

## Approches génériques du TAL pour la RI

---

Patrice Bellot — Aix-Marseille Université (AMU) / LSIS

octobre 2012  
patrice.bellot@univ-amu.fr



1

## Plan

---

- Linguistique et TAL : objets d'étude, applications...
- Un exemple bien connu : l'analyse syntaxique
- Des ressources linguistiques informatisées
- Extraction d'information
  - Reconnaissance d'entités nommées
  - Annotation automatique
- Recherche d'informations précises
  - Systèmes de questions-réponses

P. Bellot

2

## Introduction

3

## Quelques thèmes de la revue TAL

<http://atala.org/-Appels-a-soumission->

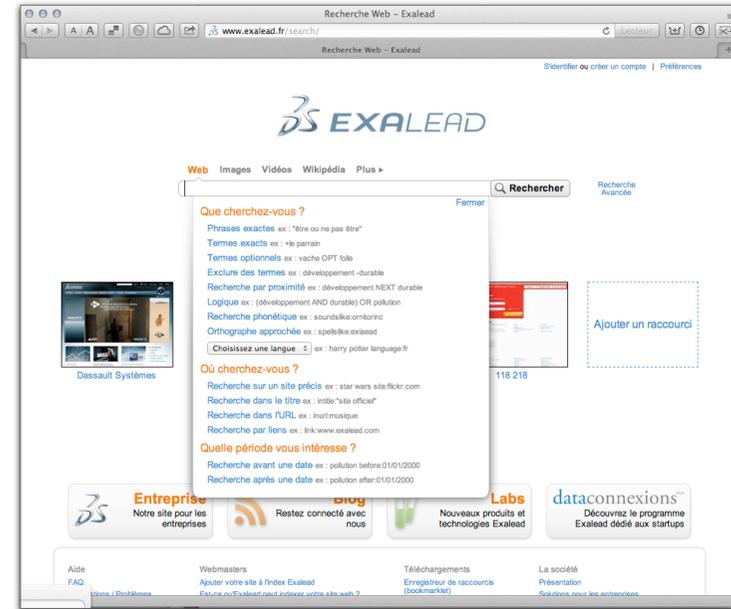
- Modèles et algorithmes pour la résolution d'anaphores.
- Systèmes de question/réponse.
- Les dictionnaires électroniques : pour les personnes, les machines ou pour les deux ?
- Le résumé automatique de texte
- Fouille de données textuelles : complexité, algorithmique et passage à l'échelle.
- Recherches actuelles en phonologie et en phonétique : interfaces avec le traitement automatique des langues.
- Opinions, sentiments et jugements d'évaluation.
- Traitement automatique des informations temporelles et spatiales en langage naturel.
- Du bruit dans le signal : gestion des erreurs en traitement automatique des langues.
- ..... **la Recherche d'information**

P. Bellot

4

## Exemples d'applications du TAL en RI

- Indexation de syntagmes nominaux
- Normalisation d'unités lexicales (mots, entités nommées...) et désambiguïsation sémantique
- Typage (annotation) des mots des documents
- Recherche de paraphrases (identification de variantes)
- RI multilingue (traduction automatique, alignement)
- Analyse de la temporalité, de l'opinion, de sentiments...
- Interactivité et systèmes de dialogue
- Correction orthographique, suggestion de requêtes
- RI multimédia (reconnaissance de la parole)
- RI à partir de documents *bruités* (OCR, tweets, SMS ...)



<http://asso-aria.org/coria/2009/19.pdf>

## stemmer et lemmatiseur

### Evaluation de diverses stratégies de désambiguïsation lexicale

Claire Fautsch, Jacques Savoy

Institut d'informatique  
Université de Neuchâtel, rue Emile Argand 11, 2009 Neuchâtel (Suisse)  
Claire.Fautsch@unine.ch, Jacques.Savoy@unine.ch

```
<NUM> C180 </NUM>
<EN-TITLE> Bankruptcy of Barings </EN-TITLE>
<EN-DESC> What was the extent of the losses in the Barings
</EN-DESC>
<EN-NARR> Relevant documents must quantify in some way the
collapse of the oldest bank in Great Britain </EN-NARR>
...
<NUM> C180 </NUM>
<EN-TITLE>
<TERM ID="10.2452/180-AH-1" LEMA = "bankruptcy" POS = "NNP">
<WF> Bankruptcy </WF>
<SYNSET SCORE = "0.4819665883771086" CODE = "10386276-n"/>
<SYNSET SCORE = "0.5180534116228914" CODE = "10386165-n"/>
<TERM ID = "10.2452/180-AH-2" LEMA = "of" POS = "IN">
<WF> of </WF> </TERM>
<TERM ID = "10.2452/180-AH-3" LEMA = "baring" POS = "NNPS">
<WF> Barings </WF>
<SYNSET SCORE = "1" CODE = "00819570-n"/> </TERM>
</EN-TITLE>
...
```

Figure 2 : Exemple d'une requête avec les indications du lemme, de sa discours, et des numéros de classe de WordNet associées

	Précision moyenne (MAP)					
	aucun	S-stemmer ‡	Porter	Lovins	SMART	lemme
Okapi	<b>0,3743</b>	0,4044†	<b>0,4150†‡</b>	<b>0,3930</b>	<b>0,4152†‡</b>	0,3988†
DFR-PL2	0,3703	0,4006†	0,4116†‡	0,3927†	0,4096†‡	<b>0,3994†</b>
DFR-I(n <sub>c</sub> )C2	0,3731	<b>0,4054†</b>	0,4141†‡	0,3894 ‡	0,4139†‡	0,3988†
LM	0,3445	0,3709†	0,3809†‡	0,3522 ‡	0,3760†‡	0,3602
<i>tf idf</i>	0,2230	0,2393†	0,2399†	0,2194 ‡	0,2431†‡	0,2308
Moyenne	0,3370	0,3641	0,3723	0,3493	0,3716	0,3576
% différence		+8,0 %	+10,5 %	+3,6 %	+10,2 %	+6,1 %

Table 1. Précision moyenne (MAP) obtenues avec différents modèles et enraineurs (284 requêtes, format T)

	Précision moyenne (MAP)		
	lemme	lemme & POS	lemme & synset
Okapi	0,3988	<b>0,4053†</b>	0,3986
DFR-PL2	<b>0,3994</b>	0,4013	0,3918
DFR-I(n <sub>c</sub> )C2	0,3988	0,4047	<b>0,4018</b>
LM	0,3602	0,3659†	0,3546
<i>tf idf</i>	0,2308	0,2315	0,2325
Moyenne	0,3576	0,3617	0,3559
% différence		+1,2 %	-0,5 %

Table 2. Précision moyenne (MAP) pour différents modèles de RI et variantes d'analyse morphologique (284 requêtes, format T)

## Combinaison d'index ?

- Thèse de F. Moreau (IRISA, 2006) «*Revisiter le couplage traitement automatique des langues et recherche d'information*»

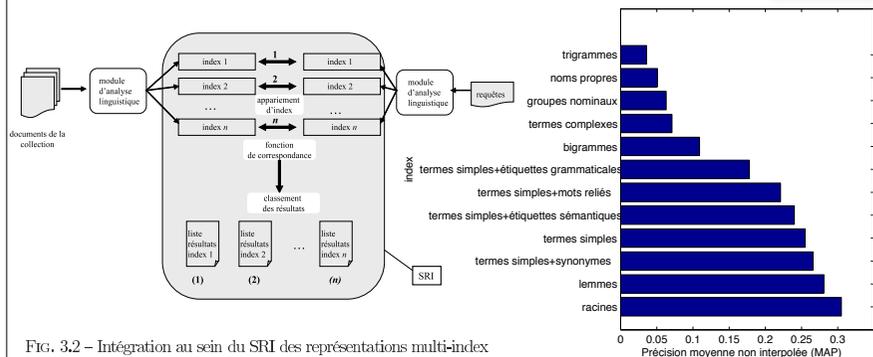


Fig. 3.2 - Intégration au sein du SRI des représentations multi-index

paire d'index (index 1 + index 2)	taux de résultats différents (en %)	documents pertinents retrouvés uniquement par le 1 <sup>er</sup> index (en %)	documents pertinents retrouvés uniquement par le 2 <sup>nd</sup> index (en %)
trigrammes + synonymes	72,22	1	99
trigrammes + étiquettes sémantiques	66,72	1	99
trigrammes + mots reliés	66,30	2	98
groupes nominaux + synonymes	61,96	0	100
noms propres + synonymes	59,83	4	96
groupes nominaux + étiquettes sémantiques	57,92	2	98
noms propres + étiquettes sémantiques	56,27	6	94
groupes nominaux + mots reliés	56,30	2	98
noms propres + mots reliés	55,97	8	92
termes complexes + synonymes	54,18	2	98
termes complexes + étiquettes sémantiques	50,50	3	97
termes complexes + mots reliés	49,53	4	96
bigrammes + synonymes	44,86	8	92
racines + trigrammes	43,31	97	3
lemmes + trigrammes	42,66	98	2
bigrammes + mots reliés	41,95	13	87
racines + groupes nominaux	41,54	99	1
bigrammes + étiquettes sémantiques	41,08	10	90
lemmes + groupes nominaux	40,60	99	1
termes simples + trigrammes	40,23	97	3
noms propres + racines	38,43	8	92
termes simples + groupes nominaux	38,22	97	3
racines + termes complexes	38,09	95	5

	Index des racines	Fusion sans infos req. (amélioration %)
Précision	27,90	26,49 (-5,04%)
Rappel	51,32	41,65 (-18,84%)
F-mesure	29,67	30,28 (+2,05%)

	Fusion sans infos req. Moyenne (cart-type)	Fusion sans infos req. Moyenne (cart-type)
Précision	29,48 (3,81)	26,49 (4)
Rappel	45,88 (3,77)	41,65 (1,72)
F-mesure	34,20 (3,92)	30,28 (3,97)

Tab. 4.10 – Moyennes et écart-type des valeurs de précision, rappel et F-mesure obtenus sur les 10 jeux de test par la méthode de fusion avec et sans prise en compte des informations sur les requêtes

Tab. 4.6 – Moyennes des précision, rappel et F-mesure obtenues par la méthode de fusion sans prise en compte des informations sur la requête et comparées aux performances d'un SRI exploitant des racines

## Exemples de requêtes TREC

### • Question-Answering Track

- What Jeopardy contestant is the biggest money-winner in television game show history?
- Into how many languages has "Harry Potter and the Goblet of Fire" been translated?

#### Utilisation de *narrative* ?

```
<?xml version="1.0" encoding="UTF-8" ?>
<topic num="68">
  <template id="2">
    What [financial relationships] exist between [DARPA] and [BBN]?
  </template>
  <narrative>
    The analyst would like to know of financial relationships between the Defense Advanced Research Projects Agency and BBN Technologies in Cambridge. Specifically, the analyst would like to know of projects funded by DARPA at BBN.
  </narrative>
</topic>
```

### • Entity Track

```
<query>
<num>109</num>
<entity_name>Airbus</entity_name>
<entity_URL>clueweb09-en0004-69-31337</entity_URL>
<target_entity>airline</target_entity>
<narrative>airlines that fly the Airbus 320 plane</narrative>
</query>
```

## Exemples de requêtes TREC (2)

### • Medical Track

- Patients diagnosed with localized prostate cancer and treated with robotic surgery
- Patients who received methotrexate for cancer treatment while in the hospital

### • Chemical IR Track

```
<top>
<num>TCS-xx</num>
<query>liquid coating composition</query>
<narr>
  I need all available data on the production of liquid coating compositions.
</narr>
</top>
```

### • Enterprise Track

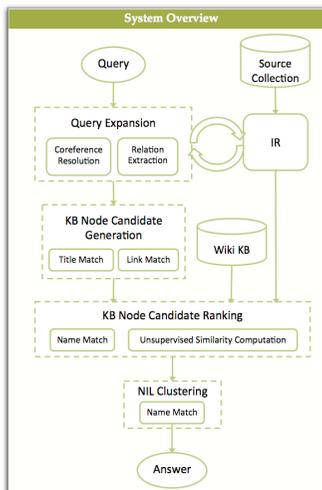
```
<top>
<num>CE-053</num>
<query>selenium soil</query>
<narr>
  Can you please provide a current e-mail address, or failing that can you please put me in contact with the group responsible for the research into the use of selenium as an additive to soils, to promote sheep productivity/health. There were some trials conducted in WA and I am looking for additional information on these.
</narr>
</top>
```

## Exemples de tâches DEFT (Défi Fouille de Textes)

- détection de l'**opinion** (positive, neutre, négative) exprimée dans un texte (débats parlementaires, critiques de livres et de jeux vidéos, relectures d'articles scientifiques)
- **classification** automatique de corpus **en genres** (Le Monde vs. Wikipédia) et en catégories (International, Politique française, Société, Economie, Sciences, Art, Littérature, Sports, Télévision)
- **détermination du parti politique d'appartenance** des parlementaires d'après leurs interventions + détection des articles globalement objectifs (**factuels**) ou **subjectifs** (porteurs d'opinion) dans les corpus de journaux
- **identification de l'année** de publication d'un extrait de journal (de 1801 à 1944) et appariement de résumés avec des articles scientifiques
- **identification des mots-clés** indexant le contenu d'un article scientifique paru en revue de SHS

## Exemple de tâche

- TAC KBP (*Knowledge Base Population*)



P. Bellot L. Bonnefoy, P. Bellot 2012

**Query Expansion**

Acronym expansion: if the source entity is an acronym (all its letters are upper-cased), named entities with first letters of each words corresponding to the acronym are considered as extended forms.

Name variations resolution:

- Coreference resolution (Stanford CoreNLP)
- Recall oriented
- All named entities in a chain are selected if at least one of them contains a part of the source entity.

$$w(v_i) = \sum_d t f(v_i, d) * conf(d)$$

Relation extraction:

- Relation of interest non specified in advance
- Tool used: Reverb
- Relations containing a variant are selected

$$w(r_i, e_j) = \sum_d t f(r_i, e_j, d) * conf(d)$$

IR:

- Retrieve documents containing at least one tuple  $v_i, r_i, e_j$ .
- Top 3 documents are selected for the next QE step
- Loop 5 times

**Candidate Generation and Ranking**

Candidate KB nodes are selected if:

- *Exact Match*: the KB node have its title matching a name variation
- *Link Match*: a link in the KB to the node matches a name variation.

Candidate KB nodes are ranked according to:

- *Title Match*:
 
$$s(t, v) = \sum_{v_i \in V} (1 + \sum_{v_j \in V} w(v) * isOfType(v_i, type_{v_j})) * w(v_j) * (1 + \log(t f(v_j, kb)))$$
- *Unsupervised Similarity Computation*:
  - Cosine similarity between two vectors of words

13

## TAL : quelles unités d'étude ?

- **Mot ? Phrase ? Paragraphe ? Texte ? Document ?**
- **Chunk** : séquence de mots + catégorie associée (groupe nominal, groupe verbal / pas d'inclusion d'autre chunk)
- **Syntagme** : nominal, verbal
- **Clause** : sujet + prédicat + ...
- **Énoncé** : prise de parole
- **Proposition** : valeur de vérité associée

P. Bellot

14

## Quelques définitions

- **Mot** : unité (empirique) du *lexique* – mots dénominatifs, désignant, grammaticaux...
- **Lemme** : forme simple d'un mot permettant d'identifier un groupe de flexions
- **Terme** : nom, entité nommée, mot de spécialité utile pour l'indexation (mot clé)
- **Thesaurus** : base de connaissance lexicale qui permet d'identifier des associations de mots (synonymie, hyperonymie....)
- **Syntaxe** : [étude] ordre des mots / des groupes – micro/macro syntaxe vs. **morphologie** (structure interne des mots : flexions, dérivations...) → **morphosyntaxe** vs. **sémantique**...
- **Syntagme** : constituant syntaxique, organisé autour d'une tête (nom pour SN, verbe pour SV, adjectif pour SA...) : **le stylo que tu m'as offert écrit très bien, je prendrai un café serré, très heureux d'être ici !**

P. Bellot

15

## Phrase, signification et sens

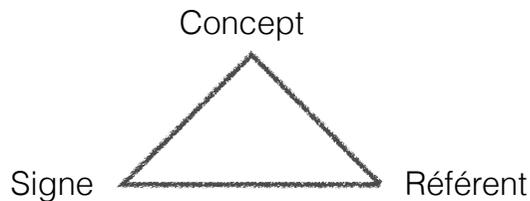
- *Le cours commence dans dix minutes.*
  - **Sémantique** : le sens indépendamment du contexte = la signification
  - **Pragmatique** : le sens en contexte (environnement, texte qui précède, voire hypothèses posées par l'auditeur) = le sens
  - Lexique : les mots sont porteurs de sens / transforment le sens = liens avec *concepts, catégories* (les entités désignées par un mot)
  - Syntaxe : indique que c'est le *cours* qui va commencer
  - Contexte : aide à l'interprétation, influence le sens des mots
  - Un *énoncé* : une phrase + informations non linguistiques (locuteur, ...)

