

Today's agenda:

- ① classic retrieval (keyword-based)
- ② neural retrieval
- ③ RAG
- ④ General examples of compound AI systems + challenges

## Classic IR / BM25

- scope: large corpus of text documents (e.g., Wikipedia - potentially split into passages)
  - input: textual query (e.g., question)
  - output: Top-K RANKING of relevant docs
- \* simple case: "single-term" queries
- build term document matrix
  - each cell says how many times a term appears in a doc (perhaps w/ reweighting)
  - if query has a single term, just return K documents w/ largest weights
- for a single term

\* "against" has off-weight for d4, d7

	d1	d2	d3	d4	d5	d6	d7	d8	d9	d10
against	0	0	0	1	0	0	3	2	3	0
age	0	0	0	1	0	3	1	0	4	0
agent	0	0	0	0	0	0	0	0	0	0
ages	0	0	0	0	0	2	0	0	0	0
ago	0	0	0	2	0	0	0	0	3	0
agree	0	1	0	0	0	0	0	0	0	0
ahead	0	0	0	1	0	0	0	0	0	0
ain't	0	0	0	0	0	0	0	0	0	0
air	0	0	0	0	0	0	0	0	0	0
aka	0	0	0	1	0	0	0	0	0	0

\* what about "multi-term" queries?

$$\text{Relevance Score}(\text{query}, \text{doc}) = \sum_{\text{term} \in \text{query}} \text{weight}_{\text{doc}, \text{term}}$$

\* core questions in classical IR:

- How do we best tokenize queries / documents?
- How do we best WEIGHT term-document pairs?

\* Intuitions behind term-document weights:

### 1) FREQUENCY

- if term  $t$  occurs frequently in document  $\delta$  - doc is more likely to be relevant for queries w/  $t$

### 2) NORMALIZATION

- if term  $t$  is rare overall, any document  $\delta$  including  $t$  is likely very relevant

- if  $\delta$  is overall short and contains  $t$  → probably important

## TF-IDF: Term-Document weighting

•  $N = |\text{collection}|$  which doc contain

•  $\text{DF}(\text{term}) = |\{ \text{doc } \in \text{collection} : \text{term } \in \text{doc} \}|$

, term freq

$\text{TF}(\text{term}, \text{doc}) = \log(1 + \text{Freq}(\text{term}, \text{doc}))$

↳ grows sublinearly w/ frequency

$$\text{IDF}(\text{term}) = \log \left( \frac{N}{\text{df}(\text{term})} \right)$$

\* if many docs contain term, IDF will be smaller

/  
inverse  
doc  
frequency

↳ grows sublinearly w/  $1/\text{IDF}$

\* accounts for normalization

$$\text{TF-IDF}(\text{term}, \text{doc}) = \text{TF}(\text{term}, \text{doc}) \times \text{IDF}(\text{term})$$

$$\text{TF-IDF}(\text{query}, \text{doc}) = \sum_{\text{term} \in \text{query}} \text{TF-IDF}(\text{term}, \text{doc})$$

\* BM25: Inspired from TF-IDF, but better:

→ change: term frequency saturates to constant value  
and penalizes LONGER doc

$$\text{IDF}(\text{term}) = \log \left( 1 + \frac{N - \text{df}(\text{term}) + 0.5}{\text{df}(\text{term}) + 0.5} \right)$$

$$\text{TF}(\text{term}, \text{doc}) = \frac{\text{freq}(\text{term}, \text{doc}) \cdot (K+1)}{\text{freq}(\text{term}, \text{doc}) + K \cdot (1 - b + b \cdot \frac{\text{len}(\text{doc})}{\text{avg doc len}})}$$

\* k, b are params

• Efficient Retrieval: Inverted Indexing

raw collection: docs → terms

term doc matrix: terms → docs

↳ could be very sparse!

