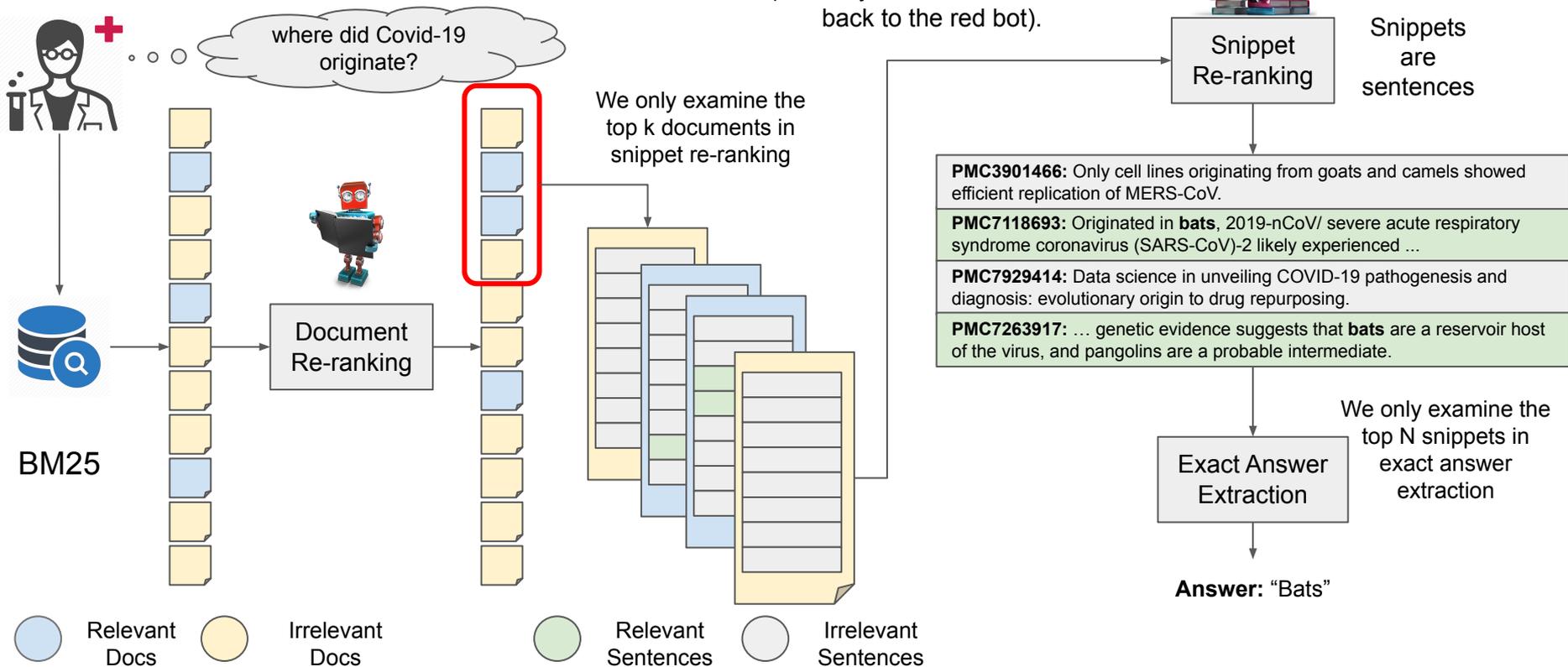


# A Neural Model for Joint Document and Snippet Ranking in Question Answering for Large Document Collections

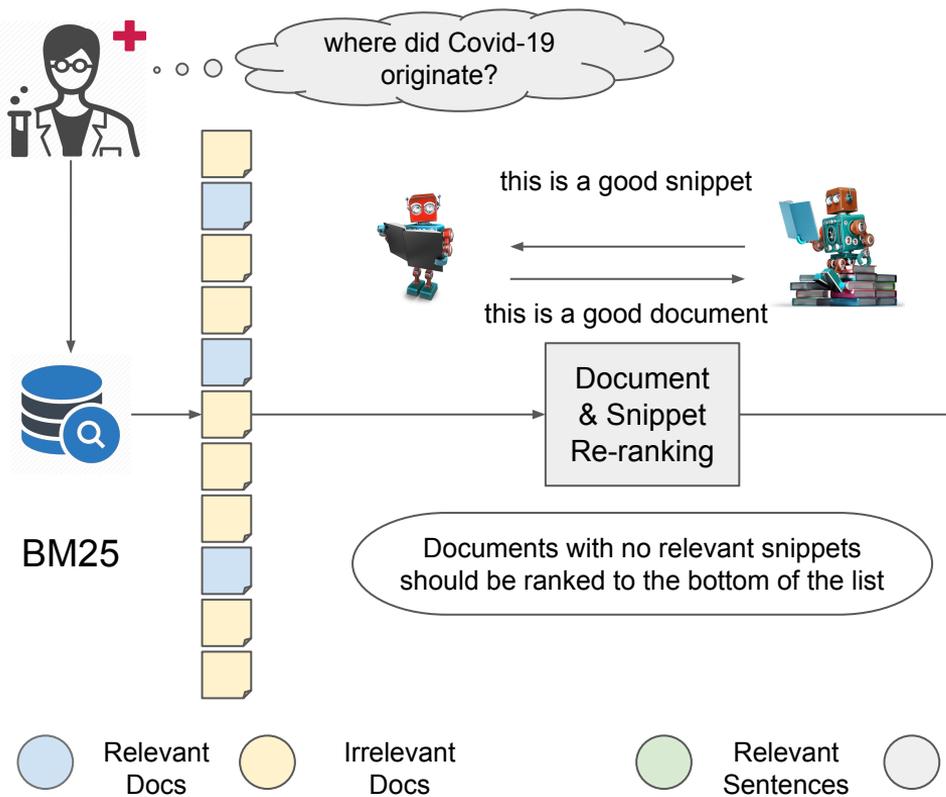
\*Dimitris Pappas, Ion Androutsopoulos



# Retrieval Pipeline



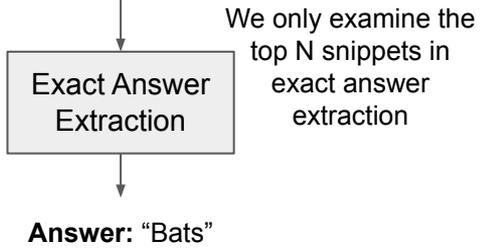
# Retrieval Pipeline



- Snippet Re-ranking handles
- relevant snippets
  - irrelevant snippets

Pipeline methods missed article PMC7113610 in document retrieval so the sentence retrieval module did not have the chance to examine the snippet

|  |
|--|
| <p><b>PMC7113610:</b> Genomic analysis revealed that SARS-CoV-2 ..., therefore <b>bats</b> could be the possible primary reservoir.</p>                |
| <p><b>PMC7118693:</b> Originated in <b>bats</b>, 2019-nCoV/ severe acute respiratory syndrome coronavirus (SARS-CoV)-2 likely experienced ...</p>      |
| <p><b>PMC7929414:</b> Data science in unveiling COVID-19 pathogenesis and diagnosis: evolutionary origin to drug repurposing.</p>                      |
| <p><b>PMC7263917:</b> ... genetic evidence suggests that <b>bats</b> are a reservoir host of the virus, and pangolins are a probable intermediate.</p> |



# Data collection and Indexing



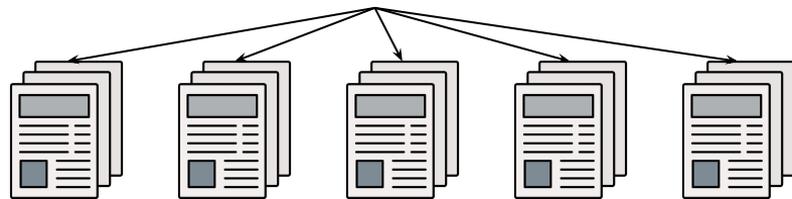
**BioA?Q** 2,747 questions

Which disease is caused by de novo VPS4A mutations?

