

INTERPRÉTATION CONTEXTUELLE ET ASSISTÉE DE FONDS D'ARCHIVES NUMÉRISÉS : APPLICATION À DES REGISTRES DE VENTES DU XVIII^e SIÈCLE

Joseph CHAZALON

Direction : Bertrand COÜASNON et Jean CAMILLERAPP

Collaborateurs : Laurent GUICHARD, Aurélie LEMAITRE, Alejandro TOSELLI



Équipe



Organismes de recherche



UNIVERSITÉ
EUROPÉENNE
DE BRETAGNE

PRES



Yvelines
Conseil général



Financement

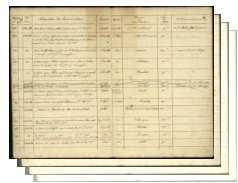
Contexte sociétal et scientifique

Dématérialisation de fonds d'archives

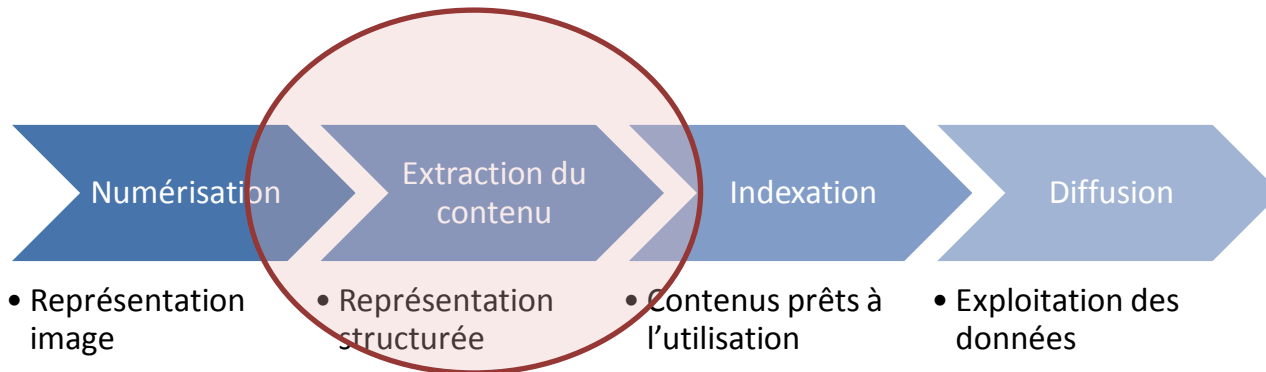
Avantages

- Protéger des contenus fragiles
- Diffuser des contenus rares
- Valoriser le patrimoine documentaire

Étapes



Documents
papier



Utilisateur

Exemple de données : séquestres révolutionnaires

Numéro
de vente

NUMÉROS des VENTES.	DATES des procès-verbaux des VENTES.	DÉSIGNATION DES OBJETS ALIÉNÉS, et de la Commune où ils sont situés.	INDICATION DE L'ANCIEN ÉTABLISSEMENT, ou de l'ancien Propriétaire.	NOM de l'Adjudicataire ou de son Command.	MONTANT de l'Adjudication.	SOMMES PAYÉES	SOMMES PRÉSENTES DUES en capital.	RÉDUCTION de ce qui reste dû au cours du jour de la vente, d'après le tableau de dépréciation du département.	INTÉRÊT de ce qui EN NUMÉRO calculé jusqu'au 2 ^e . 3
<i>Avril 1791.</i>									
188	28	98 Souches & défriches en 2 pièces, terrain de Savigny, ch ^{re} sur la fontaine	Cure de Savigny sur orge	Dubamel Couvreur à Savigny	1,400				
189	r	149 Souches & défriches à Vigney, en 4 pièces même terrain, ch ^{re} de Buisson, et de la route de la ville	Idem	Yart Citoyen de Paris g ^{re} de la rue Dumouriez aux Ornières	2,200				
190	r	Un demi arpent de même terrain, ch ^{re} du mail	Idem	Duval Boulangier à Savigny	1,150				
191	r	Un demi arpent de terrain, même terrain, ch ^{re} du Nord de la Chapelle	Idem	Olivier Vigneron à Savigny	700				
<i>Mai 1791.</i>									
192	2	4 arpents 75 Souches défriches, en 2 pièces, terrain de la commune, ch ^{re} de la Borde, et la piece au bout	Cure de la Quêre	Bucquet Cure de la Quêre	2,500				
193	r	187 Souches & défriches en une pièce, même terrain près le Pont rouge, et 100 perches de terre, ch ^{re} de la ville	Idem	Anchille Citoyen de la Borde aux lettres de la Quêre	3,050				
194	r	2 arpents de terrain, en 2 pièces, même terrain, ch ^{re} de la ville	Idem	Desvignes Cabanetier	1,550				

Ancien
propriétaire

Acquéreur

Exemple de besoin

Numéro	Quantité	Description	Propriétaire	Acquéreur	Prix
188	28	98 Souches 8 de Vignes en 3 pièces, terrain de Savigny, 1400	Cure de Savigny	Duhamel	1400
189	7	149 Souches 1/2 de Vignes en 4 pièces, même terrain et de même, 2200	Yart	Yart	2200

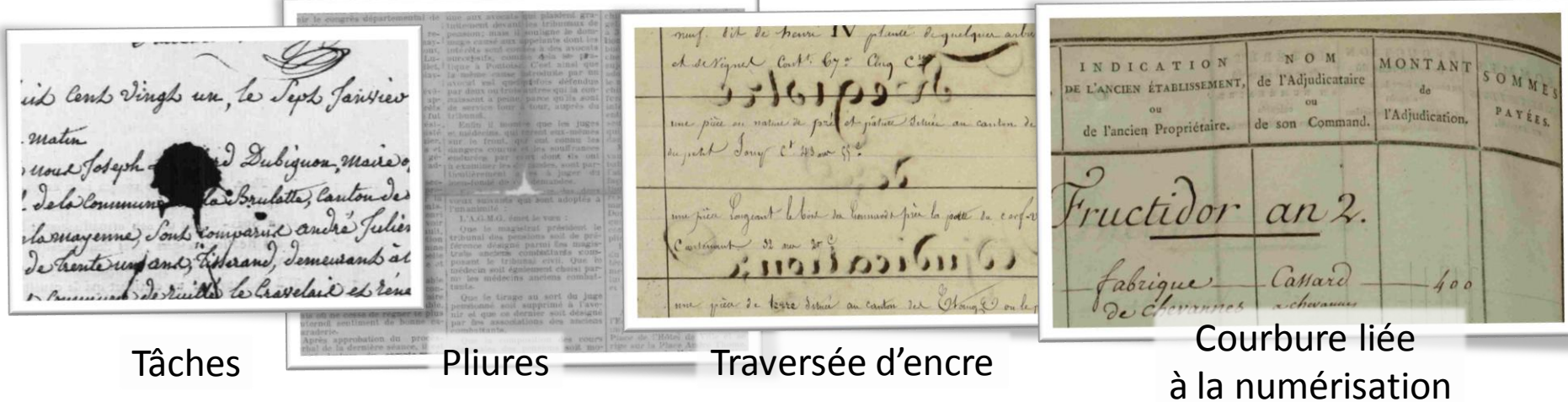
Extraction du contenu

- Localisation et reconnaissance
- Production de résultats prêts pour l'indexation
- Interprétation sémantique

Fichier image	Numéro vente	Ancien prop.	Acquéreur
00071.jpg	188	Cure de Savigny	Duhamel
00071.jpg	189	Cure de Savigny	Yart

Difficultés des fonds d'archives (1/2)

Dégradations : exemples

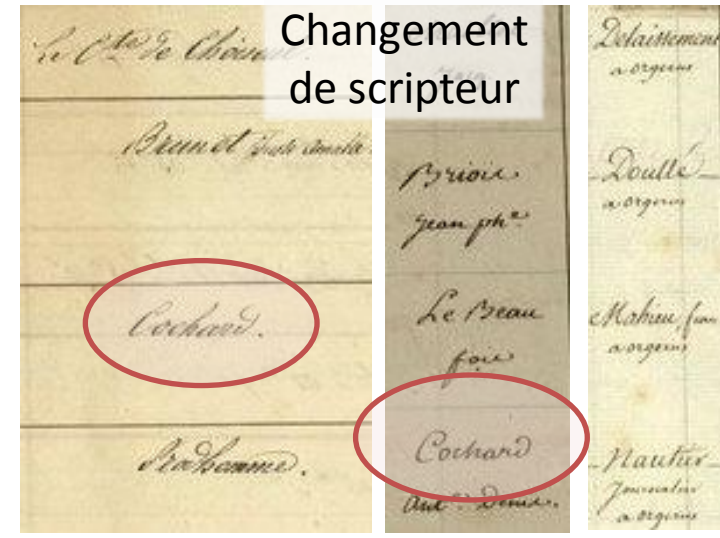
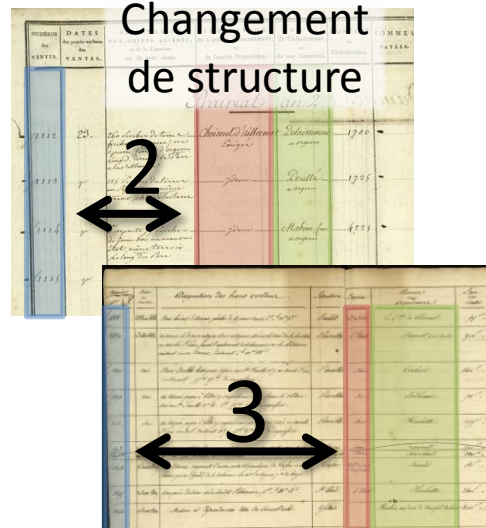
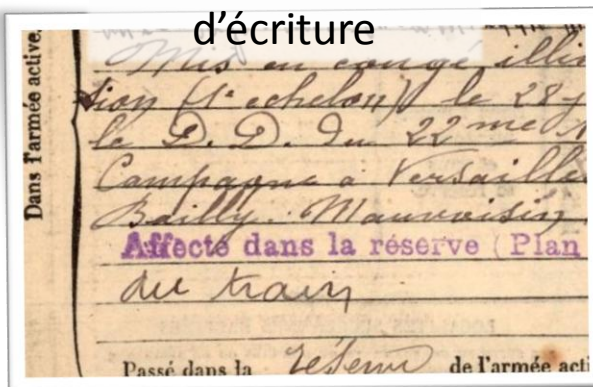


Variabilités : exemples

Mélange de styles d'écriture

Changement de structure

Changement de scripteur



Difficultés des fonds d'archives (2/2)

