

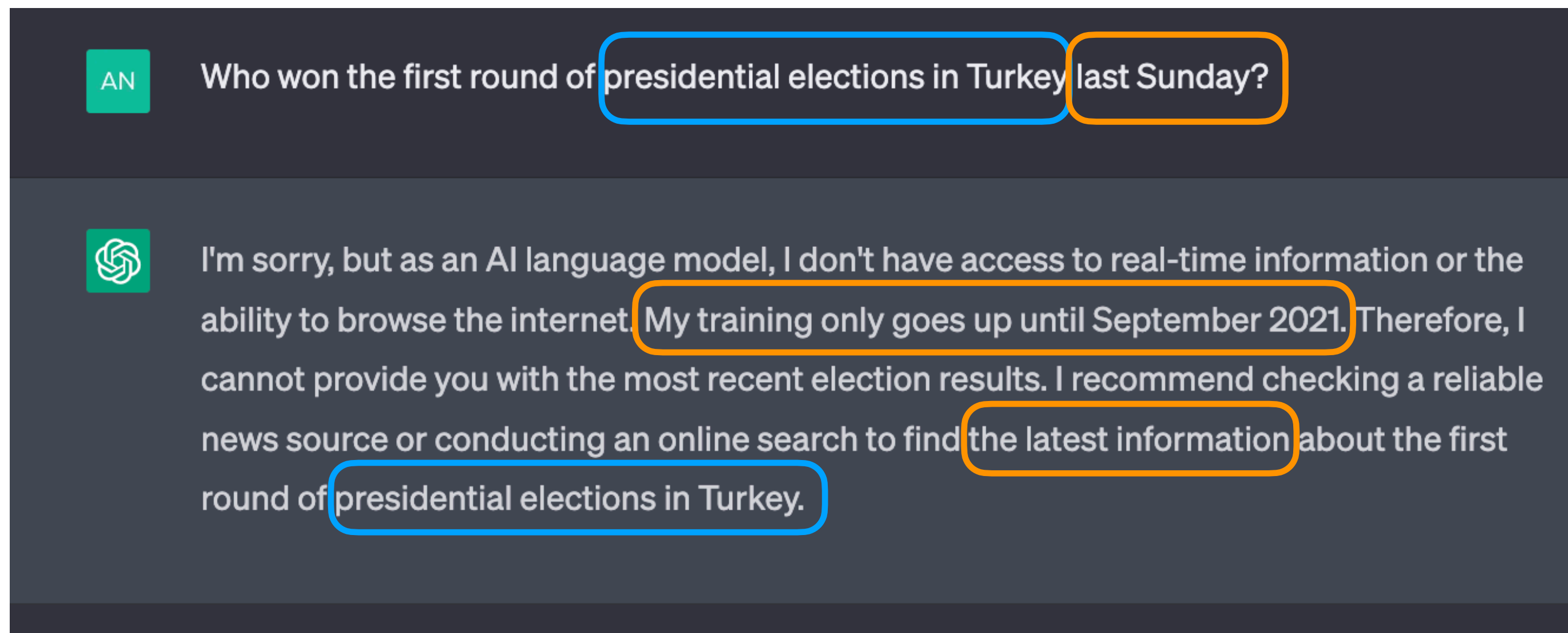
Retrieval-Augmented LMs

Antoine Bosselut

Reading Comprehension challenges

Factuality

Temporality



AN Who won the first round of presidential elections in Turkey last Sunday?

I'm sorry, but as an AI language model, I don't have access to real-time information or the ability to browse the internet. My training only goes up until September 2021. Therefore, I cannot provide you with the most recent election results. I recommend checking a reliable news source or conducting an online search to find the latest information about the first round of presidential elections in Turkey.

The screenshot shows a chat interface with a user question and an AI response. The user question is "Who won the first round of presidential elections in Turkey last Sunday?". The AI response is "I'm sorry, but as an AI language model, I don't have access to real-time information or the ability to browse the internet. My training only goes up until September 2021. Therefore, I cannot provide you with the most recent election results. I recommend checking a reliable news source or conducting an online search to find the latest information about the first round of presidential elections in Turkey." Annotations include a blue box around "presidential elections in Turkey" in the user question, an orange box around "last Sunday?" in the user question, a blue box around "presidential elections in Turkey" in the AI response, and an orange box around "My training only goes up until September 2021" in the AI response.

What are challenges of answering this question?

- Can we update the model's knowledge without updating its parameters?

Limitations of PLMs (& LLMs)

- **Hallucination** problem (e.g., factual errors)
- Long-tail knowledge (e.g., domain-specific) may not be well-represented in the model's pretraining corpus
- Cannot easily expand or update parameters after pretraining
 - knowledge learned during pretraining is **static**
- Source of information is **non-attributable**

Factual-heavy NLP tasks

Fact Verification

Claim: The Rodney King riots took place in the most populous county in the USA.

[wiki/Los Angeles Riots]

The 1992 Los Angeles riots, also known as the Rodney King riots were a series of riots, lootings, arsons, and civil disturbances that occurred in Los Angeles County, California in April and May 1992.

[wiki/Los Angeles County]

Los Angeles County, officially the County of Los Angeles, is the most populous county in the USA.

Verdict: Supported

FEVER dataset

Factual Question Answering

Huguenot numbers peaked near an estimated two million by 1562, concentrated mainly in the southern and central parts of France, about one-eighth the number of French Catholics. As Huguenots gained influence and more openly displayed their faith, Catholic hostility grew, in spite of increasingly liberal political concessions and edicts of toleration from the French crown. A series of religious conflicts followed, known as the Wars of Religion, fought intermittently from 1562 to 1598. The wars finally ended with the granting of the Edict of Nantes, which granted the Huguenots substantial religious, political and military autonomy.

Where was France's Huguenot population largely centered?

Ground Truth Answers: the southern and central parts of France southern and central parts of France, about one-eighth

SQuAD.v2 dataset

How can we tackle these limitations?

Retrieval

- Precise knowledge access mechanism
- Easy update to known knowledge (update the retrieval knowledge base)
- Neural Retrieval starting to outperform traditional IR

Limitations

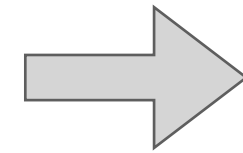
- Needs supervision or “heuristics”
- Task specific way to integrate into downstream tasks

Today's Outline

- **Lecture: Retrieval Augmented LMs**
 - **Models:** Model types, training objectives, different external knowledge
 - **Downstream tasks**
 - **Augmented LLMs:** Retrieval in the LLM era
 - **Augmentation benefits:** Modularity, Attribution, Parameter efficiency

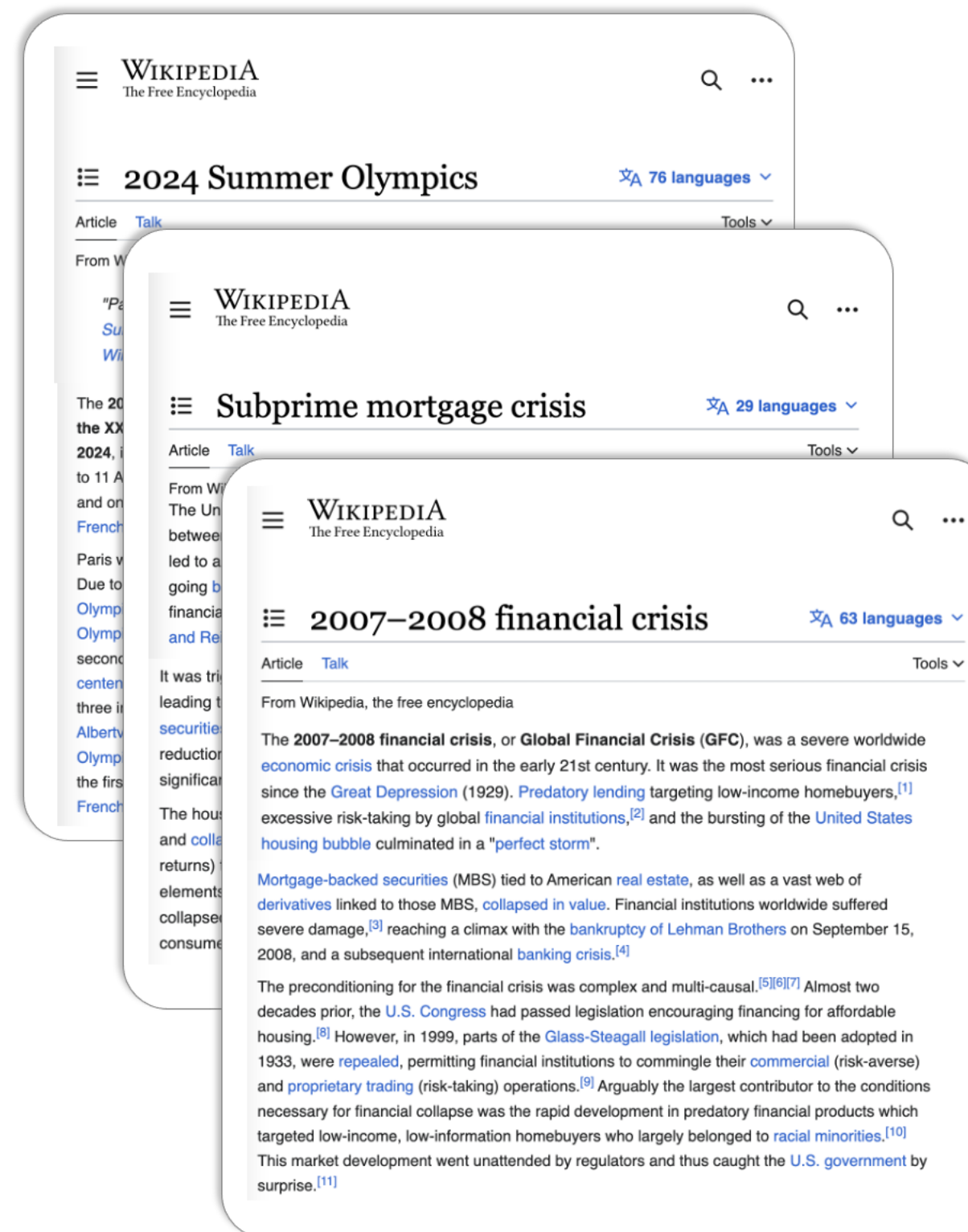
Finding the answer in 21M documents

Query



Documents

"Where the financial crisis of 2008 started?"



