

Un exemple de convergence

Apprentissage / Représentation des connaissances:

Extraction automatique de relations par ontologies et Programmation Logique Inductive

Bernard Espinasse¹ - Rinaldo Lima²
- Fred Freitas²

¹ LIS UMR CNRS 7020, Aix Marseille Université, Marseille (AMU), France
bernard.espinasse@lis-lab.fr

² Federal University of Pernambuco (UFPE), Recife, Brazil
rinaldo.jose@ufrpe.br

APSEM-Nov-2018

L'extraction d'information (EI)

Extraction Information (EI) composée de 2 tâches principales:

- **La reconnaissance d'entités nommées (REN):** extraire des instances d'entités nommées, Ex. des noms de personnes, de lieux;
- **L'extraction de relations (RE):** extraire des relations entre ces entités nommées

Soit la phrase:

“American saxophonist David Murray recruited Amidu Berry”

- **Extraction des entités nommées :**

“David Murray” et “Amidu Berry”

- **Extraction des relations :**

CITIZEN(David Murray, American) et

HIRE(David Murray, Amidu Berry)

EI et ressources sémantiques: OBIE

- Pour être plus **précis**, les systèmes d'IE doivent exploiter plus de **ressources sémantiques** (Nédellec & Nazarenko, 2005).
- Emergence de l'EI basée sur des **ontologies** - ***Ontology-Based Information Extraction – OBIE*** (Wimalasuriya et Dou, 2010) :
 - ***Ontologie en entrée*** : processus d'extraction guidé par une ontologie avec une annotation sémantique des textes à traiter
 - ***Ontologie en sortie***: utilisation d'une ontologie pour représenter et stocker les informations extraites par peuplement d'une ontologie
- L'OBIE permet aussi :
 - D'exploiter un **traitement du langage naturel en profondeur**
 - De **générer automatiquement des contenus sémantiques** pour le Web sémantique (Wu and Weld, 2008)

El et apprentissage automatique

- Pour être plus **rapidement développés** et **adaptables** à d'autres domaines d'application, les systèmes d'El utilisent des **techniques d'apprentissage automatique**:
- ***Apprentissage supervisé statistique*** largement utilisé :
 - **REN**: très bonne performance, autour de 90%,
 - **ER**: performance très nettement inférieure (Giuliano et al., 2007) (Bach et Badaskar, 2007), et peu de progrès réalisé depuis un certain temps.
 - Problème de l'explicabilité
- ***Apprentissage profond*** travaux en cours, pas encore performant
 - Usage de Word Embedding – plongement lexical)
 - Problème de l'explicabilité et l'interprétabilité
- ***Apprentissage supervisé symbolique***, avec une de ses techniques: ***la Programmation Logique Inductive (PLI)***

Programmation Logique Inductive (PLI)

- La Programmation Logique Inductive (PLI-*Inductive Logic programming*) est une **technique d'apprentissage symbolique** (Muggleton, 1991)(Lavrack&Dzeroski, 1994) - **Relational & Logical Learning**
- En **apprentissage supervisé**, la PLI utilise **les clauses du premier ordre** pour obtenir une **représentation expressive uniforme** des exemples, de la base de connaissances et des hypothèses (règles)
- Cette **représentation** :
 - Est **plus expressive** que la représentation **attribut-valeur** (*propositionnelle*) des méthodes **d'apprentissage statistiques**
 - Permet un **traitement de la langue naturelle plus profond**
 - Permet une **intégration facile** et **naturelle** de **connaissances de domaine** (ontologie, thésaurus, ...) au processus d'apprentissage
- La PLI est très **proche des ontologies et du Web sémantique**

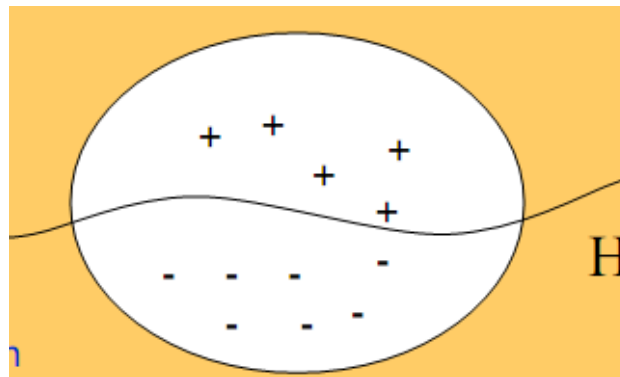
Apprentissage supervisé d'un classifieur basé sur la PLI

Entrée :

- Base de connaissance ou “Background Knowledge” (BK)
- Un ensemble E d'exemples positifs et négatifs d'apprentissage
- Une étiquette de classification c pour chaque exemple d'apprentissage

Sortie :

- Une théorie logique ou **Hypothèse H** séparant les exemples positifs des exemples négatifs



Apprentissage d'une règle en PLI

BK

Intentional:

- `parent(X, Y) :- father(X, Y).`
- `parent(X, Y) :- mother(X, Y).`

Extensional:

- `father(pat, ann).`
- `father(tom, sue).`
- `female(ann).`
- `female(eve).`
- `female(sue).`
- `male(pat).`
- `male(tom).`
- `mother(eve, sue).`
- `mother(ann, tom).`

Entrée

Examples:

Positive:

- `daughter(sue, eve).`
- `daughter(ann, pat).`

Negative:

- `daughter(tom, ann).`
- `daughter(eve, ann).`

Sortie

`daughter(D, P) :- parent(P, D), female(D).`

Le prédicat **daughter** est induit à partir des exemples positifs et négatifs, ainsi que des prédicats **parent** et **female** déclarés dans la BK

