

Question Answering / Retrieval

CSE 5525: Foundations of Speech and Natural Language
Processing

<https://shocheen.github.io/courses/cse-5525-fall-2025>



THE OHIO STATE UNIVERSITY

Logistics

- Final Project Presentations
 - Dec 5, Dec 10
 - I'll randomly assign you to one of the days, if you can't make one of them please reach out to me
 - Each presentation will be 7-8 minutes depending on total number of teams with 1-2 time for QA (I will enforce hard time stops)
 - Will post presentation format on Canvas shortly.
- Project Reports will be due: Dec 17 (absolutely no extensions).

Goal: Dive into one NLP application: Question Answering

- 📌 QA Landscape
- 📌 An Overview of Retriever-Reader & Retriever-Generator Architectures for Open-Ended QA
- 📌 Retrieval
- 📌 Answer Extractor
- 📌 Retrieval Augmented Generation (RAG)
- 📌 RAG: Overview of Retriever-Generator training options

Question Answering

We fill our information needs by talking to a virtual assistant or a chatbot, interacting with a search engine, or querying a database

No wonder that question answering has been a major task in NLP

But there isn't one QA task

Intent behind the question:

Was the person seeking information they did not have,

... or trying to test the knowledge of another person or machine?

The task of **Question Answering (QA)** in NLP is more often associated with information-seeking questions:

- The questions are real queries by users of products like Google Search, Reddit, StackOverflow, etc
- Thus, questions are often ill-specified, full of “ambiguity and presupposition”
- Less frequently involve complex reasoning
- Questions often assume no given context and are almost never poised as multiple choice
- Industrial research because research progress directly translates to improved products

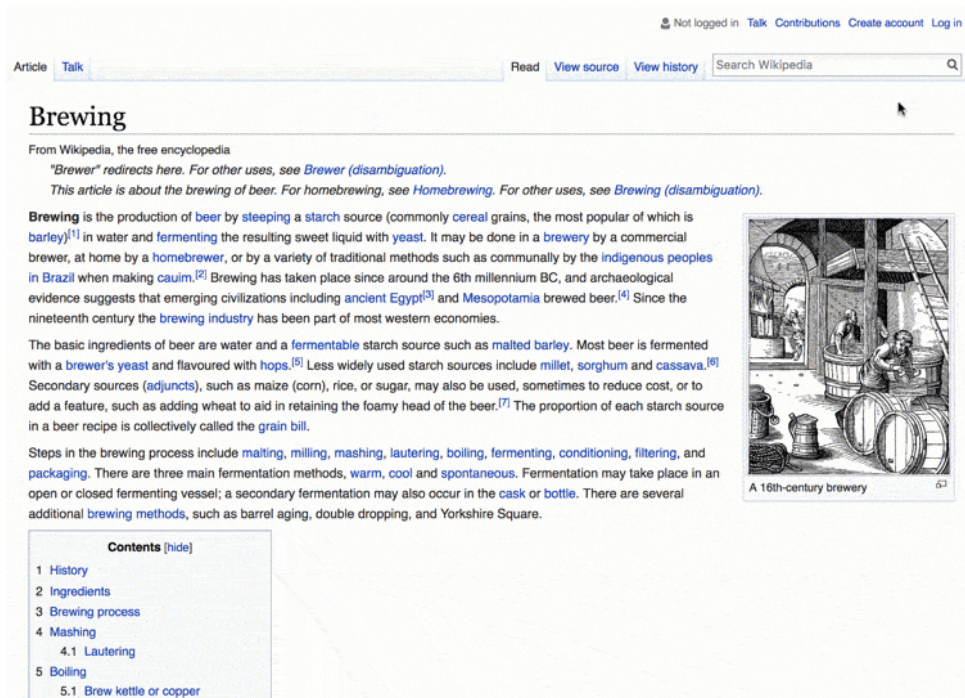
Example of such QA: NaturalQuestions [\[Kwiatkowski et al., 2019\]](#)

Question: when are hops added to the brewing process?

Short answer: The boiling process

Long answer:

After mashing , the beer wort is boiled with hops (and other flavourings if used) in a large tank known as a " copper " or brew kettle – though historically the mash vessel was used and is still in some small breweries . The boiling process is where chemical reactions take place , including sterilization of the wort to remove unwanted bacteria , releasing of hop flavours , bitterness and aroma compounds through isomerization , stopping of enzymatic processes , [...]



But there isn't one QA task

Intent behind the question:

Was the person seeking information they did not have,

... or trying to test the knowledge of another person or machine?

The task of **Question Answering (QA)** in NLP is more often associated with information-seeking questions:

- The questions are real queries by users of products like Google Search, Reddit, StackOverflow, etc
- Thus, questions are often ill-specified, full of “ambiguity and presupposition”
- Less frequently involve complex reasoning
- Questions often assume no given context and are almost never posed as multiple choice
- Industrial research because research progress directly translates to improved products

Reading Comprehension (RC) is more associated with probing questions

- A block of text that contains information relevant to the questions being asked is provided
- Probe the understanding of the given block of text

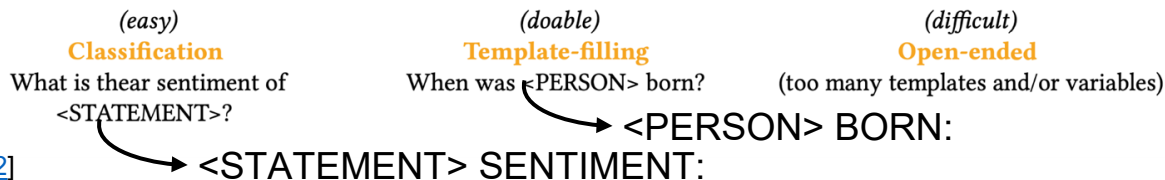
Almost any NLP task can be formulated as question answering
This is leveraged for model reuse, multi-task learning, prompting

- Example: "Is this movie review sentence negative or positive?"

In such cases, QA is **not a task but a format**: "a way of posing a particular problem to a machine, just as classification or natural language inference are formats" [\[Gardner et al., 2019\]](#)

The key distinction to keep in mind is "how easy would it be to replace the questions in a dataset with content-free identifiers?"

[\[Gardner et al., 2019\]](#)



[\[Rogers et al., 2022\]](#)

Question formats

Evidence	Format	Question	Answer	Example datasets
Einstein was born in 1879.	Questions	When was Einstein born?	1879	SQuAD [235], RACE [156]
	Queries	Which year Einstein born	1879	generated queries in BEIR [282]
	Cloze	Einstein was born in ____.	1979	CNN/Daily Mail [125], CBT [127]
	Completion	Einstein was born ...	in 1879	SWAG [319], RocStories [204]

Natural questions: questions that a human speaker could ask

- Yes/no questions (Did it rain on Monday?)
- Wh-questions (When did it rain?)

