# END-TO-END NEURAL ENTITY LINKING

*Nikolaos Kolitsas, Octavian-Eugen Ganea, Thomas Hofmann*
CoNLL 2018

Mariana Leite
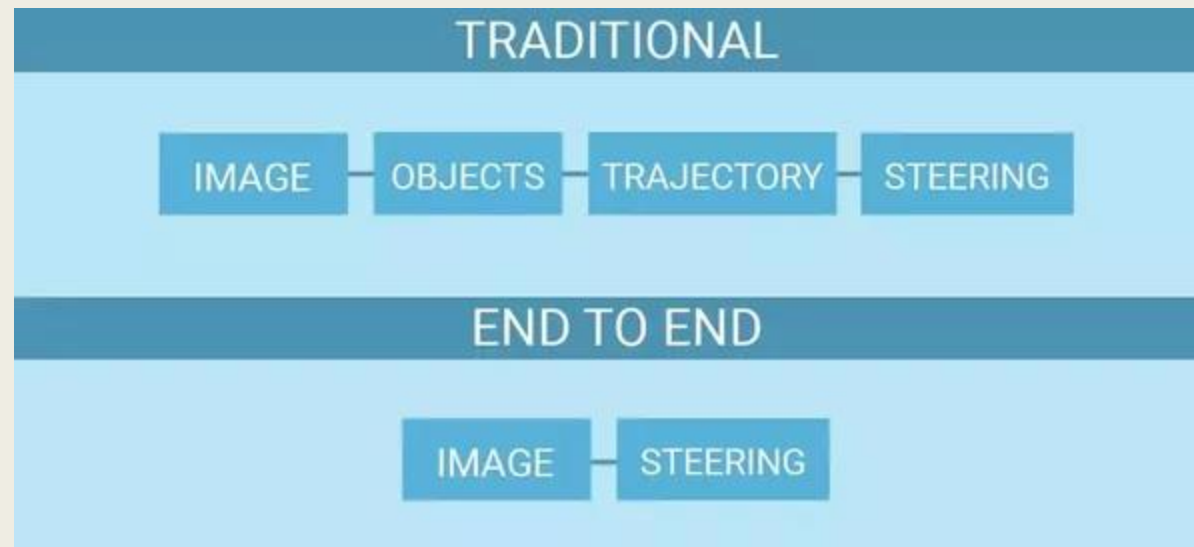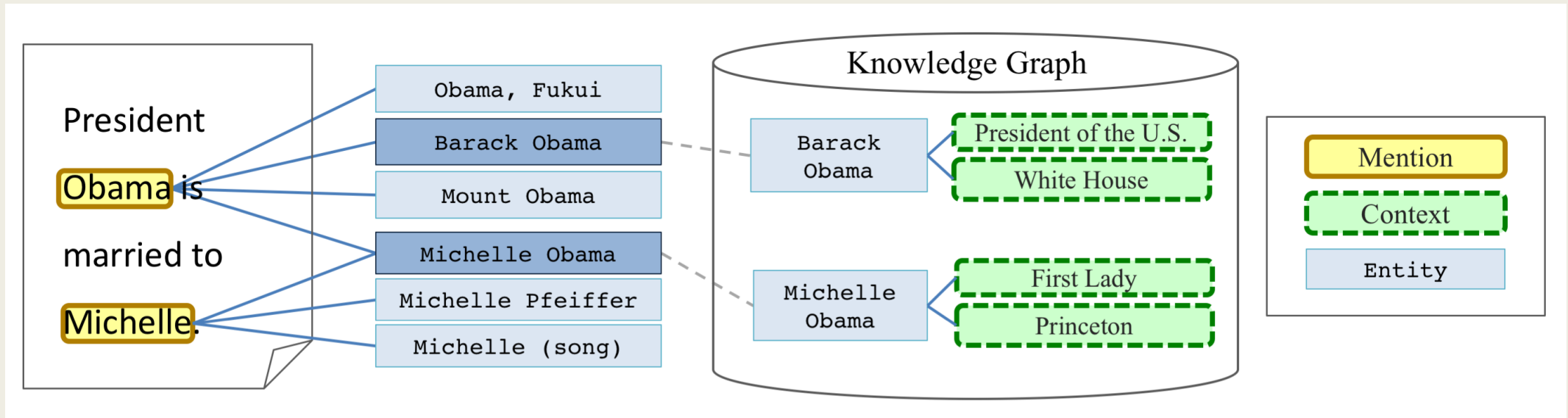me.leite@campus.fct.unl.pt
NOVA Search Reading Group

