



# Language Models as the New Search Engines

**Tianyu Gao**

Princeton Language and Intelligence (PLI)

Princeton University

**How do we conduct  
information-seeking activities?**

# The Internet era: search engines

# The Internet era: search engines



Who is Taylor Swift



All Images News Perspectives Videos More Tools

About 764,000,000 results (0.54 seconds)

## Taylor Swift

American singer-songwriter

Overview

Songs

The BRIT Awards

Events

Albums

Videos

Relationships

Listen



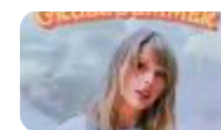
People  
All About Taylor Swift's Parents, Scott and Andrea Swift - People  
All About Taylor Swift's Parents, Scott and Andrea Swift · They raised Taylor on a...  
1 week ago

Age  
34 years  
December 13, 1989

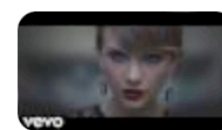
Net worth  
1.1 billion USD (2024)  
Forbes



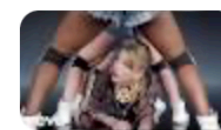
## Songs >



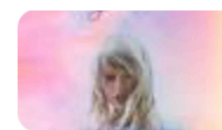
Cruel Summer  
Lover · 2019



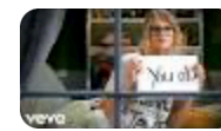
Blank Space  
1989 · 2014



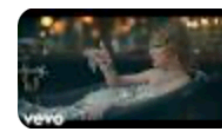
Shake it Off  
1989 · 2014



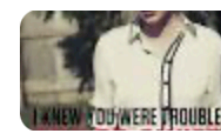
Lover  
Lover · 2019



You Belong With Me  
Fearless · 2008



Look What You Made Me Do  
reputation · 2017



I Knew You Were Trouble.  
Red · 2012



Don't Blame Me  
reputation · 2017

View 20+ more →

## Listen



YouTube



Spotify



iHeartRadio



Apple Music

## About

taylorswift.com

Taylor Alison Swift is an American singer-songwriter. Her reinventive artistry, songwriting and entrepreneurship have influenced the music industry, popular culture, and politics, while her life is a subject of widespread media coverage. Swift began professional songwriting at 14. [Wikipedia](#)



# The Internet era: search engines



Where can I find gluten-free sushi (gluten free soy sauce provided) in Princ



# The Internet era: search engines



Where can I find gluten-free sushi (gluten free soy sauce provided) in Princ



Places :

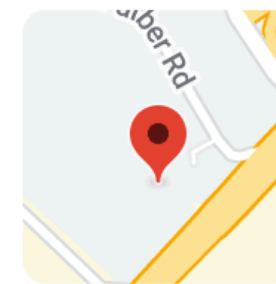


Rating ▾ Price ▾ Hours ▾

### Sponsored

#### Whole Foods Market :

4.4 ★★★★★ (1.4K) · \$\$\$ · Grocery store  
3495 U.S. 1 South  
Open · Closes 9 PM  
Delivery



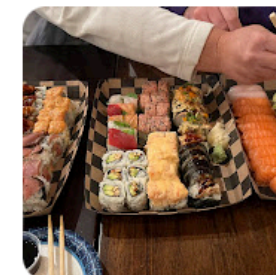
#### MTea Sushi & Dessert

4.8 ★★★★★ (213) · Sushi  
86 Nassau St  
"... thought of little details like providing small containers of **soy sauce**."



#### Masa Sushi-Princeton

4.2 ★★★★★ (1.1K) · \$20–30 · Sushi  
415 Nassau Park Blvd  
"No **gluten free soy sauce** and they forgot our Godzilla roll but other ..."



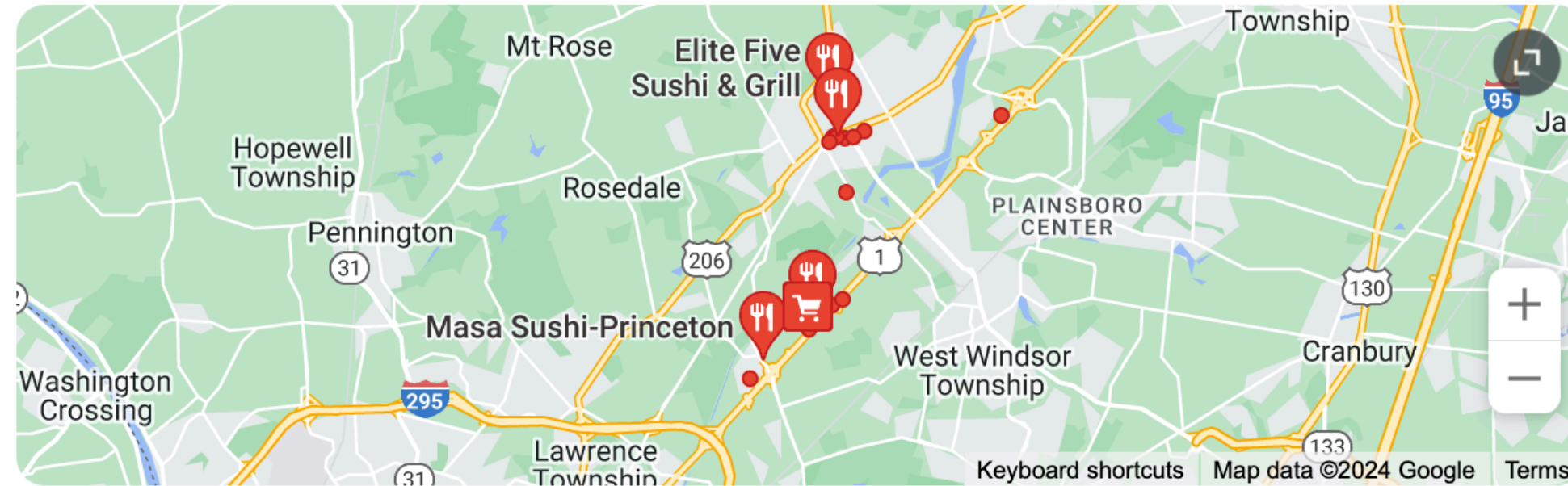
# The Internet era: search engines



Where can I find gluten-free sushi (gluten free soy sauce provided) in Princ



Places :



Rating Price Hours

Sponsored

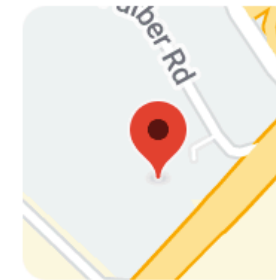
Whole Foods Market :

4.4 ★★★★★ (1.4K) · \$\$\$ · Grocery store

3495 U.S. 1 South

Open · Closes 9 PM

Delivery



MTea Sushi & Dessert

4.8 ★★★★★ (213) · Sushi

86 Nassau St

"... thought of little details like providing small containers of **soy sauce**."

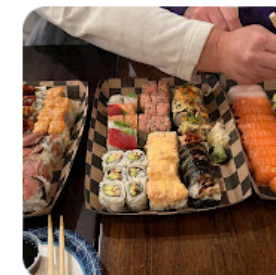


Masa Sushi-Princeton

4.2 ★★★★★ (1.1K) · \$20-30 · Sushi

415 Nassau Park Blvd

"No **gluten free soy sauce** and they forgot our Godzilla roll but other ..."



Do not directly provide answers



# The Internet era: search engines



Where can I find gluten-free sushi (gluten free soy sauce provided) in Princ



Places :



Rating Price Hours

Sponsored

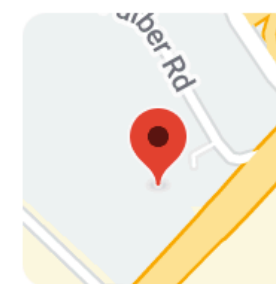
Whole Foods Market :

4.4 ★★★★★ (1.4K) · \$\$\$ · Grocery store

3495 U.S. 1 South

Open · Closes 9 PM

Delivery



MTea Sushi & Dessert

4.8 ★★★★★ (213) · Sushi

86 Nassau St

"... thought of little details like providing small containers of **soy sauce**."

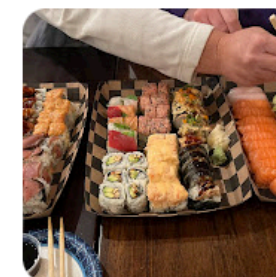


Masa Sushi-Princeton

4.2 ★★★★★ (1.1K) · \$20–30 · Sushi

415 Nassau Park Blvd

"No **gluten free soy sauce** and they forgot our Godzilla roll but other ..."



Do not directly provide answers  
Need to verify results one by one

# The Internet era: search engines

Google

🔍 In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?





# The Internet era: search engines



In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?



Johns Hopkins University  
<https://jhpce.jhu.edu> › knowledge-base › setting-a-time-li...

## Setting a time limit for your SLURM job on JHPCE

Feb 13, 2024 — The default **time limit** for **job** execution on the shared and **gpu** partitions is set to 1 day, ensuring that shorter **tasks do** not face unnecessary ...

Missing: ~~restriet~~ most



Unix & Linux Stack Exchange  
<https://unix.stackexchange.com> › questions › running-...

## Running multiple SLURM jobs on the same GPU

Nov 25, 2022 — It **will** run **jobs** in parallel if you have multiple **GPUs** that **can** run the **jobs**, otherwise it runs them in **series** as a **GPU** becomes available. I ...

[1 answer](#) · Top answer: You should use "Sharding" GRES (gres:shard) instead of gres:GPU, a...

Missing: ~~range?~~ | Show results with: [range?](#)



Slurm Documentation  
<https://slurm.schedmd.com> › faq

## Frequently Asked Questions - Slurm Workload Manager

Feb 12, 2024 — If **jobs do** not specify **time** limits, then all **jobs will** receive the same **time limit** (that associated with the partition), and the ability to ...

# The Internet era: search engines



In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?



Johns Hopkins University  
<https://jhpce.jhu.edu> › knowledge-base › setting-a-time-li...

## Setting a time limit for your SLURM job on JHPCE

Feb 13, 2024 — The default **time limit** for **job** execution on the shared and **gpu** partitions is set to 1 day, ensuring that shorter **tasks do** not face unnecessary ...

Missing: ~~restriet~~ most



Unix & Linux Stack Exchange  
<https://unix.stackexchange.com> › questions › running-...

## Running multiple SLURM jobs on the same GPU

Nov 25, 2022 — It **will** run **jobs** in parallel if you have multiple **GPUs** that **can** run the **jobs**, otherwise it runs them in **series** as a **GPU** becomes available. I ...

1 answer · Top answer: You should use "Sharding" GRES (gres:shard) instead of gres:GPU, a...

Missing: ~~range?~~ | Show results with: [range?](#)



Slurm Documentation  
<https://slurm.schedmd.com> › faq

## Frequently Asked Questions - Slurm Workload Manager

Feb 12, 2024 — If **jobs do** not specify **time** limits, then all **jobs will** receive the same **time limit** (that associated with the partition), and the ability to ...

Too many irrelevant results

# The Internet era: search engines



In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?



Johns Hopkins University  
<https://jhpce.jhu.edu> › knowledge-base › setting-a-time-li...

## Setting a time limit for your SLURM job on JHPCE

Feb 13, 2024 — The default **time limit** for **job** execution on the shared and **gpu** partitions is set to 1 day, ensuring that shorter **tasks do** not face unnecessary ...

Missing: ~~restriet~~ most



Unix & Linux Stack Exchange  
<https://unix.stackexchange.com> › questions › running-...

## Running multiple SLURM jobs on the same GPU

Nov 25, 2022 — It **will** run **jobs** in parallel if you have multiple **GPUs** that **can** run the **jobs**, otherwise it runs them in **series** as a **GPU** becomes available. I ...

1 answer · Top answer: You should use "Sharding" GRES (gres:shard) instead of gres:GPU, a...

Missing: ~~range?~~ | Show results with: [range?](#)



Slurm Documentation  
<https://slurm.schedmd.com> › faq

## Frequently Asked Questions - Slurm Workload Manager

Feb 12, 2024 — If **jobs do** not specify **time** limits, then all **jobs will** receive the same **time limit** (that associated with the partition), and the ability to ...

Too many irrelevant results  
Word-matching does not return relevant webpages



# Language models: the new “search engines”

**You**

In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?

# Language models: the new “search engines”

**You**

In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?

**ChatGPT**

To restrict the number of GPUs a user's jobs can use at most within a specific job time range in a SLURM environment, you will need to utilize SLURM's Quality of Service (QoS) and account management features. This approach allows you to define limits on resource usage that apply to users, accounts, or even specific jobs based on various criteria, including job time ranges. Here's a step-by-step guide to setting up such restrictions:

## 1. Define or Update a Quality of Service (QoS)

First, you need to define a Quality of Service (QoS) that specifies the GPU usage limits. If a suitable QoS already exists, you may update it instead.

To create a new QoS:

```
bash
```

 Copy code

```
sudo sacctmgr add qos name=gpu_limit maxtresperuser=gpu=X
```



# Language models: the new “search engines”

**You**

In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?

**ChatGPT**

To restrict the number of GPUs a user's jobs can use at most within a specific job time range in a SLURM environment, you will need to utilize SLURM's Quality of Service (QoS) and account management features. This approach allows you to define limits on resource usage that apply to users, accounts, or even specific jobs based on various criteria, including job time ranges. Here's a step-by-step guide to setting up such restrictions:

## 1. Define or Update a Quality of Service (QoS)

First, you need to define a Quality of Service (QoS) that specifies the GPU usage limits. If a suitable QoS already exists, you may update it instead.

To create a new QoS:

```
bash
```

 Copy code

```
sudo sacctmgr add qos name=gpu_limit maxtresperuser=gpu=X
```

User-friendly presentations

# Language models: the new “search engines”

**You**

In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?

**ChatGPT**

To restrict the number of GPUs a user's jobs can use at most within a specific job time range in a SLURM environment, you will need to utilize SLURM's Quality of Service (QoS) and account management features. This approach allows you to define limits on resource usage that apply to users, accounts, or even specific jobs based on various criteria, including job time ranges. Here's a step-by-step guide to setting up such restrictions:

## 1. Define or Update a Quality of Service (QoS)

First, you need to define a Quality of Service (QoS) that specifies the GPU usage limits. If a suitable QoS already exists, you may update it instead.

To create a new QoS:

```
bash
```

 Copy code

```
sudo sacctmgr add qos name=gpu_limit maxtresperuser=gpu=X
```

User-friendly presentations  
Relevant answers

# Language models: the new “search engines”

**You**

In SLURM how can I restrict how many GPUs a user's jobs can use at most in a specific job time range?

**ChatGPT**

To restrict the number of GPUs a user's jobs can use at most within a specific job time range in a SLURM environment, you will need to utilize SLURM's Quality of Service (QoS) and account management features. This approach allows you to define limits on resource usage that apply to users, accounts, or even specific jobs based on various criteria, including job time ranges. Here's a step-by-step guide to setting up such restrictions:

## 1. Define or Update a Quality of Service (QoS)

First, you need to define a Quality of Service (QoS) that specifies the GPU usage limits. If a suitable QoS already exists, you may update it instead.

To create a new QoS:

```
bash
```

 Copy code

```
sudo sacctmgr add qos name=gpu_limit maxtresperuser=gpu=X
```

User-friendly presentations

Relevant answers

Boost productivity

# Language models: the new “search engines”



**You**

Where can I find gluten-free sushi (gluten free soy sauce provided) in Princeton, NJ?



# Language models: the new “search engines”



**You**

Where can I find gluten-free sushi (gluten free soy sauce provided) in Princeton, NJ?



**ChatGPT**

One option to find gluten-free sushi in Princeton, NJ, is to try Ajihei Sushi, located at 11 Chambers St, Princeton, NJ 08542. They offer a variety of sushi options and often accommodate dietary restrictions. It's recommended to call them directly to inquire about their gluten-free options and if they provide gluten-free soy sauce.

Another option is Sakura Japanese Sushi & Steak House, located at 4437 NJ-27, Princeton, NJ 08540. They typically offer a wide range of sushi options and may be able to accommodate gluten-free requests. Again, it's best to call ahead and ask about their gluten-free options and soy sauce availability.



# Language models: the new “search engines”



**You**

Where can I find gluten-free sushi (gluten free soy sauce provided) in Princeton, NJ?



**ChatGPT**

One option to find gluten-free sushi in Princeton, NJ, is to try Ajihei Sushi, located at 11 Chambers St, Princeton, NJ 08542. They offer a variety of sushi options and often accommodate dietary restrictions. It's recommended to call them directly to inquire about their gluten-free options and if they provide gluten-free soy sauce. **Does not provide gluten-free soy sauce**

Another option is Sakura Japanese Sushi & Steak House, located at 4437 NJ-27, Princeton, NJ 08540. They typically offer a wide range of sushi options and may be able to accommodate gluten-free requests. Again, it's best to call ahead and ask about their gluten-free options and soy sauce availability.

# Language models: the new “search engines”



You

Where can I find gluten-free sushi (gluten free soy sauce provided) in Princeton, NJ?



ChatGPT

One option to find gluten-free sushi in Princeton, NJ, is to try ~~Ajihei Sushi~~, located at 11 Chambers St, Princeton, NJ 08542. They offer a variety of sushi options and often accommodate dietary restrictions. It's recommended to call them directly to inquire about their gluten-free options and if they provide gluten-free soy sauce. **Does not provide gluten-free soy sauce**

Another option is ~~Sakura Japanese Sushi & Steak House~~, located at 4437 NJ-27, Princeton, NJ 08540. They typically offer a wide range of sushi options and may be able to accommodate gluten-free requests. Again, it's best to call ahead and ask about their gluten-free options and soy sauce availability. **Does not exist?!**

# Language models: the new “search engines”



You

Where can I find gluten-free sushi (gluten free soy sauce provided) in Princeton, NJ?



ChatGPT

One option to find gluten-free sushi in Princeton, NJ, is to try Ajihei Sushi, located at 11 Chambers St, Princeton, NJ 08542. They offer a variety of sushi options and often accommodate dietary restrictions. It's recommended to call them directly to inquire about their gluten-free options and if they provide gluten-free soy sauce. **Does not provide gluten-free soy sauce**

Another option is Sakura Japanese Sushi & Steak House, located at 4437 NJ-27, Princeton, NJ 08540. They typically offer a wide range of sushi options and may be able to accommodate gluten-free requests. Again, it's best to call ahead and ask about their gluten-free options and soy sauce availability. **Does not exist?!**

“hallucination”



# Language models: the new “search engines”



You

Where can I find gluten-free sushi (gluten free soy sauce provided) in Princeton, NJ?



ChatGPT

One option to find gluten-free sushi in Princeton, NJ, is to try Ajihei Sushi, located at 11 Chambers St, Princeton, NJ 08542. They offer a variety of sushi options and often accommodate dietary restrictions. It's recommended to call them directly to inquire about their gluten-free options and if they provide gluten-free soy sauce. **Does not provide gluten-free soy sauce**

Another option is Sakura Japanese Sushi & Steak House, located at 4437 NJ-27, Princeton, NJ 08540. They typically offer a wide range of sushi options and may be able to accommodate gluten-free requests. Again, it's best to call ahead and ask about their gluten-free options and soy sauce availability. **Does not exist?!**

“hallucination”

No references for users to verify

# Language models: the new “search engines”



# Language models: the new “search engines”

How can we make LMs better information-seeking tools?

# Language models: the new “search engines”

How can we make LMs better information-seeking tools?

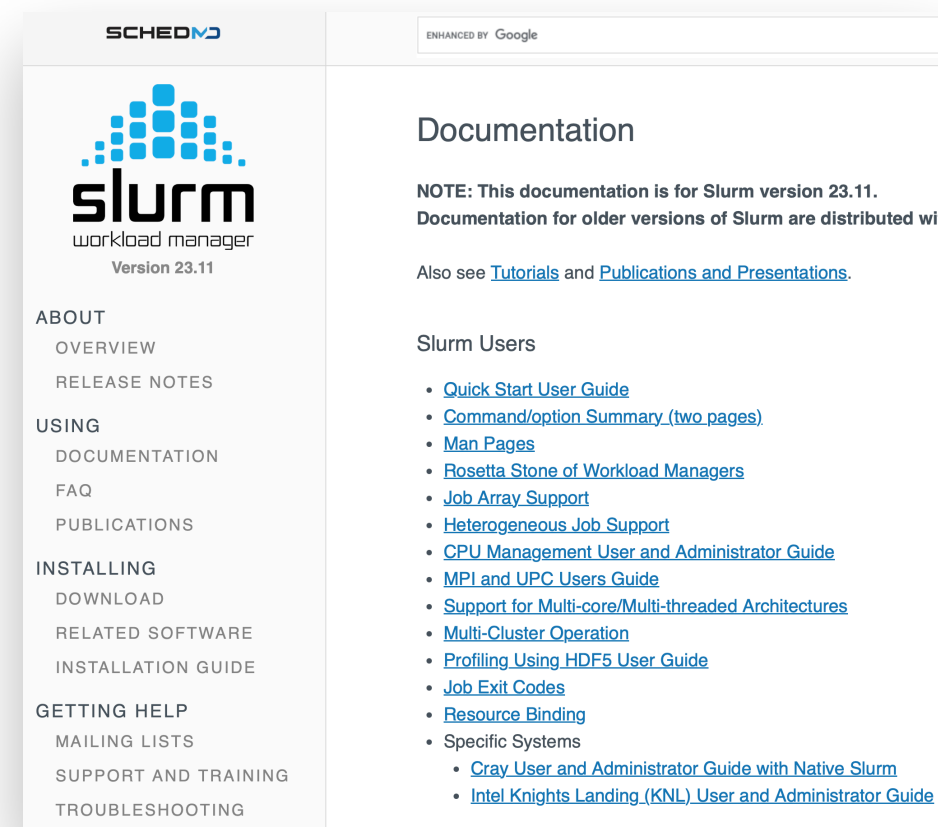
**1. The ability to find and utilize reference materials**

# Language models: the new “search engines”

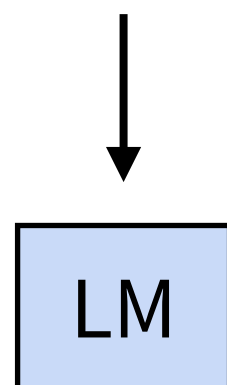
## How can we make LMs better information-seeking tools?

### 1. The ability to find and utilize reference materials

*Document-augmented*



Q: In SLURM how can i restrict ...

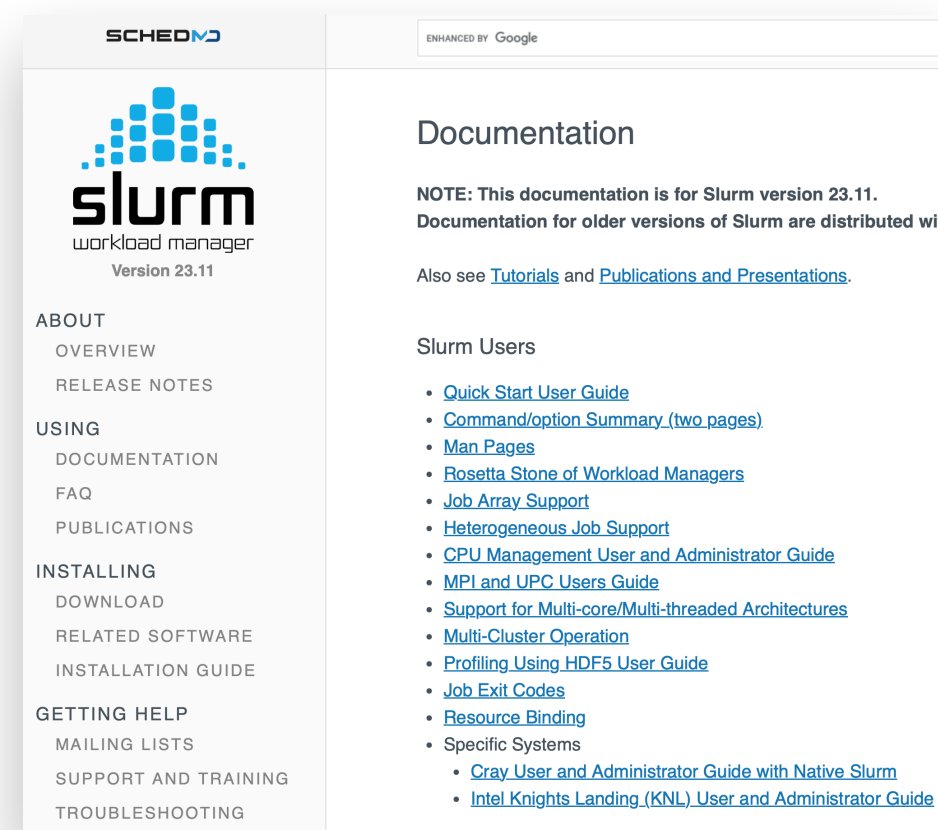


# Language models: the new “search engines”

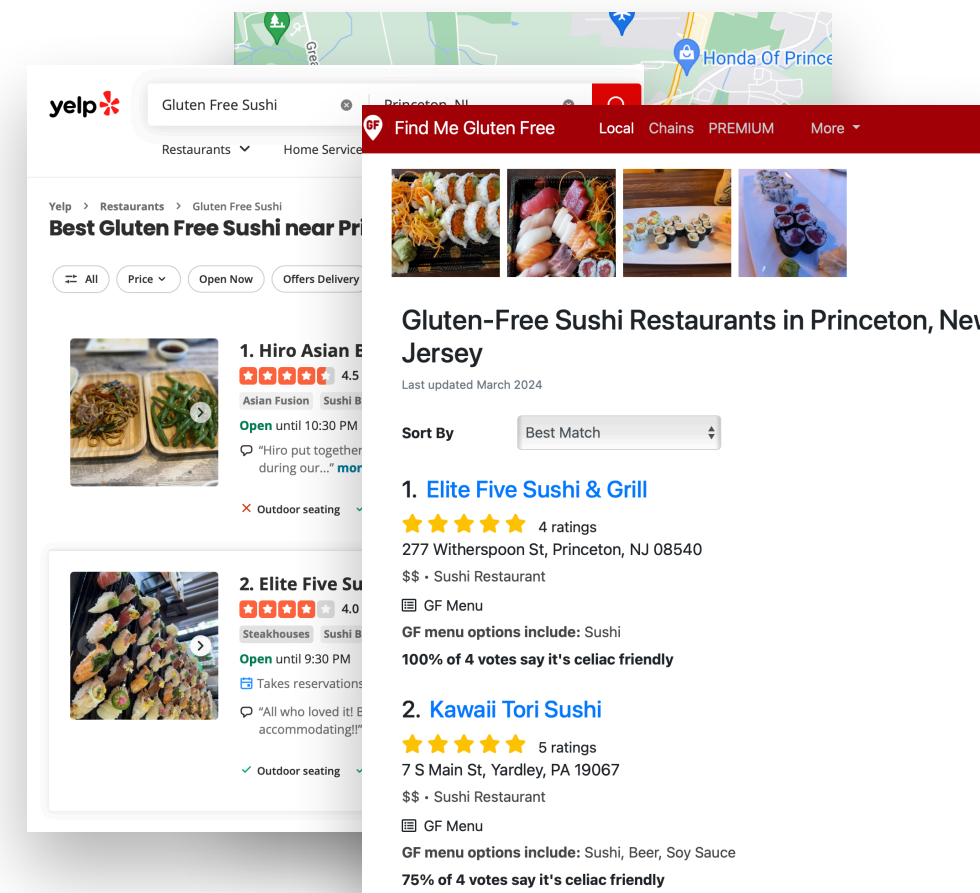
## How can we make LMs better information-seeking tools?

### 1. The ability to find and utilize reference materials

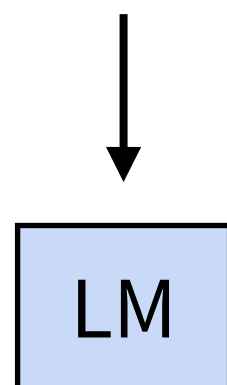
*Document-augmented*



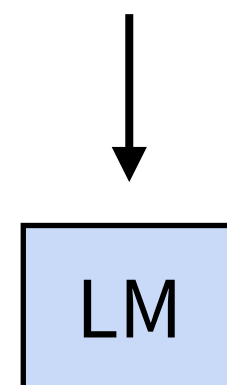
*Retrieval-augmented*



Q: In SLURM how can i restrict ...



Q: Where can I find gluten-free sushi ...





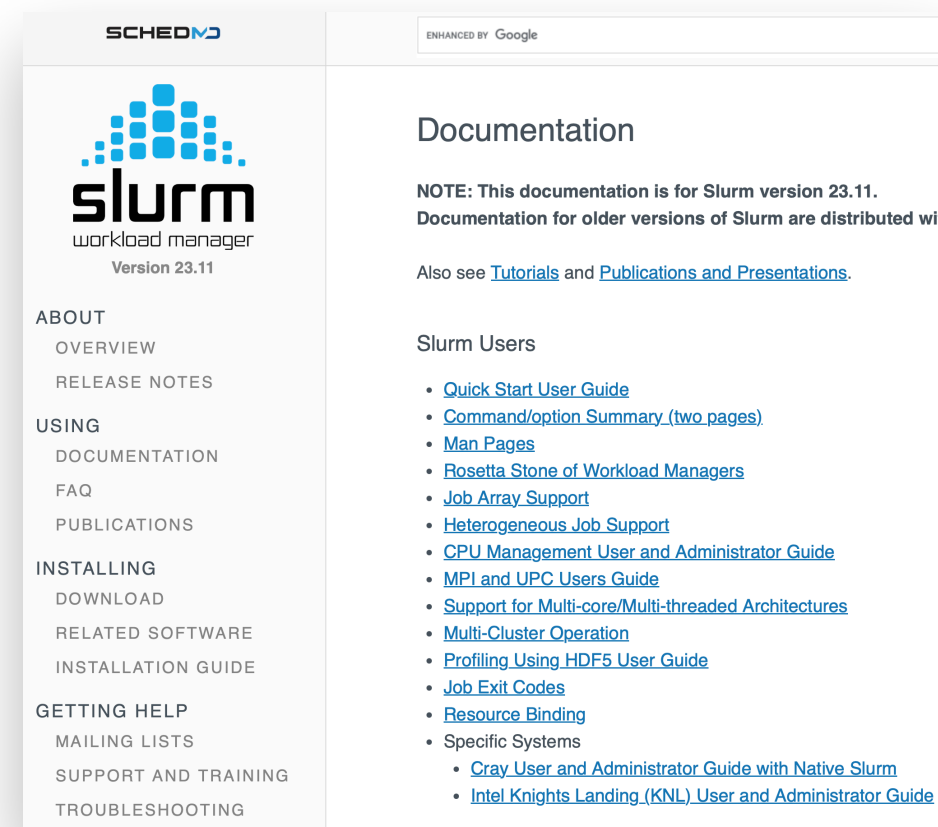
# Language models: the new “search engines”

## How can we make LMs better information-seeking tools?

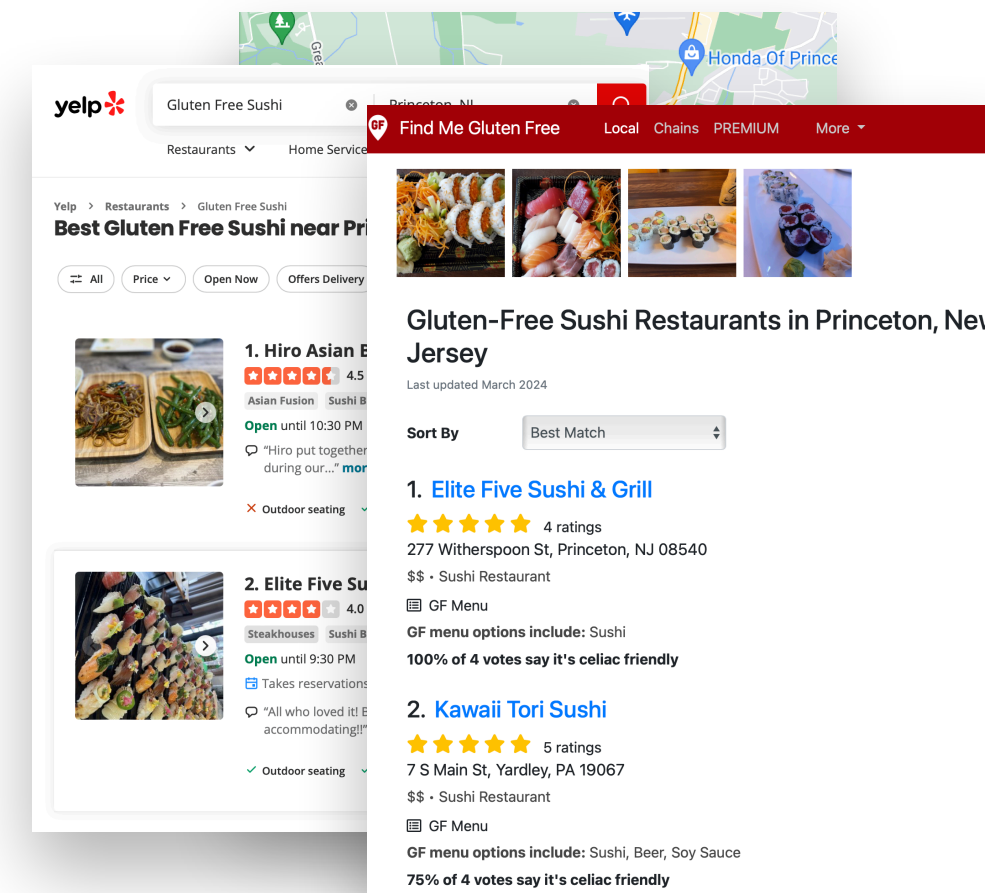
### 1. The ability to find and utilize reference materials

### 2. The ability to provide “citations”

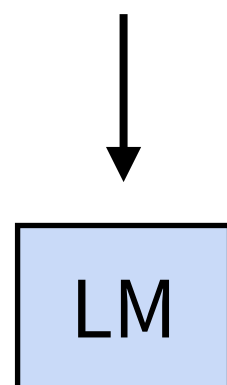
*Document-augmented*



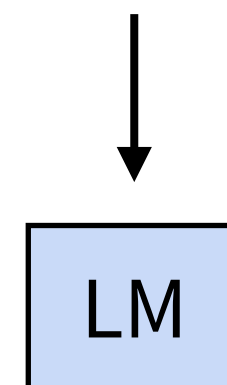
*Retrieval-augmented*



Q: In SLURM how can i restrict ...



Q: Where can I find gluten-free sushi ...

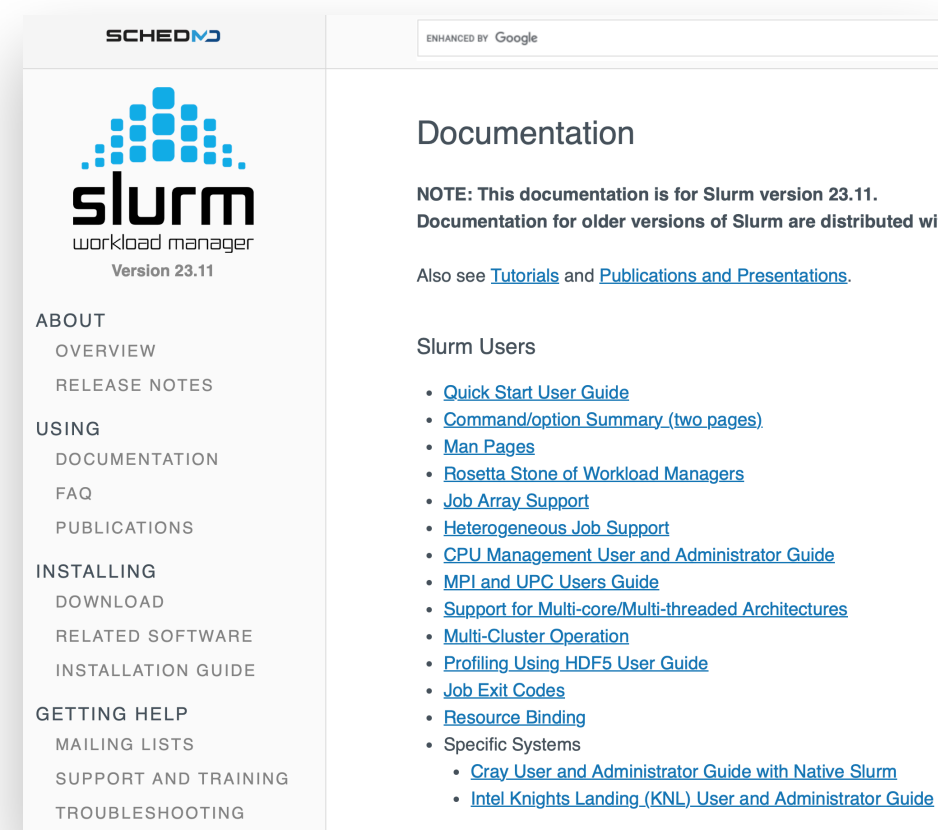


# Language models: the new “search engines”

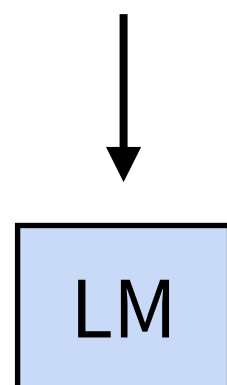
## How can we make LMs better information-seeking tools?

### 1. The ability to find and utilize reference materials

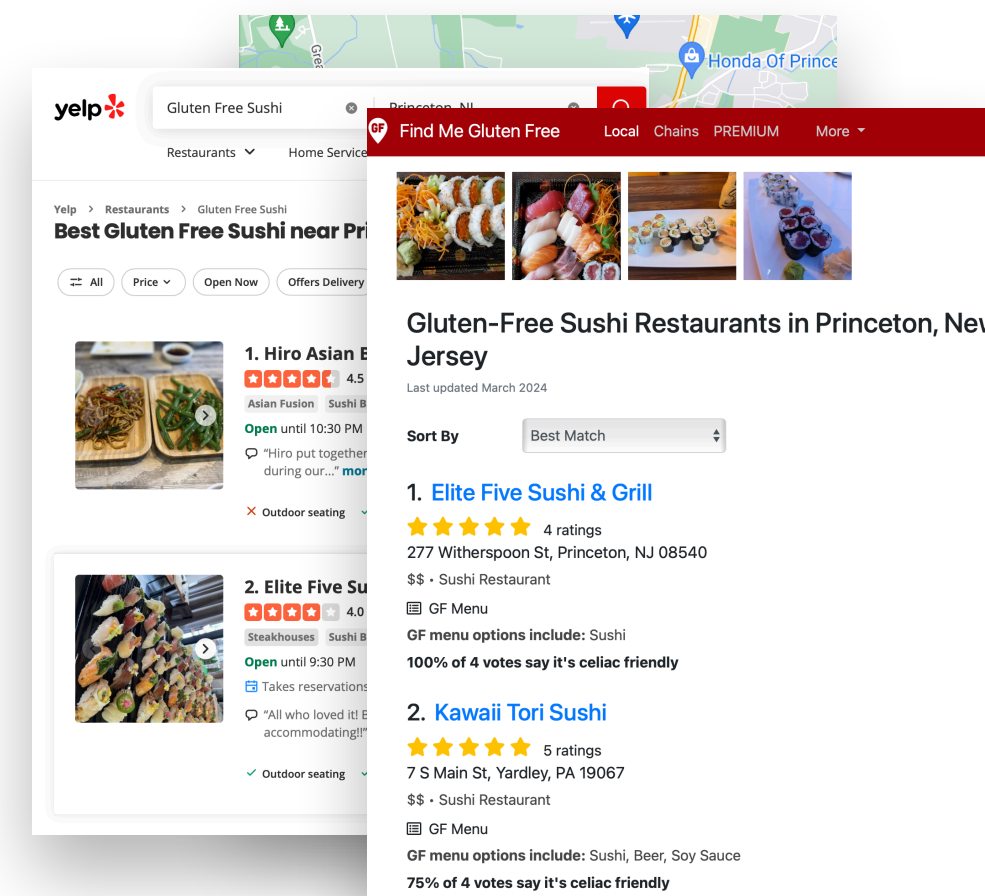
*Document-augmented*



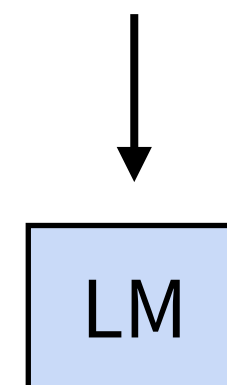
Q: In SLURM how can i restrict ...



*Retrieval-augmented*



Q: Where can I find gluten-free sushi ...



### 2. The ability to provide “citations”

Where can I find gluten-free sushi (gluten free soy sauce provided) in Princeton, NJ?

Sources

Gluten-Free Sushi Restaurants in... findme glutenf... · 1	TOP 10 BEST Gluten Free Sushi near... yelp · 2	Elite Five Sushi & Grill Gluten-Free -... findme glutenf... · 3	View 2 more
--	---	--	-------------

Answer

You can find gluten-free sushi with gluten-free soy sauce provided in Princeton, NJ at Elite Five Sushi & Grill. This restaurant offers gluten-free options, including sushi, and is known for taking celiac allergies seriously, providing GF soy sauce, and having knowledgeable staff about gluten-free entrees <sup>3</sup> <sup>1</sup>.

# Language models: the new “search engines”

# Language models: the new “search engines”

**Challenge #1: the need for better retrieval/search**



# Language models: the new “search engines”

## Challenge #1: the need for better retrieval/search

Traditional word-matching methods do not work well for “semantic search”

# Language models: the new “search engines”

## Challenge #1: the need for better retrieval/search

Traditional word-matching methods do not work well for “semantic search”

Query  
Pre-trained language models are great.

BERT obtains new SoTA results on 11 NLP tasks. (Devlin et al., 2019)

# Language models: the new “search engines”

## Challenge #1: the need for better retrieval/search

Traditional word-matching methods do not work well for “semantic search”

Query  
Pre-trained language models are great.