\* inverted index: sparse representation

of term-doc matrix (maps each unique term  $T$  to "posting list")

↳ posting list contains list of docs where term appears, + frequency

of term  $t$  in each doc

\* in terms of calculating TF, IDF:

- when we see a query:

- for each term:

→  $df(\text{term})$  stored in inverted index

→  $tf(\text{term}, \text{doc})$  stored

→ can calculate TF-IDF

\* until 2019, when BERT-based ranking

came out - BM25 was easy to deploy.

strong baseline

- JUDGING IR SYSTEMS:

1) effective (good accuracy) and efficient

various measures

- latency per query
- throughput (queries/sec)
- space for index
- HW required
- scaling to larger indexes

→ success and MRR

- let  $r_{\text{rank}} \in \{1, 2, 3, \dots, 3\}$  be position in returned top- $K$  of 1<sup>st</sup> RELEVANT DOC

success @  $K = \begin{cases} 1 & \text{if } r_{\text{rank}} \leq K \\ 0 & \text{otherwise} \end{cases}$  → if user needs returned list to contain anything in top- $K$

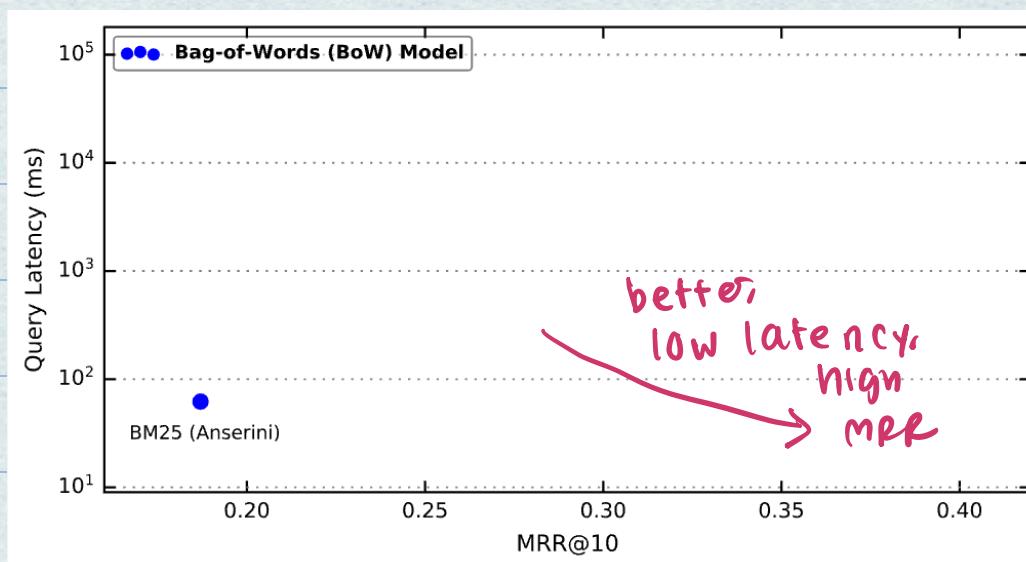
reciprocal rank at  $K = \begin{cases} 1/r_{\text{rank}} & \text{if } r_{\text{rank}} \leq K \\ 0 & \text{otherwise} \end{cases}$  → covers all positi

↳ MRR → mean of this across queries

→ look at  $K=1, 5, 100, 1000, \dots$

→ for BM25, we get:

low latency, but, low MRR?



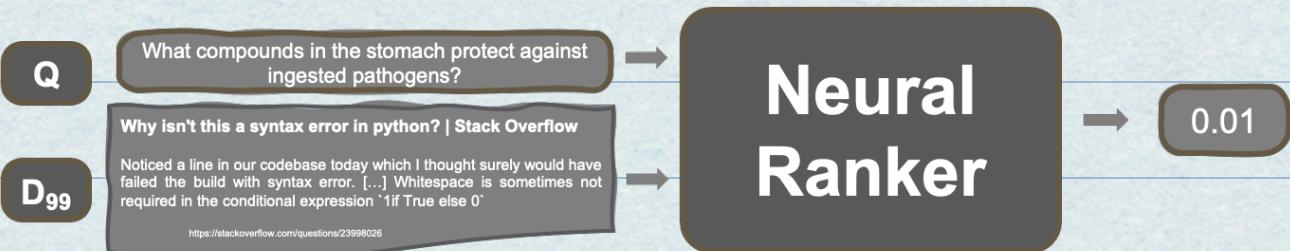
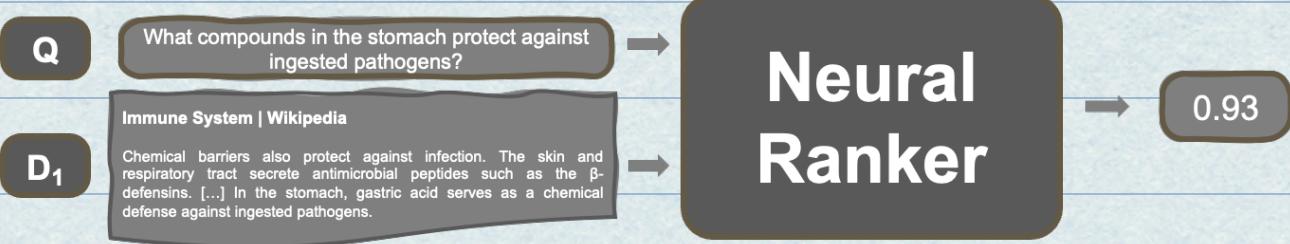
- now do we increase MRR (accuracy) for some cost of query latency?