# End-to-End?

■ End-to-end (E2E) learning refers to training a possibly complex learning system represented by a single model (specifically a Deep Neural Network) that represents the complete target system, bypassing the intermediate layers usually present in traditional pipeline designs [1].



[1] https://towardsdatascience.com/e2e-the-every-purpose-ml-method-5d4f20dafee4

# The Goal of Automatic Text Understanding...

- Models are expected to accurately **extract** ambiguous mentions of entities from a textual document and **link** them to a knowledge base.



- EL systems typically perform 2 tasks: Mention Detection (MD) & Entity Disambiguation (ED)

# Approaches

■ The usual: solve the two tasks independently.

! Important depency between the two steps is ignored;

! Errors in first step will propagate to second step without possibility of recovery.

■ This paper's approach:

– End-to-End Entity Linking (like humans do!).

– Emphasizes the **importance of the mutual dependency between MD & ED.**

# MD may benefit from ED (and viceversa)

1) MD may split a larger span into two mentions of less informative entities:

B. Obama's wife gave a speech [...]

Federer's coach [...]

2) MD may split a larger span into two mentions of incorrect entities:

Obama Castle was built in 1601 in Japan.

The Kennel Club is UK's official kennel club.

A bird dog is a type of gun dog or hunting dog.

Romeo and Juliet by Shakespeare [...]

Natural killer cells are a type of lymphocyte

Mary and Max, the 2009 movie [...]

# MD may benefit from ED (and viceversa)

3) MD may choose a shorter span,
   referring to an incorrect entity:
The Apple is played again in cinemas.
The New York Times is a popular newspaper.

4) MD may choose a longer span,
   referring to an incorrect entity:
Babies Romeo and Juliet were born hours apart.

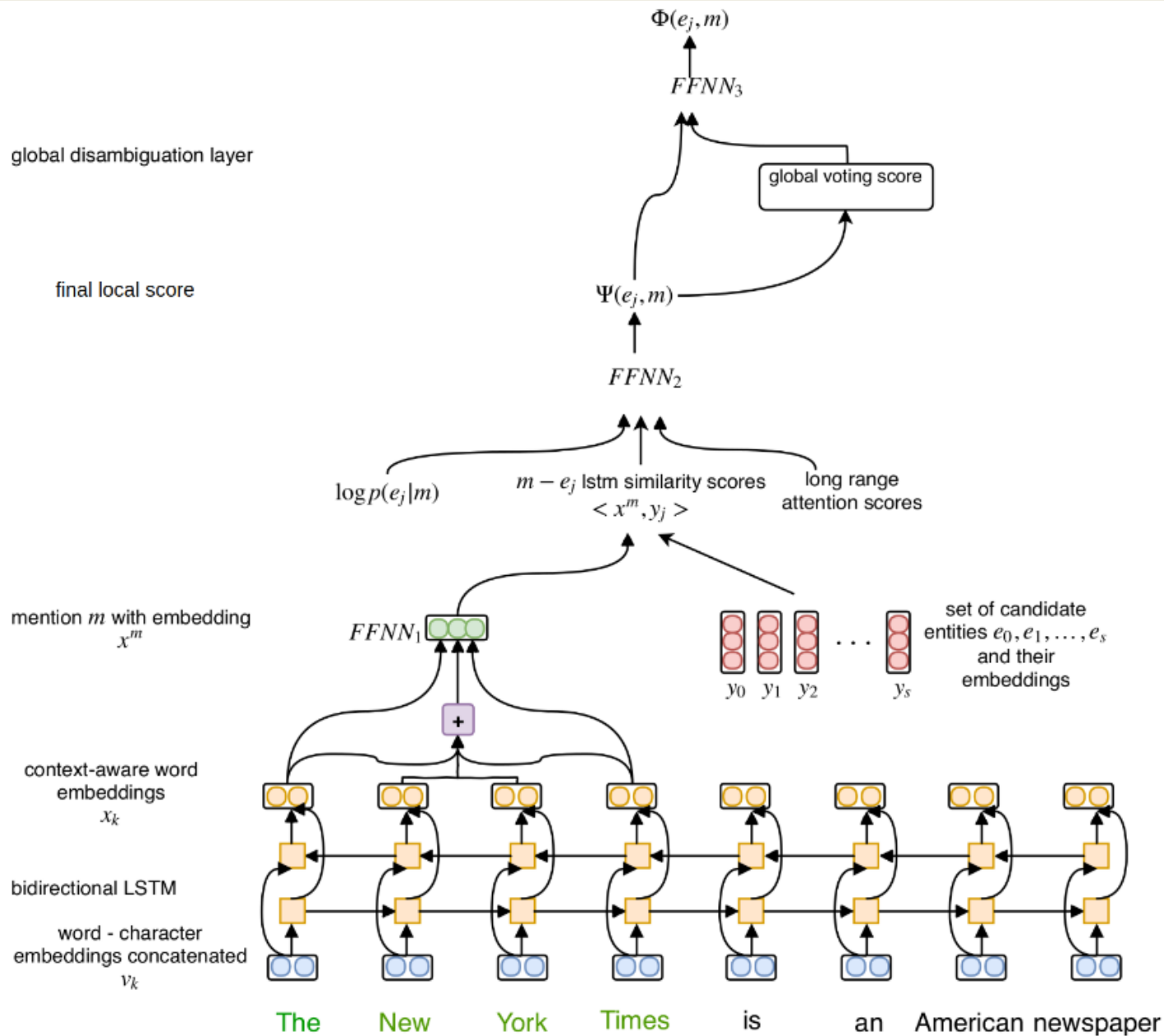# Neural Joint Mention Detection & Entity Disambiguation

