



# Modèles et algorithmes pour la simulation du contact frottant dans les matériaux complexes : application aux milieux fibreux et granulaires

Gilles Daviet

## ► To cite this version:

Gilles Daviet. Modèles et algorithmes pour la simulation du contact frottant dans les matériaux complexes : application aux milieux fibreux et granulaires. Algorithme et structure de données [cs.DS]. Université Grenoble Alpes, 2016. Français. NNT : 2016GREAM084 . tel-01684673

**HAL Id: tel-01684673**

**<https://theses.hal.science/tel-01684673>**

Submitted on 15 Jan 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## THÈSE

Pour obtenir le grade de

### DOCTEUR DE L'UNIVERSITÉ DE GRENOBLE

Spécialité : **Mathématiques-informatique**

Arrêté ministériel du 25 mai 2016

Présentée par

**Gilles Daviet**

Thèse dirigée par **Florence Bertails-Descoubes**

préparée au sein de l' **équipe-projet Bipop, Inria et Laboratoire Jean Kuntzmann**  
et de l'école doctorale "**Mathématiques, Sciences et Technologies de l'Information, Informatique**"

## Modèles et algorithmes pour la simulation du contact frottant dans les matériaux complexes

Application aux milieux fibreux et granulaires

Thèse soutenue publiquement le **15 décembre 2016**,  
devant le jury composé de :

**Dr Georges-Henri Cottet**

Professeur des universités, Université Grenoble Alpes, Président

**Dr Robert Bridson**

Adjunct Professor, University of British Columbia et

Senior Principal Researcher for Visual Effects, Autodesk, Rapporteur

**Dr Ioan Ionescu**

Professeur des universités, Université Paris 13, Rapporteur

**Dr Jean-Marie Aubry**

HDR, Visual Effects Researcher, Double Negative Visual Effects, Examineur

**Dr Pierre-Yves Lagrée**

Directeur de recherche, Institut d'Alembert, CNRS, Examineur

**Dr Pierre Saramito**

Directeur de recherche, Laboratoire Jean Kuntzmann, CNRS, Examineur

**Dr Florence Bertails-Descoubes**

Chargée de recherche, Inria Rhône-Alpes, Directrice de thèse





## Remerciements

En premier lieu, je tiens bien sûr à remercier les rapporteurs et examinateurs de cette thèse, en particulier pour leur patience quand à la lecture de ce manuscrit — qui, reconnaissons le, comporte quelques sections pouvant s'avérer arides. Je remercie également le labex PERSYVAL-Lab<sup>1</sup> de m'avoir octroyé la bourse m'ayant finalement permis d'effectuer cette thèse, après une demi-douzaine de refus auprès d'organismes variés — et évidemment, Florence et Pierre, pour ne pas avoir désespéré malgré lesdits refus. Merci en outre à Pierre pour m'avoir initié à la simulation des fluides complexes, et Florence, pour m'avoir recueilli tout fraîchement sorti de l'Ensimag et m'avoir donné goût à ce domaine bien spécifique de l'animation pour l'informatique graphique, puis m'avoir porté en tant qu'ingénieur de recherche et doctorant. Je remercie l'équipe Bipop et l'Inria de m'avoir accueilli près de six ans, et tous ceux, permanents ou non-permanents, que j'ai pu rencontré dans ce cadre, et qui m'ont transmis quelques intuitions quand à la mécanique du contact, l'optimisation, ou d'autres sujets moins scientifiques. Merci en particulier aux précédents doctorants de Bipop et de la Tour pour leur rôle de modèles exemplaires : Florent, dont le manuscrit aura été une source d'informations inestimable ; Alexandre, Romain, Sofia, mes co-bureaux m'ayant préparé à la Rédaction, vraisemblablement avec succès. Merci également à mes anciens collègues de Weta Digital, qui ont largement contribué à élargir mes horizons culturels et scientifiques.

Car cette thèse n'aurait pu s'achever sans les moments de détente qui l'ont entrecoupé, je remercie les InuIts (et apparentés) pour les sorties à Grenoble et ailleurs, les randonnées et autres bivouacs — et parmi eux mes illustres prédécesseurs thésards, pour m'avoir encouragé à continuer dans cette voie et pour la richesse de leurs discussions, en particulier en fin de Fam's.

Finalement, je remercie ma famille, non seulement pour son support inconditionnel mais pour avoir initié ma curiosité scientifique et mon goût pour la programmation (à travers le GW-BASIC !), sans lesquels je n'aurais sans doute jamais entrepris de telles études. Plus que tout, merci à Anaïs, qui m'aura supporté et encouragé tout au long de cette thèse, et m'aura accompagné au bout du monde.

---

<sup>1</sup> This work has been partially supported by the LabEx PERSYVAL-Lab (ANR-11-LABX-0025-01) funded by the French program Investissement d'avenir



# Contents

<b>Remerciements</b>	<b>3</b>
<b>Contents</b>	<b>5</b>
<b>Nomenclature</b>	<b>11</b>
<b>Introduction</b>	<b>15</b>
0.1 Motivation . . . . .	15
0.1.1 Granular materials . . . . .	15
0.1.2 Dynamics of hair and fur . . . . .	16
0.1.3 Target applications . . . . .	16
0.2 Contacts and dry friction . . . . .	18
0.2.1 Impacts . . . . .	18
0.2.2 Dry friction . . . . .	19
0.2.3 Other friction laws . . . . .	20
0.2.4 Discrete simulation of complex materials with frictional contacts . . . . .	22
0.3 Continuum modeling of dry friction . . . . .	23
0.3.1 Yield-stress flows . . . . .	23
0.3.2 Frictional yield surfaces . . . . .	23
0.3.3 Shearing granular flows . . . . .	25
0.3.4 Other complex materials . . . . .	26
0.4 Synopsis . . . . .	26
<b>I Numerical treatment of friction in discrete contact mechanics</b>	<b>29</b>
<b>1 Mathematical structure of Coulomb friction</b>	<b>31</b>
1.1 Coulomb's friction law . . . . .	31
1.1.1 Second-Order Cone . . . . .	31
1.1.2 Disjunctive formulation of the Signorini-Coulomb conditions . . . . .	32
1.1.3 Alart-Curnier function . . . . .	33
1.2 Implicit Standard Materials . . . . .	34
1.2.1 Generalized Standard Materials . . . . .	34
1.2.2 Implicit Standard Materials . . . . .	35
1.3 Application to Drucker-Prager plasticity . . . . .	38
1.3.1 Symmetric tensors . . . . .	38
1.3.2 Drucker-Prager yield surface . . . . .	39
1.3.3 Dilatancy and non-associated Drucker-Prager flow rule . . . . .	40
1.3.4 Bipotential and reformulations of the Drucker-Prager flow rule . . . . .	41
1.3.5 Viscoplasticity . . . . .	44
<b>2 Modeling contacts within the Discrete Element Method</b>	<b>47</b>
2.1 A few mechanical models for rigid and deformable bodies in finite dimension . . . . .	47
2.1.1 Rigid-body dynamics . . . . .	47
2.1.2 Lumped system . . . . .	51

2.1.3	Lagrangian mechanics	52
2.1.4	Discussion	54
2.2	Contacts	55
2.2.1	Continuous-time equations of motion with contacts	56
2.2.2	Time integration	57
2.2.3	Collision detection	60
2.3	Discrete Coulomb Friction Problem	61
2.3.1	Reduced formulation	61
2.3.2	Fixed-point algorithms and existence criterion	62
<b>3</b>	<b>Solving the Discrete Coulomb Friction Problem</b>	<b>67</b>
3.1	Global strategies	67
3.1.1	Pyramidal friction cone	67
3.1.2	Complementarity functions	68
3.1.3	Optimization-based methods	69
3.2	Interior-point methods	70
3.2.1	Second-Order Cone Programs	70
3.2.2	Discussion	71
3.3	First-order proximal methods	71
3.3.1	Proximal operator	71
3.3.2	Projected Gradient Descent	73
3.3.3	Primal–dual proximal methods	74
3.4	Splitting methods	78
3.4.1	Operator splitting	78
3.4.2	Convergence properties	79
3.4.3	Performance	80
3.4.4	Discussion	82
<b>4</b>	<b>A Robust Gauss–Seidel Solver and its Applications</b>	<b>83</b>
4.1	Hybrid Gauss–Seidel algorithm	83
4.1.1	SOC Fischer-Burmeister function	83
4.1.2	Analytical solver	85
4.1.3	Full algorithm	88
4.2	Application to hair dynamics	89
4.2.1	Hair simulation in Computer Graphics	91
4.2.2	Full-scale simulations	92
4.2.3	Friction solvers comparisons	94
4.2.4	Limitations	97
4.3	Application to cloth simulation	98
4.3.1	Nodal algorithm	98
4.3.2	Results	100
4.3.3	Limitations	100
4.4	Inverse modeling with frictional contacts	101
4.4.1	Linear case	101
4.4.2	Nonlinear case	103
<b>II</b>	<b>Continuum simulation of granular materials</b>	<b>107</b>
<b>5</b>	<b>Continuum simulation of granular flows</b>	<b>109</b>
5.1	Continuum models	109
5.1.1	Inelastic yield-stress fluids	110
5.1.2	Elasto-plastic models	111
5.2	Spatial discretization strategies	112

5.2.1	Continuous conservation equations	112
5.2.2	Mesh-based discretization	113
5.2.3	Particle-based discretization	115
5.2.4	Hybrid methods	116
5.3	Our approach	117
5.3.1	Design goals	117
5.3.2	Outline of this second part	118
<b>6</b>	<b>Dense granular flows</b>	<b>119</b>
6.1	Constitutive equations	119
6.1.1	Unilateral incompressibility	119
6.1.2	Friction	120
6.2	Creeping flow	122
6.2.1	Steady-state and boundary conditions	122
6.2.2	Variational formulation	123
6.2.3	Cadoux algorithm	124
6.3	Discretization using finite-elements	126
6.3.1	Discretization of the symmetric tensor fields	126
6.3.2	Discretization of the (bi)linear forms	128
6.3.3	Discretization of the Drucker–Prager flow rule	128
6.3.4	Considerations on $R^+$	131
6.3.5	Final discrete system	132
6.4	Solving the discrete problem	133
6.4.1	Discrete Cadoux fixed-point algorithm	133
6.4.2	Dual problem	134
6.4.3	Solving the minimization problems	135
6.5	Results	136
6.5.1	Model problems	136
6.5.2	Flow around a cylinder	139
6.5.3	Extension to inertial flows: discharge of a silo	140
6.5.4	Performance	143
<b>7</b>	<b>Dry granular flows</b>	<b>149</b>
7.1	Spatially continuous model	149
7.1.1	Constitutive equations	149
7.1.2	Energy considerations	150
7.1.3	Semi-implicit integration	152
7.1.4	Discrete-time equations	153
7.2	Discretization using finite elements	155
7.2.1	Space compability criterion	155
7.2.2	Piecewise-constant discretization	156
7.2.3	Results	160
7.3	Material Point Method	162
7.3.1	Application to our variational formulation	163
7.3.2	Grid–particles transfers	163
7.3.3	Shape functions	165
7.3.4	Numerical resolution	167
7.3.5	Overview of a time-step	169
7.4	Extensions	169
7.4.1	Rigid body coupling and frictional boundaries	170
7.4.2	Anisotropy	172
7.5	Results	173
7.5.1	Model problems	174
7.5.2	Complex scenarios	176



7.5.3	Performance	178
7.5.4	Limitations	179
7.6	Discussion	180
<b>8</b>	<b>Granular flows inside a fluid</b>	<b>181</b>
8.1	Related work	181
8.1.1	Modeling	182
8.1.2	Numerical simulations	184
8.2	Two-phase model	186
8.2.1	Base equations	186
8.2.2	Stresses and buoyancy	186
8.2.3	Drag force	188
8.2.4	Mixture conservation equations	189
8.2.5	Dimensionless equations	191
8.2.6	Particular cases	193
8.3	Numerical resolution of the two-phase equations	195
8.3.1	Time discretization	195
8.3.2	Variational formulation	196
8.3.3	Discrete system	197
8.3.4	Spatial discretization	198
8.4	Results	199
8.4.1	Rayleigh-Taylor instability	199
8.4.2	Sedimentation	199
8.4.3	Regimes	200
8.4.4	Limitations	202
8.4.5	Conclusion	203
	<b>Conclusion</b>	<b>205</b>
9.1	Key remarks and summary of contributions	205
9.2	Perspectives	206
	<b>Appendices</b>	<b>209</b>
<b>A</b>	<b>Convex analysis</b>	<b>211</b>
A.1	Operations on convex functions	211
A.1.1	Fundamental definitions	211
A.1.2	Subdifferential of a function	211
A.1.3	Convex conjugate	213
A.2	Normal and convex cones	214
A.2.1	Normal cone	214
A.2.2	Operations on normal cones	215
A.2.3	Convex cones	217
A.3	Constrained optimization	218
A.3.1	Optimality conditions	219
A.3.2	Lagrange multipliers	220
<b>B</b>	<b>Discrete Coulomb Friction Problem solvers</b>	<b>227</b>
B.1	Newton SOC Fischer–Burmeister function	227
B.1.1	SOC Fischer–Burmeister derivatives	227
B.1.2	Optimistic Newton algorithm	228
B.2	Suggested variants of the Projected Gradient Descent algorithm	228
B.3	Convergence of the out-of-order Gauss–Seidel algorithm	229
<b>C</b>	<b>Supplemental justifications related to Drucker–Prager constraints</b>	<b>233</b>

C.1	Constraints on quadrature points . . . . .	233
C.2	Frictional boundaries . . . . .	234
C.2.1	Signorini condition . . . . .	234
C.2.2	Tangential reaction . . . . .	234
C.2.3	Reverse inclusion . . . . .	235
<b>Bibliography</b>		<b>237</b>
<b>Abstract – Résumé</b>		<b>252</b>



# Nomenclature

## Abbreviations

ADMM	Alternating Direction Method of Multipliers (Section 3.3.3)
AMA	Alternating Minimization Algorithm (Section 3.3.3)
DCFP	Discrete Coulomb Friction Problem (Section 2.3)
DEM	Discrete-Element Modeling (Section 0.2.4)
FEM	Finite-Element Modeling (Section 5.2.2)
FLIP	FLuid-Implicit Particle method (Section 5.2.4)
GSM	Generalized Standard Material (Section 1.2.1)
ISM	Implicit Standard Material (Section 1.2.2)
MPM	Material Point Method (Section 5.2.4)
NSCD	Non-Smooth Contact Dynamics (Section 3.4)
PIC	Particle-in-Cell method (Section 5.2.4)
SOC	Second-Order Cone (Definition 1.1)
SOCQP	Second-Order Cone Program (Section 3.2)
SOCQP	Second-Order-Cone Quadratic program (Section 2.3.2)
SOR	Successive Over-Relaxations (Section 3.4.3)
DP	Drucker–Prager (yield surface, Section 0.3.2)
MC	Mohr–Coulomb (yield surface, Section 0.3.2)

## Constants

Bi	Bingham number
Fr	Froude number
Re	Reynolds number
St	Stokes number
$d$	Dimension of the simulation space (usually 2 or 3)
$\alpha$	(Chapter 8 only) Scaled density difference, $\alpha := \frac{\rho_g - \rho_f}{\rho_f}$
$\epsilon$	(Chapter 8 only) Ratio of scales, $\epsilon := D_g/L$
$\mathbf{e}_i$	$i^{\text{th}}$ component of the canonical basis of an Euclidean space
$\mathbb{I}_n$	Identity tensor of dimension $n$ ( $n$ might be omitted)
$\boldsymbol{\iota}_d$	Unit normal tensor for the scalar product $S_d^2 \rightarrow \mathbb{R}$ , $\boldsymbol{\sigma}, \boldsymbol{\tau} \mapsto \frac{1}{2} \cdot \boldsymbol{\sigma} : \boldsymbol{\tau}$
$\mathfrak{s}_d$	Dimension of the space $S_d$ of symmetric tensors with $d$ rows and columns, $\mathfrak{s}_d = \frac{1}{2}d(d+1)$

$\Omega$	Simulation domain (for continuum mechanics) or rigid-body
$B_D$	Part of the boundary of $\Omega$ with Dirichlet boundary conditions
$B_N$	Part of the boundary of $\Omega$ with Neumann boundary conditions
$\Delta_t$	Timestep size
$\zeta$	Dilatancy coefficient (in the context of granular flows)
$\eta$	Dynamic viscosity
$\mathbf{g}$	Gravity vector
$\mu$	Coefficient of friction
$\mathbf{n}$	Normal vector (for a given contact point)
$\phi_{\max}$	Maximal volume fraction
$\rho$	Volumetric mass
$\sigma_S$	Shear yield stress
$\tau_c$	Tensile yield stress
$c$	Cohesion coefficient
$D_g$	Average diameter of the grains
$L$	Characteristic length

#### Differential operators

$\nabla \cdot \boldsymbol{\tau}$	Divergence of a vector or tensor field $\boldsymbol{\tau}$
$\frac{\partial f}{\partial x}$	Partial derivative or subdifferential of a function $f$ w.r.t. a single variable $x$
$\frac{D\phi}{Dt}$	Total or <i>material</i> derivative of a field $\phi$ , c.f. Section 5.2.1
$\nabla f$	Gradient of a function (or of a vector or scalar field) $f$
$W(\mathbf{v})$	Skew-symmetric part of the gradient of a vector field $\mathbf{v}$ , $W(\mathbf{v}) := \frac{1}{2}\nabla\mathbf{v} - \frac{1}{2}(\nabla\mathbf{v})^T$
$D(\mathbf{v})$	Symmetric part of the gradient of a vector field $\mathbf{v}$ , $D(\mathbf{v}) := \frac{1}{2}\nabla\mathbf{v} + \frac{1}{2}(\nabla\mathbf{v})^T$
$\partial f$	Subdifferential of a function $f$ (Definition A.6)

#### Functions

$\mathcal{I}_C$	Characteristics function of the set $C$ (Definition A.8)
$\delta_j^i$	Kronecker delta
$\delta$	Dirac delta function
$f_{AC}$	Alart–Curnier function, defined in Equations (1.6, 1.27)
$f_{BK}$	Kynch batch flux density function, defined in Section 8.1.1
$f_{DS}$	De Saxcé function, defined in Equation (1.29)
$f_{DS}$	Fischer–Burmeister function, defined in Section 4.1.1

#### Mathematical operators

$\langle \mathbf{u}, \mathbf{v} \rangle$	Dot product of vectors $\mathbf{u}$ and $\mathbf{v}$
$\boldsymbol{\tau} : \boldsymbol{\sigma}$	Twice-contracted tensor product. For rank-2 tensors, $\boldsymbol{\tau} : \boldsymbol{\sigma} = \sum_{i,j} \tau_{ij} \sigma_{ji}$
$\boldsymbol{\tau} \otimes \boldsymbol{\sigma}$	Tensor (outer) product of $\boldsymbol{\tau}$ and $\boldsymbol{\sigma}$

---

$\mathbf{u} \wedge \mathbf{v}$	Cross product of 3D vectors $\mathbf{u}$ and $\mathbf{v}$
$\Pi_C$	Orthogonal projection on a set (Corollary A.6)
$M$	Mass or stiffness matrix
$W$	Delassus operator
adj	Adjugate matrix (transpose of cofactor matrix)
atan2	A two-arguments arctangent function that is robust to edge cases <sup>2</sup>
$\langle \cdot   \cdot \rangle$	Average values of a field at a discontinuity
$f^*$	Convex conjugate of a function (Definition A.7)
Conv	Convex hull of a set
Dev	Deviatoric (traceless) part of a tensor
diag	(Block-)diagonal matrix obtained by diagonal concatenation of several coefficients or blocks
dom	Effective domain of a function (Definition A.3)
epi	Epigraph of a function (Definition A.2)
Im	Image of a linear operator
relint	Relative interior of a set
$\llbracket \cdot \rrbracket$	Jump (difference in values) of a field at a discontinuity
Ker	Kernel of a linear operator
$\cdot_N$	Normal part of a vector (w.r.t. $\mathbf{n}$ ) or tensor (w.r.t. $\boldsymbol{\iota}_d$ )
$\overline{C}$	Closure of the set $C$
prox	Proximal operator (Definition 3.1)
Span	Set spanned by linear combinations of a set of vectors
$\cdot_T$	Tangential part of a vector (w.r.t. $\mathbf{n}$ ) or tensor (w.r.t. $\boldsymbol{\iota}_d$ )
Tr	Trace of a tensor
Bd	Boundary of a set
int	Interior of a set
$I_1(\boldsymbol{\sigma})$	First invariant of a tensor $\boldsymbol{\sigma}$ , $I_1(\boldsymbol{\sigma}) := \text{Tr } \boldsymbol{\sigma}$
$J_2(\boldsymbol{\sigma})$	Second invariant of the deviatoric part of a tensor $\boldsymbol{\sigma}$ , $J_2(\boldsymbol{\sigma}) := \text{Tr } \boldsymbol{\sigma}^2$

### Sets and spaces

$\mathcal{K}_\mu$	Second-Order Cone (SOC) of aperture $\mu$ (Definition 1.1)
$\mathcal{T}_{\mu, \sigma_s}$	Truncated SOC of aperture $\mu$ and base section $\sigma_s$ (Section 1.3.2)
$H^1(\Omega)$	Sobolev space $W^{1,2}$ of square integrable functions with square-integrable derivatives over a domain $\Omega$
$H_0^1(\Omega)$	Subspace of $H^1(\Omega)$ satisfying homogeneous Dirichlet boundary conditions
$L_2(\Omega)$	Space of square integrable functions over a domain $\Omega$
$T_h \subset L_2(\Omega)^{\mathbb{S}_d}$	Discrete space of symmetric tensor fields

---

<sup>2</sup>See <https://en.wikipedia.org/wiki/Atan2>

$V_h(\mathbf{0}) \subset H_0^1(\Omega)^d$  Discrete space of velocity fields

$[a, b] \subset \mathbb{R}$  Closed interval

$]a, b[ \subset \mathbb{R}$  Open interval

$\mathcal{N}_C$  Normal cone to a set  $C$  (Definition A.9)

$K^\circ$  Polar cone to a set  $K$  (Definition A.11)

$\mathcal{C}_\mu$  Set of velocity–force solutions to the Signorini–Coulomb frictional contact law (Section 1.1.2)

$\mathcal{D}\mathcal{P}$  Set of strain–stress solutions to the non-associated Drucker–Prager rheology (Section 1.3.3)

$\tilde{\mathbb{R}}$   $\mathbb{R} \cup \{-\infty, +\infty\}$

**Variables (discrete mechanics)**

$\lambda$  Lagrange multipliers associated to holonomic constraints

$q$  Generalized coordinates

$r$  Contact reaction forces

$u$  Relative velocities of contacting objects

$\tilde{u}$  Relative velocities after *de Saxcé* change of variable (Section 1.2.2)

$v$  Generalized velocities

**Variables (granular flows)**

$\beta$  Scaled mass field,  $\beta := \alpha\phi + 1$

$\dot{\epsilon}$  Strain rate tensor,  $\dot{\epsilon} := D(\mathbf{u})$

$\epsilon$  Strain tensor

$\eta_{\text{eff}}$  Effective viscosity field

$\gamma$  Affine combination of the strain rate tensor with positive divergence

$\lambda$  Opposite of contact stress tensor

$\phi$  Volume fraction field

$\pi$  Product fraction field,  $\pi := \phi(1 - \phi)$

$\sigma$  Stress tensor

$\mathbf{u}$  Velocity field

$\xi$  Effective drag field

$(\mathbf{T}_j)$  Basis of the discrete space of symmetric tensors  $T_h$

$(T_j)$  Shape functions (basis scalar fields) for stresses and strains

$(\mathbf{V}_i)$  Basis of the discrete velocity space  $V_h(\mathbf{0})$

$(\omega_i^\nu)$  Shape functions (basis scalar fields) for stresses and strains velocities

# Introduction

Complex materials can be defined as large collections of discrete constituents; rigid bodies, slender elastic objects, or anything in between. In this work, we are particularly interested in the case where the interactions between the different constituents are mainly driven by *dry frictional contact*, and more specifically, where the *Coulomb friction* law holds. As the Coulomb model is a macroscopic approximation of the fine-scale interactions occurring between contacting surfaces, our study will be restricted to systems whose constituents are above a critical size, around  $100\mu\text{m}$ . Moreover, we will focus on materials with no fixed structure — the different constituents are free to reorganize themselves at will, and their relative motion will only be impeded by frictional contact (and possibly cohesive) forces. Natural examples of such systems include the likes of sand and scree, but also animal fur and human hair; manufactured examples can be as diverse as dry food troves or ball (and more rarely coin) pools.

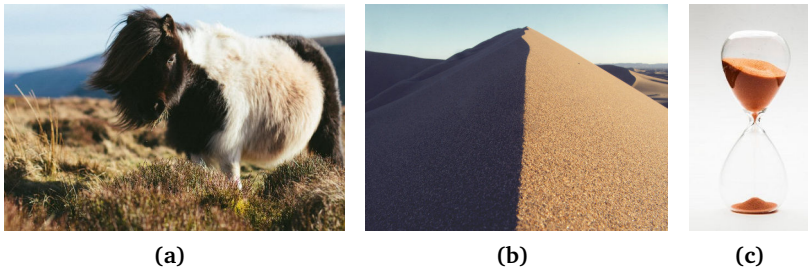
Being able to numerically reproduce the dynamics of such complex systems is important for a wide range of applications. For instance, geotechnical communities are particularly interested in the avalanching behavior of soil or gravel, while cosmetology researchers would like to assess the impact of care products on the motion of human hair. Moreover, the last decades have seen the rise of a strong demand for realism in digital special effects for feature films; the visual richness of the motion of fur, hair, or granular media have thus driven the increasing interest of the Computer Graphics community in the dynamics of complex materials.

## 0.1 Motivation

The numerical methods advocated in this dissertation were mostly motivated by two particular cases of complex materials, fiber assemblies and granular medias.

### 0.1.1 Granular materials

Granular materials (see, e.g., Andreotti et al. [2011](#) for a comprehensive description) commonly refer to a large collection of small solid grains larger than  $100\mu\text{m}$  in size — which typically



**Figure 0.1:** *Fur, herbs, and sand are examples of natural complex materials. The hourglass on the right illustrates the different dynamical regimes that can be exhibited by granular materials: liquid (above the outlet), gaseous (below the outlet), and solid (the core of the heap in the bottom compartment).*



distinguishes them from powders, made of much smaller grains. Considering this limit size, grain-grain interactions in dry granulars are mainly dictated by contact and dry friction, while air-grain interactions can be neglected. The case of immersed materials, for which interactions with the surrounding fluid can no longer be neglected, will also be treated in Chapter 8. Cohesion between grains may furthermore be considered, typically in the case of wet materials.

Being ubiquitous in outdoor environments, materials made of such grains have been heavily studied by the mechanical and geotechnical communities in the last century. They have also seen applications in a wide range of industries, including Computer Graphics. Indeed, despite their apparent simplicity, granular materials — even when constituted of rigid grains — are capable of exhibiting visually very rich dynamics. In particular, contacts and friction in such materials allow them to switch between three distinct regimes:

- a *solid* regime, when the material is maintained at rest by dry friction — for instance the core of the sand dune in Figure 0.1(b);
- a *flowing* regime, in which the material behaves like a liquid — consider the flow at the outlet of the hourglass from Figure 0.1(c), or the avalanching behavior on the outer layer of a dune;
- a *gaseous* regime, when the grains are mostly separated and only interact through sparse impacts — the flow below the hourglass' outlet, or the projections made by an impact on a granular bed.

All of these regimes (and the transitions between them) have to be properly modeled in order to produce visually convincing simulations.

### 0.1.2 Dynamics of hair and fur

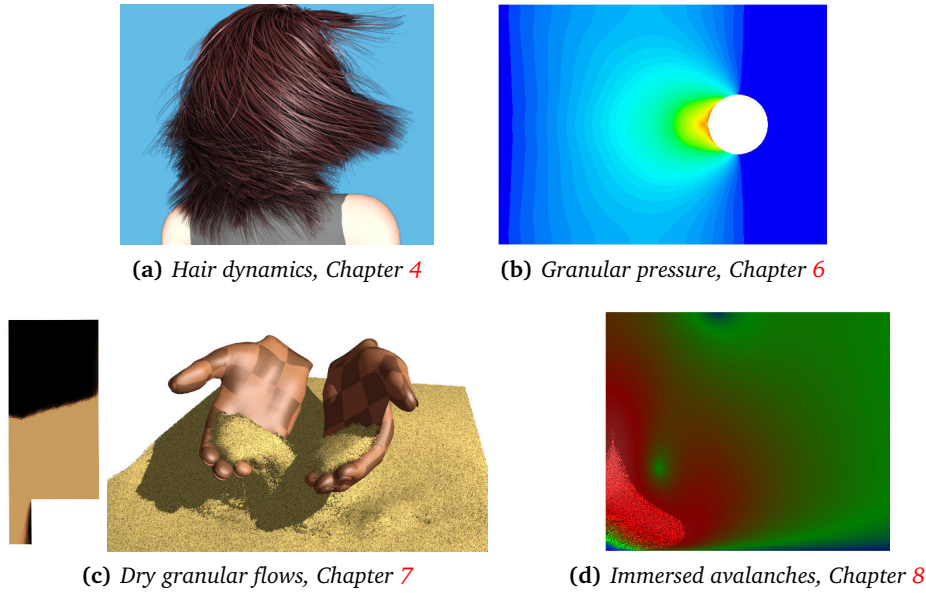
Fibrous materials feature constituents with one dimension much longer than the other ones. Driven by industrial applications, particular cases of fibrous materials have also been the subject of extensive research. For instance, the flow of polymer suspensions is critical to injection molding, and the tire engineering community is deeply interested in the study of the cords' wear by repeated small deformations. In contrast, the large-deformation dynamics of assemblies of slender elastic rods subject to frictional contacts, as is the case of hair and fur, have historically seen less interest. Yet, industrial applications such as cosmetology and digital virtual effects have recently put the spotlight on such complex materials (Ward et al. 2007).

A human head of hair consists of about 150,000 individual strands, which are very elongated, with a diameter of about  $100\mu\text{m}$  for a potential length of dozens of centimeters. Conversely, animal fur such as in Figure 0.1(a) may contain millions of (generally shorter) strands. The relative importance of contact forces compared to other interactions, such as air drag or electrostatic forces, is not well known. However, it is a certainty that contacts and friction play a huge role in the appearance of hair and fur, and proper handling of these interactions is of utmost importance for Computer Graphics applications. Indeed, frictional contacts maintain the volume of the groom, and thus the silhouette of the virtual character. Dry friction is furthermore responsible for the persistence of intricate patterns at rest, and ignoring it can lead to an uncannily tidy appearance. Despite these considerations, the work that we will present in Chapter 4 was among the first to attempt to properly capture dry friction in hair simulations.

### 0.1.3 Target applications

In this dissertation, we will not attempt to quantitatively reproduce the behavior of the simulated materials in tightly controlled settings. Instead, we will be interested in capturing the qualitative characteristics of their large-scale dynamics.

This choice is partly motivated by applications to Computer Graphics, which we will discuss in more details below. Independently of the peculiarities of this industry, we believe that convincing simulations should be based on sound physics; when possible, we will also attempt to capture in our simulations the macroscopic laws that are experimentally observed to govern the materials.

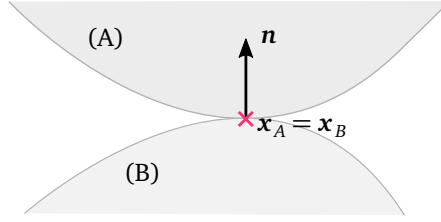


**Figure 0.2:** A few snapshots of simulations from this manuscript

Moreover, we will not solely focus on Computer Graphics; for instance, Chapters 6 and 8, in which we will study only 2D model problems, have no direct graphical application, and might be of greater interest to the mechanical engineering community.

**Computer Graphics** As already mentioned, a significant part of this work will be focused on devising simulation methods that are viable for Computing Graphics. A peculiarity of this application is that it strives to capture the emerging features that are created by the motion of the individual constituents. Indeed, while some of these features may not significantly influence the macroscopic mechanical properties of the material, they largely contribute to the visual richness of the overall phenomenon and would be extremely tedious to animate by hand. This overarching goal can thus be quite different from that of the mechanical engineering communities, so the simulation approaches will also be evaluated using different criteria. We list below a few of the virtues that numerical methods targeted at Computer Graphics should meet.

- *realism*, but not *accuracy*. While the human mind is good at pointing out things that “feel off”, it is a poor judge about whether a simulated is physically accurate or not, especially as part of a heavily stylized movie. As such, we want to be able to capture qualitative features of our complex materials (for instance, the distinct regimes of granulars), but do not necessarily want to solve our equations to a high precision.
- *artifacts-free*. The resulting simulations must be free of disturbing visual artifacts, such as *flickering*, *popping*, *privileged directions* or *creeping*. Proper modeling of dry friction is necessary to avoid this last pitfall, and care has to be taken for the underlying discretization never to be visible.
- *controllability*. While physical realism is a good thing, being pleasant to the client’s eye is a prime requirement for any Computer Graphics simulation. If a gravity-defying hair wisp is required to achieve the desired look, then the numerical method should be able to handle it. Here, we will not worry too much about this aspect. However, we will prefer models based on measurable physical parameters, and shall ensure that implementation details such as a “number of solver iterations” do not affect too much the simulated physics.
- *computational efficiency*. As providing a good set of control parameters is a hard problem (see also Sigal et al. 2015), a trial-and-error process is often required for artists to ob-



**Figure 0.3:** Two bodies (A) and (B) are in contact when two points of their respective surfaces,  $\mathbf{x}_A$  and  $\mathbf{x}_B$ , coalesce. The resulting contact normal,  $\mathbf{n}$ , is arbitrarily defined to point towards A.

tain a satisfying look. To make this process less painful, the numerical method should be reasonably efficient, and simulations for most shots should be able to finish overnight.

## 0.2 Contacts and dry friction

Proper modeling of individual constituents (for instance, in the case of fibrous materials, choosing an adequate mechanical model for individual fibers) is obviously of primary importance for computing the dynamics of any complex material. However, we will not discuss this topic in this dissertation, and will simply rely on existing models from the literature. We will instead focus on the modeling and simulation of contacts and dry friction inside the material. We provide below a brief introduction to these physical phenomena, and present the modeling choices that will underpin the remainder of this dissertation.

### 0.2.1 Impacts

We assume the distinct constituents of our complex material to be large enough that they can be considered to never overlap. Impacts happen when two previously disjoint objects come into contact; that is, when their (initially positive) relative distance drops to zero. Assuming sufficient smoothness of their boundaries, the two bodies will then share a normal direction along the contacting surfaces, and further relative motion of each pair of contacting points will be restricted to the half-space spanned by this normal. The simplest scenario, where two locally smooth and convex bodies (A) and (B) come into contact, is illustrated in Figure 0.3. Let  $\mathbf{x}_A(t)$  and  $\mathbf{x}_B(t)$  denote the position over time, for each object, of the surface point that will take part in the contact. The gap function,  $\mathbf{h}(t) := \mathbf{x}_A(t) - \mathbf{x}_B(t)$ , gives the relative position of those contacting points. The condition that objects (A) and (B) should not overlap can be written as  $\langle \mathbf{h}(t), \mathbf{n} \rangle \geq 0$ , where  $\mathbf{n}$  is the normal to (B) at  $\mathbf{x}_B$  and  $\langle \cdot, \cdot \rangle$  denotes the usual scalar product. The two objects are in contact at instant  $t$  if  $\mathbf{h}(t) = \mathbf{0}$ ; as long as this is the case, the normal relative velocity,  $\mathbf{u}_N(t) := \left\langle \frac{d\mathbf{h}}{dt}, \mathbf{n} \right\rangle$ , should remain positive. Note that  $\mathbf{u}$  can be discontinuous at the time of impact; however, following Moreau (1988), we will assume locally bounded variations of the relative velocity  $\mathbf{u}$ , i.e., the existence of a left-limit  $\mathbf{u}(t^-)$  and a right-limit  $\mathbf{u}(t^+)$  at every instant  $t$ .

At the onset of contact, a finitely-elastic body will compress, storing potential energy in the process, then restitute this energy in a second phase. Note that the amount of restituted energy does not depend on the “hardness” of the material, but rather on its internal structure; if this amount is high enough, the objects may end up separating themselves. If the elastic body is very stiff, these compression and decompression phases may happen on a time scale which is much lower than that of the studied system dynamics (Cadoux 2009, Section 1.1.1). In order to avoid having to explicitly simulate this fine time scale, one may simply model the impact as an instantaneous jump of the relative normal velocities; several models have been proposed for this purpose. The simplest of them, the empirical Newton impact law, simply states that the post-impact normal velocity should be opposed and proportional to the pre-impact one, with a

material-dependent *restitution* coefficient. On our simple example, this means that  $\mathbf{u}_N(t_i^+) = -\xi \mathbf{u}_N(t_i^-) \geq 0$ ,  $t_i$  denotes the time of impact and where  $0 \leq \xi \leq 1$  is the restitution coefficient. Note that this naive law may yield incorrect results in the presence of simultaneous contacts, and for instance will fail to reproduce the alternation in the contact points of a rigid block rocking on the ground; see (Brogliato 1999) for more discussion about impact laws.

In this work, we will focus on the simplest case of purely inelastic impacts — the energy will be instantaneously dissipated by the system. We shall thus enforce the post-impact normal velocity,  $\mathbf{u}_N(t_i^+)$ , to always be null, and refer to (Cadoux 2009, Section 1.1.5) and (Smith et al. 2012) for suggestions about how the numerical framework used throughout this dissertation may be adapted to handle Newton-like impacts.

**Signorini conditions** Physically, the inter-penetration of the two objects is prevented by the onset of a contact-force,  $\mathbf{r}$ . In the absence of friction and cohesion, this force should be colinear to the contact normal, i.e.,  $\mathbf{r} = r_N \mathbf{n}$  with  $r_N \geq 0$ . Now, suppose that the two objects are in contact at time  $t$ , i.e.,  $\mathbf{h}_N(t) = 0$ . We have already stated the following conditions:

1. The normal relative velocity should be positive as long as the points are in contact, i.e.,  $\mathbf{h}_N = 0 \implies \mathbf{u}_N \geq 0$ .
2. If  $t$  is a time of impact, the post-impact normal velocity,  $\mathbf{u}_N(t^+)$ , should vanish.
3. The normal contact force should be positive, and vanish when  $\mathbf{h}_N > 0$ .

The Signorini conditions are constructed by considering each contact in isolation, or more precisely, by assuming that the shock due to an impact does not propagate to other contacts. This is consistent with our choice of a purely inelastic impact law, which already forbids energy restitution. In the more general setting of a Newton impact law, this means that the pre-impact velocities are computed independently for each contact point, or again that the work of the normal contact force at already existing contacts should not be strictly positive, i.e.,  $\mathbf{u}_N(t^-) = 0 \implies \mathbf{u}_N(t^+) \mathbf{r}_N(t^+) \leq 0$ . Another characterization of this hypothesis is that if  $t$  is not a time of impact for the considered contact, then the contact force should locally be of bounded variations at  $t$ , i.e., should possess both left and right limits. Note that this strategy is unable to correctly model the famous Newton cradle, as the impacted balls would be incorrectly predicted to stick together (see Smith et al. 2012, Figure 3, bottom).

Combining this new implication with the three previous ones, we get that

$$\begin{cases} \mathbf{u}_N(t^+) \geq 0 \\ \mathbf{r}_N(t^+) = 0 & \text{if } \mathbf{u}_N(t^+) > 0 \\ \mathbf{r}_N(t^+) \geq 0 & \text{if } \mathbf{u}_N(t^+) = 0 \end{cases} \quad (1)$$

which together are known as the *Signorini* conditions. Remember however that these velocity-level conditions apply only when the objects are in contact at instant  $t$ , i.e., when  $\mathbf{h}_N(t) = 0$ .

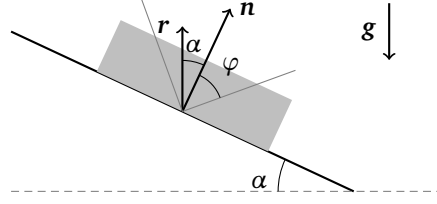
The Signorini conditions are more commonly written in a more compact manner, making the use of complementarity notation and dropping the time variable, as

$$0 \leq \mathbf{u}_N \perp \mathbf{r}_N \geq 0,$$

where the  $\mathbf{u}_N \perp \mathbf{r}_N$  notation means that the two variables should be orthogonal, i.e.,  $\mathbf{r}_N \mathbf{u}_N = 0$ .

## 0.2.2 Dry friction

The set of inequalities commonly referred to as the Coulomb friction law is actually the result of observations made by several authors and over the span of many centuries (Besson 2007). The discovery of the proportionality between the maximal tangential friction force and the normal load, as well as the irrelevance of the apparent contact surface area, are attributed to Leonard de Vinci at the end of the fifteenth century, and independently to Guillaume Amontons two hundred years later. Leonhard Euler distinguishes the *static* (or sticking) regime, when there is



**Figure 0.4:** Euler observed that as long as the angle  $\alpha$  of an inclined plane remains below the friction angle  $\varphi = \arctan \mu_s$ , the ratio of the tangential to normal components of the contact force  $\mathbf{r}$  will remain below  $\mu_s$ , and the body (gray) cannot not slide.

no relative motion between the two objects, from the dynamic regime, when one object is sliding on top of the other. Euler also introduced the notation  $\mu_s$  for the static *friction coefficient*, i.e., the maximum ratio between the tangential and normal components of the reaction force, and related this coefficient to the maximum angle  $\varphi$  at which a mass may rest on an inclined plane without sliding as  $\mu_s = \tan \varphi$  (see Figure 0.4).

In his famed manuscript, *Théorie des machines simples: en ayant égard au frottement de leurs parties et à la roideur des cordages*, Coulomb (1781) compiled the results of several experiments, validating previous theories and noting that in the dynamic regime, the friction coefficient was independent of the sliding velocity. He also observed that for most materials, the static friction coefficient,  $\mu_s$ , was higher than the dynamic one,  $\mu_d$ . Overall, Coulomb observed the relationship between the normal and tangential forces as obeying

$$\begin{cases} \|\mathbf{r}_T\| \leq \mu_s \mathbf{r}_N & \text{if } \mathbf{u}_T = \mathbf{0} \quad (\text{static regime}) \\ \|\mathbf{r}_T\| = \mu_d \mathbf{r}_N & \text{if } \mathbf{u}_T \neq \mathbf{0} \quad (\text{dynamic regime}), \end{cases} \quad (2)$$

where the  $\cdot_T$  denotes the tangential part of the reaction force and relative velocity vectors, e.g.,  $\mathbf{r}_T = \mathbf{r} - \mathbf{r}_N \mathbf{n}$ . Incidentally, the proportionality of friction to the applied load was initially postulated by Amontons to be the result of the upper object having to elevate itself above the fine-scale irregularities of the contact surface. However, investigations in the twentieth century showed that this relationship is actually caused by an increase of the microscopic-level contact area when a higher normal load is applied (Bowden and Tabor 1950).

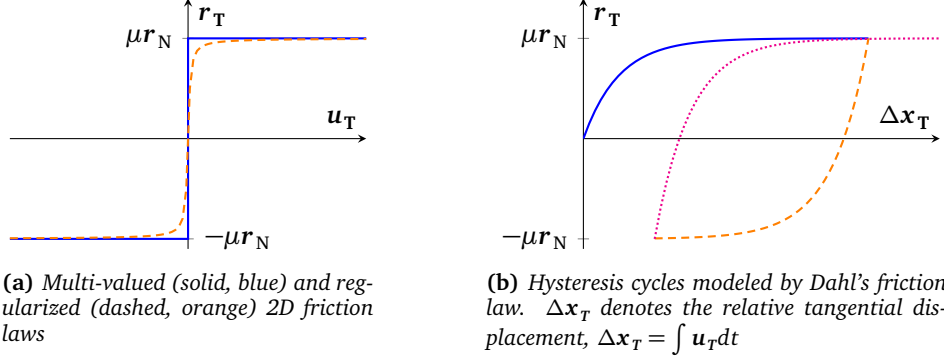
In this work, we will not distinguish between the static and dynamic friction coefficients. Indeed, the main characteristics of Coulomb friction, such as the existence of a sliding threshold that depends on the applied normal load, can already be captured without making this distinction, and we did not judge the gain in realism brought by the introduction of a distinct sliding friction coefficient worth the significant associated increase in mathematical complexity<sup>3</sup>. Taking into account the fact that the tangential friction force must oppose the sliding velocity, we will thus consider the Coulomb friction law as defined by the disjunction (3),

$$\begin{cases} \|\mathbf{r}_T\| \leq \mu \mathbf{r}_N & \text{if } \mathbf{u}_T = \mathbf{0} \\ \begin{cases} \|\mathbf{r}_T\| = \mu \mathbf{r}_N \\ \mathbf{r}_T = -\alpha \mathbf{u}_T, \alpha \in \mathbb{R}_+ \end{cases} & \text{if } \mathbf{u}_T \neq \mathbf{0}. \end{cases} \quad (3)$$

### 0.2.3 Other friction laws

While our work will be focused on Coulomb friction, we mention below a few other friction laws that are worthy of interest.

<sup>3</sup>In discrete-time numerical algorithms, the friction coefficient can always be updated explicitly at each timestep depending on the status of each contact point.



**Figure 0.5:** Some alternative friction laws: regularized (left) and Dahl's (right)

**Regularization** In contrast to a fluid (or *viscous*) friction law, which could be defined as, say,  $\mathbf{r}_T = -\eta(\mathbf{u})\mathbf{u}_T$ , the Coulomb friction law (3) is *multi-valued*. For  $\mathbf{u} = \mathbf{0}$ , the friction force is not uniquely defined, but may lie anywhere inside a ball of radius  $\mu r_N$ . As such, dealing with Coulomb friction will require devising specialized numerical method (Acary and Brogliato 2008). To avoid this complexity, one may choose to *regularize* the law, writing  $\mathbf{r}_T = -\alpha(\mathbf{u}_T, \mathbf{r}_N)\mathbf{u}_T$  with, for instance,  $\alpha(\mathbf{u}_T, \mathbf{r}_N) := \min(\mu r_N, \frac{1}{\epsilon}\|\mathbf{u}_T\|)/\|\mathbf{u}_T\|$ , or  $\alpha(\mathbf{u}_T, \mathbf{r}_N) := \mu r_N \frac{2}{\pi} (\arctan \frac{1}{\epsilon}\|\mathbf{u}_T\|)/\|\mathbf{u}_T\|$  with  $\epsilon$  small, as in Figure 0.5(a). However, such approaches may result in very stiff numerical systems, prone to flickering, and allow the lingering of creeping residual velocities; we will thus avoid this strategy.

**Tresca model** The Tresca friction law can be derived from the Coulomb law by removing the dependency of the friction force on the normal applied load,

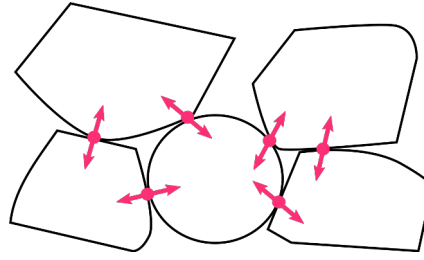
$$\begin{cases} \|\mathbf{r}_T\| \leq s & \text{if } \mathbf{u}_T = \mathbf{0} \\ \|\mathbf{r}_T\| = s & \text{if } \mathbf{u}_T \neq \mathbf{0} \\ \mathbf{r}_T = -\alpha \mathbf{u}_T, \alpha \in \mathbb{R}_+ & \end{cases}$$

where  $s$  is a positive scalar. As it reduces the coupling between tangential and normal components, the Tresca law is easier to handle numerically, yet still models a proper sliding threshold. However, the lack of proportionality of the frictional force to the applied load forbids the modeling of arbitrarily-sized stable heaps using Tresca friction, and we will thus discard this approximation.

**Dahl model** Taking inspiration from standard strain–stress diagrams, the Dahl (1968) friction model describes the evolution of the friction force w.r.t. the relative tangential displacement, with a slope “reset” at each change of sign of the tangential velocity  $\mathbf{u}_T$ . As depicted in Figure 0.5(b), this model is able to capture hysteresis cycles induced by friction in loading–unloading experiments, with a smooth reversing of the friction force, and has been especially popular to macroscopically account for the displacement of textile fibers under stretching and bending (Miguel et al. 2013; Ngo Ngoc and Boivin 2004). This smooth reversing of the friction force models slack in the sticking contact regime, and is therefore not really relevant for contacts between stiff bodies, which are close to slackless; we will thus not consider Dahl's law in the following.

**Rolling friction** Rolling friction is induced by the deformation of a wheel near its contact point with the ground, and is responsible for energy dissipation that is not captured by the tangential Coulomb friction law (3). Indeed, Equation (3) implies that the work of the friction force is non-zero only when sliding occurs. However, in the remainder of this dissertation we will focus on stiff materials and relatively small applied loads, and will thus neglect rolling friction.





**Figure 0.6:** *With Discrete Element Modeling, constituents are simulated individually, and all interactions between neighboring bodies must be taken into account.*

#### 0.2.4 Discrete simulation of complex materials with frictional contacts

The most natural way to simulate complex materials numerically would be to follow the framework of Discrete Element Modeling (DEM),<sup>4</sup> that is, simulating individually each body and its interactions with the surrounding ones, as illustrated in Figure 0.6.

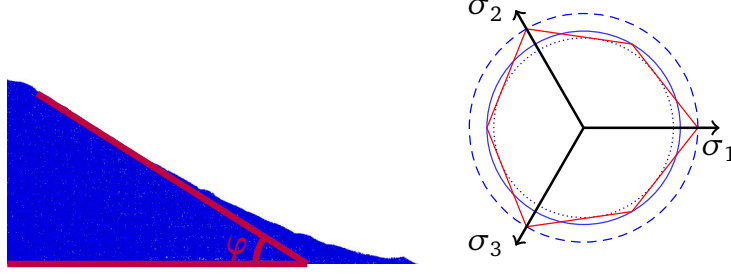
As contacts and dry friction between the grains plays a primary role in the dynamics of complex materials, special attention should be given to the numerical treatment of those phenomena. Different classes of approaches have been proposed in the literature, of which we can cite three:

- *Molecular Dynamics (MD)*, which relax the assumption that distinct bodies cannot overlap, and use nonlinear springs to model the contacts between the particles (Cundall and Strack 1979). While this approach is the simplest to implement, limiting interpenetration can require the use of very stiff springs, which introduces a time scale much smaller than that of the macroscopic dynamics. This makes stable numerical integration difficult to achieve unless very small time steps are used, and may lead to visually disturbing flickering effects. Moreover, how to handle multi-valued friction laws in the MD framework is not obvious, and as such creeping residual motion may plague this approach.
- *Constraint-based approaches*, such as the *Non-Smooth Contact Dynamics* (NSCD; Jean 1999), propose to solve for each object's dynamics while ensuring that the Signorini-Coulomb conditions (1) and (3) are satisfied. While being inherently more complex than MD, constraint-based approaches allow the use of larger timesteps, and thus still prove computationally efficient. Chapter 2 will be dedicated to this kind of approaches, with an emphasis on the timestepping scheme proposed by Jean and Moreau (1987).

**Scaling up** The first part of this dissertation is dedicated to the simulation of dry frictional contacts in the DEM framework, using a constraint-based method; Chapter 4 presents results from the application of this strategy to hair dynamics. While those results were relatively good-looking, computational performance prevented us from simulating anything close to the whole 150,000 individual fibers of a human head of hair. This motivated us to look for alternative approaches. Moreover, even though we were only simulating a small subset of the whole hair, we kept using a *fiber* model for each simulated strand, while a *wisp* model, representing the averaged behavior of several fibers, would have been more appropriate. Assuming the existence of such a model, we could as well take this approach one step further, and simulate the whole material using continuum mechanics.

As a first step towards the simulation of very large complex materials, the second part of this dissertation embraces this strategy, but focuses on simpler systems: granular materials consisting only of rigid grains. Devising a similar macroscopic simulation method for fibrous media such as hair remained out of the reach of this thesis.

<sup>4</sup>Note that the name DEM is also commonly used as a synonym for Molecular Dynamics; we use it here in a broader sense.



**Figure 0.7:** Left: visualization of the friction angle  $\varphi$  on a 2D granular heap at rest. Right: 3D Mohr-Coulomb (red) with the three different Drucker-Prager yield surfaces (blue) in the plane of constant normal stress  $\sum \sigma_i = 3$ .

Using a continuum approach for granular materials makes sense, as they can be *extremely* large systems – a cubic meter of sand contains close to a trillion individual grains. We see immediately that simulating every grain as in the DEM framework, and taking into account each of its interactions with its neighbors, is not tractable, especially on standard computers. Moreover, the scale of inhomogeneities in granular materials is usually much smaller than the material itself. For these reasons, several constitutive laws have already been proposed in the literature to model the macroscopic behavior of granulars. The next section introduces fundamental concepts for the continuum modeling of dry frictional contact in granular materials.

### 0.3 Continuum modeling of dry friction

When the discrete constituents are sufficiently small w.r.t. the scale at which a phenomenon is studied for the material to appear spatially homogeneous, averaging processes may be used to intuit constitutive equations (or *rheologies*) on macroscopic quantities such as stress and strain.

#### 0.3.1 Yield-stress flows

For instance, the presence of large molecules in so-called Bingham plastics such as mayonnaise manifest itself at the macroscopic scale by the onset of a *yield stress*  $\sigma_s$ . This means that irreversible deformations of the material will occur, i.e., the plastic strain rate  $\dot{\epsilon}$  will be non-zero, only once the norm of the deviatoric stress,  $|\text{Dev } \boldsymbol{\sigma}|$ , has reached the critical value  $\sigma_s$ . Such materials may remain indefinitely stuck in various shapes, in contrast to Newtonian fluids which will always, albeit potentially slowly, go back to a flat shape. Note that different choices can be used for the definition of the norm  $|\cdot|$  in the above expression, yielding slightly different rheologies, but an *objectivity* criterion should always be satisfied: one must ensure that the norm is invariant to changes in the reference frame. The Bingham model is classically defined using the second invariant of the deviatoric part of the stress tensor,  $J_2(\boldsymbol{\sigma}) := \frac{1}{2} \text{Tr}(\text{Dev } \boldsymbol{\sigma})^2$ , with  $|\text{Dev } \boldsymbol{\sigma}| := \sqrt{J_2(\boldsymbol{\sigma})}$ . This definition ensures the objectivity of the model.

Such plastic phenomena are commonly described with a *yield surface*, that is, a function  $F$ , objective w.r.t. the stress tensor, such that the material remains solid while  $F(\boldsymbol{\sigma}) < 0$ , and  $F(\boldsymbol{\sigma}) = 0$  corresponds to the flowing regime where irreversible deformation occurs. The yield surface for the Bingham model is given  $F^{\text{Bl}}(\boldsymbol{\sigma}) := \sqrt{J_2(\boldsymbol{\sigma})} - \sigma_s$ .

#### 0.3.2 Frictional yield surfaces

Granular materials also exhibit a yield stress, as demonstrated by their ability to form heaps that do not (systematically) collapse over time. However, just like Coulomb friction featured a sliding threshold proportional to the normal applied load, the yield stress of dry granular materials is observed to depend on the normal (or *mean*) stress, i.e., the internal pressure.



**Mohr–Coulomb criterion** The Mohr–Coulomb (MC) criterion is the continuum mechanics generalization of Euler’s inclined plane experiment, and consider the maximum angle  $\varphi$  (the so-called *rest angle*) that the slope of a granular heap can make without starting to avalanche (Figure 0.7, left). Let us first consider the 2D case, and let  $\sigma_\varphi$  and  $\tau_\varphi$  denote the normal and shear stresses acting on an inclined plane of angle  $\varphi$ . By analogy with Euler’s criterion, the stability of the granular heap up to an angle  $\varphi$  means that the material’s yield condition can be written  $\tau_\varphi \leq \sigma_\varphi \tan \varphi$ . Mohr’s circle relates  $\sigma_\varphi$  and  $\tau_\varphi$  to the principal stresses of the material  $\sigma_1 \leq \sigma_2$ , i.e., the eigenvalues of the *applied* stress tensor, as

$$\sigma_\varphi = -\frac{\sigma_1 + \sigma_2}{2} - \frac{\sigma_2 - \sigma_1}{2} \sin \varphi, \quad \tau_\varphi = \frac{\sigma_2 - \sigma_1}{2} \cos \varphi.$$

The cohesionless MC criterion thus states that a granular material with rest angle  $\varphi$  will remain stable as long as

$$\begin{aligned} \left( -\frac{\sigma_1 + \sigma_2}{2} - \frac{\sigma_2 - \sigma_1}{2} \sin \varphi \right) \tan \varphi &\geq \frac{\sigma_2 - \sigma_1}{2} \cos \varphi \\ -\frac{\sigma_1 + \sigma_2}{2} \sin \varphi &\geq \frac{\sigma_2 - \sigma_1}{2}. \end{aligned} \quad (4)$$

The 3D version of the Mohr–Coulomb criterion considers each plane of maximum shear, and can be summarized as

$$-\frac{\sigma_1 + \sigma_3}{2} \sin \varphi \geq \frac{\sigma_3 - \sigma_1}{2}, \quad (5)$$

where  $\sigma_1 \leq \sigma_2 \leq \sigma_3$  are the eigenvalues of the material’s stress. Inequation (5) written for each potential ordering of the eigenvalues defines 6 yield planes in principal stresses space; the MC yield surface is thus an hexagon-shaped convex cone centered around the hydrostatic axis  $\sigma_1 = \sigma_2 = \sigma_3$ , as illustrated in Figure 0.7, right. This hexagon degenerates to a triangle for  $\sin \varphi = 1$ , and approaches a regular (yet vanishing) hexagon for  $\sin \varphi = 0$ .

Note that while Coulomb friction is an approximation of interactions induced by the microscopic asperities of the contacting surface between grains, Mohr–Coulomb theory averages grain-sized inhomogeneities, and is thus only valid at a much bigger scale.

**Drucker–Prager yield criterion** The Mohr–Coulomb criterion (5) involves the individual eigenvalues of the stress tensor, and is numerically unwieldy. Drucker and Prager (1952) proposed a yield surface that is defined using only invariants of the stress tensor based on the Bingham model, but with a yield-stress that grows linearly with the first invariant of the stress tensor,  $I_1(\boldsymbol{\sigma}) := \text{Tr } \boldsymbol{\sigma}$ . In the cohesionless case, the Drucker–Prager (DP) criterion is thus

$$\sqrt{J_2(\boldsymbol{\sigma})} + \hat{\mu} \frac{I_1(\boldsymbol{\sigma})}{d} \leq 0, \quad (6)$$

and  $\hat{\mu}$  is called the friction coefficient.

Note that in 2D, the Mohr–Coulomb and Drucker–Prager yield surfaces coincide. Indeed,  $I_1(\boldsymbol{\sigma}) = \sigma_1 + \sigma_2$ , and  $J_2(\boldsymbol{\sigma}) = \frac{1}{4}(\sigma_1 - \sigma_2)^2$ ; Equations (6) and (4) thus become equivalent when  $\hat{\mu} = \tan \varphi$ .

However, in 3D, direct computations yield  $J_2(\boldsymbol{\sigma}) = \frac{1}{6} \sum_{i \neq j} (\sigma_i - \sigma_j)^2$ ; the Drucker–Prager yield surface is thus a convex cone spanned by a circle centered on the hydrostatic axis (i.e., a Second-Order Cone). There is no hope for the DP and MC surfaces to fully match, but one may still choose the friction coefficient  $\hat{\mu}$  using several heuristics (illustrated in Figure 0.7, right):

- $\hat{\mu} = \frac{2\sqrt{3} \sin \varphi}{3 - \sin \varphi}$ , so that DP circumscribes MC;
- $\hat{\mu} = \frac{\sin \varphi}{\sqrt{1 + \frac{1}{3} \sin^2 \varphi}}$ , so that DP inscribes MC;
- $\hat{\mu} = \frac{2\sqrt{3} \sin \varphi}{3 + \sin \varphi}$ , so that DP interpolates MC at middle vertices.

Choice between these different values is application-dependent. For instance, risk-assessment simulations may want to use the inscribed surface, so that the predicted run-out length of an avalanche with DP will always overestimate the one using MC.

**Cohesion and tensile strength** The Drucker–Prager model can be extended to the modeling of cohesive materials, modifying the yield surface as (Alejano and Bobet 2012)

$$F_{\hat{\mu}, \hat{c}}^{\text{DP}}(\boldsymbol{\sigma}) := \hat{\mu} \frac{I_1(\boldsymbol{\sigma})}{d} - \hat{c} + \sqrt{J_2(\boldsymbol{\sigma})}. \quad (7)$$

A slightly more complex yield surface, the so-called Drucker–Prager yield surface with *tension cut-off*, may also be of interest for materials such as concrete. The cut-off dictates that the material will break when the mean tensile stress exceeds a critical value  $\widehat{\tau}_c$ ,

$$F_{\hat{\mu}, \hat{c}, \widehat{\tau}_c}^{\text{DP}}(\boldsymbol{\sigma}) := \max\left(I_1(\boldsymbol{\sigma}) - d\widehat{\tau}_c, F_{\hat{\mu}, \hat{c}}^{\text{DP}}(\boldsymbol{\sigma})\right). \quad (8)$$

The cut-off stress will influence the set of admissible stresses set only when there holds simultaneously  $I_1(\boldsymbol{\sigma}) \geq d\widehat{\tau}_c$  and  $\hat{\mu}I_1(\boldsymbol{\sigma}) \leq d\hat{c}$ , which means  $\hat{\mu}\widehat{\tau}_c \leq \hat{c}$ ; the original Drucker–Prager yield surface is retrieved when  $\hat{\mu}\widehat{\tau}_c \geq \hat{c}$ . For this reason, we will prefer parameterizing the yield surface (9) with a shear yield stress,  $\sigma_s := \hat{c} - \hat{\mu}\widehat{\tau}_c$ , rather than with the cohesion coefficient  $\hat{c}$ . In the following, we will thus write the Drucker–Prager yield surface with tension cut-off as

$$F_{\hat{\mu}, \sigma_s, \widehat{\tau}_c}^{\text{DP}}(\boldsymbol{\sigma}) := \max\left(I_1(\boldsymbol{\sigma}) - d\widehat{\tau}_c, \hat{\mu} \frac{I_1(\boldsymbol{\sigma}) - d\widehat{\tau}_c}{d} - \sigma_s + \sqrt{J_2(\boldsymbol{\sigma})}\right). \quad (9)$$

Note that the Bingham yield surface is recovered when  $\widehat{\tau}_c = +\infty$ .

**Other yield surfaces** Both the Mohr–Coulomb and Drucker–Prager yield surfaces are nonsmooth; the normal to the surface is not uniquely defined everywhere, in particular for  $\boldsymbol{\sigma} = \mathbf{0}$  in the cohesionless case. As we will see in Chapter 1, this complicates the definition of a flow rule, that is, we will not be able to unambiguously express the direction of plastic displacement as a function of the stress tensor. ‘Mast (2013) presents different strategies to circumvent this difficulty, such as using a smooth cap for the Drucker–Prager cone, or prescribing the flow to be along the hydrostatic axis when  $\boldsymbol{\sigma} = \mathbf{0}$ . Another interesting option that they explore is the use of the Matzuo–Nakai yield surface, which is smooth everywhere and better matches the hexagonal shape of the Mohr–Coulomb surface than the Drucker–Prager law.

### 0.3.3 Shearing granular flows

The ‘GDR MiDi’ group (GDR MiDi 2004) studied dense granular shearing flows, and proposed a new constitutive law that was able to match experiments quantitatively, the so-called  $\mu(I)$  rheology (Jop et al. 2006). Based on the Drucker–Prager yield criterion, this rheology suggests to vary the friction coefficient with the *inertial number*  $I$ ,

$$I(\dot{\boldsymbol{\epsilon}}, \boldsymbol{\sigma}) := \frac{\sqrt{J_2(\dot{\boldsymbol{\epsilon}})} D_g}{\sqrt{I_1(\boldsymbol{\sigma})/(d\rho_g)}},$$

where  $\dot{\boldsymbol{\epsilon}}$  is the strain rate,  $D_g$  the average diameter of grains, and  $\rho_g$  their density. This dimensionless number relates the fluctuation of the velocity at the grain scale to that of the macroscopic flow. When  $I = 0$ , the material behaves like a solid, and for very high values of  $I$ , the material becomes akin to a gas; in between lies the dense flowing regime.

**Shear-hardening friction** The higher the inertial number, the more energy will be dissipated by grain–grain interactions, and the higher the friction coefficient should be. Jop et al. (2006) propose the following expression:

$$\mu(I) = \mu_s + \frac{\mu_D - \mu_s}{I_0/I + 1},$$

with  $\mu_D \geq \mu_S$  and where  $I_0$  is a material-dependent constant. In contrast to discrete friction laws which are usually taken to be slip-weakening (e.g., Coulomb friction when  $\mu_S > \mu_D$ ), the  $\mu(I)$  rheology is thus shear-hardening.

Note that the  $\mu(I)$  rheology does not influence the rest angle of the material; at rest,  $I = 0$  and  $\mu(I) = \mu_S$ .

**Dilatancy** In a similar manner, the volume fraction of grains  $\phi$ , that is, the fraction of space occupied by the granular material, can be affected by the shear rate; when this translates into an augmentation of the flow volume, this phenomenon is known as *dilatancy*. In dense shearing flow, the volume fraction has been observed to decrease with the inertial number (GDR MiDi 2004). Roux and Radjai (1998) define the dilatancy angle  $\psi$  from the ratio between the volumetric and shear strain rates,  $\frac{1}{d}I_1(\dot{\epsilon}) \tan \psi = J_2(\dot{\epsilon})$ . They relate this angle to a critical volume fraction  $\phi_c$  as  $\psi(\phi) = \psi_0(\phi - \phi_c)$ , so that positive dilatancy occurs above the critical volume fraction, but shear tends to compress the material when  $\phi < \phi_{\max}$ .

### 0.3.4 Other complex materials

This section has made clear that devising macroscopic laws for granular materials, even when assuming perfectly rigid and spherical grains, was already complex. Continuum modeling of assemblies of elastic and anisotropic objects such as fibers would require much more sophisticated models, and was thus left out of the scope of this dissertation.

## 0.4 Synopsis

This dissertation will be divided in two parts. The first one will be mostly dedicated to Coulomb friction, going from the modeling to the numerical resolution of frictional contacts between discrete bodies. Note that the first three chapters will consist mostly in a walk through standard models and numerical methods from the literature, which Chapter 4 will present original contributions.

- Chapter 1 will first look at the mathematical structure of Signorini–Coulomb conditions, presenting a few useful reformulations and showing where such law can fit in standard plasticity theory. Structural similarities with Drucker–Prager flows will be made explicit.
- Chapter 2 will go through standard modeling of contacts in discrete mechanical systems. The Moreau–Jean scheme will be presented, and we will show that each timestep can be reduced to one or more instances of a canonical problem, which we will call a *Discrete Coulomb Friction Problem* (DCFP).
- Chapter 3 will be dedicated to numerical algorithms for solving the DCFP, discussing their relative relevance for particular problem structures.
- Chapter 4 will present an original variant of the Gauss–Seidel algorithm that proved to perform very well on Computer Graphics applications, such as hair dynamics, hair inverse modeling, and cloth dynamics.

The second part will be focused on the continuum simulation of granular materials. Taking advantage of the similarities between the Coulomb friction law and the Drucker–Prager yield surface, we will show how the numerical methods devised for discrete mechanics can still be relevant in the continuum limit.

- Chapter 5 will serve as an introduction to numerical methods in continuum mechanics, and the different strategies use for the simulation of granular materials.
- Chapter 6 will present our proposed numerical method for the simulation of dense, dry granular flows, which will relax the classical incompressibility assumption and take profit of the DCFP formalism.

- Chapter 7 will extend this method to more general flows with varying volume fraction of grains. A discretization scheme and numerical resolution strategy focused on improving relevance to Computer Graphics applications will also be presented.
- Finally, Chapter 8 will discuss the extension to granular flows in which the interactions with a surrounding fluid cannot be neglected. Once again, we will show that the dynamics of such flows can be numerically solved as a sequence of DCFP.

## Software and Publications

The numerical methods presented throughout this dissertation have been accompanied by the development of associated simulation programs. Some of them have been released as open-source: a sparse block matrix linear algebra library featuring a few DCFP solvers<sup>5</sup> and a granular simulation software<sup>6</sup> based on the Material Point Method framework of Chapter 7. The approach proposed in the first part of this manuscript has also been implemented in industrial settings, both for cosmetology and visual effects applications.

Moreover, some of the contributions presented in this dissertation have already been published. These publications are listed below, together with pointers to the corresponding parts of this manuscript.

### Peer-reviewed journals

G. Daviet and F. Bertails-Descoubes (2016a). “A Semi-Implicit Material Point Method for the Continuum Simulation of Granular Materials”. In: *ACM Transactions on Graphics*. SIGGRAPH ’16 Technical Papers 35.4, p. 13. DOI: [10.1145/2897824.2925877](https://doi.org/10.1145/2897824.2925877) (Chapter 7)

G. Daviet and F. Bertails-Descoubes (2016b). “Nonsmooth simulation of dense granular flows with pressure-dependent yield stress”. In: *Journal of Non-Newtonian Fluid Mechanics* 234, pp. 15–35. ISSN: 0377-0257. DOI: <http://dx.doi.org/10.1016/j.jnnfm.2016.04.006> (Chapter 6)

A. Derouet-Jourdan, F. Bertails-Descoubes, G. Daviet, et al. (2013). “Inverse dynamic hair modeling with frictional contact”. In: *ACM Transactions on Graphics* 32.6, pp. 1–10. ISSN: 0730-0301. DOI: [10.1145/2508363.2508398](https://doi.org/10.1145/2508363.2508398) (Chapter 4, Section 4.4.1)

G. Daviet, F. Bertails-Descoubes, and L. Boissieux (2011). “A hybrid iterative solver for robustly capturing Coulomb friction in hair dynamics”. In: *ACM Transactions on Graphics* 30.6, pp. 1–12. ISSN: 0730-0301. DOI: [10.1145/2070781.2024173](https://doi.org/10.1145/2070781.2024173) (Chapter 4, Sections 4.1 and 4.2)

F. Bertails-Descoubes et al. (2011). “A nonsmooth Newton solver for capturing exact Coulomb friction in fiber assemblies”. In: *ACM Transactions on Graphics* 30 (1), 6:1–6:14. ISSN: 0730-0301. DOI: <http://doi.acm.org/10.1145/1899404.1899410> (Chapter 3, Section 3.1.2)

### Technical reports

R. Casati et al. (2016). *Inverse Elastic Cloth Design with Contact and Friction*. Research Report. Inria Grenoble Rhône-Alpes, Université de Grenoble. URL: <https://hal.archives-ouvertes.fr/hal-01309617> (Chapter 4, Section 4.4.2)

O. Bonnefon and G. Daviet (2011). *Quartic formulation of Coulomb 3D frictional contact*. Anglais. Tech. rep. INRIA - Laboratoire Jean Kuntzmann. URL: <http://hal.archives-ouvertes.fr/inria-00553859/en/> (Chapter 4, Section 4.1.2)

### Poster

G. Daviet, F. Bertails-Descoubes, and R. Casati (2015). “Fast cloth simulation with implicit contact and exact coulomb friction”. In: *Proceedings of the 14th ACM SIGGRAPH / Eurographics Symposium on Computer Animation - SCA ’15*. DOI: [10.1145/2786784.2795139](https://doi.org/10.1145/2786784.2795139) (Chapter 4, Section 4.3)

<sup>5</sup>bogus: <http://gdaviet.fr/code/bogus>

<sup>6</sup>Sand6: <http://bipop.inrialpes.fr/~gdaviet/code/sand6>



## **Part I**

# **Numerical treatment of friction in discrete contact mechanics**



# 1 Mathematical structure of Coulomb friction

In this chapter, we present a few equivalent formulations of the Coulomb law that will prove useful for the construction of numerical algorithms. Most of the properties compiled below stem from the works of Pierre Alart and G  ry de Saxc   in the early 1990's, and rely upon convex analysis tools that were largely developed by Jean Jacques Moreau and R. Tyrrell Rockafellar in the 1960's. We refer the reader to Appendix A for a brief introduction to this theory, and the enunciation of a few properties of convex cones, subdifferentials and convex conjugates that we will make frequently use of. Note that this chapter is quite heavy on notations, and not fully required to follow the remainder of this dissertation — the equivalences established here will be appropriately referred to in the following chapters. The casual reader may choose to read Sections 1.1 and 1.3.1–1.3.3, and skip the remaining derivations.

## 1.1 Coulomb's friction law

### 1.1.1 Second-Order Cone

A convenient tool for expressing the Coulomb law and Drucker–Prager yield surface in an unified manner is the formalism of the Second-Order Cone, sometimes called Lorentz cone or ice-cream cone, which is a special case of convex cone (Definition A.10).

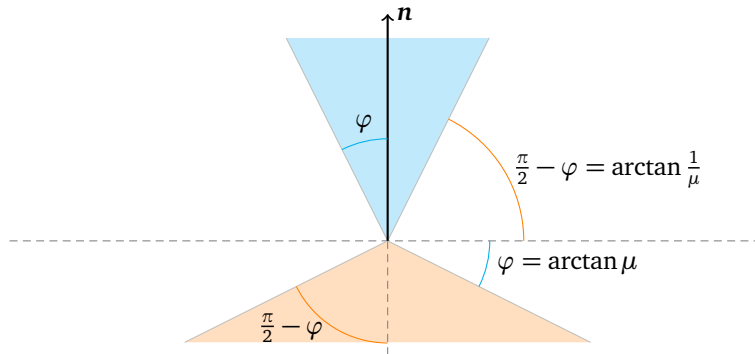
**Definition 1.1** (Second-Order Cone). *Let  $X$  be an Hilbert space,  $\langle \cdot, \cdot \rangle$  a scalar product with associated norm  $\| \cdot \| := \sqrt{\langle \cdot, \cdot \rangle}$ , and  $\mathbf{n} \neq 0 \in X$ .*

*For  $\mathbf{x} \in X$ , we note  $\mathbf{x}_N \in \mathbb{R}$  the projection of  $\mathbf{x}$  on the subspace spanned by  $\mathbf{n}$ , and  $\mathbf{x}_T \in X$  the projection on its orthogonal complement, i.e., such that  $\mathbf{x} = \mathbf{x}_N \frac{\mathbf{n}}{\|\mathbf{n}\|} + \mathbf{x}_T$  with  $\langle \mathbf{x}_T, \mathbf{n} \rangle = 0$ .*

*The Second-Order Cone of aperture  $\mu \in \mathbb{R}_+ \cup \{+\infty\}$  w.r.t.  $\mathbf{n}$  and  $\langle \cdot, \cdot \rangle$  is the closed convex cone*

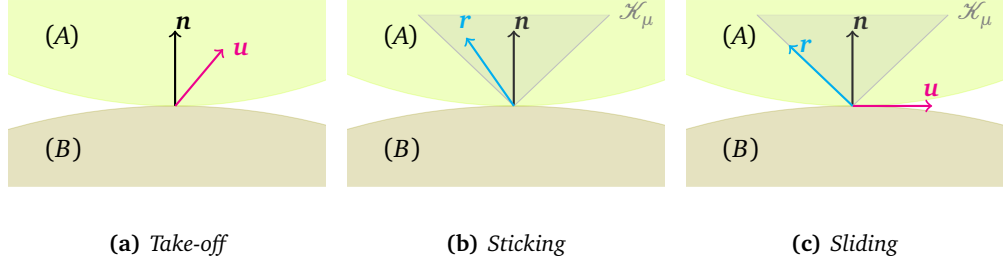
$$\mathcal{K}_\mu(\mathbf{n}) := \{ \mathbf{x} \in \mathcal{K}_\mu, \|\mathbf{x}_T\| \leq \mu \mathbf{x}_N \}.$$

*Proof.* The triangular inequality ensures that  $\mathcal{K}_\mu(\mathbf{n})$  is indeed a convex cone. □



**Figure 1.1:** The 2D Second-Order Cone of aperture  $\mu$ ,  $\mathcal{K}_\mu(\mathbf{n})$  (above, cyan), and its polar cone  $\mathcal{K}_\mu(-\mathbf{n})$  (below, orange)





**Figure 1.2:** The three cases of the disjunctive Coulomb frictional contact law

An interesting property of the class of Second-Order Cones is that it is closed w.r.t. duality, that is, the dual cone to a SOC is also a SOC. Indeed, direct geometric computations illustrated in Figure 1.1 yield that:

**Property 1.1.** The dual cone to the Second-Order Cone  $\mathcal{K}_\mu(\mathbf{n})$  is  $\mathcal{K}_{\frac{1}{\mu}}(\mathbf{n})$ , the Second Order Cone of aperture  $\frac{1}{\mu}$ .

### 1.1.2 Disjunctive formulation of the Signorini-Coulomb conditions

The Signorini (1) and Coulomb friction (3) conditions can be combined into a slightly more compact set of three cases, illustrated in Figure 1.2, which we will refer to as the *disjunctive* formulation of the Coulomb contact law (Cadoux 2009). In the *take-off* case, the two objects are separating and the reaction force vanishes; the *sticking* case holds when there is no relative motion; finally, the *sliding* corresponds to saturated friction and purely tangential relative motion. More formally, the relative velocity  $\mathbf{u}$  and the contact force  $\mathbf{r}$  should satisfy one of the following cases:

$$\left\{ \begin{array}{ll} \text{or} & \mathbf{r} = \mathbf{0} \quad \text{and} \quad u_N > 0, \quad (\text{take-off}) \\ \text{or} & \mathbf{r} \in \mathcal{K}_\mu(\mathbf{n}) \quad \text{and} \quad \mathbf{u} = \mathbf{0}, \quad (\text{sticking}) \\ \text{or} & \left\{ \begin{array}{l} \mathbf{r} \in \text{Bd } \mathcal{K}_\mu(\mathbf{n}) \\ \mathbf{r}_T = -\alpha \mathbf{u}_T, \alpha \in \mathbb{R}_+ \end{array} \right. \quad \text{and} \quad \left\{ \begin{array}{l} u_N = 0 \\ \mathbf{u}_T \neq \mathbf{0} \end{array} \right. \quad (\text{sliding}) \end{array} \right. \quad (1.1)$$

where  $\mathcal{K}_\mu$  is defined w.r.t. the usual scalar product in  $\mathbb{R}^d$ . In the following, we will denote by  $\mathcal{C}_\mu(\mathbf{n}) \subset \mathbb{R}^d \times \mathbb{R}^d$  the set of velocity–force pairs satisfying the Signorini-Coulomb condition,

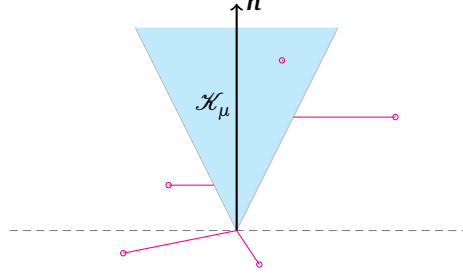
$$(\mathbf{u}, \mathbf{r}) \in \mathcal{C}_\mu(\mathbf{n}) \iff \mathbf{u} \text{ and } \mathbf{r} \text{ satisfy (1.1)} \iff \mathbf{u} \text{ and } \mathbf{r} \text{ satisfy (1) and (3)}.$$

For the sake of simplicity, we will also stop writing systematically the contact normal  $\mathbf{n}$  relative to which the normal cone and the Coulomb law solution set are defined. That is, when the precise direction of the contact normal is of no relevance, we will write  $\mathcal{K}_\mu$  instead of  $\mathcal{K}_\mu(\mathbf{n})$  and  $\mathcal{C}_\mu$  instead of  $\mathcal{C}_\mu(\mathbf{n})$ .

An analogous formulation may be derived by expressing the disjunction on  $\mathbf{r}$  instead of  $\mathbf{u}$ ,

$$(\mathbf{u}, \mathbf{r}) \in \mathcal{C}_\mu \iff \left\{ \begin{array}{ll} \text{or} & u_N \geq 0 \quad \text{and} \quad \mathbf{r} = \mathbf{0} \\ \text{or} & \mathbf{u} = \mathbf{0} \quad \text{and} \quad \mathbf{r} \in \text{int } \mathcal{K}_\mu \\ \text{or} & \left\{ \begin{array}{l} u_N = 0 \\ \mathbf{u}_T = -\alpha \mathbf{r}_T, \alpha \in \mathbb{R}_+ \end{array} \right. \quad \text{and} \quad \mathbf{r} \in \text{Bd } \mathcal{K}_\mu \setminus \{\mathbf{0}\} \end{array} \right. \quad (1.2)$$

In Section 4.1.2 we will show that for a system with a single contact and a linear relationship between  $\mathbf{u}$  and  $\mathbf{r}$ , we can find an analytical solution to the Signorini-Coulomb conditions by enumerating the cases of Equation (1.1). However, as soon as we have to deal with multiple contacts, the disjunctive formulations become cumbersome to work with. Indeed, for a system with  $n$  contact points, one would have to check for the existence of a solution in each of the  $3^n$  cases — the cost of such an enumeration would quickly become prohibitive. This



**Figure 1.3:** *The projection at the heart of the Alart–Curnier formulation*

motivates the search for alternative expressions of these conditions, i.e., ones that would lend themselves to numerical optimization methods.

### 1.1.3 Alart–Curnier function

Using the notion of normal cone (see Definition A.9), we can rewrite the Signorini and Coulomb friction conditions into a pair of root-finding problems.

Indeed, computing the normal cones of  $\mathbb{R}_+$ , the set of positive reals, and  $\mathcal{B}^d(a) \subset \mathbb{R}^d$ , the ball of radius  $\max(0, a)$  centered at  $\mathbf{0}$ , yield

$$\mathcal{N}_{\mathbb{R}_+}(x) = \begin{cases} \emptyset & \text{if } x < 0 \\ \{0\} & \text{if } x > 0 \\ \mathbb{R}_- & \text{if } x = 0, \end{cases} \quad (1.3)$$

and

$$\mathcal{N}_{\mathcal{B}^d(a)}(\mathbf{x}) = \begin{cases} \emptyset & \text{if } \|\mathbf{x}\| > a \\ \{\mathbf{0}\} & \text{if } \|\mathbf{x}\| < a \\ \{\alpha \mathbf{x}, \alpha \in \mathbb{R}_+\} & \text{if } \|\mathbf{x}\| = a \text{ and } a > 0 \\ \mathbb{R}^d & \text{if } \|\mathbf{x}\| = a \text{ and } a = 0. \end{cases} \quad (1.4)$$

Studying each case of those expressions confirms that they correspond to those of the Signorini and Coulomb friction laws. We thus get the equivalences

$$\begin{aligned} (1) &\iff \mathbf{u}_N \in -\mathcal{N}_{\mathbb{R}_+}(\mathbf{r}_N), \\ (3) &\iff \mathbf{u}_T \in -\mathcal{N}_{\mathcal{B}^{d-1}(\mu \mathbf{r}_N)}(\mathbf{r}_T). \end{aligned}$$

Finally, we can use Corollary A.6 to express the normal cone inclusions as fixed points of orthogonal projections, and obtain the functional formulation of the Coulomb contact law introduced by Alart and Curnier (1991),

$$(\mathbf{u}, \mathbf{r}) \in \mathcal{C}_\mu \iff f_{AC}(\mathbf{u}, \mathbf{r}) = \mathbf{0} \quad (1.5)$$

with, for any  $\xi \in \mathbb{R}_+^*$ ,

$$\begin{aligned} f_{AC} : \mathbb{R}^d \times \mathbb{R}^d &\rightarrow \mathbb{R}^d \\ (\mathbf{u}, \mathbf{r}) &\mapsto \begin{pmatrix} \Pi_{\mathbb{R}_+}(\mathbf{r}_N - \xi \mathbf{u}_N) \\ \Pi_{\mathcal{B}^{d-1}(\mu \mathbf{r}_N)}(\mathbf{r}_T - \xi \mathbf{u}_T) \end{pmatrix} - \mathbf{r}. \end{aligned} \quad (1.6)$$

The effect of the tangential and normal projections in Equation 1.6 is illustrated in Figure 1.3; note that their result is always in  $\mathcal{K}_\mu$ . Now, given a kinematic relationship between  $\mathbf{u}$  and  $\mathbf{r}$ , one could express the search of solutions satisfying the Coulomb law as a root-finding problem, or as a minimization problem on the norm of  $f_{AC}$ . Note however that the orthogonal projections are

not differentiable everywhere, making  $f_{AC}$  a *nonsmooth* function, and thus requiring a careful design of the minimization algorithm.

While being an effective tool for the resolution of practical problems, the Alart–Curnier function yields little insight into the mathematical structure of the Coulomb contact law. One may wonder if we could retrieve more information from the theory of plasticity, which is what we will discuss in the following section.

## 1.2 Implicit Standard Materials

### 1.2.1 Generalized Standard Materials

In its simplest form, the theory of Generalized Standard Materials (GSM) relates the evolution of a strain-like variable,  $\boldsymbol{\varepsilon} \in X$ , and a stress-like variable,  $\boldsymbol{\sigma} \in Y$ , where  $X$  and  $Y$  are reflexive Banach spaces, dual for the bilinear form  $\langle \cdot, \cdot \rangle$ . Considering an isothermal system described by  $\boldsymbol{\varepsilon}$  and another internal variable  $\boldsymbol{\gamma}$ , its evolution should satisfy (see, e.g., Saramito 2015, chapter 5)

$$\begin{cases} \boldsymbol{\sigma} \in \rho \frac{\partial \mathcal{E}}{\partial \boldsymbol{\varepsilon}}(\boldsymbol{\varepsilon}, \boldsymbol{\gamma}) + \frac{\partial \mathcal{D}}{\partial \dot{\boldsymbol{\varepsilon}}}(\dot{\boldsymbol{\varepsilon}}, \dot{\boldsymbol{\gamma}}) \\ \mathbf{0} \in \rho \frac{\partial \mathcal{E}}{\partial \boldsymbol{\gamma}}(\boldsymbol{\varepsilon}, \boldsymbol{\gamma}) + \frac{\partial \mathcal{D}}{\partial \dot{\boldsymbol{\gamma}}}(\dot{\boldsymbol{\varepsilon}}, \dot{\boldsymbol{\gamma}}) \end{cases}$$

with  $\mathcal{E}$  and  $\mathcal{D}$  proper closed convex functions,  $\rho$  is the density of the material, and  $\dot{\boldsymbol{\varepsilon}} := \frac{d\boldsymbol{\varepsilon}}{dt}$ .  $\mathcal{E}$  is called the Helmholtz free-energy function, and  $\mathcal{D}$  the dissipation potential.

Still following Saramito (2015, chapter 5), we can decompose the strain  $\boldsymbol{\varepsilon}$  as an elastic part,  $\boldsymbol{\varepsilon}_e$ , and a plastic part,  $\boldsymbol{\varepsilon}_p$ . Assuming locally small deformations, we write  $\dot{\boldsymbol{\varepsilon}} = \dot{\boldsymbol{\varepsilon}}_e + \dot{\boldsymbol{\varepsilon}}_p$ . We furthermore consider that the free energy depends only on  $\boldsymbol{\varepsilon}_e$ , and the potential of dissipation only on  $\boldsymbol{\varepsilon}_p$ . Writing again our equations for  $\boldsymbol{\gamma} := \boldsymbol{\varepsilon}_e$ , we get

$$\begin{cases} \boldsymbol{\sigma} \in \frac{\partial \mathcal{D}}{\partial \dot{\boldsymbol{\varepsilon}} - \dot{\boldsymbol{\varepsilon}}_e}(\dot{\boldsymbol{\varepsilon}} - \dot{\boldsymbol{\varepsilon}}_e) \\ \mathbf{0} \in \rho \frac{\partial \mathcal{E}}{\partial \boldsymbol{\varepsilon}_e}(\boldsymbol{\varepsilon}_e) - \frac{\partial \mathcal{D}}{\partial \dot{\boldsymbol{\varepsilon}} - \dot{\boldsymbol{\varepsilon}}_e}(\dot{\boldsymbol{\varepsilon}} - \dot{\boldsymbol{\varepsilon}}_e), \end{cases}$$

or equivalently,

$$\begin{cases} \boldsymbol{\sigma} \in \frac{\partial \mathcal{D}}{\partial \dot{\boldsymbol{\varepsilon}}_p}(\dot{\boldsymbol{\varepsilon}}_p) \end{cases} \quad (1.7)$$

$$\begin{cases} \boldsymbol{\sigma} \in \rho \frac{\partial \mathcal{E}}{\partial \boldsymbol{\varepsilon}_e}(\boldsymbol{\varepsilon}_e). \end{cases} \quad (1.8)$$

Note the purely plastic case is recovered for  $\mathcal{E} = \mathcal{J}_{\{0\}}$ , as Equation (1.8) then boils down to  $\boldsymbol{\varepsilon}_e = \mathbf{0}$ .

Using Theorem A.2, we can also write Equation (1.7) as the so-called *flow rule*,

$$\dot{\boldsymbol{\varepsilon}}_p \in \frac{\partial \mathcal{D}^*}{\partial \boldsymbol{\sigma}}(\boldsymbol{\sigma}). \quad (1.9)$$

We are interested in the particular case of yield-stress materials, where analogously to our Coulomb contact force  $\mathbf{r}$  that must lie in  $\mathcal{K}_\mu$ , the stress  $\boldsymbol{\sigma}$  must belong to an admissible set. Let  $C := \{\boldsymbol{\sigma}, F(\boldsymbol{\sigma}) \leq 0\}$  be this set of admissible stresses, with  $F(\boldsymbol{\sigma}) = 0$  the yield isosurface. The flow rule (1.9) is said to be associated if it is derived from Hill's maximum power principle (see e.g., de Saxcé and Bousshine 2002), which imposes

$$\boldsymbol{\sigma} = \arg \min_{\boldsymbol{\tau} \in C} -\langle \boldsymbol{\tau}, \dot{\boldsymbol{\varepsilon}}_p \rangle.$$

From the optimality condition of Theorem A.6, a flow rule will thus be associated if and only if  $\dot{\boldsymbol{\varepsilon}}_p \in \mathcal{N}_C(\boldsymbol{\sigma})$ , and therefore  $\mathcal{D} = \mathcal{J}_C^*$ .

**Remark** The flow rule (1.9) is often written as

$$\begin{aligned} \dot{\epsilon}_p \in \frac{\partial \mathcal{D}^*}{\partial \sigma}(\sigma) &\iff \dot{\epsilon}_p \in \alpha \frac{\partial F}{\partial \sigma}(\sigma) \text{ and } \begin{cases} F(\sigma) \leq 0 \\ \alpha = 0 & \text{if } F(\sigma) < 0 \\ \alpha \in \mathbb{R}_+ & \text{if } F(\sigma) = 0 \end{cases} \\ &\iff \dot{\epsilon}_p \in \alpha \frac{\partial F}{\partial \sigma}(\sigma) \text{ and } 0 \geq F(\sigma) \perp \alpha \geq 0 \end{aligned}$$

where  $\alpha$  is called the plastic multiplier, or consistency parameter. The equivalence – under regularity conditions for  $F$  – is granted by Property (A.16).

**Coulomb law** In our multibody contact setting,  $-\mathbf{r}$  plays the role of  $\sigma$  and  $\mathbf{u}$  that of  $\dot{\epsilon}_p$ . As a first attempt to fit Coulomb friction in the GSM framework, we can define an associated flow rule from the admissible set for the contact force,  $-\mathbf{r} \in -\mathcal{K}_\mu$ . Since  $\mathcal{K}_\mu$  is a convex cone, Theorem A.4 states that the flow rule  $\mathbf{u} \in \mathcal{N}_{-\mathcal{K}_\mu}(-\mathbf{r})$  will be equivalent to  $\mathcal{K}_\mu \ni \mathbf{r} \perp \mathbf{u} \in \mathcal{K}_{\frac{1}{\mu}}$ . This means that the power dissipated through the friction force,  $\langle \mathbf{u}, \mathbf{r} \rangle$ , will always be zero; this is physically absurd, unless  $\mu = 0$ . We therefore conclude that Coulomb friction cannot be modeled with an associated flow rule. Actually, de Saxcé and Feng (1998) showed that attempting to fit the Coulomb contact law inside the GSM framework is vain, and that Rockafellar's cyclic monotonicity criterion asserts that one cannot find a dissipation potential  $\mathcal{D}$  such that  $(\mathbf{u}, \mathbf{r}) \in \mathcal{C}_\mu \iff -\mathbf{r} \in \partial \mathcal{D}(\mathbf{u})$ . Instead, they introduce the Implicit Standard Material framework, a superset of GSM in which the dissipation potential may depend on the strain-rate-like variables.

### 1.2.2 Implicit Standard Materials

The theory of Implicit Standard Materials (ISM) was introduced by de Saxcé and Feng (1991), and extended in (Berga and de Saxcé 1994; de Saxcé 1992; de Saxcé and Feng 1998). As the standard notion of dissipation potential is not adapted to the modeling of Coulomb's law, they introduce the weaker notion of *bipotentials*, which they define as follow:

**Definition 1.2** (Bipotential). *A function  $b : X \times Y \rightarrow \bar{\mathbb{R}}$  is called a bipotential if*

1.  *$b$  is convex and closed w.r.t. each of its two parameters, and*
2.  *$b$  satisfies a Fenchel-Young-like inequality, that is*

$$b(\dot{\epsilon}, \sigma) \geq \langle \dot{\epsilon}, \sigma \rangle \quad \forall \dot{\epsilon}, \sigma \in X \times Y. \quad (1.10)$$

This definition allows for a weaker, unilateral version of Theorem A.2,

**Property 1.2.** *If  $b$  is a bipotential, then*

$$b(\dot{\epsilon}, \sigma) = \langle \dot{\epsilon}, \sigma \rangle \implies \dot{\epsilon} \in \frac{\partial b}{\partial \sigma}(\dot{\epsilon}, \sigma) \text{ and } \sigma \in \frac{\partial b}{\partial \dot{\epsilon}}(\dot{\epsilon}, \sigma).$$

*Proof.* Combining the left-hand-side with the inequality (1.10) yields

$$\begin{aligned} b(\dot{\epsilon}, \tau) - b(\dot{\epsilon}, \sigma) &\geq \langle \dot{\epsilon}, \tau \rangle - \langle \dot{\epsilon}, \sigma \rangle & \forall \tau \in Y, \text{ i.e.,} \\ b(\dot{\epsilon}, \tau) &\geq b(\dot{\epsilon}, \sigma) + \langle \dot{\epsilon}, \tau - \sigma \rangle & \forall \tau \in Y, \end{aligned}$$

which by definition of the subdifferential means  $\dot{\epsilon} \in \frac{\partial b}{\partial \sigma}(\dot{\epsilon}, \sigma)$ . The second inclusion,  $\sigma \in \frac{\partial b}{\partial \dot{\epsilon}}(\dot{\epsilon}, \sigma)$ , can be derived in a similar manner.  $\square$

This motivates the Implicit Standard Material approach, which replaces Equations (1.7) and (1.9) from GSM with

$$\boldsymbol{\sigma} \in \frac{\partial b}{\partial \dot{\boldsymbol{\varepsilon}}_p}(\dot{\boldsymbol{\varepsilon}}_p, \boldsymbol{\sigma}) \quad (1.11)$$

$$\dot{\boldsymbol{\varepsilon}}_p \in \frac{\partial b}{\partial \boldsymbol{\sigma}}(\dot{\boldsymbol{\varepsilon}}_p, \boldsymbol{\sigma}). \quad (1.12)$$

One can readily see that for  $b(\dot{\boldsymbol{\varepsilon}}, \boldsymbol{\sigma}) = \mathcal{D}(\dot{\boldsymbol{\varepsilon}}) + \mathcal{D}^*(\boldsymbol{\sigma})$ , we fall back to the GSM case and the reciprocal of Property 1.2 is also true. In the general case, the reverse implication will have to be checked manually to ensure that Equations (1.11) and (1.12) are equivalent.

**Coulomb law** In the framework of ISM, de Saxcé and Feng (1991) constructed a bipotential  $b(\mathbf{u}, -\mathbf{r})$  such that Equations (1.11) and (1.12) are equivalent to the Coulomb contact law (1.1).

In order to heuristically construct such a bipotential, we can make the following considerations:

- The force  $\mathbf{r}$  should lie in  $\mathcal{K}_\mu$
- The relative velocity  $\mathbf{u}$  should lie in  $\mathcal{K}_\infty = \mathbb{R}_+ \times \mathbb{R}^{d-1}$
- The energy dissipation should be consistent with Coulomb's law, meaning

$$b(\mathbf{u}, -\mathbf{r}) = -\langle \mathbf{u}, \mathbf{r} \rangle = -\langle \mathbf{u}_T, \mathbf{n}_T \rangle = \mu \mathbf{r}_N \|\mathbf{u}_T\| \quad (1.13)$$

The simplest function satisfying those rules is indeed the one proposed by de Saxcé and Feng (1991),

$$b(\mathbf{u}, -\mathbf{r}) = \mathcal{J}_{-\mathcal{K}_\mu}(-\mathbf{r}) + \mathcal{J}_{\mathcal{K}_\infty}(\mathbf{u}) + \mu \|\mathbf{u}_T\| \mathbf{r}_N. \quad (1.14)$$

**Property 1.3.**  $b$  defined as in (1.14) is a bipotential.

*Proof.* The closed, convex nature of  $b$  w.r.t. each variable is immediate as  $b$  is the sum of such functions. Now, for  $\mathbf{u} \in \mathcal{K}_\infty$  and  $\mathbf{r} \in \mathcal{K}_\mu$ , we have to show that

$$\begin{aligned} b(\mathbf{u}, -\mathbf{r}) &\geq -\langle \mathbf{u}, \mathbf{r} \rangle & \text{i.e.,} \\ \mu \|\mathbf{u}_T\| \mathbf{r}_N &\geq -\mathbf{u}_N \mathbf{r}_N - \langle \mathbf{u}_T, \mathbf{r}_T \rangle. \end{aligned}$$

Using the Cauchy-Schwartz inequality and the fact that  $\mathbf{r} \in \mathcal{K}_\mu$ ,

$$-\langle \mathbf{u}_T, \mathbf{r}_T \rangle \leq \|\mathbf{u}_T\| \|\mathbf{r}_T\| \leq \mu \mathbf{r}_N \|\mathbf{u}_T\|.$$

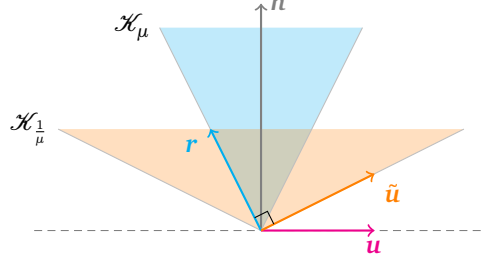
Finally,  $\mathbf{u}_N \geq 0$  and  $\mathbf{r}_N \geq 0$ , yielding  $0 \geq -\mathbf{u}_N \mathbf{r}_N$  and thus concluding the proof.  $\square$

**De Saxcé change of variable** Now, let us show that with  $b$  defined as per (1.14), Equations (1.11) and (1.12) are equivalent to  $(\mathbf{u}, \mathbf{r}) \in \mathcal{C}_\mu$ . Still following de Saxcé and Feng, we first state an intermediate result; let us introduce the auxiliary variable  $\tilde{\mathbf{u}}$ , deduced from  $\mathbf{u}$  as  $\tilde{\mathbf{u}} := \mathbf{u} + \mu \|\mathbf{u}_T\| \mathbf{n}$ .

This change of variable  $\mathbf{u} \mapsto \tilde{\mathbf{u}}$ , illustrated in Figure 1.4 and which we will refer to as the *de Saxcé change of variable*, maps the half-space  $\mathcal{K}_\infty$  to the Second-Order Cone  $\mathcal{K}_\mu^1$ , with

$$\begin{cases} \mathbf{u}_N \geq 0 & \iff \mathbf{u} + \mu \|\mathbf{u}_T\| \mathbf{n} \in \mathcal{K}_\mu^1 & \text{and} \\ \mathbf{u}_N = 0 & \iff \mathbf{u} + \mu \|\mathbf{u}_T\| \mathbf{n} \in \text{Bd } \mathcal{K}_\mu^1. \end{cases}$$

We get the following equivalences, which will prove very useful for the conception of numerical algorithms.



**Figure 1.4:** De Saxcé change of variable  $u \mapsto \tilde{u}$  in the sliding case of Coulomb friction. More generally, this change of variable maps the half-space  $\mathcal{K}_\infty$  to the dual cone  $\mathcal{K}_\mu^\perp$

**Property 1.4.** The Coulomb contact law (1.1) can be expressed in a compact manner on the variables  $r$  and  $\tilde{u}$  as

$$(u, r) \in \mathcal{C}_\mu \iff \tilde{u} \in -\mathcal{N}_{\mathcal{K}_\mu} r \iff r \in -\mathcal{N}_{\mathcal{K}_\mu^\perp} \tilde{u} \iff \mathcal{K}_\mu^\perp \ni \tilde{u} \perp r \in \mathcal{K}_\mu,$$

where the  $\perp$  notation is referring to the usual  $\mathbb{R}^d$  orthogonality.

*Proof.* Let us prove the first equivalence. Property A.14 gives us the expression of  $\mathcal{N}_{\mathcal{K}_\mu}$ ,

$$\begin{aligned} \tilde{u} \in -\mathcal{N}_{\mathcal{K}_\mu}(r) &\iff \begin{cases} r \in \mathcal{K}_\mu \\ \tilde{u} = 0 & \text{if } r \in \text{int } \mathcal{K}_\mu \\ \tilde{u} \in \mathcal{K}_\mu^\perp & \text{if } r = 0 \\ \tilde{u} \in \text{Bd } \mathcal{K}_\mu^\perp \cap \{r\}^\perp & \text{if } r \in \text{Bd } \mathcal{K}_\mu \setminus \{0\} \end{cases} \\ &\iff \begin{cases} r \in \mathcal{K}_\mu \\ u = 0 & \text{if } r \in \text{int } \mathcal{K}_\mu \\ u \in \mathcal{K}_\infty & \text{if } r = 0 \\ u \in \text{Bd } \mathcal{K}_\infty \text{ and } u_T = -\alpha r_T, \alpha \in \mathbb{R}_+ & \text{if } r \in \text{Bd } \mathcal{K}_\mu \setminus \{0\} \end{cases} \end{aligned}$$

We recognize the disjunctive formulation of the Coulomb law on  $r$ , Equation (1.2), and conclude that  $\tilde{u} \in -\mathcal{N}_{\mathcal{K}_\mu}(r) \iff (u, r) \in \mathcal{C}_\mu$ . Then, the other equivalences are direct applications of Theorem A.4, noting that  $-\mathcal{K}_\mu^\circ = \mathcal{K}_\mu^\perp$ .  $\square$

We therefore want our candidate bipotential to satisfy  $\tilde{u} \in \partial \mathcal{G}_{-\mathcal{K}_\mu}(-r) \iff u \in -\frac{\partial b}{\partial r}(u, -r)$ . This will be the case if we choose  $b$  as  $b(u, -r) = \mathcal{G}_{-\mathcal{K}_\mu}(-r) + \mu \|u_T\| r_N + f(u)$ .

**Property 1.5** (Bipotential for Coulomb law). The bipotential  $b$  defined as in Equation (1.14) models the Coulomb contact law, i.e.,

$$b(u, -r) = -\langle u, r \rangle \iff r = -\frac{\partial b}{\partial u}(u, -r) \iff u = -\frac{\partial b}{\partial r}(u, -r) \iff (u, r) \in \mathcal{C}_\mu$$

*Proof.* From Properties 1.2 and 1.3, we know that

$$b(u, -r) = -\langle u, r \rangle \implies r = -\frac{\partial b}{\partial u}(u, -r) \text{ and } u = -\frac{\partial b}{\partial r}(u, -r).$$

As the regularity conditions of Corollary A.1 on the subdifferential of a sum are satisfied, we get

$$\begin{aligned} u \in -\frac{\partial b}{\partial r}(u, -r) &\iff u \in \left( \partial \mathcal{G}_{-\mathcal{K}_\mu}(-r) - \{\mu \|u_T\| n\} \right) \\ &\iff \tilde{u} \in -\mathcal{N}_{\mathcal{K}_\mu} && \text{using Property A.9} \\ &\iff (u, r) \in \mathcal{C}_\mu. && \text{using Property 1.4} \end{aligned}$$

Now, suppose  $\mathbf{r} \in -\frac{\partial b}{\partial \mathbf{u}}(\mathbf{u}, -\mathbf{r})$ , and let us show  $b(\mathbf{u}, -\mathbf{r}) = -\langle \mathbf{u}, \mathbf{r} \rangle$ . Using Theorem A.1 on the subdifferential of a sum and Property A.14 on the normal cone to a convex cone, we have  $\mathbf{r} \in \mathcal{K}_0 \cap \{\mathbf{u}\}^\perp - \mu \mathbf{r}_N \partial \|\mathbf{u}\|$ . We always get  $\mathbf{r} \in \mathcal{K}_\mu$ , and can distinguish three cases:

1.  $\mathbf{u} = \mathbf{0}$ , then  $b(\mathbf{u}, -\mathbf{r}) = 0 = -\langle \mathbf{u}, \mathbf{r} \rangle$ .
2.  $\mathbf{u}_N \neq 0$  and  $\mathbf{u}_T = \mathbf{0}$ , then once again  $-\langle \mathbf{u}, \mathbf{r} \rangle = -\mathbf{u}_N \mathbf{r}_N = 0 = b(\mathbf{u}, -\mathbf{r})$ .
3.  $\mathbf{u}_N \neq 0$  and  $\mathbf{u}_T \neq \mathbf{0}$ , then  $\mathbf{r} + \mu \mathbf{r}_N \frac{\mathbf{u}_T}{\|\mathbf{u}_T\|} \in \{\mathbf{u}\}^\perp$ , and therefore  $-\langle \mathbf{u}, \mathbf{r} \rangle = \mu \mathbf{r}_N \|\mathbf{u}_T\| = b(\mathbf{u}, -\mathbf{r})$ .

To conclude the proof, remark that the implication  $(\mathbf{u}, \mathbf{r}) \in \mathcal{C}_\mu \implies b(\mathbf{u}, -\mathbf{r}) = -\langle \mathbf{u}, \mathbf{r} \rangle$  was ensured by our third criterion (1.13) in the heuristic construction of  $b$ .  $\square$

**Remark** We have seen that the Coulomb law is more complex than associated flow rules, in the sense that its bipotential includes a term that couples the strain and stress variables,  $\mathbf{r}_N \|\mathbf{u}_T\|$ ; we can well imagine that this will lead to a tougher numerical problem. However, one can get back to the “easier” case of an associated law by temporarily freezing either one of the variables. Indeed we will see in Section 2.3.2 that freezing  $\mathbf{r}_N$  leads to the Haslinger (1983) algorithm (successive approximations with the Tresca law, yielding minimization problems over a half-cylinder), while freezing  $\|\mathbf{u}_T\|$  leads to the Cadoux (2009) algorithm (successive minimization problems over a Second-Order Cone).

### 1.3 Application to Drucker–Prager plasticity

The Drucker–Prager yield surface (7) is structurally very similar to the Coulomb friction law. Indeed, Moreau (1965) cites this application as a motivating factor for convex analysis, together with the study of unilateral contacts with friction. More recently, Berga and de Saxcé (1994) showed that the bipotential (1.14) could be easily adapted to model non-associated Drucker–Prager plasticity; we will keep their point of view below. Alternatively, we provided a less general but more direct derivation of the equivalence between Coulomb law and a particular case of the Drucker–Prager flow rule in (Daviet and Bertails-Descoubes 2016b, Section 3.2).

In order to derive such results, we first present a few notations to reason about the space to which the stress  $\boldsymbol{\sigma}$  and displacements  $\boldsymbol{\varepsilon}$  belong, i.e., the space of  $d \times d$  symmetric tensors.

#### 1.3.1 Symmetric tensors

Let us consider the space of  $d \times d$  symmetric tensors  $S_d$ ,

$$S_d = \{(\mathbf{x} \otimes \mathbf{y} + \mathbf{y} \otimes \mathbf{x}), \mathbf{x}, \mathbf{y} \in \mathbb{R}^{2d}\} \quad (1.15)$$

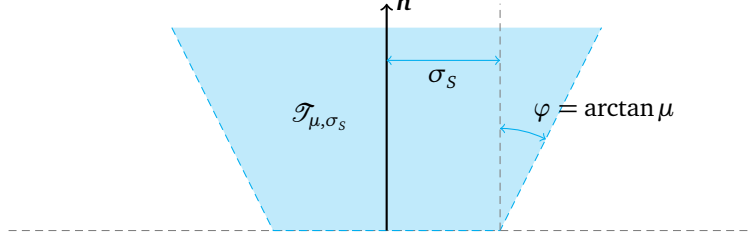
equipped with the scalar product  $\langle \cdot, \cdot \rangle$  defined from the twice-contracted tensor product,

$$\langle \boldsymbol{\sigma}, \boldsymbol{\tau} \rangle := \frac{1}{2} \boldsymbol{\sigma} : \boldsymbol{\tau} := \frac{1}{2} \sum_{i,j} \tau_{ij} \sigma_{ij} \quad \forall \boldsymbol{\sigma}, \boldsymbol{\tau} \in S_d^2.$$

We note  $|\cdot| := \sqrt{\langle \cdot, \cdot \rangle}$  the norm associated to this scalar product. Finally, we introduce the notation  $\text{Tr}$  for the linear form associating to each tensor its the trace,  $\text{Tr} : S_d \rightarrow \mathbb{R}$ ,  $\boldsymbol{\tau} \mapsto \sum_i \tau_{ii}$ , and  $\text{Dev}$  for the linear application yielding its deviatoric part,  $\text{Dev} : S_d \rightarrow S_d$ ,  $\boldsymbol{\tau} \mapsto \boldsymbol{\tau} - \frac{1}{d} \text{Tr } \boldsymbol{\tau}$ .

Using these notions, we can define a Second-Order Cone on  $S_d$ . First, notice that the element of  $S_d$  representing the linear form  $\text{Tr}$  through the scalar product  $\langle \cdot, \cdot \rangle$  is  $2\mathbb{I}_d$ , i.e.,  $\text{Tr } \boldsymbol{\tau} = \langle \boldsymbol{\tau}, 2\mathbb{I}_d \rangle$ . We note  $\boldsymbol{\iota}_d$  the corresponding unit tensor,

$$\boldsymbol{\iota}_d = \frac{2\mathbb{I}_d}{|2\mathbb{I}_d|} = \sqrt{\frac{2}{d}} \mathbb{I}_d.$$



**Figure 1.5:** The 2D truncated Second-Order Cone  $\mathcal{T}_{\mu, \sigma_S}$

Then, it can be readily seen that  $\text{Tr}(\text{Dev } \boldsymbol{\tau}) = 0$ , and that any symmetric tensor  $\boldsymbol{\tau}$  can be decomposed as the orthogonal sum  $\boldsymbol{\tau} = \langle \boldsymbol{\tau}, \boldsymbol{\iota}_d \rangle \boldsymbol{\iota}_d + \text{Dev } \boldsymbol{\tau}$ . We can therefore extend the notation of *normal part* and *tangential part* to the space of symmetric tensors with, for  $\boldsymbol{\tau} \in S_d$ ,

$$\boldsymbol{\tau}_N = \langle \boldsymbol{\tau}, \boldsymbol{\iota}_d \rangle \quad \boldsymbol{\tau}_T = \text{Dev } \boldsymbol{\tau}.$$

In the following, we will reuse the notation  $\mathcal{K}_\mu$  for the Second-Order Cone on the space of symmetric tensors defined w.r.t.  $\langle \cdot, \cdot \rangle$  and the linear form  $\text{Tr}$ ,

$$\mathcal{K}_\mu(\boldsymbol{\iota}_d) := \{ \boldsymbol{\sigma} \in S_d, |\boldsymbol{\sigma}_T| \leq \mu \boldsymbol{\sigma}_N \} = \left\{ \boldsymbol{\sigma} \in S_d, |\text{Dev } \boldsymbol{\sigma}| \leq \frac{\mu}{\sqrt{2d}} \text{Tr } \boldsymbol{\sigma} \right\}$$

### 1.3.2 Drucker–Prager yield surface

In order to study the Drucker–Prager yield surface in this framework, we have to express the invariants  $I_1(\boldsymbol{\sigma})$  and  $J_2(\boldsymbol{\sigma})$  as functions of the normal and tangential parts of the stress tensor  $\boldsymbol{\sigma}$ . We have directly  $I_1(\boldsymbol{\sigma}) = \text{Tr } \boldsymbol{\sigma} = \sqrt{2d} \boldsymbol{\sigma}_N$ , and

$$\begin{aligned} J_2(\boldsymbol{\sigma}) &= \frac{1}{2} \text{Tr}[(\text{Dev } \boldsymbol{\sigma})^2] = \frac{1}{2} \sum_i \left( \sum_j [(\text{Dev } \boldsymbol{\sigma})_{ij}]^2 \right) = |\text{Dev } \boldsymbol{\sigma}|^2 \\ &= |\boldsymbol{\sigma}_T|^2 \end{aligned}$$

The Drucker–Prager yield surface with tension cut-off is defined by Equation (9),

$$\begin{aligned} F_{\hat{\mu}, \sigma_S, \widehat{\tau}_c}^{\text{DP}}(\boldsymbol{\sigma}) &= \max \left( I_1(\boldsymbol{\sigma}) - d \widehat{\tau}_c, \hat{\mu} \frac{I_1(\boldsymbol{\sigma}) - d \widehat{\tau}_c}{d} - \sigma_S + \sqrt{J_2(\boldsymbol{\sigma})} \right) \\ &= \max \left( \sqrt{2d} (\boldsymbol{\sigma}_N - \tau_c), \mu (\boldsymbol{\sigma}_N - \tau_c) - \sigma_S + |\boldsymbol{\sigma}_T| \right) \end{aligned}$$

using the notation  $\mu := \sqrt{\frac{2}{d}} \hat{\mu}$  and  $\tau_c = \sqrt{\frac{d}{2}} \widehat{\tau}_c$ . We have thus the equivalence

$$F_{\hat{\mu}, \sigma_S, \widehat{\tau}_c}^{\text{DP}}(\boldsymbol{\sigma}) \leq 0 \iff \begin{cases} \tau_c - \boldsymbol{\sigma}_N \geq 0 \\ \mu (\tau_c - \boldsymbol{\sigma}_N) + \sigma_S \geq \|\boldsymbol{\sigma}_T\|. \end{cases}$$

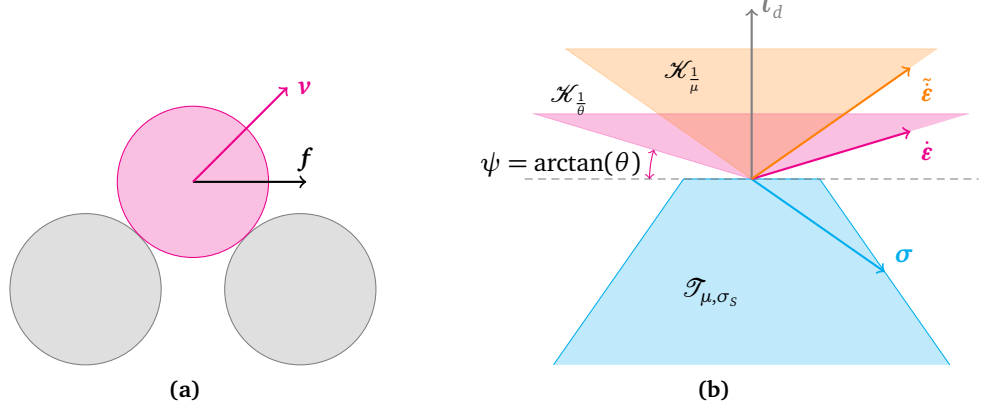
In the case  $\sigma_S = 0$ , we retrieve a set of admissible stresses similar in structure to that of the Coulomb friction, a translation of the SOC  $\mathcal{K}_\mu$ :

$$F_{\hat{\mu}, \hat{\epsilon}}^{\text{DP}}(\boldsymbol{\sigma}) \leq 0 \iff (\tau_c \boldsymbol{\iota}_d - \boldsymbol{\sigma}) \in \mathcal{K}_\mu.$$

In the general case, for  $\sigma_S \geq 0$ , the set of admissible stresses is instead a translation of a *truncated* Second-Order Cone  $\mathcal{T}_{\mu, \sigma_S}$ ,

$$\mathcal{T}_{\mu, \sigma_S} := \{ \boldsymbol{\sigma} \in S_d, \boldsymbol{\sigma}_N \geq 0 \text{ and } |\boldsymbol{\sigma}_T| \leq \tau_c + \mu \boldsymbol{\sigma}_N \}, \quad (1.16)$$





**Figure 1.6:** (a) Rationale for dilatancy: the application of a shearing force  $\mathbf{f}$  will make the “trapped” pink particle move upwards as well as horizontally. (b) De Saxcé change of variable in the flowing regime of the Drucker–Prager flow rule with dilatancy angle  $\psi$

depicted in Figure 1.5. We have the equivalence

$$F_{\hat{\mu}, \sigma_S, \hat{\tau}_c}^{\text{DP}}(\boldsymbol{\sigma}) \leq 0 \iff (\tau_c \mathbf{l}_d - \boldsymbol{\sigma}) \in \mathcal{T}_{\mu, \sigma_S}. \quad (1.17)$$

Notice that another expression for  $\mathcal{T}_{\mu, \sigma_S}$  is  $\mathcal{T}_{\mu, \sigma_S} = \mathcal{K}_{\mu} + \mathcal{B}_T(\sigma_S)$ , where  $\mathcal{B}_T(\sigma_S) \subset S_d$  is the closed ball of traceless symmetric tensors with norm lower than  $\sigma_S$ .

### 1.3.3 Dilatancy and non-associated Drucker–Prager flow rule

Drucker and Prager (1952) deduce from the theory of GSM that, using their yield surface as a dissipation potential, “volume expansion is seen to be a necessary accompaniment to shearing deformation”. As mentioned in our introduction, such a phenomenon is known as dilatancy.

Indeed, we can easily verify that the associated Drucker–Prager flow rule predicts a positive dilatancy. Defining a plasticity potential from the Drucker–Prager yield surface,  $\mathcal{D}^*(\boldsymbol{\sigma}) = \mathcal{D}_{-\mathcal{T}_{\mu, \sigma_S}}(\boldsymbol{\sigma} - \tau_c \mathbf{l}_d)$ , one gets that

$$\dot{\boldsymbol{\epsilon}} \in \frac{\partial \mathcal{D}^*}{\partial \boldsymbol{\sigma}} \iff \dot{\boldsymbol{\epsilon}} \in \mathcal{N}_{-\mathcal{T}_{\mu, \sigma_S}}(\boldsymbol{\sigma} - \tau_c \mathbf{l}_d).$$

Now if  $\sigma_S = 0$ , the normal cone to  $\mathcal{T}_{\mu, \sigma_S}$  obviously coincide with that of  $\mathcal{K}_{\mu}$ , and for  $\sigma_S > 0$ ,

$$\begin{aligned} \dot{\boldsymbol{\epsilon}} \in \mathcal{N}_{-\mathcal{T}_{\mu, \sigma_S}}(\boldsymbol{\sigma}) \\ \iff \begin{cases} \boldsymbol{\sigma} \in -\mathcal{T}_{\mu, \sigma_S} & \text{if } \boldsymbol{\sigma} \in \text{int } -\mathcal{T}_{\mu, \sigma_S} \\ \dot{\boldsymbol{\epsilon}} = \mathbf{0} & \text{if } \sigma_N = 0 \text{ and } |\boldsymbol{\sigma}_T| < \sigma_S \\ \dot{\boldsymbol{\epsilon}}_N > 0 \text{ and } \dot{\boldsymbol{\epsilon}}_T = \mathbf{0} & \text{if } \sigma_N > 0 \text{ and } |\boldsymbol{\sigma}_T| < \sigma_S \\ \dot{\boldsymbol{\epsilon}} \in \mathcal{K}_{\frac{1}{\mu}} \cap \{\boldsymbol{\sigma} - \sigma_S \frac{\boldsymbol{\sigma}_T}{|\boldsymbol{\sigma}_T|}\}^\perp & \text{if } \sigma_N > 0 \text{ and } |\boldsymbol{\sigma}_T| = \sigma_S - \mu \sigma_N \\ \dot{\boldsymbol{\epsilon}} \in \mathcal{K}_{\frac{1}{\mu}} \cap \{\alpha \mathbf{l}_d + \beta \frac{\boldsymbol{\sigma}_T}{|\boldsymbol{\sigma}_T|}, \alpha \in \mathbb{R}, \beta \in \mathbb{R}_+\} & \text{if } \sigma_N = 0 \text{ and } |\boldsymbol{\sigma}_T| = \sigma_S. \end{cases} \end{aligned} \quad (1.18)$$

In any case, we see that

$$\dot{\boldsymbol{\epsilon}} \in \mathcal{N}_{-\mathcal{T}_{\mu, \sigma_S}}(\boldsymbol{\sigma} - \tau_c \mathbf{l}_d) \implies \dot{\boldsymbol{\epsilon}} \in \mathcal{K}_{\frac{1}{\mu}}.$$

The associated flow rule therefore implies that the norm of the shear rate is limited by the normal part of the strain rate:  $\zeta |\boldsymbol{\sigma}_T| \leq \sigma_N$  with the dilatancy coefficient  $\zeta$  equal to the friction coefficient  $\mu$ . This translates to the dilatancy angle  $\psi$  being equal to the friction angle  $\varphi$ .

Dilatancy makes physical sense; indeed, a grain initially stuck between others, as in Figure 1.6(a), will first have to go upwards before being able to perform any lateral motion. However, we see from the above computation that a lower friction coefficient will mean a lower dilatancy, which is not necessarily intuitive. Actually, experiments on granular materials disproved this postulate, and the dilatancy was observed to be significantly lower than predicted by the associated flow rule (see e.g., Vermeer 1998). Moreover, the same argument that we used to reject an associated flow rule for the Coulomb friction law can be made again: for a cohesionless material,  $\sigma_S = \tau_c = 0$ , the rate of energy dissipated through the plastic strain rate,  $\dot{\epsilon} : \sigma$ , would be always zero.

Overall, this motivates keeping the dilatancy angle  $\psi$ , or equivalently our dilatancy coefficient  $\theta$ , as a independent parameter of the model. Noting that  $\dot{\epsilon} \in \mathcal{K}_\zeta^\perp \iff \dot{\epsilon} + (\mu - \zeta)|\dot{\epsilon}_T|\mathbf{l}_d \in \mathcal{K}_\mu^\perp$ , we will follow Berga and de Saxcé (1994) and use the non-associated flow rule defined by

$$\dot{\epsilon} + (\mu - \zeta)|\dot{\epsilon}_T|\mathbf{l}_d \in \mathcal{N}_{\mathcal{G}_{\mu, \sigma_S}}(\sigma - \tau_c \mathbf{l}_d). \quad (1.19)$$

For the material to be stable, plastic displacements need to dissipate energy, which means  $\dot{\epsilon} : \sigma \geq 0$  for all admissible stresses  $\sigma$ . This implies  $0 \leq \zeta \leq \mu$ , or in terms of angles,  $0 \leq \psi \leq \phi$ . Note that for  $\zeta = \mu$ , we retrieve the associated case; however, for  $\zeta < \mu$  we leave the GSM framework and obtain an Implicit Standard Material. For  $\zeta = 0$ , and a cohesionless material, we retrieve the non-associated Coulomb friction flow rule. Finally, for  $\zeta = \mu = 0$ ,  $\sigma_S > 0$  and  $\tau_c = +\infty$ , we obtain the incompressible Bingham yield surface (an infinite cylinder around the hydrostatic axis) with associated flow rule; for  $\tau_c = 0$ , we get unilateral incompressibility, that is, plastic expansion is allowed but not compression. More generally, we will see in the following section that the flow rule (1.19) enforces a maximum dissipation principle on the deviatoric component of the stress.

