# Retrieval-based Language Models

**Sewon Min**

University of Washington
shmsw25.github.io

CPSC 488/588 • Fall 2023 • Yale University

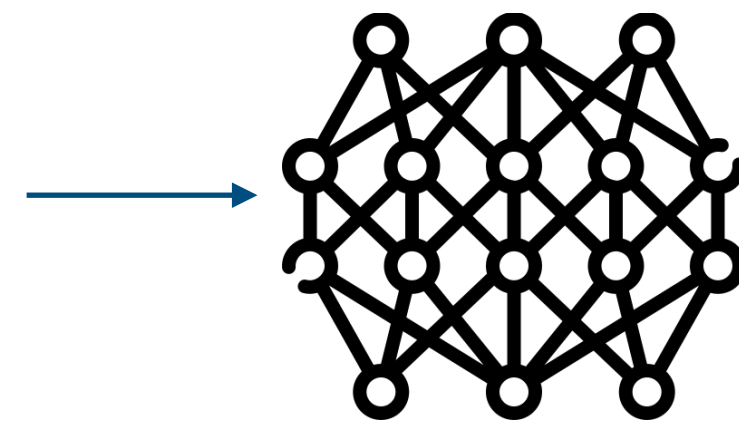Adapted from ACL 2023 Tutorial w/ Akari Asai, Zexuan Zhong, & Danqi Chen

# Language Models

$$P(x_n \mid x_1, x_2, \cdots, x_{n-1})$$

# Language Models

$$P(x_n \mid x_1, x_2, \cdots, x_{n-1})$$

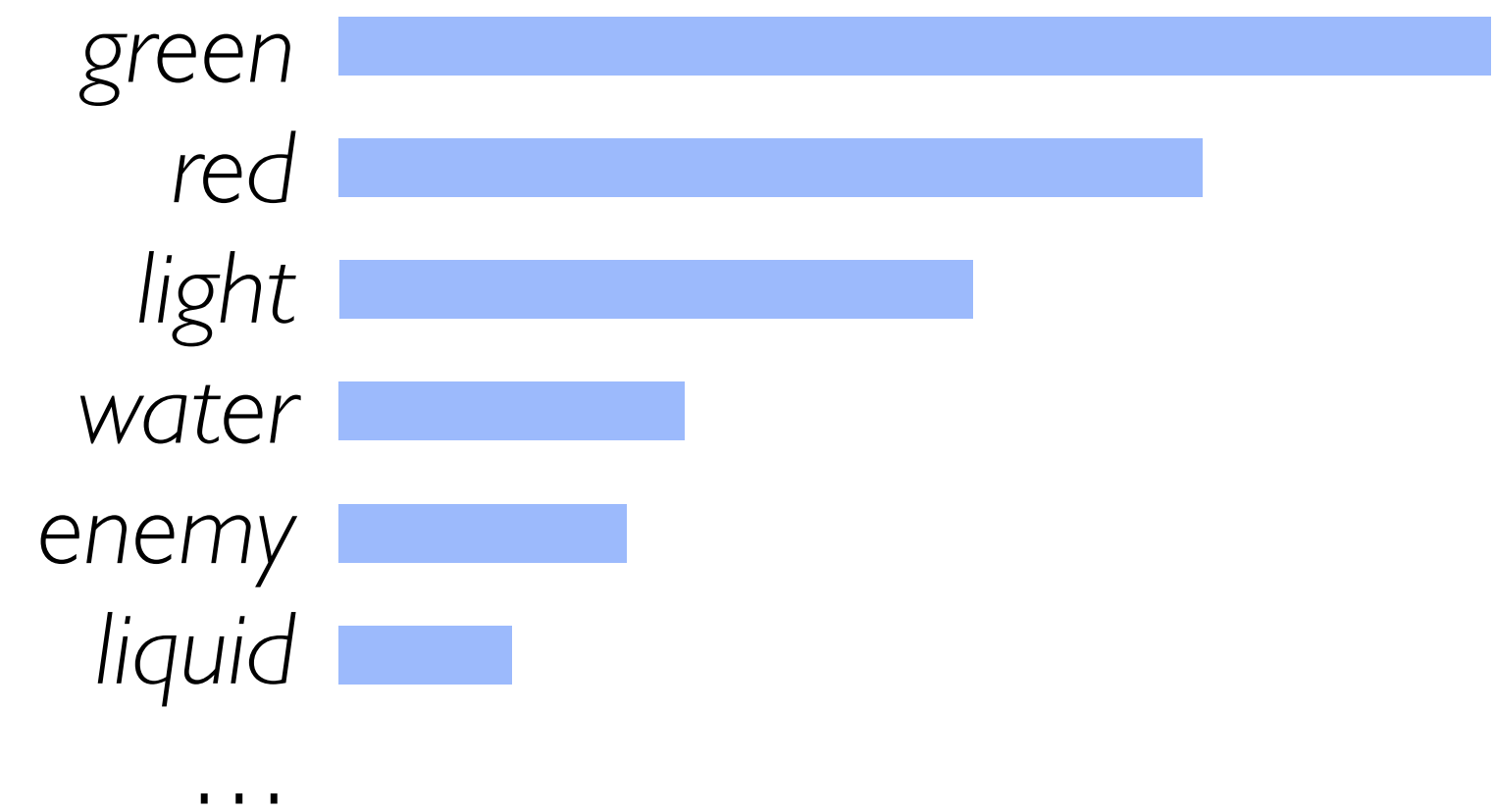Harry felt Greenback collapse against … on the floor as a jet of
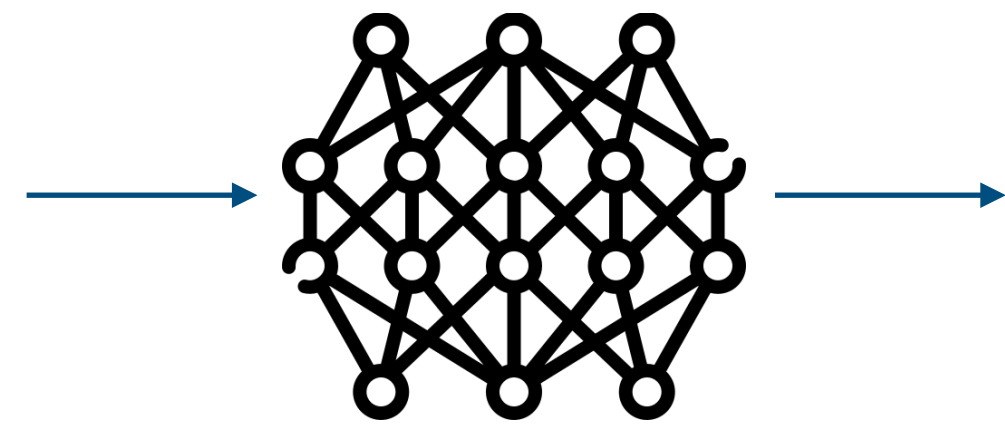
# Language Models

$$P(x_n \mid x_1, x_2, \cdots, x_{n-1})$$

# Language Models

$$P(x_n \mid x_1, x_2, \cdots, x_{n-1})$$

# Retrieval-based language models (LMs)

(also called semiparametric or nonparametric LMs)

# Retrieval-based language models (LMs)

## (also called semiparametric or nonparametric LMs)

(This figure assumes autoregressive LMs, but the idea can be broadly extended to masked LMs)

# Retrieval-based language models (LMs)

## (also called semiparametric or nonparametric LMs)

# Overview



## Why Retrieval-based LMs?

👤 Tell me about Meta Platform.

🔵 I don't have any information about a
**ChatGPT** company called Meta Platforms. It
is possible that the company is …

# Overview

# Overview

## Why Retrieval-based LMs?

Tell me about Meta Platform.

I don't have any information about a company called Meta Platforms. It is possible that the company is …

ChatGPT

## Retrieval Augmentation



$x \longrightarrow$ Retrieval $\longrightarrow$

LM $\longrightarrow y$

## New Retrieval-based LMs

$x \longrightarrow$ LM

… *"Avada Kedavra!" A jet of **green light** issued* …

… *move and a flash of **green light** and .*

… *just as a jet of **red light** blasted from Harry's*

… *is operated or driven by a jet of **water**.*

…

# Overview

## Why Retrieval-based LMs?

Tell me about Meta Platform.

I don't have any information about a company called Meta Platforms. It is possible that the company is …

*ChatGPT*

## Retrieval Augmentation



## New Retrieval-based LMs



… *"Avada Kedavra!" A jet of* **green light** *issued* …

… *move and a flash of* **green light** *and* …

… *just as a jet of* **red light** *blasted from Harry's*

… *is operated or driven by a jet of* **water**.

…

## Open Problems



Scaling **datastore** not just parameters?
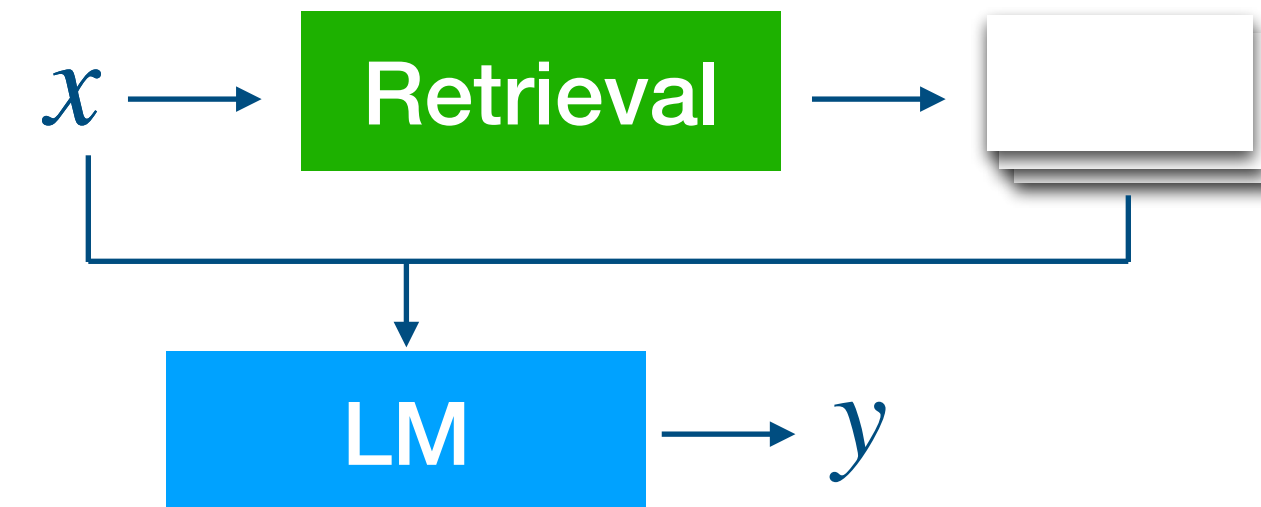
# Overview

## Why Retrieval-based LMs?

Tell me about Meta Platform.

I don't have any information about a company called Meta Platforms. It is possible that the company is …

**ChatGPT**

## Retrieval Augmentation

$x \longrightarrow$ **Retrieval** $\longrightarrow$

$\downarrow$

**LM** $\longrightarrow y$

## New Retrieval-based LMs

$x \longrightarrow$ **LM**

… *"Avada Kedavra!" A jet of* **green light** *issued* …

… *move and a flash of* **green light** *and* .

… *just as a jet of* **red light** *blasted from Harry's*

… *is operated or driven by a jet of* **water**.

…

## Open Problems

datastore

Scaling **datastore** not just parameters?

# Why Retrieval-based LMs?

New dimension in data
use & better at long-tail

Can grow & update w/o
additional training

Provide data attribution

Motivation

# Why Retrieval-based LMs?

**New dimension in data use & better at long-tail**

Can grow & update w/o additional training

Provide data attribution

# Why Retrieval-based LMs?

**New dimension in data use & better at long-tail**

Can grow & update w/o additional training

Provide data attribution



easy to memorize

modest-size LMs can't memorize

even large LMs can't memorize

Performance

Rarities of concepts/facts

# Why Retrieval-based LMs?

**New dimension in data use & better at long-tail**

Can grow & update w/o additional training

Provide data attribution

easy to memorize

modest-size LMs can't memorize

even large LMs can't memorize

Performance

Rarities of concepts/facts

Standard LMs: Need to remember everything

# Why Retrieval-based LMs?

**New dimension in data use & better at long-tail**

Can grow & update w/o additional training

Provide data attribution

easy to memorize

modest-size LMs can't memorize

even large LMs can't memorize

Performance

Rarities of concepts/facts

Standard LMs: Need to remember everything

datastore

Retrieval-based LMs: Can look-up anytime

# Why Retrieval-based LMs?

New dimension in data
use & better at long-tail

**Can grow & update
w/o additional training**

Provide data attribution



Tell me about Meta Platforms, Inc.

I'm sorry, I don't have information about a company called Meta Platforms, Inc. It's possible that the company is private or doesn't have a significant online presence. Can you provide more context or specify what information you're looking for?

# Why Retrieval-based LMs?

New dimension in data
use & better at long-tail

**Can grow & update
w/o additional training**

Provide data attribution

# Why Retrieval-based LMs?

New dimension in data use & better at long-tail

Can grow & update w/o additional training

**Provide data attribution**

# Why Retrieval-based LMs?

New dimension in data
use & better at long-tail

Can grow & update w/o
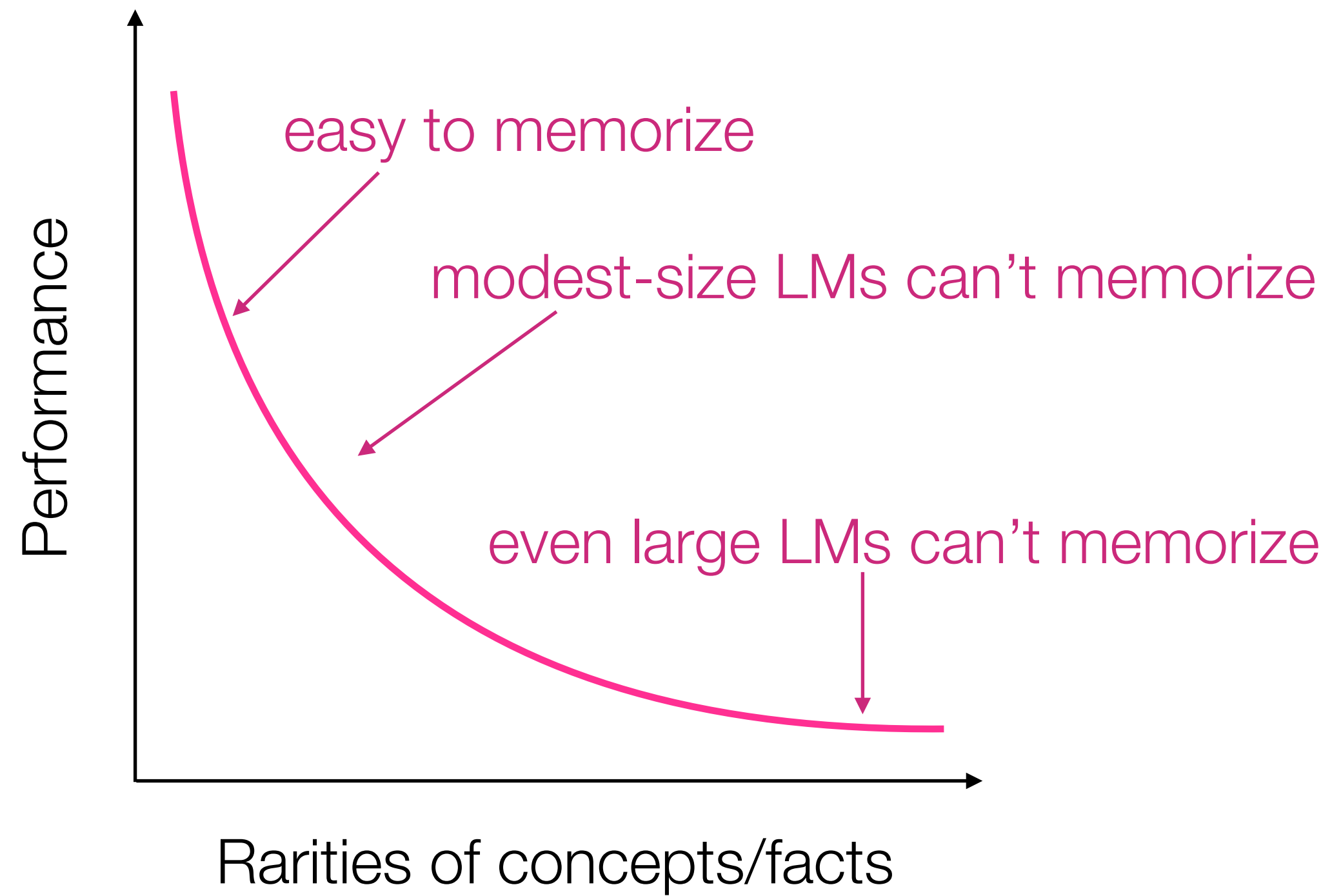additional training

**Provide data attribution**

👤 List 4 important papers authored by Geoffrey Hinton

# Why Retrieval-based LMs?

New dimension in data use & better at long-tail | Can grow & update w/o additional training | **Provide data attribution**

List 4 important papers authored by Geoffrey Hinton

**ChatGPT**

Geoffrey Hinton is a renowned computer scientist … Here are four 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,… ✖️

# Overview

## Why Retrieval-based LMs?

Tell me about Meta Platform.

I don't have any information about a company called Meta Platforms. It is possible that the company is …

ChatGPT

## Retrieval Augmentation

$x \longrightarrow$ **Retrieval** $\longrightarrow$

**LM** $\longrightarrow y$

## New Retrieval-based LMs

$x \longrightarrow$ **LM**

… *"Avada Kedavra!" A jet of* **green light** *issued* …

… *move and a flash of* **green light** *and* …

… *just as a jet of* **red light** *blasted from Harry's*

… *is operated or driven by a jet of* **water**.

…

## Open Problems

datastore

Scaling **datastore** not just parameters?

# Language Models (w/o retrieval)



| Harry felt Greenback collapse against him … on the floor as a jet of | → | **LM** |

*green*
*red*
*light*
*water*
*enemy*
*liquid*
…

# Language Models (w/ retrieval)

Harry felt Greenback collapse against him ... on the floor as a jet of

datastore ↔ **Retrieval Model**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**Most *relevant* text blocks**
(documents, passages, etc)

Guu et al. 2020. "REALM: Retrieval-Augmented Language Model Pre-Training"

# Language Models (w/ retrieval)

Harry felt Greenback collapse against him … on the floor as a jet of

datastore

**Retrieval Model**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**LM**

green
red
light
water
enemy
liquid
…

**Most *relevant* text blocks**
(documents, passages, etc)

Guu et al. 2020. "REALM: Retrieval-Augmented Language Model Pre-Training"

# Retrieval augmentation

Harry felt Greenback collapse against him … on the floor as a jet of

datastore ↔ **Retrieval Model**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**1) Retrieve stage**