\* MS MARCO: dataset widely used

for retrieval (Bing queries, 9 mil. passed)

What would input and output of neural ranking be?

- if neural retriever were a blackbox:



- input = query, document

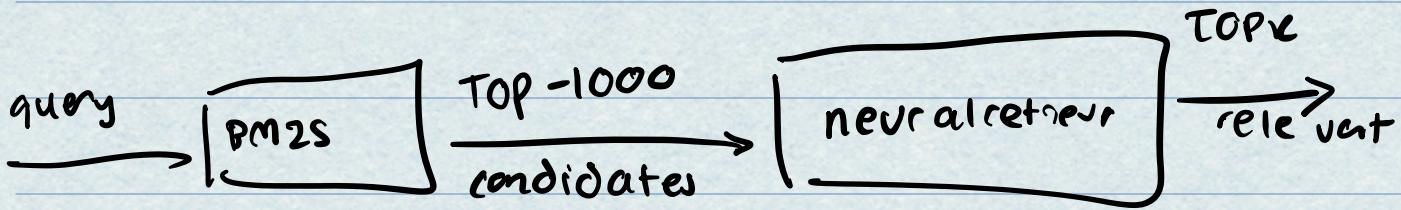
Om :

- output = score (relevance, higher = more relevant)
- TOPK docs = ranked ordering of docs according to score

\* if we have 9 billion passages, and each query, doc pair inference takes 1ns  $\Rightarrow$  still 9s per query (unacceptable!)

IDEA:

- do rerankings:



\* potential problem: if BM2S misses

some docs, we won't get them in

our ans ("artificial recall ceiling")

