

A GUIDED TOUR ON NLP APPLICATIONS FROM DOCUMENT CLASSIFICATION TO SENTENCE UNDERSTANDING

Agro-IODAA, semestre 1

Vincent Guigue



Around textual data: a tour on NLP tasks

Think as a data-scientist:

- No magic !
- For each application :
 - 1 Identify the problem class (regression, classification, ...)
 - 2 Find global Input / output
 - 3 Think about the data format (I/O)
 - 4 Find a model to deal with this data
 - 5 **Optimize all parameters & hyper-parameters**
 - **Conduct an experiment campaign**

$$\{(\mathbf{x}_i, y_i)\}_{i=1,\dots,N} \quad \Rightarrow \text{ learn:} \quad f_{\theta, \phi}(\mathbf{x}) \approx y$$

Let's have a tour on NLP applications !

Your role = identifying (x, y) for all applications



Different levels of analysis

1 At the document level

- Topic/sentiment classification
- etc...

2 At the paragraph / sentence level

- Co-reference resolution
- Named Entity Recognition / Relation Extraction
- Sentiment classification (also)

3 At the word level

- Synonymy

4 At the stream level

- Topic detection & tracking

⇒ Different levels often refer to different format of data (tabular, sequence, stream)

AT THE DOCUMENT LEVEL



Tasks

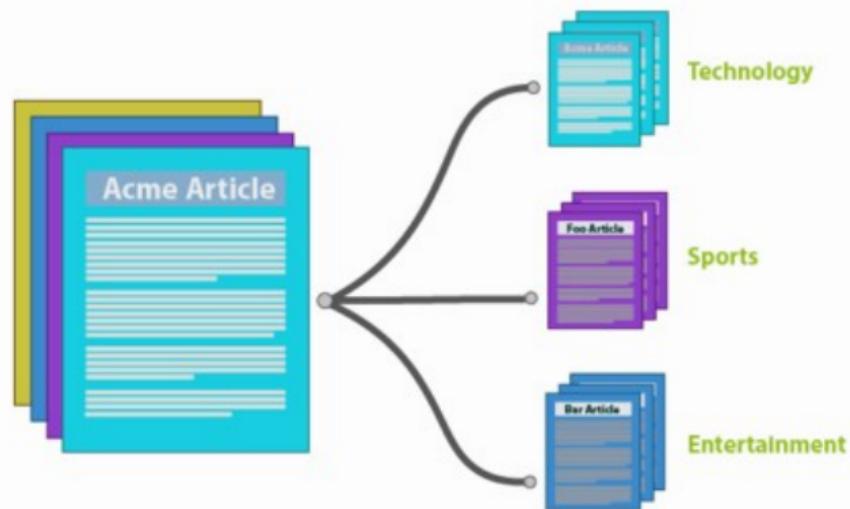
- **Indexing**
 - cf Information Retrieval
- **Classification / filtering**
 - topic
- **Counting / survey**
 - Sentiment classification
- **Clustering (topic analysis)**
 - Unsupervised formulation
 - Large corpus primary exploration
 - Lexical field extraction





Tasks

- **Indexing**
 - cf Information Retrieval
- **Classification / filtering**
 - topic
- **Counting / survey**
 - Sentiment classification
- **Clustering (topic analysis)**
 - Unsupervised formulation
 - Large corpus primary exploration
 - Lexical field extraction



Tasks

- **Indexing**
 - cf Information Retrieval
 - **Classification / filtering**
 - topic
 - **Counting / survey**
 - Sentiment classification
 - **Clustering (topic analysis)**
 - Unsupervised formulation
 - Large corpus primary exploration
 - Lexical field extraction





Tasks

- **Indexing**
 - cf Information Retrieval
- **Classification / filtering**
 - topic
- **Counting / survey**
 - Sentiment classification
- **Clustering (topic analysis)**
 - Unsupervised formulation
 - Large corpus primary exploration
 - Lexical field extraction



Still at the document level?

■ Document Segmentation

- Sentence classification
 - Link with stream
 - Break detection

■ Automated summary

- Sentence extraction
 - Generative architecture

■ Translation

- #### ■ Mostly at the sentence level



Segment Sparse Model

PISA Model



Still at the document level?

■ Document Segmentation

- Sentence classification
- Link with stream
- Break detection

■ Automated summary

- Sentence extraction
- Generative architecture

■ Translation

- Mostly at the sentence level

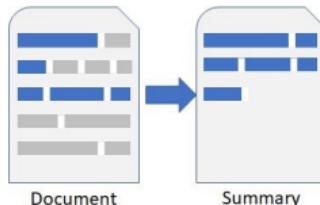
Source Document



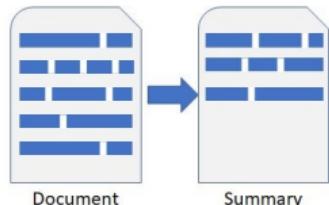
Extractive Summary

To summarize is to reduce in complexity, and hence in length, while retaining some of the essential qualities of the original. This paper focuses on document extracts, a particular kind of computed document summary. Document extracts consisting of roughly 20% of the original can be as

Extractive Summarization



Abstractive Summarization



AT THE PARAGRAPH LEVEL



At the paragraph level

■ Question Answering (QA)

- Classification formulation
- Extraction approach (Siri/Google assistant)
- An emerging task... For several other tasks!

