

## Motivation, Applications

# Géométrie d'images multiples

**Bill Triggs**

MOVI

INRIA Rhône-Alpes

## Plan d'exposé

- 1 Introduction
- 2 L'approche tensorielle à la vision projective
- 3 La reconstruction projective
- 4 L'auto-calibrage d'une caméra en mouvement
- 5 Bilan des contributions

## Extraire de l'information 3D des images multiples

- 1 Mesure 3D pour l'industrie, la santé
- 2 Reconstruction des modèles 3D d'objets et de scènes
- 3 Réalité virtuelle et augmenté, synthèse de nouvelles vues

## Notre orientation

- Une orientation théorique de base, « géométrie pure et dur »
- Formules et algorithmes, pas d'élaboration de système, pas de jolies images

## Challenges théoriques

- Méthodes de reconstruction plus flexibles — réduire l'information / calibrage préalable, extraire le maximum des données.
- Mieux comprendre la structure géométrique des images multiples

## Représentation euclidienne homogène

$\mathbf{x} = \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}$	point 3D
$\mathbf{d} = \mathbf{x} - \mathbf{y} = \begin{pmatrix} d_x \\ d_y \\ 0 \end{pmatrix}$	vecteur de déplacement/direction 3D
$\mathbf{p} = \begin{pmatrix} p_x & p_y & d \end{pmatrix}$	plan 3D: $\mathbf{p} \cdot \mathbf{x} = 0$ si $\mathbf{x}$ est sur $\mathbf{p}$
$\mathbf{T} = \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix}$	déplacement rigide
$\mathbf{x} = \begin{pmatrix} u \\ v \\ 1 \end{pmatrix}$	point 2D / image

- Toute transformation euclidienne (ou projective ...) est représentée par une simple multiplication matricielle.
- Les lois de transformation de vecteurs « point » et « direction » sont unifiées.

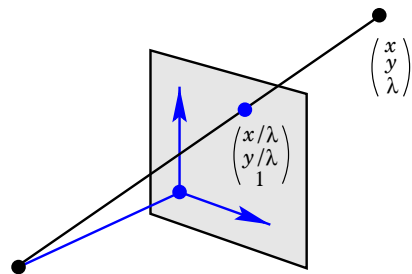
## Modèle de caméra sténopée

- Une caméra perspective est paramétrée par sa matrice de projection  $3 \times 4$  :

$P = K \begin{pmatrix} R & t \end{pmatrix}$	matrice de projection
$t, R$	position et orientation 3D de la caméra
$K = \begin{pmatrix} 1 & s & u_0 \\ 0 & a & v_0 \\ 0 & 0 & 1/f \end{pmatrix}$	matrice de calibrage interne
$f, (u_0, v_0)$	focale, point principal
$a, s$	rapport des échelles, skew ( $s \approx 0$ )

- Toute matrice  $3 \times 4$  peut être décomposée en ce forme (à une facteur d'échelle près)
- C'est un modèle maniable, mais un peu simplifié — la distorsion optique n'est pas modélisée.

## Projection perspective



- Un point 3D est projeté dans une image par multiplication par la matrice de projection, puis renormalisation d'échelle :

$$\mathbf{x} = \begin{pmatrix} x/\lambda \\ y/\lambda \\ 1 \end{pmatrix} \quad \text{où} \quad \lambda \mathbf{x} = \begin{pmatrix} x \\ y \\ \lambda \end{pmatrix} = P \mathbf{x}$$

- Tout point le long du même rayon optique projette au même point image.
- La **profondeur**  $\lambda$  — mesure de la distance le long du rayon — est perdue pendant la projection.
- La retrouver correspond à la reconstruction du point 3D dans le repère caméra ...

