

# Generative IR

@ ESSIR 2024

**Guido Zuccon**

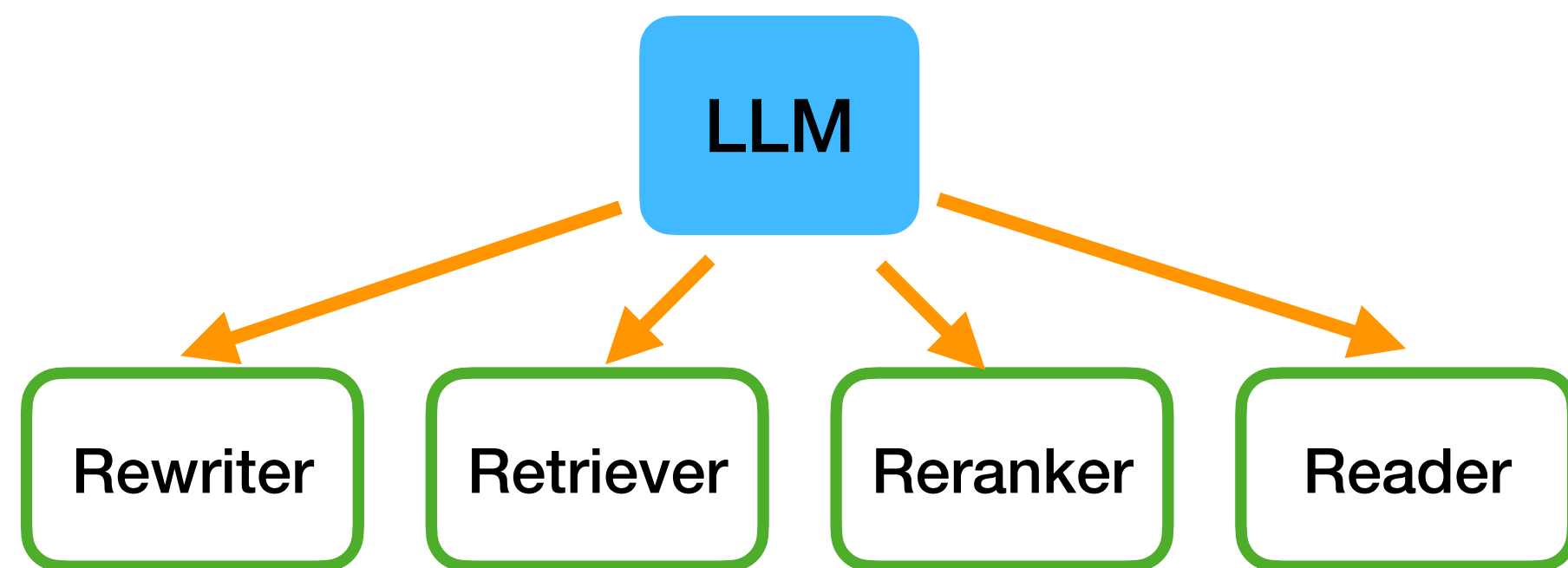
ielab, The University of Queensland, Australia

# Focus of This Lecture

- Generative IR: The use of sequence-to-sequence language models to perform IR tasks
  - Encoder-decoder LLMs (e.g., T5) and decoder-only LLMs (e.g., GPT)
- In this lecture when we say LLMs we mean any of the two type of models above
  - i.e. we do not mean BERT-style models (encoder-only)

# Three Directions in Generative IR

## Prompting LLMs for Retrieval and Ranking

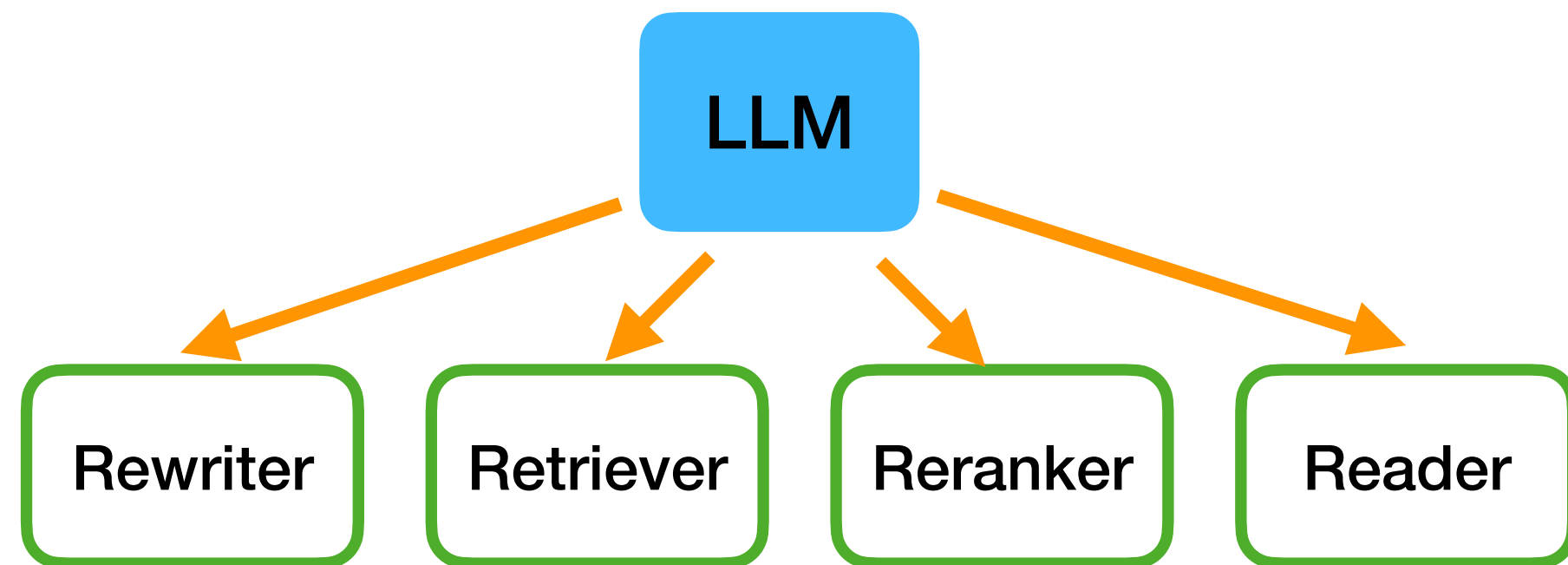


*LLM-enhanced information retrieval*

e.g. RankGPT, PromptReps

# Three Directions in Generative IR

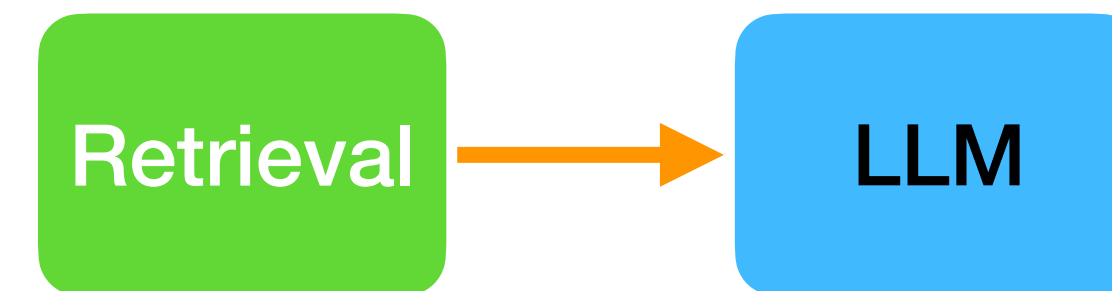
## Prompting LLMs for Retrieval and Ranking



*LLM-enhanced information retrieval*

e.g. RankGPT, PromptReps

## Retrieval Augmented Generation

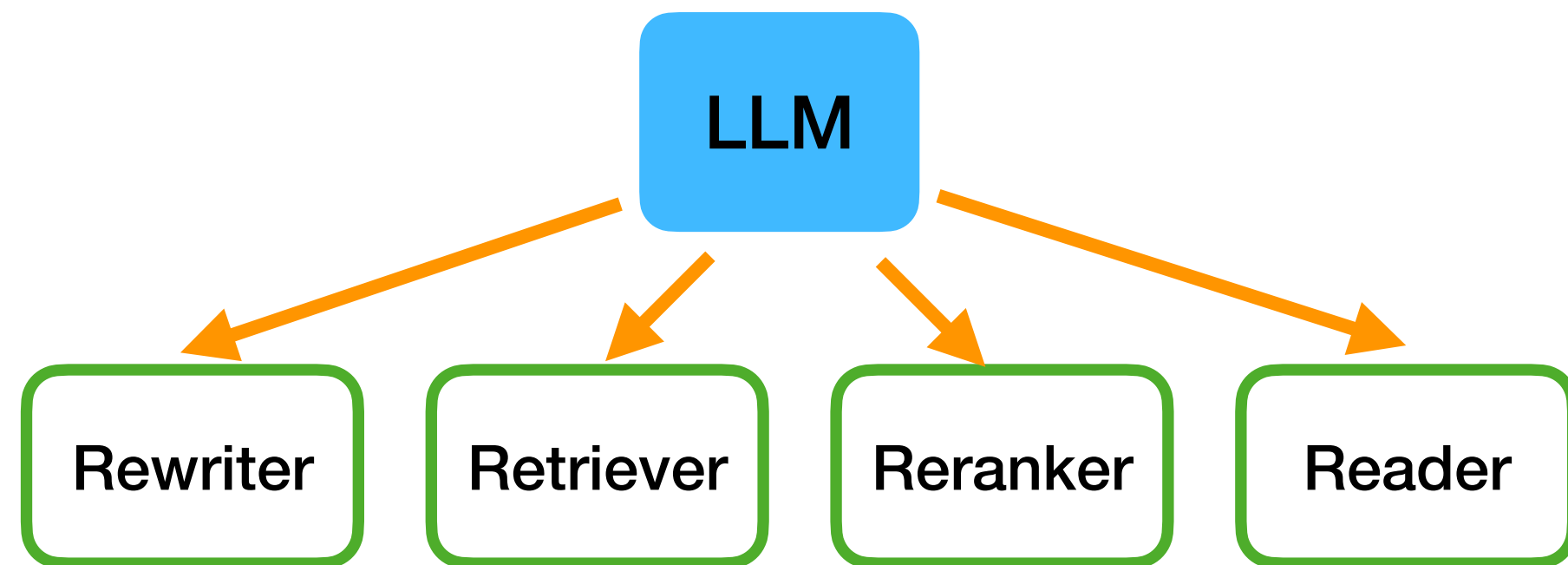


*IR4LLM,  
Retrieval-augmented LLM,  
Reliable response generation*

e.g. Self-RAG, BlendFilter, etc.

# Three Directions in Generative IR

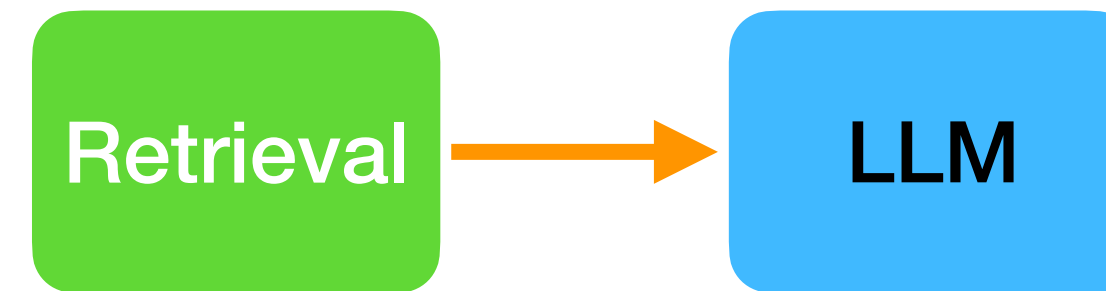
## Prompting LLMs for Retrieval and Ranking



*LLM-enhanced information retrieval*

e.g. RankGPT, PromptReps

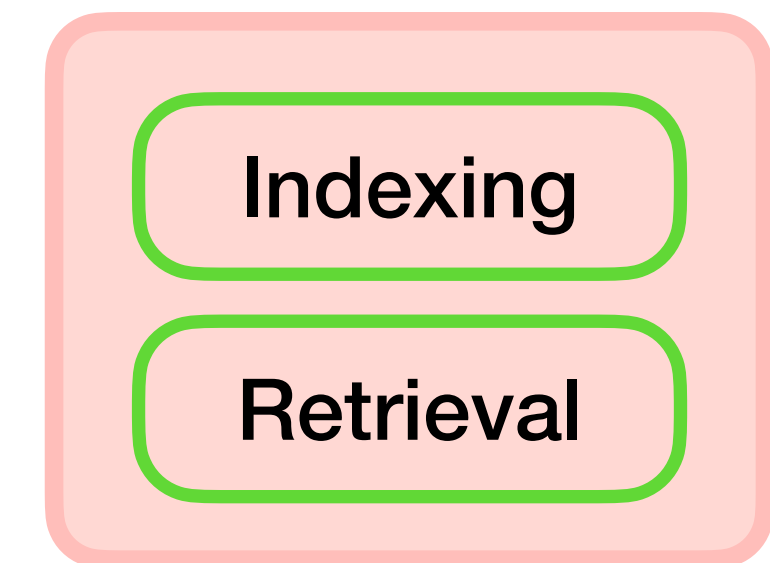
## Retrieval Augmented Generation



*IR4LLM,  
Retrieval-augmented LLM,  
Reliable response generation*

e.g. Self-RAG, BlendFilter, etc.

## Model-based IR



*LLM as Indexer and  
Retriever, Generative  
[document] Retrieval (GR)*

e.g. Differentiable Search Index (DSI)

# Disclaimer

- Would **not have time to cover all the slides**: slides are made publicly available for offline consultation + pointers to more readings
- Most of the time can **only provide a few examples** of methods, **not a comprehensive** overview of all methods:
  1. sometimes **preference to more recent work** (i.e. arXiv, not peer-reviewed)-> **not necessarily perfect work**, but examples of very fresh directions (most recent on arXiv on July 1st, 2024)
  2. Other times to more **fundamental** work
- Some methods are talked about in a **simplified/restricted/high-level** manner to give you the gist — always read the original paper to get more if interested
  - Please pardon my lingo sloppiness from time to time

# Three Directions in Generative IR

Will spend the least amount of time  
on DSI

TUTORIAL | OPEN ACCESS

## Recent Advances in Generative Information Retrieval

Authors:  [Yubao Tang](#),  [Ruqing Zhang](#),  [Weiwei Sun](#),  [Jiafeng Guo](#), and  [Maarten De Rijke](#)

WWW '24: Companion Proceedings of the ACM on Web Conference 2024 • May 2024 • Pages 1238 - 1241  
<https://doi.org/10.1145/3589335.3641239>

Published: 13 May 2024 [Publication History](#)



also coming to SIGIR 2024!

## Model-based IR

Indexing

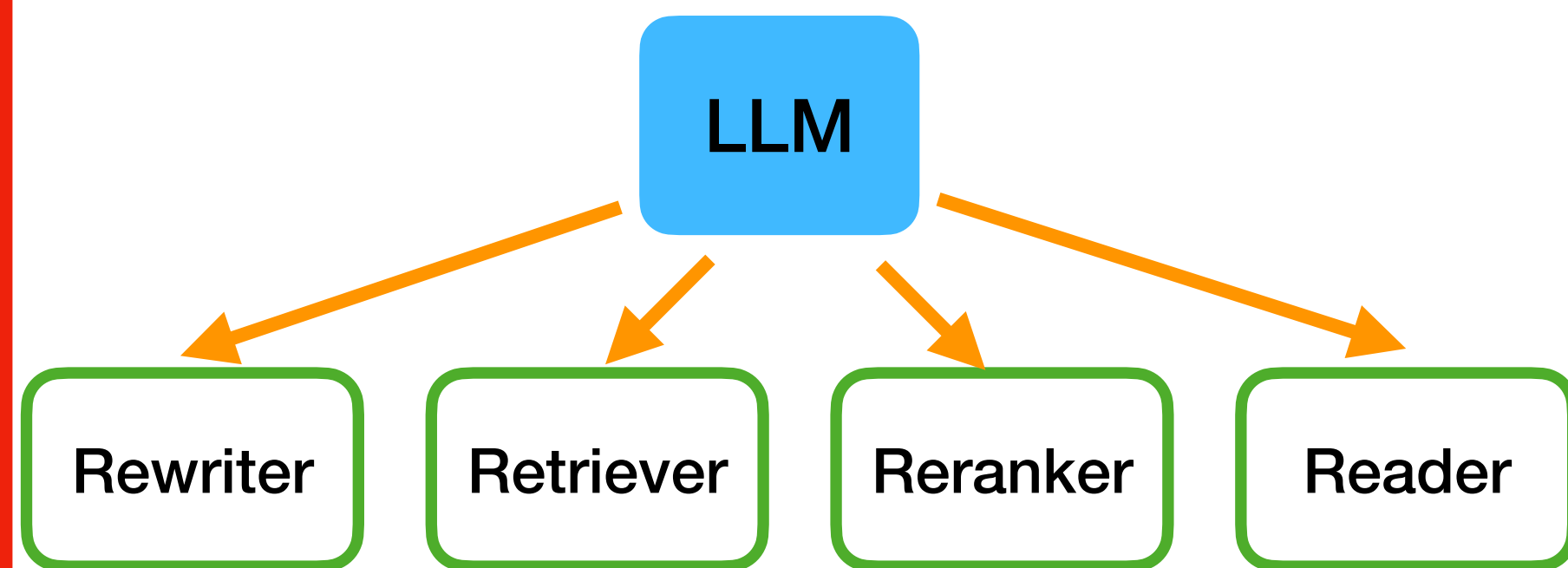
Retrieval

*LLM as Indexer and Retriever, Generative [document] Retrieval (GR)*

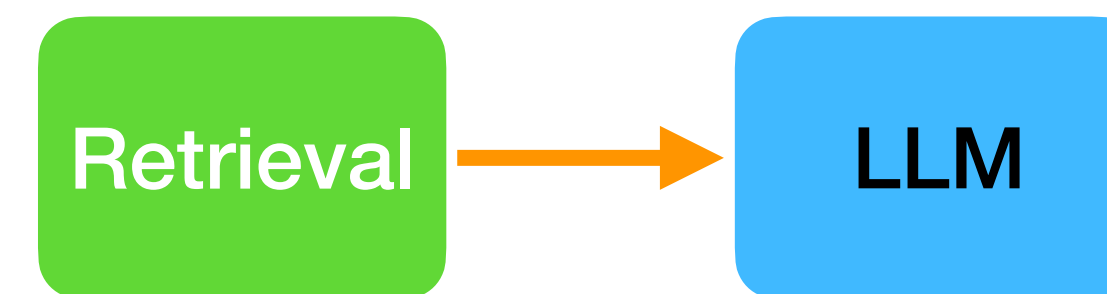
e.g. Differentiable Search Index (DSI)

# Three Directions in Generative IR

## Prompting LLMs for Retrieval and Ranking



## Retrieval Augmented Generation



In the live presentation I incorrectly attributed a great tutorial on LLMs4IR from the Chinese IR community to BAAI — the correct attribution should be made to CFF. The event was part of the CCF Advanced Disciplines Lectures series. Correct link: <https://ccf.org.cn/ADL147>; videos are not yet online but are supposed to be released at <https://dl.ccf.org.cn/>

Some great tutorial resources on other aspects of Generative IR too



<https://2024.baai.ac.cn/>



THE UNIVERSITY  
OF QUEENSLAND  
AUSTRALIA

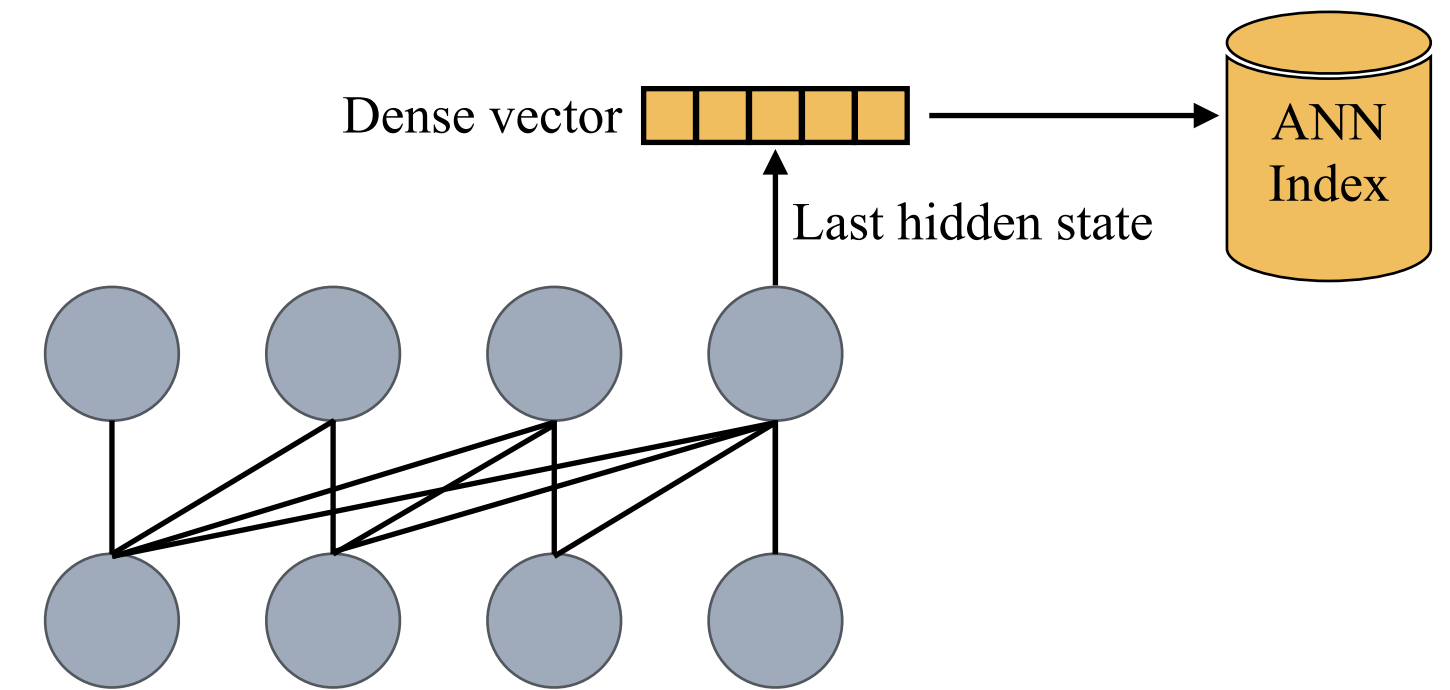
CREATE CHANGE

# Part 1: Prompting LLMs for Retrieval and Ranking

# Key Intuitions

Can we ask LLMs to:

- **Retrieve:** generate representations for a document/query?
- Then we can use representations to do matching, e.g. create embeddings and use them for dense retrieval

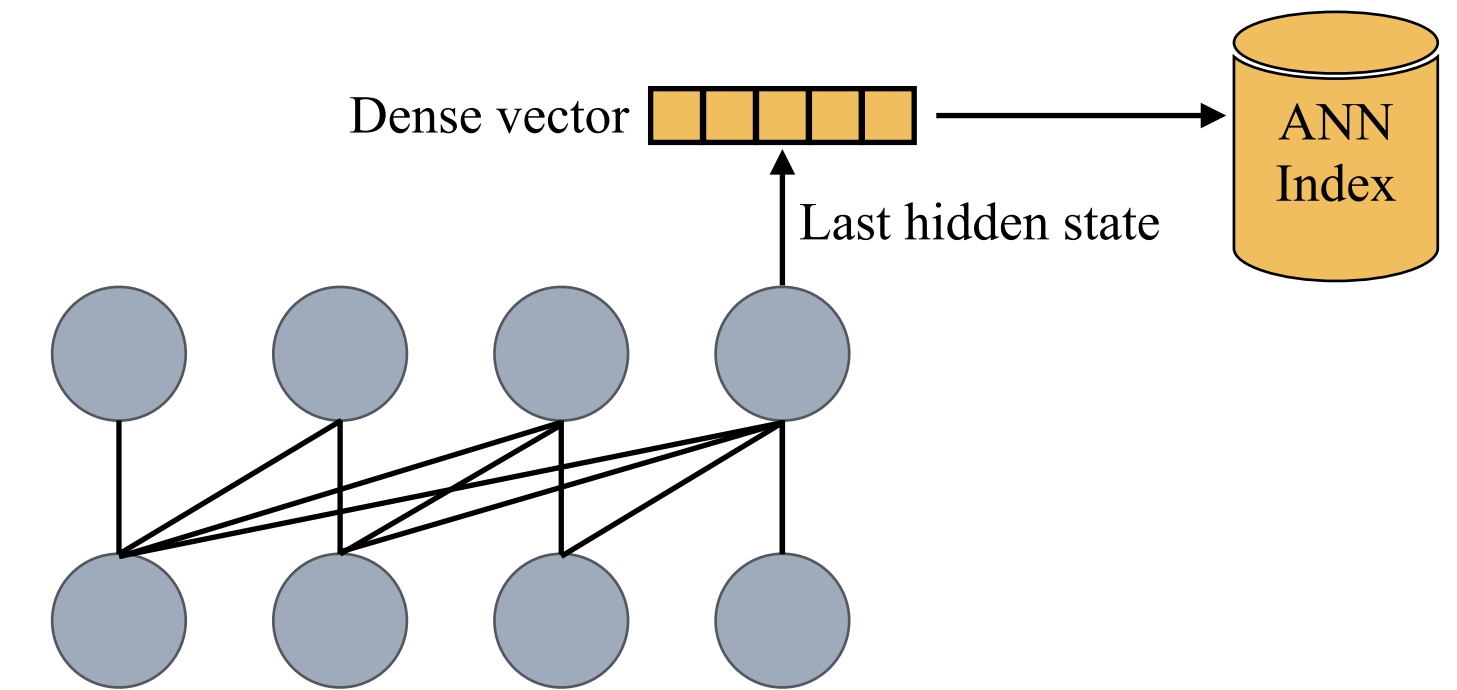


**<System>** You are an AI assistant that can understand human language.  
**<User>** Passage: “[text]”. Use one word to represent the passage in a retrieval task. Make sure your word is in lowercase.  
**<Assistant>** The word is: “

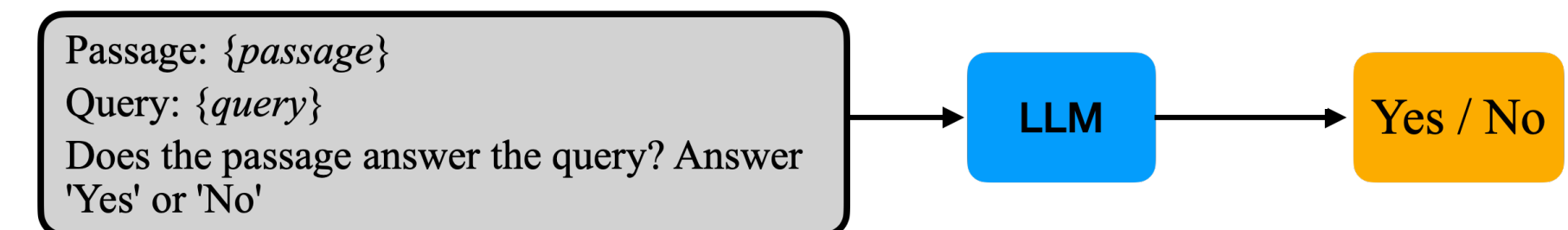
# Key Intuitions

Can we ask LLMs to:

- **Retrieve:** generate representations for a document/query?
  - Then we can use representations to do matching, e.g. create embeddings and use them for dense retrieval
- **Ranking:** tell us the relevance of a document to a query?
  - Then we can use this indication of relevance (or relative relevance of n documents) to rank documents for the query



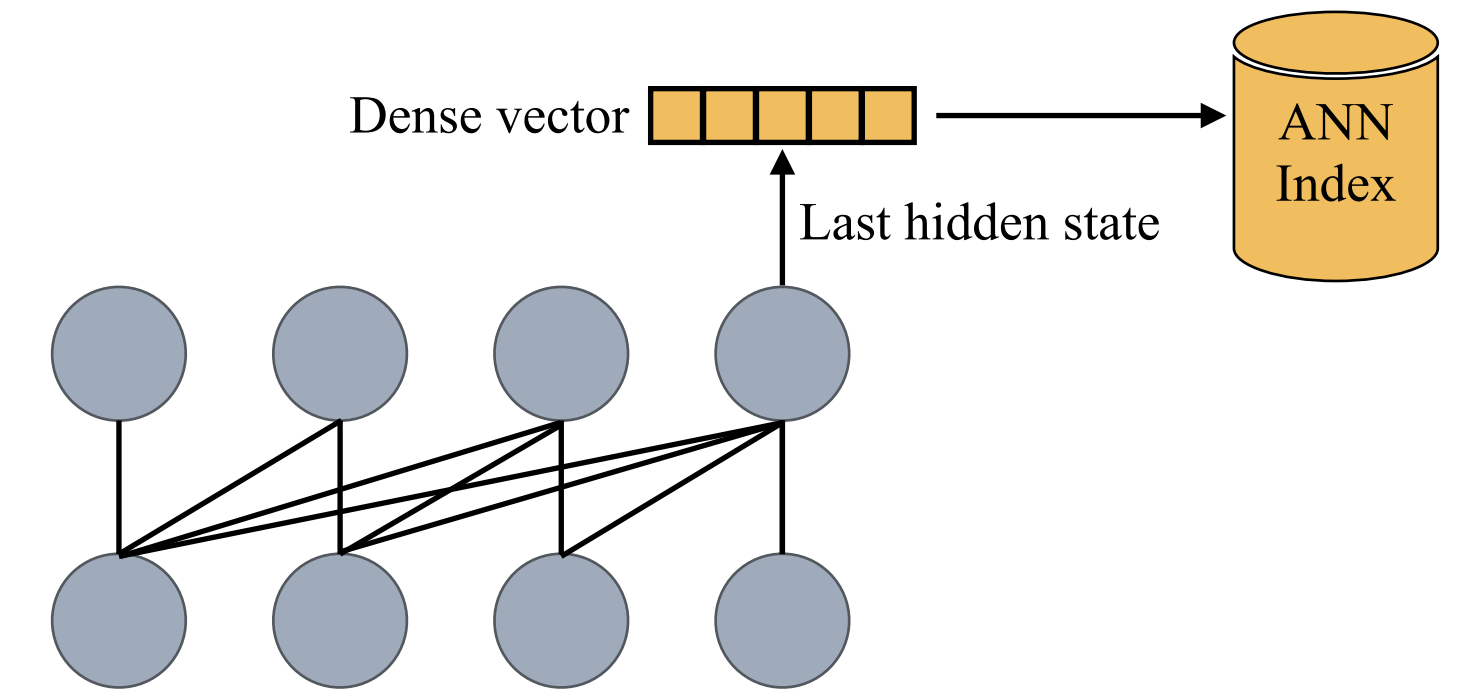
**<System>** You are an AI assistant that can understand human language.  
**<User>** Passage: “[text]”. Use one word to represent the passage in a retrieval task. Make sure your word is in lowercase.  
**<Assistant>** The word is: “



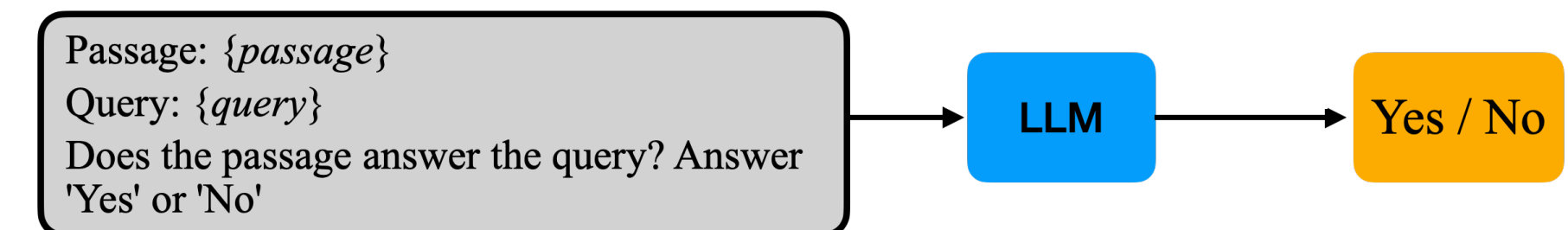
# Key Intuitions

Can we ask LLMs to:

- **Retrieve:** generate representations for a document/query?
  - Then we can use representations to do matching, e.g. create embeddings and use them for dense retrieval
- **Ranking:** tell us the relevance of a document to a query?
  - Then we can use this indication of relevance (or relative relevance of n documents) to rank documents for the query



**<System>** You are an AI assistant that can understand human language.  
**<User>** Passage: “[text]”. Use one word to represent the passage in a retrieval task. Make sure your word is in lowercase.  
**<Assistant>** The word is: “

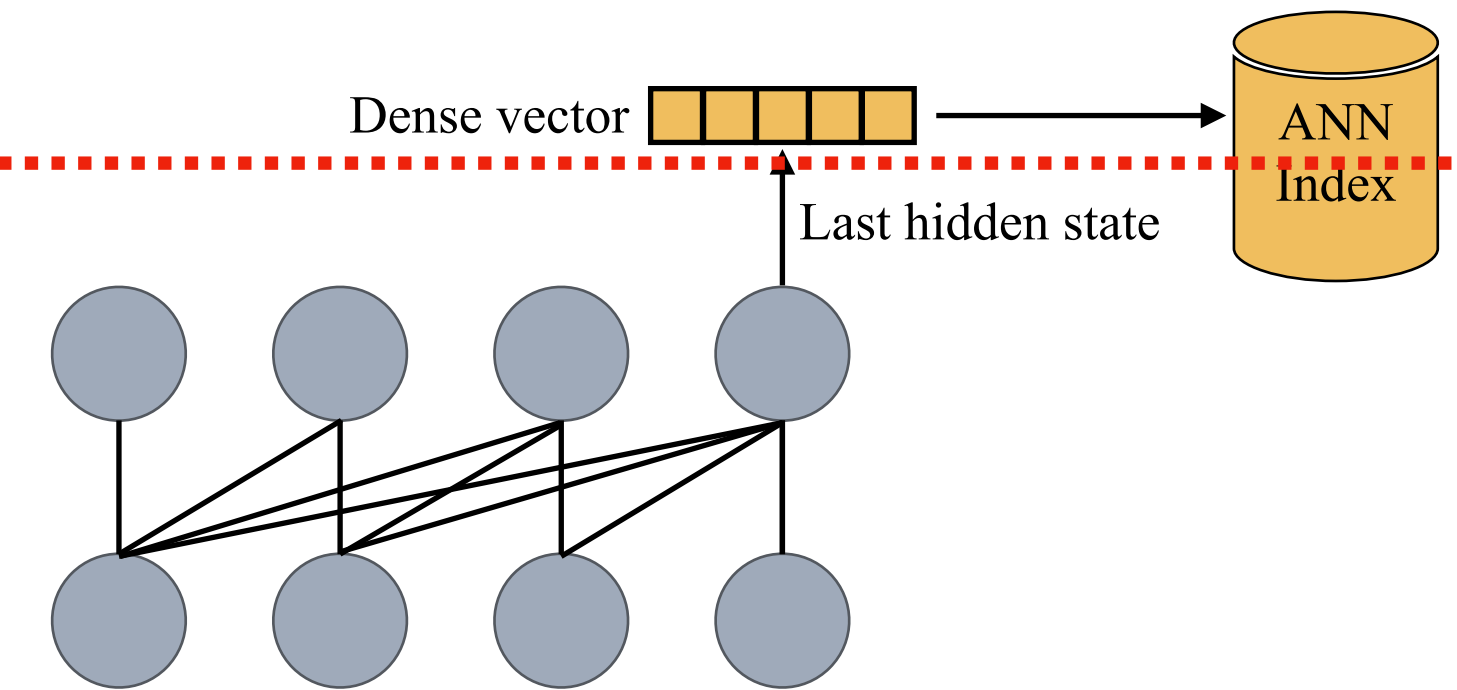


**Idea: devise prompts/instructions to tell the LLM how to perform these tasks effectively**

# Key Intuitions

Can we ask LLMs to:

- **Retrieve:** generate representations for a document/query?
  - Then we can use representations to do matching, e.g. create embeddings and use them for dense retrieval



**<System>** You are an AI assistant that can understand human language.  
**<User>** Passage: “[text]”. Use one word to represent the passage in a retrieval task. Make sure your word is in lowercase.  
**<Assistant>** The word is: “

• **Ranking:** tell

- Then we relevance query

**Plan:**

**1. LLMs for dense bi-encoding:**

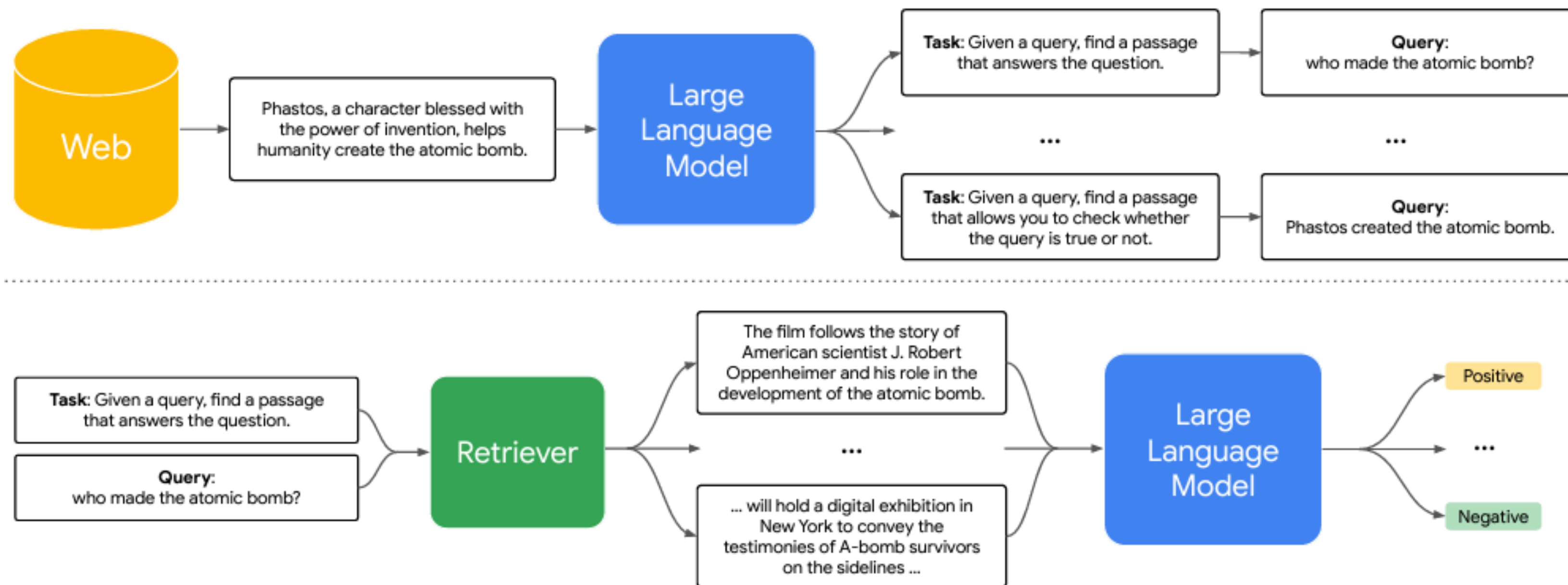
- (i) leverage LLM for data generation (e.g. Gecko, E5)
- (ii) as backbone (e.g. E5, LLM2Vec)

**2. LLMs as a zero-shot hybrid representation creation (PromptReps)**

Yes / No

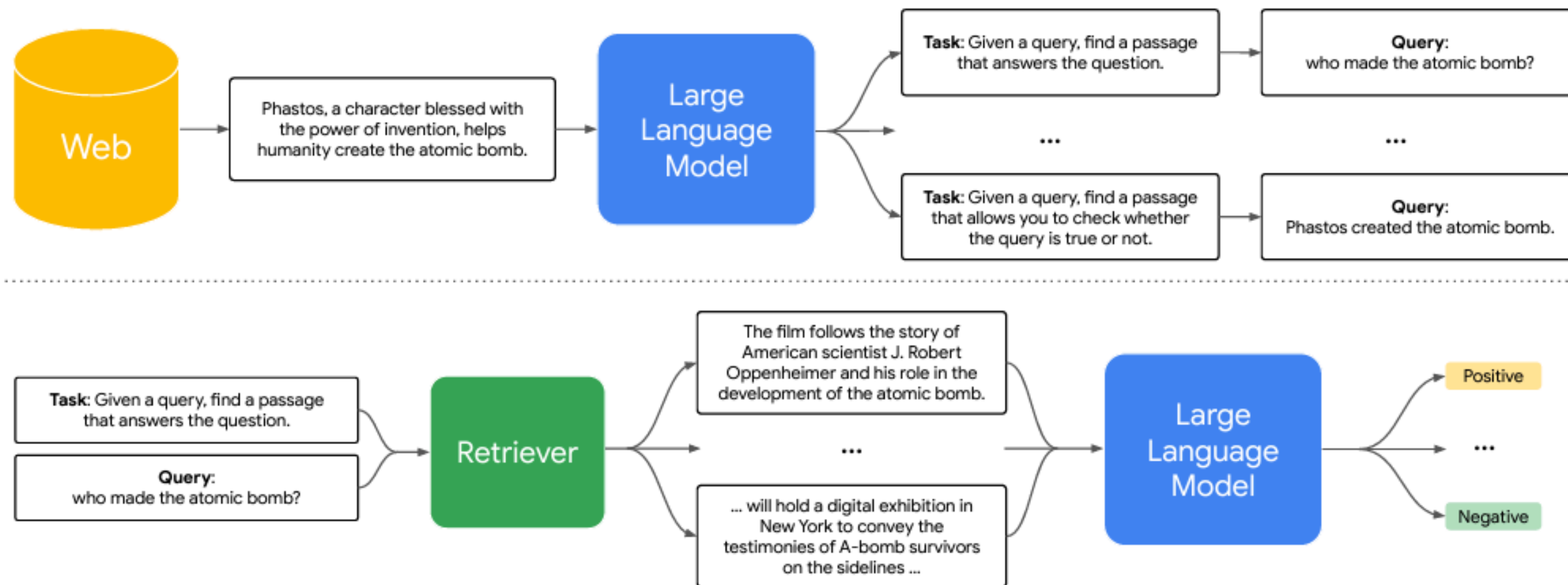
# Gecko embeddings:

- The underlying LLM architecture is not specified; likely uses Google's Matryoshka embeddings technique
- Training relies on LLMs
  - LLMs used to generate a Few-shot Prompted Retrieval dataset (FRet) for knowledge distillation from LLM into embedding through 2 tasks: (1) diverse query generation, (2) positive & negative mining



# Gecko embeddings:

- The underlying LLM architecture is not specified; likely uses Google’s Matryoshka embeddings technique
- Training relies on LLMs
  - LLMs used to generate a Few-shot Prompted Retrieval dataset (FRet) for knowledge distillation from LLM into embedding through 2 tasks: (1) diverse query generation, (2) positive & negative mining



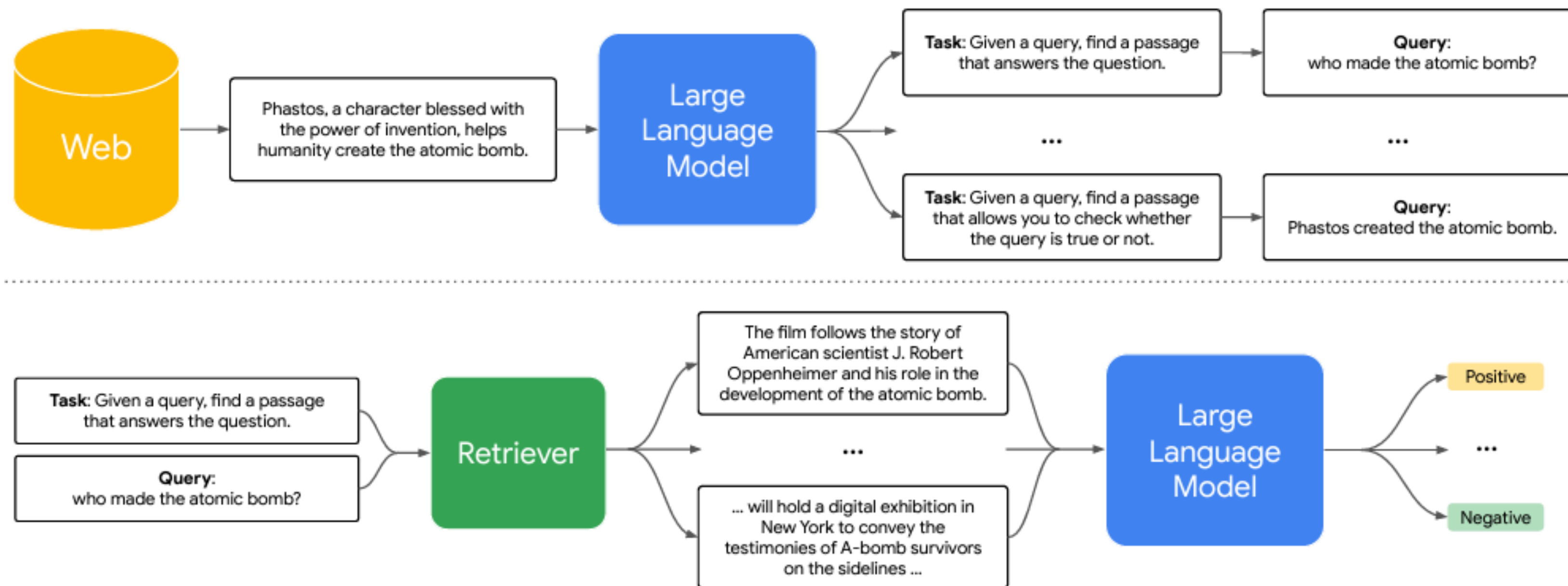
## (1) Query/task generation

Prompt for task generation:

- “Given a query, find a passage that has the answer to the query” [question answering]
- “Given a query, find a passage that allows you to check whether the query is true or not” [fact checking]

# Gecko embeddings:

- The underlying LLM architecture is not specified; likely uses Google's Matryoshka embeddings technique
- Training relies on LLMs
  - LLMs used to generate a Few-shot Prompted Retrieval dataset (FRet) for knowledge distillation from LLM into embedding through 2 tasks: (1) diverse query generation, (2) positive & negative mining



## (2) Pos/Neg Mining

Observation: generated queries focus on particular aspect of original passage

Sampling Method:

1. Use initial model trained with  $(q, p_{seed})$  pairs to retrieve top passages
2. Use an LLM to rank (2 tried: query likelihood (Sachin's UPR) and relevance classification (H. Zhuang's RG) — see later)
3. Reciprocal Rank Fusion of 1&2
4. Top ranked:  $p_+$ ; other tops or bottom  $p_-$

# E5 with LLMs

- Original E5 embeddings are BERT-based; new E5 embeddings are decoder-only LLM (Mistral-7b)
- E5-Mistral uses LLMs in two ways:

(A) **Synthetic data generation:** (1) prompt LLM (GPT4) to “brainstorm” a list of potential retrieval tasks, (2) generate <query, positive, hard negative> for each task.

You have been assigned a retrieval task: *{task}*  
Your mission is to write one text retrieval example for this task in JSON format. The JSON object must contain the following keys:


- **"user\_query"**: a string, a random user search query specified by the retrieval task.
- **"positive\_document"**: a string, a relevant document for the user query.
- **"hard\_negative\_document"**: a string, a hard negative document that only appears relevant to the query.

Please adhere to the following guidelines:

- The "user\_query" should be *{query\_type}*, *{query\_length}*, *{clarity}*, and diverse in topic.
- All documents should be at least *{num\_words}* words long.
- Both the query and documents should be in *{language}*.

... (omitted some for space)

Your output must always be a JSON object only, do not explain yourself or output anything else. Be creative!

 {"user\_query": "How to use Microsoft Power BI for data analysis",  
"positive\_document": "Microsoft Power BI is a sophisticated tool that requires time and practice to master. In this tutorial, we'll show you how to navigate Power BI ... (omitted) ",  
"hard\_negative\_document": "Excel is an incredibly powerful tool for managing and analyzing large amounts of data. Our tutorial series focuses on how you...(omitted)" }

# E5 with LLMs

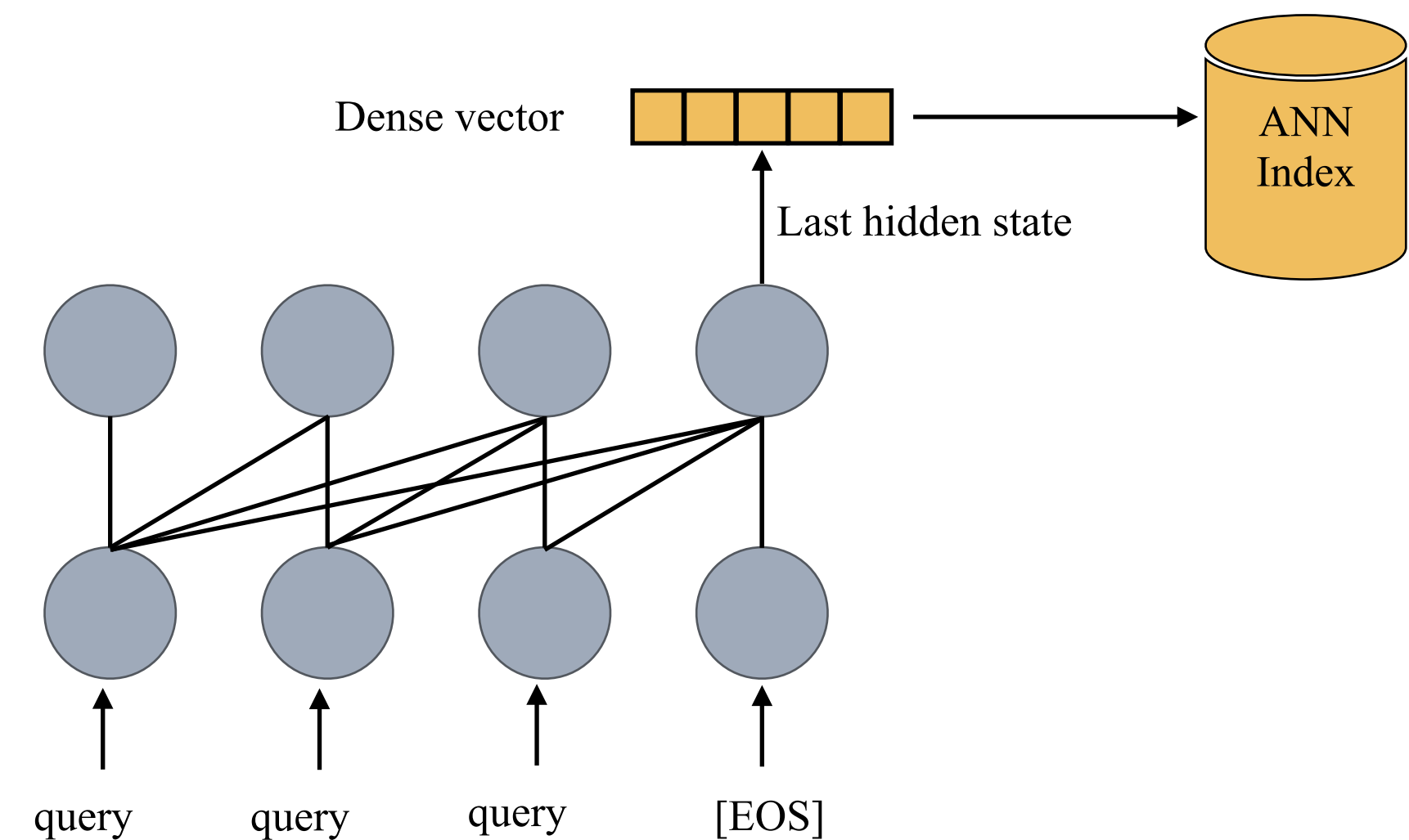
- Original E5 embeddings are BERT-based; new E5 embeddings are decoder-only LLM (Mistral-7b)
- E5-Mistral uses LLMs in two ways:

(A) **Synthetic data generation:** (1) prompt LLM (GPT4) to “brainstorm” a list of potential retrieval tasks, (2) generate <query, positive, hard negative> for each task.

(B) **Embedding backbone:** (1) modify query to:

q+ inst = Instruct: {task\_definition} \n Query: {q+}

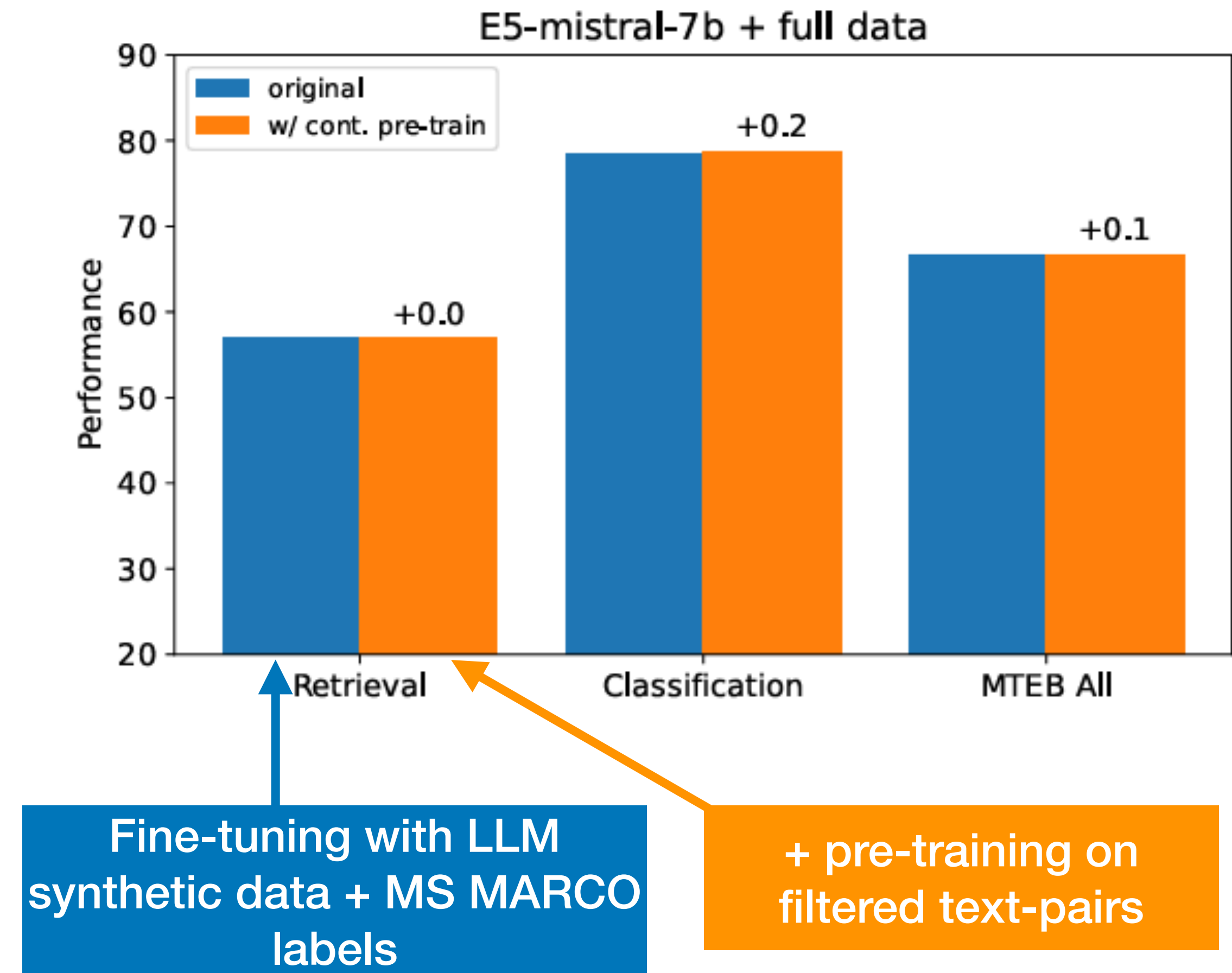
(2) append [EOS] token to end of query and document; feed them into LLM to get embeddings by taking last layer of [EOS] vector



- Embedding model then trained with standard contrastive loss over the in-batch negatives and hard negatives

# E5 with LLMs

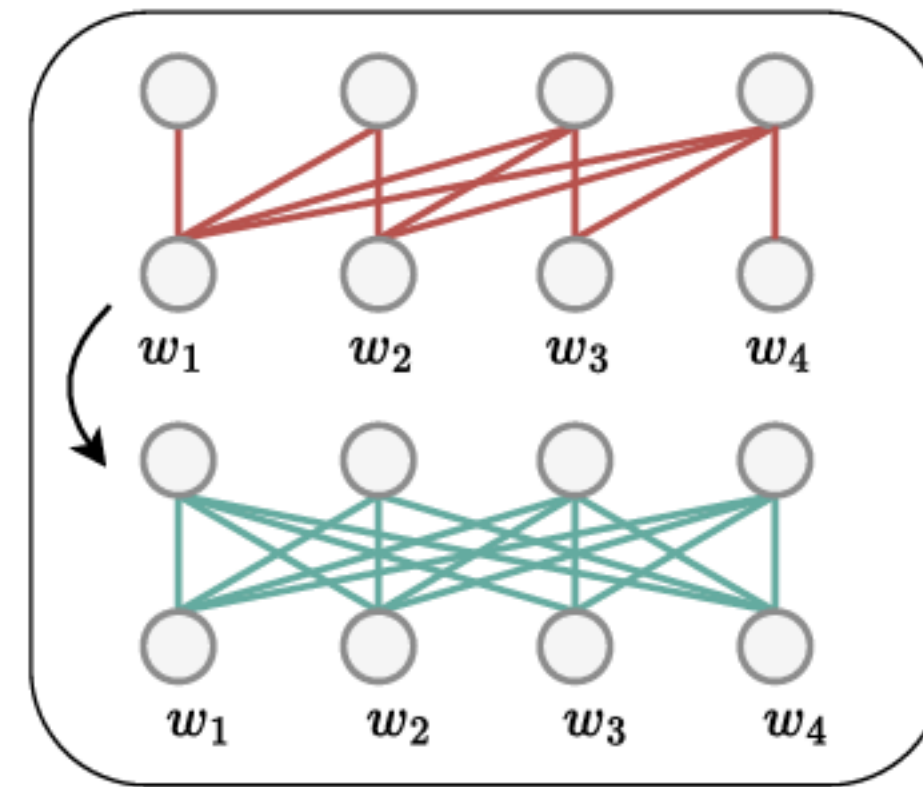
- Training in E5Mistral:
  - Contrastive pre-training: weakly supervised data from filtered text pairs (from E5 paper)
  - Contrastive fine-tuning: LLM generated synthetic data + MS MARCO labelled data [they also have version + 18 other datasets]
- Observation: contrastive *pre*-training has negligible impact on model quality.
- extensive auto-regressive pre-training enables LLMs to acquire good text representations, and only fine-tuning required to obtain effective embeddings.



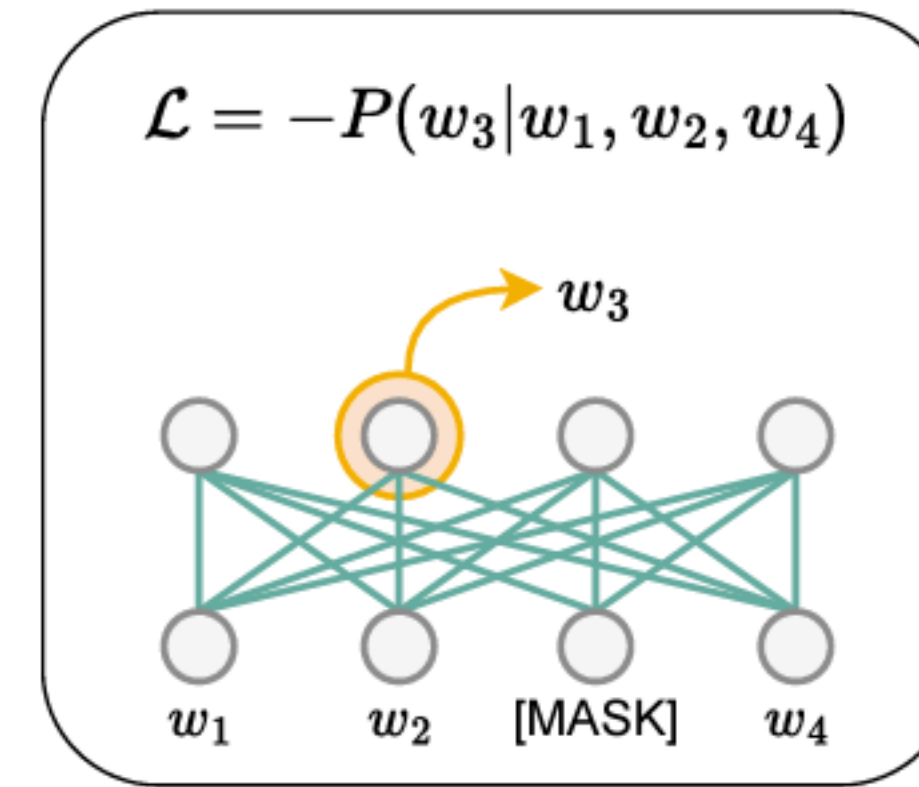
# LLM2Vec

- Transform any decoder-only LLM into a strong text encoder
- Three key steps:
  1. enable bidirectional attention,
  2. masked next token prediction,
  3. unsupervised contrastive learning
- (4) can add supervised contrastive learning

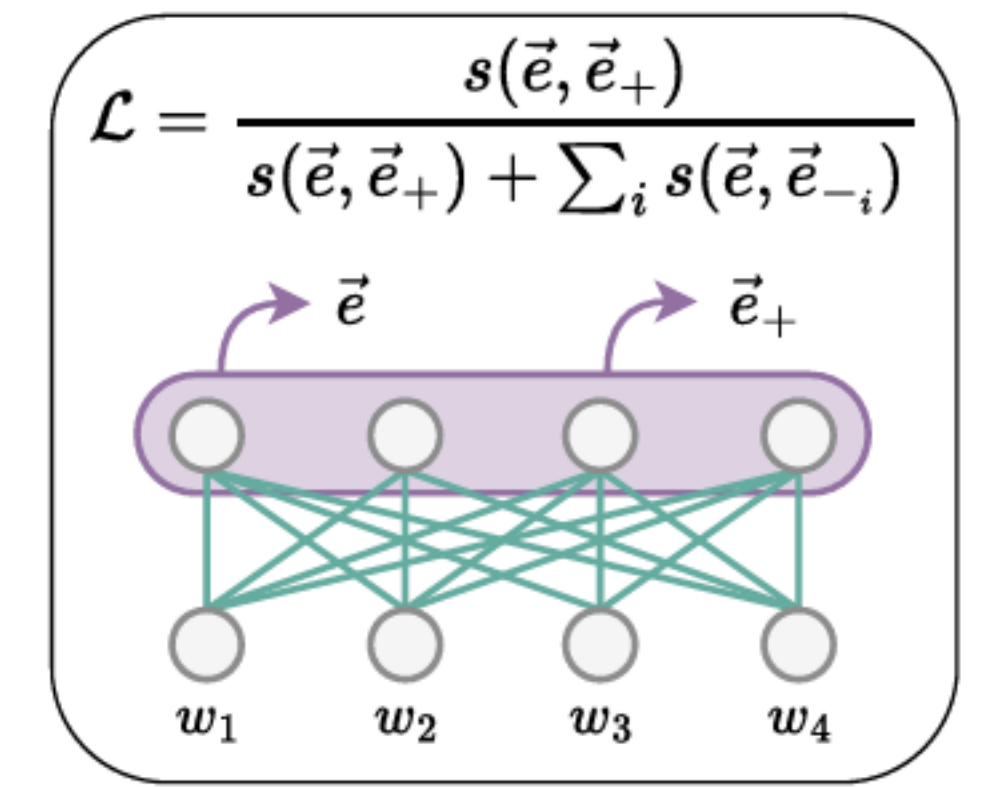
Enabling Bidirectional Attention



Masked Next Token Prediction



Unsupervised Contrastive Learning



Modify the self-attention mechanism in the decoder-only architecture by replacing the causal attention mask by an all-ones matrix

gives each token access to every other token in the sequence, converting it into a bidirectional LLM

simply enabling bidirectional attention decreases embedding performance: decoder-only LLM was not trained to attend to future tokens

BehnamGhader, P., Adlakha, V., Mosbach, M., Bahdanau, D., Chapados, N. and Reddy, S., 2024. Llm2vec: Large language models are secretly powerful text encoders. *arXiv preprint arXiv:2404.05961*.

