

Advancing state of the art of Retrieval-augmented Generation

Akari Asai

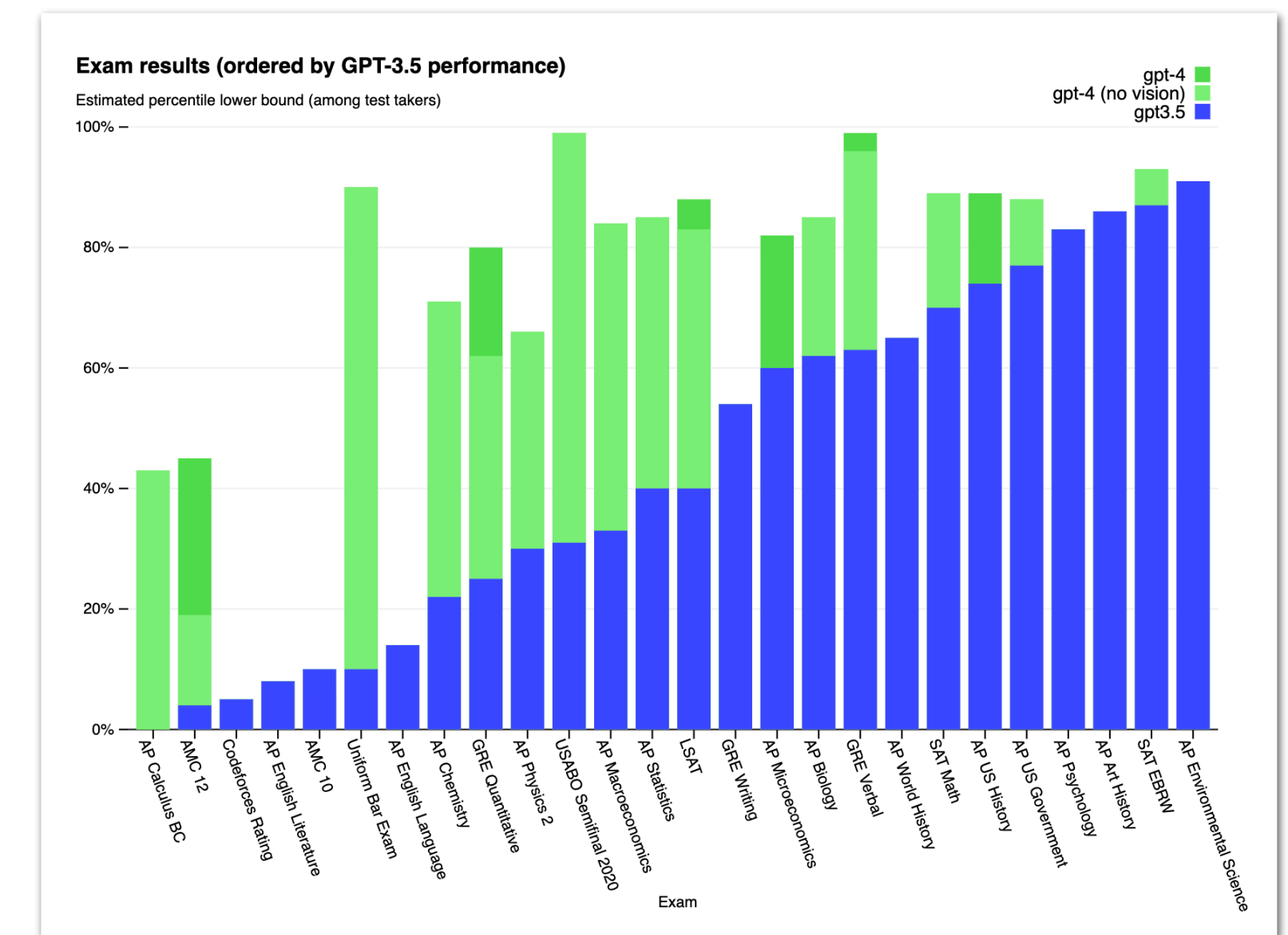
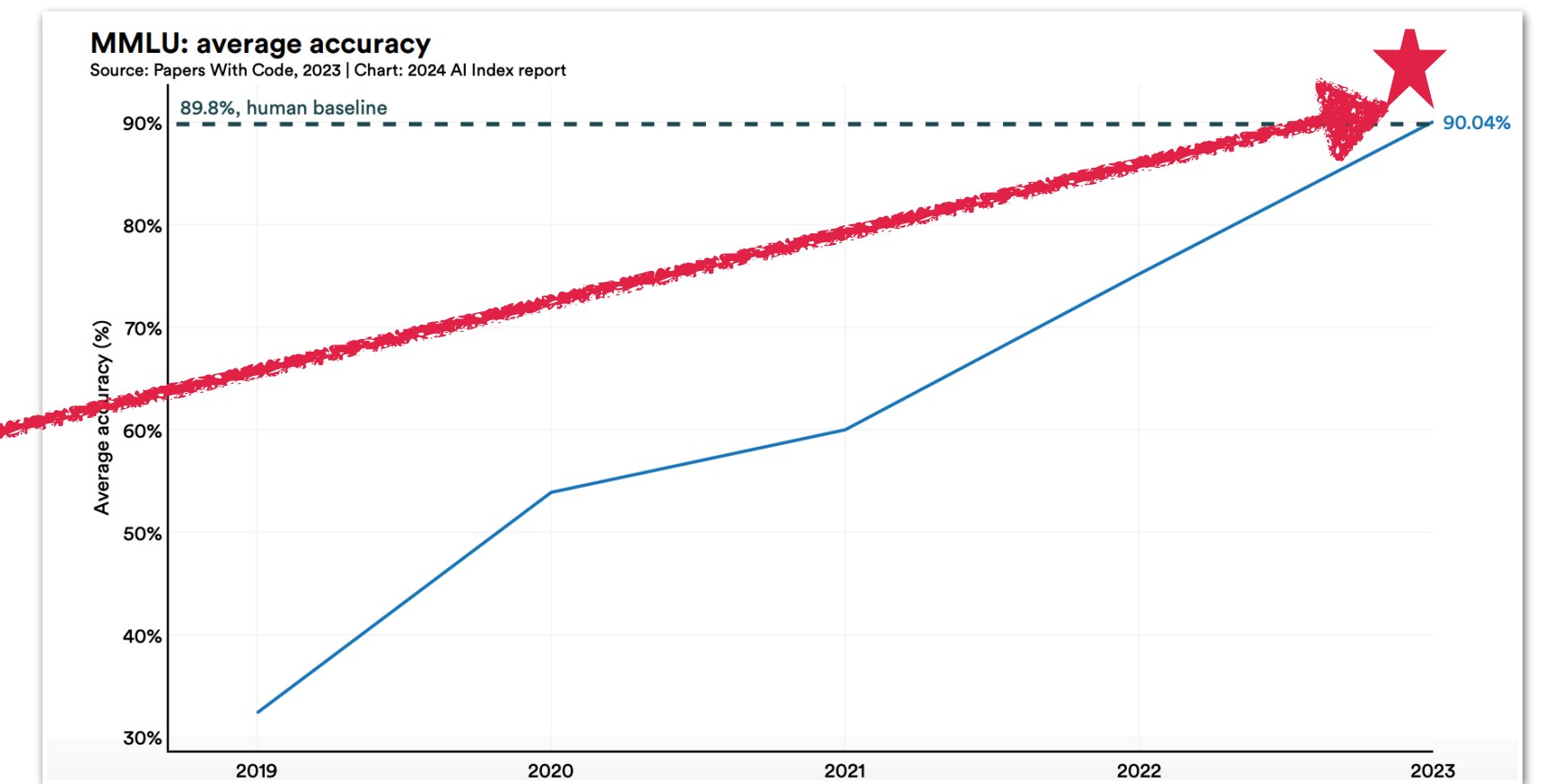
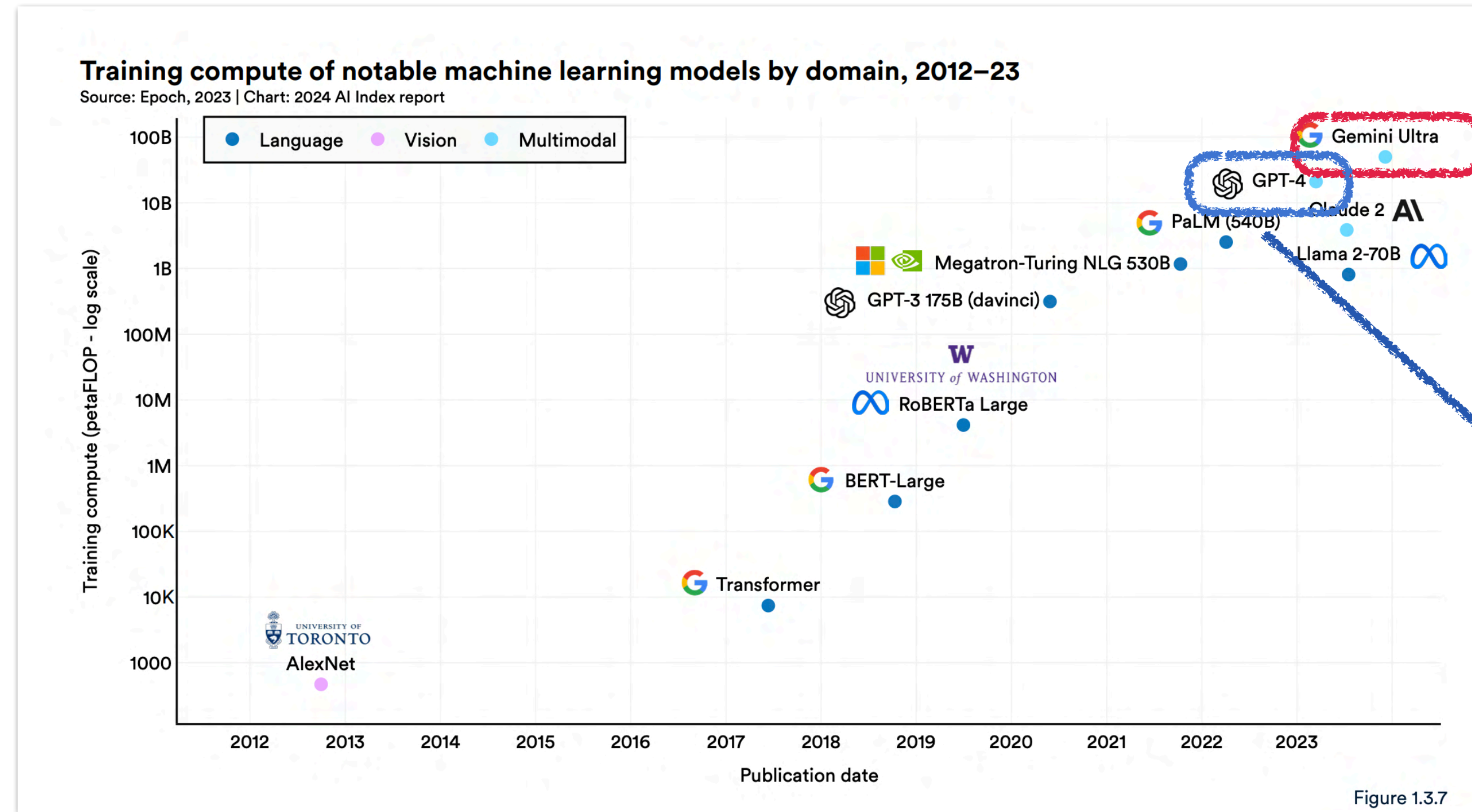
University of Washington

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



April 18 | CMU CS 11-711

LLMs got much more powerful as they scale



Stanford University HAI. Artificial Intelligence Index Report 2024.

Open AI. GPT-4 (OpenAI Blog). 2023.

Hallucinations in LM outputs (April, 2024 ver.)

AK You

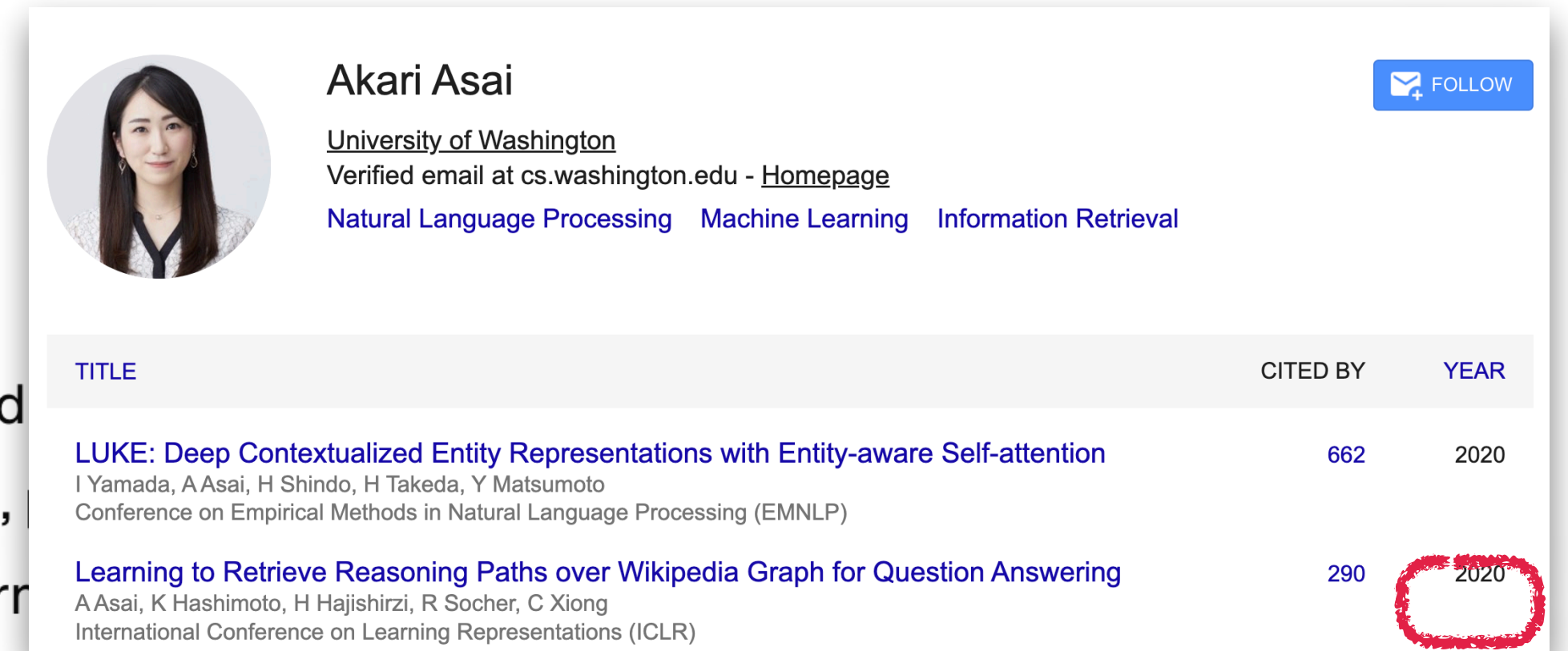
Explain Akari Asai, including her famous papers. (*April 16, 2024)

ChatGPT

Akari Asai is a prominent figure in the field of natural language processing (NLP) and machine learning. She has made significant contributions to various areas within these fields, including the development of techniques for language understanding and representation learning.

One of Asai's notable papers is "Learning to Retrieve Reasoning Paths over Wikipedia Graph for Question Answering" (2018) which she co-authored. In this paper, the authors introduced a method for question answering that utilizes the structure of Wikipedia to retrieve relevant information. They proposed a novel approach that learns to navigate the Wikipedia graph to find paths connecting relevant pieces of information, improving the accuracy of question answering systems.

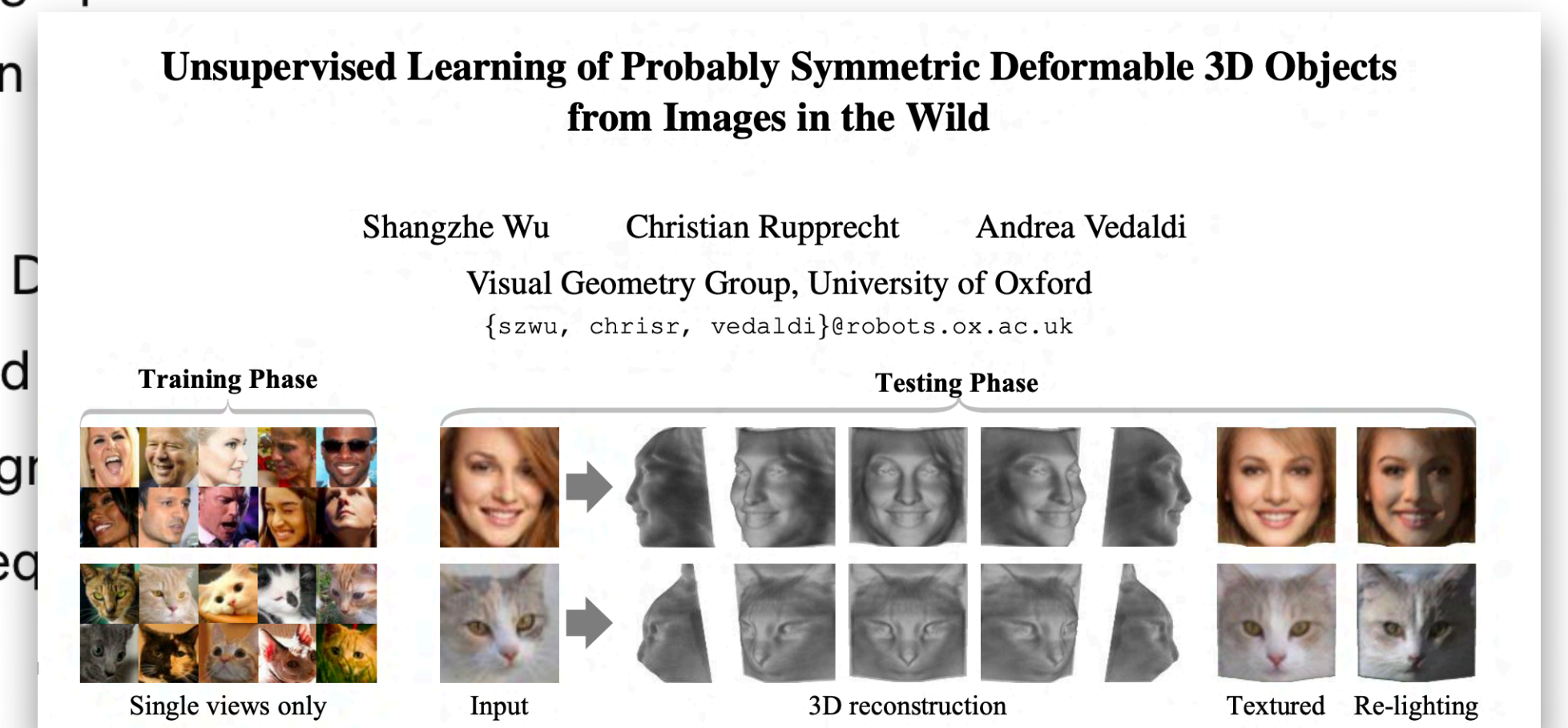
Another influential paper by Asai is "Unsupervised Learning of Probably Symmetric Deformable 3D Objects from Images in the Wild" (2019), where she and her co-authors presented a method for unsupervised learning of 3D object representations from 2D images. This work is significant because it addresses the challenge of learning 3D object representations without requiring explicit supervision, which is crucial for many computer vision tasks.



Akari Asai FOLLOW

University of Washington
Verified email at cs.washington.edu - [Homepage](#)
Natural Language Processing Machine Learning Information Retrieval

TITLE	CITED BY	YEAR
LUKE: Deep Contextualized Entity Representations with Entity-aware Self-attention I Yamada, A Asai, H Shindo, H Takeda, Y Matsumoto Conference on Empirical Methods in Natural Language Processing (EMNLP)	662	2020
Learning to Retrieve Reasoning Paths over Wikipedia Graph for Question Answering A Asai, K Hashimoto, H Hajishirzi, R Socher, C Xiong International Conference on Learning Representations (ICLR)	290	2020

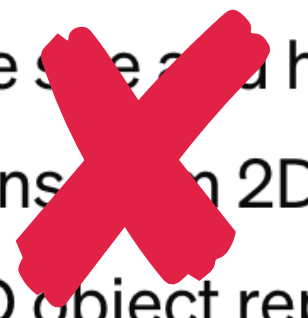
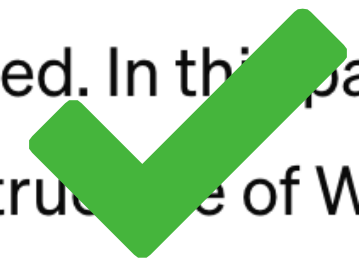


Unsupervised Learning of Probably Symmetric Deformable 3D Objects from Images in the Wild

Shangzhe Wu Christian Rupprecht Andrea Vedaldi
Visual Geometry Group, University of Oxford
{szwu, chrisr, vedaldi}@robots.ox.ac.uk

Training Phase **Testing Phase**

Single views only Input 3D reconstruction Textured Re-lighting

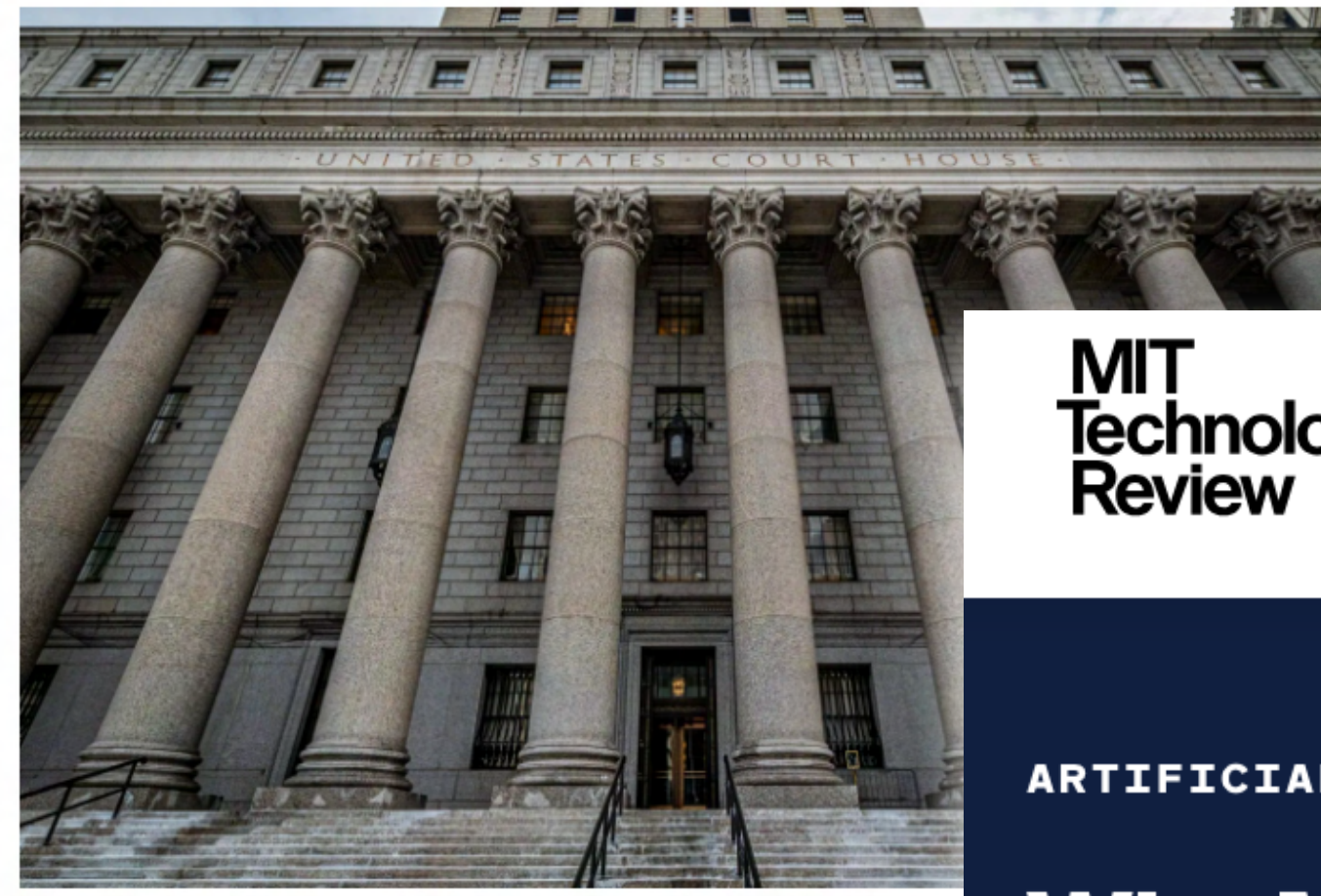


Catastrophic errors as results of LM hallucinations

TECH · LAW

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

BY RACHEL SHIN
June 23, 2023 at 9:41 AM PDT



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

MIT
Technology
Review

Featured Topics Newsletters Events Podcasts

SIGN IN

ARTIFICIAL INTELLIGENCE

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

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

By Will Douglas Heaven

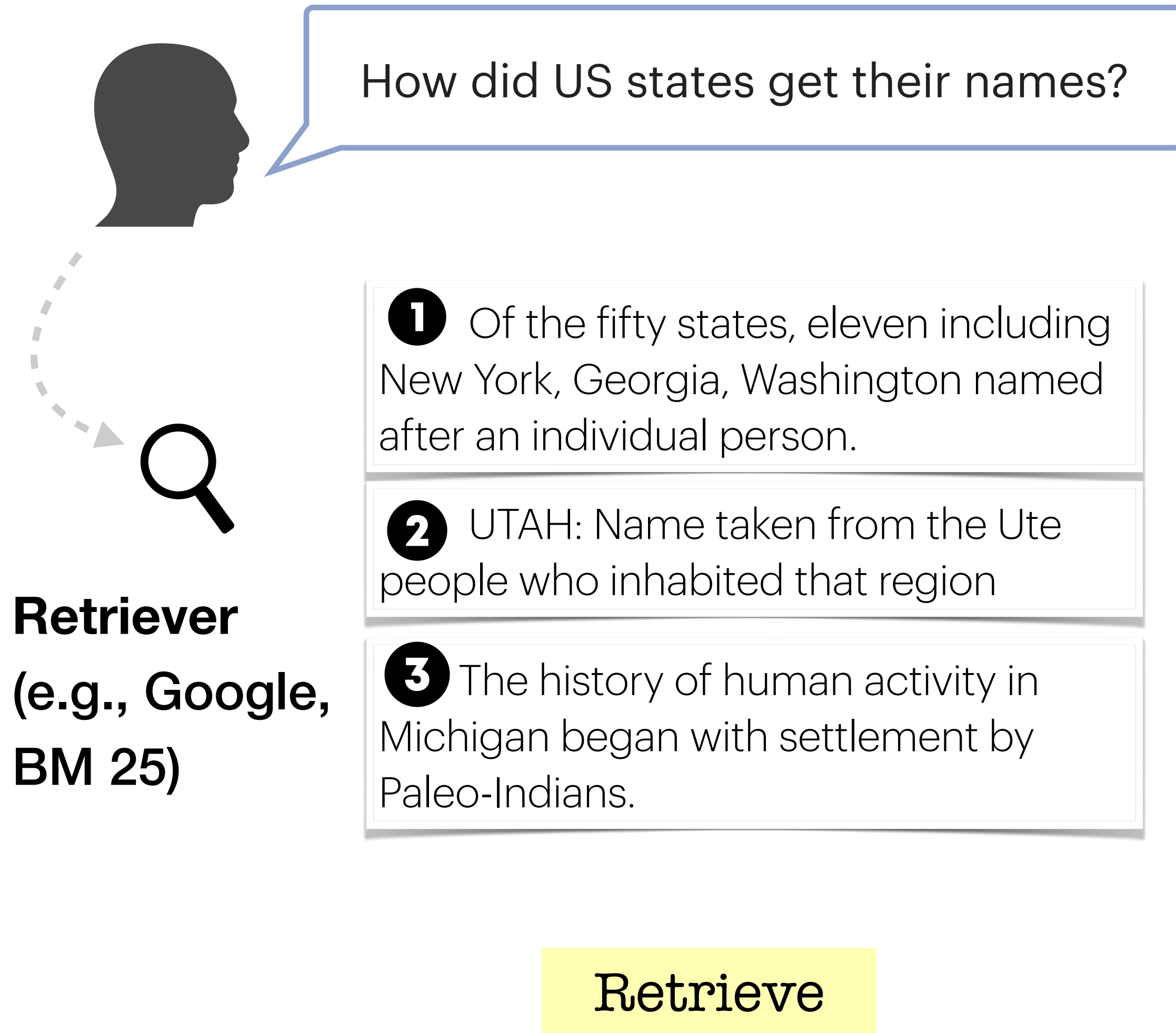
November 18, 2022

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

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

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

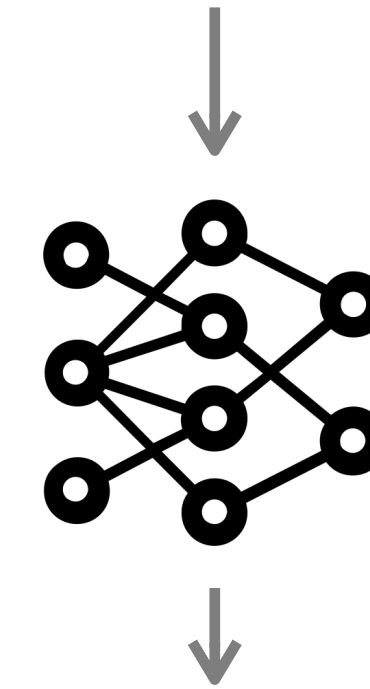
Retrieval-augmented generations (RAG)



Answer my question using references.

References: 1 2 3

Question: How did US states get their names?

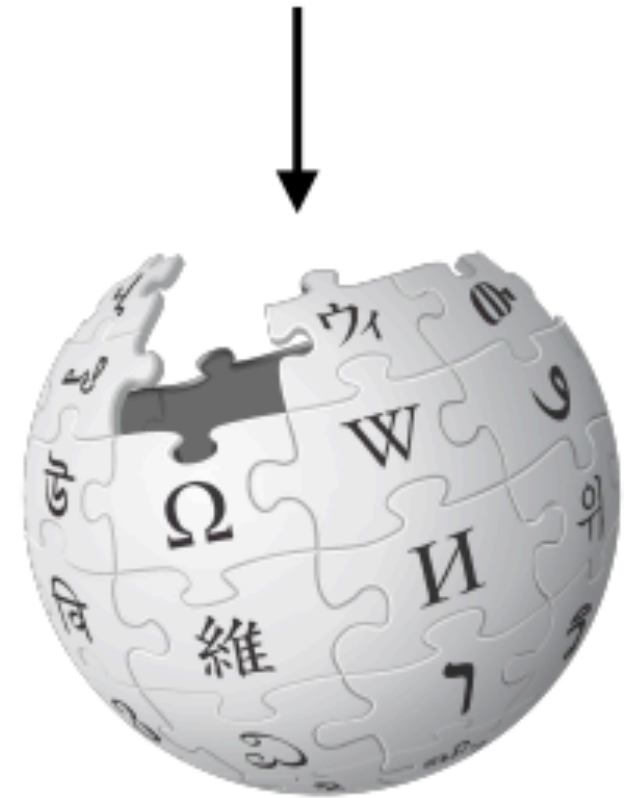


Eleven states are named after an individual person. Some states including Utah are named after native American tribe names.

Read

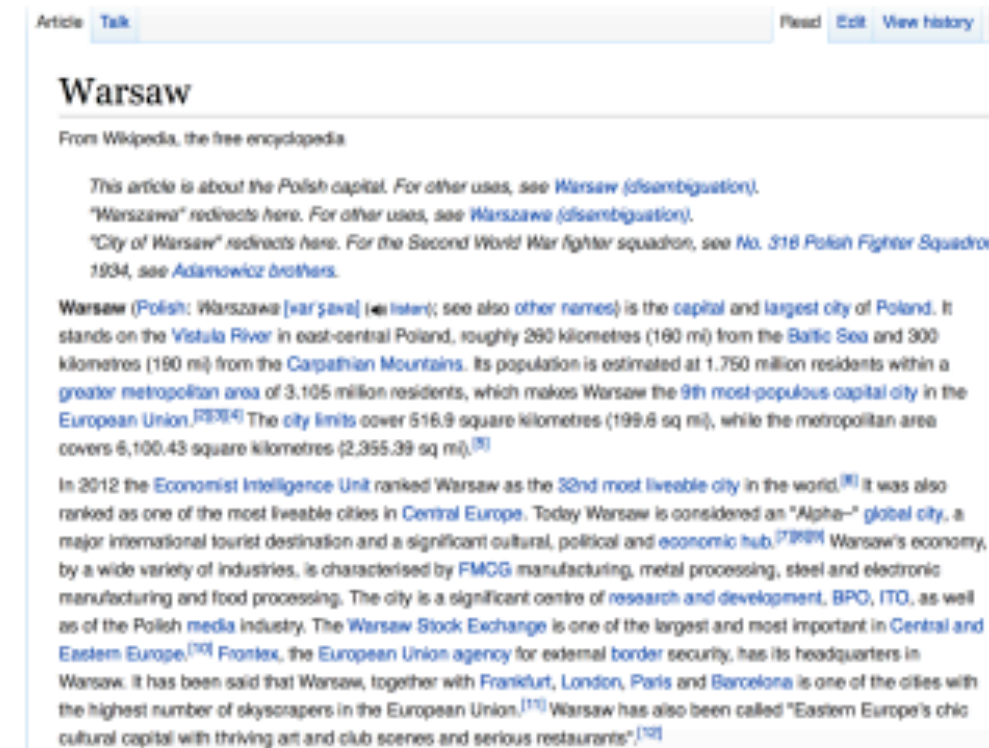
Retrieval-augmented generations (RAG)

Q: How many of Warsaw's inhabitants spoke Polish in 1933?



WIKIPEDIA
The Free Encyclopedia

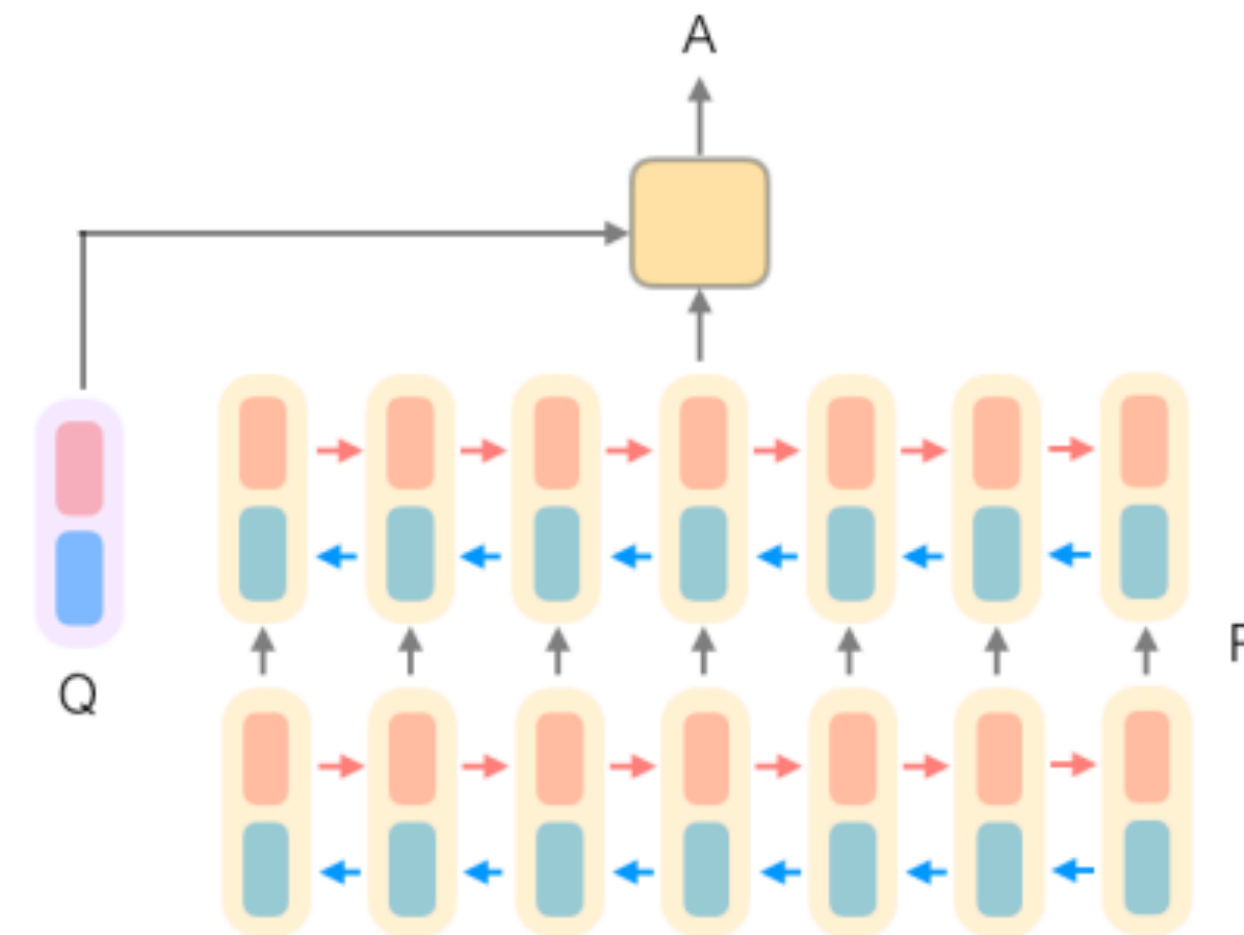
Document
Retriever



Retrieve

Document
Reader

833,500

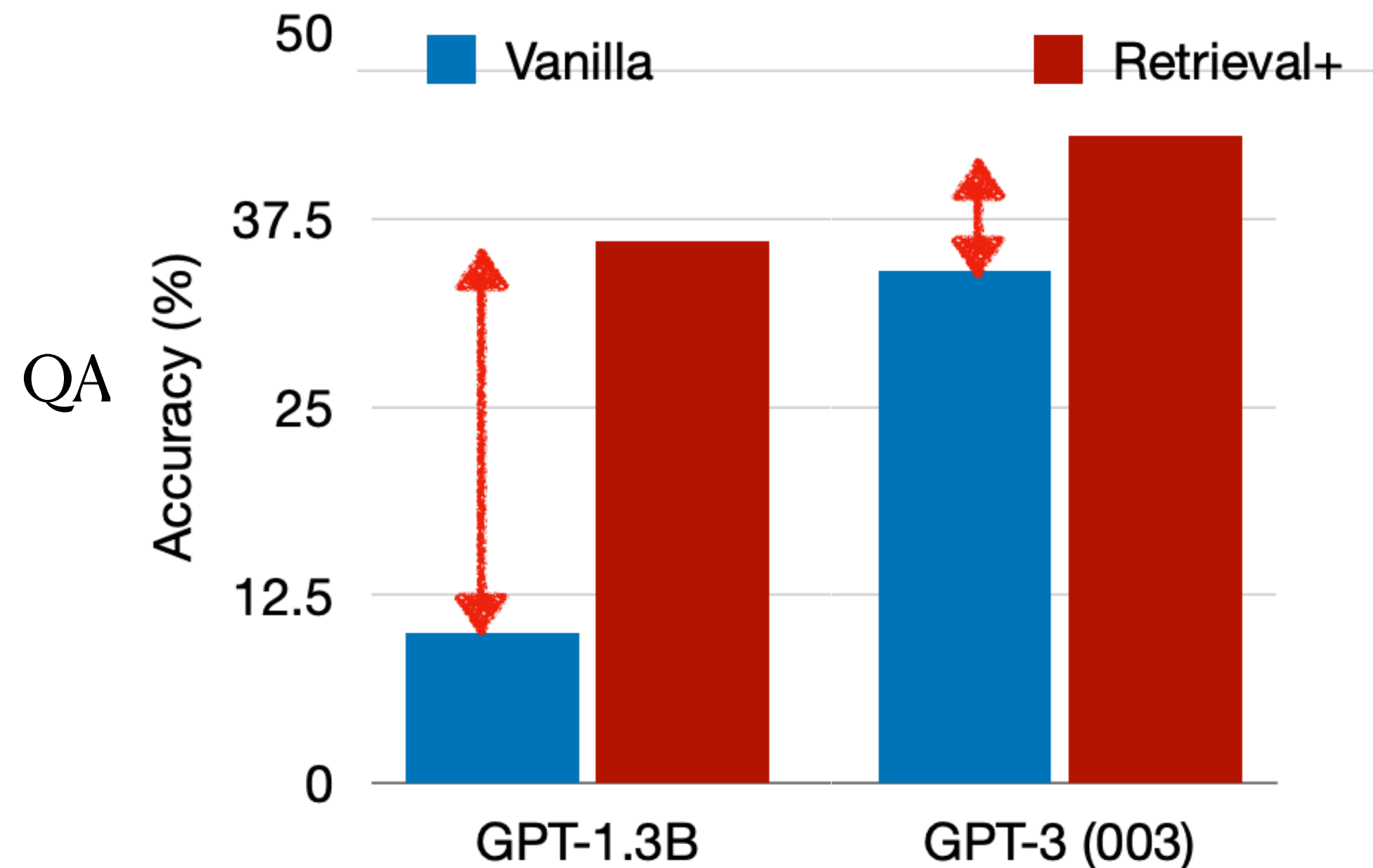


Read

RAG has been widely and successfully adopted

RAG has shown effective in many benchmarks

Widely applied to real-world production systems



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

The shift from traditional QA to *open-ended* instructions



How many of US states got their names from an individual person?

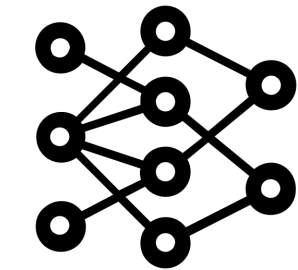
Simple factoid question



WIKIPEDIA
The Free Encyclopedia

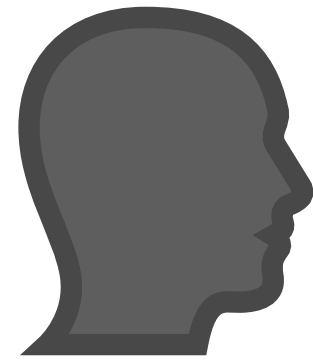
Of the fifty states, eleven are named after an individual person. Six of those are named in honor of European monarchs: the two Carolinas, the two Virginias, Georgia, and Louisiana.

Eleven



Short answer based on single document

The shift from traditional QA to *open-ended* instructions



Make a table for me summarizing how different US states got their names, grouping them together.

Group	States
Named after indigenous communities	Alabama, Alaska, Arizona, Arkansas, Colorado, Connecticut, Delaware, Florida, Hawaii, Idaho, Illinois, Iowa, Kansas, Kentucky, Louisiana, Maine, Maryland, Massachusetts, Michigan, Minnesota, Mississippi, Missouri, Montana, Nebraska, Nevada, New Mexico, North and South Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Rhode Island, Tennessee, Texas, Utah, Vermont, Virginia, Washington, West Virginia, Wisconsin, Wyoming
Named after royalty and other important figures	Georgia, Louisiana, Maryland, New York, North and South Carolina, Pennsylvania, Virginia, Washington, West Virginia

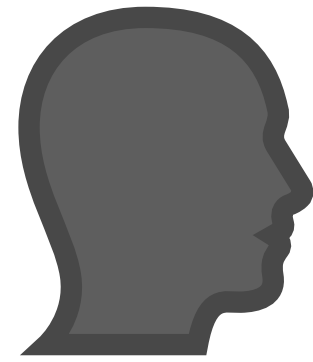
10 references

- WIKIPEDIA
LIST OF STATE AND TERRITORY NAME ETY...
- BIA
ORIGIN OF NAMES OF US STATES | INDIA...
- MENTALFLOSS
HOW ALL 50 STATES GOT THEIR NAMES | ...



Is this fully correct?

The shift from traditional QA to *open-ended* instructions

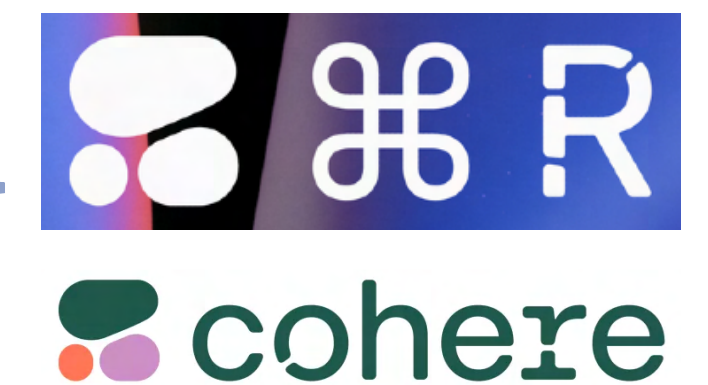


Make a table for me summarizing how different US states got their names, grouping them together.

Group	States
Named after indigenous communities	Alabama, Alaska, Arizona, Arkansas, Colorado, Connecticut, Delaware, Florida, Hawaii, Idaho, Illinois, Iowa, Kansas, Kentucky, Louisiana, Maine, Maryland, Massachusetts, Michigan, Minnesota, Mississippi, Missouri, Montana, Nebraska, Nevada, New Mexico, North and South Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Rhode Island, Tennessee, Texas, Utah, Vermont, Virginia, Washington, West Virginia, Wisconsin, Wyoming
Named after royalty and other important figures	Georgia, Louisiana, Maryland, New York, North and South Carolina, Pennsylvania, Virginia, Washington, West Virginia <i>Hint: Only 10 states here...?</i>

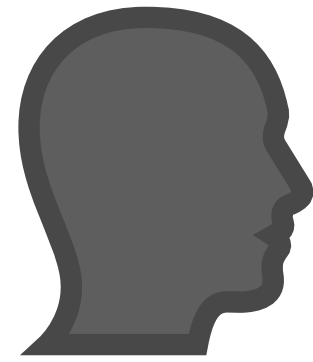
10 references

- WIKIPEDIA
LIST OF STATE AND TERRITORY NAME ETY...
- BIA
ORIGIN OF NAMES OF US STATES | INDIA...
- MENTALFLOSS
HOW ALL 50 STATES GOT THEIR NAMES | ...



Is this fully correct?

The shift from traditional QA to *open-ended* instructions



Make a table for me summarizing how different US states got their names, grouping them together.

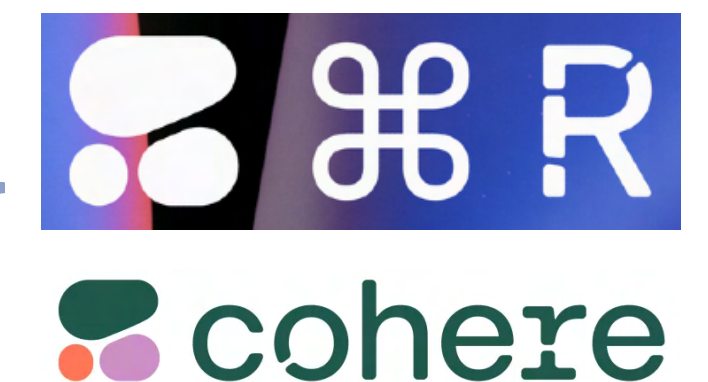
Group	States
Named after indigenous communities	Alabama, Alaska, Arizona, Arkansas, Colorado, Connecticut, Delaware, Florida, Hawaii, Idaho, Illinois, Iowa, Kansas, Kentucky, Louisiana, Maine, Maryland, Massachusetts, Michigan, Minnesota, Mississippi, Missouri, Montana, Nebraska, Nevada, New Mexico, North and South Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Rhode Island, Tennessee, Texas, Utah, Vermont, Virginia, Washington, West Virginia, Wisconsin, Wyoming
Named after royalty and other important figures	Georgia, Louisiana, Maryland, New York, North and South Carolina, Pennsylvania, Virginia, Washington, West Virginia

DELAWARE: Named for Lord De La Warr, first governor and captain-general of Virginia, who in 1630 explored the bay and river area where his name was first applied.

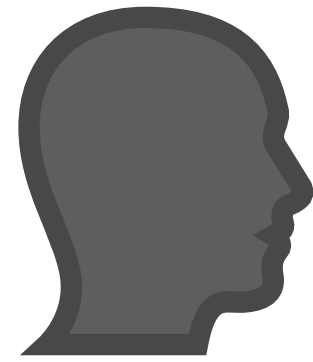
Indian Affairs (.gov)
https://www.bia.gov › as-ia › opa › online-press-release

[Origin of Names of US States | Indian Affairs](#)

- 10 references
- WIKIPEDIA
LIST OF STATE AND TERRITORY NAME ETY...
 - BIA
ORIGIN OF NAMES OF US STATES | INDIA...
 - MENTALFLOSS
HOW ALL 50 STATES GOT THEIR NAMES | ...




The shift from traditional QA to *open-ended* instructions



Make a table for me summarizing how different US states got their names, grouping them together.

