

Question:

[https://bit.ly/
akari_ralm_lec](https://bit.ly/akari_ralm_lec)



Scan me

Retrieval-augmented Language Models

Akari Asai

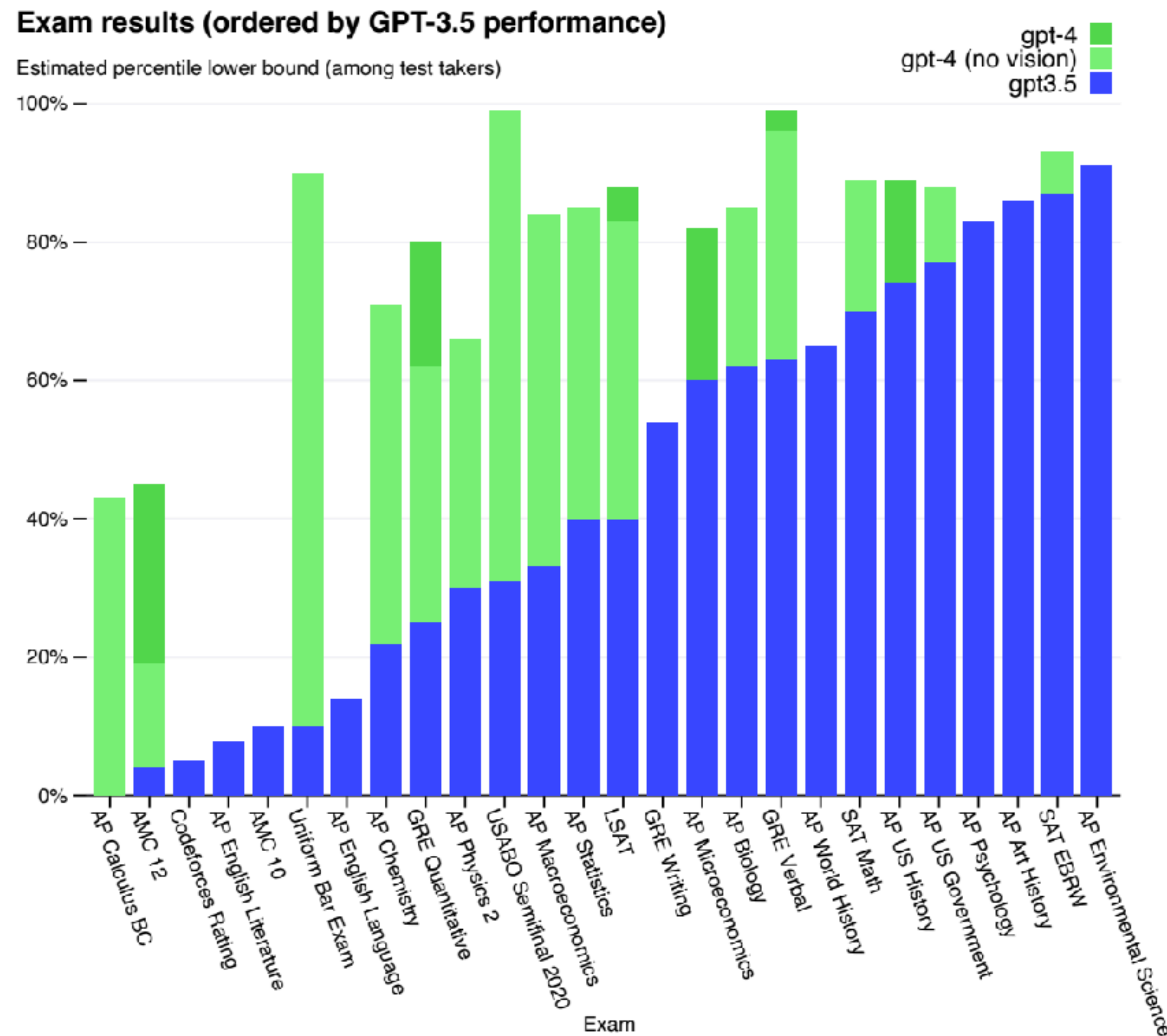
<https://akariasai.github.io/>

March 1, 2024 at UW CSE 517 / 447 NLP

Slides adapted from ACL 2023 tutorial by Akari Asai, Sewon Min, Zexuan Zhong, Danqi Chen

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

Knowledgable and versatile Language Models (LMs)



Knowledgable and versatile Language Models (LMs)



You

Write a poem about large language models, including the following three words: coffee, croissants and banana (*February 18, 2024)



ChatGPT

In the realm where thoughts converge and flow,
Resides a wonder, a grand tableau,
Where words dance in intricate array,
In the domain of the model's sway.

From the depths of bytes, they come alive,
In the language vast, where dreams contrive,
Coffee fueled, the model's might,
Awakens realms in the dead of night.

Croissants of syntax, buttered with care,
Layered meanings, rich and rare,
Each line a pastry, crisp and sweet,
In the banquet of words, where senses meet.

Banana shaped metaphors, ripe and bold,
In the orchard of thoughts, they unfold,
Symbolic fruits of the mind's terrain,
Nourishing narratives, ripe with gain.

How do such parametric LMs work?

$$P(x_n | x_1, x_2, \dots, x_{n-1})$$

Tr The capital city of Ontario is Toronto



Large-scale pre-training corpus
(e.g., 1T tokens)



Language model (Transformers)

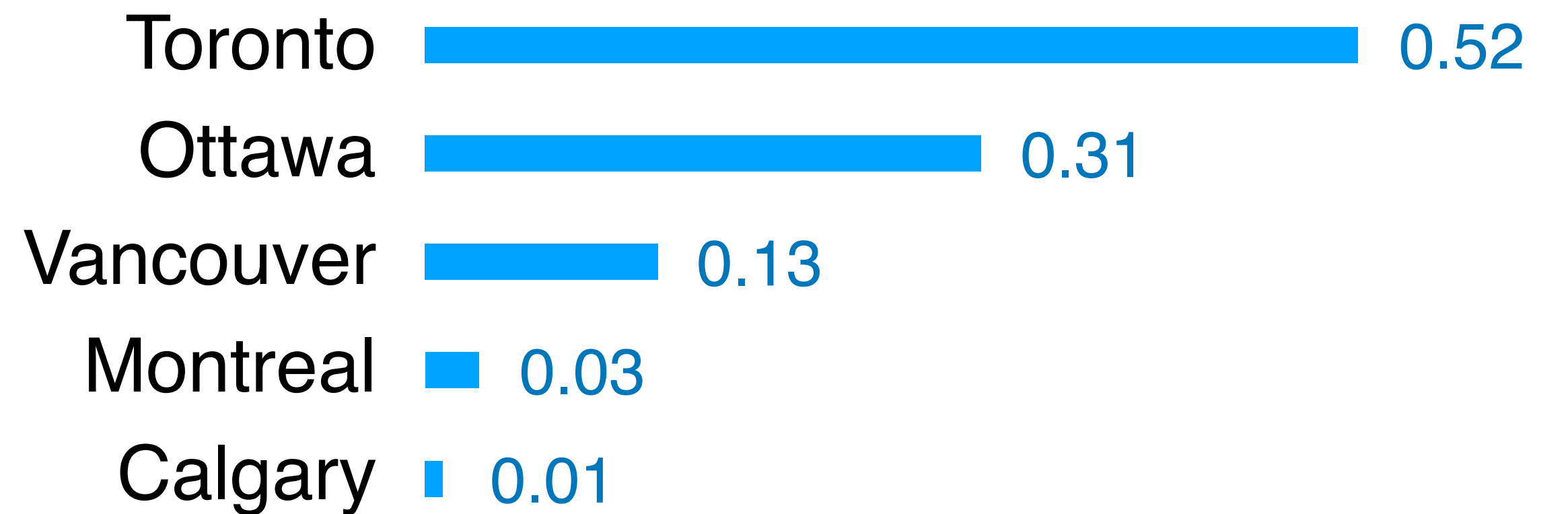
The capital city of Ontario is

x_1

x_2

...

x_{n-1}



...

How do such parametric LMs work?

$$P(x_n | x_1, x_2, \dots, x_{n-1})$$

Large-scale pre-training corpus
(e.g., 11 tokens)

No explicit access to large-scale text data



Language model (Transformers)

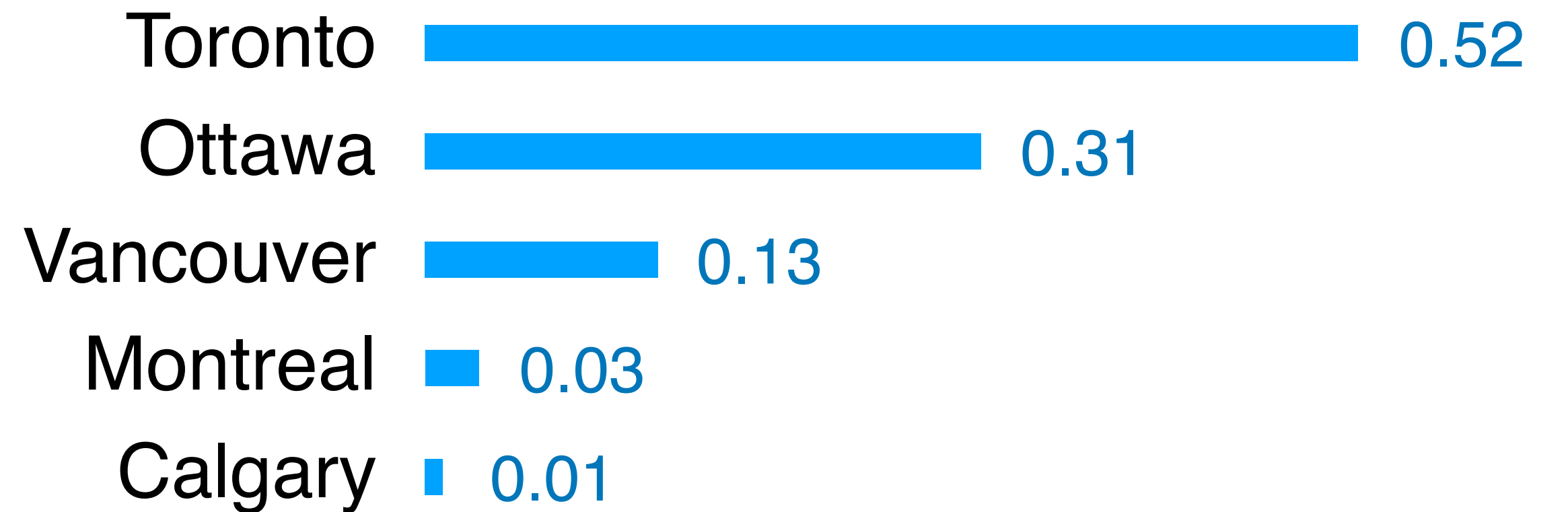
Ontario's capital city,

x_1

x_2

...

x_{n-1}



...

Hallucinations in LM outputs



You

Explain Akari Asai, including the list of her famous papers. (*Febru



ChatGPT

Akari Asai is a prominent figure in the field of natural language processing with significant contributions to various aspects of NLP research. She has made contributions particularly in the areas of question answering, dialogue systems, and machine learning.

Here are some of her famous papers:

1. "Learning to Retrieve Reasoning Paths over Wikipedia Graph for Question Answering" - In this paper, Asai and her co-authors proposed a method for learning to retrieve reasoning paths over the Wikipedia graph to improve question answering systems.
2. "Large-Scale Multi-Document Summarization Dataset and Evaluation Benchmark" - This paper introduces a large-scale dataset for multi-document summarization, along with an evaluation benchmark, providing valuable resources for advancing research in this area.

Akari Asai FOLLOW

Multi-News: a Large-Scale Multi-Document Summarization Dataset and Abstractive Hierarchical Model

Alexander R. Fabbri Irene Li
Tianwei She Suyi Li Dragomir R. Radev

Department of Computer Science
Yale University
{alexander.fabbri,irene.li,tianwei.she,suyi.li,dragomir.radev}@yale.edu

Learning to Retrieve Reasoning Paths over Wikipedia Graph for Question Answering 286 2020
A Asai, K Hashimoto, H Hajishirzi, R Socher, C Xiong
International Conference on Learning Representations (ICLR)



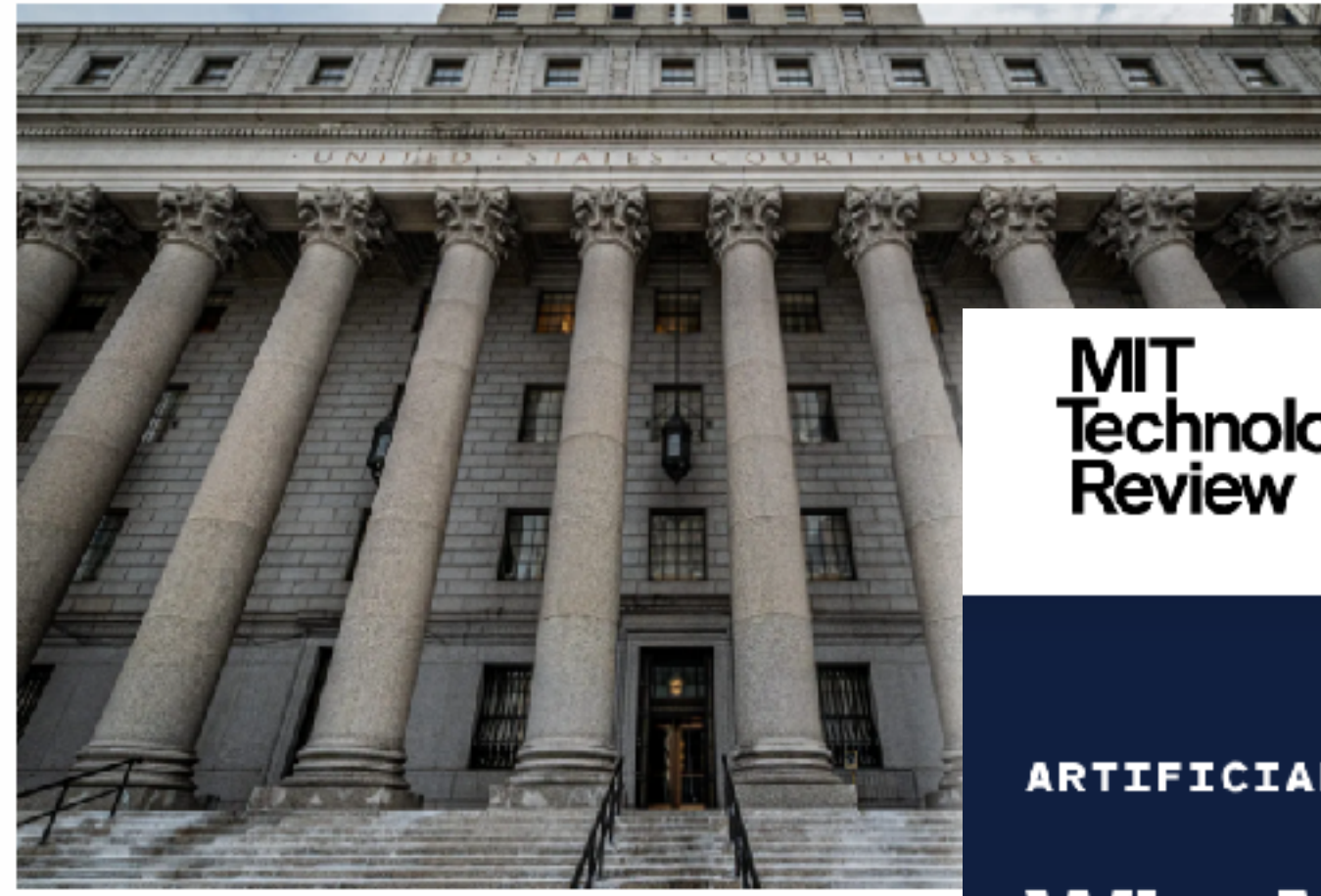
Catastrophic Errors as Results of LM Hallucinations

TECH · LAW

Humiliated lawyers fined \$5,000 for submitting ChatGPT hallucinations in court: 'I heard about this new site, which I falsely assumed was, like, a super search engine'

BY RACHEL SHIN

June 23, 2023 at 9:41 AM PDT



Lawyers who filed legal documents with false citations generated by ChatGPT have been fined \$5,000 each. (Erik McGregor - LightRocket/Getty Images)

MIT
Technology
Review

Featured Topics Newsletters Events Podcasts

SIGN IN

ARTIFICIAL INTELLIGENCE

Why Meta's latest large language model survived only three days online

Galactica was supposed to help scientists. Instead, it mindlessly spat out biased and incorrect nonsense.

By Will Douglas Heaven

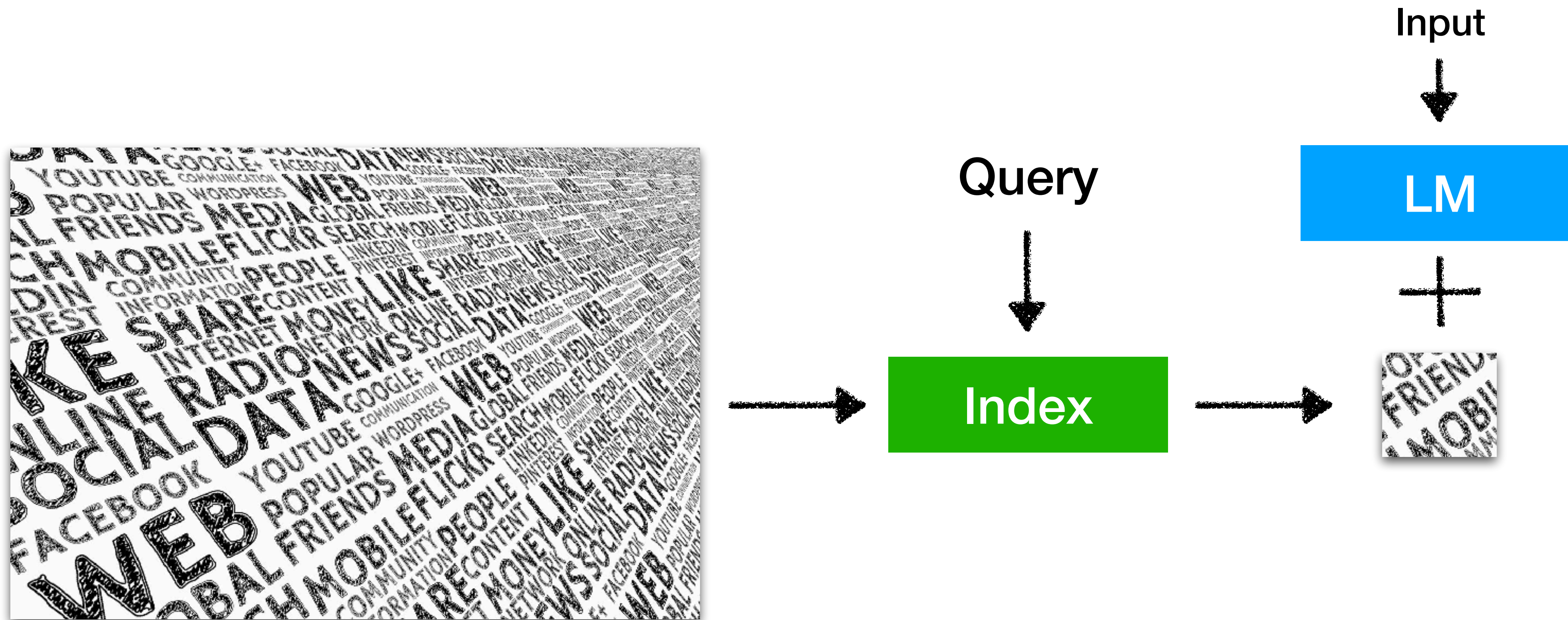
November 18, 2022

Air Canada must honor re... invented by airline's chatb...

Air Canada appears to have quietly killed its costly chatbot support.

ASHLEY BELANGER - 2/16/2024, 12:12 PM

Inference: Datastore



Datastore

Raw text corpus

At least billions~trillions of tokens

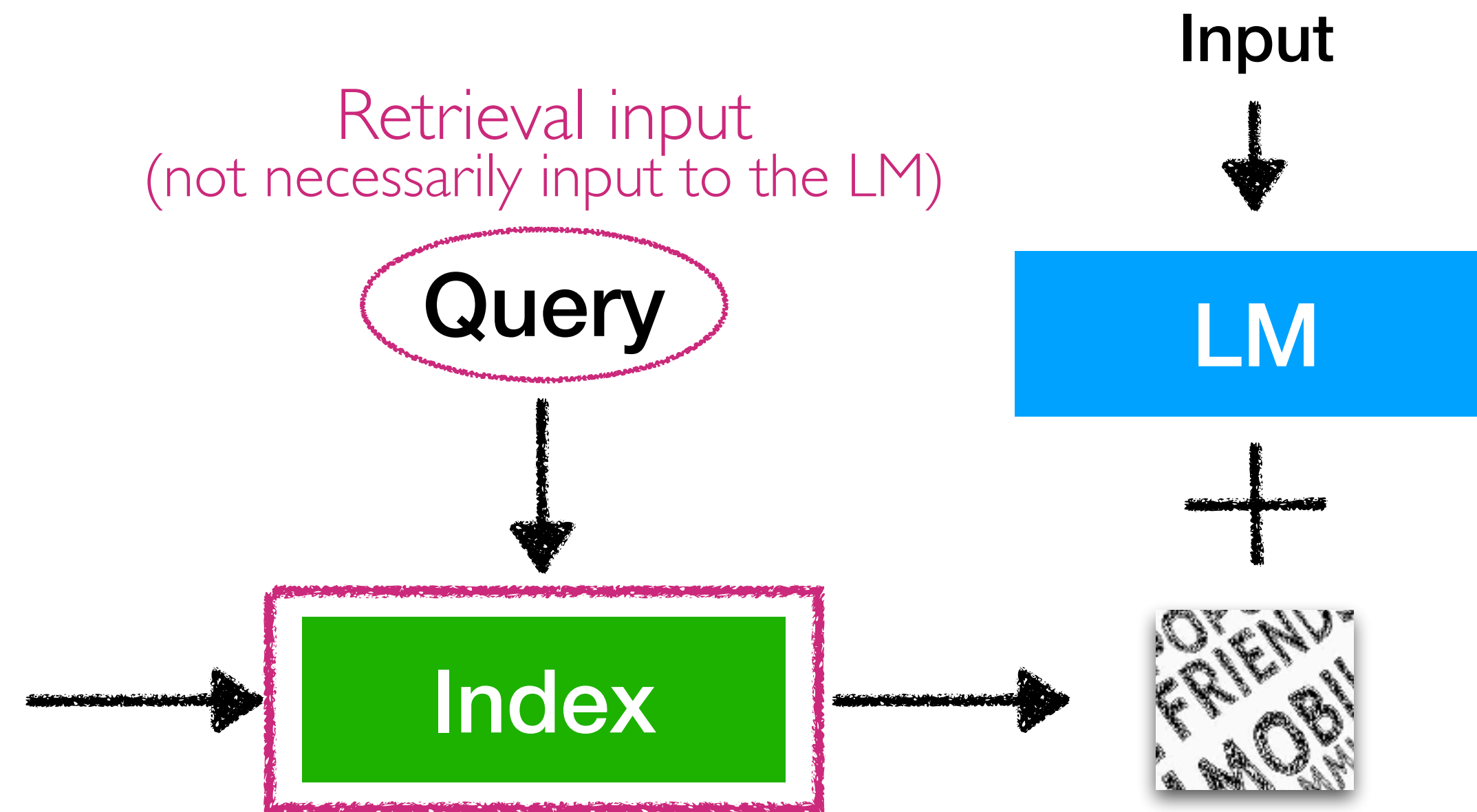
Not labeled datasets

Not structured data (knowledge bases)

Inference: Index



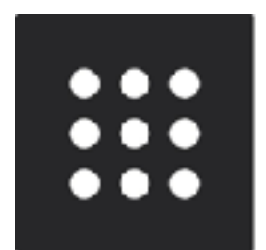
Datastore



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

Retrieval-augmented LMs are now widely used!

YOU



Perplexity



Today's outline

Question:

[https://bit.ly/
akari_ralm_lec](https://bit.ly/akari_ralm_lec)



Scan me

Why do we need retrieval-augmented LMs?

Architectures of retrieval-augmented LMs (Inference)

Training of retrieval-augmented LMs

Limitations and future directions

Today's outline

Question:

[https://bit.ly/
akari_ralm_lec](https://bit.ly/akari_ralm_lec)



Scan me

Why do we need retrieval-augmented LMs?

Architectures of retrieval-augmented LMs (Inference)

Training of retrieval-augmented LMs

Limitations and future directions

A: Because retrieval-augmented LMs **can solve many core limitations** of parametric LMs!

Core limitations of parametric LMs

Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



You

Explain Akari Asai, including the list of her famous papers. (*February 18, 2024)

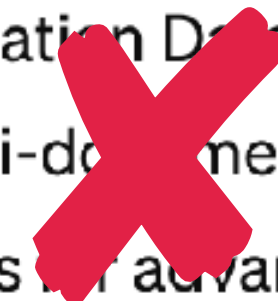


ChatGPT

Akari Asai is a prominent figure in the field of natural language processing (NLP), known for her significant contributions to various aspects of NLP research. She has made notable contributions particularly in the areas of question answering, dialogue systems, and machine learning.

Here are some of her famous papers:

1. "Large-Scale Multi-Document Summarization Dataset and Evaluation Benchmark" - This paper introduces a large-scale dataset for multi-document summarization, along with an evaluation benchmark, providing valuable resources for advancing research in this area.
2. "Large-Scale Multi-Document Summarization Dataset and Evaluation Benchmark" - This paper introduces a large-scale dataset for multi-document summarization, along with an evaluation benchmark, providing valuable resources for advancing research in this area.



Core limitations of parametric LMs

Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Explain Akari Asai, including the list of her famous papers.



Language model



Her most famous paper is “*Large-Scale Multi-Document Summarization Dataset and Evaluation Benchmark*”

Core limitations of parametric LMs

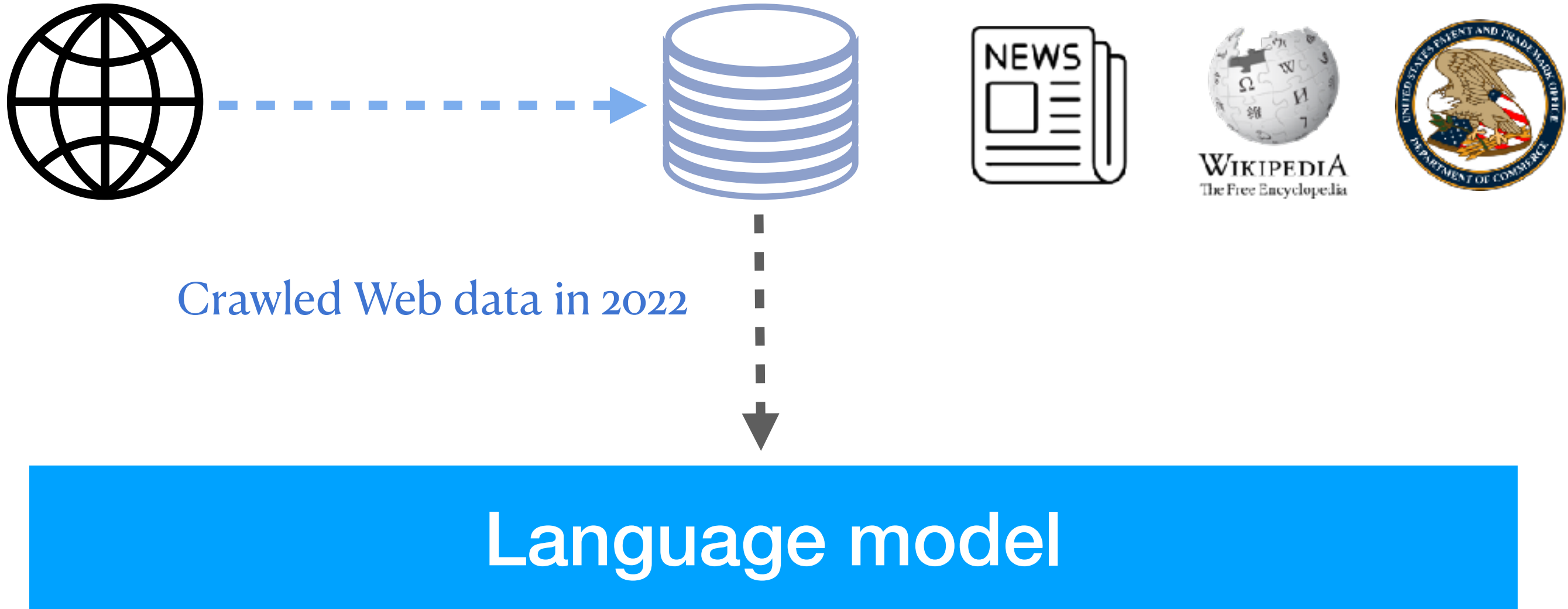
Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



ChatGPT

I'm sorry, but I don't have access to real-time information including events beyond January 2022.

Core limitations of parametric LMs

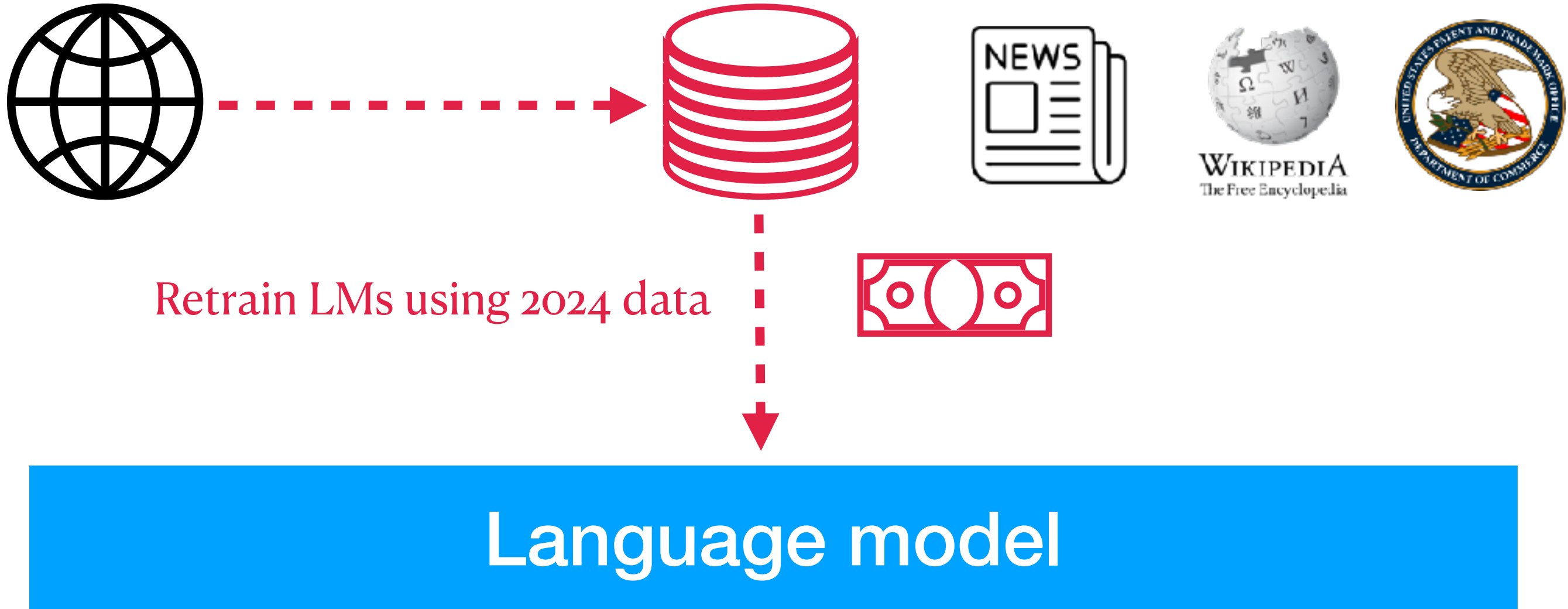
Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



ChatGPT

I'm sorry, but I don't have access to real-time information including events beyond January 2022.