# LLM2Vec

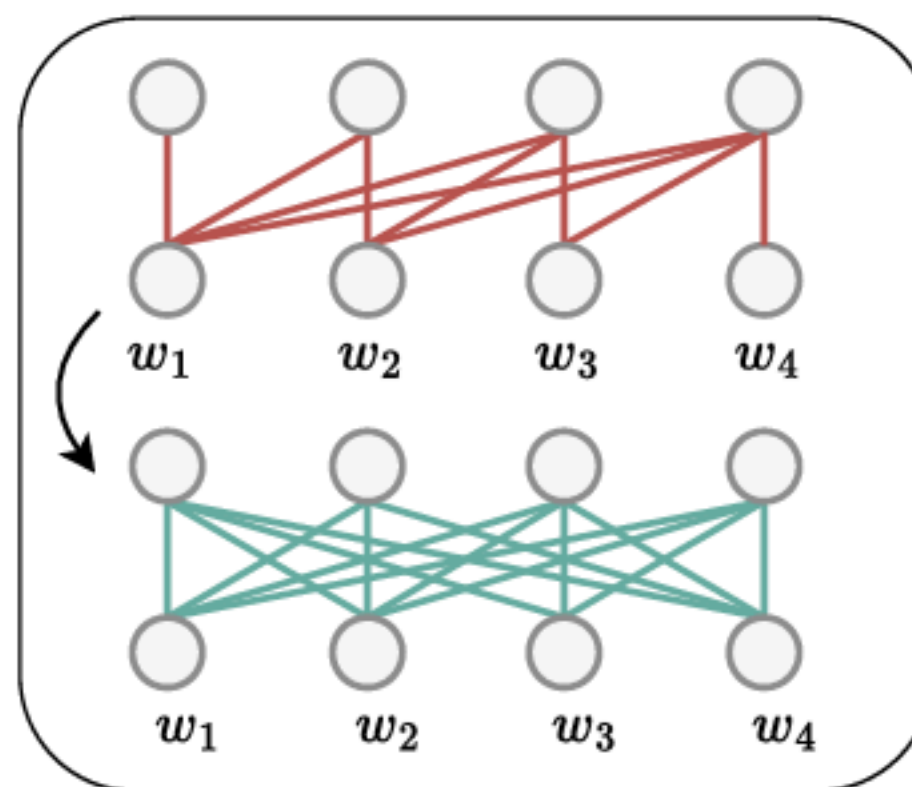
- Transform any decoder-only LLM into a strong text encoder

- Three key steps:

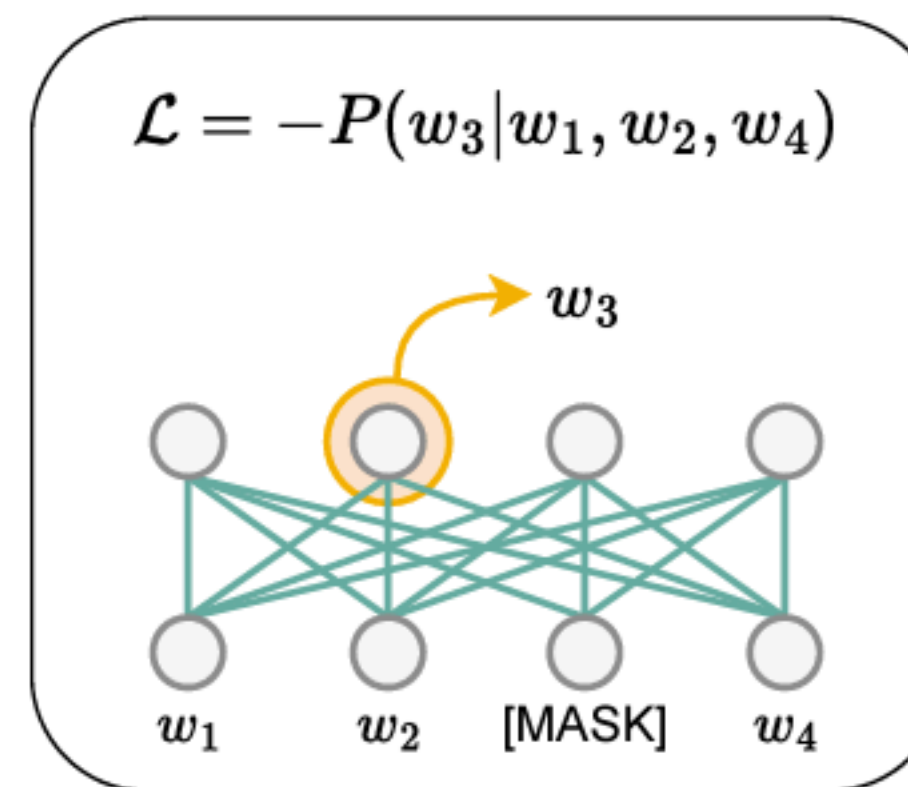
1. enable bidirectional attention,
2. masked next token prediction,
3. unsupervised contrastive learning

- (4) can add supervised contrastive learning

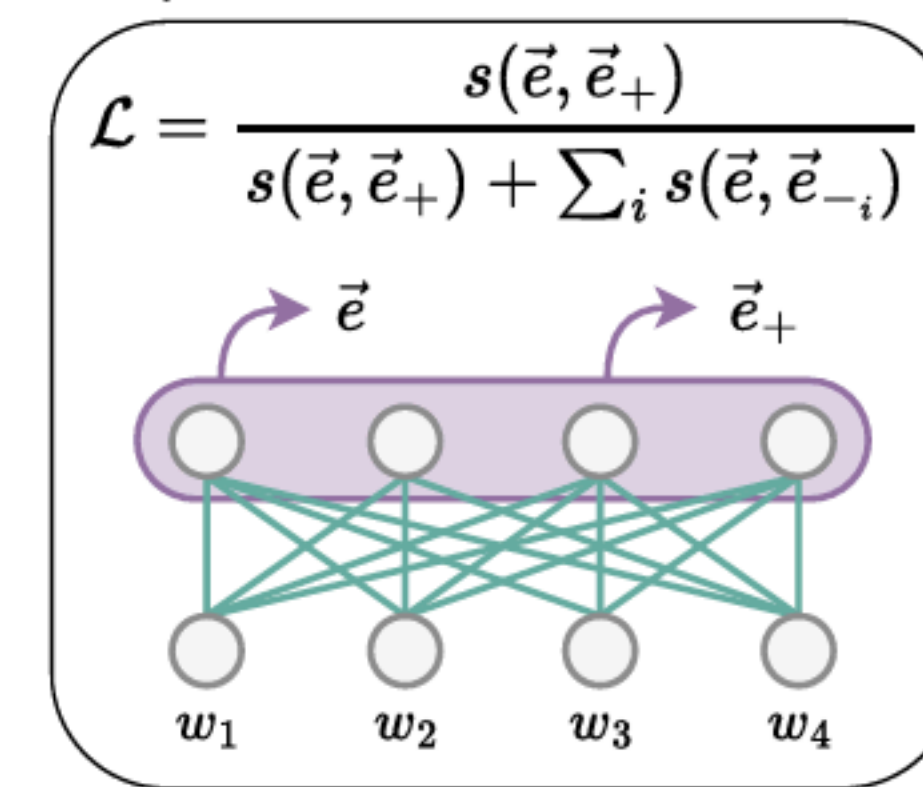
Enabling Bidirectional Attention



Masked Next Token Prediction



Unsupervised Contrastive Learning



make the model aware of its bidirectional attention by adapting it via masked next token prediction

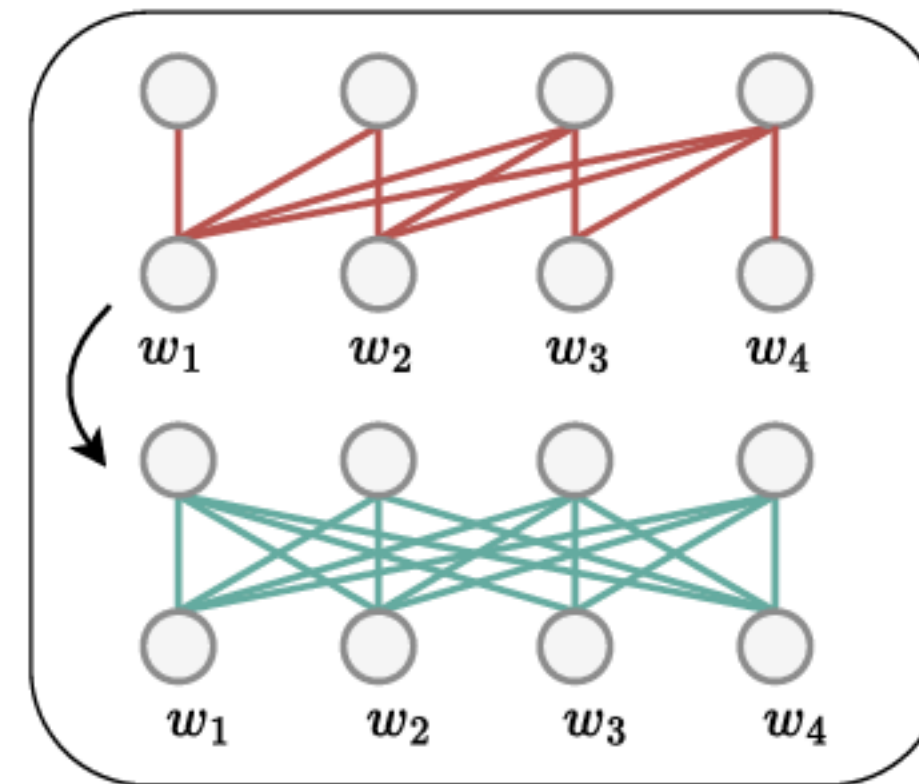
when predicting a masked token at position  $i$ , compute loss based on logits from token representation at position  $i - 1$ , not the masked position itself

No next sentence prediction objective in pre-training: not explicitly trained to capture the context of the entire sequence

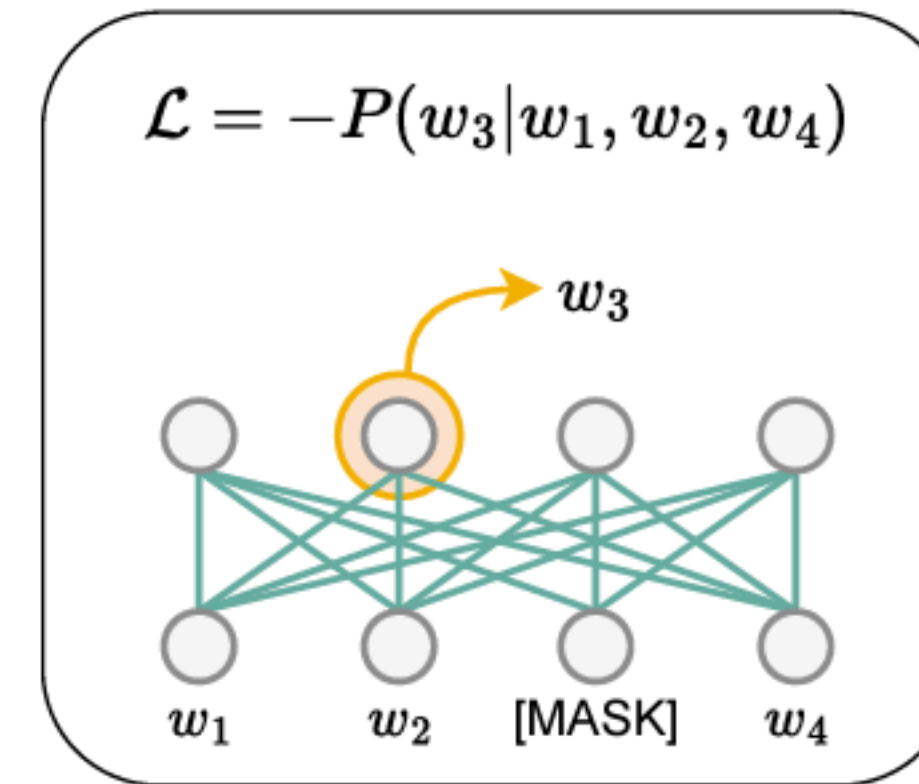
# LLM2Vec

- Transform any decoder-only LLM into a strong text encoder
- Three key steps:
  1. enable bidirectional attention,
  2. masked next token prediction,
  3. unsupervised contrastive learning
- (4) can add supervised contrastive learning

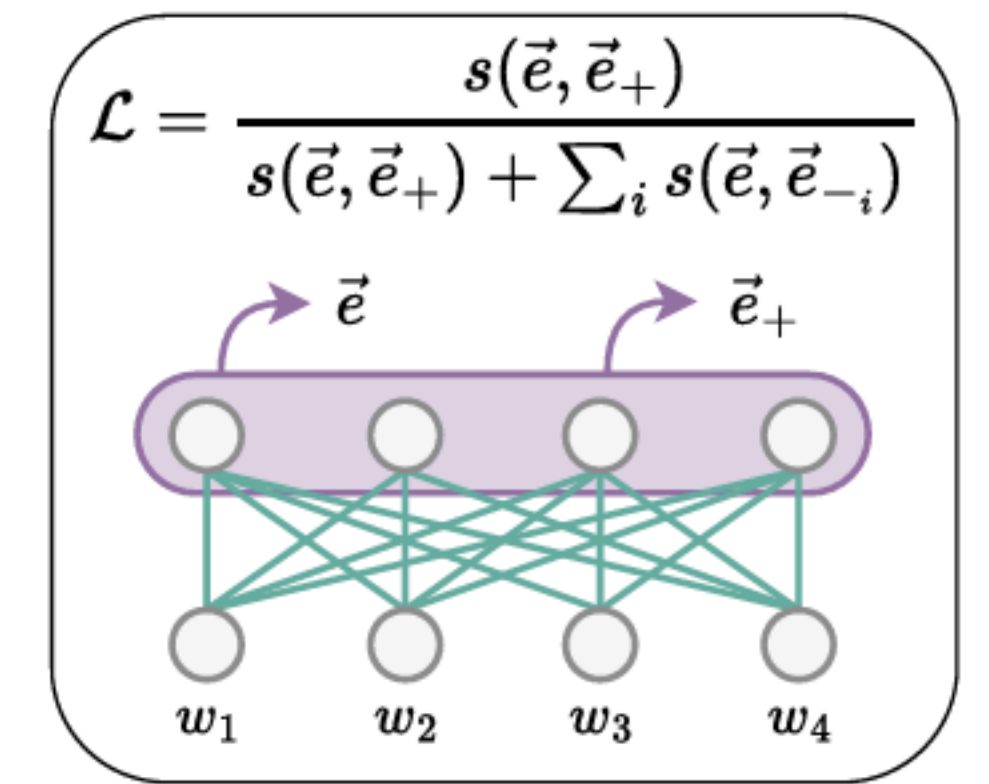
Enabling Bidirectional Attention



Masked Next Token Prediction



Unsupervised Contrastive Learning



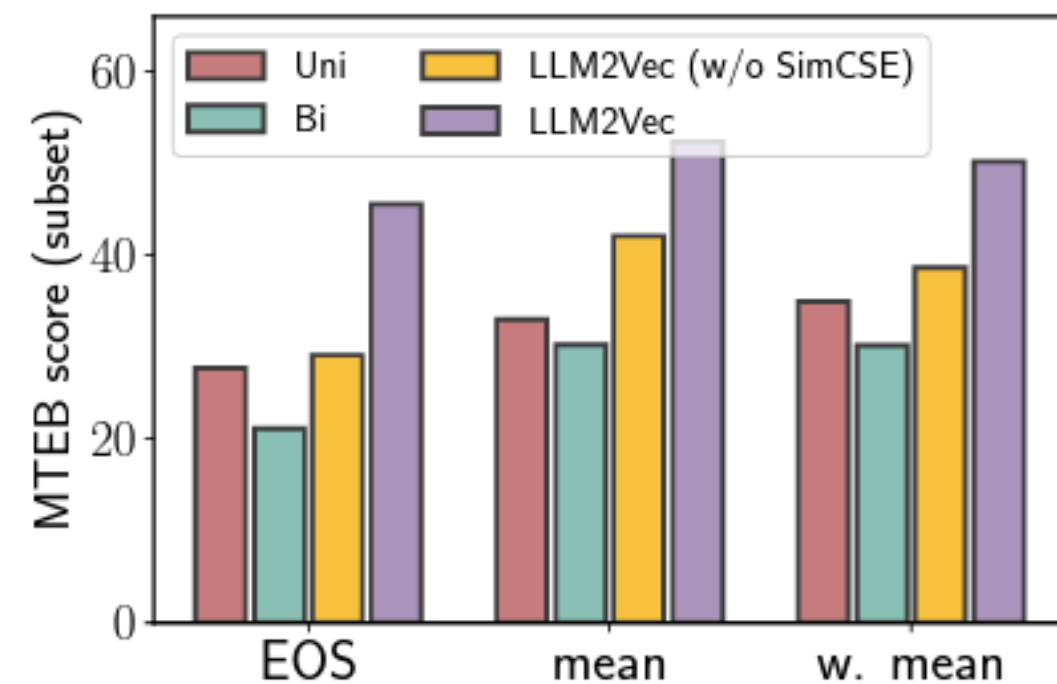
unsupervised contrastive learning via SimCSE: pass input sequence into the model twice with independently sampled dropout masks

Train model to maximise similarity between the two sequences representations and minimise similarities with others in batch

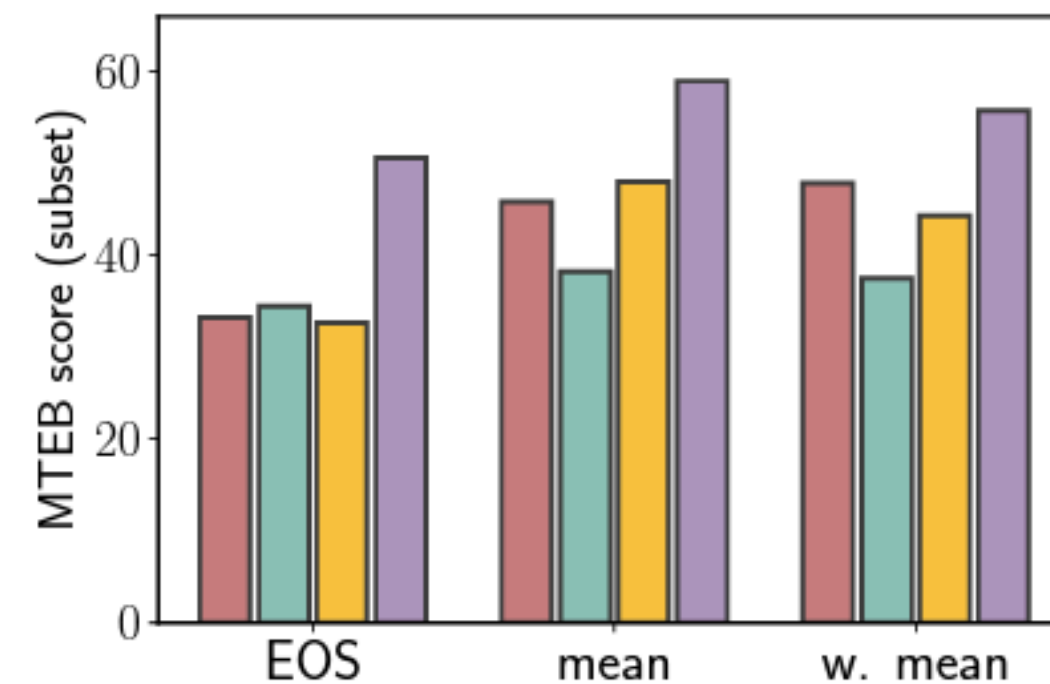
Results in two different representations for the same sequence

# LLM2Vec

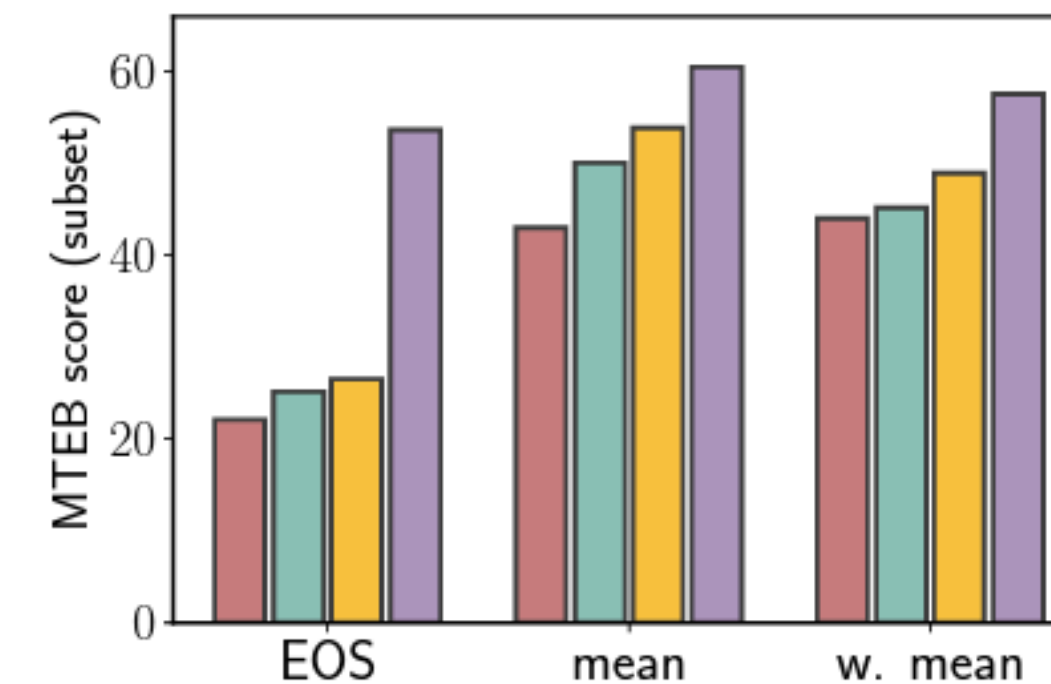
To obtain representations, can use EOS, mean pooling or weighted mean pooling



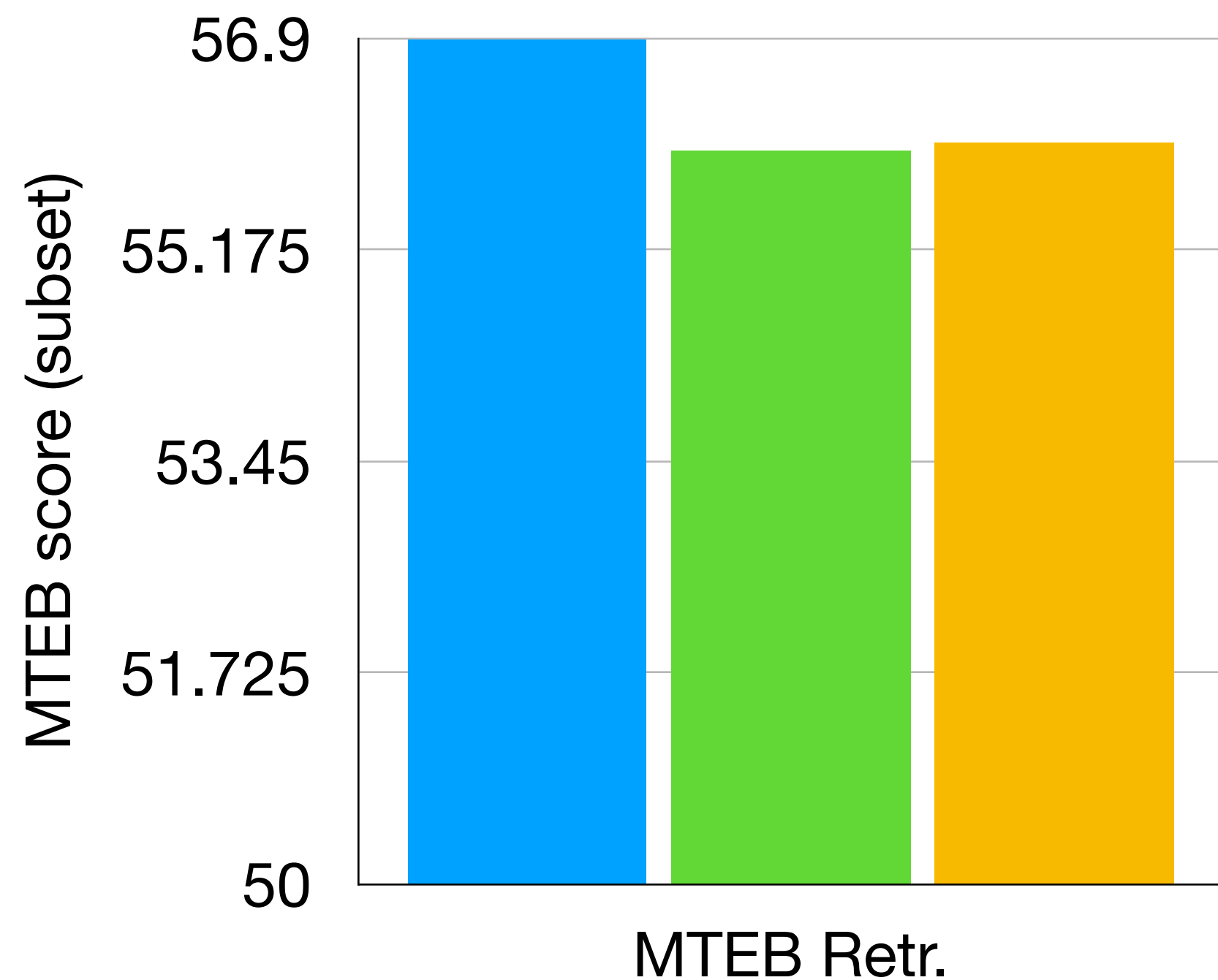
(a) S-LLaMA-1.3B



(b) LLaMA-2-7B



(c) Mistral-7B



■ E5-Mistral-7B  
■ LLM2Vec (w/o SimCSE)  
■ LLM2Vec

Effectiveness of LLM2Vec with Mistral-7B backbone & contrastive fine-tuning is lower than E5-Mistral on MTEB Retrieval datasets

(Bigger differences in the re-ranking task)

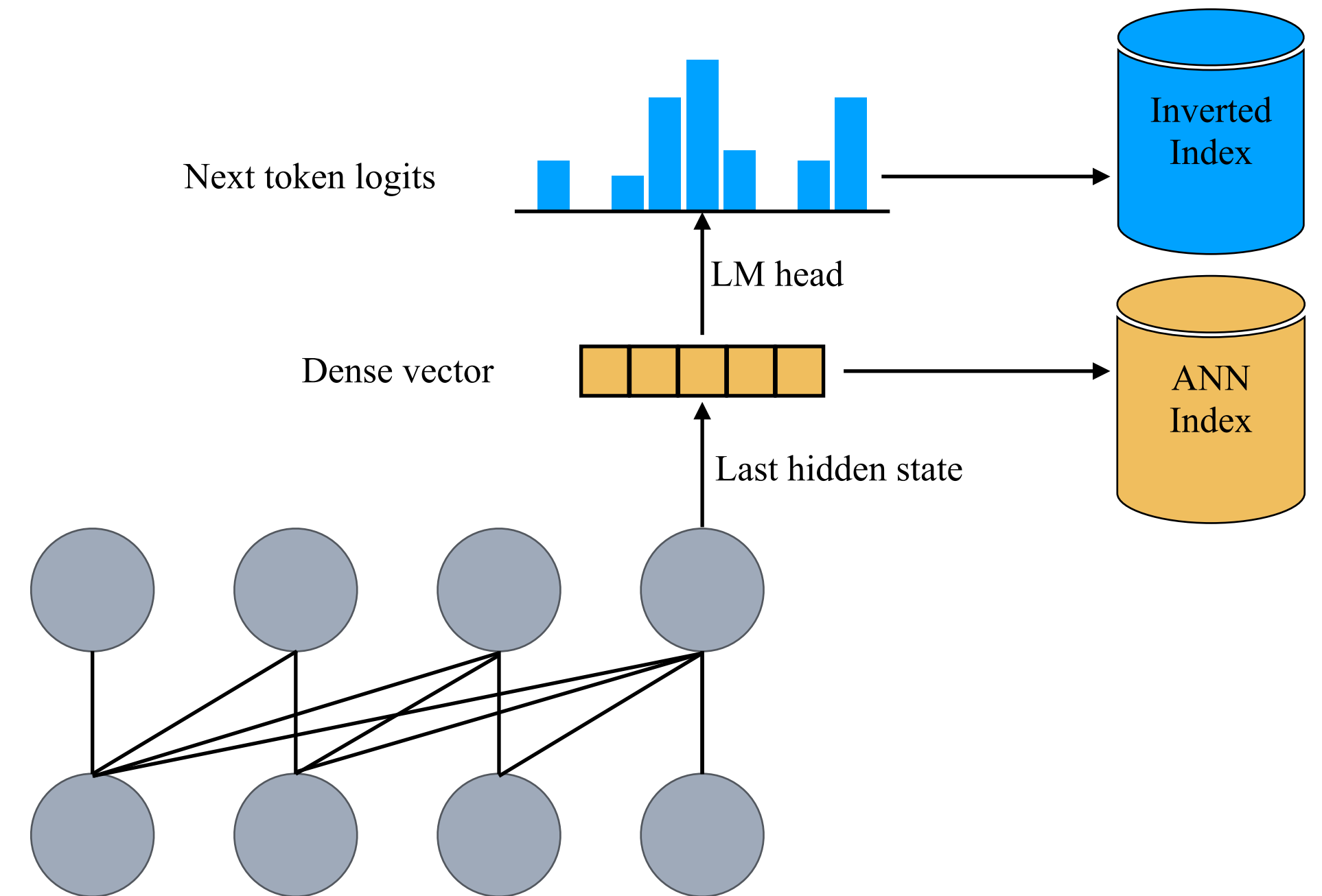
# PromptReps

- E5: gets embedding from last hidden layer of [EOS] token; requires (supervised) contrastive fine-tuning
- LLM2Vec: gets embedding from mean pooling of sequence tokens; requires (unsupervised/supervised) contrastive fine-tuning
- Can we engineer LLMs prompts to obtain an effective representation **without** need for contrastive fine-tuning?
- PromptReps! — “zero-shot” generation of representations using LLMs

# PromptReps

PromptReps at a high level:

1. **prompt the LLM** to represent given text (doc/query) using **one word**
2. Use last **token hidden layer** to obtain a **dense** representation
3. Use **logits** associated with last token hidden layer to obtain a **sparse** representation over the LLM vocabulary
4. **Combine the two** representations to get a hybrid retriever

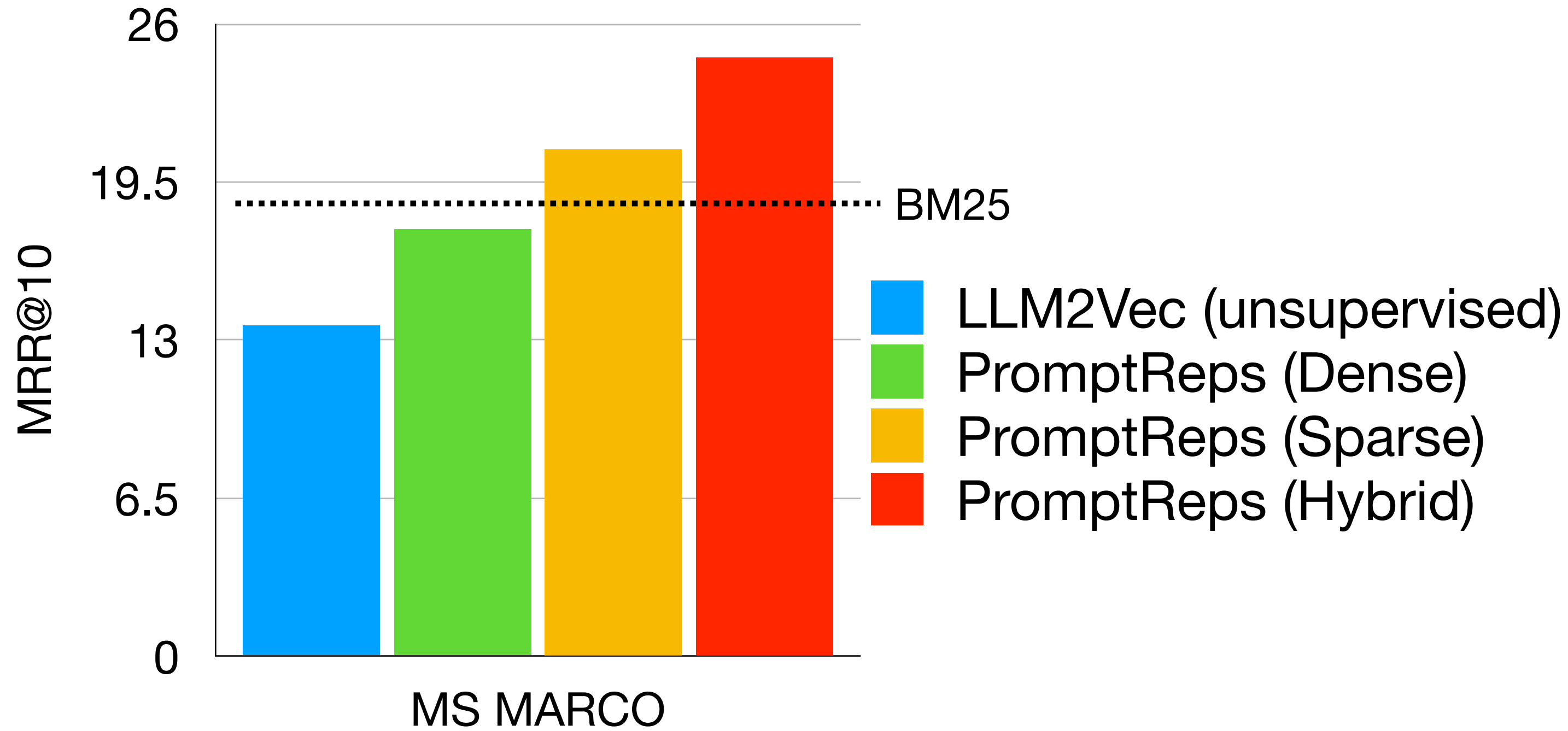


**<System>** You are an AI assistant that can understand human language.

**<User>** Passage: “[text]”. Use one word to represent the passage in a retrieval task. Make sure your word is in lowercase.

**<Assistant>** The word is: “

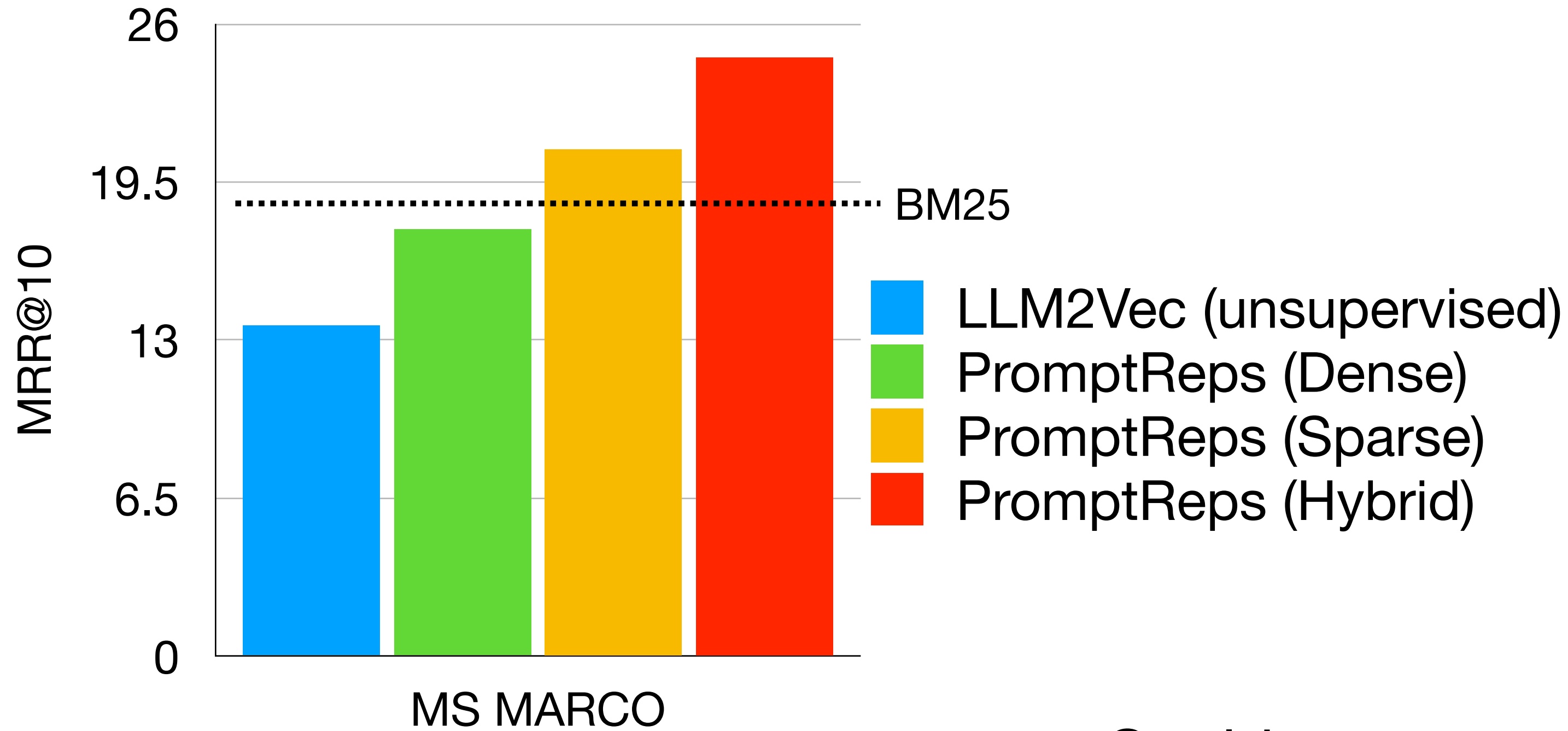
# PromptReps



Combining Dense & Sparse representation allows effective zero-shot retrieval (no training of representations)

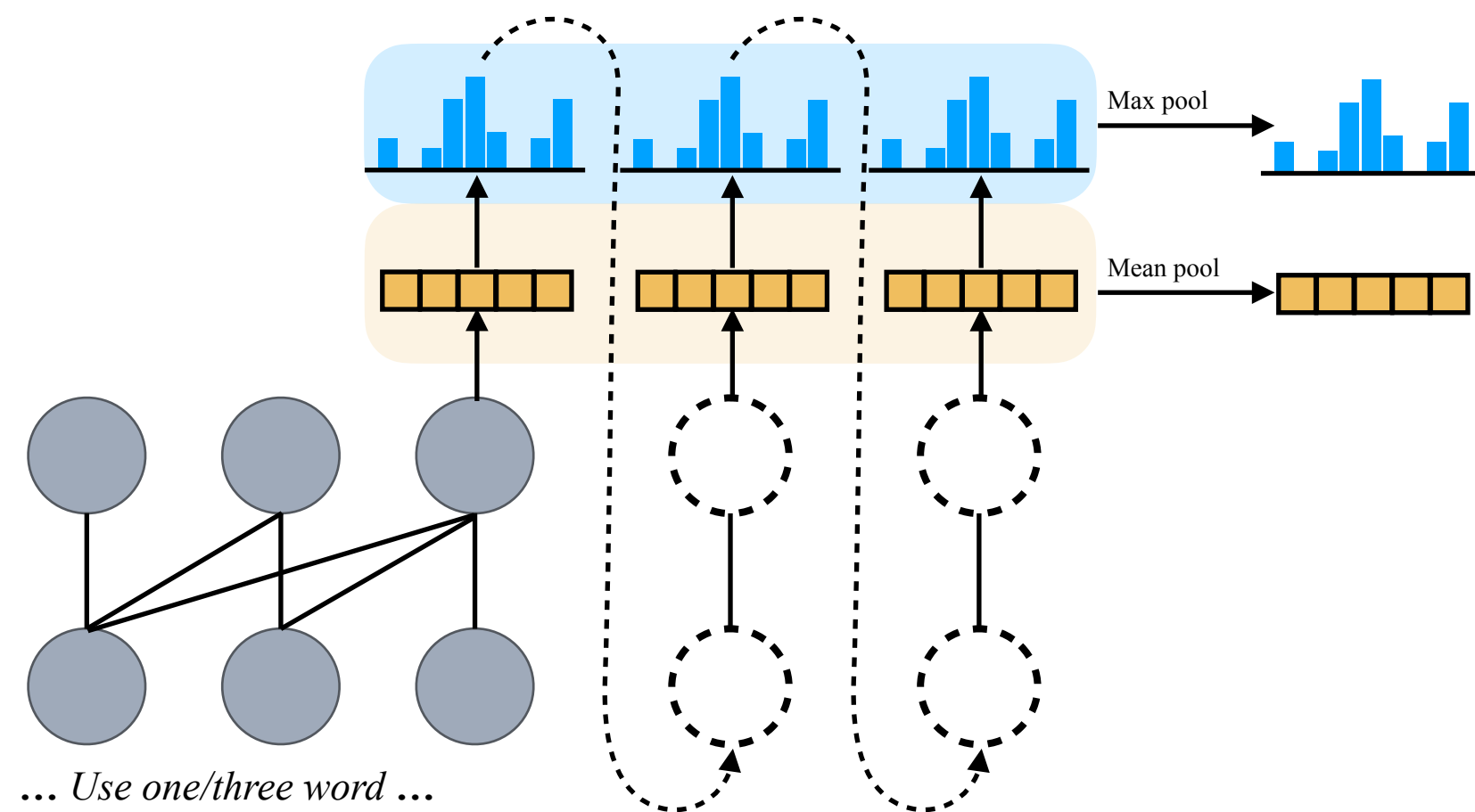
[LLM2Vec baseline here trained unsupervised with SimCSE]

# PromptReps



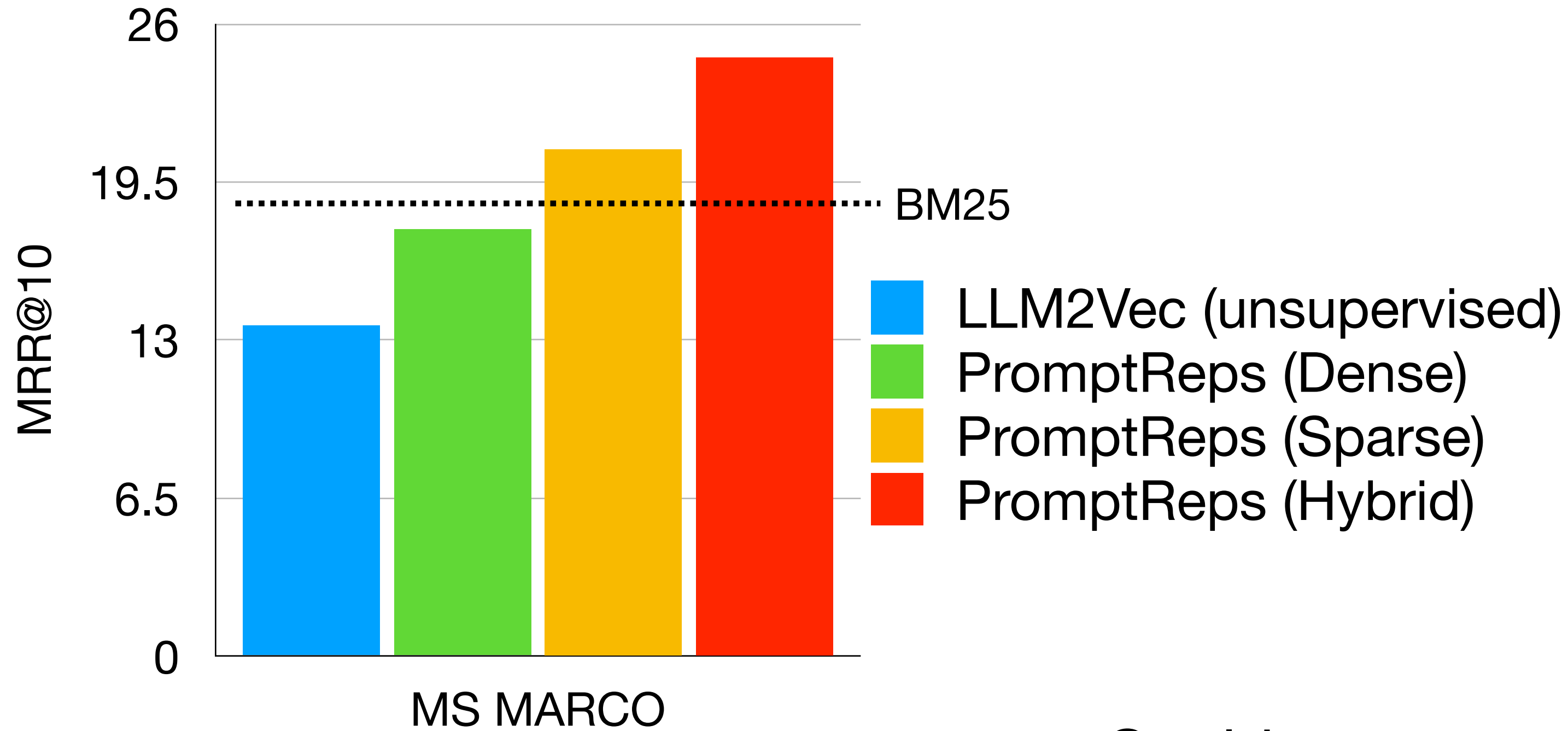
Combining Dense & Sparse representation allows effective zero-shot retrieval (no training of representations)

[LLM2Vec baseline here trained unsupervised with SimCSE]



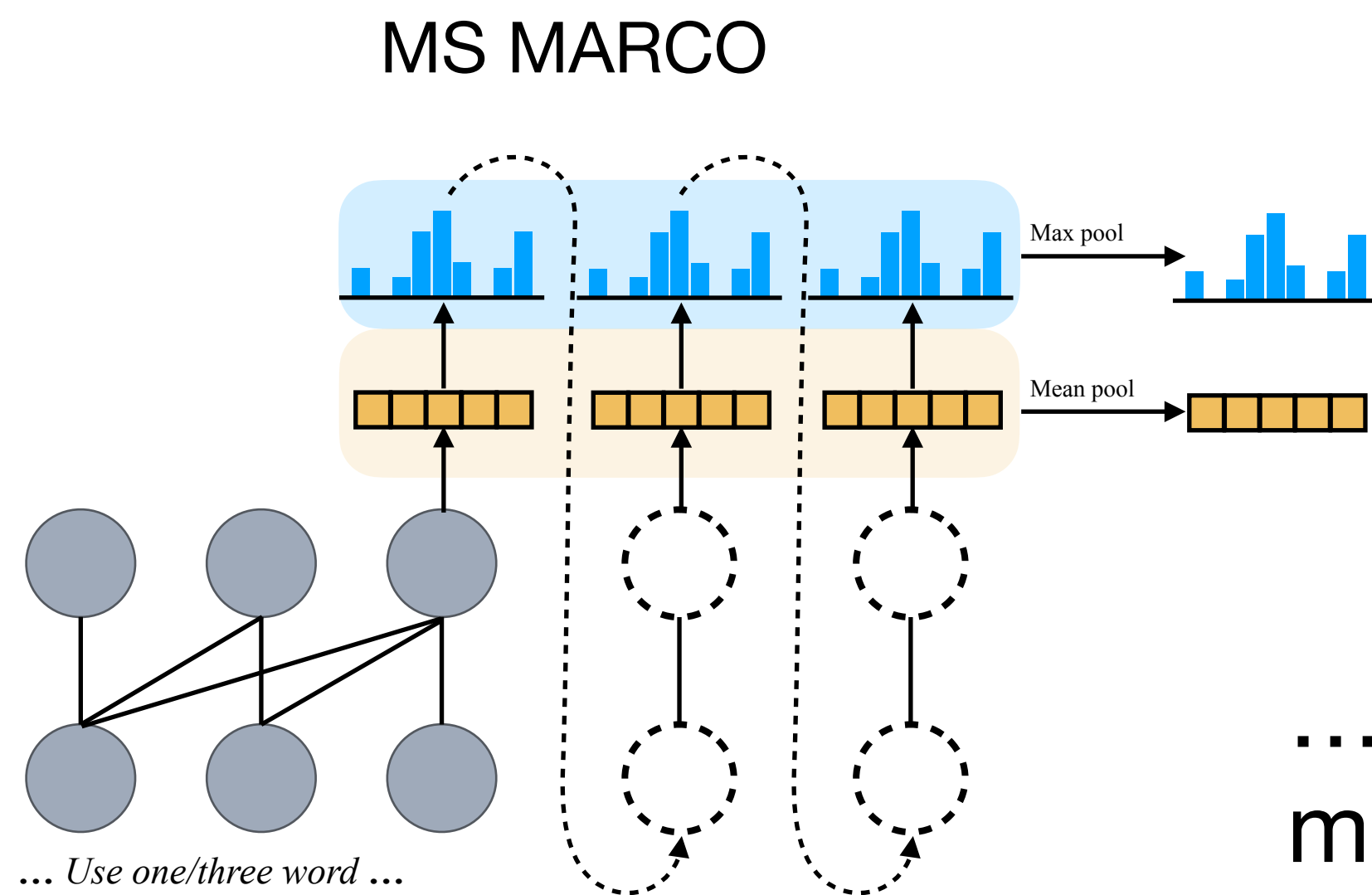
Could generate more words...

# PromptReps



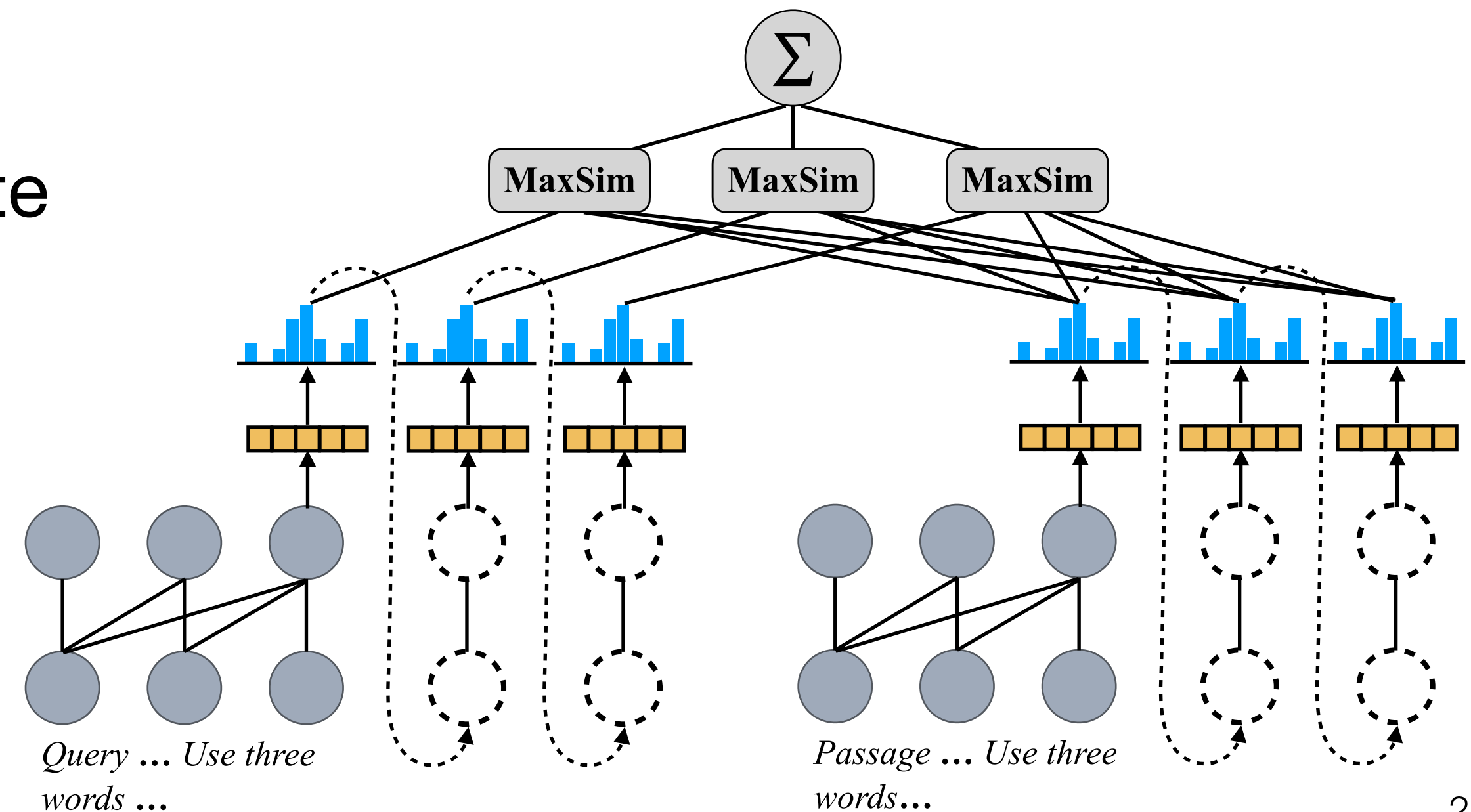
Combining Dense & Sparse representation allows effective zero-shot retrieval (no training of representations)

[LLM2Vec baseline here trained unsupervised with SimCSE]



Could generate more words...

... and use multiple tokens



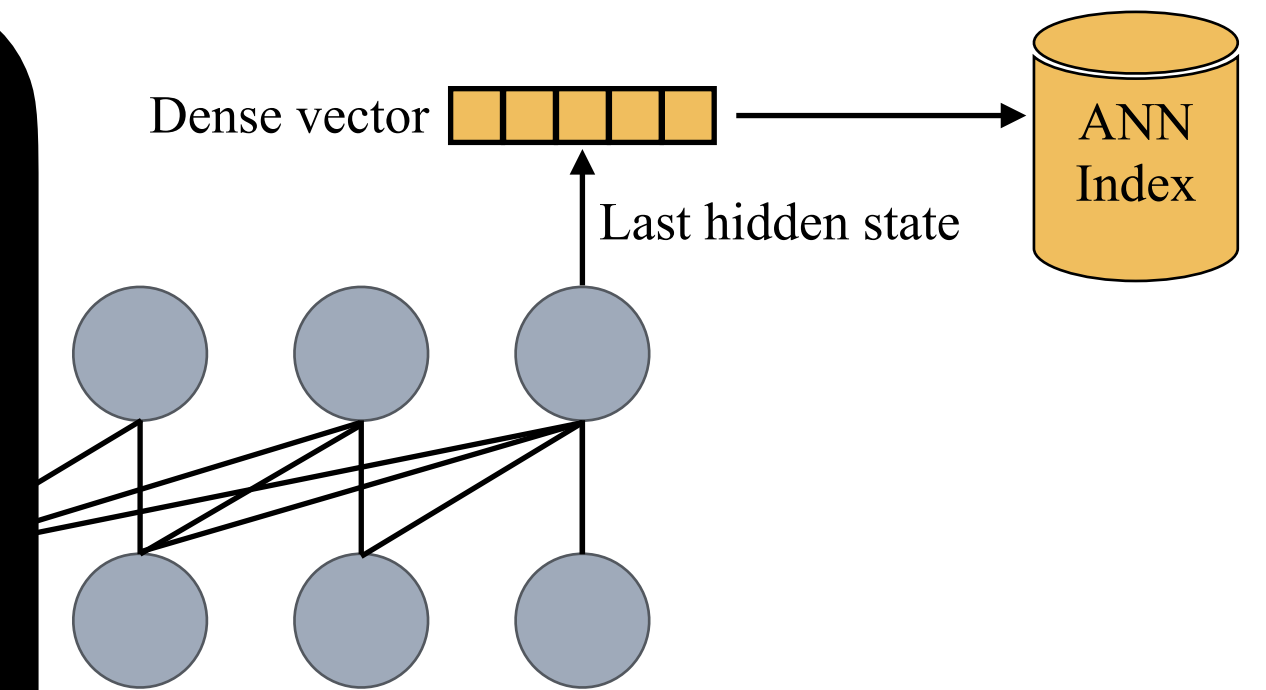
# Key Intuitions

Can we ask LLMs to:

- **Retrieve:** generate representations
- Then we can use representations to create embeddings and

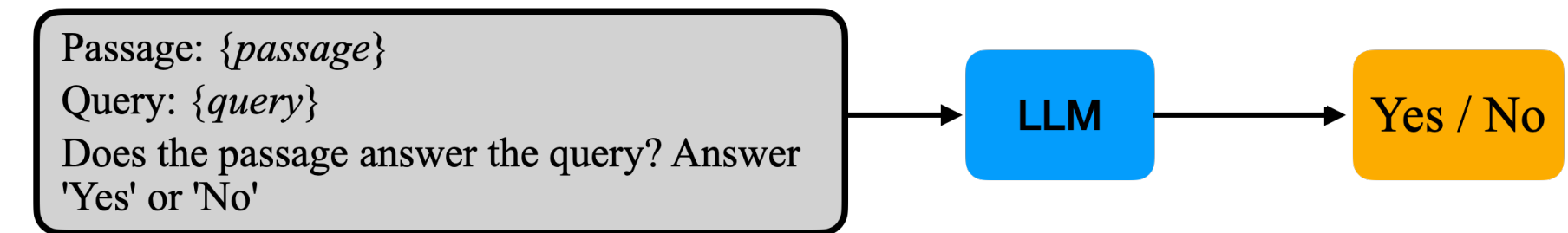
**Plan:**

- 1) Four “families” of LLM rankers
- 2) On prompt variations
- 3) On prompt optimisation



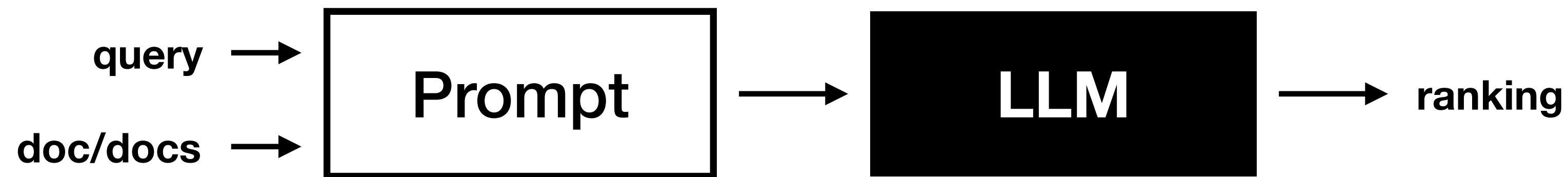
You are an AI assistant that can understand human language.  
 Passage: “[text]”. Use one word to represent the passage in a  
 retrieval task. Make sure your word is in lowercase.  
 <Assistant> The word is: “

- **Ranking:** tell us the relevance of a document to a query?
- Then we can use this indication of relevance (or relative relevance of n documents) to rank documents for the query



Idea: devise prompts/instructions to tell them LLM how to perform these tasks effectively

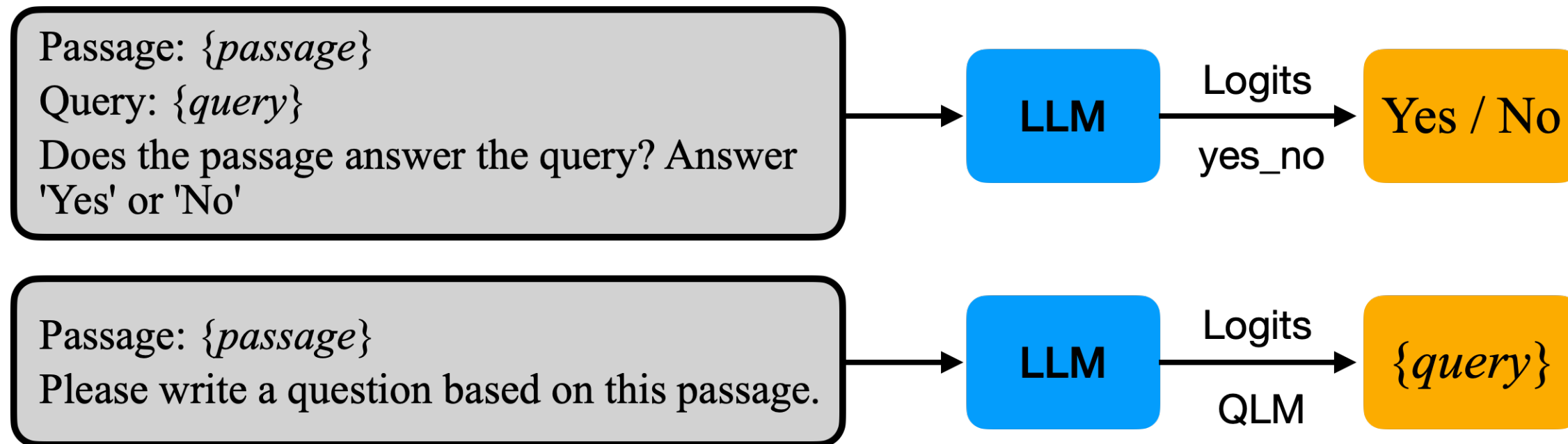
# LLMs as Rankers



- Four main families, characterised by how documents are passed in the prompt and how relevance of document to query is determined
- All are “zero-shot”: i.e. once you obtained the pre-trained, instruction tuned LLM, no need to do contrastive training

# LLMs as Rankers: Pointwise

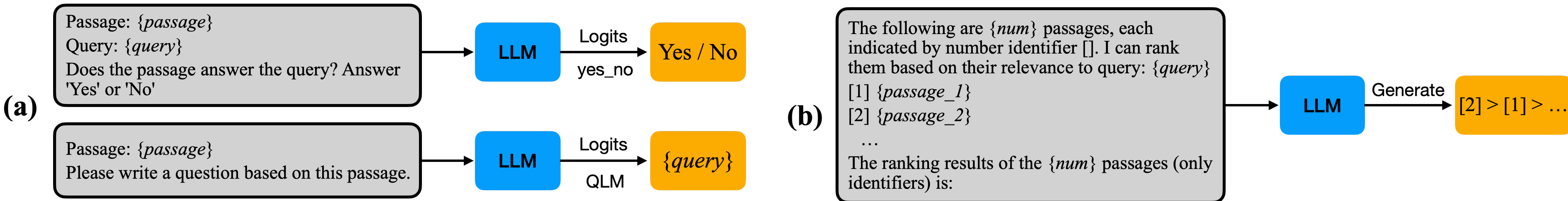
(a)



Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A. and Newman, B., 2022. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.

Sachan, D.S., Lewis, M., Joshi, M., Aghajanyan, A., Yih, W.T., Pineau, J. and Zettlemoyer, L., 2022. Improving passage retrieval with zero-shot question generation. *arXiv preprint arXiv:2204.07496*.

# LLMs as Rankers: Listwise

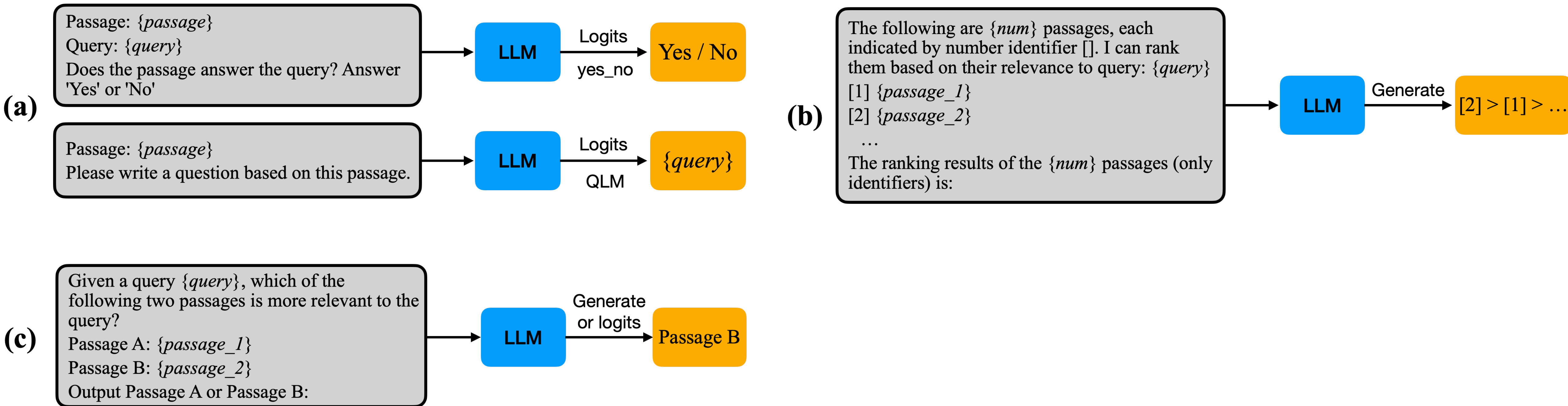


Ma, X., Zhang, X., Pradeep, R. and Lin, J., 2023. Zero-Shot Listwise Document Reranking with a Large Language Model. *arXiv preprint arXiv:2305.02156*.

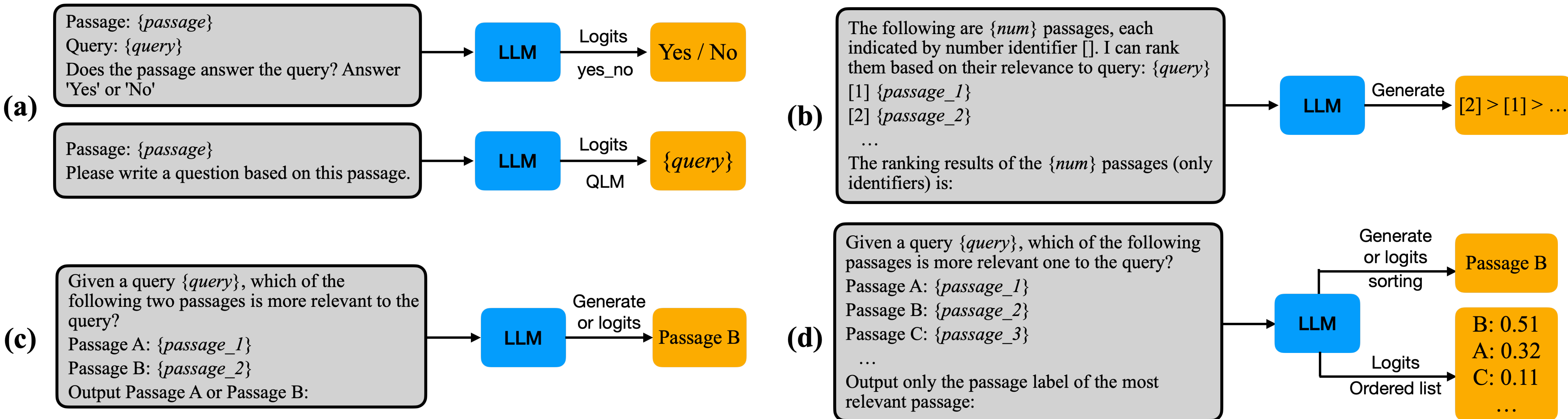
Pradeep, R., Sharifymoghaddam, S. and Lin, J., 2023. RankVicuna: Zero-Shot Listwise Document Reranking with Open-Source Large Language Models. *arXiv preprint arXiv:2309.15088*.

Sun, W., Yan, L., Ma, X., Ren, P., Yin, D. and Ren, Z., 2023. Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agent. *arXiv preprint arXiv:2304.09542*.

# LLMs as Rankers: Pairwise



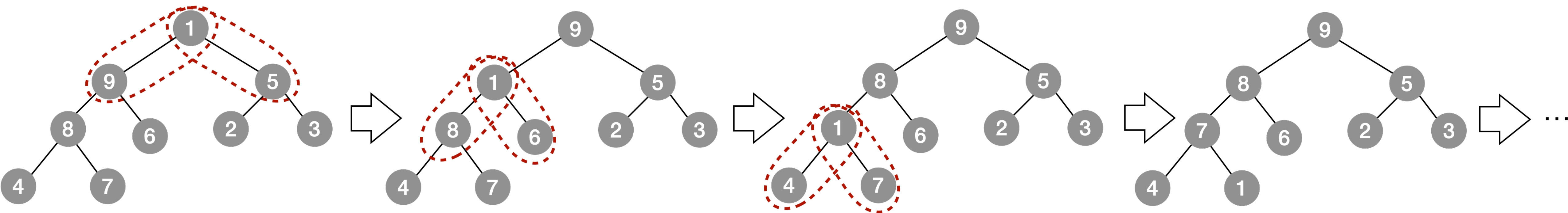
# LLMs as Rankers: Setwise



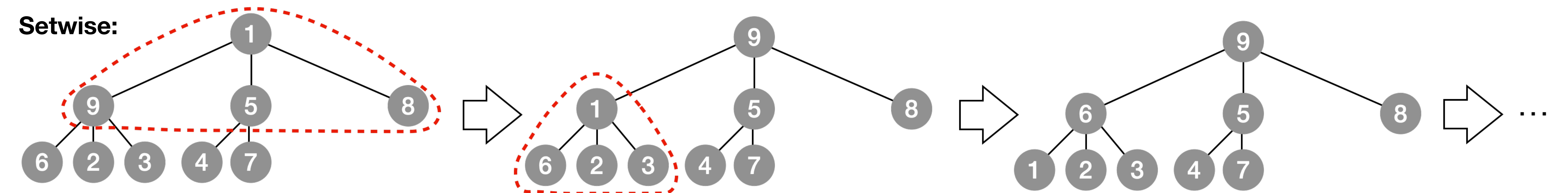
# LLMs as Rankers: Setwise

- Setwise offers two advantages
  - Compared to listwise: It can rely on logits rather than generation: faster
  - Compared to pairwise: it requires less comparisons (i.e. inferences with LLM)

Pairwise:



Setwise:



# LLMs as Rankers

Method	NDCG@10	TREC DL 2019	TREC DL 2020
setwise.bubblesort	0.6800	\$4.62	\$4.67
pairwise.allpair	0.6783	\$90.59	\$90.6
setwise.heapsort	0.6743	\$1.28	\$1.27
pointwise.yes_no	0.6398	\$0.48	\$0.49
listwise.generation	0.5929	\$3.49	\$3.75
pointwise.qlm	0.5343	\$0.46	\$0.46

Zhuang, S., Zhuang, H., Koopman, B. and Zuccon, G., 2023. A Setwise Approach for Effective and Highly Efficient Zero-shot Ranking with Large Language Models. *arXiv preprint arXiv:2310.09497*.

# LLMs as Rankers

Method	NDCG@10	TREC DL 2019	TREC DL 2020
<u>pairwise.heapsort</u>	0.6800	\$3.39	\$3.4
setwise.bubblesort	0.6800	\$4.62	\$4.67
pairwise.allpair	0.6783	\$90.59	\$90.6
listwise.likelihood	0.6745	\$2.83	\$2.86
<u>setwise.heapsort</u>	0.6743	\$1.28	\$1.27
pairwise.bubblesort	0.6550	\$12.28	\$11.89
pointwise.yes_no	0.6398	\$0.48	\$0.49
listwise.generation	0.5929	\$3.49	\$3.75
pointwise.qlm	0.5343	\$0.46	\$0.46

0.8% loss in nDCG

62% reduction in cost

Zhuang, S., Zhuang, H., Koopman, B. and Zuccon, G., 2023. A Setwise Approach for Effective and Highly Efficient Zero-shot Ranking with Large Language Models. *arXiv preprint arXiv:2310.09497*.