■ Coreference resolution

■ Information extraction

- Mainly a sentence level task...
- Sometimes at the paragraph level (with coreference / QA)

■ Dialog State Tracking

- Chatbot...

WARNING: a mix of old & new tasks

Passage Sentence

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity.

Question

What causes precipitation to fall?

Answer Candidate

gravity

- Between question and answer

cause---gravity

precipitation---gravity

fall---gravity

what---gravity

At the paragraph level

■ Question Answering (QA)

- Classification formulation
- Extraction approach (Siri/Google assistant)
- An emerging task... For several other tasks!

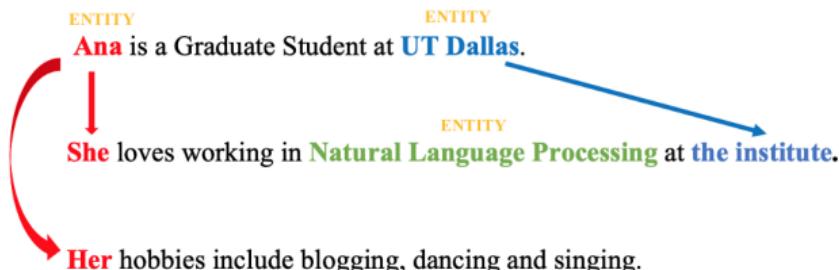
■ Coreference resolution

■ Information extraction

- Mainly a sentence level task...
- Sometimes at the paragraph level (with coreference / QA)

■ Dialog State Tracking

- Chatbot...



WARNING: a mix of old & new tasks



At the paragraph level

■ Question Answering (QA)

- Classification formulation
- Extraction approach (Siri/Google assistant)
- An emerging task... For several other tasks!

■ Coreference resolution

■ Information extraction

- Mainly a sentence level task...
- Sometimes at the paragraph level (with coreference / QA)

■ Dialog State Tracking

- Chatbot...

Sam walks into the kitchen.

Sam picks up an apple.

Sam walks into the bedroom.

Sam drops the apple.

Q: Where is the apple?

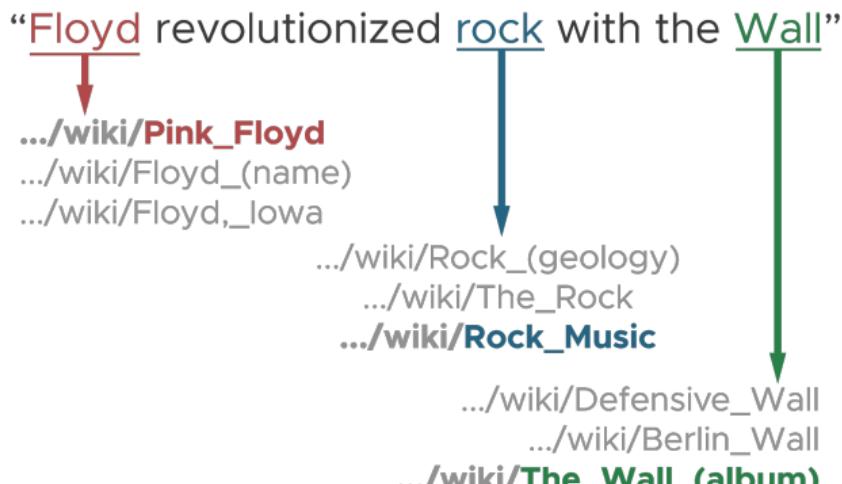
A. Bedroom

WARNING: a mix of old & new tasks

AT THE SENTENCE LEVEL

At the sentence level

- Word/entity resolution
- Part-Of-Speech tagging (POS)
- Chunking (sentence segmentation) =
Syntactic Parsing
- Semantic Role Labeling (SRL)
- Sentiment Analysis (fine grained)
- Named Entity Recognition (NER)
- Information Extraction
 - Entity resolution
 - Relation extraction
 - ⇒ Knowledge Graph
([opt] + ontologies...)
- Machine Translation



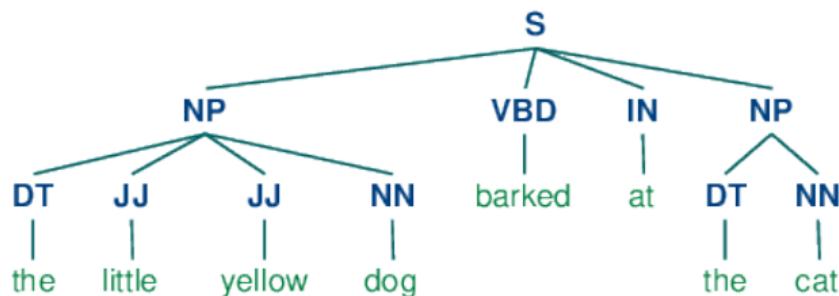
At the sentence level

- Word/entity resolution
- Part-Of-Speech tagging (POS)
- Chunking (sentence segmentation) =
Syntactic Parsing
- Semantic Role Labeling (SRL)
- Sentiment Analysis (fine grained)
- Named Entity Recognition (NER)
- Information Extraction
 - Entity resolution
 - Relation extraction
 - ⇒ Knowledge Graph
([opt] + ontologies...)
- Machine Translation