Core limitations of parametric LMs

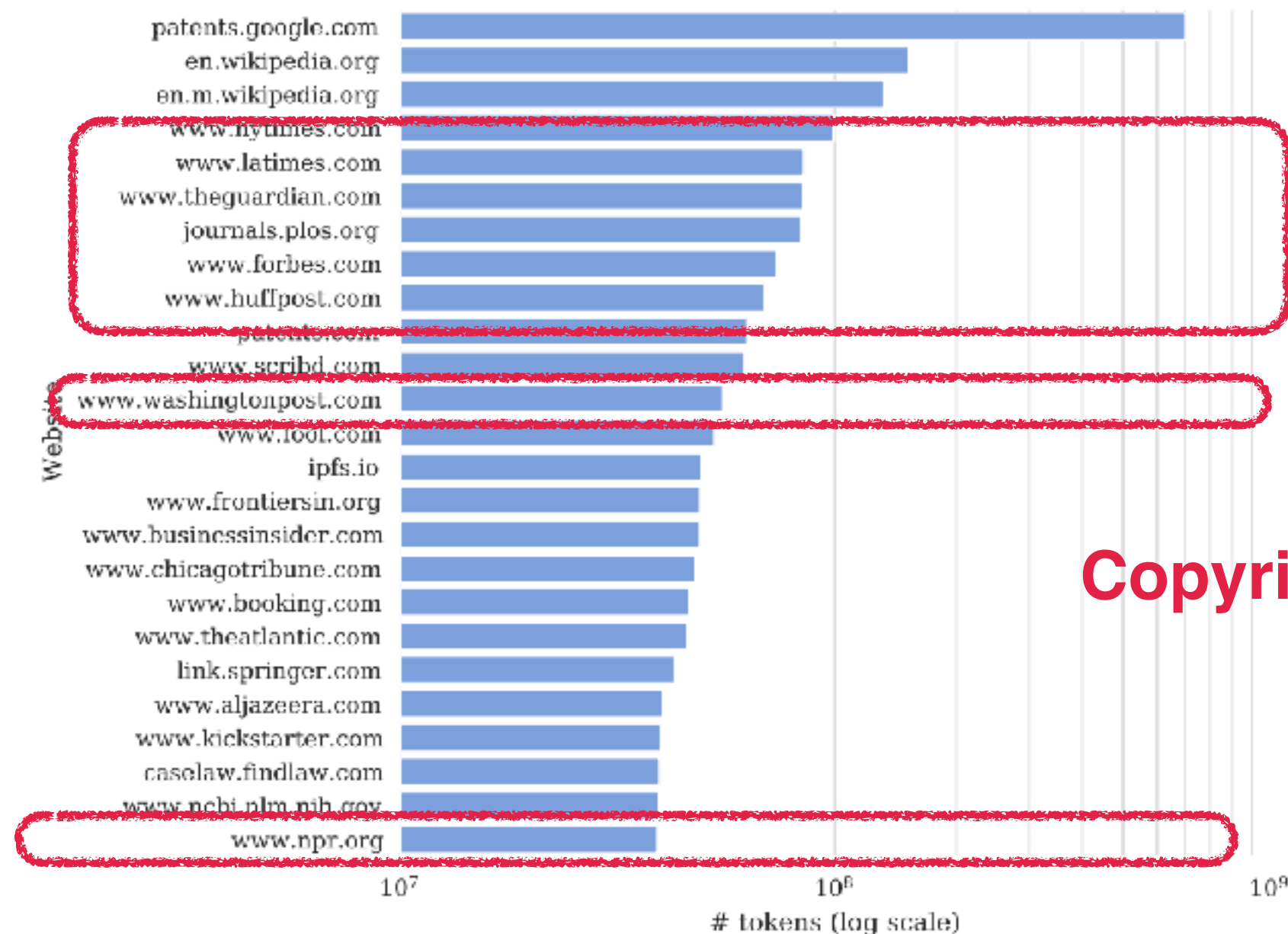
Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



Copyright-protected data?

Dodge et al., Documenting Large Webtext Corpora: A Case Study on the Colossal Clean Crawled Corpus. EMNLP 2021.

Core limitations of parametric LMs

Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

B. Defendants' GenAI Products

1. A Business Model Based on Mass Copyright Infringement

57. Despite its early promises of altruism, OpenAI quickly became a multi-billion-dollar for-profit business built in large part on the unlicensed exploitation of copyrighted works belonging to The Times and others. Just three years after its founding, OpenAI shed its exclusively

Plaintiff The New York Times Company ("The Times"), by its attorneys Susman Godfrey LLP and Rothwell, Figg, Ernst & Manbeck, P.C., for its complaint against Defendants Microsoft Corporation ("Microsoft") and OpenAI, Inc., OpenAI LP, OpenAI GP, L.P., OpenAI OpCo LLC, OPENAI CORPORATION, LLC, OAI CORPORATION, LLC, HOLDINGS, LLC, (collectively "OpenAI" and, with Microsoft, "Defendants"), alleges as follows:

I. NATURE OF THE ACTION

1. Independent journalism is vital to our democracy. It is also increasingly rare and valuable. For more than 170 years, The Times has given the world deeply reported, expert, independent journalism. Times journalists go where the story is, often at great risk and cost, to inform the public about important and pressing issues. They bear witness to conflict and disasters, provide accountability for the use of power, and illuminate truths that would otherwise go unseen. Their essential work is made possible through the efforts of a large and expensive organization that provides legal, security, and operational support, as well as editors who ensure their journalism meets the highest standards of accuracy and fairness. This work has always been important. But

New York Times lawsuits against OpenAI

Core limitations of parametric LMs

Hallucinations

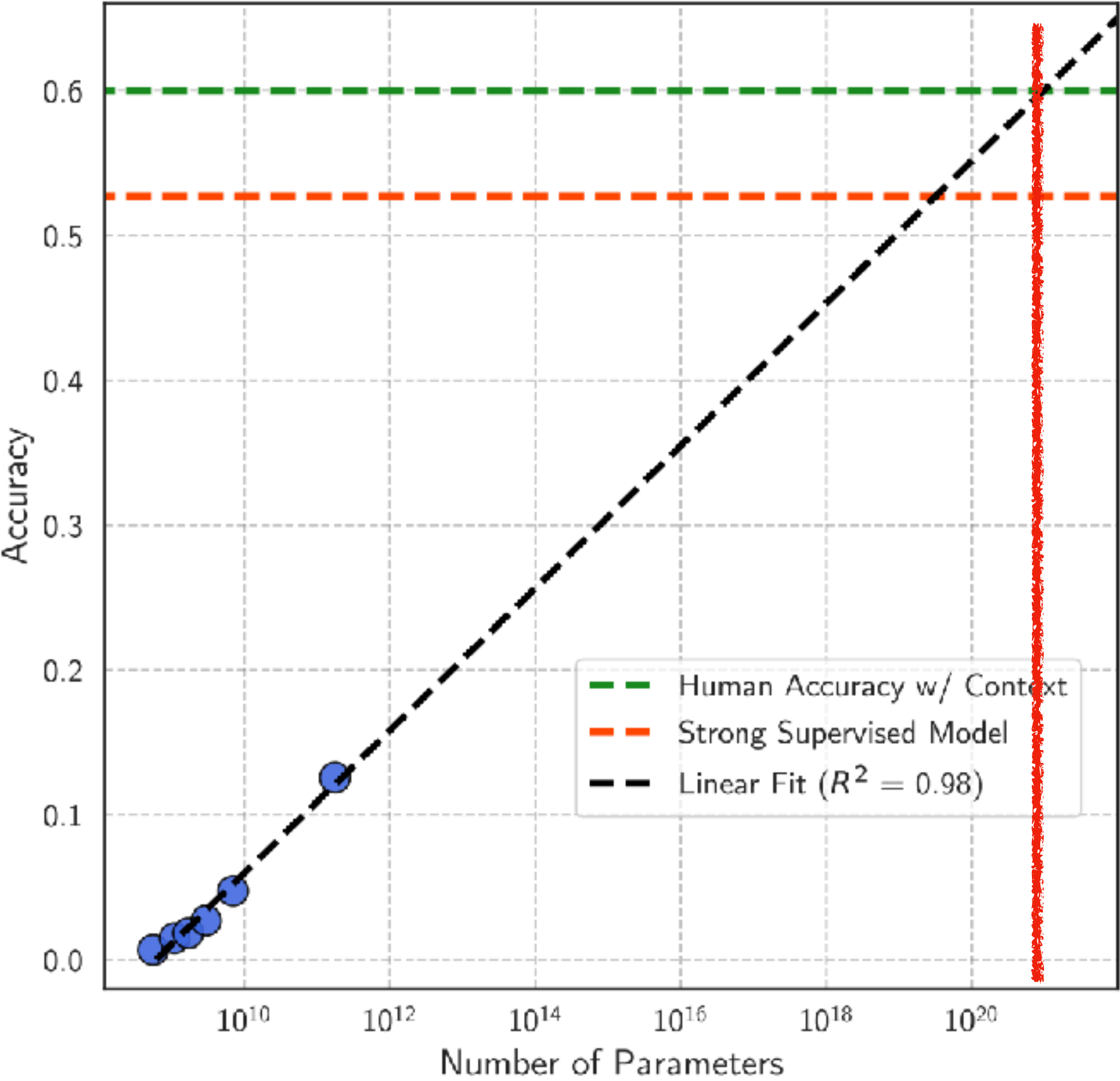
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Long-tail QA performance



100 quintillion parameters required to reach human performance

Kandpal et al., Large Language Models Struggle to Learn Long-Tail Knowledge. ICML 2023.

Q: So how can **retrieval**-augmented LMs solve those challenges?

How retrieval-augmented LMs solve the issues?

Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Significant improvements across model scale, with larger gain with smaller LMs

QA



How retrieval-augmented LMs solve the issues?

Hallucinations

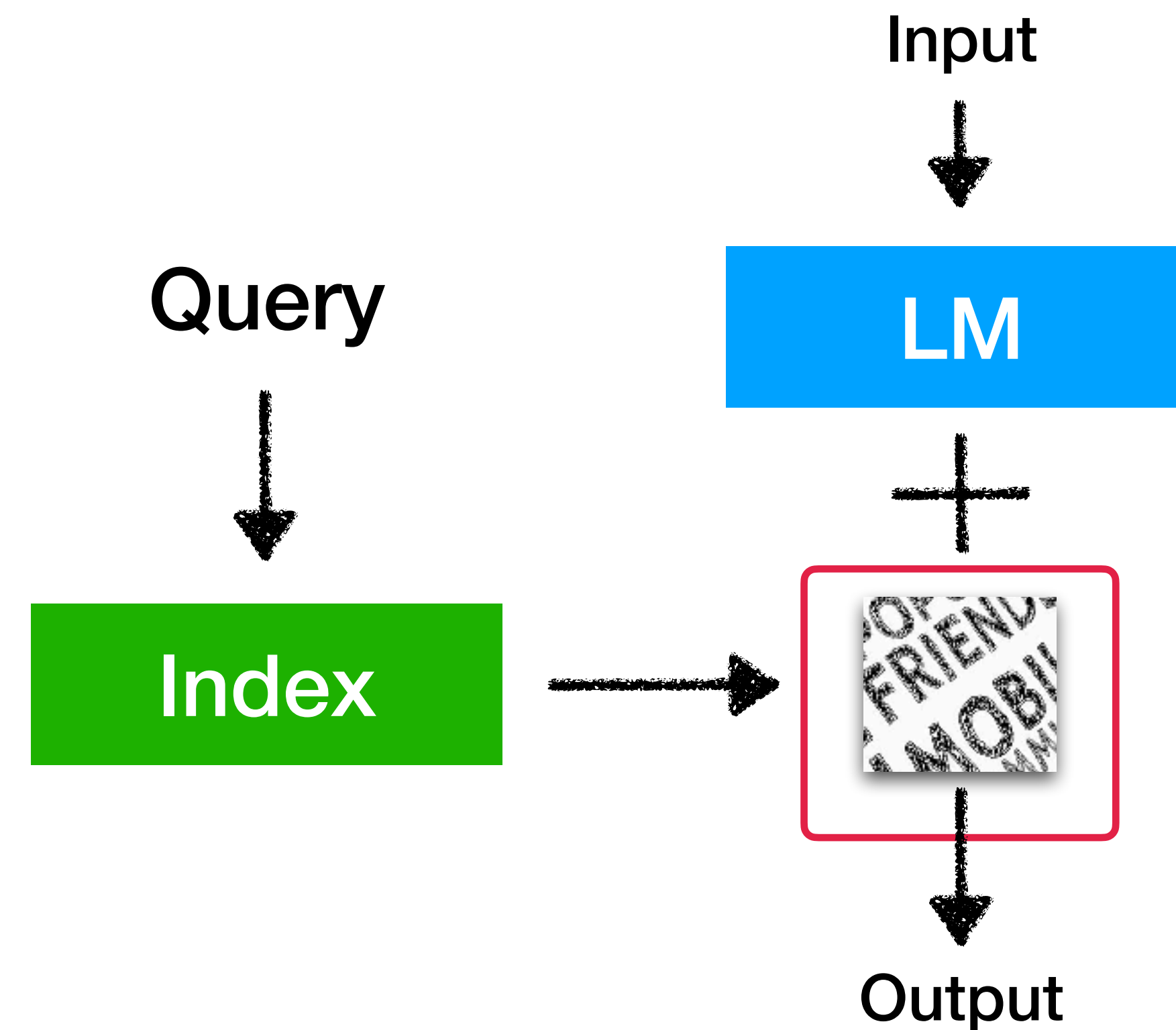
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Retrieved text can be used as attributions



How retrieval-augmented LMs solve the issues?

Hallucinations

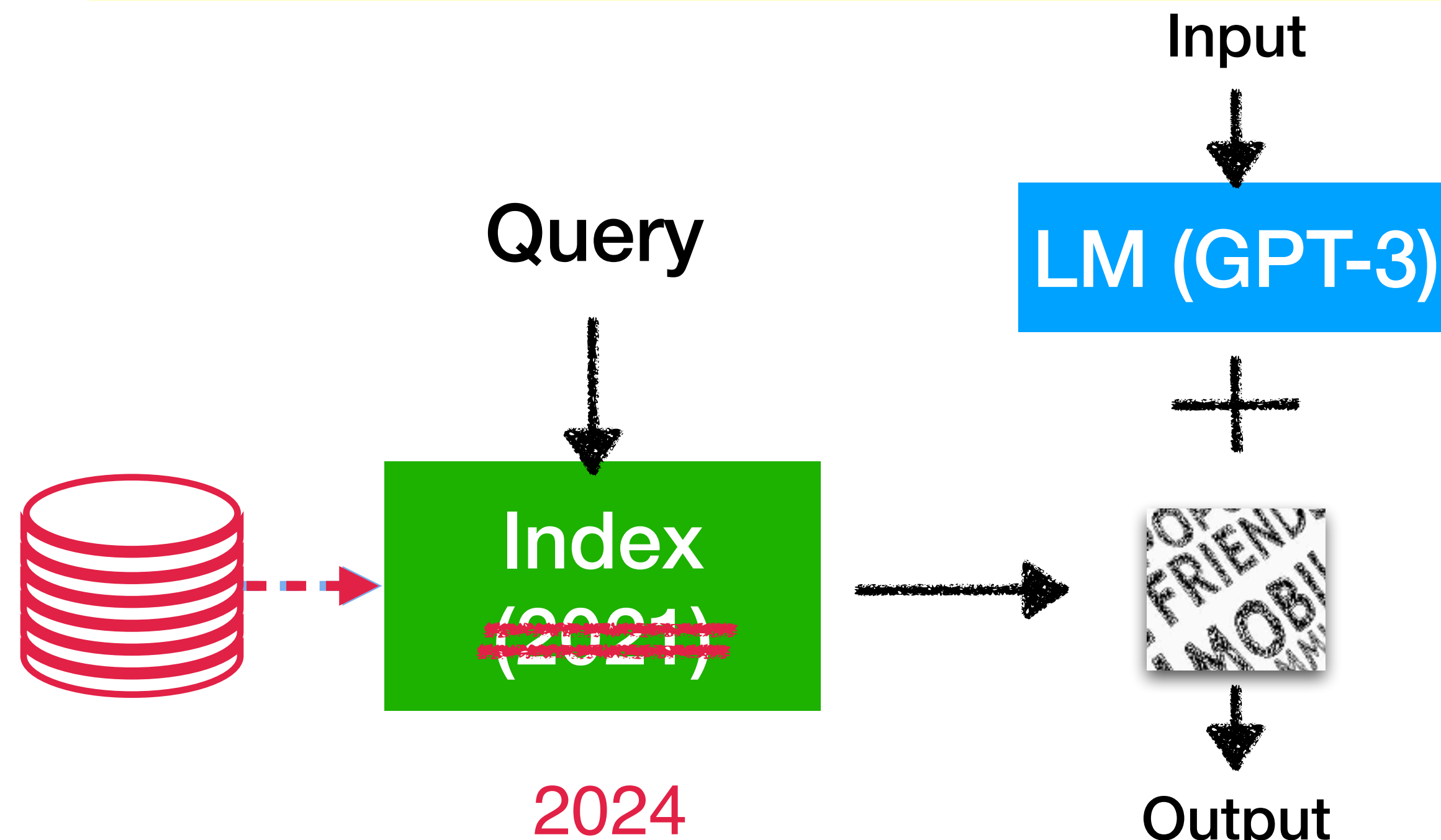
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Replacing datastores to catch up dynamically changing world without re-training



How retrieval-augmented LMs solve the issues?

Hallucinations

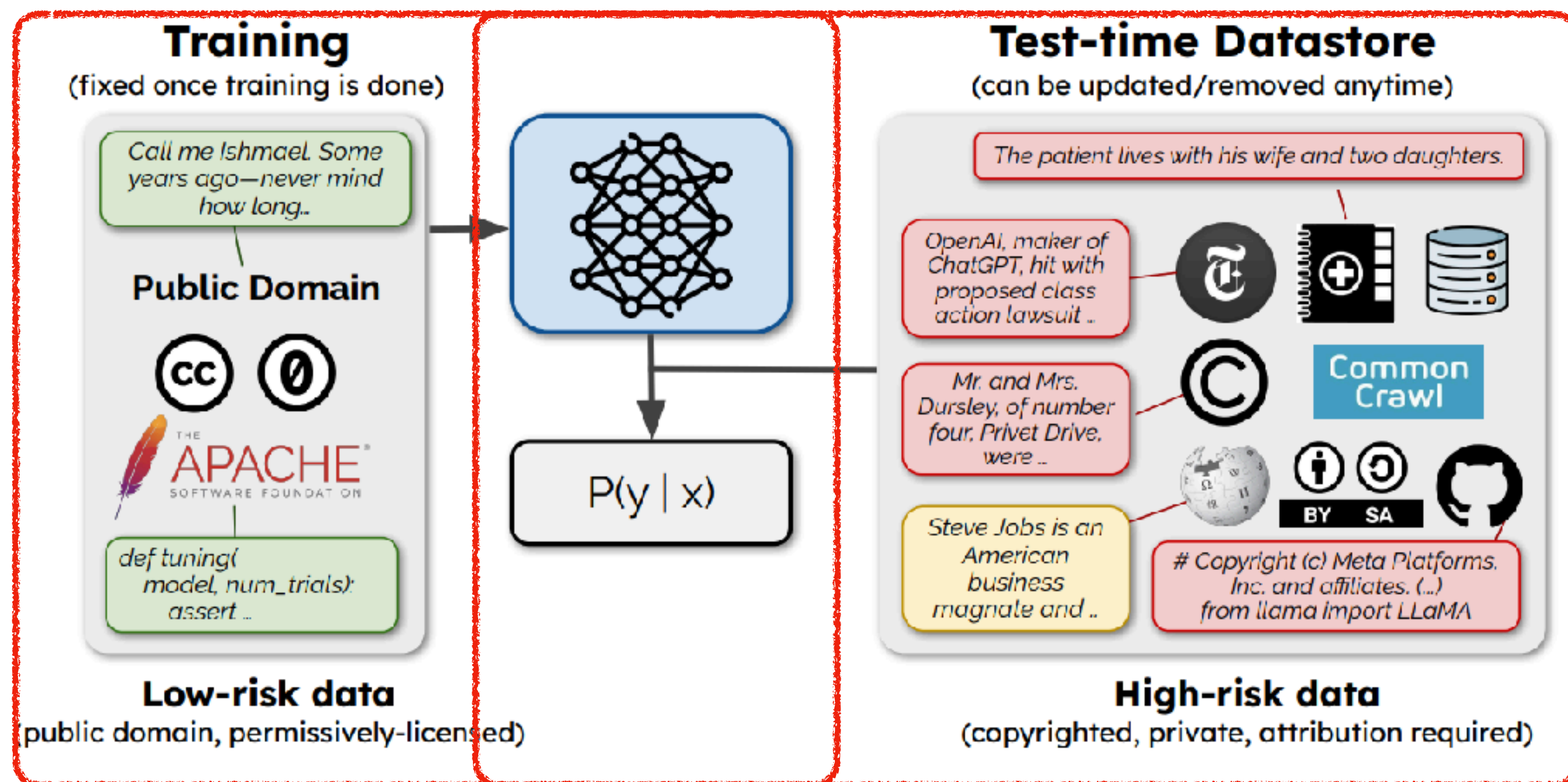
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Segregating copyright-sensitive data from pre-training data



Min* and Gururangan* et al., SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore. ICLR 2024.

How retrieval-augmented LMs solve the issues?

Hallucinations

Lack of attributions

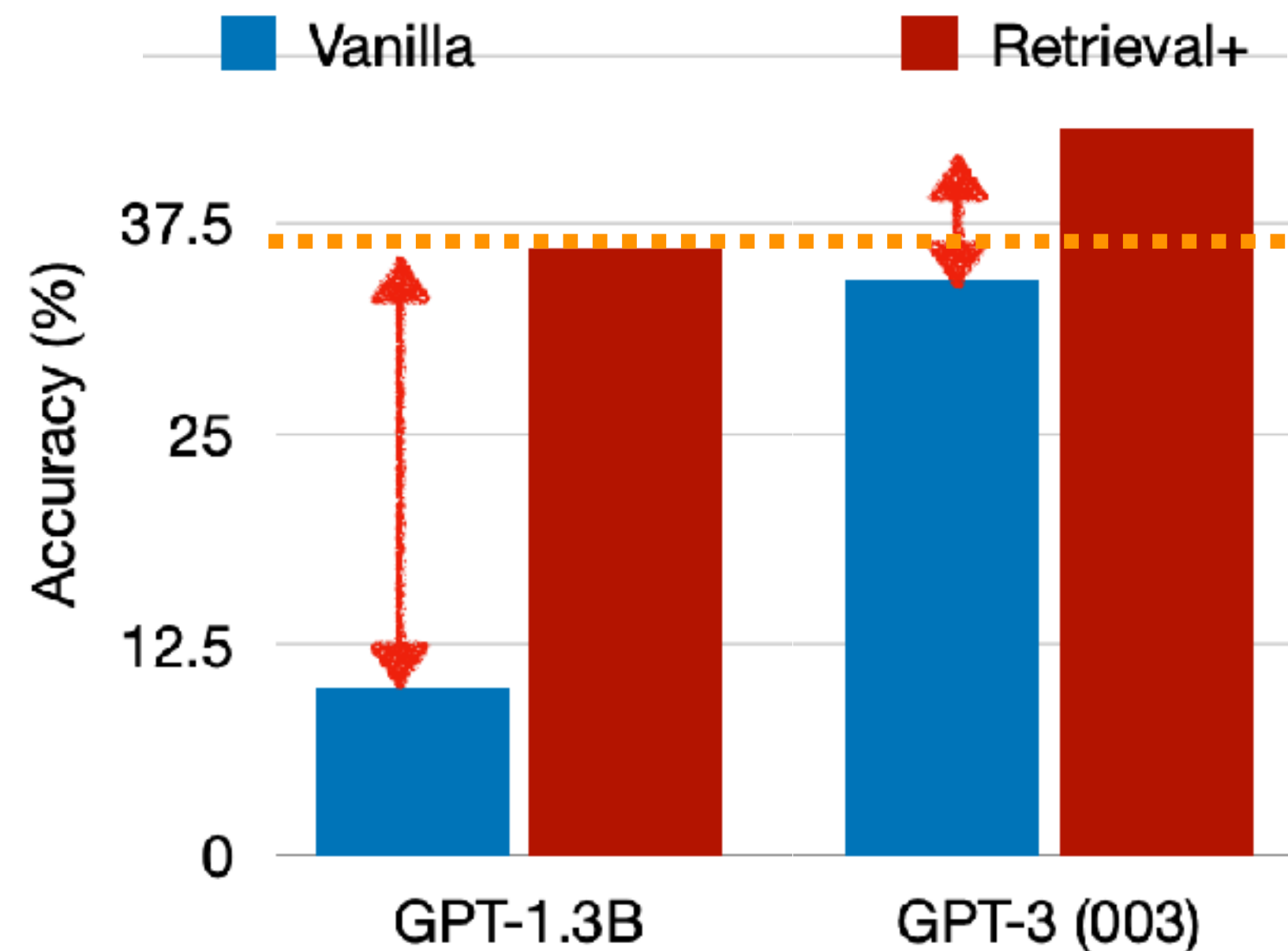
Costs of adaptations

Copyright / privacy

Large parameter size

Smaller LMs with retrieval outperform much larger LMs e.g., GPT-3

QA



Today's outline

Why do we need retrieval-augmented LMs?

Architectures of retrieval-augmented LMs (Inference)

Training of retrieval-augmented LMs

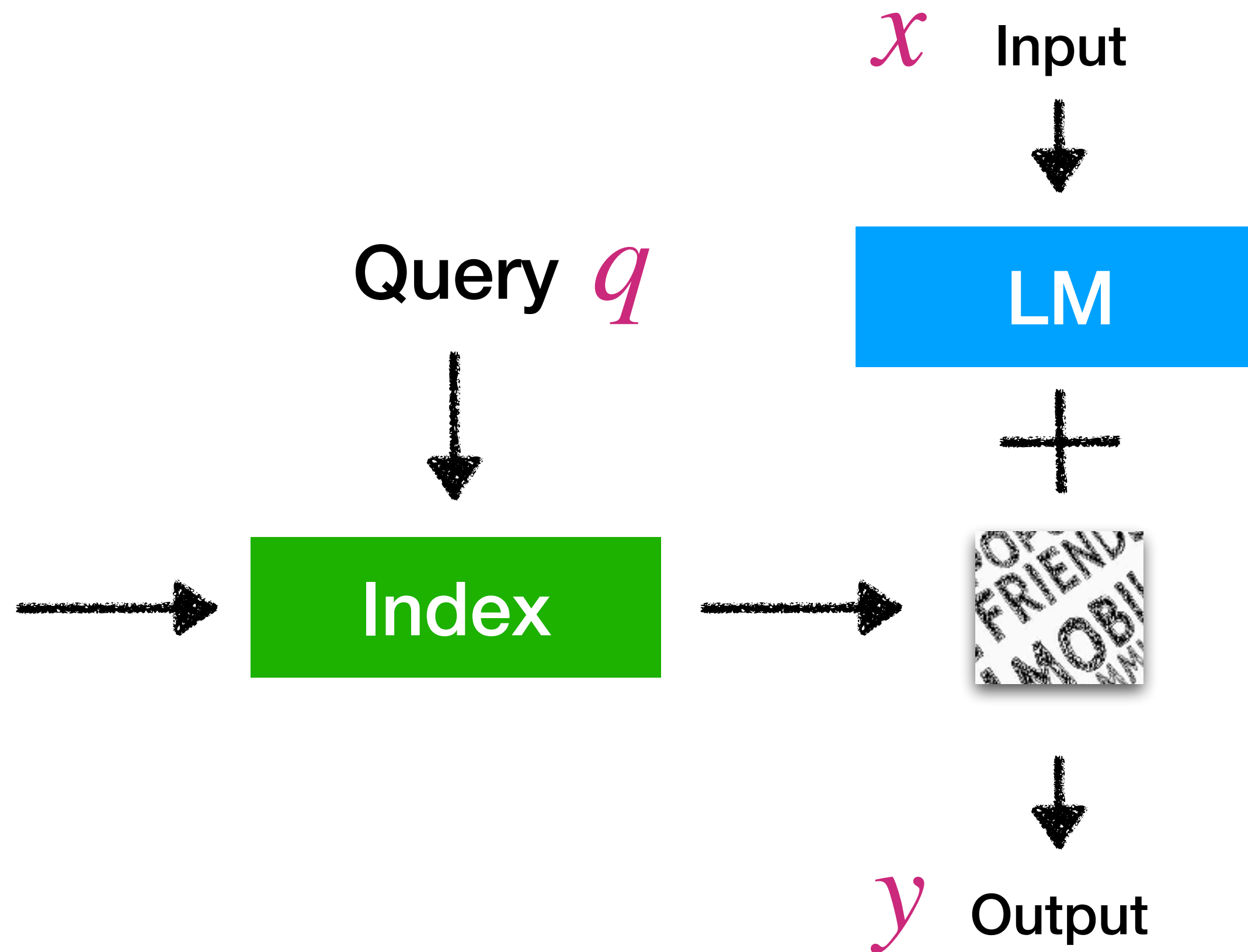
Limitations and future directions

Notations



Datastore

\mathcal{D}



Inference: Index

Goal: find a small subset of elements in a datastore that are the most similar to the query

sim: a similarity score between two pieces of text

Example $\text{sim}(i, j) = \text{tf}_{i,j} \times \log \frac{N}{\text{df}_i}$

$\text{tf}_{i,j}$ # of occurrences of i in j

N # of total docs

df_i # of docs containing i

Example $\text{sim}(i, j) = \text{Encoder}(i) \cdot \text{Encoder}(j)$

Maps the text into an h -dimensional vector

An entire field of study on how to get (or learn) the similarity function better
(We'll see some later!)

Inference: Index

Goal: find a small subset of elements in a datastore that are the most similar to the query

sim: a similarity score between two pieces of text

Can be a totally separate research area on how to do this fast & accurate

Index: given q , return $\text{argTop-}k_{d \in \mathcal{D}} \text{sim}(q, d)$ through fast nearest neighbor search

k elements from a datastore

<https://github.com/facebookresearch/faiss/wiki/>

Categorization of retrieval-augmented LMs

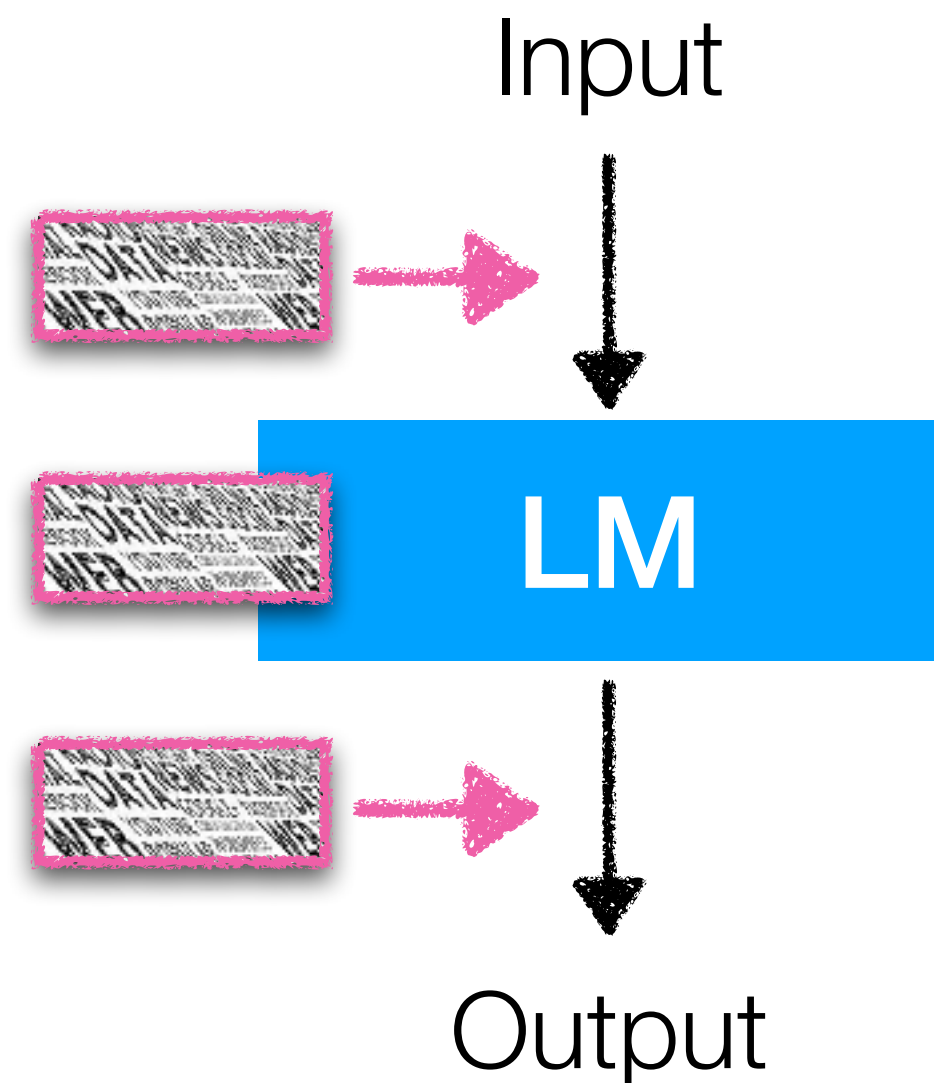
What to retrieve?