32.6M biomedical articles  
27.8M biomedical articles in English

21.8M articles with  
Title and Abstract

**PubMed**



Data collection and preprocessing



BM25

# PDRMM

And computes three cosine similarity matrices across query and document terms

PDRRM uses three kinds of word representations:

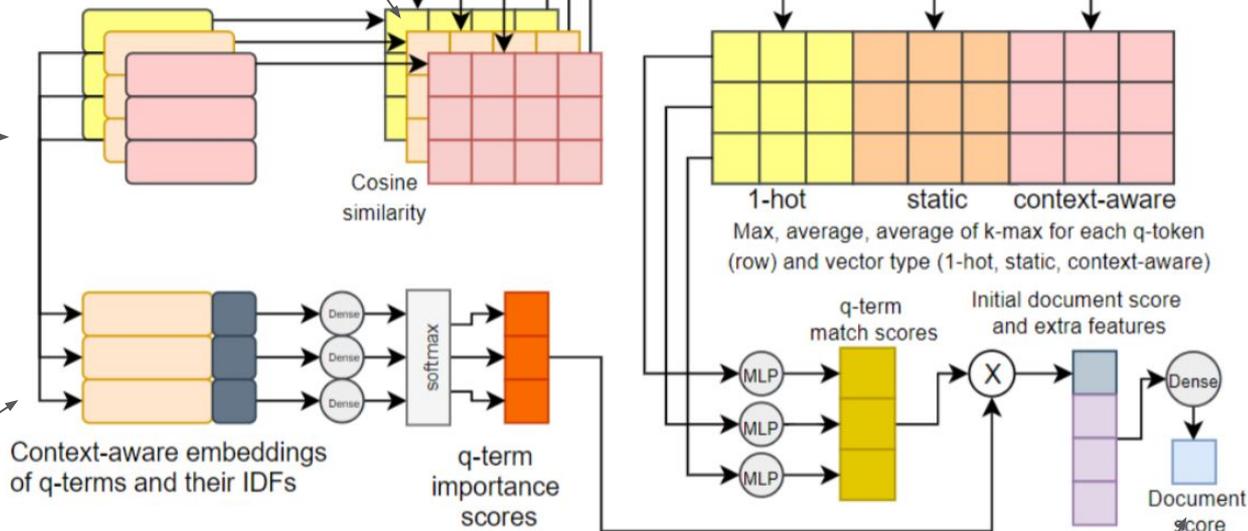
- one hot for exact match
- pretrained w2v vectors
- cnn contextual representations

PDRRM uses three kinds of word representations:

- one hot for exact match
- pretrained w2v vectors
- cnn contextual representations

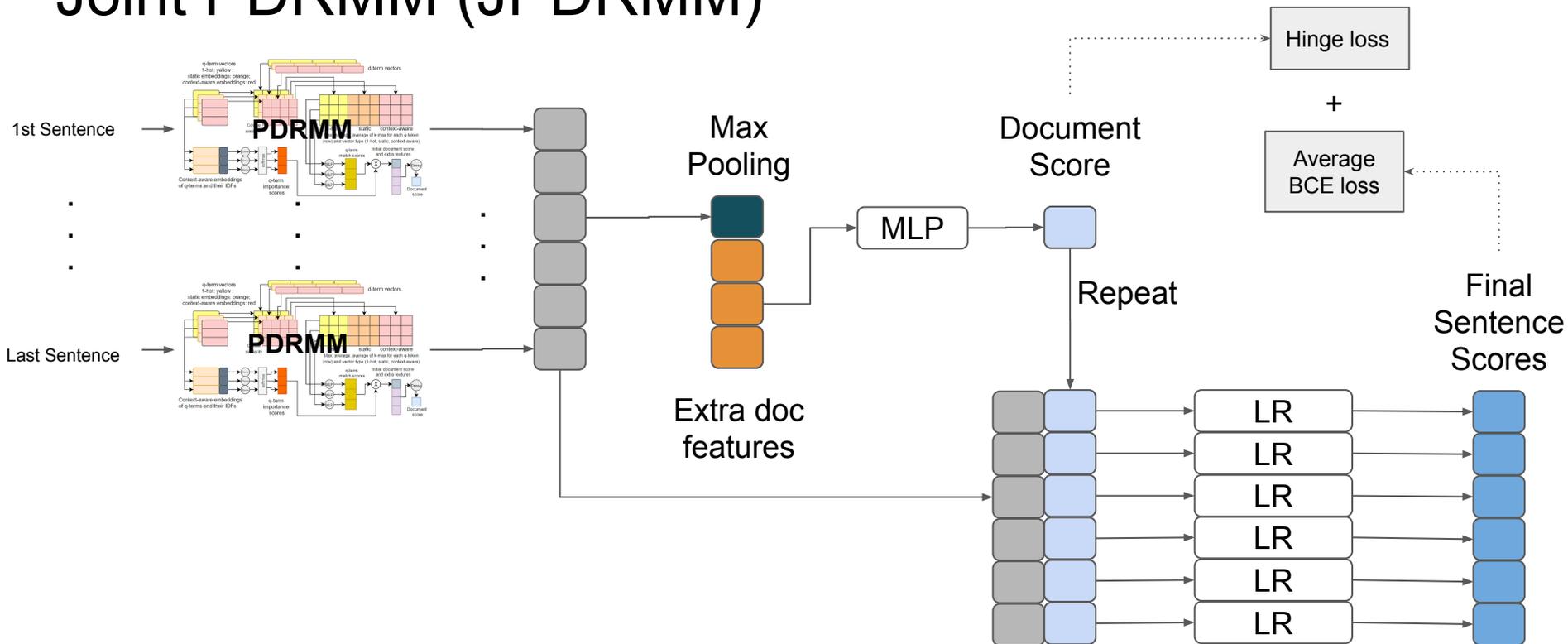
q-term vectors  
1-hot: yellow ;  
static embeddings: orange;  
context-aware embeddings: red

