

# Inférence semi-automatique et interactive de règles avec ou sans vérité terrain pour la reconnaissance de structure de documents

Cérès Carton

Encadrée par Bertrand Couïasnon et Aurélie Lemaitre

23 Mars 2016



# Analyse d'images de documents

Adaptation rapide d'un système de reconnaissance de structure de documents

- Pour le concepteur
- En temps et en effort

Exploitation efficace documents numérisés

⇒ connaissance de la structure

- Indexation
- Navigation
- Classification
- Reconnaissance de caractères
- ...

# Analyse de la structure de documents

RICHARD Patricia  
 ECLOSE N°1  
 57830 KERPICH AUX BOIS  
 Tel: 03.68.59.33.53  
 réf: FQNAO23

MANIF Assurances  
 1 PLACE FOCH  
 80490 AMIENS

le 11 août 2006

## Structure logique d'un document

- Structure physique
- Étiquetage logique
- Organisation hiérarchique

Objet: résiliation d'assurance auto.

Madame, Monsieur,

Suite à mon accident de voiture survenu le 20 janvier de cette année, il m'est apparu que votre assurance auto n'offrait pas toutes les garanties escomptées, notamment en ce qui concerne les délais de remboursement et de réparation du véhicule. C'est pourquoi je désire résilier la dite assurance.

En vous remerciant de bien vouloir prendre note de ce changement exposé ci-dessus, je vous prie de recevoir, Madame, Monsieur, mes salutations

P. RICHARD  


# Analyse de la structure de documents

## Structure logique d'un document

- Structure physique
- Étiquetage logique
- Organisation hiérarchique

coordonnées expéditeur

RICHARD Patricia  
ECLUSE N1  
57830 KERPICH AUX BOIS  
Tel: 03.68.59.33.53  
réf: FQNA023

coordonnées destinataire

MANIF Assurances  
1 PLACE FOCH  
80490 AMIENS

le 11 août 2006

date

objet

Objet: résiliation d'assurance auto.

ouverture

Madame, Monsieur,

corps de texte

Suite à mon accident de voiture survenu le 20 janvier de cette année, il m'est apparu que votre assurance auto m'offrait pas toutes les garanties escomptées, notamment en ce qui concerne les délais de remboursement et de réparation du véhicule. C'est pourquoi je désire résilier la dite assurance.

En vous remerciant de bien vouloir prendre note de ce changement exposé ci-dessus, je vous prie de recevoir, Madame, Monsieur, mes salutations

P. RICHARD  
nom  
signature



# Limites des systèmes actuels

## Difficulté croissante de la reconnaissance de structure

- Documents variés
- Documents complexes



## Adaptation d'un système de reconnaissance

⇒ modélisation de connaissances complexes

- Constitution d'un corpus étiqueté bien sélectionné
  - Très coûteux en temps
- Expression manuelle de connaissances
  - Bonne vue sur les données et leur variabilité

# Illustration 1 : corpus avec vérité terrain

- Courriers manuscrits en français
- Corpus de la compétition internationale RIMES
- Tâche de reconnaissance de structures
  - Localisation de 8 types de blocs
- Vérité terrain au niveau bloc
  - Boîtes englobantes
  - Transcription

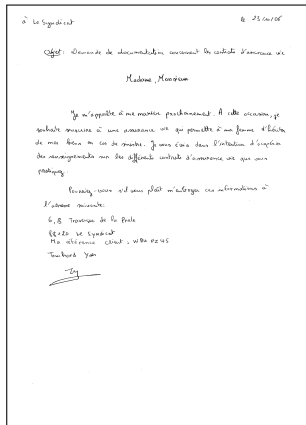
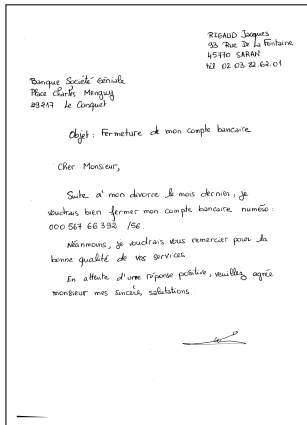
The diagram shows a handwritten letter with several colored boxes highlighting different parts:

- Sender (pink box):** Sandrine Brand, 21 rue Principale, 90340 Froidfontaine, Réf Client: R0ETS88
- Recipient (red box):** Froidfontaine, le 10 juillet 2006
- Address (cyan box):** GDF, Grande Allée de Tency, 41800 St Christophe en Brannois
- Subject (green box):** Objet: Demande de remise gracieuse
- PS (blue box):** PS: Dernière facture GDF caute chateau, kuest de garnille
- Salutation (yellow box):** Madame, Monsieur,
- Main Text (orange box):** La dernière facture GDF que j'ai reçue le 28 juin 2006 s'élève à un montant de 175,13 € (vous en trouverez une copie jointe à ce courrier). Licenciée de mon travail depuis le 1er juin 2006, je suis actuellement sans travail. Mon mari est également en recherche d'emploi et nous avons deux enfants, un de 3 ans et un autre de 11 ans. Je suis donc dans l'impossibilité de régler la somme demandée, aussi je me permets de vous solliciter afin d'obtenir, à titre gracieux, la remise ou la modération de la somme réclamée. Veuillez agréer, Madame, Monsieur, mes salutations distinguées.
- Signature (green box):** Sandrine Brand, SB Brand

# Illustration 1 : corpus avec vérité terrain

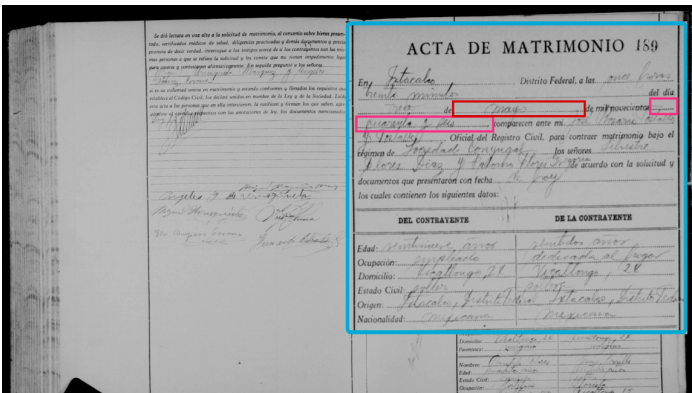
## Difficultés

- Segmentation de la page en blocs
- Règles usuelles non respectées



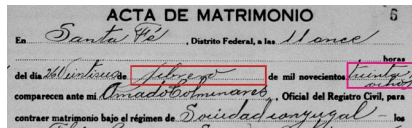
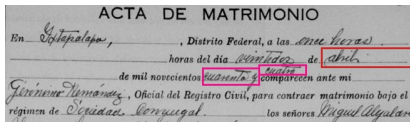
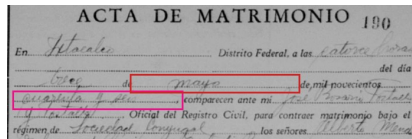
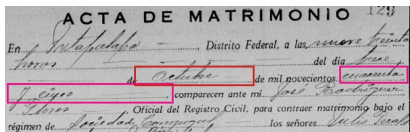
# Illustration 2 : corpus sans vérité terrain

- Registres de mariages mexicains
  - 30 000 documents
- Compétition organisée par Family Search à HIP 2013
- Localisation des champs manuscrits mois et année dans un acte pré-imprimé



## Illustration 2 : corpus sans vérité terrain

- Difficulté : mise en page variable des pré-imprimés
  - Position des champs varie
- Combien de variantes ?



- Comment détecter les variantes ?
  - Regarder un par un les 30 000 documents ?

- ① État de l'art
- ② Philosophie de notre contribution
- ③ Présentation détaillée de la méthode EWO
- ④ Validation
- ⑤ Conclusion

## ① État de l'art

Méthodes statistiques

Méthodes syntaxiques

Inférence grammaticale

## ② Philosophie de notre contribution

## ③ Présentation détaillée de la méthode EWO

## ④ Validation

## ⑤ Conclusion

# Méthodes statistiques

- Apprentissage automatique de l'étiquetage par un modèle statistique
  - Niveau bloc
  - Niveau pixel
- À partir d'un échantillon étiqueté

## Différents modèles possibles

- Champs aléatoire de Markov [Lemaitre, 2007]
- Réseaux bayésiens [Le Bourgeois, 2001]
- Champs aléatoires conditionnels [Montreuil, 2010]  
[Chaudhury, 2009]
- Réseaux de neurones [Rangoni, 2008]

## Nouveau type de documents $\Rightarrow$ base d'apprentissage étiquetée

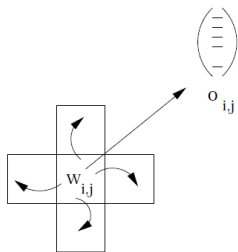
- Coût de production
- Sélection des données



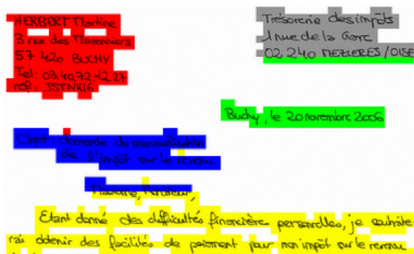
# Méthodes statistiques - exemple

Apprentissage du modèle : pour chaque pixel

- Étiquette
- Caractéristiques
  - Position
  - Texture
- Voisinage



Résultat obtenu



# Méthodes syntaxiques

Explicites des connaissances par un expert dans un modèle

- À base de règles [Fisher, 1991]
  - Difficilement adaptables
  - Peu flexibles
- À base de grammaires [Conway, 1993] [Krishnmoorthy, 1993]
  - Adapté aux structures hiérarchiques
  - Bidimensionnelles [Coüasnon, 2006]
  - Stochastiques [Tateisi, 94] [Maroneze, 2011]

Nouveau type de documents

⇒ description manuelle de nouvelles règles

- Échantillon de taille limitée
- Choix des exemples

# Méthodes syntaxiques - illustration

## Description des connaissances syntaxiques

objetCourrier ::=

AT(milieuGauche) &&

ligneLongue L &&

AT(sousLigne L) &&

ligneDecalee.

objetCourrier ::=

AT(milieuGauche) &&

ligneLongue.