- Formally for EL:
  - **Input:** $D = \{w_1, \ldots, w_n\}$, where $w_k \in W$
  - **Output:** list of mention - entity pairs $\{(m_i, e_i)\}_{i \epsilon \overline{1,T}}$
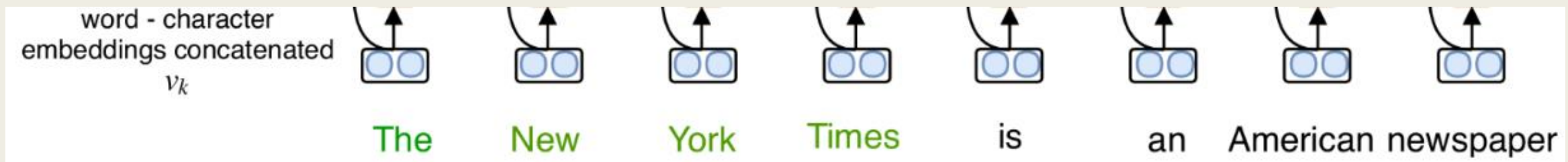
- Formally for ED:
  - **Input:** $D = \{w_1, \ldots, w_n\}$, where $w_k \in W$ + $\{m_i\}_{i \epsilon \overline{1,T}}$
  - **Output:** $\{e_i\}_{i=\overline{1,T}} \in \varepsilon^T$
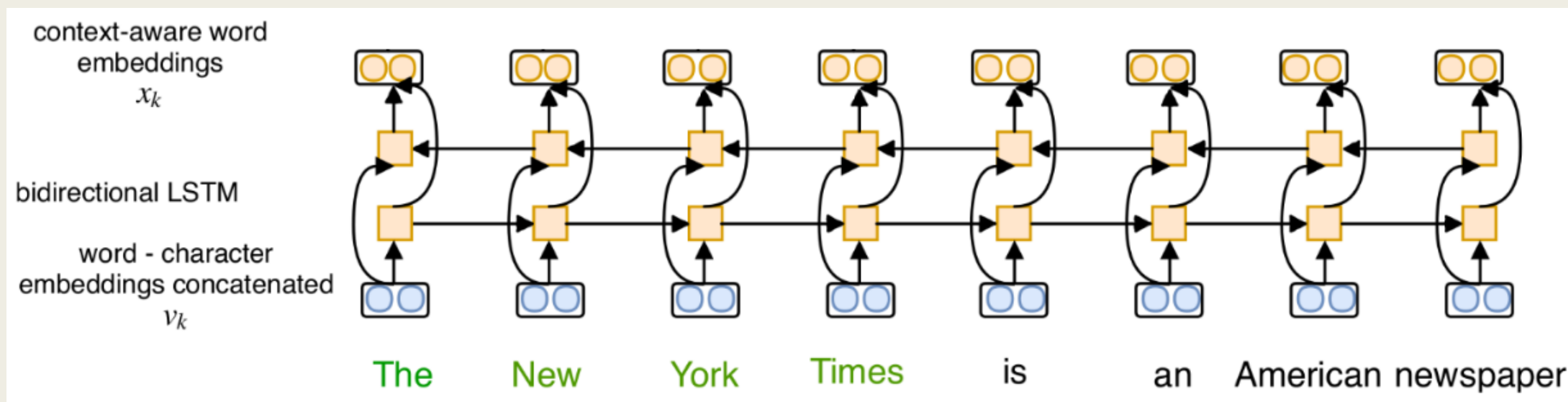
The Model Architecture

