

Natural Language Processing

IR based Question Answering
Instructor: Yangqiu Song

Open-domain QA

- **Question answering** = build computer systems that **automatically** answer questions posed by humans in a **natural language**



- **Open-domain** = deal with questions about nearly anything, usually rely on *general ontologies* and *world knowledge*
 - *Q: Where does the energy in a nuclear explosion come from? A: high-speed nuclear reaction*
 - *Q: Where is Einstein's house? A: 112 Mercer St, Princeton, NJ*
 - *Q: How many papers were accepted by ACL*

Open-domain QA

Knowledge Bases



Structured

Tabl es

Category	Structure	Country	City	Height (metres)	Height (feet)
Mixed use	Burj Khalifa	United Arab Emirates	Dubai	829.8	2,722
Self-supporting tower	Tokyo Skytree	Japan	Tokyo	634	2,080
Mixed use	Shanghai Tower	China	Shanghai	632	2,073
Clock building	Abraj Al Bait Towers	Saudi Arabia	Mecca	601	1,972
Military structure	Large masts of INS Kattabomman	India	Tirunelveli	471	1,545
Mast radiator	Lualualei VLF transmitter	United States	Lualualei, Hawaii	458	1,503
Twin towers	Petronas Twin Towers	Malaysia	Kuala Lumpur	452	1,482
Residential	432 Park Avenue	United States	New York	425.5	1,396
Chimney	Ekibastuz GRES-2 Power Station	Kazakhstan	Ekibastuz	419.7	1,377
Radar	Dimona Radar Facility	Israel	Dimona	400	1,312
Lattice tower	Kiev TV Tower	Ukraine	Kiev	385	1,263
Electricity pylon	Zhoushan Island Overhead Powerline Tie	China	Zhoushan	370	1,214

Semi-
structured

Web Documents & Wikipedia



Unstructured

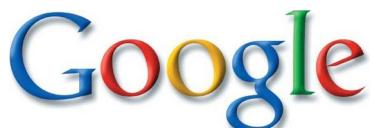
Answer Questions using Structured Data

- General problem setting
 - Information source: A “database”
 - Collection of records
 - Tables
 - Large scale DB with complex schema
 - Input: a natural language question (instead of formal “query”)
 - Output: Answer

Early Work

- Small scale & domain-specific KBs
 - Simple schema
 - Small number of entities and relations
 - Limited set of sensible questions
- Approaches
 - Ad-hoc methods (e.g., manually crafting rules) can be quite effective
 - Semantic parsing (of questions)
- Issues
 - Not clear if the methods are scalable
 - Cannot support “open-domain” question answering

Modern Large-scale Knowledge Graphs



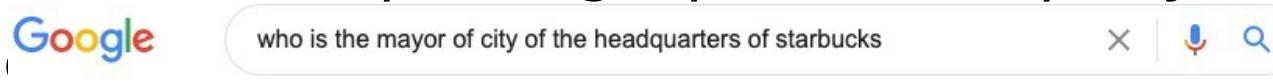
Knowledge
Graph



- Database that stores a large number of facts in an organized way
 - facts: 46m entities, 2.6b
 - facts: WikiData: 87m items
 - Most knowledge bases are curated

Open-Domain Question Answering using Knowledge Base

- Curated KB ensures the correctness of the information (answer)
- Common or "simple" questions can be answered easily
 - If semantic parsing (question->query) is done



-

The image shows a portion of a Google search results page. At the top, there's a navigation bar with 'All' selected, followed by 'News', 'Images', 'Videos', 'Maps', 'More', 'Settings', and 'Tools'. Below the bar, it says 'About 3,370,000 results (0.78 seconds)'. The first result is a card for 'Seattle / Mayor'. Inside the card, the name 'Jenny Durkan' is displayed in large text, with 'Since 2017' underneath. To the right of the text is a small portrait photo of a woman with blonde hair, smiling. The background of the card is white.

Key Challenges

- Language mismatch
 - Lots of ways to ask the same question
 - “Who played the role of Meg on Family Guy?”
 - “What is the name of the actress for Meg on Family Guy?”
 - “In the TV show Family Guy, who is the voice for Meg?”
 - Need to map questions to the predicates defined in KB
 - `tv.tv_program.regular_cast - tv.regular_tv_appearance.actor`
- Large search space
 - Some Freebase entities have $> 160,000$ immediate neighbors
- Compositionality
 - “What movies are directed by the person who won the most Academy and Golden Globe awards combined?”

Knowledge Graph is largely incomplete



Relation	Percentage unknown	
	<i>All 3M</i>	<i>Top 100K</i>
PROFESSION	68%	24%
PLACE OF BIRTH	71%	13%
NATIONALITY	75%	21%
EDUCATION	91%	63%
SPOUSES	92%	68%
PARENTS	94%	77%
CHILDREN	94%	80%
SIBLINGS	96%	83%
ETHNICITY	99%	86%

Search engines: from keyword matching to question answering with unstructured data

Google X |

All News Shopping Videos Images More Settings Tools

About 13,100,000 results (0.41 seconds)

779 Accepted Papers

ACL 2020 Announces Its 779 Accepted Papers | Synced. May 20, 2020

syncedreview.com › 2020/05/20 › acl-2020-announces... ▾

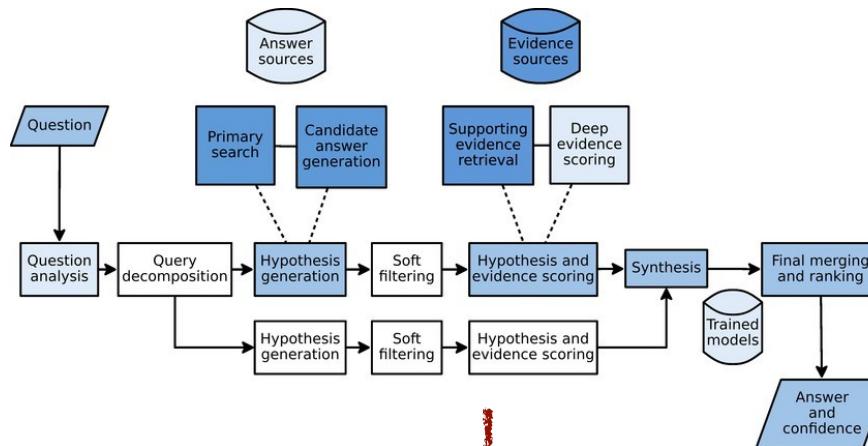
[ACL 2020 Announces Its 779 Accepted Papers | Synced](#)



[About Featured Snippets](#) [Feedback](#)

Paradigm Switching

Classical QA pipeline



Q: How many of Warsaw's inhabitants spoke Polish in 1933?

Two-stage
Retriever-
reader
approaches

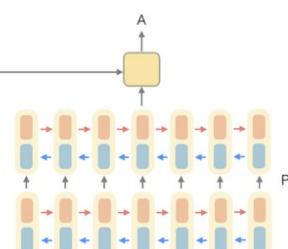


Document
Retriever



Document
Reader

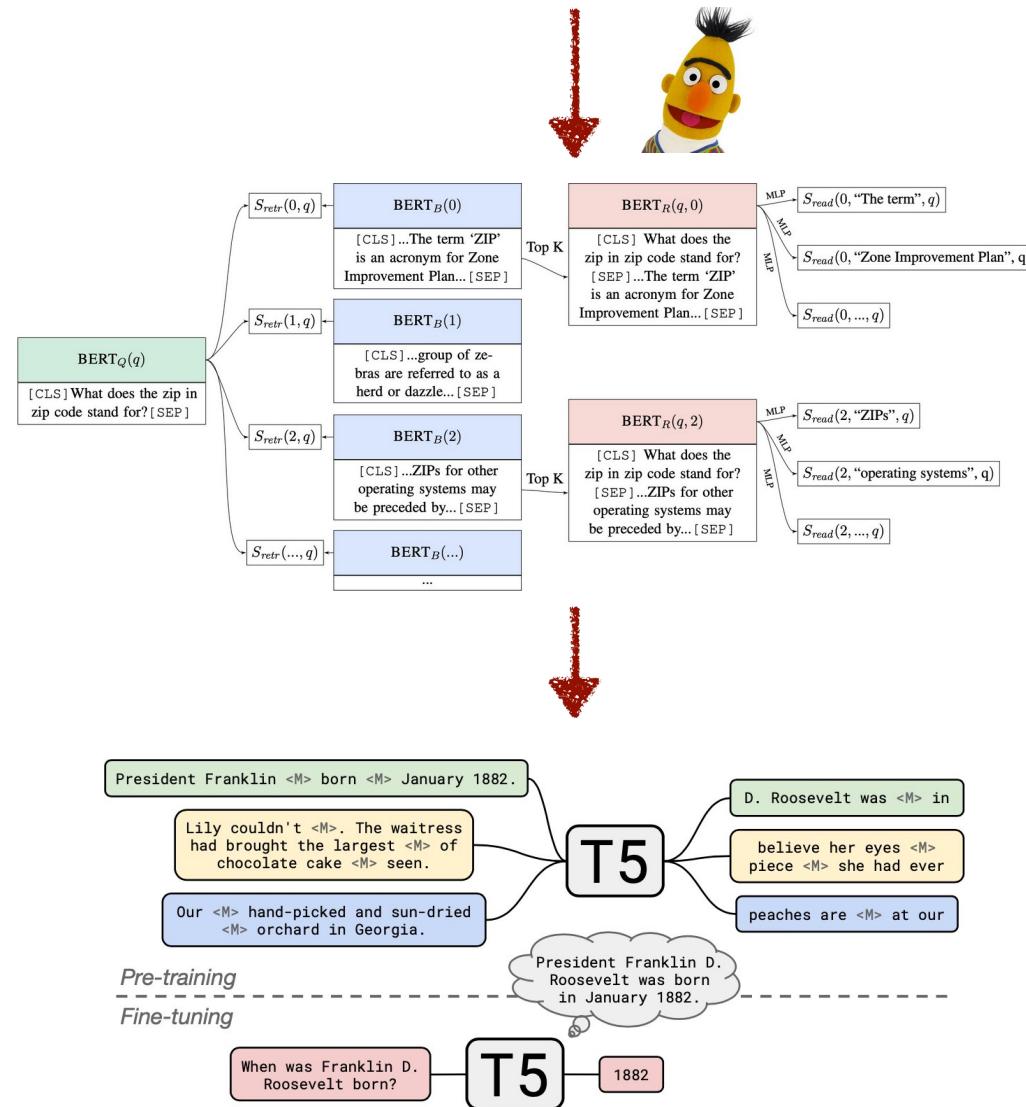
833,500



Paradigm Switching

End-to-end learning

Retrieval-free models



Outline

- History

Text Retrieval Conference (TREC)

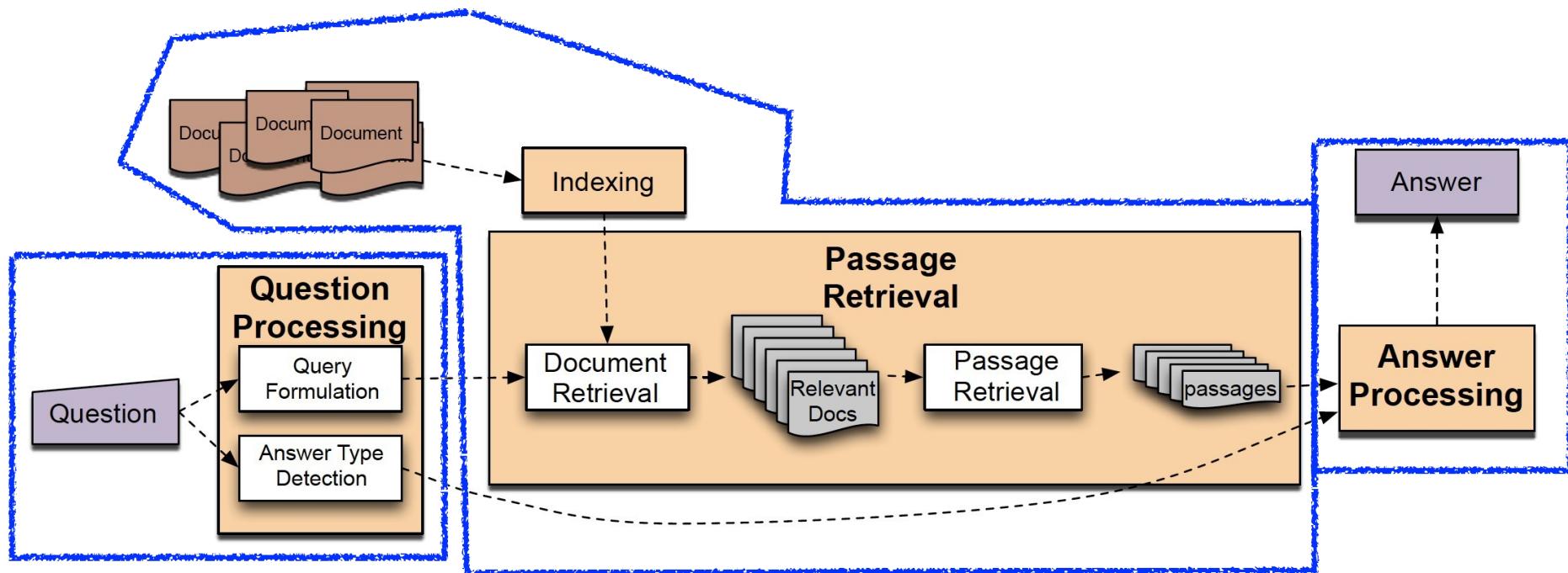
QA Tracks (1999 - 2007)