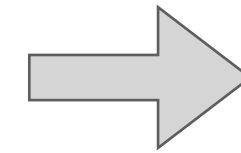
**Why can't we do this with closed-book
Extractive & Generative QA models?**

We don't have a context!

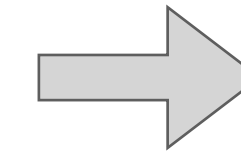
**But we can get one!
(With adaptations)**

Finding the answer in 21M documents

Query



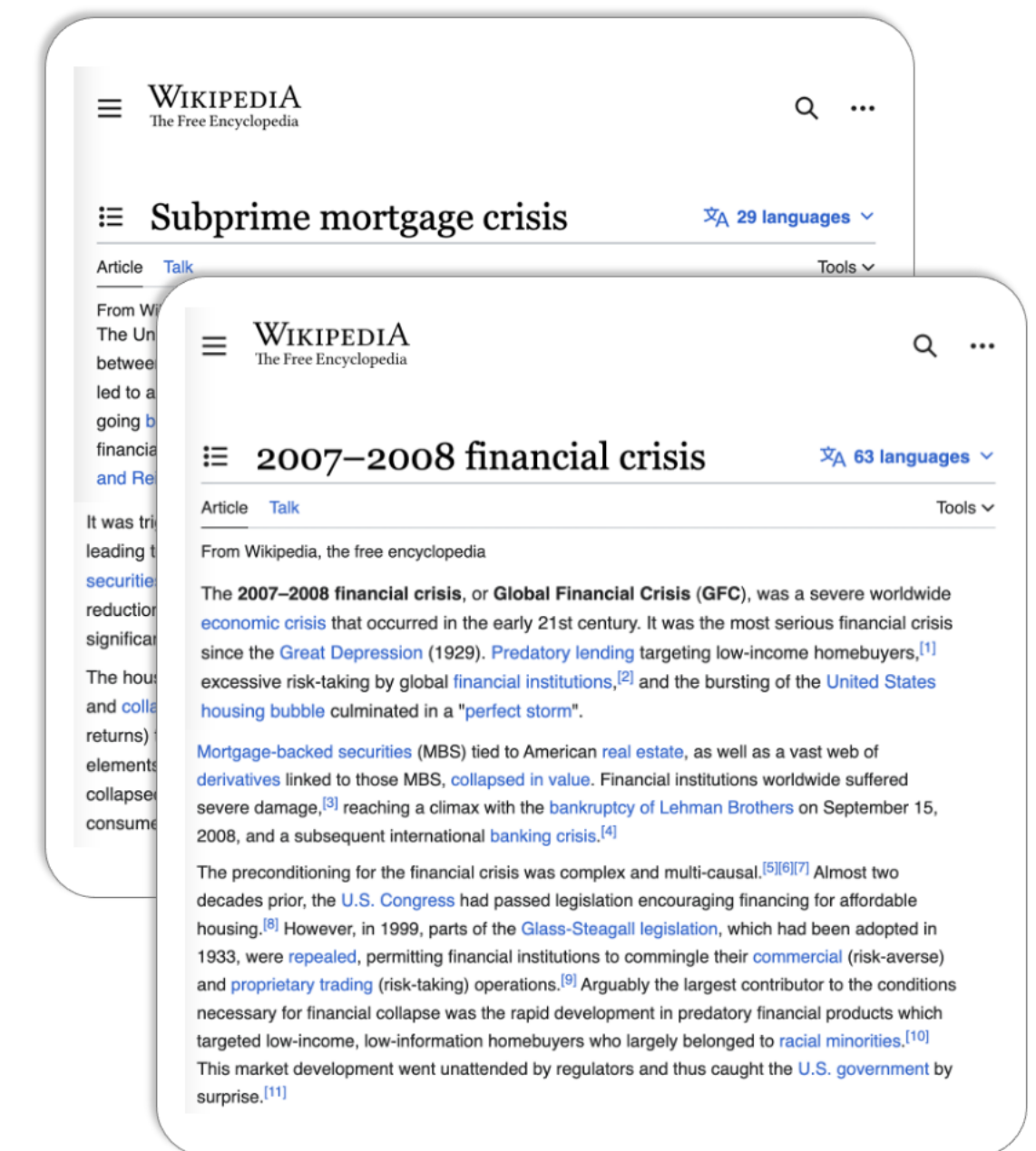
Documents



Retrieve relevant documents

That might contain the answer

"Where the financial crisis of 2008 started?"



Classical Retrieval: Okapi BM25

Repetitions of query words \rightarrow good

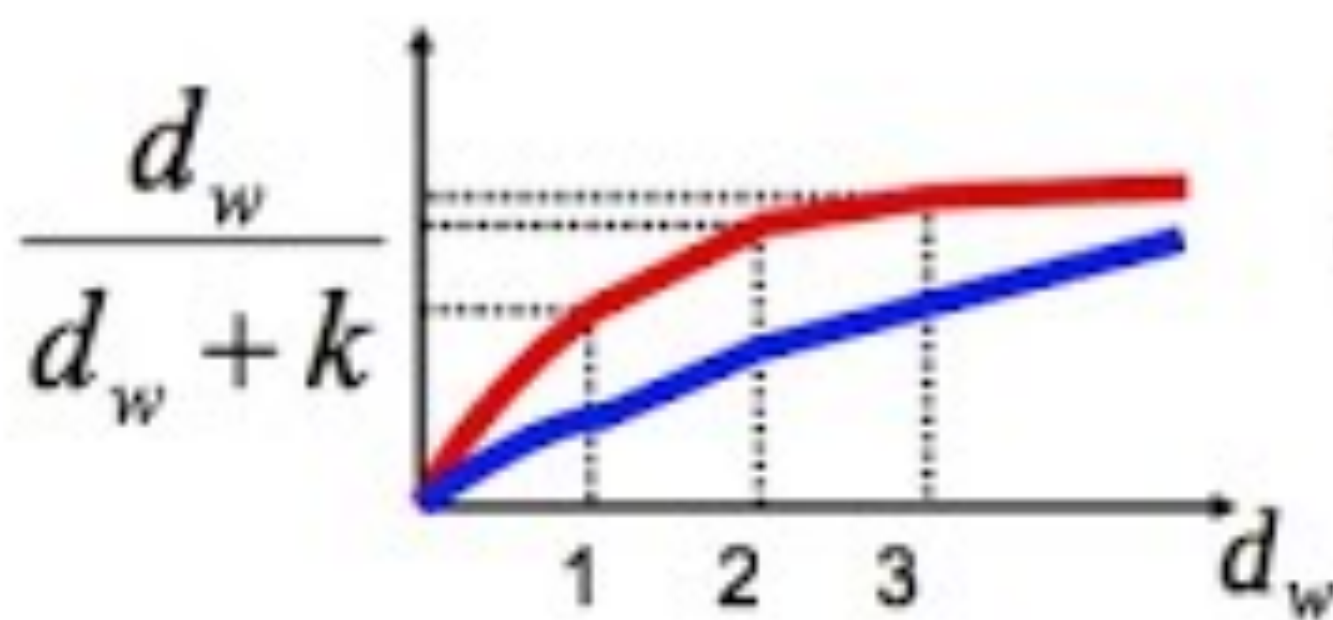
Common words less important

$$\log \frac{P(D | R=1)}{P(D | R=0)} \approx \sum_w \left(\frac{d_w(1+k)}{d_w + k((1-b) + \frac{b \cdot dl}{avg \cdot dl})} \cdot \log \frac{N - N_w + \frac{1}{2}}{N_w + \frac{1}{2}} \right)$$

More words in common with the query \rightarrow good

Repetitions less important than different query words

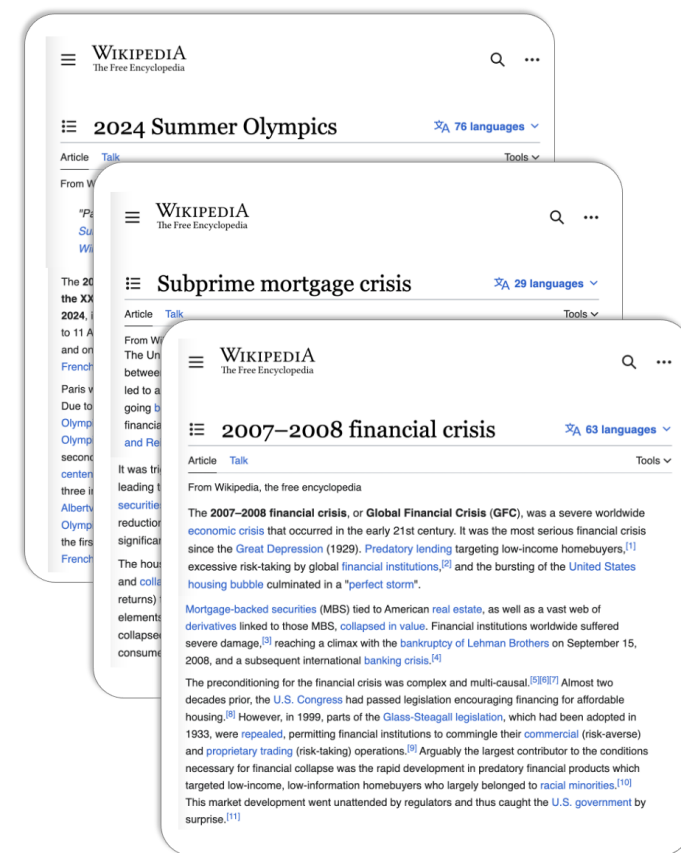
But more important if document is relatively long (wrt. average)



BM25 still works quite well for many applications, so don't ignore classical retrieval if it works better !

Dense Passage Retrieval (DPR)

Documents



Query

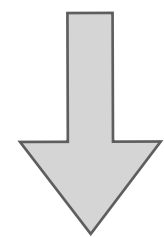
"Where the financial crisis of 2008 started?"

