

# Information Extraction

Anne-Laure Ligozat

last update: 2019

# NLP & the semantic web

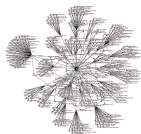
## Textual documents and knowledge bases

paper

tweets

wikipedia

web page



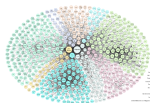
wikidata



wordnet



gene ontology



linking open data

# Textual documents and structured knowledge

interrogation et complétion des données structurées

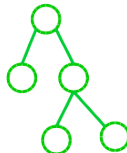
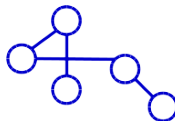
non structuré

A quelques jours d'élections professionnelles qui se dérouleront dans un climat tendu, **SUD** confirme être devenu une cible pour la direction d'**Ad**. Coup sur coup, trois représentants du syndicat ont été visés par des sanctions importantes dans le groupe fondé par **Xavier Niel** (qui détient la marque **Free**). Deux d'entre eux ont fait l'objet d'un entretien préalable au licenciement **vendredi 25 octobre**.

En **ce jeudi après-midi**, malgré une température démente, les cache-nez sont de rigueur aux abords de la **rue Hénard**, dans le **12e arrondissement de Paris**. 2 à 300 policiers « en colère » se massent devant un cordon de sécurité dressé à proximité des locaux de l'**IGPN** (Inspection générale de la police nationale). Ils sont venus soutenir leur collègue **Guillaume Lebeau**, agent de la **BAC** des **Hauts-de-Seine**, auditionné pour s'être répandu à visage découvert dans les médias lors des précédentes manifestations spontanées.

Au cours d'un deuxième débat tendu, l'ex-chef de l'Etat a servi de punching-ball à ses adversaires. Incarnation présidentielle, alliances diplomatiques douteuses, inconstance politique... Les candidats ont décidé d'user de leur droit d'inventaire sur les années **Sarkozy**.

(exemples Mediapart)



structuré

ajout de structure au texte

# Information/knowledge extraction

## Objectives

- targeted understanding of texts
- produce a structured representation of relevant information
  - relational database
  - knowledge base
- reasoning and inference

## Machine readable abstract



Subject	Relation	Object
p53	is_a	protein
Bax	is_a	protein
p53	has_function	apoptosis
Bax	has_function	induction
apoptosis	involved_in	cell_death
Bax	is_in	mitochondrial outer membrane
Bax	is_in	cytoplasm
apoptosis	related_to	caspase activation
...	...	...

5 / 73

# Targeted understanding?

**Lawrence Livermore National Laboratory**

From Wikipedia, the free encyclopedia

The **Lawrence Livermore National Laboratory (LLNL)** in **Livermore, California** is a scientific research laboratory founded by the University of California in 1952. It is funded by the United States Department of Energy (DOE) and managed by Lawrence Livermore National Security, LLC (LLNS), a partnership of the University of California, Bechtel Corporation, Babcock and Wilcox, the LBNL Corporation, and Battelle Memorial Institute. On October 1, 2007 LLNS assumed management of LLNL from the University of California, which had exclusively managed and operated the Laboratory since its inception 55 years before.

**Contents** (hide)

- 1 Background
- 2 Origins
- 3 Weapons projects
- 4 Plutonium research
- 5 National Ignition Facility and proton science
- 6 Global security program
- 7 Other programs
- 8 Key accomplishments
- 9 Unmanned facilities
- 10 Workstations, computers
- 11 Spinoffs
- 12 Directors
- 13 Organization
- 14 Facilities
- 15 References
- 16 External links and sources

**Background**

LLNL is self-described as "a premier research and development institution for science and technology applied to national security."<sup>[1]</sup> Its principal responsibility is ensuring the safety, security and reliability of the nation's nuclear weapons through the application of advanced science, engineering and technology. The Laboratory also applies its special expertise and multidisciplinary capabilities to preventing the proliferation and use of weapons of mass destruction, bolstering homeland security and solving other nationally important problems, including energy and environmental security, basic science and economic competitiveness.

LLNL is home to many unique facilities and a number of the most powerful computer systems in the world, according to the TOP500 list, including Blue Gene/L, the world's fastest computer from 2004 until Los Alamos National Laboratory's Roadrunner supercomputer surpassed it in 2008. The Lab is a leader in technical innovation: since 1978, LLNL has received a total of 118 prestigious R&D 100 Awards, including

"The **Lawrence Livermore National Laboratory (LLNL)** in **Livermore, California** is a **scientific research laboratory** founded by the **University of California** in 1952."



**LLNL** EQ **Lawrence Livermore National Laboratory**  
**LLNL** LOC-IN **California**  
**Livermore** LOC-IN **California**  
**LLNL** IS-A **scientific research laboratory**  
**LLNL** FOUNDED-BY **University of California**  
**LLNL** FOUNDED-IN 1952

# Real life IE applications









# Real life IE applications

Google  🔍

Tous Maps Actualités Images Vidéos Plus ▾ Outils de recherche

---



**Lyon** / Lieux d'intérêt

<b>Vieux Lyon</b> Renaissance, histoire, jardin		<b>Primatiale Saint-Jean de Lyon</b> Cathédrale	
<b>Basilique Notre-Dame de Fourvière</b> Église à plan basilical		<b>Parc de la Tête d'Or</b> Parc	
<b>Fourvière</b> Musée, théâtre, amphithéâtre, cathédrale, église		<b>Musée des beaux-arts de Lyon</b> Musée, musée d'art, art, architecture	

---


**Office du Tourisme de Lyon: Accueil**  
[www.lyon-france.com/](http://www.lyon-france.com/) ▾  
Only Lyon - Tourisme et Congrès ..... Sandrine, conseiller séjour, vous accueille aujourd'hui lundi 7 novembre, au pavillon du tourisme place Bellecour. Lyon ...