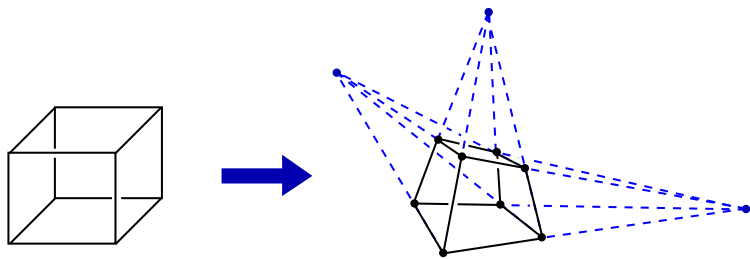
## Structure 3D projective

- Une **transformation 3D projective** de la scène est représenté par une matrice  $4 \times 4$  non-singulière  $H$

$$\mathbf{x} \longrightarrow \mathbf{x}' \simeq H \mathbf{x} \quad P \longrightarrow P' \simeq P H^{-1}$$

- “ $\simeq$ ” dénote égalité à une échelle près, de tels vecteurs homogènes sont vus comme équivalents.

- Les transformations projectives préservent les structures linéaires et d'incidence.
- Les distances, les angles, et la distinction entre vecteurs de déplacement et points finis ne sont *pas* préservés



## L'approche tensorielle à la vision projective

### Pourquoi la reconstruction projective ?

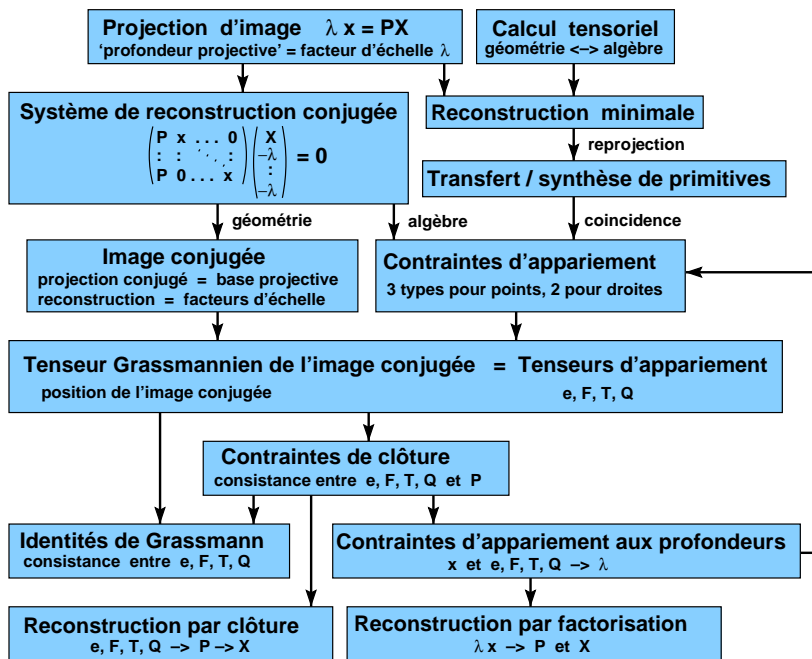
- Une transformation projective 3D  $H$  ne change pas les images :

$$P x = (P H^{-1}) (H x)$$

- Donc une reconstruction visuelle à partir d'images perspectives non-calibrées est toujours ambiguë jusqu'à (au moins) une projectivité 3D près.
- On décompose le problème de reconstruction 3D en deux étapes :
  - reconstruction projective, à partir d'images non-calibrées
  - rectification euclidienne, par moyen d'informations supplémentaires
- L'étape projective donne déjà toute la structure 3D de la scène, à 9 paramètres euclidiennes près.

### Motivation

- La géométrie 3D des caméras laisse sa trace dans les images
- Retrouver cette trace aide à
  - la mise en correspondance des primitives
  - la reconstruction projective
  - la synthèse de nouvelles vues
- Au coeur du problème: les **tenseurs d'appariement**
  - objets algébro-géométriques multi-indices et multi-images
  - ils caractérisent la géométrie 3D relative des caméras
  - ils peuvent être estimés à partir des correspondances images



## L'image conjuguée

L'ensemble d'images d'une primitive 3D caractérisent la primitive, et donnent une représentation implicite d'elle.