Mme JOUANE Colette  
16 rue de la Paix  
67240 KALTENHOUSE  
tel : 0357 53 84 95

NAIF ASSURANCES  
125 rue Nationale de la Vierge  
13015 MARSEILLE

objet : demande de catalogue des produits proposés

réf client : VRVN W67

Madame Jouane,

Je suis actuellement client auprès de vos services pour mon véhicule Ford.

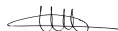
Nos contrats habitation (résidence principale et secondaires) arrivent à échéance au 31 décembre 2006. (Je suis actuellement assuré chez ARNAC ASSU)

Pourriez-vous me faire parvenir le catalogue de vos prestations relatives à l'habitation ?

Dans cette attente,

Croyez, Madame Jouane, à l'assurance de mes salutations distinguées

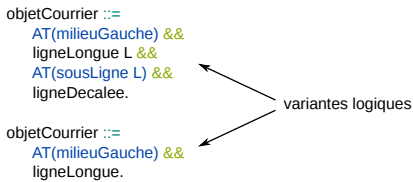
le 27.09.06



→ Type de règles à inférer

# Méthodes syntaxiques - illustration

## Description des connaissances syntaxiques



Mme JOUANE Colette  
 16 rue de la Paix  
 67240 KALTENHOUSE  
 tel : 03 57 53 84 95

NAIF ASSURANCES  
 125 rue Nationale de la Vierge  
 13015 MARSEILLE

objet : demande de catalogue des produits proposés

réf client : VRVN W67

Madame, Jorriou,

Je suis actuelle client auprès de vos services pour mon véhicule Ford.

Nos contrats habitation (résidence principale et secondaires) arrivés à échéance au 31 décembre 2006. (Je suis actuellement assurée chez ARNAC ASSU)

Pourriez-vous me faire parvenir le catalogue de vos prestations relatives à l'habitation ?

Dans cette attente,

Croyez, Madame, Jorriou, à l'assurance de mes salutations distinguées

le 27.09.06

→ Type de règles à inférer

# Méthodes syntaxiques - illustration

## Description des connaissances syntaxiques

objetCourrier ::=  
 AT(milieuGauche) &&  
 ligneLongue L &&  
 AT(sousLigne L) &&  
 ligneDecalee.

objetCourrier ::=  
 AT(milieuGauche) &&  
 ligneLongue.

ligneLongue L :-  
 segment L,  
 largeur > 1000px.

milieuGauche :-  
 Xd = 49% ImageWidth,  
 Yd = 32% ImageHeight,  
 X1 = 0,  
 Y1 = 15% ImageHeight,  
 X2 = 98% ImageWidth,  
 Y2 = 50% ImageHeight.

...

variantes logiques



→ Type de règles à inférer

Mme JOUANE Colette  
 16 rue de la Paix  
 67240 KALTENHOUSE  
 tel : 0357 53 84 95

NAIF ASSURANCES  
 125 rte Nationale de la Viste  
 13015 MARSEILLES

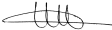
objet : demande de catalogue des produits proposés  
 réf client : VRVN W67

Madame Jorjieu,  
 Je suis actuelle client auprès de vos services pour mon  
 véhicule Ford.

Nos contrats habitation (résidence principale et secondaires)  
 arrivés à échéance au 31 décembre 2006. (Je suis actuelle-  
 -ment assurée chez ARNAC ASSU)

Pourriez-vous me faire parvenir le catalogue de vos  
 prestations relatives à l'habitation ?

Dans cette attente,  
 Cordz. Madame Jorjieu, à l'assurance de mes salutations  
 distinguées

le 27.09.06 

# Méthodes syntaxiques - illustration

## Description des connaissances syntaxiques

objetCourrier ::=

AT(milieuGauche) &&  
 ligneLongue L &&  
 AT(sousLigne L) &&  
 ligneDecalee.

objetCourrier ::=

AT(milieuGauche) &&  
 ligneLongue.

ligneLongue L :-

segment L,  
 largeur > 1000px.

milieuGauche :-

Xd = 49% ImageWidth,  
 Yd = 32% ImageHeight,  
 X1 = 0,  
 Y1 = 15% ImageHeight,  
 X2 = 98% ImageWidth,  
 Y2 = 50% ImageHeight.

...

→ Type de règles à inférer

← variantes logiques

← propriétés physiques

← opérateur de position

Mme JOUANE Colette  
 16 rue de la Paix  
 67240 KALTENHOUSE  
 tel : 0357 53 84 95

NAIF ASSURANCES  
 125 rte Nationale de la Viste  
 13015 MARSEILLES

objet : demande de catalogue des produits proposés

réf client : VRVN W67

Madame Jouane,

Je suis actuelle client auprès de vos services pour mon véhicule Ford.

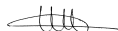
Nos contrats habitation (résidence principale et secondaires) arrivés à échéance au 31 décembre 2006. (Je suis actuellement assurée chez ARNAC ASSU)

Pourriez-vous me faire parvenir le catalogue de vos prestations relatives à l'habitation ?

Dans cette attente,

Croyez, Madame, Jouane, à l'assurance de mes salutations distinguées

le 27.09.06



# Comparaison des méthodes

- Modélisation de connaissances complexes
- Adaptation facile et rapide à un nouveau type de documents

	Statistiques	Syntaxiques
Apprentissage	automatique	manuel
Données d'apprentissage	étiquetées	non étiquetées
		taille limité
Gestion de l'incertitude	oui	selon les méthodes
Gestion des cas rares	difficile	oui
Pouvoir d'expression	limité	élevé

# Comparaison des méthodes

- Modélisation de connaissances complexes
- Adaptation facile et rapide à un nouveau type de documents

	Statistiques	Syntaxiques
Apprentissage	automatique	manuel
Données d'apprentissage	étiquetées	non étiquetées
		taille limité
Gestion de l'incertitude	oui	selon les méthodes
Gestion des cas rares	difficile	oui
Pouvoir d'expression	limité	élevé

Objectif : combiner les avantages

- Méthodes syntaxiques
  - Pouvoir d'expression et gestion des cas rares
- Méthodes statistiques
  - Apprentissage automatique → sans le coût de l'annotation



# Inférence dans les systèmes syntaxiques

## Inférence grammaticale

- Nombreux domaines d'application
  - Bioinformatique, linguistique, reconnaissance de la parole
- Nombreux obstacles à l'inférence grammaticale [de la Higuera, 2005]
  - Bidimensionnalité
  - Données bruitées
  - Données d'apprentissage étiquetées nécessaires
  - Pas de combinaison avec des connaissances a priori

# Inférence dans les systèmes syntaxiques

## Inférence grammaticale

- Nombreux domaines d'application
  - Bioinformatique, linguistique, reconnaissance de la parole
- Nombreux obstacles à l'inférence grammaticale [de la Higuera, 2005]
  - Bidimensionnalité
  - Données bruitées
  - Données d'apprentissage étiquetées nécessaires
  - Pas de combinaison avec des connaissances a priori

## Inférence grammaticale en reconnaissance de documents

[Shilman, 2005]

Utilisateur fournit	Inférence
grammaire de pages	caractéristiques
pages annotées	paramètres

# Inférence dans les systèmes syntaxiques

## Inférence grammaticale

- Nombreux domaines d'application
  - Bioinformatique, linguistique, reconnaissance de la parole
- Nombreux obstacles à l'inférence grammaticale [de la Higuera, 2005]
  - Bidimensionnalité
  - Données bruitées
  - Données d'apprentissage étiquetées nécessaires
  - Pas de combinaison avec des connaissances a priori

## Inférence grammaticale en reconnaissance de documents

[Shilman, 2005]

Utilisateur fournit	Inférence
grammaire de pages	caractéristiques
pages annotées	paramètres

⇒ Pas d'inférence grammaticale automatique possible

Utilisation de l'expertise de l'utilisateur

# Méthode Eyes Wide Open

## Objectifs

- Combinaison statistique/structurelle
- Apprentissage sans vérité terrain
- Inférence de règles

## Proposition : méthode Eyes Wide Open (EWO)

↔ Vue exhaustive sur les données

↔ Intégration de l'utilisateur

- Minimisation des interventions

Méthode	Données	Modèle
Statistiques	Utilisateur	Automatique
Syntaxiques	∅	Utilisateur
EWO	Utilisateur	

① État de l'art

② Philosophie de notre contribution

Construction progressive d'un système  
Présentation générale de la méthode EWO

③ Présentation détaillée de la méthode EWO

④ Validation

⑤ Conclusion

# Construction progressive d'un système

- Premier niveau de la grammaire connu
  - éléments à reconnaître

# Construction progressive d'un système

- Premier niveau de la grammaire connu  
→ éléments à reconnaître

```
courrier ::=  
  coordonneesExpeditEUR &&  
  coordonneesDestinataire &&  
  dateLieu &&  
  objet &&  
  ouverture &&  
  corpsDeTexte &&  
  signature &&  
  ps.
```

# Construction progressive d'un système

- Premier niveau de la grammaire connu

→ éléments à reconnaître

- Construction progressive des règles

- Réutilisation d'une approche classique pour les problèmes complexes
- Décomposition en sous-problèmes

courrier ::=

coordonneesExpéditeur &&

coordonneesDestinataire &&

dateLieu &&

objet &&

ouverture &&

corpsDeTexte &&

signature &&

ps.