Query



Text chunks (passages)?
Tokens?
Something else?

How to use retrieval?



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

The capital city of Ontario is Toronto.

Categorization of retrieval-augmented LMs

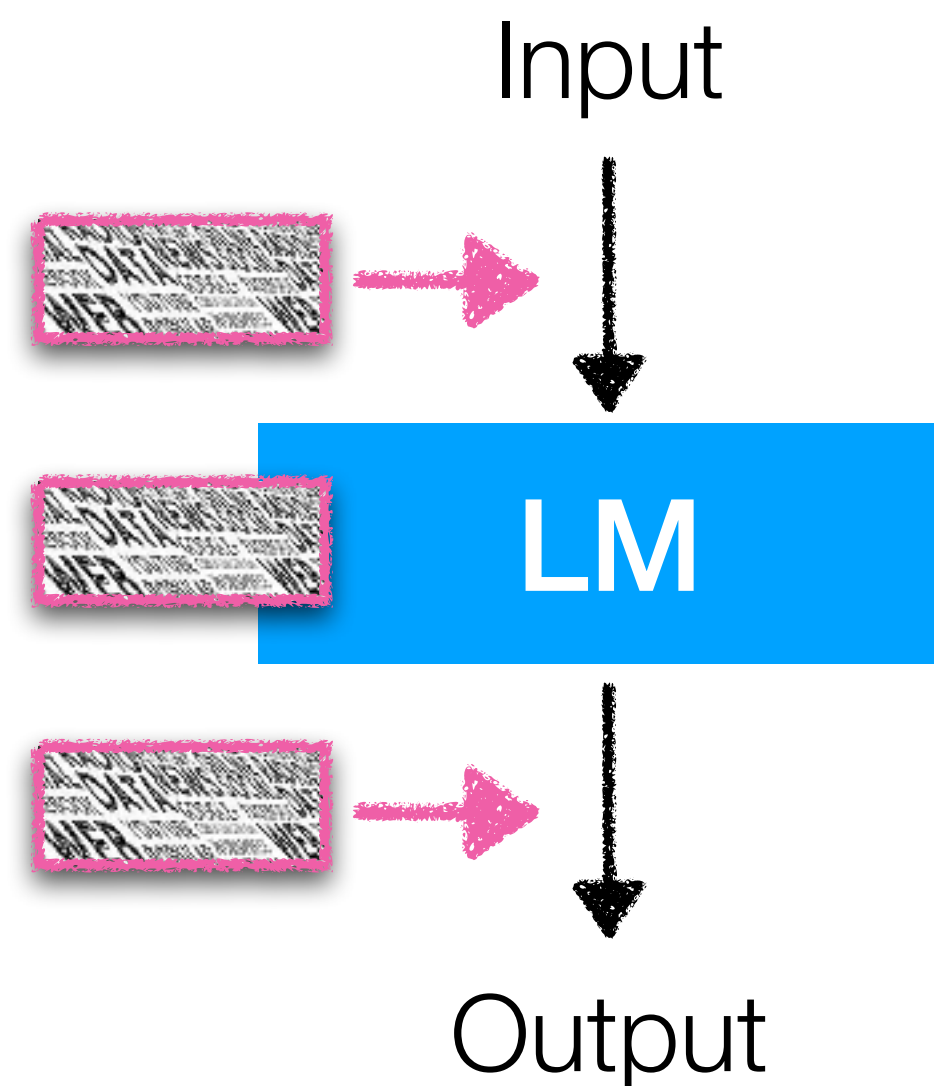
What to retrieve?

Query



Text chunks (passages)?
Tokens?
Something else?

How to use retrieval?



When to retrieve?

w/ retrieval

Today we focus on

1. What to retrieve

2. How to use retrieval

w/ retrieval w/r 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.

Three representative architectures

What: Text chunks
How: Input

Input augmentation (RAG)

What: Text chunks
How: Intermediate

Intermediate fusion

What: Tokens
How: Output

Output interpolations

More details?

- Section 3 of our tutorial (<https://acl2023-retrieval-lm.github.io/>)
- Our position paper (Asai et al., 2024; https://akariasai.github.io/assets/pdf/ralm_position.pdf)

Three representative architectures

What: Text chunks
How: Input

REALM (Guu et al., 2020)

What: Text chunks
How: Intermediate

RETRO (Borgeaud et al., 2021)

What: Tokens
How: Output

kNN-LM (Khandelwal et al., 2020)

More details?

- Section 3 of our tutorial (<https://acl2023-retrieval-lm.github.io/>)
- Our position paper (Asai et al., 2024; https://akariasai.github.io/assets/pdf/ralm_position.pdf)

Three representative architectures

What: Text chunks
How: Input

REALM (Guu et al., 2020)

What: Text chunks
How: Intermediate

RETRO (Borgeaud et al., 2021)

What: Tokens
How: Output

kNN-LM (Khandelwal et al., 2020)

More details?

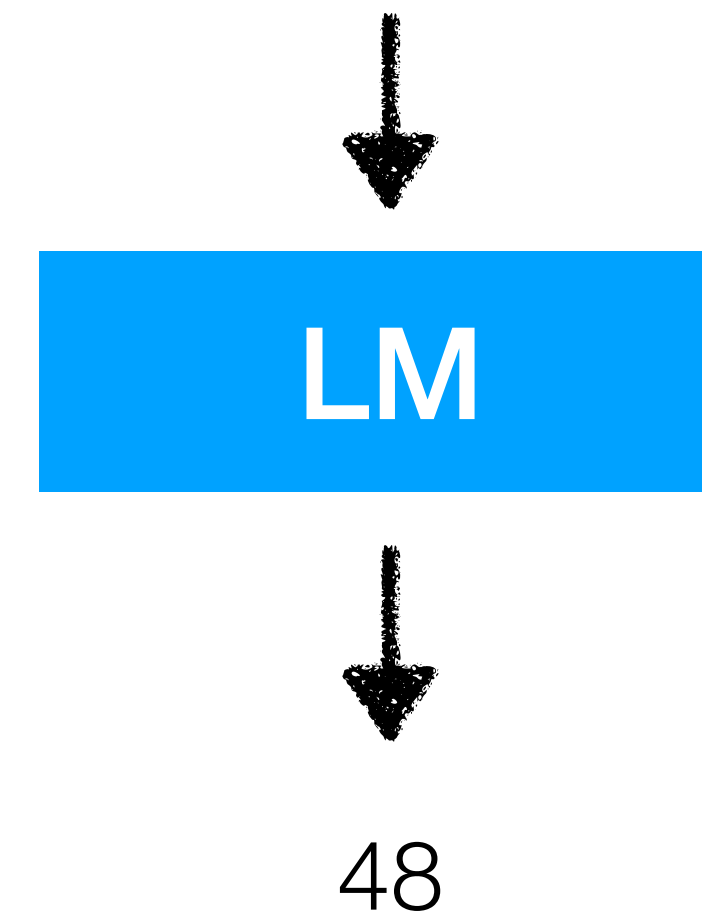
- Section 3 of our tutorial (<https://acl2023-retrieval-lm.github.io/>)
- Our position paper (Asai et al., 2024; https://akariasai.github.io/assets/pdf/ralm_position.pdf)

REALM (Guu et al 2020)



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

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



REALM (Guu et al 2020)

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

$q (=x)$



Retrieval



FIFA World Cup 2026
will expand to 48 teams.

k chunks of text
(passages)

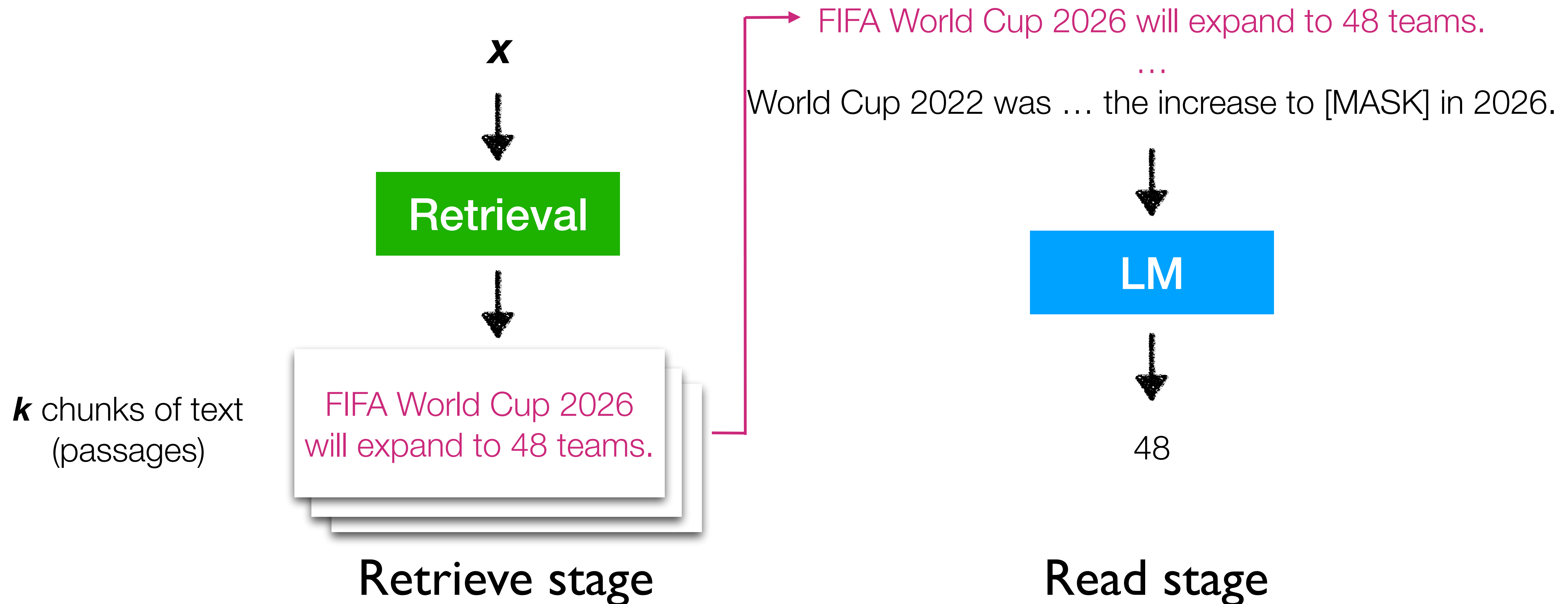
World Cup 2022 was ... the increase to [MASK] in 2026.



LM

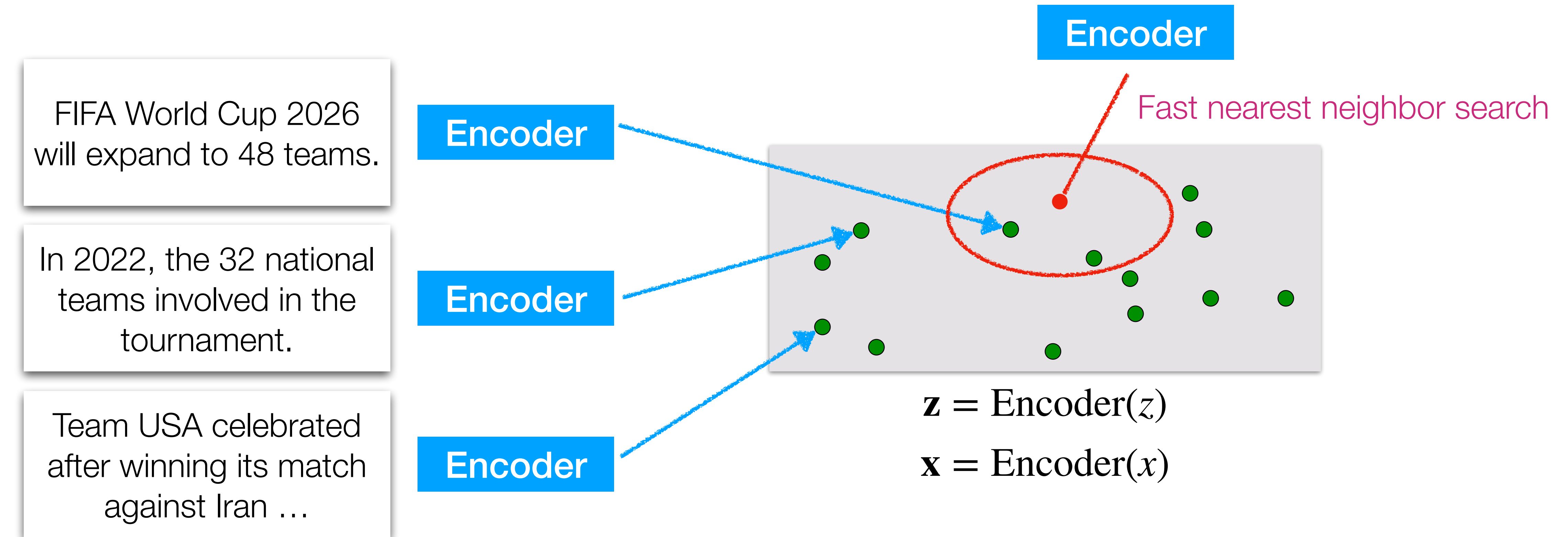
REALM (Guu et al 2020)

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



REALM: (I) Retrieve stage

\mathbf{x} = World Cup 2022 was ... the increase to [MASK] in 2026.

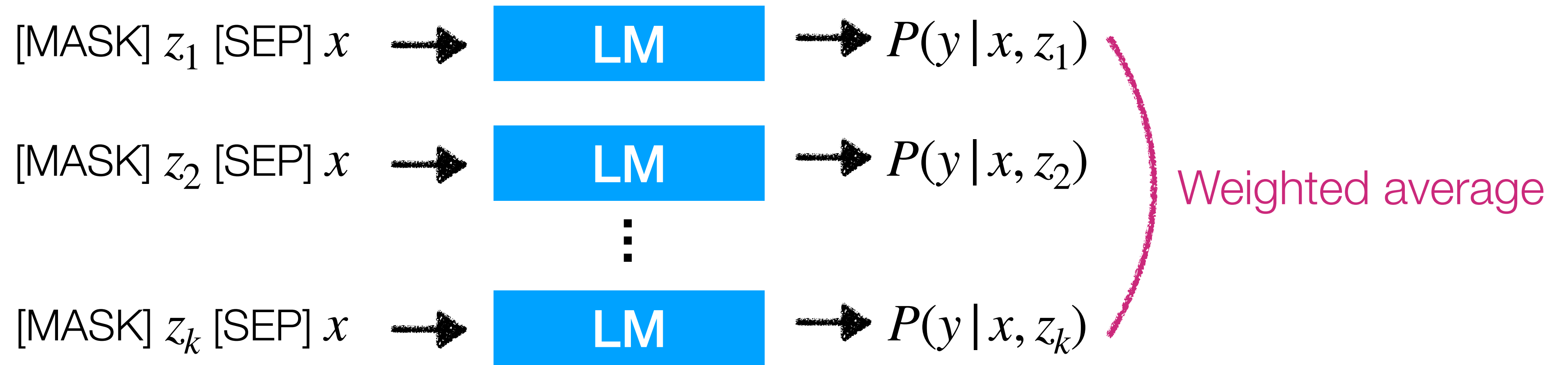


Wikipedia
13M chunks (passages)
(called *documents* in the paper)

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

k retrieved chunks

REALM: (2) Read stage



Need to approximate \rightarrow Consider top k chunks only

$$\sum_{z \in \mathcal{D}} \underbrace{P(z | x)}_{\text{from the retrieve stage}} \underbrace{P(y | x, z)}_{\text{from the read stage}}$$

Recent trend: RAG with LLMs

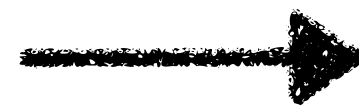
Existing parametric LMs
(e.g., GPT-3)



LM



Index



Off-the-shelf retrievers (e.g.,
Google search, BM25, DPR)

Simply combining existing models w/o
training has shown to be successful!

Three representative architectures

What: Text chunks
How: Input

REALM (Guu et al., 2020)

What: Text chunks
How: Intermediate

RETRO (Borgeaud et al., 2021)

What: Tokens
How: Output

kNN-LM (Khandelwal et al., 2020)

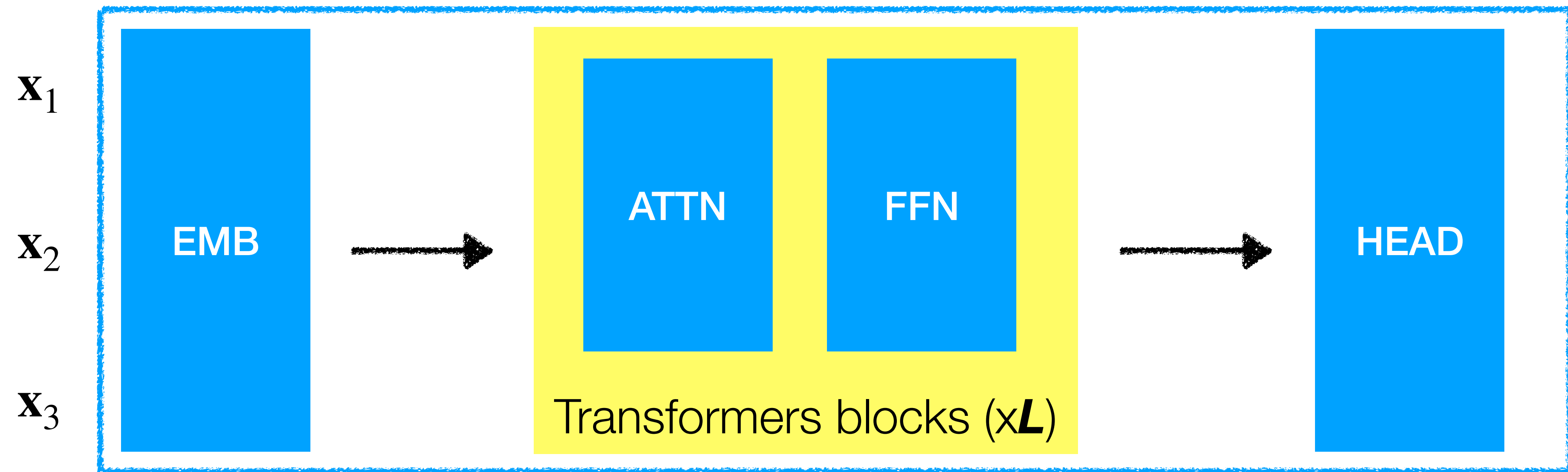
More details?

- Section 3 of our tutorial (<https://acl2023-retrieval-lm.github.io/>)
- Our position paper (Asai et al., 2024; https://akariasai.github.io/assets/pdf/ralm_position.pdf)

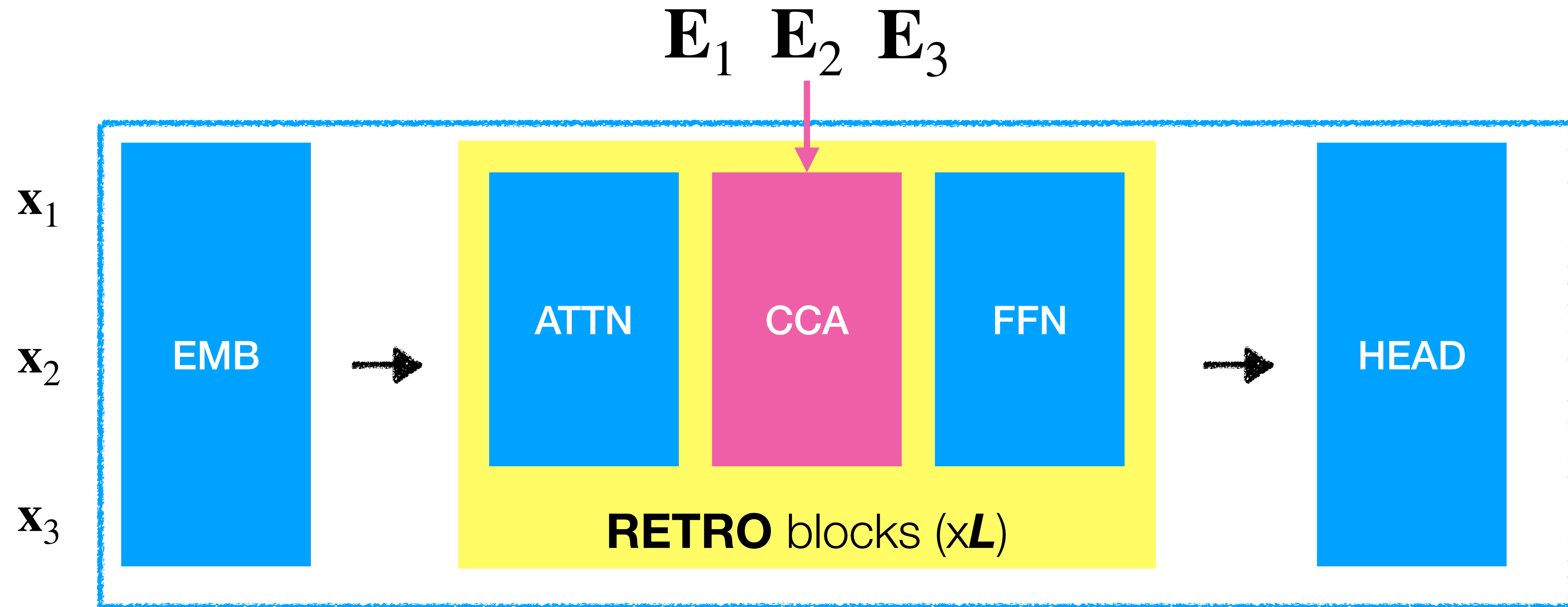
RETRO (Borgeaud et al. 2022)

- ✓ Incorporation in the “intermediate layer” instead of the “input” layer
→ designed for *many* chunks, *frequently*, more *efficiently*
- ✓ Scale the datastore (1.8T tokens)

Regular decoder

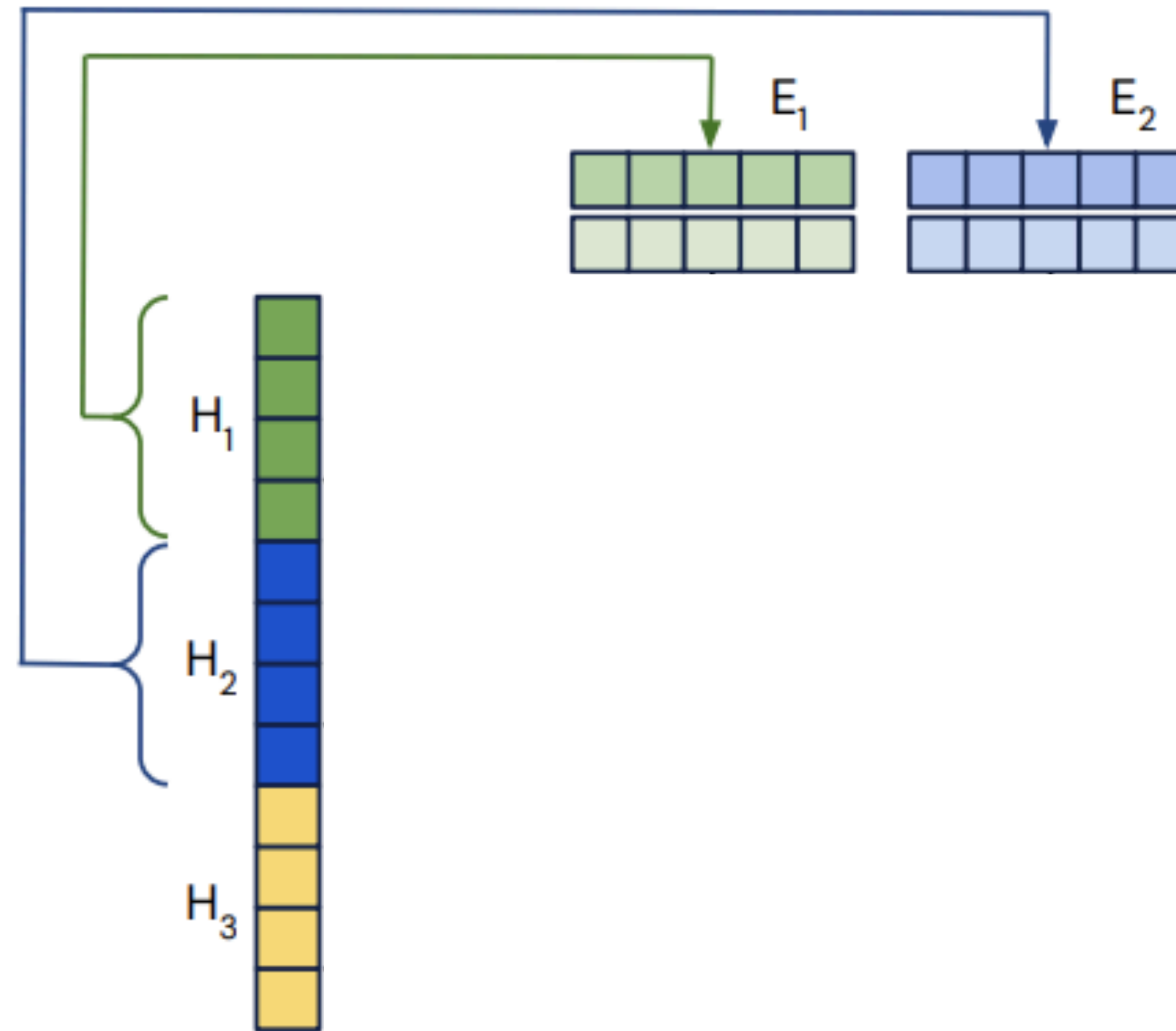


Decoder in RETRO



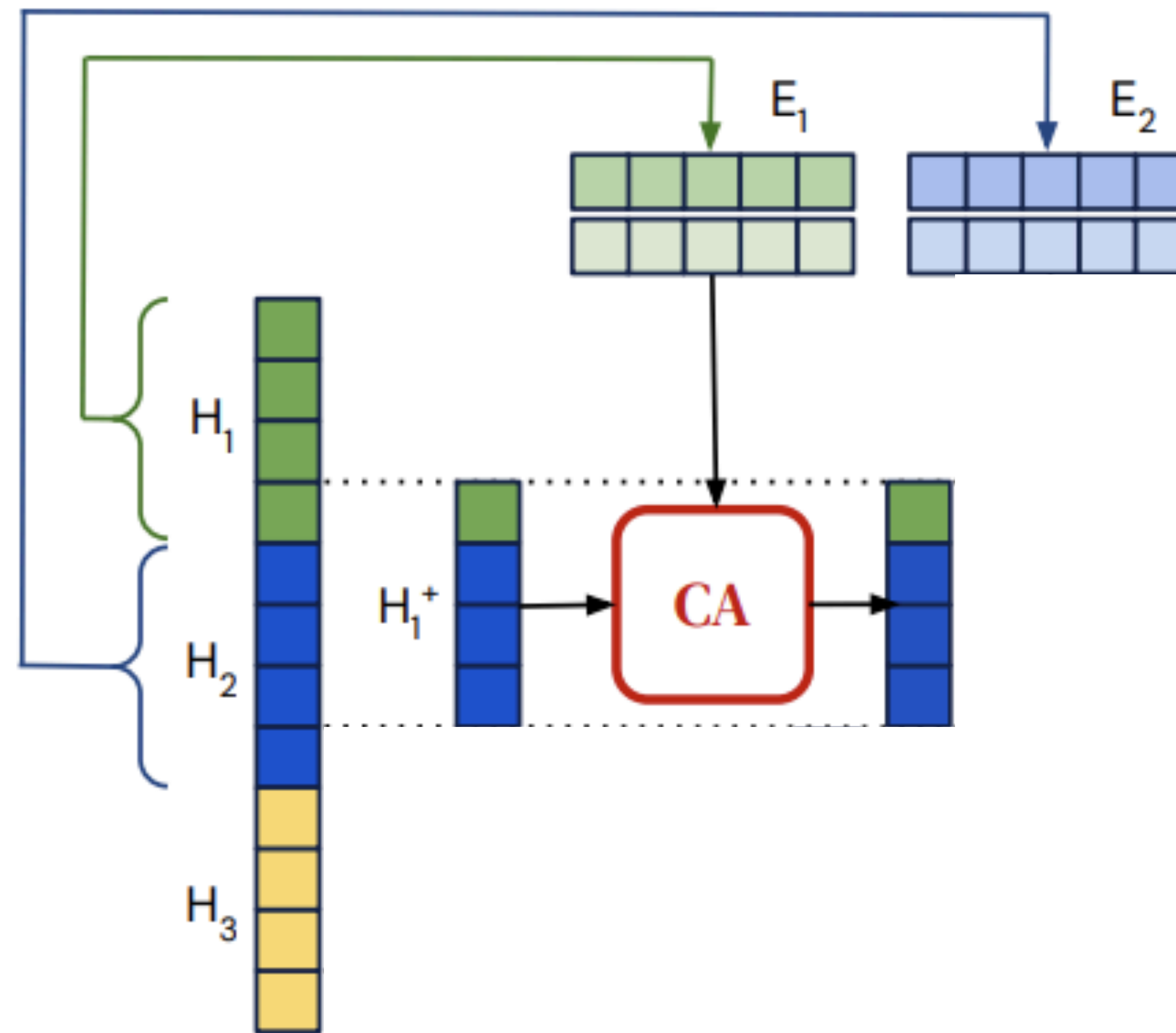
Chunked Cross Attention (CCA)

