# Natural Language Processing and Digital Humanities

## Claire Gardent

CNRS/LORIA, Nancy, France
Empire Workshop, Reno
7-8 January 2020

# Three ways of processing/visualising a text collection

# Document-Based Visualisation

- Corpus $\Rightarrow$ Graph, nodes = documents, edge = similarity
- Goal: Explore structure of document collection, Identify key topics, Weed out irrelevant documents

# Location-Based Visualisation

- Corpus $\Rightarrow$ Map, dots = locations + some information (text?)
- Goal: Geograhical representation of events

# Event-Based Visualisation

- Corpus $\Rightarrow$ Graph, nodes = entities, edge = relation
- Goal: Visualisation of the relations between entities (Social Netwok)

# Document-Based Visualisation

## GOAL

- Group together similar documents

## How ?

- Historically relevant information is extracted

    - Persons, Locations, Organisations (**NER**)

- A document is represented by a vector of named entities

    - Similar documents have similar vectors (small cosine)

- A document collection is a graph

    - Nodes are documents, edges encode similarity

- The graph can be searched (integration in an IR system)

# Named Entities

Persons

Locations

Organizations

# Named Entity Recognition

*"We have been hoping for a long time that Comrade **Brezhnev** would visit Cuba. The relations between the CPSU and the communist party of Cuba, relations between our governments and peoples are developing as well as possible."*

# Document Representation

Castro $\begin{bmatrix} 1 \\ 1 \\ 0 \\ \dots \end{bmatrix}$ Cuba $\begin{bmatrix} 1 \\ 1 \\ 1 \\ \dots \end{bmatrix}$ CPSU $\begin{bmatrix} 2 \\ 3 \\ 1 \\ \dots \end{bmatrix}$ declare $\begin{bmatrix} 6 \\ 4 \\ 3 \\ \dots \end{bmatrix}$
Breznev · · · Moscow · · · Health Ministry · · · visit
Batista · · · Volgograd · · · Columbian Army · · · relations
· · · · · · · · · · · ·

# Alias Recognition
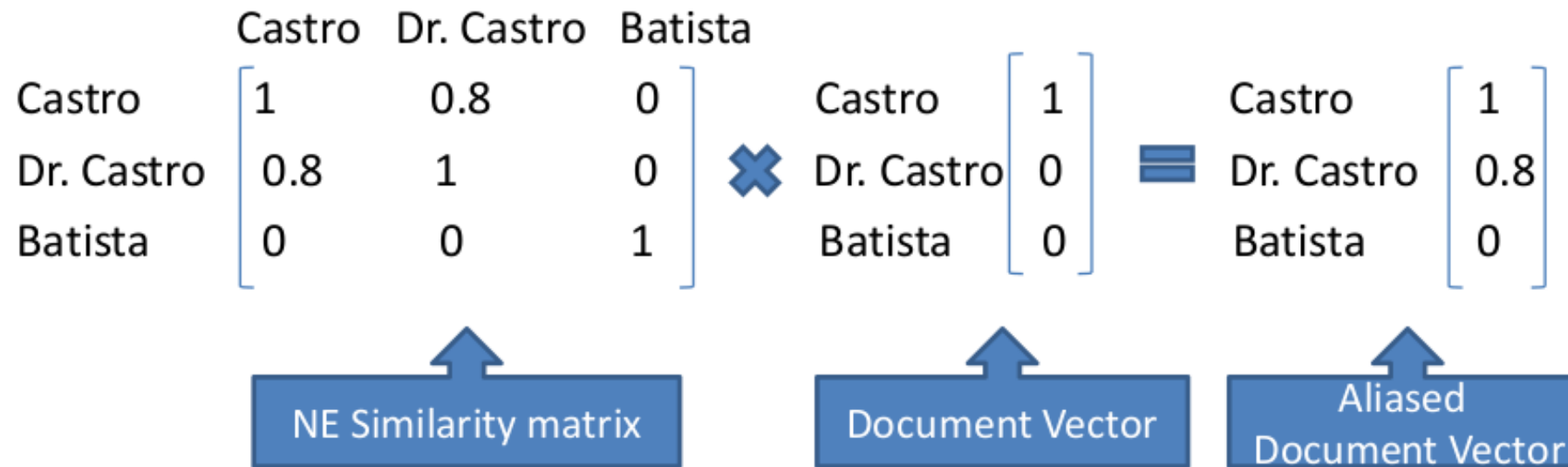


Fidel Castro
Dr. Fidel Castro
Maj. Fidel Castro
Castro

...