Dense
Embedding
Model

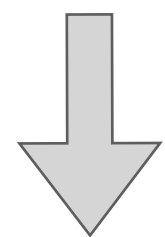
- Create the representations of documents
- Create the representation of the query
- Retrieve k documents vectors based on the query vector

Dense Passage Retrieval (DPR)

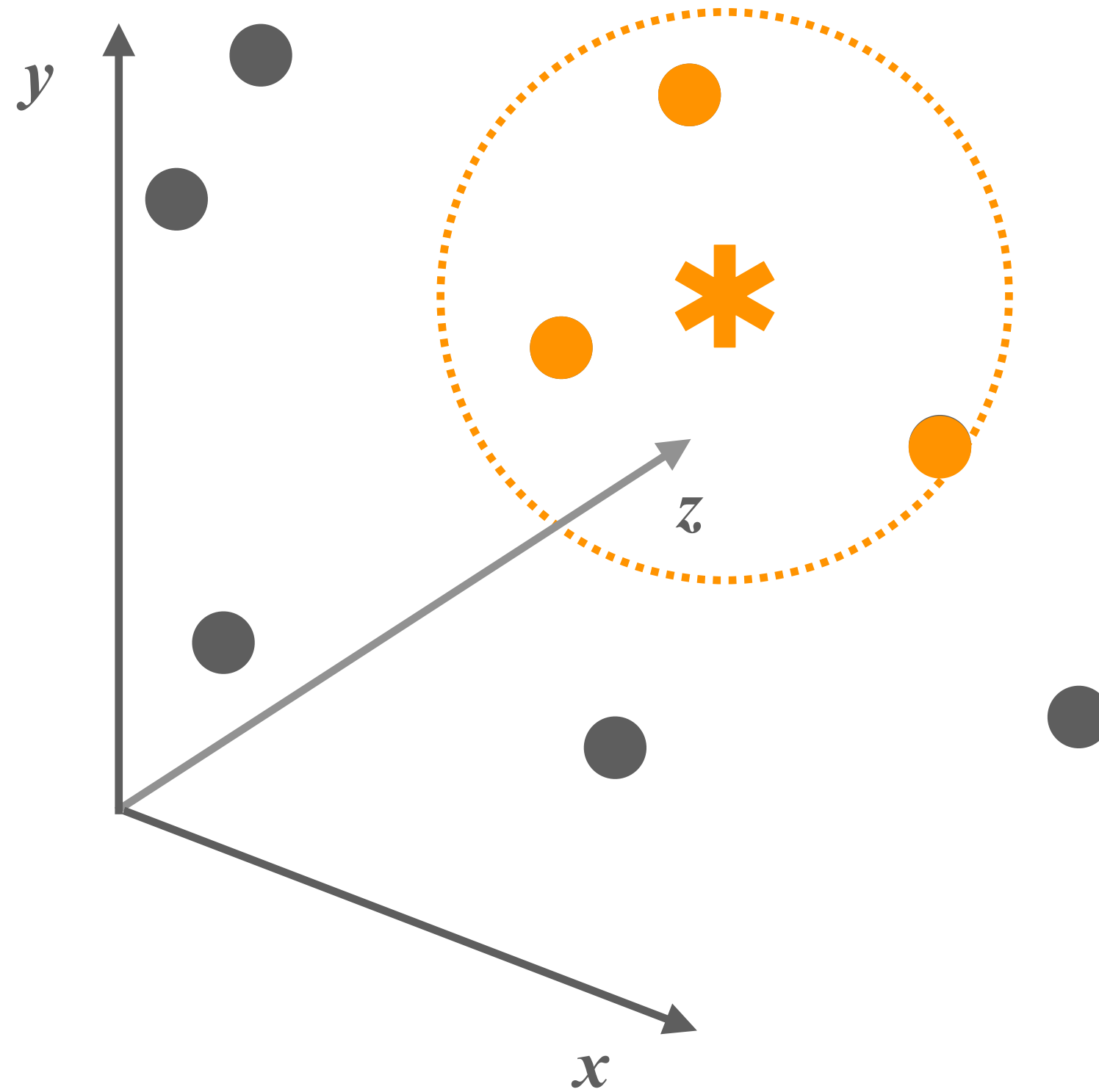
Documents



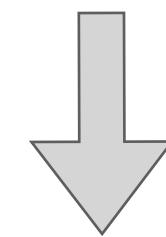
$E_P(\cdot)$



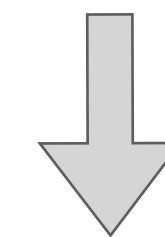
$[-0.5968882, -0.33086956, -0.32643065, -0.3670732, \dots]$



Query



$E_Q(\cdot)$



$[-0.3692328, -0.37902787, -0.12308089, -0.38124698, \dots]$

$$\text{sim}(q, p) = E_Q(q)^\top E_P(p)$$

Training DPR

"Where the financial crisis of 2008 started?"

How to create a Document-Query vector space?

Goal: **Relevant** pairs of questions-passages will have a smaller distance compared to the **irrelevant** ones.

Positive passage p^+



Negative passage p^-



DPR LOSS FUNCTION

$$L(q_i, p_i^+, p_{i,1}^-, \dots, p_{i,n}^-) = -\log \frac{e^{\text{sim}(q_i, p_i^+)}}{e^{\text{sim}(q_i, p_i^+)} + \sum_{j=1}^n e^{\text{sim}(q_i, p_{i,j}^-)}}$$

**How we can integrate a neural retriever
into a Language Model?**

Retrieval-Augmented LMs

$$p(y | x) =$$

LM

Retriever

Auto-Encoder

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

$$p(z | x)$$

Auto-Encoder

Auto-Regressive

$$\sum_{z \in \mathcal{Z}} \prod_i^N p(y_i | x, z, y_{1:i-1})$$

Trained to retrieve relevant documents (optional)

z often represented as a latent variable
(may not know what the current document is)

Trained to produce the right answer given the input query and the retrieved documents.

Retrieval-Augmented LMs - Terminology

Information that is stored in the parameters of the models used (both for the LM and the retrieval parts).

memory

knowledge

The type of external source the retriever will use.