Travaux reliés

La PLI est déjà utilisé en IE, principalement en RE:

- **(Seneviratne & Ranasinghe, 2011)**: extraction d'une seule relation (located_in) sur un petit corpus de 13 pages de Wikipédia sur les oiseaux.
- **(Smole et al., 2011)**: apprentissage de règles pour extraire des informations à partir des définitions d'entités géographiques dans le texte (en langue slovène) pour l'extraction des 5 relations les plus fréquentes en 1308 définitions d'entités spatiales
- **(Kordjamshidi et al., 2012)**: « Spatial Role Labeling – SpRL », en combinant klog (P. Frasconi et al., 2014) environnement d'apprentissage relationnel basé sur les noyaux et un classificateur SVM.

Mais les corpus traités, le nombre de relations extraites et les performances sont limités.

Le système OntoILPER

OntoILPER permet l'extraction **d'instances d'entités nommées (REN)** et de **relations binaires (RE)** de textes en **anglais**

OntoILPER repose sur :

- Un **modèle relationnel des phrases** basé sur un **graphe des dépendances** (Marneffe&Manning, 2008), traitées comme des **prédicats logiques**
- Un **processus d'apprentissage basé sur la PLI**, induisant des **règles d'extractions symboliques**
- Une **ontologie de domaine** :
 - **En entrée**: permet de **choisir les concepts** qui doivent être peuplés
 - **En sortie** : est peuplée par les instances extraites
- Une **ontologie d'annotation** :
 - Permet de stocker les annotations
 - Utilisée pour appliquer les règles d'extraction

Architecture d'OntoILPER

- **2 phases:**

- **PH1-Phase d'apprentissage:**

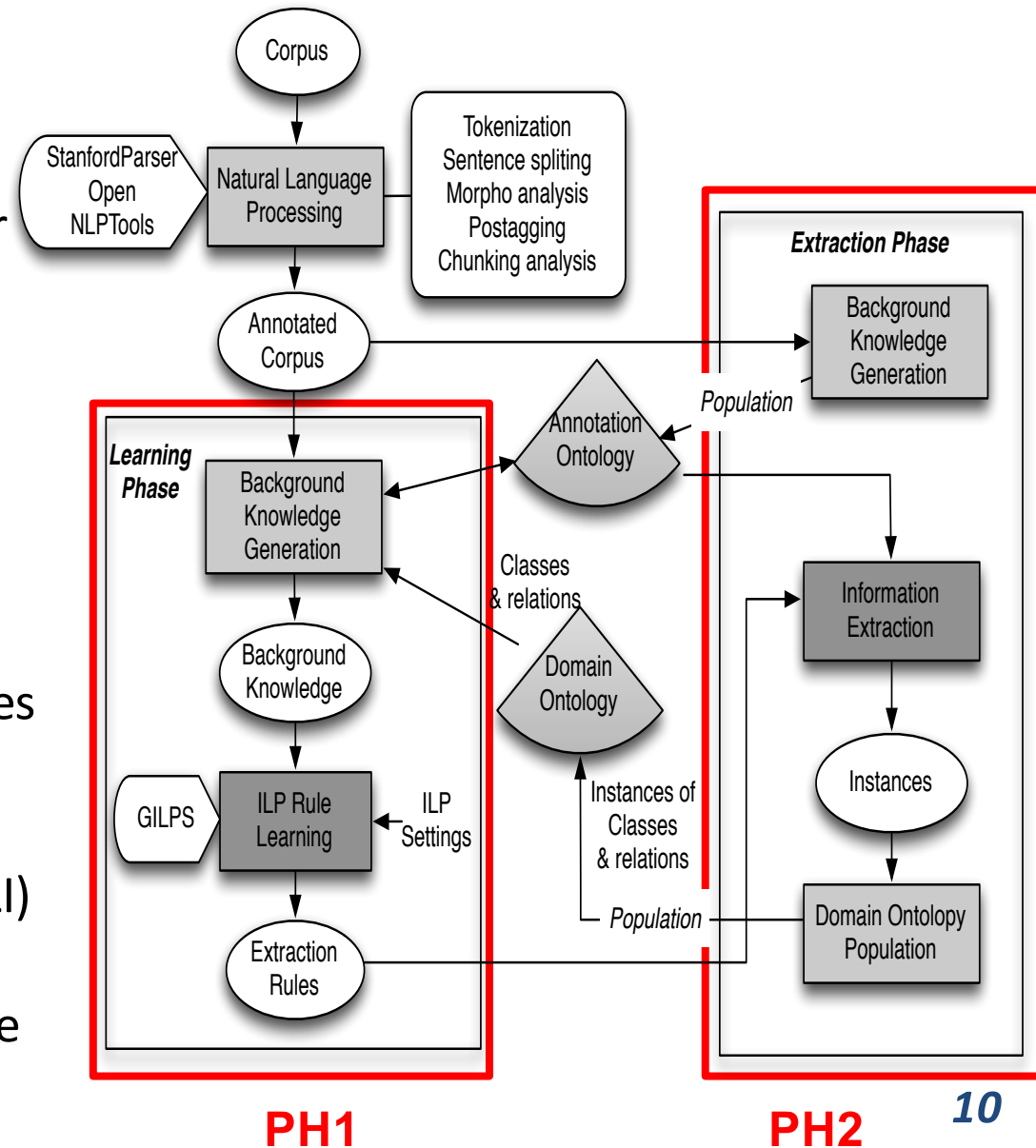
Induction de règles symboliques d'extraction par PLI