Group	States
Named after indigenous communities	Alabama, Alaska, Arizona, Arkansas, Colorado, Connecticut, Delaware, Florida, Hawaii, Idaho, Illinois, Iowa, Kansas, Kentucky, Louisiana, Maine, Maryland, Massachusetts, Michigan, Minnesota, Mississippi, Missouri, Montana, Nebraska, Nevada, New Mexico, North and South Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Rhode Island, Tennessee, Texas, Utah, Vermont, Virginia, Washington , West Virginia, Wisconsin, Wyoming
Named after royalty and other important figures	Georgia, Louisiana, Maryland, New York, North and South Carolina, Pennsylvania, Virginia, Washington , West Virginia



- 10 references
- WIKIPEDIA
LIST OF STATE AND TERRITORY NAME ETY...
 - BIA
ORIGIN OF NAMES OF US STATES | INDIA...
 - MENTALFLOSS
HOW ALL 50 STATES GOT THEIR NAMES | ...



The shift from traditional QA to *open-ended* instructions



Make a table for me summarizing how different US states got their names, grouping them together.

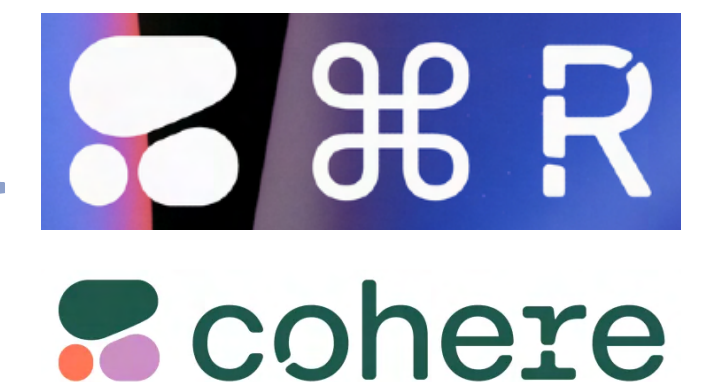
Complex instruction

Group	States	Long-form answers
Named after indigenous communities	Alabama, Alaska, Arizona, Arkansas, Colorado, Connecticut, Delaware, Florida, Hawaii, Idaho, Illinois, Iowa, Kansas, Kentucky, Louisiana, Maine, Maryland, Massachusetts, Michigan, Minnesota, Mississippi, Missouri, Montana, Nebraska, Nevada, New Mexico, North and South Dakota, Ohio, Oklahoma, Oregon, Pennsylvania, Rhode Island, Tennessee, Texas, Utah, Vermont, Virginia, Washington, West Virginia, Wisconsin, Wyoming	
Named after royalty and other important figures	Georgia, Louisiana, Maryland, New York, North and South Carolina, Pennsylvania, Virginia, Washington, West Virginia	

Requiring aggregating multiple evidence

10 references

- WIKIPEDIA
LIST OF STATE AND TERRITORY NAME ETY...
- BIA
ORIGIN OF NAMES OF US STATES | INDIA...
- MENTALFLOSS
HOW ALL 50 STATES GOT THEIR NAMES | ...



Challenges of the current naive RAG systems: reliability



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



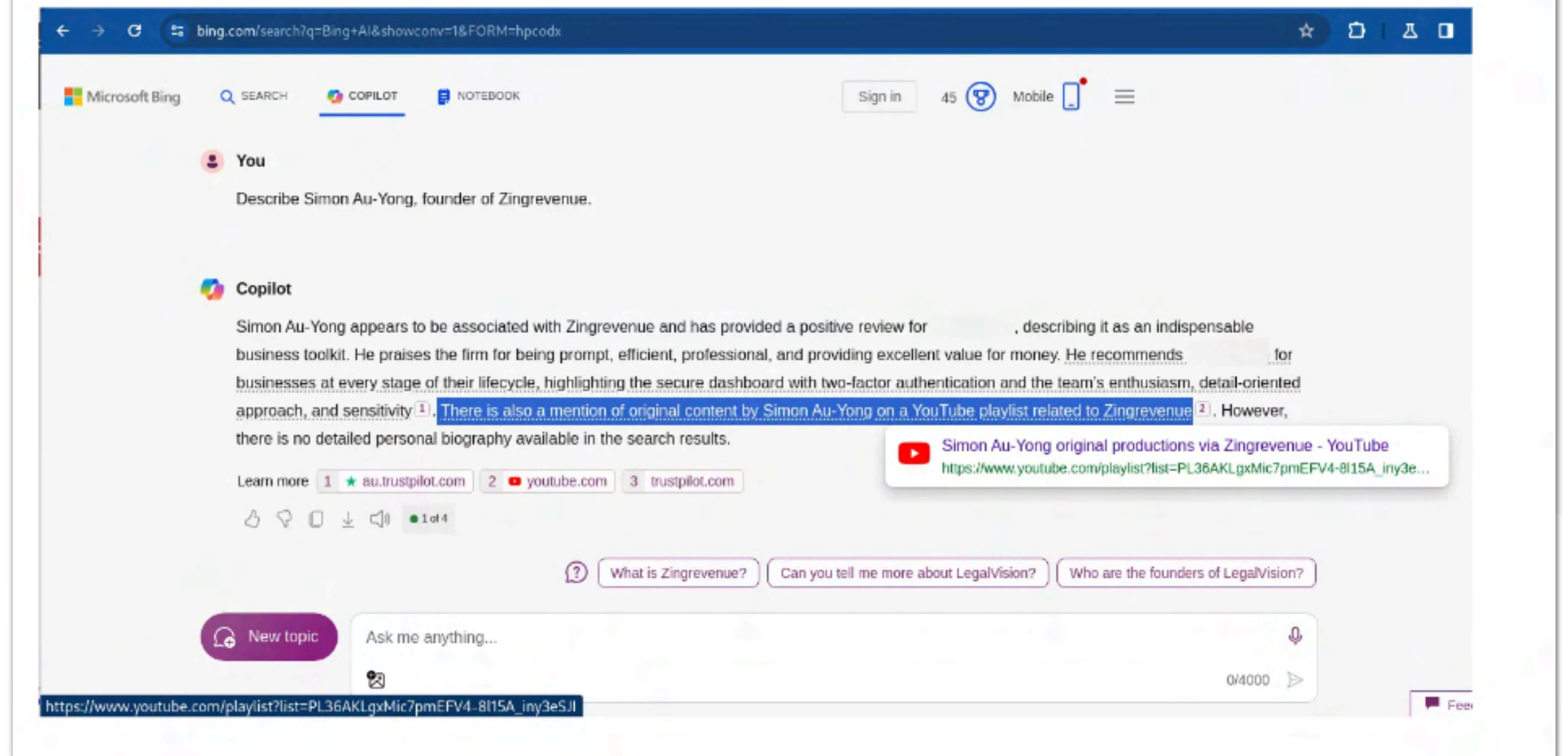
The James Webb Space Telescope is designed to peer into the dusty clouds of gas where stars and planetary systems are born. Webb has captured the first direct image of an exoplanet, and the Pillars of Creation in the Eagle Nebula[1][2]. Additionally, the telescope will be used to study the next interstellar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

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

Cited Webpages

- [1]: nasa.gov (✗ citation does not support its associated statement)
[NASA's Webb Confirms Its First Exoplanet](#)
... Researchers confirmed an exoplanet, a planet that orbits another star, using NASA's James Webb Space Telescope for the first time. ...
- [2]: cnn.com (⚠ citation partially supports its associated statement)
[Pillars of Creation: James Webb Space Telescope ...](#)
... The Pillars of Creation, in the Eagle Nebula, is a star-forming region captured in a new image (right) by the James Webb Space Telescope that reveals more detail than a 2014 image (left) by Hubble ...
- [3]: nasa.gov (✓ citation fully supports its associated statement)
[Studying the Next Interstellar Interloper with Webb](#)
... Scientists have had only limited ability to study these objects once discovered, but all of that is about to change with NASA's James Webb Space Telescope... The team will use Webb's spectroscopic capabilities in both the near-infrared and mid-infrared bands to study two different aspects of the interstellar object.

I asked Bing Copilot to describe me. It did and said that there is a mention of original content by Simon Au-Yong on a YouTube playlist related to Zingrevenue (my company). The link is at the bottom of the screenshot and there is a button that should send me to that playlist. But the playlist is made up.

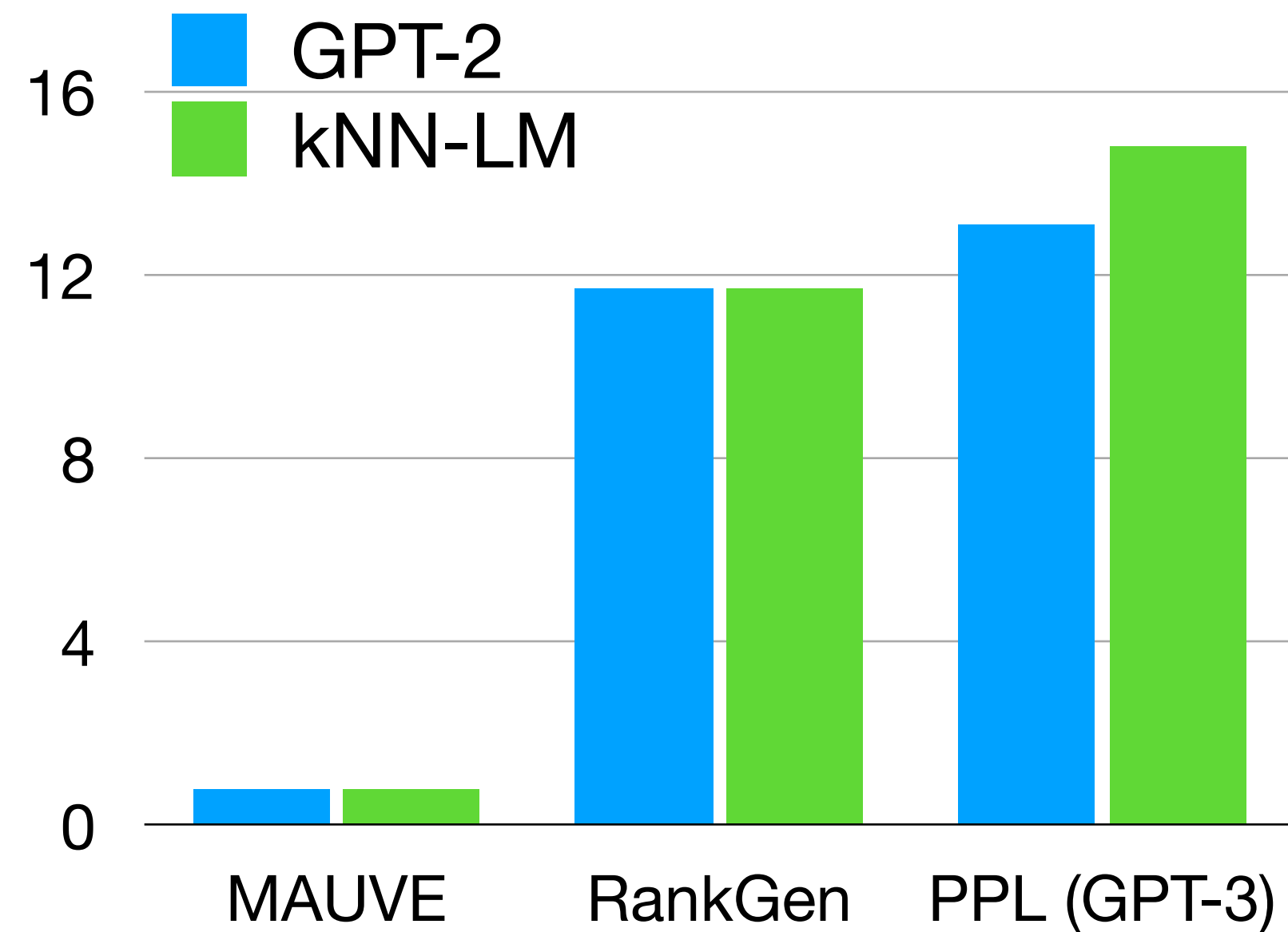


Liu et al. Evaluating Verifiability in Generative Search Engines. Findings of EMNLP 2023.

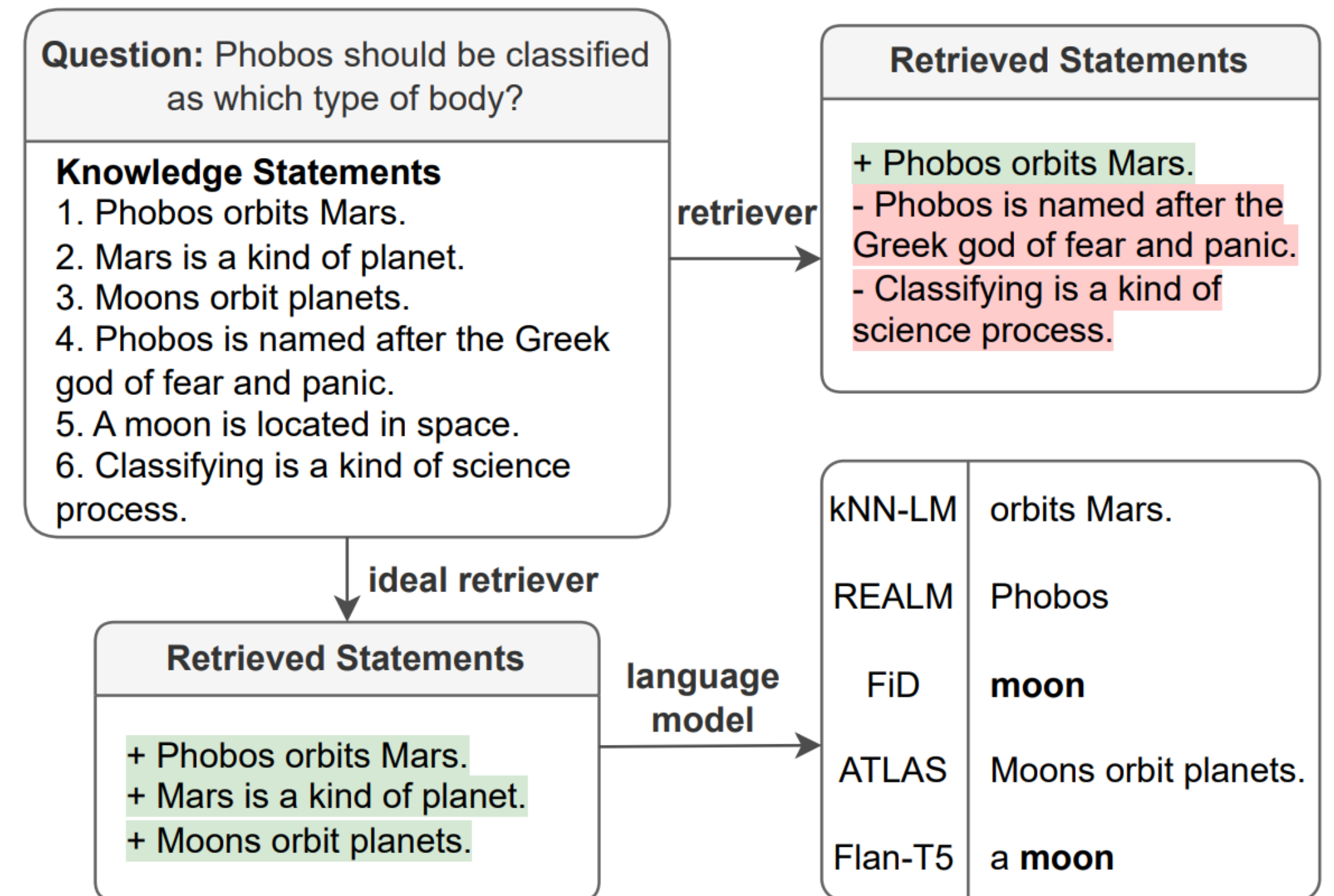
Marcus. No, RAG is probably not going to rescue the current situation. 2024.

Challenges of the current naive RAG systems: versatility

Limited effectiveness beyond information-seeking QA-like tasks



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



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

Today's lecture

Promises and limitations of retrieval-augmented LMs

Reliable inference: Self-reflective RAG with dynamic retrieval

Versatile Retriever: Intent-aware retrievers with LMs

Summary and future directions: RAG in the wild

Today's lecture

Promises and limitations of retrieval-augmented LMs

Reliable inference: Self-reflective RAG with dynamic retrieval

Versatile Retriever: Intent-aware retrievers with LMs

Summary and future directions: RAG in the wild

Q: Why do we need RAG?

A: Because **RAG** can solve many core
limitations of parametric LMs!

When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories

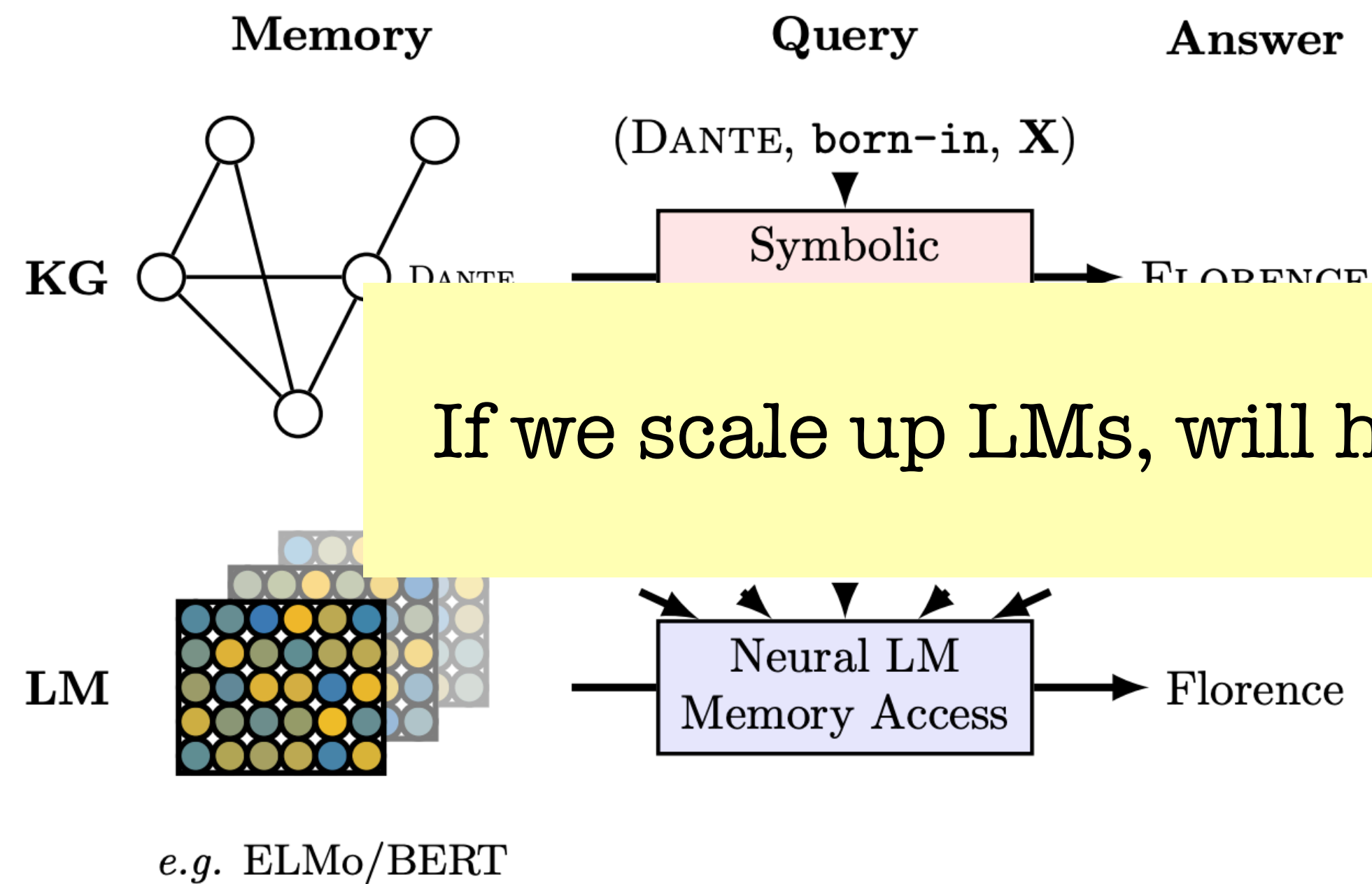
Alex Mallen*, Akari Asai*, Victor Zhong, Rajarshi Das,
Daniel Khashabi, Hannaneh Hajishirzi

* = core contributors

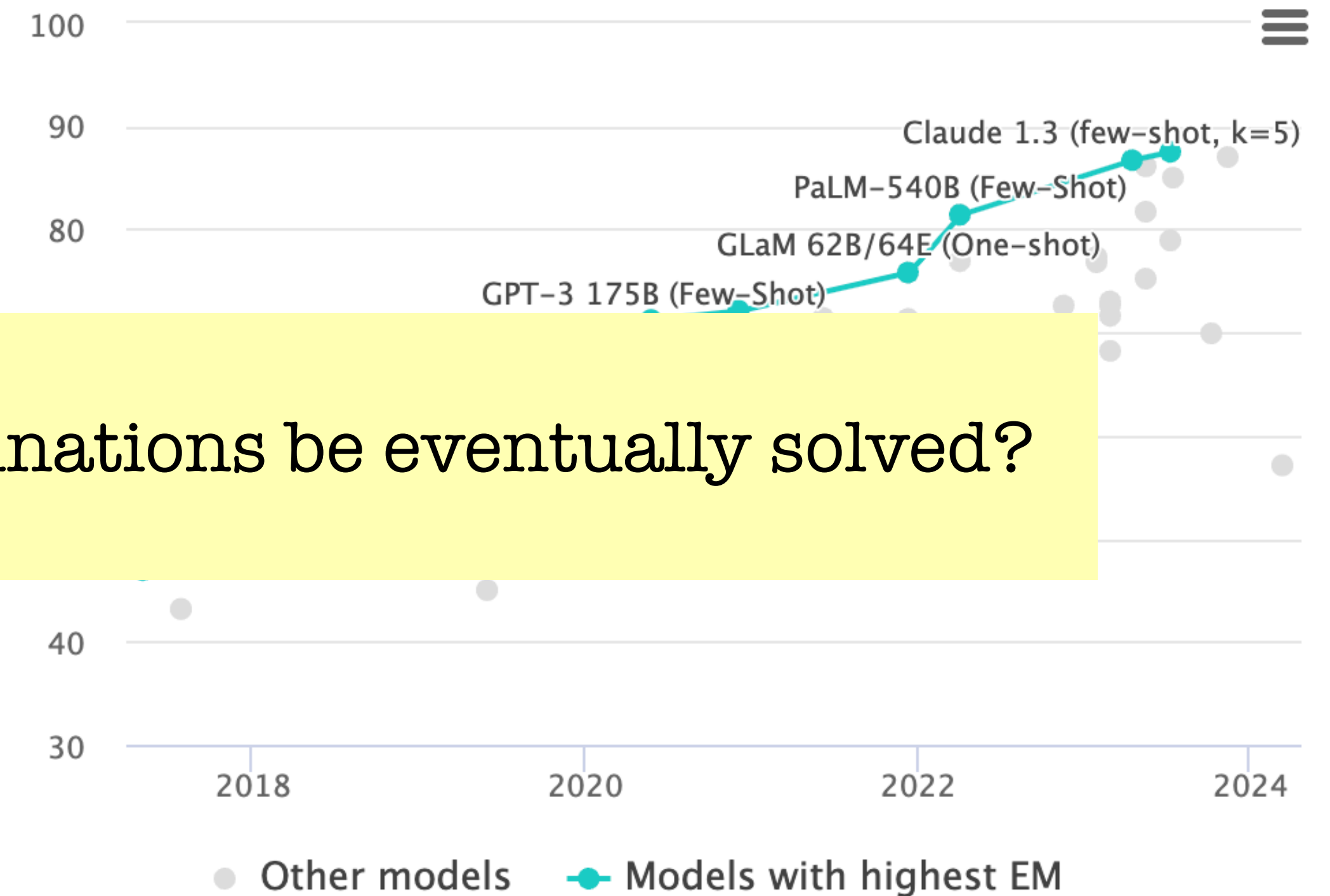


ACL 2023 (Oral, Best Video Award — most viewed)

Factual knowledge memorization in LLMs



If we scale up LMs, will hallucinations be eventually solved?



Petroni et al. Language Models as Knowledge Bases?.
EMNLP 2019.

TriviaQA (paper with code)
<https://paperswithcode.com/sota/question-answering-on-triviaqa>

Focus and task

Focus:
Factual knowledge

(Kathy Saltzman, Occupation, Politician)
Subject Relationship Object

Task:
Open-domain QA

Relationship Subject
What is the occupation of Kathy Saltzman?
A: Politician
Object

Dimensions of analysis



We created a new dataset, PopQA (17k openQA questions with fine-grained meta data). See more details in our paper!

Aspect 1:
Subject entity popularity

(Kathy Saltzman, Occupation, Politician)



(Barack Obama, Occupation, Politician)

Aspect 2:
Relationship type

(Barack Obama, Elementary School, St. Francis Assisi)



(Barack Obama, Occupation, Politician)

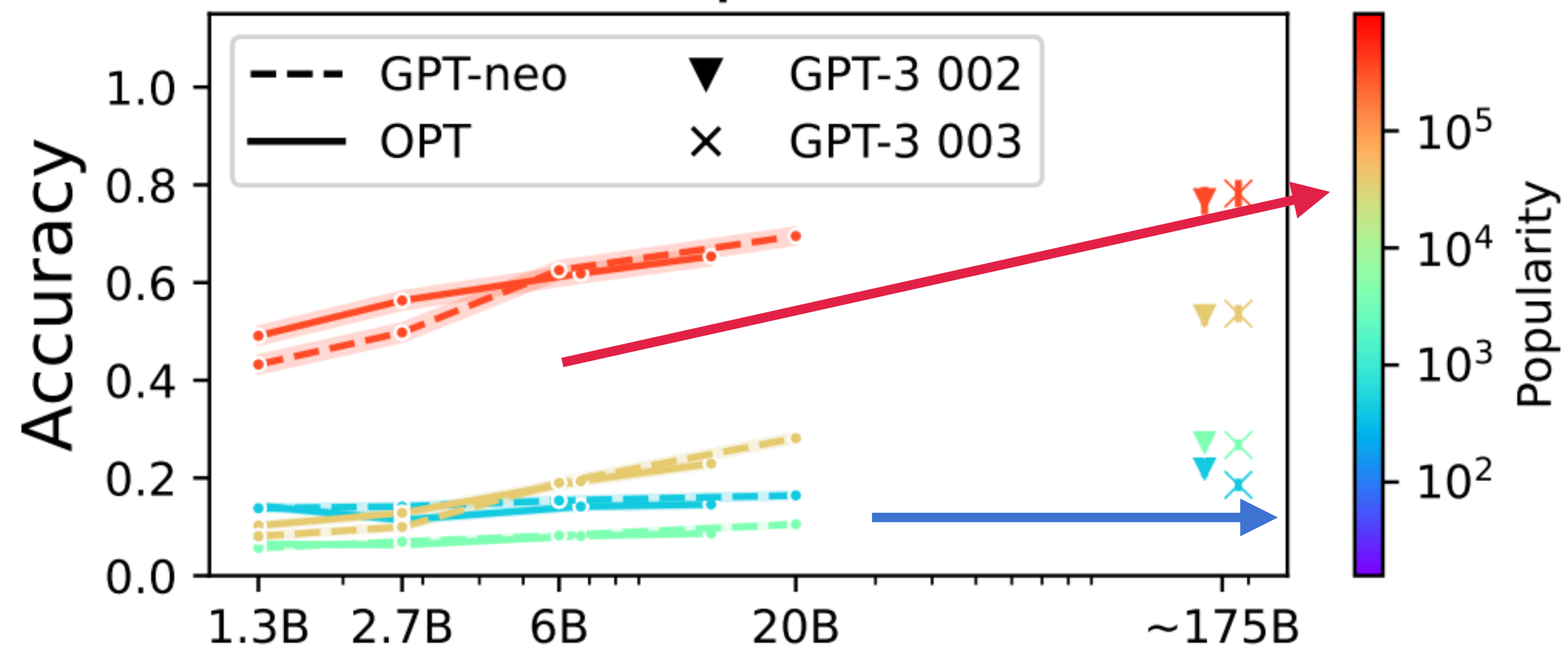
Does scaling solve memorization? **Probably not!**

On popular facts

Performance gets better as scaling.

On long-tail facts

Almost flat trends.



Scaling may not overcome hallucination in long tail!

RAG can address hallucinations in such long-tail!

On popular facts

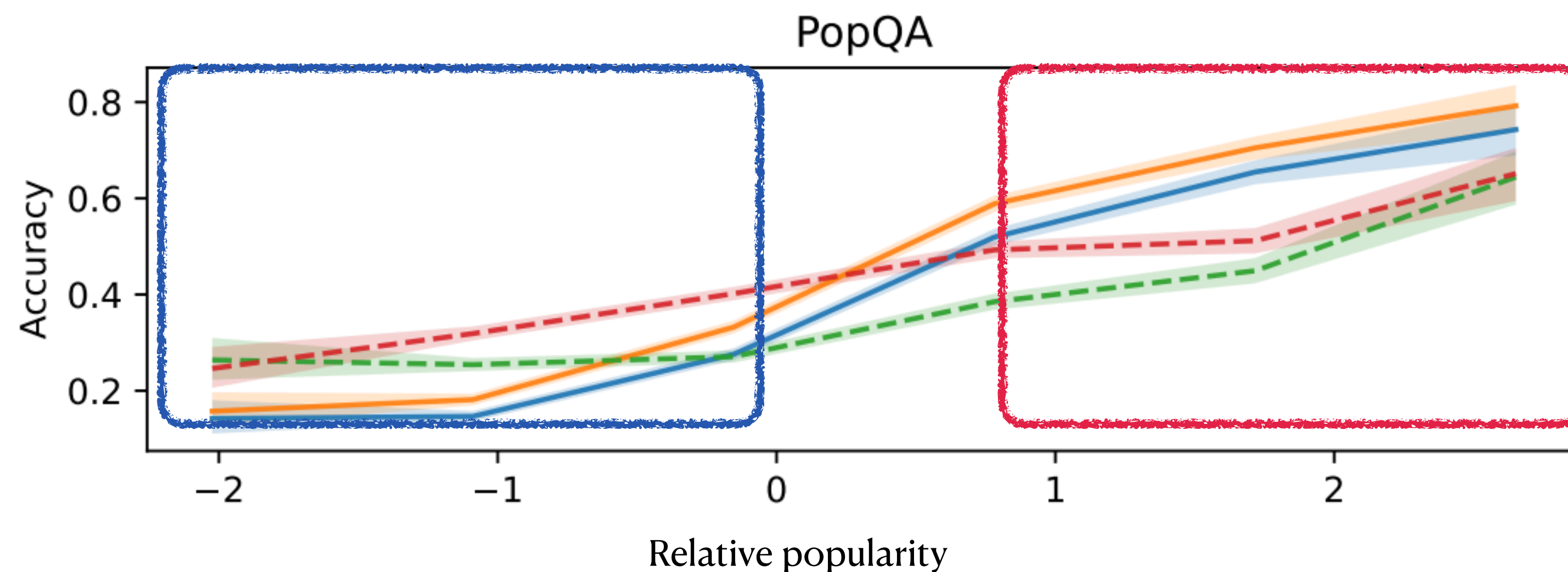
vanilla \approx RAG or vanilla $>$ RAG

will discuss why in the next section!

On long-tail facts

RAG $>$ vanilla

— Vanilla — GenRead - - - BM25 - - - Contriever



RAG gives large gain in long-tail questions!

Is RAG only useful for hallucination?

Reliable, adaptable, attributable LMs with Retrieval

Akari Asai, Zexuan Zhong, Danqi Chen, Pang Wei Koh
Luke Zettlemoyer, Hannaneh Hajishirzi, Wen-tau Yih



Arxiv 2024

<https://arxiv.org/abs/2403.03187>

Core limitations of parametric LMs

Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



You

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



ChatGPT

Akari Asai is a prominent figure in the field of natural language processing (NLP) and machine learning. She has made significant contributions to various areas within these fields, particularly in the development of techniques for language understanding and representation learning.

One of Asai's notable papers is "Learning to Retrieve Reasoning Paths over Wikipedia Graph for Question Answering" (2018), which she co-authored. In this paper, the authors introduced a method for question answering that leverages the structure of Wikipedia to retrieve relevant information. They proposed a novel approach that learns to navigate the Wikipedia graph to find paths connecting relevant pieces of information, improving the accuracy of question answering systems.

Another influential paper by Asai is "Unsupervised Learning of Probably Symmetric Deformable 3D Objects from Images in the Wild" (2019), where she and her co-authors presented a method for unsupervised learning of 3D object representations from 2D images. This work is significant because it addresses the challenge of learning 3D object representations without requiring explicit supervision, which is crucial for many computer vision tasks.

Core limitations of parametric LMs

Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

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

Language model



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

Core limitations of parametric LMs

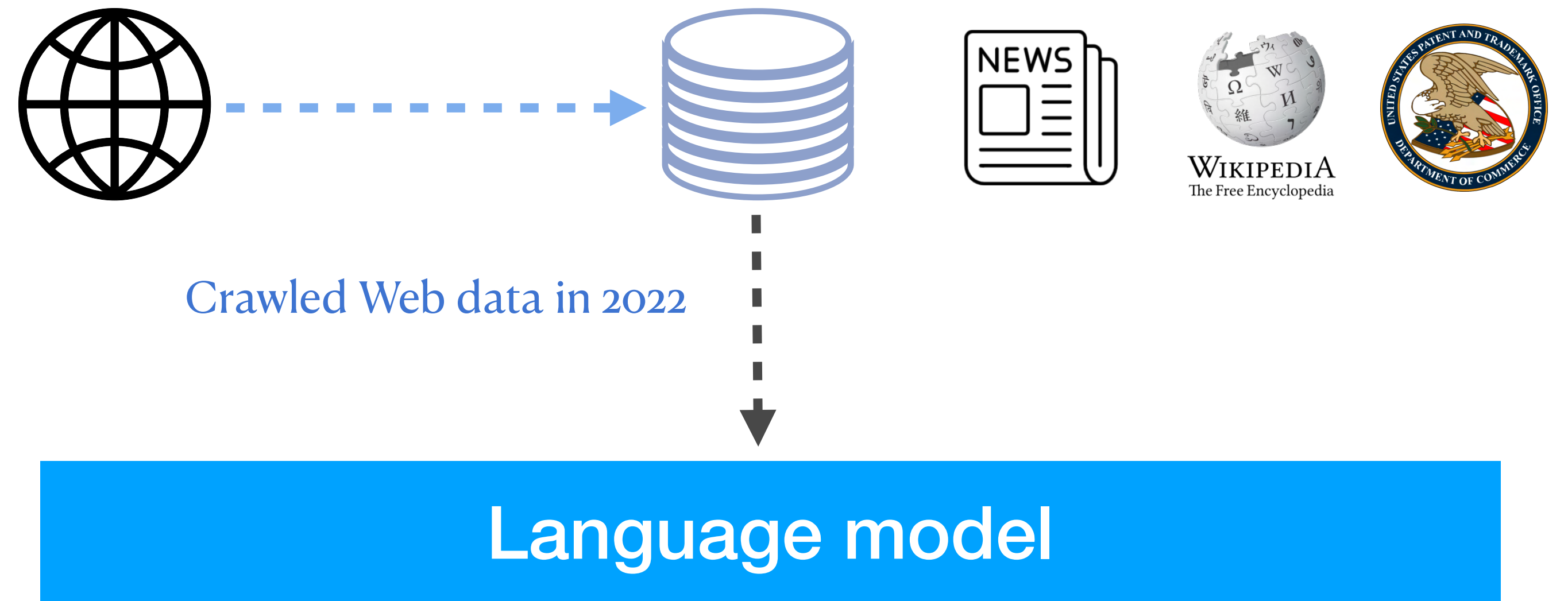
Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



ChatGPT

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

Core limitations of parametric LMs

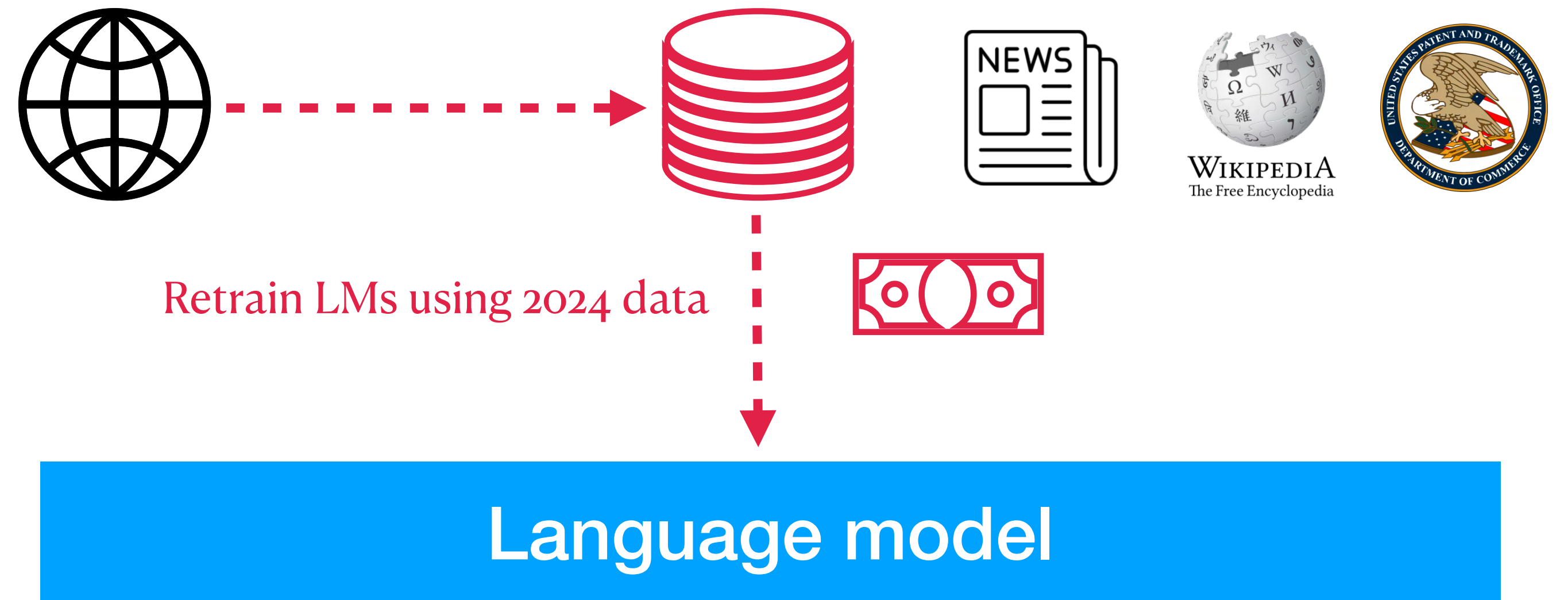
Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



ChatGPT

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

Core limitations of parametric LMs

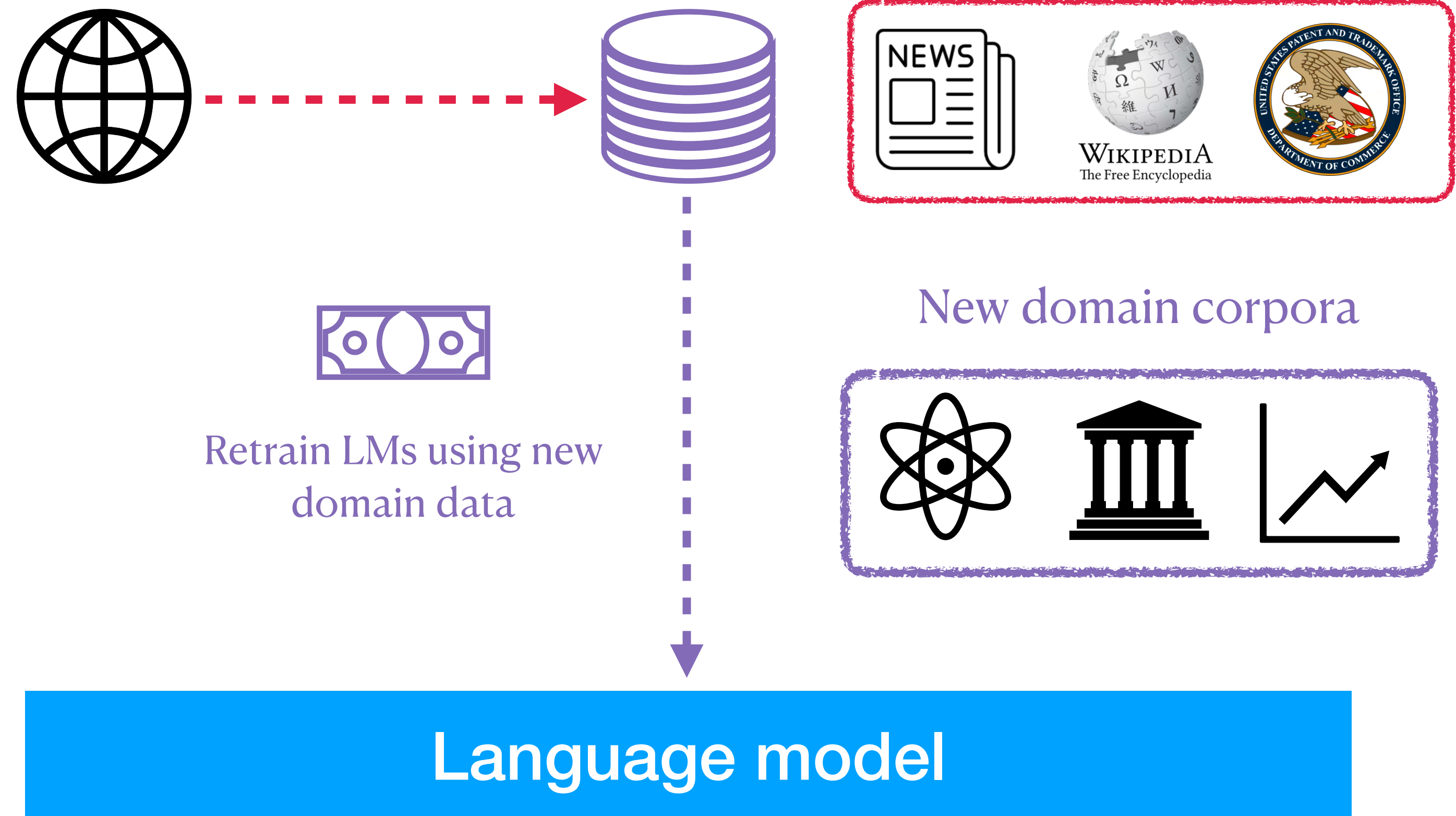
Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



Core limitations of parametric LMs

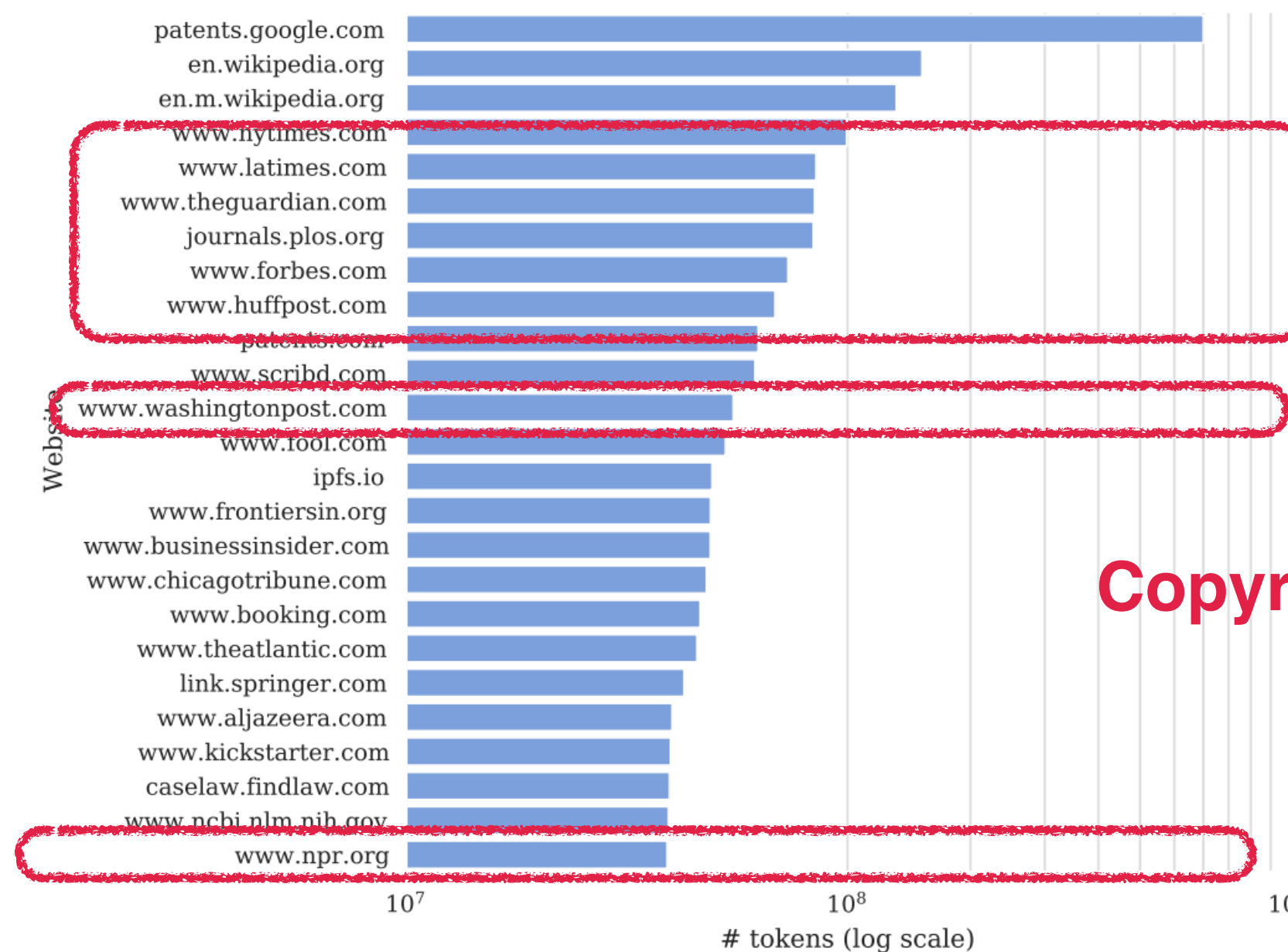
Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size



Copyright-protected data?

Core limitations of parametric LMs

Hallucinations

Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

B. Defendants' GenAI Products

1. *A Business Model Based on Mass Copyright Infringement*

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

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

I. NATURE OF THE ACTION

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

New York Times lawsuits
against OpenAI

Core limitations of parametric LMs

Hallucinations

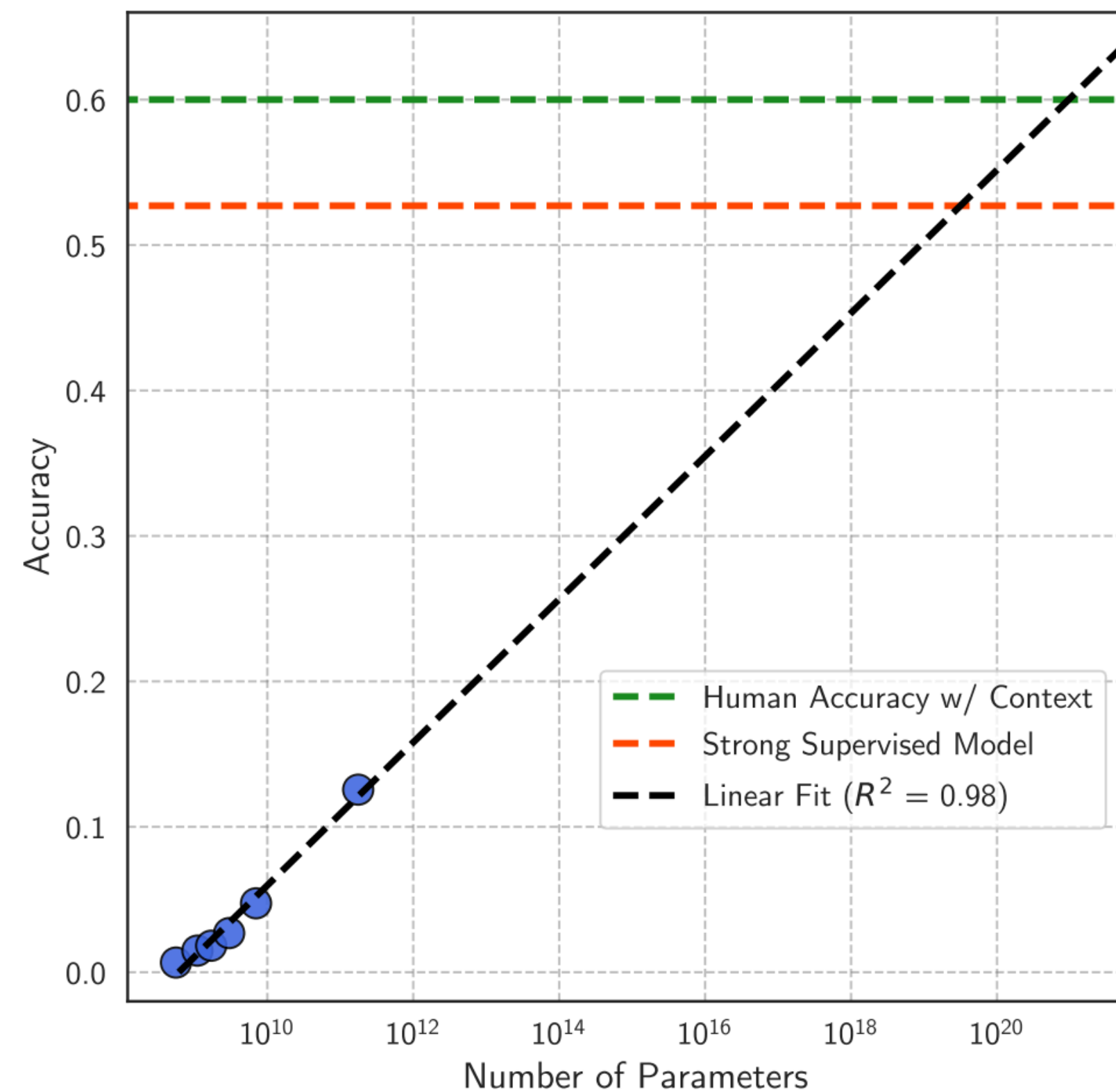
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Massive scaling for better performance



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

How do retrieval-augmented LMs address them?