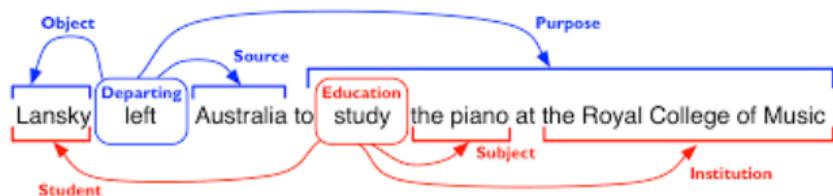
At the sentence level

- Word/entity resolution
- Part-Of-Speech tagging (POS)
- Chunking (sentence segmentation) =
Syntactic Parsing
- Semantic Role Labeling (SRL)
- Sentiment Analysis (fine grained)
- Named Entity Recognition (NER)
- Information Extraction
 - Entity resolution
 - Relation extraction
 - ⇒ Knowledge Graph
([opt] + ontologies...)
- Machine Translation



At the sentence level

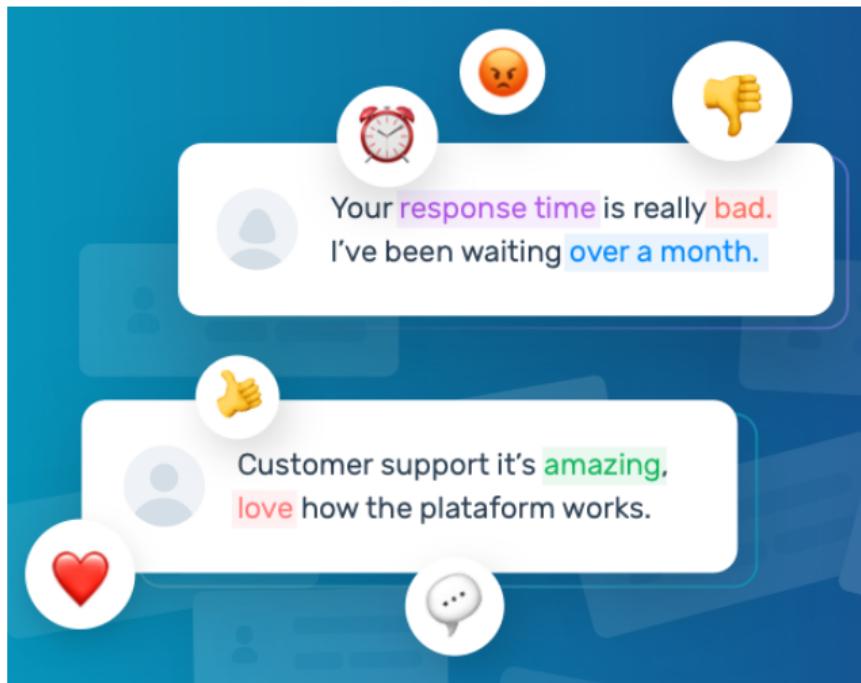
- Word/entity resolution
- Part-Of-Speech tagging (POS)
- Chunking (sentence segmentation) =
Syntactic Parsing
- Semantic Role Labeling (SRL)
- Sentiment Analysis (fine grained)
- Named Entity Recognition (NER)
- Information Extraction
 - Entity resolution
 - Relation extraction
 - ⇒ Knowledge Graph
([opt] + ontologies...)
- Machine Translation





At the sentence level

- **Word/entity resolution**
- **Part-Of-Speech** tagging (POS)
- **Chunking** (sentence segmentation) =
Syntactic Parsing
- **Semantic Role Labeling** (SRL)
- **Sentiment Analysis** (fine grained)
- **Named Entity Recognition** (NER)
- **Information Extraction**
 - Entity resolution
 - Relation extraction
 - ⇒ Knowledge Graph
 - ([opt] + ontologies...)
- **Machine Translation**



At the sentence level

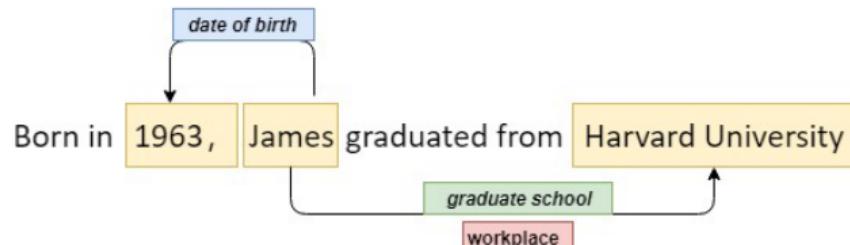
- Word/entity resolution
 - Part-Of-Speech tagging (POS)
 - Chunking (sentence segmentation) = Syntactic Parsing
 - Semantic Role Labeling (SRL)
 - Sentiment Analysis (fine grained)
 - Named Entity Recognition (NER)
 - Information Extraction
 - Entity resolution
 - Relation extraction

