

Retrieval Augmented Generation

Lecture: Information Retrieval 2

29-11-2024

Antonis Krasakis

Disclaimer

Slides are largely-based on the excellent ACL 2023 tutorial:

Retrieval-based Language Models and Applications

Akari Asai, Sewon Min, Zexuan Zhong, Danqi Chen

<https://acl2023-retrieval-lm.github.io/>

What?

Retrieval-based language models (LMs)

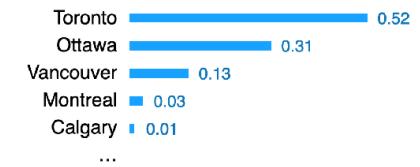
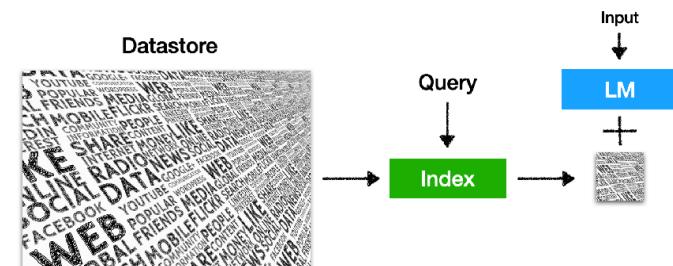
Retrieval-based LMs = Retrieval + LMs

- It is a **language model** $P(x_n | x_1, x_2, \dots, x_{n-1})$

The capital city of Ontario is _____

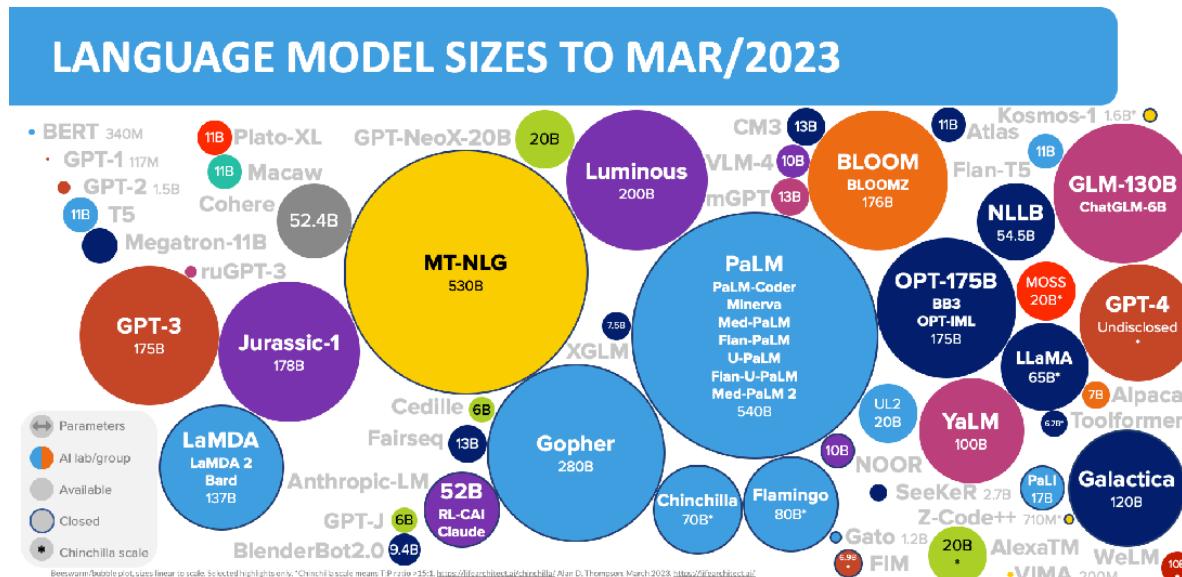
(can be broadly extended to masked language models or encoder-decoder models)

- It retrieves from an **external datastore** (at least during inference time)



(Also referred to as semiparametric and non-parametric models)

The age of large language models (LLMs)



- Transformers-based, **fully parametric**
- Trained on next-token prediction tasks (+ RLHF; not the focus today)
- **Model size ↑, data size↑**

Image: <https://lifearchitect.ai/models/>

When?

A range of target tasks

Question Answering

RETRO (Borgeaud et al., 2021)
REALM (Gu et al, 2020)
ATLAS (Izacard et al, 2023)

Fact verification

RAG (Lewis et al, 2020)
ATLAS (Izacard et al, 2022)
Evi. Generator (Asai et al, 2022)

Dialogue

BlenderBot3 (Shuster et al., 2022)
Internet-augmented generation
(Komeili et a., 2022)

Retrieval-based LMs have been extensively evaluated on knowledge-intensive tasks

A range of target tasks

Question answering

RETRO (Borgeaud et al., 2021)
REALM (Gu et al, 2020)
ATLAS (Izacard et al, 2023)

Fact verification

RAG (Lewis et al, 2020)
ATLAS (Izacard et al, 2022)
Evi. Generator (Asai et al, 2022)

Dialogue

BlenderBot3 (Shuster et al., 2022)
Internet-augmented generation
(Komeili et a., 2022)

Summarization

FLARE (Jiang et al, 2023)

Machine translation

kNN-MT (Khandelwal et al., 2020)
TRIME-MT (Zhong et al., 2022)

Code & proof generation

DocPrompting (Zhou et al., 2023)
Natural Prover
(Welleck et al., 2022)

NLI

kNN-Prompt (Shi et al., 2022)
NPM (Min et al., 2023)

Sentiment analysis

kNN-Prompt (Shi et al., 2022)
NPM (Min et al., 2023)

Commonsense reasoning

Raco (Yu et al, 2022)

More general NLP tasks

Why?

Why retrieval-based LMs?

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



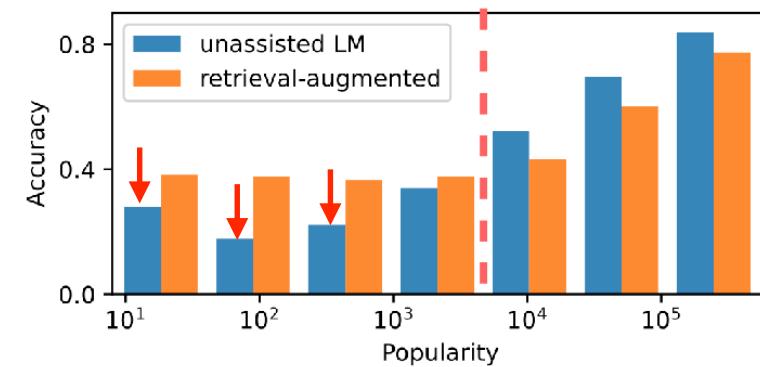
List 5 important papers authored by Geoffrey Hinton



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. "**Deep Learning**" (with Y. Bengio and A. Courville) - Published as a book in **2016**,...
4. "**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

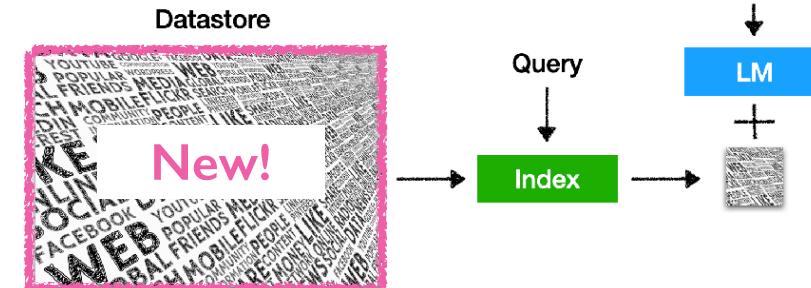
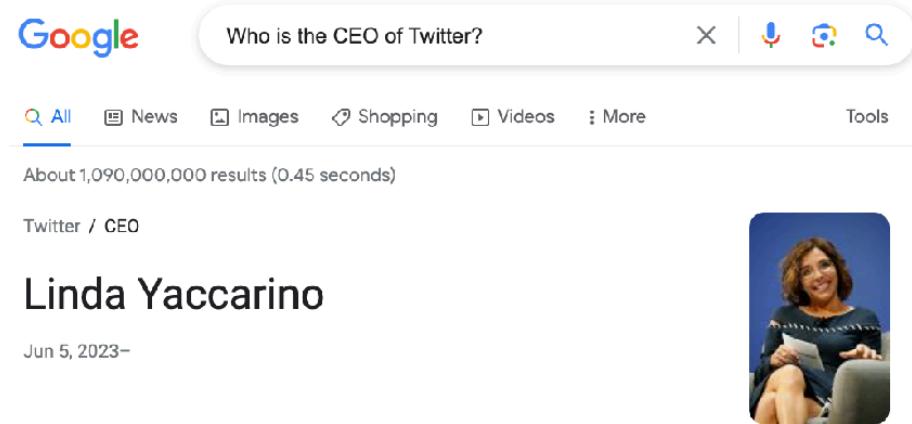


Who is the CEO of Twitter?

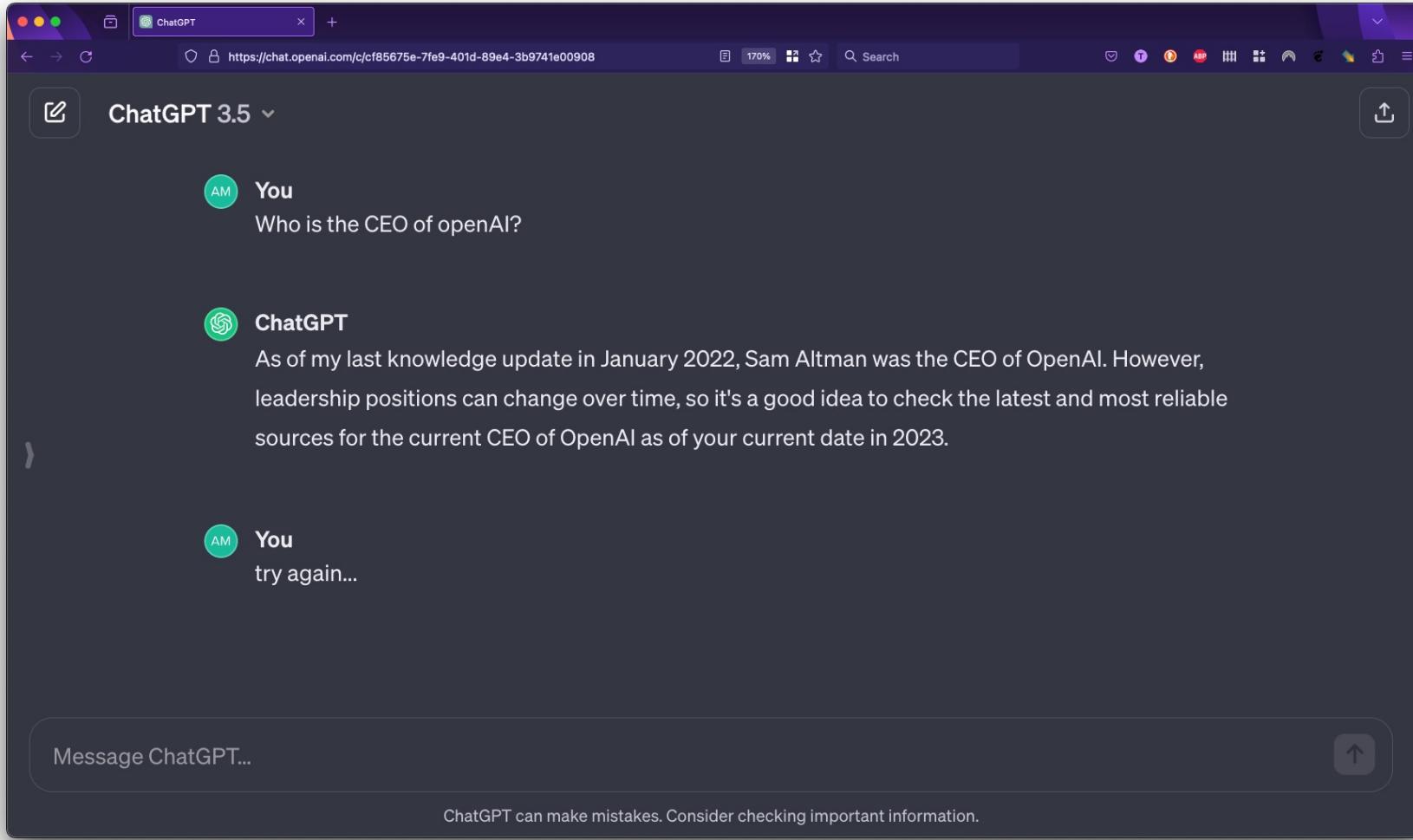


As of my **knowledge cutoff in September 2021**, the CEO of Twitter is **Jack Dorsey**....

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

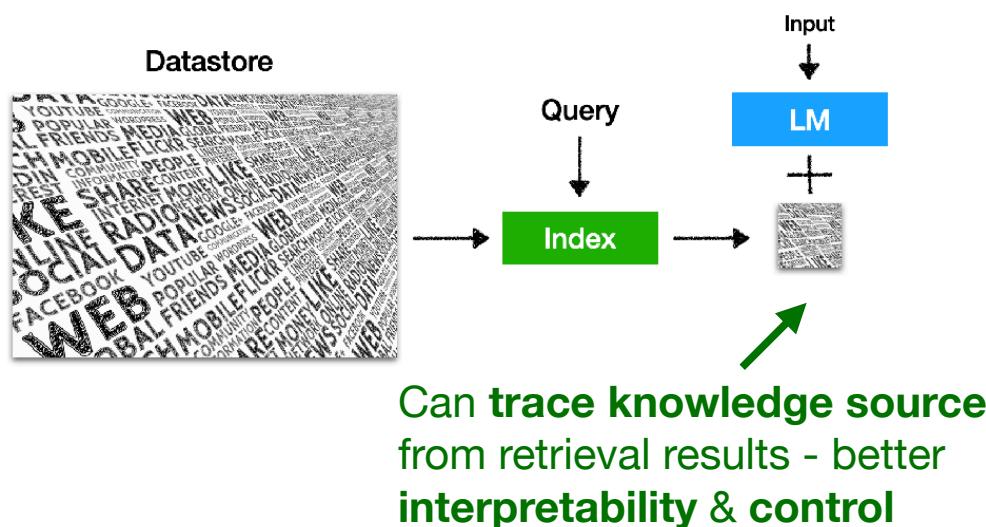


Or in fact...

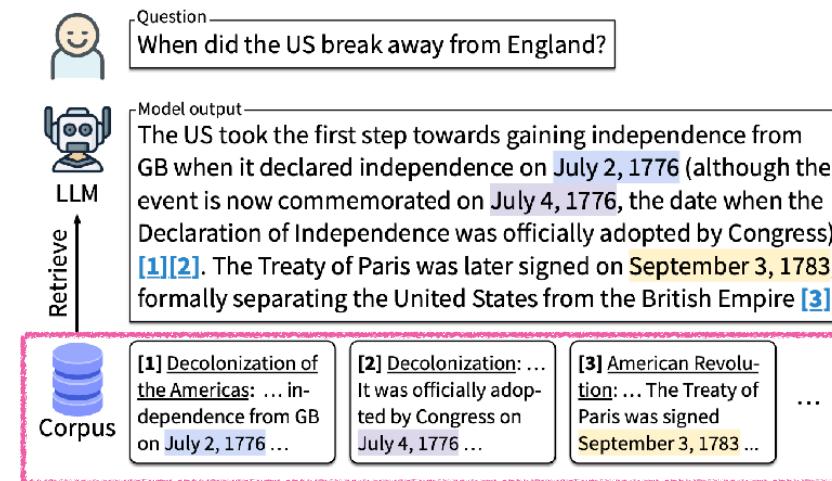


Why retrieval-based LMs?

LLMs' output is challenging to interpret and verify



Generating text with citations



(Nakano et al. 2021; Menick et al., 2022; Gao et al., 2023)

Why retrieval-based LMs?