In order to lighten notations, we will denote by  $\mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta)$  the set of  $(\dot{\epsilon}, -\sigma) \in S_d \times S_d$  such that the flow rule (1.19) is satisfied. We put the “minus” in front of the stress variable in order to be consistent with our definition of the Coulomb friction solution set  $\mathcal{C}_\mu$ . Indeed we recognize from Property 1.4 that  $\mathcal{C}_\mu \sim \mathcal{DP}(\mu, 0, 0, 0)$ , the only difference being that the former is defined on vector spaces, while the later is defined on tensor spaces. As such, the flow rule (1.19) naturally motivates the introduction of another version of the *de Saxcé change of variable* presented in Property 1.4, and illustrated in Figure 1.6(b).

**Property 1.6.** *With the change of variable  $\dot{\epsilon} \mapsto \tilde{\epsilon} := \dot{\epsilon} + (\mu - \zeta)|\dot{\epsilon}_T|\mathbf{l}_d$ , the Drucker–Prager flow rule may be written as*

$$(\dot{\epsilon}, -\sigma) \in \mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta) \iff \tilde{\epsilon} \in \mathcal{N}_{\mathcal{G}_{\mu, \sigma_S}}(\sigma - \tau_c \mathbf{l}_d).$$

When there are no ambiguities, we will omit the parameters which are taken to be zero, e.g.,  $\mathcal{DP}(\mu, \tau_c) = \mathcal{DP}(\mu, 0, \tau_c, 0)$ . Indeed, a non-zero  $\tau_c$  do not fundamentally change the nature of the solution set, so we will often discard it.

### 1.3.4 Bipotential and reformulations of the Drucker–Prager flow rule

Using the bipotential framework, we will show that  $\mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta)$  is fully characterized by the disjunctive formulation (1.20),

$$\begin{cases} \sigma_T = (\sigma_S + \mu(\tau_c - \sigma_N)) \frac{\dot{\epsilon}_T}{|\dot{\epsilon}_T|} & \text{if } \dot{\epsilon}_T \neq 0 \\ |\sigma_T| \leq (\sigma_S + \mu(\tau_c - \sigma_N)) & \text{if } \dot{\epsilon}_T = 0 \\ 0 \leq \tau_c - \sigma_N \perp \dot{\epsilon}_N - \zeta|\dot{\epsilon}_T| \geq 0. \end{cases} \quad (1.20)$$

We mentioned in the previous paragraph that for  $\sigma_S = \tau_c = \zeta = 0$ , we retrieved a flow rule similar to that of the Coulomb contact law; that is,  $\mathcal{C}_\mu \sim \mathcal{DP}(\mu)$ . This motivates looking for a bipotential with a structure similar to Equation (1.14), but adapted to our new version of the *de Saxcé change of variable* from Property 1.6.

**Property 1.7** (Bipotential for the Drucker–Prager flow rule). *The function  $b : S_d \times S_d \rightarrow \mathbb{R}$ ,*

$$b(\dot{\epsilon}, \sigma) = \mathcal{J}_{-\mathcal{T}_{\mu, \sigma_S}}(\sigma - \tau_c \iota_d) + \mathcal{J}_{\mathcal{K}_{\frac{1}{\theta}}}(\dot{\epsilon}) + \tau_c \dot{\epsilon}_N + \sigma_S |\dot{\epsilon}_T| + (\mu - \zeta)(\tau_c - \sigma_N) |\dot{\epsilon}_T|$$

*is a bipotential modeling the non-associated Drucker–Prager flow rule (1.19). That is,*

$$\begin{aligned} (\dot{\epsilon}, -\sigma) \in \mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta) &\iff b(\dot{\epsilon}, \sigma) = \langle \dot{\epsilon}, \sigma \rangle \\ &\iff \dot{\epsilon} \in \frac{\partial b}{\partial \sigma}(\dot{\epsilon}, \sigma) \iff \sigma \in \frac{\partial b}{\partial \dot{\epsilon}}(\dot{\epsilon}, \sigma) \\ &\iff \text{the disjunctive formulation (1.20) is satisfied.} \end{aligned}$$

*Proof.* It follows from the definition of  $\mathcal{DP}$  and the Drucker–Prager flow rule (1.19) that  $(\dot{\epsilon}, \sigma) \in \mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta) \iff \dot{\epsilon} \in \frac{\partial b}{\partial \sigma}(\dot{\epsilon}, \sigma)$ .

Let us now show that  $b$  is a bipotential.  $b$  is the sum of convex functions w.r.t.  $\sigma$  and  $\dot{\epsilon}$ , and is therefore convex w.r.t. each variable; indeed, we supposed that  $\mu \geq \zeta$ , and  $\sigma_N \leq \tau_c$  everywhere on the effective domain of  $b$ . Moreover, for any  $\sigma$  and  $\dot{\epsilon}$  in the effective domain of  $b$ , we have  $\dot{\epsilon}_N \geq \zeta |\dot{\epsilon}_T|$ , and  $\sigma_T \leq \mu(\tau_c - \sigma_N) + \sigma_S$ , therefore

$$\begin{aligned} \langle \dot{\epsilon}, \sigma \rangle &= \dot{\epsilon}_N \sigma_N + \langle \dot{\epsilon}_T, \sigma_T \rangle \leq \dot{\epsilon}_N \tau_c + \dot{\epsilon}_N (\sigma_N - \tau_c) + |\dot{\epsilon}_T| |\sigma_T| \\ &\leq \dot{\epsilon}_N \tau_c + \zeta |\dot{\epsilon}_T| (\sigma_N - \tau_c) + |\dot{\epsilon}_T| (\mu(\tau_c - \sigma_N) + \sigma_S) \\ &\leq \tau_c \dot{\epsilon}_N + \sigma_S |\dot{\epsilon}_T| + (\mu - \zeta)(\tau_c - \sigma_N) |\dot{\epsilon}_T| \\ &\leq b(\dot{\epsilon}, \sigma). \end{aligned}$$

Therefore,  $b$  is a bipotential.

From Property 1.2, this means that we only have to show that

$$\left\{ \begin{array}{ll} (\dot{\epsilon}, \sigma) \in \mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta) \implies b(\dot{\epsilon}, \sigma) = \langle \dot{\epsilon}, \sigma \rangle & (1.21) \\ (1.20) \implies b(\dot{\epsilon}, \sigma) = \langle \dot{\epsilon}, \sigma \rangle & (1.22) \\ (1.20) \iff \sigma \in \frac{\partial b}{\partial \dot{\epsilon}}(\dot{\epsilon}, \sigma), & (1.23) \end{array} \right.$$

to obtain the complete equivalence chain of Property 1.7. We can first check in each case of the expression of the normal cone of  $\mathcal{T}_{\mu, \sigma_S}$  given in Equation (1.18) that the implication (1.21) holds. Let  $\tilde{\epsilon}$  denote the *de Saxcé* change of variable from Property 1.6 and  $\tilde{\sigma} := \sigma - \tau_c \iota_d$ . The flow rule (1.19) means  $\tilde{\epsilon} \in \mathcal{N}_{-\mathcal{T}_{\mu, \sigma_S}}(\tilde{\sigma})$ , and then either:

1.  $\tilde{\sigma} \in \text{int-}\mathcal{T}_{\mu, \sigma_S}$ , then  $\tilde{\epsilon} = \dot{\epsilon} = \mathbf{0}$ , therefore  $b(\dot{\epsilon}, \sigma) = \langle \dot{\epsilon}, \sigma \rangle = 0$
2.  $\tilde{\sigma}_T < \sigma_S$  and  $\tilde{\sigma}_N = 0$ , then  $\dot{\epsilon}_T = \mathbf{0}$  and  $b(\dot{\epsilon}, \sigma) = \dot{\epsilon}_N \tau_c = \langle \dot{\epsilon}, \sigma \rangle$
3.  $\tilde{\sigma}_N > 0$  and  $|\sigma_T| = \sigma_S - \mu \tilde{\sigma}_N$ , then

$$\langle \tilde{\epsilon}, \tilde{\sigma} \rangle = \left\langle \tilde{\epsilon}, \tilde{\sigma} - \sigma_S \frac{\sigma_T}{|\sigma_T|} \right\rangle + \sigma_S \left\langle \tilde{\epsilon}, \frac{\sigma_T}{|\sigma_T|} \right\rangle = \sigma_S |\dot{\epsilon}_T|$$

and

$$\begin{aligned} \langle \dot{\epsilon}, \sigma \rangle &= \langle \tilde{\epsilon}, \tilde{\sigma} \rangle + \tau_c \dot{\epsilon}_N - \langle (\mu - \theta) |\dot{\epsilon}| \iota_d, \tilde{\sigma} \rangle \\ &= \sigma_S |\dot{\epsilon}_T| + \tau_c \dot{\epsilon}_N - (\mu - \theta) \tilde{\sigma}_N |\dot{\epsilon}_T| = b(\dot{\epsilon}, \sigma). \end{aligned}$$

4.  $\tilde{\sigma}_N = 0$  and  $|\sigma_T| = \sigma_S$ , then once again  $\langle \tilde{\epsilon}, \tilde{\sigma} \rangle = \sigma_S |\dot{\epsilon}_T|$ , and

$$\langle \dot{\epsilon}, \sigma \rangle = \langle \tilde{\epsilon}, \tilde{\sigma} \rangle + \tau_c \dot{\epsilon}_N - \langle (\mu - \theta) |\dot{\epsilon}| \iota_d, \tilde{\sigma} \rangle = \sigma_S |\dot{\epsilon}_T| + \tau_c \dot{\epsilon}_N = b(\dot{\epsilon}, \sigma).$$

Now, let us treat the right-hand side of the equivalence chain. First, we show the equivalence (1.23),

$$(1.20) \iff \sigma \in \frac{\partial b}{\partial \dot{\epsilon}}(\dot{\epsilon}, \sigma).$$

Indeed,

$$\begin{aligned} \sigma \in \frac{\partial b}{\partial \dot{\epsilon}}(\dot{\epsilon}, \sigma) &\iff \tilde{\sigma} \in \mathcal{N}_{\mathcal{K}_{\frac{1}{\zeta}}}(\dot{\epsilon}) + (\sigma_S - (\mu - \zeta)) \tilde{\sigma}_N \frac{\partial |\dot{\epsilon}_T|}{\partial \dot{\epsilon}} \\ &\iff \begin{cases} \dot{\epsilon} \in \mathcal{K}_{\frac{1}{\zeta}} \\ \tilde{\sigma} \in (\sigma_S - (\mu - \zeta) \tilde{\sigma}_N) \frac{\dot{\epsilon}_T}{|\dot{\epsilon}_T|} - \mathcal{K}_{\theta} \cap \{\dot{\epsilon}\}^{\perp} & \text{if } \dot{\epsilon}_T \neq 0 \\ \tilde{\sigma} \in \mathcal{B}_T((\sigma_S - (\mu - \zeta) \tilde{\sigma}_N)) - \mathcal{K}_{\theta} \cap \{\dot{\epsilon}\}^{\perp} & \text{if } \dot{\epsilon}_T = 0. \end{cases} \end{aligned} \quad (1.24)$$

Now, for any  $\tau \in S_d$ ,  $\tau - \zeta \tau_N \frac{\tau_T}{|\tau_T|} \in -\mathcal{K}_{\zeta} \iff \tau \in -\mathcal{K}_0$ , so for  $\dot{\epsilon} \in \mathcal{K}_{\frac{1}{\zeta}}$ ,

$$\tau - \zeta \tau_N \frac{\tau_T}{|\tau_T|} \in -\mathcal{K}_{\zeta} \cap \{\dot{\epsilon}\}^{\perp} \iff \tau \in -\mathcal{K}_0 \text{ and } \tau_N(\dot{\epsilon}_N - \zeta |\dot{\epsilon}_T|).$$

Moreover,  $\tau \in -\mathcal{K}_{\zeta} \cap \{\dot{\epsilon}\}^{\perp}$  implies that  $\tau_T$  is colinear with  $\dot{\epsilon}_T$  (with  $\dot{\epsilon}_T$  eventually zero). Using these results in Equation (1.24), we get

$$\sigma \in \frac{\partial b}{\partial \dot{\epsilon}}(\dot{\epsilon}, \sigma) \iff \begin{cases} \dot{\epsilon} \in \mathcal{K}_{\frac{1}{\zeta}} \\ \tilde{\sigma}_N \leq 0 \\ 0 = \tilde{\sigma}_N(\dot{\epsilon}_N - \zeta |\dot{\epsilon}_T|) \\ \sigma_T = (\sigma_S - \mu \tilde{\sigma}_N) \frac{\dot{\epsilon}_T}{|\dot{\epsilon}_T|} & \text{if } \dot{\epsilon}_T \neq 0 \\ \sigma_T \in \mathcal{B}_T((\sigma_S - \mu \tilde{\sigma}_N)) & \text{if } \dot{\epsilon}_T = 0 \end{cases} \iff (1.20).$$

It now only remains to verify (1.21), i.e., that the disjunctive formulation (1.20) implies that  $b(\dot{\epsilon}, \sigma) = \langle \dot{\epsilon}, \sigma \rangle$ . Indeed, if  $\dot{\epsilon}$  and  $\sigma$  satisfy (1.20), we have

$$\begin{aligned} \langle \dot{\epsilon}_T, \sigma_T \rangle &= (\sigma_S + \mu(\tau_c - \sigma_N)) |\dot{\epsilon}_T| \\ \dot{\epsilon}_N \sigma_N &= \tilde{\sigma}_N(\dot{\epsilon}_N - \zeta |\dot{\epsilon}_T|) + \tilde{\sigma}_N \zeta |\dot{\epsilon}_T| + \tau_c \dot{\epsilon}_N = (\sigma_N - \tau_c) \zeta |\dot{\epsilon}_T| + \tau_c \dot{\epsilon}_N, \end{aligned}$$

so  $b(\dot{\epsilon}, \sigma) = \langle \dot{\epsilon}_T, \sigma_T \rangle + \dot{\epsilon}_N \sigma_N$ . □

**Convex conjugate of the truncated SOC characteristic function** Property 1.7 implicitly gives an expression for the subdifferential of  $\mathcal{J}_{\mu, \sigma_S}^*$ . Indeed, considering the Drucker–Prager flow rule for  $\mu = \zeta$  and  $\tau_c = 0$ , we get

$$\begin{aligned} \sigma \in \partial \left( \mathcal{J}_{-\mathcal{J}_{\mu, \sigma_S}}^* \right)(\dot{\epsilon}) &\iff \dot{\epsilon} \in \mathcal{N}_{-\mathcal{J}_{\mu, \sigma_S}}(\sigma) \\ &\iff \sigma \in \frac{\partial b}{\partial \dot{\epsilon}}(\dot{\epsilon}, \sigma) \iff \sigma \in \partial \left( \mathcal{J}_{\mathcal{K}_{\frac{1}{\mu}}} + \sigma_S |\sigma_T| \right). \end{aligned}$$

As moreover  $\mathcal{J}_{-\mathcal{J}_{\mu, \sigma_S}}^*(0) = \left( \mathcal{J}_{\mathcal{K}_{\frac{1}{\mu}}} + \sigma_S |\sigma_T| \right) = 0$ , we deduce the equality (Moreau 1966–1967, Proposition 10.j)

$$\mathcal{J}_{-\mathcal{J}_{\mu, \sigma_S}}^* = \mathcal{J}_{\mathcal{K}_{\frac{1}{\mu}}} + \sigma_S |\sigma_T|. \quad (1.25)$$

**Maximum dissipation principle** The first two equations of the disjunctive formulation (1.20) may be rewritten in an equivalent manner as  $\sigma_T \in \mathcal{N}_{\mathcal{B}_T(\sigma_S + (\tau_c - \sigma_N)\mu)}(\dot{\epsilon}_T)$ , which we recognize from Theorem A.6 as the optimality conditions of the optimization problem

$$\max_{\tau \in \mathcal{B}_T(\sigma_S + (\tau_c - \sigma_N)\mu)} \langle \tau, \dot{\epsilon}_T \rangle.$$

In other words,  $\sigma_T$  maximizes the dissipated energy over the set of admissible deviatoric stresses.

**Complementarity functions** We can use the disjunctive formulation (1.20) to define an extension of the Alart–Curnier complementarity function (1.6) whose roots will coincide with the solutions of the Drucker–Prager flow rule,

$$(\dot{\mathbf{e}}, -\boldsymbol{\sigma}) \in \mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta) \iff f_{AC}(\dot{\mathbf{e}}, -\boldsymbol{\sigma}) = \mathbf{0} \quad (1.26)$$

with, for any  $\xi \in \mathbb{R}_+^*$ ,

$$\begin{aligned} f_{AC} : S_d \times S_d &\rightarrow S_d \\ (\dot{\mathbf{e}}, -\boldsymbol{\sigma}) &\mapsto \begin{pmatrix} \Pi_{\mathbb{R}_+}(\tau_c - \sigma_N - \xi(\dot{\mathbf{e}}_N - \theta|\dot{\mathbf{e}}_T|)) \\ \Pi_{\mathcal{B}_T(\sigma_S + \mu(\tau_c - \sigma_N))}(-\boldsymbol{\sigma}_T - \xi\dot{\mathbf{e}}_T) \end{pmatrix} + \boldsymbol{\sigma} - \tau_c \mathbf{l}_d. \end{aligned} \quad (1.27)$$

In both the disjunctive formulation (1.20) and the Alart–Curnier formulation (1.27), a non-zero dilatancy significantly complicates the model, yielding a nonlinear complementarity problem on the normal components of  $\dot{\mathbf{e}}$  and  $\boldsymbol{\sigma}$ . However, this coefficient has not as much influence on the structure of the bipotential. This motivates tackling problems with non-zero dilatancy with methods derived from the bipotential framework, for instance using the *de Saxcé* change of variable. As such, we can define the De Saxcé complementarity function  $f_{DS}$  satisfying

$$(\dot{\mathbf{e}}, -\boldsymbol{\sigma}) \in \mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta) \iff f_{DS}(\dot{\mathbf{e}}, -\boldsymbol{\sigma}) = \mathbf{0} \quad (1.28)$$

with, for any  $\xi \in \mathbb{R}_+^*$ ,

$$\begin{aligned} f_{DS} : S_d \times S_d &\rightarrow S_d \\ (\dot{\mathbf{e}}, -\boldsymbol{\sigma}) &\mapsto \Pi_{\mathcal{J}_{\mu, \sigma_S}}(\tau_c \mathbf{l}_d - \boldsymbol{\sigma} - \xi \tilde{\mathbf{e}}) + \boldsymbol{\sigma} - \tau_c \mathbf{l}_d \end{aligned} \quad (1.29)$$

and  $\tilde{\mathbf{e}}$  defined as the *de Saxcé* change of variable of Property 1.6,  $\tilde{\mathbf{e}} := \dot{\mathbf{e}} + (\mu - \zeta)|\dot{\mathbf{e}}_T| \mathbf{l}_d$ .

### 1.3.5 Viscoplasticity

Starting again from the GSM equations, we can model viscoplasticity by waiving the assumption that  $\mathcal{D}$  depends only on the plastic displacement  $\dot{\mathbf{e}}_p$ , and writing instead  $\mathcal{D} = \mathcal{V}(\dot{\mathbf{e}}) + \mathcal{D}_p(\dot{\mathbf{e}}_p)$ . Equations (1.7–1.8) become

$$\begin{cases} \boldsymbol{\sigma} \in \frac{\partial \mathcal{V}}{\partial \dot{\mathbf{e}}}(\dot{\mathbf{e}}) + \frac{\partial \mathcal{D}_p}{\partial \dot{\mathbf{e}}_p}(\dot{\mathbf{e}}_p) \\ \boldsymbol{\sigma} \in \frac{\partial \mathcal{V}}{\partial \dot{\mathbf{e}}}(\dot{\mathbf{e}}) + \rho \frac{\partial \mathcal{E}}{\partial \boldsymbol{\varepsilon}_e}(\boldsymbol{\varepsilon}_e). \end{cases}$$

With the viscous dissipation potential  $\mathcal{V}(\dot{\mathbf{e}}) = \eta_N \dot{\mathbf{e}}_N^2 + \eta_T \dot{\mathbf{e}}_T^2$ , we get

$$\begin{cases} \boldsymbol{\sigma} = \boldsymbol{\sigma}_e + (\eta_T \dot{\mathbf{e}} + (\eta_N - \eta_T) \dot{\mathbf{e}}_N \mathbf{l}_d) \\ \boldsymbol{\sigma}_e \in \frac{\partial \mathcal{D}_p}{\partial \dot{\mathbf{e}}_p}(\dot{\mathbf{e}}_p) \\ \boldsymbol{\sigma}_e \in \rho \frac{\partial \mathcal{E}}{\partial \boldsymbol{\varepsilon}_e}(\boldsymbol{\varepsilon}_e). \end{cases}$$

We can once again generalize this equations to the framework of Implicit Standard Materials by allowing  $\mathcal{D}_p$  to depend on  $\boldsymbol{\sigma}_e$  as well as on  $\dot{\mathbf{e}}_p$ . Our full visco-elasto-plastic model then reads

$$\begin{cases} \boldsymbol{\sigma} = \boldsymbol{\sigma}_e + (\eta_T (\dot{\mathbf{e}}_e + \dot{\mathbf{e}}_p) + (\eta_N - \eta_T) (\dot{\mathbf{e}}_e + \dot{\mathbf{e}}_p)_N \mathbf{l}_d) \\ \boldsymbol{\sigma}_e \in \frac{\partial b}{\partial \dot{\mathbf{e}}_p}(\dot{\mathbf{e}}_p, \boldsymbol{\sigma}_e) \\ \boldsymbol{\sigma}_e \in \rho \frac{\partial \mathcal{E}}{\partial \boldsymbol{\varepsilon}_e}(\boldsymbol{\varepsilon}_e). \end{cases} \quad (1.30)$$

In this thesis, we will restrict ourselves to the case where  $\mathcal{E} = \mathcal{J}_0$ , i.e., the viscoplastic case where  $\boldsymbol{\varepsilon}_e = 0$ , and will choose  $b$  to be the Drucker–Prager bipotential defined in Property 1.7.

## Summary

In this first technical chapter, we have derived several equivalent formulations of the Signorini–Coulomb contact law (1–3) and of the non-associated Drucker–Prager flow rule with dilatancy (1.19), which shall prove useful for their numerical resolution. We have seen that these laws do not fit into the framework of Generalized Standard Materials, but can be described using the more comprehensive theory of Implicit Standard Materials.

**Other cones** In this chapter, we studied the nature of laws that were intrinsically linked to the Second-Order Cone (SOC). It is natural to wonder if and how the equivalences that we obtained would translate for other kind of convex cones; for instance, the Mohr–Coulomb yield surface.

As a first step, we can define the  $p$ -order cone of aperture  $\mu$ ,  $K_\mu^{(p)}$ , as

$$K_\mu^{(p)} := \{ \mathbf{x} \in \mathbb{R}^d, \mu \mathbf{x}_N \geq \|\mathbf{x}_T\|_p \}.$$

For  $p = 1$  or  $p = +\infty$ , the  $p$ -order cone becomes a pyramid, which yields a popular way of approximating the Coulomb law in numerical solvers, as in e.g., (Klarbring 1987). Its dual cone is given by  $(K_\mu^{(p)})^* = K_{\frac{1}{\mu}}^{(q)}$ , with  $q$  such that  $\frac{1}{p} + \frac{1}{q} = 1$ . The *de Saxcé* change of variable must then be adapted to  $\tilde{\mathbf{u}}^{(p)} := \mathbf{u} + \mu \|\mathbf{u}_T\|_q \mathbf{n}$  so that the  $p$ -order Coulomb friction law reads  $\tilde{\mathbf{u}}^{(p)} \in -\mathcal{N}_{K_\mu^{(p)}}(\mathbf{r})$ , or equivalently  $K_{\frac{1}{\mu}}^{(q)} \ni \tilde{\mathbf{u}}^{(p)} \perp \mathbf{r} \in K_\mu^{(p)}$ .

**Remaining equations** The set of equations that we obtained is not sufficient to fully determine a dynamical system. For the discrete contact mechanics, we have two variables,  $\mathbf{u}$  and  $\mathbf{r}$ , but only one equation linking them. Similarly, our visco-elasto-plastic model (1.30) boasts four variables  $(\boldsymbol{\sigma}_t, \boldsymbol{\sigma}_e, \dot{\boldsymbol{\epsilon}}_p$  and  $\dot{\boldsymbol{\epsilon}}_e)$ , but only three equations. We therefore need one more relationship; it will be provided the so-called conservation equations, which describe how a local stress affects the global motion of the dynamical system.

In the next chapter, we will see how to integrate the Coulomb contact law within the framework of Discrete-Element Modeling. Starting from chapter 6, we will go back to our viscoplastic model (1.30), and apply it to the context of Finite-Element Modeling of granular flows.



## 2 Modeling contacts within the Discrete Element Method

In this chapter, we will first briefly present a few models of mechanical systems, and write their unconstrained dynamics. Then, we will see how we can account for self or external unilateral contacts with friction.

### 2.1 A few mechanical models for rigid and deformable bodies in finite dimension

Before going on with more general models, let us start with the simple case of a rigid-body,

#### 2.1.1 Rigid-body dynamics

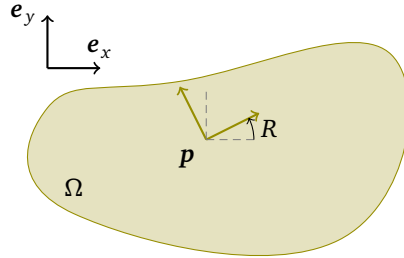
The kinematics of a rigid-body  $\Omega$  can be entirely described by a *frame*  $\mathbf{Q}$ ; that is, the combination of a world position  $\mathbf{p}$  and a rotation matrix  $R$ .

**Rotations** Formally, rotations matrices in dimension  $d$  are the elements of  $\mathcal{SO}(d)$ , the group of orthonormal matrices with positive determinant. However, such description is not very convenient in numerical programs;  $d \times d$  coefficients have to be stored, composing rotations requires computing expensive matrix-matrix products, and complex constraints, prone to drift, have to be enforced to ensure that the matrices stay in  $\mathcal{SO}(d)$ . It is thus preferable to use lower-dimensional parameterizations of the set of rotations matrices.

In 2D, the rotation matrix  $R$  can be represented simply with a scalar angle  $\theta$ ;

$$R = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} = \exp \begin{pmatrix} 0 & \theta \\ -\theta & 0 \end{pmatrix}.$$

Conversely, Euler's rotation theorem states that rotations in 3D can be described by the combination of an axis  $\mathbf{a}$  (a unit vector in  $\mathbb{R}^3$ ) and an angle  $\theta$ . This decomposition is non-unique, since the angle can be chosen modulo  $\pi$ , and the axis can be chosen arbitrarily when  $\theta = 0$  or  $\theta = \pi$ . More precisely, using Rodrigues's formula, every rotation matrix in 3D can be expressed



**Figure 2.1:** A rigid body  $\Omega$  with center of mass  $\mathbf{p}$  and rotation  $R$  w.r.t. the world's reference frame



as  $R = \exp[\mathbf{e}_\wedge]$ , where  $[\mathbf{e}_\wedge]$  is the skew-symmetric matrix corresponding to the application of the cross product  $\mathbf{e} \wedge \cdot$ ,

$$[\mathbf{e}_\wedge] := \begin{pmatrix} 0 & -e_z & e_y \\ e_z & 0 & -e_x \\ -e_y & e_x & 0 \end{pmatrix}.$$

The angle of the rotation  $R$  is then  $\theta = \|\mathbf{e}\|$ , and the axis is  $\frac{\mathbf{e}}{\|\mathbf{e}\|}$  (or any unit vector when  $\mathbf{e} = \mathbf{0}$ ). Alternatively, the set of 3D rotations can be parameterized by the three angles of the successive rotations around each axis of an orthogonal basis — the so-called Euler angles. However, while very compact, those representations are not much more satisfying for a numerical implementation; they require numerous trigonometric operations which are computationally expensive and may accumulate errors, and compositing rotations is still non-trivial. Euler angles are also subject to gimbal-locking, that is, in certain configurations they may loose a degree of freedom.

A better way of representing the matrix  $R$  in a compact yet well-defined manner is through the mean of a unit quaternion  $\mathbf{q} \in \mathcal{H}_1$ ,  $\mathcal{H}_1 = \{\mathbf{q} := [w, x, y, z] := w + x\mathbf{i} + y\mathbf{j} + z\mathbf{k}, (w, x, y, z) \in \mathbb{R}^4, \|\mathbf{q}\| = 1\}$  with  $\mathbf{i}$ ,  $\mathbf{j}$  and  $\mathbf{k}$  square roots of  $-1$ . Numerous texts have been written about the relationship between unit quaternions, matrix exponentials and rotations, we refer the reader to e.g., (Cadoux 2009, Appendix C) for a concise introduction focused on their application to rigid-body dynamics. In practice, a unit quaternion  $\mathbf{q}$  can also be seen as the combination of an angle  $\theta$  and a unit direction vector  $\mathbf{a}$ . Indeed,  $\mathbf{q} = [\cos(\frac{\theta}{2}); \sin(\frac{\theta}{2})\mathbf{a}]$  is unitary and reciprocally, every unit quaternion can be decomposed in such a way: for instance,  $\theta = 2 \operatorname{atan2}(\| [x, y, z]^T \|, w)$ , and  $\mathbf{a} = \frac{[x, y, z]^T}{\| [x, y, z]^T \|}$  if  $\sin \theta \neq 0$ , or any other unit vector otherwise. The decomposition is therefore not unique, but the rotation defined by the axis-angle couple is unique.

Now, the major advantage of using quaternions instead of rotation matrices or their angle-axis representation is that operations on quaternions are well-suited for numerical computations. Indeed,  $\mathcal{H}_1$  is a group for the canonical quaternion product  $\cdot \times \cdot$ , and  $\mathbf{q}_1 \times \mathbf{q}_2$  coincide with the composition of the rotations represented by  $\mathbf{q}_1$  and  $\mathbf{q}_2$ . At the same time, it is much cheaper and precise to compute than a matrix-matrix product, and less prone to numerical drift (even though the resulting quaternion must be kept unitary). Moreover, the application  $\mathcal{H}_1 \times \mathbb{R}^3 \rightarrow \mathbb{R}^3$ ,  $(\mathbf{q}, \mathbf{p}) \mapsto \mathbf{r}$  such that

$$[0; \mathbf{r}] = \mathbf{q}^* \times [0; \mathbf{p}] \times \mathbf{q} \quad \text{with } ([w, x, y, z]^T)^* := [w, -x, -y, -z]^T$$

coincide with the application of the rotation  $R(\mathbf{q})$  represented by the quaternion  $\mathbf{q}$  to the vector  $\mathbf{p}$ , i.e.,  $\mathbf{r} = R(\mathbf{q})\mathbf{p}$ . Compared to Rodrigues formula, this does not involve any trigonometric operation, and is therefore much cheaper to compute numerically. Finally, unit quaternions can also be expressed as the exponential (w.r.t. the quaternion product) of a purely imaginary quaternion,

$$\left[ \cos\left(\frac{\theta}{2}\right); \sin\left(\frac{\theta}{2}\right)\mathbf{a} \right] = \exp\left[0; \frac{\theta}{2}\mathbf{a}\right], \quad (2.1)$$

which naturally links them to the solutions of first-order linear differential equations, and yields an insight on their integration through time.

**Equations of motions** Without loss of generality, we will assume that our kinematic frame  $\mathbf{Q} = (\mathbf{p}; \mathbf{q})$  is positioned at the center of mass of the rigid body. Let also  $\mathbf{v}$  and  $\boldsymbol{\omega}$  be the linear and angular velocity of  $\Omega$  at its center of mass; the time derivative of the frame  $\mathbf{Q} := (\mathbf{p}; \mathbf{q})$  can be expressed as  $\dot{\mathbf{Q}} := (\mathbf{v}; \boldsymbol{\omega})$ . Indeed,  $\frac{d\mathbf{p}}{dt} = \mathbf{v}$ , and since the velocity at any point  $\mathbf{x}$  of  $\Omega$  may be computed as  $\mathbf{v}(\mathbf{x}) = \mathbf{v} + \boldsymbol{\omega} \wedge (\mathbf{x} - \mathbf{p})$ , it holds that  $\frac{dR(\mathbf{q})}{dt}\mathbf{x} = \boldsymbol{\omega} \wedge (R(\mathbf{q})\mathbf{x})$ , which leads to  $\frac{dR(\mathbf{q})}{dt} = [\boldsymbol{\omega}]_\wedge R(\mathbf{q})$ . This identity transpose to the quaternion framework as a 4-dimensional first-order linear differential equation,  $\frac{d\mathbf{q}}{dt} = \frac{1}{2}[0; \boldsymbol{\omega}] \times \mathbf{q}$ . Note that if the direction of  $\boldsymbol{\omega}$  is constant between  $t_0$  and  $t$ , i.e.,  $\boldsymbol{\omega}(t) = \dot{\theta}(t)\mathbf{a}$ , using (2.1) the solution in time is simply given by  $\mathbf{q}(t) = [\cos\frac{\theta(t)}{2}; \sin\frac{\theta(t)}{2}\mathbf{a}] \times \mathbf{q}(t_0)$ , with  $\theta(t)$  such that  $\theta(t_0) = 0$  and  $\frac{d\theta}{dt} = \dot{\theta}$ .

Supposing a Galilean frame of reference, Newton's second law over a elementary volume  $d\mathbf{x}$  around  $\mathbf{x}$  reads

$$\rho(\mathbf{x}) \frac{d\mathbf{v}(\mathbf{x})}{dt} d\mathbf{x} = \mathbf{g}(\mathbf{x}, \mathbf{v}(\mathbf{x}), t) d\mathbf{x}, \quad (2.2)$$

where  $\mathbf{g}$  is the volumetric force applied at each point of  $\Omega$ , and  $\rho(\mathbf{x})$  is the density of the rigid body. As  $\frac{d\mathbf{v}(\mathbf{x})}{dt} = \mathbf{v} + \frac{d\boldsymbol{\omega}}{dt} \wedge (\mathbf{x} - \mathbf{p}) + \boldsymbol{\omega} \wedge (\boldsymbol{\omega} \wedge (\mathbf{x} - \mathbf{p}))$ , integrating over  $\Omega$  yields the conservation of linear momentum equation,

$$m \frac{d\mathbf{p}}{dt} = \mathbf{g}(\mathbf{Q}, \dot{\mathbf{Q}}, t) := \int_{\Omega} \mathbf{g}(t, \mathbf{x}, \mathbf{v}(\mathbf{x})) d\mathbf{x}, \quad (2.3)$$

and applying the vector product  $(\mathbf{x} - \mathbf{p}) \wedge \cdot$  to each side of Equation (2.2) before integrating yields the conservation of angular momentum equation,

$$I(\mathbf{Q}) \frac{d\boldsymbol{\omega}}{dt} + \boldsymbol{\omega} \wedge (I(\mathbf{Q})\boldsymbol{\omega}) = \mathbf{l}(\mathbf{Q}, \dot{\mathbf{Q}}, t) := \int_{\Omega} (\mathbf{x} - \mathbf{p}) \wedge \mathbf{g}(t, \mathbf{x}, \mathbf{v}(\mathbf{x})) d\mathbf{x}, \quad (2.4)$$

where  $m$  is the mass and  $I$  the inertia matrix of  $\Omega$  computed as

$$\begin{aligned} m &= \int_{\Omega} \rho(\mathbf{x}) d\mathbf{x} \\ I(\mathbf{q}) &= \text{Tr}(T(\mathbf{q})) \mathbb{I}_3 - T(\mathbf{q}) \\ T(\mathbf{q}) &= \int_{\Omega} \rho(\mathbf{x}) (\mathbf{x} - \mathbf{p}) \otimes (\mathbf{x} - \mathbf{p}) d\mathbf{x}. \end{aligned}$$

We see that the inertia matrix  $I(\mathbf{q})$ , which is computed in the world's coordinate system, depends of the orientation of the frame  $\mathbf{Q}$ . However, it can easily be deduced from the computation of the inertia matrix  $I|_{\mathbf{Q}}$  in the frame attached to the rigid body,  $I(\mathbf{q}) = R(\mathbf{Q}) I|_{\mathbf{Q}} R(\mathbf{Q})^T$ .

Finally, we can group the conservation equations (2.3–2.4) into a single one,

$$M(\mathbf{Q}) \frac{d\dot{\mathbf{Q}}}{dt} = \mathbf{f}(\mathbf{Q}, \dot{\mathbf{Q}}, t), \quad (2.5)$$

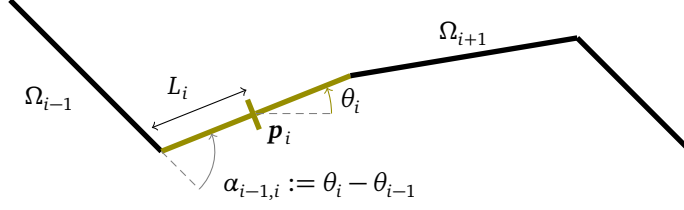
$$\text{with } M(\mathbf{Q}) := \begin{pmatrix} m\mathbb{I}_3 & 0 \\ 0 & I(\mathbf{q}) \end{pmatrix} \quad \text{and} \quad \mathbf{f}(\mathbf{Q}, \dot{\mathbf{Q}}, t) := \begin{pmatrix} \mathbf{g}(\mathbf{Q}, \dot{\mathbf{Q}}, t) \\ \mathbf{l}(\mathbf{Q}, \dot{\mathbf{Q}}, t) - \boldsymbol{\omega} \wedge (I(\mathbf{Q})\boldsymbol{\omega}) \end{pmatrix}.$$

**Multiple bodies** Now that we know how to model the dynamics of a single rigid-body, we can extend our equations to handle more complex objects, such as articulate bodies. Let  $(\Omega_i)$  be a set of  $N \in \mathbb{N}$  rigid-bodies with frames  $(\mathbf{Q}_i)$ . Concatenating the degrees of freedom, velocities and forces as  $\mathbf{Q} := [\mathbf{Q}_1; \dots; \mathbf{Q}_N]$ ,  $\dot{\mathbf{Q}} := [\dot{\mathbf{Q}}_1; \dots; \dot{\mathbf{Q}}_N]$  and assembling the  $N$  mass matrices  $M_i$  into a single block-diagonal matrix  $M := \text{diag}(M_1, \dots, M_N)$ , the dynamics of the whole system may be written as

$$\begin{cases} M(\mathbf{Q}) \frac{d\dot{\mathbf{Q}}}{dt} = \mathbf{f}(\mathbf{Q}, \dot{\mathbf{Q}}, t) + \mathbf{f}^{\text{cst}} \\ C(\mathbf{Q}) = \mathbf{0}. \end{cases}$$

where  $C : (\mathbb{R}^3 \times \mathbb{R}^4)^m \rightarrow \mathbb{R}^n$  is a function defining constraints on the system, and  $\mathbf{f}_i^{\text{cst}}$  is the sum of the forces applied on the body  $i$  through the constraints. Now, D'Alembert's principle<sup>1</sup> states that the work of the constraint forces  $\mathbf{f}_i^{\text{cst}}$  vanishes during a virtual displacement  $\delta \mathbf{Q}$  (that is, an infinitesimal displacement such that the constraints are still satisfied, i.e., such that

<sup>1</sup>Traité de dynamique, Jean le Rond d'Alembert, 1743



**Figure 2.2:** An articulated rigid-body chain. Degrees-of-freedom are  $(\mathbf{p}_i, \theta_i)$ , and joints are modeled by holonomic constraints.

$\delta \mathbf{Q} \in \text{Ker } \frac{\partial C}{\partial \mathbf{Q}}(\mathbf{Q})$ . This means that  $\mathbf{f}^{\text{cst}} \in (\text{Ker } \frac{\partial C}{\partial \mathbf{Q}}(\mathbf{Q}))^\perp = \text{Im} \left( \frac{\partial C}{\partial \mathbf{Q}}(\mathbf{Q}) \right)^\top$ , so the dynamics can be written again as

$$\begin{cases} M(\mathbf{Q}) \frac{d\dot{\mathbf{Q}}}{dt} = \mathbf{f}(\mathbf{Q}, \dot{\mathbf{Q}}, t) + \left( \frac{\partial C}{\partial \mathbf{Q}}(\mathbf{Q}) \right)^\top \boldsymbol{\lambda} \\ C(\mathbf{Q}) = 0 \end{cases} \quad (2.6)$$

with  $\boldsymbol{\lambda} \in \mathbb{R}^n$ .

For instance, one may enforce that two bodies share the same center of mass by taking  $C = (\mathbf{p}_1 - \mathbf{p}_2)$ , which will lead to  $\mathbf{f}^{\text{cst}} = (\boldsymbol{\lambda}; -\boldsymbol{\lambda})$  with  $\boldsymbol{\lambda}$  in  $\mathbb{R}^3$ .

**Chain of rigid segments** As an illustrative example, we can attempt to model a 2D flexible chain of  $N$  inextensible segments of length  $2L_i$  and width  $r_i \ll L_i$  with frames given by  $\mathbf{Q}_i = (\mathbf{p}_i, \theta_i)$ , as in Figure 2.2. The end points of  $\Omega_i$  are given by  $\mathbf{p}_i \pm L_i \mathbf{e}(\theta_i)$ , with  $\mathbf{e}(\theta) := (\cos \theta; \sin \theta)$ . The velocity of the point of abscissa  $s$  on the  $i^{\text{th}}$  segment is  $\mathbf{v}_i(s) = \mathbf{v}_i + s\omega \mathbf{e}(\theta + \frac{\pi}{2})$ . By homogeneity with the 3D equations, we will reuse the notation  $\mathbf{x} \wedge \mathbf{y}$  for the determinant of two vectors in  $\mathbb{R}^2$ ,  $\mathbf{x} \wedge \mathbf{y} := \det(\mathbf{x}, \mathbf{y})$ .

Assuming that the first link is clamped at the world's frame origin, and that the chain's end is free to move, our nonlinear constraint function  $C$  could be defined as

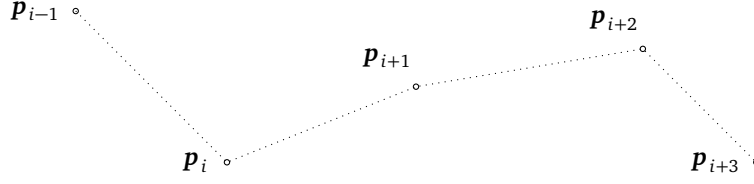
$$C(\mathbf{Q}) = \begin{pmatrix} \mathbf{p}_1 \\ \mathbf{p}_1 + L_1 \cos \theta_1 - \mathbf{p}_2 + L_2 \cos \theta_2 \\ \dots \\ \mathbf{p}_{N-1} + L_{N-1} \cos \theta_{N-1} - \mathbf{p}_N + L_N \cos \theta_N \end{pmatrix}.$$

Next, assuming a constant density  $\rho_i$  over the volume of the  $i^{\text{th}}$  segment, the mass of each cylinder is  $m_i = 2\rho_i L_i r_i$ . In 2D, the inertia matrix  $I_i$  boils down to a single scalar coefficient, in our case

$$I_i = \rho_i r_i \int_{-L_i}^{L_i} (s \mathbf{e}(\theta_i)) \wedge (s \mathbf{e}(\theta_i + \frac{\pi}{2})) = \frac{1}{3} m_i L_i^2.$$

Finally, we assume that the volumetric force (actually, lineic in our slender case) can be decomposed into the following contributions:

- The gravity, which we will assume along the  $y$  axis.  $\mathbf{g}_i^g = -\rho_i r_i g \mathbf{e}_y$ , and  $\mathbf{l}_i^g = \mathbf{0}$ .
- A viscous air friction term opposing the velocity at each point of the rod,  $\mathbf{g}^a(\mathbf{x}) = -\eta \mathbf{v}(\mathbf{x})$ ,



**Figure 2.3:** *N-tuple pendulum modeled with point masses. Degrees of freedom are limited to the positions ( $\mathbf{p}_i$ ), and each element's length is maintained by an holonomic constraint  $\|\mathbf{p}_{i+1} - \mathbf{p}_i\| = 2L_i$ .*

yielding

$$\begin{aligned} \mathbf{g}_i^a &= -\eta \int_{s=-L_i}^{L_i} \mathbf{v}_i(\mathbf{x}) ds = 2L_i \eta \mathbf{v}_i \\ l_i^a &= -\eta \int_{s=-L_i}^{L_i} (\mathbf{x}(s) - \mathbf{p}_i) \wedge \mathbf{v}_i(\mathbf{x}(s)) ds \\ &= -\eta \int_{s=-L_i}^{L_i} s \mathbf{e}(\theta_i) \wedge (\omega_i(s \mathbf{e}(\theta_i))) ds \\ &= -\frac{2}{3} L_i^3 \omega_i. \end{aligned}$$

- The force derived from a bending energy  $E_{i,j}^b$ , which, assuming a Hookean torsion spring, is quadratic in the bending angle  $\alpha_{i,j} := \theta_j - \theta_i$  between adjacent rods:  $E_{i,j}^b := \frac{1}{2} K \alpha_{i,j}^2$ . The bending energy  $E_{i,j}^b$  therefore causes the  $j^{\text{th}}$  element to induce a torque  $K \alpha_{i,j}$  on the  $i^{\text{th}}$  element. Writing the total bending force and torque applied by the  $i^{\text{th}}$  element at the link between the  $i^{\text{th}}$  and  $(i+1)^{\text{th}}$  elements gives

$$\begin{aligned} \mathbf{0} &= \int_{-L}^L \mathbf{g}_{i \rightarrow i+1}^b(s) ds \\ -K \alpha_{i,i+1} &= \int_{-L}^L (s-L) \mathbf{e}(\theta_i) \wedge \mathbf{g}_{i \rightarrow i+1}^b(s) ds = \int_{-L}^L s \mathbf{e}(\theta_i) \wedge \mathbf{g}_{i \rightarrow i+1}^b(s) ds, \end{aligned}$$

from which we deduce  $\mathbf{g}_i^b = \mathbf{0}$  and

$$l_i^b = -K (\alpha_{(i-1),i} - \alpha_{i,(i+1)}) = -K (2\theta_i - \theta_{i-1} - \theta_{i+1}).$$

To summarize, the dynamics of our articulated chain are governed by a second-order ordinary differential equation with a banded, symmetric positive semi-definite mass-matrix, non-linear in  $\mathbf{Q}$ . The force  $\mathbf{f}$  is linear in  $\mathbf{Q}$  and  $\dot{\mathbf{Q}}$ , and the constraints are non-linear in  $\mathbf{Q}$ .

### 2.1.2 Lumped system

Another popular way of modeling deformable objects is to consider them as a set of point at which all of the body's mass is concentrated, possibly with constraints on their relative positions. That is, writing  $\rho_i(\mathbf{x}) = \delta(\mathbf{x} - \mathbf{p}_i)$ , where  $\delta$  denotes the Dirac distribution such that for any field  $\mathbf{f}$  over  $\Omega$ ,  $\int_{\Omega} \mathbf{f}(\mathbf{x}) \delta(\mathbf{y}) d\mathbf{x} = \mathbf{f}(\mathbf{y})$ . We see immediately that with such a mass distribution,  $I_i = \mathbf{0}$  and  $l_i = 0$ ; therefore, we need only concern ourselves with the evolution of the position of linear velocity variables  $\mathbf{p}_i$  and  $\mathbf{v}_i$ , and can discard  $\mathbf{q}_i$  and  $\omega_i$ . The mass matrix  $M$  is then constant, diagonal and positive definite, simplifying greatly the numerical integration of the dynamical

system. However, this approach is not without drawbacks; the discretization of the physical model might be less accurate, and the forces and constraints will often be nonlinear.

A  $N$ -tuple pendulum can be modeled in this framework by positioning the point masses at the links between consecutive elements, as well as at the start and end points of the chain (Figure 2.3). For our chain of  $N$  elements, we therefore have  $(N + 1)$  point masses with coordinates  $(\mathbf{p}_i)$ , with the constraint that for  $1 \leq i \leq N$ ,  $\|\mathbf{p}_{i+1} - \mathbf{p}_i\| = L_i$ . The mass associated to the  $i^{\text{th}}$  node is then  $m_i := \rho(r_{i-1}L_{i-1} + r_iL_i)$  (with the convention  $L_0 = 0$ ). Note that this model is *not* equivalent to the chain of rigid-bodies presented in the previous section; the kinetic energy of an individual slender rod is  $\frac{m}{2}(\mathbf{v}^2 + \frac{1}{3}L^2\omega_i^2)$ , while that of the corresponding two-masses system would be  $\frac{m}{2}(\mathbf{v}^2 + L^2\omega_i^2)$ . Our goal here is simply to give an insight about the structure of the resulting equations; for the rigid-body chain version of the  $N$ -tuple pendulum, only the value (but not the structure) of the inertia matrix and air friction torque would have to be modified.

Now, the forces acting on the node  $\mathbf{p}_i$  are the sum of:

1. the gravity contribution,  $\mathbf{f}_i^g = -m_i g \mathbf{e}_y$ .
2. the force induced by air friction,  $\mathbf{f}_i^a = -(L_{i-1} + L_i)\eta \mathbf{v}_i$ .
3. the bending force derived from the bending energy  $E_{i,j}^b = \frac{1}{2}K\theta_{i,j}^2$ ,

$$\begin{aligned} \mathbf{f}_i^b &= -\alpha_{(i-2),(i-1)}^\top \frac{\partial \alpha_{(i-2),(i-1)}}{\partial \mathbf{p}_i} - \alpha_{(i-1),i}^\top \frac{\partial \alpha_{(i-1),i}}{\partial \mathbf{p}_i} - \alpha_{i,(i+1)}^\top \frac{\partial \alpha_{i,(i+1)}}{\partial \mathbf{p}_i} \\ \alpha_{i,i+1} &= \sin^{-1} \left( \frac{\|\mathbf{y}_{i,i+1}\|}{4L_{i-1}L_i} \right) \frac{\mathbf{y}_{i,i+1}}{\|\mathbf{y}_{i,i+1}\|} \\ \mathbf{y}_{i,j} &:= (\mathbf{p}_i - \mathbf{p}_{i-1}) \wedge (\mathbf{p}_{i+1} - \mathbf{p}_i), \end{aligned}$$

which is nonlinear in  $\mathbf{Q}$ .

In brief, the lumped system boasts a diagonal, positive-definite mass matrix, but suffers from highly nonlinear constraints and forces.

### 2.1.3 Lagrangian mechanics

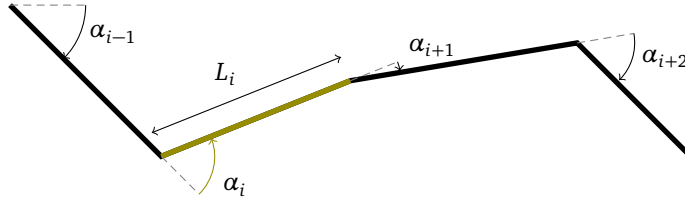
In the two previous approaches, modeling any complex system led to the introduction of constraints, which might seriously complicate the process of solving the dynamics differential equation. In contrast, Lagrangian mechanics opens the way for using alternative sets of coordinates. In particular, one may choose to directly parameterize the manifold of admissible kinematic variables with so-called *reduced* coordinates, which will eventually lead to a constraint-free system. For the sake of simplicity, we only consider this case below.

We consider a non-conservative dynamical system described by a set of  $m \in \mathbb{N}$  kinematic variables  $\mathbf{q}$  and its time-derivative  $\mathbf{v} := \frac{d\mathbf{q}}{dt}$ , the so-called generalized velocity. Assuming that  $\mathbf{q}$  are reduced coordinates, any infinitesimal change  $\delta \mathbf{q}$  yields a *virtual* displacement  $\delta \mathbf{x} = \sum \frac{\partial \mathbf{x}}{\partial q_i} \delta q_i$ . Expressing Newtons' second law of motion for any of those virtual displacements gives

$$\int_{s \in \Omega} \left( \frac{\partial \mathbf{x}}{\partial \mathbf{q}}(t, s, \mathbf{q}) \right)^\top \rho(s) \frac{d\mathbf{v}}{dt}(t, s, \mathbf{q}, \mathbf{v}) ds = \int_{s \in \Omega} \left( \frac{\partial \mathbf{x}}{\partial \mathbf{q}}(t, s, \mathbf{q}) \right)^\top \mathbf{g}(t, \mathbf{x}(t, s, \mathbf{q}), \mathbf{v}(t, s, \mathbf{q}, \mathbf{v})) ds$$

where  $\mathbf{g}$  is the external force applied at each point of  $\Omega$ . Moreover, we have also

$$\begin{aligned} & \int_{s \in \Omega} \rho(s) \frac{d\mathbf{v}}{dt}(t, s, \mathbf{q}, \mathbf{v}) \frac{\partial \mathbf{x}}{\partial \mathbf{q}}(t, s, \mathbf{q}) ds \\ &= \int_{s \in \Omega} \rho(s) \left[ \frac{d}{dt} \left( \mathbf{v}(t, s, \mathbf{q}, \mathbf{v}) \frac{\partial \mathbf{x}}{\partial \mathbf{q}}(t, s, \mathbf{q}) \right) - \mathbf{v}(t, s, \mathbf{q}, \mathbf{v}) \frac{d}{dt} \frac{\partial \mathbf{x}}{\partial \mathbf{q}}(t, s, \mathbf{q}) \right] ds \\ &= \int_{s \in \Omega} \rho(s) \left[ \frac{d}{dt} \left( \mathbf{v}(t, s, \mathbf{q}, \mathbf{v}) \frac{\partial \mathbf{v}}{\partial \mathbf{v}}(t, s, \mathbf{q}, \mathbf{v}) \right) - \mathbf{v}(t, s, \mathbf{q}, \mathbf{v}) \frac{\partial \mathbf{v}}{\partial \mathbf{q}}(t, s, \mathbf{q}, \mathbf{v}) \right] ds \\ &= \frac{1}{2} \int_{s \in \Omega} \rho(s) \left[ \frac{d}{dt} \frac{\partial \mathbf{v}^2}{\partial \mathbf{v}}(t, s, \mathbf{q}, \mathbf{v}) - \frac{\partial \mathbf{v}^2}{\partial \mathbf{q}}(t, s, \mathbf{q}, \mathbf{v}) \right] ds. \end{aligned}$$



**Figure 2.4:** With the Lagrangian approach and reduced coordinates, the angles  $(\alpha_i)$  between successive elements are the sole degrees of freedom for our 2D articulated chain, and no supplemental constraint is necessary. However, the position of each segment depends on the angle at all prior joints, yielding a dense numerical system.

Introducing the kinetic energy  $\mathcal{E}^c$  of the system,

$$\mathcal{E}^c(t, \mathbf{q}, \mathbf{v}) := \frac{1}{2} \int_{s \in \Omega} \rho(s) \mathbf{v}^2(t, s, \mathbf{q}, \mathbf{v}) ds,$$

we obtain the Euler–Lagrange dynamics equation<sup>2</sup>,

$$\frac{d}{dt} \frac{\partial \mathcal{E}^c}{\partial \mathbf{v}}(t, \mathbf{q}, \mathbf{v}) - \frac{\partial \mathcal{E}^c}{\partial \mathbf{q}}(t, \mathbf{q}, \mathbf{v}) = \int_{s \in \Omega} \left( \frac{\partial \mathbf{x}}{\partial \mathbf{q}}(t, s, \mathbf{q}) \right)^\top \mathbf{g}(t, s, \mathbf{q}, \mathbf{v}) ds. \quad (2.7)$$

It is convenient to express forces that can be derived from potentials as such; for instance, a large class of systems can be written as

$$\frac{d}{dt} \frac{\partial \mathcal{E}^c}{\partial \mathbf{v}}(t, \mathbf{q}, \mathbf{v}) - \frac{\partial \mathcal{E}^c}{\partial \mathbf{q}}(t, \mathbf{q}, \mathbf{v}) + \frac{\partial \mathcal{E}^p}{\partial \mathbf{q}}(\mathbf{q}) = - \frac{\partial \mathcal{D}}{\partial \mathbf{v}}(\mathbf{q}, \mathbf{v}) \quad (2.8)$$

where  $\mathcal{E}^p$  is the potential energy of the system and  $\mathcal{D}$  is a dissipation potential.

$\mathcal{E}^c$  being generally quadratic w.r.t.  $\mathbf{v}$ , expressing the derivatives and integral of Equation (2.7) leads to

$$M(t, \mathbf{q}) \frac{d\mathbf{v}}{dt} = \mathbf{f}(t, \mathbf{q}, \mathbf{v})$$

where  $M(t, \mathbf{q})$  is symmetric positive semi-definite, with  $M_{ij} = \int_{\Omega} \rho(s) \frac{\partial^2 \mathbf{v}}{\partial v_i \partial v_j}(t, s, \mathbf{q}, \mathbf{v}) ds$ , and the generalized force  $\mathbf{f}$  groups terms from internal, external and inertial forces. The expression of the dynamics of the system in the reduced coordinates  $\mathbf{q}$  is therefore similar to that of the unconstrained dynamics equation (2.5) in the framework of Section 2.1.1. However, as we will now see on our 2D articulated chain example, the mass-matrix  $M(t, \mathbf{q})$  will not necessarily be sparse anymore.

**Articulated chain** Let  $\mathbf{q} = (\alpha_i) \in \mathbb{R}^N$ , with  $\alpha_i$  the angle between the  $(i-1)^{\text{th}}$  and the  $i^{\text{th}}$  segments, and with  $\alpha_0$  defined w.r.t. to the clamping angle  $\theta_0$  which we will assume to be 0. The generalized velocity is  $\mathbf{v} = (\omega_i)$ , with  $\omega_i = \frac{d\alpha_i}{dt}$ .

**Potential energy** The potential energy is the sum of two terms, the bending energy  $E^b$  and the contribution of gravity  $E^g$ . We have  $E_{i-1,i}^b(\mathbf{q}) = \frac{1}{2} K \alpha_i^2$ , and therefore

$$\frac{\partial E^b}{\partial \mathbf{q}}(\mathbf{q}) = \sum_i \frac{\partial E_{i-1,i}^b}{\partial \mathbf{q}}(\mathbf{q}) = K \mathbf{q}^\top.$$

<sup>2</sup>First established by Joseph Louis Lagrange from the principle of least action in 1756, then from D'Alembert principle in 1788.

The potential energy associated to the gravity is more complex. The position of the point at abscissa  $s$  of the  $i^{\text{th}}$  element can be recovered from a recursive formula  $\mathbf{x}_i(s) = \mathbf{x}_{i-1}(L_{i-1}) + (s + L_i)(\prod_{j=0}^i R(\alpha_j))\mathbf{e}_x$ . and the clamping condition  $\mathbf{x}_0 = \mathbf{0}$ . We can therefore evaluate

$$E^g(\mathbf{q}) = \sum_i \rho_i r_i \int_{-L_i}^{L_i} \langle \mathbf{x}_i(s), \mathbf{e}_y \rangle ds = \sum_i m_i \langle \mathbf{x}_i(0), \mathbf{e}_y \rangle,$$

and the gravity generalized force,  $-\frac{\partial E^g}{\partial \mathbf{q}}(\mathbf{q})$ , can then be computed using

$$\frac{\partial \mathbf{x}_i(s)}{\partial \mathbf{q}_j}(\mathbf{q}) = \begin{cases} \left( \prod_{k=1}^{j-1} R(\alpha_k) \right) R(\alpha_j + \frac{\pi}{2}) (\mathbf{x}_i(s) - \mathbf{x}_j(-L_j)) & \text{if } j \leq i \\ \mathbf{0} & \text{if } j > i \end{cases}$$

**Air friction** The generalized force  $f^a$  due to air friction derives from the dissipation potential

$$\mathcal{D}(\mathbf{q}, \mathbf{v}) := \frac{\eta}{2} \sum_i \int_{-L_i}^{L_i} \mathbf{v}_i^2(s) ds,$$

where  $\mathbf{v}_i(s)$  can be once again be computed using a recursive formula,

$$\mathbf{v}_i(s) = \mathbf{v}_{i-1}(L) + (s + L) \left( \sum_{j=1}^i \omega_j \right) \mathbf{e} \left( \frac{\pi}{2} + \sum_{j=1}^i \alpha_j \right). \quad (2.9)$$

The expression the friction force is then

$$f_a = -\frac{\partial \mathcal{D}(\mathbf{q}, \mathbf{v})}{\partial \mathbf{v}} = -\eta \sum_i \int_{-L_i}^{L_i} \left( \frac{\partial \mathbf{x}_i}{\partial \mathbf{q}}(s) \right)^\top \mathbf{v}_i(s) ds,$$

**Inertia** Finally, it remains to consider the derivatives of the kinetic energy  $\mathcal{E}^c$  w.r.t. the generalized positions and velocity. We have

$$\mathcal{E}^c = \frac{1}{2} \sum_i \int_{-L_i}^{L_i} \mathbf{v}_i^2(s) ds = \frac{1}{2} \sum_i m_i \left( \mathbf{v}_i^2(0) + \frac{1}{3} L_i^2 \omega_i^2 \right).$$

We see from Equation (2.9) that  $\mathbf{v}_N(0)$  depends linearly on all the velocity components ( $\omega_i$ ); indeed,

$$\frac{\partial \mathbf{v}_N(0)}{\partial \omega_i} = (s + L) \mathbf{e} \left( \frac{\pi}{2} + \sum_{j=1}^N \alpha_j \right) + 2L \sum_{k=i}^{N-1} \mathbf{e} \left( \frac{\pi}{2} + \sum_{j=1}^k \alpha_j \right).$$

Therefore,  $\forall i, j \in [1, N]^2$ ,

$$\frac{\partial^2 \mathbf{v}_N^2(0)}{\partial \omega_i \partial \omega_j} = \left( \frac{\partial \mathbf{v}_N(0)}{\partial \omega_i} \right) \left( \frac{\partial \mathbf{v}_N(0)}{\partial \omega_j} \right)^\top$$

and we deduce that the mass matrix  $M(t, \mathbf{q})$  is generally dense.

#### 2.1.4 Discussion

While being far from exhaustive, we illustrated a few different approaches for deriving the equations of motion of a mechanical system with kinematic constraints. The logical follow-up question is then: which one should we use ? There is no definitive answer, but we can compare the different approaches on selected criterions.

To this aim, let us remark that all previous derivations can be coalesced into a single canonical set of equations,

$$M(t, \mathbf{q}) \frac{d\mathbf{v}}{dt} = \mathbf{f}(t, \mathbf{q}, \mathbf{v}) + \left( \frac{\partial C}{\partial \mathbf{q}}(t, \mathbf{q}) \right)^\top \boldsymbol{\lambda}$$

$$C(t, \mathbf{q}) = \mathbf{0}$$

where  $\mathbf{q}$  is a set a generalized coordinates, and  $\mathbf{v}$  a set of generalized velocities. However, the structure of the mass-matrix  $M$ , the properties of force  $\mathbf{f}$ , and the kinematic constraints  $C$  showed substantial variations depending on the nature of coordinates.

**Modeling capacity** Reduced coordinates may not always be adequate. For once, not every system can be described with solely holonomic constraints, and even then, finding a constraint-free parametrization is a hard problem — such is the case of inextensible shells (Casati 2015). Note, however that the Lagrangian approach does not prescribe the use of reduced coordinates; it is entirely possible, using again D'Alembert's principle, to add supplemental constraints to the Lagrange equation of motion (2.7), yielding an “intermediate” set of coordinates. However, the Lagrange equation (2.7) does require that the velocity variable  $\mathbf{v}$  be the derivative in time of the position variable. For instance, this prohibits using a quaternion for the position variable and an angular velocity as the velocity variable, like we did for rigid bodies.

Conversely, point masses are a very rough discretization of a mechanical system, and might not model correctly some physical quantities, such as the angular momentum in our 2D example. Getting more faithful to the continuous model then means adding more point masses and constraints, making the equations more expensive to solve numerically.

**Numerical integration** In the presence of holonomic constraints, the equations of motion boil down to an index-3 differential algebraic equation (DAE). Numerical integration of such equations have been extensively studied in the literature, we point the reader to (Haddouni 2015) for a recent review.

On the one hand, the absence of additional constraints is obviously a huge advantage of the Lagrangian approach with reduced coordinates, as it allows the use of much simpler integrators. However, in practice supplemental constraints may be required to allow artistic control of the simulation, and deriving a new reduced kinematic model for each new scene is not practical — it is safer to devise a numerical integrator that can handle kinematic constraints anyway.

On the other hand, the main advantage of the point-masses method is that the mass matrix will always be diagonal and positive-definite, and thus trivial to inverse; on the contrary, slender 3D structures may have vanishing inertia terms when modeled with the Lagrangian of articulated rigid-body approach. Moreover, the mass matrix of reduced-coordinates models can be dense, leading to more expensive computations.

Another important point deals with the structure of the internal forces; having them depend linearly on the kinematic variables may greatly simplify the numerical integration process, as we shall see later. In our 2D example, we saw that all three forces were linear for the chain of rigid-bodies, while the bending force was nonlinear of the point-masses discretization, and both gravity and air friction forces were nonlinear for the reduced coordinates model. One may thus choose one or the other model depending on the respective influence of each force term.

## 2.2 Contacts

For the sake of simplicity, we first consider a system without kinematic constraints. Using any kind of coordinates, the equation of motion may be written as

$$M(t, \mathbf{q}) \frac{d\mathbf{v}}{dt} = \mathbf{f}(t, \mathbf{q}, \mathbf{v}) + \mathbf{f}^c, \quad \mathbf{f}^c := \int_{s \in \Omega} \left( \frac{\partial \mathbf{v}}{\partial \mathbf{v}}(t, \mathbf{s}, \mathbf{q}) \right)^\top \mathbf{g}^c(\mathbf{s}) d\mathbf{s},$$

where  $\mathbf{g}^c$  is a supplemental force density due to the contacts.



### 2.2.1 Continuous-time equations of motion with contacts

First, let us discretize our contact surface as a finite number  $n$  of contact points ( $\mathbf{c}_i$ ), at which the contact normal is well-defined. Generating this set is a non trivial task, and the chosen process will largely influence the numerical resolution of the system's dynamics; however, we consider it as out of the scope of this thesis. Let  $\mathbf{u}_i$  denote the relative velocity at the  $i^{\text{th}}$  contact point. If the  $i^{\text{th}}$  contact is between two points of our mechanical system (*self* contact), then there exist  $\mathbf{s}_{i,1}$  and  $\mathbf{s}_{i,2}$  such that  $\mathbf{u}_i(t, \mathbf{q}, \mathbf{v}) = \mathbf{v}(t, \mathbf{s}_{i,1}, \mathbf{q}, \mathbf{v}) - \mathbf{v}(t, \mathbf{s}_{i,2}, \mathbf{q}, \mathbf{v})$ . Otherwise, the contact is between a point of our system and an external object (*external* contact), and  $\mathbf{u}_i(t, \mathbf{q}, \mathbf{v}) = \mathbf{v}(t, \mathbf{s}_{i,1}, \mathbf{q}, \mathbf{v}) - \mathbf{w}^{\text{ext}}(t)$ , where  $\mathbf{w}^{\text{ext}}(t)$  is the velocity of the external object at  $\mathbf{c}_i$ . Let  $\mathbf{u} = (\mathbf{u}_i)$ .

The contact force density  $\mathbf{g}^c$  can be then be expressed as a sum of  $n$  punctual forces ( $\mathbf{r}_i$ ),  $\mathbf{g}^c = \sum_i \mathbf{r}_i \delta(\mathbf{x} - \mathbf{c}_i)$ . We get

$$\begin{aligned} \mathbf{f}^c &= \sum_{i=1}^n \begin{cases} \left( \frac{\partial \mathbf{v}}{\partial \mathbf{v}}(t, \mathbf{s}_{i,1}, \mathbf{q}) \right)^\top \mathbf{r}_i & \text{if } i \text{ is a self contact} \\ \left( \frac{\partial \mathbf{v}}{\partial \mathbf{v}}(t, \mathbf{s}_{i,1}, \mathbf{q}) - \frac{\mathbf{v}}{\mathbf{v}}(t, \mathbf{s}_{i,2}, \mathbf{q}) \right)^\top \mathbf{r}_i & \text{if } i \text{ is an external contact} \end{cases} \\ &= \left( \frac{\partial \mathbf{u}}{\partial \mathbf{v}}(t, \mathbf{q}) \right)^\top \mathbf{r}. \end{aligned}$$

We will assume that  $\mathbf{u}$  is an affine function of  $\mathbf{v}$ ; that is,  $\mathbf{u}(t, \mathbf{q}) := H(t, \mathbf{q})\mathbf{v} + \mathbf{w}(t)$  with  $H(t, \mathbf{q}) = \frac{\partial \mathbf{u}}{\partial \mathbf{v}}(t, \mathbf{q})$ . Our equations of motion with unilateral contact can then be summarized as

$$\begin{cases} M(t, \mathbf{q}) \frac{d\mathbf{v}}{dt} = \mathbf{f}(t, \mathbf{q}, \mathbf{v}) + H(t, \mathbf{q})^\top \mathbf{r} \\ \mathbf{u} = H(t, \mathbf{q})\mathbf{v} + \mathbf{w}(t) \\ \mathbf{u}_i, \mathbf{r}_i \text{ satisfy the contact law for } 1 \leq i \leq n. \end{cases} \quad (2.10)$$

First, we can treat the case of frictionless contacts, or more generally, an associated contact law with the convex cone  $K$  as the set of admissible forces. The contact law is then  $K^* \ni \mathbf{u}_i \perp \mathbf{r}_i \in K$ , and we can rewrite system (2.10) as the differential inclusion

$$M(t, \mathbf{q}) \frac{d\mathbf{v}}{dt} - \mathbf{f}(t, \mathbf{q}, \mathbf{v}) \in -H(t, \mathbf{q})^\top \mathcal{N}_{K^*}(H_i(\mathbf{q})\mathbf{u}_i + \mathbf{w}_i(t)).$$

Under a regularity assumption, for instance,  $\exists \mathbf{v}$  such that  $\mathbf{u}(t, \mathbf{q}, \mathbf{v}) \in \text{int}(K^*)^n$ , we can use Property A.12 on the subdifferential of the precomposition with an affine map to write equivalently our dynamics equation with unilateral contact as

$$M(t, \mathbf{q}) \frac{d\mathbf{v}}{dt} - \mathbf{f}(t, \mathbf{q}, \mathbf{v}) \in -\mathcal{N}_{C(\mathbf{q}, t)}(\mathbf{v}), \quad (2.11)$$

with the admissible velocity set  $V(\mathbf{q}, t) := \{\mathbf{v}, H(t, \mathbf{q})\mathbf{v} + \mathbf{w}(t) \in (K^*)^n\}$ . The associated contact law can therefore be modeled with the addition of a dissipation potential  $\mathcal{D}(t, \mathbf{q}, \mathbf{v}) = \mathcal{J}_{V(\mathbf{q}, t)}(\mathbf{v})$  to the equations of motion.

This form of the equations highlights a difficulty in the definition of  $\frac{d\mathbf{v}}{dt}$ : as the right-hand-side of (2.11) may be unbounded and non-differentiable w.r.t.  $\mathbf{v}$ ,  $\mathbf{q}$ , or  $t$ ,  $\mathbf{v}$  might be discontinuous. This is physically intuitive: considering the inelastic impact of a body falling on the ground at time  $t$ , its vertical velocity goes instantaneously from a non-zero value at  $t_-$  to zero at  $t_+$ . The symbol  $\frac{d\mathbf{v}}{dt}$  should instead be understood as a function defined almost everywhere on a time interval  $[T_0, T_1]$ , such that

$$\mathbf{v}(t) = \mathbf{v}(T_0) + \int_{T_0}^t \frac{d\mathbf{v}}{dt} dt \quad \forall t \in [T_0, T_1].$$

Equation (2.11) can then be cast as a *measure differential inclusion*. We will not worry here about the existence of solutions in time to such inclusions, as we will promptly resort to a time-stepping scheme which alleviates most concerns (see e.g., Moreau 1999).

Going back to the Coulomb frictional contact law, our equations of motions read

$$\begin{cases} M(t, \mathbf{q}) \frac{d\mathbf{v}}{dt} = \mathbf{f}(t, \mathbf{q}, \mathbf{v}) + H(t, \mathbf{q})^\top \mathbf{r} \\ \mathbf{u} = H(t, \mathbf{q}) \mathbf{v} + \mathbf{w}(t) \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} \end{cases} \quad \forall 1 \leq i \leq n, \quad (2.12)$$

where  $\frac{d\mathbf{v}}{dt}$  is still to be understood in the sense mentioned above. The Coulomb law does not derive from a dissipation potential, but the equation of motions can still be written once again as a differential inclusion,

$$M(t, \mathbf{q}) \frac{d\mathbf{v}}{dt} - \mathbf{f}(t, \mathbf{q}, \mathbf{v}) \in -H(t, \mathbf{q})^\top \mathcal{N}_{\left(\Pi_{\mathcal{K}_{\frac{1}{\mu_i}}}\right)}(\tilde{\mathbf{u}}), \quad (2.13)$$

where  $\tilde{\mathbf{u}} = (\tilde{\mathbf{u}}_i) = (\mathbf{u}_i + \mu_i \|(\mathbf{u}_i)_\top\| \mathbf{n}_i)$  is the contact-wise de Saxcé change of variable defined in Property 1.4.

### 2.2.2 Time integration

We now consider the problem of finding solution in time to our equations of motion with unilateral contacts. Several algorithms have been developed; we again refer to (Haddouni 2015, chapter 2) for a complete review on this topic.

We can cite two large class of methods. First, event-driven schemes decouple the nonsmooth events (such as impacts), which happen at a countable set of instants  $(t_i)$ , from the smooth dynamics which hold in every interval between those instants. The main advantage of this approach is the ability to use arbitrarily-high order integration schemes between the nonsmooth events, reaching tight error tolerances at a reasonable computational cost. However, there are also drawbacks:

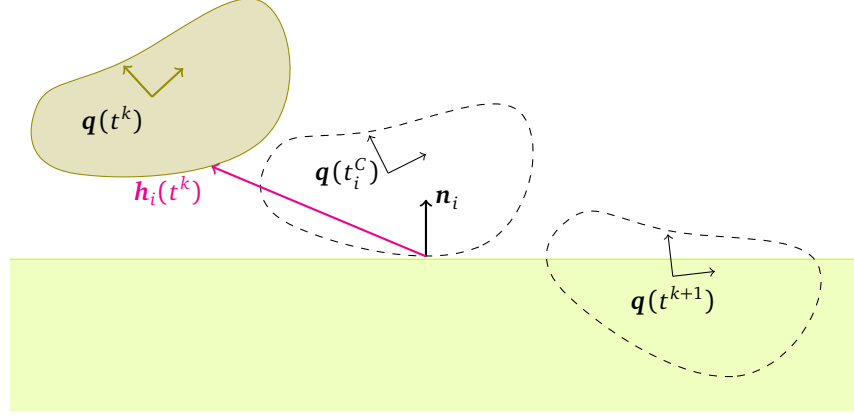
- A infinite number of events may occur in a finite time interval; this is known as a *Zeno* phenomenon, and is observed in mechanical systems as simple as rigid ball bouncing on the ground with a restitution coefficient in  $]0, 1[$  (Brogliato 1999);
- Simultaneous nonsmooth events require special care;
- Numerically, detecting events requires the introduction of several thresholds — saying that for instance, two objects will no longer be considered in contact when the gap between them exceeds a certain value. Tuning this threshold parameters is difficult to automate.

A second kind of integration method is usually preferred in practice for systems in which a large number of events may occur, such as large multi-body systems : *time-stepping* methods, which integrate the equation of motion over a time interval  $[t^k, t^{k+1} = t^k + \Delta_t]$  regardless of whether or not nonsmooth events occur during the timestep.

Time-stepping methods suppose prior knowledge of the contacts that are going to be active during the timestep; in practice, this means using proximity-based or continuous-time collision detection methods, which we will discuss in the next section. Then, the Moreau–Jean (1987, Moreau 1988) time-stepping algorithm discretizes the continuous-time measure differential equation (2.13) by enforcing that the end-of-step velocity  $\mathbf{v}^{k+1}$  — or the corresponding relative velocity,  $\mathbf{u}^{k+1}$  — satisfies the contact constraints.

**Velocity-level constraints** Consider the  $i^{\text{th}}$  contact predicted to happen during a timestep  $[t^k, t^k + \Delta_t]$ , and let  $h_i(t)$  denote the gap function, that is, the difference in position at time  $t$  of the two points that are going to come into contact at  $t_i^c \in [t^k, t^{k+1}]$ , as illustrated in Figure 2.5.

For a large  $\Delta_t$ ,  $h_{\text{IN}}(t^k)$  might be far from zero, that is, the two points of the future contact may still be far away from each other. However, the velocity constraint that we formulated in the continuous case reads  $\mathbf{u}_{\text{IN}}^{k+1} \geq 0$ , which would prevent the two points to get any closer.



**Figure 2.5:** Gap function  $h_i$  for a contact between a body with coordinates  $\mathbf{q}(t)$  and a fixed obstacle predicted to happen at time  $t_i^C \in ]t^k, t^{k+1}]$ .

Instead, the Moreau-Jean time-stepping algorithm formulates the Signorini constraint as follows:

$$\begin{cases} \mathbf{h}_i(t^{k+1})_{\text{IN}} = 0 \text{ and } 0 \leq \mathbf{u}_{\text{IN}}^{k+1} \perp \mathbf{r}_{\text{IN}} \geq 0 \\ \text{or } \mathbf{h}_i(t^{k+1})_{\text{IN}} > 0 \text{ and } \mathbf{r}_{\text{IN}} = 0, \end{cases}$$

where  $\mathbf{r}$  is the average of the contact forces over the timestep. In order to avoid introducing a new set of constraints, a modified Signorini condition,  $0 \leq \mathbf{h}_i(t^{k+1}) \perp \mathbf{r}_{\text{IN}} \geq 0$ , is commonly enforced. Using a first-order integration scheme, it is usual to write  $\mathbf{q}^{k+1} - \mathbf{q}^k = \Delta_t(\alpha \mathbf{v}^{k+1} + \alpha \mathbf{v}^k)$  with  $\alpha \in [0, 1]$ ;  $\mathbf{h}_i(t^{k+1})$  can then be approximated to the first order as  $\mathbf{h}_i(t^{k+1}) = \mathbf{h}_i(t^k) + \Delta_t \mathbf{u}_i^\alpha$ , with  $\mathbf{u}_i^\alpha = \alpha \mathbf{u}^{k+1} + (1 - \alpha) \mathbf{u}^k$ . The modified Signorini constraint written at the velocity level then reads

$$0 \leq \mathbf{u}_{\text{IN}}^{k+1} + \frac{1 - \alpha}{\alpha} \mathbf{u}_{\text{IN}}^k + \frac{1}{\alpha \Delta_t} \mathbf{h}_{\text{IN}}^k \perp \mathbf{r}_{\text{IN}} \geq 0. \quad (2.14)$$

This restricts our choice for  $\alpha$ ; indeed, we can see that for  $\alpha < 1$ , we may have  $\frac{1 - \alpha}{\alpha} \mathbf{u}_{\text{IN}}^k + \frac{1}{\alpha \Delta_t} \mathbf{h}_{\text{IN}}^k < 0$ , and therefore the modified Signorini condition (2.14) will not prevent having  $\mathbf{u}_{\text{IN}}^{k+1} > 0$  simultaneously with  $\mathbf{r}_{\text{IN}} > 0$ . In other terms, we may observe a rebound, which is contrary to our inelastic impact assumption. We will therefore restrict ourselves to  $\alpha = 1$ , and the end-of-step positions will be given by  $\mathbf{q}^{k+1} = \mathbf{q}^k + \Delta_t \mathbf{v}^{k+1}$ .

Accounting for the relative motion of the two points before the impact with the modified Signorini condition (2.14) will then just mean adding an offset  $\frac{1}{\Delta_t} \mathbf{h}(t^k)$  to the affine term  $\mathbf{w}_i(t)$  of the relative velocity  $\mathbf{u}_i$ .

**Discrete-time equations of motion** The Moreau-Jean algorithm approximates the smooth dynamics using a first-order  $\vartheta$ -scheme, and the acceleration using finite differences. For  $\vartheta \in [0, 1]$ , we note

$$\begin{aligned} t^\vartheta &= \vartheta t^{k+1} + (1 - \vartheta) t^k \\ \mathbf{v}^\vartheta &= \vartheta \mathbf{v}^{k+1} + (1 - \vartheta) \mathbf{v}^k \\ \mathbf{q}^\vartheta &= \vartheta \mathbf{q}^{k+1} + (1 - \vartheta) \mathbf{q}^k \end{aligned}$$

The differential inclusion (2.13) is then integrated over the timestep as <sup>3</sup>

$$M(t^\vartheta, \mathbf{q}^\vartheta)(\mathbf{v}^{k+1} - \mathbf{v}^k) - \mathbf{f}(t^\vartheta, \mathbf{q}^\vartheta, \mathbf{v}^\vartheta) = \Delta_t H(t^\vartheta, \mathbf{q}^\vartheta)^\top \mathbf{r}.$$

<sup>3</sup>As mentioned by Cadoux (2009),  $\mathbf{f}$  might depend on  $\mathbf{v}$  and may thus be nonsmooth; the end-of-step velocity may be considered in this case.

Enforcing holonomic constraints to be satisfied at  $t^\vartheta$ , our discrete equations of motion read

$$\left\{ \begin{array}{l} M(t^\vartheta, \mathbf{q}^\vartheta) \frac{\mathbf{v}^{k+1} - \mathbf{v}^k}{\Delta_t} = \mathbf{f}(t^\vartheta, \mathbf{q}^\vartheta, \mathbf{v}^\vartheta) + H(t^\vartheta, \mathbf{q}^\vartheta)^\top \mathbf{r} + \left( \frac{\partial C}{\partial \mathbf{q}}(t^\vartheta, \mathbf{q}^\vartheta) \right)^\top \boldsymbol{\lambda} \\ C(t^\vartheta, \mathbf{q}^\vartheta) = \mathbf{0} \\ \mathbf{u} = H(t^\vartheta, \mathbf{q}^\vartheta) \mathbf{v}^{k+1} + \mathbf{w}(t^\vartheta) + \frac{1}{\Delta_t} \mathbf{h}(t^k) \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} \quad \forall 1 \leq i \leq n. \end{array} \right. \quad (2.15)$$

**Linearized dynamics** Solving the implicit nonlinear system (2.15) can then be decomposed as a sequence of problems with explicit mass matrix, internal forces and constraints Jacobians using a quasi-Newton algorithm; see, e.g., (Jean 1999; Kaufman, Tamstorf, et al. 2014).

Indeed, let the sequence  $(\mathbf{v}^{k,l}, \mathbf{q}^{k,l})$  be recursively defined as the solution of the problem (2.16),

$$\left\{ \begin{array}{l} \mathbf{q}^{k,l+1} = \mathbf{q}^k + \Delta_t \mathbf{v}^{k,l+1} \\ \tilde{M}^l \frac{\mathbf{v}^{k,l+1} - \mathbf{v}^{k,l}}{\Delta_t} = \mathbf{f}(t^\vartheta, \mathbf{q}^{\vartheta,l}, \mathbf{v}^{\vartheta,l}) - M^l \frac{\mathbf{v}^{k,l} - \mathbf{v}^k}{\Delta_t} + H^l \mathbf{r} + G^l \boldsymbol{\lambda} \\ \tilde{M}^l := M^l - \vartheta \Delta_t \frac{\partial \mathbf{f}}{\partial \mathbf{v}}(t^\vartheta, \mathbf{q}^{\vartheta,l}, \mathbf{v}^{\vartheta,l}) - \vartheta \Delta_t^2 \frac{\partial \mathbf{f}}{\partial \mathbf{q}}(t^\vartheta, \mathbf{q}^{\vartheta,l}, \mathbf{v}^{\vartheta,l}) \\ G^l \mathbf{v}^{k,l+1} = G^l \mathbf{v}^{k,l} - C(\mathbf{q}^{\vartheta,l}) \\ \mathbf{u} = H^l \mathbf{v}^{k,l+1} + \mathbf{w}(t^\vartheta) + \frac{1}{\Delta_t} \mathbf{h}(t^k) \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} \quad \forall 1 \leq i \leq n, \end{array} \right. \quad (2.16)$$

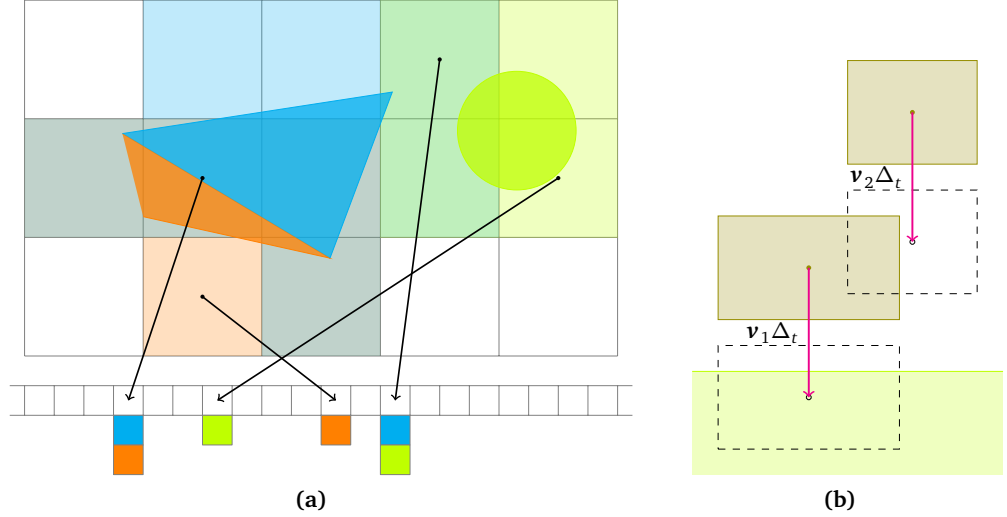
with  $M^l := M(t^\vartheta, \mathbf{q}^{\vartheta,l})$ ,  $H^l := H(t^\vartheta, \mathbf{q}^{\vartheta,l})$ ,  $G^l := \partial C \mathbf{q}(t^\vartheta, \mathbf{q}^{\vartheta,l})$ ,  $\mathbf{q}^{\vartheta,l} := \vartheta \mathbf{q}^{k,l} + (1 - \vartheta) \mathbf{q}^k$ ,  $\mathbf{v}^{\vartheta,l} := \vartheta \mathbf{v}^{k,l} + (1 - \vartheta) \mathbf{v}^k$ , and  $(\mathbf{v}^{k,0}, \mathbf{q}^{k,0}) = (\mathbf{v}^k, \mathbf{q}^k)$ . Then any value towards which this sequence converges will be a solution of the equations of motion (2.15). If  $\mathbf{f}$  derives from potentials, then  $\tilde{M}$  should remain symmetric; however  $\tilde{M}$  may non-remain positive semi-definite for large values of  $\Delta_t$ . Other approximations should then be used to estimate the Hessian, for instance the Gauss–Newton method for quadratic potentials.

In practice, this algorithm can be costly; even if the solution at iteration  $l$  is a good initial guess for the solution at  $l + 1$ , solving each subproblem may require expensive preprocessing steps. The algorithm is thus (2.16) often truncated to its first step, yielding the so-called *linearized* equations of motion. We will actually adopt this strategy for all of our applications. However, in the presence of highly nonlinear forces, this might worsen the stability or energy conservation properties of the time-stepping scheme (Kaufman, Tamstorf, et al. 2014).

**Limitations** The time-stepping approach also suffers from limitations. First, it is a low-order method, and thus requires small timesteps to approximate the continuous solution, even when no nonsmooth events occur. Note that Acary (2009) proposed a first method for improving the order of accuracy of the Moreau time-stepping scheme, which performed well on academic examples.

Then, in the above formulation, the contact constraints are formulated on the velocity, and may be subject to numerical drift. If  $\mathbf{h}_{\text{IN}}(t)$  becomes negative, then once again we may have  $\mathbf{u}_{\text{IN}}^{k+1} \mathbf{r}_{\text{IN}} > 0$ , and observe a rebound. On the other hand, if we clamp  $\mathbf{h}_{\text{IN}}(t)$  to positive values, the unilateral contact constraints may be violated.

Finally, as already mentioned, an important limitation of the time-stepping algorithm is that all contacts that may occur during the timestep must be predicted, which is hard problem, as we are going to see in the next section.



**Figure 2.6:** (a) Spatial hashing of two objects. Only simplices located in cells containing more than one object will have to be checked for collisions during the fine pass. (b) Failure case for continuous-time collision detection: the contact between the two falling bodies can only be detected once the contact between the green and the middle body has been resolved.

### 2.2.3 Collision detection

In this section, we briefly discuss a few concepts related to the detection of possible collisions within one timestep. For both a theoretical overview and practical implementation of fundamental algorithms, we refer the reader to Ericson (2004).

On the one hand, for slow enough relative velocities, or small enough timesteps, a very simple and efficient way of performing collision detection is to use a proximity criterion; two points will be considered in contact if the distance between them is below a detection radius  $r$ . As it would not be computationally tractable to compute the distance between every pair of points of the system, several approaches have been developed to speed-up the computation. Generally, they amount to a coarse pass, yielding a conservative list of possible collisions, then a finer pass checking the real distance between all the possible pairs. One such simple yet effective acceleration structure is the spatial hashing of Teschner et al. (2003). First, each simulated object is discretized as a set of simplices ( $S_i$ ), such as balls, segments or triangles. The space  $\mathbb{R}^3$  is then partitioned as a regular grid, and to each grid cell  $c$  is associated an index  $\text{hash}(c) \in \mathbb{N}$ . This mapping is used to define a hash-map  $\mathcal{H}$  associating to each cell the set of simplices that are within a radius  $r$  of this cell:  $\mathcal{H}[\text{hash}(c)] = \{i \in \mathbb{N}, d(S_i, c) < r\}$ . The fine pass simply amounts to checking for collisions between the simplices of each of the non-empty hash-map cells; that is, computing  $d(S_i, S_j) \forall (i, j) \in \mathcal{H}(k), i \neq j, \forall k \in \mathbb{N}$ , and for each pair such that  $d(S_i, S_j) < r$ , computing the closest point on each simplex and deducing the collision normal. Compared to representing the grid as a dense array, this approach has the advantage of a much lower memory footprint when the simulated objects occupy only sparse portions of the domain. However, the grid cell size has to be chosen carefully; if it is too coarse or too fine, the number of distance computations will increase.

On the other hand, for fast-moving objects or large timesteps, proximity-based collision detection will miss some collisions. In this case, the preferred approach is continuous-time collision detection, i.e., checking for collisions all along the trajectories of the simulated objects; obviously, this will require more complex and expensive methods. However, for several pairs of simple geometric shapes, collisions may be found by looking at the roots of a low-order polynomial. For instance, collisions between triangular meshes may be of two kinds: vertex–face collisions and

edge–edge collisions. Assuming a constant rigid motion for each triangle over the timestep, both correspond to the roots of a cubic polynomial. Brochu et al. (2012) proposed a geometrically-exact detection method for such intersections, avoiding the rounding errors associated with the polynomial root-finding approach.

Yet, in order to compute the collisions along a trajectory, we first need to know the trajectory. Which we don't, since the trajectory will be influenced by the contact forces, and obviously, to compute the contact forces, we first need to know the contact points. We therefore need to compromise, and will use unconstrained velocity (the solution of the unconstrained equations of motion) to define the trajectory that will be used for continuous-time collision detection. However, some contacts may not be predicted before the resolution of other collisions. For instance, consider two thin bodies A and B, with B above A, falling on the ground G at the same speed. If  $\Delta_t$  is big enough, they might impact the ground during the same time step; the continuous-time collision detection process will then predict the contacts A/G and B/G, but not the contact A/B, and the two bodies may end up in the same place. Such situations are likely to happen in the simulation of layered cloth, and may require sophisticated heuristics.

Another solution, if the computational cost is deemed acceptable, would be to perform a collision detection step at each iteration of the algorithm (2.16).

### 2.3 Discrete Coulomb Friction Problem

Let us refer to the system of equations that we have to solve at each step of the algorithm (2.16) as a *Discrete Coulomb Friction Problem* (DCFP), and let us write it using a lighter notation as

$$\begin{cases} M\mathbf{v} = \mathbf{f} + C^T\boldsymbol{\lambda} + H^T\mathbf{r} \\ C\mathbf{v} = \mathbf{k} \\ \mathbf{u} = H\mathbf{v} + \mathbf{w} \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} \end{cases} \quad \forall 1 \leq i \leq n, \quad (2.17)$$

where the unknowns are  $\mathbf{v} \in \mathbb{R}^m$ ,  $\mathbf{u} \in \mathbb{R}^{nd}$ ,  $\boldsymbol{\lambda} \in \mathbb{R}^c$  and  $\mathbf{r} \in \mathbb{R}^{nd}$ . Moreover,  $M$  will be assumed to be symmetric, positive definite, and the rows of  $C$  will generally be assumed to be linearly independent (that is, unless otherwise mentioned there are no redundant equality constraints).

#### 2.3.1 Reduced formulation

We can eliminate the variable  $\mathbf{v}$  from the DCFP (2.17), and obtain

$$\begin{cases} P\boldsymbol{\lambda} + B^T\mathbf{r} + \mathbf{c} = \mathbf{0} \\ B\boldsymbol{\lambda} + W\mathbf{r} + \mathbf{b} = \mathbf{u} \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} \end{cases} \quad \forall 1 \leq i \leq n. \quad (2.18)$$

with

$$\begin{aligned} W &:= HM^{-1}H^T & P &:= CM^{-1}C^T & B &:= HM^{-1}C^T \\ \mathbf{b} &:= HM^{-1}\mathbf{f} + \mathbf{w} & \mathbf{c} &:= CM^{-1}\mathbf{f} - \mathbf{k}. \end{aligned}$$

The symmetric, positive semi-definite matrix  $W$  is often referred to as the *Delassus* operator.

When the rows of  $C$  are linearly independent,  $C$  is surjective, therefore  $P$  is invertible and we can go one step further, writing system (2.18) with  $\mathbf{r}$  and  $\mathbf{u}$  as the only unknowns,

$$\begin{cases} \mathbf{u} = \tilde{W}\mathbf{r} + \tilde{\mathbf{b}} \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} \end{cases} \quad \forall 1 \leq i \leq n. \quad (2.19)$$

with  $\tilde{W} := W - BP^{-1}B^T$  and  $\tilde{\mathbf{b}} := \mathbf{b} - BP^{-1}\mathbf{c}$ .

While this formulation is short and convenient, in practice we will often avoid computing the  $\tilde{W}$  matrix explicitly. First, the inverse of  $M$  may be dense; this will be the case for cloth, or finite-element models, among others. In this case,  $W$  and  $P$  will be dense, and therefore expensive to

compute and less likely to fit into main memory. It is also rarely a good idea to compute  $\tilde{W}$ , as  $P^{-1}$  will often be dense and the double-inversion might lead to tremendous numerical errors.

However, for multi-body systems,  $M$  is a block-diagonal matrix; if the size of the diagonal blocks is reasonable, they can be easily factorized, for instance using a  $QLDL^T Q^T$  factorization<sup>4</sup>, and it may become fruitful to explicitly compute the Delassus operator. Even then,  $W$  might be costly to assemble. We argue in (Daviet, Bertails-Descoubes, and Boissieux 2011) that it is important to exploit the block structure of  $M$  and  $H$  to accelerate the computation. Our bogus (Daviet 2013) C++ template library provides an easy way to perform efficient linear algebra operations with sparse block matrices, including a parallelized matrix–matrix multiplication routine. Despite this, our results for hair simulations (presented in the next chapter, Table 4.2) show that a significant portion of the total time for solving the DCFP (2.17) was spent assembling the Delassus matrix  $W$ , and (Daviet, Bertails-Descoubes, and Boissieux 2011, Figure 6) implies a scaling with the square of the number of contacts, both in time and memory consumption.

More generally, a trade-off has to be made between a higher preprocessing and memory cost, or a higher solver cost. Some of the algorithms that we will present in Chapter 3 will require an explicit computation of  $W$ , others will just require being able to perform matrix–vector multiplications with  $W$ , some will use completely different formulations; their relative performance is application-dependent.

### 2.3.2 Fixed-point algorithms and existence criterion

Still following from the fact that the Coulomb friction law does not derive from a potential, the DCFP (2.17) cannot be formulated as a convex optimization problem. However, (2.17) can be characterized as the fixed-point of sequences of such problems; we just need to temporarily fix either  $\mathbf{u}_i$  or  $\mathbf{r}_i$  in the coupling term of the bipotential (1.14),  $\mu_i \mathbf{r}_{iN} \|\mathbf{u}_{iT}\|$ .

As the minimization of convex functions is a well-studied problem, for which numerous efficient algorithms have been devised, this iterative strategy may prove fruitful in practice. Moreover, this characterization of (2.17) as a fixed-point was used by Cadoux (2009) to state a sufficient criterion for the existence of solutions, which we will recall in Theorem 2.1.

**Haslinger algorithm** First, let us fix the normal force  $\mathbf{r}_{iN}$ . Consider the function  $h_s$  (which is no-longer a bipotential) derived from (1.14) as

$$h_s(\mathbf{u}_i, -\mathbf{r}_i) = \mathcal{J}_{-\mathcal{K}_{\mu_i}}(-\mathbf{r}_i) + \mathcal{J}_{\mathcal{K}_{\infty}}(\mathbf{u}_i) + s_i \|\mathbf{u}_{iT}\|.$$

A velocity–force couple satisfying  $\mathbf{r}_i \in -\partial h_s(\mathbf{u}_i)$  will be a solution to the Coulomb friction problem if and only if  $s_i = \mu_i \mathbf{r}_{iN}$ . Now, we recognize from Equation (1.25) the expression of the convex conjugate of  $\mathcal{J}_{\mathcal{T}_{0,s_i}}$ , that this,  $\mathcal{J}_{\mathcal{K}_{\infty}}(\mathbf{u}_i) + s_i \|\mathbf{u}_{iT}\| = \mathcal{J}_{\mathcal{T}_{0,s_i}}^*(\mathbf{u}_i)$  (note that  $\mathcal{T}_{0,s_i}$  is simply a semi-infinite cylinder of radius  $s_i$ ). This means, using Theorem A.2,

$$\mathbf{r}_i \in -\partial h_s(\mathbf{u}_i) \iff \mathbf{r}_i \in -\partial \left( \mathcal{J}_{\mathcal{T}_{0,s_i}}^* \right)(\mathbf{u}_i) \iff \mathbf{u}_i \in -\mathcal{N}_{\mathcal{T}_{0,s_i}}(\mathbf{r}_i).$$

As the optimality conditions of the quadratic minimization problem under cylindrical constraints (2.20),

$$\min_{\lambda \in \mathbb{R}^c, \mathbf{r} \in \Pi_{\mathcal{T}_{0,s_i}}} \frac{1}{2} (\lambda^T, \mathbf{r}^T) \begin{pmatrix} P & B^T \\ B & W \end{pmatrix} \begin{pmatrix} \lambda \\ \mathbf{r} \end{pmatrix} + (\lambda^T, \mathbf{r}^T) \begin{pmatrix} \mathbf{c} \\ \mathbf{b} \end{pmatrix}, \quad (2.20)$$

are, from Theorem (A.6),

$$\begin{cases} P\lambda + B^T \mathbf{r} + \mathbf{c} = \mathbf{0} \\ B\lambda + W\mathbf{r} + \mathbf{b} = \mathbf{u} \\ \mathbf{u}_i \in -\mathcal{N}_{\mathcal{T}_{0,s_i}}(\mathbf{r}_i) \quad \forall i, \end{cases}$$

<sup>4</sup>That is, a Cholesky factorization where the square-root of the diagonal is not explicitly computed for better numerical precision, and with a reordering step to reduce the fill-in of the triangular matrix  $L$ . The `SimplicialLDLT` class of the Eigen (Guennebaud, Jacob, et al. 2010) C++ library provides an implementation of this factorization.



we deduce that the solution of the minimization problem (2.20) will be the solution of the dual form of the DCFP (2.18) if and only if  $\mathbf{s} := (s_i) = (\mu_i \mathbf{r}_{iN})$ . This leads to the fixed-point algorithm introduced by Haslinger (1983), which consists in iteratively solving the minimization problem (2.20)<sup>5</sup> for a given parameter value  $\mathbf{s}^k$ , then computing  $\mathbf{s}^{k+1}$  from the optimal solution, and looping again.

However, the objective function of (2.20) is not strictly convex, and therefore the problem may have more than one solution. The fixed-point iteration is thus ill-defined, and studying its convergence would be tricky. In contrast, Cadoux (2009) formulated an algorithm that does not suffer from this drawback, by freezing instead the velocity part of the bipotential coupling term.

**Cadoux algorithm** Let us now consider the function derived from the bipotential by fixing  $\mathbf{u}$  in the coupling term,

$$c_s(\mathbf{u}_i, -\mathbf{r}_i) := \mathcal{G}_{-\mathcal{K}_{\mu_i}}(-\mathbf{r}_i) + \mathcal{G}_{\mathcal{K}_{\infty}}(\mathbf{u}_i) + s_i \mathbf{r}_{iN}.$$

The equivalence  $(\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} \iff \mathbf{u}_i \in \frac{\partial c_s}{\partial -\mathbf{r}_i}(-\mathbf{r}_i)$  will hold if and only if  $s_i = \mu_i \|\mathbf{u}_{iT}\|$ . Now, direct computations yield that  $\mathbf{u}_i \in \frac{\partial c_s}{\partial -\mathbf{r}_i}(-\mathbf{r}_i) \iff \mathbf{u}_i + s_i \mathbf{n}_i \in -\mathcal{N}_{\mathcal{K}_{\mu_i}}(\mathbf{r}_i)$ . The solution of the DCFP (2.18) will therefore be given by the solution of the minimization problem (2.21),

$$q(\mathbf{s}) := \min_{\lambda \in \mathbb{R}^c, \mathbf{r} \in \Pi_{\mathcal{K}_{\mu_i}}} \frac{1}{2} (\lambda^\top, \mathbf{r}^\top) \begin{pmatrix} P & B^\top \\ B & W \end{pmatrix} \begin{pmatrix} \lambda \\ \mathbf{r} \end{pmatrix} + (\lambda^\top, \mathbf{r}^\top) \begin{pmatrix} \mathbf{c} \\ \mathbf{b} + \mathbf{s} \end{pmatrix}, \quad (2.21)$$

when  $\mathbf{s} = (\mu_i \|\mathbf{u}_{iT}\| \mathbf{n}_i)$ , with  $\mathbf{u} := W\mathbf{r} + B\lambda + \mathbf{b}$

We will refer to problems such as (2.21), where a quadratic objective function is minimized subject to SOC constraints, as Second-Order Cone Quadratic Programs (SOCQP). A fixed-point algorithm can then be obtained by computing  $\mathbf{s}^{k+1}$  from the solution  $\lambda^k$  to the problem parameterized by  $\mathbf{s}^k$  as  $\mathbf{s}^{k+1} = (\mu_i \|\mathbf{u}_{iT}^k\|)$ .

At first sight, the gain over the Haslinger algorithm is not obvious; the solution to the minimization problem (2.21) is still not unique. Cadoux (2009, Theorem 3.5) proves the well-definition of the fixed-point update from a duality argument, using Fenchel's duality theorem (A.3) to show that the velocity  $\mathbf{v}$  corresponding to the optimum of (2.21) is the solution of a strictly convex minimization problem, meaning that  $\mathbf{v}$  and thus  $\mathbf{u}$  are uniquely defined for each  $\mathbf{s}$ . Here, we will show this well-definition by directly constructing the primal problem from the optimality condition of the dual (2.21).

Expressing  $W$  and  $P$  in the optimality conditions of (2.21), we can get back to a system in the same unknowns as our original DCFP (2.17),

$$\begin{cases} M\mathbf{v} = \mathbf{f} + C^\top \lambda + H^\top \mathbf{r} \\ C\mathbf{v} = \mathbf{k} \\ \mathbf{u} = H\mathbf{v} + \mathbf{w} \\ \mathcal{K}_{\frac{1}{\mu_i}} \ni \mathbf{u}_i + s_i \mathbf{n}_i \perp \mathbf{r}_i \in \mathcal{K}_{\mu_i} \quad \forall i = 1 \dots n. \end{cases} \quad (2.22)$$

Since system (2.22) derives from the optimality conditions of a convex problem, we will refer to (2.22) as a *convexified* DCFP. Now, discarding the relative velocity variable  $\mathbf{u}$ , system (2.22) can be written equivalently as

$$\begin{cases} M\mathbf{v} = \mathbf{f} + C^\top \lambda + H^\top \mathbf{r} \\ \lambda \in -\mathcal{N}_{\{\mathbf{0}_{\mathbb{R}^c}\}}(C\mathbf{v} - \mathbf{k}) \\ \mathbf{r} \in -\mathcal{N}_{\Pi_{\frac{1}{\mu_i}}} (H\mathbf{v} + \mathbf{w} + \mathbf{s}). \end{cases}$$

<sup>5</sup>Actually, the original Haslinger algorithm solves the Fenchel primal of (2.20).



Moreover, from Corollary A.4 to the Property A.12 on the subdifferential of a precomposition with an affine map, we have always

$$C^\top \mathcal{N}_{\{0_{\mathbb{R}^c}\}}(C\mathbf{v} - \mathbf{k}) + H^\top \mathcal{N}_{\Pi_{\mu_i} \mathcal{K}_\perp}(H\mathbf{v} + \mathbf{w} + \mathbf{s}) \subset \mathcal{N}_{V(\mathbf{s})}(\mathbf{v})$$

with  $V(\mathbf{s}) := \left\{ \mathbf{v} \in \mathbb{R}^m, C\mathbf{v} = \mathbf{k} \text{ and } (H\mathbf{v} + \mathbf{w} + \mathbf{s}) \in \Pi_{\mu_i} \mathcal{K}_\perp \right\}$ , and equality is achieved under the regularity condition  $(\mathcal{H}(\mathbf{s}))$ ,

$$\exists \mathbf{v} \in \mathbb{R}^m, C\mathbf{v} = \mathbf{k} \text{ and } (H\mathbf{v} + \mathbf{w} + \mathbf{s}) \in \text{int } \Pi_{\mu_i} \mathcal{K}_\perp. \quad (\mathcal{H}(\mathbf{s}))$$

This means that all the solutions to (2.22) must satisfy (2.23),

$$M\mathbf{v} - \mathbf{f} \in -\mathcal{N}_{V(\mathbf{s})}(\mathbf{v}) \quad (2.23)$$

and that they coincide when the condition  $(\mathcal{H}(\mathbf{s}))$  is satisfied. Finally, we recognize (2.23) as the optimality conditions of the strictly convex minimization problem (2.24),

$$p(\mathbf{s}) := \min_{\mathbf{v} \in V(\mathbf{s})} \frac{1}{2} \mathbf{v}^\top M \mathbf{v} - \mathbf{v}^\top \mathbf{f}. \quad (2.24)$$

Summarizing, the velocity  $\mathbf{v}$  reconstructed from a solution  $(\boldsymbol{\lambda}, \mathbf{r})$  of the dual SOCQP (2.21) as  $\mathbf{v} = M^{-1}(\mathbf{f} + C^\top \boldsymbol{\lambda} + H^\top \mathbf{r})$  will always be a solution of the primal SOCQP (2.24), and if  $\mathbf{v}$  is a solution to the primal (2.24), then the condition  $(\mathcal{H}(\mathbf{s}))$  suffices to ensure the existence of a solution  $(\boldsymbol{\lambda}, \mathbf{r})$  to (2.21).

**Property 2.1** (Cadoux fixed-point algorithm). *We introduce the mapping  $s : \mathbb{R}^{nd} \rightarrow \mathbb{R}^{nd}$ ,  $(\mathbf{u}_i) \mapsto (\mu_i \|\mathbf{u}_{iT}\| \mathbf{n}_i)$ . Let  $v(\mathbf{s})$  denote the optimum of the strictly-convex primal SOCQP (2.24), and  $u : \mathbf{v} \mapsto \mathbf{u} := H\mathbf{v} + \mathbf{w}$ . Then, if  $V(\mathbf{0}) \neq \emptyset$ , the mapping  $F : \mathbb{R}^{nd} \rightarrow \mathbb{R}^{nd}$ ,  $F := s \circ u \circ v$  is well-defined on  $\text{Im}s \supset \text{Im}F$ , and under the condition  $\mathcal{H}(\mathbf{0})$ , any fixed-point of  $F$  will correspond to a solution of the DCFP (2.17).*

*The sequence  $\mathbf{s}^0 := s \circ u(\mathbf{v}^0)$ ,  $\mathbf{s}^{k+1} := F(\mathbf{s}^k)$  will be referred to as the Cadoux fixed point algorithm.*

*Proof.* Notice that  $\forall i \in 1 \dots n$ ,  $s(\mathbf{u})_i \in \mathcal{K}_0$ ; this means that for any  $\mathbf{s}$  in the image of  $s$ ,  $\forall \mathbf{z} \in \Pi_{\mu_i} \mathcal{K}_\perp$ ,  $(\mathbf{z} + \mathbf{s}) \in \mathcal{K}_\perp$  and  $\forall \mathbf{z} \in \text{int } \Pi_{\mu_i} \mathcal{K}_\perp$ ,  $(\mathbf{z} + \mathbf{s}) \in \text{int } \mathcal{K}_\perp$ . Therefore, for all  $\mathbf{s} \in \text{Im}s$ ,  $V(\mathbf{s}) \subset V(\mathbf{0})$ , and  $\mathcal{H}(\mathbf{0}) \implies \mathcal{H}(\mathbf{s})$ .

If  $V(\mathbf{0}) \neq \emptyset$ , the primal SOCQP (2.24) is thus feasible for any  $\mathbf{s} \in \text{Im}s$ , and as it is strictly-convex it admits a unique solution.  $v$  is therefore well-defined on  $\text{Im}s$ , and so is  $F$ .

Now, let  $\bar{\mathbf{s}}$  be a fixed-point of  $F$ . Under the condition  $\mathcal{H}(\mathbf{0})$ ,  $\mathcal{H}(\bar{\mathbf{s}})$  is satisfied and therefore  $v(\bar{\mathbf{s}})$  will correspond to a solution  $(\bar{\boldsymbol{\lambda}}, \bar{\mathbf{r}})$  of the dual SOCQP (2.21) at  $\bar{\mathbf{s}}$ . As  $\bar{\mathbf{s}} = F(\bar{\mathbf{s}}) = (\mu_i \|\mathbf{u}_{iT}\| \mathbf{n}_i)$ ,  $(\bar{\boldsymbol{\lambda}}, \bar{\mathbf{r}})$  is also a solution of the dual for of the DCFP (2.17).  $\square$

Cadoux (2009) also investigated other update rules for  $\mathbf{s}$ , for instance using a Newton or quasi-Newton iteration instead of a fixed-point one. Those more complex rules were not found to perform significantly better in practice, and we will not consider them here.

In practice  $v(\mathbf{s})$  may be computed either from the solution to the primal SOCQP (2.24), or from the dual (2.21). The primal has the advantage of boasting strict convexity, but is subject to a more complex constraint, as it involves the linear application  $H$  on top of the SOC. For some values of  $\mathbf{s}$ ,  $p(\mathbf{s})$  might admit a solution while  $q(\mathbf{s})$  does not; however, an approximate solution for  $q(\mathbf{s})$  might suffice to keep the algorithm going. Actually, it is often a good idea to truncate the resolution of the intermediate SOCQP, using heuristics to refine error tolerance as the fixed-point algorithm converges, to increase the computational performance of the algorithm as a whole.

Cadoux (2009) proves furthermore — in the case without equality constraints, but the generalization is easy — the following theorem about the existence of solutions:

**Theorem 2.1** (Cadoux existence criterion). *If  $V(\mathbf{0}) \neq \emptyset$ , then the mapping  $F$  introduced in Property 2.1 admits a fixed-point  $\bar{\mathbf{s}}$ . If  $\mathcal{H}(\bar{\mathbf{s}})$ , or a fortiori if  $\mathcal{H}(\mathbf{0})$  is satisfied, then the DCFP (2.17) admits a solution.*

Note that Theorem 2.1 gives a *sufficient* existence condition, but not a *necessary* one. For instance, Cadoux (2009, Theorem 3.20) states that a solution to DCFP (2.17) will exist as soon as the fixed-point  $\bar{\mathbf{s}}$  is such that the constraints for which  $\bar{\mathbf{s}}_i = \mathbf{0}$  satisfy a qualification condition.

**Anitescu regularization** Using the regularized contact law proposed by Anitescu (2005) amounts to approximating the DCFP (2.17) with the convexified DCFP

$$\left\{ \begin{array}{l} M\mathbf{v} = \mathbf{f} + C^\top \boldsymbol{\lambda} + H^\top \mathbf{r} \\ C\mathbf{v} = \mathbf{k} \\ \mathbf{u} = H\mathbf{v} + \mathbf{w} \\ \mathcal{K}_{\frac{1}{\mu_i}} \ni \mathbf{u}_i \perp \mathbf{r}_i \in \mathcal{K}_{\mu_i} \end{array} \right. \quad \forall i = 1 \dots n,$$

which corresponds to the optimality conditions of our couple of dual SOCQP at  $\mathbf{s} = \mathbf{0}$ . Using the Drucker–Prager analogy (Section 1.3.3), this corresponds to adding a dilatancy with coefficient  $\frac{1}{\mu}$  to the Coulomb friction law, i.e., in plasticity terms, using an associated flow rule. Mazhar et al. (2015) argue that numerical simulations using this regularization are still able to match experimental data on several complex examples. However we can readily see that the approximation becomes rougher for higher values of the friction coefficient or of the tangential relative velocity, and cause the objects to separate instead of sliding.

## Summary

In this chapter, we have seen that the equations of motion of a large class of mechanical systems subject to holonomic kinematic constraints and unilateral contact with Coulomb friction can be integrated in time by solving a sequence of problems structurally similar to the Discrete Coulomb Friction Problem (2.17). The DCFP itself can then be solved as a sequence of Second-Order Cone Quadratic Programs, using either the primal (2.24) or dual (SOCQP (2.21) formulation.

In the following chapter, we will present numerical algorithms for solving either:

- the original DCFP (2.17);
- the convexified DCFP (2.22), or the corresponding primal (2.24) and dual (2.21) SOCQP.



## 3 Solving the Discrete Coulomb Friction Problem

This chapter will present an overview of different classes of methods that have been proposed in the literature to solve Discrete Coulomb Friction Problems, in both their standard and convexified versions. We briefly recall the differences between these two problems:

- In the case of the standard DCFP (2.17), each pair of local relative velocity and contact force  $(\mathbf{u}_i, \mathbf{r}_i)$  must satisfy the Coulomb law, i.e.,  $(\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i}$ . This is equivalent to  $\mathcal{K}_{\frac{1}{\mu_i}} \ni \tilde{\mathbf{u}}_i \perp \mathbf{r}_i \in \mathcal{K}_{\mu_i}$ , where  $\tilde{\mathbf{u}}_i$  is defined by the *de Saxcé* change of variable from Property 1.4,  $\tilde{\mathbf{u}}_i := \mathbf{u}_i + \mu_i \|\mathbf{u}_{iT}\| \mathbf{n}$ .
- In the convexified DCFP (2.17), each pair of local relative velocity and contact force  $(\mathbf{u}_i, \mathbf{r}_i)$  must directly satisfy the complementarity relationship  $\mathcal{K}_{\frac{1}{\mu_i}} \ni \mathbf{u}_i \perp \mathbf{r}_i \in \mathcal{K}_{\mu_i}$ . The convexified DCFP is equivalent to a minimization problem, the dual SOCQP (2.21), and, under the regularity condition  $\mathcal{H}(\mathbf{0})$ , to the primal SOCQP (2.24).

We recall also that the standard DCFP can be recast as a sequence of convexified DCFP using the Cadoux (2009) fixed-point algorithm from Section 2.3.2.

There has been a considerable amount of research dedicated to solving DCFP in the last decades, and to the best of our knowledge no comprehensive review exists. Such is not the goal of this chapter either; what follows is an opinionated presentation of the methods that we found to be of practical or historical interest, and we refer the reader to (Acary and Brogliato 2008; Cadoux 2009; Heyn 2013) for a larger covering of the literature. In the next chapter, we will also present an algorithm designed to be robust to ill-conditioned problems, and that performed well on a few applications.

To lighten notations and when no confusion is possible, we will use  $\mathcal{K}_{\mu}$  ( $\mathcal{K}_{\frac{1}{\mu}}$ ) to designate indifferently the SOC in  $\mathbb{R}^d$  and  $\Pi_i \mathcal{K}_{\mu_i}$  ( $\Pi_i \mathcal{K}_{\frac{1}{\mu_i}}$ ), the product of  $n$  SOC in  $\mathbb{R}^{nd}$ , .

### 3.1 Global strategies

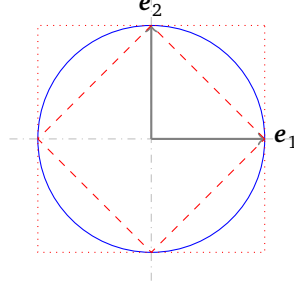
#### 3.1.1 Pyramidal friction cone

In 2D, the SOC can be described by two linear constraints. This is no longer the case in 3D, however approximating the 3D SOC with a pyramid (using a *facetting* process) is a popular way of writing the friction problem with only a set of linear constraints.

One can define an approximation of the friction cone as  $q$ -sided pyramid  $\mathcal{P}_{\mu}^{p,d}$ ,

$$\mathcal{P}_{\mu}^{p,d} := \left\{ \mathbf{r} \in \mathbb{R}^d, \left\| (\langle \mathbf{r}, \mathbf{e}_j \rangle)_{1 \leq j \leq q} \right\|_p \leq \mu \mathbf{r}_N \right\}, \quad \mathbf{e}_j := \frac{\mathbf{d}_j}{\left\| (\langle \mathbf{d}_j, \mathbf{d}_k \rangle)_{1 \leq k \leq q} \right\|_p}$$

where  $(\mathbf{d}_j)$  is a set of  $q$  unit vectors in the tangent plane, and  $p = 1$  (inner approximation) or  $+\infty$  (outer approximation). In practice, it is common to choose  $(\mathbf{d}_j)$  as  $d - 1$  orthogonal vectors. In order to minimize the effect of the linearization of the cone, it is advocated to choose the vectors in  $(\mathbf{d}_j)$  such that the faceted cone coincides with the exact one in the direction of the relative velocity at the beginning of the timestep.



**Figure 3.1:** Horizontal section of the inner (red, dashed) and outer (red, dotted) approximations with two basis vectors of the 3D Coulomb cone (blue).

The faceted problem can then be solved using a variety of methods; we will not detail them, as our main goal will be to solve the problem in 3D (or more) with the exact cone. We refer to (Erleben 2013) for a recent review such algorithms. In a non-exhaustive manner, we can mention:

- the Lemke algorithm used by Klarbring (1987), a direct method with good convergence properties, but which may become expensive for large systems;
- a conjugate-gradient algorithm on a sequence of minimization problems (Renouf and Alart 2005);
- the staggered projection approach (Kaufman, Sueda, et al. 2008), popular in Computer Graphics, which solves for the normal and tangential components in an iterative manner;
- an interior-point method (Trinkle et al. 1997);
- the Gauss-Seidel splitting algorithm (e.g., Raous et al. 1988, Jean 1999, Erleben 2007), which we will detail in Section (3.4).

This faceting approach also possesses weaknesses; for instance, Renouf, Acary, et al. (2005) argue that this approximation may lead to significant errors in the resulting trajectories.

Moreover, while the structure of the individual constraints is indeed simpler, their number is drastically increased (the simplest approximation, a 4-sided pyramid, already requires 4 linear constraints and 2 dual variables, instead of a single quadratic constraint and a single dual variable for the exact cone). The constraints also become more prone to switching between active and inactive states (the number of faces is increased), which might degrade the overall convergence of the solver algorithm. Overall, the faceted problem might therefore not be simpler nor faster to solve than the original one<sup>1</sup>, and the pertinence of this approximation should be carefully evaluated.

For these reasons, in the remainder of this chapter we will focus on the original, non-linearized problem defined with forces inside a SOC.

### 3.1.2 Complementarity functions

A natural way of solving the reduced dual formulation (2.19) is to express Coulomb law as a the root of a complementarity function, for instance the Alart–Curnier function (1.6) or the De Saxcé (1.29) complementarity functions. In the next chapter, we will also present another complementarity function, based on an adaptation of the Fischer–Burmeister function to the SOC algebra (Fukushima et al. 2002). Note that while the De Saxcé complementarity function can easily be adapted to solve the convexified DCFP rather than the DCFP (by simply not performing the eponymous change of variable), it is not as simple in the case of the Alart–Curnier function.

<sup>1</sup>Rockafellar (1993) famously said that “the great watershed in optimization isn’t between linearity and nonlinearity, but convexity and nonconvexity”.

In a general manner, let  $f : \mathbb{R}^{nd} \times \mathbb{R}^{nd} \rightarrow \mathbb{R}^{nd}$  be such that  $f(\mathbf{u}, \mathbf{r}) = \mathbf{r} \iff (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} \forall i$ . Then finding a solution to (2.19) amounts to finding  $\mathbf{r} \in \mathbb{R}^{nd}$  such that

$$f(\tilde{W}\mathbf{r} + \tilde{\mathbf{b}}, \mathbf{r}) = \mathbf{r}. \quad (3.1)$$

As it now involves solely the force variable  $\mathbf{r}$ , a variety of algorithms may be used to find a solution to (3.1); however, remember that  $f$  is in general non-convex, non-contracting, and non-smooth.

Using a simple fixed-point iteration  $\mathbf{r}^{k+1} = f(\tilde{W}\mathbf{r}^k + \tilde{\mathbf{b}}, \mathbf{r}^k)$  is tempting, and has been used in the case of a single contact (e.g., Feng et al. 2005)<sup>2</sup>, but the convergence is not always satisfying for larger systems in practice. As the points at which  $f$  is non-differentiable are generally sparse, a better solution is to use a Newton algorithm<sup>3</sup> on the function  $g$  (Alart and Curnier 1991; Dumont 2012),  $g : \mathbb{R}^{nd} \rightarrow \mathbb{R}^{nd}$ ,  $\mathbf{r} \mapsto f(\tilde{W}\mathbf{r} + \tilde{\mathbf{b}}, \mathbf{r}) - \mathbf{r}$ . Each step is defined as the solution to a linear system,

$$J_g^k(\mathbf{r}^{k+1} - \mathbf{r}^k) = -\alpha^k g(\mathbf{r}^k)$$

where  $J_g^k \in \partial g(\mathbf{r}^k)$  is any element of the subdifferential of  $g$  at  $\mathbf{r}^k$ , and  $\alpha^k$  is an optional damping term that can be computed using a line-search based on the error function  $\|g\|^2$ .

