



# **Accelerating Biomedical Entity Linking Models**

Javier Sanz-Cruzado



WORLD CHANGING GLASGOW

Information Access and Retrieval in the Al Age
Festival of Data Science & Al 2025
28<sup>th</sup> October 2025

A WORLD TOP 100 UNIVERSITY



#### Based on

Sanz-Cruzado, J. & Lever, J. **Accelerating Cross-Encoders in Biomedical Entity Linking**. BioNLP @ ACL 2025, pp. 136-147



Paper link







#### **Motivation**

- Around 167,000 people die from cancer in the UK per year
- Pandemics like Covid-19 have a large impact on our lives

Understanding diseases and developing effective treatments is fundamental for our healthcare!

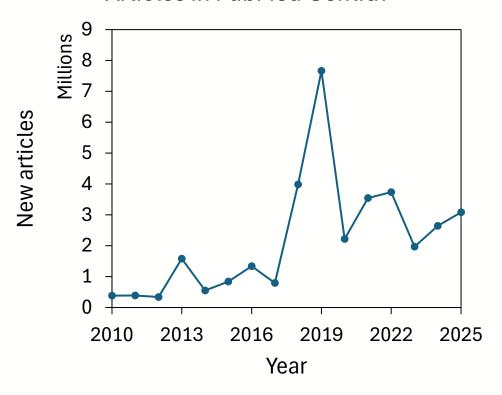
For this, it's important to know what people did before.

- What treatments were tested
- What genes / proteins are involved
- How is a disease related to others

•

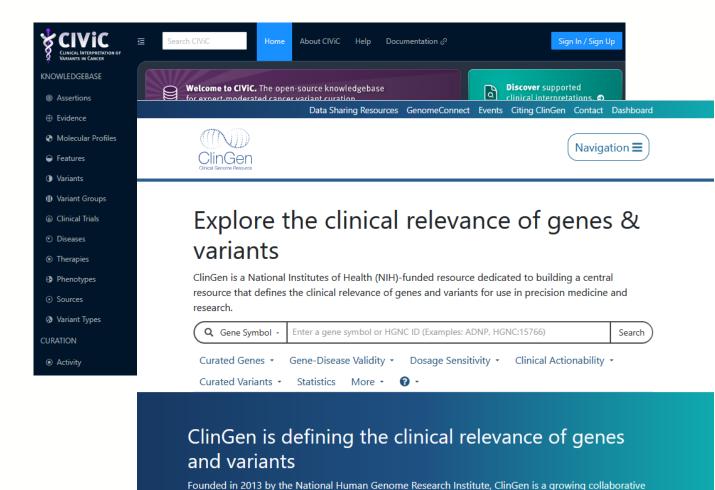
# Keeping up with prior research is becoming increasingly difficult!

#### Articles in PubMed Central





#### **Biomedical databases**



effort, involving three grants, nine principal investigators and over 2,700 contributors from more than

72 countries. Below are a series of recent updates that ClinGen has been working on.

- Used in research for precision medicine
- Examples:
  - ClinGen: Clinical Genome Resource
  - CIViC: Clinical Interpretation of Variants in Cancer
- Summarize information about
  - Diseases
  - Treatments
  - Gene variants
  - Etc.
- Information is usually manually annotated from research papers

Natural language processing can help!



# **Extracting information from biomedical texts**



Biomedical documents

Varicella is a highly contagious viral infection that causes an acute fever and blistered rash, mainly in children. Immunocompromised patients infected with the virus need intravenous treatment with the antiviral aciclovir.



#### **Entities**

- Varicella
- Aciclovir
- Fever



#### Relations

- Varicella is an infection
- Varicella causes fever
- Aciclovir treats varicella



# More complex information

 Aciclovir treats varicella via intravenous treatment



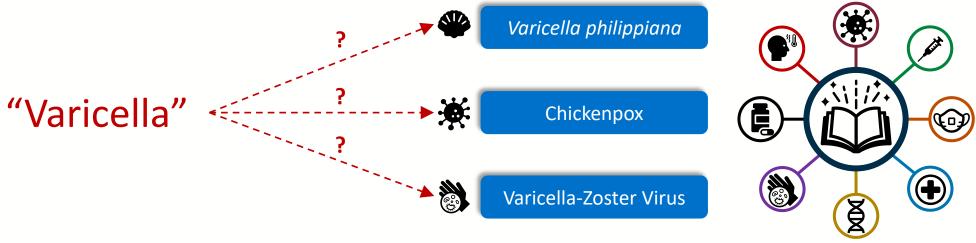
# What is biomedical entity linking?