- Originating from the IR community as the next version of search
 - Relevant documents -> short answer with support
- Shared tasks & competitions
 - Corpus: newswire (AP, WSJ, LA Times, etc.); 979k articles, 3GB
 - Test set: 500 questions from Excite, Encarta, MSNSearch, AskJeeves
 - Human judges decide the correctness of the answers from QA systems

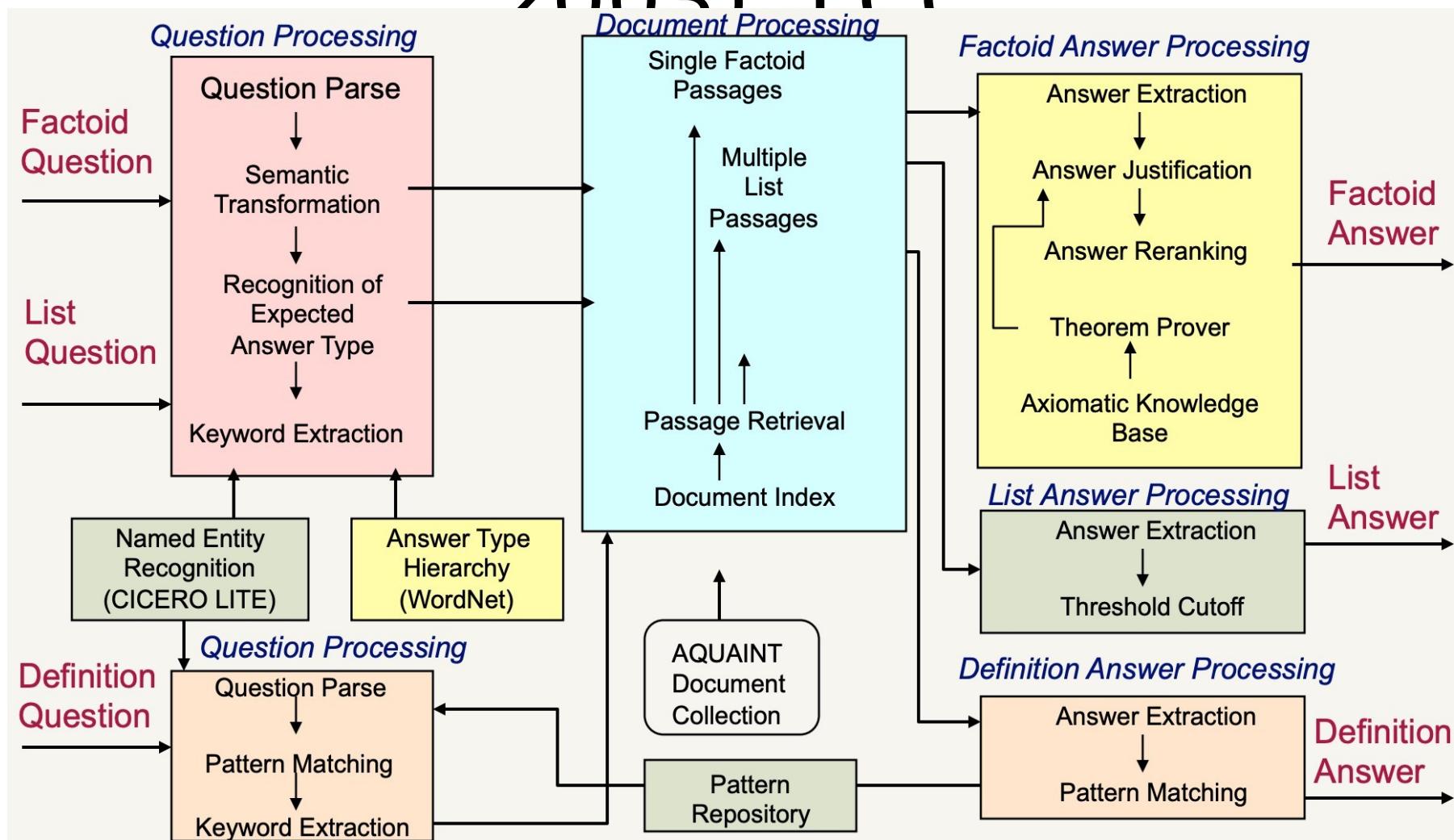


[https://
trec.nist.gov/](https://trec.nist.gov/)

Typical pipeline of TREC-QA systems



Top TREC-QA system (circa 2003) · ICC



Factoid Questions: Sometimes the question asks for a body of information instead of a fact.
abagiu and Moldovan, 2003. Question Answering.

Recent developments 2013+

- Trend: Macro-reading -> Micro-reading
- General problem setting
 - Given a question and a "context", answer the question using the "context"
 - Two different goals
 - Test machine's intelligence AI (machine reading comprehension)
 - Fulfill user's information need (answer extraction/processing stage)
- Research directions guided by development of new tasks/datasets
- Rapid progress made by new deep learning models

Stanford Question Answering Dataset (SQuAD)

[Rajpurkar et al., 2016]

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called "showers".

What causes precipitation to fall?
gravity

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?
graupel

Where do water droplets collide with ice crystals to form precipitation?
within a cloud

- (passage, question, answer) triples
- Passage is from Wikipedia, question is crowd-sourced
- Answer must be a span of text in the passage (aka. "**extractive question answering**")
- SQuAD 1.1: 100k answerable questions, SQuAD 2.0: another 50k unanswerable questions

Document Reader

$$P_{start}(i) \propto \exp(\mathbf{p}_i \mathbf{W}_s \mathbf{q})$$

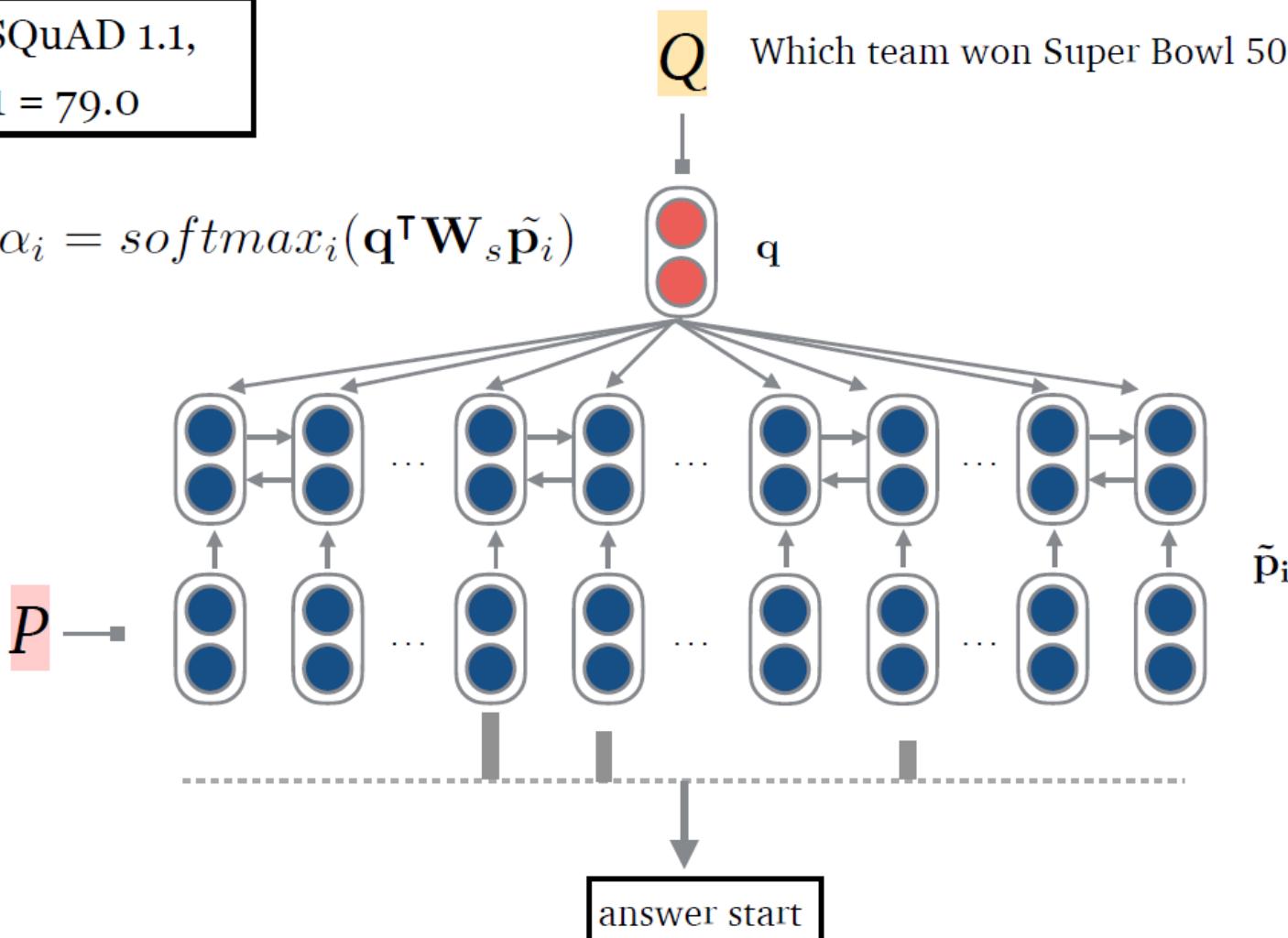
$$P_{end}(i) \propto \exp(\mathbf{p}_i \mathbf{W}_e \mathbf{q})$$