LLMs are *large* and expensive to train and run



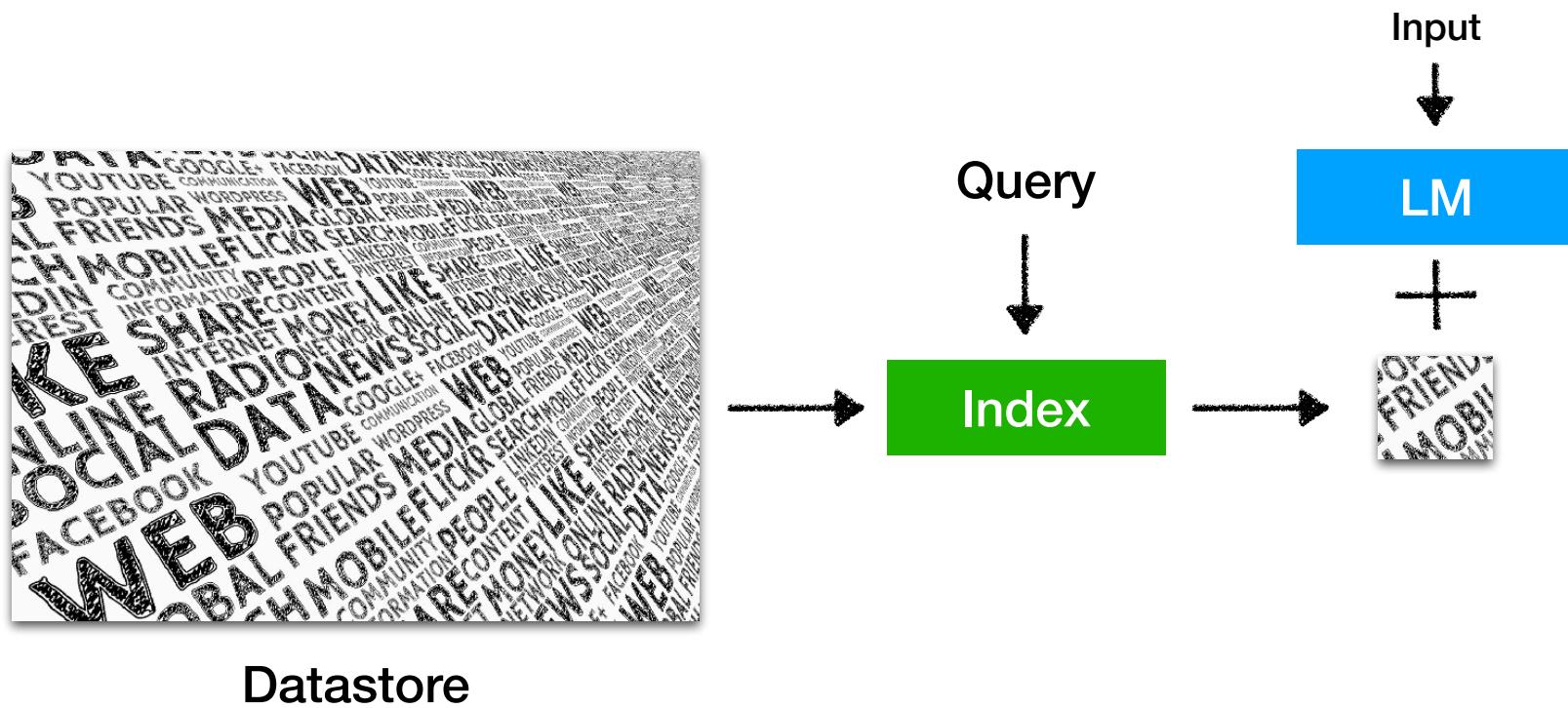
Long-term goal: can we possibly reduce the **training** and **inference costs**, and scale down the size of LLMs?

e.g., RETRO (Borgeaud et al., 2021): “obtains comparable performance to GPT-3 on the Pile, despite using **25x fewer parameters**”

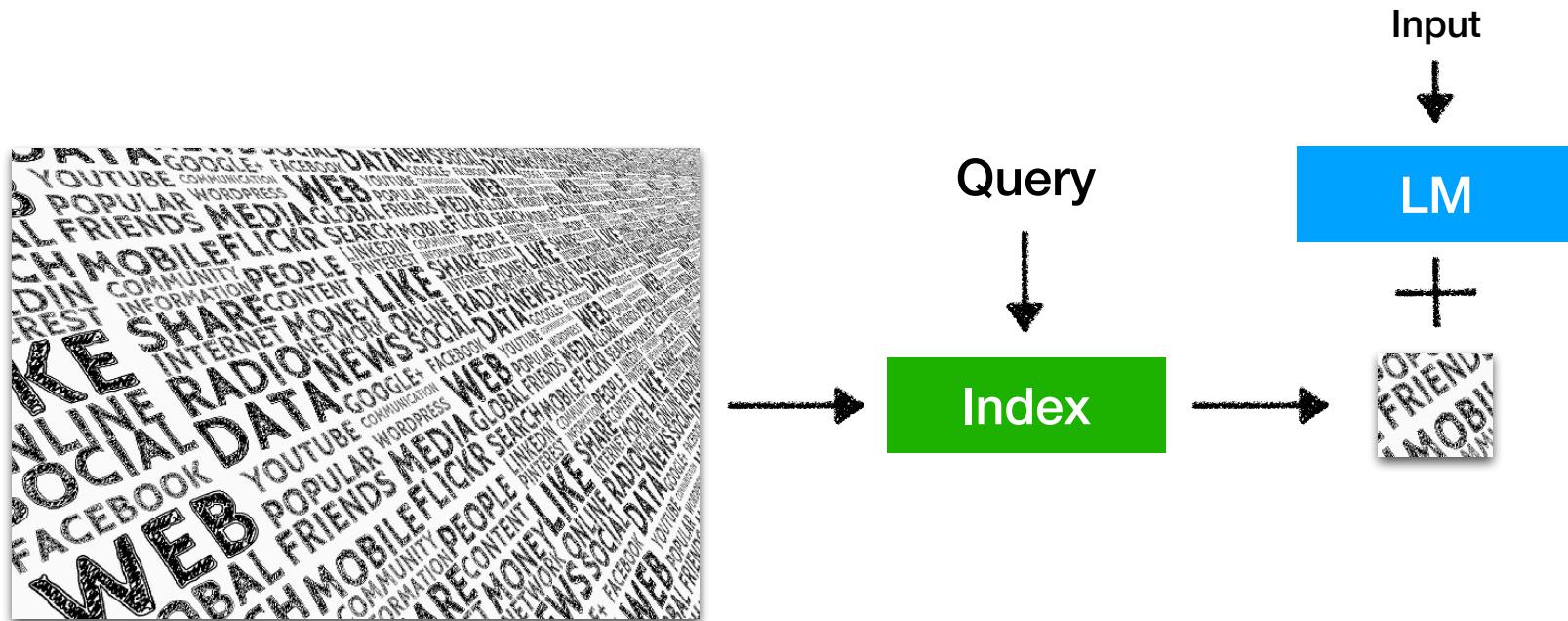


How to do it?

Inference



Inference: Datastore



Datastore

Raw text corpus

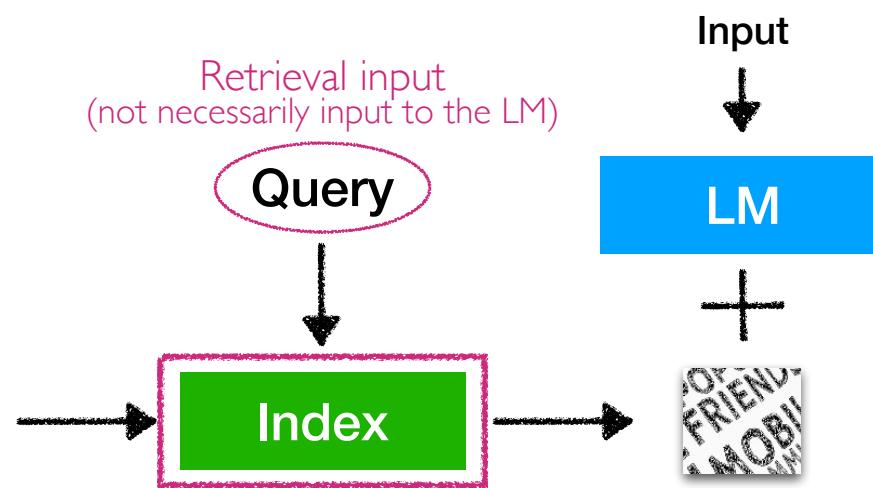
At least billions~trillions of tokens
Not labeled datasets
Not structured data (knowledge bases)

(Extending this definition is possible - some of them in a later part of Section 3)

Inference: Index



Datastore

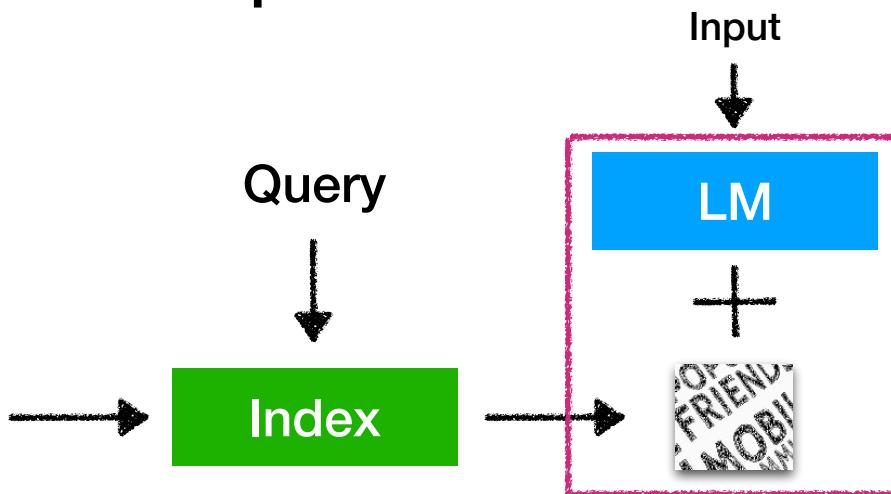


Find a small subset of elements in a datastore that are the most similar to the query

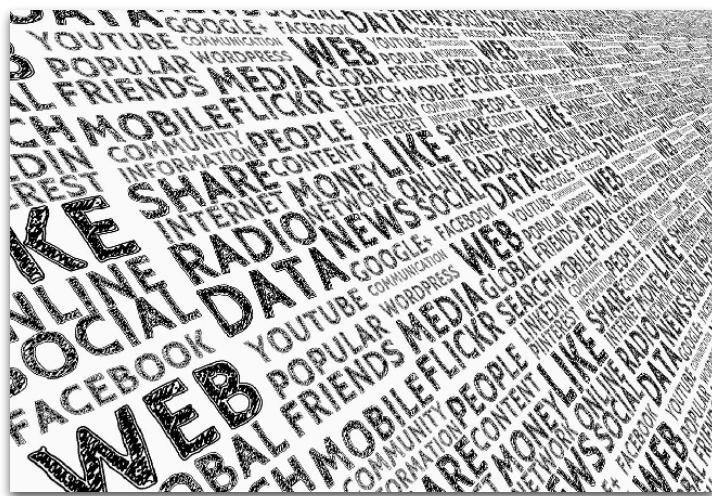
Inference: Incorporation



Datastore



Questions to answer



Datastore

What's the query & when do we retrieve?

Query

Input

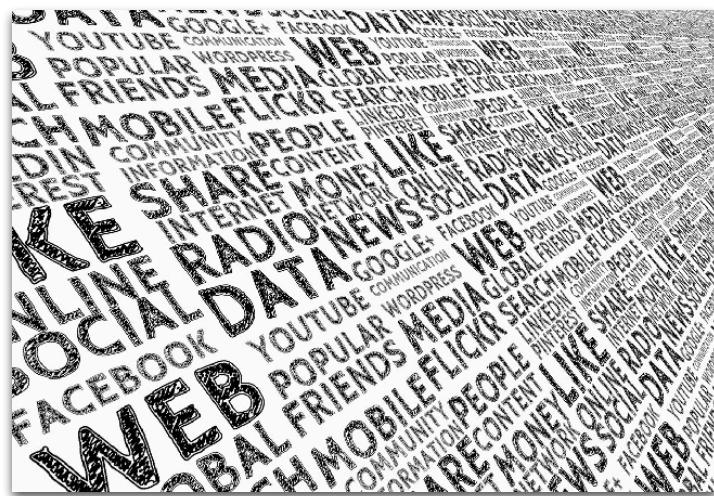
LM

+

How do we use retrieval?

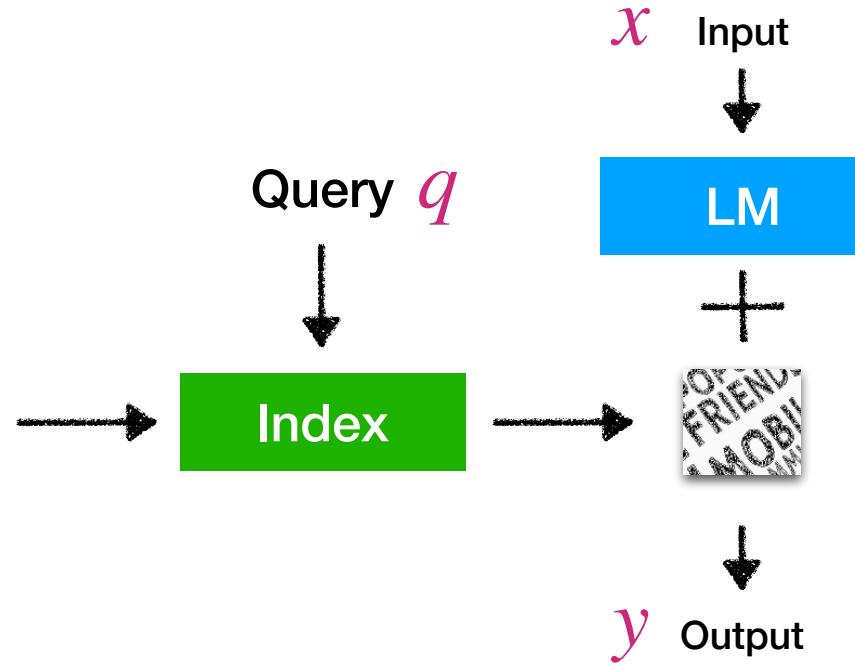
What do we retrieve?

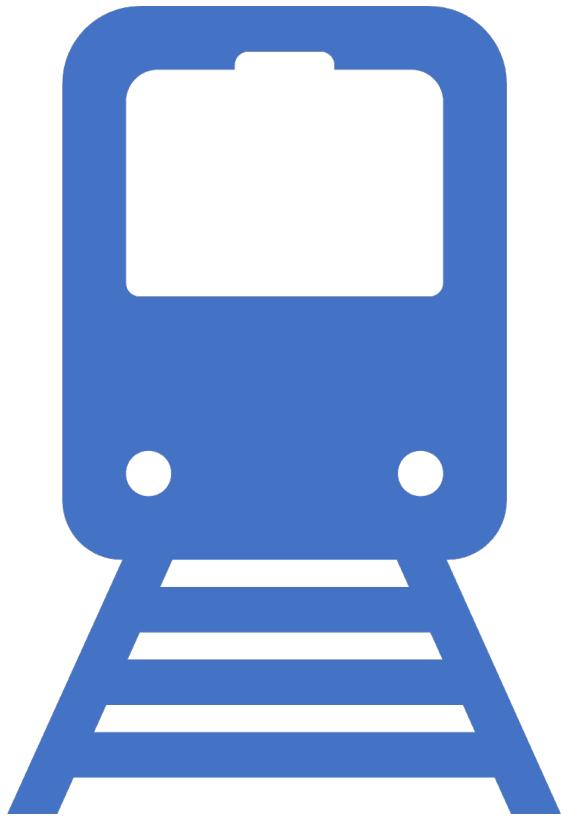
Notations



Datastore

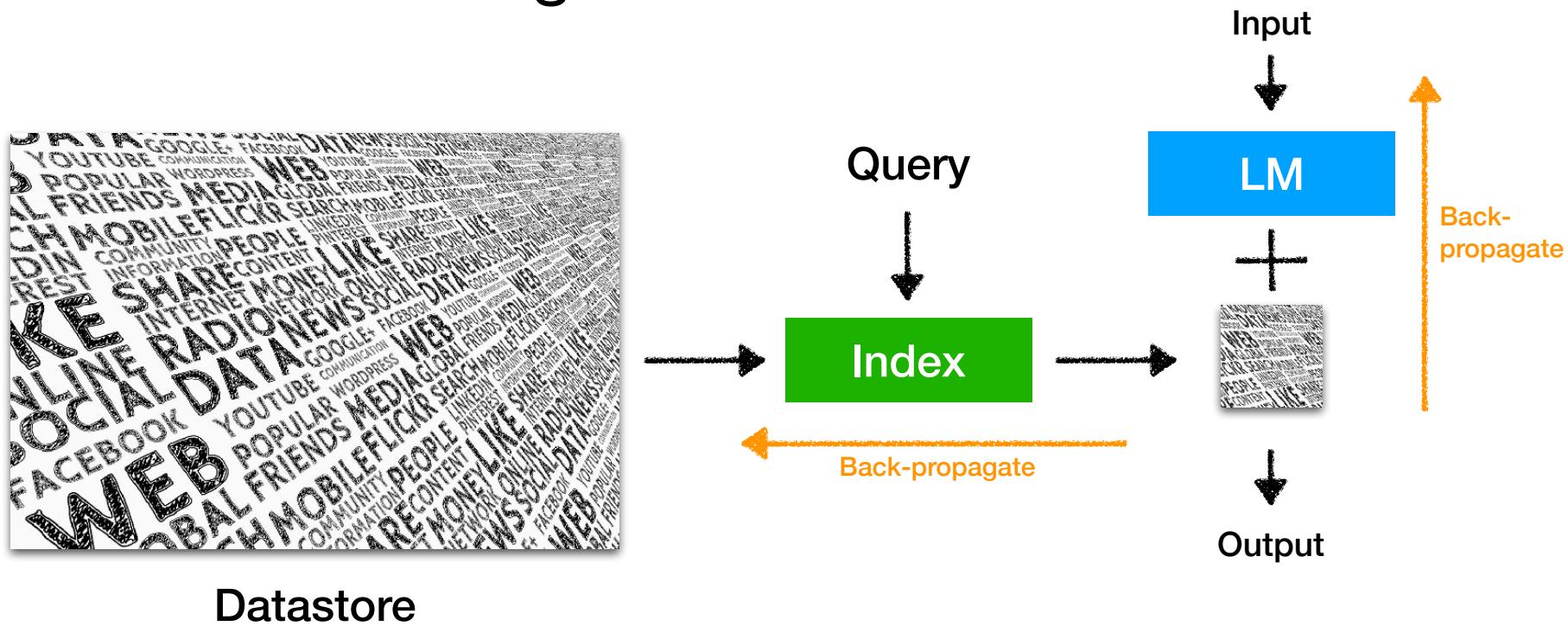
\mathcal{D}





How to train?