**Drawbacks** The left-hand-side matrix of the linear system that must be solved at each step,

$$\frac{\partial g}{\partial \mathbf{r}}(\mathbf{r}^k) = \tilde{W} \frac{\partial f}{\partial \mathbf{u}}(\tilde{W}\mathbf{r}^k + \tilde{\mathbf{b}}) + \frac{\partial f}{\partial \mathbf{r}}(\mathbf{r}^k) - \mathbb{I}_{\mathbb{R}^{nd}},$$

may not be invertible. A failsafe is required when one cannot find a step increment, for instance a gradient-descent or fixed-point iteration. In any case, the cost of solving this linear system grows quickly with the number of contacts, making the cost of this algorithm prohibitive for very large systems.

Moreover, we have observed in (Bertails-Descoubes et al. 2011) that while a Newton algorithm on the Alart-Curnier formulation yields good results for a system with a small number of contacts, the success rate of the algorithm drops quickly when a conditioning parameter,  $\mu \frac{nd}{m}$ , grows beyond 1. Indeed, in such cases the number of rows of  $H$  exceeds its number of columns, and the number of vanishing eigenvalues of  $\tilde{W}$  increases. Moreover, criterion  $\mathcal{H}(\mathbf{s})$  implies a link between the surjectivity of  $H$  and the existence of a solution to the DCFP; as the relative rank of  $H$  degrades, so does the probability that this existence criterion is satisfied.

### 3.1.3 Optimization-based methods

Since the convexified DCFP can be recast as either the primal (2.24) or dual (2.21) SOCQP, it is natural to attempt to use optimization-based algorithms to solve this problem. Moreover, as a SOCQP consists in the minimization of a convex function over a convex feasible set, theoretical convergence properties can generally be derived.

In the following sections, we will present three kinds of optimization-based methods: interior-point methods, proximal algorithms, and splitting algorithms. The latter family of methods can be interpreted as block-coordinate-descent algorithms, i.e., partial minimization over a small number of variables in turn.

Optimization-based methods can be leveraged to solve the original DCFP thanks to the Cadoux fixed-point algorithm defined in Property 2.1. Alternatively, we will see in the following that some optimization-based algorithms (e.g., projected gradient descent or Gauss–Seidel) can be repurposed to directly solve the original DCFP. However, both strategies induce the loss of any theoretical convergence property.

<sup>2</sup>When  $f$  is a projection function, i.e.,  $f \sim \Pi_C(\mathbf{r} - \rho \mathbf{u})$ , as is the case for the Alart-Curnier and De Saxcé complementarity function, the fixed-point algorithm can be interpreted as a slight variation of a proximal algorithm; see Feng et al. 2005 and Section 3.3 below.

<sup>3</sup>Quasi-Newton methods on complementarity functions are rarely observed to perform well in practice.

### 3.2 Interior-point methods

Interior-point methods have become extremely popular in optimization communities in the last decades, thanks to their superlinear convergence on many problems (Potra and Wright 2000).

At their core, they replace the characteristic function  $\mathcal{I}_C$  associated to the constraint  $C$  with a logarithmic *barrier* function,  $\mathcal{I}_{C,\beta}$ , such that  $\lim_{\beta \rightarrow 0} \mathcal{I}_{C,\beta} = \mathcal{I}_C$ . If the feasible set  $C$  is defined as  $C = \{\mathbf{x}, F(\mathbf{x}) \leq 0\}$ , one may take  $\mathcal{I}_{C,\beta}(\mathbf{x}) := -\beta \log(-F(\mathbf{x}))$ . The interior-point method then performs a Newton algorithm on the modified optimization problem, while simultaneously driving  $\beta$  to 0 in the process.

Note that high values for  $\beta$  will tend to make the successive iterates stay in the “center” of the feasible set. Iterates will thus follow an interior path to reach the solution, even when this optimal point lies on the boundary of the constraint, which is where the interior-point method’s name come from. This is in contrast to most other kinds of methods (pivoting, active-set, proximal, ...) which tend instead to follow the boundary of the constraint.

#### 3.2.1 Second-Order Cone Programs

A particular case for which several efficient solvers<sup>4</sup> have been devised is Second Order Cone Programs (SOCP), that is, minimizing a linear objective function under equality and SOC constraints. Our convexified DCFP corresponds to the optimality conditions of a (primal or dual) SOCQP, and this class of problem can easily be recast as SOCP, as shown in Cadoux (2009) (more generally, quadratic optimization problems with quadratic constraints are a subclass of SOCP, see e.g., E. Andersen and K. Andersen 2013). The trick is in the transformation of the quadratic objective function into a linear objective function with a SOC constraint. Indeed, the problem

$$\min_{\mathbf{x} \in C} \frac{1}{2} \mathbf{x}^\top M \mathbf{x} + \mathbf{x}^\top \mathbf{b}$$

is equivalent to

$$\begin{cases} \min_{\mathbf{x} \in C, t \in \mathbb{R}} t + \mathbf{x}^\top \mathbf{b} \\ 2t \geq \mathbf{x}^\top M \mathbf{x}. \end{cases} \quad (3.2)$$

Given a square root  $L$  of  $M$  (for instance, its Cholesky factorization), Problem (3.2) becomes

$$\begin{cases} \min_{\mathbf{x} \in C, t \in \mathbb{R}} t + \mathbf{x}^\top \mathbf{b} \\ L^\top \mathbf{x} = \mathbf{z} \\ 2t \geq \|\mathbf{z}\|^2. \end{cases}$$

The constraint  $t \geq \|\mathbf{z}\|^2$  can then be recognized as the rotated cone constraint  $(1, t, \mathbf{z}) \in \mathcal{R}\mathcal{H}$ , i.e.,  $2 \times 1 \times t \geq \|\mathbf{z}\|^2$ . Rotated cone constraints are handled by most SOCP solvers, or can be recast as standard SOC constraints.

The SOCP corresponding to the primal SOCQP (2.24) is therefore

$$\begin{cases} \min_{\mathbf{v} \in \mathbb{R}^m, t \in \mathbb{R}} t - \mathbf{v}^\top \mathbf{f} \\ \mathbf{z} = L^\top \mathbf{v} \\ \mathbf{k} = C \mathbf{v} \\ \mathbf{u} = H \mathbf{v} + \mathbf{x} + \mathbf{s} \\ \mathbf{u}_i \in \mathcal{K}_{\frac{1}{\mu_i}} \quad \forall i = 1 \dots n \\ (1, t, \mathbf{z}) \in \mathcal{R}\mathcal{H}, \end{cases} \quad (3.3)$$

<sup>4</sup>Including several out-of-the-box commercial packages, such as CPLEX, MOSEK, Gurobi, ...

and the dual SOCQP (2.21) can be in turn transformed into

$$\left\{ \begin{array}{ll} \min_{\mathbf{r} \in \mathbb{R}^{nd}, \boldsymbol{\lambda} \in \mathbb{R}^c, t \in \mathbb{R}} & t + \mathbf{r}^\top (\mathbf{b} + \mathbf{s}) + \boldsymbol{\lambda}^\top \mathbf{c} \\ & L^\top \mathbf{z} = H^\top \mathbf{r} + C^\top \boldsymbol{\lambda} \\ & \mathbf{r}_i \in \mathcal{K}_{\mu_i} \\ & (1, t, \mathbf{z}) \in \mathcal{R}\mathcal{K} \end{array} \right. \quad \forall i = 1 \dots n \quad (3.4)$$

with the square-root matrix  $L$  such that  $M = LL^\top$ .

Note that this technique can only be applied to the convexified DCFP; there is no hope of casting the original DCFP into a SOCP.

### 3.2.2 Discussion

**Warm-starting** A major drawback of interior-point methods is that they are hard to warm-start; that is, the solution to a similar (or even to the same) optimization problem will not be able to serve as an initial guess for a subsequent solve.

Indeed, while the solution from a previous solve may be arbitrarily close to the boundary of the constraint, most interior-point algorithms start with a relatively high value for the barrier parameter  $\beta$ . This will make the previous solution very poor (that is, yielding a very high value of the objective function), and the next iterate will go back towards the center of the constraint.

This drawback is especially serious for the Cadoux algorithm, which consists in solving a sequence of SOCQP that become more and more similar as the fixed-point algorithm converges. This motivates the investigation of algorithms with not-as-good convergence properties as interior-point methods, but whose potential slowness on single problems may be compensated by their ability to be easily warm-started.

**Heyn's algorithm** While we have not implemented nor tested this algorithm, we would like to mention that Heyn (2013, Section 4.5) also proposes a direct primal-dual interior-point method that does not require the SOCP reformulation. He observes very good convergence w.r.t. the number of iterations, but deplores that the computational cost per iteration is much higher than for proximal or splitting methods (which we will discuss below).

## 3.3 First-order proximal methods

### 3.3.1 Proximal operator

We refer the reader to (Parikh and S. P. Boyd 2014) for an excellent introduction and overview of proximal algorithms. In this manuscript we will focus on very narrow applications of this framework, and won't make use of the more abstract general setting. Yet, some knowledge about proximal methods may help in finding links between the different methods presented below.

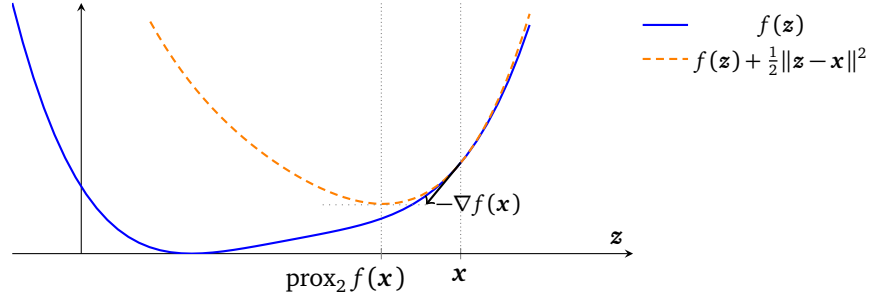
**Definition 3.1.** *The proximal operator with coefficient  $\beta > 0$  of a closed proper convex function<sup>5</sup>  $f : X \rightarrow \bar{\mathbb{R}}$ , with  $X$  a reflexive Banach space, is defined as the solution to a convex minimization problem,*

$$\begin{aligned} \text{prox}_\beta f : X &\rightarrow X \\ \mathbf{x} &\mapsto \arg \min_{\mathbf{z} \in X} f(\mathbf{z}) + \frac{1}{2\beta} \|\mathbf{x} - \mathbf{z}\|. \end{aligned} \quad (3.5)$$

For  $\beta = +\infty$ ,  $\text{prox}_\beta f$  may be multi-valued, as any point of  $X$  optimum for  $f$  satisfies the definition. However, for  $\beta < +\infty$ , the function being minimized becomes strictly convex and

<sup>5</sup>See Appendix A.1.1





**Figure 3.2:** Geometrical evaluation of the proximal operator and relationship with a gradient descent step

coercive, and therefore from Theorem A.5, the proximal operator is uniquely defined. As  $\beta$  goes towards 0, the minimization problem becomes more and more skewed towards the point at which the operator is evaluated. Alternatively, as illustrated in Figure 3.2, evaluating the proximal operator  $\text{prox}_\beta f(\mathbf{x})$  may be seen as moving  $\mathbf{x}$  towards the minimum of  $f$ ; the lower  $\beta$ , the smaller the step size. More precisely, Parikh and S. P. Boyd (2014, Section 3.3) point out that when  $f$  is twice-differentiable,

$$\text{prox}_\beta f(\mathbf{x}) = \mathbf{x} - \beta(\nabla f)(\mathbf{x}_0) + o(\beta). \quad (3.6)$$

The fixed-point algorithm defined by the induction  $\mathbf{x}^{k+1} := \text{prox}_\beta f(\mathbf{x}^k)$  can be thus seen as a modified gradient descent step. Corroborating this interpretation, it can be shown that  $\mathbf{x}_0 \in X$  is a minimum for  $f$  on  $X$  if and only if  $\text{prox}_\beta f(\mathbf{x}_0) = \mathbf{x}_0$ , i.e., if and only if  $\mathbf{x}_0$  is a fixed-point of the proximal operator (Parikh and S. P. Boyd 2014, Section 2.3).

An interesting property of the proximal operator manifests itself when  $f = \mathcal{J}_C$ , the characteristic function of  $C \subset X$ , a proper, closed, convex subset of  $X$  (Definition A.8). Then

$$\text{prox}_\beta \mathcal{J}_C(\mathbf{x}) = \arg \min_{\mathbf{z} \in C} \frac{1}{2\beta} \|\mathbf{x} - \mathbf{z}\| = \Pi_C(\mathbf{x}),$$

the orthogonal projection on  $C$ . This interpretation of the orthogonal projection as a special case of a proximal operator translates to the optimality conditions of convex optimization problems. Recall (Remark A.3) that the optimality condition for the minimization problem (A.4),

$$\min_{\mathbf{x} \in X} f(\mathbf{x}) + \mathcal{J}_C(\mathbf{x}) \quad (A.4)$$

reads, when  $f$  is differentiable on  $C$ , and for any  $\alpha \in \mathbb{R}_+^*$ ,

$$\Pi_C(\mathbf{x}_0 - \alpha(\nabla f)(\mathbf{x}_0)) = \mathbf{x}_0;$$

a similar but more general optimality condition can be stated on proximal operators (Parikh and S. P. Boyd 2014, Section 4.2).

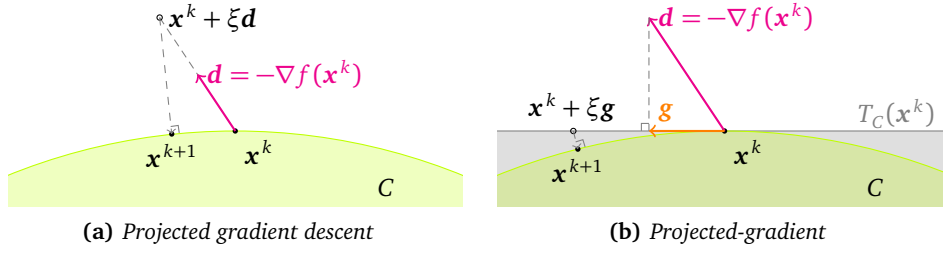
**Property 3.1.** *The optimality conditions of the minimization problem*

$$\min_{\mathbf{x} \in X} f(\mathbf{x}) + g(\mathbf{x})$$

with  $f$  and  $g$  proper closed convex and  $f$  differentiable on  $\text{dom } g$  are

$$\text{prox}_\beta g(\mathbf{x}_0 - \alpha(\nabla f)(\mathbf{x}_0)) = \mathbf{x}_0, \quad \alpha \in \mathbb{R}_+^*.$$

Any fixed-point of the sequence induced by  $\mathbf{x}^{k+1} := \text{prox}_\beta g(\mathbf{x}^k - \alpha(\nabla f)(\mathbf{x}^k))$  will thus yield a minimum of  $f + g$  over  $X$ . When  $(\nabla f)$  is Lipschitz with constant  $L$ , this fixed-point algorithm will converge as  $O(1/k)$  for  $\beta \in ]0, \frac{1}{L}]$  (Parikh and S. P. Boyd 2014) — though it is only practical as long as evaluating  $\text{prox}_\beta g$  is simple enough. The projected gradient descent method, presented below, will thus appear as a special case of fixed-point algorithms on a proximal operator where  $g = \mathcal{J}_C$ .



**Figure 3.3:** Comparison of the projected gradient descent and projected-gradient steps at a point  $\mathbf{x}^k$  in the boundary of the feasible set  $C$ .

### 3.3.2 Projected Gradient Descent

The projected gradient descent algorithm simply consists in taking a step in the direction opposite to the objective function's gradient, then projecting the result onto the feasible set.

Applied to the dual SOCQP (2.21), it reads as Algorithm 3.1,

---

**Algorithm 3.1:** Canonical Projected Gradient Descent Algorithm

---

**Result:**  $\mathbf{r}$  approximate solution to  $\arg \min_{\mathbf{r} \in \mathcal{X}_\mu} \frac{1}{2} \mathbf{r}^\top W \mathbf{r} + \mathbf{r}^\top \mathbf{b}$

**for**  $k \in \mathbb{N}$  **do**

$\mathbf{u}^k \leftarrow \tilde{W} \mathbf{r}^k + \tilde{\mathbf{b}}$ ;

    Compute step size  $\xi^{k+1}$  using a line-search procedure;

3.1.4    $\mathbf{r}^{k+1} \leftarrow \Pi_{\mathcal{X}_\mu}(\mathbf{r}^k - \xi^{k+1} \mathbf{u}^k)$ ;

**end**

---

It is easy to see that any fixed-point of this procedure will satisfy  $\Pi_{\mathcal{X}_\mu}(\mathbf{r} - \xi \mathbf{u}) = \mathbf{r}$ , and thus will be an optimum of the dual SOCQP. The iterative process can also be terminated as soon as the norm of a complementarity function (see Section 3.1.2) at  $(\mathbf{u}, \mathbf{r})$  yields an error below a chosen tolerance.

Theoretically, a similar algorithm could also be written for the primal SOCQP; however, except for a very specific structure of  $H$ , the projection on the feasible set  $V := \{\mathbf{v} \in \mathbb{R}^m, H\mathbf{v} = \mathbf{w} \in \mathcal{X}_{\frac{1}{\mu}}\}$  cannot be explicitly computed.

Algorithm 3.1 can also be repurposed to directly solve the DCFP (instead of the convexified version). Indeed, as mentioned earlier, the projection gradient descent algorithm can simply be seen as a fixed-point algorithm with a varying step size. Replacing  $\mathbf{u}^k$  in line 3.1.4 with  $\tilde{\mathbf{u}}(\mathbf{u}^k)$ , we obtain the fixed-point algorithm on the De Saxcé formulation (with a well-chosen step size); replacing  $\Pi_{\mathcal{X}_\mu}$  with  $\mathcal{T}_{0, \mu r_N}$ , the fixed-point algorithm on the Alart–Curnier formulation. However, these algorithms no longer minimize convex functions under convex constraints, and the theoretical convergence properties are lost — which does not prevent them from sometimes being quite effective in practice, as confirmed by our benchmarks presented in Chapters 4 and 7.

**Acceleration** Under regularity conditions, the projected gradient descent algorithm converges as  $O(1/k)$  (Parikh and S. P. Boyd 2014); however, incorporating a momentum term may yield convergence as  $O(1/k^2)$ . Heyn (2013) applied the Nesterov acceleration (Nesterov 1983) to SOCQP minimization, yielding the APGD method and allowing efficient solving of very large convexified DCFP (Mazhar et al. 2015).

**Projected gradient** When on the boundary and close to the solution, the gradient will be very close to the normal cone of the constraint, and therefore unless the step size is chosen to be quite

big, the projected position will be very close to the original position (Figure 3.3). The projected-gradient algorithm alleviates this problem by first projecting the gradient onto the tangent cone  $T_C$  to the constraint  $C \subset X$ ,

$$T_C(\mathbf{x}) := \overline{\{\mathbf{y} \in X, \exists \epsilon > 0, \forall 0 < t < \epsilon, (\mathbf{x} + t\mathbf{y}) \in C\}}.$$

---

**Algorithm 3.2:** Projected-Gradient Algorithm

---

```

for  $k \in \mathbb{N}$  do
     $\mathbf{u}^k \leftarrow \tilde{W}\mathbf{r}^k + \tilde{\mathbf{b}}$ ;
    3.2.3  $\mathbf{g}^{k+1} \leftarrow \Pi_{T_{\mathcal{K}_\mu}(\mathbf{r}^k)}\mathbf{u}^k$ ;
        Compute step size  $\xi^{k+1}$  using a line-search procedure;
         $\mathbf{r} \leftarrow \Pi_{\mathcal{K}_\mu}(\mathbf{r}^k - \xi^{k+1}\mathbf{g}^{k+1})$ ;
    end
```

---

If the projection tangent cone is hard to compute, it can be approximated by performing another projection on the feasible set  $C$ . Line 3.2.3 is then replaced with two steps,

---

```

 $\mathbf{x} \leftarrow \Pi_{\mathcal{K}_\mu}(\mathbf{r} - \rho\mathbf{g})$ ;
 $\mathbf{g} \leftarrow \mathbf{x} - \mathbf{r}$ ;
```

---

**Other variants** Many other variants of the projected-gradient and projected gradient descent exist. We can mention the Spectral Projected-Gradient (SPG) method, which was observed to work well on convexified DCFP by Tasora (2013). In Appendix B.2, we propose a Nesterov-accelerated, line-search free variant of this algorithm, which we coined **ASPG** and that performed consistently well on our problems (both convexified and standard). Algorithm B.2 is implemented in the bogus library, along with the APGD method, canonical algorithms, and a projected-gradient variant augmented with a simple conjugation step.

### 3.3.3 Primal–dual proximal methods

Primal–dual proximal methods are a class of first-order minimization algorithms that work by splitting the objective function, and iterating on both the primal variable (in our case, the velocity  $\mathbf{v}$ ) and the dual variable (the forces  $\mathbf{r}$ ). The methods presented below were not observed to converge very well on our DCFP without requiring significant manual tuning of the step size parameters, and we therefore won't use them in practice. However, their popularity in computational fluid mechanics make them worthy of interest.

**Alternating Directions Method of Multipliers** ADMM is another way of minimizing separable problems, introduced by Glowinski and Marroco (1975) for non-differentiable continuum flow problems and made popular by Fortin and Glowinski (1983) under the name of Augmented Lagrangian algorithm. ADMM has also been recently introduced to the Computer Graphics simulation community by Narain, Overby, et al. (2016), though in a slightly different setting.

The minimization of  $f(\mathbf{v}) + g(H\mathbf{v} + \mathbf{w})$  is written as

$$\min_{\mathbf{u}=H\mathbf{v}+\mathbf{w}} f(\mathbf{v}) + g(\mathbf{u})$$

The Lagrangian associated to this minimization problem is

$$\mathcal{L} = f(\mathbf{v}) + g(\mathbf{u}) + \mathbf{r}^\top(\mathbf{u} - H\mathbf{v} - \mathbf{w}),$$

and the augmented Lagrangian with parameter  $\gamma$  is

$$\mathcal{L}_\gamma = f(\mathbf{v}) + g(\mathbf{u}) + \mathbf{r}^\top(\mathbf{u} - H\mathbf{v} - \mathbf{w}) + \frac{\gamma}{2}\|H\mathbf{v} + \mathbf{w} - \mathbf{u}\|^2.$$

We are looking for a saddle-point of  $\mathcal{L}_\gamma$ . Minimizing over the first two variables  $\mathbf{v}$  and  $\mathbf{u}$  in turn, and taking a gradient step to maximize  $\lambda$ , we obtain Algorithm 3.3,

---

**Algorithm 3.3:** Alternating Directions Method of Multipliers

---

```

for  $k \in \mathbb{N}$  do
3.3.2    $\mathbf{v}^{k+1} \leftarrow \arg \min_{\mathbf{v}} \mathcal{L}_\gamma(\mathbf{v}, \mathbf{u}^k, \mathbf{r}^k);$ 
3.3.3    $\mathbf{u}^{k+1} \leftarrow \arg \min_{\mathbf{u}} \mathcal{L}_\gamma(\mathbf{v}^{k+1}, \mathbf{u}, \mathbf{r}^k)(\mathbf{v}^{k+1} + \mathbf{r}^k);$ 
        $\mathbf{r}^{k+1} \leftarrow \mathbf{r}^k - \gamma(H\mathbf{v} + \mathbf{w} - \mathbf{u});$ 
end
```

---

Now, Line 3.3.3 can easily be written as a proximal operator:

$$\begin{aligned} \mathbf{u}^{k+1} &:= \arg \min_{\mathbf{u}} g(\mathbf{u}) + \frac{\gamma}{2}\|H\mathbf{v}^{k+1} + \mathbf{w} - \mathbf{u} - \frac{\mathbf{r}^k}{\gamma}\| \\ &= \text{prox}_{\frac{1}{\gamma}} g(H\mathbf{v}^{k+1} + \mathbf{w} - \frac{1}{\gamma}\mathbf{r}^k) \end{aligned}$$

The first step, Line 3.3.2, is more tricky, but can be dealt with by linearizing the quadratic term around  $\mathbf{v}^k$  and adding a new quadratic regularization term, (Parikh and S. P. Boyd 2014, Section 4.4.2),

$$\begin{aligned} \frac{\gamma}{2}\|H\mathbf{v} + \mathbf{w} - \mathbf{u}^k\|^2 &\sim \gamma \langle H\mathbf{v}^k + \mathbf{w} - \mathbf{u}^k, H\mathbf{v} \rangle + \frac{\beta}{2}\|\mathbf{v} - \mathbf{v}^k\|^2 \\ &\sim \frac{\beta}{2}\|\mathbf{v} - \mathbf{v}^k + \frac{\gamma}{\beta}H^\top(H\mathbf{v}^k + \mathbf{w} - \mathbf{u}^k)\|^2. \end{aligned}$$

This yields the *linearized* ADMM algorithm, with Line 3.3.2 replaced by

$$\begin{aligned} \mathbf{v}^{k+1} &\leftarrow \arg \min_{\mathbf{v}} f(\mathbf{v}) + \frac{\beta}{2}\|\mathbf{v} - \mathbf{v}^k + \frac{\gamma}{\beta}H^\top(H\mathbf{v}^k + \mathbf{w} - \mathbf{u}^k) - \frac{1}{\beta}H^\top\mathbf{r}^k\|^2 \\ &= \text{prox}_{\frac{1}{\beta}} f\left(\mathbf{v}^k + \frac{1}{\beta}H^\top\mathbf{r}^k - \frac{\gamma}{\beta}H^\top(H\mathbf{v}^k + \mathbf{w} - \mathbf{u}^k)\right) \\ &= \text{prox}_{\frac{1}{\beta}} f\left(\mathbf{v}^k + \frac{1}{\beta}H^\top\left(\mathbf{r}^k - \gamma(H\mathbf{v}^k + \mathbf{w} - \mathbf{u}^k)\right)\right) \end{aligned} \quad (3.7)$$

The ADMM algorithm 3.3 can also be interpreted as an integral controller, where  $\mathbf{r}$  is the control  $\mathbf{u}$  the state, and the goal is to achieve  $\mathbf{u} = H\mathbf{v} + \mathbf{w}$  (Parikh and S. P. Boyd 2014).

**ADMM on primal formulation** Now, let us consider our primal SOCQP (without linear constraints for now). We have  $f(\mathbf{v}) = \frac{1}{2}\mathbf{v}^\top M\mathbf{v} - \mathbf{v}^\top \mathbf{f}$ , and  $g(\mathbf{u}) := \mathcal{G}_{\frac{1}{\mu}}(\mathbf{u})$ . Therefore

$$\begin{aligned} \text{prox}_{\frac{1}{\beta}} f(\mathbf{x}) &= \arg \min_{\mathbf{v}} \frac{1}{2}\mathbf{v}^\top M\mathbf{v} - \mathbf{v}^\top \mathbf{f} + \frac{\beta}{2}\|\mathbf{v} - \mathbf{x}\|^2 \\ &= \arg \min_{\mathbf{v}} \frac{1}{2}\mathbf{v}^\top (M + \beta\mathbb{I})\mathbf{v} - \mathbf{v}^\top (\mathbf{f} + \beta\mathbf{x}) \\ &= (M + \beta\mathbb{I})^{-1}(\mathbf{f} + \beta\mathbf{x}) \\ \text{prox}_{\frac{1}{\gamma}} g(\mathbf{y}) &= \Pi_{\mathcal{K}_{\frac{1}{\mu}}}(\mathbf{y}). \end{aligned}$$

The linearized ADMM algorithm on the primal SOCQP thus becomes Algorithm 3.4,

---

**Algorithm 3.4:** Linearized ADMM for SOCQP

---

```

for  $k \in \mathbb{N}$  do
     $\mathbf{v}^{k+1} \leftarrow (M + \beta \mathbb{I})^{-1} (\mathbf{f} + \beta \mathbf{v}^k + H^\top (\mathbf{r}^k - \gamma (H\mathbf{v}^k + \mathbf{w} - \mathbf{u}^k)))$ ;
     $\mathbf{u}^{k+1} \leftarrow \Pi_{\mathcal{K}_{\frac{1}{\mu}}} (H\mathbf{v}^{k+1} + \mathbf{w} - \frac{1}{\gamma} \mathbf{r}^k)$ ;
     $\mathbf{r}^{k+1} \leftarrow \mathbf{r}^k - \gamma (H\mathbf{v}^{k+1} + \mathbf{w} - \mathbf{u}^{k+1})$ ;
end
    
```

---

We can easily verify that any fixed point  $(\mathbf{v}, \mathbf{u}, \mathbf{r})$  of this algorithm satisfies

$$\begin{cases} M\mathbf{v} = \mathbf{f} + H^\top \mathbf{r} \\ \mathbf{u} = H\mathbf{v} + \mathbf{w} \\ \mathbf{r} \in -\mathcal{N}_{\mathcal{K}_{\frac{1}{\mu}}} \mathbf{u} \end{cases}$$

and is thus a solution of the primal SOCQP without linear constraints.

The Alternating Minimization Algorithm (AMA, see Goldstein et al. 2014) is obtained by suppressing the quadratic term in the expression of the Lagrangian when minimizing w.r.t.  $\mathbf{v}$ , that is, replacing Line 3.3.2 with

$$\mathbf{v}^{k+1} \leftarrow \arg \min_{\mathbf{v}} f(\mathbf{v}) - \langle \mathbf{r}^k, H\mathbf{v} \rangle$$

Applied to our SOCQP, we obtain Algorithm 3.5,

---

**Algorithm 3.5:** AMA for SOCQP

---

```

for  $k \in \mathbb{N}$  do
     $\mathbf{v}^{k+1} \leftarrow M^{-1} (\mathbf{f} + H^\top \mathbf{r}^k)$ ;
     $\mathbf{u}^{k+1} \leftarrow \Pi_{\mathcal{K}_{\frac{1}{\mu}}} (H\mathbf{v}^{k+1} + \mathbf{w} - \frac{1}{\gamma} \mathbf{r}^k)$ ;
     $\mathbf{r}^{k+1} \leftarrow \mathbf{r}^k - \gamma (H\mathbf{v}^{k+1} + \mathbf{w} - \mathbf{u}^{k+1})$ ;
end
    
```

---

In contrast to the original ADMM algorithm, AMA may be ill-defined when  $M$  is only positive semi-definite. Both algorithms can be easily extended to handle linear constraints; for ADMM, we get Algorithm 3.6,

---

**Algorithm 3.6:** Linearized ADMM for SOCQP with linear constraints

---

```

for  $k \in \mathbb{N}$  do
     $\mathbf{v}^{k+1} \leftarrow (M + \beta \mathbb{I})^{-1} (\mathbf{f} + \beta \mathbf{v}^k + H^\top (\mathbf{r}^k - \gamma_1 (H\mathbf{v}^k + \mathbf{w} - \mathbf{u}^k)) + C^\top (\boldsymbol{\lambda}^k - \gamma_2 (C\mathbf{v}^k - \mathbf{k})))$ ;
     $\mathbf{u}^{k+1} \leftarrow \Pi_{\mathcal{K}_{\frac{1}{\mu}}} (H\mathbf{v}^{k+1} + \mathbf{w} - \frac{1}{\gamma_1} \mathbf{r}^k)$ ;
     $\mathbf{r}^{k+1} \leftarrow \mathbf{r}^k - \gamma_1 (H\mathbf{v}^{k+1} + \mathbf{w} - \mathbf{u}^{k+1})$ ;
     $\boldsymbol{\lambda}^{k+1} \leftarrow \boldsymbol{\lambda}^k - \gamma_2 (C\mathbf{v}^{k+1} - \mathbf{k})$ ;
end
    
```

---

Like the projected-gradient algorithm on the dual SOCQP, this algorithm requires solving at each iteration a linear systems with  $M$  (or  $M + \beta \mathbb{I}$ ) as left-hand-side. However, the constraint projection is here done on  $\mathcal{K}_{\frac{1}{\mu}}$ , instead of  $\mathcal{K}_{\mu}$  for the projected gradient. One may also attempt to directly solve the original DCFP by replacing the term  $\mathbf{w}$  in each step with the expression  $\mathbf{w} + \mathbf{s}(H\mathbf{v} + \mathbf{w})$ ; however, once again the theoretical convergence properties are then lost.

Starting from dual SOCQP we can also define an AMA algorithm that will only require one multiplication by  $M$  at each iteration.

**AMA on dual formulation** We write the dual SOCQP (2.21) without linear constraints as

$$\min_{\mathbf{z}=H^T\mathbf{r}} \underbrace{\frac{1}{2}\mathbf{z}^T M^{-1}\mathbf{z} + \mathbf{z}^T M^{-1}\mathbf{f}}_{f(\mathbf{z})} + \underbrace{\mathbf{r}^T \mathbf{w} + \mathcal{J}_{\mathcal{K}_\mu}(\mathbf{r})}_{g(\mathbf{r})}.$$

Writing the AMA on this formulation yields Algorithm 3.7,

---

**Algorithm 3.7:** AMA on dual formulation

```

for  $k \in \mathbb{N}$  do
  3.7.3  $\mathbf{z}^{k+1} \leftarrow \arg \min_{\mathbf{z}} f(\mathbf{z}) - \langle \mathbf{z}, \mathbf{v}^k \rangle;$ 
   $\mathbf{r}^{k+1} \leftarrow \arg \min_{\mathbf{r}} g(\mathbf{r}) + \langle H^T \mathbf{r}, \mathbf{v}^k \rangle + \frac{\gamma}{2} \|\mathbf{z}^{k+1} - H^T \mathbf{r}\|^2;$ 
   $\mathbf{v}^{k+1} \leftarrow \mathbf{v}^k + \gamma(H^T \mathbf{r} - \mathbf{z});$ 
end
```

---

The first minimization is straightforward;

$$\begin{aligned} \mathbf{z}^{k+1} &= \arg \min_{\mathbf{z}} \frac{1}{2} \mathbf{z}^T M^{-1} \mathbf{z} + \mathbf{z}^T (M^{-1} \mathbf{f} - \mathbf{v}^k) \\ &= M(\mathbf{v}^k - \mathbf{f}). \end{aligned}$$

However the second step has to be reworked. Linearizing once again the quadratic part, we replace Line 3.7.3 with

$$\begin{aligned} \mathbf{r}^{k+1} &\leftarrow \arg \min_{\mathbf{r}} g(\mathbf{r}) + \langle H^T \mathbf{r}, \mathbf{v}^k \rangle + \gamma \langle H^T \mathbf{r}^k - \mathbf{z}, H^T \mathbf{r} \rangle + \frac{\beta}{2} \|\mathbf{r} - \mathbf{r}^k\|^2 \\ &= \arg \min_{\mathbf{r}} g(\mathbf{r}) + \frac{\beta}{2} \|\mathbf{r} - \mathbf{r}^k + \frac{1}{\beta} H(\mathbf{v}^k + \gamma(H^T \mathbf{r}^k - \mathbf{z}))\|^2 \\ &= \arg \min_{\mathbf{r} \in \mathcal{K}_\mu} \frac{\beta}{2} \|\mathbf{r} - \mathbf{r}^k + \frac{1}{\beta} \mathbf{w} + \frac{1}{\beta} H(\mathbf{v}^k + \gamma(H^T \mathbf{r}^k - \mathbf{z}^{k+1}))\|^2 \\ &= \Pi_{\mathcal{K}_\mu} \left( \mathbf{r}^k - \frac{1}{\beta} [\mathbf{w} + H(\mathbf{v}^k + \gamma(H^T \mathbf{r}^k - \mathbf{z}))] \right). \end{aligned}$$

We see easily that any fixed-point  $(\mathbf{v}, \mathbf{z}, \mathbf{r})$  of this linearized version of Algorithm 3.7 satisfies

$$\begin{cases} M\mathbf{v} = \mathbf{z} + \mathbf{f} \\ \mathbf{z} = H^T \mathbf{r} \\ H\mathbf{v} + \mathbf{w} \in -\mathcal{N}_{\mathcal{K}_\mu}(\mathbf{r}), \end{cases}$$

and thus is also a solution of the convexified DCFP. Once again, this algorithm can be extended to handle linear constraints, by driving the auxiliary variable  $\mathbf{z}$  to satisfy  $\mathbf{z} = H^T \mathbf{r} + C^T \boldsymbol{\lambda}$ . We obtain Algorithm 3.8,

---

**Algorithm 3.8:** Linearized AMA on dual formulation with linear constraints

```

for  $k \in \mathbb{N}$  do
  3.8.3  $\mathbf{z}^{k+1} \leftarrow M\mathbf{v}^k + \mathbf{f};$ 
   $\mathbf{r}^{k+1} \leftarrow \Pi_{\mathcal{K}_\mu} \left( \mathbf{r}^k - \frac{1}{\beta_1} [\mathbf{w} + H(\mathbf{v}^k + \gamma(H^T \mathbf{r}^k + C^T \boldsymbol{\lambda}^k - \mathbf{z}^{k+1}))] \right);$ 
   $\boldsymbol{\lambda}^{k+1} \leftarrow \boldsymbol{\lambda}^k - \frac{1}{\beta_2} [-\mathbf{k} + C(\mathbf{v}^k + \gamma(H^T \mathbf{r}^k + C^T \boldsymbol{\lambda}^k - \mathbf{z}^{k+1}))];$ 
   $\mathbf{v}^{k+1} \leftarrow \mathbf{v}^k + \gamma(H^T \mathbf{r}^{k+1} + C^T \boldsymbol{\lambda}^{k+1} - \mathbf{z}^{k+1});$ 
end
```

---

Just like with the primal formulation, one can directly solve the original DCFP by replacing  $\mathbf{w}$  on line 3.8.3 with  $\mathbf{w} + \mathbf{s}(H\mathbf{v} + \mathbf{w})$ .

**Variants** Goldstein et al. (2014) also define Nesterov-like accelerated versions of those algorithms. Other variants of such separation schemes exist, such as the Arrow–Hurwitz algorithm (which is itself an extension the Uzawa algorithm to the proximal setting):

---

```

for  $k \in \mathbb{N}$  do
     $\mathbf{r}^{k+1} \leftarrow -\text{prox}_{\frac{1}{\gamma}} g^* (-\mathbf{r}^k + \gamma (H\mathbf{v}^k + \mathbf{w}))$ ;
     $\mathbf{v}^{k+1} \leftarrow \text{prox}_{\frac{1}{\beta}} f(\mathbf{v}^k + \frac{1}{\beta} H^\top \mathbf{r}^{k+1})$ ;
end
    
```

---

which for our primal SOCQP is written

**Algorithm:** Arrow–Hurwitz algorithm for SOCQP

---

```

for  $k \in \mathbb{N}$  do
     $\mathbf{r}^{k+1} \leftarrow \Pi_{\mathcal{K}_\mu} (\mathbf{r}^k - \gamma (H\mathbf{v}^k + \mathbf{w}))$ ;
     $\mathbf{v}^{k+1} \leftarrow (M + \beta \mathbb{I})^{-1} (f + \beta \mathbf{v}^k + H^\top \mathbf{r}^{k+1})$ ;
end
    
```

---

Compared to our original ADMM algorithm, the constraint projection is done on the forces  $\mathbf{r}$  instead of on the relative velocities  $\mathbf{u}$ ; the Arrow–Hurwitz degrades to a projected gradient descent algorithm (without line search) when  $\beta = 0$ . The popular Chambolle and Pock (2010) algorithm is an accelerated version of the Arrow–Hurwitz algorithm.

The bogus library (Daviet 2013), implements the accelerated ADMM from (Goldstein et al. 2014) and a preconditioned and accelerated version of the AMA on the dual SOCQP. In practice, our early tests were not very satisfying; fine-tuning of the coefficients  $\beta$  and  $\gamma$  was require to achieve fast convergence.

### 3.4 Splitting methods

#### 3.4.1 Operator splitting

The most popular way to solve the dual formulation of the DCFP (2.19) is to use a splitting method: that is, treating one part of the operator  $W$  implicitly, and another part explicitly. In practice, this means decomposing  $W$  as  $W = W_1 + W_2$ , and iterating as

$$\mathbf{u}^{k+1} = W_1 \mathbf{r}^{k+1} + W_2 \mathbf{r}^k + C^\top \boldsymbol{\lambda}^k + \tilde{\mathbf{b}} \quad (3.8)$$

$$(\mathbf{u}_i^{k+1}, \mathbf{r}_i^{k+1}) \in \mathcal{C}_{\mu_i} \forall 1 \leq i \leq n. \quad (3.9)$$

$$P\boldsymbol{\lambda}^{k+1} = \mathbf{c} - B^\top \mathbf{r}^{k+1}. \quad (3.10)$$

a fixed-point for  $(\mathbf{r}, \boldsymbol{\lambda})$  is reached, that is  $(\mathbf{r}^{k+1}, \boldsymbol{\lambda}^{k+1}) = (\mathbf{r}^k, \boldsymbol{\lambda}^k)$ .

Obviously,  $W_1$  and  $W_2$  should be chosen so that the resulting problem (3.8 – 3.9) is easier to solve than the original DCFP. The convexified DCFP can also be solved in this manner, replacing the Coulomb law condition (3.9) with  $\mathcal{K}_{\frac{1}{\mu_i}} \ni \mathbf{u}_i^{k+1} \perp \mathbf{r}_i^{k+1} \in \mathcal{K}_{\mu_i}$ .

In practice,  $W$  has a sparse block structure composed of  $d \times d$  blocks  $(W_{0j})$ ,  $1 \leq i, j \leq n$ , where each block  $(W_{ij})$  describe the interaction between the  $i^{\text{th}}$  and  $j^{\text{th}}$  contacts (that is, the relative velocity change at  $i$  when a force increment is applied at  $j$ ).  $W$  can thus be decomposed as the sum of a block-diagonal matrix  $W_D := (W_{ii})$ , a lower part  $W_L := (W_{i,j})_{i>j}$  and an upper part  $W_U := (W_{i,j})_{i<j}$ .

Choosing  $W_1 = D$  and  $W_2 = L + U$  yields the Jacobi algorithm; equation (3.8) becomes separable, and result in  $n$   $d$ -dimensional problems that can be solved in parallel:

$$\begin{cases} \mathbf{u}_i^{k+1} = W_{ii} \mathbf{r}_i^{k+1} + \sum_{j \neq i} W_{i,j} \mathbf{r}_j^k + C^\top \boldsymbol{\lambda}^k + \tilde{\mathbf{b}} \\ (\mathbf{u}_i^{k+1}, \mathbf{r}_i^{k+1}) \in \mathcal{C}_{\mu_i}. \end{cases}$$

Choosing  $W_1 = D + L$  and  $W_2 = U$ , we obtain the Gauss–Seidel algorithm, which boasts better convergence but loses parallelism; the  $d$ -dimensional local problems now have to be solved in order:

$$\begin{cases} \mathbf{u}_i^{k+1} = W_{ii}\mathbf{r}_i^{k+1} + \sum_{j<i} W_{ij}\mathbf{r}_j^{k+1} + \sum_{j>i} W_{ij}\mathbf{r}_j^k + C^\top \boldsymbol{\lambda}^k + \tilde{\mathbf{b}} \\ (\mathbf{u}_i^{k+1}, \mathbf{r}_i^{k+1}) \in \mathcal{C}_{\mu_i}. \end{cases}$$

In both cases, solving our global problem can thus be reduced to solving a sequence of local problems (or *one-contact* problems),

$$\begin{cases} \mathbf{u}_i = W_{ii}\mathbf{r}_i + \tilde{\mathbf{b}}_i \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i}, \end{cases} \quad (3.11)$$

where  $W_{ii}$  is a  $d \times d$ , positive semi-definite matrix.

This one-contact problem is much easier to solve than the original one; indeed, for  $d = 2$  or  $d = 3$ , we can construct an analytical solver, as we will show in Chapter 4.

Splitting methods have been used on contact-related problems for several decades. Raous et al. (1988) applied them to the statics of an elastic body under frictional contact, and mention previous uses in the frictionless case. However, their huge popularity nowadays is largely due to the advent of the *Nonsmooth Contact Dynamics* (NSCD) method (Jean 1999; Jean and Moreau 1992; Jourdan et al. 1998), which proposed to compute the dynamics of Discrete Element Models with frictional contacts using the Gauss–Seidel algorithm and an efficient one-contact solver.

### 3.4.2 Convergence properties

**Interpretation as block-coordinate descent** When solving the convexified DCFP, the local problem (3.11) becomes

$$\begin{cases} \mathbf{u}_i = W_{ii}\mathbf{r}_i + \tilde{\mathbf{b}}_i \\ \mathbf{u}_i \in -\mathcal{N}_{\mathcal{K}_{\mu_i}}(\mathbf{r}_i), \end{cases} \quad (3.12)$$

which from Theorem A.6 we recognize as the optimality conditions of the minimization problem

$$\arg \min_{\mathbf{r}_i \in \mathcal{K}_{\mu_i}} \frac{1}{2} \mathbf{r}_i^\top W_{ii} \mathbf{r}_i + \mathbf{r}_i^\top \tilde{\mathbf{b}}_i = \arg \min_{\mathbf{r}_i \in \mathcal{K}_{\mu_i}} J_i(\mathbf{r}_i), \quad (3.13)$$

where  $J_i(\mathbf{r}_i) := \frac{1}{2} \mathbf{r}_i^\top W \mathbf{r} + \mathbf{r}^\top \mathbf{b}$  is defined by freezing every component of  $\mathbf{r}$  except for  $\mathbf{r}_i$ . As such, successively solving the local problems (3.12) amounts to iteratively minimizing the quadratic objective function  $J(\mathbf{r}) := \frac{1}{2} \mathbf{r}^\top W \mathbf{r} + \mathbf{r}^\top \mathbf{b}$  w.r.t. each contact force  $\mathbf{r}_i$ ; the Gauss–Seidel algorithm can thus be seen as a block-coordinate descent method.

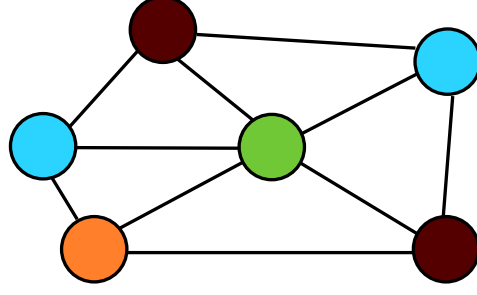
If  $W_{ii}$  is symmetric, positive-definite for all  $i$ , then each local problem admits a unique solution, and the objective function will monotonically decrease to its optimal value (Grippo and Sciandrone 2000).

**Proximal modification** However,  $W_{ii}$  may only be positive semi-definite, and in this case the minimization problem (3.13) will be ill-posed. However, we can modify it slightly as

$$\arg \min_{\mathbf{r}_i \in \mathcal{K}_{\mu_i}} \frac{1}{2} \mathbf{r}_i^\top W \mathbf{r} + \mathbf{r}^\top \mathbf{b} + \frac{\beta}{2} \|\mathbf{r}_i^k - \mathbf{r}_i\|^2 = \text{prox}_{\frac{1}{\beta}} J_i(\mathbf{r}_i^k)$$

with  $\beta > 0$ . Any fixed-point of the proximal operator will be optimal for  $J_i$ , and therefore a fixed-point of the modified Gauss–Seidel algorithm will still be a solution of the dual SOCQP. Grippo and Sciandrone (2000) show that this fixed-point will always be reached for  $W$  symmetric positive semi-definite.





**Figure 3.4:** Sample coloring of a contact graph; contacts (nodes) with the same color can be solved in parallel. Edges represent non-zero  $d \times d$  blocks of the Delassus operator.

By analogy, we define a proximal version of our splitting algorithms for the DCFP by replacing the local problem (3.11) with (3.14),

$$\begin{cases} \mathbf{u}_i = (W_{ii} + \beta \mathbb{I}) \mathbf{r}_i + \bar{\mathbf{b}}_i - \beta \mathbf{r}_i^k \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i}. \end{cases} \quad (3.14)$$

However, note that when using the Gauss–Seidel to solve the DCFP (rather than the convexified DCFP), convergence is no longer guaranteed; indeed, the algorithm may even exhibit cycles. Jourdan et al. (1998) still manage to derive a convergence proof for a particular 2D scenario.

### 3.4.3 Performance

**Accelerations** The Jacobi and Gauss–Seidel names come from the analogy with similar algorithms devised to solve linear systems, and just like their siblings, the splitting methods applied to the DCFP feature a quite slow asymptotic convergence.

Instead of using directly the solution of the local problem (3.11) as  $\mathbf{r}_i^{k+1}$ , one may store this solution as  $\mathbf{r}_i^{k+\frac{1}{2}}$ , and compute the final iterate  $\mathbf{r}_i^{k+1}$  as  $\mathbf{r}_i^{k+1} = (1 - \omega) \mathbf{r}_i^k + \omega \mathbf{r}_i^{k+\frac{1}{2}}$ , with  $\omega$  a positive coefficient. For  $\omega = 1$  we retrieve the Gauss–Seidel algorithm, and for  $\omega > 1$  the resulting method is known as the Successive Over-Relaxation (SOR) algorithm. When used to solve linear systems with a symmetric positive-definite matrix, SOR has been proved to converge for  $0 < \omega < 2$  (Saad 2003, Theorem 4.6). However, finding an  $\omega$  that will improve the convergence rate is in general a hard problem. Raous et al. (1988) proposed to accelerate convergence using the Aitken extrapolation (also known as the  $\Delta^2$  rule), which can also be interpreted as a SOR method with a well-chosen  $\omega$ . Another way to potentially accelerate the convergence is the Symmetric-SOR method, which alternates between solving the local problems from  $i = 1$  to  $n$  and from  $n$  to 1. The contact force updates are then propagated back-and-forth inside the contact chains, instead of in only one direction.

However, we did not find any of these accelerations to perform robustly in practice.

**Parallelism** While the Jacobi algorithm is embarrassingly parallel, the Gauss–Seidel is inherently sequential. On the modern many-cores architectures, this can be a significant drawback.

We can define the adirectional contact graph  $G$ ,  $G := (V, E)$  with vertices  $V = \{1 \dots n\}$  and edges  $E = \{1 \leq i < j \leq n, W_{ij} \neq \mathbf{0}\}$ . Updating the force  $\mathbf{r}_i$  will influence the right-hand-side  $\bar{\mathbf{b}}_j$  if and only if  $(i, j) \in E$ . As such, every disjoint subgraph of  $G$  can be solved in parallel using the Gauss–Seidel algorithm. However, in practice it is common to have a fully connected contact graph. A substitute solution is to use a graph partitioning software, such as the open-source SCOTCH (Pellegrini and Roman 1996), to compute a set of subgraphs  $G_i$  that are weakly connected. The contact subgraphs can then be solved for in parallel, and coupled together with

an outer Jacobi algorithm. Hopefully, the loss of convergence speed will not be too dramatic; however, the more contacts subgraphs are being used, the more degraded the convergence, as the algorithm ultimately becomes a pure Jacobi. Another solution, implemented in *bogus*, is to use a coloring approach, illustrated in Figure 3.4. When the contact graph has a grid-like structure, it is possible to color each node with a checkerboard-like pattern. Then at each Gauss–Seidel iteration, all the white nodes can be solved in parallel, then all the black nodes can be solved in parallel; this is the so-called *red-black* Gauss–Seidel algorithm. It is possible to generalize this approach to arbitrary contact graphs using more colors; however, computing an optimal coloring is a NP-hard task. A greedy algorithm can be easily implemented: each contact is considered in turn, and is either: (i) assigned one of the existing colors if it has no edge to any of the contacts in this color, or (ii) assigned a new color. However, such an algorithm will often lead to a lot of colors with very few contacts, making the overall algorithm very sequential once again. Moreover, computing the partitioning or coloring can be costly, and contacts have then to be reordered to preserve good memory locality; the gain from the recovered parallelism may not suffice to outbalance the cost of the preprocessing cost.

Another approach has been explored by Tonge et al. (2012). Each pair of contacting bodies is considered independently of the other ones, resulting in  $n_B$  distinct velocity vectors for an object in contact with  $n_B$  bodies. The  $n_B$  velocity vectors are prescribed to remain consistent using fixed-joint constraints. The algorithm then iterates between solving for each contacting pair (which can be done in parallel) and enforcing the fixed-joint constraints (for which they provide an analytical formula). This results in a very GPU-friendly method, allowing the authors to solve systems with several thousands of contacts in real-time.

Overall, we found that the most efficient solution for a number of cores that is very small w.r.t. the number of contacts is the most naive one: the set of all contacts is statically assigned to a small number of threads, and each thread is allowed to process its contacts without any synchronization with the others. This might lead to race conditions when one thread is writing a given contact force at the same time as another one is reading it; the reading thread may get back the first components of the force at iteration  $k + 1$ , and the last components at iteration  $k$ <sup>6</sup>. Yet we have not found this to be a problem in practice, and observed instead a good scaling w.r.t. the number of threads. Actually, the randomness induced by the scheduling of the threads was found to sometimes *improve* the convergence of the global algorithm. The major downside, due to this randomness, is that the algorithm is no longer deterministic, and may produce significantly different results each time that it is run on a given DCFP (using the proximal version somewhat reduces this variance).

**Delassus operator** A fundamental requisite to achieve satisfying performance is the ability to efficiently evaluate the local problem’s right-hand-side  $\tilde{\mathbf{b}}_i$ , which means multiplying  $d$  rows of the Delassus operator  $W$  with the current aggregate force vector  $\mathbf{r}$ . Using row-major block-compressed-storage (BSR) for  $W$  yields good performance, but requires explicitly assembling and storing  $W$ ; as mentioned in Section 2.3.1, this can be costly both in time and memory, and prohibitive if the inverse mass matrix  $M^{-1}$  is dense. Moreover, even when  $M$  is the identity matrix,  $HH^T$  may be much more expensive to store than  $H$ .

On the other hand, computing the  $i^{\text{th}}$  block-row-vector product  $W_{i[\cdot]} \mathbf{r} = (HM^{-1}H^T)_{i[\cdot]} \mathbf{r}$  analytically for each contact would be very light on memory, but incur a huge runtime cost. A good middle ground, when  $M$  is block-diagonal, is to consider an intermediate vector  $\mathbf{z} \in \mathbb{R}^m$  storing at each-instant the matrix-vector product  $M^{-1}H^T \mathbf{r}$ . In this case,  $M^{-1}H^T$  possesses the same block structure as  $H^T$ , and can be precomputed. Updating  $\mathbf{z}$  when a contact force  $\mathbf{r}_i$  changes then becomes very cheap. Indeed, it simply amounts to  $d$  scale-add operations on the columns of  $M^{-1}H^T$ . As the rows of  $H$  (and therefore the columns of  $H^T$ ) are usually quite sparse, this operation is inexpensive. Then, the right-hand-side  $\tilde{\mathbf{b}}_i$  can be deduced from  $\mathbf{z}$  by once again

<sup>6</sup>Note that reading and writing aligned double-precision floats is usually atomic on modern architectures, so this approach will not be subject to intra-component corruption.

a cheap multiplication with  $d$  rows of  $H$ . Overall, the inner loop of the Gauss–Seidel algorithm becomes:

---

```

for  $1 \leq i \leq n$  do
     $\tilde{\mathbf{b}}_i \leftarrow \mathbf{b}_i + H_{i[\cdot]} \mathbf{z} - W_{ii} \mathbf{r}_i^k$  ;
    Set  $\mathbf{r}_i^{k+1}$  as solution of local problem (3.11) ;
     $\mathbf{z} \leftarrow \mathbf{z} + (M^{-1} H^\top)_{[\cdot]i} (\mathbf{r}_i^{k+1} - \mathbf{r}_i^k)$  ;
end

```

---

where only the diagonal blocks  $W_{ii}$  have now to be precomputed.

When  $H$  is stored using BSR, and when its rows are quite-sparse, this algorithm remains very efficient, and is also implemented in the bogus library. It is however more prone to data races when synchronization-free parallelism is used; for this reason, we periodically recompute  $\mathbf{z}$  from scratch.

Whether one or the other of these two algorithms performs best when factoring into account the cost of computing  $W$  into account is application-dependent. In practice we usually explicitly compute  $W$  for Discrete Element Models (or generally systems with a high number of degrees of freedom w.r.t. the number of contacts), and use the matrix-free variation for continuum models (the second part of this dissertation). Note that Kaufman, Tamstorf, et al. (2014) recommends avoiding computing explicitly the Delassus operator for non-linear DEM; this seems sensible, as a sequence of DCFP with varying  $W$  but improving initial guesses has then to be solved, which mean relatively less work for the solver w.r.t. the matrix assembly time.

#### 3.4.4 Discussion

On the one hand, despite all the above theoretical drawbacks (slow or lack of convergence, sequentiality, computation of  $W$ ), the Gauss–Seidel algorithm actually performs very well in practice. This will be confirmed by the results presented over the course of the next chapters, such as in Section 4.2.3. As long as the local contact problems can be solved reliably, the global algorithm is very robust, and quickly lowers the error below a tolerance that is acceptable for a lot of purposes.

Jean (1999) states that the local problems can usually be solved using 3 or less iterations of the Newton method on the Alart–Curnier function. We have not found this to be always the case, especially on reduced-coordinates flexible models, and have observed that failure of convergence in the local solver may quickly escalate to divergence of the global solver. This has motivated the conception of a very robust local solver, to which is dedicated the following chapter.

On the other hand, the Gauss–Seidel algorithm is not adapted for problems where the inverse  $M^{-1}$  of the DCFP matrix is dense, or for massively parallel architectures. In this case, we advocate instead using proximal methods. The following chapters will show that one variant of the projected gradient descent, APSG (Algorithm B.2), performs generally well. The Dual AMA (Algorithm 3.8) is also worthy of interest as it does not require explicitly solving any linear system. However, achieving satisfying practical performance would first require devising better heuristics for choosing the step sizes.

## 4 A Robust Gauss–Seidel Solver and its Applications

The first section of this chapter presents a few modifications over the canonical Gauss–Seidel algorithm allowing to tackle problems which cause divergence of most standard implementations. The resulting algorithm performed well on a wide range of problems, was often faster than proximal or interior-points methods, and was faster and more robust than standard GS approaches.

The last three sections are dedicated to applications of this algorithm; the first one, realistic simulation hair dynamics, being the one that motivated building this solver in the first place.

Our hybrid solver (Section 4.1) was originally published in (Daviet, Bertails-Descoubes, and Boissieux 2011), in the context of hair dynamics (Section 4.3), then leveraged for the inverse design of hair (Section 4.4), and, more recently, for cloth dynamics (Section 4.3).

### 4.1 Hybrid Gauss–Seidel algorithm

As mentioned earlier, our applications — especially hair dynamics with reduced coordinate models (see Section 4.2) — caused regular failures of the local solvers commonly in use in the literature: Newton methods or fixed-point iterations over the Alart–Curnier or De Saxcé complementarity functions. These local failures could then escalate to global divergence of the frictional contact solver, leading to visible popping or complete blowup of the simulation.

Our first objective was thus to devise a very robust solver for the local problem (3.11).

#### 4.1.1 SOC Fischer–Burmeister function

One first idea is to try new complementarity functions that could be more suitable for root-finding procedures than the Alart–Curnier or De Saxcé functions.

First, we may want to look back at the scalar case, and consider the linear complementarity problem  $0 \leq u \perp r \geq 0$ . This scalar complementarity problem can be expressed as the normal cone inclusion  $u \in -\mathcal{N}_{\mathbb{R}_+}(r)$ , and thus as the zero of a complementarity function  $f : (u, r) \mapsto \Pi_{\mathbb{R}_+}(r - \xi u) - r$ , which looks somewhat similar to our Alart–Curnier and De Saxcé functions. However, for root-finding purposes this projection-based complementarity function is usually discarded in favor of the Fischer–Burmeister complementarity function (Burmeister 1985; Fischer 1992),

$$f_{\text{FB}} : \begin{cases} \mathbb{R}^2 \longrightarrow \mathbb{R} \\ (u, r) \longmapsto \sqrt{u^2 + r^2} - u - r. \end{cases}$$

The Fischer–Burmeister function is indeed “smoother”; it is only non-differentiable for  $u = r = 0$  (while the projection is not-differentiable for  $r = -\xi u$ ) and its Jacobian has been shown to be regular under reasonable conditions (Fischer 1992). Newton algorithms over the Fischer–Burmeister function have been proved to perform quite well in practice (e.g., Munson et al. 2001; Silcowitz et al. 2009).

Finding an analogous well-behaved complementarity function for SOC problems would therefore be of interest, and the Fischer–Burmeister has fortunately been extended to SOC comple-

mentarity problems by Fukushima et al. (2002). Indeed,

$$\begin{aligned} \mathcal{K}_1 \ni \mathbf{x} \perp \mathbf{y} \in \mathcal{K}_1 &\iff f_{\text{FB}}^{\mathcal{K}}(\mathbf{x}, \mathbf{y}) = \mathbf{0}, \\ f_{\text{FB}}^{\mathcal{K}} : \mathbb{R}^d \times \mathbb{R}^d &\rightarrow \mathbb{R}^d \\ (\mathbf{x}, \mathbf{y}) &\mapsto (\mathbf{u} \circ \mathbf{u} + \mathbf{r} \circ \mathbf{r})^{\frac{1}{2}} - \mathbf{x} - \mathbf{y} \end{aligned} \quad (4.1)$$

where the square-root is also defined w.r.t. the bilinear operator  $\cdot \circ \cdot : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,

$$\mathbf{x} \circ \mathbf{y} \mapsto \begin{pmatrix} \mathbf{x}^\top \mathbf{y} \\ \mathbf{x}_N \mathbf{y}_T + \mathbf{y}_N \mathbf{x}_T \end{pmatrix}.$$

The vectorial space  $\mathbb{R}^d$  equipped with standard addition and the “ $\circ$ ” product defines a Jordan (i.e., commutative, but generally not associative) algebra. The well-definition of the SOC  $f_{\text{FB}}^{\mathcal{K}}$  function will be asserted in a following paragraph.

Now, remember from Property 1.4 that the Coulomb law can be expressed equivalently as

$$(\mathbf{u}, \mathbf{r}) \in \mathcal{C}_\mu \iff \mathcal{K}_{\frac{1}{\mu}} \ni \tilde{\mathbf{u}} \perp \mathbf{r} \in \mathcal{K}_{\frac{1}{\mu}}$$

where  $\tilde{\mathbf{u}} = \mathbf{u} + \mu \|\mathbf{u}_T\|$  is defined from the De Saxcé change of variable. For  $\mu > 0$ , we introduce two new changes of variable,

$$\hat{\mathbf{r}} := \begin{pmatrix} \mu \mathbf{r}_N \\ \mathbf{r}_T \end{pmatrix} \quad \hat{\mathbf{u}} := \rho \begin{pmatrix} \tilde{\mathbf{u}}_N \\ \mu \mathbf{u}_T \end{pmatrix} \quad (4.2)$$

where  $\rho$  is a scalar introduced for conditioning purposes (Here as  $\mathbf{r}$  is an impulse and  $\mathbf{u}$  a velocity, therefore  $\rho$  should be a mass). Then,

$$\mathcal{K}_{\frac{1}{\mu}} \ni \tilde{\mathbf{u}} \perp \mathbf{r} \in \mathcal{K}_{\frac{1}{\mu}} \iff \mathcal{K}_1 \ni \hat{\mathbf{u}} \perp \hat{\mathbf{r}} \in \mathcal{K}_1$$

and thus from Equation (4.1),

$$f_{\text{FB}}^{\mathcal{K}}(\hat{\mathbf{u}}, \hat{\mathbf{r}}) = \mathbf{0} \iff (\mathbf{u}, \mathbf{r}) \in \mathcal{C}_\mu. \quad (4.3)$$

This new complementarity function for Coulomb friction is only valid for  $\mu > 0$ , but since for  $\mu = 0$  Coulomb’s law degenerates to a scalar complementarity problem, in this case we can use the original, scalar Fischer–Burmeister function. The SOC  $f_{\text{FB}}^{\mathcal{K}}$  function can also be trivially adapted to convexified DCFP, as it is sufficient to use the original velocity  $\mathbf{u}$  instead of the De Saxcé change of variable  $\tilde{\mathbf{u}}$  in Equation (4.2).

Unlike what we claimed in (Daviet, Bertails-Descoubes, and Boissieux 2011), the SOC Fischer–Burmeister function has already been applied (and discarded as not efficient) by Cadoux (2009) inside a quasi-Newton algorithm. We proposed instead to use a Newton algorithm, and obtained satisfying results.

**Numerical evaluation** Computing the SOC  $f_{\text{FB}}^{\mathcal{K}}$  function involves taking the square-root of a vector of  $\mathbb{R}^d$  w.r.t. the bilinear operator  $\cdot \circ \cdot$ ; one may wonder about the well-definedness of this operation, and its ease of computation.

The association of the vector space  $\mathbb{R}^d$  with the bilinear operator  $\circ$  forms a  $\mathbb{R}$ -algebra, tightly associated to the SOC  $\mathcal{K}_1$  through two notable properties:

- $\forall \mathbf{x} \in \mathcal{K}_1$ , there exists a unique  $\mathbf{z} \in \mathcal{K}_1$  s.t.  $\mathbf{z} \circ \mathbf{z} = \mathbf{x}$ ; we note  $\mathbf{z} = \mathbf{x}^{\frac{1}{2}}$ ;
- $\forall \mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{x} \circ \mathbf{x} \in \mathcal{K}_1$ .

As  $\mathcal{K}_1$  is a convex cone, we deduce that  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d \times \mathbb{R}^d$ ,  $\mathbf{x} \circ \mathbf{x} + \mathbf{y} \circ \mathbf{y} \in \mathcal{K}_1$  and thus the SOC Fischer–Burmeister function is well-defined.

Moreover, any vector  $\mathbf{x} \in \mathbb{R}^d$  can be decomposed as  $\mathbf{x} = \lambda_1 \mathbf{v}_1 + \lambda_2 \mathbf{v}_2$ , where

$$\lambda_i = \mathbf{x}_N + (-1)^i \|\mathbf{x}_T\| \quad (4.4)$$

$$\mathbf{v}_i = \frac{1}{2} \begin{cases} \left( 1; (-1)^i \frac{\mathbf{x}_T}{\|\mathbf{x}_T\|} \right) & \text{if } \mathbf{x}_T \neq \mathbf{0} \\ \left( 1; (-1)^i \mathbf{e} \right) & \text{if } \mathbf{x}_T = \mathbf{0} \end{cases} \quad (4.5)$$

where  $\mathbf{e}$  is any unit vector in  $\mathbb{R}^{d-1}$ .  $\lambda_1$  and  $\lambda_2$  are called the *eigenvalues* of  $\mathbf{x}$ , and  $\mathbf{v}_1$  and  $\mathbf{v}_2$  are the associated *eigenvectors*. If  $\mathbf{x} \in \mathcal{K}_1$ ,  $\lambda_1$  and  $\lambda_2$  will both be positive; we can thus define  $\mathbf{z} := \sqrt{\lambda_1} \mathbf{v}_1 + \sqrt{\lambda_2} \mathbf{v}_2$ , and check that  $\mathbf{z} \circ \mathbf{z} = \mathbf{x}$ . This yields an easy access to the square-root of any element of  $\mathcal{K}_1$ , and thus to the SOC Fischer–Burmeister function  $f_{\text{FB}}^{\mathcal{K}}$ .

**Newton algorithm** The  $f_{\text{FB}}^{\mathcal{K}}$  function is Lipschitz-continuous, as is the gradient of its squared norm,  $(\partial f_{\text{FB}}^{\mathcal{K}})^{\top} f_{\text{FB}}^{\mathcal{K}}$ , suggesting good conditions for a Newton minimization algorithm on  $\Psi_{\text{FB}} := \frac{1}{2} \|f_{\text{FB}}^{\mathcal{K}}(\hat{\mathbf{u}}(\mathbf{r}), \hat{\mathbf{r}})\|^2$ . Moreover,  $f_{\text{FB}}^{\mathcal{K}}$  is *strongly semismooth* (D. Sun and J. Sun 2005), justifying performing directly a nonsmooth root-finding Newton algorithm (Qi and J. Sun 1993) on  $f_{\text{FB}}^{\mathcal{K}}(\hat{\mathbf{u}}(\mathbf{r}), \hat{\mathbf{r}}) = \mathbf{0}$ , as this avoids computing the second derivative of  $f_{\text{FB}}^{\mathcal{K}}$ .

The Newton algorithm implemented in the bogus library is a simplification of the one that we suggested originally (Daviet, Bertails-Descoubes, and Boissieux 2011, Appendix A.2), and is given in Appendix B.1.2; the derivatives of  $f_{\text{FB}}^{\mathcal{K}}$  are provided in Appendix B.1.1.

We observed that on average, using the SOC Fischer–Burmeister complementarity function made the Gauss–Seidel algorithm perform slightly better than using the Alart–Curnier formulation (see Table 4.3). However, this was still not sufficient to achieve convergence for every local problem. Damping the Newton algorithms (by performing a line-search) improves their robustness, at the cost of being much slower in the general case, and still failing sometimes. More precisely, on our hair problems the undamped Fischer–Burmeister Newton (Algorithm B.1) was much more robust than the undamped Alart–Curnier Newton, and slightly less robust but much faster than the undamped Alart–Curnier Newton with line-search (Table 4.3 below).

Building failsafes with different combinations of root-finding algorithms on complementarity functions did not yield significantly better results. Instead, we have chosen to conceive a failsafe using a radically different approach, solving the one-contact problem analytically.

#### 4.1.2 Analytical solver

The idea of solving the one-contact friction problems analytically is not new. In 2D, it was mentioned by Klarbring (1990), and used inside a Gauss–Seidel algorithm by Jean (1999). In 3D however, this method appears to be surprisingly less popular, and we could not find any evidence of previous usage in numerical codes.

Analytical solvers attempt to solve problem (3.11) by trying to find a solution in each case of the disjunctive formulation (1.1). Here we will assume that  $W$  is symmetric positive-definite; if not, we can simply use the proximal variant (Section 3.4.2) of the Gauss–Seidel algorithm to alleviate this difficulty.

The enumerative algorithm proceeds as follow:

1. We first check the easiest case, take-off. If  $\mathbf{b}_N \geq 0$ , then it is sufficient to choose  $\mathbf{r} = \mathbf{0}$  and  $\mathbf{u} = \mathbf{b}$  for  $\mathbf{u}$  and  $\mathbf{r}$  to satisfy Coulomb’s law.
2. Then, we check the sticking case, i.e.,  $\mathbf{u} = \mathbf{0}$ . This means finding  $\mathbf{r} \in \mathcal{K}_\mu$  such that  $W\mathbf{r} + \mathbf{b} = \mathbf{0}$ . As  $W$  is positive-definite, this amounts to checking that  $W^{-1}\mathbf{b} \in -\mathcal{K}_\mu$ .
3. Finally, the sliding case is the hard one, as we will see in the next paragraph.

**Sliding case** We try to find a solution to the sliding case of the one-contact problem,

$$\begin{cases} \begin{pmatrix} 0 \\ -\alpha \mathbf{r}_T \end{pmatrix} = W\mathbf{r} + \mathbf{b} & \alpha \in \mathbb{R}_+ \\ \|\mathbf{r}_T\| = \mu \mathbf{r}_N. \end{cases} \quad (4.6)$$

In 2D, the second row of (4.6) imposes  $\mathbf{r}_T = \pm \mu \mathbf{r}_N$ . Plugging this into the first equation, we get  $(W_{NN} \pm \mu W_{TN})\mathbf{r}_N = -\mathbf{b}_N$ , and it suffices to check if any of the two possibilities yields a positive  $\alpha$ .

For the 3D case, we remark in (Bonnefon and Daviet 2011) that the null normal velocity condition  $\mathbf{u}_N = 0$  together with the cone boundary condition  $\|\mathbf{r}_T\| = \mu \mathbf{r}_N$  state that  $\mathbf{r}$  should lie in the intersection of a plane and a hollow cone, which is an ellipse. This ellipse can be parameterized by a polar angle  $\theta$ , and plugging-in the  $\mathbf{u}_T = -\alpha \mathbf{r}_T$  equation allows to reduce the potential values of  $\theta$  to the roots of a quartic polynomial, which can be enumeratively checked. Alternatively, we propose also in (Bonnefon and Daviet 2011) to look for the solutions to the sliding case in the roots of another degree-four polynomial on  $\alpha$ . This second polynomial may have a worse conditioning, but its coefficients are simpler to compute. Indeed, let us decompose  $W$  as

$$W := \begin{pmatrix} W_{NN} & W_{TN}^T \\ W_{TN} & W_{TT} \end{pmatrix}$$

We define

$$\begin{aligned} T &:= W_{TT} - \frac{1}{W_{NN}} W_{TN} W_{TN}^T & A &:= \text{Tr } T - W_{TN}^T \mathbf{b}_T & C &:= \det T - W_{TN}^T \mathbf{B} \\ b_T &:= \frac{\mathbf{b}_T}{\mathbf{b}_N} - \frac{W_{TN}}{W_{NN}} & \mathbf{B} &:= (\text{adj } T) \mathbf{b}_T & D &:= \left( \frac{W_{NN}}{\mu} \right)^2 \end{aligned}$$

where  $\text{adj } T$  designate the adjugate of  $T$ , that is, the transpose of its cofactor matrix. The solutions of the sliding problem (4.6) must then satisfy  $P(\alpha) = 0$ , with

$$P(\alpha) = \alpha^4 + 2A\alpha^3 + (2C + A^2 - D\mathbf{b}_T^2)\alpha^2 + 2(CA - D\mathbf{b}_T^T \mathbf{B})\alpha + C^2 - D\mathbf{B}^2. \quad (4.7)$$

**Convexified DCFP** The same approach can be used for the local problems of the convexified DCFP. The enumerative algorithm becomes:

1. (“take-off” case) If  $\mathbf{b} \in \mathcal{K}_{\frac{1}{\mu}}$ , then take  $\mathbf{r} = \mathbf{0}$ .
2. (“sticking” case) If  $W^{-1}\mathbf{b} \in -\mathcal{K}_{\mu}$ , take  $\mathbf{r} = -W^{-1}\mathbf{b}$ .
3. (“sliding” case) Check for solutions satisfying Equation (4.8),

$$\text{Bd } \mathcal{K}_{\frac{1}{\mu}} \ni W\mathbf{r} + \mathbf{b} \perp \mathbf{r} \in \text{Bd } \mathcal{K}_{\mu}. \quad (4.8)$$

Once again, we will be able to obtain 3D solutions by analyzing the roots of a quartic polynomial.

First, suppose that  $\mu$  is zero, then  $\mathbf{r} \in \text{Bd } \mathcal{K}_0 \implies \mathbf{r}_T = \mathbf{0}$  and therefore  $\mathbf{r} = \left[ -\frac{\mathbf{b}_N}{W_{NN}}; \mathbf{0} \right]$  is a solution since  $\mathbf{b}_N < 0$  (otherwise we would have  $\mathbf{b} \in \mathcal{K}_{\infty}$  and therefore be in the “take-off” case).

Now we assume that  $\mu > 0$ . We are looking for solutions where  $\mathbf{r}_N$  and  $\mathbf{u}_N$  are non zero, as  $\mathbf{u} = \mathbf{0}$  and  $\mathbf{r} = \mathbf{0}$  already correspond to “take-off” and “sliding” cases. Such solutions must satisfy

$$\begin{cases} \mathbf{u} &= W\mathbf{r} + \mathbf{b} \\ \|\mathbf{r}_T\| &= \mu \mathbf{r}_N \\ \|\mathbf{u}_T\| &= \frac{1}{\mu} \mathbf{u}_N \\ \exists \alpha > 0 & \text{s.t. } \mathbf{u}_T = -\alpha \mathbf{r}_T. \end{cases} \quad (4.9)$$



In the 2D case, this means  $\mathbf{r}_T = \pm \mu \mathbf{r}_N$ . We can successively look for a solution in both cases; with  $s := \text{sign}(\mathbf{r}_T)$ , this amounts to solving

$$W_{NN} + 2\mu s W_{TN} + \mu^2 W_{TT} \mathbf{r}_N = -(\mathbf{b}_N + \mu s \mathbf{b}_T) \quad (4.10)$$

and checking that the resulting  $\mathbf{r}_N$  and  $\mathbf{u}_N$  are both positive.

Now let us consider the 3D case, and write  $(\mathbf{e}_N, \mathbf{e}_{T_1}, \mathbf{e}_{T_2})$  the canonical basis of  $\mathbb{R}^3$ . Writing Equation (4.9) with cylindrical coordinates yields

$$\begin{cases} \mathbf{u} = W \mathbf{r} + \mathbf{b} \\ \mathbf{r} = r(\mathbf{e}_N + \mu(\cos \theta \mathbf{e}_{T_1} + \sin \theta \mathbf{e}_{T_2})) \\ \mathbf{u} = u(\mu \mathbf{e}_N - (\cos \theta \mathbf{e}_{T_1} + \sin \theta \mathbf{e}_{T_2})) \\ u > 0, \quad r > 0. \end{cases} \quad (4.11)$$

The first line of (4.11) reads

$$u \begin{pmatrix} \mu \\ \cos \theta \\ \sin \theta \end{pmatrix} - r \begin{pmatrix} \mu \cos \theta W_{NT_1} + \mu \sin \theta W_{NT_2} + W_{NN} \\ \mu \cos \theta W_{T_1 T_1} + \mu \sin \theta W_{T_1 T_2} + W_{NT_1} \\ \mu \cos \theta W_{T_1 T_2} + \mu \sin \theta W_{T_2 T_2} + W_{NT_2} \end{pmatrix} = \mathbf{b}. \quad (4.12)$$

We can eliminate  $u$  by plugging the first row into the second and third ones,

$$\begin{aligned} r \begin{pmatrix} -\cos^2 \theta W_{NT_1} + \cos \theta \sin \theta W_{NT_2} + (\cos \theta W_{T_1 T_1} + \sin \theta W_{T_1 T_2})\mu + \frac{\cos \theta W_{NN}}{\mu} + W_{NT_1} \\ -\cos \theta \sin \theta W_{NT_1} + \sin^2 \theta W_{NT_2} + (\cos \theta W_{T_1 T_2} + \sin \theta W_{T_2 T_2})\mu + \frac{\sin \theta W_{NN}}{\mu} + W_{NT_2} \end{pmatrix} \\ = \mathbf{b}_T + \frac{1}{\mu} \mathbf{b}_N \begin{pmatrix} \cos \theta \\ \sin \theta \end{pmatrix} \end{aligned}$$

then eliminate  $r$  by subtracting these two equations scaled by respectively  $\mathbf{b}_{T_2} + \frac{b_N \sin \theta}{\mu}$  and  $\mathbf{b}_{T_1} + \frac{b_N \cos \theta}{\mu}$ , to obtain the necessary condition (4.13),

$$\begin{aligned} 0 = & \mu(b_{T_1} W_{NT_1} - b_{T_2} W_{NT_2} - b_N W_{T_1 T_1} + b_N W_{T_2 T_2}) \cos \theta \sin \theta \\ & - ((b_{T_2} W_{T_1 T_2} - b_{T_1} W_{T_2 T_2})\mu^2 + b_{T_1} W_{NN} - b_N W_{NT_1}) \sin \theta \\ & - (b_{T_2} W_{NT_1} - b_N W_{T_1 T_2}) \cos^2 \theta \mu + (b_{T_1} W_{NT_2} - b_N W_{T_1 T_2}) \sin^2 \theta \mu \\ & - ((b_{T_2} W_{T_1 T_1} - b_{T_1} W_{T_1 T_2})\mu^2 + b_{T_2} W_{NN} - b_N W_{NT_2}) \cos \theta \\ & - (b_{T_2} W_{NT_1} - b_{T_1} W_{NT_2})\mu. \end{aligned} \quad (4.13)$$

Using the half-angle change of variable  $t = \tan \frac{\theta}{2}$ , we obtain a degree-4 polynomial in  $t$  whose roots yield a superset of the solutions to Equation (4.9). We then just need to look for any of the at most 4 roots that satisfies the original equation.

For the sake of completeness, we include below the expression of the coefficients of this polynomial, obtained thanks to the Sage<sup>1</sup> computer algebra system. Let

$$\begin{aligned} A &:= (b_{T_1} W_{T_1 T_2} - b_{T_2} W_{T_1 T_1})\mu^2 + b_N W_{N, T_2} - b_{T_2} W_{N, N} \\ B &:= (b_N W_{T_1 T_2} + b_{T_1} W_{N, T_2} - b_{T_2} W_{N, T_1})\mu \\ C &:= 2((b_{T_1} W_{T_2 T_2} - b_{T_2} W_{T_1 T_2})\mu^2 - b_N W_{N, T_1} + b_{T_1} W_{N, N}) \\ D &:= 2(b_N (W_{T_1 T_1} - W_{T_2 T_2}) - b_{T_1} W_{N, T_1} + b_{T_2} W_{N, T_2})\mu \\ E &:= -6(b_N W_{T_1 T_2} - b_{T_1} W_{N, T_2})\mu, \end{aligned}$$

then Equation (4.13) becomes equivalent to  $P(t) = 0$  with

$$P(t) := (B - A)t^4 + (C + D)t^3 + Et^2 + (C - D)t + A + B. \quad (4.14)$$

<sup>1</sup><https://www.sagemath.org>



**Quartic solvers** We need an efficient and precise way of finding the roots of the quartic polynomials from the sliding cases, (4.7) and (4.14). There exists a few analytical algorithms for finding the roots of a quartic polynomial, the first one being attributed to Lodovico Ferrari in 1540. Such algorithms are complex, as they involve accounting for a lot of potential cases; luckily, open-source implementations have been written, for instance in the GNU Scientific Library (GSL)<sup>2</sup>. If a perfectly analytical solution is not required, one may also use a general purpose polynomial root-finding algorithm, such as the famous Netlib routine `rpoly`<sup>3</sup>, or compute the eigenvalues of the companion matrix. Despite performing more floating point operations, these latter solutions are often faster in practice, thanks to a lower number of conditionals.

**Hybrid local solver** Using our new analytical solver inside the Gauss–Seidel algorithm proved to be quite robust, but once again did not fully eliminate the simulation failures (see Table 4.3). The main drawback of the enumerative solver is that, in cases where no solution is found, no approximate solution is computed either, as would be the case for iterative solvers. Even if a problem is very close to have a solution (say, for instance that  $\mathbf{b}_N$  is slightly below 0 due to numerical errors and no solution is found in either the sliding or sticking case), the analytical solver will simply state that no solution exists. In contrast, in this same scenario, a Newton solver would simply return a zero force vector and yield a very small residual. One solution could be to add tolerances when checking for each case in the enumerative algorithm; but that would just amount to slightly postponing the difficulty. We found that one more satisfying solution was to combine the two different approaches, enumerative and iterative, thus benefitting from the best of both strategies.

In (Daviet, Bertails-Descoubes, and Boissieux 2011), we computed the quartic’s roots analytically using the GSL’s algorithm. As this was quite slow, we used the enumerative solver only as a failsafe for the Fischer–Burmeister Newton algorithm. In *bogus*, a slightly different approach is used: the enumerative solver is first used, computing the quartic’s roots as the eigenvalues of the polynomial’s companion matrix. For this, we use Eigen’s built-in solver (Guennebaud, Jacob, et al. 2010), which proved to be quite fast. Then, the Newton solver is unconditionally applied. The computed eigenvalues might not be extremely precise, but will still be close enough to the solution so that the Newton algorithm will only need a very low number of iterations (usually 0 or 1) to bring the error below the required tolerance. If no solution is found by the enumerative algorithm, then the Newton algorithm is simply started from the previous iteration solution.

Both approaches work well in practice, proving themselves both more robust and more efficient than any non-hybrid solution. From a practical point of view, however, the second strategy may be more convenient to implement.

#### 4.1.3 Full algorithm

Now that we have a robust and efficient one-contact solver, let us focus on the global Gauss–Seidel algorithm. At each iteration, the canonical algorithm solves the local problem for each contact. However, we observe that the speed of convergence of the local forces is highly inhomogeneous; some contacts will reach their final value after a few iterations, with others will require dozens of them. For this reason, we would rather focus our computing power on the contacts that are slow to converge.

**Sleeping heuristics** Radjai et al. (1998) observe that the intensity of the contact forces inside a granular media follows a bimodal distribution. They describe the emergence of *force chains*, that is, a percolating load-bearing structure consisting of sticking contacts with high forces, and local dissipative regions where the forces are much smaller and where sliding occurs. Using this insight, we want to focus on getting the force chains right; the dissipative regions, being local and having less influence on the force network, should require less iterations to converge. We

<sup>2</sup><https://www.gnu.org/software/gsl/>

<sup>3</sup><http://netlib.sandia.gov/port/prop.upd/rpoly.f>



**Figure 4.1:** Comparison of the hair collective behavior between (top) real hair motion sequences and (bottom) our corresponding simulations, based on large assemblies of (up to 2,000) individual fibers with contacts and Coulomb friction.

propose a *contact “sleeping”* heuristic, which temporarily stops solving contacts at which the force remains roughly constant or is small. The contacts are then reactivated after a few iterations, 10 in our implementation. This simple heuristic leads to significant performance improvements (Figure 4.8), and we can show that it does not affect the theoretical convergence results of the Gauss–Seidel applied to the minimization of a convex function (Appendix B.3).

**Stopping criterion** The standard stopping criterion for Gauss–Seidel algorithms is to look at the difference in forces (or velocities) between one step and the next, and exit when this difference is small enough. Yet, we did not find this criterion to be always representative of the quality of the current iterate, as small steps could simply indicate failed local solves. We choose instead to make use of our Fischer–Burmeister complementarity function, and to evaluate the error as the maximum of the  $f_{\text{FB}}$  norm over all contact points. However, this requires computing the current relative velocity as  $\mathbf{u} = W\mathbf{r} + \mathbf{b}$ . This costs one matrix–vector operation, which is as much as the cost of computing  $\bar{\mathbf{b}}_i$  for the whole set of contacts. Using our error measurement at each Gauss–Seidel iteration would therefore mean almost doubling the total cost of our algorithm, which is unacceptable. Instead, we evaluate the residual only every  $N$  iterations ( $N = 25$  in our case), so that the overhead becomes negligible.

**Final algorithm** Our final algorithm for solving the DCFP or convexified DCFP (more specifically the matrix-free version, which can be easily adapted for an Delassus operator  $W$ ) is laid out in Algorithm 4.1; the exact implementation can be found in the bogus library.

## 4.2 Application to hair dynamics

Modeling hair dynamics is challenging: human hair is composed of about 150,000 individual fibers that tightly interact together, leading to a complex collective behavior. Due to the rough surface of the hair fibers, covered with microscopic scales, dry friction is substantial at the contact points, and consequently greatly influences the hair dynamics at the *macroscopic* level. As illustrated in Figure 4.2, we have identified in (Daviet, Bertails-Descoubes, and Boissieux 2011) three major hair visual features that directly emerge from those nonsmooth frictional contacts:

1. Stick-slip instabilities during motion;
2. Spontaneous splitting of hair into multiple untidy wisps and “flyers” during strong motion vs. its spontaneous grouping into a few globally coherent locks during gentle motion;
3. Appearance of complex hair patterns that can remain perfectly still at the end of the motion.

We claim that accounting for these phenomena is essential for producing realistic and compelling hair animations.

---

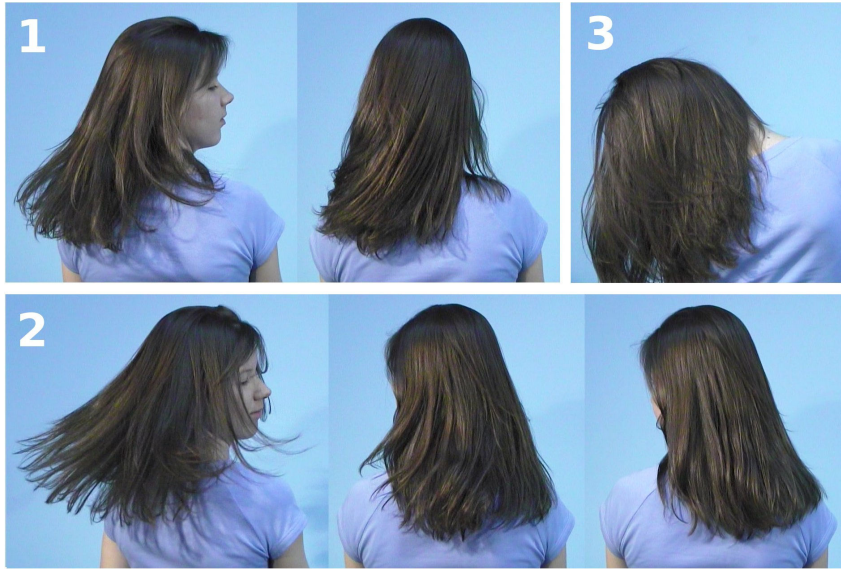
**Algorithm 4.1:** Matrix-free Gauss–Seidel algorithm with sleeping heuristics

```

input : Initial guess  $\mathbf{r}$ 
input : Matrices  $M$ ,  $H$  and  $P$ 
input : Vectors  $\mathbf{b}$  and  $\mathbf{c}$ 
input :  $\beta$  array of proximal coefficients
Result:  $\mathbf{r}$  and  $\lambda$  approximate solutions to the DCFP
 $\mathbf{l} \leftarrow (HM^{-1}C^\top)P^{-1}(\mathbf{c} - (HM^{-1}C^\top)^\top \mathbf{r})$ ; // Linear constraints
for  $i=1$  to  $\text{maxIters}$  do
    | SkipTab [ $i$ ]  $\leftarrow 0$ ; // Reset sleeping contacts counter
    |  $W_{ii} \leftarrow \sum_j^m H_{ij}M_{jj}H_{ij}^\top + \beta_i \mathbb{I}$ ; // Precompute local matrices
end
 $\mathbf{z} \leftarrow M^{-1}H^\top \mathbf{r}$ ;
for  $k \leftarrow 1$  to  $\text{maxIters}$  do // GS iteration
    | for  $i \leftarrow 1$  to  $n$  do // Loop over contacts
        | if SkipTab [ $i$ ]  $> 0$  then // Sleeping contact
            | SkipTab [ $i$ ]  $\leftarrow$  SkipTab [ $i$ ]  $- 1$ ;
            | continue;
        | end
        |  $\bar{\mathbf{b}}_i \leftarrow \mathbf{b}_i + \mathbf{l}_i + H.\text{row}(i) \mathbf{z} - W_{ii} \mathbf{r}_i$ ;
        |  $\mathbf{r}_i^{\text{prev}} \leftarrow \mathbf{r}_i$ ;
        |  $\text{solutionFound} \leftarrow \text{solveLocalProblem}(W_{ii}, \bar{\mathbf{b}}_i, \mathbf{r}_i)$ ;
        | if not  $\text{solutionFound}$  then  $\mathbf{r} \leftarrow \frac{1}{2}(\mathbf{r}_i + \mathbf{r}_i^{\text{prev}})$ ;
        |  $\mathbf{z} \leftarrow \mathbf{z} + (M^{-1}H).\text{col}(i) (\mathbf{r}_i - \mathbf{r}_i^{\text{prev}})$ ;
        | if  $\|\mathbf{r}_i\|^2 < \epsilon_{\text{small}}$  or  $\|\mathbf{r}_i - \mathbf{r}_i^{\text{prev}}\|^2 < \epsilon_{\text{converged}}$  then
            | SkipTab [ $i$ ]  $\leftarrow n\text{Skip}$ ; // Put contact to sleep
        | end
    | end
    |  $\mathbf{l} \leftarrow (HM^{-1}C^\top)P^{-1}(\mathbf{c} - (HM^{-1}C^\top)^\top \mathbf{r})$ ; // Linear constraints
    | if  $k \bmod n\text{EvalEvery} = 0$  then
        | // Evaluate residual
        |  $\mathbf{z} \leftarrow M^{-1}H^\top \mathbf{r}$ ;
        |  $\mathbf{u} \leftarrow H\mathbf{z} + \mathbf{b} + \mathbf{l}$ ;
        | if  $\max_i \|f_{FB}(\mathbf{u}_i, \mathbf{r}_i)\|^2 < \text{tol}$  then break;
    | end
end

```

---



**Figure 4.2:** Three important features of the real hair collective behavior, emerging from nonsmooth friction: (1) Stick-slip instabilities, especially visible here between the hair and the right shoulder; (2) Spontaneous splitting of hair into thin wisps and “flyers” during strong motion (left) vs. spontaneous grouping of hair into a few coherent locks during gentle motion (middle and right); (3) Complex hair patterns that may remain stable at rest. Courtesy of Sylvain Paris and Tilke Judd, MIT 2007.

#### 4.2.1 Hair simulation in Computer Graphics

Due to performance limitations, the first works that attempted to simulate hair dynamics mostly neglected hair self-interactions (Anjyo et al. 1992; Rosenblum et al. 1991), or processed them at a coarse level between a small number of predefined interacting wisps using penalty forces (Bertails et al. 2006; Choe et al. 2005; Plante et al. 2001). Alternatively, Hadap and Magnenat-Thalmann (2001) proposed a macroscopic model of the hair medium based on a fluid solver, observing that some fluid properties, such as incompressibility, could be representative of the hair collective behavior. Though interesting, their approach fails to capture the discontinuities emerging from large hair motions. Thanks to the design of realistic, robust and fast primitives for thin elastic rods (Bergou, Audoly, et al. 2010; Bergou, Wardetzky, et al. 2008; Bertails et al. 2006; Hadap 2006; Pai 2002; Selle et al. 2008; Spillmann and Teschner 2007), some approaches have been developed at the fiber lever in order to gain realism in hair simulations.

Selle et al. (2008) designed an efficient mass-spring model for an individual fiber, allowing them to simulate up to 10,000 fibers in a reasonable computational time (from a few minutes up to one hour per frame) on a quad-core architecture. Unfortunately, many self-contacts were ignored, causing the fibers to penetrate each other, and thus failing to preserve the hair volume. To resolve these issues while still retaining some discontinuous details in the simulations, McAdams et al. (2009) proposed a hybrid Eulerian/Lagrangian hair model combining a fluid model together with the explicit treatment of fiber self-contacts. With this approach, the hair volume is properly preserved while detailed interactions at the fiber level yield nice visual effects. However, nonsmooth effects due to dry friction, which play a major role in hair dynamics, are not captured.

In (Daviet, Bertails-Descoubes, and Boissieux 2011), we used the Super-Helix (Bertails et al. 2006) fiber model, which is a reduced-coordinate discretization of a Kirchhoff rod whose sole degrees of freedom are piecewise-constant curvatures and twists. This model possesses two main

Name	$N_{\text{rods}}$	$\mu_{\text{hair}}/\mu_{\text{body}}$	$\overline{N_{\text{contacts}}}$	$\max(N_{\text{contacts}})$
Free/A	1920	0.3 / 0.5	24659	46915
Free/B	1920	0.2 / 0.3	28153	36287
Pony	334	0.3 / 0.5	9850	16613
Curly	1920	0.3 / 0.5	21425	33578

Overlines indicate averaged quantities.

**Table 4.1:** *Physical properties of our hair simulations.*

Name	$\%_{>\text{tol}}^1$	$\max(\text{err})$	$\overline{\text{iters}}$	$\overline{T_{\text{GS}}} / \overline{T_{\text{solv}}}^2$ (s)	$\overline{T_{dt}}^3$ (s)	$\overline{T_{\text{frame}}}$ (min)	$T_{\text{tot}}$ (hours)
Free/A	0.056	0.004	136	2.60 / 4.06	7.15	2.09	25
Free/B	0.13	0.037	160	5.3 / 7.0	11.1	3.02	36
Pony	0.42	0.003	301	4.43 / 5.25	6.2	2.47	30
Curly	0.013	0.003	118	2.26 / 3.39	7.54	4.15	49

Overlines indicate averaged quantities.

<sup>1</sup>  $\%_{>\text{tol}}$ : Percentage of one-step problems that did not reach (global) tolerance

<sup>2</sup>  $\overline{T_{\text{GS}}} / \overline{T_{\text{solv}}}$ : Time for the Gauss–Seidel loop (Algorithm 4.1) alone / Total time for the contact solver (including the assembly of  $W$ )

<sup>3</sup>  $\overline{T_{dt}}$ : Total time for one simulation timestep ( $dt$  ranges from 1 to 4ms)

**Table 4.2:** *Performance results for our hair simulations.*

advantages: the inextensibility of the fibers is intrinsically enforced by the choice of coordinates without requiring a supplemental holonomic constraint nor stiff springs, and the bending forces are linear in the degrees of freedom, allowing for easy implicit integration. Drawbacks include having to use a limited number of degrees of freedom, as the stiffness matrices of the fibers are dense and expensive to compute, and vanishing inertia terms in straight configurations, leading to potentially degenerate systems. We treated contacts using the Moreau-Jean timestepping scheme (Section 2.2.2) and our Gauss–Seidel solver (Section 4.1). This allowed us to simulate a few thousand individual fibers at a few minutes per frame, capturing the variety of effects induced by static friction.

More recently, Iben et al. (2013) used penalties for hair–hair contacts as in (Selle et al. 2008), but did not limit the maximum number of contacts and used a contact-pruning algorithm specifically designed to maintain the volume of stylized hair. Other authors have also leveraged our hybrid solver. Aubry and Xian (2015) proposed to extend the Discrete Elastic Rods (DER; Bergou, Audoly, et al. 2010; Bergou, Wardetzky, et al. 2008) model with root constraints to perform implicit simulation of flexible trees, leveraging our hybrid Gauss–Seidel algorithm to deal with frictional contacts. Kaufman, Tamstorf, et al. (2014) devised a Newton algorithm to adaptively handle the nonlinearities of the DER model. Thanks to a matrix-free variant of our Gauss–Seidel algorithm, they were able to simulate scenes with tens of thousands of fibers and up to a million contact points per timestep. In a similar manner, Gornowicz and Borac (2015) used a Newton algorithm over a variation of the Discrete Elastic Rod model with a modified elastic energy, but with a linear approximation of the friction cone.

#### 4.2.2 Full-scale simulations

To evaluate the effectiveness of our contact solver (Algorithm 4.1) on hair simulations, we ran three kinds of experiments, summarized in Table 4.1. All are using the Super-Helix model, with 16 degrees of freedom per rod. The character was animated using 3ds Max (Autodesk 2009) by reproducing a real video-captured motion that serves as a reference (see Figure 4.1 and the accompanying video<sup>4</sup>). The hair simulation entitled “Free” models a full, unconstrained haircut, and consists of about 2000 simulated rods. It is divided into two parts, “A” and “B”, featuring a

<sup>4</sup><http://bipop.inrialpes.fr/~bertails/Papiers/Videos/hairContactSiggraphAsia2011.mp4>



head rotation motion and a head leaning motion, respectively. The third hair simulation, entitled “Pony”, contains only slightly over 300 rods, but those are tightly packed into a ponytail. All of these experiments include smooth as well as rough head motions. Finally, a last motion, “Curly”, illustrates the fact that our method can also easily handle curly hair (see Figure 4.5). It is based on the “Free/A” head motion.



**Figure 4.3:** *Simulation of a fast head movement without (top) and with (bottom) Coulomb friction. In the latter case, hair remains much more coherent and the results are visually more appealing.*

Visual results and comparisons to real hair motions are presented in Figures 4.1, 4.3, 4.4 and 4.5. Final rendering was performed using 3ds Max. Our method allows us to preserve the hair volume and to capture subtle phenomena such as stick-slip instabilities or the spontaneous appearance of transient coherent movements in hair. We also observe that a lot of energy is dissipated by Coulomb friction; capturing it accurately is essential to achieve realism. Figure 4.3 shows the effect of decreasing the hair/hair friction coefficient on the collective hair behavior. Without friction, hair looks artificially clean and light. In contrast, in the presence of friction, motion looks coherent and properly damped, while simultaneously featuring complex details at the fiber level.

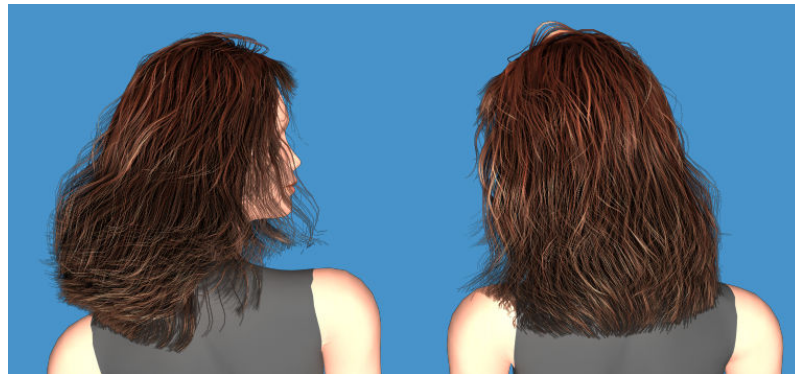
**Measure of performance and robustness** All our simulations were run on a desktop machine featuring a Intel<sup>®</sup> Xeon<sup>®</sup> W3520 processor with 8 GB of memory. Numerical results for the three experiments described above are given in Table 4.2. The tolerance for the global Gauss–Seidel was chosen so that we do not observe any visual disturbance. In our internal units, the tolerance ranges from  $10^{-4}$  when the motion is fast to  $10^{-6}$  in long static phases<sup>5</sup>. Though the solver sometimes failed to strictly reach the requested precision (in less than 2% of the cases), large errors never occur: the solver always gives an approximate solution from which the simulation can go on without exhibiting any artifact. To improve the convergence of the global algorithm, the tolerance for the local solvers is always set to a lower value, typically  $10^{-7}$ .

Each simulation ran at a rate of a few minutes per frame, that is, 25 seconds of video in about 48 hours.

<sup>5</sup>This roughly translates into an average dimensionless relative error on  $\mathbf{r}$  ranging from  $10^{-3}$  to  $10^{-4}$ .



**Figure 4.4:** *Our simulations (right) capture the complex patterns emerging from static friction in real hair (left).*



**Figure 4.5:** *Our method can also handle curly hair. Note how the volume of the hairstyle is preserved throughout the simulation, and how static friction is properly captured near the top of the head.*

#### 4.2.3 Friction solvers comparisons

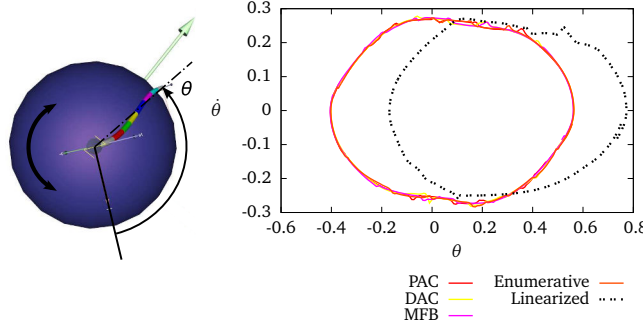
We compared our method against a variety of frictional contact solvers,

- **MFB:** Our Newton method based on the modified Fischer-Burmeister formulation. It can be used either as a global solver or as a local solver within a Gauss–Seidel loop.
- **PAC** and **DAC:** Respectively pure and damped Newton methods based on the Alart-Curnier formulation used in (Bertails-Descoubes et al. 2011). Both can be used as either global or local solvers.
- **Duriez08:** Local solver from Duriez (2008), which consists of successive iterations of the approach from Duriez et al. (2006) embedded in a fixed point loop. Unlike the original

approach, it solves the exact Coulomb friction problem when convergence is achieved.

- **4sides**: The local solver from Otaduy et al. (2009), which approximates the Coulomb friction cone with a four-sided pyramid.
- **Enum**: Our quartic enumerative local solver from (Bonnefon and Daviet 2011).

**Exact vs faceted Coulomb friction** To illustrate the influence of the choice of the local friction formulation, we created a very simple experiment: the free end of a fiber is dropped on a ball that is rotating with sinusoidal oscillations, and a non-zero friction coefficient ( $\mu = 1$ ) is set between the ball and the fiber. Figure 4.6 (top) shows the phase plots of the free end once it has reached its periodic orbit.



**Figure 4.6:** *Exact vs approximate model for Coulomb friction.*

*Periodic orbit of the free end of a rod resting on top of a rotating sphere. Frictional contact is simulated using different solvers. Solvers that model exact Coulomb friction all reach the same orbit, no matter the choice of the error function. The linearized model (dashed lines) reaches a completely different one.*

We tried two global Newton methods, **DAC** and **MFB**, and five local solvers within a Gauss–Seidel loop: **DAC** and **MFB** again, as well as **Duriez08**, **4sides** and **Enum**. For **DAC** and **MFB** based solvers, the stopping criterion relied on the norm of the corresponding objective functions, while the other methods were stopped as soon as the step size between two Gauss–Seidel iterations got below 5%.

No matter the choice of the stopping criterion, all solvers that model exact Coulomb friction converged towards the same orbit. For the linearized local solver, (**4sides**), the Gauss–Seidel algorithm still converged quickly to a fixed point, but the trajectory differed substantially — even though as suggested by Otaduy et al. (2009), we aligned the friction pyramid with the unconstrained tangential velocity.

**Performance** We could not get the full-scale simulations to complete with any local solver other than our hybrid method **MFB+Enum**. All other solvers led to the divergence of the fiber model at some point of the simulation. Still, to quantitatively evaluate the performance of our solver, we saved about three hundred one-step problems from our smallest simulation, “Pony”. We then successively ran on this benchmark all the local solvers mentioned above, with or without the help of the enumerative solver.

Convergence plots are shown in Figure 4.7 and numerical results are given in Table 4.3. We used the Fischer-Burmeister error measure for all our tests, except for the pure **4sides** solver which does not attempt to solve exact Coulomb friction<sup>6</sup>.

<sup>6</sup>In this case, we resorted to a measure of the iterates length. To avoid introducing any bias, we made sure that the tolerance was such that the timings of **MFB+Enum** were roughly the same under both error measures.



Local solver	% <sub>alt</sub> <sup>1</sup>	% <sub>fail</sub> <sup>2</sup>	% <sub>&gt;tol</sub> <sup>3</sup>	iters <sup>4</sup>	$\overline{T}_{GS}$ <sup>5</sup> (ms)
4sides	–	0	0	575	6497
Duriez08	–	2.45	50	267	4336
PAC	–	0.43	19.3	163	1265
DAC	–	0.26	0.33	60	874
MFB	–	0.13	4.9	72	484
Enum	–	0.0005	1	67	1044
4sides + Enum	38.7	0.001	0	47	763
Duriez08 + Enum	51	0.0007	0.33	90	1447
PAC + Enum	0.1	10 <sup>−5</sup>	0	62	543
DAC + Enum	0.09	10 <sup>−6</sup>	0	57	789
<b>MFB + Enum</b>	<b>0.07</b>	<b>10<sup>−6</sup></b>	<b>0</b>	<b>41</b>	<b>312</b>

<sup>1</sup> %<sub>alt</sub>: Percentage of calls to fail-safe

<sup>2</sup> %<sub>fail</sub>: Percentage of local problems that did not reach tolerance

<sup>3</sup> %<sub>>tol</sub>: Percentage of one-step problems that were not solved to (global) tolerance (tol = 10<sup>−6</sup> except for **4sides**: 5 × 10<sup>−2</sup>)

<sup>4</sup> iters: Mean number of Gauss–Seidel iterations

<sup>5</sup>  $\overline{T}_{GS}$ : Mean time in Gauss–Seidel algorithm

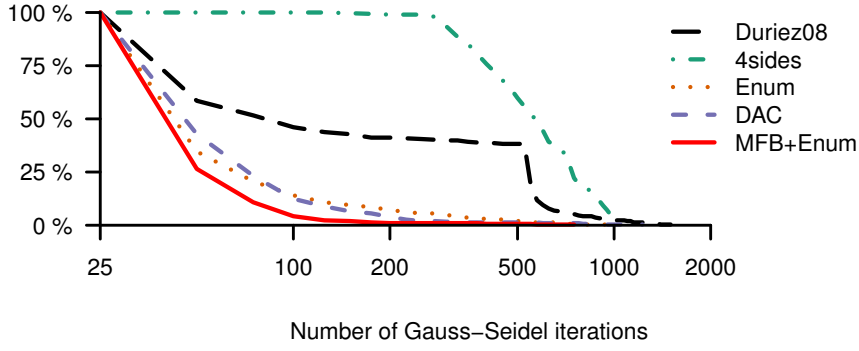
**Table 4.3:** Performance comparison of various local solvers on a set of 306 one-step problems.

From these numerical results we first note that using a linearized cone does not necessarily bring better time performance, despite a lower cost per call to the local solver. The **4sides** method was actually the one that required the highest number of Gauss–Seidel iterations to reach sufficient accuracy. We also observed that the contact freezing policy was not of much help to the **4sides** solver which, as a result, had to process a higher number of local problems per Gauss–Seidel iteration than exact friction solvers. Using a finer Coulomb friction model thus does *not necessarily* imply more costly simulations, quite the reverse in our case.

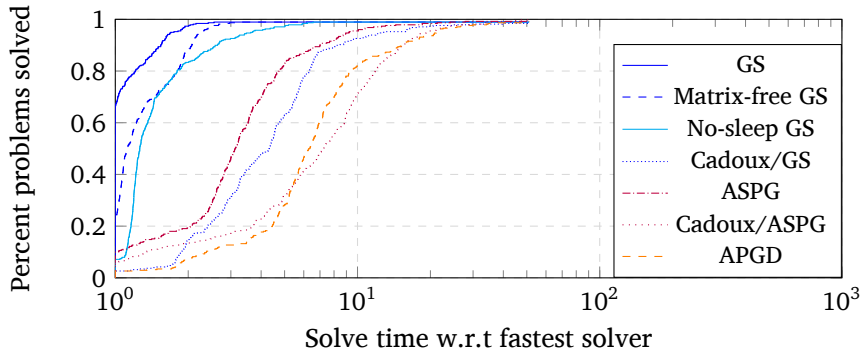
A second interesting point deals with the analysis of the role played by the **Enum** solver. For all the local solvers we tested, using the **Enum** solver as a fail-safe improved both the success rate and the computation time. In such a configuration, it turns out to be unnecessary to use a very robust primary solver. Indeed, the **MFB**, the **PAC** and more surprisingly the **4sides** method all outperformed the **DAC** solver. We also noted that the **Enum** solver requires a large number of iterations to reach the global tolerance, despite a very low rate of unsolved local problems. This is due to problems that do not admit an analytical solution, but for which an approximate solution with a low numerical error can still be found. As mentioned previously, while optimization-based solvers manage to find such acceptable solutions, the **Enum** solver remains stuck, thus spoiling the global convergence of the Gauss–Seidel algorithm. When using the **Enum** solver as a fail-safe, the rate of remaining problems without solution fortunately becomes very low, and in these rare cases, simply resetting the force to zero proved to be sufficient.

Overall, these results confirm our claim that our hybrid method **MFB+Enum** is both robust *and* efficient. The **MFB** solver alone turns out to be both faster and more robust than **PAC**, and, although slightly less robust, far much faster than **DAC**. The robustness issue becomes irrelevant when the solvers are combined with **Enum**.

**Comparison of global solvers** Finally, we compared our hybrid Gauss–Seidel solver to other approaches from Chapter 3 on a set of 378 DCFP coming from hair simulations with different fiber models (Super-Helices and Discrete Elastic Rods), and number of contacts ranging from 7000 to 27000. The results of this benchmark, shown on Figure 4.8, crowned Algorithm 4.1 as more efficient and robust than all the other solvers, and confirmed the speed-gain brought by the sleeping heuristics. Note also that computing the explicit Delassus operator proved to be more efficient than using the matrix-free algorithm. We also compared with two variants of the projected gradient descent algorithm, **APGD** from Heyn (2013), and **ASPG** adapted from Tasora



**Figure 4.7:** Percentage of one-step problems (Y axis) requiring more than a given number of Gauss-Seidel iterations (X axis) to converge, for various local solvers. Our hybrid method (solid red) induces better global convergence than previous approaches.



**Figure 4.8:** Performance profiles for a variety of solvers on our hair simulation benchmark — percentage of problems that each solver was able to solve under a multiple  $x$  of the time taken by the fastest solver for this problem. The line  $x = 1$  gives the solvers that were the most often the fastest, while  $y = 1$  shows the most robust ones.

(2013) as per Algorithm B.2. Finally, **GS** and **ASPG** were also compared to their SOCQP variants inside a Cadoux fixed-point loop.

#### 4.2.4 Limitations

Even though our hair simulations appear more convincing than those generated by previous methods, they do not perfectly match the real movement depicted in the reference videos. This is partly due to inaccuracies in the hair styling process and to the difficulty of precisely identifying the true physical parameters. The next section will present a first step towards reducing this inaccuracy by deducing the groom rest parameters from an input geometry.

Another reason is related to the actual complexity of hair interactions. While Coulomb friction is a key ingredient, further effects such as anisotropic friction, adhesion, or electrostatics are likely to influence the hair motion, depending on the hair state (clean, dirty, wet) and external conditions (dry or wet atmosphere). Moreover, air-hair friction becomes preponderant during energetic motions, and a full simulation the surrounding air would be necessary to properly capture the intricate dynamics of this tightly coupled system.



**Figure 4.9:** *Our nodal algorithm efficiently handles the frictional contacts between the mannequin and her dress during a walking cycle.*

### 4.3 Application to cloth simulation

In this section, we depart slightly from the framework of Discrete Element Models and consider instead the simulation of a single discretized piece of cloth. Several numerical models exist for unconstrained cloth, which we will not detail here; we used an implementation of the Discrete Shell model (Grinspun et al. 2003) written by R. Casati (2015).

Treatment of cloth contacts in the Computer Graphics community typically rely upon the timestepping framework of Bridson et al. (2002), as in (Brochu et al. 2012; Harmon, Vouga, Tamstorf, et al. 2008), or event-driven asynchronous integrators (Ainsley et al. 2012; Harmon, Vouga, Smith, et al. 2009). These approaches attempt to resolve contacts by moving cloth vertices, but do not consider the increase in elastic energy associated to such displacements; friction is also treated in an explicit manner. This negatively impact the stability of the simulations.

Conversely, naive implicit treatment of contacts in cloth simulations leads to performance issues, as we explain below. Most numerical cloth models (such as the Discrete Shells) use a set of vertices as their degrees of freedom; moving a single vertex will affect the internal forces only at surrounding vertices, so the stiffness matrix will be sparse. However, applying a force anywhere will induce a displacement of each of the cloth vertices; the inverse of the stiffness matrix is dense. This means that the Delassus operator will be dense as well, and our usual hybrid Gauss–Seidel will become very inefficient. We could still use other optimization-based algorithms that do not require explicitly computing the Delassus operator, such as interior points, projected gradient descent, or ADMM, see Chapter 3). Or, as proposed by Otaduy et al. (2009), embed the friction solve inside another iterative splitting algorithm, and compute the Delassus operator using only the block-diagonal part of the stiffness matrix.

In (Daviot, Bertails-Descoubes, and Casati 2015), we propose to use a different approach and present a first step towards a more efficient and physically accurate treatment of frictional contacts in cloth simulations. Note that this method suffer from severe limitations, making it totally unsuitable for real-world usage; indeed, we consider only contacts between cloth vertices and external objects (such as the character body), and assume that there is at most one contact per vertex. However, the simplicity of the formulation might be of interest, and could potentially be extended to more useful settings.

#### 4.3.1 Nodal algorithm

Consider the DCFP for a 3D system whose degrees of freedom are its  $m$  vertices, and  $n$  contacts with external objects occurring exactly at the vertices. The relative velocities are simply given by

$$\mathbf{u}_i = E_i^T(\mathbf{v}_j - \mathbf{w}_i^{ext}), \quad 1 \leq i \leq n, \quad 1 \leq j \leq m$$

where  $j$  is the index of the contacting vertex at the  $i^{\text{th}}$  contact point,  $\mathbf{w}_i^{ext}$  the velocity of the external object and  $E_i := (\mathbf{n}_i, \mathbf{t}_i, \mathbf{n}_i \wedge \mathbf{t}_i)$  a rotation matrix transforming the local contact basis coordinates into the world coordinates.

The matrix  $H$  then possesses a very simple structure: there is only one non-zero block per contact,  $H_{i,j} = E_i^T$ .  $H$  is non invertible; however, assuming that there is at most one contact per

vertex, we can easily extend it into a square  $3m \times 3m$  block-diagonal matrix  $G$ ,

$$G_{i,i} := \begin{cases} E_i^\top & \text{if vertex } i \text{ is in contact} \\ \mathbb{I} & \text{otherwise.} \end{cases} \quad (4.15)$$

The DCFP can then be equivalently written using this matrix  $G$  as finding  $\mathbf{v}, \mathbf{u}, \mathbf{r}$ , all in  $\mathbb{R}^{3m}$ , such that

$$\begin{cases} M\mathbf{v} = \mathbf{f} + G^\top \mathbf{r} \\ \mathbf{u} = G\mathbf{v} - G\mathbf{w}^{ext} \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} & \text{if vertex } i \text{ is in contact} \\ \mathbf{r}_i = \mathbf{0} & \text{otherwise.} \end{cases} \quad (4.16)$$

Now,  $G$  is orthonormal, and can be easily inverted as  $G^{-1} = G^\top$ . We can thus eliminate  $\mathbf{v}$  in system (4.16) and write its first line as

$$MG^\top(\mathbf{u} + \mathbf{w}^{ext}) = \mathbf{f} + G^\top \mathbf{r}.$$

Multiplying both sides by  $G$ , we obtain system (4.17)

$$\begin{cases} \underbrace{GMG^\top}_{\tilde{W}} \mathbf{u} + \underbrace{G(MG^\top \mathbf{w}^{ext} - \mathbf{f})}_{\tilde{\mathbf{b}}} = \mathbf{r} \\ (\mathbf{u}_i, \mathbf{r}_i) \in \mathcal{C}_{\mu_i} & \text{if vertex } i \text{ is in contact} \\ \mathbf{r}_i = \mathbf{0} & \text{otherwise.} \end{cases} \quad (4.17)$$

System (4.17) is very similar to the dual formulation of the DCFP (2.19), where the role of  $\mathbf{u}$  and  $\mathbf{r}$  have been reversed. However, the matrix  $\tilde{W}$  is now sparse, and we can therefore solve (4.17) using our Gauss–Seidel algorithm by either

- transforming (4.17) as a sequence of convexified DCFP, using a variation of the Cadoux fixed-point algorithm (see Daviet, Bertails-Descoubes, and Casati 2015), or;
- if  $\tilde{W}_{ii}$  is invertible, writing the inverse local linear relationship,  $\mathbf{u}_i = \tilde{W}_{ii}^{-1}(\mathbf{r}_i - \tilde{\mathbf{b}}_i)$  and using our original one-contact solvers, or;
- devising a new hybrid local solver for Coulomb friction with reversed linear relationship between the force and relative velocity. The Fischer–Burmeister Newton algorithm can be trivially adapted, and a corresponding quartic polynomial has been derived by a student in our group, L. Toran — but has not been published yet.

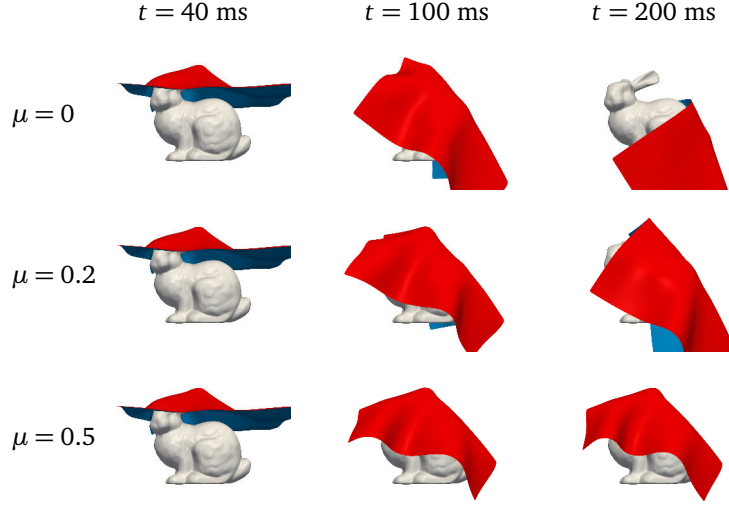
One can also take advantage of the orthonormality of the matrix  $G$  by performing a projected gradient descent algorithm on the primal formulation. Indeed, the feasible set of each intermediate primal SOCQP of the Cadoux algorithm is

$$C(\mathbf{s}) = \{\mathbf{v} \in \mathbb{R}^{3m}, G\mathbf{v} + \mathbf{w}^{ext} + \mathbf{s} \in K := K_1 \times \dots \times K_m\}$$

$$K_i := \begin{cases} \mathcal{K}_{\mu_i}^\perp & \text{if } i \text{ is in contact} \\ \mathbb{R}^d & \text{otherwise} \end{cases} \quad \forall 1 \leq i \leq m$$

and the projection on  $C(\mathbf{s})$  is now easy to compute:

$$\begin{aligned} \mathbf{y} = \Pi_{C(\mathbf{s})}(\mathbf{v}) &\iff \mathbf{v} - \mathbf{y} \in \mathcal{N}_{C(\mathbf{s})}(\mathbf{y}) \\ &\iff \langle \mathbf{v} - \mathbf{y}, \mathbf{z} - \mathbf{y} \rangle \leq 0 && \forall \mathbf{z} \in C(\mathbf{s}) \\ &\iff \langle \mathbf{v} - \mathbf{y}, G^\top(\mathbf{u} - \mathbf{s} - \mathbf{w}^{ext}) - \mathbf{y} \rangle \leq 0 && \forall \mathbf{u} \in K \\ &\iff \langle G\mathbf{v} - G\mathbf{y}, \mathbf{u} - (G\mathbf{y} + \mathbf{s} + \mathbf{w}^{ext}) \rangle \leq 0 && \forall \mathbf{u} \in K \\ &\iff (G\mathbf{y} + \mathbf{s} + \mathbf{w}^{ext}) \in \mathcal{K}_K(G\mathbf{v} + \mathbf{s} + \mathbf{w}^{ext}) \\ &\iff \mathbf{y} = G^T(\Pi_K(G\mathbf{v} + \mathbf{s} + \mathbf{w}^{ext}) - \mathbf{s} - \mathbf{w}^{ext}). \end{aligned}$$



**Figure 4.10:** When the friction coefficient  $\mu$  is high enough (bottom), the cloth falling onto the bunny no longer slides away.

Scene	$m^1$	$n^2$	Dual		Nodal algorithm	
			Build	Solve	Build	Solve
<b>Bunny</b>	6500	169	2.25	0.017	0.02	0.074
<b>Skirt</b>	5208	741	76.7	0.23	0.074	0.14
<b>Dress</b>	6060	1086	157	0.19	0.077	0.048

Times are in seconds.

<sup>1</sup> Number of cloth vertices.

<sup>2</sup> Average number of contacts.

**Table 4.4:** Performance comparison between the standard dual algorithm (with explicit Delassus operator) and the nodal algorithm for solving body–cloth contacts.

### 4.3.2 Results

**Capturing dry friction** We first want to ensure that our algorithm is indeed capable of capturing the effects of Coulomb friction. We run a simple simulation consisting of a square piece of cloth falling onto the Stanford bunny (Figure 4.10). While the friction coefficient is low, the cloth is slipping down under its own weight; however, for higher values of this coefficient, friction maintains the cloth stuck atop the bunny.