## HOW DO WE CAPTURE QUERY-DOCUMENT RELEVANCE?

1. Tokenize query and document

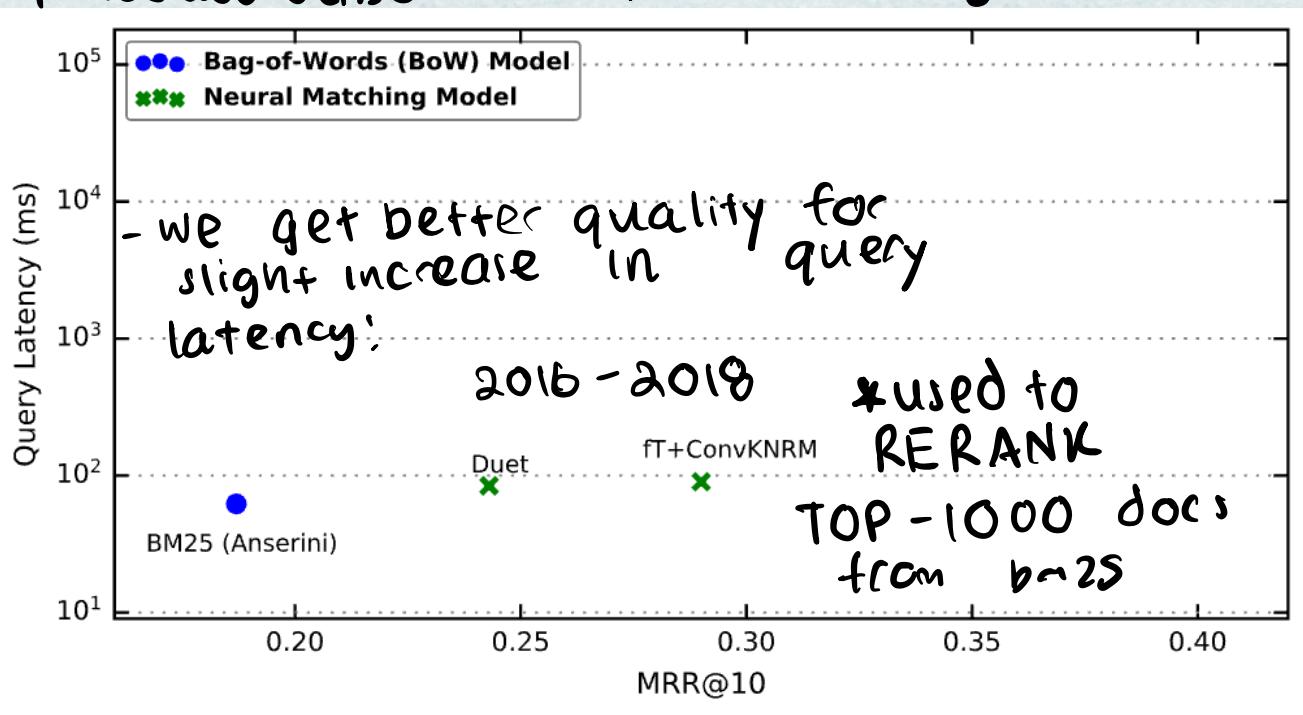
2. Embed tokens into static vector

repr

3. Build query-document interaction matrix