"Varicella" is a highly contagious "viral infection" that causes an acute "fever" and "blistered rash", mainly in children.

"Immunocompromised patients" infected with the "virus" need "intravenous treatment" with the "antiviral" "aciclovir".

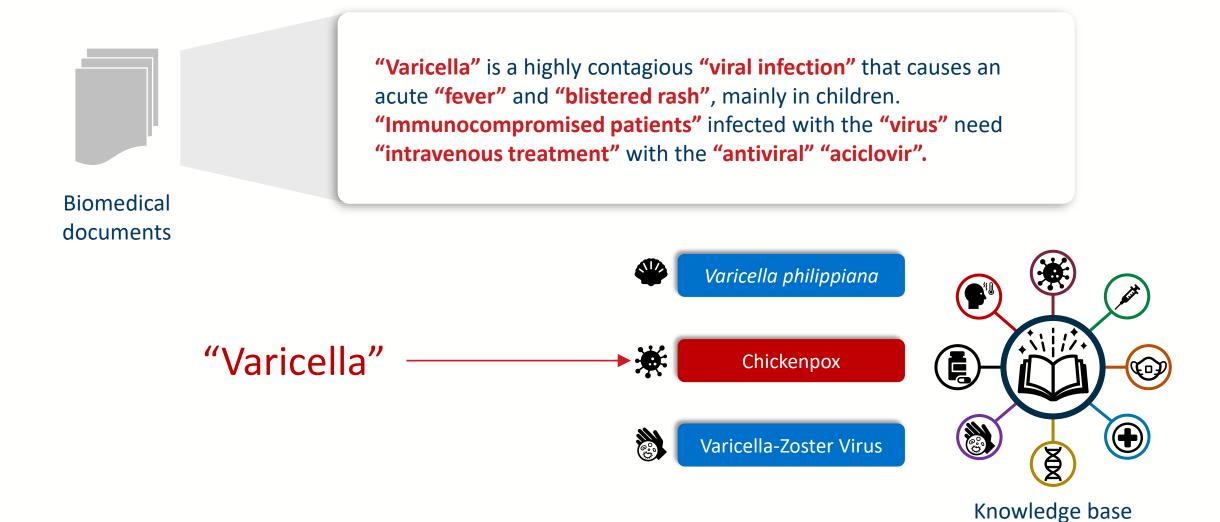
documents



Knowledge base



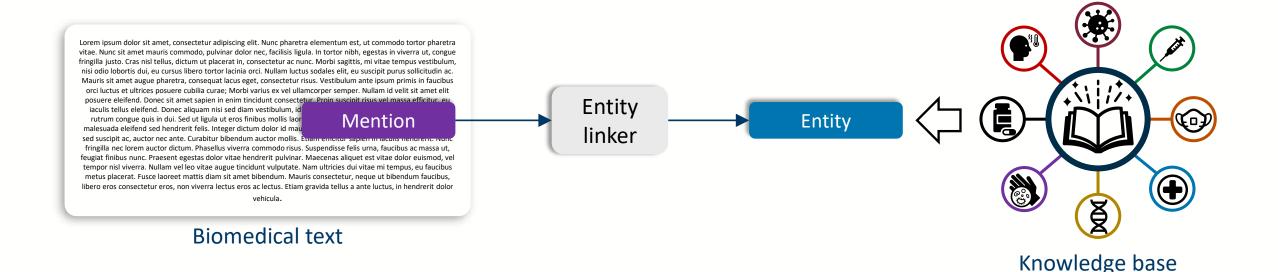
# What is biomedical entity linking?