⇒ Knowledge Graph
([opt] + ontologies...)
 - Machine Translation

1 Sebastian Thrun **PERSON** started at Google **ORG** in 2007 **DATE**, fe
him seriously. "I can tell you very senior CEOs of major American **NORP**
and turn away because I wasn't worth talking to," said Thrun **PERSON**, I
e higher education startup Udacity, in an interview with Recode **ORG**

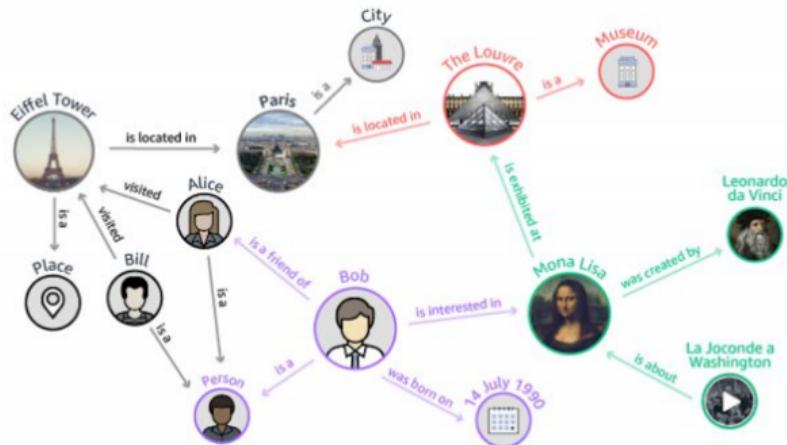
At the sentence level

- Word/entity resolution
- Part-Of-Speech tagging (POS)
- Chunking (sentence segmentation) =
Syntactic Parsing
- Semantic Role Labeling (SRL)
- Sentiment Analysis (fine grained)
- Named Entity Recognition (NER)
- Information Extraction
 - Entity resolution
 - Relation extraction
 - ⇒ Knowledge Graph
([opt] + ontologies...)
- Machine Translation



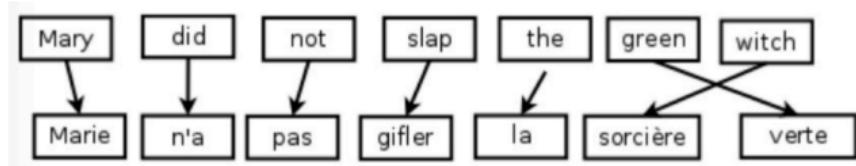
At the sentence level

- Word/entity resolution
- Part-Of-Speech tagging (POS)
- Chunking (sentence segmentation) = Syntactic Parsing
- Semantic Role Labeling (SRL)
- Sentiment Analysis (fine grained)
- Named Entity Recognition (NER)
- Information Extraction
 - Entity resolution
 - Relation extraction
 - ⇒ Knowledge Graph
([opt] + ontologies...)
- Machine Translation



At the sentence level

- Word/entity resolution
- Part-Of-Speech tagging (POS)
- Chunking (sentence segmentation) =
Syntactic Parsing
- Semantic Role Labeling (SRL)
- Sentiment Analysis (fine grained)
- Named Entity Recognition (NER)
- Information Extraction
 - Entity resolution
 - Relation extraction
 - ⇒ Knowledge Graph
([opt] + ontologies...)
- Machine Translation



AT THE WORD LEVEL

At the word level

■ Syntactic similarities

- spelling correction
- Levenshtein = DTW

■ Semantic distance (synonymy)

- WordNet (& other resources)
- Representation learning

```
e}
Correction "distnce": suggestions [fr]: distance, distancé
Calcul de Aperçu du problème (⌘F8) Aucune solution disponible dan
Calcul de distnce Sémantique
itemize
\item Ressources WordNet
itemize}
Dictionnaires en tous genres
```

At the word level

- Syntactic similarities
 - spelling correction
 - Levenshtein = DTW
- Semantic distance (synonymy)
 - WordNet (& other resources)
 - Representation learning

		H	Y	U	N	D	A	I
	0	1	2	3	4	5	6	7
H	1	0	1	2	3	4	5	6
O	2	1	1	2	3	4	5	6
N	3	2	2	2	2	3	4	5
D	4	3	3	3	3	2	3	4
A	5	4	4	4	4	3	2	3

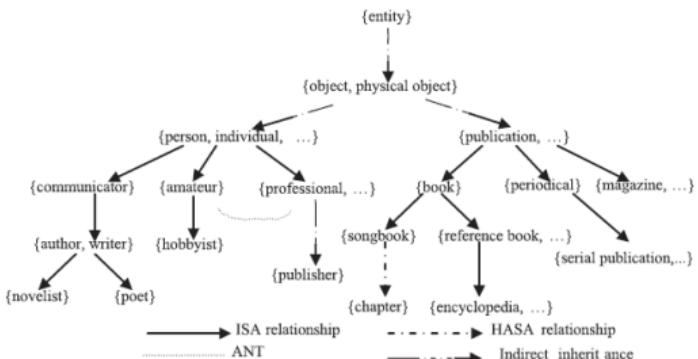
At the word level

■ Syntactic similarities

- spelling correction
- Levenshtein = DTW

■ Semantic distance (synonymy)