Training retrieval-based LMs



Training language models



Minimize $-\log P_{\text{LM}}(y|x)$



GPT



PaLM



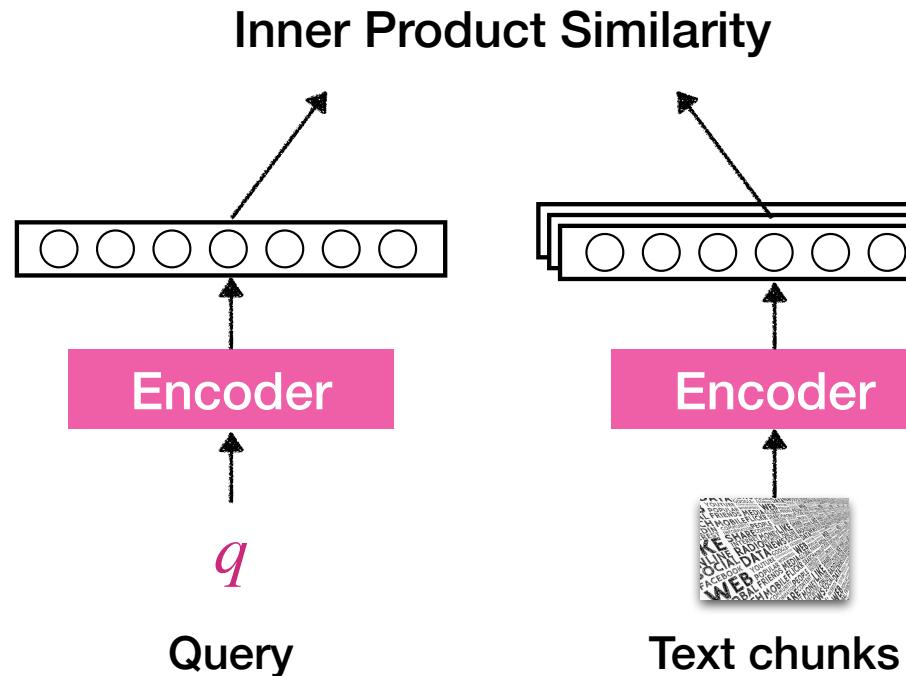
LLaMA



GPT-J

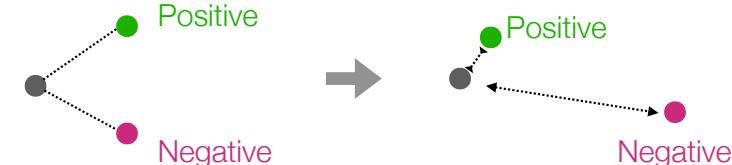
.....

Training dense retrieval models: DPR



$$L(q, p^+, p_1^-, p_2^-, \dots, p_n^-) = -\log \frac{\exp(\text{sim}(q, p^+))}{\exp(\text{sim}(q, p^+)) + \sum_{j=1}^n \exp(\text{sim}(q, p_j^-))}$$

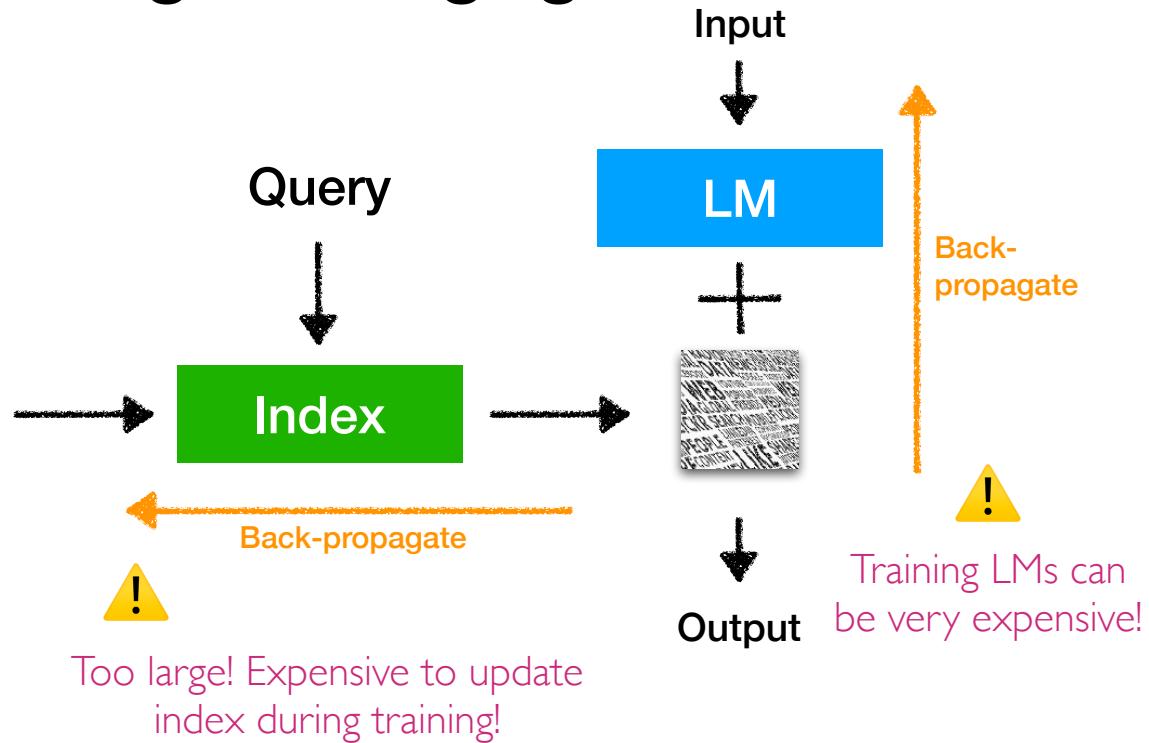
Contrastive learning



Why is training challenging?



Datastore



Training methods for retrieval-based LMs

- Independent training
- Sequential training
- Joint training w/ asynchronous index update
- Joint training w/ in-batch approximation

Training method		
Independent training (Ram et al 2023; Khandelwal et al 2020) Sequential training (Borgeaud et al 2021; Shi et al 2023) Joint training: async update (Guu et al 2020; Izacard et al 2022) Joint training: in-batch approx (Zhong et al 2022; Min et al 2023; Rubin and Berant 2023)	 <ul style="list-style-type: none"> * Easy to implement: off-the-shelf models * Easy to improve: sub-module can be separately improved * End-to-end trained — very good performance! 	 <ul style="list-style-type: none"> * Models are not end-to-end trained — suboptimal performance * Training may be complicated (overhead, batching methods, etc) * Train-test discrepancy still remains

Categorization of retrieval-based LMs

What to retrieve?

Query



Text chunks (passages)?
Tokens?
Something else?

How to use retrieval?

Input



LM

Output

When to retrieve?

w/ retrieval

The capital city of Ontario is Toronto.

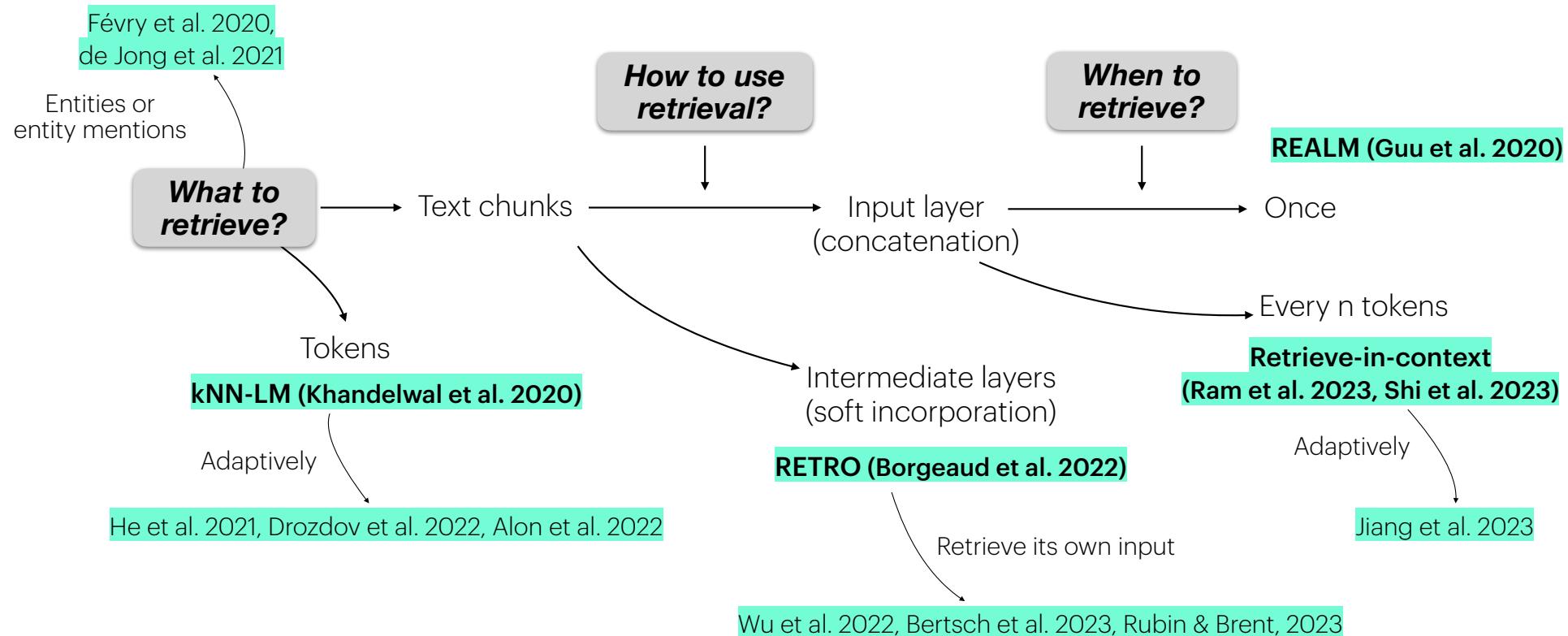
w/ retrieval w/ r w/r w/r w/r w/r w/r

The capital city of Ontario is Toronto.

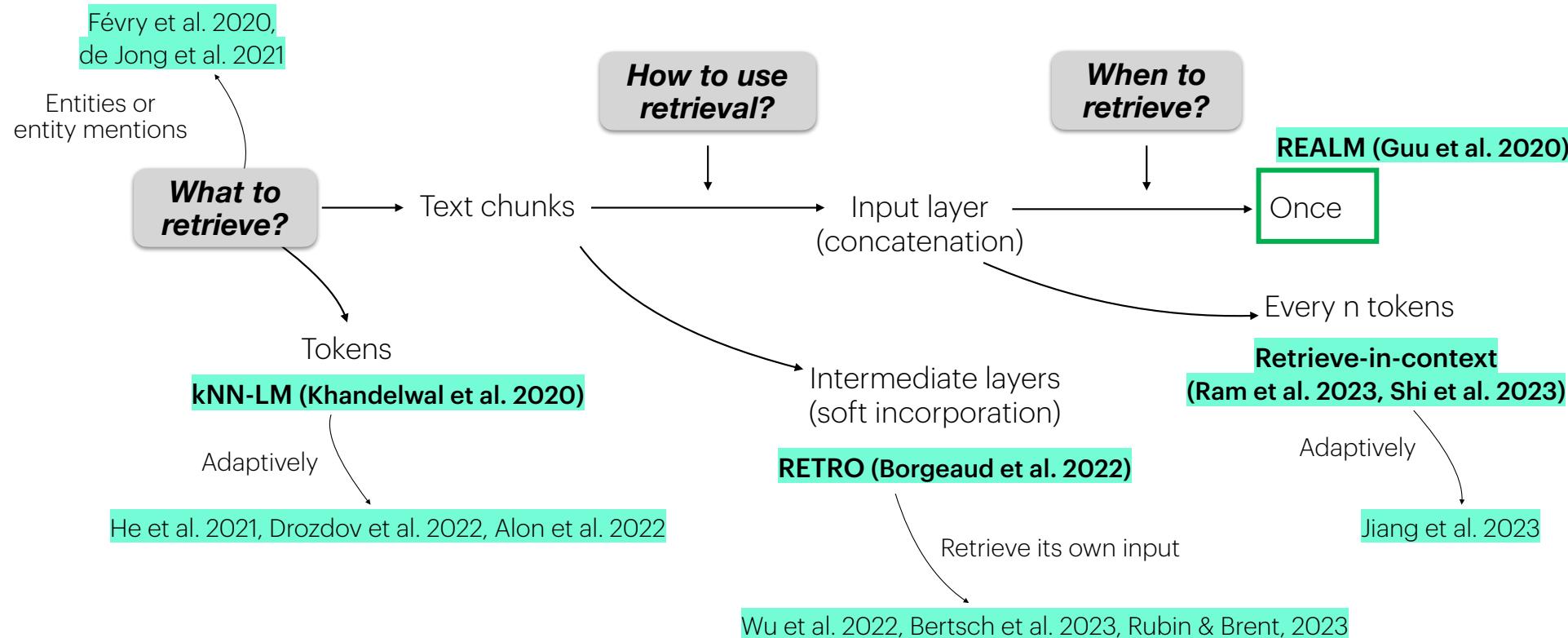
w/ retrieval w/r w/r

The capital city of Ontario is Toronto.