# What is biomedical entity linking?

Biomedical entity linking matches mentions of biomedical concepts (diseases, chemicals) in text with unique entities within a knowledge base

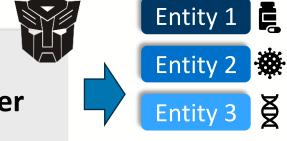




Sentence

# Two-stage entity linking

Transformers model



Reranker

**Candidate** Candidate Mention **Entities** retriever

> Selects a reduced set of potential entities for a mention



Very fast



Filters entities

Estimates probability of the mention matching the entity



Accurate



Ranks entities



Slow



# How slow are entity linking rerankers?

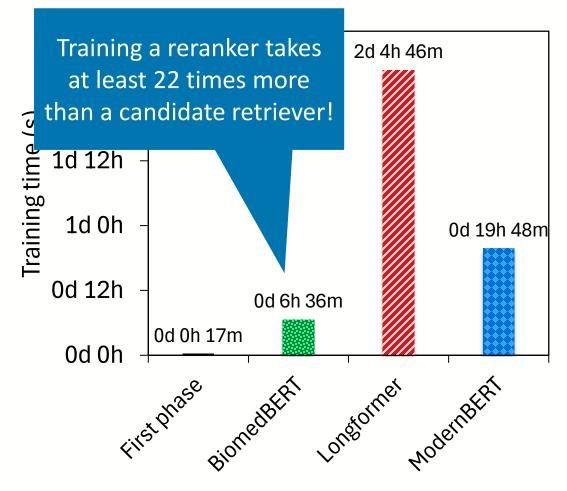
- We compare:
  - A first-stage candidate retrieval model
  - With three <u>cross-encoders</u> (state-of-the-art rerankers)
- We use the same hardware
- Dataset: MedMentions
  - Training data:
    - 2,635 biomedical paper abstracts
    - 211, 029 mentions
  - Test data:
    - 879 biomedical paper abstracts
    - 70,405 mentions

Far from what we might find on the complete PubMed!