- Rassembler dans un grand système  $3m \times 4$ , les équations de projection  $\lambda_i x_i = P_i x$  d'un point 3D dans  $m$  images:

$$X = P x \quad \text{où} \quad X \equiv \begin{pmatrix} \lambda_1 x_1 \\ \vdots \\ \lambda_m x_m \end{pmatrix} \quad \text{et} \quad P \equiv \begin{pmatrix} P_1 \\ \vdots \\ P_m \end{pmatrix}$$

- Un point 3D  $x$  est représenté par son **image conjuguée**  $X$  — l'ensemble  $3m$ -D de toutes ses coordonnées images homogènes.
- La géométrie 3D des caméras est représentée par leur **matrice de projection conjuguée**  $P$ .

## Dérivation des contraintes d'appariement

- Les equations de projection conjuguées sont linéaires en  $x$  et  $\lambda$  :

$$\begin{pmatrix} P_1 & x_1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ P_m & 0 & \dots & x_m \end{pmatrix} \begin{pmatrix} x \\ -\lambda_1 \\ \vdots \\ -\lambda_m \end{pmatrix} = 0$$

- Il y a une solution  $(x, \lambda)$ , donc toute mineur (déterminante de sous-matrice maximale) doit être nulle. Développe et simplifie.
- On trouve des contraintes multi-linéaires en les points de 2,3,4 images, dont les coefficients sont des déterminants  $4 \times 4$  de 4 lignes de  $P$ .
- Ces coefficients peuvent être rassemblées en **tenseurs** (tableaux multi-indices) inter-images.

## Tenseurs d'appariement

- Tout déterminant  $4 \times 4$  de 4 lignes de  $P$  donne une entrée d'un tenseur d'appariement.
- Il y a 4 types de tenseur, qui correspondent aux 4 façons de choisir les lignes en 2-4 images: 1+3, 2+2, 1+1+2, 1+1+1+1

$e_{12}$	épipôle de caméra 1 en image 2
$F_{12}$	matrice fondamentale
$T_1^{23}$	tenseur trifocale
$Q^{1234}$	tenseur quadrifocale

- Les tenseurs dépendent seulement de la *géométrie intrinsèque* 3D des caméras — les déterminants élimine la dépendance sur le repère 3D.

## Contraintes d'appariement

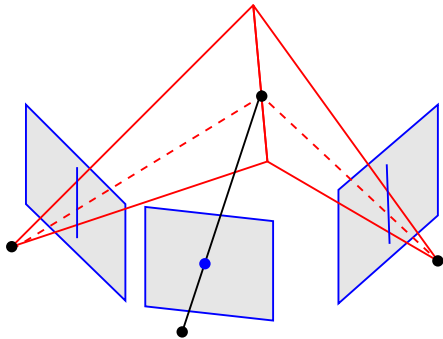
- Les différentes images d'une primitive 3D doivent vérifier des contraintes de consistance entre eux.
- Pour les points 3D, il y a 3 types de ces **contraintes d'appariement**, *p.ex.*

$$\begin{aligned} \mathbf{x}_1 \mathbf{F}_{12} \mathbf{x}_2 &= 0 && \text{épipolaire} \\ [\mathbf{x}_2]_{\times} (\mathbf{T}_1^{23} \cdot \mathbf{x}_1) [\mathbf{x}_3]_{\times} &= \mathbf{0} && \text{trifocale} \end{aligned}$$

- Il y a aussi des contraintes analogues sur les droites 3D
- **Applications**: estimation de tenseurs, aide à la mise en correspondance et « transfert » des primitives entre images.

## Interprétation géométrique

- Les faisceaux de rayons optiques 3D reprojétés des primitives images doivent s'intersecter en une primitive cohérente 3D.
- Les différentes contraintes d'appariement correspondent aux différentes façons de formuler ceci algébriquement.