BERT obtains new SoTA results on 11 NLP tasks. (Devlin et al., 2019)

*Dense retrieval* becomes more and more popular

# Language models: the new “search engines”

## Challenge #1: the need for better retrieval/search

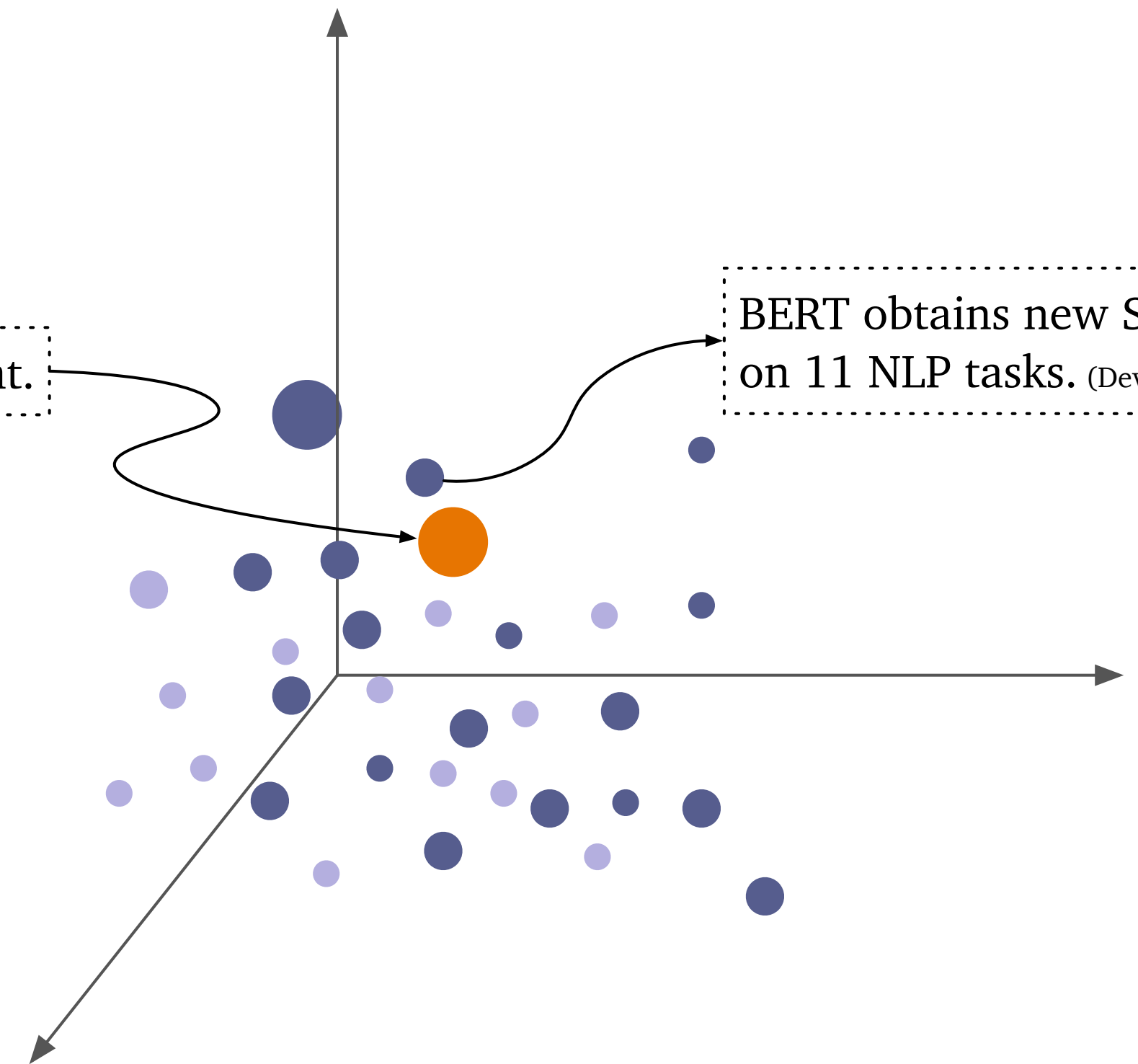
Traditional word-matching methods do not work well for “semantic search”

*Dense retrieval* becomes more and more popular

- Represent sentences/paragraphs/documents as *vectors* and perform nearest-neighbor search

Query  
Pre-trained language models are great.

BERT obtains new SoTA results on 11 NLP tasks. (Devlin et al., 2019)





# Language models: the new “search engines”

**Challenge #1: the need for better retrieval/search**

# Language models: the new “search engines”

## Challenge #1: the need for better retrieval/search

**SimCSE: Simple Contrastive Learning of Sentence Embeddings**

EMNLP 2021



# Language models: the new “search engines”

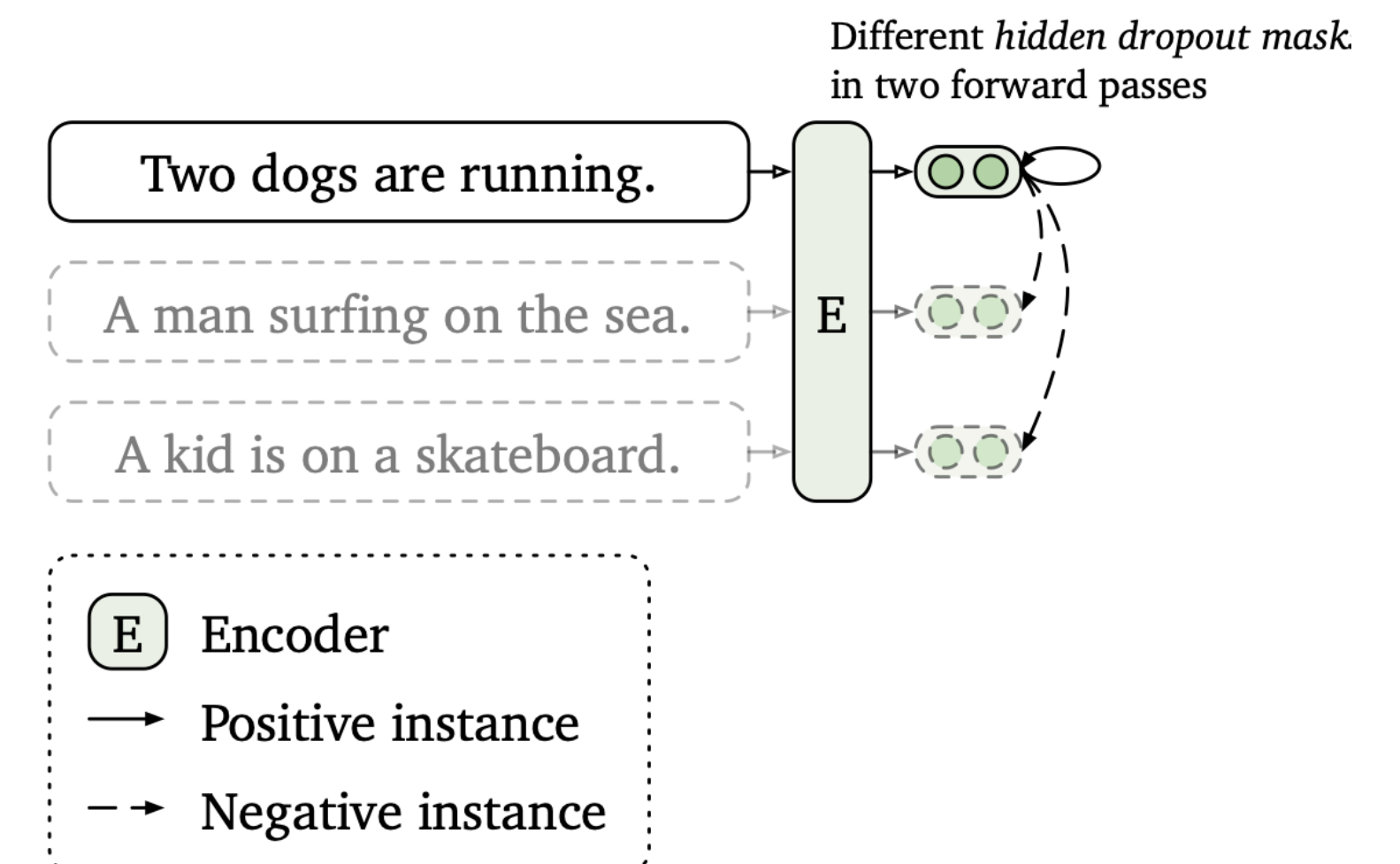
## Challenge #1: the need for better retrieval/search

### SimCSE: Simple Contrastive Learning of Sentence Embeddings

EMNLP 2021



- Propose a simple *contrastive learning* framework for sentence embeddings



# Language models: the new “search engines”

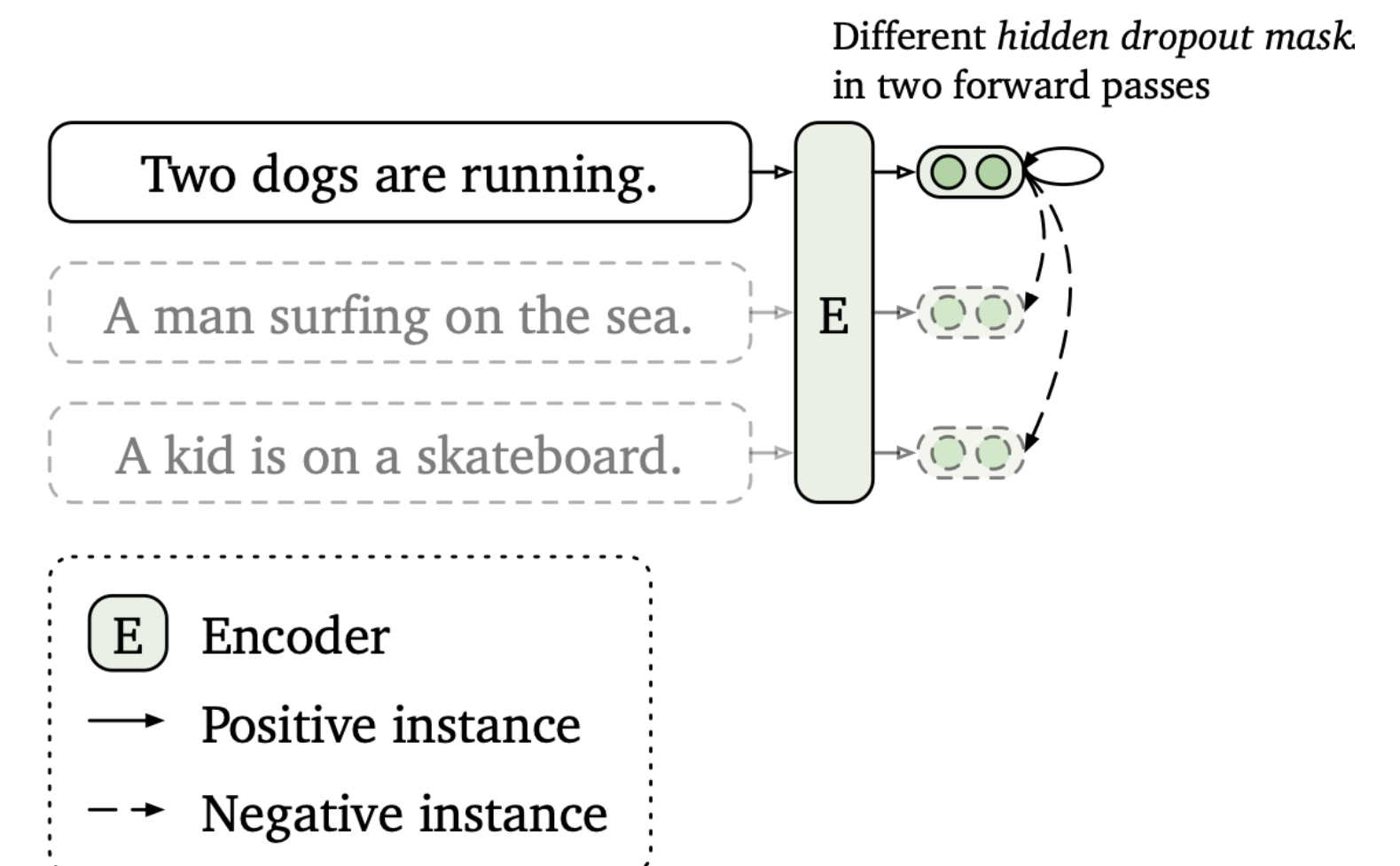
## Challenge #1: the need for better retrieval/search

### SimCSE: Simple Contrastive Learning of Sentence Embeddings

EMNLP 2021



- Propose a simple *contrastive learning* framework for sentence embeddings
  - A technique used by SoTA embedding tools (OpenAI, 2022; Su et al., 2023; Muennighoff et al., 2024)



OpenAI, 2022. Text and Code Embeddings by Contrastive Pre-Training.  
Su et al., 2023. One Embedder, Any Task: Instruction-Finetuned Text Embeddings.  
Muennighoff et al., 2024. Generative Representational Instruction Tuning.



# Language models: the new “search engines”

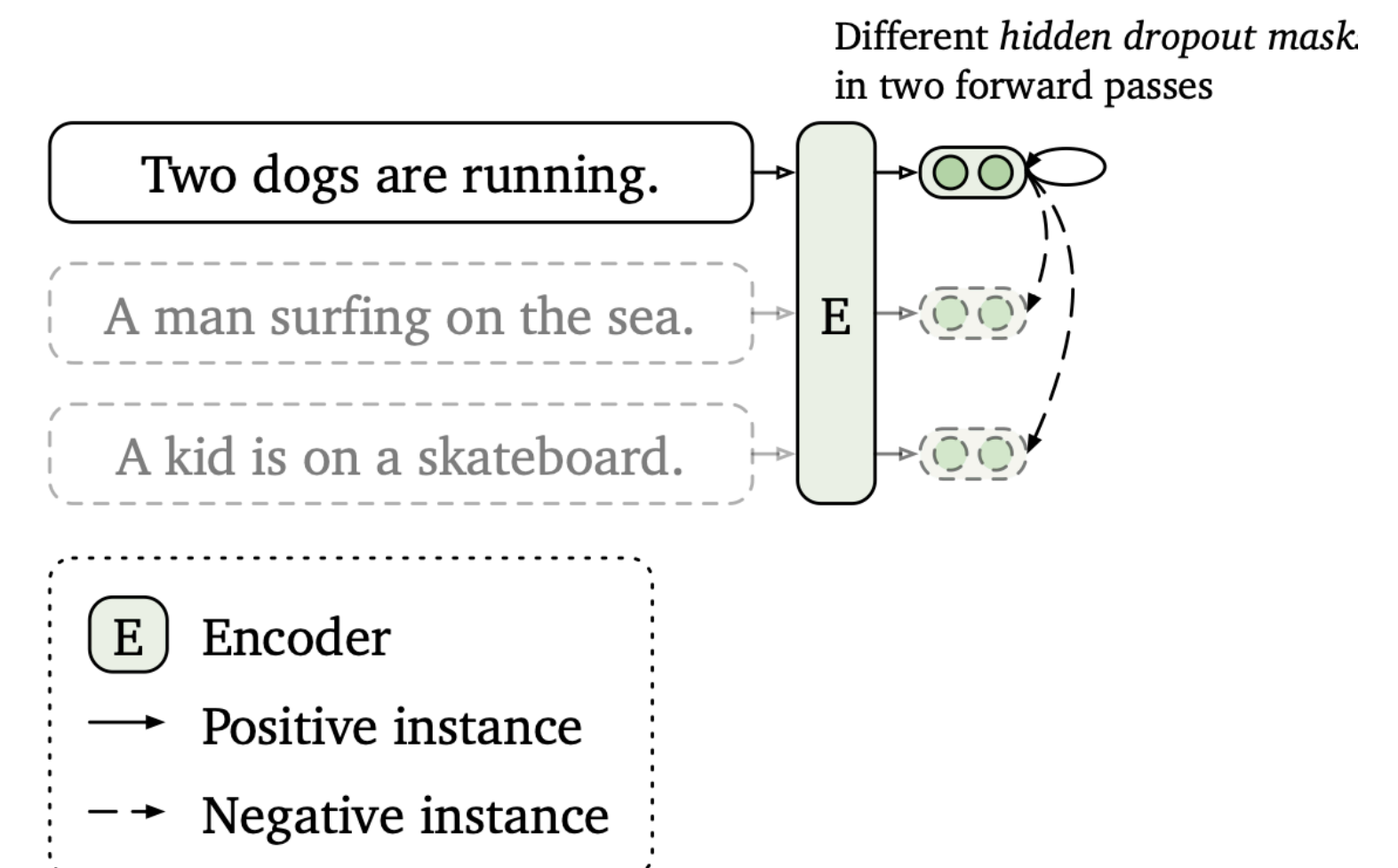
## Challenge #1: the need for better retrieval/search

### SimCSE: Simple Contrastive Learning of Sentence Embeddings

EMNLP 2021



- Propose a simple *contrastive learning* framework for sentence embeddings
  - A technique used by SoTA embedding tools (OpenAI, 2022; Su et al., 2023; Muennighoff et al., 2024)
- Downloaded by >18M times (from HuggingFace)



OpenAI, 2022. Text and Code Embeddings by Contrastive Pre-Training.  
Su et al., 2023. One Embedder, Any Task: Instruction-Finetuned Text Embeddings.  
Muennighoff et al., 2024. Generative Representational Instruction Tuning.

# Language models: the new “search engines”

# Language models: the new “search engines”

**Challenge #2: how to evaluate?**

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

- There is no good evaluation for long text generation



# Language models: the new “search engines”

## Challenge #2: how to evaluate?

- There is no good evaluation for long text generation
- There is no automatic evaluation for generations with citations

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

- There is no good evaluation for long text generation
- There is no automatic evaluation for generations with citations



Human evaluation?

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

- There is no good evaluation for long text generation
- There is no automatic evaluation for generations with citations



Human evaluation? Slow, costly, unreliable



# Language models: the new “search engines”

**Challenge #2: how to evaluate?**

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

**Enabling Large Language Models to Generate Text with Citations**

EMNLP 2023





# Language models: the new “search engines”

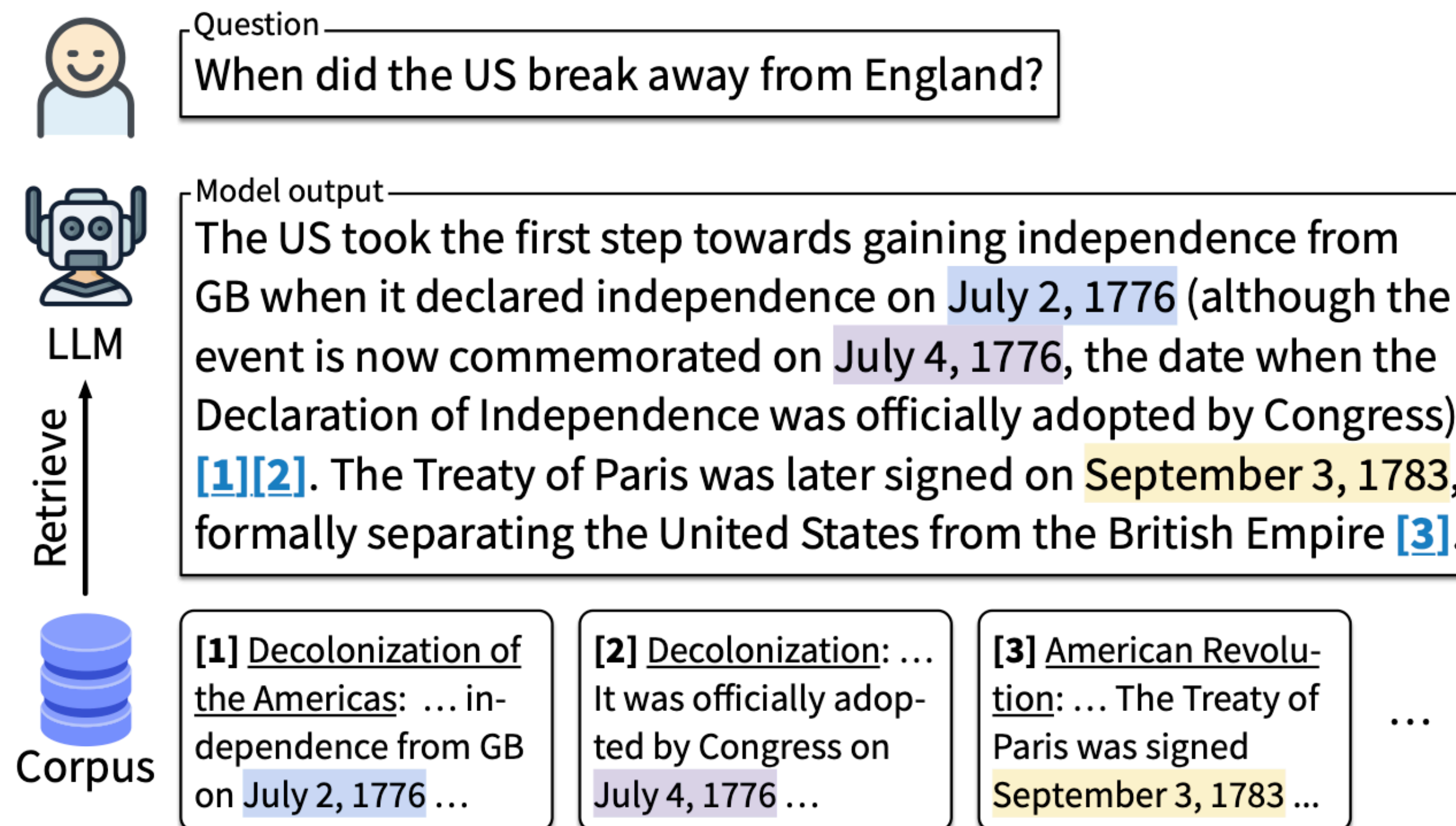
## Challenge #2: how to evaluate?

### Enabling Large Language Models to Generate Text with Citations

EMNLP 2023



### ALCE: automatic LLM citation evaluation



# Language models: the new “search engines”

## Challenge #2: how to evaluate?

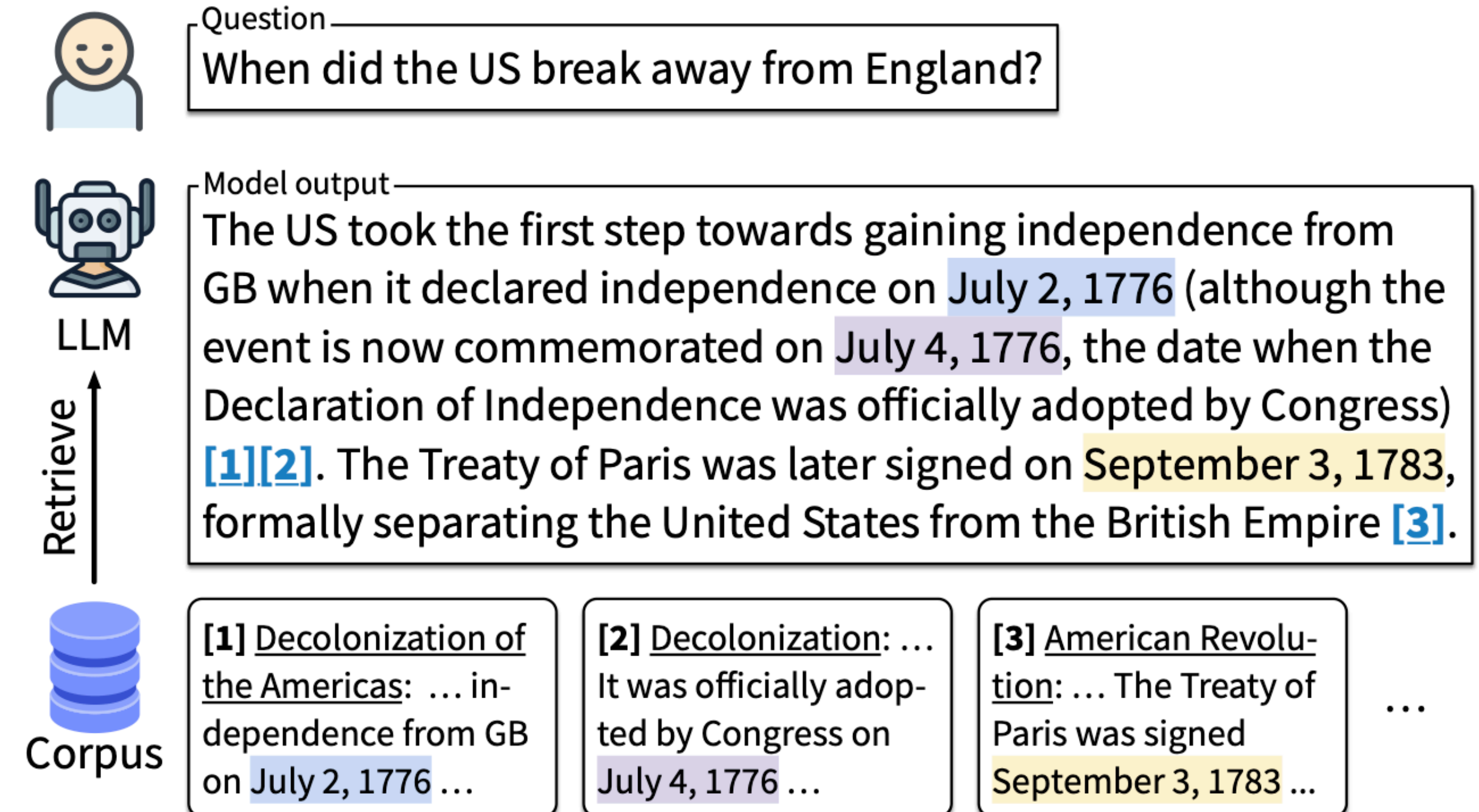
### Enabling Large Language Models to Generate Text with Citations

EMNLP 2023



#### ALCE: automatic LLM citation evaluation

- Given a **question**





# Language models: the new “search engines”

## Challenge #2: how to evaluate?

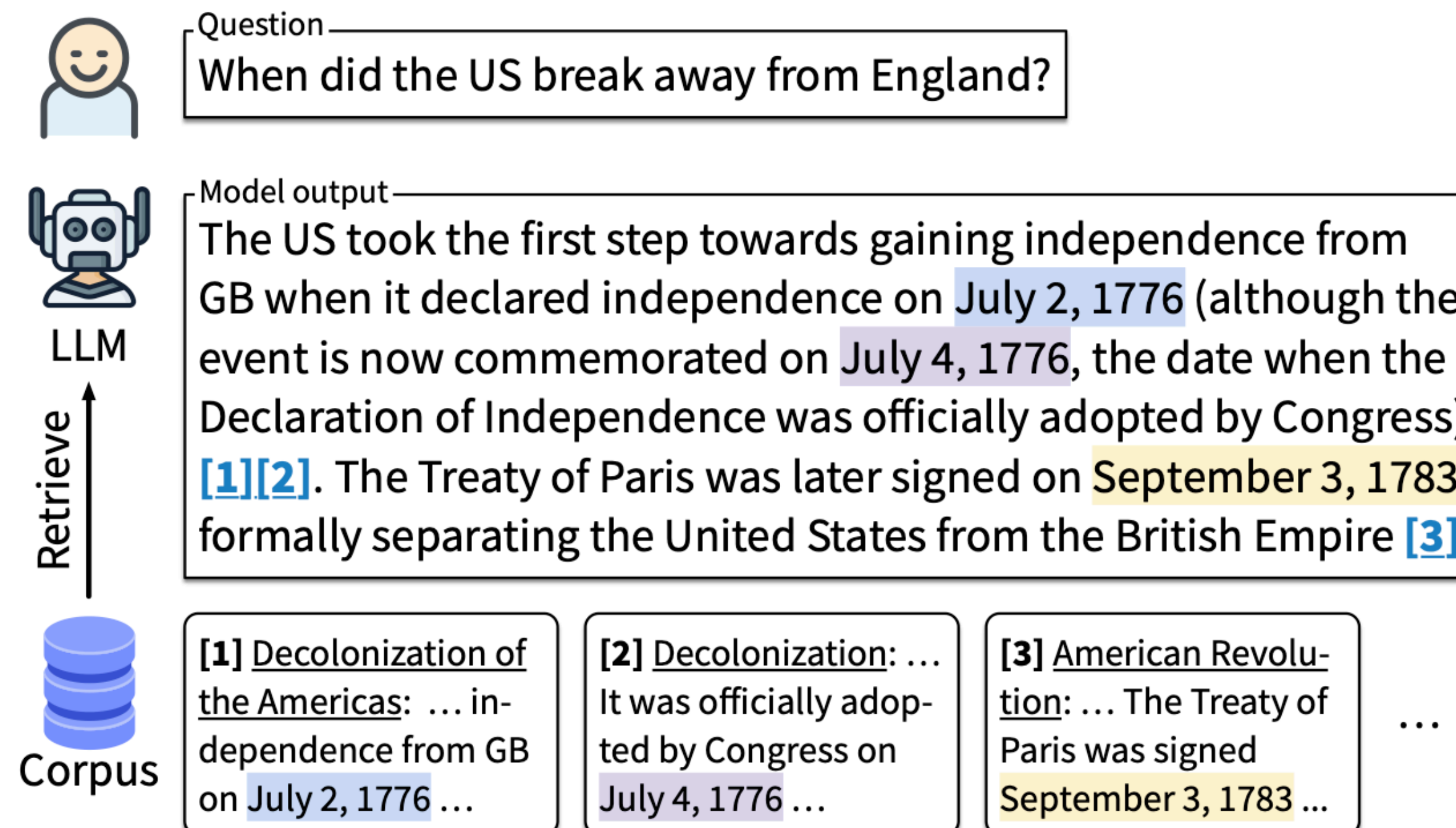
### Enabling Large Language Models to Generate Text with Citations

EMNLP 2023



#### ALCE: automatic LLM citation evaluation

- Given a **question**
- Given a **corpus**



# Language models: the new “search engines”

## Challenge #2: how to evaluate?

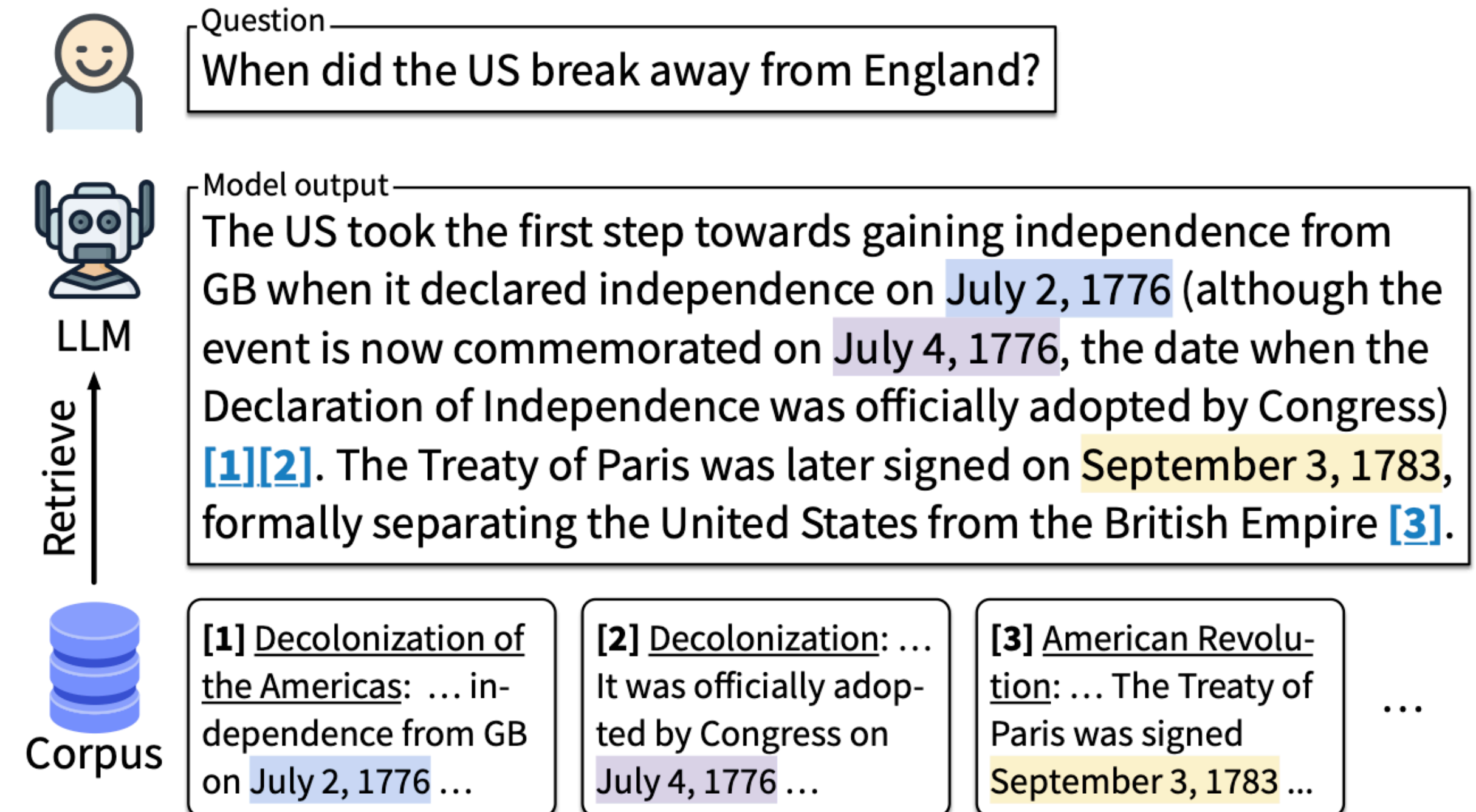
### Enabling Large Language Models to Generate Text with Citations

EMNLP 2023



#### ALCE: automatic LLM citation evaluation

- Given a **question**
- Given a **corpus**
- Requires end-to-end systems to





# Language models: the new “search engines”

## Challenge #2: how to evaluate?

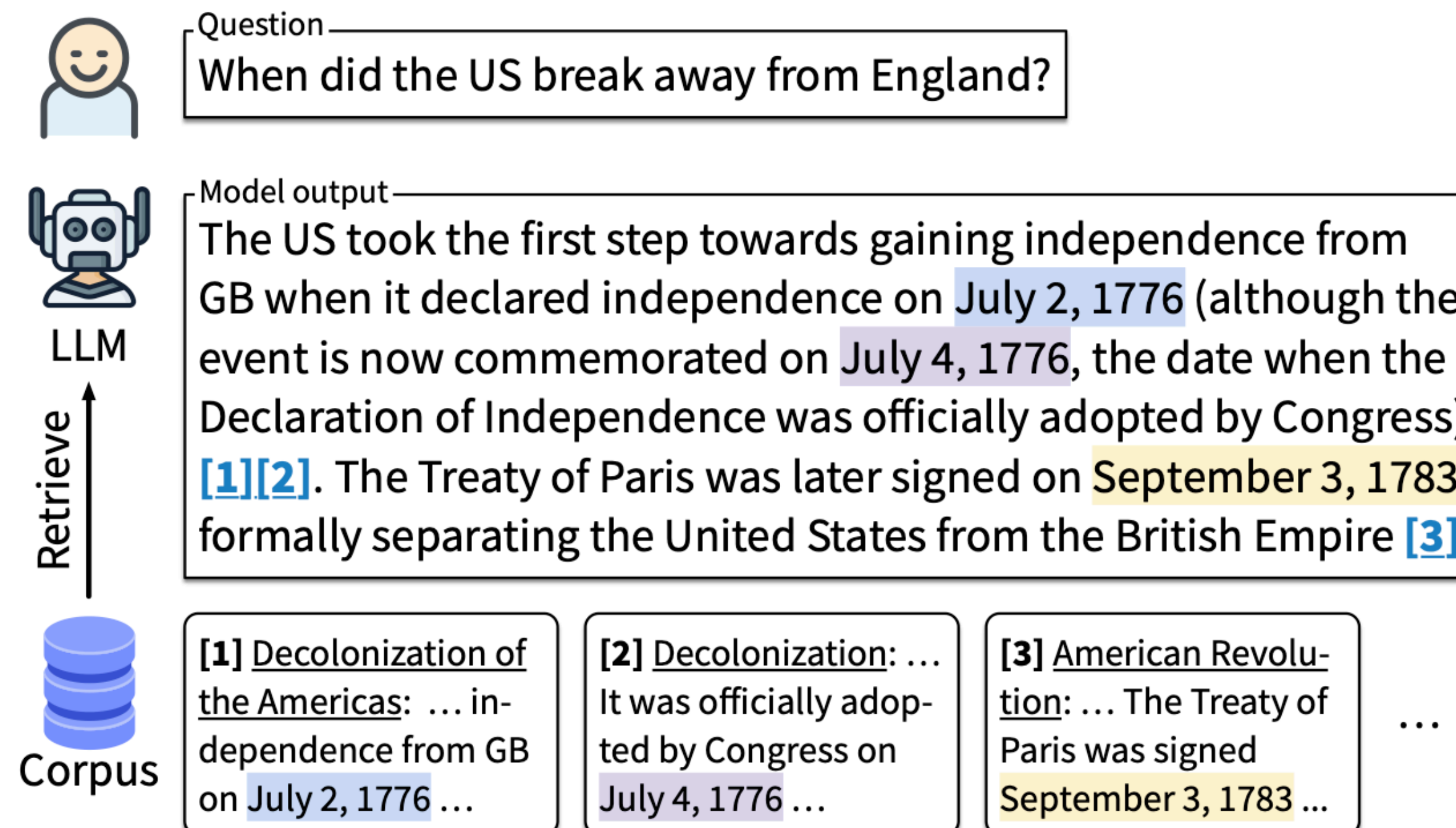
### Enabling Large Language Models to Generate Text with Citations

EMNLP 2023



#### ALCE: automatic LLM citation evaluation

- Given a **question**
- Given a **corpus**
- Requires end-to-end systems to
  - **Retrieve** passages from the corpus





# Language models: the new “search engines”

## Challenge #2: how to evaluate?

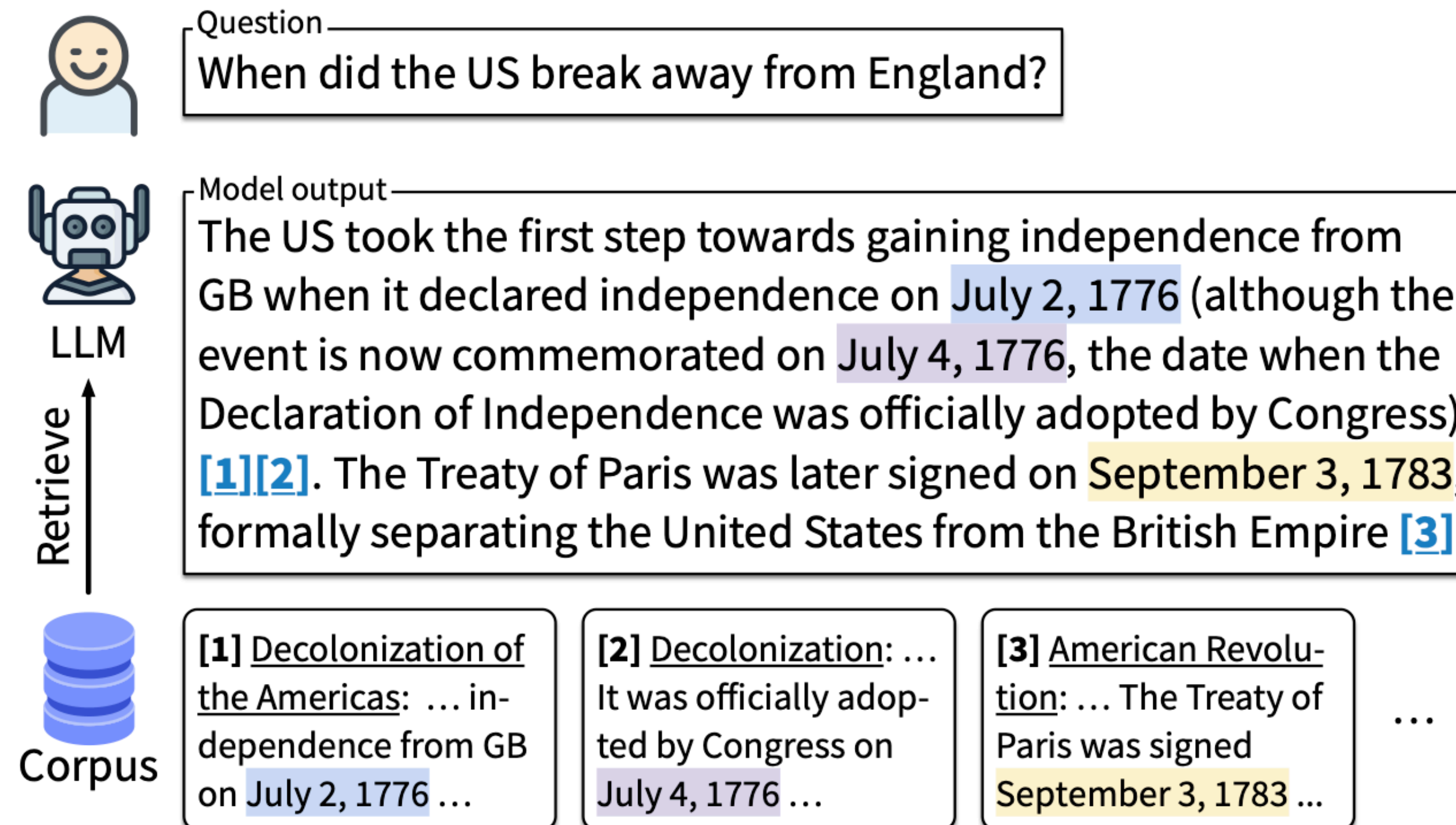
### Enabling Large Language Models to Generate Text with Citations

EMNLP 2023



#### ALCE: automatic LLM citation evaluation

- Given a **question**
- Given a **corpus**
- Requires end-to-end systems to
  - **Retrieve** passages from the corpus
  - **Generate** long-text responses



# Language models: the new “search engines”

## Challenge #2: how to evaluate?

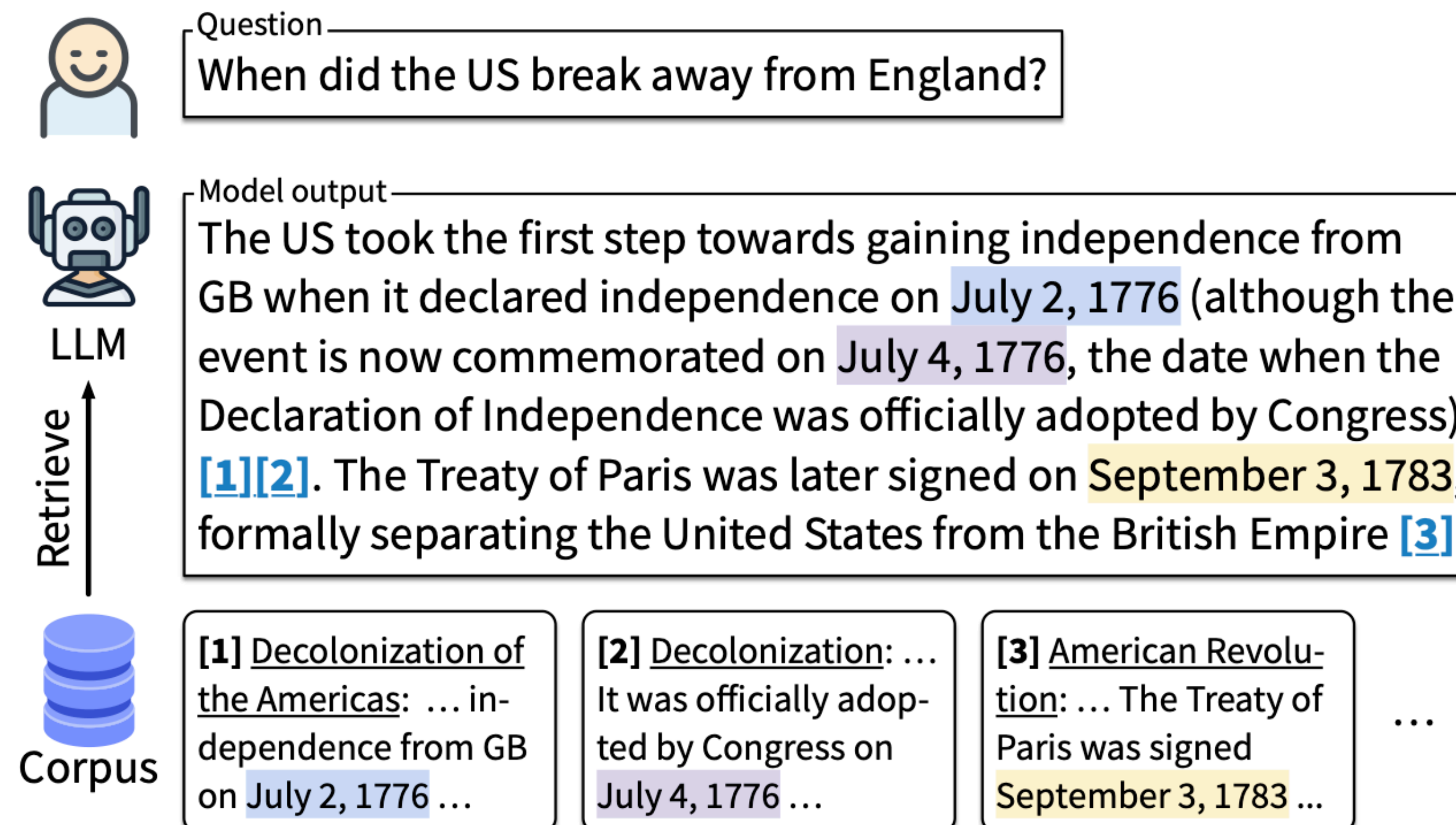
### Enabling Large Language Models to Generate Text with Citations

EMNLP 2023



#### ALCE: automatic LLM citation evaluation

- Given a **question**
- Given a **corpus**
- Requires end-to-end systems to
  - **Retrieve** passages from the corpus
  - **Generate** long-text responses
  - **Cite** supporting passages



# Language models: the new “search engines”

**Challenge #2: how to evaluate?**

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

**ALCE** automatically evaluates a long-text answer’s

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

**ALCE** automatically evaluates a long-text answer's

- Fluency



# Language models: the new “search engines”

## Challenge #2: how to evaluate?

**ALCE** automatically evaluates a long-text answer's

- Fluency
- Correctness

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

**ALCE** automatically evaluates a long-text answer’s

- Fluency
- Correctness
- Citation quality

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

**ALCE** automatically evaluates a long-text answer’s

- Fluency
- Correctness
- Citation quality

Question—  
When did the US break away from England?