# Résolution d'un sous-problème

## Cas 1

- Forte connaissance a priori
- Peu de variabilité

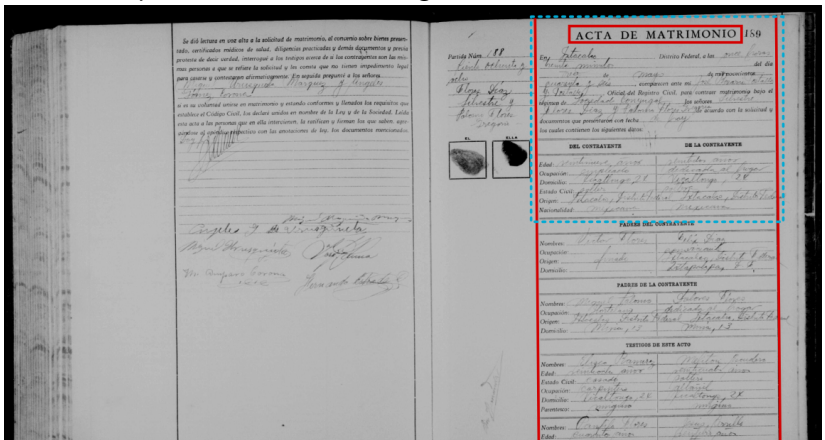
⇒ Description manuelle de la règle

# Résolution d'un sous-problème

## Cas 1

- Forte connaissance a priori
- Peu de variabilité

⇒ Description manuelle de la règle



# Résolution d'un sous-problème

## Cas 2

- Peu de connaissance a priori
- Forte variabilité

⇒ Utilisation de la méthode EWO  
Inférence semi-automatique de la règle

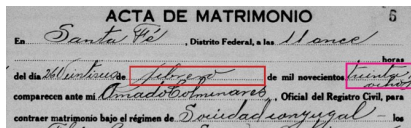
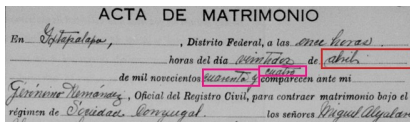
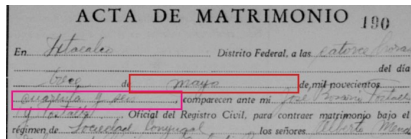
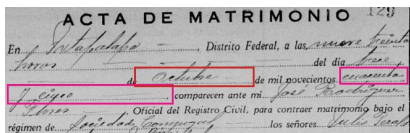
# Résolution d'un sous-problème

## Cas 2

- Peu de connaissance a priori
- Forte variabilité

⇒ Utilisation de la méthode EWO

Inférence semi-automatique de la règle



# Construction progressive d'un système

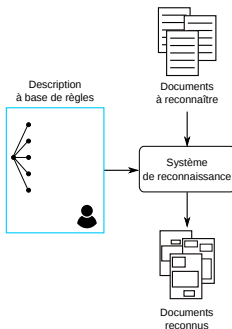
Ensemble des  
données disponibles

images de documents  
+ vérité terrain si dispo

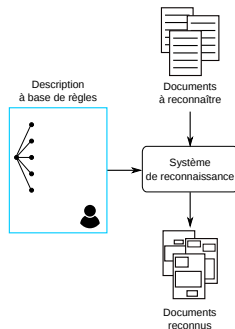
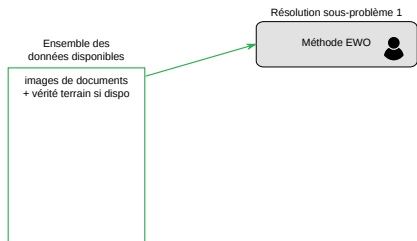
# Construction progressive d'un système

Ensemble des  
données disponibles

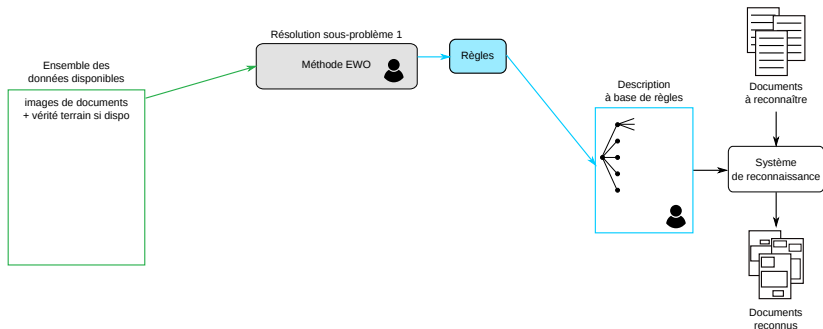
images de documents  
+ vérité terrain si dispo



# Construction progressive d'un système

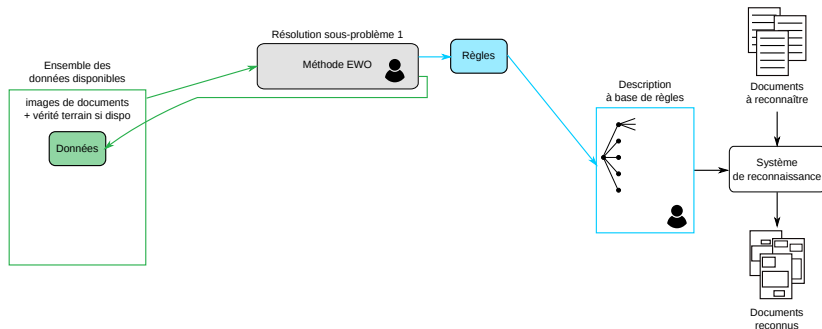


# Construction progressive d'un système

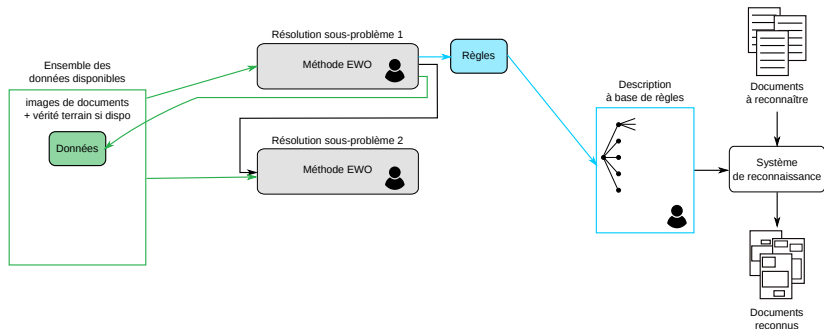




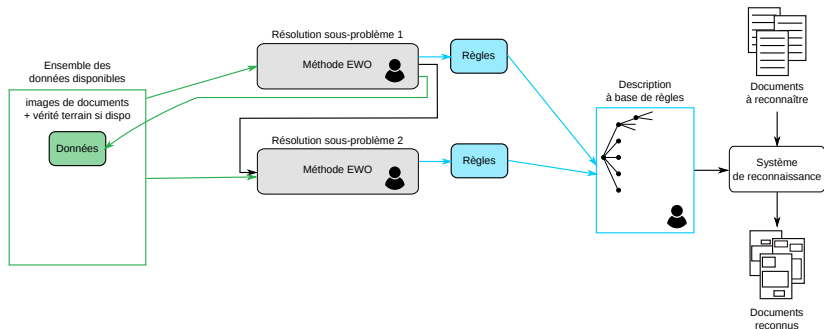
# Construction progressive d'un système



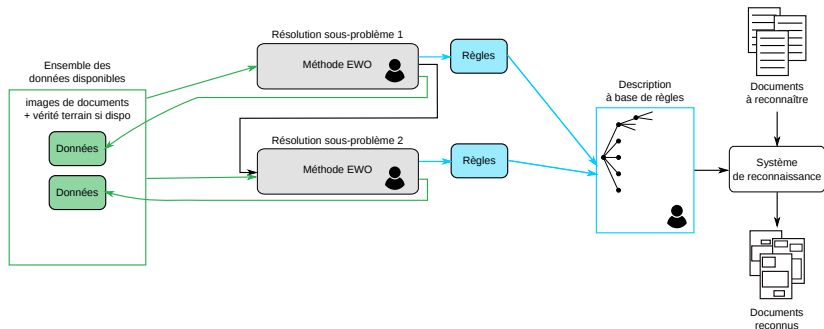
# Construction progressive d'un système



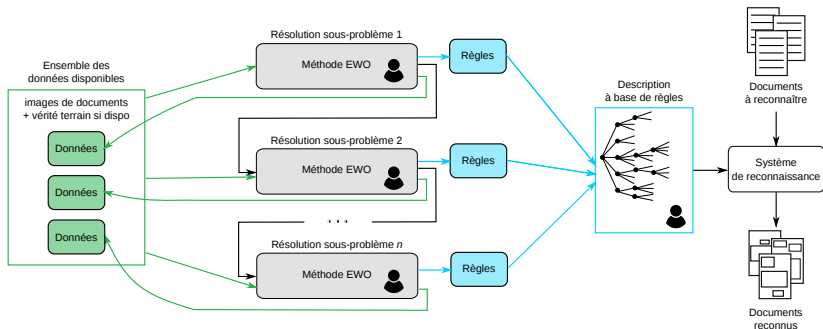
# Construction progressive d'un système



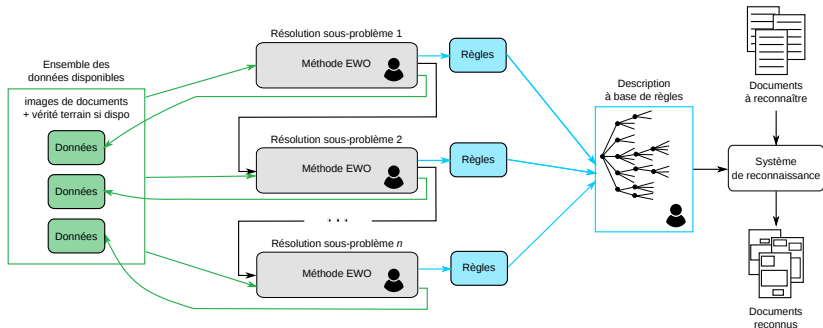
# Construction progressive d'un système



# Construction progressive d'un système



# Construction progressive d'un système



- Enrichissement progressif
  - Description grammaticale
  - Ensemble des données disponibles
- Interaction avec l'utilisateur
  - Présentation synthétique des données
    - Prise de décision efficace
  - Intégré au cœur du système