P. Bellot

16

## Le triangle sémiotique

- Mot / signe, forme orthographique / signifiant, concept / signifié
- **Triangle d'Odgen et Richards** : les mots (*signes*) désignent des *référents* (entités du monde, l'extension) via des *concepts* (l'intension)
  - Noms propres : lien direct signe / référent
  - Connecteurs, adverbes, pronoms... : pas de sens sans référent (*information procédurale*) --> résolution des références (anaphores...)
  - Expression génériques : peu importe le référent (classe conceptuelle)



## Notion de pertinence (*relevance*)

- Pertinence vs. Vérité (subjectif/objectif, vrai/faux...)
- Est pertinent ce qui crée de l'information (dans l'esprit)
- Notion d'effort / coût (adapté au contexte, auditeurs, besoin...)
- RI : constitution de référentiels (qrels), mesures de précision/rappel...

**Relevance Theory\***

DEIRDRE WILSON & DAN SPERBER

---

**Abstract**

This paper outlines the main assumptions of relevance theory (Sperber & Wilson 1985, 1995, 1998, 2002, Wilson & Sperber 2002), an inferential approach to pragmatics. Relevance theory is based on a definition of relevance and two principles of relevance: a Cognitive Principle (that human cognition is geared to the maximisation of relevance), and a Communicative Principle (that utterances create expectations of optimal relevance). We explain the motivation for these principles and illustrate their application to a variety of pragmatic problems. We end by considering the implications of this relevance-theoretic approach for the architecture of the mind.

**1 Introduction**

Relevance theory may be seen as an attempt to work out in detail one of Grice's central claims: that an essential feature of most human communication, both verbal and non-verbal, is the expression and recognition of intentions (Grice 1989: Essays 1-7, 14, 18; Retrospective Epilogue). In developing this claim, Grice laid the foundations for an inferential model of communication, an alternative to the classical code model. According to the code model, a communicator encodes her intended message into a signal, which is decoded by the audience using an identical copy of the code. According to the inferential model, a communicator provides evidence of her intention to convey a certain meaning, which is inferred by the audience on the basis of the evidence provided. An utterance is, of course, a linguistically coded piece of evidence, so that verbal comprehension involves an element of decoding. However, the linguistic meaning recovered by decoding is just one of the inputs to a non

---

\* A version of this paper will appear in L. Horn and G. Ward (eds.) *Handbook of Pragmatics* (Oxford: Blackwell), and a shortened version in *Proceedings of the Tokyo Conference on Psycholinguistics 2002*. We are grateful to Larry Horn, Tomoko Matsui, Yuji Nishiyama, Yukio Otsu and Gregory Ward for many valuable comments and suggestions.

<http://people.bu.edu/bfraser/Relevance%20Theory%20Oriented/Sperber%20&%20Wilson%20-%20ORT%20Revisited.pdf>

## Analyseurs (syntaxic)



## Analyseurs en partie du discours (*POS tagger*)

- Etiquetage morpho-syntaxique (la syntaxe influence la forme)
  - nom, verbe, article/déterminant, adjectif, pronom, adverbe, préposition, conjonction
  - masculin / féminin ( / neutre), singulier / pluriel, auxiliaire, temps, participe, voie (active / passive)...
- Méthodes :
  - règles d'étiquetage
  - modèles de Markov cachés (HMM)
- Exemples :
  - *TreeTagger* <http://www.ims.uni-stuttgart.de/projekte/corplex/TreeTagger/>  
anglais, français, allemand, italien, swahili...
  - *HMMTagger* (dans UIMA)

Sample output:

word	pos	lemma
The	DT	the
TreeTagger	NP	TreeTagger
is	VBZ	be
easy	JJ	easy
to	TO	to
use	VB	use
.	SENT	.

<http://hal.inria.fr/docs/00/49/38/47/PDF/article-taln-2010.pdf>  
Un étiqueteur de rôles grammaticaux libre pour le français intégré à Apache UIMA  
Charles Dejean Manoeil Fortun Clotilde Massot Vincent Pottier Fabien Poulard\* Matthieu Vernier



The Stanford Natural Language Processing Group

home · people · teaching · research · publications · software · events · local

<http://www-nlp.stanford.edu/software/>

- Stanford CoreNLP**  
An integrated suite of natural language processing tools for English in Java, including tokenization, part-of-speech tagging, named entity recognition, parsing, and coreference. [Online CoreNLP demo](#)
- Stanford Parser**  
Implementations of probabilistic natural language parsers, both highly optimized PCFG and dependency parsers, and a lexicalized PCFG parser in Java. Includes: [Online parser demo](#), [Stanford Dependencies page](#), and [Parser FAQ](#).
- Stanford POS Tagger**  
A maximum-entropy (CMM) part-of-speech (POS) tagger for English, Arabic, Chinese, French, and German, in Java.
- Stanford Named Entity Recognizer**  
A Conditional Random Field sequence model, together with well-engineered features for Named Entity Recognition in English and German. [Online NER demo](#)
- Stanford Word Segmenter**  
A CRF-based word segmenter in Java. Supports Arabic and Chinese.
- Stanford Classifier**  
A machine learning classifier, directed at text categorization. A conditional loglinear classifier (a.k.a. a maximum entropy or multiclass logistic regression model).
- Tregex and Tsurgeon**  
A Tregex2-style utility for matching patterns in trees, and a tree-transformation utility built on top of this matching language.
- Phrasal**  
A state-of-the-art phrase-based machine translation system.
- Stanford Biomedical Event Parser (SBEP)**  
Biomedical Event Extraction for the BioNLP 2009/2011 shared task.
- Stanford EnglishTokenizer**  
A fast tokenizer for English text (producing Penn Treebank tokenization, roughly)
- Stanford TokensRegex**  
A tool for matching regular expressions over tokens.
- Stanford Temporal Tagger (SUTime)**  
A rule-based temporal tagger for English text. [Online SUTime demo](#)

Stanford Parser

<http://nlp.stanford.edu/software/lex-parser.shtml>

- Anglais, Chinois, Arabe, Allemand (adaptable à d'autres langues)
- *The strongest rain ever recorded in India shut down the financial hub of Mumbai, snapped communication lines, closed airports and forced thousands of people to sleep in their offices or walk home during the night, officials said today.*

The/DT strongest/JJS rain/NN ever/RB recorded/VBN in/IN India/NNP shut/VBD down/VP the/DT financial/JJ hub/NN of/IN Mumbai/NNP ,/, snapped/VBD communication/NN lines/NNS ,/, closed/VBD airports/NNS and/CC forced/VBD thousands/NNS of/IN people/NNS to/TO sleep/VB in/IN their/PRP\$ offices/NNS or/CC walk/VB home/NN during/IN the/DT night/NN ,/, officials/NNS said/VBD today/NN ./.

Stanford Parser

<http://nlp.stanford.edu/software/lex-parser.shtml>

- Anglais, Chinois, Arabe, Allemand (adaptable à d'autres langues)
- *The strongest rain ever recorded in India shut down the financial hub of Mumbai, snapped communication lines, closed airports and forced thousands of people to sleep in their offices or walk home during the night, officials said today.*

The/DT strongest/JJS rain/NN ever/RB recorded/VBN in/IN India/NNP shut/VBD down/VP the/DT financial/JJ hub/NN of/IN Mumbai/NNP ,/, snapped/VBD communication/NN lines/NNS ,/, closed/VBD airports/NNS and/CC forced/VBD thousands/NNS of/IN people/NNS to/TO sleep/VB in/IN their/PRP\$ offices/NNS or/CC walk/VB home/NN during/IN the/DT night/NN ,/, officials/NNS said/VBD today/NN ./.

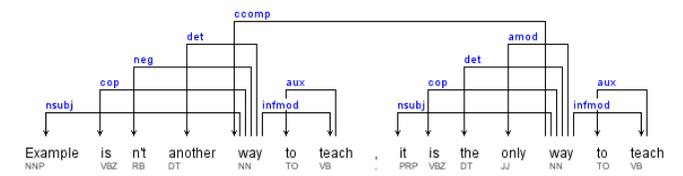
```

(ROOT
(S
(S
(NP
(NP (DT The) (JJS strongest) (NN rain))
(VP
(ADVP (RB ever))
(VBN recorded)
(PP (IN in)
(NP (NNP India))))))
(VP
(VP (VBD shut)
(PRT (RP down))
(NP
(NP (DT the) (JJ financial) (NN hub))
(PP (IN of)
(NP (NNP Mumbai))))))
(, ,)
(VP (VBD snapped)
(NP (NN communication) (NNS lines)))
(, ,)
(VP (VBD closed)
(NP (NNS airports)))
(CC and)
(VP (VBD forced)

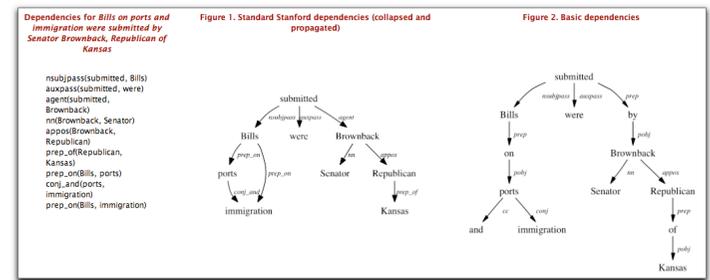
```

Stanford Parser : GUI

- <http://chaotcity.com/dependensee-a-dependency-parse-visualisation-tool/>



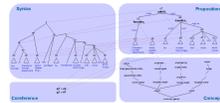
- GrammarScope / Stanford dependencies



# Stanford CoreNLP

<http://www-nlp.stanford.edu/software/coresf.shtml>

- Résolution des références (dont anaphores)
- Evaluation sur la tâche CoNLL-2011 (<http://conll.cemantix.org/2011/introduction.html>)  
She had a good suggestion and it was unanimously accepted



– données du projet OntoNotes (<http://www.bbn.com/ontonotes/>)

Stanford's Multi-Pass Sieve Coreference Resolution System at the CoNLL-2011 Shared Task  
Heeyoung Lee, Yves Peirsman, Angel Chang, Nathanael Chambers, Mihai Surdeanu, Dan Jurafsky  
Stanford NLP Group

Official: Closed track: Predicted mentions

System	MO	MUC	B-CUBED	CEAF	CEAF	BILANC	ORIG
	P	R	P	P	P	P	P
bas	78.10	52.03	68.31	66.37	45.48	73.22	37.79
napena	43.20	59.95	67.09	53.51	41.32	71.10	55.99
chang	64.26	57.15	68.79	54.63	41.64	73.21	55.96
maguik	68.06	58.61	65.46	51.45	39.52	71.11	54.53
napion	65.45	56.65	65.86	50.43	35.91	69.49	53.41
song	67.26	59.95	63.23	46.29	35.96	61.47	53.05
napionov	67.76	58.43	61.44	48.08	35.39	69.48	51.92
sohka	64.23	50.48	64.00	49.48	41.23	63.28	51.90
lobidat	61.03	53.49	60.25	47.70	39.79	62.61	51.94
zhou	62.31	48.96	64.07	47.53	39.74	64.72	50.92
charon	64.30	52.45	62.10	49.23	36.54	64.20	50.96
napie	63.93	52.31	62.32	46.55	35.33	64.63	49.99
napie	64.30	54.41	61.01	49.07	32.87	65.35	49.36
vitain	61.92	46.62	61.93	44.75	35.23	64.27	48.46
zhong	61.13	47.28	61.14	44.66	35.19	65.21	48.07
kummerfeld	62.72	42.10	60.29	45.35	35.32	59.91	47.10
zhifeng	48.29	24.58	61.46	44.83	35.75	53.77	48.43
train	26.07	19.98	50.46	31.68	25.21	51.12	31.28

Scores on various corpora described in Raghunathan et al. 2010

	MUC			B cubed			Pairwise		
	P	R	F1	P	R	F1	P	R	F1
ACE2004 dev	86.0	75.5	80.4	89.3	76.5	82.4	81.7	55.2	65.9
ACE2004 test	82.7	70.2	75.9	88.7	74.5	81.0	77.2	44.6	56.6
ACE2004 mwire	84.6	75.1	79.6	87.3	74.1	80.2	79.4	50.1	61.4
MUC6 test	90.6	69.1	78.4	90.6	63.1	74.4	89.7	57.0	69.7

# Stanford CoreNLP (2) - démo

<http://nlp.stanford.edu:8080/corenlp/process>

**Stanford CoreNLP**

Output format:

Please enter your text here:

Though Apple stock fell after Tuesday's iPad mini announcement with a \$329 starting price, analysts are largely sticking by their price targets for AAPL stock. Chris Whitmore of Deutsche Bank said pricing of the iPad mini was a "modest disappointment," but he still expects the product to perform well.

Envoyer

```

<?xml version="1.0" encoding="UTF-8"?>
<?xml:stylesheet href="CoreNLP-to-HTML.xsl" type="text/xsl"?>
<root>
<document>
<sentences>
<sentence id="1">
<tokens>
<token id="1">
<word>Though</word>
<lemma>though</lemma>
<CharacterOffsetBegin>0</CharacterOffsetBegin>
<CharacterOffsetEnd>6</CharacterOffsetEnd>
<POS>IN</POS>
<NER>O</NER>
</token>
<token id="2">
<word>Apple</word>
<lemma>Apple</lemma>
<CharacterOffsetBegin>7</CharacterOffsetBegin>
<CharacterOffsetEnd>12</CharacterOffsetEnd>
<POS>NNP</POS>
<NER>O</NER>
</token>
<token id="3">
<word>stock</word>
<lemma>stock</lemma>
<CharacterOffsetBegin>13</CharacterOffsetBegin>
<CharacterOffsetEnd>18</CharacterOffsetEnd>
<POS>NN</POS>
<NER>O</NER>
</tokens>

```

# Stanford CoreNLP (3)

Though Apple stock fell after Tuesday's iPad mini announcement with a \$329 starting price, analysts are largely sticking by their price targets for AAPL stock. Chris Whitmore of Deutsche Bank said pricing of the iPad mini was a "modest disappointment," but he still expects the product to perform well.

**Part-of-Speech:**

IN NNP NN VBD IN NNP POS NNP NNS NN IN DT DOLLAR CD VBG

1. Though Apple stock fell after Tuesday 's iPad mini announcement with a \$ 329 starting price, analysts are largely sticking by their price targets for AAPL stock.

NN NNS VBP RB VBG IN PRPS NN NNS IN NN NN

NNP NNP IN NNP NNP VBD NN IN DT NNP NNS VBD DT JJ NN CC

2. Chris Whitmore of Deutsche Bank said pricing of the iPad mini was a `modest`disappointment,`but he still expects the product to perform well.

PRP RB VBZ DT NN TO VB RB

**Named Entity Recognition:**

Date MONEY

1. Though Apple stock fell after Tuesday's iPad mini announcement with a \$329 starting price, analysts are largely sticking by their price targets for AAPL stock.

Org

Person Organization

2. Chris Whitmore of Deutsche Bank said pricing of the iPad mini was a `modest`disappointment,`but he still expects the product to perform well.

# Stanford CoreNLP (3)

Though Apple stock fell after Tuesday's iPad mini announcement with a \$329 starting price, analysts are largely sticking by their price targets for AAPL stock. Chris Whitmore of Deutsche Bank said pricing of the iPad mini was a "modest disappointment," but he still expects the product to perform well.