8

# Step 1 - Word and Char Embeddings

- $\{z_1, \dots, z_L\}$ - character vectors of word w

- $h_t^f = FWD - LSTM\left(h_{t-1}^f, z_t\right)$

- $h_t^b = BKWD - LSTM\left(h_{t+1}^b, z_t\right)$

- Character embedding of w is $[\boldsymbol{h_L^f}; \boldsymbol{h_1^b}]$

- Character embedding is concatenated with the pre-trained word embedding [2]
  - *Forms the context-independent word-character embedding of w*



[2] Mikolov et al. 2013. Distributed representations of words and phrases and their compositionality.
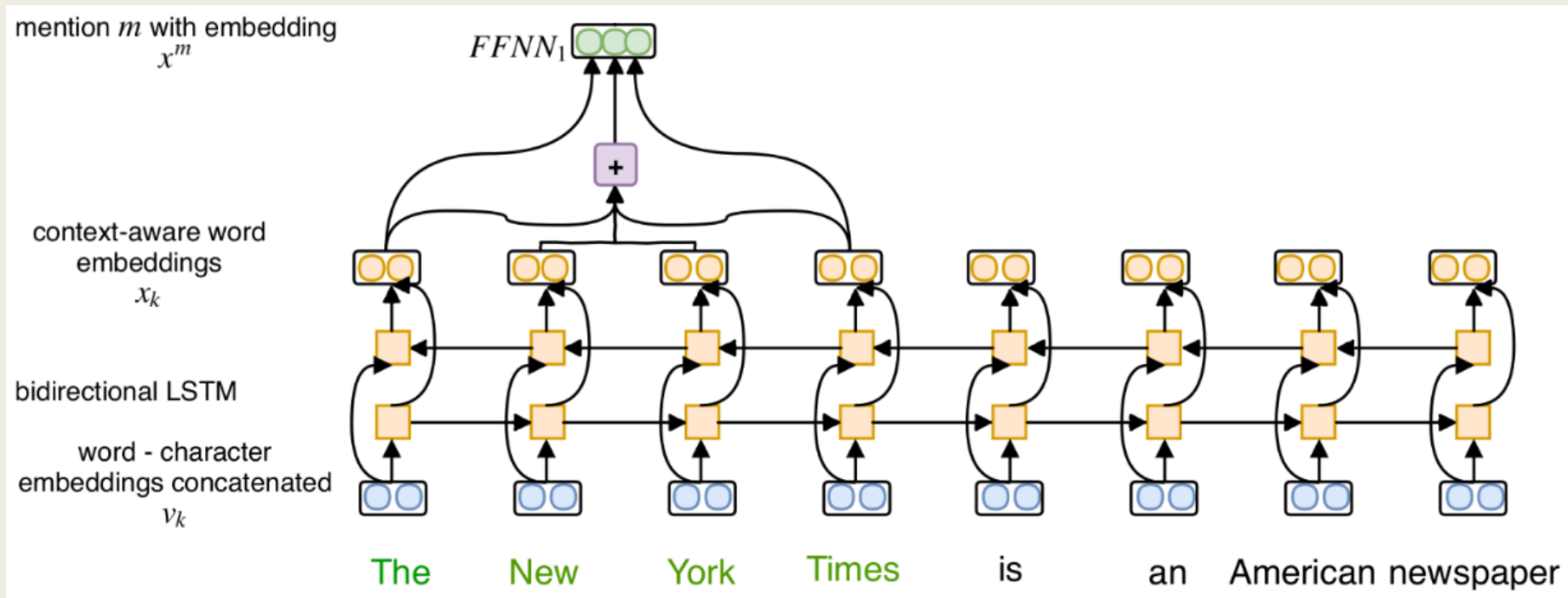
# Step 2 – Mention Representation