modalities

Implicit

Parametric

LM

vs

Explicit

Non-parametric

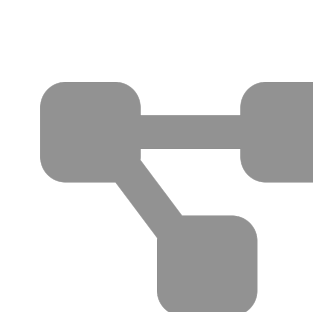
Retriever

vs

KB



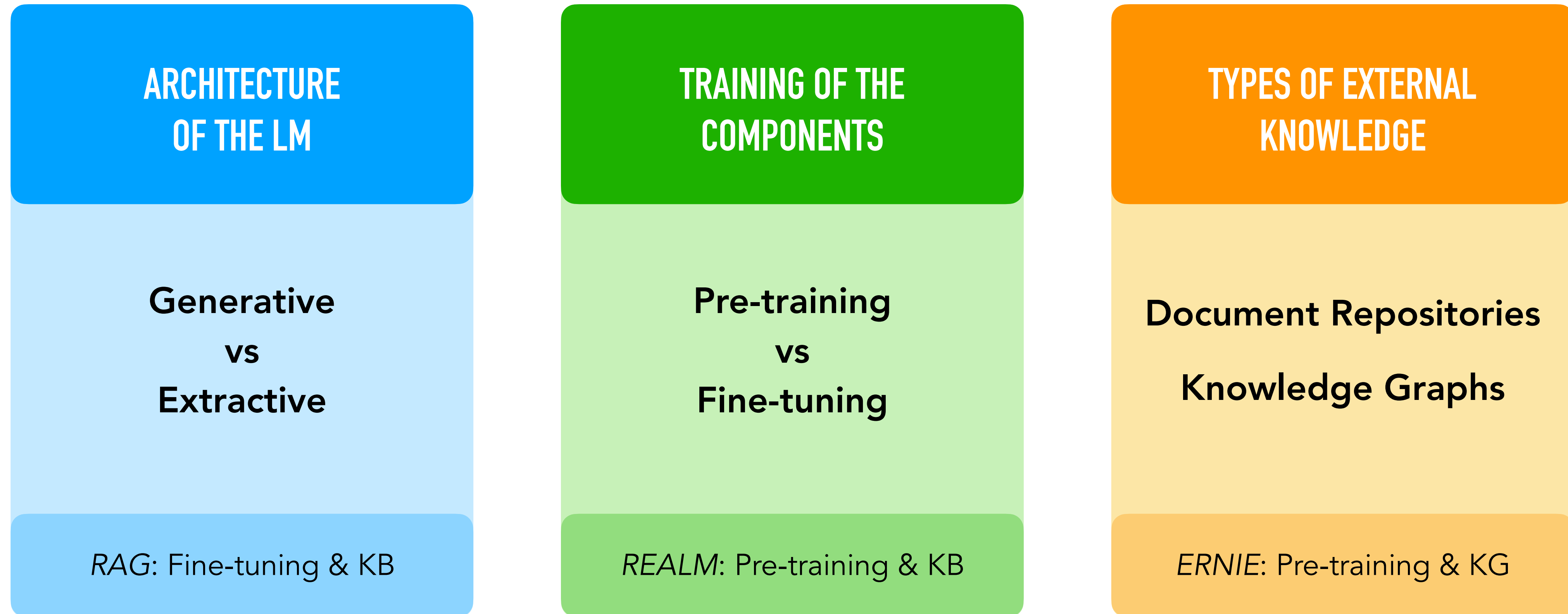
KG



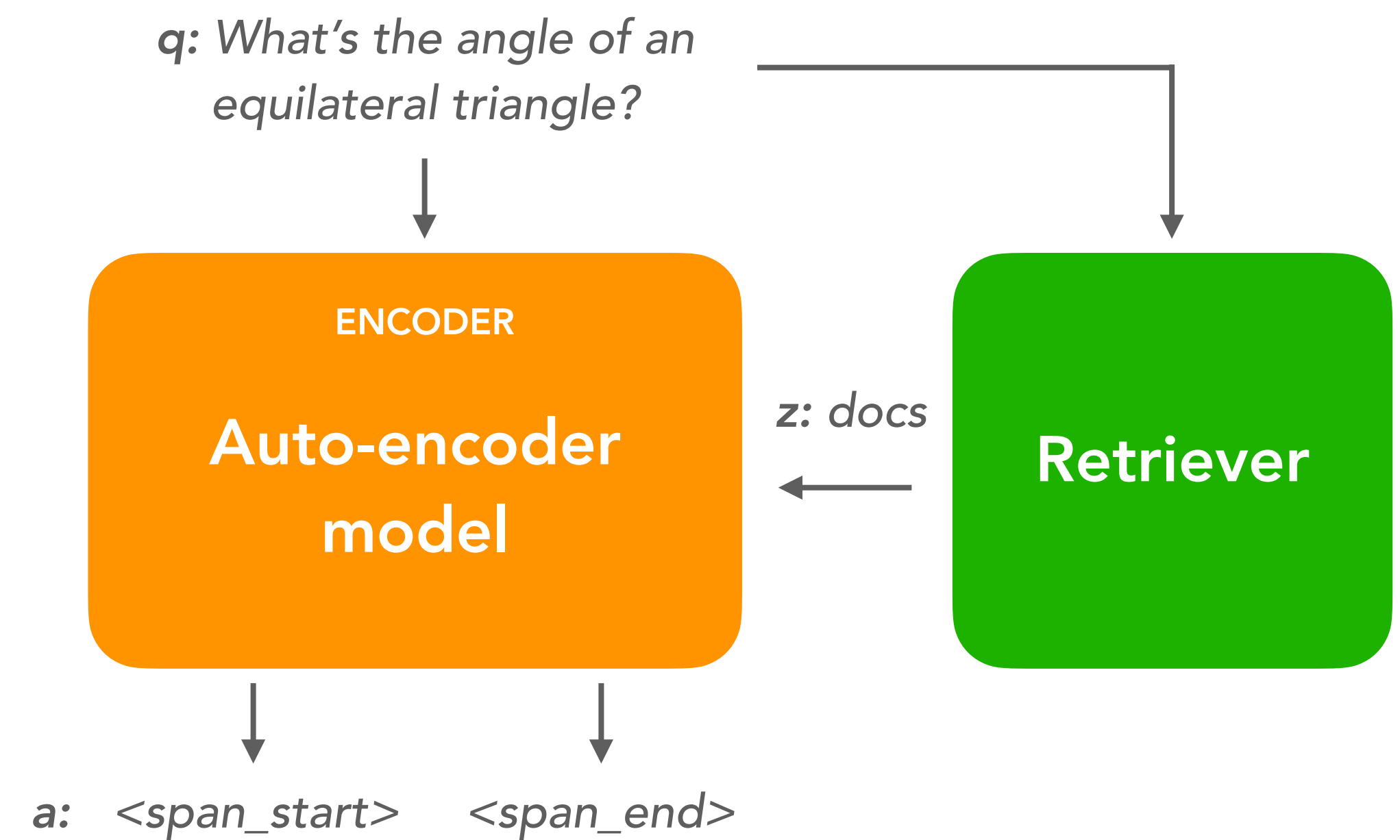
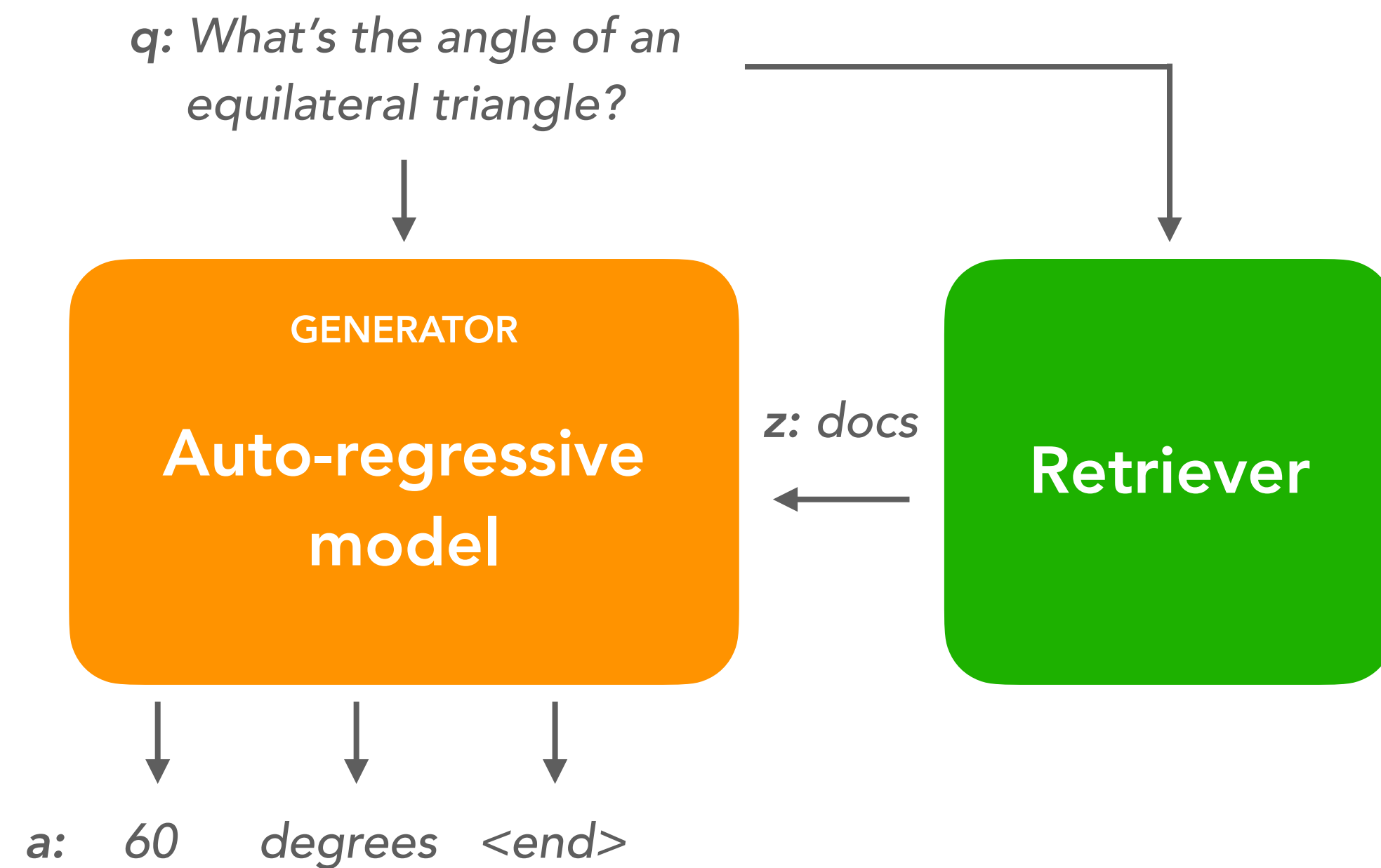
Tools



