

# Entity and Relation Extraction

Department of Information Engineering and Mathematical Sciences  
University of Siena - Siena, Italy

Andrea Zugarini



**SAILab**

December 5th, 2019

Text is a huge source of information!

**Information Extraction** (IE) is one of the most important topics in NLP, and it is about extracting **structured** information from **unstructured** text (documents).

**Goal:** align textual spans to a **Knowledge Base** KB.

**KBs** typically store **factual** data into a **triple**-based ontology.

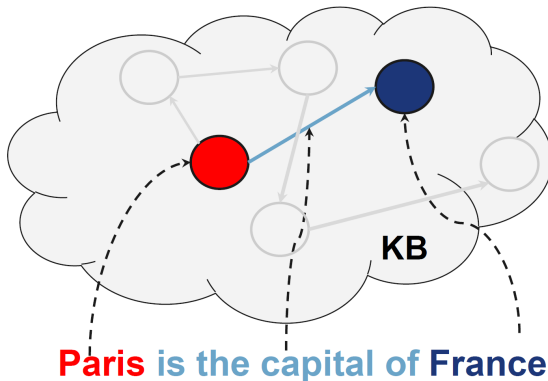
In its simplest version, a KB stores a fact as a triple of two **entities** and a **relation**:  $(e_i, r_k, e_j)$ .

Often, both entities and relations may belong to a **type** of the ontology.

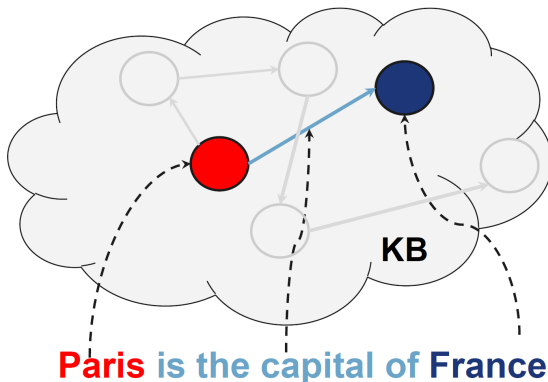
A structured representation of the information it is easier to handle automatically than plain text, allowing efficient storing and retrieval.

A knowledge base can be represented with a directed **graph**, where entities are **nodes** and relations are **edges**.

Information Extraction has to align textual spans of entities and relation to their respective nodes and edges in the KB.



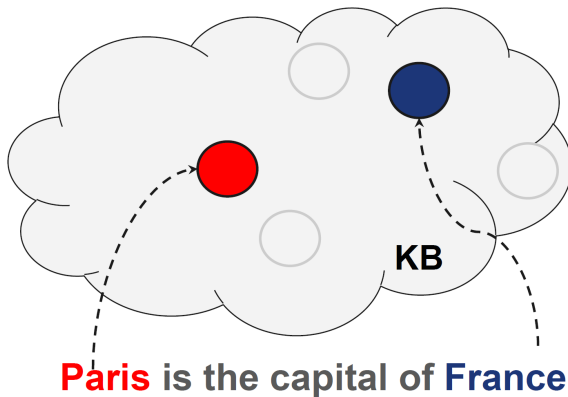
Textual spans are referred as **mentions**.



An entity is a unique instance of something in the real world.

The same entity may be referred by multiple mentions, coreferences included.

Problem related to **Named Entity Recognition**.



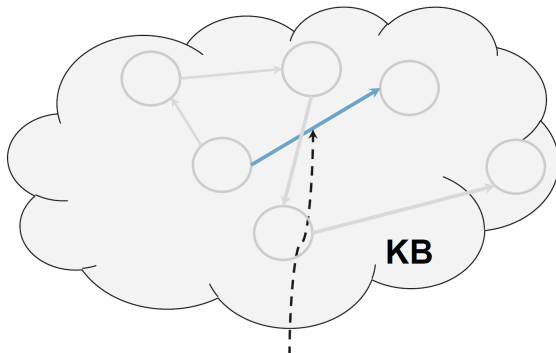
# Knowledge Base

## Triple-based

A relation is a property that connects two (or more) entities.

