CS11-711 Advanced NLP

# Retrieval and RAG

Akari Asai

https://cmu-l3.github.io/anlp-spring2026/

# Limitations of LMs: Hallucinations

(GPT 5.1, 2026/01/20)

# Limitations of LMs: Hallucinations

(GPT 5.1, 2026/01/20)

Who is Akari Asai?

# Limitations of LMs: Hallucinations



(GPT 5.1, 2026/01/20)

Who is Akari Asai?

**Akari Asai** is a prominent researcher in the field of **Natural Language Processing (NLP)** and **Artificial Intelligence (AI)**. She is particularly known for her work on **open-domain question answering**, **retrieval-augmented generation**, and **large language models**. Her research focuses on enabling machines to retrieve, understand, and reason over large corpora of unstructured text — a critical area for building systems like ChatGPT and other advanced NLP applications.

# Limitations of LMs: Hallucinations



(GPT 5.1, 2026/01/20)

**Who is Akari Asai?**

**Akari Asai** is a prominent researcher in the field of **Natural Language Processing (NLP)** and **Artificial Intelligence (AI)**. She is particularly known for her work on **open-domain question answering**, **retrieval-augmented generation**, and **large language models**. Her research focuses on enabling machines to retrieve, understand, and reason over large corpora of unstructured text — a critical area for building systems like ChatGPT and other advanced NLP applications.

📄 **Representative Papers by Akari Asai**

Here are **five representative papers** by Akari Asai, with their **publication venues**:

1. **"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"**
   *Venue:* Proceedings of NeurIPS 2020 (Conference on Neural Information Processing Systems)

2. **"Learning to Retrieve Passages without Supervision"**
   *Venue:* Proceedings of ACL 2021 (Annual Meeting of the Association for Computational Linguistics)

3. **"Task-Aware Query Reformulation with Reinforcement Learning for Open-Domain Question Answering"**
   *Venue:* Proceedings of EMNLP 2019 (Conference on Empirical Methods in Natural Language Processing)

# Limitations of LMs: Hallucinations

(GPT 5.1, 2026/01/20)

Who is Akari Asai?

**Akari Asai** is a prominent researcher in the field of **Natural Language Processing (NLP)** and **Artificial Intelligence (AI)**. She is particularly known for her work on **open-domain question answering**, **retrieval-augmented generation**, and **large language models**. Her research focuses on enabling machines to retrieve, understand, and reason over large corpora of unstructured text — a critical area for building systems like ChatGPT and other advanced NLP applications.

📄 **Representative Papers by Akari Asai**

Here are **five representative papers** by Akari Asai, with their **publication venues**:

1. ~~**"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"**~~

    *Venue:* Proceedings of NeurIPS 2020 (Conference on Neural Information Processing Systems)

2. **"Learning to Retrieve Passages without Supervision"**

    *Venue:* Proceedings of ACL 2021 (Annual Meeting of the Association for Computational Linguistics)

3. **"Task-Aware Query Reformulation with Reinforcement Learning for Open-Domain Question Answering"**

    *Venue:* Proceedings of EMNLP 2019 (Conference on Empirical Methods in Natural Language Processing)

# Limitations of LMs: Hallucinations

(GPT 5.1, 2026/01/20)

Who is Akari Asai?

**Akari Asai** is a prominent researcher in the field of **Natural Language Processing (NLP)** and **Artificial Intelligence (AI)**. She is particularly known for her work on **open-domain question answering**, **retrieval-augmented generation**, and **large language models**. Her research focuses on enabling machines to retrieve, understand, and reason over large corpora of unstructured text — a critical area for building systems like ChatGPT and other advanced NLP applications.

📄 **Representative Papers by Akari Asai**

Here are **five representative papers** by Akari Asai, with their **publication venues**:

1. ~~"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"~~        **Paper exists but not by me**

   *Venue:* Proceedings of NeurIPS 2020 (Conference on Neural Information Processing Systems)

2. **"Learning to Retrieve Passages without Supervision"**

   *Venue:* Proceedings of ACL 2021 (Annual Meeting of the Association for Computational Linguistics)

3. **"Task-Aware Query Reformulation with Reinforcement Learning for Open-Domain Question Answering"**

   *Venue:* Proceedings of EMNLP 2019 (Conference on Empirical Methods in Natural Language Processing)

# Limitations of LMs: Hallucinations

(GPT 5.1, 2026/01/20)

Who is Akari Asai?

**Akari Asai** is a prominent researcher in the field of **Natural Language Processing (NLP)** and **Artificial Intelligence (AI)**. She is particularly known for her work on **open-domain question answering**, **retrieval-augmented generation**, and **large language models**. Her research focuses on enabling machines to retrieve, understand, and reason over large corpora of unstructured text — a critical area for building systems like ChatGPT and other advanced NLP applications.

## 📄 Representative Papers by Akari Asai

Here are **five representative papers** by Akari Asai, with their **publication venues**:

1. ~~"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"~~

   *Venue:* Proceedings of NeurIPS 2020 (Conference on Neural Information Processing Systems)

2. ~~"Learning to Retrieve Passages without Supervision"~~

   *Venue:* Proceedings of ACL 2021 (Annual Meeting of the Association for Computational Linguistics)

3. **"Task-Aware Query Reformulation with Reinforcement Learning for Open-Domain Question Answering"**

   *Venue:* Proceedings of EMNLP 2019 (Conference on Empirical Methods in Natural Language Processing)

Paper exists but not by me

# Limitations of LMs: Hallucinations

(GPT 5.1, 2026/01/20)

Who is Akari Asai?

**Akari Asai** is a prominent researcher in the field of **Natural Language Processing (NLP)** and **Artificial Intelligence (AI)**. She is particularly known for her work on **open-domain question answering**, **retrieval-augmented generation**, and **large language models**. Her research focuses on enabling machines to retrieve, understand, and reason over large corpora of unstructured text — a critical area for building systems like ChatGPT and other advanced NLP applications.

## 📄 Representative Papers by Akari Asai

Here are **five representative papers** by Akari Asai, with their **publication venues**:

1. ~~"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"~~
   *Venue:* Proceedings of NeurIPS 2020 (Conference on Neural Information Processing Systems)

   `Paper exists but not by me`

2. ~~"Learning to Retrieve Passages without Supervision"~~
   *Venue:* Proceedings of ACL 2021 (Annual Meeting of the Association for Computational Linguistics)

   `Paper exists but not by me`

3. **"Task-Aware Query Reformulation with Reinforcement Learning for Open-Domain Question Answering"**
   *Venue:* Proceedings of EMNLP 2019 (Conference on Empirical Methods in Natural Language Processing)

2

# Limitations of LMs: Hallucinations

(GPT 5.1, 2026/01/20)

Who is Akari Asai?

**Akari Asai** is a prominent researcher in the field of **Natural Language Processing (NLP)** and **Artificial Intelligence (AI)**. She is particularly known for her work on **open-domain question answering**, **retrieval-augmented generation**, and **large language models**. Her research focuses on enabling machines to retrieve, understand, and reason over large corpora of unstructured text — a critical area for building systems like ChatGPT and other advanced NLP applications.

## 📄 Representative Papers by Akari Asai

Here are **five representative papers** by Akari Asai, with their **publication venues**:

1. ~~**"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"**~~

   *Venue:* Proceedings of NeurIPS 2020 (Conference on Neural Information Processing Systems)

   ==Paper exists but not by me==

2. ~~**"Learning to Retrieve Passages without Supervision"**~~

   *Venue:* Proceedings of ACL 2021 (Annual Meeting of the Association for Computational Linguistics)

   ==Paper exists but not by me==

3. ~~**"Task-Aware Query Reformulation with Reinforcement Learning for Open-Domain Question Answering"**~~

   *Venue:* Proceedings of EMNLP 2019 (Conference on Empirical Methods in Natural Language Processing)

# Limitations of LMs: Hallucinations

(GPT 5.1, 2026/01/20)

Who is Akari Asai?

**Akari Asai** is a prominent researcher in the field of **Natural Language Processing (NLP)** and **Artificial Intelligence (AI)**. She is particularly known for her work on **open-domain question answering**, **retrieval-augmented generation**, and **large language models**. Her research focuses on enabling machines to retrieve, understand, and reason over large corpora of unstructured text — a critical area for building systems like ChatGPT and other advanced NLP applications.

📄 **Representative Papers by Akari Asai**

Here are **five representative papers** by Akari Asai, with their **publication venues**:

1. ~~"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"~~

   *Venue:* Proceedings of NeurIPS 2020 (Conference on Neural Information Processing Systems)

   > Paper exists but not by me

2. ~~"Learning to Retrieve Passages without Supervision"~~

   *Venue:* Proceedings of ACL 2021 (Annual Meeting of the Association for Computational Linguistics)

   > Paper exists but not by me

3. ~~"Task-Aware Query Reformulation with Reinforcement Learning for Open-Domain Question Answering"~~

   *Venue:* Proceedings of EMNLP 2019 (Conference on Empirical Methods in Natural Language Processing)

   > Doesn't exist

2

# Limitations of LMs: Hallucinations

| Title | Average Review Rating | Paper Link | Citation Check Scan Link | Example of Verified Hallucination | Comment |
|---|---|---|---|---|---|
| TamperTok: Forensics-Driven Tokenized Autoregressive Framework for Image Tampering Localization | 8.0 | TamperTok: Forensics-Driven Tokenized Autoregressive Framework for Image Tampering Localization \| OpenReview | https://app.gptzero.me/documents/4645494f-70eb-40bb-aea7-0007e13f7179/share | Chong Zou, Zhipeng Wang, Ziyu Li, Nan Wu, Yuling Cai, Shan Shi, Jiawei Wei, Xia Sun, Jian Wang, and Yizhou Wang. Segment everything everywhere all at once. In Advances in Neural Information Processing Systems (NeurIPS), volume 36, 2023. | This paper exists, but all authors are wrong. |
| MixtureVitae: Open Web-Scale Pretraining Dataset With High Quality Instruction and Reasoning Data Built from Permissive Text Sources | 8.0 | MixtureVitae: Open Web-Scale Pretraining Dataset With High Quality Instruction and Reasoning Data Built from Permissive Text Sources \| OpenReview | https://app.gptzero.me/documents/bfd10666-ea2d-454c-9ab2-75faa8b84281/share | Dan Hendrycks, Collin Burns, Steven Basart, Andy Critch, Jerry Li, Dawn Ippolito, Aina Lapedriza, Florian Tramer, Rylan Macfarlane, Eric Jiang, et al. Measuring massive multitask language understanding. In Proceedings of the International Conference on Learning Representations (ICLR), 2021. | The paper and first 3 authors match. The last 7 authors are not on the paper, and some of them do not exist |
| Catch-Only-One: Non-Transferable Examples for Model-Specific Authorization | 6.0 | Catch-Only-One: Non-Transferable Examples for Model-Specific Authorization \| OpenReview | https://app.gptzero.me/documents/9afb1d51-c5c8-48f2-9b75-250d95062521/share | Dinghuai Zhang, Yang Song, Inderjit Dhillon, and Eric Xing. Defense against adversarial attacks using spectral regularization. In International Conference on Learning Representations (ICLR), 2020. | No Match |

# Limitations of LMs:
# Obsolete knowledge

When did GPT 5.2 come out?

(GPT 5.2, 2026/01/20)

There's **no confirmed public release date for "GPT-5.2"** in the information I was trained on.

As of mid-2024, OpenAI hadn't officially announced a version named GPT-5.2, so I don't have a specific release date for it. If it was released after that, the exact date isn't in my training.
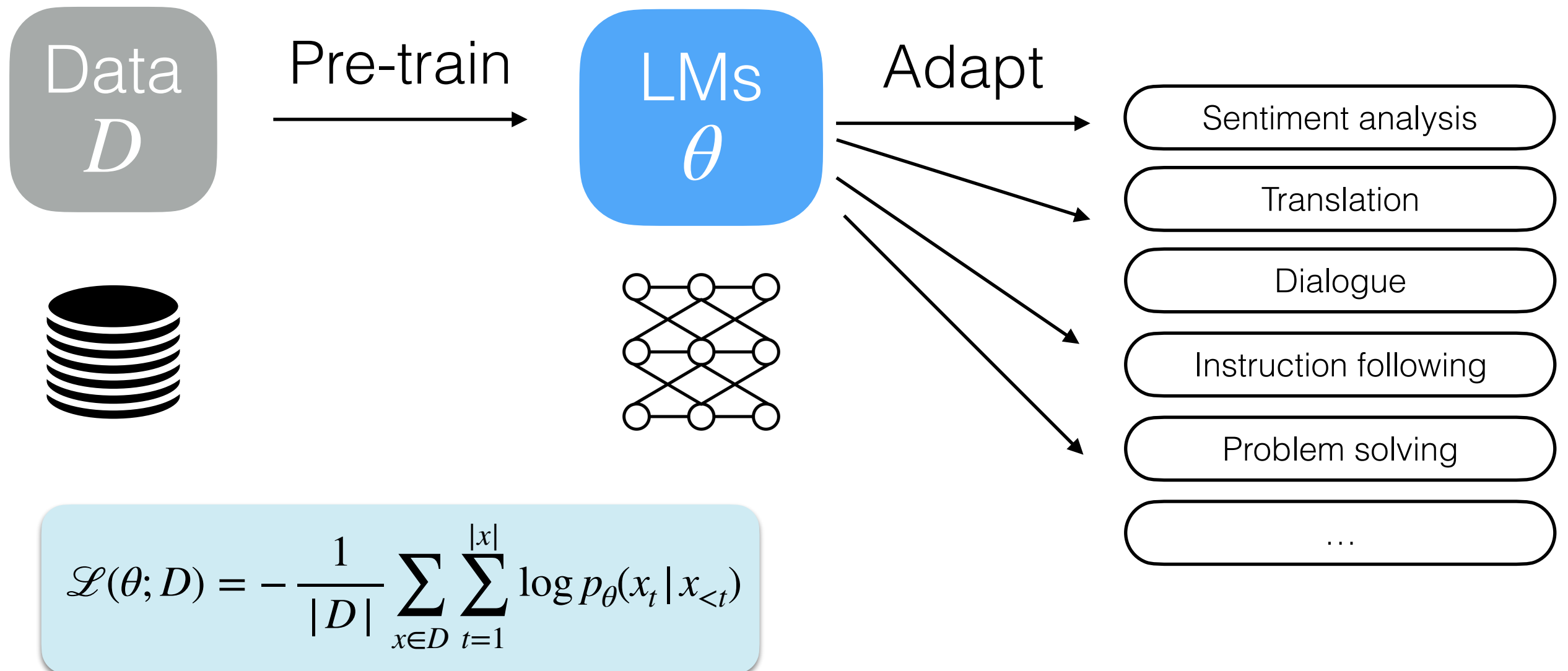
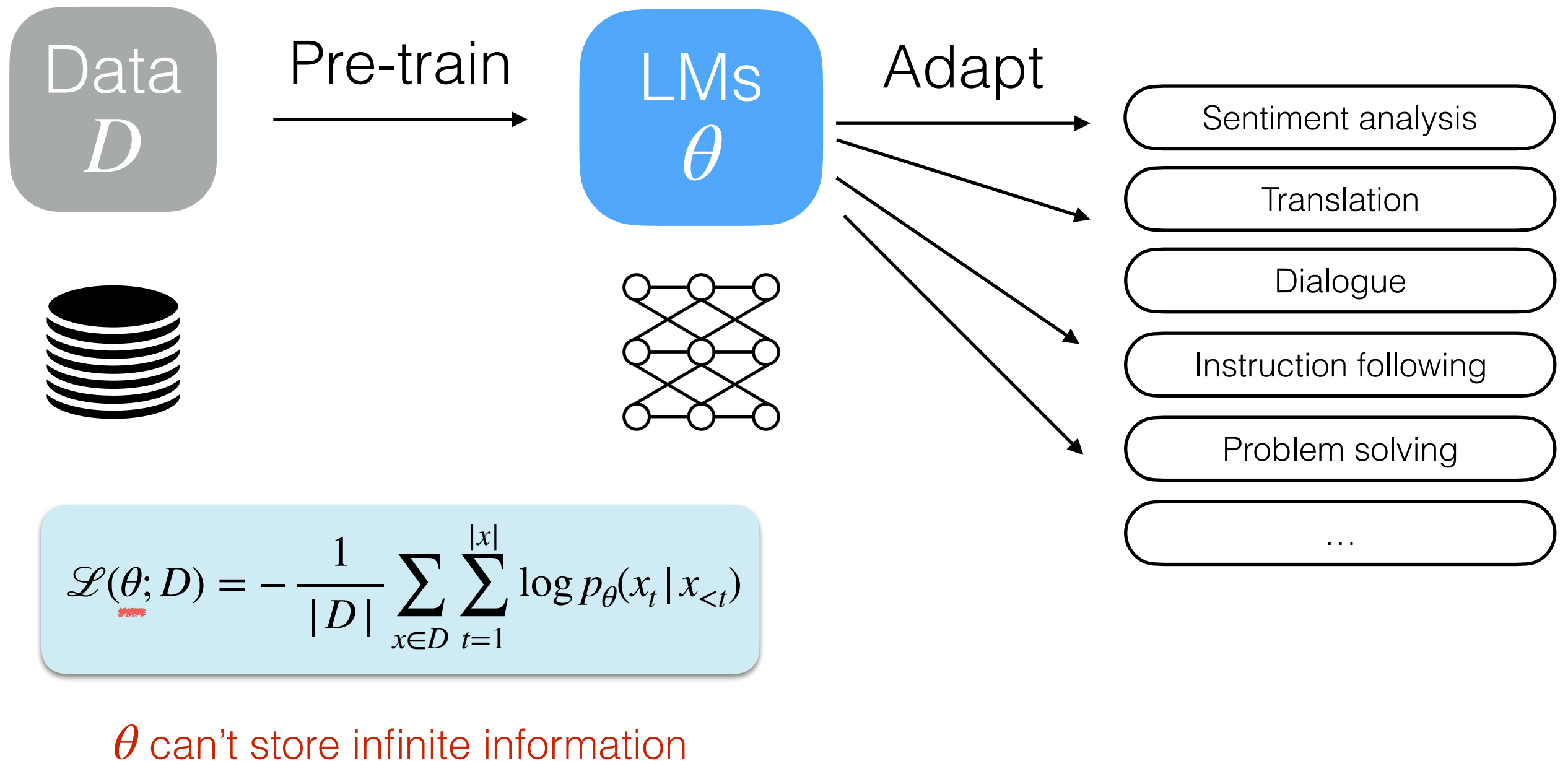OpenAI
https://openai.com › index › introducing-gpt-5-2

**Introducing GPT-5.2**

Dec 11, 2025 — **GPT-5.2** Instant is a fast, capable workhorse for everyday work improvements in info-seeking questions, how-tos and walk ... Read more

# Limitations of Monolithic LMs

Data $D$  →  **Pre-train**  →  LMs $\theta$  →  **Adapt**  →

- Sentiment analysis
- Translation
- Dialogue
- Instruction following
- Problem solving
- …

$$\mathcal{L}(\theta; D) = -\frac{1}{|D|} \sum_{x \in D} \sum_{t=1}^{|x|} \log p_\theta(x_t \mid x_{<t})$$

# Limitations of Monolithic LMs

Data $D$

Pre-train →

LMs $\theta$

Adapt →

- Sentiment analysis
- Translation
- Dialogue
- Instruction following
- Problem solving
- …

$$\mathscr{L}(\theta; D) = -\frac{1}{|D|} \sum_{x \in D} \sum_{t=1}^{|x|} \log p_\theta(x_t | x_{<t})$$

$\theta$ can't store infinite information

# Limitations of Monolithic LMs

Data
$D$

Pre-train →

LMs
$\theta$

Adapt →

- Sentiment analysis
- Translation
- Dialogue
- Instruction following
- Problem solving
- …

$$\mathcal{L}(\theta; D) = -\frac{1}{|D|} \sum_{x \in D} \sum_{t=1}^{|x|} \log p_\theta(x_t | x_{<t})$$

$\theta$ can't store infinite information

Allen-Zhu et al 2024. Physics of Language Models: Part 3.3, Knowledge Capacity Scaling Laws.

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

# Limitations of Monolithic LMs

Data $D$

Pre-train →

LMs $\theta$

Adapt →

- Sentiment analysis
- Translation
- Dialogue
- Instruction following
- Problem solving
- …

$$\mathcal{L}(\theta; D) = -\frac{1}{|D|} \sum_{x \in D} \sum_{t=1}^{|x|} \log p_\theta(x_t \mid x_{<t})$$

$\theta$ can't store infinite information

# Limitations of Monolithic LMs



$$\mathscr{L}(\theta; D) = -\frac{1}{|D|} \sum_{x \in D} \sum_{t=1}^{|x|} \log p_\theta(x_t \mid x_{<t})$$

$\theta$ can't store infinite information

$D$ only capture information as of the data cutoff time

Kasai et al 2023. RealTime QA: What's the Answer Right Now?

# Retrieval-Augmented LMs: Intuitions

# Retrieval-Augmented LMs: Intuitions

Who is <u>Akari Asai</u>?

**LM**

Japanese voice actress ✖

# Retrieval-Augmented LMs: Intuitions

Who is <u>Akari Asai</u>?

↓

Who is <u>Akari Asai</u>?

↓

**LM**

↓

Japanese voice actress ✖

# Retrieval-Augmented LMs: Intuitions

Who is <u>Akari Asai</u>?

Akari Asai is a Ph.D. student in the Paul G. Allen School of

Who is <u>Akari Asai</u>?

**LM**

Japanese voice actress ✗

# Retrieval-Augmented LMs: Intuitions

Who is _Akari Asai_?

Akari Asai is a Ph.D. student in the Paul G. Allen School of

**LM**

Ph.D. student at UW

Who is _Akari Asai_?

**LM**

Japanese voice actress ✖

# Retrieval-Augmented LMs: Intuitions



Who is <u>Akari Asai</u>?

Akari Asai is a Ph.D. student in the Paul G. Allen School of

**LM**

Ph.D. student at UW ✔

Who is <u>Akari Asai</u>?

**LM**

Japanese voice actress ✘

# Retrieval-Augmented LMs: Intuitions

Who is <u>Akari Asai</u>?

Akari Asai is a Ph.D. student in the Paul G. Allen School of

**LM**

Ph.D. student at UW ✔️

Who is <u>Akari Asai</u>?

**LM**

Japanese voice actress ❌

☺ Reduces hallucinations

# Retrieval-Augmented LMs: Intuitions

Who is <u>Akari Asai</u>?

Akari Asai is a Ph.D. student in the Paul G. Allen School of

**LM**

Ph.D. student at UW ✔

☺ Reduces hallucinations

Who is <u>Akari Asai</u>?