→ Conséquences

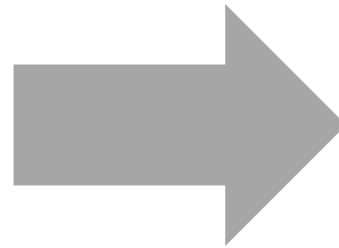
Dégradations

+

Variabilités

+

Cas particuliers

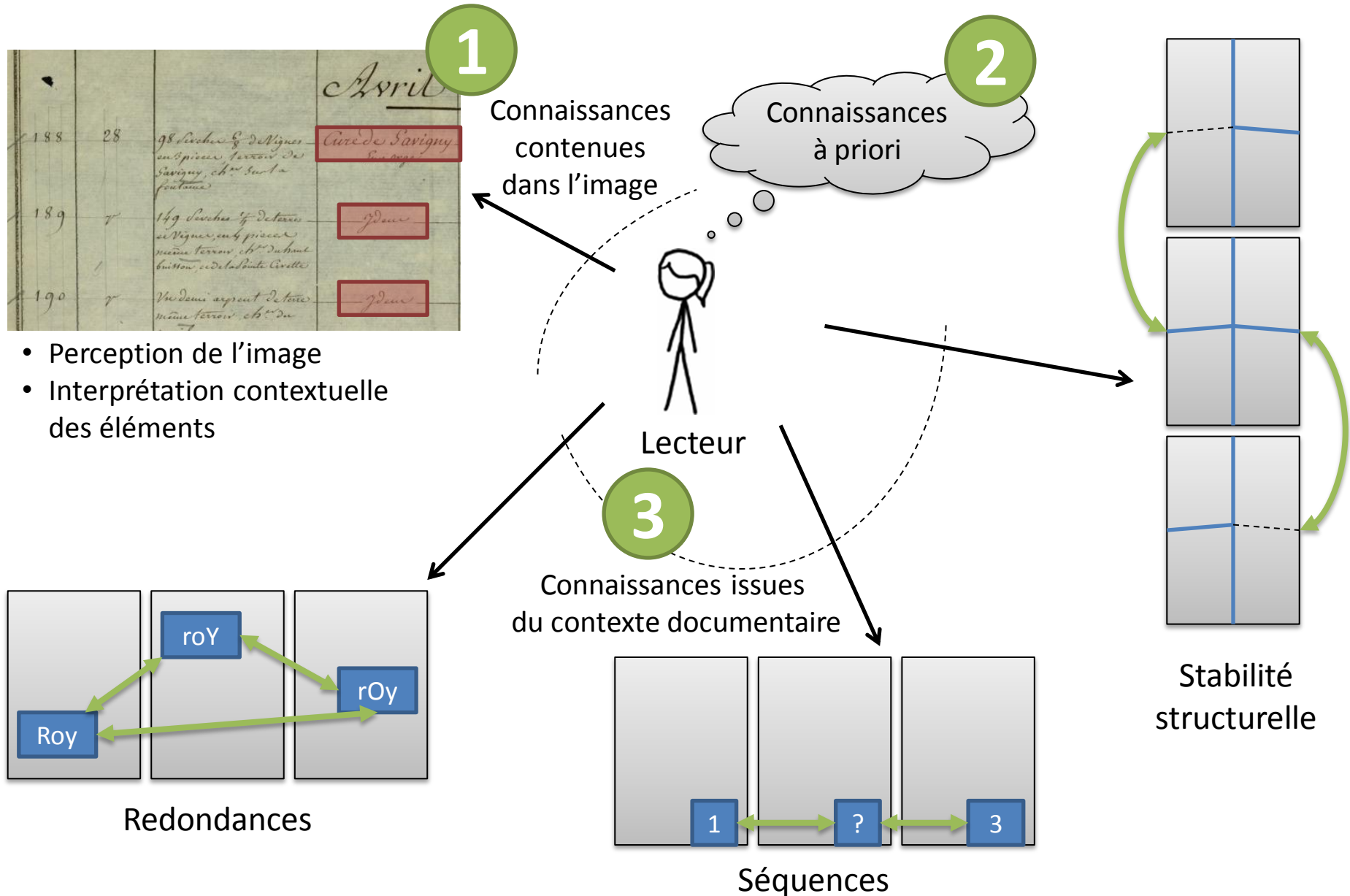


Ambiguïté

+

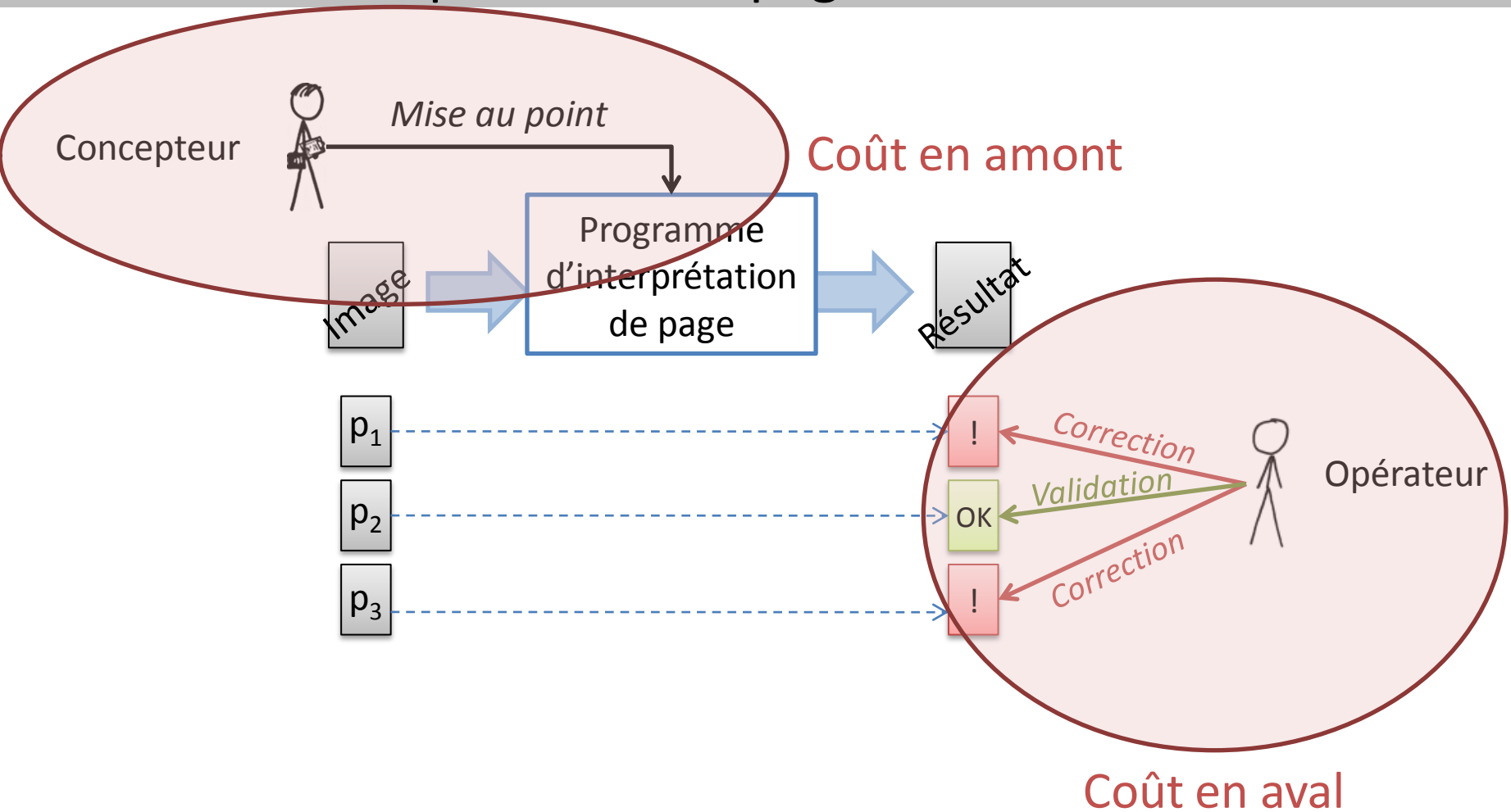
Impossibilité d'anticiper
tous les cas

Quelles connaissances les humains utilisent-ils ?



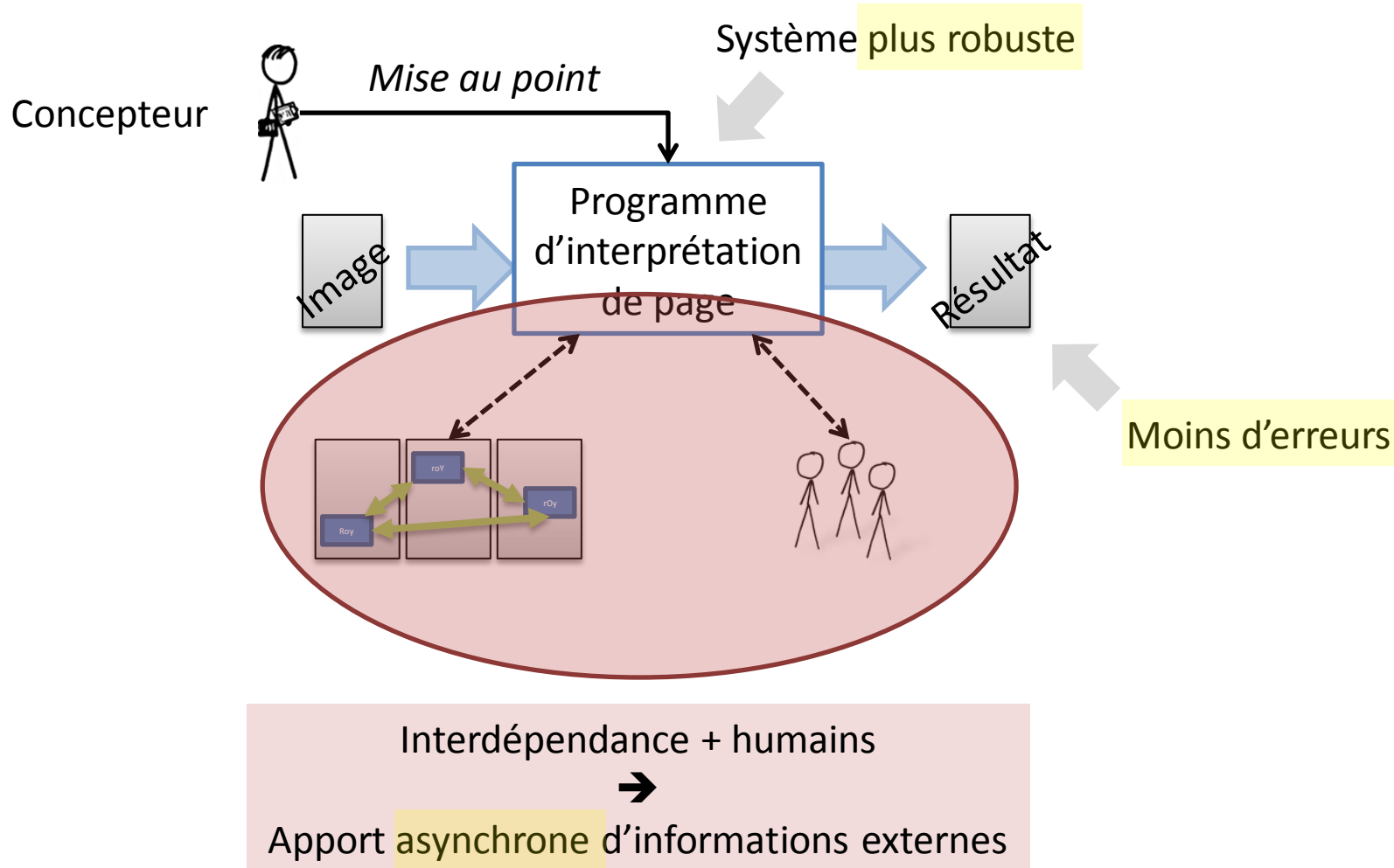
Interprétation automatique de fonds documentaires

Limites de l'interprétation de pages isolées



Notre objectif

Intégrer plus de connaissances lors de l'interprétation



Plan

1. Introduction

2. Résumé de l'état de l'art

- a. Formalisation des connaissances
- b. Exploitation du contexte documentaire
- c. Interaction avec des opérateurs humains
- d. Positionnement de notre approche

3. Contribution : interprétation itérative

4. Validation expérimentale et en production

5. Conclusion et perspectives

Formalisation des connaissances (1/2)

Approches algorithmiques

[Tang94, O'Gorman95, Clavier04, Shafait06]

```
...  
image.binariser(SEUIL);  
couche1 = image.extraireComposantes();  
couche1.appelerOCR(LEXIQUE);  
...
```

Concepteur → programme

☹ Trop figées

Approches statistiques

[LeBourgeois01, Montreuil09]

Exemples

Apprentissage

Modèle

Programme
d'interprétation

$P(\text{« zone texte »} \mid \text{« texte à gauche »}) = \dots$

Concepteur → entraînement

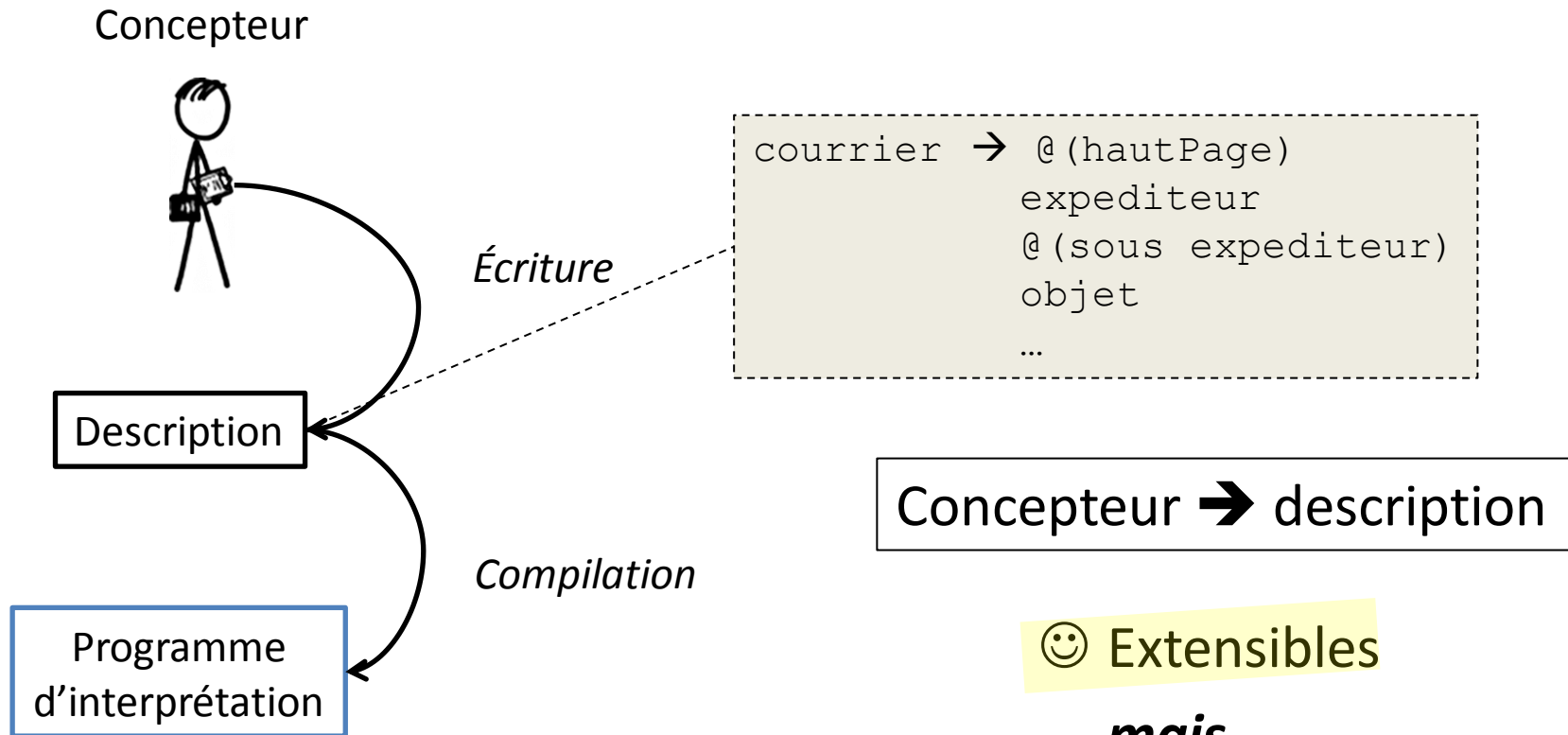
☹ Amorçage difficile

☹ Expressivité limitée

Formalisation des connaissances (2/2)