Challenging problem: there are **many** ways to express the same relations.

Very similar to **Relation Extraction**.



**Paris is the capital of France**

Can plain deep learning techniques be applied for Entity linking and Relation Extraction? Yes, however:

- ▶ End-to-end approaches (encoder-decoder) are very good at mapping text into new text, but they **do not build any concrete understanding model**
- ▶ What does “**understanding text**” mean?

UNDERSTANDING: *mapping text onto a structured (factual) knowledge base* (**Entity Linking**)

TEXT STREAMS: *we have to incrementally build and update the knowledge while we process a text stream!*



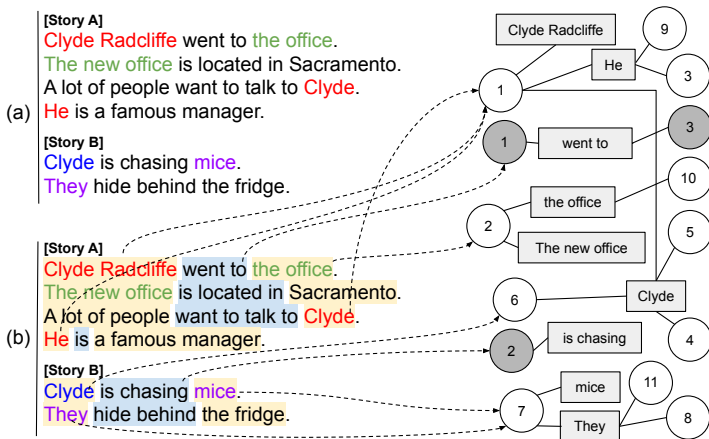
- ▶ We consider a continuous stream of text
- ▶ Groups of sentences organized into small stories about a (not-known-in-advance) set of actors/objects -  $m$  mentions and  $n$  entities/relations
- ▶ The narration is discontinuous whenever a new story begins

## Challenges

KB construction, Online Learning, Entity Discovery, Entity Linking, Multiple Stories

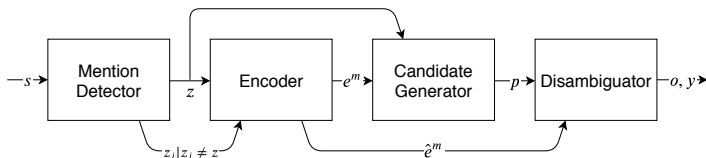
# Problem Setting

## An Example



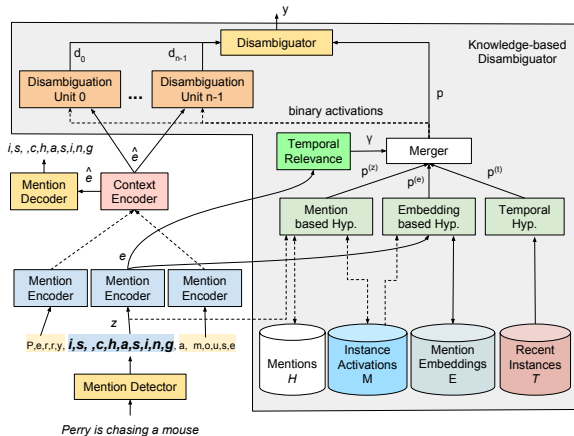
The system is the composition of multiple modules.

- ▶ **Mention Detector:** Segment each sentence in non-overlapping text fragments.
- ▶ **Encoder:** textual mention and its context are encoded into a vectorial representation.
- ▶ **Candidate Generator:** given an input mention  $z$ , the candidate generator implements memory components that are used to generate a list of compatible candidates from the KB.
- ▶ **Disambiguator:** Based on the mention context, the disambiguator is responsible of determining which candidate is the most likely.