Queries: Pieces of information that could be interpreted as a question

Cloze format: Sentences with a masked span that the model needs to predict

Story completion: The choice of the alternative endings for the passage

Answer formats

Evidence	Format	Question	Answer(s)	Example datasets
Einstein was born in 1879.	Extractive	When was Einstein born?	1879 (token 5)	SQuAD [235], NewsQA [287]
	Multi-choice	When was Einstein born?	(a) 1879, (b) 1880	RACE [156]
	Categorical	Was Einstein born in 1880?	No	BoolQ [65]
	Freeform	When was Einstein born?	1879 (generated)	MS MARCO [19], CoQA [238]

Extractive format / extractive QA

- Either get or find short text segments from the web or some other large collection of documents
- Answer a question by extracting the answer from the evidence text
 - Extracting = find the starting and ending token in the input text
 - If the answer is generated \Rightarrow **Retrieval-augmented generation (RAG)**
- The limited range of answer options means that it is easier to define what an acceptable correct answer is
- It limits the kinds of questions that can be asked to questions with answers directly contained in the text

Answer formats (cont.)

Evidence	Format	Question	Answer(s)	Example datasets
Einstein was born in 1879.	Extractive	When was Einstein born?	1879 (token 5)	SQuAD [235], NewsQA [287]
	Multi-choice	When was Einstein born?	(a) 1879, (b) 1880	RACE [156]
	Categorical	Was Einstein born in 1880?	No	BoolQ [65]
	Freeform	When was Einstein born?	1879 (generated)	MS MARCO [19], CoQA [238]

Multiple-choice format

- Questions for which a small number of answer options are given as part of the question text itself
- The answers are no longer restricted to something explicitly stated in the text
- The question writer also has full control over the available options, and therefore over the kinds of reasoning that the test subject would need to be capable of
- Writing good multi-choice questions is not easy: if the incorrect options shouldn't be easy to rule out

Answer formats (cont.)

Evidence	Format	Question	Answer(s)	Example datasets
Einstein was born in 1879.	Extractive	When was Einstein born?	1879 (token 5)	SQuAD [235], NewsQA [287]
	Multi-choice	When was Einstein born?	(a) 1879, (b) 1880	RACE [156]
	Categorical	Was Einstein born in 1880?	No	BoolQ [65]
	Freeform	When was Einstein born?	1879 (generated)	MS MARCO [19], CoQA [238]

Categorical format

- The answers come from a strictly pre-defined set of options (yes/no)

Freeform format

- Generate the answer independently rather than choose from the evidence or available alternatives
- The “gold” answer is probably not the only correct one, which makes evaluation difficult
- Evaluate on at least two axes: linguistic fluency and factual correctness; both hard

Goal: Dive into one NLP application: Question Answering

- 📌 QA Landscape
- 📌 An Overview of Retriever-Reader & Retriever-Generator Architectures for Open-Ended QA
- 📌 Retrieval
- 📌 Answer Extractor
- 📌 Retrieval Augmented Generation (RAG)
- 📌 RAG: Overview of Retriever-Generator training options

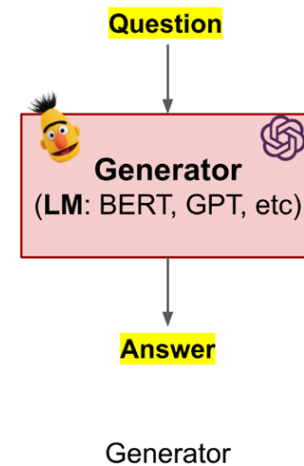
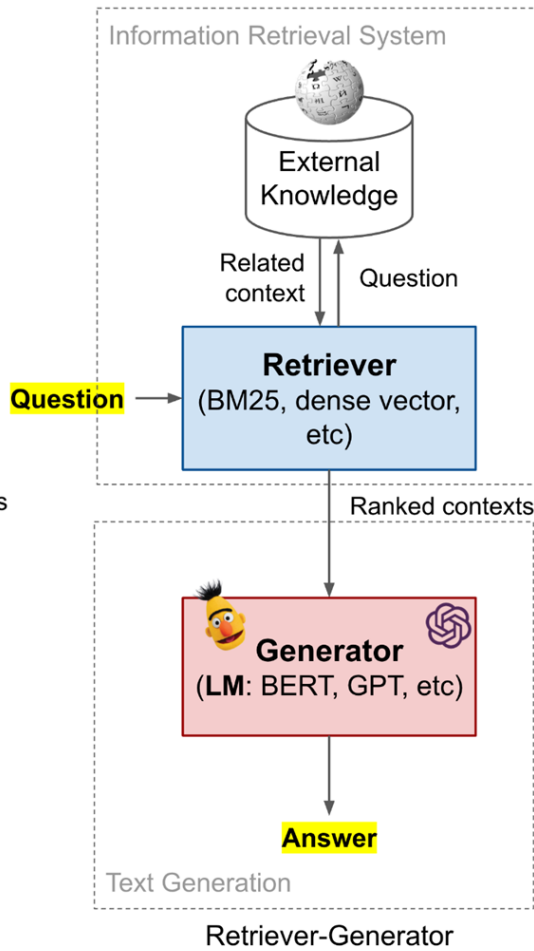
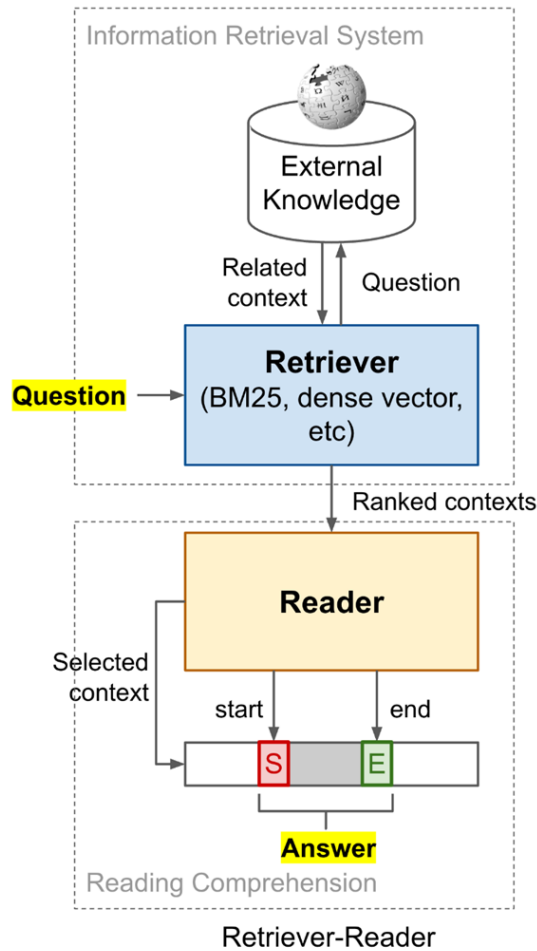
Open-domain Question Answering (ODQA): Asking a model to produce answers to, typically factoid, questions in natural language

The “open-domain” part:

- ⚡ Relevant context is not provided and needs to be found:
 - When both the question and the context are provided, the task is known as *reading comprehension*
- ⚡ Answer choices are not provided
 - When the choices are provided, the task is multiple-choice QA

General idea:

- ⚡ An LLMs’ pretraining data potentially didn’t contain relevant information, or even if it did, the model is not capable enough to memorize it and refer to it when asked to answer a related question
- ⚡ Find the relevant context where the answer is contained, then answer the question condition on the found context