**Performance** We compared the cost of solving the DCFP using the standard (dual) Gauss–Seidel algorithm with the dense Delassus operator versus running our nodal algorithm (using the Cadoux fixed-point algorithm and our hybrid Gauss–Seidel for each intermediate SOCQP). Table 4.4 summarize our results for three test cases: the cloth-on-bunny simulation of Figure 4.10, the walk of a mannequin wearing a dress of as in Figure 4.9), and a similar motion with a skirt. When taking into account the cost of computing the Delassus operator, the nodal algorithm ran about two orders of magnitude faster than the naive dual method.

### 4.3.3 Limitations

Our naive nodal algorithm allowed for a significant speed-up w.r.t. solving the standard DCFP on simple simulations, such as that of Figure 4.9 where the dress is contacting only with the underlying body. As the body mesh does not have sharp features, dealing only with contacts

at cloth vertices did not prove to be too limiting. However, the inability to account for self-contacts becomes quite disturbing when folds overlap with each other during motion. Yet, future work could probably accommodate for self-contacts while retaining the efficiency of the nodal algorithm. Indeed, the relative velocity at a self-contacting point between two vertices reads  $\mathbf{u}_i = \mathbf{v}_j - \mathbf{v}_k$ , and the correspond block-row of  $H$ ,  $H_i$ , contains only two blocks,  $H_{i,j} = E_i^\top$  and  $H_{i,k} = -E_i^\top$ . Now, as the Coulomb law is invariant w.r.t. a constant scaling factor on  $\mathbf{u}_i$ , this row could be replaced with  $H_{i,j} = \frac{1}{\sqrt{2}}E_i^\top$  and  $H_{i,k} = -\frac{1}{\sqrt{2}}E_i^\top$ . Then, assuming that the vertices  $j$  and  $k$  are not involved in any other contact, an orthonormal matrix  $G$  can still be constructed by defining its  $j^{\text{th}}$  and  $k^{\text{th}}$  rows as

$$\begin{aligned} G_{j,j} &:= \frac{1}{\sqrt{2}}E_i^\top & G_{j,k} &= -\frac{1}{\sqrt{2}}E_i^\top \\ G_{k,j} &:= \frac{1}{\sqrt{2}}E_i^\top & G_{k,k} &= \frac{1}{\sqrt{2}}E_i^\top. \end{aligned}$$

The  $j^{\text{th}}$  and  $k^{\text{th}}$  constraints should then be  $(\mathbf{u}_j, \mathbf{r}_j) \in \mathcal{C}_{\mu_j}$  and  $\mathbf{r}_k = \mathbf{0}$ , respectively.

However, this still does not address all of the limitations of our approach. The constraint of having only contacts at cloth vertices, and at most one contact per vertex is severe, though it might be alleviated by adaptively spawning vertices with kinematic constraints. Furthermore, cloth are thin objects. In contrast to 3D bodies, for which it is easy to determine what is outside and what is inside, if penetration between two layers of cloth is not resolved at one time step, it will no longer be possible in the following ones to know which layer should be on “top” of the other. Our velocity-based method does not provide this guarantee, in contrast to, for instance, (Brochu et al. 2012).

## 4.4 Inverse modeling with frictional contacts

Until now, we were interested in the direct simulation of mechanical systems; that is, given a discrete set of physical parameters  $\mathbf{p} \in \mathbb{R}^p$  describing the system (for instance, the Young modulus of the material, friction coefficients, etc), and initial positions and velocities, we wanted to compute the trajectory of the system.

Inverse modeling consists in the reverse problem: given the trajectory of the system, can we get back to its physical parameters? In this section, we present briefly two inverse modeling problems which have the particularity of being again able to be expressed as a (or a sequence of) SOCQP, and thus are good candidates for the application of our Gauss–Seidel algorithm. These works have been spearheaded by A. Derouet-Jourdan for fiber inversion and R. Casati for cloth inversion, and have already been treated in much more details in their respective dissertations (Casati 2015; Derouet-Jourdan 2013).

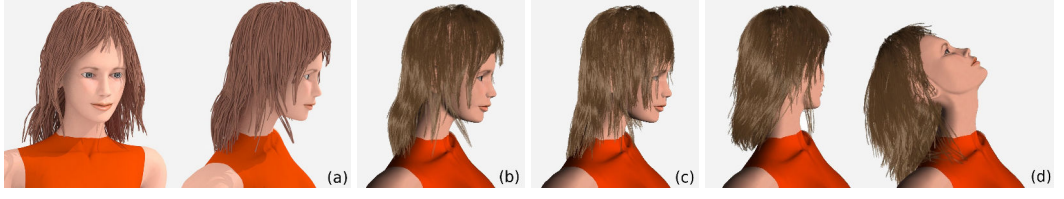
We consider a mechanical system at equilibrium (i.e.,  $\mathbf{v} = \mathbf{0}$ ) under internal, external and unilateral frictional contact forces. Writing that the sum of all forces must be zero, there holds

$$\mathbf{0} = \mathbf{f}^{\text{int}}(\bar{\mathbf{q}}; \mathbf{p}) + \mathbf{f}^{\text{ext}}(\bar{\mathbf{q}}) + H^\top(\bar{\mathbf{q}})\mathbf{r}, \quad (4.18)$$

where  $\bar{\mathbf{q}}$  denotes the generalized coordinates of the system at equilibrium, and the unknowns are  $\mathbf{p}$  and  $\mathbf{r}$  such that  $\forall i = 1 \dots n, \mathbf{r}_i \in \mathcal{K}_{\mu_i}$ .

### 4.4.1 Linear case

We can first assume that the internal forces derive from a potential energy that is quadratic in  $\mathbf{p}$ . This is the case treated in (Derouet-Jourdan, Bertails-Descoubes, Daviet, et al. 2013), where we attempt to retrieve the natural curvatures of a set of Super-Helices (Bertails et al. 2006), assuming that the fibers’ Young modulus, radii, volumetric masses and friction coefficients are known.



**Figure 4.11:** Outline of our hair inverse modeling strategy. (a) Input hair geometry, here a capture from Herrera et al. (2012). (b) Conversion to a mechanical fiber model, here a set of Super-Helices. (c) Without inversion, the hair sags at the beginning of the simulation, and the groom is lost. (d) Using our inversion procedure, the groom stays in place and behaves naturally once the head starts moving.

This addresses an important problem in the visual effects industry: when creating a digital double of a live performing actor, artists (or sufficiently-advanced software) will sculpt the groom so that it matches a reference photography. One can “easily” transform each of the groom’s strands into the rest configuration of an animatable fiber (e.g., a Super-Helix, as in Derouet-Jourdan, Bertails-Descoubes, and Thollot 2013). However, once the simulation is run with gravity, the fibers will fall from their rest configuration and significant sagging will be observed, ruining the similarity between the digital groom and its physical reference (Figure 4.11). One could simply consider fibers one by one, and find a set of rest curvatures and twists so that their deformed configuration under gravity matches the input groom (see e.g., Derouet-Jourdan, Bertails-Descoubes, and Thollot 2010). Then, no sagging will occur when the simulation is run. However, such an inversion procedure will explain all deformations, including those due to contacts (such as hair resting upon a shoulder) as merely stemming from large rest curvatures, which will inevitably lead to weird dynamics. To avoid such artifacts, one must consider the whole groom and take contacts into account when looking for the rest parameters.

**Inversion strategy** In the linear case, the inverse problem (4.18) boils down to finding  $\mathbf{p} \in \mathbb{R}^p$  and  $\mathbf{r} \in \mathbb{R}^{nd}$  such that

$$\mathbf{0} = \mathbf{f}^{\text{ext}} + K\mathbf{p} + H^\top \mathbf{r}, \quad \mathbf{r} \in \Pi_{\mathcal{K}_{\mu_i}}. \quad (4.19)$$

Even when  $K$  is symmetric positive definite, problem (4.19) is under-determined; we can always find a solution for which  $\mathbf{r} = \mathbf{0}$ , but this might not be the only one.

To circumvent this problem, we define a target parameter set  $\mathbf{p}^0 \in \mathbb{R}^p$  from physical considerations, and minimize the distance between our retrieved parameter  $\mathbf{p}$  and  $\mathbf{p}^0$ . Our problem becomes

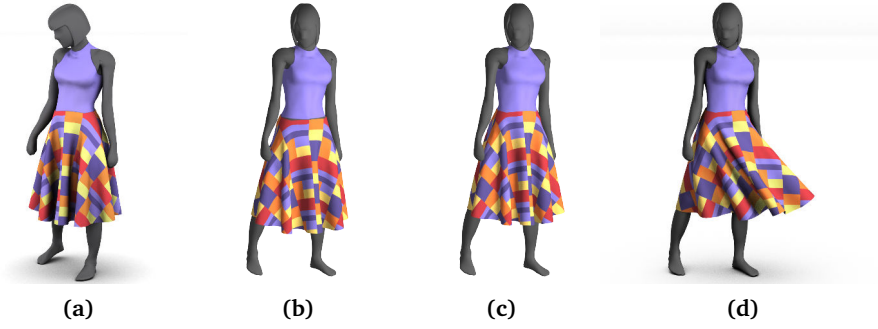
$$\begin{aligned} \min_{\mathbf{p} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{p} - \mathbf{p}^0\|^2 &= \min_{\mathbf{r} \in \Pi_{\mathcal{K}_{\mu_i}}} \frac{1}{2} \|K^{-1} (H^\top \mathbf{r} + \mathbf{f}^{\text{ext}}) + \mathbf{p}^0\|^2 \\ &= \min_{\mathbf{r} \in \Pi_{\mathcal{K}_{\mu_i}}} \frac{1}{2} \mathbf{r}^\top W \mathbf{r} + \mathbf{r}^\top \mathbf{b} + \frac{1}{2} \|K^{-1} \mathbf{f}^{\text{ext}} + \mathbf{p}^0\|^2, \end{aligned}$$

with  $W = HK^{-2}H^\top$  and  $\mathbf{b} = H(K^{-2}\mathbf{f}^{\text{ext}} + K^{-1}\mathbf{p}^0)$ , whose optimality conditions satisfy a problem structurally similar to the convexified DCFP (2.22),

$$\begin{cases} K^2 \mathbf{v} = \mathbf{f}^{\text{ext}} + K\mathbf{p}^0 + H^\top \mathbf{r} \\ \mathbf{u} = H\mathbf{v} \\ \mathcal{K}_{\mu_i}^\perp \ni \mathbf{u}_i \perp \mathbf{r}_i \in \mathcal{K}_{\mu_i} \end{cases} \quad \forall i = 1 \dots n. \quad (4.20)$$

As we saw in Section 2.3.2, problem (4.20) always admits a solution in  $\mathbf{v}$  and can be solved either as a primal or dual SOCQP.  $\mathbf{p}$  can then be deduced from a solution to (4.20) as  $\mathbf{p} = -K^{-1} (H^\top \mathbf{r} + \mathbf{f}^{\text{ext}}) = \mathbf{p}^0 - K\mathbf{v}$ .





**Figure 4.12:** An application of inverse cloth modeling. **(a)** An artist creates a skirt model by sculpting its rest pose under gravity. **(b)** Without inverse modeling, the simulated cloth sags and falls from the hips of the character. **(c)** Inverse modeling ensures that friction at the hips suffices to maintain the cloth at its original position. **(d)** The recovered cloth behaves naturally when external forces (wind) are applied.

**Application to inverse hair modeling** The above framework for taking frictional contacts into account within the inversion procedure requires defining a set of target parameters  $\mathbf{p}^0$ , and three possibilities can be considered:

- Using directly the input curvatures from the deformed configurations;
- Using the curvature computed at the end of each input curve; this assumes that each fiber possesses relatively constant natural curvatures, and that the fiber's end is the least-deformed portion.
- Using a vanishing target curvatures, assuming that the fibers are naturally straight.

Each of these strategies yields a different balance between the influence of the contact forces and that of the natural curliness, and we have not found any of them to work better than the others overall. Instead, we suggest to try each of them for any input, and manually compare the results.

The resulting SOCQP will also often have a singular  $W$  matrix, which means underdetermination of the contact forces. In practice, we propose to regularize the minimization problem by adding a penalization term  $\varepsilon \mathbb{I}_{\mathbb{R}^{nd}}$  to  $W$ , preventing the inversion procedure to consider unlikely huge contact forces. This strategy allowed us to find plausible parameters for input geometries with hundreds of fibers and thousands of contacts, at a cost similar to that of a single simulation step. The hair recovered from the inversion procedure could then either be simulated unmodified, or further edited using physically-based grooming techniques.

#### 4.4.2 Nonlinear case

In (Casati et al. 2016), we proposed a first attempt to solve the inverse problem (4.18) with a nonlinear dependency of the forces w.r.t. the parameters  $\mathbf{p}$ . This is the case for the Discrete Elastic Shell model (Grinspun et al. 2003), when the parameters to recover are the positions of the vertices of the rest shape.

We assume that all non-contact forces derive from a potential energy  $\mathcal{E}^p(\mathbf{q}, \mathbf{p})$ , i.e.,

$$\mathbf{f}^{int}(\mathbf{q}; \mathbf{p}) + \mathbf{f}^{ext}(\mathbf{q}) = -\frac{\partial \mathcal{E}^p}{\partial \mathbf{q}}(\mathbf{q}, \mathbf{p}).$$

We introduce the equilibrium function  $F : \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}^m$ ,

$$F : (\mathbf{q}, \mathbf{p}) \mapsto \frac{\partial \mathcal{E}^p}{\partial \mathbf{q}}(\mathbf{q}, \mathbf{p}) - \Pi_C \left( \frac{\partial \mathcal{E}^p}{\partial \mathbf{q}}(\mathbf{q}, \mathbf{p}) \right)$$



where  $C = H(\bar{\mathbf{q}})^\top \Pi \mathcal{K}_{\mu_i}$ , such that

$$F(\mathbf{q}, \mathbf{p}) = \mathbf{0} \iff \frac{\partial \mathcal{E}^p}{\partial \mathbf{q}}(\mathbf{q}, \mathbf{p}) \in H(\bar{\mathbf{q}})^\top \Pi \mathcal{K}_{\mu_i}.$$

The inverse problem (4.18) is thus equivalent to finding  $\mathbf{p}$  such that  $F(\bar{\mathbf{q}}, \mathbf{p}) = \mathbf{0}$ . Just like for the dynamics (Section 4.3), we assume that the degrees of freedom  $\mathbf{q}$  of the dynamical model are the positions of the shell vertices, and that contacts occur only at those points. This means that the projection on  $C$  simply amounts to projections on rotated SOC, and is thus trivial to compute.

**Draping function** Without diving too deeply into the algorithmic details, we then recast (4.18) as the minimization problem

$$\min_{\mathbf{p} \in \mathbb{R}^p} J(\mathbf{p}) := \frac{1}{2} \|\Phi(\mathbf{p}) - \bar{\mathbf{q}}\|^2 \quad (4.21)$$

where  $\Phi : \mathbb{R}^p \rightarrow \mathbb{R}^m$  is a so-called *draping* function, which associates to a parameter  $\mathbf{p}$  an equilibrium coordinate  $\mathbf{q}$ , i.e., such that  $F(\mathbf{q}, \mathbf{p}) = \mathbf{0}$ . For instance,  $\Phi$  can be defined locally through the implicit function theorem<sup>7</sup>. One solution to evaluate  $\Phi$  would be to perform a full simulation, but that would be too costly. Instead, let us now show how we can evaluate  $\Phi$  as the solution to an optimization problem.

As  $\forall \mathbf{x} \in \mathbb{R}^3$ ,  $-\mathcal{N}_{\mathcal{K}_{\frac{1}{\mu_i}}}(\mathbf{x}) \subset \mathcal{K}_{\mu_i}$ , a sufficient condition to obtain  $F(\mathbf{q}, \mathbf{p}) = \mathbf{0}$ , and thus  $\mathbf{q} = \Phi(\mathbf{p})$ , is to have  $\frac{\partial \mathcal{E}^p}{\partial \mathbf{q}}(\mathbf{q}, \mathbf{p}) \in -H(\bar{\mathbf{q}})^\top \mathcal{N}_{\Pi \mathcal{K}_{\frac{1}{\mu_i}}}(H(\bar{\mathbf{q}})(\mathbf{q} - \bar{\mathbf{q}}))$ ; the condition also becomes necessary when  $\Phi(\mathbf{p}) = \bar{\mathbf{q}}$ , i.e., when  $\mathbf{p}$  is a solution to the inverse problem. Under our nodal contact assumption  $H(\bar{\mathbf{q}})$  is surjective, so we can apply the Corollary A.3 to Property A.12 on the subdifferential of the precomposition with an affine map and obtain an equivalent condition,

$$\frac{\partial \mathcal{E}^p}{\partial \mathbf{q}}(\mathbf{q}, \mathbf{p}) \in -\mathcal{N}_V(\mathbf{q}) \quad (4.22)$$

with  $V = \left\{ \mathbf{q} \in \mathbb{R}^m, H(\mathbf{q} - \bar{\mathbf{q}}) \in \Pi \mathcal{K}_{\frac{1}{\mu_i}} \right\}$ . We recognize this last inclusion as defining a (local) minimum of an optimization problem,

$$\min_{\mathbf{q} \in V} \mathcal{E}^p(\mathbf{q}, \mathbf{p}). \quad (4.23)$$

Evaluating the draping function  $\Phi(\mathbf{p})$  therefore amounts to finding a local minimum of an optimization problem under SOC constraints, which can be cast as a sequence of convexified DCFP:

$$\begin{cases} \frac{\partial^2 E_c}{\partial^2 \mathbf{q}}(\mathbf{q}^k, \mathbf{p})(\mathbf{q}^{k+1} - \mathbf{q}^k) = -\frac{\partial E_c}{\partial \mathbf{q}}(\mathbf{q}^k, \mathbf{p}) + H(\bar{\mathbf{q}})^\top \mathbf{r} \\ \mathbf{u} = H(\bar{\mathbf{q}})(\mathbf{q}^{k+1} - \bar{\mathbf{q}}) \\ \mathcal{K}_{\frac{1}{\mu_i}} \ni \mathbf{u}_i \perp \mathbf{r}_i \in \mathcal{K}_{\mu_i} \end{cases} \quad \forall i = 1 \dots n. \quad (4.24)$$

In practice, we penalize the objective function with a proximal regularization term  $\frac{1}{2\lambda} \|\mathbf{q} - \bar{\mathbf{q}}\|^2$  in order to ensure that the Hessian in (4.24) remains reasonably positive, and avoid finding local minima that are too far away from the target position.  $\lambda$  is then progressively increased as the overall algorithm converges.

<sup>7</sup>This definition does not hold in the presence of contacts, as  $F$  is no longer differentiable, but we can still define  $\Phi$  algorithmically.

**Gradient descent** Now that we know how to compute  $\Phi(\mathbf{p})$ , let us go back to our original minimization problem (4.21). Using a gradient-descent algorithm requires evaluating the gradient of  $J$ ,  $\frac{\partial J}{\partial \mathbf{p}}(\mathbf{p}) = 2 \left( \frac{\partial \Phi}{\partial \mathbf{p}}(\mathbf{p}) \right)^\top (\Phi(\mathbf{p}) - \bar{\mathbf{q}})$ .  $\frac{\partial \Phi}{\partial \mathbf{p}}$  can be evaluated from the application of the implicit function theorem (see Casati et al. 2016 for more details) on the local optimality condition of (4.23), Equation (4.22), formulated as an orthogonal projection as per Corollary A.6,

$$\mathbf{0} = \Pi_V \left( \mathbf{q} - \frac{\partial \mathcal{E}^p}{\partial \mathbf{q}}(\mathbf{q}, \mathbf{p}) \right) - \mathbf{q}.$$

**Results** We attempted to use Algorithm 4.1 to evaluate the draping function, but obtained mitigated results. Indeed, as the stiffness matrix of the intermediate DCFP (4.24) can still have negative eigenvalues, our hybrid Gauss–Seidel algorithm does not always achieve convergence. The non-linear projected gradient descent algorithm presented in (Casati et al. 2016) proved to be more robust (though less efficient), and was used as a fail-safe. Evaluating the draping function thus constituted a major computational bottleneck of the inversion procedure.

For the global optimization process, we considered gradient descent, conjugate gradient descent, and L-BFGS strategies. Achieving convergence with any of those strategies required to equip them with a line search, and thus demanded a very high number of evaluations of the draping function. This limited the applicability of the method to simple examples for which the rest shape was not too far from the target. However, our method was still able to yield useful results, such as preventing a skirt from falling from the hips of a character, as illustrated on Figure 4.12.

## Discussion

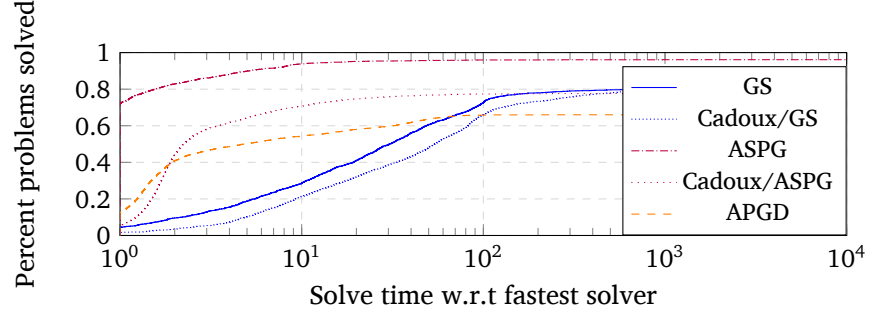
In this chapter, we presented a simple variation of a classical algorithm which performed efficiently and robustly over a variety of applications: hair dynamics and statics, cloth dynamics, and more generally, any DEM contact problem yielding a DCFP or convexified DCFP for which the inverse stiffness matrix  $M^{-1}$  is sparse. Furthermore, we have seen that some inverse modeling problems under frictional contacts also involved SOCQP, justifying again the need for efficient resolution methods.

**Performance** It may seem strange that an algorithm as naive as the Gauss–Seidel, subject to a lot of theoretical limitations, performs so well in practice. One explanation would be that a lot of domain knowledge is built into the local solver, especially when an analytical one is used. This is not as much the case for projective approaches, in which the constraint is only defined through an orthogonal projection operator. Note that Mazhar et al. (2015) argue that the Gauss–Seidel algorithm performs much more poorly than their APGD algorithm. However, they use a crippled variant of the Gauss–Seidel algorithm in which the local solver consist only in performing one projection of the forces onto the admissible cone, as in (Heyn 2013, p. 44). As each contact is also solved with a good precision, even if the asymptotic convergence of whole algorithm is slow, the local force values are never too far from an acceptable solution — this makes the Gauss–Seidel very robust in practice.

That being said, for now we have only tested our solver on hair examples consisting of a few thousands fibers only — that is, about 50 times smaller than a real human hair. Although we feel our solver could possibly resist some further scaling up (using a similar strategy, Kaufman, Tamstorf, et al. (2014) simulated one order of magnitude more fibers), we know that computational time will eventually become the main bottleneck, and concurrent algorithms appear to remain necessary for going beyond that.

The Gauss–Seidel algorithm may also not be the most adequate when considering smaller problems that must be solved to higher precisions, such as those of the FCLIB<sup>8</sup> benchmark (Acary,

<sup>8</sup><http://fclib.gforge.inria.fr/>



**Figure 4.13:** Performance profiles for a variety of solvers on the FCLIB benchmark — percentage of problems that each solver was able to solve under a multiple  $x$  of the time taken by the fastest solver for this problem. The line  $x = 1$  gives the solvers that were the most often the fastest, while  $y = 1$  shows the most robust ones.

Brémond, et al. 2014). The FCLIB-0.2 collection is a set of thousands of DCFP in reduced formulation (2.19) and without linear constraints, generated by four different kinds of simulations (stacks of rigid boxes, unstructured capsules and rigid-body chains). All problems feature less than on thousand degrees of freedom and contact points. Figure 4.13 shows performance profiles for different solvers on an arbitrary subset of 2500 problems. On this benchmark, the projected gradient variants, and especially **ASPG** (Algorithm B.2), performed much better than our Gauss–Seidel algorithm. This in part due to the fact that the combination of sleeping heuristics and naive parallelization may hamper convergence on problems with low number of contacts, as the algorithm becomes more alike to Jacobi.

**Continuum mechanics** Alternatively, one may renounce attempting to solve for each frictional contact force, and instead simulate the macroscopic interactions between the fibers using a continuum model. This is the strategy that we are going to adopt in the second part of this dissertation, though for much simpler systems consisting only of rigid spherical grains. While DCFP solvers (Chapter 3) have been initially devised to work on problems arising from contacts between discrete bodies (Chapter 2), we will see that, exploiting the similarity between the Coulomb and non-associated Drucker–Prager laws (Chapter 1), they can actually be used for a wider range of applications, including some problems from continuum mechanics. In some cases, we will once again be able to leverage the hybrid Gauss–Seidel algorithm presented in this chapter.

## **Part II**

# **Continuum simulation of granular materials**



## 5 Continuum simulation of granular flows

We already mentioned in our introduction that several attempts have been made to model the averaged behavior of granular materials. In particular, we mentioned the Mohr–Coulomb and Drucker–Prager yield surfaces, which restrict the set of admissible stresses inside granulars to take into account that:

- Granular materials offer a much higher resistance to compression than to stretching. In fact, materials like dry sand can freely dilate, but their compression Young modulus is very high (hundreds of MPa).
- Due to the Coulombic nature of the friction between grains, the resistance of granular materials to shearing increases with the local pressure.

In the following, we review a few methods that have been used in the literature to perform numerical simulations of granular materials modeled as continua.

**Notation** We will generally denote *values* with italic Greek (usually for tensors and scalars) or roman (for vectors and scalars) letters, and *fields* with the corresponding upright symbols. For instance, we shall write  $\mathbf{u} = \mathbf{u}(\mathbf{x}, t)$  the value of the velocity field at  $\mathbf{x} \in \Omega$  and at instant  $t$ , and  $\boldsymbol{\sigma} = \boldsymbol{\sigma}(\mathbf{x}, t)$  the value of the stress field.

### 5.1 Continuum models

As the numerical simulation of granular continua has found a wide range of applications, from the small-strain study of soils to the prediction of the run-out of avalanches, it is natural that different point of views have been taken in the literature.

Most of the constitutive models devised for this purpose fit in the framework of Implicit Standard Materials (ISM) presented in Chapter 1, and boil down to (omitting potential hardening variables),

$$\left\{ \begin{array}{l} \dot{\boldsymbol{\varepsilon}} = \dot{\boldsymbol{\varepsilon}}_e + \dot{\boldsymbol{\varepsilon}}_p \\ \boldsymbol{\sigma} \in \rho \frac{\partial \mathcal{E}}{\partial \boldsymbol{\varepsilon}_e}(\boldsymbol{\varepsilon}_e) \\ \boldsymbol{\sigma} \in \frac{\partial b}{\partial \dot{\boldsymbol{\varepsilon}}_p}(\dot{\boldsymbol{\varepsilon}}_p, \boldsymbol{\sigma}), \end{array} \right. \quad (5.1) \quad (5.2) \quad (5.3)$$

where  $\mathcal{E}$  is the free-energy function and  $b$  is a bipotential (Definition 1.2). In the case of Generalized Standard Materials (GSM),  $b$  is simply defined through a dissipation potential  $\mathcal{D}$  as  $b(\dot{\boldsymbol{\varepsilon}}_p, \boldsymbol{\sigma}) := \mathcal{D}(\dot{\boldsymbol{\varepsilon}}_p) + \mathcal{D}^*(\boldsymbol{\sigma})$ . However, different bodies of research have considered different expressions for  $\mathcal{E}$  and  $b$  in the context of granular continua.

A first class of methods (*flow-oriented*) is motivated by the liquid regime of granulars, thus the modeling of very large deformations, and derive from the standard framework of yield stress fluids. The grains are generally assumed to be perfectly rigid, so that the material is inelastic, leading to visco-plastic models.

Another class of methods (*elasticity-oriented*) has been historically focused on the study of soils and the solid behavior of granulars, in the fashion of the elasto-plastic “Cam Clay” model (Vermeer 1998).

### 5.1.1 Inelastic yield-stress fluids

The flow-oriented approach consists in adapting the standard framework of inelastic yield-stress fluids to account for the pressure-dependency observed in granular materials.

**Incompressible yield stress fluids** Fluids such as toothpaste behave like a solid at rest, but can flow once an external load is applied. Like Newtonian fluids, they remain inelastic and incompressible. In the framework of GSM, this means that the free-energy function is chosen as  $\mathcal{E} = \mathcal{J}_{\{0_{S_d}\}}$ , so that Equation (5.2) implies  $\varepsilon_e = \mathbf{0}$ . Moreover, the dissipation potential must be such that  $\mathcal{D}(\gamma\mathbb{I}) = \mathcal{J}_{\{0\}}(\gamma)$ , so that Equation (5.3) enforces  $\text{Tr } \dot{\varepsilon} = 0$ , i.e.,  $\nabla \cdot \mathbf{u} = 0$ .

However, in contrast to Newtonian fluids, the dissipation potential is nonsmooth w.r.t. the tangential strain  $\text{Dev } \dot{\varepsilon}$ . The simplest constitutive law for incompressible yield-stress fluid is the Bingham rheology, of which we already presented an inviscid version in Section 0.3.1. This rheology derives from a simple yet nonsmooth visco-plastic dissipation potential  $\mathcal{D}_{\text{Bi}}$ ,

$$\mathcal{D}_{\text{Bi}}(\text{Dev } \dot{\varepsilon}) := \eta(\text{Dev } \dot{\varepsilon})^2 + \sigma_s |\text{Dev } \dot{\varepsilon}|,$$

and can be written using a disjunctive formulation as

$$\begin{cases} \text{Dev } \boldsymbol{\sigma} = \eta \text{Dev } \dot{\varepsilon} + \sigma_s \frac{\text{Dev } \dot{\varepsilon}}{|\text{Dev } \dot{\varepsilon}|} & \text{if } \text{Dev } \dot{\varepsilon} \neq \mathbf{0} \\ |\text{Dev } \boldsymbol{\sigma}| \leq \sigma_s & \text{if } \text{Dev } \dot{\varepsilon} = \mathbf{0}. \end{cases} \quad (5.4)$$

This can be seen as the combination of a Newtonian viscosity with the inelastic associated flow rule deriving from the yield surface  $F_{\text{Bi}}(\boldsymbol{\sigma}) := |\text{Dev } \boldsymbol{\sigma}| - \sigma_s$ .

The Bingham rheology (5.4) is a particular case of the Herschel-Buckley rheology, which is widely used for concrete or crude oil flows. Such industrial applications have made incompressible Bingham flows the subject of a very large body of research, and several numerical methods have been devised to handle the nonsmoothness of the rheology in an implicit manner; we refer to (Saramito and Wachs 2016) for a comprehensive review. Such methods can be classified into two main categories:

- Regularizing methods, which employ diverse numerical artifacts to smooth out the singularities of the rheology (see, e.g., (Frigaard and Nouar 2005) for a review). The main benefit of a regularizing strategy is the use of classical numerical schemes dedicated to differential equations. While some regularizing methods have been shown to converge to the true solution even in the ill-conditioned inviscid case (Bouchut, Eymard, et al. 2014), this limit cannot be realized numerically, and the regularization may predict yielding where none should occur (Frigaard and Nouar 2005).
- Nonsmooth methods, the first of which being Augmented Lagrangian methods (such as Arrow–Hurwitz or ADMM) based on the framework of (Fortin and Glowinski 1983), as used by Saramito and Roquet (2001). This class of methods is able to correctly predict the yielded and unyielded zones, but historically suffered from slow numerical convergence in practice. Recently however, several authors have focused on improving the convergence rate of non-regularizing approaches. Bleyer et al. (2015) proposed a method for Herschel–Buckley fluids based on a SOCP reformulation, which they claim benefits from much faster convergence properties than Augmented Lagrangian methods. Saramito (2016) used a Newton algorithm on a complementarity function to solve for the steady-state of yield stress flows without Newtonian viscosity, and also observed improved convergence over the Augmented Lagrangian algorithm. Concurrently and in a similar fashion to our work, Treskatis et al. (2016) proposed the use of a Nesterov-accelerated proximal method on the dual minimization problem.

**Granular materials** Several works have focused on extending these numerical methods to the simulation of the non-associated Drucker–Prager or  $\mu(I)$ <sup>1</sup> rheologies in 2D, inheriting in

<sup>1</sup> See Section 0.3.3.

the process the incompressibility constraint of Bingham flows (Chambon et al. 2011; Chauchat and Médale 2014; Ionescu et al. 2015; Lagrée et al. 2011; Staron et al. 2014). In Computer Graphics, Y. Zhu and Bridson (2005) used a similar approach for the simulation of 3D granular flows, enriching an incompressible flow solver with special treatment for zones below the yield stress.

However, incompressibility of the flow means that the pressure field will only be defined up to a constant, and the usual zero-average condition will lead to the appearance of negative pressure values. This choice may prove unfortunate in certain configurations, such as in the wake of an obstacle, where special care has to be taken to ensure that the rheology remains well-defined (Chauchat and Médale 2014). Even then, using an incompressible flow model does not allow to capture the asymmetry of the pressure field experimentally observed by Seguin et al. (2016). Moreover, Barker et al. (2015) showed that the incompressible  $\mu(I)$  rheology is unstable for a wide range of parameters; they suggest that this instability may be alleviated by modeling the transitions between the different regimes of the granular material, and thus allowing dilation. In a very recent work, Heyman et al. (2016) also argue that relaxing the incompressibility constraint can resolve these stability issues.

In contrast, Narain, Golas, et al. (2010) proposed a flow-oriented approach that does not preclude the expansion of the material, and allows transitions between the solid, liquid and gaseous regimes. However, their approach was hampered by their use of a staggered time-integration procedure and a quite-rough approximation of the yield surface.

### 5.1.2 Elasto-plastic models

Another important body of research considers the plastic deformation of an elastic body. These approaches have been originally devised to study the behavior of metals outside of their elastic range, but can accommodate a much wider range of constitutive relationships, including that of granular materials (Vermeer 1998).

In contrast to flow-oriented approaches, these methods allow elastic displacement, and split the strain tensor into an elastic part  $\epsilon_e$  and a plastic part  $\epsilon_p$ , using either additive or multiplicative plasticity theory. Using the terminology of the ISM framework, the bipotential  $b$  is usually chosen to be inviscid, and is defined from a *yield surface*, — i.e., a function  $F$  such that the set of admissible stresses is defined by  $F(\sigma) \leq 0$  — and a *flow rule*. Associated flow rules remain within the framework of GSM and are obtained by choosing the dissipation potential  $\mathcal{D}$  such that  $\mathcal{D}^* = \mathcal{D}_{\{\tau, F(\tau) \leq 0\}}$ . In this case, Remark 1.2.1 states that Equation (5.3) can be written equivalently as (5.5),

$$\dot{\epsilon}_p \in \alpha \frac{\partial F}{\partial \sigma}(\sigma) \text{ and } 0 \geq F(\sigma) \perp \alpha \geq 0. \quad (5.5)$$

Non-associated flow rules are also commonly used to model granular continua, and can be devised by replacing (5.3) with (5.6),

$$\dot{\epsilon}_p \in \alpha \frac{\partial G}{\partial \sigma}(\sigma) \text{ and } 0 \geq F(\sigma) \perp \alpha \geq 0, \quad (5.6)$$

where the *plastic potential*  $G$  differs from  $F$ . The strains and stress at every point of the material are either computed explicitly — for instance by projecting the stresses onto the inside of the yield surface — or implicitly, using a so-called *return-mapping* algorithm (Simo and Hughes 2000), which usually takes the form of a root-finding Newton algorithm. Note that the non-associated Drucker–Prager flow rule, which we presented in the framework of ISM in Chapter 1, can also be written as (5.6). In this case,  $G$  is simply deduced from  $F$  by replacing the friction coefficient  $\mu$  with the dilatancy coefficient  $\zeta$ .

Mast (2013) compares different yield surfaces and flow rules (associated Drucker–Prager, Bingham or non-associated Natzo–Makai), and proposes both implicit and explicit methods for their simulation. As for Drucker–Prager the plastic potential  $G$  is not differentiable at the “tip” of the cone of admissible stresses,  $\frac{\partial G}{\partial \sigma}$  may contain several elements, which complicates the return-mapping algorithm. Mast (2013, Section 3.6) proposes to use a two-surface approximation of



the yield-surface to alleviate this difficulty. Klar et al. (2016) argue that in the case of dry sand, the projection on the Drucker–Prager yield surface should be orthogonal to the hydrostatic axis  $\sigma_1 = \sigma_2 = \sigma_3$ . This corresponds indeed to an unassociated Drucker–Prager flow-rule with null dilatancy, as the projection is then analogous to that defining the Alart–Curnier complementarity function, as illustrated in Figure 1.3. Dunatunga and Kamrin (2015) explicitly models the transition between the gaseous and dense regimes based on the local volume fraction of grains, and restrict the use of the elasto-plastic flow rule to the latter regime.

This class of methods inherently model more complete physics, as it accounts for the measurable elasticity of the macroscopic material while flow-based methods do not. However, it suffers from the fact the compression Young modulus of granular materials made of rigid grains is very high, which leads to extremely stiff numerical systems; moreover, the time-scale of the elastic response is generally much lower than that of the plastic deformation, and not necessarily relevant for our applications. Using explicit integration schemes then requires taking very small timesteps, and implicit integration yields ill-conditioned, hard to solve systems. In practice, several works artificially lower the Young modulus for numerical efficiency purposes, weakening the physical correctness argument.

## 5.2 Spatial discretization strategies

While the above dichotomy focused on the constitutive law of the material (elastoplastic vs viscoplastic), we can also sort the different granular simulation methods following the strategy that they use to discretize the so-called *conservation* equations, which we first recall below.

### 5.2.1 Continuous conservation equations

The evolution of the material state in time is driven by two conservation equations. Let  $\rho$  denote the density of the material and  $\mathbf{u}$  its velocity, then the conservation of mass reads

$$\frac{\partial \rho}{\partial t} + \nabla \cdot [\rho \mathbf{u}] = 0, \quad (5.7)$$

and the conservation of momentum (in conservative form) is

$$\frac{\partial \rho \mathbf{u}}{\partial t} + \nabla \cdot [\rho \mathbf{u} \otimes \mathbf{u} - \boldsymbol{\sigma}^{tot}] = \mathbf{f}^{ext}, \quad (5.8)$$

where  $\boldsymbol{\sigma}^{tot}$  is the total stress tensor and  $\mathbf{f}^{ext}$  the total density of external forces. These equations are derived from the application of basic physical principles over elementary volumes; we refer the interested reader to (Saramito 2013) for more details.

The momentum conservation equation (or *momentum balance*) (5.8) is typically written under another, *non-conservative* form, using the identity

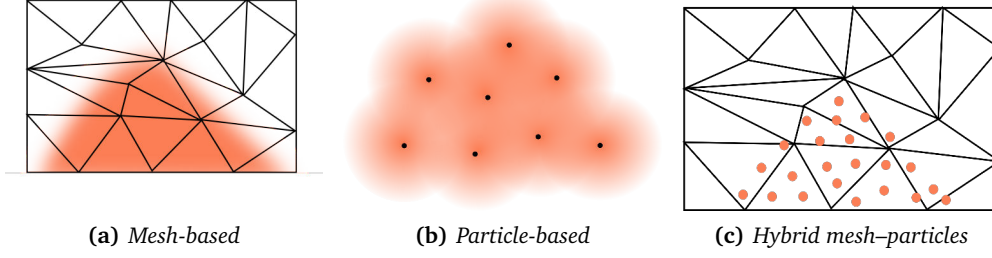
$$\begin{aligned} \nabla \cdot [\rho \mathbf{u} \otimes \mathbf{u}] &= \nabla \cdot [\rho \mathbf{u}] \mathbf{u} + (\rho \mathbf{u} \cdot \nabla) \mathbf{u} && \text{(see Remark 5.1)} \\ &= -\frac{\partial \rho}{\partial t} \mathbf{u} + \rho (\mathbf{u} \cdot \nabla) \mathbf{u} && \text{using (5.7).} \end{aligned} \quad (5.9)$$

Then using (5.9) inside the momentum balance (5.8), we get the non-conservative form,

$$\rho \left( \frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} \right) - \nabla \cdot \boldsymbol{\sigma}^{tot} = \mathbf{f}^{ext}. \quad (5.10)$$

**Remark 5.1.** We use the classical notation  $(\mathbf{u} \cdot \nabla) \mathbf{v}$ , where  $\mathbf{u}$  is a vector field (a “velocity”) and  $\mathbf{v}$  a scalar or vector field, to denote the quantity  $(\nabla \mathbf{v}) \mathbf{u}$ , that is, the variation of  $\mathbf{v}$  when traveling along the streamlines of the velocity field  $\mathbf{u}$ .

To further simplify notation, we can also introduce the notion of total derivative.



**Figure 5.1:** Different strategies for the discretization of the mass and momentum conservation equations.

**Definition 5.1.** Let  $\mathbf{u}$  be a vector field and  $\mathbf{v}$  a scalar or vector field, the total derivative of  $\mathbf{v}$  w.r.t. time and the velocity field  $\mathbf{u}$  will be denoted by

$$\frac{D_{\mathbf{u}}\mathbf{v}}{Dt} := \frac{\partial \mathbf{v}}{\partial t} + (\mathbf{u} \cdot \nabla)\mathbf{v}.$$

When no confusion is possible, we will omit the velocity field and simply write  $\frac{D\mathbf{v}}{Dt}$ .

The total derivative expresses the total variation of the quantity  $\mathbf{v}$  at a material point  $\mathbf{x}(t)$ , taking into account the fact that the material point moves with the velocity field  $\mathbf{u}$ , i.e.,  $\frac{d\mathbf{x}}{dt} = \mathbf{u}$ . Using this new notation, the mass and momentum conservation equations can be written

$$\frac{D\rho}{Dt} + \rho \nabla \cdot \mathbf{u} = 0 \quad (5.11)$$

$$\rho \frac{D\mathbf{u}}{Dt} - \nabla \cdot \boldsymbol{\sigma}^{tot} = \mathbf{f}^{ext} \quad (5.12)$$

Now, there are several ways to discretize these equations in space and express the gradient and divergence operators. We first recall briefly the existence of two points of view for looking at the evolution of a dynamical system:

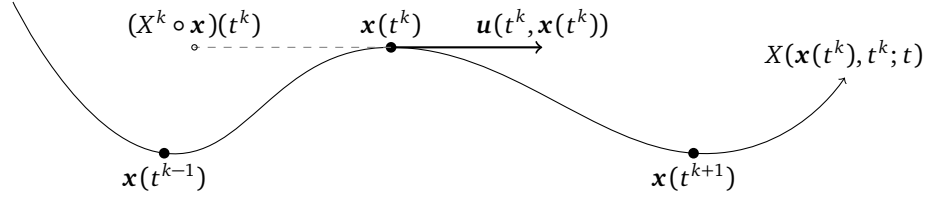
- With the *Lagrangian* point-of-view, which is classically used in DEM, the system is described w.r.t. the moving material points. As every material point knows its own history, the transport terms are easy to handle; however, the topology of the material may change over time, making the spatial differential operators difficult to evaluate.
- In contrast, the *Eulerian* point of view looks at the system from a fixed point of space. The material can be described w.r.t. a fixed domain, allowing for easy spatial differentiation, but making transport terms harder to treat.

In practice, many methods use combinations of these two points of view.

### 5.2.2 Mesh-based discretization

A first solution is, in an Eulerian fashion illustrated in Figure 5.1(a), to extrapolate the density, velocity and stress field over the whole simulation domain from their values at selected points on a mesh, the so-called *degrees of freedom*. The structure of the mesh can then be used to approximate the differential operators of the conservation equations using different strategies:

- finite-differences, where the differential operators are evaluated by computing the difference in the values at neighboring nodes of a regular grid, eventually staggered;
- finite-volumes, where the divergence operator is computed from the flux between neighboring mesh cells using Green's theorem (e.g., Lagr  e et al. 2011);



**Figure 5.2:** Characteristics morphism  $X$  with its linearized version  $X^k$

- finite-elements (FEM) — of which the two previous items can actually be seen as special cases, where the conservation equations are solved under weak form, by looking at their image through their scalar product with a finite set of test functions (e.g., Ionescu et al. 2015).

Note that shallow models, which reduce a phenomenon whose depth is much smaller than its horizontal extent to a lower-dimensional problem, have also been developed with success for granular flows. Savage and Hutter (1989) were the first to propose one such model, assuming a proportional relationship of the components of the stress tensor — thus a flowing hypothesis. Ionescu (2010, 2013) studied the onset of the flow, with and without topography. Bouchut, Ionescu, et al. (2016) take a slightly different approach and consider a two-layer model, then study the depth of the interface between the solid and liquid regimes.

**Transport terms** The presence of total derivatives (transport terms) significantly complicates the numerical treatment of mesh-based conservation equations. The momentum balance equation is thus often studied in the *creeping flow* limit, when the inertial term  $\rho(\mathbf{u} \cdot \nabla)\mathbf{u}$  is much smaller than the internal or external forces, but this strategy is not viable for all scenarios.

A convenient technique for the treatment of transport terms in timestepping schemes is the use of the *characteristics* method (see Etienne 2004), which brings insight from the Lagrangian point of view. The velocity and/or density fields are then advected following a backtracking strategy based on a low-order approximation of the trajectory of material points during the previous timestep. As in Figure 5.2, let us denote by  $X(\mathbf{x}_0, t_0; t)$  the position at every instant of the material point that is at  $\mathbf{x}_0$  at  $t_0$ ;  $X$  is the solution of the Cauchy problem

$$\begin{cases} X(\mathbf{x}_0, t_0, t) = \mathbf{x}_0 \\ \frac{\partial X}{\partial t}(\mathbf{x}_0, t_0, t) = \mathbf{u}(X(\mathbf{x}_0, t_0, t), t). \end{cases} \quad (5.13)$$

Etienne (2004, Section 2.1) argues that for a regular-enough velocity field  $\mathbf{u}$  and a small-enough time interval, the solution is unique, and thus the characteristics morphism  $X$  is well-defined. Now, composition by  $X$  allows to go back from the total derivative to the usual partial derivative. Indeed, let  $\mathbf{v}$  be a scalar or vector field, then using the notation  $(\mathbf{v} \circ X)(\mathbf{x}, t_0; t) := \mathbf{v}(X(\mathbf{x}, t_0; t), t)$ ,

$$\frac{\partial \mathbf{v} \circ X}{\partial t}(\mathbf{x}, t_0; t) = \left[ \frac{\partial \mathbf{v}}{\partial t} + \frac{\partial \mathbf{v}}{\partial \mathbf{x}} \frac{\partial X}{\partial t} \right](X(\mathbf{x}, t_0; t), t) = \frac{D\mathbf{v}}{Dt} \circ X. \quad (5.14)$$

This identity yield a simple first-order discretization of the total derivative. Consider a finite time-step  $[t^l, t^{k+1} = t^k + \Delta_t]$ , then

$$\frac{D\mathbf{v}}{Dt}(\mathbf{x}, t^{k+1}) = \frac{\mathbf{v}(\mathbf{x}, t^{k+1}) - (\mathbf{v} \circ X)(\mathbf{x}, t^{k+1}; t^k)}{\Delta_t} + O(\Delta_t).$$

Now, computing  $X(\mathbf{x}, t^{k+1}; t^k)$  amounts to solving a Cauchy problem, which is unwieldy; in practice a low-order approximation of this quantity is used, usually using a Runge-Kutta method.

At the first order,

$$\frac{D\mathbf{v}}{Dt}(\mathbf{x}, t^{k+1}) = \frac{\mathbf{v}(\mathbf{x}, t^{k+1}) - \mathbf{v}(X^k(\mathbf{x}), t^k)}{\Delta_t} + O(\Delta_t), \quad X^k(\mathbf{x}) := \mathbf{x} - \Delta_t \mathbf{u}^k(\mathbf{x}, t^k). \quad (5.15)$$

Note that this requires being able to locate arbitrary points in the mesh, which can be costly and require using a dedicated spatial index structure. In a similar manner, one may evolve the mesh after each timestep, moving its nodes along the velocity field in a Lagrangian manner and thus yielding an explicit access to the characteristics morphism  $X$ . However, for flows incurring large displacements, as is often the case with granular materials, the mesh can quickly become degenerate and require remeshing.

Other approaches discretize directly the transport operators – possibly without resorting to timestepping scheme — in the Eulerian framework, but have to take special care to ensure energy and mass conservation yet prevent overshooting (the advected quantities should be physically admissible, for instance the density should remain positive). Saramito (2013, Section 4.10) illustrates that centered methods, such as finite-differences or continuous finite-elements, lead to parasitic oscillations on transport equations. Conversely, upwind discontinuous Galerkin methods (a finite-element method which can be interpreted as an arbitrary-degree generalization of finite-volumes, and which we will discuss briefly in Section 7.2.2) have been shown to possess good stability properties when discretizing transport terms (Pietro and Ern 2011). Specialized methods have also been devised for convection-dominated equations such as the mass conservation (5.7), for instance the WENO schemes (Shu 2009).

Note that when considering an incompressible flow with a spatially constant initial density, the mass conservation equation becomes trivial, as the density simply remains constant through time. This makes mesh-based approaches especially suited for simulating granular materials as incompressible yield-stress flows. However, if the portion of space occupied by the material is moving, one still has either to deform the mesh (e.g., Ionescu et al. 2015), or advect a boundary (e.g., Lagr  e et al. 2011).

**Visualization** For Computer Graphics applications, the ultimate goal is to produce a set of grain samples that can be rendered. While it is easy to sample a density field at each individual frame, doing so in a temporally consistent manner is a hard problem; this drawback has also contributed to spark interest in particle-based approaches.

### 5.2.3 Particle-based discretization

Smoothed Particle Hydrodynamics (SPH) methods take a different approach, and remove the need for any structured mesh by making further use of the Lagrangian point of view. The material is discretized as a finite set of particles that represent clumps of material rather than individual grains. The velocity, stress and density fields are then discretized as the sum of smooth, compact, radially-symmetric kernels centered on those particles, as in Figure 5.1(b).

As such, the results are intrinsically free from grid artifacts. Moreover, this approach makes the transport terms trivial to discretize. Indeed, solving the mass conservation equations just amounts to moving the particles. Using an implicit Euler integration rule and denoting by  $\mathbf{x}_p^k(t)$  and  $\mathbf{v}_p^k(t)$  the position and velocity of the  $p^{\text{th}}$  particle at time  $t^k$ , this means

$$\mathbf{x}_p^{k+1} = \mathbf{x}_p^k + \Delta_t \mathbf{v}_p^k.$$

The particles also yield a direct access to the characteristics morphism (5.13), and the velocity's total derivative in the momentum balance equation can thus be simply discretized as

$$\frac{D\mathbf{u}}{Dt}(\mathbf{x}_p(t^{k+1}), t^{k+1}) = \frac{\mathbf{v}_p(t^{k+1}) - \mathbf{v}_p(t^k)}{\Delta_t} + O(\Delta_t).$$

Several authors have applied SPH to granular simulation (Alduán and Otaduy 2011; Alduán et al. 2009; Chambon et al. 2011; Ihmsen et al. 2013; Lenaerts and Dutré 2009). However, such methods typically require more degrees of freedom (particles) than mesh-based methods. Evaluating the differential operators require keeping track of each pair of neighboring particles, which can be costly. Enforcing constraints on the density is usually done with predictive-corrective approaches, requiring a lot of iterations with slow convergence.

### 5.2.4 Hybrid methods

Hybrid methods, illustrated on Figure 5.1(c), attempt to strike a compromise between the two aforementioned approaches. They still use particles for the representation of the material state — thus allowing easy handling of the transport terms, but augment them with a background mesh to help with the computation of differential operators and the internal forces. At each step, the new velocities are computed on the mesh, and then transferred to the particles.

The main difficulty of this approach thus lies in the particles-to-mesh and mesh-to-particles transfers; they should conserve energy, yet be stable. The first method of this kind, the Particle-in-Cell (PIC; Harlow 1963) method, used finite-differences on a regular grid for the mesh-based portion, and simple interpolation for transferring the velocities both ways.

We assume that the  $j^{\text{th}}$  component of a field  $\mathbf{g}_j$  at any point in space can be evaluated by extrapolating its values  $\mathbf{g}_{i,j}$  at the degrees of freedom ( $\mathbf{y}_i$ ) using the shape functions  $\omega_{i,j}^v$ , i.e.,  $\mathbf{g}_j(\mathbf{x}, t^k) = \sum_i \mathbf{g}_{i,j}^k \omega_{i,j}^v(\mathbf{x})$ . A any timestep  $t^k$ , PIC transfers the velocities of the particles to the grid as a weighted average,

$$\mathbf{u}_{i,j}^{p \rightarrow g, k} = \frac{\sum_p m_p \omega_{i,j}^v(\mathbf{x}_p^k) \mathbf{v}_{p,j}^k}{\sum_p m_p \omega_{i,j}^v(\mathbf{x}_p^k)}, \quad (5.16)$$

where  $m_p$  denotes the mass of the  $p^{\text{th}}$  particle.

The total velocity derivative can then be approximated to the first order as

$$\frac{D\mathbf{u}}{Dt}(\mathbf{x}, t^{k+1}) = \frac{\mathbf{u}(\mathbf{x}, t^{k+1}) - \mathbf{u}^{p \rightarrow g}(\mathbf{x}, t^k)}{\Delta_t} + O(\Delta_t).$$

Finally, particles gets their velocities from the newly computed velocity field,

$$\mathbf{v}_p^{k+1} = \mathbf{u}(\mathbf{x}_p^k, t^{k+1}) \quad (5.17)$$

are moved using a semi-implicit Euler step,

$$\mathbf{x}_p^{k+1} = \mathbf{x}_p^k + \Delta_t \mathbf{v}_p^{k+1}. \quad (5.18)$$

However, this strategy was found to quickly dissipate kinetic energy. This led to the introduction of the FLuid-Implicit-Particle (FLIP) method (Brackbill and Ruppel 1986); instead of transferring the velocities from the grid to the particles, the particle velocities are updated using the difference in grid velocity from the previous to the current timestep. That is, Equation (5.17), is replaced with (5.19),

$$\mathbf{v}_p^{k+1} = \mathbf{v}_p^k + \left( \mathbf{u}(\mathbf{x}_p^k, t^{k+1}) - \mathbf{u}^{p \rightarrow g}(\mathbf{x}_p^k, t^k) \right) \quad (5.19)$$

While this second velocity update rule leads to a much better conservation of kinetic energy, FLIP suffers from being prone to instabilities, losing the high-frequency filtering properties of PIC (Jiang, Schroeder, Selle, et al. 2015). In practice, a weighted average of the PIC and FLIP update rules is generally used. More recently, Jiang, Schroeder, Selle, et al. (2015) introduced the Affine Particle-in-Cell method (APIC). In this case, the grid-to-particles transfer use the original PIC rule (5.17), but an additional term, modeling the velocity gradients, is added to the particle-to-grid velocity transfers (5.16). This approach boasts a much better energy conservation than PIC, while preserving its stability and filtering properties.

The Material-Point Method (MPM; Sulsky 1994) and its derivatives, on which we will expand in Chapter 7, is another refinement of the PIC method that allows the usage of other mesh-based discretization strategies, such as the finite-element method. Any of the previously mentioned grid transfer strategies, PIC, FLIP and APIC, can also be used. The PIC particles-to-mesh transfer (5.16) can then be interpreted as using the *lumped* mass-matrix, corresponding to a trapezoidal approximation of the finite-element integral. As we will in Chapter 7, using an exact integration rule (*consistent* mass-matrix) with the PIC velocity update (5.17) preserves the kinetic energy, but yields a singular system when there are not enough particles in any given mesh cell. MPM has been often leveraged for the simulation of granular materials, mainly using elasto-plastic models (e.g., (S. Bardenhagen et al. 2000; Dunatunga and Kamrin 2015; Klar et al. 2016; Wieckowski, Youn, et al. 1999)) and explicit time-integration.

### 5.3 Our approach

We will dedicate the two following chapters to the devising of a simulation method for dry granular materials that meets the following design goals.

#### 5.3.1 Design goals

**Inelasticity** Following the flow-oriented approaches, we will not attempt to capture elasticity of the material, but rather claim that for most applications the solid regime may be considered as purely rigid — which is a valid assumption for sand-like materials at low confining pressures (Roux and Radjai 1998). This will avoid us having to model the elasticity time scale, and thus allow the use of rougher time discretization schemes for increased computational performance. Moreover, we will not regularize the resulting multi-valued constitutive law. Instead, we claim that tackling it in an *implicit* manner by leveraging tools from nonsmooth optimization will yield better conditioned systems and stable results, even for large timesteps. Finally, similarly as for the discrete contact mechanics setting, we shall assume impacts to be inelastic, that is, shocks will not propagate inside the granular medium.

**Dilatability and regime switching** In contrast with most flow-oriented methods and following (Narain, Golas, et al. 2010), our method will not preclude the expansion of the flow, allowing the material to transition freely from a dense regime to a gaseous regime. This strategy avoids negative pressure zones and the ill-definition of the yield stress in these regions, and alleviates the instability of the incompressible  $\mu(I)$  rheology.

A few existing approaches explicitly model abrupt transitions between the gaseous and dense (i.e., solid or liquid) regimes. They do so by considering a critical value for the density below which the grains are assumed to never be in persistent contact, and thus the contact stress vanishes. Switching between these two regimes may be done either in an explicit (Dunatunga and Kamrin 2015) or implicit (Narain, Golas, et al. 2010) manner; for numerical stability reasons we will follow the latter strategy. Moreover, in order to be consistent with our inelastic impact hypothesis, no rebound shall occur when the material suddenly enters the dense regime.

**Computers graphics** The different regimes exhibited by granular materials yield very rich dynamics, and their ubiquity in outdoor environments has made their visually plausible simulation a primary goal of the Computer Graphics community.

As we already mentioned, Computer Graphics is subject to different constraints and targets than engineering communities. While the physical accuracy criterion is less drastic in movies than risk-assessment applications, stability, robustness and numerical efficiency are prime requirements for visual effects pipelines. Obviously, any kind of visual artifacts is proscribed.

One of our design goals is to devise a method that is not only physically sound, but that should also be able to serve as a basis for Computer Graphics applications.

### 5.3.2 Outline of this second part

In Chapter 6, we will first consider dense granular materials, that is, materials that are already tightly packed together and cannot be compacted anymore. We will relax this assumption in Chapter 7, and consider the unconstrained flow of granular materials, but assume that the interactions with the surrounding fluid are negligible (which is the case for large-enough grains in the air). Finally, in Chapter 8 we will relax this second assumption, and model fully coupled flows such as powders in the air or immersed avalanches.

## 6 Dense granular flows

In this chapter, we consider a first simplified set of equations that is paradoxical, yet will allow us to retrieve classical results in a variety of scenarios, and will serve as a building block for more complex or accurate models. Most of the contents of this chapter have been published in (Daviet and Bertails-Descoubes 2016b).

### 6.1 Constitutive equations

#### 6.1.1 Unilateral incompressibility

Let  $\phi$  denote the volume fraction field, that is, the fraction of space occupied by the granular material at every point of the domain.  $\phi$  can take values in  $[0, \phi_{\max}]$ , where  $\phi_{\max} \leq 1$  is a maximum packing fraction depending on the geometry of the grains (for rigid spheres all of identical size,  $\phi_{\max}$  cannot be above 0.74). Assuming a constant density  $\rho$  for the grains, the density field  $\rho$  can be computed as  $\rho(\phi) = \rho\phi$ , and thus the mass conservation equation (5.7) reads

$$\frac{D\phi}{Dt} = -\phi \nabla \cdot \mathbf{u}. \quad (6.1)$$

**Dense flow hypothesis** We consider that the material is, everywhere on the simulation domain and at every instant in time, densely packed; that is,  $\phi = \phi_{\max}$ .

However, we want to take into account the asymmetric yielding behavior of granulars by allowing the onset of expansion, while strictly preventing compaction. Taking the derivative w.r.t. time of the maximum compaction constraint  $\phi \leq \phi_{\max}$  for an already densely packed material yields  $\frac{D\phi}{Dt} \leq 0$ , that is, using (6.1),

$$\nabla \cdot \mathbf{u} \geq 0. \quad (6.2)$$

We will call (6.2) the *unilateral incompressibility constraint*.

Actually, for the material to remain densely packed everywhere, mass conservation (6.1) dictates  $\nabla \cdot \mathbf{u} = 0$ . The paradoxical nature of our model comes from the fact that we only require positive divergence of the flow — not null divergence — yet still assume  $\phi = \phi_{\max}$  at every instant. This means that mass conservation will not hold inside the strictly dilating zones. Fortunately, this inconsistency will have little effect on the predicted flow inside the dense regions. As validated in our Section 6.5, our simplified model still remains applicable to a variety of relevant scenarios.

Now, the main advantage of this dense flow hypothesis is that it avoids having to couple the momentum and mass conservation equations, side-stepping several technical difficulties and allowing for a much more concise exposition of our implicit numerical method, which is the main contribution of this chapter. Moreover, by carefully choosing the temporal and spatial discretization strategies for the constraint  $\phi \leq \phi_{\max}$ , we will show in Chapter 7 that the framework presented here can be easily adapted to accommodate arbitrary flows.

**Complementarity condition** For standard incompressible flows, the pressure field  $p$  can be seen as a Lagrange multiplier enforcing the null-divergence condition. Here, we will instead



set the pressure to enforce the unilateral compressibility constraint (6.2), i.e., formulating the Karush-Kuhn-Tucker conditions for this inequality,

$$\begin{cases} p \geq 0 & \text{and } \nabla \cdot \mathbf{u} = 0 \\ p = 0 & \text{and } \nabla \cdot \mathbf{u} > 0, \end{cases} \quad \text{or}$$

or, using an equivalent complementarity notation,

$$0 \leq p \perp \nabla \cdot \mathbf{u} \geq 0.$$

This means that the dilating zones will feature a vanishing pressure, while those that tend to remain densely packed may exhibit positive pressure. With our dense flow hypothesis, this follows the suggestion by Drew (1983) that the pressure should vanish when  $\phi < \phi_{\max}$ . Cohesion may also be modeled by allowing negative values for the pressure field; more generally, we will use the complementarity condition

$$0 \leq p + c \perp \nabla \cdot \mathbf{u} \geq 0, \quad (6.3)$$

which states that the flow has to overcome a negative pressure of  $-c$  before starting to dilate.

Our results (Section 6.5) will show that relaxing the common incompressibility assumption  $\nabla \cdot \mathbf{u} = 0$  prevents the arising of an ill-defined rheology in some typical scenarios such as the flow in the wake of an obstacle (Chauchat and Médale 2014), and allows us to correctly retrieve a vanishing pressure field in this region.

### 6.1.2 Friction

Similarly as in Ionescu et al. (2015), we consider a rheology combining a viscosity  $\eta$  with a yield-stress  $\kappa(\sqrt{\frac{d}{2}}p)$  that depends linearly on the local pressure  $p$ ,

$$\kappa(x) := \sigma_s + \mu \left( x + \sqrt{\frac{d}{2}}c \right).$$

However, here we consider that  $\eta$  does not depend on the local velocity or pressure. We have chosen to incorporate the factor  $\sqrt{\frac{d}{2}}$  in the above definitions in order for our physical parameters to match those of the non-associated Drucker–Prager rheology defined in Section 1.3, as we will see below. Similarly, the distinct roles of the two constant coefficients in  $\kappa$ , the cohesion  $c$  and the additional stress  $\sigma_s$ , will be made clear by using this analogy.

The total stress tensor  $\boldsymbol{\sigma}^{\text{tot}}$  is defined as<sup>1</sup>

$$\begin{aligned} \boldsymbol{\sigma}^{\text{tot}} &:= 2\eta E^b(\boldsymbol{\varepsilon}) + \boldsymbol{\tau}^F - p\mathbb{I} \\ E^b(\boldsymbol{\varepsilon}) &:= \left( \text{Dev} + \frac{b}{d}\mathbb{I}\text{Tr} \right) \dot{\boldsymbol{\varepsilon}}, \end{aligned} \quad 0 < b \leq 1,$$

where  $\boldsymbol{\tau}^F$  is a traceless symmetric tensor bounded by the yield stress  $\kappa(\sqrt{\frac{d}{2}}p)$  w.r.t. the norm  $|\cdot|$  defined in Section (1.3.1),

$$|\boldsymbol{\tau}^F| = \sqrt{\langle \boldsymbol{\tau}^F, \boldsymbol{\tau}^F \rangle} = \sqrt{\frac{1}{2} \boldsymbol{\tau}^F : \boldsymbol{\tau}^F}.$$

Writing the maximum dissipation principle (i.e., the associated flow rule) for  $\boldsymbol{\tau}^F$  yields

$$\begin{cases} \boldsymbol{\tau}^F = \kappa \left( \sqrt{\frac{d}{2}}p \right) \frac{\text{Dev} \dot{\boldsymbol{\varepsilon}}}{|\text{Dev} \dot{\boldsymbol{\varepsilon}}|} & \text{if } \text{Dev} \dot{\boldsymbol{\varepsilon}} \neq \mathbf{0} \\ |\boldsymbol{\tau}^F| \leq \kappa \left( \sqrt{\frac{d}{2}}p \right) & \text{if } \text{Dev} \dot{\boldsymbol{\varepsilon}} = \mathbf{0}. \end{cases} \quad (6.4)$$

<sup>1</sup> Using  $E^b(\dot{\boldsymbol{\varepsilon}})$  rather than the usual  $\text{Dev} \dot{\boldsymbol{\varepsilon}}$  for the Newtonian viscosity models a non-zero bulk viscosity, and will ensure well-posedness of our equations (Proposition 6.2). However, in practice  $b$  can be chosen equal to 0.

**Reformulation** We will use once again the notions of *normal* and *tangential* parts of a symmetric tensor  $\tau \in S_d$  introduced in Section 1.3.1,

$$\tau_N := \langle \tau, \iota_d \rangle = \frac{1}{\sqrt{2d}} \text{Tr } \tau \quad \tau_T := \text{Dev } \tau.$$

Using these notations, the total contact stress tensor,  $\sigma^c := \tau^F - p\mathbb{I}$ , satisfies  $\sigma_T^c = \tau^F$  and  $\sigma_N^c = -\sqrt{\frac{d}{2}}p$ . The frictional stress conditions (6.4) can be written as

$$\begin{cases} \sigma_T^c = \kappa(-\sigma_N^c) \frac{\dot{\epsilon}_T}{|\dot{\epsilon}_T|} & \text{if } \dot{\epsilon}_T \neq 0 \\ |\sigma_T^c| \leq \kappa(-\sigma_N^c) & \text{if } \dot{\epsilon}_T = 0. \end{cases} \quad (6.5)$$

Moreover, since  $\nabla \cdot \mathbf{u} = \text{Tr } \dot{\epsilon} = \frac{1}{\sqrt{2d}} \dot{\epsilon}_N$  and using the fact that complementarity is invariant to strictly positive scalings, Equation (6.3) can be written equivalently as

$$0 \leq \sqrt{\frac{d}{2}}c - \sigma_N^c \perp \dot{\epsilon}_N \geq 0. \quad (6.6)$$

Putting (6.5) and (6.6) together, we can identify the disjunctive formulation (1.20) of the non-associated Drucker–Prager flow rule without dilatancy and with a tensile yield stress  $\tau_c = \sqrt{\frac{d}{2}}c$ . We will thus replace our rheology constraints (6.3–6.4) with the more general inclusion (6.7) allowing the modeling of a non-zero dilatancy  $\zeta$ ,

$$(\dot{\epsilon}, p\mathbb{I} - \tau^F) \in \mathcal{DP}(\mu, \sigma_S, \tau_c, \zeta).$$

Finally, since we mentioned earlier that  $\tau_c$  amounts to a simple translation of the set of admissible stresses along the hydrostatic axis, we will introduce the change of variable  $\lambda := \tau_c \iota_d - \sigma^c := (p + c)\mathbb{I} - \tau^F$  and write our rheology as

$$(\dot{\epsilon}, \lambda) \in \mathcal{DP}(\mu, \sigma_S, \zeta). \quad (6.7)$$

**Cases covered by our choice of rheology** Much like in (Ionescu et al. 2015), our set of parameters allows us to explore an interesting range of constitutive laws. When  $\mu = 0$ , we retrieve the viscoplastic Bingham rheology (with unilateral incompressibility – fully incompressible Bingham could be obtained in the  $\tau_c = +\infty$  limit), while when taking  $\eta = 0$ ,  $\tau_c = 0$  and  $\sigma_S = 0$ , we get a purely Coulombic plastic flow — or, in other terms, the  $\mu(I)$  rheology with identical static and dynamic friction coefficients. The numerical method presented in Section 6.2 will assume a non-zero  $\eta$  to obtain a well-posed system. However, our so-called “primal” algorithm will be able to handle a vanishing Newtonian viscosity, although some theoretical results will be lost. This constraint will be also be altogether alleviated in Section 6.5.3 when using temporal schemes. In our numerical experiments, we were also able to simulate the complete  $\mu(I)$  rheology by explicitly computing the inertial number  $I$  at each time step. As the range of values for the  $\mu(I)$  coefficient is relatively small, treating this term in an explicit fashion did not significantly degrade the stability of our simulations.

$\tau_c$  models cohesion (actually, a tensile yield strength), but the role of  $\sigma_S$  is different. For instance, a non-zero  $\sigma_S$  means that there can be a non-zero frictional stress for dilating flows. In contrast to  $\tau_c$ ,  $\sigma_S$  will induce a purely deviatoric stress. Moreover, as we saw in chapter 1, while  $\tau_c$  amounts to a simple translation of the set of admissible stresses (which in our case will mean adding a confining pressure),  $\sigma_S$  fundamentally changes the structure of this admissible set, which becomes a truncated cone.  $\sigma_S$  is probably more relevant for concrete-like materials than granular ones; however, this parameter can also be useful in the 2D case to model wall friction or effects related to the geometry of the grains, which may induce a yield stress independent of the pressure in the simulation plane.

## 6.2 Creeping flow

In this first section we assume that the flow is slow enough for its inertia to be neglected, and solve for its steady state. While inertial effects are generally significant for granular flows, this case remains relevant since the structure of the equations for each timestep of the fully dynamic case will be similar (see Section 6.5.3).

### 6.2.1 Steady-state and boundary conditions

Let us consider a domain  $\Omega \subset \mathbb{R}^d$  and decompose its boundary as  $\text{Bd } \Omega := B_D \cup B_N$ , with Dirichlet boundary conditions on  $B_D$  and homogeneous Neumann on  $B_N$ ,

$$\mathbf{u} = \mathbf{u}_D \quad \text{on } B_D \quad (6.8a)$$

$$E^b \mathcal{D}(\mathbf{u}) \mathbf{n}_\Omega = \mathbf{0} \quad \text{on } B_N \quad (6.8b)$$

where  $\mathbf{n}_\Omega$  is the outward-pointing normal to  $\Omega$  on each point of its boundary.

Moreover, let  $\boldsymbol{\sigma}^{\text{ext}}$  gather all external stresses applied onto the Neumann boundary of the domain. We thus have

$$-\lambda \mathbf{n}_\Omega = (\boldsymbol{\sigma}^{\text{ext}} - c\mathbb{I}) \mathbf{n}_\Omega \quad \text{on } B_N. \quad (6.9)$$

We assume that the sole external force is the action of gravity whose orientation is given by the “down” unit vector  $\mathbf{e}_g$ . Writing the conservation of momentum (5.12) for the steady-state gives

$$-\nabla \cdot \left[ \underbrace{2\eta E^b \dot{\boldsymbol{\varepsilon}} - \lambda}_{=\boldsymbol{\sigma}^{\text{tot}} - c\mathbb{I}} \right] = \rho g \mathbf{e}_g + \underbrace{\nabla c}_0 \quad \text{on } \Omega. \quad (6.10)$$

From Equations (6.9) and (6.10), we see that the cohesion term  $c$  simply modifies the external stress term  $\boldsymbol{\sigma}^{\text{ext}}$ ; in the following, we will gather both of those terms into  $\boldsymbol{\sigma}^{\text{ext}}$  and ignore  $c$ .

**Dimensionless equations** Let  $L$  be a characteristic dimension of the flow. We define  $U := \sqrt{gL}$  as the characteristic velocity and  $P := \rho gL$  the characteristic pressure of the flow. We consider dimensionless differential operators defined through  $\tilde{\nabla} := L \nabla$ . We furthermore introduce two dimensionless numbers, the Reynolds number  $\text{Re} := \frac{\rho UL}{\eta}$  and the Bingham number  $\text{Bi} := \frac{\sigma_s}{\rho gL} = \frac{\sigma_s}{P}$ .

Considering the dimensionless quantities  $\tilde{\mathbf{u}} := \frac{1}{U} \mathbf{u}$ ,  $\tilde{\boldsymbol{\varepsilon}} := \tilde{\mathcal{D}}(\tilde{\mathbf{u}}) = \frac{L}{U} \dot{\boldsymbol{\varepsilon}}$ ,  $\tilde{\lambda} := \frac{1}{P} \lambda$ , and  $\tilde{\mathbf{u}}_D := \frac{1}{U} \mathbf{u}_D$ , Equations (6.8 – 6.10) can be made dimensionless as

$$\left\{ \begin{array}{ll} -\tilde{\nabla} \cdot \left[ \frac{2}{\text{Re}} E^b \tilde{\boldsymbol{\varepsilon}} - \tilde{\lambda} \right] = \mathbf{e}_g & \text{on } \Omega \\ \tilde{\mathbf{u}} = \tilde{\mathbf{u}}_D & \text{on } B_D \\ E^b \tilde{\mathcal{D}}(\tilde{\mathbf{u}}) \mathbf{n}_\Omega = \mathbf{0} & \text{on } B_N \\ -\tilde{\lambda} \mathbf{n}_\Omega = \tilde{\boldsymbol{\sigma}}^{\text{ext}} \mathbf{n}_\Omega & \text{on } B_N, \end{array} \right. \quad (6.11)$$

with the dimensionless rheology

$$(\tilde{\boldsymbol{\varepsilon}}, \tilde{\lambda}) \in \mathcal{DP}(\mu, \text{Bi}, \zeta). \quad (6.12)$$

From now on we shall use the dimensionless quantities and omit the tildes.

### 6.2.2 Variational formulation

Let  $H^1(\Omega)^d$  be the usual Sobolev space containing square-integrable functions from  $\Omega \subset \mathbb{R}^d$  to  $\mathbb{R}^d$ , with square-integrable gradients. As in (Saramito 2015, Appendix A), we introduce the affine subspace  $V(\mathbf{u}_D)$  of  $H^1(\Omega)^d$  for which the Dirichlet boundary condition (6.8a) is satisfied, i.e.,

$$V(\mathbf{u}_D) := \{\mathbf{u} \in H^1(\Omega)^d; \mathbf{u} = \mathbf{u}_D \text{ on } B_D\}.$$

$V(\mathbf{0})$  is therefore the vector subspace of  $H^1(\Omega)^d$  for which the homogeneous Dirichlet boundary condition  $\mathbf{u}|_{B_D} = \mathbf{0}$  is satisfied. Let also  $T(\Omega)$  be the space of square-integrable symmetric tensor fields on  $\Omega$ .

**Proposition 6.1.** *A weak form of System (6.11 – 6.12) amounts to finding  $\mathbf{u} \in V(\mathbf{u}_D)$ ,  $\boldsymbol{\lambda} \in T(\Omega)$ , and  $\boldsymbol{\gamma} \in T(\Omega)$ , such that*

$$\begin{cases} a(\mathbf{u}, \mathbf{v}) = b(\boldsymbol{\lambda}, \mathbf{v}) + l(\mathbf{v}) & \forall \mathbf{v} \in V(\mathbf{0}) & (6.13a) \\ m(\boldsymbol{\gamma}, \boldsymbol{\tau}) = b(\boldsymbol{\tau}, \mathbf{u}) & \forall \boldsymbol{\tau} \in T(\Omega) & (6.13b) \\ (\boldsymbol{\gamma}, \boldsymbol{\lambda}) \in \mathcal{DP}(\mu, \text{Bi}, \zeta), & & (6.13c) \end{cases}$$

where  $\forall \mathbf{x}, \mathbf{y} \in H^1(\Omega)^d$  and  $\forall \boldsymbol{\sigma}, \boldsymbol{\tau} \in T(\Omega)$ ,  $a(\mathbf{x}, \mathbf{y})$  and  $m(\boldsymbol{\sigma}, \boldsymbol{\tau})$  are the symmetric positive-definite bilinear forms on  $H^1(\Omega)^d \times H^1(\Omega)^d$  and  $T(\Omega) \times T(\Omega)$ , respectively,

$$\begin{aligned} a(\mathbf{x}, \mathbf{y}) &:= \frac{2}{\text{Re}} \int_{\Omega} \text{Dev} \mathbf{D}(\mathbf{x}) : \text{Dev} \mathbf{D}(\mathbf{y}) + \frac{b}{d} (\nabla \cdot \mathbf{x}) (\nabla \cdot \mathbf{y}) \\ m(\boldsymbol{\sigma}, \boldsymbol{\tau}) &= \int_{\Omega} \boldsymbol{\sigma} : \boldsymbol{\tau}, \end{aligned}$$

$b(\boldsymbol{\tau}, \mathbf{x})$  is the bilinear form on  $T(\Omega) \times H^1(\Omega)^d$ ,

$$b(\boldsymbol{\tau}, \mathbf{x}) = \int_{\Omega} \mathbf{D}(\mathbf{x}) : \boldsymbol{\tau},$$

and  $l(\mathbf{x})$  is the linear form on  $H^1(\Omega)^d$ ,

$$l(\mathbf{x}) = \int_{\Omega} \langle \mathbf{e}_g, \mathbf{x} \rangle - \int_{B_N} \langle (\boldsymbol{\sigma}^{\text{ext}} \mathbf{n}_{\Omega}), \mathbf{x} \rangle.$$

*Proof.* First, let us consider the stress boundary condition (6.9). One solution would be to enforce it strongly, by constraining  $\boldsymbol{\lambda}$  to the subspace of  $T(\Omega)$  which satisfy (6.9). However, this may lead to difficulties in the discretization of the  $\mathcal{DP}(\mu, \text{Bi}, \zeta)$  rheology. In our proposed implementation, we choose instead to model, in the integration of the term  $\nabla \cdot \boldsymbol{\sigma}^{\text{tot}}$  over  $\Omega$ , a possibly non-zero jump  $\llbracket \boldsymbol{\sigma}^{\text{tot}} \rrbracket = -\boldsymbol{\lambda}|_{B_N} - \boldsymbol{\sigma}^{\text{ext}}$  of the stress on  $B_N$ .

Let us now derive the variational formulation for System (6.11). We assume  $\mathbf{u} \in V(\mathbf{u}_D)$  and  $\boldsymbol{\lambda} \in T(\Omega)$ , and let  $\mathbf{v} \in V(\mathbf{0})$  be a test function. Multiplying both sides of the first line of (6.11) by  $\mathbf{v}$  and integrating over  $\Omega$  yields

$$-\int_{\Omega} \left\langle \nabla \cdot \left[ \frac{2}{\text{Re}} E^b \mathbf{D}(\mathbf{u}) - \boldsymbol{\lambda} \right], \mathbf{v} \right\rangle + \int_{B_N} \langle \llbracket \boldsymbol{\sigma}^{\text{tot}} \rrbracket \mathbf{n}_{\Omega}, \mathbf{v} \rangle = \int_{\Omega} \langle \mathbf{e}_g, \mathbf{v} \rangle. \quad (6.14)$$

Using the Green formula with the Neumann (6.8b) and homogeneous Dirichlet (6.8a) boundary conditions for  $\mathbf{u}$ ,

$$\begin{aligned}
a(\mathbf{u}, \mathbf{v}) &= \frac{2}{\text{Re}} \int_{\Omega} \text{Dev} D(\mathbf{u}) : \text{Dev} D(\mathbf{v}) + \frac{b}{d} (\nabla \cdot \mathbf{v})(\nabla \cdot \mathbf{u}) \\
&= \int_{\Omega} \frac{2}{\text{Re}} E^b D(\mathbf{u}) : D(\mathbf{v}) \\
&= - \int_{\Omega} \left\langle \nabla \cdot \left[ \frac{2}{\text{Re}} E^b D(\mathbf{u}) \right], \mathbf{v} \right\rangle
\end{aligned} \tag{6.15}$$

and

$$\int_{\Omega} \langle \nabla \cdot \boldsymbol{\lambda}, \mathbf{v} \rangle + \int_{B_N} \langle \llbracket \boldsymbol{\sigma}^{\text{tot}} \rrbracket \mathbf{n}_{\Omega}, \mathbf{v} \rangle = - \int_{B_N} \langle \boldsymbol{\sigma}^{\text{ext}} \mathbf{n}_{\Omega}, \mathbf{v} \rangle - b(\boldsymbol{\lambda}, \mathbf{v}). \tag{6.16}$$

By combining (6.14 – 6.16) and the definition of  $l$ , we retrieve (6.13a). Let us now focus on the rheology  $\mathcal{DP}(\mu, \text{Bi}, \zeta)$  given in (6.12), which contains inequalities that cannot be put directly under weak form. To circumvent this difficulty, we introduce an auxiliary variable  $\boldsymbol{\gamma} \in T(\Omega)$  that weakly satisfies  $\boldsymbol{\gamma} = D(\mathbf{u})$ , i.e.,

$$\int_{\Omega} D(\mathbf{u}) : \boldsymbol{\tau} = \int_{\Omega} \boldsymbol{\gamma} : \boldsymbol{\tau} \quad \forall \boldsymbol{\tau} \in T(\Omega), \tag{6.17}$$

which is exactly equation (6.13b). We can thus express the rheology  $\mathcal{DP}(\mu, \text{Bi}, \zeta)$  under the weak form as (6.13b – 6.13c).  $\square$

**Remark** Note that we do not include any additional equation ensuring the well-posedness of our system, such as the zero-average pressure condition which is commonly used for Stokes flows. Indeed, for a given velocity field, the  $\mathcal{DP}(\mu, \text{Bi}, \zeta)$  rheology imposes the value of  $\boldsymbol{\lambda}$  in the yielded regions, serving as an intrinsic boundary condition for the stress field inside the rigid zones.

### 6.2.3 Cadoux algorithm

In this section, we adapt the algorithm of Cadoux (2009) to show that solutions to our variational problem (6.13a–6.13c) can be characterized as the fixed-point of a sequence of strictly convex and well-posed minimization problems.

**Proposition 6.2.** *Let  $s \in L_2(\Omega)$ , let us use  $\mathcal{K}_{\frac{1}{\mu}}$  as a short-hand for  $\{\boldsymbol{\tau} \in T(\Omega), \boldsymbol{\tau}(\mathbf{x}) \in \mathcal{K}_{\frac{1}{\mu}(\mathbf{x})} \text{ a.e. on } \Omega\}$ . Under the regularity condition ( $\mathcal{H}(s)$ ),*

$$\exists \mathbf{u} \in V(\mathbf{u}_D), D(\mathbf{u}) + s \boldsymbol{\iota}_d \in \text{int } \mathcal{K}_{\frac{1}{\mu}}, \tag{\mathcal{H}(s)}$$

and assuming  $\text{Re} < +\infty$ , the minimization problem

$$q(s) := \min_{\mathbf{u} \in C(s)} J(\mathbf{u}) \tag{6.18}$$

$$J(\mathbf{u}) := \frac{1}{2} a(\mathbf{u}, \mathbf{u}) - l(\mathbf{u}) + \text{Bi } g(D(\mathbf{u}))$$

$$g(\boldsymbol{\gamma}) := 2 \int_{\Omega} |\boldsymbol{\gamma}_T| \tag{6.19}$$

$$C(s) := \{\mathbf{u} \in V(\mathbf{u}_D), D(\mathbf{u}) + s \boldsymbol{\iota}_d \in \mathcal{K}_{\frac{1}{\mu}}\}$$

admits a unique solution which satisfies

$$\begin{cases} a(\mathbf{u}, \mathbf{v}) = b(\boldsymbol{\lambda}, \mathbf{v}) + l(\mathbf{v}) & \forall \mathbf{v} \in V(\mathbf{0}) \\ m(\boldsymbol{\gamma}, \boldsymbol{\tau}) = b(\boldsymbol{\tau}, \mathbf{u}) & \forall \boldsymbol{\tau} \in T(\Omega) \\ \boldsymbol{\gamma} + s \boldsymbol{\iota}_d \in -\mathcal{N}_{\mathcal{K}_{\frac{1}{\mu}}(\boldsymbol{\lambda})}. \end{cases} \tag{6.20}$$

**Remark 6.1.** We readily deduce from Property 1.6 and Equation (6.20) that a solution  $\mathbf{u}(s)$  to the minimization problem (6.18) such that  $s = (\mu - \zeta)|\gamma_T|$  will yield a solution to our original variational problem (6.13a–6.13c). We can therefore define an operator  $F : L_2(\Omega) \rightarrow L_2(\Omega)$ ,  $s \mapsto (\mu - \zeta)|D(u(s))_T|$ , with  $u : s \mapsto \arg \min_{\mathbf{u} \in C(s)} J(\mathbf{u})$ , so that any fixed-point of  $F$  will yield a solution to (6.13a–6.13c). Note that when  $\mu = \zeta$ , the fixed-point will be reached in a single step.

*Proof.* As from Korn's inequality  $J$  is coercive and strictly convex, and under  $(\mathcal{H}(s))$  the feasible set is not empty, Theorem (A.5) ensures that there exists a unique solution  $\mathbf{u}$  to (6.18), which using Theorem A.6 can be characterized by

$$\mathbf{0} \in \partial J(\mathbf{u}) + \mathcal{N}_{C(s)}(\mathbf{u}). \quad (6.21)$$

$D(\cdot) : H^1(\Omega)^d \rightarrow T(\Omega)$  and  $\cdot|_{B_D} : H^1(\Omega)^d \rightarrow L^2(B_D)^d$  are linear operators from and to Hilbert spaces (trace theorem), we can thus use Corollary A.4 to Property A.12 on the normal cone to a precomposition by an affine map. Under the regularity condition  $(\mathcal{H}(s))$ ,

$$\mathcal{N}_{C(s)}(\mathbf{u}) = D(\cdot)^T \mathcal{N}_{\mathcal{H}_{\frac{1}{\mu}}} (D(\mathbf{u}) + s\mathbf{t}_d) + (\cdot|_{B_D})^T \mathcal{N}_{\{\mathbf{0}\}}(\mathbf{u}|_{B_D} - \mathbf{u}_D) \quad (6.22)$$

where  $D(\cdot)^T$  denotes the adjoint operator to  $D(\cdot)$  for the usual scalar products of  $H^1(\Omega)^d$  and  $T(\Omega)$ . Using once again Property A.12, we have also

$$\partial J(\mathbf{u}) = a(\mathbf{u}, \cdot) - l + \text{Bi } D(\cdot)^T \partial g(D(\mathbf{u})), \quad (6.23)$$

where  $\partial g$  is also defined w.r.t. the usual scalar product of  $T(\Omega)$ ,  $\int_{\Omega} (\cdot : \cdot)$ . Replacing (6.23) and (6.22) into (6.21), taking the scalar product with a test function  $\mathbf{v} \in V(\mathbf{0})$  and noticing that  $V(\mathbf{0}) = \text{Ker}(\cdot|_{B_D}) = (\text{Im}(\cdot|_{B_D})^T)^\perp$ , we get

$$(6.21) \iff \begin{cases} a(\mathbf{u}, \mathbf{v}) = l(\mathbf{v}) + \langle D(\cdot)^T \lambda, \mathbf{v} \rangle_{H^1} & \forall \mathbf{v} \in V(\mathbf{0}) \\ \gamma = D(\mathbf{u}) \\ \lambda \in -\text{Bi } \partial g(\gamma) - \mathcal{N}_{\mathcal{H}_{\frac{1}{\mu}}}(\gamma + s\mathbf{t}_d). \end{cases} \quad (6.24)$$

Let us reformulate the last inclusion. We recognize from Section 1.3 and Equation (1.25) that for  $\tau \in S_d$ ,

$$\text{Bi}(|\tau_T|) + \mathcal{J}_{\mathcal{H}_{\frac{1}{\mu}}}(\tau) = \mathcal{J}_{-\mathcal{T}_{\mu, \text{Bi}}}^*(\tau). \quad (6.25)$$

Using once again the notation  $\mathcal{T}_{\mu, \text{Bi}}$  as a shorthand for  $\{\tau \in T(\Omega), \tau(\mathbf{x}) \in \mathcal{T}_{\mu(\mathbf{x}), \text{Bi}} \text{ a.e.}\}$ , we get from Property A.5 on subdifferentials in function spaces that

$$\begin{aligned} -\mathcal{N}_{\mathcal{T}_{\mu, \text{Bi}}}(-\sigma) &= \left\{ \tau \in T(\Omega), \tau(\mathbf{x}) \in \mathcal{N}_{-\mathcal{T}_{\mu(\mathbf{x}), \text{Bi}}}(\sigma(\mathbf{x})) \text{ a.e. on } \Omega \right\} \\ &= \left\{ \tau \in T(\Omega), \sigma(\mathbf{x}) \in \partial \mathcal{J}_{-\mathcal{T}_{\mu, \text{Bi}}}^*(\tau(\mathbf{x})) \text{ a.e.} \right\} && \text{(Theorem A.2)} \\ &= \left\{ \tau \in T(\Omega), \sigma(\mathbf{x}) \in \partial \left( \text{Bi}(|\cdot_T|) + \mathcal{J}_{\mathcal{H}_{\frac{1}{\mu}}} \right)(-\tau(\mathbf{x})) \text{ a.e.} \right\} && \text{using (6.25)} \\ &= \left\{ \tau \in T(\Omega), \sigma(\mathbf{x}) \in \left( \text{Bi} \partial(|\cdot_T|) + \mathcal{N}_{\mathcal{H}_{\frac{1}{\mu}}} \right)(\tau(\mathbf{x})) \text{ a.e.} \right\} \\ &= \left\{ \tau \in T(\Omega), \sigma \in \text{Bi} \partial g(\tau) + \mathcal{N}_{\mathcal{H}_{\frac{1}{\mu}}}(\tau) \right\}. \end{aligned}$$

Note that the “2” in the definition of  $g$  is indeed retrieved in the last equality, as the scalar product and subdifferential on  $T(\Omega)$  are defined w.r.t. the “ $\cdot : \cdot$ ” scalar product on  $S_d$ , while the subdifferential of  $|\cdot|$  on  $S_d$  is defined using the  $\langle \cdot, \cdot \rangle = \frac{1}{2} : \cdot$  scalar product.

Hence, noticing that  $\gamma_T = (\gamma + s\mathbf{t}_d)_T$ , we get the equivalence

$$-\lambda \in \text{Bi} \partial \left( \int |\cdot_T| \right)(\gamma) + \mathcal{N}_{\mathcal{H}_{\frac{1}{\mu}}}(\gamma + s\mathbf{t}_d) \iff \gamma + s\mathbf{t}_d \in -\mathcal{N}_{\mathcal{T}_{\mu, \text{Bi}}}(\lambda). \quad (6.26)$$

Finally, using that  $\langle D(\cdot)^\top \lambda, \mathbf{v} \rangle_{H^1} = \langle D(\mathbf{v}), \lambda \rangle_{L_2} = b(\lambda, \mathbf{v})$  and (6.26) to rewrite the first and last line of (6.24), respectively, we get that

$$(6.21) \iff \begin{cases} a(\mathbf{u}, \mathbf{v}) = b(\lambda, \mathbf{v}) + l(\mathbf{v}) & \forall \mathbf{v} \in V(\mathbf{0}) \\ m(\gamma, \tau) = b(\tau, \mathbf{u}) & \forall \tau \in T(\Omega) \\ \gamma + s\mathbf{t}_d \in -\mathcal{N}_{\mathcal{G}_{\mu, \text{Bi}}}(\lambda). \end{cases} \quad (6.27)$$

□

The  $F$  operator, as defined in Remark 6.1, allows us to define a Cadoux-like fixed-point algorithm.

**Remark 6.2.** We used the coercivity of the bilinear form “ $a$ ” to ensure the existence of a solution to each feasible optimization problem, and thus the well-posedness of the fixed-point algorithm. This condition is sufficient but not necessary, and in practice the method could be still be attempted for  $\text{Re} = +\infty$ .

**Existence of a solution to the variational problem ?** In the discrete case, Cadoux (2009) proves the existence of a fixed-point for the  $F$  operator under the hypothesis  $\mathcal{H}(s)$  by demonstrating the continuity and boundedness of  $F$  and applying the Brouwer fixed-point theorem. In our continuous case, one can also show relatively easily (at least for  $\text{Bi} = 0$ ) that the function  $q : L_2(\Omega) \rightarrow \mathbb{R}$  defined from the minimization problem (6.18) is continuous on  $S_+ := \{s \in L_2(\Omega), s \geq 0 \text{ a.e. on } \Omega\}$ , and that the application  $u : S_+ \rightarrow H^1(\Omega)^d$ , associating to  $s \in S_+$  the unique solution to the minimization problem (6.18), is continuous and bounded on  $S_+$ . Since  $F = f \circ u$ , with  $f : H^1(\Omega)^d \rightarrow S_+$ ,  $\mathbf{u} \mapsto (\mu - \zeta)|D(\mathbf{u})_T|$ , the existence of a solution to our flow problem under the hypothesis  $\mathcal{H}(s)$  would be ensured by the existence of a fixed-point to  $u \circ f$ . This application looks indeed like a good candidate;  $u \circ f$  is continuous and bounded on  $H^1(\Omega)^d$ , and the Rellich theorem states that its injection on  $L_2(\Omega)^d$  is compact. However, the continuity of  $u \circ f$  w.r.t. the  $L_2$ -norm does not seem obvious.

### 6.3 Discretization using finite-elements

In this section, we propose to discretize the variational problem (6.13a–6.13c) in such a way that the characterization of the solution as a fixed-point of a sequence of convex minimization problems (Remark 6.1) will still hold in the discrete case.

Moreover, we will see that the discrete version of the variational problem will have the same structure as the Discrete Coulomb Friction Problem (DCFP) that we studied in Chapters 2–4, allowing us to leverage similar algorithms.

#### 6.3.1 Discretization of the symmetric tensor fields

For the discretization of the space  $T(\Omega)$ , we shall make use of Lagrange FEM, which means that all symmetric tensor fields will be expressed as the extrapolation over  $\Omega$  of their values at  $n$  degrees of freedom  $(\mathbf{y}_i) \in \mathbb{R}^{nd}$ .

Let  $Q_h$  be a subspace of  $L_2(\Omega)$  with dimension  $n \in \mathbb{N}$ , and  $(\omega_i^\tau)_{1 \leq i \leq n}$  a basis of  $Q_h$  such that  $\omega_i^\tau(\mathbf{y}_j) = \delta_{i,j}$ . Let  $(\mathbf{S}_j)_{1 \leq j \leq s_d}$  be a basis of  $S_d$ , with  $s_d := \dim S_d = \frac{1}{2}d(d-1)$ . We may build a finite subspace  $T_h \subset T(\Omega)$  from the basis  $(\mathbf{T}_k)_{1 \leq k \leq s_d n}$  defined as

$$\mathbf{T}_{s_d(i-1)+j} := \omega_i^\tau \mathbf{S}_j \quad \text{for } 1 \leq i \leq n, 1 \leq j \leq s_d.$$

Then  $\forall \tau_h \in T_h, \exists \underline{\tau} := (\tau_k) \in \mathbb{R}^{n s_d}$  such that

$$\begin{aligned} \tau_h &= \sum_{\substack{1 \leq i \leq n \\ 1 \leq j \leq s_d}} \tau_{s_d(i-1)+j} \mathbf{T}_{s_d(i-1)+j} = \sum_{1 \leq i \leq n} \tau_h(\mathbf{y}_i) \omega_i^\tau, \\ \tau_h(\mathbf{y}_i) &= \sum_{1 \leq j \leq s_d} \tau_{s_d(i-1)+j} \mathbf{S}_j \quad \forall 1 \leq i \leq n. \end{aligned}$$

**Choice of basis for  $S_d$**  We want to leverage the analogy between the Drucker–Prager flow rule, which is expressed on tensors in  $S_d$ , and the Coulomb frictional contact law, which is expressed on vectors in  $\mathbb{R}^d$ . As such, we will now construct a basis of  $S_d$  such that the properties of the “normal” and “tangential” parts are preserved.

**Definition 6.1.** Let us introduce the morphism  $\chi$ ,

$$\begin{aligned} \chi : \mathbb{R}^{s_d} &\rightarrow S_d \\ (a; b, c) &\mapsto \begin{pmatrix} b & c \\ c & -b \end{pmatrix} + a \mathbb{I} & \text{if } d = 2 \\ (a; b, c, d, e, f) &\mapsto \begin{pmatrix} b - \frac{c}{\sqrt{3}} & d & e \\ d & -b - \frac{c}{\sqrt{3}} & f \\ e & f & \frac{2c}{\sqrt{3}} \end{pmatrix} + \frac{\sqrt{2}}{\sqrt{3}} a \mathbb{I} & \text{if } d = 3. \end{aligned}$$

The morphism  $\chi$  satisfies the following two properties:

**Property 6.1.**  $\chi$  is an orthonormal isomorphism between the two Euclidean spaces  $(\mathbb{R}^{s_d}; \cdot^T \cdot)$  and  $(S_d; \langle \cdot, \cdot \rangle)$ . This means

$$\forall (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{s_d} \times \mathbb{R}^{s_d} \quad \mathbf{x}^T \mathbf{y} = \langle \chi(\mathbf{x}), \chi(\mathbf{y}) \rangle, \quad (6.28)$$

where  $\mathbf{x}^T \mathbf{y}$  is the usual scalar product on  $\mathbb{R}^m$ ,  $m \geq 1$  and  $\langle \boldsymbol{\sigma}, \boldsymbol{\tau} \rangle = \frac{\boldsymbol{\sigma} : \boldsymbol{\tau}}{2}$  is our scalar product on  $S_d$ .

**Property 6.2.** Let  $(\cdot)_N$  and  $(\cdot)_T$  designate the normal and tangential parts of vectors and symmetric tensors as introduced in Chapter 1;

$$\begin{aligned} \mathbf{n} &:= (1, 0, \dots, 0) \in \mathbb{R}^{s_d} & \mathbf{x}_N &:= \mathbf{n}^T \mathbf{x} & \mathbf{x}_T &:= \mathbf{x} - \mathbf{x}_N \mathbf{n} & \forall \mathbf{x} \in \mathbb{R}^m \\ \boldsymbol{\iota}_d &:= \sqrt{\frac{d}{2}} \mathbb{I} \in S_d & \boldsymbol{\tau}_N &:= \langle \boldsymbol{\iota}_d, \boldsymbol{\tau} \rangle & \boldsymbol{\tau}_T &:= \boldsymbol{\tau} - \boldsymbol{\tau}_N \boldsymbol{\iota}_d = \text{Dev } \boldsymbol{\tau} & \forall \boldsymbol{\tau} \in S_d. \end{aligned}$$

Then  $\forall \mathbf{x} \in \mathbb{R}^m$ ,

$$\begin{cases} [\chi(\mathbf{x})]_N = \mathbf{x}_N \\ |[\chi(\mathbf{x})]_T| = \|\mathbf{x}_T\|. \end{cases} \quad (6.29a) \quad (6.29b)$$

It follows immediately from (6.29a–6.29b) that  $\boldsymbol{\sigma} \in \mathcal{T}_{\mu, \text{Bi}}^{(S_d)} \iff \chi^{-1}(\boldsymbol{\sigma}) \in \mathcal{T}_{\mu, \text{Bi}}^{(\mathbb{R}^{s_d})}$ . As  $\chi$  is orthogonal (Property 6.1), the inclusion in normal cones is also preserved. Finally, by linearity of  $\chi$ ,

$$\begin{aligned} &\boldsymbol{\gamma} + (\mu - \zeta) |\boldsymbol{\gamma}_T| \boldsymbol{\iota}_d \in -\mathcal{N}_{\mu, \text{Bi}}^{(S_d)}(\boldsymbol{\sigma}) \\ \iff &\chi^{-1}(\boldsymbol{\gamma} + (\mu - \zeta) |\boldsymbol{\gamma}_T| \boldsymbol{\iota}_d) \in -\mathcal{N}_{\mu, \text{Bi}}^{(\mathbb{R}^{s_d})}(\chi^{-1}(\boldsymbol{\sigma})) \\ \iff &\chi^{-1}(\boldsymbol{\gamma}) + (\mu - \zeta) \|[\chi^{-1}(\boldsymbol{\gamma})]_T\| \mathbf{n} \in -\mathcal{N}_{\mu, \text{Bi}}^{(\mathbb{R}^{s_d})}(\chi^{-1}(\boldsymbol{\sigma})). \end{aligned}$$

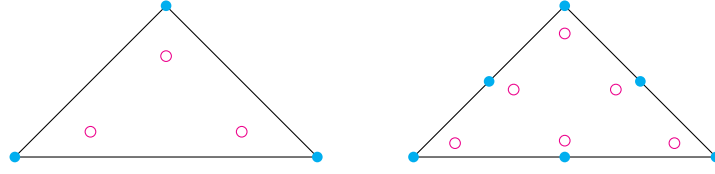
We will thus be able to write the Drucker–Prager flow rule indifferently on symmetric tensors or on their image in  $\mathbb{R}^{s_d}$  by the orthonormal isomorphism  $\chi^{-1}$ , and, with a slight abuse of notations, use the equivalence

$$(\boldsymbol{\gamma}, \boldsymbol{\lambda}) \in \mathcal{DP}(\mu, \text{Bi}, \zeta) \iff (\chi^{-1}(\boldsymbol{\gamma}), \chi^{-1}(\boldsymbol{\lambda})) \in \mathcal{DP}(\mu, \text{Bi}, \zeta). \quad (6.30)$$

Going back to the discretization of our symmetric tensor fields, the equivalence (6.30) motivates choosing  $\mathbf{S}_j := \chi(\mathbf{e}_j)$  as the basis for  $S_d$ , where  $(\mathbf{e}_j)_{1 \leq j \leq s_d}$  is the canonical basis of  $\mathbb{R}^{s_d}$ . Our global basis for  $T_h$  becomes

$$\mathbf{T}_{s_d(i-1)+j} := \omega_i^\tau \chi(\mathbf{e}_j) \quad \text{for } 1 \leq i \leq n, 1 \leq j \leq s_d.$$





**Figure 6.1:** Lagrange degrees of freedom (discs, blue) and Gauss quadrature points (circles, magenta) for two choices of triangular discretization. Left: piecewise-linear ( $\mathbb{P}_1$ ) shape function and order-2 quadrature. Right: piecewise-quadratic ( $\mathbb{P}_2$ ) shape function and order-4 quadrature.

The Drucker–Prager rheology at the degrees of freedom ( $\mathbf{y}_i$ ) can now be directly expressed on the vectors of coefficients of the strain rate and stress fields. Indeed, denoting by  $\underline{\tau}_{[i]} \in \mathbb{R}^{s_d}$  the segment of  $s_d$  coefficients of  $\underline{\tau}_h$  corresponding to the  $i^{\text{th}}$  degree of freedom,  $\underline{\tau}_{[i]} := [\tau_{s_d(i-1)+j}]_{1 \leq j \leq s_d}$ , then  $\forall 1 \leq i \leq n$ ,  $\underline{\tau}_h(\mathbf{y}_i) = \chi(\underline{\tau}_{h[i]})$ , and thus using (6.30),

$$(\gamma_h(\mathbf{y}_i), \lambda_h(\mathbf{y}_i)) \in \mathcal{DP}(\mu, \text{Bi}, \zeta) \iff (\underline{\gamma}_{[i]}, \underline{\lambda}_{[i]}) \in \mathcal{DP}(\mu, \text{Bi}, \zeta). \quad (6.31)$$

### 6.3.2 Discretization of the (bi)linear forms

Let  $U_h \subset H^1(\Omega)^d$  be a finite-dimensional vector space, and  $V_h \subset U_h$  its subspace satisfying the Dirichlet boundary conditions

$$V_h(\mathbf{u}_{D,h}) := \{\mathbf{v}_h \in U_h, \mathbf{v} = \mathbf{u}_{D,h} \text{ on } B_D\},$$

where  $\mathbf{u}_{D,h}$  discretizes  $\mathbf{u}_D$  on the FEM mesh. Let  $\mathbf{u}_{D,h}^* \in V_h(\mathbf{u}_{D,h}) \cap V_h(\mathbf{0})^\perp$  i.e., the unique function such that  $V_h(\mathbf{u}_{D,h}) = V_h(\mathbf{0}) + \{\mathbf{u}_{D,h}^*\}$ .

Unlike for the space  $T_h$ , here we make no specific assumption regarding the structure of  $V_h(\mathbf{0})$  or the construction of a corresponding basis. Let  $\mathbf{v} := \dim V_h$ . Given  $(\mathbf{V}_i)_{1 \leq i \leq \mathbf{v}}$  a basis of  $V_h(\mathbf{0})$ , we denote by  $A$ ,  $B$ , and  $M$  the matrices corresponding to the decomposition of the bilinear forms  $a$ ,  $b$ , and  $m$ , respectively, and by  $\underline{l}$  the vector corresponding to the decomposition of the linear form  $l$ . More precisely, we have  $A_{i,j} = a(\mathbf{V}_i, \mathbf{V}_j)$ ,  $M_{k,\ell} = m(\mathbf{T}_k, \mathbf{T}_\ell)$ ,  $B_{k,j} = b(\mathbf{T}_k, \mathbf{V}_j)$  and  $\underline{l}_j = l(\mathbf{V}_j)$ . Similarly, let  $\underline{u}$  be the vector of scalar coefficients corresponding to the decomposition of the projection of  $\mathbf{u}_h$  on the basis  $(\mathbf{V}_i)$  of  $V_h(\mathbf{0})$ . We have  $\mathbf{u} = \sum_{1 \leq i \leq \mathbf{v}} \underline{u}_i \mathbf{V}_i(x) + \mathbf{u}_{D,h}^*$ .

**FEM discrete system** At this point the discrete version of Equations (6.13a – 6.13b) reads: Find  $\mathbf{u}_h \in V_h(\mathbf{0})$ ,  $(\lambda_h, \gamma_h) \in T_h^2$ ,

$$\begin{cases} A \underline{u} = B^\top \underline{\lambda} + \underbrace{\underline{l} - \underline{a}_D}_{\underline{l}^{\text{tot}}} \end{cases} \quad (6.32a)$$

$$M \underline{\gamma} = B \underline{u} + \underline{k} \quad (6.32b)$$

$$(\gamma_h, \lambda_h) \in \text{discrete version of } \mathcal{DP}(\mu, \text{Bi}, \zeta) \quad (6.32c)$$

where for  $1 \leq i \leq \mathbf{v}$ ,  $\underline{a}_{D,i} = a(\mathbf{u}_{D,h}^*, \mathbf{V}_i)$  and for  $1 \leq k \leq s_d n$ ,  $\underline{k}_k = b(\mathbf{T}_k, \mathbf{u}_{D,h}^*)$ . It now remains to discretize the rheology constraint,  $(\gamma, \lambda) \in \mathcal{DP}(\mu, \text{Bi}, \zeta)$ , that is, to write an explicit expression for Equation (6.32c).

### 6.3.3 Discretization of the Drucker–Prager flow rule

As  $\mathcal{DP}(\mu, \text{Bi}, \zeta)$  is non-convex, the inclusion  $(\gamma_h, \lambda_h) \in \mathcal{DP}(\mu, \text{Bi}, \zeta)$ , where the values of  $\gamma_h$  and  $\lambda_h$  are extrapolated between degrees of freedom, will not be able to hold at every point in space. Instead, we will attempt to satisfy the rheology in a weaker sense.

Remember from Chapter 1 that the Drucker–Prager flow rule can be expressed as a root-finding problem, for instance on the De Saxcé complementarity function (1.29),

$$(\gamma, \lambda) \in \mathcal{DP}(\mu, \text{Bi}, \zeta) \iff f_{\text{DS}}(\gamma, \lambda) = \mathbf{0}.$$

We can thus write a weak form of (6.32c) as

$$\int_{\Omega} f_{\text{DS}}(\gamma_h, \lambda_h) : \tau_h = 0 \quad \forall \tau_h \in T_h,$$

which amounts to saying that the projection of  $f_{\text{DS}}(\gamma_h, \lambda_h)$  on  $T_h$  should vanish.

If  $f_{\text{DS}}$  were piecewise-polynomial, we could reduce these integrals to the evaluation of the integrand at a discrete number of quadrature points  $(\check{\mathbf{x}}_q)$ , and thus discretize the constraint as

$$f_{\text{DS}}(\gamma_h(\check{\mathbf{x}}_q), \lambda_h(\check{\mathbf{x}}_q)) = 0 \quad \forall q. \quad (6.33)$$

Obviously,  $f_{\text{DS}}$  is not piecewise-polynomial, but we shall still restrict ourselves to a discrete number of points  $\check{\mathbf{x}}_q$  at which to enforce  $(\gamma_h(\check{\mathbf{x}}_q), \lambda_h(\check{\mathbf{x}}_q)) \in \mathcal{DP}(\mu, \text{Bi}, \zeta)$  — in a sense, saying that a “projection” of  $f_{\text{DS}}$  onto a polynomial space should be zero. The remaining question is what makes a good choice for this set of points  $(\check{\mathbf{x}}_q)$ . In the following we shall see that this choice has a direct impact on the final form of the numerical system that must be solved, and thus on both the physical relevance of the discrete problem and the computational performance of solving methods.

**Discretization on  $Q_h$ ’s Lagrange degrees of freedom** One obvious choice for  $(\check{\mathbf{x}})_q$  is to consider the  $n$  degrees of freedom  $(\mathbf{y}_i)$  which served to define the  $(\omega^\tau)_i$  basis and thus the finite-dimensional spaces  $Q_h$  and  $T_h$ . Indeed, let  $\tau_h \in T_h$  such that  $\tau(\mathbf{y}_i) = \mathbf{0} \quad \forall 1 \leq i \leq n$ , then  $\tau_h = \mathbf{0}$ . From (6.31), replacing (6.32c) with

$$(\gamma_{[i]}, \lambda_{[i]}) \in \mathcal{DP}(\mu_i, \text{Bi}_i, \zeta_i) \quad \text{for } 1 \leq i \leq n.$$

thus means that any  $\tau_f \in T_h$  interpolating  $f_{\text{DS}}(\gamma_h, \lambda_h)$  at the degrees of freedom satisfies  $\tau_h = \mathbf{0}$ .

With this choice of discretization for the constraint, the discrete system (6.32a–6.32c) looks a lot like a DCFP, with one notable difference: the presence of the  $M$  matrix. Since  $A$  and  $M$  are positive-definite, we may eliminate the velocity variable  $\underline{\mathbf{u}}$  from (6.32a), and then get a linear relationship between  $\underline{\lambda}$  and  $\underline{\gamma}$  from (6.32b),  $\underline{\gamma} \propto W \underline{\lambda}$  where  $W = M^{-1}BA^{-1}B^\top$ . Continuing the DCFP analogy, the “Delassus” operator  $W$  of (6.32a–6.32c) is generally not symmetric, as  $M$  and  $BA^{-1}B^\top$  do not necessarily commute. This lack of symmetry is problematic; indeed, consider the case  $\mu = \zeta$ , the discrete problem reduces to

$$W \underline{\lambda} + \underline{\mathbf{b}} \in -\mathcal{N}_{\mathcal{D}_{\mu, \text{Bi}}}(\underline{\lambda}),$$

with  $\underline{\mathbf{b}}$  a constant vector in  $\mathbb{R}^{nd}$ . As  $W$  is asymmetric, this does not correspond to the optimality conditions of a convex minimization problem; more generally, the Cadoux algorithm that we devised in the continuous case can no longer be applied.

The discrete system (6.32) therefore lacks a fundamental symmetry property, which is key not only to guarantee physical consistency of our model, but also to the design of efficient numerical solvers. From a physical point of view, such an asymmetry in our discrete frictional contact law typically implies that the maximum dissipation principle cannot be satisfied, meaning that some anisotropy is artificially introduced through the discretization. From a purely numerical point of view, symmetry of the Delassus operator is not necessarily a prerequisite to common numerical solvers, but in our case it proves to be highly desirable for coming up with a tractable solving method. Indeed, among scalable solvers for DCFP, we basically have the choice between operator-splitting algorithms (Section (3.4)), and optimization-based algorithms (Sections 3.2 and 3.3). On the one hand, as mentioned above the asymmetric nature of

$W$  disqualifies the interpretation of as the fixed-point of a sequence of minimization problems, and thus optimization-based methods have to be discarded. On the other hand, the Gauss–Seidel splitting algorithm does not require  $W$  to be symmetric (except for the interpretation as block coordinate-descent), yet for efficiency purposes it requires the explicit knowledge of  $W$  or at least a very cheap way to compute matrix-vector products with  $A^{-1}B^\top$ . In our case  $A^{-1}$  is dense, and so is  $W$ , making the Gauss–Seidel algorithm intractable, especially for the large systems that we are going to have to deal with. For these reasons, we choose to discretize our constraints on an alternative set of points  $(\check{\mathbf{x}})_q$  that will allow us to eliminate the matrix  $M$  and thus retrieve symmetry of  $W$ . This way we shall both recover physical consistency of our model, and benefit from efficient optimization-based solving methods.

**Discretization on Gauss quadrature points** As we are using a Lagrange FEM discretization of the space  $T(\Omega)$  with polynomial interpolating bases, each integral  $M_{k,\ell} = \int \mathbf{T}_k : \mathbf{T}_\ell$  can be computed exactly using Gaussian quadrature of order  $2r$ , where  $r$  is the degree of the piecewise polynomials in  $T_h$ . That is,  $M_{k,\ell} = \sum_q w_q (\mathbf{T}_k(\hat{\mathbf{x}}_q) : \mathbf{T}_\ell(\hat{\mathbf{x}}_q))$  where the  $(\hat{\mathbf{x}}_q)$ ,  $1 \leq q \leq n_Q$  are the so-called Gauss quadrature points and  $(w_q) \in \mathbb{R}^{n_Q}$  are their corresponding weights. Similarly, let  $\sigma_h \in T_h$  such that  $\sigma_h(\hat{\mathbf{x}}_q) = \mathbf{0} \ \forall 1 \leq q \leq n_Q$ ; then  $\forall \tau_h \in T_h$ ,  $\int_\Omega \sigma_h : \tau_h = 0$ . Hence, enforcing

$$(\gamma_h(\hat{\mathbf{x}}_q), \lambda_h(\hat{\mathbf{x}}_q)) \in \mathcal{DP}(\mu, \text{Bi}, \zeta) \quad \forall 1 \leq q \leq n_Q$$

implies that any function  $\tau_h \in T_h$  interpolating  $f_{\text{DS}}(\gamma_h, \lambda_h)$  at the quadrature points  $(\hat{\mathbf{x}}_q)$  must be zero.

As shown below, defining  $(\check{\mathbf{x}})_q$ , the set of points at which the constraint is enforced, as the set of quadrature points  $(\hat{\mathbf{x}}_q)$  allows us to retrieve a symmetric Delassus operator.

Recall that for  $\tau$  a symmetric tensor field in  $T(\Omega)$ ,  $\tau_h \in T_h$  corresponds to its discretized version interpolating the values at the  $n$  degrees of freedom  $(\mathbf{y}_i)$ , with  $\tau(\mathbf{y}_i) = \chi(\underline{\tau}_{[i]})$ . Let  $R$  be the  $(s_d n_Q \times s_d n)$  matrix mapping  $\underline{\tau}$  to the interpolated values  $\hat{\underline{\tau}}$  at the quadrature points  $(\hat{\mathbf{x}}_q)$ . That is,  $R$  is such that  $\hat{\underline{\tau}} = R \underline{\tau}$ , with  $\hat{\underline{\tau}} := (\hat{\underline{\tau}}_{[q]})_{1 \leq q \leq n_Q} := \chi^{-1}(\tau(\hat{\mathbf{x}}_q))$ . The matrix  $R$  contains  $n_Q \times n$  square blocks  $\mathfrak{R}_{q,j}$  of size  $s_d \times s_d$  with  $\mathfrak{R}_{q,p} = \omega_i^\tau(\hat{\mathbf{x}}_q) \mathbb{I}_{s_d}$  for all  $1 \leq q \leq n_Q$  and  $1 \leq i \leq n$ . This translates into the following coefficient-wise expression

$$R_{(q-1)s_d+p, (i-1)s_d+j} = \omega_i^\tau(\hat{\mathbf{x}}_q) \delta_{p,j} \quad \text{for } 1 \leq p, \ell \leq s_d. \quad (6.34)$$

**Proposition 6.3.** Let  $\hat{\underline{\lambda}} := R \underline{\lambda}$ , and let  $R^\dagger$  denote the Moore–Penrose pseudoinverse of  $R$ . Then Equations (6.13a – 6.13c) can be discretized as:

Find  $\underline{\mathbf{u}} \in \mathbb{R}^v$ ,  $\hat{\underline{\lambda}}, \hat{\underline{\gamma}} \in \mathbb{R}^{n_Q s_d}$ ,

$$\begin{cases} A \underline{\mathbf{u}} = B^\top R^\dagger \hat{\underline{\lambda}} + \underline{\mathbf{l}}^{\text{tot}} \end{cases} \quad (6.35a)$$

$$\begin{cases} \hat{\underline{\gamma}} = R^{\dagger, \top} B \underline{\mathbf{u}} + R^{\dagger, \top} \underline{\mathbf{k}} \end{cases} \quad (6.35b)$$

$$\begin{cases} (\hat{\underline{\gamma}}_{[q]}, \hat{\underline{\lambda}}_{[q]}) \in \mathcal{DP}(\mu, \text{Bi}, \zeta) \end{cases} \quad \forall 1 \leq q \leq n_Q. \quad (6.35c)$$

$\hat{\underline{\gamma}}$  defined as in (6.35b) is such that for  $1 \leq q \leq n_Q$ ,  $\hat{\underline{\gamma}}_{[q]} = 2w_q \gamma_h(\mathbf{x}_q)$ . The corresponding proof is given in Appendix C.1. This time we have obtained in (6.35) a system which preserves the symmetry of the new Delassus operator  $W = R^{\dagger, \top} B A^{-1} B^\top R^\dagger$ ; we will be able to express the solutions to (6.35a–6.35c) as a sequence of convex minimization problems, which are discrete version of (6.18).

Indeed, the convexified version of (6.35a–6.35c), i.e., with  $\mu = \zeta$ , reduces to

$$A \underline{\mathbf{u}} - \underline{\mathbf{l}}^{\text{tot}} \in -B^\top R^\dagger \mathcal{N}_{\frac{1}{\mu}}(R^{\dagger, \top} (B \underline{\mathbf{u}} + \underline{\mathbf{k}})). \quad (6.36)$$

Under the hypothesis  $\exists \underline{\mathbf{u}} \in \mathbb{R}^{\mathfrak{v}}, (R^{\dagger, \top} (B \underline{\mathbf{u}} + \underline{\mathbf{k}})) \in \text{int } \mathcal{K}_{\mu}^1$ , the convexified problem (6.36) reads as the optimality conditions of the minimization problem

$$\begin{aligned} \min_{\underline{\mathbf{u}} \in C} J(\underline{\mathbf{u}}), \quad (6.37) \\ J(\underline{\mathbf{u}}) &:= \frac{1}{2} \underline{\mathbf{u}}^{\top} A \underline{\mathbf{u}} + \underline{\mathbf{u}}^{\top} \underline{\mathbf{l}} + \text{Bi} \sum_{q=1}^{n_Q} \left\| \hat{\underline{\mathbf{f}}}(\underline{\mathbf{u}})_{q\top} \right\| \\ \hat{\underline{\mathbf{f}}}(\underline{\mathbf{u}}) &:= R^{\dagger, \top} B \underline{\mathbf{u}} + R^{\dagger, \top} \underline{\mathbf{k}} \\ C &:= \left\{ \underline{\mathbf{u}} \in \mathbb{R}^{\mathfrak{v}}, \hat{\underline{\mathbf{f}}}(\underline{\mathbf{u}}) \in \mathcal{K}_{\mu}^1 \right\}. \end{aligned}$$

Problem (6.37) corresponds to a discretization of the continuous minimization problem (6.18) where the Dirichlet boundary conditions are enforced intrinsically and the strain rate  $D(\mathbf{u})$  is evaluated at the Gauss quadrature points. Indeed,

$$\sum_q \|\hat{\underline{\mathbf{f}}}(\underline{\mathbf{u}})_{q\top}\| = \sum_q |2w_q \gamma_h(\hat{\mathbf{x}}_q)_{\top}| = 2 \sum_q w_q |\Pi_{T_h}(D(\mathbf{u}_h)_{\top})(\hat{\mathbf{x}}_q)| \sim 2 \int_{\Omega} |D(\mathbf{u}_h)_{\top}|.$$

One remaining difficulty stems from the presence of the matrix  $R^{\dagger}$ , which in the general case could substantially increase the cost of solving the system. We present in below a few cases for which this difficulty vanishes.

### 6.3.4 Considerations on $R^{\dagger}$

The first observation is that if the quadrature points  $(\hat{\mathbf{x}}_q)$  were to coincide with the degrees of freedom  $(\mathbf{y}_i)$ ,  $R$  would boil down to the identity matrix and the operator  $R^{\dagger}$  would not induce any additional cost. This is actually always the case for a piecewise constant ( $\mathbb{P}_0$ ) approximation, for which both the degrees of freedom and the Gauss quadrature points are located at the barycenter of each element.

**Trapezoidal quadrature rule** For higher-order polynomial basis functions, having the  $(\hat{\mathbf{x}})_q$  coincide with the  $(\mathbf{y})_i$  amounts to computing  $m(\gamma, \tau)$  using a trapezoidal integration rule. Obviously, such an approximation induces a loss of precision — the integral being exact only for functions that are linear between the degrees of freedom. This means that the order of convergence will not increase with that of the basis functions, and using high-order discretization space ( $\mathbb{P}_2$  or higher-order polynomials) would be wasteful. However, for piecewise linear ( $\mathbb{P}_1$ ) polynomials, we found this approximation to be acceptable, and used it in practice.

**Mixed finite elements** In the case of piecewise-polynomial discontinuous basis functions, degrees of freedom are not shared between adjacent elements. When considering such a discretization of the space  $T(\Omega)$ , the matrix  $R$  becomes block-diagonal, and consequently its pseudo-inverse has a similar structure and is easy to compute. The additional cost induced by the presence of the linear operator  $R^{\dagger}$  in Problem (6.35) is therefore once again negligible.