## Contraintes de clôture

- Contraintes de consistance entre les matrices de projection, et les tenseurs d'appariement dérivés d'eux.
- Il y a 5 types, dont les deux les plus utiles:

$$\begin{aligned} \mathbf{F}_{ji} \mathbf{P}_i + [\mathbf{e}_{ij}]_{\times} \mathbf{P}_j &= \mathbf{0} && \text{clôture } \mathbf{F}\text{-}\mathbf{e} \\ \mathbf{T}_{B_j}^{A_i C_k} \mathbf{P}_a^{B_j} + \mathbf{e}_j^{A_i} \mathbf{P}_a^{C_k} - \mathbf{P}_a^{A_i} \mathbf{e}_j^{C_k} &= \mathbf{0} && \text{clôture } \mathbf{T}\text{-}\mathbf{e}^2 \end{aligned}$$

- Facteurs d'échelle consistantes pour  $\mathbf{e}$ ,  $\mathbf{F}$ ,  $\mathbf{T}$  sont nécessaires.
- **Dérivation**: équivalent aux contraintes Grassmann qui disent que la matrice de projection conjuguée est une base de l'image conjuguée
- **Applications**: la reconstruction par clôture, les contraintes suivantes ...

## Contraintes d'appariement aux profondeurs

- Contraintes qui lient les primitives images *avec leurs profondeurs projectives correctes* aux tenseurs d'appariement.
- Il y a 5 types, parallèles aux contraintes de clôture, *p.ex.* :

$$\begin{aligned} \mathbf{F}_{ji} (\lambda_{ip} \mathbf{x}_{ip}) + \mathbf{e}_{ij} \wedge (\lambda_{jp} \mathbf{x}_{jp}) &= \mathbf{0} && \text{épipolaire} \\ \mathbf{T}_{B_j}^{A_i C_k} (\lambda_j \mathbf{x}^{B_j}) - (\lambda_i \mathbf{x}^{A_i}) \mathbf{e}_j^{C_k} + \mathbf{e}_j^{A_i} (\lambda_k \mathbf{x}^{C_k}) &= \mathbf{0} && \text{trifocal} \end{aligned}$$

- Les contraintes d'appariement standards suivent par multiplication par  $\mathbf{x}$ ,  $[\mathbf{x}]_{\times}$
- **Application**: retrouver les profondeurs pendant la reconstruction.
- **Dérivation**: projeter une primitive 3D avec les matrices de projection des contraintes de clôture.

## Identités de Grassmann

- Contraintes de consistance quadratiques entre les tenseurs d'appariement, plus leurs conséquences cubiques, quartiques, ...
- 2 en 2 images, 15 en 3, ...
- Les plus simples sont:

$$\begin{aligned} \mathbf{F}_{12} \mathbf{e}_{12} &= \mathbf{0} \\ \text{Cofactor}(\mathbf{F}_{12}) + 2 \mathbf{e}_{21} \mathbf{e}_{12}^\top &= \mathbf{0} & \det(\mathbf{F}_{12}) &= 0 \\ \mathbf{F}_{32} \mathbf{e}_{12} - \mathbf{e}_{13} \wedge \mathbf{e}_{23} &= \mathbf{0} & \det(\mathbf{T}_1^{23} \cdot \mathbf{x}_1) &= 0 \quad \forall \mathbf{x}_1 \\ \mathbf{T}_1^{23} \cdot \mathbf{e}_{21} + \mathbf{e}_{12} \mathbf{e}_{23}^\top &= \mathbf{0} \end{aligned}$$

- **Dérivation**: contraintes Grassmann-Plücker de l'image conjuguée.
- Même en 3 images elles sont trop nombreuses et trop complexes à manier facilement.

## Autres travaux sur l'approche tensorielle

- Estimation optimale des tenseurs à partir des primitives images, sujet à leurs contraintes de consistance.
- Tenseurs et contraintes différentiels, pour le cas où les images sont proches les unes aux autres.