Goal: Dive into one NLP application: Question Answering

- 📌 QA Landscape
- 📌 An Overview of Retriever-Reader & Retriever-Generator Architectures for Open-Ended QA
- 📌 **Retrieval**
- 📌 Answer Extractor
- 📌 Retrieval Augmented Generation (RAG)
- 📌 RAG: Overview of Retriever-Generator training options

TF-IDF

```
t = token  
d = document (movie review)  
term_frequency(t,d) = number of times t occurs in d  
document_frequency(t) = # documents t occurs in  
N = number of documents  
inverse_document_frequency(t) = N / document_frequency(t)  
tf-idf(t,d) = tf(t,d) x inverse_document_frequency(t)  
score(query,d) = sum([tf-idf(t,d) for t in query])
```

Return documents that are the most similar to the query

TF-IDF \Rightarrow BM25 [\[Robertson et al., 1995\]](#)

BM25 is really just a more refined version of TF-IDF with **two additional hyperparameters** [\[Kamphuis et al. \(2020\)\]](#)

- k**, a knob that adjust the balance between term frequency and IDF,
- **k** \rightarrow 0, term frequency has almost no effect, IDF-driven
 - E.g., keywords in short texts like tweets/titles, repetition stylistic, special domains

b, which controls the importance of document length normalization

$$\sum_{t \in q} \overbrace{\log \left(\frac{N}{df_t} \right)}^{\text{IDF}} \overbrace{\frac{tf_{t,d}}{k \left(1 - b + b \left(\frac{|d|}{|d_{\text{avg}}|} \right) \right) + tf_{t,d}}}^{\text{weighted tf}}$$

Why could this classic IR approach to retrieval be limited?

tf-idf/BM25 algorithms work only if there is exact overlap of tokens between the question/query and document

What if a question contains a lot of synonyms of tokens in a relevant document?

Vocabulary mismatch problem: The user posing a query (or asking a question) needs to guess exactly what words the writer of the answer might have used

→ Instead of (sparse) word-count vectors, using (dense) embeddings

Two ways to dense retrieval

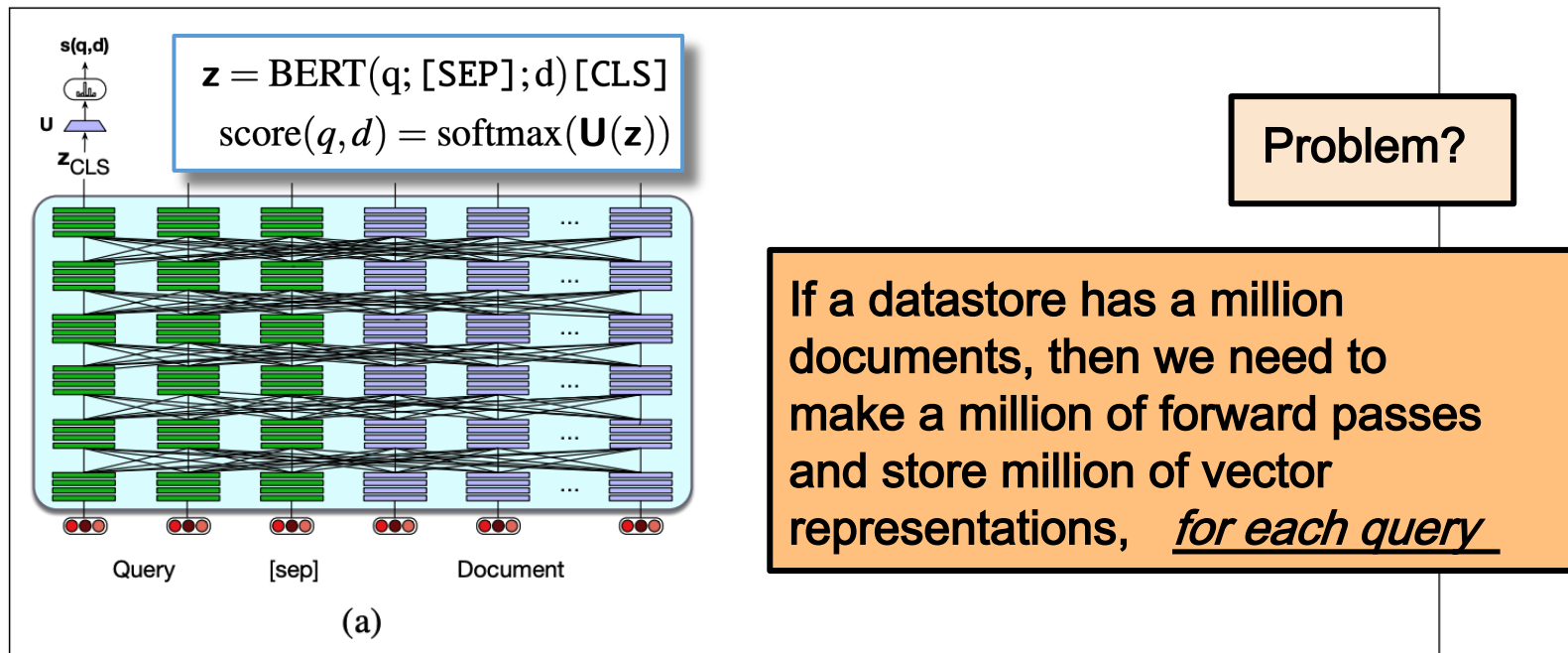


Figure 14.7 Two ways to do dense retrieval, illustrated by using lines between layers to schematically represent self-attention: (a) Use a single encoder to jointly encode query and document and finetune to produce a relevance score with a linear layer over the CLS token. This is too compute-expensive to use except in rescoring (b) Use separate encoders for query and document, and use the dot product between CLS token outputs for the query and document as the score. This is less compute-expensive, but not as accurate.