Model output—  
... it declared independence on July 3, 1776 ... The Treaty of Paris was signed on September 3, 1783 ...

Short answers (from the dataset)  
July 2, 1776 ❌ July 4, 1776 ❌ September 3, 1783 ✅

String exact match  
Recall=33.3%

An example for correctness evaluation.

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

**ALCE** automatically evaluates a long-text answer’s

- Fluency
- Correctness
- Citation quality

Question: When did the US break away from England?

Model output: ... it declared independence on July 3, 1776 ... The Treaty of Paris was signed on September 3, 1783 ...

Short answers (from the dataset): July 2, 1776 ❌ July 4, 1776 ❌ September 3, 1783 ✅

String exact match Recall=33.3%

An example for correctness evaluation.

Question: When did the US break away from England?

Model output: {statement 1} [1][2]. {statement 2} [3]. {statement 3} [2][4][5].

{statement 1} [1][2] ✅

Premise: passage [1][2] → NLI model → “entailment”  
Hypothesis: {statement 1}

An example for citation evaluation

# Language models: the new “search engines”

## Challenge #2: how to evaluate?

**ALCE** automatically evaluates a long-text answer’s

- Fluency
- Correctness
- Citation quality

We show that even **GPT-4** lacks complete citation support **50%** of the times

Question: When did the US break away from England?

Model output: ... it declared independence on July 3, 1776 ... The Treaty of Paris was signed on September 3, 1783 ...

Short answers (from the dataset): July 2, 1776 ❌ July 4, 1776 ❌ September 3, 1783 ✅

String exact match Recall=33.3%

An example for correctness evaluation.

Question: When did the US break away from England?

Model output: {statement 1} [1][2]. {statement 2} [3]. {statement 3} [2][4][5].

{statement 1} [1][2] ✅

Premise: passage [1][2] → NLI model → “entailment”  
Hypothesis: {statement 1}

An example for citation evaluation



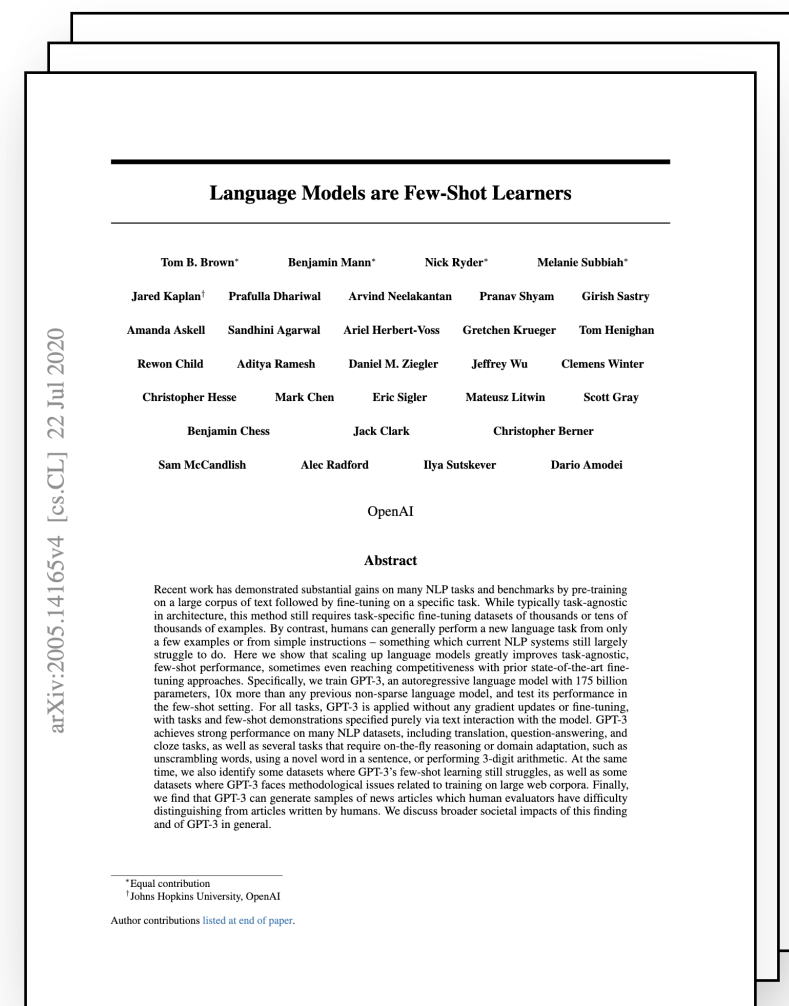
# Language models: the new “search engines”

# Language models: the new “search engines”

**Challenge #3: how to fit the context**

# Language models: the new “search engines”

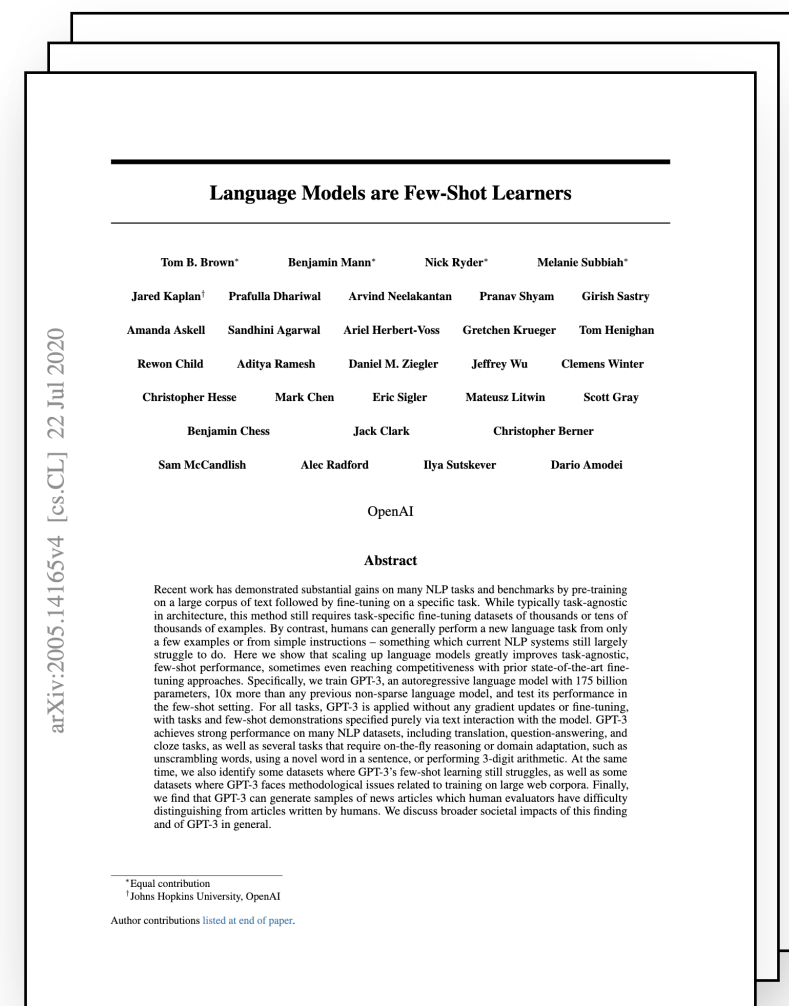
## Challenge #3: how to fit the context



The GPT-3 paper  
(~**75K** tokens)

# Language models: the new “search engines”

## Challenge #3: how to fit the context



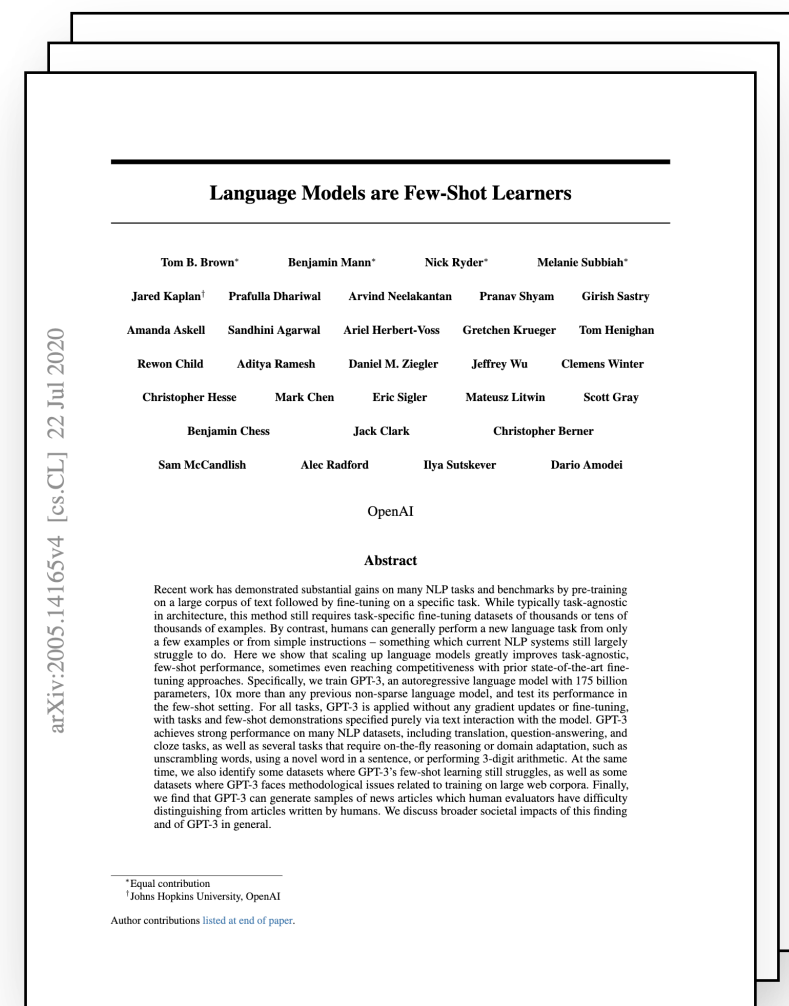
The GPT-3 paper  
(~**75K** tokens)

The Dune series  
(~**1M** words)



# Language models: the new “search engines”

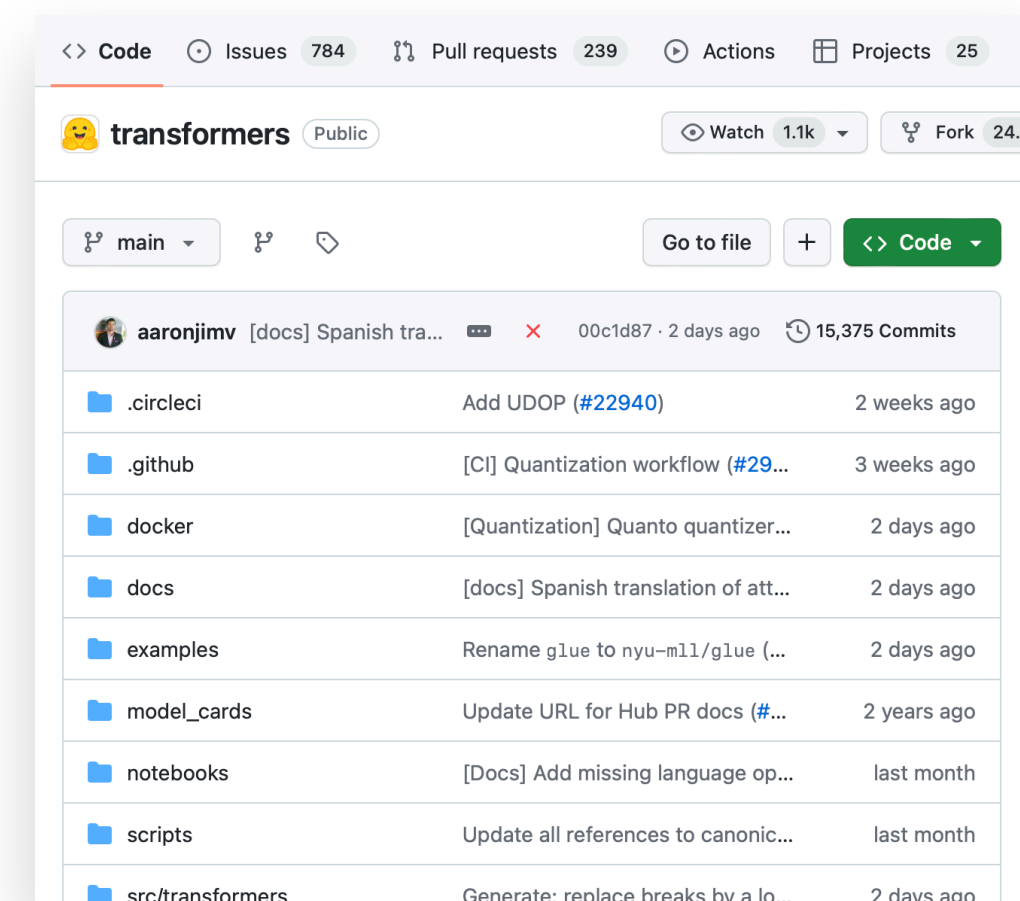
## Challenge #3: how to fit the context



The GPT-3 paper  
(~**75K** tokens)



The Dune series  
(~**1M** words)

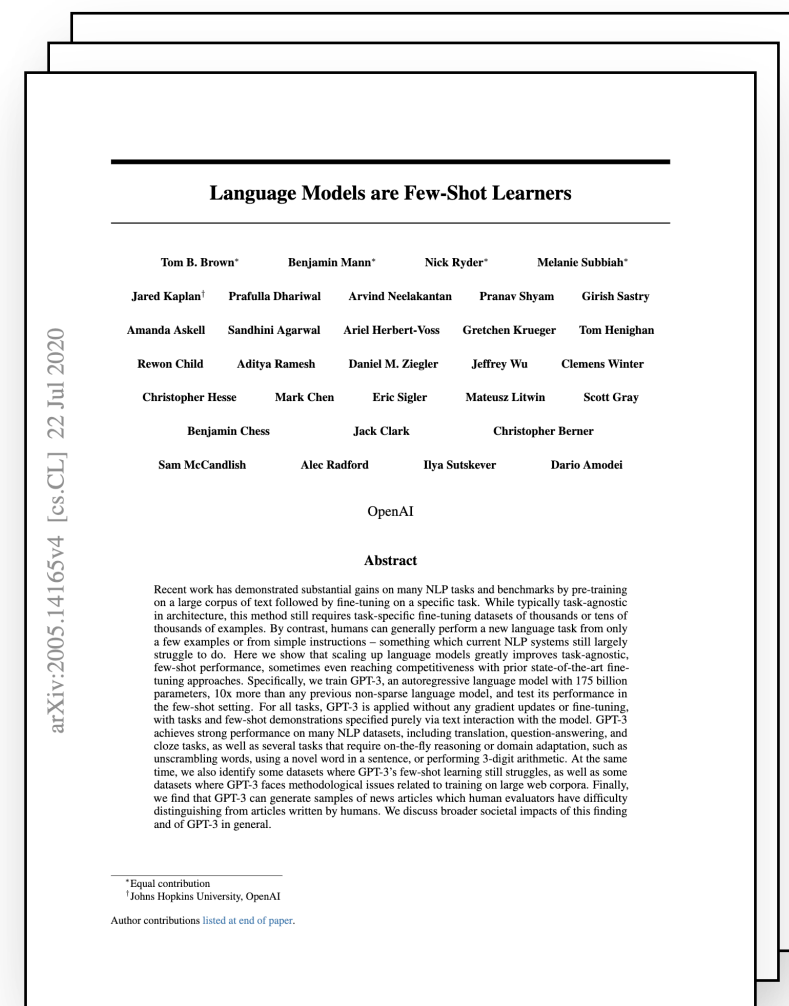


The Transformers package  
(~**10M** tokens)



# Language models: the new “search engines”

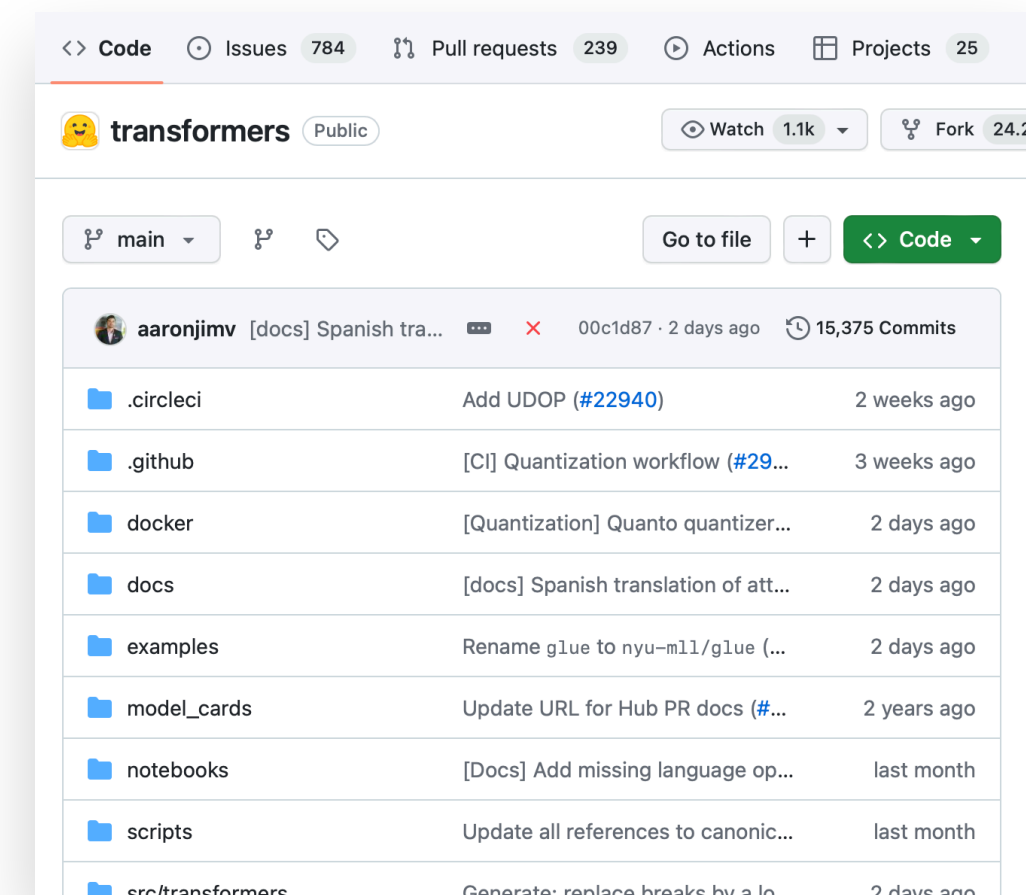
## Challenge #3: how to fit the context



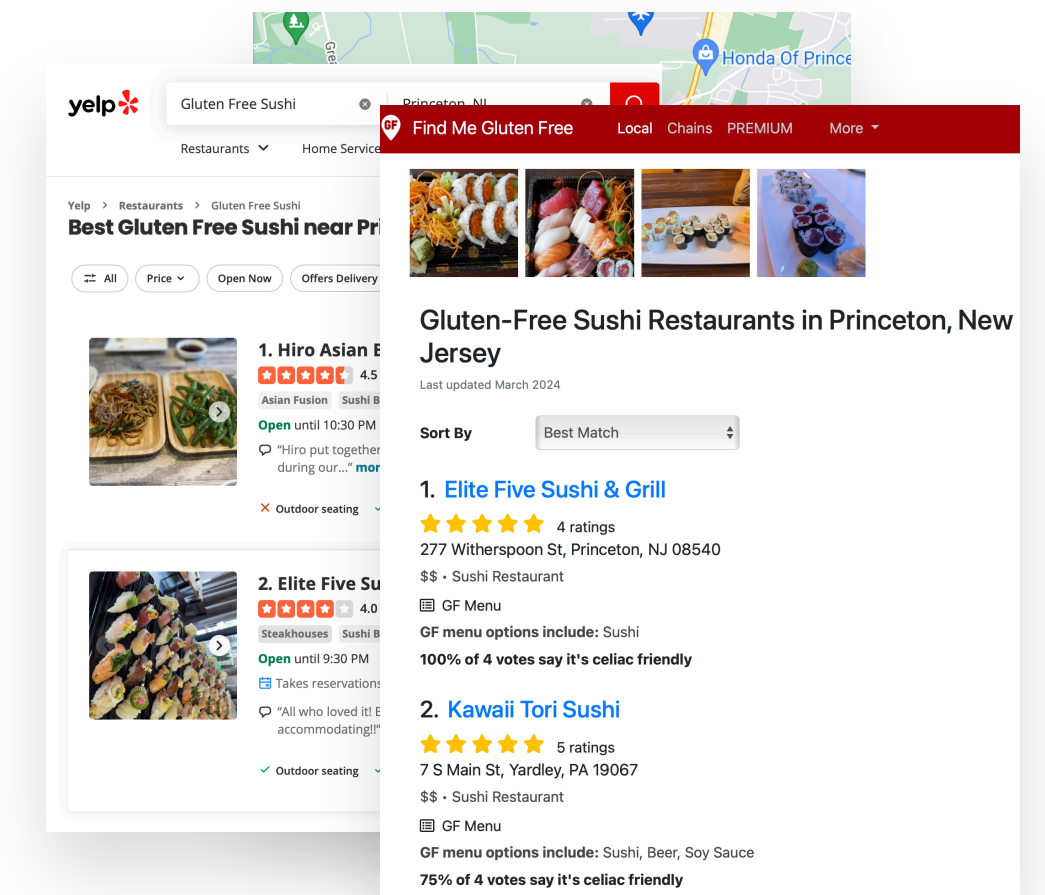
The GPT-3 paper  
(~**75K** tokens)



The Dune series  
(~**1M** words)



The Transformers package  
(~**10M** tokens)

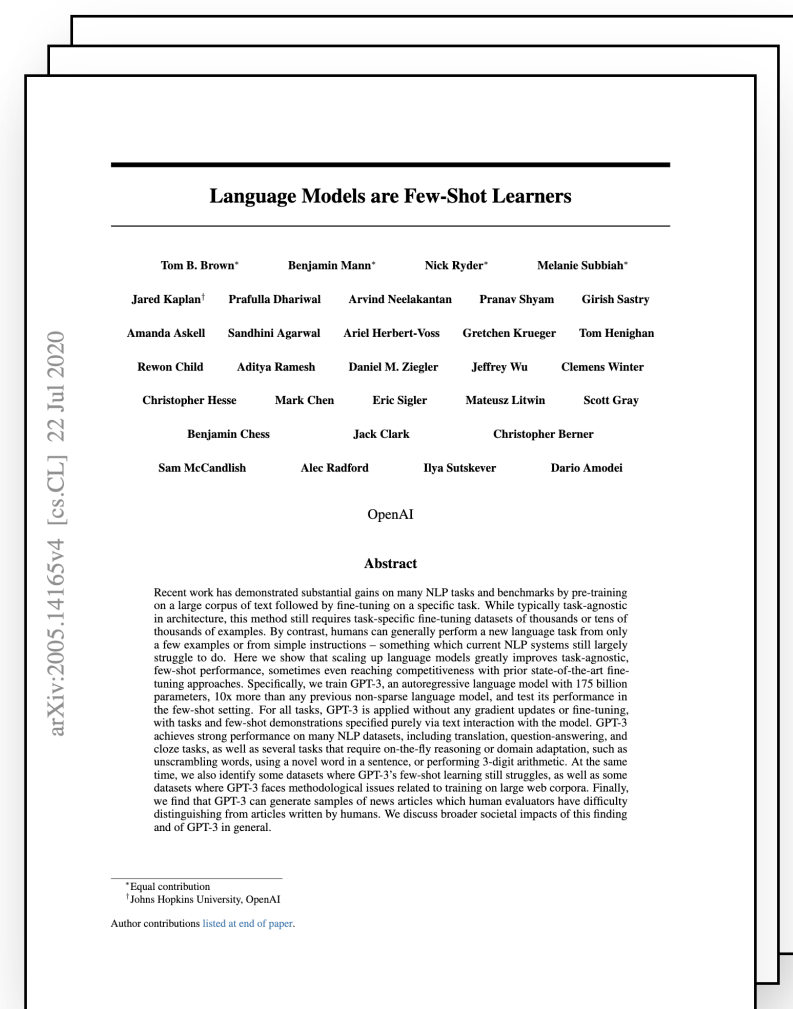


100 web pages  
(~**100K** tokens)



# Language models: the new “search engines”

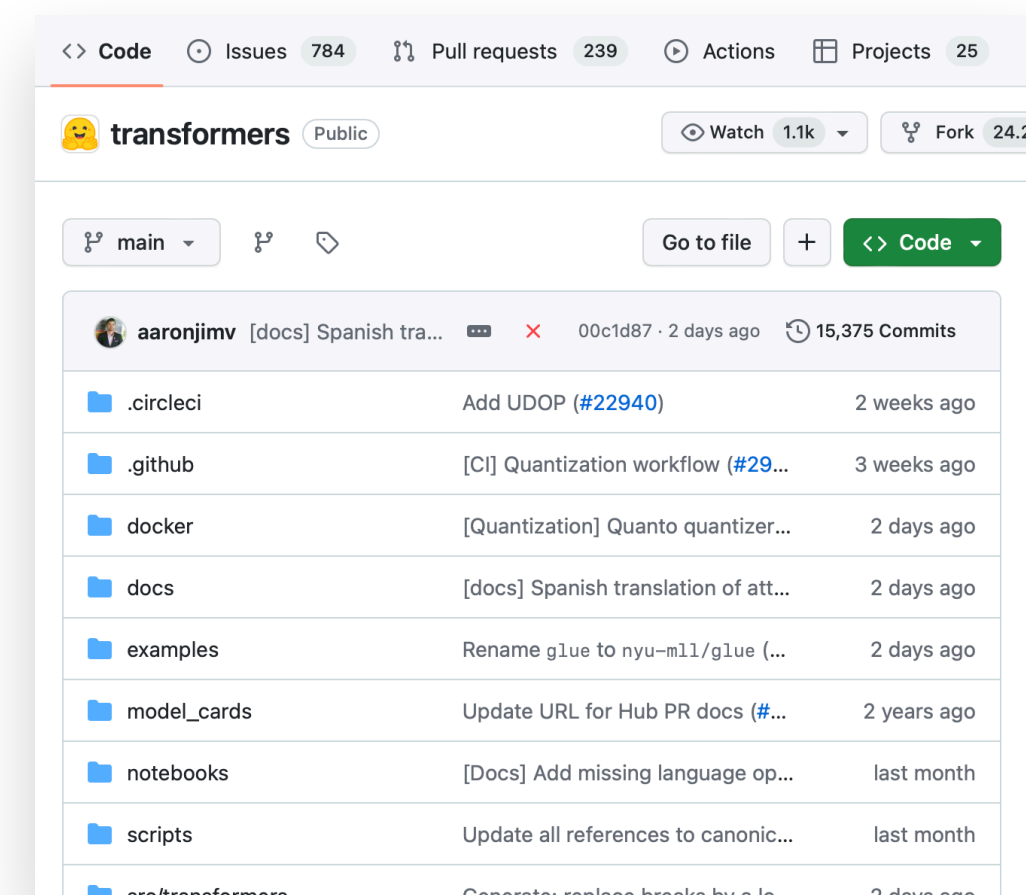
## Challenge #3: how to fit the context



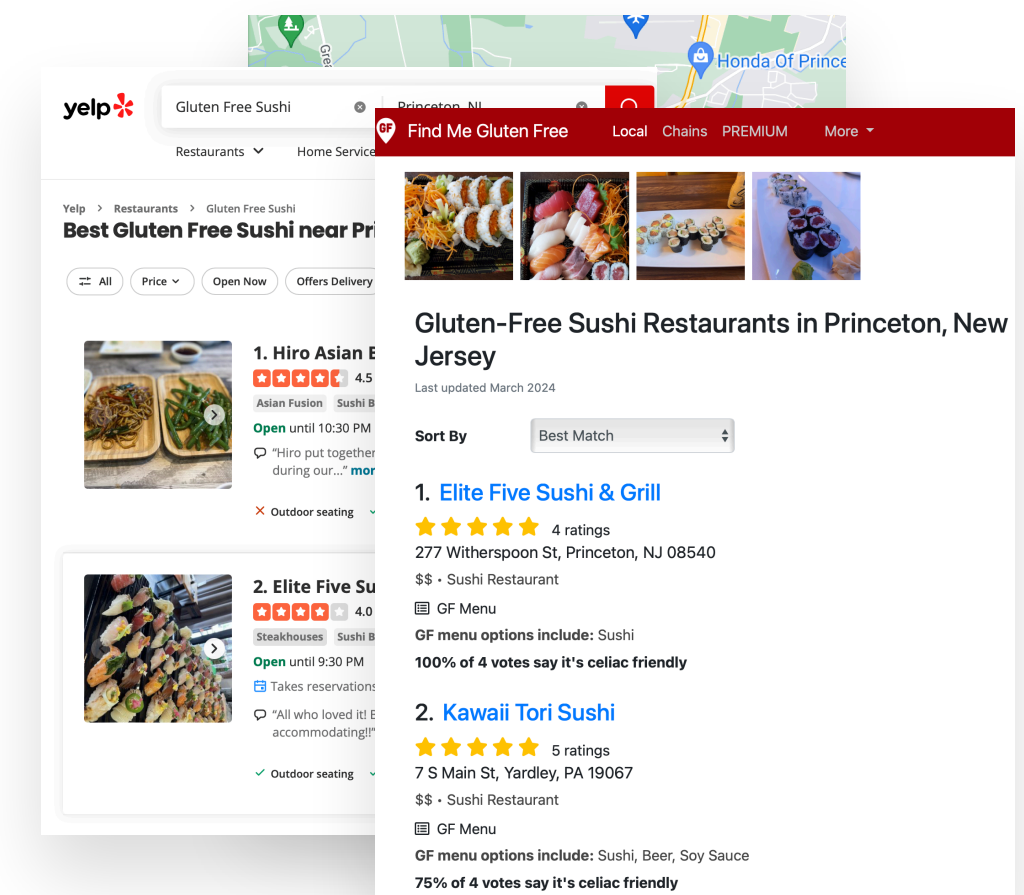
The GPT-3 paper  
(~**75K** tokens)



The Dune series  
(~**1M** words)



The Transformers package  
(~**10M** tokens)



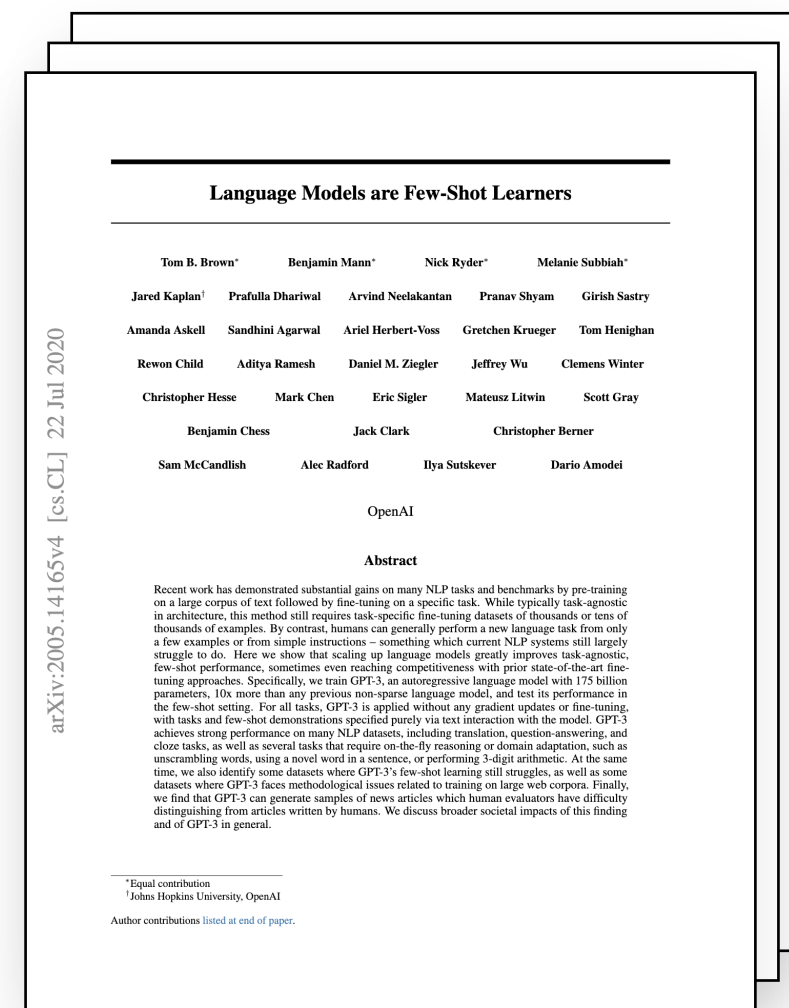
100 web pages  
(~**100K** tokens)

Meta LLaMA-2: 4K



# Language models: the new “search engines”

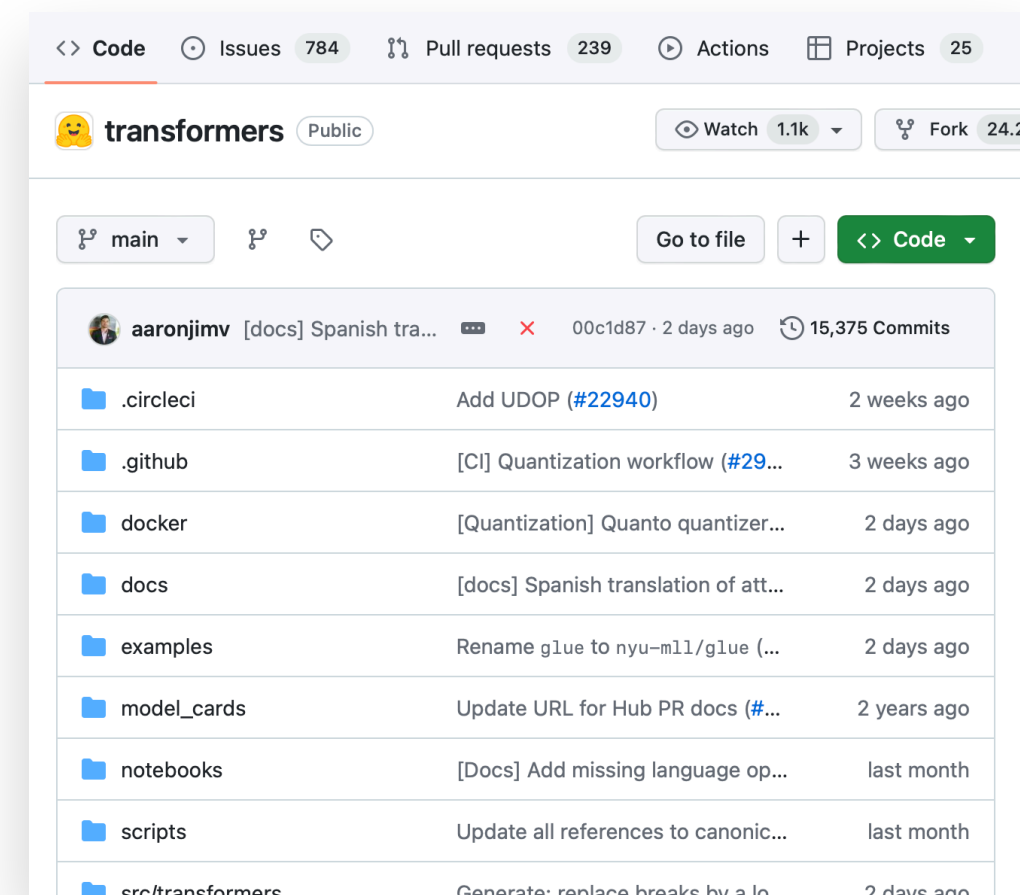
## Challenge #3: how to fit the context



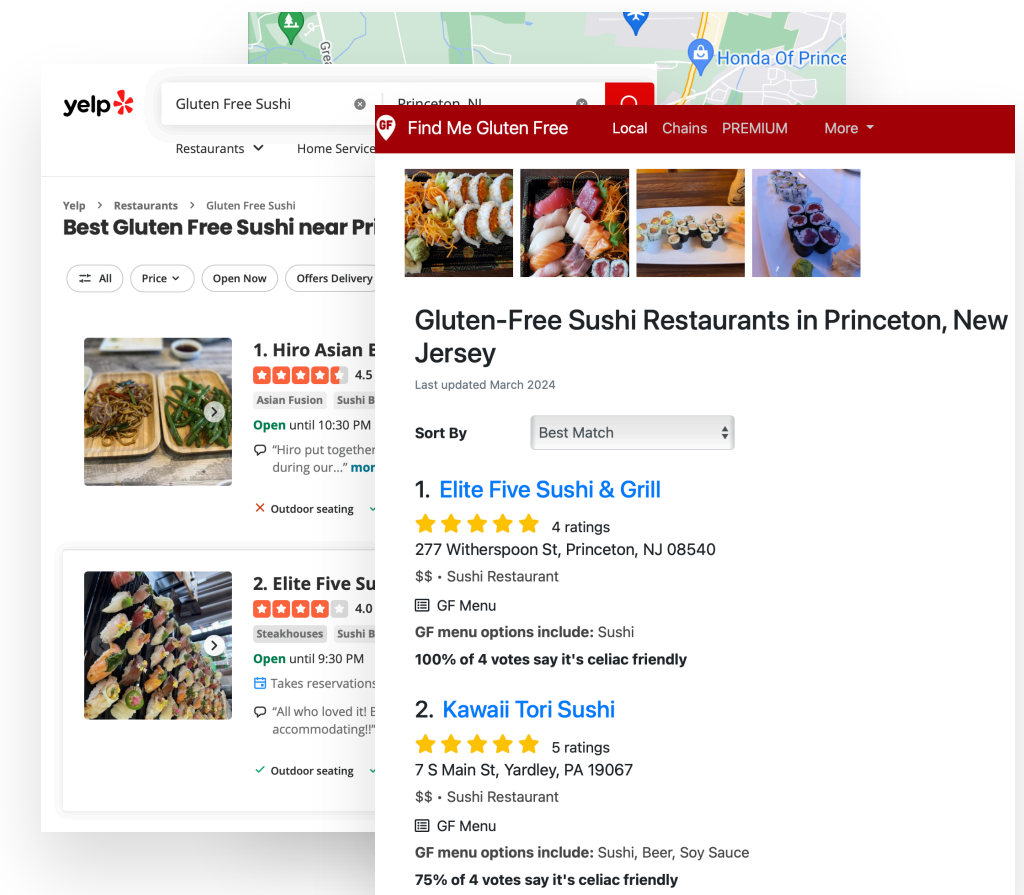
The GPT-3 paper  
(~**75K** tokens)