## Perspectives sur l'approche tensorielle

- L'image conjuguée et l'approche systématique tensorielle ont beaucoup apporté.
- Les facteurs d'échelle sont souvent plus significatives qu'on ne le pense.
- Pour points et droites, la théorie de base semble plus ou moins complète.
- Depuis 3–4 images une représentation « grande ensemble de tenseurs » devient trop lourde — il vaut mieux expliciter les matrices de projection.
- Problème ouvert: les contraintes d'appariement entre les images d'une quadrique 3D, en  $\geq 3$  images (application à l'auto-calibrage)

## Reconstruction Projective

## Reconstruction projective par factorisation

- Regrouper les équations de projection conjuguées de  $n$  points 3D en  $m$  images:

$$\begin{pmatrix} \lambda_{11} \mathbf{x}_{11} & \cdots & \lambda_{1n} \mathbf{x}_{1n} \\ \vdots & \ddots & \vdots \\ \lambda_{m1} \mathbf{x}_{m1} & \cdots & \lambda_{mn} \mathbf{x}_{mn} \end{pmatrix}_{3m \times n} = \begin{pmatrix} \mathbf{P}_1 \\ \vdots \\ \mathbf{P}_m \end{pmatrix}_{3m \times 4} \begin{pmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n \end{pmatrix}_{4 \times n}$$

- La matrice  $3m \times n$  des  $(\lambda_{ip} \mathbf{x}_{ip})$  a pour rang seulement 4.

## Factorisation

- On peut décomposer une telle matrice dans deux facteurs de la même forme par, *p.ex.* SVD (Décomposition par Valeurs Singulières)
- Toute factorisation de ce type est une *reconstruction projective de la scène et des caméras* — l'ambiguïté est celle d'un changement de repère 3D projectif.
- La factorisation concrétise une reconstruction qui était déjà implicite dans les points renormalisés  $(\lambda \mathbf{x})$ .
- En plus, elle moyenne de façon efficace le bruit dans les points images mesurés.

## Retrouver les profondeurs projectives

- Les profondeurs projectives requises sont estimées depuis les contraintes d'appariement aux profondeurs, *p.ex.*

$$\mathbf{F}_{12}(\lambda_2 \mathbf{x}_2) = \mathbf{e}_{21} \wedge (\lambda_1 \mathbf{x}_1) \implies \lambda_2 \approx \lambda_1 \frac{(\mathbf{e}_{21} \wedge \mathbf{x}_1) \cdot (\mathbf{F}_{12} \mathbf{x}_2)}{\|\mathbf{e}_{21} \wedge \mathbf{x}_1\|^2}$$

- Les tenseurs  $\mathbf{F}_{ij}$  et  $\mathbf{e}_{ij}$  sont extraites des correspondances images.
- On propage les  $\lambda$  depuis une première image, dans un réseau d'images liées par les  $\mathbf{F}$ ,  $\mathbf{e}$ .

## L'algorithme

- 1 Extraire et mettre en correspondance les points dans toutes les images
  - 2 Normaliser les coordonnées images (pour augmenter la stabilité)
  - 3 Estimer un réseau de matrices fondamentales qui relie tous les images.
  - 4 Estimer les profondeurs projectives par les équations de profondeur.
  - 5 Pour stabilité, équilibrer la matrice des points renormalisés  $(\lambda \mathbf{x})$
  - 6 Factoriser la matrice par SVD (ou autres méthodes ...)
  - 7 Extraire la structure et les matrices de projection 3D
- On peut aussi reconstruire les droites 3D par moyen d'une paramétrisation « 2 points ».