**Part-of-Speech:**

IN NNP NN VBD IN NNP POS NNP NNS NN IN DT DOLLAR CD VBG

**Coreference:**

1. Though Apple stock fell after Tuesday's iPad mini announcement with a \$329 starting price, analysts are largely sticking by their price targets for AAPL stock.

-----Coref-----Ment

2. Chris Whitmore of Deutsche Bank said pricing of the iPad mini was a `modest`disappointment,`but he still expects the product to perform well.

Ment Ment Ment

-----Coref-----

M

-----Coref-----

Person Organization

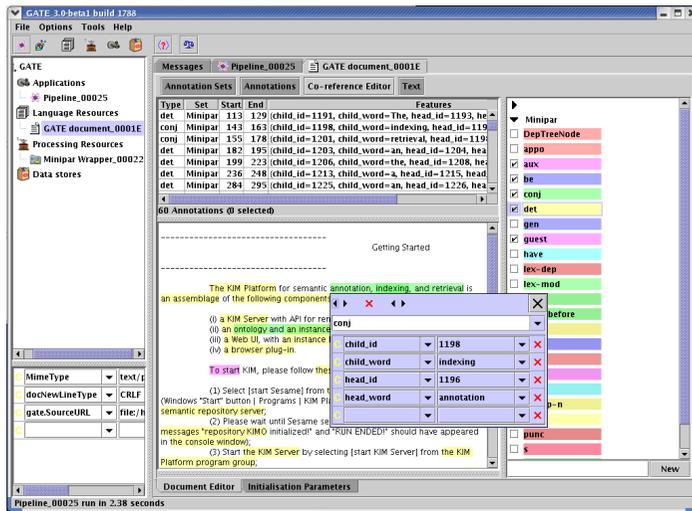
2. Chris Whitmore of Deutsche Bank said pricing of the iPad mini was a `modest`disappointment,`but he still expects the product to perform well.



# MiniPar Parser

<http://gate.ac.uk/sale/tao/splitch17.html>

- Partie du discours (POS), lemme, tête, dépendances



# LIRMM - Sygmart

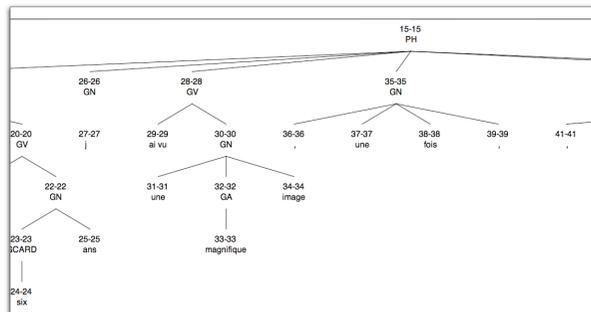
- Je relis les énoncés et trouve un exemple pour chaque utilisation.

<http://www.lirmm.fr/~chauche/ExempleAnl.html>



# LIRMM : Sygmart

- Analyse complète du Petit Prince ([http://www.sygtext.fr/Le\\_petit\\_prince.svg](http://www.sygtext.fr/Le_petit_prince.svg))



# Ressources linguistiques

PARIS DESCARTES

Le propos de cet ouvrage est de fixer le peu géographique de nos auteurs...  
**Lexique 3**  
 Un site réalisé par Boris New & Christophe Pallier et hébergé par le RISC

Menu principal  
 Accueil  
 Interroger Open Lexique  
 Télécharger  
 Documentation  
 Autour de Lexique  
 Autres outils

Base de données textuelles

- Elda: Agence pour l'évaluation et la distribution des ressources linguistiques.**  
Elda propose des bases de données (BDD) très souvent payantes (et assez chères)
- Novlex : une base de données lexicales pour des élèves de primaire**  
Novlex estime l'étendue et la fréquence lexicale du vocabulaire écrit adressé à des élèves francophones de l'enseignement primaire
- Fréquence de bigrammes en français**  
Les fréquences des bigrammes (deux lettres successives) sont présentées ici en français et calculées sur 5 types de textes
- Fréquence de journaux**  
Fréquence de journaux
- Fréquences orales**  
Ce fichier comporte les fréquences des formes (non lemmatisées) dans un corpus de français parlé d'un million de mots (Corpax, version mai 2000).
- Normes pour 299 images destinées aux protocoles en psycholinguistique**  
Ce site permet de visualiser les images et de connaître les normes correspondantes.
- Psycholinguistic norms for action photographs in French**  
Normes pour des images représentant des actions
- Norms for 866 french words**  
Fichier répertoriant les normes de concrétude, d'imagerie... de 866 mots français
- Behind the name : site très complet sur les prénoms**  
Ce site répertorie un ensemble de prénoms et les classe selon plusieurs critères (popularité du prénom, année)
- Vocolex : Base de Données Lexicales sur les similarités phonologiques entre les mots Français**  
Cette base fournit un ensemble d'indicateurs statistiques sur les similarités entre mots de la langue française
- Lexique Morphalou : Lexique ouvert des formes fléchies du français**  
Lexique ouvert des formes fléchies du français
- Lexop : BDD lexicale**  
Propositions de statistiques sur l'orthographe et la phonologie
- Omnilex : BDD sur le lexique français contemporain**  
Calculs de statistiques descriptives sur le lexique

<http://www.lexique.org>

33

EARIA-2012 ARIA

## Dictionnaire de formes fléchies : Morphalou

CNRTL Centre National de Ressources Textuelles et Lexicales

■ **Morphalou**

Le lexique *Morphalou* est un lexique ouvert des formes fléchies du français. Les données initiales de *Morphalou* proviennent du *TLFName*, la nomenclature du *Trésor de la Langue Française* qui a fourni 539.413 formes fléchies, appartenant à 68.075 lemmes. Le transfert du *TLFName* vers *Morphalou* s'est fait par une réorganisation structurale des données et une normalisation des étiquettes grammaticales, sans perte d'informations linguistiques. Le lexique résultant est un lexique à large couverture (~540.000 formes fléchies), linguistiquement valide (sous la responsabilité d'un comité éditorial) et formellement en accord avec les propositions de normalisation pour les ressources lexicales du TAL à l'ISO (TC37/SC4). Il est en accès libre à des fins de recherche et d'enseignement. Le maintien et la mise à jour du lexique sont assurés par l'ATILF.

Origine de la ressource : ATILF (Nancy Université - CNRS)  
 Nature des données : Lexique morphologique  
 Origine des données : ATILF - TLFName 0.95 (Nomenclature du TLF)  
 Soutien institutionnel : Contrat Plan-Etat Région : Ingénierie des Langues, du Document et de l'Information Scientifique, Technique et Culturelle

Documentation sur Morphalou 2.0  
 Télécharger Morphalou 1.0 | Télécharger Morphalou 2.0

■ **Fiche technique**

Version	Morphalou 2.0
Conception	Susanne Salmon-Ait, Laurent Romary, Jean-Marie Pierrrel
Responsable scientifique	Susanne Salmon-Ait
Responsable informatique	Etienne Petitjean
Contenu	Nombre de lemmes 65 810 Nombre de formes fléchies 524 725
Format	XML
Codage des caractères	UTF-8
Taille	195 Mo

34

EARIA-2012 ARIA

## Dictionnaire de formes fléchies : Morphalou

CNRTL Centre National de Ressources Textuelles et Lexicales

■ **Morphalou**

Le lexique *Morphalou* est un lexique ouvert des formes fléchies du français. Les données initiales de *Morphalou* proviennent du *TLFName*, la nomenclature du *Trésor de la Langue Française* qui a fourni 539.413 formes fléchies, appartenant à 68.075 lemmes. Le transfert du *TLFName* vers *Morphalou* s'est fait par une réorganisation structurale des données et une normalisation des étiquettes grammaticales, sans perte d'informations linguistiques. Le lexique résultant est un lexique à large couverture (~540.000 formes fléchies), linguistiquement valide (sous la responsabilité d'un comité éditorial) et formellement en accord avec les propositions de normalisation pour les ressources lexicales du TAL à l'ISO (TC37/SC4). Il est en accès libre à des fins de recherche et d'enseignement. Le maintien et la mise à jour du lexique sont assurés par l'ATILF.

Origine de la ressource : ATILF (Nancy Université - CNRS)  
 Nature des données : Lexique morphologique  
 Origine des données : ATILF - TLFName 0.95 (Nomenclature du TLF)

```

<lexicalEntry id="championne_1">
  <feminineVariantOf target="champion_1">champion</feminineVariantOf>
  <formSet>
    <lemmatizedForm>
      <orthography>championne</orthography>
      <grammaticalCategory>commonNoun</grammaticalCategory>
      <grammaticalGender>feminine</grammaticalGender>
    </lemmatizedForm>
    <inflectedForm>
      <orthography>championne</orthography>
      <grammaticalNumber>singular</grammaticalNumber>
    </inflectedForm>
    <inflectedForm>
      <orthography>championnes</orthography>
      <grammaticalNumber>plural</grammaticalNumber>
    </inflectedForm>
  </formSet>
  <originatingEntry target="TLF">CHAMPION, ONNE, subst.</originatingEntry>
</lexicalEntry>
  
```

34

EARIA-2012 ARIA

<http://www.cnrtl.fr/lexiques/prolex/>

## Lexique multilingue de noms propres : Prolex

CNRTL Centre National de Ressources Textuelles et Lexicales

■ **Prolex**

Le projet Prolex, p du traitement au journalistiques, de propres (Prolexba: La ressource Prolex • Le groupe de • L'université C

Ce projet a reçu le • De l'action T • Du programm

■ **Prolexbase**

La modélisation de pivot ne représente qui est une famille éponymique, etc.) C Il n'est pas évider sémantique et une large incluant ce souvent de la corr comme une expo proche de celle ut

■ **Fiche techniq**

Prolex

Recherche Texte Commentaire Glossaire Contact Remerciement

Votre recherche  
 italie

Choix de la langue  
 fra

Rechercher

Recherche avancée

Historique  
 italie  
 holland

Italie

Type : Pays  
 Existence : Historique  
 Nom relationnel : Italien  
 Adj. Relationnel : italien  
 Synonyme : République italienne  
 Holonyme : Europe, Organisation des nations unies  
 Accessibilité : Rome (Capitale)  
 Classifiant : pays  
 Traduction : Italien (allemand), Italia (italien), Itália (portugais), Italia (espagnol), Italië (hollandais), Italia (anglais), Italija (serbe), 이탈리아 (coréen)

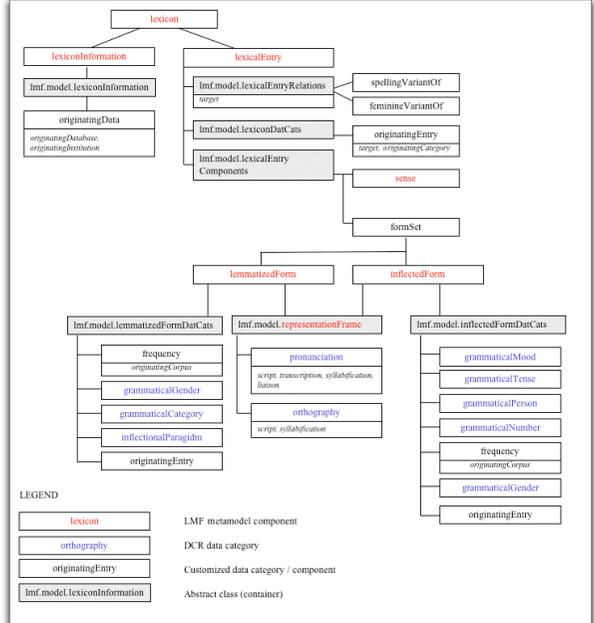
Encyclopédie : Wikipédia

Plus de détails  
 Résultat contenant la séquence

campagnes d'Italie  
 guerres d'Italie  
 Humbert Ier d'Italie

Contenu : Nombre de proxèmes : 54 774 (4 568 antroponymes, 49 789 toponymes, 175 ergonimes et 222 pragmonymes)  
 Nombre de relations : 33 512 (2 694 accessibilités, 49 970 méronymes et 649 synonymes)  
 Nombre de lemmes : 76 118  
 Nombre de formes fléchies : 124 721

35



## Lexiques liés à la subjectivité, aux opinions...

- Anglais :
  - The MPQA (Multi-Perspective Question Answering) Subjectivity Lexicon
  - SentiWordNet (valeurs >0 ou <0 aux synsets de WordNet)

#	POS	ID	PosScore	NegScore	SynsetTerms	Gloss
a		00001740	0.125	0	able#1	(usually followed by 'to') having the necessary means or [...]
a		00002098	0	0.75	unable#1	(usually followed by 'to') not having the necessary means or [...]
a		00002312	0	0	dorsal#2 abaxial#1	facing away from the axis of an organ or organism; [...]
a		00002527	0	0	ventral#2 adaxial#1	nearest to or facing toward the axis of an organ or organism; [...]
a		00002730	0	0	acroscopic#1	facing or on the side toward the apex
a		00002843	0	0	basiscopic#1	facing or on the side toward the base
a		00002956	0	0	abductive#1 abducent#1	especially of muscles; [...]
a		00003131	0	0	adductive#1 adducting#1 adductor#1	especially of muscles; [...]
a		00003356	0	0	nascent#1	being born or beginning; [...]
a		00003553	0	0	emerging#2 emergent#2	coming into existence; [...]

- Espagnol : <http://lit.csci.unt.edu/~rada/downloads/SpanishSentimentLexicons.tar.gz>

## WordNet

<http://wordnet.princeton.edu/>

POS	Unique Strings	Synsets	Total Word-Sense Pairs
Noun	117097	81426	145104
Verb	11488	13650	24890
Adjective	22141	18877	31302
Adverb	4601	3644	5720
Totals	155527	117597	207016

Table 1. Semantic Relations in WordNet

Semantic Relation	Syntactic Category	Examples
Synonymy (similar)	N, V, Adj, Av	pipe, tube rise, ascend sad, unhappy rapidly, speedily
Antonymy (opposite)	Adj, Av, (N, V)	wet, dry powerful, powerless friendly, unfriendly rapidly, slowly
Hyponymy (subordinate)	N	sugar maple, maple maple, tree tree, plant
Meronymy (part)	N	brim, hat gin, martini ship, fleet
Troponymy (manner)	V	march, walk whisper, speak
Entailment	V	drive, ride divorce, marry

Note: N = Nouns, Adj = Adjectives, V = Verbs, Av = Adverbs

WordNet: A Lexical Database for English  
George A. Miller

COMMUNICATIONS OF THE ACM November 1995/Vol. 38, No. 11

## FrameNet



- Hypothèse : les mots peuvent être compris selon la manière (structure conceptuelle) dont ils sont employés
- >> 1000 structures sémantiques conceptuelles décrivant des relations entre participants
- >> 170 000 phrases étiquetées --> semantic role labeling
- Application aux tâches de questions-réponses, d'extraction d'information, de détection d'inférence (textual entailment)

WOLF : Wordnet Libre du Français <http://alpage.inria.fr/~sagot/wolf.html>



## Entités nommées, bases de connaissances

## Conférences MUC

- MUC (Message Understanding)
  - 7 catégories d'entités  
*people, organization, location, time, date, money, percentage*
  - peu réaliste pour des tâches plus complexes (ex. questions-réponses)

[http://www.itl.nist.gov/iaui/894.02/related\\_projects/muc/proceedings/muc\\_7\\_toc.html](http://www.itl.nist.gov/iaui/894.02/related_projects/muc/proceedings/muc_7_toc.html)

**Table 2: Maximum Results Reported in MUC-3 through MUC-7 by Task**

Evaluation/Tasks	Named Entity	Coreference	Template Element	Template Relation	Scenario Template	Multilingual
MUC-3					R < 50% P < 70%	
MUC-4					F < 56%	
MUC-5					EJV F < 53% EME F < 50%	JJV F < 64% JME F < 57%
MUC-6	F < 97%	R < 63% P < 72%	F < 80%		F < 57%	
MUC-7	F < 94%	F < 62%	F < 87%	F < 76%	F < 51%	
Multilingual						
MET-1	C F < 85% J F < 93% S F < 94%					
MET-2	C F < 91% J F < 87%					

Legend:  
 R = Recall P = Precision F = F-Measure with Recall and Precision Weighted Equally  
 E = English C = Chinese J = Japanese S = Spanish  
 JV = Joint Venture ME = Microelectronics