Approches déclaratives

[Vaxivere92, Tang95, Pasternak95, Couasnon96, Marriott98]



😊 Extensibles

mais

😞 Sensibles aux
dégradations
et à la variabilité

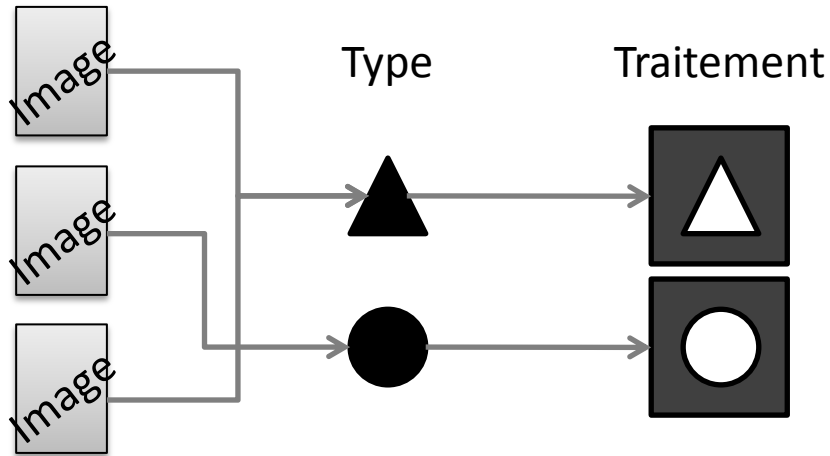
Plan

1. Introduction
2. Résumé de l'état de l'art
 - a. Formalisation des connaissances
 - b. Exploitation du contexte documentaire
 - c. Interaction avec des opérateurs humains
 - d. Positionnement de notre approche
3. Contribution : interprétation itérative
4. Validation expérimentale et en production
5. Conclusion et perspectives

Exploitation du contexte documentaire

Approches avec classification du type de document

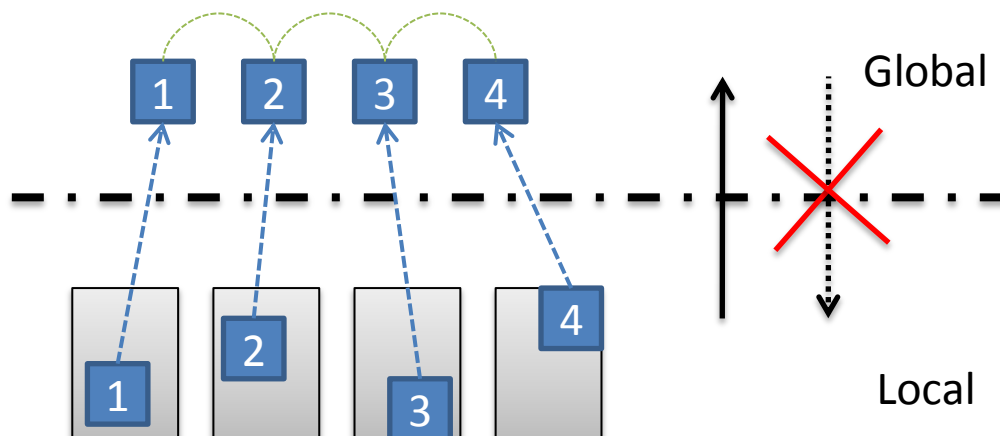
[Klein04, Lamiroy12]



☹ Pas d'exploitation globale des données

Approches à interprétation globale

[Lin97, Saund11, Xui12]



😊 Fiabilisation des résultats locaux
mais

☹ Pas de réinjection locale de l'information

Plan

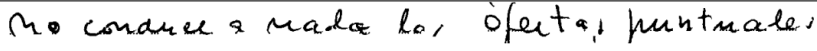
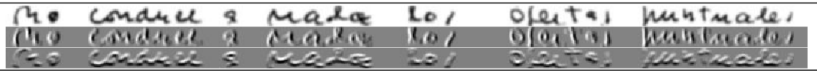
1. Introduction
2. Résumé de l'état de l'art
 - a. Formalisation des connaissances
 - b. Exploitation du contexte documentaire
 - c. Interaction avec des opérateurs humains
 - d. Positionnement de notre approche
3. Contribution : interprétation itérative
4. Validation expérimentale et en production
5. Conclusion et perspectives

Interaction avec des opérateurs humains

Traitements initiés par l'humain

[Bapst96, Ramel07, Vidal08, Llad08]

Humain → Détecte les erreurs

	Line Img	
		
		
INTER-0	(p)	()
	(\hat{s})	(no conduce o mala los afectos eventuales)
	(\hat{s}_p)	(no conduce)
INTER-1	(c)	(a)
	(p)	(no conduce a)
	(\hat{s})	(nada los afectos eventuales)
	(\hat{s}_p)	(nada)
INTER-2	(c)	(las)
	(p)	(no conduce a nada las)

😊 Résultat optimal en théorie

mais

😞 Interaction synchrone

😞 Détection manuelle des erreurs

Traitements initiés par la machine

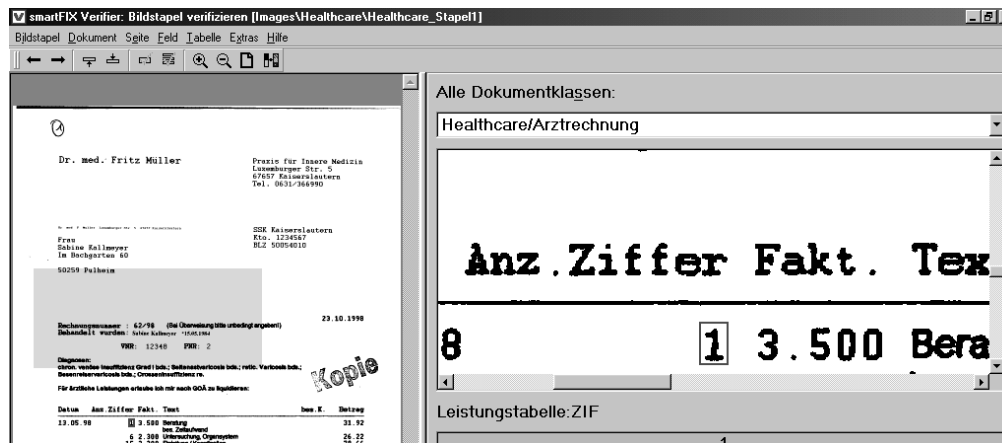
[Ogier98, Robadey01, Klein04]

Humain → Répond aux sollicitations

😊 Sollicitation asynchrone parcimonieuse

mais

😞 Détection d'erreur automatique faillible



Positionnement de notre approche

Formalisation des connaissances

- ➔ Extension des approches déclaratives

Exploitation du contexte documentaire

- ➔ Interprétation globale
- ➔ Ajout d'un mécanisme de réintégration (fusion)
global ➔ local

Interaction avec des opérateurs humains

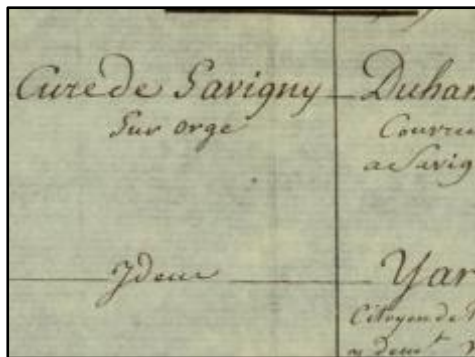
- ➔ Traitements initiés par la machine (questions)
mais possibilité d'accepter un guidage par l'humain
- ➔ Gestion automatique d'échanges asynchrones

Plan

1. Introduction
2. Résumé de l'état de l'art
3. Contribution : interprétation itérative
 - a. Interprétation itérative
 - b. Extension d'une approche déclarative
3 étapes
4. Validation expérimentale et en production
5. Conclusion et perspectives

Interprétation itérative (1/4)

Traitement page par page



Image

**Module
d'interprétation
de page**

Ancien prop.

Cure de Savigny

Julien

Résultat

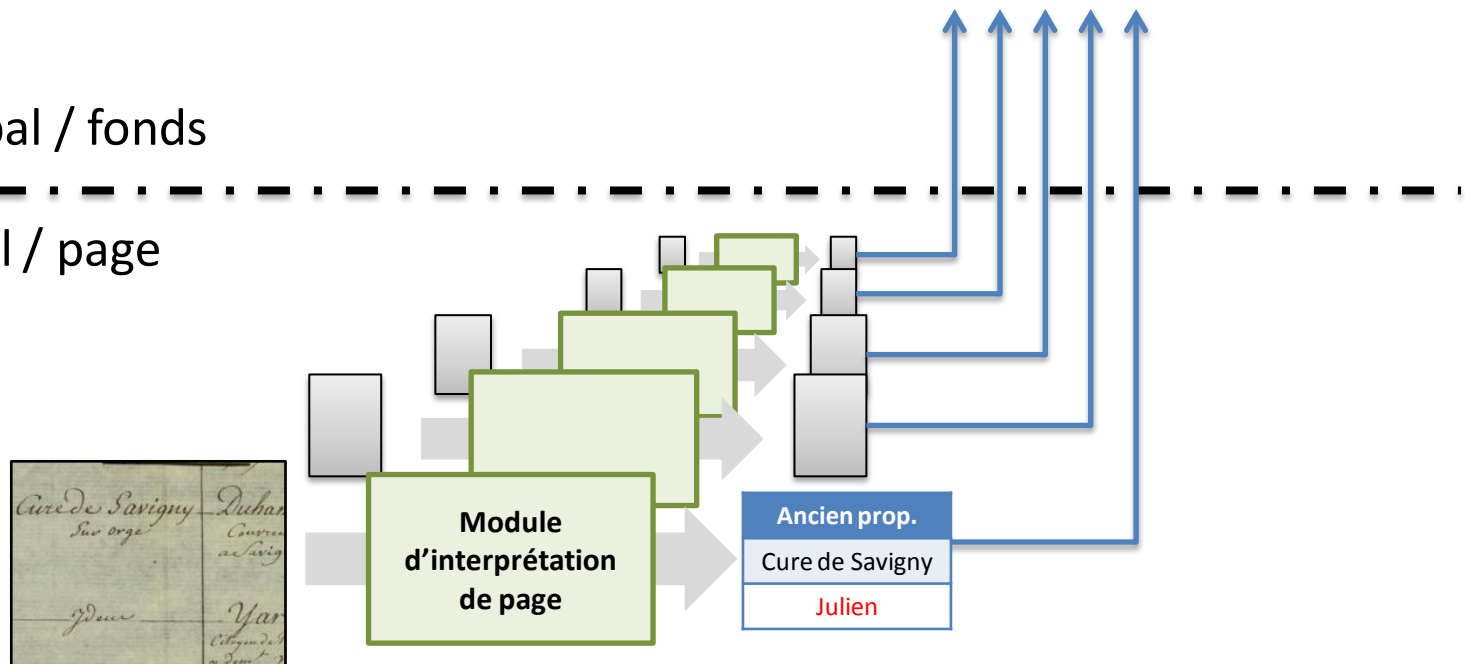
Interprétation itérative (2/4)