**L'Office du Tourisme - Office du Tourisme de Lyon**  
[www.lyon-france.com/L-Office-du-Tourisme](http://www.lyon-france.com/L-Office-du-Tourisme) ▾









# Real life IE applications

Google  

Tous Maps Actualités Images Vidéos Plus ▾ Outils de recherche   

---



**Lyon** / Lieux d'intérêt

<b>Vieux Lyon</b> Renaissance, histoire, jardin		<b>Primatiale Saint-Jean de Lyon</b> Cathédrale	
<b>Basilique Notre-Dame de Fourvière</b> Eglise à plan basilical		<b>Parc de la Tête d'Or</b> Parc	
<b>Fourvière</b> Musée, théâtre, amphithéâtre, cathédrale, église		<b>Musée des beaux-arts de Lyon</b> Musée, musée d'art, art, architecture	


---

**Office du Tourisme de Lyon: Accueil**  
[www.lyon-france.com/](http://www.lyon-france.com/) ▾  
Only Lyon - Tourisme et Congrès ..... Sandrine, conseiller séjour, vous accueille aujourd'hui lundi 7 novembre, au pavillon du tourisme place Bellecour. Lyon ...

**L'Office du Tourisme - Office du Tourisme de Lyon**  
[www.lyon-france.com/L-Office-du-Tourisme](http://www.lyon-france.com/L-Office-du-Tourisme) ▾

# Real life IE applications

Google  


[Tous](#) [Maps](#) [Actualités](#) [Images](#) [Shopping](#) [Plus ▾](#) [Outils de recherche](#)

Environ 10 300 résultats (0,47 secondes)

## Rue John Von Neumann, 91400 Orsay

LIMSI, Adresse

[Commentaires](#)



[Voir les photos](#)

**LIMSI**  
<https://www.limsi.fr/fr/> ▼

Laboratoire de recherche en informatique pluridisciplinaire, le LIMSI rassemble des chercheurs et enseignants-chercheurs relevant des Sciences de l'Ingénieur ...

**Annuaire - Limsi**  
<https://annuaire.limsi.fr/fr/annuaire> ▼

**LIMSI** ★ [Site Web](#)

Centre National de la Recherche Scientifique

**Adresse :** Rue John Von Neumann, 91400 Orsay  
**Téléphone :** 01 69 85 80 80

# Application domains

- open domain
- digital libraries (google scholar, citeseer)
- bioinformatics
- patent analysis...

# Knowledge bases

## RDF knowledge base

### Set of facts

- fact = sujet, predicate, object
  - resources = entities, concrete or abstract
  - properties = relations, such as height for a Person

## Triple example

```
<http://dbpedia.org/resource/J._K._Rowling>  
<http://dbpedia.org/ontology/notableWork>  
<http://dbpedia.org/resource/Harry_Potter>
```



# Some KB examples

## Wikidata

### In a nutshell

- partial transfer from Freebase
  - collaborative content
  - import of other sources such as MusicBrainz
- more than 20 millions items (nov 2016)



Main page  
Community portal  
Project chat  
Create a new item  
Item by title  
Recent changes  
Random item  
Query Service  
Nearby  
Help  
Donate

#### Tools

What links here  
Related changes  
Special pages  
Permanent link  
Page information  
Concept URI

Item

Discussion

Read

View history

Search Wikidata

## London (Q84)

capital of England and the United Kingdom edit

London, UK | London, United Kingdom | London, England

[In more languages](#) [Configure](#)

Language	Label	Description	Also known as
English	London	capital of England and the United Kingdom	London, UK London, United Kingdom London, England
French	Londres	capitale du Royaume-Uni	London
Occitan	Londres	No description defined	
Italian	Londra	capitale dell'Inghilterra e del Regno Unito	

[All entered languages](#)

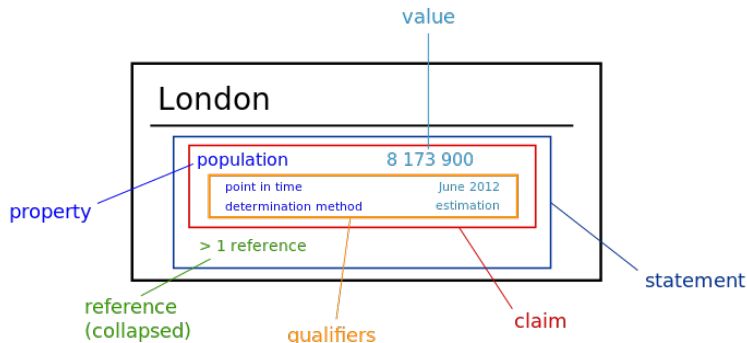
### Statements

# Some KB examples

## Wikidata

### In a nutshell

- partial transfer from Freebase
  - collaborative content
  - import of other sources such as MusicBrainz
- more than 20 millions items (nov 2016)



# Some KB examples

## YAGO

### In a nutshell

- ontology with a very rich typing system
- partly built from WordNet and Wikipedia, in particular categories
- built on an RDFS extension
- everything is an entity:
  - objects: cities, persons, URLs, numbers, words...
  - classes (hierarchy)
  - relations
  - facts = (entity, relation, entity)
    - entity < possible to give the reference
- *n*-ary relations = one main fact + other arguments in relation with this fact



# Course objective

## Information extraction/Knowledge acquisition

- from the NLP view

## Main components

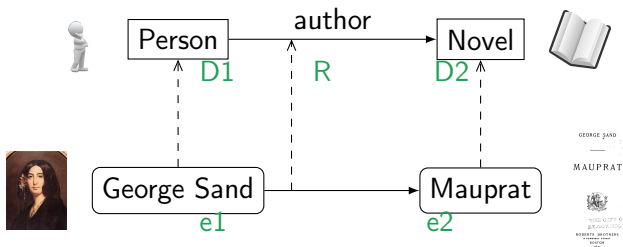
- What is it about?
  - entities: who? what? when? where?
- What is said about it?
  - relations between entities

# A few definitions

relation and types

entities and relation instance

mentions



Mauprat<sub>m1</sub>, que Sand<sub>m2</sub> écrivit<sub>mr</sub> entre 1835 et 1837, est bien un roman capital dans son œuvre.

# Example

## Objective

Find Cecilia Bartoli's first active years in order to store it in a KB (DBPedia relation `activeYearsStartYear`)

En 1985 — elle n'a que 19 ans —, Cecilia Bartoli se fait connaître en France.