The landscape of Retrieval-Augmented LMs



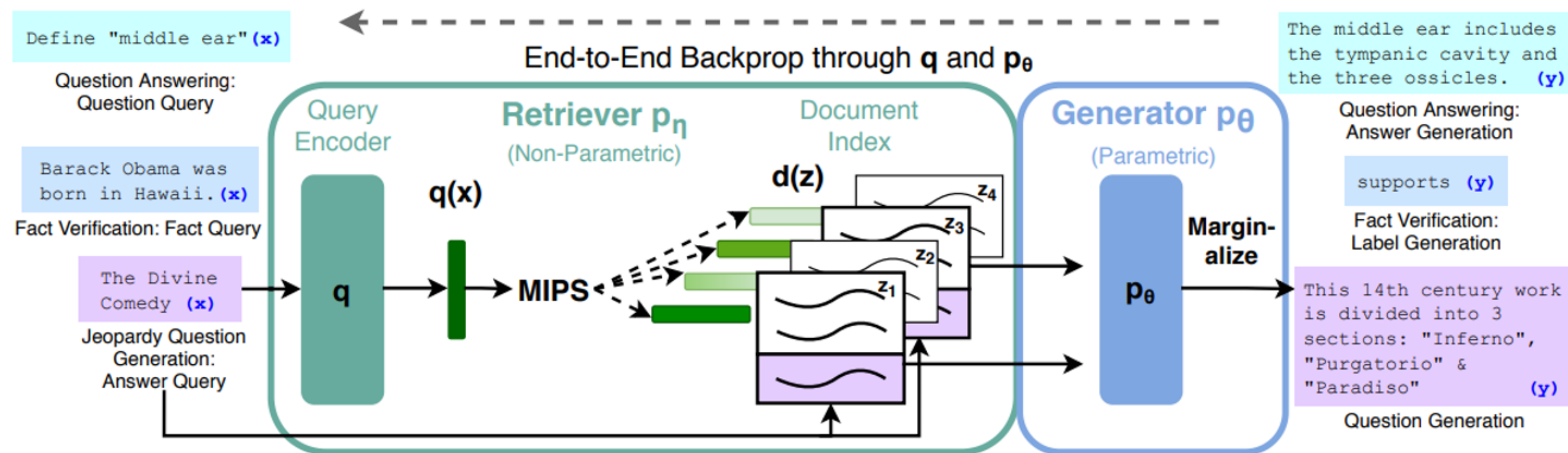
Generative vs Extractive



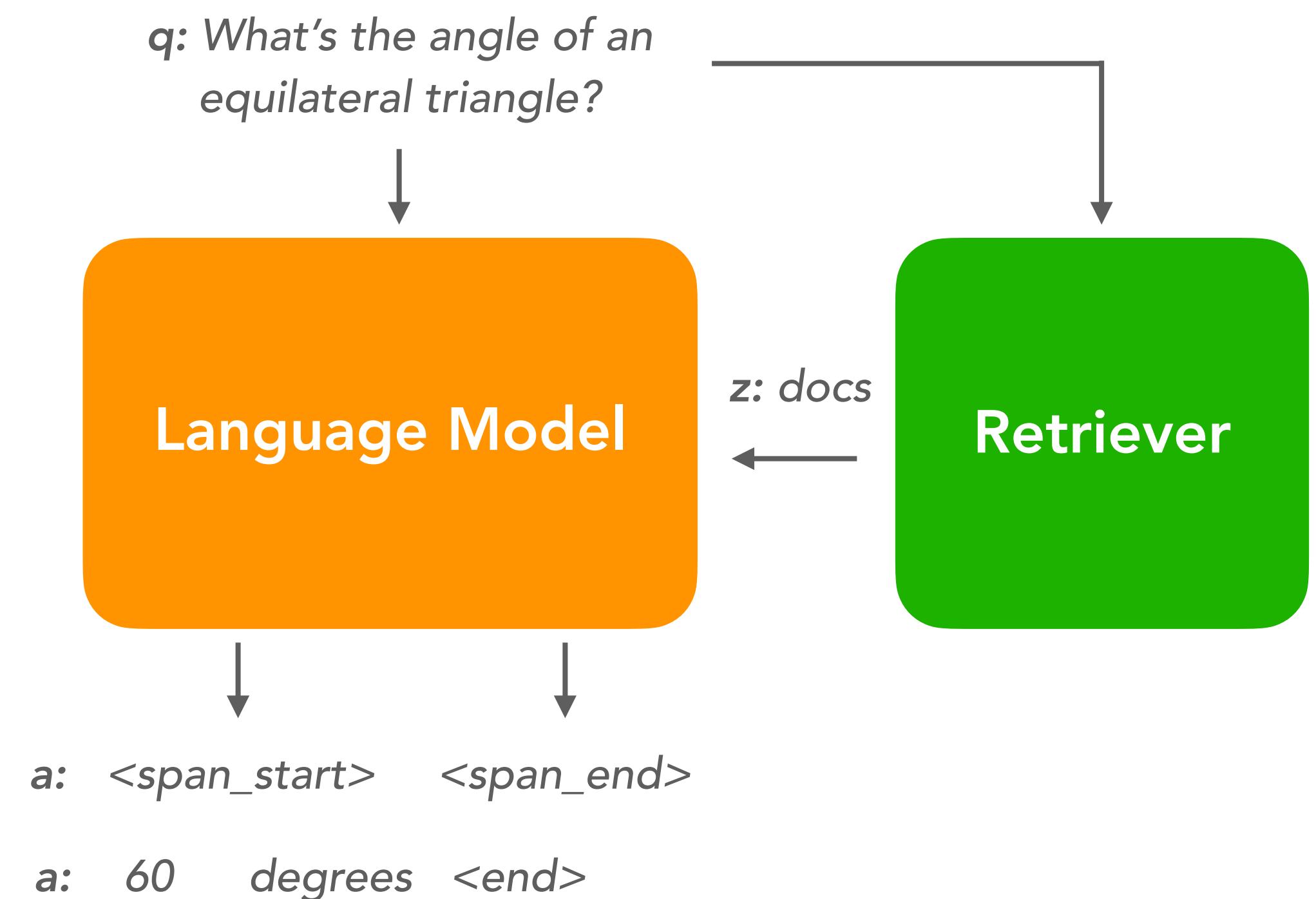
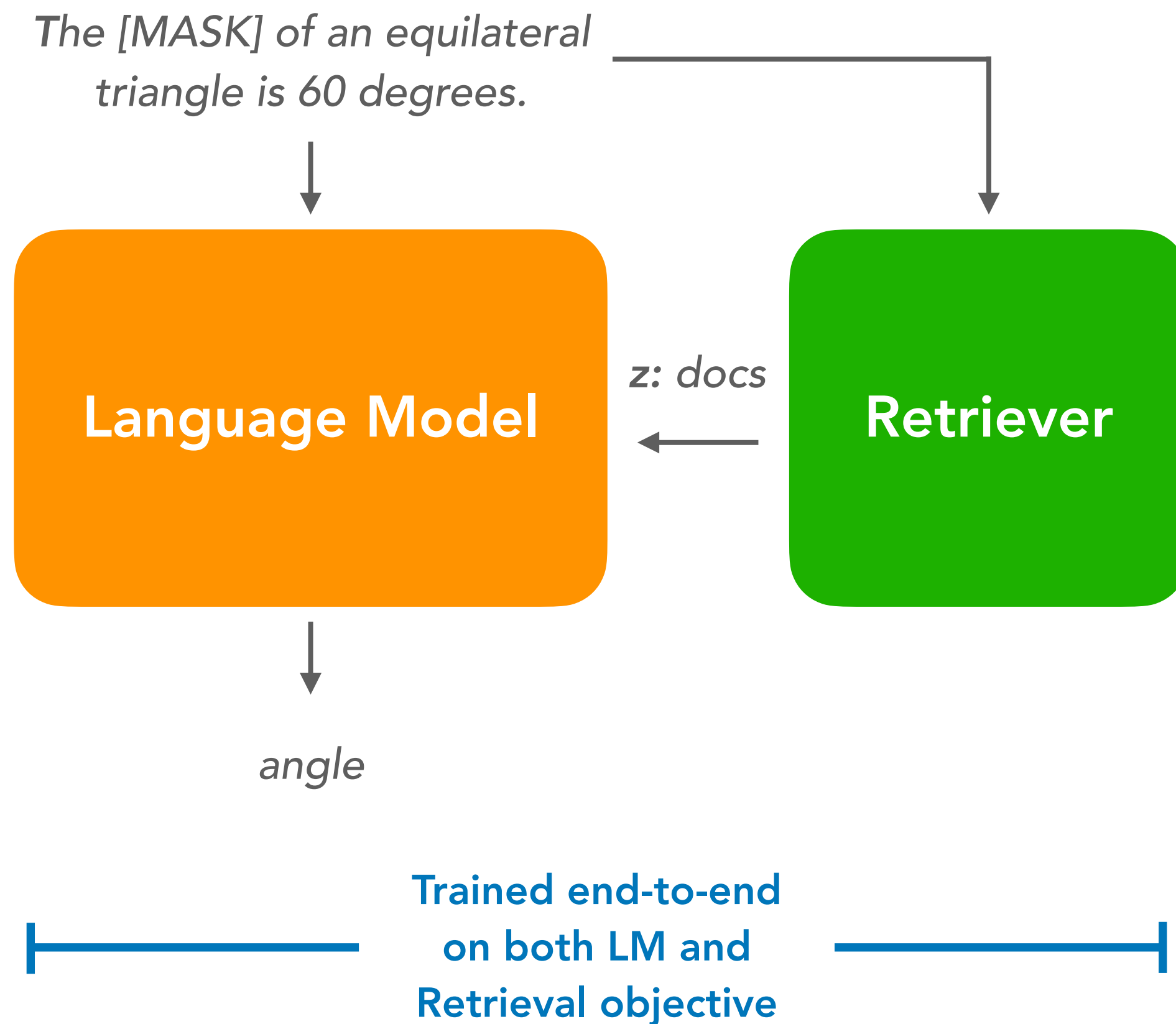
RAG: Generative Retrieval-Augmented LM

1. Pre-trained generator (e.g. BART)
2. Pre-trained retriever (e.g. DPR)
3. Indexed KB of text documents (e.g. Wikipedia)

$$p(y | x) = \sum_{z \in \mathcal{Z}} \prod_i^N p(y_i | x, z, y_{1:i-1})$$



Pre-training vs Fine-tuning



REALM: Pre-training Retrieval Augmented LMs

First Retrieve:

The retriever model is trained on what documents are relevant.

Goal: Penalise uninformative retrievals

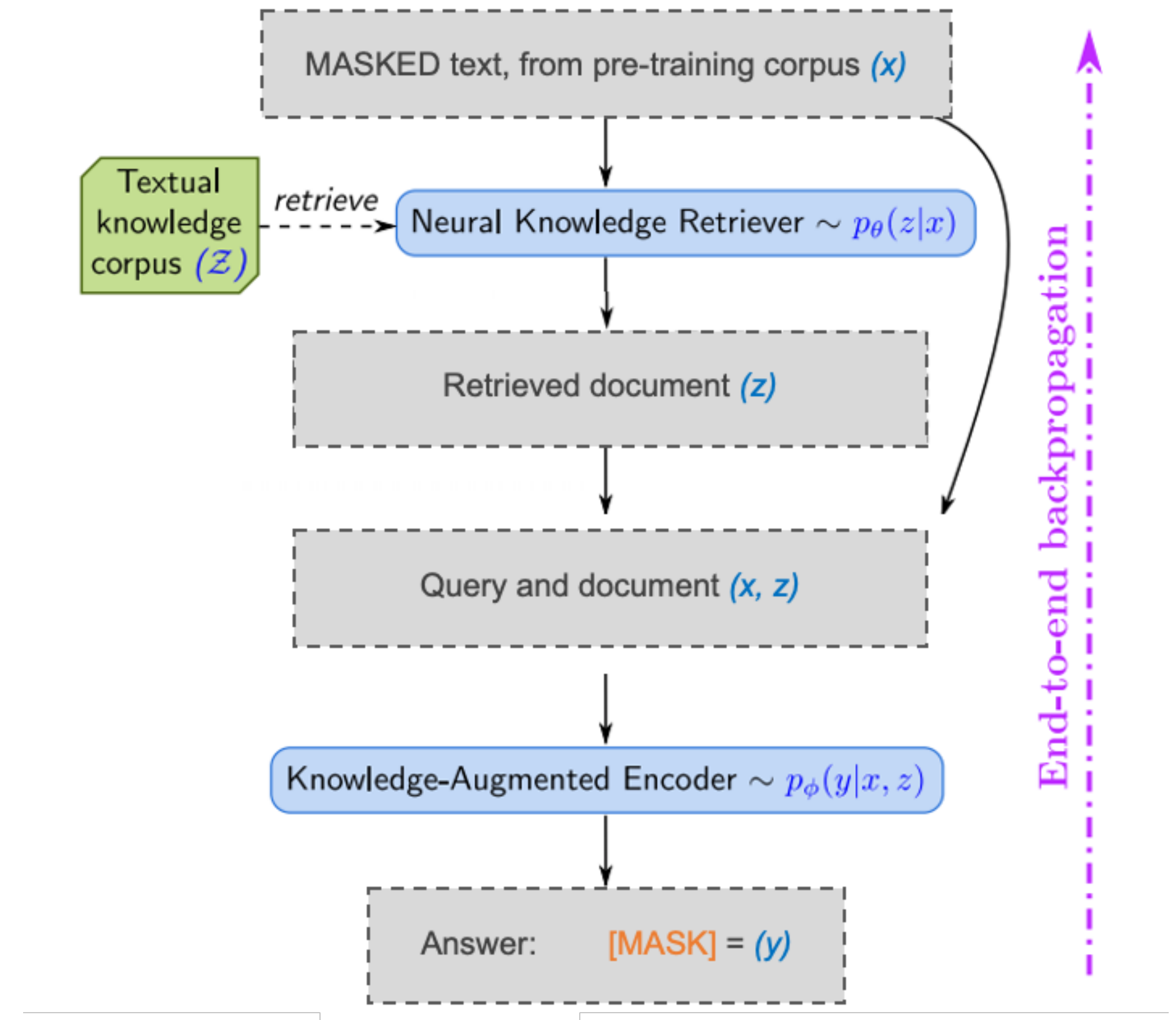
Then Predict:

The encoder model is trained to predict the original value of each masked token by attending to the input query and the retrieved documents.