Collecte globale des résultats

1. Traitement
2. Déchargement programme
3. Collecte des résultats

Niveau global / fonds

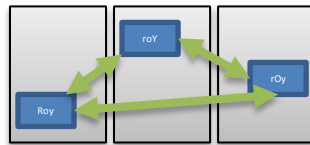
Niveau local / page



Interprétation itérative (3/4)

Interprétation globale

Exploitation du
contexte documentaire



Interaction avec des
opérateurs humains



~~Julien~~



Idem

Niveau global / fonds

Niveau local / page

Ancien prop.
Cure de Savigny
Julien

Résultat

Résultat

Résultat

Résultat

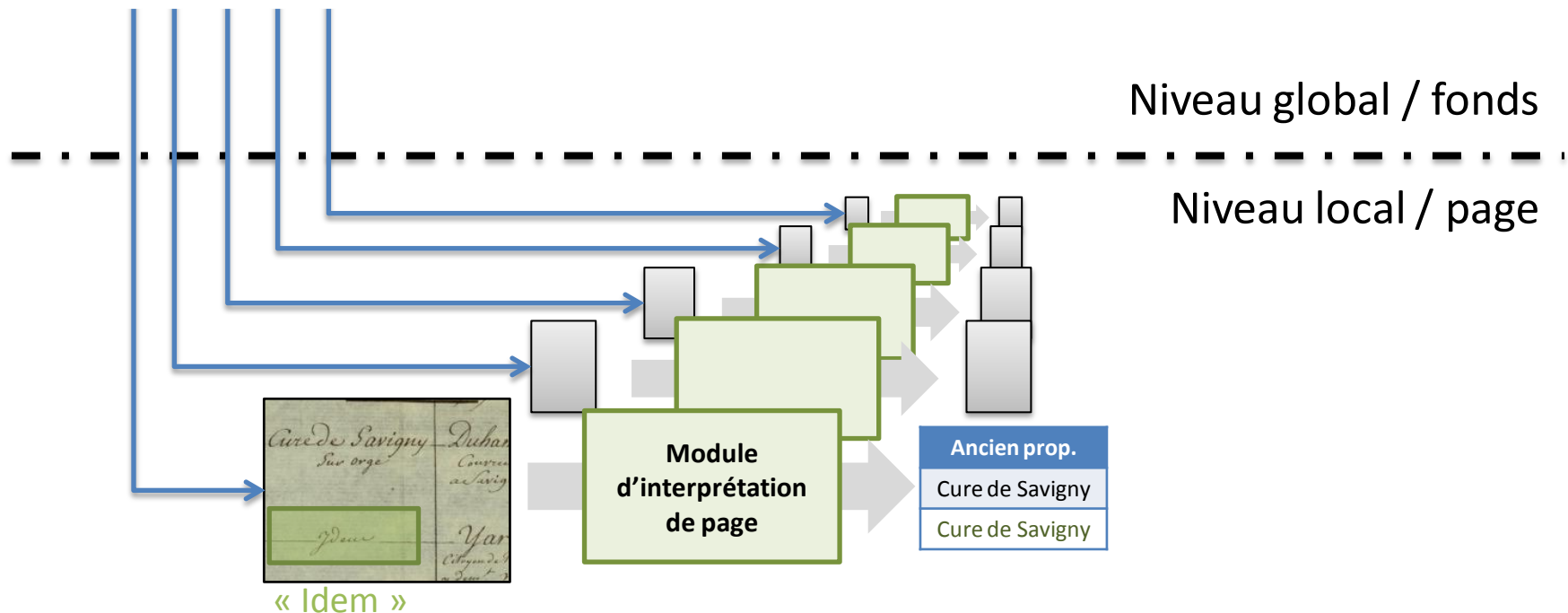
Résultat

Résultat

Interprétation itérative (4/4)

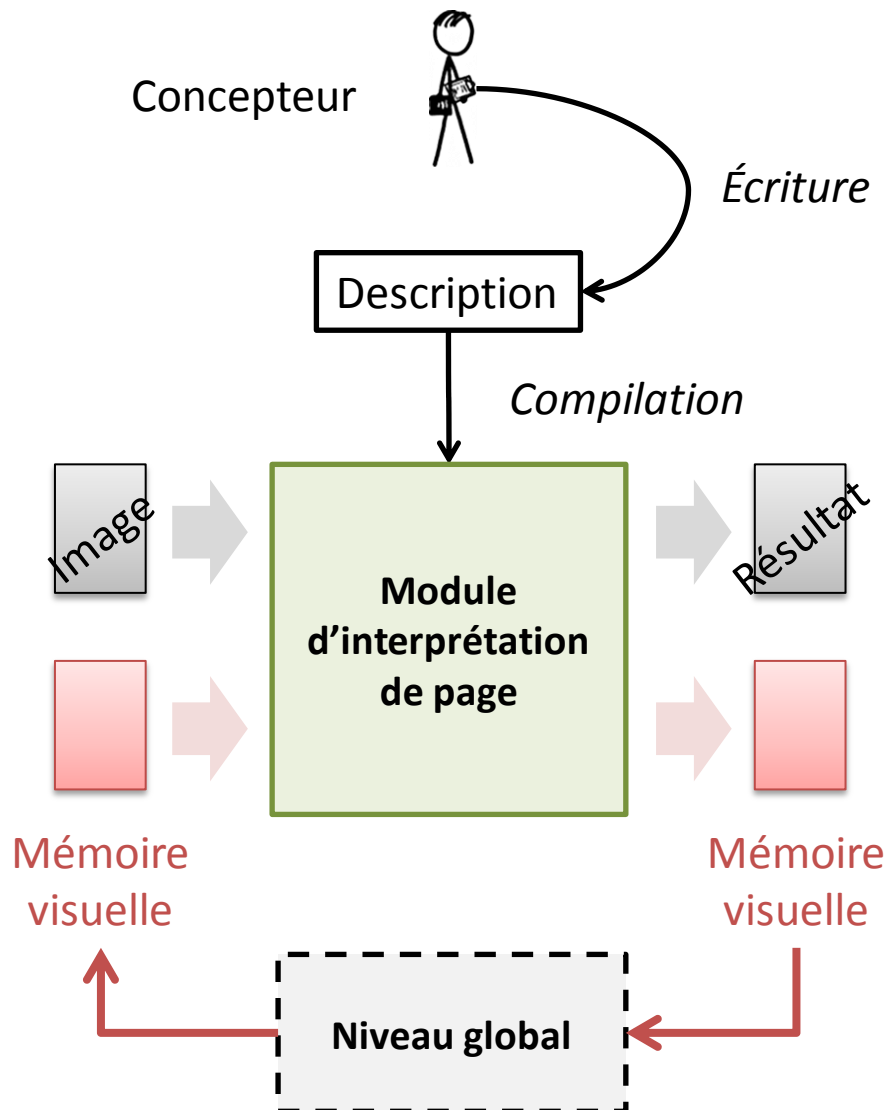
Réinterprétation locale avec des informations globales

1. Nouvelle interprétation complète avec les données externes
2. Production de nouveaux résultats



Notre mise en œuvre d'une approche itérative (1/11)

Adapter une approche déclarative



Extension en 3 étapes

1. Mémoire visuelle

- Support de la communication avec l'environnement

2. Langage de description

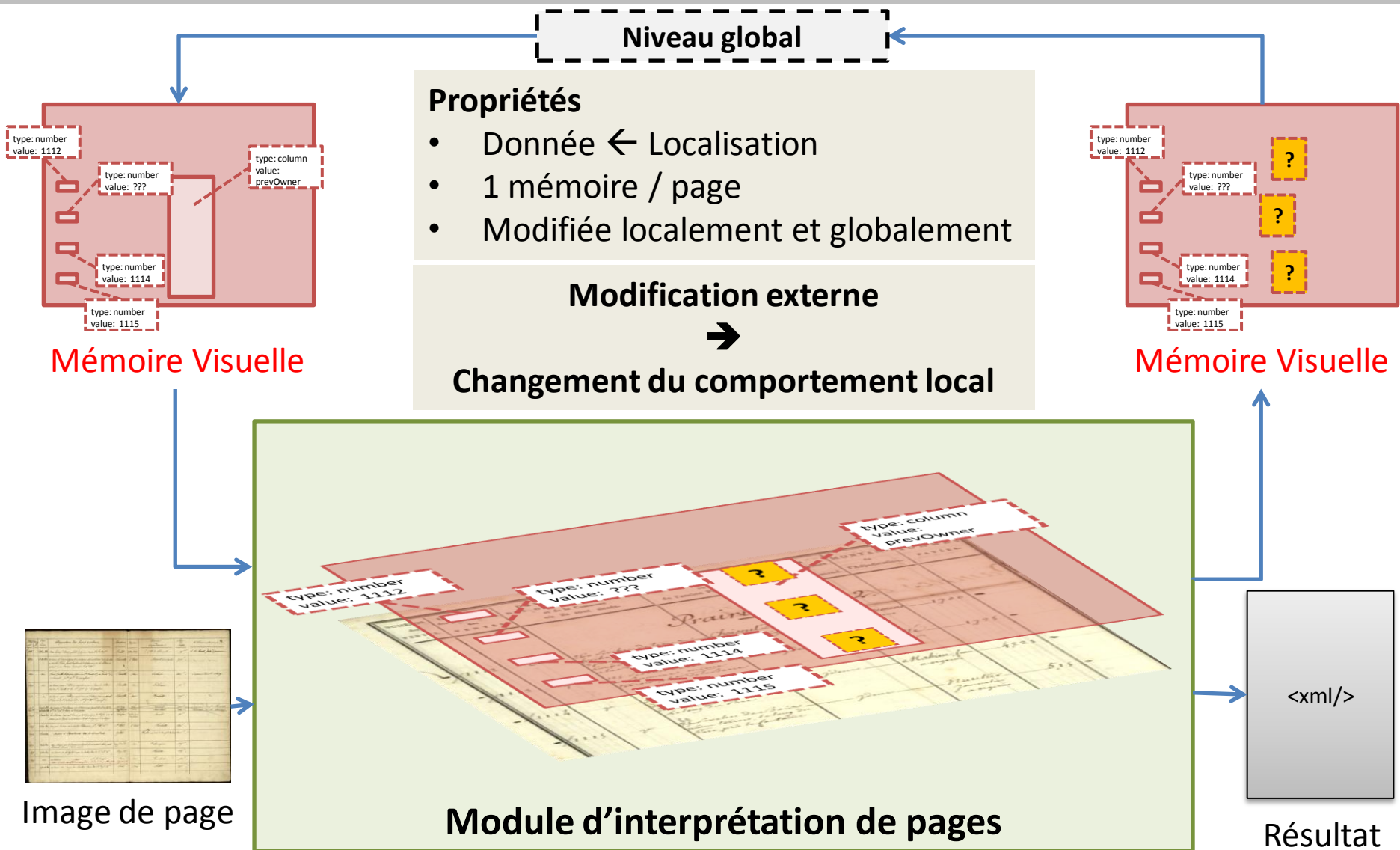
- Introduire de nouveaux opérateurs d'interaction

3. Architecture à deux niveaux

- Niveau global :
collecte et distribue les données
- Niveau local :
produit et exploite les données