# Problems with ListWise

- LLM has a **maximum context(prompt) length**, which limits the number of passages in a list
- Because ranking needs to be broken into many sublists, **listwise requires many inferences -> high latency**
- Also **generating** the string with the **ranking** can be **difficult & slow**: LLM may output wrong format, requiring to re-generate, or produces additional content (e.g. explanation), making generation slow

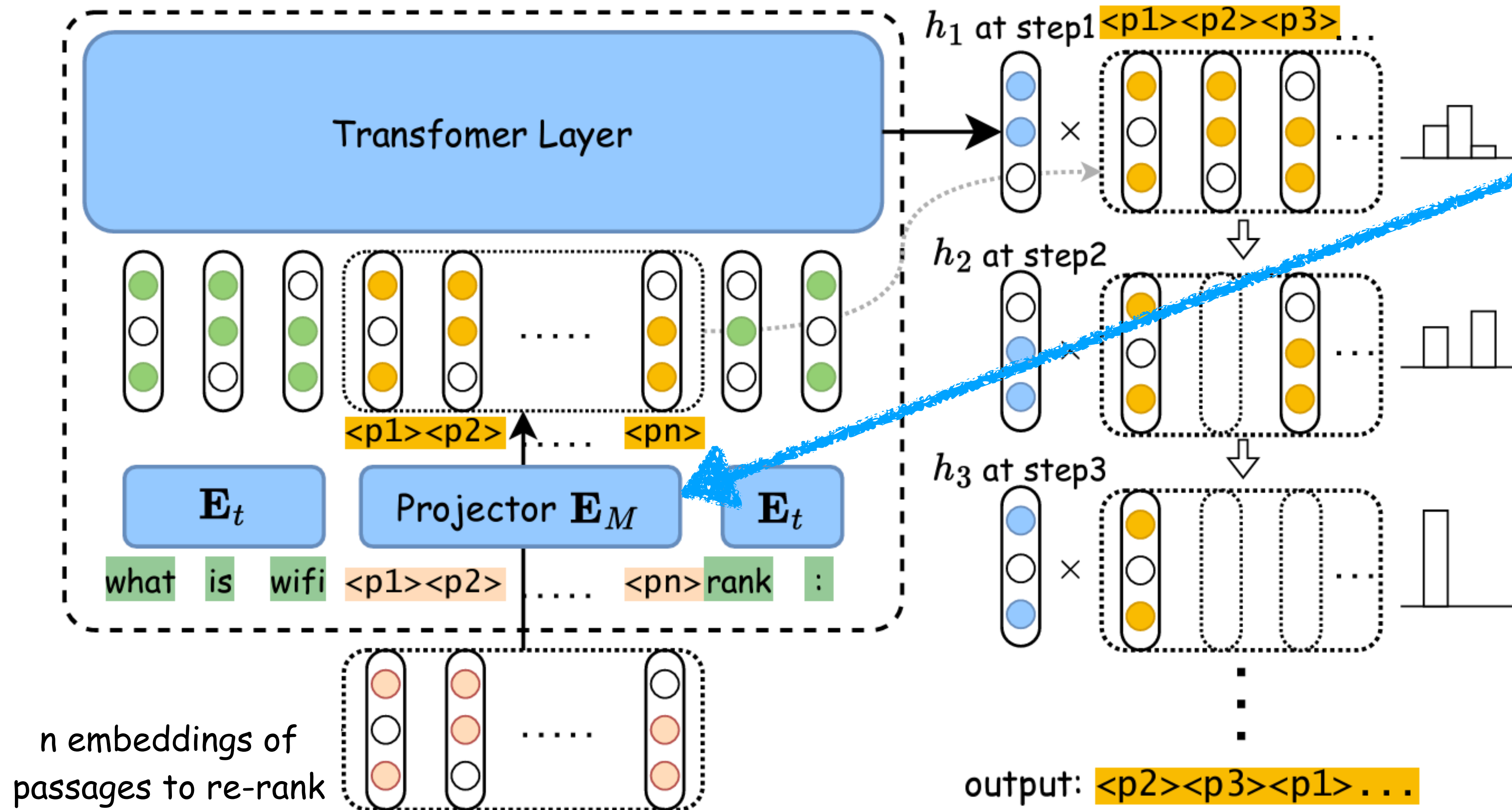
# PE-Rank: passage compression

- LLM has a **maximum context(prompt) length**, which limits the number of passages in a list
- Because ranking needs to be broken into many sublists, **listwise requires many inferences -> high latency**
- Also **generating** the string with the **ranking** can be **difficult & slow**: LLM may output wrong format, requiring to re-generate, or produces additional content (e.g. explanation), making generation slow
- PE-Rank:
  1. Treat passage as a single special token
  2. Use the passage embedding as token representation
  3. Dynamic constrain of decoding space for special tokens

**Reduce input length**

**Accelerate decoding process**

# PE-Rank: passage compression



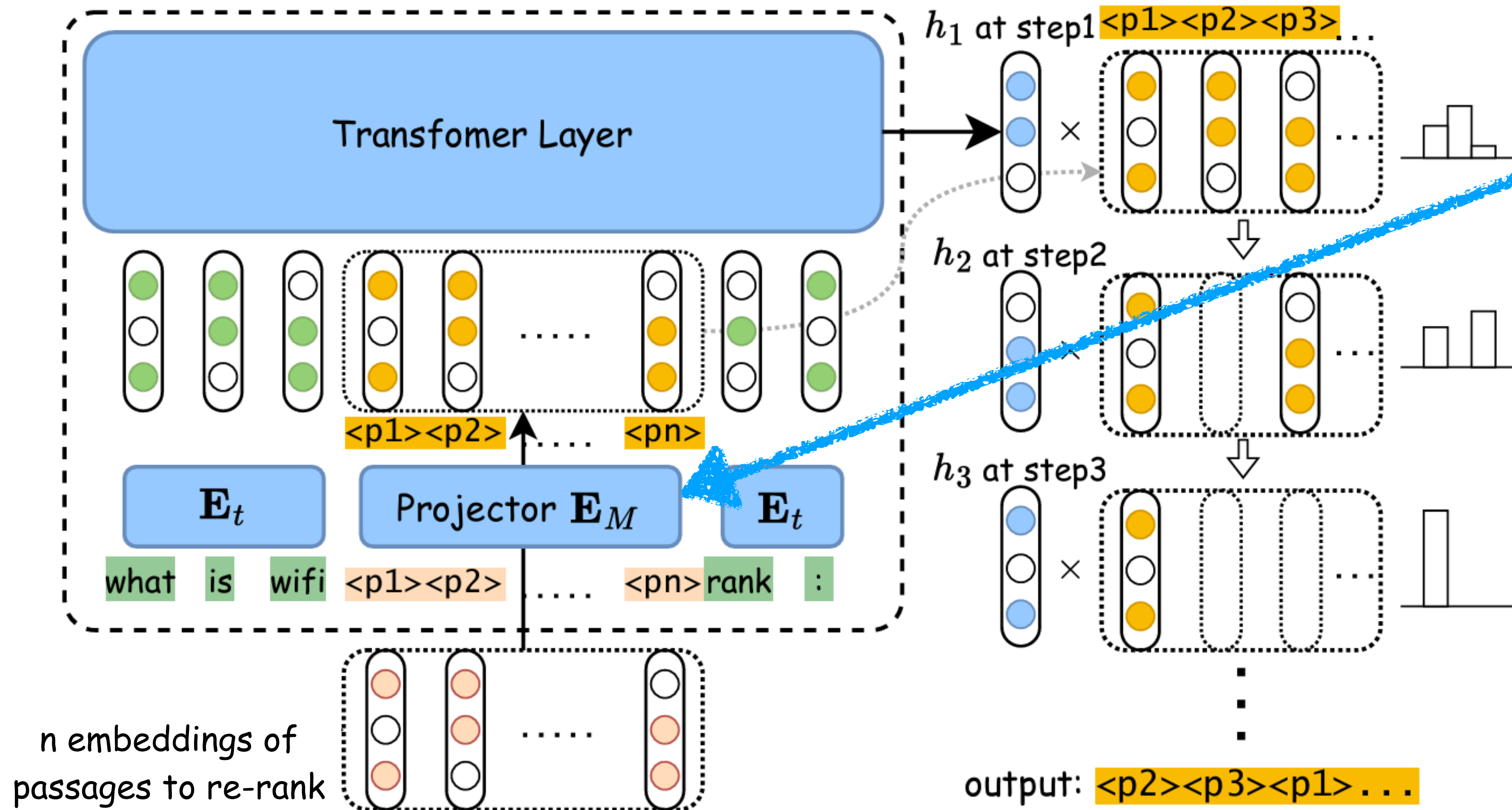
This is needed because of dimensional and distributional differences between passage embeddings and LLM token embeddings

EM is a 2-layer MLP that acts as a learned mapping function

n embeddings of passages to re-rank

output:  $\langle p2 \rangle \langle p3 \rangle \langle p1 \rangle \dots$

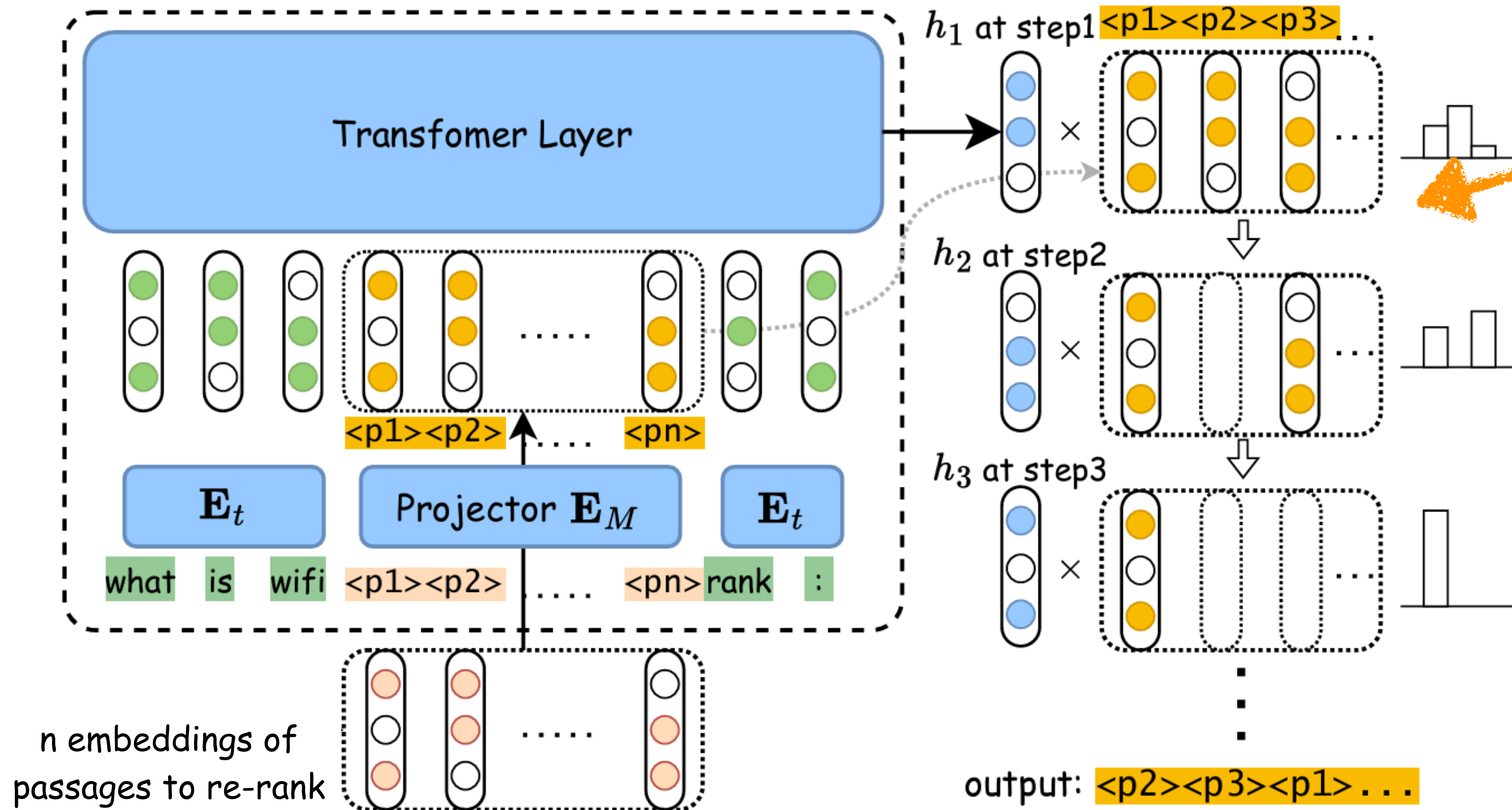
# PE-Rank: passage compression



- Training: text reconstruction task**
1. encode  $t$ , pass to MLP
  2. Take transformed embedding, pass it to LLM
  3. Prompt LLM to reconstruct text
- (Only MLP is trained, LLM & embeddings are frozen)

Training data: Wikipedia

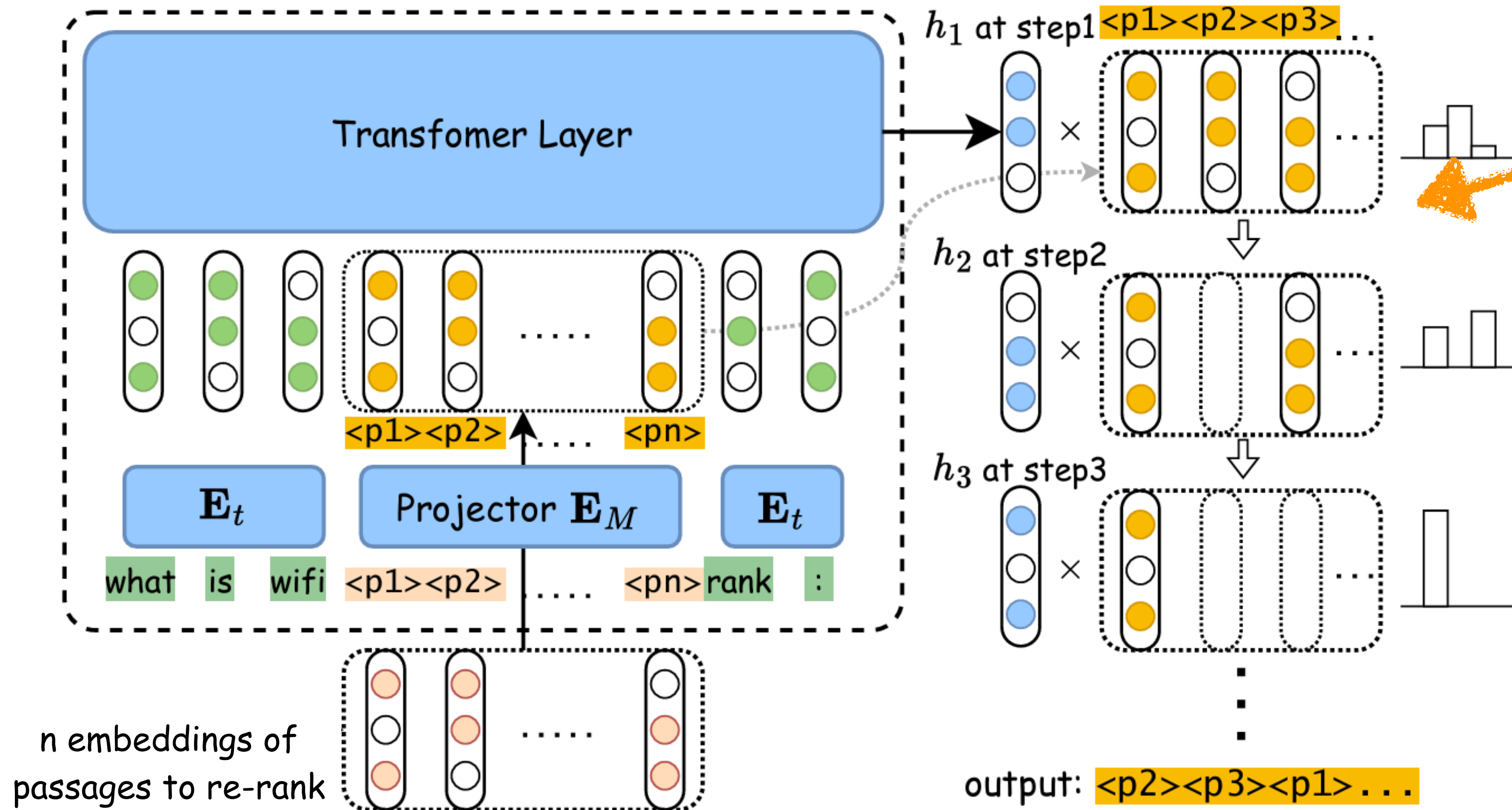
# PE-Rank: passage compression



Constrain generation to only special tokens

- \* smaller vocabulary -> faster inference
- \* remove already generated tokens -> further inference acceleration

# PE-Rank: passage compression



Training: sequential ranking learning process — at each step, provide previous decoded ranking, then maximise probability of generating next most relevant passage (i.e. ListMLE)

+ content loss + distillation loss

Training data: MS MARCO

# PE-Rank: passage compression

Common listwise

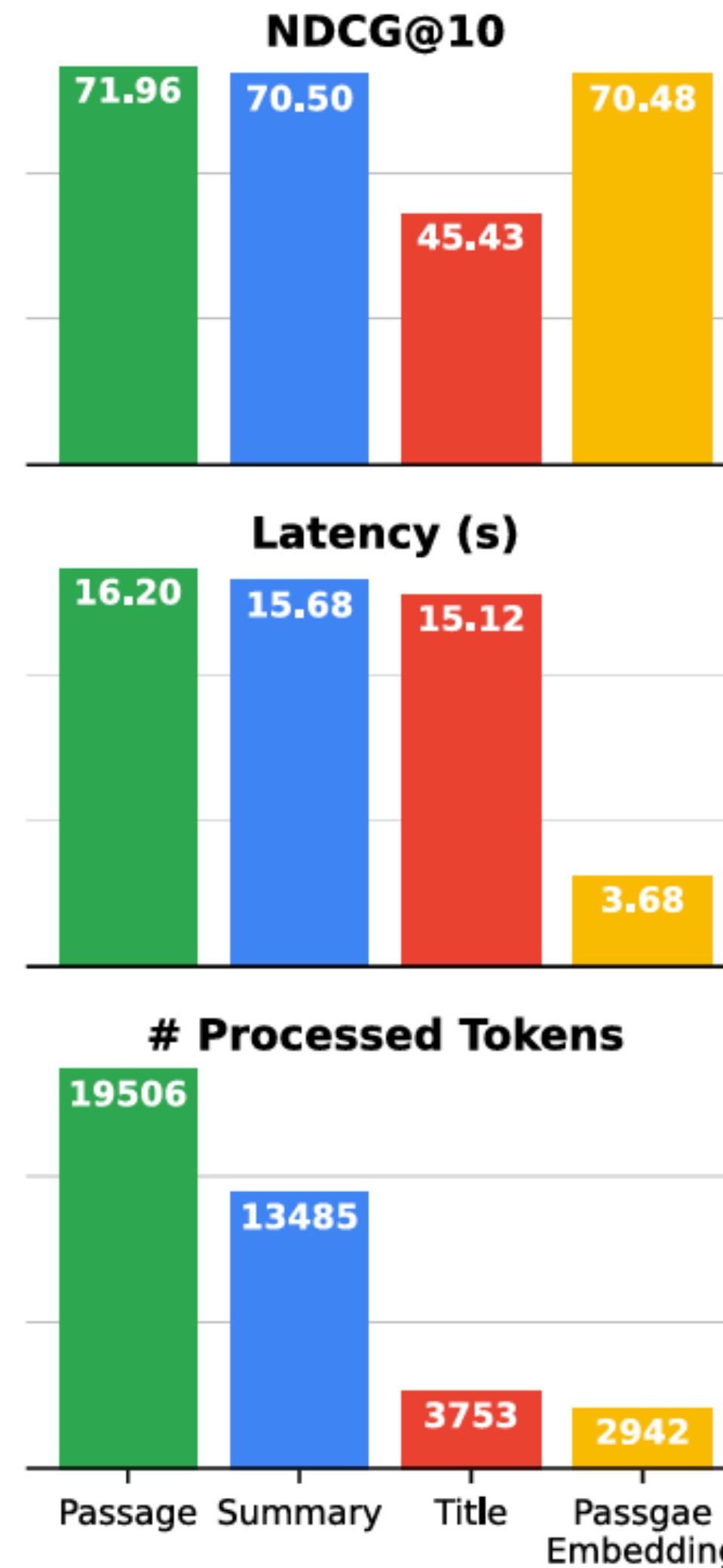
The following are passages related to query #{query}.  
 Passage 1: #{passage 1}  
 ...  
 Rank these passages based on their relevance to the query.

[2] > [3] > [1] ...

PE-Rank

The following are passages related to query #{query}, each with a special token representing the passage enclosed in [].  
 Passage 1: [<p1>]  
 ...  
 Rank these passages based on their relevance to the query.

<p2><p3><p1> ...



# PE-Rank: passage compression

Common listwise

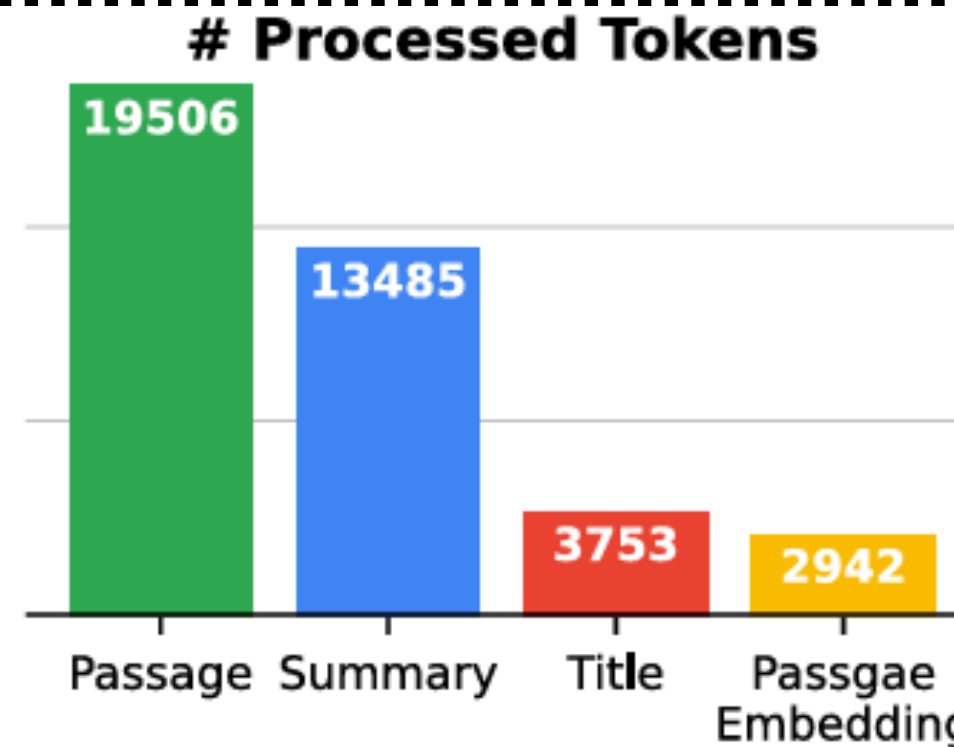
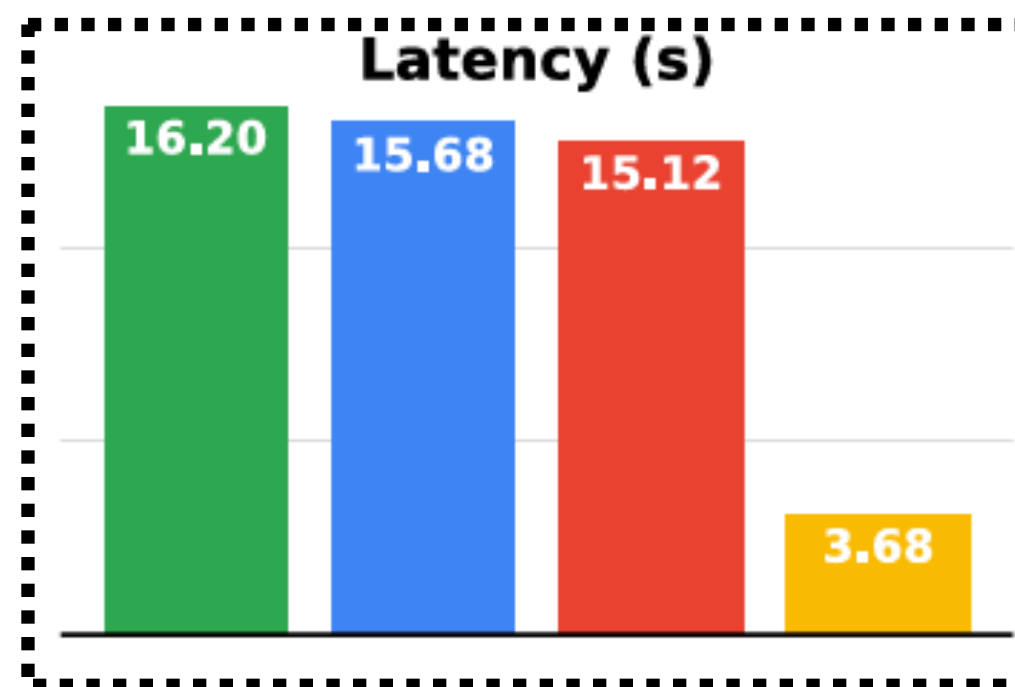
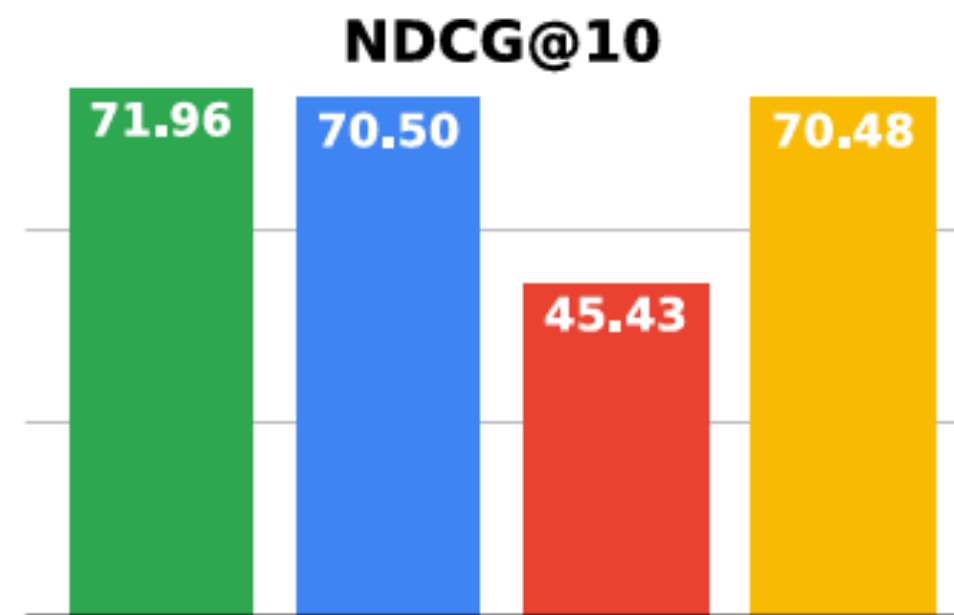
The following are passages related to query #{query}.  
 Passage 1: #{passage 1}  
 ...  
 Rank these passages based on their relevance to the query.

[2] > [3] > [1] ...

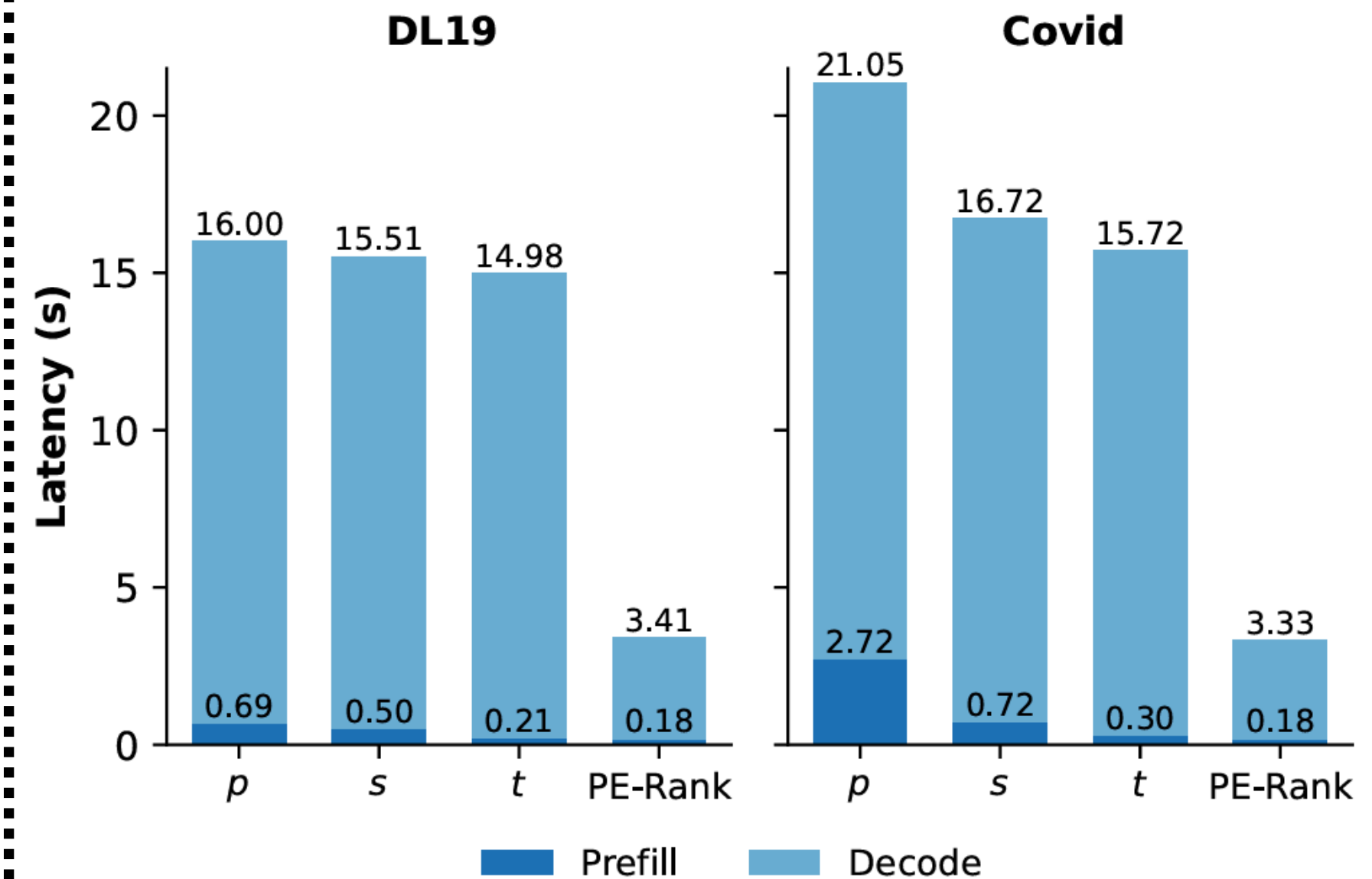
PE-Rank

The following are passages related to query #{query}, each with a special token representing the passage enclosed in [].  
 Passage 1: [<p1>]  
 ...  
 Rank these passages based on their relevance to the query.

<p2><p3><p1> ...



Most latency advantages arise from decoding, little gain in profiling phase



# Impact of Prompt Variations on LLM Rankers

- Prompts proposed for different LLM ranker methods vary largely
- Not just in terms of instruction for ranking, but also for (unrelated) additional wordings, e.g. role playing, and ordering of components (e.g. passage first, then query — or vice versa?)
- What are the effects of prompt wordings on methods? What makes a good prompt?

Method	Prompt
PRP, Qin et al.	Passage: {text} Query: {query} Does the passage answer the query?
RankGPT, Sun et al.	<p>You are RankGPT, an intelligent assistant that can rank passages based on their relevancy to the query. I will provide you with num passages, each indicated by number identifier []. Rank the passages based on their relevance to query: {query}.</p> <p>{PASSAGES}</p> <p>Search Query: {query}.</p> <p>Rank the num passages above based on their relevance to the search query. The passages should be listed in descending order using identifiers. The most relevant passages should be listed first. The output format should be [] &gt; [], e.g., [1] &gt; [2]. Only response the ranking results, do not say any word or explain.</p>

Role Playing

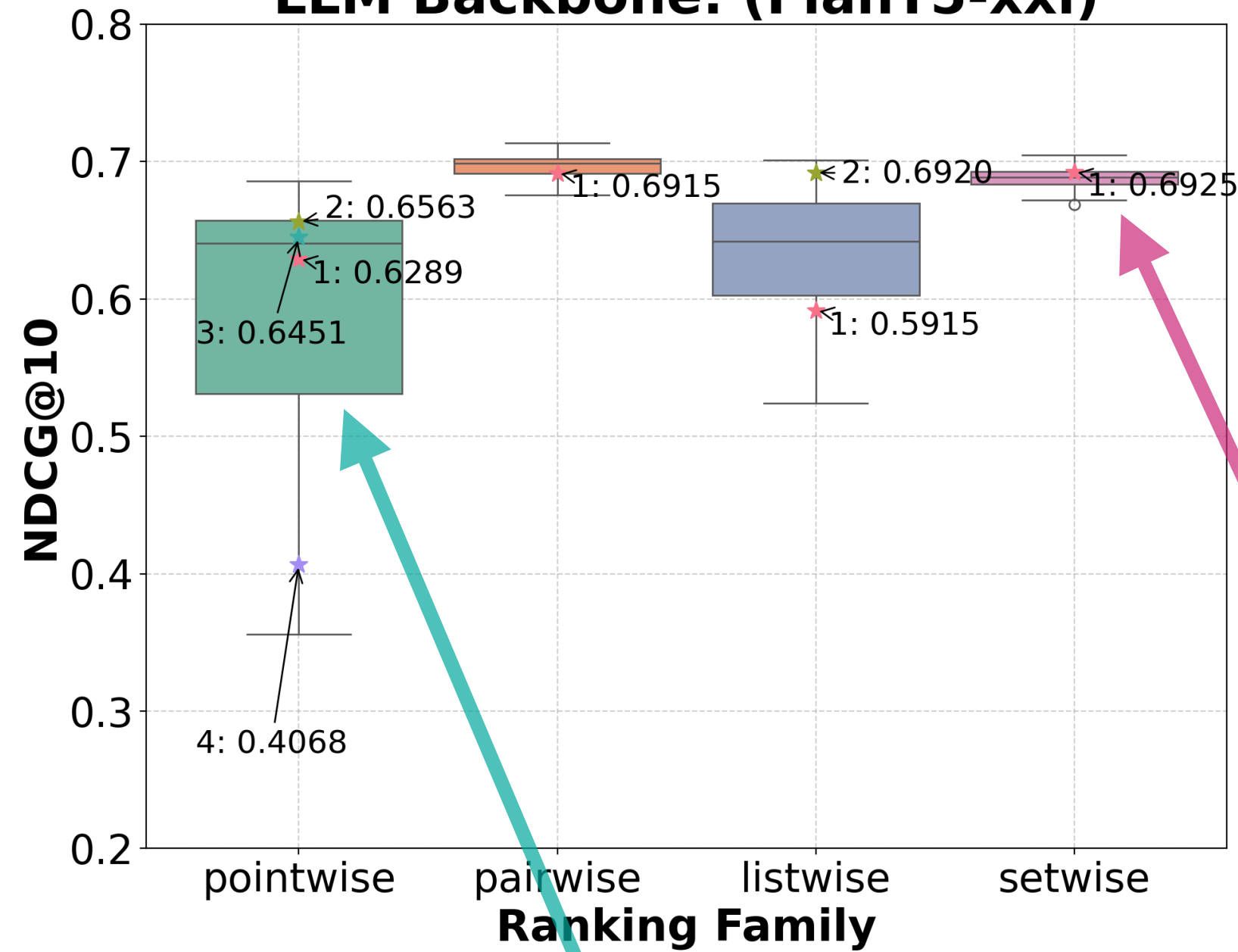
Formatting Instr.

Restriction on Output

# Impact of Prompt Variations on LLM Rankers

TREC DL 2019

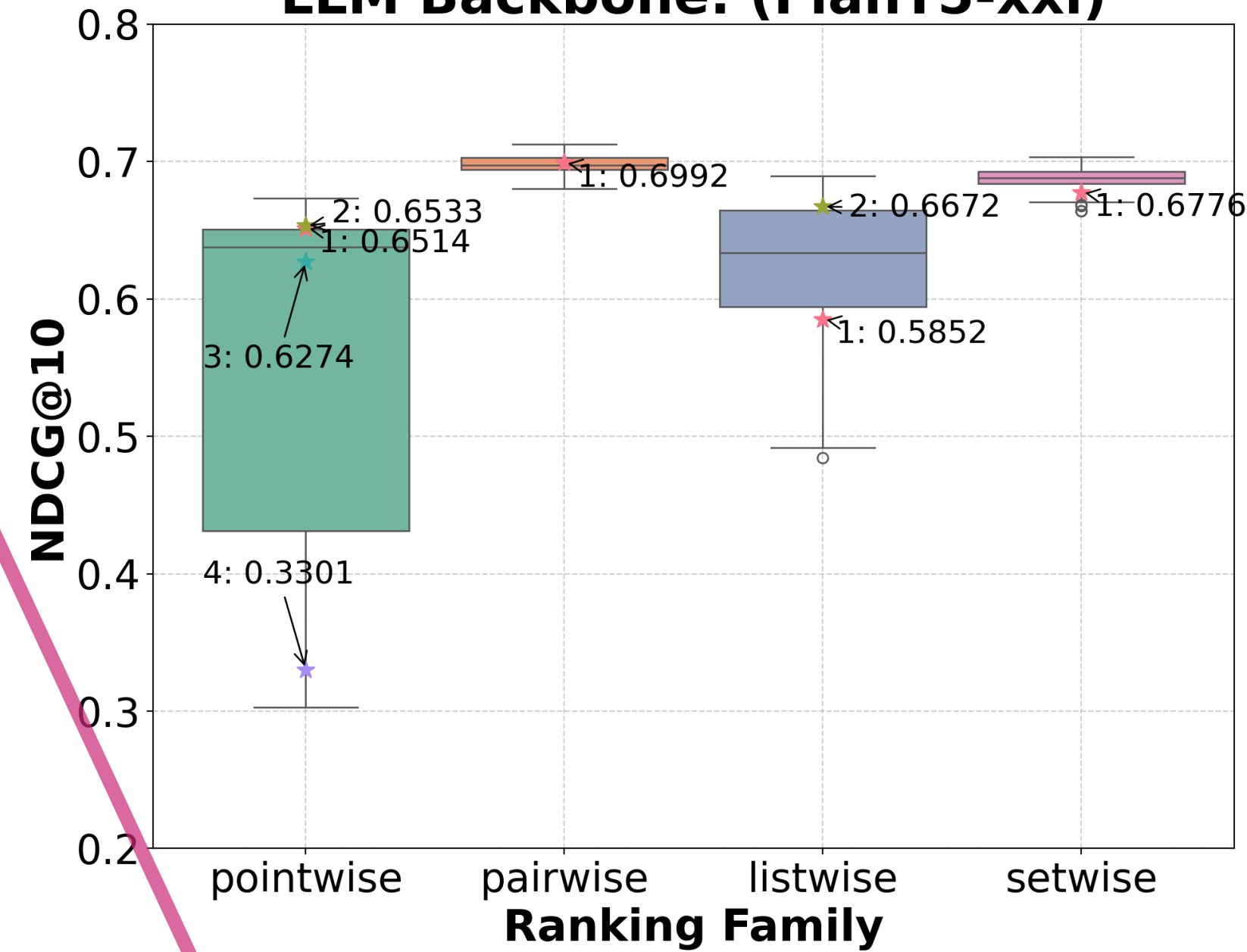
LLM Backbone: (FlanT5-xxl)



Pointwise, listwise often display large variations of effectiveness across different prompts

TREC DL 2020

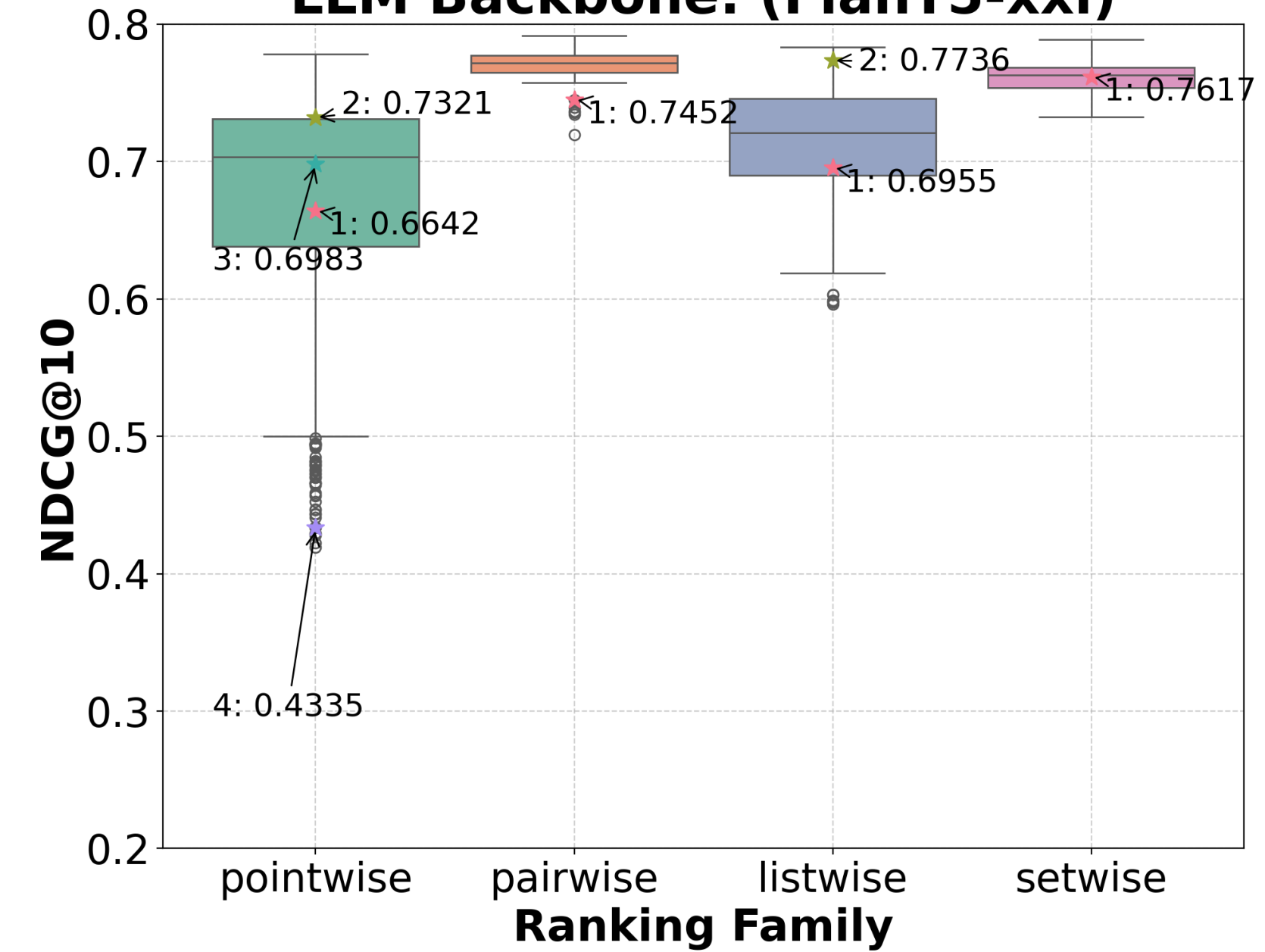
LLM Backbone: (FlanT5-xxl)



Pairwise, set wise very robust across different prompts

COVID (BEIR)

LLM Backbone: (FlanT5-xxl)

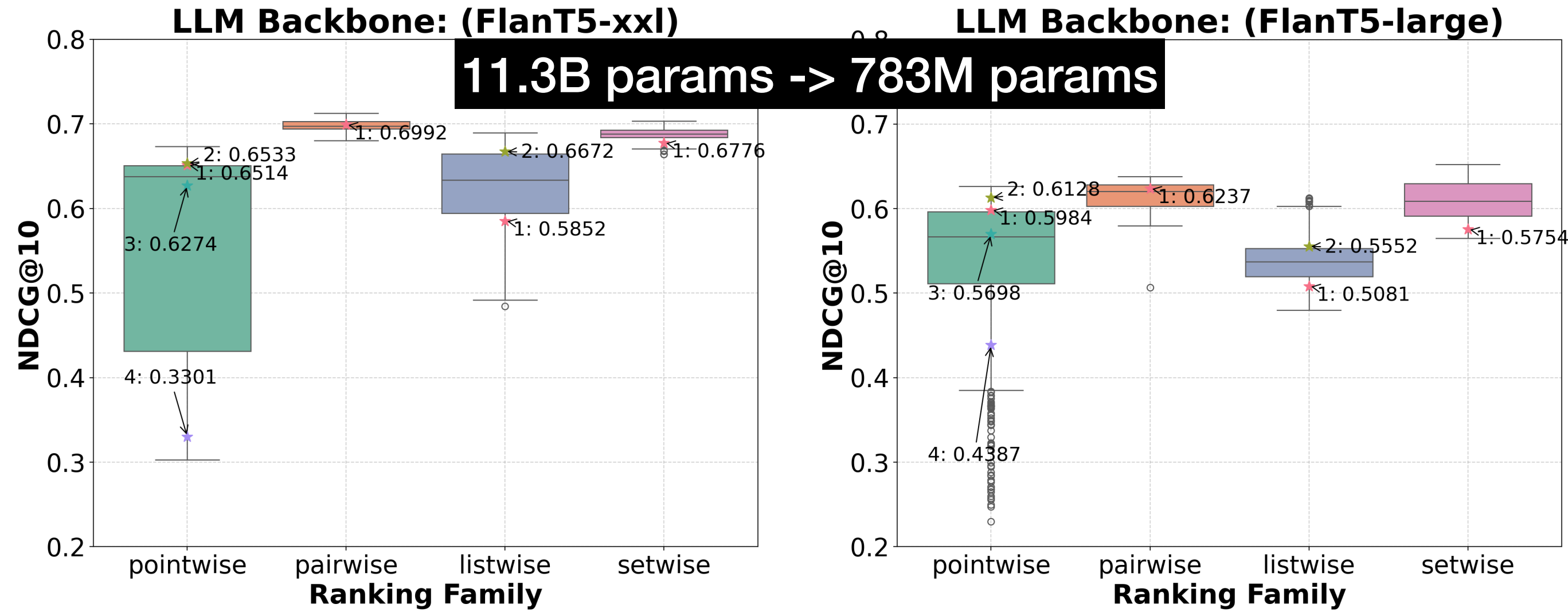


For all approaches, there are better prompts than the one in the original paper



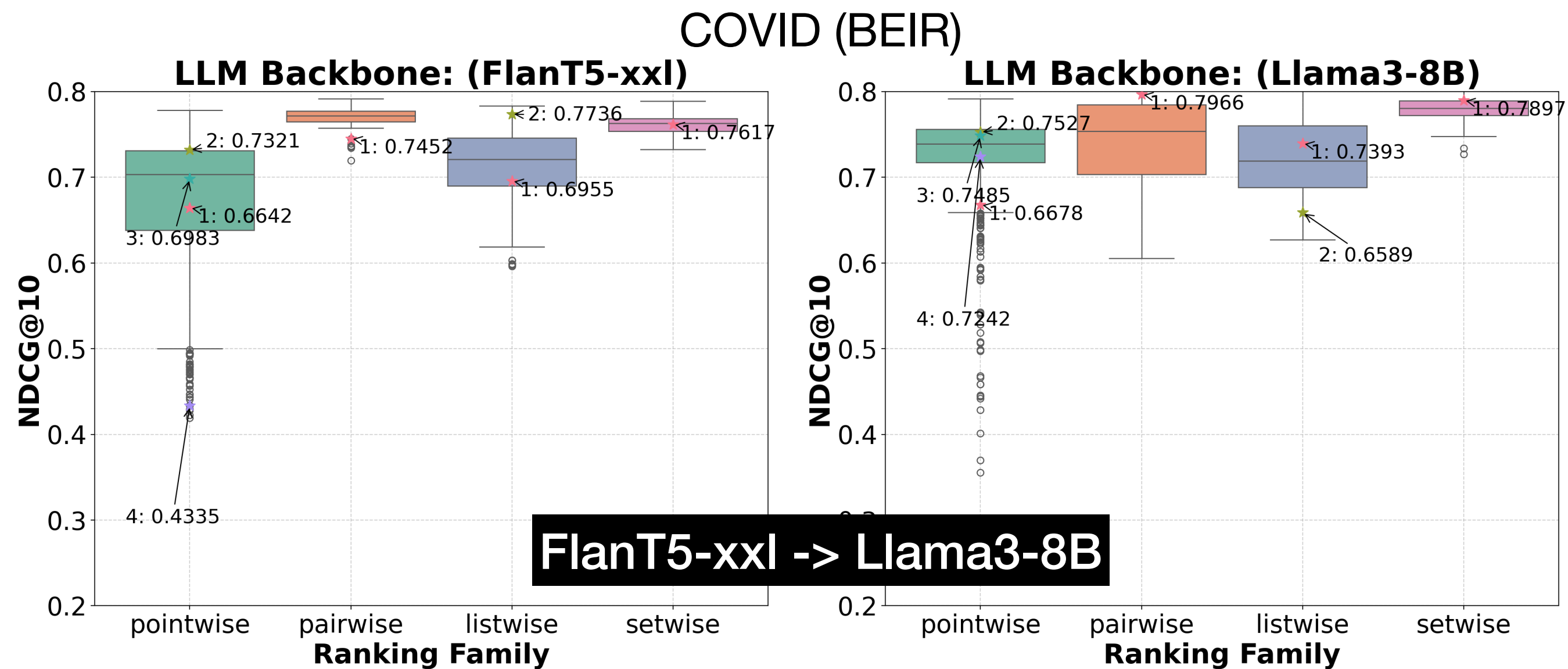
# Impact of Prompt Variations on LLM Rankers

TREC DL 2020



But specific behaviour depends also on

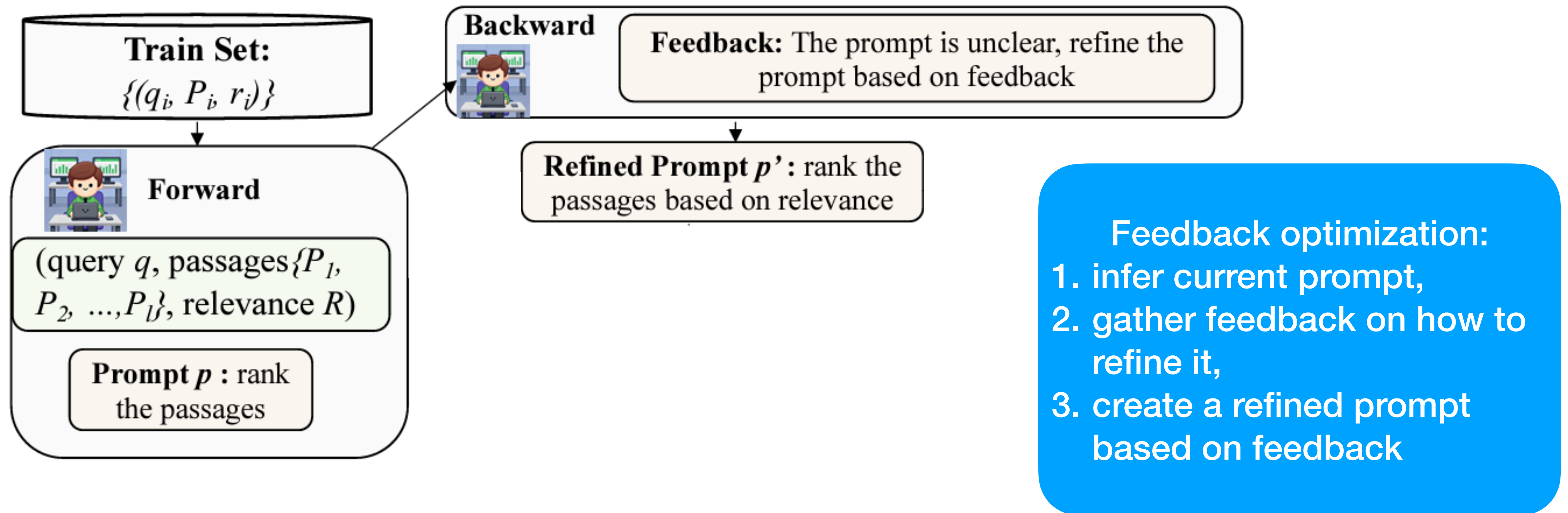
- Size of LLM
- The actual LLM architecture/training



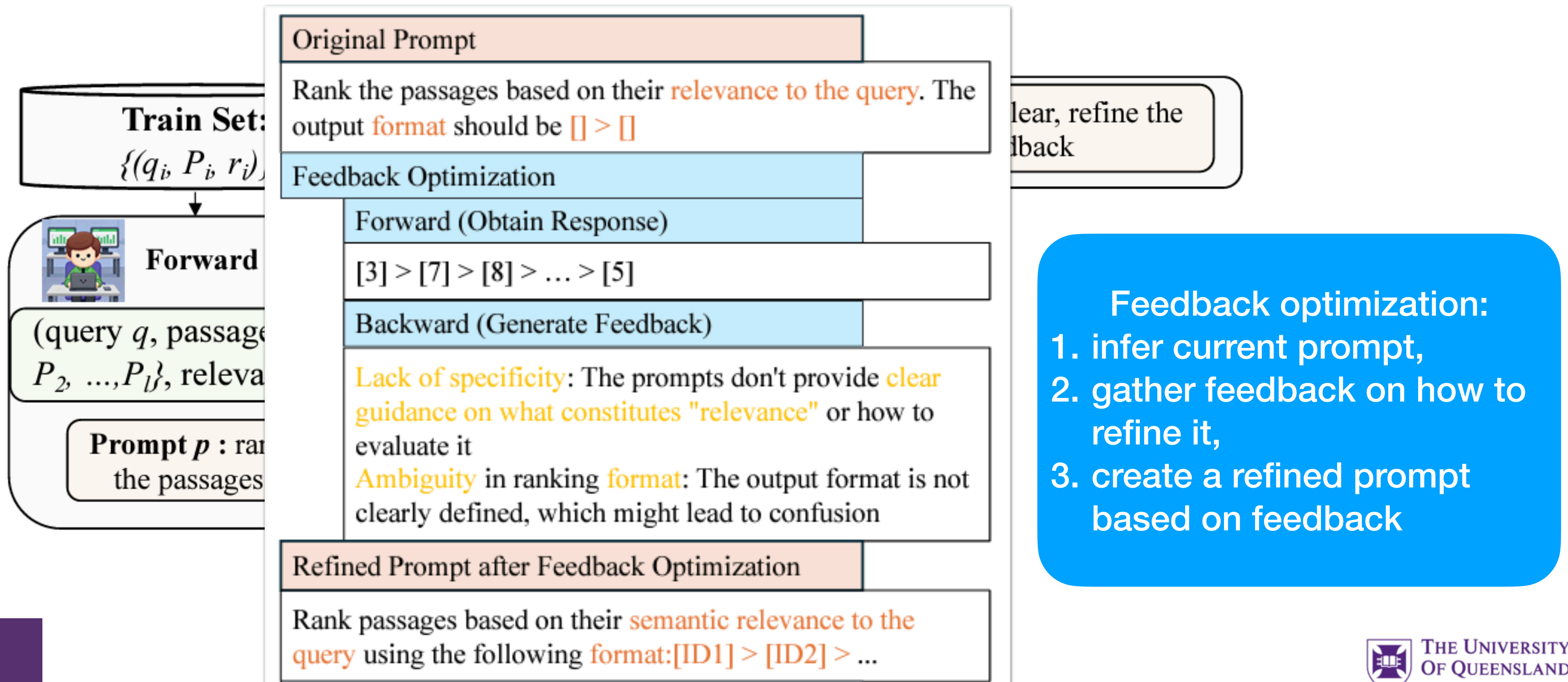
# From Manual Prompting to Automatic Prompt Engineering

- LLM rankers **effectiveness varies** across prompt **wording**
- Can we **automatically identify** create effective **prompt**?
- Direct application of current LLM automatic prompt engineering challenging cause ranking task more complex than classification task:
  1. Requires comprehension of query and of relevance of passages
  2. input-output demonstrations for relevance ranking are more complex than those for language modeling.
- Input: query, passage. Output: relevance — which is not unique to query-passage pair (others may have same relevance level)

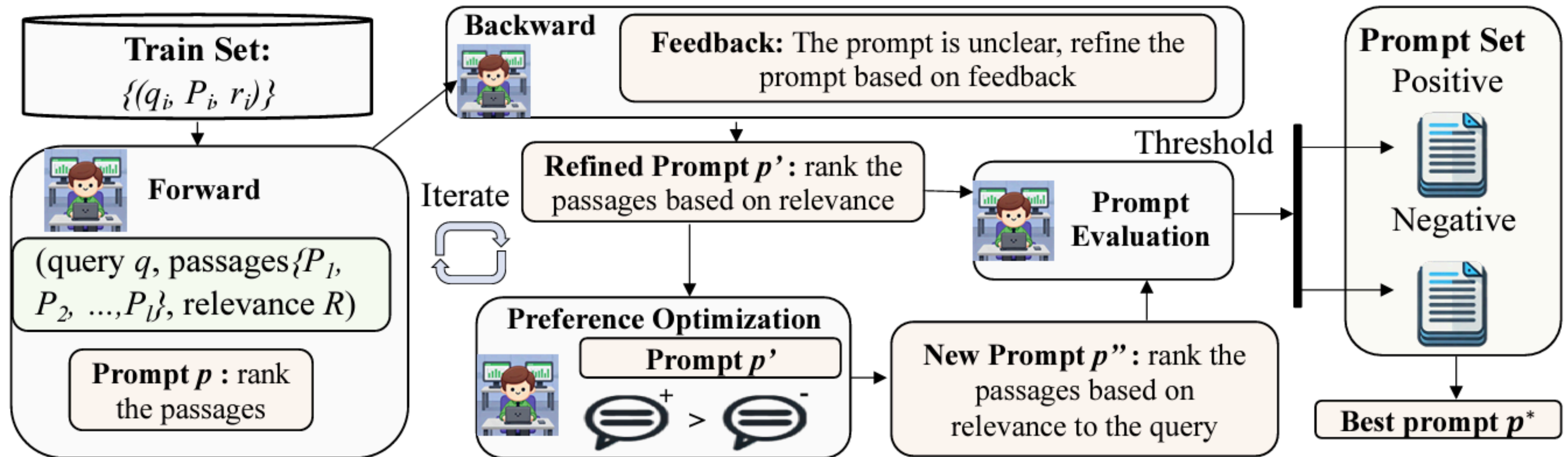
# From Manual Prompting to Automatic Prompt Engineering: APEER



# From Manual Prompting to Automatic Prompt Engineering: APEER



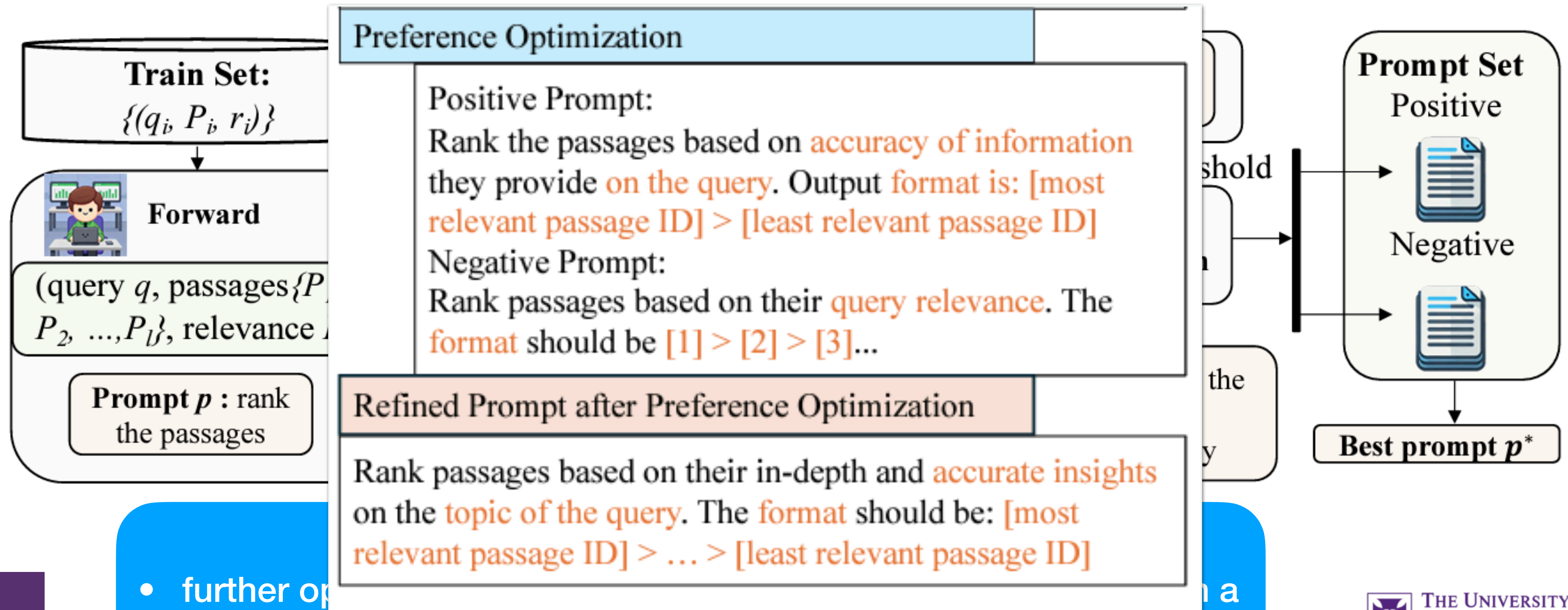
# From Manual Prompting to Automatic Prompt Engineering: APEER



## Preference optimization:

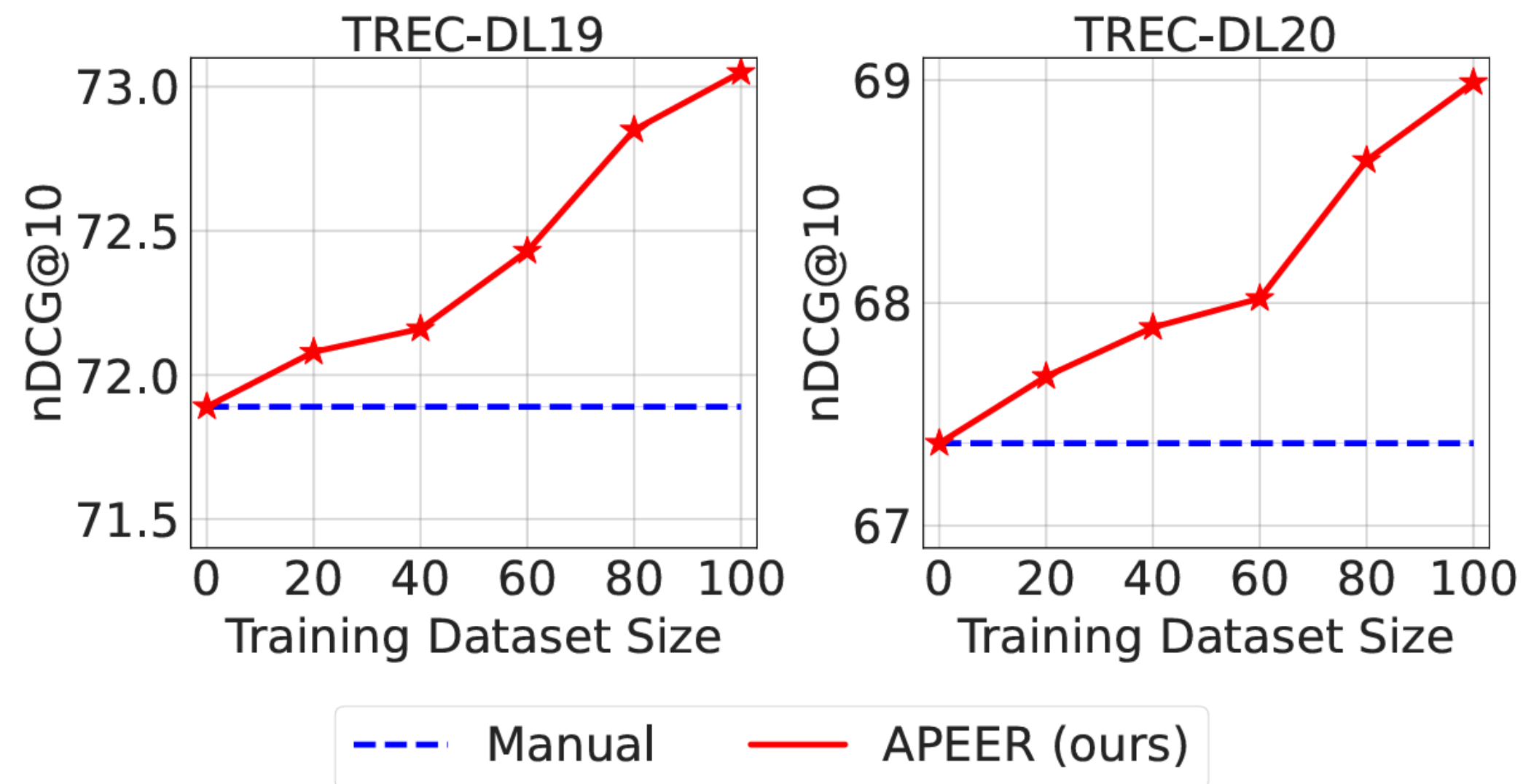
- further optimize refined prompt by learning preferences through a set of positive and negative prompt demonstrations

# From Manual Prompting to Automatic Prompt Engineering: APEER



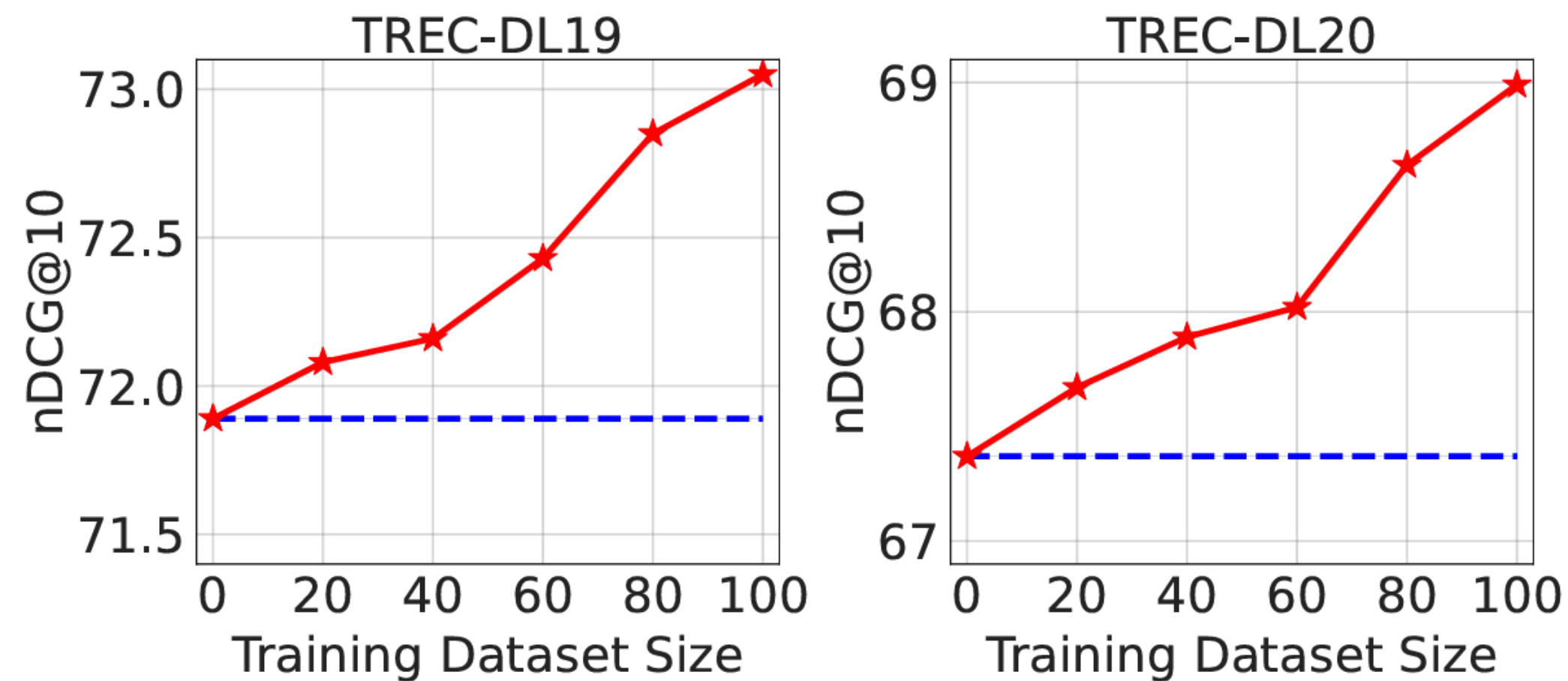
- further optimization of the prompt set based on a set of positive and negative prompt demonstrations

# From Manual Prompting to Automatic Prompt Engineering: APEER



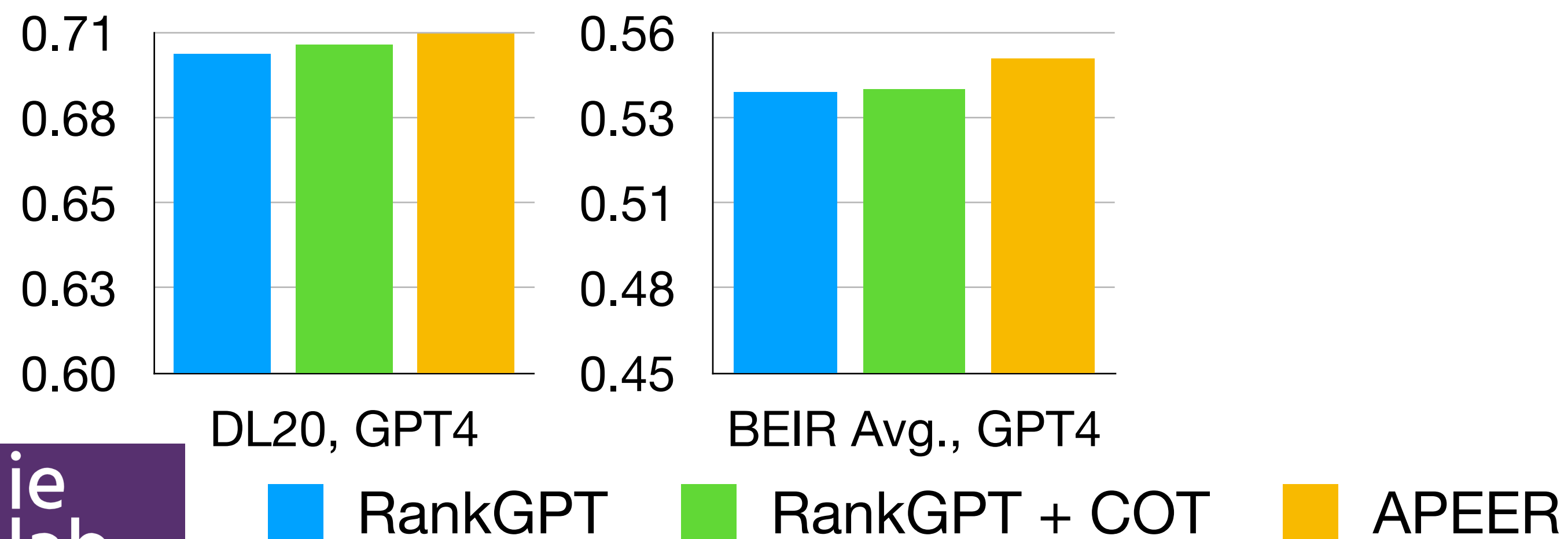
- **More training data** consistently creates **better and better prompts**
- At the cost of additional computational costs for training

# From Manual Prompting to Automatic Prompt Engineering: APEER

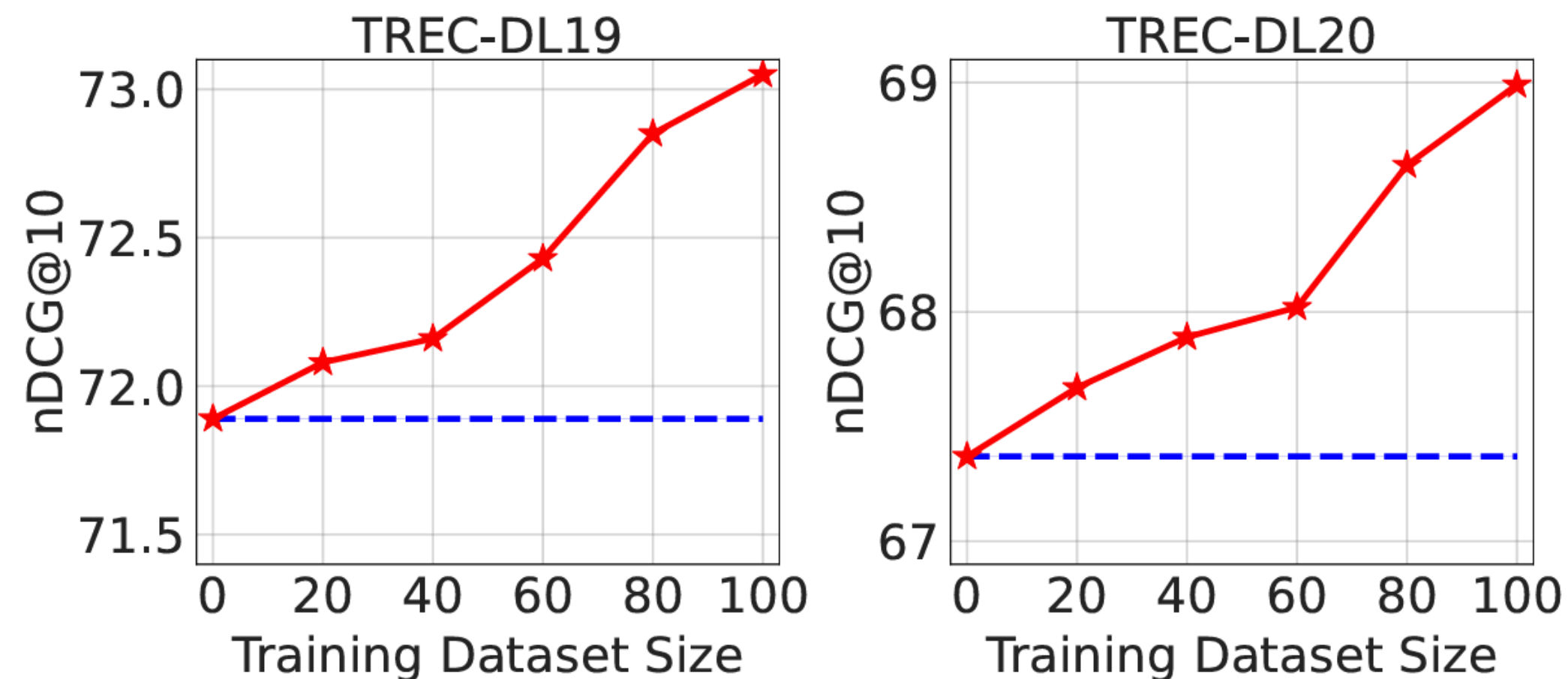


--- Manual    — APEER (ours)

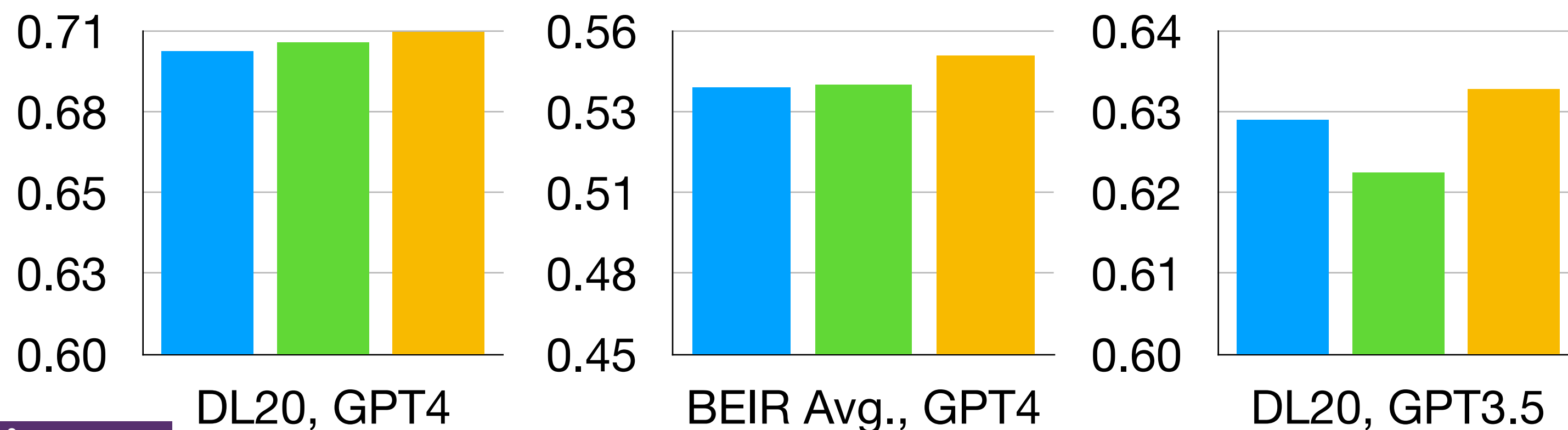
- **More training data** consistently creates **better and better prompts**
- At the cost of additional computational costs for training
- Prompts created with APEER trained on MS MARCO also **transfer to out of domain** datasets (BEIR)



# From Manual Prompting to Automatic Prompt Engineering: APEER



--- Manual    — APEER (ours)

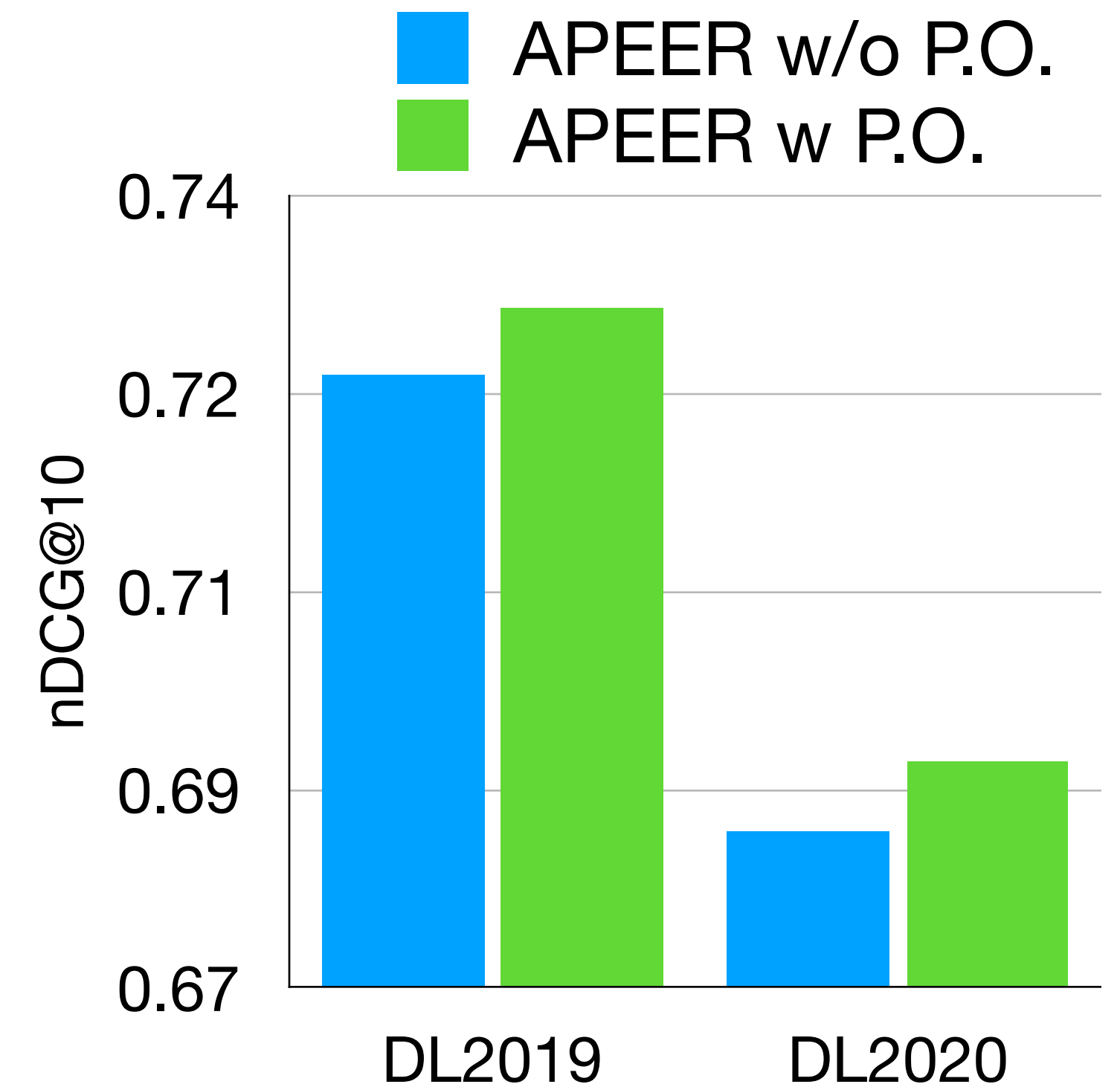


■ RankGPT    ■ RankGPT + COT    ■ APEER

- More training data consistently creates better and better prompts
- At the cost of additional computational costs for training
- Prompts created with APEER trained on MS MARCO also **transfer** to **out of domain** datasets (BEIR)
- Prompts can also **transfer** between **models**, e.g. GPT4->GPT3.5; Qwen2 -> Llama3

# More on Preference Optimisation Phase

- Catalog a collection of potential positive and negative responses within  $H_{\text{pos}}$  and  $H_{\text{neg}}$
- Utilize 2t pairs of  $p_{\text{pos}}, p_{\text{neg}}$  for demonstration
- Generate a new prompt, evaluate its performance on validation dataset relative to initial prompt, insert it in  $H_{\text{pos}}$  and  $H_{\text{neg}}$  accordingly
- Feedback Optimization: acts as local optimizer for current batch
- Preference Optimization: extends local optimization by globally aligning the local optimized prompts towards superior global prompts



**Preference Optimisation beneficial but gains minor (~0.8%)**



THE UNIVERSITY  
OF QUEENSLAND  
AUSTRALIA

CREATE CHANGE

# Part 2: Retrieval Augmented Generation

# Retrieval Augmented Generation

Reinventing search with a new AI-powered Microsoft Bing and Edge, your copilot for the web

The screenshot shows the Microsoft Bing chat interface. At the top, there are navigation links for "Microsoft Bing", "SEARCH", and "CHAT". A blue chat bubble contains the user's query: "I am planning a trip for our anniversary in September. What are some places we can go that are within a 3 hour flight from London Heathrow?". Below this, a white response bubble contains the AI's reply: "Congratulations on your anniversary! 🎉 There are many places you can go that are within a 3 hour flight from London Heathrow. Here are some suggestions based on your preferences and the best destinations in Europe in September 4 5 6 :". The response lists three suggestions: Malaga in Spain, Annecy in France, and Florence in Italy. At the bottom, there is a search bar with the placeholder text "Ask me anything..." and a blue microphone icon.