Two ways to dense retrieval

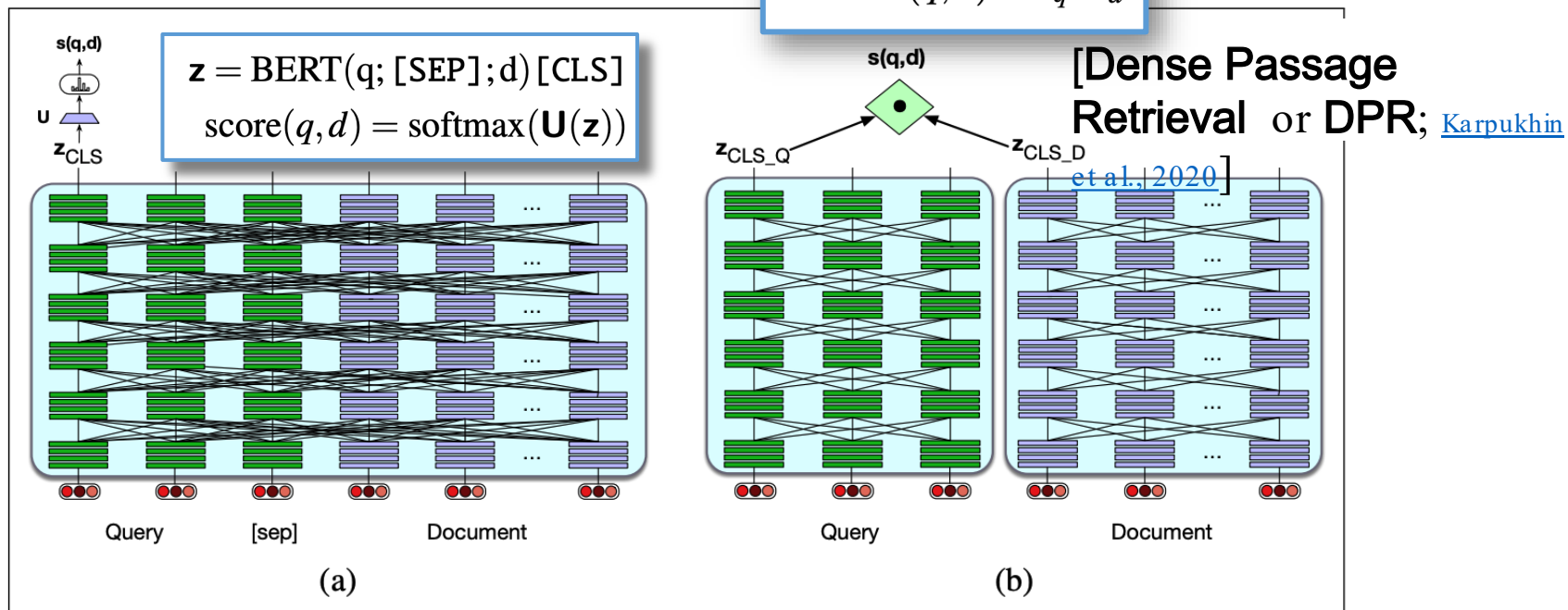


Figure 14.7 Two ways to do dense retrieval, illustrated by using lines between layers to schematically represent self-attention: (a) Use a single encoder to jointly encode query and document and finetune to produce a relevance score with a linear layer over the CLS token. This is too compute-expensive to use except in rescoring (b) Use separate encoders for query and document, and use the dot product between CLS token outputs for the query and document as the score. This is less compute-expensive, but not as accurate.

Another way to the second approach – ColBERT

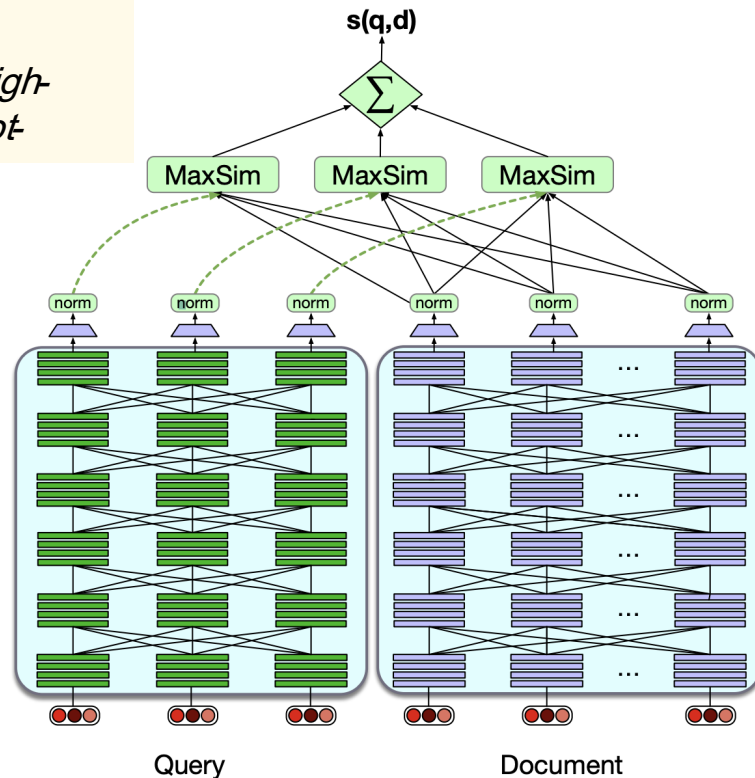
[Khattab et al., 2021]

Problem: DPR representations are relatively coarse. They encode each passage into a single high-dimensional vector & estimate relevance via one dot-product

$$\text{score}(q, d) = \sum_{i=1}^N \max_{j=1}^m \mathbf{E}_{q_i} \cdot \mathbf{E}_{d_j}$$

BERT output vectors
rescaled to unit length

Essentially, for each token in q , ColBERT finds the most contextually similar token in d , & then sums up these similarities



Training dense retrievals - Overview

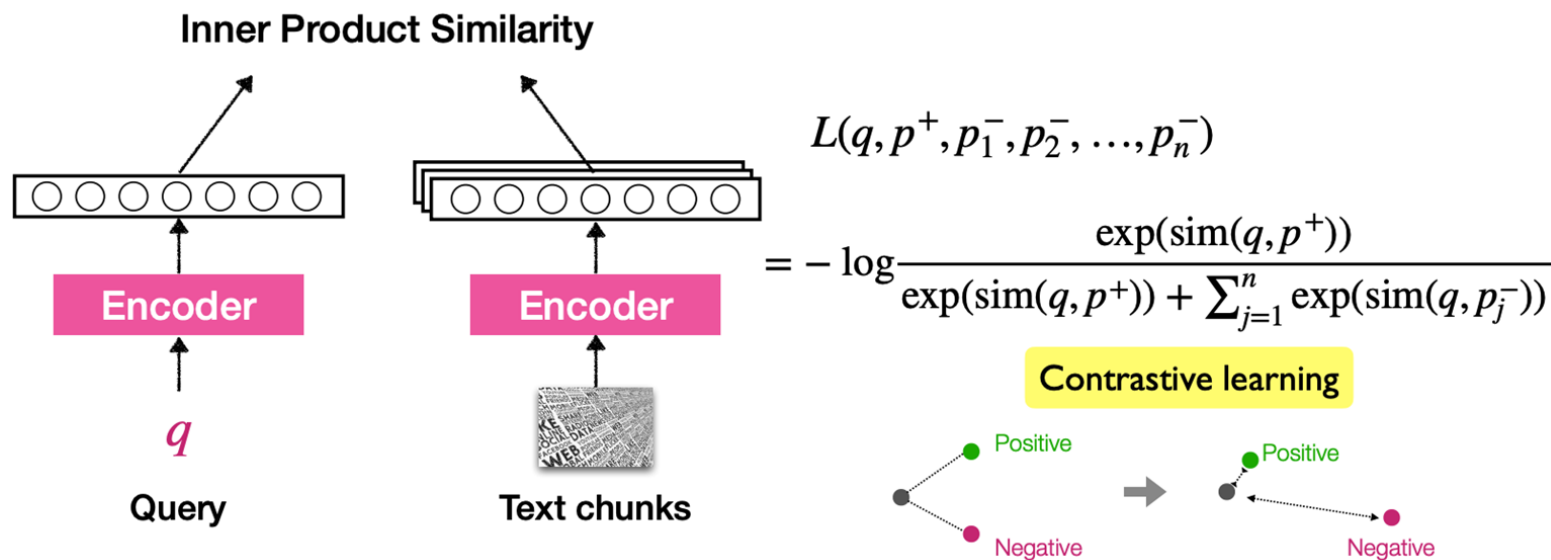
Objective: Train a model to predict the relevance of a document to a given query by optimizing a similarity score (e.g., dot product)

