

Entity Tagging & Linking

Afonso Mendes, David Nogueira, Pedro Balage, André Martins



amm@priberam.pt, david.nogueira@priberam.pt
pedro.balage@priberam.pt, afm@priberam.pt

Goals

- To develop statistical models for entity recognition and linking.
- Explore multilingual approaches for the core languages of the project (English, German, Spanish) and for Portuguese.

Named Entity Recognition (NER)

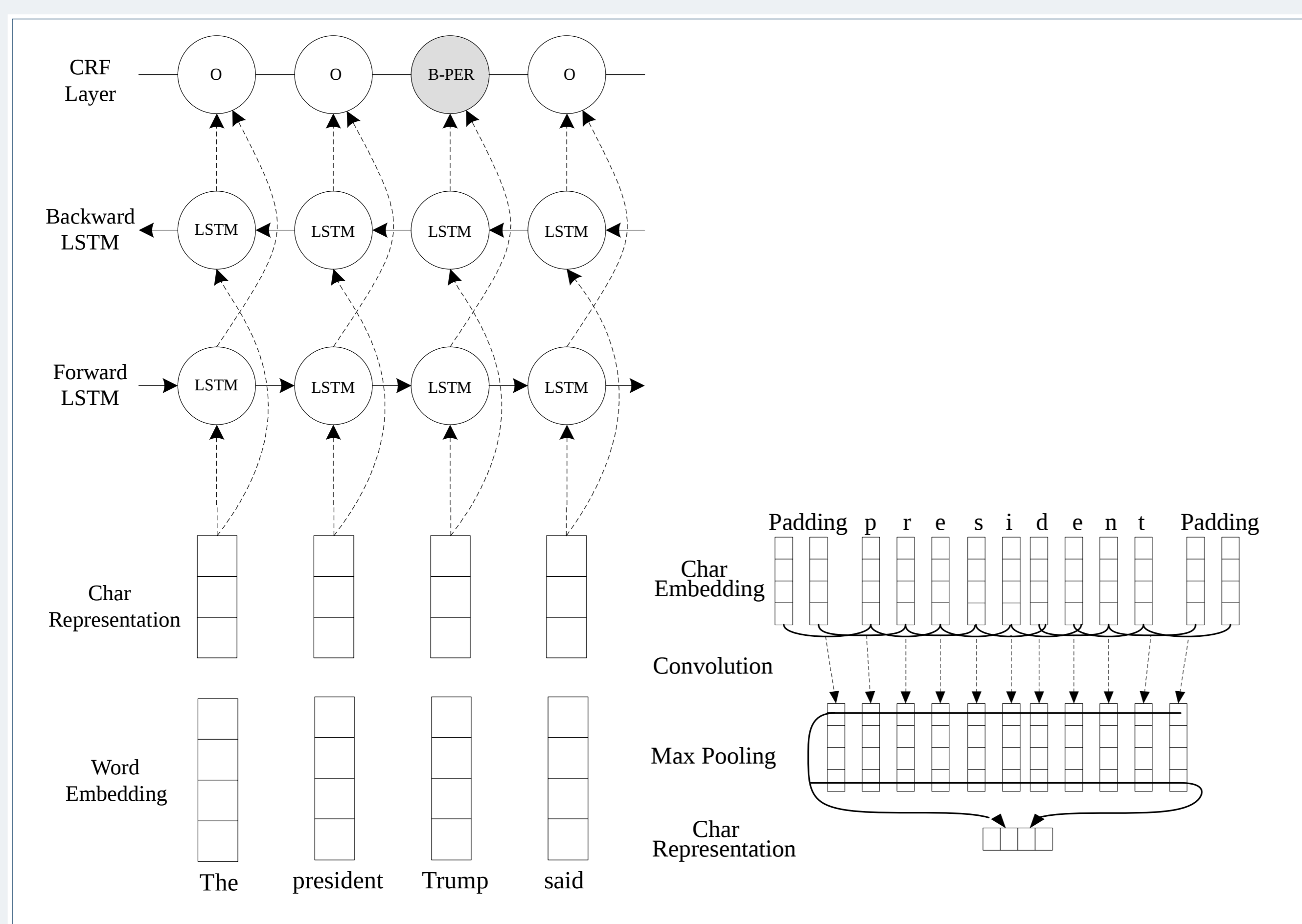
Media Item

Queries > Stories of Query: Venezuela in News > Story: Venezuela protests show no signs of ... > Media Item: Venezue