## Reconstruction par « clôture »

- La factorisation exige que tous les points soient visibles dans toutes les images — ce qui est souvent difficile en pratique.
- Avec les **contraintes de clôture**, on peut estimer les matrices de projection directement des tenseurs d'appariement
- Par exemple, la contrainte de clôture  $F-e$  donne deux contraintes linéaires entre chaque paire d'images :

$$F_{ij} P_j + [e_{ji}]_{\times} P_i = 0$$

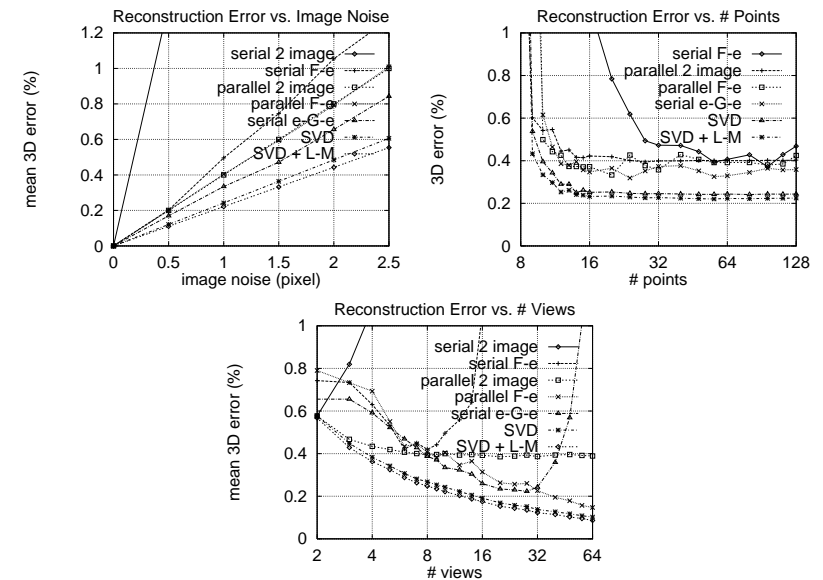
## Algorithme

- 1 Estimer un réseau de matrices fondamentales qui lie les images entre eux (chaque image lié à au moins deux autres).
- 2 Empiler les contraintes sur la matrice de projection conjuguée (les  $P_i$ ):

$$\begin{pmatrix} \vdots & \vdots \\ \cdots & F_{ij}^T & \cdots & [e_{ij}]_{\times} & \cdots \\ \vdots & \vdots \end{pmatrix} \begin{pmatrix} \vdots \\ P_i \\ \vdots \\ P_j \\ \vdots \end{pmatrix} = 0$$

- 3 Estimer l'espace nul. Toute base linéaire de cet espace donne une reconstruction des  $P_i$ , à une projectivité 3D près.
- 4 Reconstruire la structure 3D à partir des  $P_i$  par résection linéaire.

## Résultats expérimentaux



## Bilan sur la reconstruction

- La méthode par factorisation est très stable et à recommander, sauf qu'elle ne tolère pas des données manquantes, ce qui est gênant en pratique
- La méthode de clôture est moins stable, mais toujours mieux que de utiliser deux images seules.



## La quadrique absolue duale

- Une paramétrisation simple et projective de la structure euclidienne.
- Dans un repère euclidien, la **quadrique absolue duale** est la matrice

$$\Omega = \begin{pmatrix} I_{3 \times 3} & \mathbf{0} \\ \mathbf{0} & 0 \end{pmatrix}$$

- Elle est invariante par transformations euclidiennes:

$$\mathbf{T} \Omega \mathbf{T}^T = \Omega \quad \text{où} \quad \mathbf{T} = \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix}$$

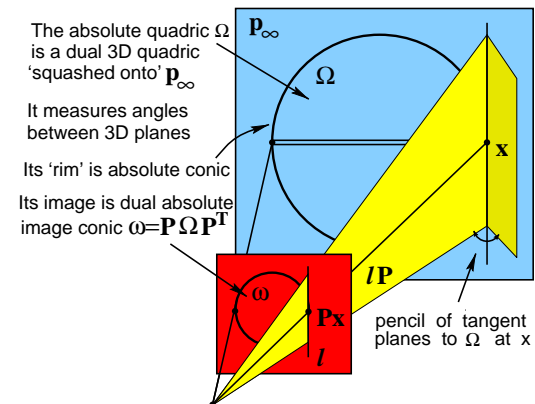
- Elle donne le produit scalaire entre les normaux des plans:

$$\begin{pmatrix} \mathbf{n}_1^T & d_1 \end{pmatrix} \Omega \begin{pmatrix} \mathbf{n}_2 \\ d_2 \end{pmatrix} = \mathbf{n}_1 \cdot \mathbf{n}_2$$

- Sa unique vecteur nulle est le plan à l'infini  $(\mathbf{0} \ 1)$ .