# Méthode EWO

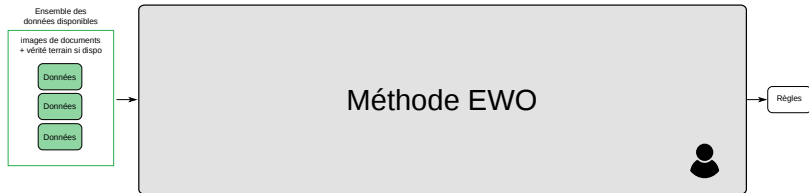
## Propriétés de la méthode

- Généricité
- Utilisable sans vérité terrain
- Intégration de connaissances a priori
- Vue exhaustive sur les données

## Moyens utilisés

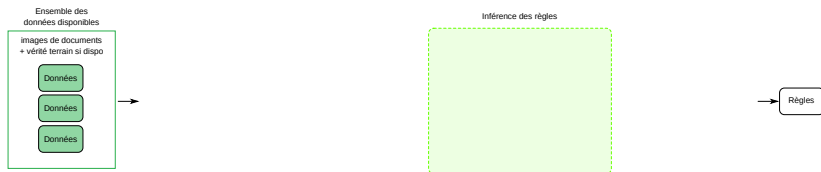
- Clustering
  - ↪ Extraction de connaissances
  - ↪ Utilisation des redondances dans les données
  - ↪ Émergence automatique de structures
- Interaction utilisateur
  - ↪ Apport de sens aux données

# Inférence semi-automatique des règles d'un sous-problème

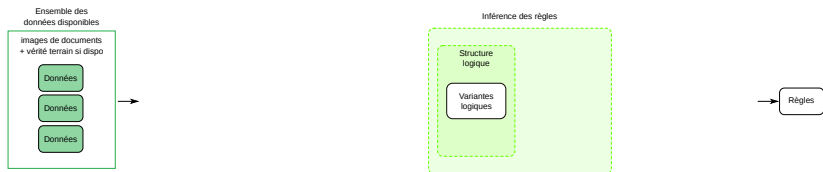




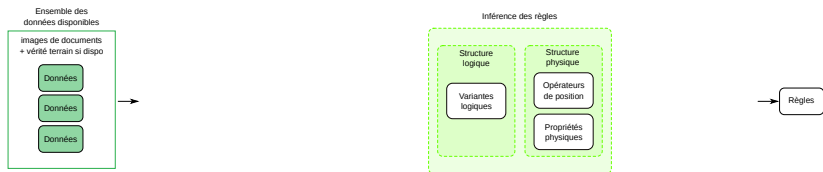
# Inférence semi-automatique des règles d'un sous-problème



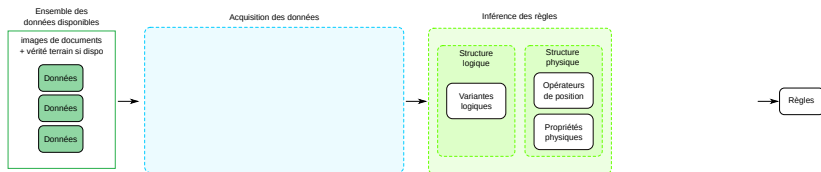
# Inférence semi-automatique des règles d'un sous-problème



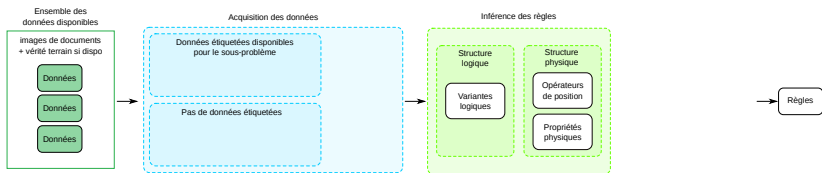
# Inférence semi-automatique des règles d'un sous-problème



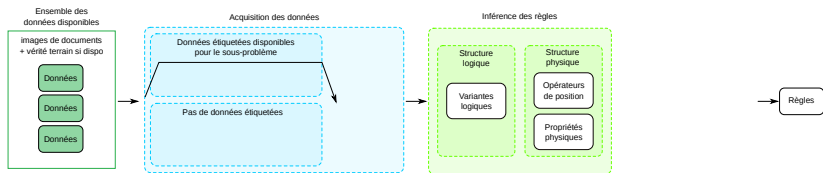
# Inférence semi-automatique des règles d'un sous-problème



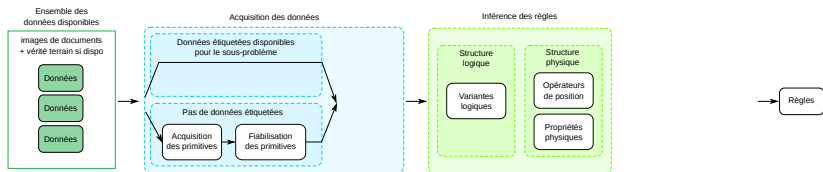
# Inférence semi-automatique des règles d'un sous-problème



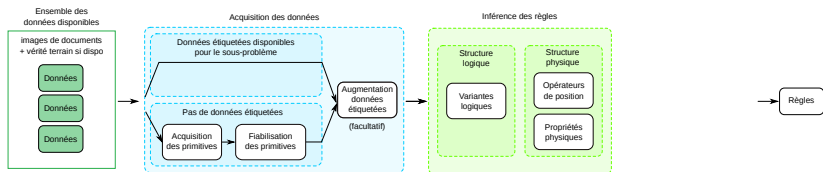
# Inférence semi-automatique des règles d'un sous-problème



# Inférence semi-automatique des règles d'un sous-problème

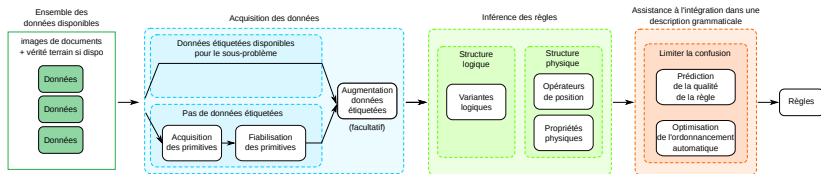


# Inférence semi-automatique des règles d'un sous-problème

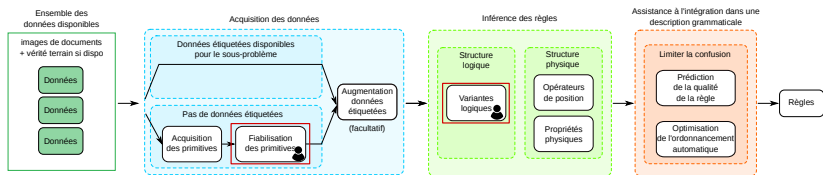




# Inférence semi-automatique des règles d'un sous-problème



# Inférence semi-automatique des règles d'un sous-problème

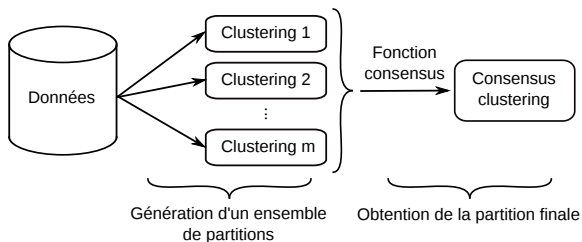


- Utilisateur intégré dans la méthode
- Vue exhaustive et synthétique
  - Prise de décision facilitée

# Clustering - quel algorithme choisir ?

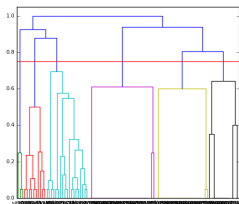
- Pas de connaissance utilisateur
  - Algorithme générique s'adaptant à n'importe quel jeu de données
  - Minimisation des paramètres à fixer par l'utilisateur
    - ↪ En particulier le nombre de clusters

Proposition : utilisation d'une méthode de clustering ensembliste



# Evidence Accumulation Clustering [Fred and Jain, 2002]

- 1 Construction de  $N$  partitions avec différents
  - Algorithmes
  - Choix de paramètres
- 2 Combinaison des différentes partitions pour générer une matrice de similarité
  - Si 2 instances souvent ensemble  $\Rightarrow$  proches
- 3 Détermination de la partition finale : algorithme de clustering hiérarchique
  - $\hookrightarrow$  *Maximum cluster lifetime* pour couper le dendrogramme