Microsoft Bing SEARCH CHAT

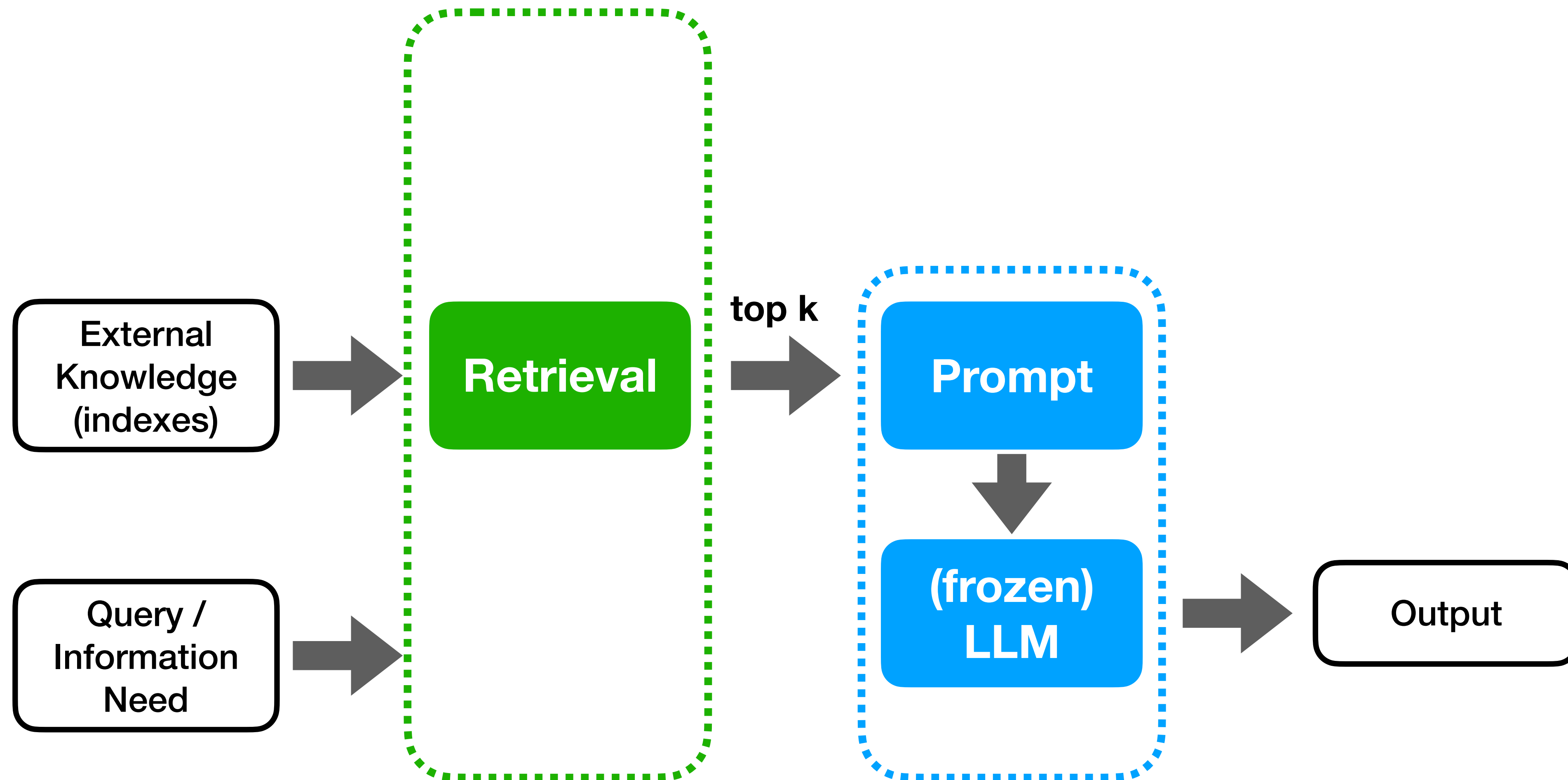
I am planning a trip for our anniversary in September. What are some places we can go that are within a 3 hour flight from London Heathrow?

Congratulations on your anniversary! 🎉 There are many places you can go that are within a 3 hour flight from London Heathrow. Here are some suggestions based on your preferences and the best destinations in Europe in September 4 5 6 :

- If you like **beaches** and **sunshine**, you can fly to **Malaga** in Spain <sup>1</sup>, which has a pleasant climate, a lively culture, and a beautiful coastline. You can enjoy the sandy beaches, the historic monuments, and the delicious tapas. Malaga is also a great base to explore other nearby attractions, such as Granada, Seville, and Cordoba.
- If you like **mountains** and **lakes**, you can fly to **Annecy** in France <sup>4</sup>, which is a charming alpine town near the Swiss border. Annecy is known as the "Venice of France" because of its canals and bridges. You can admire the stunning scenery, visit the medieval castle, and try some local cheese and wine. Annecy is also close to Geneva, Chamonix, and Lyon, if you want to see more of the region.
- If you like **art** and **history**, you can fly to **Florence** in Italy <sup>6</sup>, which is the birthplace of the Renaissance and a UNESCO World Heritage Site. Florence is a treasure trove of artistic and architectural masterpieces, such as the Duomo, the Uffizi Gallery, and the Ponte Vecchio. You can also explore the Tuscan countryside, taste the famous gelato, and shop for leather goods.

Ask me anything...

# The RAG Pipeline

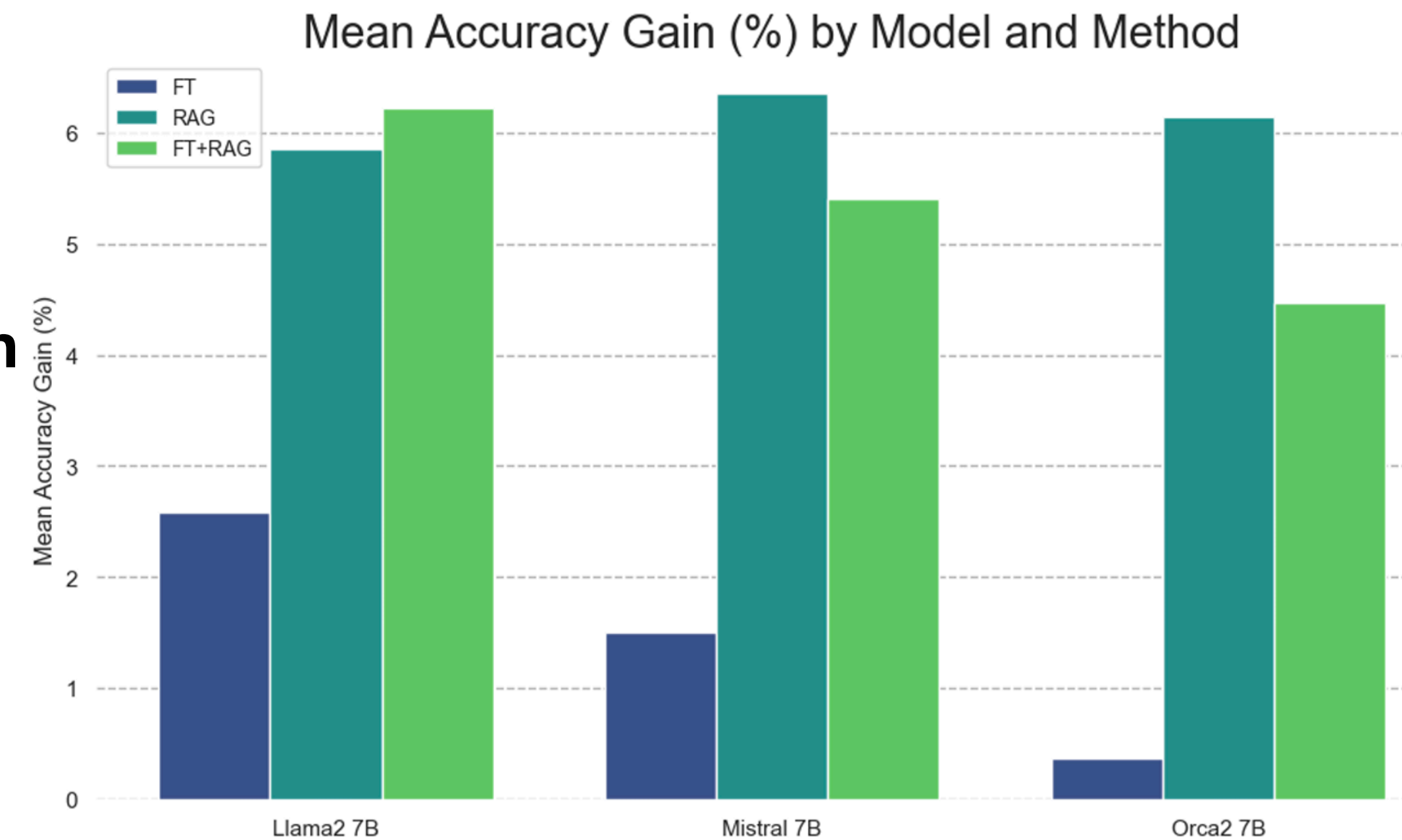


# Why RAG?

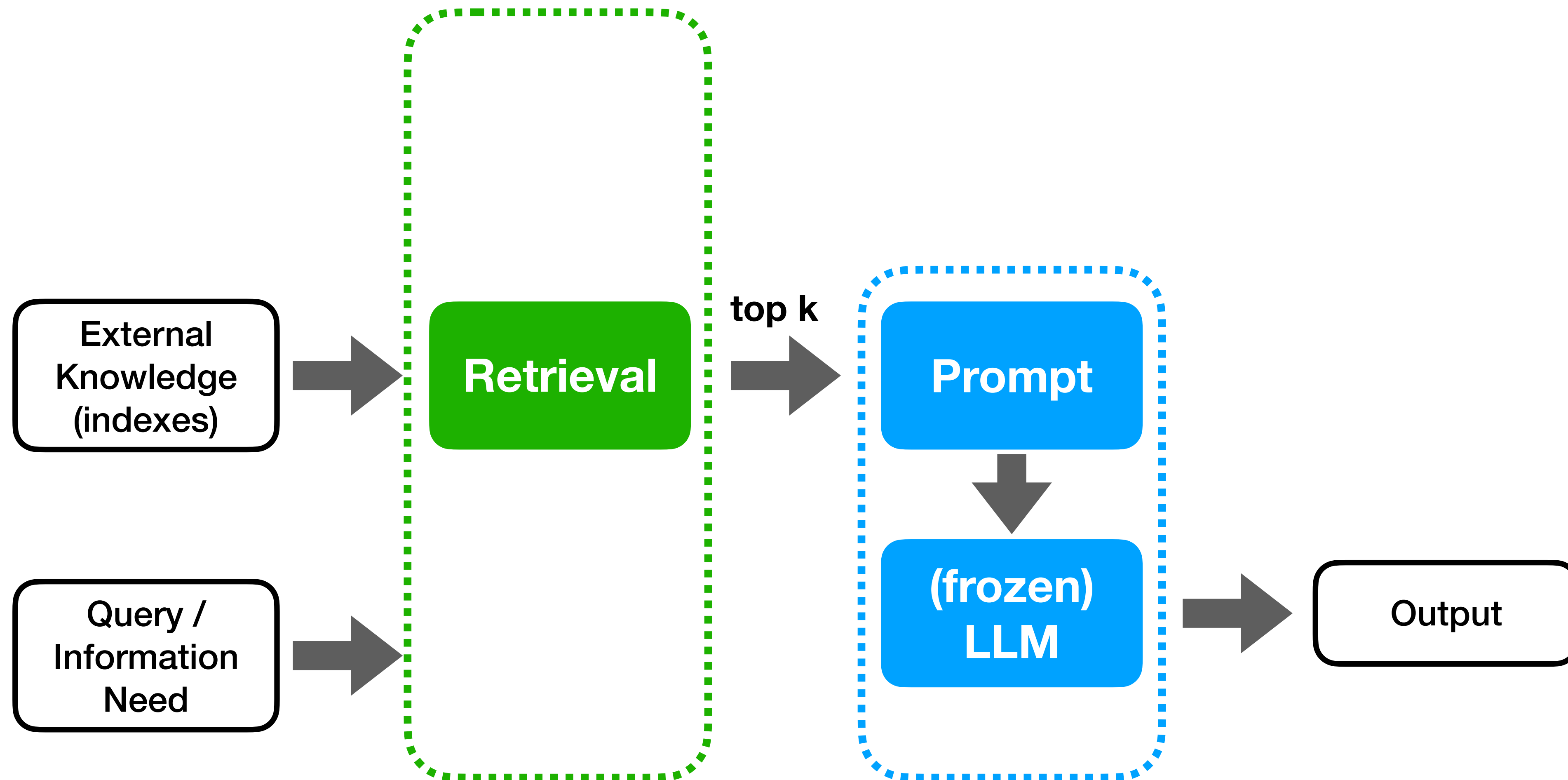
- **LLMs encapsulate** a vast amount of factual **information** within their **pre-trained weights**
- This **knowledge** is inherently **limited**, relying heavily on the characteristics of the training data
- Solution: use **external datasets** to incorporate **new information** or refine the capabilities of LLMs

# Why RAG?

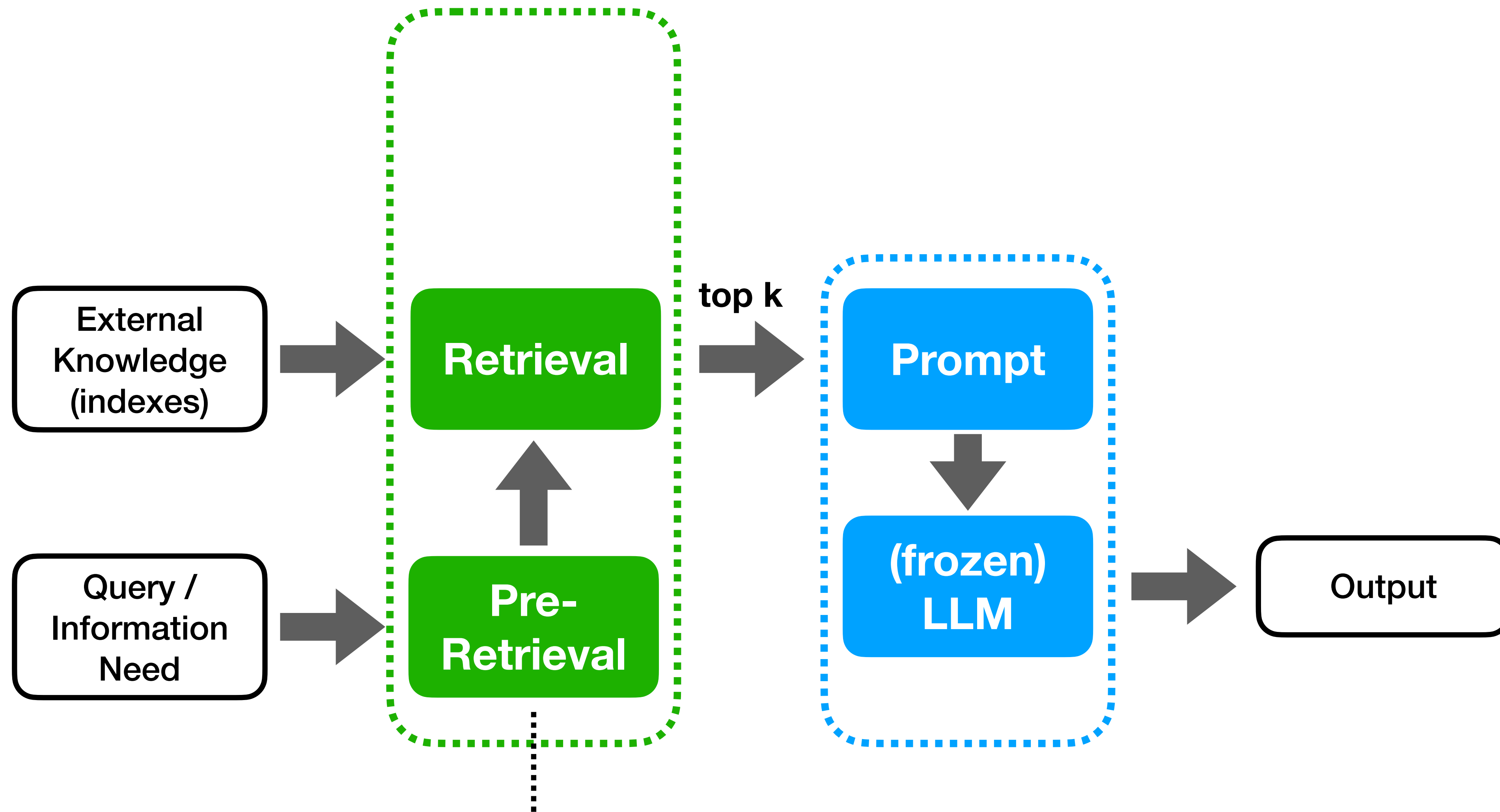
- **LLMs encapsulate** a vast amount of factual **information** within their **pre-trained weights**
- This **knowledge** is inherently **limited**, relying heavily on the characteristics of the training data
- Solution: use **external datasets** to incorporate **new information** or refine the capabilities of LLMs
- Two directions to doing this:
  1. **unsupervised fine-tuning**
  2. retrieval-augmented generation (**RAG**)
- Ovadia et al.: unsupervised fine-tuning offers some improvement, but **RAG consistently outperforms it**, both for existing knowledge encountered during training and entirely new knowledge
  - LLMs struggle to learn new factual information through unsupervised fine-tuning



# The RAG Pipeline

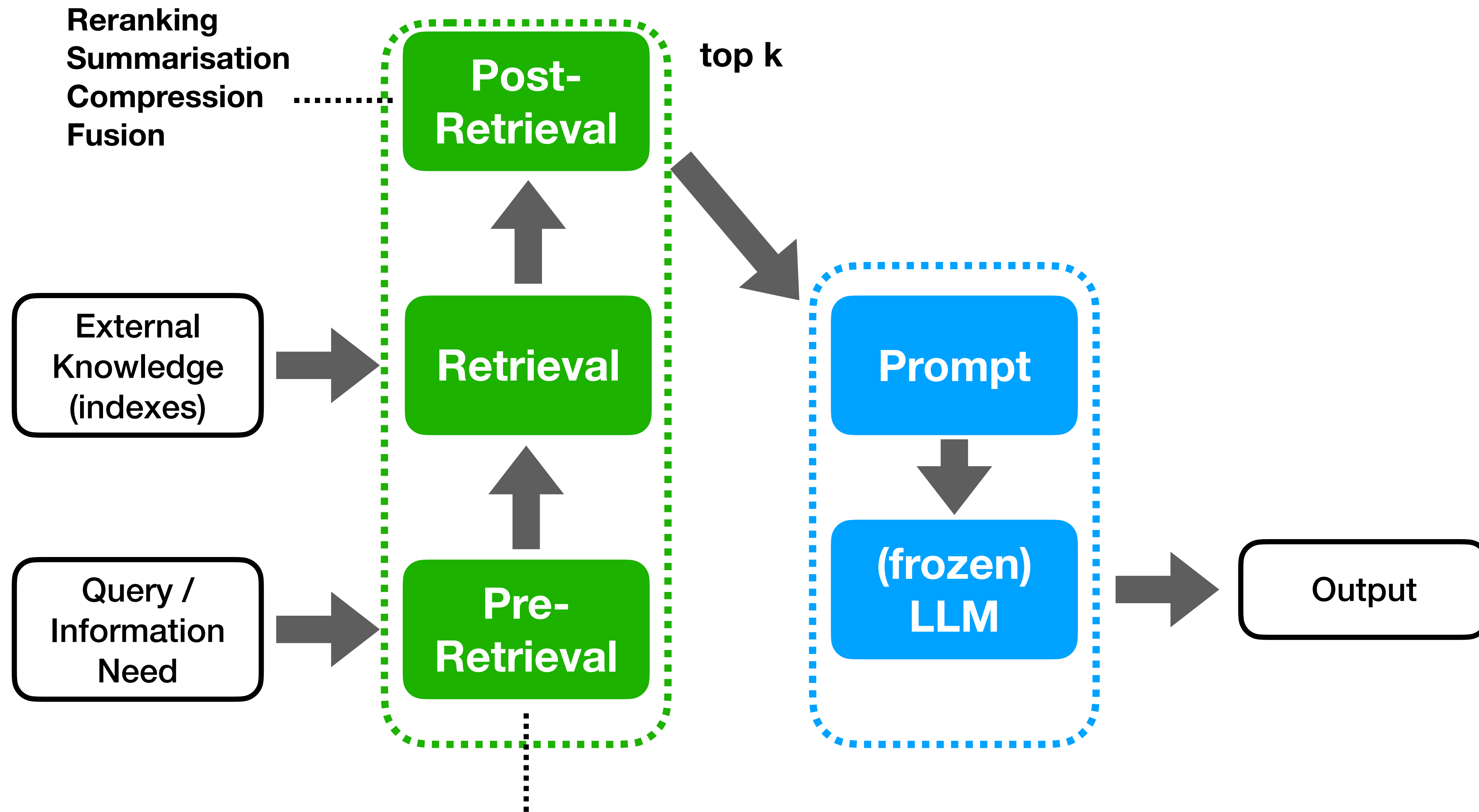


# The RAG Pipeline

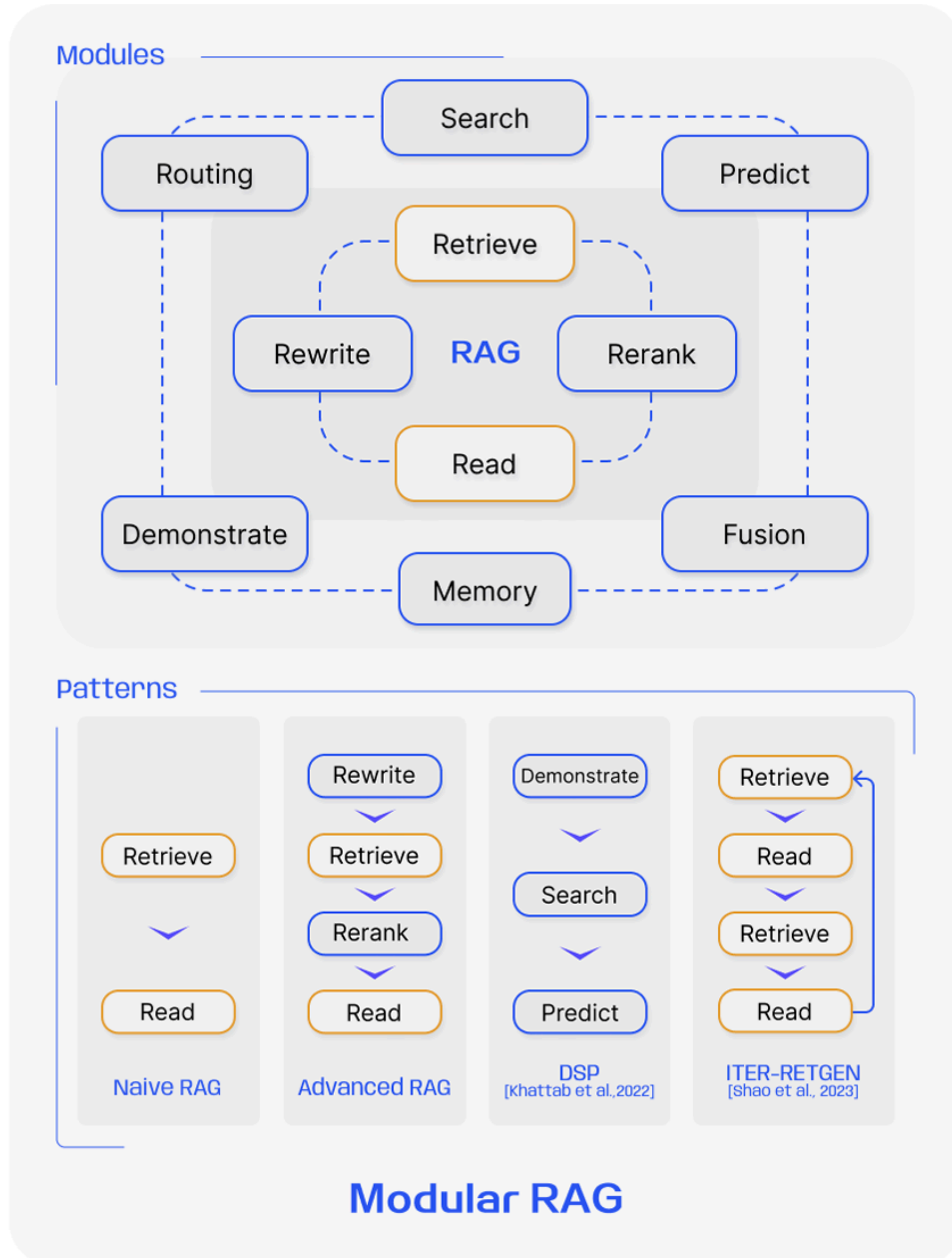


Query Routing  
Query Rewriting  
Query Expansion

# The RAG Pipeline

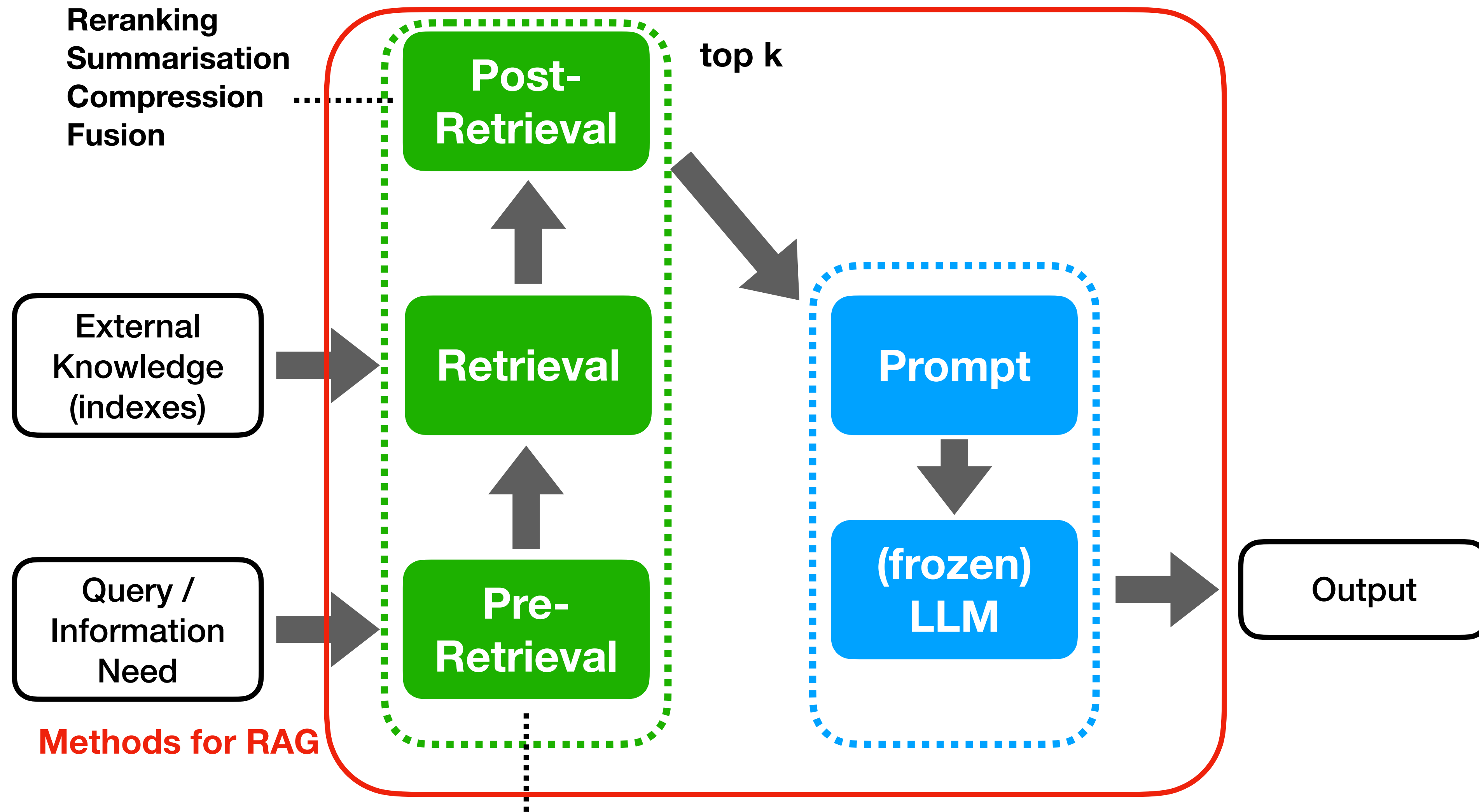


# The Modular RAG pipeline

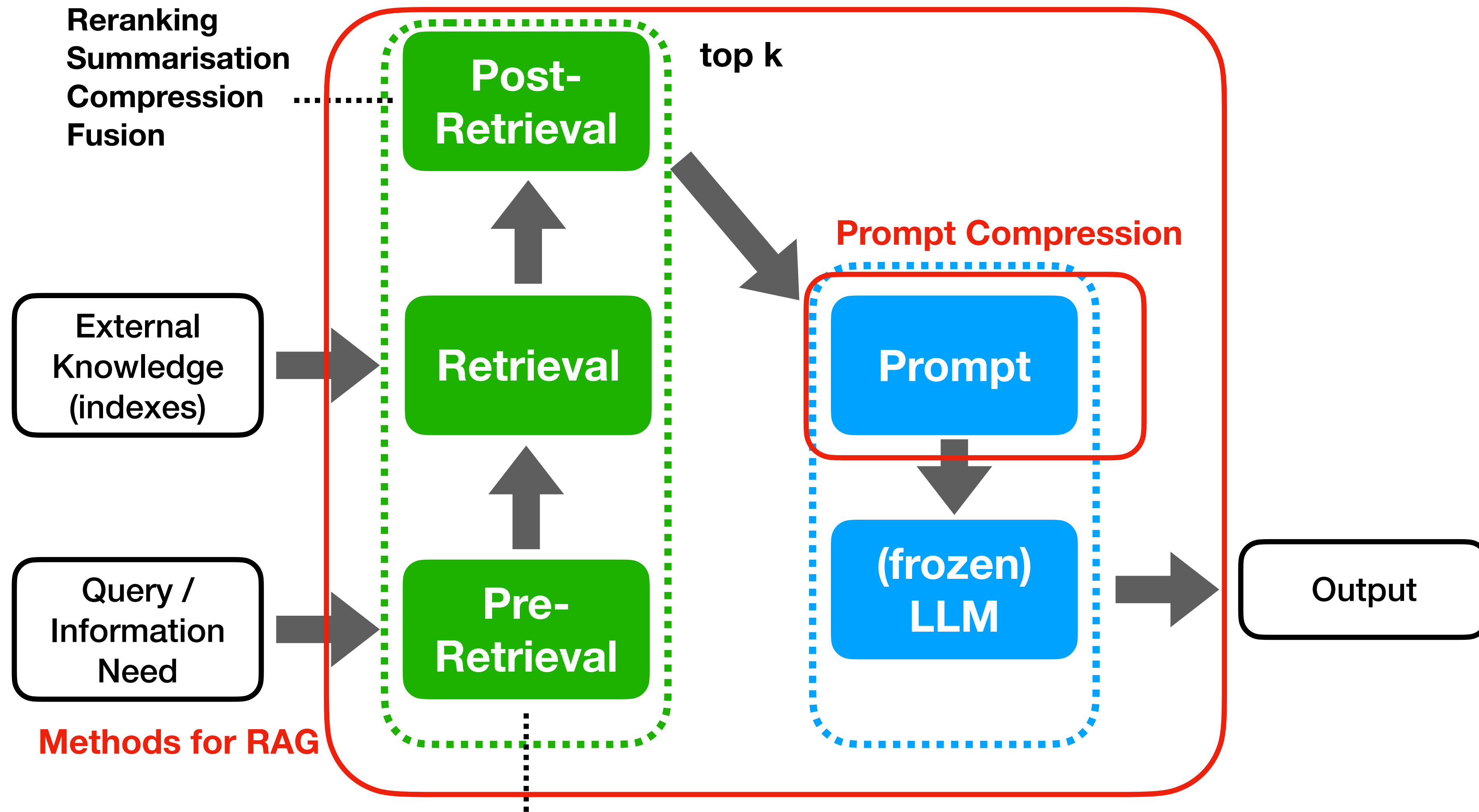


- restructured RAG modules; rearranged RAG pipelines
- Supports **sequential processing** and integrated **end-to-end training** across components
- New/Additional **modules**
- New **Patterns**: e.g. Rewrite-Retrieve-Read, Generate-Read, Recite-Read, Retrieve-Read-Retrieve-Read, Hypothetical Document Embeddings (HyDE), Demonstrate-Search-Predict (DSP)
- **Adaptive RAG**: flexible orchestration of RAG modules and flows, e.g. Self-RAG

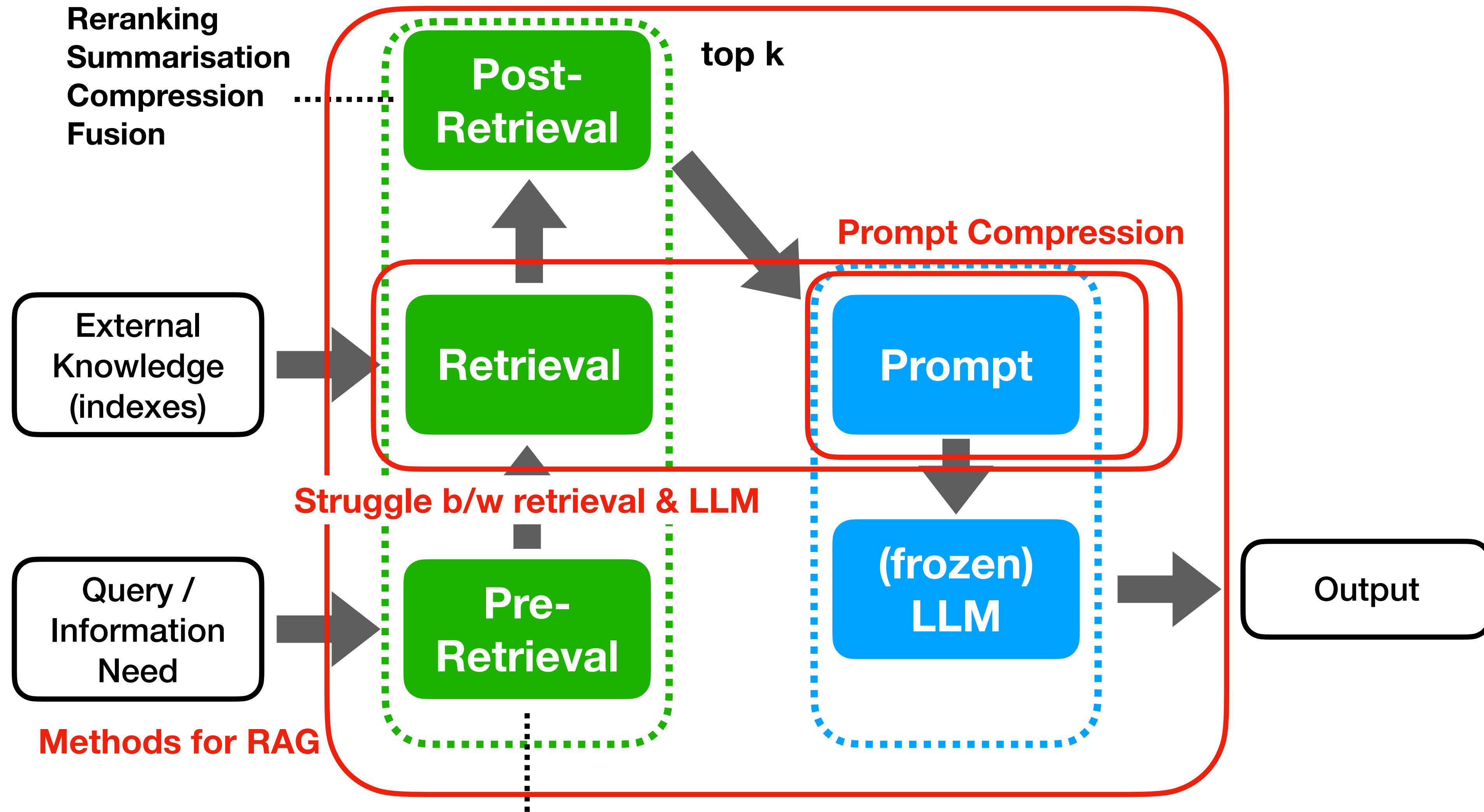
# In the remaining of this part...



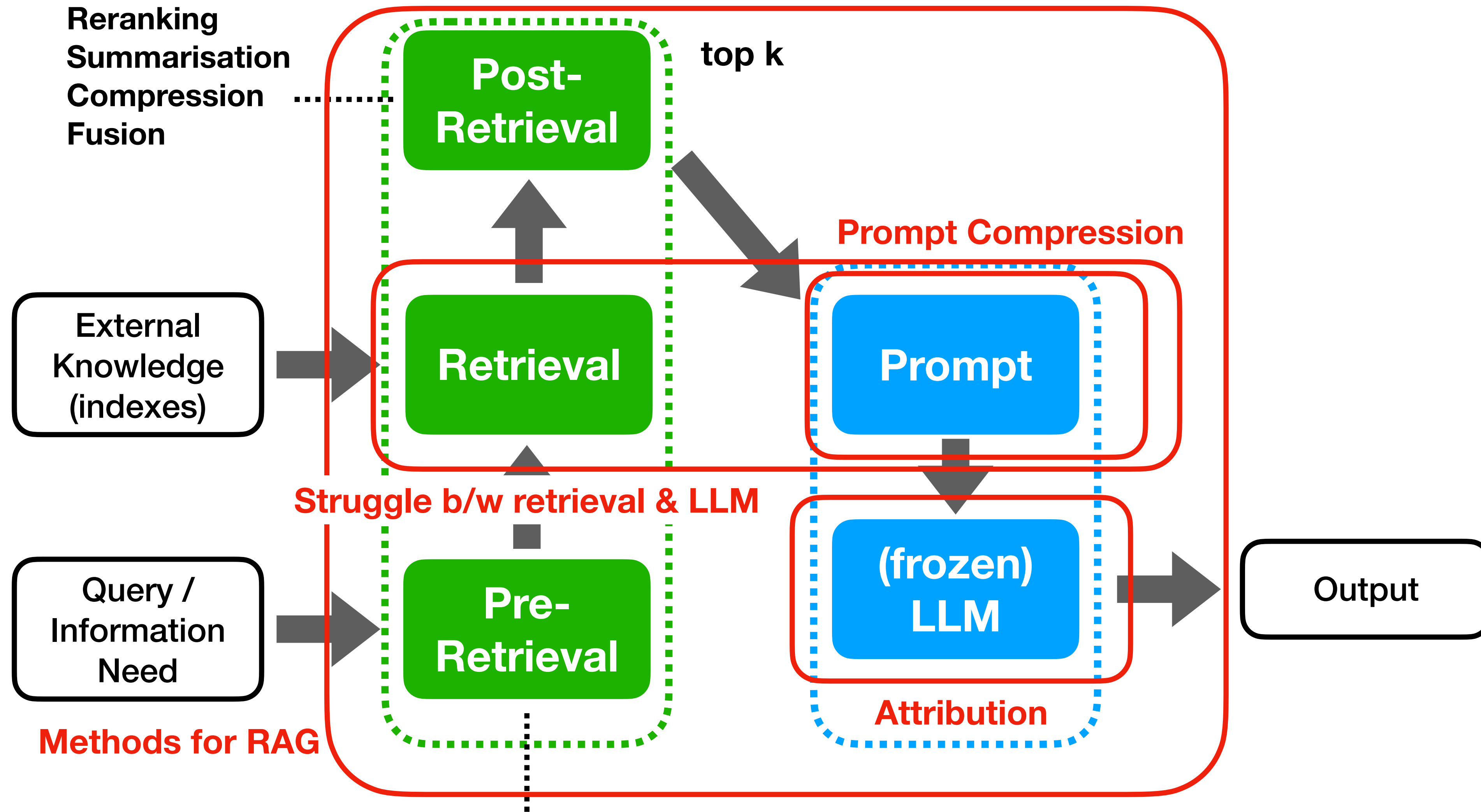
# In the remaining of this part...



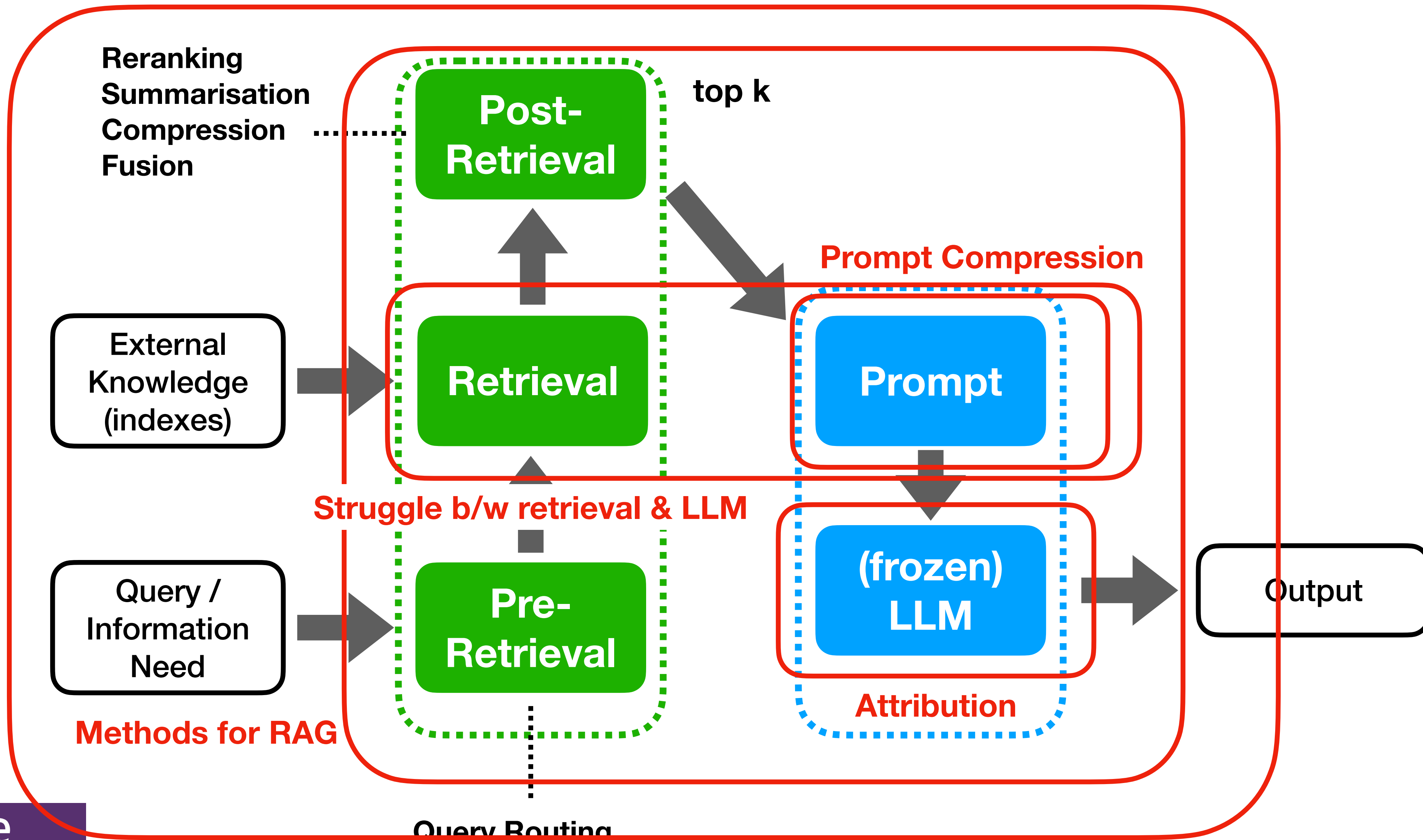
# In the remaining of this part...



# In the remaining of this part...



# In the remaining of this part...

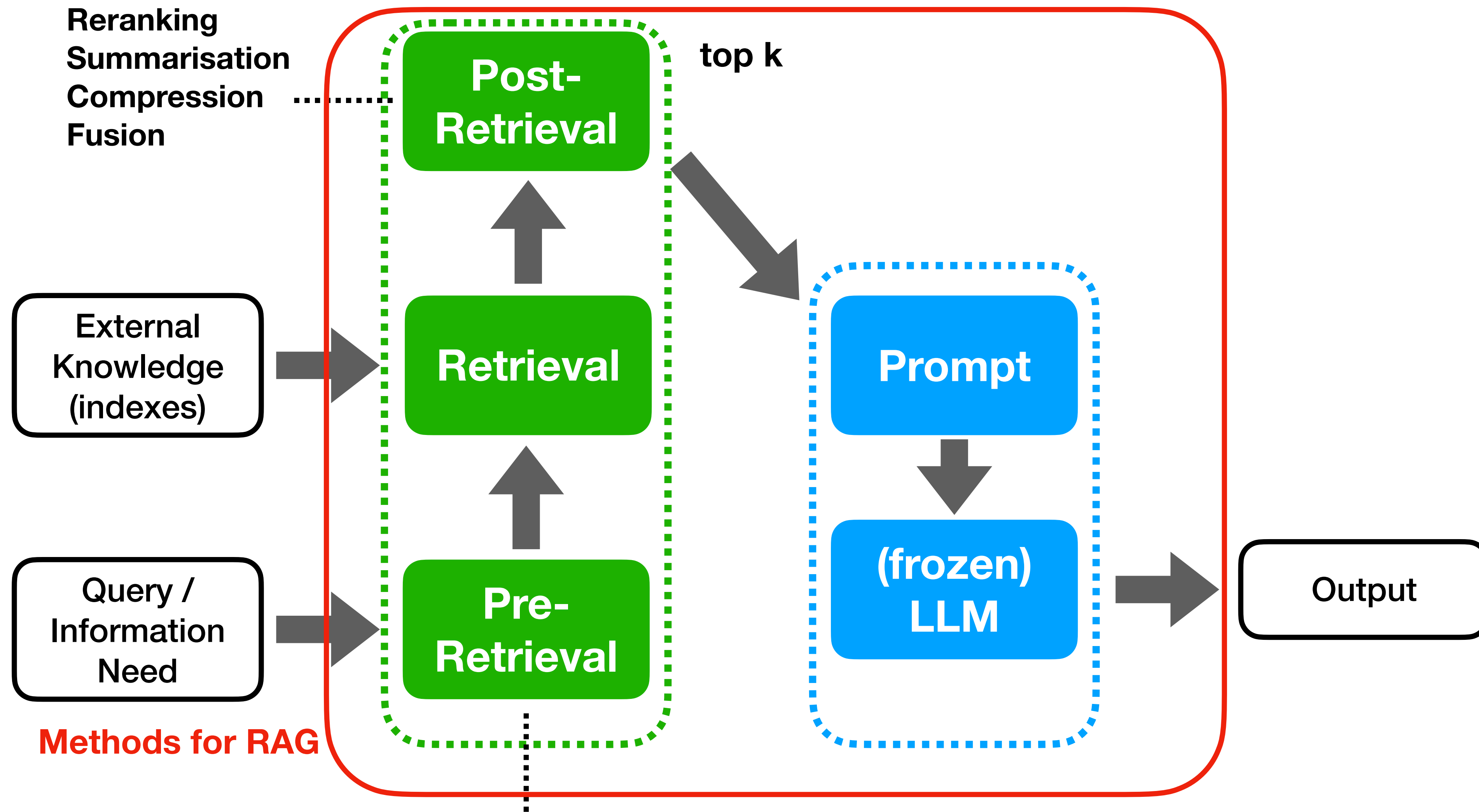


Methods for RAG

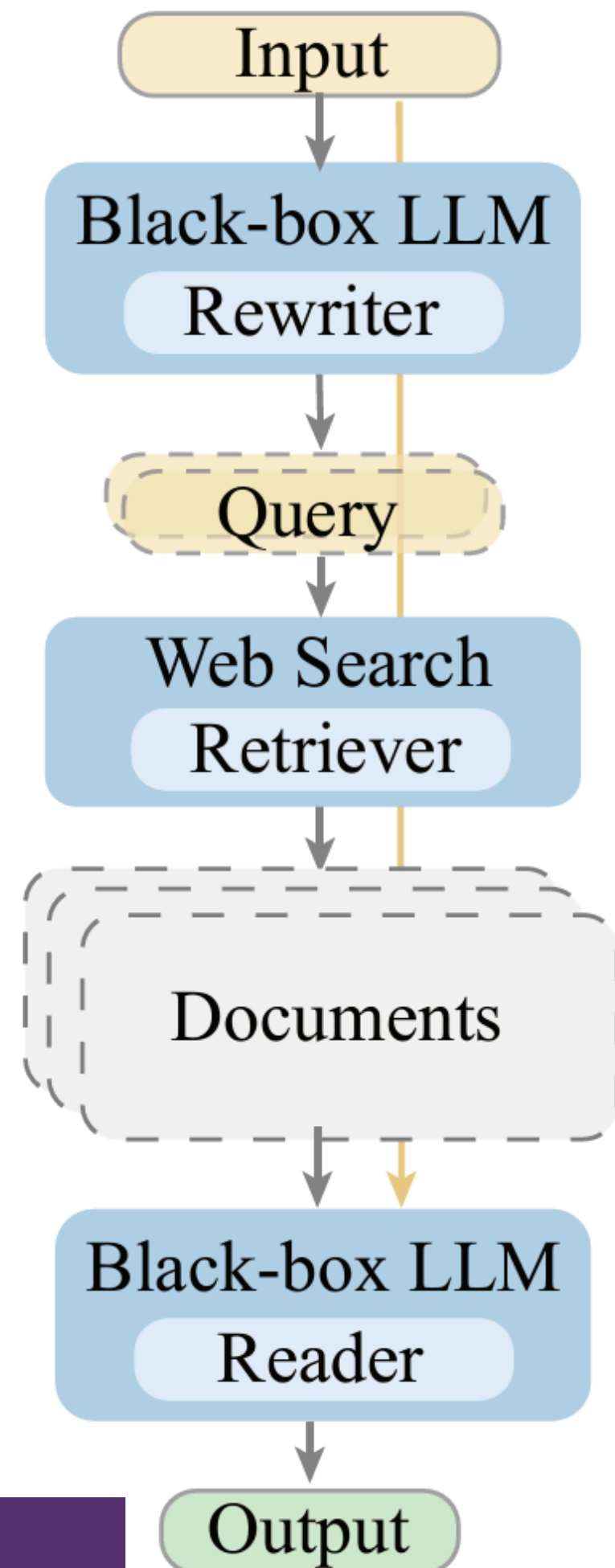
Query Routing  
Query Rewriting  
Query Expansion

Resources & Platforms for RAG

# In the remaining of this part...



# The Rewrite-Retrieve-Read Method

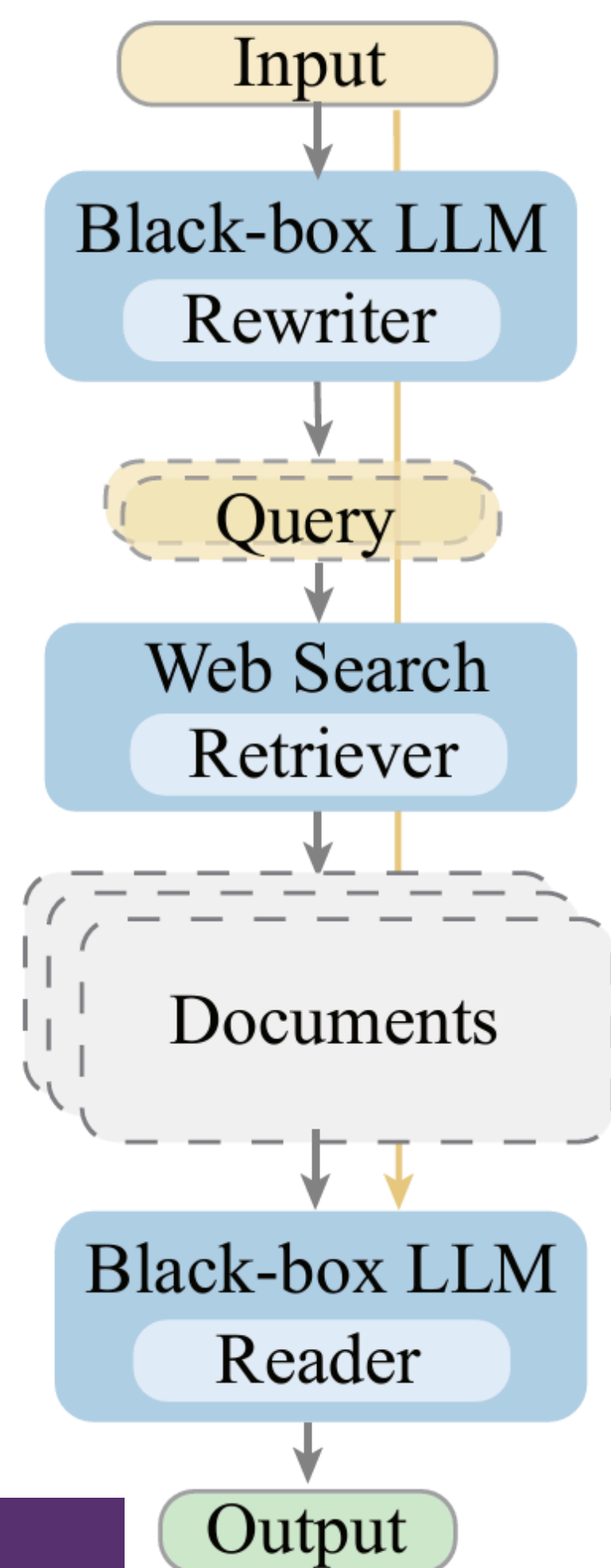


- Queries are often ambiguous and underspecified
- **Rewrite:** black-box LLM is prompted to **re-write query**
- **Retrieve:** search based on the re-written query
- **Read:** provided instruction, query, documents, generate answer

Few-shot (1-3) prompt  
in format [instruction,  
demonstrations, input]

Output can be none,  
one or more queries

# The Rewrite-Retrieve-Read Method



## Open-domain QA

*Think step by step to answer this question, and provide search engine queries for knowledge that you need. Split the queries with ';' and end the queries with '\*'. {demonstration} Question: {x} Answer:*

## Multiple-choice QA

*Provide a better search query for web search engine to answer the given question, end the queries with '\*'. {demonstration} Question: {x} Answer:*

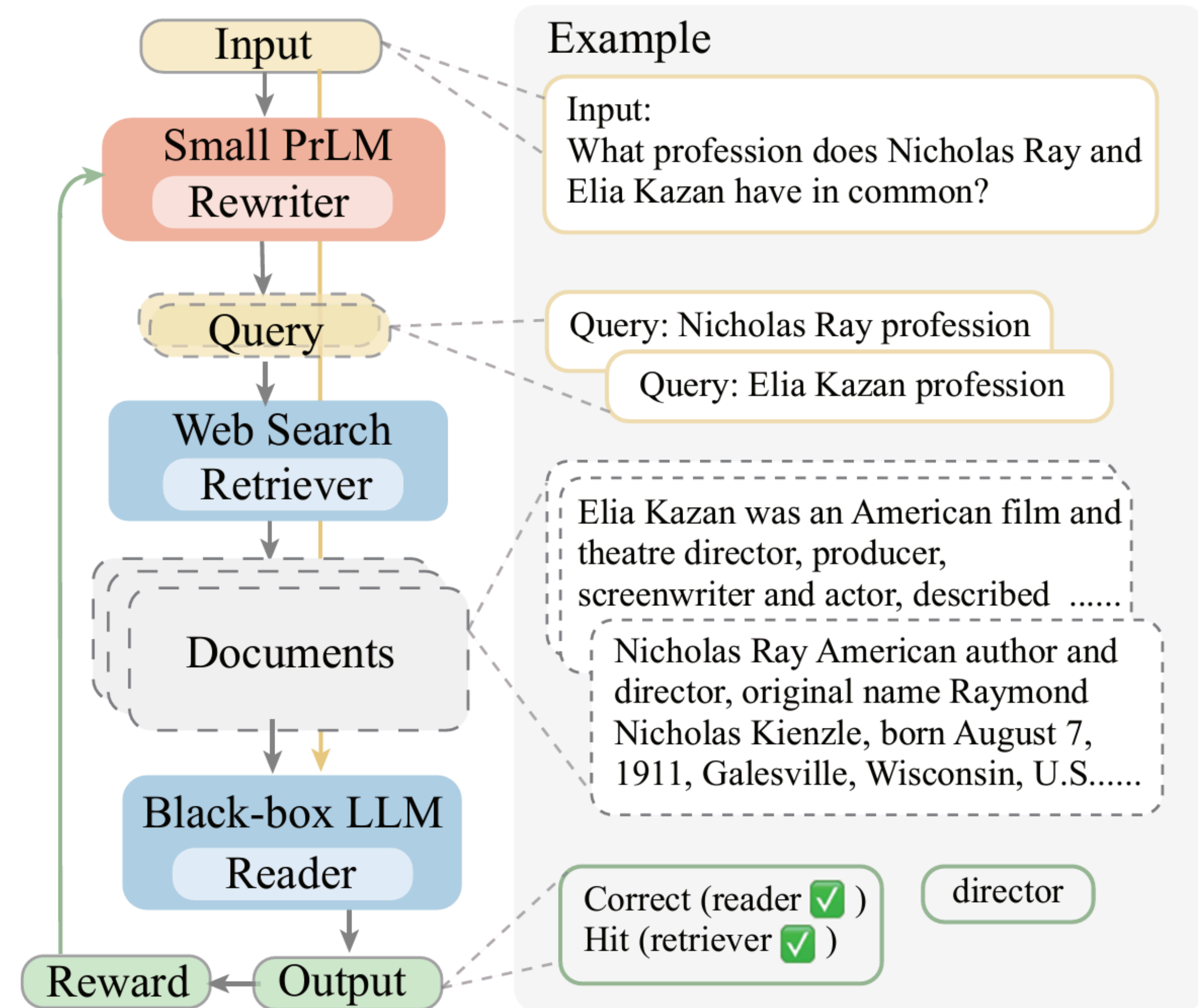
Few-shot (1-3) prompt in format [instruction, demonstrations, input]

Output can be none, one or more queries

# The Rewrite-Retrieve-Read Method

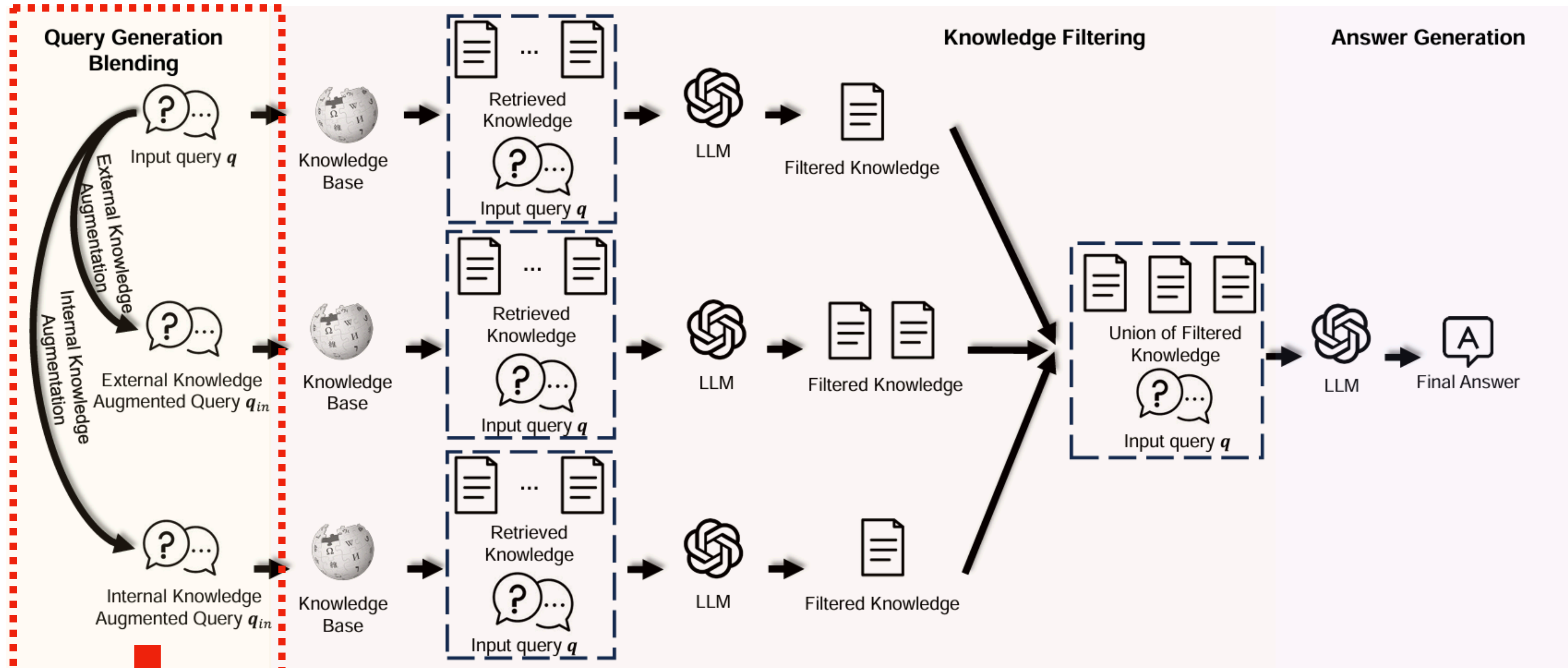
## Trainable Rewriter

- Warm-up: training on pseudo-data
  - Collect rewritten queries  $\underline{x}$ , keep only those that generate correct end-to-end answers: use to form warm-up dataset.
- Fine-tune rewriter on warm-up dataset
- Then continually trained by Reinforcement Learning with PPO
- Reward obtained from end-to-end answers compared to gold + KL regularisation to prevent model doom deviating too far





# BlendFilter

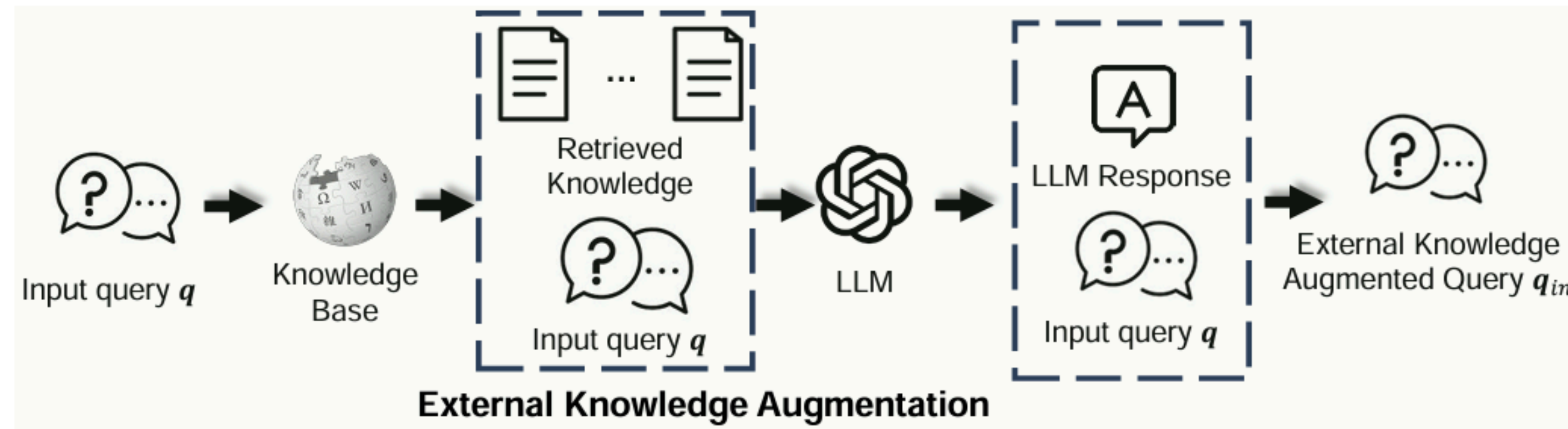
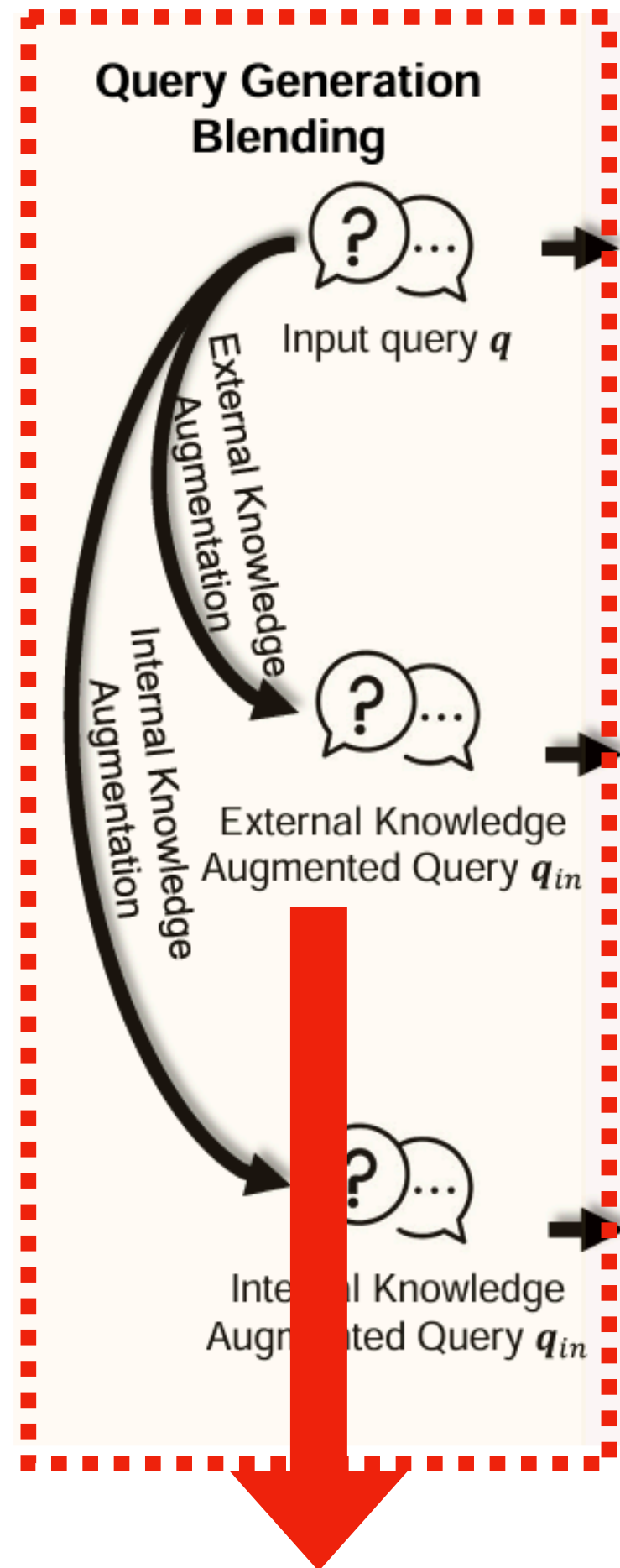


Enhance input queries through different augmentation strategies

Wang, H., Zhao, T. and Gao, J., 2024. BlendFilter: Advancing Retrieval-Augmented Large Language Models via Query Generation Blending and Knowledge Filtering. *arXiv preprint arXiv:2402.11129*.