Hallucinations

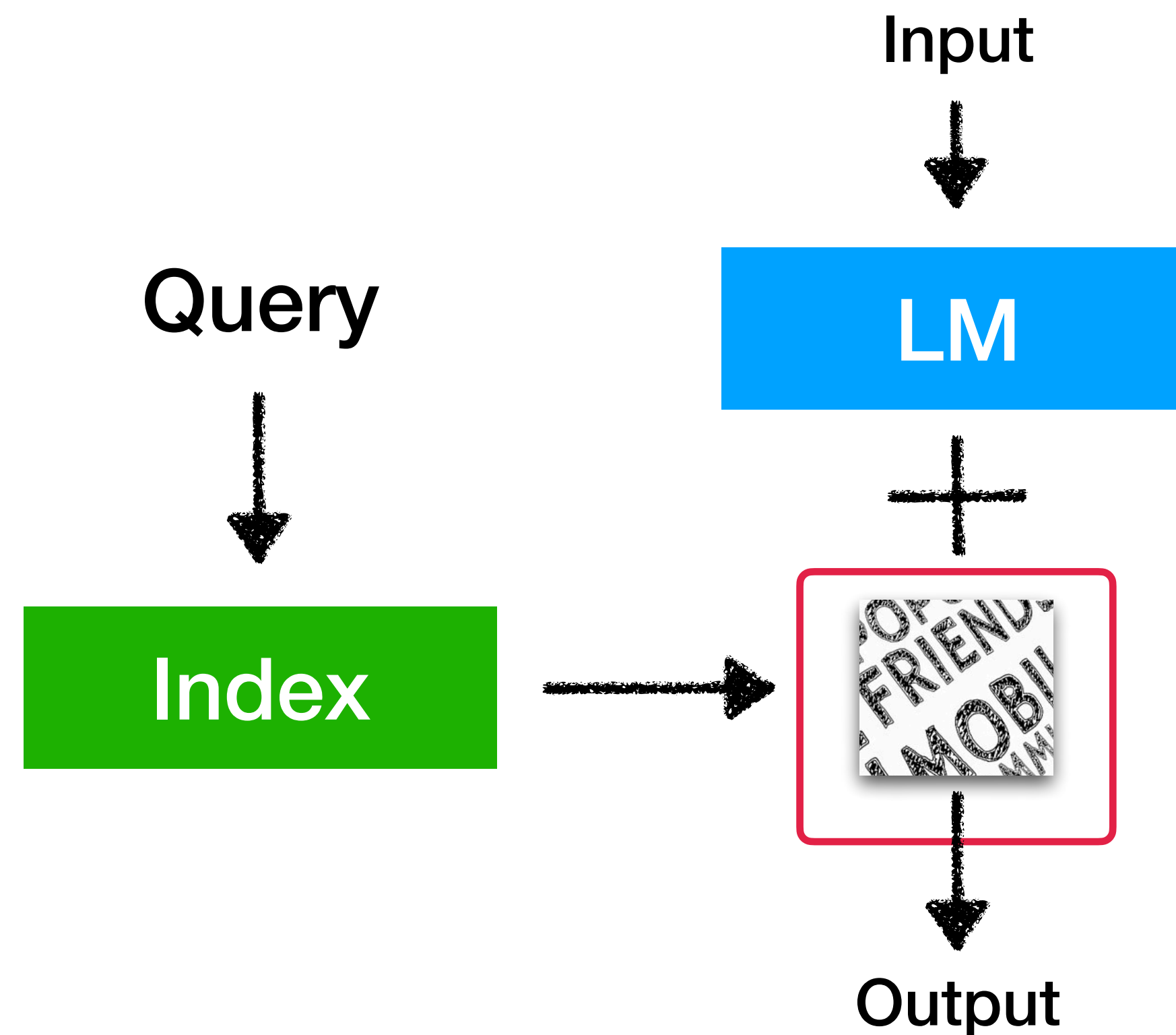
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Retrieved text can be used as attributions



How do retrieval-augmented LMs address them?

Hallucinations

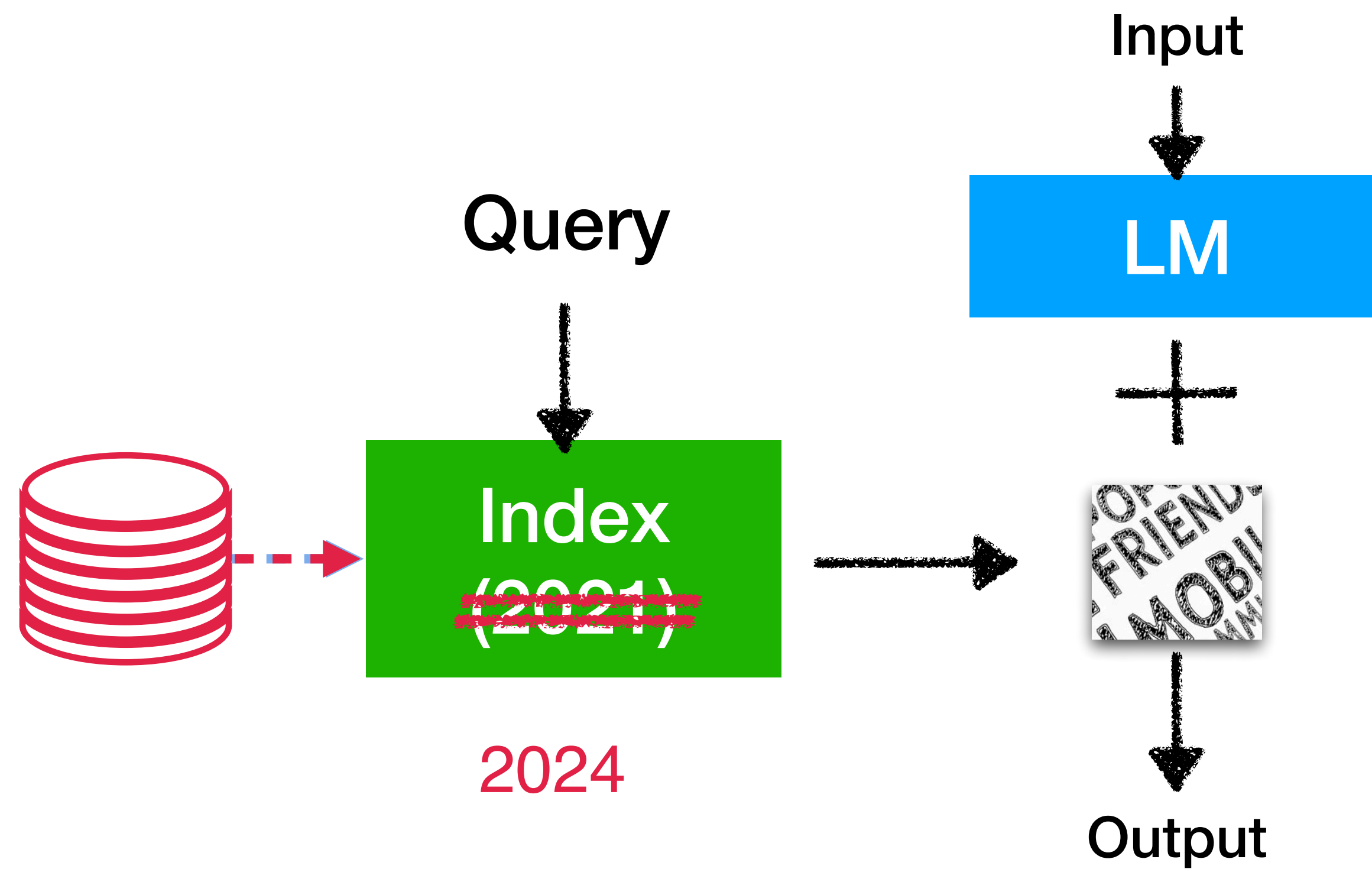
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Replacing datastore (Index) for adaptations without training



How do retrieval-augmented LMs address them?

Hallucinations

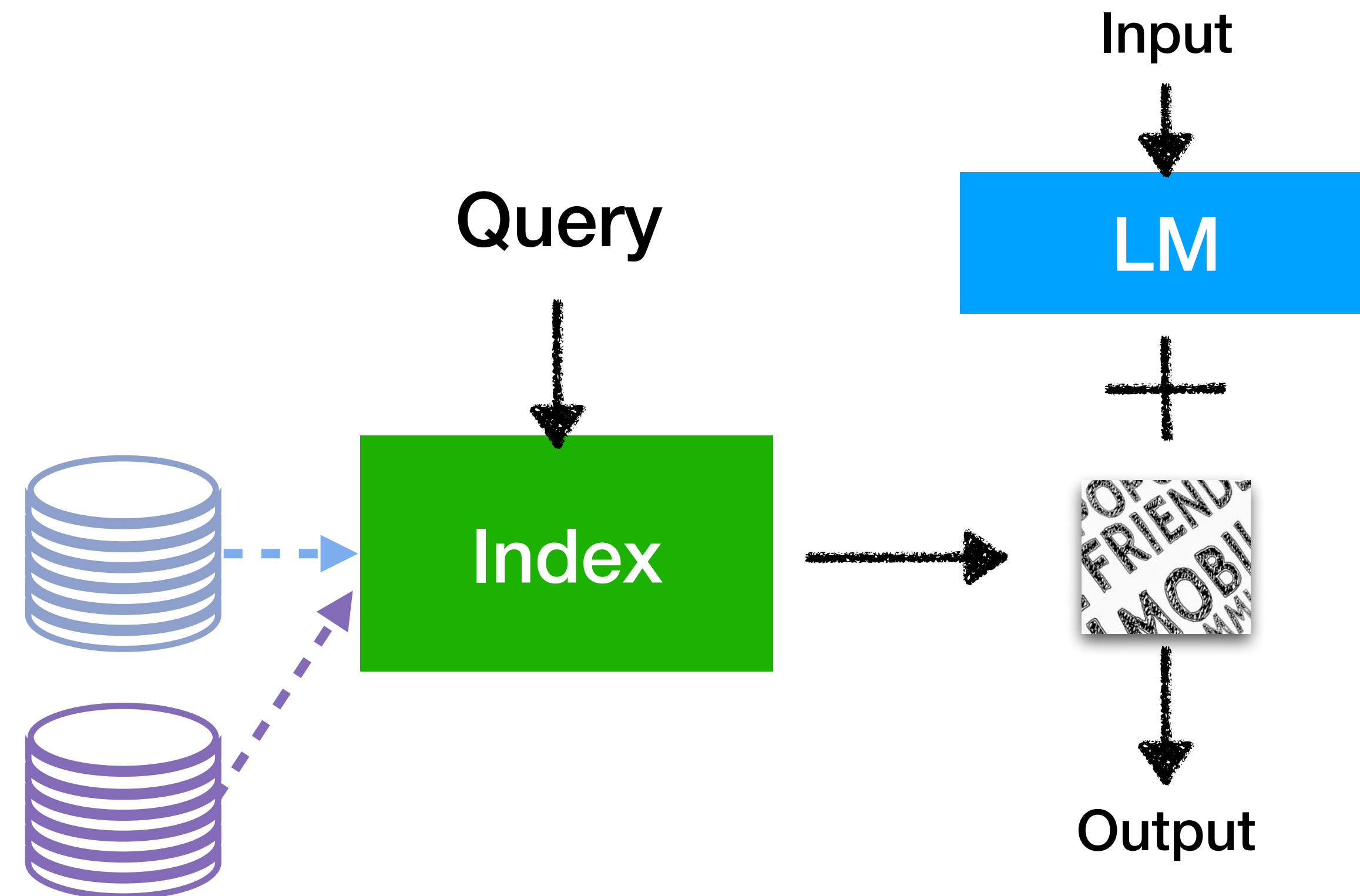
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Adding new domain corpora for domain adaptations



How do retrieval-augmented LMs address them?

Hallucinations

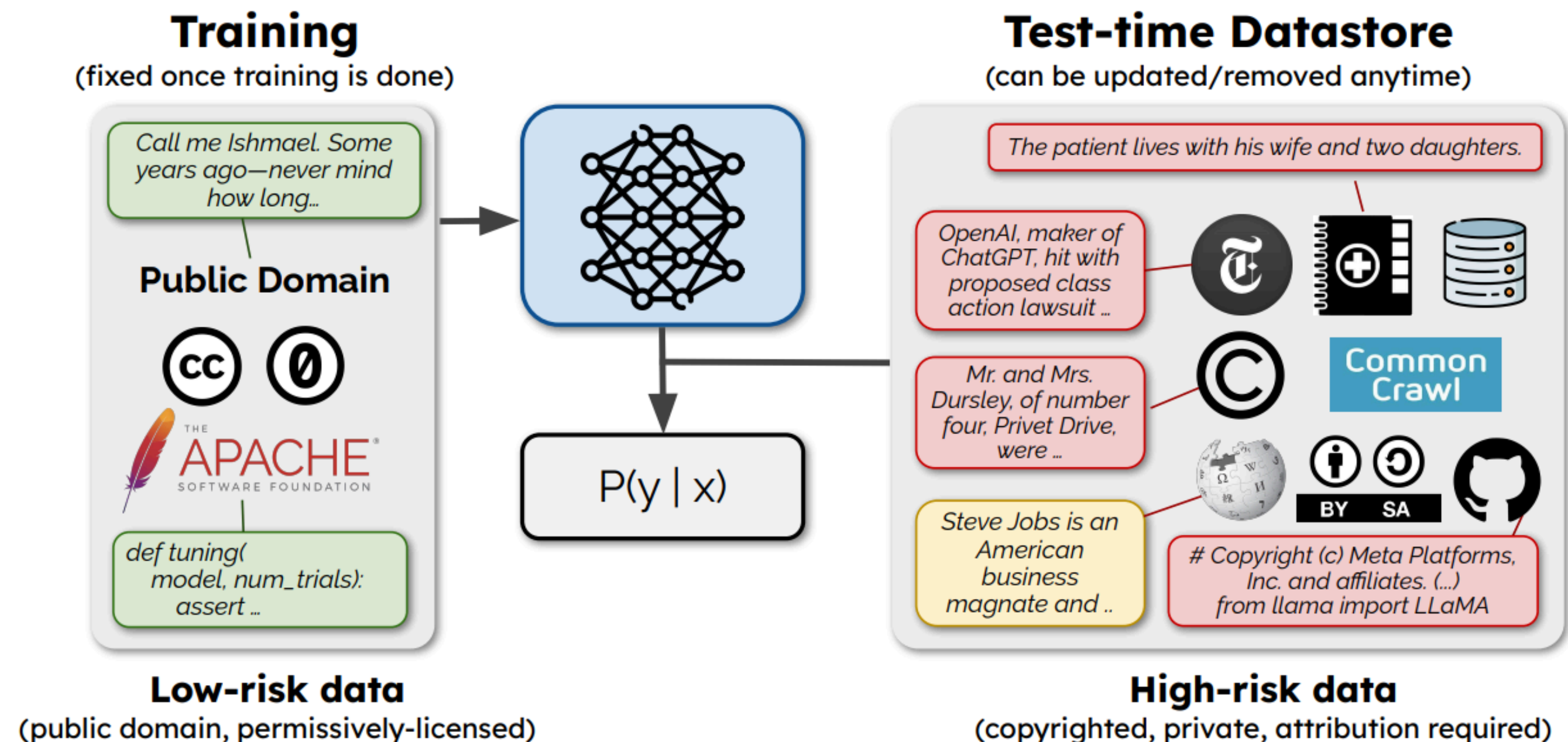
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Segregating copyright-sensitive data from pre-training data



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

How do retrieval-augmented LMs address them?

Hallucinations

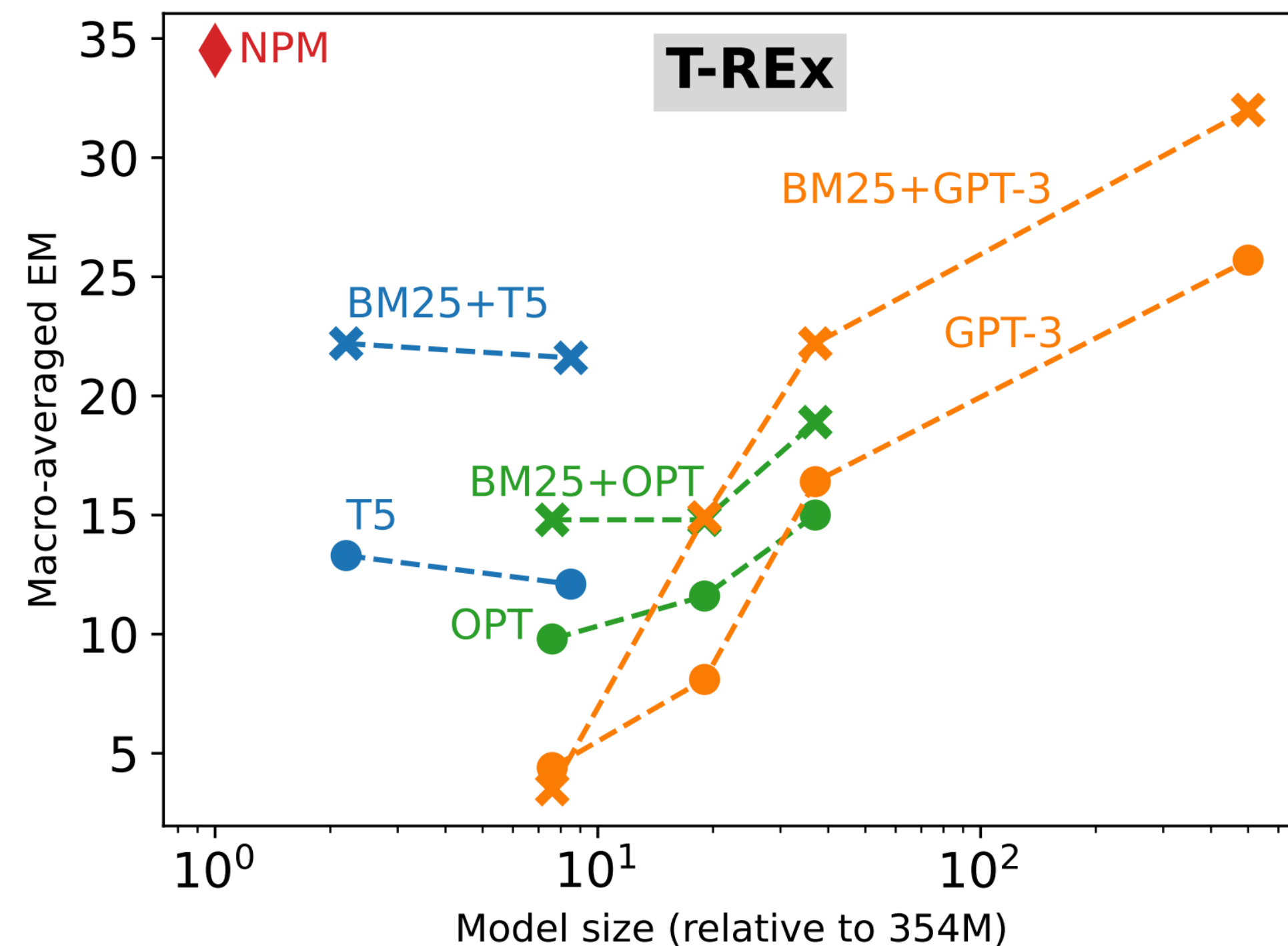
Lack of attributions

Costs of adaptations

Copyright / privacy

Large parameter size

Models with much less parameters can outperforms much larger models!



Promise and limitations of retrieval-augmented LMs

- ✓ Parametric LMs have numerous challenges to build reliable systems.
- ✓ Retrieval-augmented LMs such as RAG can effectively address them.

**When Not to Trust Language Models:
Investigating Effectiveness of Parametric
and Non-Parametric Memories (ACL 2023)**

<https://arxiv.org/abs/2403.03187>

**Reliable, adaptable and
attributable LMs with retrieval
(Arxiv 2024)**

<https://arxiv.org/abs/2403.03187>

Let's talk about how we can improve traditional RAG!

Today's lecture

Promises and Limitations of Retrieval-augmented LMs

Reliable inference: Self-reflective RAG with dynamic retrieval

Versatile Retriever: Intent-aware retrievers with LMs

Summary and Future directions: RAG in the wild

RAG hurts in popular knowledge memorized by vanilla LMs

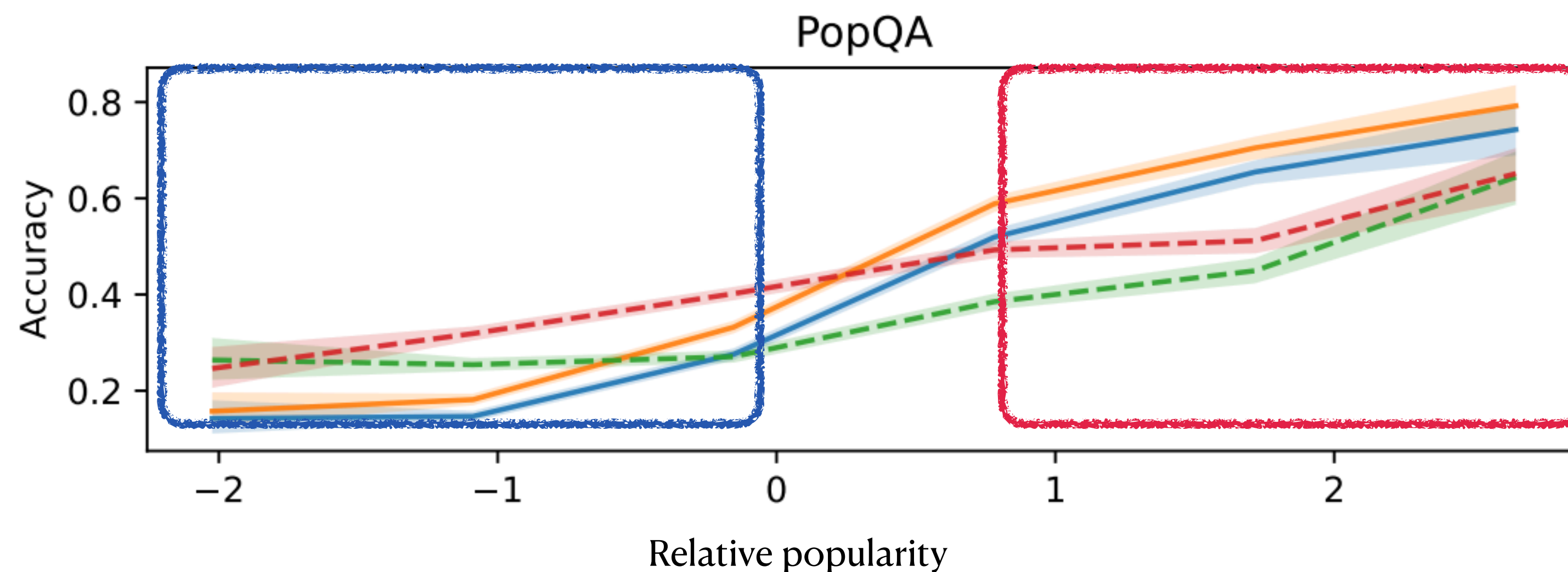
On popular facts

vanilla \approx RAG or vanilla $>$ RAG

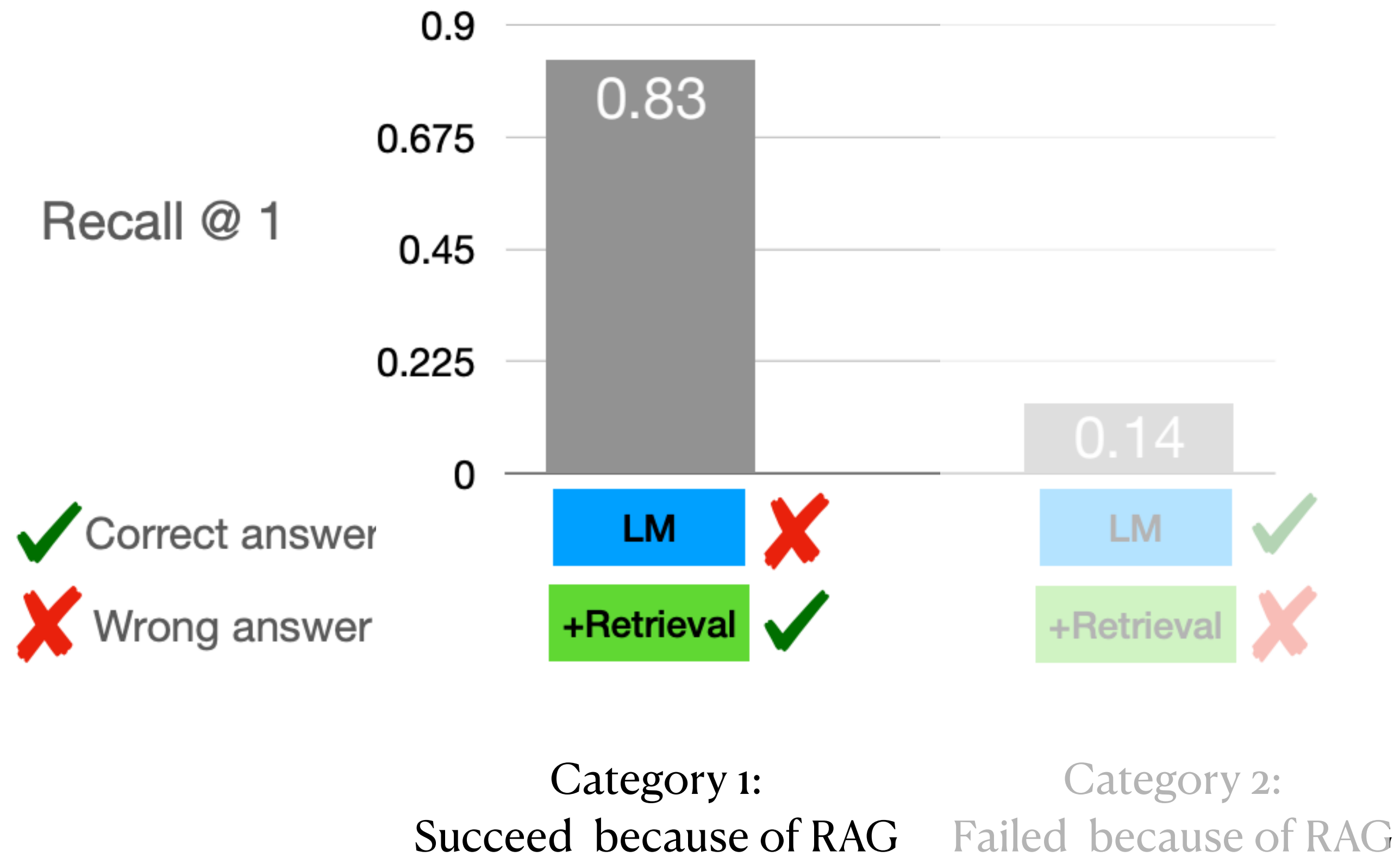
Why?

On long-tail facts

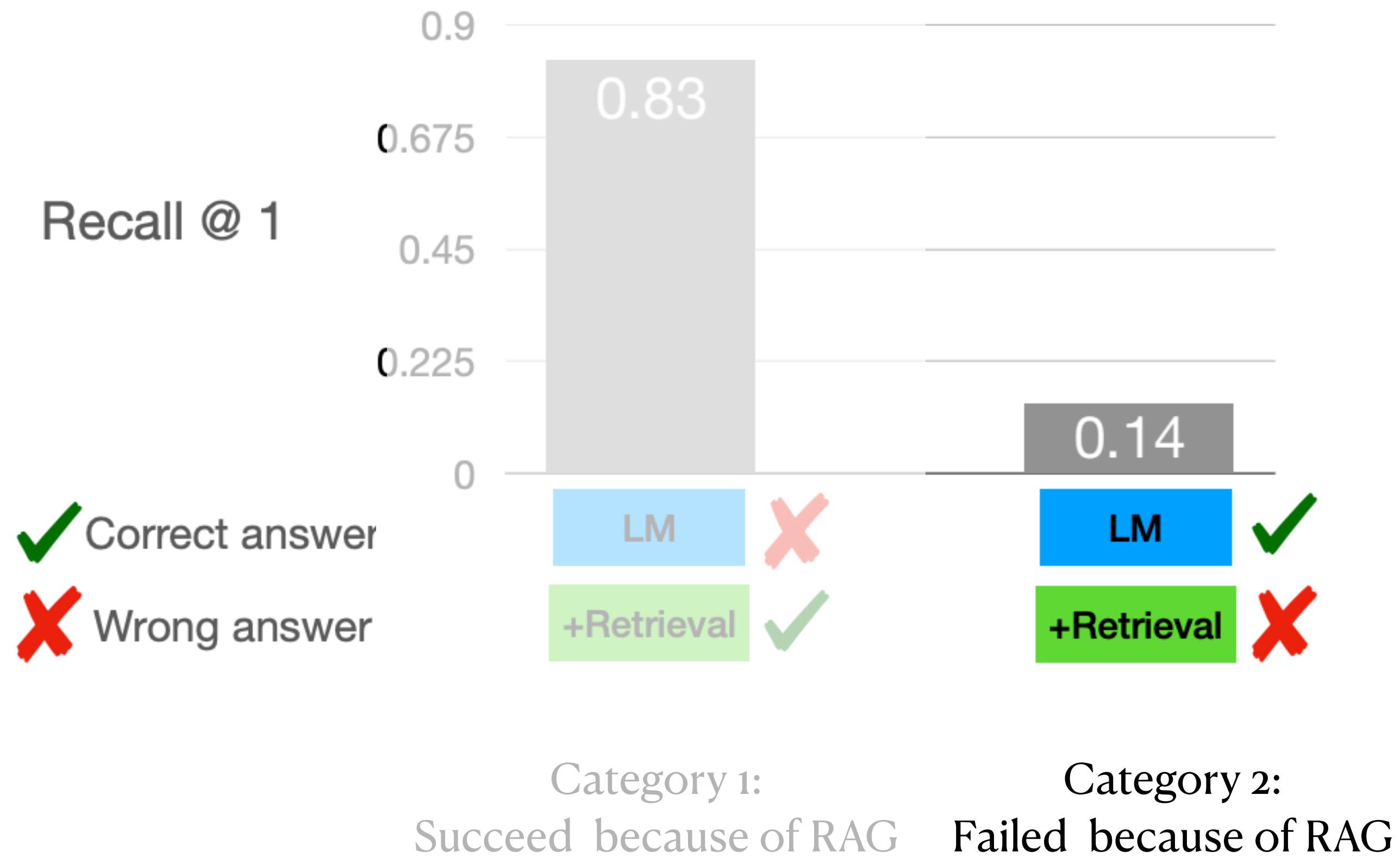
RAG $>$ vanilla



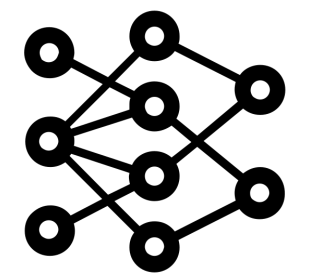
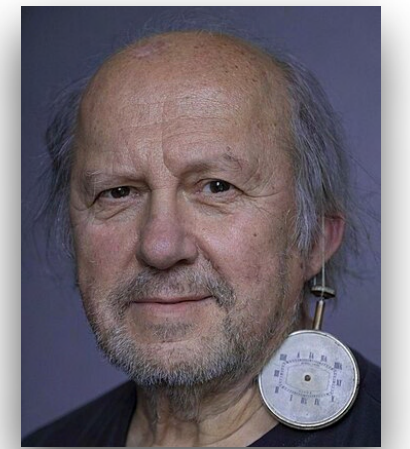
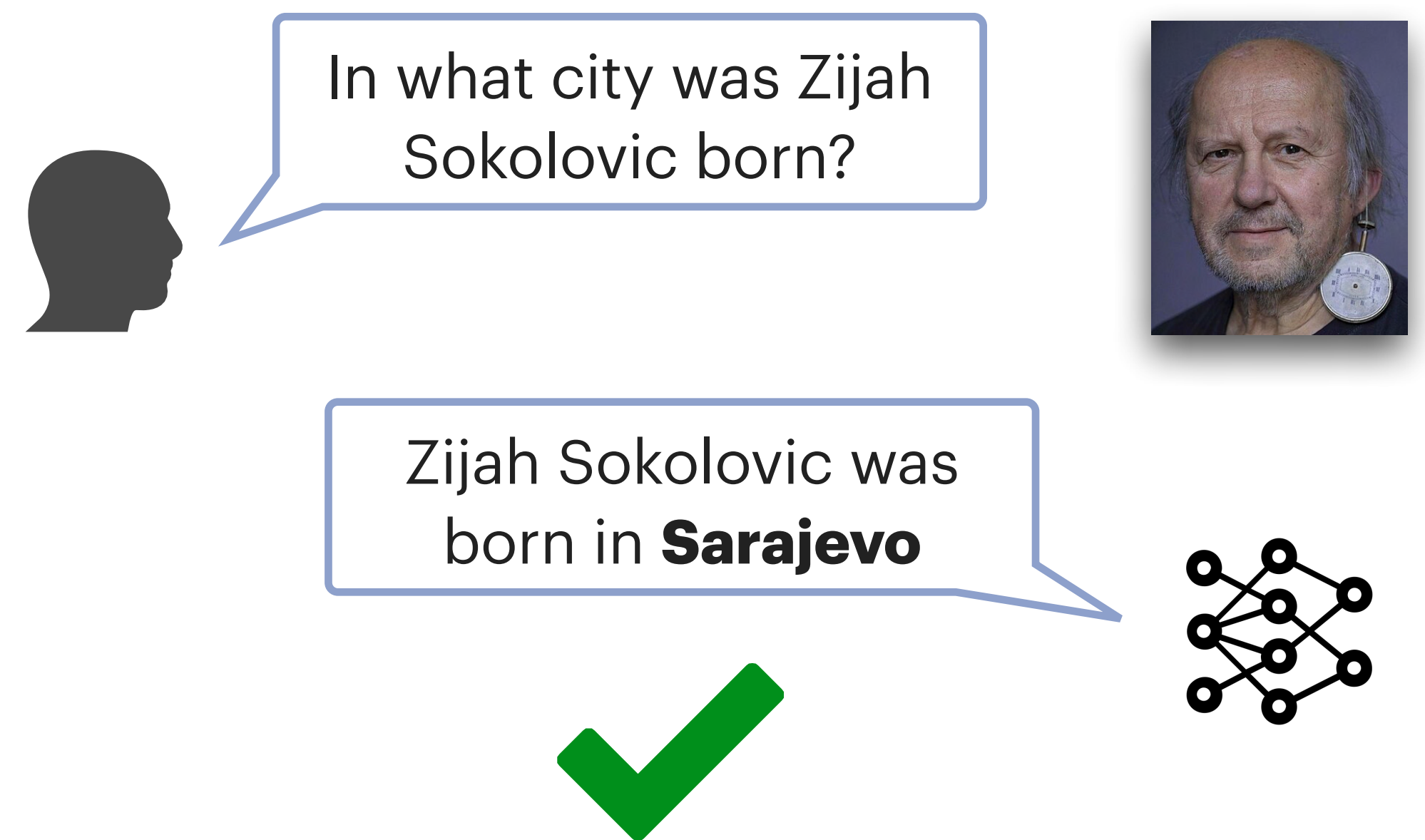
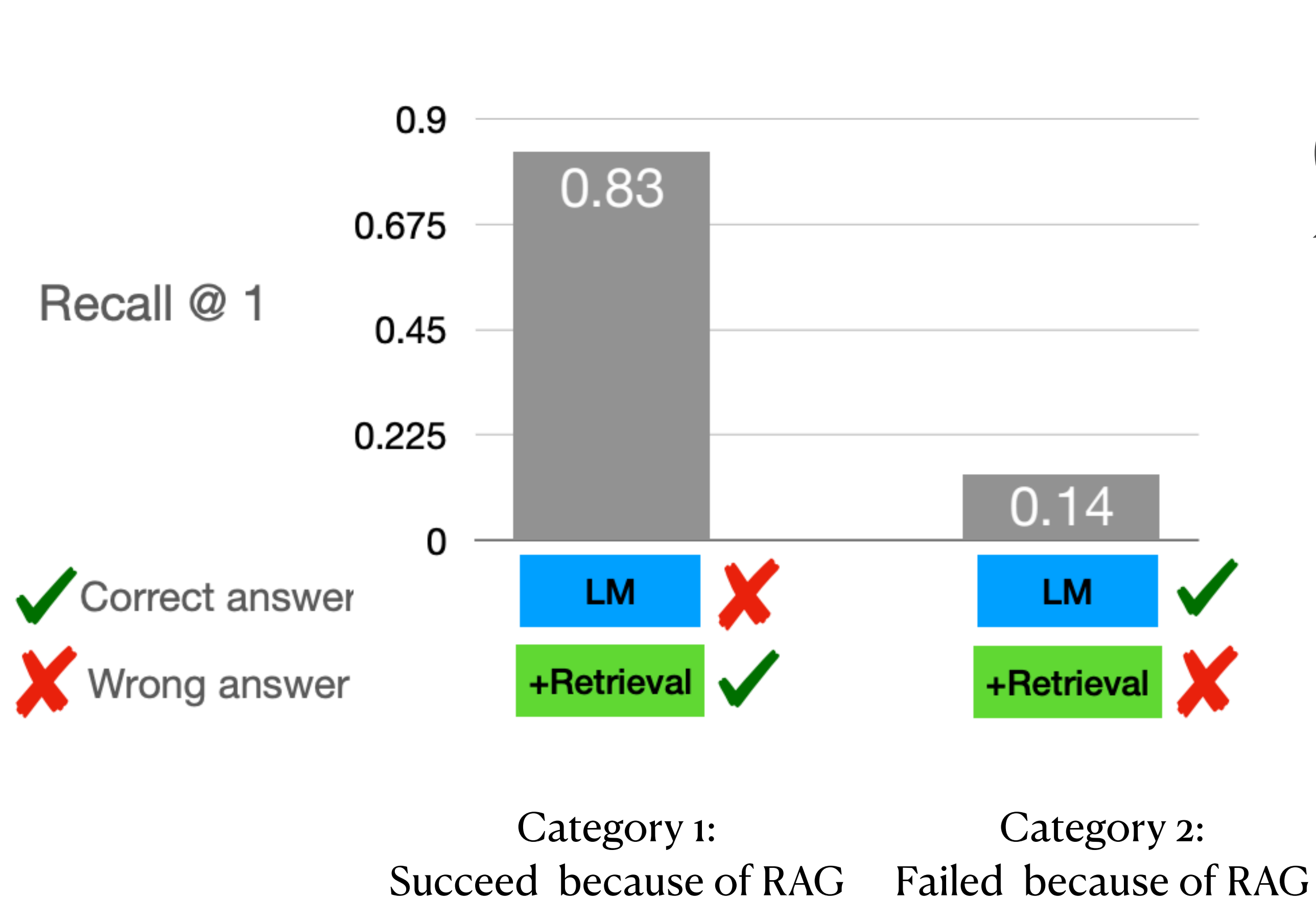
Incorrect retrieval can easily confuse LMs



Incorrect retrieval can easily confuse LMs

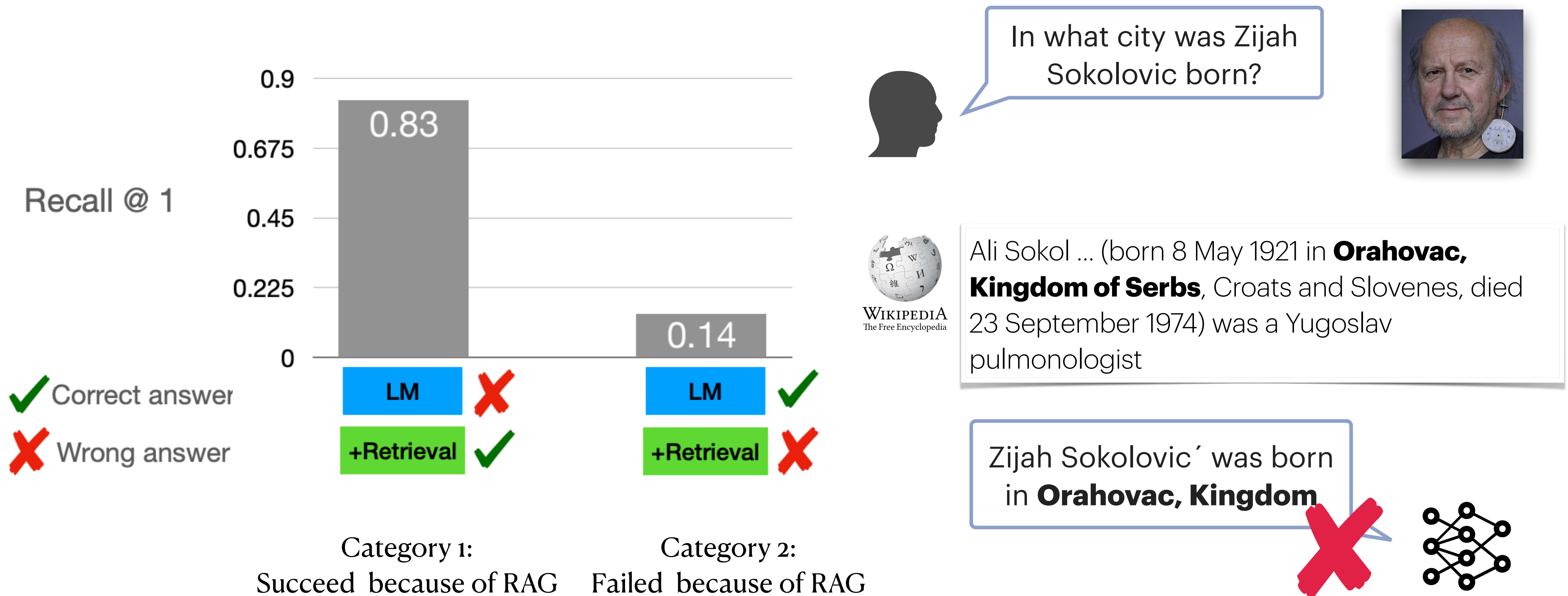


Incorrect retrieval can easily confuse LMs



Incorrect retrieval can easily confuse LMs

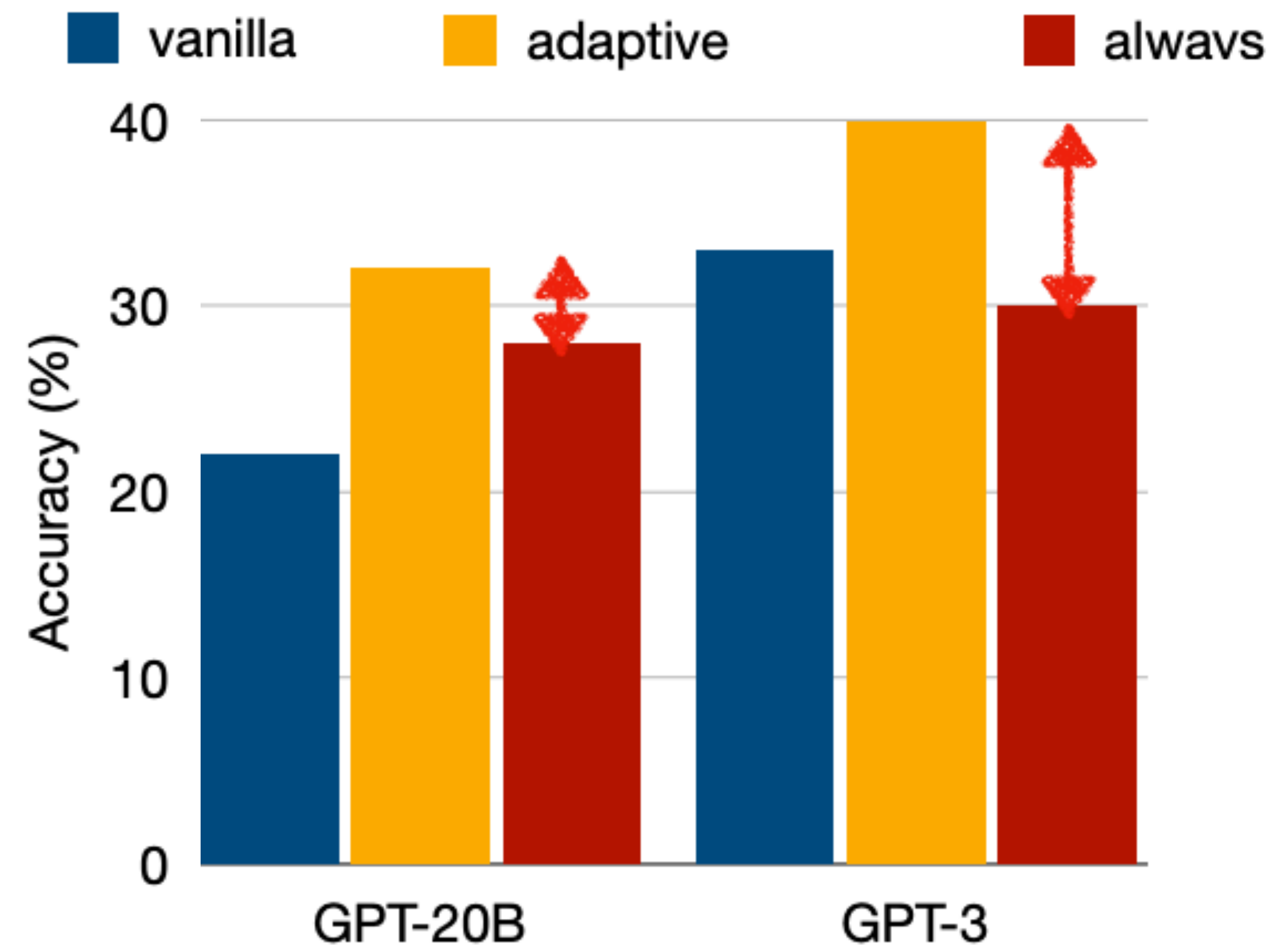
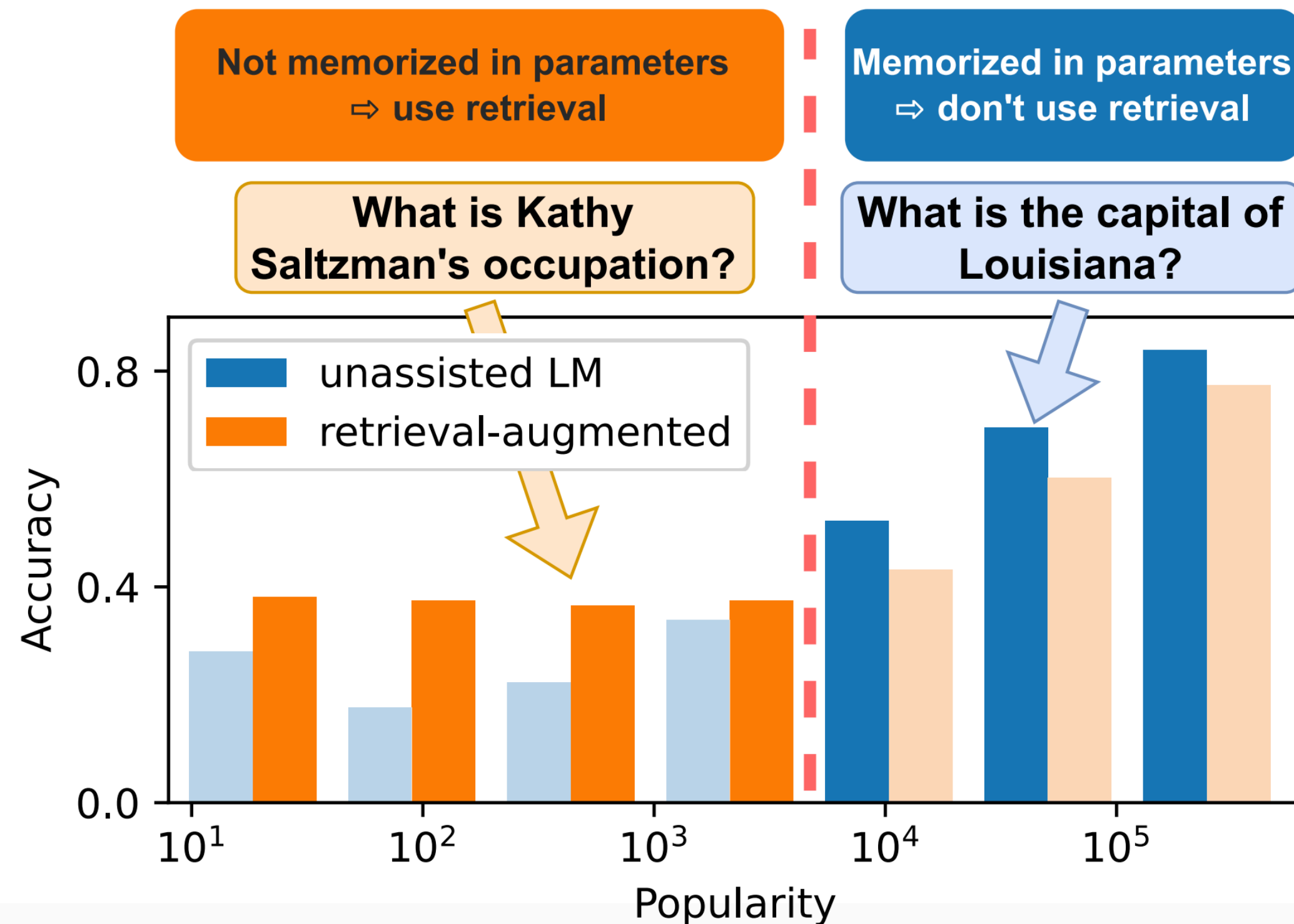
Failed RAG can counterfactually make LMs answer incorrectly



A simple hack: a threshold-based **adaptive** retrieval

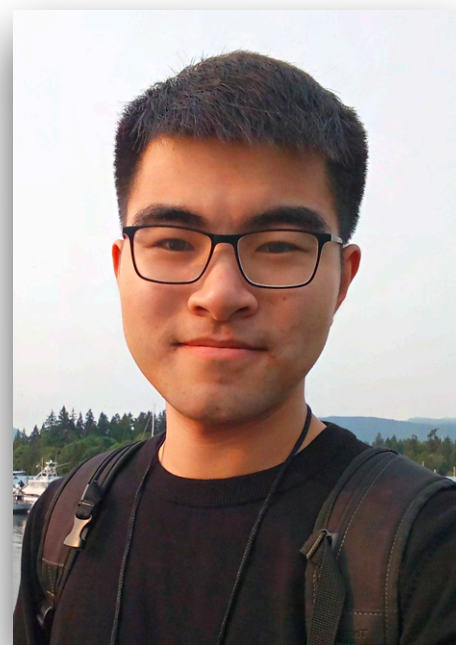
Simple adaptive RAG significantly improve RAG performance & efficiency

Popularity threshold learned on dev set



Self-RAG: **Learning to Retrieve, Generate and Critique through Self-Reflections**

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirul Sil, Hannaneh Hajishirzi



ICLR 2024 (Oral — top 1%)

Best Paper Honorable Mention at NeurIPS Instruction workshop

Standard RAG aren't (always) reliable

Step 1: Retrieve K documents

Prompt How did US states get their names?