d-term vectors

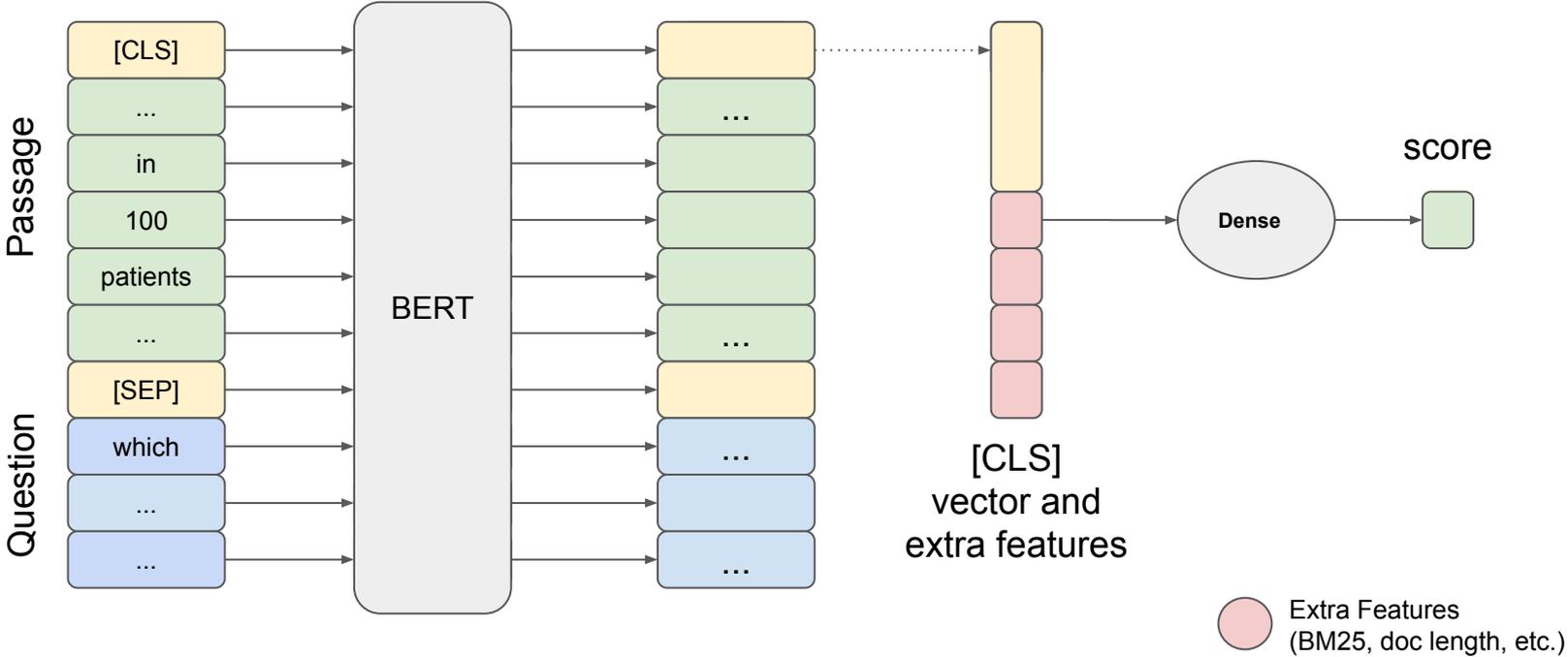


to compute a relevance score for the entire document

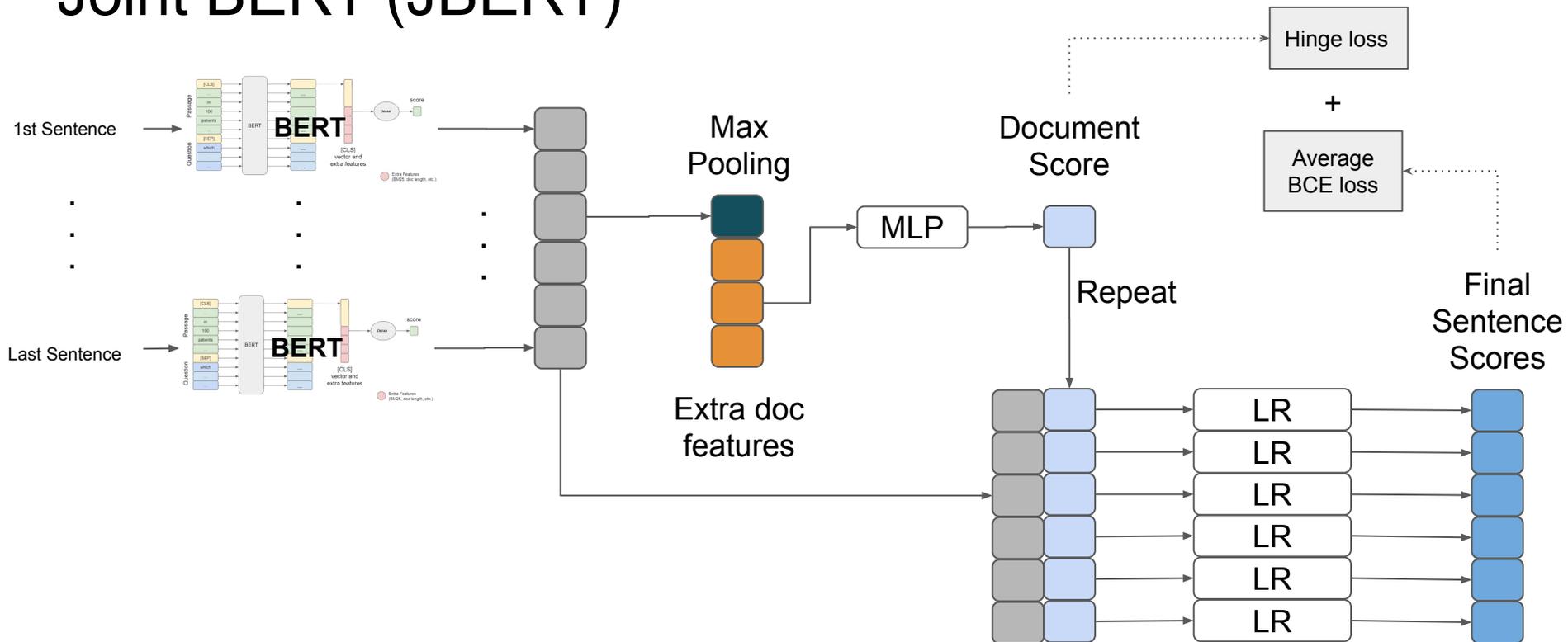
# Joint PDRMM (JPDRMM)



# BERT



# Joint BERT (JBERT)



# BioASQ 7 Results (Best of 14 methods)

Models perform better when Bert weights are not fine-tuned for the retrieval task

| Method           | Method Type | Params | Doc. MAP | Snip. MAP |
|------------------|-------------|--------|----------|-----------|
| BERT+PDRMM       | Pipeline    | 109.5M | 8.79     | 9.63      |
| JPDRMM           | Joint       | 5.79k  | 6.69     | 15.72     |
| BJPDRMM-ADAPT-NF | Joint       | 3.5M   | 7.42     | 17.35     |
| JBERT-ADAPT-NF   | Joint       | 6.3K   | 7.84     | 16.53     |
| Oracle           | n/a         | 0      | 19.24    | 25.18     |

In ADAPT models we use a linear combination of all embeddings that Bert produces throughout its layers

JPDRMM remains competitive in snippet retrieval despite using orders of magnitude fewer parameters.