## Evaluations ESTER

- Transcriptions audio (300 h. transcrites, 1600 h. non transcrites)
- 7 catégories (*persons, locations, organizations, human products, amounts, time and functions*) et 38 sous-catégories

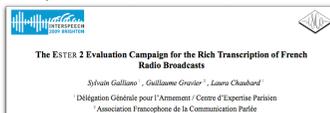
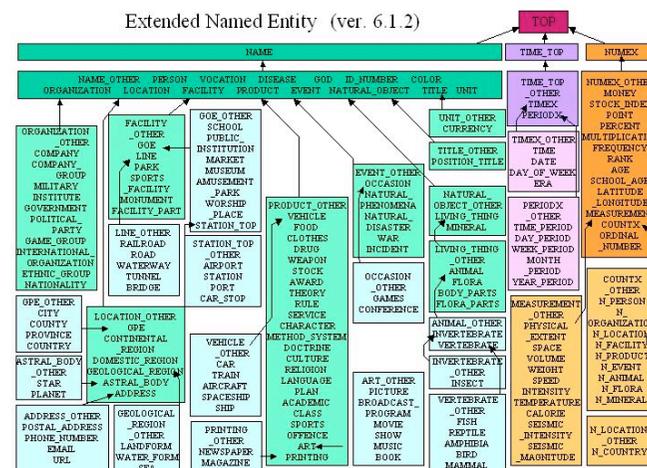


Table 6: NE task overall performance for each participating site (Slot Error Rate, Precision, Recall and F-measure)

%WER	ref. transcript			LIMSI transcript 12.11			LIA transcript 17.83			IRISA transcript 26.09		
	%S	%P	%R	%S	%P	%R	%S	%P	%R	%S	%P	%R
Sites	23.9	86.46	71.85	<b>43.4</b>	79.52	59.45	<b>51.6</b>	76.51	55.02	<b>56.8</b>	72.26	49.02
LIA	30.9	81.15	70.94	45.3	75.13	62.33	55.5	70.50	57.52	61.2	66.13	50.67
LINA	37.1	80.75	55.48	54.0	71.98	44.01	60.4	68.76	40.84	65.2	63.66	35.66
LI Tours	33.7	79.39	65.82	50.7	71.36	54.16	80.8	56.59	46.46	82.9	51.28	42.38
LSIS	35.0	82.65	73.07	55.3	70.23	58.39	86.5	70.36	28.66	88.6	67.03	25.22
Synapse	9.9	93.02	89.37	44.9	76.39	67.16	60.7	70.26	59.21	66.2	65.95	52.71
Xerox	<b>9.8</b>	<b>93.61</b>	<b>91.50</b>	44.6	58.91	70.06	—	—	—	—	—	—

## La hiérarchie de Sekine

- Plus de 200 types



<http://nlp.cs.nyu.edu/ene/>

ENE	Examples	
Name Other	Barbero, Bubbles, Max, Maggie	
Person	Bush, Michael Jackson, Elizabeth II, LeBron Raymone James	
God	Zeus, Indra, Danu, Ra	
Organization	Organization Other	the Capone Family, Department of Computer Science, CS Dept., general affairs department
	Informational Organization	UN, League of Nations, Pacific Island Forum, SEATO
	Show Organization	The Cleveland Orchestra, The Beatles, the Bolshoi Ballet troupe, Sex Pistols
	Family	The House of Hamilton, Clan Henderson, Tokugawa clan, Koga family
	Ethnic Group	Ethnic Group Other: White people, Jew, Slavic peoples, Mongoloid race, Japanese Diaspora
	Nationality	Japanese, Israeli, American, American people
	Sports Organization Other	the Breen Gym, UCLA Bruins, Ma family army, Shinagawa jogging Club
	Pro Sports Organization	New York Yankees, Seattle, NYY, Manchester United
	Sports League	NFL, National Basketball Association, Atlantic Coast Conference, National League West
	Corporation Other	Association for Computational Linguistics, National Rifle Association, NHK, BBC
Political Organization	Company Group	Tata Group, J.R. the Big Three, Big Four auditors
	Political Organization Other	Palestine Liberation Organization, Clinton Regime, Tokugawa shogunate, Ayyubid dynasty
	Government	National Security Council, Ministry of Finance, the United States Senate, USTR
	Political Party	Democratic Party, Bharatiya Janata Party, Conservative Party, LDP
	Cabinet	Tatcher's Cabinet, Major's Cabinet, Tanaka's Cabinet, Kozumi's Cabinet
	Military	Self-Defense Forces, US Air Force, Royal Navy, UN forces
	Location Other	Times Square, Ground Zero, Three Views of Japan, Garden of Eden
	Sea	Hakone Spa, Fukuchi Spa, Hakuba Spa, Yunoyama Spa
	GPE	GPE Other: Taiwain, Hong Kong, Puerto Rico, French Polynesia, Macau
	Region	City: New York City, Brooklyn, Sydney, Rio de Janeiro County: West Chester County, Madison County, Orange County, Shima District Prefecture: Osaka Prefecture, NY, Kansas, Nova Scotia, Nagorno-Karabakh Country: the United States, Japan, UK, Vatican City
Location	Region Other	Continental Region: North America, Asia, the Caribbean area, NIES Domestic Region: New England, East Coast, the South, Upper New York Geological Region Other: Grand Canyon, Alamogordo, Great Barrier Reef, Ayers Rock
	Mountain	Mount Everest, K2, Mt. Fuji, Alps
	Island	Florida Keys, Key West, Gilbert Islands, Iriomote
	River	Mississippi River, Hudson River, Yangtze River, Danube
	Lake	Lake Michigan, Lake Baikal, Dead Sea, Great Lakes
	Sea	Pacific Ocean, Sea of Japan, Sunda Strait, English Channel
	Bay	Bay of Bengal, Delaware Bay, Persian Gulf, Gulf of Guinea
	Astral Body Other	Andromeda Galaxy, Solar System, Halley's Comet, Callisto
	Astral Body	Star: Antares, Sirius, North Star, Barnard's Star Sun, Earth, Moon, Io, Juno Constellation: Taurus, Cassiopeia, Argo-Navis, Lepus
	Address Other	Address: 715 Broadway, 7th floor, New York, NY 10003 USA, 715 Broadway, 10003 Phone Number: (212) 123-4567, 911, ext445 Email: sekine@cs.nyu.edu URL: http://nlp.cs.nyu.edu/sekine/index.jp.html
Facility Other	Facility: Empire State Building, Cooper Dam, Fulmar oilfield, Eiffel Tower	
Facility Part	8th floor, room #1204, second basement, Runway 13R-31	
Archaeological Place Other	Archaeological Ruins at Moenjodaro, Cahokia Mounds State Historic Site, Angkor, Masada	
Archaeological Place	Daisen-Kofun, Zhaoqing, Ringlemere barrow, Great Pyramid of Giza	

- Plus d'un million de mots annotés
- Entités et composants

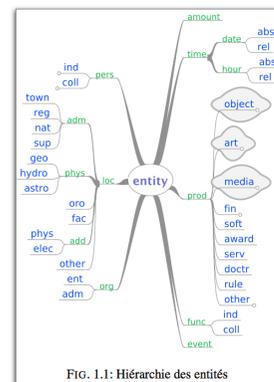


FIG. 1.1: Hiérarchie des entités

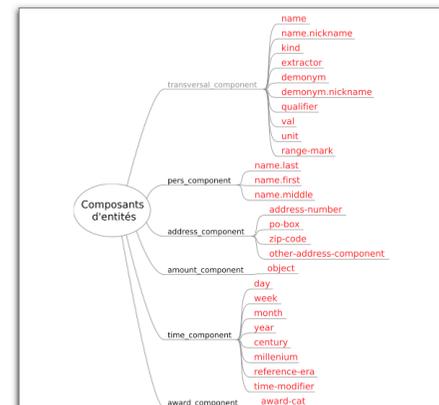


FIG. 1.2: Composants. Les noms de groupes de composants (ex : address\_component) ne sont pas utilisés dans l'annotation, ils sont seulement présents pour structurer la description des composants.

exemples de structures (analyses syntaxiques) d'entités

- ville de Paris <loc.adm.town>
  - <kind> ville
  - de <name> Paris
- la société Peugeot <org.ent>
  - <kind> société
  - <name> Peugeot
- la Mairie de Paris <org.adm>
  - de <kind> mairie
  - <loc.adm.town> Paris
  - <org.adm> de
  - <loc.adm.town> <name> Paris
- Paris <loc.adm.town> <name> Paris
- trois cuillérées de farine de châtaignes <amount>
  - <val> trois
  - de <unit> cuillérées
  - <object> de farine de c
- une douzaine de tomates <amount>
  - <val> une douzaine
  - de <object> tomates

Entités nommées structurées : guide d'annotation Quaero

Sophie Rosset, Cyril Grosin, Pierre Zweigenbaum

imsi

Centre National de la Recherche Scientifique

FIG. 1.3: Exemples d'entités.

Rosset, Grosin, Zweigenbaum (2011) <http://quaero.org/media/files/bibliographie/quaero-guide-annotation-2011.pdf>



- Approche statistique (Markov maximum entropie, CRF)

Jenny Rose Finkel, Trond Grenager, and Christopher Manning. 2005. Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling. Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL 2005), pp. 363-370.
- Evaluation dans le cadre de CoNLL *Conference on Computational Natural Language Learning* <http://www.cnts.uu.ac.be/conll2003>

types reconnus : persons, organizations, locations, times and quantities

Name	Corpus	Language	# Word Tokens		# Entities		# Features		Exact Match Score (conlleval)			Classifier
			Train	Test	Types	Instances	$\Phi(X)$	$\Lambda/(X,Y)$	Prec	Rec	F <sub>1</sub>	
CoNLL 2002	Dutch news testa (devset)		218737	37761	4	2616	838524	4192620	78.99%	77.33%	78.15%	pure CMM
CoNLL 2002	Dutch news testb		218737	68994	4	3941	838559	4192795	80.48%	78.96%	79.71%	pure CMM
CoNLL 2002	Spanish news testa (devset)		273037	52923	4	4352	776511	3882555	78.01%	76.19%	77.09%	pure CMM
CoNLL 2002	Spanish news testb		273037	51533	4	3559	776444	3882220	81.24%	81.03%	81.14%	pure CMM
CoNLL 2003	English news testa (devset)		219553	51578	4	5942	738378	3691890	91.37%	91.22%	91.29%	pure CMM
CoNLL 2003	English news testa (devset)		219554	51578	4	5942			92.15%	92.39%	92.27%	postprocessed CMM
CoNLL 2003	English news testb		219553	46666	4	5648	738378	3691890	85.65%	85.41%	85.53%	pure CMM
CoNLL 2003	English news testb		219554	46666	4	5648			86.12%	86.49%	86.31%	postprocessed CMM
CoNLL 2003	German news testa (devset)		220189	51645	4	4833	1079044	5395220	77.12%	61.37%	68.35%	pure CMM
CoNLL 2003	German news testa (devset)		220189	51645	4	4833			75.36%	60.36%	67.03%	postprocessed CMM
CoNLL 2003	German news testb		220189	52098	4	3673	1079037	5395185	79.23%	63.65%	70.59%	pure CMM
CoNLL 2003	German news testb		220189	52098	4	3673			80.38%	65.04%	71.90%	postprocessed CMM
CoNLL 2003	English news testa (devset)		219553	51578	4	5942	616918	11532202	91.64%	90.93%	91.28%	CRF (closed task)
CoNLL 2003	English news testa (devset)		219553	51578	4	5942	633786	12285708	93.28%	92.71%	92.99%	CRF (with distsm)
CoNLL 2003	English news testb		219553	46666	4	5648	633786	12285708	88.21%	87.68%	87.94%	CRF (with distsm)

# Stanford NER (2)

Feature	NER	TF
Current Word	✓	✓
Previous Word	✓	✓
Next Word	✓	✓
Current Word Character n-gram	all	length ≤ 6
Current POS Tag	✓	
Surrounding POS Tag Sequence	✓	
Current Word Shape	✓	✓
Surrounding Word Shape Sequence	✓	✓
Presence of Word in Left Window	size 4	size 9
Presence of Word in Right Window	size 4	size 9

Table 2: Features used by the CRF for the two tasks: named entity recognition (NER) and template filling (TF).

Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling  
**Jenny Rose Finkel, Trond Grenager, and Christopher Manning**  
 Computer Science Department Stanford University  
<http://nlp.stanford.edu/~manning/papers/gibbscrf3.pdf>

CoNLL					
Approach	LOC	ORG	MISC	PER	ALL
B & M LT-RMN	-	-	-	-	80.09
B & M GLT-RMN	-	-	-	-	82.30
Local+Viterbi	88.16	80.83	78.51	90.36	85.51
NonLoc+Gibbs	88.51	81.72	80.43	92.29	86.86

Table 5: F1 scores of the local CRF and non-local models on the CoNLL 2003 named entity recognition dataset. We also provide the results from Bunescu and Mooney (2004) for comparison.

U.N.	NNP	I-NP	I-ORG
official	NN	I-NP	O
Ekeus	NNP	I-NP	I-PER
heads	VBZ	I-VP	O
for	IN	I-PP	O
Baghdad	NNP	I-NP	I-LOC
.	.	O	O

# Détection des entités nommées : un système symbolique open source à la disposition de la communauté

Nathalie Friburger, Denis Maurel

## A l'origine

2000-2001 création de Casys

- Système permettant de gérer une cascade de transducteurs (aujourd'hui intégré à Unix)
- NER limitée aux personnes, lieux, organisations

2007-2010 participation de CasEN à Esterz, EPAC, Varling

- Extension aux entités : montants, événements, fonctions, productions humaines

## Pourquoi réécrire CasEN ?

- De nombreuses modifications apportées par des personnes différentes
- Maintenance difficile
- Nouvelles fonctionnalités d'Unix : morphologie, contextes à exclure ou pas
- Modification du fonctionnement de CasSys / intégration à Unix

**Refonte de CasEN**

- Architecture de la cascade
- Nommage des graphes
- Écriture des graphes
- Outils de débogage

## Différents types de transducteurs

Fonction des transducteurs identifiable par le préfixe de leur nom

- pers, loc, org, event, func, prod, amount, time : reconnaissance des entités nommées correspondant au type indiqué (peuvent assembler plusieurs transducteurs)
- ex : orgCommerceEtranger reconnaît des organisations commerciales ayant un nom étranger
- ex : amount reconnaît différentes mesures (monnaie, température, longueur, etc.)
- bool : transducteur booléen
- Ex : toolChercheSignalAvecPoints recherche les sigles contenant des points (C.G.T. au lieu de CGT aujourd'hui).
- list : transducteur de listes
- pattern : transducteurs masqués, décrivent des motifs utilisant des codes Unix pour sélectionner des séquences de mots (contraintes morphologiques, syntaxiques, lexicales, etc.)
- tag : transducteur étiquette, permettent d'insérer des étiquettes dans des entités imbriquées



**Friburger N., Maurel D.** (2004), Finite-state transducer cascade to extract named entities in texts, *Theoretical Computer Science*, vol. 313, 94-104.

**Maurel D., Friburger N., Antoine J.-Y., Eshkol-Taravella I., Nouvel D.** (2011), Cascades autour de la reconnaissance des entités nommées, *TAL* 52-1.

Ex : <pers+hum> → ((Gerhard\_N+Prénom) [Aligner\_N+nom], entity+pers+hum+grfPersPrenomNom), le {secrétaire général\_entity+fonc+pol} de l'UEFA\_entity+org+grfOrgMetier), à dans ses placards des projets de {ligue européenne\_entity+org+grfOrgDivers}.

Ex : <org> de <org> → ((Gerhard\_N+Prénom) [Aligner\_N+nom], entity+pers+hum+grfPersPrenomNom) {secrétaire général\_entity+fonc+pol} de l'UEFA\_entity+org+grfOrgMetier).