The Dune series  
(~**1M** words)



The Transformers package  
(~**10M** tokens)



100 web pages  
(~**100K** tokens)

Meta LLaMA-2: 4K

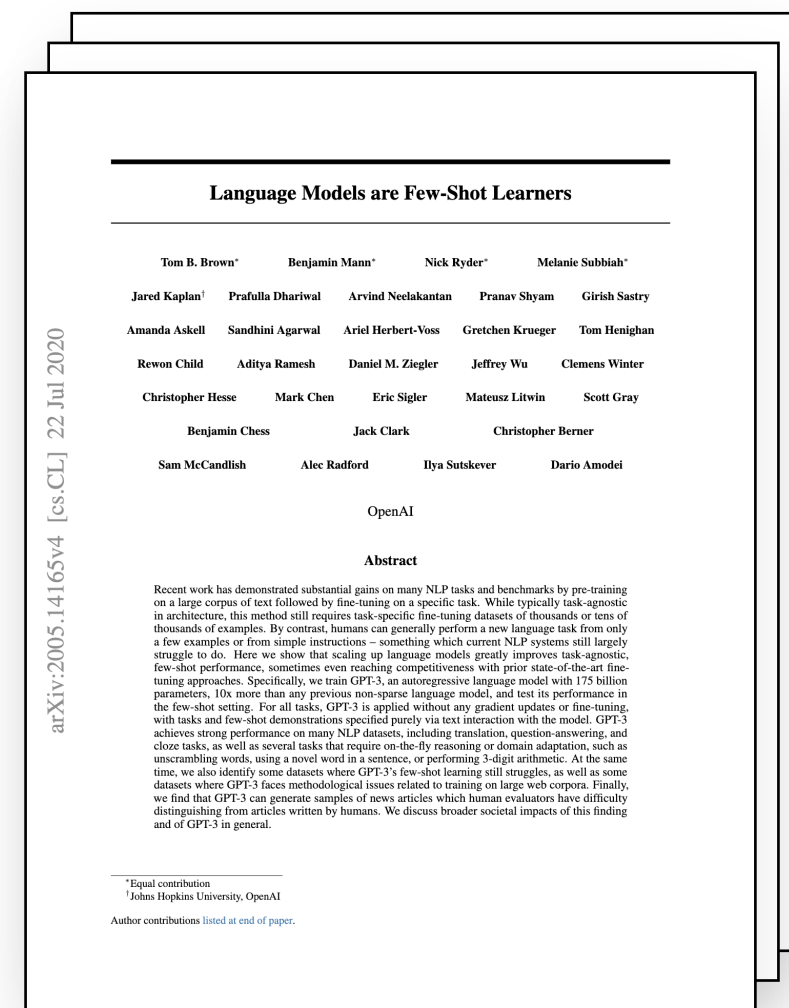


Mixtral: 32K



# Language models: the new “search engines”

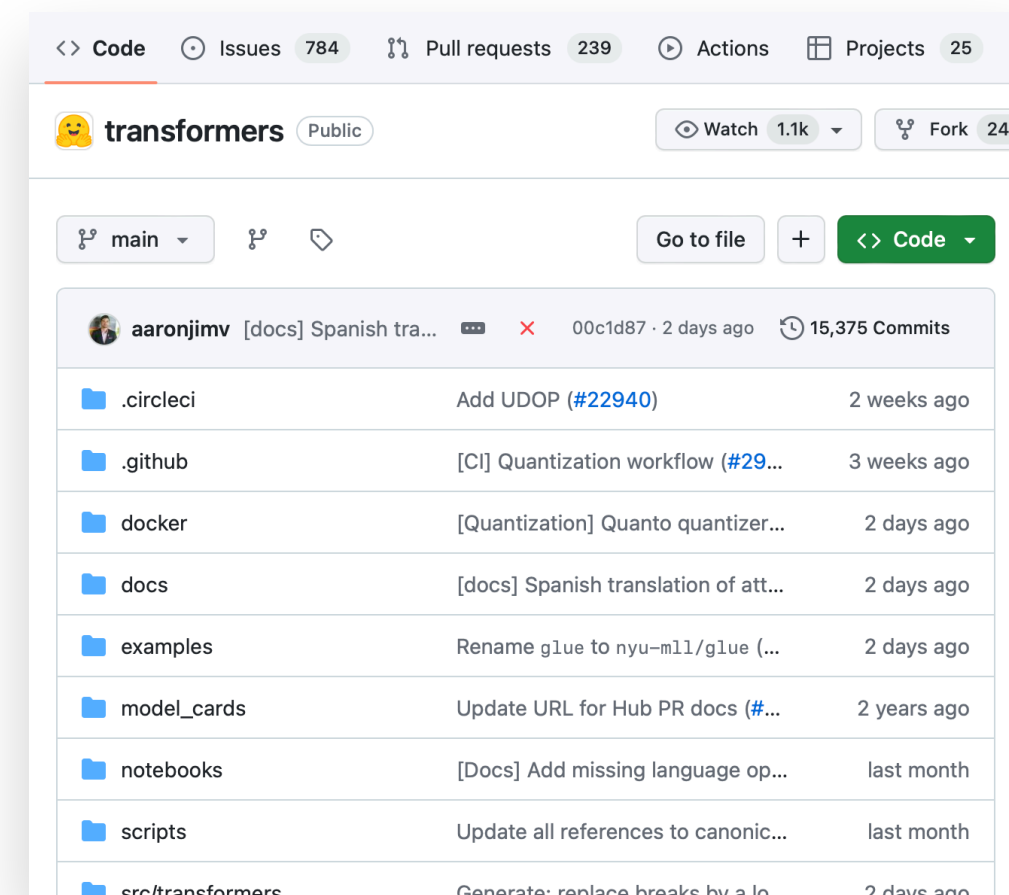
## Challenge #3: how to fit the context



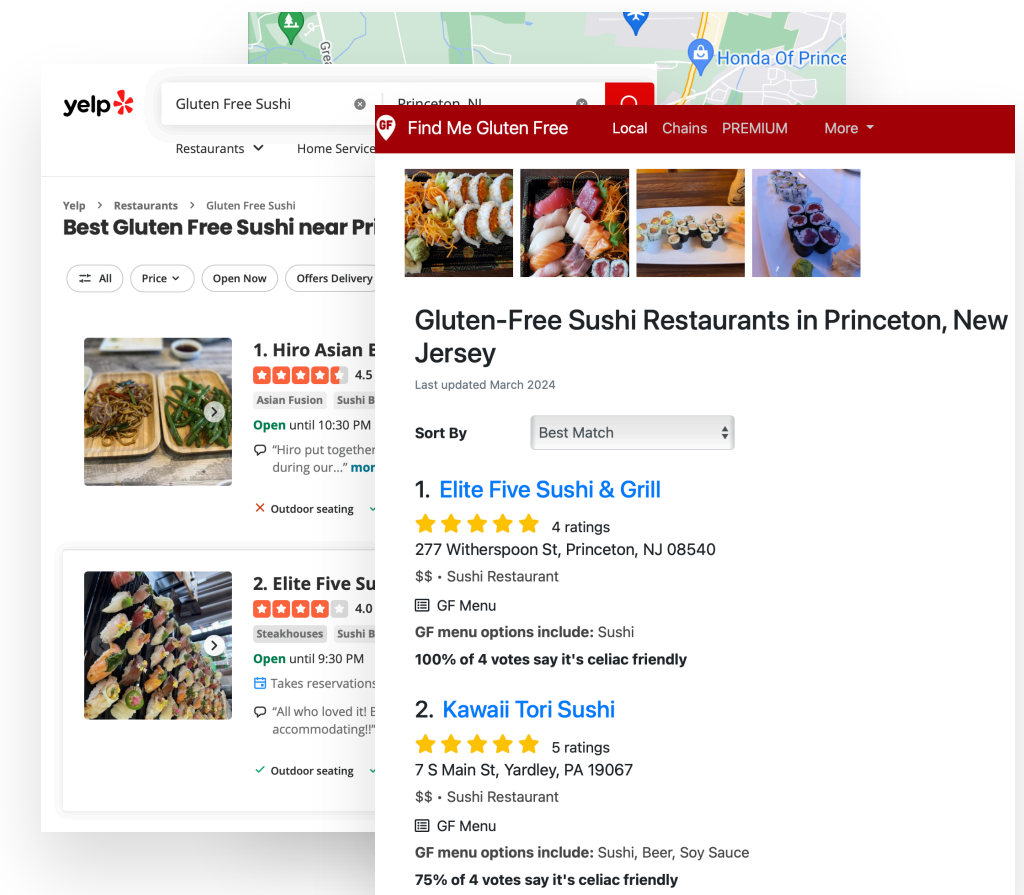
The GPT-3 paper  
(~**75K** tokens)



The Dune series  
(~**1M** words)



The Transformers package  
(~**10M** tokens)



100 web pages  
(~**100K** tokens)

Meta LLaMA-2: 4K



Mixtral: 32K



GPT-4: 128K

# Language models: the new “search engines”

## Challenge #3: how to fit the context



We need long-context language models!

The GPT-3 paper  
(~**75K** tokens)

The Dune series  
(~**1M** words)

The Transformers package  
(~**10M** tokens)

100 web pages  
(~**100K** tokens)

Meta LLaMA-2: 4K



Mixtral: 32K



GPT-4: 128K



# Long-context hurdle #1: Transformers

# Long-context hurdle #1: Transformers

**Transformers are costly (both computation and memory)**

# Long-context hurdle #1: Transformers

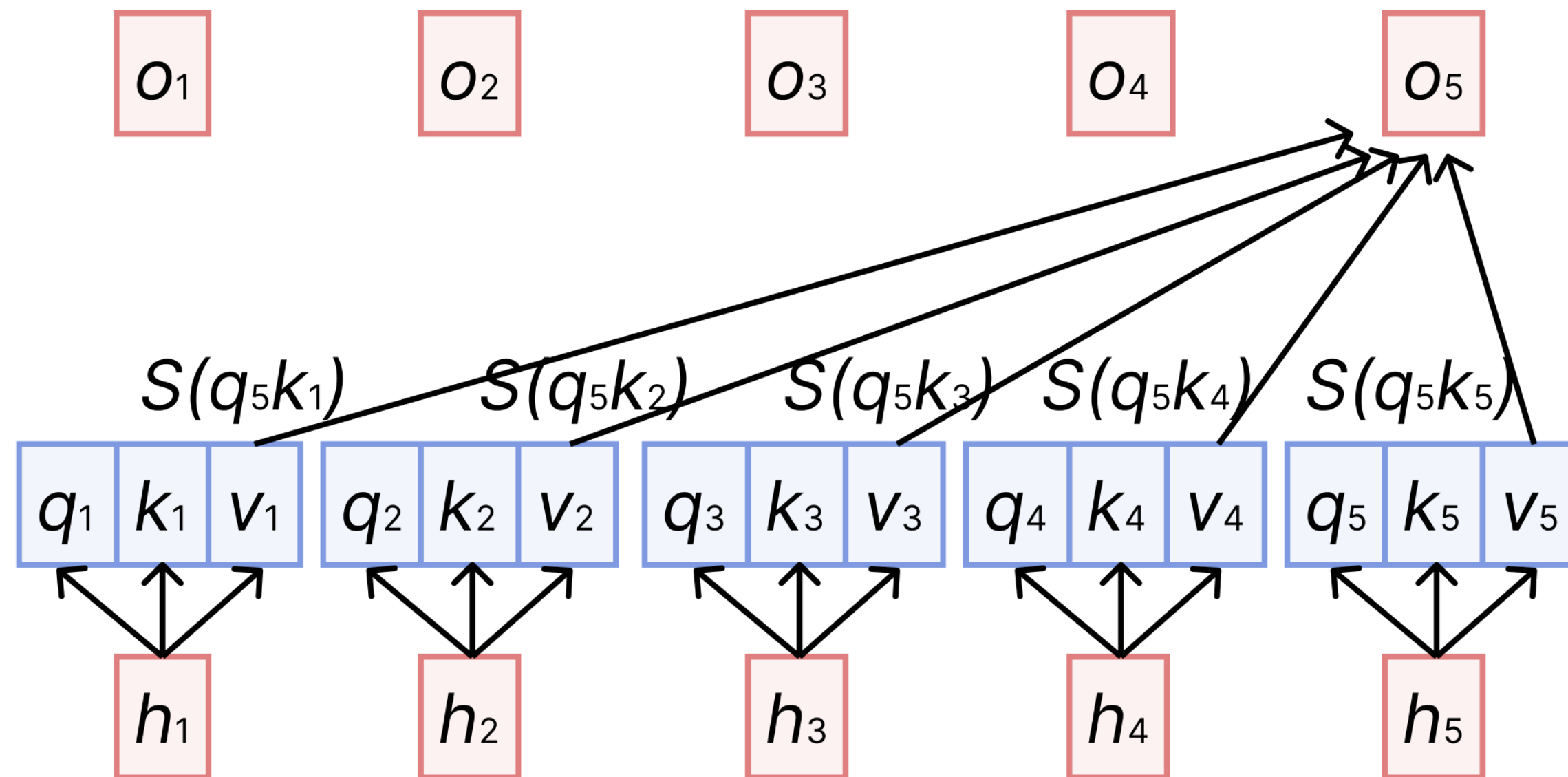
**Transformers are costly (both computation and memory)**

Multi-head attention: the representation of every word is a *weighted sum of all* previous words

# Long-context hurdle #1: Transformers

## Transformers are costly (both computation and memory)

Multi-head attention: the representation of every word is a *weighted sum of all previous words*

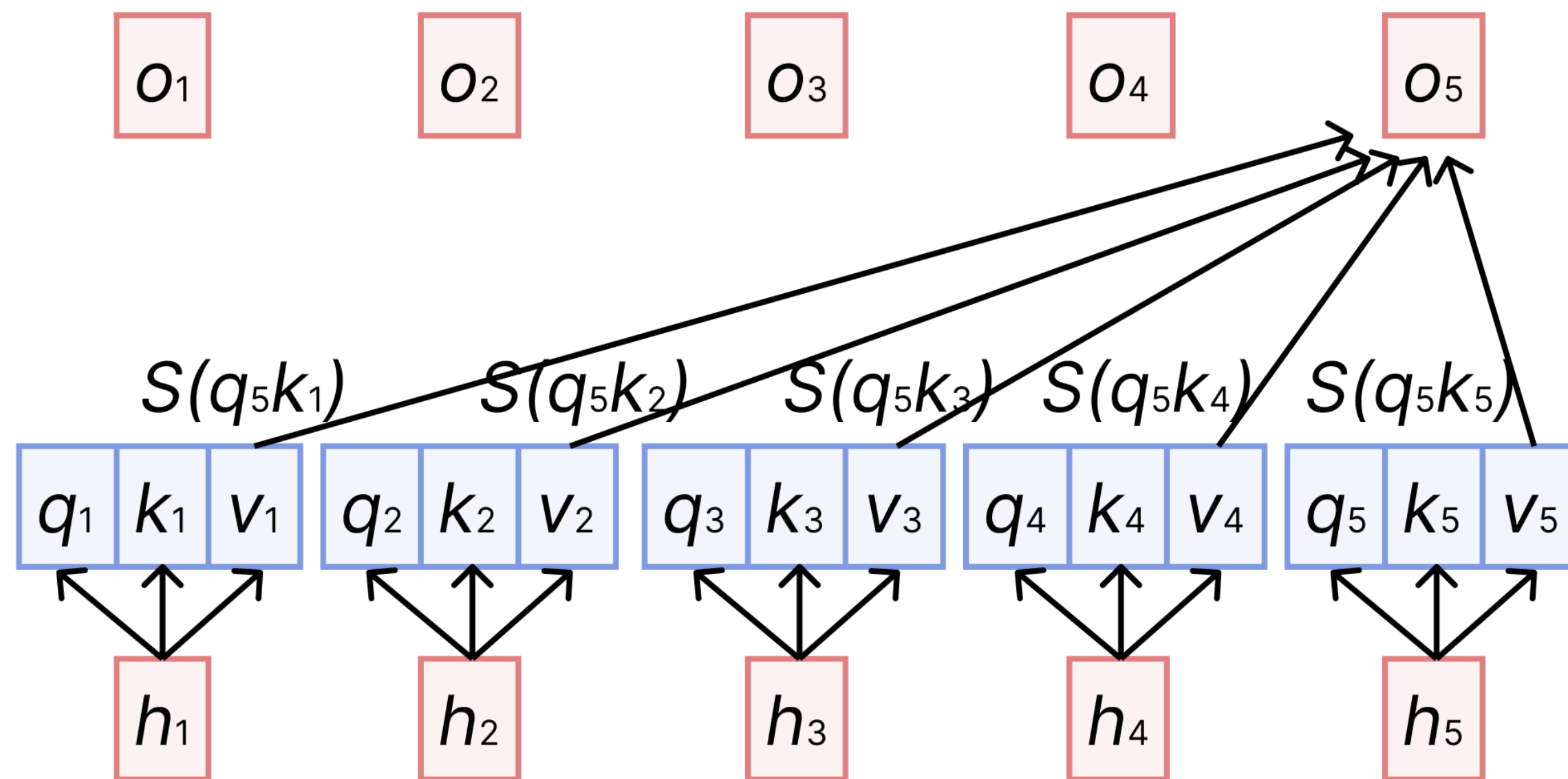




# Long-context hurdle #1: Transformers

## Transformers are costly (both computation and memory)

Multi-head attention: the representation of every word is a *weighted sum of all previous words*

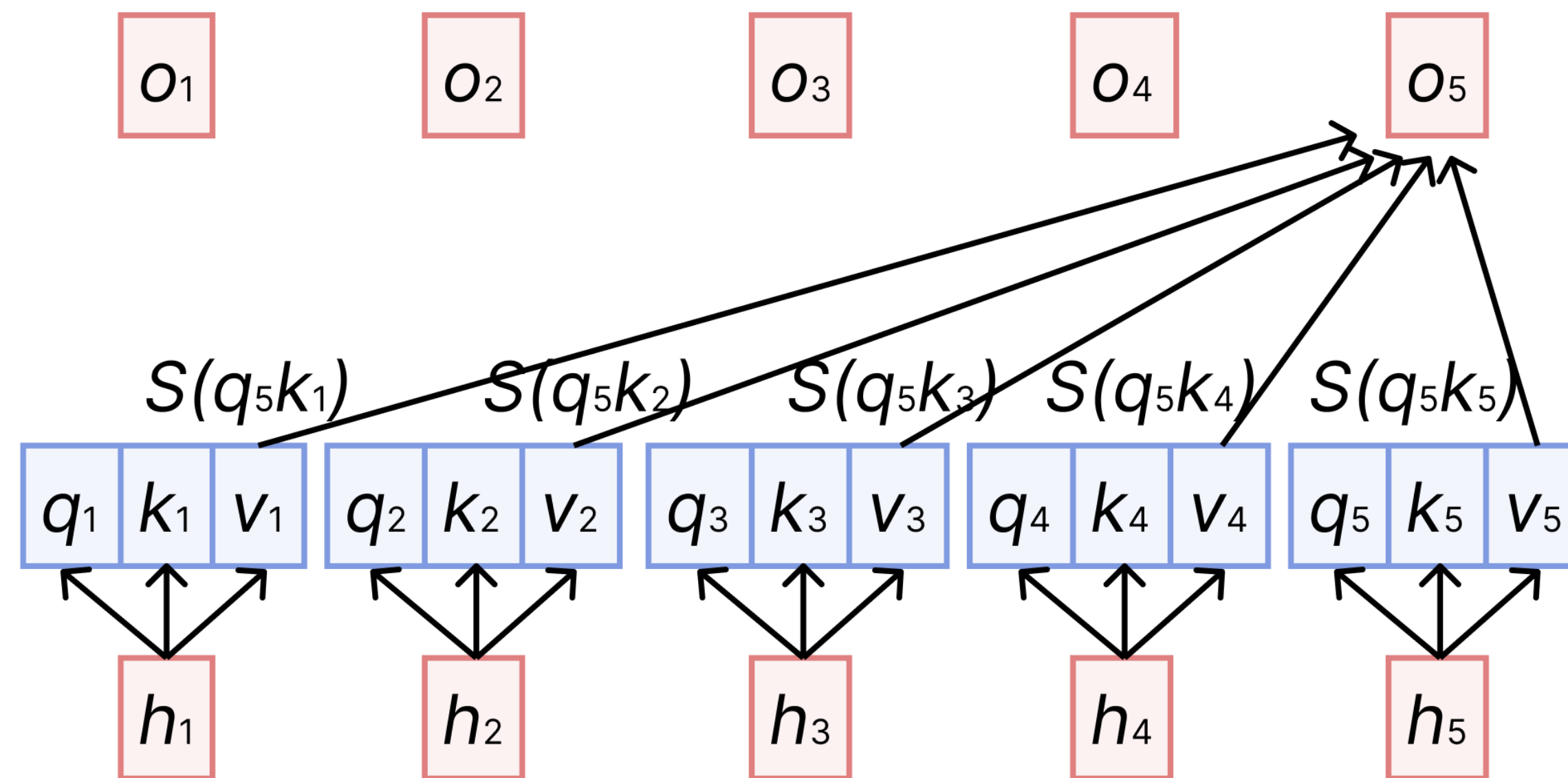


To encode a context of  $n$  words  
 $\mathcal{O}(n^2)$  compute complexity

# Long-context hurdle #1: Transformers

## Transformers are costly (both computation and memory)

Multi-head attention: the representation of every word is a *weighted sum of all previous words*



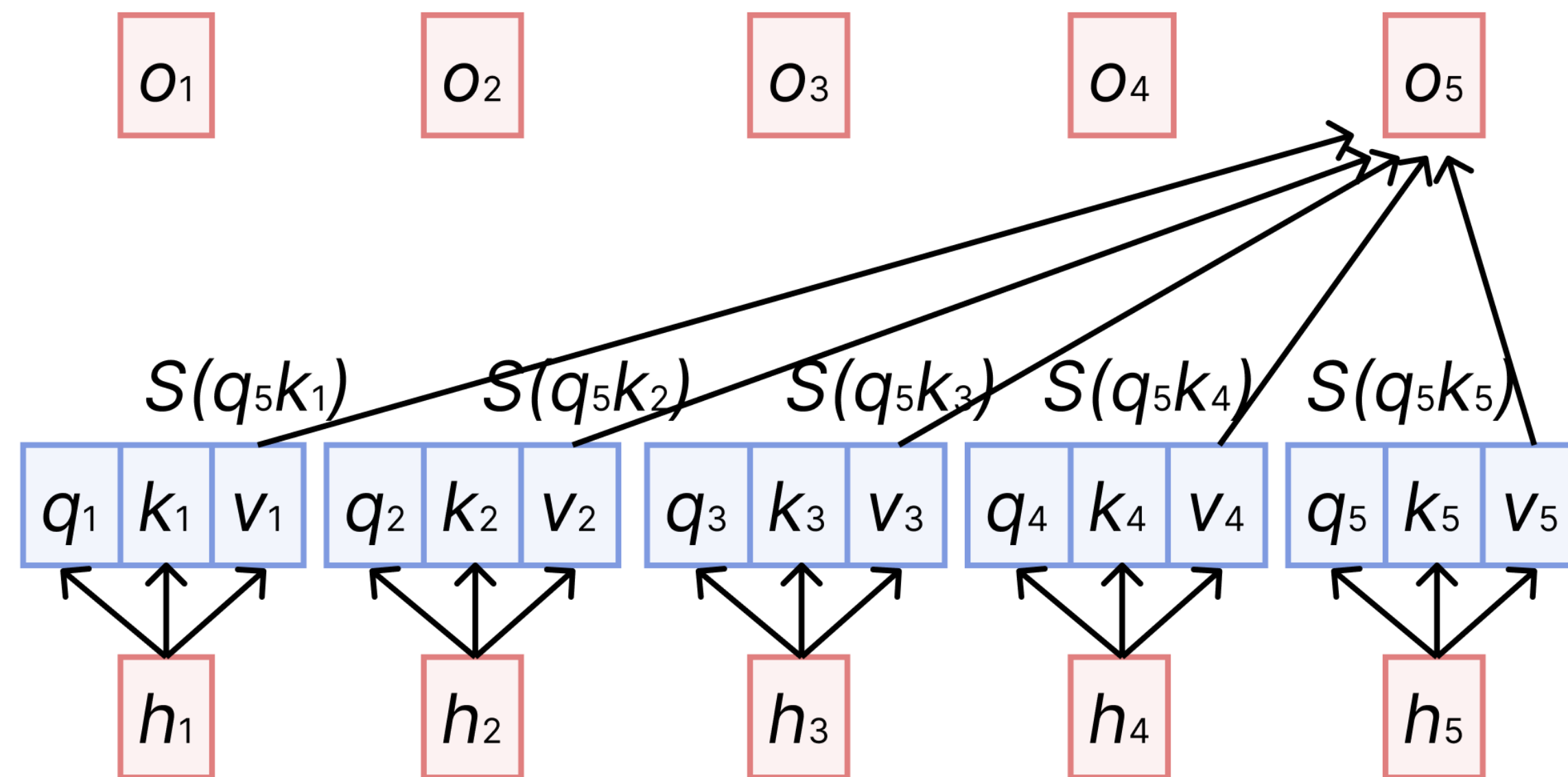
To encode a context of  $n$  words  
 $\mathcal{O}(n^2)$  compute complexity

To predict the next word  
 $\mathcal{O}(n)$  memory complexity

# Long-context hurdle #1: Transformers

## Transformers are costly (both computation and memory)

Multi-head attention: the representation of every word is a *weighted sum of all previous words*



To encode a context of  $n$  words  
 $\mathcal{O}(n^2)$  compute complexity

To predict the next word  
 $\mathcal{O}(n)$  memory complexity

A **1M** token context would cost  
**164GB** memory!  
(LLaMA-70B, FP16)

# Long-context hurdle #2: Positional encodings



# Long-context hurdle #2: Positional encodings

**Popular positional encodings are not generalizable**

# Long-context hurdle #2: Positional encodings

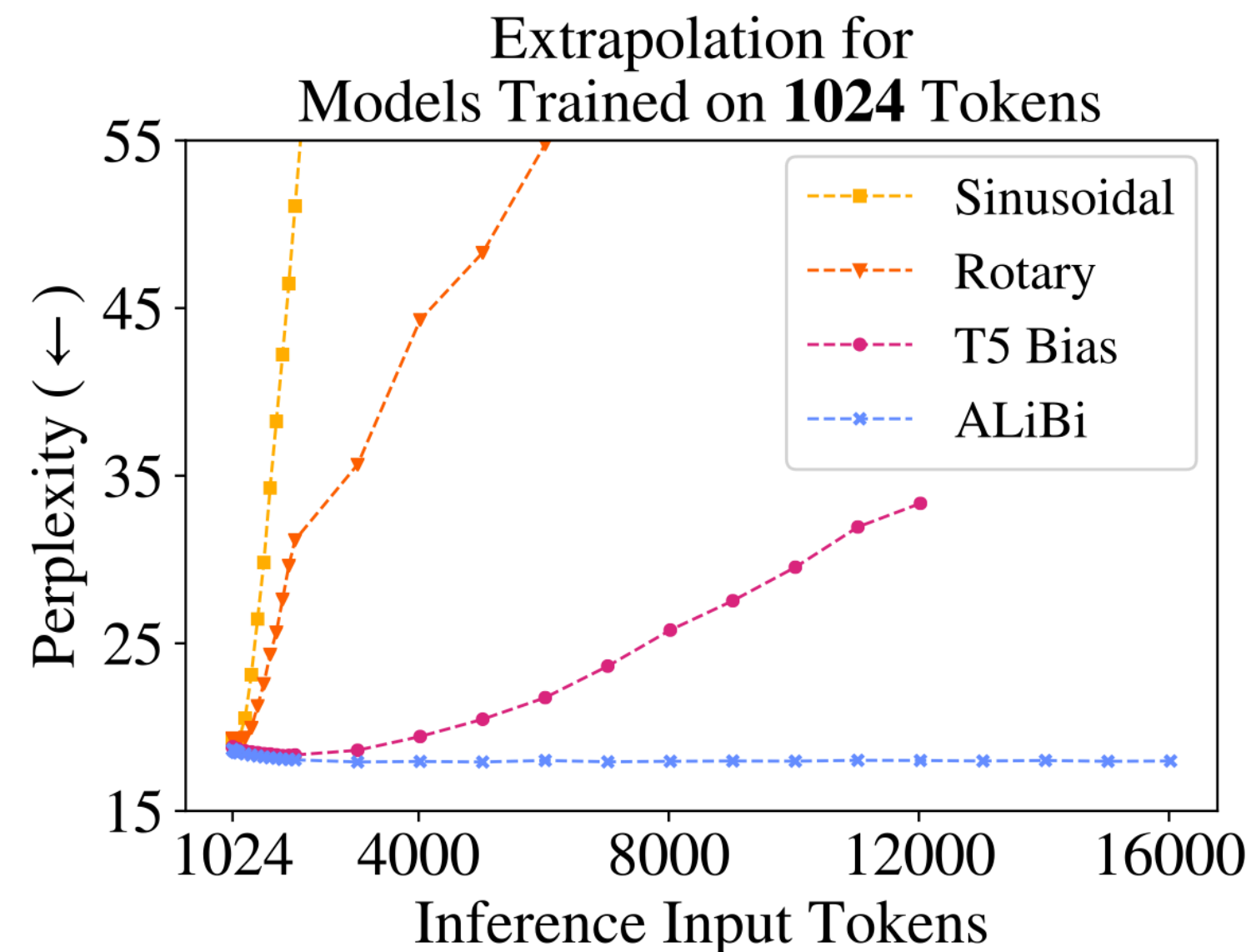
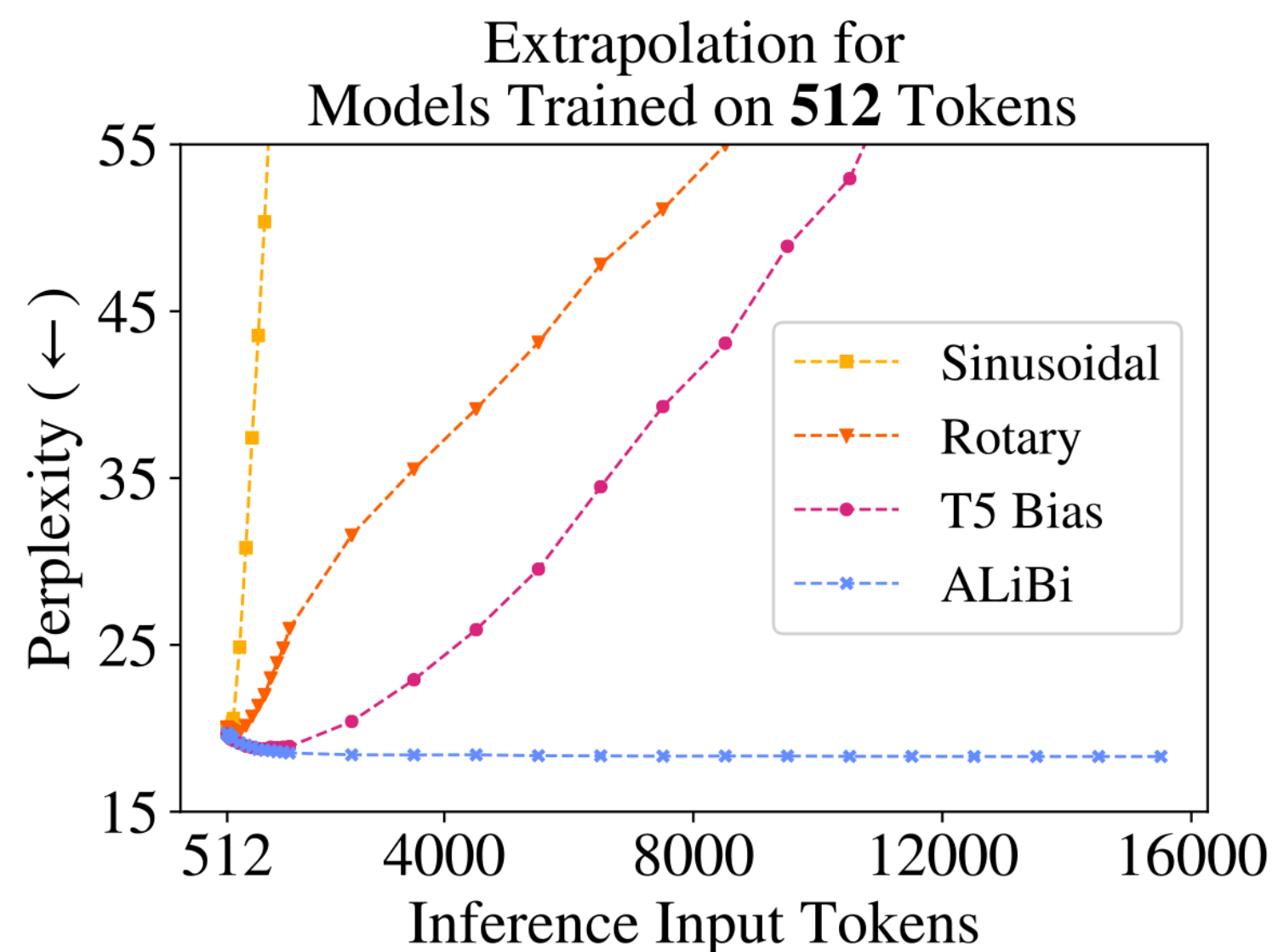
## Popular positional encodings are not generalizable

The most popular positional encoding method, RoPE (Su et al., 2021), cannot generalize beyond the training length.

# Long-context hurdle #2: Positional encodings

## Popular positional encodings are not generalizable

The most popular positional encoding method, RoPE (Su et al., 2021), cannot generalize beyond the training length.



# Long-context hurdle #3: Data scarcity



# Long-context hurdle #3: Data scarcity

**High-quality long-context data are hard to find**

# Long-context hurdle #3: Data scarcity

## High-quality long-context data are hard to find

Average length of domains from **RedPajama** (a open-source pre-training data collection)

- Wikipedia: 0.5K tokens
- C4 (webpages): 0.5K tokens
- Arxiv: 20K tokens
- Books: 147K tokens

Average length of **instruction-tuning/chat** data: <1K tokens

# Long-context hurdle #3: Data scarcity

## High-quality long-context data are hard to find

Average length of domains from **RedPajama** (a open-source pre-training data collection)

- Wikipedia: 0.5K tokens
- C4 (webpages): 0.5K tokens
- Arxiv: 20K tokens
- Books: 147K tokens

Average length of **instruction-tuning/chat** data: <1K tokens

**How can we train a model that can continually generalize to longer length?**

# Long-Context Language Modeling with Parallel Context Encoding





# Long-Context Language Modeling with Parallel Context Encoding



- We propose 🍄 **CEPE: Context Expansion with Parallel Encoding**

# Long-Context Language Modeling with Parallel Context Encoding



- We propose 🍄 **CEPE: Context Expansion with Parallel Encoding**
  - A lightweight framework that can extend any decoder-only LM's context length

# Long-Context Language Modeling with Parallel Context Encoding



- We propose 🍄 **CEPE: Context Expansion with Parallel Encoding**
  - A lightweight framework that can extend any decoder-only LM's context length
- Why CEPE

# Long-Context Language Modeling with Parallel Context Encoding



- We propose 🍄 **CEPE: Context Expansion with Parallel Encoding**
  - A lightweight framework that can extend any decoder-only LM's context length
- Why CEPE
  - **Length generalization** (train on 8K, generalize to 128K)

# Long-Context Language Modeling with Parallel Context Encoding



- We propose 🍄 **CEPE: Context Expansion with Parallel Encoding**
  - A lightweight framework that can extend any decoder-only LM's context length
- Why CEPE
  - **Length generalization** (train on 8K, generalize to 128K)
  - **Efficiency** (10x throughput, 1/256 memory for every extra token\* compared to LLaMA-7B)



# Long-Context Language Modeling with Parallel Context Encoding



- We propose 🍄 **CEPE: Context Expansion with Parallel Encoding**
  - A lightweight framework that can extend any decoder-only LM's context length
- Why CEPE
  - **Length generalization** (train on 8K, generalize to 128K)
  - **Efficiency** (10x throughput, 1/256 memory for every extra token\* compared to LLaMA-7B)
  - **Low training cost** (extending a LLaMA-7B can be done by 1 A100)

# Long-Context Language Modeling with Parallel Context Encoding



- We propose 🍄 **CEPE: Context Expansion with Parallel Encoding**
  - A lightweight framework that can extend any decoder-only LM's context length
- Why CEPE
  - **Length generalization** (train on 8K, generalize to 128K)
  - **Efficiency** (10x throughput, 1/256 memory for every extra token\* compared to LLaMA-7B)
  - **Low training cost** (extending a LLaMA-7B can be done by 1 A100)
  - Also work for **chat** models

# Long-Context Language Modeling with Parallel Context Encoding



- We propose 🍄 **CEPE: Context Expansion with Parallel Encoding**
  - A lightweight framework that can extend any decoder-only LM's context length
- Why CEPE
  - **Length generalization** (train on 8K, generalize to 128K)
  - **Efficiency** (10x throughput, 1/256 memory for every extra token\* compared to LLaMA-7B)
  - **Low training cost** (extending a LLaMA-7B can be done by 1 A100)
  - Also work for **chat** models
- CEPE achieves great performance on both **long-context** and **retrieval-augmented** applications

# How do decoder-only LMs handle long context?

# How do decoder-only LMs handle long context?



The Dune series

Q: Who betrayed the Atreides?



# How do decoder-only LMs handle long context?

Decoder-only language model



The Dune series

Q: Who betrayed the Atreides?

# How do decoder-only LMs handle long context?

Decoder-only language model



The Dune series

Q: Who betrayed the Atreides?

$\mathcal{O}(n)$  memory cost  
 $\mathcal{O}(n^2)$  computational cost



# How do decoder-only LMs handle long context?

Decoder-only language model



The Dune series

Q: Who betrayed the Atreides?

$\mathcal{O}(n)$  memory cost  
 $\mathcal{O}(n^2)$  computational cost

$n \sim 10K$  to  $1M$

# How does CEPE handle long context?



The Dune series

Q: Who betrayed the Atreides?



# How does CEPE handle long context?



The Dune series

**Additional context**

Q: Who betrayed the Atreides?



# How does CEPE handle long context?



The Dune series

**Additional context**

Q: Who betrayed the Atreides?

**Main input**

# How does CEPE handle long context?



The Dune series

## Additional context

- Can be processed by “chunks”
- Order does not matter much

Q: Who betrayed the Atreides?

## Main input



# How does CEPE handle long context?



The Dune series

## Additional context

- Can be processed by “chunks”
- Order does not matter much

Q: Who betrayed the Atreides?