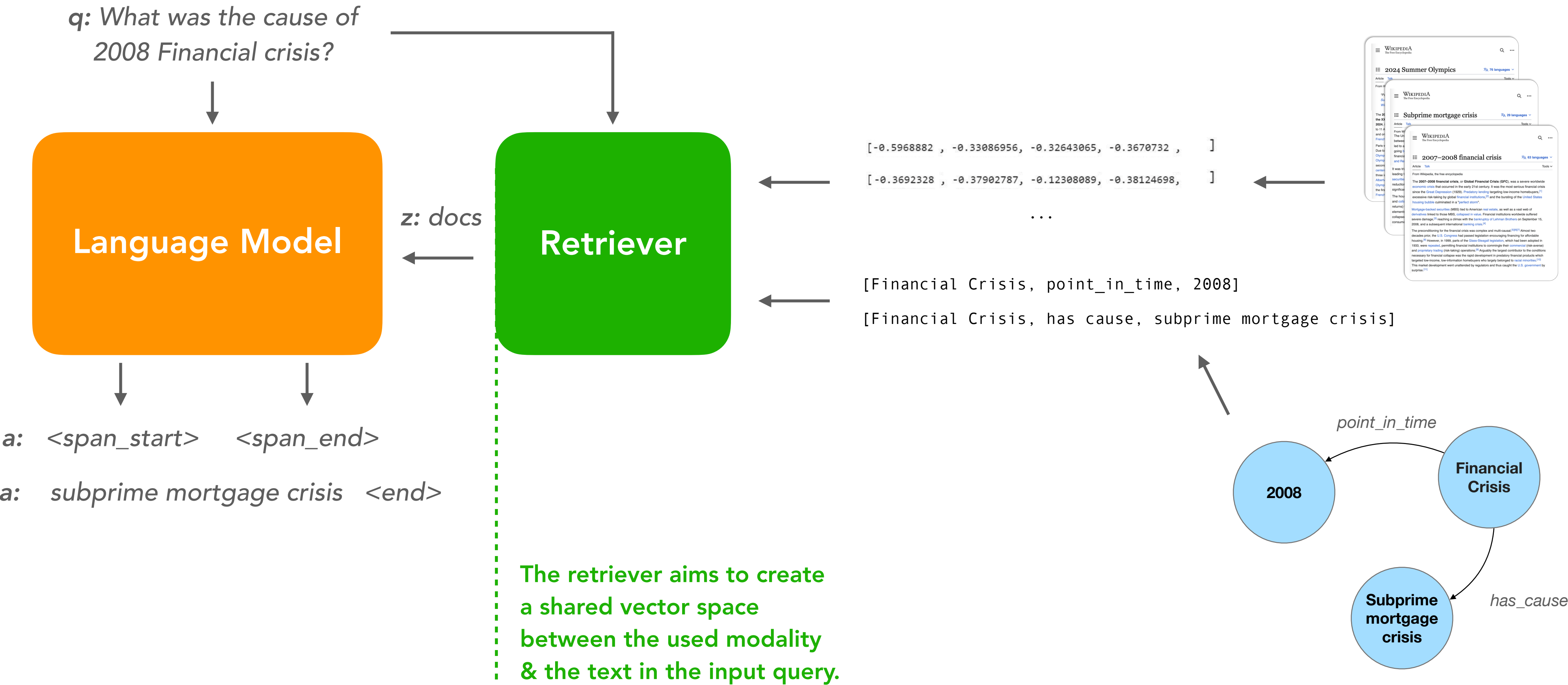
Goal: Minimise perplexity

Benefits of pre-training end-to-end

- Transferability across tasks
- Rely on information beyond lexical overlap: the model learns for itself which texts are most useful for reducing perplexity.
- Model-centric **unsupervised alignments** between text in the pre-training corpus X and knowledge corpus Z .



Different types of external knowledge

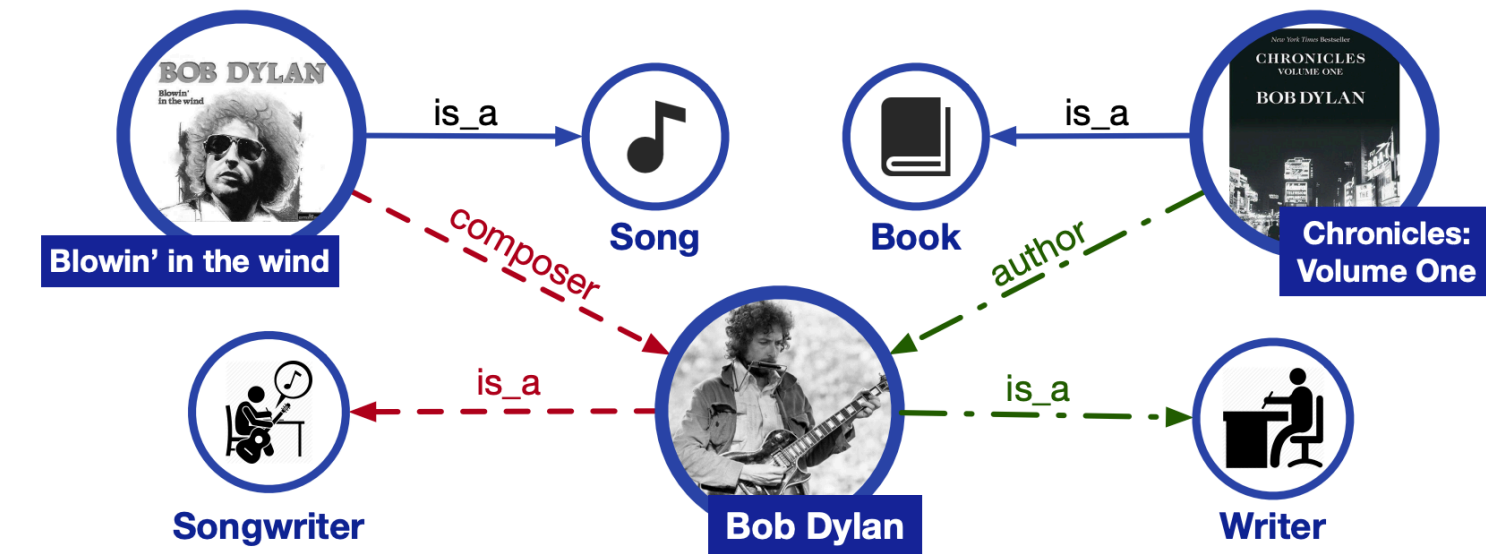


What other modalities could we use as a base ?