◆ → **LM**

*green*
*red*
*light*
*water*
*enemy*
*liquid*
…

**2) Read (or Generate) stage**

# Retrieval augmentation: Overview

- Inference

- Training

- Key results

# Retrieval augmentation: Overview

- **Inference**

- Training

- Key results

# (1) Retrieve stage



datastore

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

Voldemort's want just as a jet of red light …

"The Boy Who Lived." He saw the mouth move and a flash of green …

# (1) Retrieve stage

datastore

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

Voldemort's want just as a jet of red light …

"The Boy Who Lived." He saw the mouth move and a flash of green …

**Encoder**

**Encoder**

**Encoder**

vector space

$$\mathbf{z} = \text{Encoder}(z)$$

# (1) Retrieve stage

datastore

$\boldsymbol{x}$ = Harry felt Greenback collapse… on the floor as a jet of

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

**Encoder**

**vector space**

Voldemort's want just as a jet of red light …

**Encoder**

"The Boy Who Lived." He saw the mouth move and a flash of green …

**Encoder**

$$\mathbf{z} = \text{Encoder}(z)$$

# (1) Retrieve stage



$x$ = Harry felt Greenback collapse… on the floor as a jet of

datastore

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

Voldemort's want just as a jet of red light …

"The Boy Who Lived." He saw the mouth move and a flash of green …

**Encoder**

**Encoder**

**Encoder**

**Encoder**

vector space

$$\mathbf{z} = \text{Encoder}(z)$$

$$\mathbf{x} = \text{Encoder}(x)$$

# (1) Retrieve stage

datastore

$\boldsymbol{x}$ = Harry felt Greenback collapse… on the floor as a jet of

**Encoder**

Fast nearest neighbor search

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

**Encoder**

vector space

Voldemort's want just as a jet of red light …

**Encoder**

"The Boy Who Lived." He saw the mouth move and a flash of green …

**Encoder**

$\mathbf{z} = \text{Encoder}(z)$

$\mathbf{x} = \text{Encoder}(x)$

# (1) Retrieve stage

datastore

$x$ = Harry felt Greenback collapse… on the floor as a jet of

**Encoder**

Fast nearest neighbor search

| Voldemort cried, "Avada Kedavra!" A jet of green light issued …from … |
|---|

**Encoder**

vector space

| Voldemort's want just as a jet of red light … |
|---|

**Encoder**

| "The Boy Who Lived." He saw the mouth move and a flash of green … |
|---|

**Encoder**

$$\mathbf{z} = \text{Encoder}(z)$$

$$\mathbf{x} = \text{Encoder}(x)$$

$$z = \text{argmax}_{z \in \mathcal{Z}} \left( \text{sim}(x, z) \right)$$

# (2) Read stage

**Retrieval results** *(ranked)*
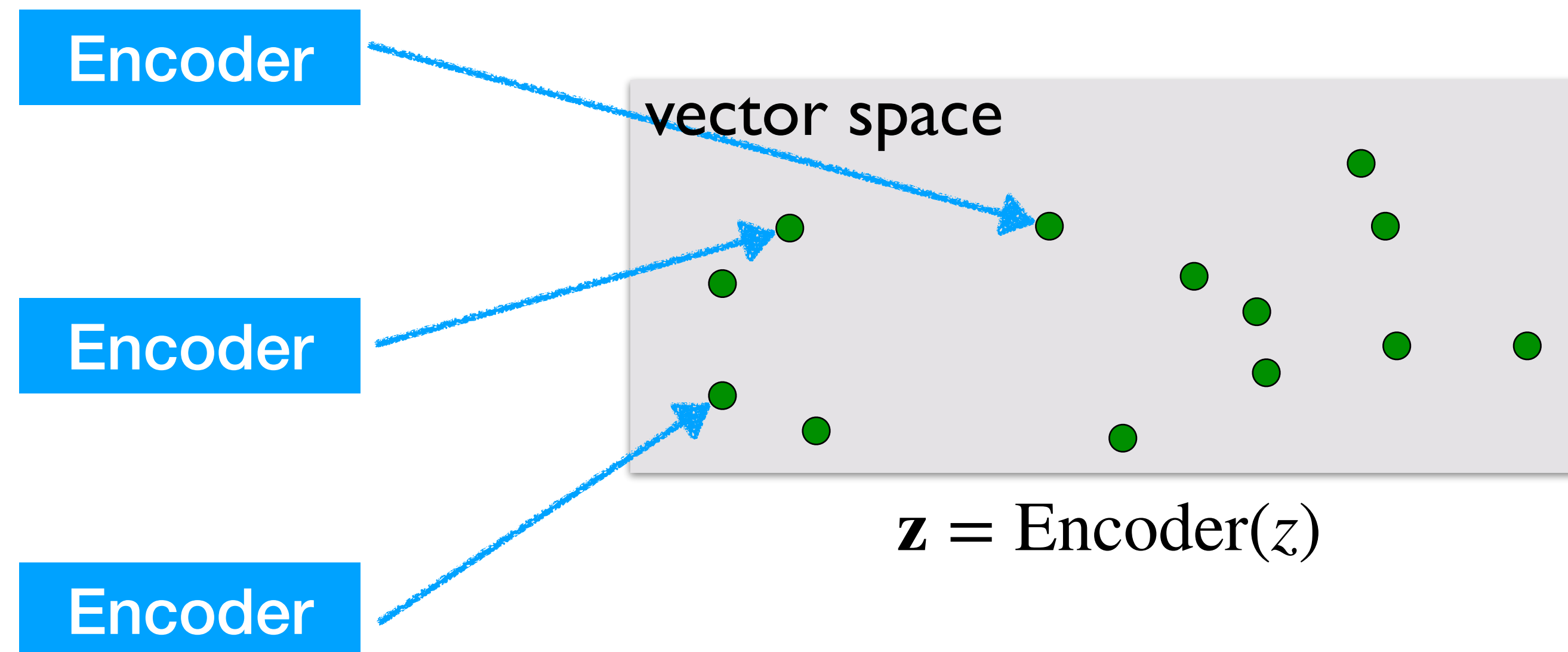
Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

Voldemort's want just as a jet of red light …

"The Boy Who Lived." He saw the mouth move and a flash of green …

# (2) Read stage

**Retrieval results** *(ranked)*

Voldemort cried, "Avada Kedavra!" A jet of green light issued from Voldemort's wand just as a jet of red light blasted from Harry's …
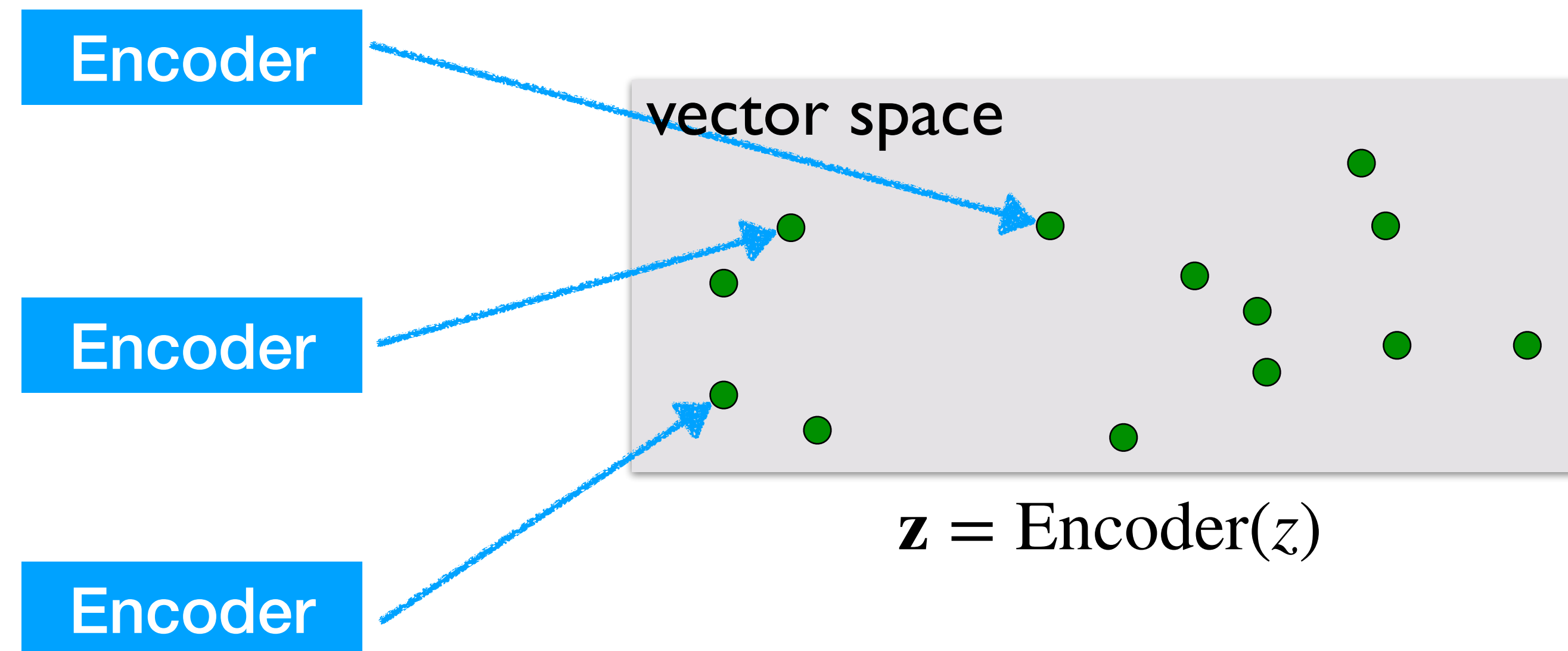
**+**

Harry felt Greenback collapse against him … a jet of

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

Voldemort's want just as a jet of red light …

"The Boy Who Lived." He saw the mouth move and a flash of green …

# (2) Read stage

**Retrieval results** *(ranked)*

Voldemort cried, "Avada Kedavra!" A jet of green light issued from Voldemort's wand just as a jet of red light blasted from Harry's …

**+**

| Voldemort cried, "Avada Kedavra!" A jet of green light issued …from … |
|---|

| Voldemort's want just as a jet of red light … |
|---|

| "The Boy Who Lived." He saw the mouth move and a flash of green … |
|---|

Harry felt Greenback collapse against him … a jet of

↓

| **LM** |
|---|

# (2) Read stage

**Retrieval results** *(ranked)*

Voldemort cried, "Avada Kedavra!" A jet of green light issued from Voldemort's wand just as a jet of red light blasted from Harry's …
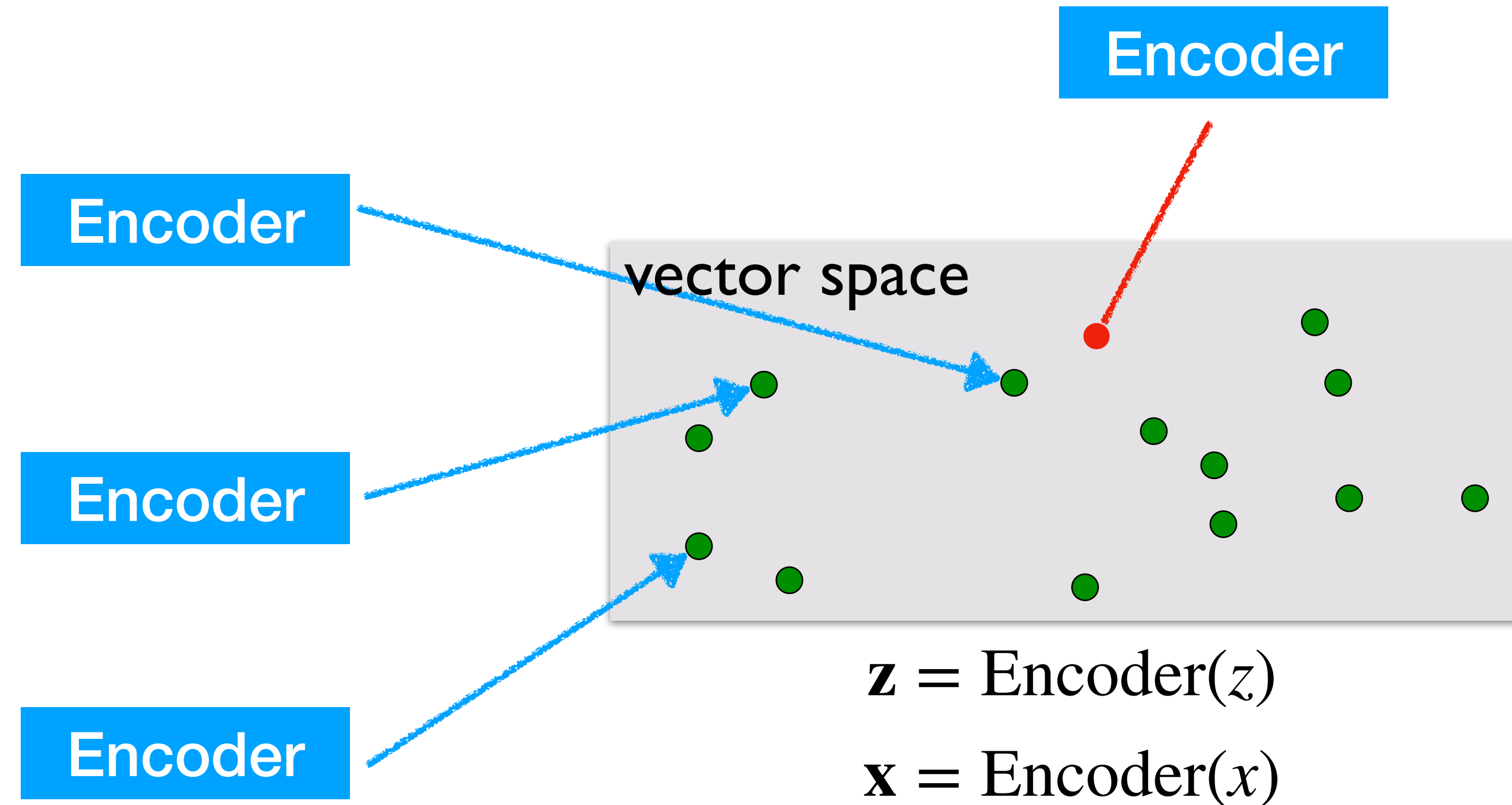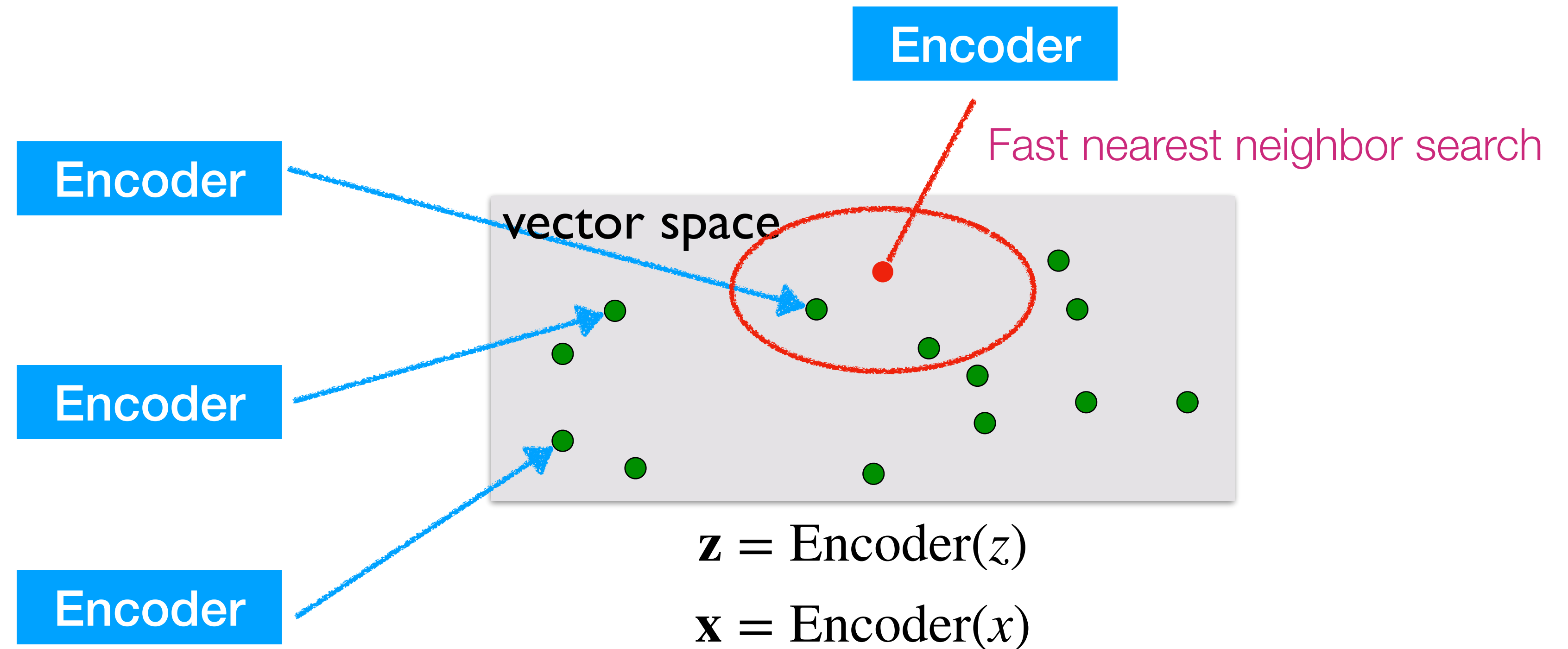
**+**

Harry felt Greenback collapse against him … a jet of

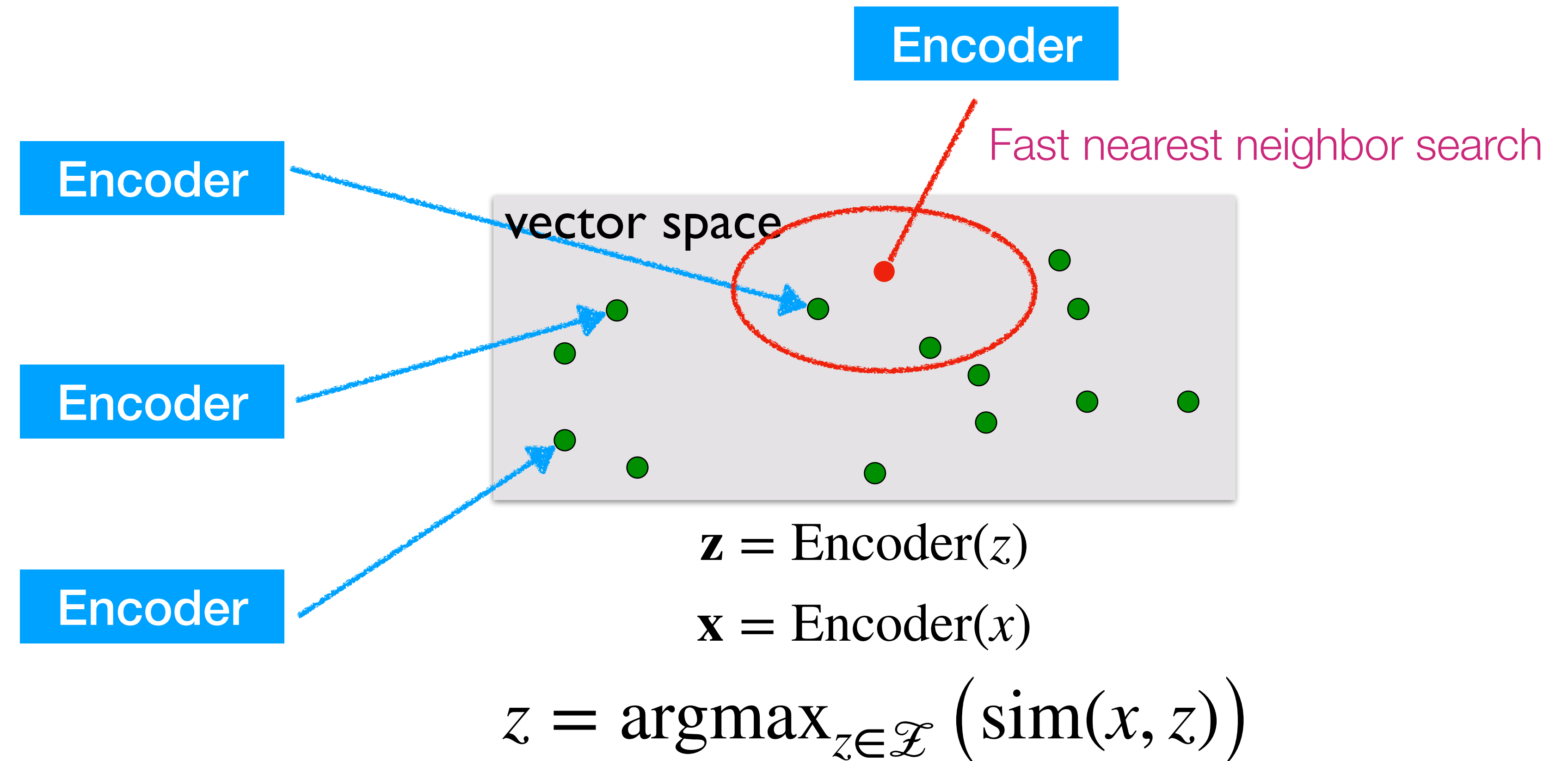Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

Voldemort's want just as a jet of red light …

"The Boy Who Lived." He saw the mouth move and a flash of green …

**LM**

*green*

*red*

*light*

*water*

*enemy*

*liquid*

…

# (2) Read stage

**Retrieval results** *(ranked)*

Voldemort cried, "Avada Kedavra!" A jet of green light issued from Voldemort's wand just as a jet of red light blasted from Harry's …

**+**

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

Harry felt Greenback collapse against him … a jet of

Voldemort's want just as a jet of red light …

**LM**

"The Boy Who Lived." He saw the mouth move and a flash of green …

*green*
*red*
*light*
*water*
*enemy*
*liquid*
…

**Very simple**
(You can use a black-box LM like an API!)

19

# (2) Read stage
## *How to use multiple text blocks?*

Voldemort cried, "Avada Kedavra!" A jet of green light issued from Voldemort's wand just as a jet of red light blasted from Harry's …

**+**

Harry felt Greenback collapse against him … a jet of

***Retrieval results*** *(ranked)*

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

Voldemort's want just as a jet of red light …

"The Boy Who Lived." He saw the mouth move and a flash of green …

**LM**

green
red
light
water
enemy
liquid
…

# (2) Read stage

*How to use multiple text blocks?*

Voldemort cried, "Avada Kedavra!" A jet of green light issued from Voldemort's wand just as a jet of red light blasted from Harry's …

**Retrieval results** *(ranked)*

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …
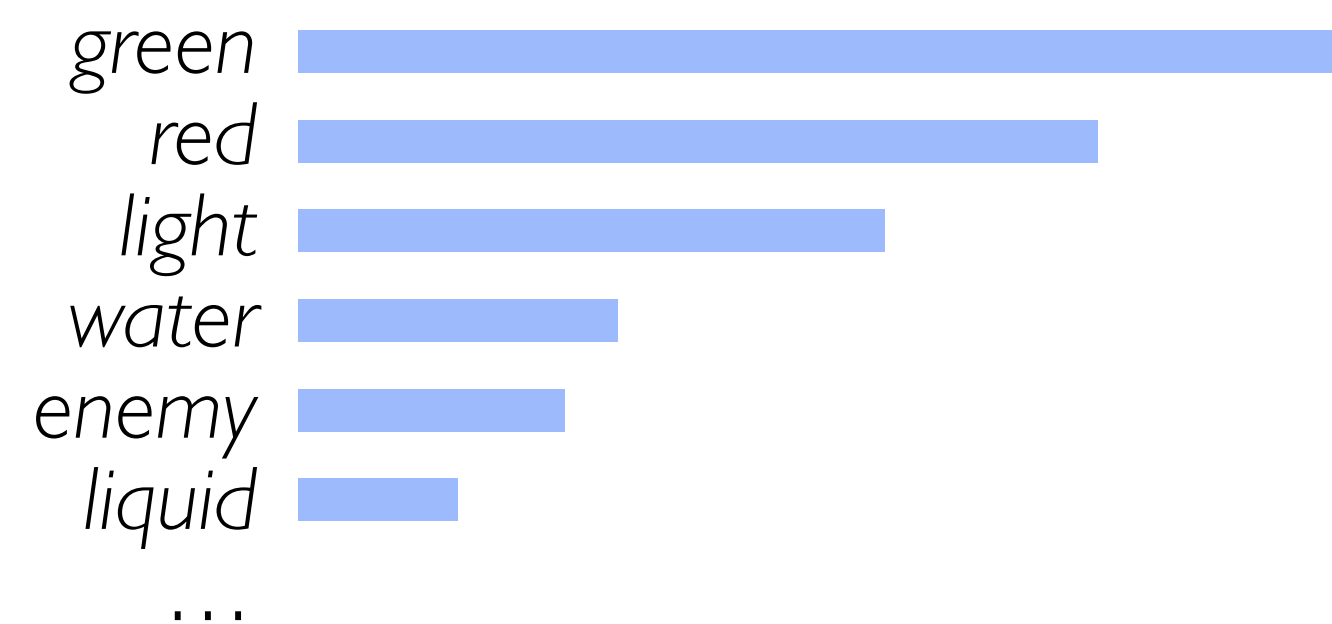
Voldemort's want just as a jet of red light …