# BlendFilter



1. Use initial query to retrieve external knowledge
2. Pass this to LLM prompt with few-shot examples and Chain of Thought
3. Generate answer
4. Concatenate answer & original query

Enhance input queries through different augmentation strategies

## Prompt for External Knowledge Augmentation on HotPotQA

Answer questions following the given format.

Knowledge:{Example Knowledge}

Question:Are It Might Get Loud and Mr. Big both Canadian documentaries?

Let's think step by step.

Mr. Big is a 2007 documentary which examines the "Mr. Big" undercover methods used by the Royal Canadian Mounted Police. However, It Might Get Loud is a 2008 American documentary film. So the answer is no.

Knowledge:{Example Knowledge}

Question:Were László Benedek and Leslie H. Martinson both film directors?

Let's think step by step.

László Benedek was a Hungarian-born film director and Leslie H. Martinson was an American film director. So the answer is yes.

Knowledge:{Example Knowledge}

Question:Lucium was confined to be an impure sample of yttrium by an English chemist who became the president of what?

Let's think step by step.

Lucium was confined to be an impure sample of yttrium by William Crookes. William Crookes is Sir William Crookes. Sir William Crookes became the president of the Society for Psychical Research. So the answer is Society for Psychical Research.

Knowledge:{Knowledge}

Question:{question}

Let's think step by step.

# Prompt for External Knowledge Augmentation

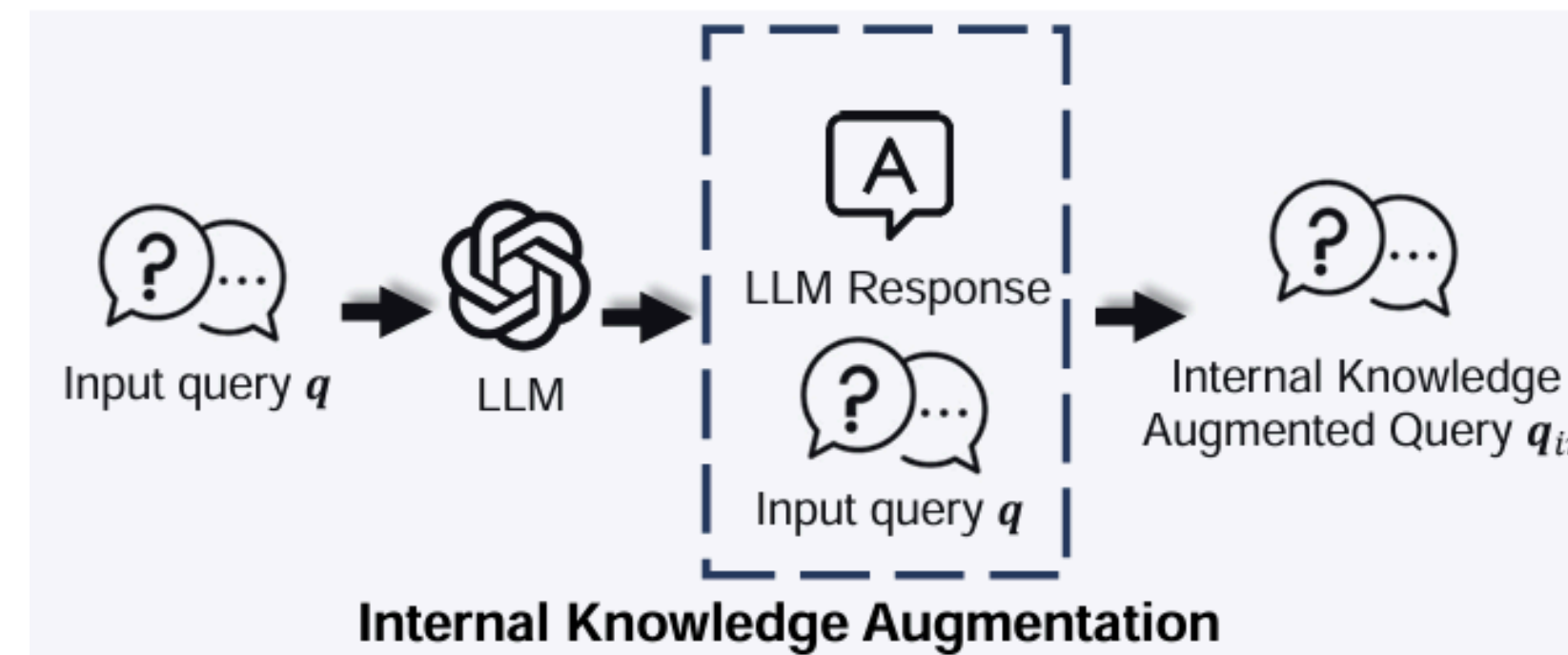
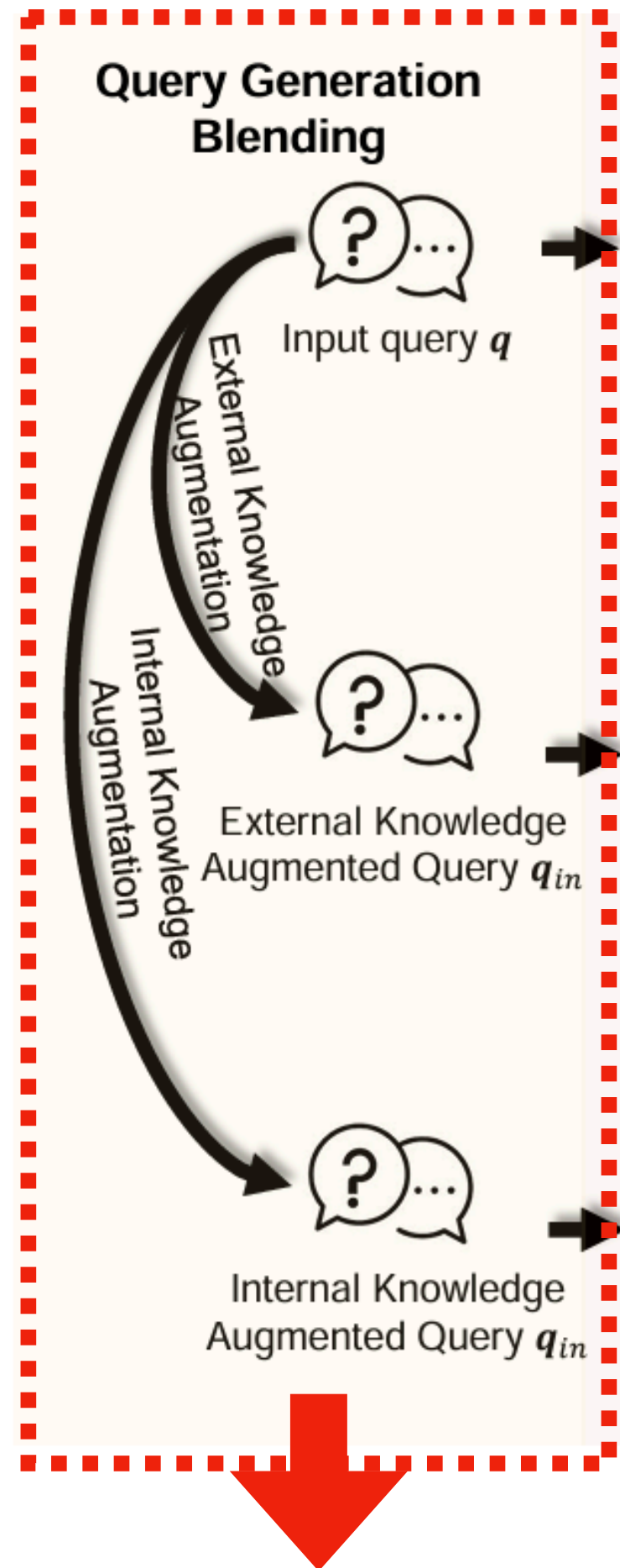
## Few Shot Examples

## Chain-of-Thought



# BlendFilter

1. Prompt LLM to respond to query based on its own knowledge only
2. Concatenate answer & query



### Prompt for Internal Knowledge Augmentation

Please write a passage to answer the question.

Question: {question}

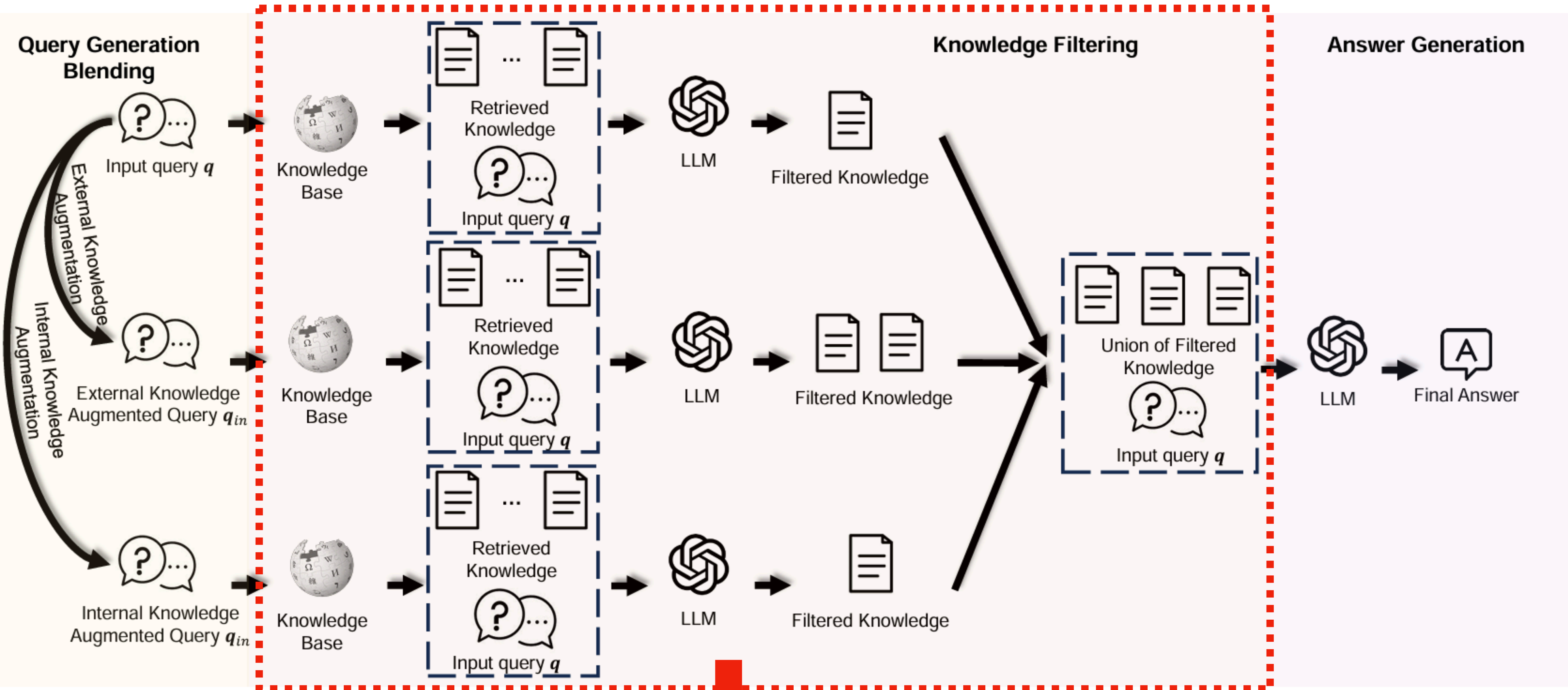
Passage:



Enhance input queries through different augmentation strategies



# BlendFilter



**Eliminate irrelevant retrieved knowledge using prompted LLM**

## Prompt for Knowledge Filtering on HotPotQA and 2Wiki-MultihopQA

What general topic is Question {question} related to?

Answer:The topic is related to

---

— forget your knowledge about {topic}. Please only consider the knowledge below.

knowledge 0 : {Retrieved\_knowledge0}

knowledge 1 : {Retrieved\_knowledge1}

knowledge 2 : {Retrieved\_knowledge2}

knowledge 3 : {Retrieved\_knowledge3}

knowledge 4 : {Retrieved\_knowledge4}

Please check the relevance between {question} and knowledges 0-4 one by one, remove the irrelevant ones and show me the relevant ones. There may be multiple relevant ones. Please take a deep breath and do it step by step.

---

— Please check the relevance between the given question and knowledges 0-4 one by one based on the given context. ONLY output the relevant knowledge ids (0-4). There may be multiple relevant ones.

Context:{LLM\_Last\_Generated\_Context}

Question:{question}

knowledge 0 : {Retrieved\_knowledge0}

knowledge 1 : {Retrieved\_knowledge1}

knowledge 2 : {Retrieved\_knowledge2}

knowledge 3 : {Retrieved\_knowledge3}

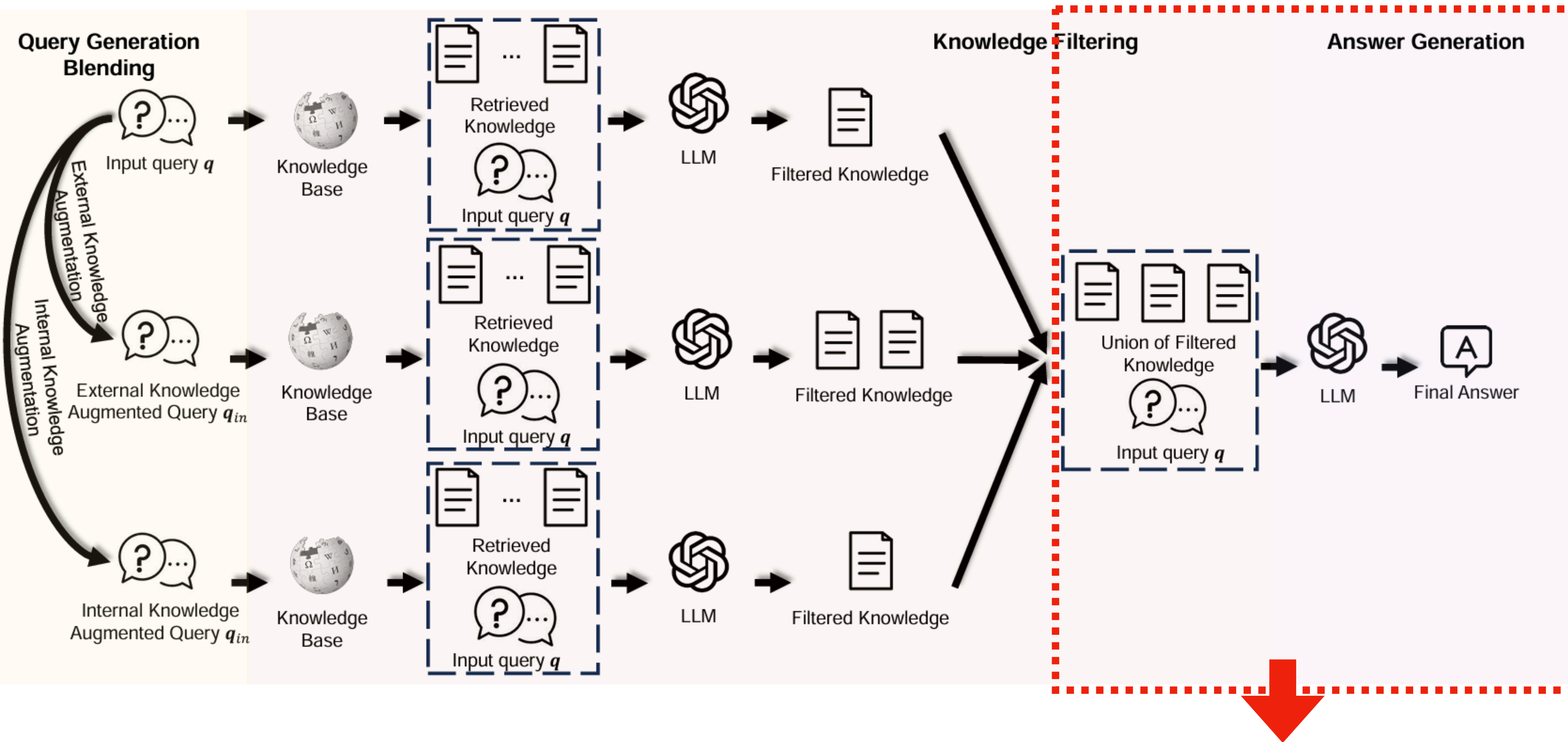
knowledge 4 : {Retrieved\_knowledge4}

Answer:

# Knowledge Filtering Prompt

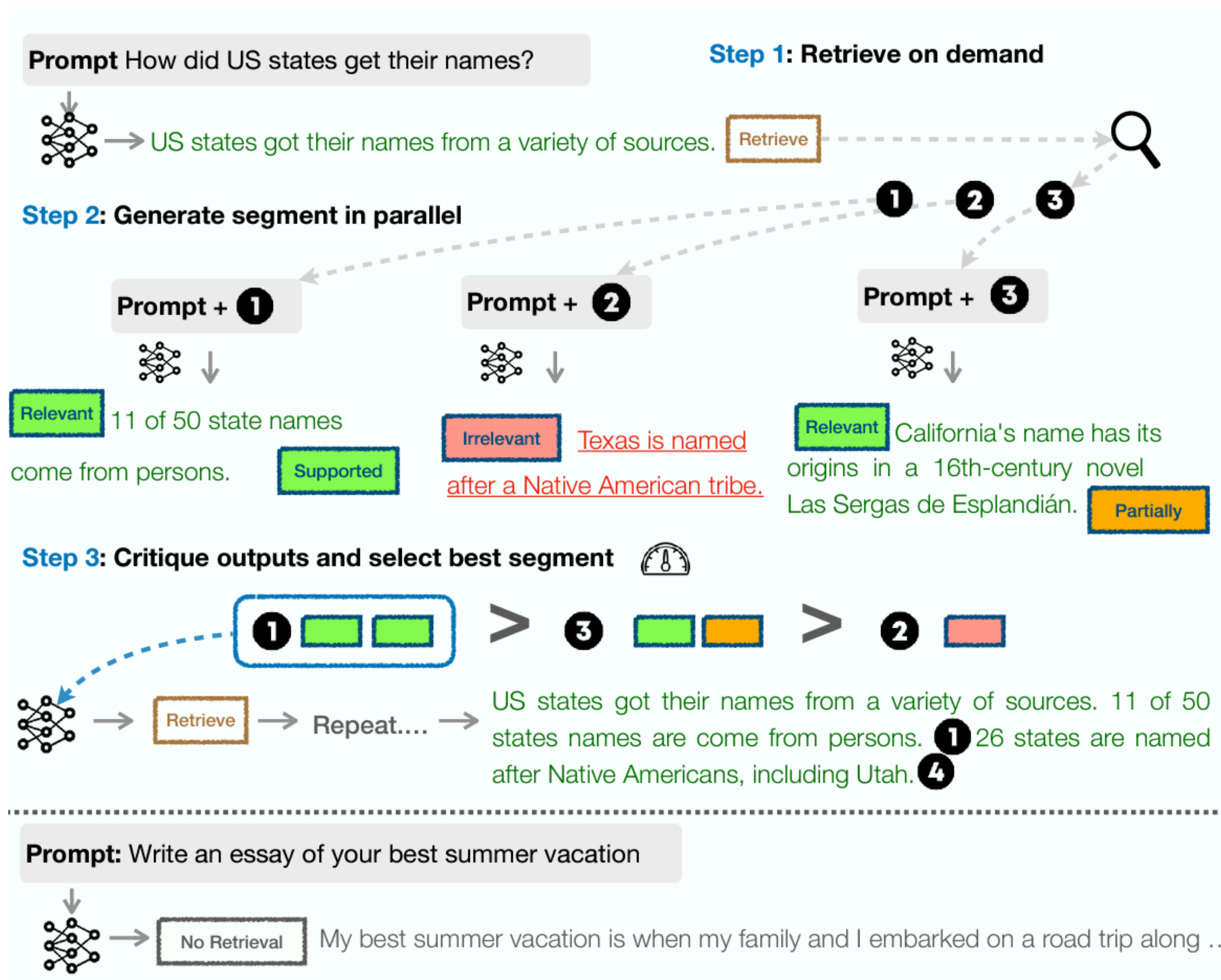


# BlendFilter



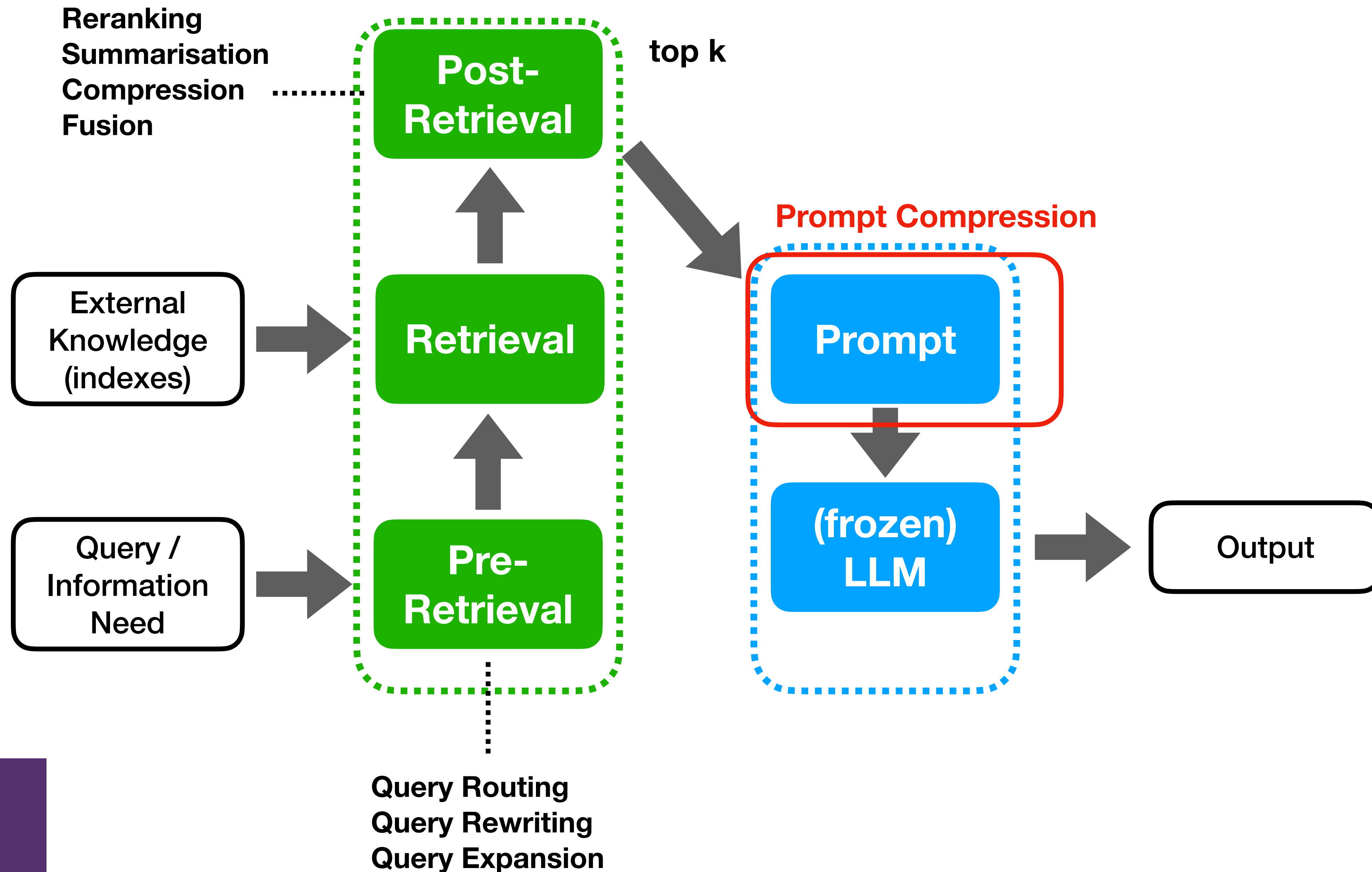
LLM generates answer based on filtered knowledge and original query

# Retrieve-Generate-Critique: Self-RAG



- **Adaptively retrieve** passages on-demand
  - Can decide to (1) **retrieve**, or (2) **not retrieve**
- LLM **reflects on retrieved passages** and its generation using special tokens (reflection tokens)
  - **Retrieval reflection:** do I need to retrieve?  
YES->output retrieval token to call retriever on demand
  - **Critique reflection:** is the generated answer good?  
SELECT best based on factuality and overall quality

# In the remaining of this part...

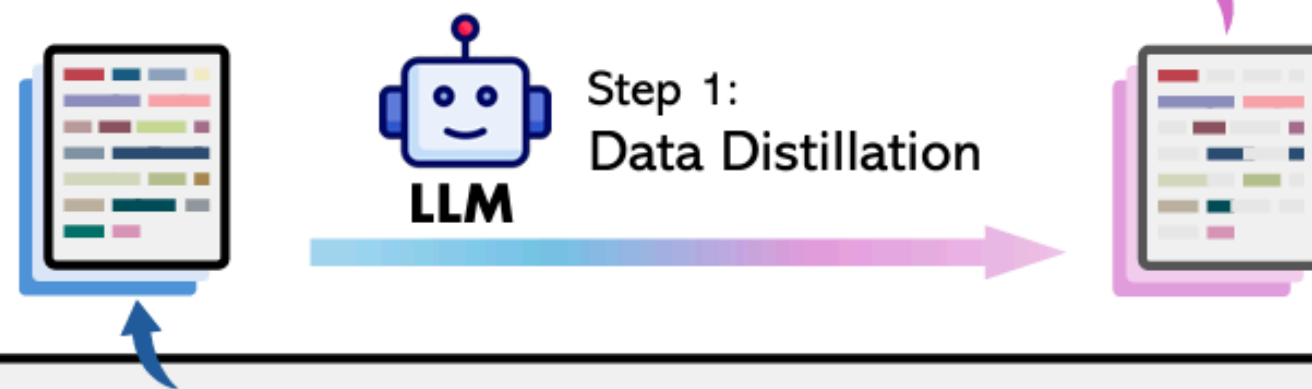


# Prompt Compression

- Reduce the prompt (aka context) to consume less tokens than in its original form.  
Why prompt compression?
- Long prompts often **confuse LLMs** and yield **lower effectiveness**
- **Latency** of decoder-based LLMs is **quadratic** with respect to **prompt length**
- Two main **categories** of approaches to prompt compression:
  1. **Lexical based:** Compress prompts by **removing tokens** according to their information entropy obtained from LLM
  2. **Embedding based:** Compress prompt into **special tokens** (short compact *memory slots*) that can be directly conditioned on by the LLM

# Lexical-based Compression via Distillation

**Compressed Text:** Item 15, City Manager Recommendation adopt three resolutions. Join Victory Pace program. Join California first program. Consent inclusion properties jurisdiction California Hero program. Emotion, motion, second, public comment. Cast vote. Public comment? Come forward. Alex Mitchell, represent Hero program. Hero program in California three half years



**Original Text:** Item 15, report from City Manager Recommendation to adopt three resolutions. First, to join the Victory Pace program. Second, to join the California first program. And number three, consenting to to inclusion of certain properties within the jurisdiction in the California Hero program. It was emotion, motion, a second and public comment. CNN. Please cast your vote. Oh. Was your public comment? Yeah. Please come forward. I thank you, Mr. Mayor. Thank you. Members of the council. My name is Alex Mitchell. I represent the hero program. Just wanted to let you know that the hero program. Has been in California for the last three and a half years.

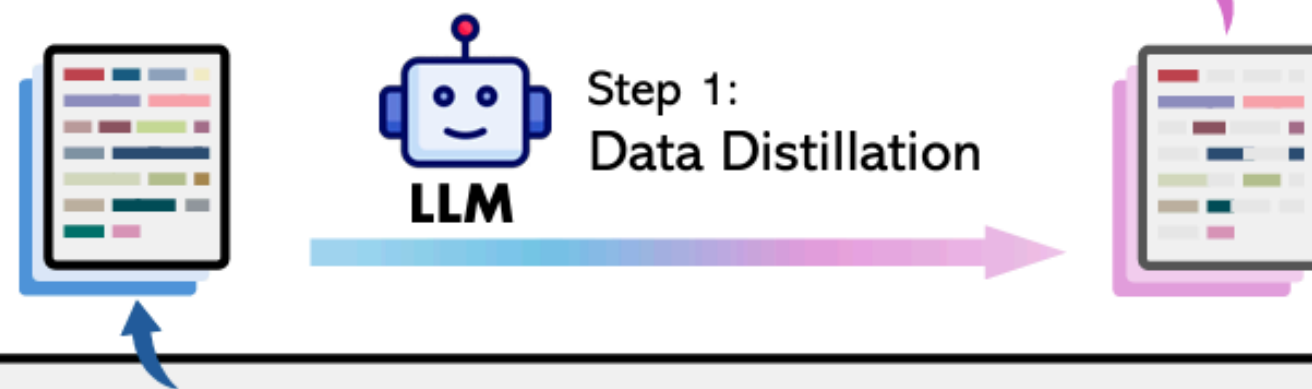
## Step 1: Data Distillation

Prompt GPT4 to generate compressed text from original text that meets criteria:

1. Token reduction
2. Informativeness
3. Faithfulness

# Lexical-based Compression via Distillation

**Compressed Text:** Item 15, City Manager Recommendation adopt three resolutions. Join Victory Pace program. Join California first program. Consent inclusion properties jurisdiction California Hero program. Emotion, motion, second, public comment. Cast vote. Public comment? Come forward. Alex Mitchell, represent Hero program. Hero program in California three half years



**Original Text:** Item 15, report from City Manager Recommendation to adopt three resolutions. First, to join the Victory Pace program. Second, to join the California first program. And number three, consenting to to inclusion of certain properties within the jurisdiction in the California Hero program. It was emotion, motion, a second and public comment. CNN. Please cast your vote. Oh. Was your public comment? Yeah. Please come forward. I thank you, Mr. Mayor. Thank you. Members of the council. My name is Alex Mitchell. I represent the hero program. Just wanted to let you know that the hero program. Has been in California for the last three and a half years.

## Our Instruction for Compression:

Compress the given text to short expressions, and such that you (GPT-4) can reconstruct it as close as possible to the original. Unlike the usual text compression, I need you to comply with the 5 conditions below:

1. You can ONLY remove unimportant words.
2. Do not reorder the original words.
3. Do not change the original words.
4. Do not use abbreviations or emojis.
5. Do not add new words or symbols.

Compress the origin aggressively by removing words only. Compress the origin as short as you can, while retaining as much information as possible. If you understand, please compress the following text: `{text to compress}`  
The compressed text is:

## Step 1: Data Distillation

Prompt GPT4 to generate compressed text from original text that meets criteria:

1. Token reduction
2. Informativeness
3. Faithfulness

Challenge: GPT-4 does not consistently follow instructions. Remedies:

- Instruction design (i.e. tweaked the prompt)
- Chunk-wise compression (i.e. segment long context into multiple 512 tokens chunks)

Pan, Z., Wu, Q., Jiang, H., Xia, M., Luo, X., Zhang, J., Lin, Q., Rühle, V., Yang, Y., Lin, C.Y. and Zhao, H.V., 2024. Lmlingua-2: Data distillation for efficient and faithful task-agnostic prompt compression. *arXiv preprint arXiv:2403.12968*.

# Lexical-based Compression via Distillation

**Compressed Text:** Item 15, City Manager Recommendation adopt three resolutions. Join Victory Pace program. Join California first program. Consent inclusion properties jurisdiction California Hero program. Emotion, motion, second, public comment. Cast vote. Public comment? Come forward. Alex Mitchell, represent Hero program. Hero program in California three half years



**Original Text:** Item 15, report from City Manager Recommendation to adopt three resolutions. First, to join the Victory Pace program. Second, to join the California first program. And number three, consenting to to inclusion of certain properties within the jurisdiction in the California Hero program. It was emotion, motion, a second and public comment. CNN. Please cast your vote. Oh. Was your public comment? Yeah. Please come forward. I thank you, Mr. Mayor. Thank you. Members of the council. My name is Alex Mitchell. I represent the hero program. Just wanted to let you know that the hero program. Has been in California for the last three and a half years.

## Step 2: Data Annotation

Assign binary label to each token in original text to determine if it should be preserved or discarded after compression.

### Challenges:

1. Ambiguity: a word in compressed text may appear multiple times in original text
2. Variation: GPT-4 may modify the original words in tense, plural form, etc. during compression
3. Reordering: order of words may be changed

# Lexical-based Compression via Distillation

**Compressed Text:** Item 15, City Manager Recommendation adopt three resolutions. Join Victory Pace program. Join California first program. Consent inclusion properties jurisdiction California Hero program. Emotion, motion, second, public comment. Cast vote. Public comment? Come forward. Alex Mitchell, represent Hero program. Hero program in California three half years



Step 1:  
Data Distillation



Step 2:  
Data Annotation

**Original Text:** Item 15, report from City Manager Recommendation to adopt three resolutions. First, to join the Victory Pace program. Second, to join the California first program. And number three, consenting to to inclusion of certain properties within the jurisdiction in the California Hero program. It was emotion, motion, a second and public comment. CNN. Please cast your vote. Oh. Was your public comment? Yeah. Please come forward. I thank you, Mr. Mayor. Thank you. Members of the council. My name is Alex Mitchell. I represent the hero program. Just wanted to let you know that the hero program. Has been in California for the last three and a half years.



Step 3:  
Quality Control  
& Filtering

## Step 3: Quality Control & Filtering

- (i) Assess quality of compressed and of automatically annotated labels
- (ii) Then filter examples by scores

Two quality control metrics:

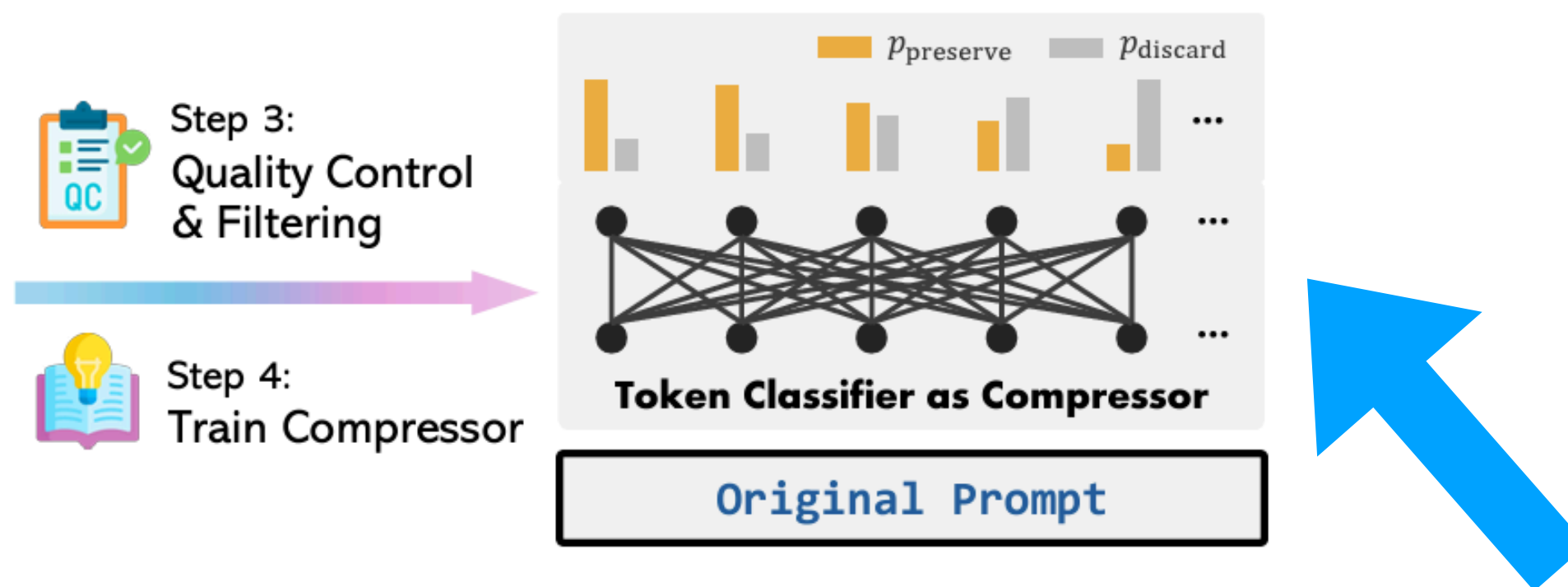
1. Variation Rate: proportion of words in compressed text that are absent in original text
2. Alignment Gap: xxx

# Lexical-based Compression via Distillation

**Compressed Text:** Item 15, City Manager Recommendation adopt three resolutions. Join Victory Pace program. Join California first program. Consent inclusion properties jurisdiction California Hero program. Emotion, motion, second, public comment. Cast vote. Public comment? Come forward. Alex Mitchell, represent Hero program. Hero program in California three half years

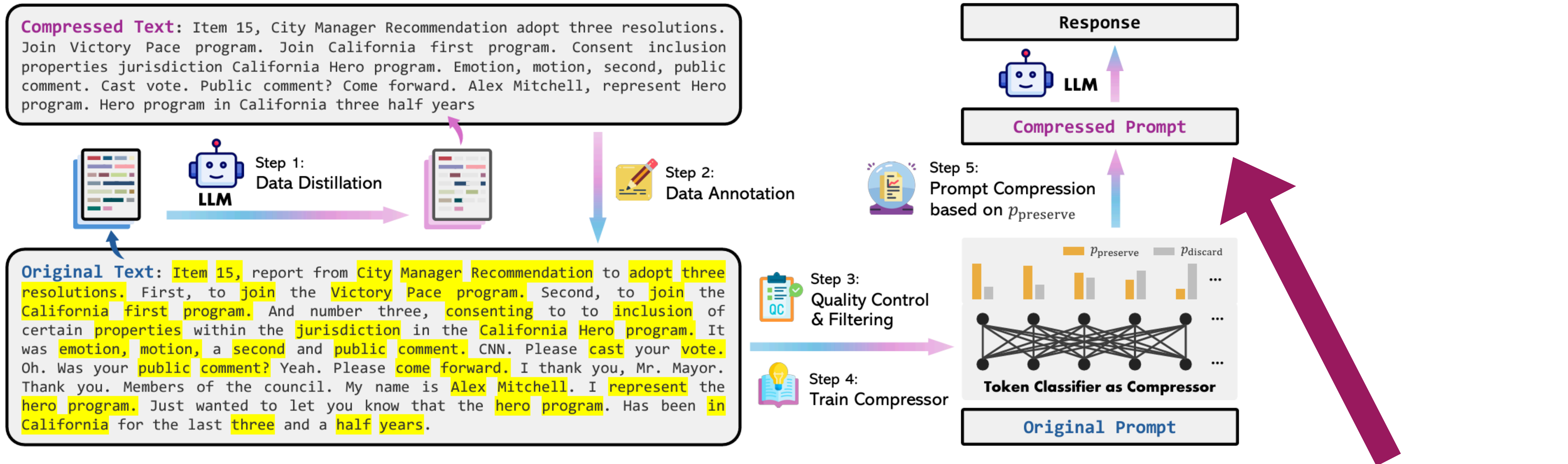


**Original Text:** Item 15, report from City Manager Recommendation to adopt three resolutions. First, to join the Victory Pace program. Second, to join the California first program. And number three, consenting to to inclusion of certain properties within the jurisdiction in the California Hero program. It was emotion, motion, a second and public comment. CNN. Please cast your vote. Oh. Was your public comment? Yeah. Please come forward. I thank you, Mr. Mayor. Thank you. Members of the council. My name is Alex Mitchell. I represent the hero program. Just wanted to let you know that the hero program. Has been in California for the last three and a half years.



- ### Step 4: Train Compressor (and then infer)
- Formulate prompt compression as binary token classification problem (i.e., preserve or discard)
  - Smaller Transformer Encoder model (e.g. mBERT) allows for lower latency than using the larger LLM

# Lexical-based Compression via Distillation



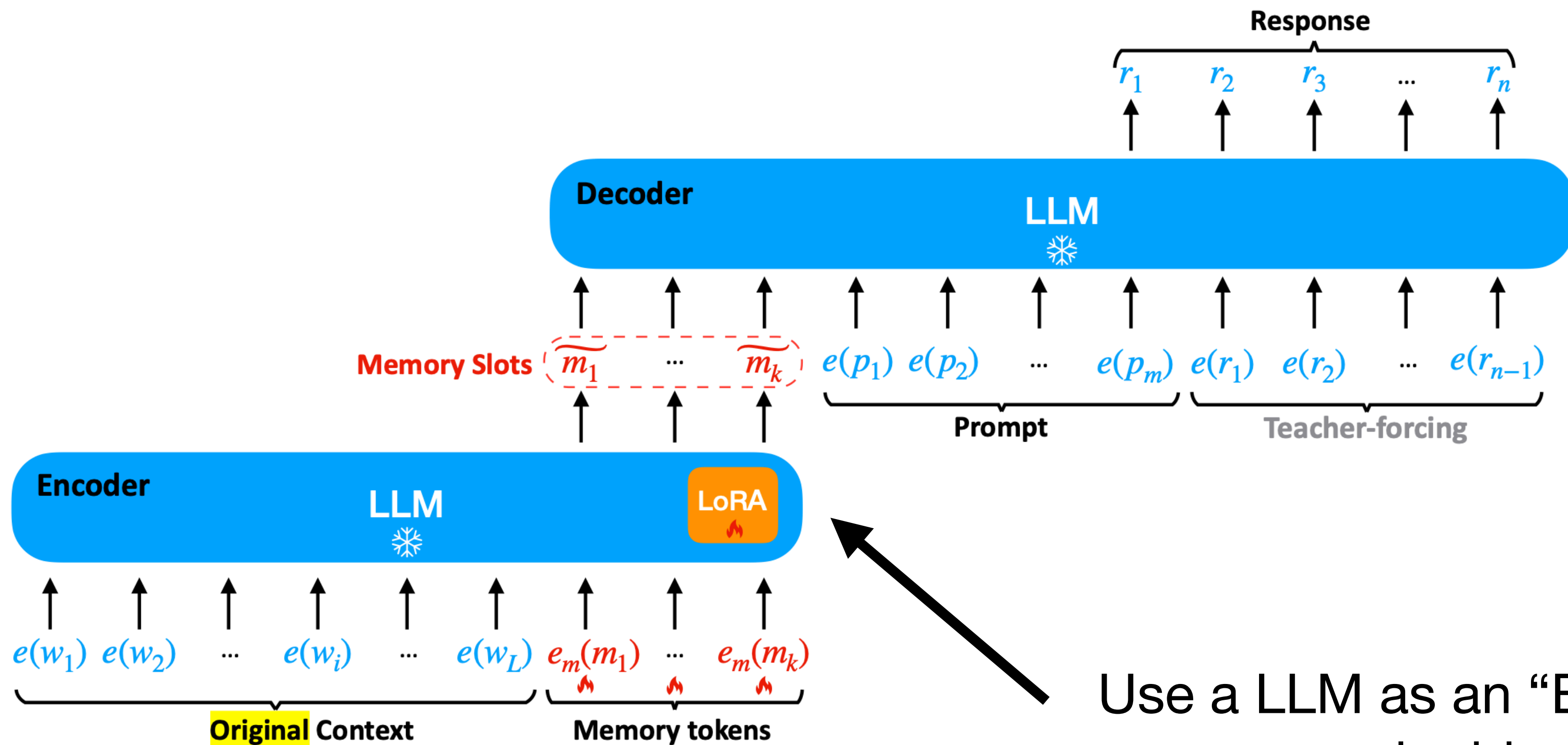
## Stage 5: Prompt Compression

1. use token classification model to predict probability of each word being labeled as preserve ( $p(\text{word}) = \text{avg}(p(\text{tokens of word}))$ )
2. retain top words in original prompt with highest preservation probability, maintain their original order

# Embedding-based Compression: In-context Autoencoder (ICAE)

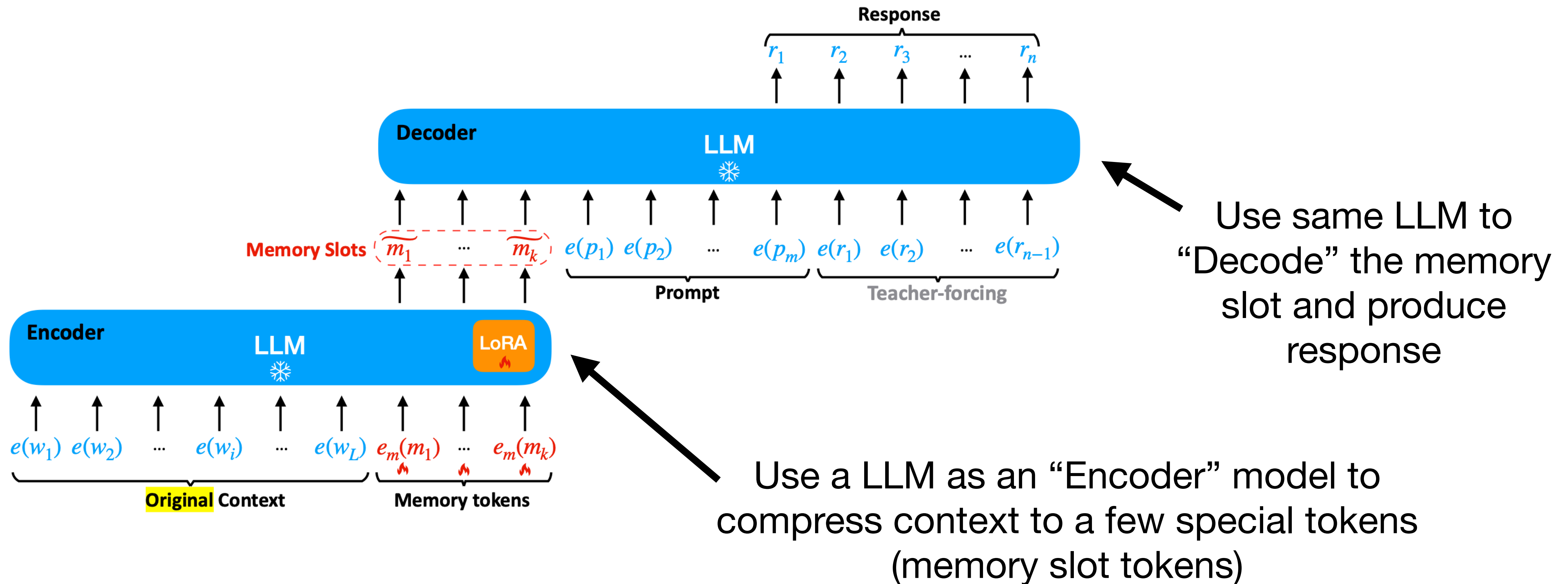
- Compress prompt into **special tokens** (short compact memory slots) that can be directly conditioned on by the LLM
- ~1% **additional parameters**, but 4× context compression
  - ➔ Reduced latency and GPU memory cost during inference

# Embedding-based Compression: In-context Autoencoder (ICAE)

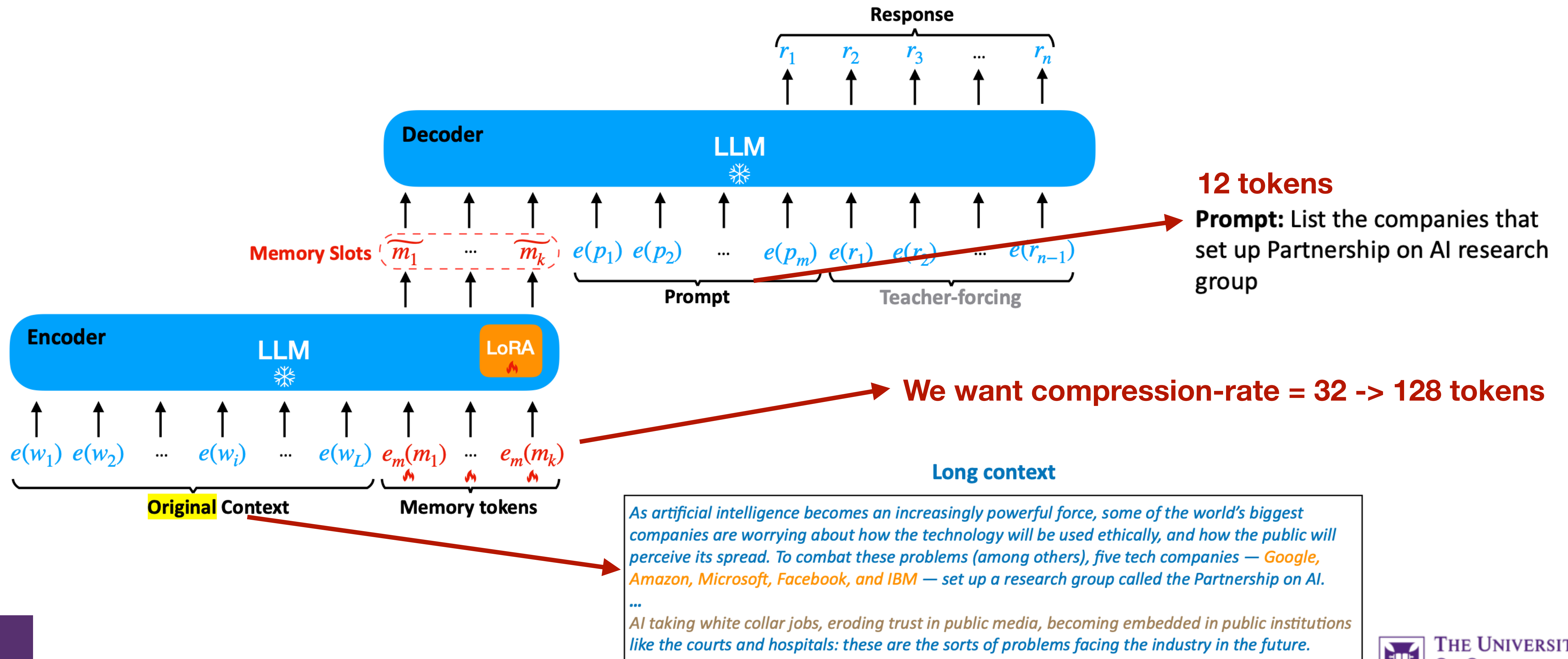


Use a LLM as an “Encoder” model to compress context to a few special tokens (memory slot tokens)

# Embedding-based Compression: In-context Autoencoder (ICAE)



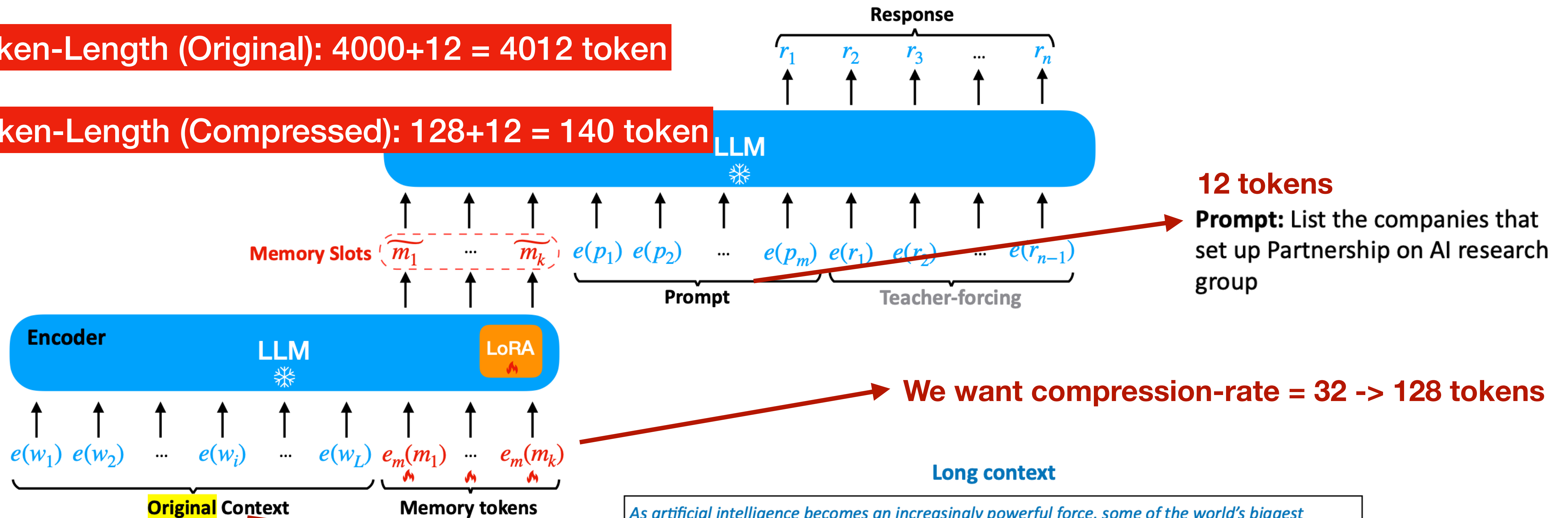
# Embedding-based Compression: In-context Autoencoder (ICAE)



# Embedding-based Compression: In-context Autoencoder (ICAE)

Token-Length (Original):  $4000+12 = 4012$  token

Token-Length (Compressed):  $128+12 = 140$  token



*As artificial intelligence becomes an increasingly powerful force, some of the world's biggest companies are worrying about how the technology will be used ethically, and how the public will perceive its spread. To combat these problems (among others), five tech companies — Google, Amazon, Microsoft, Facebook, and IBM — set up a research group called the Partnership on AI.*

...

*AI taking white collar jobs, eroding trust in public media, becoming embedded in public institutions like the courts and hospitals: these are the sorts of problems facing the industry in the future.*

# Embedding-based Compression: In-context Autoencoder (ICAE)

- Compress prompt into **special tokens** (short compact memory slots) that can be directly conditioned on by the LLM
- ~1% **additional parameters**, but 4× context compression
  - ➔ Reduced latency and GPU memory cost during inference
- Training:
  1. pre-trained using auto-encoding and language modelling objectives on massive text data -> enables to generate memory slots that accurately and comprehensively represent original context.
  2. fine-tuned on instruction data -> produces desirable responses to various prompts.



# The struggle between model knowledge and retrieval knowledge

- The output of RAG is not guaranteed to be consistent with retrieved relevant passages
- Because the models are not explicitly trained to leverage and follow facts from provided passages.

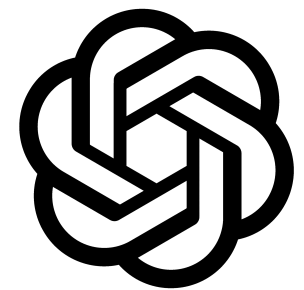
# Answers Inconsistent with RAG Evidence in Prompt

**Question** ..... *“Will drinking  
vinegar dissolve a  
stuck fish bone?”*



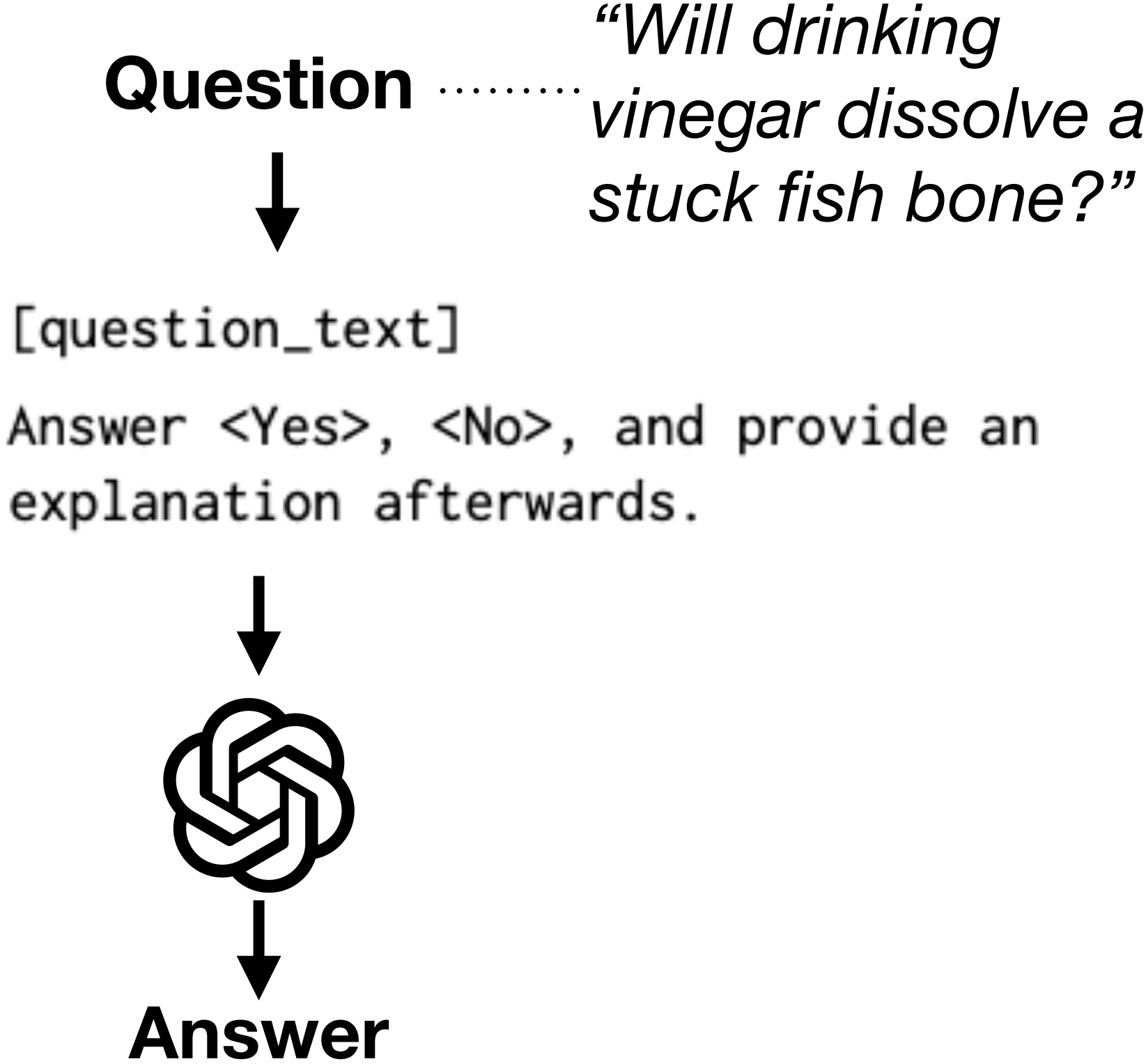
[question\_text]

Answer <Yes>, <No>, and provide an  
explanation afterwards.