❶ Of the fifty states, eleven including New York, Georgia, Washington named after an individual person.

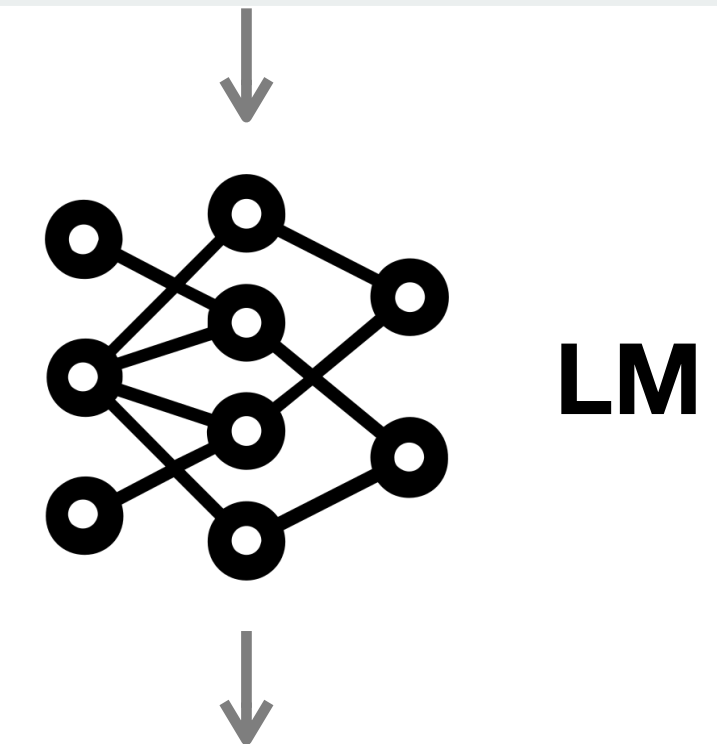
❷ UTAH: Name taken from the Ute Indians who inhabited that region

❸ The history of human activity in Michigan began with settlement by Paleo-Indians.

Easily distracted by unhelpful documents

Step 2: Prompt LM with K docs and generate

Prompt How did US states get their names? ❶ ❷ ❸



US states got their names from a variety of sources. **Eleven states such as New York and Michigan are named after an individual person.** Some states including Utah and Washington, are named after native American tribe names.

Standard RAG aren't (always) reliable

Step 1: Retrieve K documents

Prompt How did US states get their names?

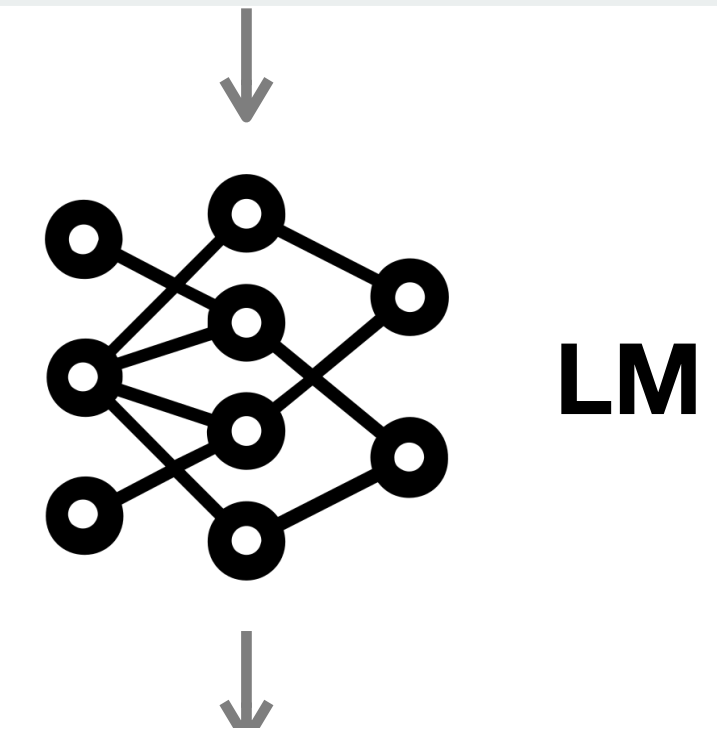


- 1 Of the fifty states, eleven including New York, Georgia, Washington named after an individual person.
- 2 UTAH: Name taken from the Ute Indians who inhabited that region
- 3 The history of human activity in Michigan began with settlement by Paleo-Indians.

LMs can also ignore retrieved context

Step 2: Prompt LM with K docs and generate

Prompt How did US states get their names? 1 2 3



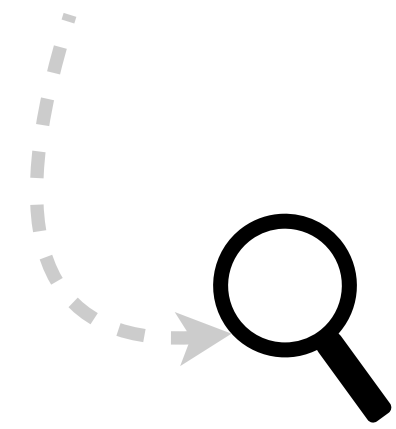
US states got their names from a variety of sources. Eleven states such as New York and Michigan are named after an individual person. **Some states including Utah and Washington**, are named after native American tribe names.



Standard RAG aren't (always) necessary

Step 1: Retrieve K documents

Write an essay of your best summer vacation



Retriever

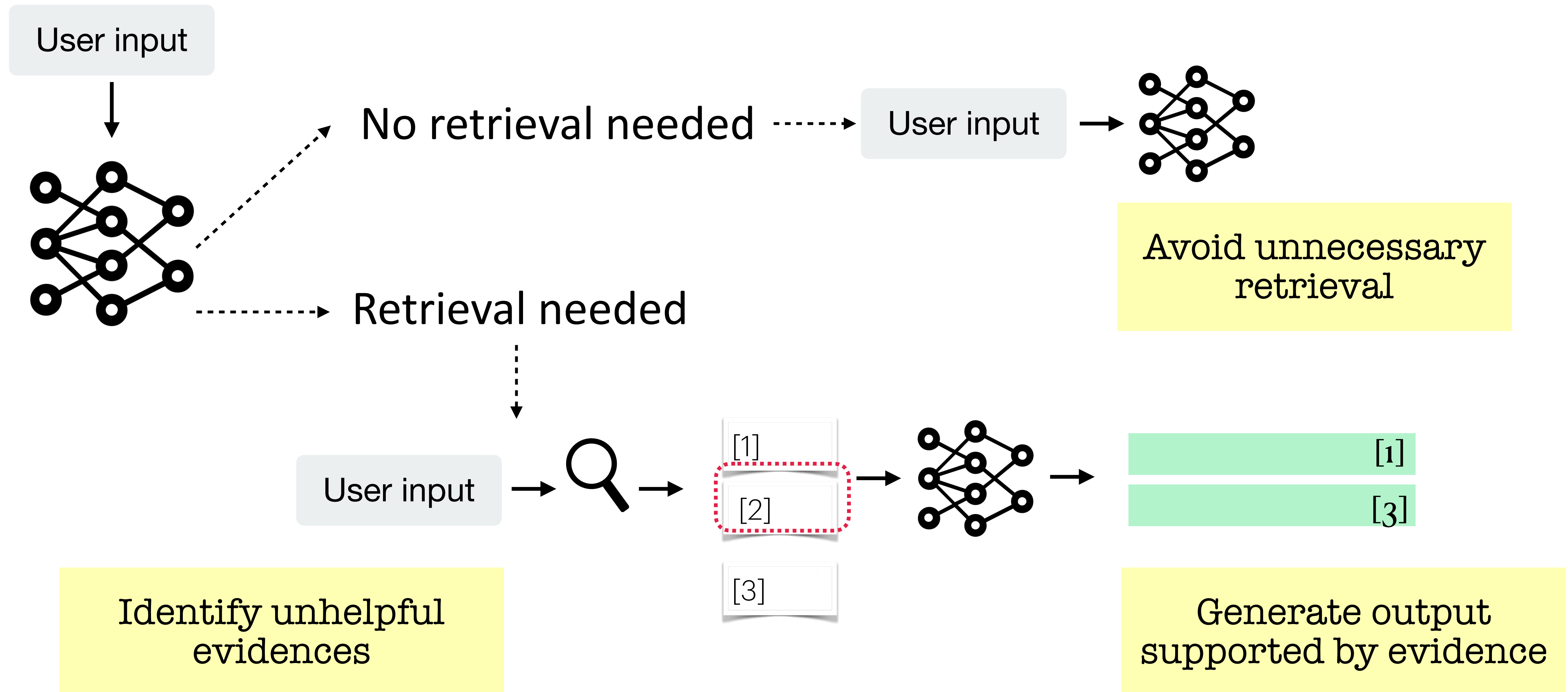
1 The term summer vacation or summer break refers to a school break in the summer between school years.

2 Summer Vacation (Korean: 여름방학) is a South Korean "home vacation" reality show

3 Summer Vacations a 2023 Spanish comedy film directed by Santiago Segura which stars Segura, Leo Harlem, Cristina Gallego, and Patricia Conde.

Always retrieving fixed number of documents is inefficient & harmful

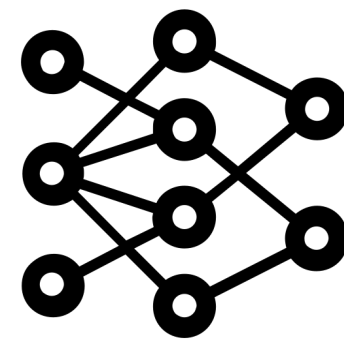
Goal of **Self-RAG**: learn to retrieve, generate and critique



Self-RAG – Self-Reflective Retrieval-Augmented Generation

Step 1: Retrieve documents on demand

How did US states get their names?



US states got their names from a variety of sources.

Retrieve

Reflection tokens



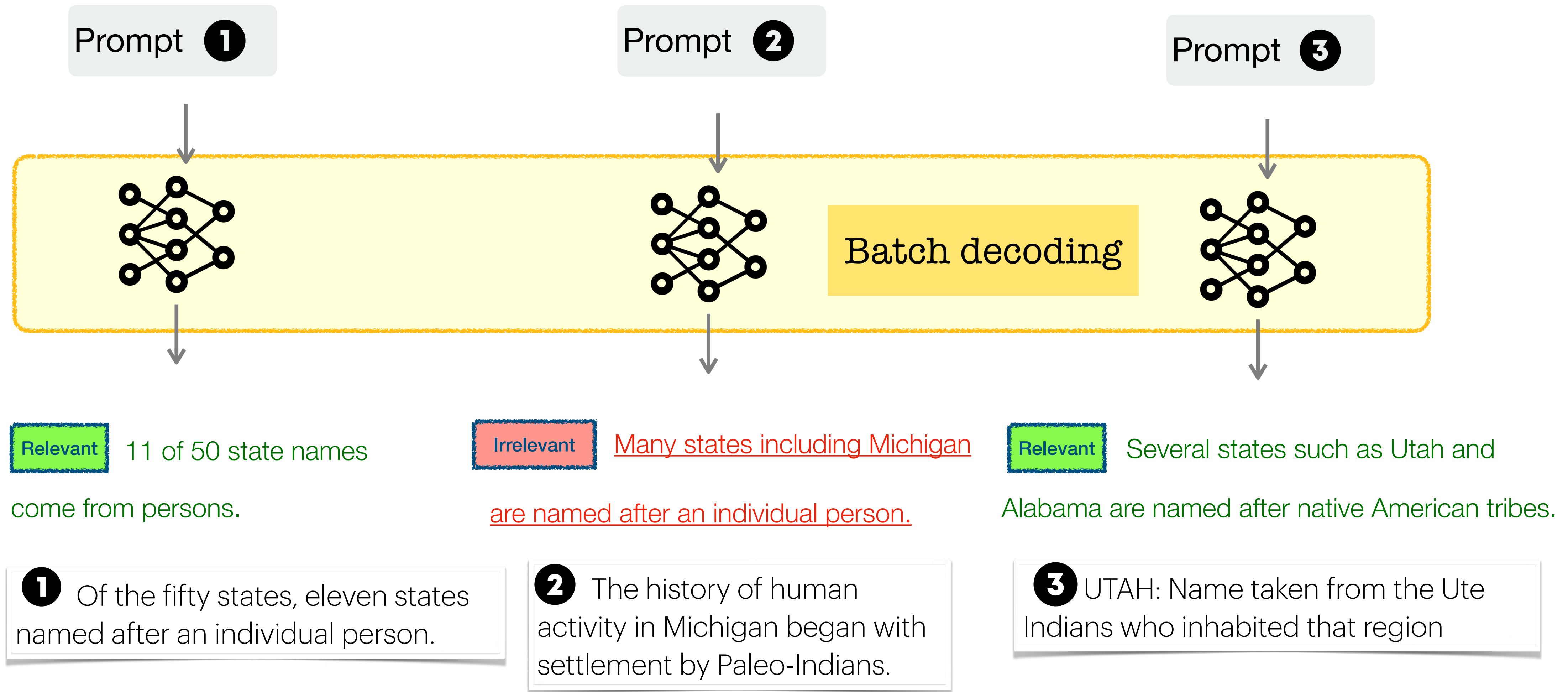
1 Of the fifty states, eleven including New York, Georgia, Washington named after an individual person.

2 The history of human activity in Michigan began with settlement by Paleo-Indians.

3 UTAH: Name taken from the Ute Indians who inhabited that region

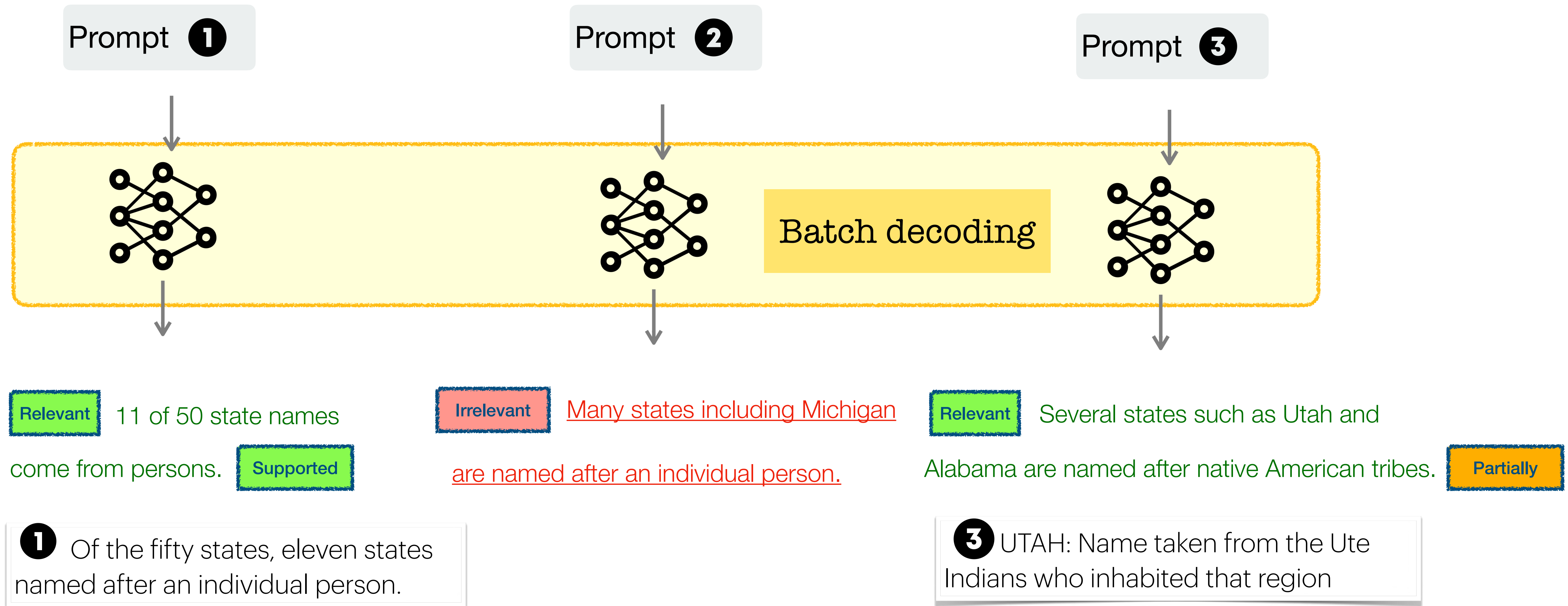
Self-RAG – Self-Reflective Retrieval-Augmented Generation

Step 2: Generate segments in *parallel*



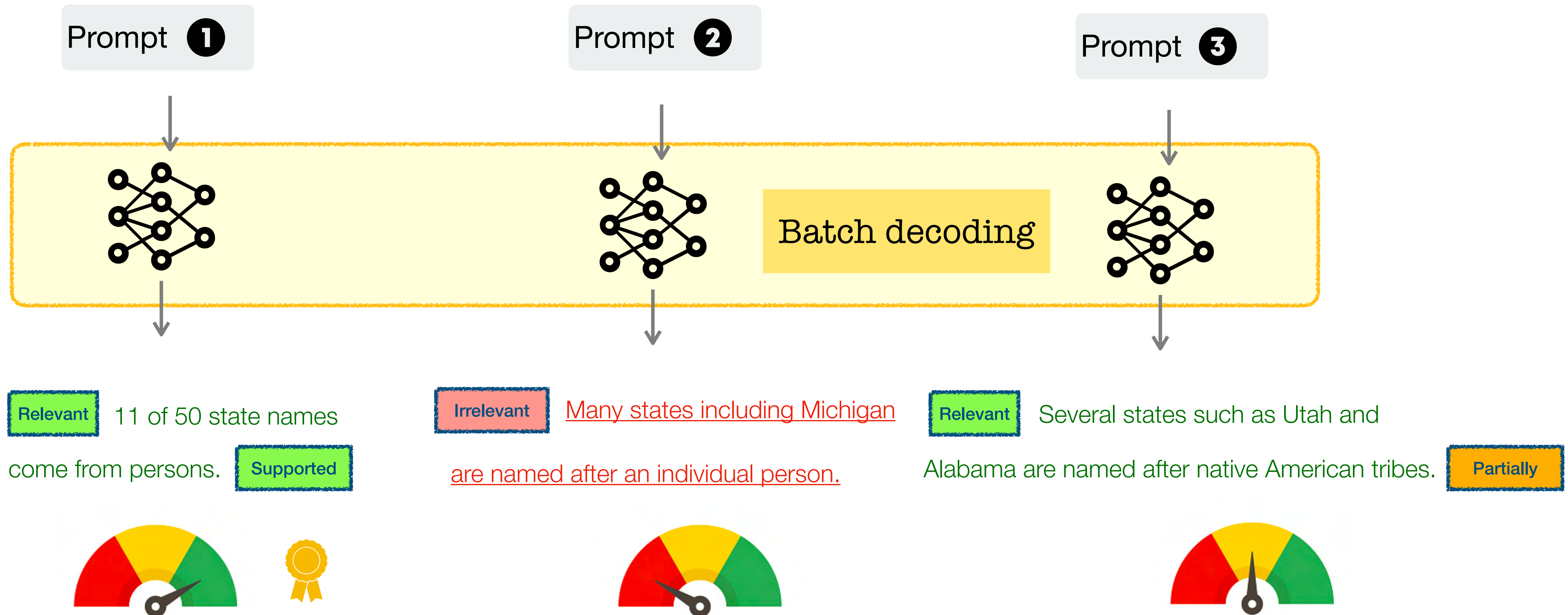
Self-RAG – Self-Reflective Retrieval-Augmented Generation

Step 2: Generate segments in *parallel*



Self-RAG – Self-Reflective Retrieval-Augmented Generation

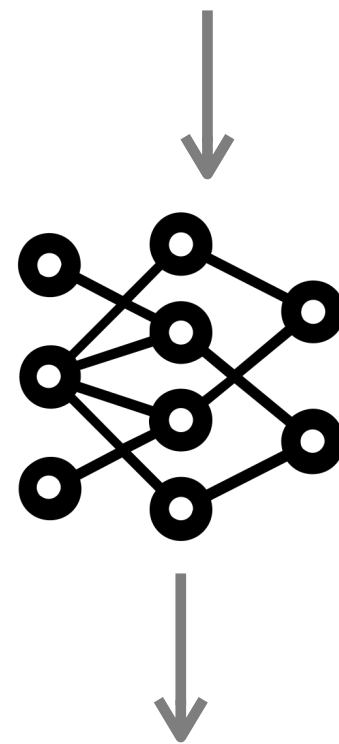
Step 2: Generate segments in *parallel*



Self-RAG – Self-Reflective Retrieval-Augmented Generation

Step 1: Generate with no retrieval

Write an essay of your best summer vacation



No Retrieval

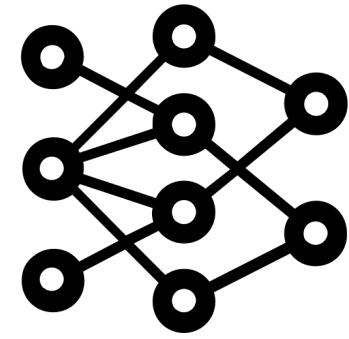
My best summer vacation was a magical escape to the coastal town of Santorini.

The azure waters, charming white-washed buildings are unforgettable.

Util:5

Reflection tokens for retrieval and critique

US states got their names from a variety of sources.



California

e.g.,

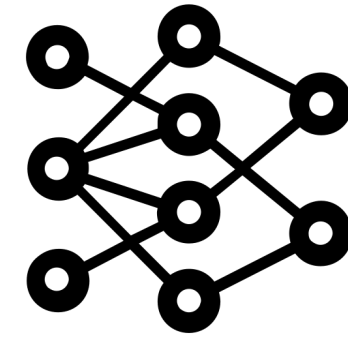
A



Original LM vocabularies

Reflection tokens for retrieval and critique

US states got their names from a variety of sources.



California

e.g.,

A

Original LM vocabularies

Retrieve

No Retrieval

Relevant

Irrelevant

Supported

No support

Useful

Not useful

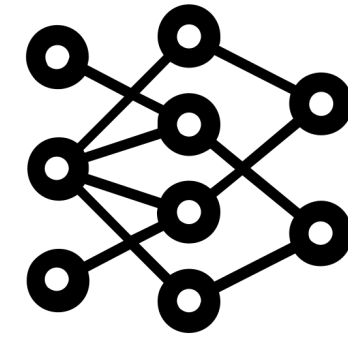
Retrieval tokens

Critique tokens

Vocabulary expanded with reflection tokens

Reflection tokens for retrieval and critique

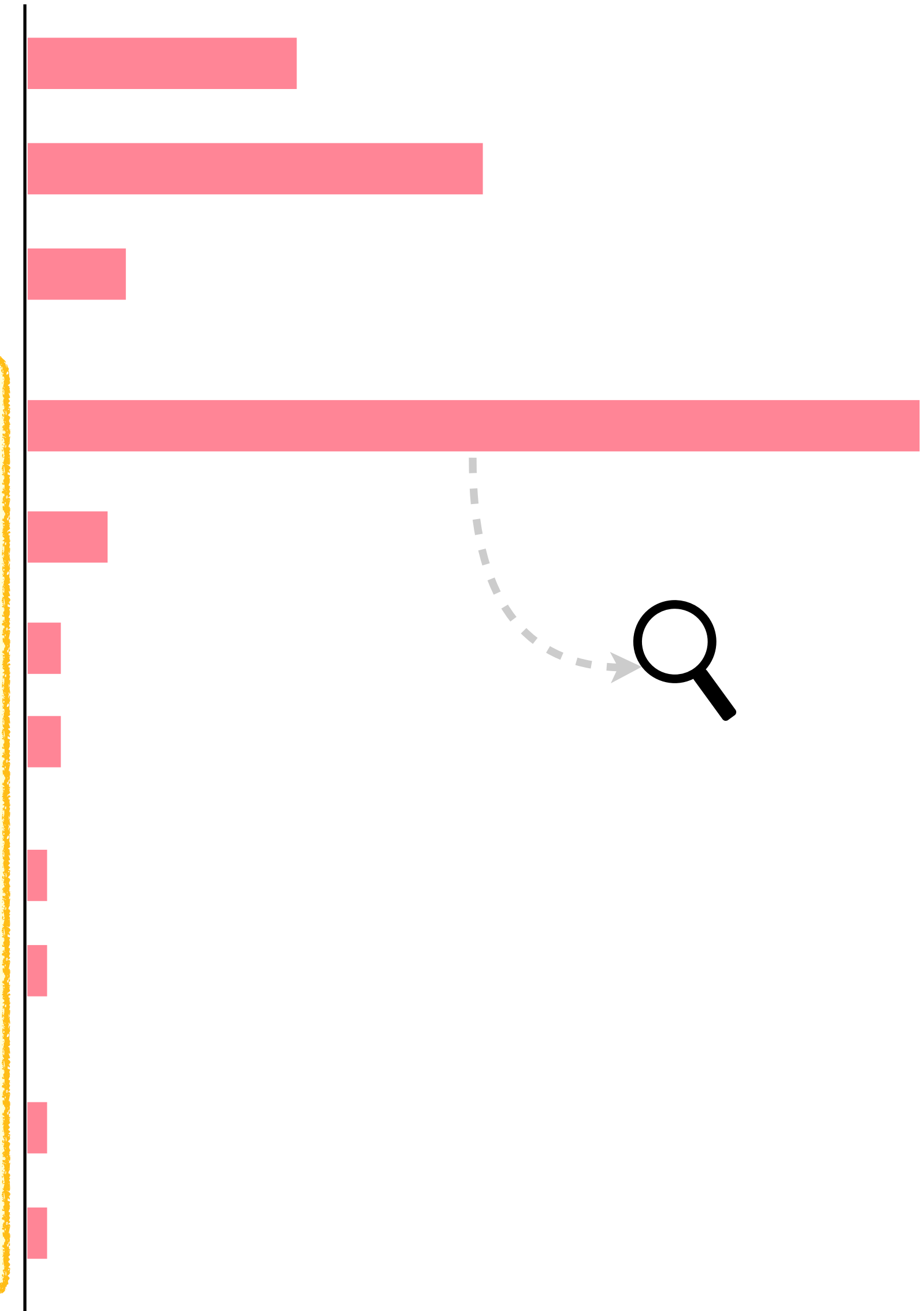
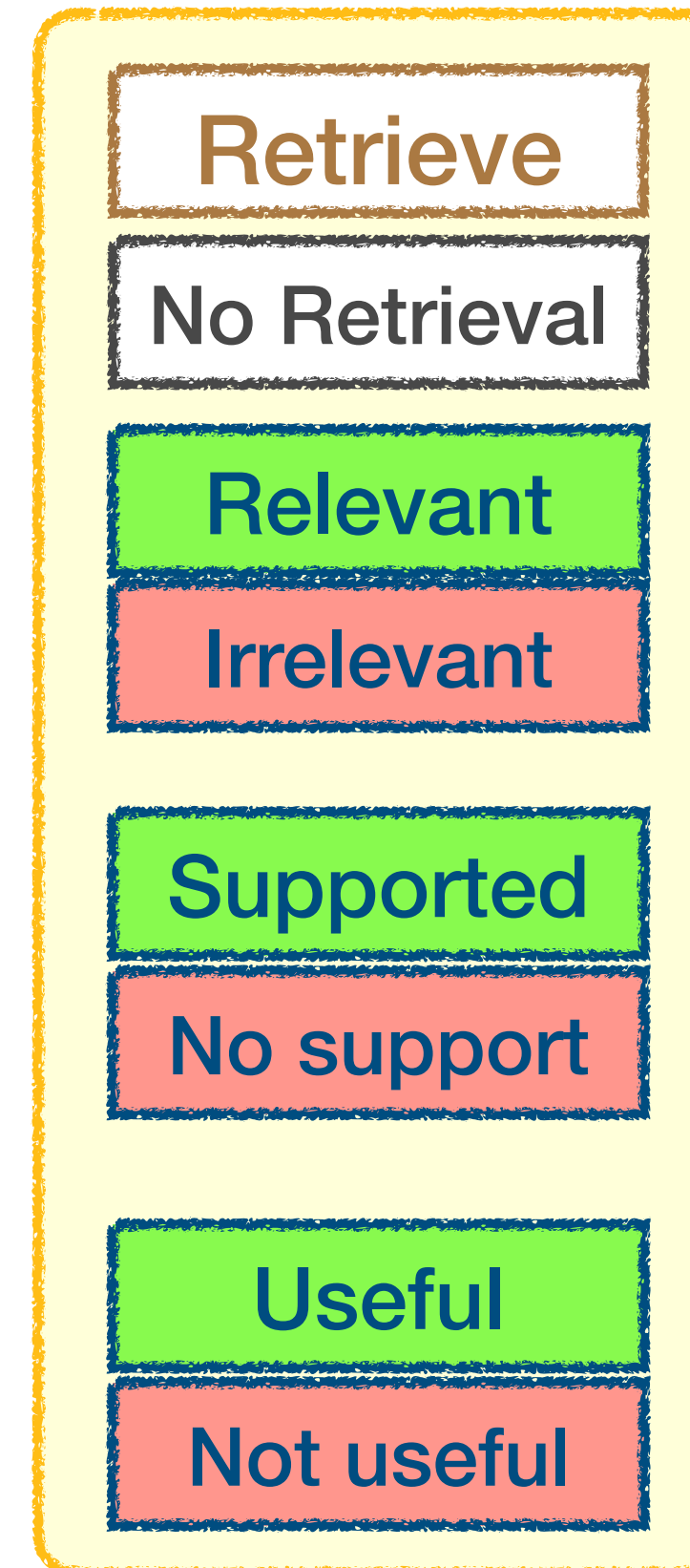
US states got their names from a variety of sources.



California

e.g.,

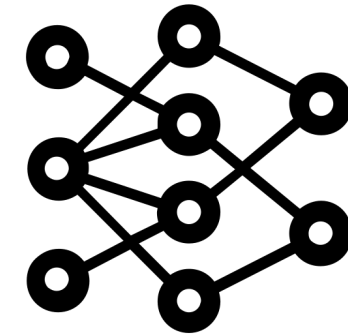
A



Vocabulary expanded with reflection tokens

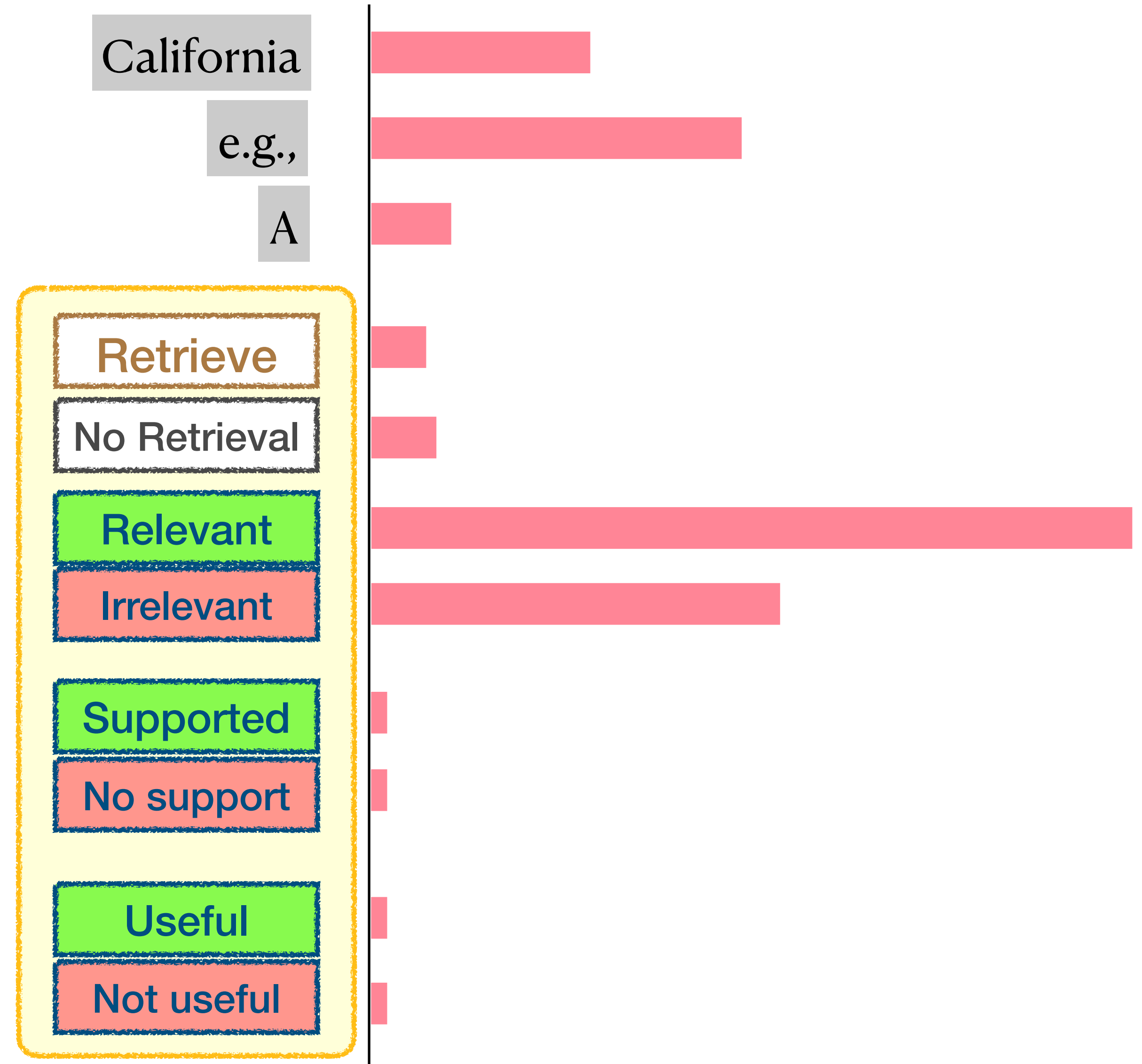
Reflection tokens for retrieval and critique

US states got their names from a variety of sources.



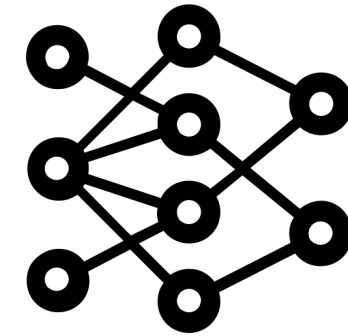
1 Of the fifty states, eleven are named after an individual person.

Vocabulary expanded with reflection tokens



Reflection tokens for retrieval and critique

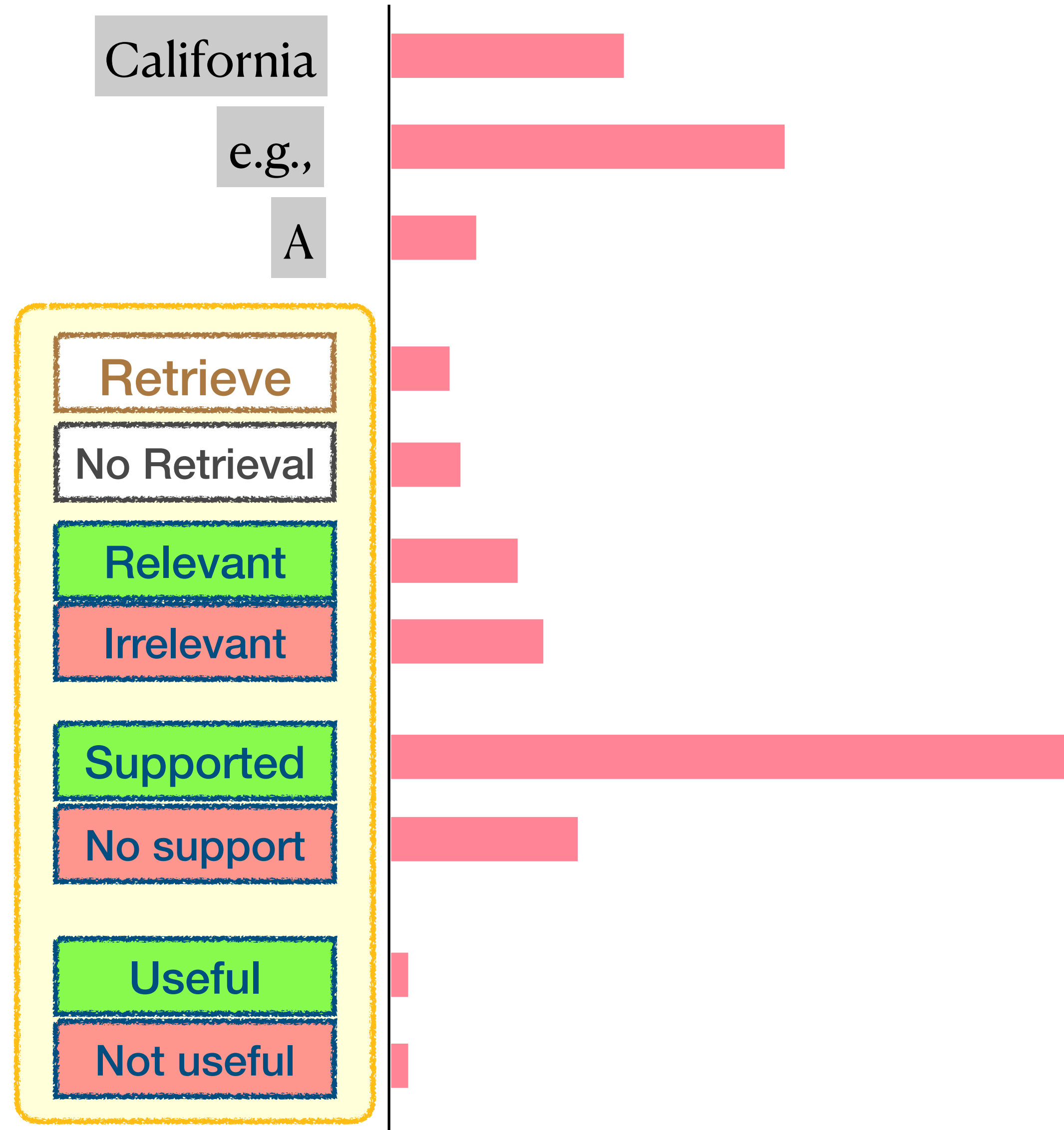
US states got their names from a variety of sources.



1 Of the fifty states, eleven are named after an individual person.

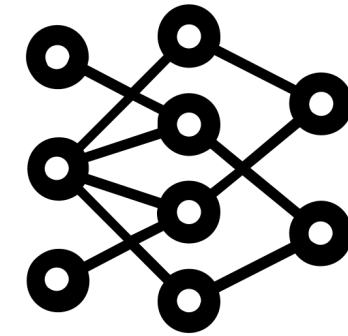
11 of 50 state names come from persons.

Vocabulary expanded with reflection tokens



Reflection tokens for retrieval and critique

US states got their names from a variety of sources.

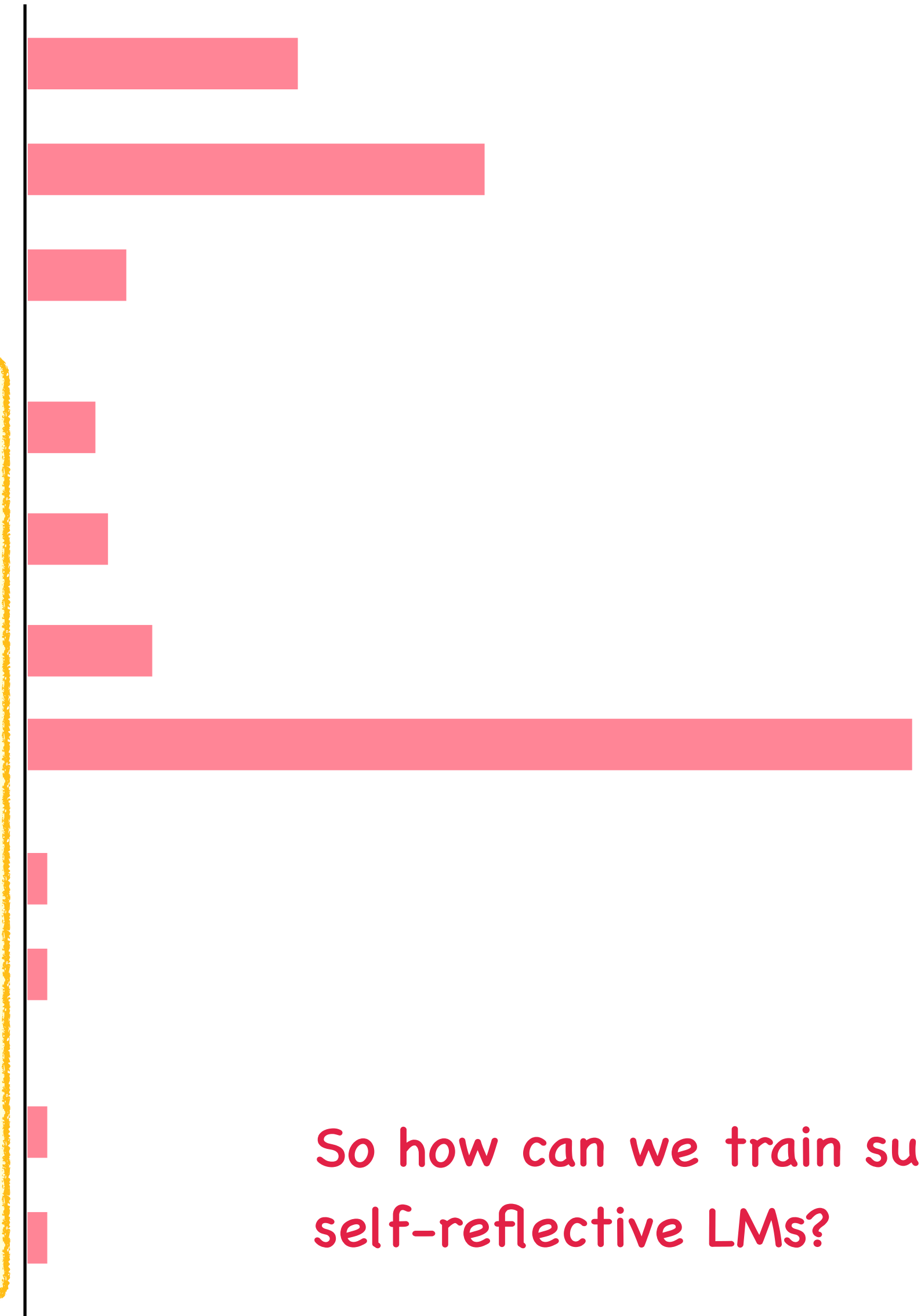
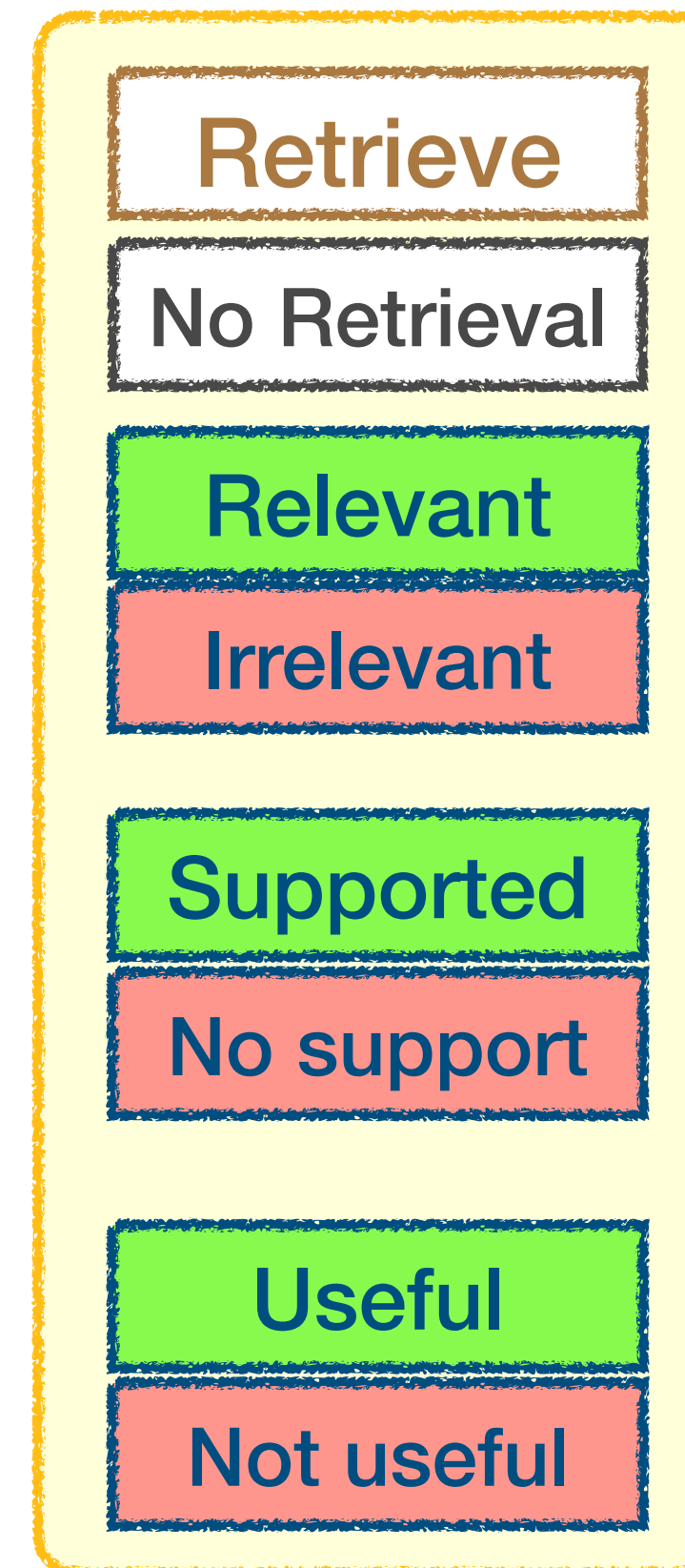


2 The history of human activity in Michigan began with settlement by Paleo-Indians.

California

e.g.,

A



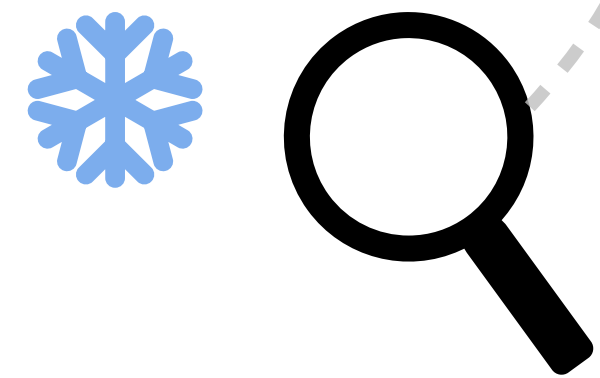
Vocabulary expanded with reflection tokens

So how can we train such self-reflective LMs?