Architectures (roadmap)



Architectures (roadmap)



**RePLUG: Retrieval-
Augmented Black-Box
Language Models**
[Shi et.al. 2023]

RePLUG: Retrieval-Augmented Black-Box Language Models

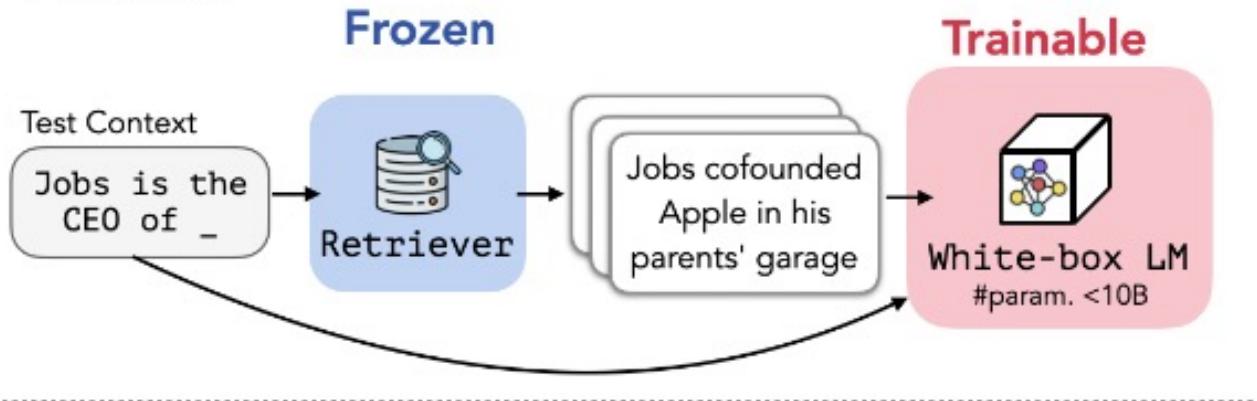
[Shi et.al. 2023]

- decoder-only LM
- adding documents in the prompt as retrieval augmentation
 - no LM finetuning needed
 - [depends on the LM]

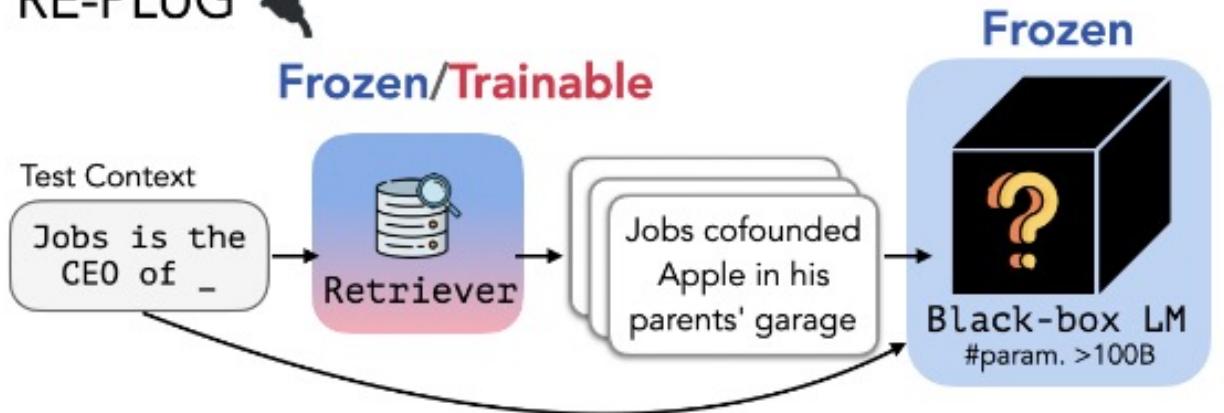
open issues:

- handling long/many documents
 - positional interpolation?
 - lost-in-middle effects
- limited controllability of LM (only via instructions/prompting)

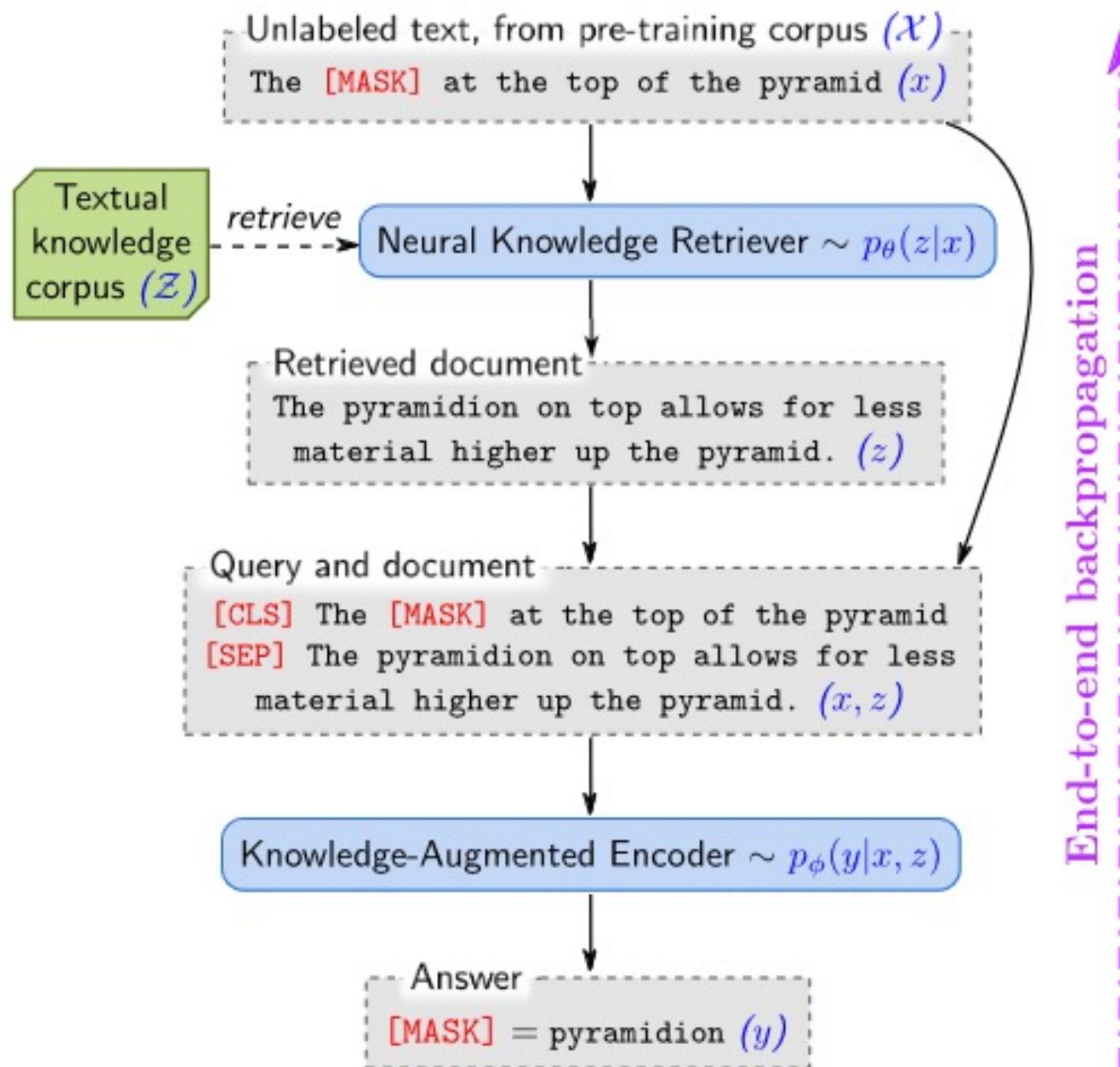
Previous



RE-PLUG ↗



REALM
[Guu et.al. 2020]



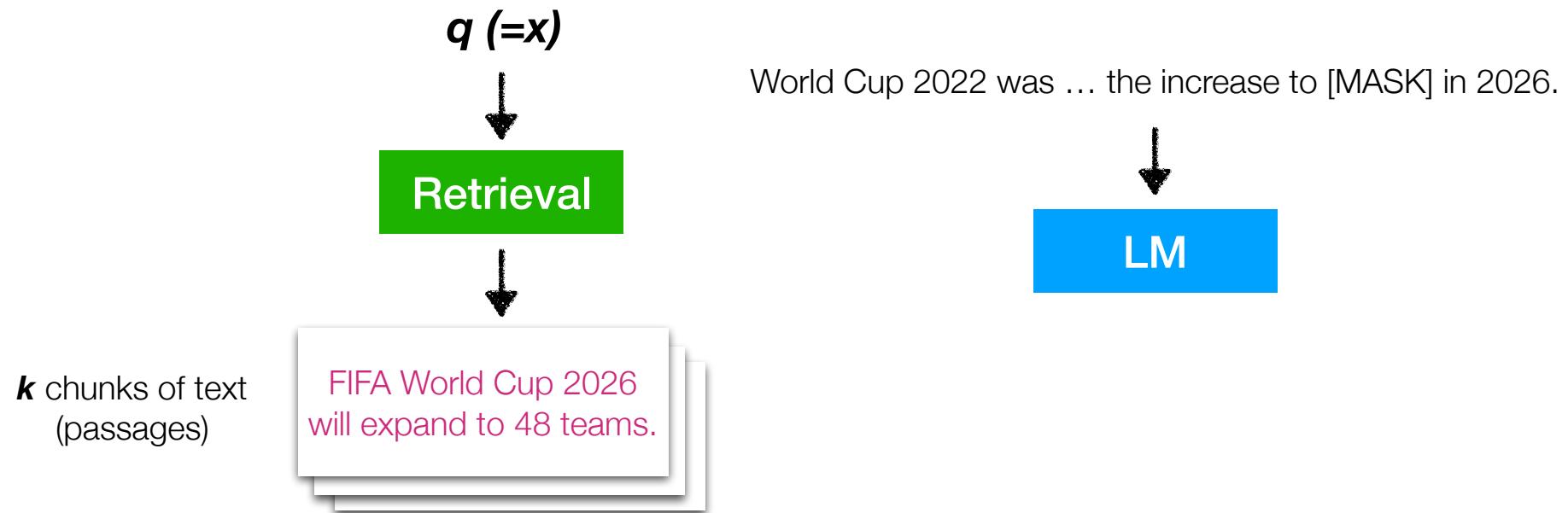
REALM

[Guu et.al. 2020]

- encoder only LM
 - BERT based
- task: QA
 - assumes answer is found as is in documents
 - predicts [answer_start], [answer_end]
- marginalizes over likelihood of document relevance and answer start/end

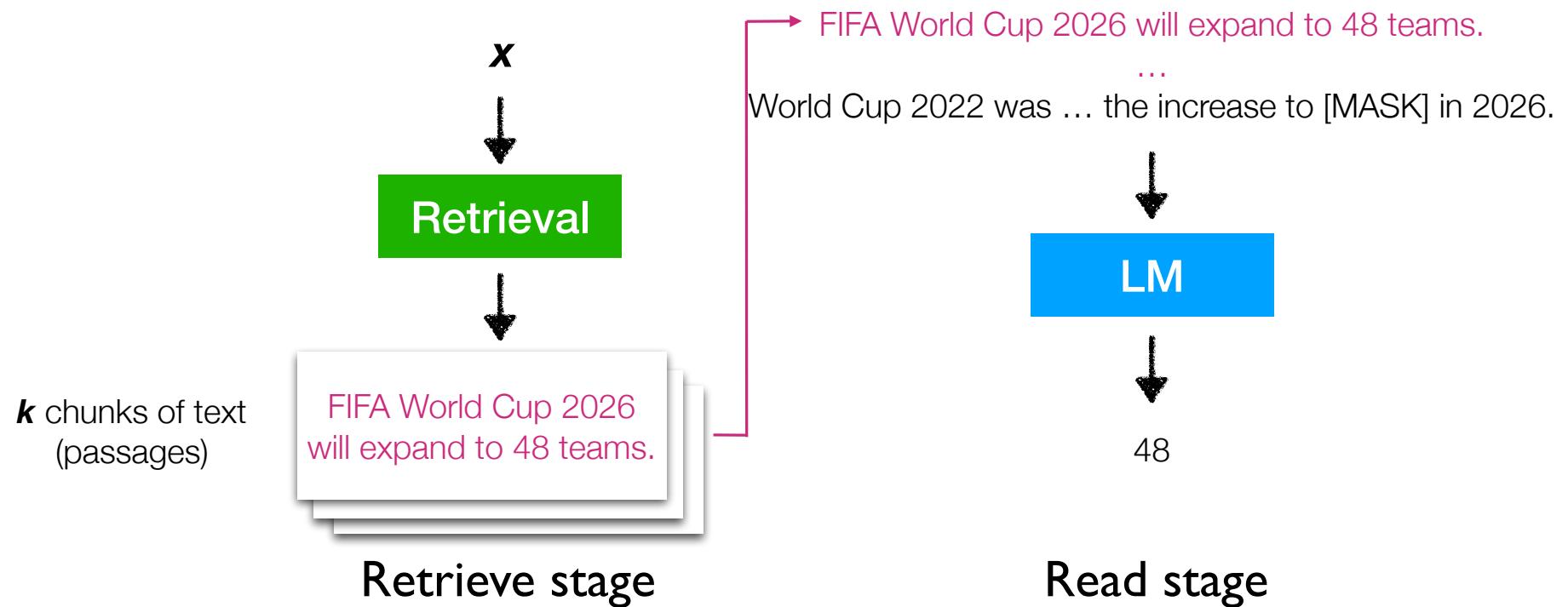
REALM (Guu et al 2020)

x = World Cup 2022 was the last with 32 teams before the increase to [MASK] in 2026.



REALM (Guu et al 2020)

x = World Cup 2022 was the last before the increase to [MASK] in the 2026 tournament.



REALM: (I) Retrieve stage

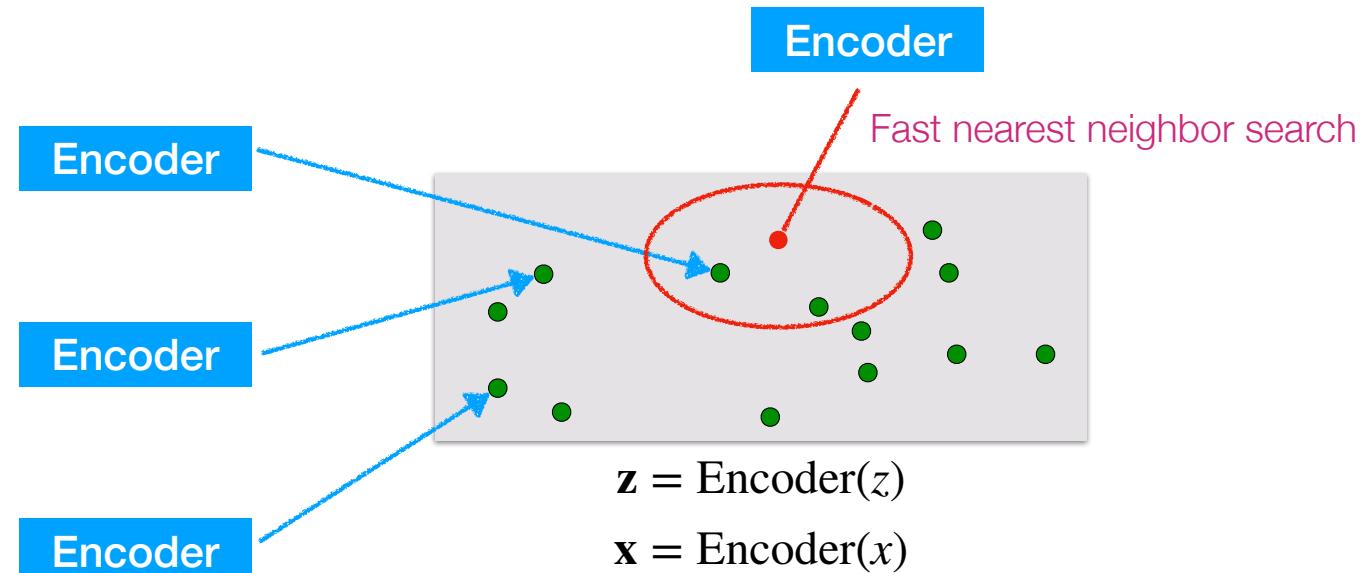
x = World Cup 2022 was ... the increase to [MASK] in 2026.

FIFA World Cup 2026 will expand to 48 teams.

In 2022, the 32 national teams involved in the tournament.

Team USA celebrated after winning its match against Iran ...