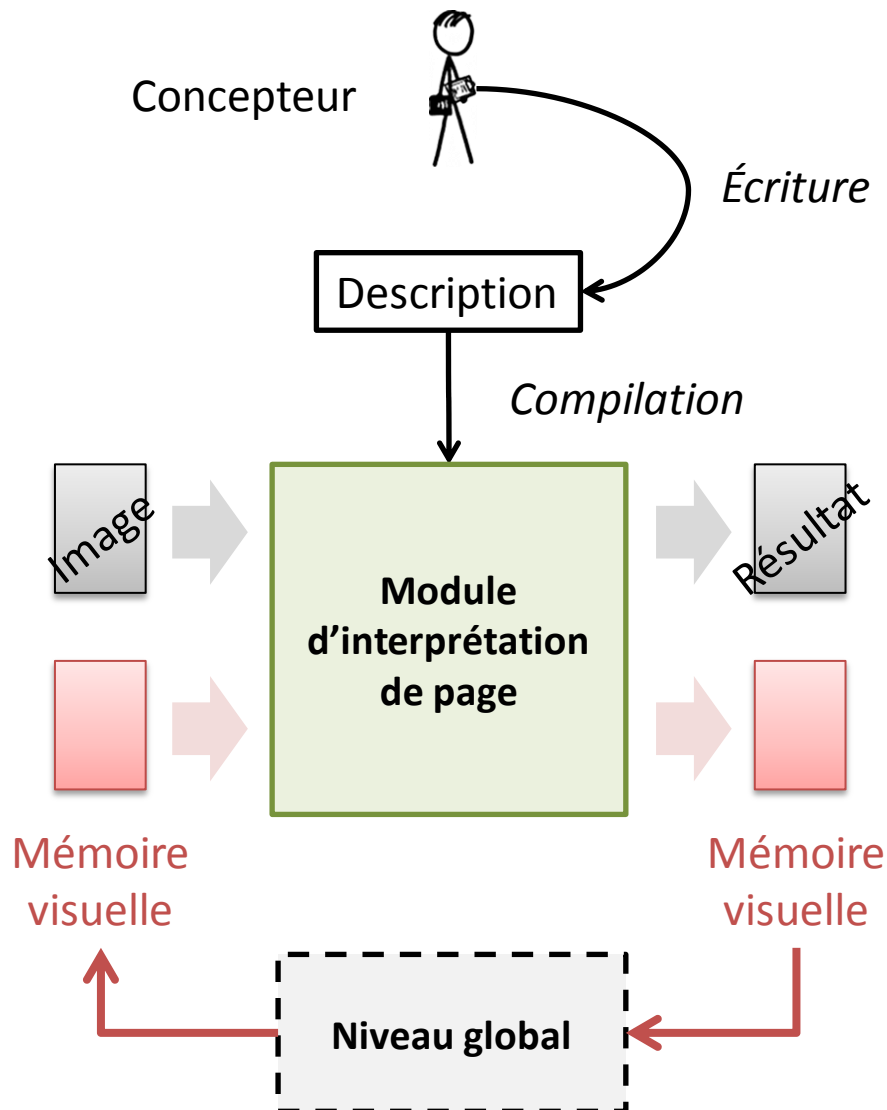
Notre mise en œuvre d'une approche itérative (2/11)

Mémoire visuelle



Notre mise en œuvre d'une approche itérative (3/11)

Adapter une approche déclarative



Extension en 3 étapes

1. Mémoire visuelle

- Support de la communication avec l'environnement

2. Langage de description

- Introduire de nouveaux opérateurs d'interaction

3. Architecture à deux niveaux

- Niveau global :
collecte et distribue les données
- Niveau local :
produit et exploite les données

Notre mise en œuvre d'une approche itérative (4/11)

Exemple de description

		Avril 1791.	
188	28	98 Souches & de Vignes en 3 pièces, terroir de Savigny, ch ^{re} sur la fontaine	Cure de Savigny Dubamel
189	r	149 Souches & de Vignes en 4 pièces, même terroir, ch ^{re} du hameau de la Pointe Crochet	Cure de Savigny Yart
190	r	Un demi arpent de terre, même terroir, ch ^{re} du mail	Cure de Savigny Duval
191	r	Un demi arpent de terre, même terroir, ch ^{re} du Nord du bien	Cure de Savigny Olivier

Description SANS interaction

```

page →
  @(colonneAncienProp)
  lireTousLesNoms("inconnu")

lireTousLesNoms(precedent) →
  @(detecterNom)
  nom = lireNom(precedent)
  @(dessous)
  LOOP(lireTousLesNoms(nom))

lireNom(precedent) →
  reco = reconnaitNom
  res = (si reco=="IDEM"
        alors precedent
        sinon reco)
  
```

Sans interaction : résultat peu fiable → propagation d'erreur

Notre mise en œuvre d'une approche itérative (5/11)

Nouveaux opérateurs de description

Détection des erreurs

RAISE_QUESTION(Question, Type) :
 Ajouter position → (Question, Type)
 en mémoire visuelle
 Poursuivre (éventuellement) l'interprétation

Reprise sur erreur

CATCH(Règle, Type) :
 Reprendre l'interprétation si une
 information de type Type est demandée dans Règle

Fusion d'informations externes

TRY(Règle, Type) :
 Si \exists donnée : Type \in mémoire visuelle,
 alors l'utiliser directement
 sinon appeler Règle
 et conserver le résultat en mémoire visuelle

Mécanisme proche
 de la gestion d'exceptions

Description AVEC interaction

```
page →
  @(colonneAncienProp)
  lireTousLesNoms("inconnu")

lireTousLesNoms(precedent) →
  @(detecterNom)
  nom = CATCH(lireNom(precedent),
              T_NOM)
  @(dessous)
  LOOP(lireTousLesNoms(nom))

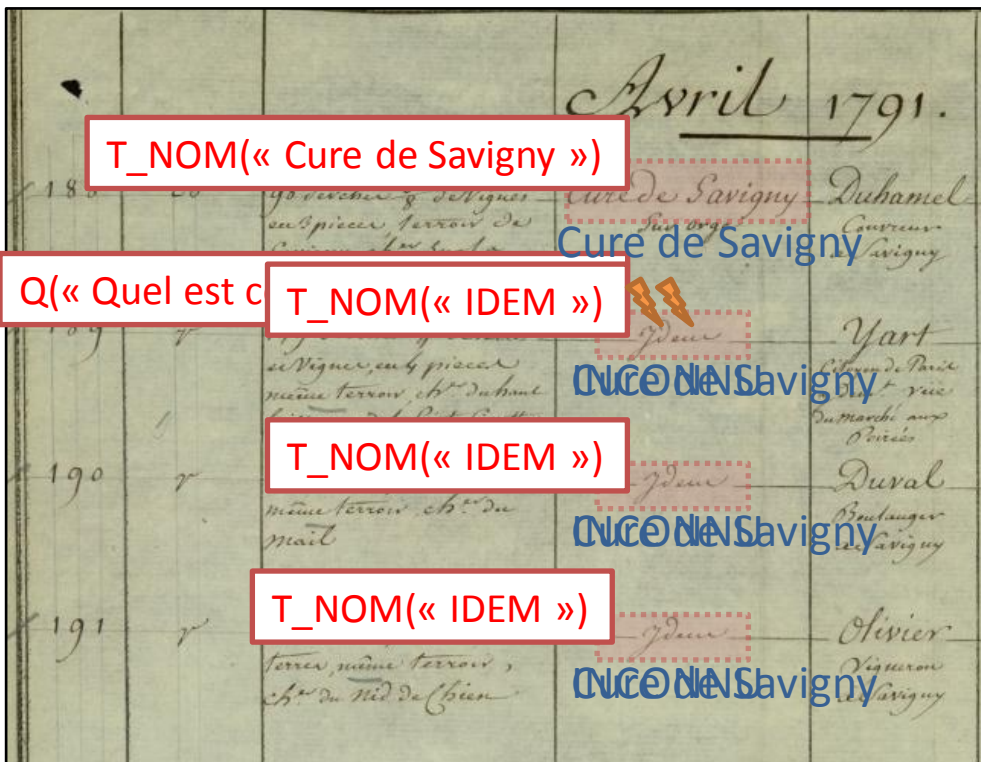
lireNom(precedent) →
  reco = TRY(reconnaitNomINT, T_NOM)
  res = (si reco=="IDEM"
        alors precedent
        sinon reco)

reconnaitNomINT →
  reconnaitNom
  ou RAISE_QUESTION("Quel est le nom
                    de cet ancien propriétaire ?",
                    T_NOM)
```

Notre mise en œuvre d'une approche itérative (6/11)

Progression au cours des itérations

Itération 1 :	Traitement	Interaction asynchrone
Itération 2 :	Traitement	...



Avec interaction :
correction des erreurs au plus tôt
➔ propagation limitée

Description **AVEC** interaction

```

page →
  @(colonneAncienProp)
  lireTousLesNoms("inconnu")

lireTousLesNoms(precedent) →
  @(detecterNom)
  nom = CATCH(lireNom(precedent),
              T_NOM)
  @(dessous)
  LOOP(lireTousLesNoms(nom))

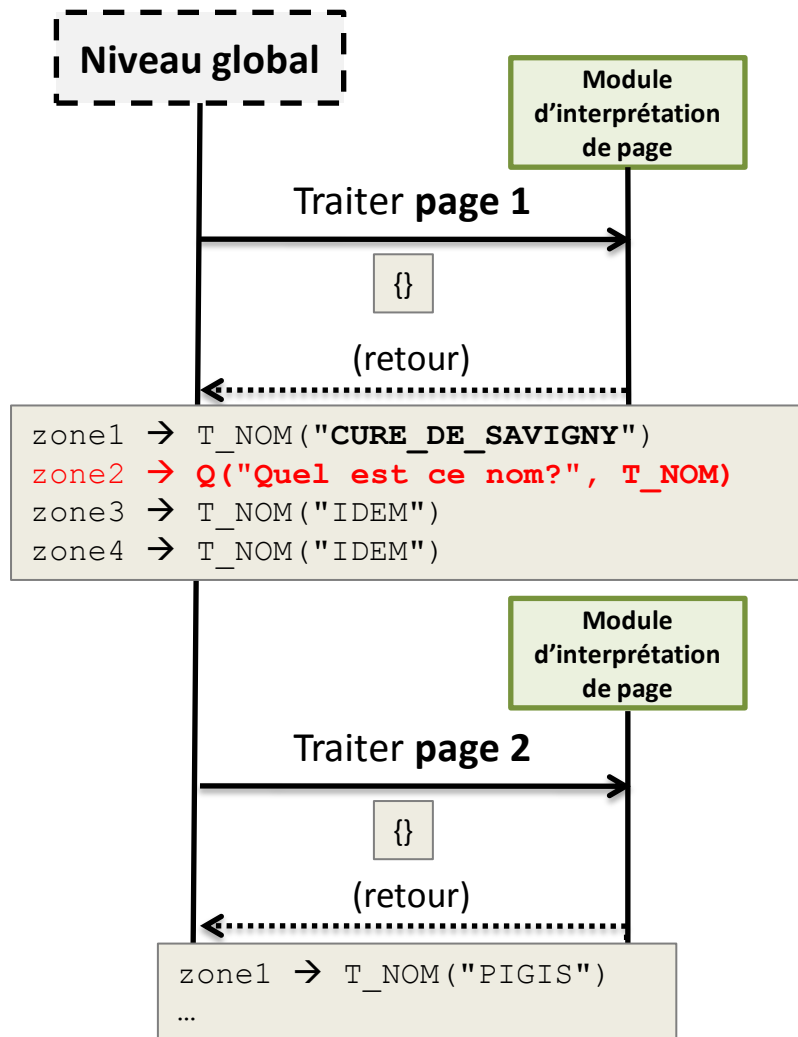
lireNom(precedent) →
  reco = TRY(reconnaitNomINT, T_NOM)
  res = (si reco=="IDEM"
        alors precedent
        sinon reco)

reconnaitNomINT →
  reconnaitNom
  ou RAISE_QUESTION("Quel est le nom
de cet ancien propriétaire ?",
T_NOM)
  
```

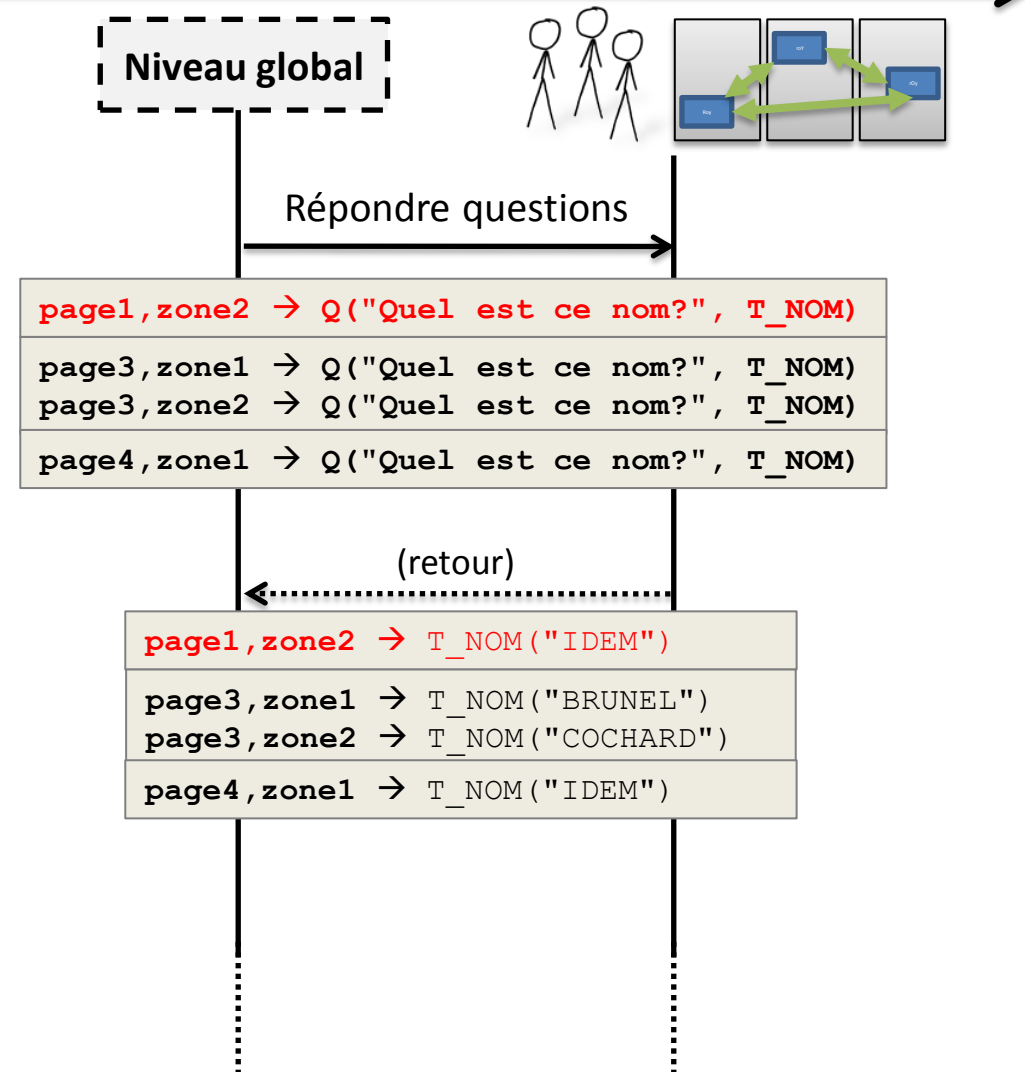

Notre mise en œuvre d'une approche itérative (7/11)