"The Boy Who Lived." He saw the mouth move and a flash of green …

**+**

Harry felt Greenback collapse against him … a jet of

**LM**

green
red
light
water
enemy
liquid
…

20

# (2) Read stage

*How to use multiple text blocks?*
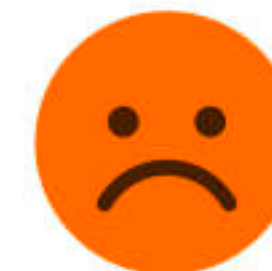
Voldemort's want just as a jet of red light …

**+**

Harry felt Greenback collapse against him … a jet of

**LM**

**Retrieval results** *(ranked)*

Voldemort's want just as a jet of red light …

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

"The Boy Who Lived." He saw the mouth move and a flash of green …

*green*
*red*
*light*
*water*
*enemy*
*liquid*
…

# (2) Read stage

*How to use multiple text blocks?* **1) Concatenation**

Voldemort's want just as a jet of red light …👎

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …👍

"The Boy Who Lived." He saw the mouth move and a flash of green …👍

# (2) Read stage

Voldemort's want just as a jet of red light … 👎

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from … 👍

"The Boy Who Lived." He saw the mouth move and a flash of green … 👍

**+**

Harry felt Greenback collapse against him … a jet of

# (2) Read stage

*How to use multiple text blocks?* **1) Concatenation**

Voldemort's want just as a jet of red light …
Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …
"The Boy Who Lived." He saw the mouth move and a flash of green …

**+**

Harry felt Greenback collapse against him … a jet of



**LM**

*green*
*red*
*light*
*water*
*enemy*
*liquid*
…

# (2) Read stage

Voldemort's want just as a jet of red light …👎
Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …👍
"The Boy Who Lived." He saw the mouth move and a flash of green …👍

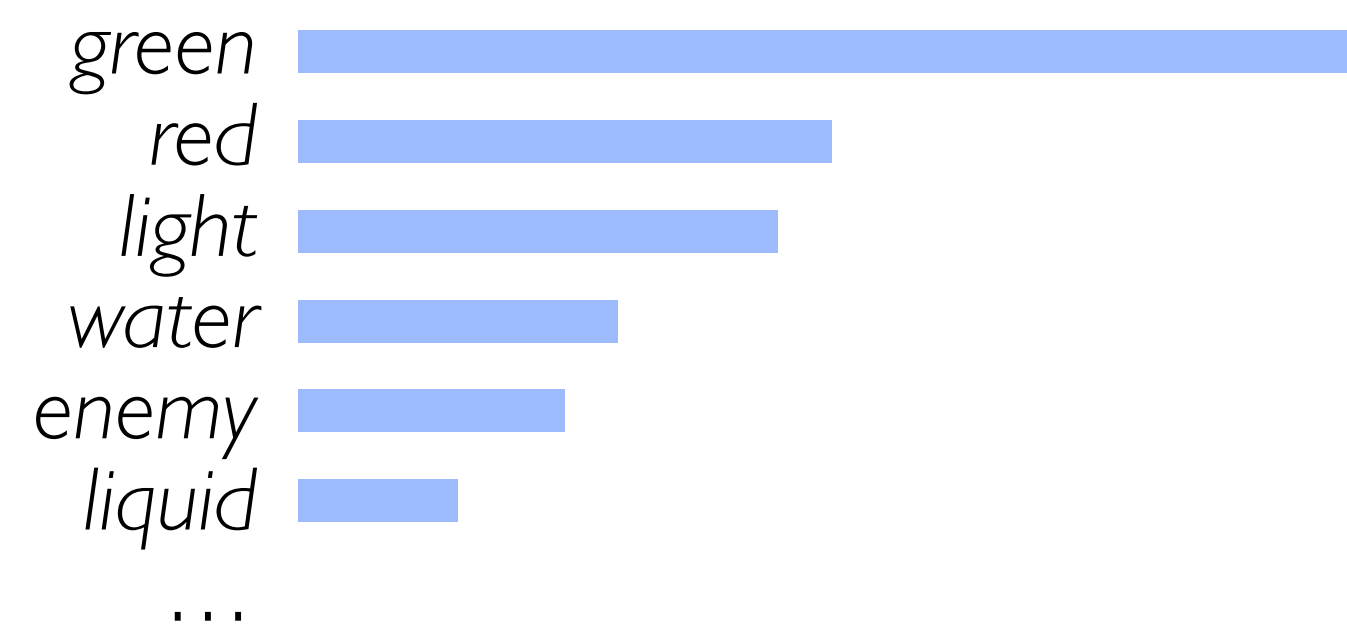**+**

Harry felt Greenback collapse against him … a jet of

↓

**LM**

↓

| | |
|---|---|
| 🙂 | Simple |
| ☹️ | Increase the inference cost & Bounded by the maximum length limit of the LM |

*green*
*red*
*light*
*water*
*enemy*
*liquid*
…

22

# (2) Read stage

*How to use multiple text blocks?* **2) Ensembling**

Voldemort's want just as a
jet of red light …
**+**
Harry felt Greenback collapse
against him … a jet of

Voldemort cried, "Avada
Kedavra!" A jet of green …
**+**
Harry felt Greenback collapse
against him … a jet of

… saw the mouth move
and a flash of green …
**+**
Harry felt Greenback collapse
against him … a jet of

Guu et al. 2020. "REALM: Retrieval-Augmented Language Model Pre-Training" * Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# (2) Read stage

*How to use multiple text blocks?* **2) Ensembling**

Voldemort's want just as a
jet of red light ... 👎

**+**

Harry felt Greenback collapse
against him ... a jet of

↓

**LM**

*green*
*red*
*light*
*water*

Voldemort cried, "Avada
Kedavra!" A jet of green ... 👍

**+**

Harry felt Greenback collapse
against him ... a jet of

↓

**LM**

*green*
*red*
*light*
*water*

... saw the mouth move
and a flash of green ... 👍

**+**

Harry felt Greenback collapse
against him ... a jet of

↓

**LM**

*green*
*red*
*light*
*water*

Guu et al. 2020. "REALM: Retrieval-Augmented Language Model Pre-Training" * Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# (2) Read stage

*How to use multiple text blocks?* **2) Ensembling**

Guu et al. 2020. "REALM: Retrieval-Augmented Language Model Pre-Training" * Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# (2) Read stage

*How to use multiple text blocks?* **2) Ensembling**

Voldemort's want just as a
jet of red light … 👎

**+**

Harry felt Greenback collapse
against him … a jet of

⬇

**LM**

*green*
*red*
*light*
*water*

Voldemort cried, "Avada
Kedavra!" A jet of green … 👍

**+**

Harry felt Greenback collapse
against him … a jet of

⬇

**LM**

*green*
*red*
*light*
*water*

… saw the mouth move
and a flash of green … 👍

**+**

Harry felt Greenback collapse
against him … a jet of

⬇

**LM**

*green*
*red*
*light*
*water*

$$P(y \mid x) = \sum_{z \in \mathcal{Z}} P_{\text{ret}}(z \mid x) P_{\text{LM}}(y \mid x, z)$$

retrieval score    LM score

*green*
*red*
*light*
*water*

Guu et al. 2020. "REALM: Retrieval-Augmented Language Model Pre-Training" * Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"
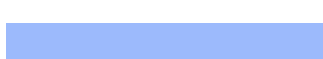
# (2) Read stage

## *How to use multiple text blocks?* **2) Ensembling**

Voldemort's want just as a jet of red light …  👎

Voldemort cried, "Avada Kedavra!" A jet of green …  👍

… saw the mouth move and a flash of green …  👍

**+**

**+**

**+**

Harry felt Greenback collapse against him … a jet of

Harry felt Greenback collapse against him … a jet of

Harry felt Greenback collapse against him … a jet of

**LM**

**LM**

**LM**

*green*
*red*
*light*
*water*

*green*
*red*
*light*
*water*

*green*
*red*
*light*
*water*

$$P(y \mid x) = \sum_{z \in \mathcal{Z}} P_{\text{ret}}(z \mid x) P_{\text{LM}}(y \mid x, z)$$

retrieval score          LM score

*green*
*red*
*light*
*water*

🙂 Not bounded by the length limit

☹️ Increase the inference cost

Guu et al. 2020. "REALM: Retrieval-Augmented Language Model Pre-Training" * Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# (2) Read stage

**Retrieval results** *(ranked)*

Voldemort's want just as a jet of red light …

👎

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

👍

"The Boy Who Lived." He saw the mouth move and a flash of green …

👍

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from …

👍

# (2) Read stage

*How to use multiple text blocks?* **3) Reranking**

**Retrieval results** *(ranked)*



Voldemort cried, "Avada Kedavra!" A jet of green light issued from Voldemort's wand just as a jet of red light blasted from Harry's ...

**+**

Harry felt Greenback collapse against him ... a jet of

Voldemort's want just as a jet of red light ...

Voldemort cried, "Avada Kedavra!" A jet of green light issued ...from ...

Voldemort cried, "Avada Kedavra!" A jet of green light issued ...from ...

"The Boy Who Lived." He saw the mouth move and a flash of green ...

**LM**

*green*
*red*
*light*
*water*
*enemy*
*liquid*
...

Ram et al. 2023. "In-Context Retrieval-Augmented Language Models"

# (2) Read stage

*How to use multiple text blocks?* **3) Reranking**

Voldemort cried, "Avada Kedavra!" A jet of green light issued from Voldemort's wand just as a jet of red light blasted from Harry's …

+

Harry felt Greenback collapse against him … a jet of

**Retrieval results** *(ranked)*

Voldemort's want just as a jet of red light … 👎

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from … 👍

Voldemort cried, "Avada Kedavra!" A jet of green light issued …from … 👍

"The Boy Who Lived." He saw the mouth move and a flash of green … 👍

**LM**

green

red

light

water

enemy

liquid

…

🙂 Not bounded by the length limit

☹️ Increase the inference cost

Ram et al. 2023. "In-Context Retrieval-Augmented Language Models"

# Key results

Ram et al. 2023. "In-Context Retrieval-Augmented Language Models"

# Key results

Perplexity: The lower the better



■ No Retrieval  ■ In-Context RALM (BM25)

| | OPT-125M | OPT-350M | OPT-1.3B | OPT-2.7B | OPT-6.7B | OPT-13B | OPT-30B | OPT-66B |
|---|---|---|---|---|---|---|---|---|
| No Retrieval | 17.4 | 14.5 | 10.4 | 9.4 | 8.4 | 8.0 | 7.5 | 7.2 |
| In-Context RALM (BM25) | 13.7 | 11.7 | 8.7 | 7.9 | 7.2 | 6.9 | 6.6 | 6.4 |

Varying sizes of LMs

# Key results



Perplexity: The lower the better

No Retrieval ■ In-Context RALM (BM25)

Varying sizes of LMs

Retrieval helps over all sizes of LMs

# Retrieval augmentation: Overview

- Inference

  - **Step 1: Retrieve**

  - **Step 2: Read (Generate)**

  - **Optionally, with multiple passages: Concatenation, Ensembling, Reranking**

- Training

- Key results

# Retrieval augmentation: Overview

- Inference

  - Step 1: Retrieve

  - Step 2: Read (Generate)

  - Optionally, with multiple passages: Concatenation, Ensembling, Reranking

- **Training**

- Key results

# How to train it?

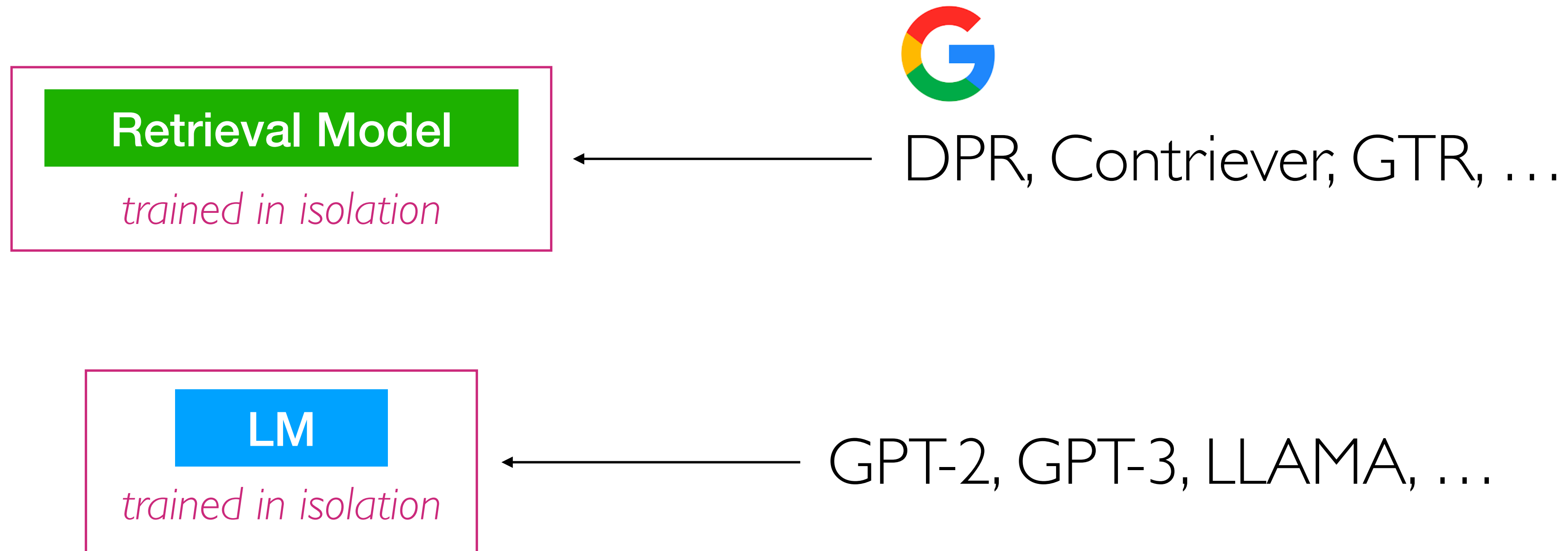**Retrieval Model**

*trained in isolation*

**LM**

*trained in isolation*

# How to train it?

Retrieval Model

*trained in isolation*

LM ← GPT-2, GPT-3, LLAMA, …

*trained in isolation*

# How to train it?

**Retrieval Model**

*trained in isolation*

DPR, Contriever, GTR, …

**LM**

*trained in isolation*

GPT-2, GPT-3, LLAMA, …

# How to train it?

Independent training

Retrieval Model

*trained in isolation*

LM

*trained in isolation*

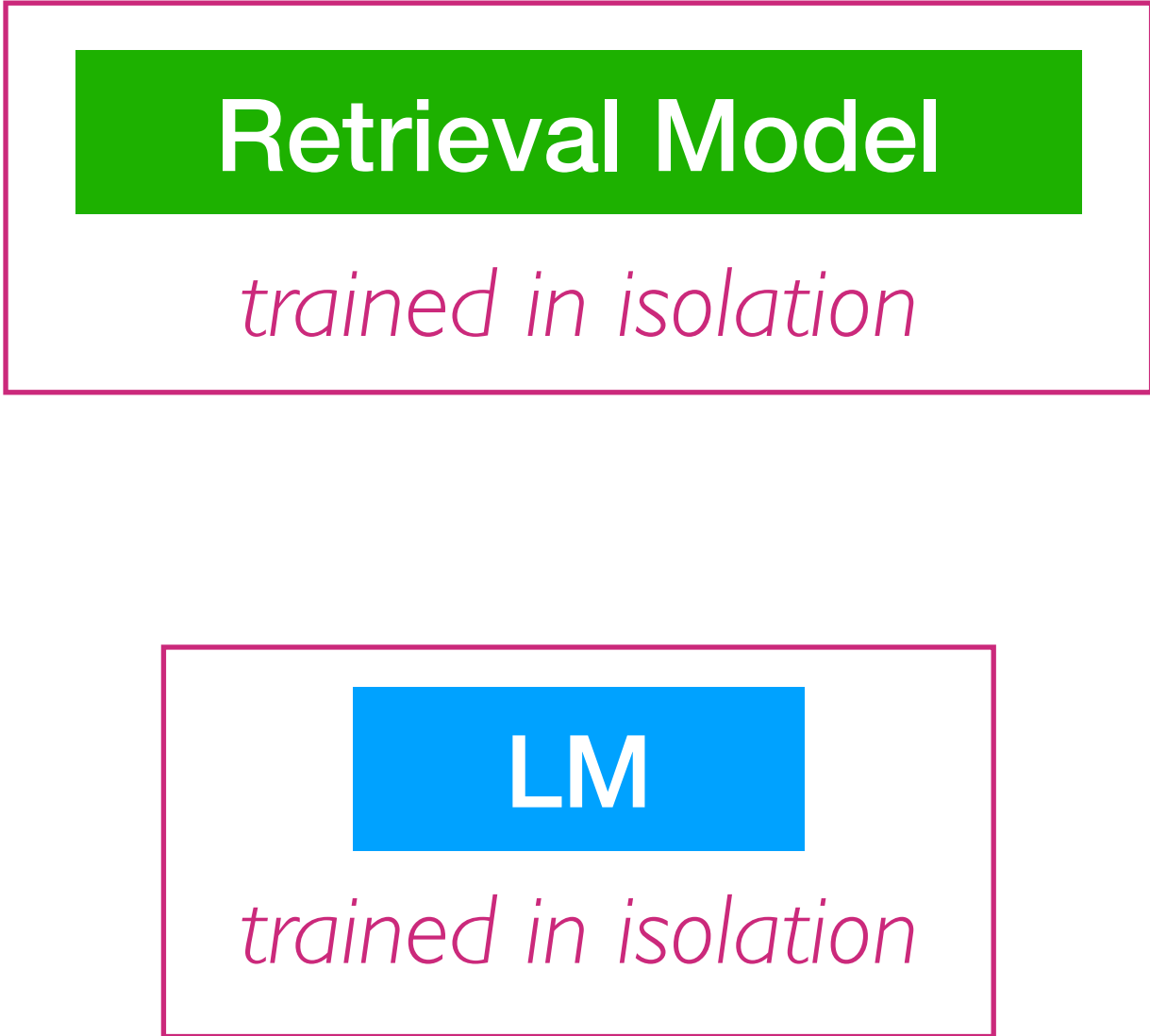# How to train it?