# Steps

En 1985<sub>DATE</sub> — elle n'a que 19 ans —, Cecilia Bartoli<sub>PERS</sub> se fait connaître en France.



Cecilia Bartoli



1985

# Steps

En 1985<sub>DATE</sub> — elle n'a que 19 ans —, Cecilia Bartoli<sub>PERS</sub> se fait connaître en France.



Cecilia Bartoli



1985

# Steps

En 1985<sub>DATE</sub> — elle n'a que 19 ans —, Cecilia Bartoli<sub>PERS</sub> se fait connaître en France.



Cecilia Bartoli

se fait connaître



1985

# Steps

En 1985<sub>DATE</sub> — elle n'a que 19 ans —, Cecilia Bartoli<sub>PERS</sub> se fait connaître en France.



## 1 Introduction

## 2 Entities

- Definitions
- NE recognition
- Entity linking

## 3 Relations

- Definitions
- Relation extraction
- Supervised methods
- Semi-supervised methods

## 4 Conclusion



# Named entity

## Definition(s)

- linguistic expression referring to a unique referent in a category in context
- typically: persons, organizations, locations
  - numerical entities often associated: dates, amounts, speed...

## Example of an annotated text

Le 27 avril 2006<sub>DATE</sub> à Washington<sub>LIEU</sub>, George Clooney<sub>PERS</sub> et Barack Obama<sub>PERS</sub> assistent à une conférence de presse sur le Darfour<sub>LIEU</sub>.

# Annotation formats

- parentheses
  - [*ORG* U.N.] official [*PERS* Ekeus] heads for [*LOC* Baghdad] .
- XML
  - `<org>U.N.</org>` official `<personne>Ekeus</personne>` heads for `<lieu>Baghdad</lieu>`.
  - `<enamel type="organisation"> U.N.</enamel>` official `<enamel type="person">Ekeus</enamel>` heads for `<enamel type="organisation">Baghdad</enamel>` . (MUC)
- BIO or variants (ex: BILOU=BIO+Last+Unique)

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

# NE recognition

## Text example

(...) et Obama assistent à (...)

# NE recognition

- recognition
  - identification

## Text example

(...) et **Obama** assistent à (...)

# NE recognition

- recognition
  - identification
  - categorization



## Text example

(...) et **<personne>** Obama **</personne>** assistent à (...)

# NE recognition

- recognition
  - identification
  - categorization



- entity linking (*désambiguïsation*)

## Text example

(...) et <personne ref="Barack\_Obama"> Obama </personne>  
assistent à (...)

# Task definition: which categories?

## Which categories?

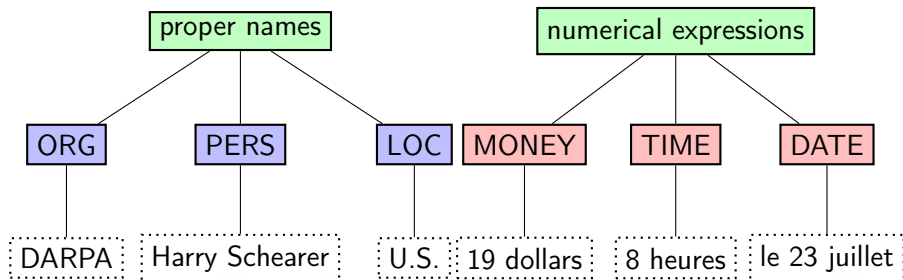
- no consensus beyond the 3 classical categories
  - category Misc in some campaigns (CoNLL, HAREM)
- dependency to the kind of targeted application
  - class granularity: length  $\neq$  height
- reference to existing datasets (evaluation campaigns)

## Categories ranges?

- which instances?
  - ☺ Matteo Renzi, la famille Kennedy
  - ☹ Zorro, Hercule, les italiens
  - ☹ Mickey, Bison futé, le Prince Charmant
- ambiguity, especially metonymy
  - la **France**<sub>ORG</sub> vote contre un traité d'interdiction des armes nucléaires (ou **France**<sub>LIEU</sub> ?)

# NE categories

MUC-6/7





# NE categories

ACE (2002-2008)

## Characteristics

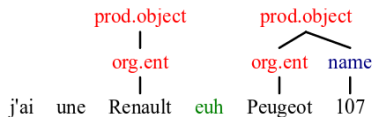
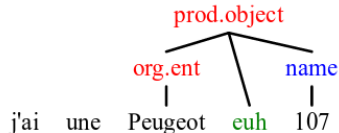
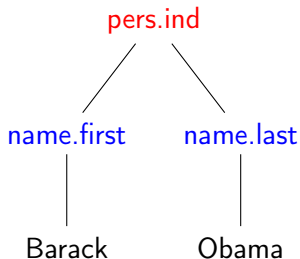
- new types of entities: nominal or pronominal mentions
- 7 types, including Person, Organization, Location and:
  - Geo-political Entity
    - **France**<sub>ORG</sub> signed a treaty with Germany last week.
    - The world leaders met in **France**<sub>LIEU</sub> yesterday.
    - **France**<sub>GPE</sub> produces better wine than New Jersey.
  - Facility (*Aéroport Charles de Gaulle*)
  - Vehicle (*les hélicoptères militaires ont...*)
  - Weapon (*des missiles sol-air ont été tirés*)
- hierarchy: subtypes  
for example for Person: Individual, Group and Indeterminate (if the context does not enable to disambiguate)

# NE categories

Quaero (2011/2012)

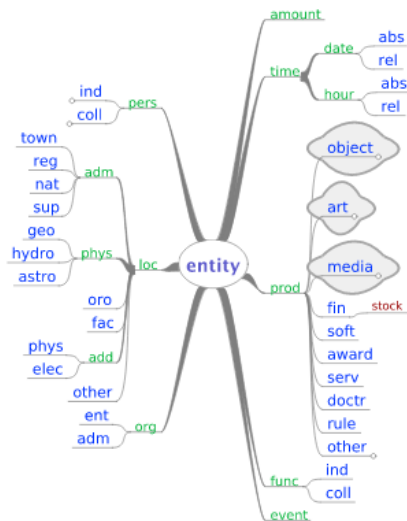
## Characteristics

- new types: products, functions
- additional structure: composition
  - metonymy taken into account: two levels of annotation
- annotation adapted to oral corpora (disfluences)