- WordNet (& other resources)
- Representation learning



TOWARD COMPOSITE OR MULTIMODAL ANALYSIS

Temporal streams

- Historical tasks: **Topic Detection and Tracking (TDT)**
- = adding a temporal axis to the topic detection
 - topic quantification
 - topic appearance
 - topic vanishing



Temporal streams

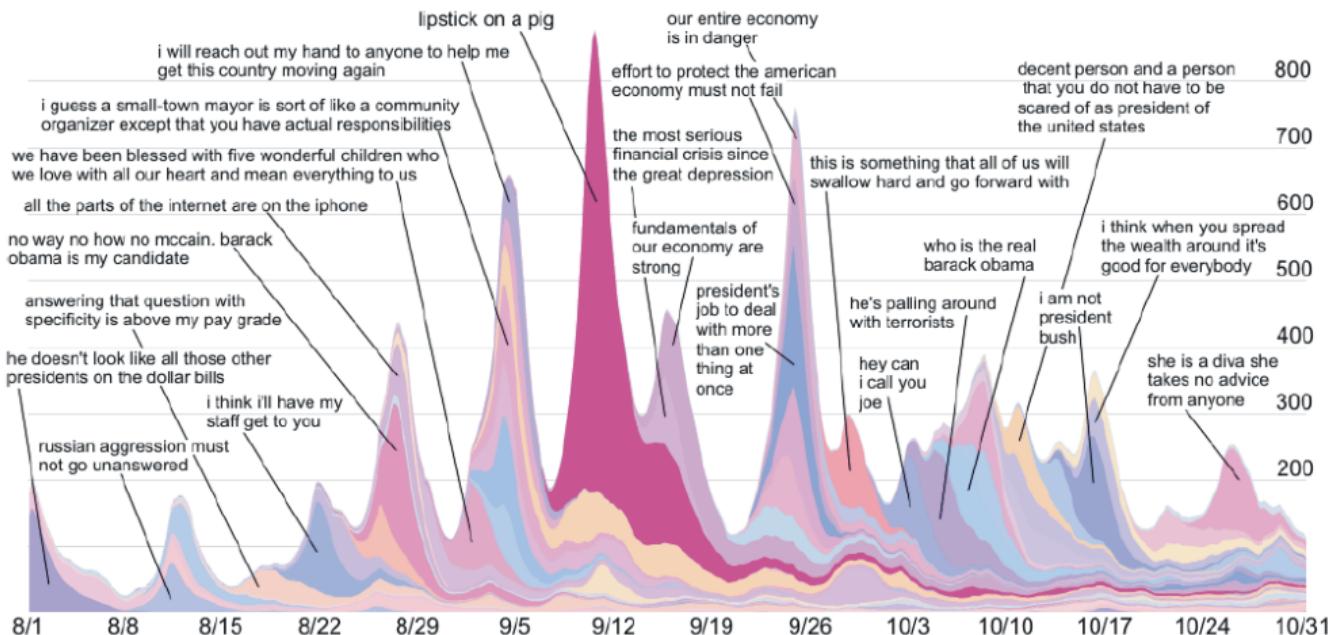


Figure 4: Top 50 threads in the news cycle with highest volume for the period Aug. 1 – Oct. 31, 2008. Each thread consists of all news articles and blog posts containing a textual variant of a particular quoted phrases. (Phrase variants for the two largest threads in each week are shown as labels pointing to the corresponding thread.) The data is drawn as a stacked plot in which the thickness of the strand corresponding to each thread indicates its volume over time. Interactive visualization is available at <http://memetracker.org>.



Profiling

Profiling = recommender system

- modern User Interface (UI)
- Information Access
 - Personalization
 - Active suggestion (user = query)

Recent proposal: mixing user interaction & textual review

- user interaction = best item similarity
- review = in depth textual description