**LM**

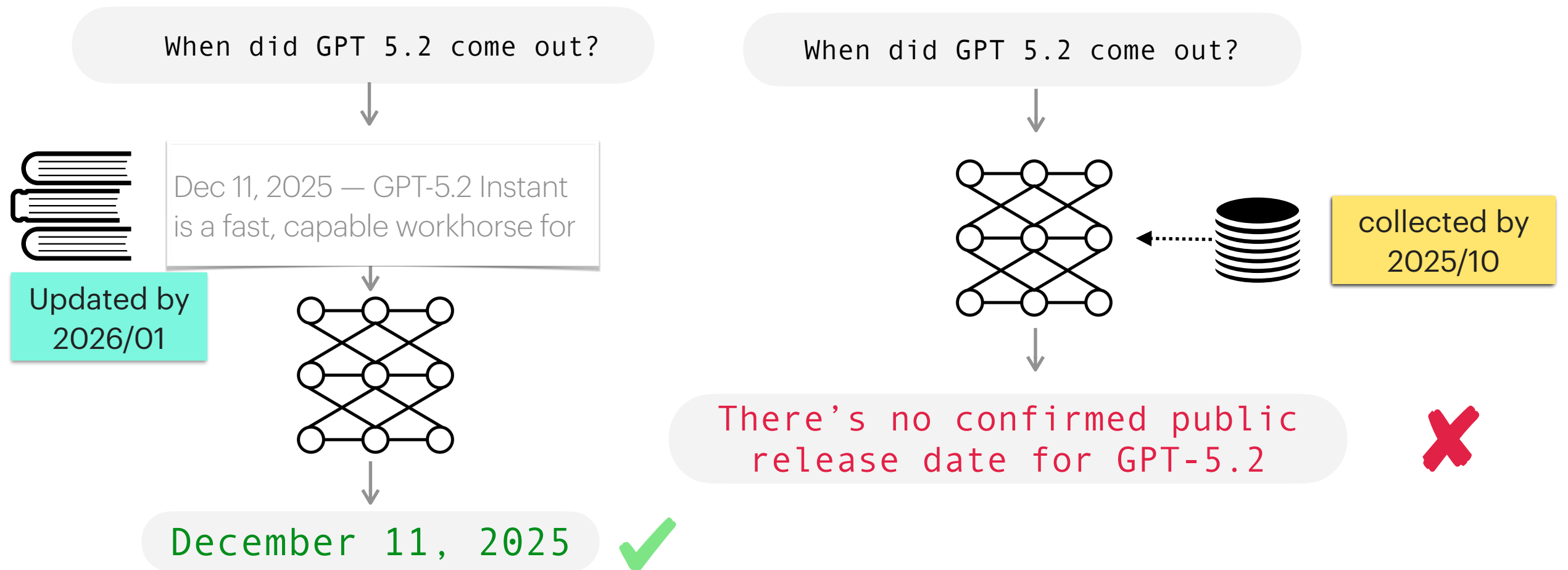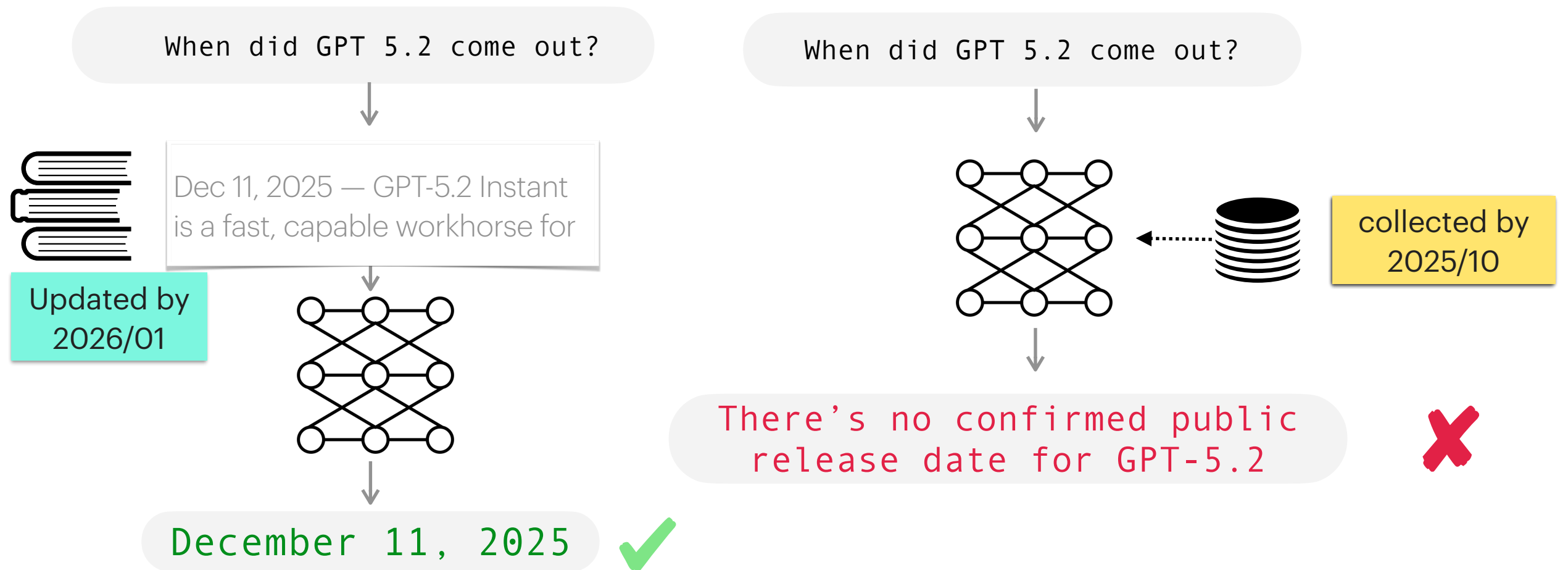Japanese voice actress ✖

☺ Parameter efficiency

# Retrieval-Augmented LMs: Intuitions

# Retrieval-Augmented LMs: Intuitions

# Retrieval-Augmented LMs: Intuitions

# Retrieval-Augmented LMs: Intuitions

# Retrieval-Augmented LMs: Intuitions
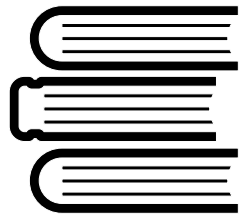
When did GPT 5.2 come out?

When did GPT 5.2 come out?

collected by 2025/10

There's no confirmed public release date for GPT-5.2 ✘

# Retrieval-Augmented LMs: Intuitions

When did GPT 5.2 come out?

Dec 11, 2025 — GPT-5.2 Instant is a fast, capable workhorse for

Updated by 2026/01

When did GPT 5.2 come out?

collected by 2025/10

There's no confirmed public release date for GPT-5.2

✖

# Retrieval-Augmented LMs: Intuitions

# Retrieval-Augmented LMs: Intuitions



When did GPT 5.2 come out?

Dec 11, 2025 — GPT-5.2 Instant is a fast, capable workhorse for

Updated by 2026/01

December 11, 2025 ✔

When did GPT 5.2 come out?

collected by 2025/10

There's no confirmed public release date for GPT-5.2 ✘

# Retrieval-Augmented LMs: Intuitions

When did GPT 5.2 come out?

Dec 11, 2025 — GPT-5.2 Instant is a fast, capable workhorse for

Updated by 2026/01

December 11, 2025 ✅

When did GPT 5.2 come out?

collected by 2025/10

There's no confirmed public release date for GPT-5.2 ❌

☺ Update knowledge w/o retraining

# Retrieval-Augmented LMs: Intuitions

When did GPT 5.2 come out?

Dec 11, 2025 — GPT-5.2 Instant is a fast, capable workhorse for

Updated by 2026/01

December 11, 2025 ✅

☺ Update knowledge w/o retraining

When did GPT 5.2 come out?

collected by 2025/10

There's no confirmed public release date for GPT-5.2 ❌

☺ Improve verifiability

# Overview

**Datastore**　　　　**Retriever**　　　　**LM**

# Overview

$x$ : When did GPT 5.2 come out?

**Datastore**

**Retriever**

**LM**

# Overview

$x$ : When did GPT 5.2 come out?

**Datastore**

**Retriever**

**LM**

$\text{Sim}(\,\cdot\,|\,x)$

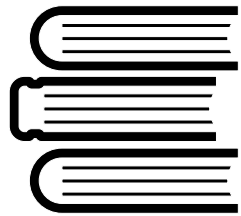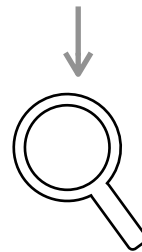| | |
|---|---|
| Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for | 0.9 |
| Ruby on Rails 5.2 Release Notes Ruby on Rails 5.2 Release Notes | 0.1 |
| GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was | 0.7 |

# Overview

$x$ : When did GPT 5.2 come out?

**Datastore**          **Retriever**          **LM**

$$D \in \mathrm{Top}_k \mathrm{Sim}( \cdot \,|\, x)$$

| | |
|---|---|
| Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for | 0.9 |
| Ruby on Rails 5.2 Release Notes Ruby on Rails 5.2 Release Notes | 0.1 |
| GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was | 0.7 |

# Overview

$x$ : When did GPT 5.2 come out?

**Datastore**            **Retriever**            **LM**

$$D \in \text{Top}_k \text{Sim}(\,\cdot\,|\,x)$$

Dec 11, 2025 - Open AIL
GPT-5.2 Instant is a fast,
capable workhorse for          0.9

Ruby on Rails 5.2 Release Notes Ruby
on Rails 5.2 Release Notes          0.1

GPT-5 (Wikipedia) Preceded in
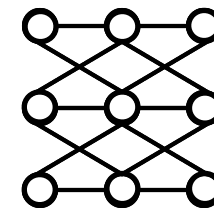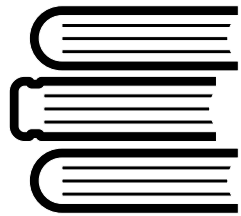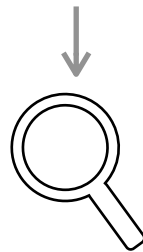the series by GPT-4, it was          0.7

# Overview



$x$ : When did GPT 5.2 come out?

$D$

**Datastore**

**Retriever**

**LM**

$$D \in \text{Top}_k \text{Sim}(\,\cdot\,|x)$$

Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for — 0.9

Ruby on Rails 5.2 Release Notes Ruby on Rails 5.2 Release Notes — 0.1

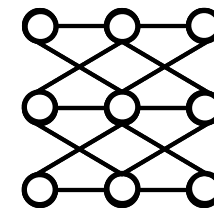GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was — 0.7

# Overview



$x$ : When did GPT 5.2 come out?

$D$

$y$ : Dec, 2026

**Datastore**

**Retriever**

**LM**

$$D \in \text{Top}_k \text{Sim}(\,\cdot\,|x)$$

Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for — 0.9

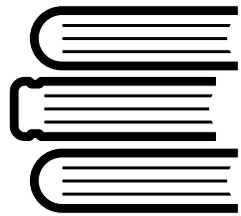Ruby on Rails 5.2 Release Notes Ruby on Rails 5.2 Release Notes — 0.1

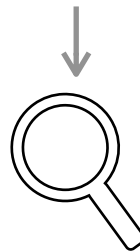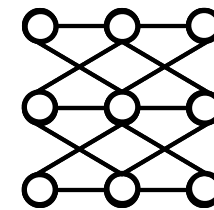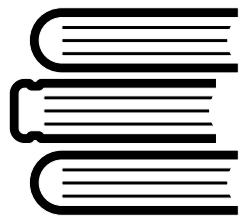GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was — 0.7

# Key Factors & Design Choices

$x$ : When did GPT 5.2 come out?



$D$

$y$ : Dec, 2026

**Datastore**

**Retriever**

**LM**

✔ Sources of datastore

✔ Processing

✔ Scaling

# Key Factors & Design Choices



$x$: When did GPT 5.2 come out?

$D$

$y$: Dec, 2026

**Datastore**

**Retriever**

**LM**

✔ Types of retrievers

✔ Training

✔ Evaluations

# Key Factors & Design Choices



$x$ : When did GPT 5.2 come out?

$D$

$y$ : Dec, 2026

**Datastore**

**Retriever**

**LM**

✔ Architectures

✔ Training

✔ Inference

# Part 1: Datastore



$x$ : When did GPT 5.2 come out?

$D$

$y$ : Dec, 2026

**Datastore**

**Retriever**

**LM**

✔ Sources of datastore

✔ Processing

✔ Scaling

# What Should be in "data store"?

# What Should be in "data store"?

$x$ : when did GPT 5.2
       come out?

# What Should be in "data store"?

$x$ : when did GPT 5.2
      come out?

**English Wikipedia**



Chen et al., 2017; Gu et al., 2020;
Asai et al., 2020; Guu et al., 2021;
Lewis et al., 2021 … etc

https://dumps.wikimedia.org/

# What Should be in "data store"?

$x$: when did GPT 5.2 come out?

$x$: How should I implement RAG using LlamaIndex?

**English Wikipedia**



Chen et al., 2017; Gu et al., 2020;
Asai et al., 2020; Guu et al., 2021;
Lewis et al., 2021 … etc

https://dumps.wikimedia.org/

# What Should be in "data store"?

$x$: when did GPT 5.2 come out?

$x$: How should I implement RAG using LlamaIndex?

**English Wikipedia**



Chen et al., 2017; Gu et al., 2020;
Asai et al., 2020; Guu et al., 2021;
Lewis et al., 2021 … etc

https://dumps.wikimedia.org/

**Code snippets**



**Documentations**



**Community forums**

# Scaling Up Datastore

# Scaling Up Datastore



**English Wikipedia** 5B tokens

# Scaling Up Datastore

# Scaling Up Datastore



**English Wikipedia** 5B tokens

**MassiveDS**
1.4 trillion tokens

# Scaling Up Datastore



**MassiveDS**
1.4 trillion tokens

## Datastore Scaling

### Language Modeling ↓
(RedPajama)

### Downstream Task ↑
(MMLU, 5-shot)



Llama-2 7B — Llama-2 13B — Llama-3 8B
● w/ Retrieval   ✖ LM-Only

# Scaling Up Datastore



**MassiveDS**
1.4 trillion tokens

Shao et al. 2024. Scaling Retrieval-Based Language Models with a Trillion-Token Datastore

# Processing Documents

# Processing Documents

# Processing Documents

# Processing Documents

## Curate and preprocess data

e.g., HTML -> Plain text

# Processing Documents

**Curate and preprocess data** $\longrightarrow$

e.g., HTML -> Plain text

# Processing Documents

## Curate and preprocess data $\longrightarrow$ Chunking

e.g., HTML -> Plain text

Paragraph-level (e.g., \n)
Every k words (e.g., 100-250)

GPT-4

Article   Talk

From Wikipedia

Generative
language mo
of GPT foun
made public
OpenAI's AP
transformer-
using both p
used to pred
tuned with re
alignment an

OpenAI

Article   Talk

From Wikipedia, the free encyclopedia

*Not to be confused with OpenAL, OpenAPI, or Open-source artificial inte*

**OpenAI** is an American Artificial Intelligence (AI) research organization
founded in December 2015 and headquartered in San Francisco,
California. Its stated mission is to develop "safe and beneficial" artificial
general intelligence (AGI), which it defines as "highly autonomous systems
that outperform humans at most economically valuable work".[5] As a
leading organization in the ongoing AI boom,[6] OpenAI is known for the
GPT family of large language models, the DALL-E series of text-to-image
models, and a text-to-video model named Sora.[7][8] Its release of
ChatGPT in November 2022 has been credited with catalyzing widespread
interest in generative AI.

GPT-4o is a pre-trained

Transformers is a series of
science fiction action films based

GPT4o was released by OpenAI in May

@I$O@

# Processing Documents

**Curate and preprocess data** → **Chunking** →

e.g., HTML -> Plain text

Paragraph-level (e.g., \n)
Every k words (e.g., 100-250)



GPT-4o is a pre-trained

Transformers is a series of science fiction action films based

GPT4o was released by OpenAI in May

@I$O@

# Processing Documents

**Curate and preprocess data** ⟶ **Chunking** ⟶ **Post-processing**

e.g., HTML -> Plain text

Paragraph-level (e.g., \n)
Every k words (e.g., 100-250)

e.g., Remove short documents

# Summary of Part 1



$x$ : When did GPT 5.2 come out?

$D$

$y$ : Dec, 2026

**Datastore**

**Retriever**

**LM**

✔ Sources of datastore

✔ Processing

✔ Scaling

- Choosing **the right datastore**
- **Chunking** and **filtering**
- **Scaling** datastores offer performance gain while adding challenges

# Part 2: Retriever

$x$ : When did GPT 5.2 come out?

$D$

$y$ : Dec, 2026

**Datastore**

**Retriever**

**LM**

✔ Types of retrievers

✔ Training

✔ Evaluations

# Types of Retrievers

$$D \in \text{Top}_k \text{Sim}(\,\cdot\,|\,x)$$

# Types of Retrievers

$$D \in \text{Top}_k \text{Sim}(\,\cdot\,|\,x)$$

**Sparse retrievers**

- **Sim:** Term-frequency based embeddings

e.g., TF-IDF, BM25

# Types of Retrievers

$$D \in \text{Top}_k \text{Sim}( \cdot \mid x)$$

## Sparse retrievers

- **Sim:** Term-
  frequency based
  embeddings

e.g., TF-IDF, BM25

## Dense retrievers

- **Sim**: dense
  embeddings
  encoded by pre-

e.g., DPR, Contriever, ColBERT

# Types of Retrievers

$$D \in \text{Top}_k \text{Sim}( \cdot \mid x)$$

**Sparse retrievers**

- **Sim:** Term-frequency based embeddings

e.g., TF-IDF, BM25

**Dense retrievers**

- **Sim**: dense embeddings encoded by pre-

e.g., DPR, Contriever, ColBERT

**Rerankers**

- **Sim**: Scores based on jointly encoded query and doc

e.g., cross-encoder reranker

# Types of Retrievers

$$D \in \text{Top}_k \text{Sim}( \cdot \, | \, x)$$

**Sparse retrievers**

- **Sim:** Term-frequency based embeddings

e.g., TF-IDF, BM25

$\$$

**Dense retrievers**

- **Sim**: dense embeddings encoded by pre-

e.g., DPR, Contriever, ColBERT

**Rerankers**

- **Sim**: Scores based on jointly encoded query and doc

e.g., cross-encoder reranker

# Types of Retrievers

$$D \in \text{Top}_k \text{Sim}(\,\cdot\mid x)$$

## Sparse retrievers

- **Sim:** Term-frequency based embeddings

e.g., TF-IDF, BM25

## Dense retrievers

- **Sim**: dense embeddings encoded by pre-

e.g., DPR, Contriever, ColBERT

## Rerankers

- **Sim**: Scores based on jointly encoded query and doc

e.g., cross-encoder reranker

# Types of Retrievers

$$D \in \text{Top}_k \text{Sim}(\,\cdot\,|x)$$

**Sparse retrievers**

- **Sim:** Term-
  frequency based
  embeddings

e.g., TF-IDF, BM25

$

**Dense retrievers**

- **Sim**: dense
  embeddings
  encoded by pre-

e.g., DPR, Contriever, ColBERT

$ $

**Rerankers**

- **Sim**: Scores based
  on jointly encoded
  query and doc

e.g., cross-encoder reranker

$ $ $

# Sparse Retrievers

q=what is nlp    d₁ = what is life ?    d₂ = nlp is an    d₃ = I like to
                 candy is life !        acronym for        do good
                                        natural language   research on

|           | q | d₁ | d₂ | d₃ |
|-----------|---|----|----|----|
| what      | 1 | 1  | 0  | 0  |
| candy     | 0 | 1  | 0  | 0  |
| nlp       | 1 | 0  | 1  | 1  |
| is        | 1 | 1  | 1  | 0  |
| language  | 0 | 0  | 0  | 0  |
| life      | 0 | 1  | 0  | 0  |
| …         | … | …  | …  | …  |

# Sparse Retrievers

| | q=what is nlp | d₁ = what is life ? candy is life ! | d₂ = nlp is an acronym for natural language | d₃ = I like to do good research on |
|---|---|---|---|---|
| what | 1 | 1 | 0 | 0 |
| candy | 0 | 1 | 0 | 0 |
| nlp | 1 | 0 | 1 | 1 |
| is | 1 | 1 | 1 | 0 |
| language | 0 | 0 | 0 | 0 |
| life | 0 | 1 | 0 | 0 |
| … | … | … | … | … |

**Check if a term appears in a document**

22

# One-hot Vector

q=what is nlp    $d_1$ = what is life ?    $d_2$ = nlp is an    $d_3$ = I like to
                 candy is life !           acronym for          do good
                                           natural language     research on

| | | | | |
|---|---|---|---|---|
| what | 1 | 1 | 0 | 0 |
| candy | 0 | 1 | 0 | 0 |
| nlp | 1 | 0 | 1 | 1 |
| is | 1 | 1 | 1 | 0 |
| language | 0 | 0 | 0 | 0 |
| life | 0 | 1 | 0 | 0 |
| … | … | … | … | … |

23

# One-hot Vector

q=what is nlp
d₁ = what is life ? candy is life !
d₂ = nlp is an acronym for natural language
d₃ = I like to do good research on

| | q | d₁ | d₂ | d₃ |
|---|---|---|---|---|
| what | 1 | 1 | 0 | 0 |
| candy | 0 | 1 | 0 | 0 |
| nlp | 1 | 0 | 1 | 1 |
| is | 1 | 1 | 1 | 0 |
| language | 0 | 0 | 0 | 0 |
| life | 0 | 1 | 0 | 0 |
| … | … | … | … | … |

**Check if a term appears in a document**

23

# Term-count Vector

|  | q=what is nlp | $d_1$ = what is life ? candy is life ! | $d_2$ = nlp is an acronym for natural language | $d_3$ = I like to do good research on |
|---|---|---|---|---|
| what | 1 | 1 | 0 | 0 |
| candy | 0 | 1 | 0 | 0 |
| nlp | 1 | 0 | 1 | 1 |
| is | 1 | 1 | 1 | 0 |
| language | 0 | 0 | 0 | 0 |
| life | 0 | 2 | 0 | 0 |
| … | … | … | … | … |

# Term-count Vector

q=what is nlp

d$_1$ = what is life ? candy is life !

d$_2$ = nlp is an acronym for natural language

d$_3$ = I like to do good research on

| | q | d$_1$ | d$_2$ | d$_3$ |
|---|---|---|---|---|
| what | 1 | 1 | 0 | 0 |
| candy | 0 | 1 | 0 | 0 |
| nlp | 1 | 0 | 1 | 1 |
| is | 1 | 1 | 1 | 0 |
| language | 0 | 0 | 0 | 0 |
| life | 0 | 2 | 0 | 0 |
| … | … | … | … | … |

**Count the number of appearances in a doc**

24

# Weighted-term Score

$$\mathrm{TF}(t, d) = \frac{\mathrm{freq}(t, d)}{\sum_{t'} \mathrm{freq}(t', d)} \qquad \mathrm{IDF}(t) = \log\left(\frac{|D|}{\sum_{d' \in D} \delta(\mathrm{freq}(t, d') > 0)}\right)$$

$$\mathrm{TF\text{-}IDF}(t, d) = \mathrm{TF}(t, d) \times \mathrm{IDF}(t)$$

$$\mathrm{BM\text{-}25}(t, d) = \mathrm{IDF}(t) \cdot \frac{\mathrm{freq}(t, d) \cdot (k_1 + 1)}{\mathrm{freq}(t, d) + k_1 \cdot \left(1 - b + b \cdot \frac{|d|}{\mathrm{avgdl}}\right)}$$

Robertson et al. 2009.  The Probabilistic Relevance Framework: BM25 and Beyond.