Chunked Cross Attention



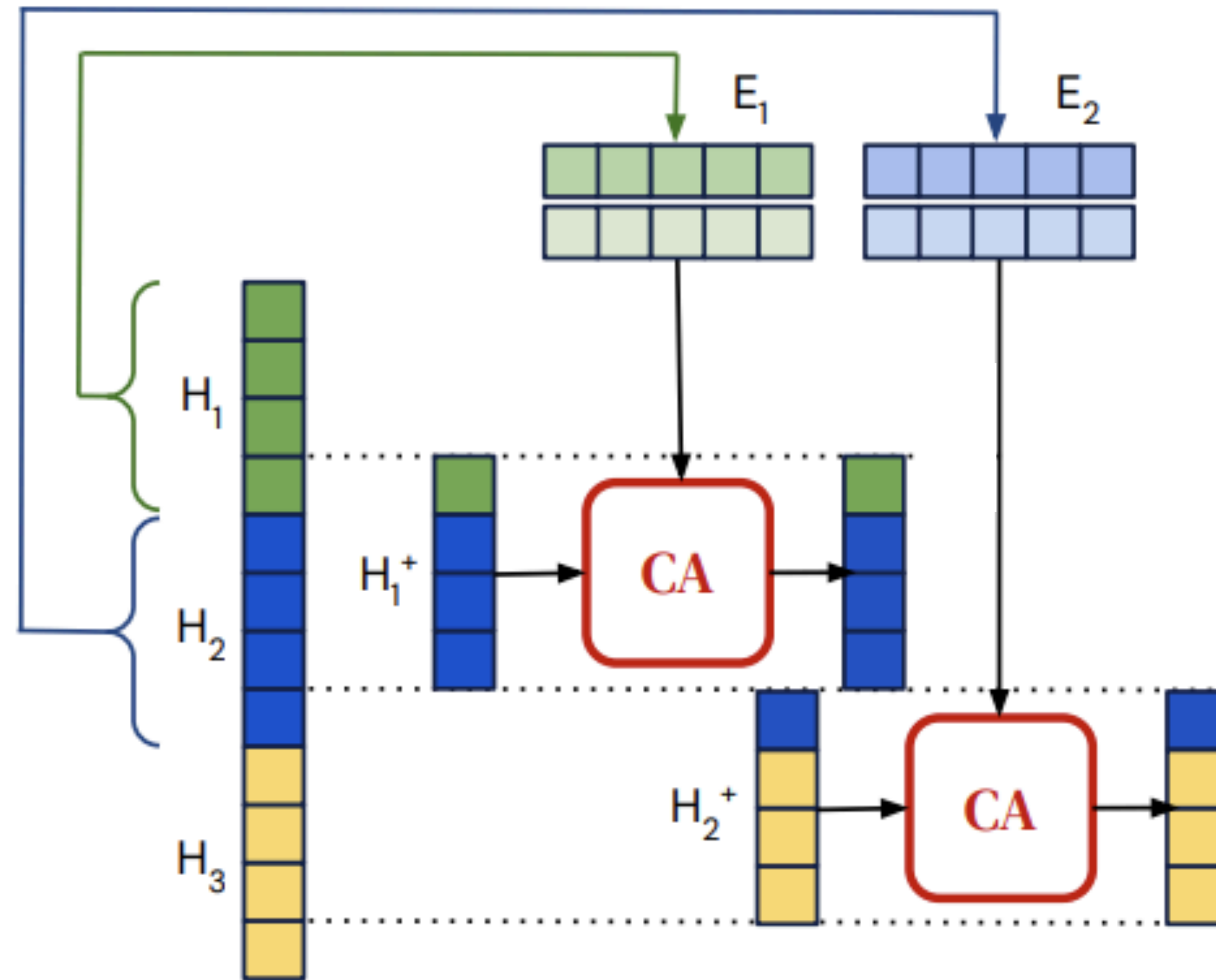
Outputs from the previous layer H

Chunked Cross Attention



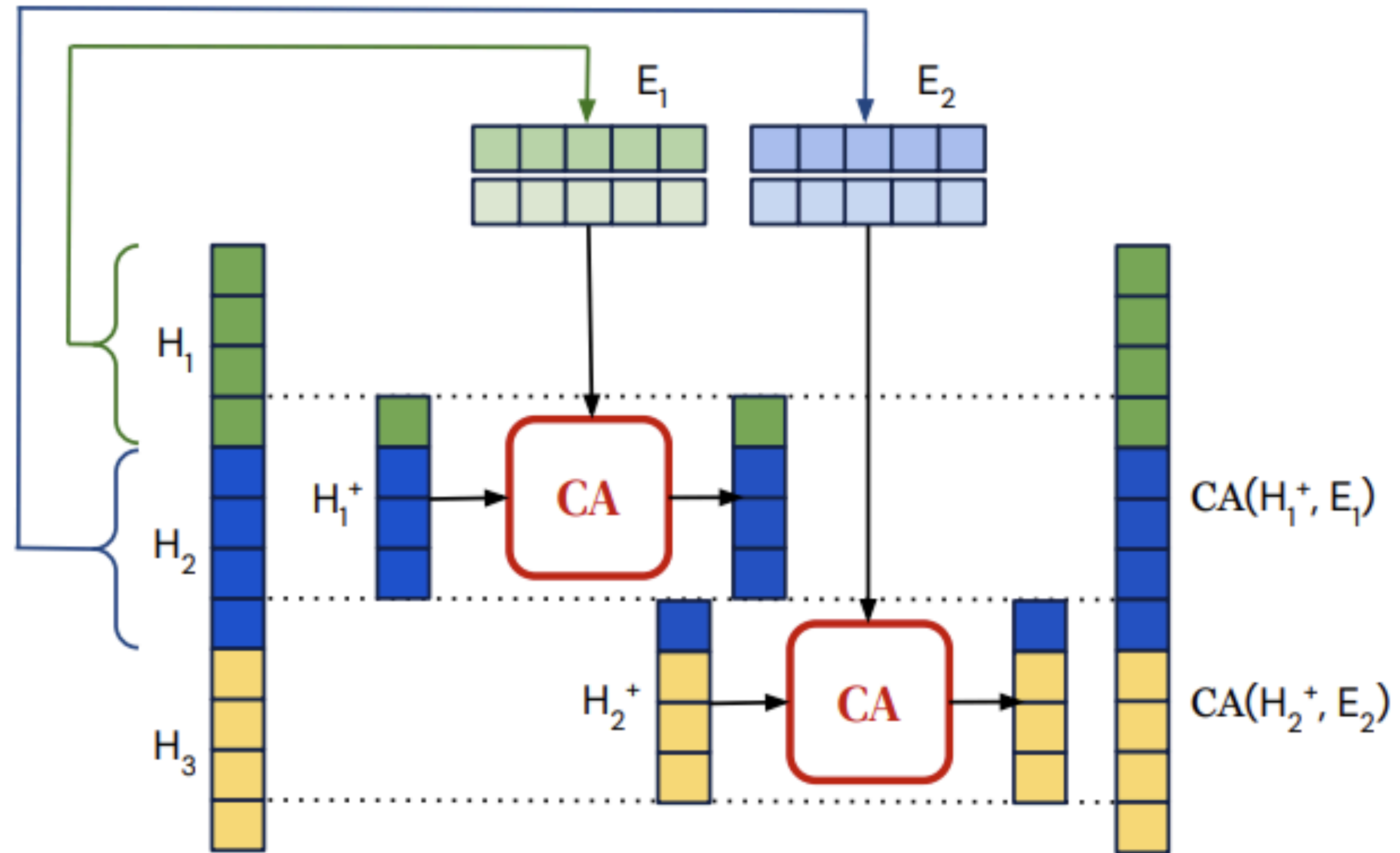
Outputs from the previous layer H

Chunked Cross Attention



Outputs from the previous layer H

Chunked Cross Attention

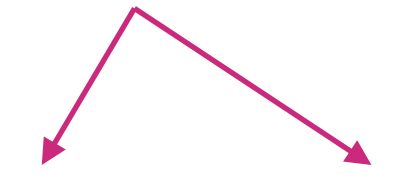


Outputs from the previous layer H

Inputs to the next layer

Results

Perplexity: The lower the better



Model	Retrieval Set	#Database tokens	#Database keys	Valid	Test
Baseline transformer (ours)	-	-	-	21.53	22.96
<i>k</i> NN-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

RETRO (w/ Wikipedia) outperforms its parametric counterpart

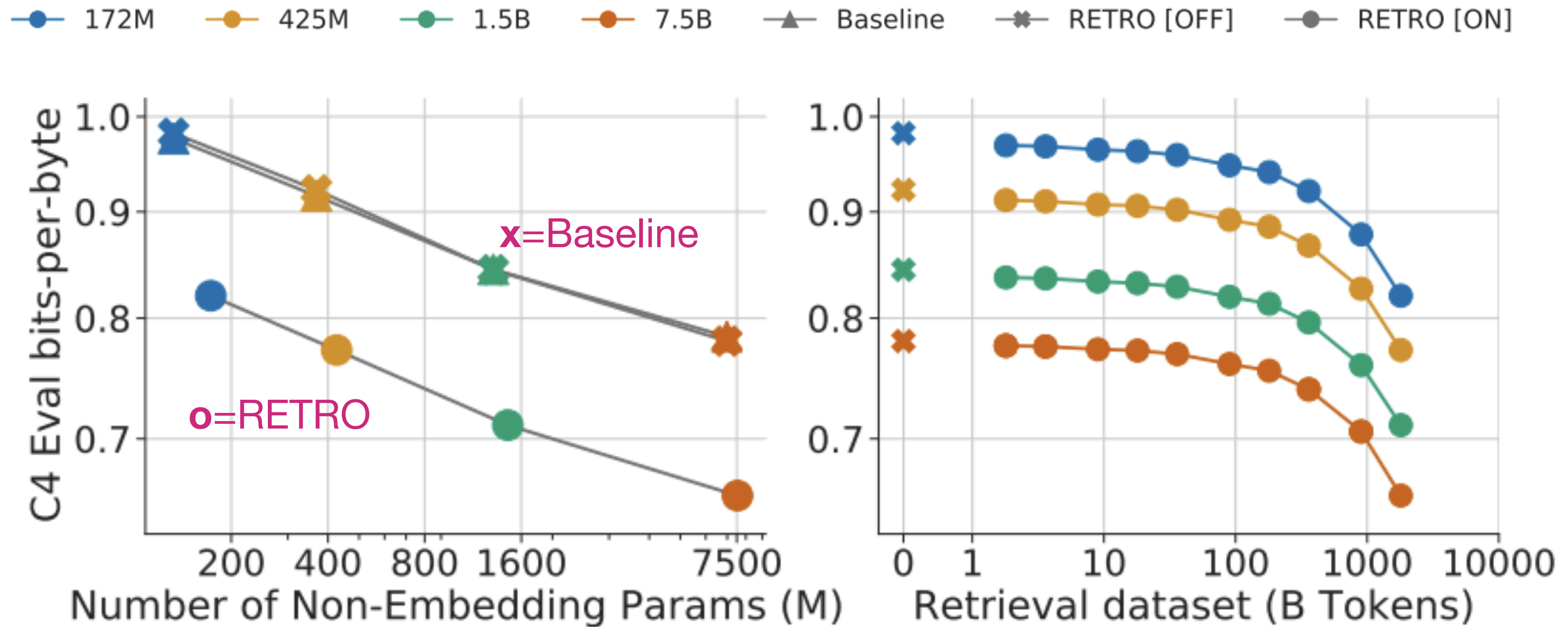
Results

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

RETRO w/ 1.8T datastores achieves SOTA

Results



Gains are constant with model scale

The larger datastore is, the better

Three representative architectures

What: Text chunks
How: Input

REALM (Guu et al., 2020)

What: Text chunks
How: Intermediate

RETRO (Borgeaud et al., 2021)

What: Tokens
How: Output

kNN-LM (Khandelwal et al., 2020)

More details?

- Section 3 of our tutorial (<https://acl2023-retrieval-lm.github.io/>)
- Our position paper (Asai et al., 2024; https://akariasai.github.io/assets/pdf/ralm_position.pdf)

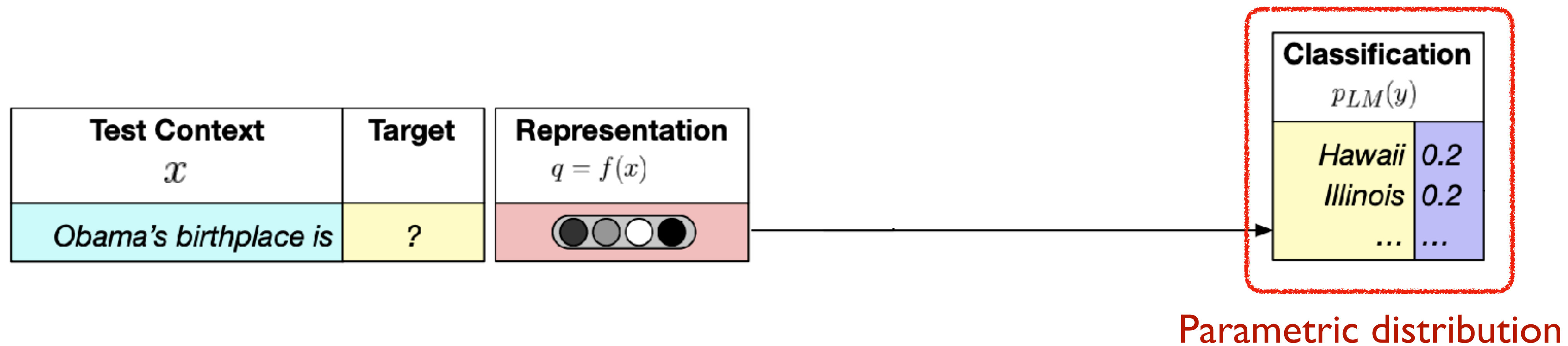
kNN-LM (Khandelwal et al. 2020)

- ✓ A different way of using retrieval, where the LM outputs a nonparametric distribution over every token in the data.
- ✓ Can be seen as an incorporation in the “output” layer

kNN-LM (Khandelwal et al. 2020)

Test Context x	Target
Obama's birthplace is	?


kNN-LM (Khandelwal et al. 2020)



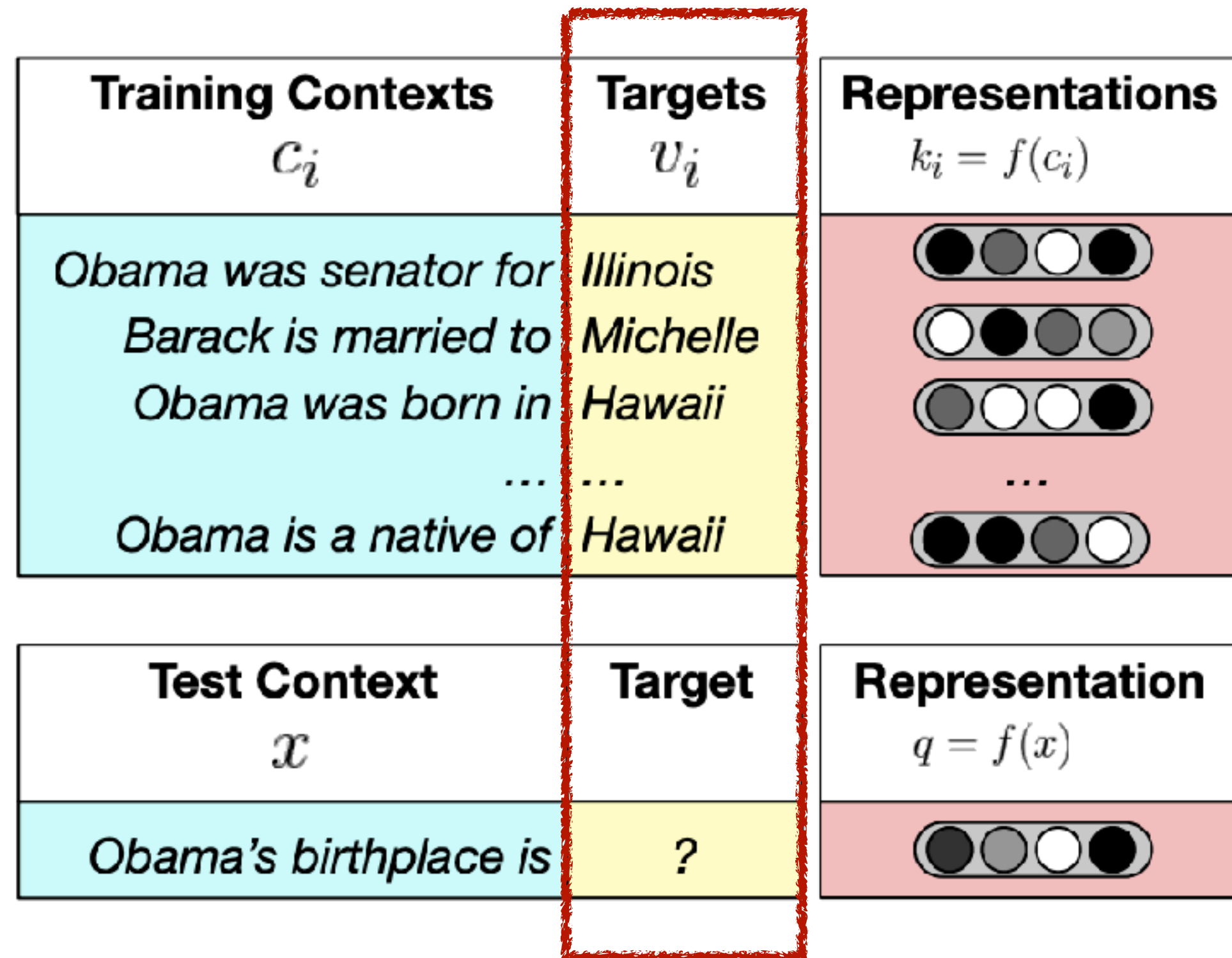
kNN-LM (Khandelwal et al. 2020)

Training Contexts c_i	Targets v_i
Obama was senator for	Illinois
Barack is married to	Michelle
Obama was born in	Hawaii
...	...
Obama is a native of	Hawaii

... Obama was senator for Illinois from 1997 to 2005, Barack is Married to Michelle and their first daughter, ... Obama was born in Hawaii, and graduated from Columbia University. ... Obama is a native of Hawaii,

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

kNN-LM (Khandelwal et al. 2020)



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

kNN-LM (Khandelwal et al. 2020)

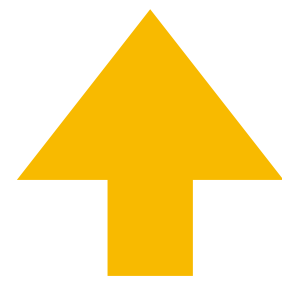
The size of the datastore = # of tokens in the corpus ($> 1B$)

Training Contexts c_i	Targets v_i	Representations $k_i = f(c_i)$
Obama was senator for	Illinois	
Barack is married to	Michelle	
Obama was born in	Hawaii	
...
Obama is a native of	Hawaii	
Test Context x	Target	Representation $q = f(x)$
Obama's birthplace is	?	

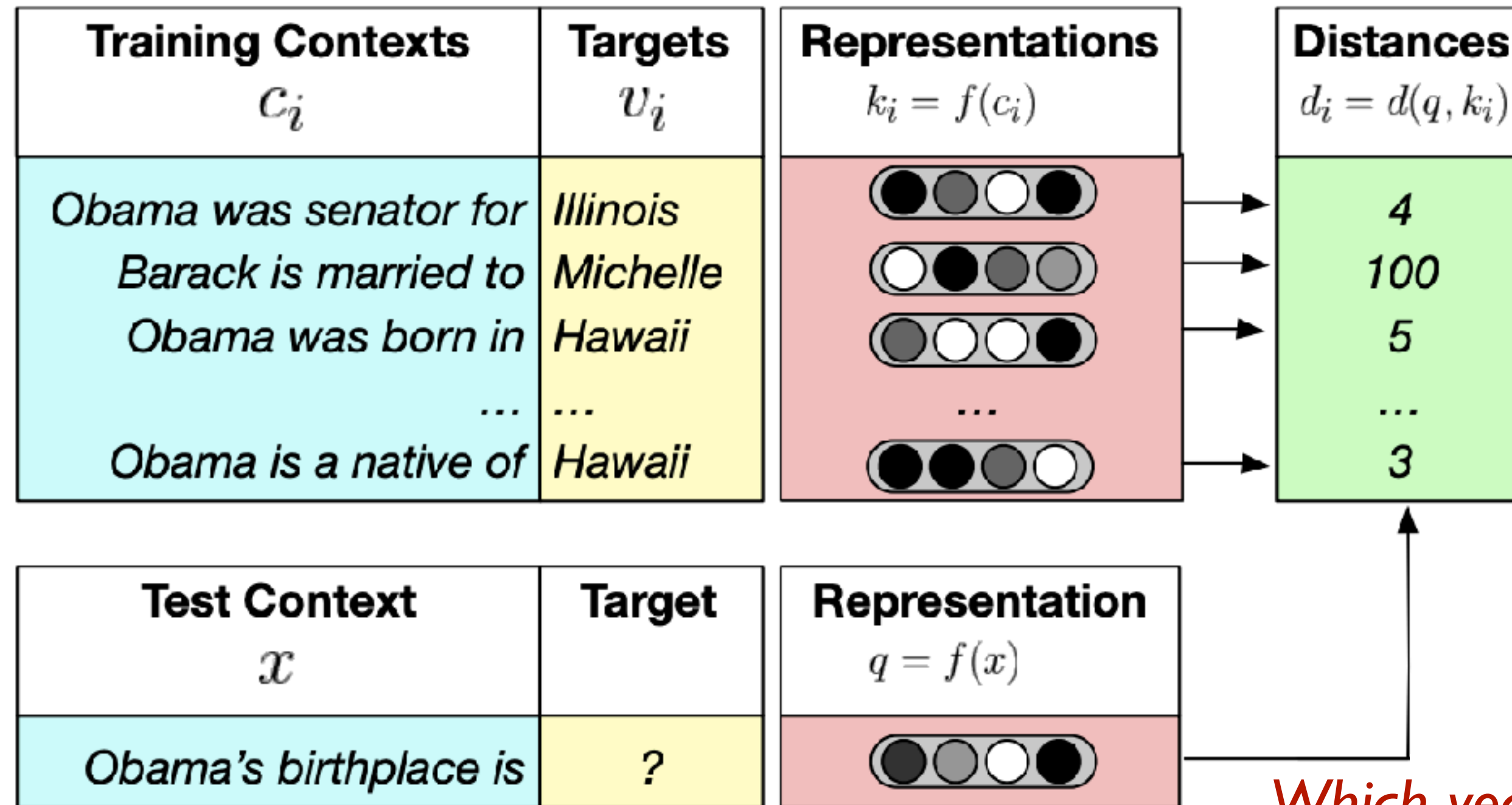
Which tokens in a datastore are close to the next token?

=

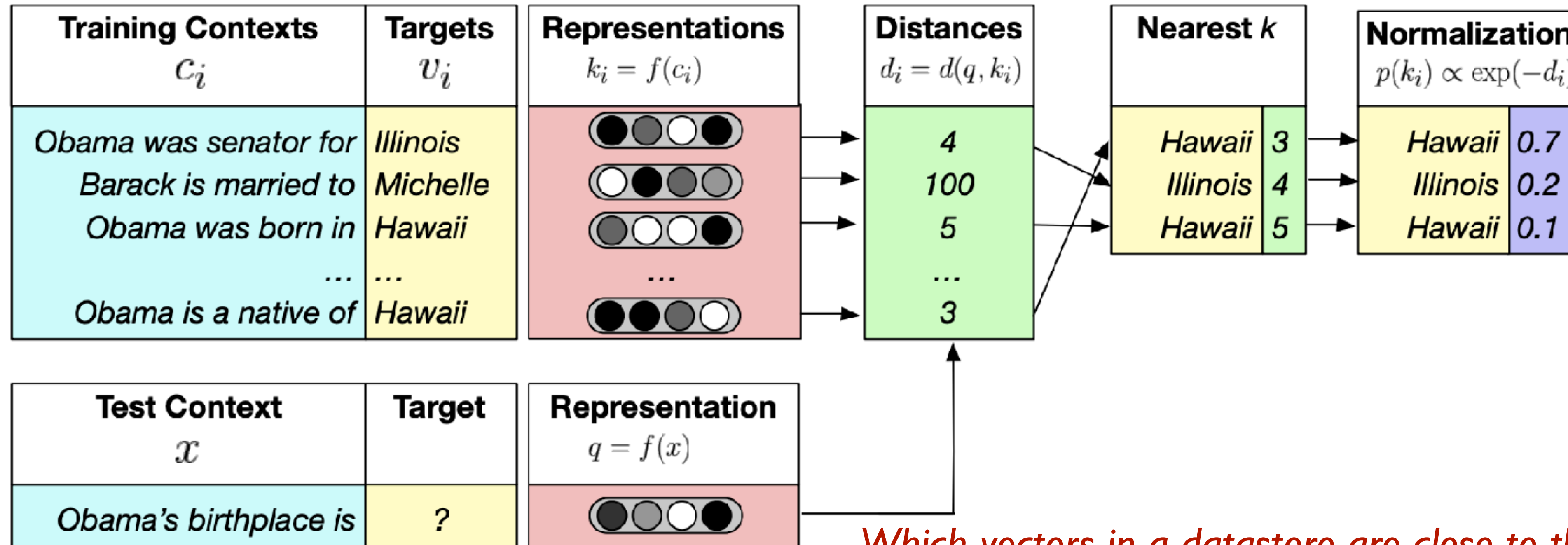
Which prefixes in a datastore are close to the prefix we have?



kNN-LM (Khandelwal et al. 2020)

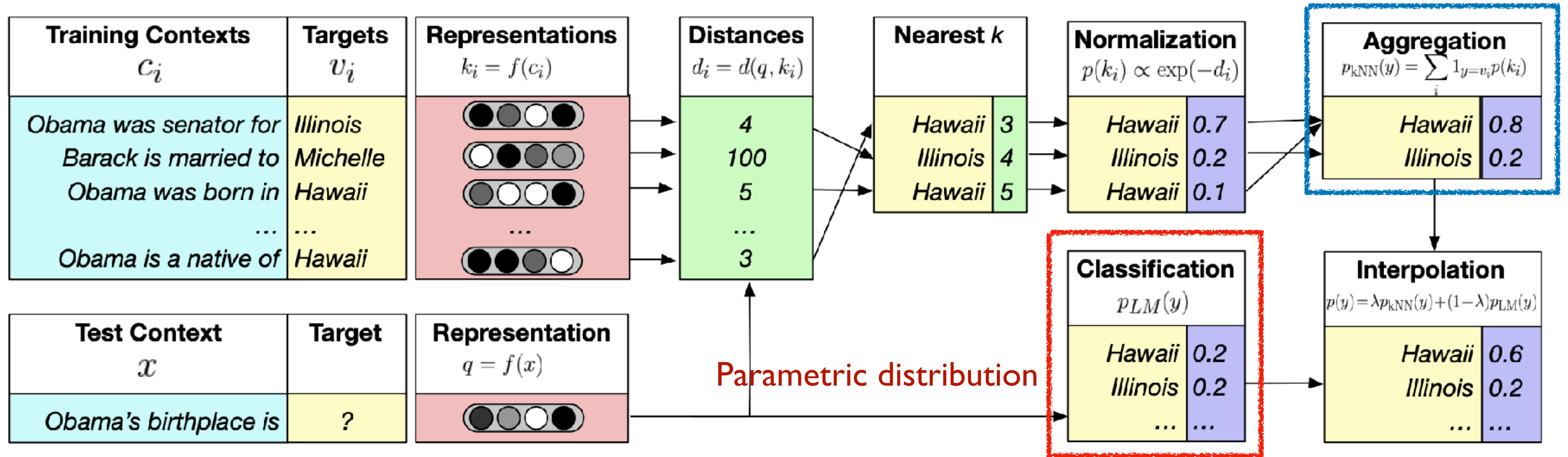


kNN-LM (Khandelwal et al. 2020)



kNN-LM (Khandelwal et al. 2020)

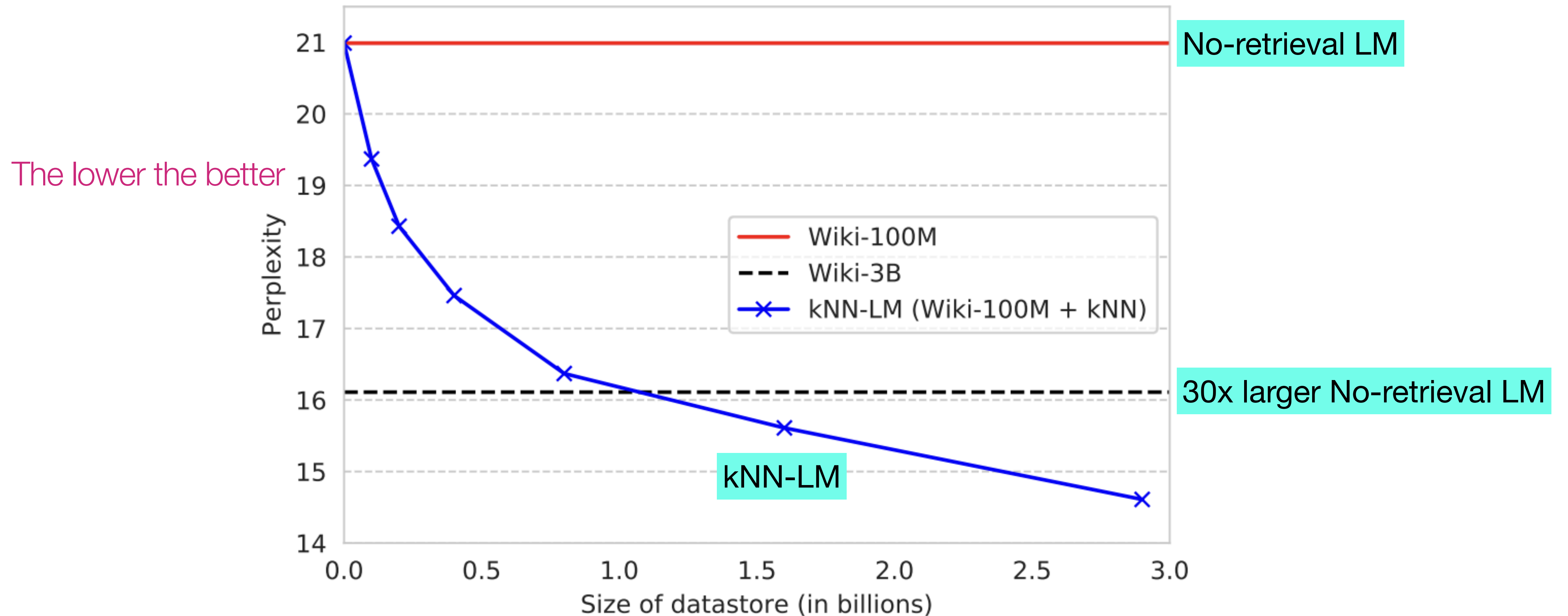
Nonparametric distribution



λ : hyperparameter

$$P_{kNN-LM}(y | x) = (1 - \lambda) P_{LM}(y | x) + \lambda P_{kNN}(y | x)$$

kNN-LM - results



Outperforms no-retrieval LM

Better with bigger datastore

Three representative architectures

What: Text chunks
How: Input

REALM (Guu et al., 2020)

What: Text chunks
How: Intermediate

RETRO (Borgeaud et al., 2021)

What: Tokens
How: Output

kNN-LM (Khandelwal et al., 2020)

More details?

- Section 3 of our tutorial (<https://acl2023-retrieval-lm.github.io/>)
- Our position paper (Asai et al., 2024; https://akariasai.github.io/assets/pdf/ralm_position.pdf)

Today's outline

Question:

[https://bit.ly/
akari_ralm_lec](https://bit.ly/akari_ralm_lec)



Scan me

Why do we need retrieval-augmented LMs?

Architectures of retrieval-augmented LMs (Inference)

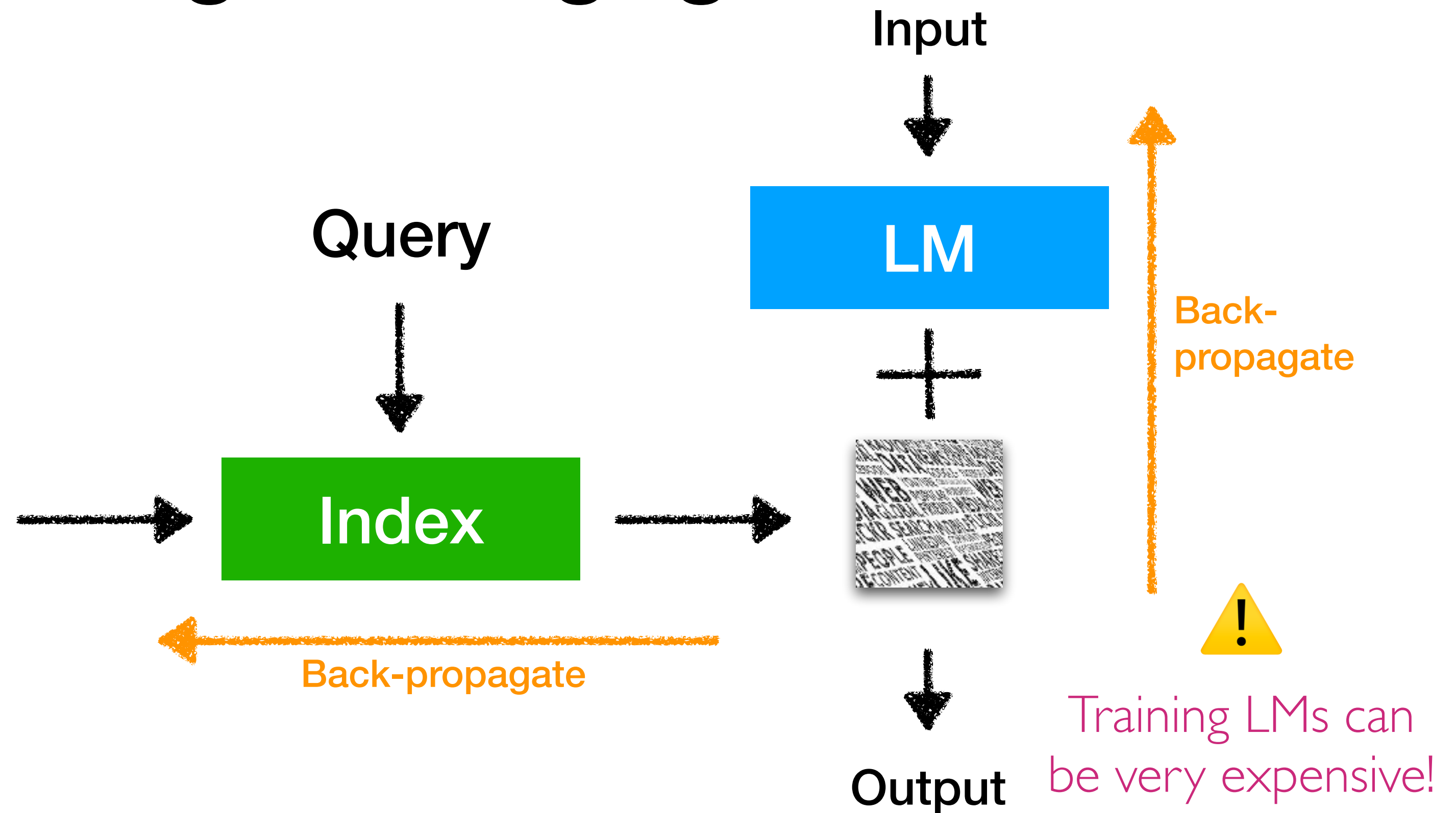
Training of retrieval-augmented LMs

Limitations and future directions

Why is training challenging?