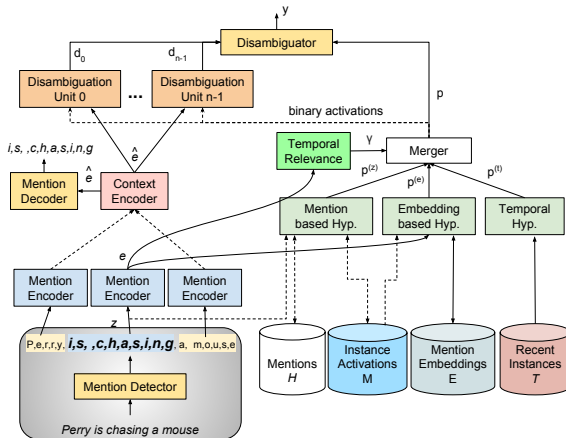
# Architecture

## In detail overview



# Architecture

## Segmentation

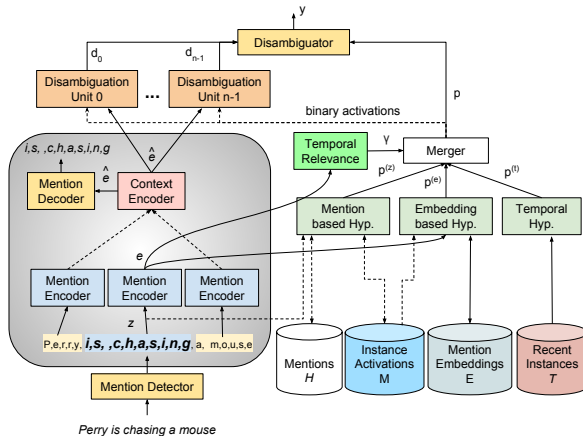


Focus the attention only on relevant text spans (mentions to entities or relations). **How?**

1. Supervised learning using syntax-based generated labels
2. Pre-trained **character-based** model to spot both entities and relations
3. Post processing of predictions to adjust misplaced markers (unclosed elements etc...)

# Architecture

## Mention and Context Encoders



We are given a sequence of *segments*, that are mentions to entities or relations

- **Represent segments:** each segment  $z_i$  is processed as a *sequence of characters* and it is embedded into  $e_i$

$$e_i = enc(z_i) = \left[ \overrightarrow{eRNN}(c_{i,1}, \dots, c_{i,|z_i|}), \overleftarrow{eRNN}(c_{i,1}, \dots, c_{i,|z_i|}) \right]$$

- **Represent contexts:** the context around  $z_i$  is embedded into the representation  $\hat{e}_i$  (that does not include  $z_i$ )

$$\hat{e}_i = enc(z_i | s - z_i) = \left[ \overrightarrow{\hat{e}RNN}(e_1, \dots, e_{i-1}), \overleftarrow{\hat{e}RNN}(e_{i+1}, \dots, e_n) \right]$$

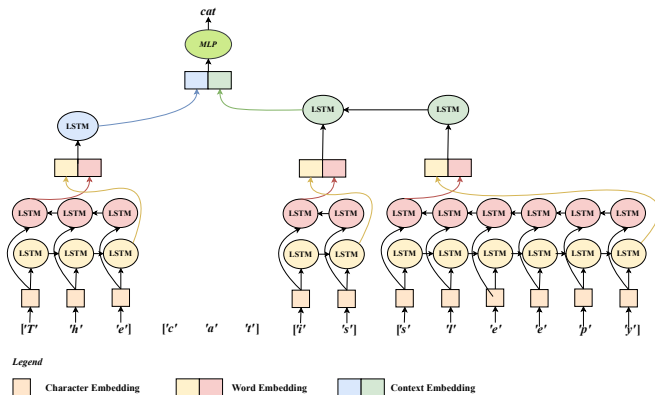
Unsupervised Learning in a encoding-decoding scheme as (CBOW)