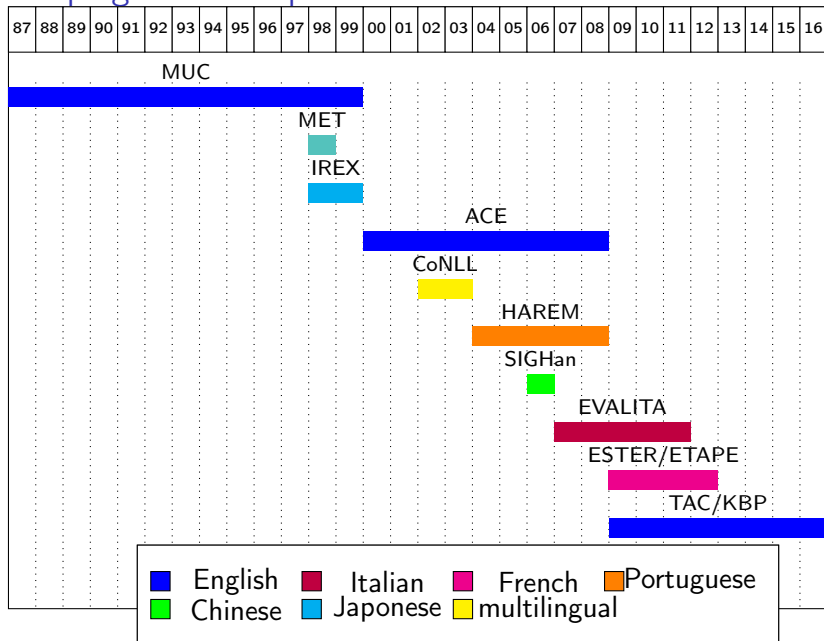


# NE categories

Quaero (2011/2012)



# Campaigns and corpora



# Task definition: which mentions?

## Annotation range

- mention forms
  - 😊 proper nouns: *Jacques Chirac*
  - 😞 nicknames, nominal phrases, pronouns: *Chichi, l'ancien président, il*
- boundaries
  - determiners: *les Rolling Stones, La Mecque, Le téléphone sonne*
  - functions: *le président Obama, l'Abbé Pierre*
  - titles: *Monsieur Fillon, Professeur Paolucci*
  - generation: *Benoît XVI, Bush Jr.*
- coordination
  - Bill and Hillary Clinton flew to Chicago last month. (ellipse partielle)
  - M. et Mme. Chirac en thalasso à Biarritz. (ellipse totale)
  - *Bill and Hillary Clinton*<sub>PERS</sub> vs *Bill*<sub>PERS</sub> and *Hillary Clinton*<sub>PERS</sub>
- imbrication
  - *Université Lyon 2, Comité Exécutif d'Orange*
  - *Université Lyon 2*<sub>ORG</sub> vs *Université Lyon*<sub>LIEU</sub> *2*<sub>ORG</sub> (structured entity)

# NE recognition

## Definition

Automatically **identify** et **classify** named entities in texts

## Examples of difficulties

- homonymy (same type or different type)  
→ JFK : person(s) or airport?, Paris
- metonymy  
→ Washington, l'Élysée : location (city) or organization?

## Which clues for NE recognition?

Laurent Courtois-Courret, délégué syndical SUD au centre Qualipel, à Vitry-sur-Seine (Val-de-Marne), a écopé de dix jours de mise à pied disciplinaire avec retenue de salaire.

# Internal features

- case
  - $mRNA = xXXX$ ,  $CPA1 = XXXd$
- character  $n$ -grams
  - *Cotrimoxazole* → drug, *Leuville-sur-Orge* → location
  - *Twilight - Chapitre 3: hésitation* → movie
- words
  - la Banque Populaire
  - l'avenue des Champs-Élysées
  - Benoît XVI (generation item)
- acronym or ampersand
  - Crédit Agricole SA
  - Standard & Poor's
  - F. Hollande
- gazettiers (first names for example), word clusters, word embeddings (*plongements lexicaux*)
  - François Hollande

# External features

- context of the entity
- additional informations or specific properties
  - **Monsieur** Hollande
  - **Mme** Michel
  - **Général** Leclerc
  - le **groupe** Sanofi
  - the Coca Cola **company**
- often given only for the first occurrence of the entity



# Symbolic systems

## Standard components

- Recognition of triggers and entities from gazetteers
- Cascaded regular expressions

## Example of a rule

*Université* + *de* + CityName  $\Rightarrow$  Organization

## Example of an entity recognized by this rule

Université de Nantes

## Limits

- low recall: gazetteers incomplete, evolutions, partial entities (*Obama*), noisy texts...
- ambiguities (homonymy and metonymy)

# Supervised learning systems

## NER as a classification problem

- training
  - create a representative corpus
    - need for many annotated examples!
  - annotate each token
  - choose features adapted to the classes and texts
  - train a classifier to predict the tags from tokens
- test
  - annotate each token
  - evaluate

token	cap	punct	firstname	pos	chunk	tag
U.N.	1	1	0	NNP	B-NP	B-ORG
official	0	0	0	NN	I-NP	O
Rolf	1	0	1	NNP	I-NP	B-PER
Ekeus	1	0	0	NNP	I-NP	I-PERS
heads	0	0	0	VBZ	I-VP	O
for	0	0	0	IN	B-PP	O
Baghdad	1	0	0	NNP	I-PP	B-LOC

# Supervised learning systems

## Standard features

- Words
  - current
  - word substrings
  - previous
  - next
- Other learned linguistic information
  - POS tags

## Annotation models

- independent tags non adapted
- annotation of tag sequences with a reading sense
  - limits: fixed window, error propagation
- sequence annotation (CRFs)

# NER today

## Objective

do without a priori knowledge and attribute selection

- deep neural networks [Collobert et al., 2011]

## Results

- [Lample et al., 2016]: LSTM-CRF, no external data
- [Guo et al., 2014, Passos et al., 2014]: CRFs + word embeddings
- $F1 \simeq 0.90$  on CoNLL 2003 data for English (PER, LOC, ORG, MISC)

# Evaluation

- $tp$  (true positives) = entities correctly recognized
- $fp$  (false positives) = entities falsely recognized
- $fn$  (false negatives) = entities not recognized

## Standard metrics

- Precision =  $\frac{tp}{tp+fp}$   
→ entities that were correctly annotated from all entities annotated by the system
- Recall =  $\frac{tp}{tp+fn}$   
→ entities correctly annotated from all entities that should have been annotated

# Evaluation

## Reference

<personne>Jean-Yves Le Drian</personne> engage ses homologues à "parler d'abord de manière européenne" sur le plan militaire.