# CasEN

[http://tln.li.univ-tours.fr/Tln\\_CasEN.html](http://tln.li.univ-tours.fr/Tln_CasEN.html)

personne (pers)	humain réel ou fictif (pers.hum) animal réel ou fictif (pers.anim)	civilité (pers.hum.titre) titre professionnel (pers+hum+func) adjectif ethnique (pers+hum+ethnic) gentilé et adjectifs toponymiques (pers+hum+gent) nationalité (pers+hum+nat) dynastie (pers+hum+dyn)
fonction (func)	politique (func.pol) militaire (func.mil) administrative (func.adm) religieuse (func.rel) artistique (func.art)	
organisation (org)	politique (org.pol) éducative (org.edu) commerciale (org.com) média & divertissement (org.div)	
lieu (loc)	géographique naturel (loc.geo) axe de circulation (loc.line) construction humaine (loc.fac) région administrative (loc.adm) ville (loc.adm.ville) adresse (loc.addr) adresse postale (loc.addr.post) téléphone et fax (loc.addr.tel) adresse électronique (loc.addr.elec)	
production humaine (prod)	produit (prod.obj) récompense (prod.award) œuvre artistique (prod.art) production documentaire (prod.doc)	
date et heure (time)	date (time.date) heure (time.hour) adverbe de date (time.advdate)	date absolue (time.date.abs) date relative (time.date.rel)
montant (amount)	valeur physique (amount.phy) valeur monétaire (amount.cur) ocets (amount.computer)	durée (amount.phy.dur) température (amount.phy.temp) longueur (amount.phy.len) surface et aire (amount.phy.area) volume (amount.phy.vol) poids (amount.phy.wgt) vitesse (amount.phy.spd) autre (amount.phy.other)
événement (event)	histoire (event+hist) célébration (event+cer) fête (event+feat) manifestation (event+manif) météorologie (event+meteo)	

# CasEN

[http://tln.li.univ-tours.fr/Tln\\_CasEN.html](http://tln.li.univ-tours.fr/Tln_CasEN.html)

personne (pers)	humain réel ou fictif (pers.hum)	civilité (pers.hum.titre)
Donnons quelques exemples de reconnaissance:		
func	« Au pire de la crise, (à l'automne dernier, entity+time+date+rel+grfTimeDateRelative), nous avons détenu jusqu'à 20 % de liquidités dans notre portefeuille », indique (( Denis_N+Prénom ) ( Remacle_N+nom ), entity+pers+hum ), ( gérant d'Amplitude Pacifique, entity+org+com ), entity+job ) , entity+pers+hum+grfPersPrenomNom ), une sicav de ( La Poste, entity+org+com+grfOrgDico ).	
org	« C'est à nos clients de décider s'ils souhaitent ou non consacrer une partie de leur patrimoine à l' (Asie, entity+loc+adm+grfLocPays ) », souligne (( Pierre_N+Prénom ) ( Cret_N+nom ), entity+pers+hum+grfPersPrenomNom ), de la ( Compagnie financière (Edmond_N+Prénom ) (de Rothschild, N+nom ), entity+pers+hum ), entity+org+com+grfOrgCommerceGauche ).	
lieu	Ils ne peuvent pas, en revanche, faire l'impasse sur la ( Bourse de Hongkong ), entity+loc+adm ), entity+org+com+grfOrgCommerceGauche ), car cette place représente près de la moitié de la capitalisation boursière de la région. (S) Pour sa part, (( Pierre-Alexis_N+Prénom ) ( Dumont_N+nom ), entity+pers+hum+grfPersPrenomNom ), de ( State Street Banque, entity+org+com+grfOrgCommerceDroite ), s'est réfugié sur le marché australien, relativement épargné par la tourmente.	
proc	( Théâtre Gérard-Philipe, entity+org+div+grfOrgDivertissementSorties ), ( 59, ( boulevard Jules-Guesde, entity+loc+line ), 93000 ( Saint-Denis, entity+loc+ville ), entity+loc+addr+post+grfLocAddr ).	
date	Selon une étude de l' ( Autorité de régulation des télécommunications, entity+org+grfOrgDivers ) (( ART, entity+org+grfOrgSuiVDeParentheses ), le taux d'équipement devrait dépasser les 50 % ( en 2002, entity+time+date+abs+grfTimeAnneeSiecle ).	
dat	Une évaluation a été réalisée sur un extrait du Journal Le Monde daté du 1er janvier 1999, 7 articles, soit 7 070 mots pour 83.2 ko. En voici les résultats :	
mor	<ul style="list-style-type: none"> <li>• Sur 582 entités nommées présentes, 432 ont été reconnues, avec, en plus, 11 fausses reconnaissances, soit : <ul style="list-style-type: none"> <li>• Précision : 97,52%</li> <li>• Rappel : 74,23%</li> </ul> </li> <li>• Sur les entités reconnues, 390 types étaient corrects et 3 types comportaient une erreur due à la métonymie (par exemple France comme toponyme au lieu d'organisation). <ul style="list-style-type: none"> <li>• Sans compter les erreurs dues à la métonymie : <ul style="list-style-type: none"> <li>• Précision : 89,71%</li> <li>• Rappel : 67,53%</li> </ul> </li> <li>• En comptant les erreurs dues à la métonymie : <ul style="list-style-type: none"> <li>• Précision : 88,04%</li> <li>• Rappel : 67,01%</li> </ul> </li> </ul> </li> <li>• Sur les entités reconnues et bien typées, 384 sous-types étaient corrects, soit : <ul style="list-style-type: none"> <li>• Précision : 85,52%</li> <li>• Rappel : 65,98%</li> </ul> </li> </ul>	
évé	Enfin, à nouveau sur les entités reconnues, 393 étaient correctement balisées, 20 avaient un début défectueux et une fin correcte, 16 un début correct et une fin défectueuse, 3 un début et une fin défectueux, soit : <ul style="list-style-type: none"> <li>• Précision : 90,97%</li> <li>• Rappel : 67,53%</li> </ul>	

# Entités et propriétés

- Infoboxes de Wikipédia
- Ontologies YAGO, SUMO...

Table 5: Size of other ontologies

Ontology	# Entities	# Facts
SUMO [114]	20,000	60,000
Ponzzetto et al. [118]	n/a	110,000
WordNet [66]	117,659	821,492
Cyc [106]	300,000	3,000,000
TextRunner [13]	n/a	7,800,000
YAGO	1,700,000	15,000,000
DBpedia [9]	1,950,000	103,000,000

Table 4: Largest relations in YAGO

Relation	# Facts	Relation	# Facts
hasUTCOffset	12724	hasWonPrize	13645
livesIn	15185	writtenInYear	16441
originatesFrom	18876	directed	18633
hasProcessor	19134	actedIn	22249
hasDuration	22652	bornInLocation	24400
hasMdb	24659	hasArea	26781
hasProductionLanguage	27840	produced	30519
hasPopulation	30731	isOfGenre	33898
hasSuccessor	46658	establishedOnDate	68629
hasNobility	20779	created	68627
locatedIn	125738	diedOnDate	168037
subclassOf	211979	bornOnDate	350613
givenNameOf	464816	familyNameOf	466969
inLanguage	2389627	isCalled	2984362
type	3057225	means	4014819

Yago: A Large Ontology from Wikipedia and WordNet

Fabian M. Suchanek, Gjergji Kasneci, Gerhard Weikum

MPI-I-2007-5-003 December 2007

Table 6: Simple queries on YAGO

Query	Result
Who was Einstein's doctoral advisor?	\$x=Alfred Kleiner
Einstein HASDOCTORALADVISOR \$x	
Who is named after a place in Africa?	\$who=Gabriel Sudan and 22 more
\$place locatedIn Africa	
\$name means \$place	
\$name familynameof \$who	



Location of Italy (dark green)  
 - in Europe (green & dark grey)  
 - in the European Union (green) — [Legend]

**Capital**  
 (and largest city) Rome  
 41°54' N 12°29' E

**Official language(s)** Italian<sup>[1]</sup>

**Demonym** Italian people

**Government** Unitary parliamentary constitutional republic

**President** Giorgio Napolitano  
**Prime Minister** Mario Monti

**Legislature** Parliament

**Upper house** Senate of the Republic  
**Lower house** Chamber of Deputies

**Formation**  
 - Unification: 17 March 1861  
 - Republic: 2 June 1946

**Area**  
 - Total: 301,338 km<sup>2</sup> (71st)  
 116,348 sq mi  
 2.4

**Water (%)**  
 - Total: 2.4

**Population**  
 - 2011 estimate: 60,813,326<sup>[2]</sup> (23rd)  
 - 2011 (preliminary results) census: 59,570,581<sup>[3]</sup>  
 - Density: 201.8/km<sup>2</sup> (51st)  
 522.7/sq mi

**GDP (PPP)**  
 - 2012 estimate: \$1,834 trillion<sup>[4]</sup> (10th)  
 - Total: \$30,118<sup>[5]</sup> (30th)  
 - Per capita: \$32,522<sup>[4]</sup> (24th)

**GDP (nominal)**  
 - 2012 estimate: \$1,560 trillion<sup>[4]</sup> (8th)  
 - Total: \$32,522<sup>[4]</sup> (24th)  
 - Per capita: \$32,522<sup>[4]</sup> (24th)

**Gini (2006)** 32<sup>[6]</sup>

**HDI (2011)** 0.874<sup>[7]</sup> (very high) (24th)

**Currency** Euro (€) (11st)

**Time zone**  
 - Summer (DST): CET (UTC+1)  
 - Winter (DST): CEST (UTC+2)

**Drives on the right**

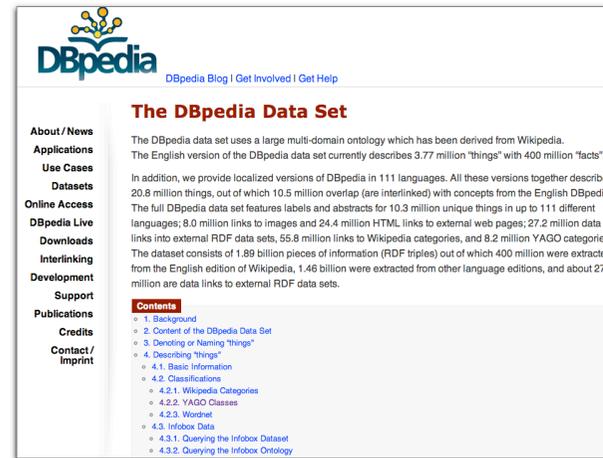
**ISO 3166 code** IT

**Internet TLD** .it

**Calling code** 39<sup>[4]</sup>

# DBpedia

http://dbpedia.aksw.org/dbpedia\_demo/dbpedia/tutorials/ranked\_keyword\_search/demo.php



**The DBpedia Data Set**

The DBpedia data set uses a large multi-domain ontology which has been derived from Wikipedia. The English version of the DBpedia data set currently describes 3.77 million "things" with 400 million "facts".

In addition, we provide localized versions of DBpedia in 111 languages. All these versions together describe 20.8 million things, out of which 10.5 million overlap (are interlinked) with concepts from the English DBpedia. The full DBpedia data set features labels and abstracts for 10.3 million unique things in up to 111 different languages; 8.0 million links to images and 24.4 million HTML links to external web pages; 27.2 million data links into external RDF data sets, 55.8 million links to Wikipedia categories, and 8.2 million YAGO categories. The dataset consists of 1.89 billion pieces of information (RDF triples) out of which 400 million were extracted from the English edition of Wikipedia, 1.46 billion were extracted from other language editions, and about 27 million are data links to external RDF data sets.

**Support**

**Publications**

**Credits**

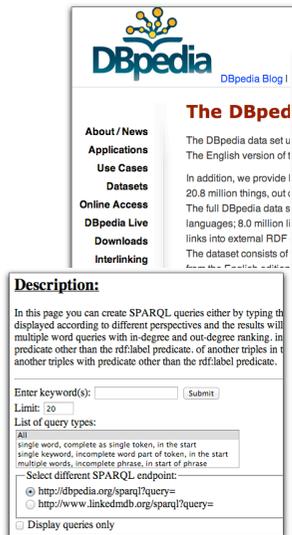
**Contact / Imprint**

**Contents**

- 1. Background
- 2. Content of the DBpedia Data Set
- 3. Denoting or Naming "Things"
- 4. Describing "Things"
  - 4.1. Basic Information
  - 4.2. Classifications
  - 4.2.1. Wikipedia Categories
  - 4.2.2. YAGO Classes
  - 4.2.3. Wordnet
  - 4.3. Infobox Data
  - 4.3.1. Querying the Infobox Dataset
  - 4.3.2. Querying the Infobox Ontology

# DBpedia

http://dbpedia.aksw.org/dbpedia\_demo/dbpedia/tutorials/ranked\_keyword\_search/demo.php



**About / News Applications**

**The DBpedia data set**  
 The English version of the DBpedia data set currently describes 3.77 million "things" with 400 million "facts".

**Use Cases**  
 In addition, we provide localized versions of DBpedia in 111 languages. All these versions together describe 20.8 million things, out of which 10.5 million overlap (are interlinked) with concepts from the English DBpedia.

**Datasets**  
 The full DBpedia data set features labels and abstracts for 10.3 million unique things in up to 111 different languages; 8.0 million links to images and 24.4 million HTML links to external web pages; 27.2 million data links into external RDF data sets, 55.8 million links to Wikipedia categories, and 8.2 million YAGO categories.

**Online Access**  
 The dataset consists of 1.89 billion pieces of information (RDF triples) out of which 400 million were extracted from the English edition of Wikipedia, 1.46 billion were extracted from other language editions, and about 27 million are data links to external RDF data sets.

**DBpedia Live**  
 Links into external RDF data sets, 55.8 million links to Wikipedia categories, and 8.2 million YAGO categories.

**Downloads**  
 The dataset consists of 1.89 billion pieces of information (RDF triples) out of which 400 million were extracted from the English edition of Wikipedia, 1.46 billion were extracted from other language editions, and about 27 million are data links to external RDF data sets.

**Interlinking**  
 The dataset consists of 1.89 billion pieces of information (RDF triples) out of which 400 million were extracted from the English edition of Wikipedia, 1.46 billion were extracted from other language editions, and about 27 million are data links to external RDF data sets.

**Description:**  
 In this page you can create SPARQL queries either by typing the displayed according to different perspectives and the results will multiple word queries with in-degree and out-degree ranking, in predicate other than the rdf:label predicate, or other triples in another triples with predicate other than the rdf:label predicate.