**Factorization of the  $B$  matrix** If the discrete velocity space  $V_h(\mathbf{0})$  is such that  $D(\mathbf{u}_h)$  is piecewise-polynomial of order less than that of  $T_h$ , then the same quadrature rule can be used to compute the bilinear form  $b$ . That is,  $b(\mathbf{T}_k, \mathbf{V}_j) = \sum_q w_q (\mathbf{T}_k)(\hat{\mathbf{x}}_q) : D(\mathbf{V}_j)(\hat{\mathbf{x}}_q)$ , and thus  $B = R^{\top} H$ , with  $H$  a  $n_Q s_d \times \mathfrak{v}$  matrix given by

$$H_{(q-1)s_d+p,j} := 2w_q \left[ \chi^{-1}(D(\mathbf{V}_j))(\hat{\mathbf{x}}_q) \right]_p.$$

*Proof.*

$$\begin{aligned}
(R^\top H)_{i s_d+k,j} &= \sum_{q=1}^{n_Q} \sum_{p=1}^{s_d} R_{(q-1)s_d+p,(i-1)s_d+k} H_{(q-1)s_d+p,j} \\
&= \sum_{q=1}^{n_Q} \sum_{p=1}^{s_d} 2w_q \alpha_i(\hat{\mathbf{x}}_q) \delta_p^k [\chi^{-1}(D(\mathbf{V}_j)(\hat{\mathbf{x}}_q))]_p \\
&= \sum_{q=1}^{n_Q} 2w_q \alpha_i(\hat{\mathbf{x}}_q) \langle \chi(\mathbf{e}_k), D(\mathbf{V}_j)(\hat{\mathbf{x}}_q) \rangle \\
&= \sum_{q=1}^{n_Q} w_q \mathbf{T}_{(i-1)s_d+k}(\hat{\mathbf{x}}_q) : D(\mathbf{V}_j)(\hat{\mathbf{x}}_q) = b(\mathbf{T}_{(i-1)s_d+k}, \mathbf{V}_j) = B_{(i-1)s_d+k,j}.
\end{aligned}$$

□

In this case, a discrete version of (6.13a–6.13c) can be written simply as:  
Find  $\underline{\mathbf{u}} \in \mathbb{R}^v$ ,  $\underline{\hat{\boldsymbol{\lambda}}}, \underline{\hat{\boldsymbol{\gamma}}} \in \mathbb{R}^{n_Q s_d}$ ,

$$\begin{cases} A \underline{\mathbf{u}} = H^\top \underline{\hat{\boldsymbol{\lambda}}} + \underline{\mathbf{l}}^{\text{tot}} \\ \underline{\hat{\boldsymbol{\gamma}}} = H \underline{\mathbf{u}} + \underline{\hat{\mathbf{k}}} \\ (\hat{\boldsymbol{\gamma}}_{[q]}, \hat{\boldsymbol{\lambda}}_{[q]}) \in \mathcal{DP}(\mu, \text{Bi}, \zeta) \quad \forall 1 \leq q \leq n_Q, \end{cases}$$

where

$$\underline{\hat{\mathbf{k}}}_{q s_d+p} := 2w_q [\chi^{-1}(D(\mathbf{u}_{D,h}^*)(\hat{\mathbf{x}}_q))]_p,$$

and the pseudo-inverse needs not be computed. However, note that unlike the previous methods which can be used in conjunction with black-box finite-element libraries, this approach requires being able to modify the assembly of matrix corresponding to the “ $b$ ” bilinear form.

### 6.3.5 Final discrete system

For brevity of notation and since there are no more ambiguities, from now on we shall drop the decorations of the variables, i.e., we shall consider Problem (6.35) written as:

Find  $\mathbf{u} \in \mathbb{R}^v$ ,  $\boldsymbol{\lambda}, \boldsymbol{\gamma} \in \mathbb{R}^{n_Q s_d}$ ,

$$\begin{cases} A \mathbf{u} = B^\top R^\top \boldsymbol{\lambda} + \mathbf{l} \\ \boldsymbol{\gamma} = R^{\top,\top} B \mathbf{u} + R^{\top,\top} \mathbf{k} \\ (\boldsymbol{\gamma}_i, \boldsymbol{\lambda}_i) \in \mathcal{DP}(\mu_i, \text{Bi}_i, \zeta_i) \quad \forall 1 \leq i \leq n_Q. \end{cases} \quad (6.38)$$

Under this form, the similarity of (6.38) with the DCFP (2.17) is clear. There remains only two differences:

- The constraint is more general; we may have to adapt our DCFP solvers to handle the dilatancy  $\zeta$  and a truncated SOC when  $\text{Bi} > 0$ ;
- In contrast with DEM methods where the global stiffness matrix was block-diagonal, and similarly to the case of cloth dynamics, here  $A$  is sparse but with a dense inverse.

**Discussion on the discretization strategy** At no point have we attempted to study the rate of convergence of our finite-element discretization; in particular, we did not specify any *inf-sup*-like condition. It is probable that one should be required; consider homogeneous Dirichlet boundary conditions, then our unilateral incompressibility constraint becomes equivalent to full incompressibility. To mitigate this consideration, in our Results section we present convergence

tests on simple Bingham Poiseuille flows for which we have an analytical solution, and observe convergence of the velocity solution in  $L_2$ - and  $H^1$ -norms.

The discrete setting will also allow us to formulate stronger existence properties than for the continuous problem, using the criterion from of Cadoux (2009).

## 6.4 Solving the discrete problem

In this section, we are concerned with the numerical resolution of System (6.38).

### 6.4.1 Discrete Cadoux fixed-point algorithm

We have already hinted in the previous section that the solutions to System (6.38) can be characterized as the fixed points of a sequence of minimization problems structurally similar to Problem (6.37); we demonstrate this assertion below.

Indeed, let  $S$  denote the subspace of  $\mathbb{R}^{n_Q \times d}$  with vanishing tangential components,  $S = \{\mathbf{s} \in \mathbb{R}^{n_Q \times d}, \mathbf{s}_{\text{IT}} = \mathbf{0} \forall 1 \leq i \leq n_Q\}$ , and  $S_+$  the positive subset of  $S$ ,  $S_+ = \{\mathbf{s} \in S, \mathbf{s}_{\text{IN}} \geq 0 \forall 1 \leq i \leq n_Q\}$ . For any  $\mathbf{s} \in S$ , consider the minimization problem (6.39),

$$\begin{aligned} q(\mathbf{s}) &:= \min_{\mathbf{u} \in C(\mathbf{s})} J(\mathbf{u}), \\ J(\mathbf{u}) &:= \frac{1}{2} \mathbf{u}^\top \mathbf{A} \mathbf{u} + \mathbf{u}^\top \mathbf{l} + \text{Bi} \sum_{i=1}^{n_Q} \|\boldsymbol{\gamma}(\mathbf{u})_{\text{IT}}\| \\ \boldsymbol{\gamma}(\mathbf{u}) &:= \mathbf{R}^{\dagger, \top} \mathbf{B} \mathbf{u} + \mathbf{R}^{\dagger, \top} \mathbf{k} \\ C(\mathbf{s}) &:= \left\{ \mathbf{u} \in \mathbb{R}^v, \boldsymbol{\gamma}(\mathbf{u}) + \mathbf{s} \in \mathcal{K}_{\frac{1}{\mu}} \right\}. \end{aligned} \quad (6.39)$$

Proceeding as in the discrete DEM (Section 2.3.2) or continuous FEM (Section 6.2.3) cases, we can show that any  $\mathbf{u} \in \mathbb{R}^v$  satisfying (6.40),

$$\mathbf{A} \bar{\mathbf{u}} - \mathbf{l} \in -\mathbf{B}^\top \mathbf{R}^\dagger \left[ \mathcal{N}_{\mathcal{K}_{\frac{1}{\mu}}}(\boldsymbol{\gamma}(\mathbf{u}) + \mathbf{s}) + \text{Bi} \partial \|\cdot\|_{\text{T}}(\boldsymbol{\gamma}(\mathbf{u})) \right] \quad (6.40)$$

is a solution of (6.39). This sufficient condition is also necessary when  $\mathcal{H}(\underline{\mathbf{s}})$  holds,

$$\exists \mathbf{u} \in \mathbb{R}^v, \boldsymbol{\gamma}(\mathbf{u}) + \mathbf{s} \in \text{int } \mathcal{K}_{\frac{1}{\mu}}. \quad (\mathcal{H}(\underline{\mathbf{s}}))$$

Moreover, Equation (6.40) can be equivalently rewritten as

$$\begin{cases} \mathbf{A} \mathbf{u} = \mathbf{B}^\top \mathbf{R}^\dagger \boldsymbol{\lambda} + \mathbf{l} \\ \boldsymbol{\gamma} = \mathbf{R}^{\dagger, \top} \mathbf{B} \mathbf{u} + \mathbf{R}^{\dagger, \top} \mathbf{k} \\ \boldsymbol{\lambda}_i \in -\mathcal{N}_{\mathcal{K}_{\frac{1}{\mu_i}}}(\boldsymbol{\gamma}_i + \mathbf{s}_i) - \text{Bi} \partial \|\cdot\|_{\text{T}}(\boldsymbol{\gamma}_i) \quad \forall 1 \leq i \leq n_Q. \end{cases} \quad (6.41)$$

Using once again Equation (1.25) to rewrite the last line of (6.41),

$$\begin{aligned} \boldsymbol{\lambda}_i \in -\mathcal{N}_{\mathcal{K}_{\frac{1}{\mu_i}}}(\boldsymbol{\gamma}_i + \mathbf{s}_i) - \text{Bi} \partial \|\cdot\|_{\text{T}}(\boldsymbol{\gamma}_i) &\iff \boldsymbol{\lambda}_i \in -\partial \left( \mathcal{G}_{\mathcal{K}_{\frac{1}{\mu_i}}, \text{Bi}}^* \right) (\boldsymbol{\gamma}_i + \mathbf{s}_i) \\ &\iff \boldsymbol{\gamma}_i + \mathbf{s}_i \in \mathcal{N}_{\mathcal{G}_{\mu_i, \text{Bi}}}(-\boldsymbol{\lambda}), \end{aligned}$$

and recalling the definition (1.19) of the Drucker–Prager flow rule, one can check that Equation (6.40) coincide with System (6.38) when  $\mathbf{s}_{\text{IN}} = (\mu - \zeta) \|\boldsymbol{\gamma}(\mathbf{u})_{\text{IT}}\|$  for any  $1 \leq i \leq n_Q$ . In other words, we can once again solve our discrete flow problem (6.38) using a fixed point algorithm.

**Property 6.3** (Cadoux fixed-point algorithm). *Let  $u : S_+ \rightarrow \mathbb{R}^v$  be a function associating to each  $s \in S_+$  a solution to the minimization problem (6.39), i.e.,  $u(s) \in C(s)$  and  $J(u(s)) = q(s)$ . Let  $s : \mathbb{R}^v \rightarrow S_+$  such that  $\forall \mathbf{u} \in \mathbb{R}^v$ ,  $\forall 1 \leq i \leq n_Q$ ,  $s(\mathbf{u})_{iN} = (\mu - \zeta) \|\gamma(\mathbf{u})_{iT}\|$ .*

*The Cadoux fixed-point algorithm, induced by  $s^{k+1} := s \circ u(s^k)$  and an initial guess  $s^0 \in S_+$ , will be well-defined as long as  $C(\mathbf{0}) \neq \emptyset$  and  $\text{Re} < +\infty$ .*

*Moreover, let  $\bar{s}$  be a fixed-point of  $s \circ u$  satisfying  $\mathcal{H}(\bar{s})$ ; then the minimization problem  $q(\bar{s})$  defined in (6.39) admits a unique solution, and its optimality conditions (6.41) realize a solution  $(\bar{\mathbf{u}}, \bar{\lambda}, \bar{\gamma})$  of the discrete flow problem (6.38).*

Indeed, if  $C(\mathbf{0})$  is non-empty, then all the minimization problems in the sequence generated by the fixed-point algorithm are feasible. If  $\text{Re} < +\infty$ , the objective function is strictly convex and coercive, hence each intermediate problem admits a unique solution, and the function  $u$  is well-defined.

The following existence property can also be directly adapted from (Cadoux 2009, Theorem 3.19).

**Property 6.4** (Cadoux existence criterion). *If  $C(\mathbf{0})$  is non-empty, then the Cadoux fixed point algorithm defined in Property 6.3 admits a fixed-point  $\bar{s}$ .*

Indeed, while Cadoux (2009) did not treat the case  $\text{Bi} > 0$ , as the strict convexity, coercivity and continuity of  $J$  are preserved, his analysis holds nevertheless without modifications.

**Homogeneous Dirichlet boundary conditions** Let us emphasize once again that this existence criterion is only sufficient, and not necessary. In the case  $\mu = 0$ , the strong hypothesis  $\mathcal{H}(\mathbf{0})$  amounts to the existence of a velocity field with strictly positive divergence everywhere – this requires outwards Dirichlet boundary conditions, or a Neumann boundary. The weaker hypothesis in Property 6.4 only requires that there exists a velocity field with nowhere strictly negative divergence. Homogeneous Dirichlet boundary conditions only satisfy the latter, weaker criterion  $C(\mathbf{0}) \neq \emptyset$ . However we can always exhibit a trivial solution in  $\mathbf{u}, \lambda$ ; for  $\mu = 0$ , it corresponds to the incompressible Stokes solution, and for  $\mu > 0$ , to the null velocity solution,  $\mathbf{u} = \mathbf{0}$ .

*Proof.* Let us use our original notations, with the underline to differentiate the concatenated coefficient vectors from the discretized fields. By construction, we have  $\underline{\mathbf{k}} = \mathbf{0}$ . Moreover, the homogeneous Dirichlet condition imposes  $\int_{\Omega} \nabla \cdot \mathbf{u}_h = \int_{\partial\Omega} \mathbf{u}_{D,h} \cdot \mathbf{n}_{\Omega} = 0$ . As the rheology requires  $\nabla \cdot \mathbf{u}_h \geq 0$ , velocity solutions must have null divergence, which means  $(R^{\dagger, \top} B \underline{\mathbf{u}})_{\mathbf{N}} = \mathbf{0}$ .

In the case  $\mu = 0$ ,  $\bar{s} := \mathbf{0}$  is a fixed-point of  $s \circ u$ , and the primal feasibility condition  $\underline{\mathbf{u}} \in C(\mathbf{0})$  boils down to  $(R^{\dagger, \top} B \underline{\mathbf{u}})_{\mathbf{N}} = \mathbf{0}$ . The optimality condition reads

$$-\nabla J \in \mathcal{N}_{C(\mathbf{0})} \underline{\mathbf{u}} = (\text{Ker}(R^{\dagger, \top} B)_{\mathbf{N}})^{\perp} = \text{Im}(R^{\dagger, \top} B)_{\mathbf{N}}^{\top},$$

so there exists  $\underline{\lambda} \in \mathbb{R}^{n_Q s_d}$  with  $\nabla J = B^T R^{\dagger} \underline{\lambda}$  and  $\underline{\lambda}_{\mathbf{T}} = \mathbf{0}$ .

In the case  $\mu > 0$ ,  $\underline{\mathbf{u}} \in C(\mathbf{0})$  boils down to  $R^{\dagger, \top} B \underline{\mathbf{u}} = \mathbf{0}$ . This means that the solution  $u(\mathbf{0})$  of the primal problem at  $\mathbf{s} = \mathbf{0}$  satisfies  $\|(B^{\top} R^{\dagger} \underline{\mathbf{u}})_{\mathbf{T}}\| = 0$ , and therefore  $\bar{s} := \mathbf{0}$  is a fixed point of  $s \circ u$ . The optimality condition yields  $\exists \underline{\lambda} \in \mathbb{R}^{n_Q s_d}$ ,  $\nabla J = B^T R^{\dagger} \underline{\lambda}$ .

In both cases, we define  $\underline{\lambda}_h^* := \underline{\lambda}_h + \Delta_p \mathbf{t}_d$ , with  $\Delta_p$  constant over all discretization points. By construction,  $\underline{\lambda} - \underline{\lambda}^* \in \text{Ker } B^T R^{\dagger}$ . We can choose  $\Delta_p$  such that  $\underline{\lambda}^* \in \mathcal{K}_{\mu}$ ; for  $\mu = 0$ , it suffices to take  $\Delta_p := -\min_i \underline{\lambda}_{iN}$ , and for  $\mu > 0$ ,  $\Delta_p := -\min_i (\underline{\lambda}_{iN} - \frac{1}{\mu} \|\underline{\lambda}_{iT}\|)$ . We have then  $A \underline{\mathbf{u}} - \underline{\mathbf{l}} = \nabla J = B^T R^{\dagger} \underline{\lambda}^*$ , and it can be easily verified that in both cases  $\underline{\gamma}^{\top} \underline{\lambda}^* = 0$ .  $\square$

#### 6.4.2 Dual problem

If the matrix  $A$  is invertible (which, once again, is achieved as soon as  $\text{Re} < +\infty$ ) then the velocity variable  $\mathbf{u}$  can be eliminated in System (6.41), so that the latter system can be written equivalently as

$$\begin{cases} \gamma = W \lambda + b \\ \gamma_i + s_i \in -\mathcal{N}_{\mathcal{T}_{\mu_i, \text{Bi}}}(\lambda_i) \quad \forall 1 \leq i \leq n_Q, \end{cases} \quad (6.42)$$

where  $W := R^\top B A^{-1} B^\top R^\top$  and  $b := R^\top B (k + B A^{-1} l)$ . As  $W$  is symmetric positive semi-definite, we recognize (6.42) as the optimality conditions of the minimization problem (6.43),

$$\min_{\lambda \in \mathcal{T}_{\mu, \text{Bi}}} \frac{1}{2} \lambda^\top W \lambda + (b + s)^\top \lambda. \quad (6.43)$$

Problem (6.43) is the dual (in the sense of Fenchel) of the primal minimization problem (6.39). The Cadoux fixed-point algorithm defined in Property 6.3 can thus be iterated by solving the intermediate problems under either primal or dual form; the next section is dedicated to the numerical resolution of those minimization problems.

### 6.4.3 Solving the minimization problems

Problems (6.39) and (6.43) are very close to SOCQP (and actually are SOCQP when  $\text{Bi} = 0$ ); for the primal, the difference is just a modification of the objective function (which is no longer quadratic), and for the dual, only the constraint is modified by replacing the SOC with a truncated SOC. It is therefore likely that the algorithms presented in Chapter 3 for solving SOCQP can be easily adapted to our flow problems. However, remember that the inverse of our matrix  $A$  is dense, therefore we will not be able to use the Gauss–Seidel solver. Moreover, as we target relatively large meshes, we will not consider performing directly root-finding methods on complementarity functions. This leaves us with two strategies, interior-points (Section 3.2) and proximal algorithms (Section 3.3).

**Interior points** When  $\text{Bi} = 0$  the SOCP formulations of the primal and dual problems given in Section 3.2 can directly be used. Moreover, when  $\text{Bi} > 0$ , our modified minimization problems (6.39) and (6.43) can still be put under SOCP form.

Indeed, the “Bi” part of our primal objective function can again be transformed into a linear contribution with SOC constraints using a trick similar to that of the quadratic part. Indeed,  $\min_{x \in \mathbb{R}^{nd}} \sum_i \|x_i\|$  can be rewritten as  $\min_{r \in \mathbb{R}^{n_{r_i}}} \sum_i r_i$  such that  $r_i \geq \|x_i\| \quad \forall i$ . Let  $L$  be a square-root of  $A$ , i.e.,  $A = LL^\top$ , then a solution of the primal minimization problem (6.39) can be found by solving the SOCP

$$\left\{ \begin{array}{l} \min_{u \in \mathbb{R}^v, t \in \mathbb{R}, r \in \mathbb{R}^{n_Q}} t + \text{Bi} \sum_{i=1}^{n_Q} r_i - u^\top l \\ z = L^\top u \\ R^\top \gamma = B u + k + R^\top s \\ \gamma_i \in \mathcal{K}_{\frac{1}{\mu_i}} \quad \forall i = 1 \dots n_Q \\ (r_i, \gamma_{iT}) \in \mathcal{K}_1 \quad \forall i = 1 \dots n_Q \\ (1, t, z) \in \mathcal{R} \mathcal{K}. \end{array} \right. \quad (6.44)$$

The objective function of the dual minimization problem has not changed, however the constraint on  $\lambda$  is no longer conical. Yet this constraint can still be expressed using two SOC constraints:

$$\lambda \in \mathcal{T}_{\mu, \text{Bi}} \iff \begin{cases} \lambda = \lambda_K + \lambda_B \\ \lambda_{Ki} \in \mathcal{K}_{\mu_i} & \forall 1 \leq i \leq n_Q \\ \lambda_{BiT} \leq \text{Bi} \text{ and } \lambda_{BiN} = 0 & \forall 1 \leq i \leq n_Q \end{cases}$$



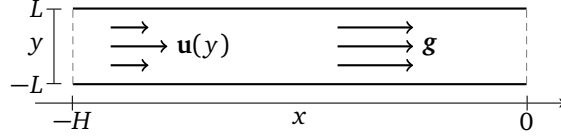


Figure 6.2: 2D channel for the Bingham Poiseuille flow

Using this insight, we can propose another SOCP for the dual minimization problem (6.43),

$$\left\{ \begin{array}{ll} \min_{\lambda_K, \lambda_B \in \mathbb{R}^{n_Q \times d}, t \in \mathbb{R}} & t + (\lambda_B + \lambda_K)^\top (b + s) \\ & Ry = \lambda_K + \lambda_B \\ & L^\top z = H^\top y \\ & \lambda_{Ki} \in \mathcal{K}_{\mu_i} & \forall i = 1 \dots n_Q \\ & (Bi, \lambda_{BiT}) \in \mathcal{K}_1 & \forall i = 1 \dots n_Q \\ & \lambda_{BiN} = 0 & \forall i = 1 \dots n_Q \\ & (1, t, z) \in \mathcal{RK}. \end{array} \right. \quad (6.45)$$

**Proximal algorithms** As computing an orthogonal projection on the truncated cone  $\mathcal{T}_{\mu, Bi}$  is not more complex than for  $\mathcal{K}_\mu$ , the projected gradient descent, projected-gradient algorithms and their variants (Section 3.3.2) for solving the dual minimization problem (2.21) can be trivially adapted to our modified dual (6.43). Note that as we do not want to explicitly assemble the Delassus operator  $W$ , each step of these algorithms will involve solving a linear system with matrix  $A$ . This can be done efficiently by precomputing a Cholesky factorization of  $A$ , or by using a truncated conjugate-gradient algorithm.

## 6.5 Results

All the finite element simulations presented in this section were performed using the open-source library Rheolef (Saramito 2015), and were run on a quad-core Intel(R) Xeon W3520 machine with 8GB memory.

### 6.5.1 Model problems

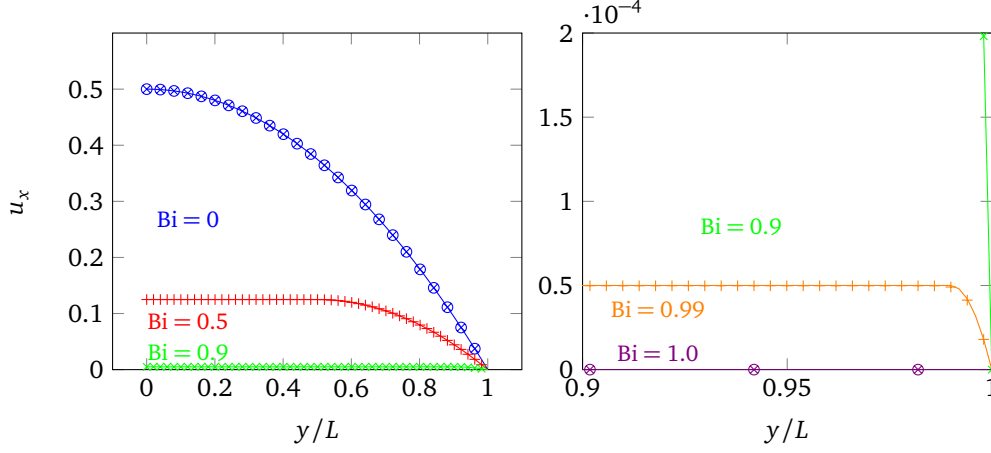
In this section, we validate our method on simple problems for which we can derive analytic solutions.

#### Bingham Poiseuille flow

We consider a Bingham Poiseuille flow as illustrated in Figure 6.2, with the following boundary conditions,

$$\begin{aligned} u(x, -L) &= u(x, L) = 0 \quad \forall x \in [-H, 0] \\ u(-H, y) &= u_x(-H, y)e_x \quad \forall y \in [-L, L] \\ u(0, y) &= u_x(0, y)e_x \quad \forall y \in [-L, L]. \end{aligned}$$

Figure 6.3 shows velocity profiles for different values of  $Bi$ , for both the analytic solution and our model using a  $\mathbb{P}_4 - \mathbb{P}_{1d}$  approximation. Using a nonsmooth solver allows us to recover the correct profile even for high values of  $Bi$  without any parameter tuning.



**Figure 6.3:** Comparison of the steady velocity profiles  $u_x(-H/2, y)$  between our numerical model using a  $\mathbb{P}_2 - \mathbb{P}_{1d}$  approximation (marks), and the analytic Bingham Poiseuille flow (lines). Plots for high values of  $Bi$  (right) are zoomed in compared to plots for lower  $Bi$  values (left).

**Convergence of spatial discretization** We study how the error between our method and the analytic solutions for  $Bi = 0.5$  and  $Bi = 0.9$  decreases as we uniformly refine a mesh with initial characteristic edge length  $h_0$ . Results for various FEM approximation orders are shown in Figure 6.4. Convergence was observed for all approximation orders, and we found that  $\mathbb{P}_2 - \mathbb{P}_1$  and  $\mathbb{P}_2 - \mathbb{P}_{1d}$  achieved a good ratio of convergence speed versus computational cost. Conversely, the higher-order approximation  $\mathbb{P}_3 - \mathbb{P}_{2d}$  performed relatively poorly on the finer meshes, which we interpret as being the result of numerically more complex quadrature rules.  $\mathbb{P}_4 - \mathbb{P}_{1d}$ , with high-order velocities but low-order stresses, yielded consistently the best results, at the cost of a very large  $A$  matrix.

### Bagnold profiles

We consider the flow of a granular layer of height  $H = 1$  on a rough infinite inclined plane with angle  $\alpha$ , as described in (Lagrée et al. 2011) and illustrated in Figure 6.5. We assume the flow to be slow enough to neglect inertial terms.

The conservation of momentum on the longitudinal  $x$ -axis and perpendicular  $y$ -axis reads

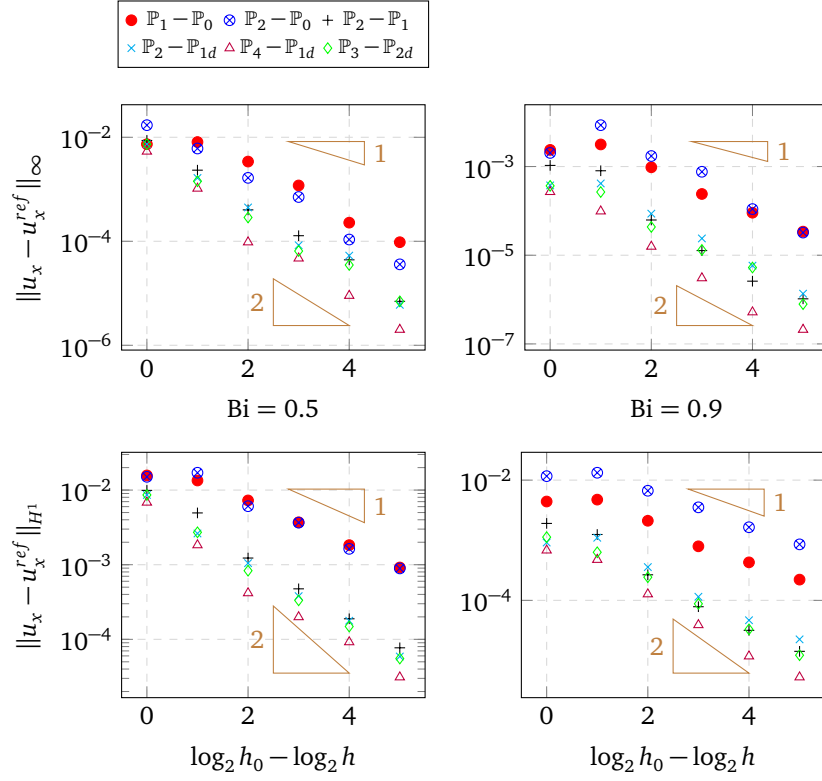
$$\begin{aligned} \frac{\partial}{\partial y} \left( \eta \frac{\partial \mathbf{u}_x}{\partial y} + \tau_{xy} \right) &= -\sin \alpha \\ \frac{\partial p}{\partial y} &= -\cos \alpha. \end{aligned}$$

Integrating the second equation with the condition that the pressure should be zero at the top of the granular layer, i.e.,  $p(H) = 0$ , gives  $p = (1 - y) \cos \alpha$ .

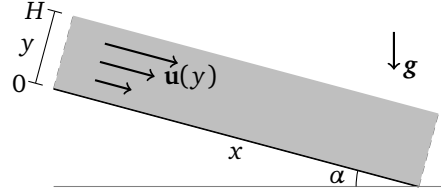
We consider the case of an avalanching flow with  $\frac{\partial \mathbf{u}_x}{\partial y} > 0$  for  $y < 1$ . This means that the friction is saturated, therefore  $\tau_{xy} = \mu(1 - y) \cos(\alpha)$ . We get

$$\frac{\partial}{\partial y} \eta \frac{\partial \mathbf{u}_x}{\partial y} = -\sin \alpha + \mu \cos \alpha.$$

A Neumann boundary condition at the interface  $\frac{\partial \mathbf{u}_x}{\partial y}(1) = 0$  imposes  $\frac{\partial}{\partial y} \eta \frac{\partial \mathbf{u}_x}{\partial y} < 0$ , which means  $\mu < \tan \alpha$ .



**Figure 6.4:**  $L_\infty$ -norm (top) and  $H^1$ -norm (bottom) convergence plots for various FEM approximation orders, for  $Bi = 0.5$  (left) and  $Bi = 0.9$  (right).  $u_x^{ref}$  denotes the analytic solution.

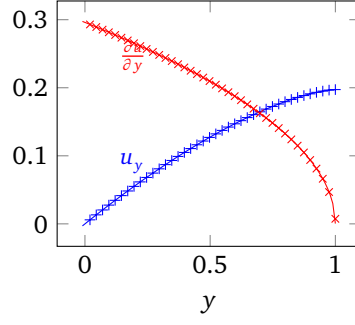


**Figure 6.5:** Flow on an infinite inclined plane

For spatially constant (Newtonian)  $\eta$  and  $\mu$ , we get  $\frac{\partial \mathbf{u}_x}{\partial y} = \frac{1-y}{\eta} \sin \alpha - (\mu \cos \alpha)$ , and  $\mathbf{u}_x(y)$  is quadratic. In order to retrieve the typical  $\frac{3}{2}$  power of the Bagnold profile, we can instead choose  $\eta(y) := |\mathbf{D}(\mathbf{u})| = \frac{1}{2} \frac{\partial \mathbf{u}_x}{\partial y}$ , which gives the analytic expressions

$$\begin{aligned} \frac{\partial \mathbf{u}_x}{\partial y} &= \sqrt{2 \sin \alpha - (\mu \cos \alpha)} (1-y)^{\frac{1}{2}} \\ \mathbf{u}_x(y) &= \frac{2}{3} \sqrt{2 \sin \alpha - (\mu \cos \alpha)} \left( 1 - (1-y)^{\frac{3}{2}} \right). \end{aligned}$$

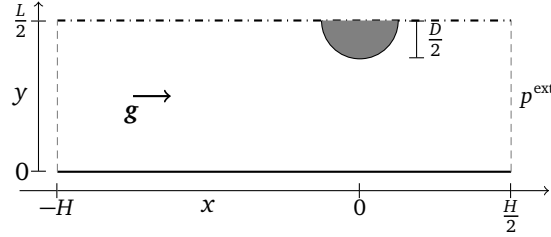
We simulated this model on a square patch using the algorithm presented in Section 6.2. The value of  $\sigma^{\text{ext}}$  on the upstream and downstream boundaries was computed using the analytic solution. In order to handle the non-constant  $\eta$ , we used a fixed-point algorithm which happened to converge very fast in practice — a dozen or so of iterations.



**Figure 6.6:** Comparison between our numerical solution (marks) and the analytic one (lines), for the velocity and shear-rate profiles of a Bagnold avalanche flow with  $\mu = 0.5$  and  $\alpha = \tan^{-1}(1.1\mu)$ .

Numerical and analytic profiles are compared in Figure 6.6. Once again, we observe a very good agreement between our simulations and the analytic solution.

### 6.5.2 Flow around a cylinder



**Figure 6.7:** Flow around a cylinder of diameter  $D$  in 2D channel of width  $L$ . The speed of the flow can be adjusted by varying the external pressure at the downstream boundary  $p^{\text{ext}}$ .

In this section, we study a gravity induced flow in a 2D channel of width  $L$  around a cylinder of diameter  $D$ , as shown in Figure 6.7. We use no-slip boundary conditions for the velocity on the sides of the channel and on the cylinder, and homogeneous Neumann conditions at the upstream and downstream boundaries.

**Equivalent drag** We attempt to reproduce the experimental setting presented in (Chehata et al. 2003). Using a regularization of the  $\mu(I)$  rheology, Chauchat and Médale (2014) performed simulations that showed a good qualitative agreement with the experimental results. However, their approach suffered from two drawbacks:

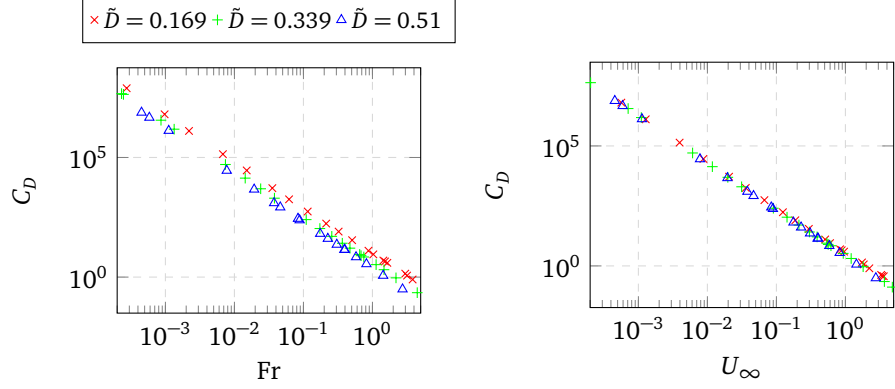
- The regularization induces a creeping flow even when none should occur;
- The  $\nabla \cdot \mathbf{u} = 0$  condition leads to negative pressures behind the obstacle. This is outside the domain of validity of the  $\mu(I)$  rheology.

Our method also suffer from approximations; we assume a constant volume fraction of grains, even if this should clearly not be the case behind the obstacle. As the flow is slow, we also neglect its inertia. We use a constant friction coefficient,  $\mu = 0.5$ , and set  $\text{Re} = 100$  and  $\text{Bi} = 0$ . We vary the average upstream velocity  $U_\infty$  by setting different values for the external pressure  $p^{\text{ext}}$  at the downstream interface; the lower this pressure, the faster the grains will flow. This emulates the varying outlet size in (Chehata et al. 2003). We set no external pressure at the upstream boundary.

The grains exercise a longitudinal drag force on the cylinder  $O$  that can be computed as  $F_D := \int_{\partial O} \boldsymbol{\lambda} \mathbf{n}_O \cdot \mathbf{e}_x$ . The Froude number is defined from the dimensionless average upstream velocity  $U_\infty$  and the length ratio  $\tilde{D} := \frac{D}{L}$  as  $\text{Fr} := \frac{U_\infty}{\sqrt{\tilde{D}}}$ .

The dimensionless equivalent drag coefficient  $C_D$  is then deduced as

$$C_D := \frac{F_D}{\frac{1}{2} U_\infty^2 \tilde{D}}.$$



**Figure 6.8:** Drag coefficient  $C_D$  versus Froude number  $\text{Fr}$  (left) and upstream velocity  $U_\infty$  (right) for the granular flow around a cylinder

Figure 6.8 shows on a logarithmic scale the evolution of the drag coefficient  $C_D$  with the Froude number (left) or average velocity (right), for different cylinder diameters. We retrieve the linear profiles observed in (Chauchat and Médale 2014; Chehata et al. 2003). The fact that the data points for the different diameters become aligned on the right plot means that the slope of  $\frac{F_D}{\frac{1}{2} \tilde{D}}$  is independent of  $\tilde{D}$ ; the drag force depends linearly on the diameter of the obstacle.

**Visualization of velocity and stress fields** We now consider a narrow channel ( $L = 4D$ ) with a free-slip boundary condition on the side walls — results for no-slip walls are depicted in (Daviet and Bertails-Descoubes 2016b, Figure 14). Figure 6.9 collects plots of the velocity and stress fields across the domain. The pressure field is of special interest as it possesses two notable features:

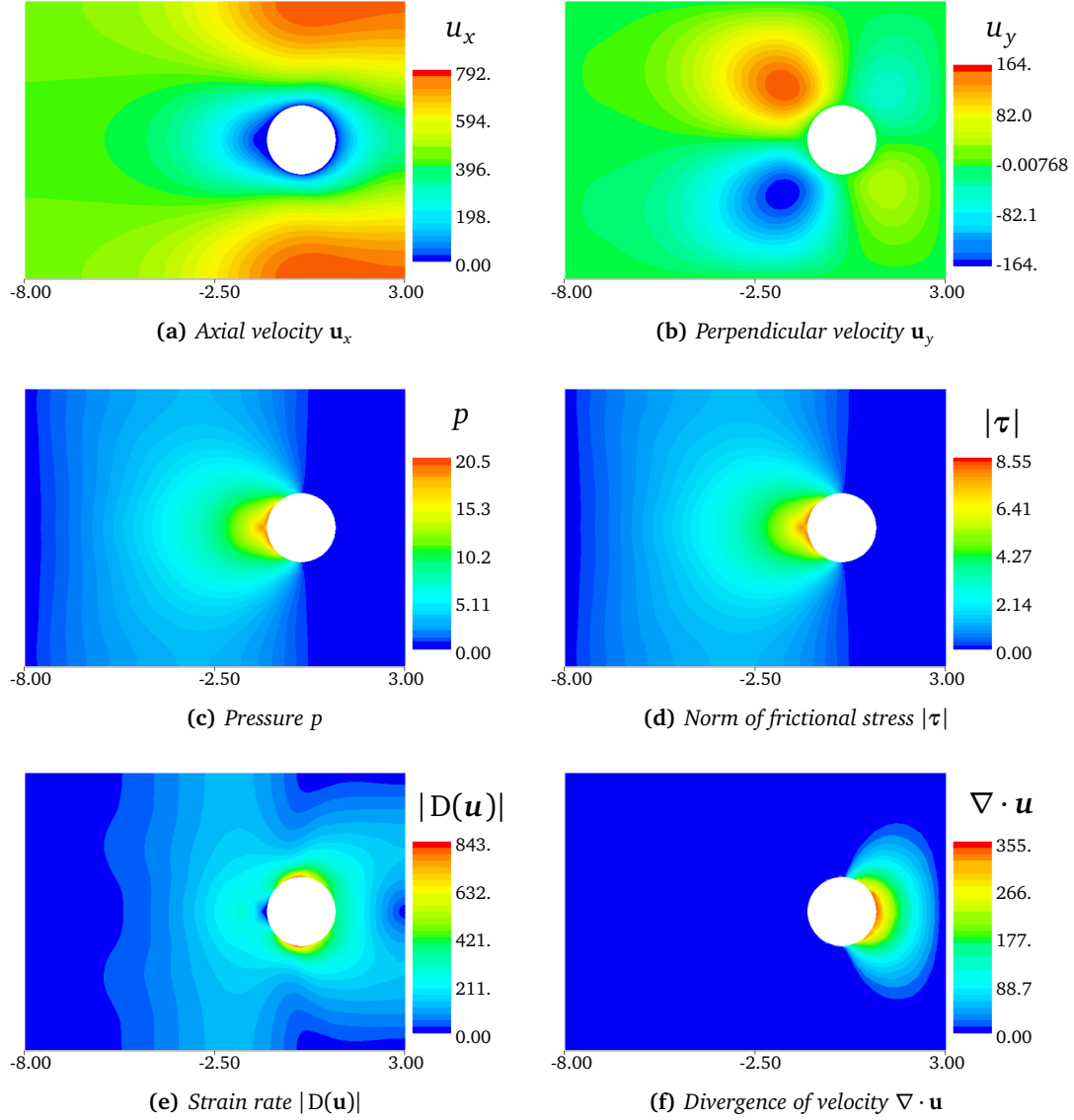
- First, it validates the benefit of allowing the dilation of the flow: the pressure does vanish in the wake of the obstacle, and the pressure intensity corresponds qualitatively with the experimental results of Seguin et al. (2016) using photoelastic grains;
- The zone of highest pressure is not located at the very front of the obstacle. Instead, we observe the formation of high-pressure arch above this point;

The shear rate plot in Figure 6.9(e) highlights the existence of a triangular rigid zone in front of the obstacle.

### 6.5.3 Extension to inertial flows: discharge of a silo

**Beverloo scaling** One of the most widely accepted macroscopic feature of the granular flow in a silo is the so-called Beverloo scaling (Beverloo et al. 1961), stating that the discharge rate  $Q$  depends on the diameter of the outlet to the power  $d - \frac{1}{2}$ ,

$$Q = C \sqrt{g} (L - kD_g)^{d - \frac{1}{2}},$$



**Figure 6.9:** Velocity and stress fields for the flow around a cylinder in a narrow channel without wall friction, with  $\mu = 0.3$ ,  $Re = 100$ ,  $Bi = 0$  and  $\tilde{D} = \frac{1}{4}$

where  $D_g$  is the diameter of a grain, and  $C$  and  $k$  are dimensionless constants depending on the silo geometry and granular properties. The number  $k$  typically lies within the range  $1 < k < 3$ .

The Beverloo phenomenon is particularly relevant for us as it has been shown that such a scaling cannot be recovered for Newtonian flows (Staron et al. 2012), nor for flows with a yield stress that does not depend on the pressure. The physical justification of the scaling involves inertia (Mankoc et al. 2007), so it is hopeless to attempt to retrieve it solely with the formulation of Section 6.2. For this problem, we therefore need to modify our equations to account for the grains inertia.

**Modification of the variational formulation** We now add an inertial term to our momentum conservation equation (6.10),

$$\rho \left( \frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} \right) - \nabla \cdot [2\eta E^b \dot{\epsilon} - \lambda] = \rho g \mathbf{e}_g,$$

which can be made dimensionless by defining  $\tilde{t}$  such that  $t = \frac{L}{U} \tilde{t} = \sqrt{\frac{L}{g}} \tilde{t}$ ,

$$\left( \frac{\partial \tilde{\mathbf{u}}}{\partial \tilde{t}} + (\tilde{\mathbf{u}} \cdot \tilde{\nabla}) \tilde{\mathbf{u}} \right) - \tilde{\nabla} \cdot \left[ \frac{2}{\text{Re}} E^b \tilde{\epsilon} - \tilde{\lambda} \right] = \mathbf{e}_g.$$

As in the previous section, we will drop the tildes from now on.

Putting directly the transport term into the variational formulation would lead to the introduction of a trilinear form, which upon discretization would ultimately yield an asymmetric matrix  $A$ . Instead, we will use a timestepping algorithm and the first-order *characteristics* method presented in Section 5.2.2. The total derivative is discretized as in (5.15),

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} = \frac{\mathbf{u}^{k+1} - \mathbf{u}^k \circ X^k}{\Delta_t} + O(\Delta_t), \quad X^k(\mathbf{x}) := \mathbf{x} - \Delta_t \mathbf{u}^k(\mathbf{x}, t^k).$$

The variational formulation for each timestep is then obtained by taking the forms defined in Section 6.2.2 and adding the following terms to  $l$  and  $a$ ,

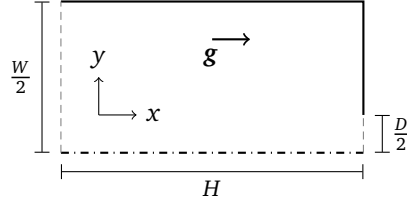
$$\begin{aligned} a(\mathbf{u}, \mathbf{v}) &:= a_{(6.2.2)}(\mathbf{u}, \mathbf{v}) + \frac{1}{\Delta_t} \int_{\Omega} \langle \mathbf{u}, \mathbf{v} \rangle \\ l(\mathbf{v}) &:= l_{(6.2.2)}(\mathbf{v}) + \frac{1}{\Delta_t} \int_{\Omega} \langle \mathbf{u}^k \circ X^k, \mathbf{v} \rangle. \end{aligned}$$

Solving for each timestep in the inertial setting is therefore equivalent to solving the problem defined in Section 6.2.

**Remark 6.3.** *Note that the timestepping scheme ensures a positive-definite form “ $a$ ” even when  $\text{Re} = +\infty$ . This allows us to extend the existence results of Property 6.4 to purely plastic flows, instead of being restricted to viscoplastic ones as in Section 6.2.*

**Results** We simulated the 2D silo shown in Figure 6.10, with  $W = H = 8L$  and an aperture size  $D$ . We made the ratio  $\tilde{D} := \frac{D}{L}$  vary from  $\frac{1}{2}$  to 2, and studied the change in the dimensionless discharge rate  $\tilde{Q}$  defined such that  $Q := \sqrt{g} L^{\frac{3}{2}} \tilde{Q}$ .

Plots for different rheologies, with the corresponding Beverloo fits when one was found, are shown in Figure 6.11. Beverloo law coefficients for  $\mu = 0.6$ , for which the fit was always acceptable except for  $\text{Re} = 1$ , are given in Table 6.1. Coefficients for  $\mu = 0.4$  are also given when the fit was deemed of sufficient quality. We remark that decreasing  $\text{Re}$ , increasing  $\mu$  or increasing  $\text{Bi}$  all contribute to a consequent increase of  $k$ .  $\mu$  and  $\text{Bi}$  have also a positive impact on  $C$ , while  $\text{Re}$  has only a small influence on this parameter.



**Figure 6.10:** Geometry of a half-silo of height  $H$ , width  $W$  and aperture size  $D$

		Re = 100		Re = 1000	Re = 10
		Bi = 0	Bi = 0.1	Bi = 0	
$\mu$	C	1.70	1.44	1.77	2.00
	k	0	0.96	0	3.34
$\mu = 0.4$	C	1.54	1.52	1.54	1.48
	k	0.87	2.49	0.38	3.07
$\mu = 0.5$	C	1.35	1.43	1.35	1.33
	k	0.83	3.20	0.22	3.56

**Table 6.1:** Beverloo law coefficients obtained for  $\mu = 0.6$  and  $\mu = 0.4$ , for different  $Bi$  and  $Re$  values. The values for  $k$  are given assuming  $L = 11.2D_g$ .

**Extension to the  $\mu(I)$  rheology** The  $\mu(I)$  (see Section 0.3.3) rheology is classically integrated into dynamics solvers by explicitly evaluating the value of the equivalent viscosity  $\eta_{\text{eff}}$  at each timestep (Lagrée et al. 2011).

In our case, we only have to explicitly evaluate the friction coefficient  $\mu(I)$  instead of  $\eta_{\text{eff}}$ . Since  $\mu(I)$  can only take values in  $[\mu_S, \mu_D]$  whereas  $\eta_{\text{eff}}$  takes values in  $\mathbb{R}^+$ , our approach has a few benefits:

- We do not have to clamp the value of  $\eta_{\text{eff}}$ , and can have fully rigid zones where the shear rate is strictly zero (which means an infinite  $\eta_{\text{eff}}$ );
- The loss of stability of the time-integration scheme due to this explicit term is much less dramatic;
- While fully implicit approaches are possible (Ionescu et al. 2015), our explicit update rule remains much simpler — and cheaper.

While the simulation frameworks are quite different — we do not take into account the air phase and use a non-zero ( $Re = 100$ ) Newtonian viscosity, we nevertheless attempted to recreate a simulation from (Staron et al. 2014, Fig 4), using the same physical parameters ( $D = 11.2D_g$ ,  $I_0 = 0.4$ ,  $\mu_S = 0.32$ ,  $\mu_D = 0.6$ ).

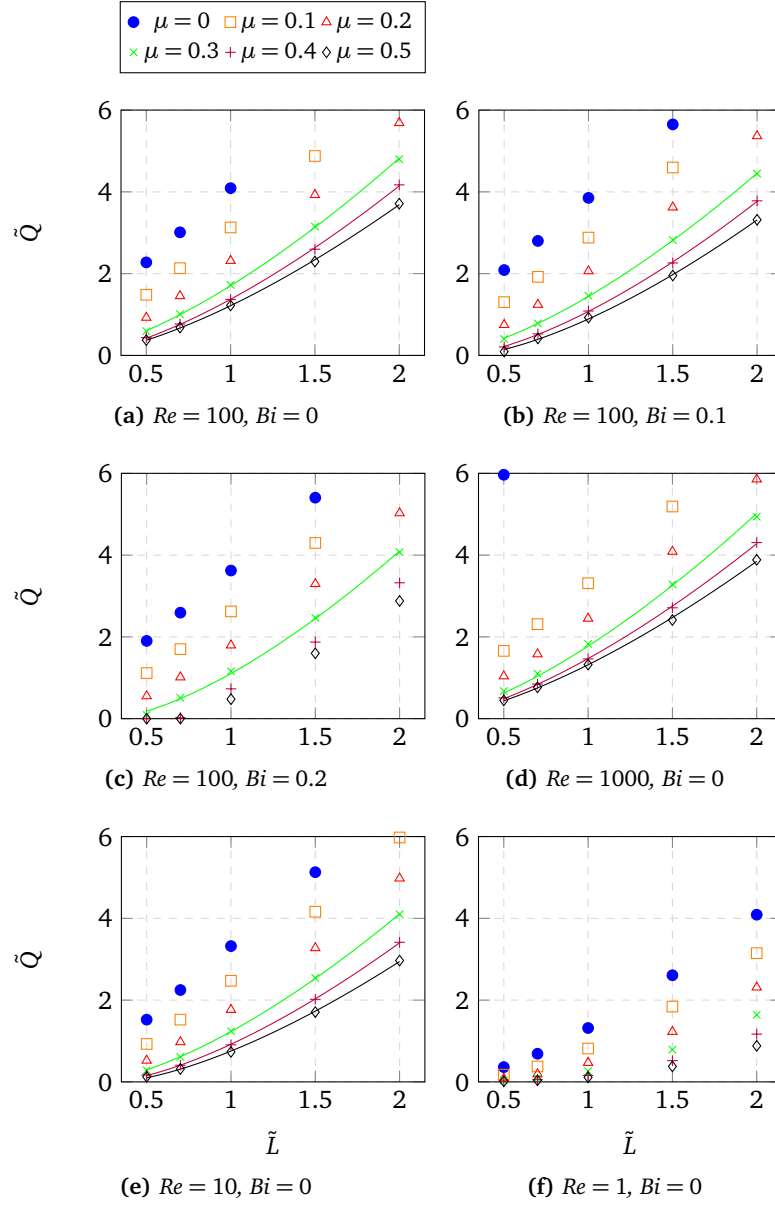
The results are shown in Figure 6.12 ; the match on the  $C$  coefficient of the Beverloo law is surprisingly good (both methods give  $C = 1.48$ ), however we retrieve a coefficient  $k$  that is significantly smaller (0.52 vs 0.73). Velocity profiles along the vertical and horizontal sections described in (Staron et al. 2014) are also presented in Figure 6.12.

#### 6.5.4 Performance

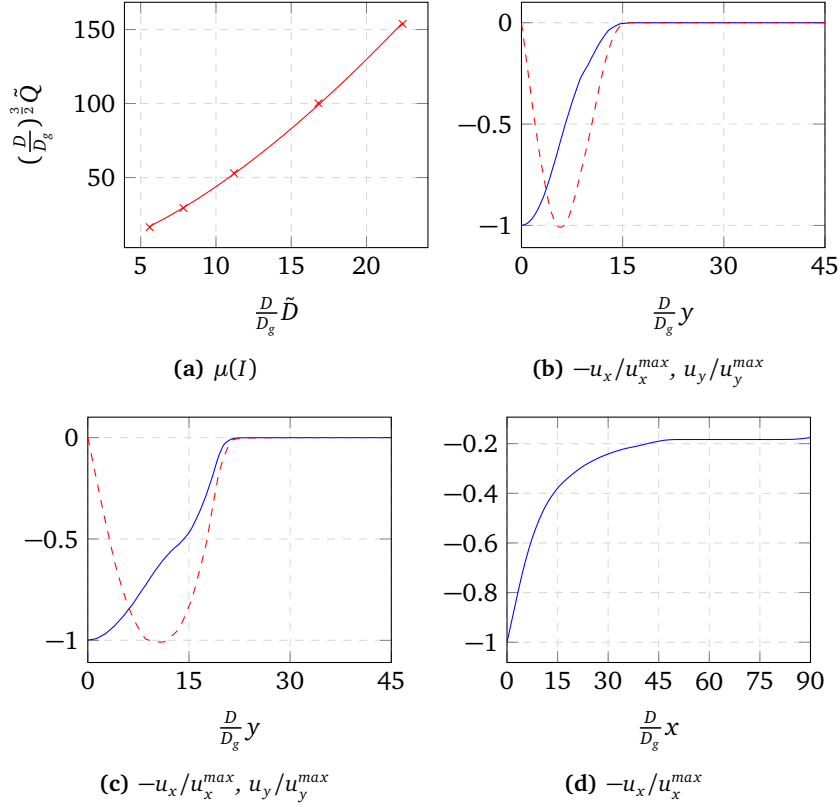
We now study the computational performance for the problem of Section 6.2 – the dynamics case is just a sequence of such problems. In the following, **PG** will denote the canonical Projected-Gradient algorithm, **APGD** the Accelerated Projected Gradient Descent from Heyn (2013) — which is roughly similar to the FISTA\* algorithm advocated in (Treskatis et al. 2016), **ASPG** the Accelerated Spectral Projected Gradient (Algorithm B.2), and **IP** the Interior-Point algorithm from the commercial package MOSEK (E. Andersen, Roos, et al. 2003).

In order to make as fair a comparison as possible, to evaluate the current error at each fixed-point iteration we have chosen not to use a residual based on a projection operator — which





**Figure 6.11:** Dimensionless discharge rate  $\tilde{Q}$  versus silo outlet diameter  $\tilde{L}$  for and different rheologies. When  $\mu$  is small, we cannot get a reasonable fit; in the case of the Newtonian flow, the discharge rate is closer to a linear law.  $Bi$  mainly influences the  $k$  parameter of the law but does suffices to obtain a fit. Large values of  $Re$  tend to linearize the discharge rate for small  $\mu$ , smaller values make it quadratic.



**Figure 6.12:** Top-left: Discharge rate as a function of the outlet size for the  $\mu(I)$  rheology. The coefficients of the Beverloo fit are, using the dimensionless unit of (Staron et al. 2014),  $C = 1.47$  and  $k = 0.063$ . Top-right and bottom: Horizontal and vertical velocity profiles for  $\tilde{D} = 1$  along section S1, S2 and H1, as defined in (Staron et al. 2014, Fig. 6)

Solver	Mesh 1			Mesh 2			Mesh 3			Mesh 4		
	$t^1$	$n^2$	$t/n$	$t$	$n$	$t/n$	$t$	$n$	$t/n$	$t$	$n$	$t/n$
IP	71	100	0.71	116	29	4	564	23	24	6281	17	369
ASPG	13	88	0.14	69	79	0.88	144	33	4.4	1098	47	23

<sup>1</sup>  $t$  : Total wall time in seconds.

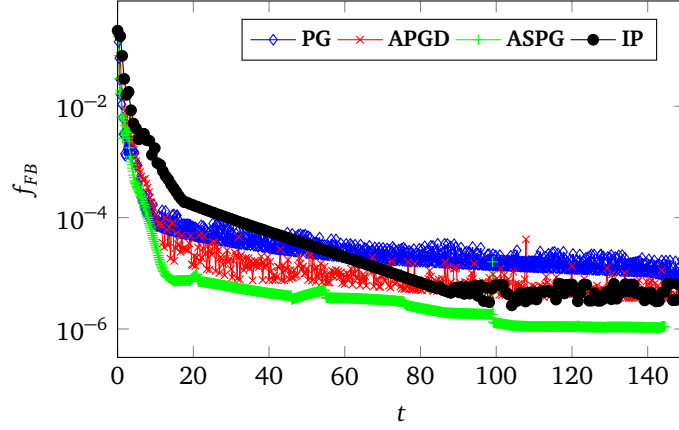
<sup>2</sup>  $n$  : Number of iterations of the fixed-point algorithm.

**Table 6.2:** Comparison of time taken to reach a give tolerance for 4 meshes with respectively 670, 2703, 10403 and 41509 vertices.

are known to be disadvantageous for interior-points — but rather the SOC Fischer-Burmeister function for frictional problems as defined e.g. in Chapter 4.

We consider the silo geometry of Figure 6.10, a  $\mathbb{P}_2 - \mathbb{P}_{1d}$  FEM approximation order, and a rheology with  $\mu = 0.5$ ,  $\text{Re} = 1$  and  $\text{Bi} = 0$ . The base mesh has a resolution of 640 vertices.

Figure 6.13 shows typical convergence plots for selected algorithms. Table 6.2 shows the evolution of the computation time in seconds as the mesh is refined, for an interior-point method and a Spectral Projected Gradient method. The latter, despite its simplicity, is therefore highly competitive.



**Figure 6.13:** *Infinity-norm of the Fischer-Burmeister function as a function of the wall time for the Projected-Gradient, Accelerated Projected Gradient Descent, Spectral Projected Gradient and Interior-Point algorithm*

## Discussion

In this chapter we have presented a macroscopic, continuous model for granular-like flows, and we have shown that it was driven by equations structurally similar to those appearing in contact mechanics between discrete elements. We have then exhibited algorithms for solving those equations on a class of spatial discretization and demonstrated that they were well-defined.

These new numerical simulation tools allowed us to capture the typical nonsmooth effects critical to the simulation of yield-stress flows, as well as to accommodate complex rheologies such as the well-known  $\mu(I)$  rheology. The main advantages over most previous approaches are the absence of any regularization or clamping, and improved efficiency over Augmented Lagrangian algorithms, thanks to formulations in the framework of Second Order Cone Programming.

**Limitations** The computational cost of our approach remains a major limitation. The SOCP formulation becomes quickly expensive when the number of required fixed-point iterations increases, while the Projected Gradient requires solving at least one linear system at each inner iteration. This makes three-dimensional scenarios only tractable for rough meshes. However, note that for the timestepping algorithm without Newtonian viscosity ( $\text{Re} = +\infty$ ), the inverse of the  $A$  matrix will be sparse for certain discretizations of  $V(\Omega)$ . The explicit computation of the Delassus operator would then become possible. In the next chapter, we will propose an approximation for large  $\text{Re}$  that makes use of this fact. This will both speed up our numerical method and allow us to use other DCFP solvers, such as our hybrid Gauss–Seidel algorithm from Chapter 4.

From a physics point of view, the fact that we do not take into account a variable density leads to invalid predictions in zones where the local divergence is strictly positive, such as regions in the wake of obstacles — or even any shearing region in the case of a non-zero dilatancy coefficient. Our simple unilateral compressibility condition  $\nabla \cdot \mathbf{u} \geq 0$  also prevents the material from recompacting after dilatation; a condition based on the current volume fraction, such as proposed by Dunatunga and Kamrin (2015), should be used instead. The next chapter will be dedicated to the treatment of a proper maximal volume fraction constraint, thus extending our approach to general scenarios. Our treatment of boundary conditions is also limited. As argued in (Domnik and Pudasaini 2012), we should consider more realistic laws, such as Coulombic ones.

Finally, as usual when dealing with Coulomb-like friction models, the theoretical existence and uniqueness criterions for our solutions are rather weak. Not being able to uniquely de-

termine the stress field for a given velocity solution is problematic for assessing the stress on structures. However, as the flowing and dilating regions enforce boundary conditions for the stress inside the rigid zones, we did not find this under-determination to be problematic for our test scenarios.



## 7 Dry granular flows

In this chapter, we relax the dense flow hypothesis from the framework presented in Chapter 6. As such, we shall be able to simulate arbitrary compacting and dilating granular flows. Note that a large part of the contents of this chapter has been published in (Daviet and Bertails-Descoubes 2016a).

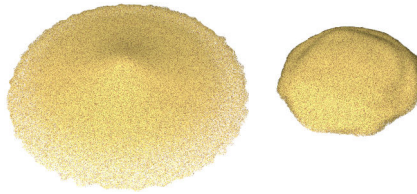
Just like in (Narain, Golas, et al. 2010) and (Dunatunga and Kamrin 2015), we shall assume that the volume fraction of grains  $\phi$  fully determines the nature of the contacts between the grains. That is, we assume the existence of a critical value  $0 < \phi_{\max} \leq 1$  such that for  $\phi < \phi_{\max}$ , we consider that the grains are mostly separated, and that they are interacting through sparse random impacts; the material is in the gaseous regime. However, for  $\phi = \phi_{\max}$ , we shall consider that the grains are in permanent contact, and that their interactions are driven by dry Coulomb-like friction. The material is then either in the liquid or the solid regime, and subject to the non-associated Drucker–Prager flow rule  $\mathcal{DP}(\mu, \sigma_s, \tau_c, \zeta)$ . In practice, we used  $\phi_{\max} = 0.6$ , corresponding to a loose 3D monodisperse packing.

As in the previous chapter, we shall also assume that the material is inelastic (the grains are fully rigid, and thus the volume fraction  $\phi$  should not exceed the critical value  $\phi_{\max}$ ), and that the interactions between the grains and the surrounding fluid can be neglected. Moreover, for the sake of simplicity we shall no longer consider non-zero values for the shear yield stress  $\sigma_s$  and the dilatancy  $\zeta$ . Following the considerations of Section 6.4, the adaptation of the numerical solvers to handle those cases will be straightforward.

### 7.1 Spatially continuous model

#### 7.1.1 Constitutive equations

As we neglect the effect of the surrounding air, the total stress of the grain–air mixture with volume fraction of grains  $\phi$  can be expressed as  $\boldsymbol{\sigma} = \phi \boldsymbol{\sigma}^g$ , where  $\boldsymbol{\sigma}^g$  is the granular phase stress. As in the previous chapter, we decompose the solid phase stress tensor as  $\boldsymbol{\sigma}^g := \eta \dot{\boldsymbol{\epsilon}} + \boldsymbol{\sigma}^C$ , where the first part corresponds to a standard Newtonian viscosity (dissipative term due to random collisions in the flowing material), and  $\boldsymbol{\sigma}^C$  is the additional stress due to the Coulombic interactions between individual grains. Note that  $\eta$  may be chosen equal to zero, in this case no internal stress is applied in the “gaseous” phase and we retrieve the constitutive law used by Narain, Golas, et al. (2010). The normal contact stress (a.k.a. “pressure”) is  $p = -\frac{1}{d} \text{Tr } \boldsymbol{\sigma}^C$ , such that  $\boldsymbol{\sigma}^C = \text{Dev } \boldsymbol{\sigma}^C - p \mathbb{I}$ .



**Figure 7.1:** Final states after the collapse of a cylindrical column with exact (Frobenius norm, left) and linearized ( $\ell_\infty$ -norm, right) Drucker-Prager yield surfaces, both for  $\mu = 0.6$ .

Just like in the previous chapter, we consider a Drucker–Prager yield criterion with friction coefficient  $\hat{\mu}$  and cohesion  $c$ , that is, a pressure-dependent yield stress  $\hat{\mu}(p + c)$ . However, following Narain, Golas, et al. (2010), the pressure no longer enforces a positive divergence of the flow, but the maximal volume fraction constraint  $\phi \leq \phi_{\max}$ . That is, we replace Equation (6.3) with the complementarity relationship

$$0 \leq p + c \perp \phi_{\max} - \phi \geq 0. \quad (7.1)$$

The maximum dissipation principle states that in the yielded regime, friction should be saturated and the frictional stress tensor should be colinear to the deviatoric part of the strain rate. The deviatoric part of  $\sigma^C$  should thus satisfy one of the two regimes,

$$\begin{cases} \text{Dev } \sigma^C = \hat{\mu}(p + c) \frac{\text{Dev } \dot{\epsilon}}{|\text{Dev } \dot{\epsilon}|} & \text{if } \text{Dev } \dot{\epsilon} \neq \mathbf{0} \text{ (yielded)} \\ |\text{Dev } \sigma^C| \leq \hat{\mu}(p + c) & \text{if } \text{Dev } \dot{\epsilon} = \mathbf{0} \text{ (unyielded)}. \end{cases} \quad (7.2)$$

Note that in contrast with (Narain, Golas, et al. 2010), we will not linearize the  $|\cdot|$  norm in Equation (7.2), thus avoiding the disturbing geometrical artifacts inherent to their method (Figure 7.1).

**Remark 7.1.** When  $\phi < \phi_{\max}$ , Equation (7.1) imposes  $p + c = 0$ , and thus Equation (7.2) dictates that the frictional stress should vanish, which is expected. However, the mixture pressure will be  $-\phi c$  in the gaseous zones; if  $c > 0$ , the grains will thus tend to attract each other even though the contacts are broken, which is unwanted. For this reason, the cohesion field  $c$  should also satisfy  $\phi < \phi_{\max} \implies c = 0$ . This may create spatial discretization issues, as we will see later.

### 7.1.2 Energy considerations

Let us study the conditions for which the flow rule defined by Equation (7.1–7.2) and the conservation equations (5.7, 5.10) lead to a dissipative system, as is expected from a granular material. We have already performed a similar study in the supplemental document to (Daviet and Bertails-Descoubes 2016a) for the discrete-time version of the rheology. We propose here to look at the continuous-time setting instead, by leveraging an analogy with impacts in discrete contact mechanics.

**Local dissipativity of contact stress** Let us consider a material point  $\mathbf{x}(t)$  of the domain, with  $\phi(t) := \phi(\mathbf{x}(t), t)$ . We will see that both the inelastic impacts hypothesis and the criterion from Remark 7.1, that cohesion should vanish in non-dense zones, are critical to ensure that the work of the contact stress is always and everywhere dissipative.

The density of energy dissipation by the contact stress is given by  $(-\phi(t)\sigma^C(t)) : \dot{\epsilon}(t)$ . For the system to be locally dissipative, it suffices that  $\phi\sigma^C : \dot{\epsilon} \geq 0$ , or, *a fortiori*, that

$$\begin{cases} \text{Dev } \sigma^C : \text{Dev}(\phi D(\mathbf{u})) \geq 0 & (7.3a) \\ \text{Tr } \sigma^C \text{Tr}(\phi D(\mathbf{u})) \geq 0. & (7.3b) \end{cases}$$

The implication (7.1–7.2)  $\implies$  (7.3a) is trivial; either  $\text{Dev}(\phi D(\mathbf{u})) = 0$ , or  $\text{Dev}(\phi D(\mathbf{u})) \neq 0$  and  $\text{Dev } \sigma^C = \hat{\mu}(p + c) \frac{\text{Dev } \dot{\gamma}}{|\text{Dev } \dot{\gamma}|}$  with  $p + c \geq 0$ .

Now, (7.3b) can be written equivalently as  $(p\phi \nabla \cdot \mathbf{u}) \leq 0$ , or again, using the mass conservation equation (6.1),  $p \frac{D\phi}{Dt} \geq 0$ . Using an analogy with discrete mechanics, we argue below that this inequality shall hold if the consideration made in Remark 7.1 is satisfied.

The maximal volume fraction constraint,  $\phi_{\max} - \phi(t) \geq 0$ , can be seen as analogous to the normal part of the *gap* function in contact mechanics. The instant at which  $\phi$  reaches  $\phi_{\max}$  constitutes an impact, and the time-derivative of the gap function (the normal relative velocity

for discrete contacts,  $\frac{D\phi}{Dt}(\mathbf{x})$  here) will be discontinuous. However, this derivative will retain both a left and a right limit at every instant.

Pursuing the analogy with nonsmooth discrete mechanics, we interpret our inelastic impacts hypothesis as preventing rebounds; that is, if  $\phi$  reach  $\phi_{\max}$  at time  $t$ , then we shall have  $\frac{D\phi}{Dt}(\mathbf{x}(t), t^+) = 0$ . Moreover, just like we did to construct the Signorini conditions, we assume that shocks do not propagate. Therefore, if  $t$  is not a time of impact for the material point  $\mathbf{x}(t)$ , then the contact pressure  $p(t) := p(\mathbf{x}(t), t)$ , should possess both left and right limits.

The volume fraction and its derivative have thus to satisfy one of the following four cases,

$$\left\{ \begin{array}{l} \phi(t) < \phi_{\max} \\ \text{or } \phi(t) = \phi_{\max} \text{ and } \frac{D\phi}{Dt}(\mathbf{x}(t), t^-) > 0 \text{ and } \frac{D\phi}{Dt}(\mathbf{x}(t), t^+) = 0 \\ \text{or } \phi(t) = \phi_{\max} \text{ and } \frac{D\phi}{Dt}(\mathbf{x}(t), t) = 0 \\ \text{or } \phi(t) = \phi_{\max} \text{ and } \frac{D\phi}{Dt}(\mathbf{x}(t), t^-) > 0 \text{ and } \frac{D\phi}{Dt}(\mathbf{x}(t), t^+) = 0. \end{array} \right. \quad \begin{array}{l} (7.4a) \\ (7.4b) \\ (7.4c) \\ (7.4d) \end{array}$$

Case (7.4a) corresponds to the gaseous regime, case (7.4b) to an impact, case (7.4c) to persistent contact and case (7.4d) to the onset of dilation.

Let us show that, under the hypothesis of Remark 7.1, i.e., that  $\phi < \phi_{\max} \implies c = 0$ , then  $p \frac{D\phi}{Dt} \geq 0$ . This is obvious when  $\frac{D\phi}{Dt} = 0$ , let us consider the remaining cases.

- Case (7.4a). The complementarity condition (7.1) requires  $p = -c = 0$ .
- Case (7.4b) at  $t^-$ . Left-continuity of  $\frac{D\phi}{Dt}$  means that there exists  $t_0$  such that  $\forall s \in [t_0, t[$ ,  $\frac{D\phi}{Dt}(s) < 0$ , and thus  $\phi(s) < \phi_{\max}$ , meaning  $c(s) = 0$ . Left-continuity of  $c$  imposes in turn that  $c(t^-) = 0$ , and condition (7.1) that  $p(t^-) \geq 0$ . We deduce  $(p \frac{D\phi}{Dt})(\mathbf{x}, t^-) \geq 0$ .
- Case (7.4d) at  $t^+$ . As  $t$  is not a time of impact, by analogy with discrete mechanics we assume right-continuity of  $p$ . There exists  $t_0$  such that  $\forall s \in ]t, t_0]$ ,  $\phi(s) < \phi_{\max}$ , thus  $c(s) = p(s) = 0$ , and  $p(t^+) = 0$ . Once again  $(p \frac{D\phi}{Dt})(\mathbf{x}(t), t^+) = 0$ .

**Global dissipativity** For the sake of simplicity, we shall consider a domain  $\Omega$  with homogeneous Dirichlet boundary conditions.

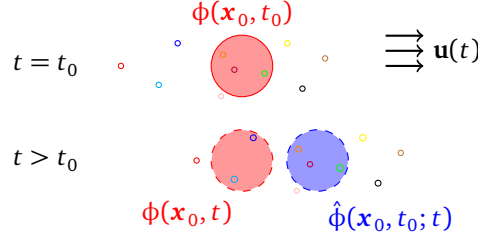
We can decompose the total energy of the system as the sum of its potential and kinetic energies, i.e.,  $\mathcal{E}^t = \mathcal{E}^p + \mathcal{E}^c$ , with

$$\mathcal{E}^p := - \int_{\Omega} \rho \phi \langle \mathbf{x}, \mathbf{g} \rangle \quad \mathcal{E}^c := \int_{\Omega} \frac{1}{2} \rho \phi \langle \mathbf{u}, \mathbf{u} \rangle.$$

Since we are using homogeneous Dirichlet boundary conditions, their evolution follows

$$\begin{aligned} \frac{d\mathcal{E}^p}{dt} &= -\rho \int_{\Omega} \frac{\partial \phi}{\partial t} \langle \mathbf{x}, \mathbf{g} \rangle = \rho \int_{\Omega} \nabla \cdot [\phi \mathbf{u}] \langle \mathbf{x}, \mathbf{g} \rangle \\ &= -\rho \int_{\Omega} \phi \langle \mathbf{u}, \mathbf{g} \rangle + \rho \underbrace{\int_{\text{Bd } \Omega} \phi \langle \mathbf{u}, \mathbf{n}_{\Omega} \rangle \langle \mathbf{x}, \mathbf{g} \rangle}_{=0}, \end{aligned}$$





**Figure 7.2:** Eulerian (in red) and Lagrangian (in blue) points of view for studying the change in volume fraction of a compressible fluid.

and

$$\begin{aligned}
 \frac{d\mathcal{E}^c}{dt} &= \int_{\Omega} \rho \phi \left\langle \frac{D\mathbf{u}}{Dt}, \mathbf{u} \right\rangle \\
 &= \int_{\Omega} \eta \langle \nabla \cdot [\phi \dot{\epsilon}], \mathbf{u} \rangle + \langle \nabla \cdot [\phi \sigma^C], \mathbf{u} \rangle + \rho \phi \langle \mathbf{g}, \mathbf{u} \rangle \\
 &= \rho \int_{\Omega} \phi \langle \mathbf{u}, \mathbf{g} \rangle - \underbrace{\eta \int_{\Omega} \phi \dot{\epsilon} : \dot{\epsilon}}_{\geq 0} - \int_{\Omega} \phi \sigma^C : \dot{\epsilon} \\
 &\quad + \underbrace{\int_{\text{Bd}\Omega} \phi \langle \mathbf{u}, (\eta \dot{\epsilon} + \sigma^C) \mathbf{n}_{\Omega} \rangle}_{=0}.
 \end{aligned}$$

Therefore  $\frac{d\mathcal{E}^c}{dt} \leq - \int_{\Omega} \phi \sigma^C : \dot{\epsilon}$ . Provided the local dissipativity of the contact stress that we discussed in the previous paragraph, the overall material is thus also dissipative.

### 7.1.3 Semi-implicit integration

Unlike in Chapter 6, our rheology involves the volume fraction field  $\phi$ ; the mass conservation and momentum balance equations are now fully coupled. Moreover, as we are using an inelastic model — i.e., in the infinite compression Young modulus limit — numerical stability considerations impose using an integration scheme with some degree of implicitity.

We propose using a two-step semi-implicit algorithm. At the  $k^{\text{th}}$  timestep,

1. we first solve the momentum balance equation (5.10) using the current volume fraction field  $\phi^k$ , and deduce the end-of-step velocity and stress fields  $\mathbf{u}^{k+1}$  and  $\sigma^{k+1}$ ;
2. we then solve the mass conservation equation (5.7) using the end-of-step velocities  $\mathbf{u}^{k+1}$  and deduce  $\phi^{k+1}$ .

However, note that Step 1. requires finding  $\mathbf{u}^{k+1}$  and  $\sigma^{k+1}$  such that the end-of-step velocity constraints are satisfied, and in particular, the end-of-step maximum volume fraction constraint,  $\phi^{k+1} \leq \phi_{\max}$ . Yet  $\phi^{k+1}$  is only computed at Step 2.! We must therefore replace the maximum volume fraction constraint in Step 1. with an approximation of  $\phi^{k+1}$  that can be computed using only  $\phi^k$  and  $\mathbf{u}^{k+1}$ , and is much cheaper to evaluate than solving the full mass conservation equation. The following section is dedicated to this choice of approximation.

**Remark 7.2.** Our algorithm could be made fully implicit by keeping iterating between the two steps until a fixed point for the volume fraction and velocity fields is reached.

**Maximum volume fraction** The mass conservation equation (5.7) describes how the volume fraction changes in time. Yet, as depicted in Figure 7.2, there are two ways to express this

change. On the one hand, the Eulerian point of view looks at the change of  $\phi$  at a fixed location in space  $\mathbf{x}_0$ , which is given by the derivative  $\frac{\partial \phi}{\partial t}$  evaluated at  $\mathbf{x}_0$ . We have

$$\begin{aligned}\phi(\mathbf{x}_0, t) &= \phi(\mathbf{x}_0, t_0) + \int_{t_0}^t \frac{\partial \phi}{\partial t}(\mathbf{x}_0, s) ds \\ &= \phi(\mathbf{x}_0, t_0) - \int_{t_0}^t \nabla \cdot [\phi \mathbf{u}](\mathbf{x}_0, s) ds, \quad \text{using (5.7).}\end{aligned}$$

On the other hand, the Lagrangian point of view follows grains as they move through space. The volume fraction at such a tracked point is given by  $\hat{\phi}(\mathbf{x}_0, t_0; t) := \phi(X(\mathbf{x}_0, t_0; t), t)$ , where  $X(\mathbf{x}_0, t_0; t)$  is the *characteristics* function defined by the Cauchy problem (5.13). Using the chain derivation rule w.r.t. time (5.14) and the mass conservation equation (5.7), we get

$$\frac{\partial \hat{\phi}}{\partial t}(\mathbf{x}_0, t_0; t) = \frac{D\phi}{Dt}(X(\mathbf{x}_0, t_0; t), t) = -\phi \nabla \cdot \mathbf{u}(X(\mathbf{x}_0, t_0; t), t).$$

We thus have

$$\begin{aligned}\hat{\phi}(\mathbf{x}_0, t_0; t) &= \hat{\phi}(\mathbf{x}_0, t_0; t_0) + \int_{t_0}^t \frac{\partial \phi}{\partial t}(\mathbf{x}_0, t_0; s) ds \\ &= \phi(\mathbf{x}_0, t_0) - \int_{t_0}^t (\phi \nabla \cdot \mathbf{u})(X(\mathbf{x}_0, t_0; s), s) ds.\end{aligned}\tag{7.5}$$

The maximum volume fraction constraint,  $\phi \leq \phi_{\max}$ , can be expressed at every instant and every point of the simulation domain using either the Eulerian or Lagrangian point of views. Studying the rate of change eventually leads to conditions on  $\nabla \cdot [\phi \mathbf{u}]$  or  $\phi \nabla \cdot \mathbf{u}$ , respectively. Both conditions are equivalent in the spatially continuous case, however they will yield different discretizations. A simple physical consideration will drive our choice: at the grain scale, the contact forces oppose the relative velocity of the particles; at the macroscopic scale, we can thus expect that the pressure will oppose the (opposite of) the divergence of the flow. Hence, we will use the Lagrangian point of view to discretize the maximum volume fraction constraint.

**First-order approximation** With the Lagrangian framework of Equation (7.5), we linearize the constraint  $\hat{\phi}^{k+1} \leq \phi_{\max}$  on a time interval  $[t^k, t^{k+1} := t^k + \Delta_t]$  as  $(\phi^k - \Delta_t \phi^k \nabla \cdot \mathbf{u}^{k+1}) \leq \phi_{\max}$ . Equation (7.1) can thus be approximated at the first order as

$$0 \leq \phi^k \nabla \cdot \mathbf{u}^{k+1} + \beta^k \perp p^{k+1} + c^k \geq 0,\tag{7.6}$$

where  $\beta^k := \frac{\phi_{\max} - \phi^k}{\Delta_t}$  is a scalar field expressing the maximum rate at which the material can compress during  $\Delta_t$ .

#### 7.1.4 Discrete-time equations

We now perform a few mathematical manipulations that will reveal a problem structure similar to that of the previous chapter. To lighten the notations, we will omit the  $k+1$  superscript; unless otherwise mentioned, the fields are evaluated at the end of the time-step.

As mentioned in Remark 7.1, the frictional contact stress field  $\text{Dev } \sigma^C$  vanishes when the material is not at maximum compaction, and in particular where there is no granular matter ( $\phi^k = 0$ ). One may thus rewrite (7.2) equivalently as

$$\begin{cases} \text{Dev } \sigma^C = \hat{\mu}(p + c^k) \frac{\phi^k \text{Dev } \dot{\epsilon}}{|\phi^k \text{Dev } \dot{\epsilon}|} & \text{if } \phi^k \text{Dev } \dot{\epsilon} \neq \mathbf{0} \\ |\text{Dev } \sigma^C| \leq \hat{\mu}(p + c^k) & \text{if } \phi^k \text{Dev } \dot{\epsilon} = \mathbf{0}. \end{cases}\tag{7.7}$$

We now consider the friction coefficient  $\mu := \sqrt{\frac{2}{d}}\hat{\mu}$ , and the symmetric tensor fields  $\lambda := c^k \mathbb{I} - \sigma^C$  and  $\gamma := \phi \dot{\epsilon} + \frac{\beta}{d} \mathbb{I}$ . We have  $\phi^k \nabla \cdot \mathbf{u} + \beta = \text{Tr } \gamma$ ,  $\phi^k \text{Dev } \dot{\epsilon} = \text{Dev}(\gamma)$ , and  $\mu \lambda_N = \hat{\mu} \sqrt{\frac{2}{d}} \left( \sqrt{\frac{d}{2}} c^k - \frac{\text{Tr } \sigma^C}{\sqrt{2d}} \right) = \hat{\mu} (c^k + p)$ , so that Equations (7.6)–(7.7) reread

$$\begin{cases} \text{Dev } \lambda = -\mu \lambda_N \frac{\text{Dev } \gamma}{|\text{Dev } \gamma|} & \text{if } \text{Dev } \gamma \neq 0 \\ |\text{Dev } \lambda| \leq \mu \lambda_N & \text{if } \text{Dev } \gamma = 0 \\ 0 \leq \gamma_N \perp \lambda_N \geq 0. \end{cases} \quad (7.8)$$

We recognize the disjunctive formulation (1.20) of the non-associated Drucker–Prager flow rule expressed on  $\gamma$  and  $\lambda$ ,

$$\gamma \text{ and } \lambda \text{ satisfy (7.8)} \iff (\gamma, \lambda) \in \mathcal{DP}(\mu).$$

**Conservation of momentum** Replacing the expression of  $\lambda$  into the momentum conservation equation (5.10), the discrete time momentum balance equation that we will have to solve at Step 1. of our semi-implicit integration algorithm reads

$$\rho \phi^k \frac{D\mathbf{u}}{Dt} - \nabla \cdot [\eta \phi^k D(\mathbf{u}) - \phi^k \lambda] = \rho \phi^k \mathbf{g} + \nabla [\phi^k c^k]. \quad (7.9)$$

The total derivative remains to be discretized, but the method for doing so will be tightly linked to our choice of spatial discretization. For now, we will assume that  $\frac{D\mathbf{u}}{Dt}$  can be approximated at the first order as

$$\frac{D\mathbf{u}}{Dt} = \frac{1}{\Delta_t} (\mathbf{u} - u(\mathbf{u}^k)) + O(\Delta_t),$$

with the mapping  $u$  to be defined later.

**Variational formulation** Similarly as in Chapter 6, the momentum conservation equation (7.9) with the non-associated Drucker–Prager flow rule can be put under variational form. Let  $V$  denote the subspace of  $H^1(\Omega)^d$  satisfying the homogeneous Dirichlet boundary conditions, and let  $T \subset L_2(\Omega)_d^s$  a subspace of square-integrable symmetric tensor fields on  $\Omega$ .

**Proposition 7.1.** *The first step of our semi-implicit integration scheme amounts to finding  $\mathbf{u} \in V$ ,  $\gamma, \lambda \in T^2$  such that*

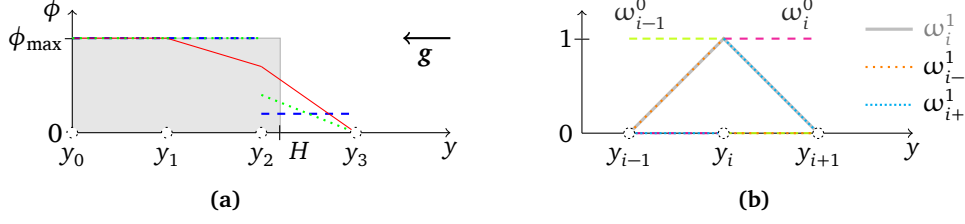
$$\begin{cases} m(\mathbf{u}, \mathbf{v}) + a(\mathbf{u}, \mathbf{v}) = b(\lambda, \mathbf{v}) + l(\mathbf{v}) & \forall \mathbf{v} \in T \\ s(\gamma, \tau) = b(\tau, \mathbf{u}) + k(\tau) & \forall \tau \in T \\ (\gamma, \tau) \in \mathcal{DP}(\mu), \end{cases} \quad (7.10)$$

where the bilinear forms  $a$ ,  $b$ ,  $m$  and linear forms  $l, k$  are given by

$$\begin{aligned} m(\mathbf{u}, \mathbf{v}) &:= \frac{\rho}{\Delta_t} \int_{\Omega} \phi^k \langle \mathbf{u}, \mathbf{v} \rangle & s(\gamma, \tau) &:= \int_{\Omega} \gamma : \tau \\ a(\mathbf{u}, \mathbf{v}) &:= \eta \int_{\Omega} \phi^k D(\mathbf{u}) : D(\mathbf{v}) & b(\tau, \mathbf{u}) &:= \int_{\Omega} \phi^k \tau : D(\mathbf{u}) \\ l(\mathbf{v}) &:= \rho \int_{\Omega} \phi^k \left( \mathbf{g} + \frac{u(\mathbf{u}^k)}{\Delta_t} \right) & k(\tau) &:= \int_{\Omega} \frac{\max(\phi_{\max} - \phi^k, 0)}{\Delta_t} \frac{\text{Tr } \tau}{d}. \end{aligned}$$

*Proof.* One can see easily that the second line of 7.10 is just the definition of  $\gamma$  put under weak form. Multiplying the conservation of momentum (7.9) by a test function  $\mathbf{v}$ , one gets

$$\frac{\rho}{\Delta_t} \int_{\Omega} \phi^k \langle \mathbf{u}, \mathbf{v} \rangle - \eta \int_{\Omega} \langle \nabla \cdot [\phi^k D(\mathbf{u})], \mathbf{v} \rangle = - \int_{\Omega} \langle \nabla \cdot [\phi^k \lambda], \mathbf{v} \rangle + \rho \int_{\Omega} \phi^k \left( \mathbf{g} + \frac{u(\mathbf{u}^k)}{\Delta_t} \right).$$



**Figure 7.3:** (a) Fraction field for a volume ( $H\phi_{\max}$ ) of grains at rest under gravity; exact (gray), and using different approximations: piecewise-linear continuous ( $\mathbb{P}_1$ , line), piecewise-constant ( $\mathbb{P}_0$ , dashed), piecewise-linear discontinuous ( $\mathbb{P}_{1d}$ , dotted). (b) Shape functions for those low-order approximations.

Then, using the Green formula and taking into account that the test function  $\mathbf{v}$  vanishes on the boundary of  $\Omega$ ,

$$b(\lambda, \mathbf{v}) = \int_{\Omega} \phi^k \lambda : D(\mathbf{v}) = - \int_{\Omega} \langle \nabla \cdot [\phi^k \lambda], \mathbf{v} \rangle$$

$$a(\mathbf{u}, \mathbf{v}) = -\eta \int_{\Omega} \langle \nabla \cdot [\phi^k D(\mathbf{u})], \mathbf{v} \rangle.$$

□

We notice that the variational formulation (7.10) is structurally similar to the one governing the dense case, System (6.13a–6.13c). From there, it would seem straightforward to use once again the discretization strategy discussed in Chapter 6. However, we will see in the next section that supplemental restrictions on the choice of elements apply.

Note also the “max” function in the expression of the bilinear form  $k$ , which at first could appear superfluous. However, it may happen that during the simulation,  $\phi$  locally grows above  $\phi_{\max}$ . Without the “max” clamping, the following timestep will attempt to correct the local volume fraction excess, and doing so may introduce energy and thus visual popping artifacts. Indeed, it can be easily verified that the dissipativity property that we exhibited in the time-continuous setting (Section 7.1.2) will remain valid in the discrete case only if  $\beta^k$ , defined such that  $\gamma = \phi \dot{\epsilon} + \beta^k \frac{1}{d}$ , is positive. This is what we ensure here by setting  $\beta^k = \frac{1}{\Delta t} \max(\phi_{\max} - \phi^k, 0)$  instead of its original definition as per Equation (7.6). The main drawback of this approach is that by not attempting to correct volume fraction excesses, the simulation may visually lose volume over time.

## 7.2 Discretization using finite elements

### 7.2.1 Space compability criterion

A simple compability consideration will greatly reduce the choice of finite-element spaces that we can use to discretize our problem. Let us consider the 1D case of a cohesionless, frictionless material. We want to be able to express the equilibrium of the material under gravity; this reads

$$\nabla[\phi p] = -\phi \rho g. \quad (7.11)$$

Moreover, the Drucker–Prager rheology imposes that  $\phi < \phi_{\max} \implies p = 0$ . Necessarily, the equilibrium (7.11) requires that every point for which  $\phi > 0$  is within the stencil of the gradient operator evaluated at a point for which  $\phi = \phi_{\max}$ . In other words,  $\phi$  should decrease “abruptly enough”. Note that this is the case of the physical solution: at equilibrium, there exists  $H$  such that  $\phi(x) = \phi_{\max}$  for  $x \leq H$ , and  $\phi(x) = 0$  for  $x > H$ .

More formally, let  $V_h$  and  $T_h$  be subspaces of  $V \subset H^1(\mathbb{R})$  and  $T \subset L_2(\mathbb{R})$ . The complementarity condition (7.1) suggests using the same discretization space for the volume fraction and stress fields, so we will assume  $\phi \in T$ . Consider the linear form  $l_h$  on  $V_h$ ,  $l_h(\mathbf{v}_h) := \int_{\Omega} \phi_h \mathbf{v}_h$ , and the bilinear form on  $T_h \times V_h$   $b_h(p_h, \mathbf{v}_h) := \int_{\Omega} \nabla[p_h] \mathbf{v}_h$ . Let  $\Lambda(\phi_h) \subset T_h$  be the linear subspace of the pressure fields that are non-zero only at mesh nodes where  $\phi_h = \phi^{max}$ . Then if  $l_h \neq 0_{V_h}$  and  $\forall p_h \in \Lambda(\phi_h), b(p_h, l_h) = 0$ , the equilibrium under gravity cannot be achieved.

Let us see if this criterion is satisfied for a few simple velocity–pressure spaces and the scenario depicted in Figure 7.3(a). We will restrict ourselves to low-order approximations (Figure 7.3(b)):

- (piecewise-constant)  $\mathbb{P}_0 := \text{Span}(\omega_i^0)$ , with  $\omega_i^0(x) := 1$  for  $x_i \leq x \leq x_{i+1}$  and zero elsewhere;
- (piecewise-linear)  $\mathbb{P}_1 := \text{Span}(\omega_i^1)$ , with  $\omega_i^1(x) := \max(0, 1 - |x - x_i|)$ ;
- (piecewise-linear discontinuous)  $\mathbb{P}_{1d} := \text{Span}(\omega_{i+}^1, \omega_{i-}^1)$ , with  $\omega_{i+}^1(x) := \omega_i^1(x)$  for  $x_i \leq x \leq x_{i+1}$  and zero elsewhere;  $\omega_{i-}^1(x) := \omega_i^1(x)$  for  $x_{i-1} \leq x \leq x_i$  and zero elsewhere.

**$T_h$  piecewise-constant** Let us choose  $\phi = \omega_0^0 + \omega_1^0 + \beta \omega_2^0$  with  $0 < \beta < \frac{1}{2}$  and where  $\omega_i^0$  is the constant unitary function over the  $i^{\text{th}}$  element.

Then  $\Lambda = \text{Span}(q_0^0, q_1^0)$ , and thus  $b_h(p_h, \mathbf{v}_h) = 0$  for any  $\mathbf{v}_h$  vanishing on  $[x_0, x_2]$ . The existence of a test function  $\mathbf{v}_h$  with positive integral on  $[x_2, x_3]$  yet vanishing for  $x \leq x_2$  would prevent the existence of an equilibrium solution.

For  $V_h = \mathbb{P}_1$ , we have  $l_h(\omega_3^1) > 0$ , and  $b_h(p_h, \omega_3^1) = 0$ . The pair  $\mathbb{P}_1 - \mathbb{P}_0$  is inconsistent. Similarly for  $V_h = \mathbb{P}_{1d}$ , we have  $l_h(\omega_{3-}^1) > 0$ , while  $b_h(p_h, \omega_{3-}^1) = 0$ . This can be generalized to higher orders.

However, the compatibility criterion does not prohibit the  $\mathbb{P}_0 - \mathbb{P}_0$  pair.

**$T_h$  piecewise-linear** We now choose  $\phi = \omega_0^1 + \omega_1^1 + \beta \omega_2^1$  with  $\frac{1}{2} < \beta < 1$ . This means  $\Lambda = \text{Span}(\omega_0^1, \omega_1^1)$ , and thus again,  $b_h(p_h, \mathbf{v}_h) = 0$  for any  $\mathbf{v}_h$  vanishing on  $[x_0, x_2]$ .

For  $V_h = \mathbb{P}_0$ ,  $l_h(\omega_2^0) > 0$  and  $b_h(p_h, \omega_2^0) = 0$ . Again for  $V_h = \mathbb{P}_1$ ,  $l_h(\omega_3^1) > 0$  and  $b_h(p_h, \omega_3^1) = 0$ . For  $V_h = \mathbb{P}_{1d}$ ,  $l_h(\omega_{3-}^1) > 0$  and  $b(p_h, \omega_{3-}^1) = 0$ . Idem for higher orders.

**$T_h$  piecewise-linear discontinuous** Let  $\phi = \omega_{0+}^1 + \omega_{1-}^1 + \omega_{1+}^1 + \omega_{2-}^1 + \beta \omega_{2+}^1$  with  $0 < \beta < 1$ . This means  $\Lambda = \text{Span}(\omega_{0+}^1, \omega_{1-}^1, \omega_{1+}^1, \omega_{2-}^1)$ , and once again,  $b_h(p_h, \mathbf{v}_h) = 0$  for any  $\mathbf{v}_h$  vanishing on  $[x_0, x_2]$ .

For  $V_h = \mathbb{P}_1$  or  $V_h = \mathbb{P}_{1d}$ ,  $l_h(\mathbf{v}_h) > 0$  and  $b_h(p_h, \mathbf{v}_h) = 0$  for  $\mathbf{v}_h := \omega_3^2$  or  $\mathbf{v}_h := \omega_{3-}^2$ , respectively. However the compatibility criterion does not prohibit  $\mathbb{P}_0$  velocities.

**Wrapping up** This incompatibility seems to arise from the fact that we are using the same discretization space for both the pressure and volume fraction fields; using a higher-order volume fraction would yield a sharper interface, and would allow this interface to stay within the range of the gradient operator stencil. However, in this case the maximal volume fraction constraint would not be able to be satisfied for every degree of freedom of the volume fraction field. In Section 7.3, we will propose to use particles to discretize the volume fraction field, and thus obtain unconditionally sharp interfaces. Alternatively, we can use one of the few discretization spaces that were found to not violate the compatibility criterion, such as the  $\mathbb{P}_0 - \mathbb{P}_0$  pair, to which the next subsection is dedicated.

## 7.2.2 Piecewise-constant discretization

In a manner similar to that of finite-volume methods, we consider only the average value (i.e., the value at the barycenter) of each mesh cell. Discontinuous velocity spaces are theoretically tricky to handle, as they are not subspaces of  $H^1(\Omega)^d$ ; we refer the reader to (Pietro and Ern 2011).

**Discontinuous Green formula** We assume that the discrete mesh  $\Omega_h$  can be decomposed as a set of disjoint polyhedras ( $K_i$ ), and denote by  $\mathbb{P}_{pd}(\Omega_h)$  the space of piecewise-polynomials that are of order less than  $p$  on each element  $K_i$ . Let  $\mathcal{F}_i$  denote the set of all interior faces of  $\Omega_h$ , which we augment with an arbitrary orientation: for any  $F \in \mathcal{F}_i$ , we number as  $K_1^F$  and  $K_2^F$  the adjacent cells, and define the normal  $\mathbf{n}^F$  as pointing towards  $K_2^F$ , i.e.,  $\mathbf{n}^F = \mathbf{n}^{K_1^F} = -\mathbf{n}^{K_2^F}$ .

For any field  $\mathbf{v} \in \mathbb{P}_{pd}(\Omega_h)$ , for any internal face  $F$  and for  $\mathbf{x} \in F$  we denote by  $\mathbf{v}_i(\mathbf{x})$ ,  $i = 1$  or  $2$  the value at  $\mathbf{x}$  of the restriction  $\mathbf{v}$  to  $K_i^F$ ,  $\mathbf{v}_{|K_i}(\mathbf{x})$ . We can then define the *jump* and *average* operators on  $\mathcal{F}_i$  as

$$[\![\mathbf{v}]\!] := \mathbf{v}_1 - \mathbf{v}_2 \quad \langle \mathbf{v} \rangle := \frac{1}{2}(\mathbf{v}_1 + \mathbf{v}_2).$$

These operators allow us to define a Stokes-like theorem for discontinuous fields. Let  $\mathbf{v} \in V_0$ , then

$$\begin{aligned} \int_{\text{Bd } \Omega_h} \langle \mathbf{v}, \mathbf{n}^\Omega \rangle &= \sum_{K \in \Omega_h} \int_{\text{Bd } K} \langle \mathbf{v}, \mathbf{n}^K \rangle - \sum_{F \in \mathcal{F}_i} \int_F \langle \mathbf{v}_1, \mathbf{n}^{K_1^F} \rangle + \langle \mathbf{v}_2, \mathbf{n}^{K_2^F} \rangle \\ &= \sum_{K \in \Omega_h} \int_K \nabla \cdot \mathbf{v} - \sum_{F \in \mathcal{F}_i} \int_F \langle [\![\mathbf{v}]\!], \mathbf{n}^F \rangle. \end{aligned}$$

Let  $V_p = \mathbb{P}_p(\Omega_h)^d$ , and  $T_p = \mathbb{P}_p(\Omega_h)^{s_d}$ . Applying our discontinuous Stokes formula to the product  $\boldsymbol{\tau} \mathbf{v}$ , with  $\boldsymbol{\tau} \in T_p$  and  $\mathbf{v} \in V_h$ , and using the identity  $[\![\boldsymbol{\tau} \mathbf{v}]\!] = [\![\boldsymbol{\tau}]\!] \langle \mathbf{v} \rangle + \langle \boldsymbol{\tau} \rangle [\![\mathbf{v}]\!]$ , we get

$$\begin{aligned} \int_{\text{Bd } \Omega_h} \langle \mathbf{v}, \boldsymbol{\tau} \mathbf{n}^\Omega \rangle + \sum_{F \in \mathcal{F}_i} \int_F \langle [\![\mathbf{v}]\!], \langle \boldsymbol{\tau} \rangle \mathbf{n}^F \rangle - \sum_{K \in \Omega_h} \int_K \text{D}(\mathbf{v}) : \boldsymbol{\tau} = \\ \sum_{K \in \Omega_h} \int_K \langle \nabla \cdot \boldsymbol{\tau}, \mathbf{v} \rangle - \sum_{F \in \mathcal{F}_i} \int_F \langle \langle \mathbf{v} \rangle, [\![\boldsymbol{\tau}]\!] \mathbf{n}^F \rangle. \end{aligned}$$

For the sake of brevity, we can extend the jump and average operators on  $\mathcal{F}_a := \mathcal{F}_i \cup \text{Bd } \Omega_h$  as  $[\![\mathbf{v}]\!]_{|\text{Bd } \Omega_h} = \mathbf{v}$  and  $\langle \mathbf{v} \rangle_{|\text{Bd } \Omega_h} = \mathbf{v}$ , and rewrite this last equation (the discontinuous Green formula) as

$$\sum_{F \in \mathcal{F}_a} \int_F \langle [\![\mathbf{v}]\!], \langle \boldsymbol{\tau} \rangle \mathbf{n}^F \rangle - \sum_{K \in \Omega_h} \int_K \text{D}(\mathbf{v}) : \boldsymbol{\tau} = \sum_{K \in \Omega_h} \int_K \langle \nabla \cdot \boldsymbol{\tau}, \mathbf{v} \rangle - \sum_{F \in \mathcal{F}_i} \int_F \langle \langle \mathbf{v} \rangle, [\![\boldsymbol{\tau}]\!] \mathbf{n}^F \rangle. \quad (7.12)$$

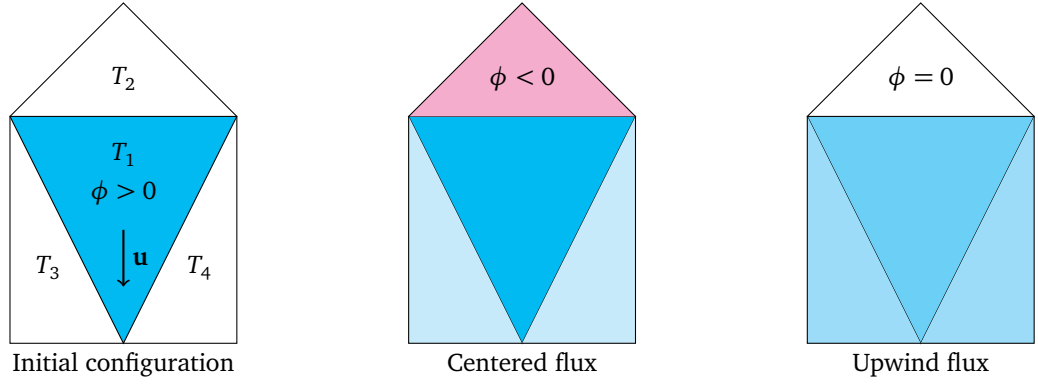
The left-hand-side of (7.12) can be interpreted as a weak expression of  $\text{D}(\mathbf{u})$  (with homogeneous Dirichlet boundary conditions), and the right-hand-side as a weak expression of  $\nabla \cdot \boldsymbol{\tau}$ .

**Discrete variational formulation** The bilinear forms  $m$  and  $s$  from the variational formulation given in Proposition 7.1 do not involve any spatial derivative, and can thus be left unmodified. The discontinuous Green formula (7.12) yields a discretization of the bilinear form  $b$  which can be easily verified to be consistent when the velocity solution is continuous, i.e., when  $[\![\mathbf{u}]\!] = \mathbf{0}$  on  $\mathcal{F}_i$ ,

$$b_h(\mathbf{u}, \boldsymbol{\tau}) = - \sum_{F \in \mathcal{F}_a} \int_F \langle [\![\mathbf{v}]\!], \langle \phi \boldsymbol{\tau} \rangle \mathbf{n}^F \rangle + \sum_{K \in \Omega_h} \int_K \phi \text{D}(\mathbf{v}) : \boldsymbol{\tau}. \quad (7.13)$$

Obviously, the rightmost term in (7.13) vanishes when considering piecewise constant polynomials, and thus  $b_h$  reduces to a sum of interfacial transfer terms.

The discretization of viscosity term, which involves a second-order derivative of the velocity, is more complex. The order of the polynomials in  $\mathbb{P}_0$  is too low to allow the application of standard discontinuous Galerkin methods, such as the Symmetric Interior Penalty method (Pietro and Ern 2011, Section 4.2). Instead, we discretize our viscosity form  $a$  from Proposition 7.1



**Figure 7.4:** The centered transport scheme (middle) will yield a negative volume fraction in  $T_2$  when  $\phi$  (cyan) and  $\mathbf{u}$  are initially non-zero only on  $T_1$ . In contrast, upwind transport (right) guarantees a positive volume fraction.

as the weak derivative of the weak derivative of  $\mathbf{u}$ . That is, we once again take profit of our discontinuous Green formula (7.12), and replace  $a(\mathbf{u}, \mathbf{v})$  with  $c(\dot{\epsilon}, \mathbf{v})$ , where

$$c(\boldsymbol{\tau}, \mathbf{v}) := - \int_{F \in \mathcal{F}_a} \langle \llbracket \mathbf{v} \rrbracket, \langle \boldsymbol{\tau} \rangle \mathbf{n}^F \rangle$$

and  $\dot{\epsilon}$  is defined as the solution to  $s(\dot{\epsilon}, \boldsymbol{\tau}) = c(\boldsymbol{\tau}, \mathbf{u}) \forall \boldsymbol{\tau} \in T_h$ . As the matrix associated to  $s$  is diagonal,  $\dot{\epsilon}$  is trivial to compute.

**Transport of the volume fraction field** We now consider the second step of our semi-implicit algorithm; given a velocity field  $\mathbf{u}$ , computing the volume fraction field at the end of the timestep.

We can list a few desirable properties for this step:

1. admissibility:  $\phi^{k+1}$  should take values in  $[0, 1]$ ;
2. global mass conservation: the integral of the volume fraction over the domain should remain constant, i.e.,  $\frac{\partial}{\partial t} \int_{\Omega} \phi = 0$ ;
3. energy conservation: in the continuous case, there is a correspondence between the work of gravity in the conservation of momentum equation and the loss of potential energy due to the motion of the volume fraction field. We would like to get a similar property in the discrete case, or at least, prevent the artificial introduction of energy.

Unfortunately, it is hard for all of those properties to simultaneously hold in the discrete setting. In our case, we prioritize the first two; in particular, if the volume fraction were to go below zero, the bilinear form  $m$  would no longer be positive and our system would be ill-posed. A correction step could be devised, at the expense of the desirable properties 2 and 3.

For stability reasons, we choose to use an implicit scheme. At the first-order, we want to solve

$$\frac{\phi^{k+1}}{\Delta_t} + \nabla \cdot [\phi^{k+1} \mathbf{u}] = \frac{\phi}{\Delta_t}. \quad (7.14)$$

Equation (7.14) can be put under weak form as  $e(\phi, \psi) = f(\psi) \forall \psi \in L_2(\Omega)$ , where, in the continuous setting,

$$\begin{aligned} e(\phi, \psi) &= \frac{1}{\Delta_t} \int_{\Omega} \phi \psi + \int_{\Omega} \nabla \cdot [\phi \mathbf{u}] \psi = \frac{1}{\Delta_t} \int_{\Omega} \psi \phi - \int_{\Omega} \phi \mathbf{u} \nabla \psi \\ f(\psi) &= \frac{1}{\Delta_t} \int_{\Omega} \phi^k \psi. \end{aligned}$$

We can check that in the continuous case, the variation of potential energy  $\mathcal{E}^p$  corresponds to the work of gravity (Desirable property number 3). Indeed, we get

$$\begin{aligned}\mathcal{E}^p(\phi) - \mathcal{E}^p(\phi^k) &= -\rho \int_{\Omega} (\phi - \phi^k) \langle \mathbf{x}, \mathbf{g} \rangle = \rho \Delta_t \int_{\Omega} \nabla \cdot [\phi \mathbf{u}] \langle \mathbf{x}, \mathbf{g} \rangle \\ &= -\rho \Delta_t \int_{\Omega} \phi \langle \mathbf{u}, \mathbf{g} \rangle.\end{aligned}$$

Now, we could discretize the bilinear form  $e$  as we did for the momentum conservation variational formulation, leading to

$$\begin{aligned}e_h(\phi, \psi) &= \frac{1}{\Delta_t} \int_{\Omega} \phi \psi - \int_{F \in \mathcal{F}_a} \langle \llbracket \phi \mathbf{u} \rrbracket, \langle \psi \rangle \mathbf{n}^F \rangle \\ &= \frac{1}{\Delta_t} \int_{\Omega} \phi \psi + \int_{F \in \mathcal{F}_i} \langle \langle \phi \mathbf{u} \rangle, \llbracket \psi \rrbracket \mathbf{n}^F \rangle.\end{aligned}$$

With this discretization, we can see that the total mass is conserved (Criterion 2); indeed

$$\frac{1}{\Delta_t} \int_{\Omega_h} \phi - \phi^k = e_h(\phi, \mathbb{1}) - \int_{F \in \mathcal{F}_i} \left\langle \langle \phi \mathbf{u} \rangle, \underbrace{\llbracket \mathbb{1} \rrbracket}_0 \right\rangle - f_h(\mathbb{1}) = 0.$$

However, looking at the difference in potential energy, we get, denoting by  $\mathbf{x}_K \in V_0$  the field associating to each mesh cell the position of its barycenter

$$\begin{aligned}\frac{\mathcal{E}^p(\phi) - \mathcal{E}^p(\phi^k)}{\rho \Delta_t} &= - \int_{\Omega} (\phi - \phi^k) \langle \mathbf{x}, \mathbf{g} \rangle = - \int_{\Omega} (\phi - \phi^k) \langle \mathbf{x}_K, \mathbf{g} \rangle \\ &= \underbrace{f_h(\langle \mathbf{x}_K, \mathbf{g} \rangle) - e_h(\phi, \langle \mathbf{x}_K, \mathbf{g} \rangle)}_0 - \int_{F \in \mathcal{F}_a} \langle \llbracket \phi \mathbf{u} \rrbracket, \langle \langle \mathbf{x}_K \rangle, \mathbf{g} \rangle \mathbf{n}^F \rangle.\end{aligned}$$

On the other hand,

$$\begin{aligned}\int_{\Omega} \phi \langle \mathbf{u}, \mathbf{g} \rangle &= \sum_K \int_K \langle \phi \mathbf{u}, \nabla \cdot [\langle \mathbf{x}, \mathbf{g} \rangle \mathbb{I}] \rangle = \sum_K \int_{\text{Bd } K} \langle \phi \mathbf{u}, \mathbf{n}^K \rangle \langle \mathbf{x}, \mathbf{g} \rangle \\ &= \sum_{F \in \mathcal{F}_a} \langle \llbracket \phi \mathbf{u} \rrbracket, \mathbf{n}^F \rangle \langle \mathbf{x}, \mathbf{g} \rangle.\end{aligned}$$

The difference between these two terms is  $\sum_{F \in \mathcal{F}_a} \langle \llbracket \phi \mathbf{u} \rrbracket, \mathbf{n}^F \rangle \langle \mathbf{x} - \langle \mathbf{x}^K \rangle, \mathbf{g} \rangle$ , and is non-zero in the general case. Criterion 3 is thus not satisfied. This discrepancy can be understood as the transported material “jumping” from one barycenter to one other, regardless of their respective proximity to the face.

Moreover, this is a centered scheme; as visible in the expression of  $e_h$ , the flux between cells is defined by the average  $\langle \phi \mathbf{u} \rangle$ . If one cell is initially empty, one of its neighbor is not, and the velocity points toward the non-empty cell, then the initially empty cell may end up with a negative value of  $\phi$ . For instance, consider the  $T_2$  cell in Figure 7.4, where  $\phi \mathbf{u} = \mathbf{e}_y$  on  $T_1$  and  $\phi = 0$  elsewhere. Criterion 1 is thus not satisfied either.

To remediate to this problem, we use an upwind scheme. The flux between cells will then be computed considering only the volume fraction of the upwind cell. As we are using an implicit scheme, if this value becomes zero the flux will vanish, and thus the volume fraction will not become negative. For any internal face  $F$ , let  $\phi^{\text{upw}}$  denote the value of  $\phi$  the value on the upstream cell, that is,  $\phi_i$  such that  $\langle \mathbf{n}^{K_i}, \langle \mathbf{u} \rangle \rangle \geq 0$ . We want the flux between cells to be  $\phi^{\text{upw}} \langle \mathbf{u} \rangle$ ;



the upwind bilinear form  $e^{\text{upw}}$  can thus be obtained as

$$\begin{aligned} e_h^{\text{upw}}(\phi, \psi) &= e_h(\phi, \psi) + \sum_{\mathcal{F}_i} \llbracket \phi \rrbracket \langle \phi^{\text{upw}}(\llbracket \mathbf{u} \rrbracket) - \langle \phi \mathbf{u} \rangle, \mathbf{n}^F \rangle \llbracket \psi \rrbracket \\ &= e_h(\phi, \psi) + \frac{1}{2} \sum_{\mathcal{F}_i} \llbracket \phi \rrbracket \left( |\langle \llbracket \mathbf{u} \rrbracket, \mathbf{n}^F \rangle| - \frac{1}{2} \langle \llbracket \mathbf{u} \rrbracket, \mathbf{n}^F \rangle \right) \llbracket \psi \rrbracket \\ &= \frac{1}{\Delta_t} \int_{\Omega} \phi \psi + \int_{F \in \mathcal{F}_i} \llbracket \psi \rrbracket \left( \langle \langle \phi \rangle \llbracket \mathbf{u} \rrbracket, \mathbf{n}^F \rangle + \frac{1}{2} \llbracket \phi \rrbracket |\langle \llbracket \mathbf{u} \rrbracket, \mathbf{n}^F \rangle| \right). \end{aligned}$$

With a similar reasoning to that of the centered case, we can verify that the upwind discretization verify Criterion 2 on the total mass conservation. However, energy conservation is not better than for the centered discretization. Criterion 1 and 2 being our priorities, we still choose to use the upwind bilinear form and take  $\phi^{k+1}$  as the solution of  $e_h^{\text{upw}}(\phi, \psi) = l(\psi) \forall \psi \in \mathbb{P}_0(\Omega_h)$ .

**Inertial terms** Finally, it remains to deal with the velocity advection term  $(\mathbf{u} \cdot \nabla) \mathbf{u}$  in the conservation of momentum equation. As the velocities are discontinuous, the characteristics method cannot be used (the morphism  $X$  is no longer locally invertible). Pietro and Ern (2011, Chapter 6) demonstrate that discretizing implicitly the transport trilinear form  $t(\mathbf{w}, \mathbf{u}, \mathbf{v})$ ,

$$t(\mathbf{w}, \mathbf{u}, \mathbf{v}) := - \int_{\Omega} \langle (\mathbf{w} \cdot \nabla) \mathbf{u}, \mathbf{v} \rangle,$$

using to the so-called *Temam* modification preserves the kinetic energy of the system. However, this modification is only valid for incompressible flows, and yields an asymmetric system; it is therefore not applicable in our case. In the end, we simply used an explicit version of the upwind transport scheme described in the previous paragraph, which we found yielded in practice a good compromise between stability and conservation of kinetic energy. Thus, we define  $u(\mathbf{u}^k)$  in the computation of the  $l$  linear form from Property 7.1 such that

$$\begin{aligned} \int_{\Omega_h} \langle u(\mathbf{u}^k), \mathbf{v} \rangle &:= \int_{\Omega_h} \langle \mathbf{u}^k, \mathbf{v} \rangle - t_h(\mathbf{u}^k, \mathbf{u}^k, \mathbf{v}) \\ t_h(\mathbf{w}, \mathbf{u}, \mathbf{v}) &:= - \sum_{F \in \mathcal{F}_i} \left\langle \llbracket \mathbf{u} \rrbracket, \langle \langle \mathbf{w} \rangle, \mathbf{n}^F \rangle \langle \mathbf{v} \rangle + \frac{1}{2} |\langle \langle \mathbf{w} \rangle, \mathbf{n}^F \rangle| \llbracket \mathbf{v} \rrbracket \right\rangle. \end{aligned}$$

**Final discrete system** The final form of the variational formulation from Property 7.1 after  $\mathbb{P}_0$ - $\mathbb{P}_0$  discretization is thus:

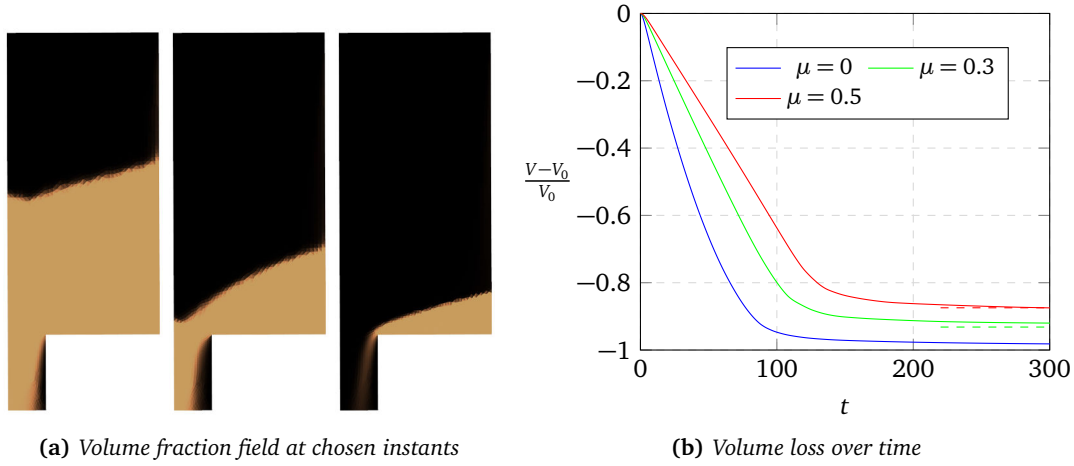
Find  $\underline{\mathbf{u}} \in \mathbb{R}^v$ ,  $\underline{\boldsymbol{\lambda}}, \underline{\boldsymbol{\gamma}} \in \mathbb{R}^{n_{s_d}}$ ,

$$\begin{cases} (M + C^\top S^{-1} C) \underline{\mathbf{u}} = B^\top \underline{\boldsymbol{\lambda}} + \underline{\mathbf{l}} \\ S \underline{\boldsymbol{\gamma}} = B \underline{\mathbf{u}} + \underline{\mathbf{k}} \\ (\underline{\boldsymbol{\gamma}}, \underline{\boldsymbol{\lambda}}) \in \mathcal{DP}(\mu), \end{cases} \quad (7.15)$$

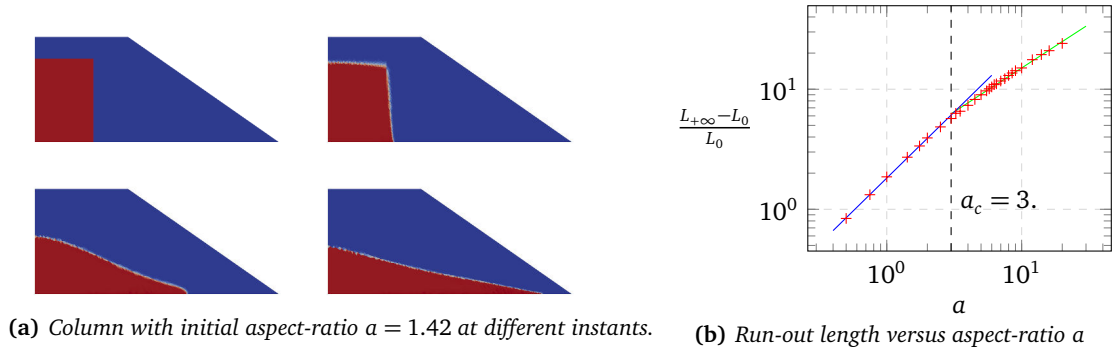
where  $M$  and  $S$  are diagonal, positive matrices. As all the integrals of the variational formulation vanish on empty mesh cells, we can only consider those in which the volume fraction is non-zero, and thus  $M$  is positive-definite. This has also the advantage of drastically reducing the size of the numerical system in some scenarios. Alternatively, we can take into account the presence of a light phase in the computation of the density in each cell, as in (Lagrée et al. 2011). In both cases, System (7.15) can be then solved in the exact same manner as for our dense flows; see Section 6.4.

### 7.2.3 Results

We attempted to validate our piecewise-constant discretization strategy by looking at some empirical laws reported in the literature. Just like in the dense case, the finite-element discretization was done with the Rheolef library (Saramito 2015).



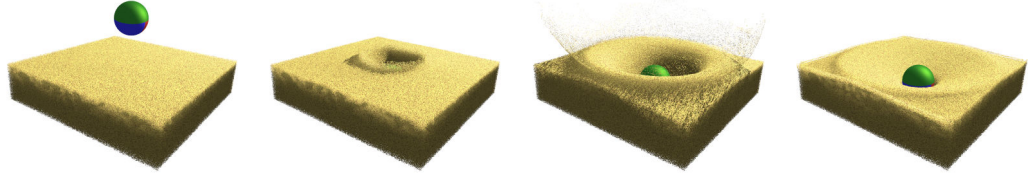
**Figure 7.5:** Evolution of the volume fraction of grains in a discharging silo.



**Figure 7.6:** Collapse of a 2D granular column using the  $\mu(I)$  rheology with for  $\mu_s = 0.32$ ,  $\mu_d = 0.6$ ,  $I_0 = 0.4$  and  $Re = 10^4$ .

**Silo discharge** The first experiment is quite similar to the one that we already studied for dense flows in Section 6.5.3. However, we are no longer interested in the steady discharge regime, but rather in the evolution of the volume of remaining matter in the 2D silo through time, as illustrated in 7.5(a). Staron et al. (2012) recall that one distinguishing feature of an hourglass is that the sand flow rate remains constant for most of the discharge process. This contrasts with clepshydras, which contain a standard Newtonian liquid and whose discharge rate decreases with time. Figure 7.5(b) demonstrates that we retrieve this behavior; for  $\mu = 0$ , the discharge rate  $\frac{dV}{dt}$  is mostly decreasing, while for  $\mu = 0.3$  and  $\mu = 0.5$ , it remains constant as long as the volume  $V$  stays in  $[0.9V_0, 0.2V_0]$ , where  $V_0$  denotes the initial volume of grains in the silo. The dashed horizontal lines in Figure 7.5(b) represent the theoretical final volume of grains computed from the Mohr–Coulomb rest angle corresponding to each friction coefficient.

**Granular column collapse** Another classical experiment that we can now run thanks to our volume fraction field concerns the 2D collapse of a granular column, as depicted in Figure 7.6(a). Note that this scenario is known to be approachable by continuum models; Ionescu et al. (2015) managed to quantitatively match experiments using the  $\mu(I)$  rheology. An empirical law relates the run-out length of the flow  $L_\infty$  (that is, the maximum abscissa reached by the grains) to the initial aspect-ratio of the column (that is, the ratio  $a$  of its height  $H_0$  to its width  $L_0$ ). The



**Figure 7.7:** A steel ball impacts a granular bed. Unlike previous approaches, our fully resolved (non-linearized) Drucker-Prager rheology allows us to retrieve a perfectly round crater.

experiments of Lajeunesse et al. (2005) led to the distinction of two regimes:

$$\frac{L_\infty - L_0}{L_0}(a) \propto \begin{cases} a & \text{if } a < a_c \\ a^{\frac{2}{3}} & \text{if } a > a_c, \end{cases}$$

with  $a_c \sim 3$ . The numerical simulations of Staron and Hinch (2005), using DEM, and Lagr  e et al. (2011), using the  $\mu(I)$  rheology, yielded similar scaling laws, although with different critical aspect ratios,  $a_c \sim 2$  and 7. Several other authors have performed similar experiments both experimentally and numerically, see (Dunatunga and Kamrin 2015, Table 3) for a summary of their results. Using our method with the physical parameters from Lagr  e et al. (2011), we obtained the scaling law plotted in Figure 7.6(b):

$$\frac{L_\infty - L_0}{L_0}(a) \sim \begin{cases} 1.82a^{1.1} & \text{if } a < 3 \\ 2.8a^{0.73} & \text{if } a > 3. \end{cases}$$

The exponents that we obtained in both regimes seem to be a bit on the high side, though they remain within the range of reported values in the numerical simulation literature, especially when taking into account the results of (Mast et al. 2014). Note that at the characteristic sizes chosen for our spatial and temporal discretization, the choice of the advection scheme for the volume fraction and velocity fields retains a non-negligible influence on the final scaling law.

**Discussion** This fully mesh-based discretization strategy suffers from several drawbacks. We observe diffusion of the volume fraction field over time, as a lot of cells retain slightly-above-zero volume fraction even once the bulk of the material has moved away. The explicit transport scheme for the inertial term also prohibits taking large timesteps, making the simulation process quite slow. In the next section we propose to use a different strategy and resort to an hybrid method, which will discretize the volume fraction field using a set of particles. These particles will also be used to transport other quantities, such as the inertial terms.