- To make word embeddings aware of **local context**: bi-LSTM layer!

- Hidden states of the bi-LSTM (corresponding to each word) are concatenated into context-aware word embeddings.
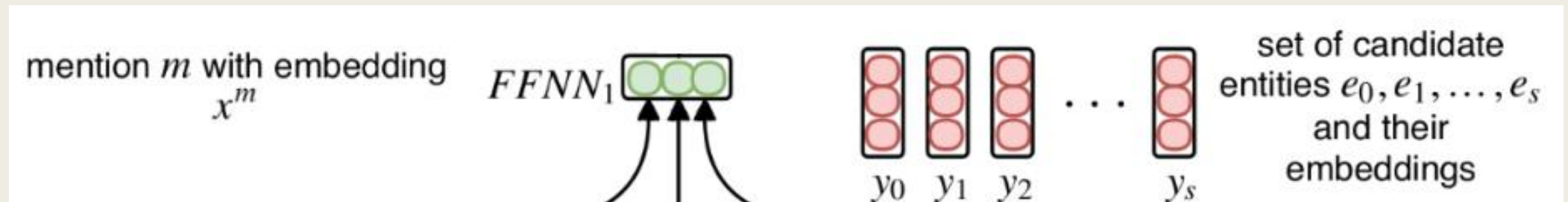
# Step 2 – Mention Representation

■ For each possible mention $m = w_q, \dots, w_r$, concatenate first, last and "soft head" words

  – *fixed size representation* $g^m = \left[ x_q; x_r; \hat{x}^m \right]$

■ To learn non-linear interactions between the component word vectors:

  – *project $g^m$ to a final mention representation $x^m = FFNN_1(g^m)$*
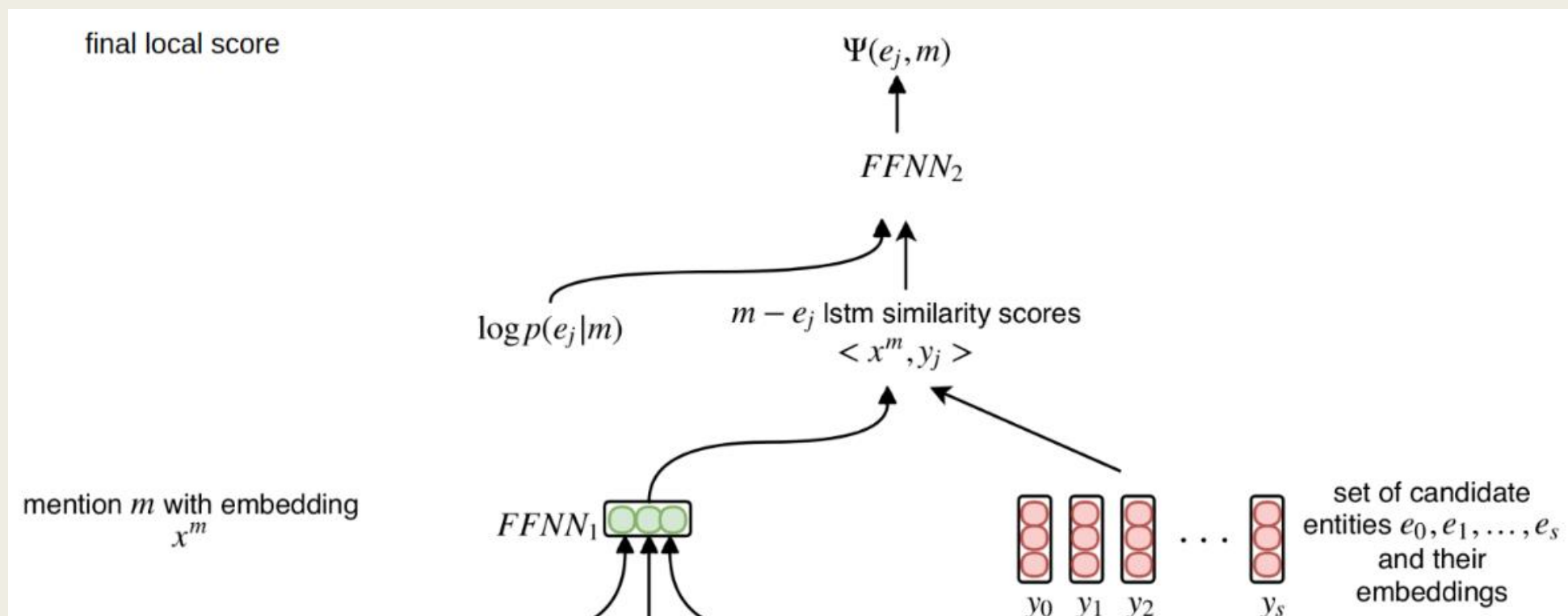
# Step 3 - Entity Embeddings
# Step 4 - Candidate Selection