① État de l'art

② Philosophie de notre contribution

③ **Présentation détaillée de la méthode EWO**

Acquisition des données

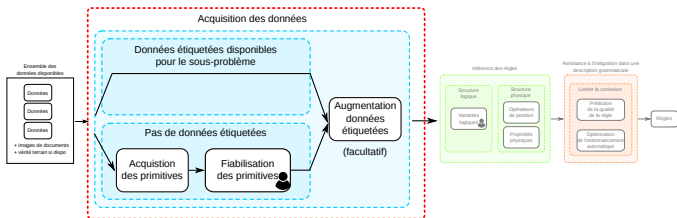
Inférence des règles

Intégration dans la description complète

④ Validation

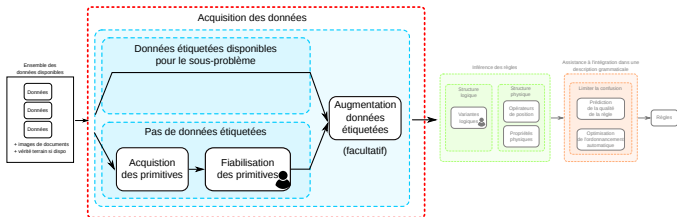
⑤ Conclusion

# Acquisition des données



- Cas simple : les données étiquetées existent
- Sinon : construction d'une pseudo vérité terrain
  - Redondances
  - Grands volumes de documents

# Acquisition des données



- Cas simple : les données étiquetées existent
- Sinon : construction d'une pseudo vérité terrain
  - Redondances
  - Grands volumes de documents
- Données utiles
  - Minimum : boîtes englobantes des éléments
  - Mais aussi : transcription, langue, ordre de lecture, etc.

Coordonnées expéditeur

*Pauline Janvier*  
*13 Jeanne d'arc*  
*57370 Darnley St Quentin*

Date, lieu

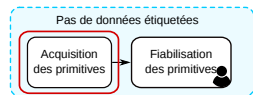
*Paris le 28 Juin 2006*

Coordonnées destinataire

*Mmees Lucie Godeau*  
*40 rue St Louis*

# Pseudo vérité terrain : acquisition des primitives

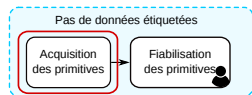
- Extraction automatique
- Sources variées : OCR, word spotting, systèmes de reconnaissance (lignes de texte, segments, etc.)
- Déterminées par l'utilisateur





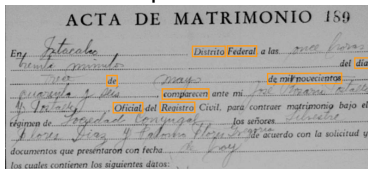
# Pseudo vérité terrain : acquisition des primitives

- Extraction automatique
- Sources variées : OCR, word spotting, systèmes de reconnaissance (lignes de texte, segments, etc.)
- Déterminées par l'utilisateur



## Extraction de la position de 8 mots-clés

- POI : descripteurs locaux sur la variation du gradient



- Autre méthode possible : OCR
  - Faibles performances
  - Mauvaise qualité + interaction manuscrit/imprimé

# Pseudo vérité terrain : extraction des primitives

Cas idéal : 8 mots clés par document

Cas réel : erreurs d'extraction

- Oublis, fausses reconnaissances

ACTA DE MATRIMONIO 189

En Itzamal Distrito Federal a las once horas del día uno de Mayo de mil novecientos veinte y seis comparecen ante mí José María Estela Oficial del Registro Civil, para contraer matrimonio bajo el régimen de Soledad Conjugal los señores Sebastián Pérez Díaz y Fabiana Pérez de acuerdo con la solicitud y documentos que presentaron con fecha de mayo los cuales contienen los siguientes datos:

ACTA DE MATRIMONIO

En Veracruz Distrito Federal, a las once horas del día uno de enero de mil novecientos veintinueve y seis comparecen ante mí Pedro Ruiz Oficial del Registro Civil, para contraer matrimonio bajo el régimen de Soledad Conjugal los señores Francisco Aguilar Prudencio Espinoza Benito Justo de acuerdo con la solicitud y documentos que presentaron con fecha de hoy los cuales contienen los siguientes datos:

⇒ Données bruitées

# Pseudo vérité terrain : fiabilisation des primitives

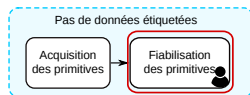
But : supprimer le bruit à moindre coût

## ① Clustering des occurrences

- EAC clustering
- Étape automatique

## ② Suppression des occurrences par cluster

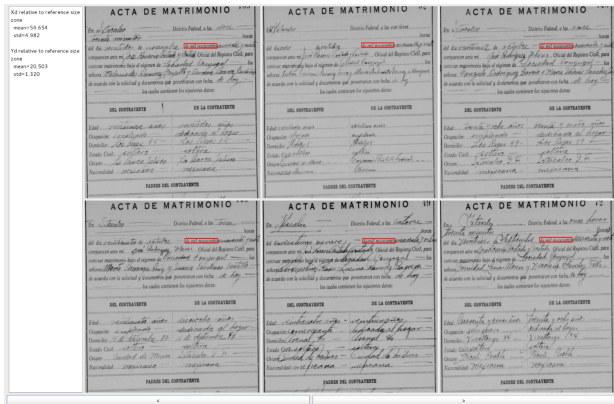
- Étape interactive
- Visualisation de quelques exemples par cluster
  - Sélectionnés automatiquement
  - Proches du centroïde
- Décision de supprimer ou conserver le cluster par l'utilisateur
  - Basée sur les exemples
  - Indicateurs de dispersion pour aider l'utilisateur



Résultat : pseudo vérité terrain construite semi-automatiquement

# Pseudo vérité terrain : fiabilisation des primitives

- Exemple : mot clé « de mil novecientos »
- Visualisation des clusters par l'utilisateur
  - ↳ quelques exemples représentatifs



1192 occurrences

# Pseudo vérité terrain : fiabilisation des primitives

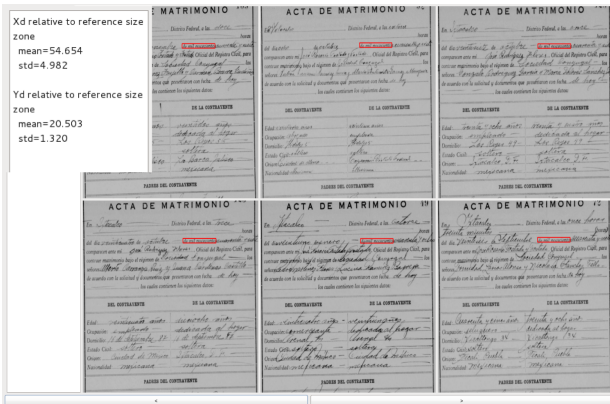
- Exemple : mot clé « de mil novecientos »
- Visualisation des clusters par l'utilisateur
  - ↳ quelques exemples représentatifs

The image displays six examples of Spanish marriage certificates (ACTA DE MATRIMONIO) arranged in a 2x3 grid. Each document is a scanned form with handwritten entries. Red rectangular boxes highlight the phrase "de mil novecientos" (of one thousand nine hundred) in various locations across the forms, such as in the date field or the names of the officiating authorities. The forms include fields for the date and location of the ceremony, the names and details of the bride and groom, and the names of the witnesses and the officiant. The text is in Spanish, and the handwriting is in cursive.

1192 occurrences

# Pseudo vérité terrain : fiabilisation des primitives

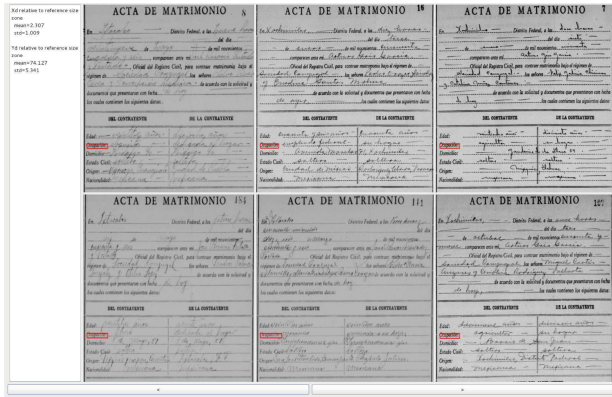
- Exemple : mot clé « de mil novecientos »
- Visualisation des clusters par l'utilisateur
  - ↳ quelques exemples représentatifs



1192 occurrences → décision : cluster conservé

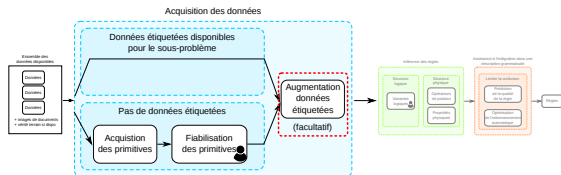
# Pseudo vérité terrain : fiabilisation des primitives

- Exemple : mot clé « de mil novecientos »
- Visualisation des clusters par l'utilisateur
  - ↳ quelques exemples représentatifs



2114 occurrences → décision : cluster rejeté (mot clé ocupación)

# Augmentation de la vérité terrain

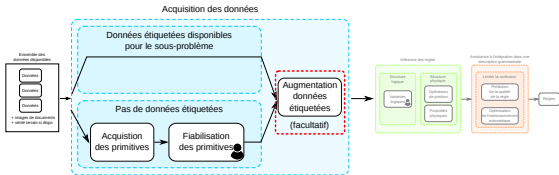


## Ajout d'éléments

- Nécessaire à la reconnaissance
- Automatiquement
  - À partir de systèmes de reconnaissances existants
  - Utilisation de la position des éléments



# Augmentation de la vérité terrain



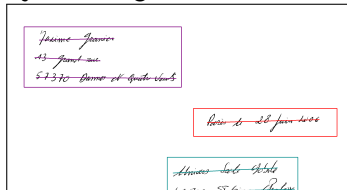
## Ajout d'éléments

- Nécessaire à la reconnaissance
- Automatiquement
  - À partir de systèmes de reconnaissances existants
  - Utilisation de la position des éléments

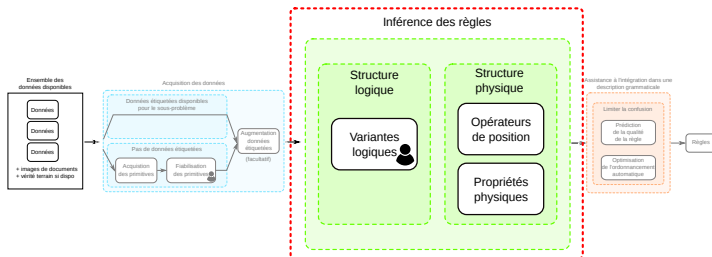
## Ajout des composantes connexes



## Ajout des lignes de texte



# Inférence des règles



## Grammaire de départ courrier manuscrit courrier ::=

coordonneesExpéditeur &&  
 coordonneesDestinataire &&  
 dateLieu &&  
 objet &&  
 ouverture &&  
 corpsDeTexte &&  
 signature &&  
 ps.