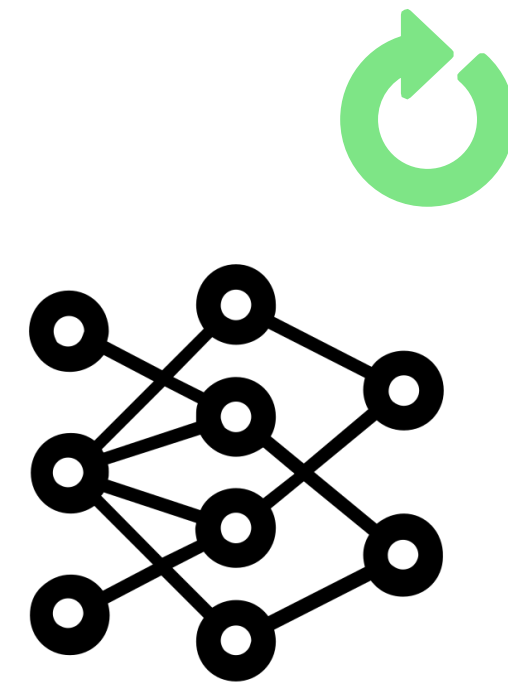
Self-RAG training

How did US states get their names?

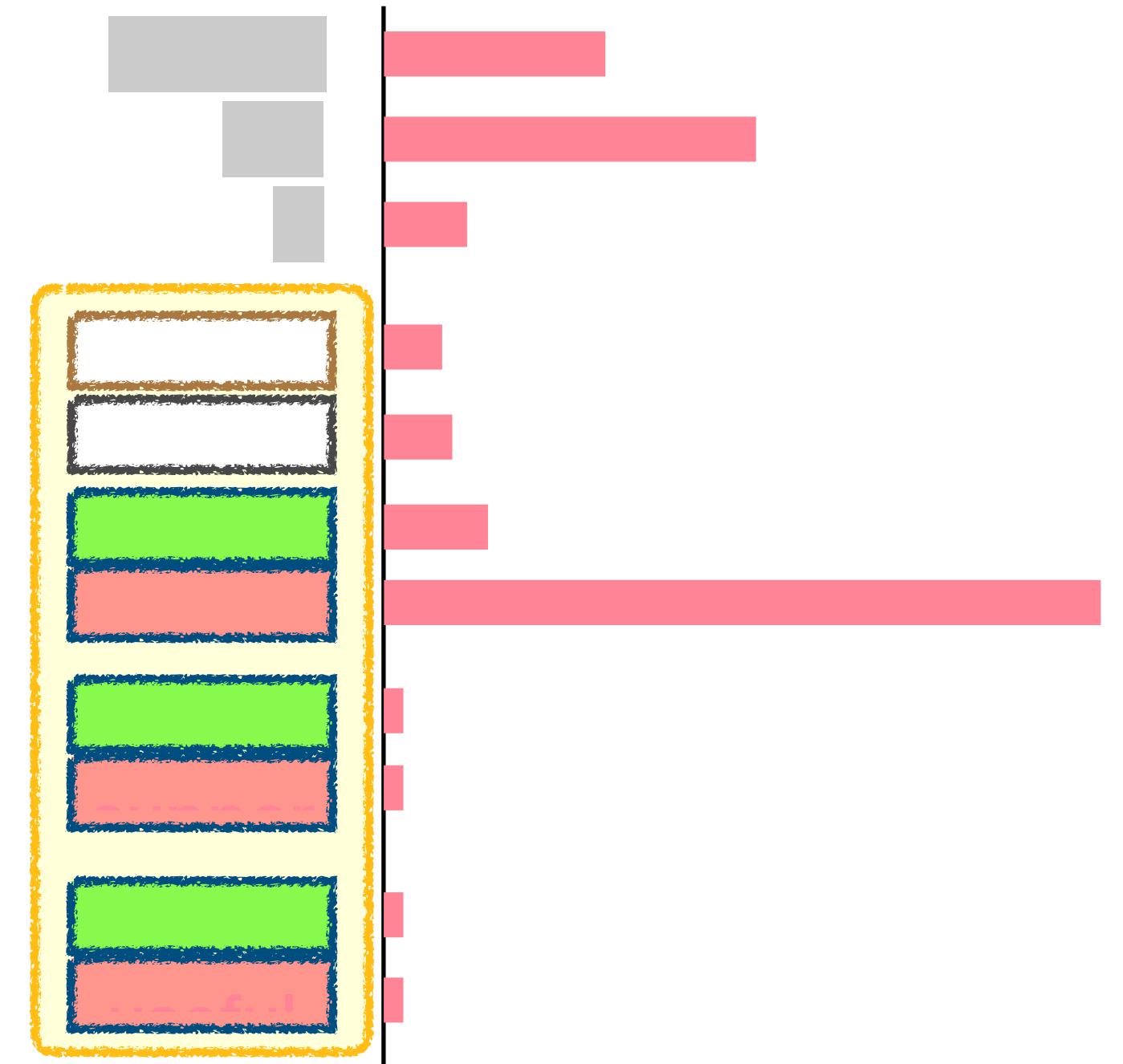
Of the fifty states, eleven are named after an individual person.



Retriever



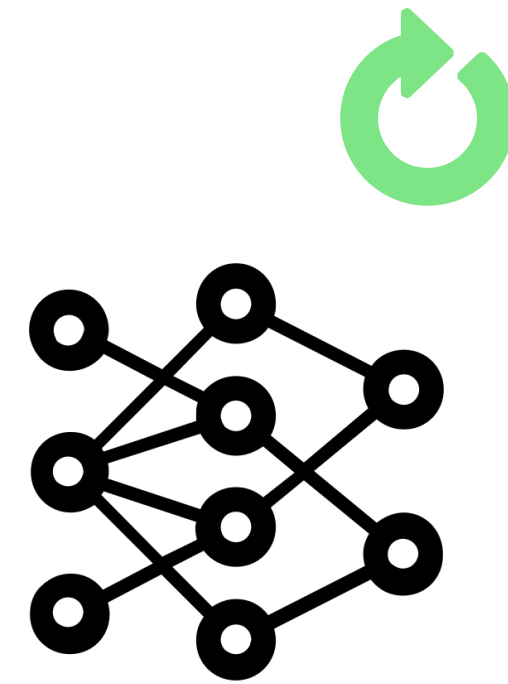
Generator LM



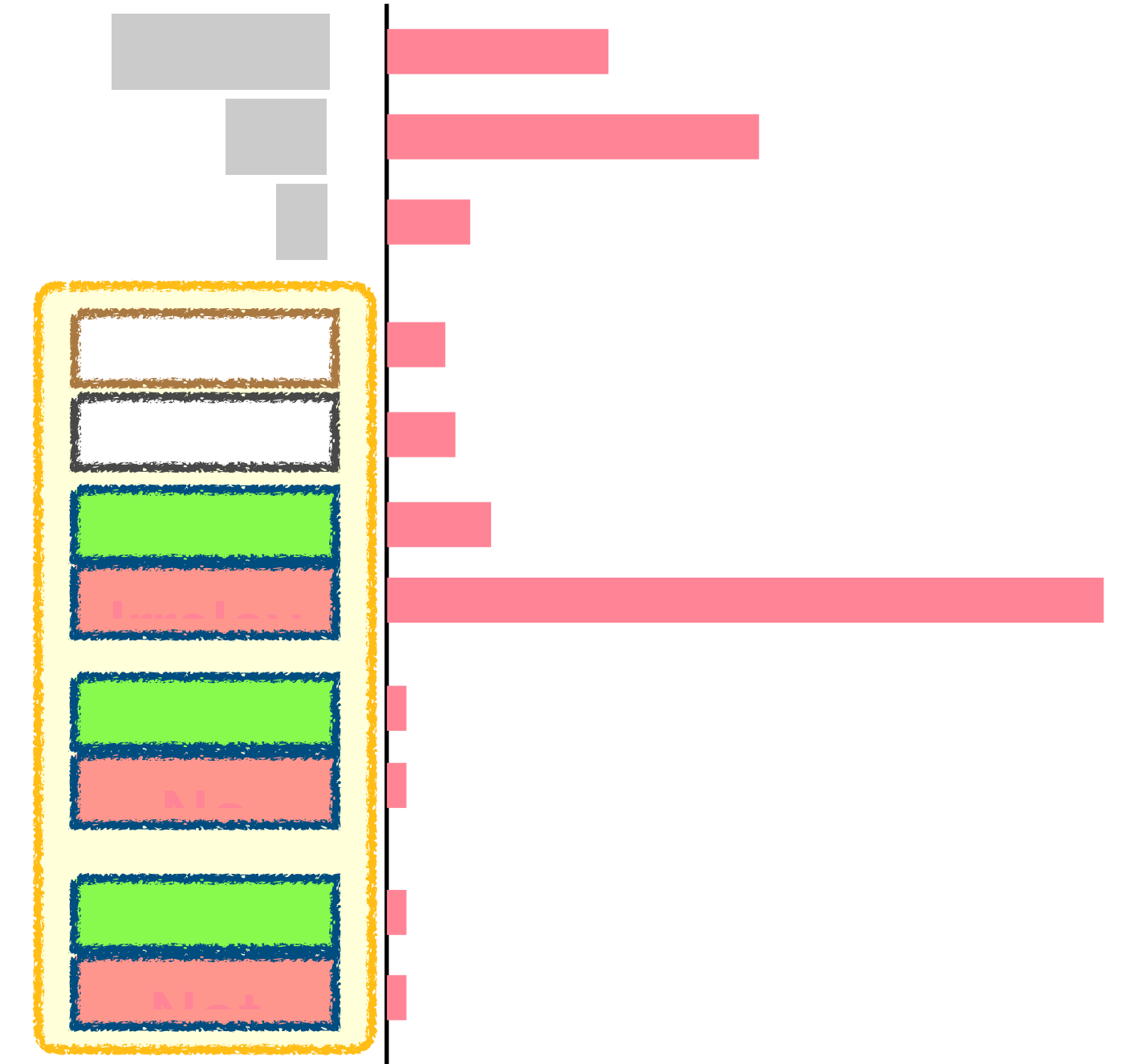
Self-RAG training

How did US states get their names?

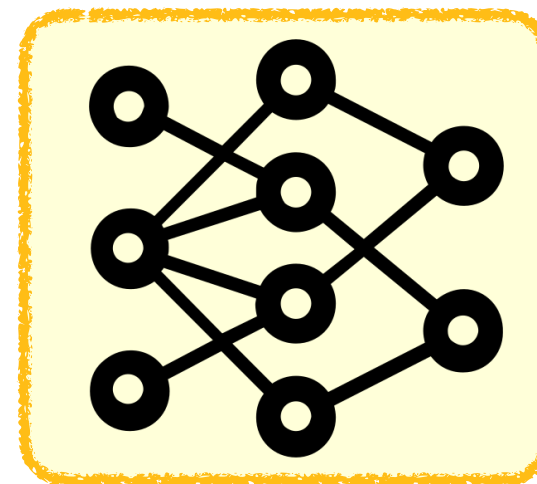
Of the fifty states, eleven are named after an individual person.



Generator LM



Critic LM teaches Generator LM to predict reflection tokens



Training time only Critic LM

Self-RAG training - critic LM

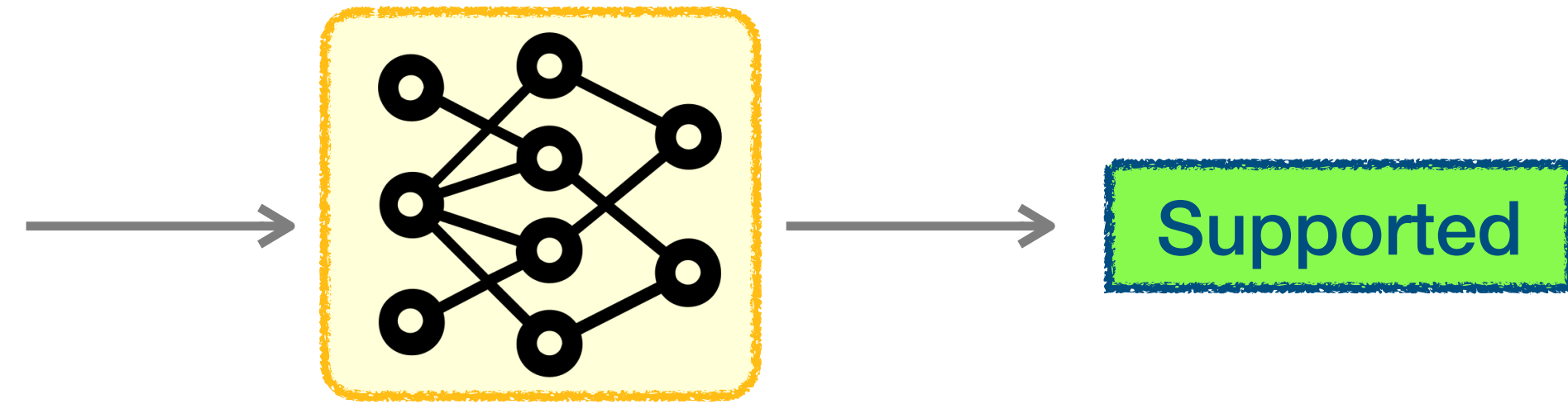
Evaluate if the output y to an input x is supported by retrieved doc d

x How did US states get their names?

d Of the fifty states, eleven are named after an individual person.

y 11 of 50 state names come from persons.

$$\max_{\mathcal{C}} \mathbb{E}_{((x,y,d),r) \sim \mathcal{D}_{critic}} \log p_{\mathcal{C}}(r | x, y, d)$$



Critic LM



85-90% acc. on validation set



30k fine-grained feedback, align with human in 90%

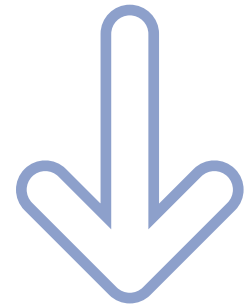


$$\mathcal{D}_{critic} = \{(x, y, d), r\}$$

Self-RAG training - augmented Instruction-tuning data



150k (input, output) instances from diverse instruction-following data



Input: Write an essay of your best summer vacation

Output: My best summer vacation was a magical escape to the coastal town of Santorini. The azure waters, charming white-washed building are unforgettable.

Augmented Output: No Retrieval My best summer vacation was a magical escape to the coastal town of Santorini. No Retrieval The azure waters, charming white-washed building are unforgettable experience. Util: 5



Input: How did US states get their names?

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



Augmented Output: Retrieve **1**

Of the fifty states, eleven are named after an individual person

. Relevant 11 of 50 states' names come from person. Supported Retrieve **2**

LOUISIANA: Named in honor of Louis XIV of France.

. Relevant For instance, Louisiana was named after King Louis XIV, and Georgia was named after King George II. Partially Util: 5



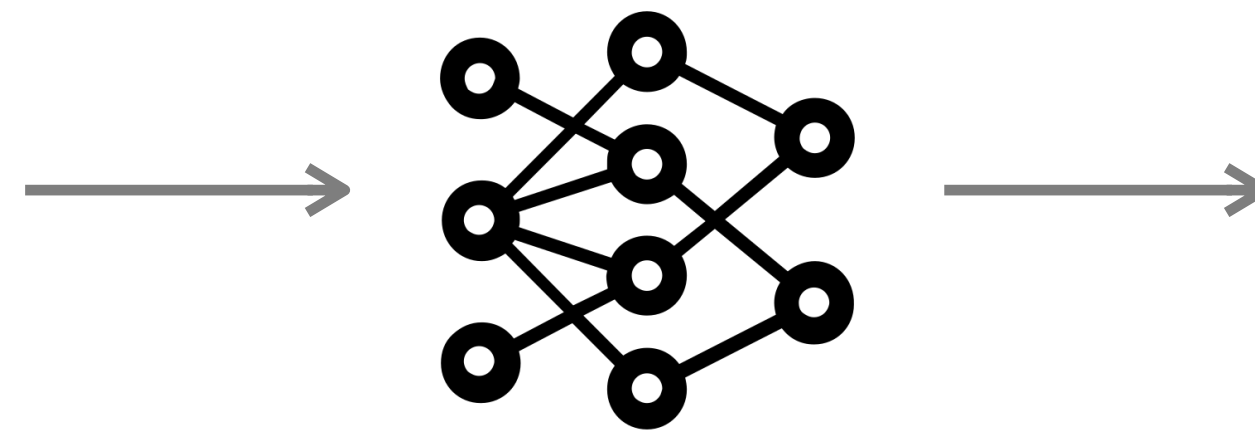
Retriever

Self-RAG training - generator LM

$$\max_{\mathcal{M}} \mathbb{E}_{(x,d,y,r) \sim \mathcal{D}_{gen}} \log p_{\mathcal{M}}(y, r | x, d).$$

x How did US states get their names?

d Of the fifty states, eleven are named after an individual person.



y 11 of 50 state names come from persons.

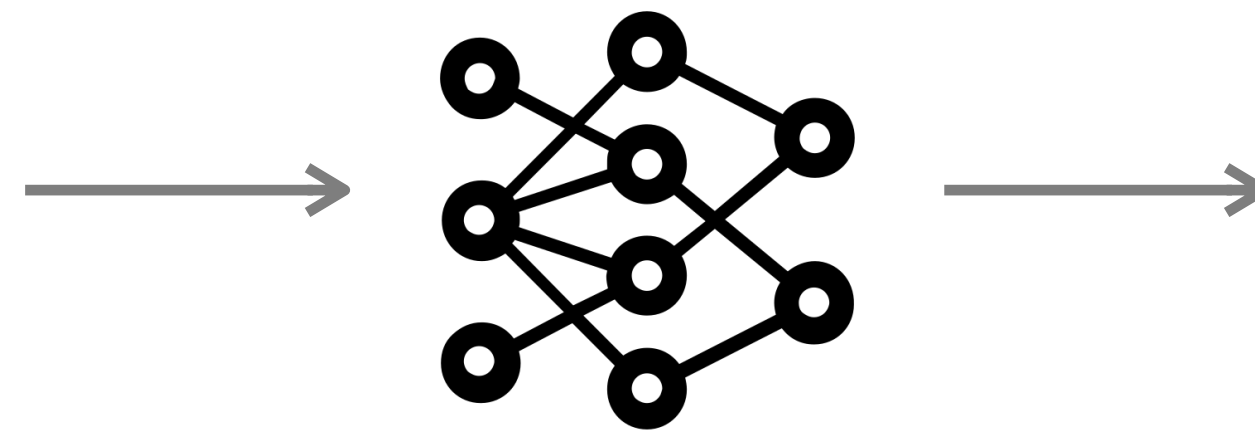
r **Supported**

Train with a standard next token objective with expanded vocabulary

Self-RAG training - generator LM

$$\max_{\mathcal{M}} \mathbb{E}_{(x,d,y,r) \sim \mathcal{D}_{gen}} \log p_{\mathcal{M}}(y, r | x, d).$$

x How did US states get their names?



y 11 of 50 state names come from persons.

d Of the fifty states, eleven are named after an individual person.

r **Supported**

Generator LM

Memory-efficient & stable training

Easily applied to new pre-trained LM

Customize & control via reflection tokens **How?**

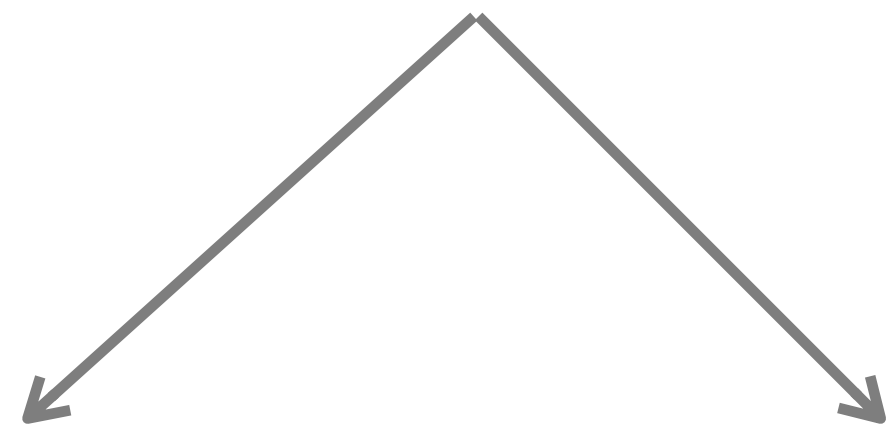
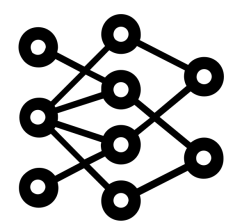
Self-reflection-guided decoding

Conduct segment-level beam search to find top k segments

Prompt 1

$f(\text{Relevant Supported})$

0.9



Prompt 2

$f(\text{Irrelevant})$

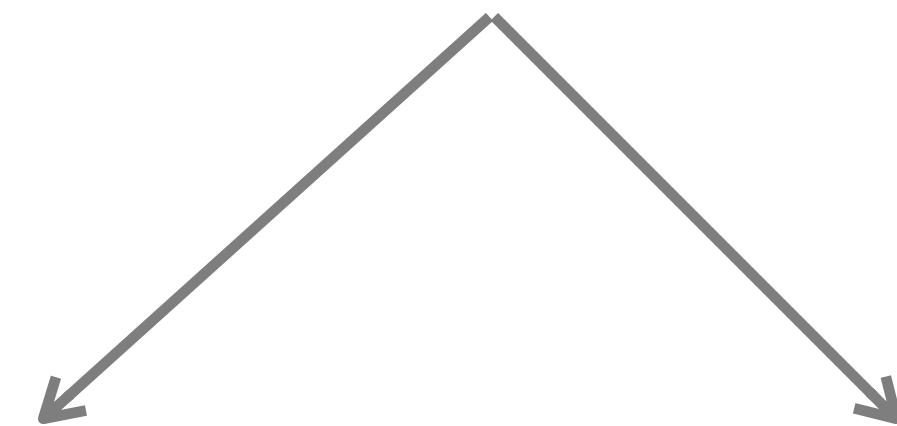
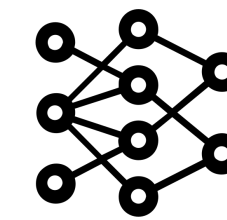
0.1

$$f(y_t, d, \text{Critique}) = p(y_t | x, d, y_{<t}) + \mathcal{S}(\text{Critique}), \text{ where}$$
$$\mathcal{S}(\text{Critique}) = \sum_{G \in \mathcal{G}} w^G s_t^G \text{ for } \mathcal{G} = \{\text{ISREL}, \text{ISSUP}, \text{ISUSE}\},$$

Prompt 3

$f(\text{Relevant Partially})$

0.4



Self-reflection-guided decoding

Enable simple model customization by changing weights

Prompt 1

$f(\text{Relevant Supported})$

0.9

Prompt 2

$f(\text{Irrelevant})$

0.1

Prompt 3

$f(\text{Relevant Partially})$

0.4

$f(y_t, d, \text{Critique}) = p(y_t|x, d, y_{<t}) + \mathcal{S}(\text{Critique})$, where

$\mathcal{S}(\text{Critique}) = \sum_{G \in \mathcal{G}} w^G s_t^G$ for $\mathcal{G} = \{\text{ISREL}, \text{ISSUP}, \text{ISUSE}\}$,

Experimental details

Tasks and datasets

- **Closed-set tasks** (classifications, multiple-choice QA)
 - ARC-Challenge (Clark et al., 2018)
 - PubHealth (Zhang et al., 2023)
- **Short-form generation**
 - OpenQA - PopQA
 - Trivia QA (Joshi et al., 2017)
- **Long-form generation**
 - ASQA-ALCE [fluency, citation accuracy, correctness] (Gao et al., 2023)
 - Bio generations [FactScore] (Min et al., 2023)

Experimental details

More details of training & test are in our paper!

Training details

- **Critic training data:** 4k-20k instances for each type
- **Generator training data:** 150k instruction-following datasets
 - ShareGPT
 - OpenAssistant
 - Alpaca
 - FLANV2
 - Natural Questions
- **Base LMs:** Llama2-7B, 13B (Touvron et al., 2023)
- **Computation:** 4*A100 (15 hours)

Experimental details

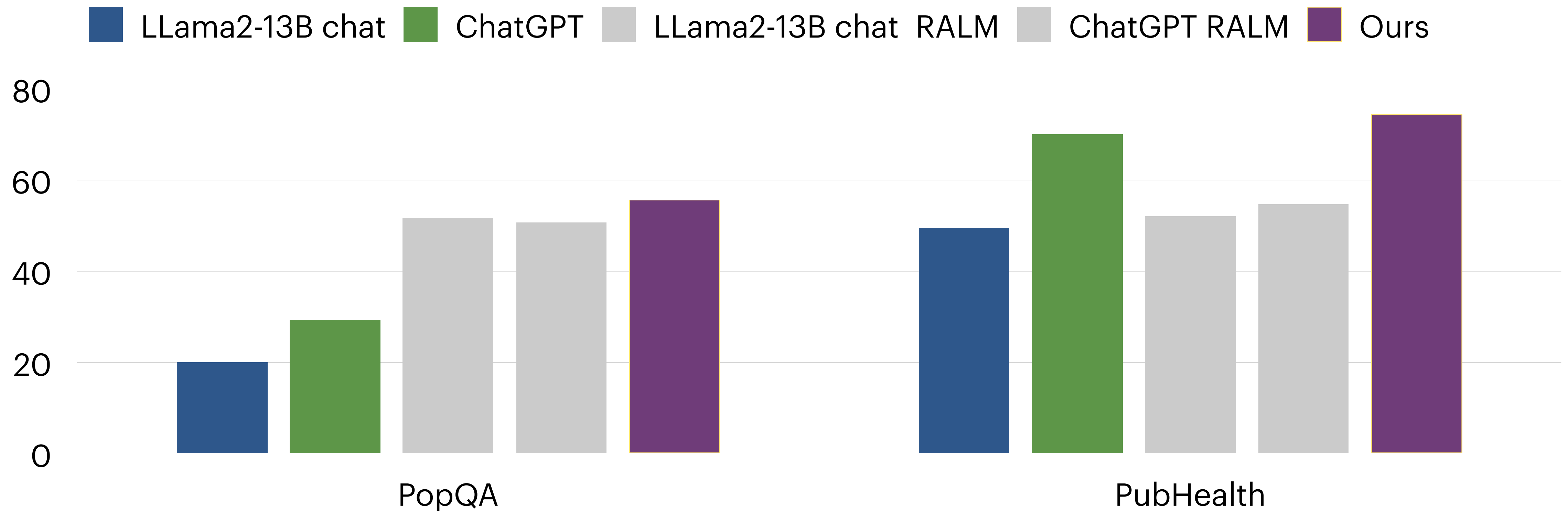
More details of training & test are in our paper!

Inference-time details

- **Retriever Encoder:** Contriever-MS MARCO (Izacard et al., 2022)
- **Index:** HNSW Index (0.1 sec / query) and FLAT Index (5 sec / query)
- **Efficient LM inference:** vllm (Kwon et al., 2023)
- **Tree decoding configuration:** max 200 tokens per depth, max depth of 7

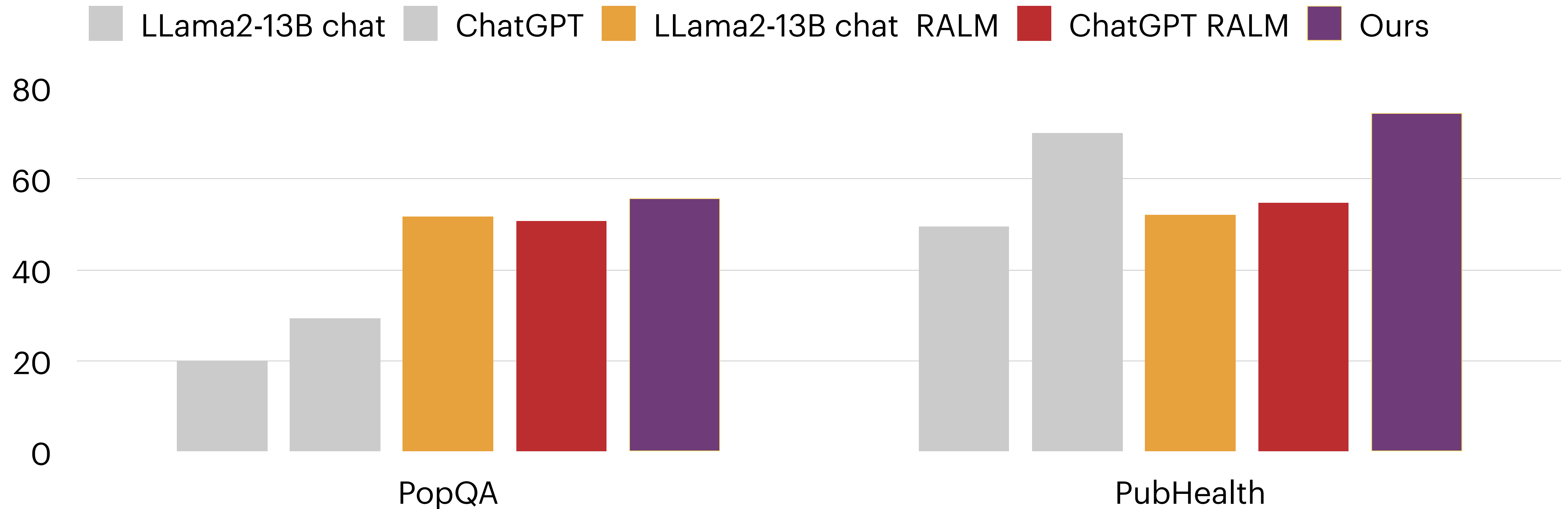
Experimental results (short-form & closed)

Self-RAG outperforms vanilla LMs incl. ChatGPT



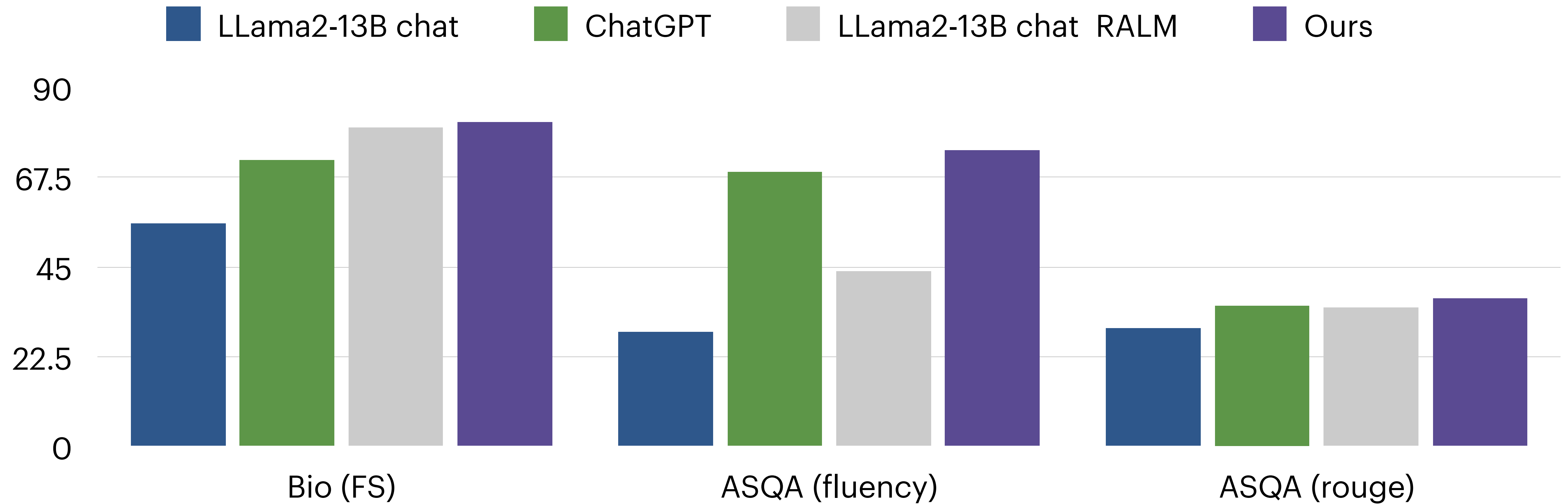
Experimental results (short-form & closed)

Self-RAG outperforms standard RAG + LLMs



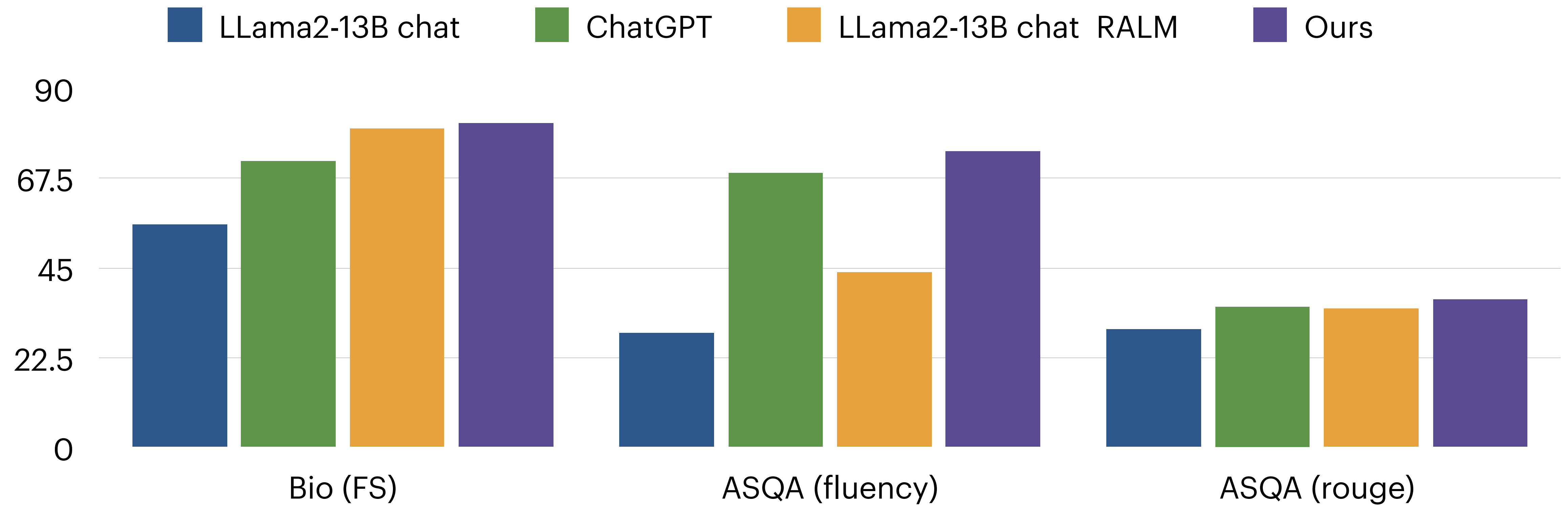
Experimental results (long-form)

Outperforms other LMs in terms of factuality & fluency correctness



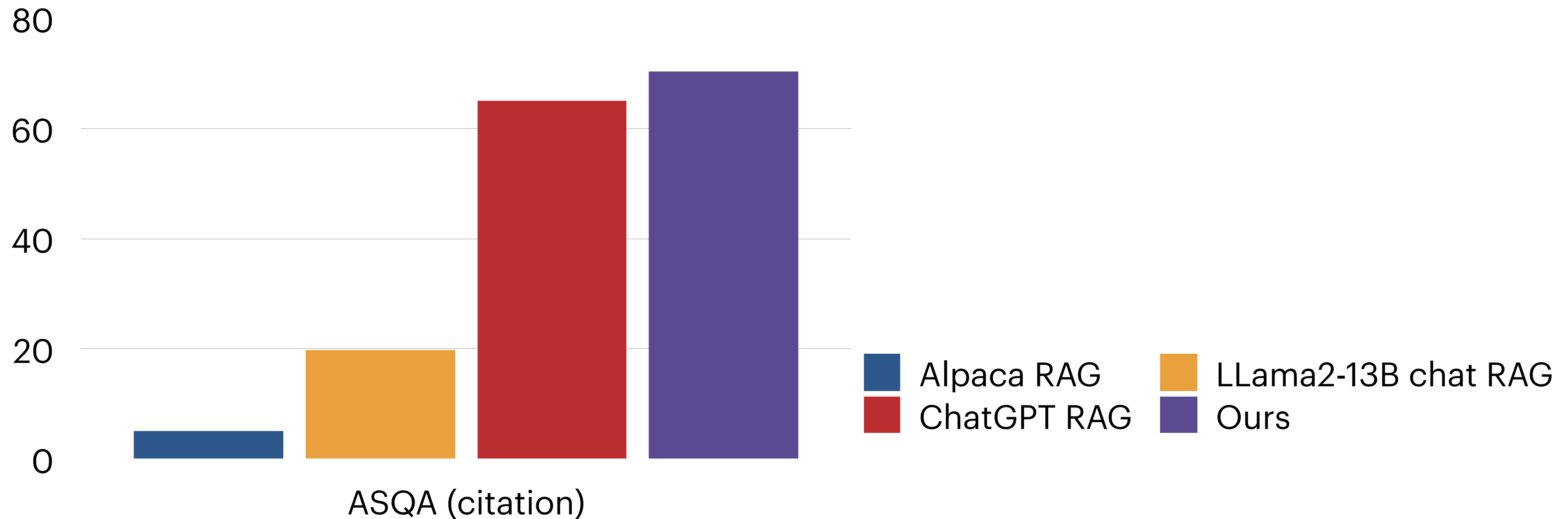
Experimental results (long-form)

Outperforms other LMs in terms of factuality & fluency correctness



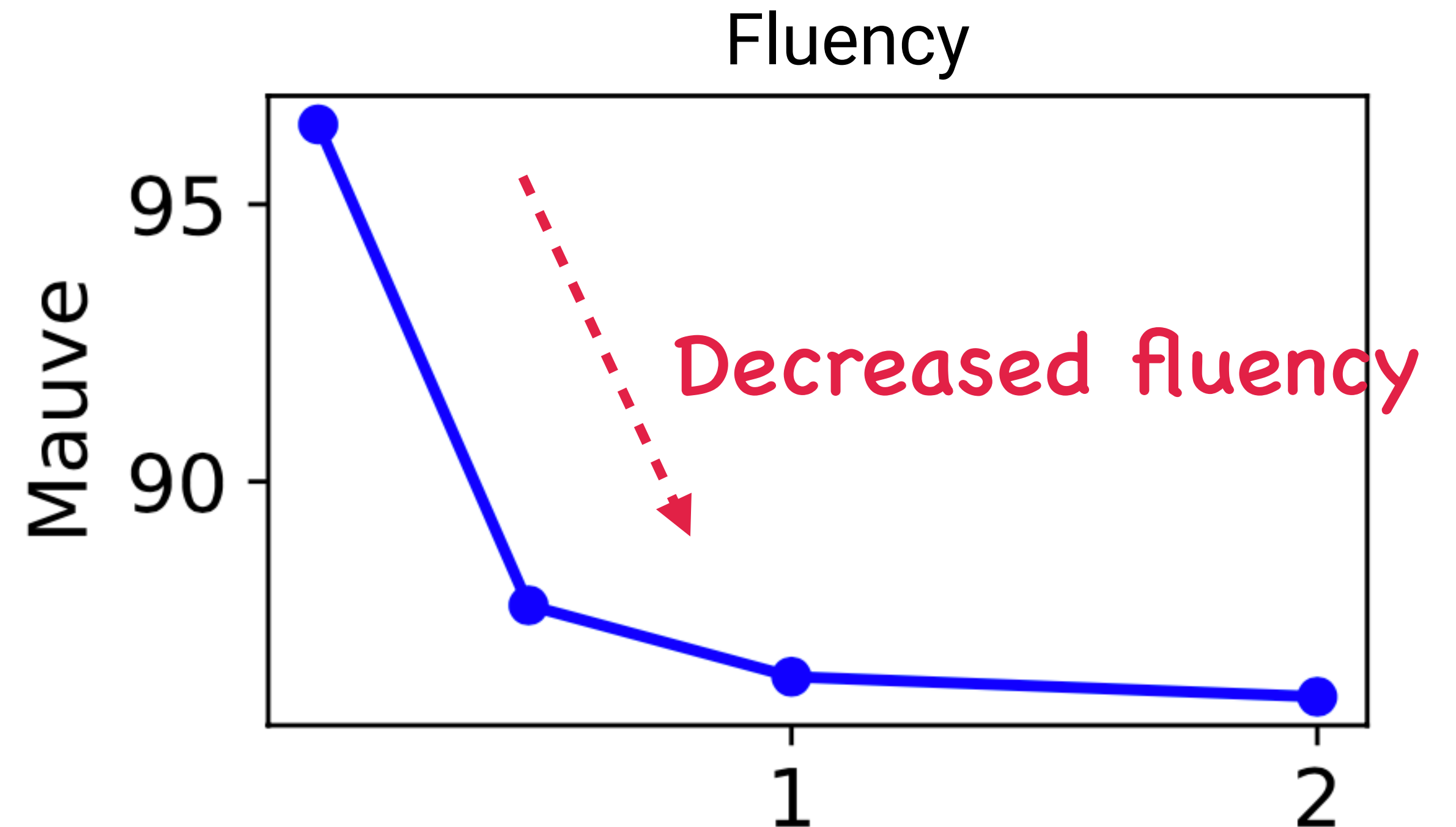
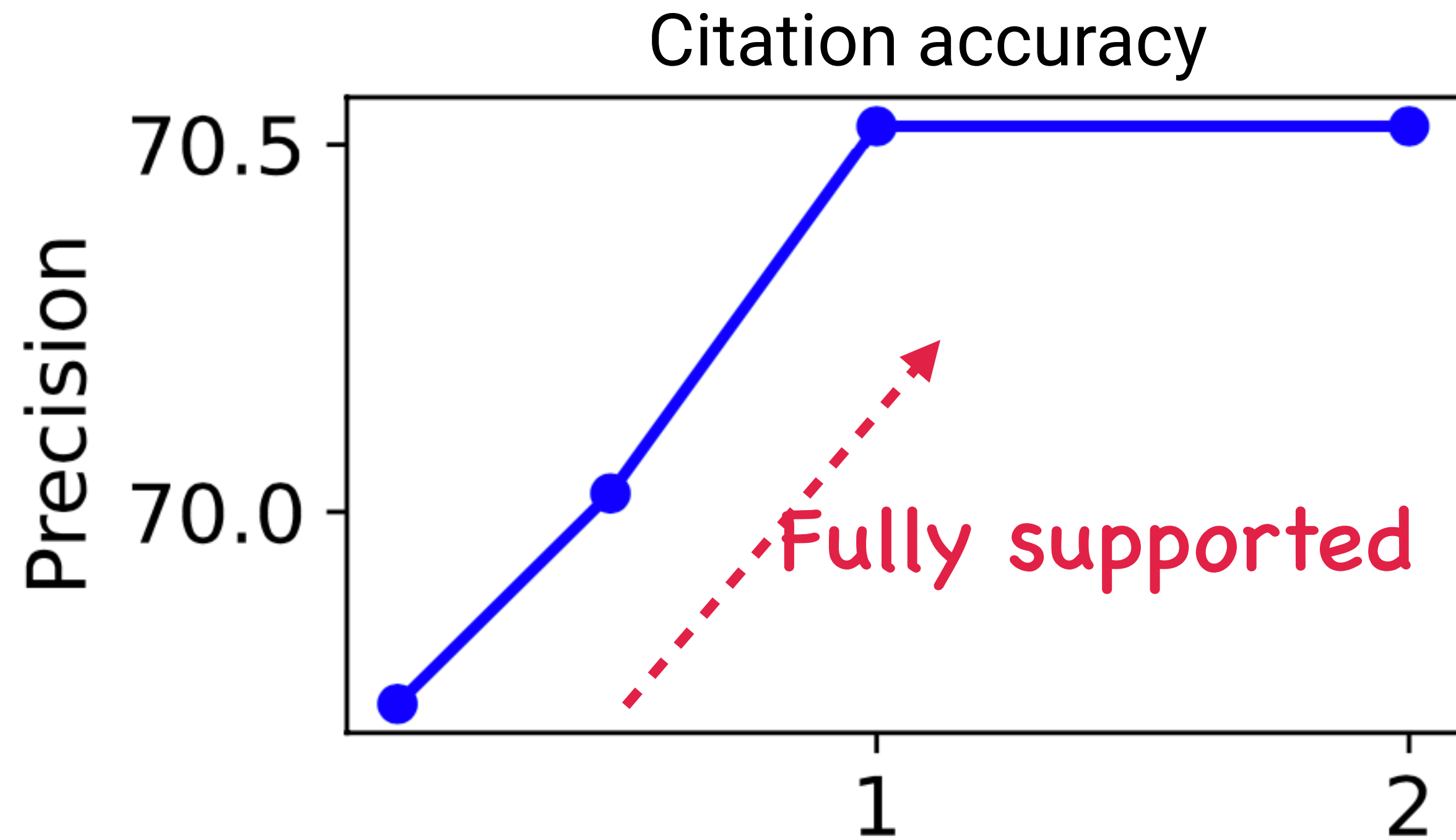
Experimental results (long-form citation precisions)

Significantly improves llama2-13B citation accuracy, matching ChatGPT



Inference-time customization via **self-reflection**

Decoding-time control via reflection tokens change model behaviors

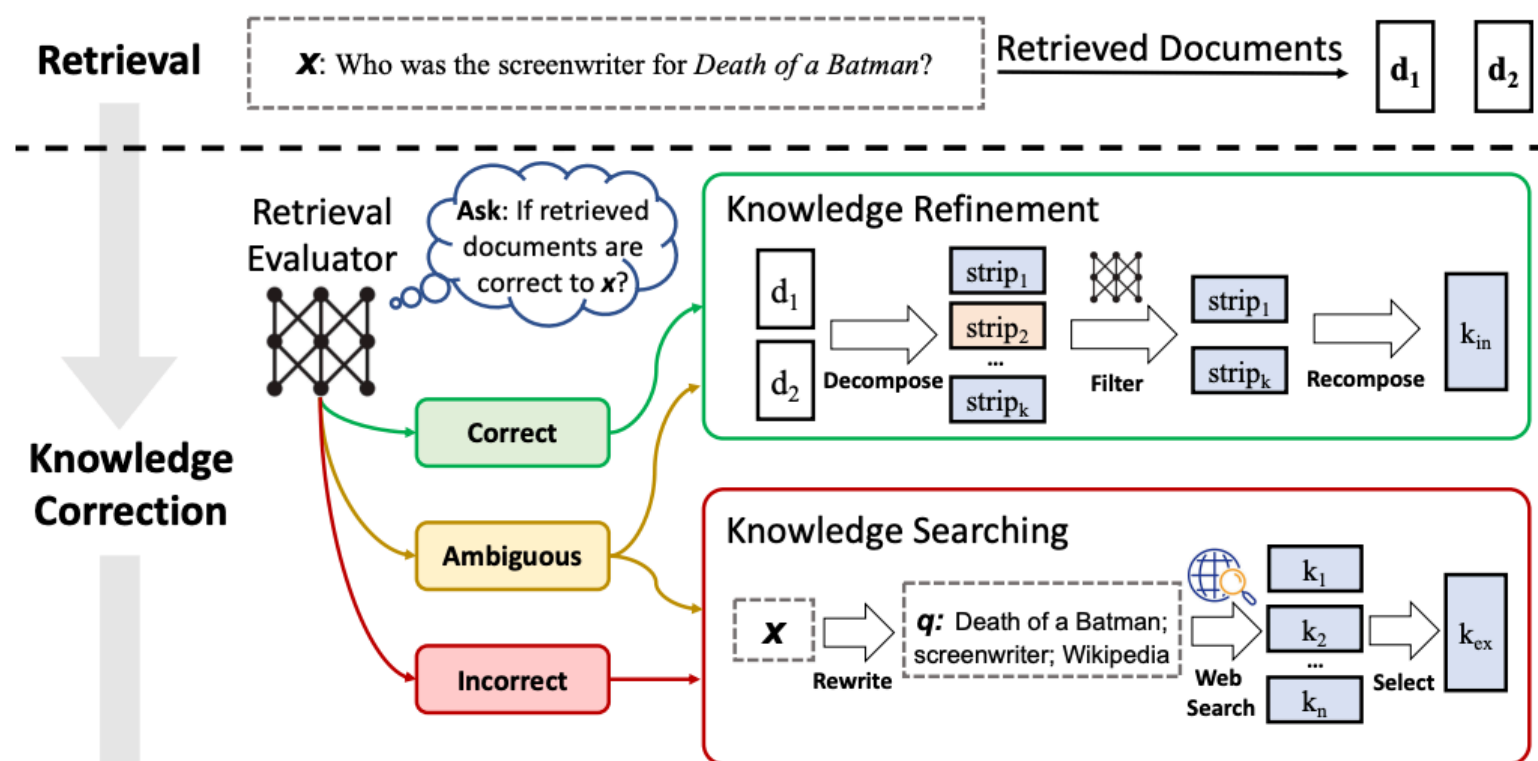


x axis – Weight for **Supported** (larger \rightarrow more emphasis on supported)