Training Data:

- ⤵ Requires queries paired with relevant (positive) and irrelevant (negative) documents
- ⤵ **Positive examples:** From datasets containing human-annotated relevant documents
- ⤵ **Negative examples:**
 - Typically sampled from the top-1000 results retrieved by an existing **retriever** (e.g., BM25)
 - This ensures that the negatives are challenging \Rightarrow Discourage the model from learning only trivial distinctions between relevant and irrelevant documents
- ⤵ If there are no labeled positive examples, use **relevance-guided supervision** [[Khattab et al., 2021](#)]:
 - Extract the top-ranked documents from an existing IR system that contain relevant answer strings (treat these as positive examples)
 - Use documents from the same top-ranked results that lack answer strings as negative examples
 - Train a new retriever and iterate

Training dense retrievals (cont.)

Train the model to maximize the score for positive documents and minimize the score for negative ones:



Maximum Inner Product Search (MIPS)

MIPS: The problem of determining a small subset of items in a datastore whose vector representations have the highest dot product with a given query vector

Although there is an **obvious linear-time implementation $O(\# \text{datastore items} \cdot \text{vector size})$** , it is generally **too slow** to be used on practical problems because the datastores are huge (millions or billions of vectors, each hundreds or thousands of dimensions)

To optimize the retrieval speed, the common choice is the approximate nearest neighbors (ANN) algorithm to return approximately top-k nearest neighbors to **trade off a little accuracy lost for a huge speedup**; <https://ann-benchmarks.com/>

Index

During the **indexing phase**, each document/passage is encoded into a vector

All these document vectors are stored in the **index**, typically in a structure that allows for fast similarity search

FAISS (Facebook AI Similarity Search) <https://github.com/facebookresearch/faiss>

- A library developed by Meta (Facebook) to perform fast similarity search on large-scale dense vectors
- FAISS builds an index by dividing the dense space into clusters using methods like k-means clustering
- At search time, the query vector is compared only with the most relevant clusters, significantly reducing the number of comparisons

Retrieval Evaluation

Each document returned by the IR system is either relevant to our purposes or not relevant

How to evaluate the retrieval then?

- **Precision:** # relevant documents among all documents retrieved (to be relevant)
- **Recall:** # truly relevant documents that are retrieved
- **F1:** as always, a weighted harmonic mean of the precision and recall

Issue: It doesn't really measure the performance of a system that ranks the documents

	What do retrieve?	How to use retrieval?	When to retrieve?
REALM (Guu et al 2020)	Text chunks	Input layer	Once
Retrieve-in-context LM (Shi et al 2023, Ram et al 2023)	Text chunks	Input layer	Every n tokens
RETRO (Borgeaud et al. 2021)	Text chunks	Intermediate layers	Every n tokens
kNN-LM (Khandelwal et al. 2020)	Tokens	Output layer	Every token
FLARE (Jiang et al. 2023)	Text chunks	Input layer	Every n tokens (<i>adaptive</i>)
Adaptive kNN-LM (He et al 2021, Alon et al 2022, etc)	Tokens	Output layer	Every n tokens (<i>adaptive</i>)
Entities as Experts (Fevry et al. 2020), Mention Memory (de Jong et al. 2022)	Entities or entity mentions	Intermediate layers	Every entity mentions
Wu et al. 2022, Bertsch et al. 2023, Rubin & Berant. 2023	Text chunks <i>from the input</i>	Intermediate layers	Once or every n tokens

Goal: Dive into one NLP application: Question Answering

- 📍 QA Landscape
- 📍 An Overview of Retriever-Reader & Retriever-Generator Architectures for Open-Ended QA
- 📍 Dense Retrieval
- 📍 **Answer Extractor**
- 📍 Retrieval Augmented Generation (RAG)
- 📍 RAG: Overview of Retriever-Generator training options

Reader algorithms: Answer span extractor

Goal: Given retrieved text (e.g., passage) and the question, predict which token is the start of the answer span and which token is its end

S, **E** ... vectors of the size of the final token representations, randomly initialized & trained

Concatenate question & passage, and delineate them:

- With BERT-like models, use [SEP]
- With today's LM use text like "Passage:"

Reader algorithms: Answer span extractor

Minimizing a negative log-likelihood loss encourages the model to maximize the probabilities of the true start (i^*) and end (j^*) positions

The probability of correctly predicting the answer span involves:

1. Predicting the correct start position
2. Predicting the correct end position

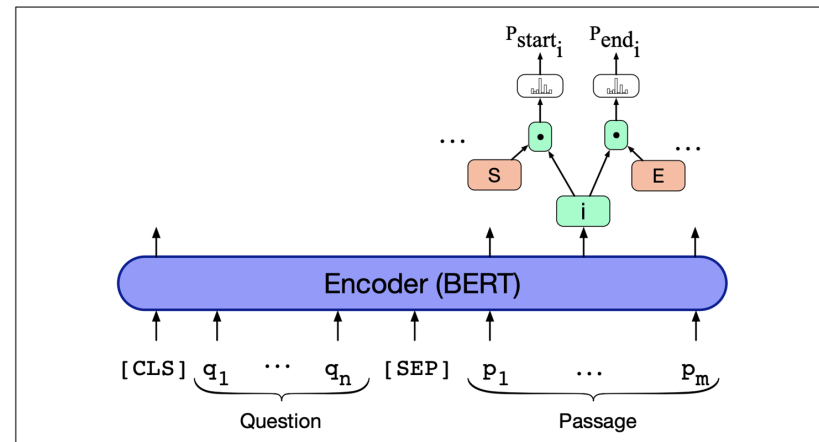


Figure 14.12 An encoder model (using BERT) for span-based question answering from reading-comprehension-based question answering tasks.