- Use the fixed continuous **pre-trained entity representations** of Ganea and Hofmann [3].

- For each span $m$ select the top $s$ entity candidates.
  - *Candidate set $C(m)$*

- Top is based on empirical probabilistic entity – map $p(\text{e}|m)$ [3].
  - *From Wikipedia, Crosswikis and YAGO dictionaries*

- $C(m)$ is used at training and test time.



mention $m$ with embedding $x^m$  —  $FFNN_1$  —  $y_0$ $y_1$ $y_2$ $\cdots$ $y_s$  —  set of candidate entities $e_0, e_1, \ldots, e_s$ and their embeddings

[3] Octavian-Eugen Ganea and Thomas Hofmann. 2017. Deep joint entity disambiguation with local neural attention.
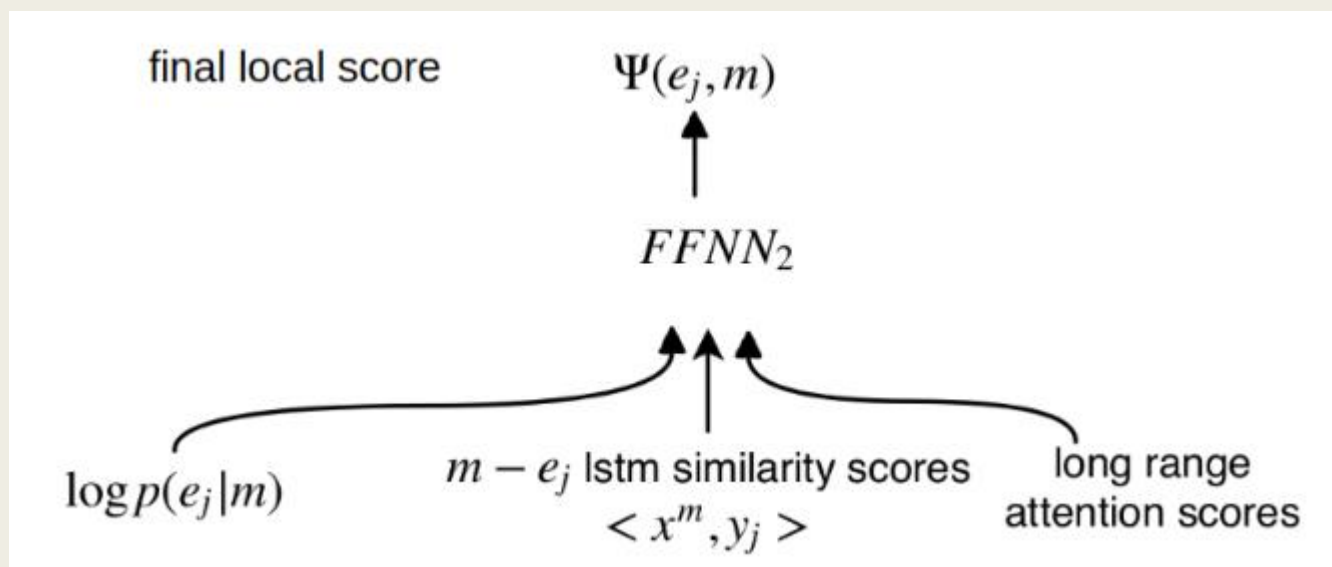
# Step 5 – Final Local Score