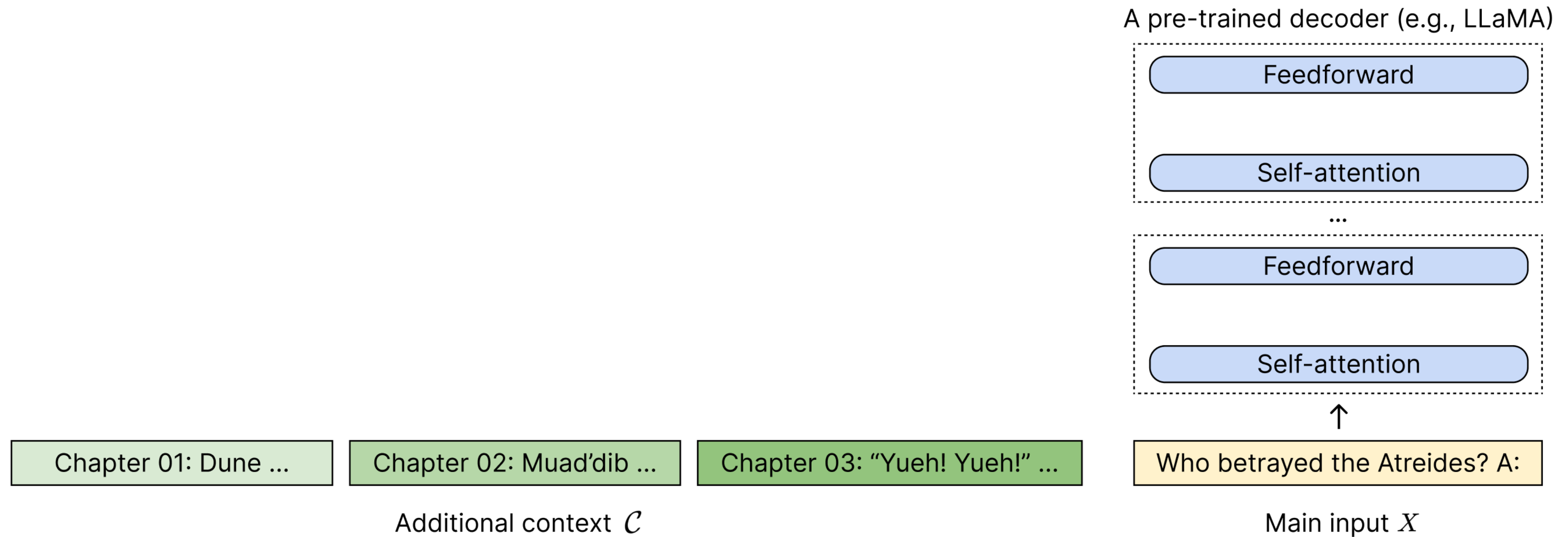
## Main input

- Should be processed with additional context
- Directly related to the generation

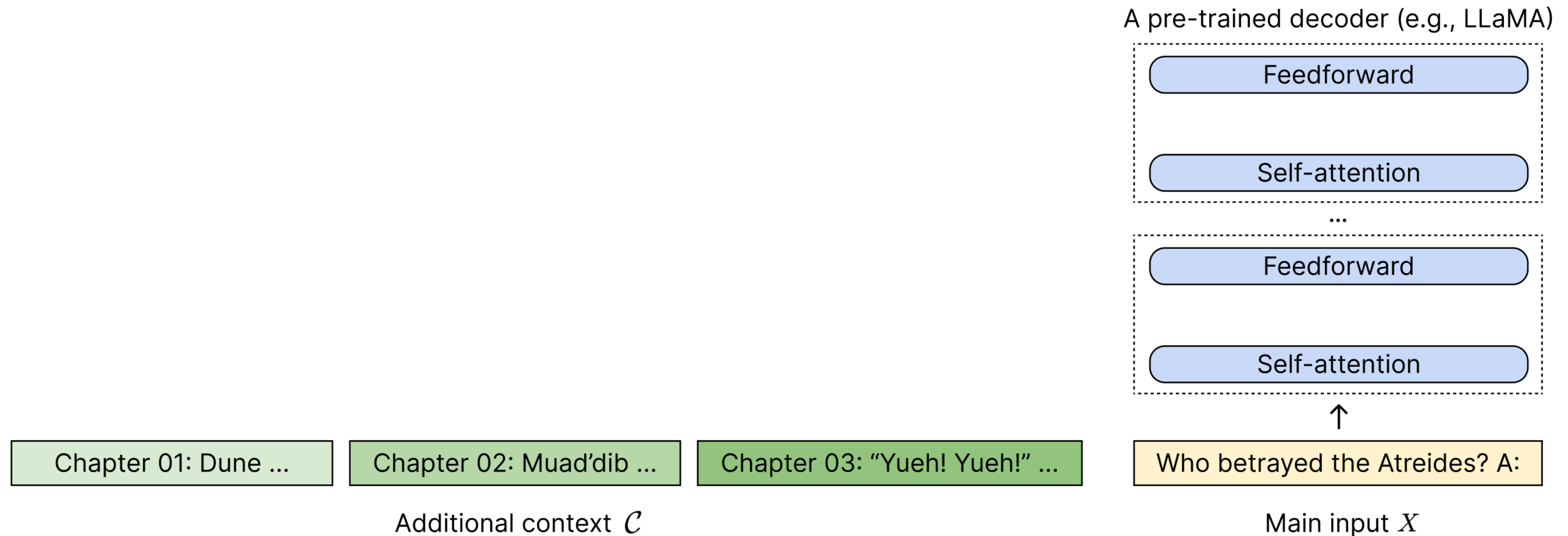
# Context Expansion with Parallel Encoding (CEPE)



# Context Expansion with Parallel Encoding (CEPE)

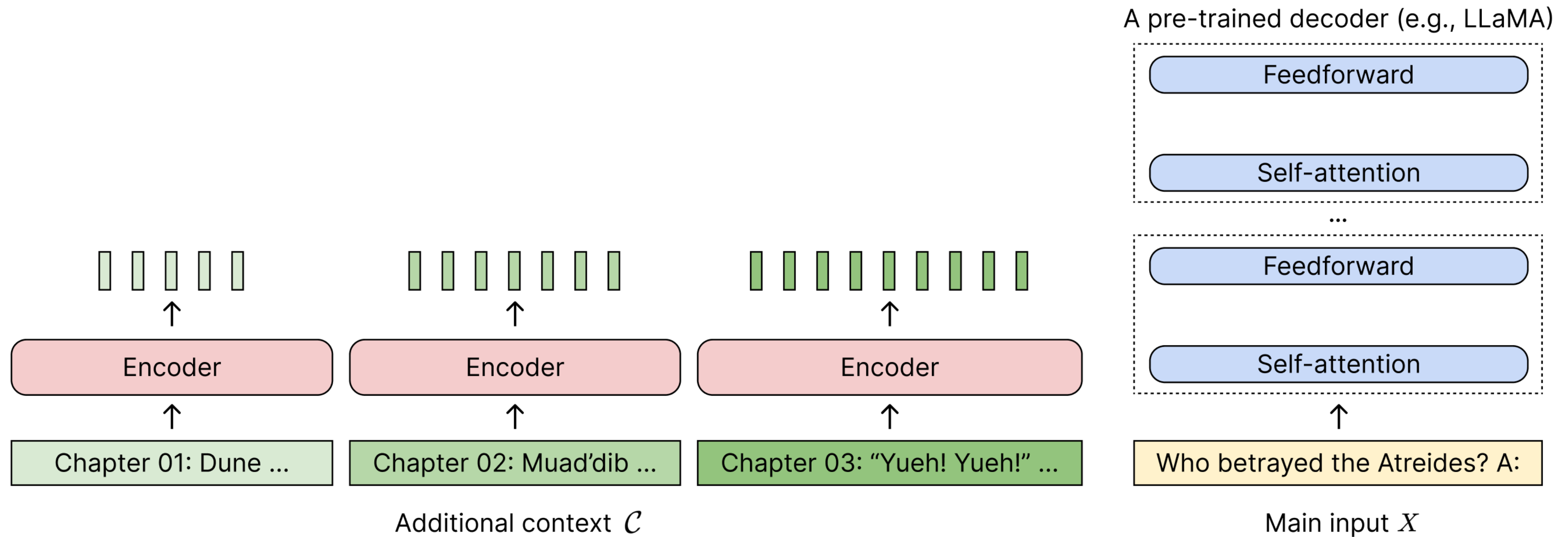


# Context Expansion with Parallel Encoding (CEPE)



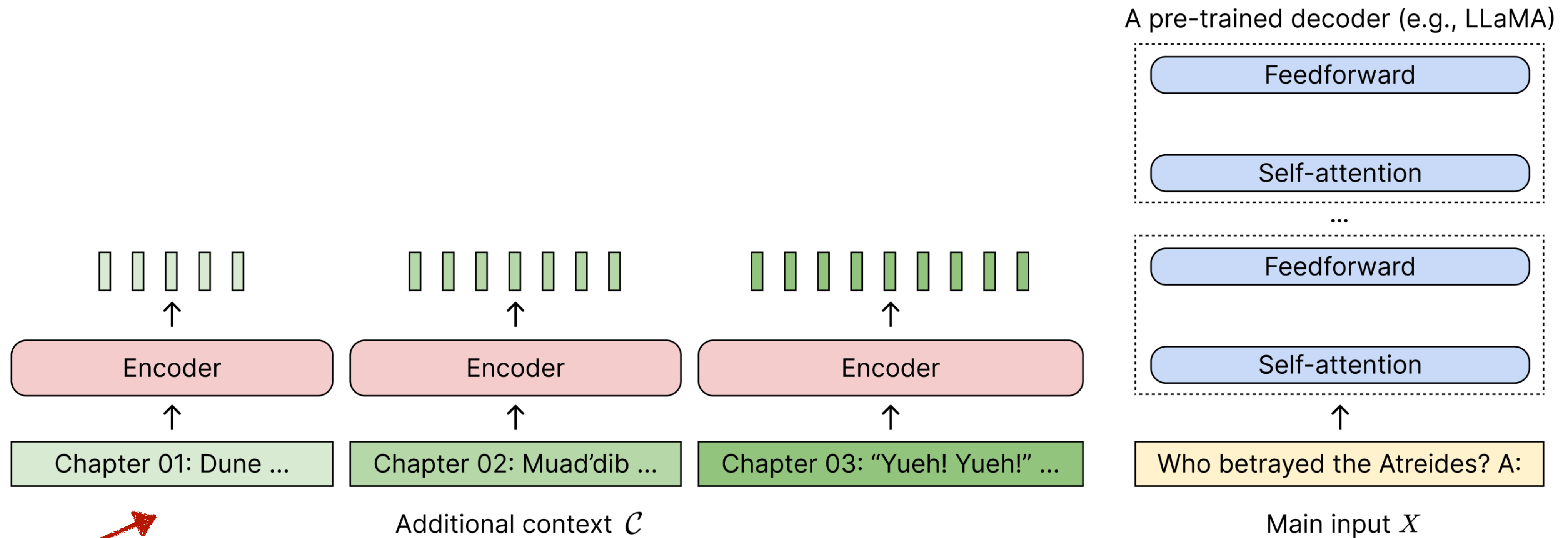
We use an existing decoder-only model (e.g., LLaMA-7B) to process the main input.

# Context Expansion with Parallel Encoding (CEPE)



We use a small bidirectional encoder (435M) to encode the additional context *by chunks*.

# Context Expansion with Parallel Encoding (CEPE)



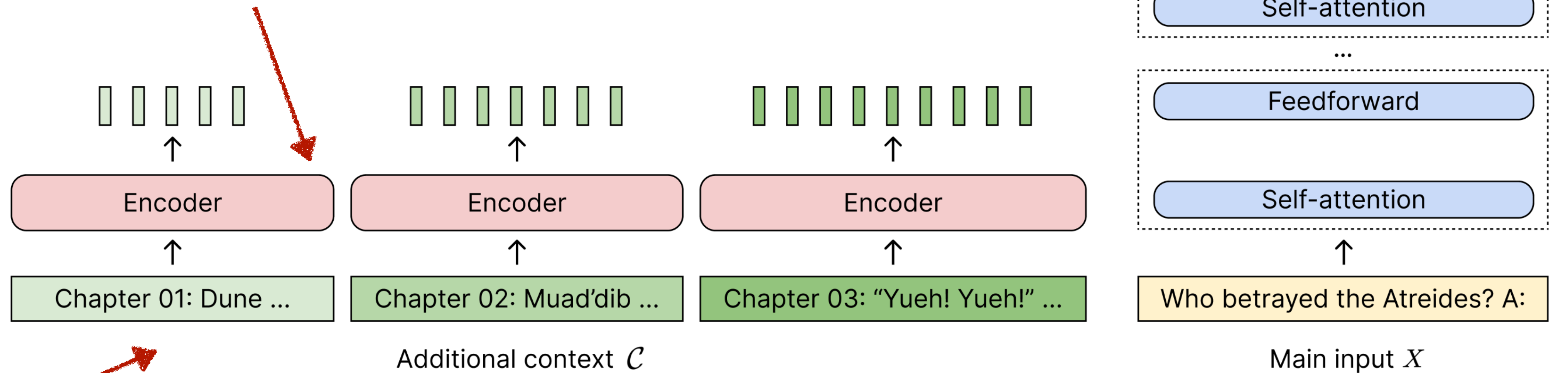
Each chunk has at most 256 tokens

We use a small bidirectional encoder (435M) to encode the additional context *by chunks*.



# Context Expansion with Parallel Encoding (CEPE)

Much faster compared to the decoder  
Bidirectional → better representation



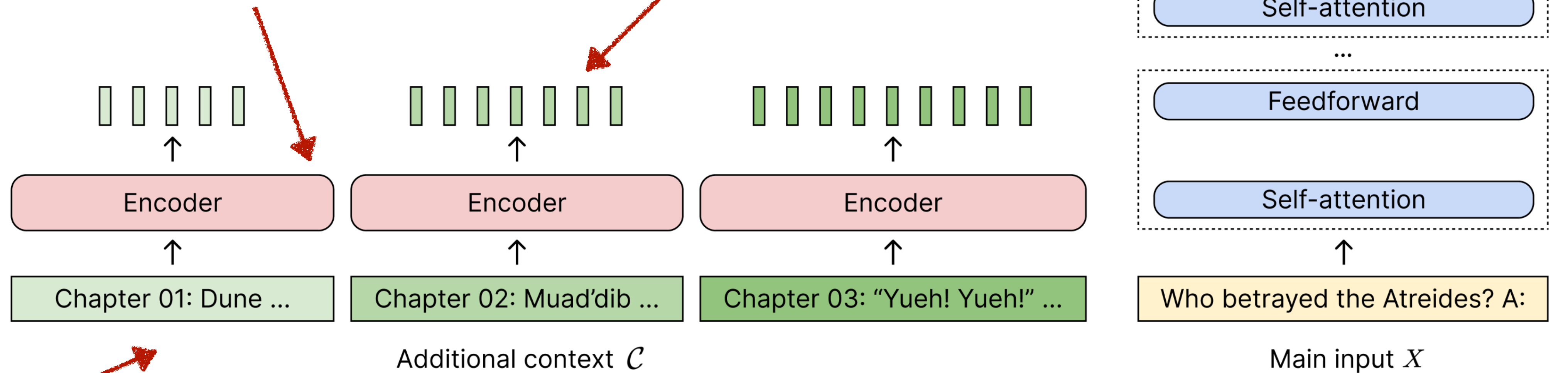
Each chunk has at most 256 tokens

We use a small bidirectional encoder (435M) to encode the additional context *by chunks*.

# Context Expansion with Parallel Encoding (CEPE)

Much faster compared to the decoder  
Bidirectional → better representation

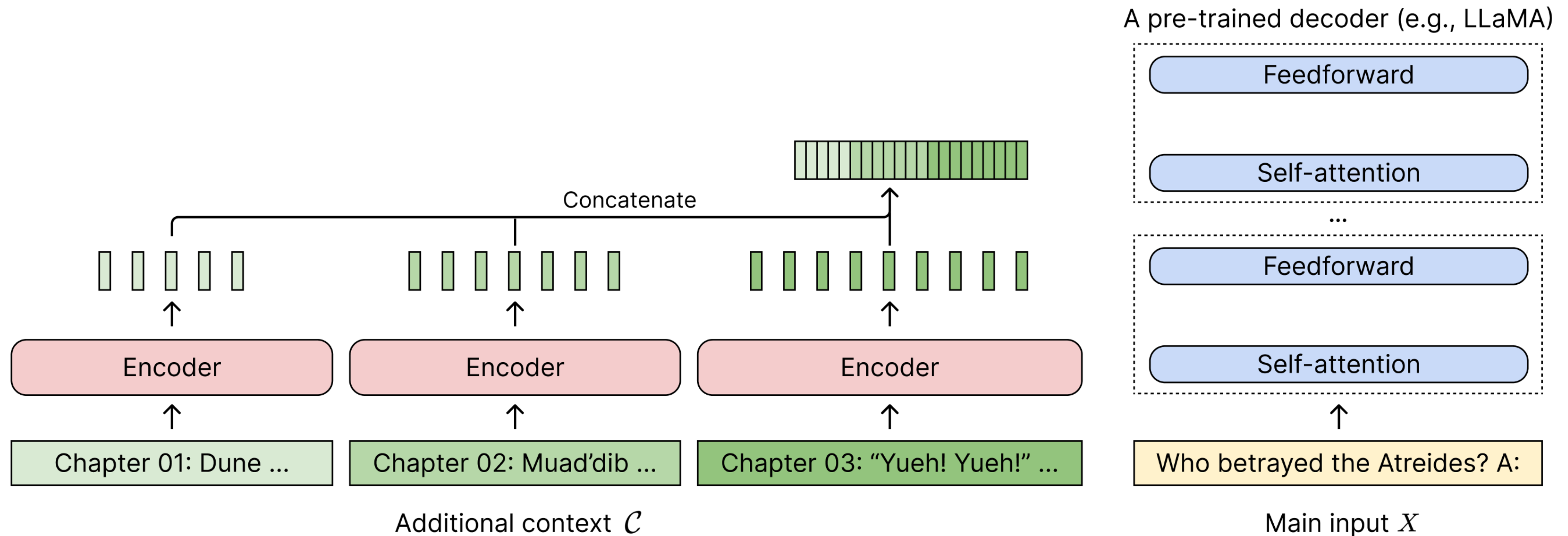
We keep the last layer's output



Each chunk has at most 256 tokens

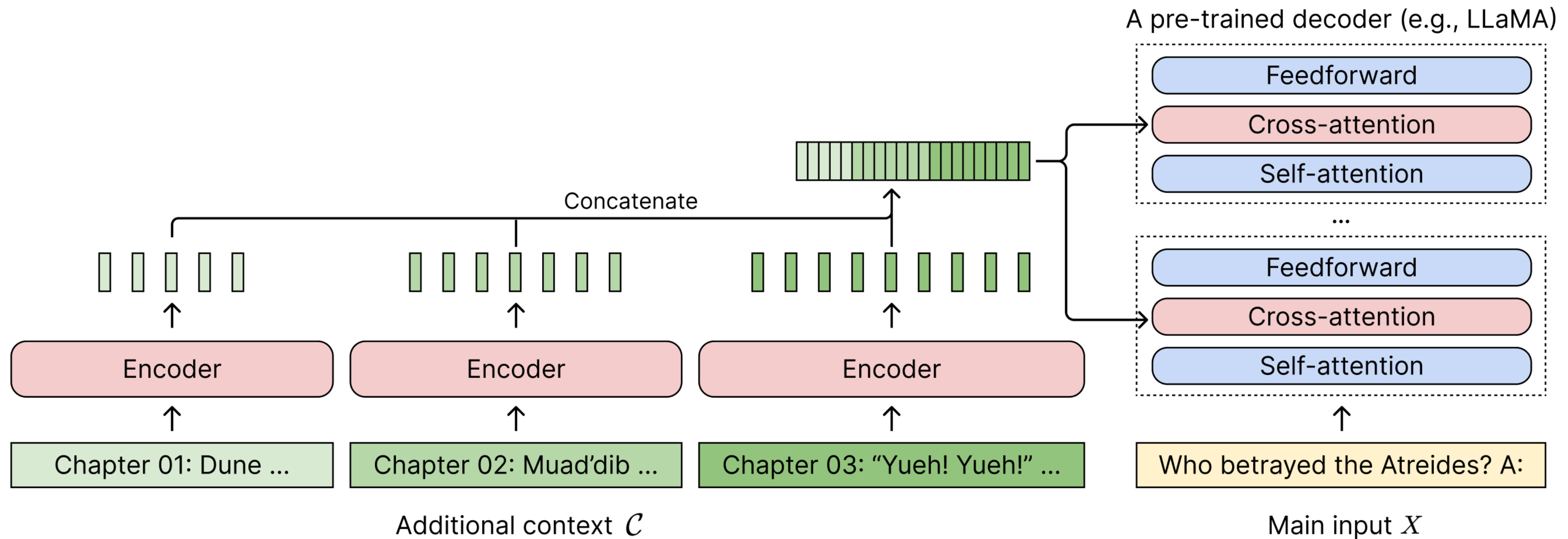
We use a small bidirectional encoder (435M) to encode the additional context *by chunks*.

# Context Expansion with Parallel Encoding (CEPE)



All the encoder outputs are concatenated as the representation for the additional context.

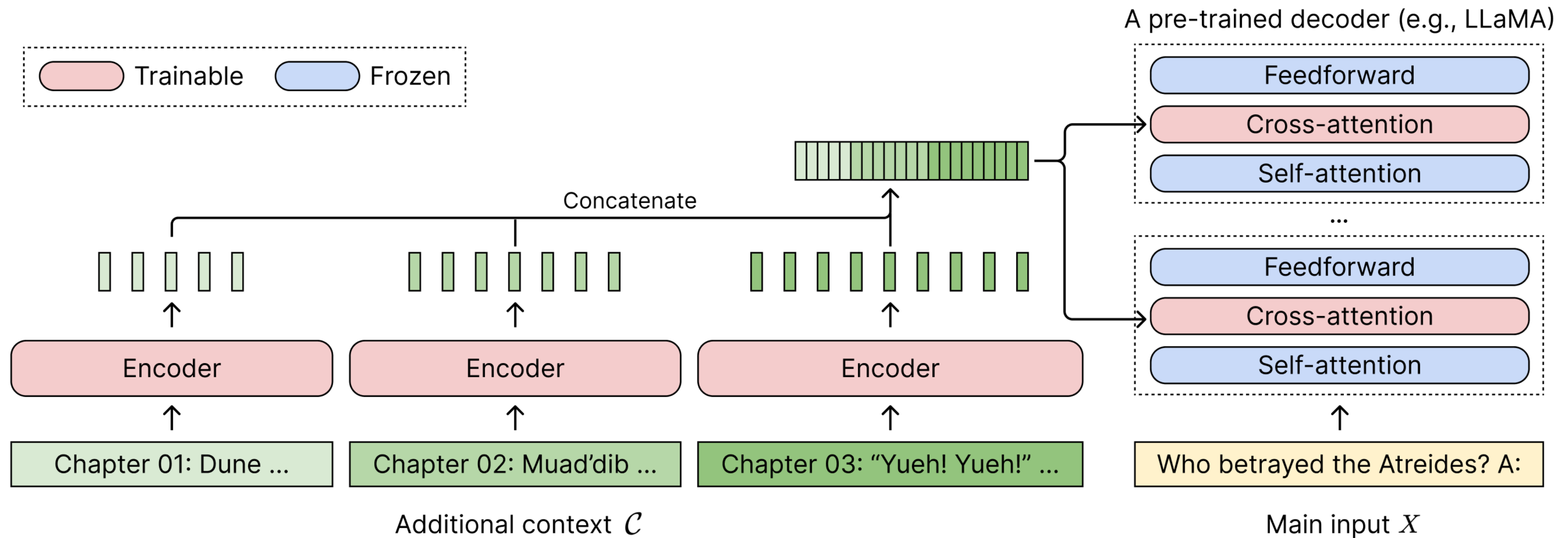
# Context Expansion with Parallel Encoding (CEPE)



We insert *cross-attention* into every layer of the decoder, which attends to the additional context.



# Context Expansion with Parallel Encoding (CEPE)



We freeze the decoder and only tune the small encoder and the cross-attention modules.

# Benefit #1: Length generalization

Su et al., 2021. RoFormer: Enhanced Transformer with Rotary Position Embedding.  
Press et al., 2021. Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation.  
Chen et al., 2023. Extending Context Window of Large Language Models via Positional Interpolation.  
Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

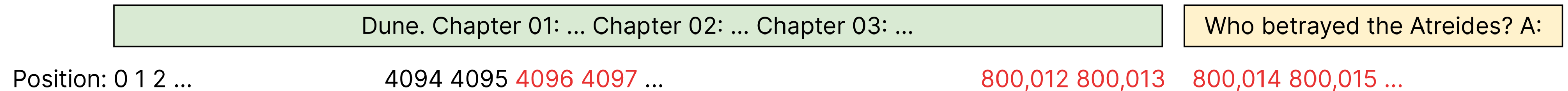
# Benefit #1: Length generalization

For the decoder-only model (assume trained on **4K** length):

Su et al., 2021. RoFormer: Enhanced Transformer with Rotary Position Embedding.  
Press et al., 2021. Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation.  
Chen et al., 2023. Extending Context Window of Large Language Models via Positional Interpolation.  
Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Benefit #1: Length generalization

For the decoder-only model (assume trained on **4K** length):

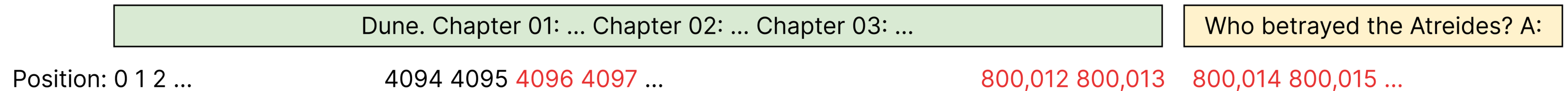


Su et al., 2021. RoFormer: Enhanced Transformer with Rotary Position Embedding.  
Press et al., 2021. Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation.  
Chen et al., 2023. Extending Context Window of Large Language Models via Positional Interpolation.  
Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.



# Benefit #1: Length generalization

For the decoder-only model (assume trained on **4K** length):

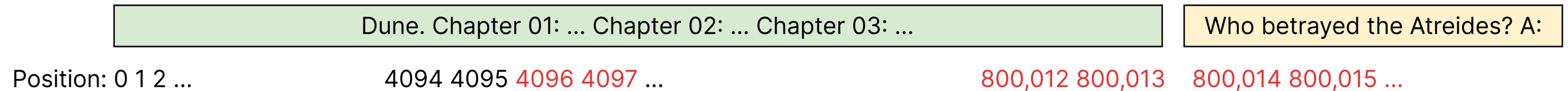


The model **does not generalize well beyond its training length**

Su et al., 2021. RoFormer: Enhanced Transformer with Rotary Position Embedding.  
Press et al., 2021. Train Short, Test Long: Attention with Linear Biases Enables Input Length Extrapolation.  
Chen et al., 2023. Extending Context Window of Large Language Models via Positional Interpolation.  
Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Benefit #1: Length generalization

For the decoder-only model (assume trained on **4K** length):

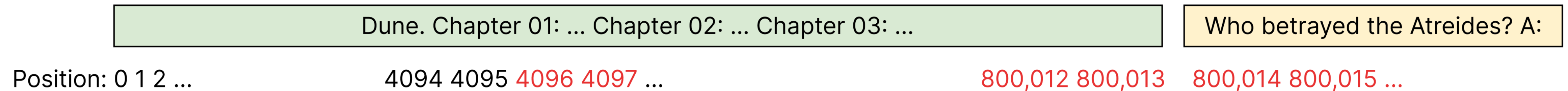


The model **does not generalize well beyond its training length**

- RoPE embeddings can generalize infinitely in theory (Su et al., 2021)

# Benefit #1: Length generalization

For the decoder-only model (assume trained on **4K** length):

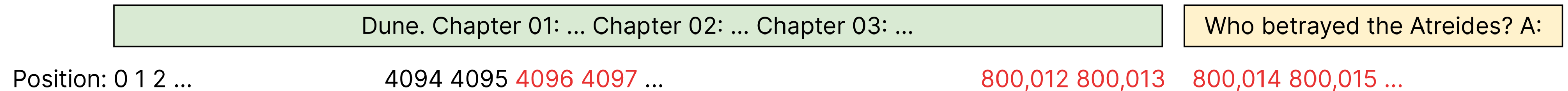


The model **does not generalize well beyond its training length**

- RoPE embeddings can generalize infinitely in theory (Su et al., 2021)
- ... but the model does not perform well on beyond-training lengths (Press et al., 2021)

# Benefit #1: Length generalization

For the decoder-only model (assume trained on **4K** length):

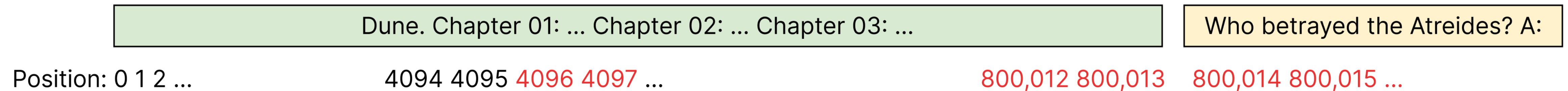


The model **does not generalize well beyond its training length**

- RoPE embeddings can generalize infinitely in theory (Su et al., 2021)
- ... but the model does not perform well on beyond-training lengths (Press et al., 2021)
- There exists positional interpolation methods (Chen et al., 2023)

# Benefit #1: Length generalization

For the decoder-only model (assume trained on **4K** length):



The model **does not generalize well beyond its training length**

- RoPE embeddings can generalize infinitely in theory (Su et al., 2021)
- ... but the model does not perform well on beyond-training lengths (Press et al., 2021)
- There exists positional interpolation methods (Chen et al., 2023)
- ... but good performance on long context still requires extensive training on the target length (Fu et al., 2024)



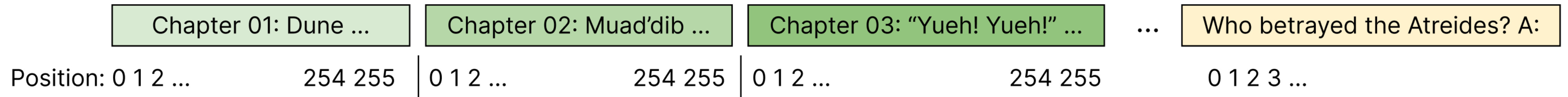
# **Benefit #1: Length generalization**

# Benefit #1: Length generalization

For **CEPE**

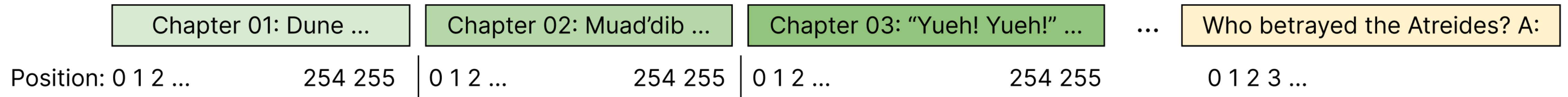
# Benefit #1: Length generalization

For **CEPE**



# Benefit #1: Length generalization

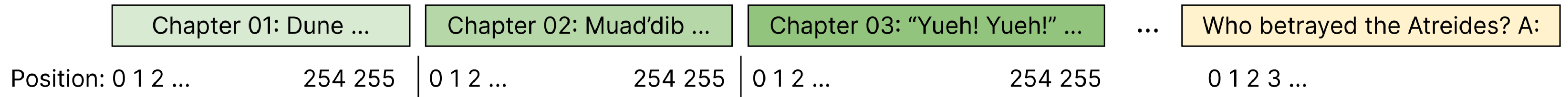
For **CEPE**



Each encoder chunk/the decoder window has its own positional encodings.

# Benefit #1: Length generalization

For **CEPE**



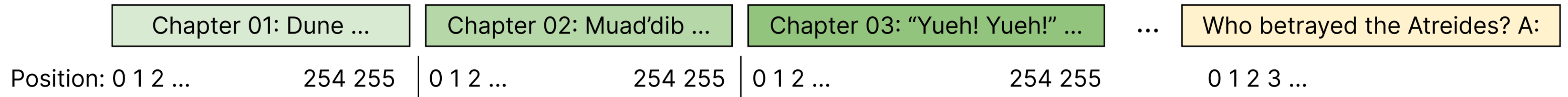
Each encoder chunk/the decoder window has its own positional encodings.

- Generalize to longer length → generalize to *more chunks*



# Benefit #1: Length generalization

For **CEPE**



Each encoder chunk/the decoder window has its own positional encodings.

- Generalize to longer length → generalize to *more chunks*
- Trained on 16 chunks, CEPE can generalize to (at least) 128 chunks

# **Benefit #2: Efficiency**

# Benefit #2: Efficiency

Throughput

# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,

# Benefit #2: Efficiency

- Throughput**      #tokens  
                         ↓
- Full attention:  $\mathcal{O}(m^2)$ ,



# Benefit #2: Efficiency

## Throughput

#tokens



- Full attention:  $\mathcal{O}(m^2)$ ,
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$

# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   
↓
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$   
↓

# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   $\downarrow$   $\#chunks$  length of the chunk
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$

# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   $\downarrow$   $\#chunks$  length of the chunk
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$
- Using a small encoder to encode most context

# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   $\downarrow$   $\#chunks$  length of the chunk
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$
- Using a small encoder to encode most context

## Memory



# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   $\downarrow$   $\#chunks$  length of the chunk
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$
- Using a small encoder to encode most context

## Memory

- Additional context:  $m$  tokens. Main input:  $n$  tokens.

# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   $\downarrow$   $\#chunks$  length of the chunk
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$
- Using a small encoder to encode most context

## Memory

- Additional context:  $m$  tokens. Main input:  $n$  tokens.
- Decoder-only:  $\mathcal{O}((m + n)Ld_{dec})$

# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   $\downarrow$   $\#chunks$  length of the chunk
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$
- Using a small encoder to encode most context

## Memory

- Additional context:  $m$  tokens. Main input:  $n$  tokens.
- Decoder-only:  $\mathcal{O}((m + n)Ld_{dec})$   
 $\#decoder\ layers$   $\swarrow$

# Benefit #2: Efficiency

## Throughput

- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   $\downarrow$
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $\#chunks$   $\swarrow$  length of the chunk  $\searrow$ ,  $m = kc$
- Using a small encoder to encode most context

## Memory

- Additional context:  $m$  tokens. Main input:  $n$  tokens.
- Decoder-only:  $\mathcal{O}((m + n)Ld_{dec})$
- Ours:  $\mathcal{O}(md_{enc} + nLd_{dec})$   $\swarrow$   $\#decoder\ layers$

# Benefit #2: Efficiency

## Throughput

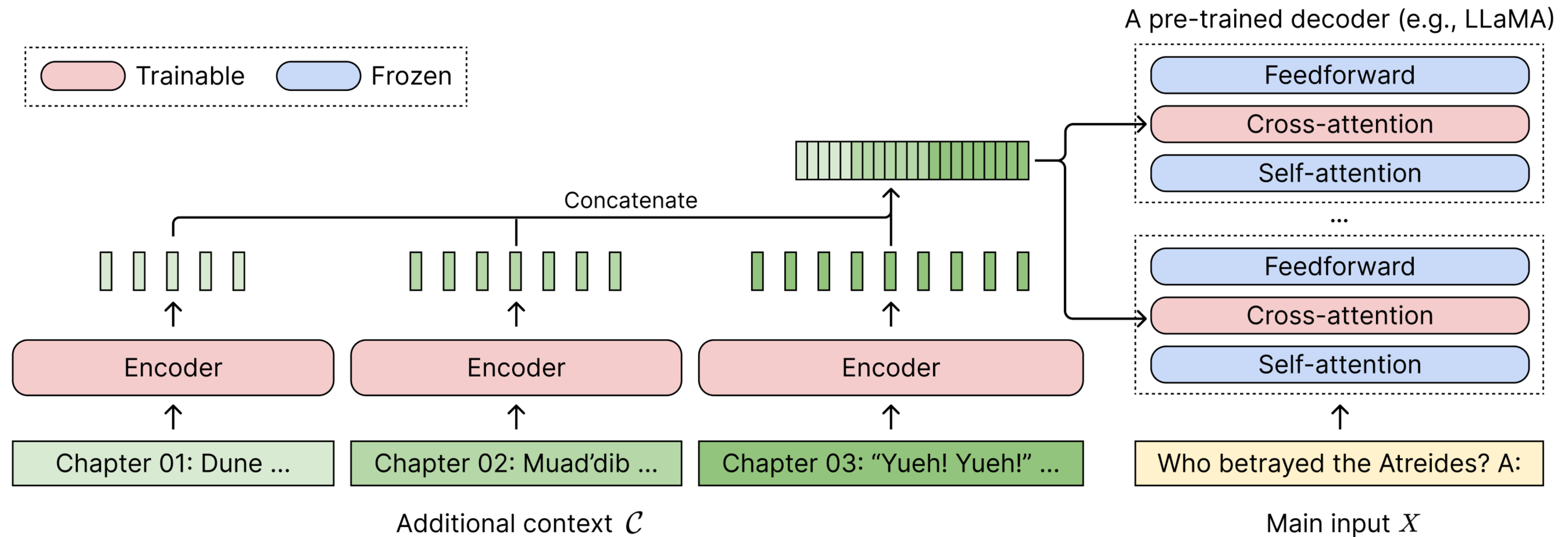
- Full attention:  $\mathcal{O}(m^2)$ ,  $\#tokens$   $\downarrow$   $\#chunks$  length of the chunk
- Attention by chunks:  $\mathcal{O}(kc^2)$ ,  $m = kc$
- Using a small encoder to encode most context

## Memory

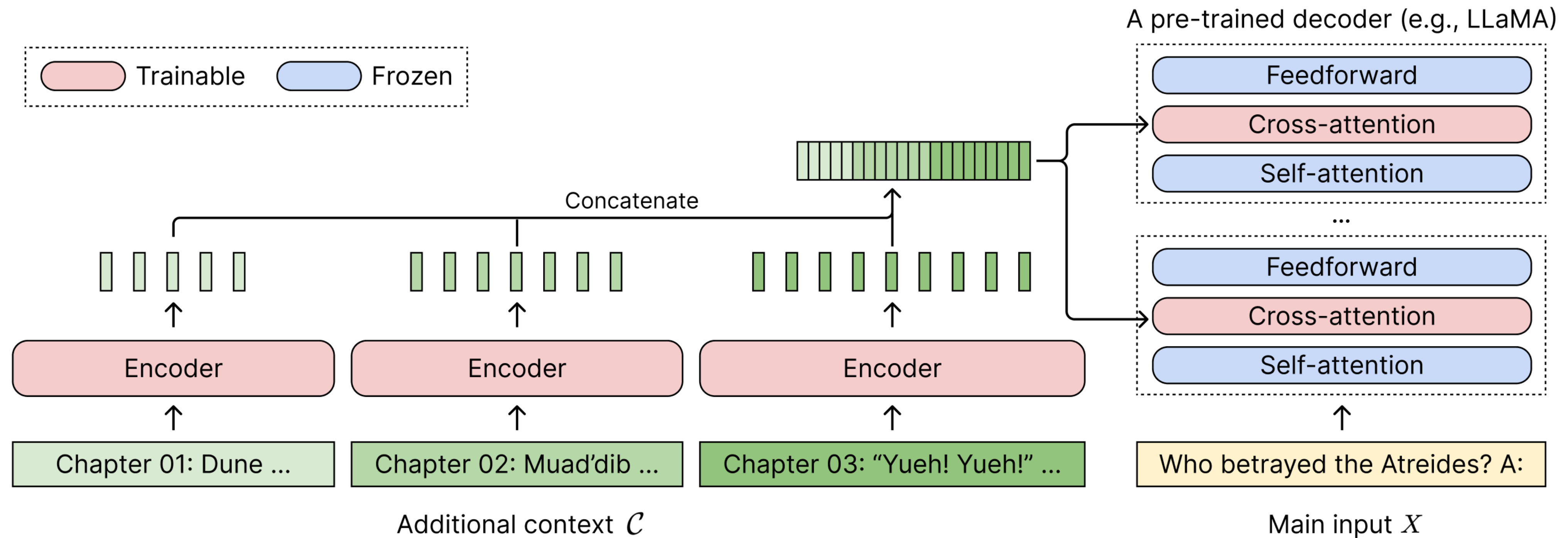
- Additional context:  $m$  tokens. Main input:  $n$  tokens.
- Decoder-only:  $\mathcal{O}((m + n)Ld_{dec})$
- Ours:  $\mathcal{O}(md_{enc} + nLd_{dec})$   $\swarrow$   $\#decoder\ layers$
- $m \gg n, d_{dec} \gg d_{enc}$



# Benefit #3: Low training cost

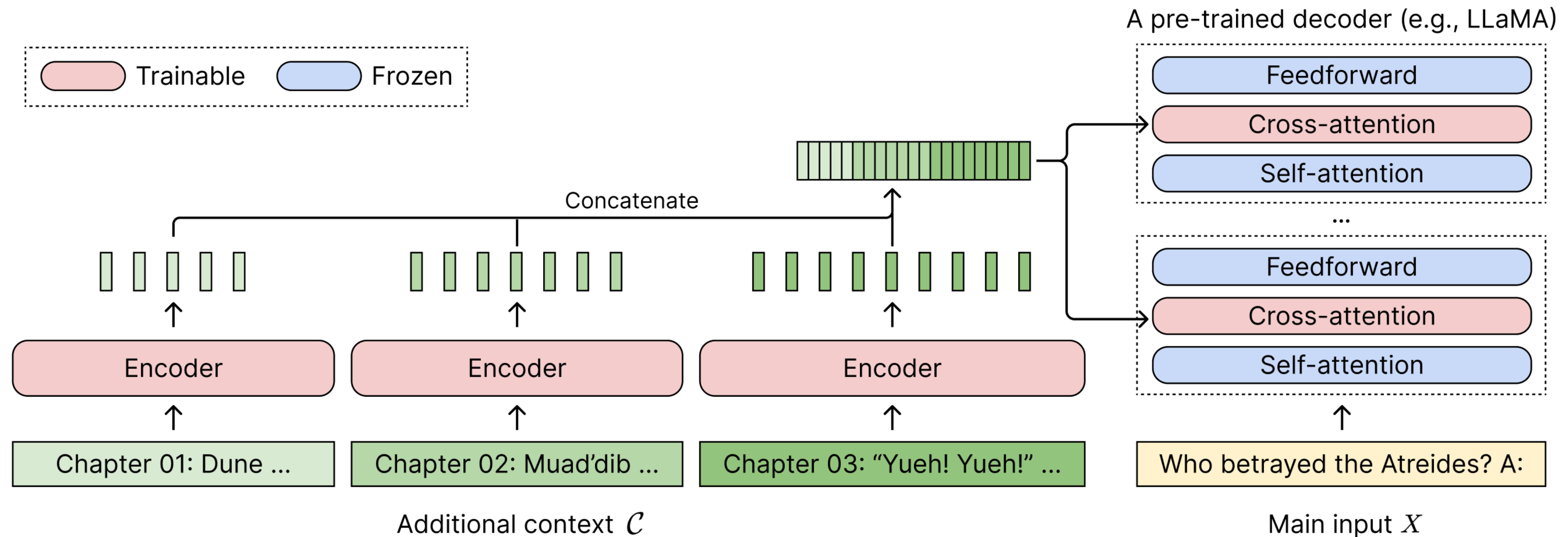


# Benefit #3: Low training cost



We only fine-tune the encoder and the cross-attention → can fit on 1 80GB GPU

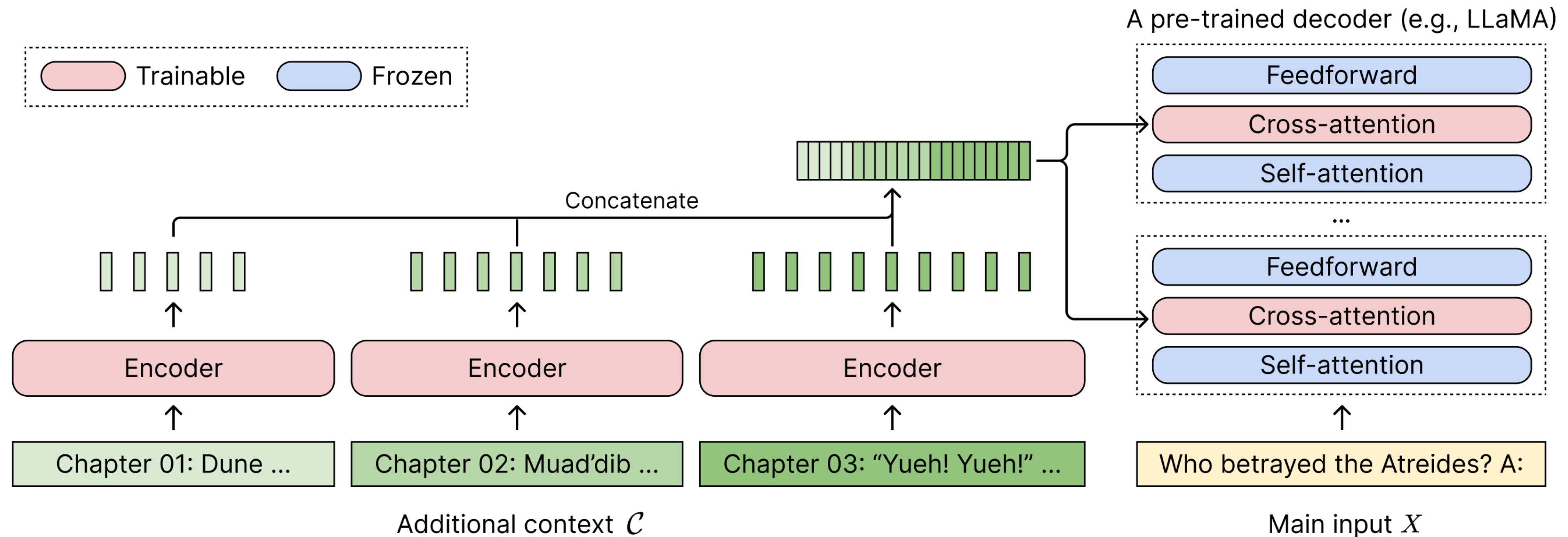
# Benefit #3: Low training cost



We only fine-tune the encoder and the cross-attention → can fit on 1 80GB GPU

- Fine-tuning a 7B decoder → at least 4 80GB GPUs

# Benefit #3: Low training cost



We only fine-tune the encoder and the cross-attention → can fit on 1 80GB GPU

- Fine-tuning a 7B decoder → at least 4 80GB GPUs

We only fine-tune on 8K sequence length → generalize to 128K

# Training data

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.