$P_{start}(i) \times P_{end}(i')$ is maximized

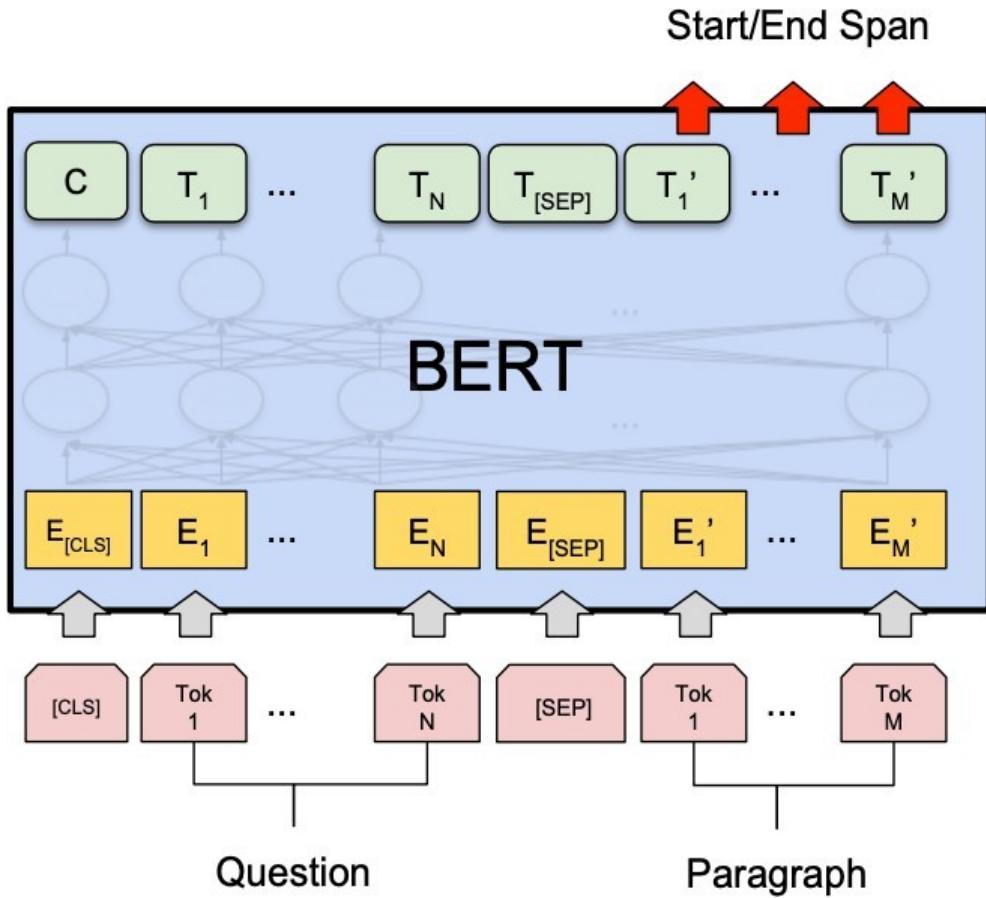
On SQuAD 1.1,

- F1 = 79.0

$$\alpha_i = softmax_i(\mathbf{q}^\top \mathbf{W}_s \tilde{\mathbf{p}}_i)$$



Document Reader



Question = Segment A

Passage = Segment B

Answer = predicting two endpoints in segment B

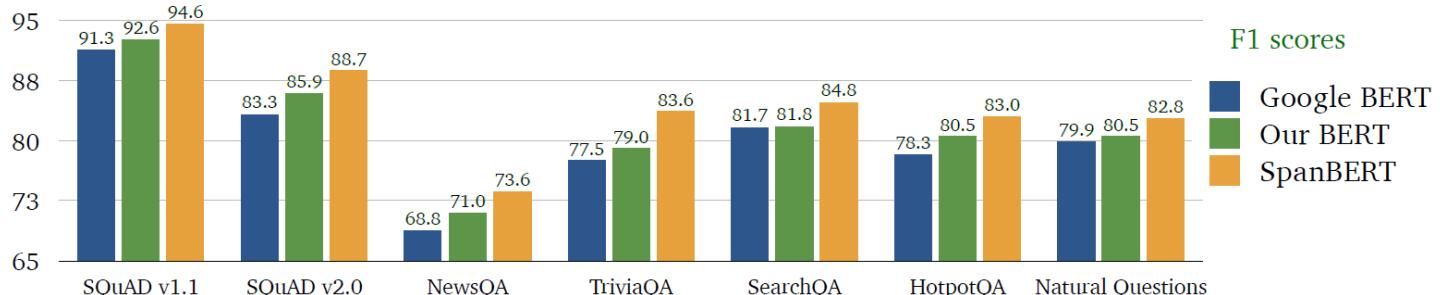
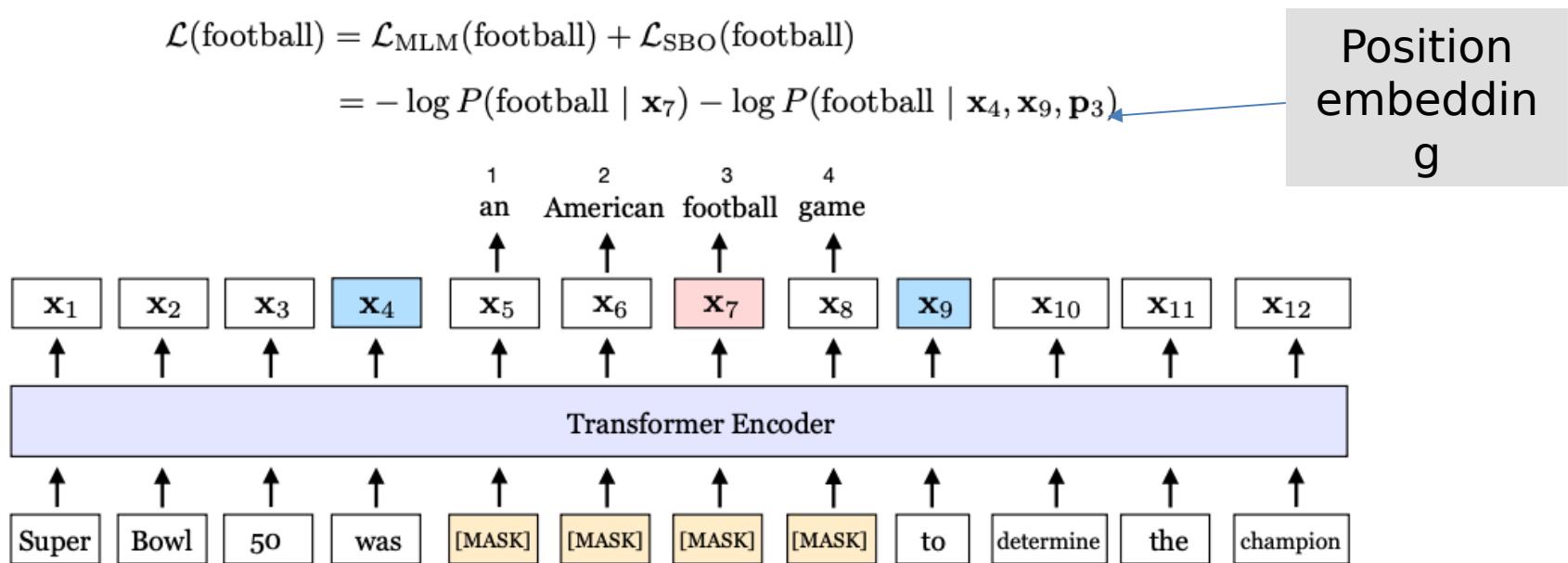
On SQuAD

1.1, : F1 =

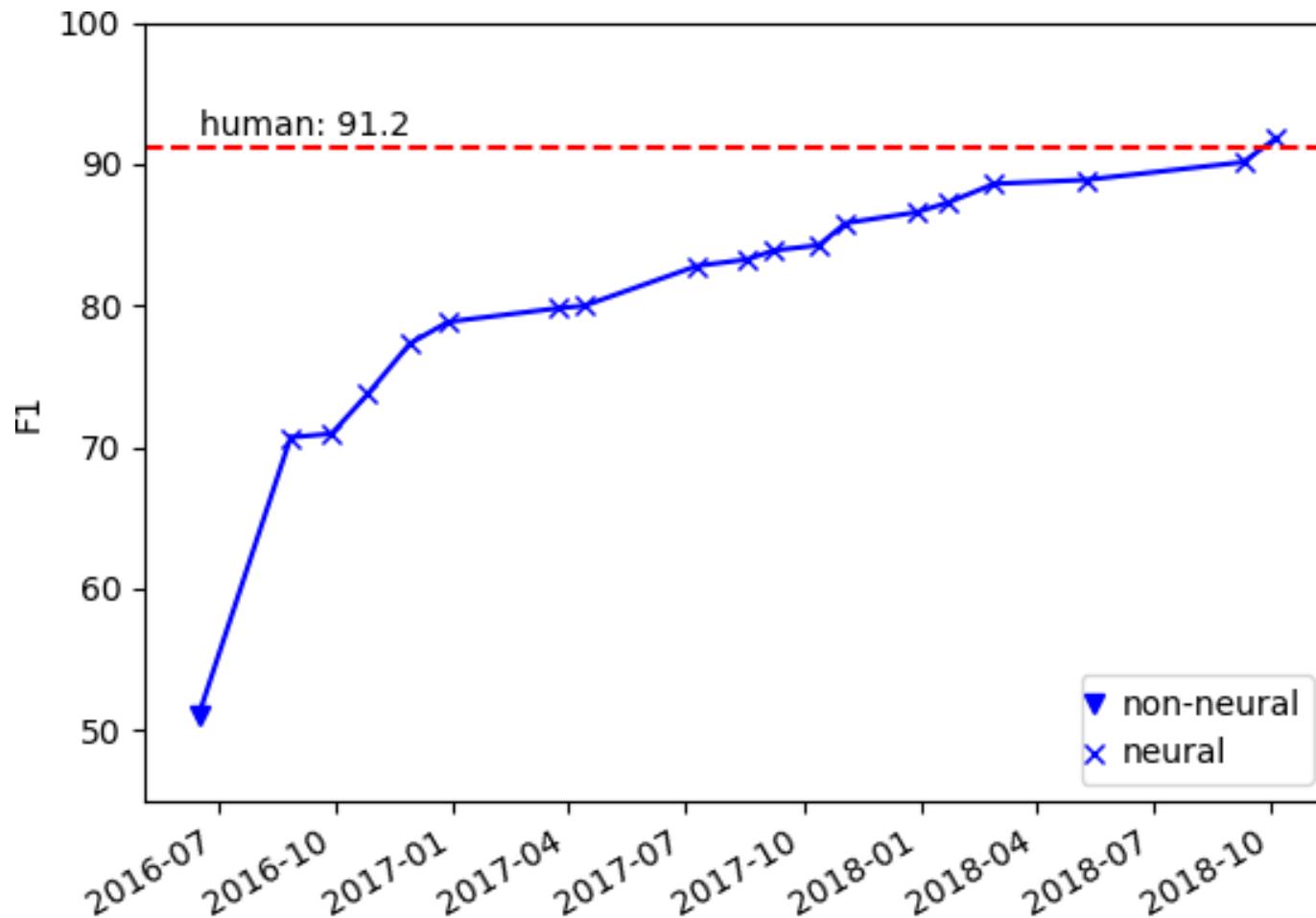
• RoBERTa: F1 = 90.9
94.6

SpanBERT

An illustration of SpanBERT training. The span an American football game is masked. The SBO (span-boundary objective) uses the output representations of the boundary tokens



Stanford Question Answering Dataset (SQuAD)



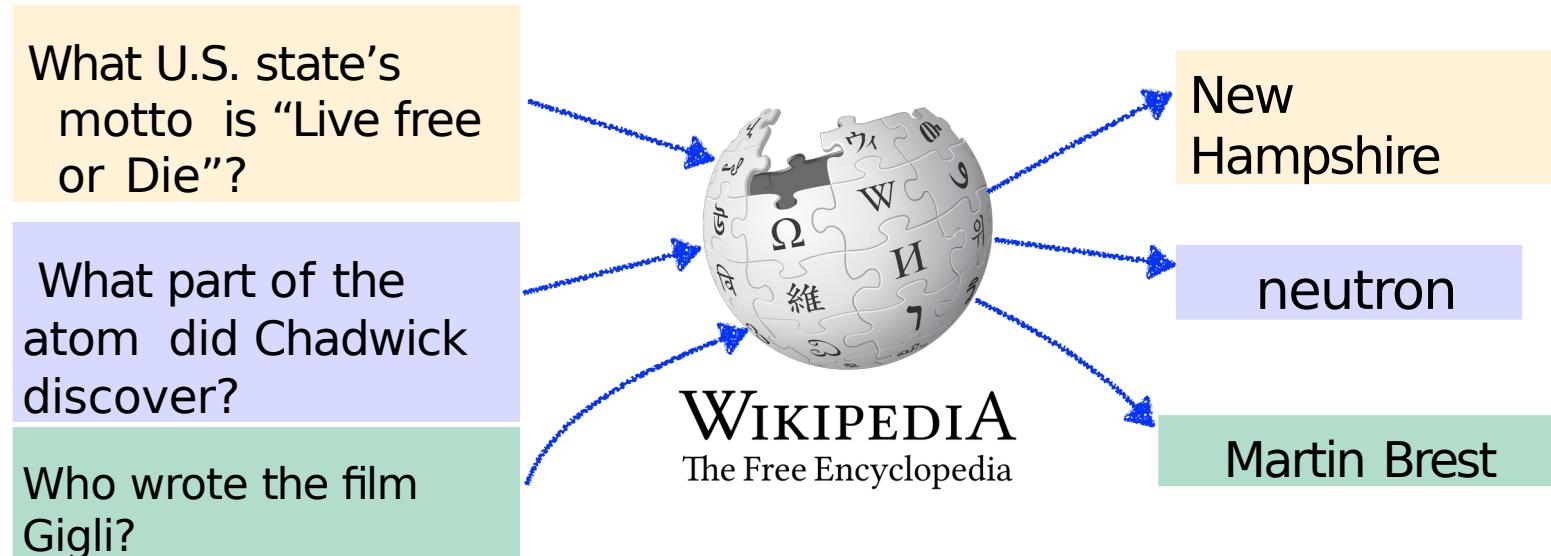
Rajpurkar et al., 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text

Outline

- Two-stage **retriever-reader** approaches

Problem setup

- **Input:** question Q , D = English Wikipedia (~ 5 million documents)
- **Output:** answer A



“Machine Reading at Scale”

Why is Wikipedia?