# BioASQ 7 Results (post-contest exper

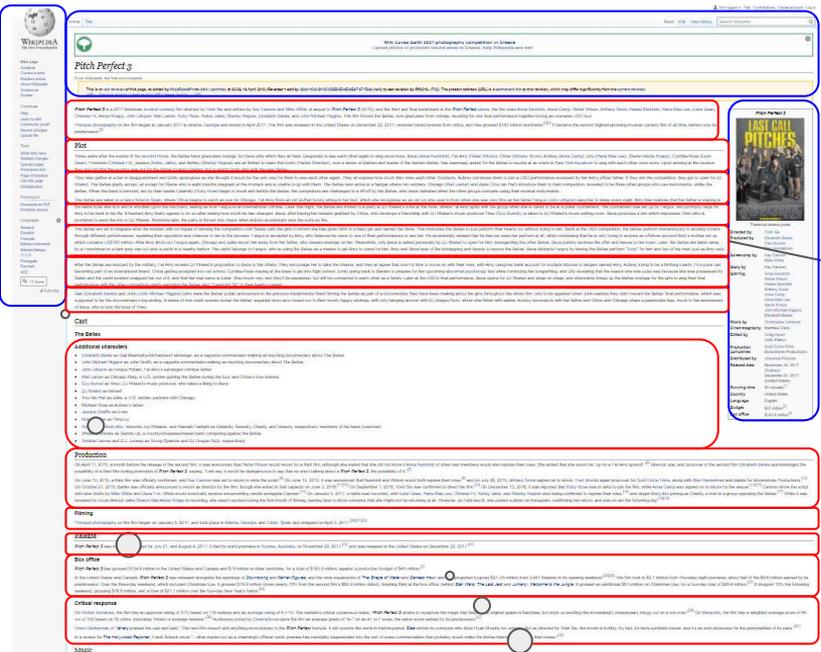
Human annotators examine all results submitted by contestants to identify unseen gold data

| Method                   | Method Type | Parameters  | Before expert inspection   |   | After expert inspection   |   |
|--------------------------|-------------|---|--|---|---|---|
|                          |             |   | Doc. MAP   | Snip. MAP   | Doc. MAP  | Snip. MAP   |
| BERT+PDRMM               | Pipeline    | 109.5M  | 7.29  | 7.58  | 14.86   | 15.61   |
| JPDRMM                   | Joint       | 5.79K  | 5.16   | 12.45  | 16.55  | 21.98   |
| BJPDRMM-NF               | Joint       | 3.5M  | 6.18   | 13.89  | 14.65   | 23.96  |
| Best BIOASQ 7 competitor | n/a         | n/a   | n/a  | n/a   | 13.18  | 14.98  |

BIOASQ 7 test batches 4 and 5 before and after post-contest

# Original Natural Questions Dataset

T. Kwiatkowski et al. 2019



Question: where do the bellas go in pitch perfect 3

**HTML content:** `<h2><span class="mw-headline" id="Plot">Plot</span></h2><p>The Bellas are taken to a fancy hotel in`

Answer: Spain

URL: [https://en.wikipedia.org/?title=Pitch\\_Perfect\\_3&oldid=837158847](https://en.wikipedia.org/?title=Pitch_Perfect_3&oldid=837158847)

Boilerplate and tables (blue color)

Paragraphs (red color)

We handle each paragraph of the entire webpage as a document which we index in our database