↳ similarity scores b/w each pair of word embeddings across query / doc, such as cosine similarity

4. Reduce dense matrix to a single score



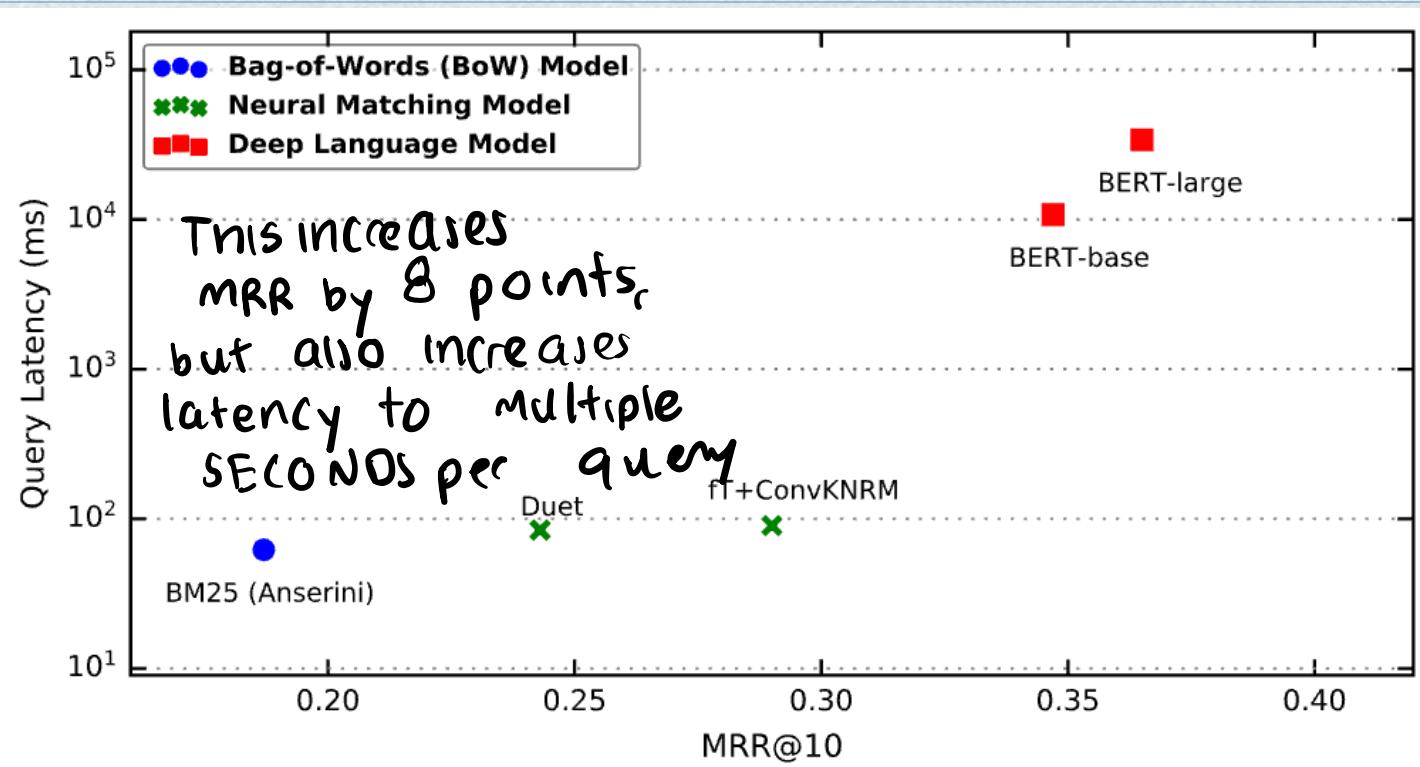
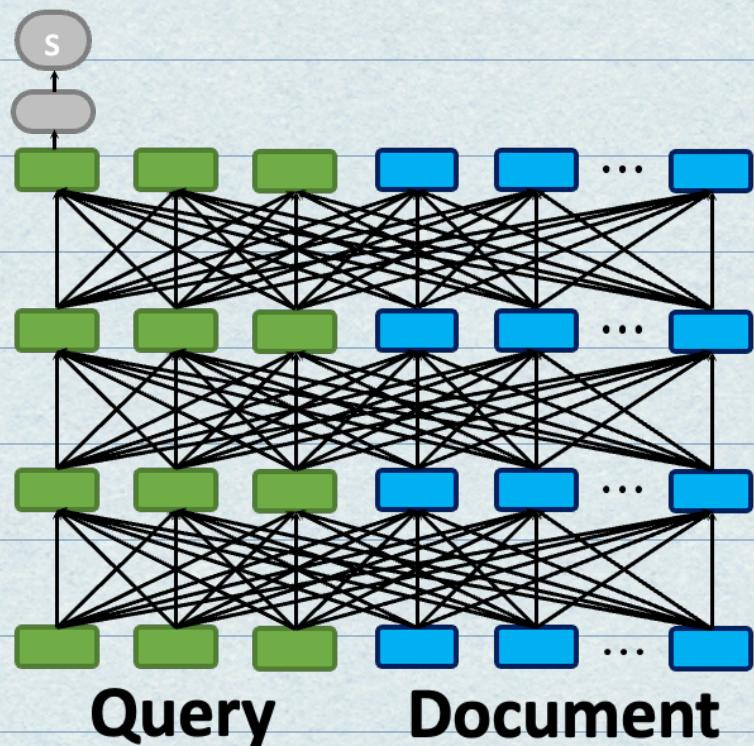
# POST 2018: All-to-All interaction w/ BERT