# Weighted-term Score

d$_1$ = <u>what</u> is life ?
  candy is life !

$$\text{TF}(t, d) = \frac{\text{freq}(t, d)}{\sum_{t'} \text{freq}(t', d)} \qquad \text{IDF}(t) = \log \left( \frac{|D|}{\sum_{d' \in D} \delta(\text{freq}(t, d') > 0)} \right)$$

t$_1$ = what

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \text{IDF}(t)$$

$$\text{BM-25}(t, d) = \text{IDF}(t) \cdot \frac{\text{freq}(t, d) \cdot (k_1 + 1)}{\text{freq}(t, d) + k_1 \cdot \left( 1 - b + b \cdot \frac{|d|}{\text{avgdl}} \right)}$$

Robertson et al. 2009. The Probabilistic Relevance Framework: BM25 and Beyond.

25

# Weighted-term Score

d₁ = <u>what</u> is life ?
candy is life !

$$\text{TF}(t, d) = \frac{\text{freq}(t, d)}{\sum_{t'} \text{freq}(t', d)}$$

t₁ = what

\# of documents

$$\text{IDF}(t) = \log \left( \frac{|D|}{\sum_{d' \in D} \delta(\text{freq}(t, d') > 0)} \right)$$

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \text{IDF}(t)$$

$$\text{BM-25}(t, d) = \text{IDF}(t) \cdot \frac{\text{freq}(t, d) \cdot (k_1 + 1)}{\text{freq}(t, d) + k_1 \cdot \left( 1 - b + b \cdot \frac{|d|}{\text{avgdl}} \right)}$$

Robertson et al. 2009.  The Probabilistic Relevance Framework: BM25 and Beyond.

# Weighted-term Score

d$_1$ = <u>what</u> is life ?
   candy is life !

# of documents

$$\text{TF}(t,d) = \frac{\text{freq}(t,d)}{\sum_{t'} \text{freq}(t',d)} \qquad \text{IDF}(t) = \log\left(\frac{|D|}{\sum_{d' \in D} \delta(\text{freq}(t,d') > 0)}\right)$$

t$_1$ = what

# of documents **where term t appears**

$$\text{TF-IDF}(t,d) = \text{TF}(t,d) \times \text{IDF}(t)$$

$$\text{BM-25}(t,d) = \text{IDF}(t) \cdot \frac{\text{freq}(t,d) \cdot (k_1 + 1)}{\text{freq}(t,d) + k_1 \cdot \left(1 - b + b \cdot \frac{|d|}{\text{avgdl}}\right)}$$

Robertson et al. 2009.  The Probabilistic Relevance Framework: BM25 and Beyond.

# Weighted-term Vectors

q=what is nlp    $d_1$ = what is life ?
candy is life !

$d_2$ = nlp is an
acronym for
natural language

$d_3$ = I like to
do good
research on

| | q | $d_1$ | $d_2$ | $d_3$ |
|---|---|---|---|---|
| what | 0.36 | 0.18 | 0 | 0 |
| candy | 0 | 0.18 | 0 | 0 |
| nlp | 0.13 | 0 | 0.05 | 0.05 |
| is | 0.13 | 0.13 | 0.05 | 0 |
| language | 0 | 0 | 0.13 | 0 |
| life | 0 | 0.36 | 0 | 0 |
| … | … | … | … | … |

Compute TF-IDF weights to build weighted vectors

# Weighted-term Vectors

q=what is nlp     d₁ = what is life ?     d₂ = nlp is an     d₃ = I like to
                   candy is life !          acronym for        do good
                                          natural language   research on

| | q | d₁ | d₂ | d₃ |
|---|---|---|---|---|
| what | 0.36 | 0.18 | 0 | 0 |
| candy | 0 | 0.18 | 0 | 0 |
| nlp | 0.13 | 0 | 0.05 | 0.05 |
| is | 0.13 | 0.13 | 0.05 | 0 |
| language | 0 | 0 | 0.13 | 0 |
| life | 0 | 0.36 | 0 | 0 |
| … | … | … | … | … |

**Compute TF-IDF weights to build weighted vectors**

# Weighted-term Vectors

q=what is nlp     d$_1$ = what is life ? candy is life !     d$_2$ = nlp is an acronym for natural language     d$_3$ = I like to do good research on

| | q | d$_1$ | d$_2$ | d$_3$ |
|---|---|---|---|---|
| what | 0.36 | 0.18 | 0 | 0 |
| candy | 0 | 0.18 | 0 | 0 |
| nlp | 0.13 | 0 | 0.05 | 0.05 |
| is | 0.13 | 0.13 | 0.05 | 0 |
| language | 0 | 0 | 0.13 | 0 |
| life | 0 | 0.36 | 0 | 0 |
| … | … | … | … | … |

**Compute TF-IDF weights to build weighted vectors**

# Weighted-term Vectors

q=what is nlp    d$_1$ = what is life ?    d$_2$ = nlp is an    d$_3$ = I like to
                 candy is life !           acronym for        do good
                                           natural language   research on

| | q | d$_1$ | d$_2$ | d$_3$ |
|---|---|---|---|---|
| what | 0.36 | 0.18 | 0 | 0 |
| candy | 0 | 0.18 | 0 | 0 |
| nlp | 0.13 | 0 | 0.05 | 0.05 |
| is | 0.13 | 0.13 | 0.05 | 0 |
| language | 0 | 0 | 0.13 | 0 |
| life | 0 | 0.36 | 0 | 0 |
| … | … | … | … | … |

**Compute TF-IDF weights to build weighted vectors**

# Compute cosine similarity

|  | q=what is nlp | $d_1$ = what is life ? candy is life ! | $d_2$ = nlp is an acronym for natural language | $d_3$ = I like to do good research on |
|---|---|---|---|---|
| what | 0.36 | 0.18 | 0 | 0 |
| candy | 0 | 0.18 | 0 | 0 |
| nlp | 0.13 | 0 | 0.05 | 0.05 |
| is | 0.13 | 0.13 | 0.05 | 0 |
| language | 0 | 0 | 0.13 | 0 |
| life | 0 | 0.36 | 0 | 0 |
| … | … | … | … | … |

$$q * d_1 = 0.44 \qquad q * d_2 = 0.21 \qquad q * d_3 = 0.32$$

27

# Compute cosine similarity

q=what is nlp      $d_1$ = what is life ?      $d_2$ = nlp is an      $d_3$ = I like to
                   candy is life !      acronym for      do good
                   natural language      research on

| | $d_1$ | $d_2$ | $d_3$ |
|---|---|---|---|
| what | 0.36 | 0.18 | 0 | 0 |
| candy | 0 | 0.18 | 0 | 0 |
| nlp | 0.13 | 0 | 0.05 | 0.05 |
| is | 0.13 | 0.13 | 0.05 | 0 |
| language | 0 | 0 | 0.13 | 0 |
| life | 0 | 0.36 | 0 | 0 |
| … | … | … | … | … |

$$q * d_1 = 0.44 \qquad q * d_2 = 0.21 \qquad q * d_3 = 0.32$$

# Compute cosine similarity

q=what is nlp    $d_1$ = what is life ? candy is life !    $d_2$ = nlp is an acronym for natural language    $d_3$ = I like to do good research on

what
candy
nlp
is
language
life
…

$$\begin{pmatrix} 0.36 \\ 0 \\ \\ \\ 0 \\ 0 \\ \\ … \end{pmatrix} \begin{pmatrix} 0.18 \\ 0.18 \\ \\ \\ 0 \\ 0.36 \\ \\ … \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ \\ \\ 0.13 \\ 0 \\ \\ … \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0.05 \\ 0 \\ 0 \\ 0 \\ \\ … \end{pmatrix}$$

Can't fully capture semantic similarities

$q * d_1 = 0.44$          $q * d_2 = 0.21$          $q * d_3 = 0.32$

# Dense Retrievers: Overview

FIFA World Cup 2026 will expand to

In 2022, the 32 national teams

Team USA celebrated after

# Dense Retrievers: Overview



FIFA World Cup 2026 will expand to → **Encoder**

In 2022, the 32 national teams → **Encoder**

Team USA celebrated after → **Encoder**

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

# Dense Retrievers: Overview

$x$ = How many teams will participate in FIFA World



FIFA World Cup 2026 will expand to

**Encoder**

In 2022, the 32 national teams

**Encoder**

Team USA celebrated after

**Encoder**

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

# Dense Retrievers: Overview

$x$ = How many teams will participate in FIFA World



$$\mathbf{z} = \mathrm{Encoder}(z)$$
$$\mathbf{x} = \mathrm{Encoder}(x)$$

29

# Dense Retrievers: Overview

**x** = How many teams will participate in FIFA World

FIFA World Cup 2026 will expand to

**Encoder**

In 2022, the 32 national teams

**Encoder**

Team USA celebrated after

**Encoder**

**Encoder**

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

29

# Dense Retrievers: Overview

$x$ = How many teams will participate in FIFA World

**Encoder**

**Encoder**

FIFA World Cup
2026 will expand to

**Encoder**

In 2022, the 32
national teams

**Encoder**

Team USA
celebrated after

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

$k$ retrieved chunks $\quad z_1, \ldots, z_k = \text{argTop-}k\,(\mathbf{x} \cdot \mathbf{z})$

# Dense Retrievers: Embeddings

- Use output vector of $\mathbb{R}^d$ [CLS] in masked LMs

  e.g., DPR



Pre-training          Fine-Tuning

- Mean / Max pooling of output vectors $\mathbb{R}^{N \times d}$ (can be applied to autoregressive LMs)

  e.g., SBERT, SGPT, Qwen Embeddings

Karpukhin et al 2020. Dense Passage Retrieval for Open-Domain Question Answering.
Muennighoff 2022. SGPT: GPT Sentence Embeddings for Semantic Search.

30

# Training Dense Retrievers

# Training Dense Retrievers



**Dot Product Similarity**

Encoder

Encoder

$q$

Query

Text chunks

**Contrastive learning**

Positive

Negative

Positive

Negative

# Training Dense Retrievers

## Dot Product Similarity



**Encoder**

$q$

**Query**

**Encoder**

**Text chunks**

## Contrastive learning



Positive

Negative

Positive

Negative

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

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

# Training Dense Retrievers

**Dot Product Similarity**

Positive

Negative

Positive

Negative

Encoder

Encoder

$q$

**Query**

**Text chunks**

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

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

# Training Dense Retrievers

**Dot Product Similarity**

**Encoder**

**Encoder**

$q$

**Query**

**Text chunks**

**Contrastive learning**

Positive

Negative

Positive

Negative

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

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

# Fast Nearest Neighbor Search

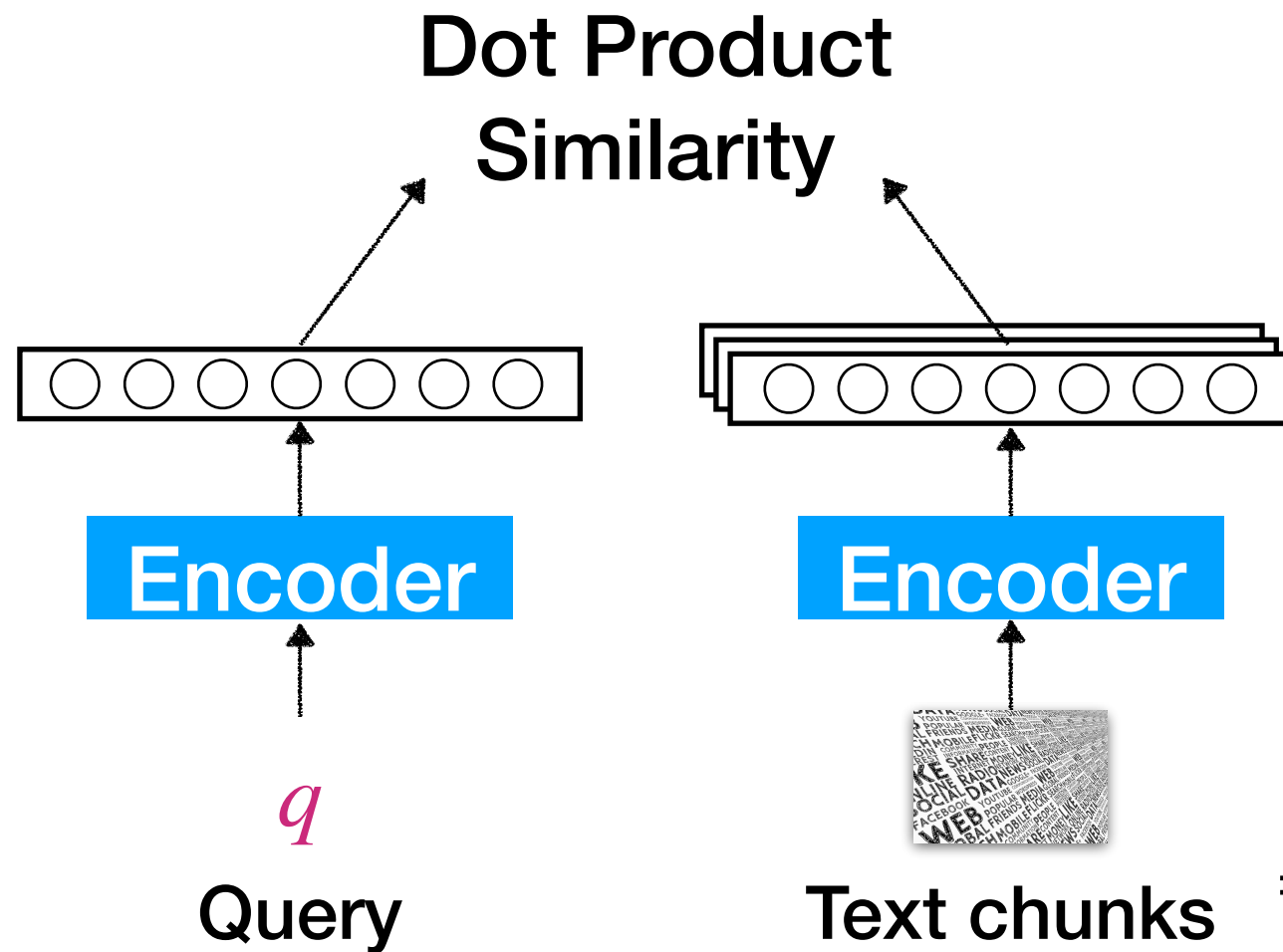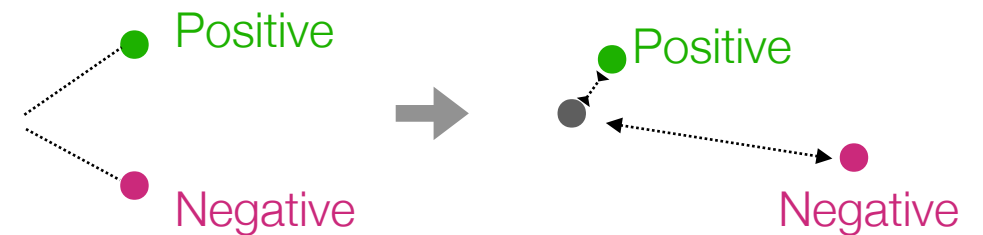| Method | Class name | index_factory | Main parameters | Bytes/vector | Exhausti |
|---|---|---|---|---|---|
| Exact Search for L2 | IndexFlatL2 | "Flat" | d | 4*d | yes |
| Exact Search for Inner Product | IndexFlatIP | "Flat" | d | 4*d | yes |
| Hierarchical Navigable Small World graph exploration | IndexHNSWFlat | "HNSW,Flat" | d, M | 4*d + x * M * 2 * 4 | no |
| Inverted file with exact post-verification | IndexIVFFlat | "IVFx,Flat" | quantizer, d, nlists, metric | 4*d + 8 | no |
| Locality-Sensitive Hashing (binary flat index) | IndexLSH | - | d, nbits | ceil(nbits/8) | yes |
| Scalar quantizer (SQ) in flat mode | IndexScalarQuantizer | "SQ8" | d | d | yes |
| Product quantizer (PQ) in flat mode | IndexPQ | "PQx", "PQ"M"x"nbits | d, M, nbits | ceil(M * nbits / 8) | yes |
| IVF and scalar quantizer | IndexIVFScalarQuantizer | "IVFx,SQ4" "IVFx,SQ8" | quantizer, d, nlists, qtype | SQfp16: 2 * d + 8, SQ8: d + 8 or SQ4: d/2 + 8 | no |
| IVFADC (coarse quantizer+PQ on residuals) | IndexIVFPQ | "IVFx,PQ"y"x"nbits | quantizer, d, nlists, M, nbits | ceil(M * nbits/8)+8 | no |

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

https://speakerdeck.com/matsui_528/cvpr20-tutorial-billion-scale-approximate-nearest-neighbor-search   (CVPR 2020 Tutorial)

32

# Fast Nearest Neighbor Search

| Method | Class name | index_factory | Main parameters | Bytes/vector | Exhausti... |
|---|---|---|---|---|---|
| Exact Search for L2 | IndexFlatL2 | "Flat" | d | 4*d | yes |
| Exact Search for Inner Product | IndexFlatIP | "Flat" | d | 4*d | yes |
| Hierarchical Navigable Small World graph exploration | IndexHNSWFlat | "HNSW,Flat" | d, M | 4*d + x * M * 2 * 4 | no |
| Inverted file with exact post-verification | IndexIVFFlat | "IVFx,Flat" | quantizer, d, nlists, metric | 4*d + 8 | no |
| Locality-Sensitive Hashing (binary flat index) | IndexLSH | - | d, nbits | ceil(nbits/8) | yes |
| Scalar quantizer (SQ) in flat mode | IndexScalarQuantizer | "SQ8" | d | d | yes |
| Product quantizer (PQ) in flat mode | IndexPQ | "PQx", "PQ"M"x"nbits | d, M, nbits | ceil(M * nbits / 8) | yes |
| IVF and scalar quantizer | IndexIVFScalarQuantizer | "IVFx,SQ4" "IVFx,SQ8" | quantizer, d, nlists, qtype | SQfp16: 2 * d + 8, SQ8: d + 8 or SQ4: d/2 + 8 | no |
| IVFADC (coarse quantizer+PQ on residuals) | IndexIVFPQ | "IVFx,PQ"y"x"nbits | quantizer, d, nlists, M, nbits | ceil(M * nbits/8)+8 | no |

**Exact search (still fast for $10^6 \sim 10^7$ scale)**

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

https://speakerdeck.com/matsui_528/cvpr20-tutorial-billion-scale-approximate-nearest-neighbor-search (CVPR 2020 Tutorial)

# Fast Nearest Neighbor Search

| Method | Class name | index_factory | Main parameters | Bytes/vector | Exhausti |
|--------|-----------|---------------|-----------------|--------------|----------|
| Exact Search for L2 | IndexFlatL2 | "Flat" | d | 4*d | yes |
| Exact Search for Inner Product | IndexFlatIP | "Flat" | d | 4*d | yes |
| Hierarchical Navigable Small World graph exploration | IndexHNSWFlat | "HNSW,Flat" | d, M | 4*d + x * M * 2 * 4 | no |
| Inverted file with exact post-verification | IndexIVFFlat | "IVFx,Flat" | quantizer, d, nlists, metric | 4*d + 8 | no |
| Locality-Sensitive Hashing (binary flat index) | IndexLSH | - | d, nbits | ceil(nbits/8) | yes |
| Scalar quantizer (SQ) in flat mode | IndexScalarQuantizer | "SQ8" | d | d | yes |
| Product quantizer (PQ) in flat mode | IndexPQ | "PQx", "PQ"M"x"nbits | d, M, nbits | ceil(M * nbits / 8) | yes |
| IVF and scalar quantizer | IndexIVFScalarQuantizer | "IVFx,SQ4" "IVFx,SQ8" | quantizer, d, nlists, qtype | SQfp16: 2 * d + 8, SQ8: d + 8 or SQ4: d/2 + 8 | no |
| IVFADC (coarse quantizer+PQ on residuals) | IndexIVFPQ | "IVFx,PQ"y"x"nbits | quantizer, d, nlists, M, nbits | ceil(M * nbits/8)+8 | no |

Exact search (still fast for 10^6~10^7 scale)

Approximate search (faster but more memory)

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

32

https://speakerdeck.com/matsui_528/cvpr20-tutorial-billion-scale-approximate-nearest-neighbor-search   (CVPR 2020 Tutorial)

# Fast Nearest Neighbor Search

| Method | Class name | index_factory | Main parameters | Bytes/vector | Exhausti |
|--------|-----------|---------------|-----------------|--------------|----------|
| Exact Search for L2 | IndexFlatL2 | "Flat" | d | 4*d | yes |
| Exact Search for Inner Product | IndexFlatIP | "Flat" | d | 4*d | yes |
| Hierarchical Navigable Small World graph exploration | IndexHNSWFlat | "HNSW,Flat" | d , M | 4*d + x * M * 2 * 4 | no |
| Inverted file with exact post-verification | IndexIVFFlat | "IVFx,Flat" | quantizer , d , nlists , metric | 4*d + 8 | no |
| Locality-Sensitive Hashing (binary flat index) | IndexLSH | - | d , nbits | ceil(nbits/8) | yes |
| Scalar quantizer (SQ) in flat mode | IndexScalarQuantizer | "SQ8" | d | d | yes |
| Product quantizer (PQ) in flat mode | IndexPQ | "PQx" , "PQ"M"x"nbits | d , M , nbits | ceil(M * nbits / 8) | yes |
| IVF and scalar quantizer | IndexIVFScalarQuantizer | "IVFx,SQ4" "IVFx,SQ8" | quantizer , d , nlists , qtype | SQfp16: 2 * d + 8, SQ8: d + 8 or SQ4: d/2 + 8 | no |
| IVFADC (coarse quantizer+PQ on residuals) | IndexIVFPQ | "IVFx,PQ"y"x"nbits | quantizer , d , nlists , M , nbits | ceil(M * nbits/8)+8 | no |

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

Exact search (still fast for $10^6$~$10^7$ scale)

Approximate search (faster but more memory)

Reduce index size with quantization

https://speakerdeck.com/matsui_528/cvpr20-tutorial-billion-scale-approximate-nearest-neighbor-search (CVPR 2020 Tutorial)

# Reranking with Cross Encoders

**Bi-Encoder**



Encoder

Encoder

FIFA World Cup
2026 will expand to

How many teams will
participate in FIFA

# Reranking with Cross Encoders

**Bi-Encoder**                    **Cross-Encoder**



Encoder          Encoder

FIFA World Cup          How many teams will
2026 will expand to          participate in FIFA