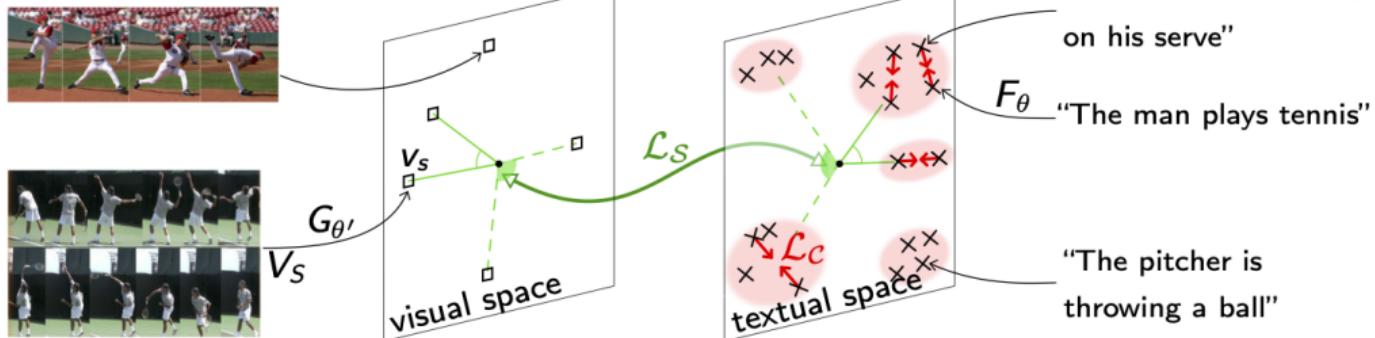
McAuley & Leskovec 2013



Text & Image

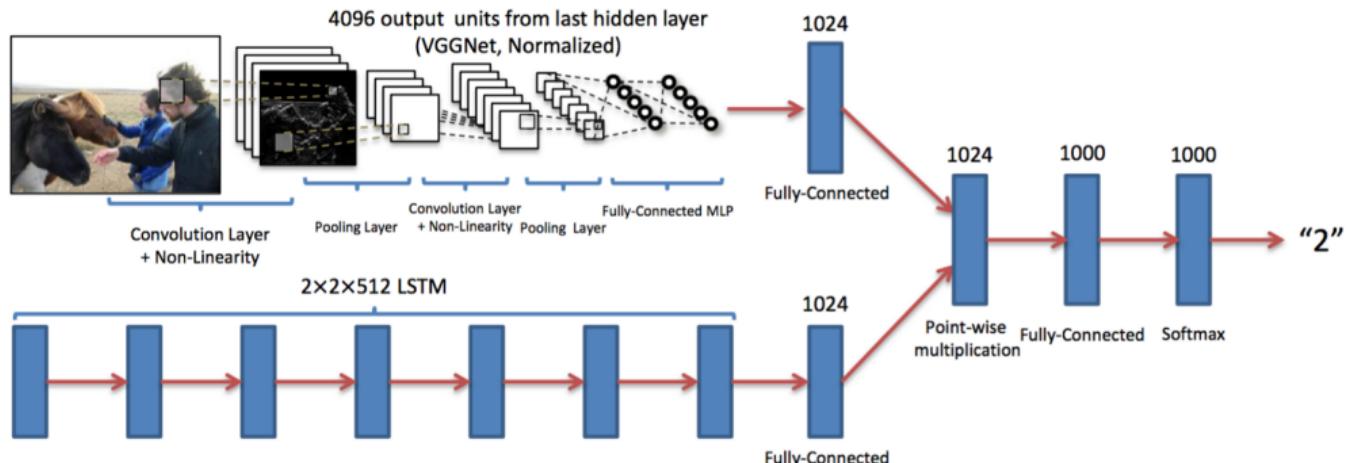
- Visual grounding
- Visual Question Answering
- Captioning
- Image Generation

Word	Teraword	Knext	Word	Teraword	Knext
Spoke	11,577,917	372,042	Hugged	610,040	11,453
Laughed	3,904,519	179,395	Blinked	390,692	21,973
Murdered	2,843,529	16,890	Was late	368,922	31,168
Inhaled	984,613	5,617	Exhaled	168,985	4,052
Breathed	725,034	41,215	Was on time	23,997	14



A Text & Image

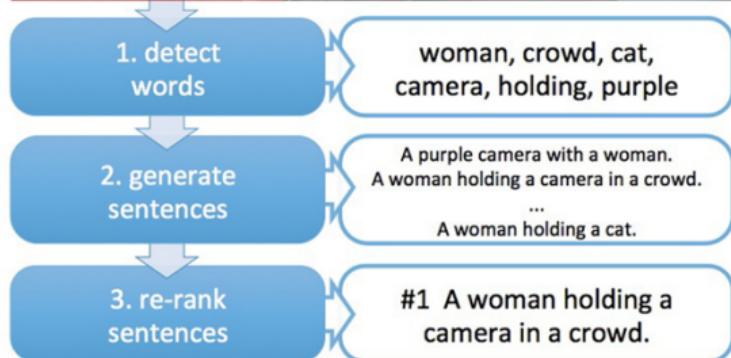
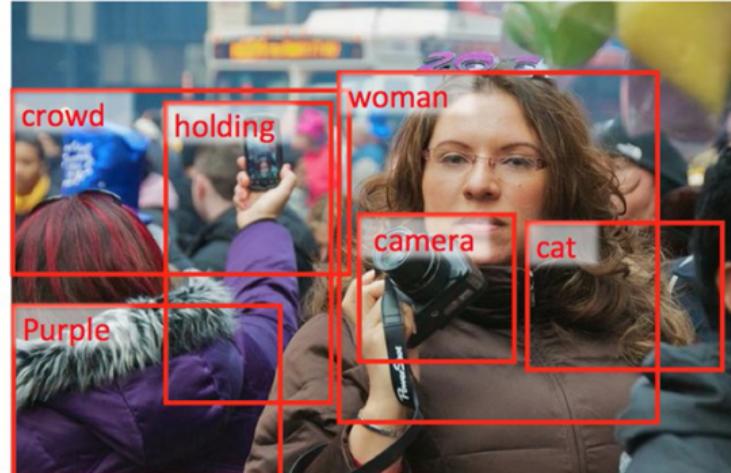
- Visual grounding
- Visual Question Answering
- Captioning
- Image Generation