1. Feed BERT the query and document as one sentence w/ two segments

2. Run this through BERT layers

3. Extract final output embedding

- use custom linear layer to reduce to single score



How do we INCREASE quality but keep latency low?

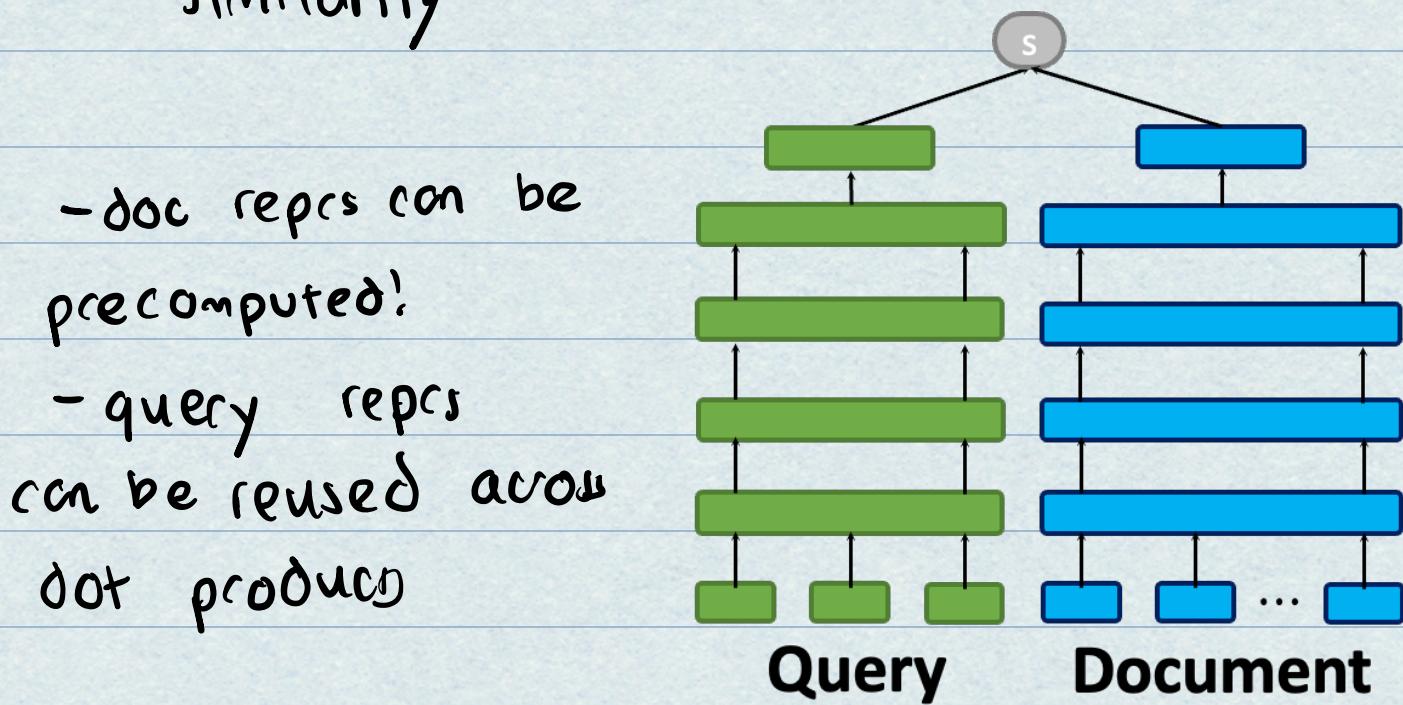
\* Bert rankers slow as computations are redundant:

- repr. query (1000 times for 1000 doc)
- repr. doc (for EACH inference)
- matching

\* can we precompute document representations and "cache" these for use across queries?

## REPRESENTATION SIMILARITY:

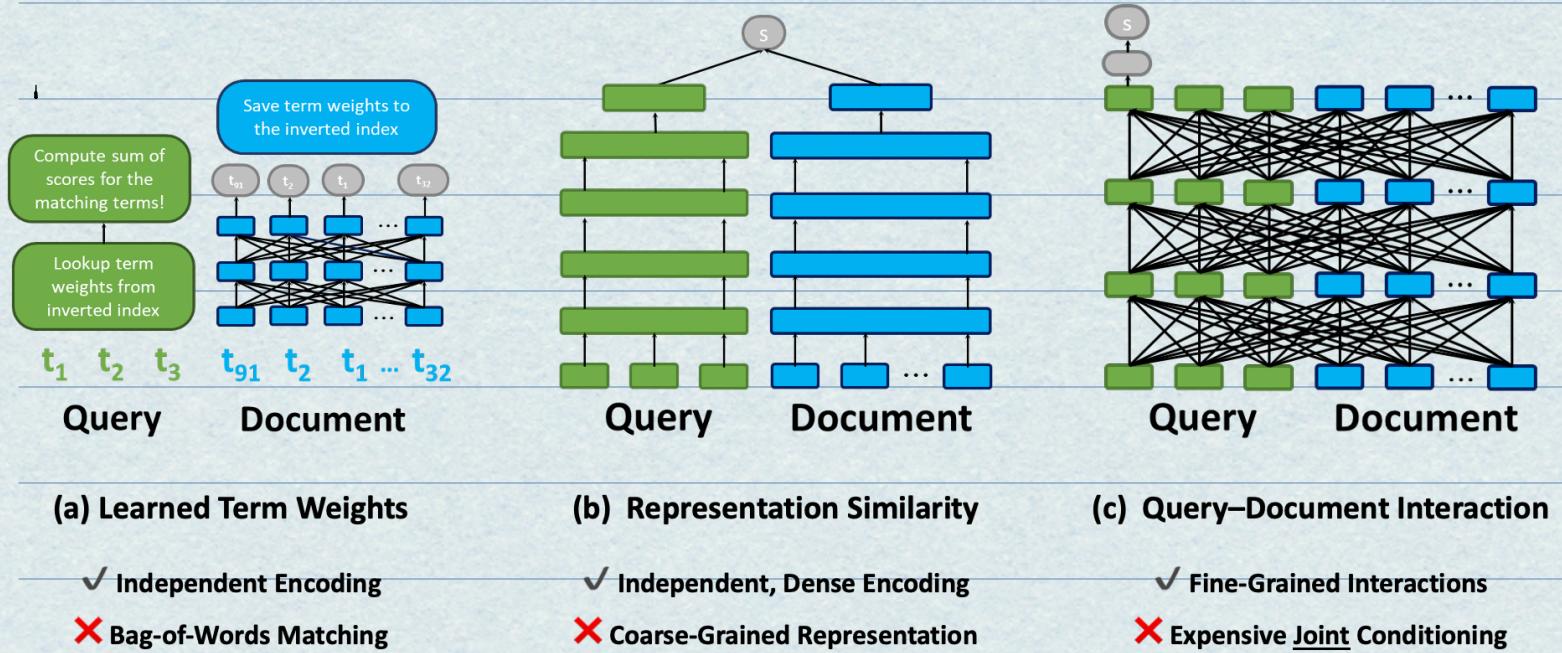
- tokenize query / doc
- independently encode query / doc
- use dot product to estimate similarity