- Sous une transformation projective, la quadrique absolue duale devient une matrice  $4 \times 4$  rang 3 symétrique arbitraire.

- Interprétation géométrique: elle est une quadrique 3D complexe (pas de points réels), et qui a dégénéré sur une disque dans le plan à l'infini.



## Auto-calibrage d'une caméra en mouvement

### Structures euclidiennes en espace projectif

- Sous les déformations projectives 3D:
  - Les points infiniment distants (*vecteurs de direction*) peuvent devenir finis, donc ne peuvent pas être distingués des points finis.
  - Les directions orthogonales peuvent devenir non-orthogonales.

- Pour convertir une reconstruction projective dans une reconstruction euclidienne, il faut « localiser »:

1 Le **plan à l'infini** — le sphère de rayon infini qui contient tous les points de fuite.

2 L'orthogonale: une **base de directions orthogonales**, la **conique absolue** ou la **quadrique absolue duale**.

## Bases de directions orthogonales

- Une **base de directions orthogonales** est une matrice  $4 \times 3$ , dont les 3 colonnes donnent 3 directions orthogonales.
- Dans un repère euclidien,  $D = \begin{pmatrix} R \\ \theta \end{pmatrix}$ , où  $R$  est une rotation  $3 \times 3$ .
- Le lien avec la quadrique duale absolue:

$$\Omega = DD^T$$

## Contrainte de base de l'auto-calibrage

La projection calibrée d'une base de directions est orthogonale

$$K^{-1}PD \simeq \text{rotation } 3 \times 3$$

- $K$  = matrice  $3 \times 3$  de calibration interne de la caméra
- $P$  = matrice  $3 \times 4$  de projection
- $D$  = matrice  $4 \times 3$  de base des directions

- Ceci implique la **contrainte de projection de la quadrique duale absolue**

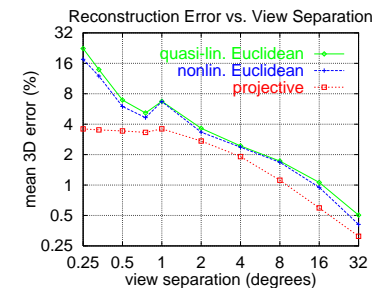
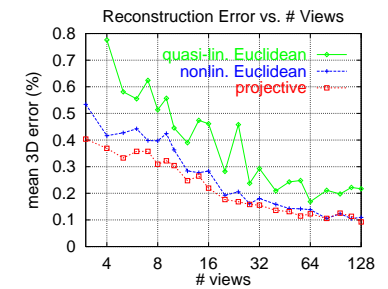
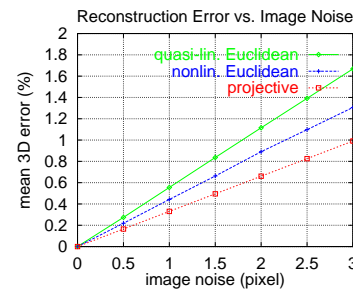
$$\omega \simeq P_i \Omega P_i \quad \text{où} \quad \omega \equiv KK^T$$

- Si la calibration interne  $K$  de la caméra est inconnue mais constante,  $\omega = KK^T$  est aussi constante.
- D'autres informations sur  $K$  peuvent aussi être intégrées.