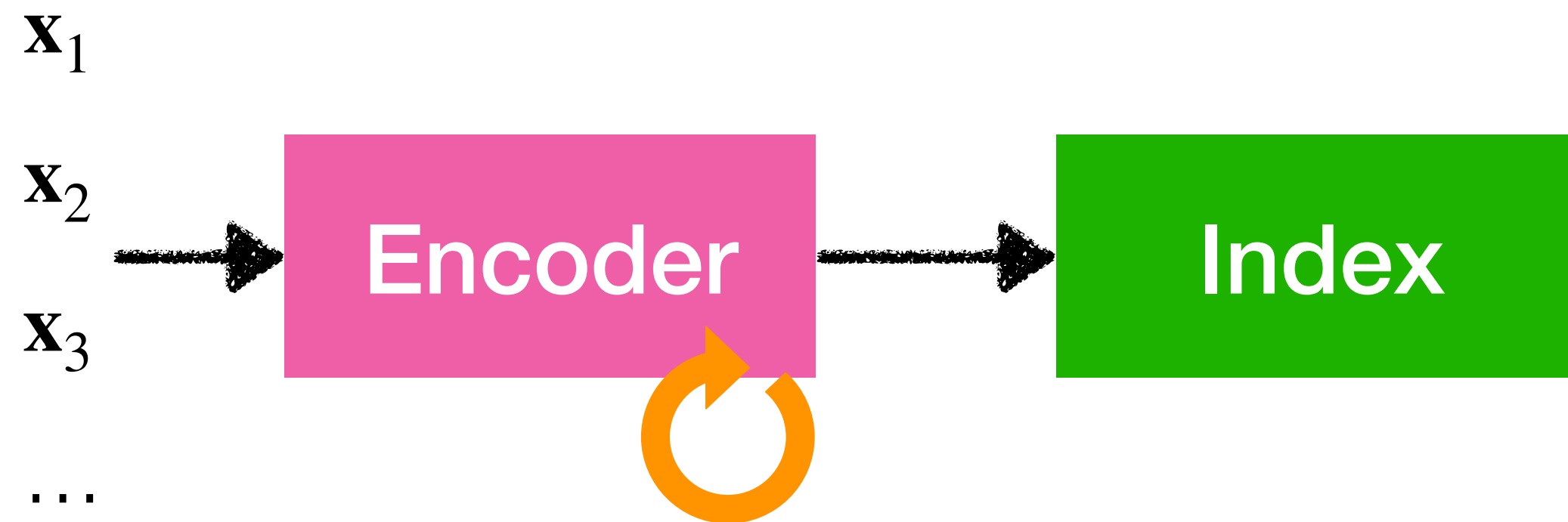
Datastore



Challenges of updating retrieval models



Datastore



During training, we will update the encoder

Training methods for retrieval-augmented LMs

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

Training methods for retrieval-augmented LMs

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

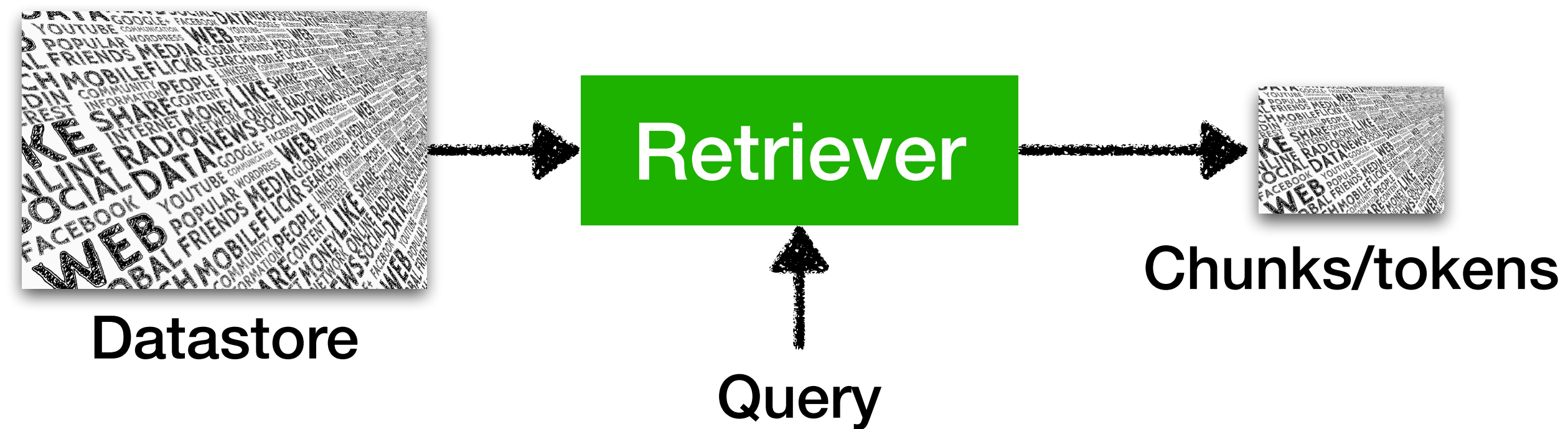
Independent training

Retrieval models and language models are trained **independently**

- Training language models



- Training retrieval models



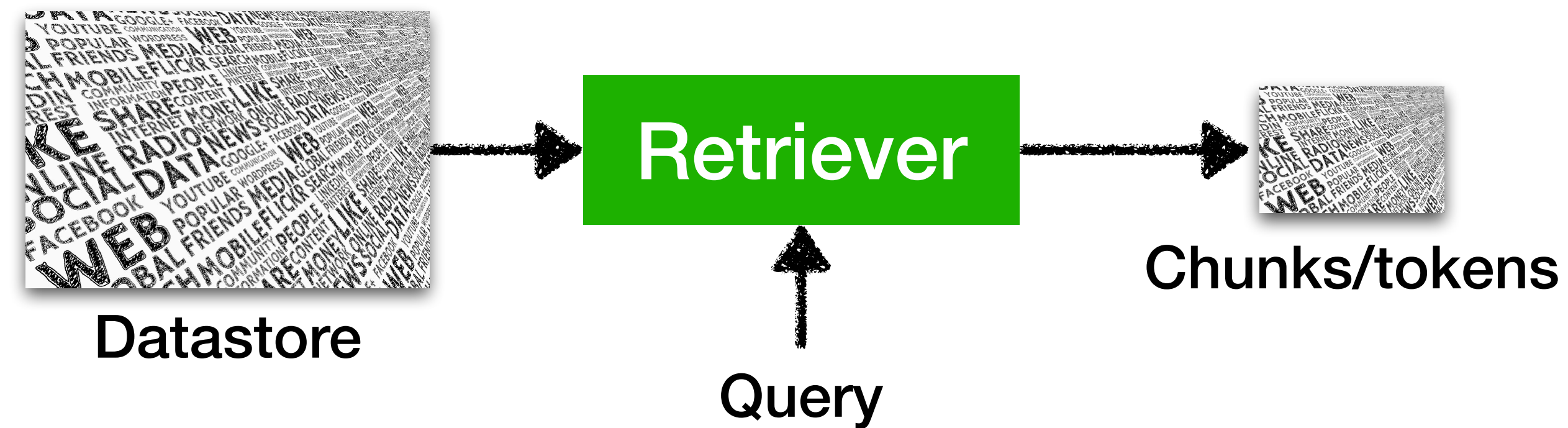
Independent training

Retrieval models and language models are trained **independently**

- Training language models



- Training retrieval models



Sparse retrieval models: TF-IDF / BM25

In 1997, **Apple** merged with NeXT,
and Steve **Jobs** became **CEO** of ...

Jobs returned to **Apple** as **CEO**
after the company's acquisition ...

[0, 0, 0.4, 0, 0.8, 0.7, ...]

[0, 1.2, 0.4, 0, 0.8, 0, ...]

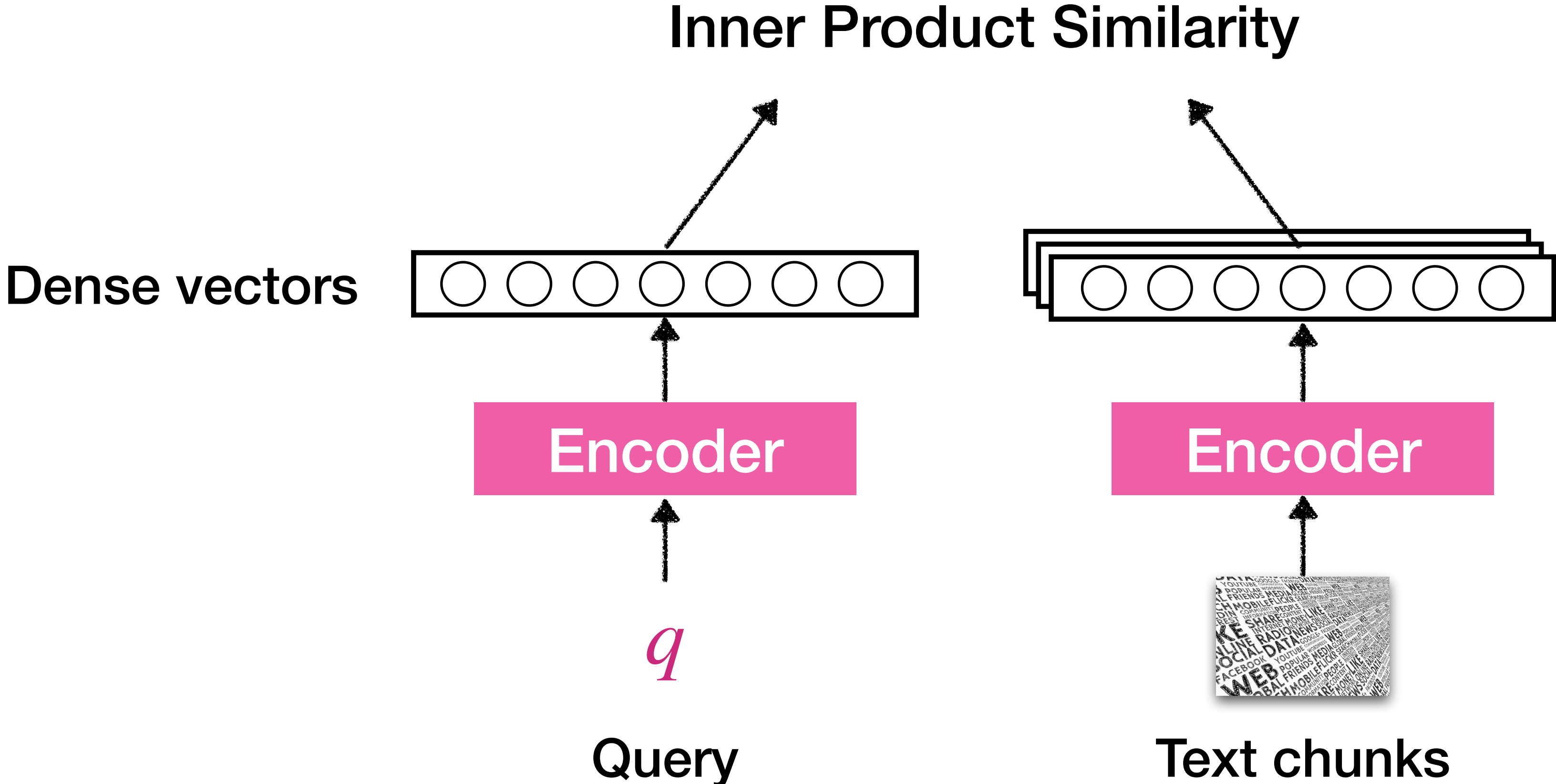
Lexical overlap

Text chunks

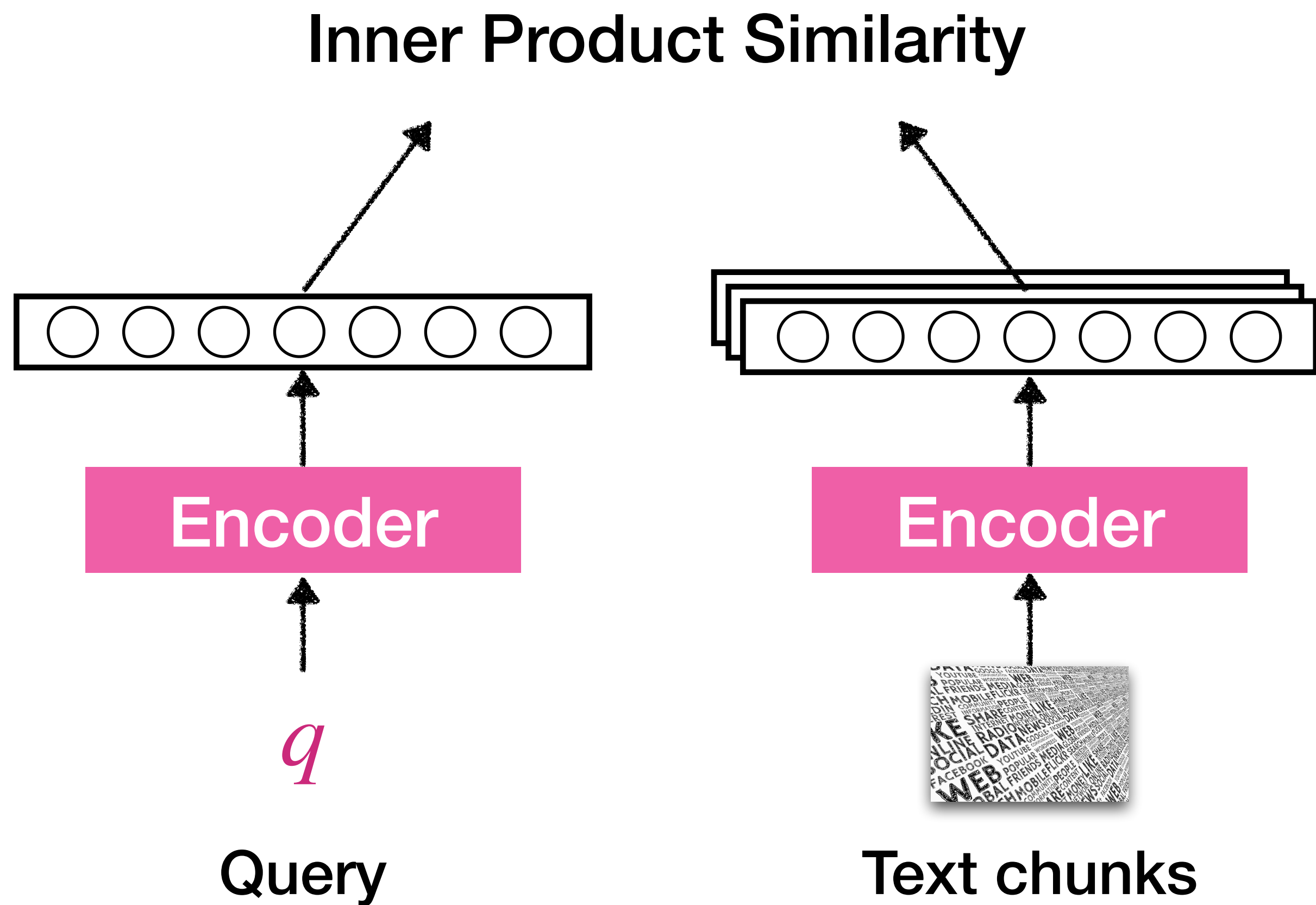
Sparse vectors

No training needed!

Dense retrieval models: DPR (Karpukhin et al. 2020)

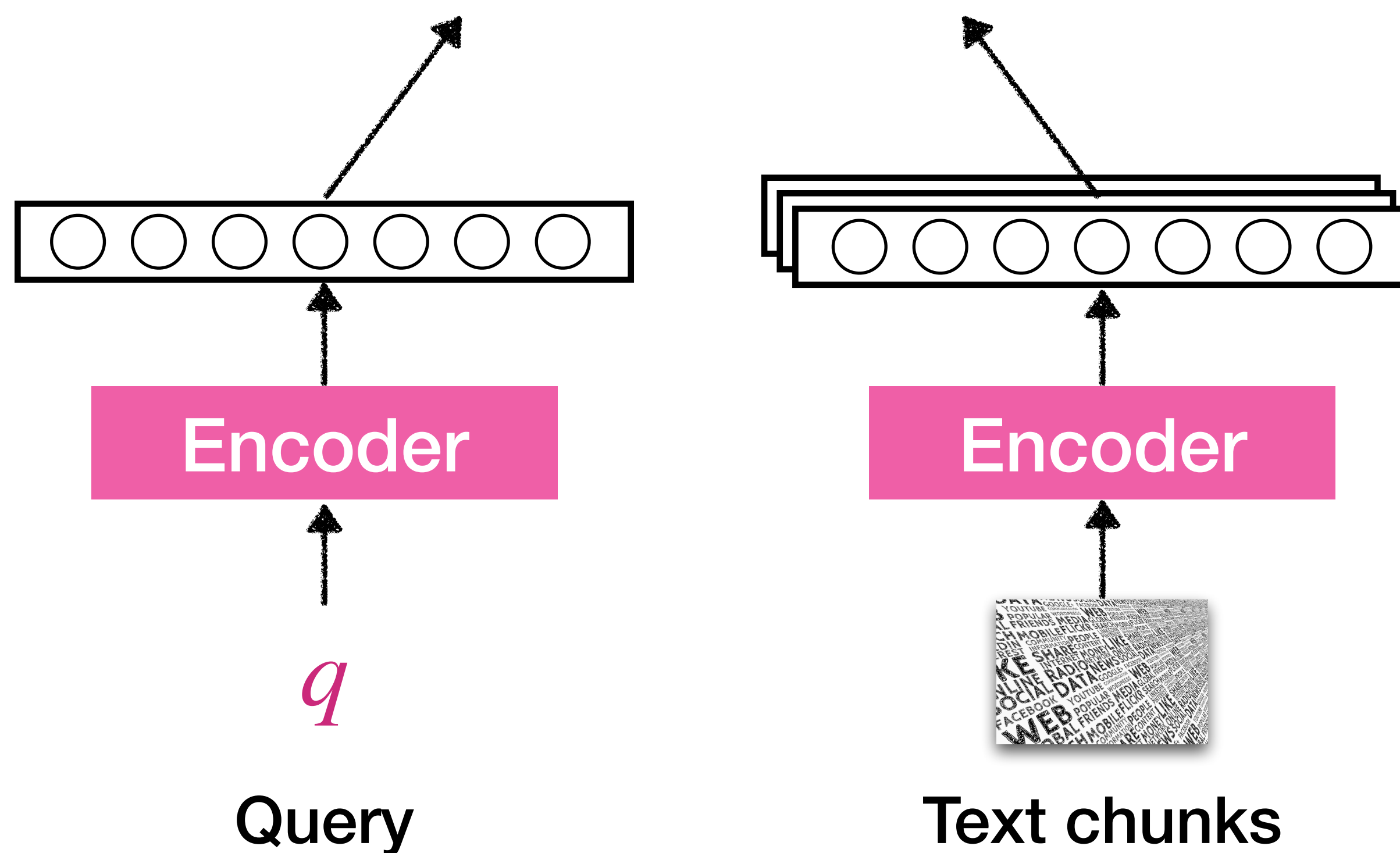


Training dense retrieval models: DPR



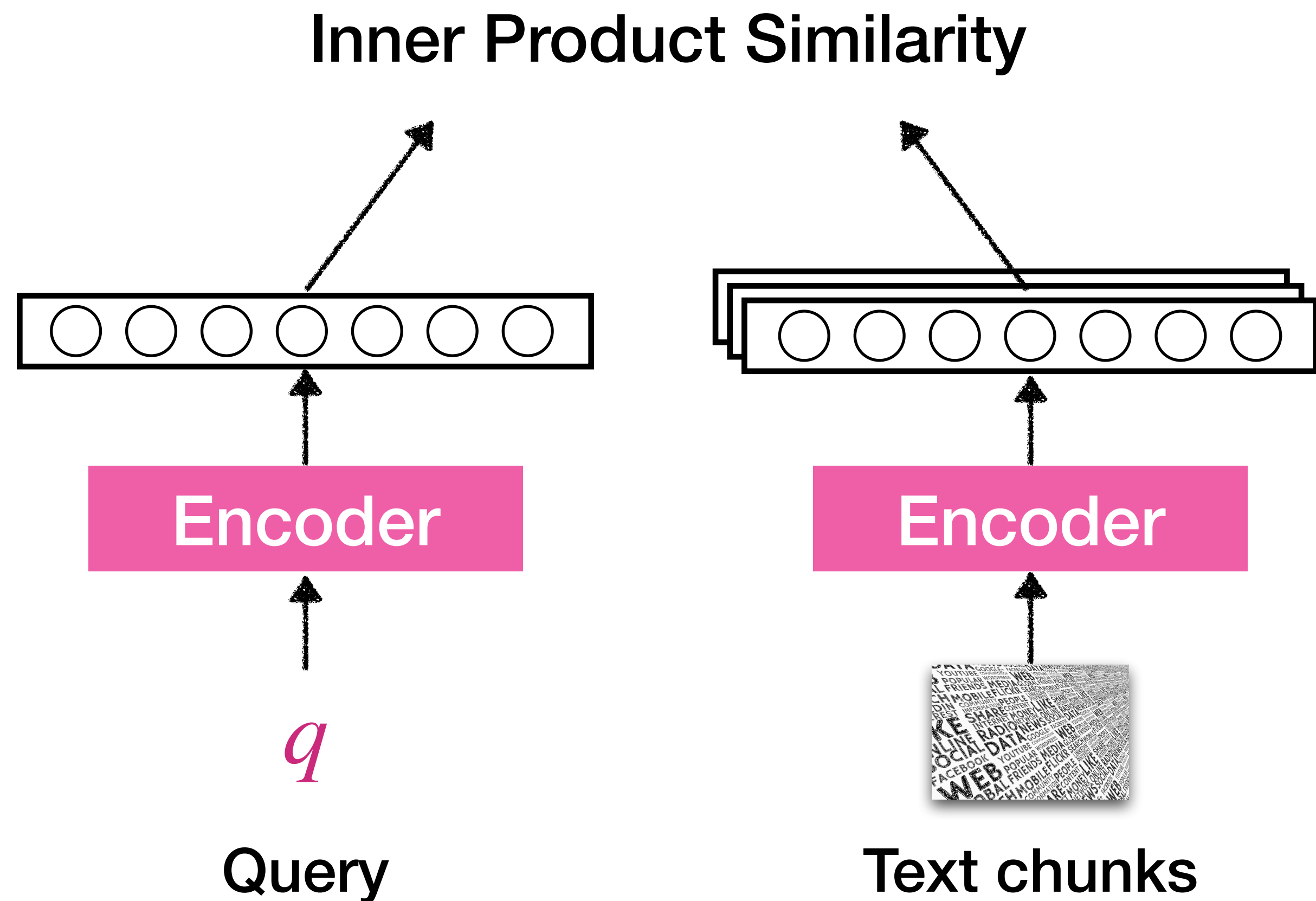
Training dense retrieval models: DPR

Inner Product Similarity



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

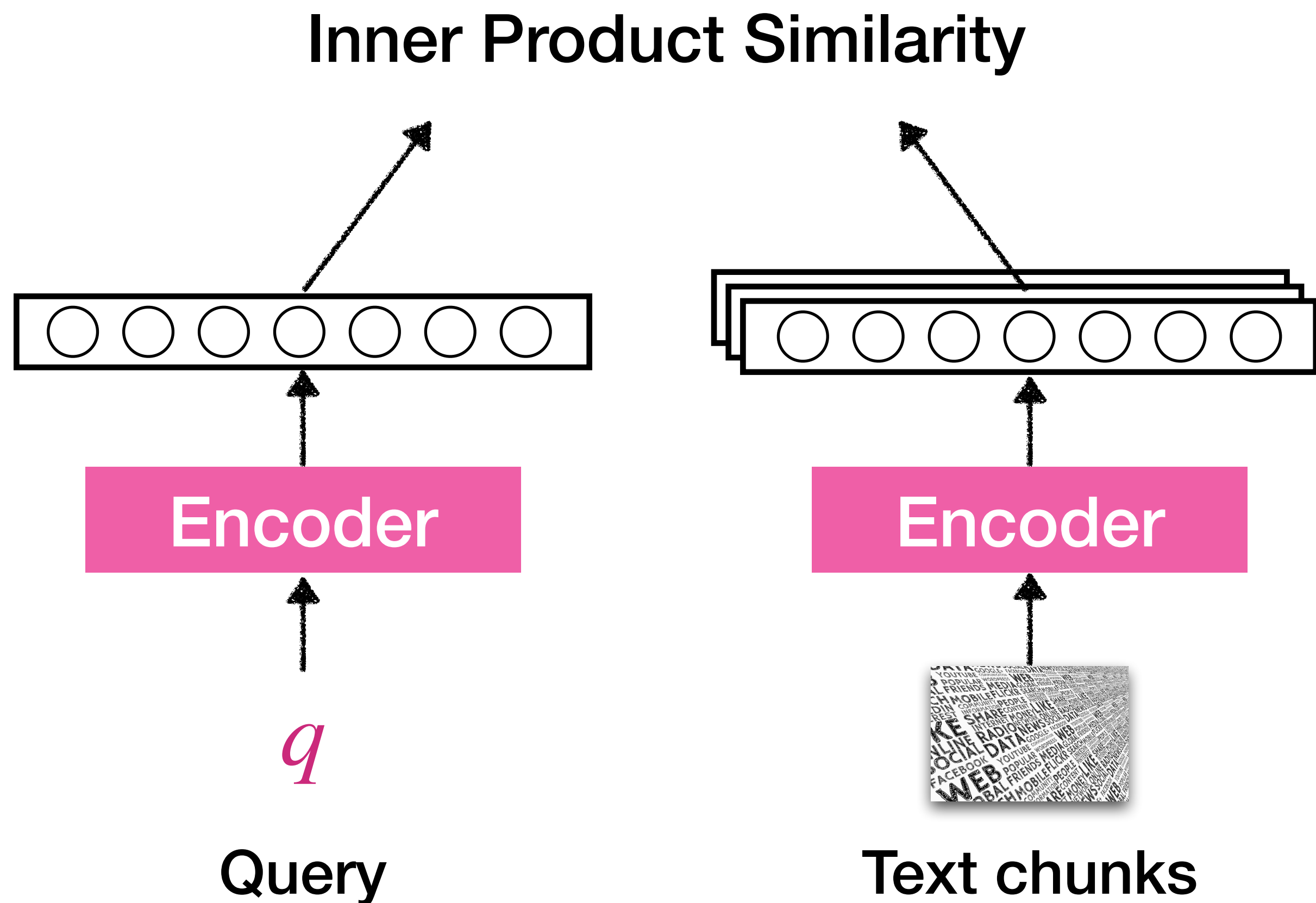
Training dense retrieval models: DPR



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

Contrastive learning

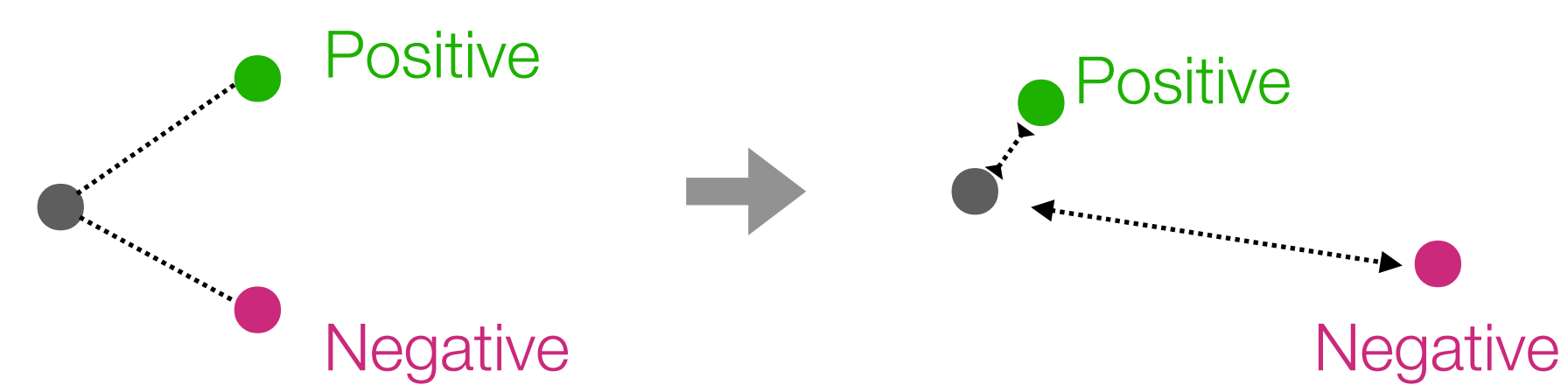
Training dense retrieval models: DPR



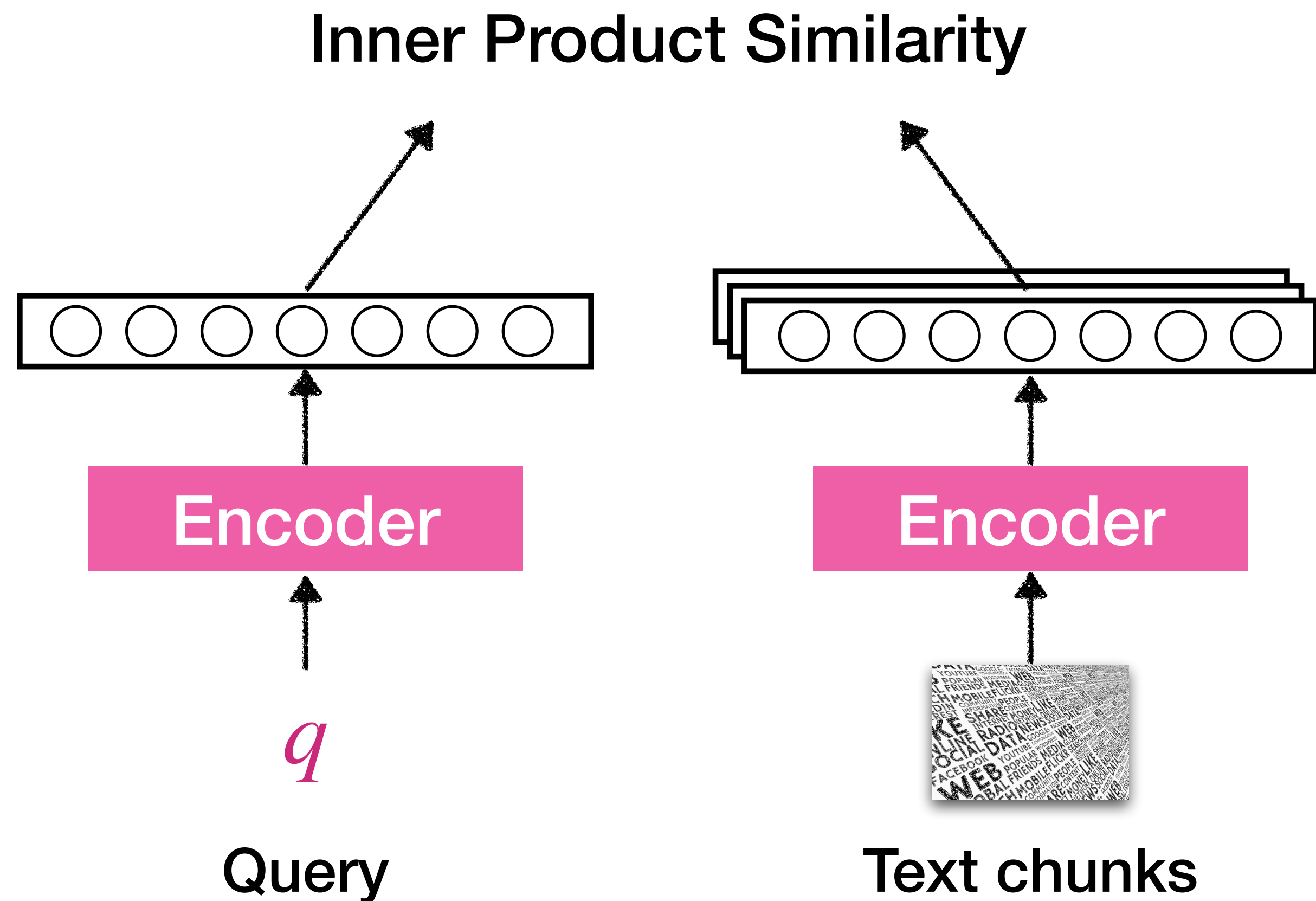
$$L(q, p^+, p_1^-, p_2^-, \dots, p_n^-)$$