- We treat Wikipedia as a **generic** collection of articles and don't consider its internal **graph structure** in this setting easy to extend to any collection of documents.
- The search problem is challenging and realistic while its scale is still **manageable**, especially **for academic research**.
- Wikipedia contains a wealth of information of real-world facts. We don't consider the **redundancy** problem too much here (vs. Web documents or news corpora).

DrQA: a first neural open-domain QA system

[Chen et al., 2017]

Q: How many of Warsaw's inhabitants spoke Polish in 1933?



WIKIPEDIA

Document
Retriever



Information
Retrieval

This article is about the Polish capital. For other uses, see Warsaw (disambiguation).
"Warszawa" redirects here. For other uses, see Warsaw (disambiguation).
"City of Warsaw" redirects here. For the Second World War fighter squadron, see No. 316 Polish Fighter Squadron. F 1934, see Adamowicz brothers.

Warsaw (Polish: Warszawa [varˈʂava] ();[1] listen);[2] also see other names) is the capital and largest city of Poland. It stands on the Vistula River in central Poland, roughly 260 kilometres (160 mi) from the Baltic Sea and 300 kilometres (190 mi) from the Carpathian Mountains. Its population is estimated at 1,750 million residents within a greater metropolitan area of 3.105 million residents, which makes Warsaw the 8th most populous capital city in the European Union.[3][4] The city limits cover 151.9 square kilometres (59.0 sq mi), while the metropolitan area covers 6,100.43 square kilometres (2,355.39 sq mi).[5]

In 2012 the Economist Intelligence Unit ranked Warsaw as the 32nd most liveable city in the world.[6] It was also ranked as one of the most liveable cities in Central Europe. Today Warsaw is considered an "Alpha+" global city, a major international business destination and a significant cultural, political and economic hub.[7][8] Warsaw's economy, by a wide margin, is the largest in Central Europe. It includes a diversified industrial base, steel, automotive, manufacturing and food processing. The city is a significant centre of research and development, IPOS, ITC, as well as of the Polish media industry. The Warsaw Stock Exchange is one of the largest and most important in Central and Eastern Europe.[9][10] Frontex, the European Union agency for external border security, has its headquarters in Warsaw. It has been said that Warsaw, together with Frankfurt, London, Paris and Barcelona is one of the cities with the highest number of skyscrapers in the European Union.[11] Warsaw has also been called "Eastern Europe's chic cultural capital with thriving art and club scenes and serious restaurants".[12]

Document Reader → 833,500



Reading
Comprehension

Document Retriever

A TF-IDF weighted term vector model over unigrams/bigrams:

$\text{tf} = \text{term frequency}$, $\text{idf} = \text{inverse document frequency}$

t : term (uni/bi), d : document (= one Wiki. article), D : corpus (= Wikipedia)

$$\text{tf-idf}(t, d, D) = \text{tf}(t, d) \times \text{idf}(t, D)$$

$$\text{tf}(t, d) = \log(1 + \text{freq}(t, d))$$

$$\text{idf}(t, D) = \log\left(\frac{|D|}{|\{d \in D : t \in d\}|}\right)$$

- This retriever is not *trainable*.
- Retriever at document level instead of paragraph level.

Document Reader

paragraph, document,
arbitrary-length text blocks.

Cast as a *reading comprehension* problem:

- **Input** is a passage P and a question Q
- **Output** is an answer A

A restricted setting is that A needs
to be a segment of text in $P \Rightarrow$
“**extractive question answering**”

Stanford Question Answering
Dataset [Rajpurkar et al., 2016]

Super Bowl 50 was an American football game to determine the champion of the National Football League (NFL) for the 2015 season. The American Football Conference (AFC) champion Denver Broncos defeated the National Football Conference (NFC) champion Carolina Panthers 24–10 to earn their third Super Bowl title. The game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at Santa Clara, California.

Question: Which NFL team won Super Bowl 50?

Answer: Denver Broncos

Question: What does AFC stand for?

Answer: American Football Conference

Question: What year was Super Bowl 50?

Answer: 2016

How to train the model?

- Using an existing reading comprehension dataset (e.g., SQuAD)!

$$D_{rc} = \{(P_i, Q_i, A_i)\}$$

Problem: very different distribution with real-world QA data.

- How about other QA datasets (e.g., WebQuestions, TREC)?

$$D_{QA} = \{(Q_i, A_i)\}$$

- Solution: create **distantly-supervised** examples using our **retriever**!

$$(Q, A) \longrightarrow (P, Q,$$

A)
if passage is retrieved and answer can be found in
passage

Similar to distant supervision

in information extraction [Mintz et al., 2009]

Chen et al., 2017. Reading Wikipedia to Answer Open-domain Questions

How to train the model?

Question: What U.S. state's motto is "Live free or Die"?

Answer: New Hampshire
Passage

Live Free or Die

From Wikipedia, the free encyclopedia

"**Live Free or Die**" is the official motto of the U.S. state of **New Hampshire**, adopted by the state in 1945.^[1] It is possibly the best-known of all [state mottos](#), partly because it conveys an assertive [independence](#) historically found in [American political philosophy](#) and partly because of its contrast to the milder sentiments found in other state mottos.

Putting together

Training

time:

Document retriever: not trained

- Document reader: a LSTM-based neural reading comprehension model trained on SQuAD + distantly-supervised data generated from QA datasets

Inference

time:

The document retriever returns *top 5 documents*

- The reader reads every (natural) paragraph in these 5 documents and predicts an answer and its span score.

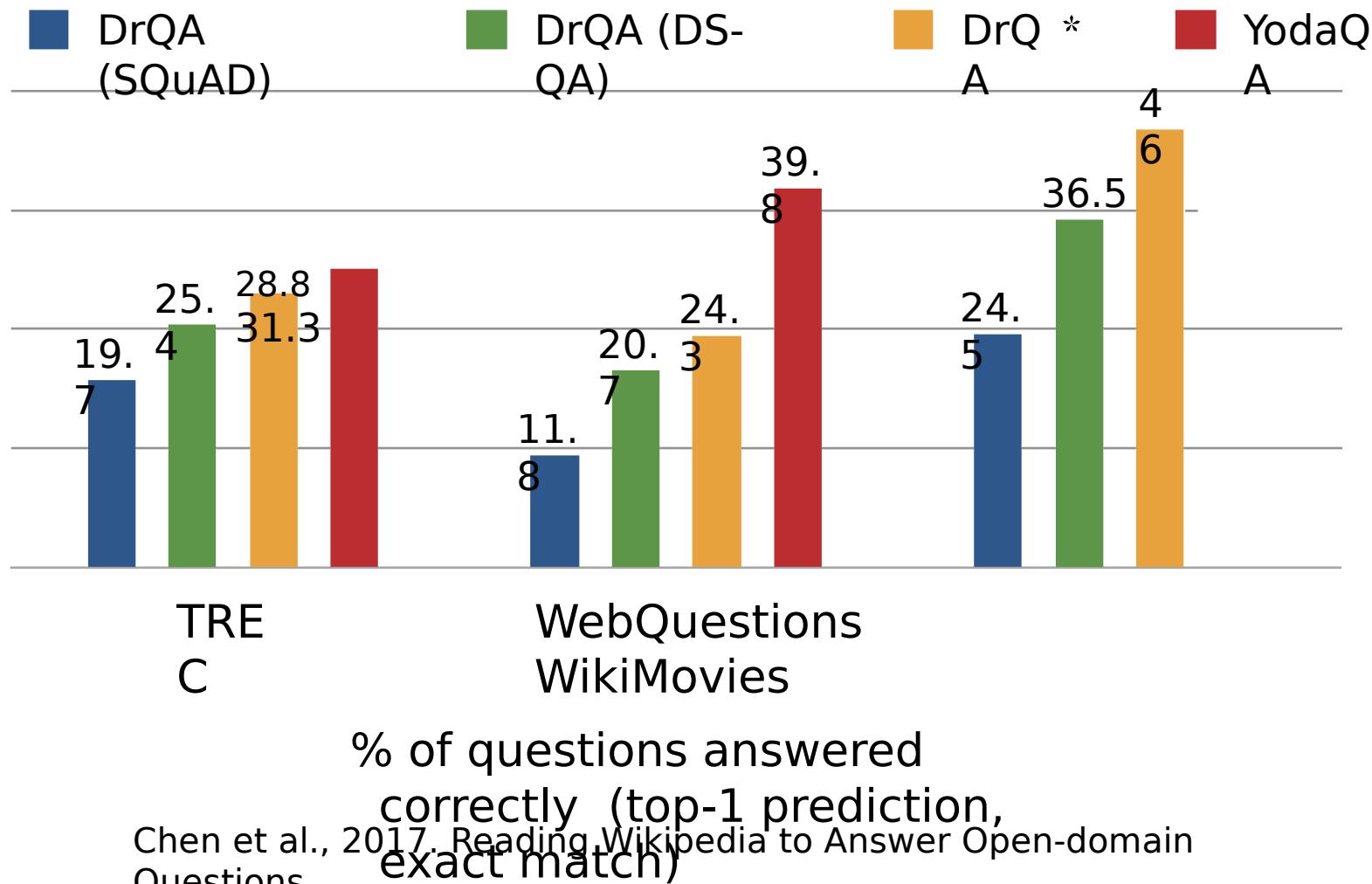
The system finally returns the answer with the highest (unnormalized) span score.

Experiments

See [Raison et al., 2018]

Trained reader on full documents

Hybrid system on KBs, text, ontologies etc
<http://ailao.eu/yodaqa/>



How can we do better?

- DrQA considers retrieval at document level.
Does paragraph-level retriever work better?
- Answers in the retrieved passages might not be directly comparable at inference time.
Does multi-passage training help?
- The importance of each passage has been omitted.
Can we use passage retriever scores or even train a better ranker?
- The retriever is not trained!
The focus of the next part.

Outline

- **Dense** Retriever and **End-to-end** Training

Key questions

Q: How many of Warsaw's inhabitants spoke Polish in 1933?



WIKIPEDIA

Document
Retriever



Article | Talk | Read | Edit | View history | 5s

Warsaw

From Wikipedia, the free encyclopedia

This article is about the Polish capital. For other uses, see Warsaw (disambiguation).
"Warszawa" redirects here. For other uses, see Warszawa (disambiguation).
"City of Warsaw" redirects here. For the Second World War fighter squadron, see No. 316 Polish Fighter Squadron. F 1934, see Adamowicz brothers.

Warsaw (Polish: Warszawa [varˈʂava] ();[1] ;[2] see also other names) is the capital and largest city of Poland. It stands on the Vistula River in central Poland, roughly 260 kilometres (160 mi) from the Baltic Sea and 300 kilometres (190 mi) from the Carpathian Mountains. Its population is estimated at 1,750 million residents within a greater metropolitan area of 3.105 million residents, which makes Warsaw the 8th most populous capital city in the European Union.[3][4] The city limits cover 151.9 square kilometres (59.0 sq mi), while the metropolitan area covers 6,100.43 square kilometres (2,355.39 sq mi).[5]

In 2012 the Economist Intelligence Unit ranked Warsaw as the 32nd most liveable city in the world.[6] It was also ranked as one of the most liveable cities in Central Europe. Today Warsaw is considered an "Alpha+" global city, a major international tourist destination and a significant cultural, political and economic hub.[7][8] Warsaw's economy, by a wide margin, is the largest in Central Europe, and includes a diversified industrial base, particularly in automotive manufacturing and food processing. The city is a significant centre of research and development, IPOS, ITCs, as well as of the Polish media industry. The Warsaw Stock Exchange is one of the largest and most important in Central and Eastern Europe.[9] Frontex, the European Union agency for external border security, has its headquarters in Warsaw. It has been said that Warsaw, together with Frankfurt, London, Paris and Barcelona is one of the cities with the highest number of skyscrapers in the European Union.[10] Warsaw has also been called "Eastern Europe's chic cultural capital with thriving art and club scenes and serious restaurants".[11]

Document
Reader

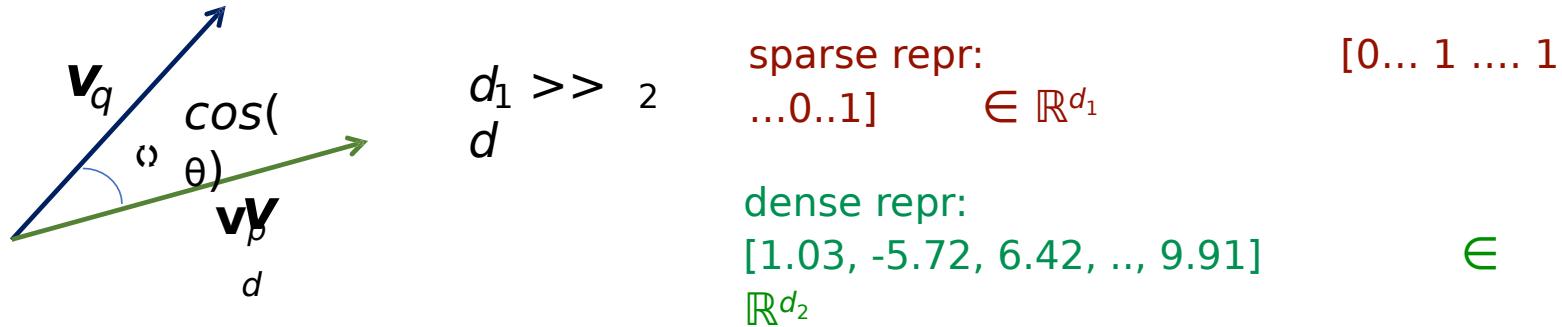
833,500



Can we also train the retriever component?