Venezuela protests show no signs of letting up... Added: 4 days ago (2017-05-04 19:17 UTC)
Changed: 3 days ago (2017-05-05 14:26 UTC)
Source: DW English Articles

Together with her supporters and family, **Lilian Tintori**, the wife of jailed Venezuelan opposition leader **Leopoldo Lopez**, stood in front of the hilltop **Ramo Verde** jail on **Thursday** and demanded to see her husband. **Luis Almagro**, **Organization of American States (OAS)** wrote via **Twitter**: "The **Venezuela** government has refused to confirm the health of political prisoner **Leopoldo Lopez**. Family and lawyers have not seen him in more than a month." "I demand to visit **Leopoldo Lopez** based on the commitments that **Venezuela** has with the **Inter-American System of Human Rights**," Almagro wrote. **Lopez**, a former mayor, was sentenced to nearly 14 years in jail in **2015** following the last major anti-government

- Strong industrial implementation** that uses a sequence labelling model with gazetteers.
- New approaches using **deep neural networks** with bidirectional LSTMs and character- and word-level embeddings can achieve higher scores.



Architecture for the Neural Named-Entity Recognition

Text Analysis Conference (TAC)

- The goal of TAC Knowledge Base Population (KBP) is to develop and evaluate technologies for populating knowledge bases (KBs) from unstructured text.
- TAC is sponsored by the U.S. National Institute of Standards and Technology (NIST) since 2008.
- SUMMA participated at TAC 2017** in the track of Entity Discovery and Linking (EDL).
- From the 22 teams that participated, SUMMA **ranked as first for English EDL** and third for Spanish EDL.

Entity Linking System	NERLC	CEAFmC	Rank
SUMMA	79.4	83.1	1st
Competitor B	79.0	82.5	2nd
Competitor C	78.2	82.4	3rd

Results for English EDL 2017 with named-mentions.

NERLC evaluates the detection, type classification and linking of a mention against a knowledge-base. CEAFmC evaluates the detection and type classification of mention clustering.