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

Contrastive learning



Training dense retrieval models: DPR

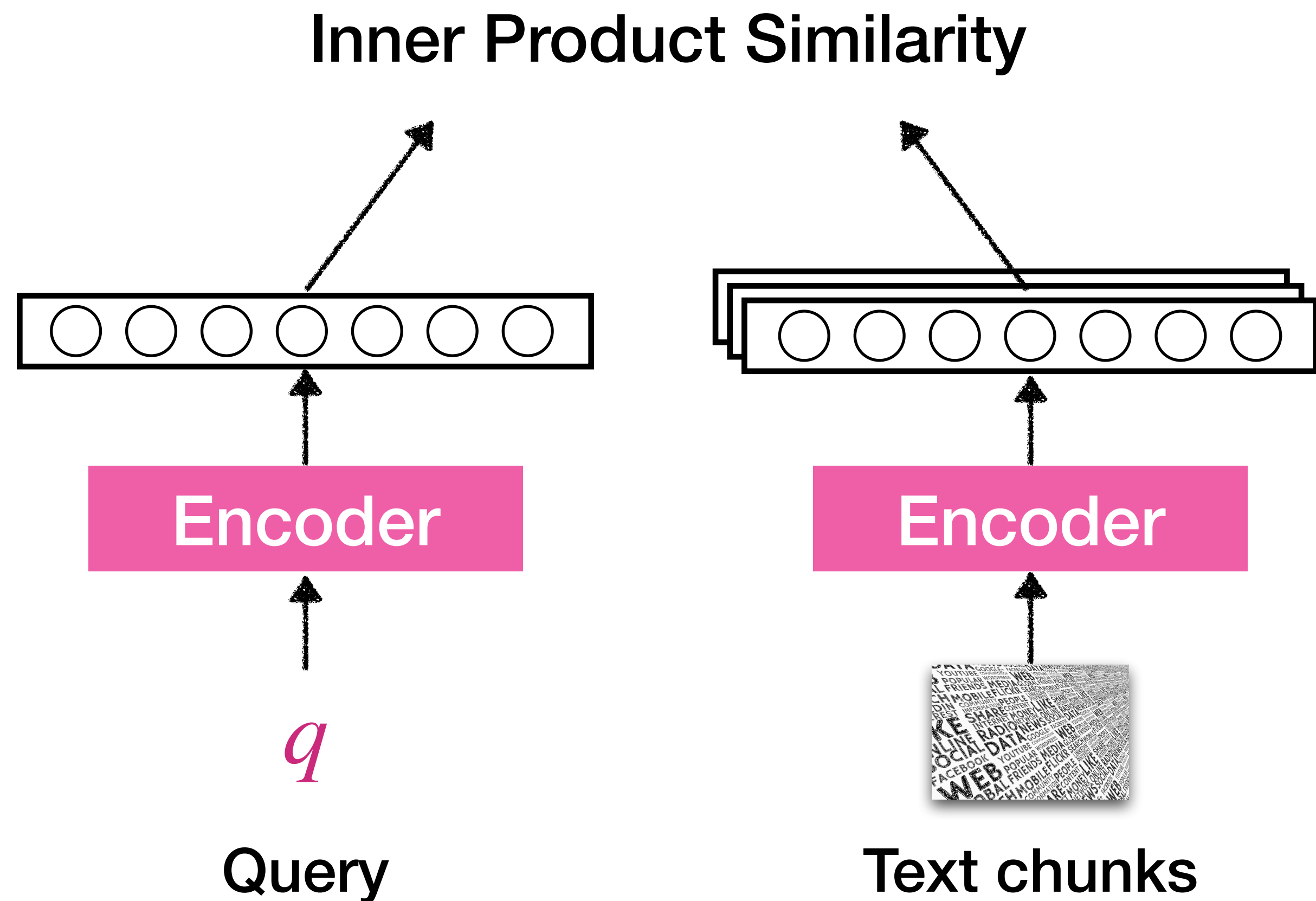


$$L(q, p^+, p_1^-, p_2^-, \dots, p_n^-)$$

Positive passage

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

Training dense retrieval models: DPR



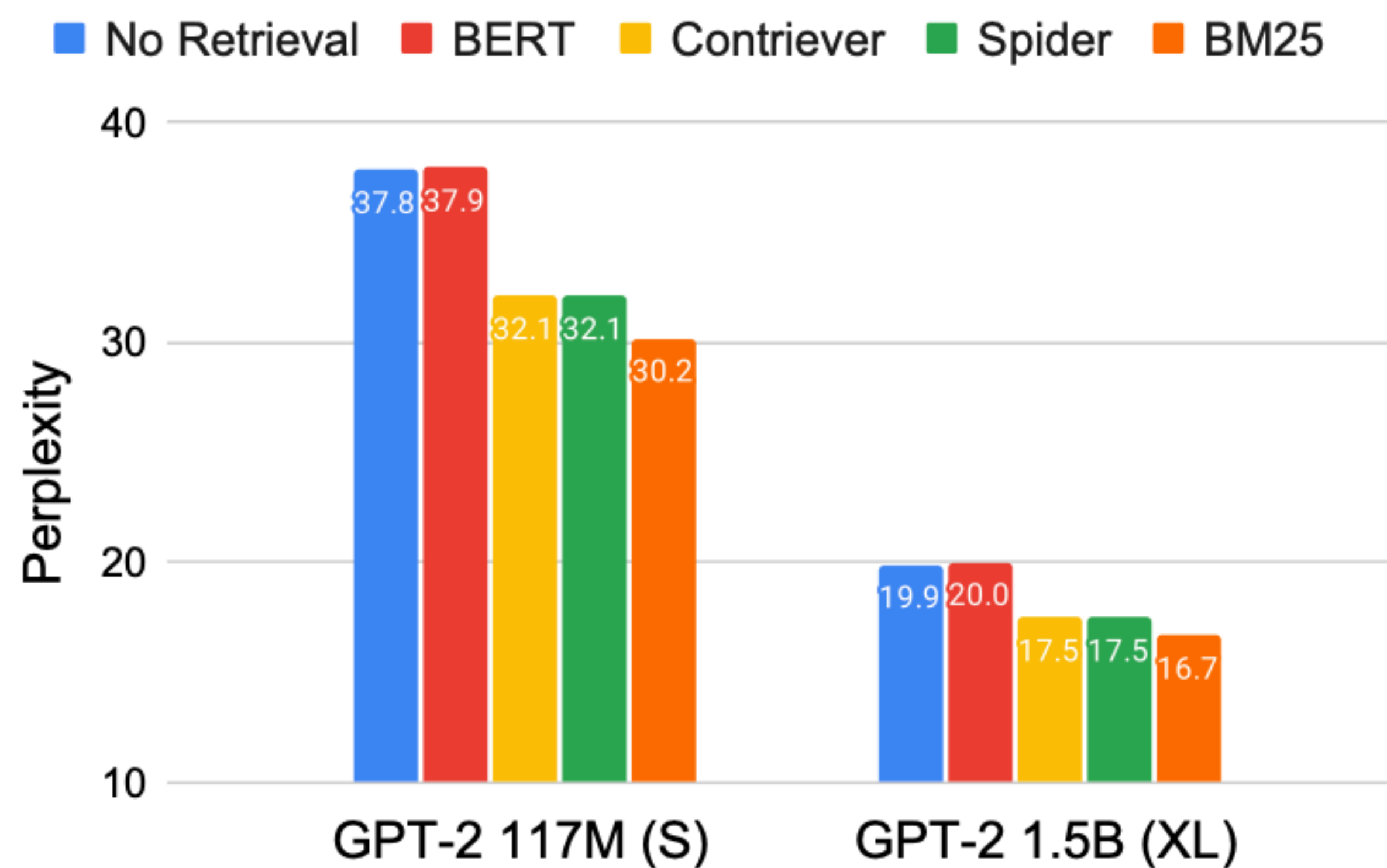
Negative passages
Too expensive to consider all negatives!

$$L(q, p^+, p_1^-, p_2^-, \dots, p_n^-)$$

Positive passage

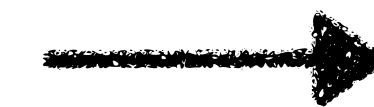
$$= -\log \frac{\exp(\text{sim}(q, p^+))}{\exp(\text{sim}(q, p^+)) + \sum_{j=1}^n \exp(\text{sim}(q, p_j^-))}$$

RAG with LMs using different retrievers



Better **retrieval model**

Better **base LMs**



Better **retrieval-based LMs**

Each component can be improved separately

Independent training



Work with off-the-shelf models (no extra training required)



Each part can be improved independently

Independent training



Work with off-the-shelf models (no extra training required)



Each part can be improved independently



LMs are not trained to leverage retrieval



Retrieval models are not optimized for LM tasks/domains

Training methods for retrieval-augmented LMs

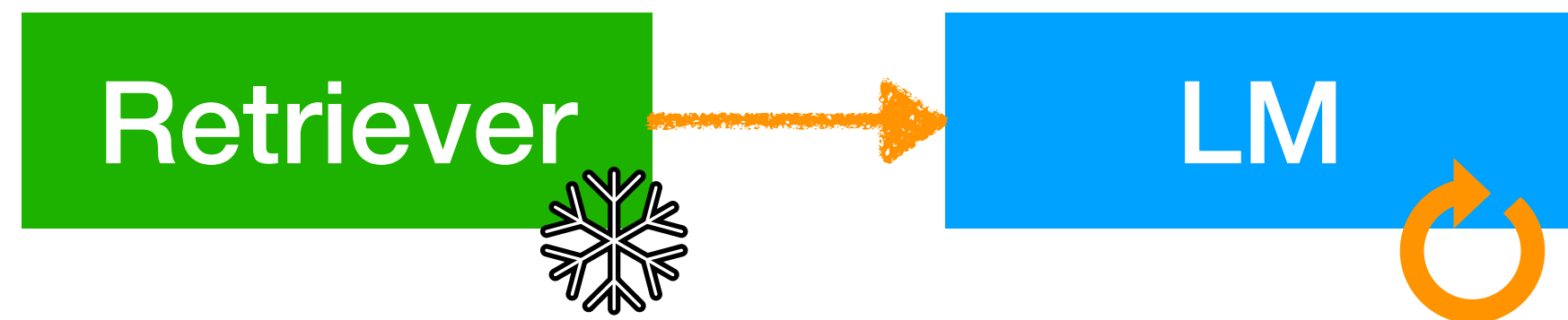
- Independent training
- **Sequential training**
- Joint training w/ asynchronous index update
- Joint training w/ in-batch approximation

Sequential training

- **One component** is first trained independently and then fixed
- **The other component** is trained with an objective that depends on the **first one**

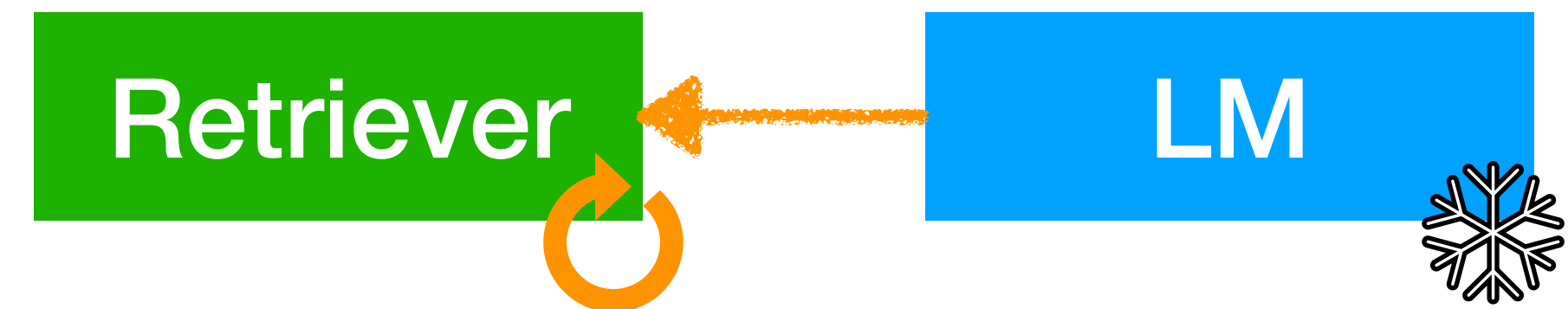
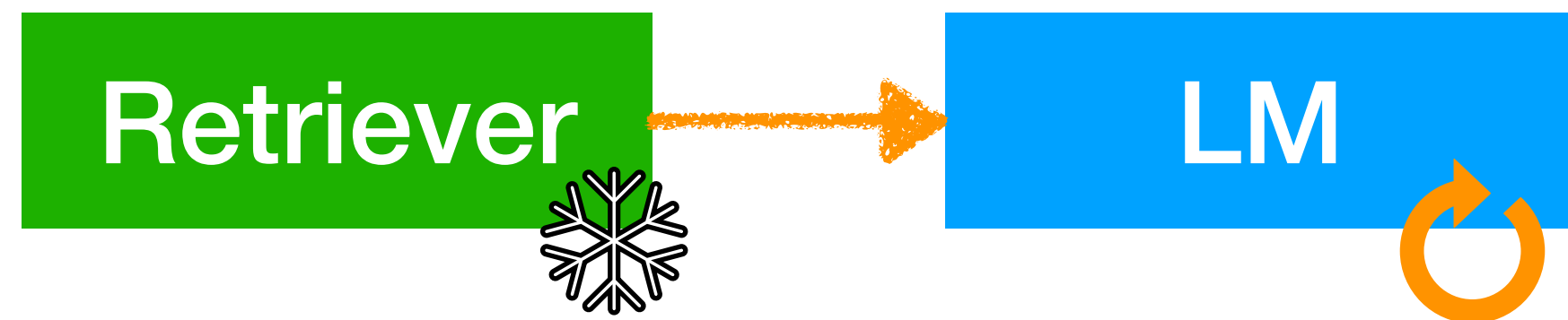
Sequential training

- One component is first trained independently and then fixed
- The other component is trained with an objective that depends on the first one



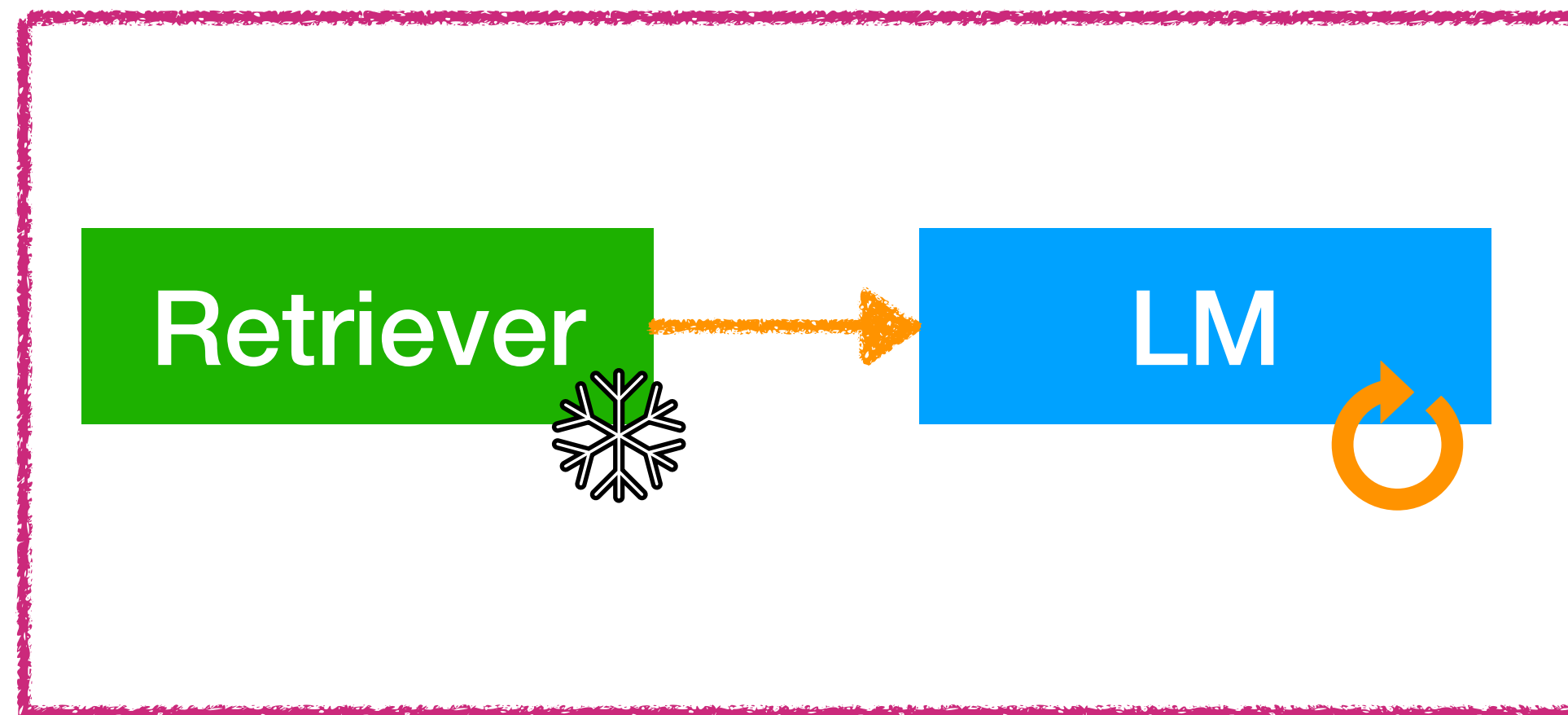
Sequential training

- One component is first trained independently and then fixed
- The other component is trained with an objective that depends on the first one

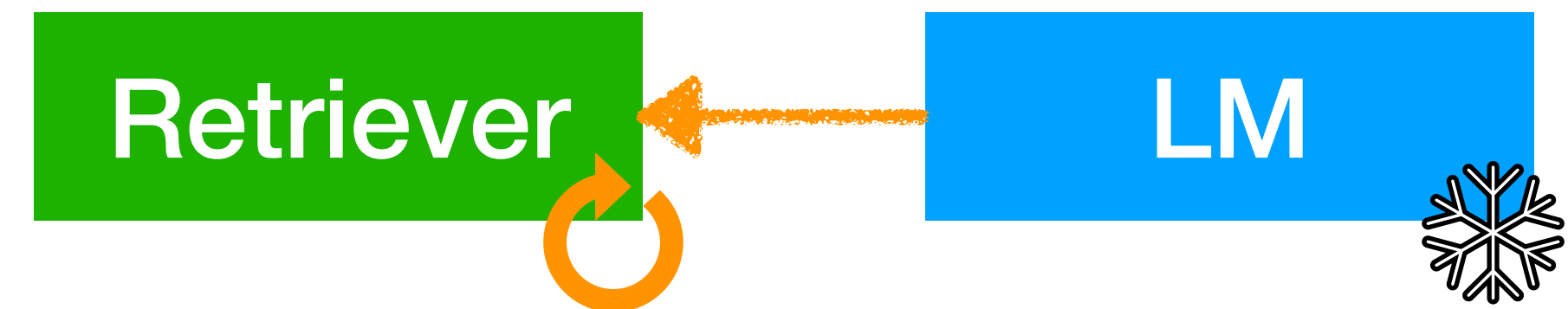


Sequential training

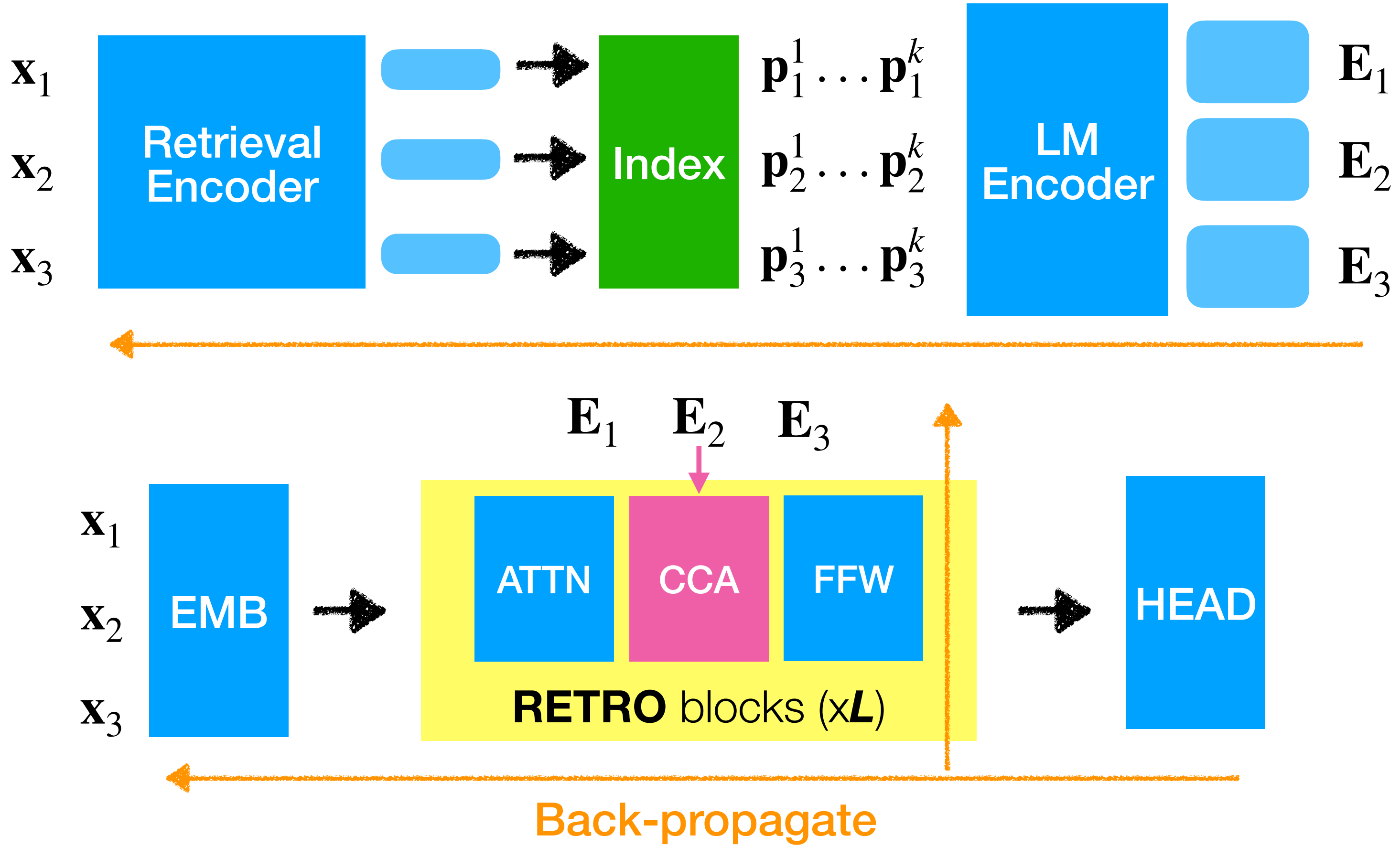
- One component is first trained independently and then fixed
- The other component is trained with an objective that depends on the first one



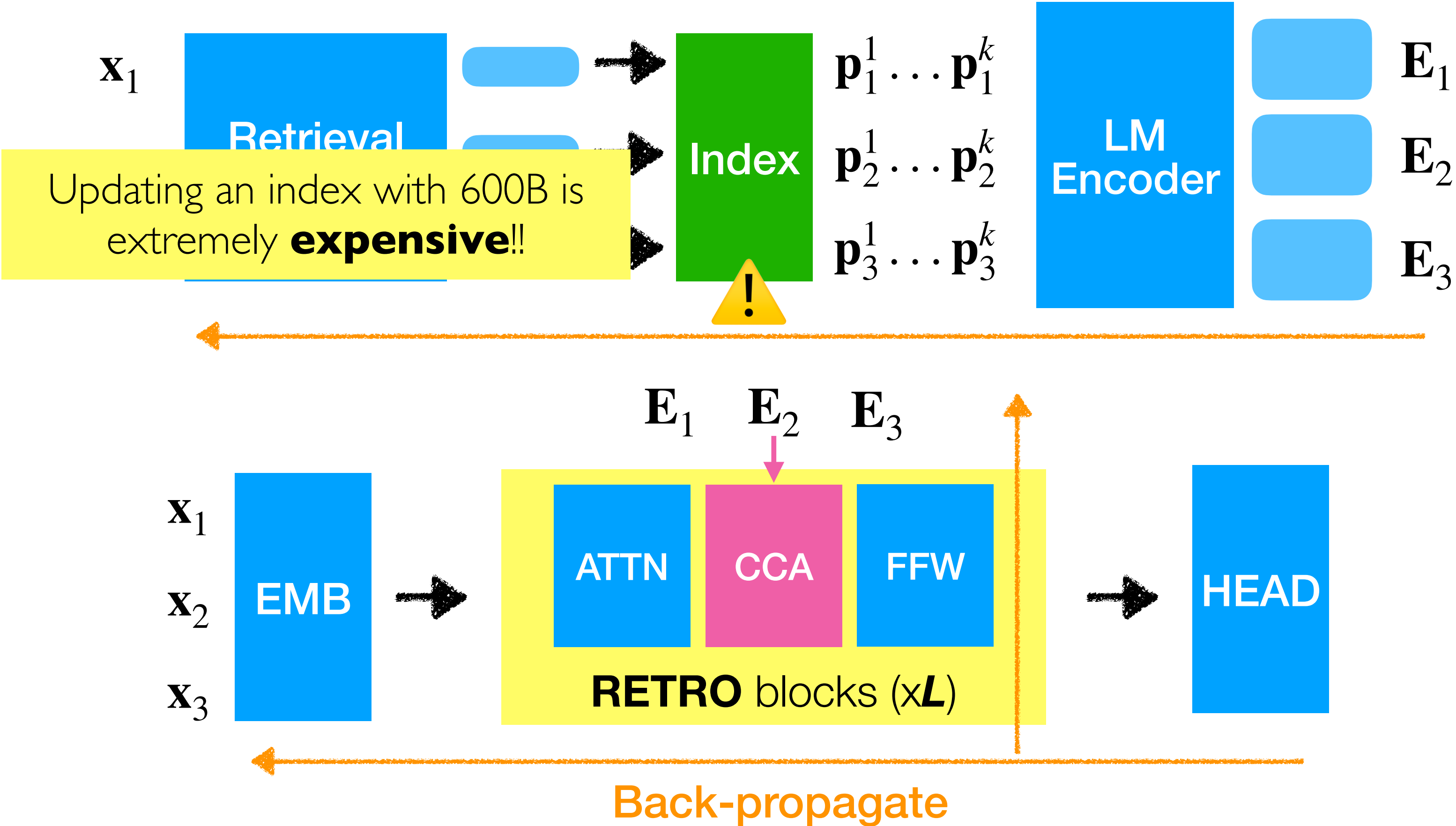
e.g., RETRO; WebGPT



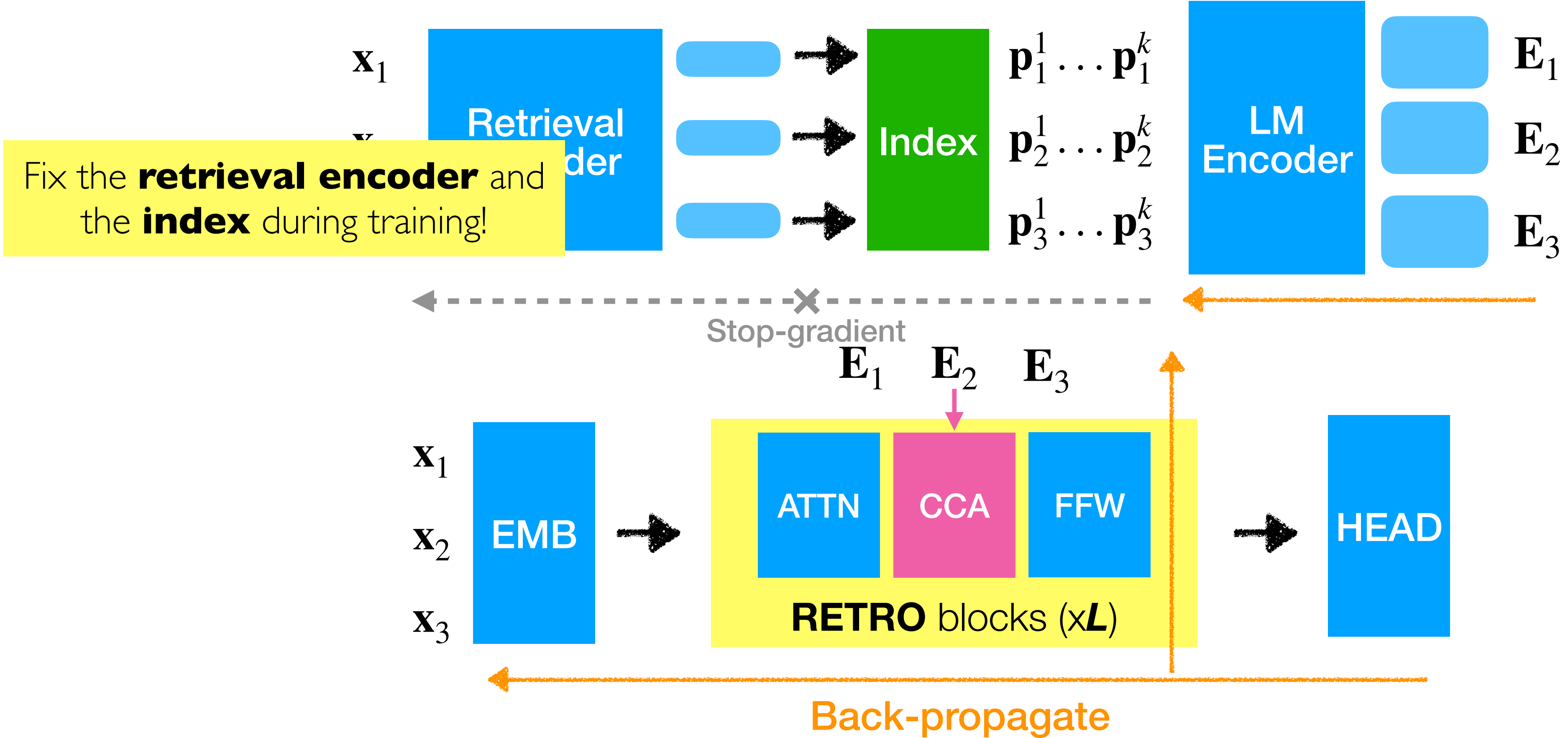
RETRO: Training



RETRO: Training



RETRO: Training



Sequential training



Work with off-the-shelf components (either a large index or a powerful LM)



LMs are trained to effectively leverage retrieval results



Retrievers are trained to provide text that helps LMs the most



One component is still fixed and not trained

Sequential training



Work with off-the-shelf components (either a large index or a powerful LM)



LMs are trained to effectively leverage retrieval results



Retrievers are trained to provide text that helps LMs the most



One component is still fixed and not trained

Let's jointly train retrieval models and LMs!