## Hypothesis (system output)

<personne>Jean-Yves</personne> Le Drian engage ses homologues à "parler d'abord de manière européenne" sur le plan militaire.

## Disadvantage of these metrics for named entities

- *Jean-Yves* falsely recognized as an entity  
→ false positive
- *Jean-Yves Le Drian* not recognized  
→ false negative

## Adapted metrics

- R : # entities in the reference
- H : # entities in hypothesis (= system output)
- C : # correct entities (= true positives)
- T : # entities with correct boundaries but wrong category
- F : # entities with correct category but wrong boundaries
- TF : # entities with wrong type and boundaries
- I : # inserted entities insérées (= false positives)
- D : # forgotten entities (= false negatives)

### Adapted metrics

→ partial recognition = half correct

- Precision =  $\frac{C+0.5 \times (T+F)}{H}$
- Recall =  $\frac{C+0.5 \times (T+F)}{R}$
- Slot Error Rate =  $\frac{D+I+TF+0.5 \times (T+F)}{R}$

# Evaluation example

## Reference (manual annotation)

<personne>Bertrand Delanoë</personne> a été élu maire de  
<lieu>Paris</lieu>.

## Hypothesis 1 (system 1)

<personne>Bertrand Delanoë</personne> a été élu maire de  
<personne>Paris</personne>.

$$\text{SER} = (0 + 0 + 0 + 0,5 * (1 + 0)) / 2 = 0,25$$

## Hypothesis 2 (system 2)

<personne>Bertrand</personne> Delanoë a été élu maire de  
<personne>Paris</personne>.

$$\text{SER} = (0 + 0 + 0 + 0,5 * (1 + 1)) / 2 = 0,5$$



## My example

En 1985<sub>DATE</sub> — elle n'a que 19 ans —, Cecilia Bartoli<sub>PERS</sub> se fait connaître en France.



Cecilia Bartoli



1985

# Entity linking - definition

## Entity linking (*désambiguïsation/résolution/liaison*)

Given a knowledge base, chose the entity corresponding to the mention (*referent*)

### Text to analyze

In a grim preview of the discontent that may cloud at least the outset of the next president's term, Hillary Clinton and Donald J. Trump are seen by a majority of voters as unlikely to bring the country back together after this bitter election season.

With more than eight in 10 voters saying the campaign has left them repulsed rather than excited, the rising toxicity threatens the ultimate victor. Mrs. Clinton, the Democratic candidate, and Mr. Trump, the Republican nominee, are seen as dishonest and viewed unfavorably by a majority of voters.

# Entity linking - definition

## Entity linking (*désambiguïsation/résolution/liaison*)

Given a knowledge base, chose the entity corresponding to the mention (*referent*)

### Expected result

In a grim preview of the discontent that may cloud at least the outset of the next president's term, **Hillary Clinton**<sub>*Hillary\_Clinton*</sub> and **Donald J. Trump**<sub>*Donald\_Trump*</sub> are seen by a majority of voters as unlikely to bring the country back together after this bitter election season. With more than eight in 10 voters saying the campaign has left them repulsed rather than excited, the rising toxicity threatens the ultimate victor. **Mrs. Clinton**<sub>*Hillary\_Clinton*</sub>, the **Democratic**<sub>*Democratic\_Party\_(United\_States)*</sub> candidate, and **Mr. Trump**<sub>*Donald\_Trump*</sub>, the **Republican**<sub>*Republican\_Party\_(United\_States)*</sub> nominee, are seen as dishonest and viewed unfavorably by a majority of voters.

# Entity linking - definition

## Entity linking (*désambiguïsation/résolution/liaison*)

Given a knowledge base, chose the entity corresponding to the mention (*referent*)

In a grim preview of the discontent that may cloud at least the outset of the next president 's term , Hillary Clinton and Donald J. Trump are seen by a majority of voters as unlikely to bring the country back together after this bitter election season .

With more than List\_of\_neighborhoods\_of\_the\_District\_of\_Columbia\_by\_ward (10) eight in 10 voters saying the campaign has left them repulsed rather than excited , the rising toxicity threatens the ultimate victor .

Hill Clinton , the Democratic Party (United States) candidate , and Mr. Donald Trump , the Republican Party (United States) nominee , are seen as dishonest and viewed unfavorably by a majority of voters .

CoreNLP

# Entity linking - definition

## Entity linking (*désambiguïsation/résolution/liaison*)

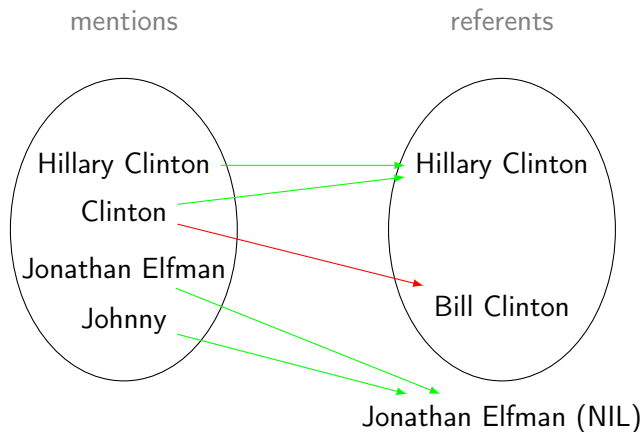
Given a knowledge base, chose the entity corresponding to the mention (*referent*)

In a grim preview of the discontent that may cloud at least the outset of the next president's term, [Hillary Clinton](#) and Donald J. Trump are seen by a majority of voters as unlikely to bring the country back together after this bitter election season.