- **PH2- Phase d'extraction:**

Extraction d'information par application de ces règles

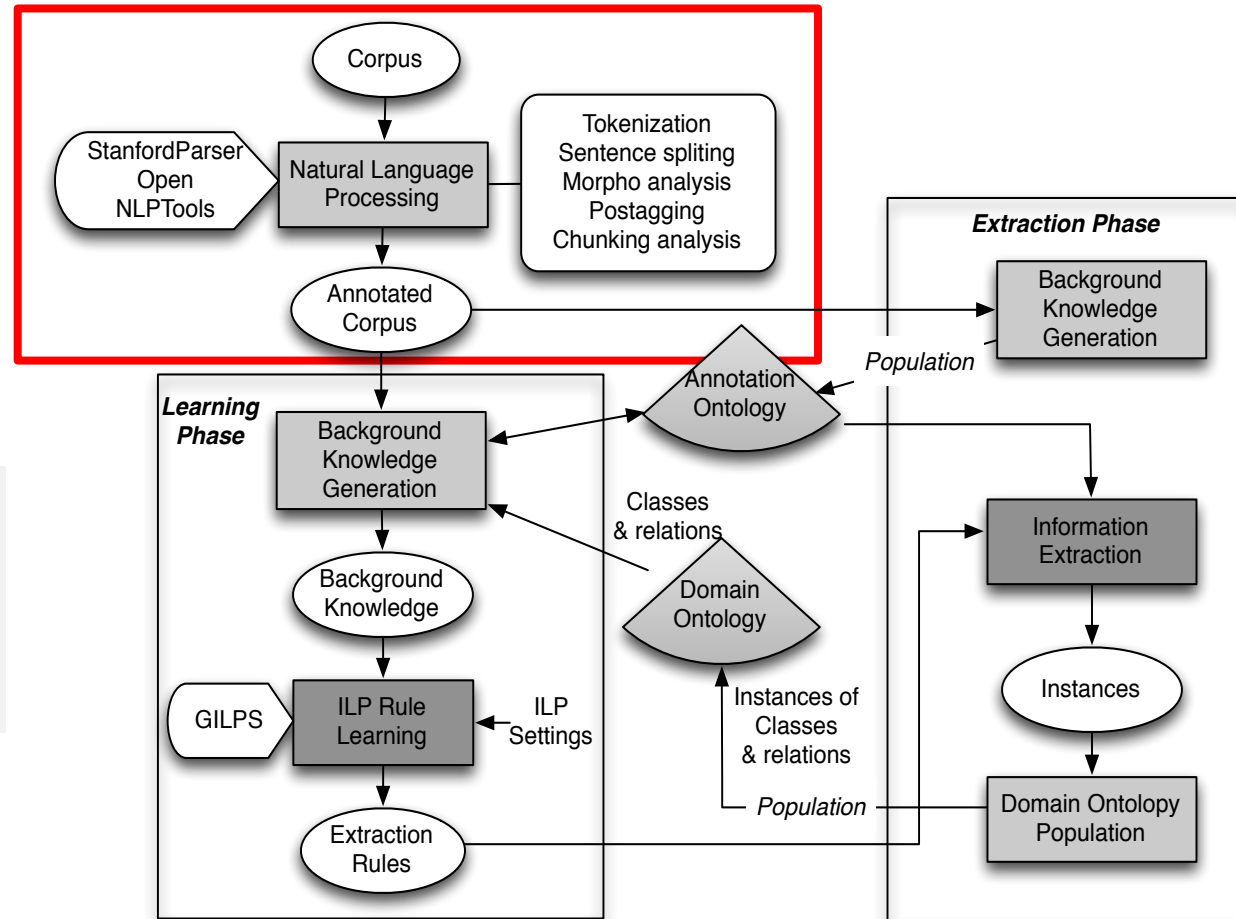
- **Composants majeurs:**

- TAL
- Génération des connaissances de base (BK – Background Knowledge)
- Apprentissage des règles (PLI)
- Application des règles
- Peuplement de l'ontologie de domaine



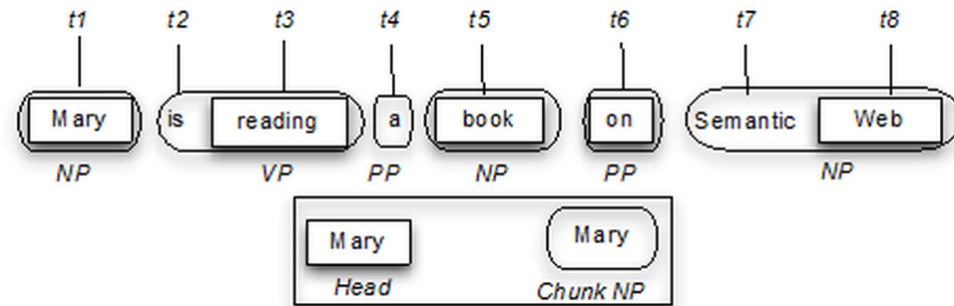
PH1 & PH2: Composant TAL (1)

Composant de traitement automatique du langage naturel (TAL)



PH1 & PH2 : Composant TAL (2)