# Reranking with Cross Encoders

**Bi-Encoder**

**Cross-Encoder**



Encoder

Encoder

FIFA World Cup 2026 will expand to

How many teams will participate in FIFA

FIFA World Cup 2026 will expand to

How many teams will participate in FIFA World

# Reranking with Cross Encoders

**Bi-Encoder**

**Cross-Encoder**



Classifier

Encoder

Encoder

Encoder

FIFA World Cup 2026 will expand to

How many teams will participate in FIFA

FIFA World Cup 2026 will expand to

How many teams will participate in FIFA World
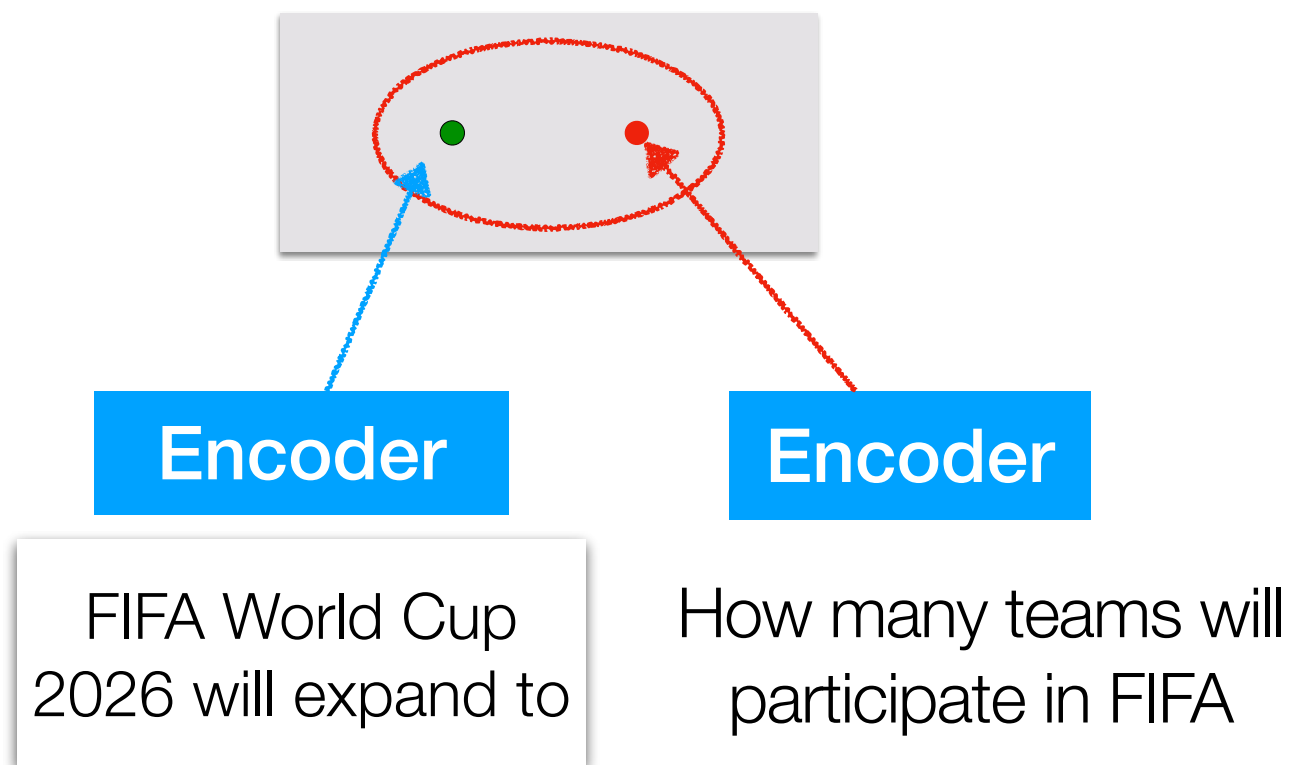
# Reranking with Cross Encoders

**Bi-Encoder**

**Cross-Encoder**

0.9

| Classifier |

| Encoder |

| Encoder | | Encoder |

FIFA World Cup 2026 will expand to

How many teams will participate in FIFA

FIFA World Cup 2026 will expand to

How many teams will participate in FIFA World

# Evaluation Metrics

# Evaluation Metrics

Evaluation of **unranked** retrieval sets

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} \qquad \text{Recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})}$$

# Evaluation Metrics

Evaluation of **unranked** retrieval sets

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} \qquad \text{Recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})}$$

Evaluation of **ranked** retrieval sets

$$\text{MAP}(Q) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{m_j} \sum_{k=1}^{m_j} \text{Precision}(R_{jk}) \qquad \text{NDCG}(Q,k) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} Z_{kj} \sum_{m=1}^{k} \frac{2^{R(j,m)} - 1}{\log_2(1+m)}$$

# Evaluation Metrics

Evaluation of **unranked** retrieval sets

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} \qquad \text{Recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})}$$

Evaluation of **ranked** retrieval sets

$$\text{MAP}(Q) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{m_j} \sum_{k=1}^{m_j} \text{Precision}(R_{jk}) \qquad \text{NDCG}(Q,k) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} Z_{kj} \sum_{m=1}^{k} \frac{2^{R(j,m)} - 1}{\log_2(1+m)}$$

nDCG@10 is widely used (e.g., BEIR)

# Retrieval Benchmark: MTEB



**Fact Checking**

FEVER
QUERY  Natural Claim
DOCS  Wikipedia Articles
*Wiki*

Climate-FEVER
QUERY  Climate-based Claim
DOCS  Wikipedia Articles
*Wiki*

SciFact
QUERY  Scientific claim
DOCS  PubMed Articles
*Scientific*

**Citation-Prediction**

SCIDOCS
QUERY  Article Title
DOCS  PubMed Articles
*Scientific*

**Dup. Question Retrieval**

Quora
QUERY  Query Title
DOCS  Quora Questions
*Quora*

CQADupStack
QUERY  Query Title
DOCS  Query Title + Body
*StackEx.*

**Argument Retrieval**

Tóuche-2020
QUERY  Controversial Query
DOCS  Args.me Arguments
*Misc.*

ArguAna
QUERY  Argument
DOCS  Idebate Arguments
*Misc.*

9 Tasks

18 Datasets

Beir
Benchmarking IR

**News Retrieval**

TREC-NEWS
QUERY  News Headline
DOCS  News Articles
*News*

Robust04
QUERY  News Query
DOCS  News Articles
*News*

**Question-Answering**

NQ
QUERY  Natural Query
DOCS  Wikipedia Articles
*Wiki*

HotpotQA
QUERY  Multi-Hop Query
DOCS  Wikipedia Articles
*Wiki*

FiQA-2018
QUERY  Financial Query
DOCS  Investment Articles
*Finance*

**Tweet Retrieval**

Signal-1M
QUERY  News Headline
DOCS  Twitter Tweets
*Twitter*

**Bio-Medical IR**

TREC-COVID
QUERY  COVID-19 Query
DOCS  CORD-19 Articles
*Scientific*

BioASQ
QUERY  Bio-Medical Query
DOCS  PubMed Articles
*Scientific*

NFCorpus
QUERY  Nutrition Facts
DOCS  PubMed Articles
*Scientific*

**Entity Retrieval**

DBPedia
QUERY  Entity-based Query
DOCS  DBPedia Articles
*Wiki*

Thakur et al. 2021.  BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models.

# BEIR Results

|              | BM25  | BM25+CE |
|--------------|-------|---------|
| MS MARCO     | 22.8  | 41.3    |
| Trec-COVID   | 65.6  | **75.7** |
| NFCorpus     | 32.5  | **35.0** |
| NQ           | 32.9  | 53.3    |
| HotpotQA     | 60.3  | 70.7    |
| FiQA         | 23.6  | 34.7    |
| ArguAna      | 31.5  | 31.1    |
| Touche-2020  | **36.7** | 27.1 |
| CQADupStack  | 29.9  | 37.0.   |
| Quora        | 78.9  | 82.5    |
| DBPedia      | 31.3  | 40.9    |
| Scidocs      | 15.8  | 16.6    |
| FEVER        | 75.3  | **81.9** |
| Climate-FEVER| 21.3  | 25.3    |
| Scifact      | 66.5  | 68.8    |
| Avg. w/o CQA | 44.0  | 49.5    |
| Avg.         | 43.0  | 48.6    |
| Best on      | 1     | 3       |

# BEIR Results

| | BM25 | BM25+CE |
|---|---|---|
| MS MARCO | 22.8 | 41.3 |
| Trec-COVID | 65.6 | **75.7** |
| NFCorpus | 32.5 | **35.0** |
| NQ | 32.9 | 53.3 |
| HotpotQA | 60.3 | 70.7 |
| FiQA | 23.6 | 34.7 |
| ArguAna | 31.5 | 31.1 |
| Touche-2020 | **36.7** | 27.1 |
| CQADupStack | 29.9 | 37.0. |
| Quora | 78.9 | 82.5 |
| DBPedia | 31.3 | 40.9 |
| Scidocs | 15.8 | 16.6 |
| FEVER | 75.3 | **81.9** |
| Climate-FEVER | 21.3 | 25.3 |
| Scifact | 66.5 | 68.8 |
| Avg. w/o CQA | 44.0 | 49.5 |
| Avg. | 43.0 | 48.6 |
| Best on | 1 | 3 |

Adding CE (cross-encoder) helps

Izacard et al. TMLR 2022. Unsupervised Dense Information Retrieval with Contrastive Learning.

# BEIR Results

|  | BM25 | BM25+CE | DPR |
|---|---|---|---|
| MS MARCO | 22.8 | 41.3 | 17.7 |
| Trec-COVID | 65.6 | **75.7** | 33.2 |
| NFCorpus | 32.5 | **35.0** | 18.9 |
| NQ | 32.9 | 53.3 | 47.4 |
| HotpotQA | 60.3 | 70.7 | 39.1 |
| FiQA | 23.6 | 34.7 | 11.2 |
| ArguAna | 31.5 | 31.1 | 17.5 |
| Touche-2020 | **36.7** | 27.1 | 13.1 |
| CQADupStack | 29.9 | 37.0. | 15.3 |
| Quora | 78.9 | 82.5 | 24.8 |
| DBPedia | 31.3 | 40.9 | 26.3 |
| Scidocs | 15.8 | 16.6 | 7.7 |
| FEVER | 75.3 | **81.9** | 56.2 |
| Climate-FEVER | 21.3 | 25.3 | 14.8 |
| Scifact | 66.5 | 68.8 | 31.8 |
| Avg. w/o CQA | 44.0 | 49.5 | 26.3 |
| Avg. | 43.0 | 48.6 | 25.5 |
| Best on | 1 | 3 | 0 |

Adding CE (cross-encoder) helps

Izacard et al. TMLR 2022. Unsupervised Dense Information Retrieval with Contrastive Learning.

# BEIR Results

|  | BM25 | BM25+CE | DPR |
|---|---|---|---|
| MS MARCO | 22.8 | 41.3 | 17.7 |
| Trec-COVID | 65.6 | **75.7** | 33.2 |
| NFCorpus | 32.5 | **35.0** | 18.9 |
| NQ | 32.9 | 53.3 | 47.4 |
| HotpotQA | 60.3 | 70.7 | 39.1 |
| FiQA | 23.6 | 34.7 | 11.2 |
| ArguAna | 31.5 | 31.1 | 17.5 |
| Touche-2020 | **36.7** | 27.1 | 13.1 |
| CQADupStack | 29.9 | 37.0. | 15.3 |
| Quora | 78.9 | 82.5 | 24.8 |
| DBPedia | 31.3 | 40.9 | 26.3 |
| Scidocs | 15.8 | 16.6 | 7.7 |
| FEVER | 75.3 | **81.9** | 56.2 |
| Climate-FEVER | 21.3 | 25.3 | 14.8 |
| Scifact | 66.5 | 68.8 | 31.8 |
| Avg. w/o CQA | 44.0 | 49.5 | 26.3 |
| Avg. | 43.0 | 48.6 | 25.5 |
| Best on | 1 | 3 | 0 |

Adding CE (cross-encoder) helps

Dense retrievers could struggle in OOD

Izacard et al. TMLR 2022. Unsupervised Dense Information Retrieval with Contrastive Learning.

# BEIR Results

|  | BM25 | BM25+CE | DPR | Ours | Ours+CE |
|---|---|---|---|---|---|
| MS MARCO | 22.8 | 41.3 | 17.7 | 40.7 | **47.0** |
| Trec-COVID | 65.6 | **75.7** | 33.2 | 59.6 | 70.1 |
| NFCorpus | 32.5 | **35.0** | 18.9 | 32.8 | 34.4 |
| NQ | 32.9 | 53.3 | 47.4 | 49.8 | **57.7** |
| HotpotQA | 60.3 | 70.7 | 39.1 | 63.8 | **71.5** |
| FiQA | 23.6 | 34.7 | 11.2 | 32.9 | **36.7** |
| ArguAna | 31.5 | 31.1 | 17.5 | 44.6 | 41.3 |
| Touche-2020 | **36.7** | 27.1 | 13.1 | 23.0 | 29.8 |
| CQADupStack | 29.9 | 37.0. | 15.3 | 34.5 | **37.7** |
| Quora | 78.9 | 82.5 | 24.8 | **86.5** | 82.4 |
| DBPedia | 31.3 | 40.9 | 26.3 | 41.3 | **47.1** |
| Scidocs | 15.8 | 16.6 | 7.7 | 16.5 | **17.1** |
| FEVER | 75.3 | **81.9** | 56.2 | 75.8 | **81.9** |
| Climate-FEVER | 21.3 | 25.3 | 14.8 | 23.7 | **25.8** |
| Scifact | 66.5 | 68.8 | 31.8 | 67.7 | 69.2 |
| Avg. w/o CQA | 44.0 | 49.5 | 26.3 | 47.5 | 51.2 |
| Avg. | 43.0 | 48.6 | 25.5 | 46.6 | 50.2 |
| Best on | 1 | 3 | 0 | 1 | 9 |

Adding CE (cross-encoder) helps

Dense retrievers could struggle in OOD

Izacard et al. TMLR 2022. Unsupervised Dense Information Retrieval with Contrastive Learning.

# BEIR Results

|  | BM25 | BM25+CE | DPR | Ours | Ours+CE |
|---|---|---|---|---|---|
| MS MARCO | 22.8 | 41.3 | 17.7 | 40.7 | **47.0** |
| Trec-COVID | 65.6 | **75.7** | 33.2 | 59.6 | 70.1 |
| NFCorpus | 32.5 | **35.0** | 18.9 | 32.8 | 34.4 |
| NQ | 32.9 | 53.3 | 47.4 | 49.8 | **57.7** |
| HotpotQA | 60.3 | 70.7 | 39.1 | 63.8 | **71.5** |
| FiQA | 23.6 | 34.7 | 11.2 | 32.9 | **36.7** |
| ArguAna | 31.5 | 31.1 | 17.5 | 44.6 | 41.3 |
| Touche-2020 | **36.7** | 27.1 | 13.1 | 23.0 | 29.8 |
| CQADupStack | 29.9 | 37.0. | 15.3 | 34.5 | **37.7** |
| Quora | 78.9 | 82.5 | 24.8 | **86.5** | 82.4 |
| DBPedia | 31.3 | 40.9 | 26.3 | 41.3 | **47.1** |
| Scidocs | 15.8 | 16.6 | 7.7 | 16.5 | **17.1** |
| FEVER | 75.3 | **81.9** | 56.2 | 75.8 | **81.9** |
| Climate-FEVER | 21.3 | 25.3 | 14.8 | 23.7 | **25.8** |
| Scifact | 66.5 | 68.8 | 31.8 | 67.7 | 69.2 |
| Avg. w/o CQA | 44.0 | 49.5 | 26.3 | 47.5 | 51.2 |
| Avg. | 43.0 | 48.6 | 25.5 | 46.6 | 50.2 |
| Best on | 1 | 3 | 0 | 1 | 9 |

Adding CE (cross-encoder) helps

Dense retrievers could struggle in OOD

Unsupervised training helps in OOD

Izacard et al. TMLR 2022. Unsupervised Dense Information Retrieval with Contrastive Learning.

# Retrieval Benchmark: MTEB

Thakur et al. 2021. BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models.
Muennighoff et al. 2022. MTEB: Massive Text Embedding Benchmark.

# Retrieval Benchmark: MTEB

Thakur et al. 2021.  BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models.
Muennighoff et al. 2022. MTEB: Massive Text Embedding Benchmark.

# Retrieval Benchmark: MTEB

## Clustering

ArxivP2P · ArxivS2S · BiorxivP2P · BiorxivS2S

MedrxivP2P · MedrxivS2S · Reddit · RedditP2P

StackExchange · StackExchangeP2P

TwentyNewsgroup

## MTEB
**Massive Text Embedding Benchmark**

8 Tasks

58 Datasets

## Bitext Mining

BUCC · Tatoeba

## Retrieval

Beir

ArguAna · ClimateFEVER · DBPedia

CQADupstackRetrieval · FEVER · FiQA2018

HotpotQA · MSMARCO · NFCorpus · NQ · Quora

SCIDOCS · SciFact · Touche2020 · TRECCOVID

## STS

BIOSESS · SICK-R

STS11 · STS12 · STS13

STS14 · STS15 · STS16

STS17 · STS22 · STSB

## Summarization

SummEval

## Classification

AmazonCounterfactual · AmazonPolarity

AmazonReviews · Banking77 · Emotion

Imdb · MassiveIntent · MassiveScenario

MTOPDomain · MTOPIntent

ToxicConversations · TweetSentimentExtraction

## Pair Classification

SprintDuplicateQuestions · TwitterSemEval2015

TwitterURLCorpus

## Reranking

AskUbuntuDupQuestions · MindSmallReranking

SciDocsRR · StackOverFlowDupQuestions

Thakur et al. 2021. BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models.
Muennighoff et al. 2022. MTEB: Massive Text Embedding Benchmark.

# MTEB Leaderboard

| Rank (Bor… | Model | Zero-shot | Memory Us… | Number of P… | Embedding D… | Max Tokens |
|---|---|---|---|---|---|---|
| **1** | KaLM-Embedding-Gemma3-12B-2511 | 73% | 44884 | 11.8 | 3840 | 32768 |
| 2 | llama-embed-nemotron-8b | 99% | 28629 | 7.5 | 4096 | 32768 |
| 3 | Qwen3-Embedding-8B | 99% | 14433 | 7.6 | 4096 | 32768 |
| 4 | gemini-embedding-001 | 99% | | | 3072 | 2048 |
| 5 | Qwen3-Embedding-4B | 99% | 7671 | 4.0 | 2560 | 32768 |
| 6 | Octen-Embedding-8B | 99% | 14433 | 7.6 | 4096 | 32768 |
| 7 | Seed1.6-embedding-1215 | 89% | | | 2048 | 32768 |
| 8 | Qwen3-Embedding-0.6B | 99% | 1136 | 0.596 | 1024 | 32768 |
| 9 | gte-Qwen2-7B-instruct | ⚠ NA | 29040 | 7.6 | 3584 | 32768 |
| 10 | Linq-Embed-Mistral | 99% | 13563 | 7.1 | 4096 | 32768 |
| 11 | multilingual-e5-large-instruct | 99% | 1068 | 0.560 | 1024 | 514 |

https://huggingface.co/spaces/mteb/leaderboard

Muennighoff et al. 2022. MTEB: Massive Text Embedding Benchmark.

38

# MTEB Leaderboard

# MTEB Leaderboard



https://huggingface.co/spaces/mteb/leaderboard

Muennighoff et al. 2022. MTEB: Massive Text Embedding Benchmark.

39

# MTEB Leaderboard

https://huggingface.co/spaces/mteb/leaderboard

Muennighoff et al. 2022. MTEB: Massive Text Embedding Benchmark.

# Summary of Part 2

$x$: When did GPT 5.2 come out?

$D$

$y$: Dec, 2026

**Datastore**

**Retriever**

**LM**

✔ Types of retrievers

✔ Training

✔ Evaluations

- Different types pf retrievers
- Training with contrastive loss
- Common metrics: NDCG@10 … etc
- Performance v.s. cost trade off

# Key Factors & Design Choices

$x$ : When did GPT 5.2 come out?

$D$

$y$ : Dec, 2026

**Datastore**

**Retriever**

**LM**

- ✔ Architectures
- ✔ Training
- ✔ Inference

# How to Use Retrieval

# How to Use Retrieval

**Input Augmentation**

Input



LM

- Augment input of LMs
- Easy to apply (w/o training) & effective
- Difficulty of using many D

e.g., RAG

# How to Use Retrieval

## Input Augmentation

Input

LM

- Augment input of LMs
- Easy to apply (w/o training) & effective
- Difficulty of using many D

e.g., RAG

## Intermediate Fusion

Input

LM

- Modify LMs to incorporate D in intermediate layers
- Scalable to many passages
- Requires retraining

e.g., RETRO, InstructRETRO

42

# How to Use Retrieval

**Input Augmentation**

Input

LM

- Augment input of LMs
- Easy to apply (w/o training) & effective
- Difficulty of using many D

e.g., RAG

**Intermediate Fusion**

Input

LM

- Modify LMs to incorporate D in intermediate layers
- Scalable to many passages
- Requires retraining

e.g., RETRO, InstructRETRO

**Output Interpolation**

LM

Output

- Directly manipulate output token distributions
- No training required*
- Limited effectiveness on tasks

e.g., kNNLM

42

# How to Use Retrieval

**Input Augmentation**

Input

LM

- Augment input of LMs
- Easy to apply (w/o training) & effective
- Difficulty of using many D

e.g., RAG

**Intermediate Fusion**

Input

LM

- Modify LMs to incorporate D in intermediate layers
- Scalable to many passages
- Requires retraining

e.g., RETRO, InstructRETRO

**Output Interpolation**

LM

Output

- Directly manipulate output token distributions
- No training required*
- Limited effectiveness on tasks

e.g., kNNLM

# RAG (Lewis et al., 2020)



$x$ : When did GPT 5.2 come out?

**Datastore**

**Retriever**

**LM**

$$D \in \text{Top}_k \text{Sim}( \cdot \,|\, x)$$

| | |
|---|---|
| Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for | 0.9 |
| GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was | 0.7 |

# RAG (Lewis et al., 2020)



$x$ : When did GPT 5.2 come out?

**Datastore**

**Retriever**

$D$

**LM**

$$D \in \text{Top}_k \text{Sim}(\,\cdot\mid x)$$

| | |
|---|---|
| Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for | 0.9 |
| GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was | 0.7 |

# RAG (Lewis et al., 2020)

$x$ : When did GPT 5.2 come out?



**LM**

$$D \in \text{Top}_k \text{Sim}( \cdot \,|\, x)$$

Dec 11, 2025 - Open AIL
GPT-5.2 Instant is a fast,
capable workhorse for

GPT-5 (Wikipedia) Preceded in
the series by GPT-4, it was

Lewis et al. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks.

# RAG (Lewis et al., 2020)

$x$ : When did GPT 5.2 come out?

$D$

**LM**

$$D \in \text{Top}_k \text{Sim}( \cdot \,|\, x)$$

Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for

GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was

**Question:** When did GPT 5.2 come out?
**References**:

Lewis et al. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks.

# RAG (Lewis et al., 2020)

$x$ : When did GPT 5.2 come out?

$D$

**LM**

$$D \in \text{Top}_k \text{Sim}( \cdot \,|\, x)$$

**Question:** When did GPT 5.2 come out?
**References**:

> Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for

> GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was

Lewis et al. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks.

45

# RAG (Lewis et al., 2020)

$x$ : When did GPT 5.2 come out?

$D$

**LM**

$$D \in \mathrm{Top}_k \mathrm{Sim}(\,\cdot\,|x)$$

**Question:** When did GPT 5.2 come out?
**References**:

Dec 11, 2025 - Open AIL GPT-5.2 Instant is a fast, capable workhorse for

GPT-5 (Wikipedia) Preceded in the series by GPT-4, it was

$y$ : Dec, 2026

Lewis et al. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks.

# Training RAG

# Training RAG



## Independent training

- DPR (Karpukhin et al., 2020)

- DRQA (Chen et al., 2017)

# Training RAG



## Independent training
- DPR (Karpukhin et al., 2020)
- DRQA (Chen et al., 2017)

## Sequential training
- Evidentiality Generator (Asai et al., 2023)
- REPLUG (Shi et al., 2023)

# Training RAG



## Independent training

- DPR (Karpukhin et al., 2020)
- DRQA (Chen et al., 2017)

## Sequential training

- Evidentiality Generator (Asai et al., 2023)
- REPLUG (Shi et al., 2023)

## Joint training

- RAG (Lewis et al., 2021)
- REALM (Guu et al., 2021)

# End-to-end Training for RAG



$$p_{\text{RAG}} \approx \prod_i \sum_{z \in \text{top-k}(p(\cdot|x))} \colorbox{green}{$p_\eta(z\,|\,x)$}\colorbox{cyan}{$p_\theta(y_i\,|\,x, z, y_{1:i-1})$}$$

Retriever score      Generator score

Lewis et al. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks.