# Training data

We use RedPajama (RP; Together, 2023).

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)
- Good diversity

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)
  - Good diversity
  - Lack long-range dependency

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)
  - Good diversity
  - Lack long-range dependency

Domain	Avg #tokens
Arxiv	16K
Book	142K
C4 (clean web)	0.5K
CC (web)	2K
Github	2K
StackExchange	0.5K
Wikipedia	0.5K

Different domains from RP.

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.



# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)
  - Good diversity
  - Lack long-range dependency
- **Length-filter recipe:** long documents from Arxiv and Book (RP-filter)

Domain	Avg #tokens
Arxiv	16K
Book	142K
C4 (clean web)	0.5K
CC (web)	2K
Github	2K
StackExchange	0.5K
Wikipedia	0.5K

Different domains from RP.

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)
  - Good diversity
  - Lack long-range dependency
- **Length-filter recipe:** long documents from Arxiv and Book (RP-filter)
  - Strong long-range dependency

Domain	Avg #tokens
Arxiv	16K
Book	142K
C4 (clean web)	0.5K
CC (web)	2K
Github	2K
StackExchange	0.5K
Wikipedia	0.5K

Different domains from RP.

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)
  - Good diversity
  - Lack long-range dependency
- **Length-filter recipe:** long documents from Arxiv and Book (RP-filter)
  - Strong long-range dependency
  - Lack diversity

Domain	Avg #tokens
Arxiv	16K
Book	142K
C4 (clean web)	0.5K
CC (web)	2K
Github	2K
StackExchange	0.5K
Wikipedia	0.5K

Different domains from RP.

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)
  - Good diversity
  - Lack long-range dependency
- **Length-filter recipe:** long documents from Arxiv and Book (RP-filter)
  - Strong long-range dependency
  - Lack diversity

Domain	Avg #tokens
Arxiv	16K
Book	142K
C4 (clean web)	0.5K
CC (web)	2K
Github	2K
StackExchange	0.5K
Wikipedia	0.5K

Different domains from RP.

**We use a 2:1 mixture of RP-filter and RP-concat,**

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.

# Training data

We use RedPajama (RP; Together, 2023).

- **Standard recipe:** concatenate the documents to the target length (RP-concat)
  - Good diversity
  - Lack long-range dependency
- **Length-filter recipe:** long documents from Arxiv and Book (RP-filter)
  - Strong long-range dependency
  - Lack diversity

Domain	Avg #tokens
Arxiv	16K
Book	142K
C4 (clean web)	0.5K
CC (web)	2K
Github	2K
StackExchange	0.5K
Wikipedia	0.5K

Different domains from RP.

**We use a 2:1 mixture of RP-filter and RP-concat,**  
a design choice echoing Together, 2023; Fu et al., 2024.

Together, 2023. RedPajama, a project to create leading open-source models, starts by reproducing LLaMA training dataset of over 1.2 trillion tokens.

Together, 2023. Preparing for the era of 32K context: Early learnings and explorations.

Fu et al., 2024. Data Engineering for Scaling Language Models to 128K Context.



# Training details

# Training details

Encoder

# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling

# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

## Training CEPE



# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

## Training CEPE

- We train CEPE with LLaMA-2-7B for 20B tokens (1% of LLaMA-2's pre-training budget)

# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

## Training CEPE

- We train CEPE with LLaMA-2-7B for 20B tokens (1% of LLaMA-2's pre-training budget)
- Decoder input = **4K** tokens

# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

## Training CEPE

- We train CEPE with LLaMA-2-7B for 20B tokens (1% of LLaMA-2's pre-training budget)
- Decoder input = **4K** tokens
- Encoder input = **16 x 256** tokens

# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

## Training CEPE

- We train CEPE with LLaMA-2-7B for 20B tokens (1% of LLaMA-2's pre-training budget)
- Decoder input = **4K** tokens
- Encoder input = **16 x 256** tokens

## Warmup

# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

## Training CEPE

- We train CEPE with LLaMA-2-7B for 20B tokens (1% of LLaMA-2's pre-training budget)
- Decoder input = **4K** tokens
- Encoder input = **16 x 256** tokens

## Warmup

We first train the cross-attention to learn to copy



# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

## Training CEPE

- We train CEPE with LLaMA-2-7B for 20B tokens (1% of LLaMA-2's pre-training budget)
- Decoder input = **4K** tokens
- Encoder input = **16 x 256** tokens

## Warmup

We first train the cross-attention to learn to copy

- We set the encoder and the decoder input to be the same (for 131M tokens)

# Training details

## Encoder

- A train-from-scratch **435M** bidirectional encoder using masked language modeling
- Same vocabulary as LLaMA

## Training CEPE

- We train CEPE with LLaMA-2-7B for 20B tokens (1% of LLaMA-2's pre-training budget)
- Decoder input = **4K** tokens
- Encoder input = **16 x 256** tokens

## Warmup

We first train the cross-attention to learn to copy

- We set the encoder and the decoder input to be the same (for 131M tokens)

# **CEPE for instruction-tuned models**

# CEPE for instruction-tuned models

Instruction-tuning / chat-tuning

# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)

# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)

Who is the president of the United States?

A: Donald Trump

Q: What are three scientific advancements?

A: Airline travel, cars, and space travel

Q: Where were the first humans found?

A: East Africa

Q: Where did humans emigrate to prior to finding the new world?

A: Asia



# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)

Who is the president of the United States?

A: Donald Trump

Q: What are three scientific advancements?

A: Airline travel, cars, and space travel



Q: Where were the first humans found?

A: East Africa

Q: Where did humans emigrate to prior to finding the new world?

A: Asia

# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)

Who is the president of the United States?

A: Donald Trump

Q: What are three scientific advancements?

A: Airline travel, cars, and space travel

Q: Where were the first humans found?

A: East Africa

Q: Where did humans emigrate to prior to finding the new world?

A: Asia



USER

Who is the president of the United States?

ASSISTANT

As of my last update in April 2023, Joe Biden is the President of the United States.

# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)
- How to turn a long-context LM to a long-context **chat** LM?

Who is the president of the United States?

A: Donald Trump

Q: What are three scientific advancements?

A: Airline travel, cars, and space travel

Q: Where were the first humans found?

A: East Africa

Q: Where did humans emigrate to prior to finding the new world?

A: Asia



USER

Who is the president of the United States?

ASSISTANT

As of my last update in April 2023, Joe Biden is the President of the United States.

# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)
- How to turn a long-context LM to a long-context **chat** LM?
  - Tune a chat LM on long-context data

Who is the president of the United States?

A: Donald Trump

Q: What are three scientific advancements?

A: Airline travel, cars, and space travel

Q: Where were the first humans found?

A: East Africa

Q: Where did humans emigrate to prior to finding the new world?

A: Asia



USER

Who is the president of the United States?

ASSISTANT

As of my last update in April 2023, Joe Biden is the President of the United States.

# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)
- How to turn a long-context LM to a long-context **chat** LM?
  - Tune a chat LM on long-context data → **lose chat abilities**

Who is the president of the United States?

A: Donald Trump

Q: What are three scientific advancements?

A: Airline travel, cars, and space travel

Q: Where were the first humans found?

A: East Africa

Q: Where did humans emigrate to prior to finding the new world?

A: Asia



USER

Who is the president of the United States?

ASSISTANT

As of my last update in April 2023, Joe Biden is the President of the United States.

# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)
- How to turn a long-context LM to a long-context **chat** LM?
  - Tune a chat LM on long-context data → **lose chat abilities**
  - Tune a long-context LM on chat data

Who is the president of the United States?

A: Donald Trump

Q: What are three scientific advancements?

A: Airline travel, cars, and space travel

Q: Where were the first humans found?

A: East Africa

Q: Where did humans emigrate to prior to finding the new world?

A: Asia



USER

Who is the president of the United States?

ASSISTANT

As of my last update in April 2023, Joe Biden is the President of the United States.



# CEPE for instruction-tuned models

## Instruction-tuning / chat-tuning

- Most useful models are fine-tuned on *chat-like data* (often proprietary)
- How to turn a long-context LM to a long-context **chat** LM?
  - Tune a chat LM on long-context data → **lose chat abilities**
  - Tune a long-context LM on chat data → **no proprietary data; no long-context chat data**

Who is the president of the United States?

A: Donald Trump

Q: What are three scientific advancements?

A: Airline travel, cars, and space travel

Q: Where were the first humans found?

A: East Africa

Q: Where did humans emigrate to prior to finding the new world?

A: Asia



USER

Who is the president of the United States?

ASSISTANT

As of my last update in April 2023, Joe Biden is the President of the United States.

# **CEPE for instruction-tuned models**

# CEPE for instruction-tuned models

CEPE-Distilled (CEPED)

# CEPE for instruction-tuned models

## CEPE-Distilled (CEPED)

- Use only *unsupervised* long-context data, we can turn a chat model to a *long-context* model

# CEPE for instruction-tuned models

## CEPE-Distilled (CEPED)

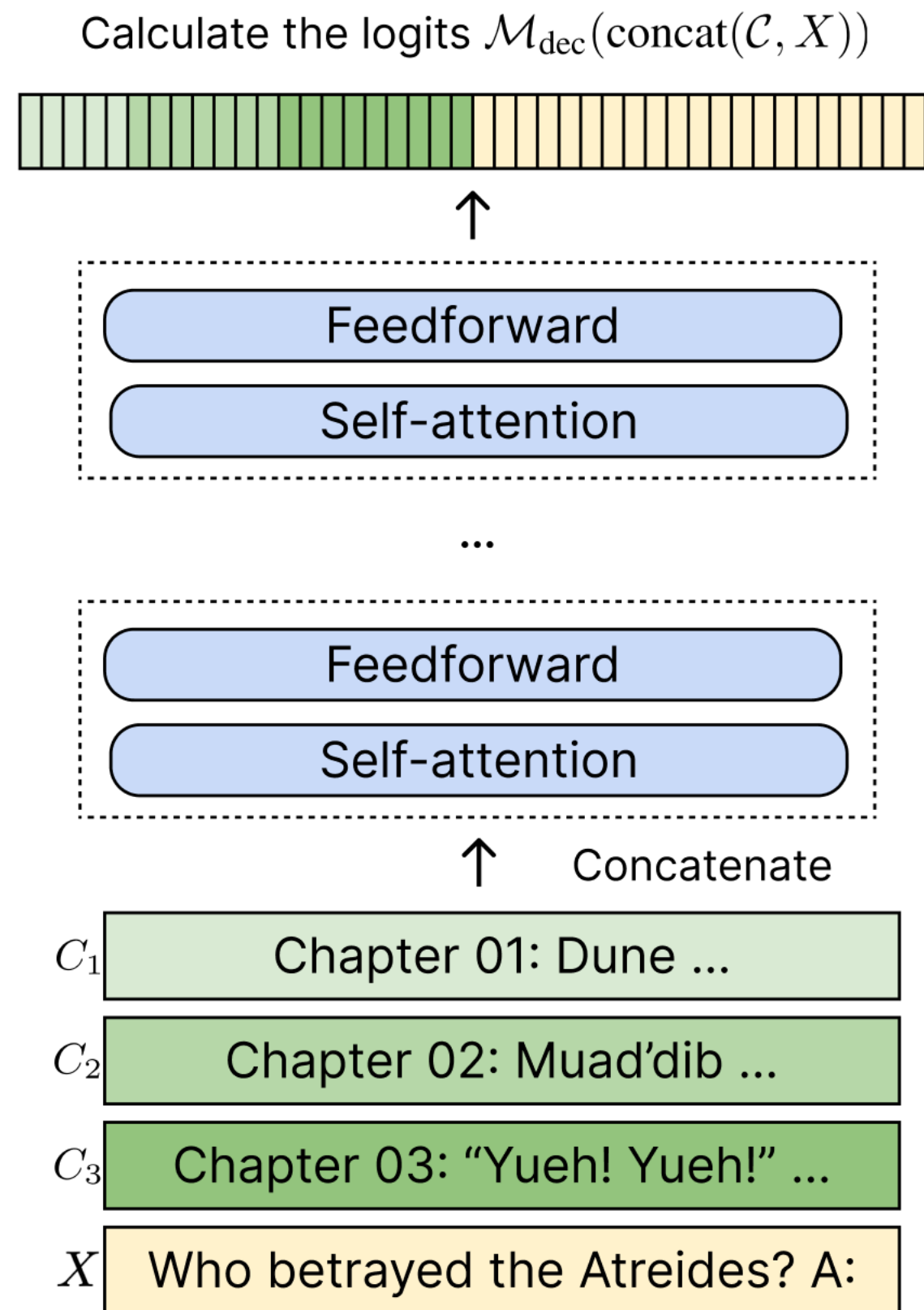
- Use only *unsupervised* long-context data, we can turn a chat model to a *long-context* model
- Add **an auxiliary distillation loss** to maintain the chat model's chat ability

# CEPE-Distilled (CEPED)

$C_1$	Chapter 01: Dune ...
$C_2$	Chapter 02: Muad'dib ...
$C_3$	Chapter 03: "Yueh! Yueh!" ...
$X$	Who betrayed the Atreides? A:

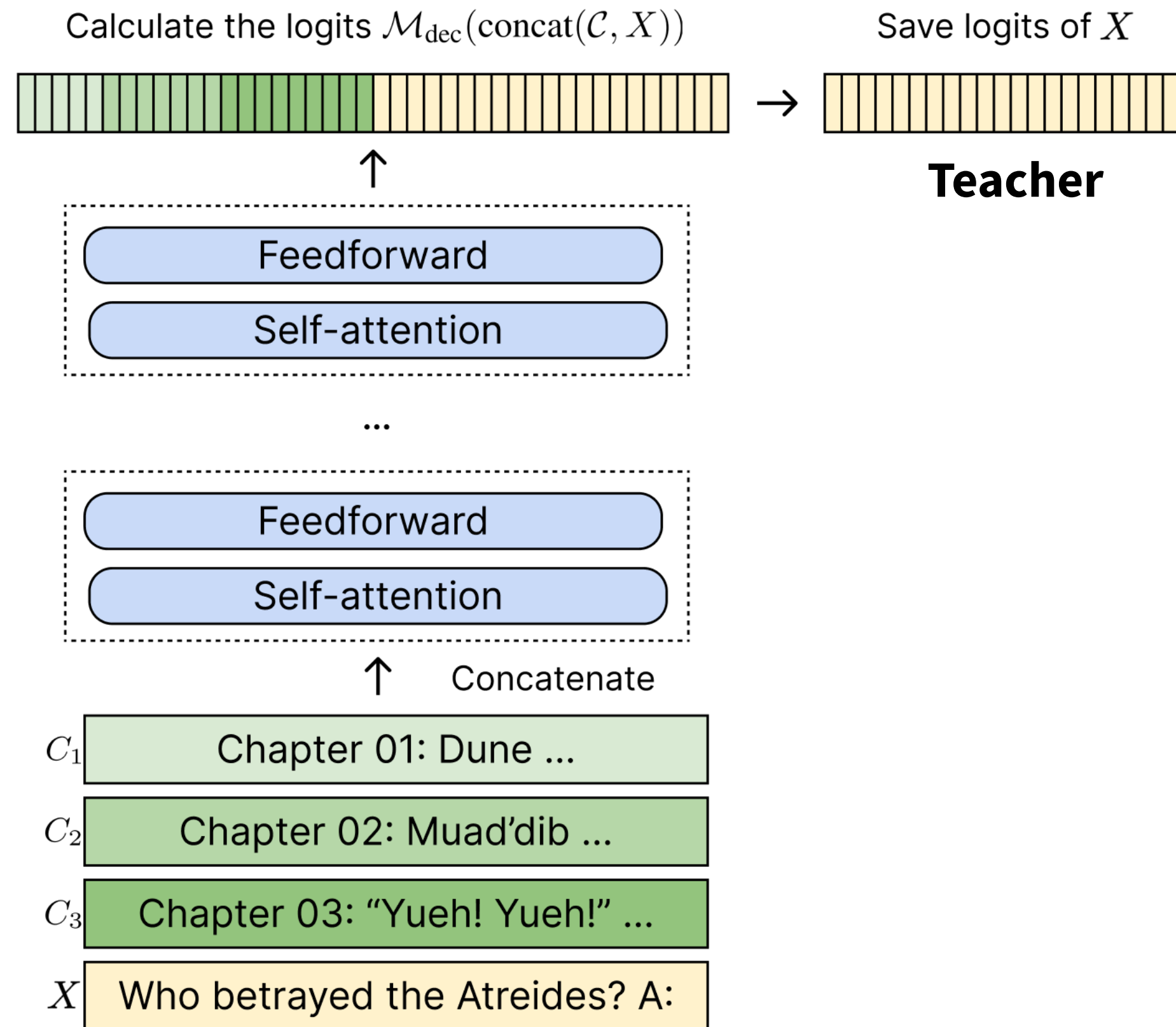


# CEPE-Distilled (CEPED)



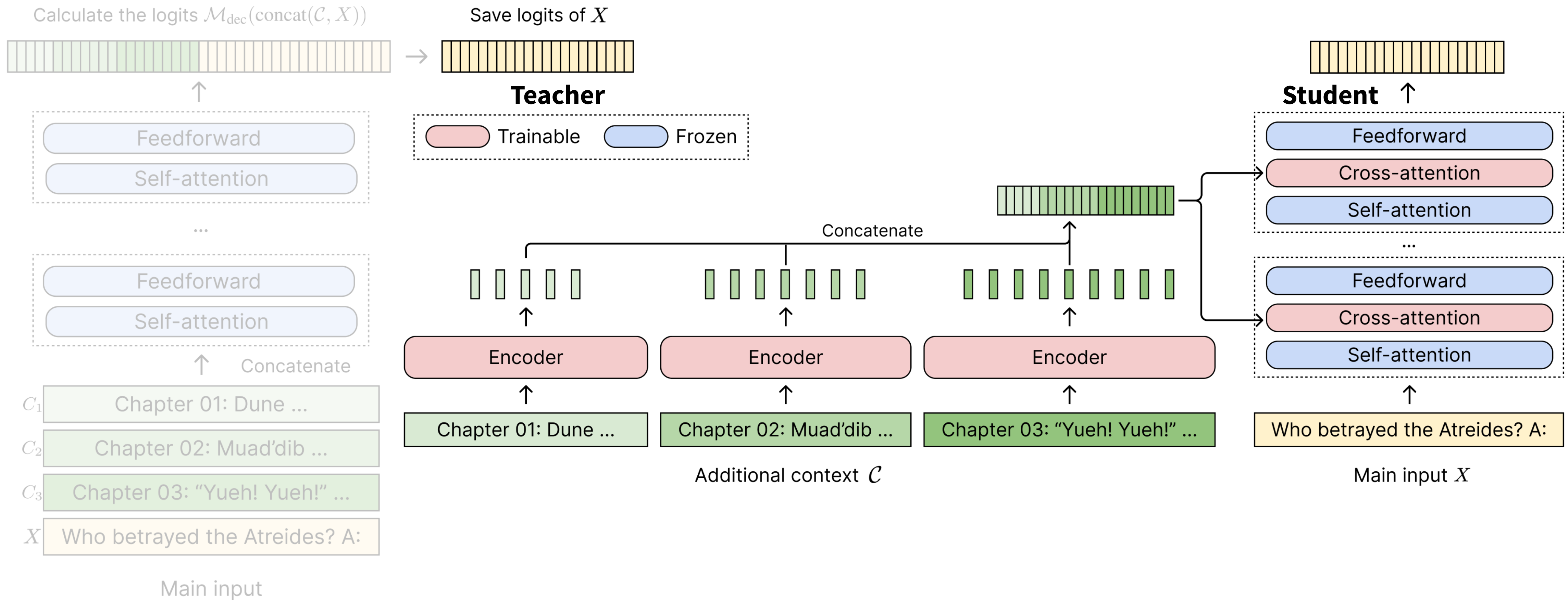
1. Run forward passes w/ the original chat model

# CEPE-Distilled (CEPED)



1. Run forward passes w/ the original chat model

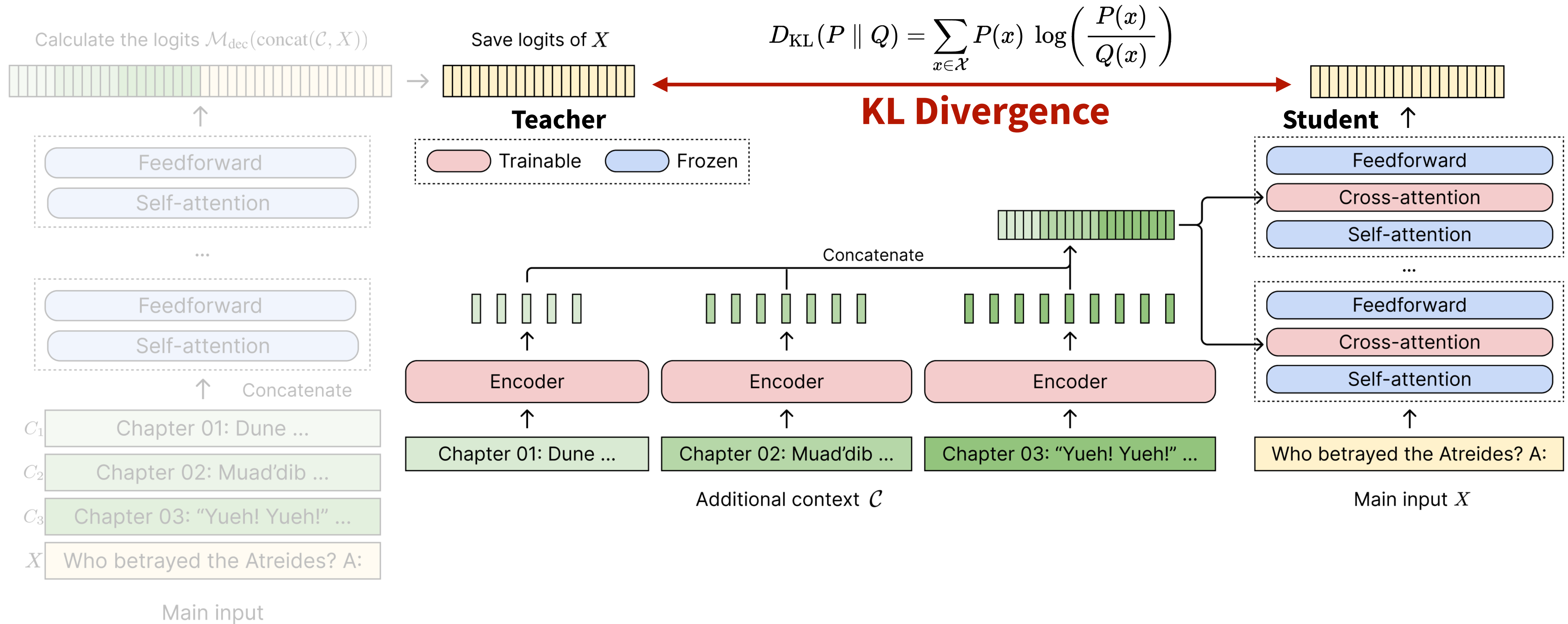
# CEPE-Distilled (CEPED)



1. Run forward passes w/ the original chat model

2. Run forward passes w/ the CEPE model

# CEPE-Distilled (CEPED)



1. Run forward passes w/ the original chat model

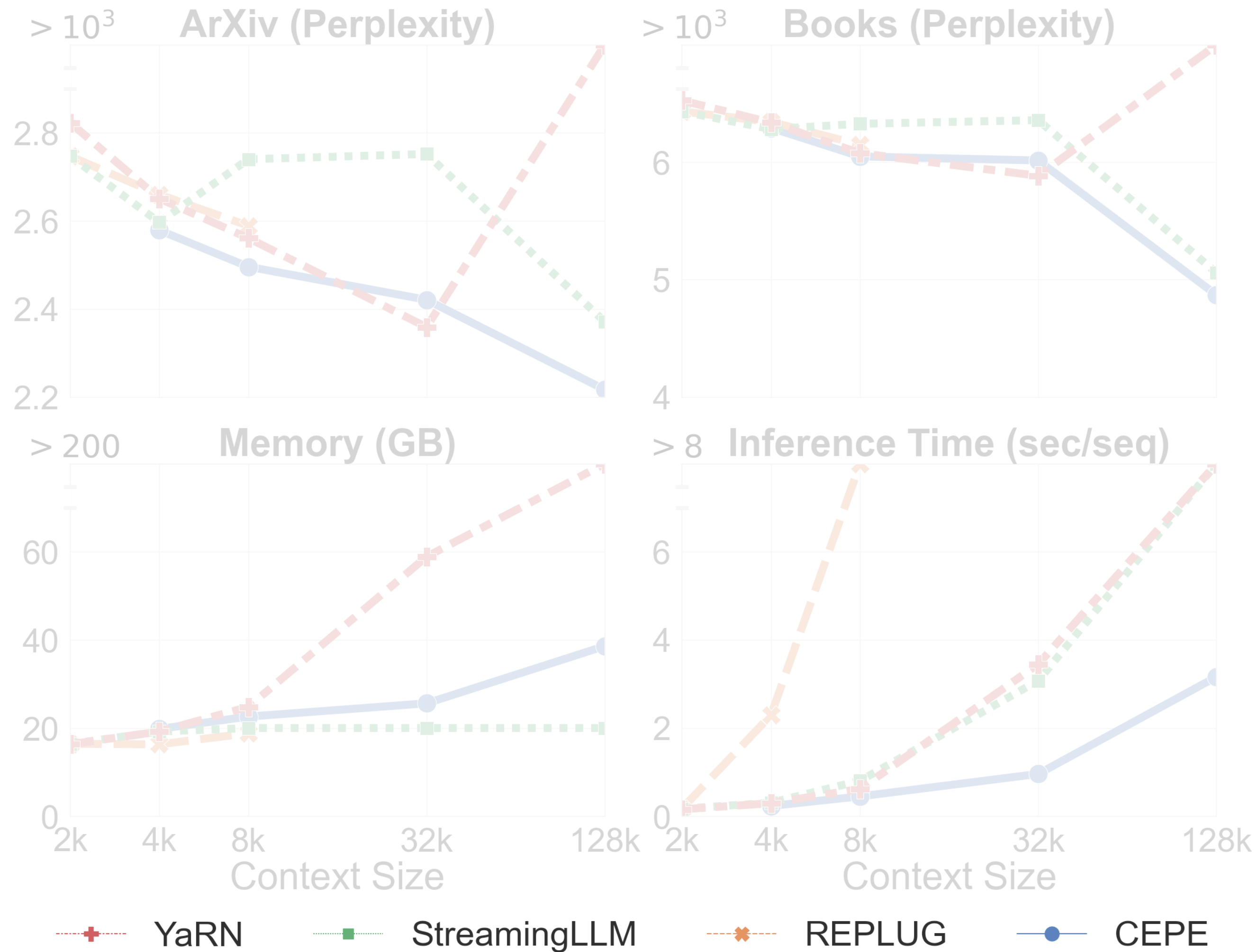
2. Run forward passes w/ the CEPE model

3. Train with KL Divergence loss + cross-entropy loss

# Evaluation

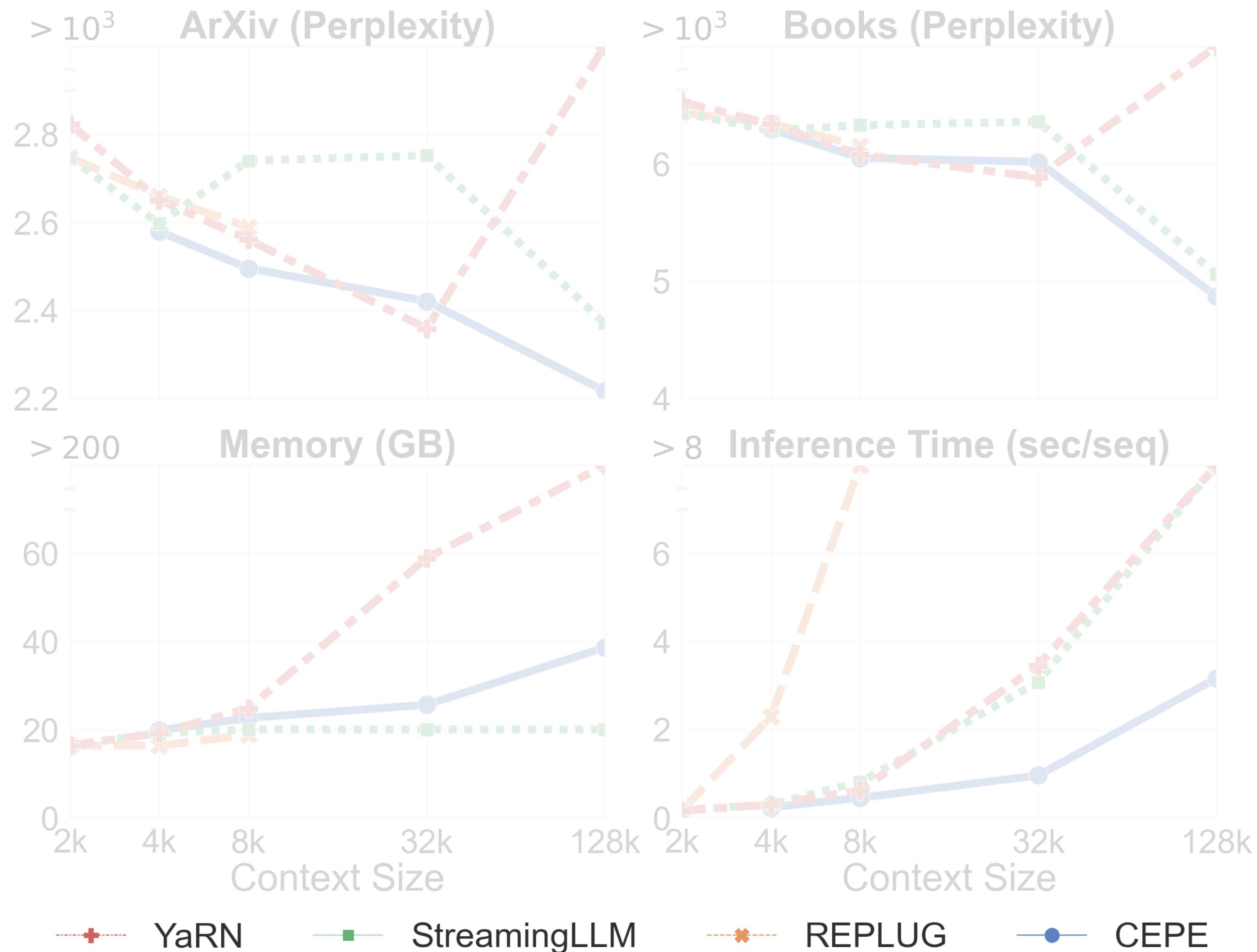
- Long-context language modeling
- Retrieval-augmented applications
- In-context learning
- Chat model evaluation

# Evaluation - long-context language modeling



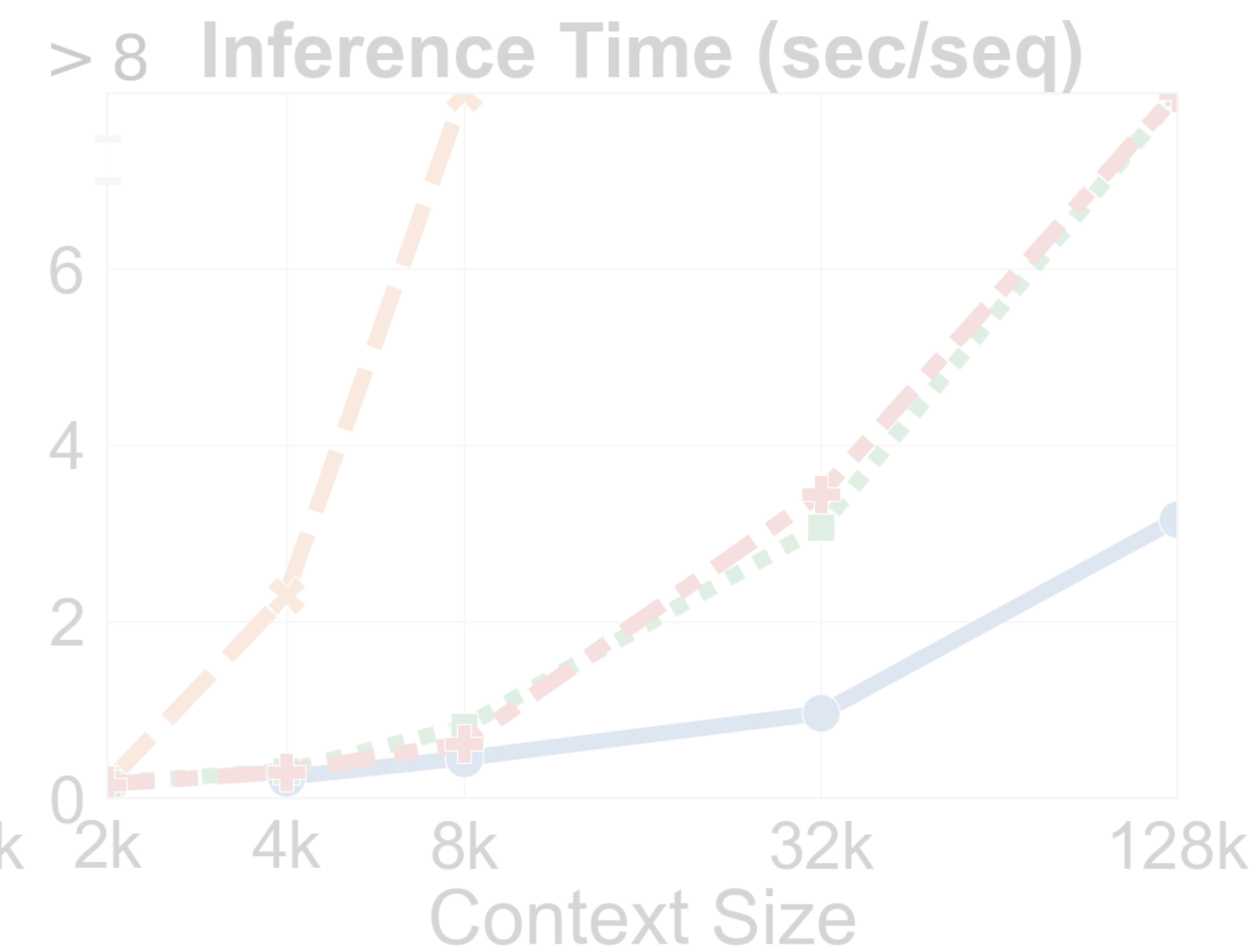
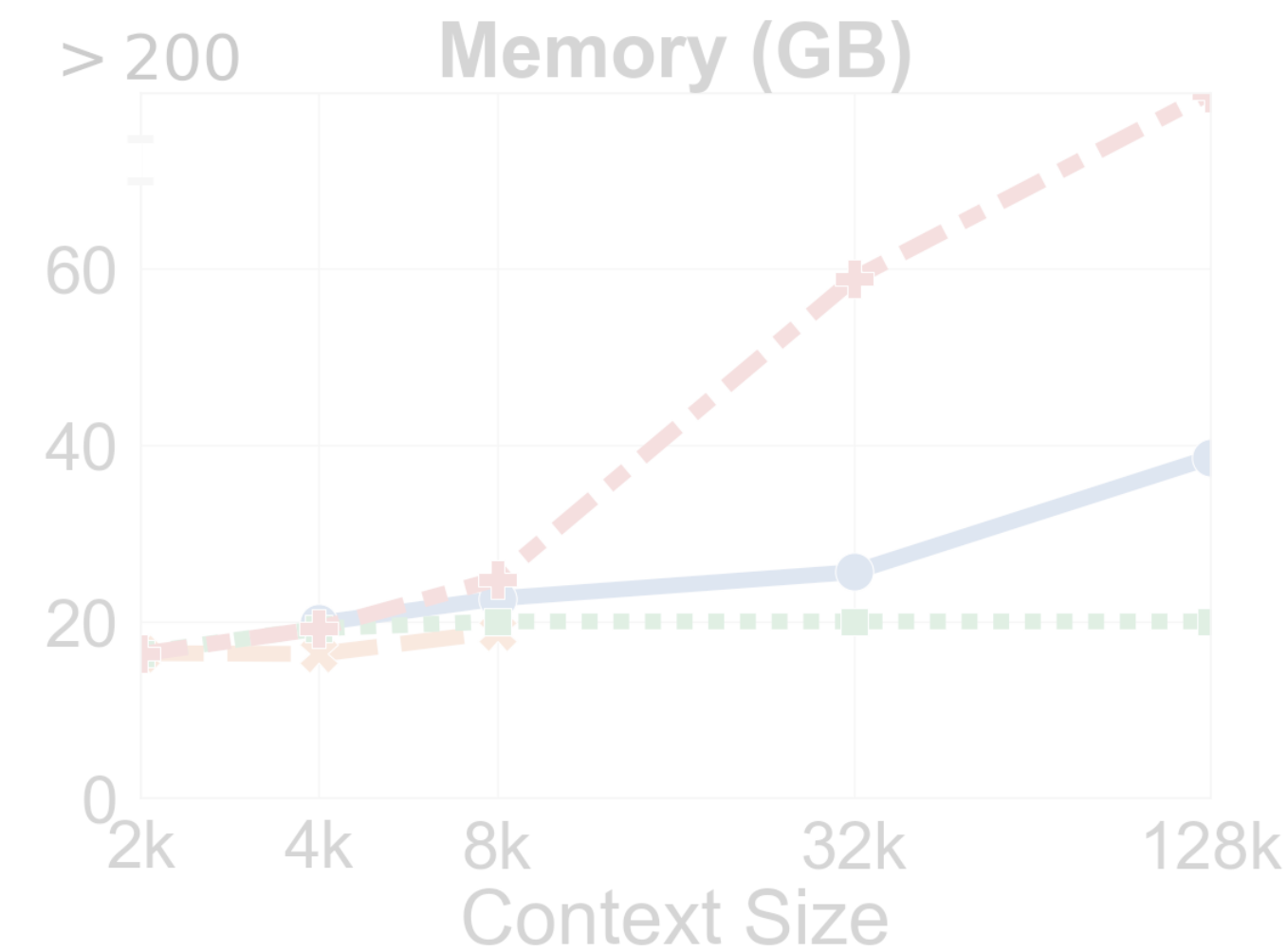
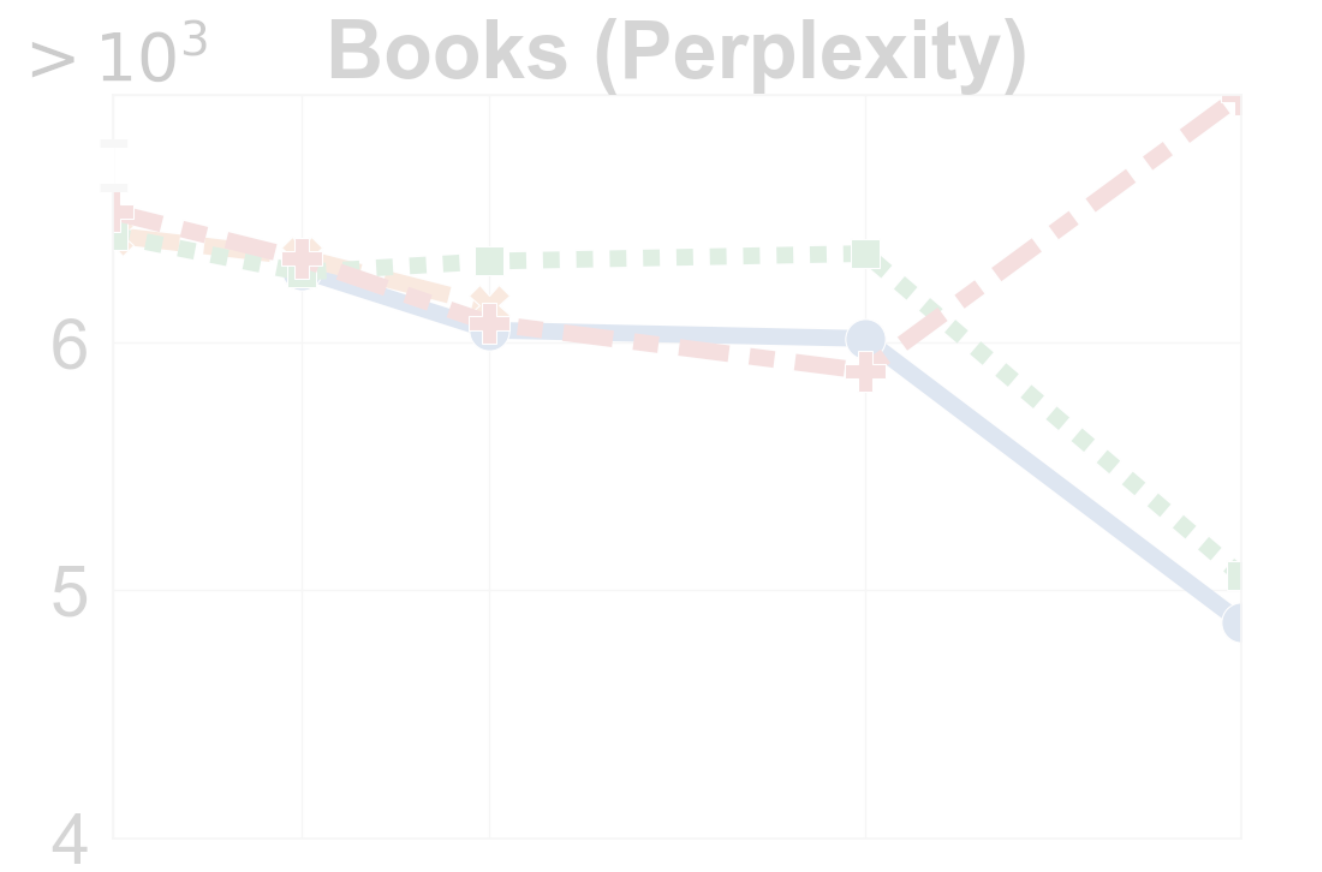
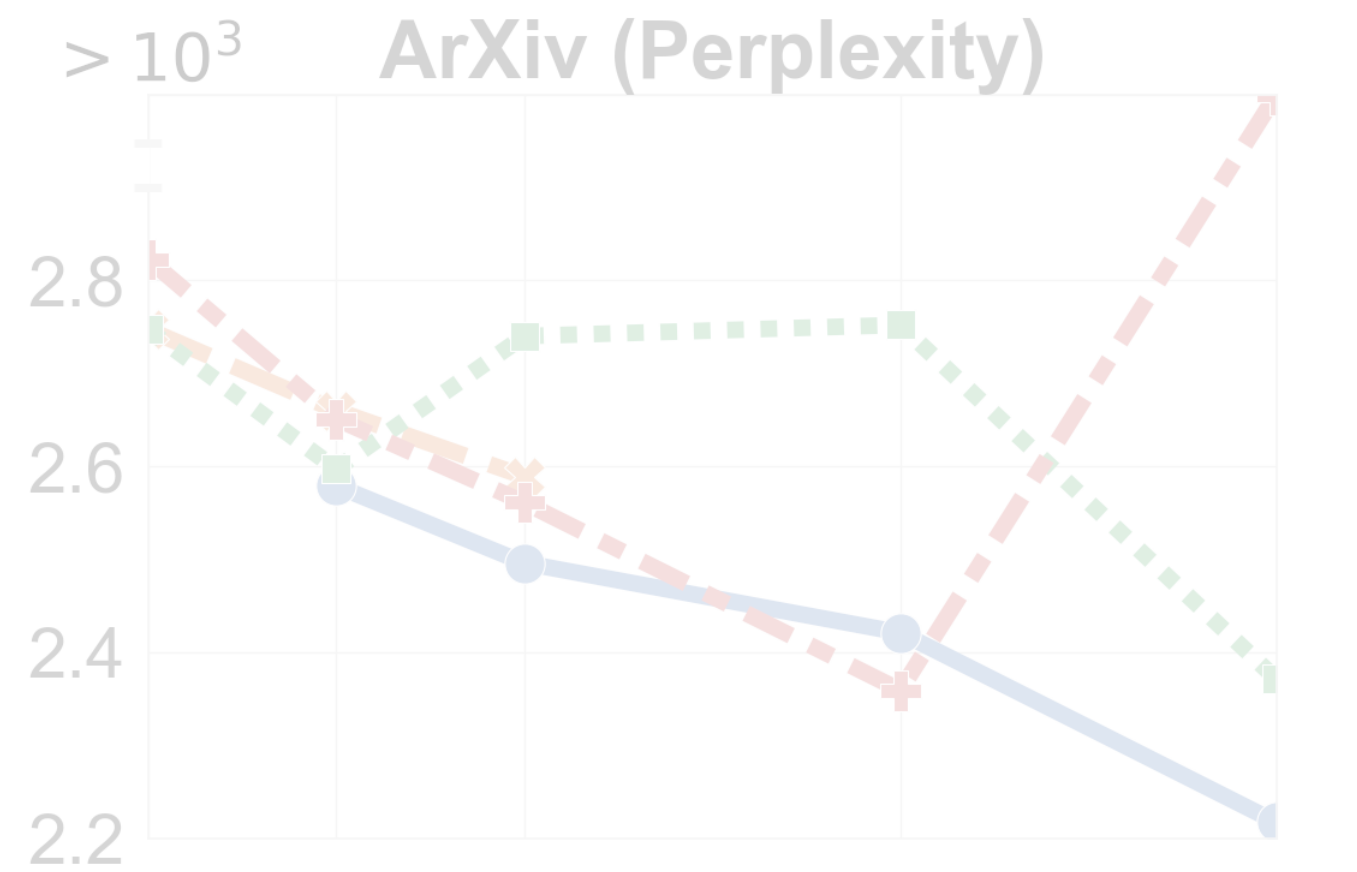