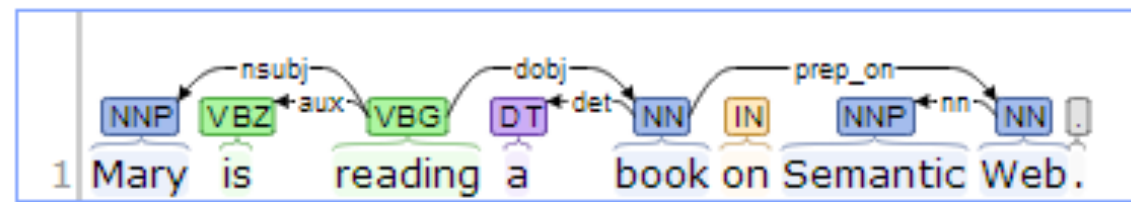
Chunking Analysis



Ce component utilise les outils Stanford CoreNLP & Open NLP

Collapsed CC-processed dependencies:

Dependency Graph



Il r alise: s eparation de phrases, tokenisation,  tiquetage morphosyntaxique (POS tagging, lemmatisation, analyse de chunks), extraction d'entit s nomm es (NER) et analyse des d pendances.

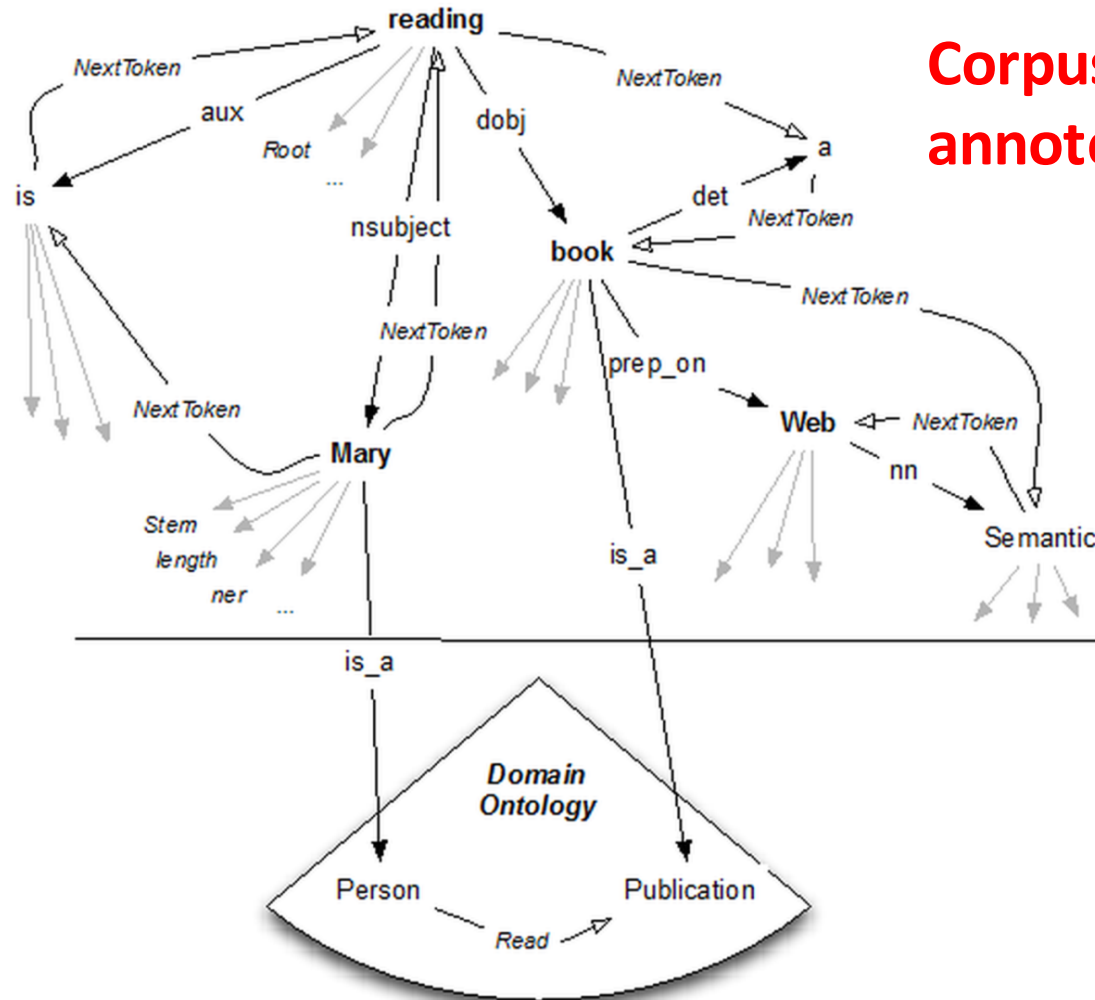
PH1 & PH2: Composant TAL (3)

Corpus
annot 

Final Graph-based
Model of
sentences

Ground
Predicates in BK

nsubj (reading, Mary)
aux (reading, is)
det (book, a)
head (Mary, NP)
nextToken (Mary, is)
nextToken (is, reading)
length (Mary, 4)
ner(Marry, person)