**Answer**

# Answers Inconsistent with RAG Evidence in Prompt

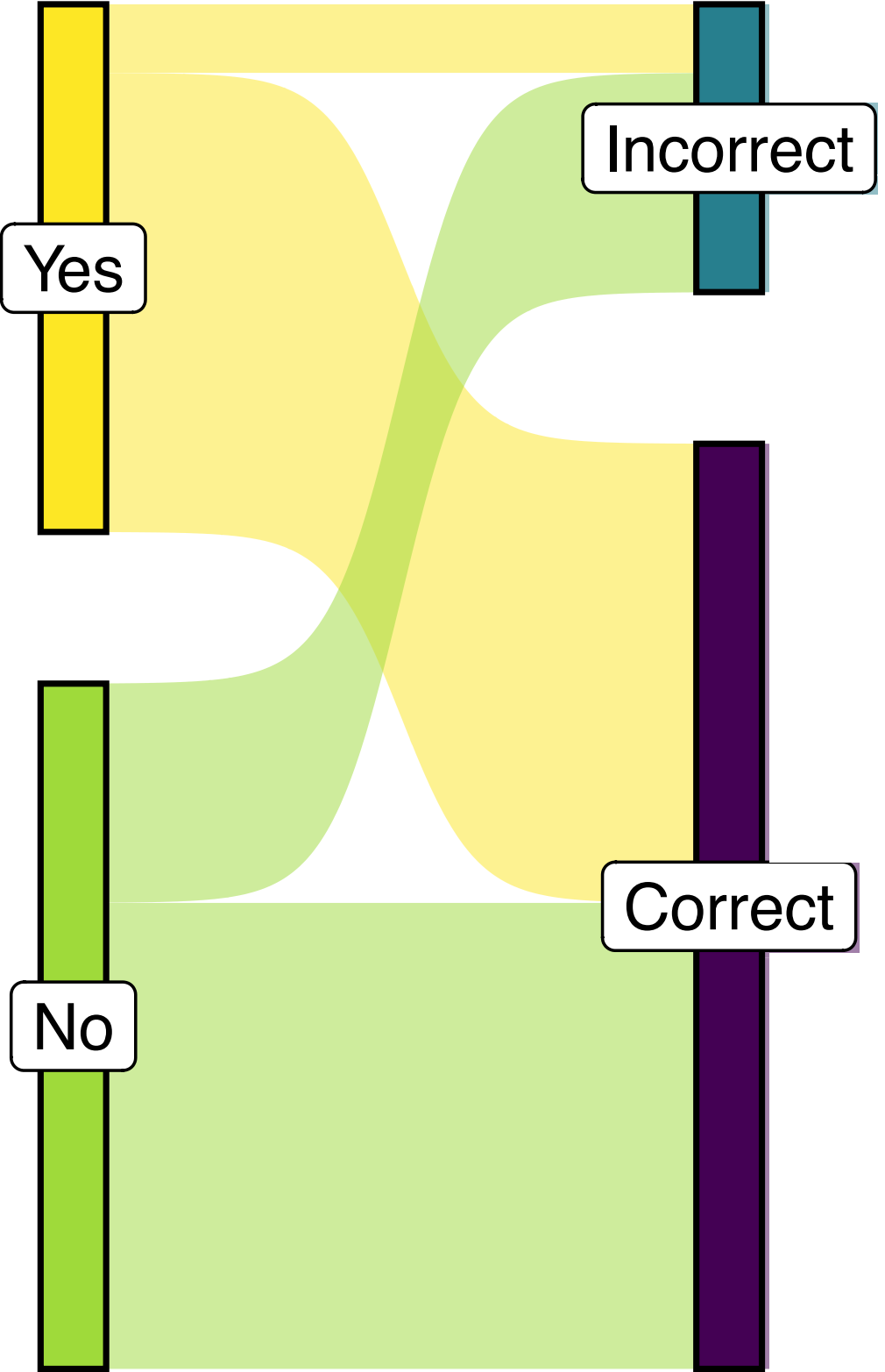
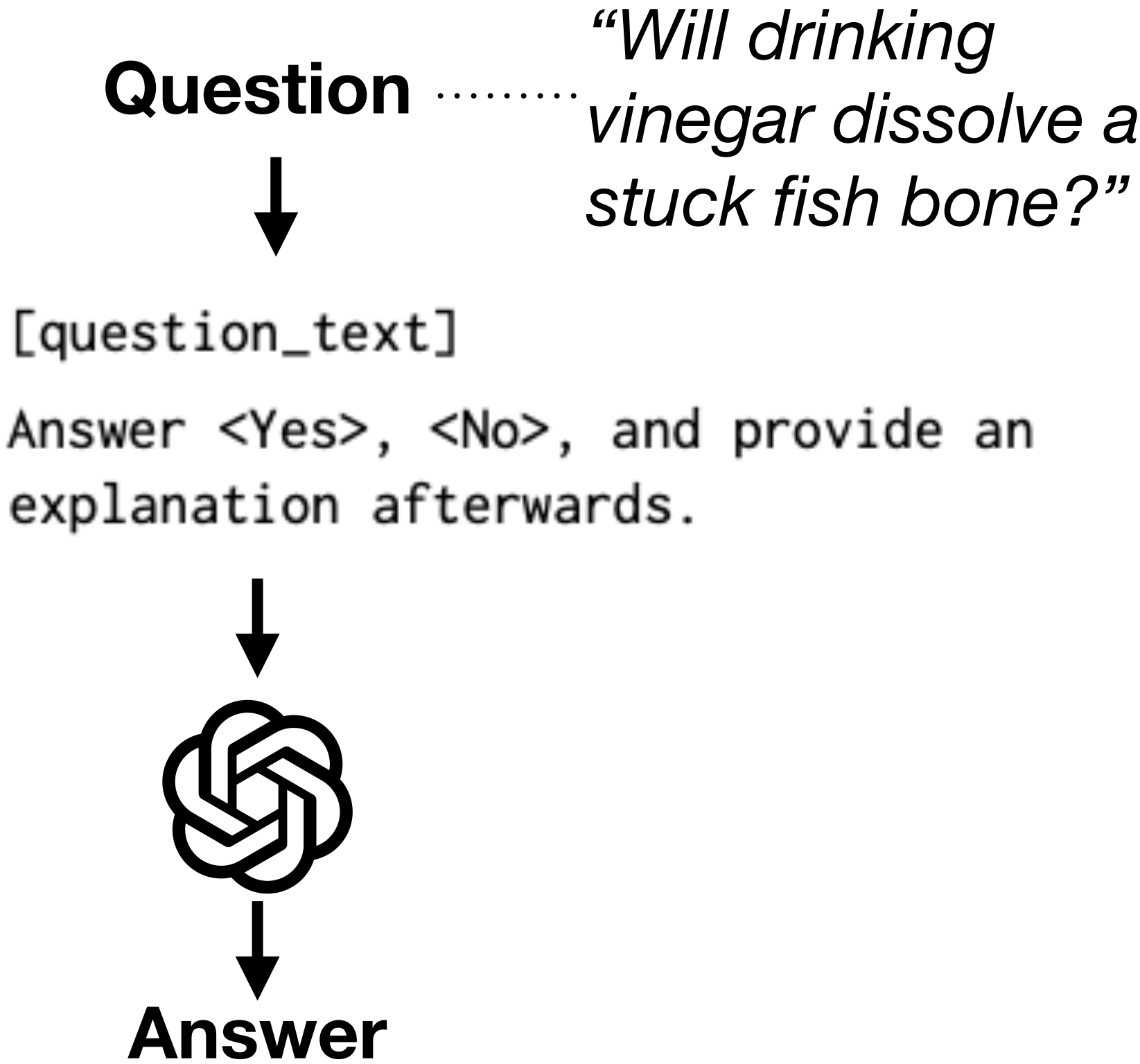


Yes

No

Ground truth answer

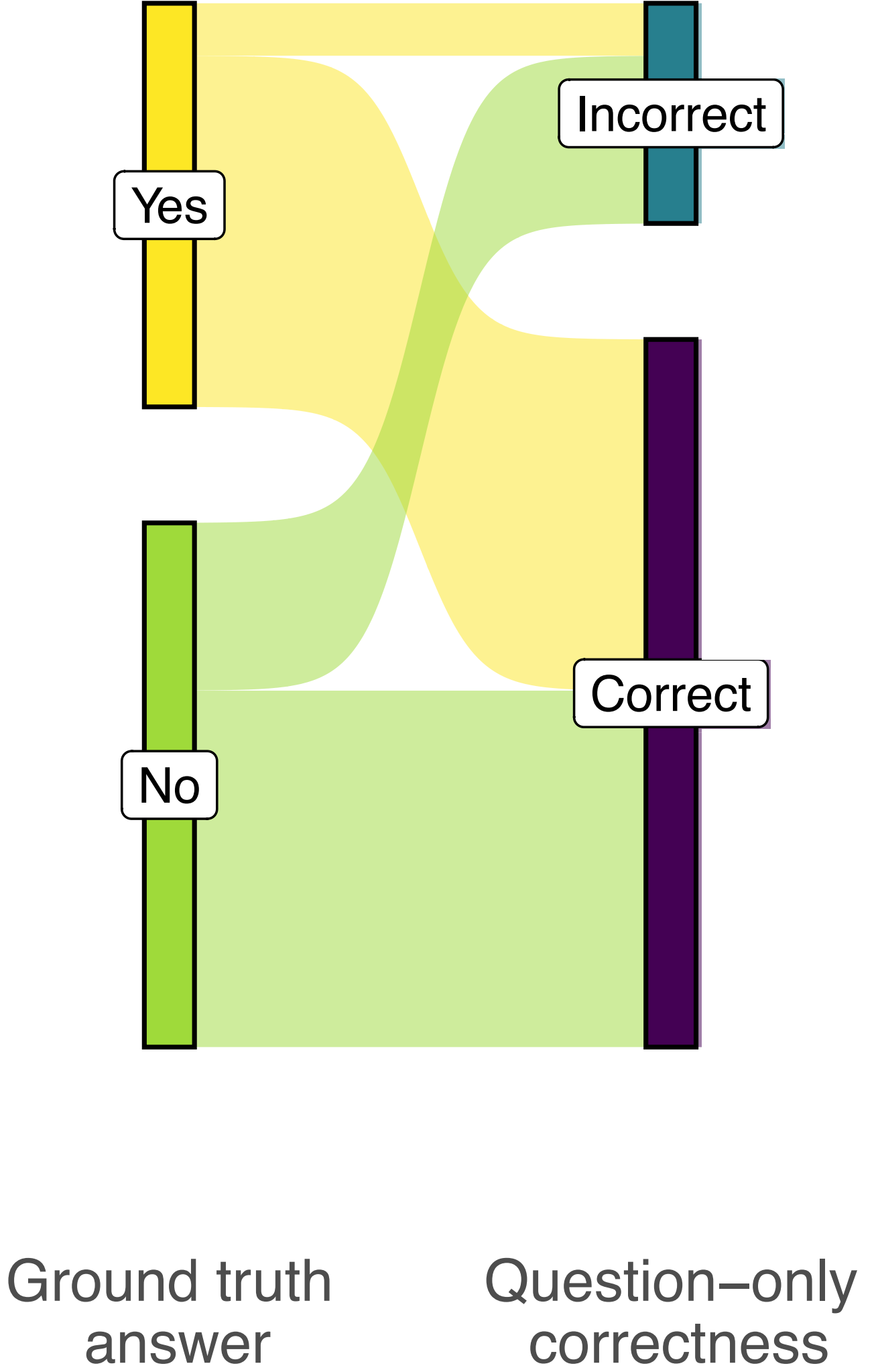
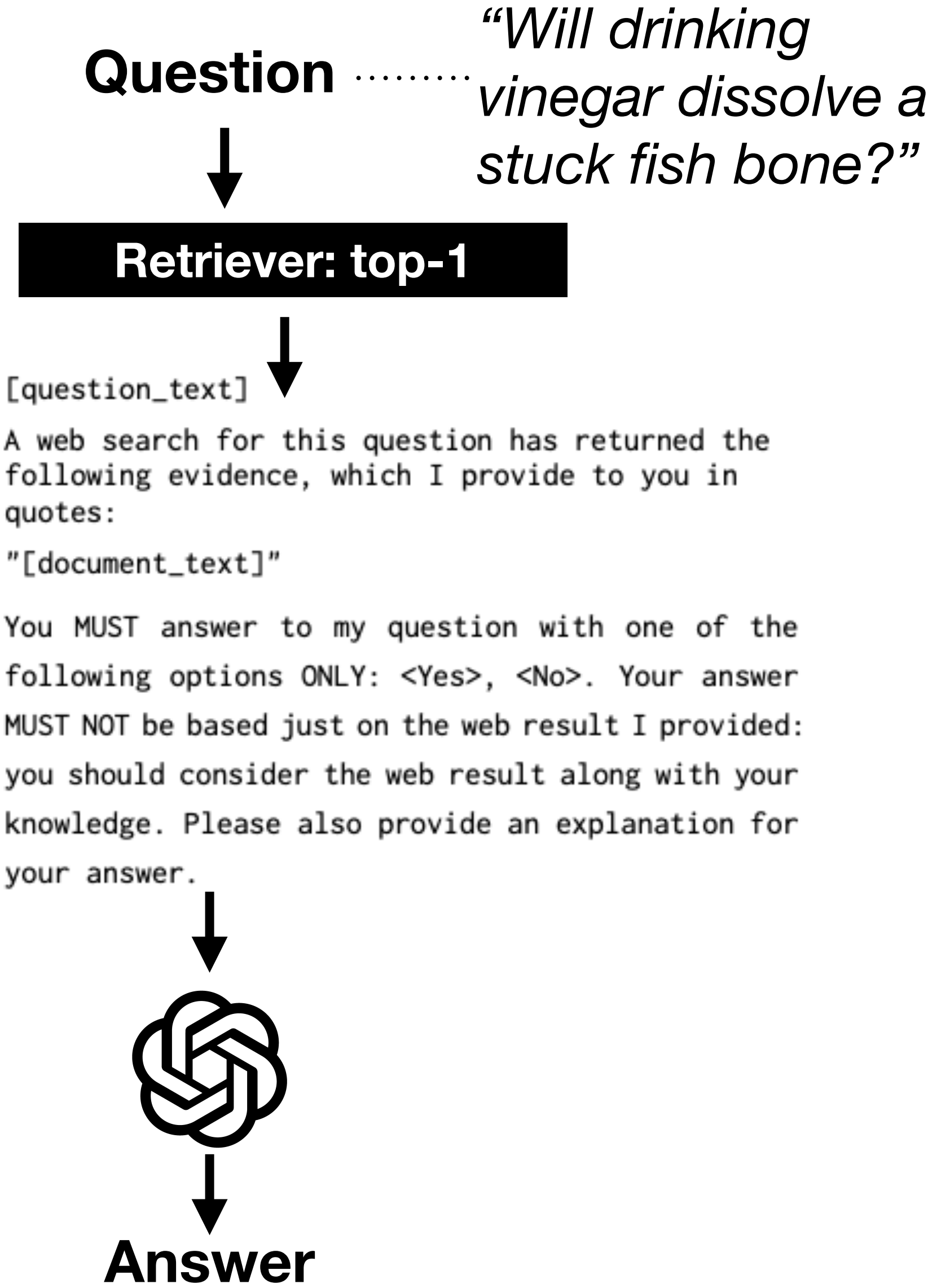
# Answers Inconsistent with RAG Evidence in Prompt



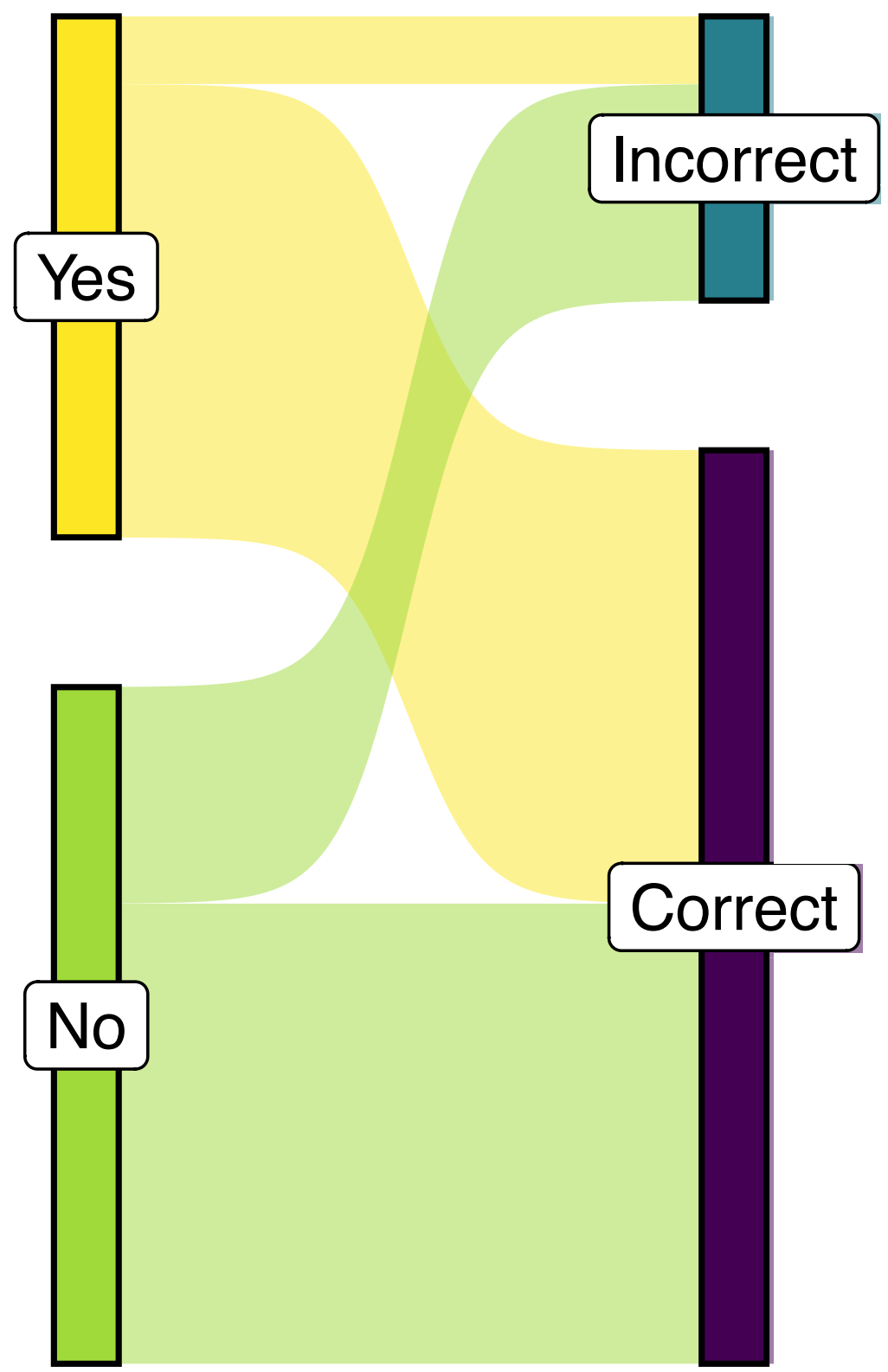
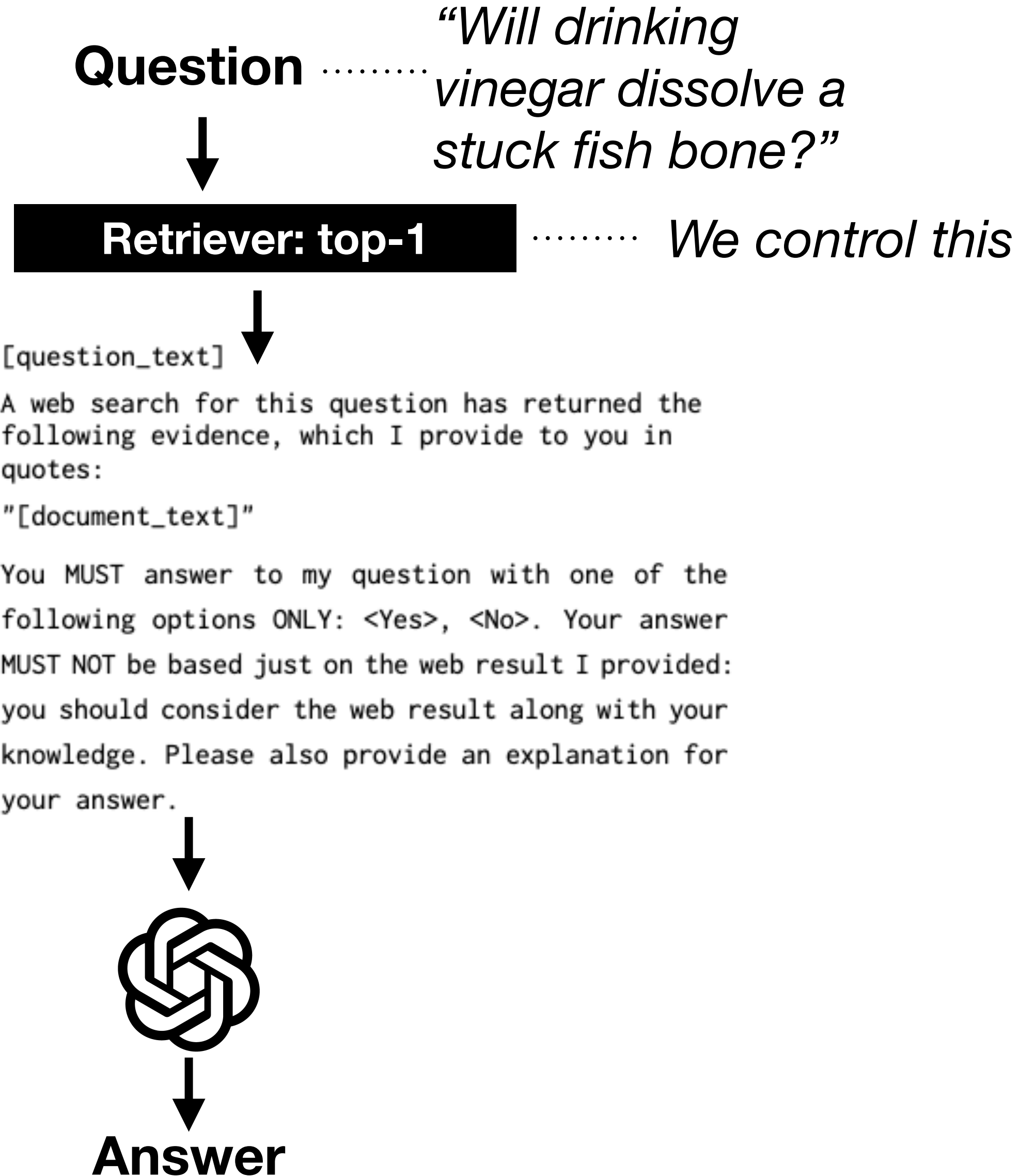
Ground truth answer

Question-only correctness

# Answers Inconsistent with RAG Evidence in Prompt



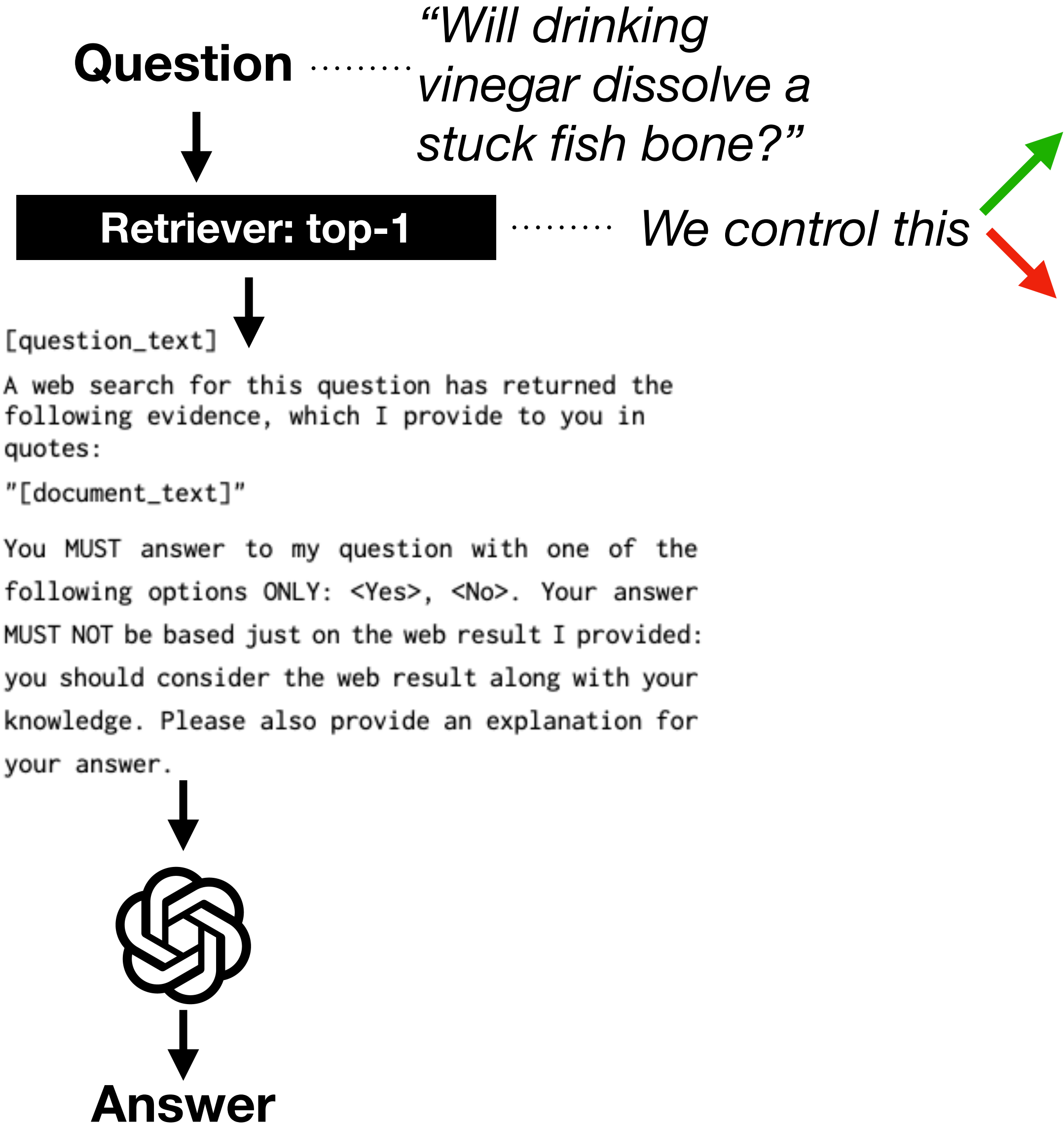
# Answers Inconsistent with RAG Evidence in Prompt



Ground truth answer

Question-only correctness

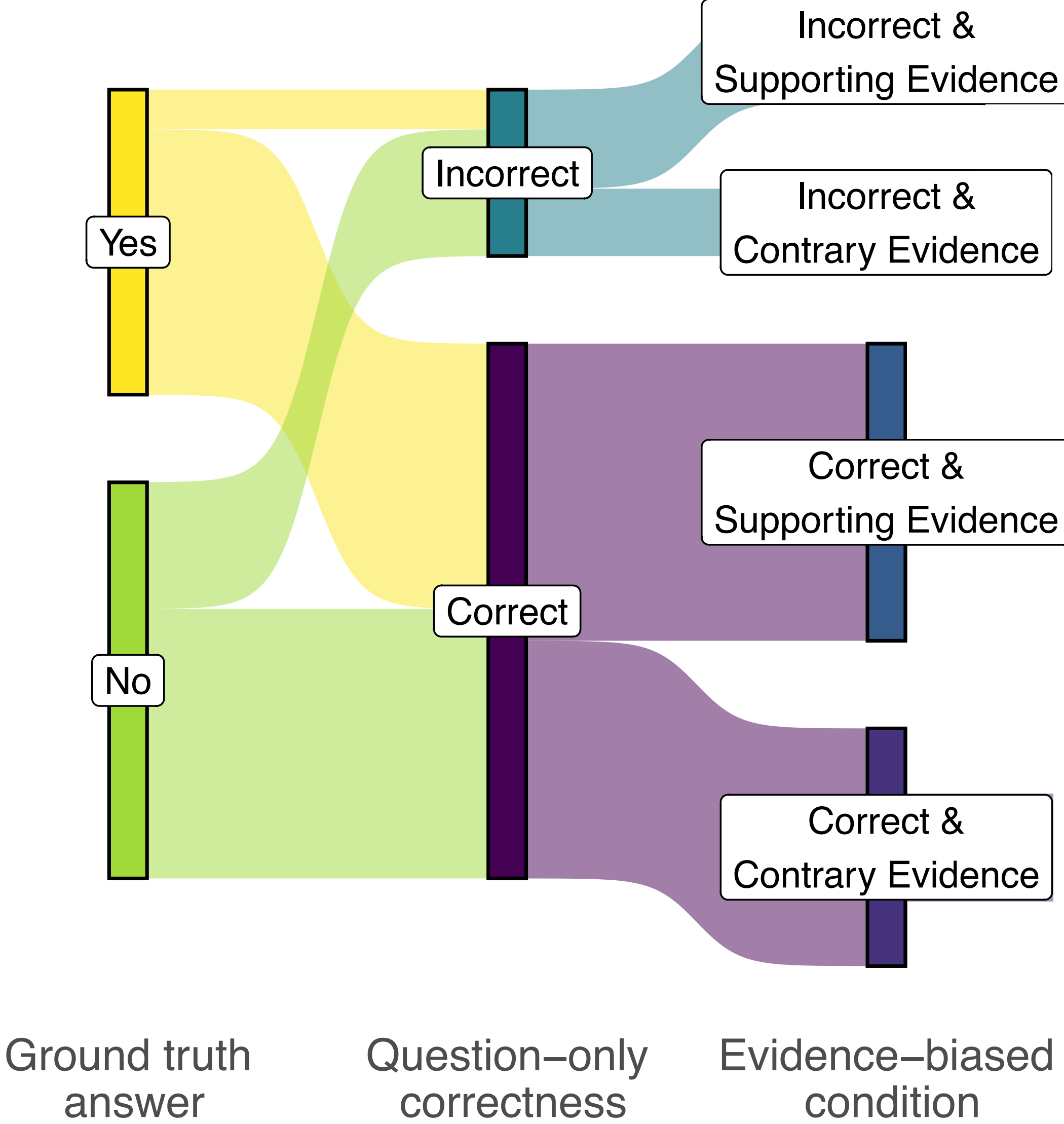
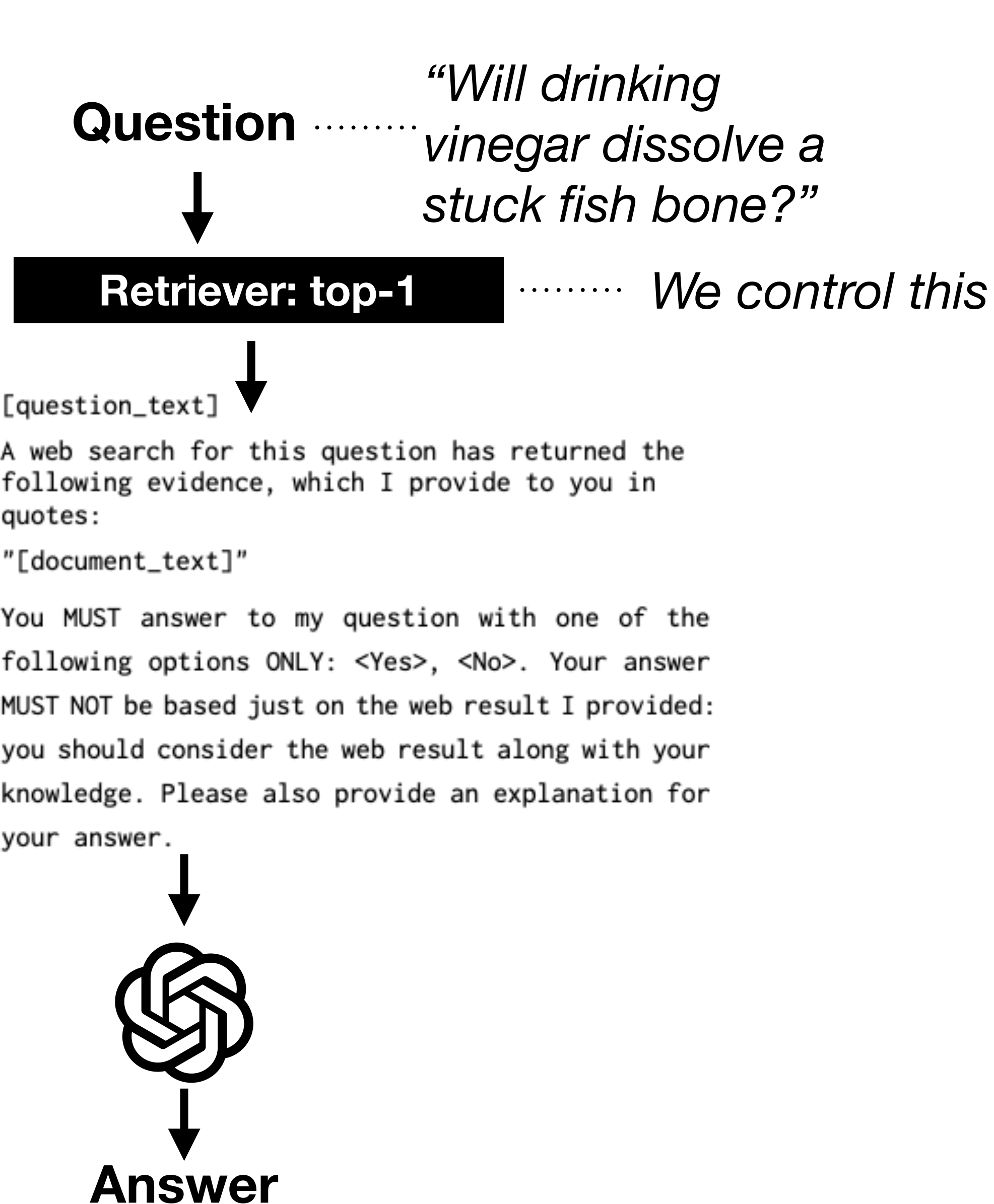
# Answers Inconsistent with RAG Evidence in Prompt



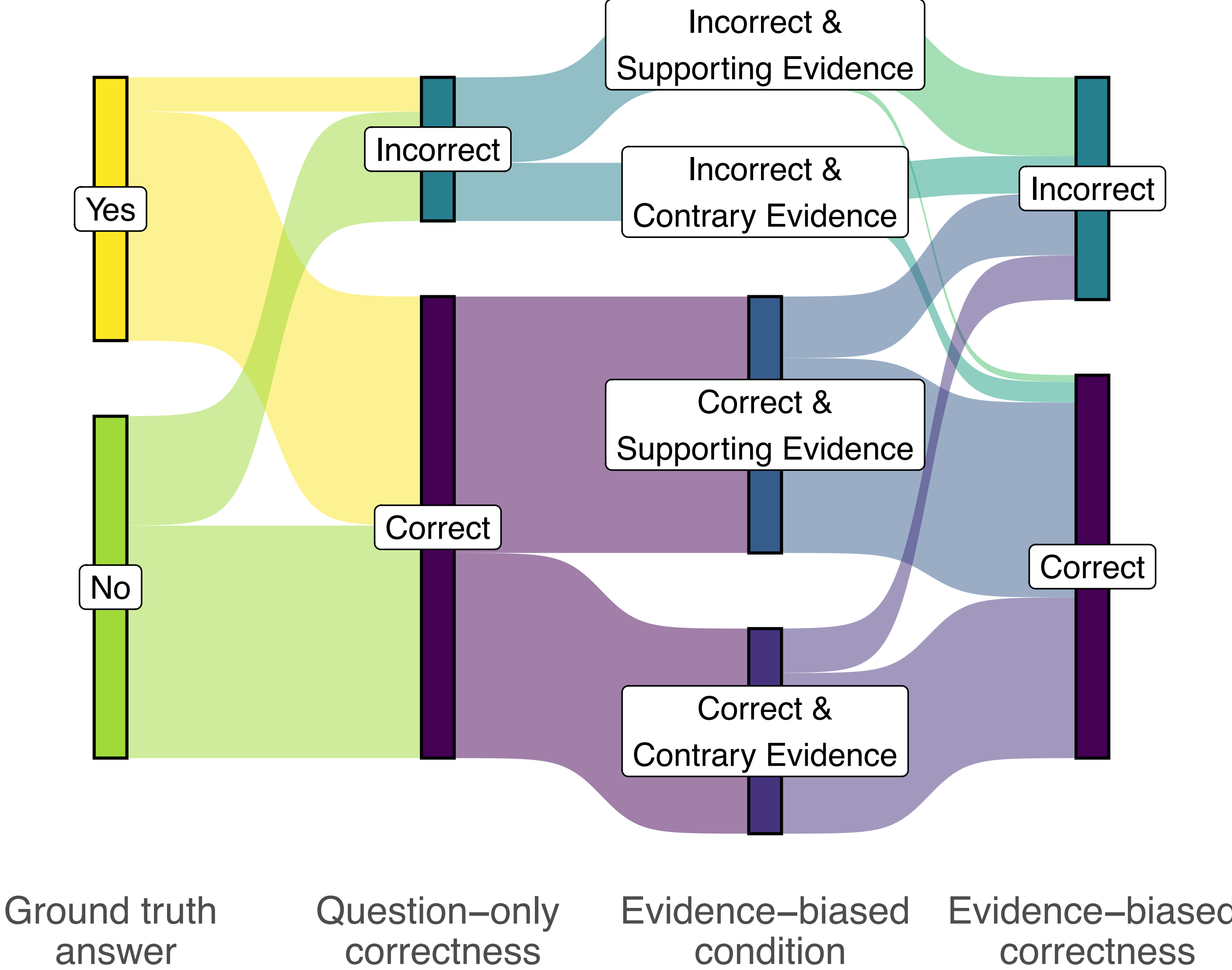
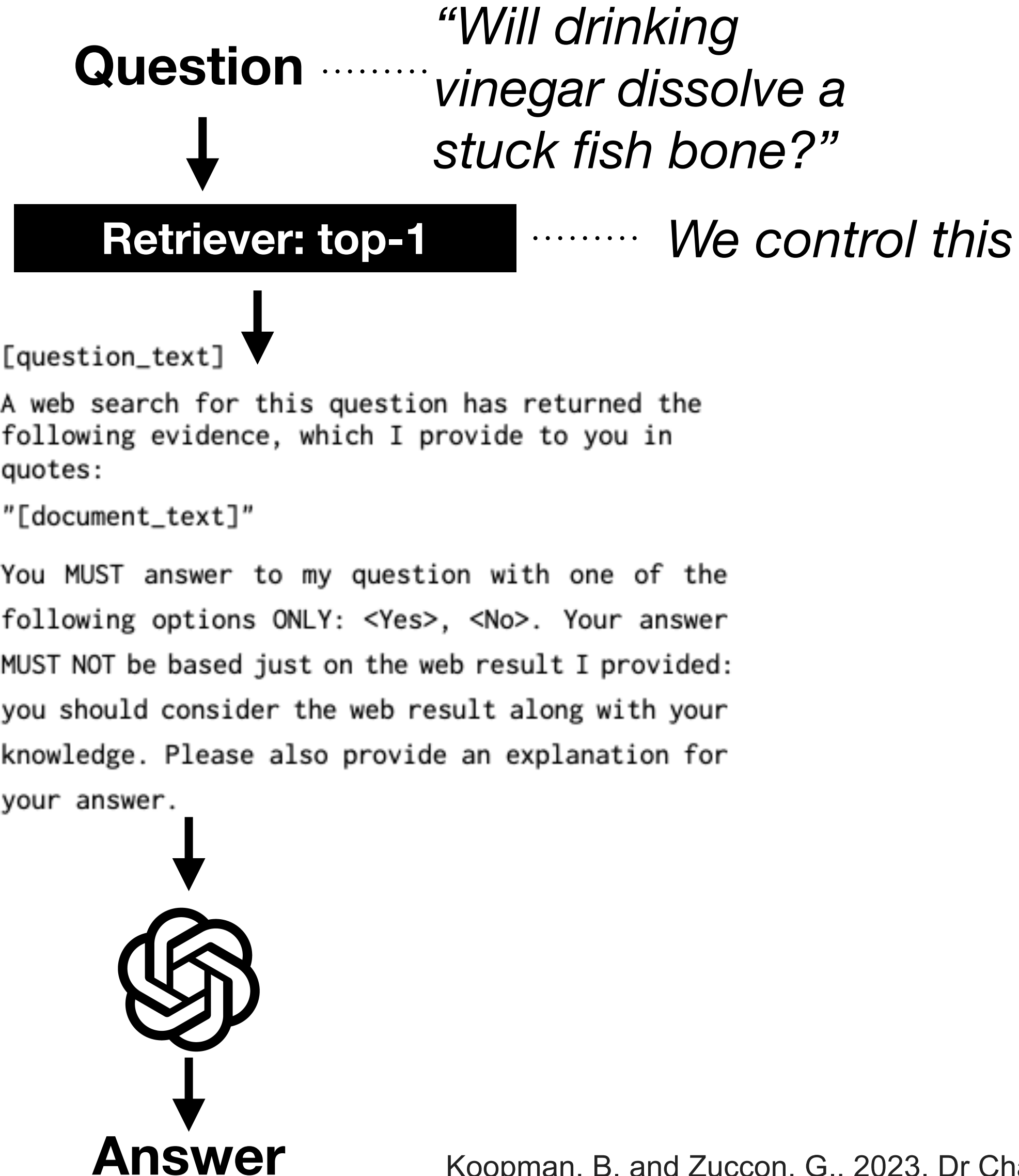
**Supporting Evidence: provides the same answer as the ground-truth**

**Contrary Evidence: provides answer contrary to the ground-truth**

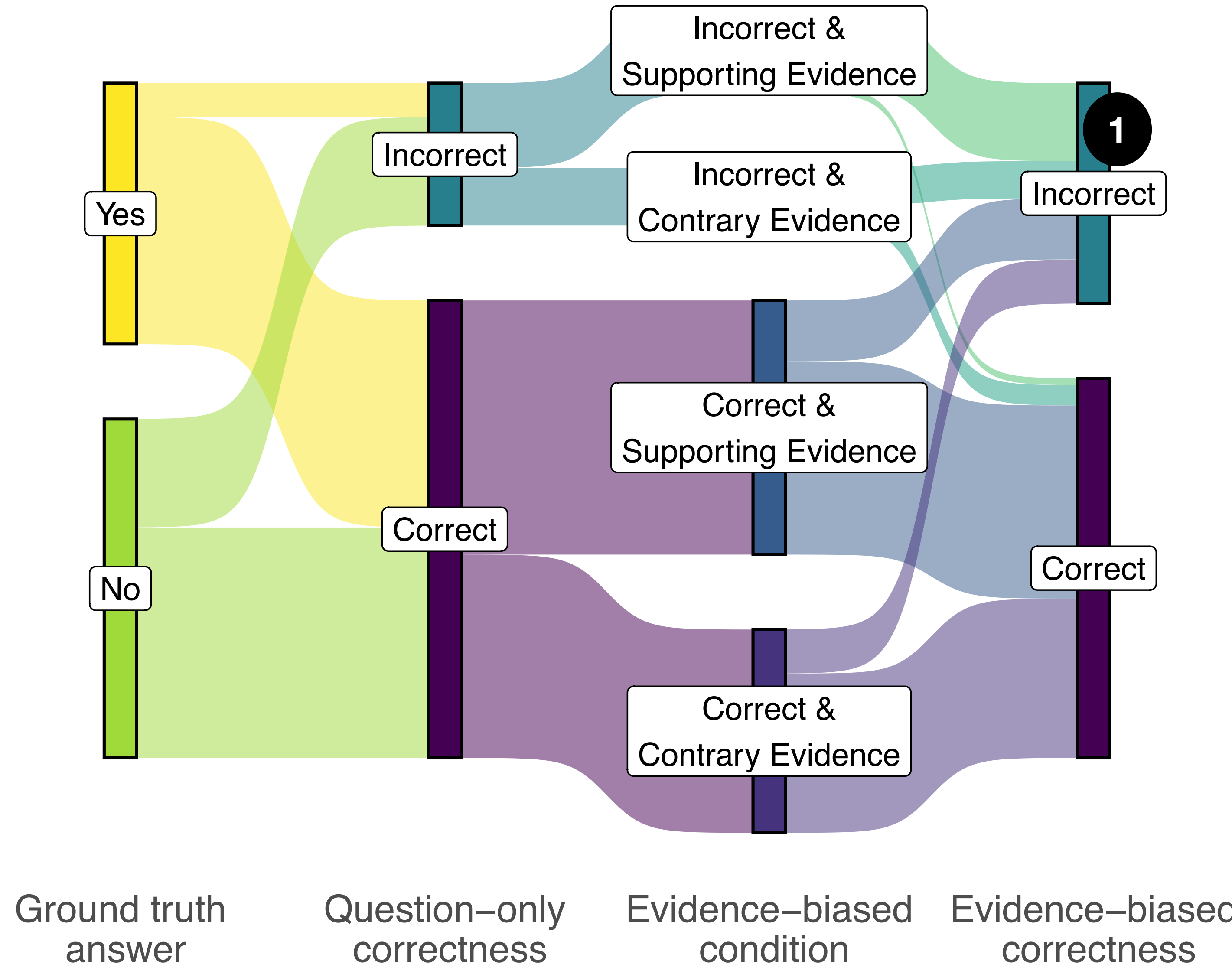
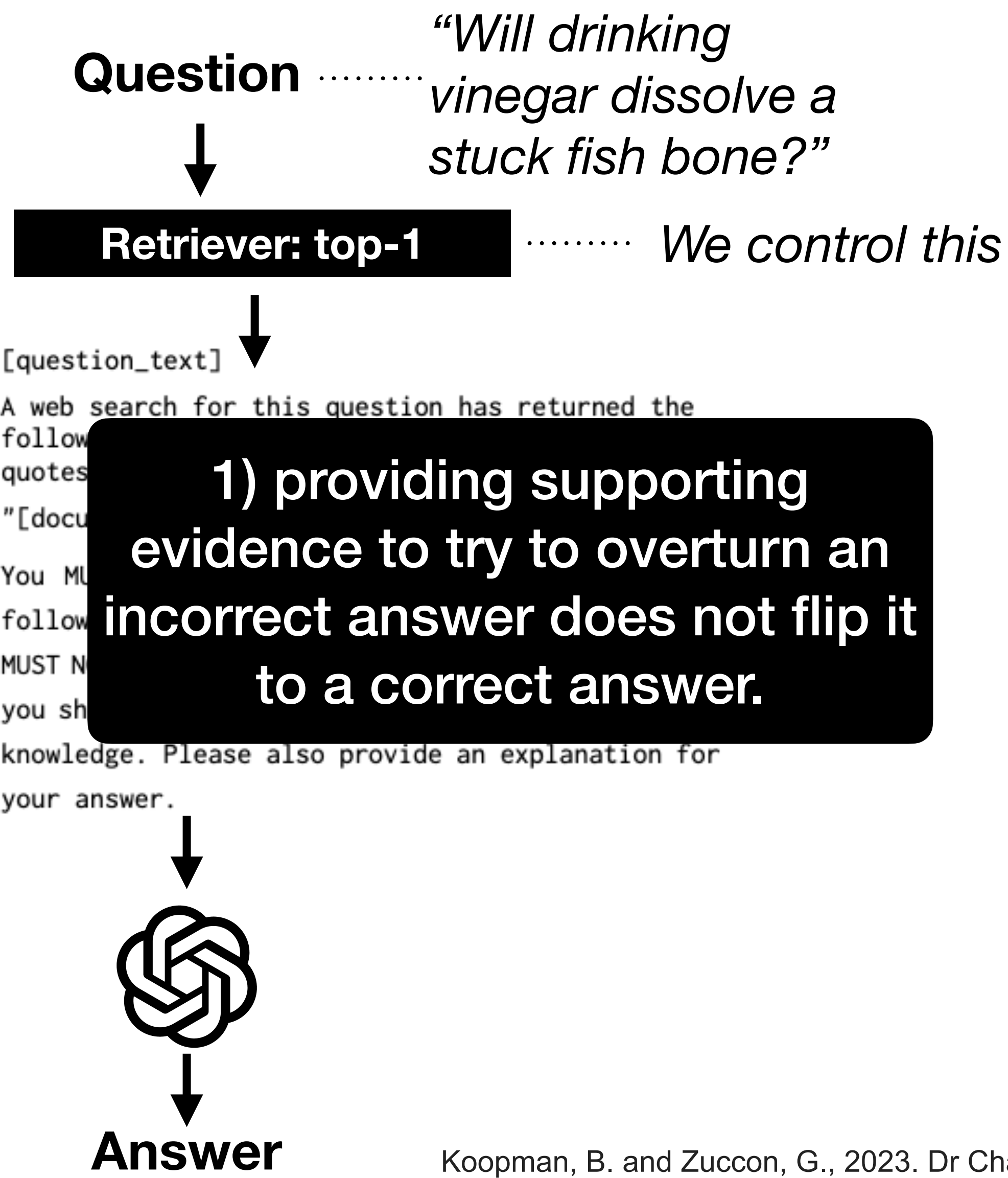
# Answers Inconsistent with RAG Evidence in Prompt



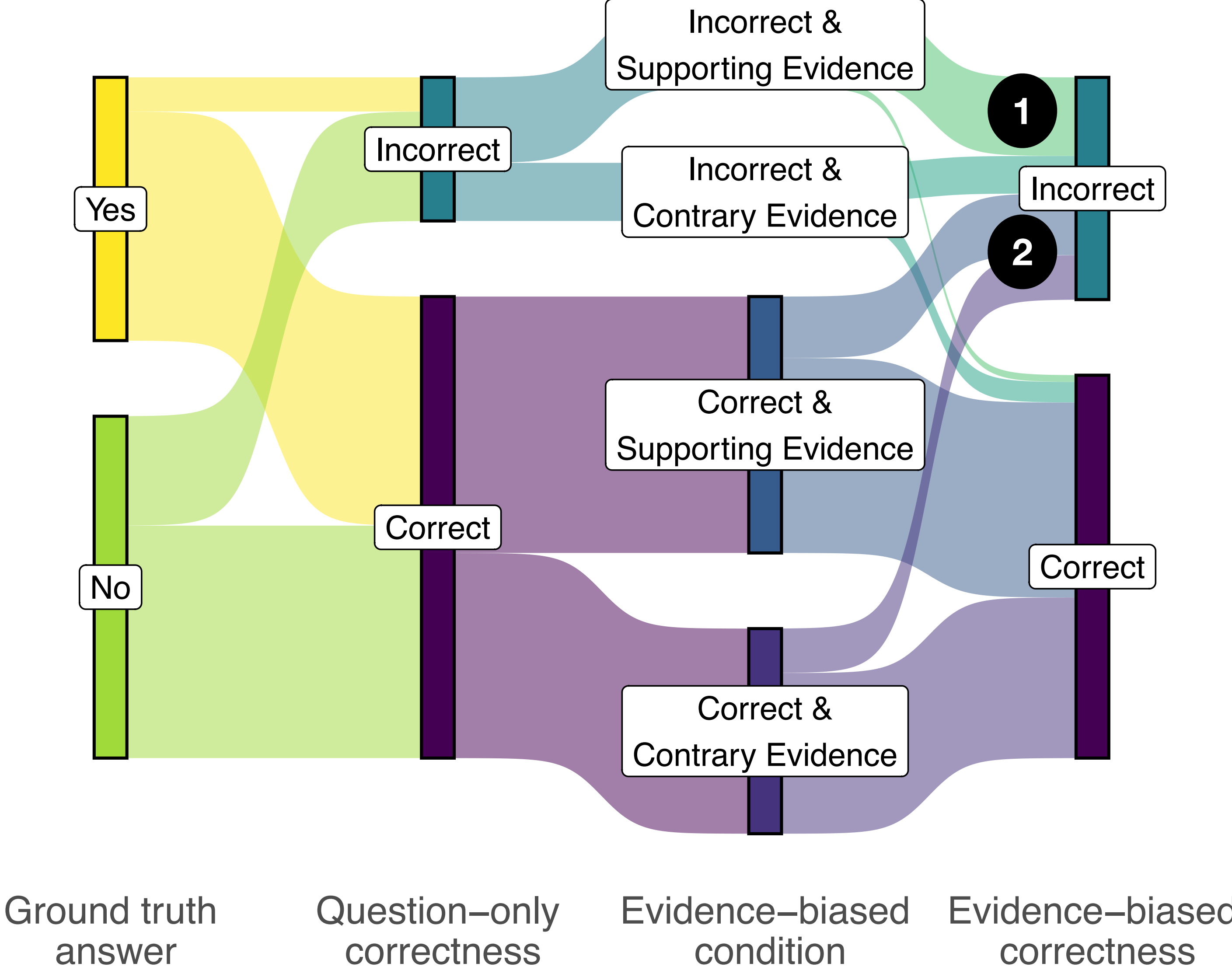
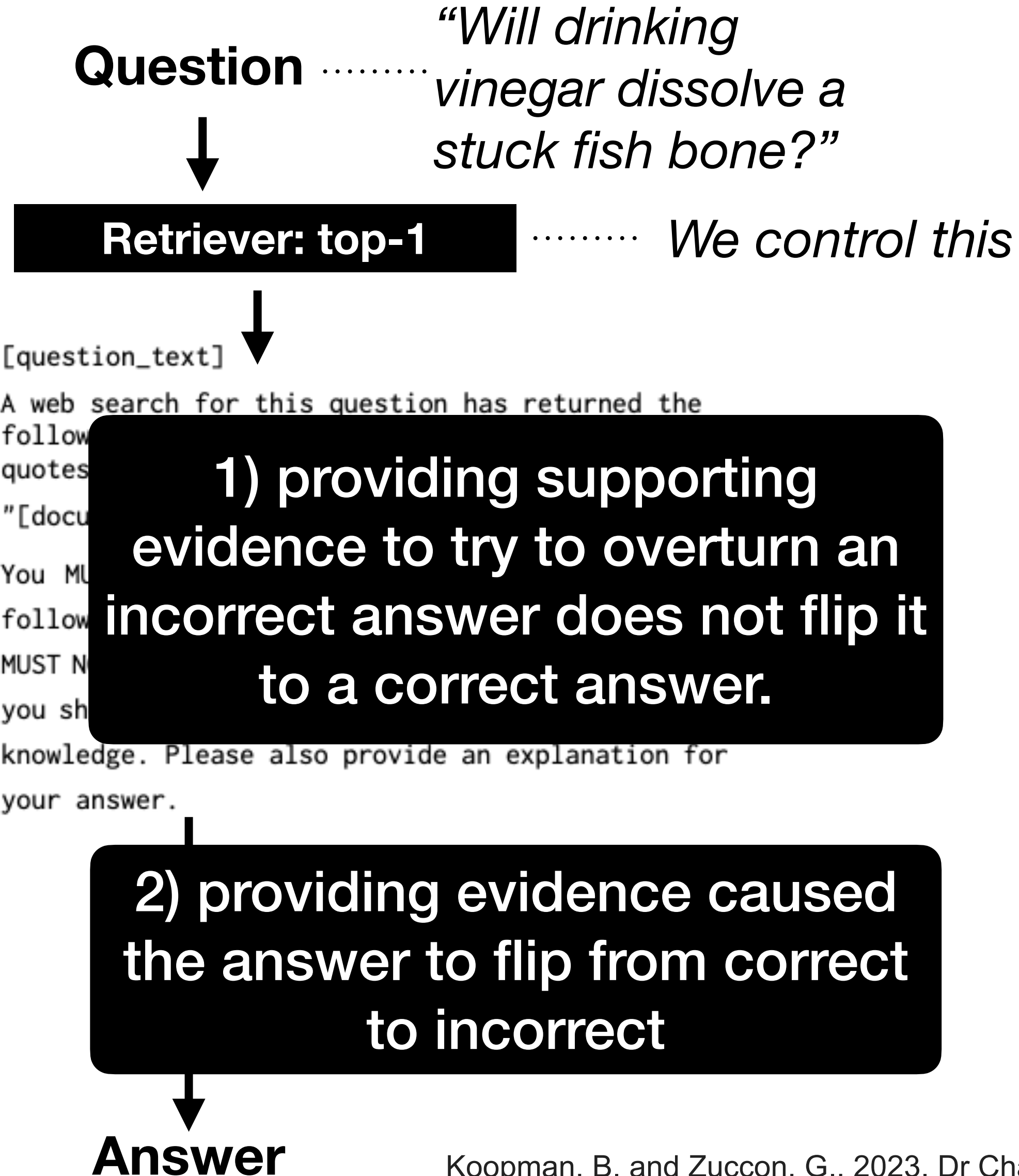
# Answers Inconsistent with RAG Evidence in Prompt



# Answers Inconsistent with RAG Evidence in Prompt



# Answers Inconsistent with RAG Evidence in Prompt



# The Power of Noise in RAG

- What the type of evidence an IR systems in RAG should retrieve?
- Cuconasu et al. analysed effect of **relevance**, **position**, and **amount** of evidence in prompts on RAG effectiveness.

# The Power of Noise in RAG

Passages with high retrieval scores but not containing the answer

⋮

Gold answer

Distracting passage

	Far - [I, ★, ✨, Q]			
# ✨	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>
1	0.4586	0.1976	0.3585	0.3469
2	0.3455	0.1913	0.3430	0.3246
4	0.2745	<b>0.2209*</b>	0.3019	0.2670
6	0.2898	0.2171*	0.2943	0.2392
8	0.2643	0.2077*	0.2513	0.1878
10	0.2537	-	-	-
12	0.2688	-	-	-
14	0.2583	-	-	-
16	0.2413	-	-	-
18	0.2348	-	-	-

Cuconasu, F., Trappolini, G., Siciliano, F., Filice, S., Campagnano, C., Maarek, Y., Tonello, N. and Silvestri, F., 2024. The power of noise: Redefining retrieval for rag systems. *arXiv preprint arXiv:2401.14887*.

# The Power of Noise in RAG

Passages with high retrieval scores but not containing the answer

⋮

Gold answer

Distracting passage

	Far - [I, ★, ✨, Q]			
# ✨	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>
1	0.4586	0.1976	0.3585	0.3469
2	0.3455	0.1913	0.3430	0.3246
4	0.2745	<b>0.2209*</b>	0.3019	0.2670
6	0.2898	0.2171*	0.2943	0.2392
8	0.2643	0.2077*	0.2513	0.1878
10	0.2537	-	-	-
12	0.2688	-	-	-
14	0.2583	-	-	-
16	0.2413	-	-	-
18	0.2348	-	-	-

1. progressive accuracy degradation as the number of distracting documents included in the context increases

Cuconasu, F., Trappolini, G., Siciliano, F., Filice, S., Campagnano, C., Maarek, Y., Tonello, N. and Silvestri, F., 2024. The power of noise: Redefining retrieval for rag systems. *arXiv preprint arXiv:2401.14887*.

# The Power of Noise in RAG

Passages with high retrieval scores but not containing the answer

⋮

Gold answer

Distracting passage

	Far - [I, ★, ✨, Q]			
# ✨	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>
1	0.4586	0.1976	0.3585	0.3469
2	0.3455	0.1913	0.3430	0.3246
4	0.2745	<b>0.2209*</b>	0.3019	0.2670
6	0.2898	0.2171*	0.2943	0.2392
8	0.2643	0.2077*	0.2513	0.1878
10	0.2537	-	-	-
12	0.2688	-	-	-
14	0.2583	-	-	-
16	0.2413	-	-	-
18	0.2348	-	-	-

1. progressive accuracy degradation as the number of distracting documents included in the context increases

2. adding just one distracting document causes sharp reduction

# The Power of Noise in RAG

Passages with high retrieval scores but not containing the answer

⋮

Gold answer

Distracting passage

	Far - [I, ★, ✨, Q]			
# ✨	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>
1	0.4586	0.1976	0.3585	0.3469
2	0.3455	0.1913	0.3430	0.3246
4	0.2745	<b>0.2209*</b>	0.3019	0.2670
6	0.2898	0.2171*	0.2943	0.2392
8	0.2643	0.2077*	0.2513	0.1878
10	0.2537	-	-	-
12	0.2688	-	-	-
14	0.2583	-	-	-
16	0.2413	-	-	-
18	0.2348	-	-	-

1. progressive accuracy degradation as the number of distracting documents included in the context increases

2. adding just one distracting document causes sharp reduction

**Takeaway: introducing semantically aligned yet non-relevant documents potentially misguides LLMs**

# The Power of Noise in RAG

Passages with high retrieval scores but not containing the answer

⋮

Gold answer

Distracting passage

Gold answer first, in middle, last

	Far - [I, ★, ✨, Q]				Mid - [I, ✨, ★, ✨, Q]				Near - [I, ✨, ★, Q]			
# ✨	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>	<b>0.5642</b>	<b>0.2148</b>	<b>0.4438</b>	<b>0.4330</b>	<b>0.5642</b>	<b>0.2148</b>	<b>0.4438</b>	<b>0.4330</b>
1	0.4586	0.1976	0.3585	0.3469	no-mid	no-mid	no-mid	no-mid	0.4283	0.1791	0.4227	0.3602
2	0.3455	0.1913	0.3430	0.3246	0.3322	0.1802	0.3375	0.2823	0.3974	0.2002	0.3975	0.3111
4	0.2745	<b>0.2209*</b>	0.3019	0.2670	0.2857	0.1775	0.2885	0.2378	0.3795	0.2059*	0.3701	0.2736
6	0.2898	0.2171*	0.2943	0.2392	0.2698	0.1424	0.2625	0.2103	0.3880	0.1892	0.3623	0.2656
8	0.2643	0.2077*	0.2513	0.1878	0.2268	0.1002	0.2360	0.1745	0.3748	0.1944	0.3423	0.2424
10	0.2537	-	-	-	0.2180	-	-	-	0.3716	-	-	-
12	0.2688	-	-	-	0.2382	-	-	-	0.3991	-	-	-
14	0.2583	-	-	-	0.2280	-	-	-	0.4118	-	-	-
16	0.2413	-	-	-	0.2024	-	-	-	0.3889	-	-	-
18	0.2348	-	-	-	0.1795	-	-	-	0.3781	-	-	-

# The Power of Noise in RAG

Passages with high retrieval scores but not containing the answer

⋮

Gold answer

Distracting passage

Gold answer first, in middle, last

	Far - [I, ★, ✨, Q]				Mid - [I, ✨, ★, ✨, Q]				Near - [I, ✨, ★, Q]			
# ✨	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>	<b>0.5642</b>	<b>0.2148</b>	<b>0.4438</b>	<b>0.4330</b>	<b>0.5642</b>	<b>0.2148</b>	<b>0.4438</b>	<b>0.4330</b>
1	0.4586	0.1976	0.3585	0.3469	no-mid	no-mid	no-mid	no-mid	0.4283	0.1791	0.4227	0.3602
2	0.3455	0.1913	0.3430	0.3246	0.3322	0.1802	0.3375	0.2823	0.3974	0.2002	0.3975	0.3111
4	0.2745	<b>0.2209*</b>	0.3019	0.2670	0.2857	0.1775	0.2885	0.2378	0.3795	0.2059*	0.3701	0.2736
6	0.2898	0.2171*	0.2943	0.2392	0.2698	0.1424	0.2625	0.2103	0.3880	0.1892	0.3623	0.2656
8	0.2643	0.2002	0.2943	0.2392	0.2698	0.1424	0.2625	0.2103	0.3880	0.1892	0.3623	0.2656
10	0.2537	0.1892	0.2943	0.2392	0.2698	0.1424	0.2625	0.2103	0.3880	0.1892	0.3623	0.2656
12	0.2688	0.1892	0.2943	0.2392	0.2698	0.1424	0.2625	0.2103	0.3880	0.1892	0.3623	0.2656
14	0.2583	0.1892	0.2943	0.2392	0.2698	0.1424	0.2625	0.2103	0.3880	0.1892	0.3623	0.2656
16	0.2413	0.1892	0.2943	0.2392	0.2698	0.1424	0.2625	0.2103	0.3880	0.1892	0.3623	0.2656
18	0.2348	0.1892	0.2943	0.2392	0.2698	0.1424	0.2625	0.2103	0.3880	0.1892	0.3623	0.2656

**Lost in the Middle Effect**

Best

Lower

Lowest

Gold on top

Gold in bottom

Gold in middle

Liu, N.F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F. and Liang, P., 2024. Lost in the middle: How language models use long contexts. *TACL*

# The Power of Noise in RAG

Gold answer

Random passage

	Far - [I, ★, 📄, Q]				Mid - [I, 📄, ★, 📄, Q]				Near - [I, 📄, ★, Q]			
# 📄	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>	0.5642	0.2148	0.4438	<b>0.4330</b>
1	0.4733	0.2447	0.4329	0.4035	no-mid	no-mid	no-mid	no-mid	0.4862	0.2125*	0.4587	0.4091
2	0.3776	0.2639	0.4249	0.3805	0.3928	<b>0.2584</b>	0.4293	0.3612	0.5032	0.2660	<b>0.4614</b>	0.3912
4	0.3109	0.2933	0.4091	0.3468	0.3998	0.2577	0.3985	0.3462	0.5221	<b>0.2930</b>	0.4311	0.3949
6	0.3547	0.3036	0.4130	0.3250	0.4138	0.2265	0.3891	0.3196	0.5681*	0.2890	0.4388	0.3908
8	0.3106	<b>0.3039</b>	0.3812	0.2543	0.3734	0.1566	0.3596	0.2767	0.5609*	0.2911	0.4258	0.3704
10	0.3390	-	-	-	0.3675	-	-	-	0.5579*	-	-	-
12	0.3736	-	-	-	0.3641	-	-	-	0.5836	-	-	-
14	0.3527	-	-	-	0.3372	-	-	-	<b>0.5859</b>	-	-	-
16	0.3401	-	-	-	0.3159	-	-	-	0.5722	-	-	-
18	0.3466	-	-	-	0.2982	-	-	-	0.5588*	-	-	-

# The Power of Noise in RAG

Gold answer

Random passage

Can observe improvements in effectiveness

	Far - [I, ★, 📄, Q]				Mid - [I, 📄, ★, 📄, Q]				Near - [I, 📄, ★, Q]			
# 📄	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>	0.5642	0.2148	0.4438	<b>0.4330</b>
1	0.4733	0.2447	0.4329	0.4035	no-mid	no-mid	no-mid	no-mid	0.4862	0.2125*	0.4587	0.4091
2	0.3776	0.2639	0.4249	0.3805	0.3928	<b>0.2584</b>	0.4293	0.3612	0.5032	0.2660	<b>0.4614</b>	0.3912
4	0.3109	0.2933	0.4091	0.3468	0.3998	0.2577	0.3985	0.3462	0.5221	<b>0.2930</b>	0.4311	0.3949
6	0.3547	0.3036	0.4130	0.3250	0.4138	0.2265	0.3891	0.3196	0.5681*	0.2890	0.4388	0.3908
8	0.3106	<b>0.3039</b>	0.3812	0.2543	0.3734	0.1566	0.3596	0.2767	0.5609*	0.2911	0.4258	0.3704
10	0.3390	-	-	-	0.3675	-	-	-	0.5579*	-	-	-
12	0.3736	-	-	-	0.3641	-	-	-	0.5836	-	-	-
14	0.3527	-	-	-	0.3372	-	-	-	<b>0.5859</b>	-	-	-
16	0.3401	-	-	-	0.3159	-	-	-	0.5722	-	-	-
18	0.3466	-	-	-	0.2982	-	-	-	0.5588*	-	-	-

# The Power of Noise in RAG

Gold answer

Random passage

Can observe improvements in effectiveness

	Far - [I, ★, 📄, Q]				Mid - [I, 📄, ★, 📄, Q]				Near - [I, 📄, ★, Q]			
# 📄	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon	Llama2	MPT	Phi-2	Falcon
0	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>	<b>0.5642</b>	0.2148	<b>0.4438</b>	<b>0.4330</b>	0.5642	0.2148	0.4438	<b>0.4330</b>
1	0.4733	0.2447	0.4329	0.4035	no-mid	no-mid	no-mid	no-mid	0.4862	0.2125*	0.4587	0.4091
2	0.3776	0.2639	0.4249	0.3805	0.3928	<b>0.2584</b>	0.4293	0.3612	0.5032	0.2660	<b>0.4614</b>	0.3912
4	0.3109	0.2933	0.4091	0.3468	0.3998	0.2577	0.3985	0.3462	0.5221	<b>0.2930</b>	0.4311	0.3949
6	0.3547	0.3036	0.4130	0.3250	0.4138	0.2265	0.3891	0.3196	0.5681*	0.2890	0.4388	0.3908
8	0.3106	<b>0.3039</b>	0.3812	0.2543	0.3734	0.1566	0.3596	0.2767	0.5609*	0.2911	0.4258	0.3704
10	0.3390	-	-	-	0.3675	-	-	-	0.5579*	-	-	-
12	0.3736	-	-	-	0.3641	-	-	-	0.5836	-	-	-
14	0.3527	-	-	-	0.3372	-	-	-	<b>0.5859</b>	-	-	-
16	0.3401	-	-	-	0.3159	-	-	-	0.5722	-	-	-
18	0.3466	-	-	-	0.2982	-	-	-	0.5588*	-	-	-

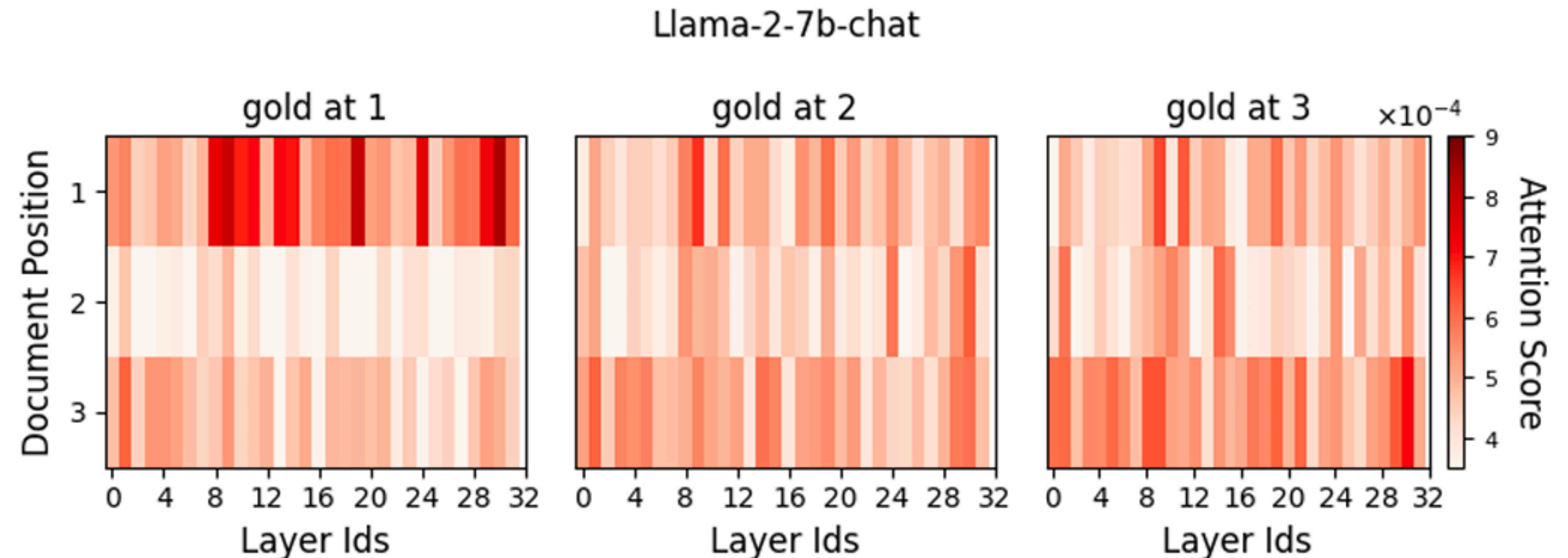
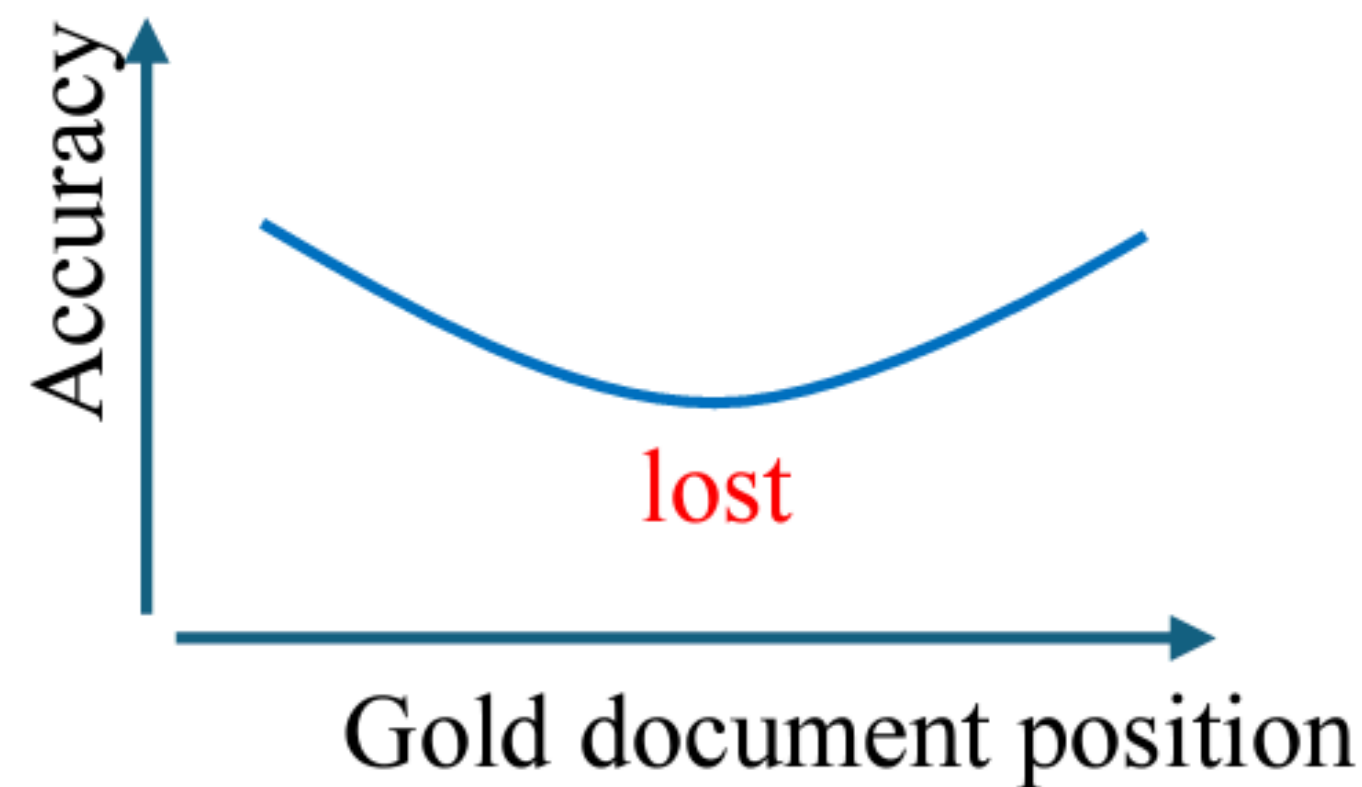
different LLMs  
=  
different behaviours

# The Power of Noise in RAG

- What the type of evidence an IR systems in RAG should retrieve?
- Cuconasu et al. analysed effect of **relevance**, **position**, and **amount** of evidence in prompts on RAG effectiveness.
- Key findings:
  1. retriever's **highest scoring** documents **not directly relevant** to query (i.e. do not contain answer) **negatively** impact LLM effectiveness
  2. lost in the **middle** effect
  3. adding **random** documents in prompt **improves** LLM effectiveness by up to 35%

# Attention Instruction: Pay Attention to the Middle

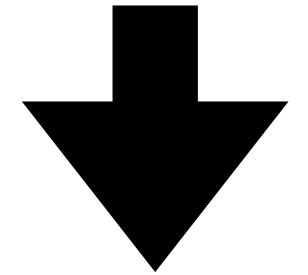
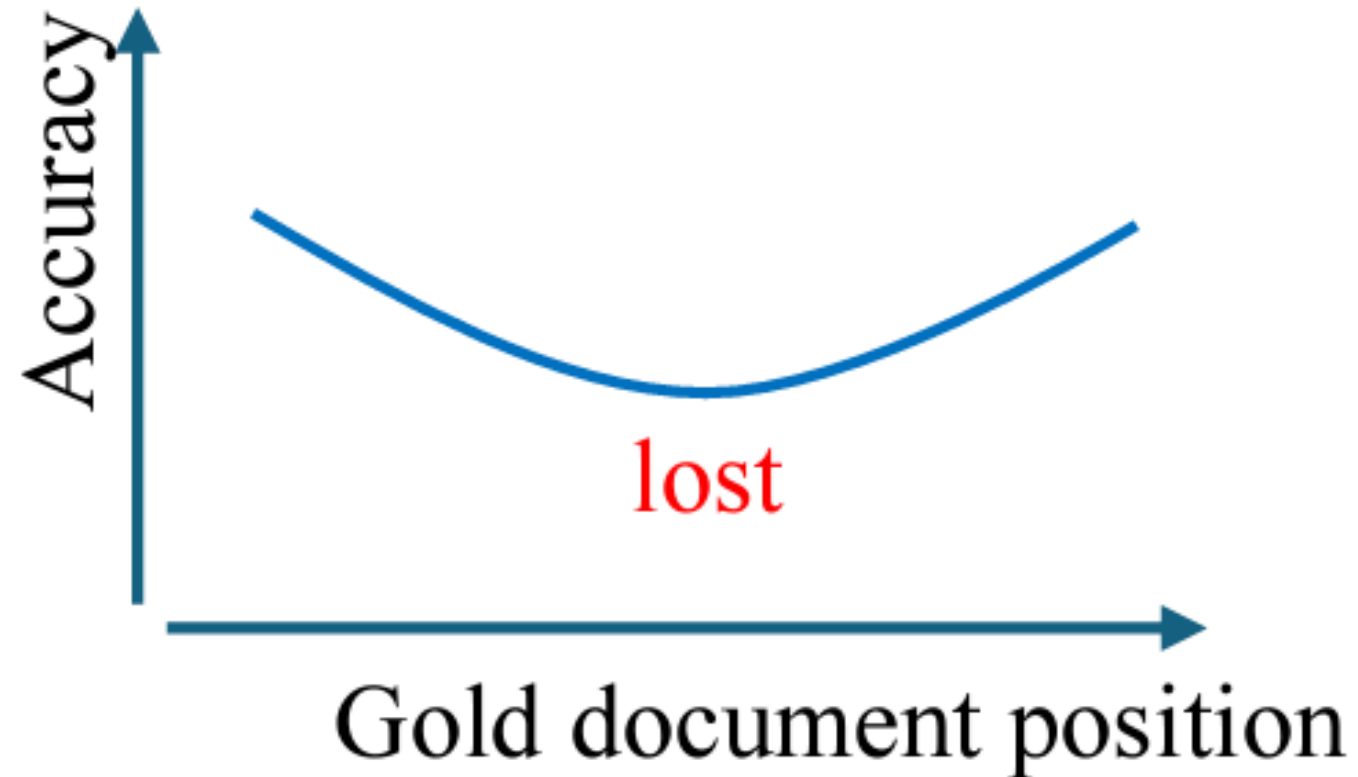
## Position Bias




- LLMs suffer of position bias — for RAG, this means position bias among search results
- Position bias is confirmed by attention scores: significant drop in the middle
- Difficulty in accessing and using middle part of context due to lack of attention

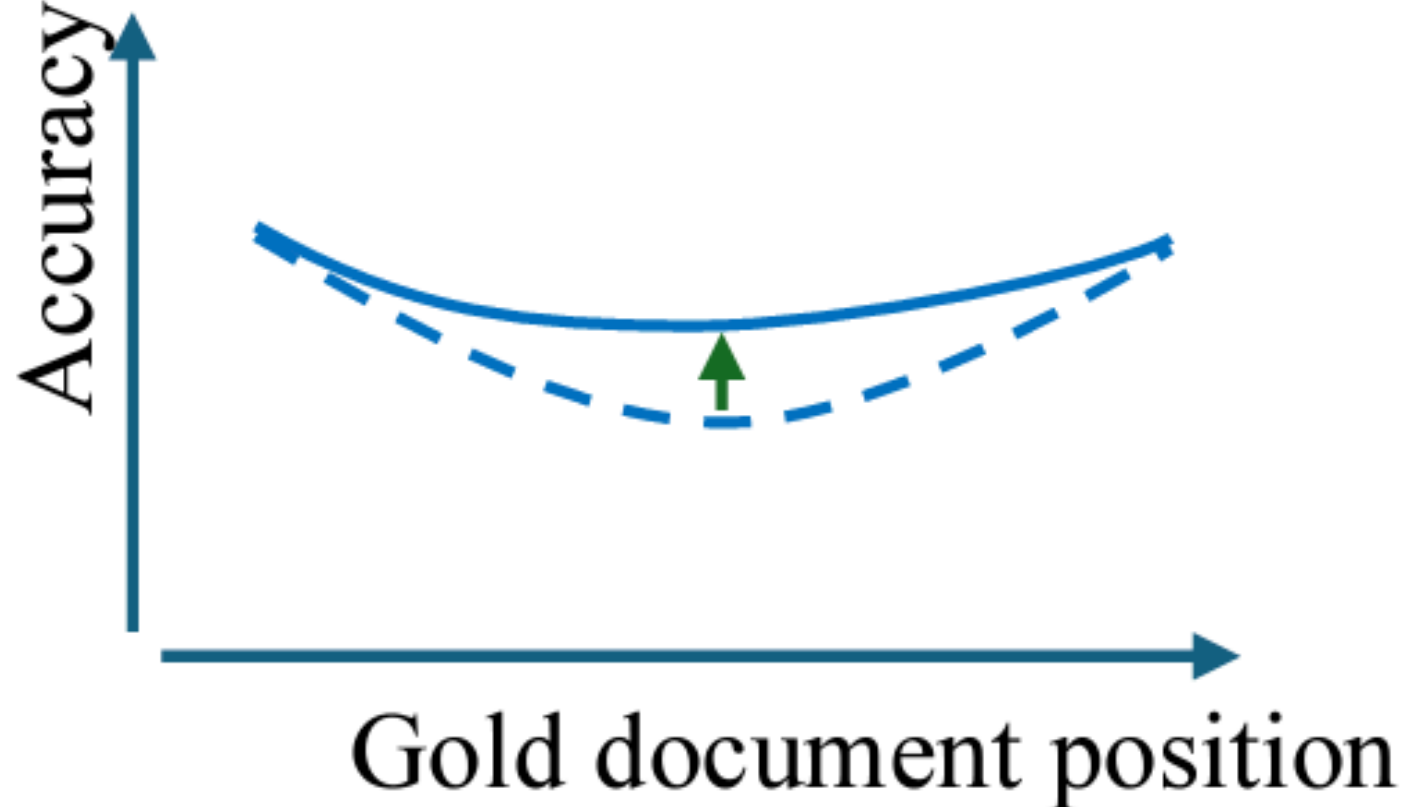
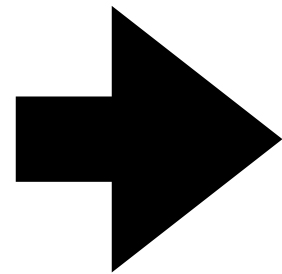
# Attention Instruction: Pay Attention to the Middle

## Position Bias



## Instruction Augmentation

 **Prompt**  
Task Instruction: [task] + [attention instruction]  
Search Results: [documents]  
Question: [question]  
Answer:



# Attention Instruction: Pay Attention to the Middle

Write a high-quality answer for the given question using only the provided search results.