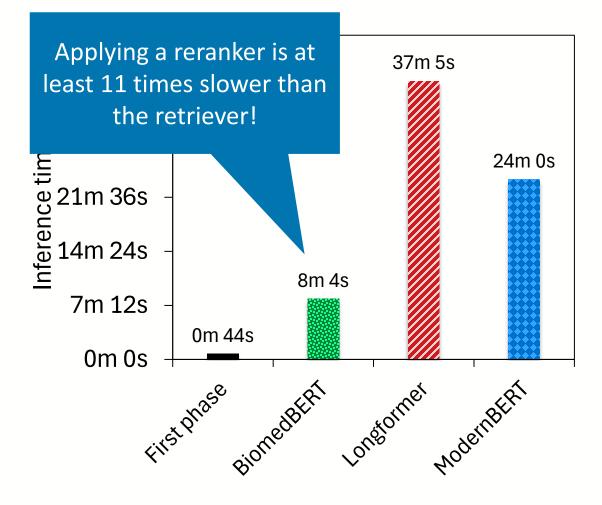


# How slow are entity linking rerankers?





#### Inference time

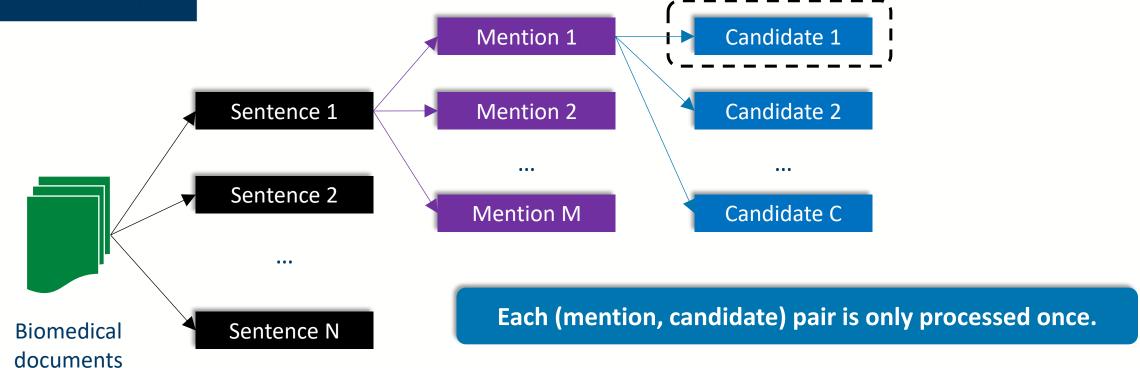


## Research question

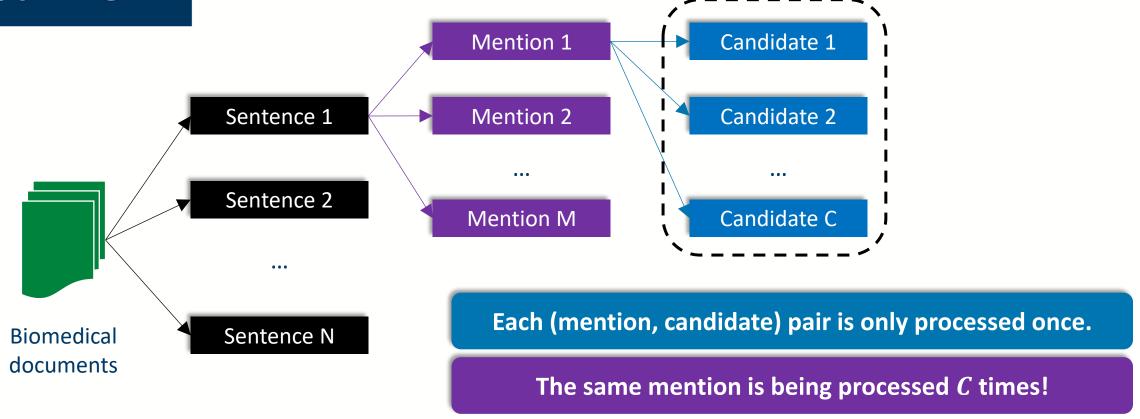
Cross-encoder rerankers are very slow...

Can we make them faster without harming performance?

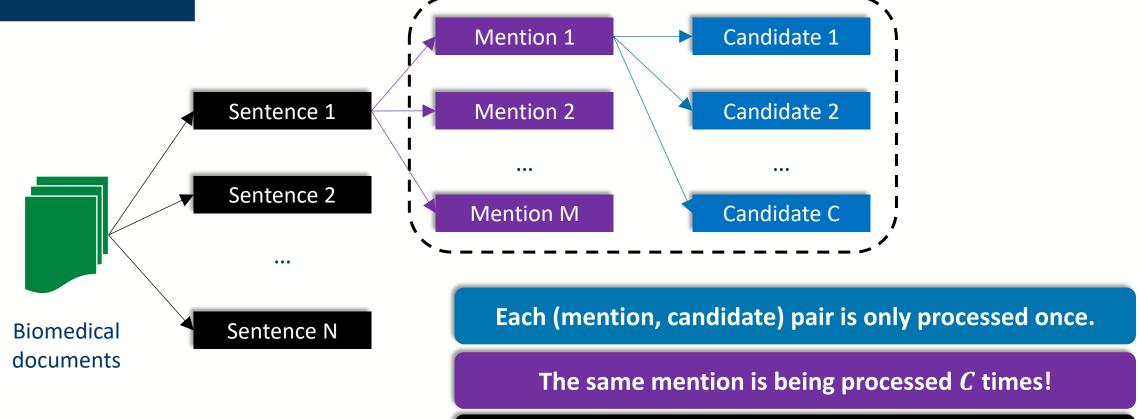






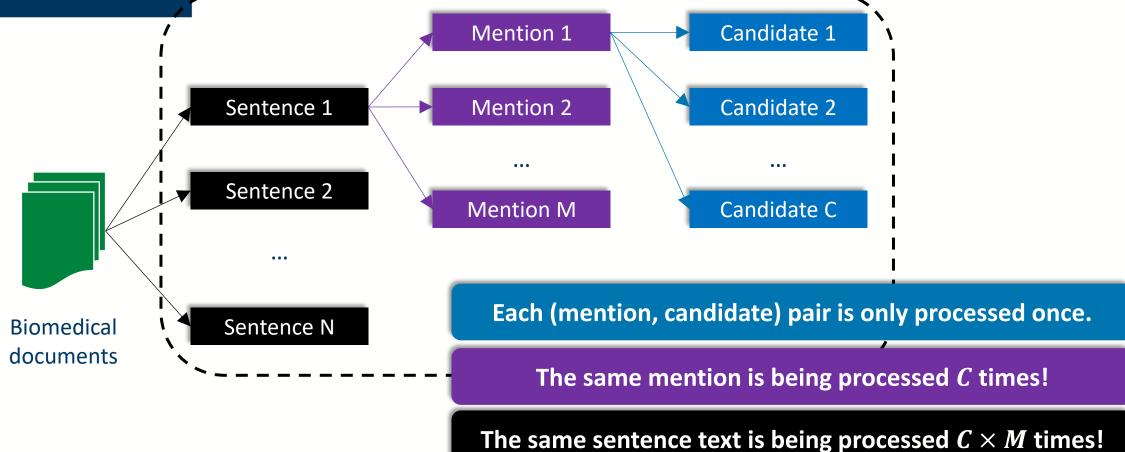






The same sentence text is being processed  $C \times M$  times!