# Evaluation - long-context language modeling



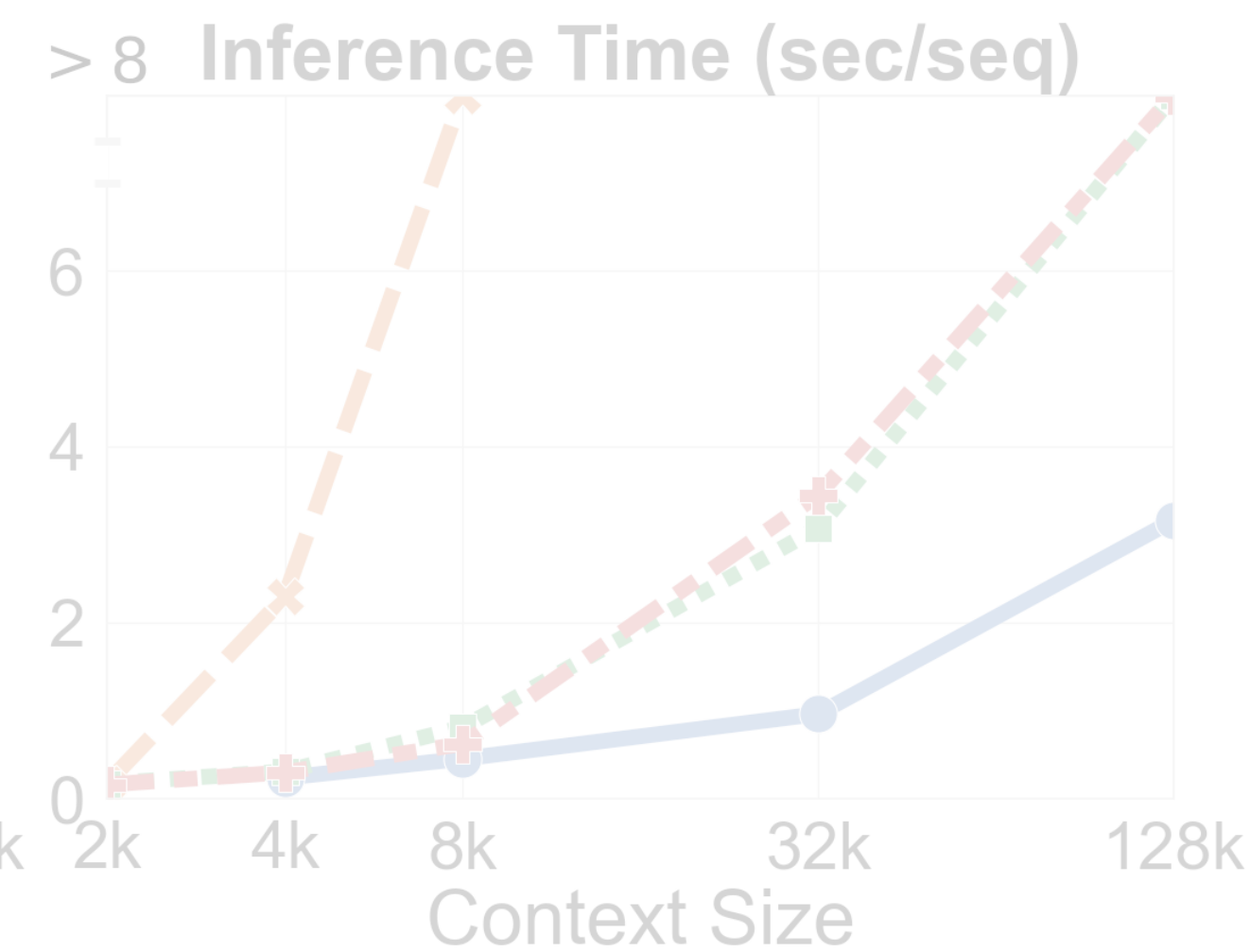
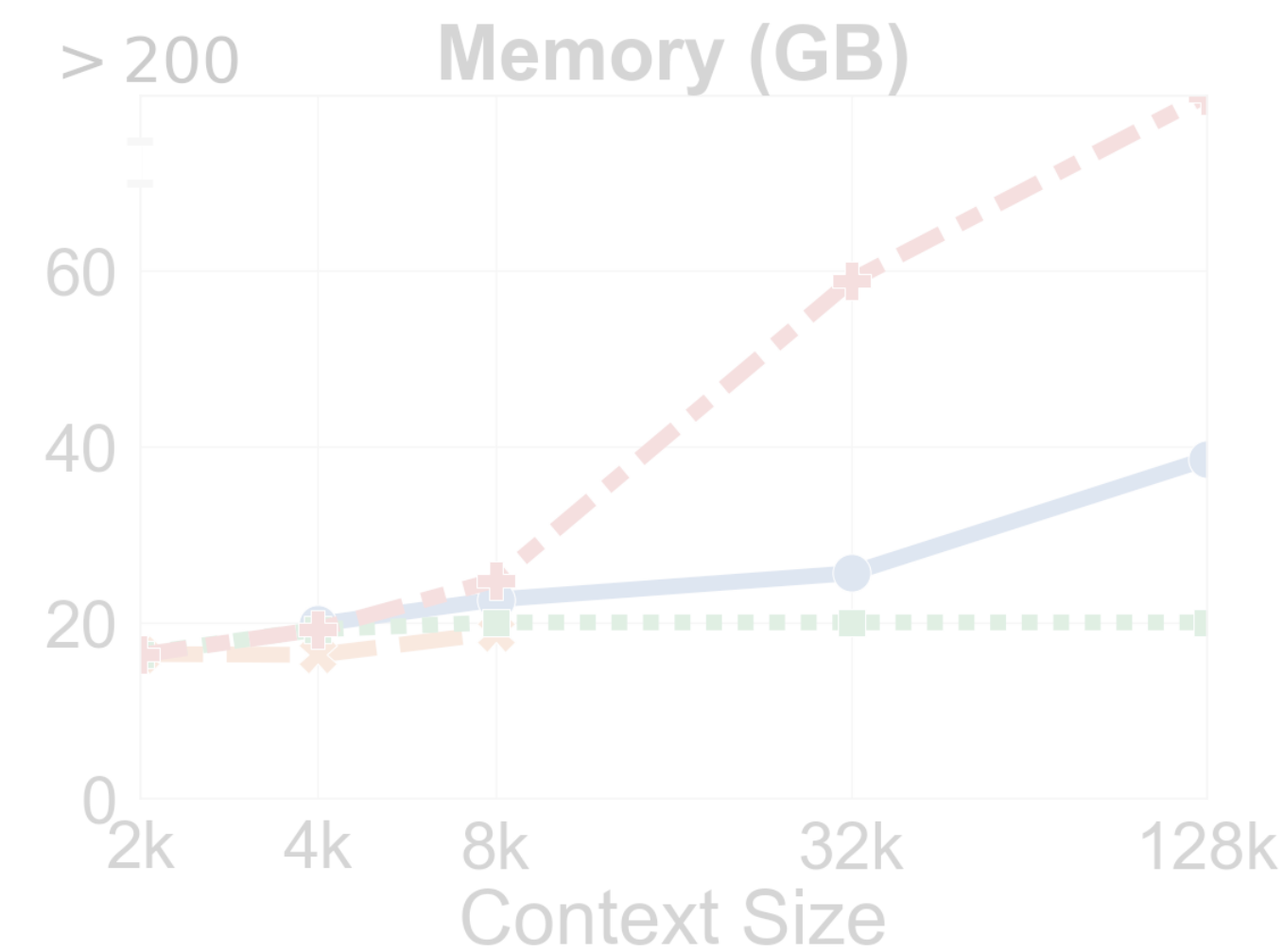
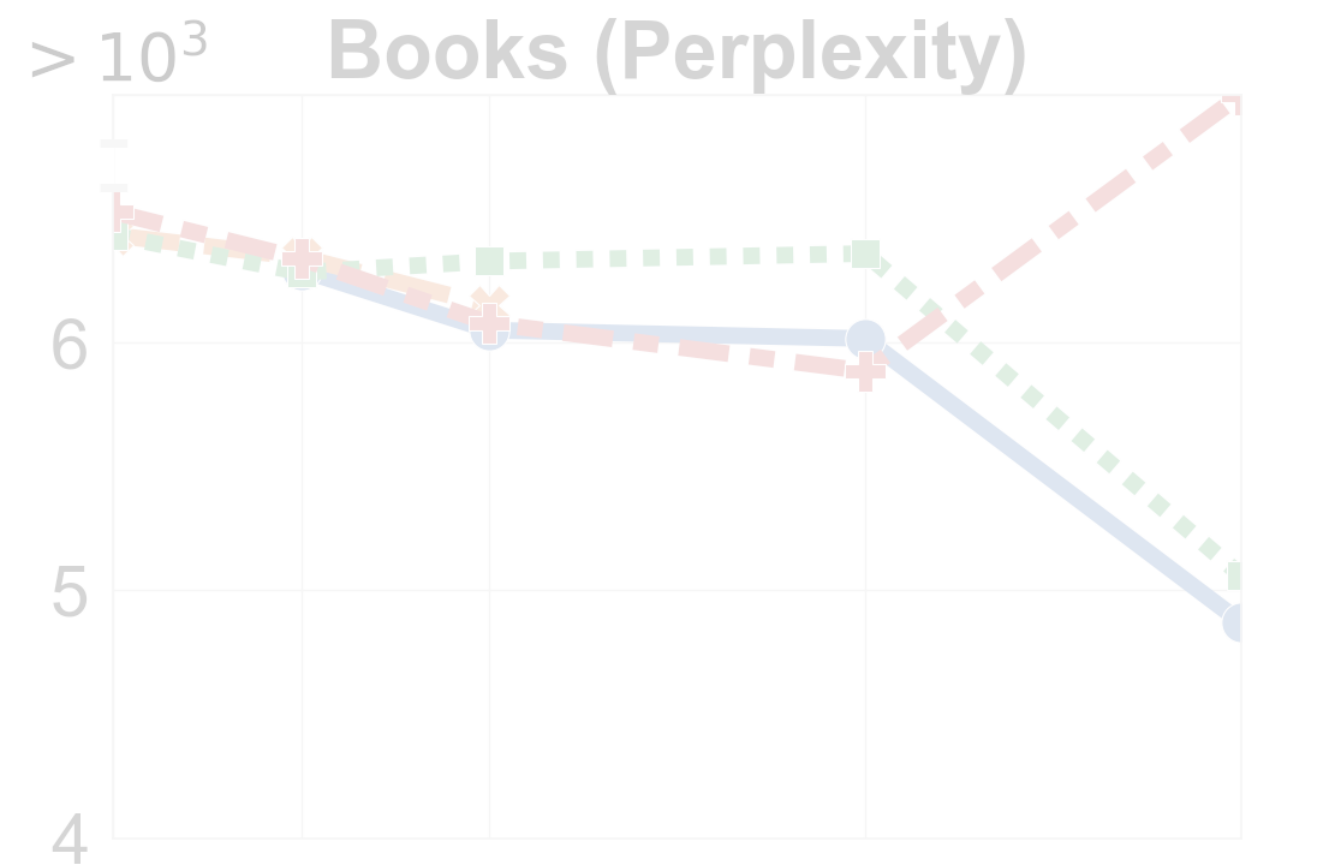
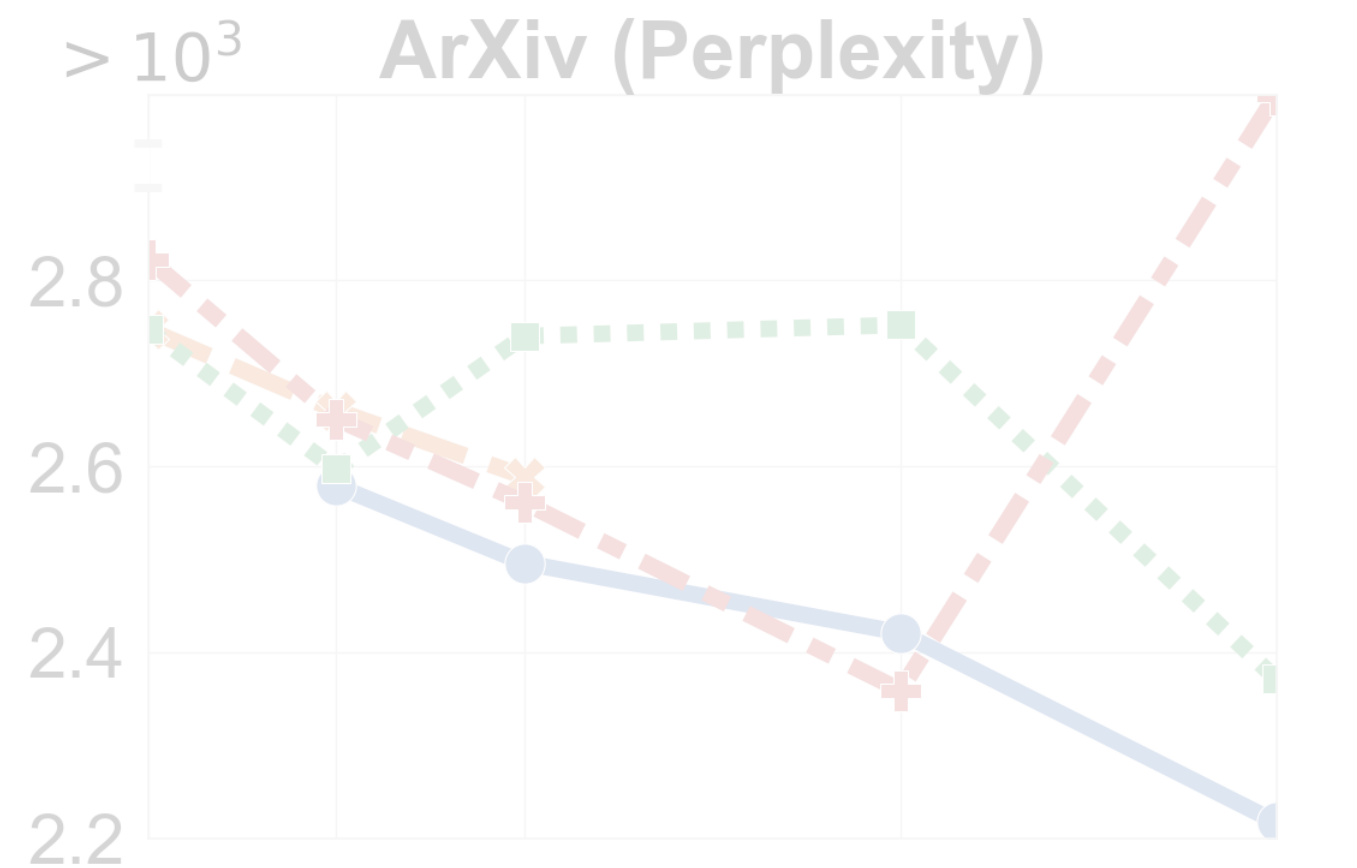
Positional interpolation

# Evaluation - long-context language modeling



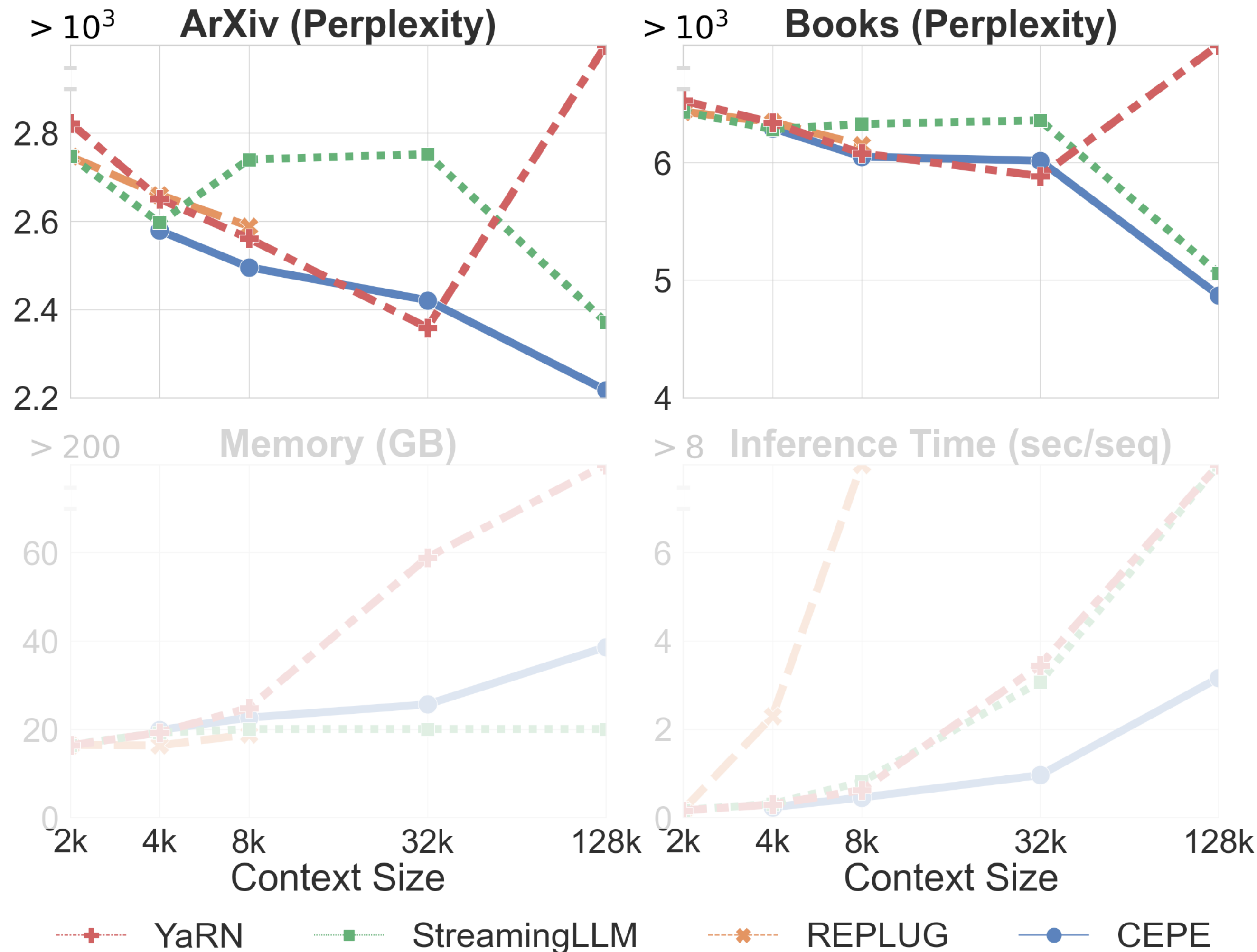
- - + - - YaRN      - - ■ - - StreamingLLM  
- - \* - - REPLUG      - - ● - - CEPE  
Positional interpolation      KV-cache dropping

# Evaluation - long-context language modeling



- - + - - YaRN      - - ■ - - StreamingLLM      - - \* - - REPLUG      - - ● - - CEPE  
Positional interpolation      KV-cache dropping      Retrieval for long-context

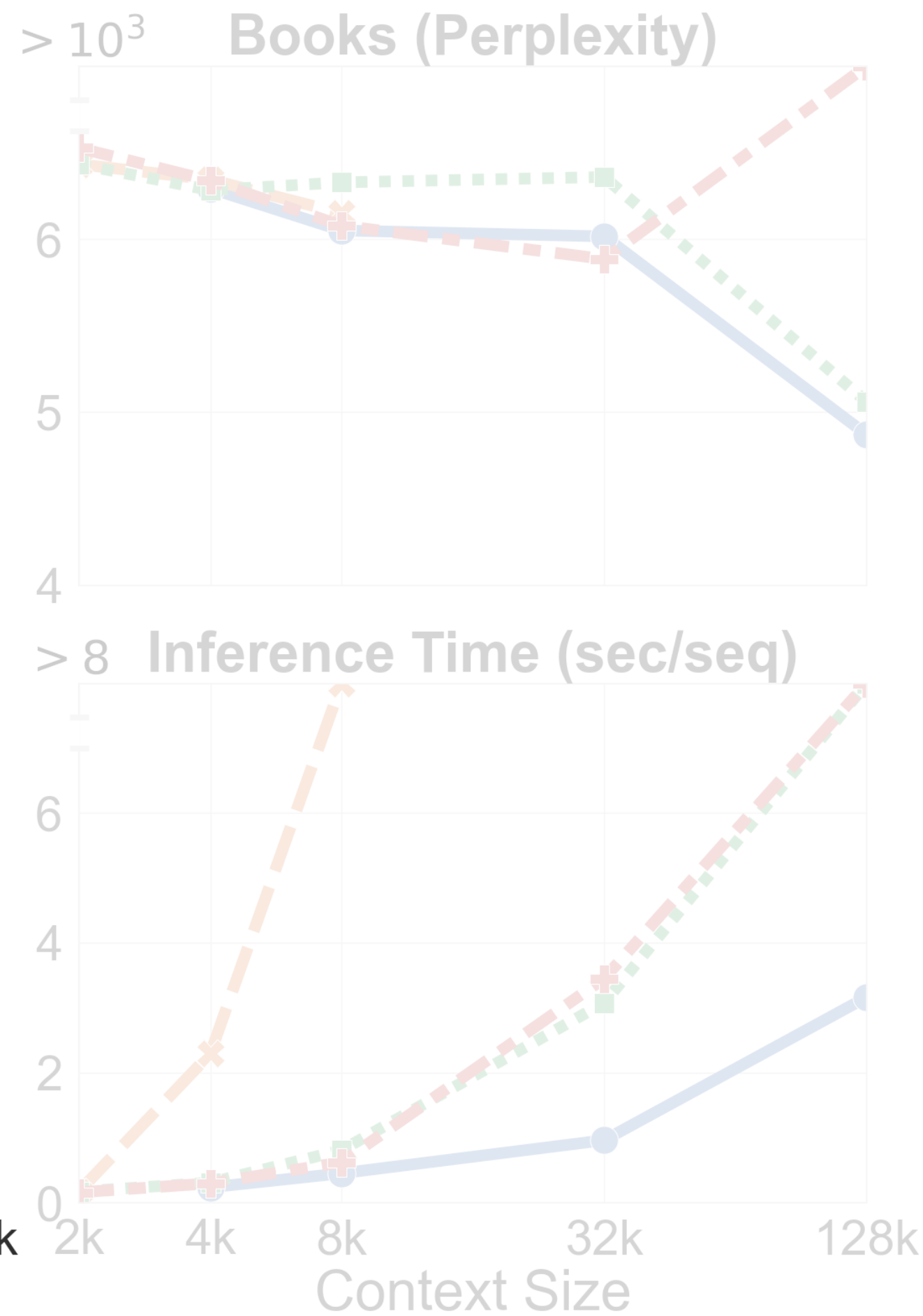
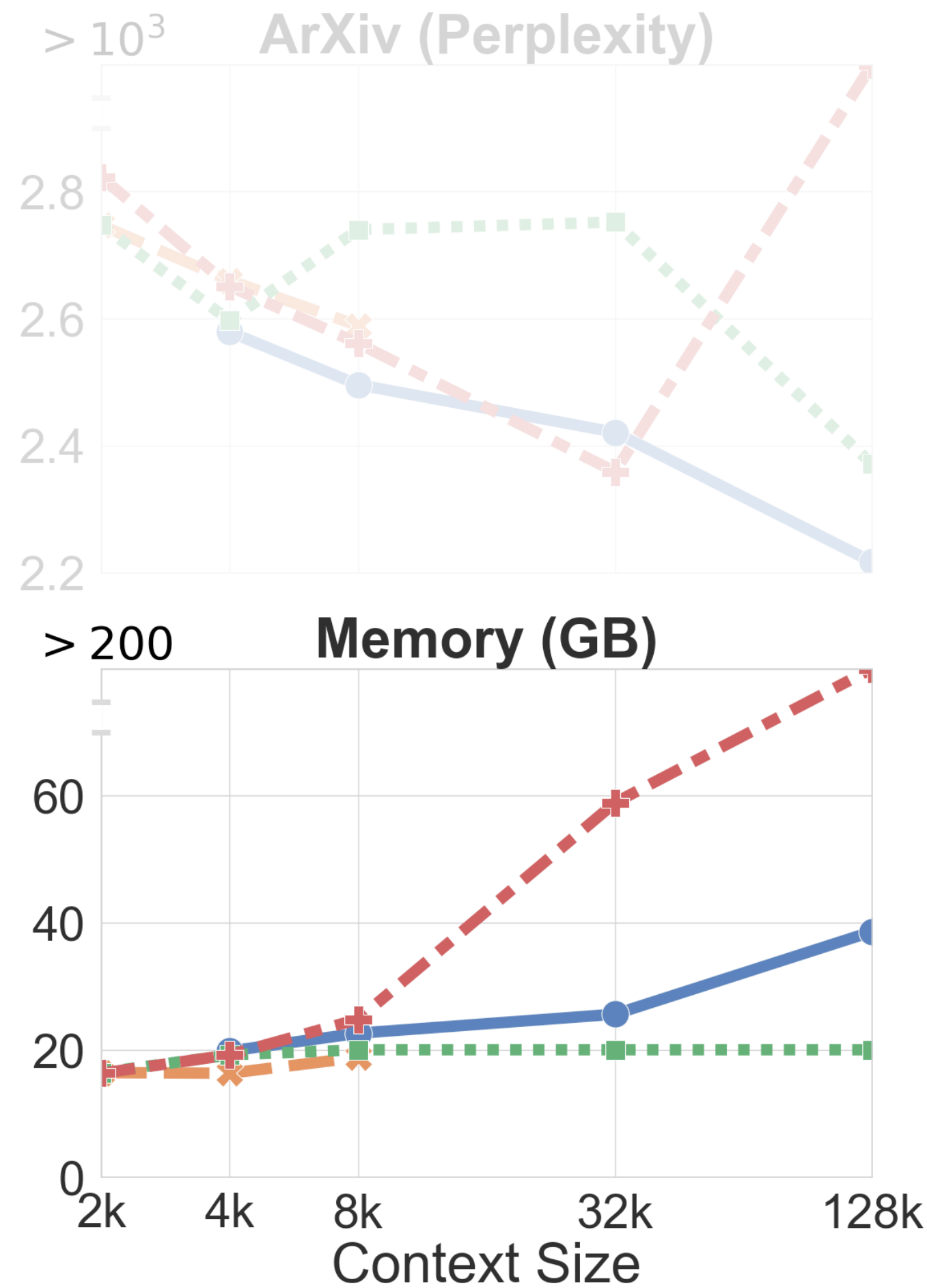
# Evaluation - long-context language modeling



## Performance

CEPE continues to improve perplexity with more context (only trained on 8K)

# Evaluation - long-context language modeling



## Performance

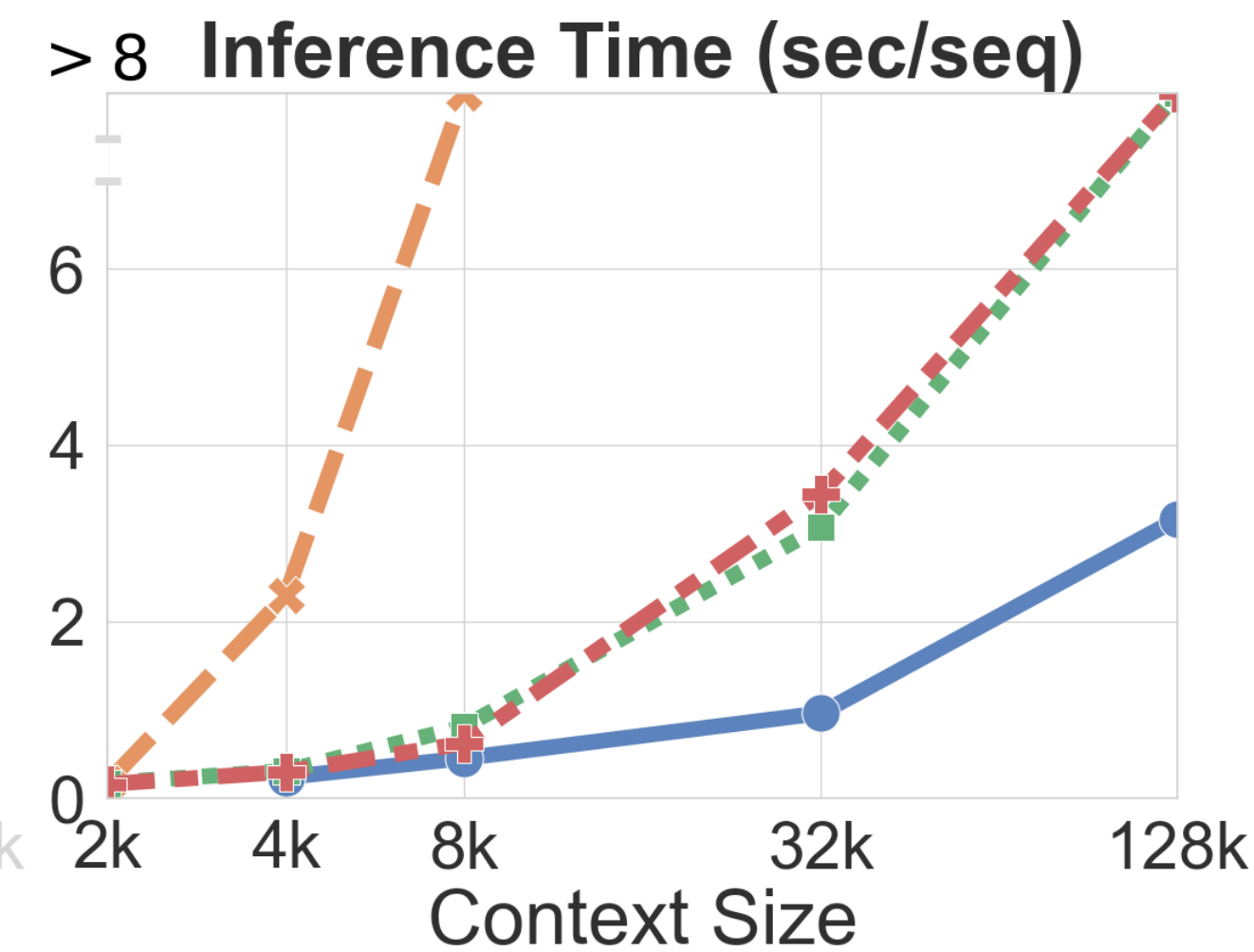
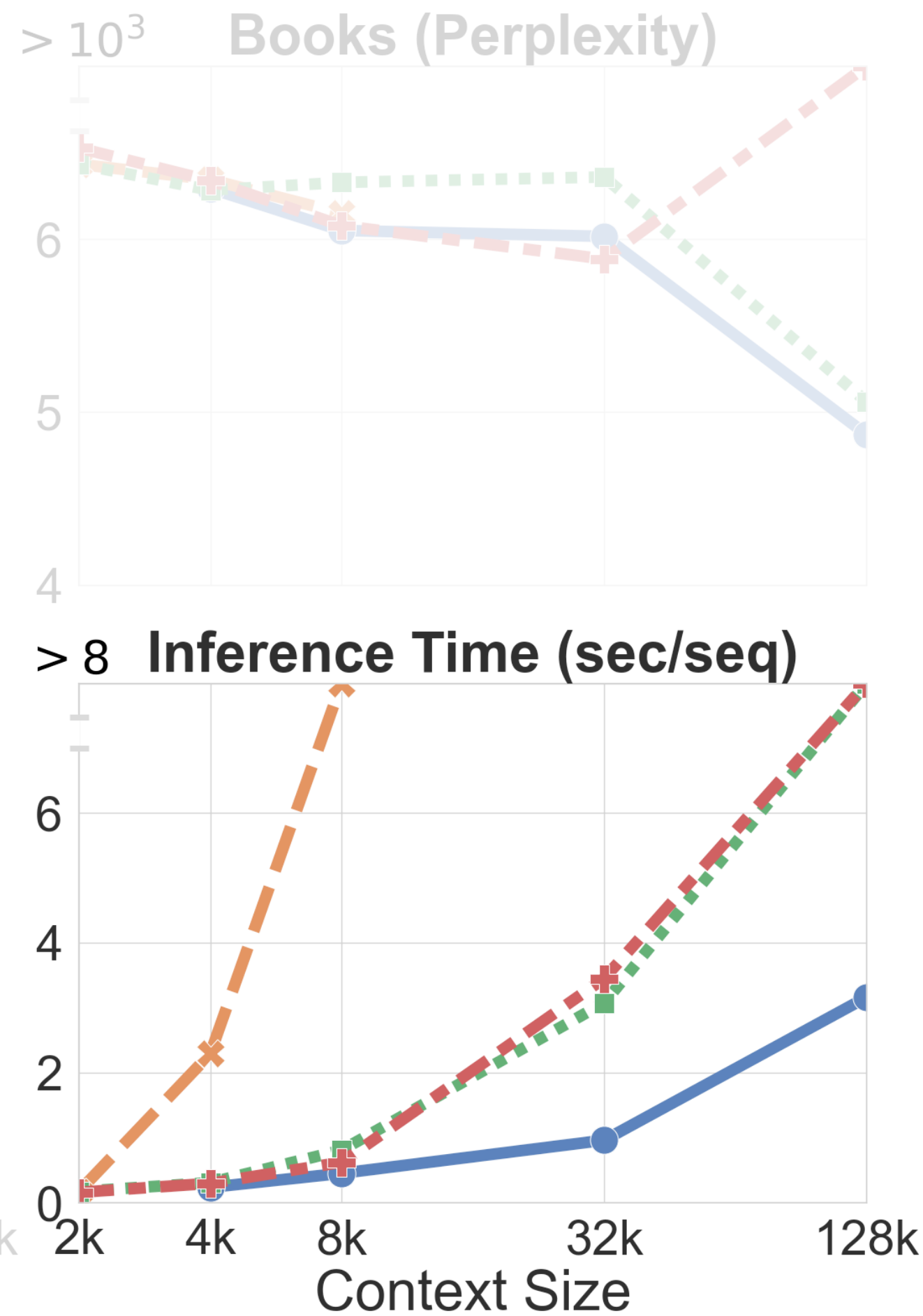
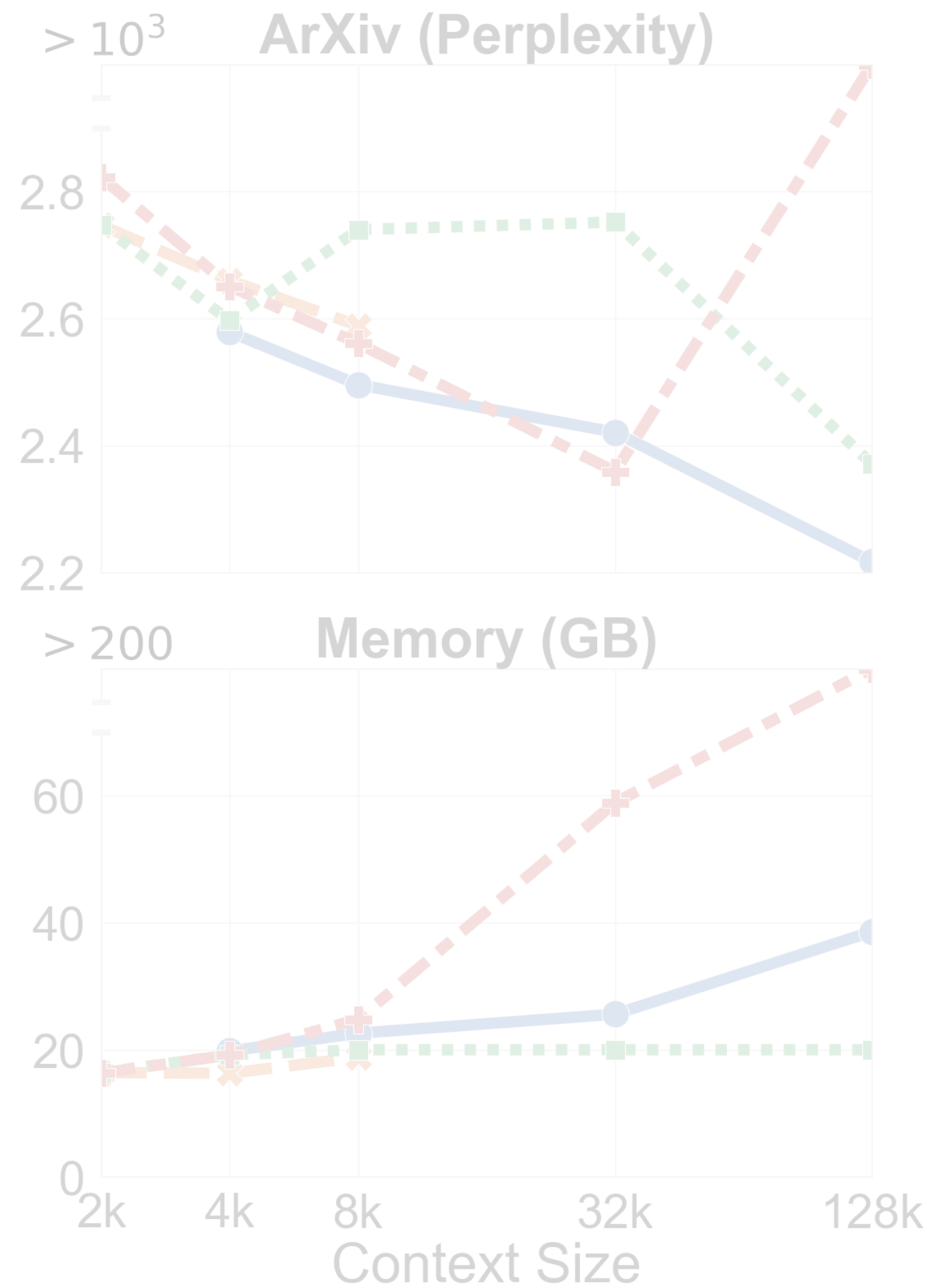
CEPE continues to improve perplexity with more context (only trained on 8K)

## Memory

CEPE maintains a low memory usage (1/6 of full attention)

-+ - YaRN   
 -■ - StreamingLLM   
 -\* - REPLUG   
 -● - CEPE

# Evaluation - long-context language modeling



## Performance

CEPE continues to improve perplexity with more context (only trained on 8K)

## Memory

CEPE maintains a low memory usage (1/6 of full attention)

## Throughput

CEPE achieves the highest throughput (10x of full attention)

---+--- YaRN    -.-■-.- StreamingLLM    -.-\*--- REPLUG    -●- CEPE



# Evaluation - retrieval-augmented applications

# Evaluation - retrieval-augmented applications

**Open-domain question answering**

# Evaluation - retrieval-augmented applications

## Open-domain question answering

Passage k: ...