Text & Image

- Visual grounding
- Visual Question Answering
- Captioning
- Image Generation





Text & Image

- Visual grounding
- Visual Question Answering
- Captioning
- Image Generation

TEXT DESCRIPTION

An astronaut Teddy bears A bowl
of soup

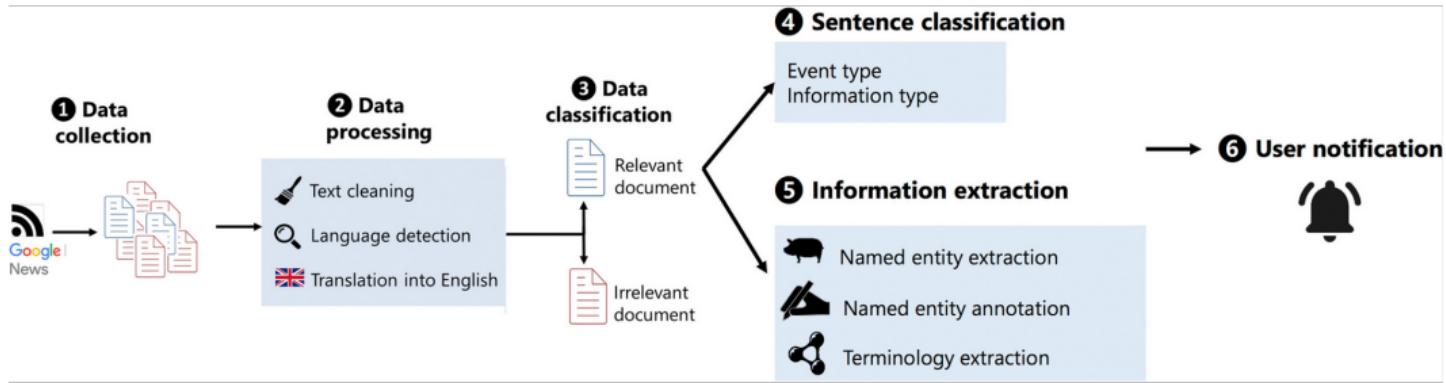
riding a horse lounging in a tropical
resort in space playing basketball
with cats in space

in a photorealistic style in the style
of Andy Warhol as a pencil drawing

DALL-E 2



Epidemiology



Valentin S. et al., PADI-web 3.0: framework for extracting fine-grained information from news for animal disease surveillance.



Event detection

Name
@username
Sairapet is scary. The Bridge is flooding & it's bringing cylinders, fridges etc [ChennaiFloods](#) [Chennaiains](#)



9:21 AM · Dec 2, 2015 · Twitter for Android

Bridge flooding

Name
@username
@ChennaiFloods: Ola opens temporary homes for residents - Times of India bit.ly/1NuIxt9

5:31 PM · Dec 4, 2015 · Buffer

Shelters available

Name
@username
these roads are closed for traffic. Pls pas 1 Sholinganallur to Siruseri closed 2 Thuraipakkam AKDR Golf course to Toll Gate Closed

7:49 PM · Dec 1, 2015 · Twitter for Android

Closed roads

Name
@username
Dear Friends, Pl help by sending boat to 54 and 58 Vivekananda Nagar Street Nesapakkam Chennai..

8:59 AM · Dec 2, 2015 · Twitter Web Client

Help request

Name
@username
Heavy rains continue to batter Chennai Airport shut and trains diverted.



Airport shutdown

Name
@username
#ChennaiFloods
Guys I don't know if ts mi8 help.But Im willing 2 recharge ur mobile no if req. Pls msg back... my no 9952343236 -Karthi

3:03 PM · Dec 3, 2015 · Twitter for Android

Volunteering response

Suwaileh R., Elsayed T., Imran M. (2023). IDRISI-RE: A generalizable dataset with benchmarks for location mention recognition on disaster tweets. *Information Processing & Management*, <https://doi.org/10.1016/j.ipm.2023.103340>



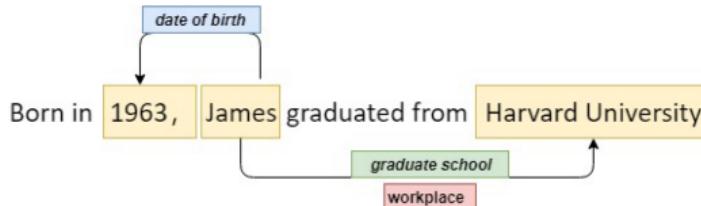
From text (NLP) to Knowledge Bases (KB)

Text = lot of resources ⇒ Universal knowledge?