## Algorithme d'auto-calibrage

- 1 Reconstruction projective pour retrouver les matrices de projection  $P_i$  dans un repère projectif inconnu.
- 2 Minimisation non-linéaire contrainte de l'erreur résiduelle de  $\omega \simeq P_i \Omega P_i$ 
  - Les variables sont  $\omega$  (calibrage de la caméra) et  $\Omega$  (la structure euclidienne de la scène).
  - Une initialisation aléatoire ou défaute suffit.
  - La contrainte est  $\det(\Omega) = 0$  (car  $\Omega$  est de rang 3).
  - La méthode numérique est la programmation quadratique séquentielle SQP.
- 3 Extraire la matrice de calibration  $K$  et la structure euclidienne de  $\omega$  et  $\Omega$ .
  - Au moins 3 images avec rotations générales sont nécessaires.
  - En générale la méthode donne de très bons résultats.

## Expériences



## Méthode plane

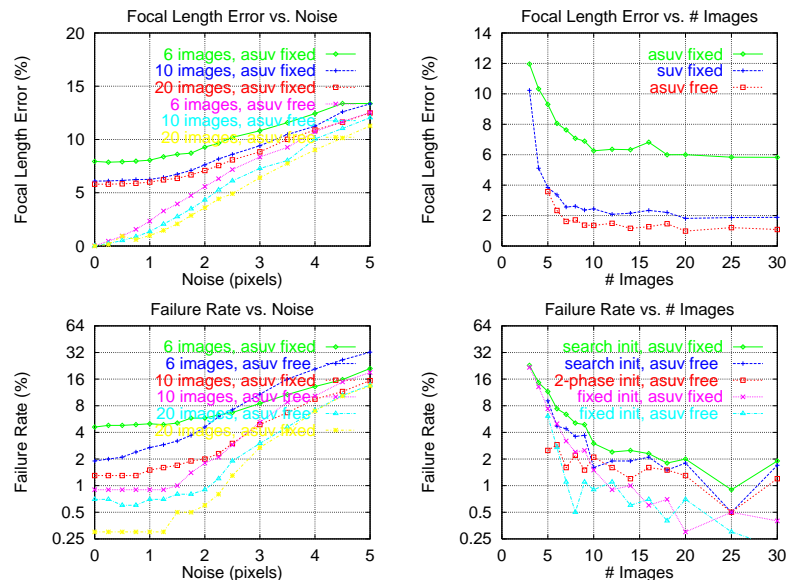
- Contrainte: les 2 directions orthogonales dans le plan 3D doivent être projetées aux directions orthogonales dans les images calibrées.
- Formulation algébrique:
  - Le plan 3D est représenté par une de ses images (disons 1).
  - Les directions orthogonales sont représentées par deux points inconnus  $x, y$  dans le plan.
  - Les autres images du plan sont représentées par les homographies  $H_{i1}$
  - La contrainte devient:

$$\begin{aligned} \|u_i\|^2 &= \|v_i\|^2 & \text{où} & \quad (u_i, v_i) \equiv K^{-1} H_{i1}(x, y) \\ u_i \cdot v_i &= 0 \end{aligned}$$

- On estime  $(x, y)$  et  $K$  à partir des  $H_{i1}$ , par optimisation non-linéaire.
- Au moins 5 images sont nécessaires pour une  $K$  générale.

## Bilan des contributions

## Expériences



## L'approche tensorielle à la vision projective

- L'image conjuguée, et le lien entre les coordonnées Grassmann-Plücker de l'image conjuguée et les tenseurs d'appariement.
- Les contraintes de clôture, d'appariement aux profondeurs, et de consistance inter-tenseur, et les applications qu'elles donnent.
- Estimation optimale des tenseurs sous contraintes de consistance.
- Contraintes d'appariement différentielles.

## Méthodes de reconstruction projective

- Méthode de factorisation, basée sur les contraintes d'appariement aux profondeurs (avec P. Sturm).
- Méthode de « clôture ».

## Auto-calibrage

- La quadrique duale absolue et le lien avec les repères de direction.
- Les méthodes quasi-linéaire et SQP, d'auto-calibrage à partir d'une reconstruction projective.
- La méthode d'auto-calibrage des scènes planes, à partir de homographies.