# Inférence des règles - structure logique

## Structure logique

- Apprentissage de la hiérarchie de la grammaire
- Grande variabilité possible par élément
- Nécessite une vision exhaustive des données

# Inférence des règles - structure logique

## Structure logique

- Apprentissage de la hiérarchie de la grammaire
- Grande variabilité possible par élément
- Nécessite une vision exhaustive des données

## Apprentissage des variantes logiques des éléments

- Détection automatique des variantes
  - Clustering
- Apport de sens aux données
  - Interaction utilisateur
  - Visualisation de quelques exemples
  - Validation des variantes détectées
  - Attribution d'un nom significatif à chaque variante
    - Description grammaticale interprétable

# Inférence des règles - structure logique

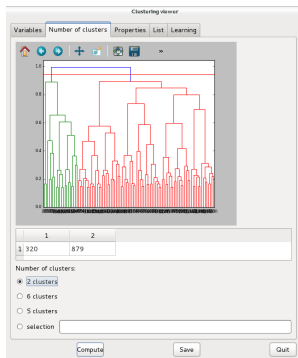
## Coordonnées expéditeur

- Recherche des variantes logiques
- Variables : hauteur et largeur
- Partition en deux clusters proposée par EWO

# Inférence des règles - structure logique

## Coordonnées expéditeur

- Recherche des variantes logiques
- Variables : hauteur et largeur
- Partition en deux clusters proposée par EWO



# Inférence des règles - structure logique

## Coordonnées expéditeur

- Recherche des variantes logiques
- Variables : hauteur et largeur
- Partition en deux clusters proposée par EWO

Cluster viewer

Cluster label:   Delete cluster occurrences

<p>height relative to page size mean=4.003 std=1.114</p> <p>width relative to page size mean=34.459 std=9.808</p> <p>TEBEC Société 22 rue Dahan 52 rue GRAND 76 1012 91 16 74</p> <p>EBI Prouvost 20 Boulevard 75 102 188 1</p> <p>Chang le Pêcheur Snc</p> <p>Objet: répartition Tribunal de Commerce</p> <p>Monsieur, Monsieur, Je vous adresse les 1/2 parts de la société en partiel de biens, y ont été des par vos copies par je soude et le soude au titre de la répartition.</p> <p>Je suis en ce document les parts en soude et par un soude par soude.</p> <p>Voilà, vous, Monsieur, Monsieur, une répartition définitive.</p> <p>TEBEC Société</p>	<p>Administration (société) 107 chemin de la Justice 76000 SAINT-DENIS Tel : 02 34 52 32 14</p> <p>Robert &amp; Associés SAF Marché City of Paris 92000 VANVRES</p> <p>Objet: répartition Tribunal de Commerce</p> <p>Monsieur, Monsieur, Je vous adresse les 1/2 parts de la société en partiel de biens, y ont été des par vos copies par je soude et le soude au titre de la répartition.</p> <p>Je suis en ce document les parts en soude et par un soude par soude.</p> <p>Voilà, vous, Monsieur, Monsieur, une répartition définitive.</p> <p>TEBEC Société</p>	<p>Objet: répartition Tribunal de Commerce</p> <p>Monsieur, Monsieur, Je vous adresse les 1/2 parts de la société en partiel de biens, y ont été des par vos copies par je soude et le soude au titre de la répartition.</p> <p>Je suis en ce document les parts en soude et par un soude par soude.</p> <p>Voilà, vous, Monsieur, Monsieur, une répartition définitive.</p> <p>TEBEC Société</p>
<p>POURQUENOT Noun 17 rue de la République 97000 Pointe-à-Pitre Tel: 01 98 42 92 32</p> <p>Objet: répartition Tribunal de Commerce</p> <p>Monsieur, Monsieur, Suite à son changement d'adresse, je vous ai envoyé un courrier. Je suis heureux que vous ayez pu recevoir le courrier. Je vous prie de m'excuser pour le retard de l'envoi. Je vous prie de m'excuser pour le retard de l'envoi. Je vous prie de m'excuser pour le retard de l'envoi.</p> <p>Voilà, vous, Monsieur, Monsieur, une répartition définitive.</p> <p>TEBEC Société</p>	<p>Monsieur LEBLANC 10 rue de la République 97000 Pointe-à-Pitre Tel: 01 98 42 92 32</p> <p>Objet: répartition Tribunal de Commerce</p> <p>Monsieur, Monsieur, Suite à son changement d'adresse, je vous ai envoyé un courrier. Je suis heureux que vous ayez pu recevoir le courrier. Je vous prie de m'excuser pour le retard de l'envoi. Je vous prie de m'excuser pour le retard de l'envoi. Je vous prie de m'excuser pour le retard de l'envoi.</p> <p>Voilà, vous, Monsieur, Monsieur, une répartition définitive.</p> <p>TEBEC Société</p>	<p>Objet: répartition Tribunal de Commerce</p> <p>Monsieur, Monsieur, Suite à son changement d'adresse, je vous ai envoyé un courrier. Je suis heureux que vous ayez pu recevoir le courrier. Je vous prie de m'excuser pour le retard de l'envoi. Je vous prie de m'excuser pour le retard de l'envoi. Je vous prie de m'excuser pour le retard de l'envoi.</p> <p>Voilà, vous, Monsieur, Monsieur, une répartition définitive.</p> <p>TEBEC Société</p>

# Inférence des règles - structure logique

## Coordonnées expéditeur

- Recherche des variantes logiques
- Variables : hauteur et largeur
- Partition en deux clusters proposée par EWO

Cluster viewer

Cluster label:   Delete cluster occurrences

height relative to page size  
mean=4.003  
std=1.114

width relative to page size  
mean=34.459  
std=9.808

<p><b>TELECOM Junifer</b> 22 rue Grandin 62 500 CHAMBY TEL: 03 32 47 34 59</p> <p><b>NWA Arrosses</b> 26 Boulevard Victor 75 251 PARIS</p> <p>Cherry le Penche 2008</p> <p>Chart: responsable civik <b>reference client: DIVERSE</b></p> <p>Maitres, Messieurs, Y'ai récemment lu le petit article d'un magazine en parlant de travaux</p>	<p>Madame, Monsieur,</p> <p>Je suis ravi de constater que vous avez accepté de participer à notre étude.</p> <p>Vous recevrez prochainement un questionnaire par la poste.</p> <p>Si vous avez des questions, n'hésitez pas à nous contacter.</p> <p>Bonne nuit.</p> <p>Dr. Jean-Luc L...</p>
<p>POURQUOI? Parce qu'il n'y a pas de solution pour résoudre ce problème.</p> <p><b>Madame, Monsieur,</b></p> <p>Tout d'abord, je tiens à vous remercier pour votre participation à notre étude.</p> <p>Vous recevrez prochainement un questionnaire par la poste.</p> <p>Si vous avez des questions, n'hésitez pas à nous contacter.</p> <p>Bonne nuit.</p> <p>Dr. Jean-Luc L...</p>	<p>Monsieur L...</p> <p>Je suis ravi de constater que vous avez accepté de participer à notre étude.</p> <p>Vous recevrez prochainement un questionnaire par la poste.</p> <p>Si vous avez des questions, n'hésitez pas à nous contacter.</p> <p>Bonne nuit.</p> <p>Dr. Jean-Luc L...</p>
<p>Cher Monsieur,</p> <p>Je suis ravi de constater que vous avez accepté de participer à notre étude.</p> <p>Vous recevrez prochainement un questionnaire par la poste.</p> <p>Si vous avez des questions, n'hésitez pas à nous contacter.</p> <p>Bonne nuit.</p> <p>Dr. Jean-Luc L...</p>	<p>Monsieur L...</p> <p>Je suis ravi de constater que vous avez accepté de participer à notre étude.</p> <p>Vous recevrez prochainement un questionnaire par la poste.</p> <p>Si vous avez des questions, n'hésitez pas à nous contacter.</p> <p>Bonne nuit.</p> <p>Dr. Jean-Luc L...</p>