Enter keyword(s):  Submit

Limit: [20]

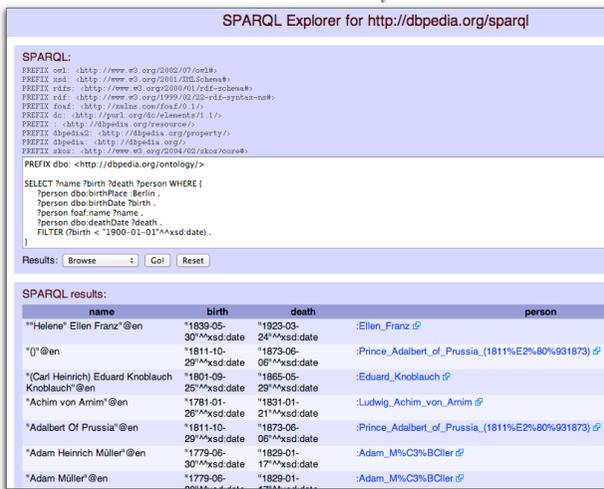
List of query types:

- Single word, complete as single token, in the start
- Single keyword, incomplete word part of token, in the start
- Multiple words, incomplete phrase, in start of phrase

Select different SPARQL endpoint:

- http://dbpedia.org/sparql?query=
- http://www.linkedmdb.org/sparql?query=

Display queries only



SPARQL Explorer for http://dbpedia.org/sparql

**SPARQL:**

```

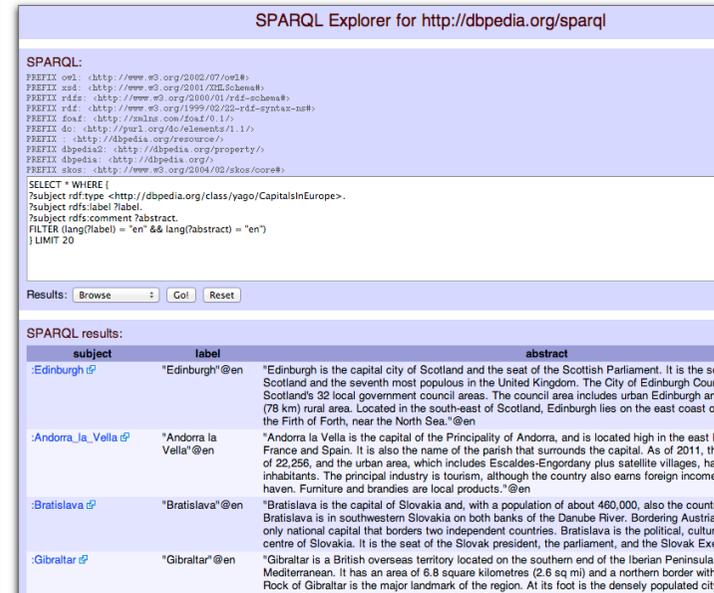
PREFIX owl: <http://www.w3.org/2002/07/owl#>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX rdf:type <http://www.w3.org/2000/01/rdf-schema#>
PREFIX rdfs: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX dbpedia: <http://www.w3.org/2003/06/daos/dbpedia/>
PREFIX dbpedia2: <http://dbpedia.org/property/>
PREFIX dbpedia3: <http://dbpedia.org/>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
PREFIX dbo: <http://dbpedia.org/ontology/>

SELECT ?name ?birth ?death ?person WHERE {
  ?person dbo:birthPlace .Berlin.
  ?person dbo:birthDate ?birth .
  ?person foaf:name ?name .
  ?person dbo:deathDate ?death .
  FILTER (?birth < "1900-01-01"^^xsd:date) .
}
    
```

Results: Browse : Go! Reset

**SPARQL results:**

name	birth	death	person
"Helen" Elien Franz @en	*1830-05-30^^xsd:date	*1823-03-24^^xsd:date	:Elien_Franz @en
"Q" @en	*1811-10-29^^xsd:date	*1873-06-06^^xsd:date	:Prince_Adalbert_of_Prussia_(1811%E2%80%931873) @en
"(Carl Heinrich) Eduard Knoblauch Knoblauch" @en	*1801-09-25^^xsd:date	*1865-06-29^^xsd:date	:Eduard_Knoblauch @en
"Achim von Arnim" @en	*1781-01-26^^xsd:date	*1831-01-21^^xsd:date	:Ludwig_Achim_von_Arnim @en
"Adalbert Of Prussia" @en	*1811-10-29^^xsd:date	*1873-06-06^^xsd:date	:Prince_Adalbert_of_Prussia_(1811%E2%80%931873) @en
"Adam Heinrich Müller" @en	*1779-06-30^^xsd:date	*1829-01-17^^xsd:date	:Adam_M%C3%B4Cler @en
"Adam Müller" @en	*1779-06-30^^xsd:date	*1829-01-17^^xsd:date	:Adam_M%C3%B4Cler @en



SPARQL Explorer for http://dbpedia.org/sparql

**SPARQL:**

```

PREFIX owl: <http://www.w3.org/2002/07/owl#>
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX rdf:type <http://www.w3.org/2000/01/rdf-schema#>
PREFIX rdfs: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
PREFIX dbpedia: <http://www.w3.org/2003/06/daos/dbpedia/>
PREFIX dbpedia2: <http://dbpedia.org/property/>
PREFIX dbpedia3: <http://dbpedia.org/>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>

SELECT * WHERE {
  ?subject rdf:type <http://dbpedia.org/class/Yago/CapitalInEurope> .
  ?subject rdfs:label ?label .
  ?subject rdfs:comment ?abstract .
  FILTER (lang(?label) = "en" && lang(?abstract) = "en")
  } LIMIT 20
    
```

Results: Browse : Go! Reset

**SPARQL results:**

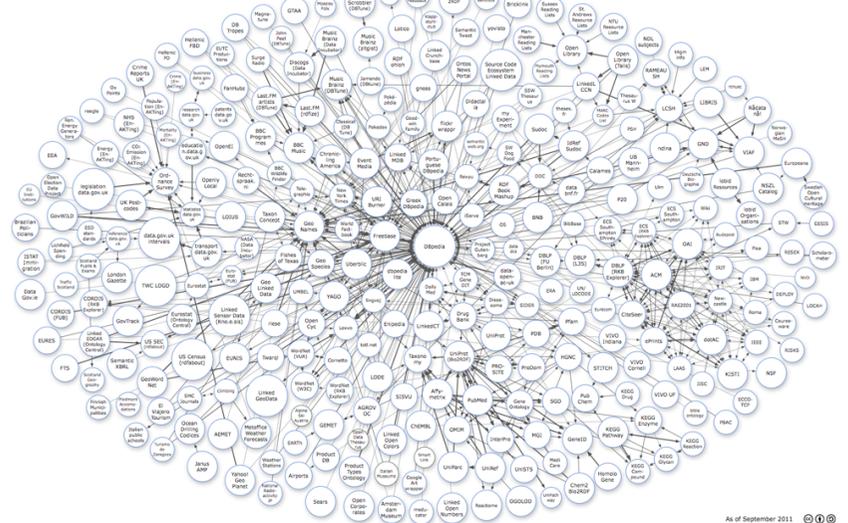
subject	label	abstract
:Edinburgh @en	"Edinburgh" @en	"Edinburgh is the capital city of Scotland and the seat of the Scottish Parliament. It is the second largest city in Scotland and the seventh most populous in the United Kingdom. The City of Edinburgh Council is the local government council area. The council area includes urban Edinburgh and a (78 km) rural area. Located in the south-east of Scotland, Edinburgh lies on the east coast of the Firth of Forth, near the North Sea." @en
:Andorra_la_Vella @en	"Andorra la Vella" @en	"Andorra la Vella is the capital of the Principality of Andorra, and is located high in the east Pyrenees mountains between France and Spain. It is also the name of the parish that surrounds the capital. As of 2011, the city has a population of 22,256, and the urban area, which includes Escaldes-Engordany plus satellite villages, has over 30,000 inhabitants. The principal industry is tourism, although the country also earns foreign income from gambling. Furniture and handicrafts are local products." @en
:Bratislava @en	"Bratislava" @en	"Bratislava is the capital of Slovakia and, with a population of about 460,000, also the country's largest city. It is situated in southwestern Slovakia on both banks of the Danube River. Bordering Austria and Hungary, it is the only national capital that borders two independent countries. Bratislava is the political, cultural, and economic centre of Slovakia. It is the seat of the Slovak president, the parliament, and the Slovak Excutive Council." @en
:Gibraltar @en	"Gibraltar" @en	"Gibraltar is a British overseas territory located on the southern end of the Iberian Peninsula at the southern tip of the Rock of Gibraltar, which has an area of 6.8 square kilometres (2.6 sq mi) and a northern border with Andalusia. The Rock of Gibraltar is the major landmark of the region. At its foot is the densely populated city of Gibraltar, which is home to 30,000 Gibraltarians and other nationalities." @en

SPARQL:		
PREFIX owl: <http://www.w3.org/2002/07/owl#>		
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>		
PREFIX rdfs: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>		
PREFIX foaf: <http://xmlns.com/foaf/0.1/>		
PREFIX owl: <http://www.w3.org/2002/07/owl#>		
PREFIX rdfs: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>		
PREFIX dbpedia: <http://dbpedia.org/resource/>		
PREFIX owl: <http://www.w3.org/2002/07/owl#>		
SELECT * WHERE {		
?subject rdfs:type <http://www.w3.org/1999/02/22-rdf-syntax-ns#label> ?label.		
?subject rdfs:comment ?c.		
FILTER (lang(?label) = "en")		
LIMIT 20		

dbpedia2:partType	"Districts"@en
dbpedia2:highestElevation	514
dbpedia2:codeType	:Vehicle_registration_plates_of_Slovakia
dbpedia2:map1Locator	"Bratislava Region"@en
dbpedia2:map1Size	128
dbpedia2:imageCaption	"Bratislava Montage"@en
dbpedia2:symbolType	"Coat of arms"@en
dbpedia2:utcOffset	"+1"@en
dbpedia2:mapCaption	"Location in Slovakia"@en
dbpedia2:singleLine	"yes"@en
dbpedia2:population	462603
dbpedia2:postalCode	8
dbpedia2:janPrecipitationMm	42
dbpedia2:febPrecipitationMm	37
dbpedia2:junPrecipitationMm	61
dbpedia2:augPrecipitationMm	52
dbpedia2:latNs	"N"@en
dbpedia2:longD	17
dbpedia2:longM	6
dbpedia2:longS	35
dbpedia2:major	:Milan_Ft%C3%A1%C4%8Dnik
dbpedia2:populationDate	2012
dbpedia2:populationDensity	"auto"@en
dbpedia2:establishedType	"First mentioned"@en
dbpedia2:areaUrban	853
dbpedia2:areaMetro	2053
dbpedia2:countryFlag	"true"@en
dbpedia2:elevation	134
dbpedia2:statistics	"[http://www.statistics.sk/mosmis/eng/prvav2.jsp?txtUroven&#61;410190&istMOS/MI/S]"@en
dbpedia2:government	"City council"@en
dbpedia2:river	:Little_Danube
dbpedia2:river	:Danube
dbpedia2:river	:Morava_(river)

<http://linkeddata.org>



Extraction d'informations

## Exemple

Annotation de références bibliographiques avec des automates à états finis

# Approches symboliques avec Unifex

http://igm.univ-mlv.fr/~unifex/

Exemple : Annotation de documents OCRisés (UNESCO)

M. Softi, P. Bellot et al., 2011

2686-26914	Royaume-Uni-United Kingdom
2686	Bhask, Madan; A PRAKASHANAM DIVYANVA INSA. - Henry A. ... London: World, 1952, 192, 316. Dt: <i>Indefatigable</i>
2687	1951. - <i>Chant de l'Inde</i> . - London: Burns & Oates, 1951, 85, 6. Fr: <i>Vivants clairs</i>
2688	Burns & Oates, 1951, 85, 6. Fr: <i>Vivants clairs</i>
2689	Carrière, Hervé. - THE SOCIOLOGY OF RELIGIOUS BELIEFS. - Arthur J. ... London: Duckworth, 1965, 200, 90. Fr: <i>Psychosociologie de l'appartenance religieuse</i>
2690	Comar, Yves. DIALOGUE BETWEEN CHRISTIANITY. - P. ... London: G. Chapman, 1965, 476, 50. Fr: <i>Christianisme et dialogue</i>
2691	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2692	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2693	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2694	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2695	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2696	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2697	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2698	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2699	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2700	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2701	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2702	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2703	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2704	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2705	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2706	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2707	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2708	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2709	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2710	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2711	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2712	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2713	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>
2714	1951. - <i>John Wesley</i> . - London: Chapman, 1951, 22, 21. Fr: <i>John Wesley</i>

# Approches symboliques avec Unifex

http://igm.univ-mlv.fr/~unifex/

Exemple : Annotation de documents OCRisés (UNESCO)

M. Softi, P. Bellot et al., 2011

26881 ECUMENICAL DIALOGUE IN EUROPE. - Fletcher Fleet. - London: Lutterworth, 83, 12/6. Dt: *Dialogue œcuménique, les rencontres œcuméniques des Dombes*  
 26882 Emery, Pierre Y. : THE COMMUNION OF SAINTS. - D. J. & M. Watson. - London: Faith P. XIII, 256. Fr: *L'unité des*



# Etiquetage sorties OCR (2)

```
<!--XML OUTPUT--><OCR Output
<pageinfo><br>159-171 INDEX TRANSLATIONUM</pageinfo><br>
<country>ESPAGNE</country></br>

<italic>Traduccions anunciades dans la « Bibliografia española e hispano-americana » du juillet, août et septembre 1935.
Translations listed in the « Bibliografia española e hispano-americana » for the months of July, August and September 1935.</italic>
<category>Philosophie, Religion,Philosophy, Religion,</category></br>

<br>159. AGUSTYŃ (San), Kempis Agustinián.Máximas de..... sobre la Vida cristiana (Recopiladasdes de Obras por el P. Antonio Tonna-Barthet). -- P.Francisco Mier. -- Barcelona, Ed. Litúrgica Española,1935, 8x, 635.

<br>160. BAETEMAN, J., Breves Respuestas a lasObjeciones contra la Religión. -- Juan Nada. -- Barcelona, Ed. Poliglota, 1935, 8o, 48, Ptas. 0.80.
<italic>Francis, Français, French,</italic>

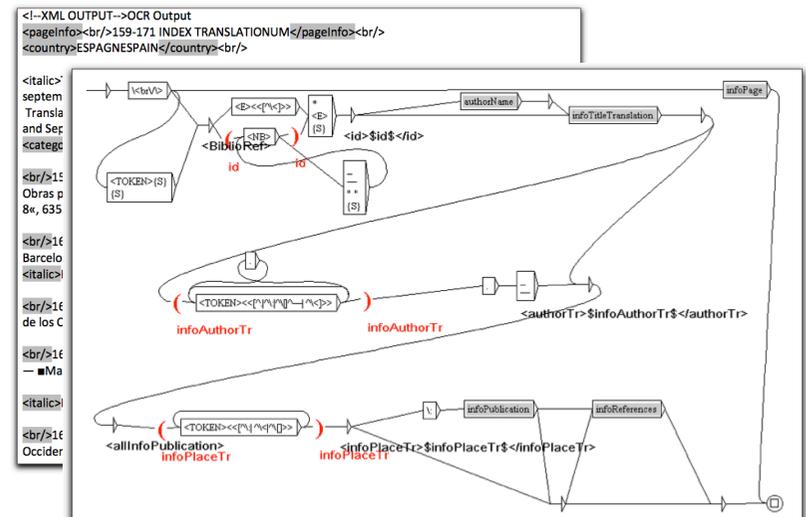
<br>161. BERDIAEFF, Nicolás, El Cristianismo y la Lucha de Clases. Dignidad del Cristianismo e indignidad de los Cristianos. -- María Cardona. --Madrid, Espasa Cal pe, 1935, 8°, 208, Ptas. 5.

<br>162. DELMAS, F. A., y BOLL, M., La Persona+shy;llidad humana. Su Análisis. -- J. Albiñana Mompó. -- Madrid, Aguilar, 1935, 8o, 270, Ptas. 6 (Biblioteca deIdeas y Estudios contemporáneos).

<italic>Francés, Français, French [La Personnalité humaine :Son Analyse ],</italic>

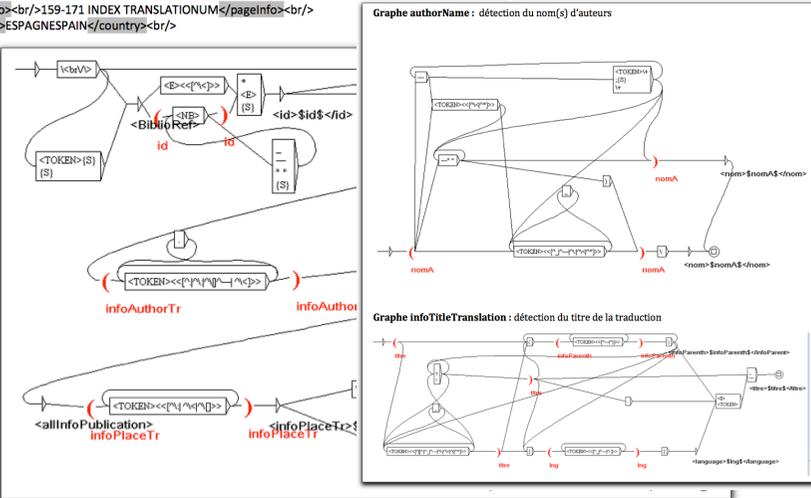
<br>163. DESCARTES, Reglas para la Direccióndel Espíritu. -- Manuel Mindán. -- Madrid, Revistade Occidente, 1935, 8°, 179, Ptas. 6 (Textos filosóficos).
```

# Etiquetage sorties OCR (2)



## Etiquetage sorties OCR (2)

<!--XML OUTPUT-->OCR Output  
 <pageInfo><br/>159-171 INDEX TRANSLATIONUM</pageInfo><br/>  
 <country>ESPAGNE</country></pageInfo>