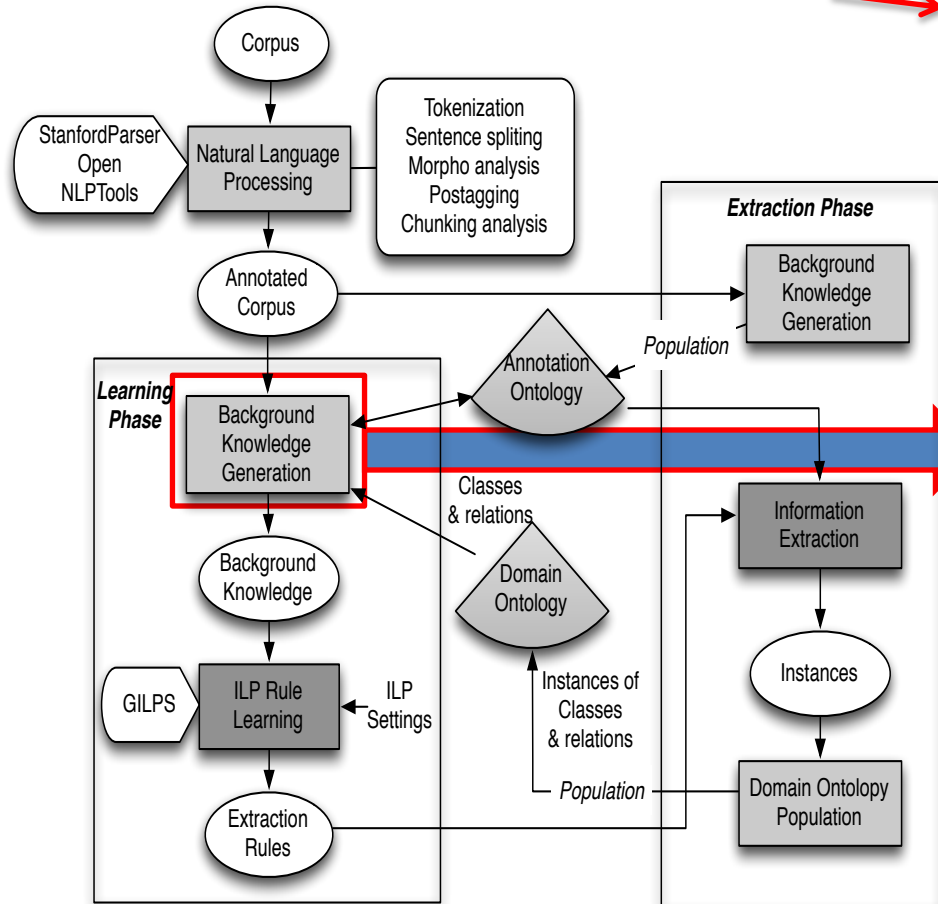
Phrase: « Mary is reading a book on Semantic Web »

PH1: Génération des connaissances de base (BK)

Pour le token « Mary » (T_1)

PROLOG PREDICATES FOR THE TOKEN "MARY" (T_1)

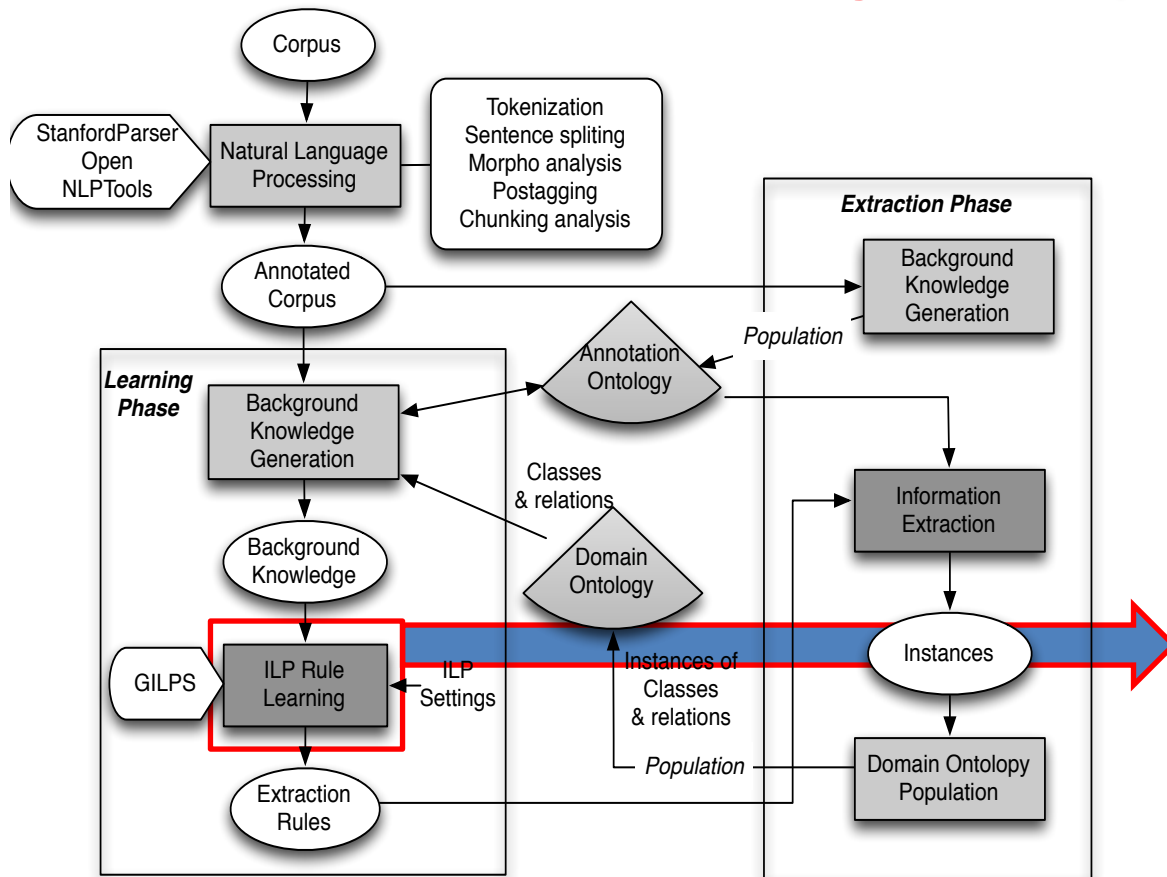
Predicates Generated	Meaning
<i>token (t_1)</i>	t_1 is the token identifier
<i>t_dep (nsubj, t_3, t_1)</i>	there is a noun subject dependency between token t_3 and t_1
<i>t_next (t_1, t_2)</i>	token t_2 follows token t_1
<i>t_stem (t_1, "Mary")</i>	the stemming of the token t_1 is "Mary"
<i>t_length (t_1, 4)</i>	t_1 has length of 4
<i>t_orth (t_1, upperInitial)</i>	t_1 has an initial uppercase letter
<i>t_type(t_1, word)</i>	t_1 is a word
<i>t_pos (t_1, nnp)</i>	t_1 is a singular proper noun
<i>t_ner (t_1, person)</i>	t_1 is a person entity
<i>t_root (t_3)</i>	t_3 is the root of the dependency graph
<i>t_bigposbef (t_n, ...)</i>	POS tag bigram of the tokens after t_n
<i>t_bigposaft (t_1, vbz-vbg)</i>	POS tag bigram of the tokens before t_1
<i>t_trigposbef (t_n, ...)</i>	POS tag trigram of the tokens before t_n
<i>t_trigposaft (t_1, vbz-vbg-dt)</i>	POS tag trigram of the tokens after t_1
<i>t_isHeadNP (t_1)</i>	t_1 is the head of the nominal chunking
<i>t_isHeadVP (t_n...)</i>	t_n is the head of the verbal chunking
<i>t_isHeadPP (t_n...)</i>	t_n is the head of the prepositional chunking
<i>t_ck_tag (t_1, NP)</i>	t_1 is part of a nominal chunking