## Independent training

Retrieval Model
*trained in isolation*

LM
*trained in isolation*

## Joint training

Retrieval Model

LM

*trained jointly*

# How to train it?

## Independent training

**Retrieval Model**

*trained in isolation*

**LM**

*trained in isolation*

## Joint training

**Retrieval Model**

**LM**

*trained jointly*

## Sequential training

*trained in isolation*

**Retrieval Model**

**LM**

*trained conditionally*

# How to train it?

**Independent training**

Retrieval Model

*trained in isolation*

LM

*trained in isolation*

**Joint training**

Retrieval Model

LM

*trained jointly*

**Sequential training**

*trained in isolation*

Retrieval Model

LM

*trained conditionally*

or

*trained conditionally*

Retrieval Model

LM

*trained in isolation*

# How to train it?

## Independent training

**Retrieval Model**

*trained in isolation*

**LM**

*trained in isolation*

## Joint training
(Skipping details)

**Retrieval Model**

**LM**

*trained jointly*

## Sequential training

*trained in isolation*

**Retrieval Model**

↓

**LM**

*trained conditionally*

or

*trained conditionally*

**Retrieval Model**

↑

**LM**

*trained in isolation*

# Sequential training: freeze LM, tune retrieval

Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# Sequential training: freeze LM, tune retrieval

Harry felt Greenback collapse against him … on the floor as a jet of ⟶

**Frozen**

**LM**

*green*
*red*
*light*
*water*
*enemy*
*liquid*
…

# Sequential training: freeze LM, tune retrieval

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

**Frozen**

LM

LM

Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# Sequential training: freeze LM, tune retrieval

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

datastore ↔ **Retrieval Model**

**Frozen**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**LM**

Voldemort's want just as a jet of red light …

**LM**

Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# Sequential training: freeze LM, tune retrieval

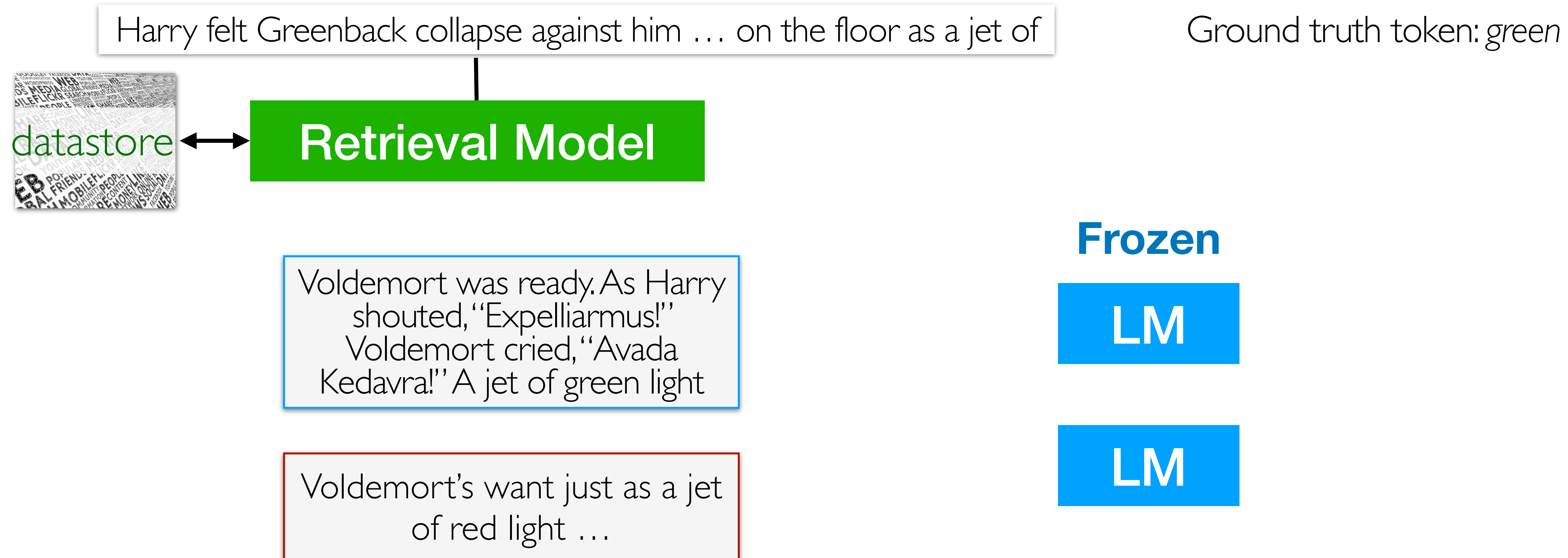Harry felt Greenback collapse against him … on the floor as a jet of
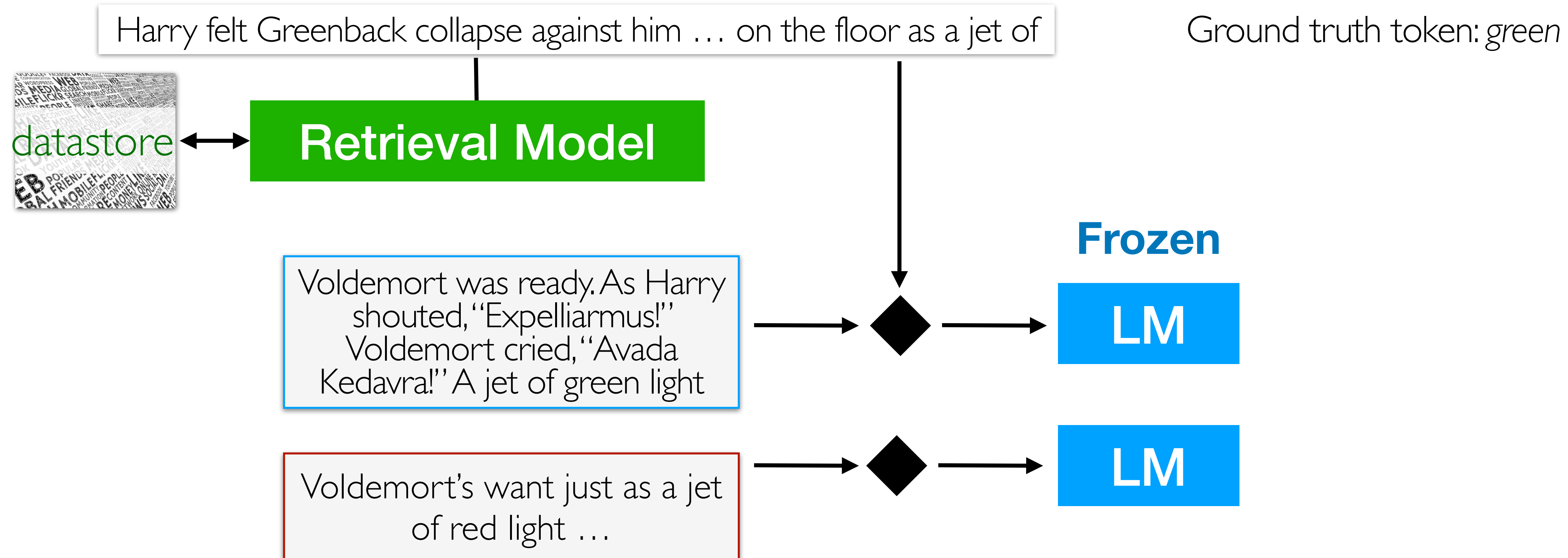
Ground truth token: *green*

datastore ↔ **Retrieval Model**

**Frozen**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

◆ → **LM**

Voldemort's want just as a jet of red light …

◆ → **LM**

# Sequential training: freeze LM, tune retrieval

Harry felt Greenback collapse against him … on the floor as a jet of
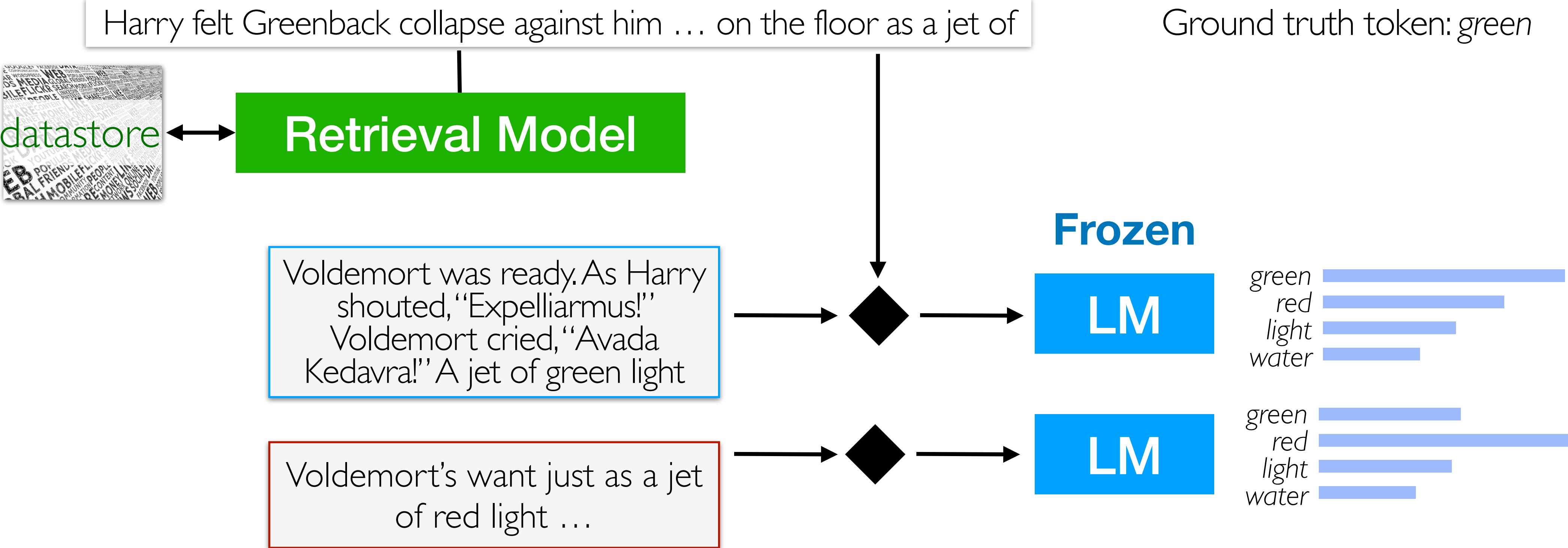
Ground truth token: *green*

datastore ↔ **Retrieval Model**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**Frozen**

LM

*green*
*red*
*light*
*water*

Voldemort's want just as a jet of red light …

LM

*green*
*red*
*light*
*water*

# Sequential training: freeze LM, tune retrieval

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

datastore ↔ **Retrieval Model**

**Frozen**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**LM**

*green*
*red*
*light*
*water*
👍

Voldemort's want just as a jet of red light …

**LM**

*green*
*red*
*light*
*water*
👎

# Sequential training: freeze LM, tune retrieval

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

datastore

**Retrieval Model**

**Frozen**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**LM**

green
red
light
water

Voldemort's want just as a jet of red light …

**LM**

green
red
light
water

Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# Sequential training: freeze LM, tune retrieval

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

datastore

**Retrieval Model** **Updated**

**Frozen**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

LM

*green*
*red*
*light*
*water*

Voldemort's want just as a jet of red light …

LM

*green*
*red*
*light*
*water*

# Sequential training: freeze LM, tune retrieval

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

datastore ⟷ **Retrieval Model** **Updated**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

Voldemort's want just as a jet of red light …

**Frozen**

LM

| | |
|---|---|
| *green* | |
| *red* | |
| *light* | |
| *water* | |

LM

| | |
|---|---|
| *green* | |
| *red* | |
| *light* | |
| *water* | |

**Updated**

$$\text{Maximize } P(y \mid x) = \sum_{z \in \mathcal{Z}} P_{\text{ret}}(z \mid x) P_{\text{LM}}(y \mid x, z)$$
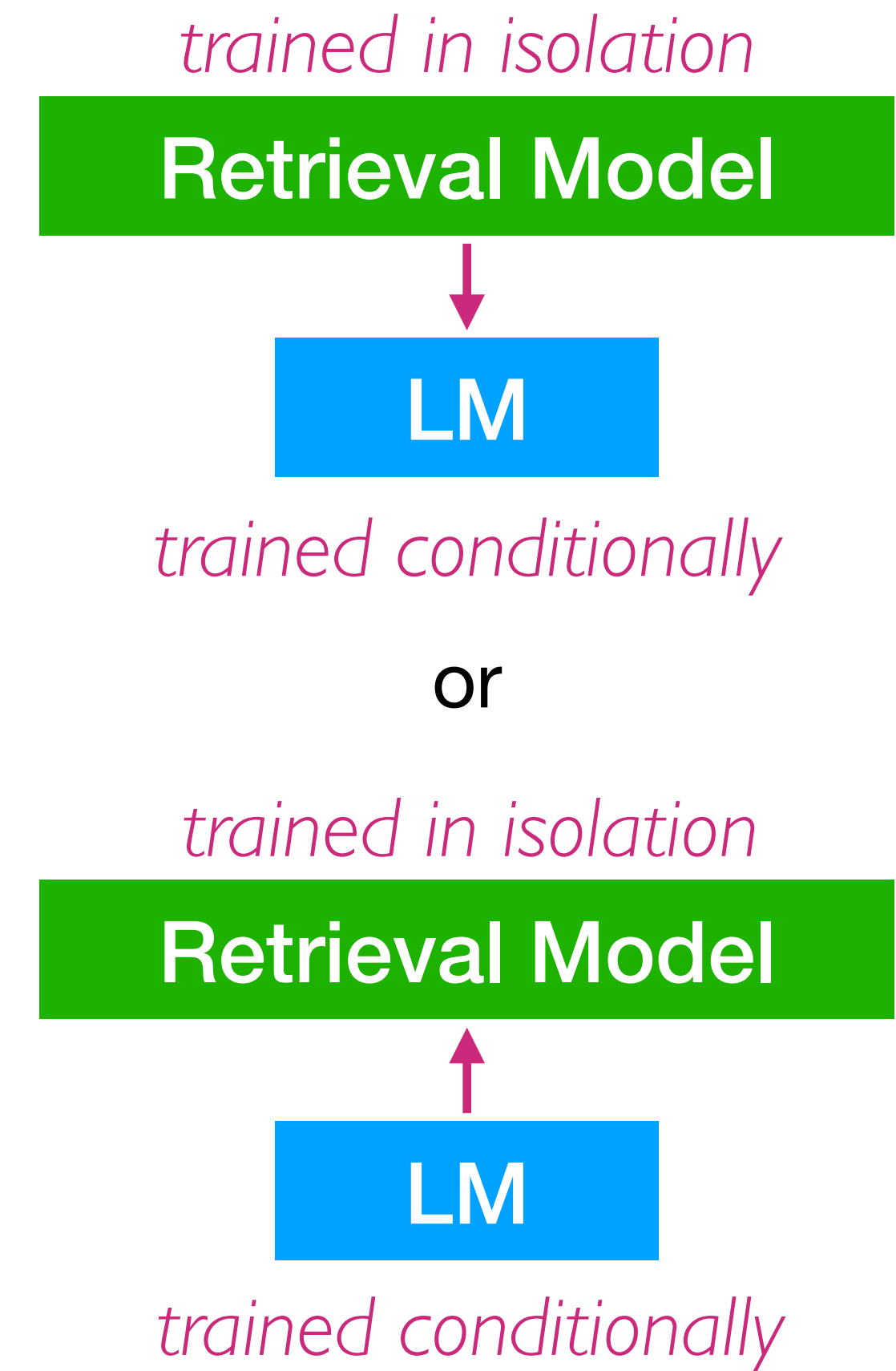
Shi et al. 2023. "REPLUG: Retrieval-Augmented Black-Box Language Models"

# Sequential training: freeze retrieval, tune LM

Shi et al. 2023. "In-Context Pretraining: Language Modeling Beyond Document Boundaries"

# Sequential training: freeze retrieval, tune LM

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

Shi et al. 2023. "In-Context Pretraining: Language Modeling Beyond Document Boundaries"

# Sequential training: freeze retrieval, tune LM

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

datastore ↔ **Retrieval Model** **Frozen**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

# Sequential training: freeze retrieval, tune LM



Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

datastore

**Retrieval Model** **Frozen**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**Updated**

LM

green
red
light
water
enemy
liquid
…

# Sequential training: freeze retrieval, tune LM

Harry felt Greenback collapse against him … on the floor as a jet of

Ground truth token: *green*

datastore ⟷ **Retrieval Model**  **Frozen**

Voldemort was ready. As Harry shouted, "Expelliarmus!" Voldemort cried, "Avada Kedavra!" A jet of green light

**Updated**

**LM**

*green*
*red*
*light*
*water*
*enemy*
*liquid*
…

**Updated**

$$\text{Maximize } P(y \mid x) = \sum_{z \in \mathcal{Z}} P_{\text{ret}}(z \mid x) \boxed{P_{\text{LM}}(y \mid x, z)}$$

# Summary: Training

**Independent training**

Retrieval Model

*trained in isolation*

LM

*trained in isolation*

**Joint training**
(Skipping details)

Retrieval Model

LM

*trained jointly*

**Sequential training**

*trained in isolation*

Retrieval Model

↓

LM

*trained conditionally*

or

*trained in isolation*

Retrieval Model

↑

LM

*trained conditionally*

# Summary: Training

## Independent training

Retrieval Model

*trained in isolation*

LM

*trained in isolation*

## Joint training
(Skipping details)

Retrieval Model

LM

*trained jointly*

Quite difficult, essentially iterative sequential training

## Sequential training

*trained in isolation*

Retrieval Model

↓

LM

*trained conditionally*

or

*trained in isolation*

Retrieval Model

↑

LM

*trained conditionally*

Gus et al. 2020. "REALM: Retrieval-Augmented Language Model Pre-Training"
Izcard et al. 2022. "Atlas: Few-shot Learning with Retrieval Augmented Language Models"

# Summary: Training

## Independent training

Retrieval Model

*trained in isolation*

LM

*trained in isolation*

**Good enough if you want minimal effort**

## Joint training
(Skipping details)

Retrieval Model

LM

*trained jointly*

**Principle way but still open question**

## Sequential training

*trained in isolation*

Retrieval Model

↓

LM

*trained conditionally*

or

*trained in isolation*

Retrieval Model

↑

LM

*trained conditionally*

**Good middle ground**

# Retrieval augmentation: Overview

- Inference

    - Step 1: Retrieve

    - Step 2: Read (Generate)

    - Optionally, with multiple passages: Concatenation, Ensembling, Reranking

- Training

    - **Independent training, Joint training, Sequential training**

- Key results

# Retrieval augmentation: Overview

- Inference
  - Step 1: Retrieve
  - Step 2: Read (Generate)
  - Optionally, with multiple passages: Concatenation, Ensembling, Reranking
- Training
  - Independent training, Joint training, Sequential training
- **Key results**

# Question Answering

# Question Answering

Izcard et al. "Atlas: Few-shot Learning with Retrieval Augmented Language Models"

# Question Answering



ATLAS largely outperforms 7x larger LMs in few-shot

- Chinchilla (70B)
- ATLAS (Few; 11B)
- ATLAS (Full; 11B)

# Question Answering



Full-shot fine-tuning further improves performance

Legend:
- Chinchilla (70B) — blue
- ATLAS (Few; 11B) — green
- ATLAS (Full; 11B) — orange with pink border

Izcard et al. "Atlas: Few-shot Learning with Retrieval Augmented Language Models"

# Question Answering

What is Kathy Saltzman's occupation?



Mallen et al. 2023. "When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories"

# Question Answering

What is Kathy Saltzman's occupation?



Gains increase as the rarity increases (even over GPT-3!)

Mallen et al. 2023. "When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories"

# Reasoning (MMLU)

# Reasoning (MMLU)



Large performance gain from base LM

- Base LM (CodeX)
- + REPLUG LSR

# Code generation

TLDR (NL —> bash)

BLEU



- CodeT5
- + DocPrompting
- CodeX
- + DocPrompting

Zhou et al. 2023. "DocPrompting: Generating Code by Retrieving the Docs"

# Code generation

## TLDR (NL —> bash)

Large gains over both CodeT5 & CodeX



BLEU

- CodeT5
- + DocPrompting
- CodeX
- + DocPrompting

Zhou et al. 2023. "DocPrompting: Generating Code by Retrieving the Docs"

# Can update effectively

# Can update effectively

Izcard et al. "Atlas: Few-shot Learning with Retrieval Augmented Language Models"

# Can update effectively



Test- 2017    Test - 2020

T5-Fine-tuned on 2017 data

Train 2017, DS 2017
Train 2017, DS 2020

# Can update effectively



T5-Fine-tuned on 2017 data

Legend: Test- 2017, Test - 2020

Legend: Train 2017, DS 2017; Train 2017, DS 2020

Swapping test datastore only gives strong performance

Izcard et al. "Atlas: Few-shot Learning with Retrieval Augmented Language Models"

# Instruction-tuning

# Instruction-tuning



**Retriever Fine-tuning**

$p_R(c_1 \mid x)$   min $KL$

**Retriever** → $c_1$

$p_{LSR}(c_1 \mid x, y)$

→ $c_2$

$p_R(c_2 \mid x)$   min $KL$

$p_{LSR}(c_2 \mid x, y)$

$x$: Why does my GPU keep failing with Xid 79 fallen off the bus?

**Retrieval-augmented Instruction Tuning**

1 **Background:** I assume that the BGA chip has damage to the substrate level $\cdots$ \n\n**Q:** Why does my GPU keep failing with Xid 79 fallen off the bus? **A:** ✓

2 **Background:** Microsoft should withdraw from the hardware market $\cdots$\n\n **Question:** Why does my GPU keep failing with Xid 79 fallen off the bus? **Answer:** ✗

**LLM**

max $p_{LM}(y \mid c_1 \circ x)$

$y$: Do not use Y-splitters or 6-pin to 8-pin converters in any of the PCIe$\cdots$

max $p_{LM}(y \mid c_2 \circ x)$

$y$: Do not use Y-splitters or 6-pin to 8-pin converters in any of the PCIe$\cdots$

# Retrieval augmentation: Summary

- Inference

    - Step 1: Retrieve

    - Step 2: Read (Generate)

    - Optionally, with multiple passages: Concatenation, Ensembling, Reranking

- Training

    - Independent training, Joint training, Sequential training

- Key results

    - QA, Reasoning, Code generation, etc

    - Update effectively

    - You can also do instruction tuning

# Retrieval Augmented LMs are already being used!

# Retrieval Augmented LMs are already being used!

Liu et al. 2023. "Evaluating Verifiability in Generative Search Engines"
Min et al. 2023. "FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation"

# Retrieval Augmented LMs are already being used!



Toronto is a vibrant and diverse city with many attractions and activities to suit different tastes and preferences. Here are some suggestions for creating an itinerary for exploring Toronto over two days:

- **Day 1**: Start your day with a visit to the **CN Tower**, the iconic landmark that offers panoramic views of the city and beyond [1]. You can also experience the thrill of walking on the glass floor or the EdgeWalk, the world's highest hands-free walk [1]. Next, head to the nearby **St. Lawrence Market**, one of the world's best food markets, where you can sample a variety of cuisines and local specialties [2]. After lunch, take a stroll along **Queen West**, a trendy neighborhood with eclectic shops, galleries, cafes and street art [3]. In the evening, enjoy a show at one of the many theaters or comedy clubs in the **Entertainment District**, or catch a game at the **Scotiabank Arena** if you're a sports fan.