Cuban Communist Party
Communist Party of Cuba

...

- String Kernel to measure string similarity
- Helps recognizing linguistic variants of names in a collection of documents

# Alias Boosting



|  | Castro | Dr. Castro | Batista |
|---|---|---|---|
| Castro | 1 | 0.8 | 0 |
| Dr. Castro | 0.8 | 1 | 0 |
| Batista | 0 | 0 | 1 |

**NE Similarity matrix**

| | |
|---|---|
| Castro | 1 |
| Dr. Castro | 0 |
| Batista | 0 |

**Document Vector**

| | |
|---|---|
| Castro | 1 |
| Dr. Castro | 0.8 |
| Batista | 0 |

**Aliased Document Vector**

- Aliases of names that appear in the documents receive an additional weight (more reliable document similarity measure)
- A similar expansion is used for queries that contain NEs (increased querying flexibility)

# Graphical User Interface



- Query for keywords
- Results shown as interactive graph and as table

# Benefits for Historians

## Search + Visualisation

- Search: to express a specific historical question
- Visualisation: to help finding answers and formulating new questions

## Helps ...

- Identify important topics (dense graph regions) and the corresponding documents
- Discover NEs that play an important role in a given topic
- Separating relevant from irrelevant documents

## Clustering ...

- Can be used to further analyse similar documents

# Relevant vs. Irrelevant



- Two groups of interconnected documents.

- The bigger group contains documents concerning health-care in Cuba in the 80s

- The smaller group is about greetings, wishing each other "good health"

# Clustering



(a)

349
864
274
337
772 345
1057
334
213
647

(b)

349
864
772
334
213
274 337 647 345
1057

**Figure 6.** The resulting graphs for the query "Giron Kennedy" before clustering (a), and after clustering (b)

## Clusters

- Documents directly related to the 1961 Bay of Pigs Invasion, e.g. speeches on the anniversary of the event or victory speeches

- Documents mentioning the invasion but not directly related to the event

# NLP Tools used

## Named Entity Recognition

- Stanford Named Entity Recognizer

## Alias Recognition

- String Kernel

## Clustering

- Chinese Whispers Clustering Algorithm

# Named Entity Recognition (NER)

## Which mentions ?

a. *portavoce della* Villa Médicis a Roma
'spokesperson of the Villa Médicis in Rome'

b. *club de pelota vasca de la* ciudad de San Sebastián
'Basque pelota club of San Sebastian city'

c. *dépression karstique dans le*
territoire aride au sud de la *région* d' Aragon
'karstic depression in the arid land south of the Aragon Region'

Gaio and Moncla, IARA'17

# Named Entity Recognition (NER)

## Which types ?

| Named entity | Tag | Description |
|---|---|---|
| Personal names | **p** | first names, surnames, artistic names, (academic) titels, (royal) family names |
| Institutions | **i** | names of institutions, organizations, clubs, companies, names of historical collectives (e. g. religious orders) |
| Geographical names | **g** | names of continents, states, territorial-administrative units, streets and public places, natural monuments including local names |
| Time expressions | **t** | date, days, hours, month, years, centuries, names of epochs, holidays and important days, historic events |
| Artifact names / Objects | **o** | names of documents, artworks, products, books, newspapers, buildings, currency |
| Ambiguous | **a** | used in case the annotator is not sure which of the types above is correct |

Czech Named Entity Corpus 2 . 0, Helena Hubkova, Pavel Kraal, Eva Pettersson LREC'18

# Location-Based Visualisation



## Goal

- Ground information to location

## How ?

- Toponym Resolution: location names are resolved to a geographic reference

# Recognising vs. Resolving Named Entities

*Hamilton is in the North-Island*

## Recognising

- Named Entities

  Hamilton, TYPE: LOC

- Aliases Recognition

  Fidel Castro, Maj. Fidel Castro, Castro, Batista ...