With more than eight in 10 voters saying the campaign has left them repulsed rather than excited, the rising toxicity threatens the ultimate victor. Mrs. Clinton, the [Democratic](#) candidate, and Mr. Trump, the [Republican](#) nominee, are seen as dishonest and viewed unfavorably by a majority of voters.

## DBPedia Spotlight

# Entity linking - difficulties



# Entity linking - principles

## Steps

- ① possible mention detection
  - often based on NE recognition
- ② candidate selection
  - graphical proximity to labels, links texts, queries that lead to the Wikipedia pages, Wikipedia disambiguation pages
- ③ candidate ranking  
WSD / Wikipedia
  - mention: distance to the referents' labels
  - referent: popularity (most frequent, Wikipedia page with most links...)
  - local context of the mention: textual similarity with Wikipedia pages, links...
  - global context of the mention (document): other entities (collective disambiguation), coreference

## Tâche Entity Discovery and Linking

- Discovery: detect and annotate mentions
  - classes: LOC, ORG, PER, FAC, GPE;
  - mentions: EN, noms, posts authors
- Linking: attach mention clusters to a KB
- Difficulties (KBP 2015):
  - detection of common names and acronyms
  - rare entities
  - popularity bias
  - general knowledge
  - informal language
  - lack of coherence between NE type and referent
- $F1 \simeq 0.60$  for EL in English 2015



## My example

En 1985<sub>DATE</sub> — elle n'a que 19 ans —, Cecilia Bartoli<sub>PERS</sub> se fait connaître en France.



Cecilia Bartoli



1985

## 1 Introduction

## 2 Entities

- Definitions
- NE recognition
- Entity linking

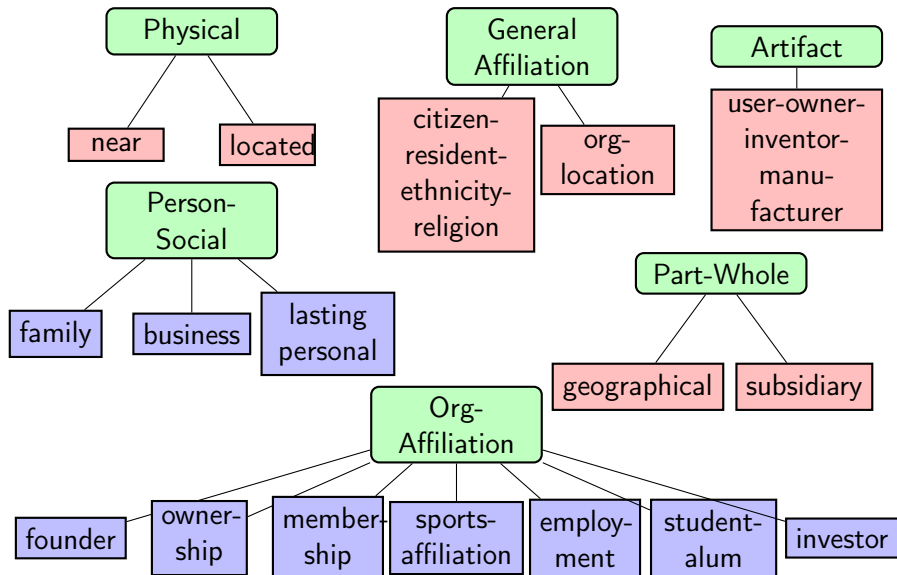
## 3 Relations

- Definitions
- Relation extraction
- Supervised methods
- Semi-supervised methods

## 4 Conclusion

# A few sets of relations

ACE 2005



# A few sets of relations

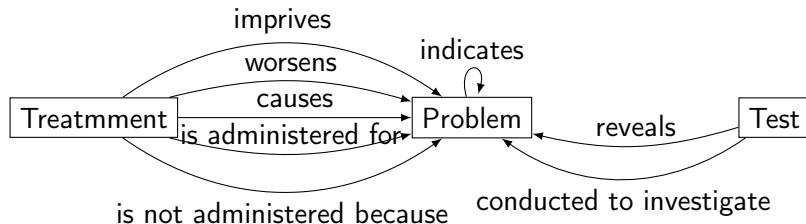
## SemEval 2010 task 8

Type	Example
Cause-Effect	The <i>news</i> brought about a <i>commotion</i> in the office.
Instrument-Agency	<i>Carpenters</i> build many things from <i>wood</i> .
Product-Producer	The <i>government</i> built 10,000 new <i>homes</i> .
Content-Container	I emptied the <i>wine bottle</i> into my glass.
Entity-Origin	It involves a spectator choosing a <i>card</i> from the <i>deck</i> .
Entity-Destination	He sent his <i>painting</i> to an <i>exhibition</i> .
Component-Whole	Feel free to download the first <i>chapter</i> of the <i>book</i> .
Member-Collection	A person who is serving on a <i>jury</i> is known as <i>juror</i> .
Message-Topic	Mr Cameron asked a <i>question</i> about tougher <i>sentences</i> for people carrying knives.

# A few sets of relations

i2b2 2010

## Natural Language Processing Challenge for Clinical Records



# A few sets of relations

## Freebase

### Most frequent Freebase relations

- /people/person/nationality
- /location/location/contains
- /people/person/location
- /people/person/place\_of\_birth
- /dining/restaurant/cuisine
- /business/business\_chain/location
- /biology/organism\_classification\_rank
- /film/film/genre
- /film/film/language
- /biology/organism\_higher\_classification
- /film/film/country
- /film/writer/film

# Relation

## Relation characteristics

- between concepts or concept instances
- hierarchical or not
- include events or only binary relations
- “real world” relations or factivity taken into account

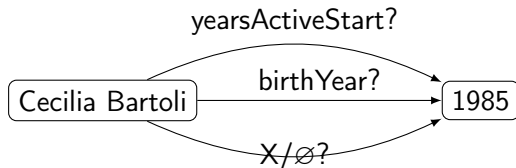
# Relation extraction

## Definition

Given two (or more) entities, determine

- whether there is a relation between them
- what kind of relation

En 1985<sub>DATE</sub> — elle n'a que 19 ans —, Cecilia Bartoli<sub>PERS</sub> se fait connaître en France.