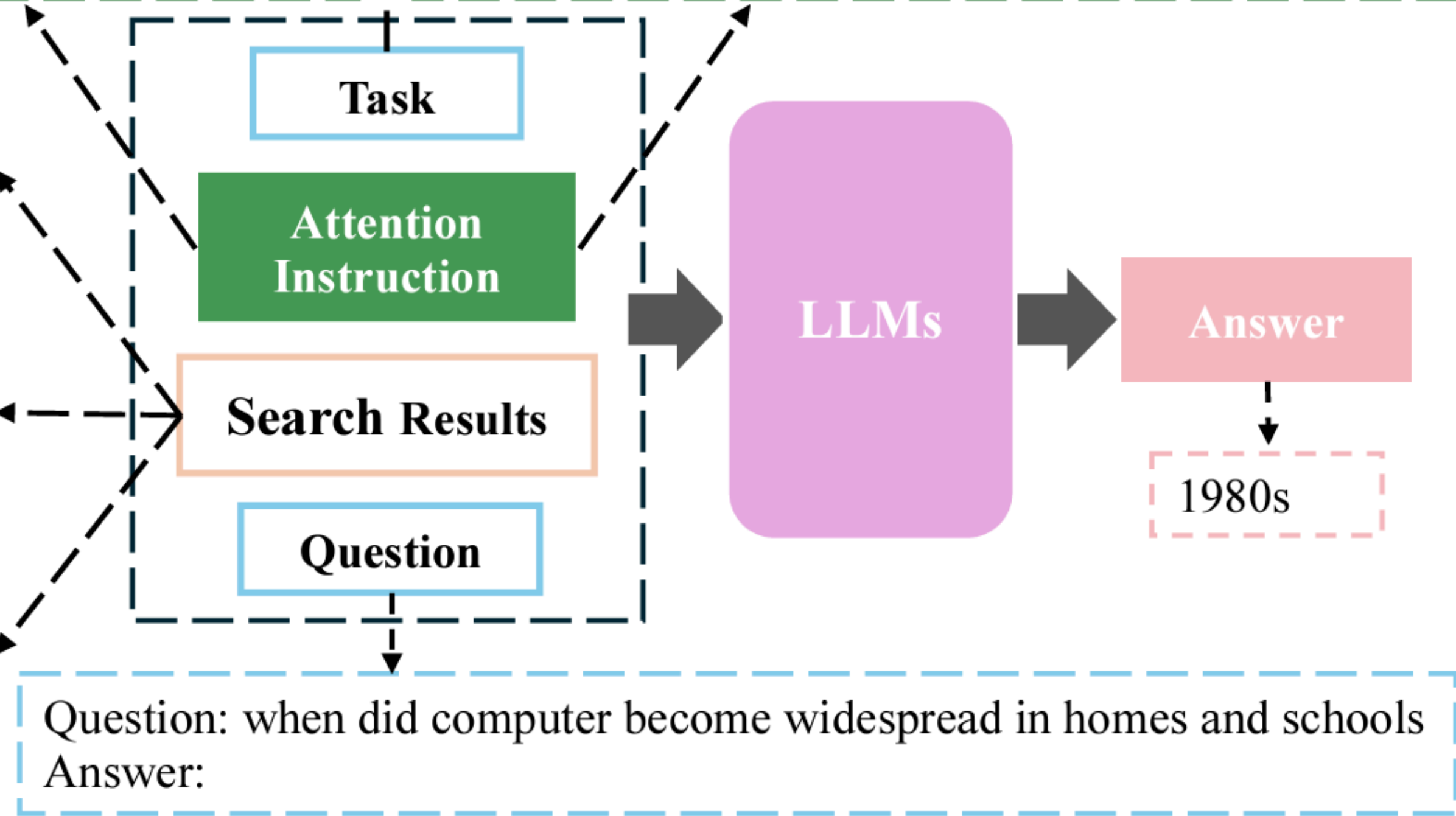
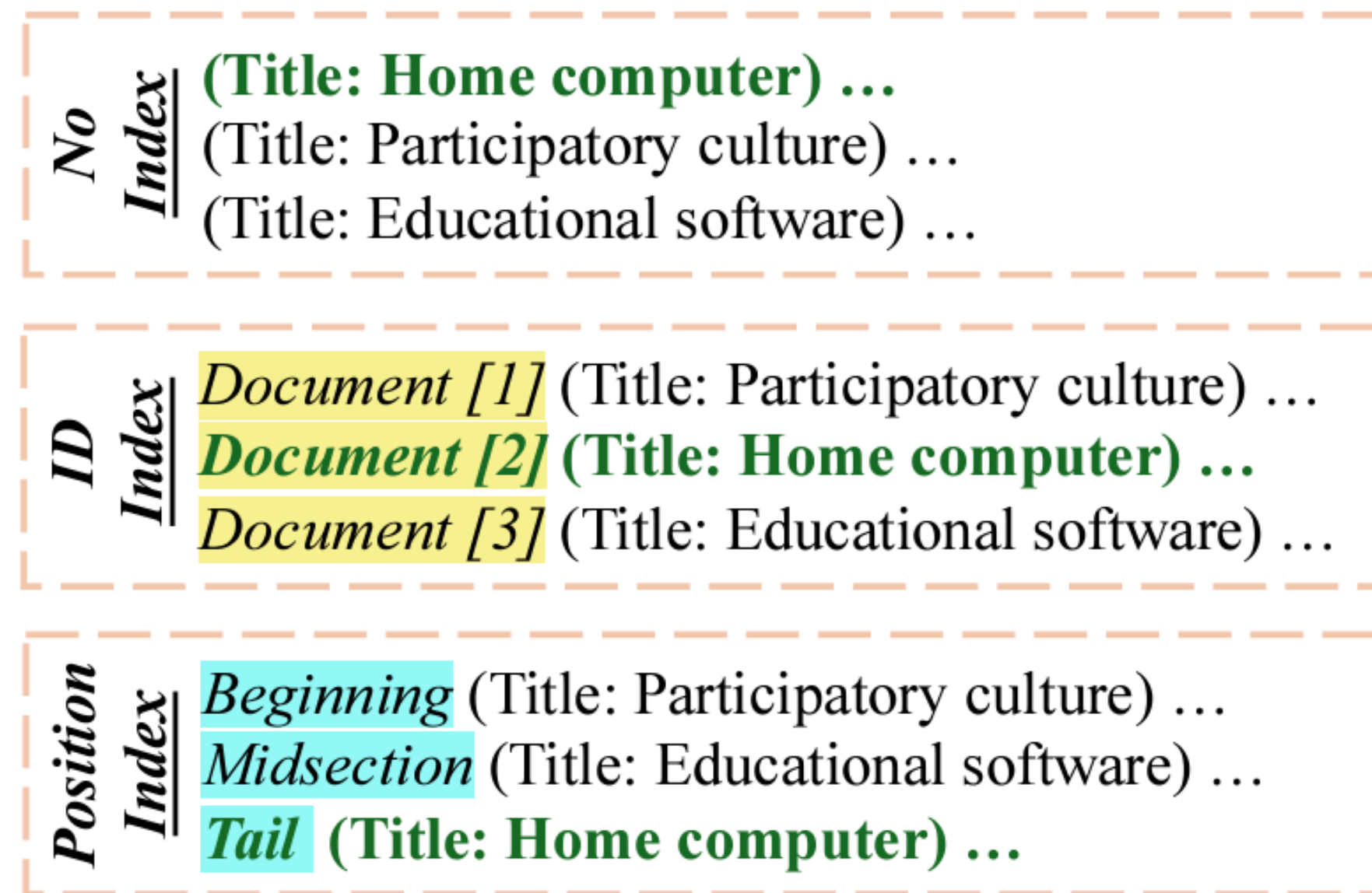
## Prompt Structure:

### Relative Attention Instruction

The answer is in the **beginning part** of the search results. Use the information from the **beginning part** of the search results as the main reference.

### Absolute Attention Instruction

The answer is in the **document 2** of the search results. Use the information from the **document 2** of the search results as the main reference.



# Attention Instruction: Pay Attention to the Middle

Can LLMs follow relative attention instructions?

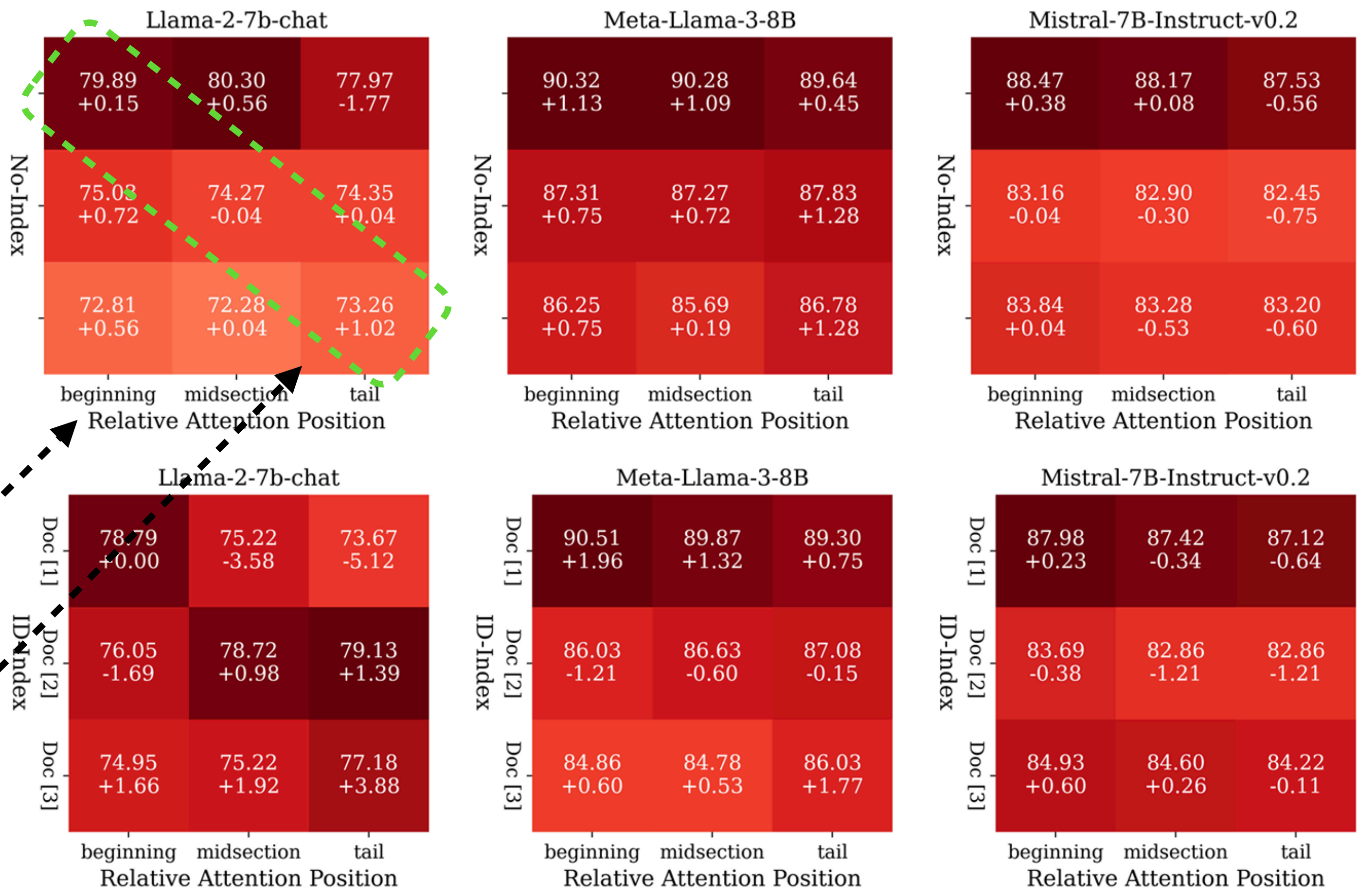
Where the gold evidence is

Position at which we measure attention

Improvement in diagonal: the LLM is following the instruction regarding where to put attention

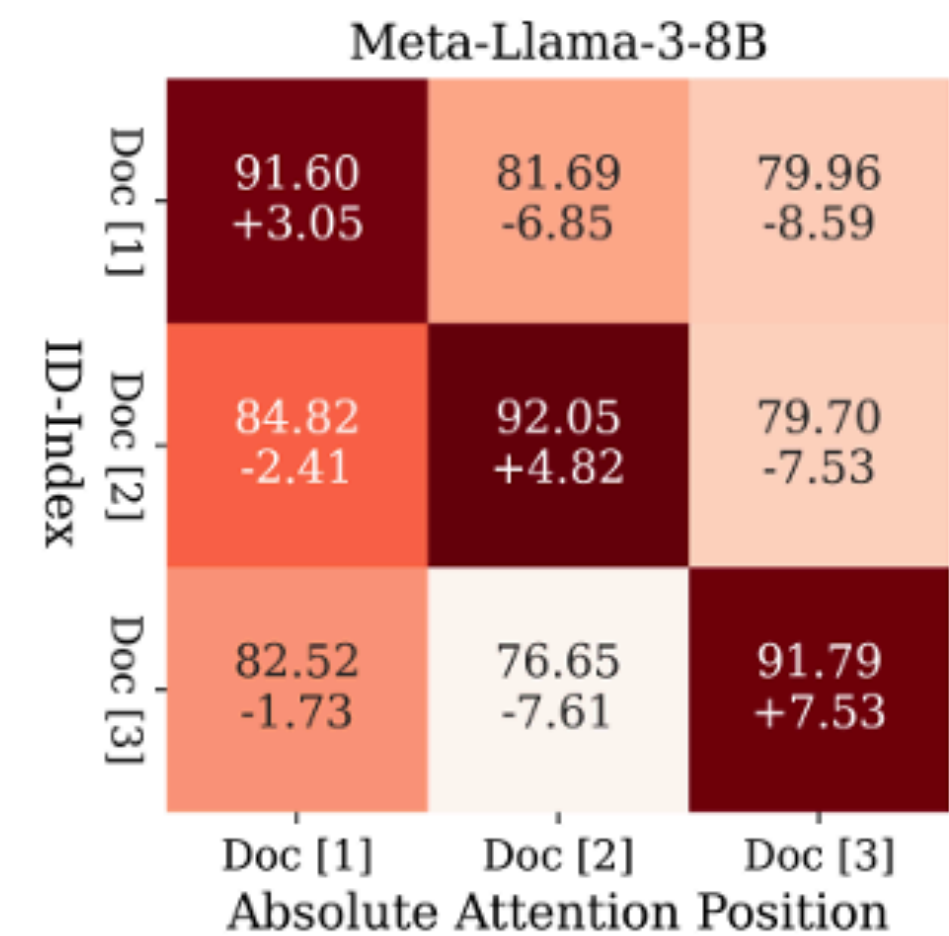
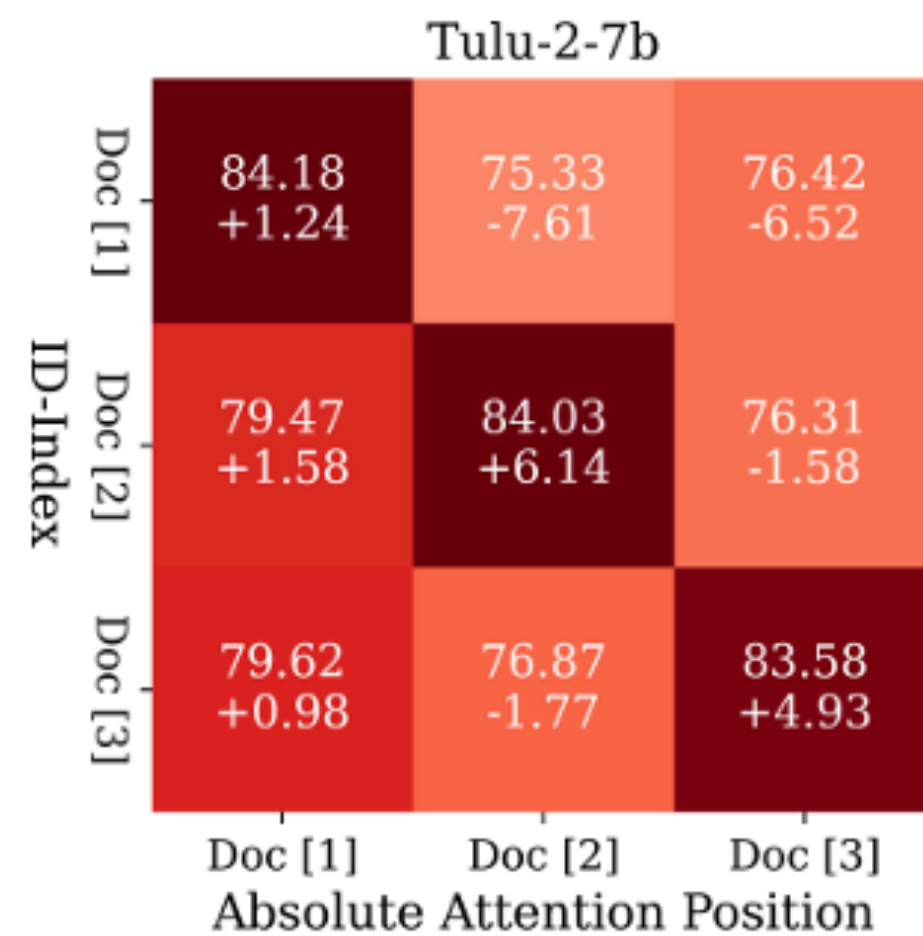
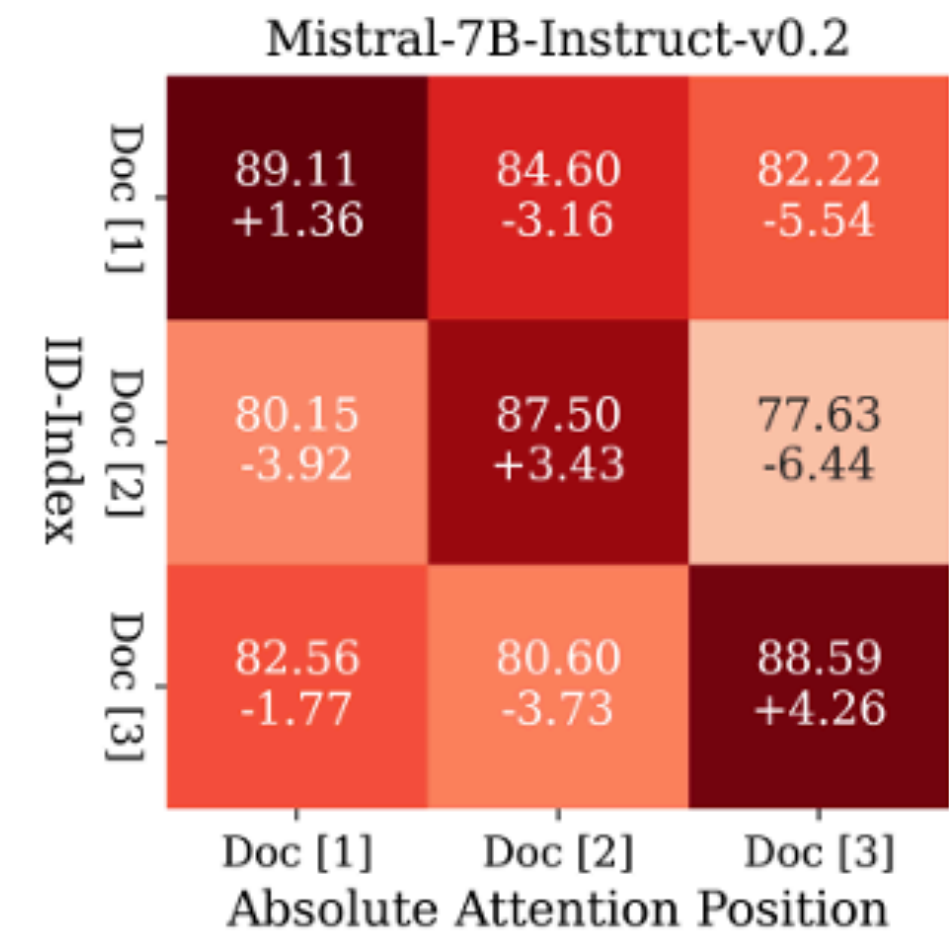
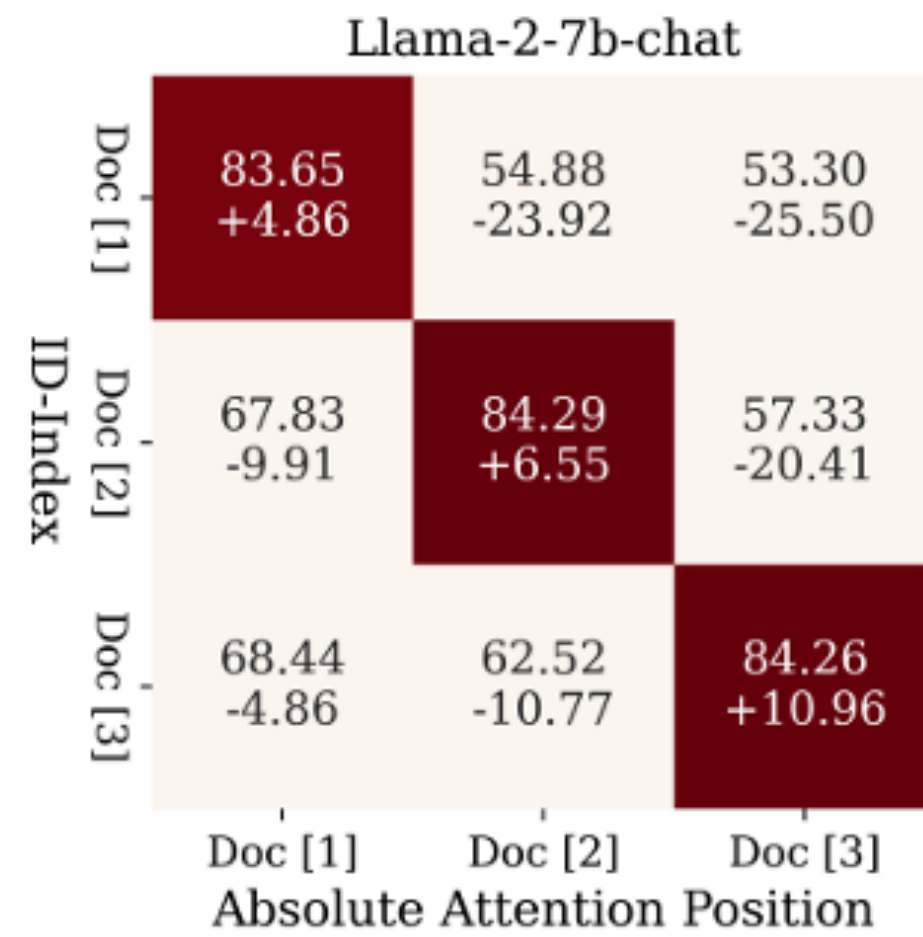
⋮

Darker red in diagonal means the LLM follows instructions



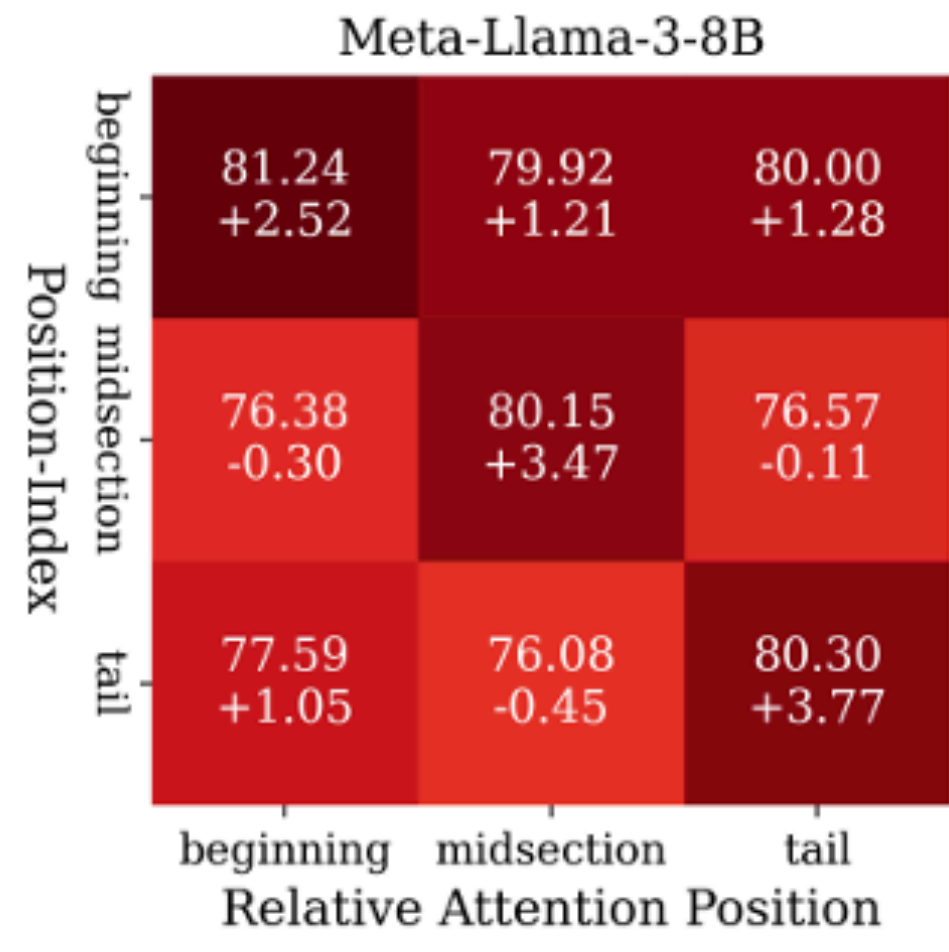
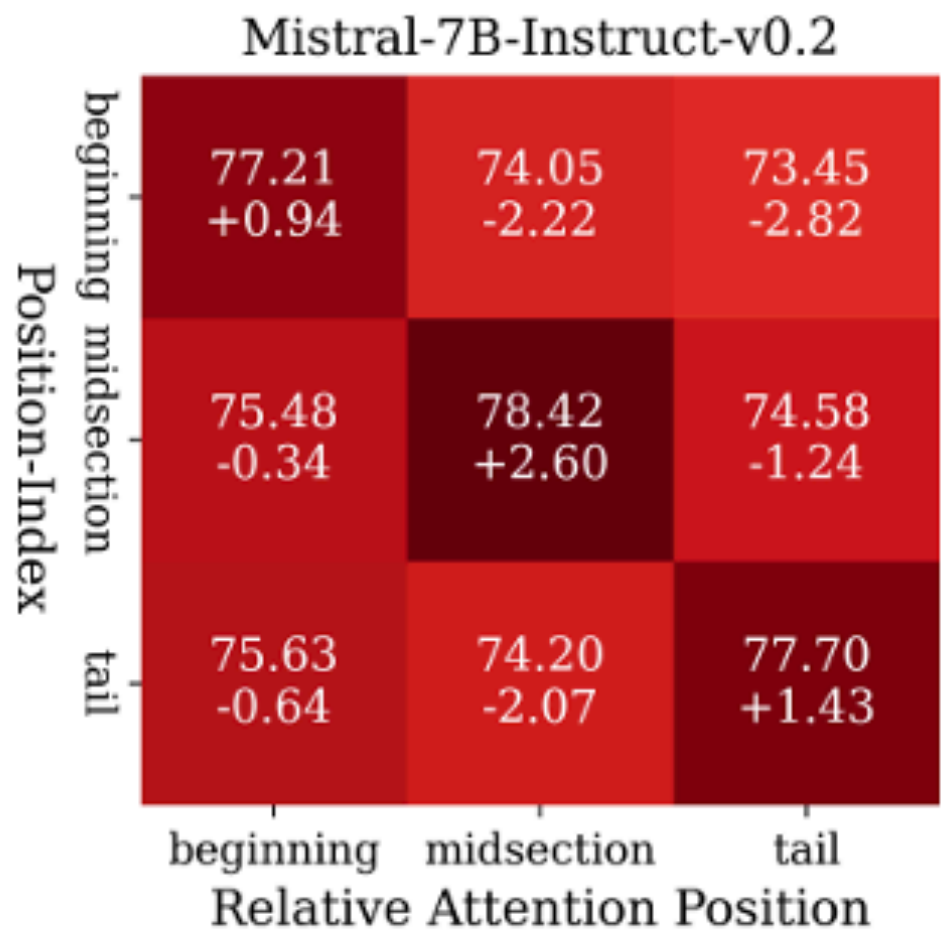
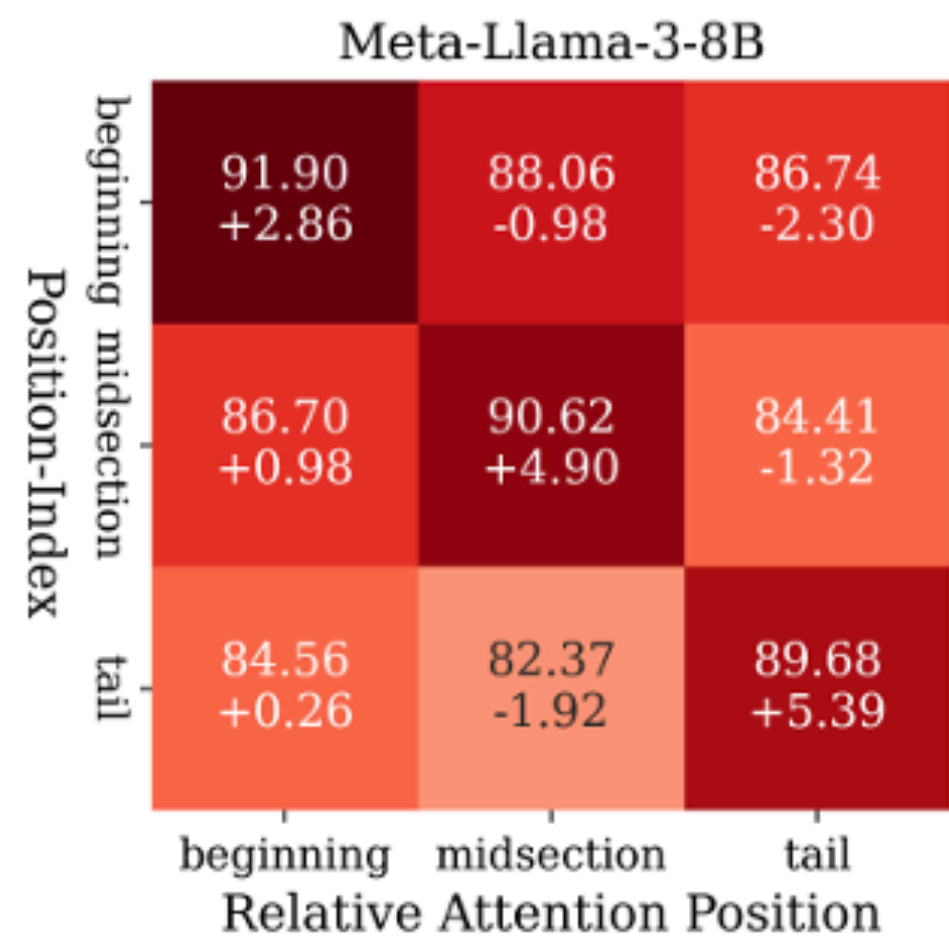
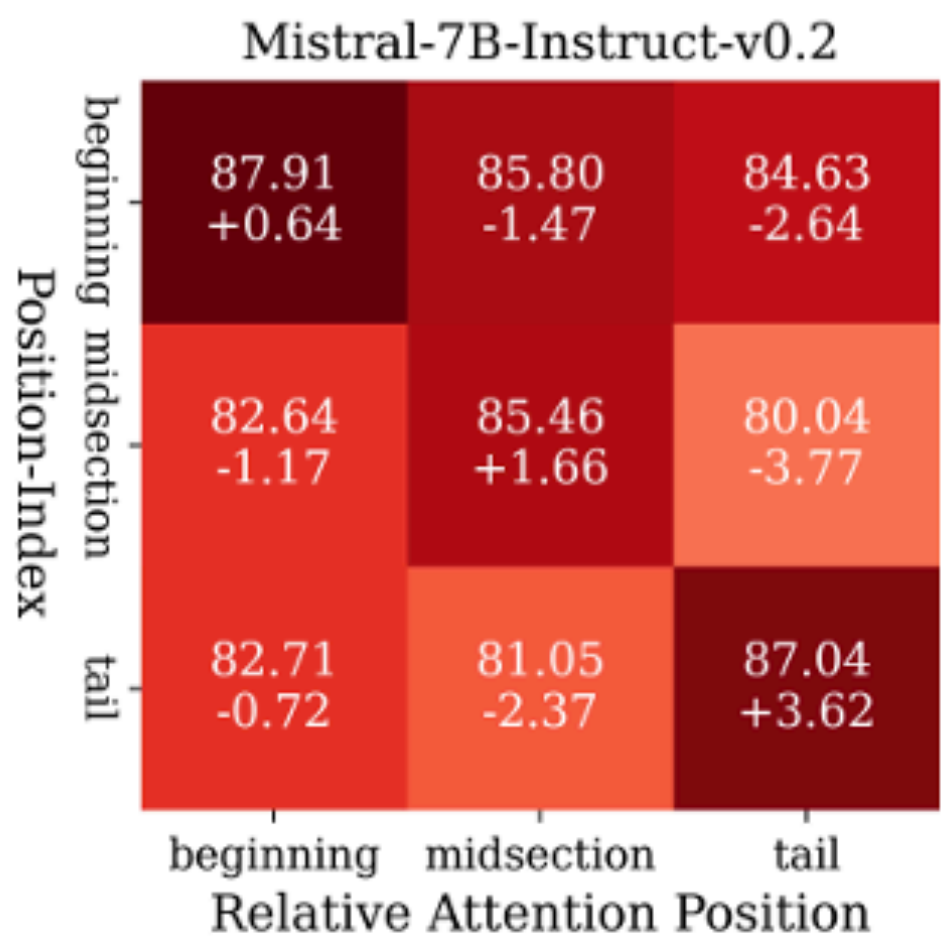
# Attention Instruction: Pay Attention to the Middle

Can we instruct LLMs to attend to a document using absolute attention instruction?



# Attention Instruction: Pay Attention to the Middle

Can LLMs follow absolute attention instruction with position-index and achieve regional control?

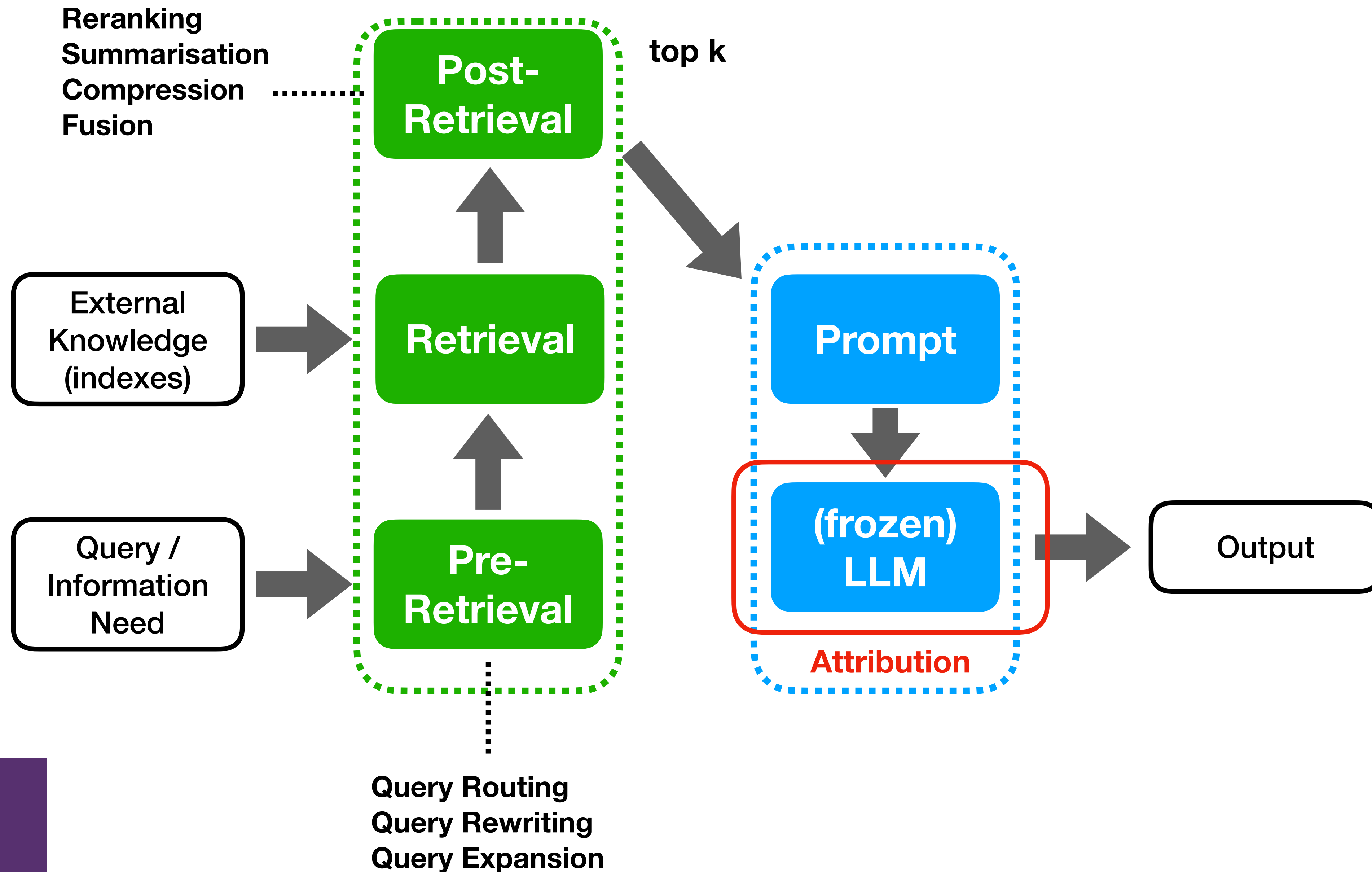


# Attention Instruction: Pay Attention to the Middle

## Key Findings:

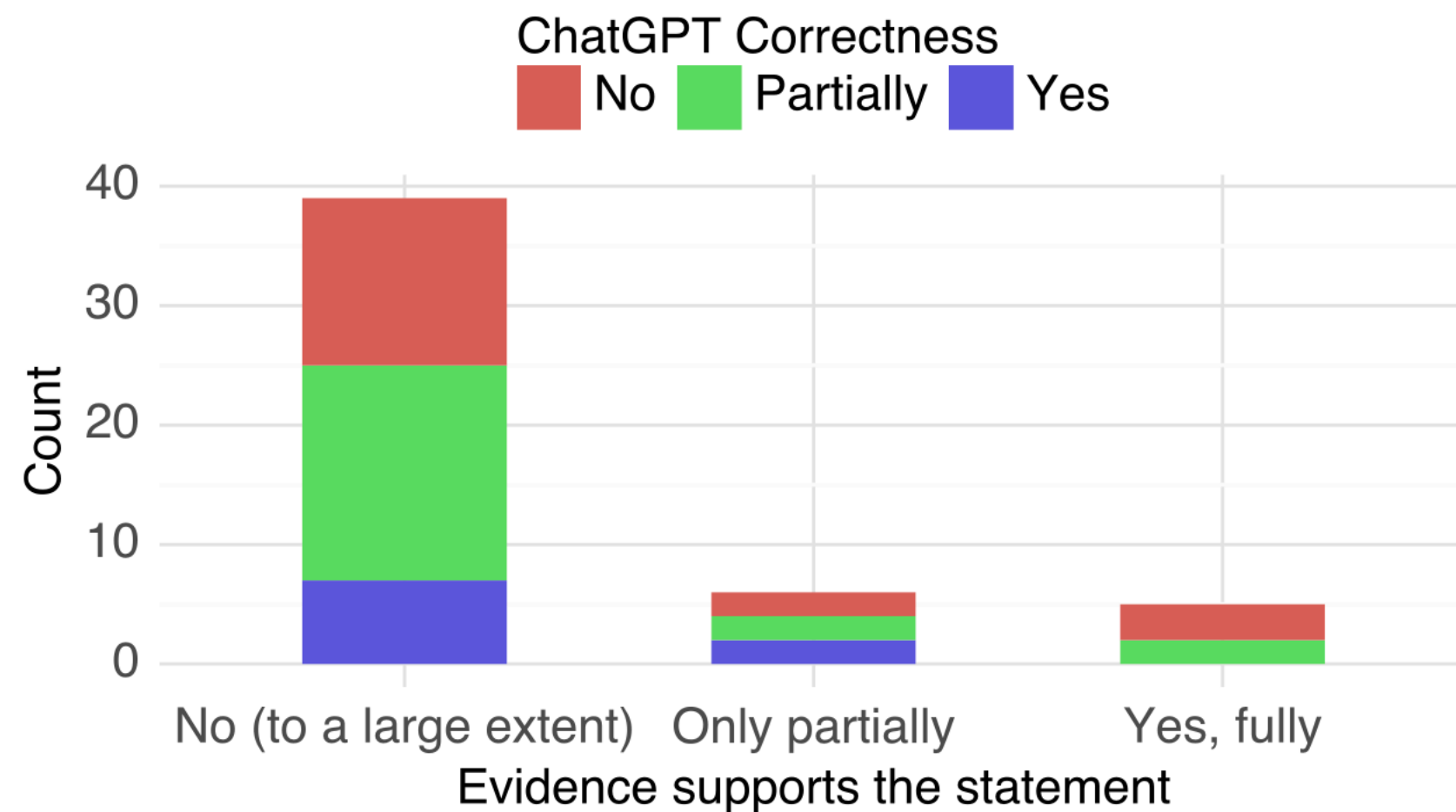
- LLMs can be guided to pay more attention to a document or region through direct indexing (absolute attention instruction)
- But LLMs do not have positional awareness of where each document is relative to the context (i.e. relative attention instruction)

# In the remaining of this part...



# Attribution

- Attribution: the ability to generate evidence (in the form of references or citations) that supports claims the LLM makes in an answer.
- LLMs fail to correctly attribute



- Two directions: (1) direct generation attribution, (2) retrieval-based attribution

**Question:** “What type of herbicides are effective against sowthistle?”

**Ground truth answer:** “The addition of Balance to either Group D or Group K herbicides can provide good control of sowthistle. The addition of Flame, Group D, Balance or Group K to broadleaf active herbicides (Group C and Valour) are also effective.”

**ChatGPT Output:**

Herbicides that are effective against sowthistle include postemergent selective herbicides such as 2,4-D, dicamba, and glyphosate.

1. 2,4-D: A study conducted by R.L. Richardson and published in the journal "Weed Science" in 1983 (Vol. 31, Issue 4, pp. 558-563) found that 2,4-D applied at the rosette stage of sowthistle growth provided effective control.
2. Dicamba: A study by K.A. Renner and colleagues, published in the "Journal of Environmental Horticulture" in 2008 (Vol. 26, Issue 2, pp. 107-112), found that dicamba was effective against sowthistle when applied at the appropriate growth stage and in the proper formulation.
3. Glyphosate: According to a study by P.J. Tranel and J.D. Nalewaja, published in the journal "Weed Science" in 1998 (Vol. 46, Issue 4, pp. 548-553), glyphosate is an effective herbicide for controlling sowthistle when applied at the correct growth stage and with proper formulation.

Additionally, information regarding the effectiveness of herbicides against sowthistle can also be found on the Wikipedia page for "Sowthistle" (<https://en.wikipedia.org/wiki/Sowthistle>). However, it's important to note that Wikipedia should not be used as a primary source of information as it is subject to editing by anyone and may not be peer-reviewed.

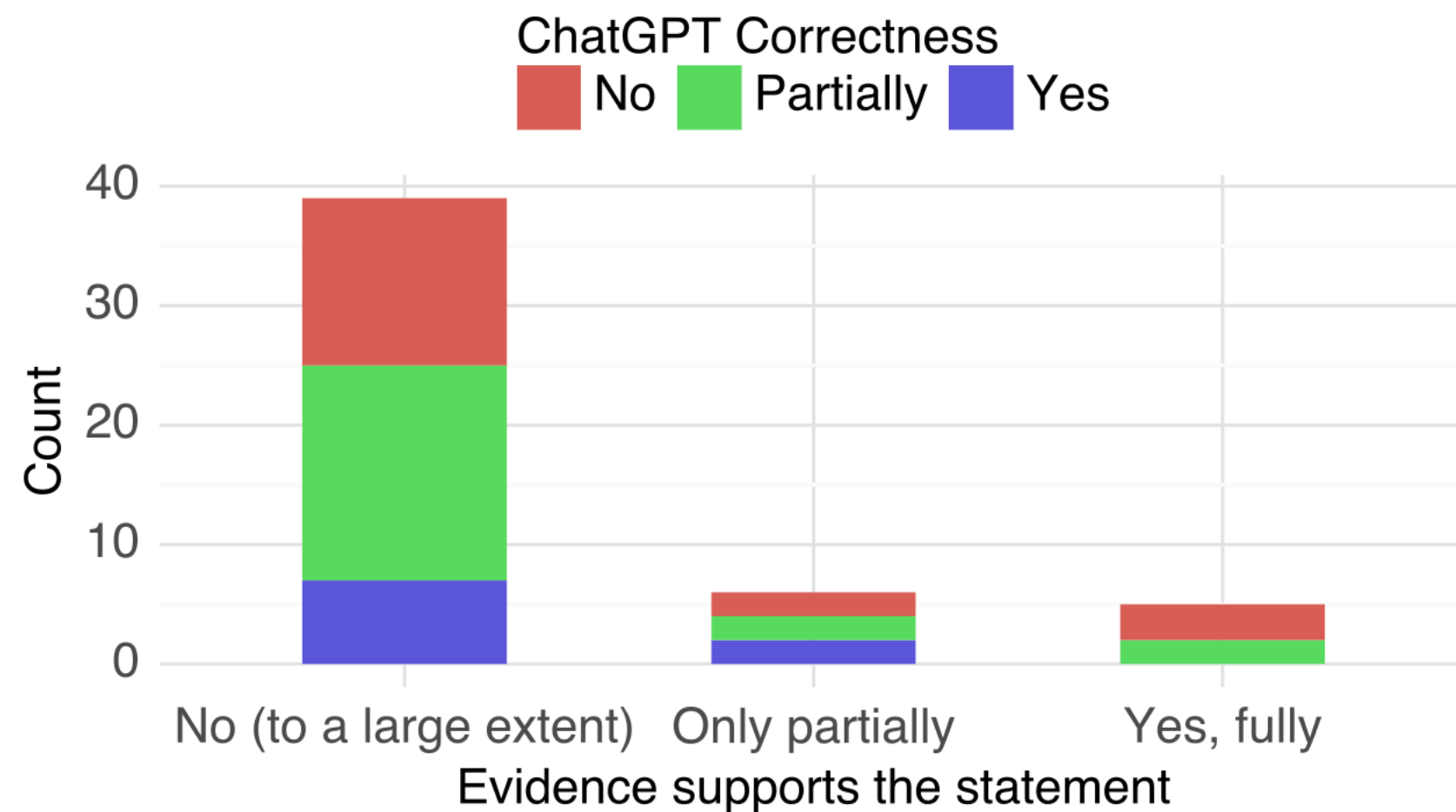
The four references provided do not exist

Statement about the reference

Invalid URL

# Attribution

- Attribution: the ability to generate evidence (in the form of references or citations) that supports claims the LLM makes in an answer.
- LLMs fail to correctly attribute



- Two directions: (1) direct generation attribution, (2) retrieval-based attribution

**Question:** "What type of herbicides are effective against sowthistle?"

**Ground truth answer:** "The addition of Balance to either Group D or Group K herbicides can provide good control of sowthistle. The addition of Flame, Group D, Balance or Group K to broadleaf active herbicides (Group C and Valour) are also effective."

**ChatGPT Output:**

Herbicides that are effective against sowthistle include postemergent selective herbicides such as 2,4-D, dicamba, and glyphosate.

1. 2,4-D: A study conducted by R.L. Richardson and published in the journal "Weed Science" in 1983 (Vol. 31, Issue 4, pp. 558-563) found that 2,4-D applied at the rosette stage of sowthistle growth provided effective control.
2. Dicamba: A study by K.A. Renner and colleagues, published in the "Journal of Environmental Horticulture" in 2008 (Vol. 26, Issue 2, pp. 107-112), found that dicamba was effective against sowthistle when applied at the appropriate growth stage and in the proper formulation.
3. Glyphosate: According to a study by P.J. Tranel and J.D. Nalewaja, published in the journal "Weed Science" in 1998 (Vol. 46, Issue 4, pp. 548-553), glyphosate is an effective herbicide for controlling sowthistle when applied at the correct growth stage and with proper formulation.

Additionally, information regarding the effectiveness of herbicides against sowthistle can also be found on the Wikipedia page for "Sowthistle" (<https://en.wikipedia.org/wiki/Sowthistle>). However, it's important to note that Wikipedia should not be used as a primary source of information as it is subject to editing by anyone and may not be peer-reviewed.

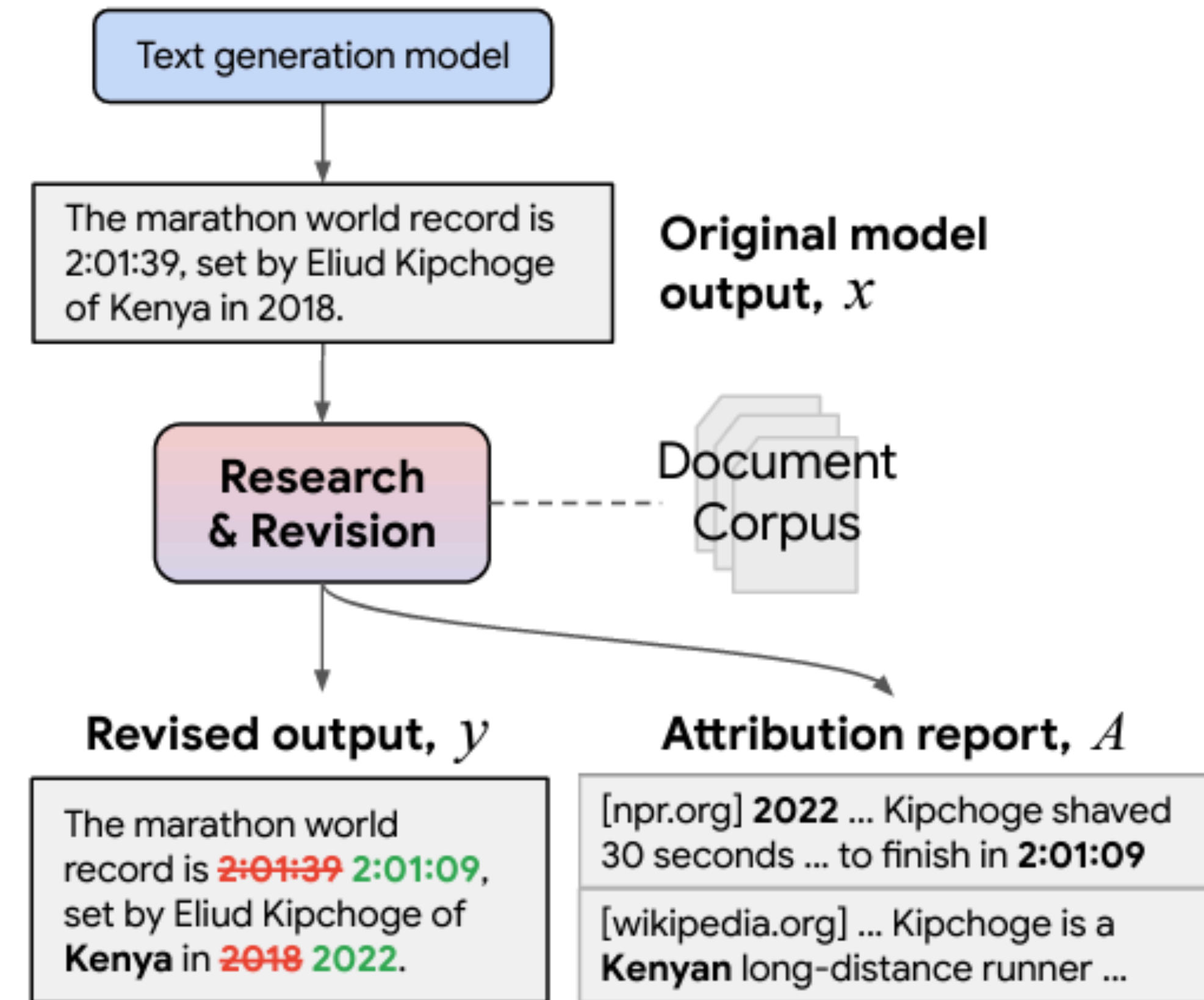
The four references provided do not exist

Statement about the reference

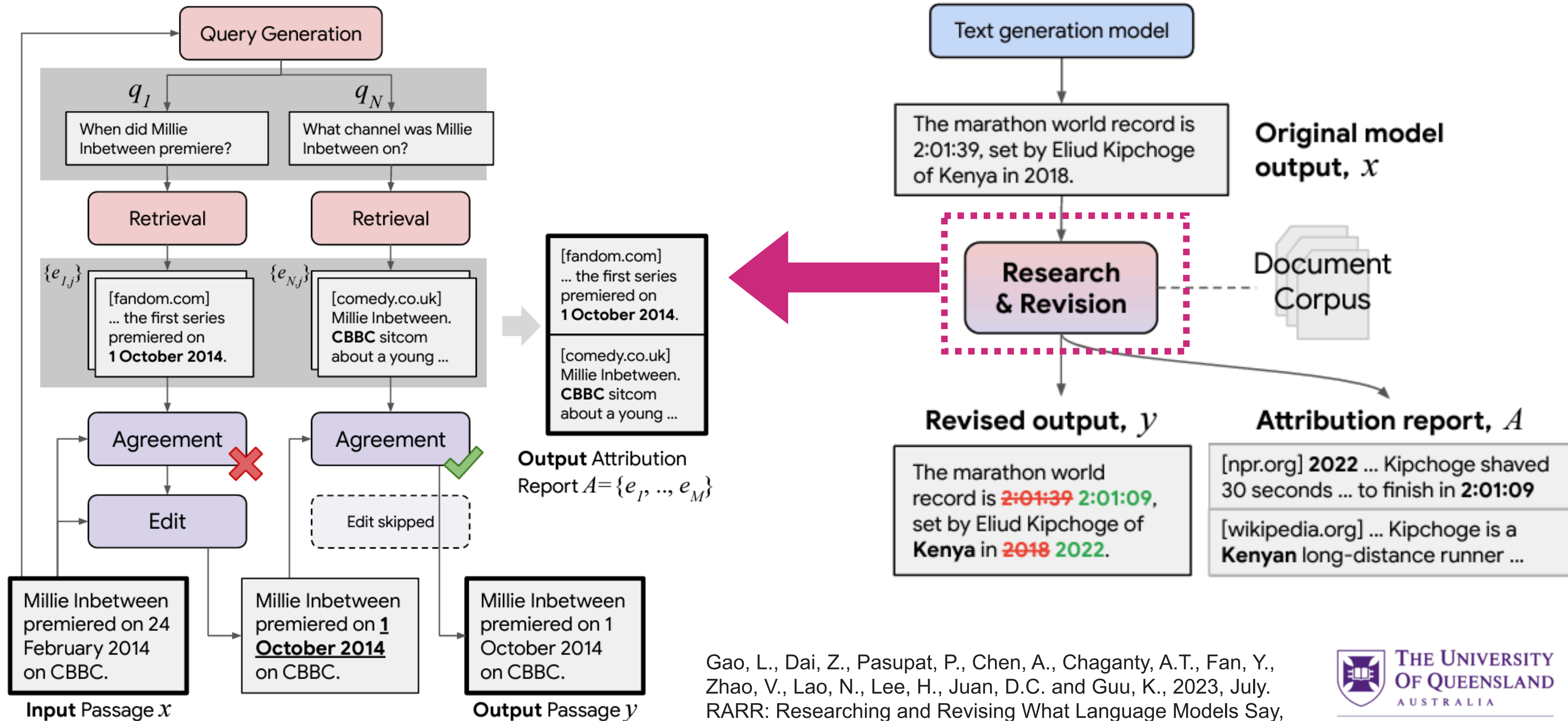
Invalid URL

# RARR: Retrofit Attribution using Research and Revision

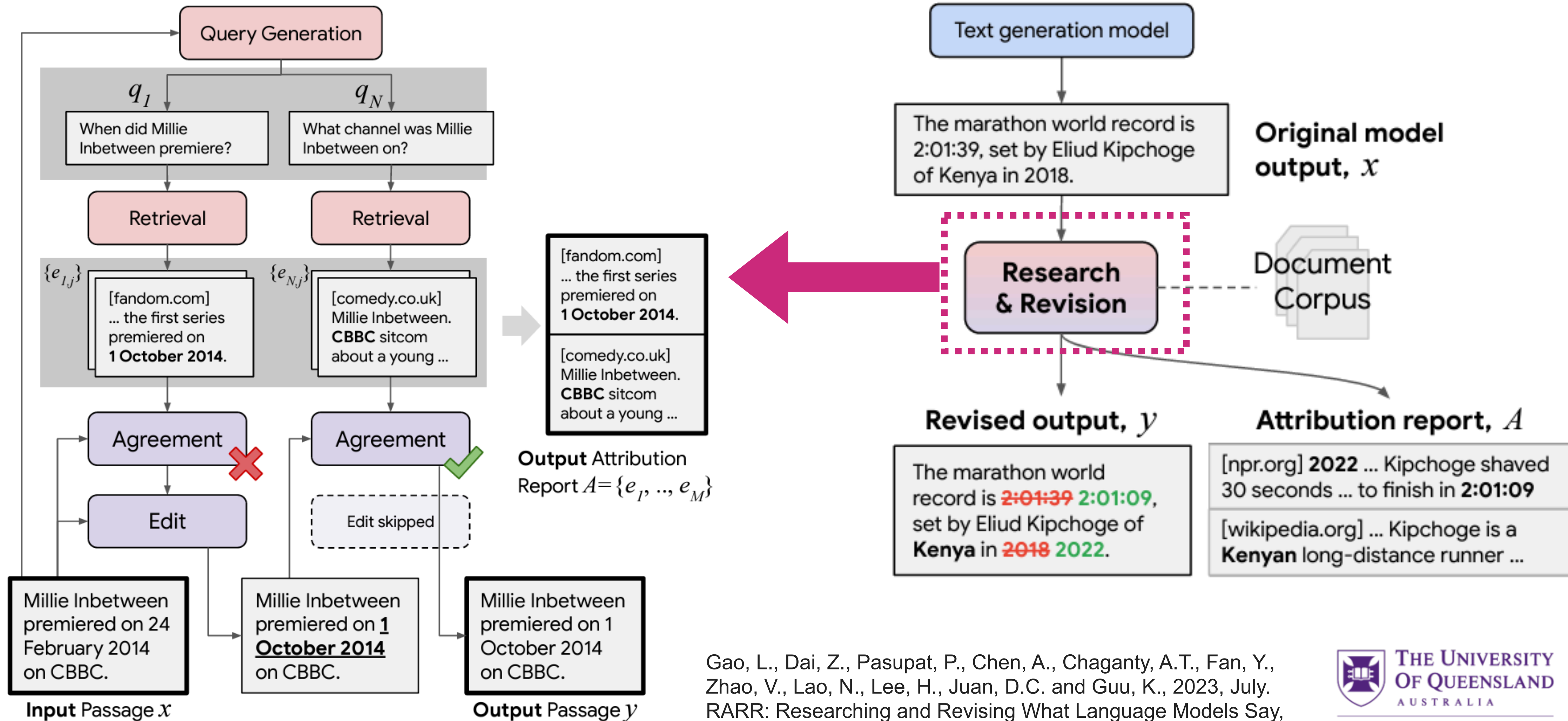
- automatically finds attribution for LLM output
- post-edits the output to fix unsupported content while preserving original output as much as possible.