- **Day 2**: Explore the history and culture of Toronto at the **Royal Ontario Museum**, one of the largest and most comprehensive museums in North America, featuring exhibits on art, natural history, world cultures and more [4]. Then, hop on a ferry to the **Toronto Islands**, a group of islands that offer a relaxing escape from the city, with beaches, parks, trails and amusement rides [3] [5]. You can also rent a bike or kayak to explore the islands at your own pace. For dinner, head to **Chinatown**, one of the largest and most vibrant in North America, where you can find a variety of Asian cuisines and shops [3].
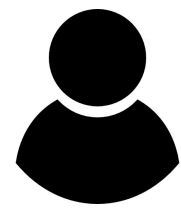
I hope this helps you plan your trip to Toronto. Have fun! 😊

Learn more:

1. cntower.ca    2. travel.usnews.com    3. bing.com

4. rom.on.ca    5. tripadvisor.com

Liu et al. 2023. "Evaluating Verifiability in Generative Search Engines"
Min et al. 2023. "FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation"

# Overview

## Why Retrieval-based LMs?
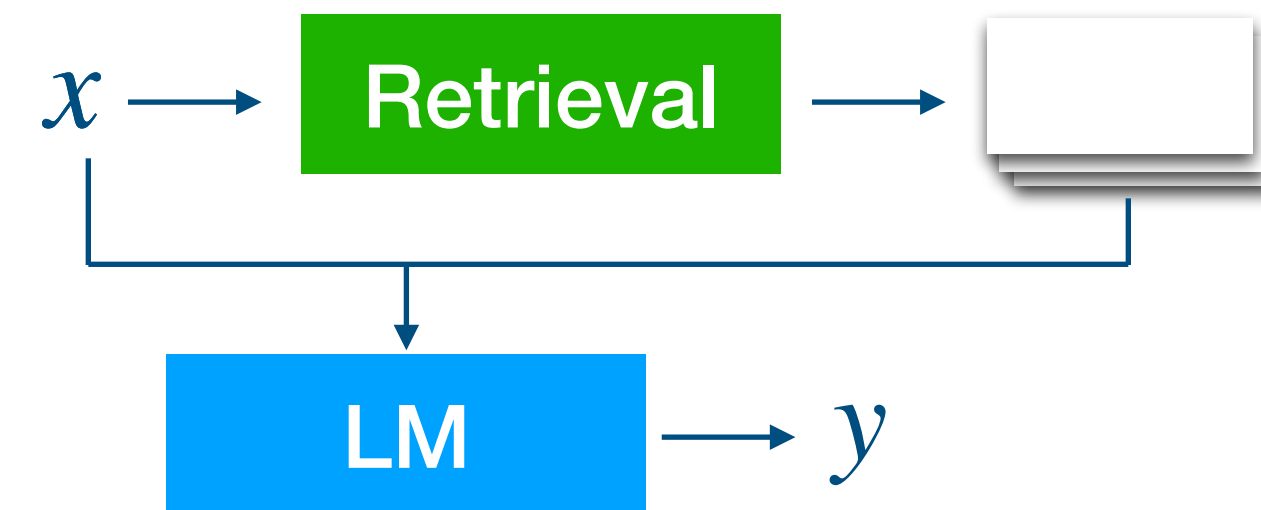
Tell me about Meta Platform.

I don't have any information about a company called Meta Platforms. It is possible that the company is …

**ChatGPT**

## Retrieval Augmentation

$x \longrightarrow$ **Retrieval** $\longrightarrow$

**LM** $\longrightarrow y$

## New Retrieval-based LMs

$x \longrightarrow$ **LM**

… *"Avada Kedavra!" A jet of* **green light** *issued* …

… *move and a flash of* **green light** *and .*

… *just as a jet of* **red light** *blasted from Harry's*

… *is operated or driven by a jet of* **water**.

…

## Open Problems

datastore

Scaling **datastore** not just parameters?

# New Retrieval-based LMs

- New Methodology 1 — Designing a new Transformer

- New Methodology 2 — Designing a new Softmax

- New LM Design — Mitigating fairness & legality issues

# New Retrieval-based LMs

> 1. How to overcome sequence length limit issue?
>
> 2. How to overcome efficiency issue when retrieving **many** blocks, **frequently**?

- **New Methodology 1 — Designing a new Transformer**

- New Methodology 2 — Designing a new Softmax

- New LM Design — Mitigating fairness & legality issues

# RETRO (Borgeaud et al. 2021)

# RETRO (Borgeaud et al. 2021)

New Transformers layers, designed to read *many* text blocks, *frequently*, more *efficiently*

# RETRO (Borgeaud et al. 2021)

*x* = World Cup 2022 was the last with 32 teams, before the increase to
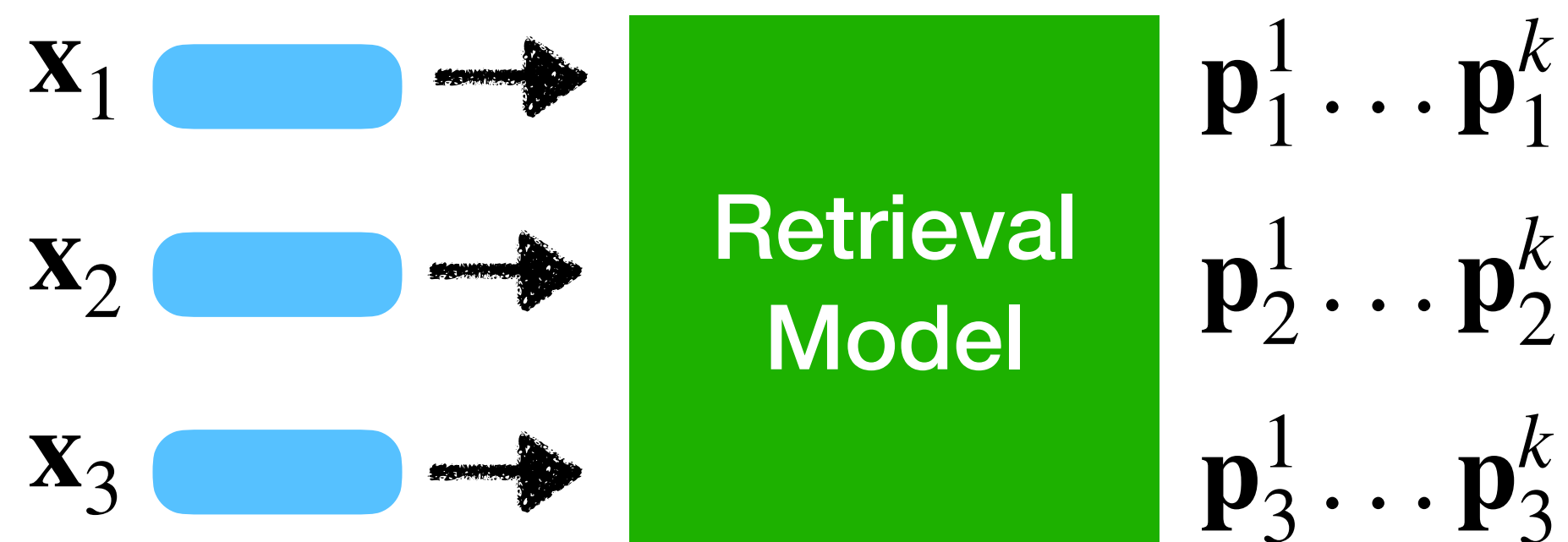
# RETRO (Borgeaud et al. 2021)

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

$$\mathbf{x}_1 \qquad\qquad \mathbf{x}_2 \qquad\qquad \mathbf{x}_3$$

# RETRO (Borgeaud et al. 2021)

$x$ = World Cup 2022 was the last with 32 teams, before the increase to

$$\mathbf{x}_1 \qquad\qquad \mathbf{x}_2 \qquad\qquad \mathbf{x}_3$$

($k$ text blocks per split)



$\mathbf{x}_1$ → Retrieval Model → $\mathbf{p}_1^1 \cdots \mathbf{p}_1^k$

$\mathbf{x}_2$ → $\mathbf{p}_2^1 \cdots \mathbf{p}_2^k$
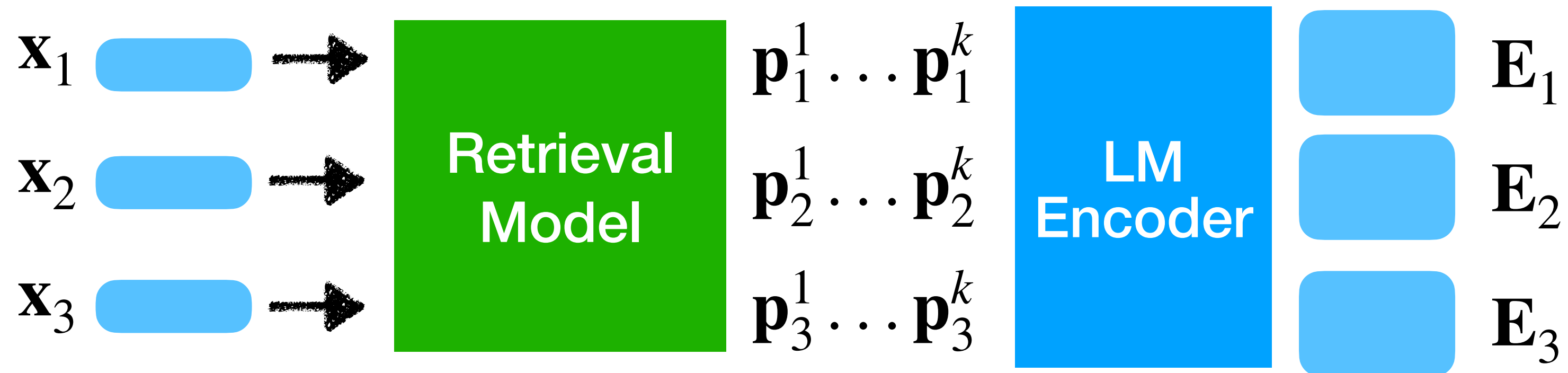
$\mathbf{x}_3$ → $\mathbf{p}_3^1 \cdots \mathbf{p}_3^k$

# RETRO (Borgeaud et al. 2021)

$x$ = World Cup 2022 was the last with 32 teams, before the increase to

$$\mathbf{x}_1 \qquad\qquad \mathbf{x}_2 \qquad\qquad \mathbf{x}_3$$

($k$ text blocks per split)



$$\mathbf{x}_1 \rightarrow \boxed{\substack{\text{Retrieval} \\ \text{Model}}} \begin{array}{l} \mathbf{p}_1^1 \cdots \mathbf{p}_1^k \\ \mathbf{p}_2^1 \cdots \mathbf{p}_2^k \\ \mathbf{p}_3^1 \cdots \mathbf{p}_3^k \end{array} \boxed{\substack{\text{LM} \\ \text{Encoder}}}$$

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"
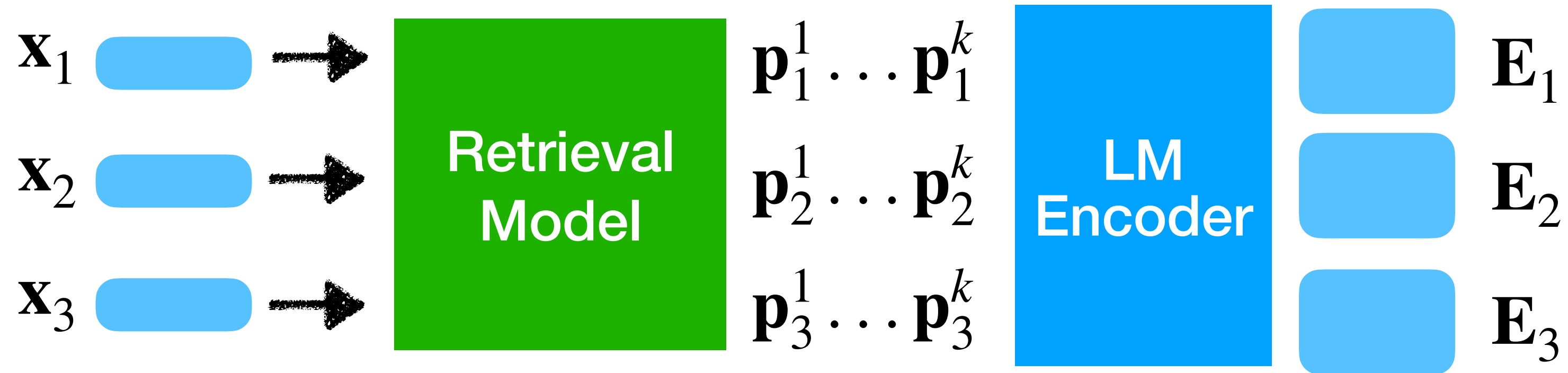
# RETRO (Borgeaud et al. 2021)

$x$ = World Cup 2022 was the last with 32 teams, before the increase to

$\mathbf{x}_1$            $\mathbf{x}_2$            $\mathbf{x}_3$
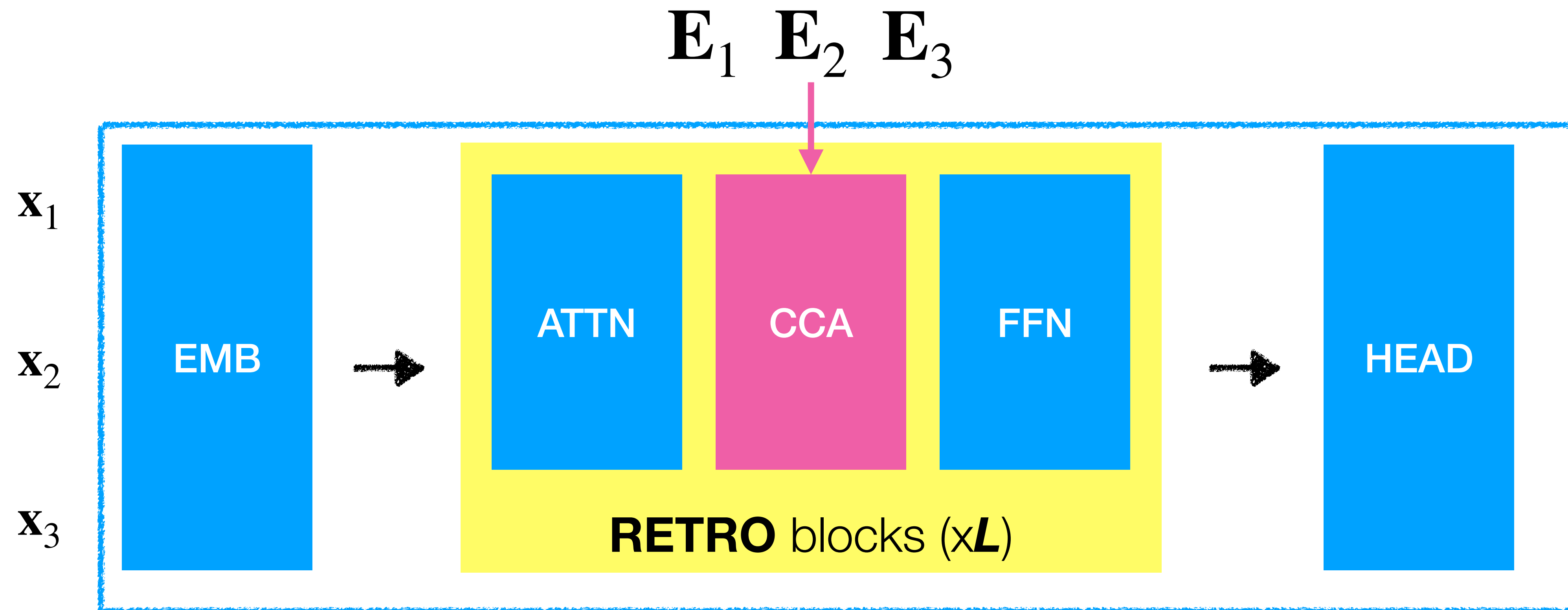
($k$ text blocks per split)



$\mathbf{x}_1$ → Retrieval Model → $\mathbf{p}_1^1 \cdots \mathbf{p}_1^k$ → LM Encoder → $\mathbf{E}_1$

$\mathbf{x}_2$ → $\mathbf{p}_2^1 \cdots \mathbf{p}_2^k$ → $\mathbf{E}_2$

$\mathbf{x}_3$ → $\mathbf{p}_3^1 \cdots \mathbf{p}_3^k$ → $\mathbf{E}_3$

# RETRO (Borgeaud et al. 2021)

$x$ = World Cup 2022 was the last with 32 teams, before the increase to

$$\mathbf{x}_1 \qquad\qquad \mathbf{x}_2 \qquad\qquad \mathbf{x}_3$$

($k$ text blocks per split)



How to incorporate them into Transformers?

 Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# Regular Transformers



$\mathbf{x}_1$

$\mathbf{x}_2$   EMB   →   ATTN   FFN   →   HEAD

$\mathbf{x}_3$   Transformers blocks (x**L**)

# RETRO Transformers

$$\mathbf{E}_1 \quad \mathbf{E}_2 \quad \mathbf{E}_3$$



**Chunked Cross Attention (CCA)**

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# Chunked Cross Attention

# Chunked Cross Attention



Outputs from the previous layer  H

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# Chunked Cross Attention

# Chunked Cross Attention

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"
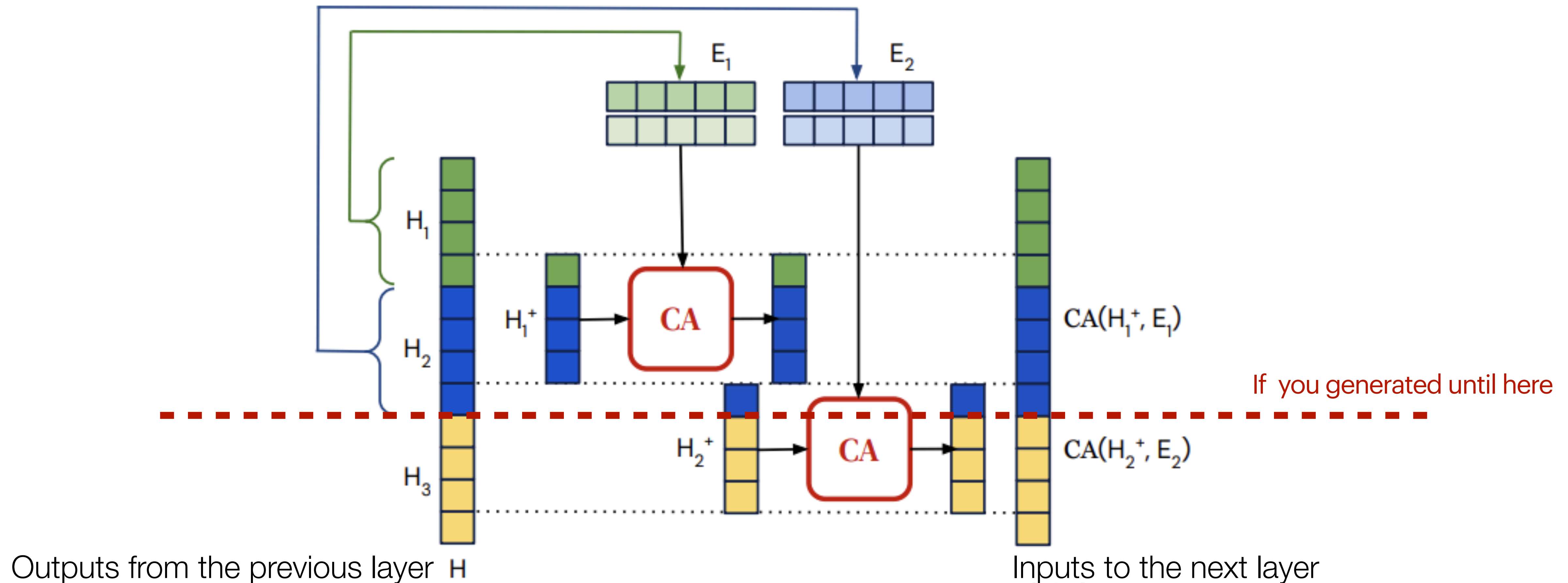
# Chunked Cross Attention



Outxuts from the previous layer  H

Inputs to the next layer

✓ Cross-attention can be computed *in parallel, and be re-used*

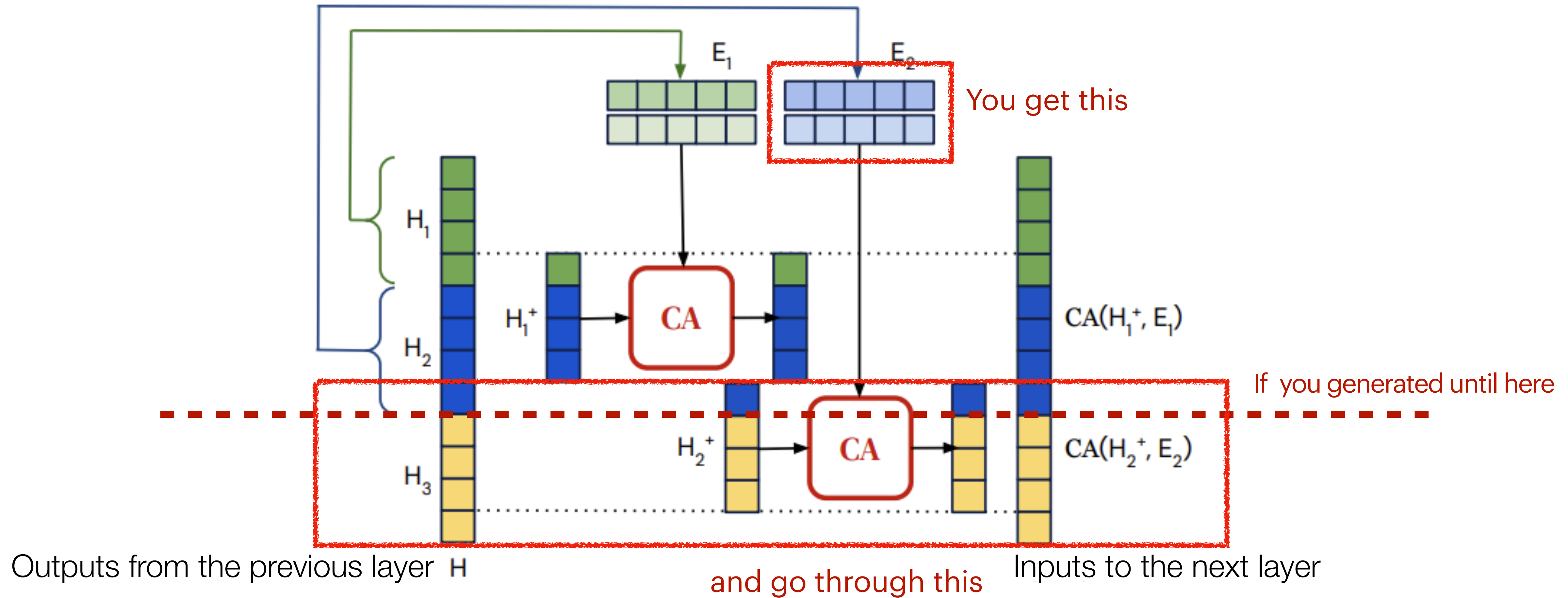Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# Chunked Cross Attention