- Downsides of JUST doing representation similarity:

- single-vector repr could be coarse-grained
- no FINE-GRAINED interaction
  - ↳ e.g. b/w individual tokens in query to docs

## SUMMARY SO FAR

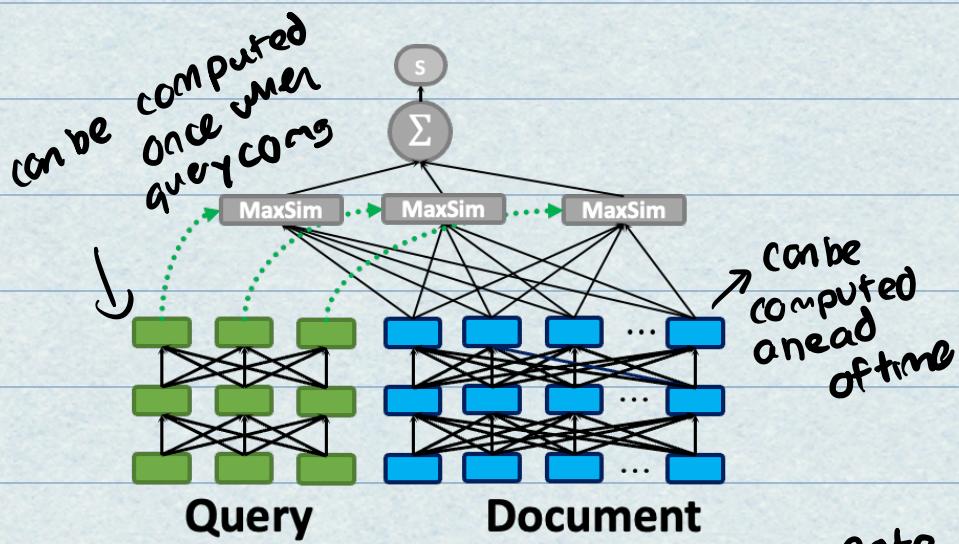


- we didn't look at this specifically but this is like bmm2s w/ better weights

LATE INTERACTION: keep precomputation, but have fine-grained interactions (CoBERT, Khattab & Zaharia)

- insight: feed BERT into query / each doc, but KEEP output embeddings corresponding to each token (don't feed in output to final linear layer)  
↳ recall model produces  $1 \times d$  vector per token