47

# End-to-end Training for RAG

Define "middle ear"**(x)**

Question Answering:
Question Query

The middle ear includes
the tympanic cavity and
the three ossicles. **(y)**

End-to-End Backprop through **q** and **p$_\theta$**

Question Answering:
Answer Generation

Query
Encoder

**Retriever p$_\eta$**
(Non-Parametric)

Document
Index

**Generator p$_\theta$**
(Parametric)

Barack Obama was
born in Hawaii.**(x)**

Fact Verification: Fact Query

**q(x)**

**d(z)**

$z_4$

$z_3$

$z_2$

supports **(y)**

Fact Verification:
Label Generation

**Margin-
alize**

Question Generation

**Now people often combine retrieval with off-the-shelf LMs**

Widely referred to as *RAG*

$$p_{\text{RAG}} \approx \prod_i \sum_{z \in \text{top}-\text{k}(p(\cdot|x))} p_\eta(z\,|\,x) p_\theta(y_i\,|\,x, z, y_{1:i-1})$$

Retriever score    Generator score

Lewis et al. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks.

# Effectiveness of In-context RAG

RAG constantly gives performance improvements in long-tail



Mallen*, Asai* (contributed equally) et al. 2023. When Not to Trust Language Models:
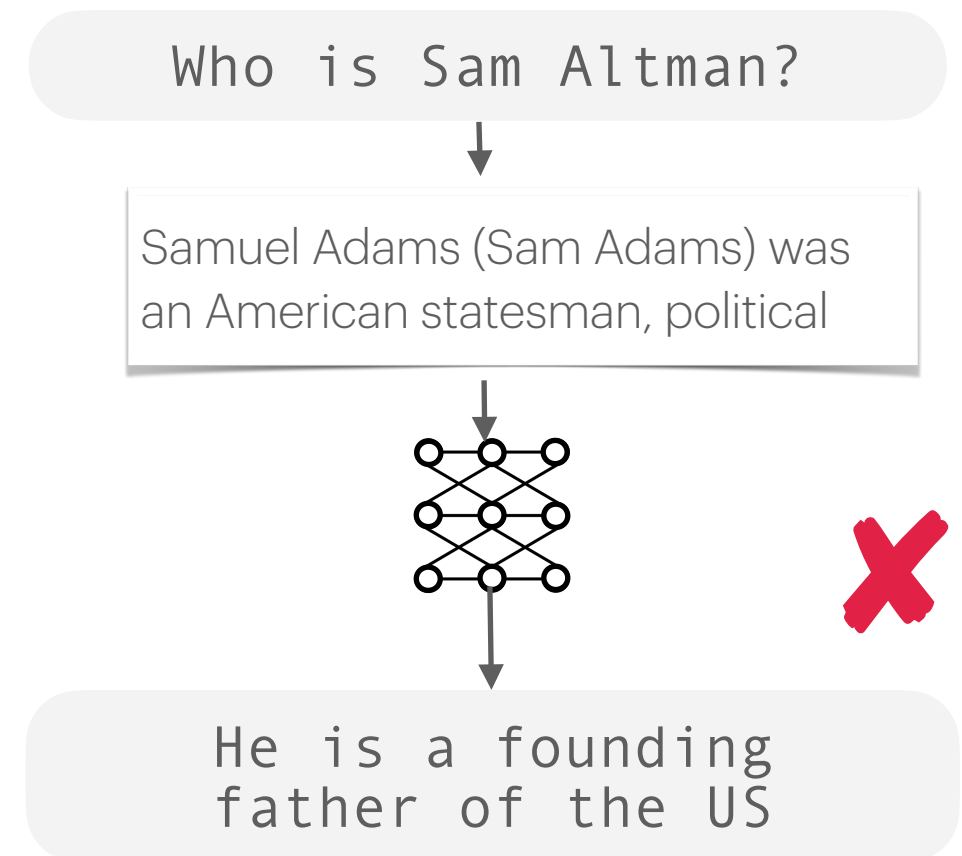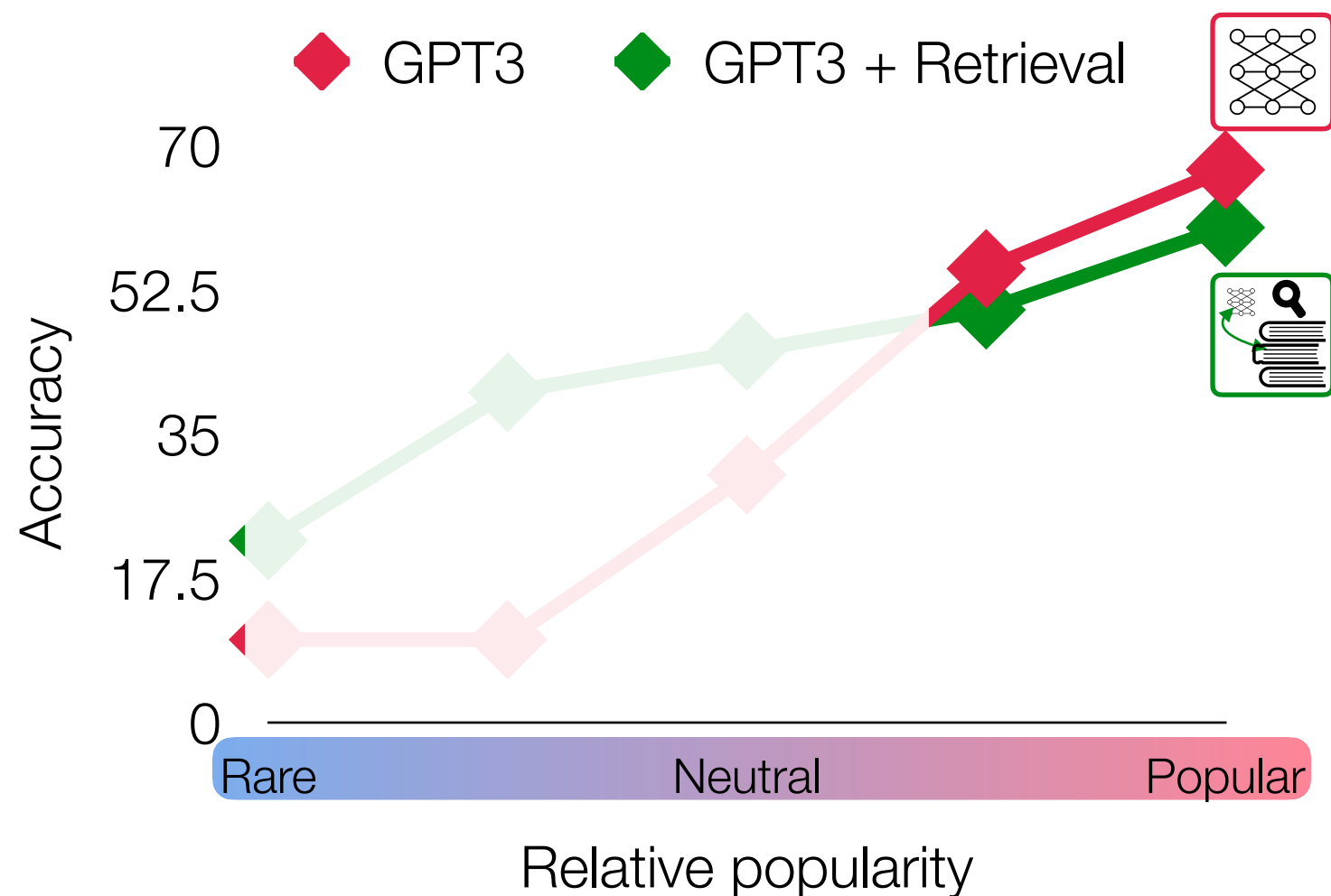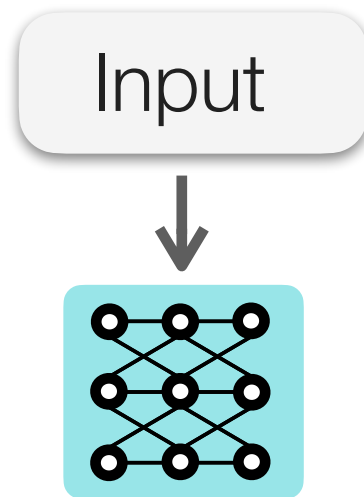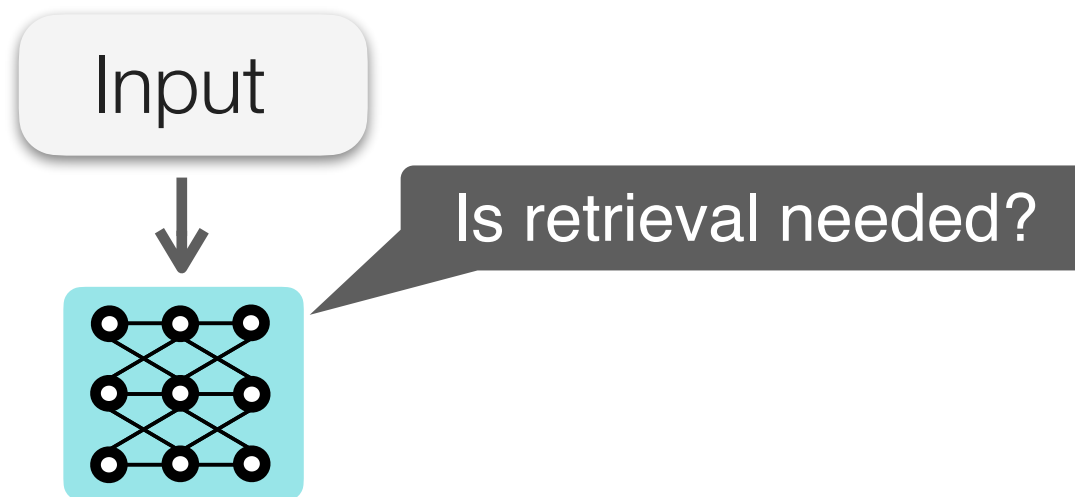Investigating Effectiveness of Parametric and Non-Parametric Memories

# Limitations of In-context RAG



Mallen*, Asai* (contributed equally) et al. 2023. When Not to Trust Language Models:
Investigating Effectiveness of Parametric and Non-Parametric Memories

# Limitations of In-context RAG
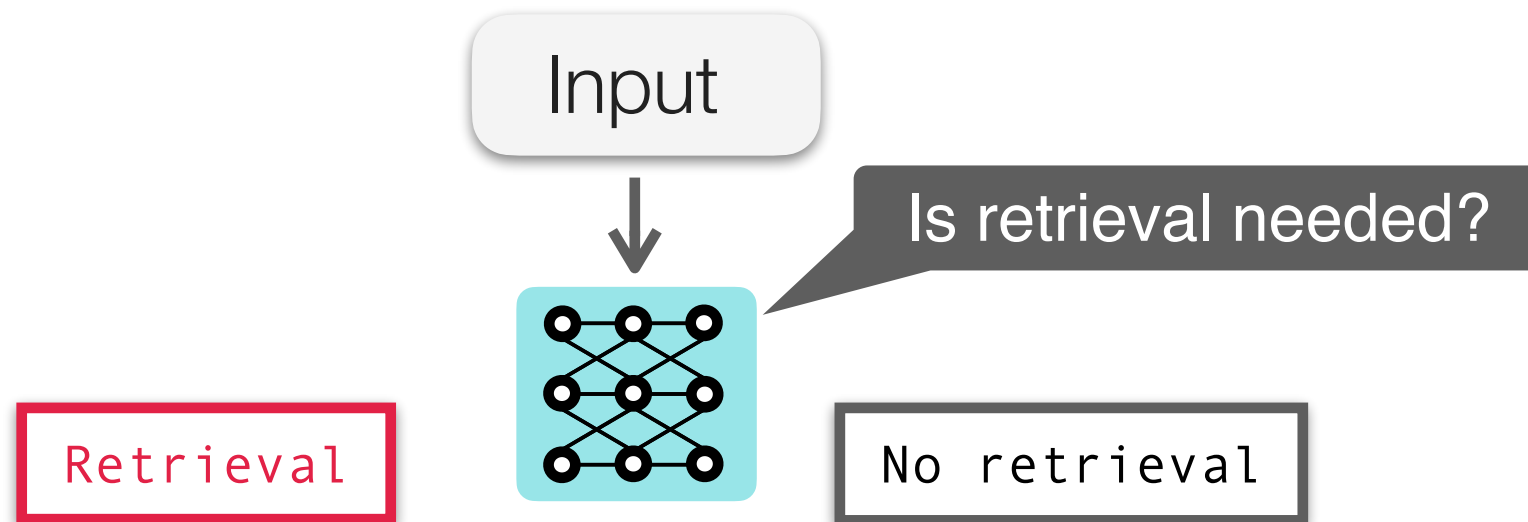


GPT3 GPT3 + Retrieval

Who is Sam Altman?

Samuel Adams (Sam Adams) was an American statesman, political

Mallen*, Asai* (contributed equally) et al. 2023. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories

# Limitations of In-context RAG



Mallen*, Asai* (contributed equally) et al. 2023. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories

# Limitations of In-context RAG



Mallen*, Asai* (contributed equally) et al. 2023. When Not to Trust Language Models:
Investigating Effectiveness of Parametric and Non-Parametric Memories

# Limitations of In-context RAG

**Inflexibility and lack of robustness to unhelpful docs**



GPT3 ◆     GPT3 + Retrieval ◆

Accuracy: 70, 52.5, 35, 17.5, 0

Relative popularity: Rare — Neutral — Popular

Who is Sam Altman?

Samuel Adams (Sam Adams) was an American statesman, political

He is a founding father of the US

Mallen*, Asai* (contributed equally) et al. 2023. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories

# Self-RAG: Adaptive Retrieval

Input

Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

# Self-RAG: Adaptive Retrieval


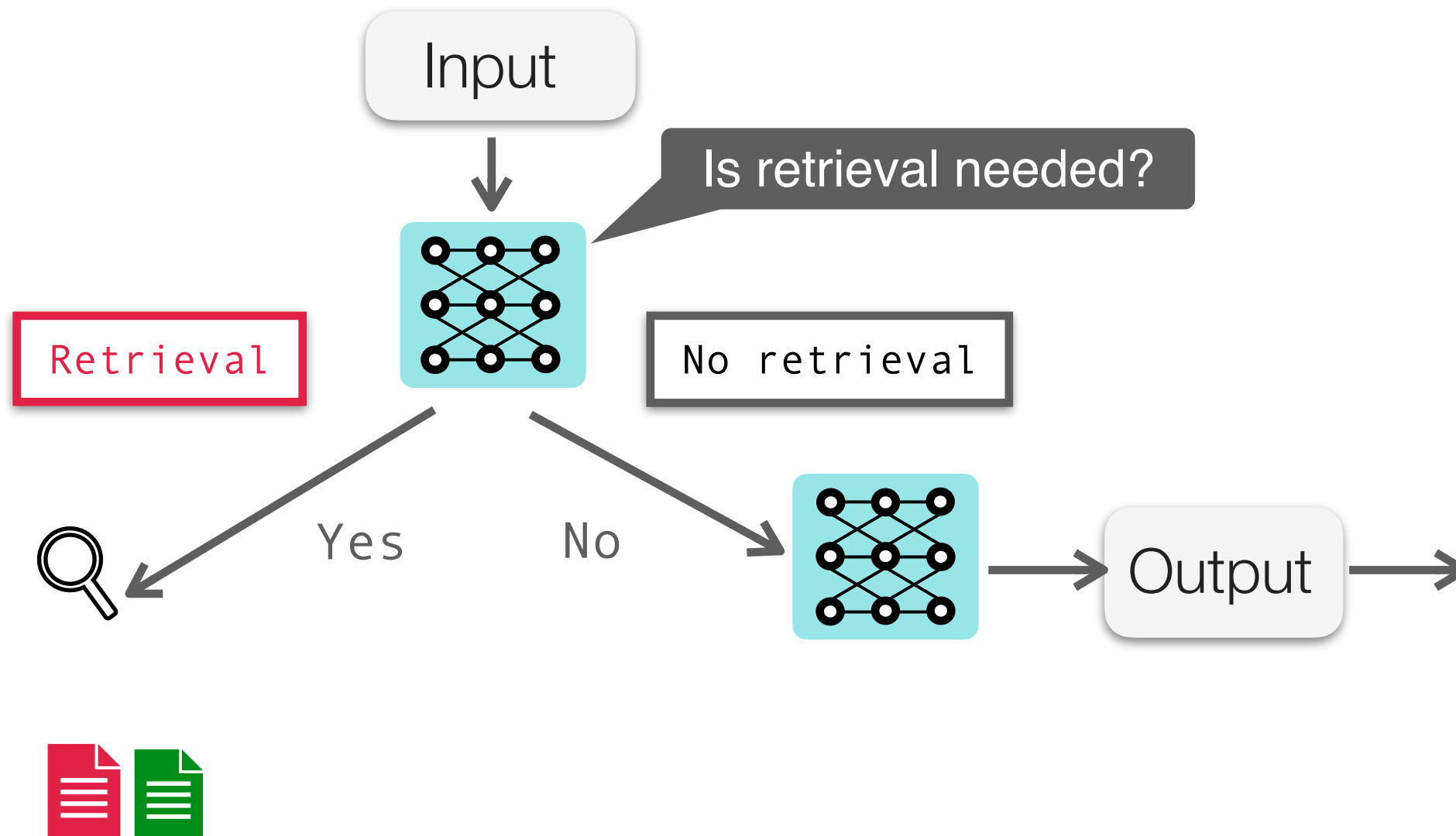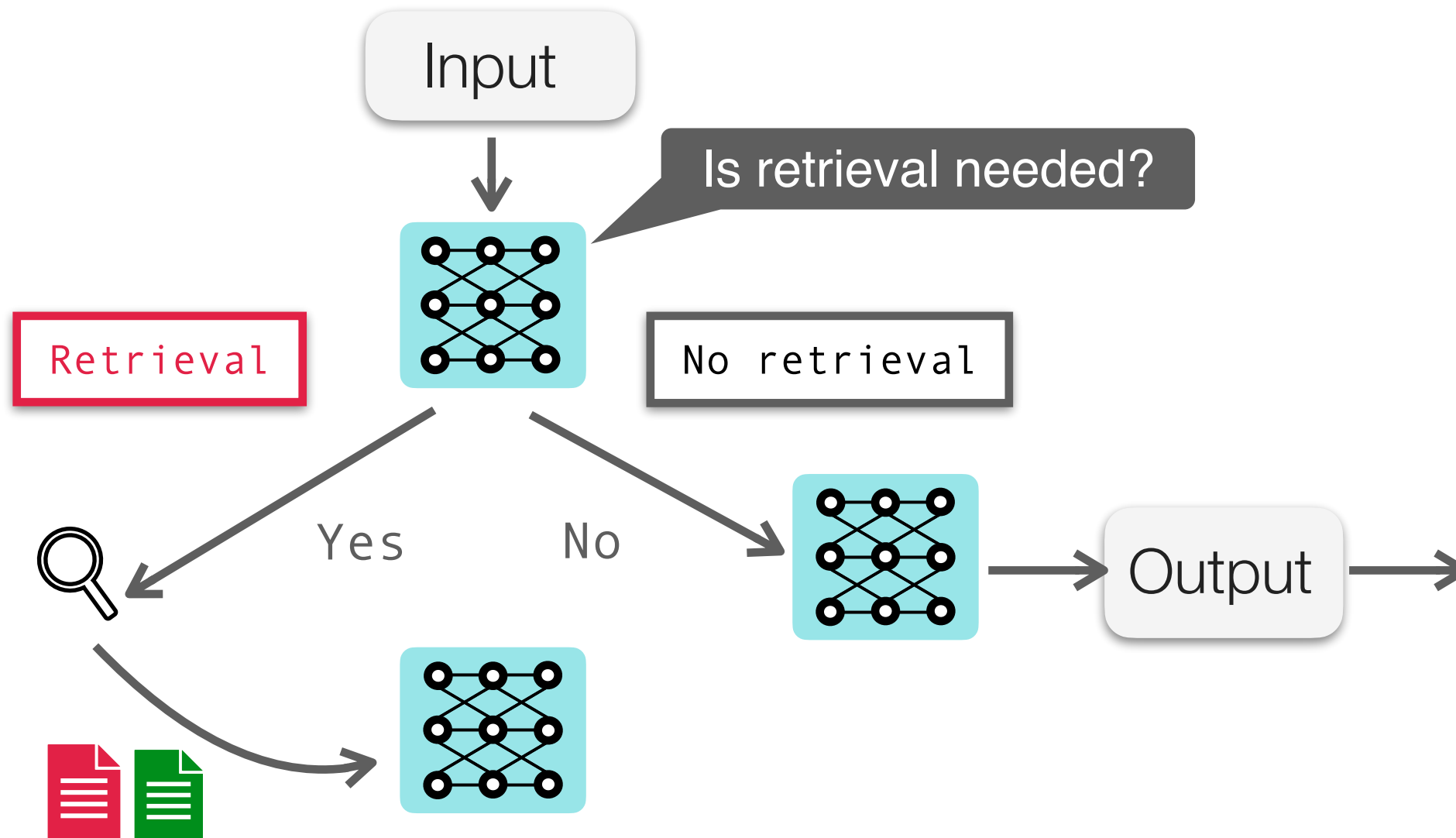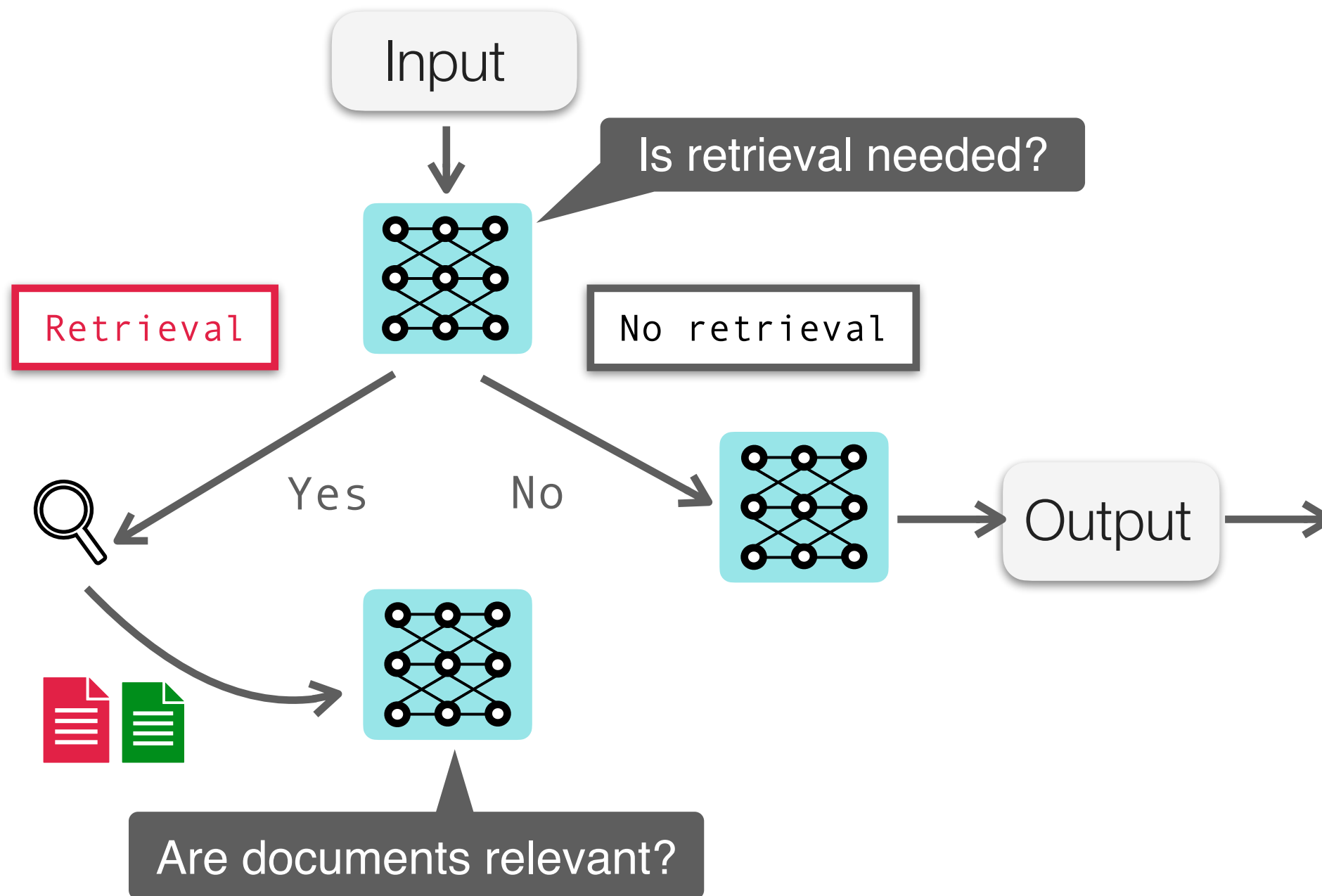
Input

Is retrieval needed?

Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

51

# Self-RAG: Adaptive Retrieval



Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

51

# Self-RAG: Adaptive Retrieval



Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.
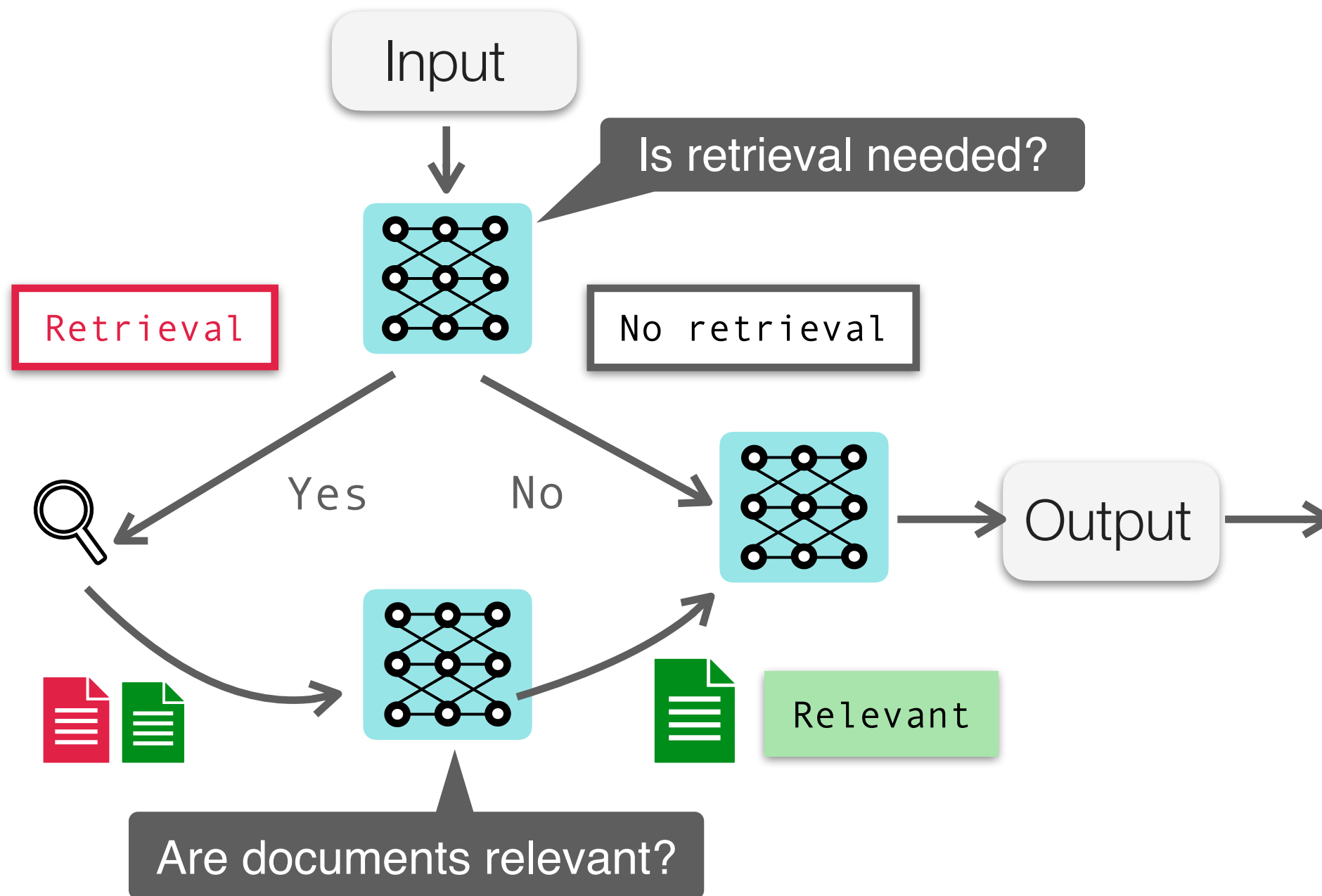
# Self-RAG: Adaptive Retrieval



Input

Is retrieval needed?

Retrieval

No retrieval

No

Output

Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

# Self-RAG: Adaptive Retrieval



Input

Is retrieval needed?
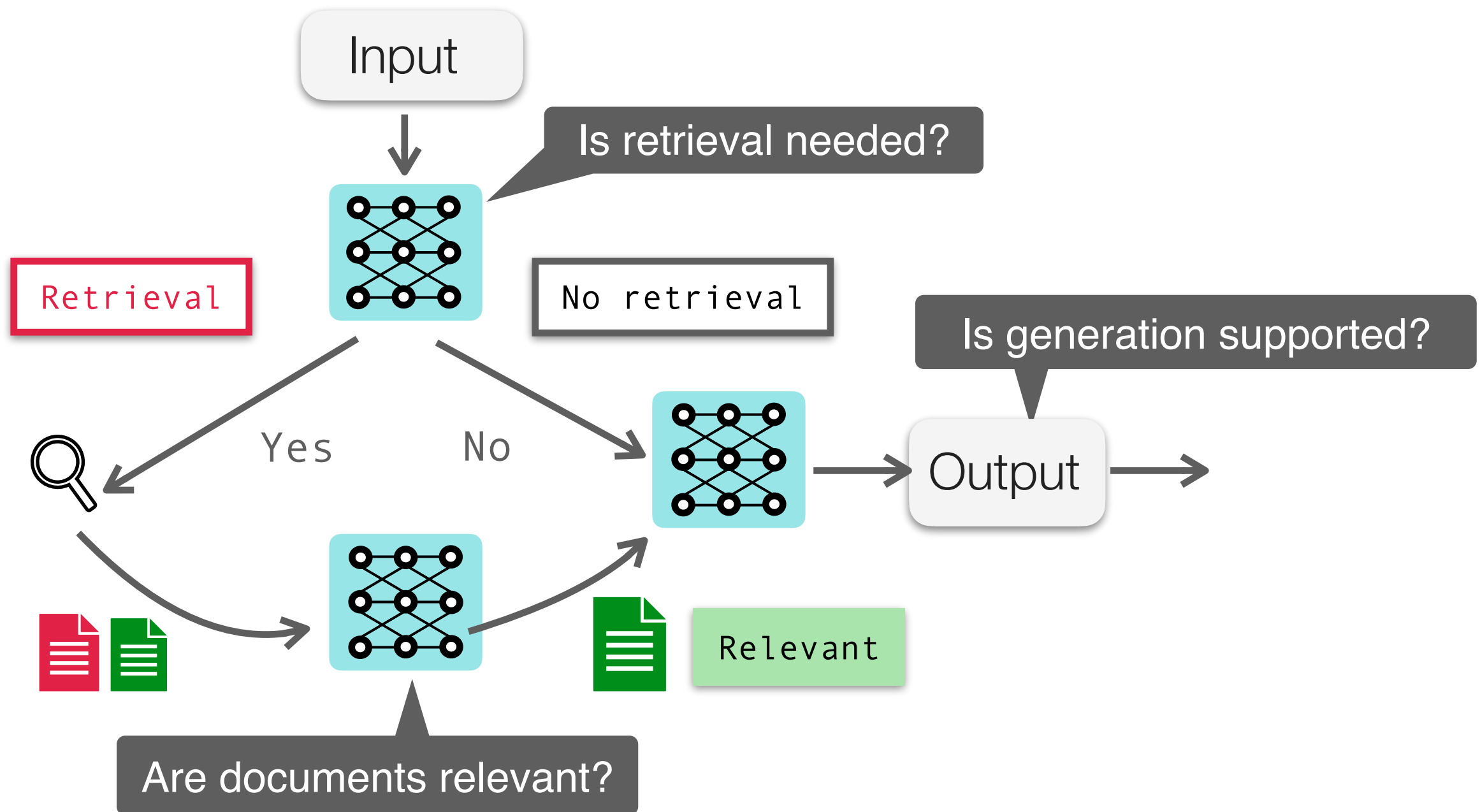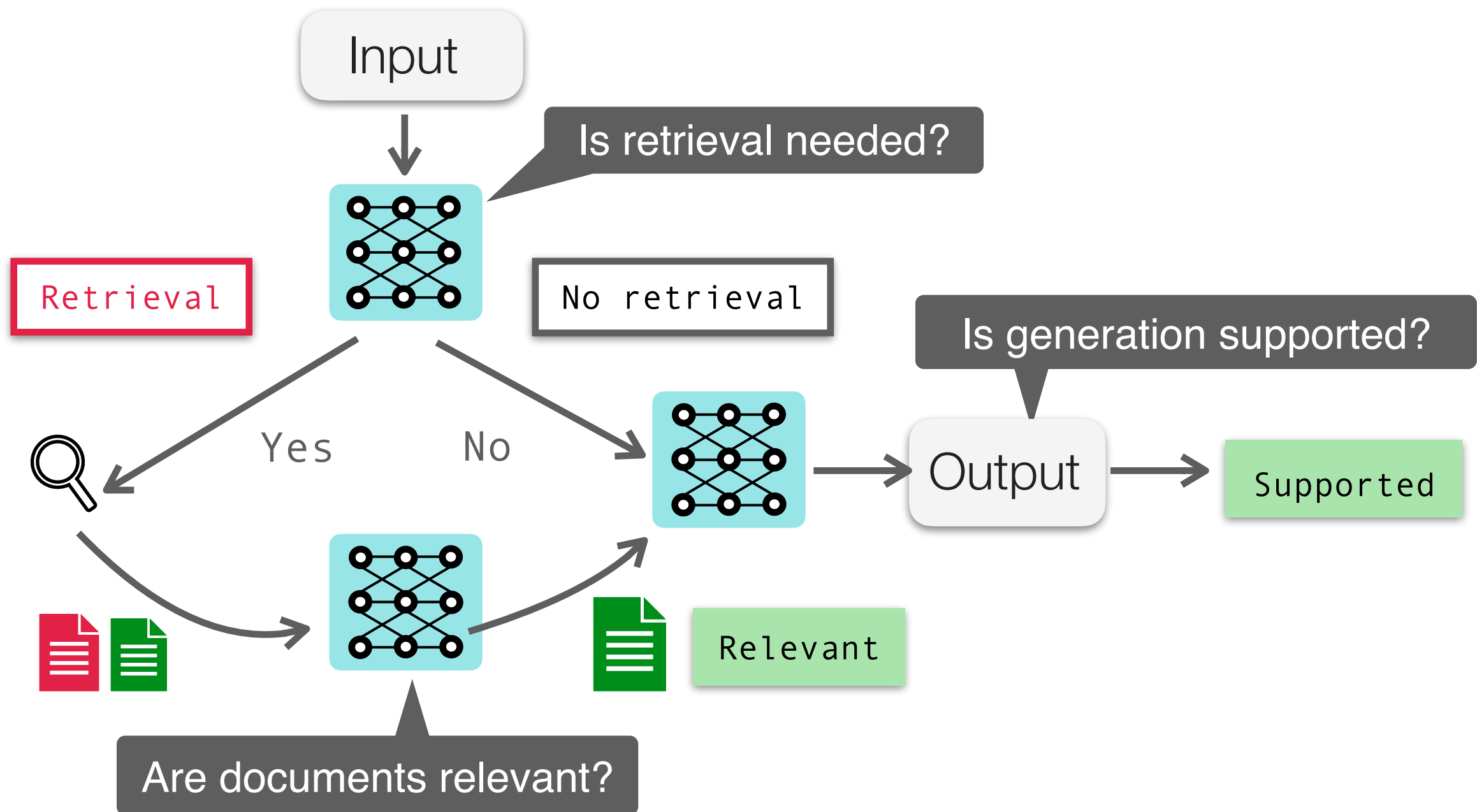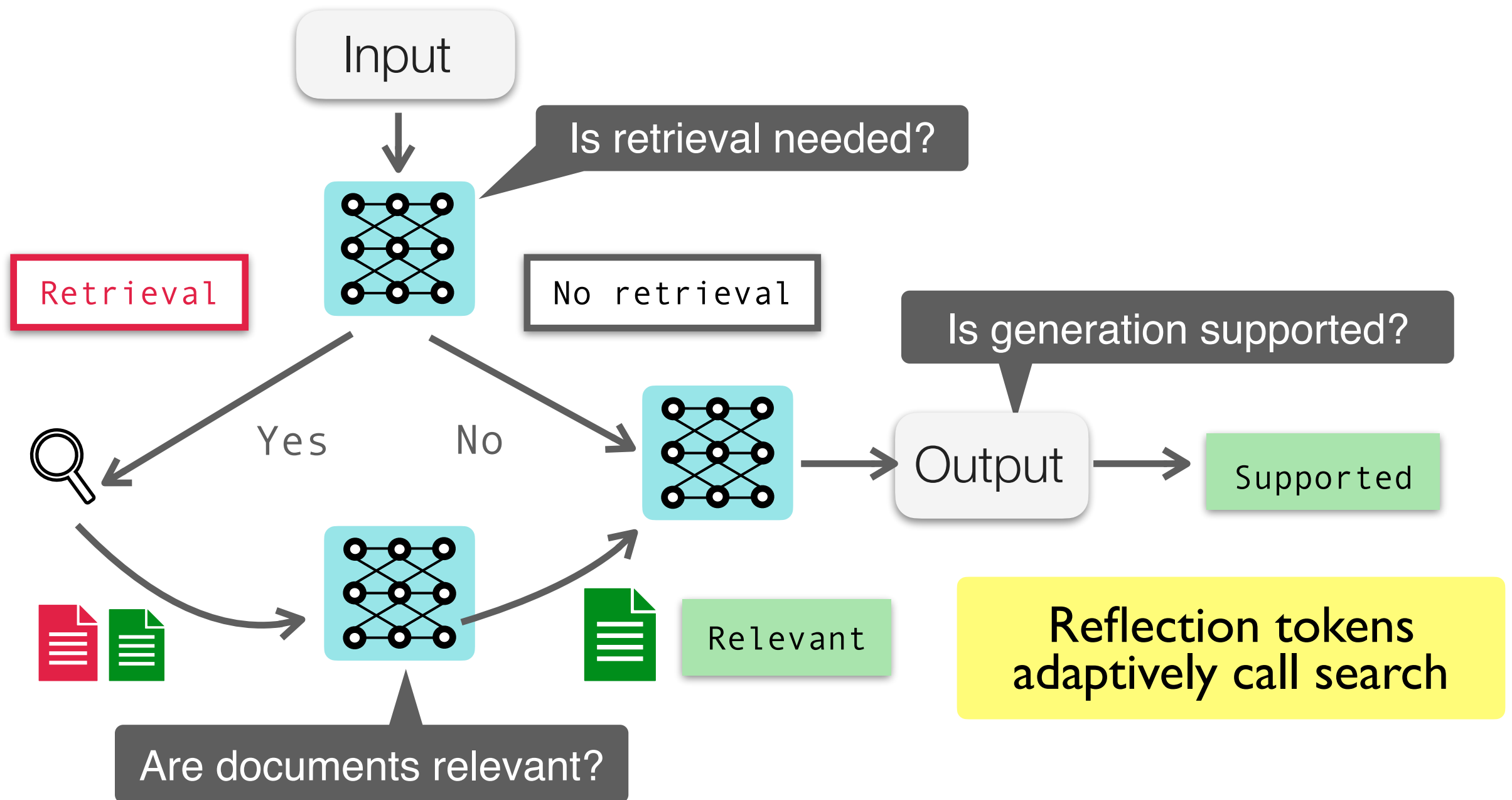
Retrieval

No retrieval

Yes

No

Output

Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

# Self-RAG: Adaptive Retrieval



Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate,
and Critique through Self-Reflection.

# Self-RAG: Adaptive Retrieval



Input

Is retrieval needed?

Retrieval

No retrieval

Yes

No

Output

Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

51

# Self-RAG: Adaptive Retrieval



Input

Is retrieval needed?

Retrieval

No retrieval

Yes    No

Output

Are documents relevant?

Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

# Self-RAG: Adaptive Retrieval

Input

Is retrieval needed?

Retrieval

No retrieval

Yes    No

Output

Relevant

Are documents relevant?

Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

# Self-RAG: Adaptive Retrieval



Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

# Self-RAG: Adaptive Retrieval



Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.

51

# Self-RAG: Adaptive Retrieval



Asai et al. 2024..Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection.
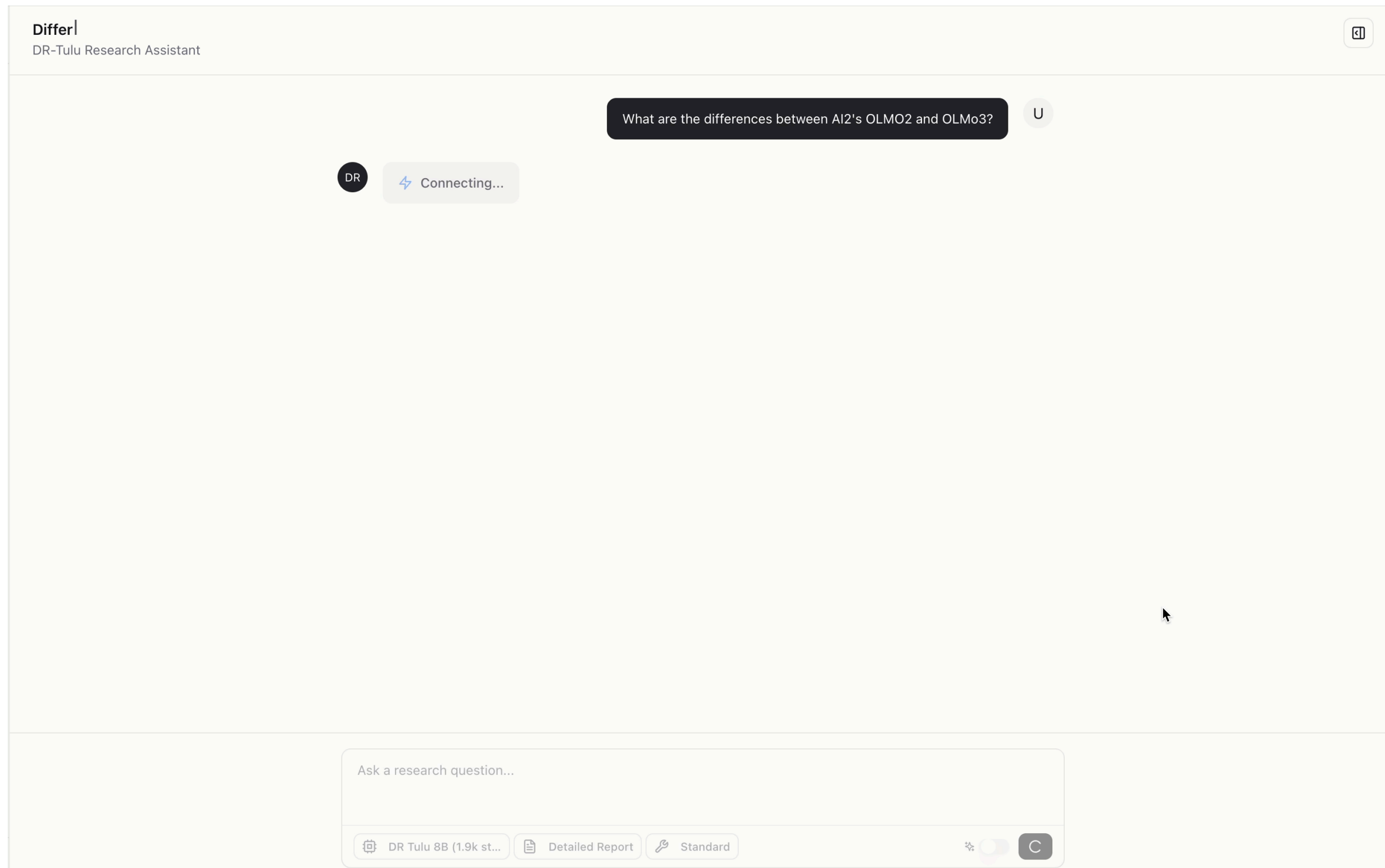
51

# Tool-augmented LMs

- Training LMs to *adaptively* and *iteratively* use external tools at inference time

- LMs can use diverse set of tools, not just retrieval

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.
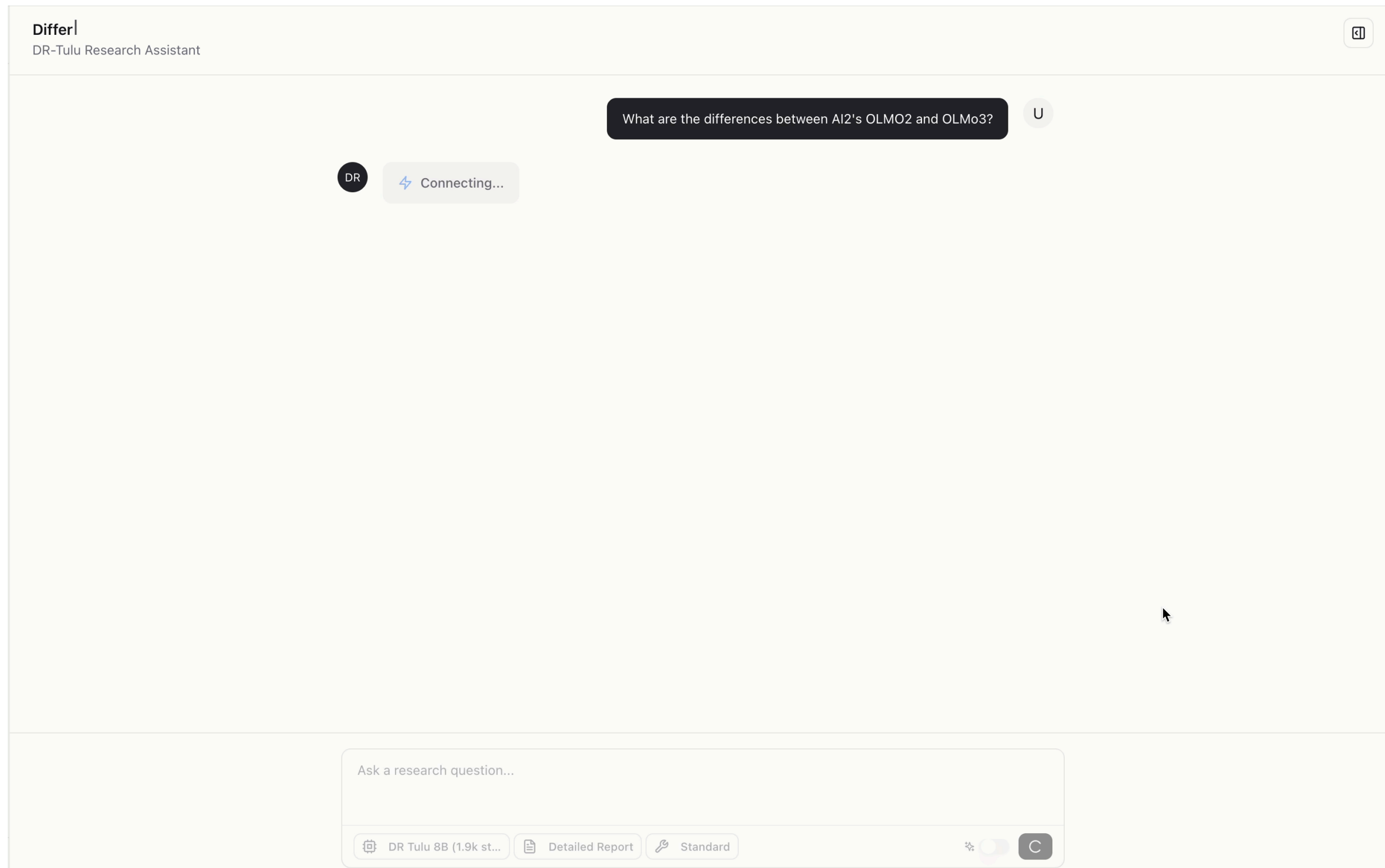
Schick et al. 2023. Toolformer: Language Models Can Teach Themselves to Use Tools.