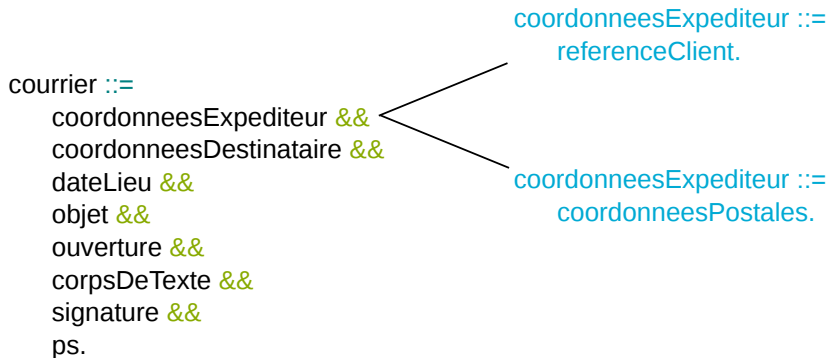






## Inférence des règles - structure logique

- Détection automatique de deux variantes logiques
- Production de nouvelles règles



# Détection des valeurs extrêmes

- Cas limite de ce qui doit être reconnu
- Difficile à détecter par un humain
- Détectées automatiquement et présentées à l'utilisateur

## ① Cas rares

- Définition de variantes supplémentaires

## ② Erreurs dans les données étiquetées

- Suppression ou correction

⇒ Nécessaire pour la robustesse de l'analyse

# Détection des valeurs extrêmes

- Cas limite de ce qui doit être reconnu
- Difficile à détecter par un humain
- Détectées automatiquement et présentées à l'utilisateur

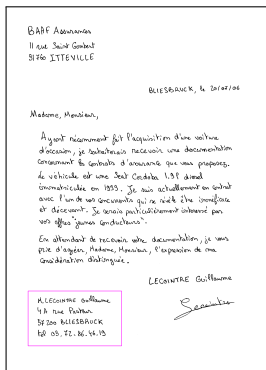
## ① Cas rares

- Définition de variantes supplémentaires

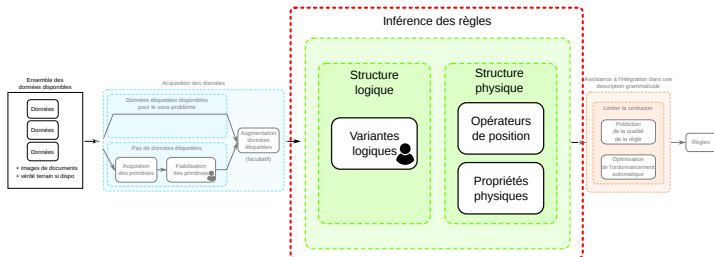
## ② Erreurs dans les données étiquetées

- Suppression ou correction

⇒ Nécessaire pour la robustesse de l'analyse



# Inférence des règles - structure physique



Comment le système peut reconnaître les éléments de la structure logique ?

- Positionnement des éléments
- Propriétés physiques des éléments

Nécessite une vue exhaustive sur les données

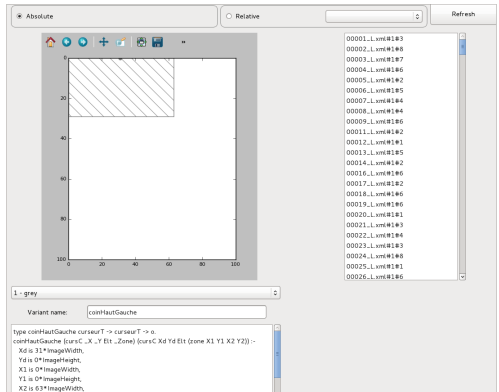
# Positionnement des éléments

- Détermine la zone de recherche d'un élément
  - Gestion de la combinatoire
  - Maximisation du rappel



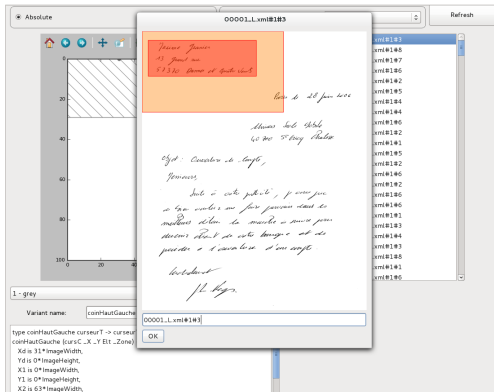
# Positionnement des éléments

- Détermine la zone de recherche d'un élément
  - Gestion de la combinatoire
  - Maximisation du rappel
- Choix utilisateur
  - positionnement absolu ou relatif
- Détection automatique des variantes
- Frontières de la zone de recherche
- Ordre de parcours des éléments
- Nom choisi par l'utilisateur



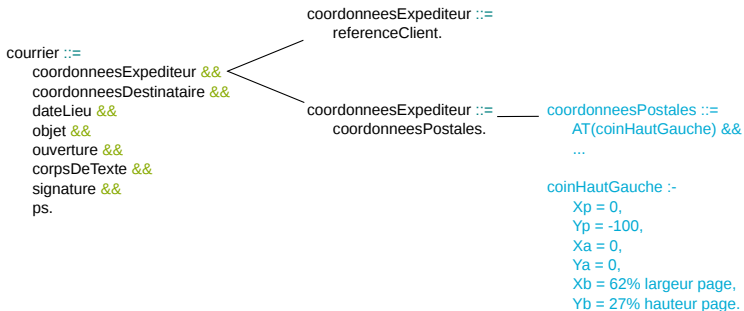
# Positionnement des éléments

- Détermine la zone de recherche d'un élément
  - Gestion de la combinatoire
  - Maximisation du rappel
- Choix utilisateur
  - positionnement absolu ou relatif
- Détection automatique des variantes
- Frontières de la zone de recherche
- Ordre de parcours des éléments
- Nom choisi par l'utilisateur



# Positionnement des éléments

## Production de nouvelles règles



# Propriétés physiques

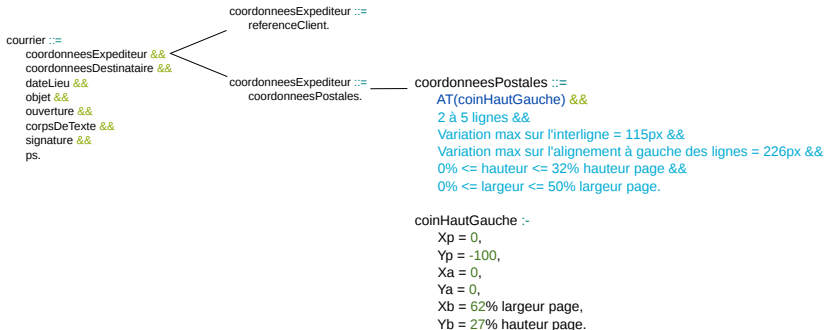
Détermination des caractéristiques propres de l'élément

- Nécessite une vision globale sur les données
- Permet la segmentation
- Statistiques descriptives

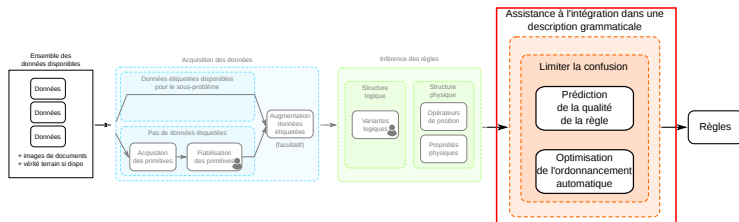
# Propriétés physiques

## Détermination des caractéristiques propres de l'élément

- Nécessite une vision globale sur les données
- Permet la segmentation
- Statistiques descriptives



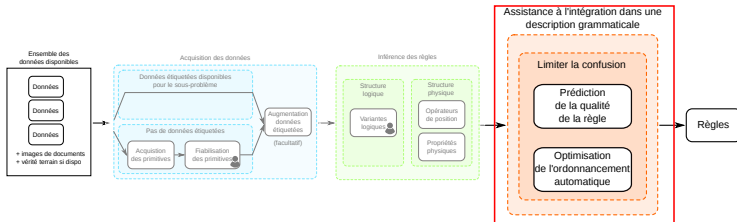
# Intégration de la règle inférée dans la description complète



But : minimiser la confusion entre les éléments

- Vérifier que la règle est suffisamment spécifiée
- Ordonner les règles
  - Recherche des éléments les plus fiables en premiers

# Intégration de la règle inférée dans la description complète



But : minimiser la confusion entre les éléments

- Vérifier que la règle est suffisamment spécifiée
- Ordonner les règles
  - Recherche des éléments les plus fiables en premiers

Comment ? Fonctionnement sous forme de requêtes au sein de la méthode EWO

- Approximation du système de reconnaissance
- Réduction du temps de calcul par rapport au système complet
- Action en temps réel de l'utilisateur avec le système

# Intégration dans la description complète

## Grammaire finale, inférence de

- 34 règles et variantes
- 14 opérateurs de position
- 60 propriétés physiques

## Grammaire de base

courrier ::=

coordonneesExpediteur &&  
 coordonneesDestinataire &&  
 dateLieu &&  
 objet &&  
 ouverture &&  
 corpsDeTexte &&  
 signature &&  
 ps.

## Grammaire finale - ordre optimisé

courrier ::=

ouverture &&  
 corpsDeTexte &&  
 signature &&  
 coordonneesPostales &&  
 coordonneesDestinataire &&  
 ps &&  
 dateLieu &&  
 objet &&  
 referenceClient.