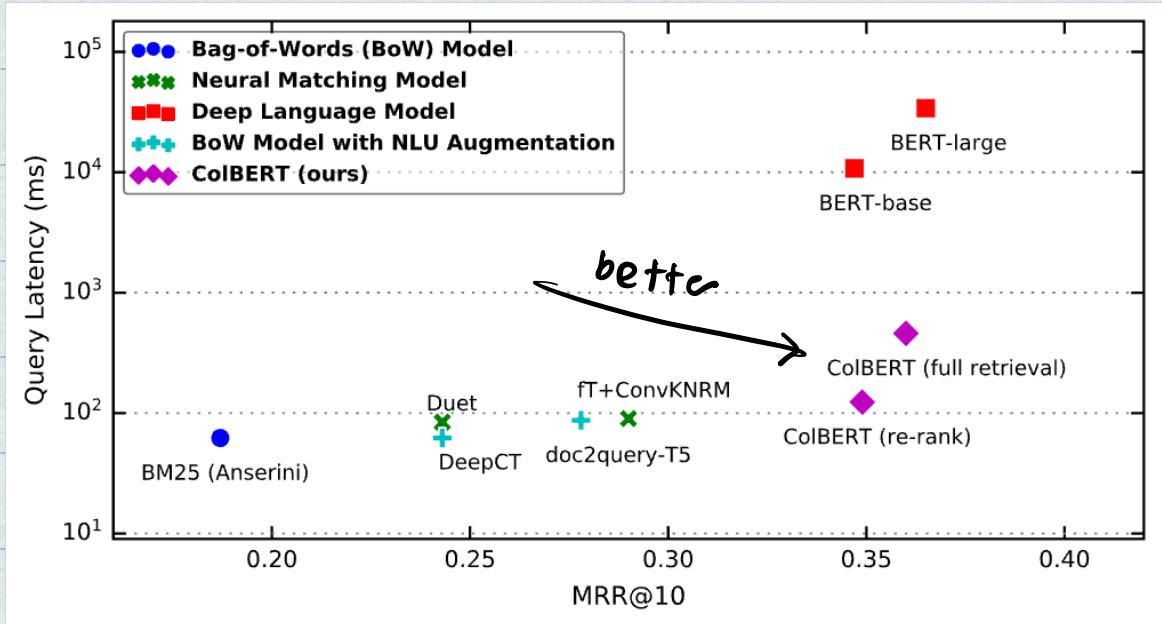
- now for each query token embedding, compute max similarity score across all doc embeddings
  - sum all these partial scores.



- note that colbert
- (d) Late Interaction represents doc as  $N \times d$  matrix (rather than a single vector)

Omar Khattab and Matei Zaharia. "CoBERT: Efficient and effective passage search via contextualized late interaction over BERT." SIGIR'20.

- w/ late interaction, we can resolve mismatches from before!



RAG: LLMs struggle on knowledge-intensive tasks

1. retrieve docs relevant to input from an external DB (either retriever OR vector DB)
2. Augment LLM input w/ retrieved docs.
3. Generate response on augmented input

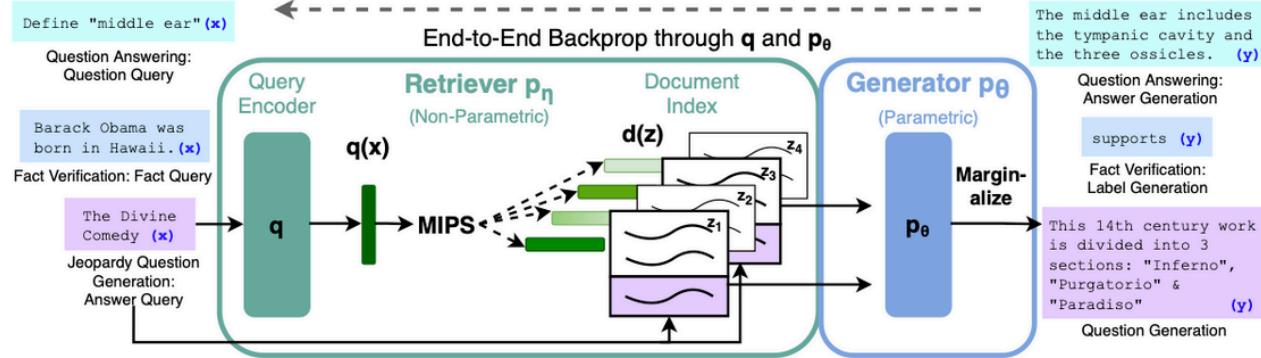
- we need:

- LM
- retrieval model (turn query into embeddings)

- DB w/ representations of docs - to search on query with

# Early RAG approaches: fine-tune pretrained LM and query encoder: (doc encoder frozen)

During fine-tuning, jointly train retriever (query encoder) and language model



Lewis et al '21 <https://arxiv.org/pdf/2005.11401>

- can change weights of LM OR  
query encoder via backprop