# Natural Questions Modification

**Question:** where do the bellas go in pitch perfect 3

**HTML content:** `<h2><span class="mw-headline" id="Plot">Plot</span></h2><p>The Bellas are taken to a fancy hotel in <a href="/wiki/Spain" title="Spain">Spain</a>, where Chloe begins to catch an eye for Chicago. Fat Amy finds an old stuffed bunny sitting in her bed, which she recognizes as an old toy she used to hold when she was very little as her father Fergus (<a href="/wiki/John_Lithgow" title="John Lithgow">John Lithgow</a>) sang her to sleep every night. Amy then realizes that her father is staying in the same hotel she is in and is shocked upon the discovery, seeing as how Fergus is an international criminal. Later that night, the Bellas are invited to a party at DJ Khaled's suite at the hotel, where Fat Amy splits with the group when she is called to be at a poker tournament. The tournament was set up by Fergus, who promptly begs for Amy to be back in his life. A hesitant Amy finally agrees to do so after seeing how much he has changed. Beca, after having her breasts grabbed by Chloe, who develops a friendship with DJ Khaled's music producer Theo (<a href="/wiki/Guy_Burnet" title="Guy Burnet">Guy Burnet</a>), is taken to DJ Khaled's music editing room. Beca produces a mix which impresses Theo who is prompted to send the mix to DJ Khaled. Moments later, the party is thrown into chaos when Aubrey accidentally sets the suite on fire.</p>`

Plain text  
extraction

**Question:** where do the bellas go in pitch perfect 3

**Sentences:**

The Bellas are taken to a fancy hotel in **Spain**, where Chloe begins to catch an eye for Chicago.

Fat Amy finds an old stuffed bunny sitting in her bed, which she recognizes as an old toy she used to hold when she was very little as her father Fergus ( John Lithgow) sang her to sleep every night. Amy then realizes that her father is staying in the same hotel she is in and is shocked upon the discovery, seeing as how Fergus is an international criminal.

Later that night, the Bellas are invited to a party at DJ Khaled's suite at the hotel, where Fat Amy splits with the group when she is called to be at a poker tournament.

The tournament was set up by Fergus, who promptly begs for Amy to be back in his life.

A hesitant Amy finally agrees to do so after seeing how much he has changed.

Beca, after having her breasts grabbed by Chloe, who develops a friendship with DJ Khaled's music producer Theo (Guy Burnet), is taken to DJ Khaled's music editing room.

Beca produces a mix which impresses Theo who is prompted to send the mix to DJ Khaled.

Moments later, the party is thrown into chaos when Aubrey accidentally sets the suite on fire.

We handle each paragraph of the entire webpage as a document which we index in our database

"Spain" is the answer therefore the yellow sentence is annotated as relevant.

# Modified Natural Questions Results

| Method                | Document Retrieval<br>(paragraph retrieval) |          |          | Snippet Retrieval<br>(sentence retrieval) |          |          |
|-----------------------|---|----------|----------|---|----------|----------|
|                       | MRR   | Recall@1 | Recall@2 | MRR                                       | Recall@1 | Recall@2 |
| Pipeline: BM25+BM25   | 30.18                                       | 16.50    | 29.75    | 8.19                                      | 3.75     | 7.13     |
| Pipeline: PDRMM+PDRMM | 40.33                                       | 28.25    | 38.50    | 22.86                                     | 13.75    | 22.75    |
| Joint: JPDRMM         | 36.50                                       | 24.50    | 36.00    | 26.92                                     | 19.00    | 25.25    |

We use Recall@1 and Recall@2 because there are at most two relevant documents and two relevant snippets

- 👎 joint model does not outperform pipelines in document retrieval
- 👍 Joint model outperforms pipelines in snippet retrieval
- 👍 Both models outperform BM25

# Conclusions

- All the neural pipelines and joint models we considered improve the traditional BM25 ranking on both datasets.
- Joint models vastly outperform the corresponding pipelines in snippet retrieval.
- The Joint PDRMM model that does not use BERT is competitive with BERT-based models.
- We surpassed all competitive systems in BioASQ 7.

# Future Work

- Other pre-trained models on biomedical data (e.g. BioBERT) or QA data (e.g. SQuAD)
- Extend our joint models to cover selecting exact answers too (Document Retrieval, Snippet extraction, Factoid QA)
- Use external knowledge (e.g. graph embeddings)



<http://nlp.cs.aueb.gr/publications.html>



@dvpappas, @AUEBNLPGroup



pappasd@aueb.gr