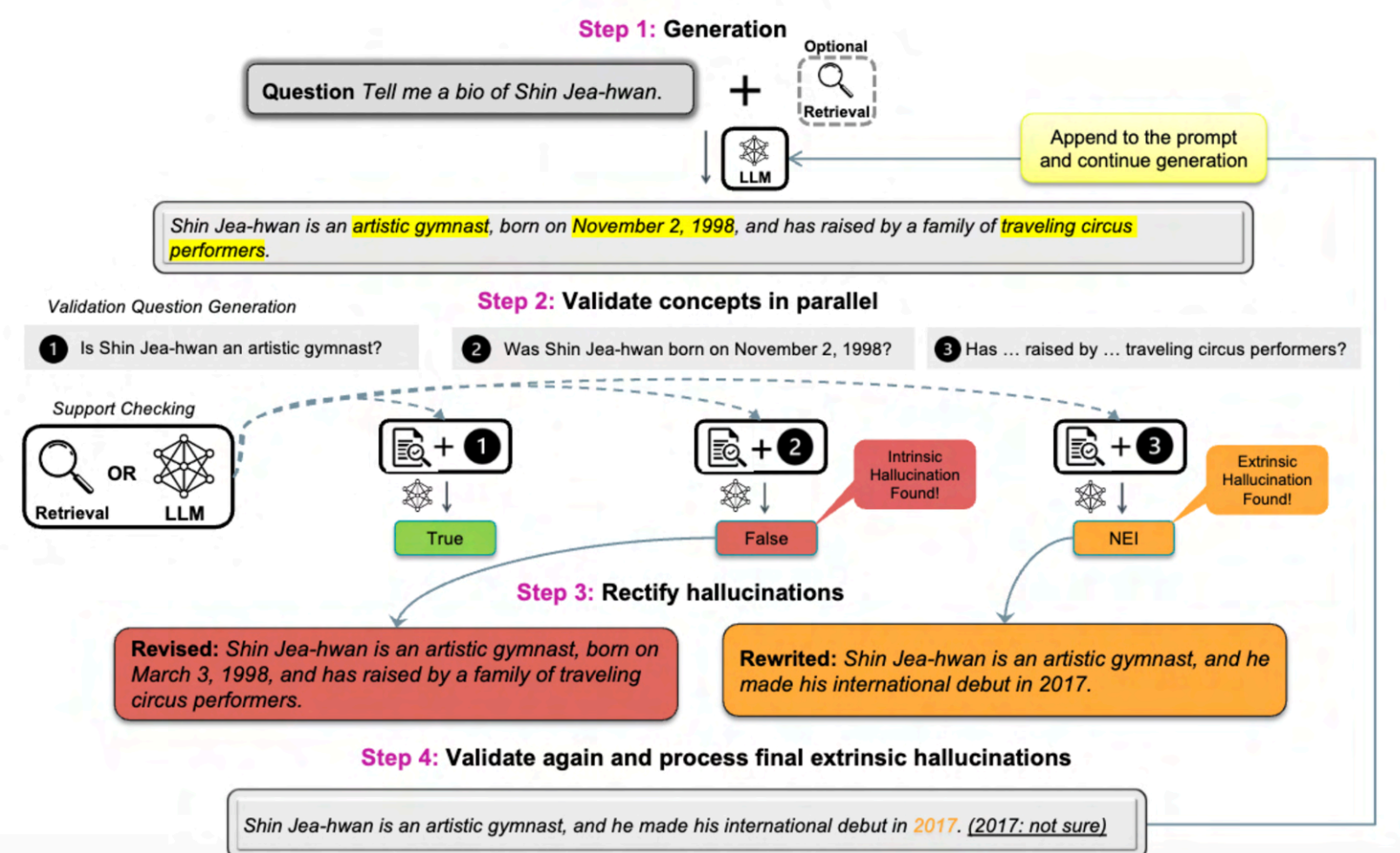
Impacts on academic communities and applications

New advanced RAG methods with reflections inspired by Self-RAG

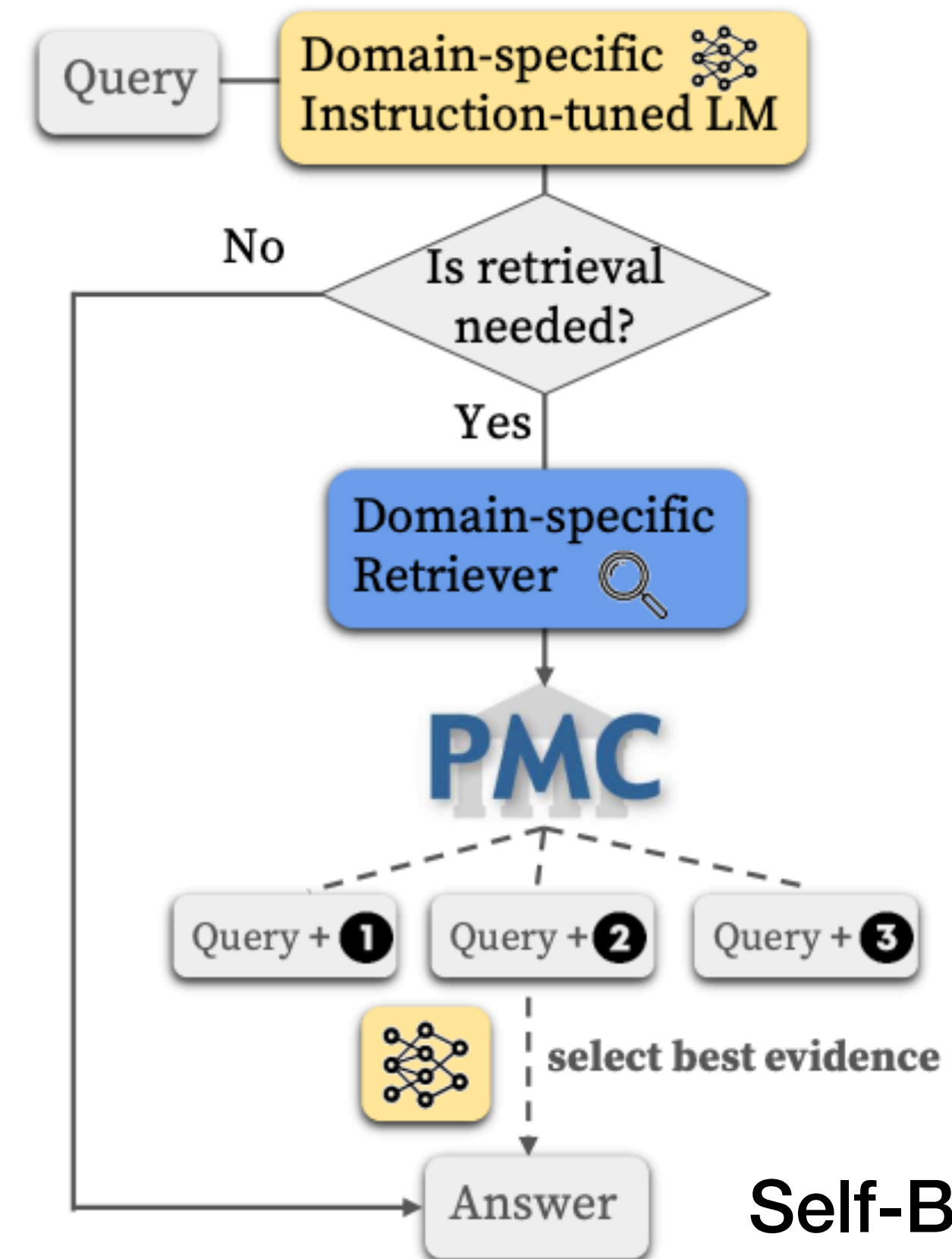
Applications to expert domains (e.g., Biomedical)



CRAG
(Yang et al., 2024)



EVER
(Kang et al., 2024)



Self-BioRAG
(Jeong et al., 2024)

Impacts on academic communities and applications

The screenshot shows the top of a LangChain website article. The header includes the LangChain logo and navigation links: 'By LangChain', 'Release Notes', 'Case Studies', 'LangChain', 'GitHub', 'Docs', 'Sign in', and a 'Subscribe' button. The article title is 'Self-Reflective RAG with LangGraph', with a sub-header '(REFLECTION AND SELF-CORRECTION)'. A diagram illustrates a workflow: 'Question' leads to a 'Retrieve (Node)' (blue circle), then to a 'Grade (Node)' (green circle), then to a decision diamond 'Any doc relevant?'. If 'Yes', it goes to a 'Generate (Node)' (purple circle) and then to 'Answer'. If 'No', it goes to a 'Re-write query (Node)' (green circle), which then loops back to the 'Retrieve (Node)'. The article is dated '6 MIN READ FEB 7, 2024'.

Downloads last month
6,146



Self-RAG has been integrated into LangChain, LlamaIndex ... etc

LlamaIndex • Feb 13, 2024

LlamaIndex Newsletter 2023-02-13

Newsletter LlamaIndex AI Rag LLM

👑 The highlights:

1. **Self-RAG:** Introducing Self-RAG, now part of LlamaIndex as a LlamaPack. Boosts LLM training and RAG workflows with dynamic capabilities. [Notebook](#), [Tweet](#).

Self-RAG – Self-Reflective Retrieval-Augmented Generation

- ✓ An LM learns to retrieve, generate and critique
- ✓ Instruction-tuned LMs trained with fine-grained reflection tokens
- ✓ Outperforms other LMs in six tasks, improving citation accuracy



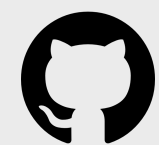
<https://selfrag.github.io/>



<https://arxiv.org/abs/2212.10511>



https://huggingface.co/selfrag/selfrag_llama2_7b (13b)



<https://github.com/AkariAsai/self-rag> (1.4k ★!)

Today's lecture

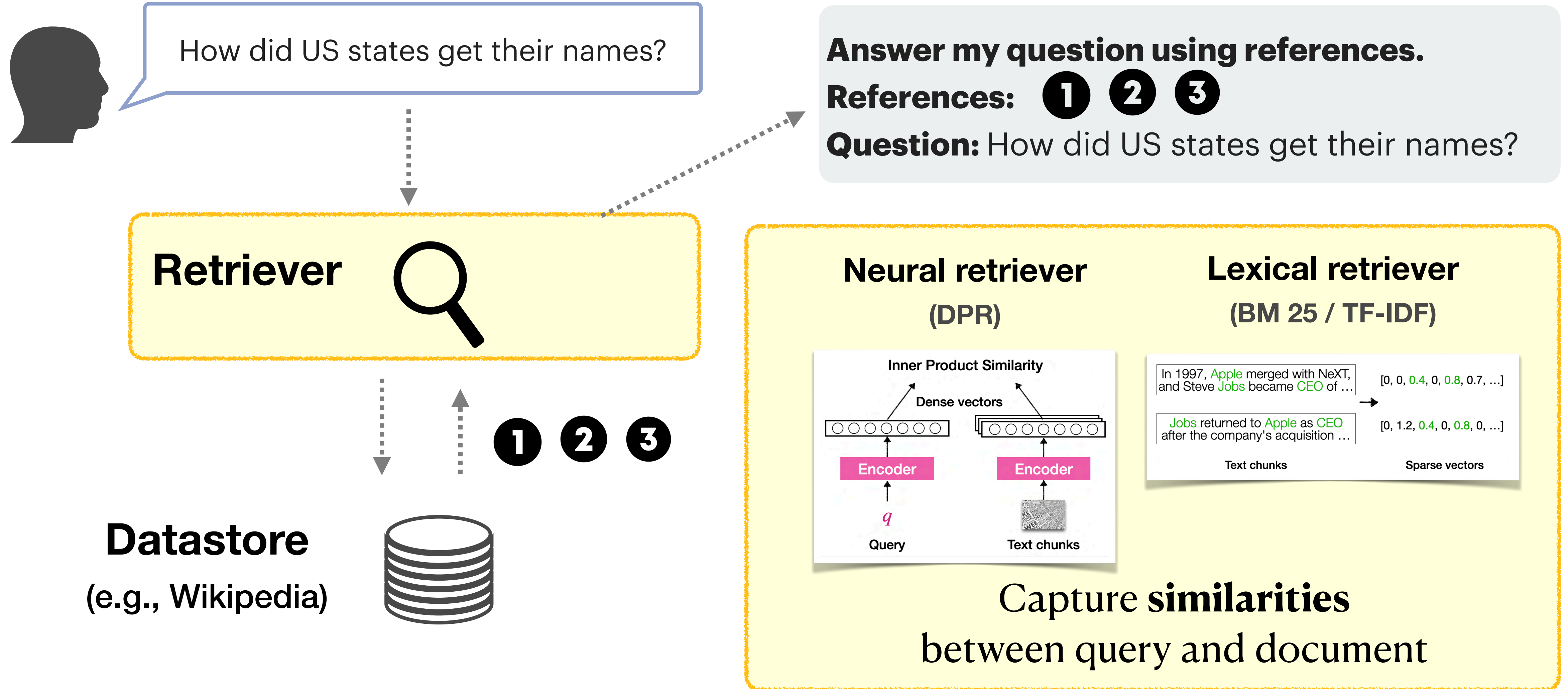
Promises and Limitations of Retrieval-augmented LMs

Reliable inference: Self-reflective RAG with dynamic retrieval

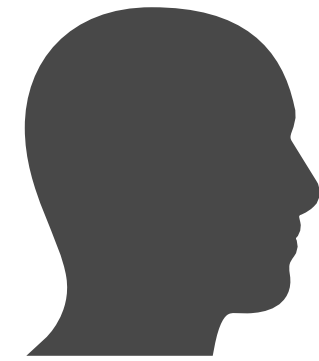
Versatile retriever: Intent-aware retrievers with LMs

Summary and future directions: RAG in the wild

Standard RAG originally designed for a single task (e.g., QA)



Modern RAG systems have been used in diverse scenarios



Implement batch normalization in python



We're planning to travel to Japan. Create a feasible 7 days travel plan to visit Tokyo, Kyoto, Sapporo, Osaka, and Okinawa.

Modern RAG systems have been used in diverse scenarios

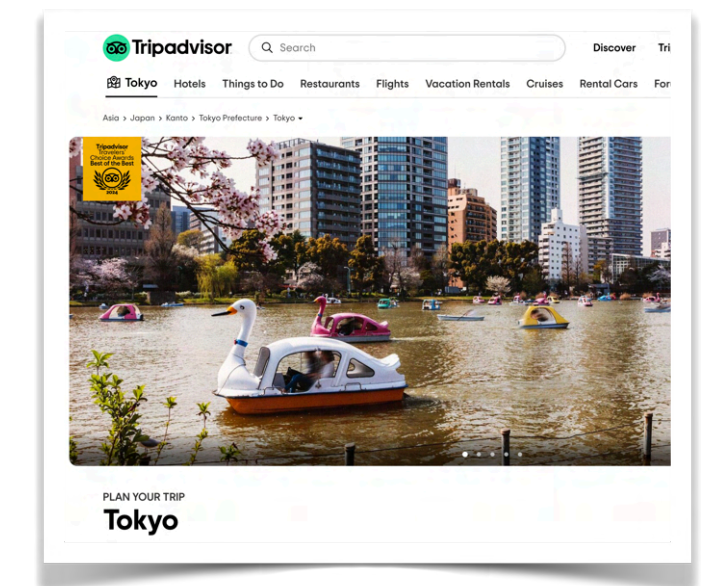
Implement batch normalization in python

We're planning to travel to Japan. Create a feasible 7 days travel plan to visit Tokyo, Kyoto, Sapporo, Osaka, and Okinawa.

Retriever



```
Code Blame 301 lines (246 loc) · 9.46 KB
78
79
80 def batchnorm_forward(x, gamma, beta, bn_param):
81     """
82     Forward pass for batch normalization.
83
84     During training the sample mean and (uncorrected) sample variance are
85     computed from minibatch statistics and used to normalize the incoming data.
86     During training we also keep an exponentially decaying running mean of the mean
87     and variance of each feature, and these averages are used to normalize data
88     at test-time.
```



Datastore

(e.g., Github codes)



Datastore

(e.g., Tripadvisor, Wikipedia pages)



Modern RAG systems have been used in diverse scenarios

Implement batch normalization in python

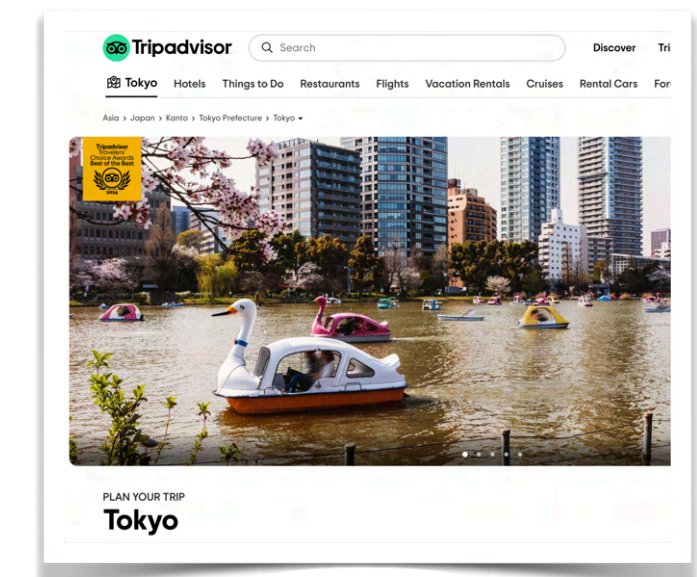
We're planning to travel to Japan. Create a feasible 7 days travel plan to visit Tokyo, Kyoto, Sapporo, Osaka, and Okinawa.

Retrieve

Different information needs often require different information retrieval systems

```
Code Blame 301 lines (246 loc) · 9.46 KB
78
79
80 def batchnorm_forward(x, gamma, beta, bn_param):
81     """
82     Forward pass for batch normalization.
83
84     During training the sample mean and (uncorrected) sample variance are
85     computed from minibatch statistics and used to normalize the incoming data.
86     During training we also keep an exponentially decaying running mean of the mean
87     and variance of each feature, and these averages are used to normalize data
88     at test-time.
```

Kyoto (/kjoʊtoʊ/ Japanese: 京都, *Kyōto* [koˈto]), officially **Kyoto City** (京都市, *Kyōto-shi*, [koˈtoˈɕi]), is the capital city of **Kyoto Prefecture** in the Kansai region of Japan's largest and most populous island of Honshu. As of 2020, the city had a population of 1.46 million, making it the **ninth-most populous** city in Japan. The city is the cultural anchor of the substantially larger **Greater Kyoto**, a metropolitan statistical area (MSA) home to a census-estimated 3.8 million people. It is also part of the even larger **Keihanshin metropolitan area**, along with **Osaka** and **Kobe**.



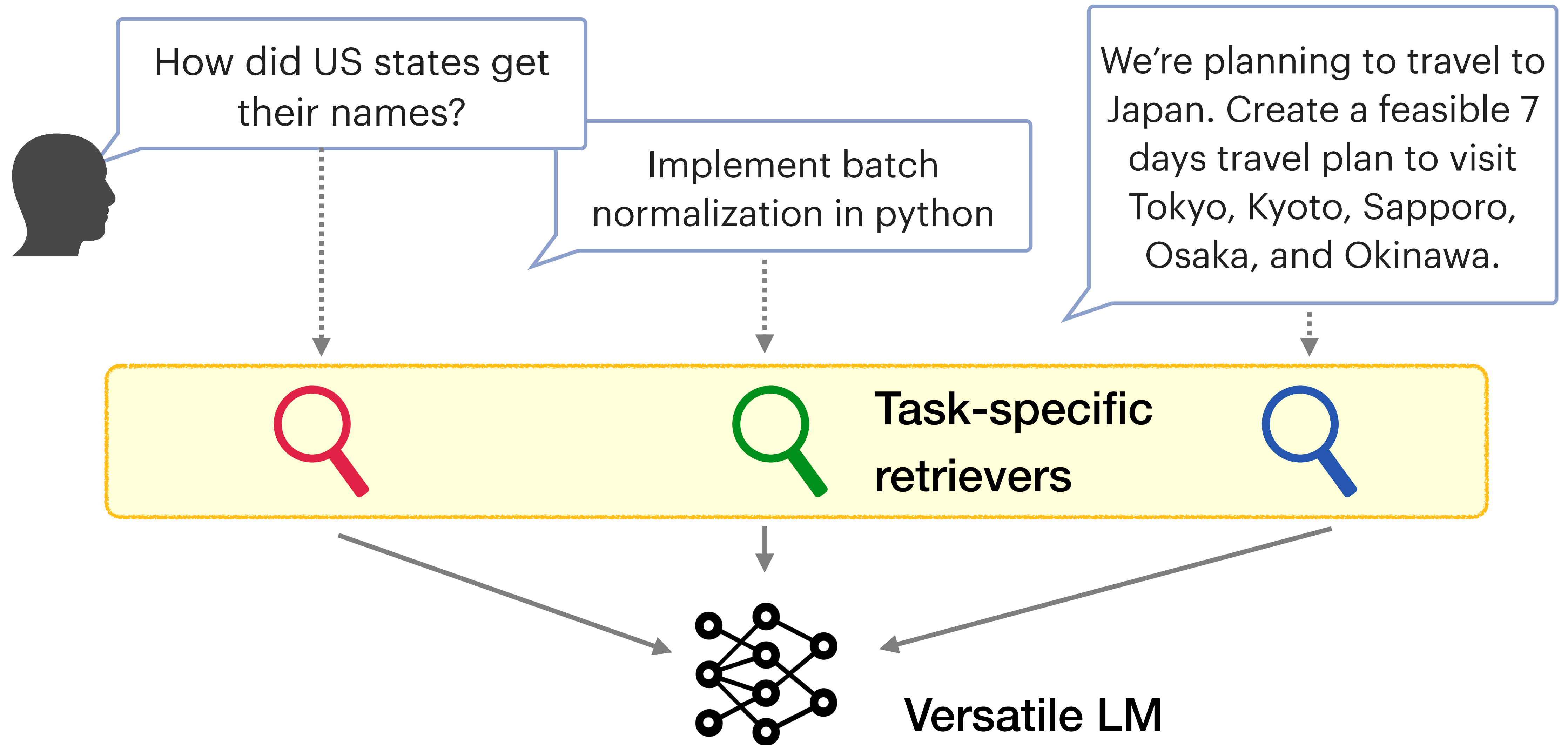
Datastore
(e.g., Github codes)



Datastore
(e.g., Tripadvisor, Wikipedia pages)



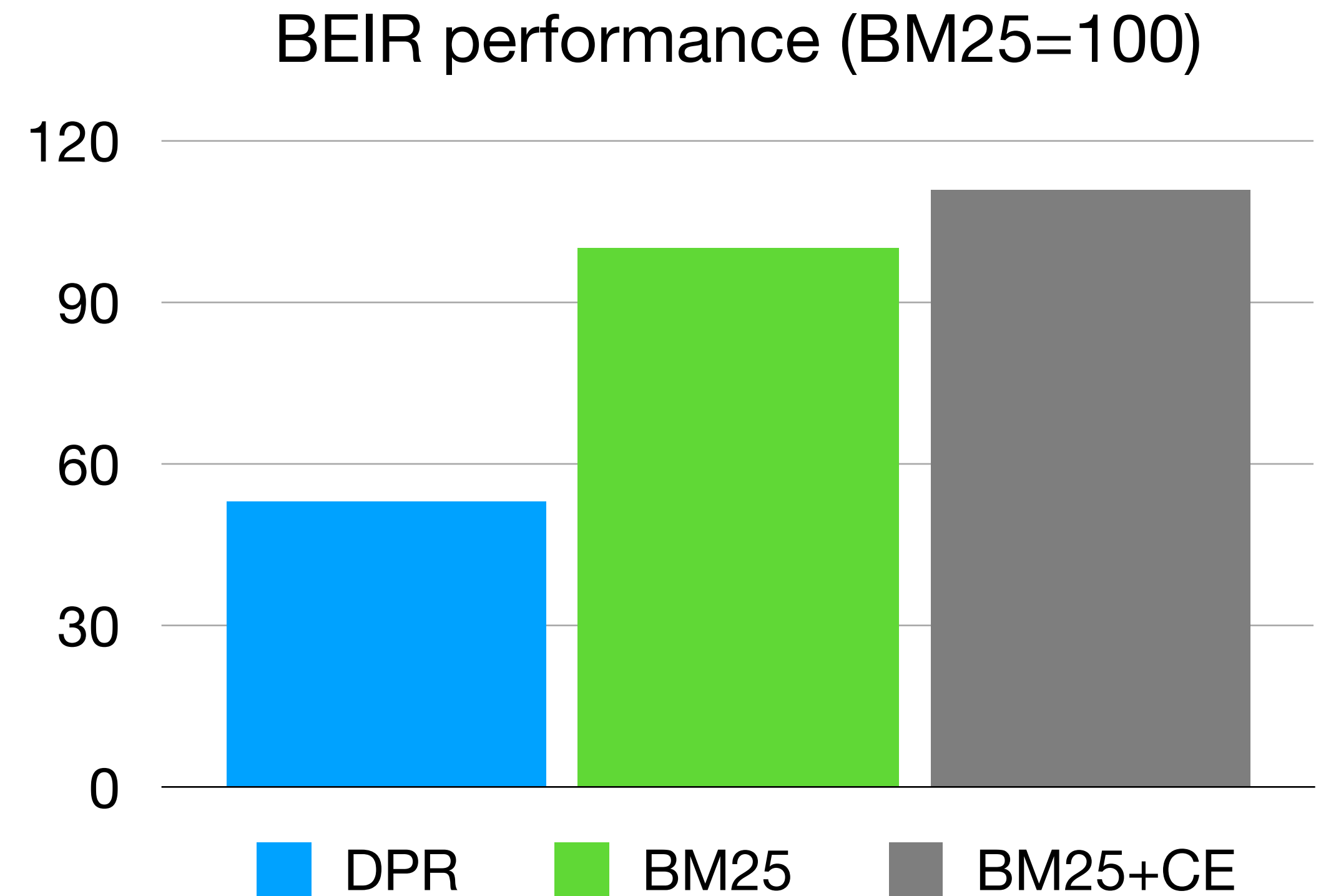
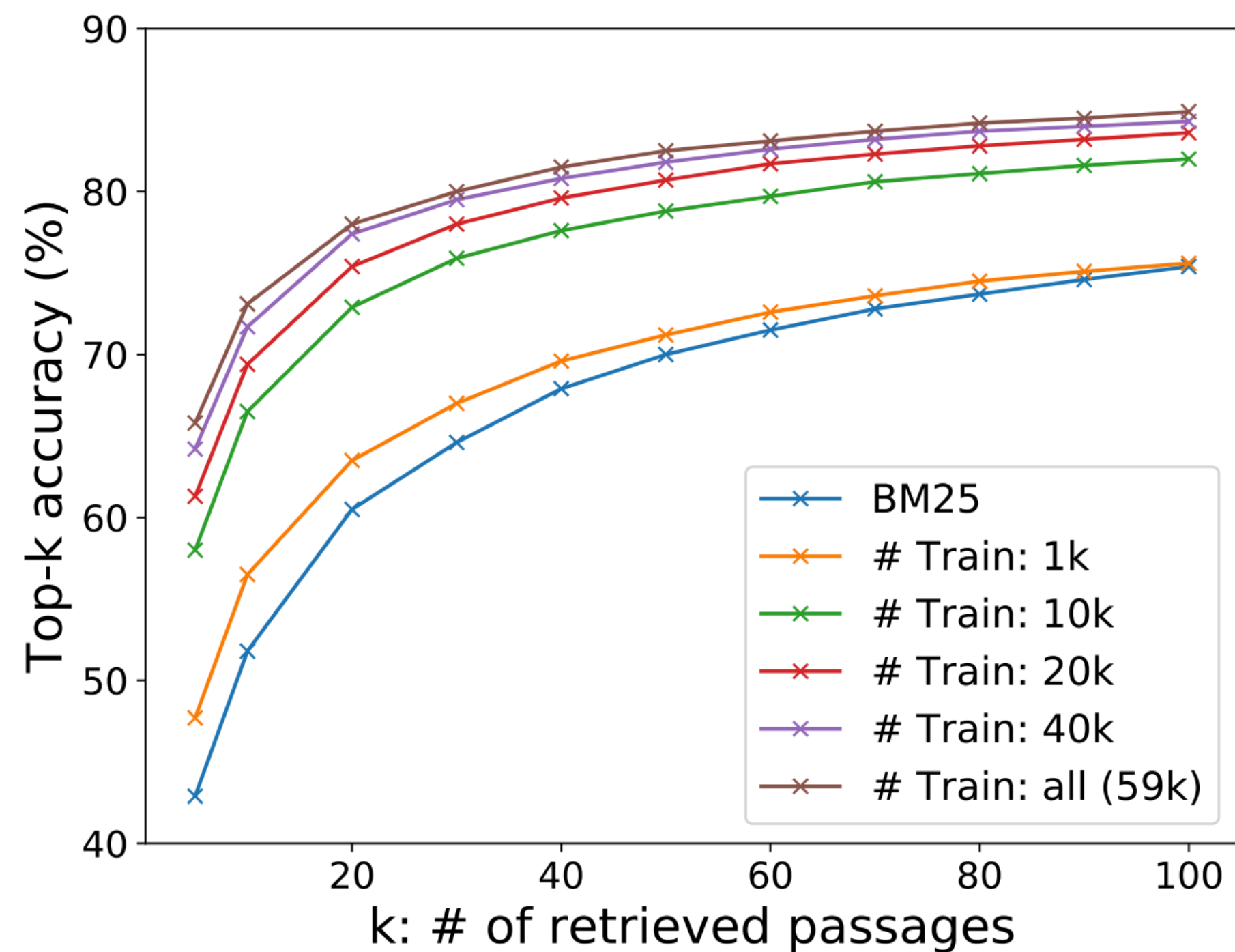
Versatility of RAG systems limited by conventional retrieval



Trained retrieval systems typically struggle in OOD

Perform well with more task training data

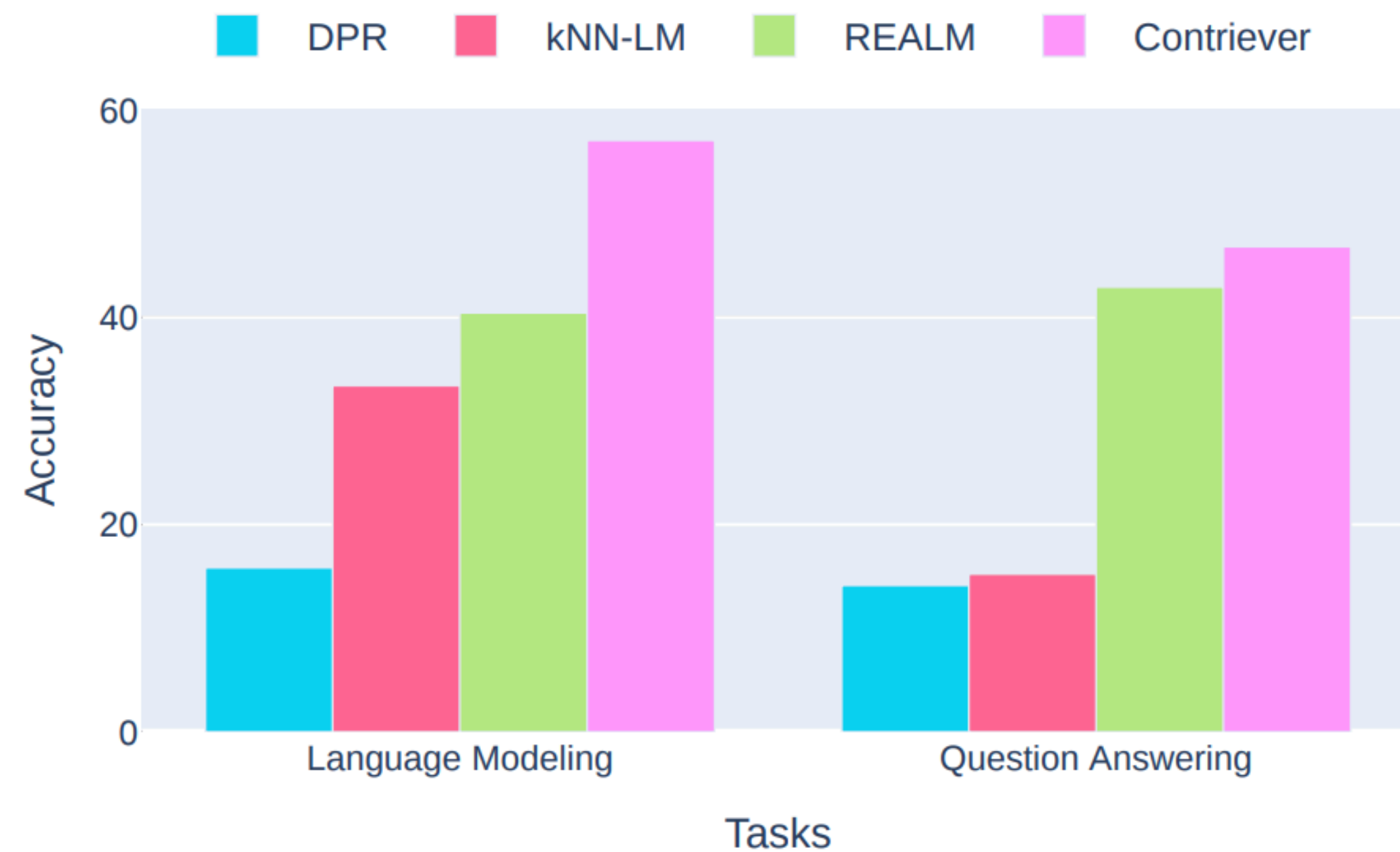
Struggle in new tasks & domains



Thakur, Nandan, et al. "BEIR: A heterogenous benchmark for zero-shot evaluation of information retrieval models." *In NeurIPS (Benchmark) 2021*.

“Similar” documents may not be always helpful

Helpful documents for some tasks aren't necessary similar to queries



In a zoo located in a warm region, what should be included in the polar bear exhibit? (Multi-hop reasoning task)

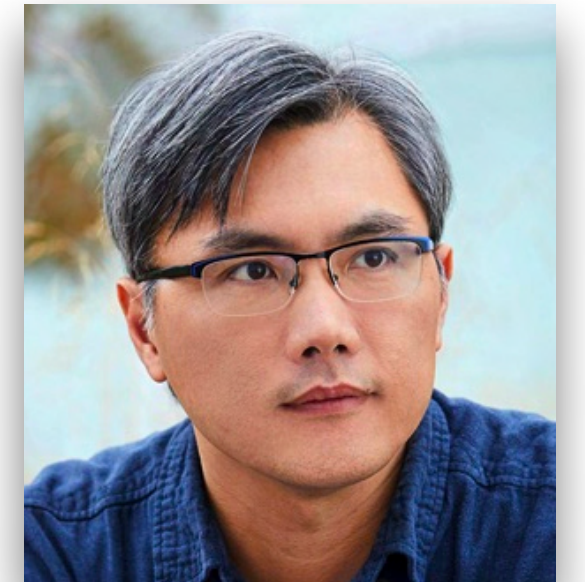
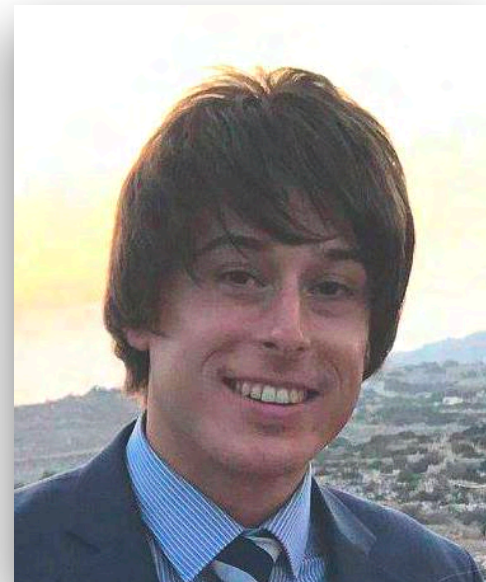
If an animal lives a certain environment then that animal usually requires that kind of environment.

Polar bears live in cold environments

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

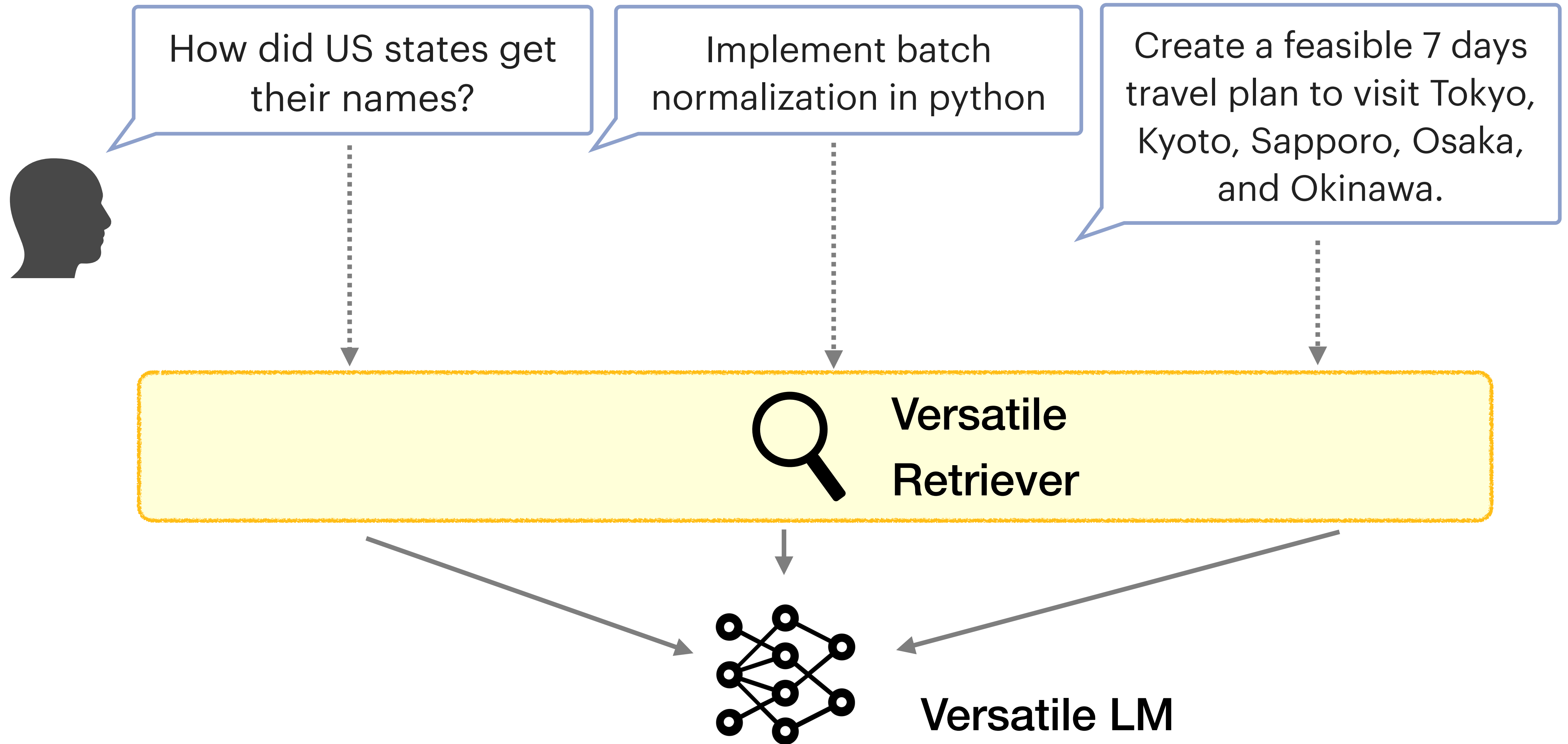
Task-aware Retrieval with Instructions

Akari Asai, Timo Schick, Patrick Lewis, Xilun Chen, Gautier Izacard,
Sebastian Riedel, Hannaneh Hajishirzi, Wen-tau Yih

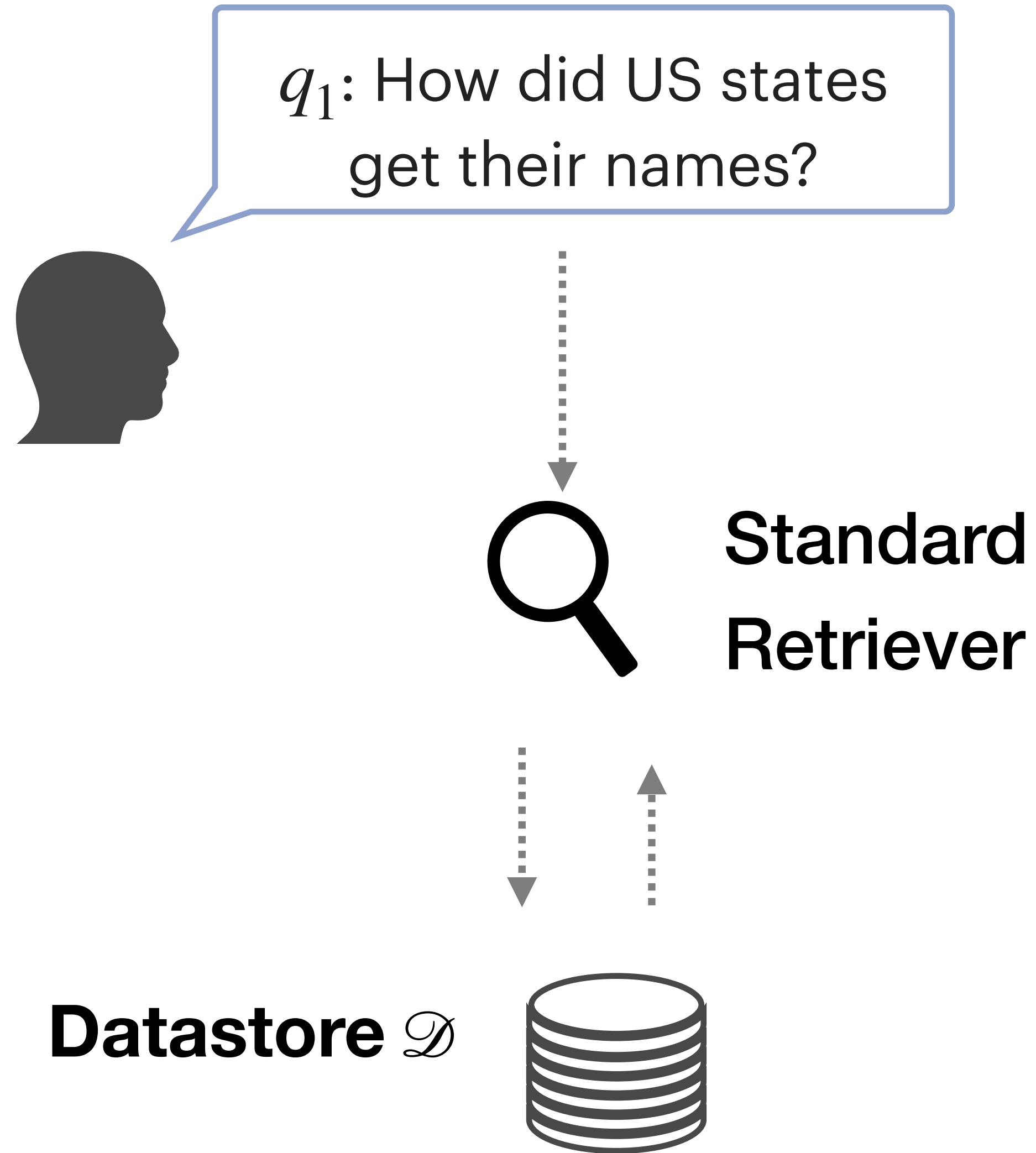


ACL Findings 2023

Goal of **Tsk-Aware ReTrieve** (**TART**)

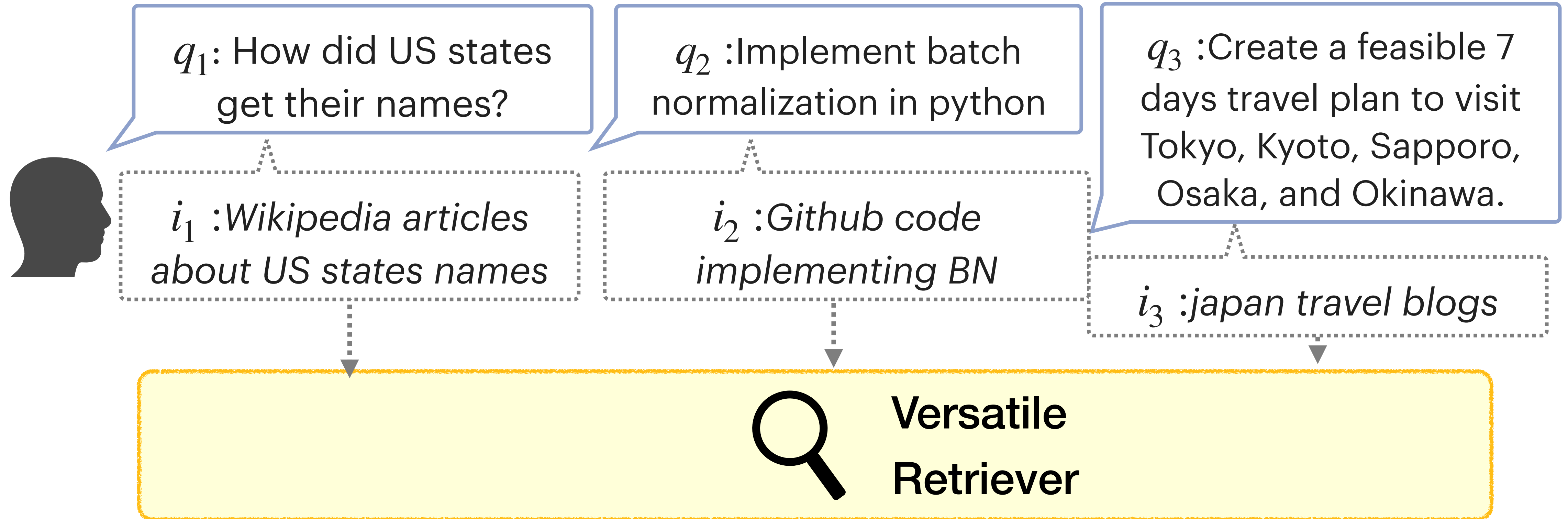


Normal retrieval task



$$\text{Topk } \arg\text{Topk}_{d \in \mathcal{D}} \text{sim}(q, d)$$

New task: Retrieval with Instruction



$$\text{Topk } \arg \text{Topk}_{d \in \mathcal{D}} \text{sim}(q, d, i)$$

New task: Retrieval with Instruction

q_2 : Implement batch normalization in python



i_1 : similar stack overflow questions

i_2 : Github code implementation of BN

How to implement batch normalization merging in python?

Asked 1 year, 5 months ago Modified 1 year, 5 months ago Viewed 382 times

stackoverflow
Do you use coding assistant tools?
Take our Coding Assistant & Productivity Survey

1

- The first layer of the model is a linear layer and there is **no bias**.
- The second layer of the model is a batch normalization and there is no weight and bias (**affine is false**)
- The third layer of the model is a linear layer.