✓ Cross-attention can be computed *in parallel, and be re-used*

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# Chunked Cross Attention



You get this

If you generated until here

$CA(H_1^+, E_1)$

$CA(H_2^+, E_2)$

Outputs from the previous layer  H

Inputs to the next layer

✓ Cross-attention can be computed *in parallel, and be* **re-used**

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# Chunked Cross Attention



✓ Cross-attention can be computed *in parallel*, *and be re-used*

# Chunked Cross Attention



Outputs from the previous layer  H

Inputs to the next layer

This part can be re-used

If you generated until here

✓ Cross-attention can be computed *in parallel, and be re-used*

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# 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** |

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# Results

| 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** |

Borgeaud et al. 2021. "Improving language models by retrieving from trillions of tokens"

# 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**
(We'll talk more about the importance of the **datastore size** later)

# 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**
(We'll talk more about the importance of the **datastore size** later)

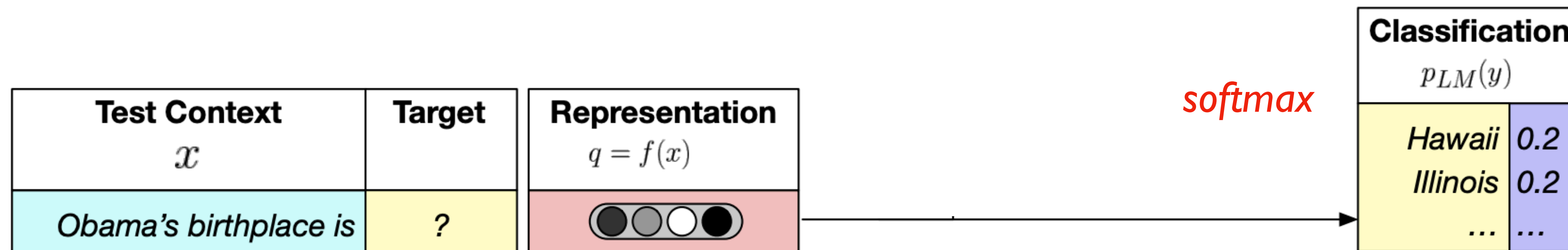# New Retrieval-based LMs: Overview

• New Methodology 1 — Designing  a new Transformer

  • **New attention layers to incorporate more blocks (RETRO)**

    • Possibly combine with long-range Transformers

• New Methodology 2 — Designing a new Softmax

• New LM Design — Mitigating fairness & legality issues

# New Retrieval-based LMs: Overview

- New Methodology 1 — Designing a new Transformer
  - New attention layers to incorporate more blocks (RETRO)
  - **Possibly combine with long-range Transformers**
- New Methodology 2 — Designing a new Softmax
- New LM Design — Mitigating fairness & legality issues

*Solve length limit issue in retrieval augmentation*
*(and probably simpler than RETRO?!)*

# New Retrieval-based LMs: Overview

- New Methodology 1 — Designing a new Transformer
  - New attention layers to incorporate more blocks (RETRO)
  - Possibly combine with long-range Transformers
- **New Methodology 2 — Designing a new Softmax** ◀ *Nonparametric softmax?*
- New LM Design — Mitigating fairness & legality issues

# kNN-LM

| Test Context $x$ | Target |
|---|---|
| Obama's birthplace is | ? |

# kNN-LM

# kNN-LM



datastore

... 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

| 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 |

datastore

… 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

| 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 | ? |  |

Khandelwal et al. 2020. "Generalization through Memorization: Nearest Neighbor Language Models"

# kNN-LM

*# of vectors = # of tokens in the corpus (>1B)*

Khandelwal et al. 2020. "Generalization through Memorization: Nearest Neighbor Language Models"

# kNN-LM



| 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 <u>in a datastore</u> are close to the next token?*

# kNN-LM



*Which tokens __in a datastore__ are close to the next token?*

**=**

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

Khandelwal et al. 2020. "Generalization through Memorization: Nearest Neighbor Language Models"

# kNN-LM

# kNN-LM

# kNN-LM

*Nonparamatric softmax*

# kNN-LM

$$P_{k\text{NN}}(y \mid x) \propto \sum_{(k,v) \in \mathcal{D}} \mathbb{I}[v = y] e^{\text{sim}(k,x)}$$

Khandelwal et al. 2020. "Generalization through Memorization: Nearest Neighbor Language Models"

# kNN-LM

$$P_{k\text{NN}}(y \mid x) \propto \sum_{(k,v)\in\mathscr{D}} \mathbb{I}[v = y]e^{\text{sim}(k,x)} \qquad \text{sim}(k, x) = -d(\text{Enc}(k), \text{Enc}(x))$$
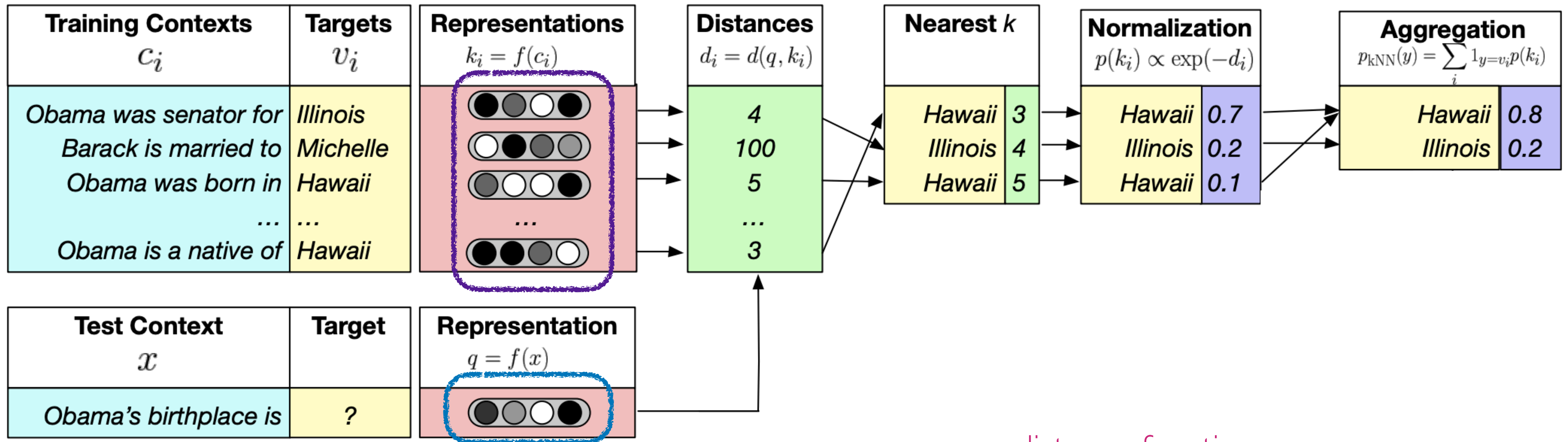
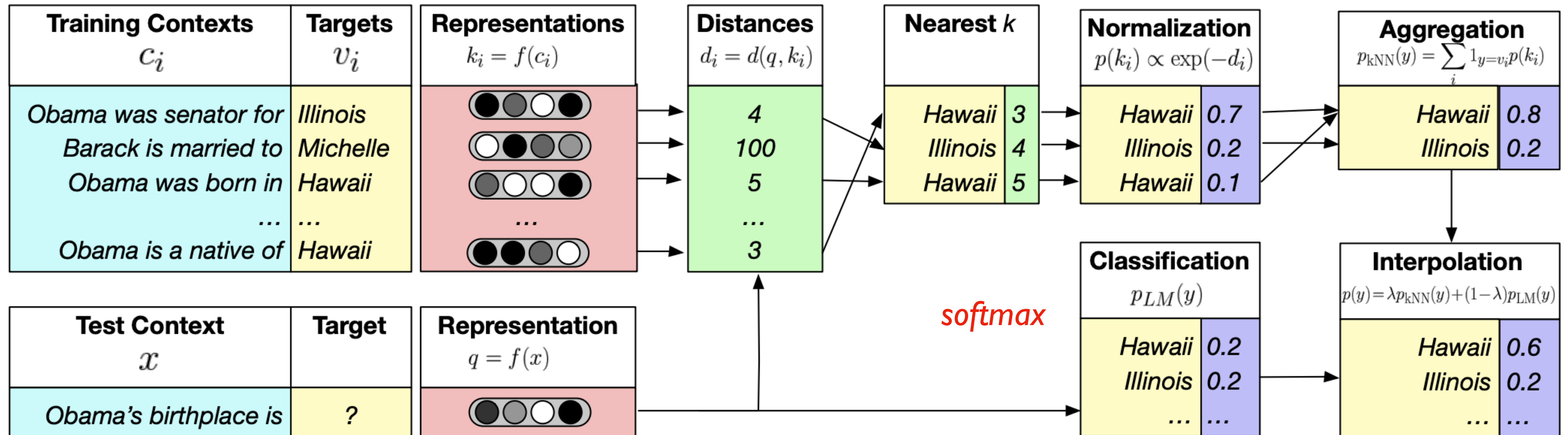Khandelwal et al. 2020. "Generalization through Memorization: Nearest Neighbor Language Models"

# kNN-LM

*Nonparamatric softmax*



$$P_{k\mathrm{NN}}(y \mid x) \propto \sum_{(k,v)\in\mathscr{D}} \mathbb{1}[v = y]e^{\mathrm{sim}(k,x)} \qquad \mathrm{sim}(k, x) = -d(\mathrm{Enc}(k), \underline{\mathrm{Enc}(x)})$$

Khandelwal et al. 2020. "Generalization through Memorization: Nearest Neighbor Language Models"

# kNN-LM

*Nonparamatric softmax*



$$P_{k\text{NN}}(y \mid x) \propto \sum_{(k,v) \in \mathscr{D}} \mathbb{I}[v = y] e^{\text{sim}(k,x)} \qquad \text{sim}(k, x) = -d(\text{Enc}(k), \text{Enc}(x))$$

# kNN-LM

*Nonparametric softmax*



| Training Contexts $c_i$ | Targets $v_i$ | Representations $k_i = f(c_i)$ | Distances $d_i = d(q, k_i)$ |
|---|---|---|---|
| Obama was senator for | Illinois | ⚫⚫⚪⚫ | 4 |
| Barack is married to | Michelle | ⚪⚫⚫⚫ | 100 |
| Obama was born in | Hawaii | ⚫⚪⚫⚪ | 5 |
| … | … | … | … |
| Obama is a native of | Hawaii | ⚫⚫⚪⚪ | 3 |

| Nearest $k$ | |
|---|---|
| Hawaii | 3 |
| Illinois | 4 |
| Hawaii | 5 |

| Normalization $p(k_i) \propto \exp(-d_i)$ | |
|---|---|
| Hawaii | 0.7 |
| Illinois | 0.2 |
| Hawaii | 0.1 |

| Aggregation $p_{\mathrm{kNN}}(y) = \sum_i 1_{y=v_i} p(k_i)$ | |
|---|---|
| Hawaii | 0.8 |
| Illinois | 0.2 |

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

distance function

$$P_{k\mathrm{NN}}(y \mid x) \propto \sum_{(k,v) \in \mathscr{D}} \mathbb{I}[v = y] e^{\mathrm{sim}(k,x)} \qquad \mathrm{sim}(k, x) = -\, d(\mathrm{Enc}(k), \mathrm{Enc}(x))$$

# kNN-LM



*Nonparamatric softmax*

*softmax*

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

# kNN-LM



*Nonparamatric softmax*

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

# kNN-LM



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

Khandelwal et al. 2020. "Generalization through Memorization: Nearest Neighbor Language Models"

# kNN-LM



$$P_{k\text{NN}-\text{LM}}(y \mid x) = (1 - \lambda)P_{\text{LM}}(y \mid x) + \lambda P_{k\text{NN}}(y \mid x) \quad \lambda: \text{hyperparameter}$$

# Why nonparametric softmax?

| 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* |

# Why nonparametric softmax?

| 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* |

# Why nonparametric softmax?

| 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* |

# Why nonparametric softmax?

| 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* |

# Why nonparametric softmax?

Dense vector space

| 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 |

... affordable

... nice

... good

... bad

... poor

... terrible

# Why nonparametric softmax?

Dense vector space

| 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* |

10/10, would buy this **cheap**

… affordable

… nice

… good

… bad

… poor

… terrible

# Why nonparametric softmax?

Dense vector space

| 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* |

10/10, would buy this **cheap**

… affordable

… nice

… good

… bad

… poor

… terrible

Item delivered broken. Very **cheap**

# Why nonparametric softmax?

| 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** |

# Why nonparametric softmax?

| 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** |

# Why nonparametric softmax?

| 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** |

# Why nonparametric softmax?

| 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** |

# Why nonparametric softmax?

Dense vector space

| 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** |



... machine

... computer

... tool

... fire

... infection

... pregnancy

# Why nonparametric softmax?

Dense vector space

| 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** |

... permitted to bring a **torch**

PyTorch, you can use **torch**

... machine

... tool

... computer

... fire

... infection

... pregnancy

... a group of infections ... **torch**

# Nonparametric-only, Phrase-level (NPM)
## (If you can train the model…)

Min et al. 2023. Nonparametric Masked Language Modeling

# Nonparametric-only, Phrase-level (NPM)

## (If you can train the model…)

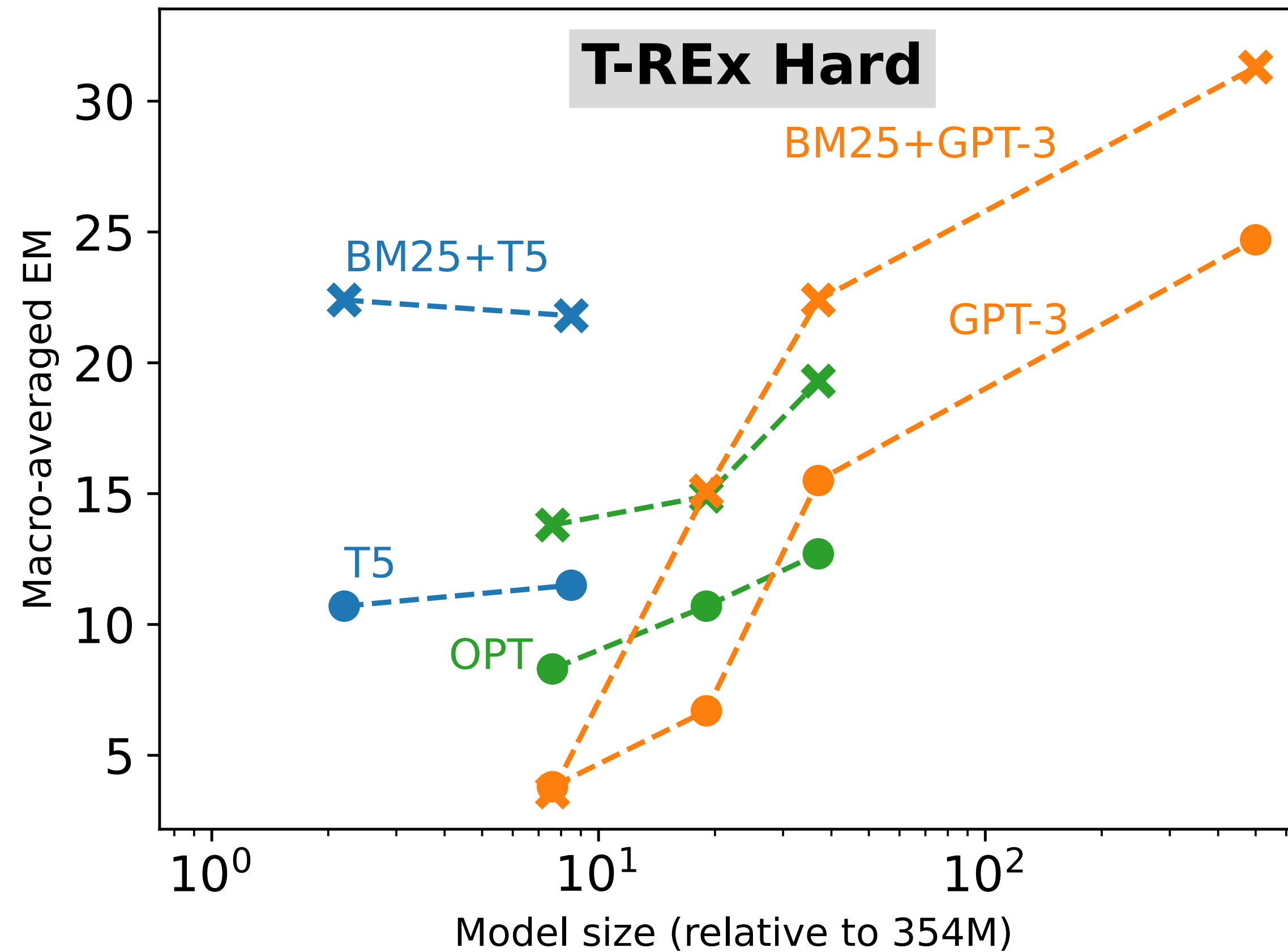# Nonparametric-only, Phrase-level (NPM)

## (If you can train the model…)

datastore

just as a jet of **red light** blasted from Harry's …

Voldemort cried, "Avada Kedavra!" A jet of **green light** issued …

"The Boy Who Lived." He saw the mouth move and a flash of **green light**, and everything was gone.

… is operated or driven by a jet of **water**.

Pick up a flat rock, skip it across **Green** River

# Nonparametric-only, Phrase-level (NPM)

## (If you can train the model…)

datastore

just as a jet of **red light** blasted from Harry's …

**Vector space**

… is operated or driven by a jet of **water**.

Voldemort cried, "Avada Kedavra!" A jet of **green light** issued …

Pick up a flat rock, skip it across **Green** River

Encoder

"The Boy Who Lived." He saw the mouth move and a flash of **green light**, and everything was gone.

Min et al. 2023. Nonparametric Masked Language Modeling

# Nonparametric-only, Phrase-level (NPM)

## (If you can train the model…)



just as a jet of **red light** blasted from Harry's …

**Vector space**

… is operated or driven by a jet of **water**.

datastore

Voldemort cried, "Avada Kedavra!" A jet of **green light** issued …

Pick up a flat rock, skip it across **Green** River

Encoder

"The Boy Who Lived." He saw the mouth move and a flash of **green light**, and everything was gone.

Harry felt Greenback collapse against him … on the floor as a jet of _____ came flying toward him.

Min et al. 2023. Nonparametric Masked Language Modeling

# Nonparametric-only, Phrase-level (NPM)

## (If you can train the model…)



just as a jet of **red light** blasted from Harry's …

… is operated or driven by a jet of **water**.

**Vector space**

datastore

Voldemort cried, "Avada Kedavra!" A jet of **green light** issued …

Pick up a flat rock, skip it across **Green** River

Encoder

"The Boy Who Lived." He saw the mouth move and a flash of **green light**, and everything was gone.

Harry felt Greenback collapse against him … on the floor as a jet of _____ came flying toward him.