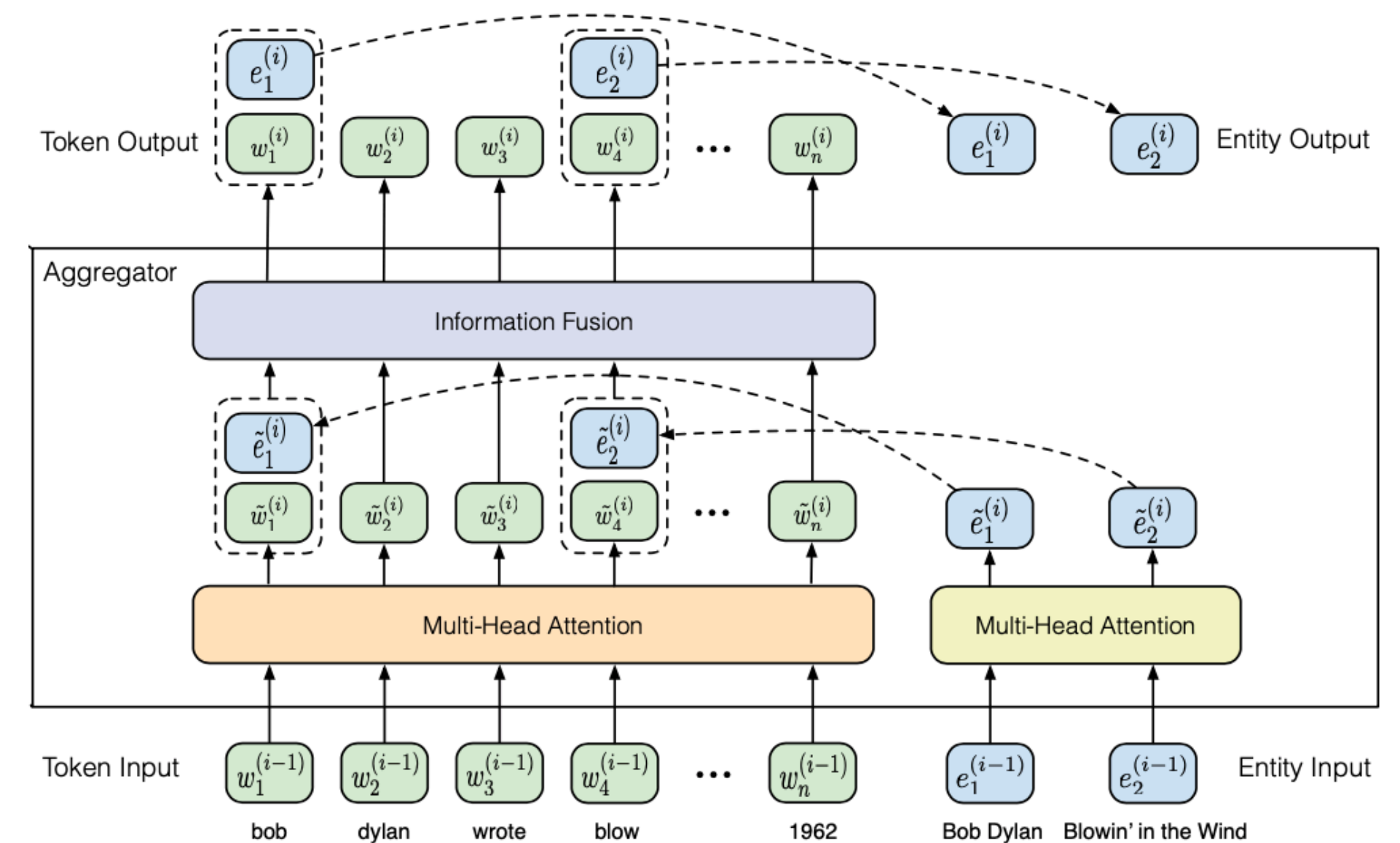
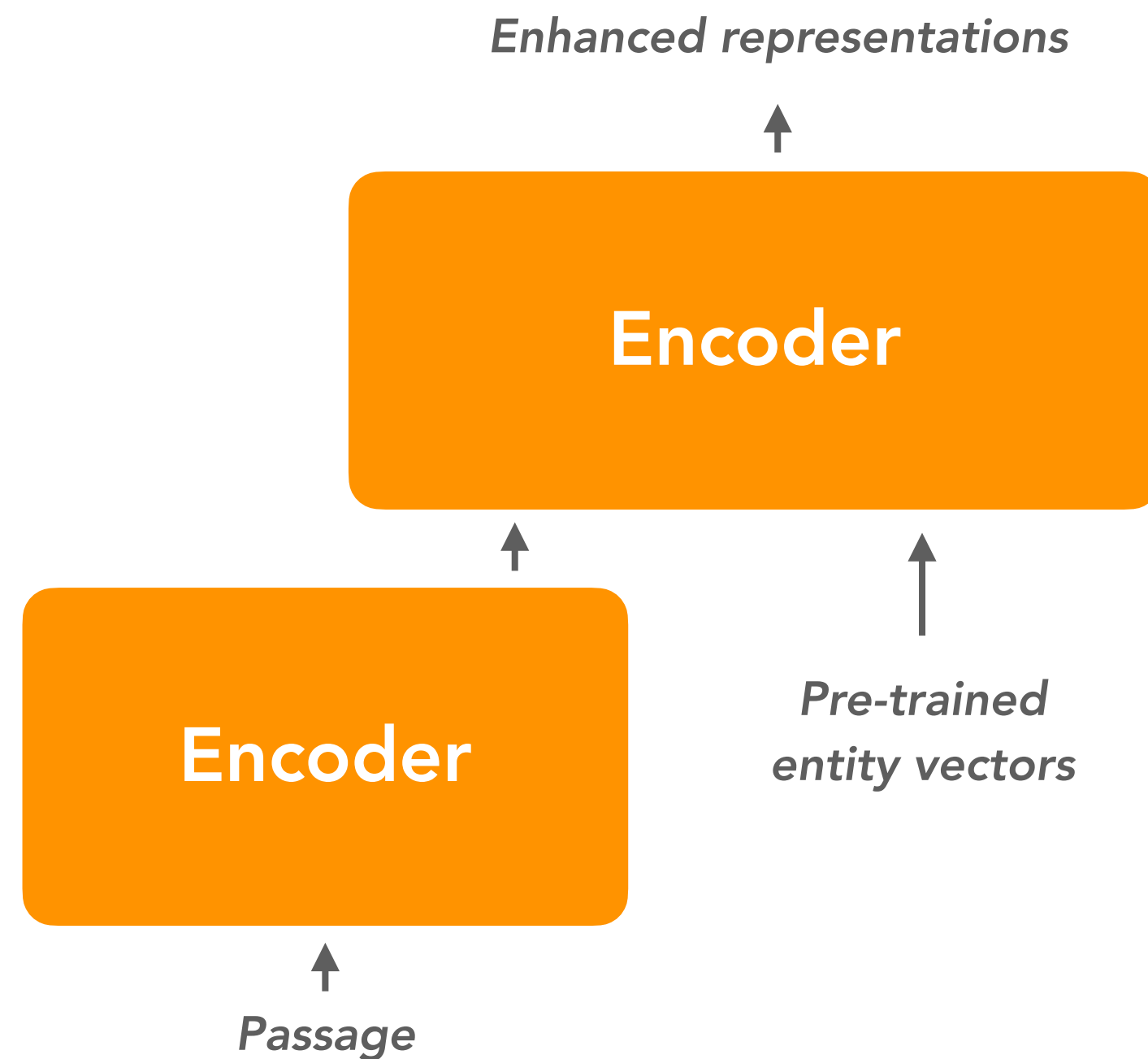
How would we integrate these modalities ?

ERNIE: Infuse KG knowledge

1. Extracts the named entity mentions in the text
2. Aligns these mentions to their corresponding entities in KGs.
3. Gets the graph pre-trained entity embeddings for each named entity.
4. Integrates the entity representations in the Encoder model.



Bob Dylan wrote **Blowin' in the Wind** in 1962, and wrote **Chronicles: Volume One** in 2004.



Bob Dylan wrote **Blowin' in the Wind** in 1962

We looked at retrieving information from document bases and knowledge graphs.

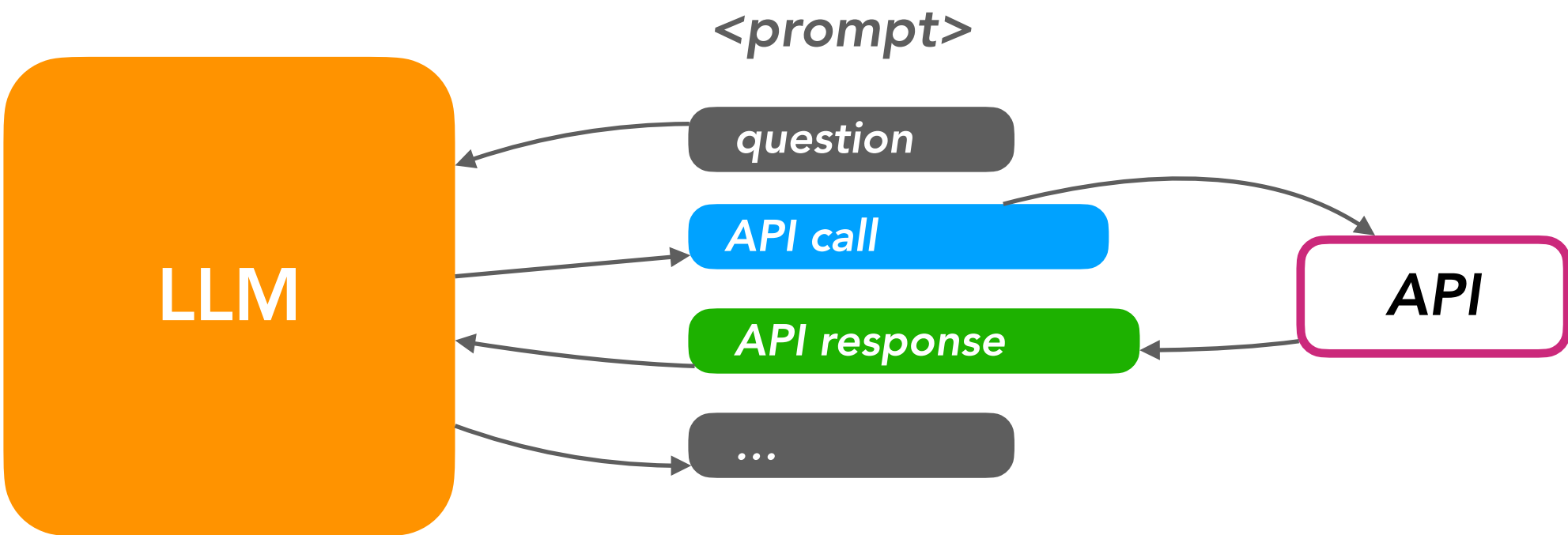
What else can be retrieved by a model ?

Augmented LLMs

Retrieval-Augmented Prompts

Retrieve from tools & APIs

Equip language models with the ability to use different tools by means of API calls



Act

Question	What is the elevation range for the area that the eastern sector of the Colorado orogeny extends into?
Action 1	Search[Colorado orogeny]
Observation 1	The Colorado orogeny was an episode of mountain building (an orogeny) in Colorado and surrounding areas.
Action 2	Lookup[eastern sector]
Observation 2	(Result 1 / 1) The eastern sector extends into the High Plains and is called the Central Plains orogeny.
Action 3	Search[High Plains]
Observation 3	High Plains refers to one of two distinct land regions:
Action 4	Search[High Plains (United States)]
Observation 4	The High Plains are a subregion of the Great Plains. From east to west, the High Plains rise in elevation from around 1,800 to 7,000 ft (550 to 2,130 m) .[3]
Action 5	Finish[1,800 to 7,000 ft]

ReAct (Yao et al. 2024)

The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

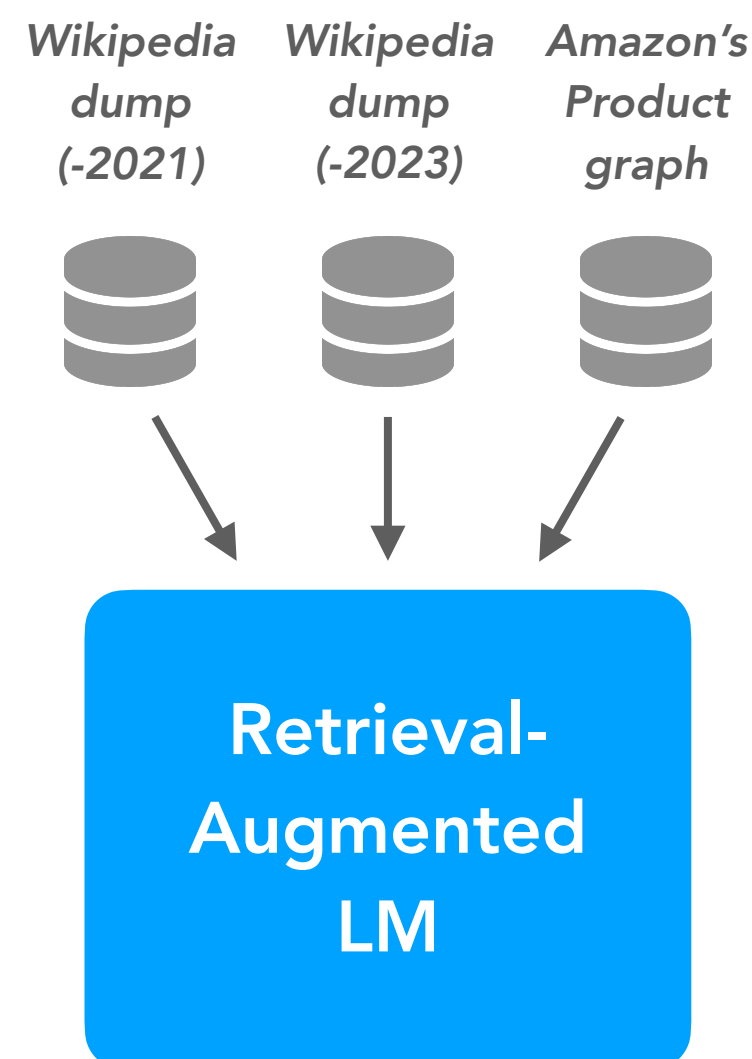
Toolformer (Schick et al. 2023)

What are some benefits of augmentation?