**Graphe authorName : détection du nom(s) d'auteurs**

**Graphe infoTitle : détection du titre de la traduction**

61

## Etiquetage sorties OCR (3)

```
<!--XML OUTPUT-->
{S}OCR Output
<pageInfo>421_464 Allemagne-Germany</pageInfo><br/>
{S}421 — <italic>N.T. : </italic>
<BiblioRef>
  <id></id>
  <titre>Die Bibel.{S} Die Heilige Schrift d.{S} Neuen Bundes.</titre>
  <allInfoPublication>
    <infoPlaceTr>Freiburg i.{S} Br.{S}</infoPlaceTr>
    <infoPublisher>Herders</infoPublisher>
    <publicationInfo>1968 (4.{S} Aufl.) vni, 276, 48,2.90.</publicationInfo>
  </allInfoPublication>
  <references>Gr</references>
</BiblioRef>
<BiblioRef>
  <id>422</id>
  <nom>— </nom>
  <language>Niederdt</language>
  <titre>Dat Nie Testament in unse Moderspraak.</titre>
  <authorTr>Johannes Jessen</authorTr>
  <allInfoPublication>
    <infoPlaceTr>Göttingen</infoPlaceTr>
    <infoPublisher>Vandenhoeck</infoPublisher>
    <publicationInfo>+ Ruprecht,1968 (6.{S} Aufl.) 494, 10.80.</publicationInfo>
  </allInfoPublication>
  <references>Gr</references>
```

## Exemple

Annotation de références bibliographiques par apprentissage automatique

## Approches probabilistes pour l'annotation

- **Modèles de Markov Cachés (Hidden Markov Models - HMM)**
  - modèles génératifs (apprentissage des paramètres qui maximisent les probabilités «observations + étiquettes»)
  - hypothèses d'indépendance (les observations sont indépendantes)
  - tient compte des états futurs (algo. Forward-Backward) (même si l'état suivant dépend seulement de l'état actuel)
- **Modèles de Markov à Maximum d'Entropie (MaxEnt Markov Models - MEMM)**
  - modèles discriminants (apprentissage des paramètres qui maximisent les probabilités des étiquettes sachant les observations)
  - ne pose pas l'hypothèse d'indépendance des observations (mots)
  - ne tient pas compte des états «futurs»
- **Champs Conditionnels Aléatoires (Conditional Random Fields - CRF)**
  - modèles discriminants tenant compte de *features* très diverses
  - ne pose pas l'hypothèse d'indépendance
  - tient compte des états «futurs»

# Conditional Random Fields (CRF)

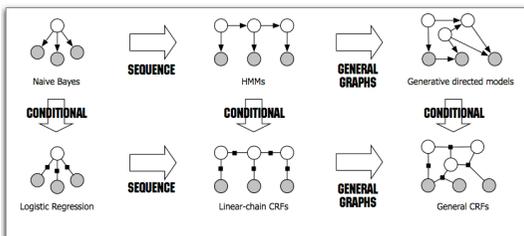


Figure 1.2 Diagram of the relationship between naive Bayes, logistic regression, HMMs, linear-chain CRFs, generative models, and general CRFs.

**1 An Introduction to Conditional Random Fields for Relational Learning**

Charles Sutton  
Department of Computer Science  
University of Massachusetts, USA  
csutton@cs.umass.edu  
<http://www.cs.umass.edu/~csutton>

Andrew McCallum  
Department of Computer Science  
University of Massachusetts, USA  
mccallum@cs.umass.edu  
<http://www.cs.umass.edu/~mccallum>

<http://people.cs.umass.edu/~mccallum/papers/crf-tutorial.pdf>

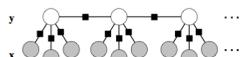


Figure 1.3 Graphical model of an HMM-like linear-chain CRF.

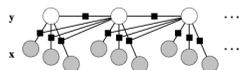


Figure 1.4 Graphical model of a linear-chain CRF in which the transition score depends on the current observation.

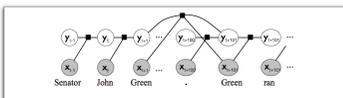


Figure 1.5 Graphical representation of a skip-chain CRF. Identical words are connected because they are likely to have the same label.

# Conditional Random Fields (CRF)

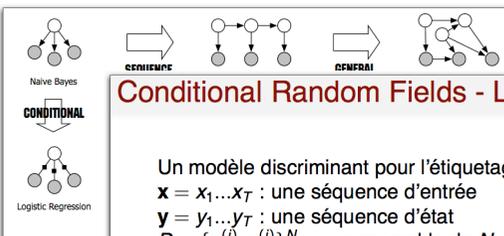


Figure 1.2 Diagram of the relationship between naive Bayes, logistic regression, HMMs, linear-chain CRFs, generative models, and general CRFs.

## Conditional Random Fields - Lafferty et al. (2001)

Un modèle discriminant pour l'étiquetage de données séquentielles  
 $\mathbf{x} = x_1 \dots x_T$  : une séquence d'entrée  
 $\mathbf{y} = y_1 \dots y_T$  : une séquence d'état  
 $D = \{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}_{i=1}^N$  : un ensemble de  $N$  exemples

La probabilité conditionnelle d'un CRF

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp\left\{\sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t)\right\}, \quad (1)$$

$\theta = \{\theta_k\} \in R^K$  : un vecteur de paramètre  
 $\{f_k(y_t, y_{t-1}, \mathbf{x}_t)\}_{k=1}^K$  : un ensemble des fonctions de caractéristiques  
 $Z(\mathbf{x})$  : une fonction de normalisation

Figure 1.3 Graphical model of a linear-chain CRF in which the transition score depends on the current observation.

Figure 1.4 Graphical model of a linear-chain CRF in which the transition score depends on the current observation.

Figure 1.5 Graphical representation of a skip-chain CRF. Identical words are connected because they are likely to have the same label.

# BILBO - Analyse automatique de documents

<http://bilbo.hypotheses.org/>



## BILBO - Annotation automatique des références bibliographiques



### Exemples

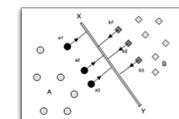
- Non bibl note**  
 26. La nature euro-centrée du projet était encore plus apparente dans la version originale du texte, qui, comme nous l'avons déjà mentionné, était appelé " Europe élargie ".
- Multi bibl note**  
 27. " Une Europe sûre dans un monde meilleur. Stratégie européenne de sécurité ", Bruxelles, 12 décembre 2003. Pour un commentaire critique, voir Toje A., " The 2003 European security strategy: A critical appraisal ", *European Foreign Affairs Review*, vol. 10, n°1, 2005, pp. 117-134.
- Part. info.**  
 31. Voir par exemple la communication de la Commission relative au " Renforcement de la politique européenne de voisinage ", COM (2006) 726 final.

# BILBO - Données d'apprentissage

- Token, features, tags -> Données d'apprentissage
- Environnement MALLET <http://mallet.cs.umass.edu/>
- Combinaison de SVM (segmentation du texte et classification) et de CRF (annotation)

```

<?xml version="1.0" encoding="UTF-8" ?>
<document>
  <title>Medical essays and observations</title>
  <author>
    <name>Galen</name>
  </author>
  <publisher>Society in Edinburgh</publisher>
  <date>1753</date>
  <link target="http://www.biodidacte.com" />
  <text>
    Research in karst hydrogeology and geomorphology, which was carried out in the 1970s, 1980s and 1990s
  </text>
  </document>
  
```



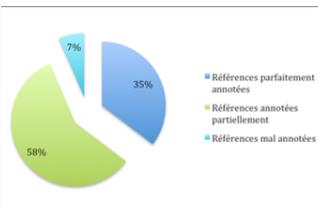
# BILBO - Evaluation détaillée

**Table 2: Evaluation on UMich Library level 1 data including all files (upper table) and excluding 'ark' and 'mbot' (lower table).**

Tool	Precision			Recall			F-measure		
	total	author	title	total	author	title	total	author	title
BILBO	<b>0.68</b>	<b>0.86</b>	<b>0.84</b>	<b>0.69</b>	0.86	0.73	<b>0.68</b>	<b>0.86</b>	<b>0.78</b>
ParsCit	0.57	0.79	0.73	0.49	0.63	0.52	0.53	0.70	0.63
Biblio	0.57	0.72	0.62	0.55	0.57	<b>0.83</b>	0.56	0.64	0.71
Freecite	0.65	0.76	0.76	0.62	<b>0.94</b>	0.68	0.63	0.83	0.72
Grobid	-	-	-	-	-	-	-	-	-

Tool	Precision			Recall			F-measure		
	total	author	title	total	author	title	total	author	title
BILBO	0.63	<b>0.79</b>	<b>0.85</b>	<b>0.64</b>	0.86	0.70	<b>0.64</b>	<b>0.82</b>	<b>0.77</b>
ParsCit	0.51	0.69	0.71	0.41	0.53	0.51	0.45	0.60	0.59
Biblio	0.57	0.72	0.60	0.59	0.63	<b>0.88</b>	0.58	0.67	0.71
Freecite	0.59	0.71	0.75	0.57	<b>0.89</b>	0.68	0.58	0.78	0.71
Grobid	<b>0.64</b>	0.67	0.79	0.62	0.80	0.63	0.63	0.73	0.70



Y.M. Kim, P. Bellot et al. (ACM DocEng 2012, LREC 2012)

# Systèmes de questions-réponses

# Les moteurs de questions-réponses

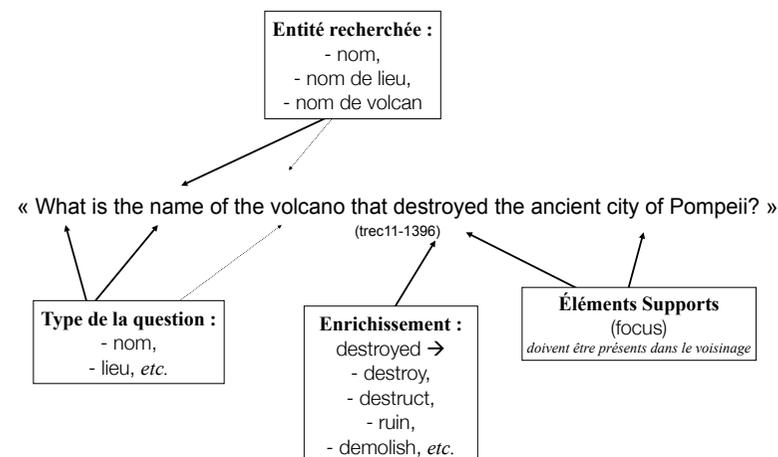
- Questions factuelles (*quelle est la capitale ....*)
- Questions booléennes (*est-ce que .... ?*)
- Questions "liste" (*quelles sont .... ?*)
- Les questions complexes : "comment ?", "pourquoi ?" etc.

Référence	Libellé de la question
GF222	Qui est le maire de Baska ?
GF30	Combien de personnes souffrent d'acné en Suisse ?
GF266	Quelle est la monnaie nationale en Hongrie ?
GF219	A combien de kilomètres de Paris se trouve la gare de Tours ?
GF178	Comment s'appelle le président Tchadéane ?
GF6	Quel âge a l'abbé Pierre ?
GF232	Combien y a t'il de chômeurs en Europe ?
GF245	Qui est la frère de la princesse Leia ?
GF17	Combien y a t'il d'habitants en Lettonie ?
GF29	Quel grade occupe Juan Carlos Ralon dans la marine ?
GF273	Quand est mort Kurt Cobain ?
GF105	Quel est le nom du roi du Maroc ?
GF298	Quelle est la capitale de Terre Neuve ?
GF147	En quelle année Hitler est arrivé au pouvoir ?
GF176	Qui est le président d'Aérospatiale ?
GF99	Combien de personnes sont mortes dans des accidents de la route en 1997 ?
GF132	Où se situe San Cristobal de Las Casas ?
GF206	Quand a été votée la loi Esin ?
GF84	En combien de langues a été publié le Petit Prince ?
GF78	Quel journal public-chaque année le top 50 des personnalités ?

Questions évaluation EQUER 2005 (Grau et al.)

# Analyses complexes

- **Objectif** : extraire de la question le + d'informations → localiser la réponse



**Population Of France**  
The population of France, as estimated in late 2006, is **60,742,000**.  
The Capital City of France is **Paris**.

**A: WikiAnswers**  
**Answer**  
65 million people. (this was in January)

**France This Way**  
The population of France at July 2006 is estimated at 65.8 million people, of a total world population of 6.1 billion people. This makes it the 23rd most populated country in the world, 18.4% of the population of France are 65 years old or more.

**Population**  
(January 1, 2008 estimate)  
Total<sup>[1]</sup> 64,473,140<sup>[5]</sup> (20th)  
- Metropolitan France 61,875,822<sup>[4]</sup> (20th)  
- Density<sup>[9]</sup> 114/km² (89th)  
295/sq mi

## Classification des questions

- Objectifs :
  - déterminer le type de la question (factuelle, procédurale, définitoire...)
  - déterminer le type de réponse attendue (nom propre, nom de lieu, quantité numérique...)
- Méthodes :
  - à base de règles et de patrons syntaxiques
  - à base d'apprentissage automatique (approche bayésienne, SVM...)
  - méthodes mixtes
  - exploitation des POS, de bigrammes, des lemmes, de ressources...
- Performances :
  - env. 90% F-Mesure sur quelques dizaines de classes

## Recherche documentaire dans QR

- Pour la recherche de documents : peu importe le modèle (booléen, vectoriel, probabiliste...)
- Pour la recherche de phrases : c'est différent ! [Radev, 2002]
- Combinaison de modèles de langage meilleure que Okapi (ordonnancement des phrases trouvées par Google, Alltheweb, Northernlight)
- Pour 132 questions (sur 200), la réponse est dans les 20 premières phrases

♦ TREC 2000  
542 questions ont une réponse juste dans les 50 premiers documents (576 dans les 200 premiers)

♦ TREC 2001 (500 questions) : redondance  
- 37 questions ont leur réponse dans au moins 25 documents de la collection complète  
- 138 questions ont leur réponse dans au moins 10 documents de la collection complète

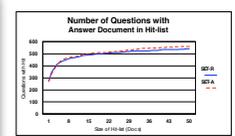


Figure 1. Comparison of number of questions with a "correct" document in the hit list against hit-list depth, for our search engine (SET-R) and AT&T's (SET-A).

Prager et al.,  
Guru / SMART

**En moyenne : document contenant la bonne réponse en position 2 (MRR = 0,49)**

## Robustesse et suggestions de réécritures

- Evaluation sur des questions "réelles" du système SQUALIA (Gillard et al.)
- 20 questions X 19 utilisateurs  
(9 Français langue maternelle, 6 étrangers vivant en France, 2 dyslexiques)

Comment s'appelle le président Tchèque ?  
Quel est le nom du Président tchèque ?  
Quel est le nom du président de la tchetchenie ?  
Comment s'appelle le président de la Tchetchénie ?  
Quel est le nom du président de la Tchetchénie ?  
Comment s'appelle le président tchèque ?  
quel est le nom du président de la tchetchénie  
Comment s'appelle le président de le Tchetchenie ?  
qui est le président de la tchetchenie  
quel est le nom du président  
Qui est le Président de la Tchetchénie ?  
Quelle est le nom du président de la Cherchenie ?  
Quelle est le nom du président Tchèque ?  
Qui est le président de la tchetchénie ?  
qui est le président de la tchetchénie  
Qui est le président de la Tchetchenia ?  
Quel est le nom du président de la Tchetchenie ?

SITBON, BELLOT, BLACHE, (LREC, 2008)

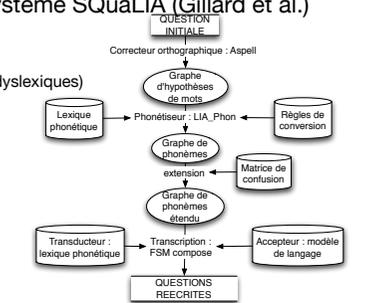


FIG. 410: La procédure de correction orthographique conduisant à la génération d'hypothèses de questions utilise un correcteur orthographique standard, un phonétiseur (LIA\_Phon) et un système de transcription automatique (figure reprise de Sitbon et al. (2008)).