① État de l'art

② Philosophie de notre contribution

③ Présentation détaillée de la méthode EWO

④ **Validation**

Corpus homogène : courriers manuscrits

Corpus hétérogène : MAURDOR

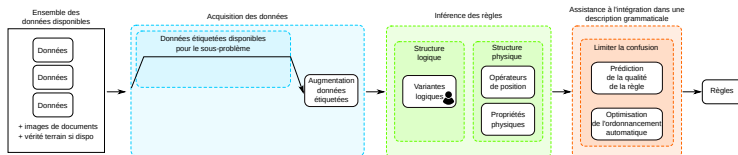
Corpus sans vérité terrain : mariages mexicains

⑤ Conclusion

# Introduction

- Validation complète
  - Différentes parties séparément
  - Fonctionnement global
- Corpus variés
  - RIMES : corpus de courriers manuscrits avec vérité terrain
  - MAURDOR : corpus hétérogène avec vérité terrain
  - FamilySearch HIP 2013 : registres de mariages mexicains sans vérité terrain
- Système syntaxique utilisé : DMOS [Coüasnon, 2006]

# RIMES



- Grammaire complète décrite avec EWO
- Métrique du concours
  - Taux de pixels noirs mal étiquetés
- Données
  - Apprentissage : 900 courriers
  - Test : base concours de 100 courriers
- Inférence
  - 34 règles et variantes
  - 14 opérateurs de position
  - 60 paramètres pour les propriétés physiques

# RIMES - Résultats

Système	Taux d'erreur
Syntaxique (DMOS) [Lemaitre, 2008]	8,97
Syntaxique (DMOS) localement stochastique [Maroneze, 2011]	5,53
Statistique basé sur un MRF [Lemaitre, 2007]	8,53
Statistique basé sur des CRF [Montreuil, 2010]	6,33
Méthode EWO	5,82

## Résultats comparables

- Meilleur système syntaxique
  - Gain de temps avec la méthode EWO
- Meilleur système statistique
  - Meilleure gestion des cas rares avec EWO
  - Homogénéité dans les blocs construits

# MAURDOR

## Tâche 5 de la compétition sur zones segmentées

- Étiquetage logique
- Ordre de lecture
- Groupe

Summit Caroline  
5 rue des Allées  
93300 Montfermeil  
03-33-55-51-60

coordinates

date + location  
À Montfermeil le  
30/04/2012

text\_section

Madame, Monsieur,

Samredi 28 juillet 2012, j'ai prouvé  
un comantelage qui a été encrei deux  
automobiles de particuliers de première  
immatriculée AB-234-CD, est changi sa  
numéro sans aucune autorisation  
AA-555-AA, et est changi ses deux fois  
avant je voudrais savoir si je puis bénéficier  
de l'aide de mon assurance automobile Vaia  
pour rembourser les réparations nécessaires.

Dans leattente et sans réponse de votre part  
je vous prie, madame, monsieur de lui rendre  
mon expression de mon plus profond respect

**Employeur**

Entreprise (raison sociale) : Sabpro

Nom du responsable de l'entreprise : Sarnier Stéphanie

Adresse (numéro et nom) : Bois du Servan

Commune de l'entreprise : POISSASSIERE

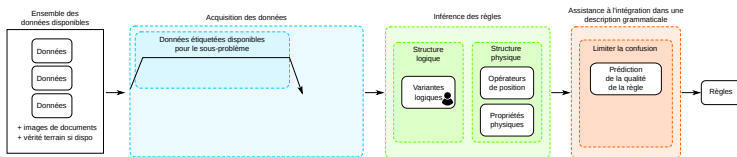
Sera présent(e) :

le matin	<input checked="" type="checkbox"/>	<input type="checkbox"/>
au déjeuner	<input type="checkbox"/>	<input checked="" type="checkbox"/>
l'après-midi	<input checked="" type="checkbox"/>	<input type="checkbox"/>

## Documents variés

- Factures, formulaires, courriers, tableaux, plans, etc.
- Langues : français, anglais et arabe
- Manuscrit et dactylographié

# MAURDOR - Résultats



- Données
  - Apprentissage : 2000 documents
  - Test : 1000 documents
- Deux systèmes
  - Syntaxique avec la méthode EWO
  - Second participant :
    - type : classifieur (SVM)
    - ordre + groupe : règles décrites manuellement

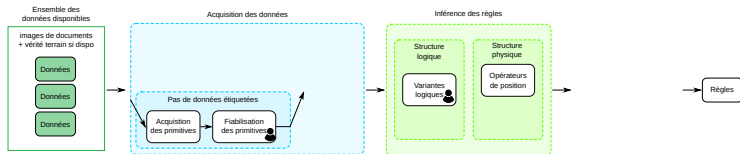
# MAURDOR - Résultats

- Métrique du concours
  - Score  $\in [-100; 100]$
  - Score  $> 0 \Leftrightarrow$  apport de plus d'informations que d'erreurs

Système	Type	Ordre	Groupe
Méthode EWO	55	45	60
Participant 2	69	28	61

- Type : zones déjà segmentées  $\Rightarrow$  tâche de classification
  - Meilleur résultat du SVM
- Ordre
  - Meilleurs résultats pour les règles définies avec EWO
  - Validation de l'apport de la vue exhaustive sur les données
- Détection de 164 erreurs dans la vérité terrain

# FamilySearch HIP2013

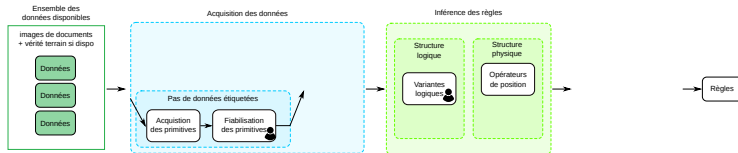


Jeu de données : registres de mariages de la compétition HIP2013  
FamilySearch

- Apprentissage : 7000 documents
  - Sans vérité terrain
- Test : 2000 documents annotés manuellement



# FamilySearch HIP2013



Jeu de données : registres de mariages de la compétition HIP2013  
FamilySearch

- Apprentissage : 7000 documents
  - Sans vérité terrain
- Test : 2000 documents annotés manuellement

Évaluation de la correspondance spatiale de zones

- Métrique introduite par Garris (Garris, 1995)
  - Recouvrement entre zone attendue et zone obtenue
- Reconnaissance complète
  - Au moins 95% de la largeur et au moins 75% de la hauteur
- Reconnaissance partielle
  - 80 à 95% de la largeur et au moins 75% de la hauteur

# FamilySearch HIP2013 - Constitution de la pseudo vérité terrain

- 7000 documents
- 54 141 zones
- 276 actions utilisateur
  - Annotation manuelle (8 actions par document)
    - 34,5 documents annotés
  - Coût divisé par 200



# FamilySearch HIP2013 - Résultats

Deux systèmes évalués :

- 4 modèles définis manuellement [Lemaitre 2013]
- 11 modèles inférés automatiquement

Modèle	1	2	3	4	5	6	7	8	9	10	11
Nombre	1448	822	740	652	566	470	359	123	92	33	25

# FamilySearch HIP2013 - Résultats

Deux systèmes évalués :

- 4 modèles définis manuellement [Lemaitre 2013]
- 11 modèles inférés automatiquement

Modèle	1	2	3	4	5	6	7	8	9	10	11
Nombre	1448	822	740	652	566	470	359	123	92	33	25

		Modèles	
		Inférés	Manuels
Zone	Reconnaissance complète	91.4%	89.7%
	Reconnaissance partielle	6.2%	4.0%
	Manquant	2.4%	6.3%
Taux de reconnaissance du document		89.8%	78.9%

# FamilySearch HIP2013 - Résultats

Deux systèmes évalués :

- 4 modèles définis manuellement [Lemaitre 2013]
- 11 modèles inférés automatiquement

Modèle	1	2	3	4	5	6	7	8	9	10	11
Nombre	1448	822	740	652	566	470	359	123	92	33	25

		Modèles	
		Inférés	Manuels
Zone	Reconnaissance complète	91.4%	89.7%
	Reconnaissance partielle	6.2%	4.0%
	Manquant	2.4%	6.3%
Taux de reconnaissance du document		89.8%	78.9%

- Amélioration des taux de reconnaissance
  - Inférence semi-automatique des modèles
  - Meilleure couverture du corpus
    - Construction rapide d'une pseudo vérité terrain

- ① État de l'art
- ② Philosophie de notre contribution
- ③ Présentation détaillée de la méthode EWO
- ④ Validation
- ⑤ Conclusion**

# Conclusion

Méthode facilitant l'adaptation d'un système de reconnaissance de structure combinant

- Apprentissage statistique
- Système syntaxique
  - Pouvoir d'expression et gestion des cas rares
- Interaction avec l'utilisateur
  - Au cœur du système

Méthode EWO permet


- Inférence des règles
  - Structures logique et physique
- Gestion de l'absence de vérité terrain
- Vision exhaustive et synthétique des données

# Perspectives

- Extension des cadres applicatifs



# Perspectives

- Extension des cadres applicatifs
    - Séparateurs d'articles d'archives de presse
    - Grammaire déjà existante
- 
- Grande variabilité
  - Création d'une pseudo vérité terrain
  - Utilisation d'éléments horizontaux combinant
    - OCR
    - Lignes de texte détectées avec une grammaire
    - Séparateurs déjà détectés

# Perspectives

- Utilisation dans d'autres contextes
  - Autres méthodes de reconnaissance de structures
  - Construction rapide de bases étiquetées

# Perspectives

- Utilisation dans d'autres contextes
  - Autres méthodes de reconnaissance de structures
  - Construction rapide de bases étiquetées
- Problématique de l'exhaustivité
  - Exploitation de très grandes quantités de données
  - Comportement de l'EAC clustering
  - Détection des cas rares
  - Minimisation des interventions utilisateur
    - Sélection des primitives
    - Sélection automatique des propriétés physiques