*Voldemort cried, "Avada Kedavra!" A jet of **green light** issued …*
*"The Boy Who Lived." … a flash of **green light** and everything was gone.*
*Voldemort's wand just as a jet of **red light** blasted from Harry's*
*… is operated or driven by a jet of **water**.*
…

# NPM: Fact probing

# NPM: Fact probing



**T-REx Hard**

Macro-averaged EM

T5

OPT

GPT-3

Model size (relative to 354M)

No-retrieval LMs are better as they get larger

# NPM: Fact probing



Retrieval augmentation helps

Min et al. 2023. Nonparametric Masked Language Modeling

# NPM: Fact probing



NPM is more parameter efficient

Min et al. 2023. Nonparametric Masked Language Modeling

# NPM: Predicting rare entities

Min et al. 2023. Nonparametric Masked Language Modeling

# NPM: Predicting rare entities

Min et al. 2023. Nonparametric Masked Language Modeling

# NPM: Predicting rare entities

Min et al. 2023. Nonparametric Masked Language Modeling

# NPM: Predicting rare entities



NPM outperforms by a larger margin as the rarity increases

Min et al. 2023. Nonparametric Masked Language Modeling

# New Retrieval-based LMs: Overview

- New Methodology 1 — Designing a new Transformer

  - New attention layers to incorporate more blocks (RETRO)

  - Possibly combine with long-range Transformers

- New Methodology 2 — Designing a new Softmax

  - **Two softmaxes together: kNN-LM**

  - **Nonparametric softmax only, phrase-level: NPM**

- New LM Design — Mitigating fairness & legality issues

# New Retrieval-based LMs: Overview

- New Methodology 1 — Designing a new Transformer
  - New attention layers to incorporate more blocks (RETRO)
  - Possibly combine with long-range Transformers
- New Methodology 2 — Designing a new Softmax
  - Two softmaxes together: kNN-LM
  - Nonparametric softmax only, phrase-level: NPM
- **New LM Design — Mitigating fairness & legality issues**

# Common practice

*Web crawl*



*Training*

$x \longrightarrow$  $\longrightarrow y$

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# Common practice



Permissively-licensed            Copyrighted            Private

*Training*

$x \longrightarrow$    $\longrightarrow y$

# Common practice



Permissively-licensed                    Copyrighted                    Private

*Training*

$x \longrightarrow$ [neural network] $\longrightarrow y$

☹ Legal risk in training on copyrighted data

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# Common practice



Permissively-licensed      Copyrighted      Private

*Training*

$x \longrightarrow$   $\longrightarrow y$

☹ Legal risk in training on copyrighted data    ☹ Failure in crediting to data creators

# New proposal: SILO



Permissively-licensed

$x \longrightarrow$ [neural network] $\longrightarrow y$

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# New proposal: SILO



Very low legal risk,
but poor performance
(small-size data, domain shift)

Permissively-licensed

*Training*

$x \longrightarrow$ [neural network] $\longrightarrow y$

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# New proposal: SILO



Very low legal risk,
but poor performance
(small-size data, domain shift)

Permissively-licensed

Significantly improve generalization

*Training*

$x \longrightarrow$ $\longrightarrow y$

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# New proposal: SILO



Very low legal risk,
but poor performance
(small-size data, domain shift)

Permissively-licensed

Significantly improve generalization

*Training*

$x \longrightarrow$ [neural network] $\longrightarrow y$

☑ Can trace inherent attribution

☑ Can modify the datastore at any time

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# New proposal: SILO



Permissively-licensed

Very low legal risk,
but poor performance
(small-size data, domain shift)

Significantly improve generalization

*Training*

$x \longrightarrow$ [neural network] $\longrightarrow y$

☑ Can trace inherent attribution
- Likely defense *fair use*
- Provide copyright notice
- Allow credits (or payment) to data creators

☑ Can modify the datastore at any time

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# New proposal: SILO



Very low legal risk,
but poor performance
(small-size data, domain shift)

Permissively-licensed

Significantly improve generalization

*Training*

☑ Can trace inherent attribution
- Likely defense *fair use*
- Provide copyright notice
- Allow credits (or payment) to data creators

$x \longrightarrow$ ☐ $\longrightarrow y$

☑ Can modify the datastore at any time
- Support removal of data at any time
- Better alignment with *GDPR*

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# SILO Attribution Example

**Test input:**
include '../lib/admin.defines.php';
include '../lib/admin.module.access.php';
include '../lib/admin.smarty.php';
if (! has_right (

*Continuation:* **[AC]**X_BILLING)) { Header …

# SILO Attribution Example

**Test input:**
include '../lib/admin.defines.php';
include '../lib/admin.module.access.php';
include '../lib/admin.smarty.php';
if (! has_right (

**Continuation:** **[AC]**X_BILLING)) { Header ...

# SILO Attribution Example

**Test input:**
include '../lib/admin.defines.php';
include '../lib/admin.module.access.php';
include '../lib/admin.smarty.php';
if (! has_right (

**Continuation: [AC]**X_BILLING)) { Header …

**Top-1 retrieved token (in kNN-LM):**
*You should have received a copy of the GNU Affero General Public License
* along with this program. If not, see <http://www.gnu.org/licenses/>.
*
*
**/
if (! has_right (
        **[AC]**X_ACCESS)) { Header …

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# SILO Attribution Example

**Test input:**
include '../lib/admin.defines.php';
include '../lib/admin.module.access.php';
include '../lib/admin.smarty.php';
if (! has_right (

*Continuation:* **[AC]**X_BILLING)) { Header …

***Top-1 retrieved token (in kNN-LM):***
\*You should have received a copy of the GNU Affero General Public License
\* along with this program. If not, see <http://www.gnu.org/licenses/>.
\*
\*
\*\*/
if (! has_right (
         **[AC]**X_ACCESS)) { Header …

# New Retrieval-based LMs: Summary

- New Methodology 1 — Designing a new Transformer

  - New attention layers to incorporate more blocks (RETRO)

- New Methodology 2 — Designing a new Softmax

  - Two softmaxes together: kNN-LM

  - Nonparametric softmax only, phrase-level: NPM

- New LM Design — Mitigating fairness & legality issues

  - Train on permissive text → place copyrighted text into a datastore

# Overview

## Why Retrieval-based LMs?

Tell me about Meta Platform.

I don't have any information about a company called Meta Platforms. It is possible that the company is …

**ChatGPT**

## Retrieval Augmentation

$x \longrightarrow$ **Retrieval** $\longrightarrow$

**LM** $\longrightarrow y$

## New Retrieval-based LMs

$x \longrightarrow$ **LM**

… *"Avada Kedavra!" A jet of* **green light** *issued* …

… *move and a flash of* **green light** *and* …

… *just as a jet of* **red light** *blasted from Harry's*

… *is operated or driven by a jet of* **water**.

…

## Open Problems

datastore

Scaling **datastore** not just parameters?

# Summary

**What?**                    **How?**                    **Why?**

# Summary

**What?**                    **How?**                    **Why?**



training

(Typical LMs)

$x$: test input
$y$: model prediction to $x$

# Summary

**What?**  **How?**  **Why?**



$x$: test input
$y$: model prediction to $x$

# Summary

**What?**



$x$: test input
$y$: model prediction to $x$

**How?**

Retrieval augmentation

New Transformers

Nonparametric Softmax

**Why?**

# Summary

**What?**



$x$: test input
$y$: model prediction to $x$

**How?**

Retrieval augmentation

New Transformers

Nonparametric Softmax

**Why?**

New dimension in improving LMs!

# Summary



slide 86 ↓

slide 89 ↓

**Why?**

New dimension in improving LMs!

101

# Summary

**What?**



$x$: test input
$y$: model prediction to $x$

**How?**

Retrieval augmentation

New Transformers

Nonparametric Softmax

**Why?**

New dimension in improving LMs!

Update & scale without additional training

# Summary



slide 45 ↓

**Test- 2017**
**Test - 2020**

T5-Fine-tuned on 2017 data

**Train 2017, DS 2017**
**Train 2017, DS 2020**

## Why?

New dimension in improving LMs!

Update & scale without additional training

# Summary

slide 96 ↓

**Test input:**
include '../lib/admin.defines.php';
include '../lib/admin.module.access.php';
include '../lib/admin.smarty.php';
if (! has_right (

*Continuation:* **[AC]**X_BILLING)) { Header …

**Top-1 retrieved context:**
*You should have received a copy of the GNU Affero General Public License
* along with this program. If not, see <http://www.gnu.org/licenses/>.
*
*
**/
if (! has_right (
        **[AC]**X_ACCESS)) { Header …

## Why?

New dimension in improving LMs!

Update & scale without additional training

Provide data attribution

# Summary

slide 96 ↓

**Test input:**
include '../lib/admin.defines.php';
include '../lib/admin.module.access.php';
include '../lib/admin.smarty.php';
if (! has_right (

*Continuation:* **[AC]**X_BILLING)) { Header …

**Top-1 retrieved context:**
*You should have received a copy of the GNU Affero General Public License
* along with this program. If not, see <http://www.gnu.org/licenses/>.
*
*
**/
if (! has_right (
        **[AC]**X_ACCESS)) { Header …

## Why?

New dimension in improving LMs!

Update & scale without additional training

Provide data attribution

New opportunities in fairness & legality

# Summary

**What?**



training

$x \longrightarrow$ datastore $\longrightarrow y$

$x$: test input
$y$: model prediction to $x$

**How?**

Retrieval augmentation

New Transformers

Nonparametric Softmax

**Why?**

New dimension in improving LMs!

Update & scale without additional training

Provide data attribution

New opportunities in fairness & legality

# Open questions

# Open question: Scaling retrieval-based LMs

# Open question: Scaling retrieval-based LMs

A small LM + a large datastore ≈ a large (no-retrieval) LM?



vs.

# Open question: Scaling retrieval-based LMs

## A small LM + a large datastore ≈ a large (no-retrieval) LM?



A new dimension in scaling!

Min et al. 2023. "SILO Language Models: Isolating Legal Risk In a Nonparametric Datastore"

# Open question: Scaling retrieval-based LMs

A small LM + a large datastore ≈ a large (no-retrieval) LM?



vs.

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

# Open question: Scaling retrieval-based LMs

A small LM + a large datastore ≈ a large (no-retrieval) LM?



vs.

|  | 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 | ~30B |
| **RETRO** (Borgeaud et al., 2021) | 7B | 2T |
| **REPLUG** (Shi et al., 2023) | $\leq$ 175B | ~5B |

# Open question: Scaling retrieval-based LMs

Scaling law?

# Open question: Scaling retrieval-based LMs

Scaling law?



Loss as a function of:

• Training data size

• # model parameters

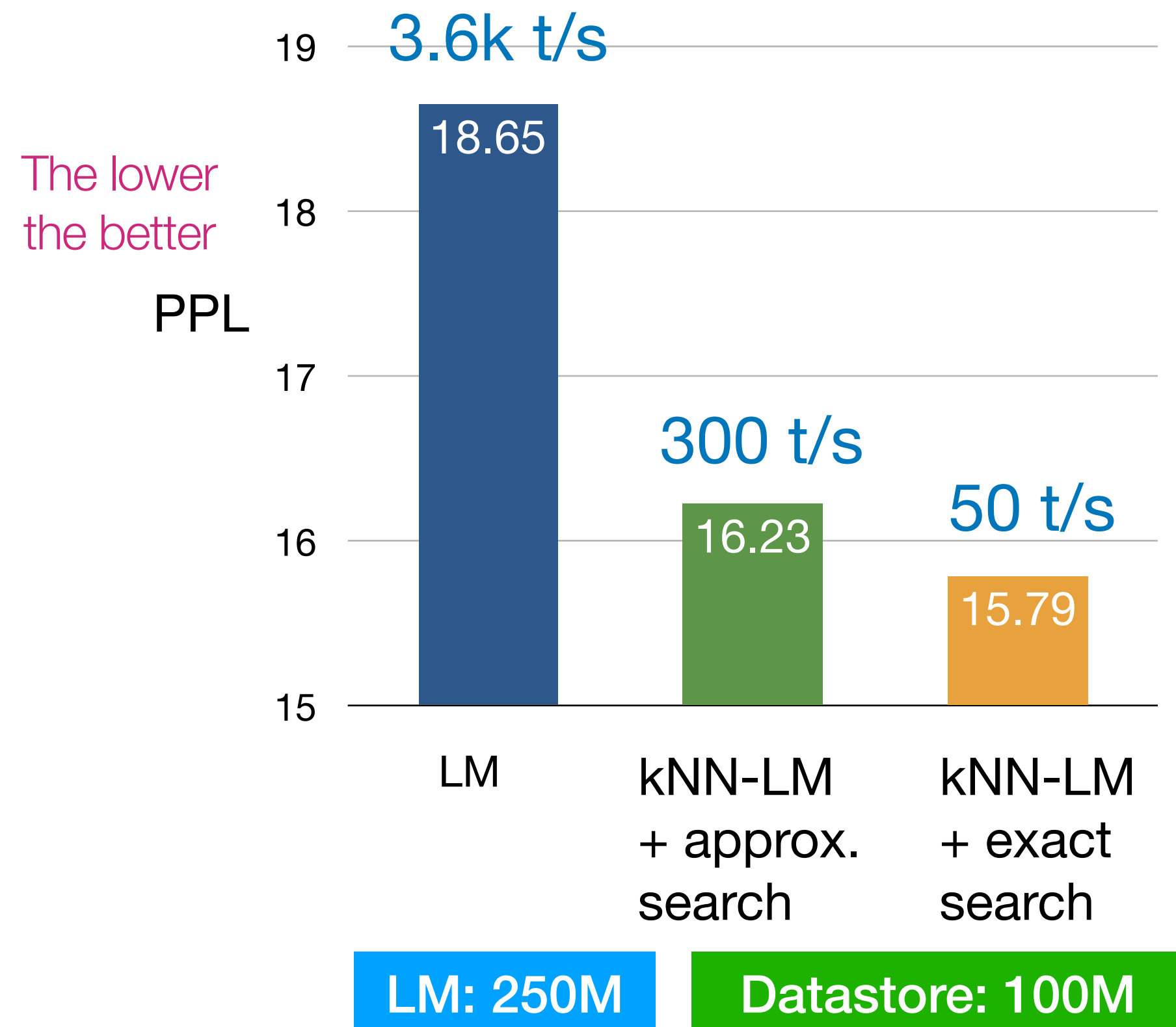Scaling law for parametric LMs (Kalpan et al., 2020; Hoffman et al., 2022)

# Open question: Scaling retrieval-based LMs

## Scaling law?



Loss as a function of:

• Training data size

• # model parameters

+ Datastore sizes?

Scaling law for parametric LMs (Kalpan et al., 2020; Hoffman et al., 2022)

# Open question: Runtime efficiency

Efficiency of similarity search

Guo et al. 2020. "Accelerating Large-Scale Inference with Anisotropic Vector Quantization"

# Open question: Runtime efficiency

## Efficiency of similarity search

Measured on NVIDIA RTX 3090 GPU (Zhong et al., 2022)
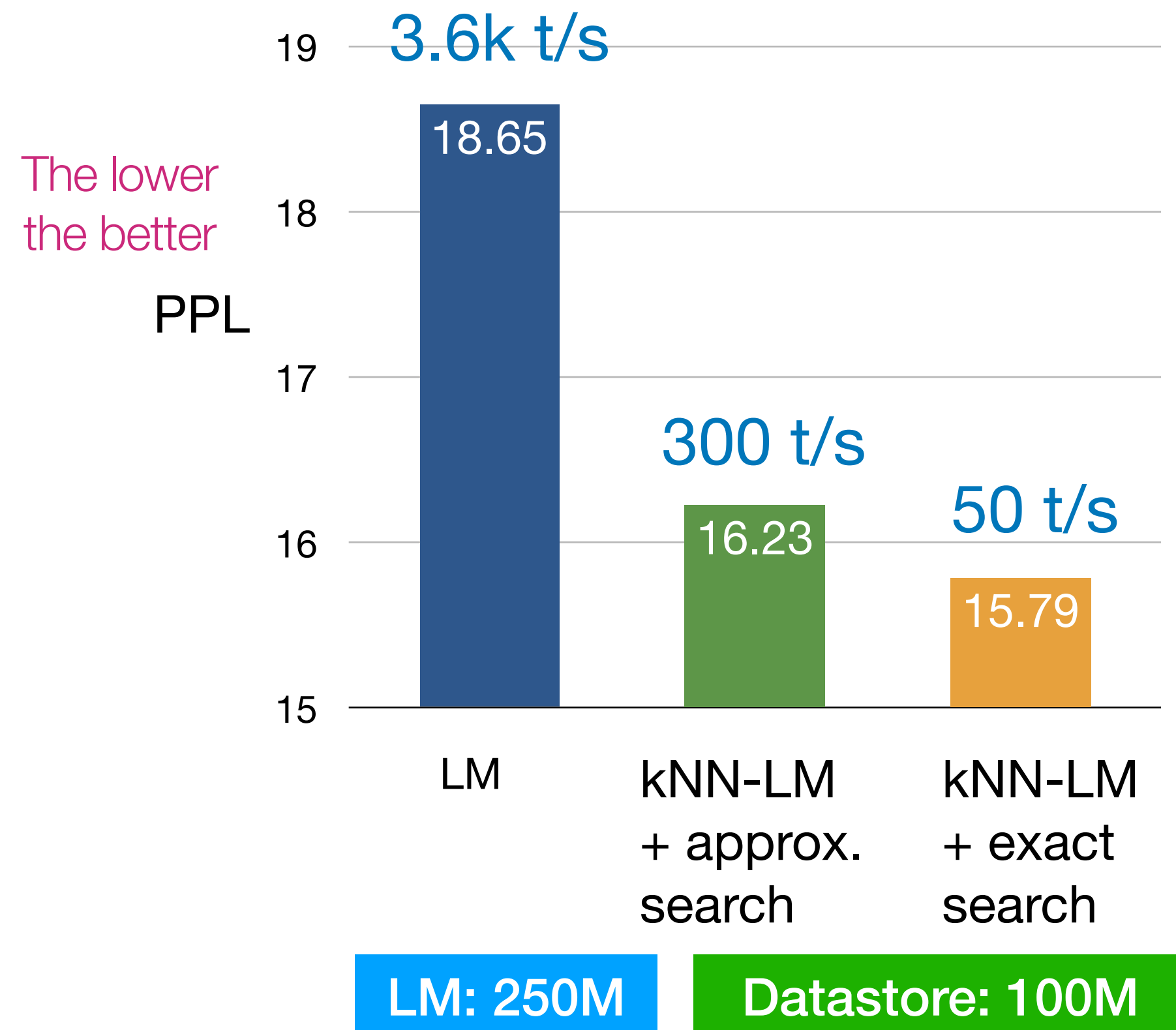with a FAISS indexer (Johnson et al., 2021) with 32 CPUs

# Open question: Runtime efficiency

## Efficiency of similarity search

Measured on NVIDIA RTX 3090 GPU (Zhong et al., 2022)
with a FAISS indexer (Johnson et al., 2021) with 32 CPUs



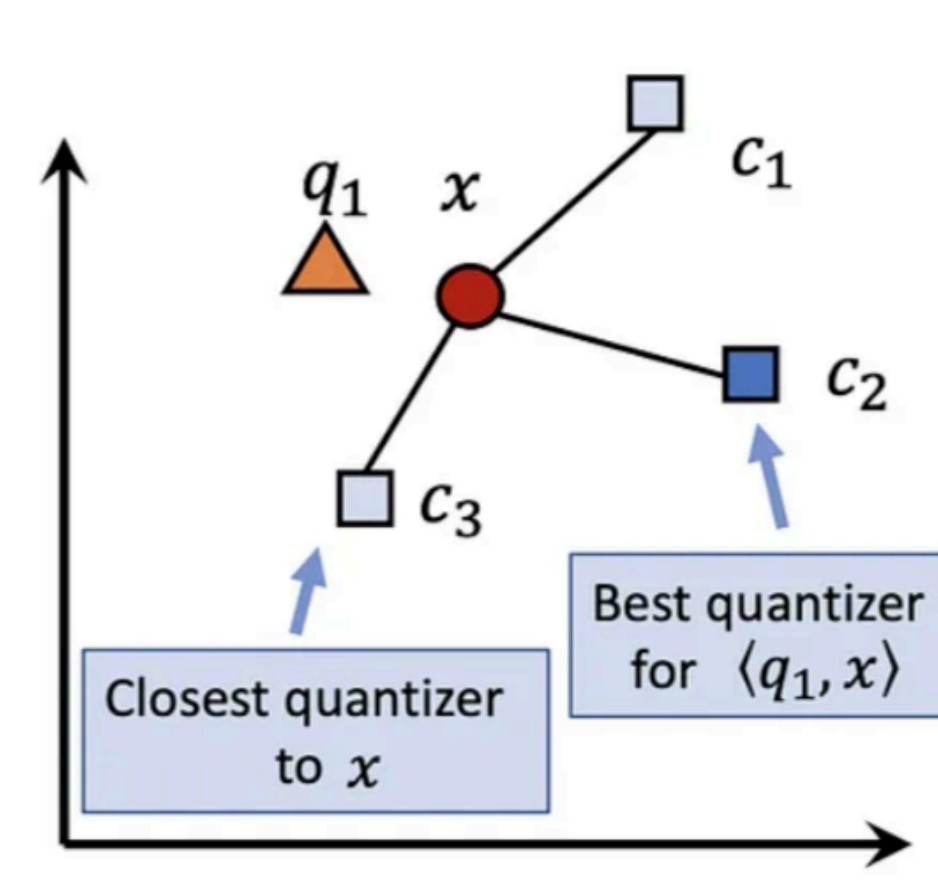- >12 times slower with **approximate** nearest neighbor search