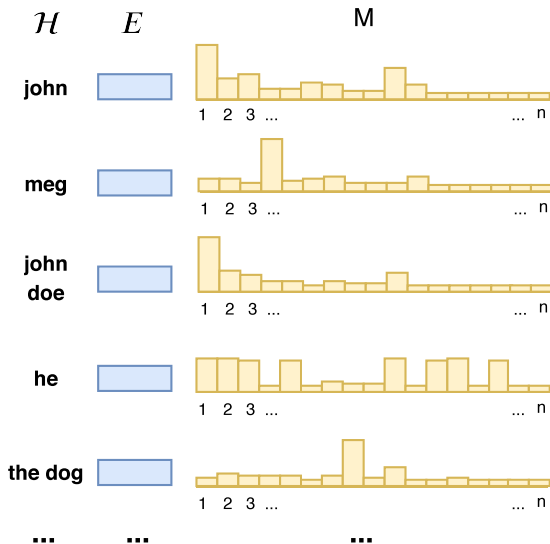


# Encodings

## Learning Example

Sketch of mention and context encoding architecture while processing the sentence “*The cat is sleepy*” with target word *cat*





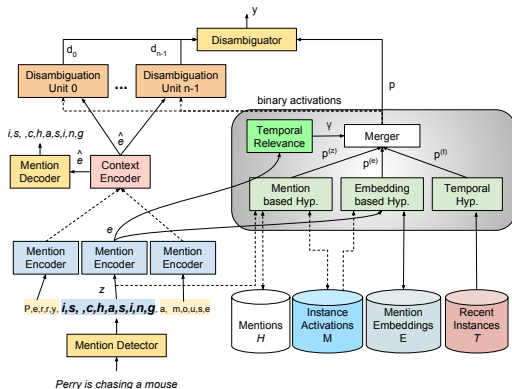
Our **Knowledge Base** is organized in 4 memory components

1.  $\mathcal{H}$  is the set of mentions (raw text)
2.  $E$  is the matrix of the embeddings of each mention
3.  $\mathcal{T}$  buffers the last disambiguated instances, resets at the beginning of a story
4.  $M$  is a matrix where is row is associated to a mention  $z \in \mathcal{H}$  and  $\sigma(M_{\mathcal{H}(z)})$  provides the activation scores of currently known instances

Anytime a new element is encountered the **KB** is updated

# Architecture

## Hypotheses Formulation



Given a mention  $z$  and its embedding  $e$  at time  $t$ , three hypotheses are formulated

► **Mention-based:**  $p^{(z)} = \sigma(M_{\mathcal{H}(z)})$

► **Embedding-based:**

$$p^{(e)} = \left( \left[ \frac{\cos(e, E_i) + 1}{\sum_{j=1}^m \cos(e, E_j) + m} \right]_{i=1}^m \right)' \cdot \sigma(M)$$

► **Time-based:**

$$p^{(t)} = \frac{[u(i, \mathcal{T})]_{i=1}^n}{\max [u(j, \mathcal{T})]_{j=1}^n}$$

ENSEMBLER

$$p = (1 - \gamma) \cdot \left( p^{(z)} + (1 - p^{(z)})p^{(e)} \right) + \gamma \cdot p^{(t)}$$

Given a mention  $z$  and its embedding  $e$  at time  $t$ , three hypotheses are formulated

► **Mention-based:**  $\mathbf{p}^{(z)} = \sigma(M_{\mathcal{H}(z)})$

► **Embedding-based:**

$$\mathbf{p}^{(e)} = \left( \left[ \frac{\cos(e, E_i) + 1}{\sum_{j=1}^m \cos(e, E_j) + m} \right]_{i=1}^m \right)' \cdot \sigma(M)$$