# Difficulties

## Variability of expression of relations

- En 1985 — elle n'a que 19 ans —, Cecilia Bartoli se fait connaître en France lors d'un concert organisé par l'Opéra de Paris en hommage à Maria Callas.
- C'est déjà une longue carrière que celle de Cecilia Bartoli. Elle débute en 1985, à Rome. Elle a dix-neuf ans et incarne la pétulante Rosina du « Barbier de Séville ».
- En 1985, une tournée en Allemagne de l'Est et un gala télévisé à Paris en hommage à Maria Callas suffisent à attirer l'attention de tous – y compris celle de chefs d'orchestre prestigieux comme Daniel Barenboim, Claudio Abbado, Simon Rattle, Herbert von Karajan – sur cette jeune cantatrice.

wikipédia, les échos et encyclopédie universalis examples

# Simple methods

## Cooccurrence

but ambiguity

- person - date: start date of the career, birthdate, other?
- treatment - illness: cures? prevents? side effect?

## Lexico-syntactic patterns

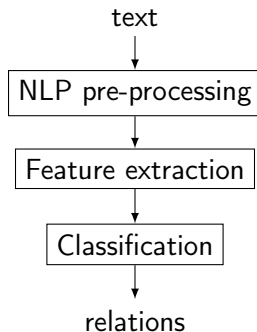
- for example for birthyear relation:
  - *PERSON*, born in *DATE*
  - *PERSON* (*DATE*-)
  - *PERSON* is a NP born in *DATE*
- have to be manually written for each relation
  - possible automatic acquisition
  - bootstrapping
- oriented recall or precision

# Supervised relation extraction

## Classification problem

- binary or multiclass classification
- positive and negative training examples

# Basic supervised method



# Context

En 1985 — elle n'a que 19 ans —, Cecilia Bartoli se fait connaître en France

# Features

- words (or lemmss)
  - from different context parts
  - bag of word and n-grams
  - syntactic head and concatenation
- entity types
  - entity types and concatenation
- syntactic information
  - constituents path
  - dependency path
- external resources
  - country list, triggers...

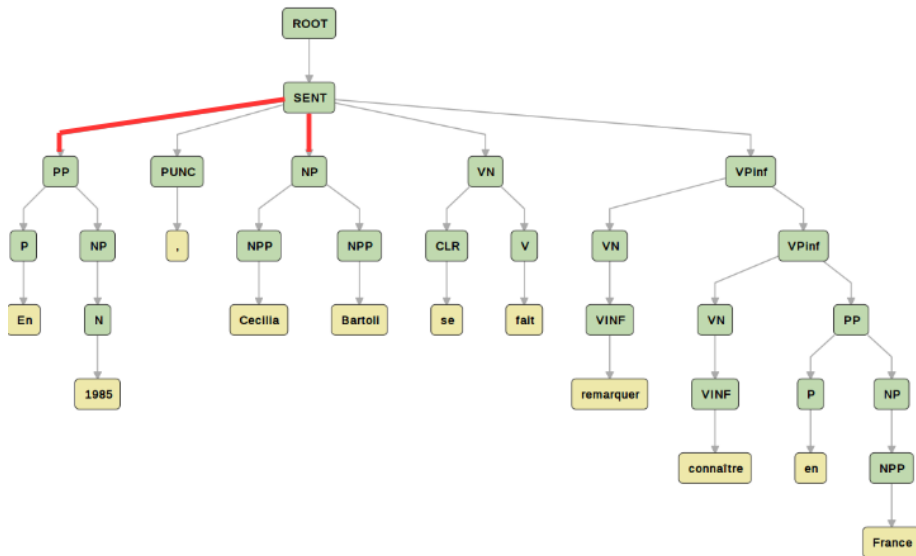
# Feature example

En 1985 — elle n'a que 19 ans —, Cecilia Bartoli se fait connaître en France

- words
  - before  $e_1$ : {En}
  - between entities (bow): {elle, n', a, que 19, ans}
  - after  $e_2$ : {se, fait, connaître, en, France}
  - head  $e_1$ : 1985
  - head  $e_2$ : Bartoli
- types
  - type  $e_1$ : DATE
  - type  $e_2$ : PERSON
  - concatenation: DATEPERSON
- syntax
  - constituents: PP - SENT - NP
  - dependencies: nmod -subj

# Example

## Constituents parse tree





# Structured representations

## Which attributes?

- intuition
- experiments

## Using structured representations

Definition of adapted similarity metrics: parse tree kernels

## Experiments

- constituent tree [Zelenko et al., 2003]
- dependency tree [Culotta and Sorensen, 2004]
- shortest path between entities [Bunescu and Mooney, 2005]

# Limits of these approaches

## Disadvantages of previous methods

- classify quality strongly dependent on pre-processing
- large annotated corpora
  - even if crowdsourcing possible [Liu et al., 2016]
- corpora imbalance
- lack of generalization