### 7.3 Material Point Method

The basic idea behind MPM is to consider that the whole mass of the material is condensed at a finite number of material points. The volume fraction field  $\phi$  can therefore be expressed as  $\phi(\mathbf{x}, t) = \sum_p V_p \delta(\mathbf{x} - \mathbf{x}_p(t))$ , where  $\delta$  is the Dirac distribution, and  $V_p$ ,  $\mathbf{x}_p$  are the volume and position of the  $p^{\text{th}}$  particle, respectively. The total volume of granular reads  $V_{\text{tot}} = \int_\Omega \phi(t) = \sum_p V_p$ .

This representation is particularly adapted to the FEM formalism, where quantities are evaluated in weak-form, by multiplying them by a test function and integrating over a domain. For instance, for  $\mathbf{v} \in L_2(\Omega)$ , we have  $\int_\Omega \phi(t) \mathbf{v} = \sum_p V_p \mathbf{v}(\mathbf{x}_p(t))$ . In this sense, the Material Point Method can be seen as a quadrature rule for which the quadrature points are given by the particles positions, and the corresponding weights by the particles volumes.

**Related work** The MPM method was introduced by Sulsky (1994) for the simulation of history-dependent elastic materials. Storing information on particles made tracking the position of each material point through time trivial, avoiding solving complex transport equations.

MPM was quickly recognized as relevant for the simulation of granular materials. Wieckowski, Youn, et al. (1999) proposed a first explicit algorithm for the integration of an elasto-plastic material subject to the non-associated Drucker–Prager flow rule<sup>1</sup> in 2D, and recently extended it to the 3D setting (Wieckowski and Pawlak 2015). S. Bardenhagen et al. (2000) proposed to use MPM for simulating the dynamics of many 2D deformable grains with sliding and rolling contacts. Konagai and Johansson (2001) coupled the Mohr–Coulomb yield criterion with a nodal finite-difference scheme to simulate 2D granular flows. More recently, Mast et al. (2014) performed 2D granular column collapse simulations using MPM with the Matzuo–Nakai yield surface. Dunatunga and Kamrin (2015) proposed a MPM method with two distinct regimes depending on whether the volume fraction is above or below a critical value, again for 2D flows.

Stomakhin, Schroeder, Chai, et al. (2013) popularized MPM to the Computer Graphics community, coupling this method with an elasto-plastic model to produce visually-striking snow simulations. This gave rise to a variety of related methods modeling a wide range of physical phenomenons, such as visco-elasto-plastic foams (Ram et al. 2015; Yue et al. 2015), phase changes (Stomakhin, Schroeder, Jiang, et al. 2014), granular materials (e.g., concurrently to the work presented below, Klar et al. 2016), and more (Jiang, Schroeder, Teran, et al. 2016). Klar et al. (2016) follow the approach of Wieckowski, Youn, et al. (1999) with one notable difference: while both prevent plastic compression, Klar et al. (2016) allow plastic expansion, performing Alart–Curnier-like projections onto the Drucker–Prager set of admissible stresses.

Most MPM methods are explicit (Dunatunga and Kamrin 2015; Wieckowski, Youn, et al. 1999; Yue et al. 2015), or semi-implicit with respect to the elastic energy, but not w.r.t. the plasticity (Ram et al. 2015; Stomakhin, Schroeder, Chai, et al. 2013). Klar et al. (2016) and Mast et al. (2014) propose to use a “return-mapping” algorithm to perform semi-implicit integration taking into account the plasticity. In practice, this means a Newton algorithm prone to falling into local minima. In contrast, our semi-implicit integration procedure will leverage our DCFP solvers that showed good empirical convergence.

Extensions of MPM using different discretization strategies for the volume fraction field have also been proposed; for instance, the Generalized Interpolation Material Point Method (GIMP; S. G. Bardenhagen and Kober 2004) use arbitrary indicator functions for the particles.

### 7.3.1 Application to our variational formulation

Taking advantage of the expression of  $\phi$  as a sum of Dirac distributions, we can now rewrite the bilinear and linear forms from the variational formulation of Property 7.1 as:

$$\begin{aligned} m(\mathbf{u}, \mathbf{v}) &:= \sum_p \frac{\rho}{\Delta_t} V_p \langle \mathbf{u}(\mathbf{x}_p^k), \mathbf{v}(\mathbf{x}_p^k) \rangle & s(\boldsymbol{\gamma}, \boldsymbol{\tau}) &:= \int_{\Omega} \boldsymbol{\gamma} : \boldsymbol{\tau} \\ a(\mathbf{u}, \mathbf{v}) &:= \sum_p V_p \eta D(\mathbf{u})(\mathbf{x}_p^k) : D(\mathbf{v})(\mathbf{x}_p^k) & b(\boldsymbol{\tau}, \mathbf{u}) &:= \sum_p V_p \boldsymbol{\tau}(\mathbf{x}_p^k) : D(\mathbf{u})(\mathbf{x}_p^k) \\ l(\mathbf{v}) &:= \sum_p \rho V_p \left\langle \left( \frac{u(\mathbf{u}^k)}{\Delta_t} + \mathbf{g} \right), \mathbf{v}(\mathbf{x}_p^k) \right\rangle & k(\boldsymbol{\tau}) &:= \int_{\Omega} \frac{\phi_{\max}}{\Delta_t} \frac{\mathbb{I}}{d} : \boldsymbol{\tau} - \sum_p \frac{V_p}{\Delta_t} \frac{\mathbb{I}}{d} : \boldsymbol{\tau}(\mathbf{x}_p) \end{aligned}$$

and where  $\mathbf{x}_p^k$  denotes the position of the  $p^{\text{th}}$  particle at timestep  $k$ .

### 7.3.2 Grid–particles transfers

In the context of MPM, the semi-implicit integration algorithm reads as follow:

<sup>1</sup>which is expressed as the association of the Drucker–Prager yield surface with a plastic incompressibility criterion, that is,  $\hat{\boldsymbol{\varepsilon}}_p \in \mathcal{N}_{\mathcal{B}_T(-\mu\sigma_N)}(\boldsymbol{\sigma}_T)$

1. solve the variational formulation (7.10), where the advected velocity field  $u(\mathbf{u}^k)$  is reconstructed from the particles;
2. deduce the new particle velocities  $\mathbf{v}_p^{k+1}$  from the velocity field  $\mathbf{u}^{k+1}$ , and move the particles accordingly.

Two steps remain to be defined; first, the computation of  $u(\mathbf{u}^k)$  (particles-to-grid transfer), and the computation of the particles velocities (grid-to-particle transfer). We mentioned in Section 5.2.4 different strategies for doing these transfers, such as PIC and FLIP. Both of those methods compute  $\mathbf{u}^{p \rightarrow g} := u(\mathbf{u}^k)$  at each mesh node as a weighted interpolation of the velocities of the surrounding particles, as per Equation (5.16). As we are using the FEM formalism however, the natural way to compute  $\mathbf{u}^{p \rightarrow g}$  would be to solve

$$\frac{\Delta_t}{\rho} m(\mathbf{u}^{p \rightarrow g}, \mathbf{v}) = \int_{\Omega} \phi(\mathbf{u}, \mathbf{v}) = \sum_p V_p \langle \mathbf{v}_p^k, \mathbf{v}(\mathbf{x}_p^k) \rangle,$$

leading to a linear system  $M \underline{\mathbf{u}}^{p \rightarrow g} = \underline{\mathbf{v}}$ , where

$$M_{i,j} := \sum_p V_p \langle \mathbf{v}_i(\mathbf{x}_p^k), \mathbf{v}_j(\mathbf{x}_p^k) \rangle$$

and  $\mathbf{v}_i := \sum_p V_p \mathbf{v}_p^k \mathbf{V}_i(\mathbf{x}_p^k)$ .  $M$  is called the *consistent mass matrix*.

In contrast, the standard PIC/FLIP update (5.16) can be understood as solving a diagonal system  $\tilde{M} \underline{\mathbf{u}}^{p \rightarrow g} = \underline{\mathbf{v}}$ , with

$$\tilde{M}_{i,j} := \sum_p V_p \langle \mathbf{v}_i(\mathbf{x}^k + p), \mathbb{1} \rangle \delta_i^j \quad (7.16)$$

Assume that  $\forall \mathbf{x} \in \Omega$ ,  $\sum \mathbf{V}_j(\mathbf{x}) = \mathbb{1}_{\mathbb{R}^d}$ , as is the case for standard Lagrange FEM. Then  $\tilde{M}_{i,j} = \sum_j M_{i,j}$ , and  $\tilde{M}$  is called the *diagonalized* or *lumped* mass matrix. Alternatively,  $\tilde{M}$  can be understood as the result of using a trapezoidal approximation of the integral in the bilinear form  $m$ .

Computing  $\mathbf{u}^{p \rightarrow g}$  using  $M$  instead of  $\tilde{M}$  yields much better kinetic energy conservation, alleviating the main drawback of PIC (Burgess et al. 1992). Indeed, the kinetic energy discrepancy is  $O(\Delta_t |\mathbf{u}|)$  for the consistent mass matrix, instead of  $O(|\mathbf{u}|)$  for the lumped mass matrix. However, we can see from the interpretation of the particles as quadrature points that  $M$  will be indefinite if there are not enough particles in one given mesh cell. Using the consistent mass matrix is therefore rarely possible in practice.

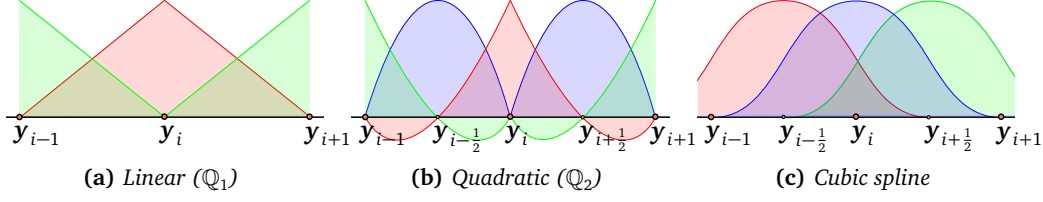
Moreover, FLIP is not convenient in our case. Indeed, consider a particular falling on a flat granular bed; we want the particle to stop as soon as it reaches the ground. As the granular bed is not moving, the velocity delta ( $\mathbf{u}^{k+1} - \mathbf{u}^k$ ) will be zero inside the ground; the FLIP update rule (5.19) will thus not modify the falling particle's velocity, and this particle will dive inside the ground. Using a weighted PIC/FLIP update rule does not fully correct this problem.

Instead, we will opt for the strategy recently proposed by Jiang, Schroeder, Selle, et al. (2015), the Affine Particle-in-Cell (APIC) method. The original PIC velocity update rule (5.17) is then used; in our falling particle scenario, this means that the particle velocity will become equal to the ground velocity ( $\mathbf{0}$ ) as soon as it reaches it. The kinetic energy loss due to the lumped-mass matrix is alleviated by storing information about the velocity gradient on the particles, and using this information during the particle-to-grid transfer step. Instead of computing the velocity at mesh nodes as a weighted average of particle velocities, they are computed as the weighted average of the node velocities *as seen by the particles*.

The exact procedure depends on the shape functions chosen for the velocity field; in the case of  $d$ -linear shape functions, with  $\mathbf{V}_{d(i-1)+k} = \omega_i^v \mathbf{e}_k$ , the expression boils down to

$$\underline{\mathbf{u}}^{p \rightarrow g} = \tilde{M}^{-1} \underline{\mathbf{v}}^{\text{APIC}} \quad (7.17)$$

$$\mathbf{v}_{3(i-1)+k}^{\text{APIC}} = \sum_p V_p (\mathbf{v}_{p,k}^k \omega_i^v(\mathbf{x}_p^k) + \mathbf{c}_{p,k}^k (\mathbf{y}_i - \mathbf{x}_p^k)) \quad (7.18)$$



**Figure 7.8:** A few possible choices for MPM velocity shape functions

where  $y_i$  denotes the position of grid node  $i$  and  $\mathbf{c}_{p,k}$  the  $k^{\text{th}}$  row of  $\mathbf{c}_p^k$ , the  $3 \times 3$  velocity gradient matrix computed at the end of the previous step as

$$\mathbf{c}_p^k = (\nabla \mathbf{u})(t^k, \mathbf{x}_p^k) = \sum_i \mathbf{u}(t^k, y_i) (\nabla \omega_i^v)(\mathbf{x}_p^k). \quad (7.19)$$

See (Jiang, Schroeder, Selle, et al. 2015, Sections 5 and 6) for more details.

**Adaptive resampling** Even when using the lumped mass matrix to avoid indefinite systems, for the MPM discretization to be justified the particles needs to remain well-distributed over the mesh cells. However, if no resampling is performed, expansion of the material will inevitably lead to a scarcity of particles. Conversely, if too many particles become concentrated in a single cell, the system might become overconstrained when using particle-based stress shape functions.

For these reasons, we follow the split/merge heuristics of Narain, Golas, et al. (2010, Section 3). To each particle is associated an ellipsoid representing its occupied volume. Such an ellipsoid with semi-axis  $(\mathbf{a}_{p,i})_{1 \leq i \leq d}$  defines a symmetric, positive-definite tensor  $F_p := \sum \mathbf{a}_{p,i} \mathbf{a}_{p,i}^T$ , whose change in time satisfies

$$\frac{dF_p}{dt} = (D(\mathbf{u}) + W(\mathbf{u}))(\mathbf{x}_p(t))F_p + F_p(D(\mathbf{u}) - W(\mathbf{u}))(\mathbf{x}_p(t)). \quad (7.20)$$

At each timestep, we assume  $\mathbf{x}_p(t)$  constant and compute the solution to (7.20) using matrix exponentiation.

Particles are then split if one axis of the ellipsoid becomes much longer than the others. The reverse operation is performed when two close particles can be merged into a more isotropic one.

Still following Narain, Golas, et al. (2010), these frames are useful not only at simulation time, but also at the rendering stage. Indeed, they define volumes from which the passively-advectioned grain samples are not allowed to escape. This ensures that a rendering sample cannot drift to a place where no simulated particle is present, i.e., to a place where the velocity field has no physical meaning.

### 7.3.3 Shape functions

The Material Point Method prescribes how the volume fraction field should be discretized, but suitable discretizations for the velocity and stress fields remain to be chosen.

**Background mesh** The background mesh used for the computation of the momentum conservation equation is typically chosen to be a regular grid (which allows efficiently locating particles), but triangle or tetrahedral-based meshes have also been used (e.g., Wieckowski and Pawlak 2015; Wieckowski, Youn, et al. 1999). Stomakhin, Schroeder, Jiang, et al. (2014) use a staggered grid to simplify the handling of incompressibility constraints.

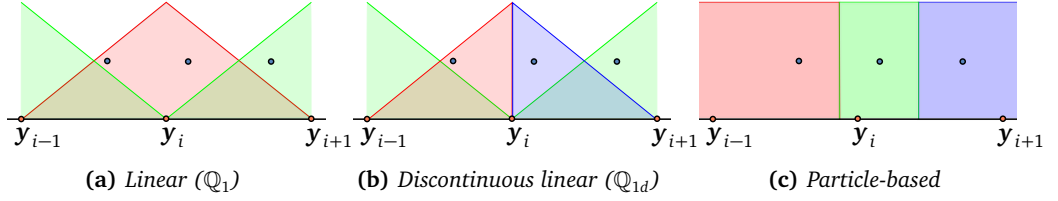


Figure 7.9: A few possible choices for MPM stress shape functions

**Velocity field** Historically, linear ( $\mathbb{P}_1$ ) or multi-linear ( $\mathbb{Q}_1$ ) elements (Figure 7.8(a)) have been favored for the discretization of the velocity field (S. Bardenhagen et al. 2000; Dunatunga and Kamrin 2015; Mast et al. 2014; Sulsky 1994; Wieckowski and Pawlak 2015; Wieckowski, Youn, et al. 1999). However, higher-order shape functions have also been proposed:  $d$ -quadratic ( $\mathbb{Q}_2$ ; Figure 7.8(b)) elements (S. Andersen and L. Andersen 2010), which feature only  $C^0$  continuity, or quadratic and cubic splines (Steffen et al. 2008; Figure 7.8(c)), which boast  $C^1$  continuity but span over several grid elements, have non-zero values at more than one node. Due to their smoothness, cubic splines have been especially popular for Computer Graphics applications, and are used for instance in (Klar et al. 2016; Ram et al. 2015; Stomakhin, Schroeder, Chai, et al. 2013; Yue et al. 2015).

Steffen et al. (2008) argue that  $C^1$  continuity is required to avoid so-called *cell-crossing artifacts*. Indeed, without this property, the test velocity gradient field  $D(\mathbf{v})$  is discontinuous, and the integration error increases drastically when a particle crosses a cell-boundary during a timestep. They showed that on some simple test cases linear shape functions did not yield convergence when the spatial resolution was increased, while quadratic and cubic splines performed better. However, for rough discretizations (i.e., while still within the convergence regime of linear shape functions), using higher-order splines did not provide much more accuracy. S. Andersen and L. Andersen (2010) compared  $\mathbb{Q}_1$ ,  $\mathbb{Q}_2$  elements and cubic splines on scenarios less prone to the cell-crossing instabilities, and observed once again that cubic splines were not much more accurate than linear elements, yet helped to better resolve collisions.  $\mathbb{Q}_2$  (and serendipity) elements yielded much lower integration error, but are subject to a major issue impeding their use in practice: as the  $\mathbb{Q}_2$  (or  $\mathbb{P}_2$ ) shape functions can take negative values, when there are not enough particles per cell they may yield negative coefficients in the lumped mass matrix  $\tilde{M}$ .

As our primary application for this hybrid method is Computer Graphics, we will prioritize computational efficiency over spatial convergence properties. This will drive our choice to use  $d$ -linear shape functions, as this discretization strategy will lead to much sparser matrices (there is less overlap between the different shape functions than for, say, quadratic or cubic splines), and thus much more efficient numerical solvers.

**Stress field** In contrast, MPM is pretty opinionated about the way the stress field should be discretized, and a large majority of the methods from the literature choose to make the stress shape functions coincide with the particles. That is, the basis  $(\omega^\tau)$  is chosen so that  $\omega_q^\tau(\mathbf{x}_p) = \delta_q^p$ , which amounts to associating to each particle its own stress  $\sigma_p$ , as in Figure 7.9(c). In most methods, the stress field is evaluated only at the particle positions, so there are no need to precisely define the shape of the stress basis elsewhere. However, there are a few reasons for looking for alternatives in our case.

- First, a peculiarity of our method is that we define a critical volume fraction  $\phi_{\max}$ . This manifests itself by the presence in the linear form  $k$  of the variational formulation of the integral  $V(\boldsymbol{\tau}) := \int_{\Omega} \frac{\mathbb{I}}{d} : \boldsymbol{\tau}$ . Computing  $k(\mathbf{T}_p)$  thus requires being able to evaluate the “volume”  $V(\omega_p^\tau)$  of each stress shape function. This volume has an actual physical meaning: if it is larger than  $V_p$ , the flow may still compress, while if it is equal to or below  $V_p$ , the critical density has already been reached. We must thus be able to evaluate the volume



of each  $(\omega_p^\tau)$ , the amount of space available for each particle to fill-up. A first solution is to compute these volumes from the particles positions; for instance, by constructing a Voronoi diagram, or by computing the distance of each particle to its  $k$ -nearest-neighbors. Another solution, more in line with the MPM philosophy, and proposed by Dunatunga and Kamrin (2015), is to track these volumes over time. Actually, this quantity is readily available as the volume of the frames used for the particle resampling strategy of Narain, Golas, et al. (2010) presented in the previous section; indeed,  $V(\omega_p^\tau) = \sqrt{\det F_p}$ , where  $F_p$  is given by Equation (7.20).

- A second drawback concerns the number of constraints of the system. As we saw in Chapter 6 for dense flows, the number of stress shape functions can be directly linked to the number of constraints in our system. A higher number of constraints leads to a more expensive system to solve, so it might be beneficial to use less than one stress shape function per particle. Moreover, when using multilinear shape functions for the velocities, the number of degrees of freedom of the system is reduced to  $d$  per grid node, and using one constraint per particle may lead to an overconstrained system. Indeed, Mast et al. (2014) describes a kinematic locking phenomenon for nearly incompressible systems when using multilinear velocities and particle-based stresses. A similar phenomenon can manifest itself for our unilateral incompressibility criterion; it suffices for a single particle in a cell to be maximally compacted (i.e.,  $V(\omega_p^\tau) = V_p$ ) to prevent negative divergence of the velocity inside the cell, even if the other particles of the cell have still room for compression.

We thus advocate considering multilinear shape functions for the stresses as well; in many cases, this will lead to much faster solving. Note however that this choice can also create artifacts in some scenarios.

- Stable configurations maintained by cohesion are hard to achieve. Indeed, as Remark 7.1 imposes that cohesion vanishes at nodes where  $\phi < \phi_{\max}$ , a structure maintained by cohesion will progressively erode unless particles are in a very special configuration.
- Another failure case for multilinear stresses is when a big lump of material falls in a neighboring cell to light particles. The momentum of the falling heavy particles will be transmitted to the light ones, resulting in a disturbing kicking effect (Figure 7.10). This scenario is prone to appear in the bottom compartment of an hourglass, for instance.

In both cases, reverting to particle-based stresses will alleviate the problem. An alternative solution is to consider multilinear-discontinuous  $\mathbb{Q}_{1d}$  stress shape functions, as in Figure 7.9(b), at the cost of having to solve  $2^d$  constraints per (active) grid cell. As this space does not depend on particle history and maintains a stable constraints–degrees-of-freedom ratio,  $\mathbb{Q}_{1d}$  is less prone to kinematic-locking than particle-based stresses, but may still induce grid artifacts in the presence of cohesion. Figure 7.11 illustrates that using discontinuous multilinear shape functions allows simulating the bottom compartment of an hourglass.

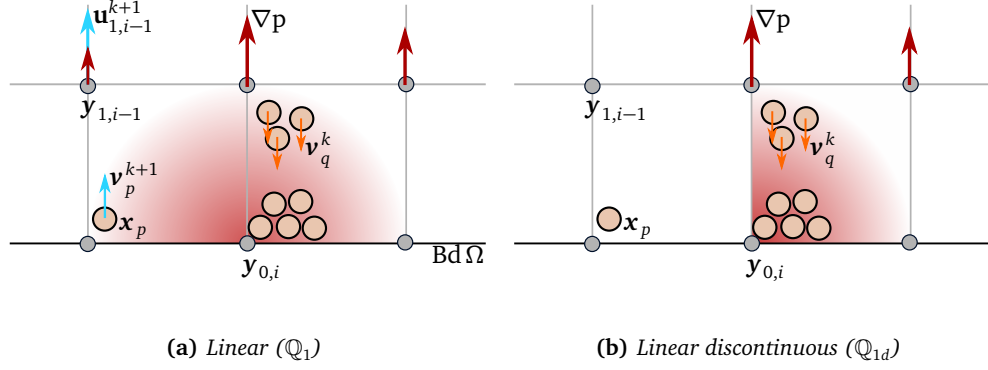
**Compability criterion** Finally, the choice of discrete basis for the velocity and stress still has to satisfy a compability criterion to allow the expression of equilibria, as argued in Section 7.2.1. We found  $\mathbb{Q}_1$ – $\mathbb{Q}_1$ ,  $\mathbb{Q}_1$ – $\mathbb{Q}_{1d}$  and  $\mathbb{Q}_1$ –*particle* to behave well, however  $\mathbb{P}_1$ – $\mathbb{P}_{1d}$  on a regular tetrahedral grid failed to reach an equilibrium state after a column collapse. In this case, switching to piecewise discontinuous stresses,  $\mathbb{P}_1$ – $\mathbb{P}_{1d}$ , was sufficient to restore stability, as the cost of a much increased number of constraints.

#### 7.3.4 Numerical resolution

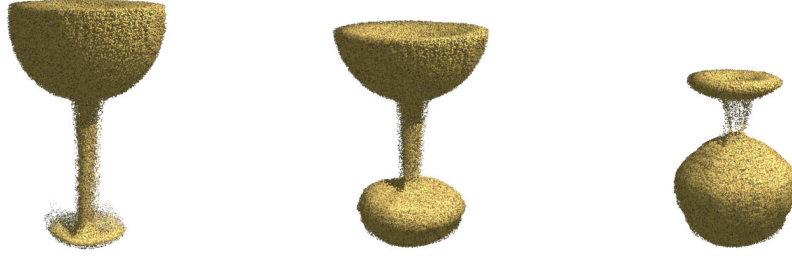
Once our velocity and stress basis functions have been chosen, the procedure described in Section 6.3 can be applied so that we end up with the discrete system (7.21),

$$\begin{cases} (\tilde{M} + A)\underline{u} = B^\top \underline{\lambda} + \underline{l} \\ \underline{\gamma} = B\underline{u} + \underline{k} \\ (\underline{\gamma}_{[i]}, \underline{\lambda}_{[i]}) \in \mathcal{DP}(\mu) \quad \forall 1 \leq i \leq n. \end{cases} \quad (7.21)$$





**Figure 7.10:** A clump of particles falls on the ground near an isolated one. (a) When using multilinear stresses, the pressure (red) at the node  $y_{0,i}$  induces an upwards force and velocity at the node  $y_{1,i-1}$ , kicking the isolated particle  $x_p$  away. (b) When using discontinuous multilinear stresses, this force vanishes.



**Figure 7.11:** Hourglass simulation using discontinuous stress shape functions

The solvers presented in Section 6.4 can then be used; however, they were found to be quite slow in 2D, and thus unlikely to scale up to three-dimensional problems. Moreover, as we have taken the shear yield stress of the Drucker–Prager law to be zero, System (7.21) is now exactly a DCFP, though one were the inverse of the stiffness matrix,  $(\tilde{M} + A)^{-1}$ , is dense. We will thus make an approximation that will allow us to leverage faster DCFP solvers.

**Two-step algorithm** We make the assumption that in the solid and liquid zones, the Newtonian viscosity is small w.r.t. the frictional contact forces; that is, as we are mainly considering gravity-driven flows, that  $\text{Re} := \frac{\rho g^{\frac{1}{2}} L^{\frac{3}{2}}}{\eta}$  is high. We then proceed to solve an approximated version of (7.21) in two steps.

We first solve the unconstrained momentum balance (7.22) using a conjugate-gradient algorithm,

$$(\tilde{M} + A)\underline{u}^* = \underline{l}. \quad (7.22)$$

The frictional response of the material is then computed by neglecting the effect of the change in Newtonian stress due to the addition of contact stresses. That is, we write  $\underline{u} = \underline{u}^* + \Delta \underline{u}$ , and solve a new DCFP:

Find  $\Delta \underline{\mathbf{u}}, \underline{\boldsymbol{\lambda}}, \underline{\boldsymbol{\gamma}}$  s.t.

$$\begin{cases} \tilde{\mathbf{M}} \Delta \underline{\mathbf{u}} = \mathbf{B}^\top \underline{\boldsymbol{\lambda}} \\ \underline{\boldsymbol{\gamma}} = \mathbf{B} \Delta \underline{\mathbf{u}} + \mathbf{B} \underline{\mathbf{u}}^* + \underline{\mathbf{k}} \\ (\underline{\boldsymbol{\gamma}}_{[i]}, \underline{\boldsymbol{\lambda}}_{[i]}) \in \mathcal{DP}(\mu) \quad \forall 1 \leq i \leq n. \end{cases} \quad (7.23)$$

Doing so, we fully preserve the impact of the Newtonian viscosity in the gaseous regime, while neglecting its effect within zones that are dominated by static friction. Our DCFP's stiffness matrix is now diagonal, positive-definite and therefore trivial to inverse. This positively impacts the performance in two ways:

- We get a significant speedup using the solvers of Section 6.4, as the linear system solve at each iteration of the inner problem is now replaced with a single multiplication by a diagonal matrix.
- The Delassus operator of (7.23) is now sparse and easy to assemble, which opens the way to a wider range of algorithms.

In practice, we leverage our Gauss–Seidel algorithm from Chapter 4. Note though that in 3D, as the dimension of the constraint is 6, the analytical polynomial-based local solver cannot be used. However, we found that the Fischer–Burmeister-based solver worked well enough for the hybrid strategy not to be required. We also found the matrix-free variant of the algorithm to be more efficient than the original one, especially for  $\mathbb{Q}_{1d}$  and particle-based stresses where the number of constraints is much higher than the number of velocity degrees of freedom. Attempting to assemble the Delassus operator with such discretization strategies could also quickly overwhelm our amount of available memory, even for relatively small test cases.

### 7.3.5 Overview of a time-step

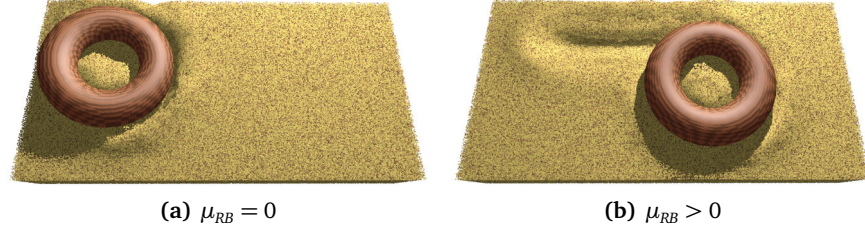
The full MPM algorithm can be summarized as follows:

1. Compute the lumped mass matrix  $\tilde{\mathbf{M}}$  with Eq (7.16);
2. Recover nodal velocities  $\underline{\mathbf{u}}^{p \rightarrow g}$  from System (7.17);
3. Use Section 6.3 to assemble the matrices  $\mathbf{A}$  and  $\mathbf{B}$  and vectors  $\underline{\mathbf{l}}$  and  $\underline{\mathbf{k}}$  corresponding to the forms given in Section 7.3.1;
4. Solve for the unconstrained velocity  $\underline{\mathbf{u}}^*$  with Eq (7.22);
5. Solve DCFP (7.23) and get the total velocity  $\underline{\mathbf{u}} = \underline{\mathbf{u}}^* + \Delta \underline{\mathbf{u}}$ ;
6. Update particles frames with Eq (7.20), then split or merge them according to Section 7.3.2;
7. Compute the velocity gradient matrix  $\mathbf{c}_p$  using Eq (7.19);
8. Update particle positions and velocities as per Eq (5.17, 5.18);
9. Proceed to next time-step.

The quantities that need to be stored on particles are therefore  $V_p$ ,  $\mathbf{x}_p$ ,  $\mathbf{v}_p$ ,  $\mathbf{c}_p$  and  $F_p$ . We also store the stress field  $\boldsymbol{\lambda}$  between timesteps, in order to warm-start the DCFP solver.

## 7.4 Extensions

While collapsing granular columns may be deemed interesting in themselves, Computer Graphics applications often require the simulated material to interact with other external objects, such as an animated character. In this section, we present two extensions of our framework that could potentially allow artists to drive the simulation, or experiment with different look-and-feels by tuning the anisotropy of the flow. We insist on the fact that these “hacks” are designed to allow the visual enrichment of the granular flow, and have not been validated against any physical experiment.



**Figure 7.12:** A slightly tilted wheel, on which a constant torque is applied, is dropped on a sandy ground. (a) Without wheel-sand friction, the wheel falls down immediately. (b) Otherwise, it manages to roll on for some distance.

#### 7.4.1 Rigid body coupling and frictional boundaries

In the following, we describe the modifications that have to be made to the variational formulation from Property 7.1 and the DCFP (7.23) to perform two-way coupling with rigid bodies. Frictional boundary conditions will then be imposed through the means of rigid bodies with infinite inertia.

**Hard boundary** For the sake of simplicity, we consider here the case where the rigid body  $\Omega_{RB}$  is defined by an exact boundary,  $\text{Bd } \Omega_{RB}$ . An heuristic derivation for the case of a diffuse interface is given in (Daviet and Bertails-Descoubes 2016a).

We consider an additional stress field inside the rigid body,  $-\lambda_{RB} \in T_{RB} \subset L_2(\Omega_{RB})^{s_d}$ . The reaction force applied by the rigid body onto the granular material through the boundary is given by  $\mathbf{r} = \lambda_{RB} \mathbf{n}_{RB}$ , with  $\mathbf{n}_{RB}$  the normal to  $\text{Bd } \Omega_{RB}$ . Let  $\hat{\mathbf{v}}$  denote the relative velocity between the granular material and the rigid body; we want  $\mathbf{r}$  and  $\hat{\mathbf{v}}$  to follow a Coulombic relationship. A convenient way to enforce this in our framework is to add the condition  $(\lambda_{RB}, \gamma_{RB}) \in \mathcal{DP}(\mu_{RB})$ , where  $\gamma_{RB} := \frac{1}{2}(\hat{\mathbf{v}} \mathbf{n}_{RB}^T + \mathbf{n}_{RB} \hat{\mathbf{v}}^T)$  and the coefficient  $\mu_{RB}$  sets the intensity of friction between the granular material and the rigid body (see Figure 7.12). The rationale behind this constraint is exposed in Appendix C.2. This formulation with two stress fields allows us to separate the constraints associated to grains–grains contacts from those associated to grains–rigid-body contacts, and thus to set distinct friction coefficients, as illustrated in Figure 7.12.

**Proposition 7.2.** Let  $\mathbf{v}_{RB} \in \mathbb{R}^6$  be the generalized velocity vector of the rigid body, i.e., the concatenation of the linear and angular velocities. The coupled granular–rigid body system then satisfies the variational formulation:

Find  $\mathbf{u} \in V$ ,  $(\lambda, \gamma) \in T$ ,  $(\lambda_{RB}, \gamma_{RB}) \in T_{RB}$ ,

$$\left\{ \begin{array}{ll} m(\mathbf{u}, \mathbf{v}) + a(\mathbf{u}, \mathbf{v}) = b(\lambda, \mathbf{v}) + b_{RB}(\lambda_{RB}, \mathbf{v}) + l(\mathbf{v}) & \forall \mathbf{v} \in V \\ \mathbf{w}^T M_{RB} \mathbf{v}_{RB} = \mathbf{w}^T \mathbf{f} + c_{RB}(\mathbf{w}, \lambda_{RB}) & \forall \mathbf{w} \in \mathbb{R}^6 \\ s(\gamma, \mathbf{u}) = b(\tau, \mathbf{u}) + k(\tau) & \forall \tau \in T \\ s(\gamma_{RB}, \mathbf{u}) = b_{RB}(\tau, \mathbf{u}) + c_{RB}(\tau, \mathbf{v}_{RB}) & \forall \tau \in T_{RB} \\ (\gamma, \lambda) \in \mathcal{DP}(\mu) & \\ (\gamma_{RB}, \lambda_{RB}) \in \mathcal{DP}(\mu_{RB}) & \end{array} \right. \quad (7.24)$$

where  $M_{RB}$  and  $\mathbf{f}$  are the inertia matrix and vector of external forces of the rigid body (see Section 2.1.1), and where we have introduced two new bilinear forms,  $b_{RB} : T_{RB} \times V \rightarrow \mathbb{R}$  and

$$c_{RB} : T_{RB} \times \mathbb{R}^6 \rightarrow \mathbb{R},$$

$$\begin{aligned} b_{RB}(\boldsymbol{\tau}, \mathbf{u}) &:= \frac{1}{2} \int_{\text{Bd } \Omega_{RB}} (\mathbf{u} \mathbf{n}_{RB}^\top + \mathbf{n}_{RB} \mathbf{u}^\top) : \boldsymbol{\tau} \\ c_{RB}(\boldsymbol{\tau}, \mathbf{v}_{RB}) &:= -\frac{1}{2} \int_{\text{Bd } \Omega_{RB}} (J(\mathbf{v}_{RB}) \mathbf{n}_{RB}^\top + \mathbf{n}_{RB} J(\mathbf{v}_{RB})^\top) : \boldsymbol{\tau}, \end{aligned}$$

with  $J(\mathbf{v}_{RB})$  the velocity field spawned by a rigid body with linear and angular velocities  $\mathbf{v}_{RB}$  and center of mass  $\mathbf{x}_{RB}$ ,

$$J(\mathbf{v}_{RB})(\mathbf{x}) = J \begin{pmatrix} \mathbf{v}_{RB}^{lin} \\ \mathbf{v}_{RB}^{ang} \end{pmatrix}(\mathbf{x}) = \mathbf{v}_{RB}^{lin} + \mathbf{v}_{RB}^{ang} \wedge (\mathbf{x} - \mathbf{x}_{RB}).$$

Discretizing System (7.24) once again yields a DCFP, which can be solved as usual.

*Proof.* By definition  $\hat{\mathbf{v}} = \mathbf{u} - J(\mathbf{v}_{RB})$ , so the weak definition of  $\boldsymbol{\gamma}_{RB}$  as per the fourth line of (7.24) is consistent.

We will now justify that the terms  $b_{RB}(\boldsymbol{\lambda}_{RB}, \mathbf{v})$  and  $c_{RB}(\boldsymbol{\lambda}_{RB}, \mathbf{w})$  added to the momentum conservation equations of the granular material and rigid body (i.e., the first two lines of (7.24)) do correspond to the interaction between the two media.

We recall the identity, for any symmetric tensor  $\boldsymbol{\tau} \in S_d$ , and for  $(\mathbf{v}, \mathbf{n}) \in \mathbb{R}^d \times \mathbb{R}^d$ ,

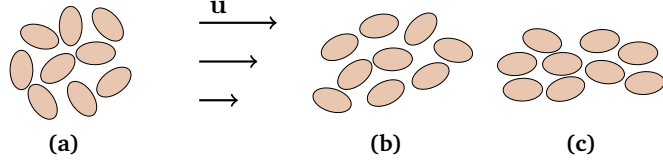
$$\mathbf{v}^\top \boldsymbol{\tau} \mathbf{n} = \sum_i \mathbf{v}_i \sum_j \tau_{ij} \mathbf{n}_j = \sum_{i,j} \mathbf{v}_i \mathbf{n}_j \tau_{i,j} = \boldsymbol{\tau} : \mathbf{v} \otimes \mathbf{n} = \frac{1}{2} \boldsymbol{\tau} : (\mathbf{v} \otimes \mathbf{n} + \mathbf{n} \otimes \mathbf{v}). \quad (7.25)$$

The force  $\mathbf{f}^C$  and torque  $\mathbf{L}^C$  applied onto the rigid body by the granular material are given by (Batty et al. 2007):

$$\mathbf{f}^C = \int_{\Omega_{RB}} \nabla \cdot [-\boldsymbol{\lambda}_{RB}] \quad \mathbf{L}^C = \int_{\Omega_{RB}} (\mathbf{x} - \mathbf{x}_{RB}) \wedge \nabla \cdot [-\boldsymbol{\lambda}_{RB}].$$

For any  $\mathbf{r} \in \mathbb{R}^3$ , there holds

$$\begin{aligned} \mathbf{r}^T \mathbf{f}^C &= \mathbf{r}^T \int_{\Omega_{RB}} \nabla \cdot [-\boldsymbol{\lambda}_{RB}] = - \int_{\text{Bd } \Omega_{RB}} \mathbf{r}^T \boldsymbol{\lambda}_{RB} \mathbf{n}_{RB} \\ &= - \int_{\text{Bd } \Omega_{RB}} \frac{1}{2} \boldsymbol{\lambda}_{RB} : (\mathbf{n}_{RB} \otimes \mathbf{r} + \mathbf{r} \otimes \mathbf{n}_{RB}) \quad \text{using (7.25)} \\ \mathbf{r}^T \mathbf{L}^C &= \mathbf{r}^T \int_{\Omega_{RB}} (\mathbf{x} - \mathbf{x}_{RB}) \wedge \nabla \cdot [-\boldsymbol{\lambda}_{RB}] = - \int_{\Omega_{RB}} \mathbf{r}^T [\mathbf{x} - \mathbf{x}_{RB}]_\wedge \nabla \cdot \boldsymbol{\lambda}_{RB} \\ &= \int_{\Omega_{RB}} ([\mathbf{x} - \mathbf{x}_{RB}]_\wedge \mathbf{r})^\top \nabla \cdot \boldsymbol{\lambda}_{RB} \\ &= \int_{\text{Bd } \Omega_{RB}} ([\mathbf{x} - \mathbf{x}_{RB}]_\wedge \mathbf{r})^\top \boldsymbol{\lambda}_{RB} \mathbf{n}_{RB} - \underbrace{\int_{\Omega_{RB}} D([\mathbf{x} - \mathbf{x}_{RB}]_\wedge \mathbf{r}) : \boldsymbol{\lambda}_{RB}}_{=0} \\ &= \int_{\text{Bd } \Omega_{RB}} \frac{1}{2} \boldsymbol{\lambda}_{RB} : (\mathbf{n}_{RB} \otimes ([\mathbf{x} - \mathbf{x}_{RB}]_\wedge \mathbf{r}) + ([\mathbf{x} - \mathbf{x}_{RB}]_\wedge \mathbf{r}) \otimes \mathbf{n}_{RB}) \quad \text{using (7.25)} \\ &= - \int_{\text{Bd } \Omega_{RB}} \frac{1}{2} \boldsymbol{\lambda}_{RB} : (\mathbf{n}_{RB} \otimes (\mathbf{r} \wedge [\mathbf{x} - \mathbf{x}_{RB}]) + (\mathbf{r} \wedge [\mathbf{x} - \mathbf{x}_{RB}]) \otimes \mathbf{n}_{RB}), \end{aligned}$$



**Figure 7.13:** Rationale for our handling of anisotropy. (a) (b) A shear flow induces a privileged orientation. (c) Aligned particles yield anisotropic macroscopic friction: reduced in the horizontal direction, more intense in the vertical one.

which means,  $\forall \mathbf{w} \in \mathbb{R}^6$ ,

$$\mathbf{w}^\top \begin{pmatrix} \mathbf{f}^C \\ \mathbf{L}^C \end{pmatrix} = c_{\text{RB}}(\boldsymbol{\lambda}_{\text{RB}}, \mathbf{w}),$$

and hence justifies the second line of (7.24).

Concerning the momentum conservation of the granular material, using (7.25),

$$b(\boldsymbol{\lambda}_{\text{RB}}, \mathbf{v}) = \int_{\text{Bd } \Omega_{\text{RB}}} \langle \boldsymbol{\lambda}_{\text{RB}} \mathbf{n}_{\text{RB}}, \mathbf{v} \rangle,$$

which does correspond to the force applied by the rigid body onto the granular material through the boundary.  $\square$

#### 7.4.2 Anisotropy

Many granular materials are composed of anisotropic grains that are thinner along one direction. This is typically the case for corn flakes, or, more common in the Computer Graphics imaginary, piles of gold coins (see Figure 7.23). It is noteworthy that anisotropy at the grain scale does play a role on the collective granular behavior. Indeed, while there is no reason to favor any particular direction when the grains are randomly oriented, when all their normals are aligned the macroscopic friction becomes anisotropic: it has much less dissipative effect when grains are sliding on top of each other, rather than when their relative displacement is along their common normal (see Figure 7.13 for an illustration).

However, remember that each one of our particles does not represent a single grain, but a collection of them, and thus it may contain different orientations. Rather than mapping a normal to each particle, we should instead store a probability distribution function (PDF)  $\psi(\mathbf{n})$ . For efficiency purposes we store only its second moment, the symmetric tensor  $\boldsymbol{\nu}_2 = \int_{S^2} \mathbf{n} \mathbf{n}^\top \psi(\mathbf{n}) d\mathbf{n}$ . In the following, we propose a way to model the influence of  $\boldsymbol{\nu}_2$  on our  $\mathcal{DP}(\mu)$  rheology, then construct an heuristic evolution equation for this tensor.

**Modification of  $\mathcal{DP}(\mu)$**  Anisotropy can be included in the rheology by simply replacing the norm  $|\cdot|$  with the norm  $|\cdot|_{\mathcal{A}}$  associated to the scalar product  $\langle \cdot, \cdot \rangle_{\mathcal{A}} := \frac{1}{2}(\cdot : \mathcal{A}) : (\mathcal{A} : \cdot)$ , where  $\mathcal{A}$  is a fourth-order symmetric tensor. In order for the maximum dissipation principle to remain satisfied,  $\gamma$  should also be replaced with  $\mathcal{A} : \gamma$  on the first line of Equation (7.8).

We choose to construct  $\mathcal{A}$  such as  $\mathcal{A} : \boldsymbol{\tau} = N \boldsymbol{\tau} N$ , i.e.,  $\mathcal{A}_{ijkl} = N_{ik} N_{jl}$ , where the symmetric tensor  $N$  is deduced from the normal orientation tensor  $\boldsymbol{\nu}_2$  following  $N^{-1} = (1 - \alpha)\mathbb{I} + d\alpha\boldsymbol{\nu}_2$ , where  $0 \leq \alpha \leq 1$  is a dimensionless coefficient parameterizing the amount of anisotropy in the frictional law. Using this formula, the effective friction coefficient will remain equal to  $\mu$  for isotropic orientations (the  $d$  eigenvalues of  $\boldsymbol{\nu}_2$  equal to  $\frac{1}{d}$ ). If all the normals are oriented in the same direction (a single non-zero eigenvalue equal to 1), the effective friction coefficient will be  $(1 + (d - 1)\alpha)^2 \mu$  in the normal direction, and  $(1 - \alpha)^2 \mu$  in the tangential ones.

Instead of modifying our friction solver, we can introduce the linear operator  $\underline{N} : \underline{\lambda} \mapsto \frac{\text{Tr} \underline{\lambda}}{d} \mathbb{I} + N \text{Dev} \underline{\lambda} N$ , make the changes of variable  $\underline{\tilde{\lambda}} = \underline{N} \underline{\lambda}$ ,  $\underline{\tilde{\gamma}} = \underline{N} \underline{\gamma}$  and replace the DCFP (7.23) with

$$\begin{aligned} \tilde{M} \Delta \underline{u} &= B^T \underline{N}^{-1} \underline{\tilde{\lambda}} \\ \underline{\tilde{\gamma}} &= \underline{N}^{-1} B \Delta \underline{u} + \underline{N}^{-1} (B \underline{u}^* + \underline{k}) \\ (\underline{\gamma}_{[i]}, \underline{\lambda}_{[i]}) &\in \mathcal{DP}(\mu) \quad \forall 1 \leq i \leq n. \end{aligned}$$

It now remains to describe how the second moment of the orientation normal PDF,  $\nu_2$ , changes in time.

**Evolution of  $\nu_2$**  The evolution of the orientation of rigid ellipsoids in a shearing Newtonian flow has been described by Jeffery (1922). Integrating it over the unit sphere of all possible orientations (see e.g., Folgar and Tucker 1984) yields the evolution of the second moment tensor  $\pi_2$  of the orientation PDF,

$$\frac{D\pi_2}{Dt} = W(\underline{u})\pi_2 + \pi_2 W(\underline{u}) + \ell(D(\underline{u})\pi_2 + \pi_2 D(\underline{u}) - 2D(\underline{u}) : \pi_4) \quad (7.26)$$

where  $0 \leq \ell \leq 1$  is a coefficient describing the elongation of the ellipsoid,  $\ell = \frac{L/W-1}{L/W+1}$ , and  $\pi_4$  the fourth moment of the orientation PDF. The parameter  $\ell$  affects the tendency of the ellipsoids to align with the flow. Note that this evolution equation does not directly apply to our problem:

- Jeffery's equation consider the velocity of a surrounding Newtonian fluid. In our case, the velocity field  $\underline{u}$  is that of the granular matter itself.
- This model is only valid for dilute suspensions. Many authors have postulated laws for extending it to higher particle concentrations (e.g., Folgar and Tucker (1984) added an additional dissipative term modeling random collisions between particles in the semi-concentrated regime).

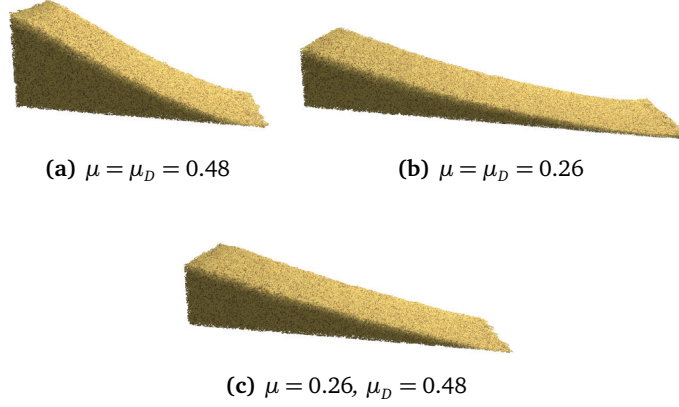
We will nonetheless take inspiration from this evolution equation and propose the following algorithm for the evolution of our normal orientation tensor  $\nu_2$ , which yields good enough results for our purposes — for high values of  $\ell$ , flat particles in a shearing flow do tend to become parallel to each other, which corresponds to the expected behavior (see Figures 7.13 and 7.15).

1. Deduce  $\pi_2(t)$  from  $\nu_2(t)$  as  $\pi_2 = \frac{1}{d-1} P(\mathbb{I} - D)P^T$ , where  $P$  and  $D$  are given by the eigen decomposition of  $\nu_2$ .
2. Explicitly compute  $\pi_2(t + \Delta_t)$  using Equation (7.26) and the quadratic approximation of the  $\pi_4$  tensor,  $\pi_4 \sim \pi_2 \otimes \pi_2$  — the outer tensor product of  $\pi_2$  with itself.
3. Deduce  $\nu_2(t + \Delta_t)$  from  $\pi_2(t + \Delta_t)$ , and normalize it using  $\nu_2 \leftarrow \frac{\nu_2}{\|\nu_2\|_1}$ .

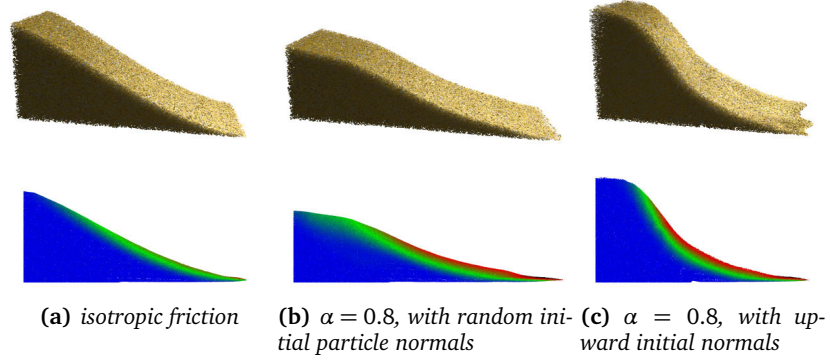
## 7.5 Results

We used this hybrid framework to perform 3D simulations of granular materials. Rendering was done following the method of Narain, Golas, et al. (2010), yielding several samples for each simulation particle, but using rasterization instead of ray-tracing. All simulation code, including the finite-element discretization, relied on a purposely-built library, that has since then been released as open-source<sup>2</sup>. Unless otherwise mentioned, trilinear shape functions were used for the stresses. The DCFP were solved using the bogus library. As usual, all our simulations and benchmarks were run on Intel<sup>®</sup> Xeon<sup>®</sup> quad-core workstations.

<sup>2</sup><http://bipop.inrialpes.fr/~gdaviet/code/sand6>



**Figure 7.14:** Collapsed granular columns with different friction coefficients. Remark that the dynamic friction coefficient  $\mu_D$  of the  $\mu(I)$  rheology influences the run-out length, but not the rest angle.



**Figure 7.15:** Snapshot and corresponding velocity field of the anisotropic column collapse with  $\mu = 0.48$  at  $t = 2.72$ .

### 7.5.1 Model problems

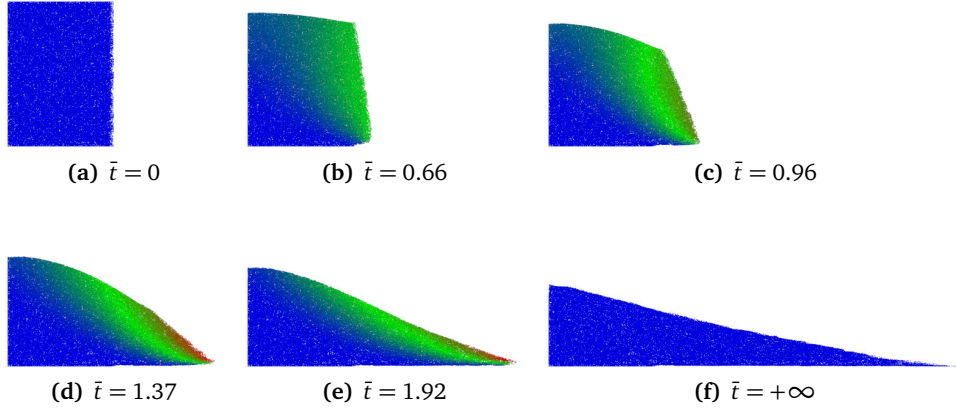
Here we study the influence of several parameters of our model on the collapse of a 3D rectangular column — the experiment setting as for the  $\mathbb{P}_0$  discretization in Section 7.2.3, though with finite (and constant) depth. We refer the reader to the accompanying video <sup>3</sup> for a more comprehensive view of their dynamical effects.

**Influence of friction coefficients** Figure 7.14 depicts the final, stable states following the collapse for the constant  $\mu$  and  $\mu(I)$  rheologies. While  $\mu$  is directly related to the slope of the final granular heap (higher  $\mu$  implies steeper slope),  $(\mu_D - \mu)$  has an effect on the dynamic regime, and therefore on the horizontal spread of the heap. Cranking up the Newtonian viscosity parameter  $\eta$  also reduces this spread, but gives a muddier look to the simulation.

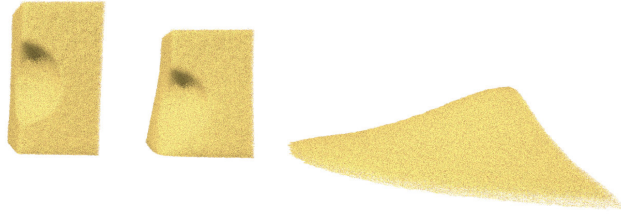
**Influence of anisotropy** Figure 7.15 compares the profiles and velocities of the column collapse for different anisotropy settings. The anisotropic collapse with random initial orientations yields a longer run-out, as the particles align with the flow and minimize friction. The distinct frictional responses of the yielded and unyielded zones are even more visible when initial orientations are uniform.

<sup>3</sup><http://bipop.inrialpes.fr/~gdaviet/files/mpm/mpmGranular.mp4>





**Figure 7.16:** Orthographic views of a collapsing column with aspect-ratio 1.4 at selected instants in time, to be compared with Lagrée et al. 2011, Figure 9. Colors denote particles velocities (blue slowest, black fastest), and can be compared with a similar visualization in Ionescu et al. 2015, Figure 2.



**Figure 7.17:** Revisiting Y. Zhu and Bridson 2005's column collapse.

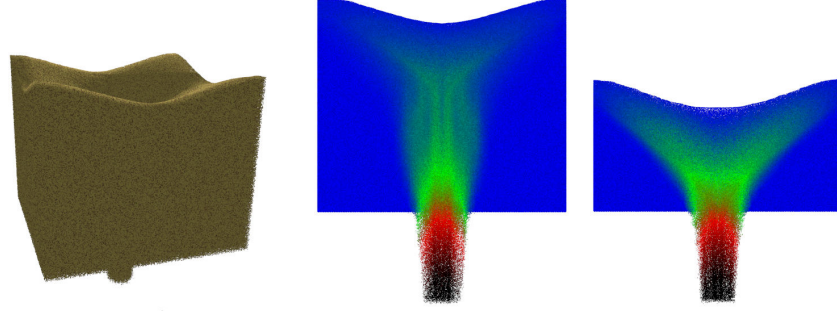
**Comparisons** Continuum simulations using the  $\mu(I)$  rheology were recently validated against DEM simulations (Lagrée et al. 2011) and real experiments (Ionescu et al. 2015) of a 2D granular column collapse. We have reproduced the experiment of Lagrée et al. (2011, Figure 9). The friction parameters  $\mu = 0.26$ ,  $\mu_D = 0.48$  of our 3D simulation were chosen so as to match those from (Lagrée et al. 2011) using the midpoint Drucker-Prager surface (see Section 0.3.2). Figure 7.16 shows that we retrieve the correct profiles throughout time. While our final relative height,  $H_\infty/H_0 = 0.56$ , quantitatively matches the 2D experiments, our final dimensionless run-out length,  $\frac{L_\infty - L_0}{L_0} \sim 2.98$ , is slightly slower.

Figure 7.17 depicts representative frames of our simulation of the sand column collapse scenario introduced by Y. Zhu and Bridson (2005). Our results are comparable to those later obtained by Narain, Golas, et al. (2010), where the incompressibility constraint was relaxed. However, Figure 7.1 illustrates that replacing the Frobenius norm of the Drucker-Prager law with  $\ell_\infty$ , as done in (Narain, Golas, et al. 2010), yields anisotropic artifacts, as the effective friction coefficient then ranges from  $\mu$  to  $\sqrt{5}\mu$  depending on the flow direction.

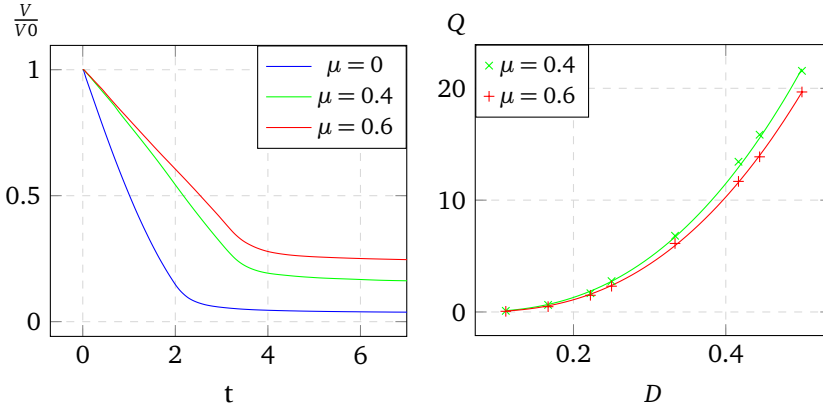
**Retrieving the empirical laws of silo discharge** Focusing our hybrid method on the simulation of a 3D rectangular granular silo with round outlet (Figure 7.18), we attempt to retrieve the empirical laws that were captured by our 2D simulations, the constant discharge rate (Section 7.2.3) and the Beverloo scaling (Section 6.5.3).

Figure 1.16(a), shows that a high enough friction coefficient is necessary to observe the constant discharge rate, while Figure 1.16(b) demonstrates that we do retrieve the power  $\frac{5}{2}$  scaling of the 3D Beverloo law,  $Q = C(D - k)^{\frac{5}{2}}$  where  $C$  and  $k$  are properties of the material and





**Figure 7.18:** Perspective view and orthographic projection of the particle velocities at different instants of the silo discharge.



**(a)** Normalized remaining volume over time **(b)** Simulated discharge rate vs outlet width (marks), and Beverloo fits (lines)

**Figure 7.19:** Constant discharge rate and Beverloo scaling for a 3D granular silo.

silo geometry. However due to our relatively rough discretization, the resulting  $k$  corresponds to about two grid cells – instead of a few grain diameters in reported experiments.

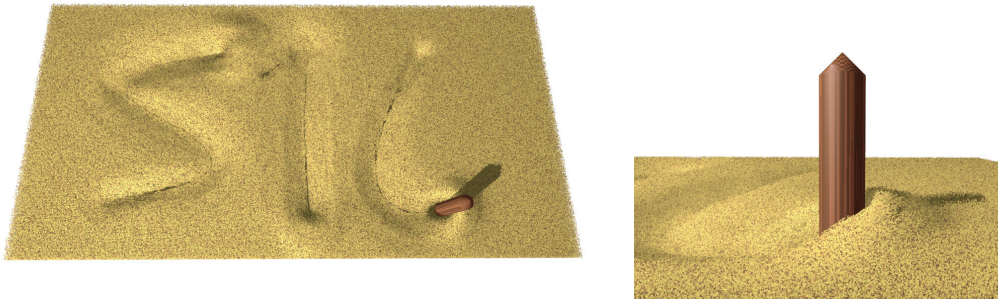
### 7.5.2 Complex scenarios

**Free-flowing material** We simulated two scenarios in which sand is manipulated through scripted rigid-body motions. The first one, reproduced in Figure 7.20, features a small cylinder being dragged across the ground. We capture the formation of a typical mound at the front of the cylinder, and the permanent marks that are due to the frictional nature of the material. In the second scenario (Figure 7.21), a handful of sand is picked up before being let to flow freely, illustrating the transitions between the gaseous, liquid and solid regimes of the material.

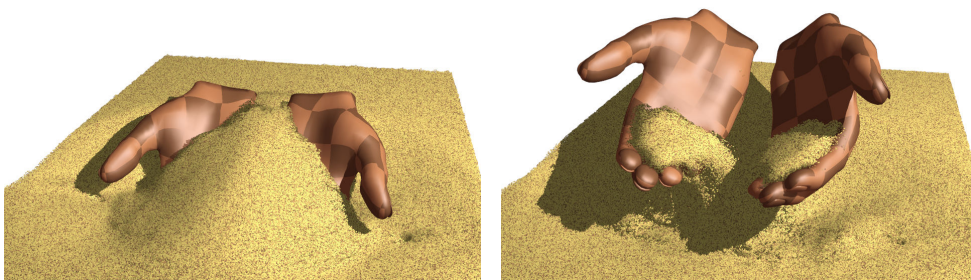
**Impact** We ran two simulations reproducing impacts of a fully-coupled rigid sphere on granular media, with and without cohesion. While Narain, Golas, et al. (2010) were able to capture the dynamics of the impact of a two-inch tungsten ball on a granular bed<sup>4</sup>, their simulation suffered from visible artifacts due to their linearization of the Drucker-Prager law. In contrast, our simulation of the same scenario yields a perfectly round crater (Figure 7.7), matching much more closely that of their reference video<sup>5</sup>.

<sup>4</sup><https://youtube.com/watch?v=ZoZ0ZAzr6eg#t=90>

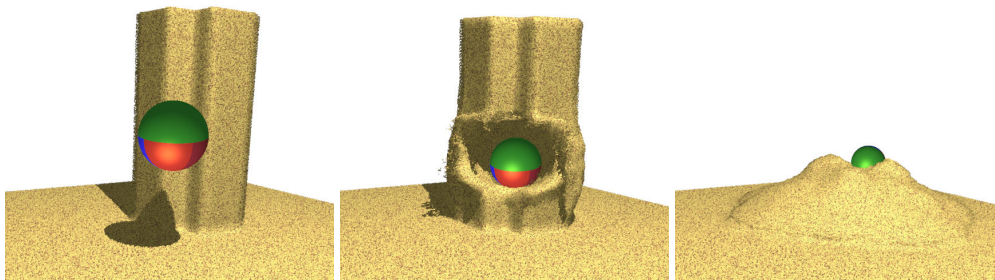
<sup>5</sup><http://dsc.discovery.com/videos/time-warp-deep-impact>



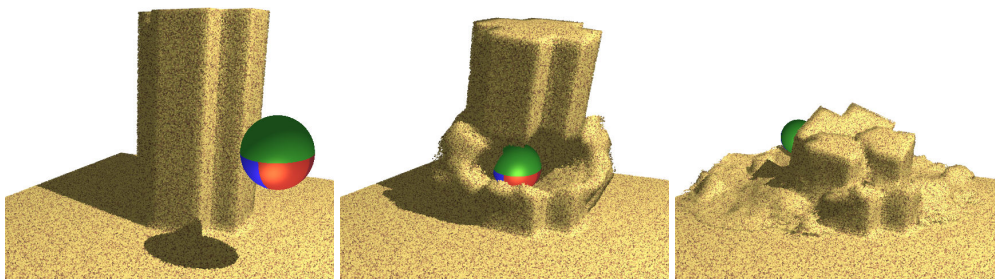
**Figure 7.20:** Letters drawn by dragging a stick in the sand. Right: a typical mound grows at the front of the stick.



**Figure 7.21:** Picking-up sand and letting it flow away.

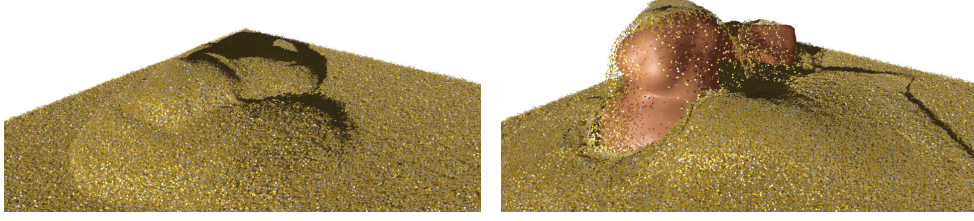


**(a)** Trilinear stress shape functions and cohesion decay



**(b)** Particle-based stress shape functions without cohesion decay

**Figure 7.22:** A sphere impacting a sand tower initially standing up thanks to a high cohesion coefficient.



**Figure 7.23:** Treasures made of gold coins are a typical example of highly anisotropic granular materials.

Example	FPS	Grid	$n_p$ <sup>1</sup>	$n_N$ <sup>1</sup>	$\bar{t}$ <sup>2</sup>	% <sub>b</sub>	% <sub>s</sub> <sup>3</sup>
<b>Collapse</b>	96	40×10×20	5.0 10 <sup>4</sup>	2755	1.47	19	72
<b>ZB05</b>	96	100×100×50	5.1 10 <sup>5</sup>	34112	23.5	21	71
<b>Wheel</b>	96	42×28×14	6.5 10 <sup>4</sup>	10152	9.62	4	92
<b>Silo</b>	96	36×36×72	1.3 10 <sup>6</sup>	49774	61.5	41	41
<b>Writing</b>	60	96×48×24	4.9 10 <sup>5</sup>	61864	48.4	18	65
<b>Cohesion</b>	240	80×80×40	3.3 10 <sup>5</sup>	42205	6.8	32	52
<b>Digging</b>	96	64×64×32	2.2 10 <sup>6</sup>	75283	32.3	23	65
<b>Crater</b>	600	75×75×50	8.3 10 <sup>5</sup>	110887	89.6	19	65
<b>Treasure</b>	96	70×70×50	3.3 10 <sup>6</sup>	121143	83.6	27	59

<sup>1</sup> Maximum numbers of particles ( $n_p$ ) and active grid nodes ( $n_N$ )

<sup>2</sup> Average simulation time in seconds per rendering frame  $\bar{t}$

<sup>3</sup> Percentage of time for building and solving the DCFP

**Table 7.1:** Sizes and simulation time of our examples

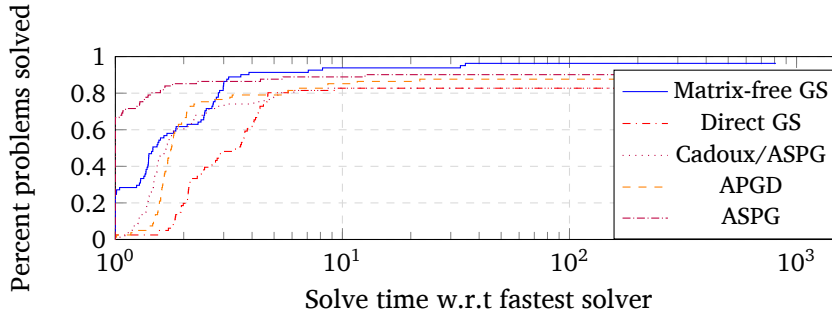
**Cohesion** As mentioned in Section 7.3.3, cohesion fits in very nicely with the continuous model, but creates tremendous discretization issues. When using trilinear shape functions for the stresses, we can only obtain a stable cohesive medium when the particles are aligned with the underlying grid. One solution to mitigate this limitation is to model debonding, for instance using a simple decay model linear in the norm of the strain rate  $|\dot{\epsilon}|$ ,  $\frac{dc}{dt} = -\xi|\dot{\epsilon}|c$ . Figure 7.22 reproduces the destruction of a granular tower initially standing thanks to cohesion using either this approach or particle-based stress shape functions. The latter option allows to retain much sharper features in the final state.

**Anisotropy** Scenes featuring large piles of gold coins are common in movies; we have simulated the Stanford bunny waking up under a large heap of coins, triggering anisotropic avalanches (see Figure 7.23).

### 7.5.3 Performance

**Simulation timings** Table 7.1 provides the sizes and simulation timings for all our examples. The target framerate (FPS) is dictated by the desired playback speed of our accompanying video — simulation frames were then subdivided using either a fixed number of substeps, or an adaptive criterion based on the fastest particle speed to grid cell size ratio. The tolerance for the DCFP solutions was set to  $10^{-3}$  times the typical stress  $\rho g L$ , and was typically reached within 25 to 250 iterations of the Gauss–Seidel algorithm warm-started with the solution of the previous timestep.

The computational bottleneck of our method lies in the DCFP solving. Note that with trilinear stress basis functions, this cost increases with the number of active grid nodes, and has very little dependence on the number of simulated particles. We can therefore use a high numbers of particles per cell, such as in the Silo simulation, without incurring too much overhead.



**Figure 7.24:** Performance profiles for a variety of solvers. The line  $x = 1$  gives the solvers that were the most often the fastest, while  $y = 1$  shows the most robust ones.

**Comparisons between Coulomb friction solvers** As mentioned previously, we use the open-source bogus library, which implements a few algorithms for solving DCFP, either tackling it directly or decomposing it as a sequence of optimization problems using the Cadoux algorithm (Section 2.3.2). In both cases, one may use our Gauss–Seidel (GS) algorithm from Chapter 4, or a few variants of the Projected-Gradient (PG) algorithm, including the APGD algorithm (Mazhar et al. 2015) and a line-search-free implementation of the ASPG method (Algorithm B.2).

Figure 7.24 shows the performance profiles for those solvers on a variety of problems extracted from our simulations — that is, the percentage of problems solved under a certain multiple of the time taken by the best-performing solver for each problem. To avoid favoring a given algorithm, we used the Alart-Curnier complementarity function to evaluate all the residuals equally. We found that the matrix-free versions of both PG and GS algorithms performed better than when explicitly assembling the Delassus operator, especially when the number of constraints (i.e., the number of stress degrees of freedom) was much larger than the velocity degrees of freedom. Note that the improved robustness of the matrix-free Gauss–Seidel shown in Figure 7.24 simply comes from the fact that for the bigger problems, the Delassus operator did not fit into main memory. We also found that using the Cadoux algorithm did not improve efficiency nor robustness on our benchmark problems. Overall, we found the matrix-free GS to perform the best on our quad-core setup and for our chosen tolerance. Note however that PG algorithms should theoretically scale more easily to a higher number of processors.

#### 7.5.4 Limitations

While our hybrid method is applicable to various challenging scenarios, it still suffers from a number of limitations. Apart from the those inherent to the use of a continuum model — which we will discuss later — the choice of the discrete function space remain the main limitation of the method, and compromises have to be made. Trilinear stress shape functions are computationally efficient, but may lead to the creation of artifacts. In particular, as mentioned in Section 7.3.3, cohesive materials have to be aligned with the underlying grid in order to remain stable, which drastically reduces their relevance. Moreover, disturbing artifacts may happen when visually disjoint clusters of particles in neighboring cells react together. Particle-based stresses mitigate these artifacts, but may lead to a much higher number of constraints and potential kinematic locking. Discontinuous trilinear stresses are a good compromise between the two approaches, but they are still expensive to solve. Moreover, enforcing the rheology to be satisfied only at a few discrete points in space may induce an overall loss of volume. A volume correction step may be applied, as in (Narain, Golas, et al. 2010), at the risk of introducing artificial energy in the system.

## 7.6 Discussion

In this chapter, we presented an extension of the numerical method of Chapter 6 that is applicable to general granular flows with varying volume fraction of grains. We saw that in order to obtain a symmetric system, the maximal volume fraction constraint has to be discretized using the Lagrangian point of view — i.e., looking at positions that are moving with the grains. The main difficulty in our approach consists in finding appropriate discretization spaces for the velocity, stress and volume fraction fields; we proposed two different choices, piecewise-constant discretization or the Material Point Method with low-order velocities and stresses. However, both of those choices suffer from limitations; finding better discretization spaces remain one of our priority for future work. Despite this, our method was still able to capture qualitative characteristics of granular flows, and to reduce the visual artifacts of previous approaches such as (Narain, Golas, et al. 2010) at a roughly similar computational cost. Moreover, in contrast to the cited approach, our numerical method derives from a proper, consistent, variational formulation rather than an ad-hoc sequence of corrective steps. This will allow us to accommodate the more complex case of diphasic flows in the next chapter.

The continuous model itself is also subject to several limitations; due to our continuum approximation of granulars, the range of possible simulations is restricted to homogeneous materials where all grains share common features, such as their size. Our method is thus not appropriate for simulating polydisperse media with various grain shapes. Moreover, being a macroscopic law, our  $\mathcal{DP}(\mu)$  rheology cannot model specific arrangements of grains, such as the formation of arches clogging the silo outlet. Our definition of a constant maximal volume fraction is also flawed; it is well known that different packings are possible over time for a given granular medium, exhibiting different volume fractions. Moreover, Roux and Radjai (1998) state that the packing volume fraction has an influence on the dilatancy angle; our model is unable to capture this phenomenon. Finally, our equations are based on an inelastic impact assumption, while a non-zero restitution coefficient might be necessary to properly model some classes of granular flows.

A last limitation, which we are going to address in the next chapter, is that we completely ignored the presence of a fluid around the grains. While the influence of air on large-enough grains is negligible, this hypothesis prevents us from tackling interesting scenarios, such as immersed avalanches.



## 8 Granular flows inside a fluid

This chapter is dedicated to the simulation of a granular material inside a Newtonian fluid. In contrast with the previous chapter, the surrounding fluid is no longer considered to be much lighter than the grains. Our main objective is to capture the qualitative influence of the surrounding fluid on the dynamics of the granular material; for instance, looking at the dynamics of an ash cloud, or at how water will affect the avalanching behavior of a column of sand. However, we still want said avalanche to eventually come to rest. For this reason, we still assume the existence of a critical volume fraction of grains  $\phi_{\max}$  such that when  $\phi$  reaches  $\phi_{\max}$ , the grains interact together through dry frictional contact, and obey the non-associated Drucker–Prager flow rule. We also suppose that no mass transfer occurs between the two phases (i.e., no chemical reaction).

Even more than in the dry case, immersed granular materials can exhibit different regimes, depending for instance on whether the inertia of the grains or that of the fluid is preponderant (Topin, Monerie, et al. 2012); we will attempt to capture this behavior in our continuum simulations. As we will see, the ratio of the particle size w.r.t. the domain will be of utmost importance. To the best of our knowledge, the method presented in this chapter is the first to combine fully-coupled two-phase equations for immersed granular flows with implicit, nonsmooth treatment of the Drucker–Prager rheology. An interesting property of our method is that it will ultimately lead to equations very similar to that of the dry case discussed in Chapter 7. The main differences will be the presence of a second velocity field (the *fluctuation*, modeling the velocity difference between the two phases), and linear constraints enforcing the conservation of the local volume of matter.

**Notations** Quantities associated with the granular phase will be denoted with a “g” subscript (e.g.,  $\mathbf{u}_g$  will be the velocity field of the grains) while those associated with the surrounding fluid will be denoted with a “f”. As in the previous chapters,  $\phi$  will denote the grain volume fraction, and we will assume that the fluid occupies the entirety of the volume where there are no grains. That is, the density of the granular phase will be  $\rho_g(\phi) = \rho_g \phi$ , and the density of the fluid phase will be  $\rho_f(\phi) = \rho_f(1 - \phi)$ . The total density of the mixture will thus be  $\rho(\phi) = \rho_f(\phi) + \rho_g(\phi)$ . We also define two velocities for the mixture: the mass-averaged velocity,  $\mathbf{u}_m$ , such that  $\rho(\phi)\mathbf{u}_m = \rho_g(\phi)\mathbf{u}_g + \rho_f(\phi)\mathbf{u}_f$ , and the volume-averaged velocity,  $\mathbf{u}_v := \phi\mathbf{u}_g + (1 - \phi)\mathbf{u}_f$ .

### 8.1 Related work

As the behavior of dry granular materials has yet to be fully understood, it is not surprising that no general model exists for immersed granular materials. However, they have been the subject of extensive studies in the second-half of the last century, in a variety of domains. Different models for different applications evolved largely in parallel, but cross-pollinated over time.

The study of the influence of dispersed grains in a Newtonian fluid goes as far as Einstein (1906), who postulated that the effective viscosity  $\eta_{\text{eff}}$  of a *slurry* — a dilute suspension of spherical particles — followed the law  $\eta_{\text{eff}}(\phi) = \eta_f(1 + \frac{5}{2}\phi)$ , where  $\eta_f$  is the viscosity of the pure Newtonian fluid, and  $\phi$  the volume fraction of particles.

Several applications have motivated subsequent work on related problems. In rough chronological order, we can mention:

- Civil engineering and soil mechanics, with the work of Terzaghi (1936) on the reduction of the effective stress in a saturated porous material; later, with the poroelasticity theory of Biot (1955);
- Sedimentation and consolidation, of practical interest for the mining industry, with the seminal work of Kynch (1952);
- Fluidized granular beds, of special importance in chemical engineering, following Jackson (1963);
- Particulate gravity currents — that is, gravity-induced flows that mostly spread in the horizontal direction, such as avalanches (Parker et al. 1986);
- Sediment transport in river beds (Hanes and Bowen 1985);
- Finally, Computer Graphics (Lenaerts and Dutré 2009; Rungtirananon et al. 2008).

To construct our numerical model, we will mainly follow the ideas of Jackson (2000). However, we propose below to walk through the different approaches that have been explored for the continuum simulation of grains in a Newtonian fluid.

### 8.1.1 Modeling

**Kynch batch flux density** Batch sedimentation study the gravity-induced settling process of initially dilute particles in a tank. Concha and Bürger (2002) relate the advances of sedimentation theory from the early works of Kynch (1952), who postulated that the evolution of the volume fraction of sediments in a horizontally homogeneous vertical tank obeyed a 1D equation,

$$\frac{\partial \phi}{\partial t} + \frac{\partial f_{\text{BK}}(\phi)}{\partial z} = 0,$$

where  $f_{\text{BK}}$  is called the Kynch *batch flux density* function. As mass conservation states that  $\frac{\partial \phi}{\partial t} = -\frac{\partial \phi u_z}{\partial z}$ , this is equivalent to saying that the vertical velocity depends only on the local volume fraction of sediments. An underlying assumption is that the volume fraction of grains is monotonically decreasing with altitude. Richardson and Zaki (1954) proposed an expression for the batch flux function, leading the vertical velocity to decrease with the volume fraction of fluid,

$$u_z(\phi) = w_\infty (1 - \phi)^\nu, \quad (8.1)$$

where  $w_\infty$  denotes the settling velocity of a single particle, and the exponent  $\nu$  is generally observed to be greater than 3. Michaels and Bolger (1962) refined this expression to ensure that the vertical velocity vanishes at the bottom of the tank, where the particles have already reached a critical volume fraction  $\phi_{\text{max}}$  and cannot be compacted anymore. They use

$$u_z(\phi) = w_\infty \left(1 - \frac{\phi}{\phi_{\text{max}}}\right)^\nu.$$

**Full two-phase models** In parallel, Jackson (1963) started the study of fluidized granular beds by modeling this phenomenon as the interactions between two continua. However, he did not include any stress preventing compression of the granular phase, and his model was found to be unstable, disagreeing with observations. Anderson and Jackson (1967) derived the governing equations of the two continua using an averaging process. In particular, they demonstrated the necessity of including a generalized buoyancy contribution in the momentum transfer term between the two phases. They modeled the drag applied by the fluid on the particles using two terms: one term opposing the relative velocity of the two phases, determined from the sedimentation experiments of Richardson and Zaki (1954); and a second term opposing the acceleration between the two phases, a so-called virtual mass. Their work has ultimately led to the popular two-phase framework presented in (Jackson 2000).

Drew and Lahey (1979) proposed a general averaging framework for the derivation of thermodynamically admissible constitutive equations, and Drew (1983) proposed closures for common two-phase flows, specifying in particular many potential contributions to the interfacial momentum transfer term. Drew (1983) also suggested the use of an (unilaterally) incompressible model for the granular phase, stating that the particulate pressure should vanish while the maximum concentration is not reached.

Batchelor (1988) argues that sediment and granular bed fluidization are two facets of the same problem, and that both should be able to be described by a single model. He attributes the difference in the two communities approaches to the fact that “Fluidized beds have mostly been studied with larger particles, such that the Reynolds number of the flow about a particle is well above unity, whereas sedimentation processes in practice usually involve a liquid continuous phase and smaller particles for which the Reynolds number is small.” Moreover, he argues that the work of Drew is hardly usable in practice; “One can go part way towards finding an equation of motion for one of the phases which formally resembles an equation for a continuum by taking an average over the volume occupied instantaneously by that phase in the manner described by Drew (1983), but his procedure achieves rigor at the cost of introducing the intractable problem of closure of averages of a complicated kind.” Instead, assuming a one-dimensional mean velocity, he proposes a model whose parameters could all be measured from physical experiments. In particular, he proposes once again to deduce the drag force from sedimentation experiments, and models the inter-particle contacts as a diffusive phenomenon, with a force proportional to the gradient of the density of grains and a supplemental viscosity.

Harris and Crighton (1994) postulates the stress in the granular phase to be the sum of a Newtonian viscosity and a particulate “pressure”  $p^C$  growing to infinity when the volume fraction reaches  $\phi_{\max}$ ,  $p^C := P_0 \phi \frac{1}{\phi - \phi_{\max}}$ . The drag force is once again modeled using the sedimentation velocity given by Richardson and Zaki (1954). Bürger (2000) proposes different models for the viscosity of each phase, and uses a virtual mass and a batch flux density function to define the drag forces. Hsu et al. (2004) propose complex expressions for the drag — based on the Richardson and Zaki (1954) model — and inter-particle pressure — diverging for a volume fraction slightly above random close packing — targeted at the modeling of sediment transport.

**Mixtures** *Mixture* theories attempt to avoid fully modeling the two phases and their interactions through the use of higher-level closures. For instance, Frankel and Acrivos (1967) and Krieger (1972) propose to extend Einstein’s law to denser concentrations of particles, making sure that the viscosity goes to infinity when the critical volume fraction of grains  $\phi_{\max}$  is reached. Frankel and Acrivos (1967) write the effective mixture viscosity

$$\eta_{\text{eff}}(\phi) = \eta_f \frac{9}{8} \left( \left( \frac{\phi_{\max}}{\phi} \right)^{\frac{1}{3}} - 1 \right)^{-1}, \quad (8.2)$$

while Krieger (1972) suggests

$$\eta_{\text{eff}}(\phi) = \eta_f \left( 1 - \frac{\phi}{\phi_{\max}} \right)^{-\frac{5}{2}\phi_{\max}}, \quad \nu > 0. \quad (8.3)$$

so that Einstein’s law is retrieved in the dilute limit. Boyer et al. (2011) empirically attempt to unify the  $\mu(I)$  rheology for dense and dry granular materials with mixture theory. They find an expression for the tangential stress which like that of Krieger (1972), diverges when  $\phi$  reaches  $\phi_{\max}$  and coincides with Einstein’s viscosity in the dilute limit. The pressure-to-frictional-stress ratio is found to follow a  $\mu(I_v)$  relationship, where  $I_v$  is called the *viscous number*.

The mixture theory developed for gravity currents use different kinds of closures. A large portion of natural gravity currents involves phases with similar densities; in this case, many works rely on the Boussinesq approximation, which consists in considering the varying density solely in the buoyancy term (through the introduction of a “reduced” gravity). Ungarish (2009) considers this approximation to remain valid for density differences of up to 10% or



20%. However, neither granular avalanches in water, nor powder avalanches in the air, fall into this range. For higher density ratios, mixture theory (see Manninen et al. 1996 for a high-level introduction, or Bedford and Drumheller 1983 for a derivation from thermodynamical and averaging considerations) propose an alternative approximation of the multiphase continuum by considering closure equations for the diffusion velocity of the different phases, for instance following Fick's law (e.g., Etienne 2004). Recalling that  $\mathbf{u}_m$  denotes the mass-averaged velocity,  $\mathbf{u}_m := \frac{1}{\rho(\phi)} (\rho_g(\phi)\mathbf{u}_g + \rho_f(\phi)\mathbf{u}_f)$ , the flux density  $\mathbf{q}$  for each phase is defined such that

$$\begin{aligned} \frac{D_{\mathbf{u}_m} \phi}{D_{\mathbf{u}_m} t} + \phi \nabla \cdot \mathbf{u}_m &= -\nabla \cdot [\mathbf{q}_g] \\ \frac{D_{\mathbf{u}_m} (1 - \phi)}{D_{\mathbf{u}_m} t} + (1 - \phi) \nabla \cdot \mathbf{u}_m &= -\nabla \cdot [\mathbf{q}_f]. \end{aligned}$$

Expressing those flux densities using mass conservation would involve the velocities of individual phases velocities;  $\mathbf{q}_g = (\mathbf{u}_g - \mathbf{u}_m)\phi$  and  $\mathbf{q}_f = (\mathbf{u}_f - \mathbf{u}_m)(1 - \phi)$ . Fick's law instead postulates that the flux densities are proportional to the gradient of the volume fraction:  $\mathbf{q}_g := q_g(\phi)\nabla\phi$ ,  $\mathbf{q}_f := q_f(\phi)\nabla\phi$ . Note that the rates  $q_f$  and  $q_g$  cannot be chosen independently, as total mass conservation requires  $\rho_f q_f + \rho_g q_g = 0$ . Remark that Fick's law amount to considering that the drift velocities  $\mathbf{u}_g - \mathbf{u}_m$  and  $\mathbf{u}_f - \mathbf{u}_m$  depend only on the volume fraction, which is reminiscent of Kynch's theory in a different, higher-dimensional setting. The evolution of the volume fraction field can then be expressed as an unsteady convection-diffusion equation,

$$\frac{D_{\mathbf{u}_m} \phi}{D_{\mathbf{u}_m} t} + \phi \nabla \cdot \mathbf{u}_m + \nabla \cdot [q_f(\phi)\nabla\phi] = 0.$$

However, this diffusive model is not satisfying for our purposes, as we want to be able to simulate not only gravity currents but also sedimentation processes (such as sand falling at the bottom of a tank). As such, and since batch sedimentation theory is mostly concerned with 1D problems and is thus also too restrictive, our numerical model will be based on full two-phase equations.

### 8.1.2 Numerical simulations

These different modeling approaches have led to a variety of numerical methods for the simulation of immersed granular or powder flows. We mention below a few techniques that are representative of the diversity of strategies employed for this purpose, but the reader should keep in mind that this list is far from exhaustive.

**Sedimentation and fluidization** Andrews and O'Rourke (1996) extended the Particle-In-Cell hybrid method to the simulation of multiphase flow, using the particle interaction stress from (Harris and Crighton 1994). Apte et al. (2003) used a similar approach with a slightly different expression for the particulate phase pressure, using a non-unit exponent for the hyperbolic divergence of the stress when  $\phi$  approaches  $\phi_{\max}$ . Chauchat and Médale (2010) model sediment transport in a laminar river bed by decomposing the flow into two (fixed) domains: the upper one where a mixture is considered, and the bottom one where the grains are assumed to be dense and a full two-phases continuum is considered. In the upper domain (the river), Einstein's viscosity for dilute suspensions is considered. The lower domain (the granular bed) is modeled following Jackson (2000), and the grains are subject to a pressure-dependent yield-stress. The volume fraction of grains is assumed to remain constant, and thus all flows are incompressible. The equations are then solved using FEM with regularization of the yield stress. Revil-Baudard and Chauchat (2013) study the transport of sediments by a turbulent fluid, using again a mixture approach for the upper, turbulent layer and a two-phase model for the bottom one. This granular layer is modeled using the  $\mu(I_v)$  rheology from Boyer et al. (2011), and solved numerically

by regularizing of the yield stress. Chauchat, Guillou, et al. (2013) simulate 1D sedimentation processes by starting from the two-phase equations of Jackson (2000) with Frankel–Acrivos viscosity (8.2), and providing closures for the drag force and effective granular stress for cohesive and non-cohesive particles. They found their results to be in agreement with experiments in both cases.

**Gravity currents** Mixture theory has been especially popular for the numerical simulation of avalanches; Etienne (2004) and Birman et al. (2005) apply this approach to lock-exchange problems with high density ratios. Iverson (1997), and Denlinger and Iverson (2001) extend the depth-averaged Savage–Hutter model for dry granular flows to accommodate the presence of an interstitial fluid. Pitman and Le (2005) propose a new depth-averaged model by building upon the two-phase equations of Jackson (2000), and demonstrate how the pore pressure induced by the fluid can reduce basal friction, and increase the run-out length of the avalanche. Bouchut, Fernández-Nieto, et al. (2016) construct a two-phase two-layers shallow model taking into account dilatancy effects with the model from Roux and Radjai (1998).

The granular collapse simulations by Lagrée et al. (2011) were also done using a multiphase solver, Gerris<sup>1</sup>; however, the granular phase was considered as fully incompressible, while here we will attempt to model unilateral incompressibility.

**Computer graphics** Müller et al. (2005) proposed a first approach for simulating interacting fluids within the SPH framework, but considered a single velocity field for the different phases. H. Zhu et al. (2006) presented simulations of general mixtures using a Lattice-Boltzmann approach. L. Boyd and Bridson (2012) introduced a FLIP extension for multiphase immiscible flows able to take into account a distinct velocity for each phase. Nielsen and Østerby (2013) simulated the coupled motion of droplets of water in air, with a convection-diffusion process for the water phase. Ren et al. (2014) developed an extension of SPH targeted at mixtures.

More relevant to our applications, Lenaerts and Dutré (2009) demonstrated a first approach for the simulation of the interactions between water and a granular continuum. To do so, they coupled a granular SPH model with the framework of (Lenaerts, Adams, et al. 2008), which simulates the interaction of a fluid with porous materials by adaptively absorbing and recreating fluid particles, with a diffusive model for the transport of the fluid inside the porous material. Their approach was able to capture the transition of a granular material into mud, but cannot reproduce sedimentation processes.

**Discrete simulations** Finally, while our goal is to simulate systems with a large enough number of grains to motivate the continuum approach, we can mention that several works have focused on coupled fluid–granular simulations using diverse discrete models for the particles.

Choi and Joseph (2001) performed direct numerical simulations of the fluidization of a few hundred of particles, adapting the simulation mesh to their boundaries. Pignatelli et al. (2009) simulated the low-velocity, gravity-driven injection of particles inside a Newtonian fluid using particle sampling and compared their results to both experiments and a continuum model using the effective viscosity from Krieger (1972) to model the stress in the particulate phase. Topin, Dubois, et al. (2011) and Topin, Monerie, et al. (2012) couple the Non-Smooth Contact Dynamics method for the granular phase with a Newtonian fluid simulation using the fictitious domain approach (that is, extending the fluid domain to the interior of the particles and enforcing rigid motion there). Maurin et al. (2015) used an hybrid approach, using mixture theory for the surrounding fluid yet computing the interactions of the particles with a molecular dynamics approach similar to that of (Cundall and Strack 1979).

**Our approach** Despite the variety of methods mentioned above, none of them proposes a fully coupled continua simulation while simultaneously enforcing the non-associated Drucker–Prager

<sup>1</sup><http://gfs.sourceforge.net>

flow rule for the granular phase, without regularization of the yield stress. Yet, such a nonsmooth treatment is necessary to avoid the lingering of residual velocities (at the end of an avalanche, for instance). To achieve this goal, we will extend the framework of Chapter 7 to account for the Newtonian fluid. More precisely, we will make use of the two-phase model of Anderson and Jackson (1967) with the drag expression of Richardson and Zaki (1954). To construct our system of dimensionless equations, we will also take inspiration from Boyer (2001). Indeed, while this latter work deals with the simulation of immiscible fluids with surface tension, it uses a two-phase approach that is relevant for us.

## 8.2 Two-phase model

In this section, we shall propose equations for the motion of a granular material subject to a Drucker–Prager or  $\mu(I)$  yield stress inside a Newtonian fluid, following the approach of Anderson and Jackson (1967) and Jackson (2000). Our overview of the modeling literature has shown that the expression of the stresses acting on the two phases is not fully known, and that several alternatives have been proposed even for the study of similar phenomena. Here we will attempt to make reasonable choices, in the hope that they will prove relevant enough to capture the qualitative behavior of the materials, but are fully aware that comparisons with experiments remain necessary in order to determine definitive expressions for those stresses.

### 8.2.1 Base equations

We start from the momentum and mass conservation equations for the granular and fluid continua. We assume that the sole external force applied our diphasic medium is gravity.

$$\rho_g \phi \frac{D_{\mathbf{u}_g} \mathbf{u}_g}{D_{\mathbf{u}_g} t} - \nabla \cdot [\boldsymbol{\sigma}_g] = \rho_g \phi \mathbf{g} + \mathbf{f}_{f \rightarrow g} \quad (8.4)$$

$$\rho_f (1 - \phi) \frac{D_{\mathbf{u}_f} \mathbf{u}_f}{D_{\mathbf{u}_f} t} - \nabla \cdot [\boldsymbol{\sigma}_f] = \rho_f (1 - \phi) \mathbf{g} - \mathbf{f}_{f \rightarrow g} \quad (8.5)$$

$$\frac{\partial \phi}{\partial t} + \nabla \cdot [\phi \mathbf{u}_g] = 0 \quad (8.6)$$

$$\frac{\partial (1 - \phi)}{\partial t} + \nabla \cdot [(1 - \phi) \mathbf{u}_f] = 0 \quad (8.7)$$

Expressions for the stresses  $\boldsymbol{\sigma}_{f,g}$  and interfacial momentum transfer  $\mathbf{f}_{f \rightarrow g}$  term remain to be written. As we have seen in our brief literature review, these have been subject to thorough debate in the last decades. Moreover, different physical considerations may lead to attribute certain contributions to one term or the other, but still lead to the same equations. Here we shall follow Anderson and Jackson (1967), and state that the interfacial momentum transfer should consist of a drag term  $\mathbf{f}_{f \rightarrow g}^d$  and a generalized buoyancy contribution  $\mathbf{f}_{f \rightarrow g}^b$ ,

$$\mathbf{f}_{f \rightarrow g} := \mathbf{f}_{f \rightarrow g}^d + \mathbf{f}_{f \rightarrow g}^b.$$

Stresses will include a pressure term, a contribution from the Newtonian mixture viscosity, and the stress induced by the contact between particles.

### 8.2.2 Stresses and buoyancy

We assume the following expression for the stresses:

$$\begin{aligned} \boldsymbol{\sigma}_g &= -p_g \mathbb{I} + \boldsymbol{\sigma}_g^v + \phi \boldsymbol{\sigma}^C \\ \boldsymbol{\sigma}_f &= -p_f \mathbb{I} + \boldsymbol{\sigma}_f^v. \end{aligned}$$

$p_{f,g}$  corresponds to the partial pressure in each phase,  $\sigma_{f,g}^v$  to Newtonian viscosities, and  $\sigma^C$  to the supplemental stress due to the contacts between grains (the effective stress, in Terzaghi's terms).

Our first assumption, which follows (Chauchat, Guillou, et al. 2013; Einstein 1906; Jackson 2000), will be that the Newtonian viscous forces should be relative to the volume-average velocity  $\mathbf{u}_v$ . Moreover, the viscosity shall depend only on the local volume fraction, i.e.,  $\sigma_{f,g} = \eta_{f,g}(\phi)D(\mathbf{u}_v)$ .

**Pore stress** We first consider the stress conveyed at the interface between the two phases; this is the case for the pressure  $p_{f,g}$  and viscous stresses  $\sigma_{f,g}^v$ . Following Anderson and Jackson (1967), these stresses generate a generalized buoyancy force density,

$$\mathbf{f}_{f \rightarrow g}^b = -\phi \nabla \cdot [\sigma_f^b] + (1 - \phi) \nabla \cdot [\sigma_g^b],$$

where  $\sigma_{f,g}^b := -p_{f,g}\mathbb{I} + \sigma_{f,g}^v$ . Note that this is consistent with, and a generalization of, the classical Archimede's principle.

Now, as mentioned by B rger (2000),  $\sigma_{f,g}^b$  are mathematical quantities that cannot be experimentally measured, in contrast to the *pore* (or interstitial) stress  $\sigma^b$ , which can. Indeed, following from Terzaghi's principle,  $\sigma^b$  can be computed as the difference of the total stress and the effective stress due to contacts, i.e.,  $\sigma^b = (\sigma_g + \sigma_f) - \phi \sigma^C$ . When the material is at rest,  $\sigma^b$  simply corresponds to the hydrostatic pressure. We can make two kinds of hypothesis to relate  $\sigma_{f,g}^b$  to  $\sigma^b$ , which, when taking into account buoyancy, will ultimately lead to the same set of equations.

- We can suppose that the viscosity and pressure belong fully to the fluid phase, as done by e.g., Pitman and Le (2005). This means  $\sigma_g^b = \mathbf{0}$ , and  $\sigma_f^b = \sigma^b$ . In this case, the generalized buoyancy is simply  $\mathbf{f}_{f \rightarrow g}^b = -\phi \nabla \cdot [\sigma^b]$ .
- Alternatively, we can suppose that the phase stresses  $\sigma_{f,g}^b$  are given by pore stress  $\sigma_{f,g}^b$  scaled with their relative volume fractions, that is  $\sigma_g^b = \phi \sigma^b$  and  $\sigma_f^b = (1 - \phi) \sigma^b$ . Then the generalized buoyancy is  $\mathbf{f}_{f \rightarrow g}^b = -\phi \nabla \cdot [(1 - \phi) \sigma^b] + (1 - \phi) \nabla \cdot [\phi \sigma^b] = \nabla \cdot [\phi \sigma^b] - \phi \nabla \cdot [\sigma^b] = \sigma^v \nabla \phi$ .

In both cases, the pore stress can be decomposed as  $\sigma^b = p\mathbb{I} - \eta_{\text{eff}}(\phi)D(\mathbf{u}_v)$ , where  $p$  is called the pore pressure, and we get:

$$\begin{aligned} \nabla \cdot [\sigma_g^b] - \mathbf{f}_{f \rightarrow g}^b &= \phi \sigma^b = \phi (\eta_{\text{eff}}(\phi)D(\mathbf{u}_v) - p\mathbb{I}) \\ \nabla \cdot [\sigma_f^b] + \mathbf{f}_{f \rightarrow g}^b &= (1 - \phi) \sigma^b = (1 - \phi) (\eta_{\text{eff}}(\phi)D(\mathbf{u}_v) - p\mathbb{I}). \end{aligned}$$

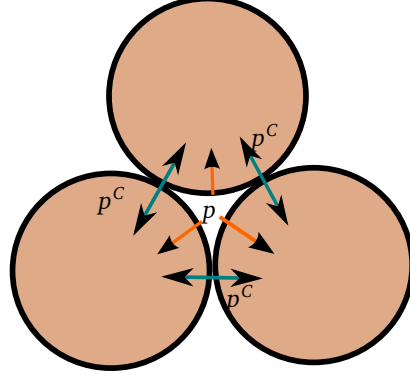
Equations (8.4–8.5) can thus be rewritten as

$$\rho_g \phi \frac{D_{\mathbf{u}_g} \mathbf{u}_g}{D_{\mathbf{u}_g} t} - \phi \nabla \cdot [\eta_{\text{eff}}(\phi)D(\mathbf{u}_v) - p\mathbb{I}] - \nabla \cdot [\phi \sigma^C] = \rho_g \phi \mathbf{g} + \mathbf{f}_{f \rightarrow g}^d \quad (8.8)$$

$$\rho_f (1 - \phi) \frac{D_{\mathbf{u}_f} \mathbf{u}_f}{D_{\mathbf{u}_f} t} - (1 - \phi) \nabla \cdot [\eta_{\text{eff}}(\phi)D(\mathbf{u}_v) - p\mathbb{I}] = \rho_f (1 - \phi) \mathbf{g} - \mathbf{f}_{f \rightarrow g}^d \quad (8.9)$$

It remains to choose an expression for  $\eta_{\text{eff}}(\phi)$ . We can for instance suppose  $\eta_{\text{eff}}$  to be constant, or use Einstein's expression,  $\eta_{\text{eff}}(\phi) = \eta_f \left(1 + \frac{5}{2}\phi\right)$ .

**Contact stress** Our second main assumption is that the contacts between grains are subject to dry friction. Just like in the dry case of Chapter 7, the grains will be considered to be self-contacting, with no layer of fluid between them, as soon as the maximal volume fraction is reached.



**Figure 8.1:** Two kinds of competing normal stresses acting on the particles: the pore pressure  $p$  and the effective stress due to contacts  $p^C = -\frac{1}{d} \text{Tr } \sigma^C$

As in the dry case, we will want the grains volume fraction  $\phi$ , contact stress  $\sigma^C$ , and strain rate  $\dot{\epsilon}_g$ , to follow the maximal volume fraction constraint with non-associated Drucker–Prager flow rule, i.e.,

$$\begin{cases} 0 \leq p^C \perp \phi_{\max} - \phi \geq 0 \\ \text{Dev } \sigma^C = \hat{\mu} p^C \frac{\text{Dev } \dot{\epsilon}_g}{|\text{Dev } \dot{\epsilon}_g|} & \text{if } \text{Dev } \dot{\epsilon}_g \neq \mathbf{0} \text{ (yielded)} \\ |\text{Dev } \sigma^C| \leq \hat{\mu} p^C & \text{if } \text{Dev } \dot{\epsilon}_g = \mathbf{0} \text{ (unyielded),} \end{cases} \quad (8.10)$$

where  $p^C := -\frac{1}{d} \text{Tr } \sigma^C$ .

This solid-contacts assumption also means that the stress due to contacts between grains will be propagated at grain–grain interfaces only, and exist solely inside the granular phase. As such, there is no associated buoyancy force. This is confirmed by a simple thought experiment: consider a dry, compact heap of rigid grains at rest, supporting a given mass  $m$ . Modifying the intensity of the contact forces between the grains, for instance by changing  $m$ , will not affect the air phase nor the interstitial pressure  $p$ .

**Remark 8.1.** We can already remark that as illustrated on Figure 8.1, the normal component of the total stress of the mixture will be the sum of two terms,  $p$  and  $p^C = -\frac{1}{d} \text{Tr } \sigma^C$ . Now, as per the Drucker–Prager criterion (8.10), increasing  $p^C$  increases the yield stress of the material; for a given applied load, the pore pressure  $p$  will thus have a weakening effect.

### 8.2.3 Drag force

Our third core hypothesis will be to consider grains that are thin enough, or with velocity (w.r.t. the surrounding fluid) small enough that Stokes’ drag law apply. That is, we consider the case where the grain-level Reynolds number,  $\text{Re}_g := \frac{\rho_f W D_g}{\eta_f}$ , where  $D_g$  is the diameter of the grains and  $W$  a characteristic relative velocity, is small. This choice has been made by numerous authors, e.g., (Anderson and Jackson 1967; Bürger 2000; Chauchat and Médale 2010; Harris and Crighton 1994; Pitman and Le 2005), while others consider the drag force to be quadratic in the relative velocity (Batchelor 1988; Chauchat, Guillou, et al. 2013; Drew 1983). We will stick with simple linear drag, knowing that in any case we can account for the nonlinearity by explicitly adjusting the drag coefficient in the numerical method, and remember that we may overestimate the relative velocity in the analysis below.

We will thus express the drag force as  $\mathbf{f}_{f \rightarrow g} = -\hat{\xi}(\phi) \mathbf{w}$ , where  $\mathbf{w}$  is the relative velocity between the two phases,  $\mathbf{w} := \mathbf{u}_g - \mathbf{u}_f$ , and  $\hat{\xi}$  is a scalar function of the volume fraction. In the

two-phase literature,  $\hat{\xi}(\phi)$  is commonly deduced from sedimentation experiments, for instance using the empirical law of Richardson and Zaki (1954), such as in (Chauchat, Guillou, et al. 2013). The quasistatic conservation of momentum for the sedimenting particles is written as

$$0 = \phi(\rho_g - \rho_f)g - \hat{\xi}(\phi)u_z(\phi),$$

yielding

$$\hat{\xi}(\phi) = \frac{\phi(\rho_g - \rho_f)g}{u_z(\phi)}.$$

Equation (8.1) gives  $u_z(\phi) = w_\infty(1 - \phi)^\nu$ ; let us compute the value of  $w_\infty$  under the hypothesis  $Re_g \leq 1$ . Let  $w$  denote the relative velocity of a single spherical grain falling through the fluid. Stokes' law state that the intensity of the drag on the particle is  $-3\pi\eta_f D_g w$ , meaning that the equilibrium is reached for  $(\rho_g - \rho_f)(\frac{4}{3}\pi(\frac{D_g}{2})^3)g = 3\pi\eta_f D_g w$ , i.e.,

$$w_\infty = \frac{(\rho_g - \rho_f)g D_g^2}{18\eta_f}.$$

Overall, we conclude

$$\hat{\xi}(\phi) = \eta_f \frac{18}{D_g^2} \phi(1 - \phi)^{-\nu}.$$

Note that this expression increases monotonically with the volume fraction of grains, which at first may appear strange, as one expects the average drag force to decrease for denser objects. However, this weakening of the drag force is due to a smaller relative velocity, as the fluid is lugged by the higher density of grains, rather than a decrease in the friction coefficient. Note also that this expression makes  $\hat{\xi}$  grow to infinity when the fluid vanishes, thus requiring the relative velocity  $\mathbf{w}$  to vanish for purely granular materials. As our maximal volume fraction constraint requires  $\phi \leq \phi_{\max} < 1$ , we need not worry about this degenerate case, however. Other closures are possible; for instance, Chauchat and Médale (2010) use the so-called *Carman–Kozeny* relationship. For now, we will not assume a precise expression for  $\hat{\xi}$ , but simply write that

$$\mathbf{f}_{f \rightarrow g}^d = -\phi(1 - \phi)\hat{\xi}(\phi)\mathbf{w}, \quad (8.11)$$

and remember that  $\hat{\xi}(\phi)$  in  $\text{Pa.s.m}^{-2}$  is similar in order of magnitude to  $g \frac{\rho_g - \rho_f}{w_\infty(1 - \phi)^{\nu+1}} = \frac{\eta_f}{D_g^2} \frac{18}{(1 - \phi)^{\nu+1}}$ .

**Summary of our hypothesis** At this point, we have made the following hypothesis:

- no mass is transferred between the two phases;
- the Newtonian viscous stress is proportional to  $D(\mathbf{u}_v)$ ;
- the linear Stokes drag law applies;
- the contacts between grains are subject to dry friction and occur when  $\phi = \phi_{\max} < 1$ .

#### 8.2.4 Mixture conservation equations

The conservation equations (8.4–8.7) written with the phase velocities are not very convenient; Boyer (2001) propose to express them instead as functions of two other velocity variables, the volume-average velocity  $\mathbf{u}_v$  and the relative velocity  $\mathbf{w}$ . This choice does have some nice properties; the Newtonian viscosity is already expressed as a function of  $\mathbf{u}_v$ , and  $\mathbf{u}_v$  is the only velocity variable that is divergence-free. Indeed, summing Equations (8.6) and (8.7), we get

$$\nabla \cdot \mathbf{u}_v = -\frac{\partial 1}{\partial t} = 0. \quad (8.12)$$

However, using  $\mathbf{u}_v$  also requires Boyer (2001) to deal with rather inconvenient inertial terms. Here, we will choose instead to use  $\mathbf{u}_m$  and  $\mathbf{w}$  as our velocity variables.

**Notations** Following Etienne (2004), we introduce the scaled density difference

$$\alpha := \frac{\rho_g - \rho_f}{\rho_f}.$$

Note that  $\rho_g/\rho_f = (\alpha + 1)$ . Let  $\beta(\phi) := (1 + \alpha\phi)$ , so that the total density of the mixture is given by  $\rho(\phi) = \phi\rho_g + (1 - \phi)\rho_f = \beta(\phi)\rho_f$ . Let  $\pi(\phi) := \phi(1 - \phi)$ .

The total mass conservation of the mixture, obtained by summing (8.6) scaled by  $\rho_g$  and (8.7) scaled by  $\rho_f$ , reads

$$\frac{\partial \beta}{\partial t} + \nabla \cdot [\beta \mathbf{u}_m] = 0. \quad (8.13)$$

Finally, note that phase and volume-averaged velocities can also be expressed as function of the mass-averaged and relative velocities:

$$\begin{aligned} \mathbf{u}_g &= \mathbf{u}_m + \frac{(1 - \phi)}{\beta} \mathbf{w} & \mathbf{u}_f &= \mathbf{u}_m - (\alpha + 1) \frac{\phi}{\beta} \mathbf{w} \\ \mathbf{u}_v &= \mathbf{u}_m - \alpha \frac{\pi}{\beta} \mathbf{w} \end{aligned}$$

**Mixture momentum conservation** To obtain the conservation equation for the total momentum of the mixture, we just have to sum Equations (8.8) and (8.9). It is simpler to consider the conservative expression<sup>2</sup> for the inertial terms,

$$\begin{aligned} \phi \frac{D_{\mathbf{u}_g} \mathbf{u}_g}{D_{\mathbf{u}_g} t} &= \frac{\partial \phi \mathbf{u}_g}{\partial t} + \nabla \cdot [\phi \mathbf{u}_g \otimes \mathbf{u}_g] \\ (1 - \phi) \frac{D_{\mathbf{u}_f} \mathbf{u}_f}{D_{\mathbf{u}_f} t} &= \frac{\partial (1 - \phi) \mathbf{u}_f}{\partial t} + \nabla \cdot [(1 - \phi) \mathbf{u}_f \otimes \mathbf{u}_f]. \end{aligned}$$

Remark first that

$$\begin{aligned} &\rho_g \phi \mathbf{u}_g \otimes \mathbf{u}_g + \rho_f (1 - \phi) \mathbf{u}_f \otimes \mathbf{u}_f \\ &= \rho_f [(\alpha + 1) \phi \mathbf{u}_g \otimes \mathbf{u}_m + (1 - \phi) \mathbf{u}_f \otimes \mathbf{u}_m] \\ &\quad + \rho_f [(\alpha + 1) \phi \mathbf{u}_g \otimes (\mathbf{u}_g - \mathbf{u}_m) + (1 - \phi) \mathbf{u}_f \otimes (\mathbf{u}_f - \mathbf{u}_m)] \\ &= \rho_f \left[ \beta \mathbf{u}_m \otimes \mathbf{u}_m + (\alpha + 1) \frac{\phi(1 - \phi)}{\beta} \mathbf{u}_g \otimes \mathbf{w} - (\alpha + 1) \frac{\phi(1 - \phi)}{\beta} \mathbf{u}_f \otimes \mathbf{w} \right] \\ &= \rho_f \left[ \beta \mathbf{u}_m \otimes \mathbf{u}_m + (\alpha + 1) \frac{\pi}{\beta} \mathbf{w} \otimes \mathbf{w} \right]. \end{aligned}$$

Therefore,

$$\begin{aligned} &\rho_g \phi \frac{D_{\mathbf{u}_g} \mathbf{u}_g}{D_{\mathbf{u}_g} t} + \rho_f (1 - \phi) \frac{D_{\mathbf{u}_f} \mathbf{u}_f}{D_{\mathbf{u}_f} t} \\ &= \rho_f \left[ \frac{\partial \beta \mathbf{u}_m}{\partial t} + (\beta \mathbf{u}_m \cdot \nabla) \mathbf{u}_m + \underbrace{\left\langle \nabla \cdot [\beta \mathbf{u}_m], \mathbf{u}_m \right\rangle}_{-\frac{\partial \beta}{\partial t}} + (\alpha + 1) \nabla \cdot \left[ \frac{\pi}{\beta} \mathbf{w} \otimes \mathbf{w} \right] \right] \\ &= \rho_f \left[ \beta \frac{D_{\mathbf{u}_m} \mathbf{u}_m}{D_{\mathbf{u}_m} t} + (\alpha + 1) \nabla \cdot \left[ \frac{\pi}{\beta} \mathbf{w} \otimes \mathbf{w} \right] \right]. \end{aligned}$$

The conservation equation for the total momentum of the mixture thus reads

$$\begin{aligned} &\rho_f \left[ \beta \frac{D_{\mathbf{u}_m} \mathbf{u}_m}{D_{\mathbf{u}_m} t} + (\alpha + 1) \nabla \cdot \left[ \frac{\pi}{\beta} \mathbf{w} \otimes \mathbf{w} \right] \right] + \nabla p - \nabla \cdot \left[ \eta_{\text{eff}} D(\mathbf{u}_m - \alpha \frac{\pi}{\beta} \mathbf{w}) \right] - \nabla \cdot [\phi \boldsymbol{\sigma}^C] \\ &= \rho_f \beta \mathbf{g} \quad (8.14) \end{aligned}$$

<sup>2</sup>see Section 5.2.1

**Fluctuation momentum conservation** Another interesting linear combination of the momentum conservation consists in taking the difference of (8.8) scaled by  $\rho_f(1-\phi)$ , and (8.9) scaled by  $\phi\rho_g$ . Processing term by term, we have

- Inertial terms:

$$\begin{aligned}\rho_g\rho_f\pi\frac{D_{\mathbf{u}_g}\mathbf{u}_g}{D_{\mathbf{u}_g}t} - \rho_g\rho_f\pi\frac{D_{\mathbf{u}_f}\mathbf{u}_f}{D_{\mathbf{u}_f}t} &= \rho_g\rho_f\pi\left[\frac{\partial\mathbf{w}}{\partial t} + (\mathbf{u}_g \cdot \nabla)\mathbf{u}_g - (\mathbf{u}_f \cdot \nabla)\mathbf{u}_f\right] \\ &= \rho_g\rho_f\pi\left[\frac{\partial\mathbf{w}}{\partial t} + (\langle\mathbf{u}_{f,g}\rangle \cdot \nabla)\mathbf{w} - (\mathbf{w} \cdot \nabla)\langle\mathbf{u}_{f,g}\rangle\right], \\ \text{with } \langle\mathbf{u}_{f,g}\rangle &:= \frac{\mathbf{u}_f + \mathbf{u}_g}{2} = \mathbf{u}_m + \frac{1-2\phi-\alpha\phi}{2\beta}\mathbf{w};\end{aligned}$$

- Mixture stress:

$$-\rho_f\pi\nabla \cdot [\boldsymbol{\sigma}^b] + \rho_g\pi\nabla \cdot [\boldsymbol{\sigma}^b] = \rho_f\alpha\pi\nabla \cdot [\boldsymbol{\sigma}^b];$$

- Drag force:

$$\rho_f(1-\phi)f_{f \rightarrow g}^d + \rho_g\phi f_{f \rightarrow g}^d = \rho_f\beta f_{f \rightarrow g}^d = -\rho_f\beta\pi\xi\mathbf{w};$$

- Gravity:

$$\rho_f(1-\phi)\phi\rho_g\mathbf{g} - \rho_g\phi\rho_f(1-\phi)\mathbf{g} = \mathbf{0}.$$

The conservation of the fluctuation momentum, where we have divided each term by  $(1-\phi)\rho_f$ <sup>3</sup>, thus reads

$$\begin{aligned}(\alpha+1)\rho_f\phi\left[\frac{\partial\mathbf{w}}{\partial t} + (\langle\mathbf{u}_{f,g}\rangle \cdot \nabla)\mathbf{w} + (\mathbf{w} \cdot \nabla)\langle\mathbf{u}_{f,g}\rangle\right] + \beta\phi\xi\mathbf{w} \\ - \alpha\phi\nabla p + \alpha\phi\nabla \cdot \left[\eta_{\text{eff}}D(\mathbf{u}_m - \alpha\frac{\pi}{\beta}\mathbf{w})\right] - \nabla \cdot [\phi\boldsymbol{\sigma}^c] = \mathbf{0}. \quad (8.15)\end{aligned}$$

**Remark** The velocity transport terms in Equations (8.14) and (8.15) seem relatively unwieldy, and, as we will see in the next section, may not be negligible for certain applications. However, their expression will be drastically simplified when using a timestepping algorithm with an explicit velocity-advection scheme, such as the characteristics method or particle-based transport.

### 8.2.5 Dimensionless equations

**Characteristic quantities** To get an idea of the respective importance of the different terms in Equation (8.14) and (8.15), we need to define characteristic values for our variables. There are several ways to do this, and none will hold for the whole range of regimes that can be exhibited by our two-phase material. As we will mostly be studying gravity-driven flows, we define the characteristic mixture velocity as  $U := \sqrt{gL}$ , where  $L$  is a characteristic length of the studied phenomenon. For the same reason, we choose  $P := \rho_f gL$  as the characteristic pore pressure,  $(\alpha+1)P = \rho_g gL$  as the characteristic contact stress,  $T = \frac{L}{U}$  as the characteristic time. Finally, at the risk of overestimating  $W$ , we define the characteristic relative velocity as the sedimentation velocity of a single grain,  $w_\infty$ .

The dimensionless variables will be denoted with a tilde in the equations below, but we will quickly drop their decoration to lighten notation. We also introduce dimensionless versions of the effective viscosity and drag fields,

$$\eta_{\text{eff}} = \eta_f\tilde{\eta}_{\text{eff}} \quad \xi = g\frac{\rho_g - \rho_f}{w_\infty}\Xi\tilde{\eta}_{\text{eff}},$$

where the dimensionless number  $\Xi$  denote the order of magnitude of the geometry dependent term in  $\xi$ , i.e.,  $\Xi \sim \xi\frac{w_\infty}{\rho_g - \rho_f}g \sim (1-\phi)^{-\nu-1}$  when using the Richardson and Zaki (1954) closure. Note that as  $\nu$  is usually taken to be greater than 3, the choice of  $\Xi$  will be highly dependent of the target volume fraction.

<sup>3</sup>We supposed that  $\phi_{\text{max}} < 1$  and thus  $\phi = 1$  is never reached



**Mixture momentum conservation** The dimensionless version of (8.14) is, after division by  $\rho_f$ ,

$$\frac{1}{L} \left[ \beta U^2 \frac{D\tilde{\mathbf{u}}_m}{D\tilde{t}} + (\alpha + 1) U W \tilde{\nabla} \cdot \left[ \frac{\pi}{\beta} \tilde{\mathbf{w}} \otimes \tilde{\mathbf{w}} \right] \right] - (\alpha + 1) g \tilde{\nabla} \cdot [\phi \tilde{\sigma}^C] + g \tilde{\nabla} p - \frac{\eta_f}{\rho_f L^2} \tilde{\nabla} \cdot \left[ \eta_{\text{eff}} \tilde{D}(U \tilde{\mathbf{u}}_m - W \alpha \frac{\pi}{\beta} \tilde{\mathbf{w}}) \right] = g \beta \mathbf{e}_g.$$

This motivates the introduction of two dimensionless number; the Reynolds number of the fluid,  $\text{Re}$ , and a Stokes number,  $\text{St}$ , relating the importance of the mixture kinetic energy to the dissipation by drag forces:

$$\text{Re} := \frac{\rho_f U L}{\eta_f} \quad \text{St} := \frac{W}{U} = \frac{(\rho_g - \rho_f) g D_g^2}{18 U \eta_f} = \frac{\rho_g U D_g^2}{18 L \eta_f}.$$

The relationship between these dimensionless numbers is given by ratios of density and length,

$$\text{St} = \frac{\alpha + 1}{18} \frac{\rho_f U D_g^2}{L \eta_f} = \frac{\alpha + 1}{18} \text{Re} \left( \frac{D_g}{L} \right)^2 = \frac{\alpha + 1}{18} \epsilon^2 \text{Re},$$

where  $\epsilon := \frac{D_g}{L}$  denote the ratio of the grains diameter to the characteristic simulation length.

Dividing our equation by  $g(\alpha + 1)$ , we get the dimensionless momentum conservation equation for the mixture,

$$\frac{\beta}{\alpha + 1} \frac{D\tilde{\mathbf{u}}_m}{D\tilde{t}} + \text{St}^2 \tilde{\nabla} \cdot \left[ \frac{\pi}{\beta} \tilde{\mathbf{w}} \otimes \tilde{\mathbf{w}} \right] - \tilde{\nabla} \cdot [\phi \tilde{\sigma}^C] - \frac{1}{\alpha + 1} \tilde{\nabla} \cdot \left[ \frac{\eta_{\text{eff}}}{\text{Re}} \tilde{D}(\tilde{\mathbf{u}}_m - \text{St} \alpha \frac{\pi}{\beta} \tilde{\mathbf{w}}) - \tilde{p} \mathbb{I} \right] = \frac{\beta}{\alpha + 1} \mathbf{e}_g. \quad (8.16)$$

**Fluctuation momentum conservation** Rewriting (8.15) using our dimensionless variables yield

$$(\alpha + 1) \rho_f \phi \frac{W}{L} \left[ U \frac{\partial \tilde{\mathbf{w}}}{\partial \tilde{t}} + U I(\tilde{\mathbf{u}}_m, \tilde{\mathbf{w}}) + W I \left( \frac{1 - 2\phi - \alpha\phi}{2\beta} \tilde{\mathbf{w}}, \tilde{\mathbf{w}} \right) \right] + \beta \phi g \frac{\rho_g - \rho_f}{w_\infty} \Xi W \tilde{\xi} \tilde{\mathbf{w}} = \alpha \phi \left( g \rho_f \nabla \tilde{p} - \frac{\eta_f}{L^2} \nabla \cdot \left[ \eta_{\text{eff}} \tilde{D}(U \tilde{\mathbf{u}}_m - \alpha W \frac{\pi}{\beta} \tilde{\mathbf{w}}) \right] \right) + (\alpha + 1) \rho_f g \tilde{\nabla} \cdot [\phi \tilde{\sigma}^C],$$

where  $I(\mathbf{v}, \mathbf{w}) := (\mathbf{v} \cdot \tilde{\nabla}) \mathbf{w} + (\mathbf{w} \cdot \tilde{\nabla}) \mathbf{v}$ . Dividing both sides by  $\rho_g g$ , and remarking that  $WU/L = \text{St}g$ , we get

$$\phi \text{St} \left[ \frac{\partial \tilde{\mathbf{w}}}{\partial \tilde{t}} + I(\tilde{\mathbf{u}}_m, \tilde{\mathbf{w}}) + \text{St} I \left( \frac{1 - 2\phi - \alpha\phi}{2\beta} \tilde{\mathbf{w}}, \tilde{\mathbf{w}} \right) \right] + \frac{\alpha\beta}{\alpha + 1} \phi \Xi \tilde{\xi} \tilde{\mathbf{w}} = \frac{\alpha\phi}{\alpha + 1} \left( \tilde{\nabla} \tilde{p} - \frac{1}{\text{Re}} \tilde{\nabla} \cdot \left[ \eta_{\text{eff}} \tilde{D}(\tilde{\mathbf{u}}_m - \alpha \text{St} \frac{\pi}{\beta} \tilde{\mathbf{w}}) \right] \right) + \tilde{\nabla} \cdot [\phi \tilde{\sigma}^C]. \quad (8.17)$$

**Mass conservation** The dimensionless mass conservation equations read:

$$\tilde{\nabla} \cdot \left[ \tilde{\mathbf{u}}_m - \text{St} \frac{\alpha\pi}{\beta} \tilde{\mathbf{w}} \right] = 0 \quad (8.18)$$

$$\frac{\partial \beta}{\partial \tilde{t}} - \tilde{\nabla} \cdot [\beta \tilde{\mathbf{u}}_m] = 0. \quad (8.19)$$

Alternatively, Equation (8.19) can be replaced with (8.20),

$$\frac{\partial \phi}{\partial \tilde{t}} - \tilde{\nabla} \cdot [\phi \tilde{\mathbf{u}}_g] = 0, \quad (8.20)$$

where  $\tilde{\mathbf{u}}_g := \frac{1}{U} \mathbf{u}_g$ .