- Corefering mentions

  *Hamilton is in the North Island.*
  *It has a nice beach.*
  *Weta is a big thing in that city*

- Cross-document Coreferring Mentions

## Resolving

- Entity Linking: NE $\Rightarrow$ **Knowledge Base Entity**

  Hamilton (Ontario)
  Hamilton (New Zealand)
  Hamilton (Ohio)
  Hamilton (Bermuda)

- Geo-tagging: Toponym $\Rightarrow$ **Geographic Reference**

  Hamilton, Lat: -37.46, Lon: 175.16

# Cross-Document Coreference Resolution

## Document 1

UAW president Stephen Yokich then met separately for at least an hour with chief executives Robert Eaton of Chrysler Corp., Alex Trotman of Ford Motor Co. and finally with John Smith Jr. of General Motors Corp.
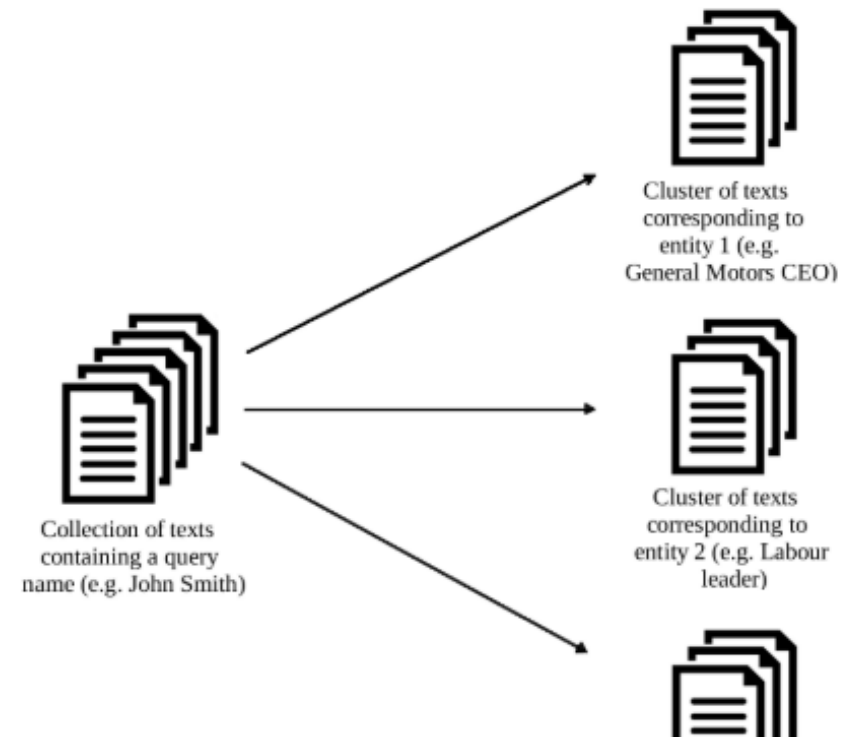
## Document 2

Blair became Labour leader after the sudden death of his successor John Smith in 1994 and since then has steadily purged the party of its high-spend and high-tax policies and its commitment to national ownership of industrial assets.

## Document 3

Two years ago, Powell switched coaches from Randy Huntington to John Smith, who is renowned for his work with sprinters from 100 to 400 meters.

John Smith $\Rightarrow$ 3 document clusters

- the CEO of General Motors
- the Labour Party leader
- an athletics coach



Collection of texts containing a query name (e.g. John Smith)

Cluster of texts corresponding to entity 1 (e.g. General Motors CEO)

Cluster of texts corresponding to entity 2 (e.g. Labour leader)

# Why is NE Resolution hard ?

## Language

### Aliasing

- The same entity can have different names
- Location names often change over time