And the document has N different sentences!



## **Accelerating cross-encoders**

Each (mention, candidate) pair is only processed once.

The same mention is being processed C times!

The same sentence text is being processed  $C \times M$  times!

And the document has N different sentences!

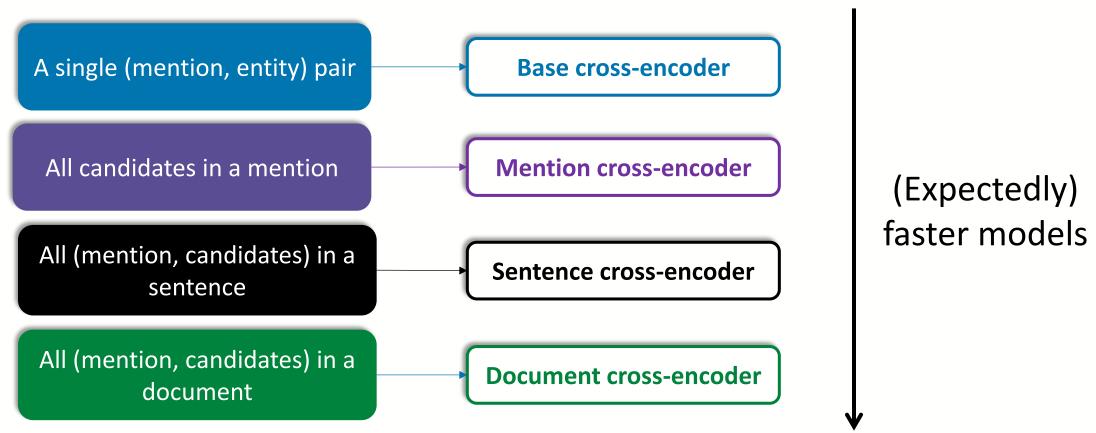
Both in training and inference

**Idea:** What if, instead of showing one (mention, candidate) pair on each call, we show several simultaneously?



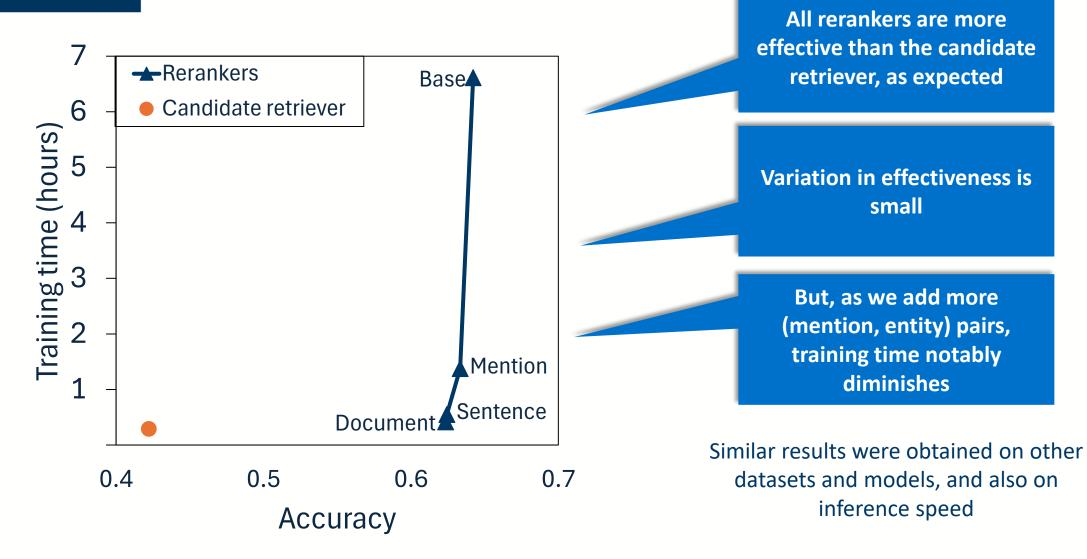
## **Accelerating cross-encoders**

We build re-rankers where we provide different amounts of (mention, candidate) pairs to process





#### How do these models work?





#### **Conclusions**

Processing more (mention, entity) pairs simultaneously has the following effects on entity linking cross-encoder rerankers

#### **Small variations in accuracy**

-3.42 to 2.76 % differences with base model

Major improvements in training speed

2.68x – 36.97x faster training than base model

Major improvements in inference speed

3.8x – 26.47x faster inference than base model

We can have a major training / inference speed boost at a small accuracy cost!