Entity Linking

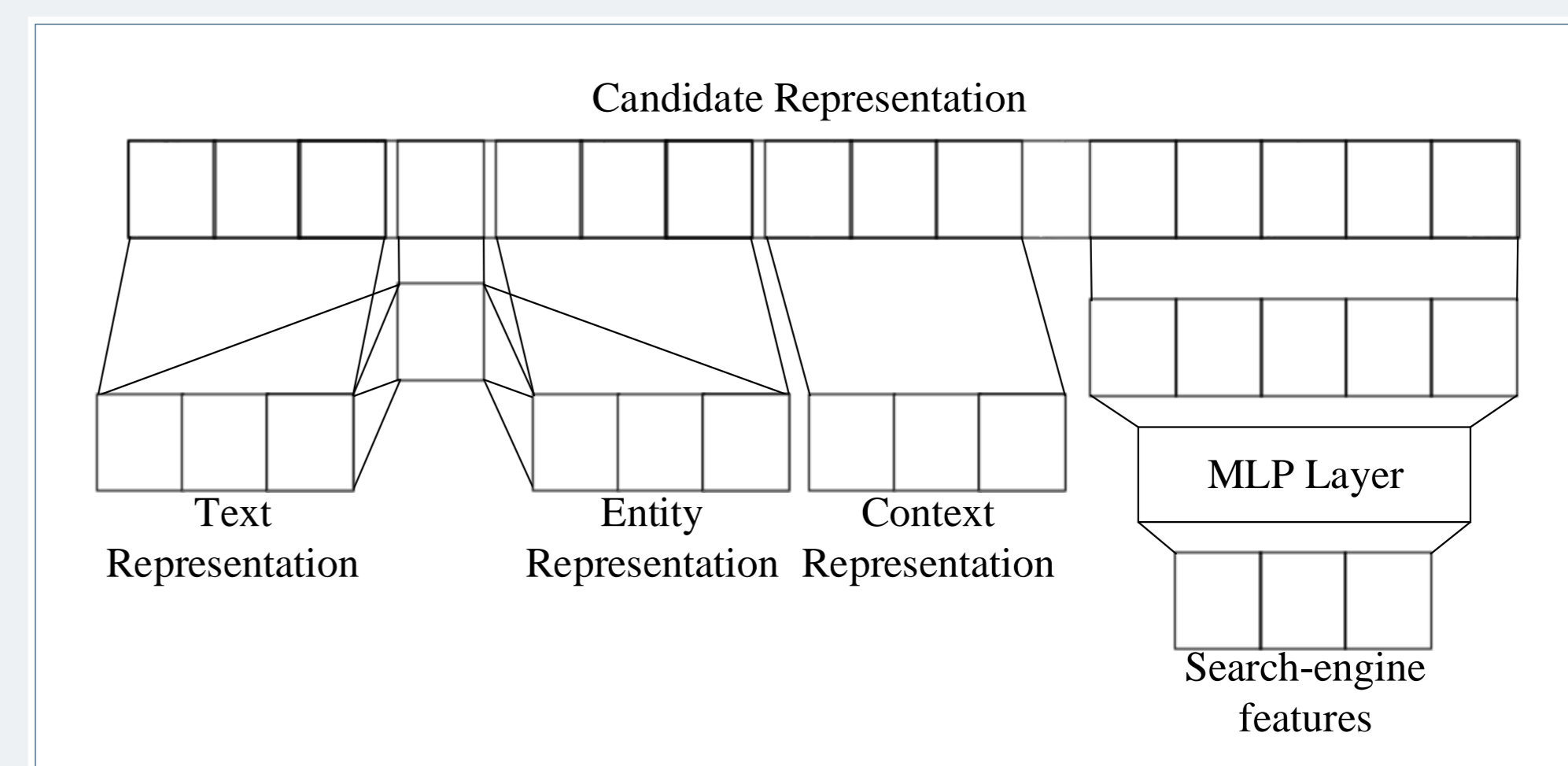
- State-of-the-art industrial language-agnostic implementation**, with document-level and corpora-level coherence steps.
- The model has a small number of parameters, being easily deployed in other languages.
- New model using **distributed representations** boosted the results even further.

The screenshot shows a media item page for Emmanuel Macron. The text contains several entities that have been linked to their respective knowledge bases. For example, "Emmanuel Macron" is linked to a profile page showing his role as President of France, his taking office date (14 May 2017), and his predecessors and successors. Other entities like "Lilian Tintori" and "Leopoldo Lopez" are also linked to their respective profiles.

Algorithm for Entity Linking:

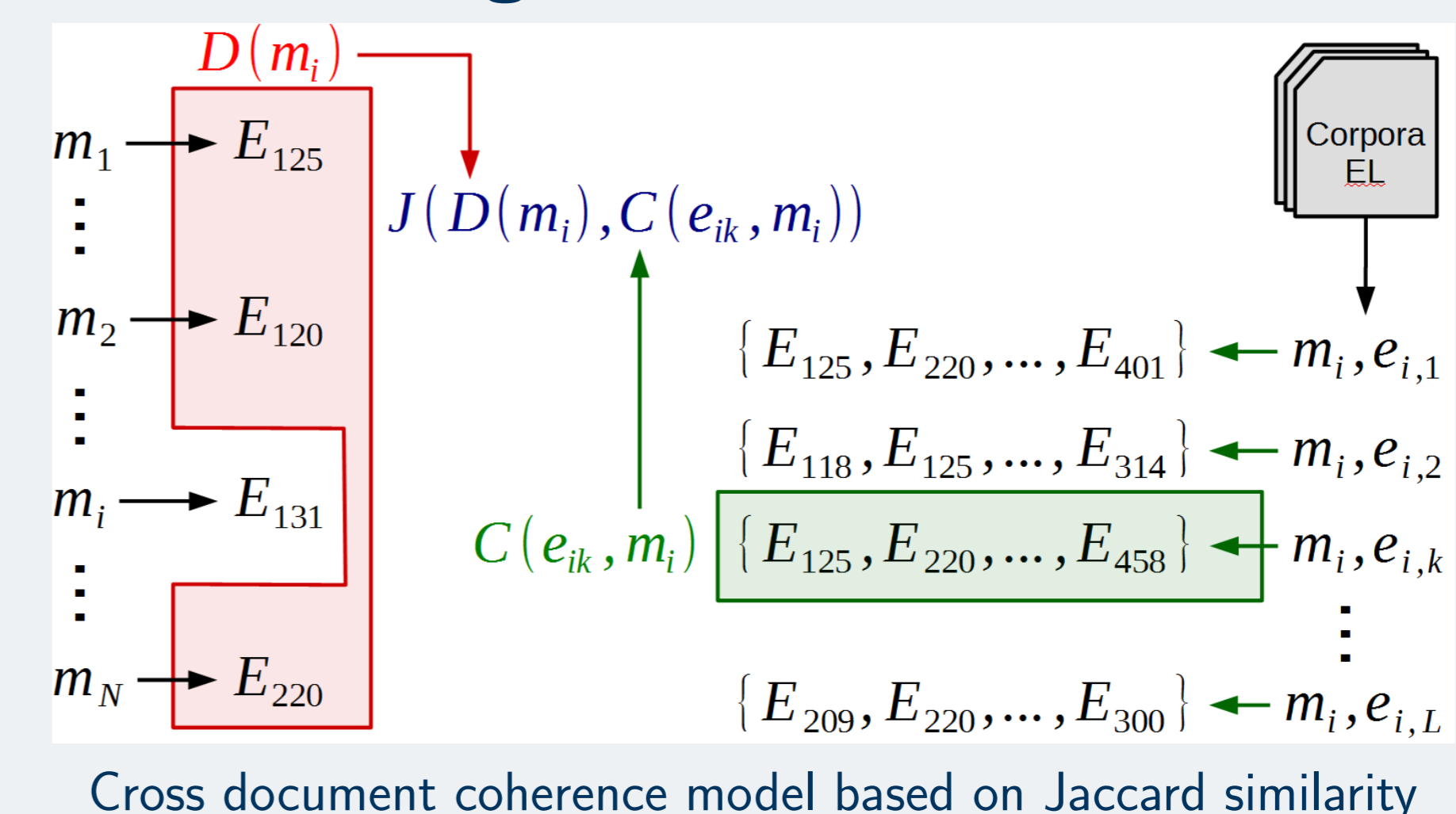
- High-precision **sub-string match mention coreference**
- Candidate generation**
- Features generation: information retrieval engine (KNN algorithm) + prior statistics + mentions candidates co-occurrences
- Distributed representation neural network disambiguation**

A Multi-layer Perceptron (MLP) classifier receives as input the learned representations obtained from the text, each candidate entity representation, context representation and the information retrieval features.



Candidate representation in the neural network model

- Nominal mentions disambiguation**
- Global NIL clustering and cross-document coherence**



Cross document coherence model based on Jaccard similarity

What are the Challenges?

- Named entity recognition performance is directly connected with the amount of data used for training.
- Ensemble methods for named entity recognition.
- Learn suitable cross-lingual entity representations.