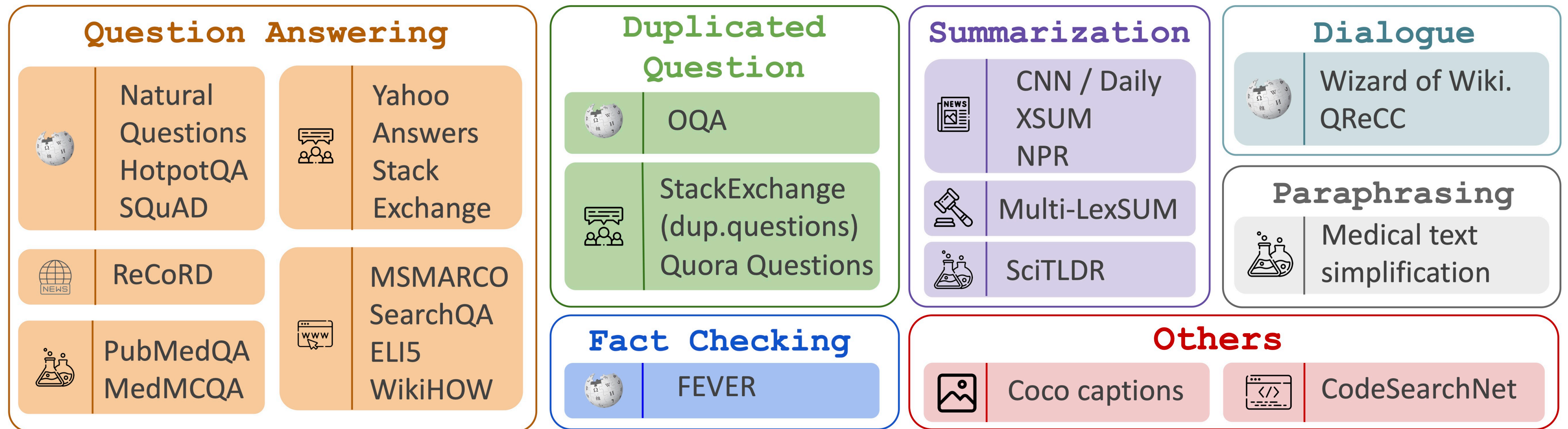
The variables named **new_weight** and **new_bias** are the weight and bias of the newly created linear layer, respectively.

```
78
79
80 def batchnorm_forward(x, gamma, beta, bn_param):
81     """
82     Forward pass for batch normalization.
83
84     During training the sample mean and (uncorrected) sample variance are
85     computed from minibatch statistics and used to normalize the incoming data.
86     During training we also keep an exponentially decaying running mean of the mean
87     and variance of each feature, and these averages are used to normalize data
88     at test-time.
```

$$\text{Topk } \arg \text{Topk}_{d \in \mathcal{D}} \text{sim}(q, d, i)$$

BERRI: first large-scale retrieval dataset with instructions

Curated 50 tasks with expert annotation instructions across domains



Instruction-scheme for retrieval tasks

Propose effective instruction scheme for retrieval tasks
and annotate instructions

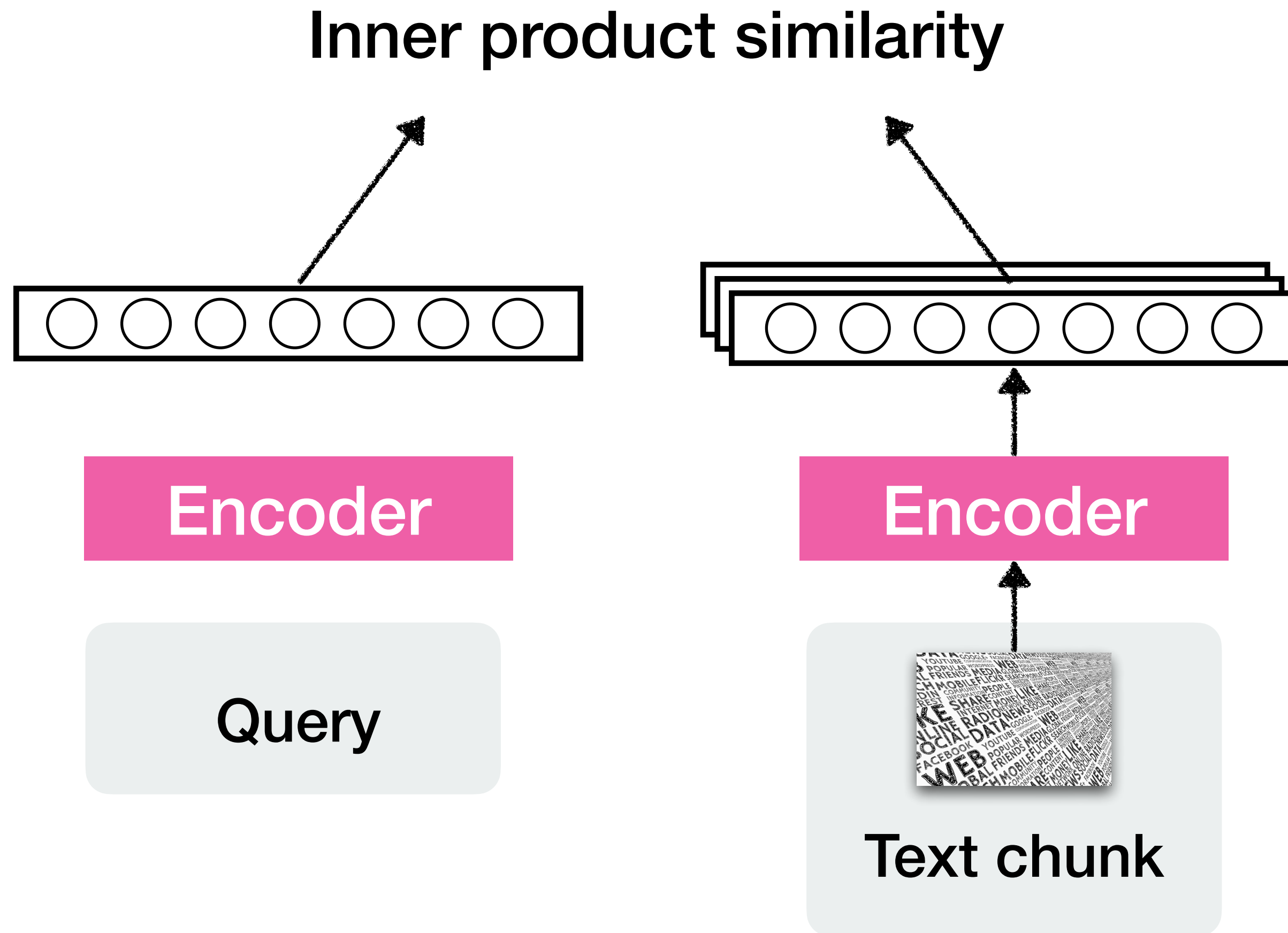
Dataset	Instruction
NQ	Retrieve a Wikipedia paragraph that answers this question .
QReCC	Find a dialogue response from dialogue history to answer the user's question .
Arguana	Retrieve a paragraph from an argument website that argues against the following argument .
SciFact	Find a sentence from a scientific paper to check if the statement is correct or not .
MultiLexSum	I want to find the one-sentence summary of this legal case .

Intent

Domain

Unit

Instruction-aware bi-encoder retriever (TART-dual)



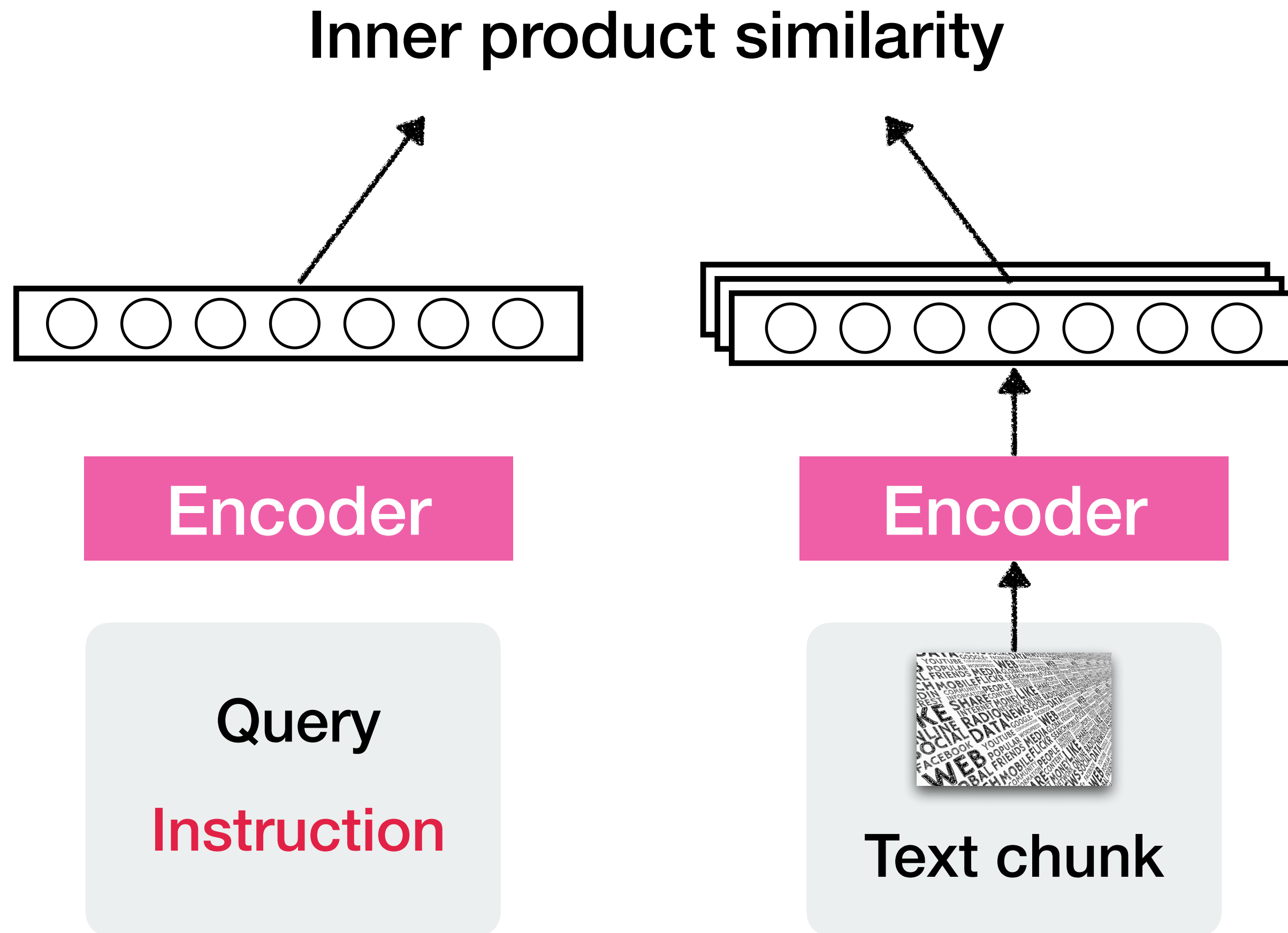
Positive paragraph Negative paragraphs

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

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

Contrastive learning

Instruction-aware bi-encoder retriever (TART-dual)



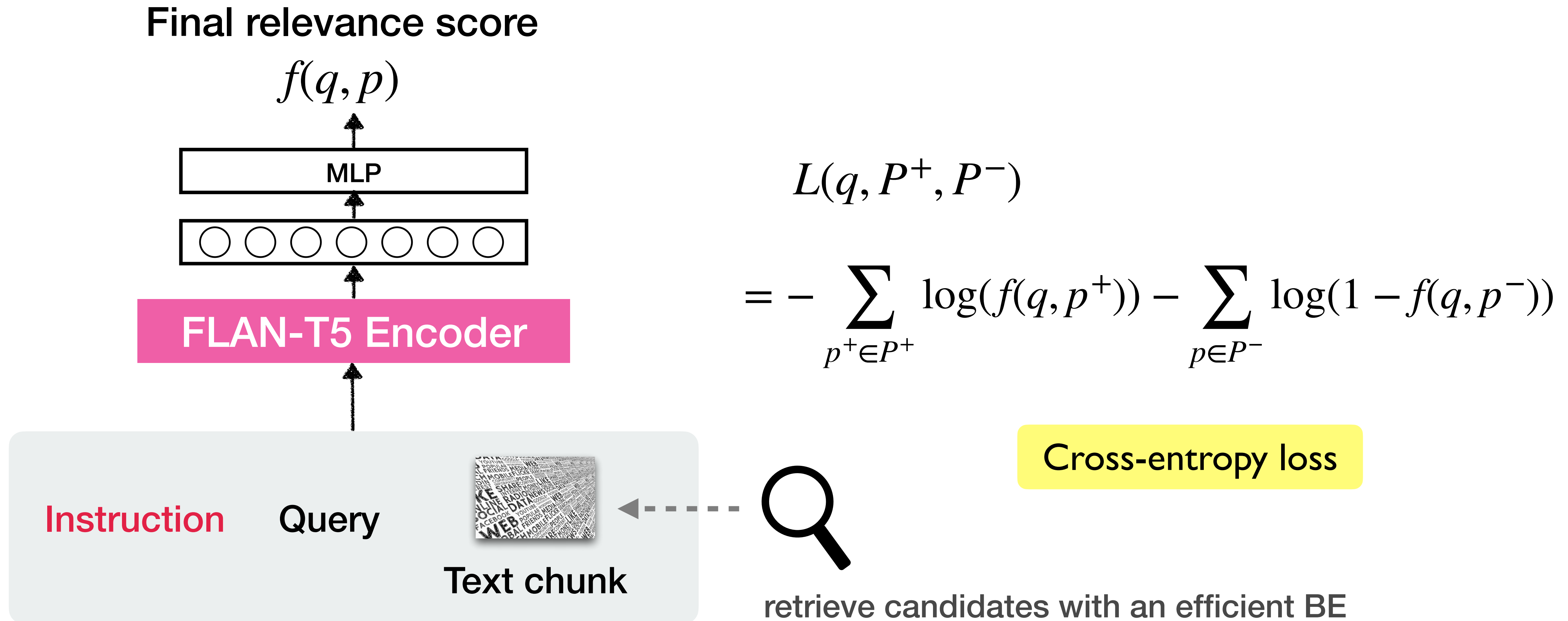
Positive paragraph Negative paragraphs

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

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

Contrastive learning

Instruction-aware cross-encoder retriever (TART-full)



New negative samples: instruction *un*-following samples

Carefully defining negative samples is a key for successful retrievers

Dup. Question Retrieval

t_1 : Retrieve a question asked in StackOverflow similar to this

q: How to compute square root in iOS?

How can we calculate the square root in Objective C or Swift?
StackOverflow Question

Hard negative documents d^{HD}

Which python function can I use to compute sq root?
StackOverflow Question

Instruction-unfollowing negatives d^{UF}

You can just use the Objective C or Swift's `sqrt` function
StackOverflow Answer

Dialogue Response Retrieval

t_1 : Find an informative dialogue response to this user's conversation

q: Are armadillos native to a Spanish-speaking part of the world?

Yes, they are most commonly found in North, Central, and South America.
Dialogue Response

I love animals and think armadillos are awesome with their leathery shell.
Dialogue Response

Armadillos are medium-sized mammals found in North, Central, and South America
Wikipedia Paragraph

Tasks

Gold documents d^+

Negative documents d^-

Follow instruction?



Relevant to the query?



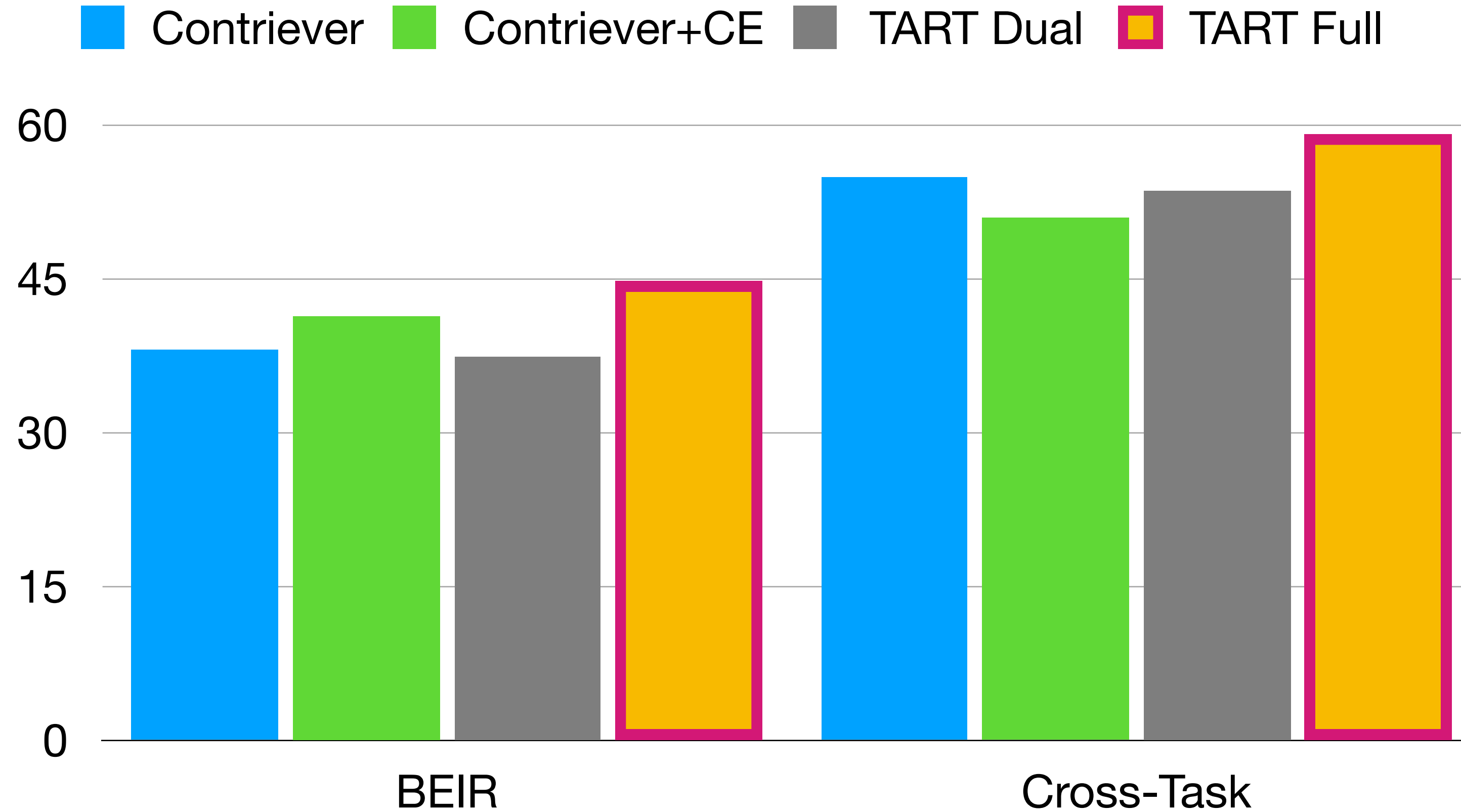
Experimental settings

More details of training & test are in our paper!

- Our models
 - **TART-dual (bi-encoder)** - Instruction-tuning of Contriever-MS MARCO 110M
 - **TART-full (cross-encoder)** - instruction-tuning FLAN-T5 3B Encoder (1.5B)
- Evaluations
 - **Zero-shot retrieval:** generalize to new retrieval task via instruction
 - **BEIR** (Thakur et al. 2021)
 - **LOTTE** (Khattab et al., 2022)
 - **Cross-task cross-domain retrieval:** synthetically combine retrieval tasks to test instruction following by pairing two relevant tasks (e.g., QA & question retrieval)

Better Generalization and Instruction Following

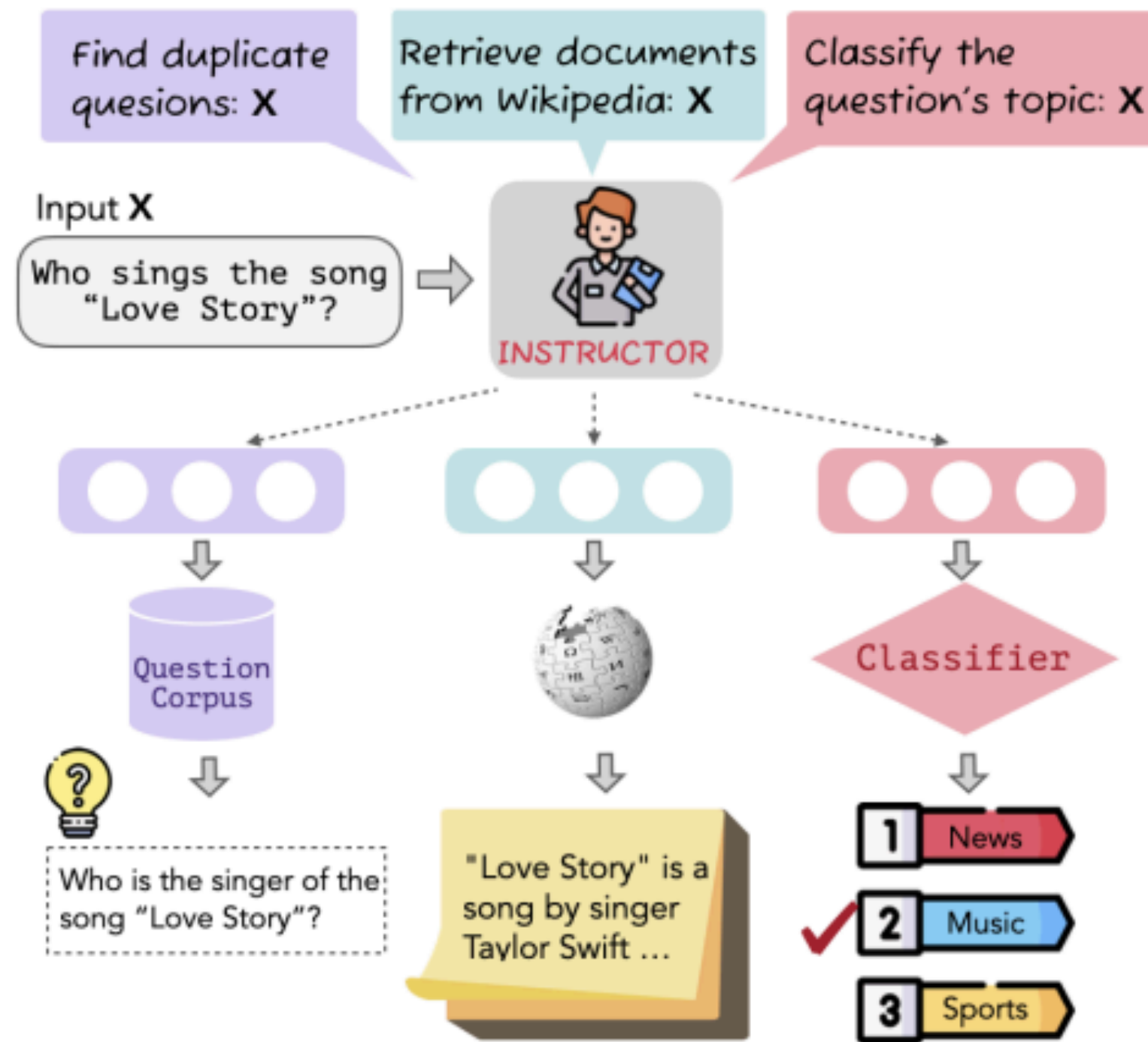
Instruction-tuning for retrieval (3B) shows effectiveness



Increasing number of instruction-following retrievers

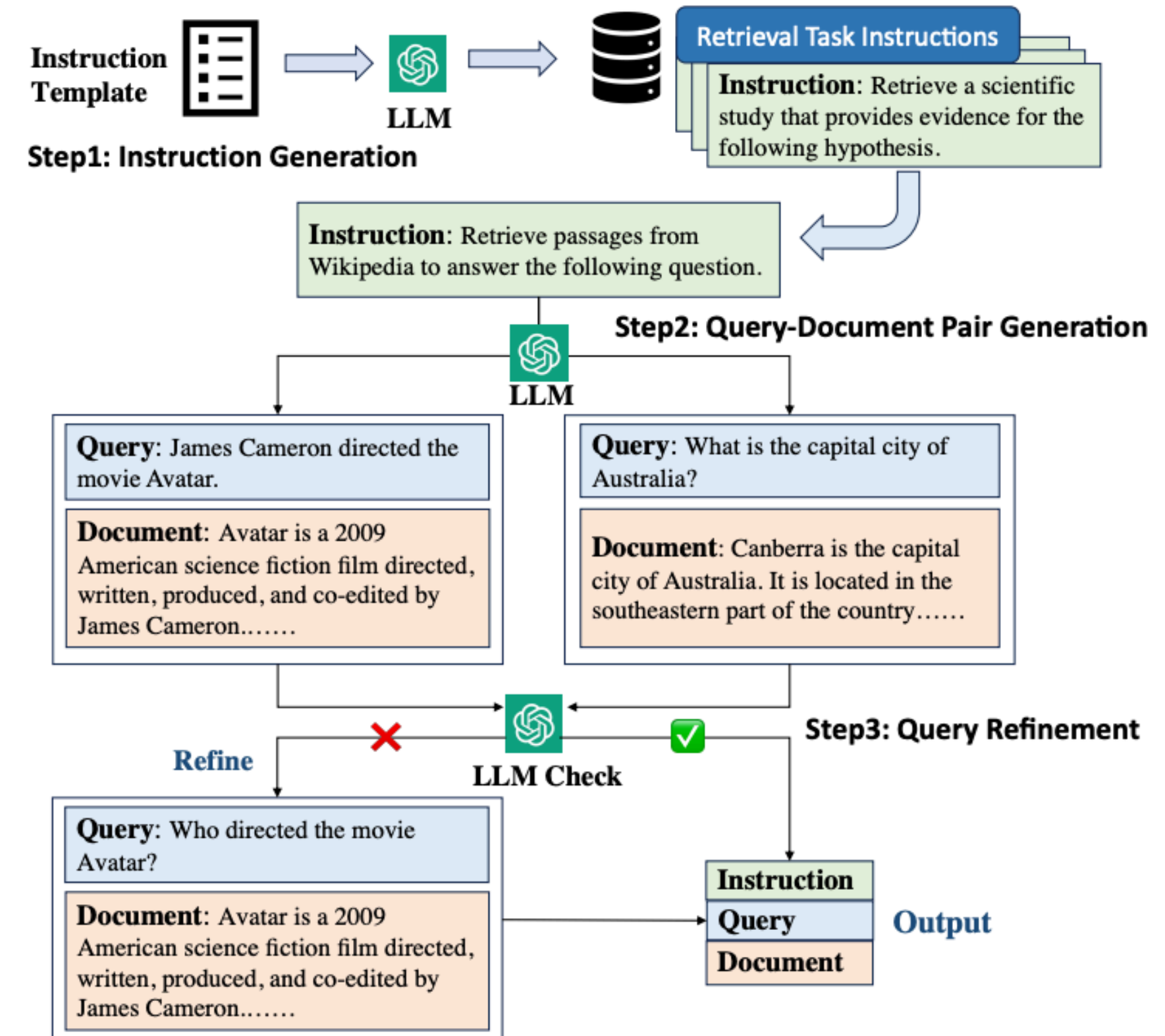
Instructor (Su et al., 2023)

Concurrent work focusing on embeddings trained on existing datasets



ControlRetriever (Pang et al., 2023)

Leverage LLMs to generate training data
E5 MISTRAL Instruct (Wang et al., 2024)



Increasing number of instruction-following retrievers

Instruction following retrievers are now dominating retrieval tasks!

mteb/leaderboard like 2.49k Running on CPU UPGRADE App Files

Massive Text Embedding Benchmark (MTEB) Leaderboard. To submit, refer to the [MTEB GitHub repository](#). Refer to the [MTEB paper](#) for details on metrics, tasks and models.

Search Bar (separate multiple queries with `;`)
Search for a model and press enter...

Model types
 Open Proprietary Sentence Transformers

Model sizes (in number of parameters)
 <100M 100M to 250M 250M to 500M 500M to 1B >1B

Overall Bitext Mining Classification Clustering Pair Classification Reranking Retrieval STS Summarization

English Chinese French Polish

Overall MTEB English leaderboard 🏆

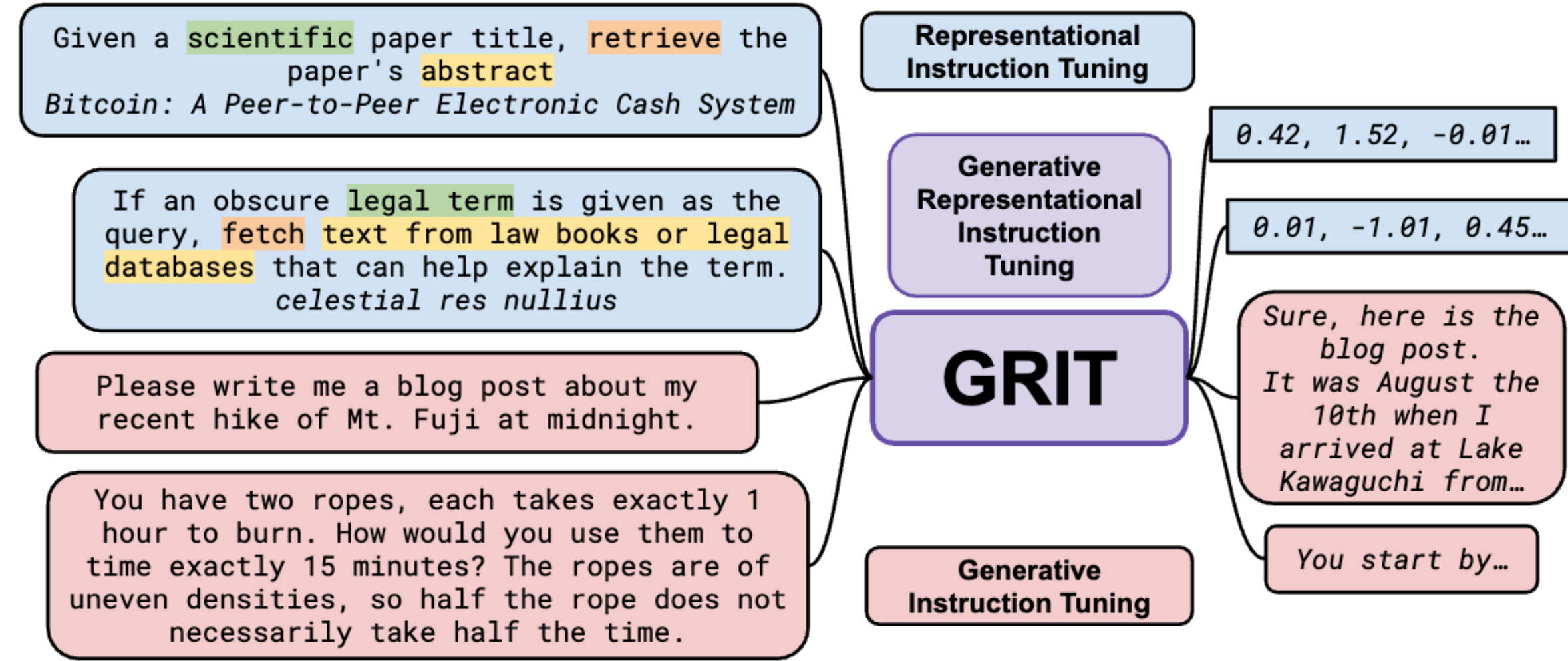
- Metric: Various, refer to task tabs
- Languages: English

Rank	Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens	Average (56 datasets)	Classification Average (12 datasets)	Clustering Average (11 datasets)	Pair Classification Average (3 datasets)	Reranking Average (4 datasets)
1	SFR-Embedding-Mistral	7111	26.49	4096	32768	67.56	78.33	51.67	88.54	60.64
2	voyage-lite-02-instruct	1220	4.54	1024	4000	67.13	79.25	52.42	86.87	58.24
3	GritLM-7B	7242	26.98	4096	32768	66.76	79.46	50.61	87.16	60.49
4	e5-mistral-7b-instruct	7111	26.49	4096	32768	66.63	78.47	50.26	88.34	60.21

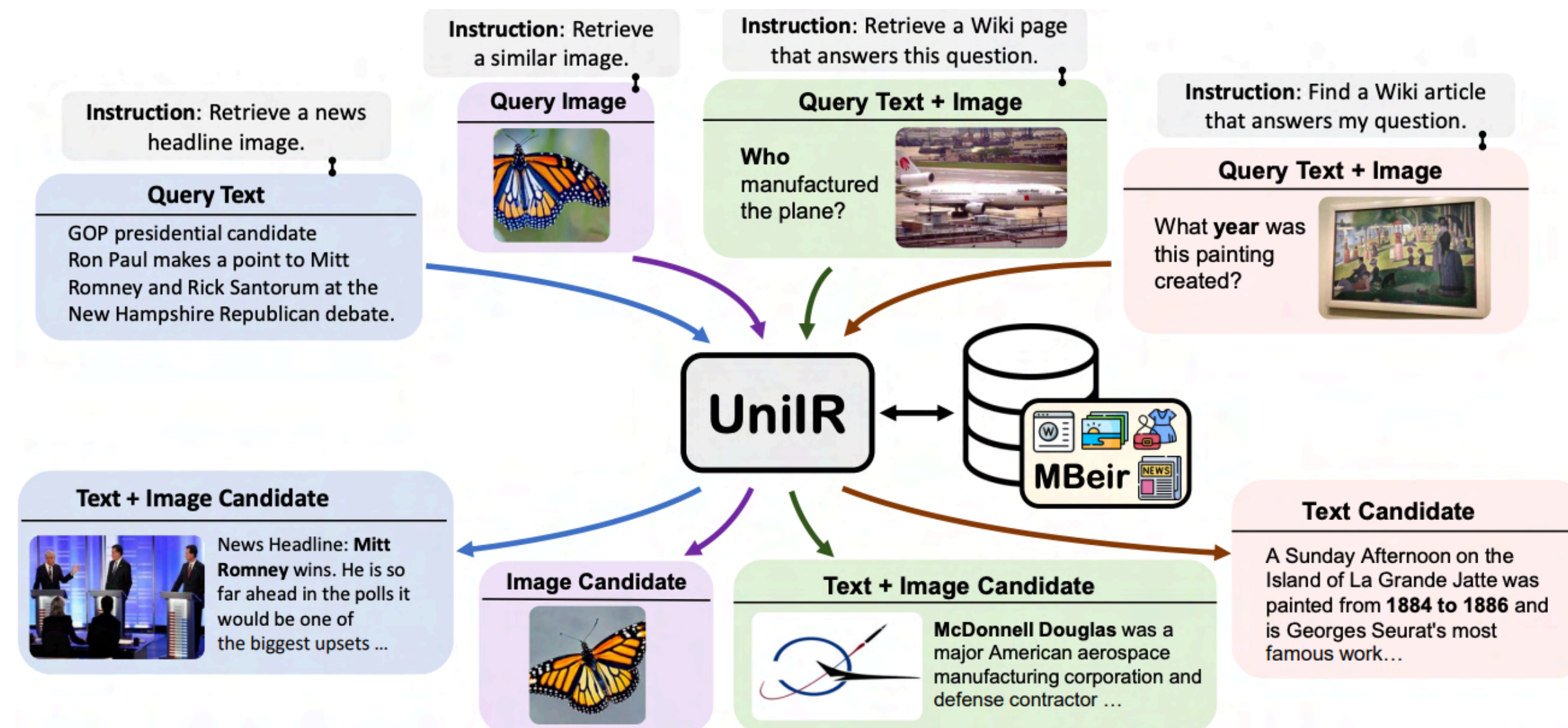
<https://huggingface.co/spaces/mteb/leaderboard>

Versatile retrievers improves RAG in text and multi-modal

GRIT (Muennighoff et al., 2024)



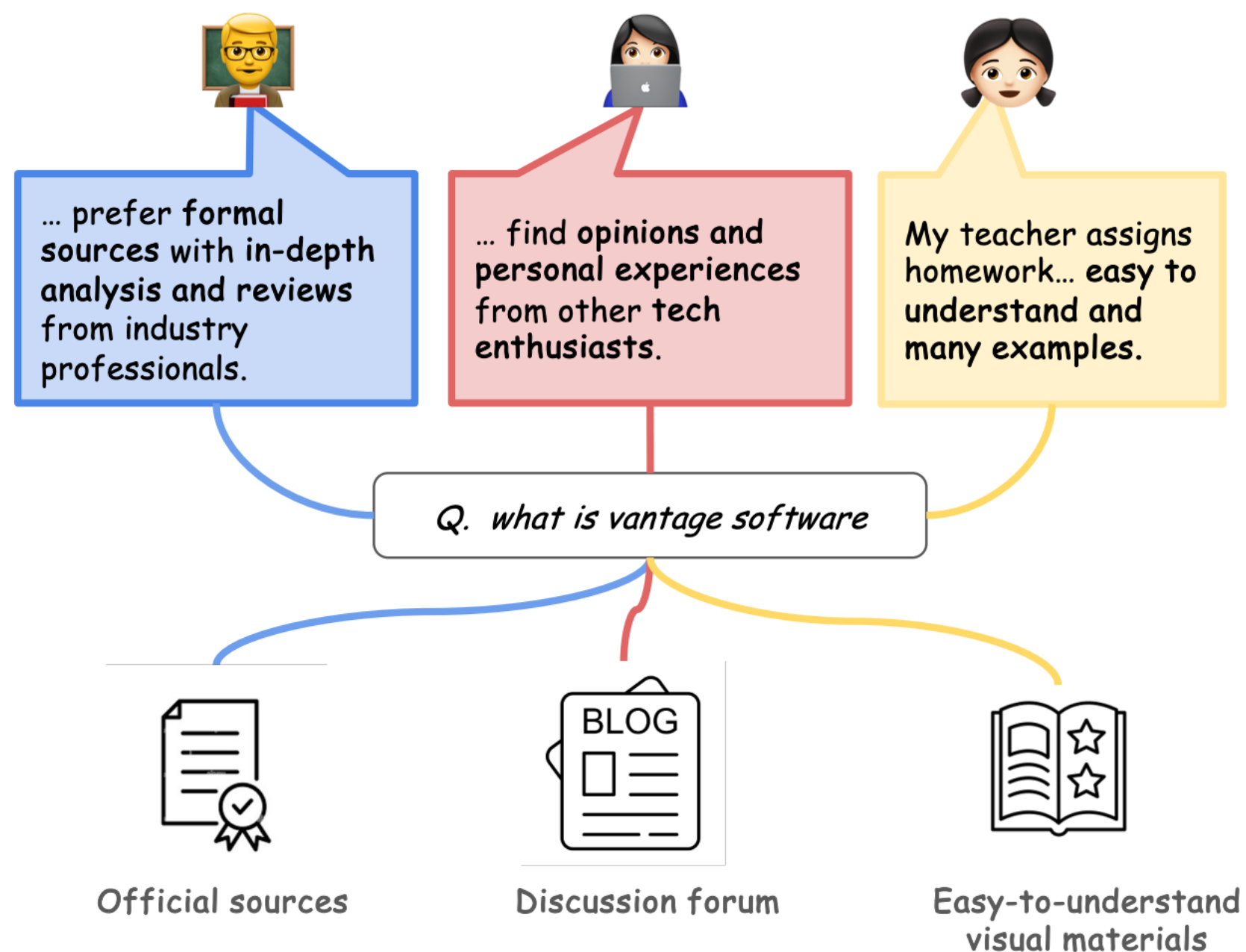
UniLR (Wei et al., 2023)



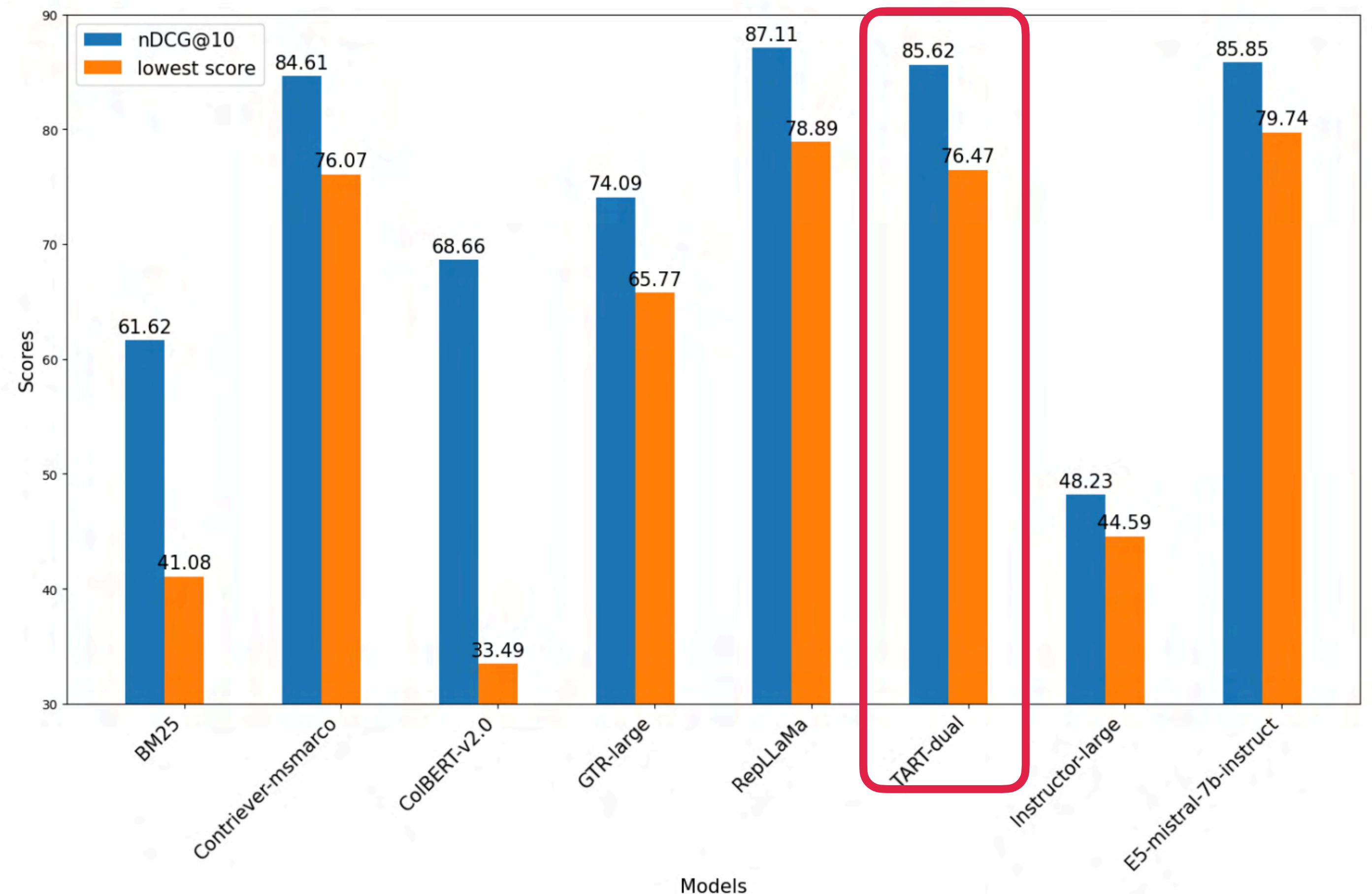
Are those retrievers following instructions?

Instructir (Oh et al., 2024)

Create a natural instruction-following retrieval dataset with LLMs



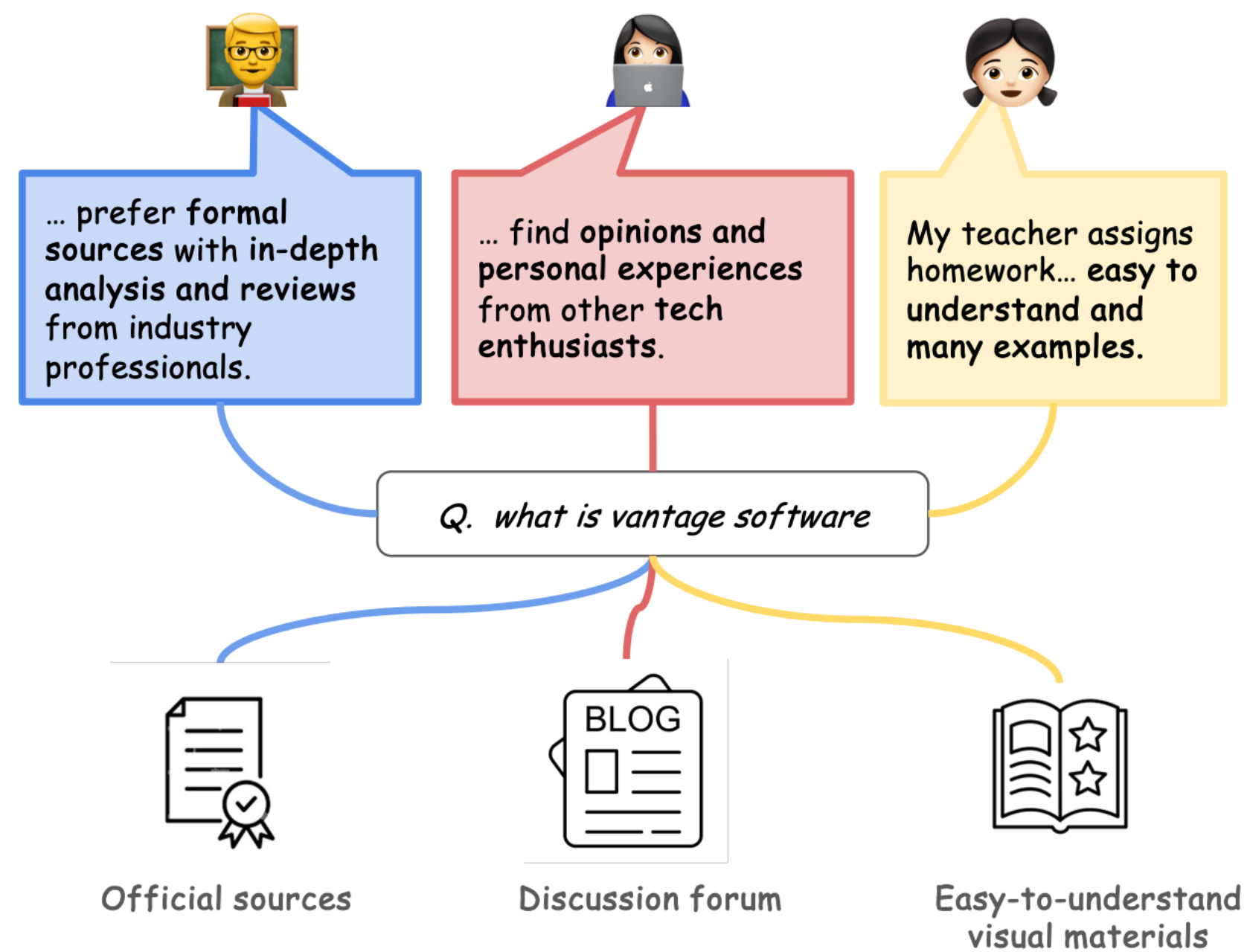
TART (110M) outperforms other <7B models by large margins



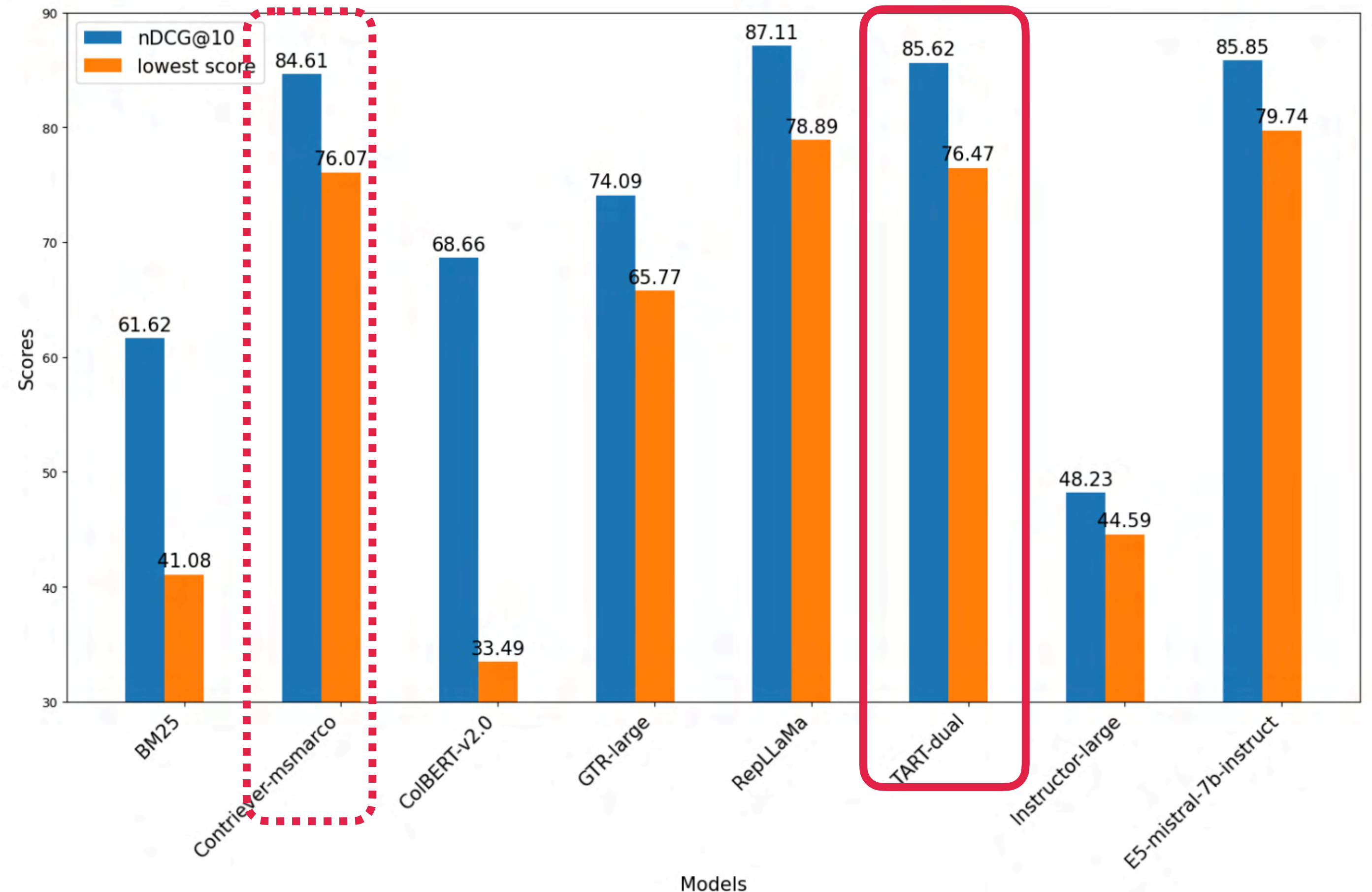
Are those retrievers following instructions?