Our solutions are suitable for environments where speed is crucial (or data is huge)



# **Questions?**



Al4BioMed Group, University of Glasgow

javier.sanz-cruzadopuig@glasgow.ac.uk



X JavierSanzCruza



Javiersanzcruza.bsky.social



in https://www.linkedin.com/in/javier-sanzcruzado-puig/





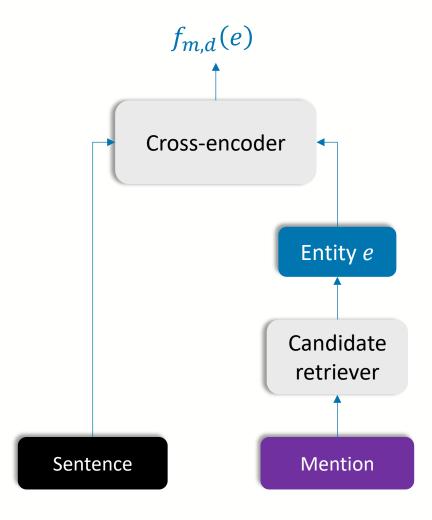


Each (mention, candidate) pair is only processed once.

The same mention is being processed *C* times!

The same sentence text is being processed  $C \times M$  times!

And the document has N different sentences!



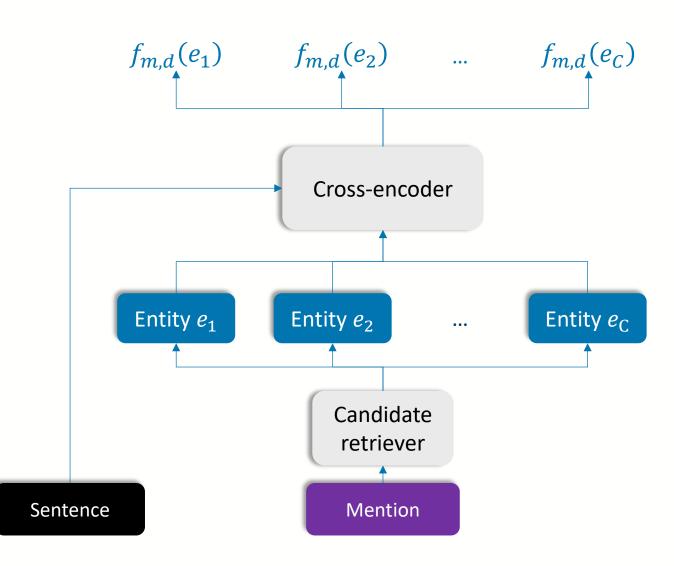


Each (mention, candidate) pair is only processed once.

The same mention is being processed *C* times!

The same sentence text is being processed  $C \times M$  times!

And the document has N different sentences!





Each (mention, candidate) pair is only processed once.

The same mention is being processed *C* times!

The same sentence text is being processed  $C \times M$  times!

And the document has N different sentences!

While the cross encoder uses template for each candidate:

Text [SEP] Mention [MASK] Entity e name



Each (mention, candidate) pair is only processed once.

The same mention is being processed *C* times!

The same sentence text is being processed  $C \times M$  times!

And the document has N different sentences!

• The mention cross-encoder receives a sentence using the following a template for each mention:

```
Text [SEP] Mention [MASK] Entity e_1 name [SEP] Mention [MASK] Entity e_2 name ... [SEP] Mention [MASK] Entity e_C name
```

• Therefore, the score of the entity  $e_i$  is the probability of its [MASK] token taking value 1



Each (mention, candidate) pair is only processed once.

The same mention is being processed once.

The same sentence text is being processed  $C \times M$  times!

And the document has N different sentences!

• The mention cross-encoder receives a sentence using the following a template for each mention:

```
Text [SEP] Mention [MASK] Entity e_1 name [SEP] Mention [MASK] Entity e_2 name ... [SEP] Mention [MASK] Entity e_C name
```

• Therefore, the score of the entity  $e_i$  is the probability of its [MASK] token taking value 1

And so on for the other models