# Open question: Runtime efficiency

## Efficiency of similarity search

Measured on NVIDIA RTX 3090 GPU (Zhong et al., 2022)
with a FAISS indexer (Johnson et al., 2021) with 32 CPUs
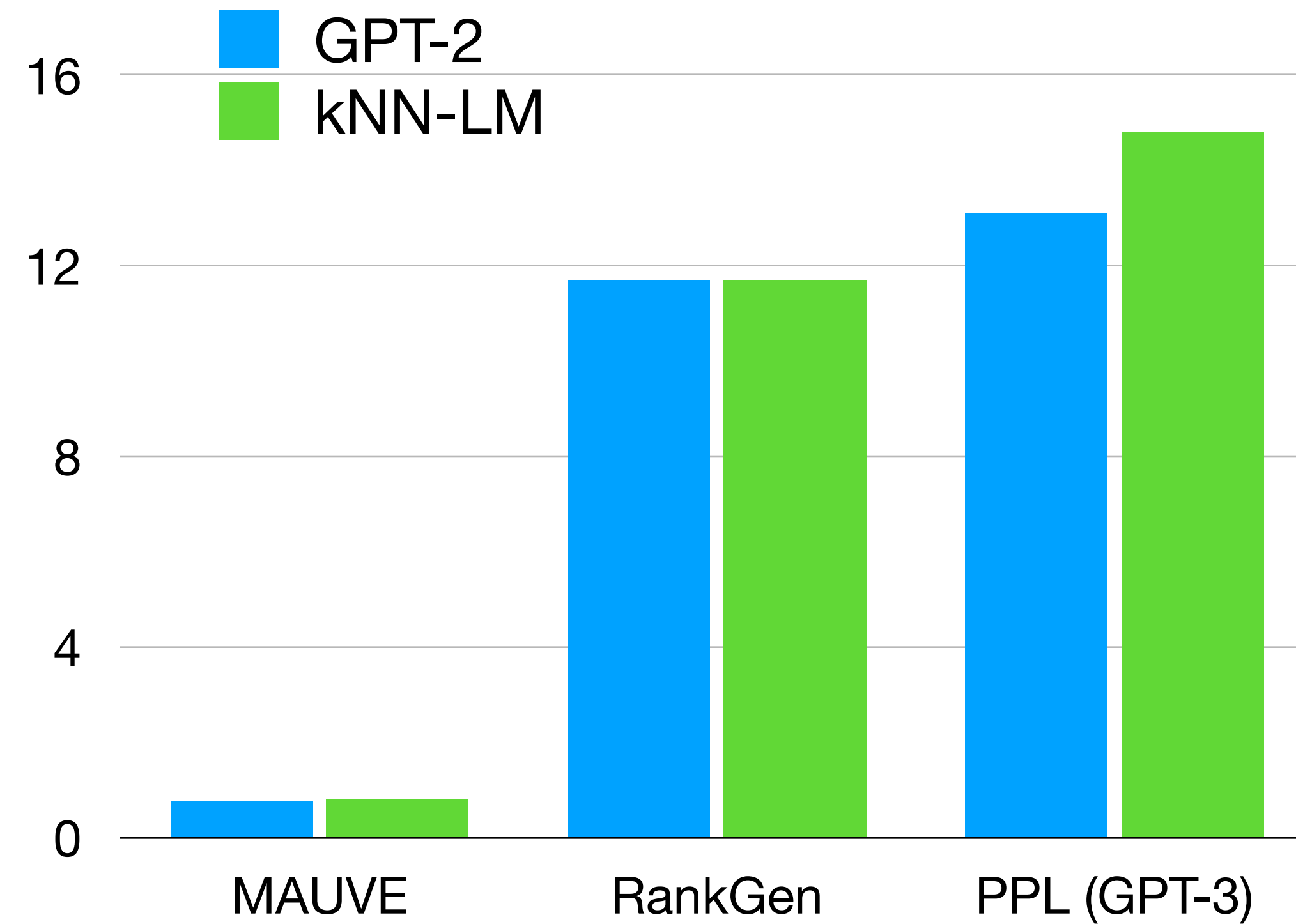


The lower the better

- >12 times slower with **approximate** nearest neighbor search

- Efficient similarity search is an active research area (in conjunction with **systems**, **databases**, & **algorithms**)



Guo et al. 2020. "Accelerating Large-Scale Inference with Anisotropic Vector Quantization"

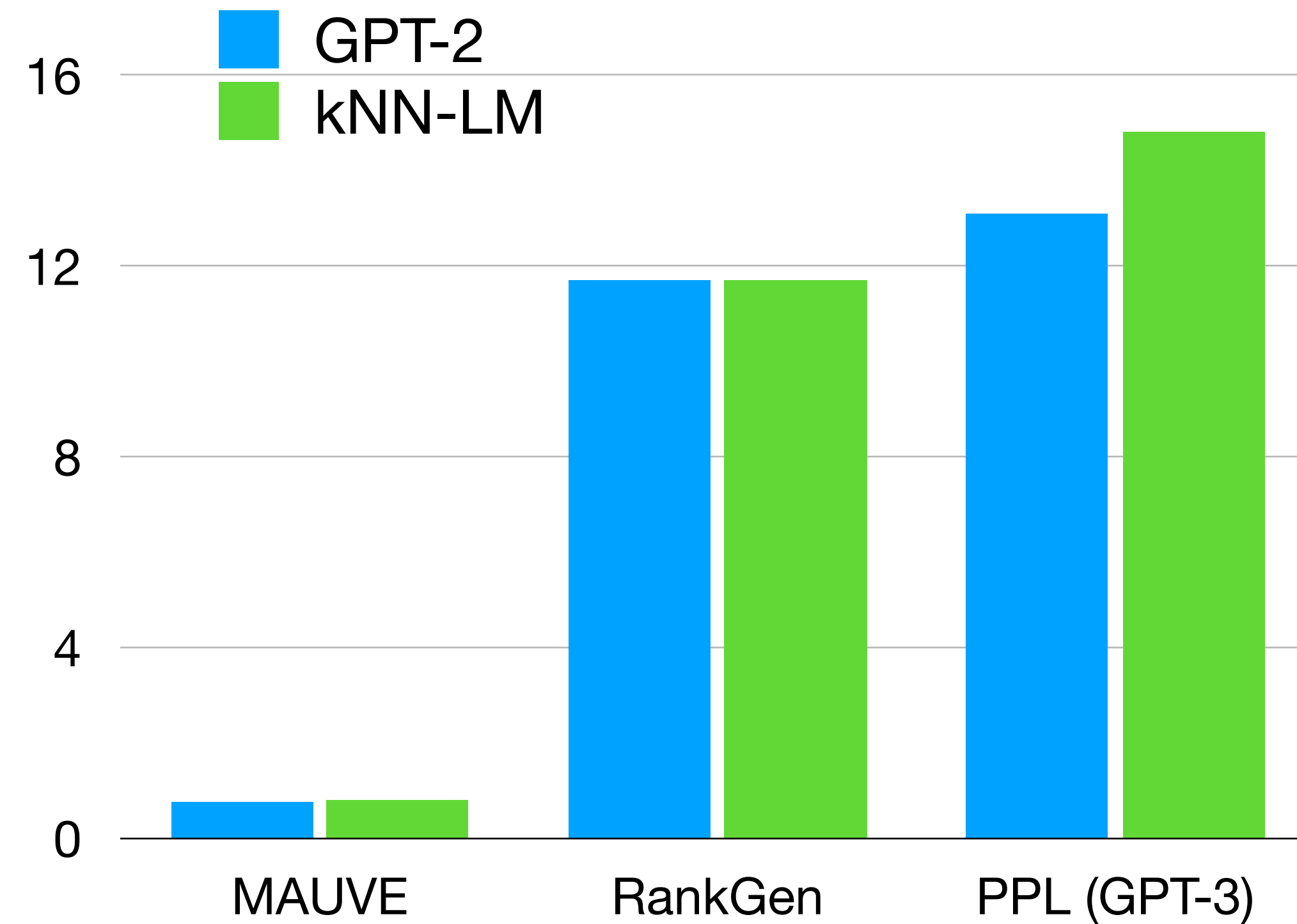# Open question: Retrieval-based LMs for applications

# Open question: Retrieval-based LMs for applications

Open-ended text generation?

# Open question: Retrieval-based LMs for applications

Open-ended text generation?



Better decoding algorithms? Better adaptation methods?

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

# Open questions: Summary

# Open questions: Summary

- What is the best **architecture & training method** for retrieval-based LMs in practice?

# Open questions: Summary

- What is the best **architecture & training method** for retrieval-based LMs in practice?

- How to **scale the datastore** to trillions of tokens?
  [Scaling law]

# Open questions: Summary

- What is the best **architecture & training method** for retrieval-based LMs in practice?

- How to **scale the datastore** to trillions of tokens? [Scaling law]

- How to improve **runtime efficiency**?

# Open questions: Summary

- What is the best **architecture & training method** for retrieval-based LMs in practice?

- How to **scale the datastore** to trillions of tokens? [Scaling law]

- How to improve **runtime efficiency**?

- How to design **new decoding** or **adaptation methods** for downstream tasks (e.g., open-ended text generation)!

# Open questions: Summary

- What is the best **architecture & training method** for retrieval-based LMs in practice?

- How to **scale the datastore** to trillions of tokens? [Scaling law]

- How to improve **runtime efficiency**?

- How to design **new decoding** or **adaptation methods** for downstream tasks (e.g., open-ended text generation)!

# Q & A

## Thank you for listening!

Check out ACL 2023 Tutorial on this topic (3-hour): https://acl2023-retrieval-lm.github.io/

Please leave feedback at tinyurl.com/sewon-min-talk
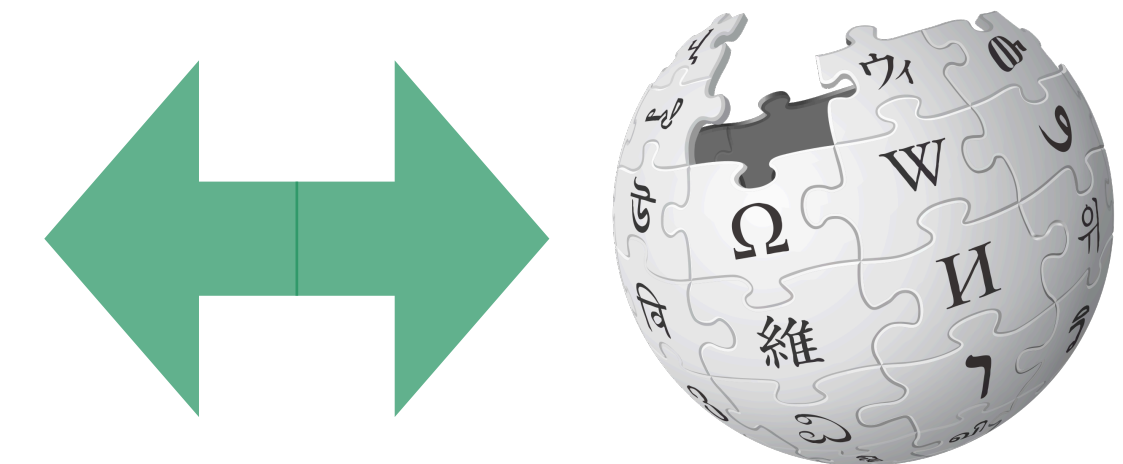
# Extra slides (from QnA)

# Validating Model Output to be Factual

Bridget Moynahan is an American actress, model and producer. She is best known for her roles in Grey's Anatomy, I, Robot and Blue Bloods. She studied acting at the American Academy of Dramatic Arts, and ...
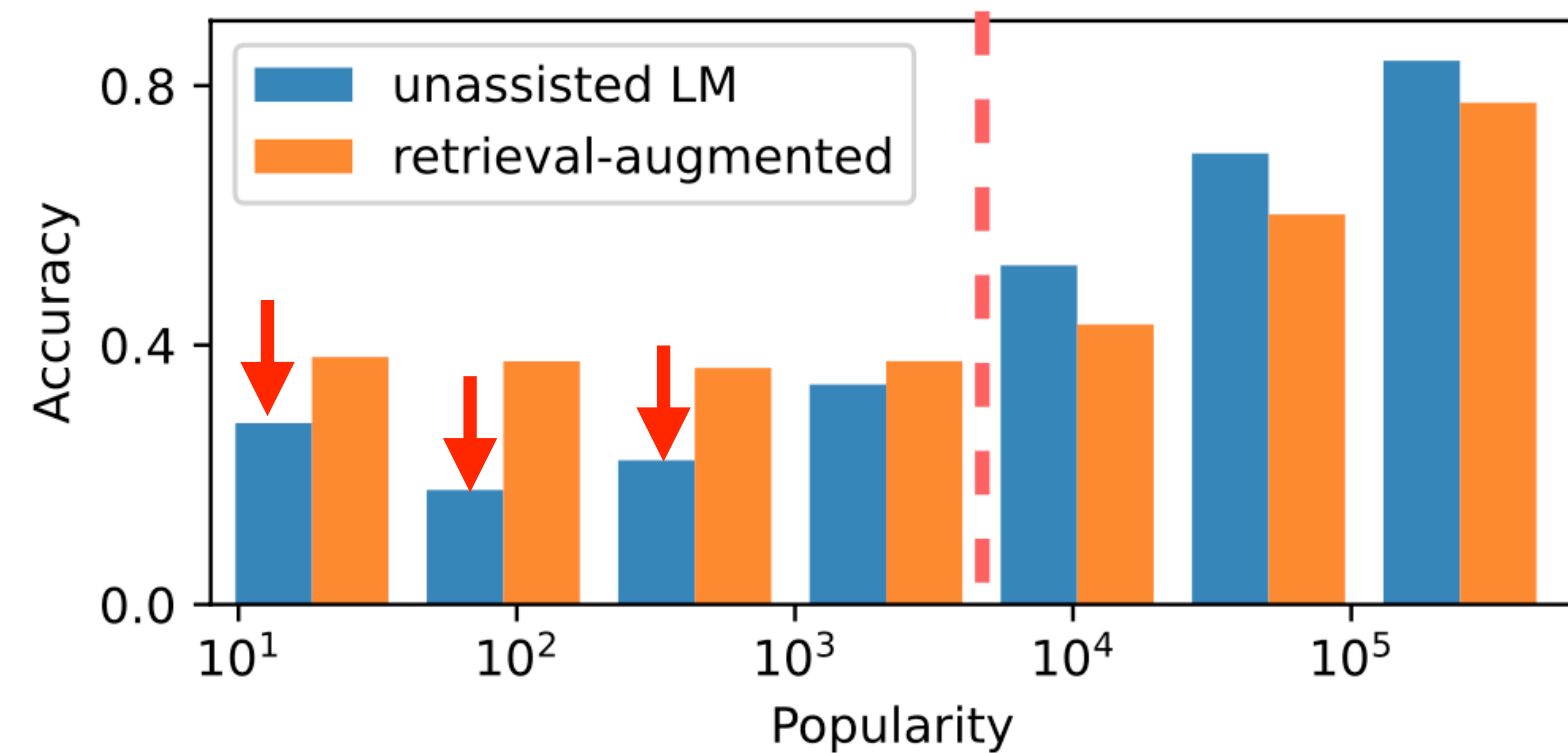
## Atomic facts

- Bridget Moynahan is American. ✔
- Bridget Moynahan is an actress. ✔
- Bridget Moynahan is a model. ✔
- Bridget Moynahan is a producer. ✘
- She is best known for her roles in Grey's Anatomy. ✘
- She is best known for her roles in I, Robot. ✔
- She is best known for her roles in Blue Bloods. ✔
- She studied acting. ✔
- She studied at the American Academy of Dramatic Arts. ✘
- ...

**66.7%**

Min et al. 2023. "FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation"
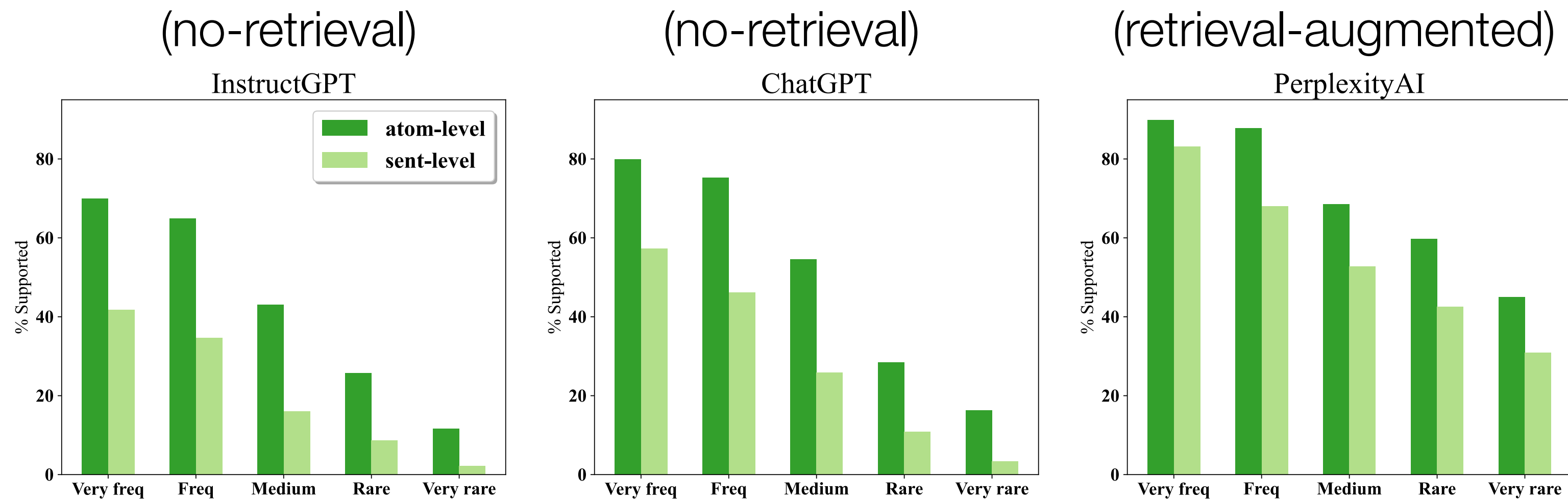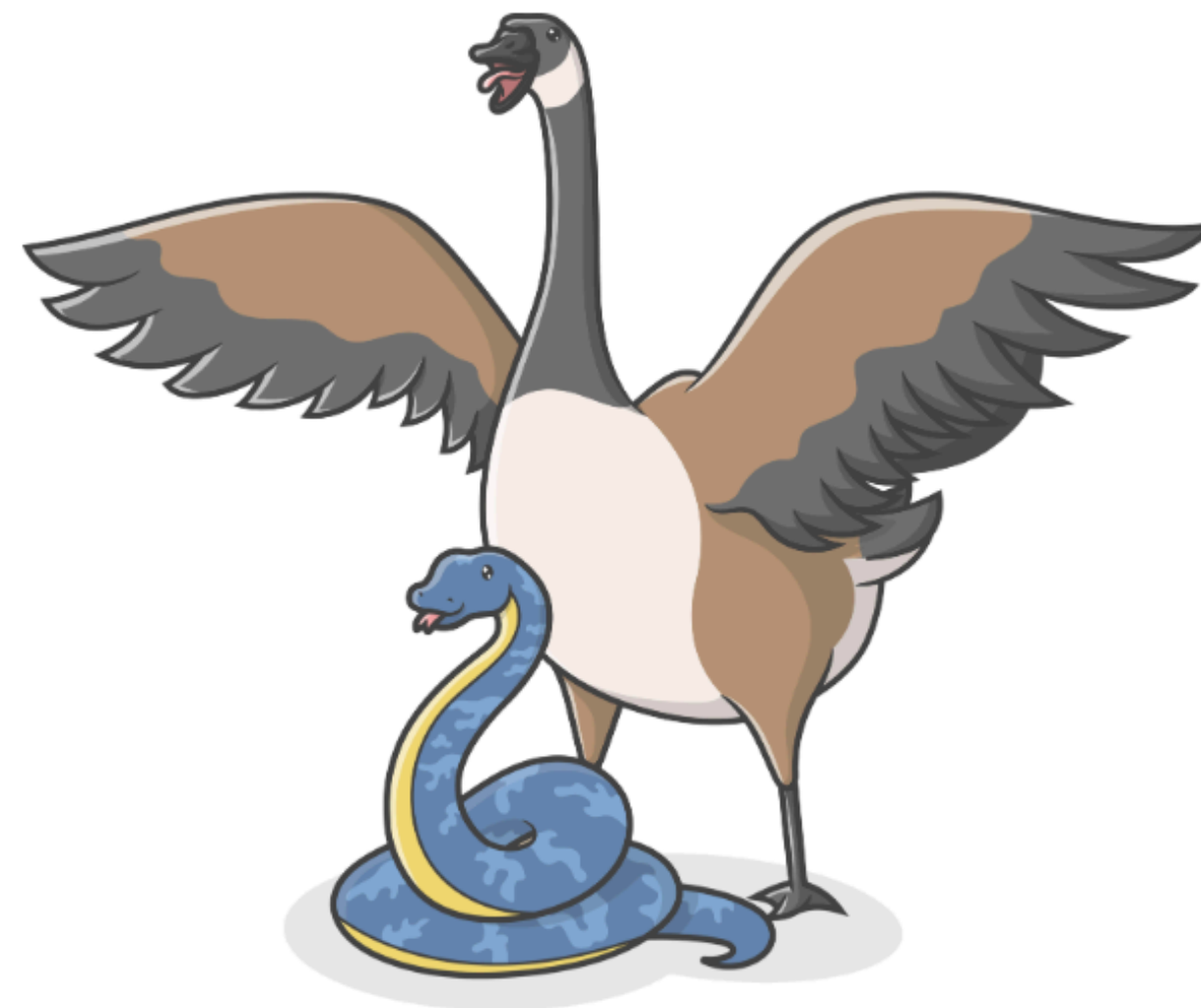
# Gains from retrieval w.r.t. frequency



There has been mixed results about whether retrieval hurts when it comes to popular entities/facts, e.g., the top graph shows it does hurt in (short-form) question answering, and the bottom graph shows retrieval always help even with frequent entities in long-form text generation. These results are likely to depend on exact setup, e.g., the task, base LMs, and datastore, etc.

Mallen et al. 2023. "When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories"

## (no-retrieval)          (no-retrieval)          (retrieval-augmented)



Min et al. 2023. "FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation"

# Research on information retrieval



Retrieval—including training the encoder, getting embeddings and indexing—is an active area of research. Recommend Pyserini (https://github.com/castorini/pyserini) for a set of references and also try some of them out easily.

# State-of-the-art retrieval-based LMs?

- If you want the model that you can use right now — retrieval-augmentation
  - Partially because you can leverage the state-of-the-art models that industry built with no modification
  - You should use state-of-the-art retrieval (BM25, Contriever or GTR) and state-of-the-art LM (LLAMA, ChatGPT)
  - Easiest: with "independent training", optionally with reranking
- Doesn't mean retrieval-augmentation is the "best" under the scenario of fair comparison, e.g., when the model has exact same parameters & is trained on the exactly same data
  - The SILO paper shows kNN-LM (kNN in the graph) outperforms retrieval-augmentation (RiC in the graph), both when training data==datastore (right) and when training data!=datastore (left)
    - However, this is based on language modeling perplexity. Downstream task eval is still an open Q.