Instructir (Oh et al., 2024)

Create a natural instruction-following retrieval dataset with LLMs



Improvements from base retrievers are limited



Are those retrievers following instructions?

FollowIR (Welleer et al., 2024)

Model		Robust04		News21		Core17		Average	
		mAP	<i>p</i> -MRR	nDCG	<i>p</i> -MRR	mAP	<i>p</i> -MRR	Score	<i>p</i> -MRR
No-Instruction IR	BM25	12.2	-2.8	21.3	+2.5	8.1	-0.7	13.9	-0.3
	E5-base-v2	14.5	-6.8	21.6	-4.1	14.0	-2.9	16.7	-4.6
	E5-large-v2	18.1	-4.1	24.9	-2.2	17.0	+0.1	20.0	-2.1
	Contriever	20.3	-6.1	24.0	-1.8	15.3	-2.5	19.9	-3.5
	MonoBERT	21.5	-9.7	26.3	-4.4	18.4	-1.3	22.1	-5.1
	MonoT5-base	16.3	-5.8	11.9	-1.2	12.2	-3.5	13.5	-3.5
	MonoT5-3B	27.8	+5.6	18.6	+7.5	18.1	+1.7	21.5	+4.9
Instruction-IR	BGE-base	17.5	-6.4	23.8	-0.2	14.6	-2.7	18.6	-3.1
	BGE-large	18.1	-7.8	26.4	+0.1	15.0	+0.1	19.8	-2.5
	TART-Contriever	14.1	-7.8	21.9	+0.0	12.4	-1.3	16.1	-3.0
	INSTRUCTOR-base	14.4	-5.6	16.3	-2.5	14.7	-2.2	15.1	-3.4
	INSTRUCTOR-xl	15.5	-2.1	14.6	-4.3	14.4	-0.6	14.8	-2.3
	TART-FLAN-T5-xl	25.2	-0.8	20.3	-1.1	17.0	+2.8	20.8	+0.3
	GritLM-7B	29.0	-1.4	25.2	+2.1	20.8	+2.6	25.0	+1.1
APIs	Cohere v3 English	22.9	-3.3	23.6	-3.1	20.6	+2.7	22.4	-1.2
	OpenAI v3 Large	27.9	-5.7	30.0	-3.3	21.4	-0.2	26.4	-3.1
Instruct LMs	FLAN-T5-base	6.8	+5.0	2.2	+1.1	6.5	-3.2	5.2	+1.0
	FLAN-T5-large	15.1	+4.0	8.5	+7.7	11.5	+1.2	11.7	+4.3
	Llama-2-7B-chat	6.9	+1.6	13.3	+2.1	5.4	+3.6	8.5	+2.4
	Mistral-7B-instruct	24.1	+12.2	22.9	+10.5	19.6	+13.4	22.2	+12.0
	FollowIR-7B	25.9	+13.6	25.7	+10.8	20.0	+16.3	23.9	+13.6

All small BE models struggles to follow instructions

Table 3: Evaluating instruction-following on FOLLOWIR. *p*-MRR is a new pairwise evaluation metric measuring instruction following when instructions change, ranging from -100 to 100 (higher is better). We see that the only models that show any success at following instructions are large models (3B+ parameters) or instruction-tuned LLMs that haven't been trained on retrieval tasks.

Are those retrievers following instructions?

FollowIR (Welleer et al., 2024)

	Model	Robust04		News21		Core17		Average	
		mAP	<i>p</i> -MRR	nDCG	<i>p</i> -MRR	mAP	<i>p</i> -MRR	Score	<i>p</i> -MRR
No-Instruction IR	BM25	12.2	-2.8	21.3	+2.5	8.1	-0.7	13.9	-0.3
	E5-base-v2	14.5	-6.8	21.6	-4.1	14.0	-2.9	16.7	-4.6
	E5-large-v2	18.1	-4.1	24.9	-2.2	17.0	+0.1	20.0	-2.1
	Contriever	20.3	-6.1	24.0	-1.8	15.3	-2.5	19.9	-3.5
	MonoBERT	21.5	-9.7	26.3	-4.4	18.4	-1.3	22.1	-5.1
	MonoT5-base	16.3	-5.8	11.9	-1.2	12.2	-3.5	13.5	-3.5
	MonoT5-3B	27.8	+5.6	18.6	+7.5	18.1	+1.7	21.5	+4.9
Instruction-IR	BGE-base	17.5	-6.4	23.8	-0.2	14.6	-2.7	18.6	-3.1
	BGE-large	18.1	-7.8	26.4	+0.1	15.0	+0.1	19.8	-2.5
	TART-Contriever	14.1	-7.8	21.9	+0.0	12.4	-1.3	16.1	-3.0
	INSTRUCTOR-base	14.4	-5.6	16.3	-2.5	14.7	-2.2	15.1	-3.4
	INSTRUCTOR-xl	15.5	-2.1	14.6	-4.3	14.4	-0.6	14.8	-2.3
	TART-FLAN-T5-xl	25.2	-0.8	20.3	-1.1	17.0	+2.8	20.8	+0.3
	GritLM-7B	29.0	-1.4	25.2	+2.1	20.8	+2.6	25.0	+1.1
APIs	Cohere v3 English	22.9	-3.3	23.6	-3.1	20.6	+2.7	22.4	-1.2
	OpenAI v3 Large	27.9	-5.7	30.0	-3.3	21.4	-0.2	26.4	-3.1
Instruct LMs	FLAN-T5-base	6.8	+5.0	2.2	+1.1	6.5	-3.2	5.2	+1.0
	FLAN-T5-large	15.1	+4.0	8.5	+7.7	11.5	+1.2	11.7	+4.3
	Llama-2-7B-chat	6.9	+1.6	13.3	+2.1	5.4	+3.6	8.5	+2.4
	Mistral-7B-instruct	24.1	+12.2	22.9	+10.5	19.6	+13.4	22.2	+12.0
	FollowIR-7B	25.9	+13.6	25.7	+10.8	20.0	+16.3	23.9	+13.6

CE models or 7B BE models have shown superiority

Introduce significant inference latency!

Table 3: Evaluating instruction-following on FOLLOWIR. *p*-MRR is a new pairwise evaluation metric measuring instruction following when instructions change, ranging from -100 to 100 (higher is better). We see that the only models that show any success at following instructions are large models (3B+ parameters) or instruction-tuned LLMs that haven't been trained on retrieval tasks.

Today's lecture

Promises and Limitations of Retrieval-augmented LMs

Advanced Retriever: Intent-aware LM-based retrievers

Advanced RAG: Self-reflective LMs with dynamic Retrievals

Summary and **future directions**: RAG in the wild

Summary

- **Understanding** retrieval-augmented LMs (Asai et al., 2024b; Mallen*, Asai et al., 2023)
 - Retrieval-augmented LMs can alleviate many issues in parametric LMs
 - More fundamental improvements for architectures or training is necessary

Summary

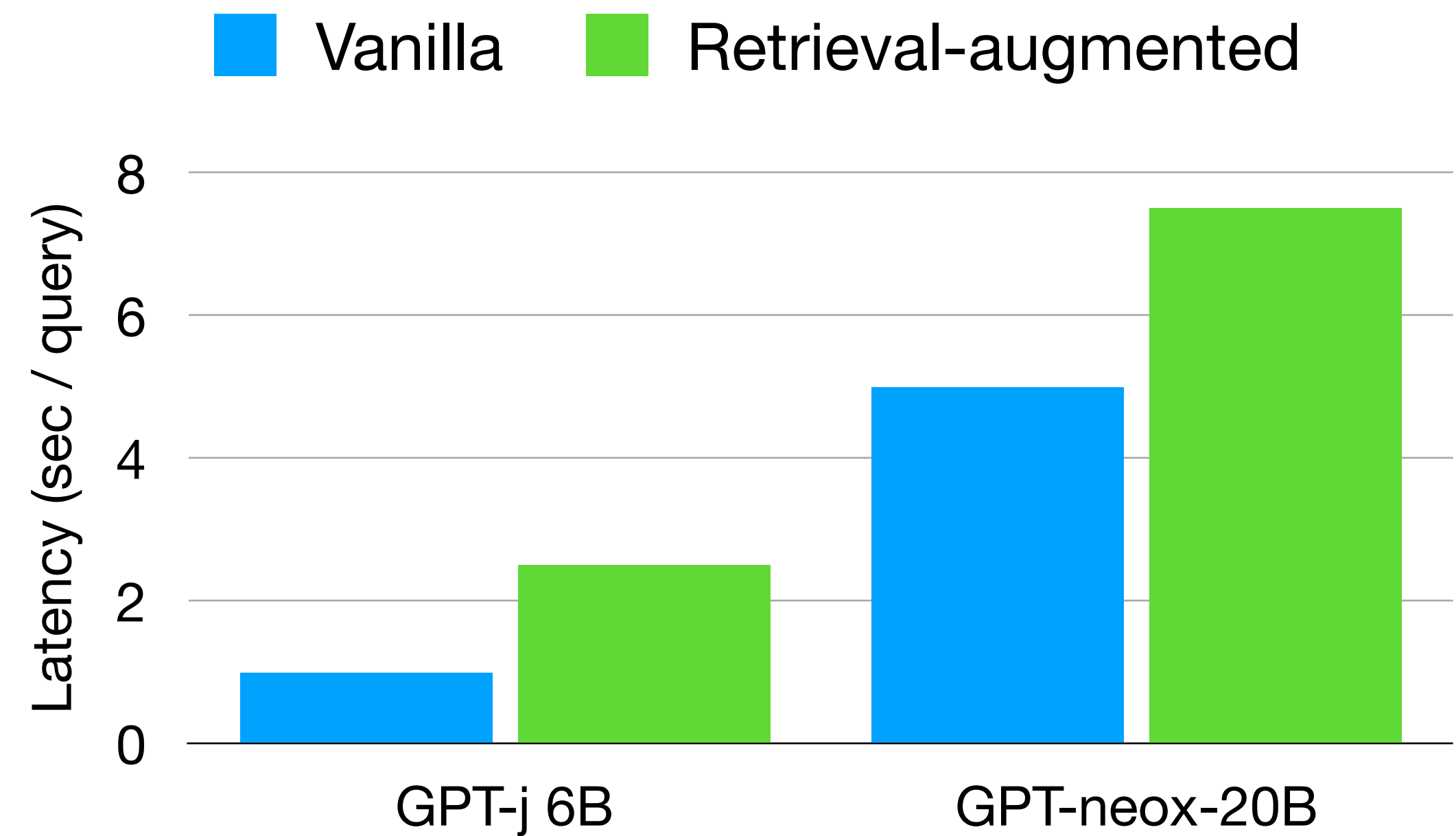
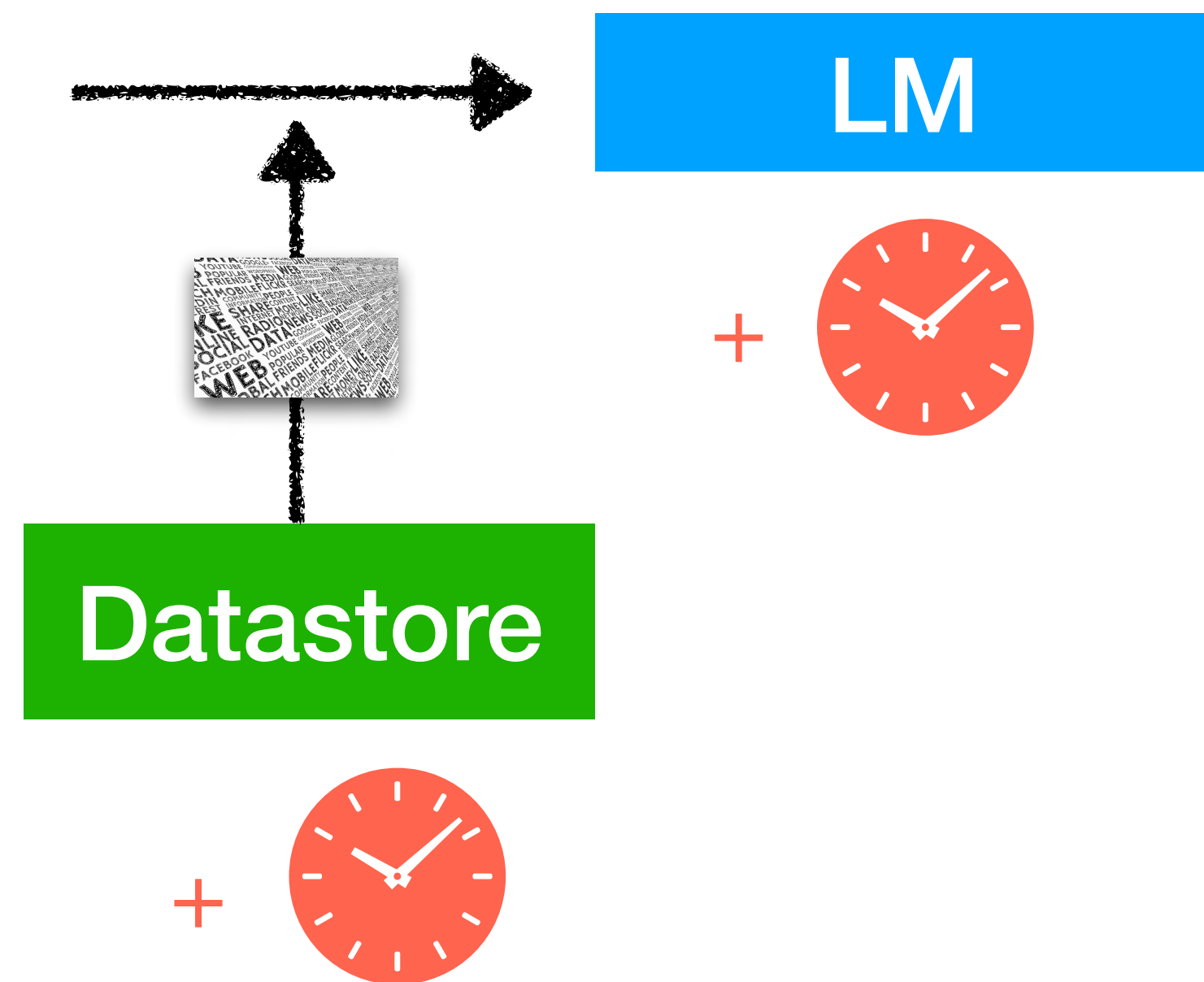
- **Understanding** Retrieval-augmented LMs (Asai et al., 2024b; Mallen*, Asai et al., 2023)
 - Retrieval-augmented LMs can alleviate many issues in parametric LMs.
 - More fundamental improvements for architectures or training is necessary
- **Advancing RAG** (Asai et al., 2024; Asai et al., 2023)
 - **Self-RAG** to build versatile retrieval-augmented LMs addressing issues in RAG
 - **Task-aware retrievals** to build versatile RAG systems

Reliable RAG in the wild: improving efficiency

Efficiency

Effective datastore

New application



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

Reliable RAG in the wild: efficient algorithms / arc. for RAG

Efficiency

Effective datastore

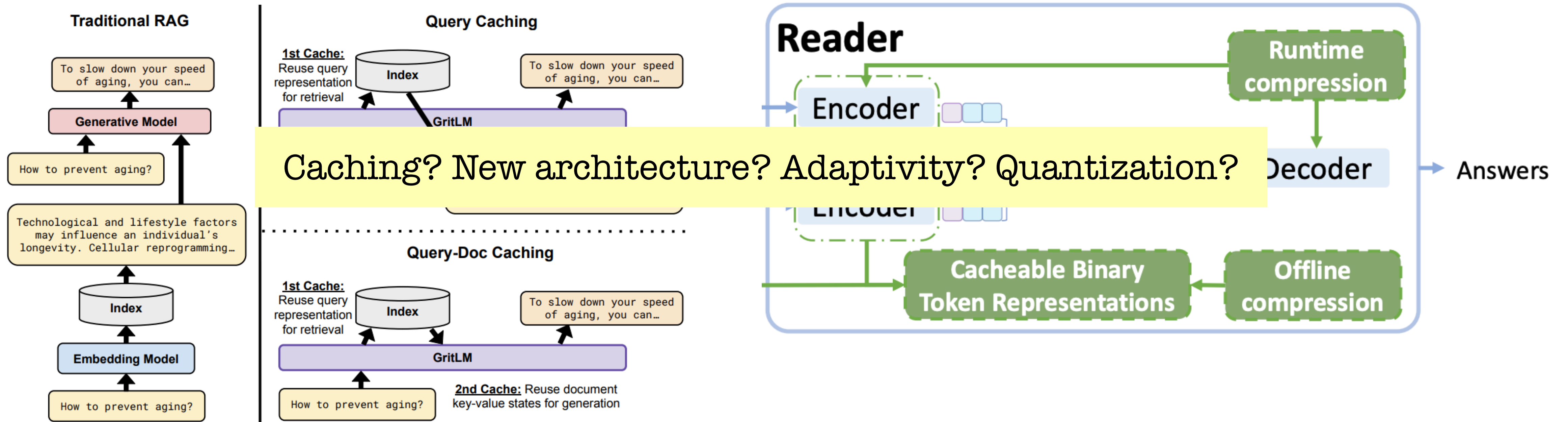
New application

GRIT

(Muennighoff et al., 2024)

BTR

(Cao et al., 2024)

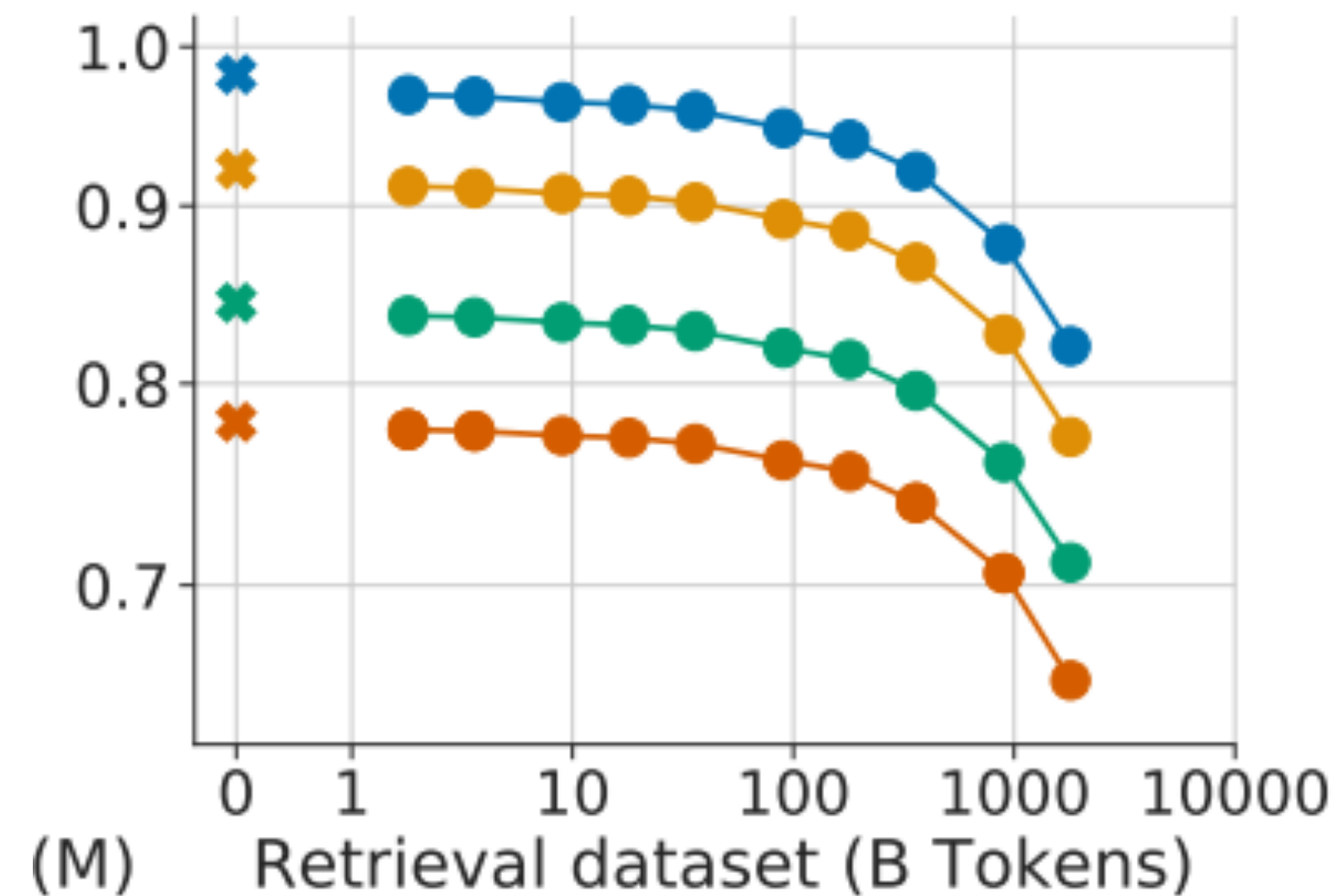
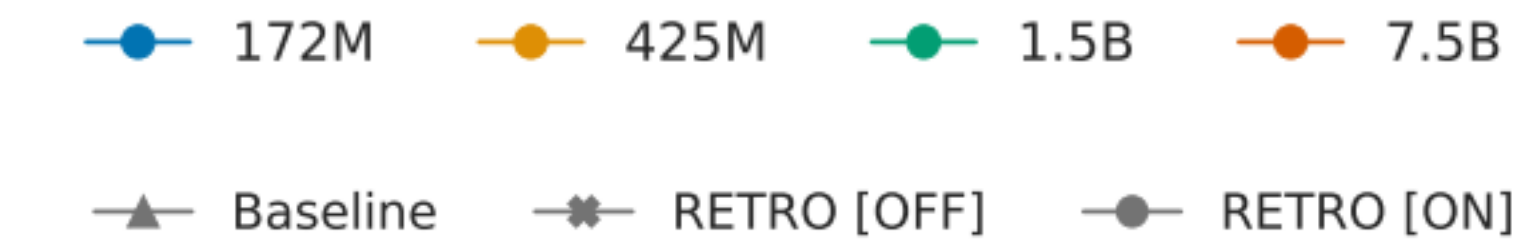
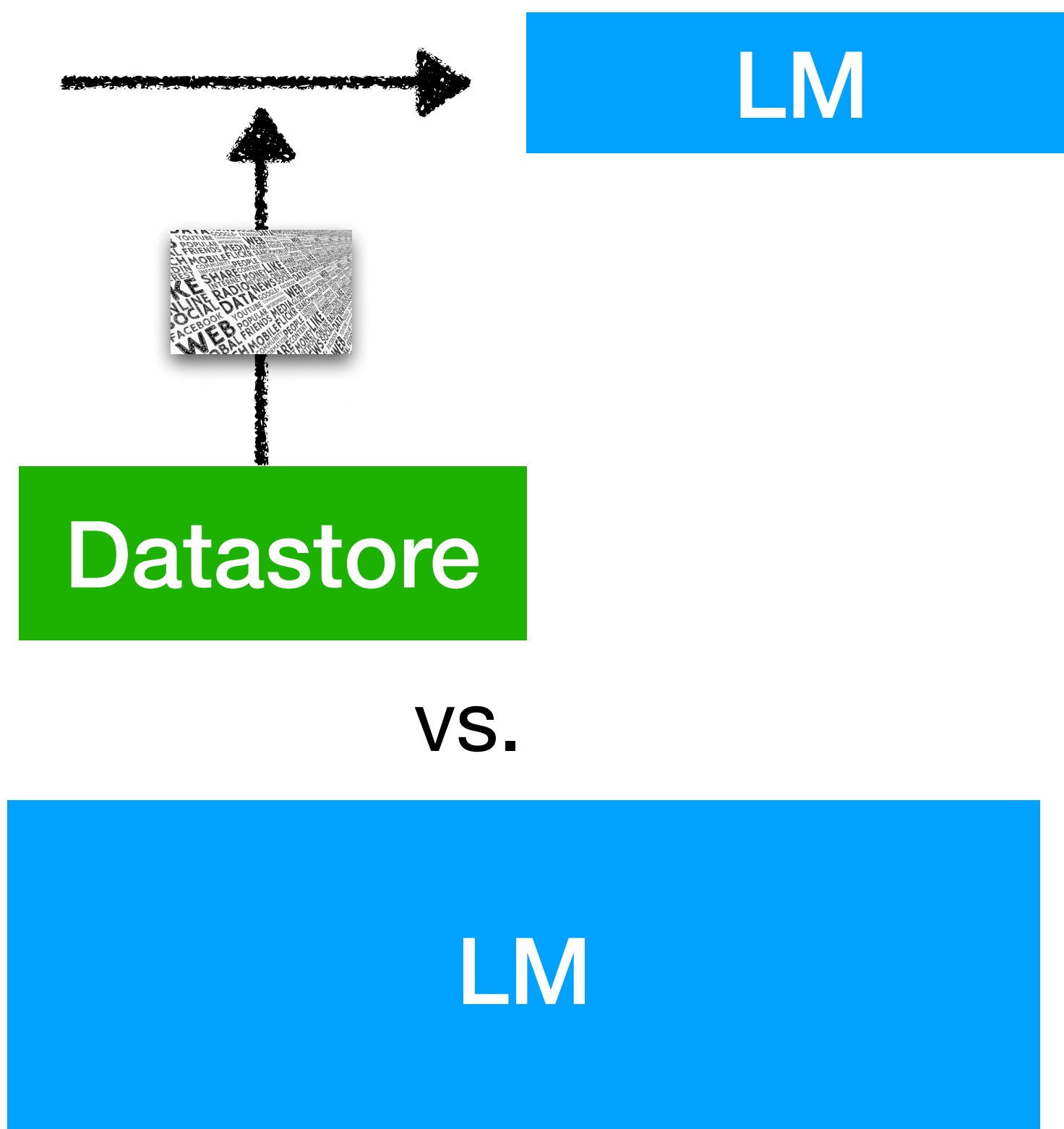


Reliable RAG in the wild: effective datastore (scale)

Efficiency

Effective datastore

New application



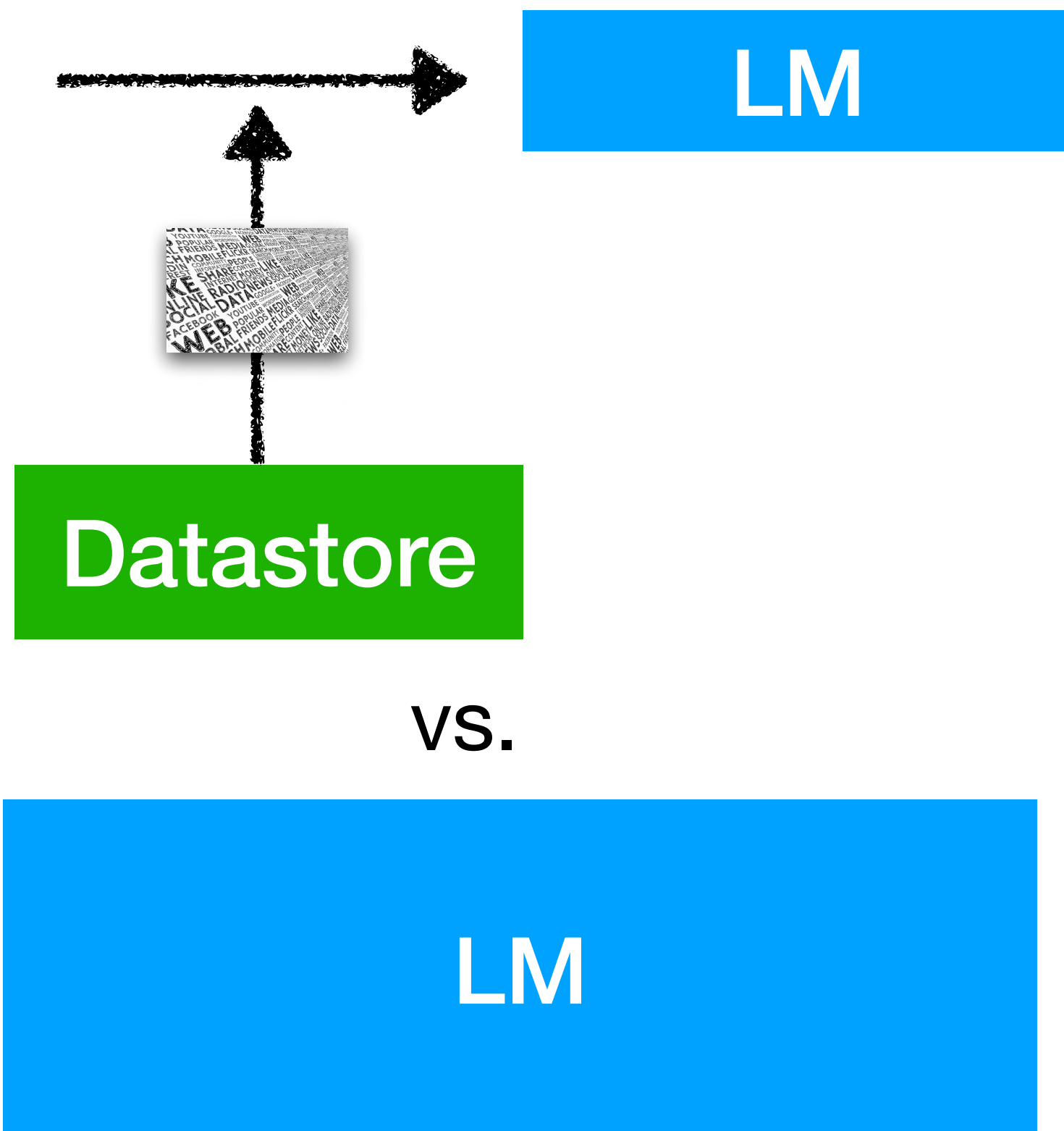
RETRO
(Borgeaud et al., 2021)

Reliable RAG in the wild: effective datastore (scale)

Efficiency

Effective datastore

New application



How should we scale RAG?

LM

Datastore

kNN-LM (Khandelwal et al., 2020)

of parameters

of tokens

250M

$\leq 3B$

NPM (Min et al., 2023)

350M

1B

Atlas (Izacard et al., 2022)

11B

$\sim 30B$

RETRO (Borgeaud et al., 2021)

7B

2T

REPLUG (Shi et al., 2023)

$\leq 175B$

$\sim 5B$

Reliable RAG in the wild: effective datastore (quality)

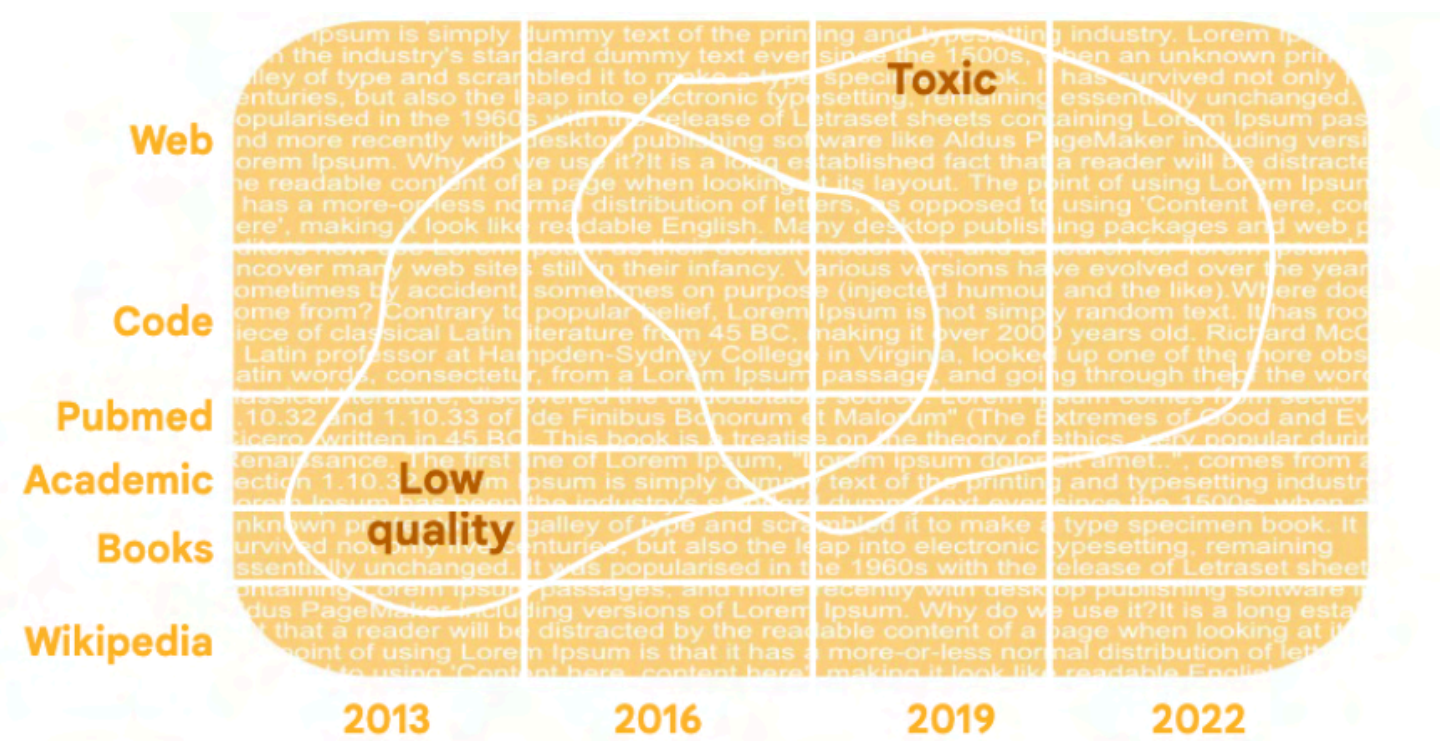
Efficiency

Effective datastore

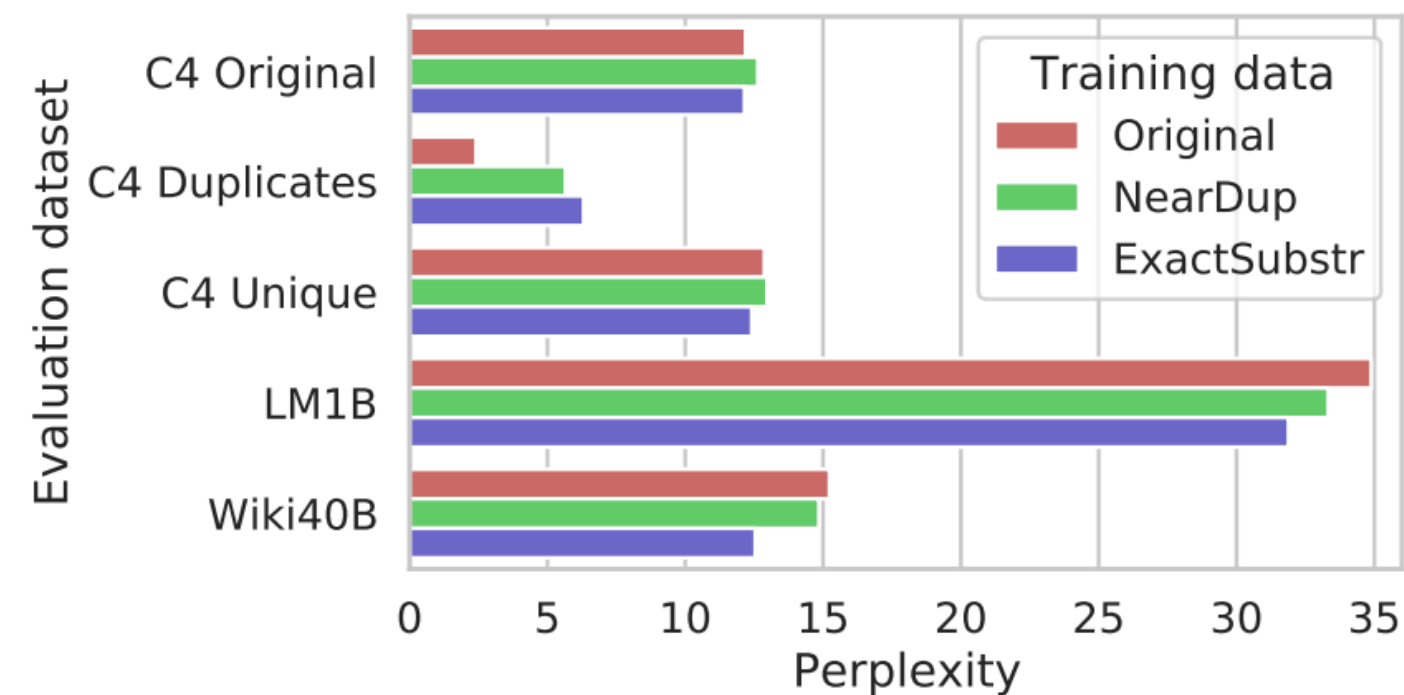
New application

Data-centric approaches to build effective large-scale datastore?

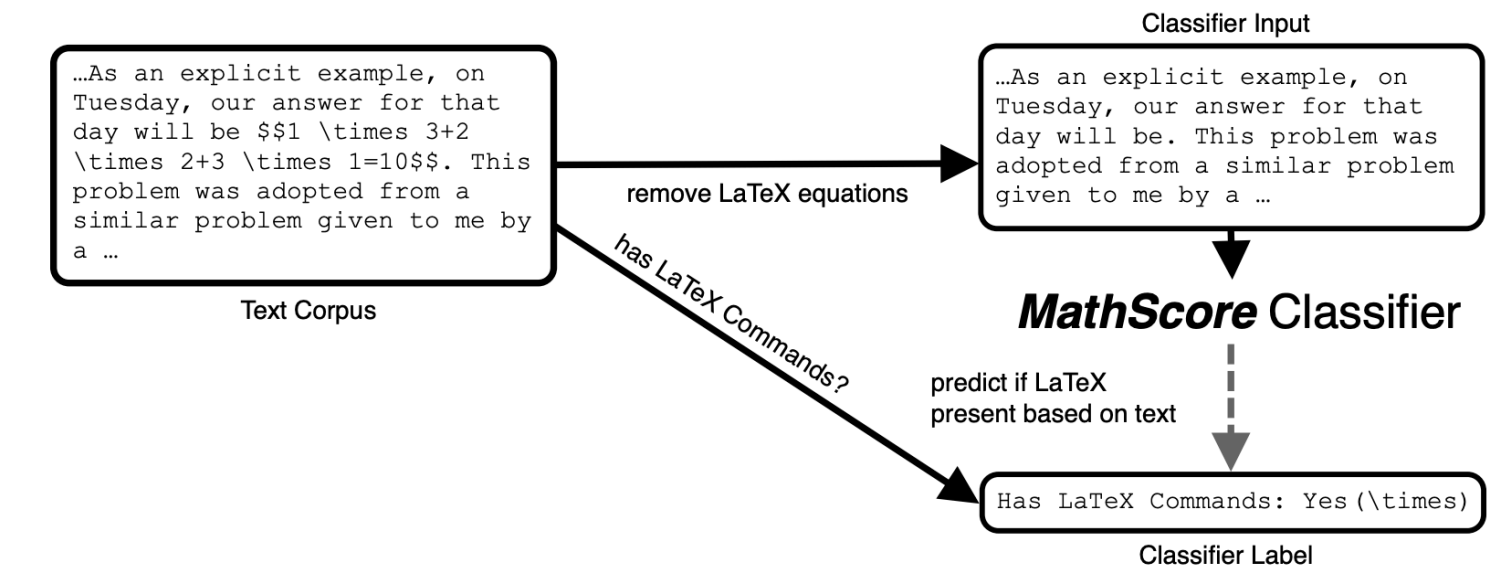
Quality & Composition (Longpre et al., 2023)



Deduplication (Lee et al., 2023)



Data Filtering (Paster et al., 2023)



Reliable RAG in the wild: scaling datastore

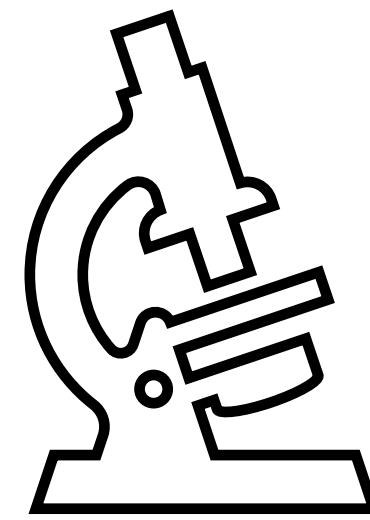
Efficiency

Scaling datastores

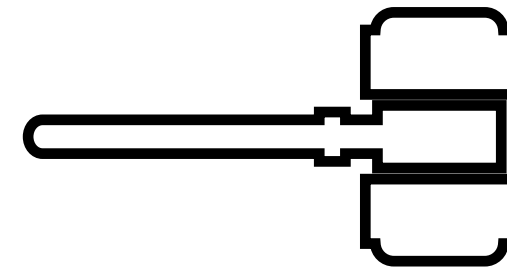
New application



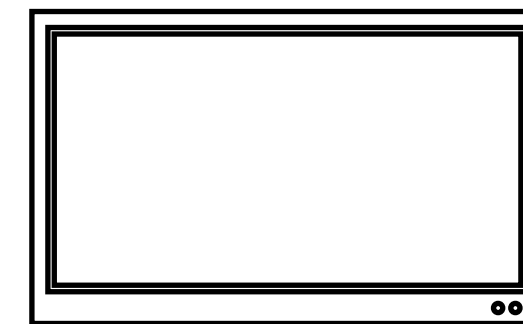
Medical



Science



Legal



Code



Finance

Reliable RAG in the wild: scaling datastore

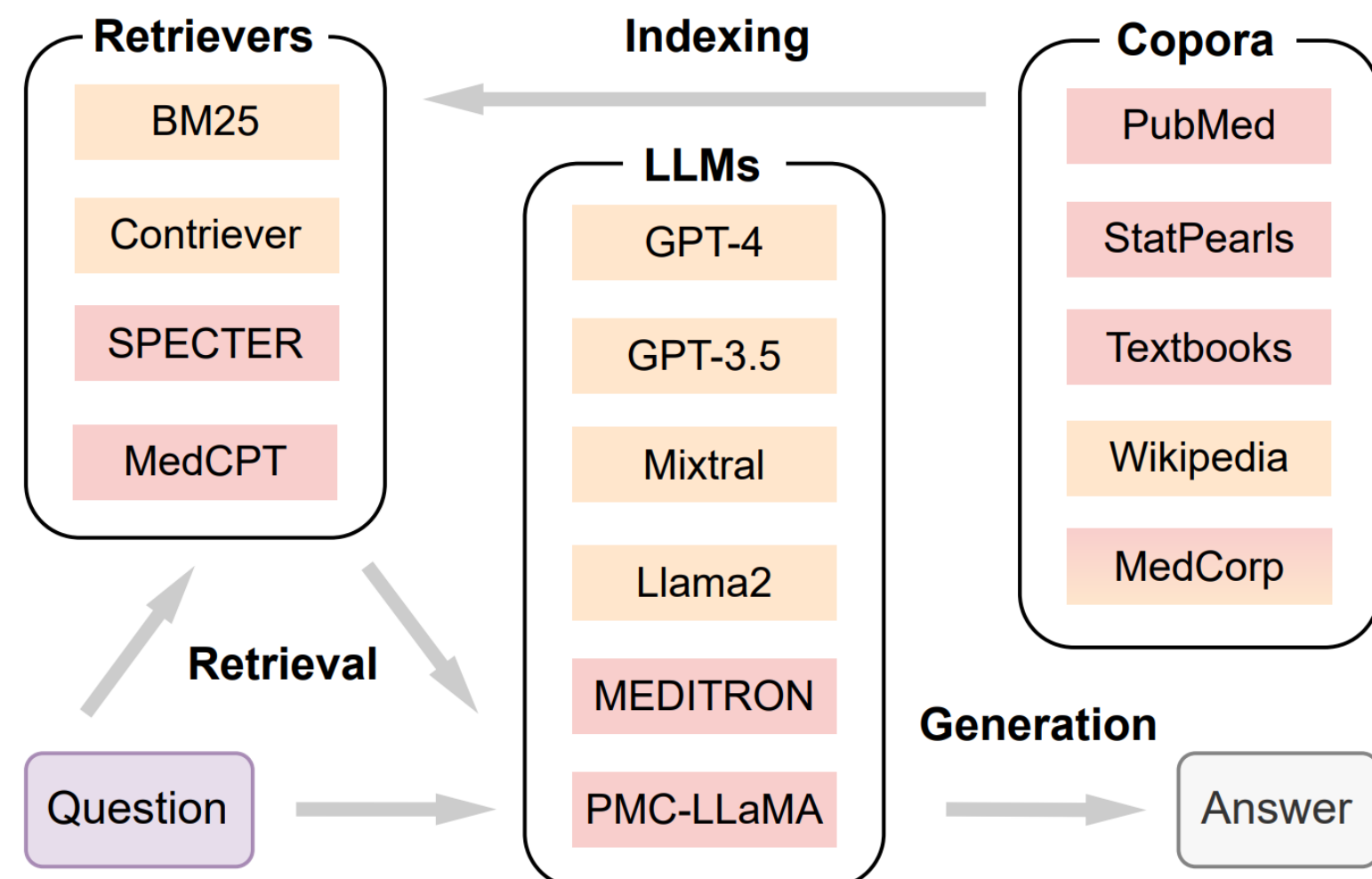
Efficiency

Scaling datastores

New application



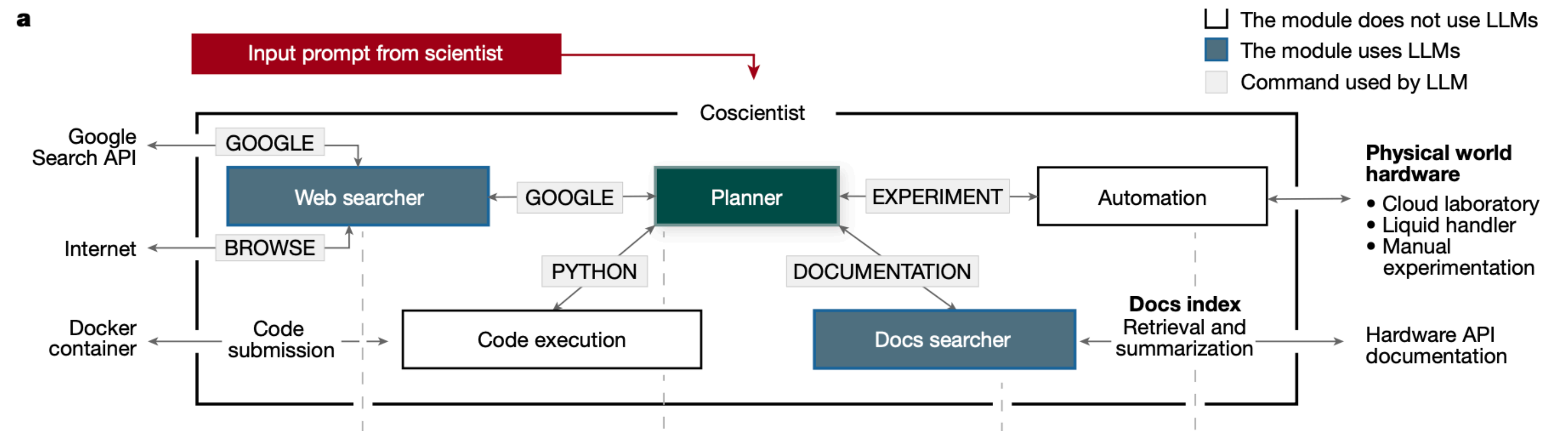
MedRAG (Xiong et al., 2024)



Realistic evaluation?



Autonomous chemistry experiments (Boiko et al., 2024)



RAG pipelines optimized for expert tasks?

Thanks for listening :)

- **Understanding RAG** (Asai et al., 2024b; Mallen*, Asai* et al., 2023)
 - Retrieval-augmented LMs can alleviate many issues in parametric LMs.
 - More fundamental improvements for architectures or training is necessary
- **Advancing RAG** (Asai et al., 2024a; Asai et al., 2023)
 - **Self-RAG** to build versatile retrieval-augmented LMs addressing issues in RAG
 - **Task-aware retrievals** to build versatile RAG systems

ACL 2023 tutorial: <https://acl2023-retrieval-lm.github.io/> by Akari, Sewon, Zexuan and Danqi
RAG survey: Retrieval-augmented Generation for Large Language Models: A Survey (Gao et al., 2024)

Contact: akari@cs.washington.edu
Website: <https://akariasai.github.io/>
Twitter: @AkariAsai
Public OH: Friday 6pm