### 8.2.6 Particular cases

**Equilibrium** First, let us look at static equilibrium, which means  $\mathbf{u}_m = \mathbf{w} = \mathbf{0}$ . The dimensionless conservation equations read simply

$$\begin{cases} -\tilde{\nabla} \cdot [\phi \tilde{\sigma}^C] + \frac{1}{\alpha+1} \tilde{\nabla} \tilde{p} = \frac{\beta}{\alpha+1} \mathbf{e}_g \\ -\tilde{\nabla} \cdot [\phi \tilde{\sigma}^C] = \frac{\alpha\phi}{\alpha+1} \tilde{\nabla} \tilde{p}, \end{cases}$$

from which we deduce easily

$$\begin{cases} \tilde{\nabla} \tilde{p} = \mathbf{e}_g \\ -\tilde{\nabla} \cdot [\phi \tilde{\sigma}^C] = \frac{\alpha\phi}{\alpha+1} \mathbf{e}_g. \end{cases}$$

Hence, up to a constant field, the pore pressure satisfies  $p = \rho_f g z$ , and corresponds to the expected hydrostatic pressure. We also have  $-\nabla \cdot [\phi \sigma_g] = \phi \alpha \rho_f \mathbf{g} = \phi \rho_g \mathbf{g} - \phi \nabla p$ ; the contact force thus oppose the action of gravity on the grains reduced by the buoyancy force given by Archimede's principle. Again, this is the expected result.

**Incompressible flow** The divergence-free condition for the volume-averaged velocity  $\mathbf{u}_v$  is expressed by Equation (8.18).

In the limit case where  $St = 0$ , the mass-averaged velocity  $\mathbf{u}_m$  becomes also divergence-free. In a similar fashion, if  $\Xi = +\infty$ , then (8.17) imposes  $\mathbf{w} = \mathbf{0}$  and then once again the flow becomes incompressible.

**Single phase limit** If we consider that the surrounding fluid is massless and inviscid,  $\alpha = +\infty$ , and  $\frac{\beta}{\alpha+1} = \phi$ . Moreover, if the grains are small enough w.r.t. the characteristic length, i.e.,  $\epsilon \ll 1$ , then  $\frac{St}{\alpha+1} \ll 1$ .

Overall, the mixture momentum conservation equation (8.16) becomes

$$\phi \frac{D_{\tilde{\mathbf{u}}_g} \tilde{\mathbf{u}}_g}{D_{\tilde{\mathbf{u}}_g} \tilde{t}} - \tilde{\nabla} \cdot [\phi \tilde{\sigma}^C] = \phi \mathbf{e}_g,$$

which is the equation that we used for the dry case in Chapter 7 in the case of a vanishing Newtonian granular viscosity. However, in the vanishing grain size limit  $St$  will also get close to zero, and the flow will be incompressible, thus departing from the single phase equations.

**Orders of magnitude** We now look at the order of magnitude of the different terms in the dimensionless conservation equations for different physical parameters. Numbers of special interest are  $St^2$  — the importance of the  $\mathbf{w}$  transport term in the mixture momentum balance equation,  $\frac{St}{Re} = \frac{\alpha\epsilon^2}{18}$  — the viscosity term on  $\mathbf{w}$ , and  $\frac{St}{\Xi}$  and  $\frac{St^2}{\Xi}$  — the inertial terms in the fluctuation equations.

Since as we already mentioned, we know that we overestimate  $\mathbf{w}$ , and since when using the Richardson-Zaki closure the term  $\beta\Xi\tilde{\xi}$  will be greater than 1 for most regimes,  $St < 10^{-1}$  appears to be a satisfying criterion for neglecting all inertial forces in  $\mathbf{w}$ .

- We first consider sand grains in water, for which  $\rho_f \sim 10^3 \text{ kg.m}^{-3}$ ,  $\rho_g \sim 2.5 \times 10^3 \text{ kg.m}^{-3}$ , and  $\eta_f \sim 10^{-3} \text{ Pa.s}$ . We have  $Re = 3 \times 10^6 L^{\frac{3}{2}}$ ,  $\alpha \sim 1$  and  $St = 10^5 L^{\frac{3}{2}} \epsilon^2 = 10^5 D_g^2 L^{-\frac{1}{2}}$ . This scaling means that when the size of the particles doubles, the size of the domain has to be multiplied by 16 to recover the same ratio for the mass-averaged and fluctuation velocities, and by 256 to get a similar influence for the inertial terms in  $\mathbf{w}$ . For grains of diameter  $D_g = 1 \text{ mm}$ , getting  $St < 0.1$  requires  $L > 1 \text{ m}$ . However, for  $D_g = 1 \text{ cm}$ , we must have  $L > 10 \text{ km}$ . As  $\frac{St}{Re} \ll 1$  as soon as  $L > 10 D_g$ , the viscous forces on  $\mathbf{w}$  can be neglected in both of those scenarios.

- Now, consider the same material in air, i.e.,  $\rho_f \sim 1\text{kg.m}^{-3}$ , and  $\eta_f \sim 10^{-3}\text{Pa.s}$ . We have  $\text{Re} = 3 \times 10^5 L^{\frac{3}{2}}$ ,  $\alpha \sim 2500$  and  $\text{St} = 10^7 L^{\frac{3}{2}} \epsilon^2 = 10^7 D_g^2 L^{-\frac{1}{2}}$ . For grains of diameter  $D_g = 1\text{mm}$ , getting  $\text{St} < 0.1$  requires  $L > 10^4\text{m}$ . However, for  $D_g = 100\mu\text{m}$ , it suffices to have  $L > 1\text{m}$ . Moreover  $\frac{\text{St}}{\text{Re}} \ll 1$  as soon as  $L > 10^2 D_g$ , thus in both cases the viscous forces on  $\mathbf{w}$  can be neglected.

Note that these estimates are quite conservative, and that these values for the characteristic length can be reduced in regimes where the volume fraction of grains is not near zero.

**Simplified model** Now, supposing that we can neglect the inertial and viscous terms in  $\mathbf{w}$ , our conservation equations boil down to

$$\begin{cases} \frac{\beta}{\alpha+1} \frac{D\tilde{\mathbf{u}}_m}{D\tilde{t}} - \tilde{\nabla} \cdot [\phi \tilde{\boldsymbol{\sigma}}^C] - \frac{1}{\alpha+1} \tilde{\nabla} \cdot \left[ \frac{\eta_{\text{eff}}}{\text{Re}} \tilde{D}(\tilde{\mathbf{u}}_m) - \tilde{p}\mathbb{I} \right] = \frac{\beta}{\alpha+1} \mathbf{e}_g, \\ \frac{\alpha\phi}{\alpha+1} \left( \tilde{\nabla} \tilde{p} - \frac{1}{\text{Re}} \nabla \cdot [\eta_{\text{eff}} \tilde{D}(\tilde{\mathbf{u}}_m)] \right) + \tilde{\nabla} \cdot [\phi \tilde{\boldsymbol{\sigma}}^C] = \beta \phi \Xi \tilde{\xi} \tilde{\mathbf{w}}. \end{cases}$$

We are primarily interested in the grain velocity, which is given by  $\tilde{\mathbf{u}}_g = \tilde{\mathbf{u}}_m + \tilde{\text{St}} \frac{\mathbf{w}}{\beta}$ . The influence of the remaining viscous term in the simplified fluctuation conservation equation on  $\mathbf{u}_g$  is thus of order  $\frac{\text{St}}{\text{Re}}$ , which we already assumed to be negligible. Our simplified model will therefore be defined as the solution to

$$\begin{cases} \frac{\beta}{\alpha+1} \frac{D\tilde{\mathbf{u}}_m}{D\tilde{t}} - \tilde{\nabla} \cdot [\phi \tilde{\boldsymbol{\sigma}}^C] - \frac{1}{\alpha+1} \tilde{\nabla} \cdot \left[ \frac{\eta_{\text{eff}}}{\text{Re}} \tilde{D}(\tilde{\mathbf{u}}_m) - \tilde{p}\mathbb{I} \right] = \frac{\beta}{\alpha+1} \mathbf{e}_g \\ \frac{\alpha\phi}{\alpha+1} \tilde{\nabla} \tilde{p} + \tilde{\nabla} \cdot [\phi \tilde{\boldsymbol{\sigma}}^C] = \frac{\alpha}{\alpha+1} \phi \beta \Xi \tilde{\xi} \tilde{\mathbf{w}} \\ \tilde{\nabla} \cdot \left[ \tilde{\mathbf{u}}_m - \text{St} \alpha \frac{\pi}{\beta} \tilde{\mathbf{w}} \right] = 0 \\ \frac{\partial \beta}{\partial \tilde{t}} + \tilde{\nabla} \cdot [\beta \mathbf{u}_m] = 0, \end{cases}$$

with the supplemental condition that  $\dot{\epsilon}_g = D(\tilde{\mathbf{u}}_m + \text{St} \frac{1-\phi}{\beta} \tilde{\mathbf{w}})$  and  $\tilde{\boldsymbol{\sigma}}^C$  should follow the Drucker-Prager flow rule of Equation (8.10). Since we have now an explicit expression for  $\mathbf{w}$ , we can eliminate this variable and get

$$\begin{cases} \frac{\beta}{\alpha+1} \frac{D\tilde{\mathbf{u}}_m}{D\tilde{t}} - \frac{1}{\alpha+1} \tilde{\nabla} \cdot \left[ \frac{\eta_{\text{eff}}}{\text{Re}} \tilde{D}(\tilde{\mathbf{u}}_m) - \tilde{p}\mathbb{I} \right] + \tilde{\nabla} \cdot [\phi \tilde{\boldsymbol{\lambda}}] = \frac{\beta}{\alpha+1} \mathbf{e}_g \\ \frac{\tilde{\nabla} \cdot [\tilde{\mathbf{u}}_m]}{\alpha+1} - \frac{\text{St}}{\Xi} \tilde{\nabla} \cdot \left[ \frac{\alpha}{\alpha+1} \frac{\pi}{\beta^2 \xi} \tilde{\nabla} \tilde{p} - \frac{1-\phi}{\beta^2 \xi} \tilde{\nabla} \cdot [\phi \tilde{\boldsymbol{\lambda}}] \right] = 0 \\ \phi \tilde{D}(\tilde{\mathbf{u}}_m) - \phi \frac{\text{St}}{\Xi} \tilde{D} \left( \frac{1-\phi}{\beta^2 \xi} \tilde{\nabla} \tilde{p} - \frac{\alpha+1}{\alpha} \frac{1-\phi}{\phi \beta^2 \xi} \tilde{\nabla} \cdot [\phi \tilde{\boldsymbol{\lambda}}] \right) = \tilde{\gamma}, \end{cases} \quad (8.21)$$

with once again our usual frictional rheology on  $\gamma := \phi D(\mathbf{u}_f)$  and  $\boldsymbol{\lambda} := -\boldsymbol{\sigma}^C$ . Each row in (8.21) as been scaled so as to highlight the symmetry of the system.

This model is interesting for a few reasons:

- We still retrieve an incompressible flow for  $\text{St} = 0$  or  $\Xi = \infty$ ;
- While this does not fall *at all* in the range of validity of our hypothesis, we also retrieve the dry granular equations of Chapter 7 when taking the  $\alpha = \infty$ ,  $\text{St} = \infty$  limit;
- Ignoring the grain-grain contact stress, the fluctuation equations now reads exactly as Darcy's law;

- Still ignoring grain-grain interactions, the system we have to solve is deeply similar to that of near-incompressible flows, i.e., a penalized Navier-Stokes problem,

$$\left\{ \begin{array}{l} \frac{\beta}{\alpha+1} \frac{D\tilde{\mathbf{u}}_m}{D\tilde{t}} - \frac{1}{\alpha+1} \tilde{\nabla} \cdot \left[ \frac{\eta_{\text{eff}}}{\text{Re}} \tilde{D}(\tilde{\mathbf{u}}_m) - \tilde{p} \mathbb{I} \right] = \frac{\beta}{\alpha+1} \mathbf{e}_g \\ \frac{\tilde{\nabla} \cdot [\tilde{\mathbf{u}}_m]}{\alpha+1} - \frac{\alpha}{\alpha+1} \frac{\text{St}}{\Xi} \tilde{\nabla} \cdot \left[ \frac{\pi}{\beta^2 \xi} \tilde{\nabla} \tilde{p} \right] = 0 \\ \frac{\partial \beta}{\partial \tilde{t}} + \tilde{\nabla} \cdot [\beta \mathbf{u}_m] = 0. \end{array} \right.$$

The reduced diphasic model (8.21) is very similar to the single phase model of Chapter 7, except that the degrees of freedom are now the combination of the mass-averaged velocity  $\mathbf{u}_m$  and the pressure  $p$ . It could be solved in a similar fashion, but one would have to be wary that the stiffness matrix is no longer guaranteed to be positive-definite, and that another term, coming from the nonlinear Laplacian of  $\lambda$  in the definition of  $\gamma$ , contributes to the Delassus operator. Since  $\mathbf{w}$  does not need to be discretized, this reduced problem could be cheaper to solve than the original one, and may still be able to capture a wide range of phenomena. Following the approach of the dry case, the simplified model can also be shown to be dissipative. Despite these potential advantages, in the following section we will go back to the general case and propose a numerical method to solve our original equations. Adapting this method to the simplified model should be straightforward.

### 8.3 Numerical resolution of the two-phase equations

We consider our original dimensionless system of four equations — the two conservation equations, (8.16) and (8.17), and the two mass conservation equations (8.18) and (8.20) — plus our frictional rheology constraint (8.10). We now drop the tildes to lighten notations.

#### 8.3.1 Time discretization

Just like in the dry case presented in Chapter 7, we will use a semi-implicit timestepping scheme, that first solves the momentum conservation equations, then advect the volume fraction field. That is, we will first solve equations (8.16), (8.17) subject to the volume-averaged incompressibility (8.18) and the frictional rheology (8.10) to get end-of-step values for  $\mathbf{u}, \mathbf{w}, p$  and  $\sigma^C$ , then solve the volume fraction transport equation (8.20) using these new velocities.

**Frictional contacts** Once again, we thus have to linearize the end-of-step maximal volume fraction constraint, following

$$(\phi(t + \Delta_t) \leq \phi_{\max}) \sim (\phi^k + \phi^k \Delta_t D(\mathbf{u}_g) \leq \phi_{\max}).$$

Just as in Chapter 7, we thus define

$$\gamma := \phi^k D(\mathbf{u}_g) + \frac{\phi_{\max} - \phi}{\Delta_t} \mathbb{I} \quad \text{and} \quad \lambda := -\sigma^C, \quad (8.22)$$

such that our frictional rheology constraint is expressed as  $(\gamma, \lambda) \in \mathcal{DP}(\mu)$ .

**Material velocity derivatives** We suppose that the material derivative for the velocity of each phase can be approximated to the first order as

$$\frac{D_{\mathbf{u}_{f,g}} \mathbf{u}_{f,g}}{D_{\mathbf{u}_{f,g}} t} = \frac{\mathbf{u}_{f,g}^{k+1} - \mathbf{u}_{f,g}^k}{\Delta_t} + O(\Delta_t).$$

For instance,  $u$  can be computed using the characteristics morphism, or by recovering the velocity from particles in the case of hybrid schemes.

Using this approximation, we can get simpler expressions for the inertial terms of the mixture and fluctuation momentum conservation equations. Indeed,

$$\begin{aligned}\phi(\alpha+1)\frac{D_{\mathbf{u}_g}\mathbf{u}_g}{D_{\mathbf{u}_g}t} + (1-\phi)\frac{D_{\mathbf{u}_f}\mathbf{u}_f}{D_{\mathbf{u}_f}t} &\sim \frac{1}{\Delta_t} \left[ \beta\mathbf{u}_m - (\alpha+1)\phi u_g(\mathbf{u}_g^k) - (1-\phi)u_f(\mathbf{u}_f^k) \right] \\ \frac{D_{\mathbf{u}_g}\mathbf{u}_g}{D_{\mathbf{u}_g}t} - \frac{D_{\mathbf{u}_f}\mathbf{u}_f}{D_{\mathbf{u}_f}t} &\sim \frac{1}{\Delta_t} \left[ \text{St} \mathbf{w} + u_f(\mathbf{u}_f^k) - u_g(\mathbf{u}_g^k) \right].\end{aligned}$$

Using this insight, we rewrite our momentum conservation equations as

$$\begin{aligned}\frac{\beta}{\alpha+1}\frac{\mathbf{u}_m}{\Delta_t} + \nabla \cdot [\phi\lambda] - \frac{1}{\alpha+1}\nabla \cdot \left[ \frac{\eta_{\text{eff}}}{\text{Re}} D(\mathbf{u}_m - \text{St}\alpha\frac{\pi}{\beta}\mathbf{w}) - \tilde{\mathbf{p}}\mathbb{I} \right] \\ = \frac{\beta}{\alpha+1}\mathbf{e}_g + \frac{1}{\Delta_t} \left[ \phi u_g(\mathbf{u}_g^k) + \frac{1-\phi}{\alpha+1}u_f(\mathbf{u}_f^k) \right]\end{aligned}\quad (8.23)$$

$$\begin{aligned}\phi \left( \frac{\text{St}}{\Delta_t} + \frac{\alpha\beta}{\alpha+1}\Xi\xi \right) \mathbf{w} - \frac{\alpha\phi}{\alpha+1} \left( \nabla p - \frac{1}{\text{Re}} \nabla \cdot \left[ \eta_{\text{eff}} D(\mathbf{u}_m - \alpha\text{St}\frac{\pi}{\beta}\mathbf{w}) \right] \right) \\ + \nabla \cdot [\phi\lambda] = \frac{\phi}{\Delta_t} (u_g(\mathbf{u}_g^k) - u_f(\mathbf{u}_f^k))\end{aligned}\quad (8.24)$$

### 8.3.2 Variational formulation

We first perform a change of variable which will allow us to obtain a symmetric system, and introduce the velocity field  $\hat{\mathbf{w}} := \sqrt{\text{St}}\frac{1-\phi}{\beta}\mathbf{w}$ . The velocity of each phase can be reconstructed from  $\mathbf{u}_m$  and  $\hat{\mathbf{w}}$  as  $\mathbf{u}_g = \mathbf{u}_m + \sqrt{\text{St}}\hat{\mathbf{w}}$  and  $\mathbf{u}_f = \mathbf{u}_m - (\alpha+1)\sqrt{\text{St}}\frac{\phi}{1-\phi}\hat{\mathbf{w}}$ .

Let us consider a simulation domain  $\Omega$ , with, for the sake of simplicity, homogeneous Dirichlet boundary conditions. Let as usual  $V$  denote a subspace of  $H^1(\Omega)^d$  satisfying the boundary conditions, and  $T(\Omega) \sim L_2(\Omega)^{5d}$  be the space of square-integrable symmetric tensor fields on  $\Omega$ . Our two-phase flow is described by the solution to the following variational formulation:

Find  $\mathbf{u}_m, \hat{\mathbf{w}} \in V^2$ ,  $p \in L_2(\Omega)$  and  $\gamma, \lambda \in T^2$  such that

$$\begin{aligned}a(\mathbf{u}_m, \mathbf{v}) + e(\hat{\mathbf{w}}, \mathbf{v}) - b(p, \mathbf{v}) - g(\lambda, \mathbf{v}) &= l(\mathbf{v}) \quad \forall \mathbf{v} \in V \\ e(\mathbf{u}_m, \mathbf{z}) + r(\hat{\mathbf{w}}, \mathbf{z}) - c(p, \mathbf{z}) - h(\lambda, \mathbf{z}) &= f(\mathbf{z}) \quad \forall \mathbf{z} \in V \\ -b(\mathbf{q}, \mathbf{u}_m) - c(\mathbf{q}, \hat{\mathbf{w}}) &= 0 \quad \forall \mathbf{q} \in L_2(\Omega) \\ -g(\tau, \mathbf{u}_m) - h(\tau, \hat{\mathbf{w}}) + m(\gamma, \tau) &= k(\tau) \quad \forall \tau \in T \\ (\gamma, \lambda) &\in \mathcal{DP}(\mu).\end{aligned}$$

where  $a = a_1 + a_2$ ,  $r = r_1 + r_2$ ,

$$\begin{aligned}a_1(\mathbf{u}, \mathbf{v}) &:= \frac{1}{\Delta_t} \int_{\Omega} \frac{\beta}{\alpha+1} \langle \mathbf{u}, \mathbf{v} \rangle & a_2(\mathbf{u}, \mathbf{v}) &:= \frac{1}{(\alpha+1)\text{Re}} \int_{\Omega} \eta_{\text{eff}} D(\mathbf{u}) : D(\mathbf{v}) \\ r_1(\mathbf{w}, \mathbf{z}) &:= \int_{\Omega} \frac{\phi\beta}{1-\phi} \left( \frac{\text{St}}{\Delta_t} + \frac{\alpha\beta}{\alpha+1}\Xi\xi \right) & r_2(\mathbf{w}, \mathbf{z}) &:= \frac{\text{St}}{\text{Re}} \frac{\alpha^2}{\alpha+1} \int_{\Omega} \eta_{\text{eff}} D(\phi\mathbf{w}) : D(\phi\mathbf{z}) \\ e(\mathbf{w}, \mathbf{z}) &:= -\frac{\alpha\sqrt{\text{St}}}{\alpha+1} \int_{\Omega} \frac{\eta_{\text{eff}}}{\text{Re}} D(\phi\mathbf{w}) : D(\mathbf{z}) & m(\gamma, \tau) &:= \int_{\Omega} \gamma : \tau \\ b(p, \mathbf{v}) &:= \frac{1}{\alpha+1} \int_{\Omega} p \nabla \cdot \mathbf{z} & c(p, \mathbf{z}) &:= \frac{\alpha\sqrt{\text{St}}}{\alpha+1} \int_{\Omega} \phi \langle \nabla p, \mathbf{z} \rangle \\ g(\tau, \mathbf{u}) &:= \int_{\Omega} \phi \tau : D(\mathbf{u}) & h(\tau, \mathbf{u}) &:= \sqrt{\text{St}} \int_{\Omega} \phi \tau : D(\mathbf{w})\end{aligned}$$

$$\begin{aligned}
 l(\mathbf{v}) &:= \int_{\Omega} \frac{\beta}{\alpha+1} \langle \mathbf{e}_g, \mathbf{v} \rangle + \frac{1}{\Delta_t} \int_{\Omega} \langle (\alpha+1)\phi u_g(\mathbf{u}_g^k) + (1-\phi)u_f(\mathbf{u}_f^k), \mathbf{v} \rangle \\
 k(\boldsymbol{\tau}) &:= \int_{\Omega} \frac{\phi_{\max} - \phi}{\Delta_t} \frac{\text{Tr } \boldsymbol{\tau}}{d} \\
 f(\mathbf{z}) &:= \frac{\sqrt{\text{St}}}{\Delta_t} \int_{\Omega} \phi \langle u_g(\mathbf{u}_g^k) - u_f(\mathbf{u}_f^k), \mathbf{z} \rangle.
 \end{aligned}$$

*Proof.* Scaling the fluctuation momentum conservation equation (8.24) by  $\sqrt{\text{St}}$ , dividing the volume-averaged incompressibility (8.18) by  $(\alpha+1)$  and transcribing the definition of  $\hat{\mathbf{w}}$  into them, we obtain:

$$\begin{aligned}
 \frac{\phi\beta}{1-\phi} \left( \frac{\text{St}}{\Delta_t} + \frac{\alpha\beta}{\alpha+1} \Xi \xi \right) \hat{\mathbf{w}} - \frac{\alpha\phi}{\alpha+1} \left( \sqrt{\text{St}} \nabla p - \frac{1}{\text{Re}} \nabla \cdot [\eta_{\text{eff}} D(\sqrt{\text{St}} \mathbf{u}_m - \alpha \text{St} \phi \hat{\mathbf{w}})] \right) \\
 + \sqrt{\text{St}} \nabla \cdot [\phi \boldsymbol{\lambda}] = \sqrt{\text{St}} \frac{\phi}{\Delta_t} (\mathbf{u}_g(\mathbf{u}_g^k) - \mathbf{u}_f(\mathbf{u}_f^k)), \\
 \frac{1}{\alpha+1} \nabla \cdot \mathbf{u}_m - \sqrt{\text{St}} \frac{\alpha}{\alpha+1} \nabla \cdot [\phi \hat{\mathbf{w}}] = 0,
 \end{aligned}$$

which, put under weak form and using Green formulas, correspond to the second and third lines of our variational formulation. The first and fourth line are obtained through a similar process on the mixture momentum conservation equation (8.24) and the definition of  $\boldsymbol{\gamma}$ , Equation (8.22).  $\square$

### 8.3.3 Discrete system

Choosing adequate discrete spaces for our variables, and expressing the constraints on quadrature points as we did in Chapters 6 and (7), we ultimately obtain a system which has the form:

$$\begin{cases}
 A\mathbf{u} + E^T \mathbf{w} = \mathbf{l} + B^T \mathbf{p} + G^T \boldsymbol{\lambda} \\
 E\mathbf{u} + R\mathbf{w} = \mathbf{f} + C^T \mathbf{p} + H^T \boldsymbol{\lambda} \\
 \mathbf{0} = B\mathbf{u} + C\mathbf{w} \\
 \boldsymbol{\gamma} = \mathbf{k} + G\mathbf{u} + H\mathbf{w} \\
 (\boldsymbol{\lambda}_{[i]}, \boldsymbol{\gamma}_{[i]}) \in \mathcal{DP}(\mu) \quad \forall i.
 \end{cases} \quad (8.25)$$

System (8.25) is a DCFP, with linear constraints and whose degrees of freedom are given by the couple  $(\mathbf{u}, \mathbf{w})$ . Concatenating the system matrices as

$$\bar{A} := \begin{pmatrix} A & E^T \\ E & R \end{pmatrix} \quad \bar{H} := (G, H) \quad \bar{C} := (B, C),$$

System (8.25) can be rewritten in a more familiar manner as

$$\begin{cases}
 \bar{A}(\mathbf{u}; \mathbf{w}) = \begin{pmatrix} \mathbf{l} \\ \mathbf{f} \end{pmatrix} + \bar{C}^T \mathbf{p} + \bar{H}^T \boldsymbol{\lambda} \\
 \bar{C}(\mathbf{u}; \mathbf{w}) = \mathbf{0} \\
 \boldsymbol{\gamma} = \mathbf{k} + \bar{H}(\mathbf{u}; \mathbf{w}) \\
 (\boldsymbol{\lambda}_{[i]}, \boldsymbol{\gamma}_{[i]}) \in \mathcal{DP}(\mu) \quad \forall i.
 \end{cases} \quad (8.26)$$

The matrix  $A$  is always symmetric positive-definite, and discretizing the fluctuation field only where  $\hat{\phi} > 0$  — which is sensible, as  $\hat{\mathbf{w}}$  should vanish elsewhere —  $R$  is also positive-definite; therefore  $\bar{A}$  is symmetric positive-positive as well.

**Numerical resolution** Chapters 3 and 4 were dedicated to solving numerically the DCFP with linear constraints. However, notice that System 8.26 features two peculiarities that restrict the choice of numerical solvers.

- First, the inverse of the stiffness matrix,  $\bar{A}^{-1}$ , is dense. This difficulty can be dealt with by either using the numerical algorithms proposed in Chapter 6, or making the high Reynolds hypothesis of Chapter 7 and using the two-steps algorithm of Section 7.3.4.
- The matrix  $\bar{C}$  is not surjective — the pressure field is simply defined up to a constant — and thus the Schur complement of the linear constraints may not be invertible, discarding some of the methods that we proposed to deal with linear constraints in a DCFP. If the discretization spaces for the velocities and pressure,  $V_h$  and  $Q_h$ , satisfy an *inf-sup* criterion, this problem can be solved by enforcing a supplemental zero-average constraint for the pressure  $p$ . Otherwise, we can assume that slightly relaxing the incompressibility constraint of the volume-averaged velocity is acceptable, and penalize the pressure variable.

In practice, as we put more of an emphasis on solving efficiency rather than on the precision of the fluid velocity field, we use the two-step approximate algorithm together with the relaxation of the incompressibility constraints.

That is, we first solve a Stokes-like problem, with either the zero-average constraint or a penalization strategy, depending on whether our choice for  $V_h$  and  $Q_h$  satisfies the *inf-sup* condition, to obtain a candidate velocity  $(\underline{u}; \underline{w})^*$ , e.g.,

$$\begin{pmatrix} A & E^\top & -B^\top \\ E & R & -C^\top \\ -B & -C & \mathbb{1}^\top \end{pmatrix} \begin{pmatrix} \underline{u}^* \\ \underline{w}^* \\ \underline{p} \\ \kappa \end{pmatrix} = \begin{pmatrix} \underline{l} \\ \underline{f} \\ \underline{0} \\ 0 \end{pmatrix} \quad \text{or} \quad \begin{pmatrix} A & E^\top & -B^\top \\ E & R & -C^\top \\ -B & -C & -c\mathbb{I} \end{pmatrix} \begin{pmatrix} \underline{u}^* \\ \underline{w}^* \\ \underline{p} \end{pmatrix} = \begin{pmatrix} \underline{l} \\ \underline{f} \\ \underline{0} \end{pmatrix},$$

where  $\kappa$  is the Lagrange multiplier associated to the zero-pressure average condition.

Then, we solve a lumped DCFP to get the velocity increment satisfying the frictional constraints,

$$\begin{cases} \tilde{\bar{A}} \Delta(\underline{u}; \underline{w}) = \bar{C}^\top \underline{p} + \bar{H}^\top \underline{\lambda} \\ \bar{C} \Delta(\underline{u}; \underline{w}) = -c \underline{p} \\ \underline{\gamma} = \bar{H}(\underline{u}; \underline{w}) + \underline{k} + \bar{H}(\underline{u}; \underline{w})^* \\ (\underline{\lambda}_{[i]}, \underline{\gamma}_{[i]}) \in \mathcal{DP}(\mu) \quad \forall i, \end{cases} \quad (8.27)$$

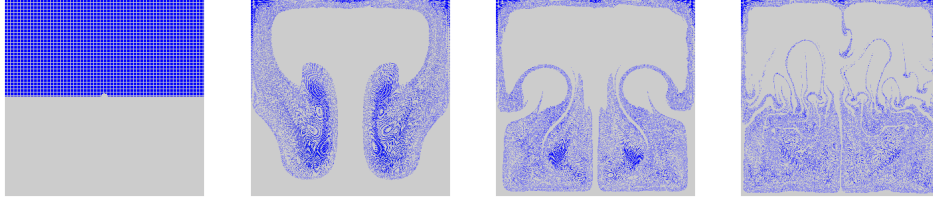
where  $\tilde{\bar{A}} = \text{diag}(\tilde{\bar{A}}_1, \tilde{\bar{R}}_1)$  is the concatenation of the lumped mass matrices for the mixture and fluctuation inertia. In practice, we use the matrix-free Gauss–Seidel algorithm presented in Chapter 4, where the linear constraints are accounted for every few iterations by solving the symmetric positive-definite system

$$(\bar{C} \tilde{\bar{A}}^{-1} \bar{C}^\top + c\mathbb{I}) \underline{p} = -\bar{C} \tilde{\bar{A}}^{-1} \bar{H}^\top \underline{\lambda}^k.$$

### 8.3.4 Spatial discretization

In the previous paragraph we assumed that suitable discretization spaces for the velocity, pressure, grain volume fraction and stress fields had been chosen. We did not pursue this investigation in details yet, and simply adapted the discretization strategies devised for the dense case. We considered either:

- Piecewise-constant approximations on a triangular mesh for the velocity, volume fraction, and stress field, with a piecewise-linear pressure field. As in Section 7.2.2, advection of the volume fraction and velocity fields was done using an upwind scheme.
- MPM discretization of the volume fraction field, as in Section 7.3, with trilinear velocities and trilinear, trilinear-discontinuous, or particle-based stresses. For the sake of simplicity



**Figure 8.2:** *Rayleigh–Taylor instability. A slight perturbation is introduced at the center of the interface between an upper, heavy fluid and a lighter one (here  $\alpha = 1000$ ). The fluids are supposed fully immiscible ( $St = 0$ ) and viscous ( $Re = 3 \cdot 10^{-2}$ ).*

we also used trilinear pressure field, which has very poor theoretical properties yet yielded satisfying results when relaxing the volume-incompressibility constraint. The advected velocity of the granular phase is recovered from the particles, while the characteristics method is used for the fluid velocities. Particle-based advection could also be used for the fluid, but one would then have to perform resampling to ensure that the sum of volume fractions of both phases remain always equal to 1.

## 8.4 Results

All the results presented below were obtained with an extension of our MPM simulation framework from Chapter 7 to diphasic flows.

### 8.4.1 Rayleigh-Taylor instability

We performed a first test to check that independently of our treatment of the granular rheology, our diphasic simulation code was behaving as expected. Figure 8.2 shows the development of a Rayleigh–Taylor instability between two Newtonian fluids, the upper one being much heavier than the bottom one. We capture the characteristic plumes associated to this instability (compare with a similar simulation in e.g., Boyer 2001, Figure II.4.12), though our particle-based discretization leads to a rather messy end state.

### 8.4.2 Sedimentation

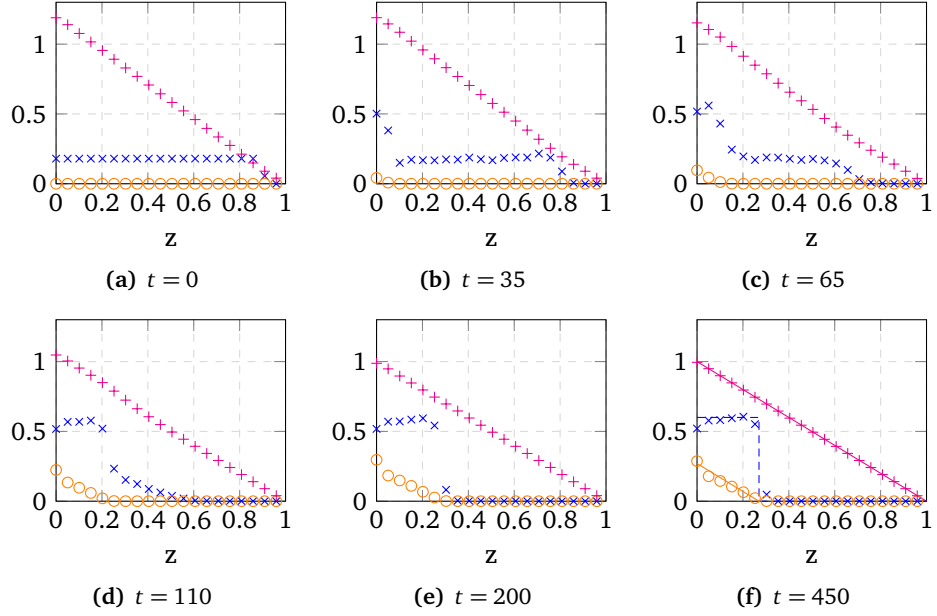
We consider a 2D water container of width and height  $W = H = 1\text{m}$ , rigid grains of volumetric mass  $\rho_g = 2500\text{kg.m}^{-3}$ , and a maximal volume fraction  $\phi_{\max} = 0.6$ . The initial volume fraction of grains  $\phi^0(x, z)$  is such that

$$\phi^0(x, z) = \begin{cases} 0.3\phi_{\max} & \text{if } z \leq H_0 = 0.9H \\ 0 & \text{if } z > H_0. \end{cases}$$

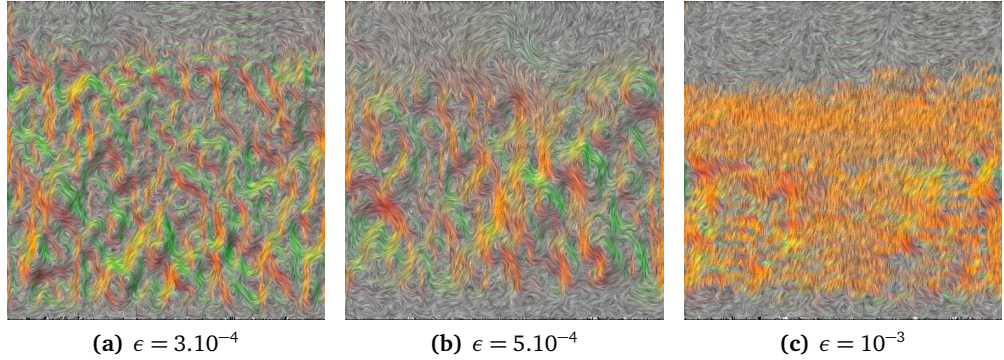
Figure 8.3 shows the evolution of the volume fraction, pore pressure  $p$  and particulate normal stress  $\text{Tr} \boldsymbol{\lambda}$  fields averaged over the width of the container for grains of diameter  $D_g = 0.5\text{mm}$ . Following the considerations of Section 8.2.6, the pressure has been made dimensionless w.r.t. the characteristic pressure  $\rho_f g H$ , and the particulate normal stress w.r.t.  $\alpha \rho_f g H$ . We observe that the final state of our simulated system does correspond to the analytical predictions of Section 8.2.6.

**Influence of grain diameters** The diameter of the grains (or rather, the ratio of length scales  $\epsilon := D_g/H$ ) has a dramatic influence on the sedimentation process. However, the grains diameters has no influence on the final state of the system.





**Figure 8.3:** Horizontally-averaged volume fraction of grains, dimensionless pore pressure and effective stress at different instants of a simulated sedimentation process. In the lower-right graph, lines correspond to the theoretical values from Section 8.2.6.



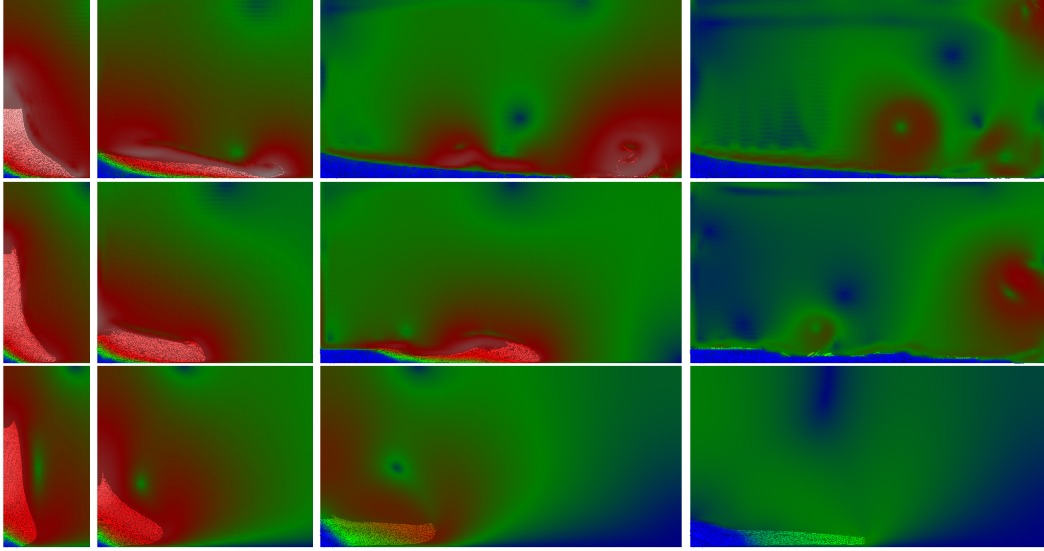
**Figure 8.4:** Line-integral convolution of the velocity field at the onset of the sedimentation process for different length ratios  $\epsilon = D_g/H$ . Hue indicates direction, brightness amplitude.

For grains of diameter 0.3mm, the dimensionless duration of the sedimentation process (scaled by  $H/\sqrt{\frac{\alpha}{\alpha+1}gH}$ ) was 30. It dropped to 17 for  $D_g = 0.5\text{mm}$  and to 6.5 for  $D_g = 5\text{mm}$ . This is explained by a much more turbulent velocity field for lower grain diameters, as illustrated in Figure 8.4.

### 8.4.3 Regimes

We attempted to reproduce the immersed column collapses of Topin, Monerie, et al. (2012).<sup>4</sup>

<sup>4</sup><http://www.irsn.fr/EN/Research/publications-documentation/Publications/PSN-RES/Pages/2012-Topin-collapse-dynamics-run-out-dense-granular-materials-fluid-videos.aspx>,



**Figure 8.5:** Snapshots at identical instants of the collapse of a granular column in three different fluids: air (top), water (middle), and viscous (bottom). Colors indicate the particle and fluid velocities, on the same scale (blue is slowest, white fastest).

We consider a 2D granular column of aspect-ratio  $a = 8$ , with grains of diameter  $D_g = 1\text{mm}$  and volumetric mass  $\rho_g = 2600\text{kg.m}^{-3}$ . The column width is  $W = 11.5D_g$ . Topin, Monerie, et al. (2012) couple the NSCD method with a fluid simulation; we use instead the  $\mu(I)$  rheology to model the frictional contacts, with  $\mu_s = 0.32$ ,  $\mu_D = 0.6$ . Three choices are considered for the surrounding fluid:

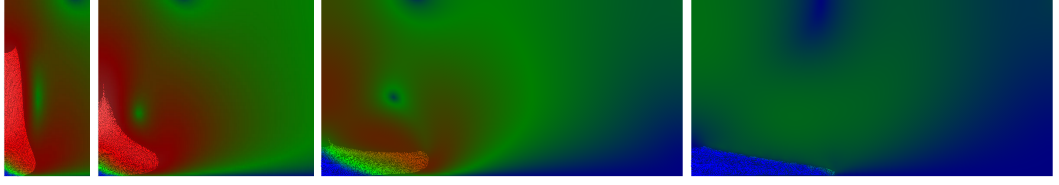
1. *air*:  $\rho_f = 1\text{kg.m}^{-3}$  and  $\eta_f = 10^{-5}\text{Pa.s}$ ;
2. *water*:  $\rho_f = 1000\text{kg.m}^{-3}$  and  $\eta_f = 10^{-3}\text{Pa.s}$ ;
3. *viscous*:  $\rho_f = 1000\text{kg.m}^{-3}$  and  $\eta_f = 1\text{Pa.s}$ .

Topin, Monerie, et al. (2012) observe that each of these choices yield a different collapse regime, which they coin respectively *grain-inertial*, *fluid-inertial*, and *viscous*. In the latter case, the collapse is simply slowed-down by the fluid, and the run-out length is much shorter than in the dry (i.e., grain-inertial) case. However, in the fluid-inertial regime, the kinetic energy initially transferred from the grains to the fluid is transferred back to the grains in the later stage of the collapse, maintaining an horizontal velocity for a much longer time than in the two other regimes. As such, the final run-out length in the fluid-inertial regime can surpass that of the dry case.

While our experiments, reproduced in Figure 8.5, allow us to retrieve the grain-inertial and fluid-inertial regimes for the choices of parameters (1.) and (2.), choice (3.) remains mostly in the fluid-inertial regime; albeit much slower than the collapse in water, the collapse in the viscous fluid achieve a barely shorter run-out length. While initially surprising, this can be explained by looking at the characteristic numbers of our model.

Indeed, in the later stage of the collapse, the grains can be either slowed down or driven by the drag forces, depending on which of the contact stress or pore pressure is preponderant. The total normal stress of the mixture is  $p + \text{Tr } \phi \lambda$ ; the fact that the velocity is mostly horizontal, and thus that the grains are not collapsing under their own weight, means that the total normal stress is higher there than in the front of the avalanche. When the grains are weakly contacting, and

<http://www.irsnn.fr/EN/Research/publications-documentation/Publications/PSN-RES/Pages/2012-Topin-collapse-dynamics-run-out-dense-granular-materials-fluid-videos-bis.aspx>



**Figure 8.6:** Granular column collapse in the viscous (bottom) fluid of Figure 8.5 with triple-sized grains.

thus the weight of the mixture is mostly supported by the pore pressure  $p$ , the relative velocity  $\mathbf{w}$  is along  $\nabla p$ . The drag force on the grains  $\xi \mathbf{w}$  will be oriented towards the front of the flow, and will sustain it — this is the flow-inertial regime. On the other hand, when the weight of the granular phase is supported by grain–grain forces,  $\mathbf{w}$  will be along  $-\nabla \text{Tr } \phi \lambda$ , and the drag forces will tend to slow down the grains — this is the viscous regime.

To determine which regime will drive the flow, we can thus look at the ratio between two timescales: the one for the viscous collapse,  $T_v := \eta_f / \rho_f g L$ , where  $L$  is the height of the basin, and the one for the grains to come into contact,  $T_c := D_g / W$ . We have

$$\frac{T_v}{T_c} = \frac{\eta_f W}{\rho_f g L D_g} = \frac{\eta_f \text{St} U}{\rho_f U^2 L \epsilon} = \frac{\text{St}}{\text{Re} \epsilon} = \frac{\alpha + 1}{18} \epsilon.$$

More than the Reynolds number  $\text{Re}$ , the  $(\alpha + 1)\epsilon$  will thus determine the regime of the collapse; the fluid-inertial one will occur when  $T_v/T_c$  is sufficiently small. For our column,  $\epsilon = 10^{-2}$ ; in the dry case,  $\alpha = 2599$  and thus  $T_v/T_c \sim 1$ ; the collapse will therefore be quickly stabilized by the contact forces. However, this number remains constant for the choice of parameters (2.) and (3.),  $\alpha \epsilon \sim 10^{-3}$ . One solution to model a viscous collapse is thus to increase  $\epsilon$ , i.e., consider larger grains. With our simulator, starting from the choice of parameters (3.) and simply tripling the diameters of the particles allows us to significantly reduce the duration and run-out length of the collapse (Figure 8.6).

## Discussion

### 8.4.4 Limitations

It appears that our modeling of the drag term in the interfacial momentum transfer term is insufficient to retrieve the transition between the fluid-inertial and viscous regimes predicted by the discrete simulations. A nonlinear dependency of the drag force on the fluid viscosity for higher concentrations of particles might be necessary.

Another ingredient missing from our model is a dependency on the initial volume fraction of grains, which Pailha et al. (2008) showed had a tremendous influence on the onset of the flow. However, note that when using the Richardson–Zaki drag model, a higher volume fraction will mean that the fluid will oppose more resistance to the initial motion of the grains, and thus will make stronger compact granular heaps.

Our approach is also intrinsically limited to fully immersed granular materials, and thus cannot be applied to the simulation of grains–air–water mixtures, such as wet sand. Finally, our method is not adequate for the simulation of phenomena which happen at a time scale much longer than that of the transport of individual grains, such as the transport of dunes by the wind.

On the numerical side, our diphasic simulation framework obviously suffers from all the discretization issues of the dry case discussed in Chapter 7, with the additional difficulty of having to find an adequate space for the discretization of the pore pressure field — which we sidestep in practice by relaxing the volume-average incompressibility constraint.

#### 8.4.5 Conclusion

Despite all those drawbacks, our approach is still able to recover the qualitative dynamics of immersed granular flows in different regimes. We are able to capture the duality of the fluid role, which may either lubricate the flow or dampen it, following the relative importance of the pore pressure and contact forces. Moreover, we were able to cast the numerical problem that must be solved at each timestep as a standard DCFP with linear constraints, allowing us to leverage once again the large body of research devoted to such systems.

Future work will be focused on better understanding the role of the drag function  $\xi$ , which discriminates between the inertial and viscous regime. Moreover, as already mentioned for the dry case, the search for better discretization spaces constitutes an important research direction for us.



# Conclusion

This dissertation focused on the simulation of complex materials featuring a large number of inelastic contacts with friction. Our main argument is that dealing with these contacts using methods derived from nonsmooth analysis is an efficient alternative to seemingly simpler elasticity- or regularization-based strategies. Indeed, nonsmooth methods allow to decouple the time scale of the macroscopic phenomenon from that of the contact dynamics themselves, the latter being often much smaller for stiff materials and of little practical interest. For instance, one avoids having to explicitly simulate the compression and decompression phases after an impact, and can thus get away with much larger timesteps.

We have also seen that continuum-based simulation methods can provide an efficient manner of simulating very large systems with small inhomogeneities, while still allowing to capture qualitative flow characteristics. Moreover, some numerical methods that have been developed for discrete mechanics can be translated in a straightforward manner to continuum materials. In particular, granular materials may be simulated either at the grain level with the Coulomb friction law or in a macroscopic fashion using the non-associated Drucker–Prager flow rule. Both approaches can leverage the exact same numerical solvers — even though the latter constitutive law is expressed on tensorial rather than vectorial quantities. This legitimizes further the focus that we have put on devising a robust and efficient solver for DCFP.

Finally, we have advocated in very different settings the use of *hybrid* methods, which can take insight from distinct point of views to alleviate the drawbacks of each of them. We have highlighted the benefits of combining an optimization-based approach with an analytical one for solving local problems arising when using splitting methods on DCFP. Similarly, we have seen that using simultaneously mesh-based and particle-based discretization strategies can simplify the treatment of both temporal and spatial differential operators.

## 9.1 Key remarks and summary of contributions

We recall below with more details the main points that we brought up over the course of this thesis, and emphasize the original contributions that we have introduced.

Following the works of Alart and de Saxcé, Chapter 1 presented the Signorini–Coulomb frictional contact law and the non-associated Drucker–Prager flow rule with tension cut-off in an unified manner. It included different ways of characterizing their sets of solutions, either with normal cone inclusions or as roots of complementarity functions. In Chapter 2, we recalled the Moreau–Jean timestepping algorithm, which allows to compute the dynamics of a mechanical system with bilateral constraints and unilateral contact with friction as a sequence of Discrete Coulomb Friction Problems (DCFP). We presented the Cadoux fixed-point algorithm, which yields both a criterion for the existence of solutions to a DCFP, and a practical way of solving this problem as a sequence of *convexified* DCFP (i.e., equivalently, as quadratic minimization problems under conical constraints). Chapter 3 then outlined standard algorithms for solving DCFP and convexified DCFP. Our first original contribution was presented in Chapter 4, and consisted in improving a numerical method for solving DCFP, targeted at ill-conditioned problems with a high number of contacts and for which reaching very high precision was not a requirement. To this aim, we proposed to equip the Gauss–Seidel algorithm with an hybrid local solver combining a novel Newton algorithm on the SOC Fischer–Burmeister function, and an enumerative analytical solver. We derived the expressions of the coefficients of the quartic polynomial that has

to be solved for both DCFP and convexified DCFP. We showed how this algorithm was relevant for applications to the direct and inverse simulation of hair and cloth, and provided benchmarks confirming its good practical performance. However, we noted that on smaller problems, a variant of the projected gradient descent algorithm which we coined ASPG also performed very well in practice. We also proposed a primal-dual algorithm “Dual AMA” requiring only multiplications by the stiffness matrix, but which failed to demonstrate consistent-enough convergence. These different solvers have been released as part of an open-source library developed alongside this thesis, *bogus* (Daviet 2013).

The second part of this dissertation built on our second contribution, presented in Chapter 6, a numerical strategy allowing to cast continuum granular flows problems into the framework of DCFP. Inspired by the work of Cadoux (2009), we showed how finding the solution to the non-associated Drucker–Prager flow rule over a continuous domain amounted once again to a sequence of minimization problems over a SOC in the space of symmetric tensor fields. However, we were only able to state a very weak existence criterion. We then demonstrated that by carefully choosing the basis for the discretization of the strain and stress fields, the discrete flow equations amounted to solving a DCFP (although with higher-dimensional constraints), and that we were able to leverage the solvers designed for standard discrete contact mechanics. We showed that this approach allowed us to capture qualitative features of dense granular flows (some of which being hard to achieve with standard methods, such as the vanishing pressure field in the wake of an obstacle), then extended our method to more general flows in Chapters 7 and 8. In particular, we proposed in Chapter 7 an approximation thanks to which our method proved to be efficient-enough for Computer Graphics applications, allowing us to significantly reduce visual artifacts w.r.t. a state-of-the-art method at a similar computational cost. However, we noted that the choice of suitable discretization spaces remained a notable difficulty. On the other hand, the soundness of our method allowed for easy generalization to more complex settings, be it anisotropic friction, two-way coupling with rigid bodies, or interactions with a surrounding fluid. This last point was the subject of our third main contribution, presented in Chapter 8. Starting from the framework of Jackson (2000), we proposed a two-phase, two-velocity model for the coupled simulation of a Drucker–Prager granular material with a Newtonian fluid. The main interesting feature of our approach is that, once again, we are able to cast the flow dynamics as a DCFP—though, in contrast to the single-phase case, one that includes linear equality constraints. While very limited at the moment, we presented a first few steps towards the qualitative validation of this model, including the capture of the different dynamical regimes governing the collapse of an immersed granular column reported in the literature.

## 9.2 Perspectives

Each step of our numerical method could probably be substantially improved: the choice of discretization spaces, the coupling between the transport and momentum conservation equations, the algorithm for solving the DCFP, *et cætera*. One of our initial goals that was finally excluded from the scope of this thesis was the devising of a massive parallelization strategy taking profit of the modern many-cores architectures. We have good hope that the proximal-based approaches such as ASPG or Dual AMA could scale quite well, even though sparse matrix–vector product is not the most parallelization-friendly routine. Indeed, our first tests in this domain indicated that for regular grids, the variance in the number of non-zeros blocks per row in the DCFP’s matrices was quite small, and that storage schemes such as NVIDIA<sup>®</sup> cuSPARSE’s HYB format<sup>1</sup> were particularly relevant in this case.

From a modeling point of view, we showed that it was possible to express the flow of a granular material, even inside a Newtonian fluid, in a compact, consistent manner and without any regularization as a sequence of numerical problems with a well-known structure. We demonstrated that that approach was able to capture some qualitative features of these coupled dynamics. Despite these first encouraging results, priority for future work should be put towards

<sup>1</sup><http://docs.nvidia.com/cuda/cusparse/#hybrid-format-hyb>

additional validation w.r.t. experimental data or DEM simulations, and towards studying the influence of the several dimensionless numbers than appear in our diphasic model. Moreover, extending the numerical method of Chapter 8 to 3D simulations remains to be done before contemplating applications to Computer Graphics — which could include reproducing underwater sand avalanches or the dynamics of volcanic ash during an eruption.

Finally, while we were able to simulate the dynamics of fibrous materials consisting of a few thousand strands, larger systems still elude our DEM strategy. Moreover, air drag greatly influences the dynamics of hair, and modeling the coupled dynamics of air and hair would probably enhance realism. As coupling a fluid solver with hundreds of thousands of very fine DEM fibers would be extremely costly, continuum modeling of this interaction would seem more appropriate. An interesting research direction would thus be to attempt to extend our macroscopic numerical method to the simulation of hair and fur; however, several key ingredients of such media would still remain to be modeled. First, the bending elasticity of the fibers. One could imagine writing a macroscopic bending energy based on the total derivative the fibers' tangents field, or computing stresses on a discrete number of fiber samples in a MPM fashion. Stretching energy would probably have to be modeled in a similar manner, as the material cannot be macroscopically approximated as inextensible. Fibrous materials are also generally very anisotropic; the basic considerations that we laid down for coin-shaped grains in Chapter 7 would have to be generalized, and the distribution of fiber orientations should influence a much wider range of parameters, from the drag coefficient to the maximal volume fraction. Finally, our granular model cannot capture the natural entanglement of human hair, which induces a kind of *normal* (instead of *shearing*) friction. Applying Dahl's frictional model to this phenomenon would be an interesting starting point. Continuum modeling of slender, elastic fiber assemblies would thus require the introduction of several additional hypothesis and closures. Before attempting such an undertaking, it is therefore essential to get a good understanding of the simpler granular model.





## Appendices



# A Convex analysis

In this section we recall a few fundamental results from convex analysis, stemming for a large part from the early works of Jean-Jacques Moreau and R. Tyrrell Rockafellar. In the following, we mainly refer to the monograph “Fonctionnelles convexes”, as it lists several results in very general settings — e.g., without assumption of finite dimension, even though here we will restrict ourselves to Banach or Hilbert spaces.

## A.1 Operations on convex functions

### A.1.1 Fundamental definitions

Let  $X$  be a reflexive Banach space on  $\mathbb{R}$ . Let  $\bar{\mathbb{R}} := \mathbb{R} \cup \{-\infty, +\infty\}$  denote the extended-real set, and  $f : X \rightarrow \bar{\mathbb{R}}$ .

**Definition A.1.**  $f$  is said to be convex if

$$f(\alpha x + (1 - \alpha)z) \leq \alpha f(x) + (1 - \alpha)f(z) \quad \forall x, z, \alpha \in X \times X \times ]0, 1[$$

where the addition is extended to  $\bar{\mathbb{R}}$  with the commutative convention  $(-\infty) + (+\infty) = +\infty$ .  $f$  is said to be strictly convex if the above definition holds with strict inequality for all  $x \neq z$ .

**Definition A.2** (Epigraph). The epigraph of  $f$  is the set  $\text{epi } f := \{(x, \alpha) \in X \times \mathbb{R}, f(x) \leq \alpha\}$ . It holds that  $f$  is convex if and only if  $\text{epi } f$  is convex.

**Definition A.3** (Effective domain). The effective domain of  $f$  is the set  $\text{dom } f := \{x \in X, f(x) < +\infty\}$ .  $f$  is said to be proper if  $\text{dom } f \neq \emptyset$  and  $\forall x \in X, f(x) > -\infty$ .

**Definition A.4** (Closed function).  $f$  is said to be closed if  $\text{epi } f$  is closed.

**Definition A.5.** [Lower semi-continuity]  $f$  is said to be lower semi-continuous at  $x_0 \in X$  if for all  $\epsilon \in \mathbb{R}_+^*$ , there exists a neighborhood  $V$  of  $x_0$  such that  $\forall x \in V, f(x) \geq f(x_0) - \epsilon$ .

**Property A.1.** The following propositions are equivalent (Moreau 1966–1967, paragraph 4.a):

1.  $f$  is closed
2.  $f$  is lower semi-continuous on  $X$
3.  $\forall \alpha \in \mathbb{R}$ , the set  $\{x \in X, f(x) \leq \alpha\}$  is closed.

**Property A.2.** Any proper closed convex function  $f$  is continuous on the interior of its effective domain (Moreau 1966–1967, paragraph 5.f).

### A.1.2 Subdifferential of a function

Let  $Y$  be the topological dual space of  $X$  (the space of continuous linear forms  $y : X \rightarrow \mathbb{R}$ ), with the bilinear form  $\langle \cdot, \cdot \rangle : X \times Y \rightarrow \mathbb{R}$  such that  $Y$  can be identified to  $\{\langle \cdot, y \rangle, y \in Y\}$  and  $X$  to  $\{\langle x, \cdot \rangle, x \in X\}$ .

**Definition A.6** (Subdifferential). *If there exists  $\mathbf{x}_0 \in X$  and  $\mathbf{y} \in Y$  such that*

$$\forall \mathbf{x} \in X, \quad f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle$$

*then  $f$  is said to be subdifferentiable at  $\mathbf{x}_0$  and  $\mathbf{y}$  is a subgradient of  $f$  at  $\mathbf{x}_0$ . We denote by  $\partial f(\mathbf{x}_0)$  the subdifferential of  $f$  at  $\mathbf{x}_0$ , defined as the set of all subgradients at  $f$  at  $\mathbf{x}_0$ .*

**Remark A.1.** *If  $f$  is proper,  $\partial f(X \setminus \text{dom } f) = \emptyset$ .*

It follows from Definition A.6 that:

**Property A.3** (Minimum of a convex function). *The global minimum of  $f$  is attained at  $\mathbf{x}_0 \in X$  if and only if  $\mathbf{0} \in \partial f(\mathbf{x}_0)$ . Moreover, if  $f$  is strictly convex, this global minimum is attained at at most one point.*

*Proof.*  $\mathbf{0} \in \partial f(\mathbf{x}_0) \iff f(\mathbf{x}) \geq f(\mathbf{x}_0) \forall \mathbf{x} \in X$ , which is the definition of a global minimum. To prove the second assertion, suppose that  $f$  is strictly convex and attains its global minimum at two different points  $\mathbf{x}_1 \neq \mathbf{x}_2$ . Then  $f(\frac{1}{2}(\mathbf{x}_1 + \mathbf{x}_2)) < \frac{1}{2}(f(\mathbf{x}_1) + f(\mathbf{x}_2)) = f(\mathbf{x}_1)$ , which yields a contradiction.  $\square$

The following property suggests that the subdifferential can be seen as a generalization of the concept of the gradient of a convex function at points where it is not differentiable.

**Property A.4.** *Let  $f$  be a proper convex function on  $X$ , continuous at  $\mathbf{x}_0 \in \text{dom}(f)$ . The following statements are equivalent (Moreau 1966–1967, paragraphs 10.f and 10.g)*

- *$f$  is Gâteaux-differentiable at  $\mathbf{x}_0$  with gradient  $(\nabla f)(\mathbf{x}_0) \in Y$ , i.e.,*

$$\lim_{\alpha \rightarrow 0_+} \frac{f(\mathbf{x}_0 + \alpha \mathbf{x}) - f(\mathbf{x}_0)}{\alpha} = \langle \mathbf{x}, (\nabla f)(\mathbf{x}_0) \rangle \quad \forall \mathbf{x} \in X;$$

- $\partial f(\mathbf{x}_0) = \{(\nabla f)(\mathbf{x}_0)\}$ .

**Remark A.2.** *If  $f$  is Fréchet-differentiable, then it is continuous and Gâteaux-differentiable as well, and the gradients from both definitions coincide. The converse is not true; intuitively Fréchet ensures that the differential is well-defined along all continuous paths converging to  $\mathbf{x}_0$ , while Gâteaux considers only straight lines.*

**Sum of differentials** It always hold that  $\partial f + \partial g \subset \partial(f + g)$ , but some regularity (a.k.a. qualification) conditions are required for the equality to be achieved. Moreau (1966–1967, Proposition 10.7) states the following theorem:

**Theorem A.1** (Subdifferential of a sum). *For  $f$  and  $g$  convex  $X \rightarrow \bar{\mathbb{R}}$ , if there exists  $\mathbf{x}_0 \in \text{dom } f \cap \text{dom } g$  such that either  $f$  or  $g$  is continuous at  $\mathbf{x}_0$ , then  $\partial(f + g) = \partial f + \partial g$  on  $X$ .*

As a proper convex function is continuous on the interior of its effective domain (Property A.2), it follows:

**Corollary A.1.** *Let  $f$  and  $g$  be convex functions  $X \rightarrow \bar{\mathbb{R}}$ , with  $f$  closed and proper. If  $(\text{int dom } f) \cap \text{dom } g \neq \emptyset$ , then  $\partial(f + g) = \partial f + \partial g$  on  $X$ .*

**Special cases** The following property will prove useful in continuum mechanics.

**Property A.5** (Subdifferential in function spaces). *Let  $\Omega$  be a subset of  $\mathbb{R}^d$ ,  $d \in \mathbb{N}$ , and let  $L_2(\Omega)$  denote the space of square-integrable functions on  $\Omega$ . Let  $\langle \cdot, \cdot \rangle$  be a scalar product on  $\mathbb{R}^n$ ,  $n \in \mathbb{N}$ ; we equip the space  $V := L_2(\Omega)^n$  with the scalar product  $\langle \cdot, \cdot \rangle_V$ ,*

$$\langle \mathbf{u}, \mathbf{v} \rangle_V := \int_{\Omega} \langle \mathbf{u}(\mathbf{x}), \mathbf{v}(\mathbf{x}) \rangle d\mathbf{x} \quad \forall (\mathbf{u}, \mathbf{v}) \in V^2.$$

Finally, consider a function  $g : V \rightarrow \bar{\mathbb{R}}, \mathbf{u} \mapsto \int_{\Omega} f(\mathbf{u}(\mathbf{x})) d\mathbf{x}$ , where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . The subdifferential of  $g$  is given by, for any  $\mathbf{u} \in V$ ,

$$\partial g(\mathbf{u}) = \{\mathbf{v} \in V, \mathbf{v}(\mathbf{x}) \in \partial f(\mathbf{u}(\mathbf{x})) \text{ a.e. on } \Omega\}.$$

*Proof.*  $V$  is an Hilbert space, thus  $\partial g \subset V$ . Let  $\mathbf{u}, \mathbf{v} \in V^2$ ,  $\Omega_{\mathbf{v}}^+ := \{\mathbf{x} \in \Omega, \mathbf{v}(\mathbf{x}) \in \partial f(\mathbf{u}(\mathbf{x}))\}$  and  $\Omega_{\mathbf{v}}^- := \Omega - \Omega_{\mathbf{v}}^+$ . By definition of the subdifferential,  $\mathbf{x} \in \Omega_{\mathbf{v}}^+ \iff \forall \mathbf{z} \in \mathbb{R}^d, f(\mathbf{z}) - f(\mathbf{u}(\mathbf{x})) \geq \langle \mathbf{v}(\mathbf{x}), \mathbf{z} - \mathbf{u}(\mathbf{x}) \rangle$ .

For any  $\mathbf{w} \in V$ , we thus have

$$\begin{aligned} g(\mathbf{w}) - g(\mathbf{u}) &= \int_{\Omega_{\mathbf{v}}^+} f(\mathbf{w}(\mathbf{x})) - f(\mathbf{u}(\mathbf{x})) + \int_{\Omega_{\mathbf{v}}^-} f(\mathbf{w}(\mathbf{x})) - f(\mathbf{u}(\mathbf{x})) \\ &\geq \int_{\Omega_{\mathbf{v}}^+} \langle \mathbf{v}(\mathbf{x}), \mathbf{w}(\mathbf{x}) - \mathbf{u}(\mathbf{x}) \rangle + \int_{\Omega_{\mathbf{v}}^-} f(\mathbf{w}(\mathbf{x})) - f(\mathbf{u}(\mathbf{x})). \end{aligned}$$

As  $f$  takes its values in  $\mathbb{R}$ ,  $\Omega_{\mathbf{v}}^-$  being of measure zero implies that  $g(\mathbf{w}) - g(\mathbf{u}) \geq \langle \mathbf{v}, \mathbf{w} - \mathbf{u} \rangle_V$ , and thus  $\mathbf{v} \in \partial g(\mathbf{u})$ . Reciprocally, suppose that  $\Omega_{\mathbf{v}}^-$  is not of measure zero. By definition of this set, for all  $\mathbf{x} \in \Omega_{\mathbf{v}}^-$  there exists  $\mathbf{z}_{\mathbf{x}} \in \mathbb{R}^d$  such that  $f(\mathbf{z}_{\mathbf{x}}) - f(\mathbf{u}(\mathbf{x})) < \langle \mathbf{v}(\mathbf{x}), \mathbf{z}_{\mathbf{x}} - \mathbf{u}(\mathbf{x}) \rangle$ . Let us define  $\mathbf{w} \in V$  as

$$\mathbf{w}(\mathbf{x}) := \begin{cases} \mathbf{z}_{\mathbf{x}} & \text{on a relatively compact subset of } \Omega_{\mathbf{v}}^- \text{ with strictly positive measure} \\ \mathbf{u}(\mathbf{x}) & \text{elsewhere.} \end{cases}$$

Then  $g(\mathbf{w}) - g(\mathbf{u}) = \int_{\Omega_{\mathbf{v}}^-} \mathbf{z}_{\mathbf{x}} - f(\mathbf{u}(\mathbf{x})) < \int_{\Omega_{\mathbf{v}}^-} \langle \mathbf{v}(\mathbf{x}), \mathbf{z}_{\mathbf{x}} - \mathbf{u}(\mathbf{x}) \rangle = \langle \mathbf{v}, \mathbf{w} - \mathbf{u} \rangle_V$ , i.e.,  $\mathbf{v} \notin \partial g(\mathbf{u})$ .  $\square$

Finally, we state a property of subdifferentials that will not be directly useful in this dissertation, but allows the easy extension of results derived for a scalar constraint to larger dimensional ones. The proof is given for a finite-dimensional space in (Hiriart-Urruty and Lemaréchal 1993, VI, Corollary 4.3.2), and extended to infinite sequences in (Hiriart-Urruty and Lemaréchal 1993, VI, Theorem 4.4.2).

**Property A.6** (Subdifferential of a maximum). *Let  $f_i, i = 1 \dots n$  be a finite number of convex functions on  $X$ . Then the subdifferential of their point-wise maximum at a point  $\mathbf{x}_0$  is the convex hull of the subdifferentials of the  $f_i$  that are active at  $\mathbf{x}_0$ , i.e.,*

$$\partial(\max(f_i))(\mathbf{x}_0) = \text{Conv} \bigcup_{j \in J} \partial f_j \text{ with } J := \{j, f_j(\mathbf{x}_0) = \max_i f_i(\mathbf{x}_0)\}.$$

### A.1.3 Convex conjugate

**Definition A.7** (Convex conjugate). *The convex conjugate (or Fenchel–Legendre transform) of  $f$  is the application*

$$\begin{aligned} f^* : Y &\rightarrow \bar{\mathbb{R}} \\ \mathbf{y} &\mapsto \sup_{\mathbf{x} \in X} (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x})). \end{aligned}$$

Moreau (1966–1967, paragraphs 6.b and 6.d) gives two fundamental properties:

**Property A.7.** *The convex conjugate of  $f$  is always convex and closed. If  $f$  is proper, convex and closed, then it is equal to its biconjugate, i.e.,  $f = f^{**}$ .*

**Theorem A.2** (Subnormality). *It always hold that  $f(\mathbf{x}) + f^*(\mathbf{y}) \geq \langle \mathbf{x}, \mathbf{y} \rangle$  (Fenchel–Young inequality). Moreover, if  $f$  is proper, closed and convex, then*

$$\mathbf{y} \in \partial f(\mathbf{x}) \iff f(\mathbf{x}) + f^*(\mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle \iff \mathbf{x} \in \partial f^*(\mathbf{y}).$$

*Proof.* From the definition of the convex conjugate,  $\forall \mathbf{x}, \mathbf{y} \in X \times Y$ ,

$$\begin{aligned} f(\mathbf{x}) + f^*(\mathbf{y}) &= f(\mathbf{x}) + \sup_{\mathbf{z} \in X} (\langle \mathbf{z}, \mathbf{y} \rangle - f(\mathbf{z})) \\ &\geq f(\mathbf{x}) + (\langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x})) = \langle \mathbf{x}, \mathbf{y} \rangle \end{aligned}$$

Then,

$$\begin{aligned} f(\mathbf{x}) + f^*(\mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle &\iff f^*(\mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}) \\ &\iff \langle \mathbf{z}, \mathbf{y} \rangle - f(\mathbf{z}) \leq \langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}) \quad \forall \mathbf{z} \in X \\ &\iff f(\mathbf{x}) + \langle \mathbf{z} - \mathbf{x}, \mathbf{y} \rangle \leq f(\mathbf{z}) \quad \forall \mathbf{z} \in X \\ &\iff \mathbf{y} \in \partial f(\mathbf{x}) \end{aligned}$$

The second part of the equivalence stems from  $f^{**} = f$  (Property A.7).  $\square$

We transcribe below a slightly weakened version of the Fenchel duality theorem as given in (Borwein and Luke 2011, theorem 1.2),

**Theorem A.3** (Fenchel duality). *Let  $X$  and  $E$  be Banach spaces, with  $E^*$  the dual of  $E$ . Let  $L$  be a bounded linear map.*

*Then the extrema of the optimization problems (A.1) and (A.2),*

$$p = \inf_{\mathbf{x} \in X} f(\mathbf{x}) + g(L\mathbf{x}) \quad (\text{A.1})$$

$$d = \sup_{\mathbf{z} \in E^*} -f^*(L^\top \mathbf{z}) - g^*(-\mathbf{z}) \quad (\text{A.2})$$

*satisfy weak duality, i.e.,  $p \geq d$ .*

*Moreover if  $f$  and  $g$  are convex and if there exists  $\mathbf{x} \in L \text{ dom } f$  such that  $g$  is finite and continuous at  $\mathbf{x}_0$ , strong duality holds:  $p = d$ , and if they are finite, then the dual problem (A.2) attains its maximum.*

**Corollary A.2.** *Extension to affine maps;  $L$  in the theorem above is replaced by  $A : \mathbf{x} \mapsto L\mathbf{x} + \mathbf{e}$ , and the optimization problems (A.1) and (A.2) with*

$$\begin{aligned} p &= \inf_{\mathbf{x} \in X} f(\mathbf{x}) + (g \circ A)(\mathbf{x}) \\ d &= \sup_{\mathbf{z} \in E^*} -f^*(L^\top \mathbf{z}) - g^*(-\mathbf{z}) - \langle \mathbf{e}, \mathbf{z} \rangle \end{aligned}$$

*Weak duality always hold, and strong duality is achieved when  $f$  and  $g$  are convex and there exists  $\mathbf{x} \in A \text{ dom } f$  such that  $g$  is finite and continuous at  $\mathbf{x}_0$ .*

*Proof.* By application of the theorem to  $g_e : \mathbf{z} \mapsto g(\mathbf{z} + \mathbf{e})$ ; indeed  $g_e(L\mathbf{x}) = g \circ A(\mathbf{x})$ , and  $(^*g_e)(-\mathbf{z}) = (^*g)(-\mathbf{z}) + \langle \mathbf{e}, \mathbf{z} \rangle$ .  $\square$

## A.2 Normal and convex cones

### A.2.1 Normal cone

**Definition A.8** (Characteristic function). *The characteristic function  $\mathcal{J}_C$  of a set  $C \subset X$  is defined as vanishing on  $C$  and taking infinite values outside of  $C$ ,*

$$\mathcal{J}_C := \begin{cases} 0 & \text{if } \mathbf{x} \in C \\ +\infty & \text{if } \mathbf{x} \notin C. \end{cases}$$

From  $\text{epi } \mathcal{J}_C = C \times \mathbb{R}^+$  and  $\text{dom } \mathcal{J}_C = C$  it follows:

- $\mathcal{J}_C$  is convex if and only if  $C$  is convex.

- $\mathcal{J}_C$  is not empty if and only if  $C$  is proper.
- $\mathcal{J}_C$  is closed if and only if  $C$  is closed.

**Definition A.9** (Normal cone). *Let  $C$  be a convex subset of  $X$ . The normal cone to  $C$  at  $\mathbf{x}_0 \in X$  is defined as*

$$\mathcal{N}_C(\mathbf{x}_0) = \begin{cases} \{\mathbf{y} \in Y, \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle \leq 0 \quad \forall \mathbf{x} \in C\} & \text{if } \mathbf{x}_0 \in C \\ \emptyset & \text{if } \mathbf{x}_0 \notin C \end{cases}$$

As it can easily be seen that for any point  $\mathbf{x}_0$  in  $C$ ,  $\mathbf{0} \in \mathcal{N}_C(\mathbf{x}_0)$ , we have the equivalence:

**Property A.8.**  $\mathbf{0} \in \mathcal{N}_C(\mathbf{x}) \iff \mathbf{x} \in C$

**Property A.9.** *If  $C$  is not empty, there holds  $\mathcal{N}_C = \partial \mathcal{J}_C$ .*

*Proof.* As  $\mathcal{J}_C$  is proper, from Remark A.1 we have  $\mathbf{x}_0 \notin C \iff \partial \mathcal{J}_C(\mathbf{x}_0) = \emptyset = \mathcal{N}_C(\mathbf{x}_0)$ . For  $\mathbf{x}_0 \in C$ , Definition A.6 reads

$$\begin{aligned} \mathbf{y} \in \partial \mathcal{J}_C(\mathbf{x}_0) &\iff \mathcal{J}_C(\mathbf{x}) \geq \mathcal{J}_C(\mathbf{x}_0) + \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle & \forall \mathbf{x} \in X \\ &\iff \mathcal{J}_C(\mathbf{x}) \geq \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle & \forall \mathbf{x} \in X \\ &\iff 0 \geq \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle & \forall \mathbf{x} \in C \\ &\iff \mathbf{y} \in \mathcal{N}_C(\mathbf{x}_0) \end{aligned}$$

□

We deduce immediately that the expression of the normal cone at an interior point.

**Property A.10.**  $\mathbf{x} \in \text{int } C \implies \mathcal{N}_C(\mathbf{x}) = \{\mathbf{0}\}$ .

*Proof.* For  $\mathbf{x} \in \text{int } C$ , there exists a closed ball with strictly positive radius centered at  $\mathbf{x}$  on which  $\mathcal{J}_C$  is constant. Therefore  $\mathcal{J}_C$  is continuous and Gâteaux-differentiable at  $\mathbf{x}$ , with  $(\nabla \mathcal{J}_C)(\mathbf{x}_0) = \mathbf{0}$ . Properties A.4 and A.9 conclude the proof. □

### A.2.2 Operations on normal cones

**Property A.11** (Normal cone to an intersection). *For any two subsets  $C_1$  and  $C_2$  of  $X$ ,  $\mathcal{N}_{C_1 \cap C_2} \supset \mathcal{N}_{C_1} + \mathcal{N}_{C_2}$ .*

*Moreover, if  $C_1$  and  $C_2$  are convex, with  $C_1$  closed and  $(\text{int } C_1) \cap C_2 \neq \emptyset$ , then  $\mathcal{N}_{C_1 \cap C_2} = \mathcal{N}_{C_1} + \mathcal{N}_{C_2}$ .*

*Proof.* If  $C_1$  and  $C_2$  are empty, the inclusion is trivial. Otherwise, we have  $\mathcal{J}_{C_1 \cap C_2} = \mathcal{J}_{C_1} + \mathcal{J}_{C_2}$ , therefore with Property A.9,

$$\mathcal{N}_{C_1 \cap C_2} = \partial(\mathcal{J}_{C_1} + \mathcal{J}_{C_2}) \supset \partial \mathcal{J}_{C_1} + \partial \mathcal{J}_{C_2} = \mathcal{N}_{C_1} + \mathcal{N}_{C_2}.$$

Now, if  $C_1$  and  $C_2$  are convex with  $C_1$  closed and  $(\text{int } C_1) \cap C_2 \neq \emptyset$ , the regularity conditions of Corollary A.1 on the subdifferential of a sum are satisfied, and the equality holds. □

Corollaries of the following property will often prove useful for the practical computation of normal cones.

**Property A.12** (Subdifferential of the precomposition by an affine map). *Let  $X$  and  $E$  be Hilbert spaces, and  $A$  an affine map from  $X$  to  $E$ , i.e.,  $A(\mathbf{x}) = L\mathbf{x} + \mathbf{e}$  with  $L$  linear and  $\mathbf{e} \in E$ . We note  $L^\top$  the linear adjoint of  $L$ . Let  $f$  be a convex function on  $E$ . There holds, for all  $\mathbf{x} \in X$ ,*

$$\partial(f \circ A)(\mathbf{x}) \supset L^\top \partial f(A(\mathbf{x})).$$

*Moreover, if  $f$  is proper and closed, and  $\exists \mathbf{x}_0 \in \text{dom}(f \circ A)$  such that  $f$  is continuous at  $A(\mathbf{x}_0)$ , then for all  $\mathbf{x} \in X$ ,*

$$\partial(f \circ A)(\mathbf{x}) = L^\top \partial f(A(\mathbf{x})).$$



*Proof.* First, let  $\mathbf{y} \in L^\top \partial f(A(\mathbf{x}_0))$ . There exists  $\mathbf{g} \in E$  such that  $\mathbf{y} = L^\top \mathbf{g}$  and

$$f(\mathbf{z}) \geq f(A(\mathbf{x}_0)) + \langle \mathbf{z} - A(\mathbf{x}_0), \mathbf{g} \rangle \quad \forall \mathbf{z} \in E.$$

In particular,

$$f(A(\mathbf{x})) \geq f(A(\mathbf{x}_0)) + \langle A(\mathbf{x}) - A(\mathbf{x}_0), \mathbf{g} \rangle \quad \forall \mathbf{x} \in X.$$

Since  $\langle A(\mathbf{x}) - A(\mathbf{x}_0), \mathbf{g} \rangle = \langle L\mathbf{x} - L\mathbf{x}_0, \mathbf{g} \rangle = \langle \mathbf{x} - \mathbf{x}_0, L^\top \mathbf{g} \rangle$ ,

$$f(A(\mathbf{x})) \geq f(A(\mathbf{x}_0)) + \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle \quad \forall \mathbf{x} \in X,$$

which means that  $\mathbf{y} \in \partial(f \circ A)(\mathbf{x}_0)$ .

Now, let us prove the reverse inclusion under our supplemental regularity conditions. Let  $\mathbf{y} \in \partial(f \circ A)(\mathbf{x}_0)$ , and let  $\mathbf{x}^* \in \text{dom}(f \circ A)$ , which is not empty under our hypothesis. Then,  $\forall \mathbf{x} \in \text{Ker } L$ ,

$$\begin{aligned} f(A(\mathbf{x}^* + \mathbf{x})) &\geq f(A(\mathbf{x}_0)) + \langle \mathbf{x}^* + \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle \\ &\quad - \langle \mathbf{x}, \mathbf{y} \rangle \geq f(A(\mathbf{x}_0)) - f(A(\mathbf{x}^*)) + \langle \mathbf{x}^* - \mathbf{x}_0, \mathbf{y} \rangle. \end{aligned}$$

As the right-hand side is finite, for the inequality to hold for all  $\mathbf{x} \in \text{Ker } L$ , we must have  $\mathbf{y} \in (\text{Ker } L)^\perp = \text{Im } L^\top$ , i.e.,  $\exists \mathbf{g} \in E$ ,  $\mathbf{y} = L^\top \mathbf{g}$ . Going back to the definition of the subdifferential, we must also have

$$\begin{aligned} f(A(\mathbf{x})) &\geq f(A(\mathbf{x}_0)) + \langle \mathbf{x} - \mathbf{x}_0, L^\top \mathbf{g} \rangle & \forall \mathbf{x} \in X \\ f(A(\mathbf{x})) &\geq f(A(\mathbf{x}_0)) + \langle A(\mathbf{x}) - A(\mathbf{x}_0), \mathbf{g} \rangle & \forall \mathbf{x} \in X \\ f(\mathbf{z}) &\geq f(A(\mathbf{x}_0)) + \langle \mathbf{z} - A(\mathbf{x}_0), \mathbf{g} \rangle & \forall \mathbf{z} \in \text{Im } A \\ (f + \mathcal{J}_{\text{Im } A})(\mathbf{z}) &\geq (f + \mathcal{J}_{\text{Im } A})(A(\mathbf{x}_0)) + \langle \mathbf{z} - A(\mathbf{x}_0), \mathbf{g} \rangle & \forall \mathbf{z} \in E, \end{aligned}$$

which means  $\mathbf{g} \in \partial(f + \mathcal{J}_{\text{Im } A})(A(\mathbf{x}_0))$ . Under our regularity conditions, we can apply Theorem A.1 and get equivalently  $\mathbf{g} \in \partial f(A(\mathbf{x}_0)) + \partial \mathcal{J}_{\text{Im } A}(A(\mathbf{x}_0))$ . As  $\partial \mathcal{J}_{\text{Im } A}(A(\mathbf{x}_0)) = \mathcal{N}_{\text{Im } A}(A(\mathbf{x}_0)) = \text{Ker } L^\top$ , we get

$$\mathbf{y} \in L^\top (\partial f(A(\mathbf{x}_0)) + \text{Ker } L^\top) = L^\top \partial f(A(\mathbf{x}_0)).$$

□

We can get the following corollary by choosing  $f = \mathcal{J}_C$  in the above property:

**Corollary A.3.** *Let  $C$  be a closed convex subset of  $E$  such that  $\text{int } C \cap \text{Im } A \neq \emptyset$ . Then the normal cone to  $\Gamma := \{\mathbf{x} \in X, A(\mathbf{x}) \in C\}$  is*

$$\mathcal{N}_\Gamma(\mathbf{x}) = L^\top \mathcal{N}_C(A(\mathbf{x}))$$

The extension to linear constraints stated below will also prove useful for mechanical systems.

**Corollary A.4.** *Let  $A : X \rightarrow E$  be an affine map defined as in Property A.12, i.e.,  $A : \mathbf{x} \mapsto L\mathbf{x} + \mathbf{e}$ ,  $B : X \rightarrow F$  a linear operator from  $X$  to  $F$  Hilbert space, and for any  $\mathbf{f} \in F$ ,  $V(\mathbf{f}) := \{\mathbf{x} \in X, B\mathbf{x} = \mathbf{f}\}$ . Let  $C$  be a closed convex subset of  $E$ , then the normal cone to  $\Gamma(\mathbf{f}) := \{\mathbf{x} \in V(\mathbf{f}), A(\mathbf{x}) \in C\}$  satisfies*

$$\mathcal{N}_{\Gamma(\mathbf{f})}(\mathbf{x}) \supset L^\top \mathcal{N}_C(A(\mathbf{x})) + B^\top \mathcal{N}_{\{\mathbf{0}_F\}}(B\mathbf{x} - \mathbf{f}),$$

with equality under the regularity condition

$$\text{int } C \cap A(V(\mathbf{f})) \neq \emptyset. \quad (\text{A.3})$$

*Proof.* The inclusion being trivial, we focus on the equality. The regularity condition (A.3) requires that  $V(\mathbf{f})$  is not empty; let  $\mathbf{v}_f \in V(\mathbf{f})$ , then  $V(\mathbf{f}) = V(\mathbf{0}_F) + \mathbf{v}_f$ .

$V(\mathbf{0}_F) = \text{Ker } B$  is a linear subspace of  $X$  and is therefore a Hilbert space. As the regularity condition (A.3) can be written

$$\text{int } C \cap (LV(\mathbf{0}_F) + \mathbf{e} + L\mathbf{v}_f) \neq \emptyset$$

we can apply Corollary A.3 to the linear map  $A_0 : V(\mathbf{0}_F) \rightarrow E$ ,  $\mathbf{x} \mapsto L\mathbf{x} + (\mathbf{e} + L\mathbf{v}_f)$ , and get that the normal cone to  $\Gamma_0(f) \subset V(\mathbf{0}_F) := \{\mathbf{x} \in V(\mathbf{0}_F), A_0(\mathbf{x}) \in C\}$  is

$$\mathcal{N}_{\Gamma_0(f)}(\mathbf{x}) = L^\top \mathcal{N}_C(A_0(\mathbf{x})).$$

Now, suppose  $\mathbf{y} \in \mathcal{N}_{\Gamma(f)}(\mathbf{x})$ . Necessarily  $\mathbf{x} \in \Gamma(f)$ , and thus  $\mathbf{x} - \mathbf{v}_f \in V(\mathbf{0}_F)$ . Moreover, we can decompose  $\mathbf{y}$  as  $\mathbf{y} = \mathbf{y}_{|V(\mathbf{0}_F)} + \mathbf{y}_\perp$ , with  $\mathbf{y}_{|V(\mathbf{0}_F)} \in V(\mathbf{0}_F)$  and  $\mathbf{y}_\perp \in \text{Im } B^\top$ . Therefore,

$$\begin{aligned} \forall \mathbf{z} \in \Gamma(f), & \quad \langle \mathbf{y}, \mathbf{z} - \mathbf{x} \rangle \leq 0 \\ \forall \mathbf{z} \in \Gamma_0(f), & \quad \langle \mathbf{y}, \mathbf{z} - (\mathbf{x} - \mathbf{v}_f) \rangle \leq 0 \\ \forall \mathbf{z} \in \Gamma_0(f), & \quad \langle \mathbf{y}_{|V(\mathbf{0}_F)}, \mathbf{z} - (\mathbf{x} - \mathbf{v}_f) \rangle \leq 0, \end{aligned}$$

and necessarily  $\mathbf{y} \in \mathcal{N}_{\Gamma_0}(\mathbf{x} - \mathbf{v}_f) + B^\top \mathcal{N}_{\{0\}}(B\mathbf{x} - f)$ . As the reverse inclusion always holds from Property A.11 on the normal cone to an intersection and Corollary A.3 applied to the linear constraint, we get

$$\begin{aligned} \mathcal{N}_{\Gamma(f)}(\mathbf{x}) &= L^\top \mathcal{N}_C(A_0(\mathbf{x} - \mathbf{v}_f)) + B^\top \mathcal{N}_{\{0\}}(B\mathbf{x} - f) \\ &= L^\top \mathcal{N}_C(A(\mathbf{x})) + B^\top \mathcal{N}_{\{0\}}(B\mathbf{x} - f). \end{aligned}$$

□

### A.2.3 Convex cones

The name *normal cone* is coined from the fact that its belong to the class of *convex cones*, which possesses interesting properties w.r.t. the convex conjugate.

**Definition A.10** (Convex cone). *A subset  $K \subset X$  is a convex cone if*

$$\mathbf{x}, \mathbf{y} \in K^2 \implies (\alpha\mathbf{x} + \beta\mathbf{y}) \in K, \quad \forall \alpha, \beta \in \mathbb{R}_+^2.$$

**Property A.13.** *There holds:*

1. *For any  $\mathbf{x} \in C \subset X$ ,  $\mathcal{N}_C(\mathbf{x})$  is a convex cone.*
2. *If  $K$  is a convex cone, then  $\mathbf{0} \in K \iff K \neq \emptyset$ .*

**Definition A.11** (Dual and polar cones). *The polar cone of  $C \subset X$  is the set*

$$C^\circ := \{\mathbf{y} \in Y, \langle \mathbf{x}, \mathbf{y} \rangle \leq 0 \quad \forall \mathbf{x} \in C\}.$$

*The dual cone of  $C$  is then defined as the negative of the polar cone,  $C^* := -C^\circ$ .*

**Property A.14** (Normal cone to a convex cone). *If  $K$  is a convex cone, then*

$$\mathcal{N}_K(\mathbf{x}) = \begin{cases} \emptyset & \text{if } \mathbf{x} \notin K \\ K^\circ \cap \{\mathbf{x}\}^\perp & \text{if } \mathbf{x} \in K. \end{cases}$$

*Proof.* The first case holds for any normal cone (Property A.8). Let us prove the second one.

First, we choose  $\mathbf{x} \in K$  and  $\mathbf{y} \in \mathcal{N}_K(\mathbf{x})$ , i.e.,  $\forall \mathbf{z} \in K, \langle \mathbf{z} - \mathbf{x}, \mathbf{y} \rangle \leq 0$ . Since  $\mathbf{0} \in K$ ,  $\langle \mathbf{0} - \mathbf{x}, \mathbf{y} \rangle \leq 0$ , and since  $2\mathbf{x} \in K$  as well,  $\langle \mathbf{x}, \mathbf{y} \rangle \leq 0$ , therefore  $\langle \mathbf{x}, \mathbf{y} \rangle = 0$ , i.e.,  $\mathbf{y} \in \{\mathbf{x}\}^\perp$ .

Moreover,  $\forall \mathbf{z} \in K, (\mathbf{x} + \mathbf{z}) \in K$ , meaning that  $\forall \mathbf{z} \in K, \langle \mathbf{z}, \mathbf{y} \rangle \leq 0$ , i.e.,  $\mathbf{y} \in K^\circ$ . Necessarily, we therefore have  $\mathcal{N}_K(\mathbf{x}) \subset K^\circ \cap \{\mathbf{x}\}^\perp$ . Let us now prove that the inclusion is not strict. Let  $\mathbf{y} \in K^\circ \cap \{\mathbf{x}\}^\perp$ , and  $\mathbf{z} \in K$ . Then  $\langle \mathbf{z} - \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{z}, \mathbf{y} \rangle - \langle \mathbf{x}, \mathbf{y} \rangle \leq 0$  from the definitions of the polar cone and of the orthogonal subspace. Since this is true  $\forall \mathbf{z} \in K, \mathbf{y} \in \mathcal{N}_K(\mathbf{x})$ . Therefore,  $\forall \mathbf{x} \in K, \mathcal{N}_K(\mathbf{x}) = K^\circ \cap \{\mathbf{x}\}^\perp$ . □

**Property A.15.** *For  $K$  a non-empty convex cone,  $(\mathcal{N}_K)^* = \mathcal{N}_{K^\circ}$ .*

*Proof.* Following from the definition of the convex conjugate,

$$\begin{aligned}\mathcal{J}_K^*(y) &= \sup_{x \in X} (\langle x, y \rangle - \mathcal{J}_K(x)) \\ &= \sup_{x \in K} (\langle x, y \rangle)\end{aligned}$$

As  $K$  is a cone,  $(0x) \in K$  and therefore  $\mathcal{J}_K^*(y) \geq 0$ . Moreover if  $\exists x_0 \in K, \langle x_0, y \rangle > 0$ , then  $\forall \beta \in \mathbb{R}_+$ , we can choose  $\alpha = \frac{\beta}{\langle x_0, y \rangle}$  so that  $(\alpha x_0) \in K$  and  $\langle \alpha x_0, y \rangle \geq \beta$ , i.e.,  $\mathcal{J}_K^*(y) = +\infty$ .

This means

$$\mathcal{J}_K^*(y) = \begin{cases} +\infty & \text{if } \exists x \in K, \langle x, y \rangle > 0 \\ 0 & \text{otherwise} \end{cases} = \begin{cases} +\infty & \text{if } x \notin K^\circ \\ 0 & \text{otherwise} \end{cases} = \mathcal{J}_{K^\circ}.$$

□

**Corollary A.5.** *A non-empty closed convex cone  $K$  is the dual (resp., polar) cone of its dual (resp., polar) cone, i.e.,  $K^{**} = K$  and  $K^{\circ\circ} = K$ .*

*Proof.*

$$\begin{aligned}x \in K^{\circ\circ} &\iff x \in \mathcal{N}_{K^\circ}(0) && \text{using Property A.14} \\ &\iff x \in \partial \mathcal{J}_{K^\circ}(0) && \text{using Property A.9} \\ &\iff x \in \partial \mathcal{J}_K^*(0) && \text{using Property A.15} \\ &\iff 0 \in \partial \mathcal{J}_K^{**}(x) && \text{using Theorem A.2} \\ &\iff 0 \in \partial \mathcal{J}_K(x) && \text{using Property A.7} \\ &\iff 0 \in \mathcal{N}_K(x) && \text{using Property A.9} \\ &\iff x \in K && \text{using Property A.8.}\end{aligned}$$

□

**Theorem A.4** (Conic complementarity). *For  $K$  a non-empty closed convex cone, there holds*

$$y \in \mathcal{N}_K(x) \iff K \ni x \perp y \in K^\circ \iff x \in \mathcal{N}_{K^\circ}(y)$$

where the “ $\cdot \perp \cdot$ ” is used to mean  $\langle \cdot, \cdot \rangle = 0$ .

*Proof.* From Theorem A.2 and Properties A.9 and A.15,

$$\begin{aligned}y \in \mathcal{N}_K(x) &\iff y \in \partial \mathcal{J}_K(x) \\ &\iff \mathcal{J}_K(x) + \mathcal{J}_K^*(y) = \langle x, y \rangle \\ &\iff \mathcal{J}_K(x) + \mathcal{J}_{K^\circ}(y) = \langle x, y \rangle\end{aligned}$$

The equality is only possible on the effective domain of the left-hand side, on which  $\mathcal{J}_K(x) + \mathcal{J}_{K^\circ}(y)$  is constant and equal to zero. This means

$$\mathcal{J}_K(x) + \mathcal{J}_{K^\circ}(y) = \langle x, y \rangle \iff x \in K \text{ and } y \in K^\circ \text{ and } \langle x, y \rangle = 0$$

Then, the rightmost equivalence follows from Corollary A.5.

□

### A.3 Constrained optimization

We consider the minimization problem

$$p = \min_{x \in C} f(x) \tag{A.4}$$

with  $f$  closed and convex, and  $C$  a closed convex subset of  $X$ .

Note that the constrained minimization problem (A.4) is equivalent to the unconstrained one (A.5),

$$p = \min_{x \in X} (f + \mathcal{I}_C)(x) \quad (\text{A.5})$$

We first recall a sufficient condition for the existence of a solution to problem (A.5) (see e.g., Barbu and Precupanu 2012, p. 72, Theorem 2.11, Proposition 2.10 and Remark 2.13).

**Theorem A.5** (Existence of a solution). *Let  $f : X \rightarrow \bar{\mathbb{R}}$  be convex and lower semi-continuous, and  $C \subset X$  is convex and closed. If the minimization problem (A.4) is feasible, that is  $\text{dom } f \cap C \neq \emptyset$ , and  $f$  is coercive on  $C$ , that is,*

$$\lim_{x \in C, \|x\| \rightarrow +\infty} f(x) = +\infty,$$

*then  $f$  attains its global minimum on  $C$ .*

### A.3.1 Optimality conditions

We now derive sufficient conditions for checking that a point  $x_0 \in X$  is a solution to (A.4), i.e.,  $x_0 \in C$  and  $f(x_0) = p$ .

**Theorem A.6** (Fundamental theorem of convex optimization). *Suppose  $f$  a proper closed convex function on  $X$  and  $C$  a closed convex subset of  $X$ , such that  $f$  is real-valued and continuous on  $C$ . The optimality condition of (A.4) reads*

$$\partial f(x_0) \cap -\mathcal{N}_C(x_0) \neq \emptyset.$$

*If  $f$  is moreover Gâteaux-differentiable on  $C$ , the optimality condition becomes*

$$(\nabla f)(x_0) \in -\mathcal{N}_C(x_0).$$

*Proof.* First, if  $C$  is empty the problem is not feasible,  $\mathcal{N}_C$  is always empty and the equivalence is trivial. Let us now assume  $C$  not empty. From Property A.3, the optimality condition of (A.5) and therefore of (A.4) is  $\mathbf{0} \in \partial g(x_0)$ .

We can apply Theorem A.1 to get  $\partial g = \partial \mathcal{I}_C + \partial f$ . Expressing  $\partial \mathcal{I}_C$  with Property A.9, this means  $\mathbf{0} \in \partial f(x_0) + \mathcal{N}_C(x_0)$ . If  $f$  is Gâteaux-differentiable at  $x_0$ , Property A.4 states that its subdifferential contains only its gradient, so the optimality condition becomes  $\mathbf{0} \in \{(\nabla f)(x_0)\} + \mathcal{N}_C(x_0)$ .  $\square$

**Corollary A.6.** *If  $X$  is an Hilbert space, then for  $x, y, \alpha \in X \times X \times \mathbb{R}_+^*$ ,*

$$y \in \mathcal{N}_C(x) \iff \Pi_C(x + \alpha y) = x$$

*where  $\Pi_C$  denotes the orthogonal projection on the closed convex subset  $C$  of  $X$ ,*

$$\Pi_C(z) = \arg \min_{x \in C} \frac{1}{2} \langle z - x, z - x \rangle$$

*Proof.* For any  $z \in X$ , let  $f_z(x) := \frac{1}{2} \langle z - x, z - x \rangle$ , which a real-valued, closed, convex, and differentiable function on  $X$ .

Applying Theorem A.6 to the minimization of  $f_z$  on  $C$ , we get

$$x_0 = \Pi_C(z) \iff \nabla f_z(x_0) \in -\mathcal{N}_C(x_0) \iff x_0 - z \in -\mathcal{N}_C(x_0).$$

For  $z := x_0 + \alpha y$  with  $\alpha \in \mathbb{R}_+^*$ , this reduces to

$$x_0 = \Pi_C(x_0 + \alpha y) \iff \alpha y \in \mathcal{N}_C(x_0) \iff y \in \mathcal{N}_C(x_0).$$

$\square$

**Remark A.3.** Combining Theorem A.6 with Corollary A.6, the optimality condition of (A.4) becomes

$$\Pi_C(\mathbf{x}_0 - \alpha(\nabla f)(\mathbf{x}_0)) = \mathbf{x}_0, \quad \alpha \in \mathbb{R}_+^*,$$

yielding the rationale for the Projected Gradient Descent algorithm.

**Corollary A.7** (Karush–Kuhn–Tucker conditions). *The optimality conditions of the minimization problem*

$$\min_{\mathbf{x} \in K} f(\mathbf{x})$$

with  $K$  is a closed convex cone, and  $f$  a proper closed convex function on  $X$  and real-valued, continuous and Gâteaux-differentiable on  $K$ , are

$$K \ni \mathbf{x}_0 \perp (\nabla f)(\mathbf{x}_0) \in K^*$$

and are called the first-order Karush–Kuhn–Tucker conditions.

*Proof.* This is a direct application of Theorem A.6 combined with the equivalences from Property A.9 and Theorem A.4:

$$(\nabla f)(\mathbf{x}_0) \in -\mathcal{N}_K(\mathbf{x}_0) \iff (\nabla f)(\mathbf{x}_0) \in -\partial \mathcal{J}_K(\mathbf{x}_0) \iff K \ni \mathbf{x}_0 \perp (\nabla f)(\mathbf{x}_0) \in K^*$$

□

### A.3.2 Lagrange multipliers

A last application concerns the structure of the normal cone at the boundary of an implicitly-defined convex set.

**Property A.16** (Normal cone at the boundary). *Consider a set  $C$  defined as  $C := \{\mathbf{x} \in X, F(\mathbf{x}) \leq 0\}$  with  $F$  a proper closed convex function on  $X$ , continuous on  $C$ . For  $\mathbf{x}_0 \in \text{Bd } C$ , there always holds*

$$\mathcal{N}_C(\mathbf{x}_0) \supset \{\alpha \mathbf{g}, \mathbf{g} \in \partial F(\mathbf{x}_0), \alpha \in \mathbb{R}_+\}, \quad (\text{A.6})$$

and equality is achieved if  $\inf F < 0$ . In this latter case, we also have the equivalence  $F(\mathbf{x}) = 0 \iff \mathbf{x}_0 \in \text{Bd } C$ , so that for any  $\mathbf{x}_0$  such that  $F(\mathbf{x}) < 0$ ,  $\mathcal{N}_C(\mathbf{x}_0) = \{\mathbf{0}\}$ .

*Proof.*  $C$  is convex (from the convexity of  $F$ ) and closed (Property A.1). Moreover, as  $F$  is continuous on  $C$ ,  $F(\mathbf{x}) < 0 \implies \mathbf{x} \in \text{int } C$  and therefore  $\mathbf{x} \in \text{Bd } C \implies F(\mathbf{x}) = 0$ .

We can easily show that  $\{\alpha \mathbf{g}, \mathbf{g} \in \partial F(\mathbf{x}_0), \alpha \in \mathbb{R}_+\} \subset \mathcal{N}_C(\mathbf{x}_0)$ . Indeed, let  $\mathbf{x}_0 \in \text{Bd } C$  and  $\mathbf{g} \in \partial F(\mathbf{x}_0)$ . This means

$$\begin{aligned} F(\mathbf{x}) &\geq F(\mathbf{x}_0) + \langle \mathbf{x} - \mathbf{x}_0, \mathbf{g} \rangle & \forall \mathbf{x} \in X \\ 0 &\geq \langle \mathbf{x} - \mathbf{x}_0, \mathbf{g} \rangle & \forall \mathbf{x}, F(\mathbf{x}) \leq 0 \end{aligned}$$

and therefore  $\mathbf{y} \in \mathcal{N}_C(\mathbf{x}_0)$ .

Now assume  $\inf F < 0$ , and let us prove the converse inclusion. Let  $\mathbf{y} \in \mathcal{N}_C(\mathbf{x}_0)$ ; if  $\mathbf{y} = \mathbf{0}$ , we have  $\exists \mathbf{g} \in \partial F(\mathbf{x}_0), \mathbf{0} \mathbf{g} = \mathbf{y}$  as  $\partial F(\mathbf{x}_0)$  is not empty (Moreau 1966–1967, paragraph 10.c). Now, let us assume  $\mathbf{y} \neq \mathbf{0}$ , and consider the constrained optimization problem

$$\min_{\mathbf{x} \in C} -\langle \mathbf{x}, \mathbf{y} \rangle. \quad (\text{A.7})$$

From Theorem A.6,  $\mathbf{x}_0$  satisfies the optimality condition of (A.7). Now, suppose that the minimum is also reached at  $\mathbf{x} \in \text{int } C$ ; this means  $\mathbf{y} \in \mathcal{N}_C(\mathbf{x}) = \{\mathbf{0}\}$  from Property A.10, which is a contradiction. Therefore,  $\forall \mathbf{x} \in \text{int } C, -\langle \mathbf{x}, \mathbf{y} \rangle > -\langle \mathbf{x}_0, \mathbf{y} \rangle$ . As  $\mathbf{x}_0 \in \text{Bd } C$ , we get

$$\mathbf{x} \in \text{int } C \implies F(\mathbf{x}) < \underbrace{F(\mathbf{x}_0)}_0 \implies \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle < 0$$

or equivalently

$$\langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle \geq 0 \implies F(\mathbf{x}) \geq F(\mathbf{x}_0). \quad (\text{A.8})$$

Let us denote by  $\mathcal{H}_{\mathbf{y}}(\mathbf{x}_0)$  the half-space

$$\mathcal{H}_{\mathbf{y}}(\mathbf{x}_0) := \{\mathbf{x} \in X, \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle \geq 0\}.$$

Then Equation (A.8) means that  $\mathbf{x}_0$  satisfies the optimality conditions of the minimization problem

$$\min_{\mathbf{x} \in \mathcal{H}_{\mathbf{y}}(\mathbf{x}_0)} F(\mathbf{x}),$$

which are, using Property (A.3) then Theorem A.1 ( $F$  is continuous and both terms are convex and finite at  $\mathbf{x}_0$ ),

$$\begin{aligned} \mathbf{0} &\in \partial(F + \mathcal{I}_{\mathcal{H}_{\mathbf{y}}(\mathbf{x}_0)})(\mathbf{x}_0) = \partial F(\mathbf{x}_0) + \mathcal{N}_{\mathcal{H}_{\mathbf{y}}(\mathbf{x}_0)}(\mathbf{x}_0) \\ \mathbf{0} &\in \partial F(\mathbf{x}_0) + \{-\alpha \mathbf{y}, \alpha \in \mathbb{R}_+\}. \end{aligned}$$

As we supposed that  $\inf F < 0$ ,  $\mathbf{0} \notin \partial F(\mathbf{x}_0)$ , and therefore  $\exists \alpha \in \mathbb{R}_+$  and  $\mathbf{g} \in \partial F(\mathbf{x}_0)$  such that  $\mathbf{y} = \alpha \mathbf{g}$ .

It remains to show that  $\mathbf{x} \in \text{int } C \implies F(\mathbf{x}) < 0$ . Suppose  $F(\mathbf{x}) = 0$ . As  $\inf F < 0$ ,  $\mathbf{0} \notin \partial F(\mathbf{x}_0)$ , and we can find a non-zero element  $\mathbf{y} \in \partial F$ , which is non-empty. From the duality of  $X$  and  $Y$ , there exists  $\mathbf{x} \in X$ ,  $\langle \mathbf{x}, \mathbf{y} \rangle > 0$ . For every  $\epsilon > 0$ , let  $\mathbf{x}_\epsilon = \mathbf{x}_0 + \frac{\epsilon}{\|\mathbf{x}\|} \mathbf{x}$ ;  $\|\mathbf{x}_0 - \mathbf{x}_\epsilon\| \leq \epsilon$  and  $F(\mathbf{x}_\epsilon) > 0$ , i.e.,  $\mathbf{x}_0 \in \text{Bd } C$ .  $\square$

The qualification hypothesis  $\inf F < 0$  is not necessary, however the equality would then require supplemental constraints on the structure of  $\partial F$ . Intuitively, we require  $F$  to go to zero with a strictly positive slope in all directions; the problem is indeed similar to that of the existence of Lagrange Multipliers, which require a surjective differential for the application of the Implicit Function Theorem.

An example of a function for which inclusion (A.6) is strict while the subdifferential is nowhere equal to  $\{\mathbf{0}\}$  is

$$F : \mathbb{R} \rightarrow \mathbb{R}, \quad x \mapsto \begin{cases} x^2 & \text{if } x \leq 0 \\ x & \text{if } x > 0. \end{cases}$$

Indeed,  $\partial F(0) = [0, 1]$  while  $\partial \mathcal{I}_C = \mathbb{R}$ .

Conversely, an example for which  $\min F = 0$  yet the equality holds is the case of affine equality constraints:

**Property A.17.** Suppose  $A : X \rightarrow E$ ,  $\mathbf{x} \mapsto L\mathbf{x} + \mathbf{e}$  an affine map between two Hilbert spaces. Let  $F : \mathbf{x} \mapsto \|A(\mathbf{x})\|$ , and  $C := \{\mathbf{x} \in X, F(\mathbf{x}) = 0\}$ . Then for  $\mathbf{x}_0$  such that  $A(\mathbf{x}_0) = \mathbf{0}$ ,

$$\mathcal{N}_C(\mathbf{x}_0) = \text{Im } L^\top = \{\alpha \mathbf{g}, \mathbf{g} \in \partial F(\mathbf{x}_0), \alpha \in \mathbb{R}_+\}.$$

*Proof.* First, suppose  $\mathbf{y} \in \mathcal{N}_C(\mathbf{x}_0)$ .  $\forall \mathbf{x} \in C, \langle \mathbf{x} - \mathbf{x}_0, \mathbf{y} \rangle \leq 0$ . As  $A(\mathbf{x}_0) = \mathbf{0}$ ,  $F(\mathbf{x}) = 0 \iff A(\mathbf{x}) - A(\mathbf{x}_0) = \mathbf{0} \iff L(\mathbf{x} - \mathbf{x}_0) = \mathbf{0}$ . Therefore,  $\mathbf{y}$  must obey that  $\forall \mathbf{x} \in \text{Ker } L, \langle \mathbf{x}, \mathbf{y} \rangle \leq 0$ , which means  $\mathbf{y} \in (\text{Ker } L)^\perp$ . We thus have  $\mathcal{N}_C(\mathbf{x}_0) \subset \text{Im } L^\top$ .

Now, let us prove that

$$\text{Im } L^\top \subset \{\alpha \mathbf{g}, \mathbf{g} \in \partial F(\mathbf{x}_0), \alpha \in \mathbb{R}_+\},$$

which will achieve the demonstration, the converse inclusion being granted by Property A.16. Property A.12 states that we always have  $\partial F(\mathbf{x}_0) \supset L^\top \partial \|\cdot\|(A(\mathbf{x}_0))$ . Since  $A(\mathbf{x}_0) = \mathbf{0}$ , using Theorem (A.2),

$$\mathbf{g} \in \partial \|\cdot\|(A(\mathbf{x}_0)) \iff \mathbf{0} \in \partial \|\cdot\|^\star(\mathbf{g}).$$

Direct computations yield  $\|\cdot\|^\star = \mathcal{I}_{\mathcal{B}(0,1)}$ , the characteristic function of the unit closed ball centered at  $\mathbf{0}$ . This means

$$\mathbf{g} \in \partial \|\cdot\|(A(\mathbf{x}_0)) \iff \mathbf{0} \in \mathcal{N}_{\mathcal{B}(0,1)}(\mathbf{g}) \iff \|\mathbf{g}\| \leq 1.$$

It results that  $\partial F(\mathbf{x}_0) \supset L^\top \mathcal{B}(\mathbf{0}, 1)$ . Then, for all  $\mathbf{y} \in \text{Im } L^\top$ ,  $\exists \alpha \in \mathbb{R}_+$ ,  $\mathbf{y} = L^\top \alpha \mathbf{g}$  with  $\|\mathbf{g}\| \leq 1$ , and in particular,  $(L^\top \mathbf{g}) \in \partial F(\mathbf{x}_0)$ .  $\square$

**Property A.18** (KKT multiplier). *Let  $C := \{\mathbf{x} \in X, F(\mathbf{x}) \leq 0 \text{ and } A(\mathbf{x}) = 0\}$  with  $F$  convex, closed and continuous on  $C$ , and  $A$  an affine map,  $A : X \rightarrow E$ ,  $\mathbf{x} \mapsto L\mathbf{x} + \mathbf{e}$ . We consider the optimization problem (A.4),*

$$\min_{\mathbf{x} \in C} f(\mathbf{x}), \quad (\text{A.4})$$

for  $f$  closed, convex, continuous and real-valued on  $C$ .

If there exists  $\mathbf{x} \in X$  such that  $A(\mathbf{x}) = 0$  and  $F(\mathbf{x}) < 0$ , then the optimality conditions of (A.4) read

$$\begin{cases} \emptyset \neq \partial f(\mathbf{x}_0) \cap (-\lambda \partial F(\mathbf{x}_0) + L^\top \mathbf{r}) \\ \mathbf{0} = A(\mathbf{x}_0) \\ 0 \leq F(\mathbf{x}_0) \perp \lambda \geq 0. \end{cases} \quad (\text{A.9})$$

$\lambda \in \mathbb{R}$  is the Karush–Kuhn–Tucker multiplier associated to the constraint  $F(\mathbf{x}) \leq 0$ , and  $\mathbf{r} \in E^*$  is the vector of Lagrange multipliers associated to the constraint  $A(\mathbf{x}) = 0$ .

*Proof.* Property A.6 gives the optimality condition  $\partial f(\mathbf{x}_0) \cap \mathcal{N}_C \neq \emptyset$ . Let us express  $\mathcal{N}_C = \mathcal{N}_{C_A \cap C_F}$ ,  $C_A := \{A(\mathbf{x}) = 0\}$ ,  $C_F := \{F(\mathbf{x}) \leq 0\}$ .

Our hypothesis,  $\exists \mathbf{x}, A(\mathbf{x}) = 0$  and  $F(\mathbf{x}) < 0$ , implies that  $C_A \cap \text{int } C_F \neq \emptyset$ ; we can therefore use the Corollary A.1 on the subdifferential of a sum, and obtain with Property A.9 that  $\mathcal{N}_{C_A \cap C_F} = \mathcal{N}_{C_A} + \mathcal{N}_{C_F}$ .

$\mathcal{N}_{C_A}$  is straightforward to compute, and Property A.16 gives us the expression of the normal cone  $\mathcal{N}_{C_F}(\mathbf{x}_0)$ ,

$$\begin{aligned} \mathcal{N}_{C_A}(\mathbf{x}_0) &= \begin{cases} \emptyset & \text{if } A(\mathbf{x}_0) \neq \mathbf{0} \\ \text{Im } L^\top & \text{if } A(\mathbf{x}_0) = \mathbf{0} \end{cases} \\ \mathcal{N}_{C_F}(\mathbf{x}_0) &= \begin{cases} \emptyset & \text{if } F(\mathbf{x}_0) > 0 \\ \{\mathbf{0}\} & \text{if } F(\mathbf{x}_0) < 0 \\ \{\lambda \mathbf{g}, \lambda \in \mathbb{R}_+, \mathbf{g} \in \partial F(\mathbf{x}_0)\} & \text{if } F(\mathbf{x}_0) = 0. \end{cases} \end{aligned}$$

$\square$

**Remark A.4.** Property A.18 can easily be extended to the case of multiple constraints thanks to Property A.6, with  $F := \max_{i=1, \dots, n} F_i(\mathbf{x})$ . The optimality conditions are then, for  $f$  Gâteaux-differentiable, and in the absence of linear constraints,

$$\begin{cases} (\nabla f)(\mathbf{x}_0) = \sum_{i=1}^n \lambda_i \mathbf{g}_i, \mathbf{g}_i \in \partial F_i \\ \mathbb{R}_+^n \ni (F_i(\mathbf{x}_0)) \perp (\lambda_i) \in \mathbb{R}_+^n. \end{cases}$$

**Theorem A.7** (Lagrange duality). *With the notations of Property A.18, let  $p$  be the minimum of the constrained minimization problem (A.4),*

$$p = \min_{\mathbf{x} \in C} f(\mathbf{x}). \quad (\text{A.4})$$

Under the qualification hypothesis of A.18, if there exists  $\mathbf{x}_0 \in C$  such that  $p = f(\mathbf{x}_0)$ , then there exists  $(\lambda_0, \mathbf{r}_0) \in \mathbb{R}_+ \times E$  such that the tuple  $(\mathbf{x}_0, \lambda_0, \mathbf{r}_0)$  is a saddle points of the Lagrangian  $\mathcal{L}(\mathbf{x}, \lambda, \mathbf{r})$ ,

$$\mathcal{L}(\mathbf{x}, \lambda, \mathbf{r}) := f(\mathbf{x}) + \lambda F(\mathbf{x}) + \langle \mathbf{r}, A(\mathbf{x}) \rangle.$$

That is,  $(\lambda_0, \mathbf{r}_0)$  is a solution to the dual maximization problem,

$$d = \max_{\mathbf{r} \in E, \lambda \geq 0} J(\lambda, \mathbf{r}), \quad J(\lambda, \mathbf{r}) := \min_{\mathbf{x} \in X} (\mathcal{L}(\mathbf{x}, \lambda, \mathbf{r})). \quad (\text{A.10})$$

with  $J(\lambda_0, \mathbf{r}_0) = \mathcal{L}(\mathbf{x}_0, \lambda_0, \mathbf{r}_0)$ .

Moreover,  $p = d$ , and if  $(\lambda_0, \mathbf{r}_0) \in \mathbb{R}_+ \times E^*$  is a solution to the dual problem (A.10), any solution  $\mathbf{x}_0 \in X$  to the primal problem (A.4) must satisfy  $d = J(\lambda_0, \mathbf{r}_0) = \mathcal{L}(\mathbf{x}_0, \lambda_0, \mathbf{r}_0)$ .

*Proof.* First, let us show that  $J$  is concave (i.e.,  $-J$  is convex). For  $\alpha \in ]0, 1[$ ,

$$\begin{aligned} \alpha J(\lambda_1, \mathbf{r}_1) + (1 - \alpha)J(\lambda_2, \mathbf{r}_2) &= \min_{\mathbf{x} \in X} (\alpha \mathcal{L}(\mathbf{x}, \lambda_1, \mathbf{r}_1)) + \min_{\mathbf{x} \in X} ((1 - \alpha) \mathcal{L}(\mathbf{x}, \lambda_2, \mathbf{r}_2)) \\ &\geq \min_{\mathbf{x} \in X} (\alpha \mathcal{L}(\mathbf{x}, \lambda_1, \mathbf{r}_1) + (1 - \alpha) \mathcal{L}(\mathbf{x}, \lambda_2, \mathbf{r}_2)) \\ &\geq \min_{\mathbf{x} \in X} \mathcal{L}(\mathbf{x}, \alpha \lambda_1 + (1 - \alpha) \lambda_2, \alpha \mathbf{r}_1 + (1 - \alpha) \mathbf{r}_2) \\ &\geq J(\alpha \lambda_1 + (1 - \alpha) \lambda_2, \alpha \mathbf{r}_1 + (1 - \alpha) \mathbf{r}_2) \end{aligned}$$

This means that we can use Property A.6 to get an optimality condition for the dual problem (A.10),

$$\emptyset \neq -\partial(-J)(\lambda, \mathbf{r}) \neq (\mathcal{N}_{\mathbb{R}_+}(\lambda) \times \mathbf{0}_{E^*}). \quad (\text{A.11})$$

Suppose that  $(\mathbf{x}_0)$  is a solution to the primal problem (A.4). From Property (A.18), this means that there exists  $(\lambda_0, \mathbf{r}_0)$  such that

$$\begin{cases} \emptyset \neq \partial f(\mathbf{x}_0) \cap (-\lambda_0 \partial F(\mathbf{x}_0) + L^\top \mathbf{r}_0) \\ \mathbf{0} = A(\mathbf{x}_0) \\ 0 \leq F(\mathbf{x}_0) \perp \lambda_0 \geq 0. \end{cases} \quad (\text{A.9})$$

The first equation ensures that  $\mathbf{x}_0$  is a solution to the minimization problem  $J(\lambda_0, \mathbf{r}_0)$ , i.e.,  $J(\lambda_0, \mathbf{r}_0) = \mathcal{L}(\mathbf{x}_0, \lambda_0, \mathbf{r}_0)$ . It remains to show that  $(\lambda_0, \mathbf{r}_0)$  satisfies the optimality condition (A.11) of the dual problem.

As  $J$  is a minimum over  $\mathbf{x} \in X$ ,  $\forall \lambda, \mathbf{r} \in \mathbb{R} \times E$ ,

$$\begin{aligned} J(\lambda, \mathbf{r}) &\leq \mathcal{L}(\mathbf{x}_0, \lambda, \mathbf{r}) \\ &\leq F(\mathbf{x}_0) + \lambda F(\mathbf{x}_0) + \langle \mathbf{r}, A(\mathbf{x}_0) \rangle \\ &\leq J(\mathbf{x}_0) + (\lambda - \lambda_0)F(\mathbf{x}_0) + \langle \mathbf{r} - \mathbf{r}_0, A(\mathbf{x}_0) \rangle \end{aligned}$$

which by definition means that  $(F(\mathbf{x}_0), A(\mathbf{x}_0)) \in -\partial(-J)$ . We can easily check from the conditions (A.9) that  $(F(\mathbf{x}_0), A(\mathbf{x}_0)) \in \mathcal{N}_{\mathbb{R}_+}(\lambda_0) \times \mathbf{0}_{E^*}$ , meaning that the optimality conditions of the dual problem (A.10) are satisfied.

Notice that  $\lambda_0 F(\mathbf{x}_0) = 0$  and  $\langle \mathbf{r}_0, A(\mathbf{x}_0) \rangle = 0$ , i.e.,  $\mathcal{L}(\mathbf{x}_0, \lambda_0, \mathbf{r}_0) = f(\mathbf{x}_0)$ , and therefore  $p = d$ .

We have proved that if there exists a solution to the primal problem (A.4), then  $p = d$ . This is actually always true, as we will show below; but for now, assume that  $p = d$ , and let us show that, for  $(\lambda_0, \mathbf{r}_0)$  solution to the dual problem (A.10), any primal solution  $\mathbf{x}_0$  is such that  $\mathcal{L}(\lambda_0, \mathbf{r}_0, \mathbf{x}_0) = J(\lambda_0, \mathbf{r}_0)$ . Indeed, as  $\mathbf{x}_0 \in C_F$ ,  $\mathcal{L}(\mathbf{x}_0, \lambda_0, \mathbf{r}_0) = f(\mathbf{x}_0) + \lambda_0 F(\mathbf{x}_0)$ , with  $\lambda_0 F(\mathbf{x}_0) \leq 0$ . Therefore

$$d = J(\lambda_0, \mathbf{r}_0) \leq \mathcal{L}(\mathbf{x}_0, \lambda_0, \mathbf{r}_0) \leq f(\mathbf{x}_0) = p = d,$$

from which we deduce that  $\mathcal{L}(\mathbf{x}_0, \lambda_0, \mathbf{r}_0) = J(\lambda_0, \mathbf{r}_0)$ .

In the remainder of this proof, we show that  $p = d$  regardless of the existence of a primal solution. It always hold that

$$d = \max_{\lambda \geq 0, \mathbf{r} \in E} \min_{\mathbf{x} \in X} \mathcal{L}(\mathbf{x}, \lambda, \mathbf{r}) \leq \max_{\lambda \geq 0, \mathbf{r} \in E} \min_{\mathbf{x} \in C} \mathcal{L}(\mathbf{x}, \lambda, \mathbf{r}) \leq \max_{\lambda \geq 0, \mathbf{r} \in E} \min_{\mathbf{x} \in C} f(\mathbf{x}) \leq p.$$

We therefore have to show that  $p \leq d$ . Let us introduce the affine map  $\hat{A} : X \rightarrow X \times E$ ,  $\mathbf{x} \mapsto \begin{bmatrix} \mathbf{x} \\ A(\mathbf{x}) + \mathbf{e} \end{bmatrix}$ , and the characteristic function  $g : X \times E \rightarrow \mathbb{R}$ ,  $(\mathbf{x}, \mathbf{y}) \mapsto \mathcal{I}_{C_F \times \{\mathbf{0}_E\}}$ .