Exemple d'échange entre le niveau local et le niveau global

Itération 1 : traitement



Itération 1 : interaction *asynchrone*

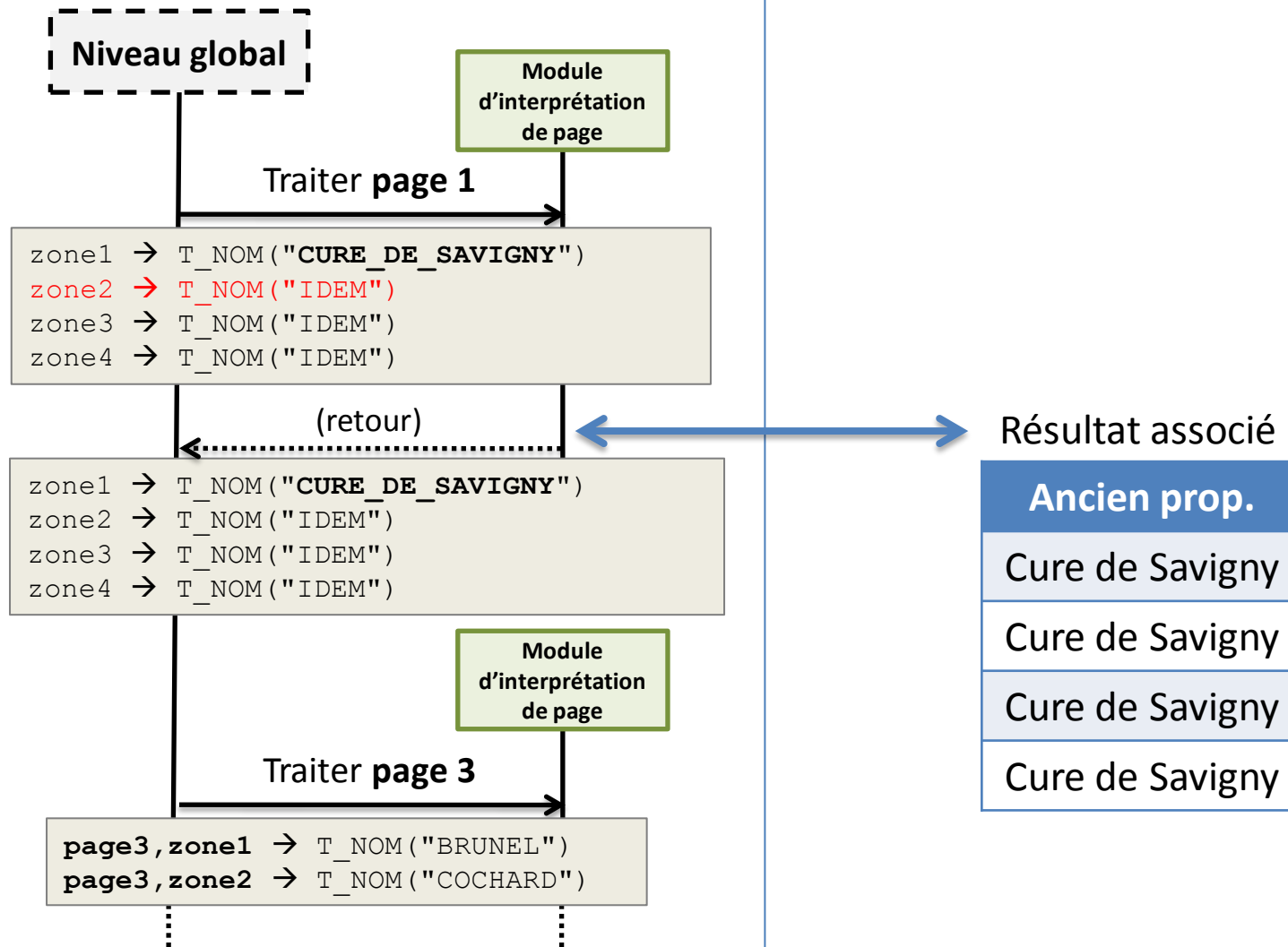


Notre mise en œuvre d'une approche itérative (8/11)

Exemple d'échange entre le niveau local et le niveau global

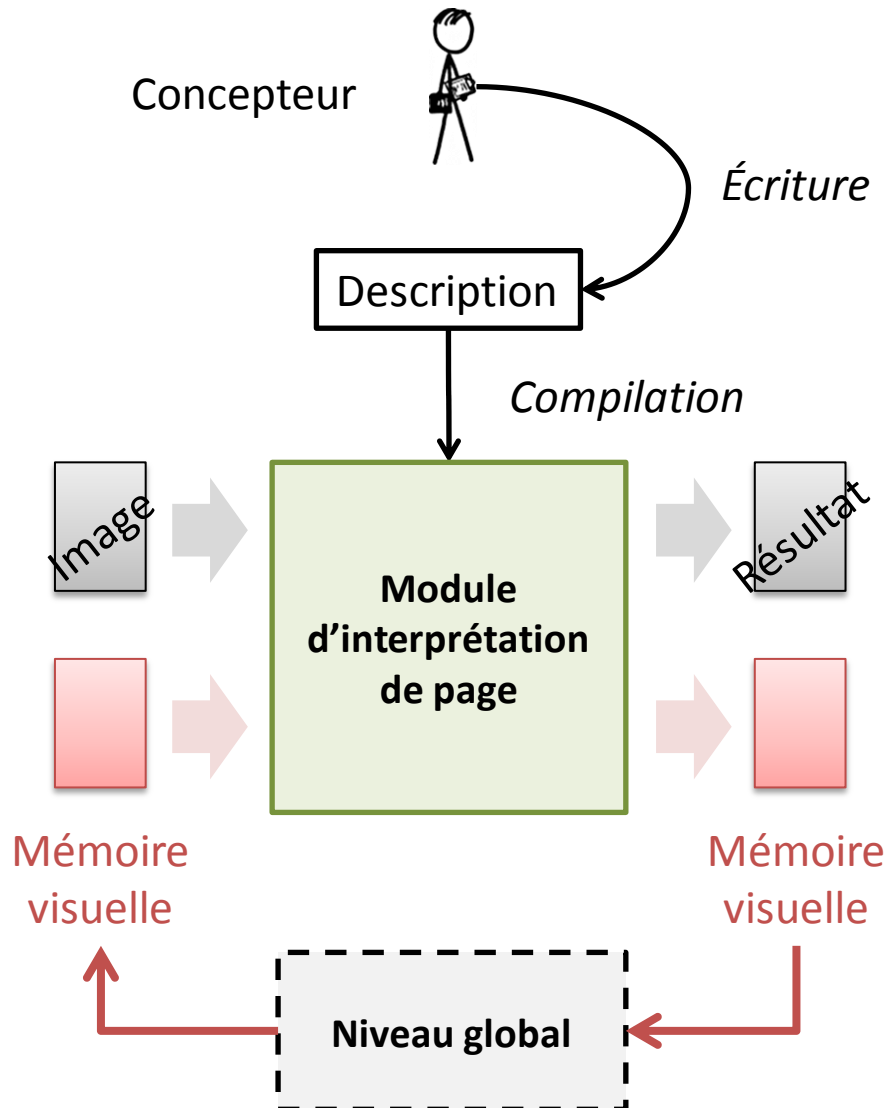
Itération 2 : traitement

etc.



Notre mise en œuvre d'une approche itérative (9/11)

Adapter une approche déclarative



Extension en 3 étapes

1. Mémoire visuelle

- ➔ Support de la communication avec l'environnement

2. Langage de description

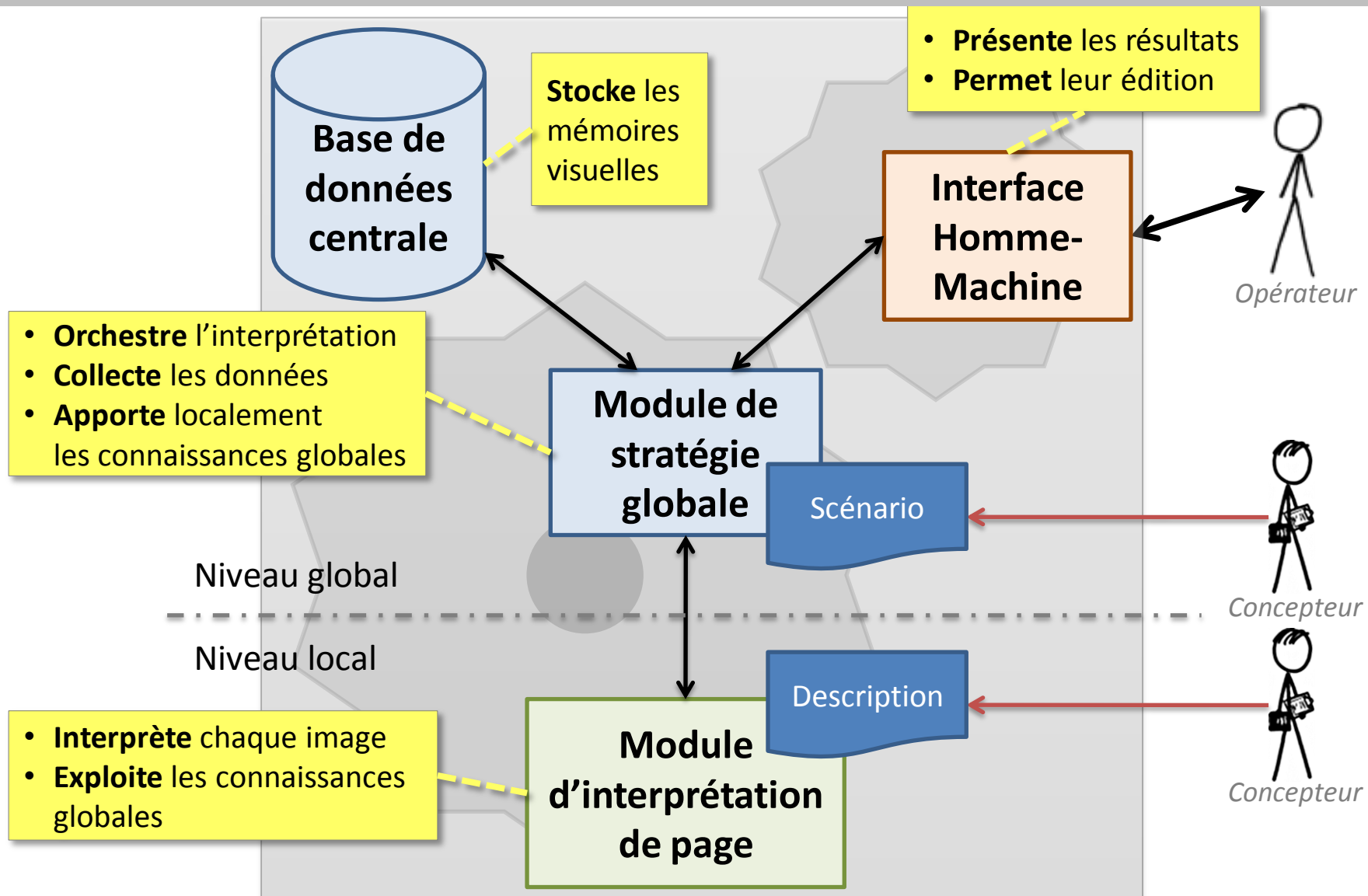
- ➔ Introduire de nouveaux opérateurs d'interaction