Base de faits Prolog

PH1: Composant PLI d'apprentissage de r gles

Base de faits Prolog



```
% MODE DECLARATIONS
:-modeb(1, located_in(+token, +token)).
```

```
:-modeb(*, t_hasDep(#dep, +token, -token)).
:-modeb(1, t_next(+token, -token)).
```

Relational features

```
:-modeb(1, t_root(+token)).
:-modeb(1, t_stem(+token, #string)).
:-modeb(1, t_pos(+token, #postag)).
:-modeb(1, t_length(+token, #int)).
:-modeb(1, t_orth(+token, #orth)).
:-modeb(1, t_type(+token, #type)).
:-modeb(1, t_ck_ot(+token, #ck_tag)).
:-modeb(1, t_ner(+token, #ner)).
:-modeb(1, t_ck_tag_ot(+token, #string)).
:-modeb(1, t_gpos(+token, #gpos)).
```

Morpho syntactical features

```
:-modeb(1, t_isHeadNP(+token)).
:-modeb(1, t_isHeadVP(+token)).
:-modeb(1, t_isHeadPP(+token)).
```

Chunking features

```
:-modeb(1, t_bigPosBef(+token, #bigposbef)).
:-modeb(1, t_bigPosAft(+token, #bigposaft)).
:-modeb(1, t_trigPosBef(+token, #trigposbef)).
:-modeb(1, t_trigPosAft(+token, #trigposaft)).
```

N-grams features

Utilise de syst me noyau de PLI « GILPS »
(Santos 2010)

Exemples de r gles d'extraction induites

R gle pour la relation "located_in":

•located_in (A,B):- t_class(A, loc), t_next(A, B), t_class(B, loc).

→ Cette r gle caract rise le patron "City, Country", ex. "Marseille, France"

•located_in (A,C) :- t_next(A,B), t_next(B, C), t_ner(A, org), t_class(C, loc).

→ Cette r gle caract rise le patron "ORG, [at|in|on] LOC", ex. "White House in USA"

...

R gles pour d'autres relations:

•part_whole (A,B):- t_gpos(A, nn), t_next(A, B), t_subtype(B, state-or-province).

•part_w(A,B):- t_next(A,B), t_pos(A, nnp), t_ne_type(B, gpl), t_subtype(A, pop-center).

•live_in(A,B):- t_pos(A, nn), t_class(A, person), t_hasDep(amod, B, C), t_next(C, B), t_class(B, loc), t_isHeadNP(B).

Ces r gles sont compr hensibles (et modifiables) par des humains ...

Protocole expérimental

- 2 corpus de références : reACE 2004/2005 datasets (broadcast news):

Relation
Type &
Subtype

reACE 2004 - Relation Type/Subtype Hierarchy	Freq
Employee-Membership-Subsidiary (EMP ORG)	
Employee-Staff	303
Employee-Executive	220
Member-of-Group	80
General-Affiliation (GEN_AFF)	
Located	352
Citizen-Resident-Religion-Ethnic	98
Part-Whole (PRT_WHOLE)	
Part-Whole	174
Subsidiary	100
Personal-Social (PER_SOC)	
Business	35
Family	15
Total	1377

reACE 2005 - Relation Type/Subtype Hierarchy	Freq
Organization-Affiliation (ORG_AFF)	
Employment	228
Membership	36
General-Affiliation (GEN_AFF)	
Located	280
Citizen-Resident-Religion-Ethnic	39
Part-Whole (PRT_WHOLE)	
Geographical	119
Subsidiary	47
Personal-Social (PER_SOC)	
Business	16
Family	42
Total	807

- Métriques d'évaluation : **Precision (P)**, **Recall (R)** et **F1-measure (F1)**
- Theory Compression Ratio (*mesure de généralité de la règle pour éviter le sur-apprentissage*):

$$TCR = \frac{\text{Nb de règles dans la théorie apprise}}{\text{nb d'exemples positifs dans l'ensemble d'apprentissage}}$$

- 5 validations croisées

Quelle combinaison de BK linguistique est la meilleure ?