- Noisy (both at the syntactic & semantic levels)
- Hard to query & exploit (search, knowledge inference,...)

(RDF) knowledge

- Def: Entity 1 (subject/head) Relation (predicate) Entity 2 (target/tail)
- Simpler version: Key / value
- Easy to handle



From text (NLP) to Knowledge Bases (KB)

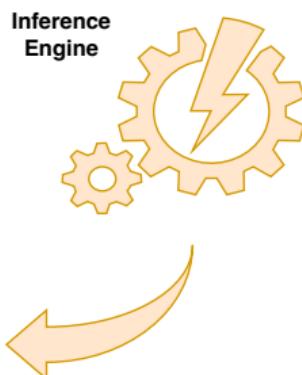
Text = lot of resources \Rightarrow Universal knowledge?

- Noisy (both at the syntactic & semantic levels)
- Hard to query & exploit (search, knowledge inference,...)

From RDF to ontology: adding hierarchies + reasoning

Fact base	
	Barack Obama was born in Honolulu
	Honolulu is the capital of Hawaii
	Hawaii is a state of the USA
	...

Rule base	
	capital => part of
	state of => part of
	born in => citizen of [USA, France]
	...



Barack Obama was born in Hawaii
Barack Obama is American
...



From text (NLP) to Knowledge Bases (KB)

Text = lot of resources ⇒ Universal knowledge?

- Noisy (both at the syntactic & semantic levels)
- Hard to query & exploit (search, knowledge inference,...)

And advanced queries based on this structure:

```
1 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
2 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 PREFIX taxon: <http://purl.uniprot.org/taxonomy/>
4 PREFIX up: <http://purl.uniprot.org/core/>
5 SELECT ?protein ?organism ?isoform ?sequence
6 WHERE
7 {
8     ?protein a up:Protein .
9     ?protein up:organism ?organism .
10    # Taxon subclasses are materialized, do not use rdfs:subClassOf+
11    ?organism rdfs:subClassOf taxon:83333 .
12    ?protein up:sequence ?isoform .
13    ?isoform rdf:value ?sequence .
```

OVERVIEW



Language Technology

making good progress

mostly solved

Spam detection

Let's go to Agra!



Buy V1AGRA ...



Part-of-speech (POS) tagging

ADJ ADJ NOUN VERB ADV

Colorless green ideas sleep furiously.

Named entity recognition (NER)

PERSON ORG LOC

Einstein met with UN officials in Princeton

Sentiment analysis

Best roast chicken in San Francisco!



The waiter ignored us for 20 minutes.



Coreference resolution

Carter told Mubarak he shouldn't run again.



Word sense disambiguation (WSD)

I need new batteries for my **mouse**.



Parsing



I can see Alcatraz from the window!

Machine translation (MT)

第13届上海国际电影节开幕...



The 13th Shanghai International Film Festival...

1

Information extraction (IE)

You're invited to our dinner party, Friday May 27 at 8:30



still really hard

Question answering (QA)

Q. How effective is ibuprofen in reducing fever in patients with acute febrile illness?

Paraphrase

XYZ acquired ABC yesterday

ABC has been taken over by XYZ

Summarization

The Dow Jones is up

The S&P500 jumped

Housing prices rose



Economy is good

Dialog

Where is Citizen Kane playing in SF?



Castro Theatre at 7:30. Do you want a ticket?





Why else is natural language understanding difficult?

non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #neversaynever & you yourself should never give up either♥

segmentation issues

the New York-New Haven Railroad
the New York-New Haven Railroad

idioms

dark horse
get cold feet
lose face
throw in the towel

neologisms

unfriend
Retweet
bromance

world knowledge

Mary and Sue are sisters.
Mary and Sue are mothers.

tricky entity name

Where is *A Bug's Life* playing?
Let It Be was recorded ...
... a mutation on the *for* gen-

DATA FORMATS & PROCESSING CHAIN

What is textual data?

- A series of letters

```
the_cat_is ...
```

- A series of words

```
the cat is ...
```

- A set of words

in alphabetical order

cat
is
the
...

- N-gram dictionary:

```
BEG_the the_cat cat_is is....
```

Standard processing chain

0. Reading : watch the encoding (utf8, latin1, ...) / file format (txt, pdf, ...)

1. Preprocessing

- encoding (latin, utf8, ...)
- punctuation
- stemming
- lemmatization
- tokenization
- capitals/lower case
- regex
- ...

2. Formatting

- Dictionary
- + reversed Index
- Vectorial format
- Sequence format

3. Learning

- Doc / sentence / word classification
- Semantic
...
- Perceptron, HMM, RNN,
Transformer?

4. Hyper-parameter optimization

Conclusion

- Large number of high value applications
 - (One of the) biggest community (both academic & industrial)
 - A community of its own
 - (very) Specific data
 - Specific pre-processing
 - Specific models

⇒ Specific companies, conferences, ...

until the end of 2010 decade !

NLP = Very high potential, with a significant cost