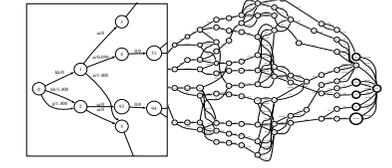


FIG. 412: Exemple d'hypothèses phonétiques (phonème/épave, nombre de syllabes) produites par LIA\_Phon à partir des hypothèses graphématiques du correcteur Aspell et de la question : Qui est le Président de la Tchetchenie ? — figure reprise de Sitbon et al. (2008).

## Quelques thèses récentes...

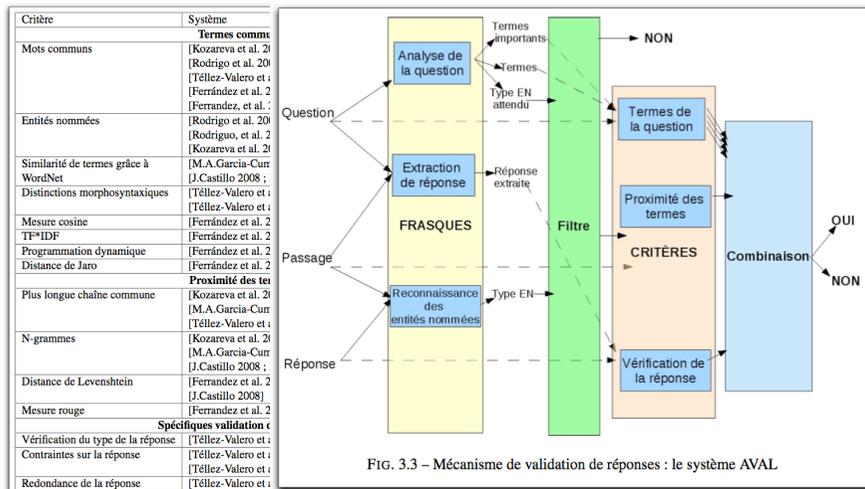
- « Réordonnement d'hypothèses dans un système de questions-réponses », Guillaume BERNARD, LIMSI Paris-Sud, 2011
- « Extraction d'information générique à partir de textes fondée sur une analyse linguistique profonde », Ludovic JEAN-LOUIS, CEA LVIC, Paris-Sud, 2011
- « Validation de réponses dans un système de questions réponses », Arnaud GRAPPY, LIMSI Paris-Sud, 2011
- « Désignations nominales des événements – Etude et extraction automatique dans les textes », Béatrice ARNULPHY, LIMSI Paris-Sud, 2012
- « Extraction automatique de segments textuels, détection de rôles, de sujets et de polarités », Rémi LAVALLEY, LIA Avignon, 2012
- ...

## Validation de réponses (A. Grappy, 2011)

Critère	Système
<b>Termes communs</b>	
Mots communs	[Kozareva et al. 2006 ; Herrera et al. 2006] [Rodrigo et al. 2006 ; M.A.Garcia-Cumbreas et al. 2007] [Téllez-Valero et al. 2007 ; Téllez-Valero et al. 2008] [Ferrández et al. 2008 ; J.Castillo 2008] [Ferrandez, et al. 2007]
Entités nommées	[Rodrigo et al. 2006 ; Herrera et al. 2006] [Rodrigo, et al. 2007 ; Ferrández et al. 2008] [Kozareva et al. 2006 ; Téllez-Valero et al. 2008]
Similarité de termes grâce à WordNet	[M.A.Garcia-Cumbreas et al. 2007] [J.Castillo 2008 ; Ferrández et al. 2008]
Distinctions morphosyntaxiques	[Téllez-Valero et al. 2007 ; Ferrández et al. 2008] [Téllez-Valero et al. 2008]
Mesure cosin	[Ferrández et al. 2008 ; J.Castillo 2008]
T <sup>F</sup> *IDF	[Ferrández et al. 2008 ; J.Castillo 2008]
Programmation dynamique	[Ferrández et al. 2008]
Distance de Jaro	[Ferrández et al. 2008]
<b>Proximité des termes</b>	
Plus longue chaîne commune	[Kozareva et al. 2006 ; Bosma & Callison-Bursh 2006] [M.A.Garcia-Cumbreas et al. 2007] [Téllez-Valero et al. 2007 ; Ferrandez et al. 2007]
N-grammes	[Kozareva et al. 2006 ; Herrera et al. 2006] [M.A.Garcia-Cumbreas et al. 2007 ; Ferrandez et al. 2007] [J.Castillo 2008 ; Rodrigo et al. 2006]
Distance de Levenshtein	[Ferrandez et al. 2007 ; Ferrández et al. 2008] [J.Castillo 2008]
Mesure rouge	[Ferrandez et al. 2007]
<b>Spécifiques validation de réponses</b>	
Vérification du type de la réponse	[Téllez-Valero et al. 2007 ; Téllez-Valero et al. 2008]
Contraintes sur la réponse	[Téllez-Valero et al. 2007 ; Ferrández et al. 2008]
Redondance de la réponse	[Téllez-Valero et al. 2008]

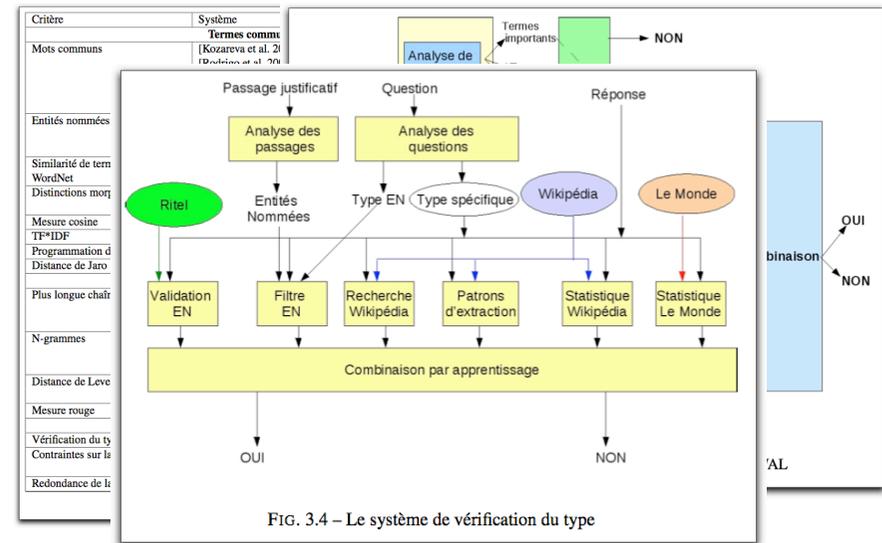
TAB. 1.1 – Critères de validation de réponses

## Validation de réponses (A. Grappy, 2011)



TAB. 1.1 – Critères de validation de réponses

## Validation de réponses (A. Grappy, 2011)



## Liens entre entités - (L. Jean-Louis, 2011)

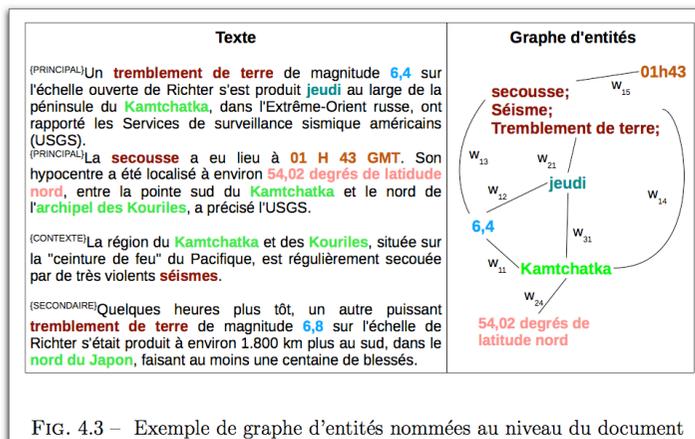


FIG. 4.3 – Exemple de graphe d'entités nommées au niveau du document

## Liens entre entités - (L. Jean-Louis, 2011)

Description des <i>features</i>	FEAT-BASE	FEAT-LEX	FEAT-NOLEX
Types des entités E1 et E2	✓	✓	✓
POS des entités E1 et E2	✓	✓	✓
Mots des entités E1 et E2	✓	✓	✓
Bigrammes de mots de E1 et E2	✓	✓	✓
Mots entre E1 et E2	✓	✓	✓
Bigrammes de mots entre E1 et E2	✓	✓	✓
POS des mots entre E1 et E2	✓	✓	✓
Nb mots entre E1 et E2	✓	✓	✓
Bigrammes de POS de mots entre E1 et E2	✓	✓	✓
Nb relations syntaxiques entre E1 et E2	✓	✓	✓
Chemin syntaxique entre E1 et E2	✓	✓	✓
Position relative et POS <sup>1</sup>	✓	✓	✓
Nb d'entités entre E1 et E2	✓	✓	✓
Nb mentions d'événements entre E1 et E2	✓	✓	✓
POS des deux mots avant/après E1	✓	✓	✓
POS des deux mots avant/après E2	✓	✓	✓

FIG. 4.3 – Ex

TAB. 1 – *Features* pour la classification de relations

## Liens entre entités - (L. Jean-Louis, 2011)

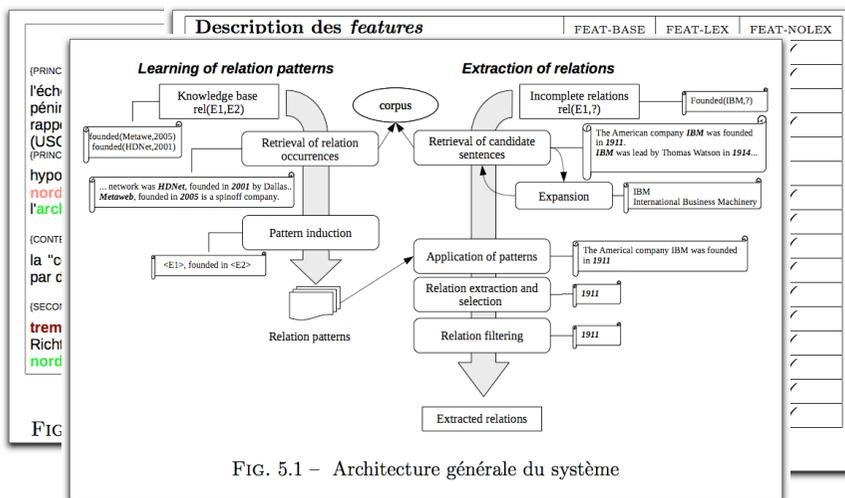


FIG. 5.1 – Architecture générale du système

## Plateformes logicielles

Infrastructure for Human  
Language Technology

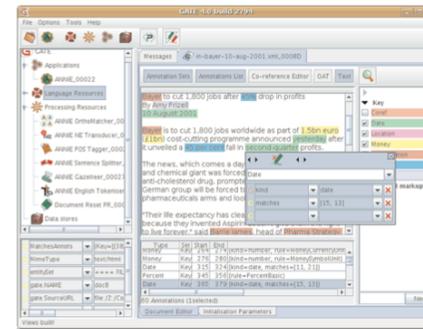
**GATE**

<http://gate.ac.uk>



**FREE**  
**Open source**, licensed under LGPL allowing unrestricted commercial use, hosted on SourceForge.  
**100% JAVA**  
 Runs on **any platform** supporting Java 5 or later. Developed and tested daily on Linux, Windows, and Mac OS X.  
**MATURE AND ACTIVELY SUPPORTED**  
 In development **since 1996**; now at version 5.0; around 20 active developers.  
**COMPREHENSIVE**  
 Support for manual annotation, performance evaluation, information extraction, [semi-]automatic semantic annotation, and many other tasks.  
 Over **50 plugins** included with the standard distribution, containing over 70 resource types. Many others available from independent sources.

**STANDARDS-BASED**  
 Reference implementation in ISO TC37/SC4 LIRICs project; supports XCES, ACE, TREC etc. formats; founder member of OASIS/UIMA committee.  
**EFFICIENT**  
 Optimisations included with the latest version provide a 20 to 40% speed and memory usage improvement. Highly efficient finite state text processing engine; many plugins with linear execution time.  
**POPULAR**  
 Assessed as "outstanding" and "internationally leading" by an anonymous EPSRC peer review. Used at thousands of sites: companies, universities and research laboratories, all over the world. ~35,000 downloads/year.  
 Rolling funding for more than 15 staff at the University of Sheffield.



**INTEGRATION**

Leveraging the power of other projects such as:

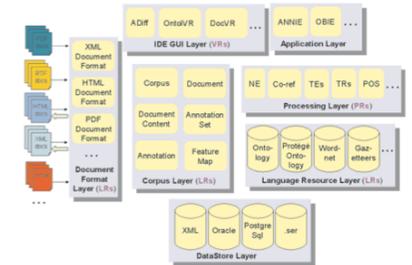
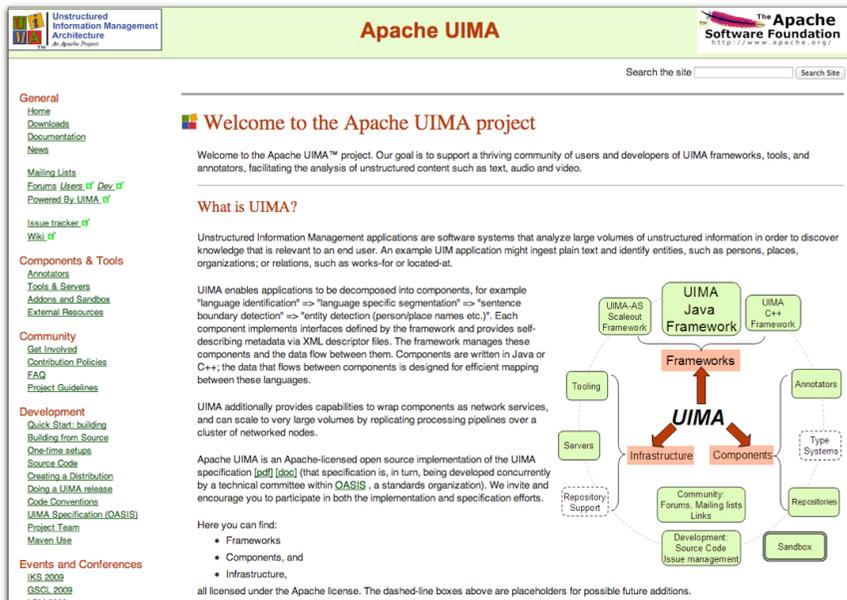
- **Information Retrieval:** Lucene (Nutch, Solr), Google and Yahoo search APIs, MG4;
- **Machine Learning:** Weka, MaxEnt, SVMlight, etc.;
- **Ontology Support:** Sesame and OWLIM;
- **Parsing:** RASP, Minipar, and SUPPLE;
- **Other:** UIMA, Wordnet, Snowball, etc.

**COMMUNITY AND SUPPORT**  
 Friendly and active community of developers and users offers efficient help. Commercial support available from Ontotext and Matrixware.

**DATA MANAGEMENT**  
 Pluggable input filters with out of the box support for XML, HTML, PDF, MS Word, email, plain text, etc.  
 Common in-memory data model built around stand-off annotation, documents and corpora.  
 Persistent storage layer with support for XML, Oracle, PostgreSQL, or Java serialisation. I/O interoperation with many other systems.

**STANDARD ALGORITHMS**  
 Ready made implementations for many typical NLP tasks such as tokenisation, POS tagging, sentence splitting, named entity recognition, co-reference resolution, machine learning, etc.

**USER INTERFACE**  
 Comprehensive tool set for data editing and visualisation, rapid application development, manual annotation, ontology management.


## Conclusions...

- La RI sur des textes (et audio) fait partie du TAL
- Le TAL utilise des approches symboliques / numériques, des méthodes d'apprentissage automatique, des ressources et bases de connaissances
- Tâches complexes à effectuer et à évaluer (y compris pour un humain)
  - au cœur de nombreuses (toutes ?) campagnes d'évaluation en RI
- Traitements souvent lourds : question de la *pertinence* de les utiliser
- Les améliorer avant de les utiliser en RI ? (complexité, robustesse...)
- Apprendre à mieux les intégrer au sein des modèles de RI ?
  - quelle information ? à quel moment ?
  - architecture informatique (intégration)
  - combinaison de scores, confiance...

## Quelques livres...