ID	Features	reACE 2004			reACE 2005		
		P	R	F1	P	R	F1
1	Baseline	81.09	39.81	53.40	60.53	25.12	35.52
2	+C	80.17	47.13	59.36	75.05	34.03	46.80
3	+D	81.01	46.93	59.43	72.91	36.51	48.65
4	+D+C	89.01	54.40	67.53	74.81	38.14	50.48
5	+D+C+P	91.16	62.04	73.83	81.75	44.24	57.37
6	+D+C+P+Cr	93.30	66.68	77.77	83.68	50.43	62.91
7	+D+C+P+Cr+N	93.04	67.12	77.99	80.59	51.39	62.68
8	+D+C+P+Cr+A	92.20	71.13	80.31	83.03	63.38	71.86
9	+D+C+P+Cr+A+N	92.91	73.07	81.80	82.30	61.85	70.62

P = Precision
R = Rappel
F1= F1-measure

Relation
subtypes

Baseline = Morphological + Next

C = Nominal and verbal chunkings

D = Dependencies

P = POS tagging

Cr = Chunking-related features

N = NER

A = reACE Corpus types

Structural features

Attributive features

Semantic features

Résultats de classification sur les sous-types de relations et les relations-types (1)

PERFORMANCE RESULTS OF RELATION SUBTYPES ON BOTH DATASETS

Relation subtypes

Rel. Type	Rel. Subtype	reACE 2004			reACE 2005			
		P	R	F1	Rel. Subtype	P	R	F1
EMP_ORG	Employ-Staff	78.10	86.90	82.27	Employ	89.60	86.22	87.88
	Employ-Exec	95.49	77.00	85.25	-	-	-	-
	Member	92.18	76.82	83.80	Member	94.30	71.03	81.03
GEN_AFF	Citizen-Resident	98.81	69.58	81.66	Citizen-Resident	100.00	61.10	75.85
	Located	83.28	80.09	81.65	Located	86.00	84.10	85.04
PERS_SOC	Business	100.00	69.42	81.95	Business	0.00	0.00	0.00
	Family	100.00	39.11	56.23	Family	92.70	57.70	71.13
PRT_WHL	Part-Whole	93.20	83.38	88.02	Geo	100.00	62.10	76.62
	Subsidiary	95.10	75.30	84.05	Subsidiary	95.80	72.51	82.54
Avg		92.91	73.07	81.80	Avg	82.30	61.85	70.62

Il est plus difficile de classifier à un niveau détaillé (sous-type de relation)

Relation types

CLASSIFICATION RESULTS OF RELATION TYPES ON THE REACE 2004

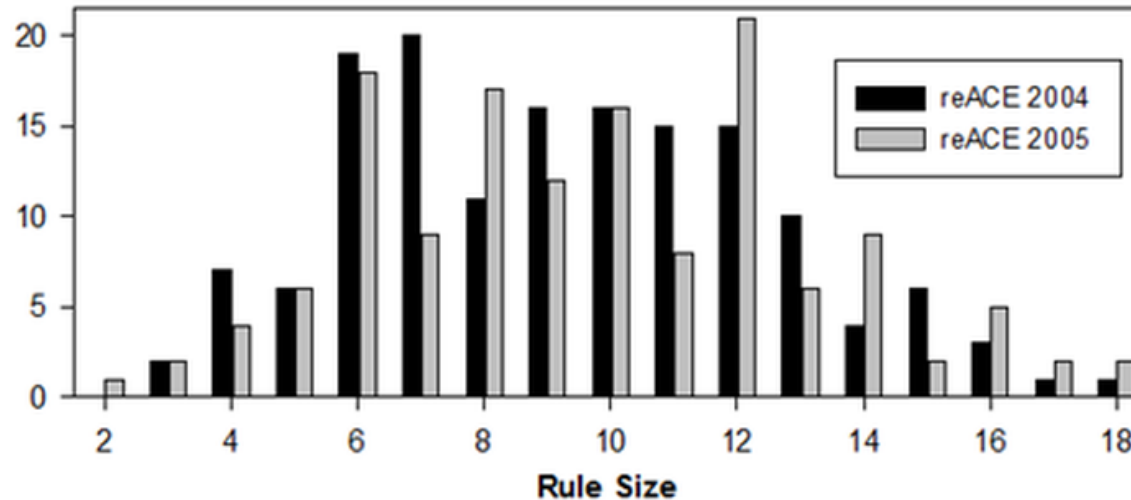
Rel. Type	#E+	#Rules	TCR	P	R	F1
EMP_ORG	603	65	0.11	86.00	84.00	84.99
GEN_AFF	450	51	0.11	86.90	78.90	82.71
PER_SOC	50	18	0.36	100.00	64.40	78.35
PRT_WHL	274	33	0.12	91.00	81.60	86.04
Total	1377	167				
Avg			0.18	90.98	77.23	83.54

CLASSIFICATION RESULTS OF RELATION TYPES ON THE REACE 2005

Rel. Type	#E+	#Rules	TCR	P	R	F1
ORG_AFF	264	38	0.14	88.70	77.80	82.89
GEN_AFF	319	60	0.19	94.40	70.40	80.65
PER_SOC	58	19	0.33	100.00	58.30	73.66
PRT_WHL	166	35	0.21	87.60	72.30	79.22
Total	807	152				
Avg			0.22	92.68	69.70	79.56