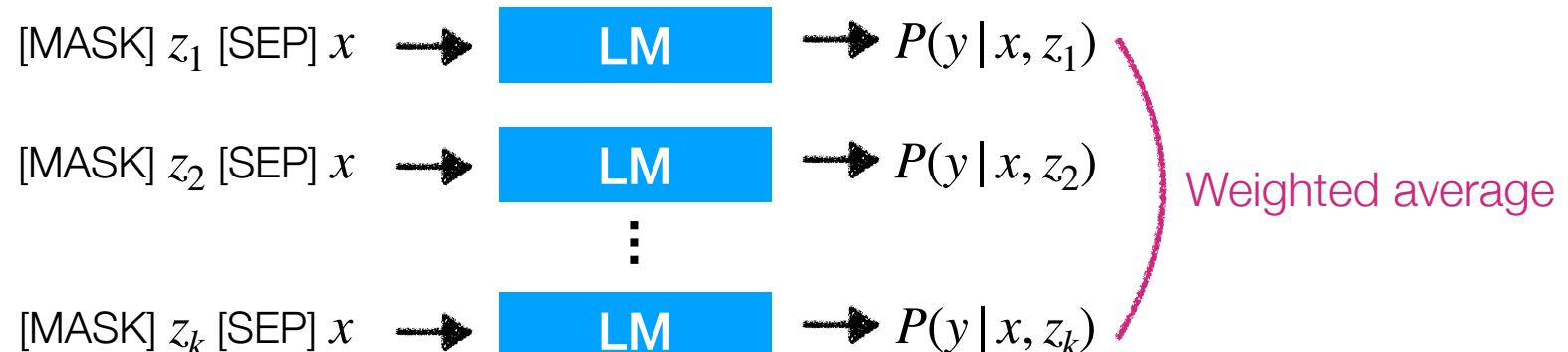
Wikipedia
13M chunks (passages)
(called *documents* in the paper)



$$z_1, \dots, z_k = \arg\top-k(\mathbf{x} \cdot \mathbf{z})$$

k retrieved blocks

REALM: (2) Read stage



Need to approximate
→ Consider top k chunks only

$$\sum_{z \in \mathcal{D}} P(z | x) P(y | x, z)$$

0 if not one of top k

from the retrieve stage from the read stage

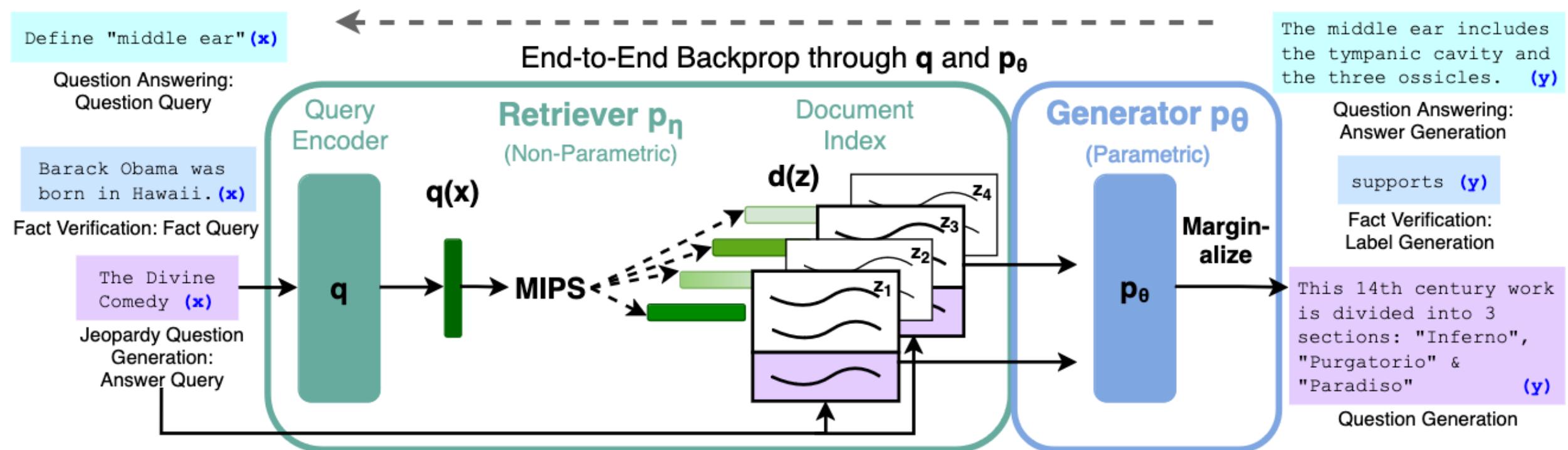
Retrieval-Augmented- Generation

[Lewis et.al. 2020]

Retrieval-Augmented-Generation

[Lewis et.al. 2020]

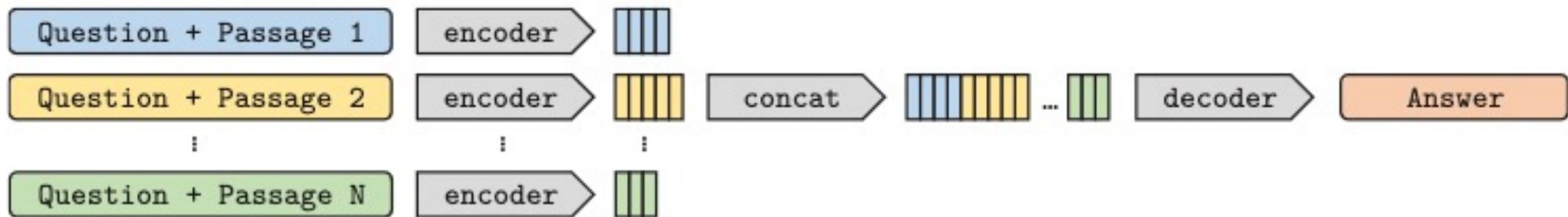
- encoder-decoder model (BART) for generation
- training:
 - retriever (DPR): index frozen; training only the query encoder, based on the probability the GT answer is generated given this document
 - LM: next token prediction



Fusion-in-Decoder (FiD) [Izacard et.al. 2020]

Fusion-in-Decoder (FiD) [Izacard et.al. 2020]

- similar to RAG, but
 - T5 instead of BART; retriever is static
 - encode **all documents at once**, pass document embeddings to decoder
 - outperforms RAG on QA
 - why??



**Fusion-in-Decoder
Knowledge-Distillation
(FiD-KD) [Izacard et.al.
2021]**

How to train the retriever?

Marie Curie was a physicist...

She was born in Warsaw...

Curie won the 1911 Nobel...

Marie Curie died in 1934...

In 1898, they discovered...

Curie, born in Warsaw in 18...

Where was Marie Curie born?

Retriever

Retriever

Retriever

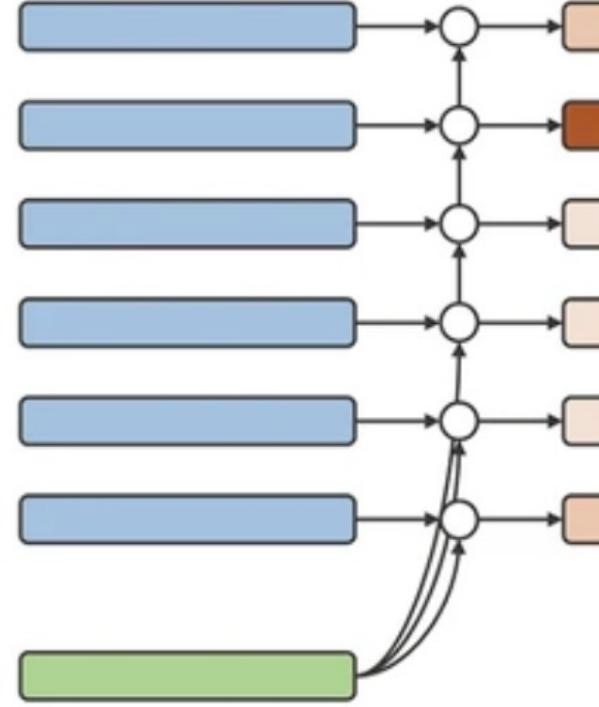
Retriever

Retriever

Retriever

Retriever

Retriever



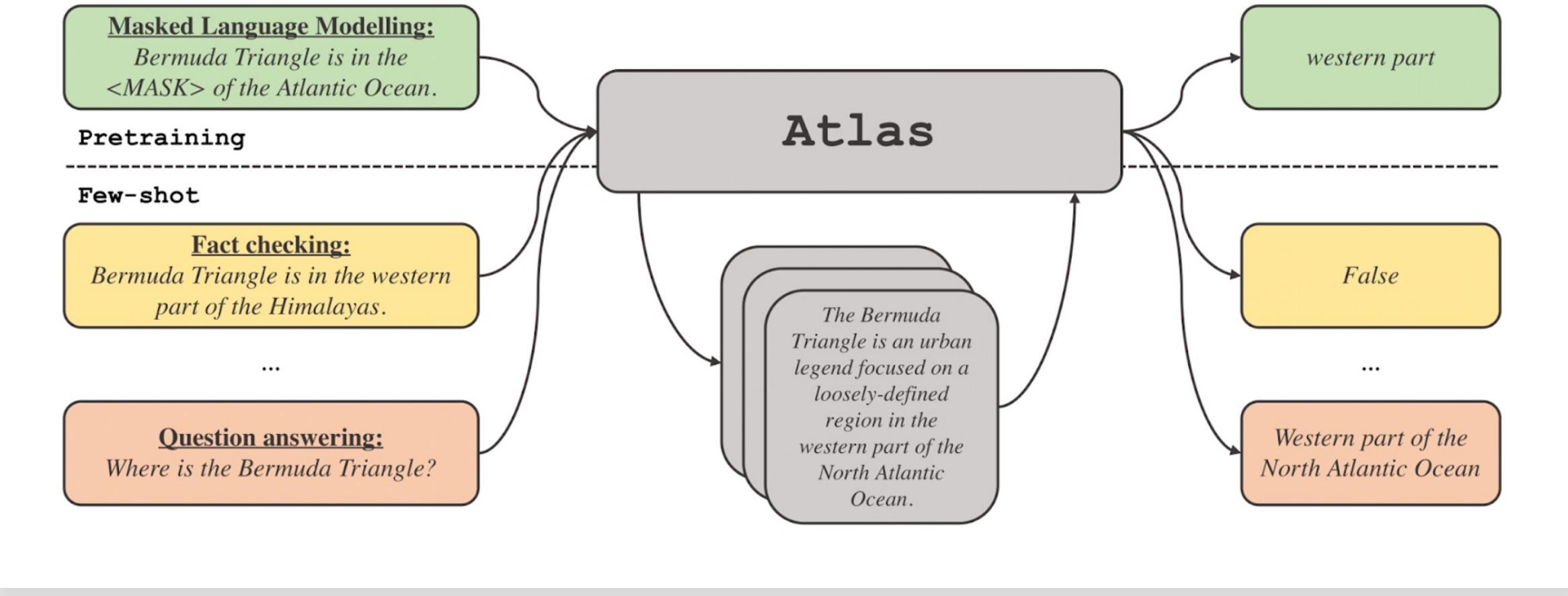
Fusion-in-Decoder Knowledge-Distillation (FiD-KD) [Izacard et.al. 2021]

How to train the retriever?

- no document relevance judgements needed; only QA pairs
- distil training signal (reader → retriever) :
 - Measure “contribution” of each document while generating answer (LM)

Few-shot learning with RAG models (Atlas) [Izacard et.al. 2022]

joint retriever/reader pre-training & few shot capabilities



Few-shot learning with RAG models (Atlas) [Izacard et.al. 2022]

- joint retriever/reader pre-training & few shot capabilities*
- joint unsupervised pretraining (ICT) of retriever & reader
 - instead of iteratively freezing & training each model
 - few-shot (64 examples) finetuning: decent performance

Atlas: Retriever training

Similarity based on retrieval encoder

$$P_{\text{retr}}(z | q) = \frac{\exp(s(z, q))}{\sum_{k=1}^K \exp(s(z_k, q))}$$

How likely each document is retrieved

KL Divergence

Prob of the gold labels if augmenting this text chunk

$$P_{\text{ppl}}(z | q, y) = \frac{\exp(\log P_{\text{LM}}(y | q, z))}{\sum_{k=1}^K \exp(\log P_{\text{LM}}(y | q, z_k))}$$

How much each document improves the ppl



Perplexity Distillation

Retrieve the text that can help LM encoders improve perplexity

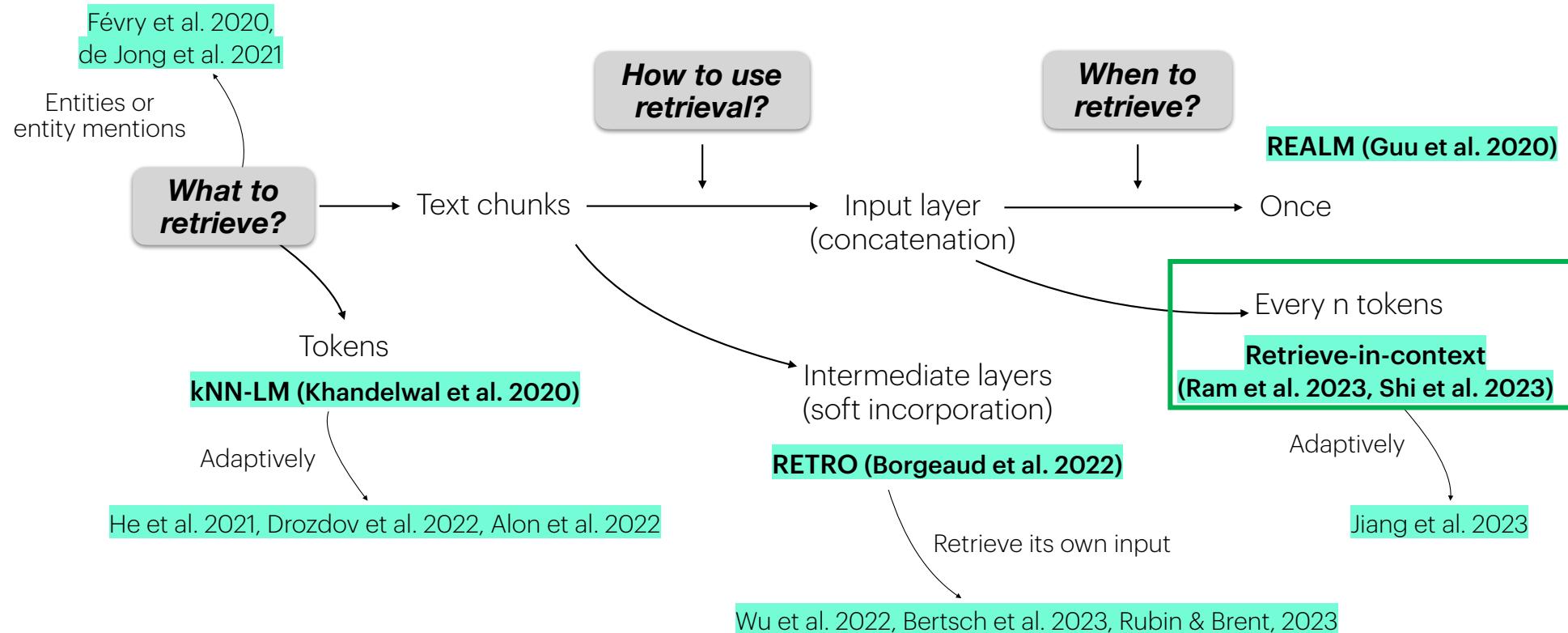
Model	NQ		TriviaQA filtered		TriviaQA unfiltered	
	64-shot	Full	64-shot	Full	64-shot	Full
GPT-3 (Brown et al., 2020)	29.9	-	-	-	71.2	-
Gopher (Rae et al., 2021)	28.2	-	57.2	-	61.3	-
Chinchilla (Hoffmann et al., 2022)	35.5	-	64.6	-	72.3	-
PaLM (Chowdhery et al., 2022)	39.6	-	-	-	81.4	-
RETRO (Borgeaud et al., 2021)	-	45.5	-	-	-	-
FiD (Izacard & Grave, 2020)	-	51.4	-	67.6	-	80.1
FiD-KD (Izacard & Grave, 2021)	-	54.7	-	73.3	-	-
R2-D2 (Fajcik et al., 2021)	-	55.9	-	69.9	-	-
ATLAS	42.4	60.4	74.5	79.8	84.7	89.4

Where do we stand?

Steps that seem to help:

- improving retriever (reader distillation, DPR → Contriever)
- iterative/joint pretraining

Architectures (roadmap)



Retrieval-in-context in LM

How frequent should retrieval be?

World Cup 2022 was the last with



Retrieval



The 2022 FIFA World Cup (...) 32 national teams involved in the tournament. World Cup 2022 was the last with



LM



32 teams before the increase to 48 in the 2026 tournament.

explained by retrieval

not really covered

Retrieval-in-context in LM

How frequent should retrieval be?

World Cup 2022 was the last with



Retrieval



The 2022 FIFA World Cup (...) 32 national teams involved in the tournament. World Cup 2022 was the last with



LM



32 teams before the increase

Retrieval-in-context in LM

How frequent should retrieval be?