► **Time-based:**

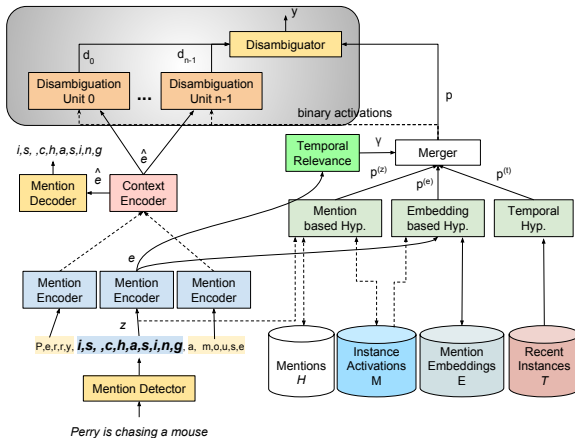
$$\mathbf{p}^{(t)} = \frac{[u(i, \mathcal{T})]_{i=1}^n}{\max [u(j, \mathcal{T})]_{j=1}^n}$$

ENSEMBLER

$$\mathbf{p} = (1 - \gamma) \cdot \left( \mathbf{p}^{(z)} + (1 - \mathbf{p}^{(z)})\mathbf{p}^{(e)} \right) + \gamma \cdot \mathbf{p}^{(t)}$$

# Architecture

## Disambiguator



Hypotheses outputs potential candidates, **disambiguation** resolves ambiguities by selecting the correct one(s). **How?**

It looks at the **context** to find most compatible mention wrt  $n$  disambiguation units

- ▶ **Disambiguation Unit:** given the context  $\hat{e}$ , predicts the activation of the instance (Predictors have a local support around  $\kappa$  centres)

$$d_i(\hat{e}) = \frac{1}{2} + \frac{1}{2} \max_{j=1}^{\kappa} \cos(\hat{e}, \hat{w}_{ij})$$

Final output of the system is  $\mathbf{o}$

$$\mathbf{o} = \delta(\mathbf{p} > \tau_r) \cdot (\eta \cdot \mathbf{p} + (1 - \eta) \cdot \mathbf{d})$$



**Online Learning** process accordingly to either supervision or **self learning**

When no supervision is provided, we distinguish among three cases:

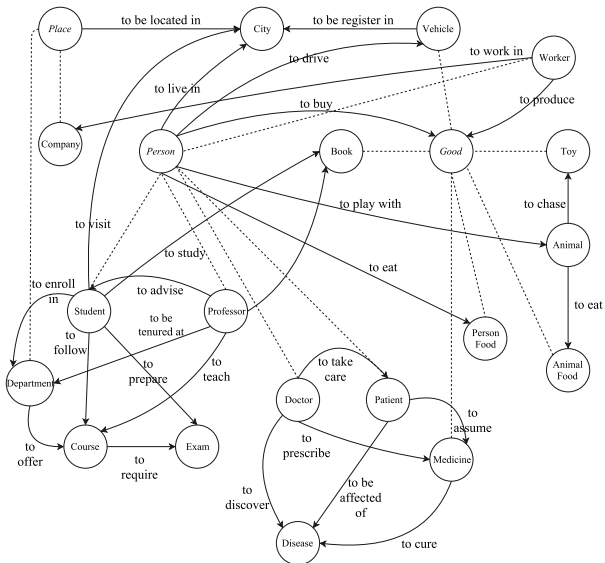
- i.*  $\max \mathbf{o} \geq \tau_a$  : RECOGNIZED SOME INSTANCES
- ii.*  $\max \mathbf{p} > \tau_r \wedge \max \mathbf{o} < \tau_a$  : UNCERTAINTY
- iii.*  $\max \mathbf{p} \leq \tau_r$  : UNKNOWN INSTANCE

**Learnable parameters:**  $M$ , disambiguation units  $d$  and  $\gamma$

- ▶ Collection of 10k sentences organized in 564 stories
- ▶ A story is a list of not repeated facts mostly focussed on a certain entity, also called **main entity**
- ▶ 130 **entity** and 27 **relation** instances, belonging to a pre-designed ontology
- ▶ Overall there are 2176 single word tokens, 1528 and 288 mentions to entities and relations, about 6830 ambiguous mentions.