We can use Fenchel's duality theorem (A.3) to get that

$$\begin{aligned}
 p &= \inf_{\mathbf{x} \in X} (f(\mathbf{x}) + g \circ \hat{A}g) = \sup_{\mathbf{y} \in Y, \mathbf{r} \in E^*} J_f(\mathbf{r}, \mathbf{y}), \\
 J_f(\mathbf{r}, \mathbf{y}) &= -f^*(\mathbf{y} + L^T \mathbf{r}) - g^*(-\mathbf{y}, -\mathbf{r}) - \langle \mathbf{e}, \mathbf{r} \rangle \\
 &= \sup_{\mathbf{y} \in Y, \mathbf{r} \in E^*} - \sup_{\mathbf{x} \in X} (\langle \mathbf{x}, \mathbf{y} + L^T \mathbf{r} \rangle - f(\mathbf{x})) - g^*(-\mathbf{y}, -\mathbf{r}) - \langle \mathbf{e}, \mathbf{r} \rangle \\
 &= \sup_{\mathbf{y} \in Y, \mathbf{r} \in E^*} - \sup_{\mathbf{x} \in X} (\langle \mathbf{x}, \mathbf{y} \rangle + \langle A(\mathbf{x}), \mathbf{r} \rangle - f(\mathbf{x})) - g^*(-\mathbf{y}, -\mathbf{r}) \\
 &= \sup_{\mathbf{y} \in Y, \mathbf{r} \in E^*} \inf_{\mathbf{x} \in X} (\langle \mathbf{x}, \mathbf{y} \rangle + \langle A(\mathbf{x}), \mathbf{r} \rangle + f(\mathbf{x})) - g^*(\mathbf{y}, \mathbf{r}).
 \end{aligned}$$

As  $p$  is a supremum, we need only consider the set of  $\mathbf{y}, \mathbf{r}$  for which  $g^*(\mathbf{y}, \mathbf{r}) < +\infty$ . Let us express  $g^*$ :

$$\begin{aligned}
 g^*(\mathbf{y}, \mathbf{r}) &= \sup_{\mathbf{x} \in X, \mathbf{z} \in E} (\langle \mathbf{y}, \mathbf{x} \rangle - \mathcal{J}_{C_F}(\mathbf{x}) + \langle \mathbf{r}, \mathbf{z} \rangle - \mathcal{J}_{0_E}(\mathbf{z})) \\
 &= \sup_{\mathbf{x} \in X} (\langle \mathbf{y}, \mathbf{x} \rangle - \mathcal{J}_{C_F}(\mathbf{x})) + \sup_{\mathbf{r} \in E} (\langle \mathbf{r}, \mathbf{z} \rangle - \mathcal{J}_{0_E}(\mathbf{z})) \\
 &= - \inf_{\mathbf{x} \in C_F} (-\langle \mathbf{y}, \mathbf{x} \rangle).
 \end{aligned}$$

The optimality condition for  $\inf_{\mathbf{x} \in C_F} (-\langle \mathbf{y}, \mathbf{x} \rangle)$  is  $\exists \mathbf{x}_0 \in X, \mathbf{y} \in \mathcal{N}_{C_F}(\mathbf{x}_0)$ . Moreover, if  $g^*(\mathbf{y}, \mathbf{r})$  is finite but the optimum is non attained, then  $\forall \epsilon > 0$ , the penalized optimization problem

$$q_\epsilon(\mathbf{y}) = \min_{\mathbf{x} \in C_F} (-\langle \mathbf{y}, \mathbf{x} \rangle + \epsilon \|\mathbf{x}\|) \quad (\text{A.12})$$

attains its minimum, and  $\lim_{\epsilon \rightarrow 0_+} q_\epsilon(\mathbf{y}) = -g^*(\mathbf{y}, \mathbf{r})$ . Indeed, the function  $h_\epsilon : X \rightarrow \bar{\mathbb{R}}$ ,

$$h_\epsilon : \mathbf{x} \mapsto -\langle \mathbf{y}, \mathbf{x} \rangle + \epsilon \|\mathbf{x}\|$$

is convex, lower semi-continuous and coercive on a reflexive Banach space, and therefore, from Theorem A.5, attains its minimum on  $C_F$ . Moreover  $\forall \delta > 0, \exists \mathbf{x}_0 \in C_F, \langle \mathbf{x}_0, \mathbf{y} \rangle > g^*(\mathbf{y}, \mathbf{r}) - \delta$ , therefore  $\exists \epsilon_0$  such that  $\forall \epsilon < \epsilon_0, \langle \mathbf{x}_0, \mathbf{y} \rangle - \epsilon \|\mathbf{x}_0\| > g^*(\mathbf{y}, \mathbf{r}) - \delta$ . i.e.,  $\lim_{\epsilon \rightarrow 0_+} (q_\epsilon(\mathbf{y}) + g^*(\mathbf{y}, \mathbf{r})) < \delta$ . Since there also holds  $q_\epsilon \geq -g^*(\mathbf{y}, \mathbf{r})$ ,  $\lim_{\epsilon \rightarrow 0_+} (q_\epsilon(\mathbf{y}) + g^*(\mathbf{y}, \mathbf{r})) = 0$ .

Let us denote by  $V$  the subset of  $Y$  for which  $g^*(\mathbf{y}, \mathbf{r})$  is finite. For any  $\mathbf{y} \in V$ , for any  $\epsilon > 0$ , we note  $\mathbf{x}_\epsilon$  the element of  $C_F$  at which the optimal value of  $q_\epsilon$  is reached. We have

$$-g^*(\mathbf{y}, \mathbf{r}) \leq \lim_{\epsilon \rightarrow 0_+} \langle \mathbf{x}_\epsilon, -\mathbf{y} \rangle \leq \lim_{\epsilon \rightarrow 0_+} (\langle \mathbf{x}_\epsilon, -\mathbf{y} \rangle + \epsilon \|\mathbf{x}_\epsilon\|) = \lim_{\epsilon \rightarrow 0_+} q_\epsilon(\mathbf{y}) = -g^*(\mathbf{y}, \mathbf{r}),$$

therefore  $\lim_{\epsilon \rightarrow 0_+} \epsilon \|\mathbf{x}_\epsilon\| = 0$ . As for all  $\mathbf{y}_\epsilon \in \mathcal{B}(\mathbf{y}, \epsilon)$ ,

$$\lim_{\epsilon \rightarrow 0_+} (\langle \mathbf{x}_\epsilon, -\mathbf{y} \rangle + \epsilon \|\mathbf{x}_\epsilon\|) \geq \lim_{\epsilon \rightarrow 0_+} \langle \mathbf{x}_\epsilon, -\mathbf{y}_\epsilon \rangle \geq \lim_{\epsilon \rightarrow 0_+} (\langle \mathbf{x}_\epsilon, -\mathbf{y} \rangle - \epsilon \|\mathbf{x}_\epsilon\|),$$

we deduce that

$$\lim_{\epsilon \rightarrow 0_+} \langle \mathbf{x}_\epsilon, -\mathbf{y}_\epsilon \rangle = \lim_{\epsilon \rightarrow 0_+} (\langle \mathbf{x}_\epsilon, -\mathbf{y} \rangle) = -g^*(\mathbf{y}, \mathbf{r}).$$

Now, the optimality condition of (A.12) is  $\exists \mathbf{x}_\epsilon \in X, \mathbf{y} \in \mathcal{B}(\mathbf{0}, \epsilon) + \mathcal{N}_{C_F}(\mathbf{x}_\epsilon)$ . Going back to our original problem, we can now write that  $\forall \mathbf{y} \in V, \forall \mathbf{r} \in E^*$ ,

$$\begin{aligned}
 J_f(\mathbf{y}, \mathbf{r}) &= \lim_{\epsilon \rightarrow 0_+} \inf_{\mathbf{x} \in X} (\langle \mathbf{x}, \mathbf{y} \rangle + \langle A(\mathbf{x}), \mathbf{r} \rangle + f(\mathbf{x})) - \langle \mathbf{x}_\epsilon, \mathbf{y}_\epsilon \rangle & \mathbf{y}_\epsilon \in \mathcal{N}_{C_F}(\mathbf{x}_\epsilon) \\
 &= \lim_{\epsilon \rightarrow 0_+} \inf_{\mathbf{x} \in X} (\langle \mathbf{x} - \mathbf{x}_\epsilon, \mathbf{y}_\epsilon \rangle + \langle \mathbf{x}, \mathbf{y} - \mathbf{y}_\epsilon \rangle + \langle A(\mathbf{x}), \mathbf{r} \rangle + f(\mathbf{x})) & \mathbf{y}_\epsilon \in \mathcal{N}_{C_F}(\mathbf{x}_\epsilon) \\
 &\leq \lim_{\epsilon \rightarrow 0_+} \inf_{\mathbf{x} \in X} (\langle \mathbf{x} - \mathbf{x}_\epsilon, \mathbf{y}_\epsilon \rangle + \epsilon \|\mathbf{x}\| + \langle A(\mathbf{x}), \mathbf{r} \rangle + f(\mathbf{x})) & \mathbf{y}_\epsilon \in \mathcal{N}_{C_F}(\mathbf{x}_\epsilon)
 \end{aligned}$$

From Property A.16, we know that

$$\mathcal{N}_{C_F}(\mathbf{x}) = \{\mathbf{y} \in Y, \mathbf{y} = \alpha \mathbf{g}, \mathbf{g} \in \partial F(\mathbf{x}) \text{ and } 0 \leq \alpha \perp F(\mathbf{x}) \geq 0\}.$$

Therefore, for  $\mathbf{y}_\epsilon \in \mathcal{N}_{C_F}(\mathbf{x}_\epsilon)$ , and  $\forall \mathbf{x} \in X$ ,

$$\langle \mathbf{x} - \mathbf{x}_\epsilon, \mathbf{y}_\epsilon \rangle = \langle \mathbf{x} - \mathbf{x}_\epsilon, \alpha \mathbf{g} \rangle + \alpha F(\mathbf{x}_\epsilon) \leq \alpha F(\mathbf{x})$$

as  $\mathbf{g} \in \partial F(\mathbf{x}_\epsilon)$ , and then

$$\begin{aligned} J_f(\mathbf{y}, \mathbf{r}) &\leq \lim_{\epsilon \rightarrow 0_+} \sup_{\alpha \in R_+} \inf_{\mathbf{x} \in X} (\alpha F(\mathbf{x}) + \epsilon \|\mathbf{x}\| + \langle A(\mathbf{x}), \mathbf{r} \rangle + f(\mathbf{x})) \\ &\leq \sup_{\alpha \in R_+} \inf_{\mathbf{x} \in X} (\alpha F(\mathbf{x}) + \langle A(\mathbf{x}), \mathbf{r} \rangle + f(\mathbf{x})) \\ &\leq \sup_{\alpha \in R_+} J(\alpha, \mathbf{r}). \end{aligned}$$

We conclude that

$$p = \sup_{\mathbf{y} \in Y, \mathbf{r} \in E^*} J_f(\mathbf{y}, \mathbf{r}) = \sup_{\mathbf{y} \in V, \mathbf{r} \in E^*} J_f(\mathbf{y}, \mathbf{r}) \leq \sup_{\mathbf{r} \in E^*} \sup_{\alpha \in R_+} J(\alpha, \mathbf{r}) \leq d.$$

□



## B Discrete Coulomb Friction Problem solvers

### B.1 Newton SOC Fischer–Burmeister function

#### B.1.1 SOC Fischer–Burmeister derivatives

Let us recall that for any pair of vectors  $(\hat{r}, \hat{u}) \in \mathbb{R}^d \times \mathbb{R}^d$ , the vector  $\mathbf{z}^2 := \hat{r} \circ \hat{r} + \hat{u} \circ \hat{u}$  is in the SOC  $\mathcal{K}_1$ , and its square root can be computed as  $\mathbf{z} = \sqrt{\lambda_1} \mathbf{v}_1 + \sqrt{\lambda_2} \mathbf{v}_2$ ,

$$\begin{aligned} \mathbf{w} &:= \hat{r}_N \hat{r}_T + \hat{u}_N \hat{u}_T \\ \lambda_i &= \|\hat{r}\|^2 + \|\hat{u}\|^2 + (-1)^i \|2\mathbf{w}\| \\ \mathbf{v}_i &= \frac{1}{2} \begin{cases} \left(1; (-1)^i \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) & \text{if } \mathbf{w} \neq \mathbf{0} \\ \left(1; (-1)^i \mathbf{e}\right) & \text{if } \mathbf{w} = \mathbf{0}. \end{cases} \end{aligned}$$

The SOC Fischer–Burmeister function can thus be evaluated as  $f_{\text{FB}}^{\mathcal{K}}(\hat{r}, \hat{u}) = \sqrt{\lambda_1} \mathbf{v}_1 + \sqrt{\lambda_2} \mathbf{v}_2 - \hat{r} - \hat{u}$ . To compute the gradient of  $f_{\text{FB}}^{\mathcal{K}}$ , let us recall two properties of the SOC algebra (Chen and Tseng 2005, Property 1):

- If  $\mathbf{x} \in \text{int } \mathcal{K}_1$ ,  $\mathbf{x}$  is invertible<sup>1</sup> and  $\mathbf{x}^{-1} \in \text{int } \mathcal{K}_1$ ;
- $\mathbf{x} \in \mathbb{R} \times \mathbb{R}^{d-1}$  is invertible  $\iff \det \mathbf{x} = \mathbf{x}_N^2 - \|\mathbf{x}_T\|^2 \neq 0$ .

Since we always have  $\mathbf{z} \in \mathcal{K}_1$  and  $\det \mathbf{z} = \sqrt{\lambda_1 \lambda_2}$ , we get  $\mathbf{z}$  invertible  $\iff \mathbf{z} \in \text{int } \mathcal{K}_1 \iff \lambda_1 \lambda_2 \neq 0$ . We can then distinguish three cases. We have (Chen and Tseng 2005, Proposition 1):

- If  $\lambda_2 = 0$ , then  $\hat{r} = \hat{u} = \mathbf{0}$ , and thus  $f_{\text{FB}}^{\mathcal{K}}(\hat{r}, \hat{u}) = \mathbf{0}$ ; we have already found a solution to our problem.
- Otherwise, if  $\lambda_1 = 0$ , then  $\mathbf{z} \in \text{Bd } \mathcal{K}_1 \setminus \mathbf{0}$ , and we can arbitrarily choose one element of the generalized Jacobian of  $f_{\text{FB}}^{\mathcal{K}}$ . We take<sup>2</sup>:

$$\frac{\partial f_{\text{FB}}^{\mathcal{K}}}{\partial \hat{r}}(\hat{r}, \hat{u}) = \left( \frac{\hat{r}_N}{\sqrt{\hat{r}_N^2 + \hat{u}_N^2}} - 1 \right) \mathbb{I} \text{ and } \frac{\partial f_{\text{FB}}^{\mathcal{K}}}{\partial \hat{u}}(\hat{r}, \hat{u}) = \left( \frac{\hat{u}_N}{\sqrt{\hat{r}_N^2 + \hat{u}_N^2}} - 1 \right) \mathbb{I}.$$

- Otherwise, we have  $\mathbf{z} \in \text{int } \mathcal{K}_1$ , and both  $\mathbf{z}^{-1}$  and the Jacobian of  $f_{\text{FB}}^{\mathcal{K}}$  are uniquely defined. For  $\mathbf{x} \in \mathbb{R}^d$ , let us denote  $L_{\mathbf{x}}$  the  $d \times d$  matrix that satisfies  $\forall \mathbf{y} \in \mathbb{R}^d \ L_{\mathbf{x}} \mathbf{y} = \mathbf{x} \circ \mathbf{y}$ .

$$\begin{aligned} L_{\mathbf{z}^{-1}} &= \frac{1}{\sqrt{\lambda_1 \lambda_2}} \begin{pmatrix} \mathbf{z}_N & -\mathbf{z}_T^\top \\ -\mathbf{z}_T & \frac{\sqrt{\lambda_1 \lambda_2}}{\mathbf{z}_N} \mathbb{I}_{d-1} + \frac{\mathbf{z}_T \mathbf{z}_T^\top}{\mathbf{z}_N} \end{pmatrix}. \\ \text{Then } \frac{\partial f_{\text{FB}}^{\mathcal{K}}}{\partial \hat{r}}(\hat{r}, \hat{u}) &= \mathbb{I} - L_{\mathbf{z}^{-1}} L_{\hat{r}} \text{ and } \frac{\partial f_{\text{FB}}^{\mathcal{K}}}{\partial \hat{u}}(\hat{r}, \hat{u}) = \mathbb{I} - L_{\mathbf{z}^{-1}} L_{\hat{u}}. \end{aligned}$$

Finally, using the chain rule,

$$\frac{df_{\text{FB}}^{\mathcal{K}}}{d\mathbf{r}}(\mathbf{r}) = \frac{\partial f_{\text{FB}}^{\mathcal{K}}}{\partial \hat{r}}(\hat{r}, \hat{u}) \frac{d\hat{r}}{d\mathbf{r}} + \frac{\partial f_{\text{FB}}^{\mathcal{K}}}{\partial \hat{u}}(\hat{r}, \hat{u}) \frac{d\hat{u}}{d\mathbf{r}}.$$

Note that  $\frac{d\hat{u}}{d\mathbf{r}}$  is not strictly defined when  $\mathbf{u}_T = \mathbf{0}$ , because  $\|\mathbf{u}_T\|$  is not differentiable at  $\mathbf{0}$ . We can however take any element of its generalized Jacobian, such as  $I_d$ .

<sup>1</sup>With respect to the Jordan product, i.e.,  $\exists \mathbf{y}$  s.t.  $\mathbf{x} \circ \mathbf{y} = \mathbf{y} \circ \mathbf{x} = [1, \mathbf{0}]$ .

<sup>2</sup> $\hat{r}_N^2 + \hat{u}_N^2 > 0$ . Otherwise we would have  $\hat{r}_N = \hat{u}_N = 0$ , which would mean that  $\mathbf{z}_T = \mathbf{0}$  and either  $\mathbf{z} = \mathbf{0}$  or  $\mathbf{z} \in \text{int } \mathcal{K}_1$ .

### B.1.2 Optimistic Newton algorithm

The algorithm below is used as part of the hybrid strategy to solve one-contact problems in the Gauss–Seidel algorithm from Chapter 4. As the initial guess from the enumerative solver is usually quite good, this algorithm needs not be extremely robust; for this reason, we only use a very slight damping. Note that this algorithm was still observed to perform quite well without the hybrid strategy, for instance for solving 6-dimensional problems from continuum granular simulations (Chapters 7 and 8).

---

**Algorithm B.1:** Optimistic Non-Smooth Newton algorithm

---

```

input : Initial guess  $\mathbf{r}^0$ 
input : Matrix  $W$ , vector  $\mathbf{b}$ , friction coefficient  $\mu$ 
input : Curvature criterion  $\sigma > 0$  (typically  $10^{-2}$ )
Result:  $\mathbf{r}_{\text{best}}$  approximate solution to  $f_{\text{FB}}^{\mathcal{K}}(\hat{\mathbf{r}}(\mathbf{r}), \hat{\mathbf{u}}(\mathbf{r})) = 0$ 
 $\Phi_{\text{best}} \leftarrow +\infty$ ;
for  $k = 0$  to  $\text{maxIters}$  do
    // Compute  $f_{\text{FB}}^{\mathcal{K}}$  and its Jacobian at  $\mathbf{r}^k$  (see Appendix B.1.1)
     $(f_{\text{FB}}, J_{\text{FB}}) \leftarrow \text{computeFBandJacFB}(W, \mathbf{b}, \mu, \mathbf{r}^k)$ ;
     $\phi^k \leftarrow \frac{1}{2} \|f_{\text{FB}}\|^2$ ; // Compute value of error function
    if  $\phi^k < \phi_{\text{best}}$  then // Check quality of current estimate
         $\mathbf{r}_{\text{best}} \leftarrow \mathbf{r}^k$ ;
        if  $\phi^k < \text{tol}$  then
            break; // Error low enough, exit algorithm
        end
    end
     $\nabla \Phi \leftarrow J_{\text{FB}}^T f_{\text{FB}}$ ; // Compute gradient of objective function
     $\Delta \mathbf{r} \leftarrow -J_{\text{FB}}^{-1} f_{\text{FB}}$ ; // Compute new step  $\Delta \mathbf{r}$ 
    if  $\langle \Delta \mathbf{r}, \nabla \Phi \rangle > -\sigma \|\Delta \mathbf{r}\| \|\nabla \Phi\|$  then // Check for bad descent direction
         $\Delta \mathbf{r} \leftarrow \frac{1}{2} \Delta \mathbf{r}$ ;
    end
     $\mathbf{r}^{k+1} \leftarrow \mathbf{r}^k + \Delta \mathbf{r}$ ;
end

```

---

### B.2 Suggested variants of the Projected Gradient Descent algorithm

The variant of the Projected Gradient Descent algorithm given below is implemented in the bogus library, and performed consistently well on our frictional contact and continuum problems (w.r.t. other algorithms of the same kind). However, theoretical convergence properties have not been studied.

Algorithm B.2 is a line-search free but Nesterov-accelerated variant of the Spectral Projected Gradient algorithm from Tasora (2013), and aims to derive an optimal descent step size from the previous iterates.

---

**Procedure** NesterovInertia( $\theta_{\text{prev}}$ )

---

```

 $\Delta \leftarrow \theta_{\text{prev}} \sqrt{\theta_{\text{prev}}^2 + 4}$ ;
 $\theta \leftarrow \frac{1}{2} (\Delta - \theta_{\text{prev}}^2)$ ;
 $\beta \leftarrow \frac{\theta_{\text{prev}} (1 - \theta_{\text{prev}})}{\theta_{\text{prev}}^2 + \theta}$ ;
return  $(\theta, \beta)$ ;

```

---

**Algorithm B.2:** Accelerated Spectral Projected Gradient Descent

---

```

input : Initial guess  $\mathbf{r}^0$ 
input : Matrix  $W$ , vector  $\mathbf{b}$ , orthogonal projectors  $\Pi_C$ 
input : Min and max step sizes  $\xi_{\min}, \xi_{\max}$  (typically  $10^{-6}$  and  $10^6$ )
Result:  $\mathbf{r}_{\text{best}}$  approximate solution to  $\min_{\mathbf{r} \in C} \frac{1}{2} \mathbf{r}^T W \mathbf{r} + \mathbf{r}^T \mathbf{b}$ 
 $\mathbf{r}^0 \leftarrow \Pi_C(\mathbf{r}^0)$ ;
 $\mathbf{u}^0 \leftarrow W \mathbf{r}^0 + \mathbf{b}$ ; // Gradient of objective function
 $\hat{\mathbf{r}}^0 \leftarrow \mathbf{r}^0$ ;
 $\theta^k \leftarrow 1$ ; // Inertia
 $\xi^0 \leftarrow \text{guessStepSize}(W)$ 
 $\Phi_{\text{best}} \leftarrow +\infty$ ;
for  $k = 0$  to  $\text{maxIters}$  do
     $\phi^k \leftarrow \text{evaluateError}(\mathbf{r}^k, \mathbf{u}^k)$ ; // For instance using  $f_{\text{FB}}^{\mathcal{K}}$ 
    if  $\phi^k < \phi_{\text{best}}$  then
         $\mathbf{r}_{\text{best}} \leftarrow \mathbf{r}^k$ ;
        if  $\phi^k < \text{tol}$  then break; // Error low enough, exit algorithm
    end
     $\mathbf{d}^k \leftarrow -\mathbf{u}^k$ ;
     $\hat{\mathbf{r}}^{k+1} \leftarrow \Pi_C(\mathbf{r}^k + \xi^k \mathbf{d}^k)$ ;
    if  $\langle \hat{\mathbf{r}}^{k+1} - \hat{\mathbf{r}}^k, \mathbf{d}^k \rangle > 0$  then // Check for descent direction
         $(\theta^{k+1}, \beta) \leftarrow \text{NesterovInertia}(\theta^k)$ ;
         $\mathbf{r}^{k+1} \leftarrow \hat{\mathbf{r}}^{k+1} + \beta(\hat{\mathbf{r}}^{k+1} - \hat{\mathbf{r}}^k)$ ;
    else
         $\mathbf{r}^{k+1} \leftarrow \hat{\mathbf{r}}^{k+1}$ ;
         $\theta^{k+1} \leftarrow 1$ ; // Reset inertia
    end
     $\mathbf{u}^{k+1} \leftarrow W \mathbf{r}^{k+1} + \mathbf{b}$ ; // Gradient of objective function
     $\xi^{k+1} \leftarrow \frac{\langle \mathbf{r}^{k+1} - \mathbf{r}^k, \mathbf{u}^{k+1} - \mathbf{u}^k \rangle}{\|\mathbf{u}^{k+1} - \mathbf{u}^k\|^2}$ ; // Eq (5) from Barzilai and Borwein (1988)
     $\xi^{k+1} \leftarrow \min(\xi_{\max}, \max(\xi_{\min}, \xi^{k+1}))$ ; // Clamp next step size
end

```

---

Algorithm B.2 can also be used to solve directly a DCFP by replacing line B.2.13 with While this works relatively well in practice, the algorithm can no longer be interpreted as a

---


$$\mathbf{d}^k \leftarrow -\mathbf{u}^k - \mathbf{s}(\mathbf{u}^k);$$


---

minimization algorithm.

### B.3 Convergence of the out-of-order Gauss–Seidel algorithm

We study the convergence of the proximal, out-of-order Gauss–Seidel algorithm applied to the constrained minimization of a convex function outlined in Algorithm (B.3), where  $\beta_i \geq 0 \forall 1 \leq i \leq n$  and the sequence  $(j^k)_{k \in \mathbb{N}}$  denotes the index of the block that will be treated at the  $k^{\text{th}}$  iteration of the algorithm.

Algorithm 4.1 applied to the minimization of a SOCQP belongs to this class.

We make a few more hypothesis (which, again, are verified for Algorithm 4.1)

1.  $f$  is bounded from below;
2.  $f$  and  $\nabla f$  are continuous;

---

**Algorithm B.3:** Out-of-order Proximal-Point algorithm
 

---

**Input:**  $\mathbf{x}^0 \in C \subset X$ 
**for**  $k \in \mathbb{N}$  **do**

$$\mathbf{x}_{j^k}^{k+1} = \arg \min_{\mathbf{y} \in C_{j^k}} f_{j^k}(\mathbf{x}_1^k, \dots, \mathbf{x}_{j^k-1}^k, \mathbf{y}, \mathbf{x}_{j^k+1}^k, \dots, \mathbf{x}_n^k) + \frac{1}{2} \beta_{j^k} \|\mathbf{y} - \mathbf{x}_{j^k}^k\|^2;$$

$$\mathbf{x}^{k+1} = (\mathbf{x}_1^k, \dots, \mathbf{x}_{j^k}^{k+1}, \dots, \mathbf{x}_n^k)$$

**end**


---

3.  $C$  is closed, convex and non-empty;
4. the maximum number of iterations between a local minimization for a given block and the next one for this same block is finite. Formally,

$$\exists M \in \mathbb{N} \text{ s.t. } \forall 1 \leq i \leq n, \forall k \in \mathbb{N}, \exists k^* \in (k, k+M) \text{ s.t. } j^{k^*} = i. \quad (\text{B.1})$$

We will show that Algorithm B.3 converges to a solution of (B.2),

$$\arg \min_{\mathbf{x} \in C} f(\mathbf{x}), \quad (\text{B.2})$$

as long as either  $f$  is strictly convex or  $\beta_i > 0 \forall i$ .

The methodology below is adapted to our out-of-order setting from (Grippo and Sciandrone 2000).

**Lemma B.1.** *The sequence  $(f(\mathbf{x}^k))_k$  generated by Algorithm B.3 converges to a limit  $\bar{f}$ .*

*Proof.* We first introduce the function  $g^k : X \rightarrow \mathbb{R}, \mathbf{z} \mapsto f(\mathbf{z}) + \frac{1}{2} \beta_{j^k} \|\mathbf{z} - \mathbf{x}^k\|^2$ . By construction of Algorithm B.3,

$$\forall k \in \mathbb{N}, f(\mathbf{x}^{k+1}) \leq g^k(\mathbf{x}^{k+1}) \leq g^k(\mathbf{x}^k) = f(\mathbf{x}^k).$$

The sequence  $(f(\mathbf{x}^k))_k$  is thus decreasing and bounded from below, hence converges to a limit  $\bar{f}$ .  $\square$

**Property B.1.** *If for any  $1 \leq i \leq n$ , either  $\beta_i > 0$  or  $f$  is strictly convex and coercive w.r.t. the  $i^{\text{th}}$  block, then the sequence  $(\mathbf{x}^k)$  generated by Algorithm B.3 converges.*

In order to adapt the proof of Grippo and Sciandrone (2000), we first introduce a few supplemental notations.

**Definition B.1.** *Consider an infinite subset  $K \subset \mathbb{N}$ . For  $m \in \mathbb{N}$ , we define  $I_K^m$  as the set of all block indices that will be solved for infinitely many times at iteration  $k+m$  for  $k \in K$ . Formally*

$$I_K^m := \{i \in [1, n] \text{ s.t. } \forall k \in K, \exists k^* \in K, k^* \geq k \text{ and } j^{k^*+m} = i\}. \quad (\text{B.3})$$

We also define for each block the subset  $J_K^m(i) \subset K$  that contains all indices  $k \in K$  such that the block  $i$  will be solved for at iteration  $k+m$ ,

$$J_K^m(i) := \{k \in K \text{ s.t. } j^{k+m} = i\}.$$

Note that the set  $J_K^m(i)$  will be infinite if and only if  $i \in I_K^m$ .

**Lemma B.2.** *Let  $\overline{J}_K^m$  denote the union of the sets  $J_K^m(i)$  for  $i \in I_K^m$ ,*

$$\overline{J}_K^m := \bigcup_{i \in I_K^m} J_K^m(i) \subset K.$$

*For  $m \in \mathbb{N}$ , the subset  $S_K^m := K \setminus \overline{J}_K^m$  is finite.*

*Proof.* Suppose this was not the case, then the integer sequence  $\{j^{k+m}\}_{k \in S_K^m} \subset [1, n]$  has a limit point  $i^*$  which is also in  $I_K^m$  as it satisfies Equation B.3. By definition of  $i^*$  there exists  $k^* \in S_K^m$  such that  $j^{k^*+m} = i^*$ , which means  $k^* \in J_K^m(i^*) \cap S_K^m = \emptyset$ , which is a contradiction.  $\square$

We should now be able to prove Proposition B.1.

*Proof.* Let us denote by  $B \subset [1, n]$  the set for which  $\beta_i = 0$ .

**Existence of a limit point** First, let us show that from  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  we can extract a converging subsequence. From Lemma B.1, we have  $\lim_{k \in \mathbb{N}} f(\mathbf{x}^k) = \lim_{k \in \mathbb{N}} g^k(\mathbf{x}^{k+1}) = \bar{f}$ , therefore  $\forall 1 \leq i \leq n$ ,  $\lim_{k \in \mathbb{N}} \beta_i \|\mathbf{x}^k - \mathbf{x}^{k+1}\| = 0$ . Each block  $\mathbf{x}_i^k$  for which  $i \notin B$  thus admits a limit; let us define

$$\hat{\mathbf{x}}_i := \begin{cases} \lim \mathbf{x}_i^k & \text{if } i \notin B \\ \mathbf{0} & \text{if } i \in B. \end{cases}$$

The function  $\hat{g} : \mathbf{z} \mapsto f(\mathbf{z}) + \frac{\beta_i}{2} \|\mathbf{z} - \hat{\mathbf{x}}\|^2$  is coercive on  $X$  and  $\hat{g}(\mathbf{x}^k) \rightarrow \bar{f}$ , therefore  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  is bounded. Hence there exists a subsequence  $(\mathbf{x}^k)_{k \in K}$  that converges to  $\bar{\mathbf{x}} \in C$ , and by continuity of  $f$ ,  $f(\bar{\mathbf{x}}) = \lim_{k \in \mathbb{N}} f(\mathbf{x}^k) = \bar{f}$ . Necessarily,  $\bar{\mathbf{x}}_i = \hat{\mathbf{x}}_i \forall i \notin B$ .

**Convergence of translated subsequences** For  $m \in \mathbb{N}$ , we define the predicate  $\Pi(m)$ ,

$$\Pi(m) := \lim_{k \in K} \mathbf{x}^{k+m} = \bar{\mathbf{x}}$$

$\Pi(0)$  hold by definition of  $K$ ; now we suppose that  $\Pi(m)$  is true, and we will show that this implies  $\Pi(m+1)$ .

For any  $i \in I_K^{m+1} \cap B$ , by definition the set  $J_K^{m+1}(i) \subset K$  is infinite and thus  $(\mathbf{x}^{k+m})_{k \in J_K^{m+1}(i)} \rightarrow \bar{\mathbf{x}}$ . We can use (Grippe and Sciandrone 2000, Proposition 4) by identifying their  $\{y^k\}$  with  $(\mathbf{x}^{k+m})_{k \in J_K^{m+1}(i)}$  and their  $\{v^k\}$  with  $(\mathbf{x}^{k+m+1})_{k \in J_K^{m+1}(i)}$ . We get

$$\lim_{k \in J_K^{m+1}(i)} \|\mathbf{x}_i^{k+m} - \mathbf{x}_i^{k+m+1}\| = 0$$

and therefore from  $\Pi(m)$ ,

$$\lim_{k \in J_K^{m+1}(i)} \mathbf{x}_i^{k+m+1} = \bar{\mathbf{x}}_i. \quad (\text{B.4})$$

Moreover, Equation (B.4) also holds for any  $i \in I_K^{m+1} \setminus B$  since  $(\mathbf{x}_i^k)_{k \in \mathbb{N}} \rightarrow \hat{\mathbf{x}}_i = \bar{\mathbf{x}}_i$ . Now by definition of  $J_K^{m+1}(i)$  and Algorithm B.3,  $\forall i \in I_K^{m+1}$ ,  $\forall \ell \neq i$ ,  $\forall k \in J_K^{m+1}(i)$ ,  $\mathbf{x}_\ell^{k+m} = \mathbf{x}_\ell^{k+m+1}$ . Thus we get also from  $\Pi(m)$  that  $\forall i \in I_K^{m+1}$ ,

$$\lim_{k \in J_K^{m+1}(i)} \mathbf{x}_\ell^{k+m+1} = \lim_{k \in J_K^{m+1}(i)} \mathbf{x}_\ell^{k+m} = \bar{\mathbf{x}}_\ell \quad \forall \ell \neq i. \quad (\text{B.5})$$

Putting together (B.4) and (B.5), we get

$$\lim_{k \in J_K^{m+1}(i)} \mathbf{x}^{k+m+1} = \bar{\mathbf{x}}.$$

This being true for any  $i \in I_K^{m+1}$ , we deduce that

$$\lim_{k \in J_K^{m+1}} \mathbf{x}^{k+m+1} = \bar{\mathbf{x}}$$

and with Lemma B.2,

$$\lim_{k \in K} \mathbf{x}^{k+m+1} = \bar{\mathbf{x}}$$

which proves  $\Pi(m+1)$ . By induction  $\Pi(m)$  holds for all  $m \in \mathbb{N}$ , which means that  $(\mathbf{x}^k)_{k \in \mathbb{N}} \rightarrow \bar{\mathbf{x}}$ .  $\square$



**Property B.2.** Under the hypothesis of Property B.1, the limit  $\bar{\mathbf{x}}$  of the sequence of iterates  $(\mathbf{x}^k)$  from Algorithm B.3 is a critical point of the minimization problem (B.2).

Again, let us first introduce an useful definition.

**Definition B.2** ( Successor function ). Given Assumption B.1, we can define for each block  $1 \leq i \leq n$  the function

$$\text{succ}_i : \begin{cases} \mathbb{N} \longrightarrow \mathbb{N} \\ k \longmapsto \min \{k^* \in (k, k + M), j^{k^*} = i\} \end{cases}$$

which indicates the next iteration at which a given block will be minimized over.

*Proof.* For any  $1 \leq i \leq n$ , the sequence  $(\text{succ}_i(k))_{k \in \mathbb{N}}$  is infinite and thus

$$\lim_{k \in \mathbb{N}} \mathbf{x}^{\text{succ}_i(k)} = \bar{\mathbf{x}}. \quad (\text{B.6})$$

Moreover by construction of Algorithm B.3,  $(\nabla_i g^k)(\mathbf{x}^{\text{succ}_i(k)}) \in -\mathcal{N}_{C_i}(\mathbf{x}_i^{\text{succ}_i(k)})$ , i.e.,

$$\langle (\nabla_i g^k)(\mathbf{x}^{\text{succ}_i(k)}), \mathbf{z} - \mathbf{x}_i^{\text{succ}_i(k)} \rangle \geq 0 \quad \forall \mathbf{z} \in C_i.$$

This means

$$\langle (\nabla_i f^k)(\mathbf{x}^{\text{succ}_i(k)}) + \beta_i(\mathbf{x}_i^{\text{succ}_i(k)} - \mathbf{x}_i^{\text{succ}_i(k)-1}), \mathbf{z} - \mathbf{x}_i^{\text{succ}_i(k)} \rangle \geq 0 \quad \forall \mathbf{z} \in C_i$$

Taking into account Equation B.6 and the continuity assumption on  $\nabla_i f$ ,

$$\langle (\nabla_i f)(\bar{\mathbf{x}}), \mathbf{y} - \bar{\mathbf{x}}_i \rangle \geq 0 \quad \forall \mathbf{y} \in C_i$$

i.e.,  $(\nabla_i f)(\bar{\mathbf{x}}) \in -\mathcal{N}_{C_i}(\bar{\mathbf{x}}_i)$ . As this is true for any  $1 \leq i \leq n$ ,  $\bar{\mathbf{x}}$  is an optimal point of Problem B.2.  $\square$

## C Supplemental justifications related to Drucker–Prager constraints

### C.1 Constraints on quadrature points

The goal of this section is to justify Proposition 6.3, that Equations (6.13a — 6.13c) can be discretized as System (6.35).

*Proof.* First, it is noteworthy that the matrix  $M$  may be factorized as  $M = R^T \text{diag}(S)R$ , where  $\text{diag}(S)$  is a diagonal matrix of size  $s_d n_Q$  mapping to each quadrature point  $q$  the diagonal weight block  $\mathfrak{S}_{q,q} := 2 w_q I_{s_d}$ . That is, the diagonal coefficients of  $S$  are given by

$$S_{(q-1)s_d+k} = 2 w_q \quad \text{for } 1 \leq k \leq s_d \quad \text{and } 1 \leq q \leq n_Q.$$

Indeed, for all  $1 \leq i, j \leq n$  and for all  $1 \leq k, \ell \leq n$  let  $r(i, k) := (i-1)s_d + k$  and  $c(j, \ell) := (j-1)s_d + \ell$ . Using the quadrature rule, we have

$$\begin{aligned} M_{r,c} &= \sum_{q=1}^{n_Q} w_q \left( \mathbf{T}_{r(i,k)}(\hat{\mathbf{x}}_q) : \mathbf{T}_{c(j,\ell)}(\hat{\mathbf{x}}_q) \right) \\ &= 2 \sum_{q=1}^{n_Q} w_q \left\langle \omega_i^\tau(\hat{\mathbf{x}}_q) \chi(e_k), \omega_j^\tau(\hat{\mathbf{x}}_q) \chi(e_\ell) \right\rangle \\ &= 2 \sum_{q=1}^{n_Q} w_q \omega_i^\tau(\hat{\mathbf{x}}_q) \omega_j^\tau(\hat{\mathbf{x}}_q) (\mathbf{e}_k^\top \mathbf{e}_\ell) && \text{from (6.28)} \\ &= 2 \sum_{q=1}^{n_Q} w_q \left( \omega_i^\tau(\hat{\mathbf{x}}_q) \omega_j^\tau(\hat{\mathbf{x}}_q) \delta_{k,\ell} \right) \\ &= 2 \sum_{q=1}^{n_Q} \sum_{p=1}^{s_d} w_q \left( \omega_i^\tau(\hat{\mathbf{x}}_q) \omega_j^\tau(\hat{\mathbf{x}}_q) \delta_{k,p} \delta_{p,\ell} \right) \\ &= \sum_{q=1}^{n_Q} \sum_{p=1}^{s_d} S_{(q-1)s_d+p} \left( R_{(q-1)s_d+p,r(i,k)} R_{(q-1)s_d+p,c(j,\ell)} \right) && \text{from (6.34),} \end{aligned}$$

which we recognize as the product  $M = R^T \text{diag}(S)R$ .

Furthermore, since  $M$  is invertible, we have  $\text{rank}(M) = s_d n$ . From the rank inequality on product of matrices, we get

$$\text{rank}(M) \leq \min(\text{rank}(\text{diag}(S)), \text{rank}(R)) \leq \text{rank}(R),$$

meaning that the rank of  $R$  is equal to its number of columns  $s_d n$ . The pseudoinverse  $R^\dagger$  can therefore be computed as  $R^\dagger = (R^T R)^{-1} R^T$ , and  $R^\dagger$  plays the role of a left inverse, i.e., we have  $R^\dagger R = \mathbb{I}$ .

We thus have  $R^\dagger \hat{\underline{\lambda}} = R^\dagger R \underline{\lambda} = \underline{\lambda}$ , so (6.35a) is directly equivalent to (6.32a). Moreover, since for any strictly positive scalar  $\xi$ ,  $(\gamma, \lambda) \in \mathcal{DP}(\mu, \text{Bi}, \zeta) \iff (\xi\gamma, \lambda) \in \mathcal{DP}(\mu, \text{Bi}, \zeta)$ , the

system is not affected by a multiplication of the left-hand-side of (6.35b) with the matrix  $\text{diag}(S)$ . Multiplying then both sides of (6.35b) by  $R^\top$ , we get

$$R^\top \text{diag}(S) \hat{\underline{\gamma}} = R^\top R^{\dagger, \top} B \underline{u} + R^\top R^{\dagger, \top} \underline{k}$$

Since  $R^\top R^{\dagger, \top} = (R^\dagger R)^\top = \mathbb{I}$ , we retrieve (6.32b).

Finally, as  $\hat{\underline{\lambda}}_{[q]} = \chi^{-1}(\underline{\lambda}(\hat{\underline{x}}_q))$  and  $\hat{\underline{\gamma}}_{[q]} = \chi^{-1}(\underline{\gamma}(\hat{\underline{x}}_q))$ , we have (6.35c)  $\iff (\underline{\gamma}(\hat{\underline{x}}_q), \underline{\lambda}(\hat{\underline{x}}_q)) \in \mathcal{DP}(\mu, \text{Bi}, \zeta)$ , that is, we ensure that  $\forall \tau \in T_h$  the integral (6.33) is zero.  $\square$

## C.2 Frictional boundaries

Suppose that  $(\underline{\gamma}_{\text{RB}}, \underline{\lambda}_{\text{RB}}) \in \mathcal{DP}\mu_{\text{RB}}$ , with  $\underline{\gamma}_{\text{RB}} := \frac{1}{2}(\hat{\underline{v}}\underline{n}_{\text{RB}}^\top + \underline{n}_{\text{RB}}\hat{\underline{v}}^\top)$ , and  $\|\underline{n}_{\text{RB}}\| = 1$ .

The force induced by a stress  $\underline{\sigma}$  through a plane with normal  $\underline{n}$  is computed as  $\underline{\sigma}\underline{n}$ ; the reaction force induced by the material on the frictional boundary is therefore  $\underline{r} = \underline{\lambda}_{\text{RB}}\underline{n}_{\text{RB}}$ . In the following, we investigate the relationship between  $\underline{r}$  and  $\hat{\underline{v}}$ , the relative velocity of the boundary w.r.t. the granular material.

### C.2.1 Signorini condition

First, remark that  $\text{Tr } \underline{\gamma}_{\text{RB}} = \langle \hat{\underline{v}}, \underline{n}_{\text{RB}} \rangle = \hat{v}_N$ ,

$$\underline{r}_N = \frac{1}{3}(\text{Tr } \underline{\lambda}_{\text{RB}}) + (\text{Dev } \underline{\lambda}_{\text{RB}}) : (\underline{n}_{\text{RB}}\underline{n}_{\text{RB}}^\top)$$

and

$$\begin{aligned} (\text{Dev } \underline{\lambda}_{\text{RB}}) : (\underline{n}_{\text{RB}}\underline{n}_{\text{RB}}^\top) &\leq 2 |\text{Dev } \underline{\lambda}_{\text{RB}}| |\underline{n}_{\text{RB}}\underline{n}_{\text{RB}}^\top| \\ &\leq \frac{\mu_{\text{RB}}}{\sqrt{6}} \text{Tr } \underline{\lambda}_{\text{RB}} \times \frac{2}{\sqrt{2}} \\ &\leq \frac{\mu_{\text{RB}}}{\sqrt{3}} \text{Tr } \underline{\lambda}_{\text{RB}}. \end{aligned}$$

Therefore,

$$\left( \frac{1 - \sqrt{3}\mu_{\text{RB}}}{3} \right) \text{Tr } \underline{\lambda}_{\text{RB}} \leq \underline{r}_N \leq \left( \frac{1 + \sqrt{3}\mu_{\text{RB}}}{3} \right) \text{Tr } \underline{\lambda}_{\text{RB}}.$$

For  $\mu_{\text{RB}} < \frac{1}{\sqrt{3}}$ , we thus have

$$0 \leq \text{Tr } \underline{\gamma}_{\text{RB}} \perp \text{Tr } \underline{\lambda}_{\text{RB}} \geq 0 \implies 0 \leq \hat{v}_N \perp \underline{r}_N \geq 0,$$

i.e., the Signorini condition is satisfied.

### C.2.2 Tangential reaction

Notice that

$$\begin{aligned} \text{Dev}(\underline{\gamma}_{\text{RB}})\underline{n}_{\text{RB}} &= \frac{1}{2}\hat{\underline{v}} + \left(\frac{1}{2} - \frac{1}{3}\right)\langle \hat{\underline{v}}, \underline{n}_{\text{RB}} \rangle \underline{n}_{\text{RB}} \\ &= \frac{1}{2}\hat{\underline{v}}_T + \frac{2}{3}\hat{\underline{v}}_N \underline{n}. \end{aligned}$$

**Sliding case** First suppose that  $\hat{v}_T \neq 0$ , therefore  $\text{Dev } \underline{\gamma}_{\text{RB}} \neq 0$ , and  $\mathcal{DP}\mu_{\text{RB}}$  imposes that  $\text{Dev } \underline{\lambda}_{\text{RB}} = -\alpha \text{Dev } \underline{\gamma}_{\text{RB}}$ ,  $\alpha > 0$ . Since  $\underline{r} = \frac{1}{3} \text{Tr } \underline{\lambda}_{\text{RB}} \underline{n} + \text{Dev } \underline{\lambda}_{\text{RB}} \underline{n}$ , we can identify that

$$\underline{r}_T = -\frac{1}{2}\alpha \hat{\underline{v}}_T$$

i.e. the tangential friction force is opposed to the tangential relative velocity. Now, let us show that  $\underline{r}$  lies on the boundary of the second-order cone of aperture  $\sqrt{\frac{3}{2}}\mu_{\text{RB}}$ , i.e.,  $\|\underline{r}_T\| = \sqrt{\frac{3}{2}}\mu_{\text{RB}}\underline{r}_N$ .

From the Signori condition,  $\hat{\mathbf{v}}_N > 0$  implies  $\mathbf{r}_N = 0$ , and therefore  $\text{Tr } \boldsymbol{\lambda}_{\text{RB}} = 0$ . This means  $|\text{Dev } \boldsymbol{\lambda}_{\text{RB}}| = 0$ , and consequently  $\|\mathbf{r}_T\|$ . Our relation  $\|\mathbf{r}_T\| = \sqrt{\frac{3}{2}}\mu_{\text{RB}}\mathbf{r}_N$  is trivially satisfied.

We now have to study the case  $\hat{\mathbf{v}}_N = 0$ .  $\text{Dev } \boldsymbol{\lambda}_{\text{RB}} = -\alpha \text{Dev } \boldsymbol{\gamma}_{\text{RB}}$  means that

$$\|\mathbf{r}_T\| = \|\text{Dev } \boldsymbol{\lambda}_{\text{RB}} \mathbf{n}\| = |\text{Dev } \boldsymbol{\lambda}_{\text{RB}}| \frac{\|\text{Dev } \boldsymbol{\gamma}_{\text{RB}} \mathbf{n}_{\text{RB}}\|}{|\text{Dev } \boldsymbol{\gamma}_{\text{RB}}|}.$$

Since  $\hat{\mathbf{v}} \cdot \mathbf{n}_{\text{RB}} = 0$ ,

$$\|\text{Dev}(\boldsymbol{\gamma}_{\text{RB}}) \mathbf{n}_{\text{RB}}\| = \frac{1}{2} \|\hat{\mathbf{v}}_T\|$$

and

$$\begin{aligned} |\text{Dev}(\boldsymbol{\gamma}_{\text{RB}})|^2 &= |\boldsymbol{\gamma}_{\text{RB}}|^2 = |\hat{\mathbf{v}} \mathbf{n}_{\text{RB}}^\top|^2 - \frac{1}{4} |\hat{\mathbf{v}} \mathbf{n}_{\text{RB}}^\top - \mathbf{n}_{\text{RB}} \hat{\mathbf{v}}^\top|^2 \\ &= \frac{1}{2} \|\hat{\mathbf{v}}\|^2 - \frac{1}{4} \|\hat{\mathbf{v}} \wedge \mathbf{n}_{\text{RB}}\|^2 = \frac{1}{2} \|\hat{\mathbf{v}}\|^2 - \frac{1}{4} \|\hat{\mathbf{v}}_T\|^2 \\ &= \frac{1}{4} \|\hat{\mathbf{v}}_T\|^2 = \|\text{Dev}(\boldsymbol{\gamma}_{\text{RB}}) \mathbf{n}_{\text{RB}}\|^2. \end{aligned}$$

This means

$$\|\mathbf{r}_T\| = |\text{Dev } \boldsymbol{\lambda}_{\text{RB}}| = \frac{\mu_{\text{RB}}}{\sqrt{6}} \text{Tr } \boldsymbol{\lambda}_{\text{RB}} = \sqrt{\frac{3}{2}} \mu_{\text{RB}} \mathbf{r}_N.$$

The sliding case thus satisfies the Coulomb law with coefficient  $\sqrt{\frac{3}{2}}\mu_{\text{RB}}$ .

**Sticking** When  $\hat{\mathbf{v}} = 0$ , we cannot conclude without more information about the relationship between  $\boldsymbol{\lambda}_{\text{RB}}$  and  $\boldsymbol{\gamma}_{\text{RB}}$ . Indeed, we can only verify that

$$\begin{aligned} \|\mathbf{r}_T\| &\leq \|\text{Dev } \boldsymbol{\lambda}_{\text{RB}} \mathbf{n}\| \leq \sqrt{2} |\text{Dev } \boldsymbol{\lambda}_{\text{RB}}| \\ &\leq \frac{\mu_{\text{RB}}}{\sqrt{3}} \text{Tr } \boldsymbol{\lambda}_{\text{RB}} \leq \frac{\sqrt{3}\mu_{\text{RB}}}{1 - \sqrt{3}\mu_{\text{RB}}} \mathbf{r}_N \end{aligned}$$

i.e the reaction force has to lie inside a second-order cone of aperture  $\frac{\sqrt{3}\mu_{\text{RB}}}{1 - \sqrt{3}\mu_{\text{RB}}}$ .

This last bound does not correspond to the one derived for the sliding case (except when  $\mu_{\text{RB}} = 0$ ), but nevertheless models a coupling between the tangential and normal reaction forces.

### C.2.3 Reverse inclusion

For any  $(\hat{\mathbf{v}}, \mathbf{r}) \in \mathcal{C}_{\sqrt{\frac{3}{2}}\mu}(\mathbf{n}_{\text{RB}})$  — i.e., satisfying the 3D Coulomb law with friction coefficient  $\sqrt{\frac{3}{2}}\mu$  — we can construct a symmetric tensor  $\boldsymbol{\lambda}_{\text{RB}}$  such that  $(\boldsymbol{\gamma}_{\text{RB}}, \boldsymbol{\lambda}_{\text{RB}}) \in \mathcal{DP}(\mu)$ . Indeed, let

$$\boldsymbol{\lambda}_{\text{RB}} := (\mathbf{r}_T \mathbf{n}_{\text{RB}}^\top + \mathbf{n}_{\text{RB}} \mathbf{r}_T^\top) + \mathbf{r}_N \mathbb{I}$$

We have

$$\begin{aligned} \text{Tr } \boldsymbol{\lambda}_{\text{RB}} &= 3\mathbf{r}_N = \sqrt{6} \sqrt{\frac{3}{2}} \mathbf{r}_N \\ \boldsymbol{\lambda}_{\text{RB}} \mathbf{n}_{\text{RB}} &= \mathbf{r}_T + \mathbf{n}_{\text{RB}} (\mathbf{r}_N) = \mathbf{r} \\ \text{Dev}(\boldsymbol{\lambda}_{\text{RB}}) \mathbf{n}_{\text{RB}} &= \mathbf{r}_T \\ |\text{Dev}(\boldsymbol{\lambda}_{\text{RB}})| &= |\mathbf{r}_T \mathbf{n}_{\text{RB}}^\top + \mathbf{n}_{\text{RB}} (\mathbf{r} - \mathbf{r}_N \mathbf{n}_{\text{RB}})^\top| \\ &= \|\mathbf{r}_T\| \end{aligned}$$

It can be easily verified that for any case of the  $\mathcal{C}_{\sqrt{\frac{3}{2}}\mu}(\mathbf{n}_{\text{RB}})$  disjunctive formulation satisfied by  $\mathbf{r}$  and  $\hat{\mathbf{v}}$ , the corresponding case of  $\mathcal{DP}(\mu)$  is satisfied by  $(\boldsymbol{\gamma}_{\text{RB}}, \boldsymbol{\lambda}_{\text{RB}})$ .



# Bibliography

- Acary, V. and B. Brogliato (2008). *Numerical methods for nonsmooth dynamical systems*. English. Vol. 35. Lecture Notes in Computational and Applied Mechanics. Springer.
- Acary, V. (2009). *Toward higher order event-capturing schemes and adaptive time-step strategies for nonsmooth multibody systems*. Research Report RR-7151, p. 31. URL: <https://hal.inria.fr/inria-00440771>.
- Acary, V., M. Brémond, T. Koziara, and F. Périçon (2014). *FCLIB: a collection of discrete 3D Frictional Contact problems*. Technical Report RT-0444. INRIA, p. 34. URL: <https://hal.inria.fr/hal-00945820>.
- Ainsley, S., E. Vouga, E. Grinspun, and R. Tamstorf (2012). “Speculative Parallel Asynchronous Contact Mechanics”. In: *ACM Trans. Graph.* 31.6, 151:1–151:8. ISSN: 0730-0301. DOI: [10.1145/2366145.2366170](https://doi.org/10.1145/2366145.2366170).
- Alart, P. and A. Curnier (1991). “A mixed formulation for frictional contact problems prone to Newton like solution methods”. In: *Computer Methods in Applied Mechanics and Engineering* 92.3, pp. 353–375. ISSN: 0045-7825. DOI: [10.1016/0045-7825\(91\)90022-x](https://doi.org/10.1016/0045-7825(91)90022-x).
- Alduán, I. and M. A. Otaduy (2011). “SPH Granular Flow with Friction and Cohesion”. In: *Proceedings of the 2011 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. SCA ’11. Vancouver, British Columbia, Canada: ACM, pp. 25–32. ISBN: 978-1-4503-0923-3. DOI: [10.1145/2019406.2019410](https://doi.org/10.1145/2019406.2019410).
- Alduán, I., A. Tena, and M. A. Otaduy (2009). “Simulation of high-resolution granular media”. In: *Proc. of Congreso Español de Informática Gráfica*. Vol. 1. 4.
- Alejano, L. R. and A. Bobet (2012). “Drucker–Prager Criterion”. In: *Rock Mechanics and Rock Engineering* 45.6, pp. 995–999. ISSN: 1434-453X. DOI: [10.1007/s00603-012-0278-2](https://doi.org/10.1007/s00603-012-0278-2).
- Andersen, E. and K. Andersen (2013). *MOSEK modeling manual*.
- Andersen, E., C. Roos, and T. Terlaky (2003). “On implementing a primal-dual interior-point method for conic quadratic optimization”. In: *Mathematical Programming* 95.2, pp. 249–277. ISSN: 1436-4646. DOI: [10.1007/s10107-002-0349-3](https://doi.org/10.1007/s10107-002-0349-3).
- Andersen, S. and L. Andersen (2010). “Analysis of spatial interpolation in the material-point method”. In: *Computers & Structures* 88.7-8, pp. 506–518. ISSN: 0045-7949. DOI: [10.1016/j.compstruc.2010.01.004](https://doi.org/10.1016/j.compstruc.2010.01.004).
- Anderson, T. B. and R. Jackson (1967). “Fluid Mechanical Description of Fluidized Beds. Equations of Motion”. In: *Industrial & Engineering Chemistry Fundamentals* 6.4, pp. 527–539. ISSN: 1541-4833. DOI: [10.1021/i160024a007](https://doi.org/10.1021/i160024a007).
- Andreotti, B., Y. Forterre, and O. Pouliquen (2011). *Granular media: between fluid and solid*. Cambridge University Press.
- Andrews, M. and P. O’Rourke (1996). “The multiphase particle-in-cell (MP-PIC) method for dense particulate flows”. In: *International Journal of Multiphase Flow* 22.2, pp. 379–402. ISSN: 0301-9322. DOI: [10.1016/0301-9322\(95\)00072-0](https://doi.org/10.1016/0301-9322(95)00072-0).
- Anitescu, M. (2005). “Optimization-based simulation of nonsmooth rigid multibody dynamics”. In: *Math. Program.* 105.1, pp. 113–143. ISSN: 1436-4646. DOI: [10.1007/s10107-005-0590-7](https://doi.org/10.1007/s10107-005-0590-7).
- Anjyo, K., Y. Usami, and T. Kurihara (1992). “A Simple Method for Extracting the Natural Beauty of Hair”. In: *Computer Graphics Proceedings (Proceedings of the ACM SIGGRAPH’92 conference)*, pp. 111–120.

- Apte, S., K. Mahesh, and T. Lundgren (2003). "An Eulerian-Lagrangian Model to Simulate Two-Phase/Particulate Flows in Complex Geometries". In: *APS Division of Fluid Dynamics Meeting Abstracts*, A1.
- Aubry, J.-M. and X. Xian (2015). "Fast Implicit Simulation of Flexible Trees". In: *Mathematics for Industry*, pp. 47–61. ISSN: 2198-3518. DOI: [10.1007/978-4-431-55483-7\\_5](https://doi.org/10.1007/978-4-431-55483-7_5).
- Autodesk (2009). *3ds Max Hair & Fur*. URL: <http://usa.autodesk.com/>.
- Barbu, V. and T. Precupanu (2012). "Convex Functions". In: *Convexity and Optimization in Banach Spaces*, pp. 67–151. ISSN: 1439-7382. DOI: [10.1007/978-94-007-2247-7\\_2](https://doi.org/10.1007/978-94-007-2247-7_2).
- Bardenhagen, S. G. and E. M. Kober (2004). "The Generalized Interpolation Material Point Method". In: *CMES: Computer Modeling in Engineering & Sciences* 5.6, pp. 477–496. DOI: [10.3970/cmes.2004.005.477](https://doi.org/10.3970/cmes.2004.005.477).
- Bardenhagen, S., J. Brackbill, and D. Sulsky (2000). "The material-point method for granular materials". In: *Computer Methods in Applied Mechanics and Engineering* 187.3-4, pp. 529–541. ISSN: 0045-7825. DOI: [10.1016/S0045-7825\(99\)00338-2](https://doi.org/10.1016/S0045-7825(99)00338-2).
- Barker, T., D. G. Schaeffer, P. Bohorquez, and J. M. N. T. Gray (2015). "Well-posed and ill-posed behaviour of the -rheology for granular flow". In: *Journal of Fluid Mechanics* 779, pp. 794–818. ISSN: 1469-7645. DOI: [10.1017/jfm.2015.412](https://doi.org/10.1017/jfm.2015.412).
- Barzilai, J. and J. M. Borwein (1988). "Two-Point Step Size Gradient Methods". In: *IMA Journal of Numerical Analysis* 8.1, pp. 141–148. ISSN: 1464-3642. DOI: [10.1093/imanum/8.1.141](https://doi.org/10.1093/imanum/8.1.141).
- Batchelor, G. K. (1988). "A new theory of the instability of a uniform fluidized bed". In: *Journal of Fluid Mechanics* 193.-1, p. 75. ISSN: 1469-7645. DOI: [10.1017/S002211208800206X](https://doi.org/10.1017/S002211208800206X).
- Batty, C., F. Bertails, and R. Bridson (2007). "A fast variational framework for accurate solid-fluid coupling". In: *ACM Transactions on Graphics* 26.3, p. 100. ISSN: 0730-0301. DOI: [10.1145/1276377.1276502](https://doi.org/10.1145/1276377.1276502).
- Bedford, A. and D. Drumheller (1983). "Theories of immiscible and structured mixtures". In: *International Journal of Engineering Science* 21.8, pp. 863–960. ISSN: 0020-7225. DOI: [10.1016/0020-7225\(83\)90071-X](https://doi.org/10.1016/0020-7225(83)90071-X).
- Berga, A. and G. de Saxcé (1994). "Elastoplastic Finite Element Analysis of Soil Problems with Implicit Standard Material Constitutive Laws". In: *Revue Européenne des Éléments* 3.3, pp. 411–456. DOI: [10.1080/12506559.1994.10511137](https://doi.org/10.1080/12506559.1994.10511137). eprint: <http://www.tandfonline.com/doi/pdf/10.1080/12506559.1994.10511137>.
- Bergou, M., B. Audoly, E. Vouga, M. Wardetzky, and E. Grinspun (2010). "Discrete Viscous Threads". In: *ACM Transactions on Graphics (Proceedings of the ACM SIGGRAPH'10 Conference)*.
- Bergou, M., M. Wardetzky, S. Robinson, B. Audoly, and E. Grinspun (2008). "Discrete elastic rods". In: *ACM Transactions on Graphics (Proceedings of the ACM SIGGRAPH'08 conference)* 27.3, pp. 1–12. ISSN: 0730-0301. DOI: [http://doi.acm.org/10.1145/1360612.1360662](https://doi.org/10.1145/1360612.1360662).
- Bertails, F., B. Audoly, M.-P. Cani, B. Querleux, et al. (2006). "Super-helices for predicting the dynamics of natural hair". In: *ACM Transactions on Graphics – SIGGRAPH 2006*. DOI: [10.1145/1179352.1142012](https://doi.org/10.1145/1179352.1142012).
- Bertails-Descoubes, F., F. Cadoux, G. Daviet, and V. Acary (2011). "A nonsmooth Newton solver for capturing exact Coulomb friction in fiber assemblies". In: *ACM Transactions on Graphics* 30 (1), 6:1–6:14. ISSN: 0730-0301. DOI: [http://doi.acm.org/10.1145/1899404.1899410](https://doi.org/10.1145/1899404.1899410).
- Besson, U. (2007). "Du frottement à la tribologie: survol historique". In: *Bulletin de l'Union des physiciens* 899, pp. 1137–1154.
- Beverloo, W., H. Leniger, and J. van de Velde (1961). "The flow of granular solids through orifices". In: *Chemical Engineering Science* 15.3-4, pp. 260–269. ISSN: 0009-2509. DOI: [10.1016/0009-2509\(61\)85030-6](https://doi.org/10.1016/0009-2509(61)85030-6).
- Biot, M. A. (1955). "Theory of Elasticity and Consolidation for a Porous Anisotropic Solid". In: *Journal of Applied Physics* 26.2, p. 182. ISSN: 0021-8979. DOI: [10.1063/1.1721956](https://doi.org/10.1063/1.1721956).

- Birman, V. K., J. E. Martin, and E. Meiburg (2005). “The non-Boussinesq lock-exchange problem. Part 2. High-resolution simulations”. In: *Journal of Fluid Mechanics* 537.-1, p. 125. ISSN: 1469-7645. DOI: [10.1017/s0022112005005033](https://doi.org/10.1017/s0022112005005033).
- Bleyer, J., M. Maillard, P. de Buhan, and P. Coussot (2015). “Efficient numerical computations of yield stress fluid flows using second-order cone programming”. In: *Computer Methods in Applied Mechanics and Engineering* 283, pp. 599–614. ISSN: 0045-7825. DOI: [10.1016/j.cma.2014.10.008](https://doi.org/10.1016/j.cma.2014.10.008).
- Bonnefon, O. and G. Daviet (2011). *Quartic formulation of Coulomb 3D frictional contact*. Anglais. Tech. rep. INRIA - Laboratoire Jean Kuntzmann. URL: <http://hal.archives-ouvertes.fr/inria-00553859/en/>.
- Borwein, J. M. and D. R. Luke (2011). “Duality and convex programming”. In: *Handbook of Mathematical Methods in Imaging*. Springer, pp. 229–270.
- Bouchut, F., R. Eymard, and A. Prignet (2014). “Convergence of conforming approximations for inviscid incompressible Bingham fluid flows and related problems”. In: *Journal of Evolution Equations* 14.3, pp. 635–669. ISSN: 1424-3202. DOI: [10.1007/s00028-014-0231-9](https://doi.org/10.1007/s00028-014-0231-9).
- Bouchut, F., E. D. Fernández-Nieto, A. Mangeney, and G. Narbona-Reina (2016). “A two-phase two-layer model for fluidized granular flows with dilatancy effects”. In: *Journal of Fluid Mechanics* 801, pp. 166–221. ISSN: 1469-7645. DOI: [10.1017/jfm.2016.417](https://doi.org/10.1017/jfm.2016.417).
- Bouchut, F., I. R. Ionescu, and A. Mangeney (2016). “An analytic approach for the evolution of the static/flowing interface in viscoplastic granular flows”. working paper or preprint. URL: <https://hal-upec-upem.archives-ouvertes.fr/hal-01081213>.
- Bowden, F. P. and D. Tabor (1950). *The friction and lubrication of solids*. Vol. 1. Oxford university press.
- Boyd, L. and R. Bridson (2012). “MultiFLIP for energetic two-phase fluid simulation”. In: *ACM Transactions on Graphics* 31.2, pp. 1–12. ISSN: 0730-0301. DOI: [10.1145/2159516.2159522](https://doi.org/10.1145/2159516.2159522).
- Boyer, F. (2001). “Ecoulements diphasiques de type Cahn-Hilliard”. Thèse de doctorat dirigée par Fabrie, Pierre Mathématiques appliquées Bordeaux 1 2001. PhD thesis. Université Bordeaux 1, 191 p. URL: <http://www.theses.fr/2001BOR10509>.
- Boyer, F., É. Guazzelli, and O. Pouliquen (2011). “Unifying Suspension and Granular Rheology”. In: *Physical Review Letters* 107.18. ISSN: 1079-7114. DOI: [10.1103/physrevlett.107.188301](https://doi.org/10.1103/physrevlett.107.188301).
- Brackbill, J. and H. Ruppel (1986). “FLIP: A method for adaptively zoned, particle-in-cell calculations of fluid flows in two dimensions”. In: *Journal of Computational Physics* 65.2, pp. 314–343. ISSN: 0021-9991. DOI: [10.1016/0021-9991\(86\)90211-1](https://doi.org/10.1016/0021-9991(86)90211-1).
- Bridson, R., R. Fedkiw, and J. Anderson (2002). “Robust treatment of collisions, contact and friction for cloth animation”. In: *Proceedings of SIGGRAPH '02*. ISSN: 0730-0301. DOI: [10.1145/566570.566623](https://doi.org/10.1145/566570.566623).
- Brochu, T., E. Edwards, and R. Bridson (2012). “Efficient geometrically exact continuous collision detection”. In: *ACM Transactions on Graphics (TOG)* 31.4, p. 96.
- Brogliato, B. (1999). *Nonsmooth Mechanics: Models, Dynamics and Control (Communications and Control Engineering)*. Springer. ISBN: 1447111613.
- Bürger, R. (2000). “Phenomenological foundation and mathematical theory of sedimentation-consolidation processes”. In: *Chemical Engineering Journal* 80.1, pp. 177–188.
- Burgess, D., D. Sulsky, and J. Brackbill (1992). “Mass matrix formulation of the FLIP particle-in-cell method”. In: *Journal of Computational Physics* 103.1, pp. 1–15. ISSN: 0021-9991. DOI: [10.1016/0021-9991\(92\)90323-q](https://doi.org/10.1016/0021-9991(92)90323-q).
- Burmeister, W. (1985). “Private communication, Dresden”.
- Cadoux, F. (2009). “Optimization and convex analysis for nonsmooth dynamics”. Theses. Université Joseph-Fourier - Grenoble I. URL: <https://tel.archives-ouvertes.fr/tel-00440798>.
- Casati, R. (2015). “Quelques contributions à la modélisation numérique de structures élançées pour l’informatique graphique”. PhD thesis. Université Grenoble Alpes.



- Casati, R., G. Daviet, and F. Bertails-Descoubes (2016). *Inverse Elastic Cloth Design with Contact and Friction*. Research Report. Inria Grenoble Rhône-Alpes, Université de Grenoble. URL: <https://hal.archives-ouvertes.fr/hal-01309617>.
- Chambolle, A. and T. Pock (2010). “A First-Order Primal-Dual Algorithm for Convex Problems with Applications to Imaging”. In: *J Math Imaging Vis* 40.1, pp. 120–145. ISSN: 1573-7683. DOI: [10.1007/s10851-010-0251-1](https://doi.org/10.1007/s10851-010-0251-1).
- Chambon, G., R. Bouvarel, D. Laigle, and M. Naaim (2011). “Numerical simulations of granular free-surface flows using smoothed particle hydrodynamics”. In: *Journal of Non-Newtonian Fluid Mechanics* 166.12-13, pp. 698–712. ISSN: 0377-0257. DOI: [10.1016/j.jnnfm.2011.03.007](https://doi.org/10.1016/j.jnnfm.2011.03.007).
- Chauchat, J., S. Guillou, D. Pham Van Bang, and K. Dan Nguyen (2013). “Modelling sedimentation-consolidation in the framework of a one-dimensional two-phase flow model”. In: *Journal of Hydraulic Research* 51.3, pp. 293–305. ISSN: 1814-2079. DOI: [10.1080/00221686.2013.768798](https://doi.org/10.1080/00221686.2013.768798).
- Chauchat, J. and M. Médale (2010). “A three-dimensional numerical model for incompressible two-phase flow of a granular bed submitted to a laminar shearing flow”. In: *Computer Methods in Applied Mechanics and Engineering* 199.9-12, pp. 439–449. ISSN: 0045-7825. DOI: [10.1016/j.cma.2009.07.007](https://doi.org/10.1016/j.cma.2009.07.007).
- (2014). “A three-dimensional numerical model for dense granular flows based on the  $\mu(I)$  rheology”. In: *Journal of Computational Physics* 256, pp. 696–712. ISSN: 0021-9991. DOI: [10.1016/j.jcp.2013.09.004](https://doi.org/10.1016/j.jcp.2013.09.004).
- Chehata, D., R. Zenit, and C. R. Wassgren (2003). “Dense granular flow around an immersed cylinder”. In: *Physics of Fluids* 15.6, p. 1622. ISSN: 1070-6631. DOI: [10.1063/1.1571826](https://doi.org/10.1063/1.1571826).
- Chen, J.-S. and P. Tseng (2005). “An unconstrained smooth minimization reformulation of the second-order cone complementarity problem”. In: *Math. Program.* 104 (2), pp. 293–327. ISSN: 0025-5610. DOI: [10.1007/s10107-005-0617-0](https://doi.org/10.1007/s10107-005-0617-0).
- Choe, B., M. Choi, and H.-S. Ko (2005). “Simulating complex hair with robust collision handling”. In: *Proceedings of the 2005 ACM SIGGRAPH/Eurographics symposium on Computer animation*. SCA '05. Los Angeles, California: ACM, pp. 153–160. ISBN: 1-59593-198-8. DOI: <http://doi.acm.org/10.1145/1073368.1073389>.
- Choi, H. G. and D. D. Joseph (2001). “Fluidization by lift of 300 circular particles in plane Poiseuille flow by direct numerical simulation”. In: *Journal of Fluid Mechanics* 438. ISSN: 1469-7645. DOI: [10.1017/s0022112001004177](https://doi.org/10.1017/s0022112001004177).
- Concha, F. and R. Bürger (2002). “A century of research in sedimentation and thickening”. In: *Kona* 20, pp. 38–70.
- Coulomb, C.-A. (1781). *Théorie des machines simples: en ayant égard au frottement de leurs parties et à la roideur des cordages*.
- Cundall, P. A. and O. D. L. Strack (1979). “A discrete numerical model for granular assemblies”. In: *Géotechnique* 29.1, pp. 47–65. ISSN: 1751-7656. DOI: [10.1680/geot.1979.29.1.47](https://doi.org/10.1680/geot.1979.29.1.47).
- Dahl, P. R. (1968). *A solid friction model*. Tech. rep. The Aerospace Corporation.
- Daviet, G. (2013). *So-bogus, a C++ linear algebra and solvers library for sparse block matrices*. <http://gdaviet.fr/code/bogus>.
- Daviet, G. and F. Bertails-Descoubes (2016a). “A Semi-Implicit Material Point Method for the Continuum Simulation of Granular Materials”. In: *ACM Transactions on Graphics*. SIGGRAPH '16 Technical Papers 35.4, p. 13. DOI: [10.1145/2897824.2925877](https://doi.org/10.1145/2897824.2925877).
- (2016b). “Nonsmooth simulation of dense granular flows with pressure-dependent yield stress”. In: *Journal of Non-Newtonian Fluid Mechanics* 234, pp. 15–35. ISSN: 0377-0257. DOI: <http://dx.doi.org/10.1016/j.jnnfm.2016.04.006>.
- Daviet, G., F. Bertails-Descoubes, and L. Boissieux (2011). “A hybrid iterative solver for robustly capturing Coulomb friction in hair dynamics”. In: *ACM Transactions on Graphics* 30.6, pp. 1–12. ISSN: 0730-0301. DOI: [10.1145/2070781.2024173](https://doi.org/10.1145/2070781.2024173).
- Daviet, G., F. Bertails-Descoubes, and R. Casati (2015). “Fast cloth simulation with implicit contact and exact coulomb friction”. In: *Proceedings of the 14th ACM SIGGRAPH / Eurographics Symposium on Computer Animation - SCA '15*. DOI: [10.1145/2786784.2795139](https://doi.org/10.1145/2786784.2795139).

- De Saxcé, G. (1992). “Une généralisation de l’inégalité de Fenchel et ses applications aux lois constitutives”. In: *Comptes rendus de l’Académie des sciences. Série 2, Mécanique, Physique, Chimie, Sciences de l’univers, Sciences de la Terre* 314.2, pp. 125–129.
- De Saxcé, G. and Z.-Q. Feng (1991). “New Inequality and Functional for Contact with Friction: The Implicit Standard Material Approach”. In: *Mechanics of Structures and Machines* 19.3, pp. 301–325. DOI: [10.1080/08905459108905146](https://doi.org/10.1080/08905459108905146). eprint: <http://dx.doi.org/10.1080/08905459108905146>.
- (1998). “The bipotential method: A constructive approach to design the complete contact law with friction and improved numerical algorithms”. In: *Mathematical and Computer Modelling* 28.4-8, pp. 225–245. ISSN: 0895-7177. DOI: [10.1016/s0895-7177\(98\)00119-8](https://doi.org/10.1016/s0895-7177(98)00119-8).
- De Saxcé, G. and L. Bousshine (2002). “Implicit Standard Materials”. In: *Inelastic Behaviour of Structures under Variable Repeated Loads*, pp. 59–76. DOI: [10.1007/978-3-7091-2558-8\\_4](https://doi.org/10.1007/978-3-7091-2558-8_4).
- Denlinger, R. P. and R. M. Iverson (2001). “Flow of variably fluidized granular masses across three-dimensional terrain: 2. Numerical predictions and experimental tests”. In: *Journal of Geophysical Research* 106.B1, p. 553. ISSN: 0148-0227. DOI: [10.1029/2000jb900330](https://doi.org/10.1029/2000jb900330).
- Derouet-Jourdan, A. (2013). “Dynamic curves : from geometrical shape capture to deformable objects animation.” Theses. Université de Grenoble. URL: <https://tel.archives-ouvertes.fr/tel-01135185>.
- Derouet-Jourdan, A., F. Bertails-Descoubes, G. Daviet, and J. Thollot (2013). “Inverse dynamic hair modeling with frictional contact”. In: *ACM Transactions on Graphics* 32.6, pp. 1–10. ISSN: 0730-0301. DOI: [10.1145/2508363.2508398](https://doi.org/10.1145/2508363.2508398).
- Derouet-Jourdan, A., F. Bertails-Descoubes, and J. Thollot (2010). “Stable inverse dynamic curves”. In: *ACM Transactions on Graphics (Proceedings of the ACM SIGGRAPH Asia’10 conference)*. DOI: [10.1145/1882262.1866159](https://doi.org/10.1145/1882262.1866159).
- (2013). “Floating tangents for approximating spatial curves with  $G^1$  piecewise helices”. In: *Computer Aided Geometric Design* 30.5, pp. 490–520. ISSN: 0167-8396. DOI: [10.1016/j.cagd.2013.02.007](https://doi.org/10.1016/j.cagd.2013.02.007).
- Domnik, B. and S. P. Pudasaini (2012). “Full two-dimensional rapid chute flows of simple viscoplastic granular materials with a pressure-dependent dynamic slip-velocity and their numerical simulations”. In: *Journal of Non-Newtonian Fluid Mechanics* 173-174, pp. 72–86. ISSN: 0377-0257. DOI: [10.1016/j.jnnfm.2012.03.001](https://doi.org/10.1016/j.jnnfm.2012.03.001).
- Drew, D. A. (1983). “Mathematical Modeling of Two-Phase Flow”. In: *Annual Review of Fluid Mechanics* 15.1, pp. 261–291. ISSN: 1545-4479. DOI: [10.1146/annurev.fl.15.010183.001401](https://doi.org/10.1146/annurev.fl.15.010183.001401).
- Drew, D. A. and R. T. Lahey (1979). “Application of general constitutive principles to the derivation of multidimensional two-phase flow equations”. In: *International Journal of Multiphase Flow* 5.4, pp. 243–264. ISSN: 0301-9322. DOI: [10.1016/0301-9322\(79\)90024-7](https://doi.org/10.1016/0301-9322(79)90024-7).
- Drucker, D. and W. Prager (1952). “SOIL MECHANICS AND PLASTIC ANALYSIS OR LIMIT DESIGN”. In: *Quarterly of Applied Mathematics* 10.2, pp. 157–165. ISSN: 0033569X, 15524485.
- Dumont, S. (2012). “On enhanced descent algorithms for solving frictional multicontact problems: application to the discrete element method”. In: *Int. J. Numer. Meth. Engng* 93.11, pp. 1170–1190. ISSN: 0029-5981. DOI: [10.1002/nme.4424](https://doi.org/10.1002/nme.4424).
- Dunatunga, S. and K. Kamrin (2015). “Continuum modelling and simulation of granular flows through their many phases”. In: *Journal of Fluid Mechanics* 779, pp. 483–513. ISSN: 1469-7645. DOI: [10.1017/jfm.2015.383](https://doi.org/10.1017/jfm.2015.383).
- Duriez, C., F. Dubois, A. Kheddar, and C. Andriot (2006). “Realistic Haptic Rendering of Interacting Deformable Objects in Virtual Environments”. In: *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 1.12, pp. 36–47.
- Duriez, C. (2008). “Rendering of Frictional Contact with Deformable Environments”. In: *Haptic Rendering: Foundations, Algorithms and Applications*. Natick, MA, USA: A. K. Peters, Ltd., pp. 421–442. ISBN: 1568813325.
- Einstein, A. (1906). “Eine neue Bestimmung der Moleküldimensionen”. In: *Annalen der Physik* 324.2, pp. 289–306. ISSN: 1521-3889. DOI: [10.1002/andp.19063240204](https://doi.org/10.1002/andp.19063240204).

- Ericson, C. (2004). *Real-time collision detection*. CRC Press.
- Erleben, K. (2007). “Velocity-based shock propagation for multibody dynamics animation”. In: *ACM Transaction on Graphics* 26.2.
- Erleben, K. (2013). “Numerical methods for linear complementarity problems in physics-based animation”. In: *ACM SIGGRAPH 2013 Courses on - SIGGRAPH’13*. DOI: [10.1145/2504435.2504443](https://doi.org/10.1145/2504435.2504443).
- Etienne, J. (2004). “Simulation numérique d’écoulements gravitaires à fortes différences de densité : application aux avalanches”. Thèse de doctorat dirigée par Saramito, Pierre et Hopfinger, Emil Mathématiques appliquées Grenoble, INPG 2004. PhD thesis, 163 p. URL: <http://www.theses.fr/2004INPG0066>.
- Feng, Z.-Q., P. Joli, J.-M. Cros, and B. Magnain (2005). “The bi-potential method applied to the modeling of dynamic problems with friction”. In: *Comput Mech* 36.5, pp. 375–383. ISSN: 1432-0924. DOI: [10.1007/s00466-005-0663-8](https://doi.org/10.1007/s00466-005-0663-8).
- Fischer, A. (1992). “A special newton-type optimization method”. In: *Optimization* 24.3-4, pp. 269–284. ISSN: 1029-4945. DOI: [10.1080/02331939208843795](https://doi.org/10.1080/02331939208843795).
- Folgar, F. and C. L. Tucker (1984). “Orientation Behavior of Fibers in Concentrated Suspensions”. In: *Journal of Reinforced Plastics and Composites* 3, pp. 98–119. DOI: [10.1177/073168448400300201](https://doi.org/10.1177/073168448400300201).
- Fortin, M. and R. Glowinski (1983). *Augmented Lagrangian Methods: Applications to the Numerical Solution of Boundary-Value Problems (Studies in Mathematics and Its Applications)*. Elsevier Science Ltd. ISBN: 0444866809.
- Frankel, N. A. and A. Acrivos (1967). “On the viscosity of a concentrated suspension of solid spheres”. In: *Chem. Eng. Sci.* 22, pp. 847–852.
- Frigaard, I. and C. Nouar (2005). “On the usage of viscosity regularisation methods for viscoplastic fluid flow computation”. In: *Journal of Non-Newtonian Fluid Mechanics* 127.1, pp. 1–26. ISSN: 0377-0257. DOI: [10.1016/j.jnnfm.2005.01.003](https://doi.org/10.1016/j.jnnfm.2005.01.003).
- Fukushima, M., Z.-Q. Luo, and P. Tseng (2002). “Smoothing Functions for Second-Order-Cone Complementarity Problems”. In: *SIAM Journal on Optimization* 12 (2), pp. 436–460. ISSN: 1052-6234. DOI: [http://dx.doi.org/10.1137/S1052623400380365](https://dx.doi.org/10.1137/S1052623400380365).
- GDR MiDi (2004). “On dense granular flows”. In: *The European Physical Journal E* 14.4, pp. 341–365. ISSN: 1292-895X. DOI: [10.1140/epje/i2003-10153-0](https://doi.org/10.1140/epje/i2003-10153-0).
- Glowinski, R. and A. Marroco (1975). “Sur l’approximation, par éléments finis d’ordre un, et la résolution, par pénalisation-dualité d’une classe de problèmes de Dirichlet non linéaires”. fre. In: *ESAIM: Mathematical Modelling and Numerical Analysis* 9.R2, pp. 41–76. URL: <http://eudml.org/doc/193269>.
- Goldstein, T., B. O’Donoghue, S. Setzer, and R. Baraniuk (2014). “Fast Alternating Direction Optimization Methods”. In: *SIAM Journal on Imaging Sciences* 7.3, pp. 1588–1623. ISSN: 1936-4954. DOI: [10.1137/120896219](https://doi.org/10.1137/120896219).
- Gornowicz, G. and S. Borac (2015). “Efficient and Stable Approach to Elasticity and Collisions for Hair Animation”. In: *Proceedings of the 2015 Symposium on Digital Production*. DigiPro ’15. Los Angeles, California: ACM, pp. 41–49. ISBN: 978-1-4503-3718-2. DOI: [10.1145/2791261.2791271](https://doi.org/10.1145/2791261.2791271).
- Grinspun, E., A. N. Hirani, M. Desbrun, and P. Schröder (2003). “Discrete Shells”. In: *Proceedings of the 2003 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. SCA ’03. San Diego, California: Eurographics Association, pp. 62–67. ISBN: 1-58113-659-5.
- Grippo, L. and M. Sciandrone (2000). “On the convergence of the block nonlinear Gauss–Seidel method under convex constraints”. In: *Operations Research Letters* 26.3, pp. 127–136. ISSN: 0167-6377. DOI: [10.1016/S0167-6377\(99\)00074-7](https://doi.org/10.1016/S0167-6377(99)00074-7).
- Guennebaud, G., B. Jacob, et al. (2010). *Eigen v3, a C++ template library for linear algebra*. <http://eigen.tuxfamily.org>.
- Hadap, S. (2006). “Oriented strands - dynamics of stiff multi-body system”. In: *ACM SIGGRAPH - EG Symposium on Computer Animation (SCA’06)*. ACM-EG SCA, pp. 91–100.
- Hadap, S. and N. Magnenat-Thalmann (2001). “Modeling Dynamic Hair as a Continuum”. In: *Computer Graphics Forum* 20.3. Proceedings of Eurographics’01, pp. 329–338.

- Haddouni, M. (2015). “Algorithmes de résolution de la dynamique du contact avec impact et frottement”. Thèse de doctorat dirigée par Brogliato, Bernard Matériaux, mécanique, génie civil, électrochimie Grenoble Alpes 2015. PhD thesis. URL: <http://www.theses.fr/2015GREAI022>.
- Hanes, D. M. and A. J. Bowen (1985). “A granular-fluid model for steady intense bed-load transport”. In: *Journal of Geophysical Research* 90.C5, p. 9149. ISSN: 0148-0227. DOI: [10.1029/jc090ic05p09149](https://doi.org/10.1029/jc090ic05p09149).
- Harlow, F. H. (1963). “The particle-in-cell method for numerical solution of problems in fluid dynamics”. In: *Proceedings of Symposia in Applied Mathematics*, pp. 269–288. ISSN: 0160-7634. DOI: [10.1090/psapm/015/9942](https://doi.org/10.1090/psapm/015/9942). URL: <http://dx.doi.org/10.1090/psapm/015/9942>.
- Harmon, D., E. Vouga, B. Smith, R. Tamstorf, and E. Grinspun (2009). “Asynchronous Contact Mechanics”. In: *ACM Trans. Graph.* 28.3, 87:1–87:12. ISSN: 0730-0301. DOI: [10.1145/1531326.1531393](https://doi.org/10.1145/1531326.1531393).
- Harmon, D., E. Vouga, R. Tamstorf, and E. Grinspun (2008). “Robust treatment of simultaneous collisions”. In: *ACM Transactions on Graphics* 27.3, p. 1. ISSN: 0730-0301. DOI: [10.1145/1360612.1360622](https://doi.org/10.1145/1360612.1360622).
- Harris, S. E. and D. G. Crighton (1994). “Solitons, solitary waves, and voidage disturbances in gas-fluidized beds”. In: *Journal of Fluid Mechanics* 266, pp. 243–276. ISSN: 1469-7645. DOI: [10.1017/S0022112094000996](https://doi.org/10.1017/S0022112094000996).
- Haslinger, J. (1983). “Approximation of the signorini problem with friction, obeying the coulomb law”. In: *Mathematical Methods in the Applied Sciences* 5.1, pp. 422–437. ISSN: 0170-4214. DOI: [10.1002/mma.1670050127](https://doi.org/10.1002/mma.1670050127).
- Herrera, T. L., A. Zinke, and A. Weber (2012). “Lighting Hair From The Inside: A Thermal Approach To Hair Reconstruction”. In: *ACM Transactions on Graphics (Proceedings SIGGRAPH Asia 2012)*. Vol. 31. 6. DOI: [10.1145/2366145.2366165](https://doi.org/10.1145/2366145.2366165).
- Heyman, J., R. Delannay, H. Tabuteau, and A. Valance (2016). “Compressibility regularizes the  $\mu(I)$  rheology for granular flows”. In: eprint: [1609.01502](https://arxiv.org/abs/1609.01502).
- Heyn, T. D. (2013). “On the modeling, simulation, and visualization of many-body dynamics problems with friction and contact”. PhD thesis. THE UNIVERSITY OF WISCONSIN-MADISON.
- Hiriart-Urruty, J.-B. and C. Lemaréchal (1993). *Convex analysis and minimization algorithms I: Fundamentals*. Vol. 305. Springer Science & Business Media.
- Hsu, T.-J., J. T. Jenkins, and P. L.-F. Liu (2004). “On two-phase sediment transport: sheet flow of massive particles”. In: *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 460.2048, pp. 2223–2250. ISSN: 1471-2946. DOI: [10.1098/rspa.2003.1273](https://doi.org/10.1098/rspa.2003.1273).
- Iben, H., M. Meyer, L. Petrovic, O. Soares, et al. (2013). “Artistic Simulation of Curly Hair”. In: *Proceedings of the 12th ACM SIGGRAPH/Eurographics Symposium on Computer Animation. SCA '13*. Anaheim, California: ACM, pp. 63–71. ISBN: 978-1-4503-2132-7. DOI: [10.1145/2485895.2485913](https://doi.org/10.1145/2485895.2485913).
- Ihmsen, M., A. Wahl, and M. Teschner (2013). “A Lagrangian framework for simulating granular material with high detail”. In: *Computers & Graphics* 37.7, pp. 800–808. ISSN: 0097-8493. DOI: [10.1016/j.cag.2013.04.010](https://doi.org/10.1016/j.cag.2013.04.010).
- Ionescu, I. R. (2010). “Onset and dynamic shallow flow of a viscoplastic fluid on a plane slope”. In: *Journal of Non-Newtonian Fluid Mechanics* 165.19-20, pp. 1328–1341. ISSN: 0377-0257. DOI: [10.1016/j.jnnfm.2010.06.016](https://doi.org/10.1016/j.jnnfm.2010.06.016).
- (2013). “Augmented Lagrangian for shallow viscoplastic flow with topography”. In: *Journal of Computational Physics* 242, pp. 544–560. ISSN: 0021-9991. DOI: [10.1016/j.jcp.2013.02.029](https://doi.org/10.1016/j.jcp.2013.02.029).
- Ionescu, I. R., A. Mangeney, F. Bouchut, and O. Roche (2015). “Viscoplastic modeling of granular column collapse with pressure-dependent rheology”. In: *Journal of Non-Newtonian Fluid Mechanics* 219, pp. 1–18. ISSN: 0377-0257. DOI: [10.1016/j.jnnfm.2015.02.006](https://doi.org/10.1016/j.jnnfm.2015.02.006).
- Iverson, R. M. (1997). “The physics of debris flows”. In: *Reviews of Geophysics* 35.3, pp. 245–296. ISSN: 8755-1209. DOI: [10.1029/97rg00426](https://doi.org/10.1029/97rg00426).



- Jackson, R. (1963). “The mechanics of fluidised beds. Part I. The stability of the state of uniform fluidisation”. In: *Transactions of the Institution of Chemical Engineers*.
- Jackson, R. (2000). *The dynamics of fluidized particles*. Cambridge University Press.
- Jean, M. (1999). “The non-smooth contact dynamics method”. In: *Computer Methods in Applied Mechanics and Engineering* 177.3-4, pp. 235–257. ISSN: 0045-7825. DOI: [10.1016/S0045-7825\(98\)00383-1](https://doi.org/10.1016/S0045-7825(98)00383-1).
- Jean, M. and J.-J. Moreau (1987). “Dynamics in the Presence of Unilateral Contacts and Dry Friction: A Numerical Approach”. In: *Unilateral Problems in Structural Analysis — 2: Proceedings of the Second Meeting on Unilateral Problems in Structural Analysis, Prescudin, June 17–20, 1985*. Ed. by G. Del Piero and F. Maceri. Vienna: Springer Vienna, pp. 151–196. ISBN: 978-3-7091-2967-8. DOI: [10.1007/978-3-7091-2967-8\\_10](https://doi.org/10.1007/978-3-7091-2967-8_10).
- (1992). “Unilaterality and dry friction in the dynamics of rigid bodies collections”. In: *Proceedings of the Contact Mechanics International Symposium*. Ed. by A. Curnier, pp. 31–48.
- Jeffery, G. B. (1922). “The Motion of Ellipsoidal Particles Immersed in a Viscous Fluid”. In: *Proceedings of the Royal Society of London. Series A* 102, pp. 161–179. DOI: [10.1098/rspa.1922.0078](https://doi.org/10.1098/rspa.1922.0078).
- Jiang, C., C. Schroeder, A. Selle, J. Teran, and A. Stomakhin (2015). “The Affine Particle-in-cell Method”. In: *ACM Transactions on Graphics* 34.4, pp. 1–10. ISSN: 0730-0301. DOI: [10.1145/2766996](https://doi.org/10.1145/2766996).
- Jiang, C., C. Schroeder, J. Teran, A. Stomakhin, and A. Selle (2016). “The material point method for simulating continuum materials”. In: *ACM SIGGRAPH 2016 Courses on - SIGGRAPH '16*. DOI: [10.1145/2897826.2927348](https://doi.org/10.1145/2897826.2927348).
- Jop, P., Y. Forterre, and O. Pouliquen (2006). “A constitutive law for dense granular flows”. In: *Nature* 441.7094, pp. 727–730. ISSN: 1476-4687. DOI: [10.1038/nature04801](https://doi.org/10.1038/nature04801).
- Jourdan, F., P. Alart, and M. Jean (1998). “A Gauss-Seidel like algorithm to solve frictional contact problems”. In: *Computer Methods in Applied Mechanics and Engineering* 155, pp. 31–47. DOI: [10.1016/S0045-7825\(97\)00137-0](https://doi.org/10.1016/S0045-7825(97)00137-0).
- Kaufman, D. M., S. Sueda, D. James, and D. Pai (2008). “Staggered Projections for Frictional Contact in Multibody Systems”. In: *ACM Transactions on Graphics (Proceedings of the ACM SIGGRAPH Asia'08 conference)* 27.5, 164:1–164:11.
- Kaufman, D. M., R. Tamstorf, B. Smith, J.-M. Aubry, and E. Grinspun (2014). “Adaptive non-linearity for collisions in complex rod assemblies”. In: *ACM Transactions on Graphics* 33.4, pp. 1–12. ISSN: 0730-0301. DOI: [10.1145/2601097.2601100](https://doi.org/10.1145/2601097.2601100).
- Klar, G., T. Gast, A. Pradhana, C. Fu, et al. (2016). “Drucker-Prager Elastoplasticity for Sand Animation”. In: *ACM Trans. Graph.* 35.4.
- Klarbring, A. (1987). “Contact Problems with Friction by Linear Complementarity”. In: *Unilateral Problems in Structural Analysis – 2*, pp. 197–219. DOI: [10.1007/978-3-7091-2967-8\\_11](https://doi.org/10.1007/978-3-7091-2967-8_11).
- (1990). “Examples of non-uniqueness and non-existence of solutions to quasistatic contact problems with friction”. In: *Ingenieur-Archiv* 60.8, pp. 529–541. ISSN: 1432-0681. DOI: [10.1007/BF00541909](https://doi.org/10.1007/BF00541909).
- Konagai, K. and J. Johansson (2001). “Two dimensional Lagrangian Particle Finite Difference Method for modeling large soil deformations”. In: *Journal of Structural Mechanics and Earthquake Engineering* 18, pp. 105–110.
- Krieger, I. M. (1972). “Rheology of monodisperse latices”. In: *Advances in Colloid and Interface Science* 3.2, pp. 111–136. ISSN: 0001-8686. DOI: [10.1016/0001-8686\(72\)80001-0](https://doi.org/10.1016/0001-8686(72)80001-0).
- Kynch, G. J. (1952). “A theory of sedimentation”. In: *Transactions of the Faraday society* 48, pp. 166–176.
- Lagrée, P.-Y., L. Staron, and S. Popinet (2011). “The granular column collapse as a continuum: validity of a two-dimensional Navier–Stokes model with a  $\mu(I)$ -rheology”. In: *Journal of Fluid Mechanics* 686, pp. 378–408. ISSN: 1469-7645. DOI: [10.1017/jfm.2011.335](https://doi.org/10.1017/jfm.2011.335).
- Lajeunesse, E., J. B. Monnier, and G. M. Homsy (2005). “Granular slumping on a horizontal surface”. In: *Physics of Fluids* 17.10, p. 103302. ISSN: 1070-6631. DOI: [10.1063/1.2087687](https://doi.org/10.1063/1.2087687).
- Lenaerts, T., B. Adams, and P. Dutré (2008). “Porous flow in particle-based fluid simulations”. In: *ACM SIGGRAPH 2008 papers on - SIGGRAPH '08*. DOI: [10.1145/1399504.1360648](https://doi.org/10.1145/1399504.1360648).

- Lenaerts, T. and P. Dutré (2009). “Mixing fluids and granular materials”. In: *Computer Graphics Forum*. Vol. 28. 2, pp. 213–218. DOI: [10.1111/j.1467-8659.2009.01360.x](https://doi.org/10.1111/j.1467-8659.2009.01360.x).
- Mankoc, C., A. Janda, R. Arévalo, J. M. Pastor, et al. (2007). “The flow rate of granular materials through an orifice”. In: *Granular Matter* 9.6, pp. 407–414. ISSN: 1434-7636. DOI: [10.1007/s10035-007-0062-2](https://doi.org/10.1007/s10035-007-0062-2).
- Manninen, M., V. Taivassalo, and S. Kallio (1996). “On the Mixture model for Multiphase Flow”. In: *VTT Publications* 288. URL: <http://www.vtt.fi/inf/pdf/publications/1996/P288.pdf>.
- Mast, C. M. (2013). “Modeling landslide-induced flow interactions with structures using the material point method”. PhD thesis. University of Washington.
- Mast, C. M., P. Arduino, P. Mackenzie-Helnwein, and G. R. Miller (2014). “Simulating granular column collapse using the Material Point Method”. In: *Acta Geotechnica* 10.1, pp. 101–116. ISSN: 1861-1133. DOI: [10.1007/s11440-014-0309-0](https://doi.org/10.1007/s11440-014-0309-0).
- Maurin, R., J. Chauchat, B. Chareyre, and P. Frey (2015). “A minimal coupled fluid-discrete element model for bedload transport”. In: *Phys. Fluids* 27.11, p. 113302. ISSN: 1089-7666. DOI: [10.1063/1.4935703](https://doi.org/10.1063/1.4935703).
- Mazhar, H., T. Heyn, D. Negrut, and A. Tasora (2015). “Using Nesterov’s Method to Accelerate Multibody Dynamics with Friction and Contact”. In: *ACM Transactions on Graphics* 34.3, pp. 1–14. ISSN: 0730-0301. DOI: [10.1145/2735627](https://doi.org/10.1145/2735627).
- McAdams, A., A. Selle, K. Ward, E. Sifakis, and J. Teran (2009). “Detail preserving continuum simulation of straight hair”. In: *ACM Transactions on Graphics Proceedings of the SIGGRAPH’09 conference*. Association for Computing Machinery, p. 1. ISBN: 9781605587264. DOI: [10.1145/1576246.1531368](https://doi.org/10.1145/1576246.1531368).
- Michaels, A. S. and J. C. Bolger (1962). “Settling Rates and Sediment Volumes of Flocculated Kaolin Suspensions”. In: *Industrial & Engineering Chemistry Fundamentals* 1.1, pp. 24–33. ISSN: 1541-4833. DOI: [10.1021/i160001a004](https://doi.org/10.1021/i160001a004).
- Miguel, E., R. Tamstorf, D. Bradley, S. C. Schwartzman, et al. (2013). “Modeling and Estimation of Internal Friction in Cloth”. In: *ACM Trans. Graph.* 32.6, 212:1–212:10. ISSN: 0730-0301. DOI: [10.1145/2508363.2508389](https://doi.org/10.1145/2508363.2508389).
- Moreau, J.-J. (1965). “Proximité et dualité dans un espace hilbertien”. fre. In: *Bulletin de la Société Mathématique de France* 93, pp. 273–299. URL: <http://eudml.org/doc/87067>.
- (1966–1967). “Fonctionnelles convexes”. fre. In: *Séminaire Jean Leray* 2, pp. 1–108. URL: <http://eudml.org/doc/112529>.
- (1988). “Unilateral Contact and Dry Friction in Finite Freedom Dynamics”. In: *Nonsmooth Mechanics and Applications*. Ed. by J.-J. Moreau and P. D. Panagiotopoulos. Vienna: Springer Vienna, pp. 1–82. ISBN: 978-3-7091-2624-0. DOI: [10.1007/978-3-7091-2624-0\\_1](https://doi.org/10.1007/978-3-7091-2624-0_1).
- (1999). “Numerical aspects of the sweeping process”. In: *Computer Methods in Applied Mechanics and Engineering* 177.3–4, pp. 329–349. ISSN: 0045-7825. DOI: [http://dx.doi.org/10.1016/S0045-7825\(98\)00387-9](https://doi.org/10.1016/S0045-7825(98)00387-9).
- Müller, M., B. Solenthaler, R. Keiser, and M. Gross (2005). “Particle-based fluid-fluid interaction”. In: *Proceedings of the 2005 ACM SIGGRAPH/Eurographics symposium on Computer animation - SCA ’05*. DOI: [10.1145/1073368.1073402](https://doi.org/10.1145/1073368.1073402).
- Munson, T. S., F. Facchinei, M. C. Ferris, A. Fischer, and C. Kanzow (2001). “The semismooth algorithm for large scale complementarity problems”. In: *INFORMS Journal on Computing* 13.4, pp. 294–311.
- Narain, R., A. Golas, and M. C. Lin (2010). “Free-flowing granular materials with two-way solid coupling”. In: *ACM Transactions on Graphics* 29.6, pp. 1–10. DOI: [10.1145/1882261.1866195](https://doi.org/10.1145/1882261.1866195).
- Narain, R., M. Overby, and G. E. Brown (2016). “ADMM  $\supseteq$  Projective Dynamics: Fast Simulation of General Constitutive Models”. In: *Eurographics/ACM SIGGRAPH Symposium on Computer Animation*. Ed. by L. Kavan and C. Wojtan. The Eurographics Association. ISBN: 978-3-03868-009-3. DOI: [10.2312/sca.20161219](https://doi.org/10.2312/sca.20161219).
- Nesterov, Y. (1983). “A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ”. In: *Soviet Mathematics Doklady* 27.2, pp. 372–376.

- Ngo Ngoc, C. and S. Boivin (2004). *Nonlinear Cloth Simulation*. Tech. rep. Inria.
- Nielsen, M. B. and O. Østerby (2013). “A two-continua approach to Eulerian simulation of water spray”. In: *ACM Transactions on Graphics* 32.4, p. 1. ISSN: 0730-0301. DOI: [10.1145/2461912.2461918](#).
- Otaduy, M. A., R. Tamstorf, D. Steinemann, and M. Gross (2009). “Implicit Contact Handling for Deformable Objects”. In: *Computer Graphics Forum* 28.2, pp. 559–568. ISSN: 1467-8659. DOI: [10.1111/j.1467-8659.2009.01396.x](#).
- Pai, D. (2002). “Strands: Interactive Simulation of Thin Solids using Cosserat Models”. In: *Computer Graphics Forum* 21.3. Proceedings of Eurographics’02, pp. 347–352.
- Pailha, M., M. Nicolas, and O. Pouliquen (2008). “Initiation of underwater granular avalanches: Influence of the initial volume fraction”. In: *Physics of Fluids* 20.11, p. 111701. ISSN: 1070-6631. DOI: [10.1063/1.3013896](#).
- Parikh, N. and S. P. Boyd (2014). “Proximal Algorithms.” In: *Foundations and Trends in optimization* 1.3, pp. 127–239.
- Parker, G., Y. Fukushima, and H. M. Pantin (1986). “Self-accelerating turbidity currents”. In: *Journal of Fluid Mechanics* 171.-1, p. 145. ISSN: 1469-7645. DOI: [10.1017/s0022112086001404](#).
- Pellegrini, F. and J. Roman (1996). “Scotch: A software package for static mapping by dual recursive bipartitioning of process and architecture graphs”. In: *Lecture Notes in Computer Science*, pp. 493–498. ISSN: 1611-3349. DOI: [10.1007/3-540-61142-8\\_588](#).
- Pietro, D. A. D. and A. Ern (2011). *Mathematical Aspects of Discontinuous Galerkin Methods (Mathématiques et Applications)*. Springer. ISBN: 3642229794.
- Pignatelli, F., M. Nicolas, É. Guazzelli, and D. Saintillan (2009). “Falling jets of particles in viscous fluids”. In: *Physics of Fluids* 21.12, p. 123303. ISSN: 1070-6631. DOI: [10.1063/1.3276235](#).
- Pitman, E. B. and L. Le (2005). “A two-fluid model for avalanche and debris flows”. In: *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 363.1832, pp. 1573–1601. ISSN: 1471-2962. DOI: [10.1098/rsta.2005.1596](#).
- Plante, E., M.-P. Cani, and P. Poulin (2001). “A Layered Wisp Model for Simulating Interactions inside Long Hair”. In: *EG workshop on Computer Animation and Simulation (EG CAS’01)*. Computer Science. Springer, pp. 139–148.
- Potra, F. A. and S. J. Wright (2000). “Interior-point methods”. In: *Journal of Computational and Applied Mathematics* 124.1-2, pp. 281–302. ISSN: 0377-0427. DOI: [10.1016/s0377-0427\(00\)00433-7](#).
- Qi, L. and J. Sun (1993). “A nonsmooth version of Newton’s method”. In: *Mathematical Programming* 58, pp. 353–367.
- Radjai, F., D. Wolf, M. Jean, and J.-J. Moreau (1998). “Bimodal Character of Stress Transmission in Granular Packings”. In: *Phys. Rev. Lett.* 80 (1), pp. 61–64. DOI: [10.1103/PhysRevLett.80.61](#).
- Ram, D., T. Gast, C. Jiang, C. Schroeder, et al. (2015). “A material point method for viscoelastic fluids, foams and sponges”. In: *Proceedings of the 14th ACM SIGGRAPH / Eurographics Symposium on Computer Animation - SCA ’15*. DOI: [10.1145/2786784.2786798](#).
- Raous, M., P. Chabrand, and F. Lebon (1988). “Numerical methods for frictional contact problems and applications”. In: *Journal de Mécanique Théorique et Appliquée* 7.1.
- Ren, B., C. Li, X. Yan, M. C. Lin, et al. (2014). “Multiple-Fluid SPH Simulation Using a Mixture Model”. In: *ACM Transactions on Graphics* 33.5, pp. 1–11. ISSN: 0730-0301. DOI: [10.1145/2645703](#).
- Renouf, M., V. Acary, and G. Dumont (2005). “3D Frictional Contact and Impact Multibody Dynamics. A Comparison of Algorithms Suitable for Real-time Applications”. In: *ECCOMAS*.
- Renouf, M. and P. Alart (2005). “Conjugate gradient type algorithms for frictional multi-contact problems: applications to granular materials”. In: *Computer Methods in Applied Mechanics and Engineering* 194.18-20, pp. 2019–2041. ISSN: 0045-7825. DOI: [10.1016/j.cma.2004.07.009](#).
- Revil-Baudard, T. and J. Chauchat (2013). “A two-phase model for sheet flow regime based on dense granular flow rheology”. In: *Journal of Geophysical Research: Oceans* 118.2, pp. 619–634. ISSN: 2169-9275. DOI: [10.1029/2012jc008306](#).

- Richardson, J. F. and W. N. Zaki (1954). “Sedimentation and fluidisation: Part I. Transactions of Institution of Chemical Engineers”. In: *Transactions of Institution of Chemical Engineers* 32, pp. 35–53.
- Rockafellar, R. T. (1993). “Lagrange Multipliers and Optimality”. In: *SIAM Review* 35.2, pp. 183–238. ISSN: 00361445. DOI: [10.2307/2133143](https://doi.org/10.2307/2133143).
- Rosenblum, R., W. Carlson, and E. Tripp (1991). “Simulating the Structure and Dynamics of Human Hair: Modeling, Rendering, and Animation”. In: *The Journal of Visualization and Computer Animation* 2.4, pp. 141–148.
- Roux, S. and F. Radjai (1998). “Texture-Dependent Rigid-Plastic Behavior”. In: *Physics of Dry Granular Media*, pp. 229–236. DOI: [10.1007/978-94-017-2653-5\\_13](https://doi.org/10.1007/978-94-017-2653-5_13).
- Rungjiratananon, W., Z. Szego, Y. Kanamori, and T. Nishita (2008). “Real-time Animation of Sand-Water Interaction”. In: *Computer Graphics Forum* 27.7, pp. 1887–1893. ISSN: 1467-8659. DOI: [10.1111/j.1467-8659.2008.01336.x](https://doi.org/10.1111/j.1467-8659.2008.01336.x).
- Saad, Y. (2003). *Iterative Methods for Sparse Linear Systems, Second Edition*. Society for Industrial and Applied Mathematics. ISBN: 0898715342.
- Saramito, P. (2013). “Méthodes numériques en fluides complexes : théorie et algorithmes”. Lecture. Grenoble. URL: <https://cel.archives-ouvertes.fr/cel-00673816>.
- (2015). “Efficient C++ finite element computing with Rheolef”. Lecture. Grenoble, France, France. URL: <https://cel.archives-ouvertes.fr/cel-00573970>.
- (2016). “A damped Newton algorithm for computing viscoplastic fluid flows”. In: *Journal of Non-Newtonian Fluid Mechanics*. ISSN: 0377-0257. DOI: [10.1016/j.jnnfm.2016.05.007](https://doi.org/10.1016/j.jnnfm.2016.05.007).
- Saramito, P. and N. Roquet (2001). “An adaptive finite element method for viscoplastic fluid flows in pipes”. In: *Computer Methods in Applied Mechanics and Engineering* 190.40-41, pp. 5391–5412. ISSN: 0045-7825. DOI: [10.1016/s0045-7825\(01\)00175-x](https://doi.org/10.1016/s0045-7825(01)00175-x).
- Saramito, P. and A. Wachs (2016). “Progress in numerical simulation of yield stress fluid flows”. working paper or preprint. URL: <https://hal.archives-ouvertes.fr/hal-01375720>.
- Savage, S. B. and K. Hutter (1989). “The motion of a finite mass of granular material down a rough incline”. In: *Journal of Fluid Mechanics* 199.-1, p. 177. ISSN: 1469-7645. DOI: [10.1017/s0022112089000340](https://doi.org/10.1017/s0022112089000340).
- Seguin, A., C. Coulais, F. Martinez, Y. Bertho, and P. Gondret (2016). “Local rheological measurements in the granular flow around an intruder”. In: *Physical Review E* 93.1. ISSN: 2470-0053. DOI: [10.1103/physreve.93.012904](https://doi.org/10.1103/physreve.93.012904).
- Selle, A., M. Lentine, and R. Fedkiw (2008). “A mass spring model for hair simulation”. In: *ACM Transactions on Graphics (Proceedings of the ACM SIGGRAPH’08 conference)* 27.3, pp. 1–11. ISSN: 0730-0301. DOI: [http://doi.acm.org/10.1145/1360612.1360663](https://doi.org/http://doi.acm.org/10.1145/1360612.1360663).
- Shu, C.-W. (2009). “High Order Weighted Essentially Nonoscillatory Schemes for Convection Dominated Problems”. In: *SIAM Review* 51.1, pp. 82–126. ISSN: 1095-7200. DOI: [10.1137/070679065](https://doi.org/10.1137/070679065).
- Sigal, L., M. Mahler, S. Diaz, K. McIntosh, et al. (2015). “A perceptual control space for garment simulation”. In: *ACM Transactions on Graphics* 34.4, 117:1–117:10. ISSN: 0730-0301. DOI: [10.1145/2766971](https://doi.org/10.1145/2766971).
- Silcowitz, M., S. Niebe, and K. Erleben (2009). “Nonsmooth Newton Method for Fischer function reformulation of contact force problems for interactive rigid body simulation”. In: pp. 105–114. DOI: [10.2312/PE/vriphys/vriphys09/105-114](https://doi.org/10.2312/PE/vriphys/vriphys09/105-114).
- Simo, J. and T. Hughes (2000). *Computational Inelasticity (Interdisciplinary Applied Mathematics)* (v. 7). Springer. ISBN: 0387975209.
- Smith, B., D. M. Kaufman, E. Vouga, R. Tamstorf, and E. Grinspun (2012). “Reflections on Simultaneous Impact”. In: *ACM Trans. Graph.* 31.4, 106:1–106:12. ISSN: 0730-0301. DOI: [10.1145/2185520.2185602](https://doi.org/10.1145/2185520.2185602).
- Spillmann, J. and M. Teschner (2007). “CoRdE: Cosserat rod elements for the dynamic simulation of one-dimensional elastic objects”. In: *ACM SIGGRAPH - EG Symposium on Computer Animation (SCA’07)*. ACM-EG SCA, pp. 63–72.



- Staron, L. and E. J. Hinch (2005). “Study of the collapse of granular columns using two-dimensional discrete-grain simulation”. In: *Journal of Fluid Mechanics* 545.-1, p. 1. ISSN: 1469-7645. DOI: [10.1017/s0022112005006415](https://doi.org/10.1017/s0022112005006415).
- Staron, L., P.-Y. Lagrée, and S. Popinet (2012). “The granular silo as a continuum plastic flow: The hour-glass vs the clepsydra”. In: *Physics of Fluids* 24, 103301. ISSN: 1070-6631. DOI: [10.1063/1.4757390](https://doi.org/10.1063/1.4757390).
- (2014). “Continuum simulation of the discharge of the granular silo”. In: *The European Physical Journal E* 37.1. ISSN: 1292-895X. DOI: [10.1140/epje/i2014-14005-6](https://doi.org/10.1140/epje/i2014-14005-6).
- Steffen, M., R. M. Kirby, and M. Berzins (2008). “Analysis and reduction of quadrature errors in the material point method (MPM)”. In: *Int. J. Numer. Meth. Engng* 76.6, pp. 922–948. ISSN: 1097-0207. DOI: [10.1002/nme.2360](https://doi.org/10.1002/nme.2360).
- Stomakhin, A., C. Schroeder, L. Chai, J. Teran, and A. Selle (2013). “A material point method for snow simulation”. In: *ACM Transactions on Graphics* 32.4, pp. 1–10. DOI: [10.1145/2461912.2461948](https://doi.org/10.1145/2461912.2461948).
- Stomakhin, A., C. Schroeder, C. Jiang, L. Chai, et al. (2014). “Augmented MPM for phase-change and varied materials”. In: *ACM Transactions on Graphics* 33.4, pp. 1–11. ISSN: 0730-0301. DOI: [10.1145/2601097.2601176](https://doi.org/10.1145/2601097.2601176).
- Sulsky, D. (1994). “A particle method for history-dependent materials”. In: *Computer Methods in Applied Mechanics and Engineering* 118, pp. 179–196. DOI: [10.1016/0045-7825\(94\)00033-6](https://doi.org/10.1016/0045-7825(94)00033-6).
- Sun, D. and J. Sun (2005). “Strong Semismoothness of the Fischer-Burmeister SDC and SOC Complementarity Functions”. In: *Math. Program.* 103.3, pp. 575–581. ISSN: 1436-4646. DOI: [10.1007/s10107-005-0577-4](https://doi.org/10.1007/s10107-005-0577-4).
- Tasora, A. (2013). “Efficient simulation of contacts, friction and constraints using a modified spectral projected gradient method”. In: *Poster Proceedings of WSCG 2013*, pp. 69–72.
- Terzaghi, K. (1936). “The Shearing Resistance of Saturated Soils”. In: *Proc. 1st International Conference on Soil Mechanics and Foundation Engineering*.
- Teschner, M., B. Heidelberger, M. Müller, D. Pomeranerts, and M. Gross (2003). “Optimized Spatial Hashing for Collision Detection of Deformable Objects”. In: *Vision, Modeling, Visualization (VMV 2003)*, pp. 47–54.
- Tonge, R., F. Benevolenski, and A. Voroshilov (2012). “Mass splitting for jitter-free parallel rigid body simulation”. In: *ACM Transactions on Graphics* 31.4, pp. 1–8. ISSN: 0730-0301. DOI: [10.1145/2185520.2185601](https://doi.org/10.1145/2185520.2185601).
- Topin, V., F. Dubois, Y. Monerie, F. Perales, and A. Wachs (2011). “Micro-rheology of dense particulate flows: Application to immersed avalanches”. In: *Journal of Non-Newtonian Fluid Mechanics* 166.1-2, pp. 63–72. ISSN: 0377-0257. DOI: [10.1016/j.jnnfm.2010.10.006](https://doi.org/10.1016/j.jnnfm.2010.10.006).
- Topin, V., Y. Monerie, F. Perales, and F. Radjaï (2012). “Collapse Dynamics and Runout of Dense Granular Materials in a Fluid”. In: *Physical Review Letters* 109.18. ISSN: 1079-7114. DOI: [10.1103/physrevlett.109.188001](https://doi.org/10.1103/physrevlett.109.188001).
- Treskatis, T., M. A. Moyers-González, and C. J. Price (2016). “An accelerated dual proximal gradient method for applications in viscoplasticity”. In: *Journal of Non-Newtonian Fluid Mechanics*. ISSN: 0377-0257. DOI: [10.1016/j.jnnfm.2016.09.004](https://doi.org/10.1016/j.jnnfm.2016.09.004).
- Trinkle, J. C., J.-S. Pang, S. Sudarsky, and G. Lo (1997). “On Dynamic Multi-Rigid-Body Contact Problems with Coulomb Friction”. In: *Z. angew. Math. Mech.* 77.4, pp. 267–279. ISSN: 1521-4001. DOI: [10.1002/zamm.19970770411](https://doi.org/10.1002/zamm.19970770411).
- Ungarish, M. (2009). *An introduction to gravity currents and intrusions*. CRC Press.
- Vermeer, P. A. (1998). “Non-Associated Plasticity for Soils, Concrete and Rock”. In: *Physics of Dry Granular Media*, pp. 163–196. DOI: [10.1007/978-94-017-2653-5\\_10](https://doi.org/10.1007/978-94-017-2653-5_10).
- Ward, K., F. Bertails, T.-Y. Kim, S. Marschner, et al. (2007). “A Survey on Hair Modeling: Styling, Simulation, and Rendering”. In: *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 13.2, pp. 213–34. URL: <http://www-evasion.imag.fr/Publications/2007/WBKMcL07>.
- Wieckowski, Z. and M. Pawlak (2015). “Material point method in three-dimensional problems of granular flow”. In: *6th European Conference on Computational Fluid Dynamics*.

- Wieckowski, Z., S.-K. Youn, and J.-H. Yeon (1999). “A particle-in-cell solution to the silo discharging problem”. In: *International Journal for Numerical Methods in Engineering* 45.9, pp. 1203–1225. ISSN: 1097-0207. DOI: [10.1002/\(sici\)1097-0207\(19990730\)45:9<1203::aid-nme626>3.0.co;2-c](https://doi.org/10.1002/(sici)1097-0207(19990730)45:9<1203::aid-nme626>3.0.co;2-c).
- Yue, Y., B. Smith, C. Batty, C. Zheng, and E. Grinspun (2015). “Continuum Foam”. In: *ACM Transactions on Graphics* 34.5, pp. 1–20. ISSN: 0730-0301. DOI: [10.1145/2751541](https://doi.org/10.1145/2751541).
- Zhu, H., X. Liu, Y. Liu, and E. Wu (2006). “Simulation of miscible binary mixtures based on lattice Boltzmann method”. In: *Computer Animation and Virtual Worlds* 17.3-4, pp. 403–410. ISSN: 1546-427X. DOI: [10.1002/cav.143](https://doi.org/10.1002/cav.143).
- Zhu, Y. and R. Bridson (2005). “Animating sand as a fluid”. In: *ACM Transactions on Graphics* 24.3, pp. 965–972. DOI: [10.1145/1073204.1073298](https://doi.org/10.1145/1073204.1073298).





### Abstract

This dissertation focuses on the numerical simulation of mechanical systems that consist of a large number of discrete pieces interacting with each other through contacts and dry friction. Examples of such systems — for instance, sand or human hair — are common in natural environments; being able to predict their dynamics is therefore of importance for diverse applications, ranging from geotechnical considerations to engaging visual effects for feature films.

A major difficulty that complicates the numerical simulation of such complex systems stems from the nonsmoothness of their dynamics. Indeed, from a macroscopic viewpoint, the velocity of the individual constituents may exhibit instantaneous jumps, for instance when impacts occur. The first part of this manuscript will thus be dedicated to the establishment of efficient algorithms for the numerical simulation of discrete mechanical systems subject to contacts and Coulomb friction. We advocate using a Gauss–Seidel algorithm with an hybrid local solver, and show that this strategy performs robustly in the challenging context of virtual hair simulation.

The second part of this dissertation will focus on much bigger systems, consisting of millions or billions of grains. As computing every force between pairs of contacting grains quickly becomes intractable, we embrace a continuum viewpoint instead. We show how the numerical methods devised for the simulation of discrete mechanical systems can be adapted to the simulation of flows governed by the Drucker–Prager rheology — a constitutive relationship between the material’s stress and strain rate that macroscopically models the action of frictional contact forces. This approach allows us to capture experimentally-observed qualitative features of granular flows. Finally, we propose a new numerical model for the coupled simulation of a granular continuum with a surrounding fluid, and show that once again, we are able to leverage algorithms from discrete contact mechanics.

### Résumé

Cette thèse traite de la simulation numérique de systèmes composés de nombreux objets distincts, et dont le principal mécanisme d’interaction consiste en des contacts inélastiques avec frottement sec. On trouve de nombreuses occurrences de tels systèmes dans la nature, par exemple sous la forme de sable ou d’une chevelure humaine ; aussi la reproduction numérique de leur dynamique trouve des applications diverses, allant de considérations géotechniques à la production d’effets spéciaux réalistes.

Une difficulté majeure pour la simulation de tels systèmes concerne la *non-régularité* de leur dynamique ; à une échelle de temps macroscopique, on observe des sauts dans les vitesses des constituants, par exemple lors d’impacts. La première partie de ce manuscrit est ainsi dédiée à la conception d’algorithmes efficaces permettant de prendre en compte les contacts avec frottement de Coulomb lors de la simulation de systèmes mécaniques discrets. La méthode que nous proposons, basée sur un algorithme de type Gauss–Seidel avec stratégie hybride, s’avère robuste et performante sur le problème délicat de la simulation virtuelle de chevelures.

La seconde partie de ce manuscrit est consacrée à l’étude de systèmes à une échelle beaucoup plus grande, au delà du million de grains. Puisque le calcul de toutes les forces de contacts pour chaque paire de grains s’avérerait trop coûteux, on adopte un point de vue macroscopique. On propose ainsi d’adapter les méthodes développées pour la simulation de systèmes discrets à la résolution de la rhéologie dite de Drucker–Prager, une relation entre la contrainte et le cisaillement du matériau exprimant l’influence moyennée des forces de frottement. On montre que cette approche nous permet de retrouver le comportement qualitatif de matériaux granulaires secs observé expérimentalement. Finalement, nous proposons un nouveau modèle numérique pour l’étude des dynamiques couplées d’un matériau granulaire immergé dans un fluide Newtonien, et montrons une nouvelle fois que les algorithmes développés pour la mécanique discrète s’avèrent également pertinents dans le cas continu.