# Tool-augmented LMs

- Training LMs to *adaptively* and *iteratively* use external tools at inference time

- LMs can use diverse set of tools, not just retrieval

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

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

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

Schick et al. 2023. Toolformer: Language Models Can Teach Themselves to Use Tools.

52

# Deep Research (DR) Agents

https://www.dr-tulu.org/chat

Shao*, Asai* et al. 2025. DR Tulu: Reinforcement Learning with Evolving Rubrics for Deep Research.

# Deep Research (DR) Agents

53

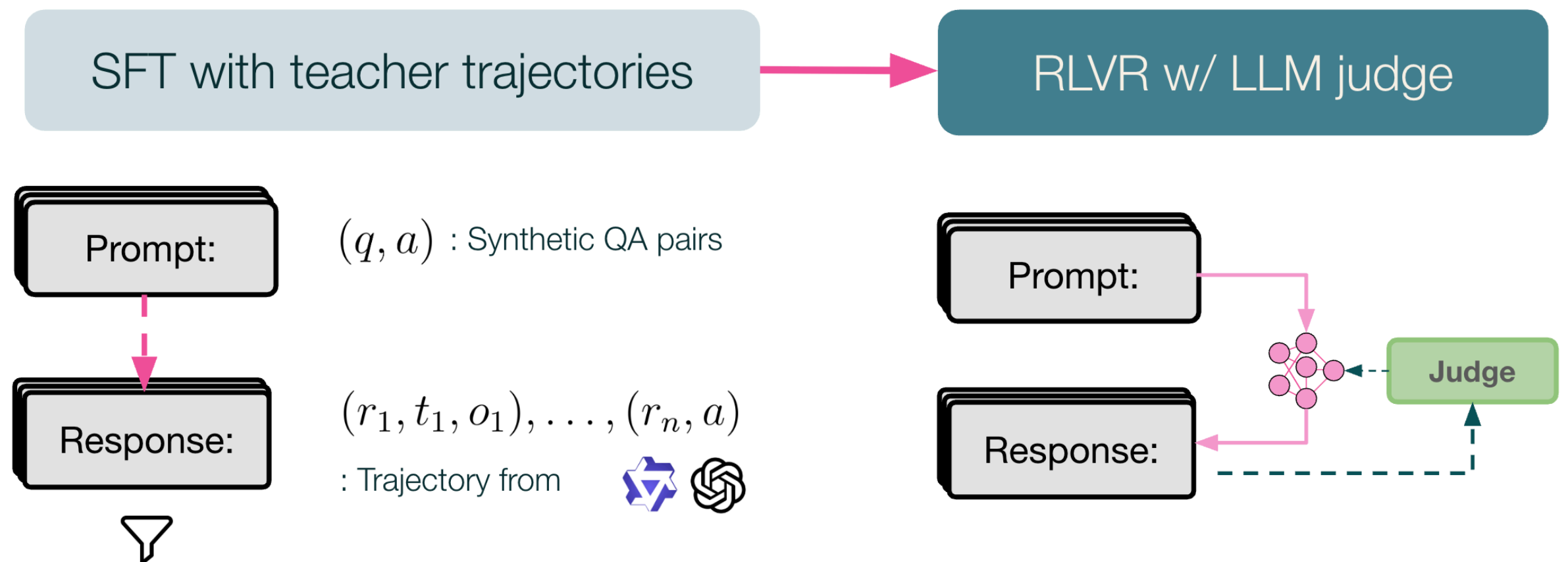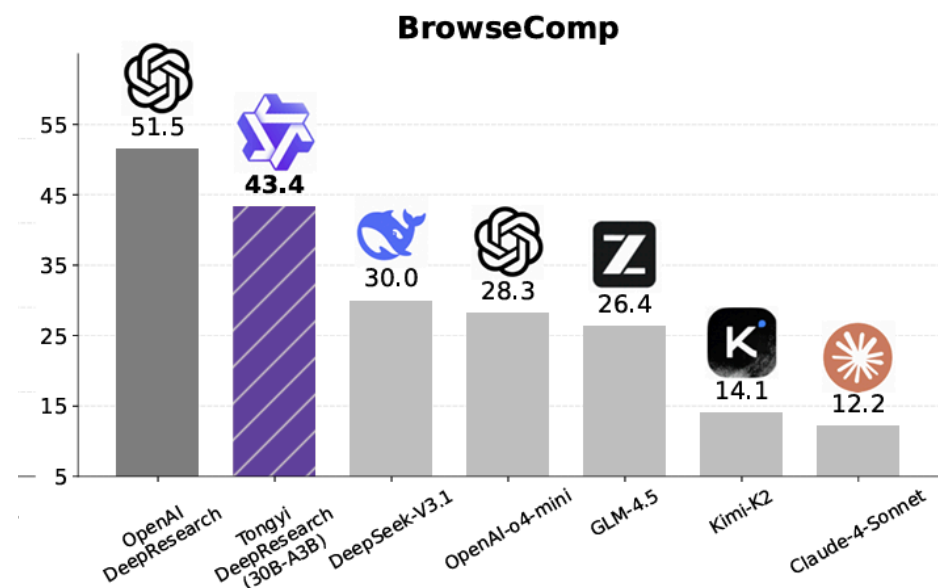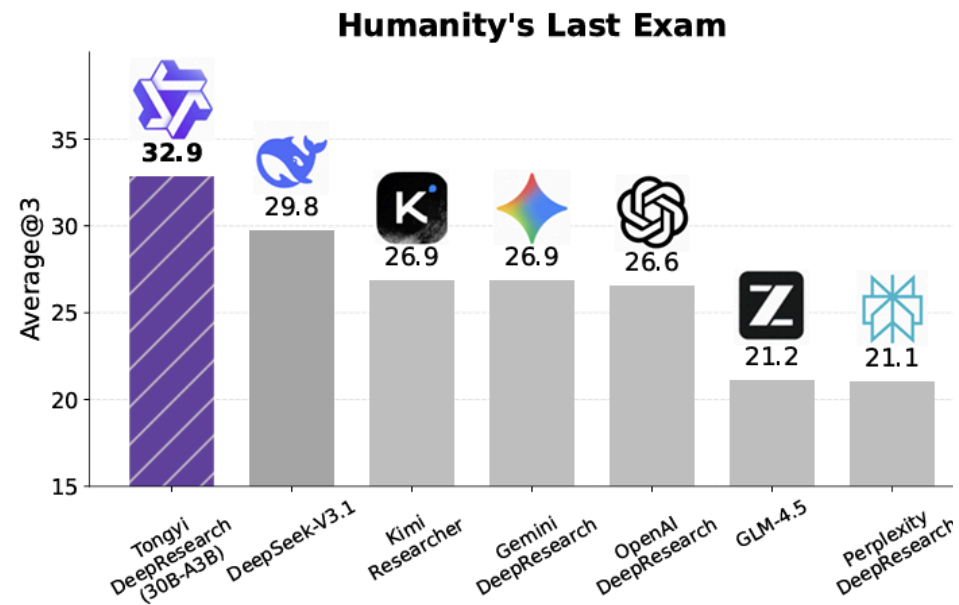Shao*, Asai* et al. 2025. DR Tulu: Reinforcement Learning with Evolving Rubrics for Deep Research.

# Training for DR Agents

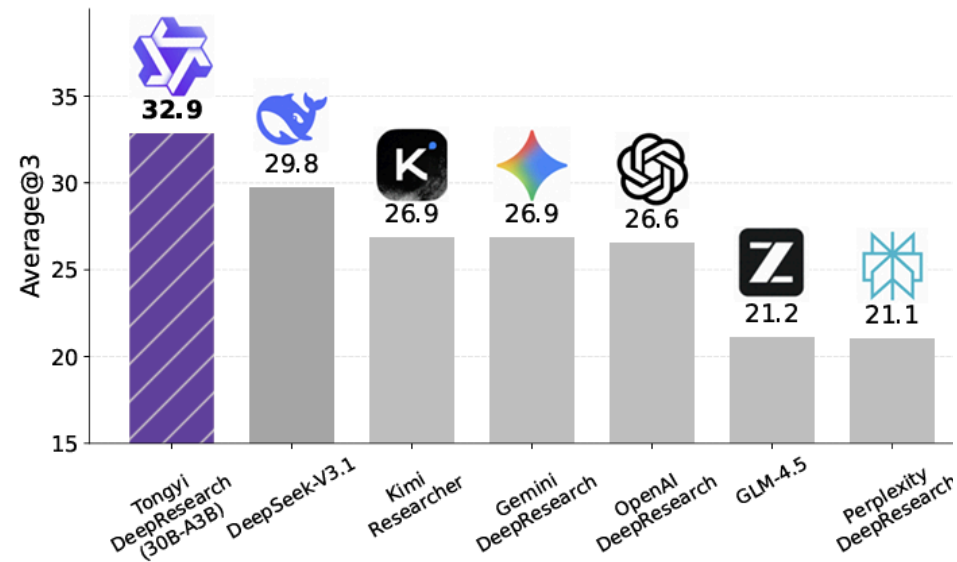- Large-scale SFT (w. Rejection sampling) followed by RLVR using answer matching as reward



SFT with teacher trajectories → RLVR w/ LLM judge
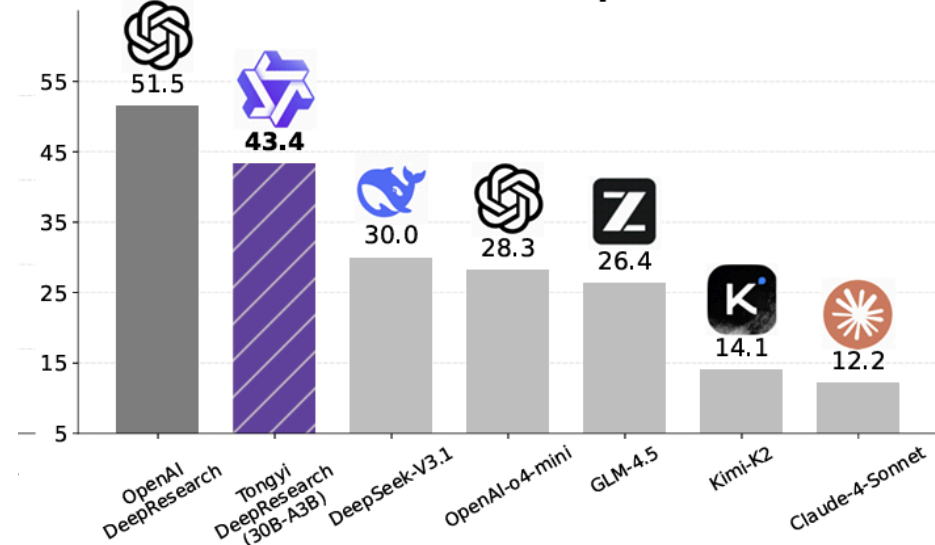
$(q, a)$ : Synthetic QA pairs

$(r_1, t_1, o_1), \ldots, (r_n, a)$
: Trajectory from

Prompt:
Response:
Judge

Tongyi Research. 2025. Tongyi Deep Research Technical Report.

# Training for DR Agents



Tongyi Research. 2025. Tongyi Deep Research Technical Report.
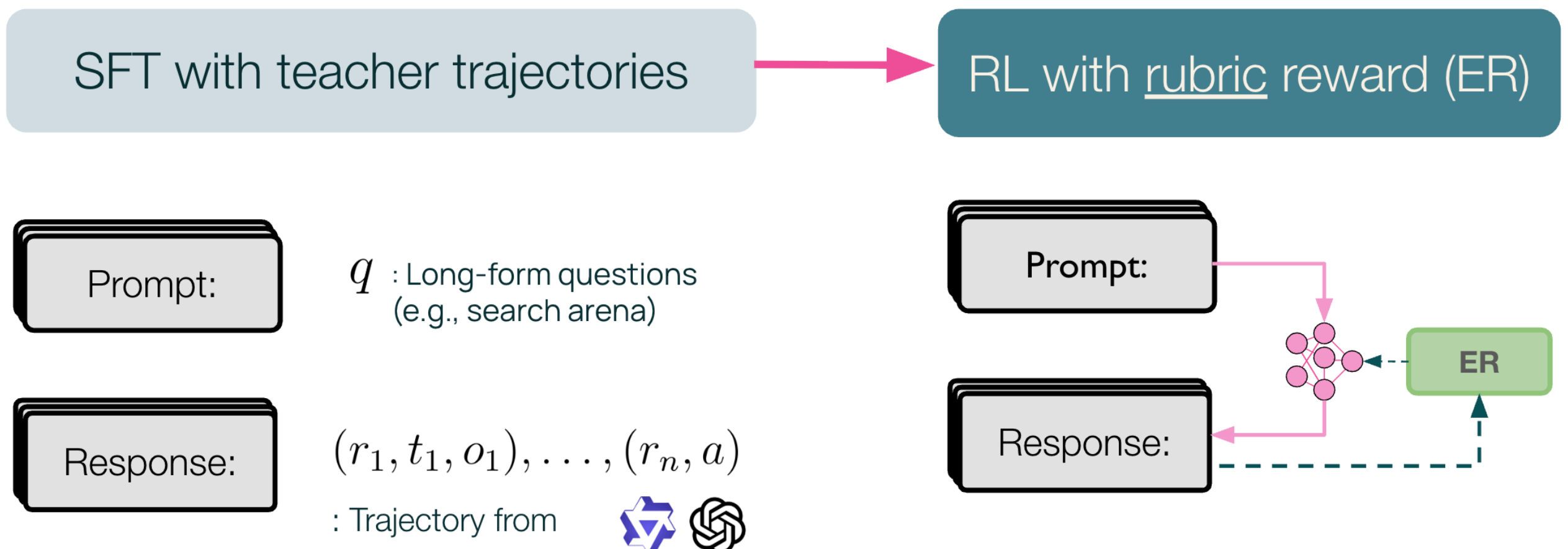
# Training for DR Agents

## Humanity's Last Exam



## BrowseComp



## Deep Research Bench (Long-form DR)

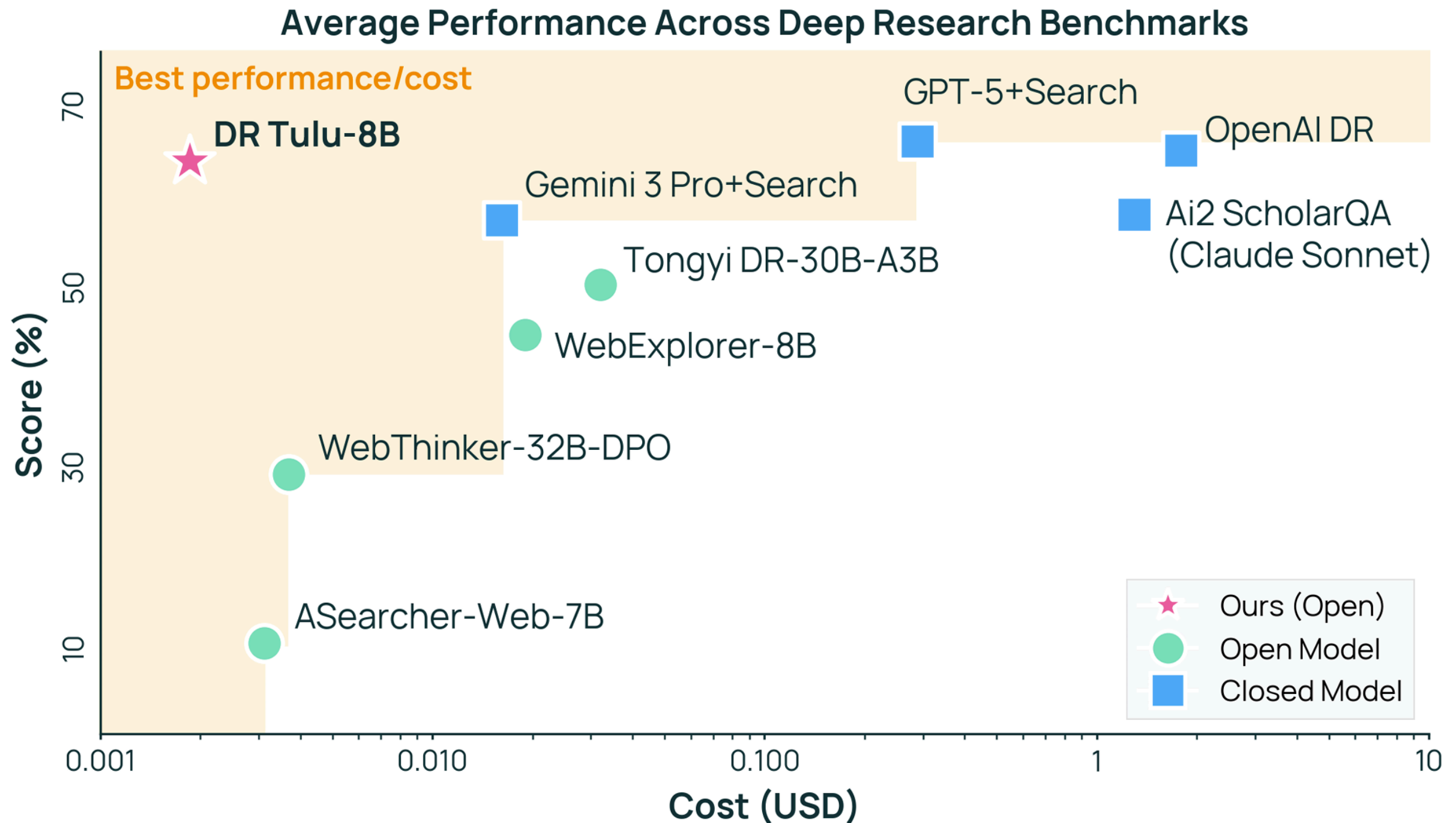| Rank | model | overall |
|------|-------|---------|
| 1 🥇 | 🚀 cellcog | 54.54 |
| 2 🥈 | 🚀 Qianfan-DeepResearch Pro | 54.22 |
| 3 🥉 | 🚀 Qianfan-DeepResearch | 53.02 |
| 4 | 🚀 tavily-research | 52.44 |
| 5 | 🚀 thinkdepthai-deepresearch | 52.43 |
| 6 | 🚀 salesforce-air-deep-research | 50.65 |
| 7 | 🚀 langchain-open-deep-research(GPT-5,with gensee search) | 50.6 |
| 8 | 🚀 gemini-2.5-pro-deepresearch | 49.71 |
| 9 | 🚀 langchain-open-deep-research(GPT-5,with Tavily) | 49.33 |
| 10 | 🚀 openai-deepresearch | 46.45 |
| ... | | |
| 17 | 🚀 tongyi-deepresearch-30B-A3B | 40.46 |

Tongyi Research. 2025. Tongyi Deep Research Technical Report.

# Training for DR Agents with Rubrics

- Long-form responses are "hard-to-verify"

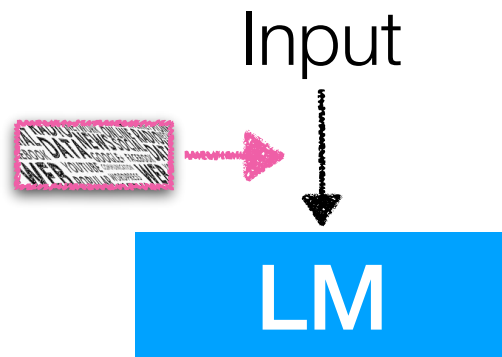- Rubric-reward based RL for DR agent training



SFT with teacher trajectories

Prompt:

$q$ : Long-form questions (e.g., search arena)

Response:

$(r_1, t_1, o_1), \ldots, (r_n, a)$

: Trajectory from

RL with <u>rubric</u> reward (ER)

Prompt:

ER

Response:

Shao*, Asai* et al. 2025. DR Tulu: Reinforcement Learning with Evolving Rubrics for Deep Research.