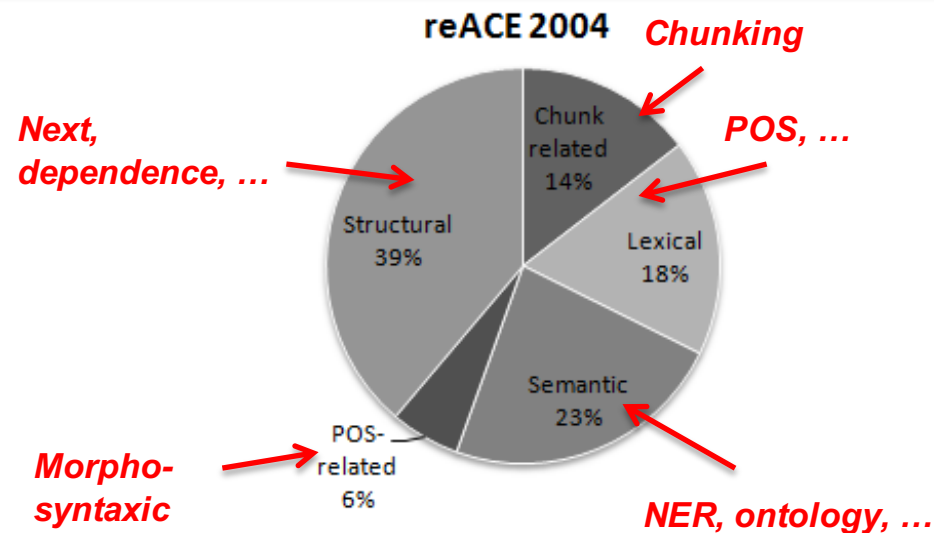
Aspects qualitatifs des règles apprises

Distribution
de la taille des
règles



80% des règles
ont moins de 11
prédicats

Ratio des
types des
prédicats dans
l'ensemble des
règles induites



Conclusion - Perspectives

Optimisation :

- Utilisation de **méthodes d'ensemble**, pour améliorer les performances du **processus d'apprentissage (PH1)**
- Utilisation d'un **triple-store** et des **techniques du Web Sémantique** pour améliorer l'application des règles symboliques dans le **processus d'extraction (PH2)**
- Usage de nouveaux algorithmes de PLI (ParProGolem, ...)
- Usage du parallélisme (Prolog parallèle, ...), Map-Reduce, cartes GPU, ...

Extensions :

- **Intégration de plus de connaissances dans la BK** au niveau du prétraitement, pour prendre en compte les **aspects sémantiques** (prédicats sur les synonymes, hypernoms / hyponymes, rôles sémantiques ...)
- **Extraction d'évènements** : passage de l'extraction de relations binaires à l'extraction de relations n-aires (**Event Extraction**)
- Utilisation de représentations plus riches : **Abstract Meaning Representation, ...**

Pour aller plus loin: publications associées ...

Reuves :

- R. Lima, R.; Espinasse, B.; Freitas, F. (2018), « A Logic-based Relational Learning Approach to Relation Extraction: the OntoLPER System ». Journal of Engineering Applications of Artificial Intelligence - Elsevier (**EAAI**) (*accepted for publication*).
- R. Lima, B. Espinasse, F. Freitas (2017), « An Ontology-and inductive logic programming-based system to extract entities and relations from text ». Knowledge And Information System (**KAIS**) Journal, Springer-Verlag, v. 54, p. 1-33, Springer-Verlag, 2017. <https://doi.org/10.1007/s10115-017-1108-3>.
- B. Espinasse, R. Lima, F. Freitas (2016), « Extraction automatique d'entités et de relations par ontologies et programmation logique inductive », in: Revue d'Intelligence Artificielle (**RIA**), Vol. 30 (n° 6/2016), dec 2016 (Répertoriée Scopus et DBLP).

Conférences :

- B. Espinasse, Lima R., Magdy D., « Extraction automatique d'entités et de relations par ontologies et programmation logique inductive », Journée Francophones sur les Ontologies - **JFO 2016**, 13-14 Octobre 2016, Bordeaux, France..
- R. Lima, S. B. Espinasse, F. Freitas « Relation Extraction from Texts with Symbolic Rules Induced by Inductive Logic Programming », IEEE International Conference on Tools with Artificial Intelligence, **IEEE-ICTAI 2015**, Vietri sul Mar, Italy, 9-11 nov. 2015.
- R. Lima, B. Espinasse, H. Oliveira, F. Freitas « Ontology Population from the Web: an Inductive Logic Programming-Based Approach », 11th International Conference on Information Technology: New Generations, **ITNG 2014**, Las Vegas, Nevada, USA, April 7-9, 2014.
- R. Lima, B. Espinasse, H. Oliveira, L. Pentagrossa, F. Freitas, « Information Extraction from the Web: An Ontology-Based Method using Inductive Logic Programming », IEEE International Conference on Tools with Artificial Intelligence, **IEEE-ICTAI 2013**, Washington DC, USA, November 4-6, 2013.
- R. Lima, B. Espinasse, H. Oliveira, R. Ferreira, L. Cabral, F. Freitas, R. Gadelha, « An Inductive Logic Programming-Based Approach for Ontology Population from the Web », **DEXA 2013**, Prague, Czech Republic, August 26-29, 2013.