Retrieval-in-context in LM

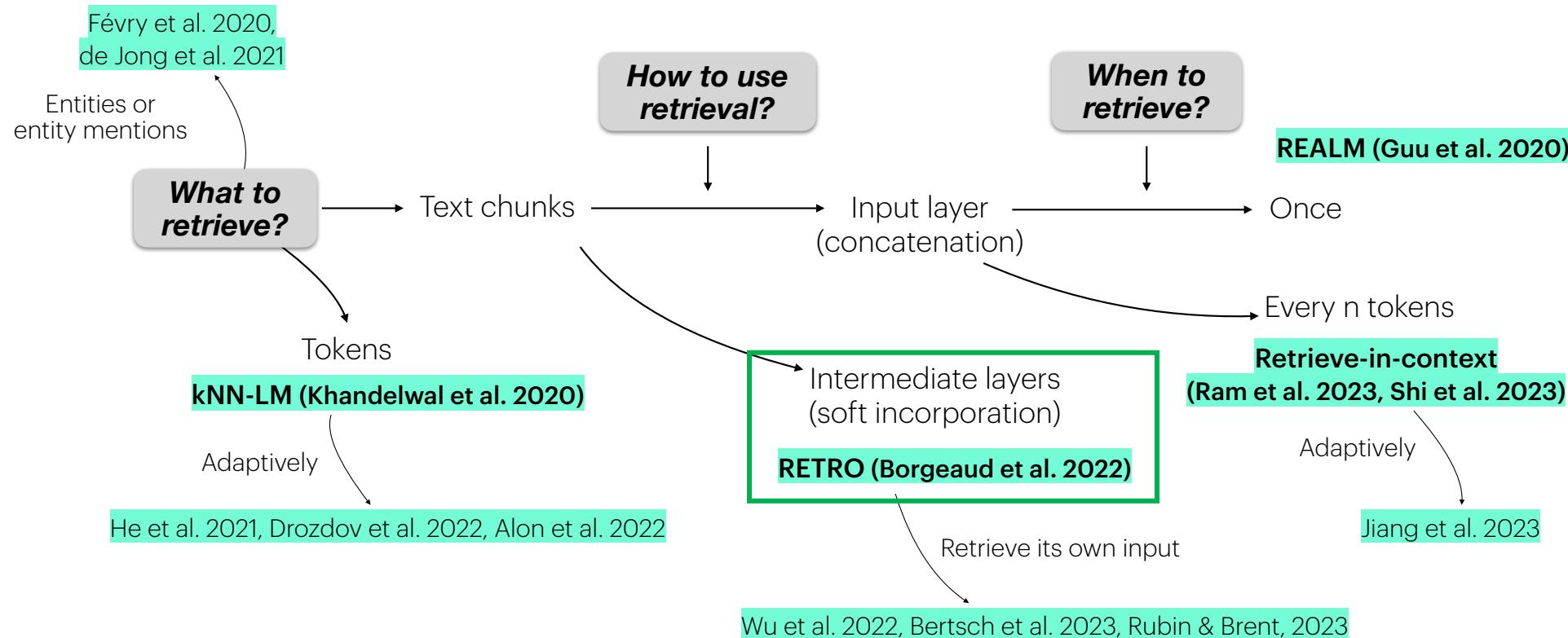
How frequent should retrieval be?



World Cup 2022 was the last with 32 teams before the increase



Architectures (roadmap)



RETRO

Beyond simple concatenation of retrieved passages to the LM

RETRO (Borgeaud et al. 2021)

- ✓ Incorporation in the “intermediate layer” instead of the “input” layer
→ designed for many blocks, frequently, more efficiently
- ✓ Scale the datastore to retrieve from

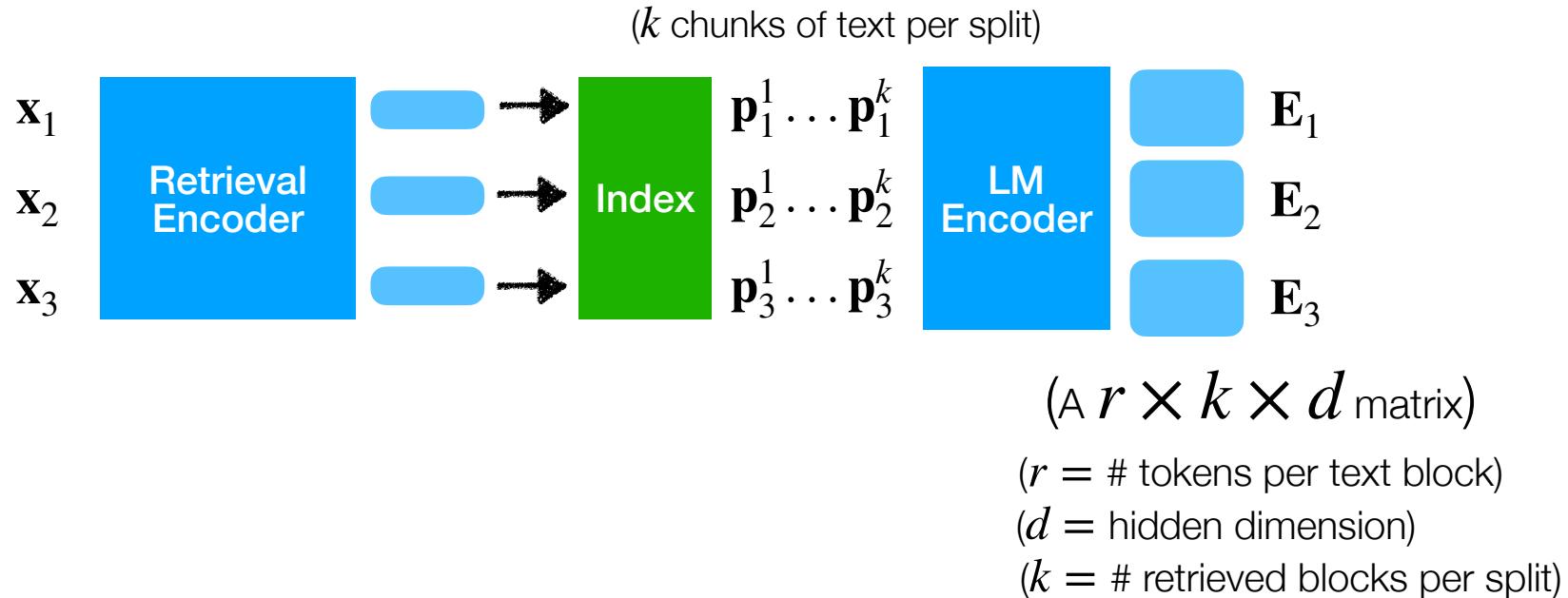
RETRO (Borgeaud et al. 2021)

\mathbf{x} = World Cup 2022 was / the last with 32 teams, / before the increase to