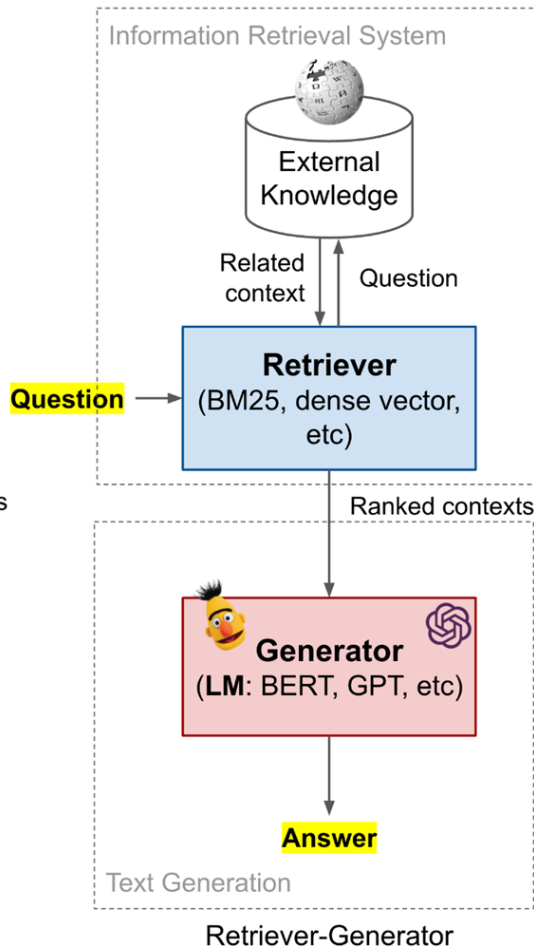
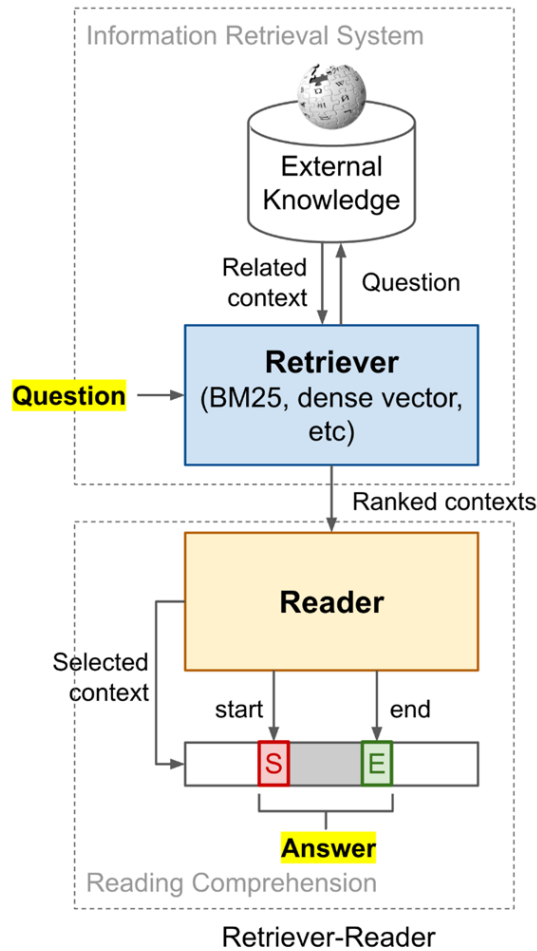
These two events are treated as independent events, the joint probability of both predictions being correct is the product of the two probabilities

Therefore the loss is:

$$L = -\log(P_{start_{i^*}} \cdot P_{end_{j^*}}) = -\log P_{start_{i^*}} - \log P_{end_{j^*}}$$

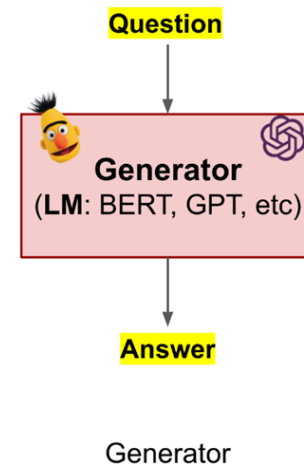
Goal: Dive into one NLP application: Question Answering

- 📍 QA Landscape
- 📍 An Overview of Retriever-Reader & Retriever-Generator Architectures for Open-Ended QA
- 📍 Dense Retrieval
- 📍 Answer Extractor
- 📍 Retrieval Augmented Generation (RAG)
- 📍 RAG: Overview of Retriever-Generator training options



A LM's *parametric knowledge* :
The information that the model has encoded within its parameters/weights during training that it can then use to do tasks for which that knowledge is required

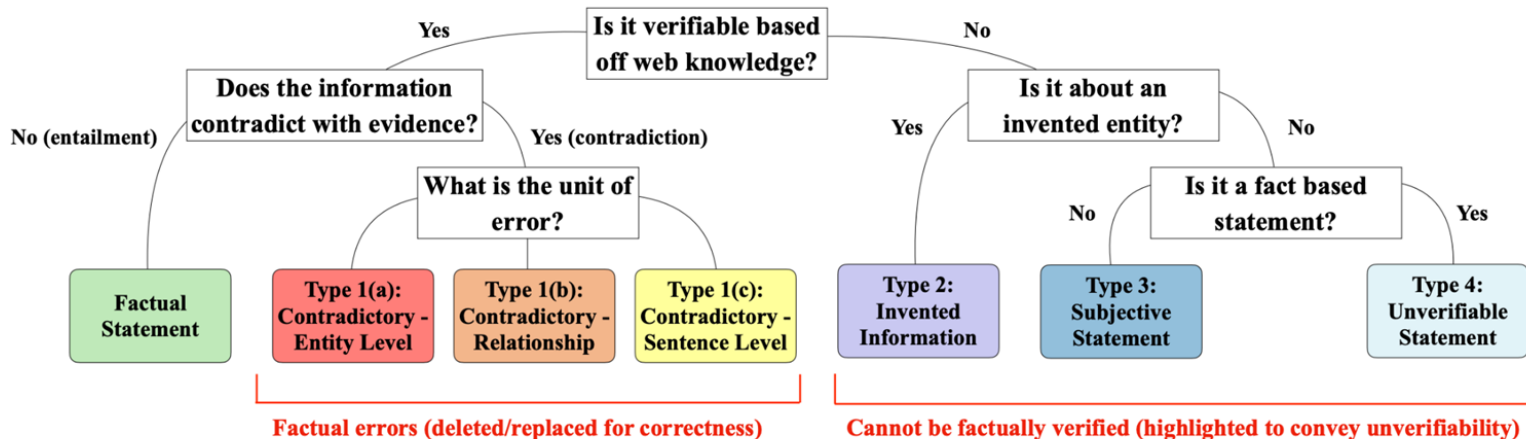
So, why LLMs need retrieval?



Hallucination & Factuality

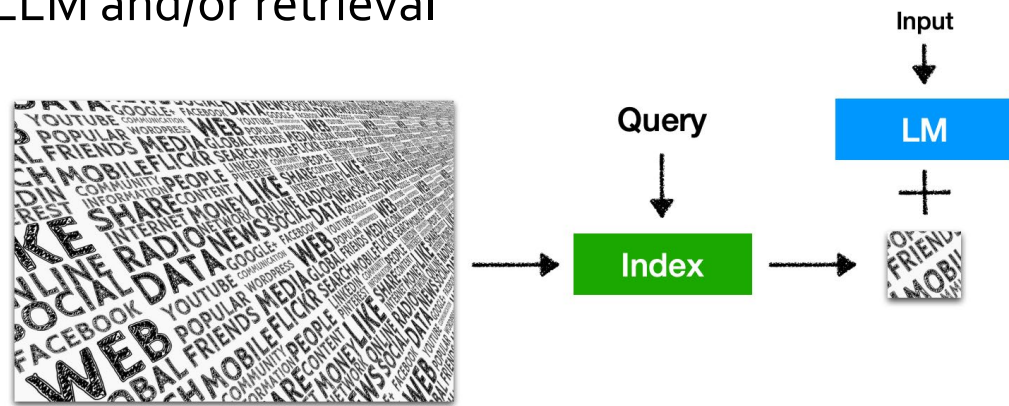
A hallucination is a response that is not faithful to the facts of the world