Can we use dense representations for

Sparse vs dense representations for retrieval



They capture complementary information; Dense representations have never been shown to outperform sparse representations in open-domain QA before 2019!

Why dense retrieval now?

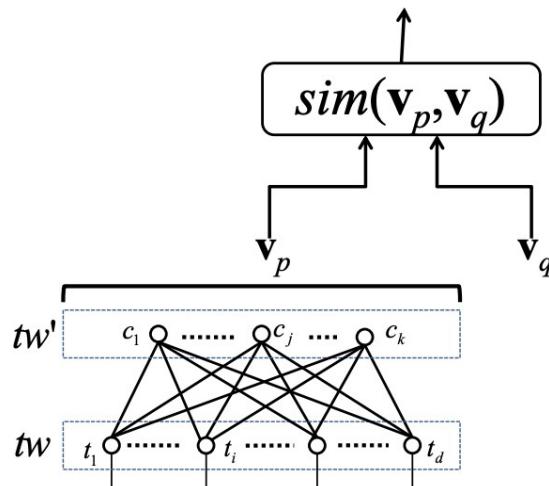
It is a difficult problem: we need to *encode*, *index* and *search* from **5M** documents or **30M** paragraphs or **60B** phrases.

Learning Discriminative Projections for Text Similarity Measures

Wen-tau Yih Kristina Toutanova John C. Platt Christopher Meek

Microsoft Research
One Microsoft Way
Redmond, WA 98052, USA

CoNLL'20
11
Best
paper



- Cross-lingual document retrieval
- Ad relevance prediction
- Web search ranking

Why dense retrieval now?

It is actually not easy to make these dense models “work”.

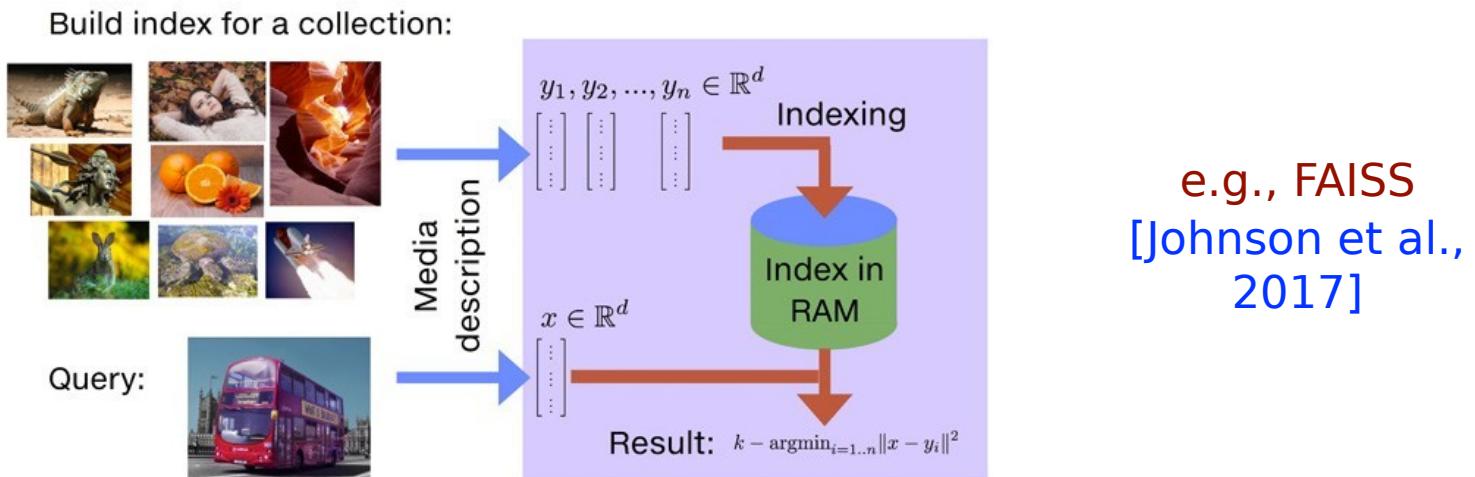
- Needs large enough labeled data (e.g., 82M query-doc pairs from user clicks).

We have pre-trained
models now!



Now we have much better techniques and tools to support fast maximum inner product search (MIPS):

- In-memory data structure and indexing schemes



ORQA: Open-Retriever Question Answering

[Lee et al., 2019]

First model to learn retriever and reader jointly

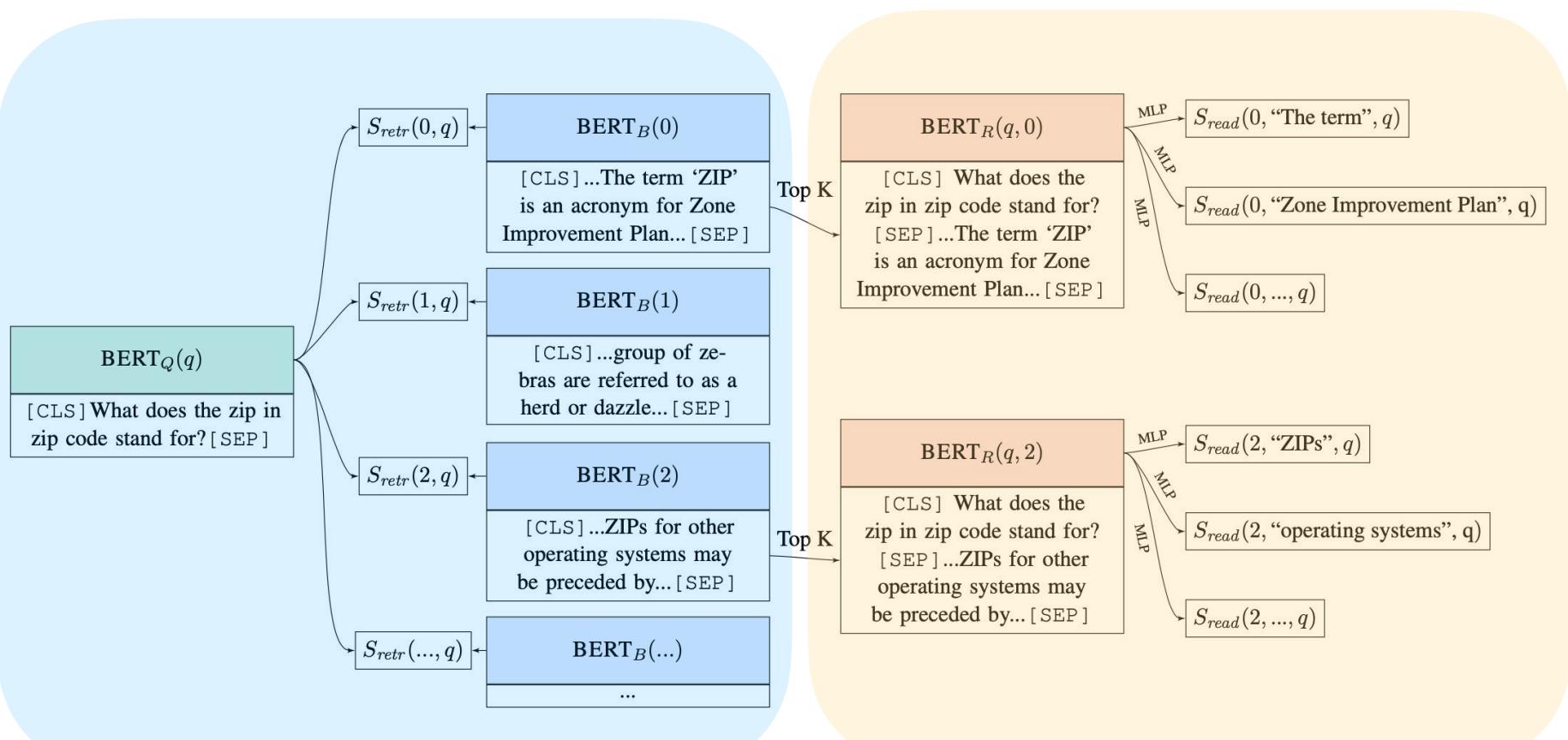
Key contributions:

- Both retriever and reader are learnable with NNs (= BERT)
- Only learned from question-answer pairs; No reading comprehension datasets!

$$D_{QA} = \{(Q_i, A_i)\}$$

- A new pre-training task called **Inverse Cloze Task (ICT)** to address the challenging retrieval problem.

ORQA: Open-Retriever Question Answering



Information
Retrieval

Reading
Comprehension

ORQA: Overview

$$S(b, s, q)$$

Notations:

- b - text block,
- s - a span of text within b ,
- q - question

Fixed-length blocks as “passages”: 288 wordpieces => 13M
in total

Each block has 2000 possible answer spans

Modeling

$$S(b, s, q) = S_{retr}(b, q) + S_{read}(b, s, q)$$

Inference

$$a^* = \text{TEXT}(\underset{b,s}{\operatorname{argmax}} S(b, s, q))$$

Retriever score:

Question

What does the zip in zip code stand for?



$BERT_Q$



oooooooooooo W_q



$$S_{\text{retr}}(b, q) = \frac{h_q^\top h_b}{b}$$

$$S_{\text{retr}}(b, q) = h_q^\top h_b$$

$$h_q = \mathbf{W}_q \text{BERT}_Q(q)[\text{CLS}]$$

$$h_b = \mathbf{W}_b \text{BERT}_B(b)[\text{CLS}]$$

$$S_{\text{retr}}(b, q) = h_q^\top h_b$$

All of Wikipedia: select top K

Each evidence block

$BERT_B$



oooooooooooo W

b



$$\text{Evidence Block 1 } S_{\text{retr}}(b_1, q)$$

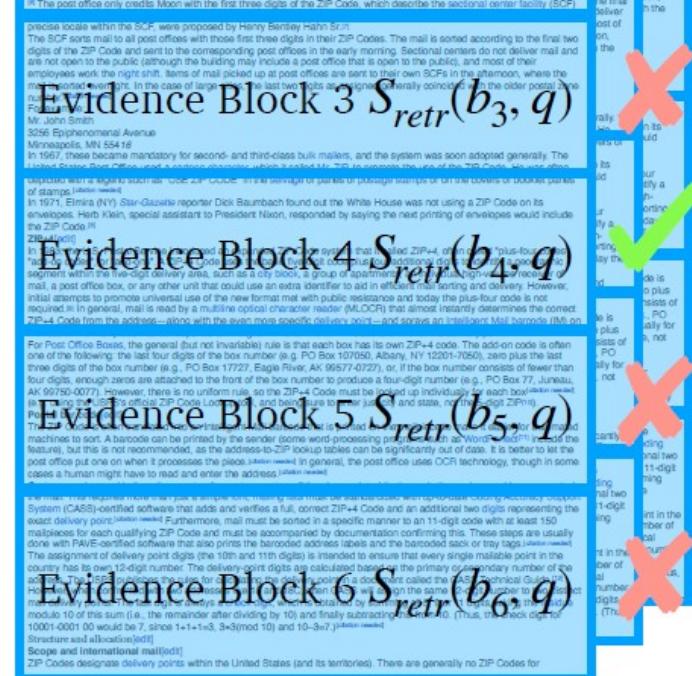
$$\text{Evidence Block 2 } S_{\text{retr}}(b_2, q)$$