\mathbf{x}_1

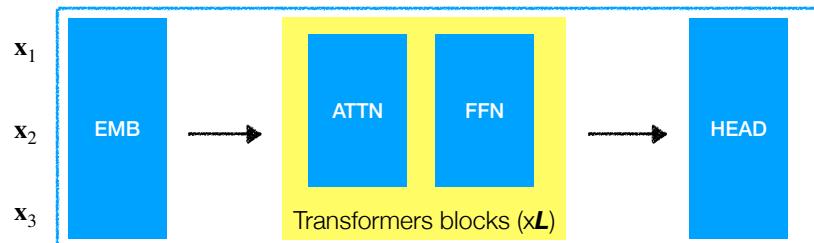
\mathbf{x}_2

\mathbf{x}_3

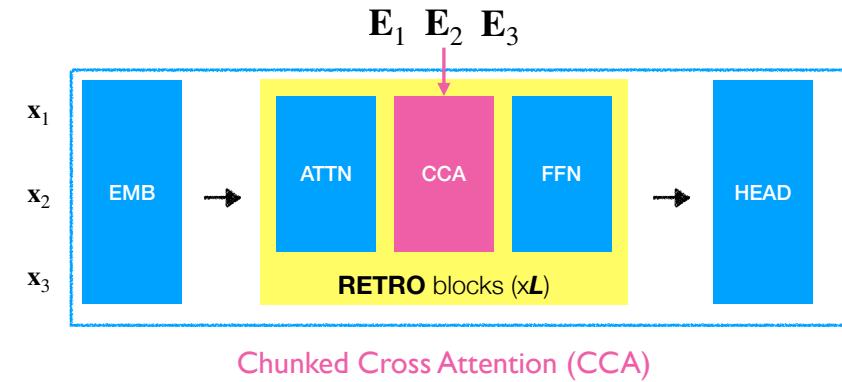


RETRO Chunked Cross Attention architecture

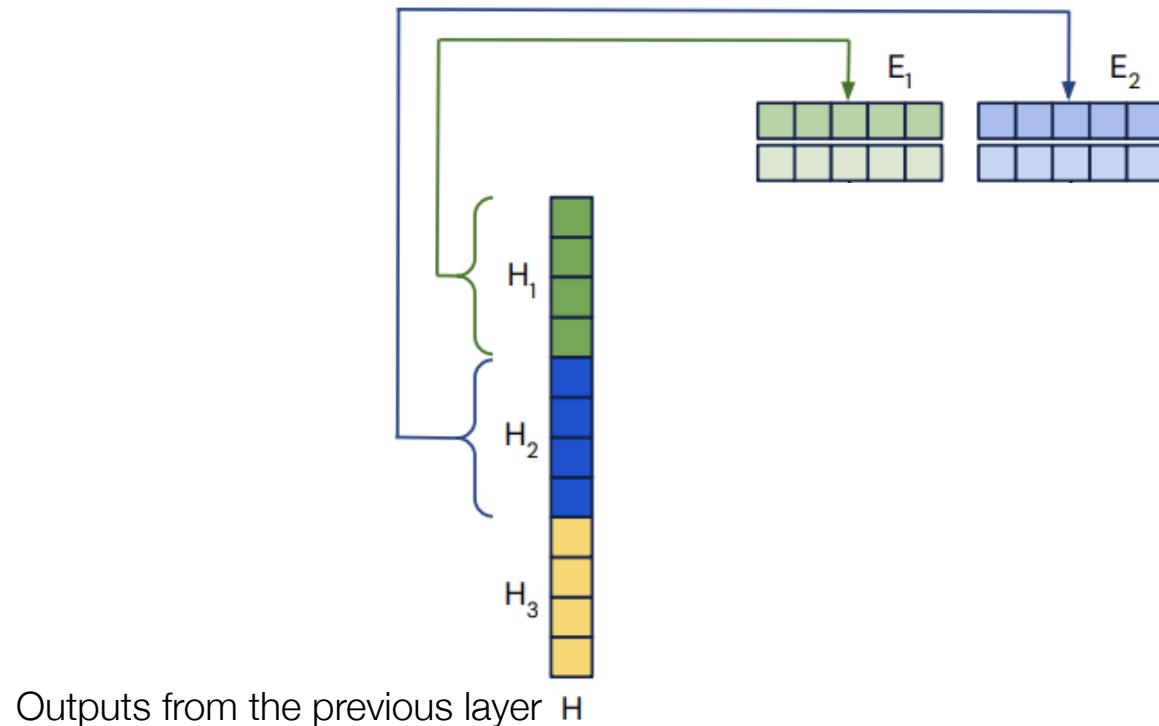
Regular decoder



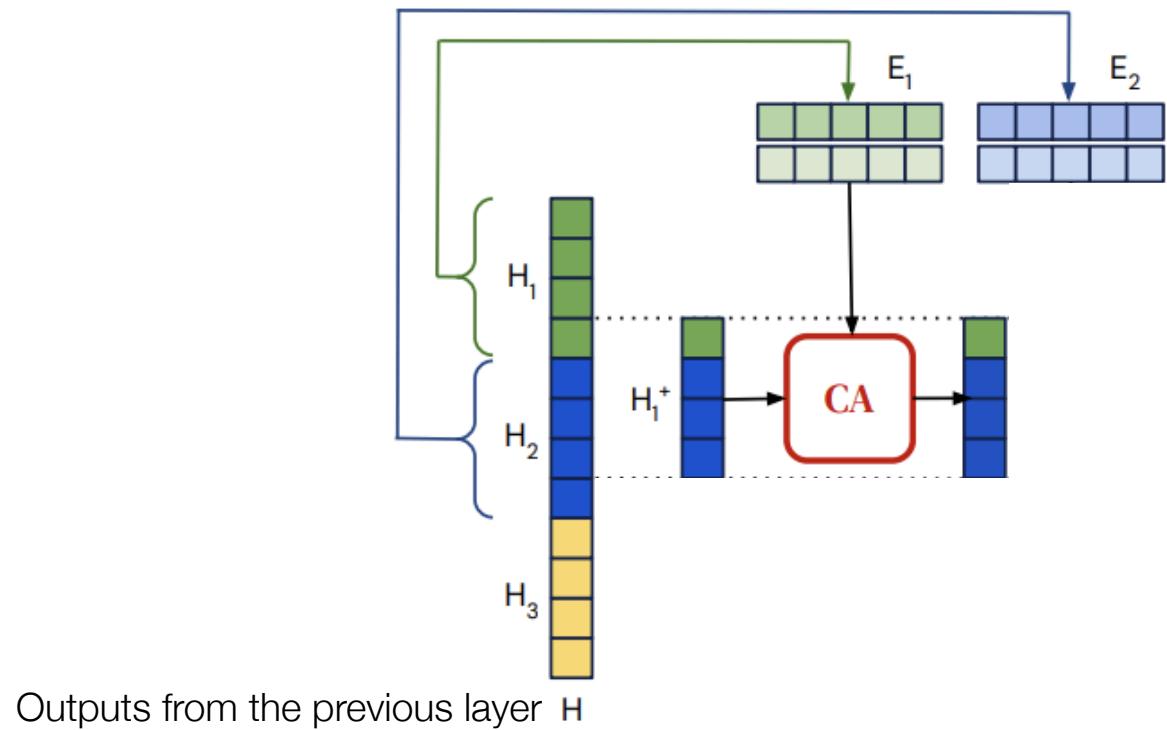
Decoder in RETRO



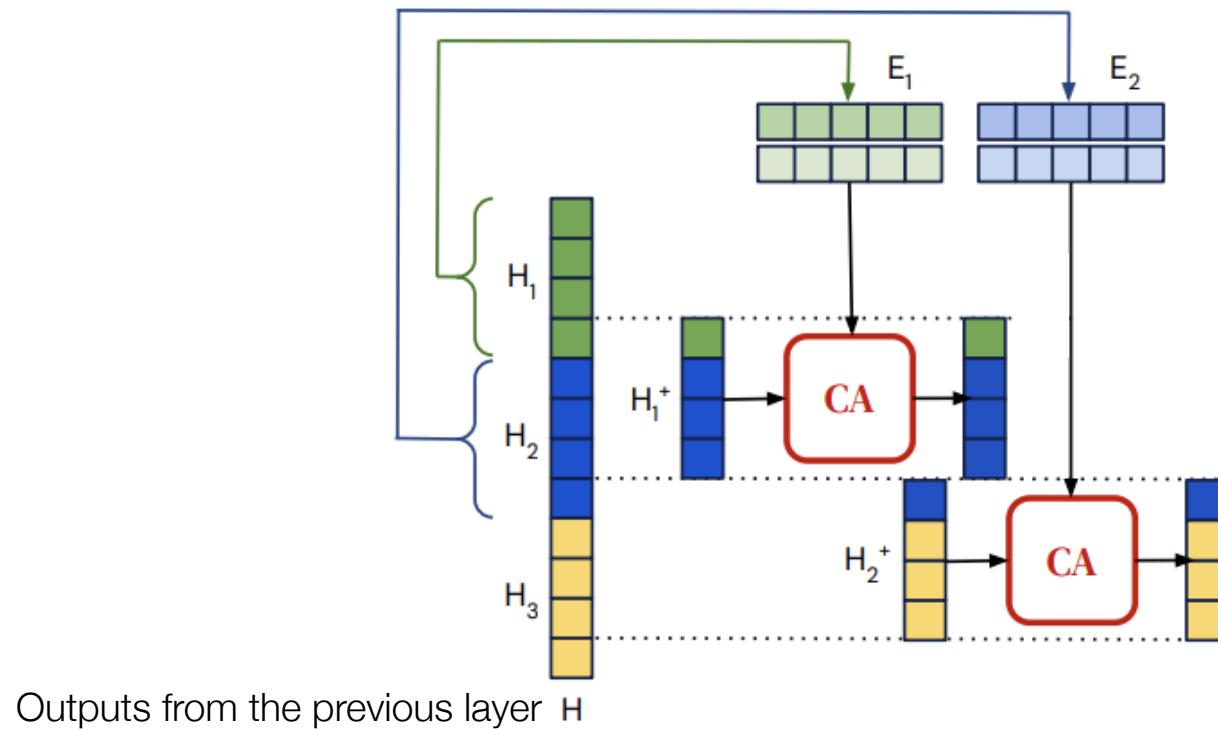
Chunked Cross Attention



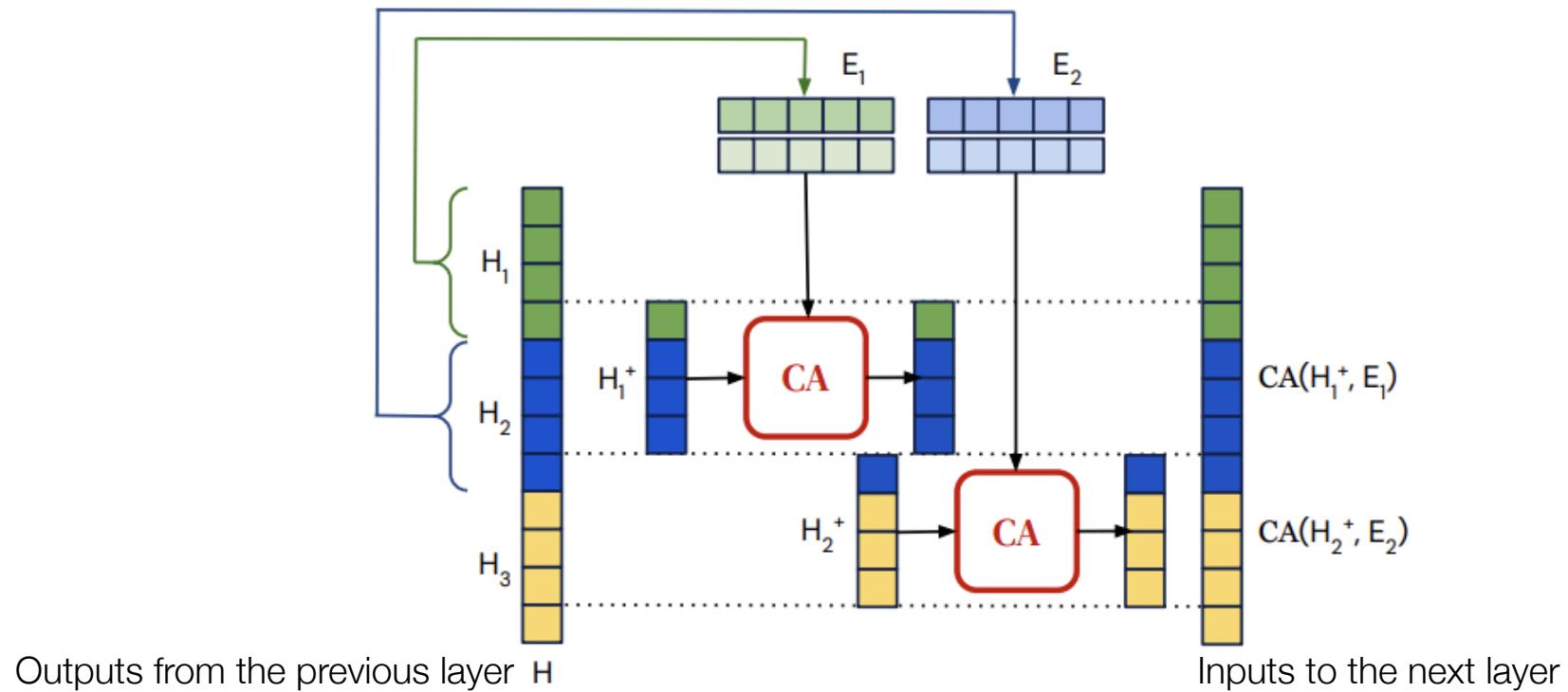
Chunked Cross Attention



Chunked Cross Attention



Chunked Cross Attention



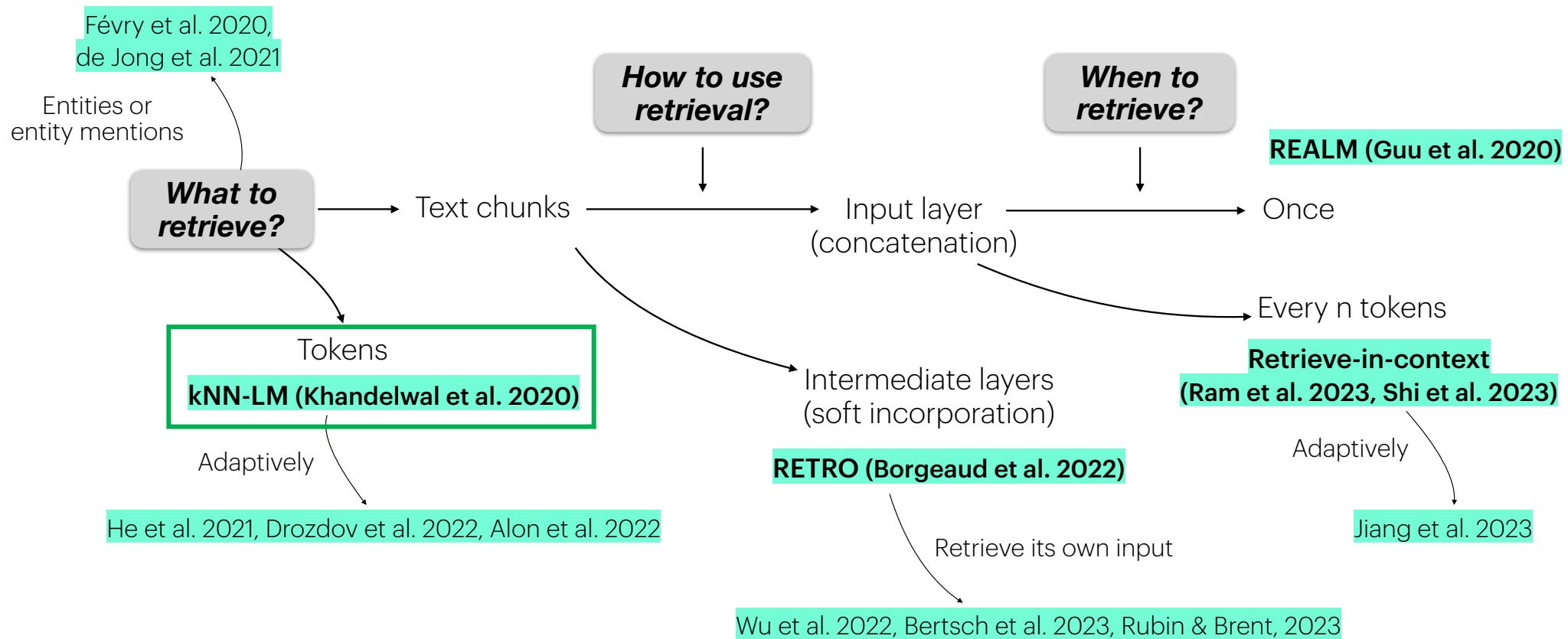
Results

Perplexity: The lower the better

Model	Retrieval Set	#Database tokens	#Database keys	Valid	Test
Adaptive Inputs (Baevski and Auli, 2019)	-	-	-	17.96	18.65
SPALM (Yogatama et al., 2021)	Wikipedia	3B	3B	17.20	17.60
kNN-LM (Khandelwal et al., 2020)	Wikipedia	3B	3B	16.06	16.12
Megatron (Shoeybi et al., 2019)	-	-	-	-	10.81
Baseline transformer (ours)	-	-	-	21.53	22.96
kNN-LM (ours)	Wikipedia	4B	4B	18.52	19.54
RETRO	Wikipedia	4B	0.06B	18.46	18.97
RETRO	C4	174B	2.9B	12.87	10.23
RETRO	MassiveText (1%)	18B	0.8B	18.92	20.33
RETRO	MassiveText (10%)	179B	4B	13.54	14.95
RETRO	MassiveText (100%)	1792B	28B	3.21	3.92

Significant improvements by retrieving from 1.8 trillion tokens

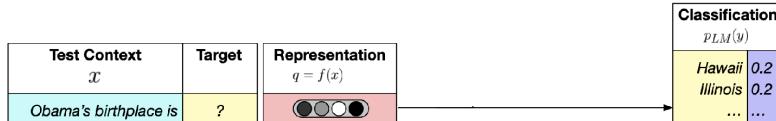
Architectures (roadmap)



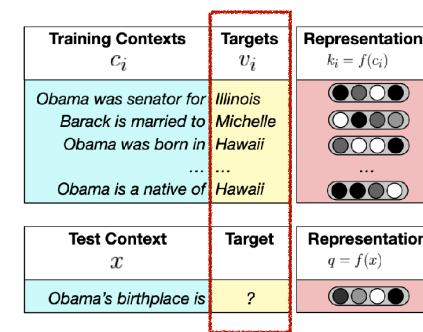
kNN-LM:

Generating by retrieving (or “memorizing”?)

Language Model



kNN-LM (Khandelwal et al. 2020)



Which tokens in a datastore are close to the next token?

kNN-LM (Khandelwal et al. 2020)

Training Contexts c_i	Targets v_i	Representations $k_i = f(c_i)$
<i>Obama was senator for</i>	<i>Illinois</i>	
<i>Barack is married to</i>	<i>Michelle</i>	
<i>Obama was born in</i>	<i>Hawaii</i>	
...
<i>Obama is a native of</i>	<i>Hawaii</i>	

Test Context x	Target	Representation $q = f(x)$
<i>Obama's birthplace is</i>	?	



Which tokens in a datastore are close to the next token?

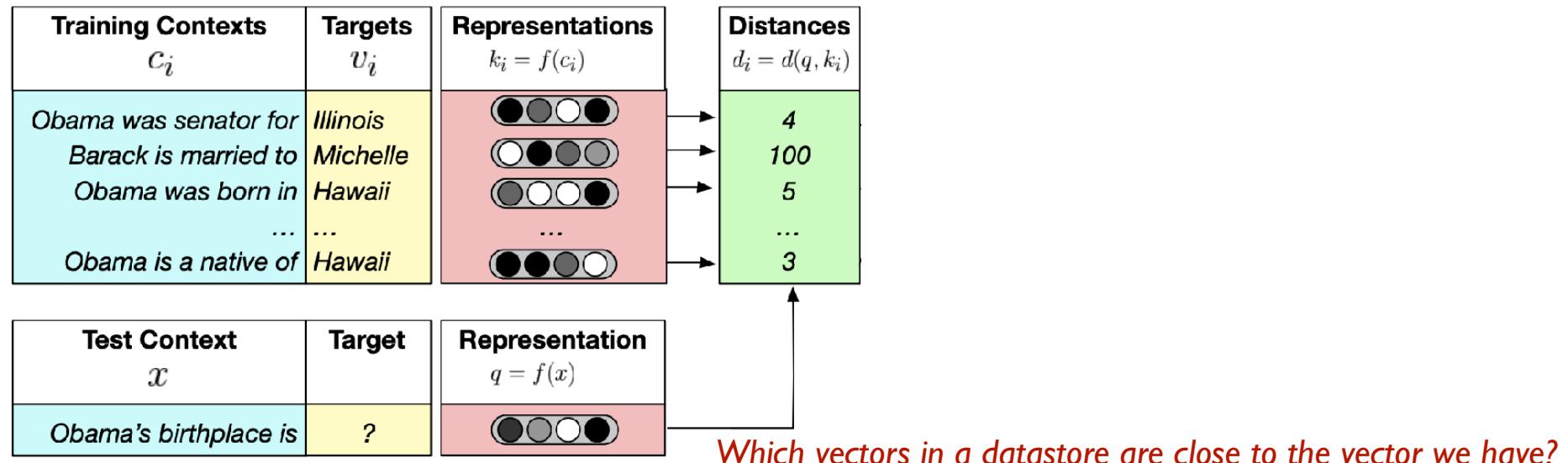
=

Which prefixes in a datastore are close to the prefix we have?

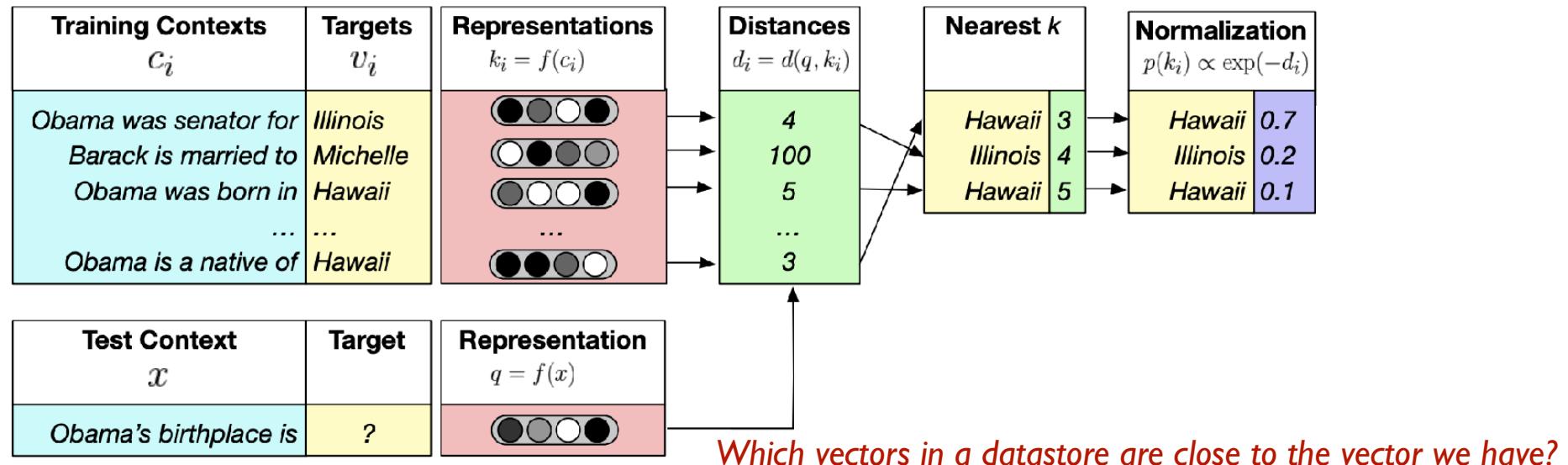
=

Which vectors in a datastore are close to the vector we have?