3. Architecture à deux niveaux

- ➔ Niveau global :
collecte et distribue les données
- ➔ Niveau local :
produit et exploite les données

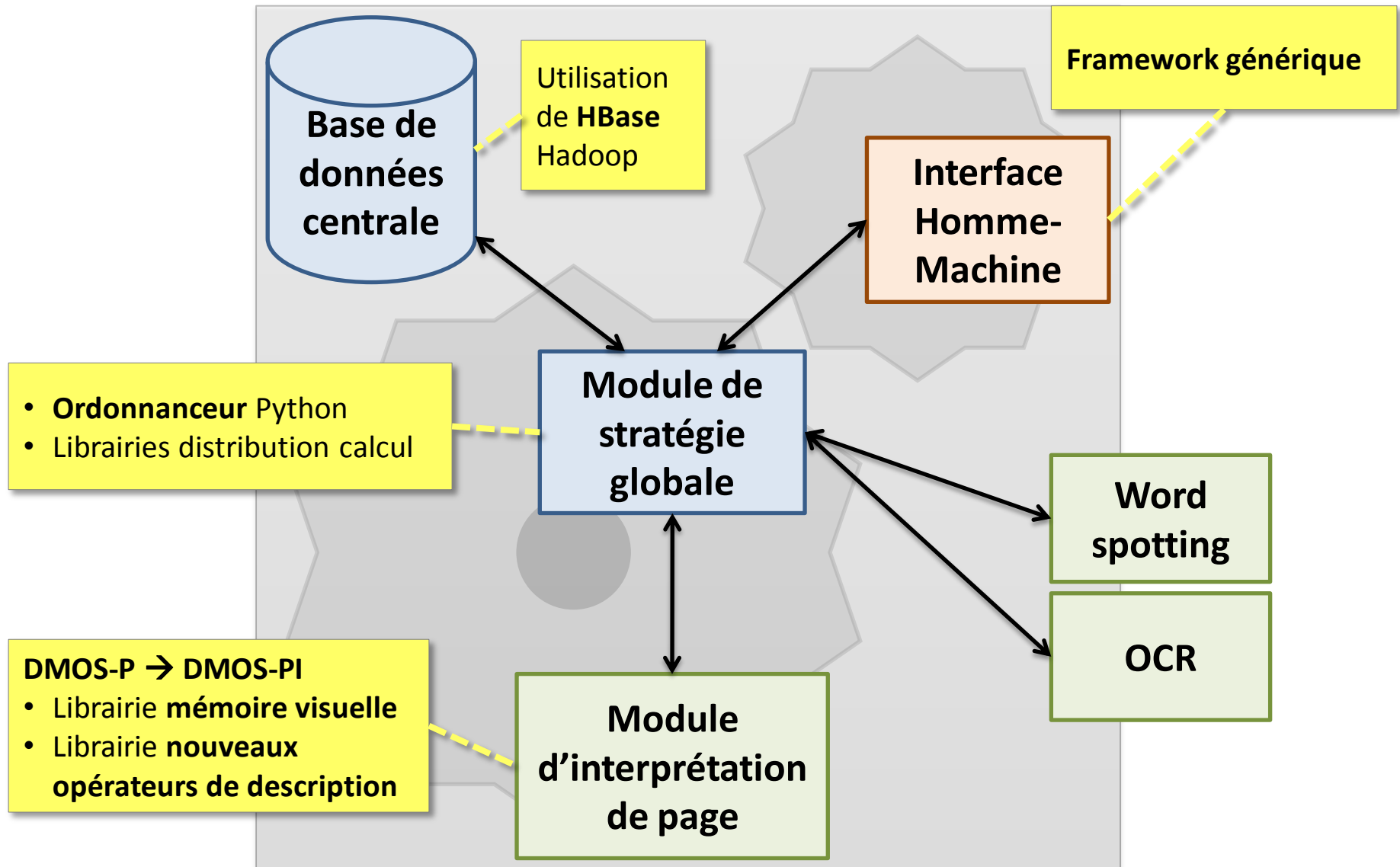
Notre mise en œuvre d'une approche itérative (10/11)

Architecture à deux niveaux



Notre mise en œuvre d'une approche itérative (11/11)

Implémentation de notre architecture



Bilan des contributions

Méthode pour rendre une approche déclarative itérative

1. Mémoire visuelle
2. Trois nouveaux opérateurs d'interaction
3. Architecture à deux niveaux

Avantages

Automatisation

- Au niveau local
 - Fusion des données produites par l'environnement
 - Échange asynchrone d'informations
- Au niveau global
 - Collecte et circulation de l'information
 - Interprétation par itération

Conception simple

- Description de la page
 - Comme si l'information externe était déjà disponible
 - Prévoir la gestion de l'incertain
- Définition d'un scénario global
 - Enchaîne les traitements
 - Ne se préoccupe pas de la production locale des données

Plan

1. Introduction
2. Résumé de l'état de l'art
3. Contribution : Interprétation itérative
4. Validation expérimentale et en production
 - a. Que valide-t-on ?
 - b. Expérimentation 1 : transcription de patronymes
 - c. Expérimentation 2 : correction de sous-segmentation
 - d. Utilisation du système en production
5. Conclusion et perspectives

Que valide-t-on ?

Intérêt des outils proposés

- Conception simple
- Possibilité de comparer différents scénarios

Intérêt d'une interprétation contextuelle et assistée

- Gain en qualité ou en coût de correction
- Évaluation du coût associé à un scénario délicat
 - Ne pas être dépendant de l'ergonomie des outils de correction
 - Éviter de comparer des actions de types différent
 - ➔ Isoler des fragments de scénarios réels

Plan

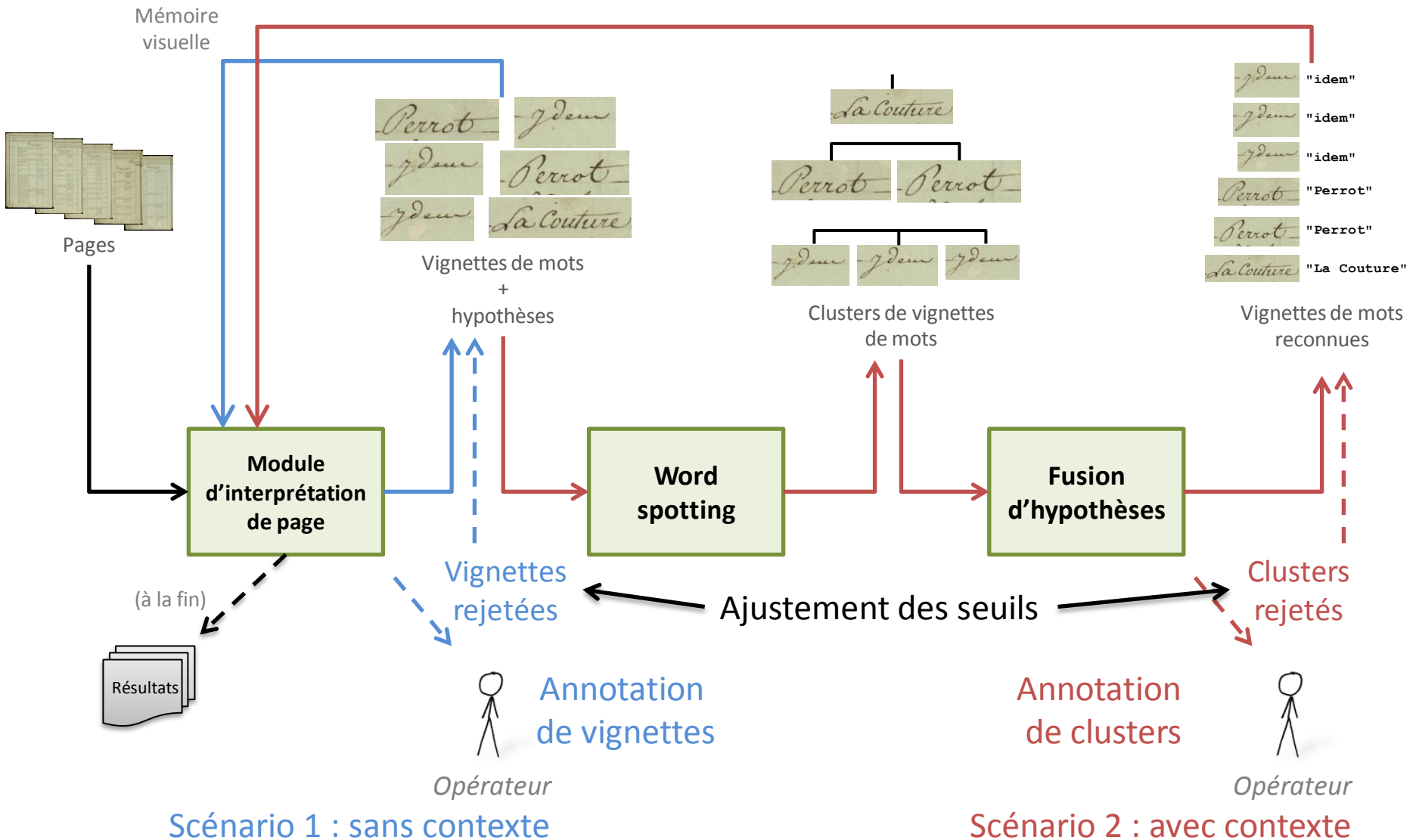
1. Introduction
2. Résumé de l'état de l'art
3. Contribution : Interprétation itérative
4. Validation expérimentale et en production
 - a. Que valide-t-on ?
 - b. Expérimentation 1 : transcription de patronymes
 - c. Expérimentation 2 : correction de sous-segmentation
 - d. Utilisation du système en production
5. Conclusion et perspectives

Expérimentation 1 : transcription de patronymes (1/4)

NUMÉROS des VENTES.	DATES des procès-verbaux des VENTES.	DESIGNATION DES OBJETS ALIÉNÉS, et de la Commune où ils sont situés.	INDICATION DE L'ANCIEN ÉTABLISSEMENT, ou de l'ancien Propriétaire.	NOM de l'Adjudicataire ou de son Command.	MONTANT de l'Adjudication.	SOMME PAYÉES, QUI RESTENT D	SOMME en capital.
<u>Ferrier 1791.</u>							
7.	3	5 arpent de pré terroir de St Germain - les arpaizon, ch ^e de la Boitelle	Curie de St Germain Les Arpaizon	Perrot m. de Boi à Arpaizon	5,25		
8	7	3 arpent de terre et hermine, en 2 pièces même terroir, ch ^e de grand St. m. de	Religieuses de St. Eutrope les chanceloup	Deliot culteur de Batim. à Arpaizon	2,075		
9	7	12 arpent 84 perches de terre, en une pièce même terroir, en chanceloup	Idem	Berson m. de Mouffettes à Marolles en chanceloup	10,600		
10	7	4 arpent 25 perches de vignes, en 2 pièces même terroir, devant au Parc de chanceloup, et au chemin de la Rocession	Idem	La Couture Bourgeois et Genard Garnier de devant le Parc à Arpaizon	4,550		
11	7	2 arpent de terre, même terroir, ch ^e de la Ceinture	Idem	Soret Patrimoine tracteur à Arpaizon	1,725		
12	7	3 arpent 75 perches de pré, même terroir	Picure de St. Guenard de Corbeil	Perrot m. de Boi à Arpaizon	3,300		

Expérimentation 1 : transcription de patronymes (2/4)

Comparaison de deux scénarios



Expérimentation 1 : transcription de patronymes (3/4)

Données traitées : environ 11 000 vignettes de patronymes

Évaluation

Approche	Taux d'annotation manuelle TM	Taux d'annotation automatique TA
Sans contexte	1 action / vignette	Vignettes dont la valeur est <ul style="list-style-type: none"> correcte déterminée automatiquement
Avec contexte (word spotting)	1 action / cluster	