### Ambiguity

- The same entity name can be used for different entities

## Tools and Resources

### Error Propagation
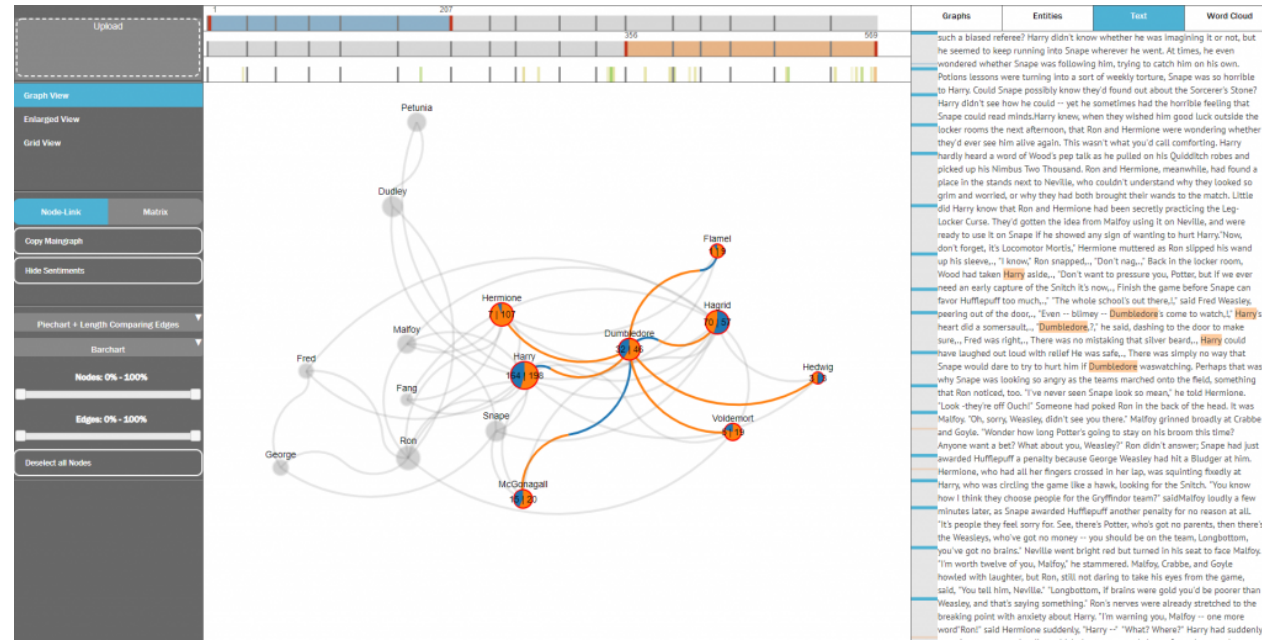
- NER errors propagate to NE disambiguation

### Domain Adpatation

- Most tools trained on contemporary Language

### Incomplete or inadequate KB/Gazeteer

- The entity may not be present in the K/GazeteerB

# Event-Based Visualisation



## Goal

- Identify entity networks (social networks)

## How ?

- NER and Text Segmentation
- Entities that cooccur within a segment are groupe withing a network

# Social Network Visualisation

## Interactive, semi-automatic Expert-Driven Data Analysis

- Automatic tools are imperfect and their results can lead e.g., to highly skewed impressions of the relative importance of characters in a text

- Interaction with domain experts helps mitigate these issues and support the development of better, more adapted NLP tools

## Generic Workflow

The workflow has been applied to

- narrative (modern and medieval) texts
- theoretical philosophical texts, with the goal of establishing relations between philosophical networks
- parliamentary debates, with the goal of connecting po- litical parties to political issues

Blessing, Echelmeyer, John and Reiter. 2017

# Workflow

Text Segmentation

NE linking

Create networks of entities that co-occur within a segment

- E.g., the characters that take part in a great feast

Manual exploration of networks for validation

# Semi-Automatic Processing

## Automatic pre-processing

- tokenization, sentence segmentation, POS tagging

## NER

- Manual annotation of NEs (persons and locations)

- Training a NER (language specific gazetteers and POS tagging for the features)

- NER evaluation using cross validation

- Additional manual annotations (to improve recall)

## NE Resolution

- Manual entity grounding (to a pre-defined list of characters)

# Semi-Automatic Processing (Ct'd)

## Segmentation

- Automatic segmentation (sentence, paragraphs)

## Web-Based Exploration

- Web-based tool for close and distant reading
    - User can view the text passages of the selected entities
    - Network graphs are created with Gephi (Bastian et al., 2009), which provides various layout algorithms, offers statistics and network metrics

# Questions

## Which visualisation method ?

- Explore vs. Display

- One or several

## Named Entities

- Recognition or Disambiguation ?

- Mentions and Types ?

## Domain Adaptation

- Language: Contemporary $\rightarrow$ Historical (Text Normalisation vs. Tool Adpatation)

- Domain: National Newspaper, blogs $\rightarrow$ Regional journal, Encyclopedia

# Thanks