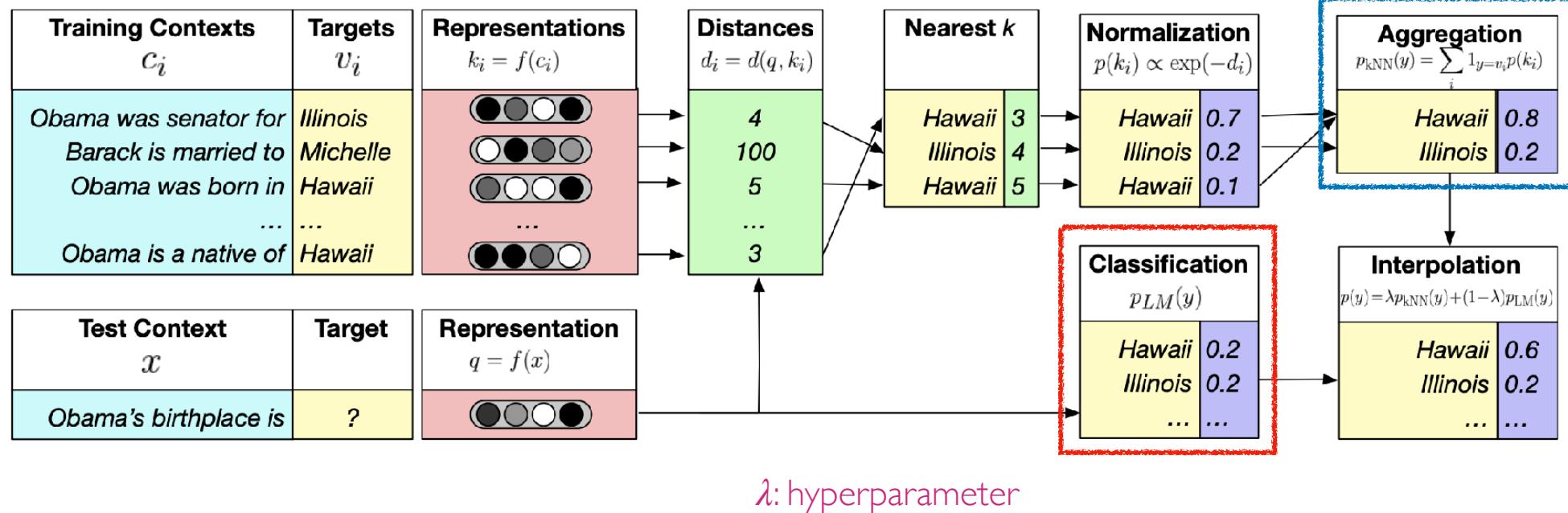
kNN-LM (Khandelwal et al. 2020)



kNN-LM (Khandelwal et al. 2020)



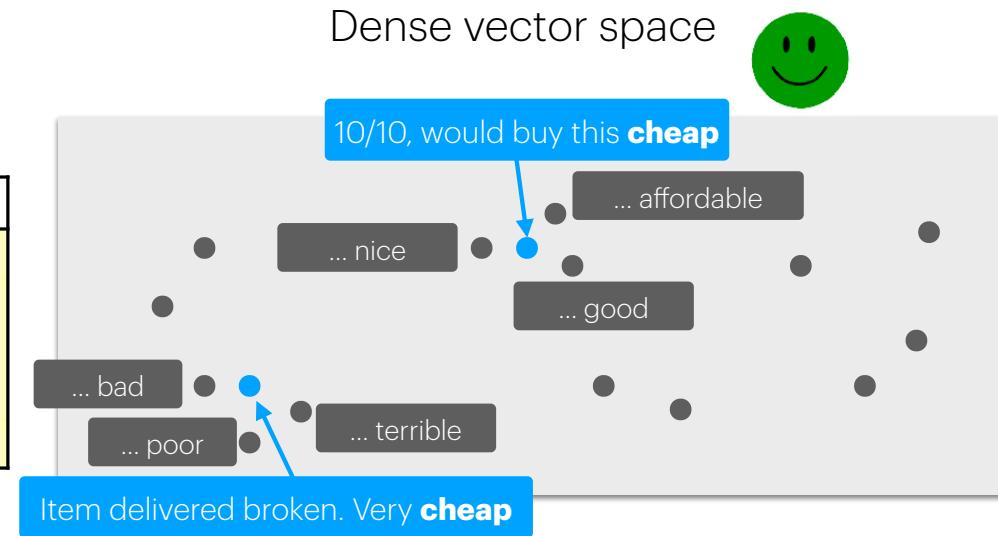
kNN-LM (Khandelwal et al. 2020)



$$P_{kNN-LM}(y|x) = \underbrace{(1 - \lambda)P_{LM}(y|x)}_{\text{Classification}} + \underbrace{\lambda P_{kNN}(y|x)}_{\text{Aggregation}}$$

kNN-LM - why?

Training contexts	Targets
10/10, would buy this	cheap
Item delivered broken. Very	cheap
To check the version of PyTorch, you can use	torch
You are permitted to bring a	torch
A group of infections ... one of the	torch



Summary

	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

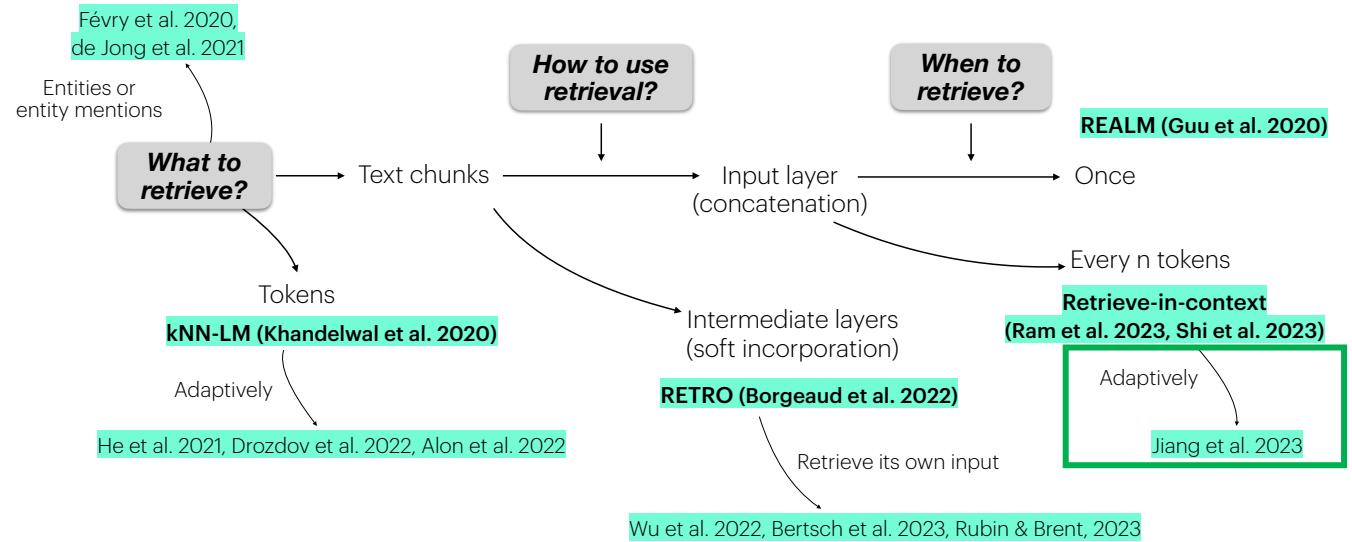


More fine-grained; Can be better at rare patterns & *out-of-domain*
Can be very efficient (as long as kNN search is fast) **(Wikipedia) 13M vs. 4B**



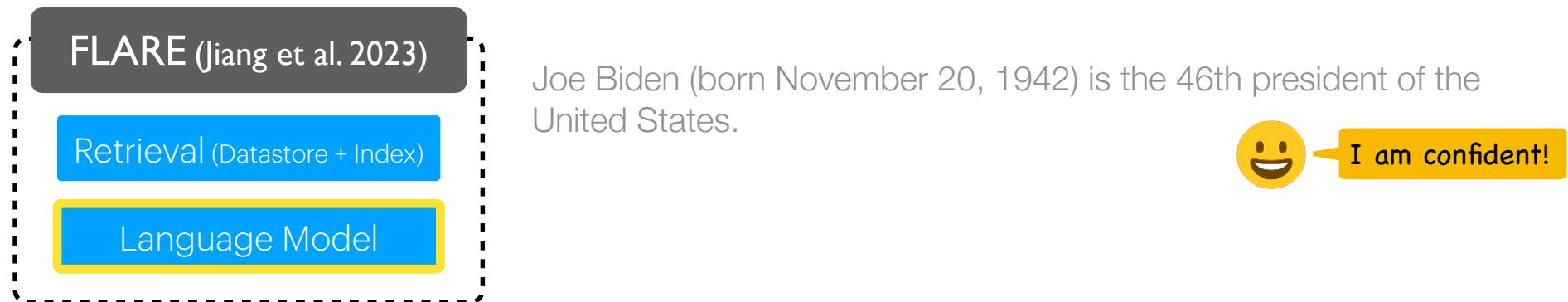
Datastore is expensive in space: given the same data, # text chunks vs. # tokens
No cross attention between input and retrieval results

Can we retrieve only when needed?



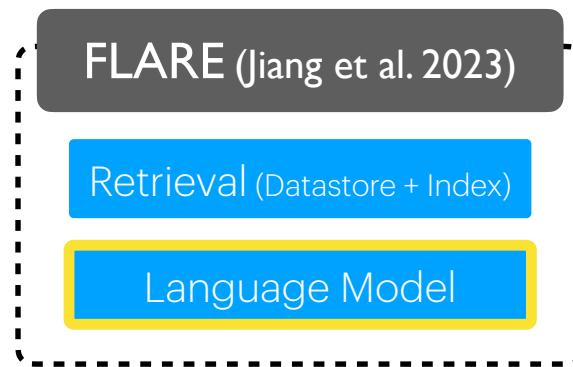
Adaptive retrieval of chunks

- Judge necessity



Adaptive retrieval of chunks

- Judge necessity



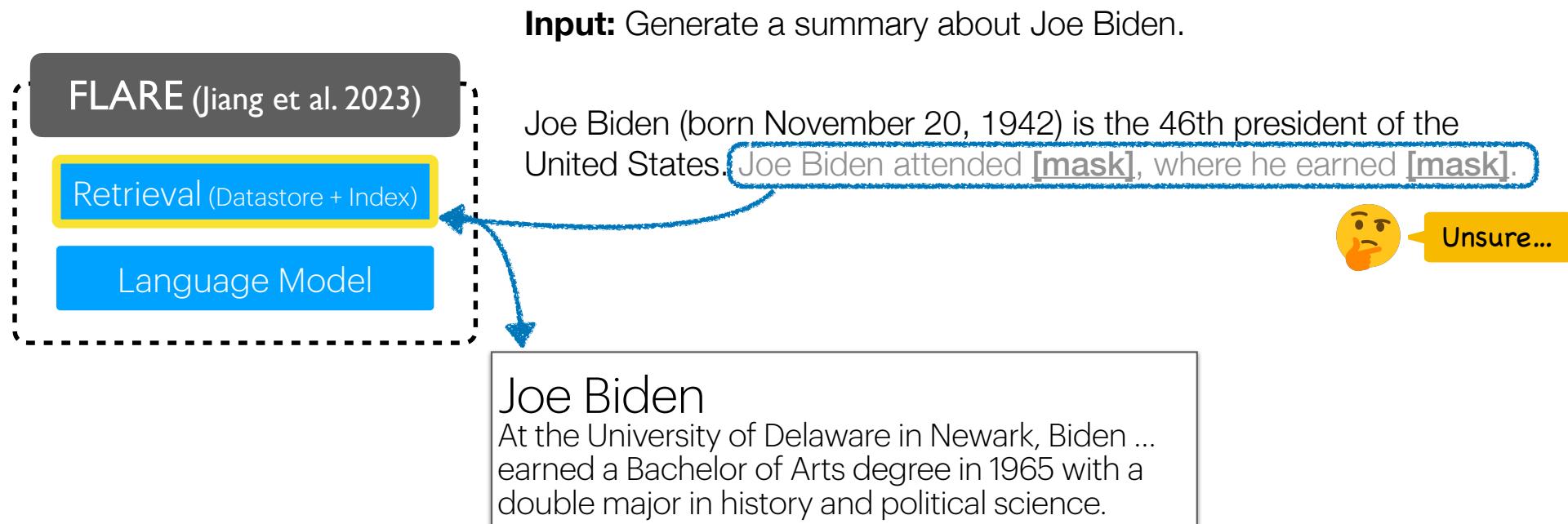
Input: Generate a summary about Joe Biden.

Joe Biden (born November 20, 1942) is the 46th president of the United States. Joe Biden attended the University of Pennsylvania, where he earned a law degree.



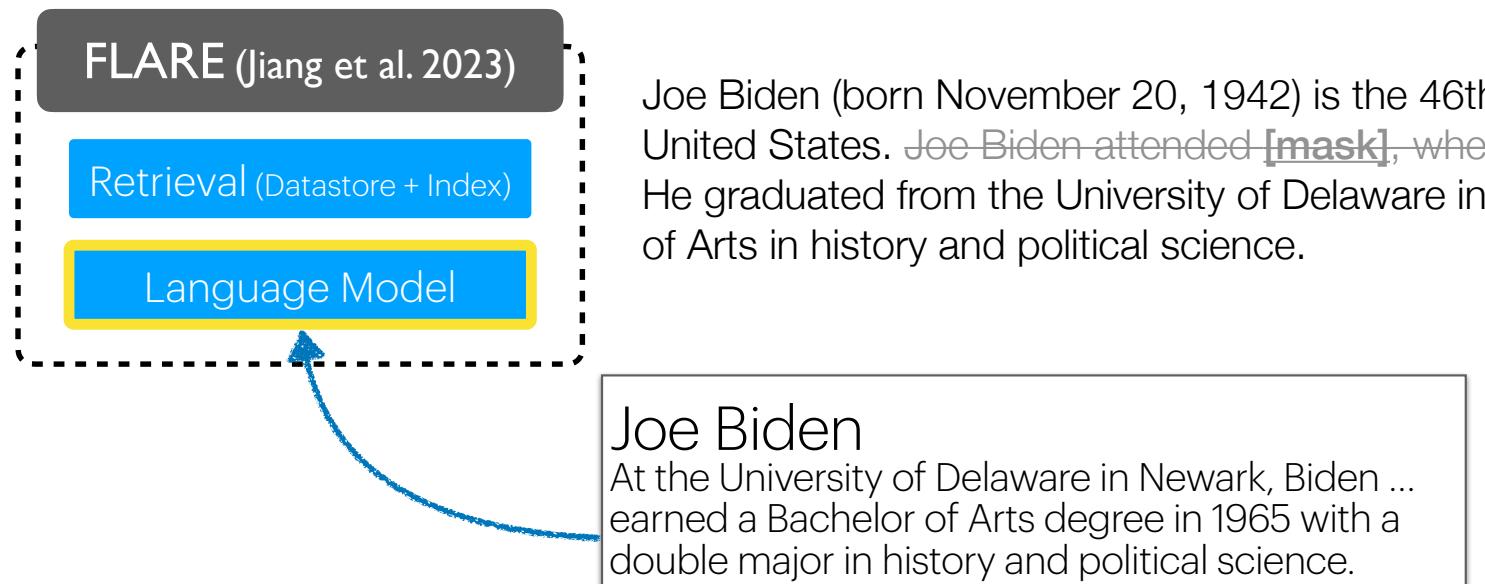
Adaptive retrieval of chunks

- Judge necessity



Adaptive retrieval of chunks

- Judge necessity

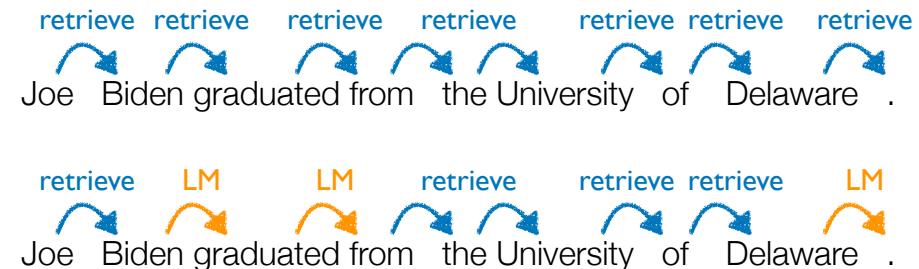


Also with kNN-LM!

[kNN-LM] does it in a per-token basis instead of sentence/chunk level

Adaptive retrieval of *tokens*

- Judge necessity



Also with kNN-LM!

Adaptive retrieval of *tokens*

- Judge necessity



$$P_{k\text{NN-LM}}(y|x) = \underline{(1 - \lambda(x))} P_{\text{LM}}(y|x) + \underline{\lambda(x)} P_{k\text{NN}}(y|x)$$

A function of the input x

$\rightarrow \lambda = 0$ if $\lambda < \gamma$

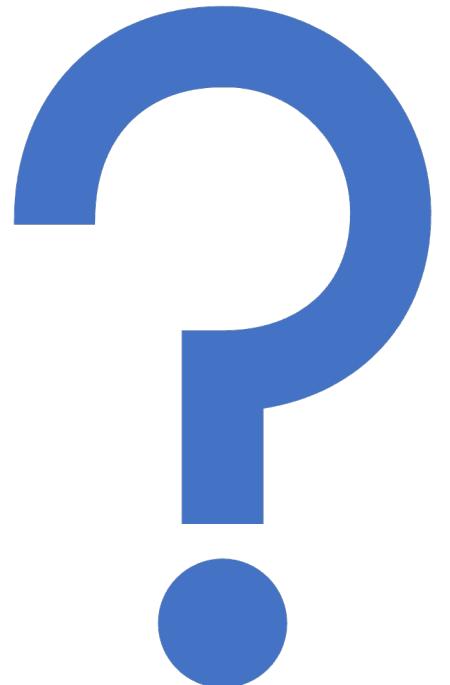
He et al. 2021. "Efficient Nearest Neighbor Language Models"

103

predict the interpolation parameter $\lambda(x)$:
trade-off between memorization (knn-LM) and LM predictions

Summary

	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 (adaptive)
Adaptive kNN-LM (He et al 2021, Alon et al 2022, etc)	Tokens	Output layer	Every n tokens (adaptive)



Questions?