■ Compute **similarity score** using embedding dot-product for each candidate e

– *When $|C(m)| \geq 1$*

■ Combine it with the log-prior probability using a FFNN

# Step 5 (Extra!) – Long Range Context Attention

- Improves model in some cases by explicitly capturing long context dependencies.
- Using an attention model [3] collect one context embedding per mention
  - *based on informative context words*
  - *related to at least of the candidate entities*

# Step 6 - Training

■ Assuming corpus with documents and gold entity - mention pairs $G = \{(m_i, e_i^*)\}_{i=\overline{1,K}}$

■ Collect set M of all (potentially overlapping) token spans m for which $|C(m)| \geq 1$

■ Train the parameters using the minimization procedure:

$$\theta^* = \arg\min_{\theta} \sum_{m \in M} \sum_{e \in C(m)} V(\Psi_\theta(e, m))$$

■ V enforces the scores of gold pairs to be linearly separable from scores of negative pairs

# Extra Step – Global Disambiguation

- Currently, the model performs disambiguation of each candidate span independently.

- Add extra layer that promotes coherence among linked and disambiguated entities inside the same document.

  1. Define the set of mention-entity pairs that can participate in the global disambiguation voting (that have a **high score**)

  2. Calculate **final "global" score** $G(e_j, m) = \cos\left(y_{e_j}, y_G^m\right)$

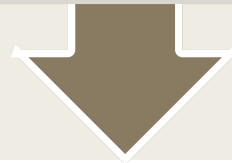  3. Combine this with the local score:

# Coreference Resolution Heuristic

It is important to be able to solve simple coreference resolution cases

*("Alan" referring to "Alan Shearer")*

these cases are difficult to handle by the candidate selection strategy

Adopt simple heuristic [3]:

Observed between 0.5% and 1% improvement on all datasets

The Model Architecture

# Experiments

- Used Wikipedia 2014 as KB.

- Conducted experiments on the most important public EL datasets using the **Gerbil** platform.

- For training, used the biggest publicly available EL dataset: AIDA/CoNLL

    *Training Set*

    – *18,448 linked mentions*

    – *946 documents*

    Validation Set

    – *4,791 mentions*

    – *216 documents*

    Test Set

    – *4,485 mentions*

    – *231 documents*

# Results

- 4 different models used
  - **Base model**: only uses the mention local score and the log-prior.
  - **Base model + att**: the Base Model plus Long Range Context Attention.
  - **Base model + att + global**: the Global Model
  - **ED base model + att + global StanfordNER:** ED Global model that runs on top of the detected mentions of the Stanford NER system [4].

[4] Finkel et al. 2005. Incorporating non-local information into information extraction systems by gibbs sampling.