$$\text{Evidence Block 3 } S_{\text{retr}}(b_3, q)$$

$$\text{Evidence Block 4 } S_{\text{retr}}(b_4, q)$$

$$\text{Evidence Block 5 } S_{\text{retr}}(b_5, q)$$

$$\text{Evidence Block 6 } S_{\text{retr}}(b_6, q)$$



Reader score:

$$S_{read}(b, s, q)$$

Top-K evidence
blocks concatenated
with question

Question

Top Evidence Block 1

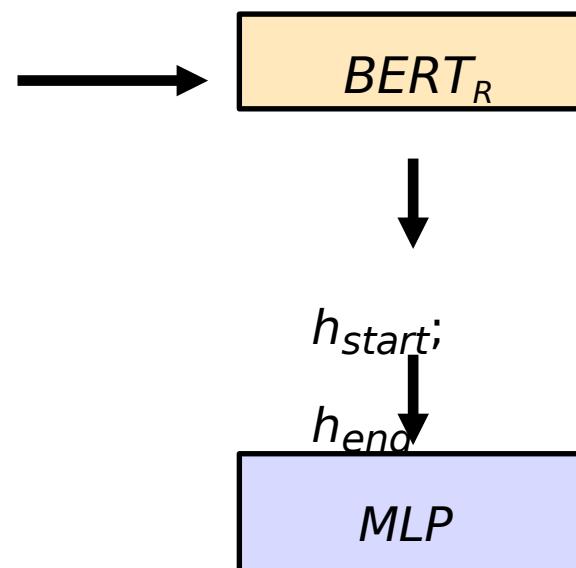
Question

Top Evidence Block 2

$$\begin{aligned} h_{start} &= \text{BERT}_R(q, b)[\text{START}(s)] \\ h_{end} &= \text{BERT}_R(q, b)[\text{END}(s)] \\ S_{read}(b, s, q) &= \text{MLP}([h_{start}; h_{end}]) \end{aligned}$$

Notations:

- b - text block,
- s - a span of text within b,
- q - question



$$S_{read}(b, s, q)$$

How is this model learned?

$$P(b, s|q) = \frac{\exp(S(b, s, q))}{\sum_{b' \in \text{TOP}(k)} \sum_{s' \in b'} \exp(S(b', s', q))}$$

Notations:

- b - text block,
- s - a span of text within b ,
- q - question

$k=5$, $\text{TOP}(k)$ = the top k retrieved blocks according to
 $S_{\text{retr}}(b, q)$

Loss function

$$L_{\text{full}}(q, a) = -\log \sum_{b \in \text{TOP}(k)} \sum_{s \in b, a = \text{TEXT}(s)} P(b, s | q)$$

How is this model learned?

Early learning: consider a larger set of evidence blocks

$c =$
5000

$$P_{\text{early}}(b|q) = \frac{\exp(S_{\text{retr}}(b, q))}{\sum_{b' \in \text{TOP}(c)} \exp(S_{\text{retr}}(b', q))}$$

$$L_{\text{early}}(q, a) = -\log \sum_{b \in \text{TOP}(c), a \in \text{TEXT}(b)} P_{\text{early}}(b|q)$$

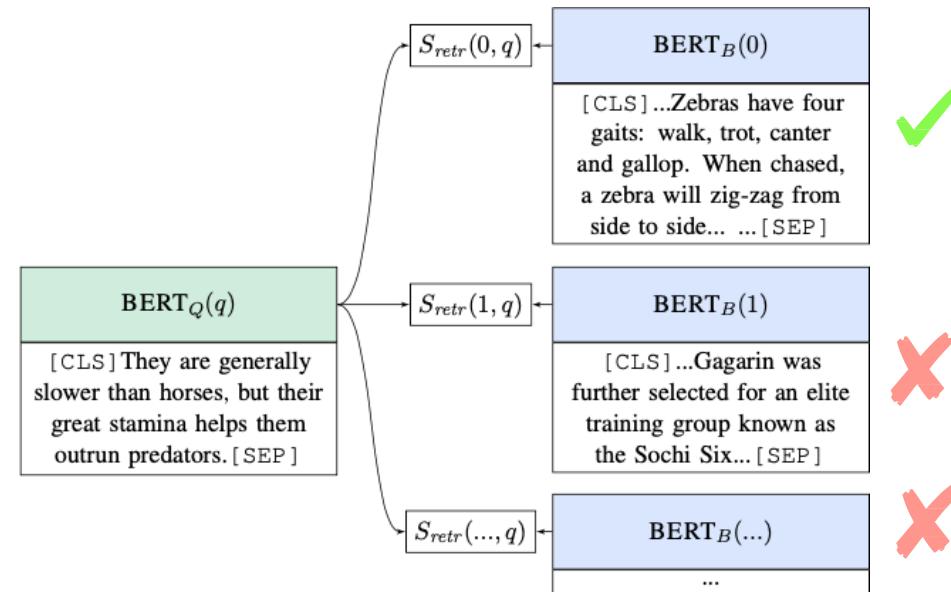
Final loss

$$L(q, a) = L_{\text{early}}(q, a) + L_{\text{full}}(q, a)$$

Pre-training: Inverse Cloze Task (ICT)

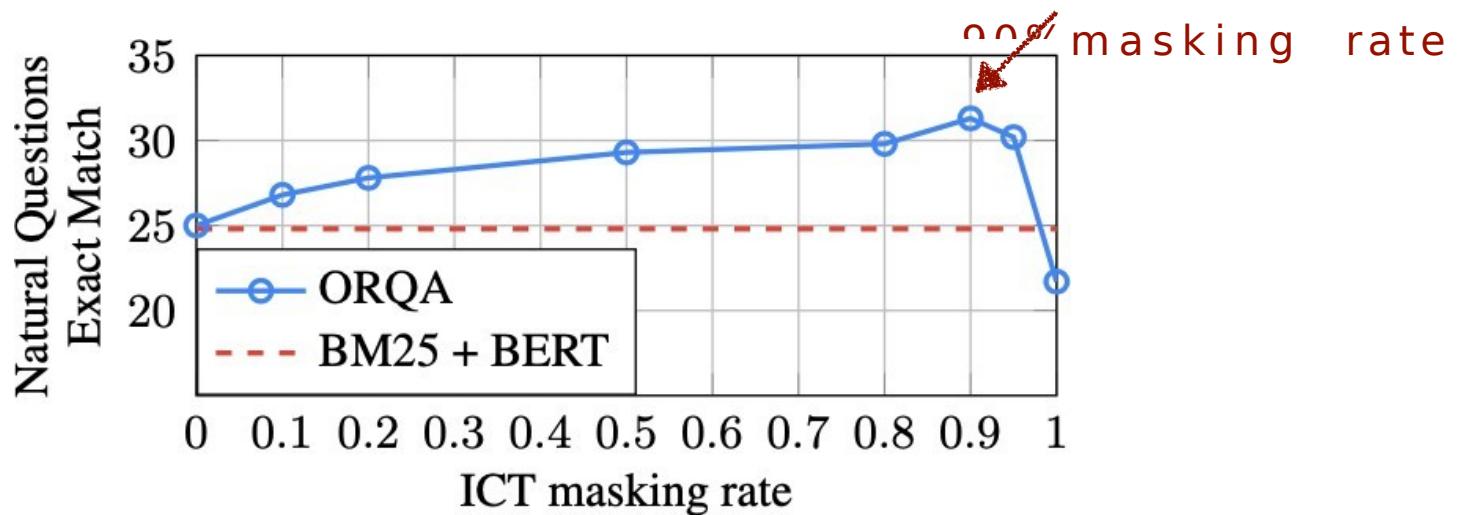
Key idea: a sentence is treated as a *pseudo-question*, and its context is treated as *pseudo-evidence*. The goal is to predict the correct context among a set of other random options.

...Zebras have four gaits: walk, trot, canter and gallop. **They are generally slower than horses, but their great stamina helps them outrun predators.** When chased, a zebra will zig-zag from side to side...



Pre-training: Inverse Cloze Task (ICT)

Still encourages the model to learn word matching — only remove sentences from its context in **90%** of the examples!



Important - After pre-training, **BERT_B is fixed** and all the block representations can be pre-computed and search efficiently using existing tools (e.g., Locality Sensitive Hashing).

REALM: Retrieval-augmented Language Model

[Guu et al., 2020]

Pre-training on both retriever and reader!

$$p(y | x) = \sum_{z \in \mathcal{Z}} p(y | z, x) p(z | x)$$

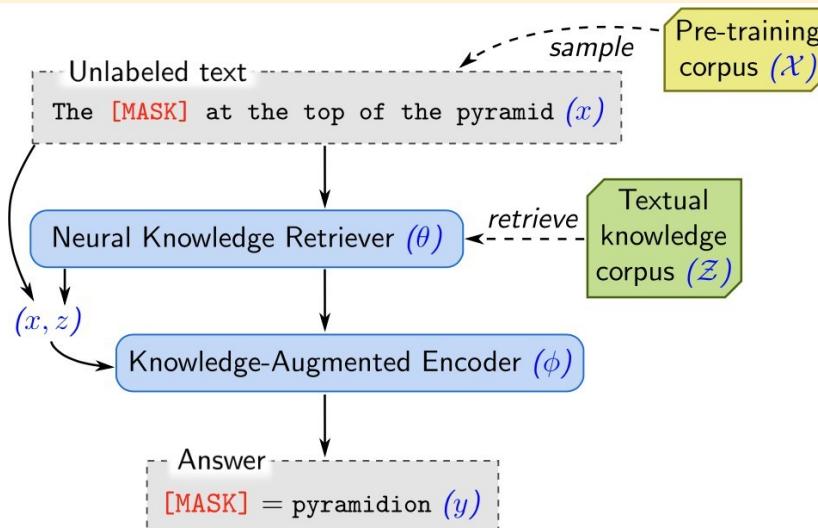
$Z = \text{Top}(K)$

Compared to

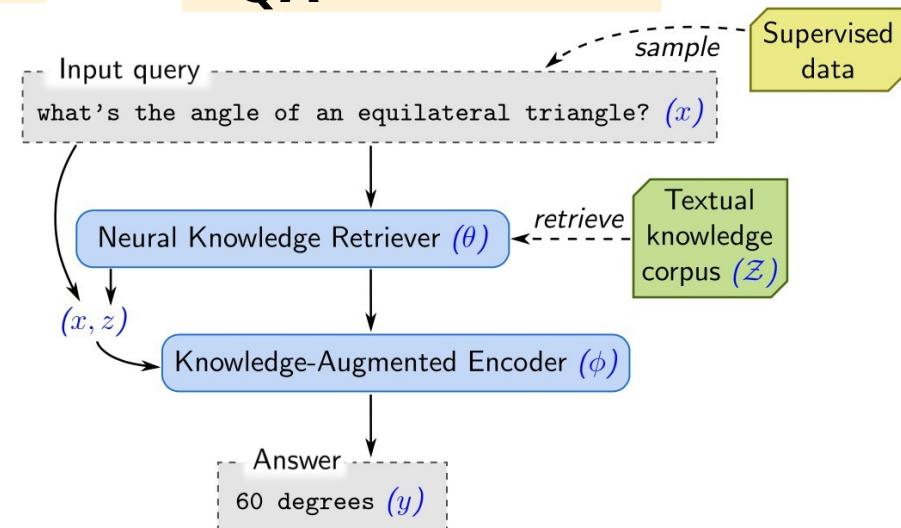
$$S(b, s, q) = S_{\text{retr}}(b, q) + S_{\text{read}}(b, s, q)$$

ORQA:

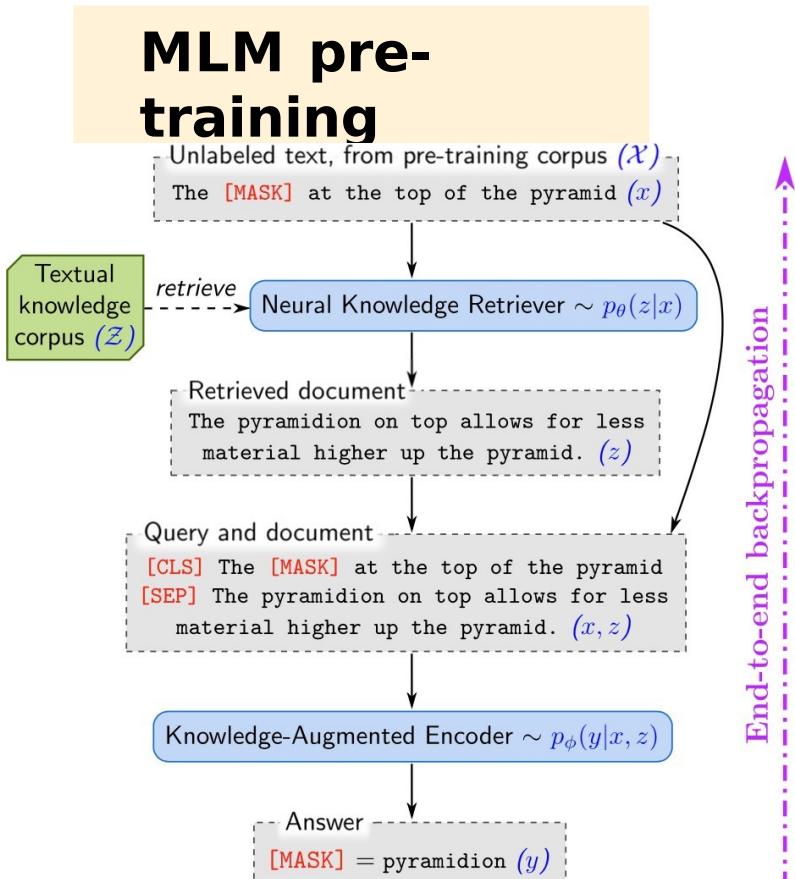
Pre-training: masked language model (MLM)



Fine-tuning:
QA



REALM: MLM pre-training



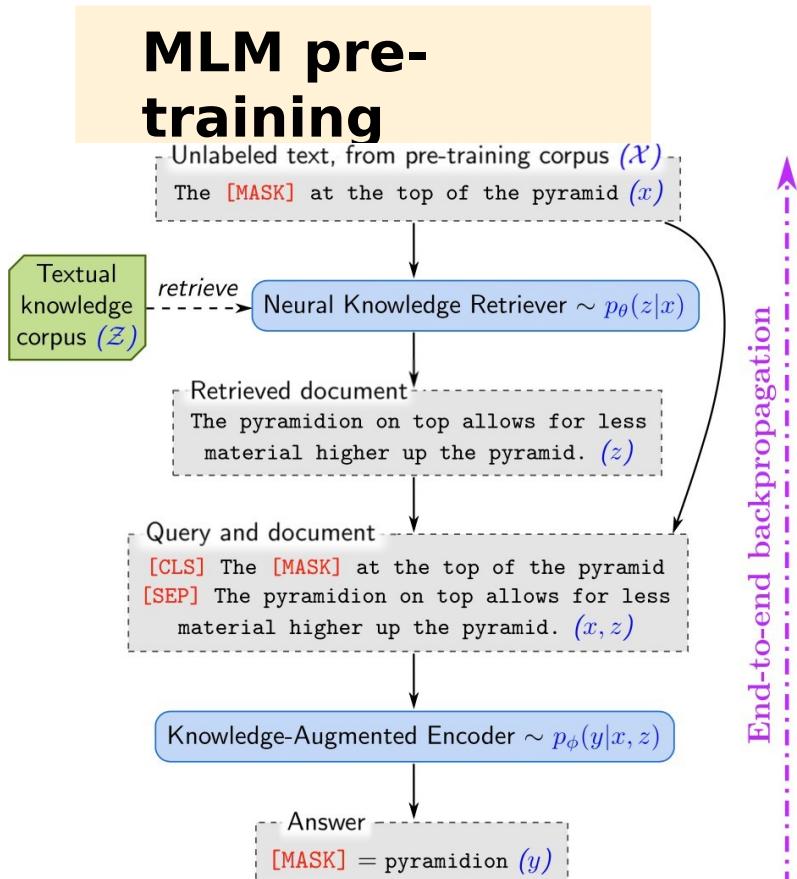
- Cold-start is hard => Use ICT as a first-stage pre-training

- **Salient span masking:** similar to span masking [Joshi et al., 2020] but only mask named **Evaluation** data.

NQ

Random **uniform** masks:
32.3 Random **span** masks: 35.3 **Salient** span masks: 38.2

REALM: MLM pre-training



- Can use corpora even larger than Wikipedia for pre-training
 - CC-News vs Wikipedia: **40.4** vs **39.2**
 - Can allow updating evidence encoder (**REDT_**)

```
graph LR; IB[Index builder (stale theta')] -- MIPS index of Z --> MLM[MLM trainer (fresh theta)]; MLM -- "Updates theta' <- theta" --> IB;
```
- MIPS: Maximum Inner Product Search
- ORQA vs REALM: 33.3 vs 40.4 on NQ

Take-aways from ORQA/REALM

- It is possible to jointly train the retriever and reader, without any sparse IR components.
- The pre-training makes this retrieval process feasible and pre-training strategies matter.
- However, pre-training is very expensive.
“We pre-train for 200k steps on 64 Google Cloud TPUs, with a batch size of 512”

Next question: Is pre-training really necessary?

Retrieval-free QA: Probing Language Models

- Key question: can we use **pre-trained language models** to act as “knowledge storage”?
- Instead of explicitly storing all the text and searching among their *dense* or *sparse* representations, can we query the LMs to obtain the answer directly?
- The LMs were pre-trained on Wikipedia (and other textual corpora) so they should be able to memorize a fair amount of information.

Language Models as Knowledge Bases

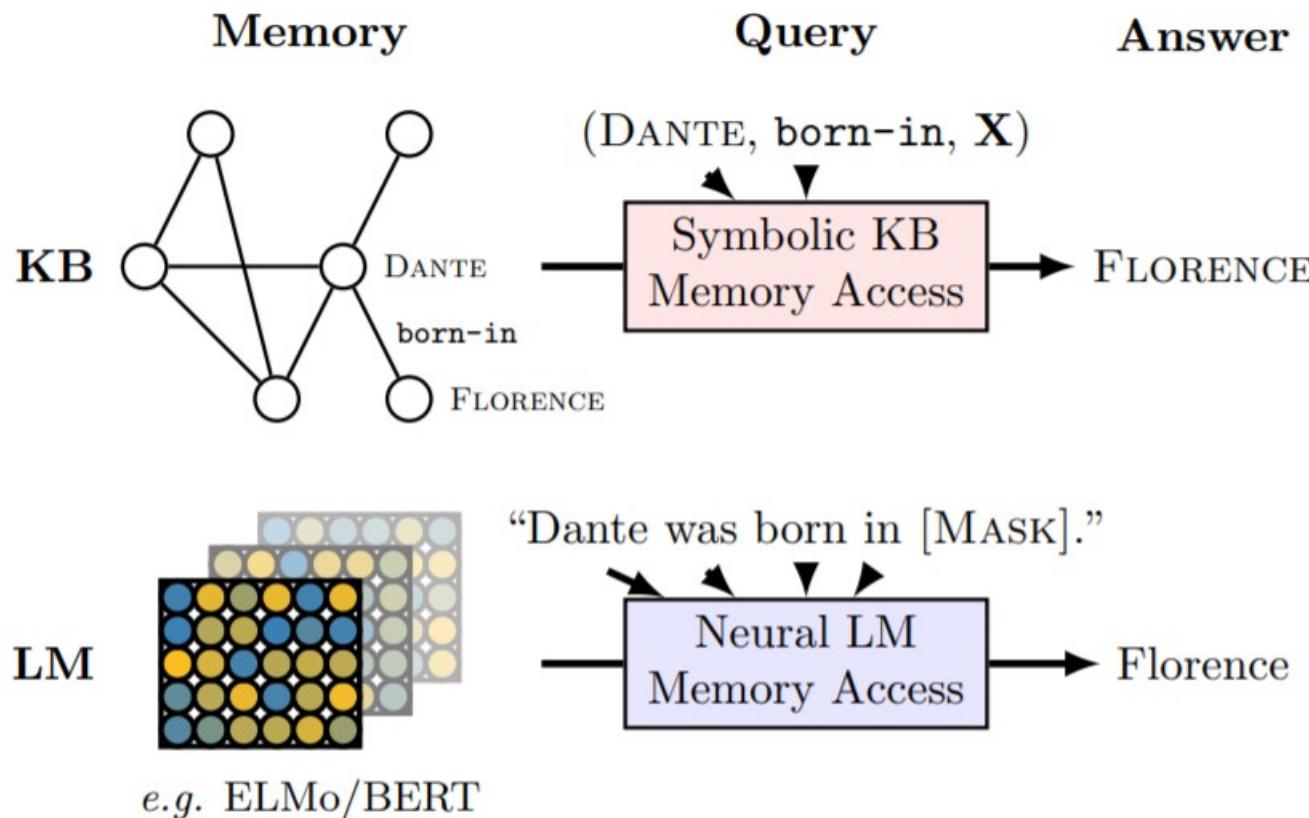


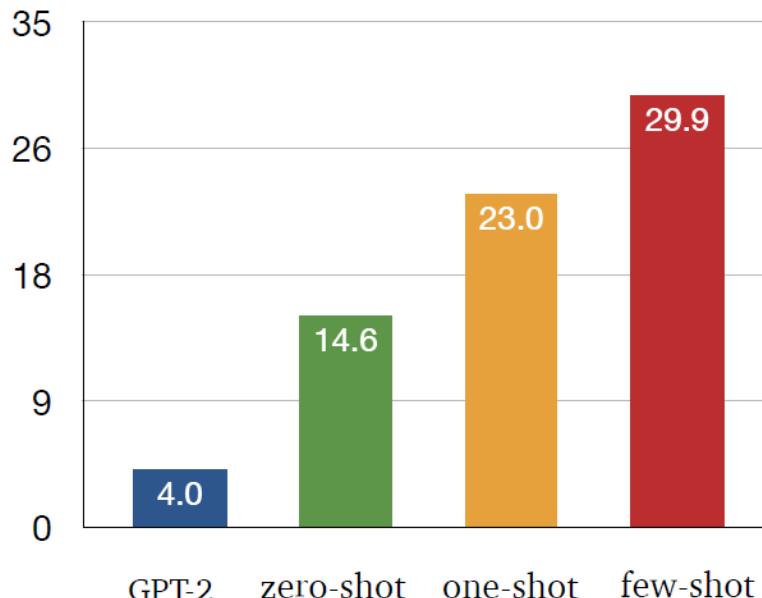
Figure 1: Querying knowledge bases (KB) and language models (LM) for factual knowledge.

GPT-3: Few-shot Learner [Brown et al., 2020]

96 layers, hidden size 12288, **175B** parameters

Larger corpora:
Common Crawl
+ WebText + Books
+ **English Wikipedia**

Evaluated on Natural



Few-shot learner

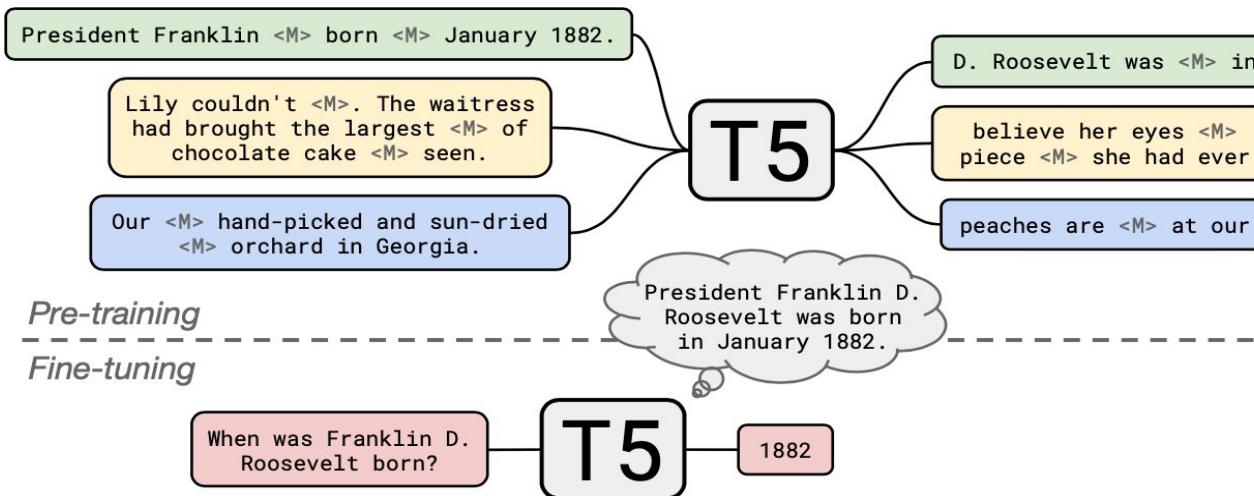
- $Q_1, A_1, Q_2, A_2, \dots, Q_k, A_k, Q ?$
- One-shot setting is a special case when only **one** example is given.