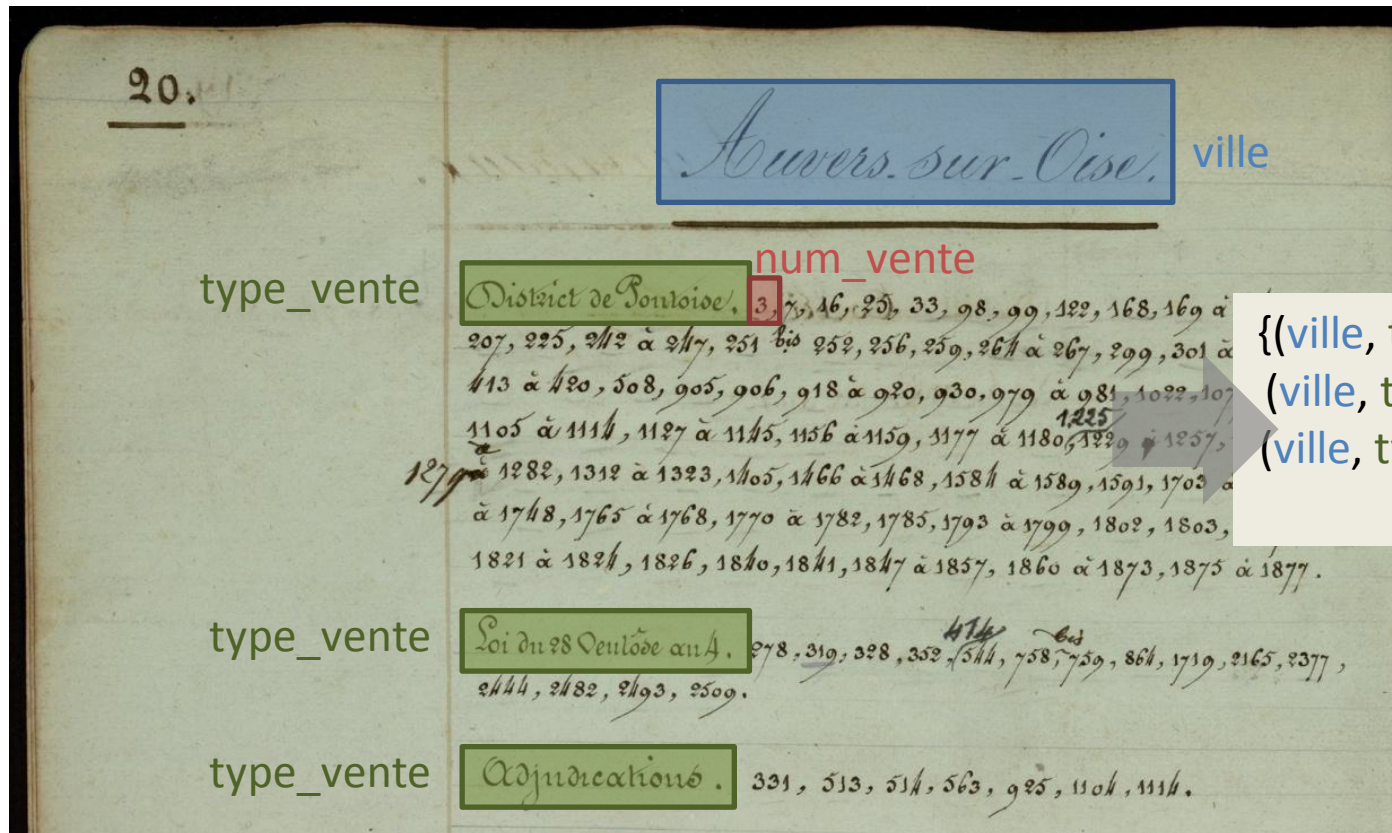
Résultats

Approche	Production de vérité terrain Tolérance : 1% d'erreur	Tâche d'indexation Tolérance : 20% d'erreur
Sans contexte	TM = 78,8% TA = 20,2%	TM = 20,2% TA = 59,8%
Avec contexte (word spotting)	TM = 55,6% TA = 43,4%	TM = 10,7% TA = 69,3%

Plan

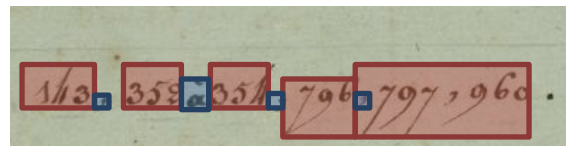
1. Introduction
2. Résumé de l'état de l'art
3. Contribution : Interprétation itérative
4. Validation expérimentale et en production
 - a. Que valide-t-on ?
 - b. Expérimentation 1 : transcription de patronymes
 - c. Expérimentation 2 : correction de sous-segmentation
 - d. Utilisation du système en production
5. Conclusion et perspectives

Expérimentation 2 : correction de sous-segmentation (1/4)



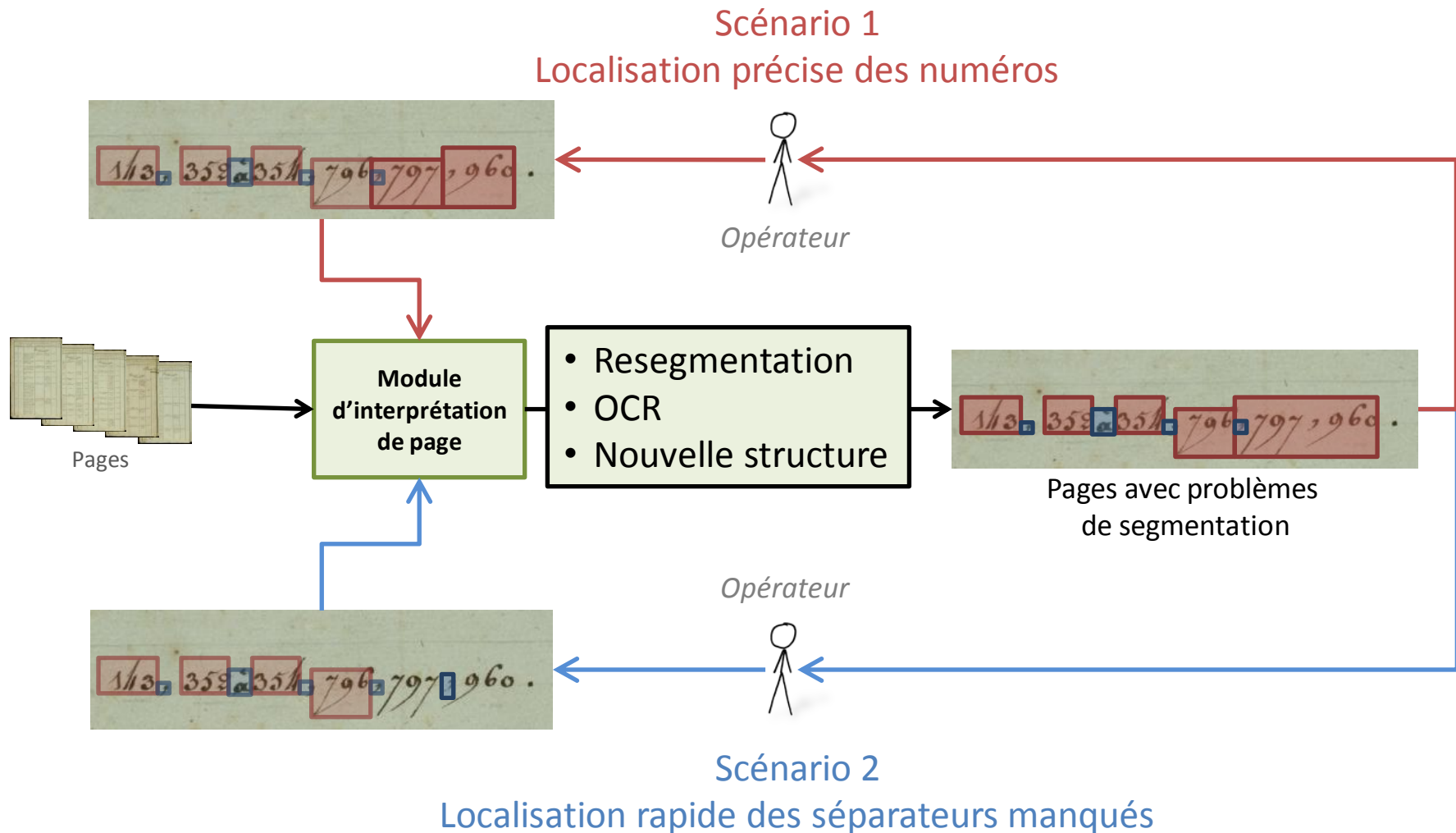
```
{(ville, type_vente, num_vente),
 (ville, type_vente, num_vente),
 (ville, type_vente, num_vente),
 ...}
```

Problème principal : sous-segmentation



Expérimentation 2 : correction de sous-segmentation (2/4)

Comparaison de deux scénarios



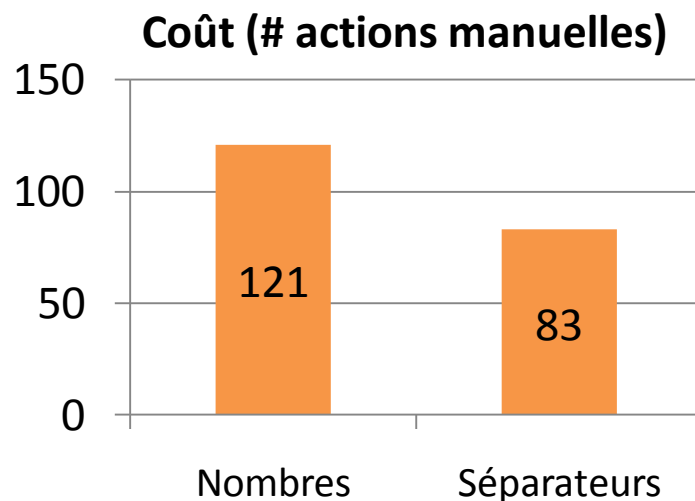
Expérimentation 2 : correction de sous-segmentation (3/4)

Données traitées : 50 pages (1637 vignettes de nombres)

Évaluation

- Quantité de travail manuel requis pour corriger les 121 nombres mal segmentés
- 1 action manuelle = ...
 - 1 localisation précise de nombre
 - 1 localisation précise de séparateur

Résultats



Gain :
30% d'actions en moins
en corrigeant la cause des erreurs
(détection du séparateur)

Expérimentation 2 : correction de sous-segmentation (4/4)

Description de la page

- 1 description pour les 2 scénarios
→ Possibilité de les combiner
- Basée sur des questions optionnelles
→ Opérateur **OPT(Type, Règle)**

```
extraireNumVenteLigne →  
  @(positionLigne)  
  lstsep = OPT(T_SEP, detecterTousSep)  
  OPT(T_NUM, extraireNbreEntreSep(lstsep))  
  
detecterTousSep →  
  sep = TRY(T_SEP, separateur)  
  res = sep :: LOOP(detecterTousSep)  
  
extraireNbreEntreSep (lstsep) →  
  @(entreseparateur lstsep)  
  TRY(T_NUM, extraireNombre)  
  LOOP(extraireNbreEntreSep(lstsep))
```

À retenir

- **Correction de la cause des erreurs et non des conséquences**
- Possibilité de laisser la charge de la détection d'erreur à l'opérateur humain
- Possibilité de poser plusieurs questions de types différents

Plan

1. Introduction
2. Résumé de l'état de l'art
3. Contribution : Interprétation itérative
4. Validation expérimentale et en production
 - a. Que valide-t-on ?
 - b. Expérimentation 1 : transcription de patronymes
 - c. Expérimentation 2 : correction de sous-segmentation
 - d. Utilisation du système en production
5. Conclusion et perspectives

Utilisation du système en production

Données traitées (partenariat avec les Archives des Yvelines)

- Environ 1200 pages de documents d'archives
- 6700 lignes de ventes
- 11 000 vignettes de patronymes

Techniques mises en œuvre en conditions réelles

- Au niveau local
 - Descriptions avec plusieurs types de données échangées
 - Possibilité de détection d'erreur manuelle
 - Remise en cause en profondeur des résultats
 - ➔ Avec des descriptions simples
- Au niveau global
 - Gestion centralisée du compromis automatique/manuel
 - Optimisation de séquences, regroupement de mots visuellement similaires
 - ➔ Grâce à la séparation global/local facilitant la collaboration en conception

Plan

1. Introduction
2. Résumé de l'état de l'art
3. Contribution : Interprétation itérative
4. Validation expérimentale et en production
5. Conclusion et perspectives

Conclusion

Notre contribution

Méthode pour mettre en œuvre une interprétation itérative

- Basée sur une approche déclarative
- Extension en 3 étapes
 1. **Mémoire visuelle**
 2. **Langage de description**
 3. **Architecture à 2 niveaux**

Avantages

- **Des outils proposés**
 - **Réintégration locale** de connaissances globales
 - Interaction **asynchrone**
 - Conception **centrée page**
- **De l'approche**
 - **Réutilisable** : basée sur des mécanismes standards
 - **Validée** : en théorie, en pratique, sur les expérimentations et en production

Perspectives

- **Consolidation : vers l'apprentissage de règles**
 - Difficulté d'amorçage du système réduite
 - Mécanisme itératif favorable à la remise en cause du modèle
- **Extension : vers une exploitation en consultation**
 - Utilisateurs actifs (crowdsourcing + traitements automatiques)