# Results

| F1@MA F1@MI | AIDA A | AIDA B | MSNBC | OKE-2015 | OKE-2016 | N3-Reuters-128 | N3-RSS-500 | Derczynski | KORE50 |
|---|---|---|---|---|---|---|---|---|---|
| FREME | 23.6 | 23.8 | 15.8 | 26.1 | 22.7 | 26.8 | 32.5 | 31.4 | 12.3 |
|  | 37.6 | 36.3 | 19.9 | 31.6 | 28.5 | 30.9 | 27.8 | 18.9 | 14.5 |
| FOX | 54.7 | 58.1 | 11.2 | 53.9 | 49.5 | 52.4 | 35.1 | 42.0 | 28.3 |
|  | 58.0 | 57.0 | 8.3 | 56.8 | 50.5 | 53.3 | 33.8 | 38.0 | 30.8 |
| Babelfy | 41.2 | 42.4 | 36.6 | 39.3 | 37.8 | 19.6 | 32.1 | 28.9 | 52.5 |
|  | 47.2 | 48.5 | 39.7 | 41.9 | 37.7 | 23.0 | 29.1 | 29.8 | 55.9 |
| Entityclassifier.eu | 43.0 | 42.9 | 41.4 | 29.2 | 33.8 | 24.7 | 23.1 | 16.3 | 25.2 |
|  | 44.7 | 45.0 | 42.2 | 29.5 | 32.5 | 27.9 | 22.7 | 16.9 | 28.0 |
| Kea | 36.8 | 39.0 | 30.6 | 44.6 | 46.3 | 17.5 | 22.7 | 31.3 | 41.0 |
|  | 40.4 | 42.3 | 30.9 | 46.2 | 46.4 | 18.1 | 20.5 | 26.5 | 46.8 |
| DBpedia Spotlight | 49.9 | 52.0 | 42.4 | 42.0 | 41.4 | 21.5 | 26.7 | 33.7 | 29.4 |
|  | 55.2 | 57.8 | 40.6 | 44.4 | 43.1 | 24.8 | 27.2 | 32.2 | 34.9 |
| AIDA | 68.8 | 71.9 | 62.7 | 58.7 | 0.0 | 42.6 | 42.6 | 40.6 | 49.6 |
|  | 72.4 | 72.8 | 65.1 | 63.1 | 0.0 | 46.4 | 42.4 | 32.6 | 55.4 |
| WAT | 69.2 | 70.8 | 62.6 | 53.2 | 51.8 | 45.0 | 45.3 | 44.4 | 37.3 |
|  | 72.8 | 73.0 | 64.5 | 56.4 | 53.9 | 49.2 | 42.3 | 38.0 | 49.6 |
| **Best baseline** | **69.2** | **71.9** | **62.7** | **58.7** | **51.8** | **52.4** | **45.3** | **44.4** | **52.5** |
|  | **72.8** | **73.0** | **65.1** | **63.1** | **53.9** | **53.3** | **42.4** | **38.0** | **55.9** |
| base model | 86.6 | 81.1 | 64.5 | 54.3 | 43.6 | 47.7 | 44.2 | 43.5 | 34.9 |
|  | 89.1 | 80.5 | 65.7 | 58.2 | 46.0 | 49.0 | 38.8 | 38.1 | 42.0 |
| base model + att | 86.5 | 81.9 | 69.4 | 56.6 | 49.2 | 48.3 | 46.0 | 47.9 | 36.0 |
|  | 88.9 | 82.3 | 69.5 | 60.7 | 51.6 | 51.1 | 40.5 | 42.3 | 42.2 |
| base model + att + global | 86.6 | 82.6 | 73.0 | 56.6 | 47.8 | 45.4 | 43.8 | 43.2 | 26.2 |
|  | 89.4 | 82.4 | 72.4 | 61.9 | 52.7 | 50.3 | 38.2 | 34.1 | 35.2 |
| ED base model + att + global using Stanford NER mentions | 75.7 | 73.3 | 71.1 | 62.9 | 57.1 | 54.2 | 45.9 | 48.8 | 40.3 |
|  | 80.3 | 74.6 | 71.0 | 66.9 | 58.4 | 54.6 | 42.2 | 42.3 | 46.0 |

# Results

Highlighted the best and second best models for the EL task.

Metrics computed in the **strong matching** setting:

➢ Requires the exact prediction of the gold mention boundaries and their entity annotations.

# Main Goal Achieved!

- The joint EL offers the best model!
    - If enough training data is available with the same characteristics as the test data.
    - True not only when training on AIDA, but also for other types of datasets (queries, tweets)

- When testing data has different statistics than the training data:
    - Method works best with a state-of-the-art NER system.

# Conclusion

Presented the first neural end-to-end entity linking model

Showed the benefit of jointly optimizing entity recognition and linking

Proved that engineered features can be almost completely replaced by modern neural networks

Q&A