Fine-grained hallucination taxonomy [Mishra et al., 2024]:



Retrieval Augmented Generation (RAG)

1. Use a retriever to identify top-k most relevant documents for the user's prompt
2. Add the retrieved documents to the input
3. Ask the LLM to generate the answer based on this rich context
4. Finetune the answer-generator LLM and/or retrieval



Datastore

Raw text corpus

At least billions~trillions of tokens
Not labeled datasets

Not structured data (knowledge bases)

Reminder: Input Context Length

The attention matrix is quadratic in the maximum sequence length

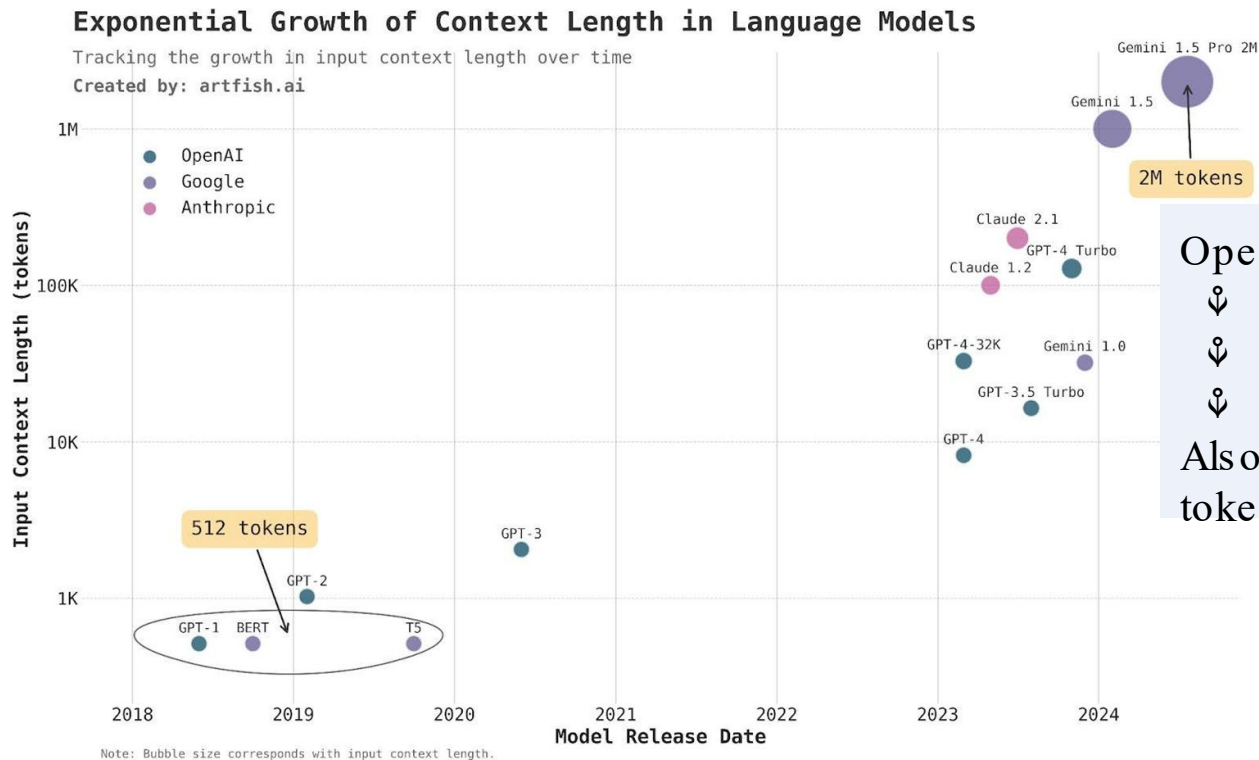
- If the length of an input sequence doubles, the amount of memory required quadruples
- Training an LLM on sequence lengths of 128k will require $\sim 1024\times$ the memory compared to training on sequence lengths of 4k
- GPU memory

Poor generalization due to position encoding:

- [RoPE](#) is the current choice (aims to preserve the relative distance between tokens)
- Performance quickly breaks down for sequence lengths significantly longer than the model has seen before [[Press et al., 2022](#)]

You can process sequences of arbitrary lengths, but you shouldn't expect a good performance for sequences longer than what's used for pretraining/post-training because of RoPE, & creators of LLMs are prevented from increasing the sequence length drastically due to the GPU memory limits

Input Context Length



Open-weight LLMs:

- ⌵ LLaMA 3.1
- ⌵ Qwen 2.5
- ⌵ Mistral-Large

Also fit 128K tokens!

Google's Gemini 1.5 can (almost) fit the entire Harry Potter + Lord of the Ring series in its 2 million context window

2M

Gemini 1.5 2M
(June 2024)



200K

Claude 2.1
(July 2023)



128K

GPT-4 Turbo
(March 2023)



GPT-3.5 Turbo
(March 2022)



Practical Considerations

1. Longer the input context length, the more potentially relevant documents we can squeeze, but:

⚡ [Liu et al., 2023]: “models are **better at using relevant information** that occurs **at the very beginning** (primacy bias) **or end of its input context** (recency bias), and performance degrades significantly when models must access and use information located in the middle of its input context”

1. [Ying et al., 2024]: Conflicts between internal/parametric knowledge and knowledge in a given context

⚡ Curation of the context is thus important

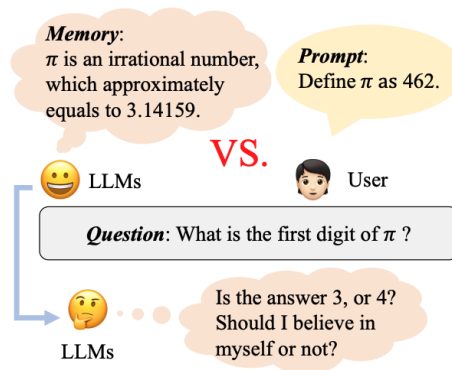
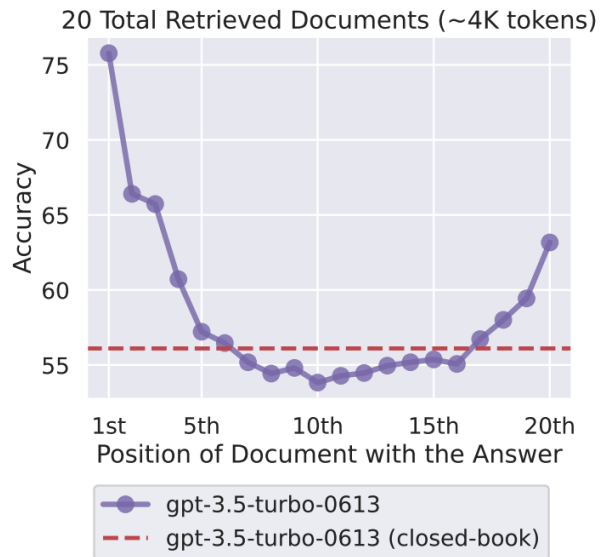