T5: Fine-tuning leads to improved performance

[Roberts et al., 2020]

Text-to-Text Transfer Transformer [Raffel et al., 2019]

Original text
Thank you ~~for inviting~~ me to your party last week.
Inputs
Thank you <X> me to your party <Y> week.
Targets
<X> for inviting <Y> last <Z>



*: Pre-trained on a multitask mixture including an **unsupervised “span corruption” task** on unlabeled text as well as supervised translation, summarization, classification, and reading comprehension tasks

T5: Fine-tuning leads to improved performance

S	Natural Question	WebQuestionTriviaQA		
		NQ	WQ	TQA
Chen et al. (2017)		–	20.7	–
Lee et al. (2019)		33.3	36.4	47.1
Min et al. (2019a)		28.1	–	50.9
Min et al. (2019b)		31.8	31.6	55.4
Asai et al. (2019)		32.6	–	–
Ling et al. (2020)		–	–	35.7
Guu et al. (2020)		40.4	40.7	–
Févry et al. (2020)		–	–	53.4
Karpukhin et al. (2020)		41.5	42.4	57.9
220	T5-Base	27.0	29.1	29.1
M	T5-Large	29.8	32.2	35.9
770	T5-3B	32.1	34.9	43.4
M	T5-11B	34.5	37.4	50.1
3	T5-11B + SSM	36.6	44.7	60.5

BERT-base = 110M parameters

11
B

SSM: salient span masking,
pre-training data
proposed in REALM

Why retrieval-based LMs?

LLMs can't memorize all (long-tail) knowledge in their parameters



List 5 important papers authored by Geoffrey Hinton

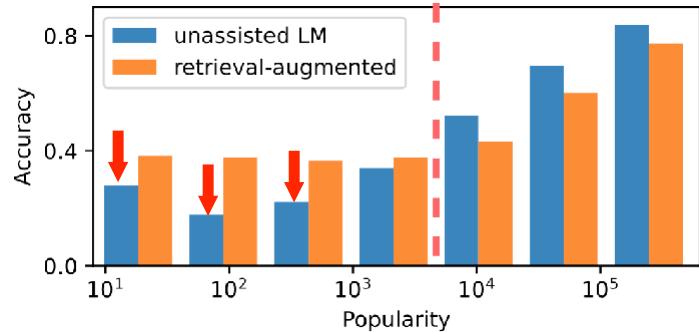


ChatGPT

Geoffrey Hinton is a renowned computer scientist ... Here are five important papers authored by him:

1. "**Learning Internal Representations by Error Propagation**" (with **D. E. Rumelhart** and **R. J. Williams**) - This paper, published in **1986**, .. ✓
2. "**Deep Boltzmann Machines**" (with **R. Salakhutdinov**) - Published in **2009**, .. ✗
3. ... ✗
4. "**Deep Learning**" (with Y. Bengio and A. Courville) - Published as a book in **2016**,... ✗
5. "**Attention Is All You Need**" (with V. Vaswani, N. Shazeer, et al.) - Published in **2017**, this paper introduced the Transformer model,...

What is Kathy Saltzman's occupation?



(Mallen et al., 2023)

GPT-3 davinci-003: 20%-30% accuracy

Why retrieval-based LMs?

LLMs' knowledge is easily outdated and hard to update



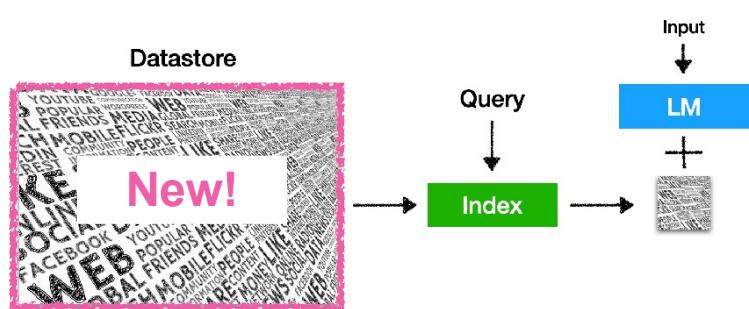
Who is the CEO of Twitter?



As of my **knowledge cutoff** in September 2021, the CEO of Twitter is Jack Dorsey....

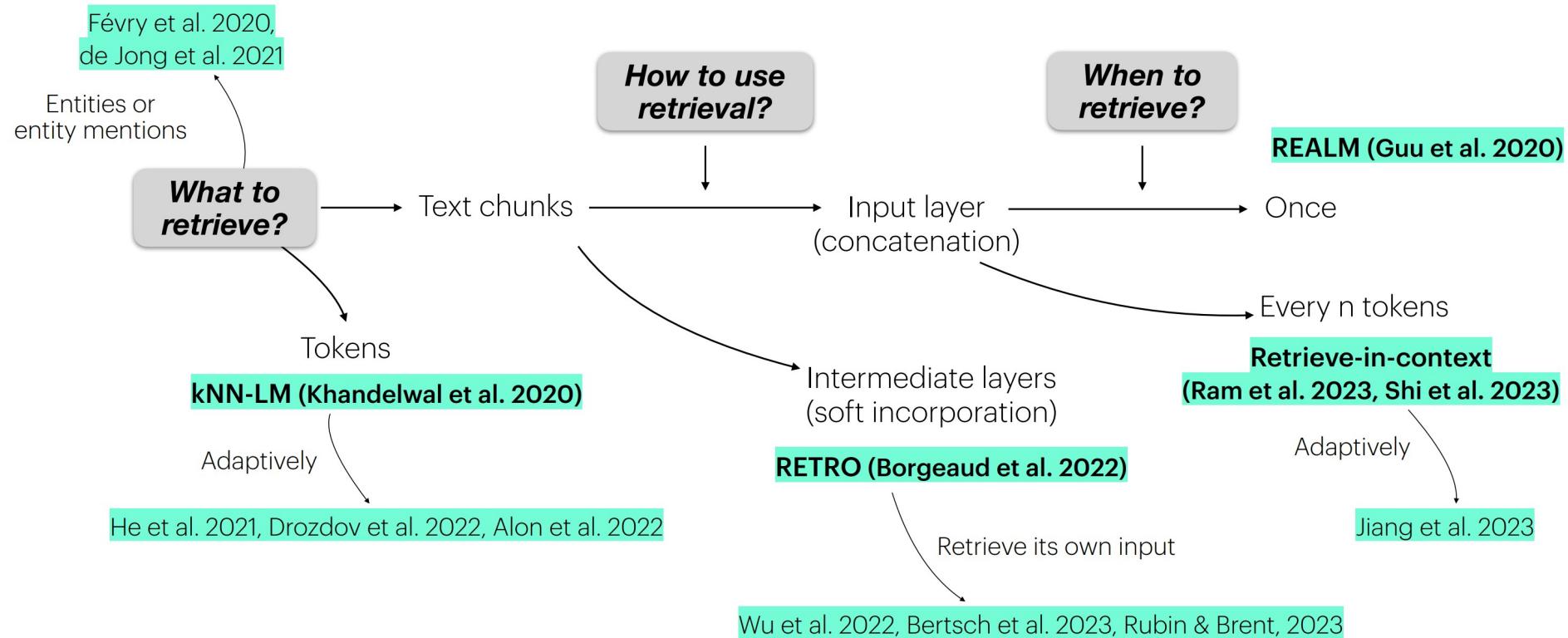
The screenshot shows a Google search results page. The query "Who is the CEO of Twitter?" is entered in the search bar. Below the search bar, there are navigation links for All, News, Images, Shopping, Videos, and More. It also shows that there are about 1,090,000,000 results found in 0.45 seconds. The top result is a snippet for Linda Yaccarino, dated Jun 5, 2023. To the right of the snippet is a small portrait photo of Linda Yaccarino.

- Existing **knowledge editing** methods are still NOT scalable (**active research!**)
- The datastore can be easily **updated** and **expanded** - even without retraining!



Retrieval Augmented Generation

Roadmap

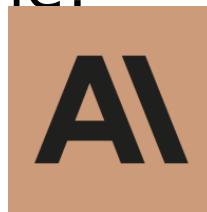


REPLUG

- **Title:** REPLUG: Retrieval-Augmented Black Box Language Models
- **Authors:** Shi et al. (University of Washington, Stanford, KAIST, Meta AI)
- **Publication:** ArXiv 2023
- **Link:** <https://arxiv.org/abs/2301.12652>

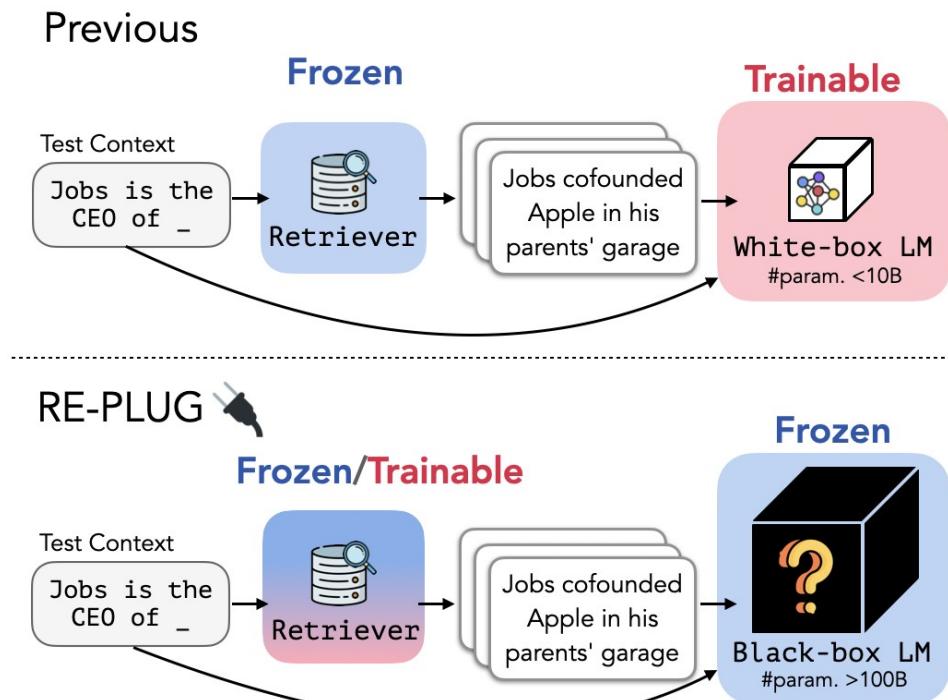
Motivation: RAG for Black-Box LMs

- **Issue:** Fine-tuning LMs with retrieval requires white-box access (i.e., access to the LM parameters)
 - To train the model
 - To index the datastore
- **In Practice:** Many LLMs are only accessible via API
 - Can't access model parameters!
 - Can't fine-tune!



Introducing REPLUG

- **Focus:** Retrieval-augmentation in the **black-box setting** (no access to model parameters)
 - LM is a black box
 - Retriever is optionally tunable
- **How does it work?**
 1. Get documents from retriever
 2. Prepend documents to LM input
 3. Feed input to LM



Implementation: Retrieval

- **Given:**
 - Corpus of documents
 - Input context
- **Output:**
 - The documents in \mathcal{D} that are most relevant to \mathbf{x} , denoted
- **How:**
 - Precompute embedding for each document (construct FAISS index)
 - Encode input \mathbf{x} to obtain embedding
 - Find top- k most similar documents to \mathbf{x} using cosine similarity function

Hang On, Can We Fit Documents?

- **Context Window:** The input text for an LM
- **Context Window Size:** Number of tokens an LM can process as input
- **Problem:** Context window may not fit documents
 - The context window size varies by LM
 - The size of each document could vary
 - The number of retrieved documents could vary

Solution: Ensemble Inference Scheme

- **Idea:** Do separate LM predictions!
 1. Concatenate each with
 2. Ensemble the output probabilities from inference calls