Training methods for retrieval-augmented LMs

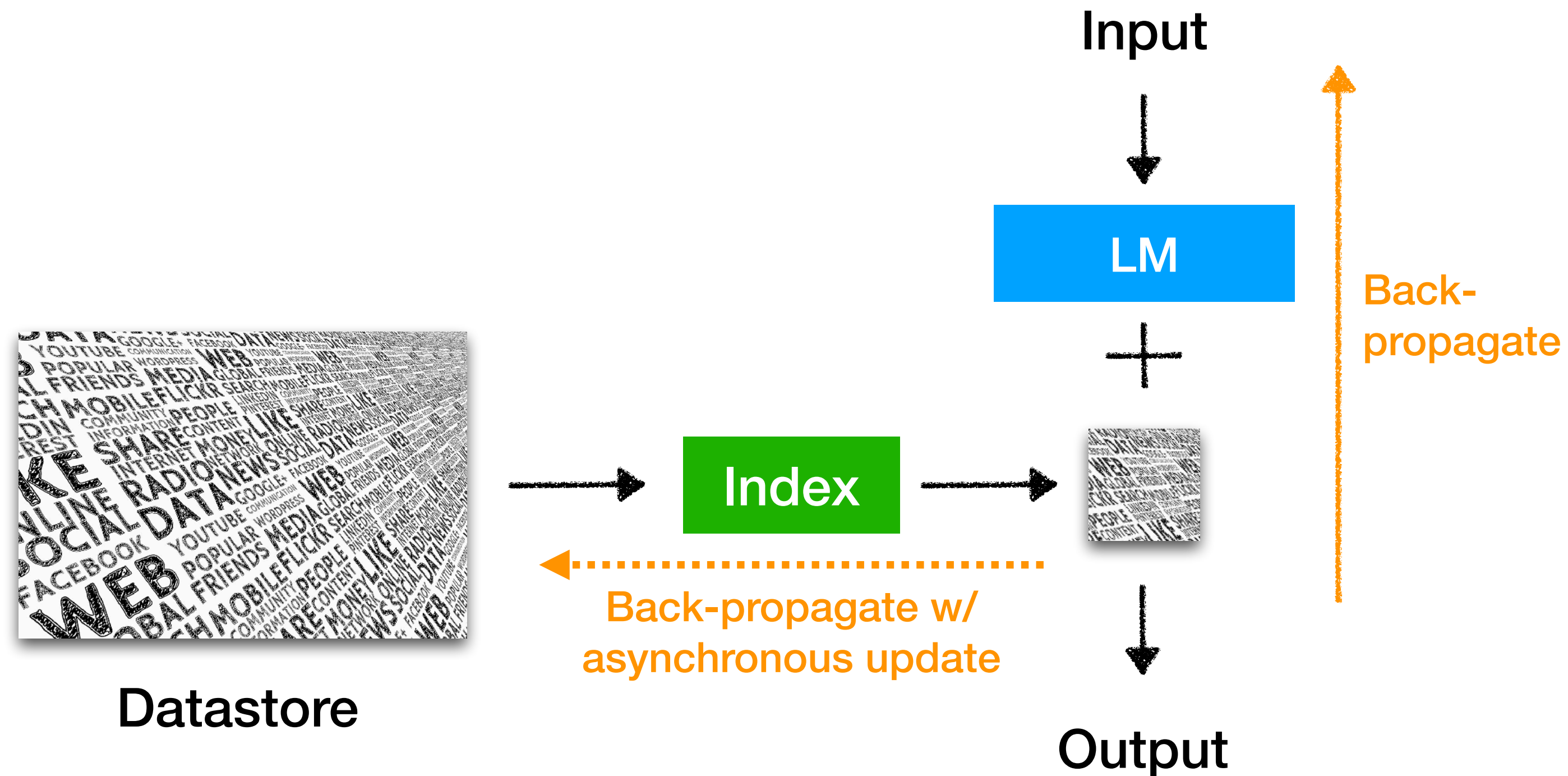
- Independent training
- Sequential training
- **Joint training w/ asynchronous index update**
- **Joint training w/ in-batch approximation**

Training methods for retrieval-augmented LMs

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

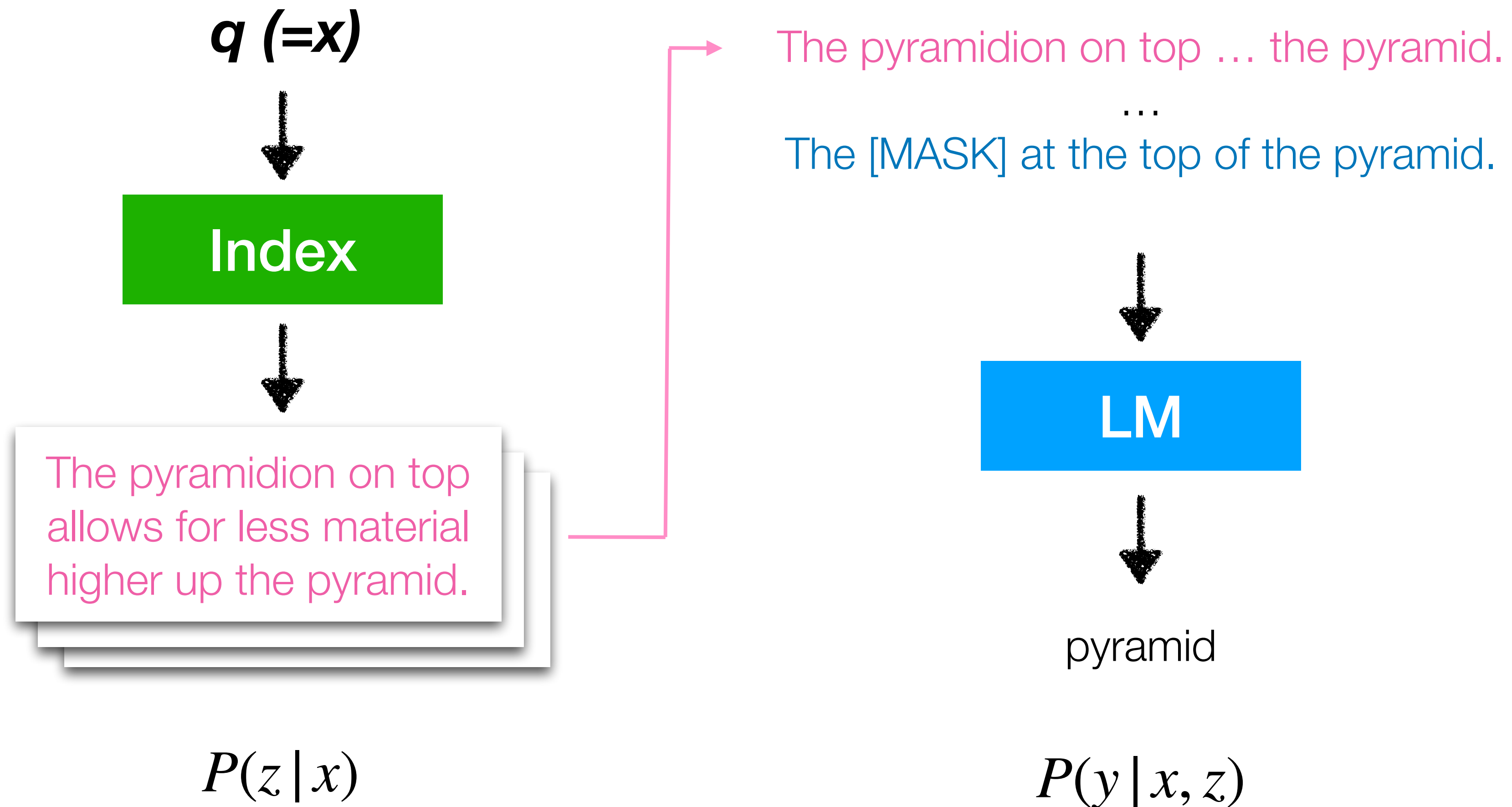
Joint training w/ asynchronous index update

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



REALM (Guu et al. 2020)

x = The [MASK] at the top of the pyramid.



REALM: Training

Objective: maximize $\sum_{z \in \mathcal{L}_\theta} P_\theta(z | q) P_\theta(y | q, z)$

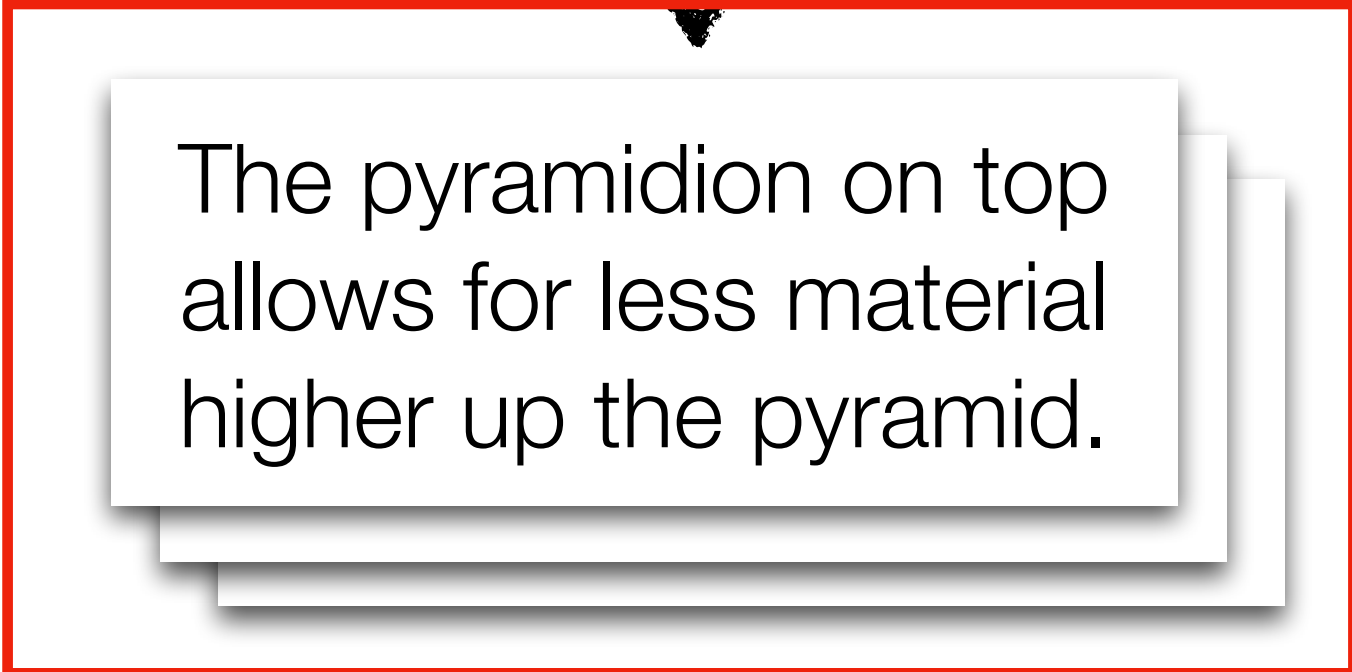
$q (=x)$



Index



\mathcal{L}_θ : top-K retrieved chunks



$P_\theta(z | x)$

The pyramidion on top ... the pyramid.
 ...
 The [MASK] at the top of the pyramid.



LM



pyramid

$P_\theta(y | x, z)$

REALM: Training

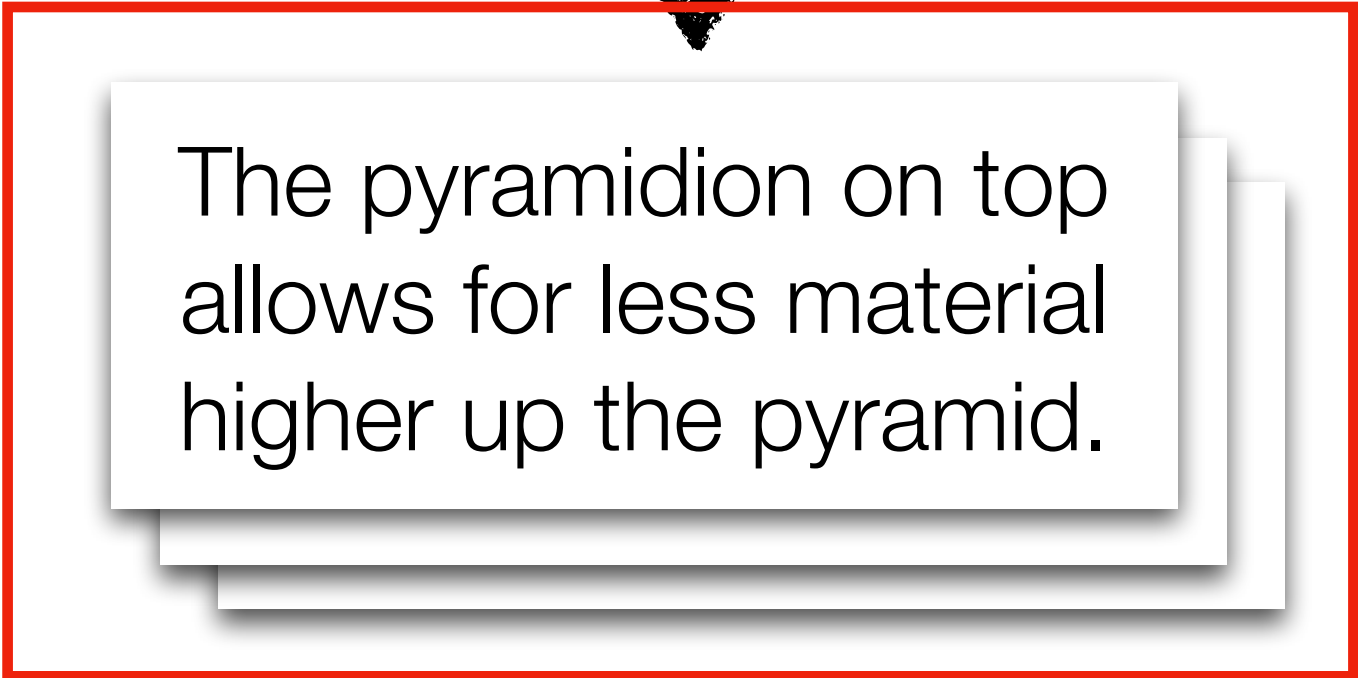
Objective: maximize $\sum_{z \in \mathcal{L}_\theta} P_\theta(z | q) P_\theta(y | q, z)$

$q (=x)$

Index

The pyramidion on top ... the pyramid.
 ...
 The [MASK] at the top of the pyramid.

\mathcal{L}_θ : top-K retrieved chunks



Back-propagation



LM

pyramid

$P_\theta(z | x)$

$P_\theta(y | x, z)$

REALM: Training

Objective: maximize $\sum_{z \in \mathcal{L}_\theta} P_\theta(z | q) P_\theta(y | q, z)$

$q (=x)$



Index

Stale index;
Update every T steps

\mathcal{L}_θ : top-K retrieved chunks



The pyramidion on top allows for less material higher up the pyramid.

The pyramidion on top ... the pyramid.
...
The [MASK] at the top of the pyramid.



LM



pyramid

$P_{\theta_{\text{new}}}(z | x)$

$P_{\theta_{\text{new}}}(y | x, z)$

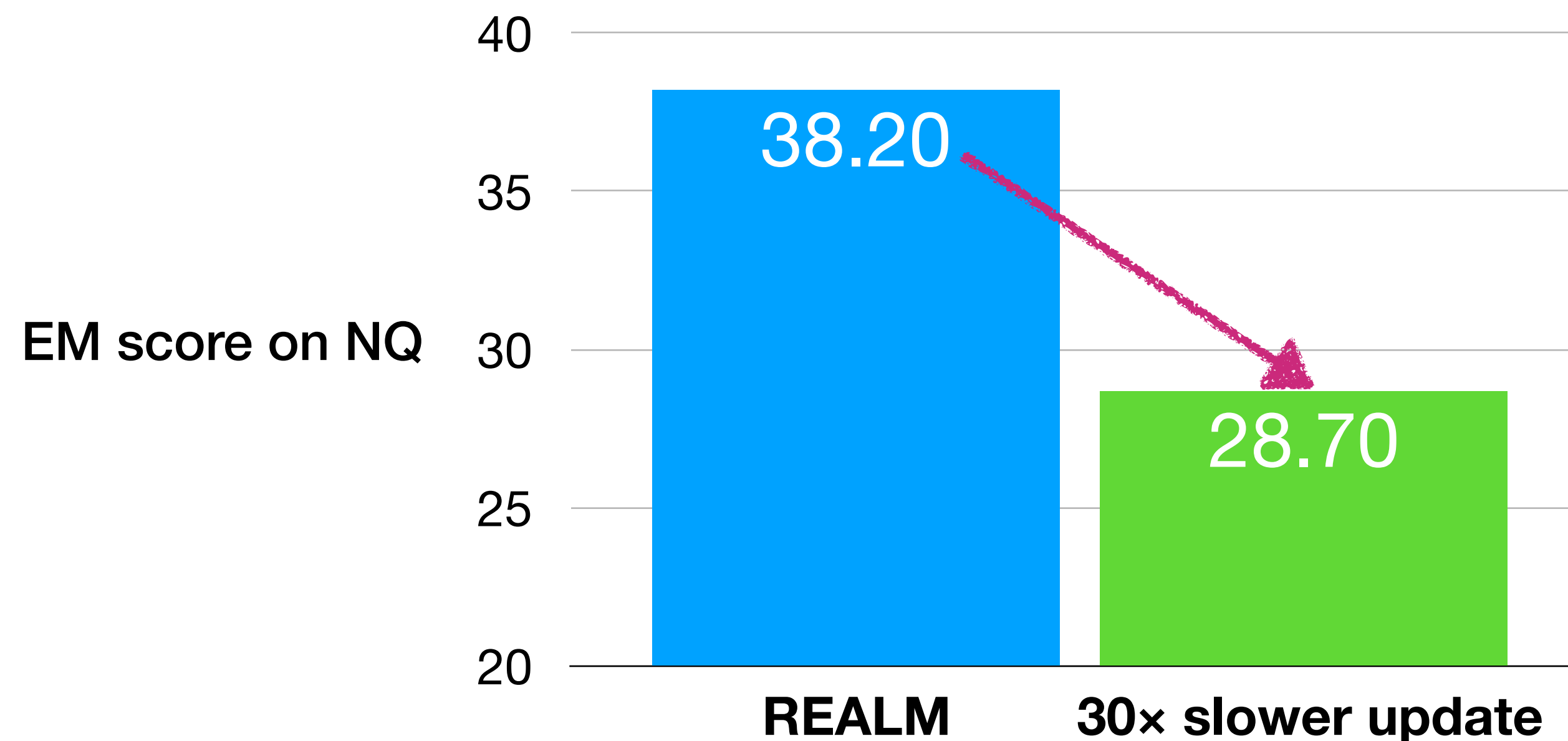
Up-to-date parameters

REALM: Index update rate

How often should we update the retrieval index?

- Frequency too high: expensive
- Frequency too slow: out-dated

REALM: updating the index every 500 training steps



Joint training



End-to-end trained — each component is optimized



Good performance



Training is more complicated
(async update, overhead, data batching, etc)



Train-test discrepancy still remains

Today's outline

Question:

[https://bit.ly/
akari_ralm_lec](https://bit.ly/akari_ralm_lec)



Scan me

Why do we need retrieval-augmented LMs?

Architectures of retrieval-augmented LMs (Inference)

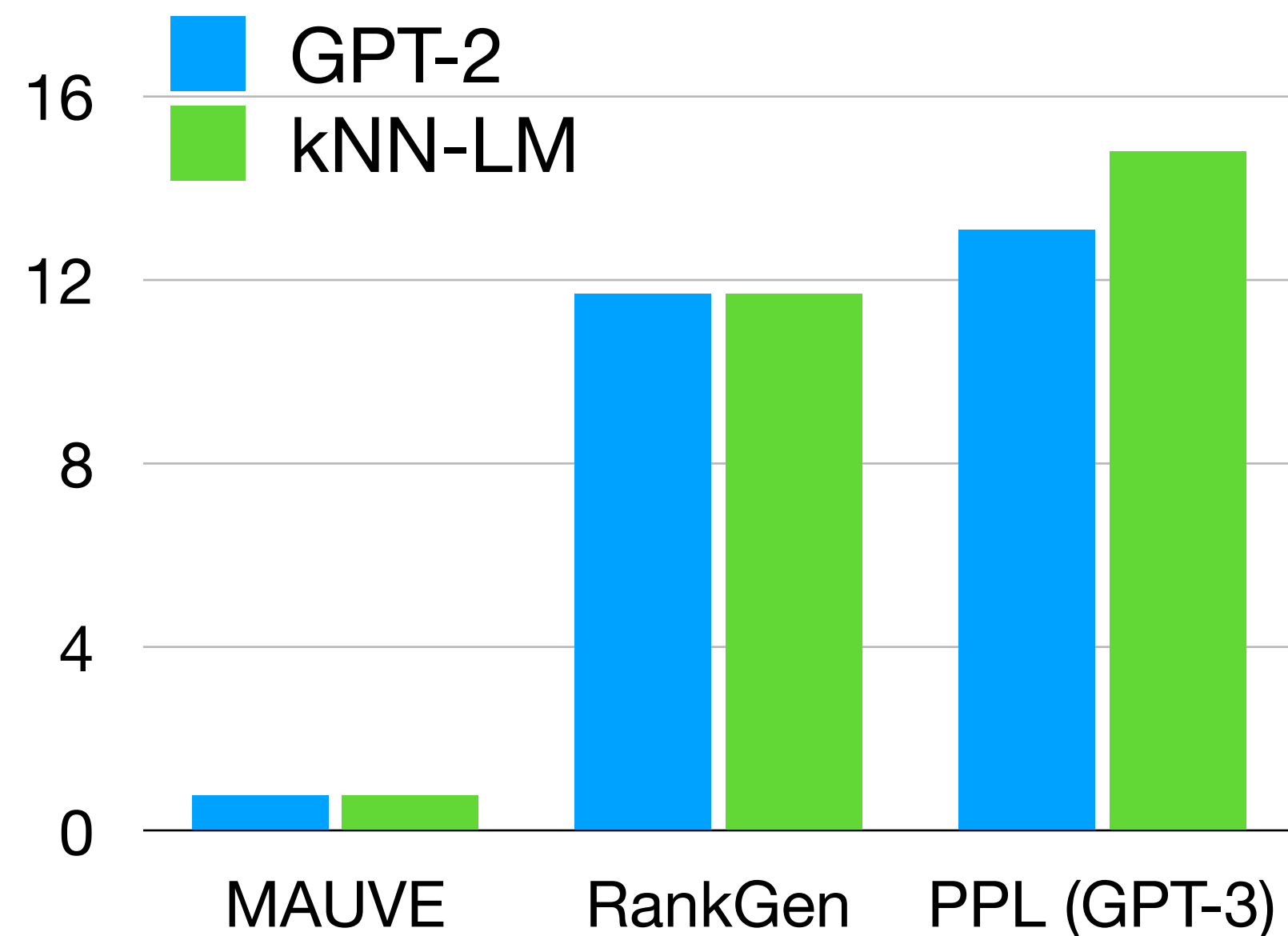
Training of retrieval-augmented LMs

Limitations and future directions

Challenge: retrieval-augmented LMs for applications

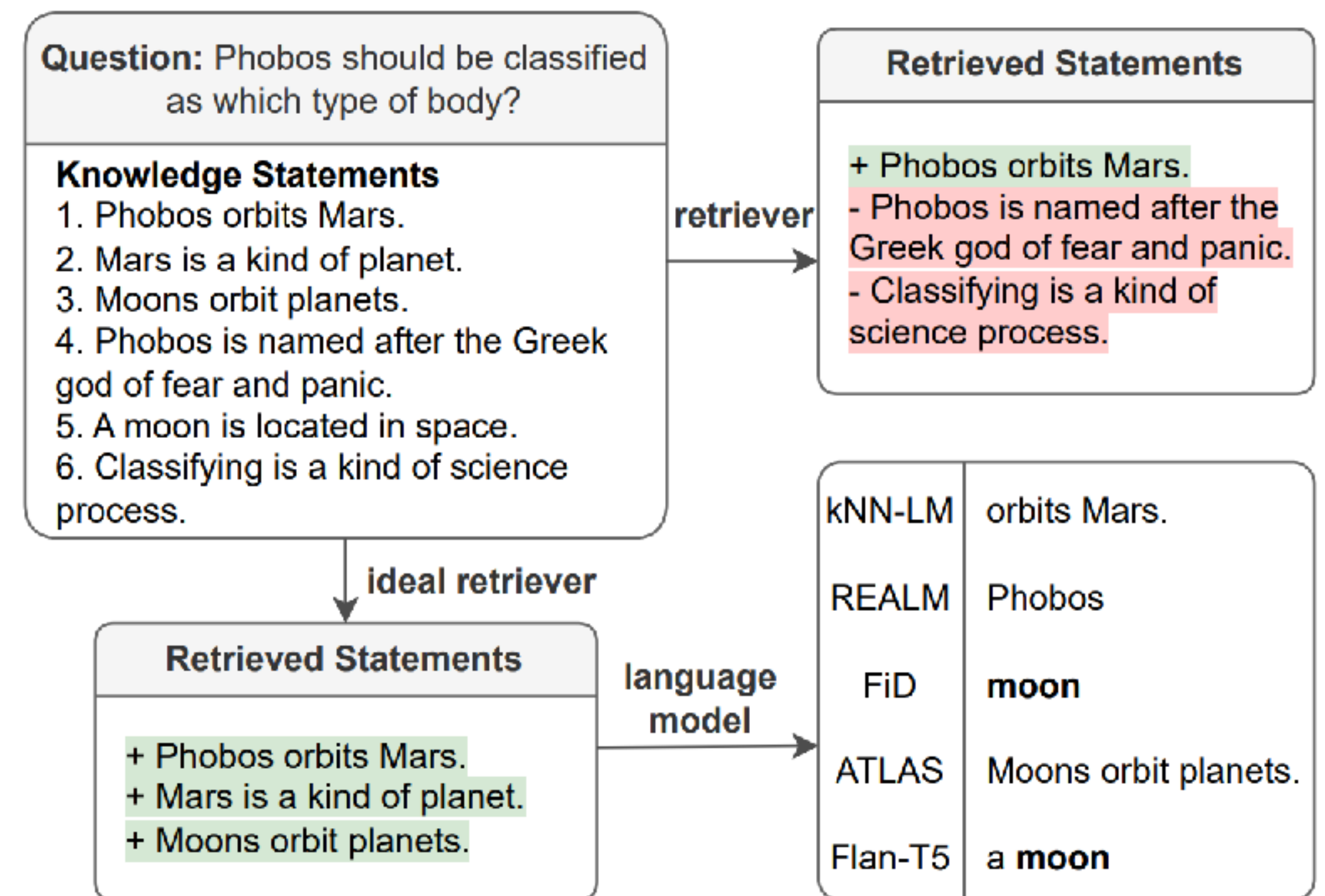
Open-ended text generation? Reasoning?

Doesn't improve open-ended generation



Wang et al. kNN-LM Does Not Improve Open-ended Text Generation. ACL 2023.

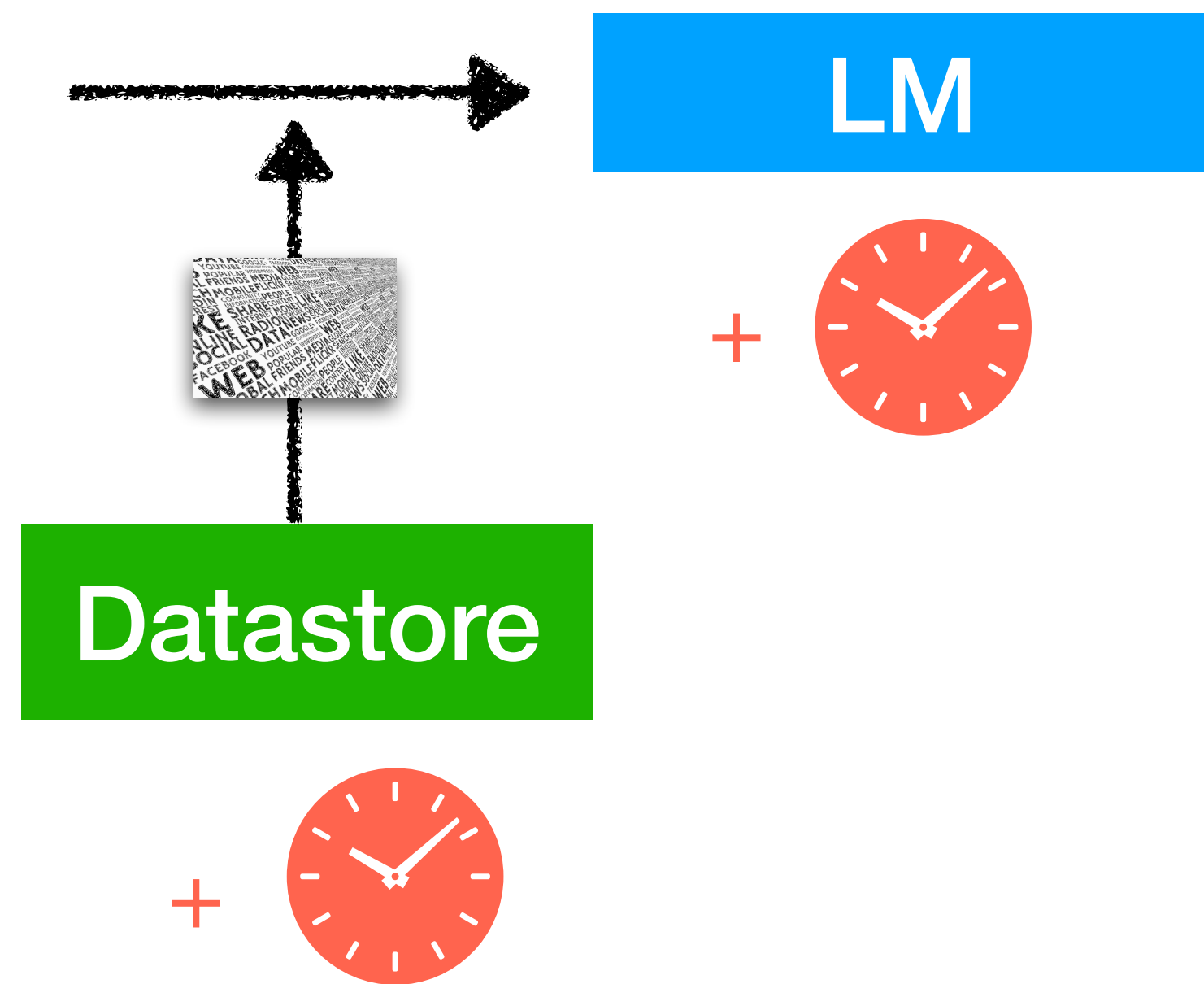
Failure of retrieval in reasoning task



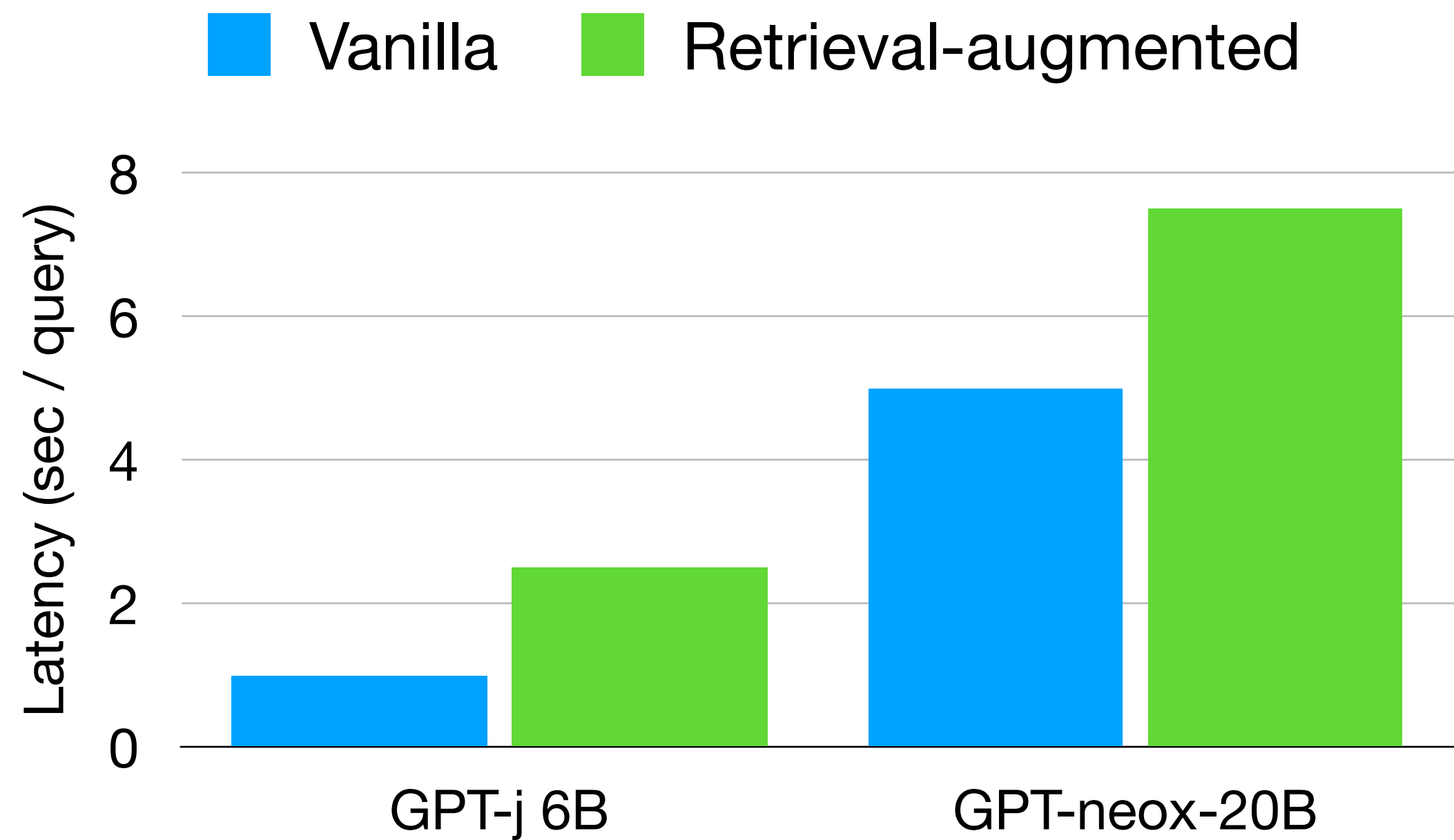
BehnamGhader et al. Can Retriever-Augmented Language Models Reason? The Blame Game Between the Retriever and the Language Model. EMNLP Findings 2023.