# Experimental Environment

## Dataset Ontology



- ▶ Wikipedia pages are loosely aligned with Freebase triples
- ▶ Composed of a collection of summaries, each being a description of a certain entity. Each summary is considered as a story.
- ▶ About  $560k$  **entities** extracted from  $10k$  pages, we took a sub-portion of 1112 pages.
- ▶ We marked text between two entities as relation.

Each story split into two parts: a supervised and an unsupervised one.

Accuracy on each prediction is measured at the same time when the prediction is made.

Two results reported:

- ▶ All the unsupervised sentences of a story (**ALL**)
- ▶ Only the last sentence of the story (**LAST**)

## RULE-BASED

An informed model that buffer statistics on the supervisions received up to time  $t$ .

**Already seen mention:** predicts the most common supervision

**Never seen mention:** responds with the most frequent supervision of the story

## DEEP-RNN

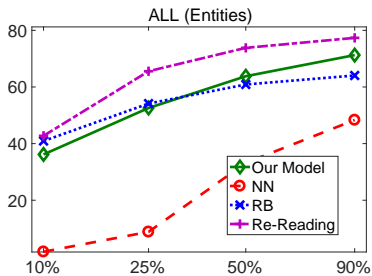
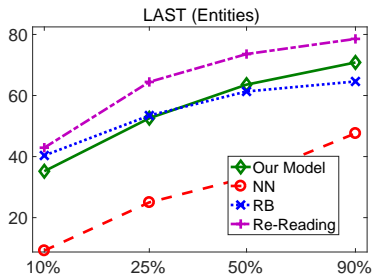
A simple neural mention classifier

- ▶ Input  $[e, \hat{e}]$
- ▶ 1 hidden layer of size 600
- ▶ Softmax activation in the output layer

**NB:** Both models always predicts on ground truth mentions!

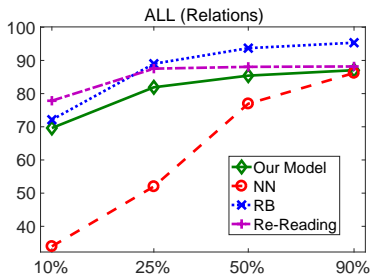
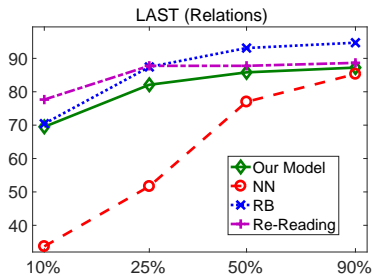
# Results

## Entities



# Results

## Relations





	Model	10%	25%	50%	90%
<b>All</b>	RB	16.84	40.44	48.28	49.55
	Deep-RNN	0.6	3.01	12.34	21.78
	Our Model	39.25	<b>54.57</b>	<b>69.64</b>	<b>75.45</b>
	Re-Reading	<b>44.75</b>	<b>54.66</b>	66.88	70.55
<b>Last</b>	RB	17.28	40.87	48.04	49.37
	Deep-RNN	0.6	3.25	12.11	21.37
	Our Model	37.44	52.93	<b>67.45</b>	<b>75.37</b>
	Re-Reading	<b>43.41</b>	<b>53.38</b>	65.13	70.39

We presented an end-to-end model for entity/relation mentions discovery and disambiguation in text streams by constantly updating an interpretable KB

## Next Steps

- ▶ Entity and Relation Types introduction
- ▶ Higher-level reasoning
- ▶ Dynamic KB re-organization (pruning, merging etc..)

Thank You !!!