Additional benefits of Augmented LMs

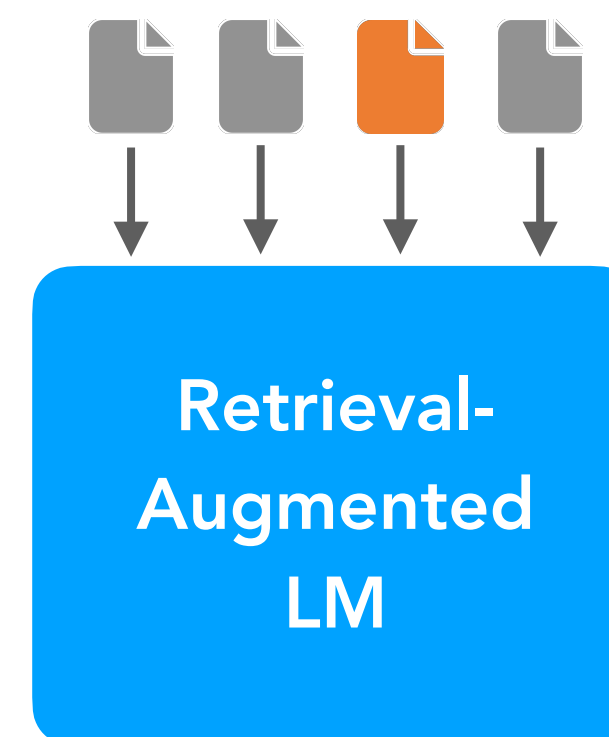
Modularity

We can change external memory and update the model's knowledge on test time.



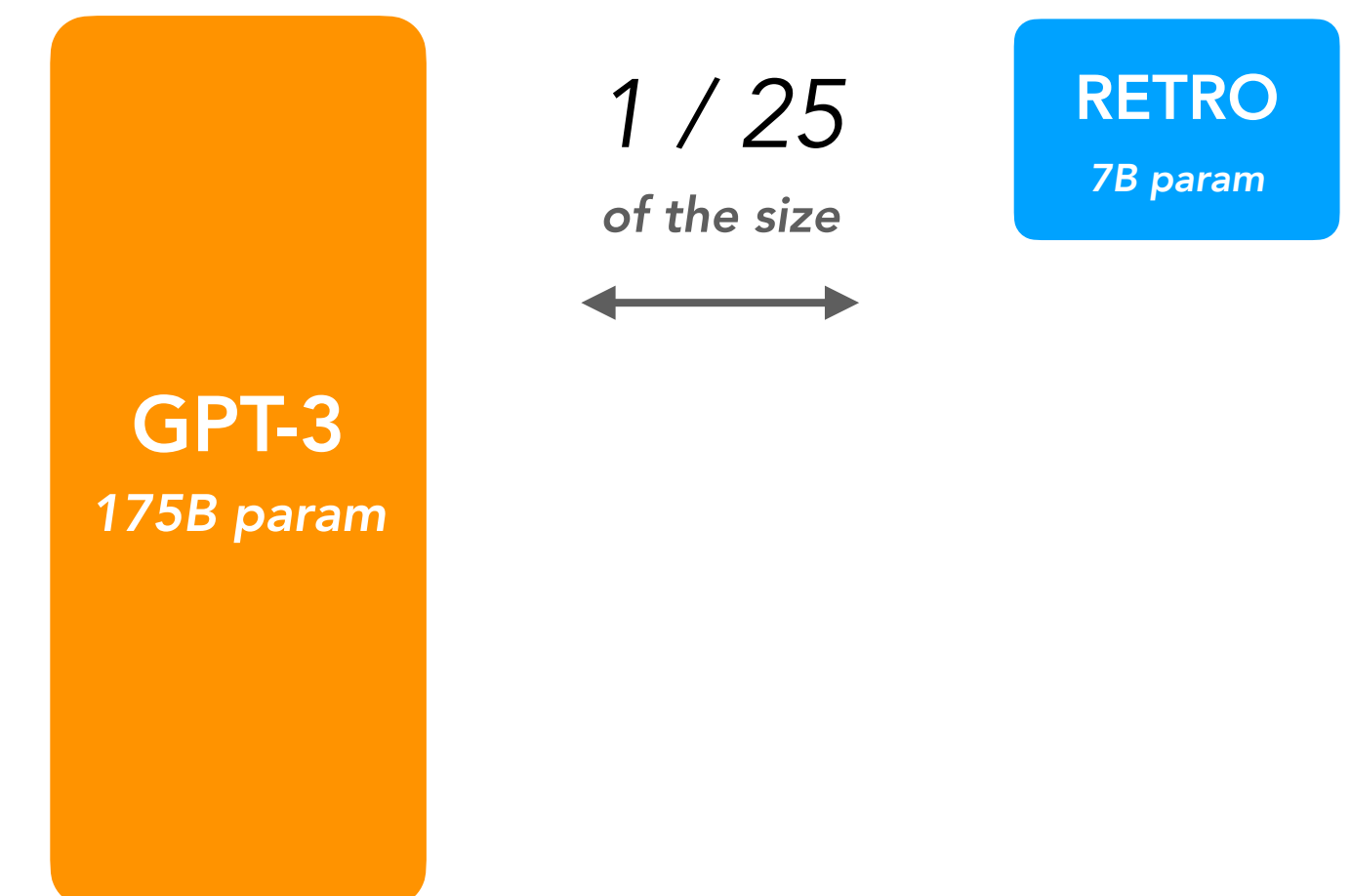
Attribution

We can trace back the information (documents) that the generated answer is based on.



Parameter efficiency

We can leverage external memory to reduce the number of implicit parameters of the LM without compromising performance.



Closed-book vs. Retrieval

Rank	Model	EM	Paper	Code	Result	Year	Tags
1	Atlas (full, Wiki-dec-2018 index)	64.0	Atlas: Few-shot Learning with Retrieval Augmented Language Models	Code	Result	2022	
2	Atlas (full, Wiki-dec-2021+CC index)	60.4	Atlas: Few-shot Learning with Retrieval Augmented Language Models	Code	Result	2022	
3	FiE	58.4	FiE: Building a Global Probability Space by Leveraging Early Fusion in Encoder for Open-Domain Question Answering		Result	2022	
4	R2-D2 (full)	55.9	R2-D2: A Modular Baseline for Open-Domain Question Answering	Code	Result	2021	
5	ReAtt	54.7	Retrieval as Attention: End-to-end Learning of Retrieval and Reading within a Single Transformer	Code	Result	2022	
6	FiD-KD (full)	54.7	Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering	Code	Result	2020	
7	EMDR^2	52.5	End-to-End Training of Multi-Document Reader and Retriever for Open-Domain Question Answering	Code	Result	2021	
8	FID (full)	51.4	Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering	Code	Result	2020	
9	RETRO + DPR (full)	45.5	Improving language models by retrieving from trillions of tokens	Code	Result	2021	
10	code-davinci-002 175B + REPLUG LSR (few-shot)	45.5	REPLUG: Retrieval-Augmented Black-Box Language Models	Code	Result	2023	
11	Atlas (few-shot, k=64, Wiki-Dec-2018 index)	45.1	Atlas: Few-shot Learning with Retrieval Augmented Language Models	Code	Result	2022	few-shot
12	code-davinci-002 175B + REPLUG (few-shot)	44.7	REPLUG: Retrieval-Augmented Black-Box Language Models	Code	Result	2023	
13	RAG	44.5	Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks	Code	Result	2020	
14	Blended RAG	42.63	Blended RAG: Improving RAG (Retriever-Augmented Generation) Accuracy with Semantic Search and Hybrid Query-Based Retrievers	Code	Result	2024	
15	Atlas (few-shot, k=64, Wiki-dec-2021+CC index)	42.4	Atlas: Few-shot Learning with Retrieval Augmented Language Models	Code	Result	2022	few-shot
16	DPR	41.5	Dense Passage Retrieval for Open-Domain Question Answering	Code	Result	2020	
17	REALM	40.4	REALM: Retrieval-Augmented Language Model Pre-Training	Code	Result	2020	
18	LLaMA 65B (few-shot, k=64)	39.9	LLaMA: Open and Efficient Foundation Language Models	Code	Result	2023	

- All leading models use retrieval
- Retrieval models often have to the order $O(10^9)$ parameters
- Much more efficient than largest LLMs

Recap

- **Retrieval-Augmented language models:**
 - Infuse knowledge from external sources into LMs.
 - Suitable for knowledge-intensive tasks where factual accuracy is needed.
- **Main components:** type of external knowledge, type of the LM, type of training.
- Using external knowledge can allow us to reduce the # of parameters of LMs, making them smaller in size without compromising performance.
- **In the LLMs era:**
 - Retrieval aims to augment the prompt.
 - Models are interacting with various tools and APIs to enhance their reasoning capabilities.

References

- Karpukhin, Vladimir, et al. "Dense passage retrieval for open-domain question answering." arXiv preprint arXiv:2004.04906 (2020).
- Guu, Kelvin, et al. "Retrieval augmented language model pre-training." International conference on machine learning. PMLR, 2020.
- Lewis, Patrick, et al. "Retrieval-augmented generation for knowledge-intensive nlp tasks." Advances in Neural Information Processing Systems 33 (2020): 9459-9474.
- Schick, Timo, et al. "Toolformer: Language models can teach themselves to use tools." arXiv preprint arXiv:2302.04761 (2023).
- Yao, Shunyu, et al. "React: Synergizing reasoning and acting in language models." arXiv preprint arXiv:2210.03629 (2022).
- Mialon, Grégoire, et al. "Augmented language models: a survey." arXiv preprint arXiv:2302.07842 (2023).