Challenge: efficiency retrieval-augmented LMs

Additional costs from retrieval augmentation

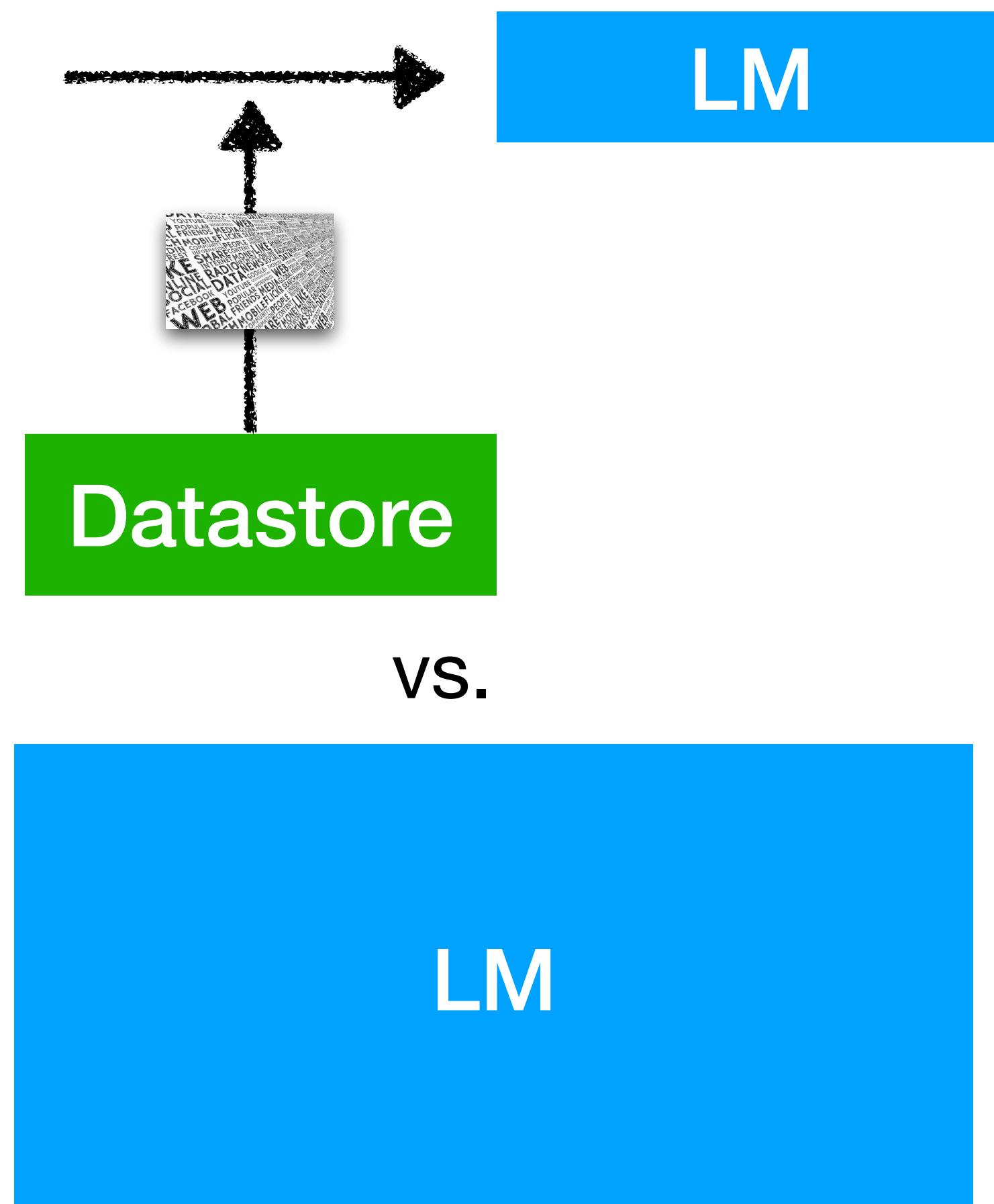


Retrieval-augmented LMs add inference costs



Challenge: scaling retrieval-augmented LMs


A small LM + a large datastore \approx a large parametric LM?




	LM	Datastore
	# of parameters	# of tokens
kNN-LM (Khandelwal et al., 2020)	250M	$\leq 3B$
NPM (Min et al., 2023)	350M	1B
Atlas (Izacard et al., 2022)	11B	$\sim 30B$
RETRO (Borgeaud et al., 2021)	7B	2T
REPLUG (Shi et al., 2023)	$\leq 175B$	$\sim 5B$

Challenge: robustness and controllability


Retrieval-augmented LMs can still hallucinate


 What are the latest discoveries from the James Webb Space Telescope?


 The James Webb Space Telescope is designed to peer into the dusty clouds of gas where stars and planetary systems are born. Webb has captured the first direct image of an exoplanet, and the Pillars of Creation in the Eagle Nebula[1][2]. Additionally, the telescope will be used to study the next interstellar interloper[3].

*(*Some generated statements may not be fully supported by citations, while others are fully supported.)*

Cited Webpages

[1]:  nasa.gov (✗ citation does not support its associated statement)
[NASA's Webb Confirms Its First Exoplanet](#)
... Researchers confirmed an exoplanet, a planet that orbits another star, using NASA's James Webb Space Telescope for the first time. ...

[2]:  cnn.com (⚠ citation partially supports its associated statement)
[Pillars of Creation: James Webb Space Telescope ...](#)
... The Pillars of Creation, in the Eagle Nebula, is a star-forming region captured in a new image (right) by the James Webb Space Telescope that reveals more detail than a 2014 image (left) by Hubble ...

[3]:  nasa.gov (✓ citation fully supports its associated statement)
[Studying the Next Interstellar Interloper with Webb](#)
...Scientists have had only limited ability to study these objects once discovered, but all of that is about to change with NASA's James Webb Space Telescope...The team will use Webb's spectroscopic capabilities in both the near-infrared and mid-infrared bands to study two different aspects of the interstellar object.

Liu et al. Evaluating Verifiability in Generative Search Engines. Findings of EMNLP 2023.

Roadmap to advance retrieval-augmented LMs

Rethink Retrieval and Datastore

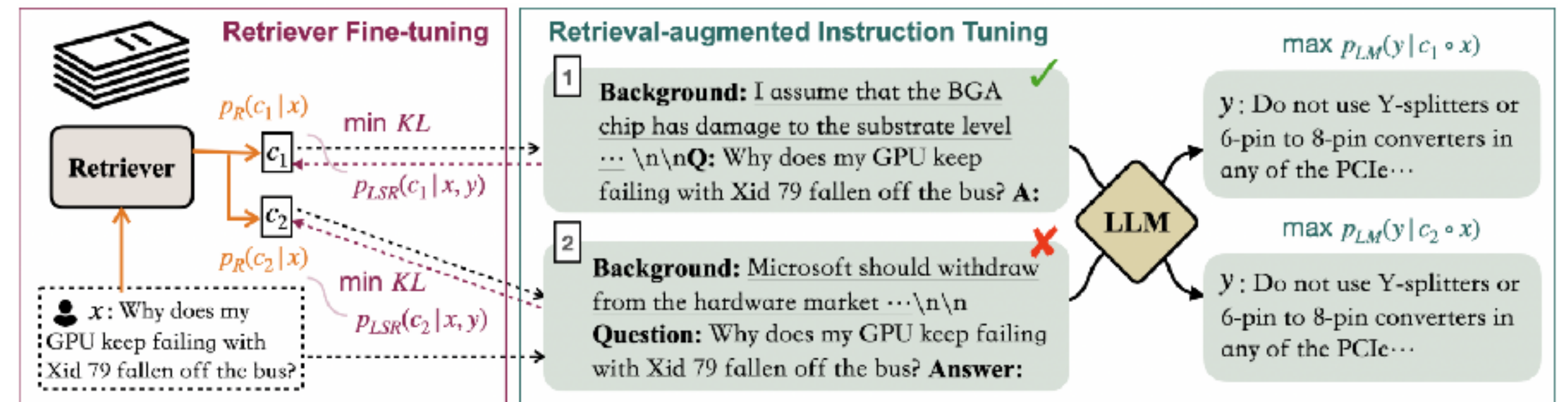
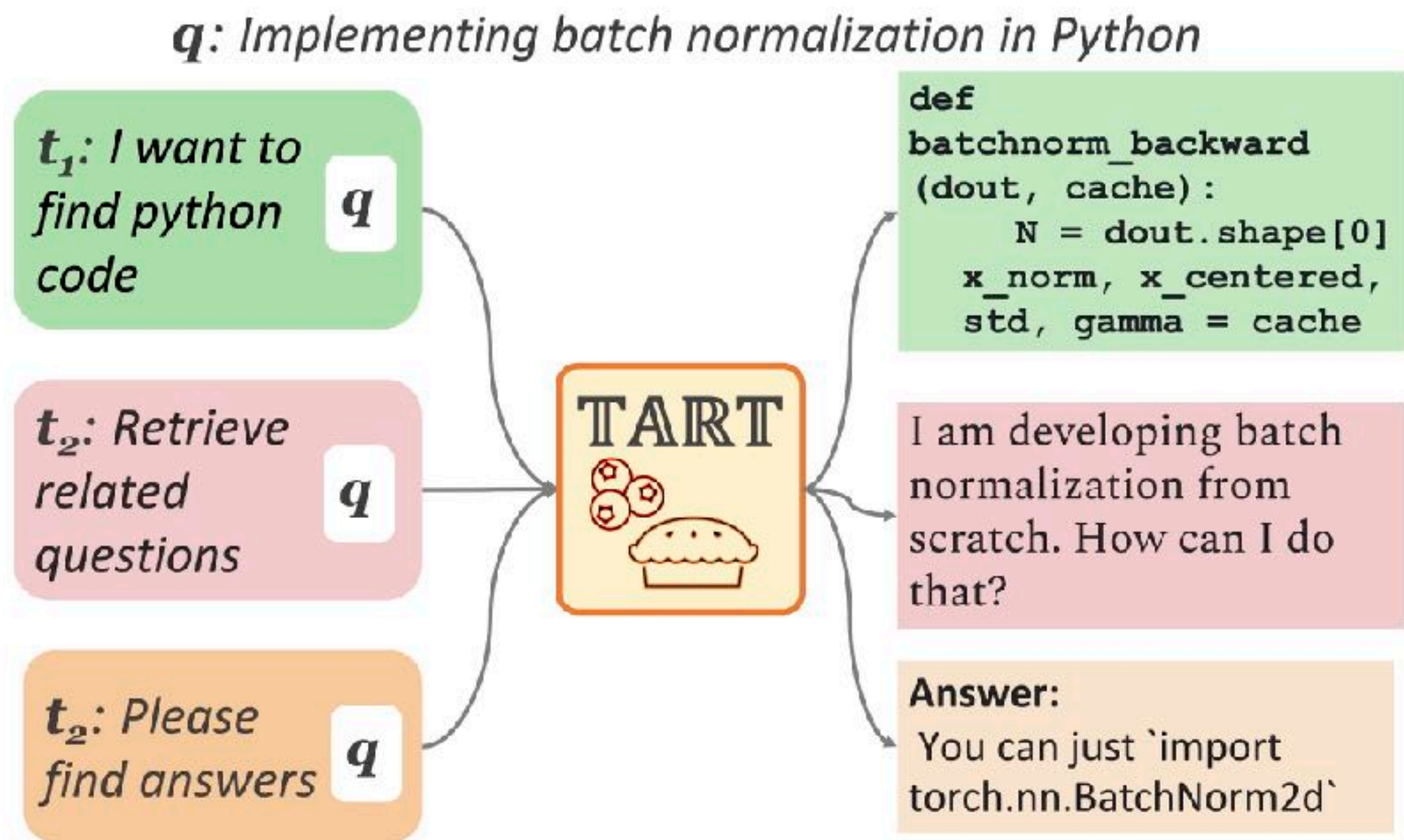


Advance Architectures & Retrieval-aware Training

Investment Infrastructures for Training and Inference at Scale

Beyond semantic and lexical-similarity based search

Training retrievers to optimize end-to-end retrieval-augmented LM performance in diverse tasks



<i>0-shot</i>	BoolQ	PIQA	SIQA	HellaSwag	WinoGrande
LLAMA 65B	85.3	82.8	52.3	84.2	77.0
RA-DIT 65B w/o retrieval	86.7	83.7	57.9	85.1	79.8
RA-DIT 65B	85.6	84.4	58.4	85.4	80.0

Asai et al., Task-aware Retrieval with Instruction. Findings of ACL 2023.

Lin et al., RA-DIT: Retrieval-Augmented Dual Instruction Tuning. ICLR 2024.

Roadmap to advance retrieval-augmented LMs

Rethink Retrieval and Datastore

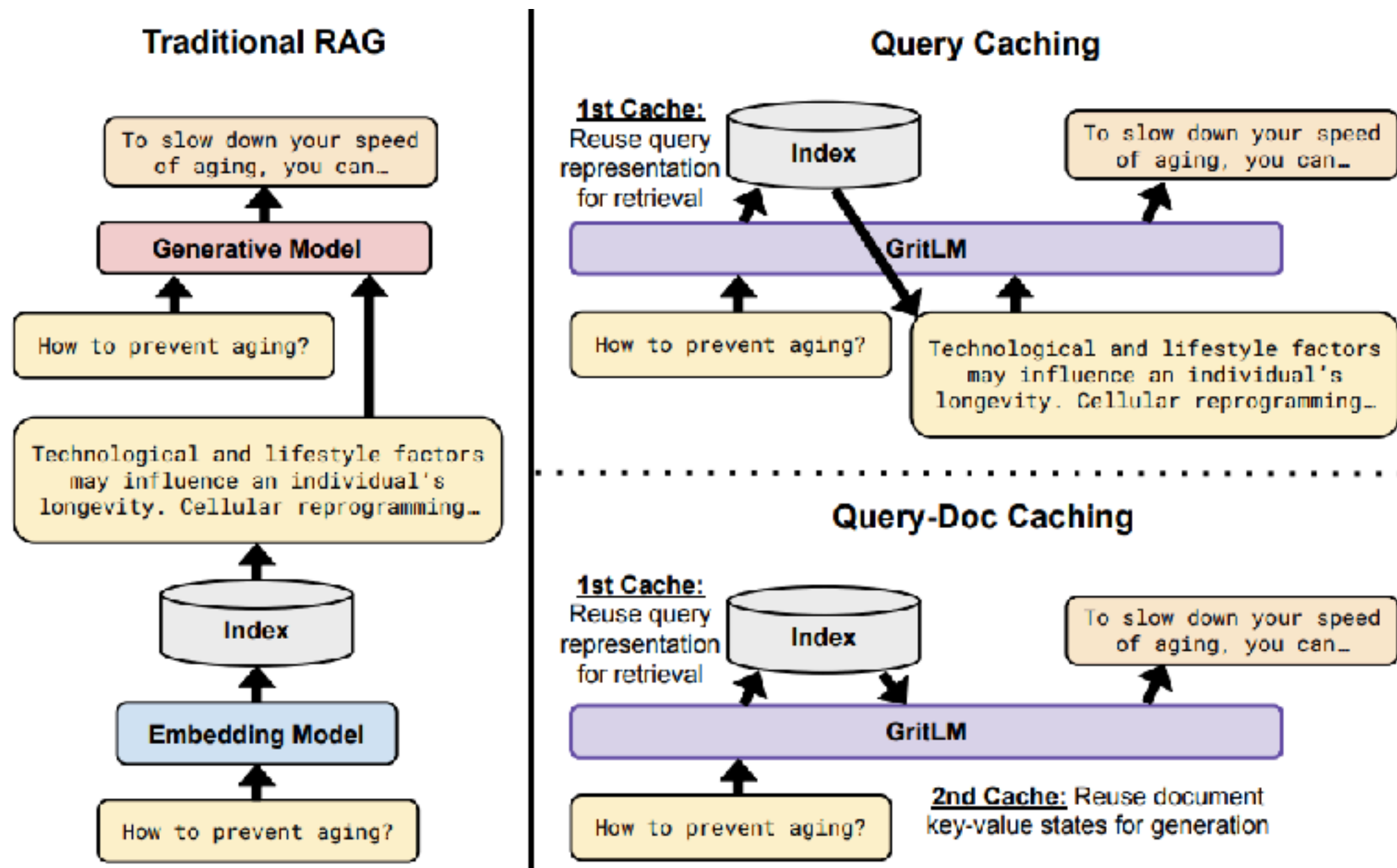


Advance Architectures & Retrieval-aware Training

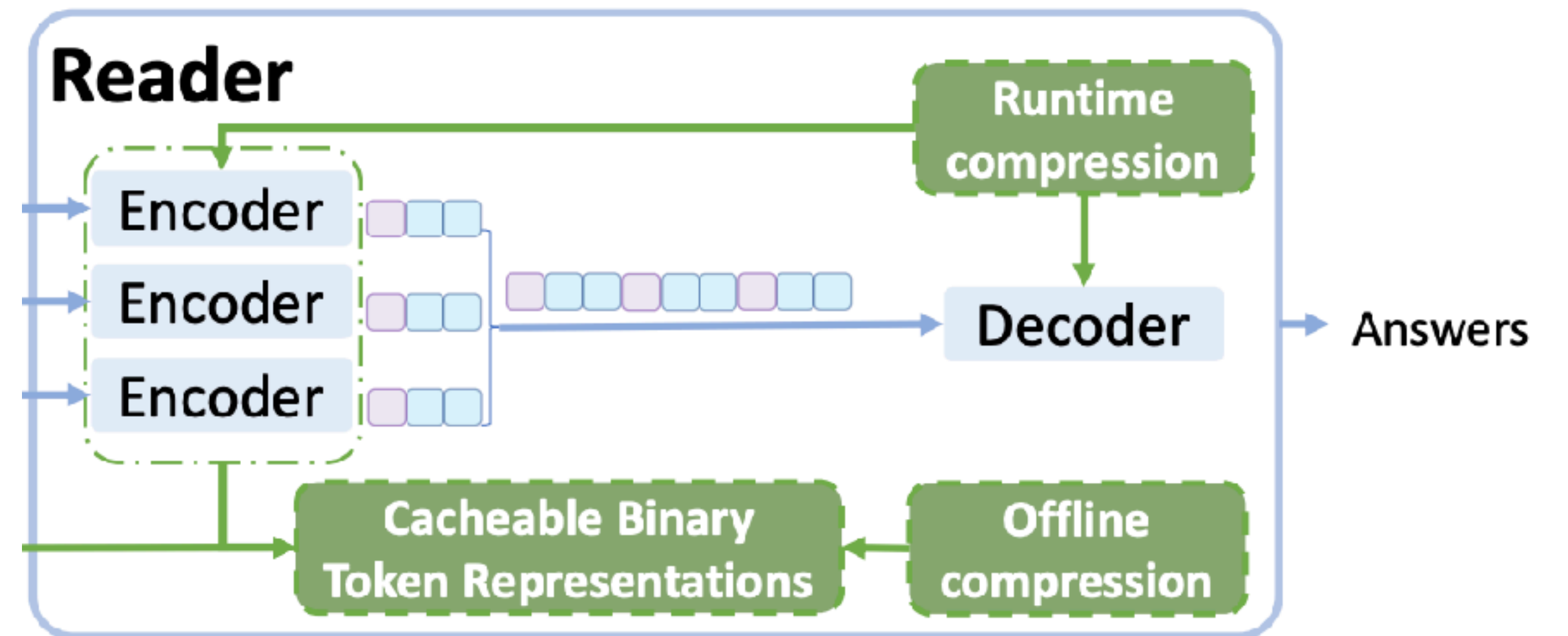
Investment Infrastructures for Training and Inference at Scale

New architectures for performance and efficiency

Further explorations of unified architectures & caching



Muennighoff et al. Generative Representational Instruction Tuning. 2024.



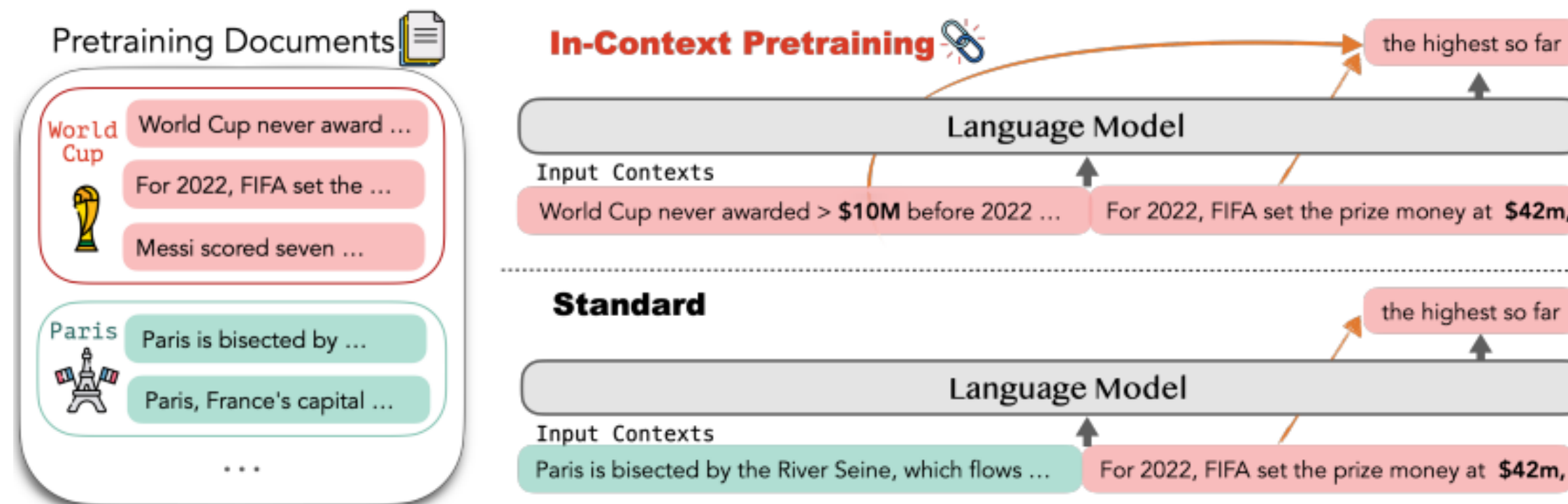
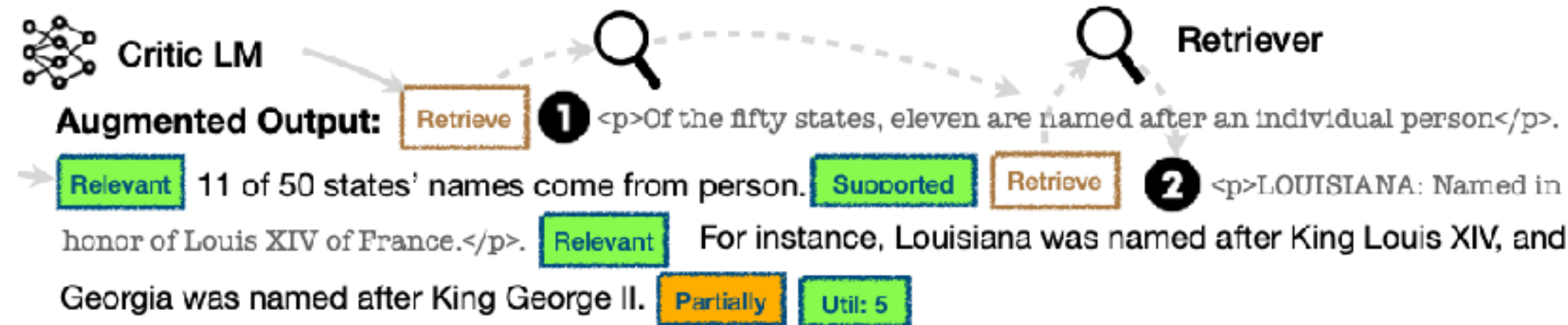
Cao et al. BTR: Binary Token Representations for Efficient Retrieval Augmented Language Models. ICLR 2024.

Training LMs with Retrieval

Training LMs to learn to use retrieval during pre-training or instruction-tuning

Input: How did US states get their names?

Output: 1 of 50 states names come from persons. For instance, Louisiana was named in honor of King Louis XIV of France and Georgia was named after King George II.



Instruction-tuning with retrieval

Asai et al. Self-RAG: Learning to Retrieve, Generate and Critique with Retrieval. ICLR 2024.

Retrieval-aware pre-training

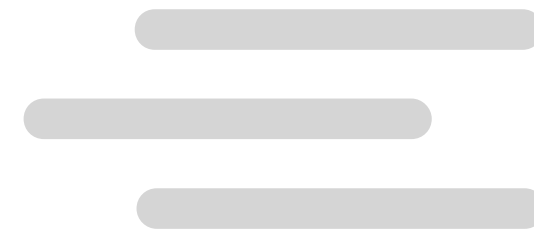
Shi. et al. In-Context Pretraining: Language Modeling Beyond Document Boundaries. ICLR 2024.

Roadmap to advance retrieval-augmented LMs

Rethink Retrieval and Datastore

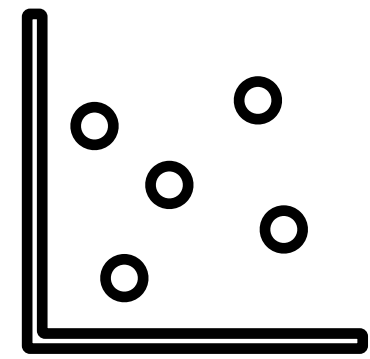


Advance Architectures & Retrieval-aware Training

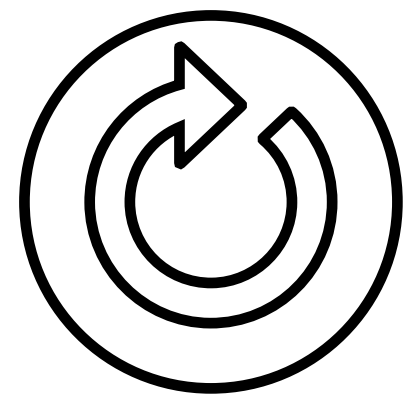


Investment Infrastructures for Training and Inference at Scale

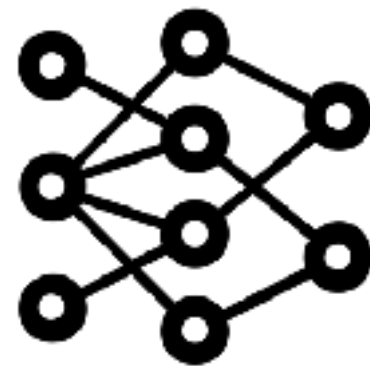
More open-sourced and collaborative opportunities



System / Algorithmic improvements for **massive Datastore**



Standardized implementations for **efficient training**



Fast inference algorithms for retrieval-augmented LMs

Summary & QA

Question:

[https://bit.ly/
akari_ralm_lec](https://bit.ly/akari_ralm_lec)



Scan me

Retrieval-augmented LMs can solve many issues e.g., hallucinations

Various architectures (not just RAG) exist with different pros&cons

Jointly training retrieval-augmented LMs is important but hard

Many interesting research opportunities — let's work together!

ACL 2023 tutorial: <https://acl2023-retrieval-lm.github.io/>

Position paper: [https://akariasai.github.io/assets/
pdf/ralm_position.pdf](https://akariasai.github.io/assets/pdf/ralm_position.pdf)

Contact: akari@cs.washington.edu

Website: <https://akariasai.github.io/>

Twitter: @AkariAsai

References (I)

Akari Asai, Sewon Min, Zexuan Zhong, Danqi Chen. Retrieval-based Language Models and Applications. ACL Tutorial 2023.

Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, Matt Gardner. Documenting Large Webtext Corpora: A Case Study on the Colossal Clean Crawled Corpus. EMNLP 2021.

Alex Mallen*, Akari Asai*, Victor Zhong, Rajarshi Das, Daniel Khashabi, Hannaneh Hajishirzi. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories. ACL 2023.

Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, Colin Raffel. Large Language Models Struggle to Learn Long-Tail Knowledge. ICML 2023.

Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A. Smith, Yejin Choi, Kentaro Inui. RealTime QA: What's the Answer Right Now?. NeurIPS (Dataset & Benchmark) 2023.

References (2)

Sewon Min, Suchin Gururangan, Eric Wallace, Hannaneh Hajishirzi, Noah A. Smith, Luke Zettlemoyer. SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore. ICLR 2024.

Juan Ramos. Using TF-IDF to Determine Word Relevance in Document Queries. The first instructional conference on machine learning 2003.

Robertson and Zaragoza. The Probabilistic Relevance Framework: BM25 and Beyond. Foundations and Trends® in Information Retrieval 2009.

Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, Wen-tau Yih. Dense Passage Retrieval for Open-Domain Question Answering. EMNLP 2020.

Akari Asai, Zexuan Zhong, Danqi Chen, Pang Wei Koh, Luke Zettlemoyer, Hannaneh Hajishirzi, Wen-tau Yih. Reliable, Adaptable, and Attributable Language Models with Retrieval. Arxiv 2024.

References (3)

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, Ming-Wei Chang. REALM: Retrieval-Augmented Language Model Pre-Training. ICML 2020.

Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, Yoav Shoham. In-Context Retrieval-Augmented Language Models. arXiv 2023.

Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, Wen-tau Yih. REPLUG: Retrieval-Augmented Black-Box Language Models. arXiv 2023.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, Laurent Sifre. Improving language models by retrieving from trillions of tokens. arXiv 2021.

Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, Mike Lewis. Generalization through Memorization: Nearest Neighbor Language Models. ICLR 2020.

References (4)

Shufan Wang, Yixiao Song, Andrew Drozdov, Aparna Garimella, Varun Manjunatha, Mohit Iyyer. KNN-LM Does Not Improve Open-ended Text Generation. ACL 2023.

Parishad BehnamGhader, Santiago Miret, Siva Reddy. Can Retriever-Augmented Language Models Reason? The Blame Game Between the Retriever and the Language Model. Findings of EMNLP 2023.

Nelson F. Liu, Tianyi Zhang, Percy Liang. Evaluating Verifiability in Generative Search Engines. Findings of EMNLP 2023.

Akari Asai, Timo Schick, Patrick Lewis, Xilun Chen, Gautier Izacard, Sebastian Riedel, Hannaneh Hajishirzi, Wen-tau Yih. Task-aware Retrieval with Instructions. Findings of ACL 2023.

Xi Victoria Lin, Xilun Chen, Mingda Chen, Weijia Shi, Maria Lomeli, Rich James, Pedro Rodriguez, Jacob Kahn, Gergely Szilvassy, Mike Lewis, Luke Zettlemoyer, Scott Yih. RA-DIT: Retrieval-Augmented Dual Instruction Tuning. ICLR 2024.

References (5)

Niklas Muennighoff, Hongjin Su, Liang Wang, Nan Yang, Furu Wei, Tao Yu, Amanpreet Singh, Douwe Kiela. Generative Representational Instruction Tuning. Arxiv 2024.

Qingqing Cao, Sewon Min, Yizhong Wang, Hannaneh Hajishirzi. BTR: Binary Token Representations for Efficient Retrieval Augmented Language Models. ICLR 2024.

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, Hannaneh Hajishirzi. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. ICLR 2024.

Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou, Margaret Li, Rich James, Xi Victoria Lin, Noah A. Smith, Luke Zettlemoyer, Scott Yih, Mike Lewis. In-Context Pretraining: Language Modeling Beyond Document Boundaries. ICLR 2024.