Figure 1: In conflict situation, LLMs may depend on the prompt or intuitively answer based on memory. 41

Goal: Dive into one NLP application: Question Answering

- 📍 QA Landscape
- 📍 An Overview of Retriever-Reader & Retriever-Generator Architectures for Open-Ended QA
- 📍 Dense Retrieval
- 📍 Answer Extractor
- 📍 Retrieval Augmented Generation (RAG)
- 📍 RAG: Overview of Retriever-Generator training options

Retrieval-based LMs: Training

Training challenges:

- ⚡ External datastore is huge
 - ⇒ Expensive to update the index = recompute dense vectors for all documents
- ⚡ LLMs are large
 - ⇒ Expensive to finetune an LLM to generate answers

Option 1 – Independent Training:

Retrieval models and language models are trained independently

Problem with this training option

We ignore that a retriever and a generator work together

We want to improve the RAG system such that:

- The retriever learns to select passages that make the generator's job easier
- The generator learns to adapt to whatever the retriever gives

$$L = -\log \sum_{d \in \text{Docs}} p_{\text{retriever}}(d|q; \theta_r) p_{\text{generator}}(y|q, d; \theta_g)$$

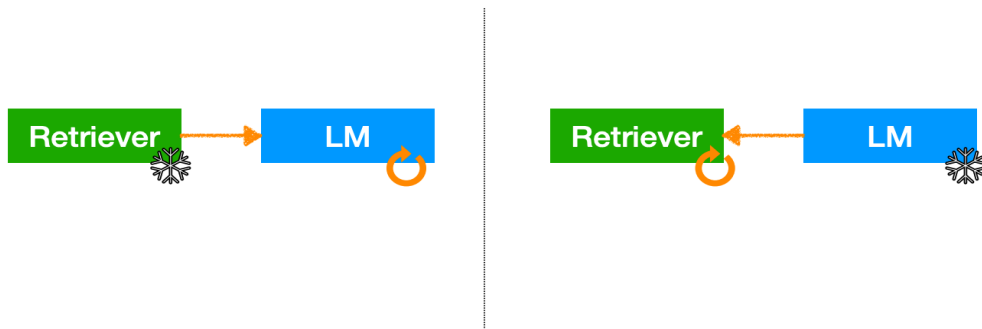
Millions or billions \Rightarrow Too
slow

Retrieval-based LMs: Training

Option 2 – Sequential Training:

One component is first trained independently & then fixed, the other component is trained with an objective that depends on the first one

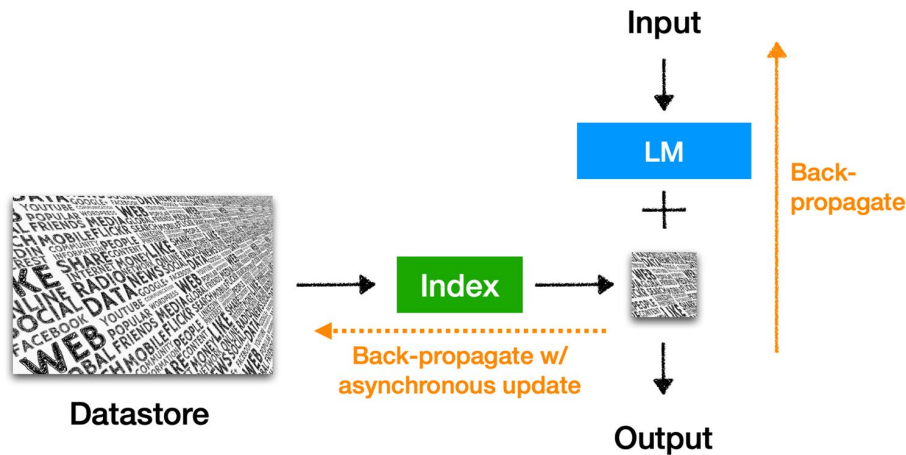
- ✚ Retrievers are trained to provide text that helps LMs the most



Retrieval-based LMs: Training

Option 3 – Joint training w/ asynchronous index update:

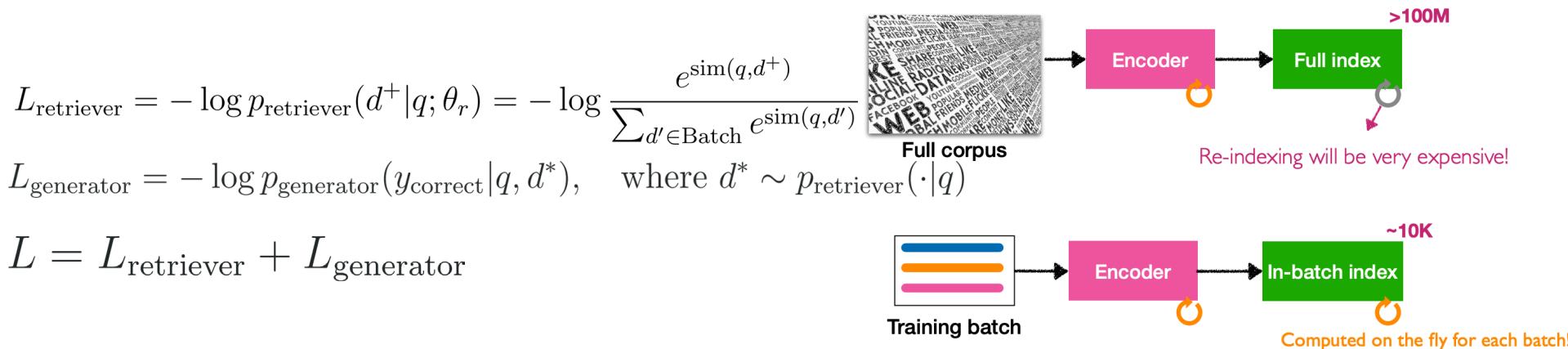
- Retrieval models and language models are trained jointly
- Allow the index to be “stale”; rebuild the retrieval index every T steps





Retrieval-based LMs: Training

Option 4 – Joint training w/ in-batch approximation:

- Retrieval models and language models are trained jointly
- Use “in-batch index” instead of full index
- Treat each other document in the batch as a negative example for a given query



Training method		
Independent training (Ram et al 2023; Khandelwal et al 2020)	* Easy to implement: off-the-shelf models	* Models are not end-to-end trained — suboptimal performance
Sequential training (Borgeaud et al 2021; Shi et al 2023)	* Easy to improve: sub-module can be separately improved	
Joint training: async update (Guu et al 2020; Izacard et al 2022)		* Training may be complicated (overhead, batching methods, etc)
Joint training: in-batch approx (Zhong et al 2022; Min et al 2023; Rubin and Berant 2023)	* End-to-end trained — very good performance!	* Train-test discrepancy still remains

How do retrieval-based language models perform on downstream tasks? → **Section 5!**

ODQA “Developer Guide”

<https://drive.google.com/file/d/1vh8S13V-LvgdhTlBrcvoXUiPYuzLCSTA/view?usp=sharing>