# Training for DR Agents with Rubrics



**Average Performance Across Deep Research Benchmarks**

Best performance/cost

DR Tulu-8B

GPT-5+Search

OpenAI DR

Gemini 3 Pro+Search

Ai2 ScholarQA (Claude Sonnet)

Tongyi DR-30B-A3B

WebExplorer-8B

WebThinker-32B-DPO

ASearcher-Web-7B

Legend:
- ★ Ours (Open)
- ● Open Model
- ■ Closed Model

Score (%) vs Cost (USD)

# How to Use Retrieval

**Input Augmentation**

Input



**LM**

- Augment input of LMs
- Easy to apply (w/o training) & effective
- Difficulty of using many D

e.g., RAG

☹ Not scalable to many documents (needs context engineering)

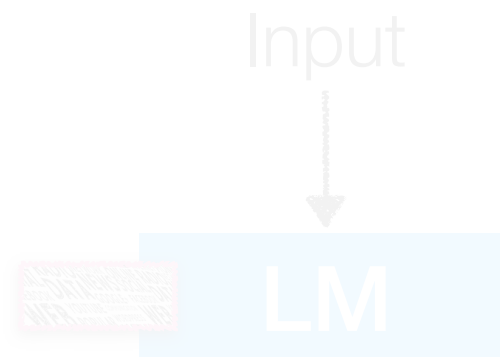☹ Not strictly grounded

# How to Use Retrieval

**Input Augmentation**

Input

LM

- Augment input of LMs
- Easy to apply (w/o training) & effective
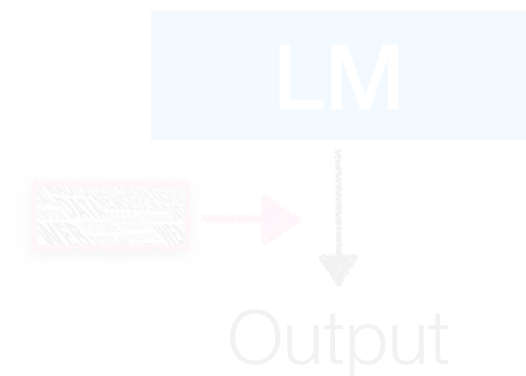- Difficulty of using many D

e.g., RAG

**Intermediate Fusion**

Input

LM

- Modify LMs to incorporate D in intermediate layers
- Scalable to many passages
- Requires retraining

e.g., RETRO, InstructRETRO
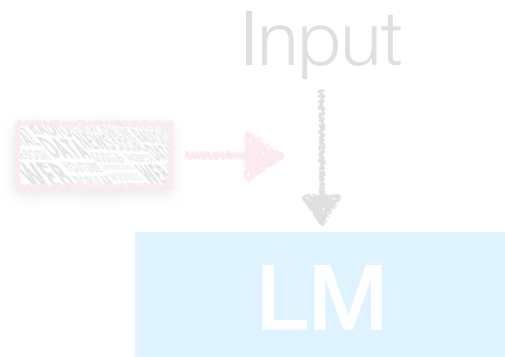
**Output Interpolation**

LM

Output

- Directly manipulate output token distributions
- No training required*
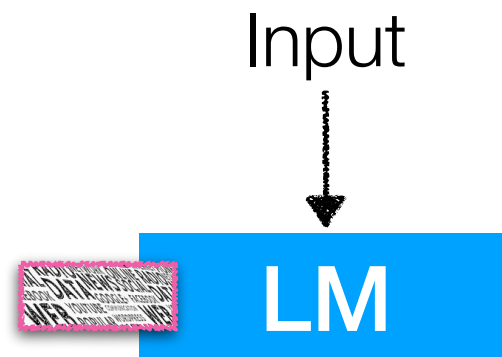- Limited effectiveness on tasks

e.g., kNNLM

# How to Use Retrieval



**Input Augmentation**

Input

LM

- Augment input of LMs
- Easy to apply (w/o training) & effective
- Difficulty of using many D

e.g., RAG

**Intermediate Fusion**

Input

LM
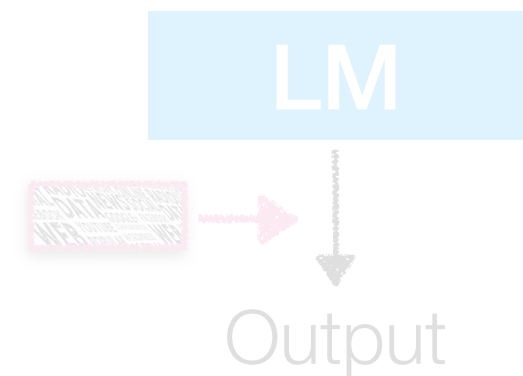
- Modify LMs to incorporate D in intermediate layers
- Scalable to many passages
- Requires retraining

e.g., RETRO, InstructRETRO

**Output Interpolation**
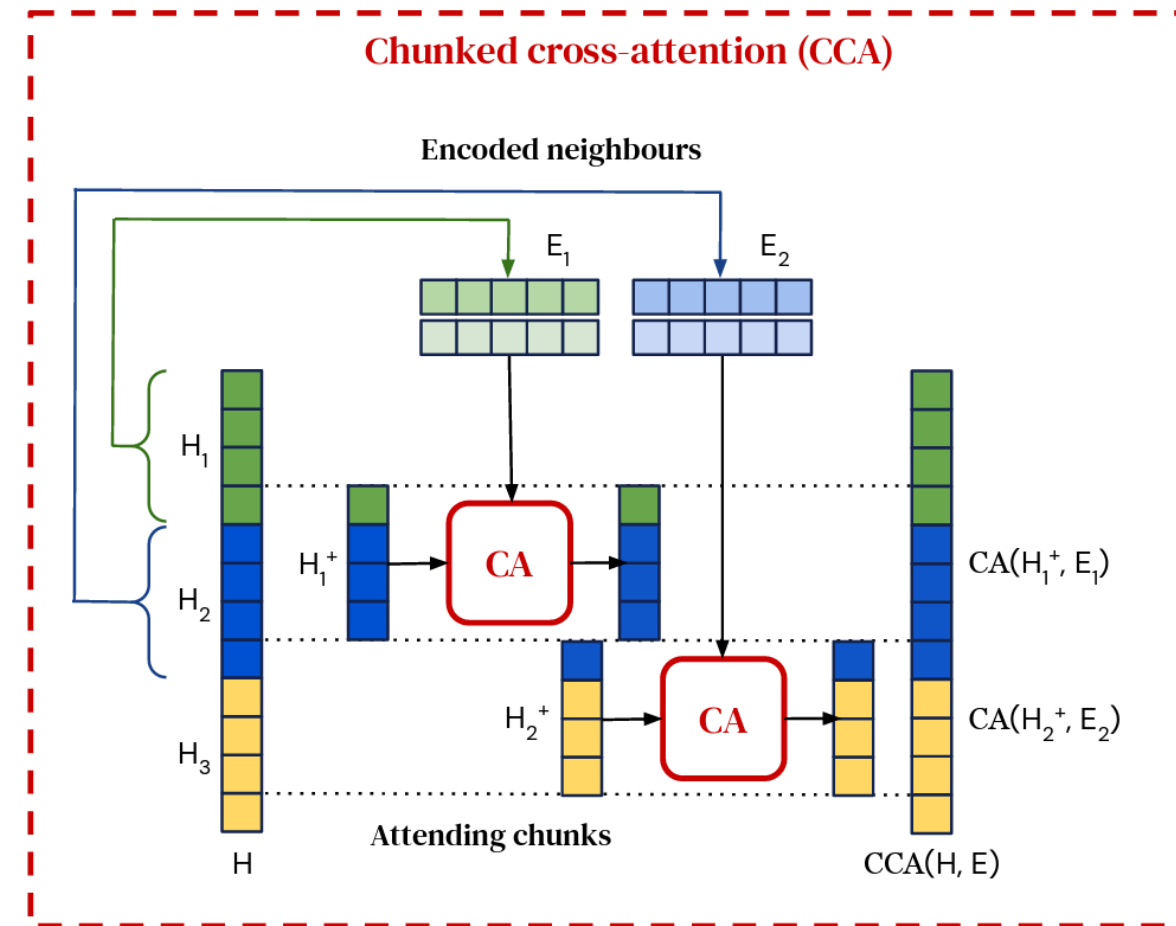
LM

Output

- Directly manipulate output token distributions
- No training required*
- Limited effectiveness on tasks

e.g., kNNLM

59

# RETRO



Borgeaud et al. 2022. Improving language models by retrieving from trillions of tokens.

60

# How to Use Retrieval

## Input Augmentation

Input

LM

- Augment input of LMs
- Easy to apply (w/o training) & effective
- Difficulty of using many D

e.g., RAG

## Intermediate Fusion

Input

LM

- Modify LMs to incorporate D in intermediate layers
- Scalable to many passages
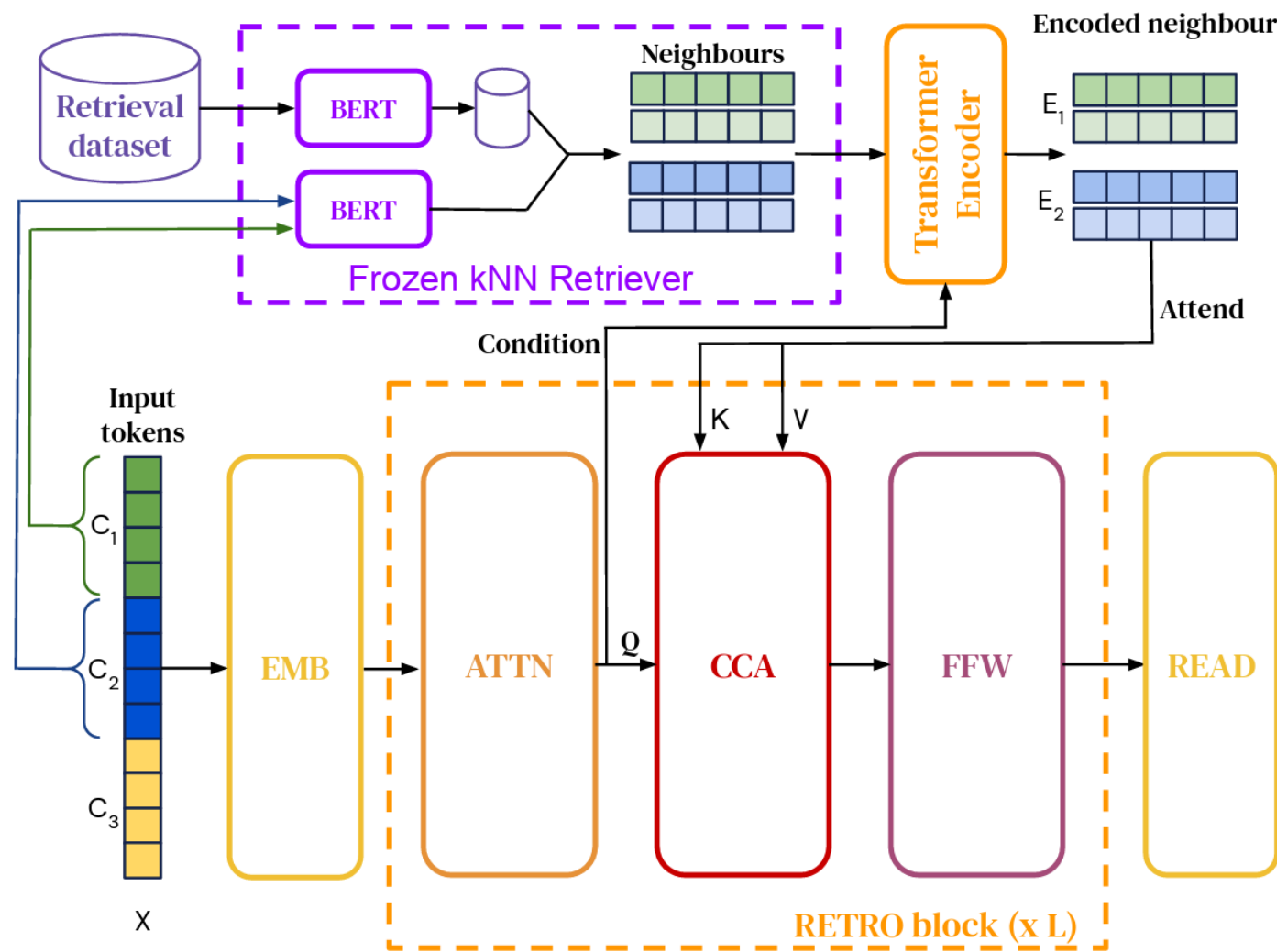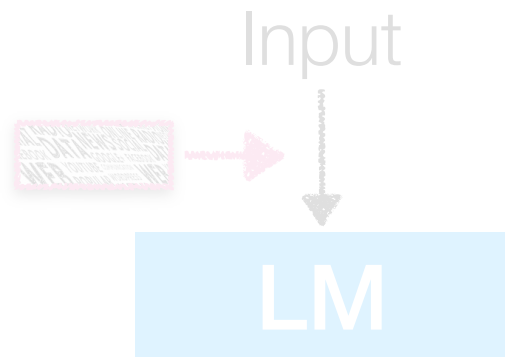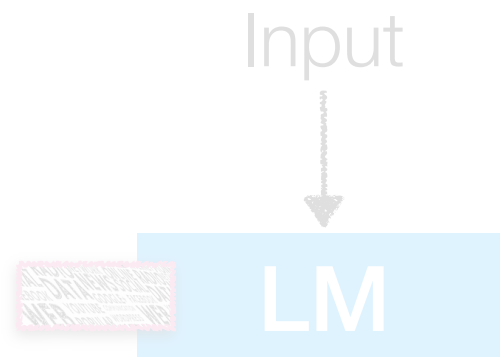- Requires retraining

e.g., RETRO, InstructRETRO

## Output Interpolation

LM

Output

- Directly manipulate output token distributions
- No training required*
- Limited effectiveness on tasks

e.g., kNNLM

# kNN-LM



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

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

# kNN-LM



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

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

# kNN-LM



Nonparametric distribution

Parametric distribution

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

Khandelwal et al. 2020. Generalization through Memorization: Nearest Neighbor Language Models.
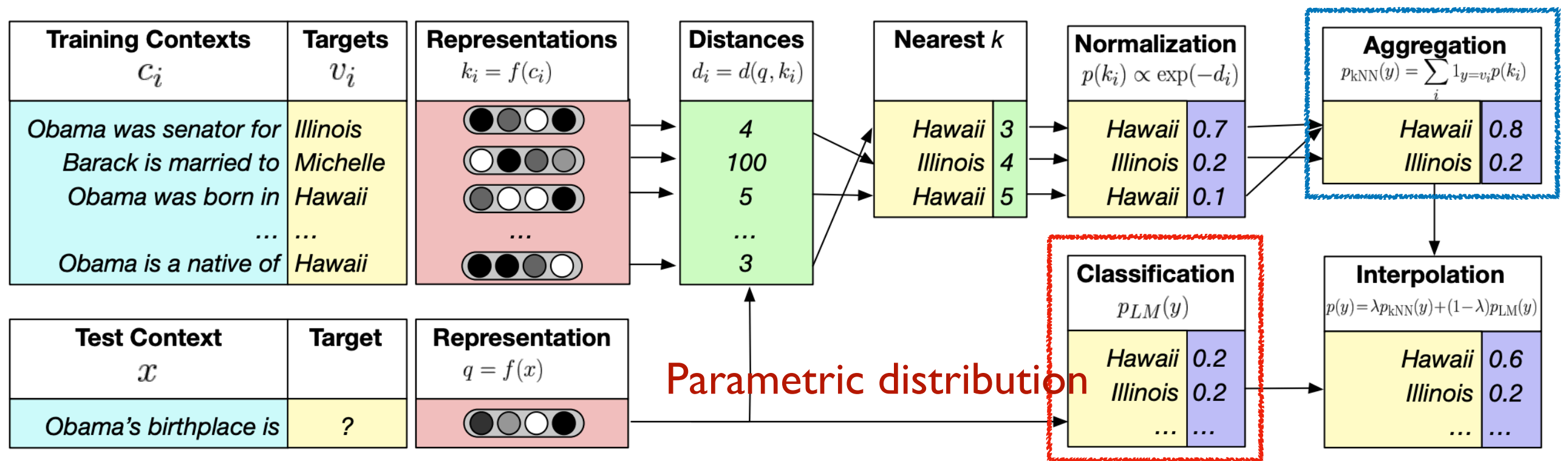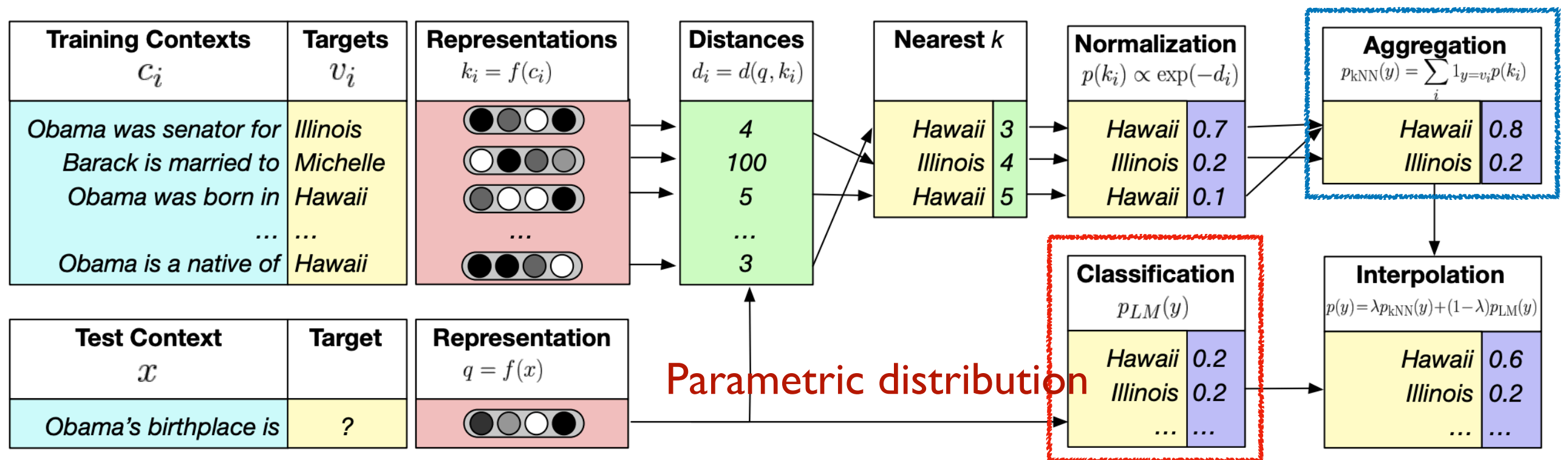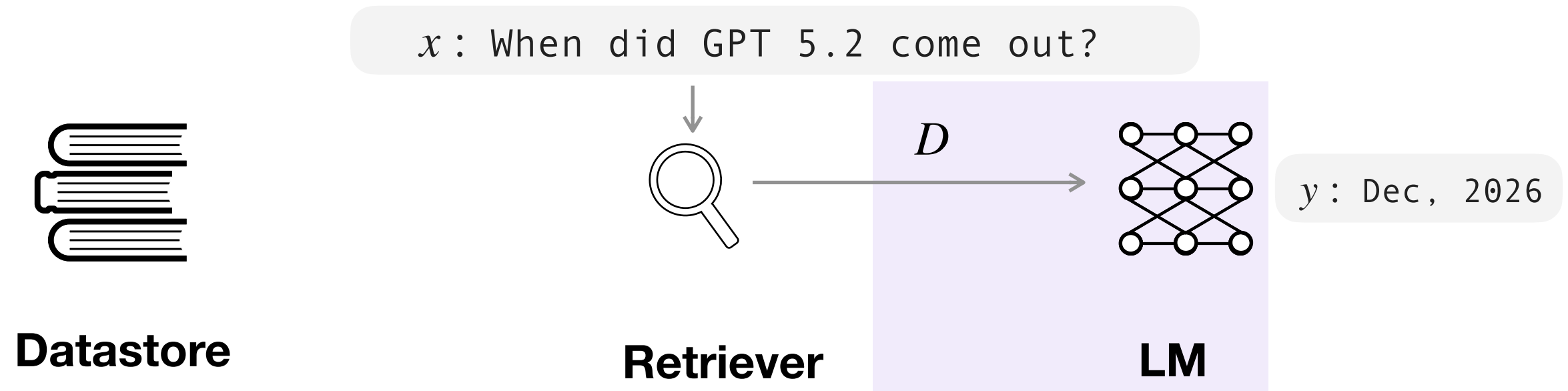
62

# kNN-LM



Nonparametric distribution

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

Representations $k_i = f(c_i)$

Distances $d_i = d(q, k_i)$

| | |
|---|---|
| 4 | |
| 100 | |
| 5 | |
| … | |
| 3 | |

Nearest $k$

| Hawaii | 3 |
|---|---|
| Illinois | 4 |
| Hawaii | 5 |

Normalization $p(k_i) \propto \exp(-d_i)$

| Hawaii | 0.7 |
|---|---|
| Illinois | 0.2 |
| Hawaii | 0.1 |

Aggregation $p_{\text{kNN}}(y) = \sum_i \mathbb{1}_{y=v_i} p(k_i)$

| Hawaii | 0.8 |
|---|---|
| Illinois | 0.2 |

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

Representation $q = f(x)$

Parametric distribution

Classification $p_{LM}(y)$

| Hawaii | 0.2 |
|---|---|
| Illinois | 0.2 |
| … | … |

Interpolation $p(y) = \lambda p_{\text{kNN}}(y) + (1-\lambda)p_{\text{LM}}(y)$

| Hawaii | 0.6 |
|---|---|
| Illinois | 0.2 |
| … | … |

$\lambda$ : hyperparameter

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

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

62

# Summary of Part 3

$x$ : When did GPT 5.2 come out?

$D$

$y$ : Dec, 2026

**Datastore**

**Retriever**

**LM**
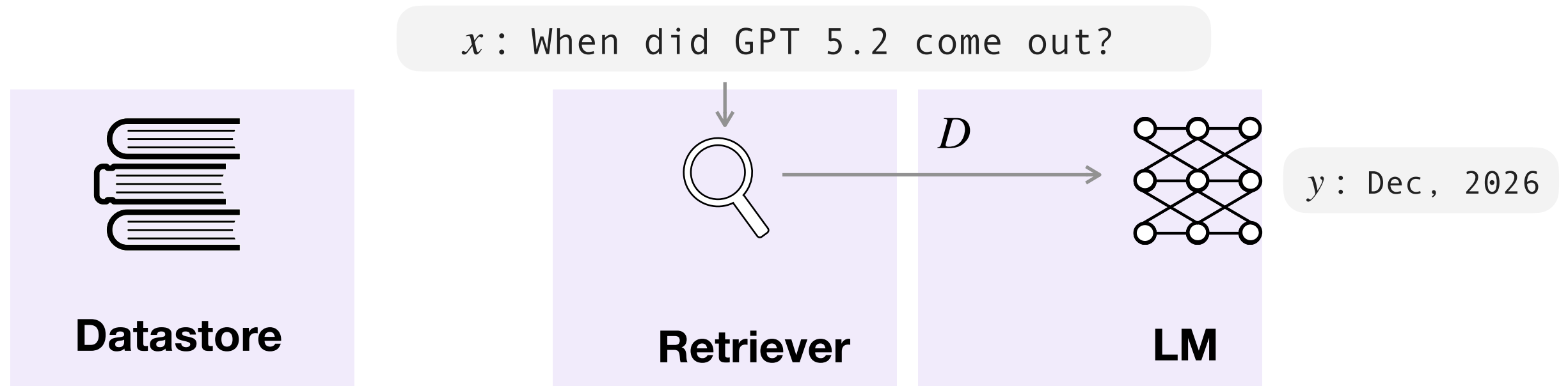
✔ Architectures

✔ Training

✔ Inference

- RAG is widely used but several limitations
- Recent progress to overcome such shortcomings e.g., Deep Research
- Other architectures: intermediate incorporation or output interpolation gain while adding challenges

# Retrieval & RAG

$x$ : When did GPT 5.2 come out?

**Datastore**

**Retriever**

$D$

**LM**

$y$ : Dec, 2026

✔ Sources of datastore

✔ Processing

✔ Scaling

✔ Types of retrievers

✔ Training

✔ Evaluations

✔ Architectures

✔ Training

✔ Inference

🌐 https://akariasai.github.io/          ✉ aasai@andrew.cmu.edu