## Getting rid of pre-processing

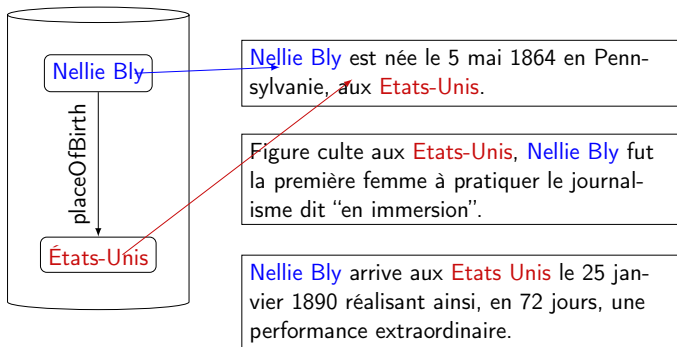
### Deep neural networks

- input: words,  $n$ -grams + positions + word embeddings
- network: RNN or CNN

# Distant supervision

## Objective

- automatically annotate training examples
- ← KBs
- then standard methods



relations dbpedia

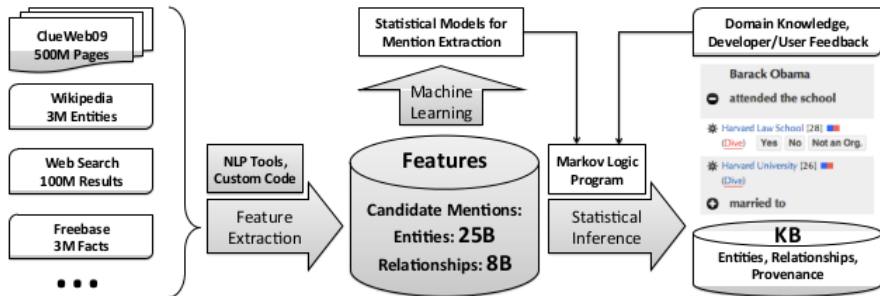
# Distant supervision

## Hypotheses

- any sentence that contains a pair of entities that participate in a known KB relation is likely to express that relation in some way [Mintz et al., 2009, Wu and Weld, 2007, Niu et al., 2012]
- multi-instances learning problem [Riedel et al., 2010]: at least one of the sentences contains a relation mention
- Several relations may exist between two entities [Hoffmann et al., 2011]

- 😊 (mostly) domain independent
- 😊 scales up well
- 😞 only for context free relations
- 😞 depends on the quality of NER

# DeepDive [Niu et al., 2012]



# Open information extraction

## Principle

Keep the relation expression from the text

Ada Lovelace was one of the  
earliest computer scientists.

The second tunnel boring machine will be named Ada after  
Ada Lovelace who was one of  
the earliest computer scientists.



Ada Lovelace

was one of

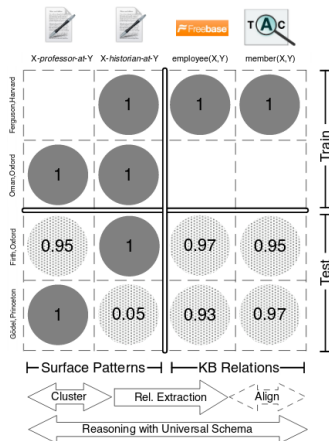
the earliest computer scientist

(exemples de <http://openie.allenai.org/> )

# Open information extraction

## Limits

- Non normalized relations
  - [Angeli et al., 2015]: cooccurrences relations OIE and KBP in corpus
  - [Riedel et al., 2013, Verga et al., 2016]: implications between relations



# Current difficulties

- rare relations (in texts)
- contextual relations
- factivity
- source: reliability, fiction...
- common sense knowledge
- NL query: several relations



## My example

En 1985<sub>DATE</sub> — elle n'a que 19 ans —, Cecilia Bartoli<sub>PERS</sub> se fait connaître en France.



# Conclusion

## A few points

- more and more explicit knowledge
- virtuous circle between IE and semantic annotation
- remain complementary
  - query both types of resources
  - real interaction between reasoning on texts and knowledge

# References

- Cours de Christopher Manning et Dan Jurafsky sur le traitement automatique des langues (Natural Language Processing)
- Livre Les entités nommées pour le traitement automatique des langues, Damien Nouvel, Maud Ehrmann et Sophie Rosset, 2015

# Références I



Angeli, G., Premkumar, M. J., and Manning, C. D. (2015).  
Leveraging linguistic structure for open domain information extraction.  
*In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing.*



Bunescu, R. and Mooney, R. (2005).  
A shortest path dependency kernel for relation extraction.  
*In Proceedings of the conference on human language technology and empirical methods in natural language processing.*



Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011).  
Natural language processing (almost) from scratch.  
*Journal of Machine Learning Research*, 12(Aug):2493–2537.



Culotta, A. and Sorensen, J. (2004).  
Dependency tree kernels for relation extraction.  
*In Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 423. Association for Computational Linguistics.



Guo, J., Che, W., Wang, H., and Liu, T. (2014).  
Revisiting embedding features for simple semi-supervised learning.  
*In EMNLP*, pages 110–120.

# Références II



Hoffmann, R., Zhang, C., Ling, X., Zettlemoyer, L., and Weld, D. S. (2011). Knowledge-based weak supervision for information extraction of overlapping relations. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 541–550. Association for Computational Linguistics.



Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., and Dyer, C. (2016). Neural architectures for named entity recognition. In *Proceedings of NAACL-HLT 2016*.



Liu, A., Soderland, S., Bragg, J., Lin, C. H., Ling, X., and Weld, D. S. (2016). Effective Crowd Annotation for Relation Extraction. In *Proceedings of NAACL-HLT 2016*.



Mintz, M., Bills, S., Snow, R., and Jurafsky, D. (2009). Distant supervision for relation extraction without labeled data. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2-Volume 2*, pages 1003–1011. Association for Computational Linguistics.



Niu, F., Zhang, C., Ré, C., and Shavlik, J. W. (2012). DeepDive: Web-scale Knowledge-base Construction using Statistical Learning and Inference. *VLDS*, 12:25–28.

# Références III



Passos, A., Kumar, V., and McCallum, A. (2014).  
Lexicon infused phrase embeddings for named entity resolution.  
*In Proceedings of the Eighteenth Conference on Computational Language Learning.*



Riedel, S., Yao, L., and McCallum, A. (2010).  
Modeling Relations and Their Mentions without Labeled Text.  
*In Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 148–163. Springer.



Riedel, S., Yao, L., McCallum, A., and Marlin, B. M. (2013).  
Relation extraction with matrix factorization and universal schemas.  
*In Proceedings of NAACL-HLT 2013.*



Verga, P., Belanger, D., Strubell, E., Roth, B., and McCallum, A. (2016).  
Multilingual relation extraction using compositional universal schema.  
*In Proceedings of NAACL-HLT 2016.*



Wu, F. and Weld, D. S. (2007).  
Autonomously semantifying wikipedia.  
*In Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*, pages 41–50. ACM.



Zelenko, D., Aone, C., and Richardella, A. (2003).  
Kernel methods for relation extraction.  
*Journal of machine learning research*, 3(Feb):1083–1106.