Passage k-1: ...

...

Passage 1: ...

*Question: Who is the lead actor of “Dune: Part Two”*

# Evaluation - retrieval-augmented applications

## Open-domain question answering

Retrieved from Wikipedia based on the question

(using a dense retriever)

Passage k: ...

Passage k-1: ...

...

Passage 1: ...

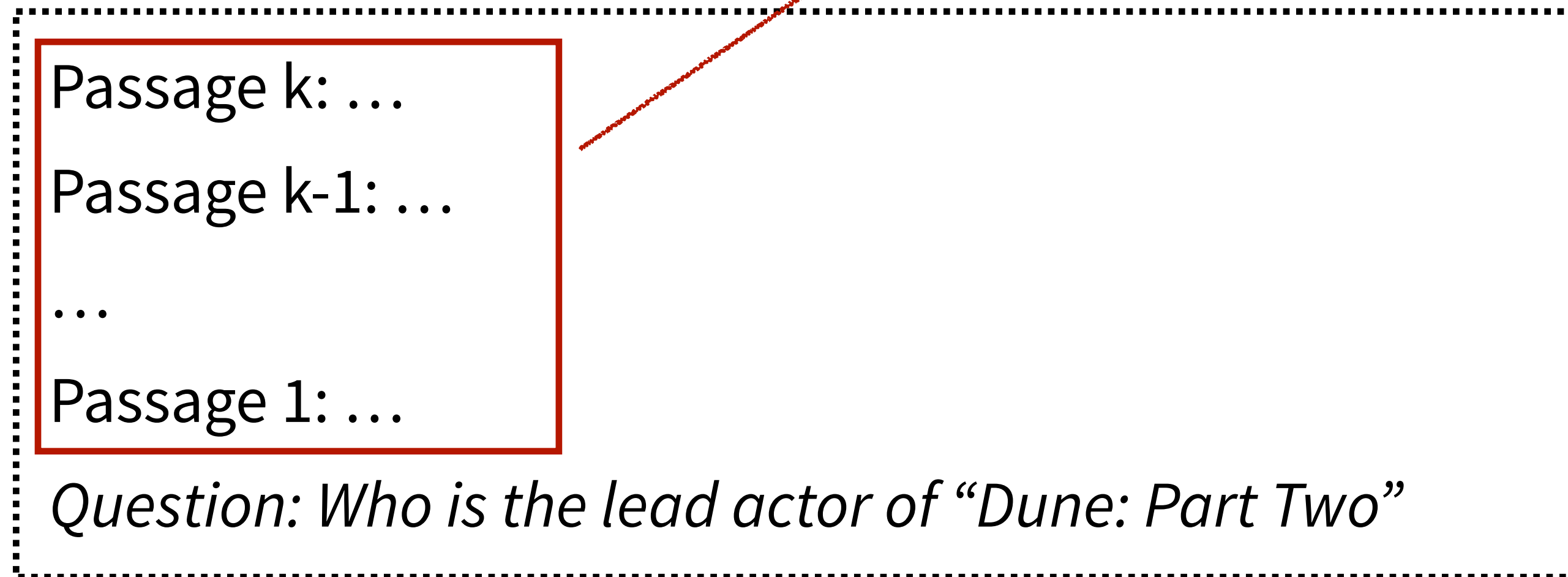
*Question: Who is the lead actor of "Dune: Part Two"*

# Evaluation - retrieval-augmented applications

## Open-domain question answering

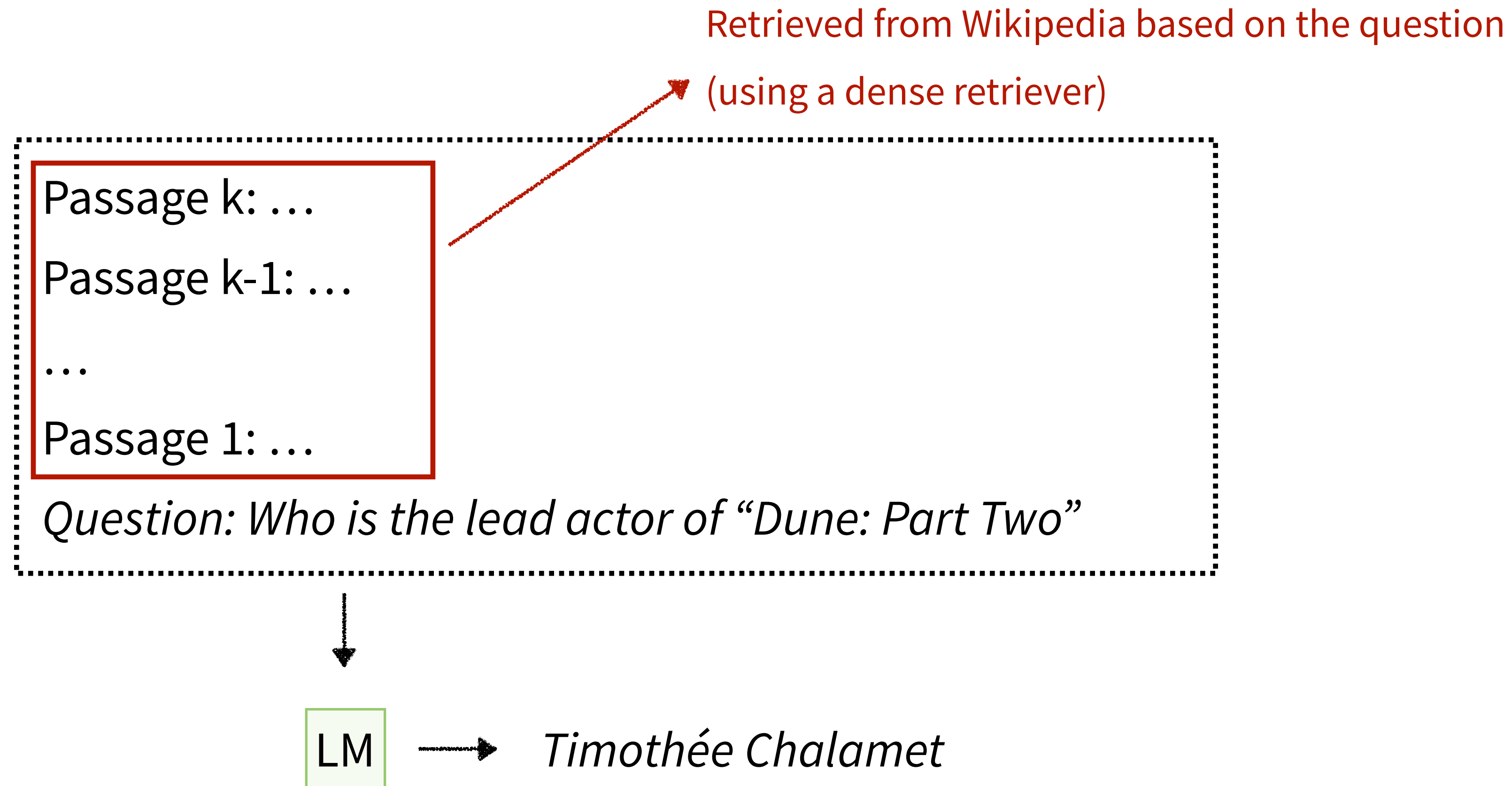
Retrieved from Wikipedia based on the question

(using a dense retriever)



# Evaluation - retrieval-augmented applications

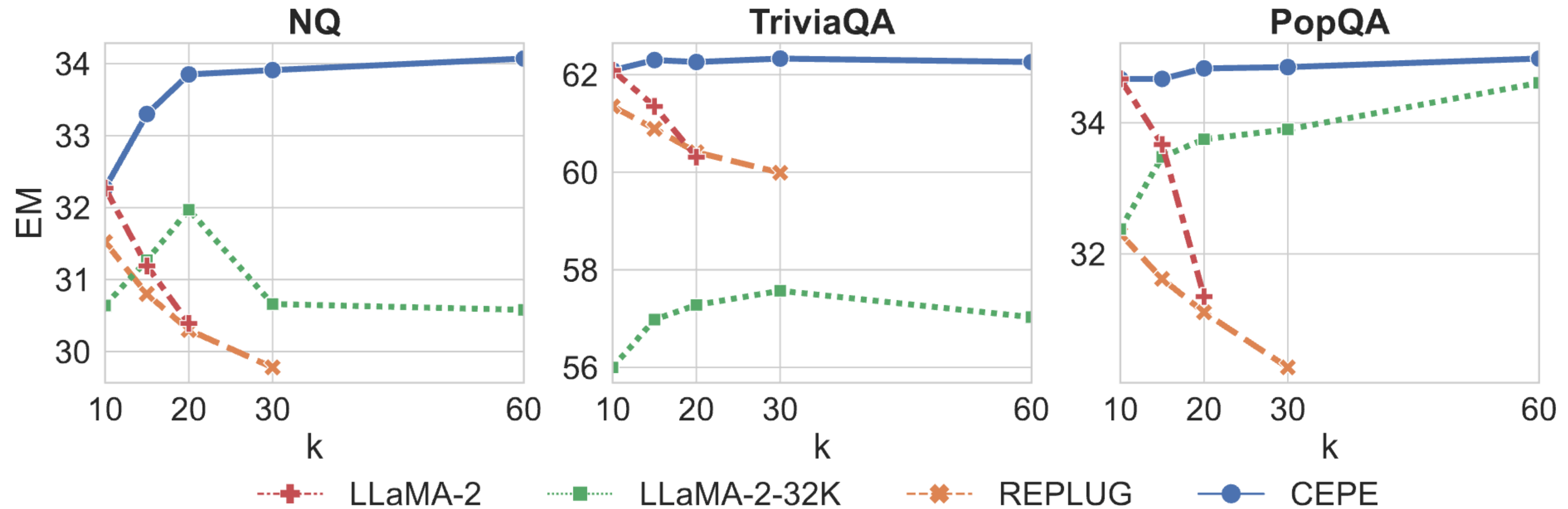
## Open-domain question answering



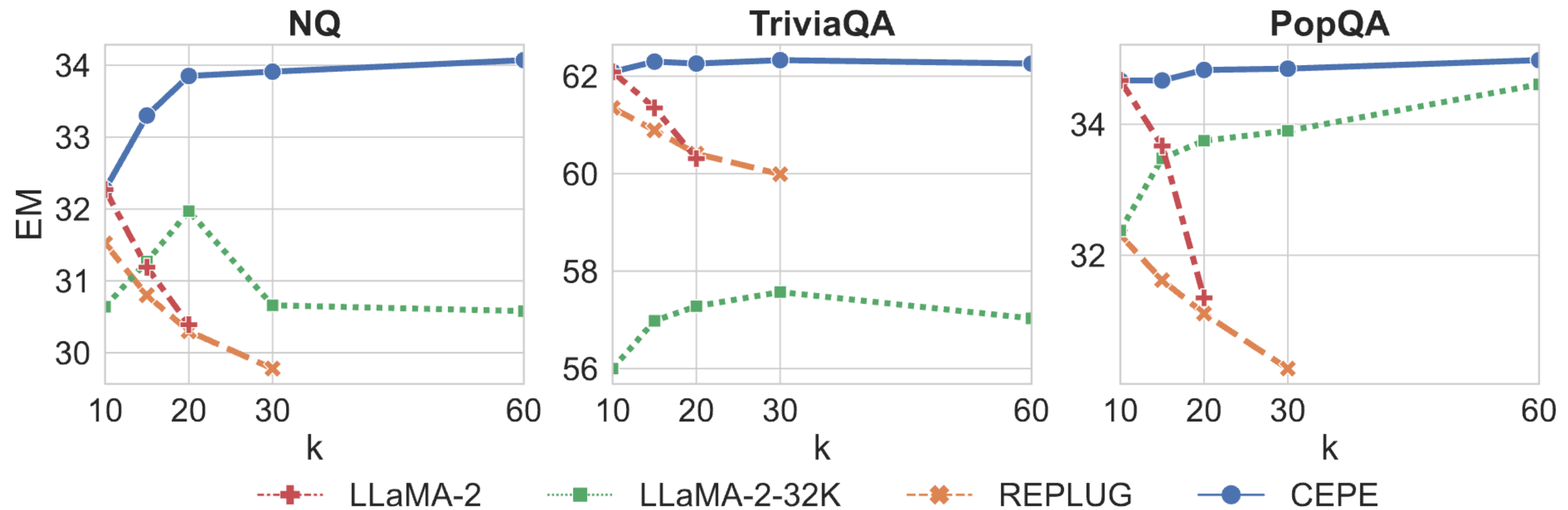


# Evaluation - retrieval-augmented applications

# Evaluation - retrieval-augmented applications



# Evaluation - retrieval-augmented applications



CEPE extrapolates well to more retrieved passages without getting distracted  
(also more efficient)

# Evaluation — in-context learning

Brown et al., 2020. Language Models are Few-Shot Learners.

Rong, 2021. Extrapolating to Unnatural Language Processing with GPT-3's In-context Learning: The Good, the Bad, and the Mysterious.

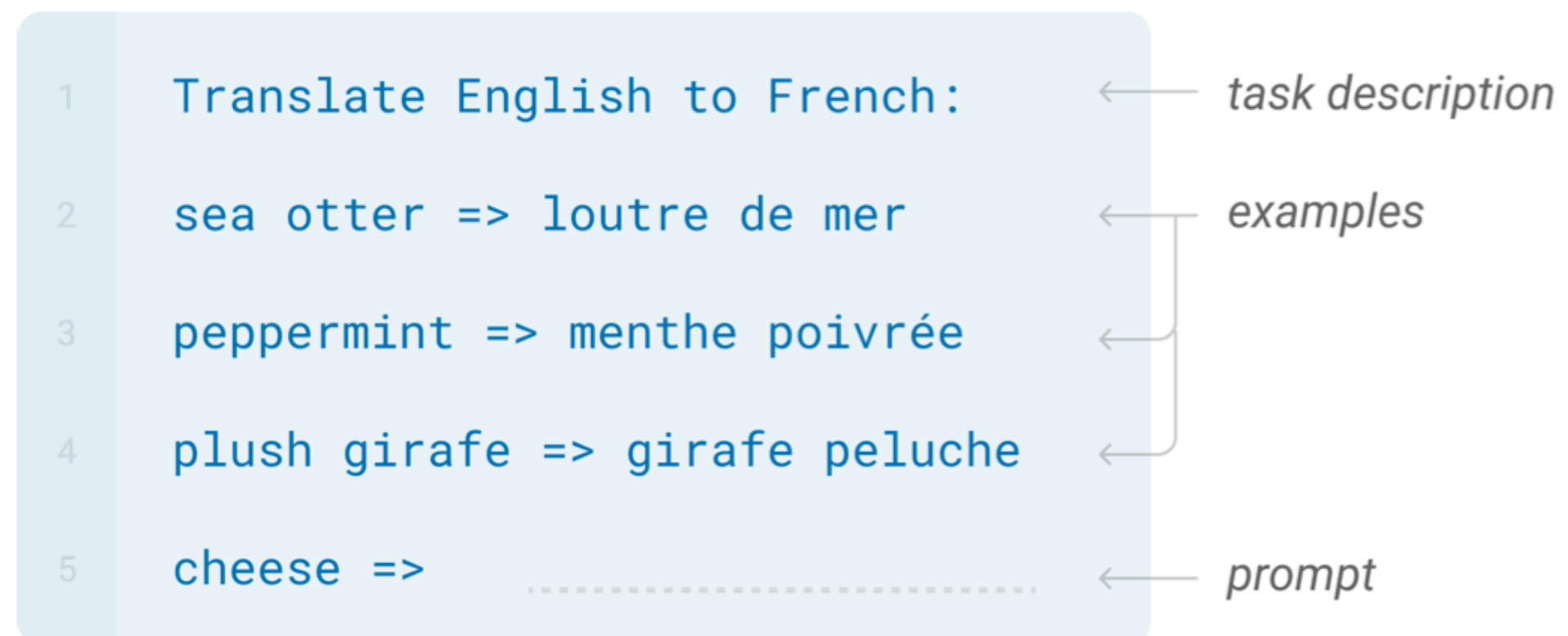
# Evaluation — in-context learning

An “emerging” ability of large language models: **in-context learning**

Brown et al., 2020. Language Models are Few-Shot Learners.

Rong, 2021. Extrapolating to Unnatural Language Processing with GPT-3's In-context Learning: The Good, the Bad, and the Mysterious.

# Evaluation — in-context learning



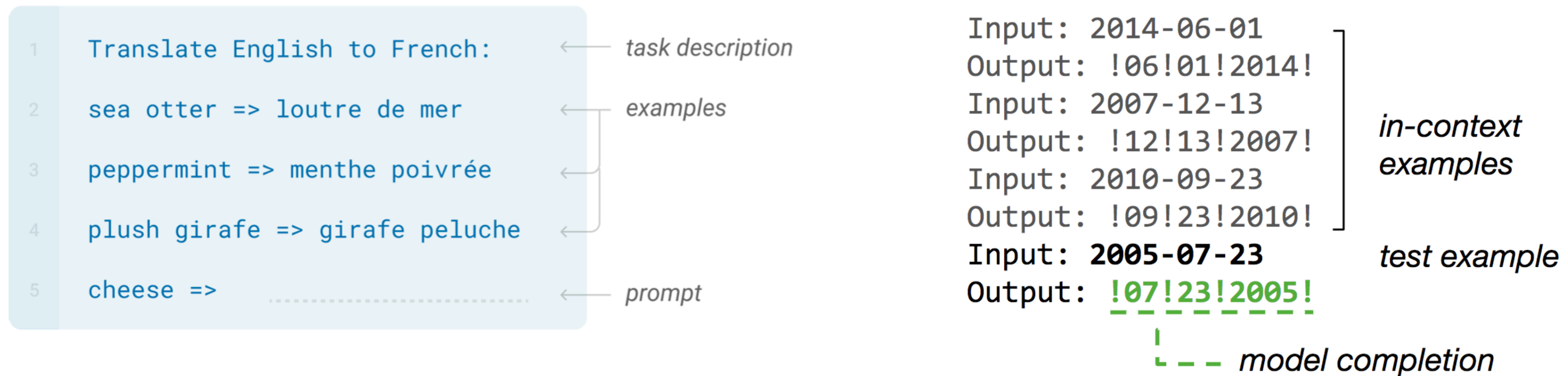
An “emerging” ability of large language models: **in-context learning**

Brown et al., 2020. Language Models are Few-Shot Learners.

Rong, 2021. Extrapolating to Unnatural Language Processing with GPT-3's In-context Learning: The Good, the Bad, and the Mysterious.

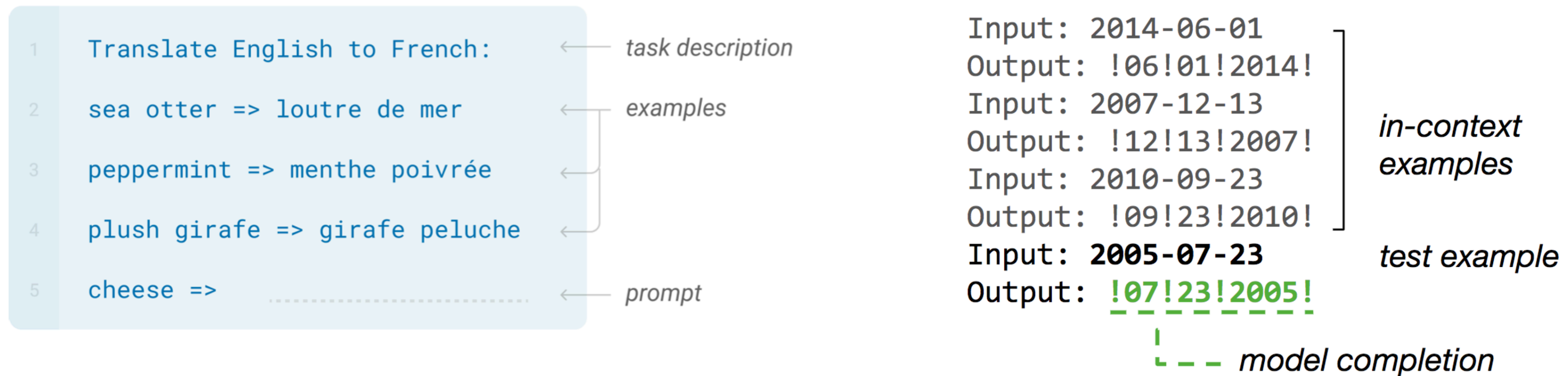


# Evaluation — in-context learning



An “emerging” ability of large language models: **in-context learning**

# Evaluation — in-context learning



An “emerging” ability of large language models: **in-context learning**

Is the *cross-attention* good enough for using in-context examples?

# Evaluation — in-context learning

	$k$	SST2	MR	AGNews	SST5	TREC	TREC-F	DBPedia	NLU-S	NLU-I	BANKING	CLINIC
LLAMA-2	2	89.1	96.7	72.7	3.9	<b>48.0</b>	16.7	<b>94.0</b>	42.3	22.3	38.4	59.1
+ CEPE	2 + 18	90.7	<b>98.4</b>	71.9	<b>46.7</b>	47.1	22.8	<b>94.0</b>	<b>48.9</b>	30.4	42.5	62.4
	2 + 38	<b>92.9</b>	98.0	<b>73.2</b>	45.5	47.5	<b>25.1</b>	93.3	48.8	<b>31.6</b>	<b>46.0</b>	<b>62.8</b>
LLAMA-2 <sup>†</sup>	40	94.3	98.7	74.7	52.3	87.7	54.8	95.1	76.7	62.1	50.4	72.0

# Evaluation — in-context learning

#examples in the decoder

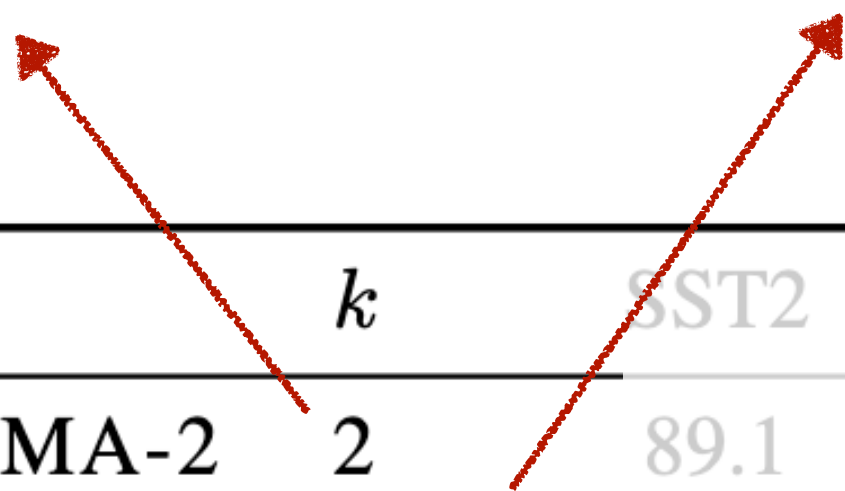


	$k$	SST2	MR	AGNews	SST5	TREC	TREC-F	DBPedia	NLU-S	NLU-I	BANKING	CLINIC
LLAMA-2	2	89.1	96.7	72.7	3.9	<b>48.0</b>	16.7	<b>94.0</b>	42.3	22.3	38.4	59.1
+ CEPE	2 + 18	90.7	<b>98.4</b>	71.9	<b>46.7</b>	47.1	22.8	<b>94.0</b>	<b>48.9</b>	30.4	42.5	62.4
	2 + 38	<b>92.9</b>	98.0	<b>73.2</b>	45.5	47.5	<b>25.1</b>	93.3	48.8	<b>31.6</b>	<b>46.0</b>	<b>62.8</b>
LLAMA-2 <sup>†</sup>	40	94.3	98.7	74.7	52.3	87.7	54.8	95.1	76.7	62.1	50.4	72.0

# Evaluation — in-context learning

#examples in the decoder

#examples in the encoder



	$k$	SST2	MR	AGNews	SST5	TREC	TREC-F	DBPedia	NLU-S	NLU-I	BANKING	CLINIC
LLAMA-2	2	89.1	96.7	72.7	3.9	<b>48.0</b>	16.7	<b>94.0</b>	42.3	22.3	38.4	59.1
+ CEPE	2 + 18	90.7	<b>98.4</b>	71.9	<b>46.7</b>	47.1	22.8	<b>94.0</b>	<b>48.9</b>	30.4	42.5	62.4
	2 + 38	<b>92.9</b>	98.0	<b>73.2</b>	45.5	47.5	<b>25.1</b>	93.3	48.8	<b>31.6</b>	<b>46.0</b>	<b>62.8</b>
LLAMA-2 <sup>†</sup>	40	94.3	98.7	74.7	52.3	87.7	54.8	95.1	76.7	62.1	50.4	72.0

# Evaluation — in-context learning

	$k$	SST2	MR	AGNews	SST5	TREC	TREC-F	DBPedia	NLU-S	NLU-I	BANKING	CLINIC
LLAMA-2	2	89.1	96.7	72.7	3.9	<b>48.0</b>	16.7	<b>94.0</b>	42.3	22.3	38.4	59.1
+ CEPE	2 + 18	90.7	<b>98.4</b>	71.9	<b>46.7</b>	47.1	22.8	<b>94.0</b>	<b>48.9</b>	30.4	42.5	62.4
	2 + 38	<b>92.9</b>	98.0	<b>73.2</b>	45.5	47.5	<b>25.1</b>	93.3	48.8	<b>31.6</b>	<b>46.0</b>	<b>62.8</b>
LLAMA-2 <sup>†</sup>	40	94.3	98.7	74.7	52.3	87.7	54.8	95.1	76.7	62.1	50.4	72.0

CEPE can perform in-context learning using demonstrations in the encoder



# Evaluation — in-context learning

	$k$	SST2	MR	AGNews	SST5	TREC	TREC-F	DBPedia	NLU-S	NLU-I	BANKING	CLINIC
LLAMA-2	2	89.1	96.7	72.7	3.9	<b>48.0</b>	16.7	<b>94.0</b>	42.3	22.3	38.4	59.1
+ CEPE	2 + 18	90.7	<b>98.4</b>	71.9	<b>46.7</b>	47.1	22.8	<b>94.0</b>	<b>48.9</b>	30.4	42.5	62.4
	2 + 38	<b>92.9</b>	98.0	<b>73.2</b>	45.5	47.5	<b>25.1</b>	93.3	48.8	<b>31.6</b>	<b>46.0</b>	<b>62.8</b>
LLAMA-2 <sup>†</sup>	40	94.3	98.7	74.7	52.3	87.7	54.8	95.1	76.7	62.1	50.4	72.0

CEPE can perform in-context learning using demonstrations in the encoder

... though the performance still lags behind putting demonstrations in the decoder

# Evaluation – chat model evaluation

	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3

# Evaluation – chat model evaluation

Long context (books, government report, papers)



	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3

# Evaluation – chat model evaluation

Long context (books, government report, papers)

Zero-shot (no training; no in-context examples)



	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3

# Evaluation – chat model evaluation

Long context (books, government report, papers)

Zero-shot (no training; no in-context examples)

The model needs to understand “instructions”



	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3

# Evaluation – chat model evaluation

	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3



# Evaluation – chat model evaluation

Tokens in the decoder



	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3

# Evaluation – chat model evaluation

Tokens in the decoder      Tokens in the encoder

	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3

# Evaluation – chat model evaluation

	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3

CEPE can utilize the long context and boost the QA/summarization performance

# Evaluation – chat model evaluation

	Total tokens	Question Answering			Summarization		
		NQA	Qspr	QALT	GvRp	SSFD	QMSum
LLAMA-2-CHAT	2K	17.1	14.6	28.6	16.0	16.4	19.3
+ CEPED	2K + 2K	19.5	<b>20.5</b>	<b>30.2</b>	<b>16.5</b>	16.4	<b>19.6</b>
	2K + 30K	21.6	19.9	29.6	15.8	<b>16.7</b>	19.5
	2K + All	<b>21.9</b>	19.9	29.6	15.9	<b>16.7</b>	19.5
LLAMA-2-32K INSTRUCT	32K	12.2	18.1	41.6	19.9	10.0	10.3

Compared to a full-attention decoder, CEPE’s performance is more stable across different tasks

# Conclusion

# Conclusion

- CEPE is a lightweight framework for extending the context window of any decoder-only LMs



# Conclusion

- CEPE is a lightweight framework for extending the context window of any decoder-only LMs
- CEPE is **length-generalizable**, **efficient**, and **easy to train**

# Conclusion

- CEPE is a lightweight framework for extending the context window of any decoder-only LMs
- CEPE is **length-generalizable**, **efficient**, and **easy to train**
- CEPE performs well on both **long-context modeling** and **retrieval-augmented** applications

# Conclusion

- CEPE is a lightweight framework for extending the context window of any decoder-only LMs
- CEPE is **length-generalizable**, **efficient**, and **easy to train**
- CEPE performs well on both **long-context modeling** and **retrieval-augmented** applications
- CEPE works well with instruction-tuned/chat models too (with only unlabeled data)

**What's next**

# What's next

# What's next

- Long-context instruction-tuning



# What's next

- Long-context instruction-tuning
- Reduce training cost

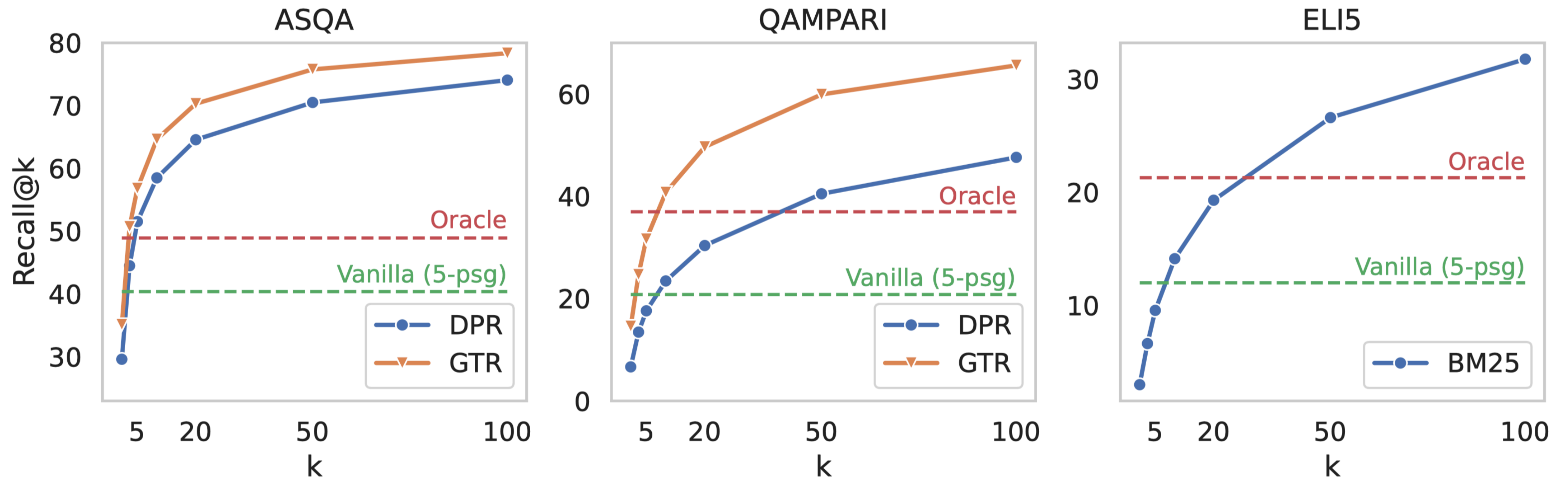
# What's next

- Long-context instruction-tuning
- Reduce training cost
- Reduce inference cost

# What's next

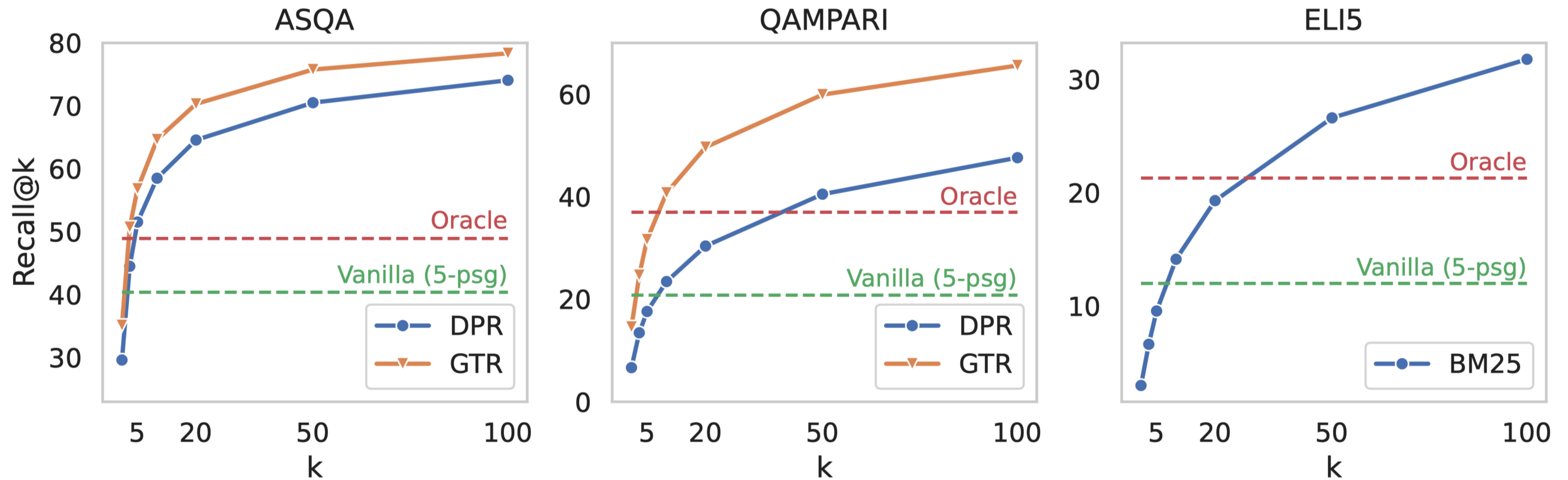
- Long-context instruction-tuning
- Reduce training cost
- Reduce inference cost
- **Improve retrieval-augmented applications**

# ... Back to ALCE



More passages encode more information to answer the question.

# ... Back to ALCE



More passages encode more information to answer the question.

Can LLMs use them effectively?

# ... Back to ALCE

	<u>Fluency</u>	<u>Correct.</u>	<u>Citation</u>	
	(MAUVE)	(EM Rec.)	Rec.	Prec.
<b>ChatGPT-16K (max #tokens=16K)</b>				
ChatGPT (5-psg)	60.3	36.1	76.2	76.5
ChatGPT (10-psg)	56.3	36.7	75.3	75.0
ChatGPT (20-psg)	56.7	36.1	73.7	73.5
<b>GPT-4 (max #tokens=8K)</b>				
GPT-4 (5-psg)	67.1	41.3	68.5	75.6
GPT-4 (10-psg)	71.5	43.1	72.0	75.5
GPT-4 (20-psg)	64.9	44.4	73.0	76.5



# ... Back to ALCE

	<b>Fluency</b>	<b>Correct.</b>	<b>Citation</b>	
	(MAUVE)	(EM Rec.)	Rec.	Prec.
<b>ChatGPT-16K (max #tokens=16K)</b>				
ChatGPT (5-psg)	60.3	36.1	76.2	76.5
ChatGPT (10-psg)	56.3	36.7	75.3	75.0
ChatGPT (20-psg)	56.7	36.1	73.7	73.5
<b>GPT-4 (max #tokens=8K)</b>				
GPT-4 (5-psg)	67.1	41.3	68.5	75.6
GPT-4 (10-psg)	71.5	43.1	72.0	75.5
GPT-4 (20-psg)	64.9	44.4	73.0	76.5

- More passages do not lead to better performance

# ... Back to ALCE

	<b>Fluency</b>	<b>Correct.</b>	<b>Citation</b>	
	(MAUVE)	(EM Rec.)	Rec.	Prec.
<b>ChatGPT-16K (max #tokens=16K)</b>				
ChatGPT (5-psg)	60.3	36.1	76.2	76.5
ChatGPT (10-psg)	56.3	36.7	75.3	75.0
ChatGPT (20-psg)	56.7	36.1	73.7	73.5
<b>GPT-4 (max #tokens=8K)</b>				
GPT-4 (5-psg)	67.1	41.3	68.5	75.6
GPT-4 (10-psg)	71.5	43.1	72.0	75.5
GPT-4 (20-psg)	64.9	44.4	73.0	76.5

- More passages do not lead to better performance
- A stronger model utilizes the more information better

# ... Back to ALCE

	<b>Fluency</b>	<b>Correct.</b>	<b>Citation</b>	
	(MAUVE)	(EM Rec.)	Rec.	Prec.
<b>ChatGPT-16K (max #tokens=16K)</b>				
ChatGPT (5-psg)	60.3	36.1	76.2	76.5
ChatGPT (10-psg)	56.3	36.7	75.3	75.0
ChatGPT (20-psg)	56.7	36.1	73.7	73.5
<b>GPT-4 (max #tokens=8K)</b>				
GPT-4 (5-psg)	67.1	41.3	68.5	75.6
GPT-4 (10-psg)	71.5	43.1	72.0	75.5
GPT-4 (20-psg)	64.9	44.4	73.0	76.5
Oracle (20-psg)		>70		

- More passages do not lead to better performance
- A stronger model utilizes the more information better

# ... Back to ALCE

	<b>Fluency</b>	<b>Correct.</b>	<b>Citation</b>	
	(MAUVE)	(EM Rec.)	Rec.	Prec.
<b>ChatGPT-16K (max #tokens=16K)</b>				
ChatGPT (5-psg)	60.3	36.1	76.2	76.5
ChatGPT (10-psg)	56.3	36.7	75.3	75.0
ChatGPT (20-psg)	56.7	36.1	73.7	73.5
<b>GPT-4 (max #tokens=8K)</b>				
GPT-4 (5-psg)	67.1	41.3	68.5	75.6
GPT-4 (10-psg)	71.5	43.1	72.0	75.5
GPT-4 (20-psg)	64.9	44.4	73.0	76.5
Oracle (20-psg)		>70		

- More passages do not lead to better performance
- A stronger model utilizes the more information better
- LLMs are not good at synthesizing information from long context

# ... Back to ALCE

	Fluency (MAUVE)	Correct. (EM Rec.)	Citation Rec. Prec.	
GPT-4 (10-psg)	71.5	43.1	72.0	75.5
GPT-4 (20-psg)	64.9	44.4	73.0	76.5
Oracle (20-psg)		>70		

More passages do not lead to better performance

Improve LLMs' ability to **retrieve and synthesize** multiple pieces of information from long-context

from long context



Contact: [tianyug@princeton.edu](mailto:tianyug@princeton.edu)

 @gaotianyu1350

 @princeton\_nlp

 @PrincetonPLI