# RARR: Retrofit Attribution using Research and Revision

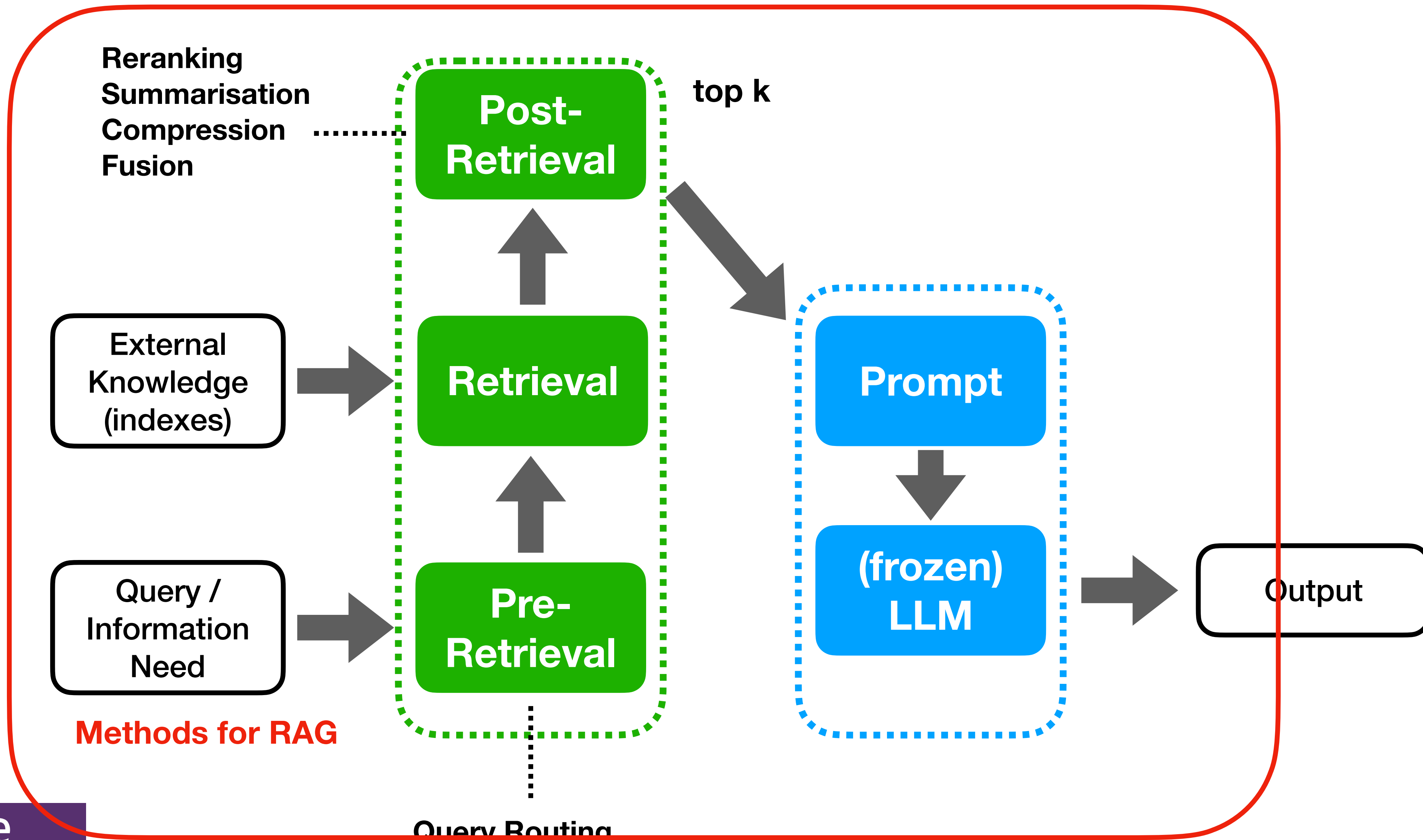


# RARR: Retrofit Attribution using Research and Revision



Gao, L., Dai, Z., Pasupat, P., Chen, A., Chaganty, A.T., Fan, Y., Zhao, V., Lao, N., Lee, H., Juan, D.C. and Guu, K., 2023, July. RARR: Researching and Revising What Language Models Say, Using Language Models. ACL 2023

# In the remaining of this part...



Methods for RAG

Query Routing  
Query Rewriting  
Query Expansion

Resources & Platforms for RAG

# Resources & Platforms for RAG

- Many **platforms/frameworks** for RAG
  - Industrial-oriented, e.g. Llama-Index, LangChain
    - Some provided only/also as a service, e.g. Nuclia AI
  - Research-oriented: **Ragnarök**, FlashRAG, **BERGEN**
- **Datasets**
  - Many datasets, not necessarily focusing on RAG
  - Few recent datasets specifically designed for RAG, e.g.
    - **TREC RAG 2024**: currently on-going [Topics released: August 4th; Submission deadline: August 11th]; Feb4Rag

# BERGEN

- end-to-end library for reproducible research standardizing RAG experiments.
- focus on QA
- Implements different state-of-the-art retrievers, rerankers, and LLMs.

```
python bergen.py dataset=kilt_nq generator=gemma-7b  
retriever=splade-v3 reranker=miniLM6
```

Natural Questions, Trivia QA, HotpotQA, WoW, ELI5,  
WikiQA, TruthfulQA, PopQA, ASQA, SCIQ

BM25, RetroMAE v2,  
RepLlama, SPLADE,  
DeBERTa-v3, MiniLM, ...

Llama-2-7B, 13B, 70B  
SOLAR-10B,  
Mixtral-8x7B,  
Gemma-2B, 7B, ...

F1, (Exact) Match, Rouge-\*, BEM, LLMEval, ...

**NAVER LABS** Europe

July, 2024

## BERGEN: A Benchmarking Library for Retrieval-Augmented Generation

David Rau\* Hervé Déjean Nadezhda Chirkova Thibault Formal Shuai Wang\* Vassilina Nikoulina  
Stéphane Clinchant

NAVER LABS Europe

\* Work performed while at Naver Labs Europe.

<https://github.com/naver/bergen>

Rau D., Déjean H., Chirkova N., Formal T., Wang S.,  
Nikoulina V., Clinchant, S.. 2024. BERGEN: A  
Benchmarking Library for Retrieval-Augmented  
Generation, arXiv:2407.01102v1



# BERGEN

Also large scale analysis of RAG components (state-of-the-art

- end retrievers, rerankers, LLMs -> 500+ experiments).

rese  
Key findings:

- foc

1. need to go beyond commonly used surface-matching metrics (e.g. exact match, F1, Rouge-L, etc.)

- Imp retr

2. retrieval quality matters for RAG response generation

**NAVER**

3. need to improve current knowledge-intensive benchmarks to use them in RAG:

**BERGEN**

**Genera**

- datasets evaluating general knowledge might not be suitable for RAG, as LLMs have acquired most such knowledge from Web/Wikipedia

David Rau\* I  
Stéphane Clin

NAVER LABS Eur

\* Work performed w

<https://github.com/naver/bergen>

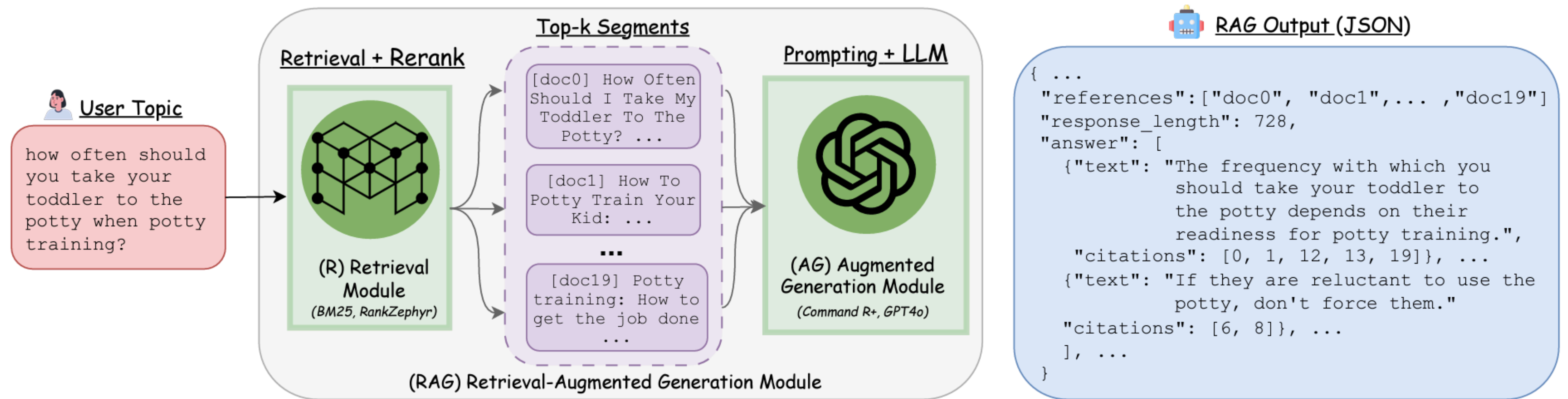
emma-7b

ELI5,

,70B

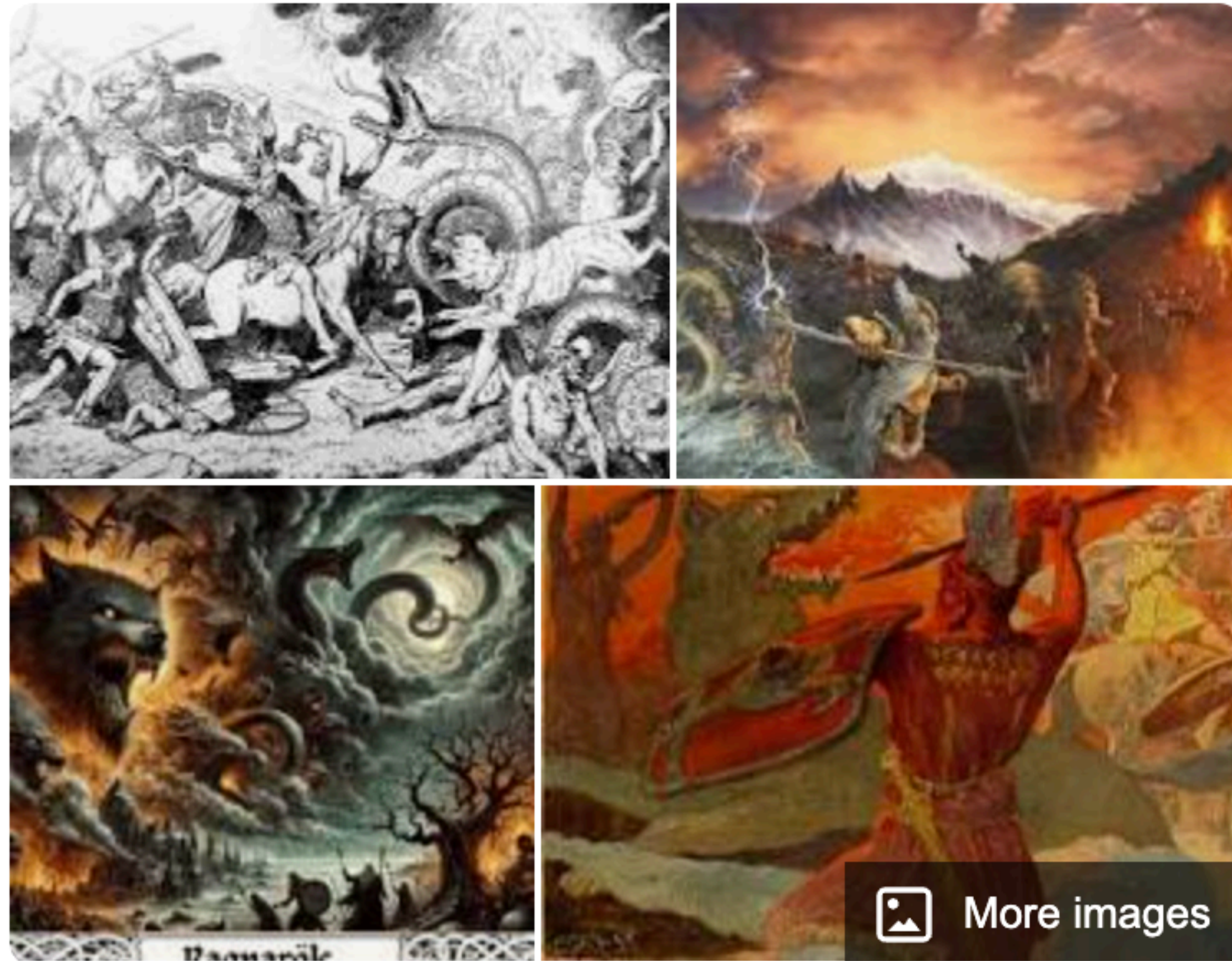
...

# TREC RAG 2024 & Ragnarök



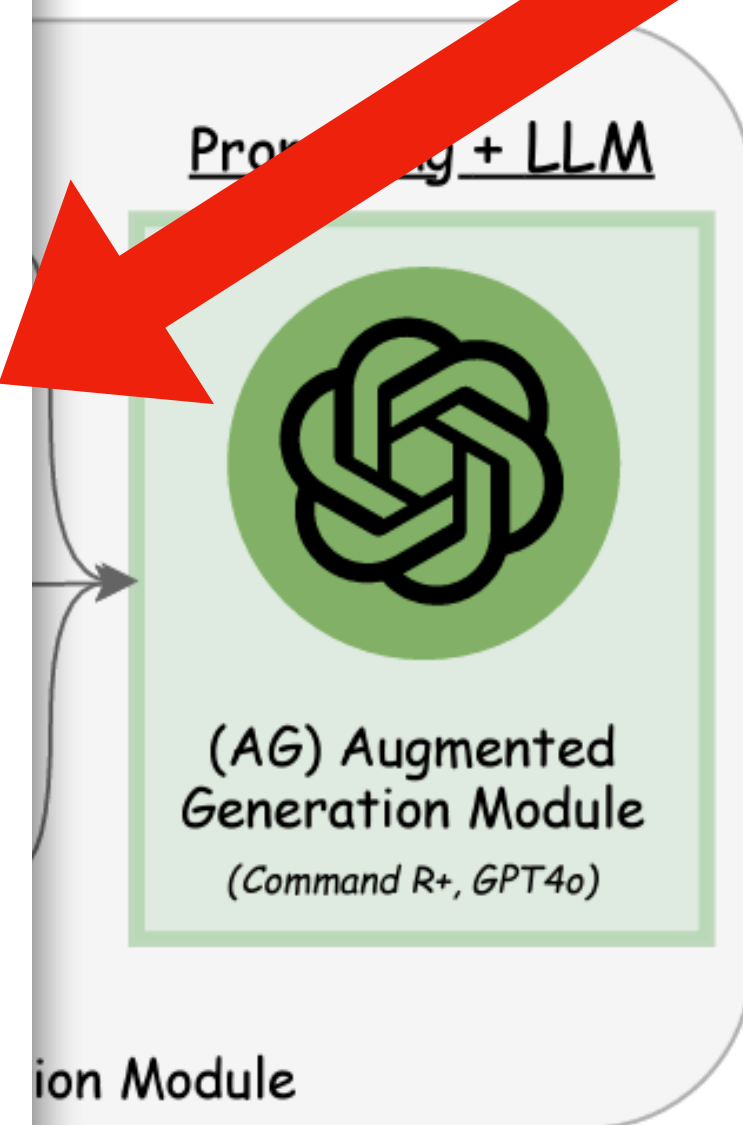
- TREC RAG 2024: dataset for RAG evaluation. Three tasks: Retrieve (R), Augmented Generation (AG), Retrieval Augmented Generation (RAG)
- Ragnarök: open-source, reproducible, reusable framework implementing RAG pipeline, with 2 sequential modules: (1) R, (2) AG

# Ragnarök :



In Norse mythology, Ragnarök is a foretold series of impending events, including a great battle in which numerous great Norse mythological figures will perish; it will entail a catastrophic series of natural disasters, including the burning of the world, and culminate in the submersion of the world underwater. [Wikipedia](#)

# 2024 & **Ragnarök**



RAG Output (JSON)

```
{ ...
  "references":["doc0", "doc1",... ,"doc19"]
  "response_length": 728,
  "answer": [
    {"text": "The frequency with which you
      should take your toddler to
      the potty depends on their
      readiness for potty training.",
      "citations": [0, 1, 12, 13, 19]}, ...
    {"text": "If they are reluctant to use the
      potty, don't force them."
      "citations": [6, 8]}, ...
  ], ...
}
```

evaluation. Three tasks: Retrieve (R), Augmented Generation (RAG)

**Really?**

framework implementing RAG

# Many More Datasets

Task	Sub Task	Dataset	Task	Sub Task	Dataset
QA	Single-hop	Natural Question(NQ) [111]	Dialog	Dialog Generation	Wizard of Wikipedia (WoW) [133]
		TrivialQA(TQA) [113]		Personal Dialog	KBP [134]
		SQuAD [114]		Task-oriented Dialog	DuleMon [136]
		Web Questions(WebQ) [115]	Recommendation	CamRest [137]	
	PopQA [116]			Amazon(Toys,Sport,Beauty) [138]	
	MS MARCO [117]				
	Multi-hop	HotpotQA [118]	IE	Event Argument Extraction	WikiEvent [139]
		2WikiMultiHopQA [119]		Relation Extraction	RAMS [140]
		MuSiQue [120]			T-REx [141],ZsRE [142]
	Long-form QA	ELI5 [121]	Reasoning	Commonsense Reasoning	HellaSwag [143]
NarrativeQA(NQA) [122]		CoT Reasoning		CoT Reasoning [144]	
ASQA [124]		Complex Reasoning		CSQA [145]	
QMSum(QM) [125]					
Domain QA	Qasper [126]	Others	Language Understanding	MMLU [146]	
	COVID-QA [127]		Language Modeling	WikiText-103 [147]	
	CMB [128],MMCU_Medical [129]		Fact Checking/Verification	StrategyQA [148]	
Multi-Choice QA	QuALITY [130]		Text Generation	FEVER [149]	
	ARC [131]		Text Summarization	PubHealth [150]	
	CommonsenseQA [132]		Text Classification	Biography [151]	
Graph QA	GraphQA [84]		Sentiment	WikiASP [152]	
			Code Search	XSum [153]	
			Robustness Evaluation	VioLens [154]	
			Math	TREC [155]	
		Machine Translation	SST-2 [156]		
				CodeSearchNet [157]	
				NoMIRACL [56]	
				GSM8K [158]	
				JRC-Acquis [159]	

From: Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J. and Wang, H., 2023. Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997.

# Many Open Directions

RAG vs. Long Context

RAG Evaluation

Scaling Laws in RAG

RAG Robustness

End-to-End Optimisation

Hybrid Approaches

Production-Ready RAGS

And many more: temporality,  
trade-offs



THE UNIVERSITY  
OF QUEENSLAND  
AUSTRALIA

CREATE CHANGE

# Part 3: Model-based IR

# Model-based IR

## Rethinking Search: Making Domain Experts out of Dilettantes\*

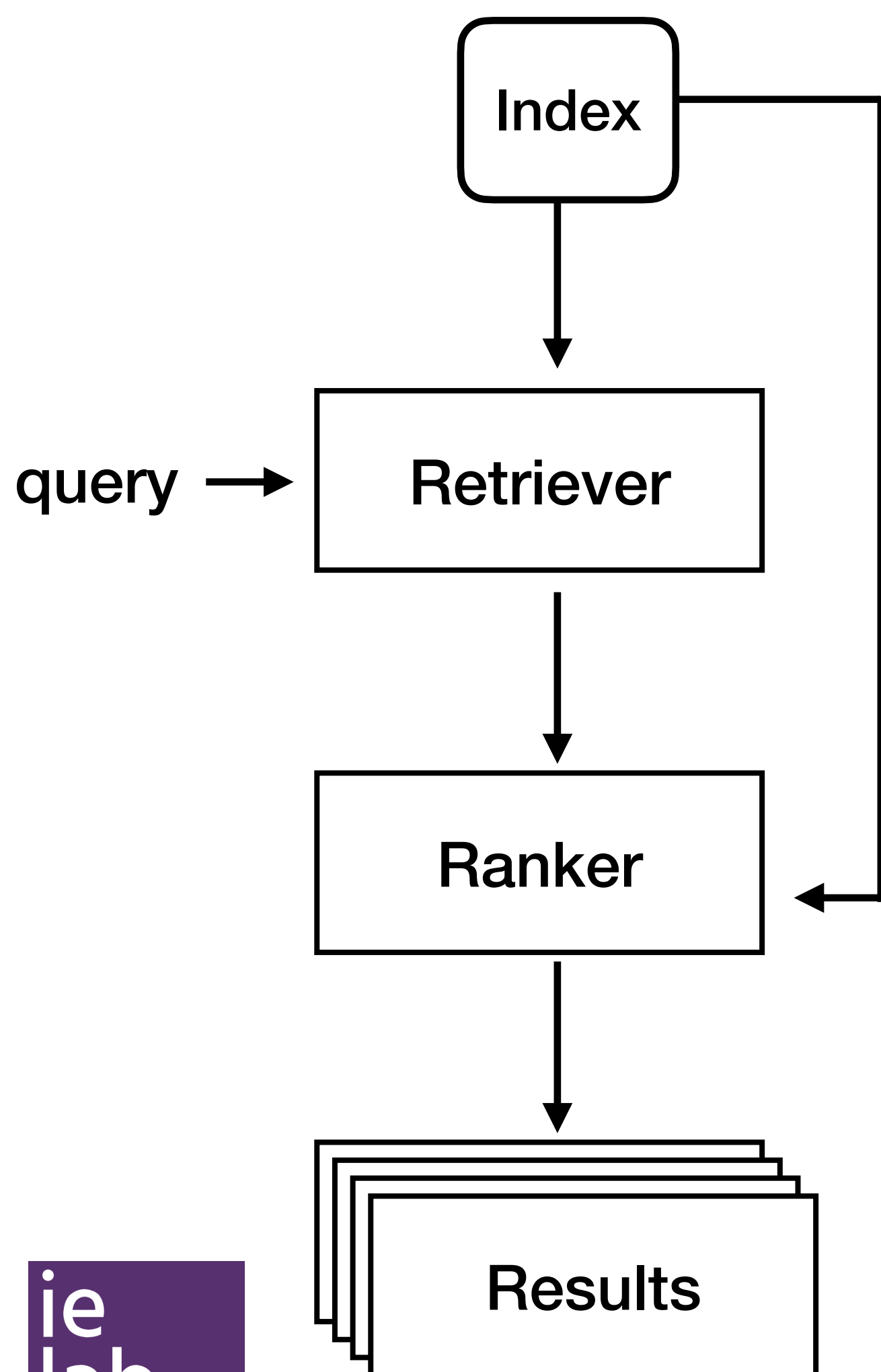
Donald Metzler  
Google Research  
*metzler@google.com*

Yi Tay  
Google Research  
*yitay@google.com*

Dara Bahri  
Google Research  
*dbahri@google.com*

Marc Najork  
Google Research  
*najork@google.com*

# Model-based IR



## Rethinking Search: Making Domain Experts out of Dilettantes\*

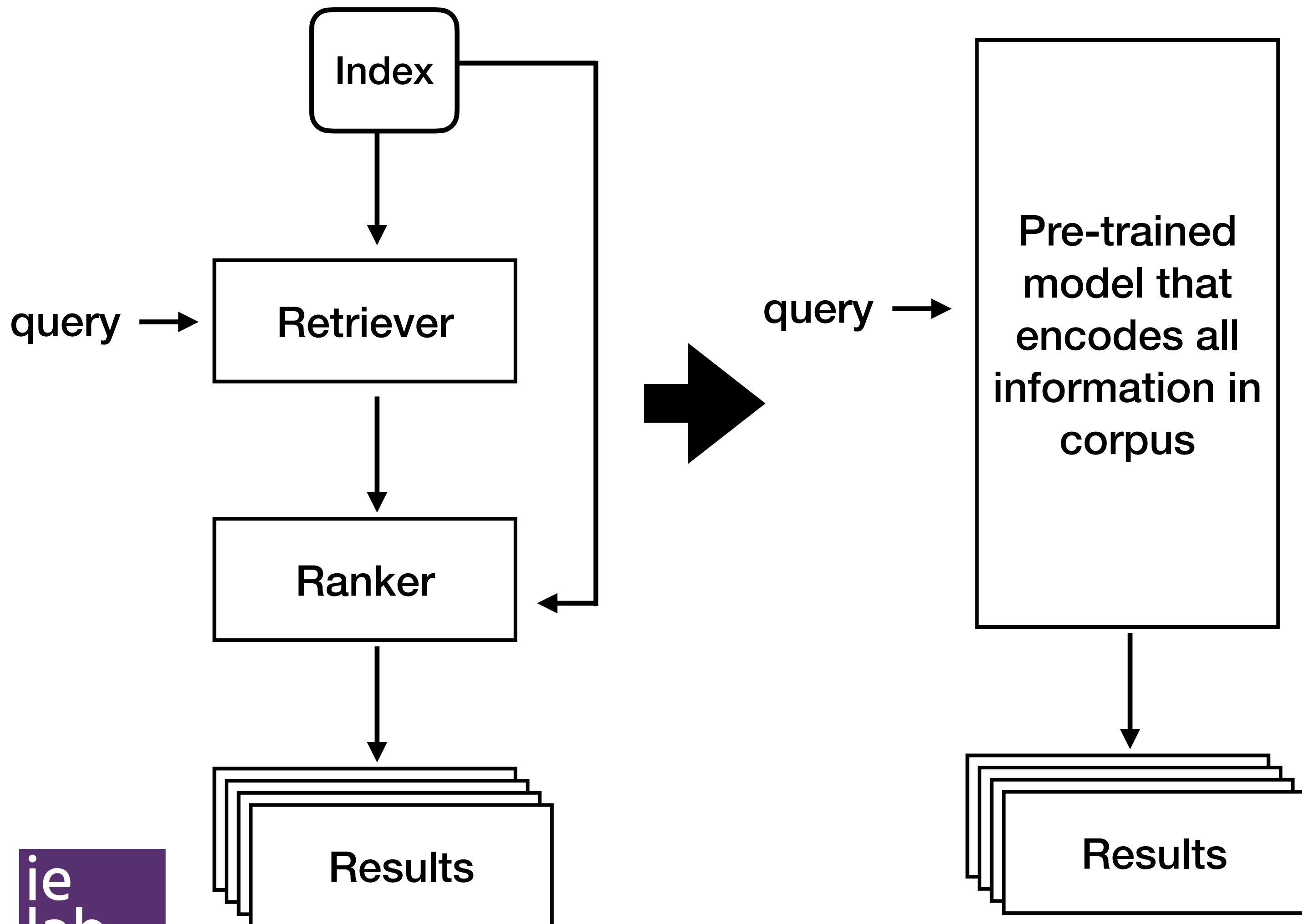
Donald Metzler  
Google Research  
*metzler@google.com*

Yi Tay  
Google Research  
*yitay@google.com*

Dara Bahri  
Google Research  
*dbahri@google.com*

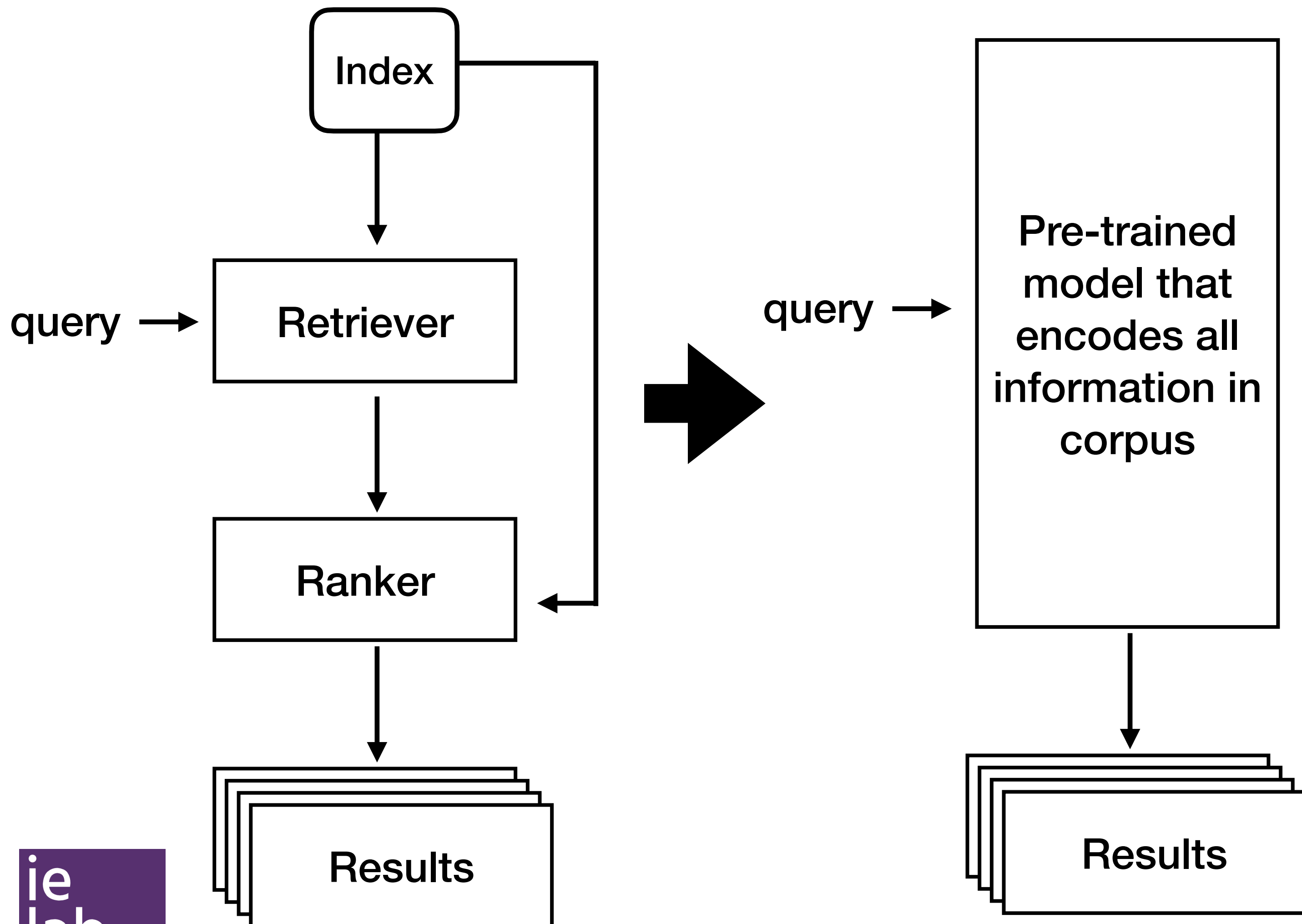
Marc Najork  
Google Research  
*najork@google.com*

# Model-based IR

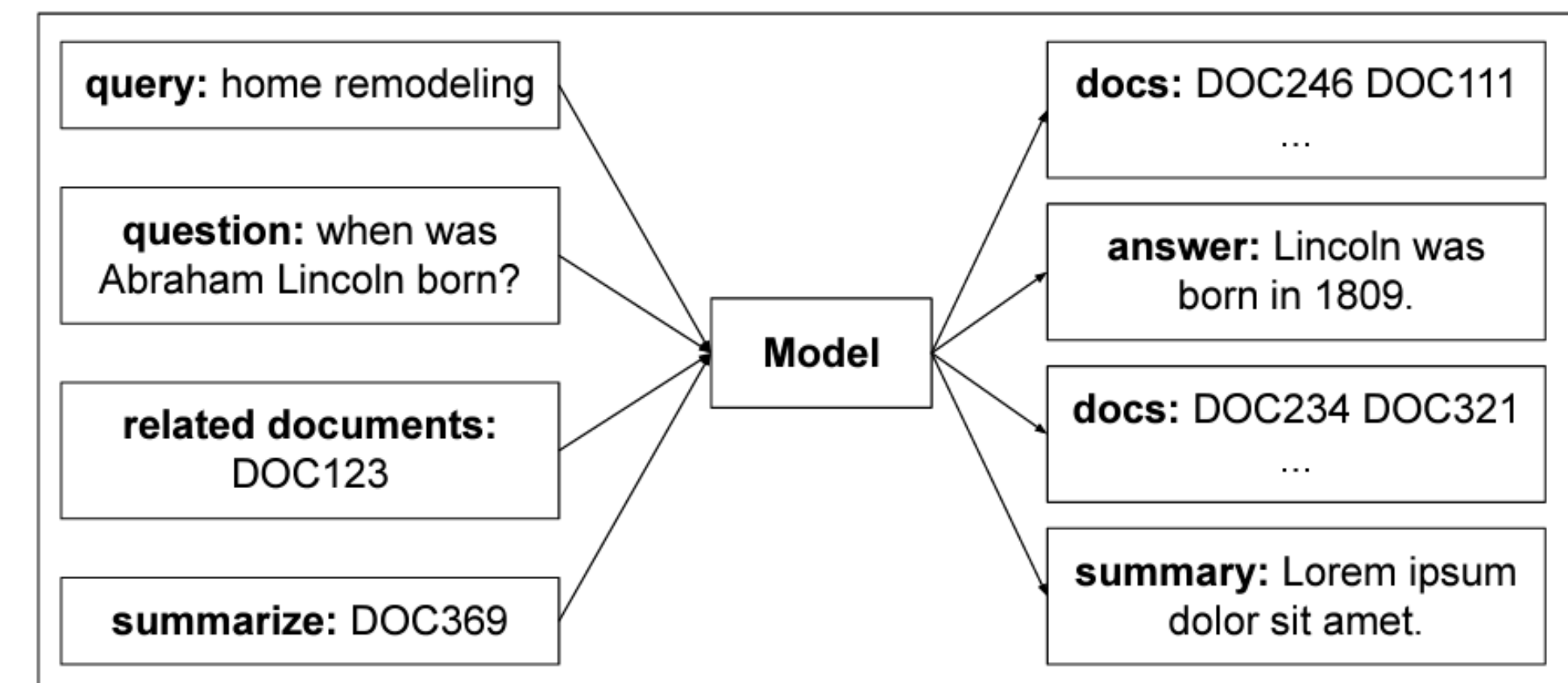


**Key Intuition:**  
Replace pipelined architecture with a single consolidated model (Unified Retrieve-and-Rank)

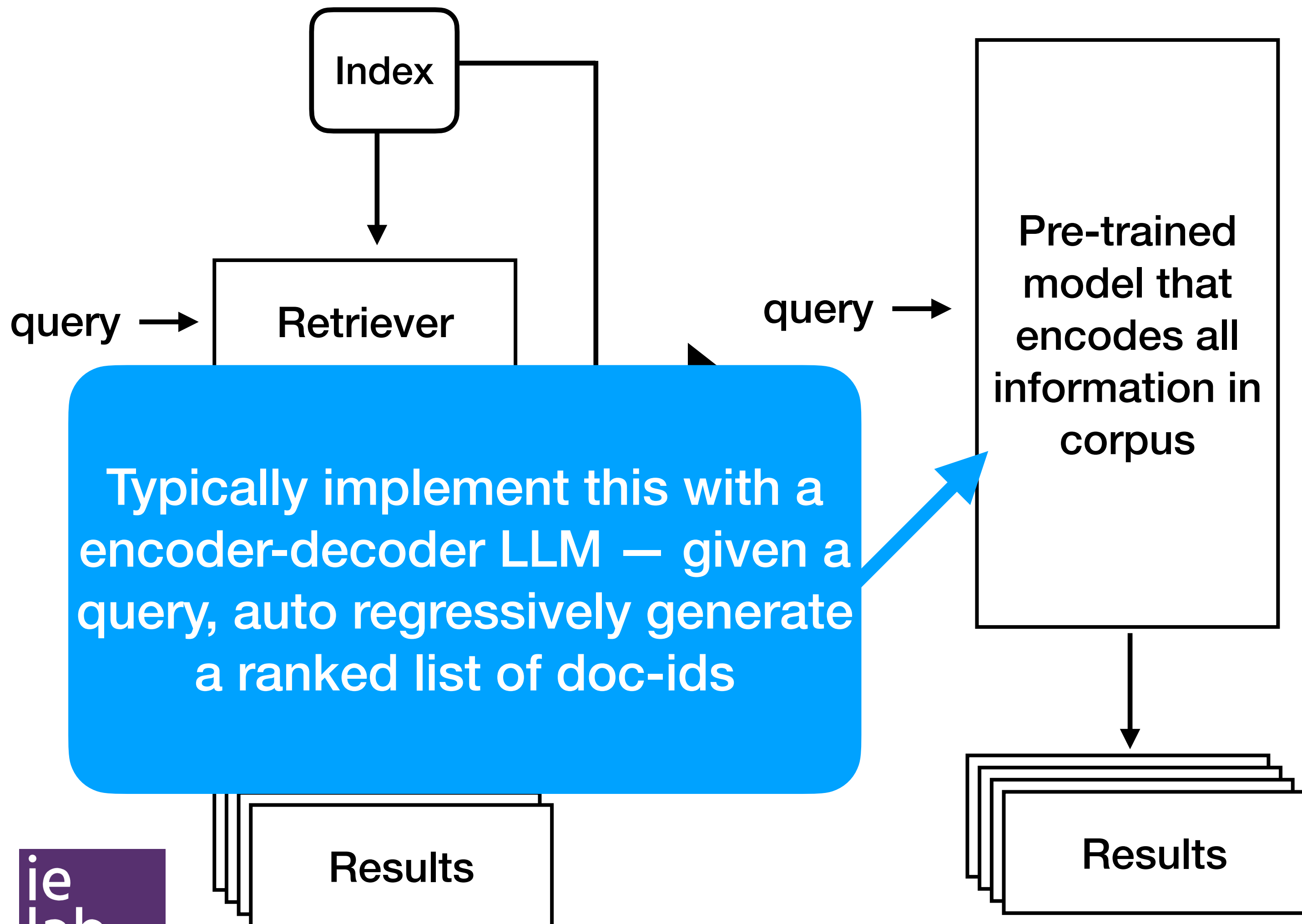
# Model-based IR



**Key Intuition:**  
Replace pipelined architecture with a single consolidated model (Unified Retrieve-and-Rank)

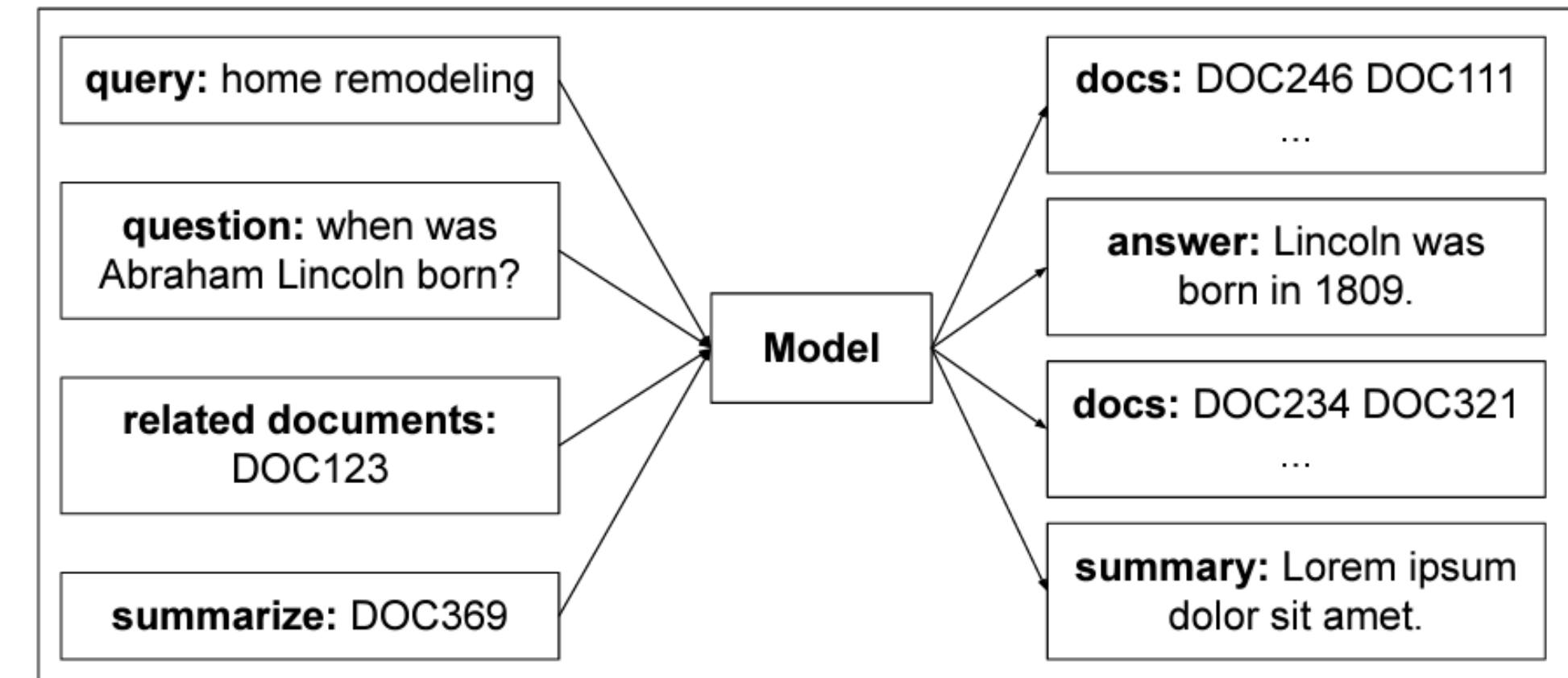


# Model-based IR



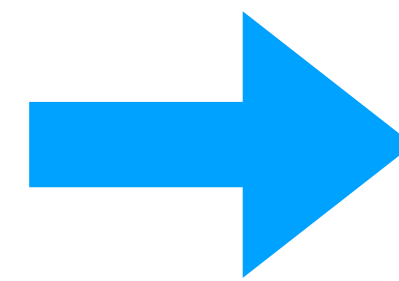
**Key Intuition:**

Replace pipelined architecture with a single consolidated model (Unified Retrieve-and-Rank)

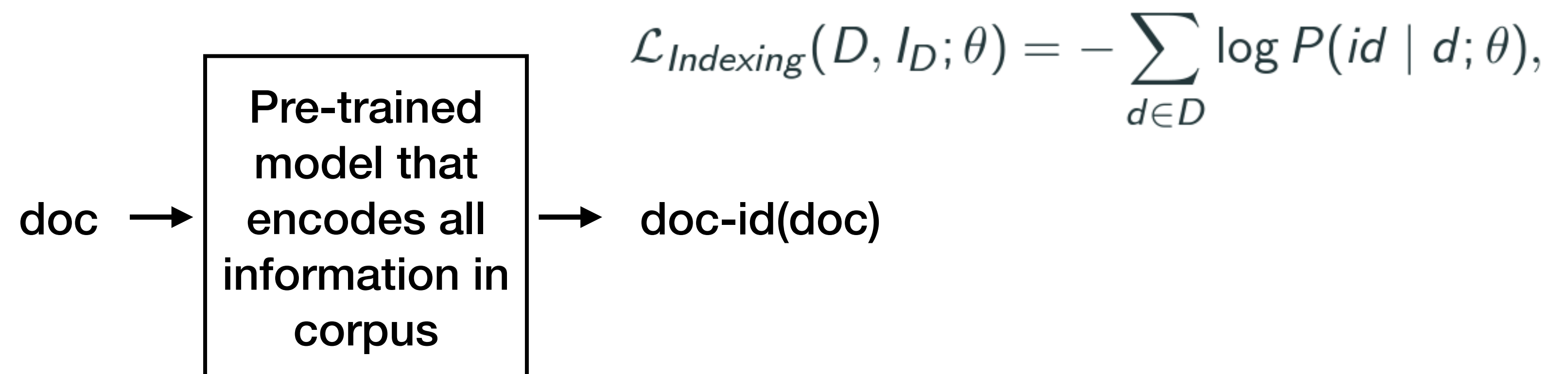


# “Indexing” in Distributed Search Index (DSI)

the model needs to memorise the content/information of each document

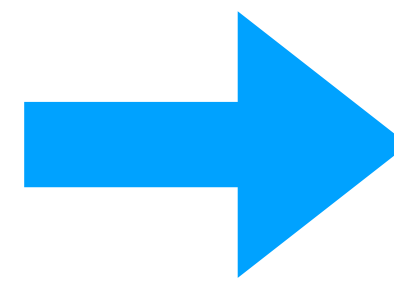


learn to associate the content of each document with its corresponding doc-id

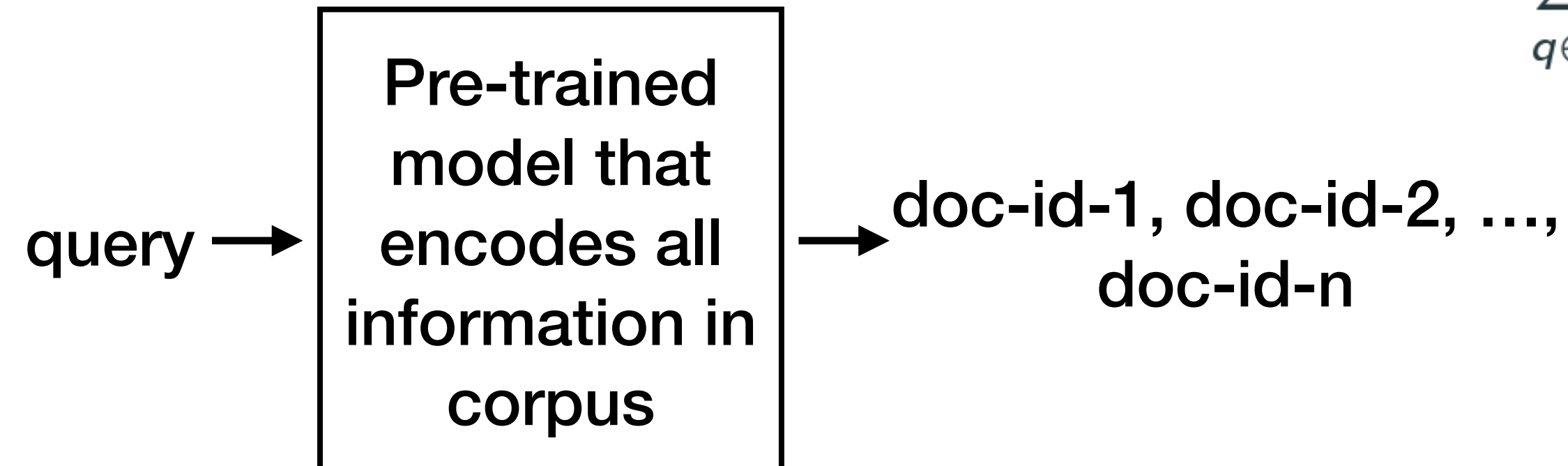


# “Retrieval” in Distributed Search Index (DSI)

the model needs to answer a query

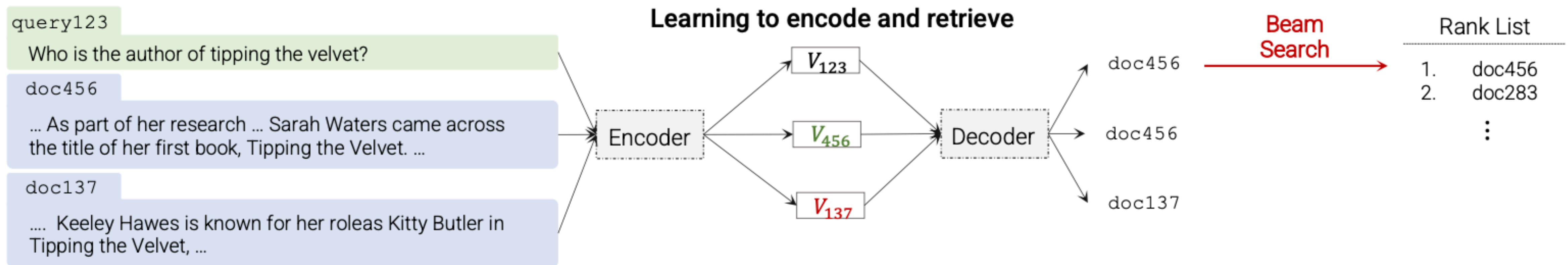


learn to generate a ranked list of candidate docids



$$\mathcal{L}_{\text{Retrieval}}(Q, I_Q; \theta) = - \sum_{q \in Q} \sum_{id^q \in I_Q} \log P(id^q | q; \theta),$$

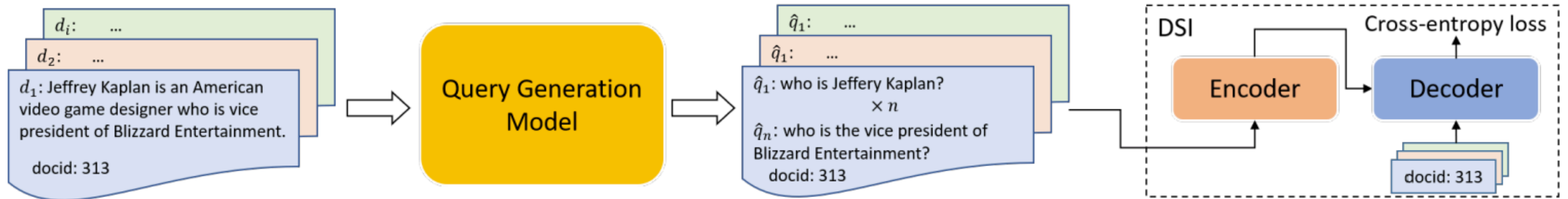
# Jointly Learn to Index and Retrieve



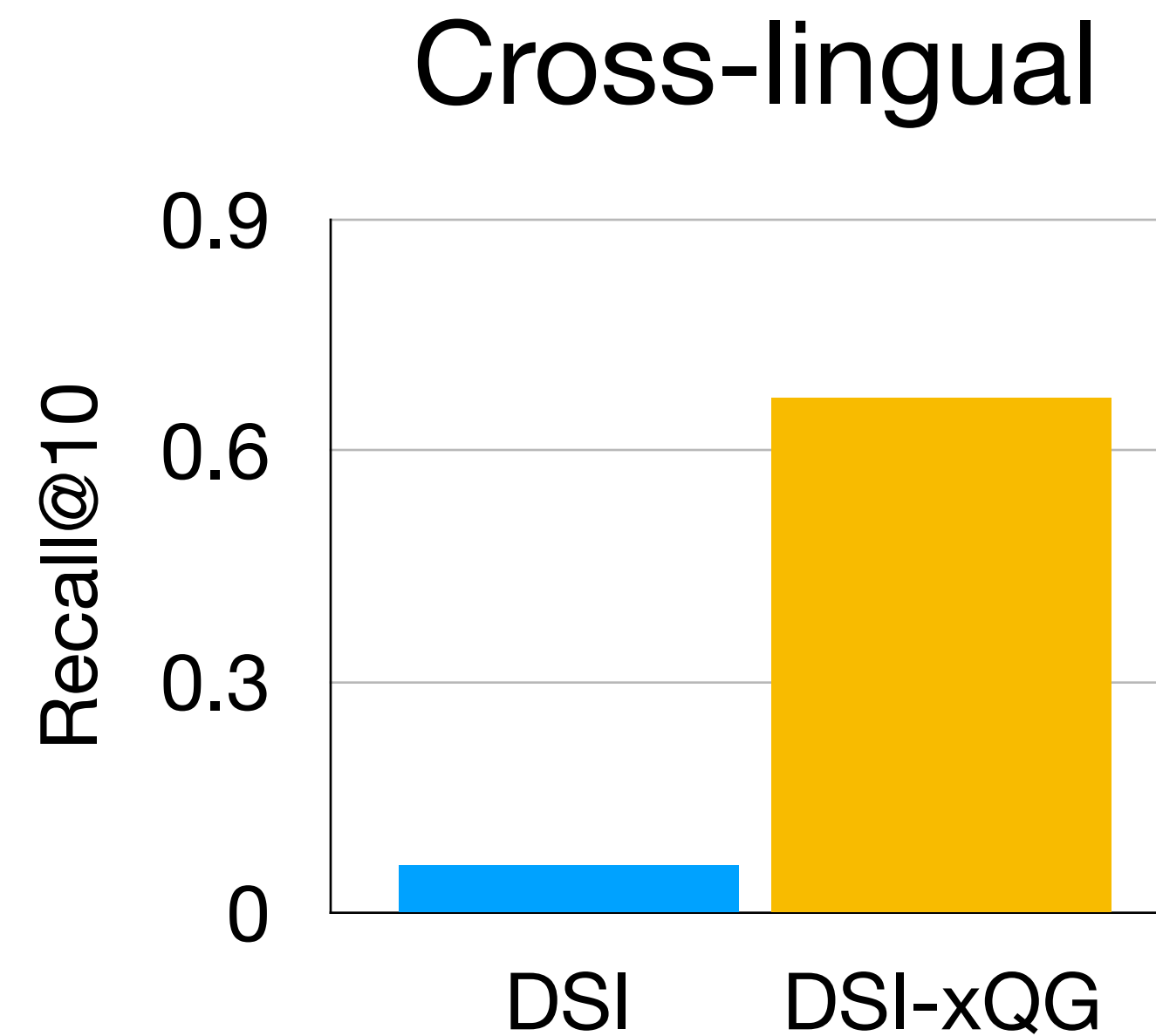
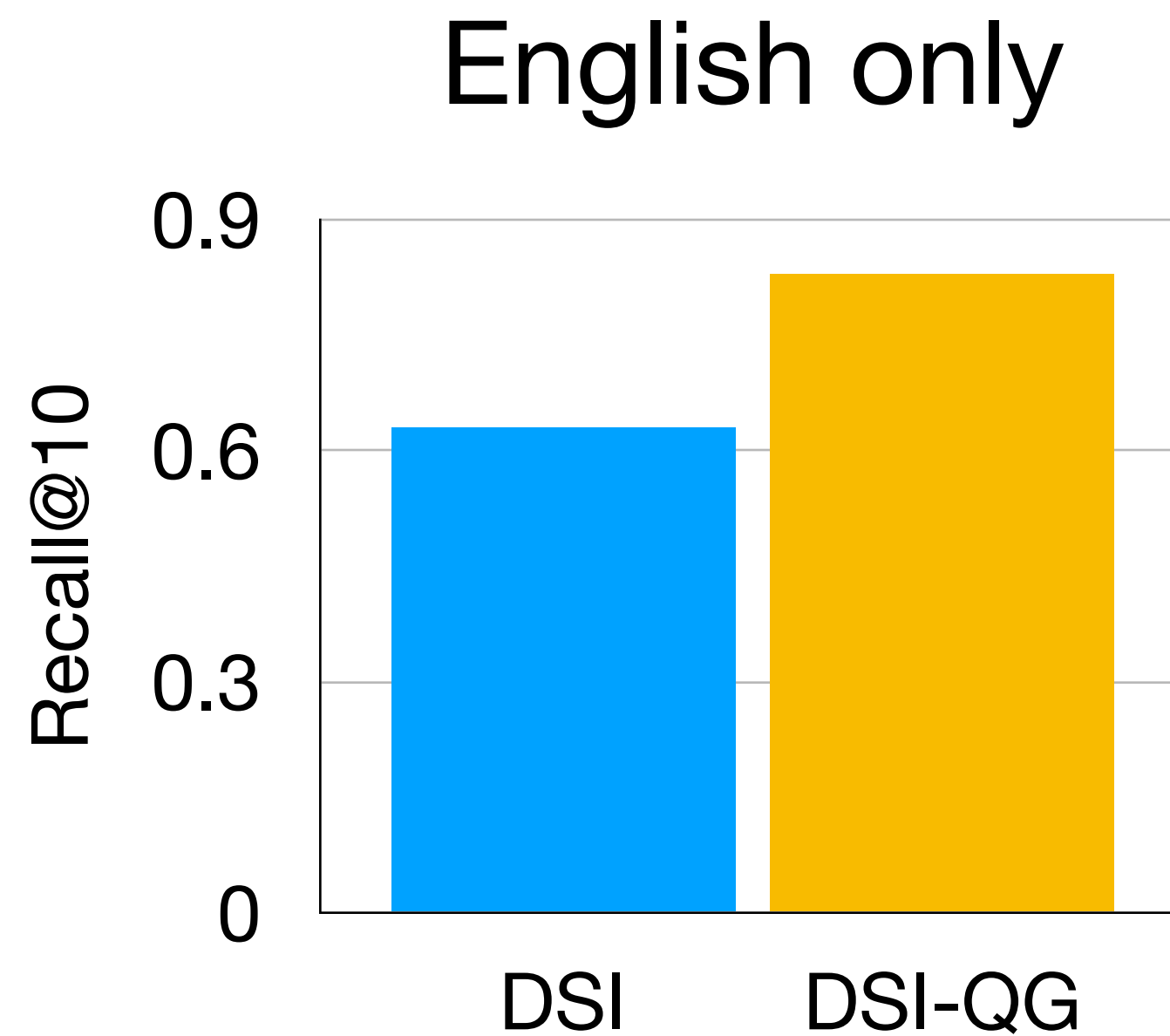
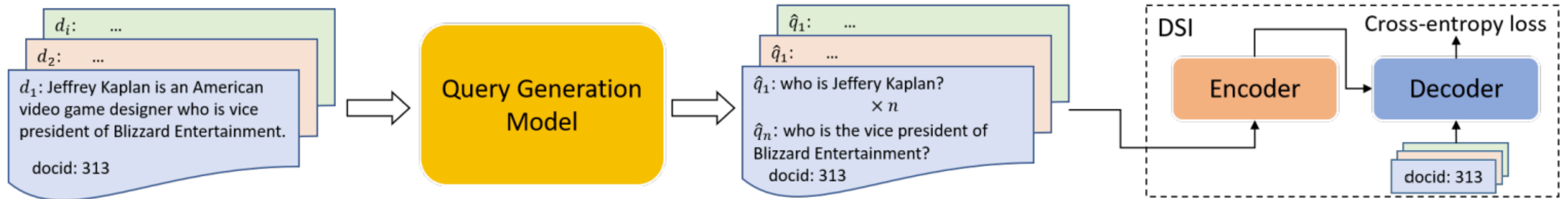
Optimise a single model directly in an end-to-end fashion towards global objective:

$$\mathcal{L}_{Global}(Q, D, I_D, I_Q; \theta) = \mathcal{L}_{Indexing}(D, I_D; \theta) + \mathcal{L}_{Retrieval}(Q, I_Q; \theta)$$

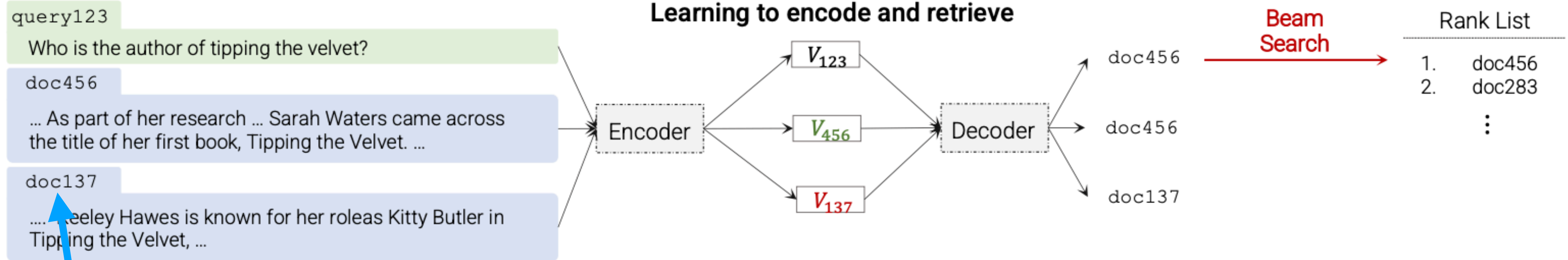
# Differentiable Search Index with Query Generation



# Differentiable Search Index with Query Generation



# Questions

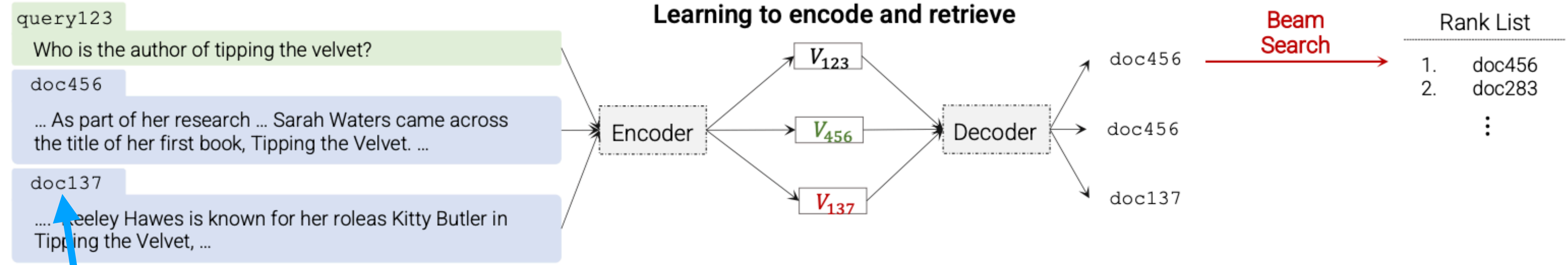


How should doc-ids look like?

# Questions

How to handle heterogeneous tasks (data distributions, objective functions)?

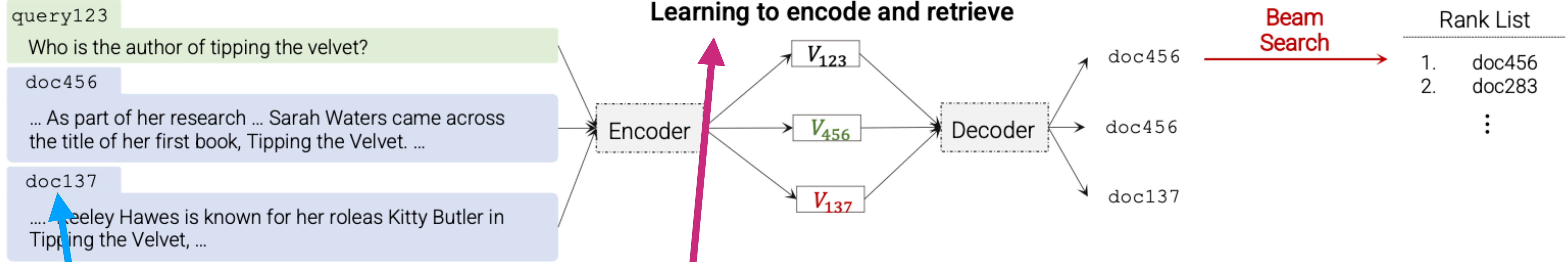
## Learning to encode and retrieve



How should doc-ids look like?

# Questions

How to handle heterogeneous tasks (data distributions, objective functions)?



How should doc-ids look like?

How to memorise

- lots of data
- new/updating data
- efficiently

# Questions

How to handle heterogeneous tasks (data distributions, objective functions)?

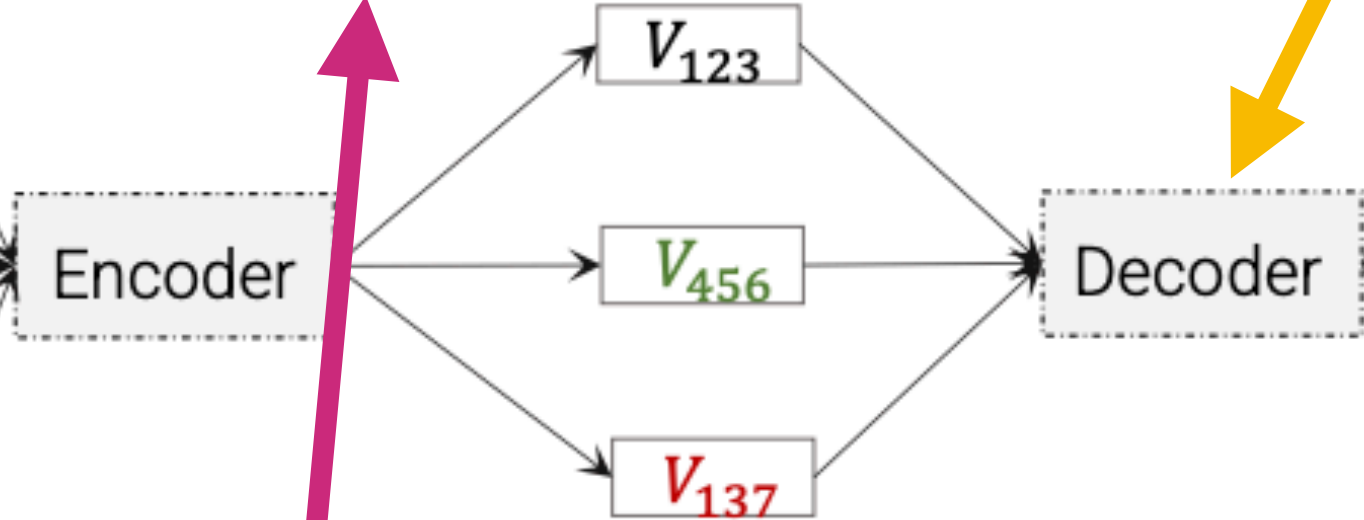
How to generate valid doc-ids?

query123  
Who is the author of tipping the velvet?

doc456  
... As part of her research ... Sarah Waters came across the title of her first book, Tipping the Velvet. ...

doc137  
... Keeley Hawes is known for her roles as Kitty Butler in Tipping the Velvet, ...

Learning to encode and retrieve



Beam Search

Rank List

1.	doc456
2.	doc283
	⋮

How should doc-ids look like?

How to memorise

- lots of data
- new/updating data
- efficiently

# Questions

How to handle heterogeneous tasks (data distributions, objective functions)?

How to generate valid doc-ids?

Learning to encode and retrieve

query123

Who is the author of tipping the velvet?

doc456

$V_{123}$

doc456

Beam Search

Rank List

- 1. doc456
- 2. doc283

Why doing all this?

... Kelly Hawes is known for her roles as Kitty Butler in Tipping the Velvet, ...

How should doc-ids look like?

- How to memorise
  - lots of data
  - new/updating data
  - efficiently





# Why doing DSI

## 1) Efficiency at inference

- **Memory:**

- Inverted index: corpus grows -> index grows
- DSI: corpus grows -> model size stays the same (+ small overhead for doc-id storage)

- **Latency:** light generative process over the vocabulary of identifiers

	Dense retrieval	Generative retrieval
Memory size (MS MARCO 300K)	 GTR 1430MB	 GenRet 860MB
Online latency	 GTR 1.97s	 GenRet 0.16s

## 2) Dramatic **simplification** of whole indexing/retrieval/generation process

- Less engineering/maintenance
- Imaging of instead of just generating doc-ids, you could generate answers, with doc-ids (attribution) ~> RAG



THE UNIVERSITY  
OF QUEENSLAND  
AUSTRALIA

CREATE CHANGE

**Bonus!**

**Extreme Model-based IR**



# Reflections on: “Can Long-Context Language Models Subsume Retrieval, RAG, SQL, and More?”

## Can Long-Context Language Models Subsume Retrieval, RAG, SQL, and More?

---

Jinhyuk Lee\* Anthony Chen\* Zhuyun Dai\*

Dheeru Dua Devendra Singh Sachan Michael Boratko Yi Luan

Sébastien M. R. Arnold Vincent Perot Siddharth Dalmia Hexiang Hu

Xudong Lin Panupong Pasupat Aida Amini Jeremy R. Cole

Sebastian Riedel Iftekhar Naim Ming-Wei Chang Kelvin Guu

Google DeepMind

### Abstract

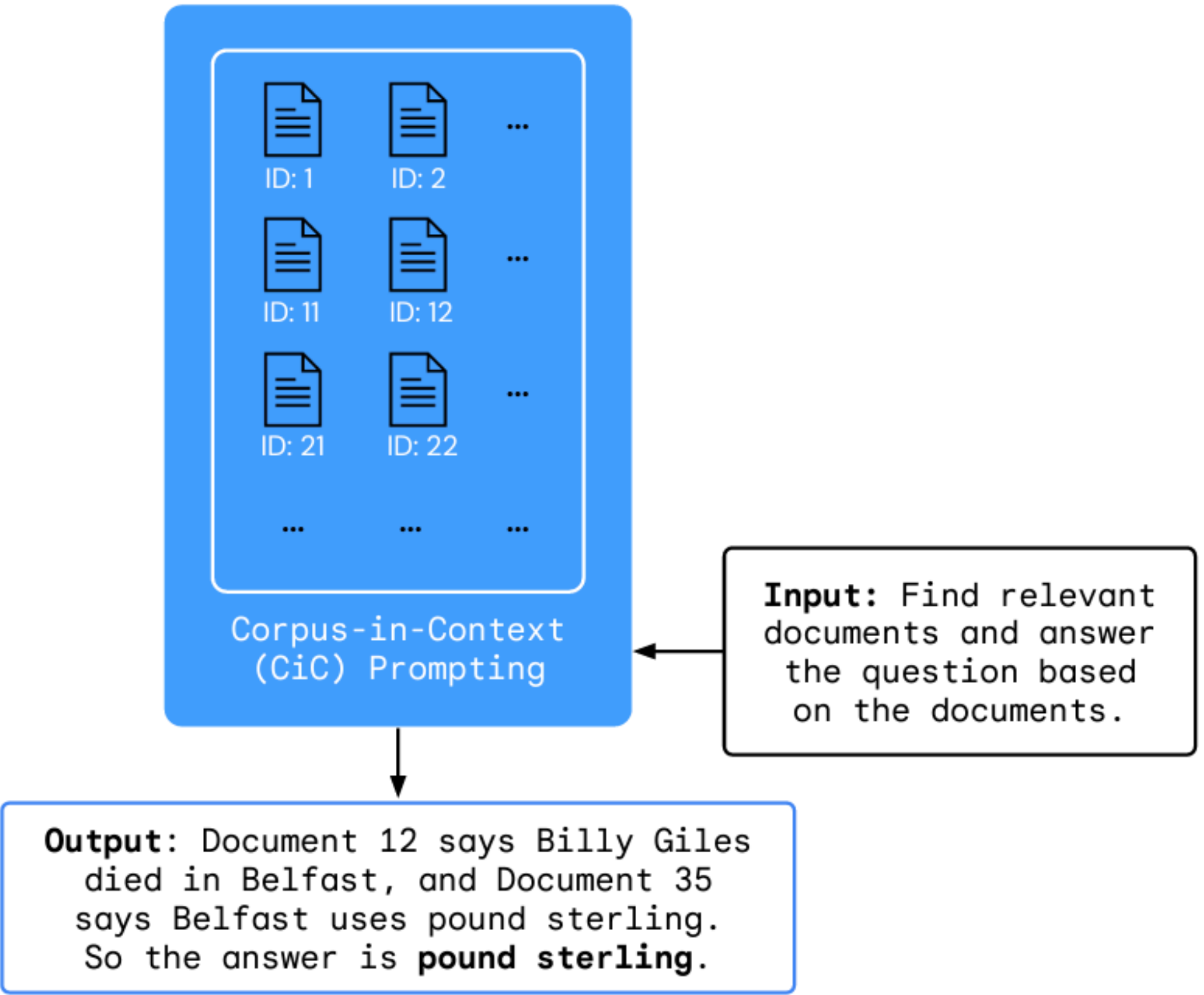
Long-context language models (LCLMs) have the potential to revolutionize our approach to tasks traditionally reliant on external tools like retrieval systems or databases. Leveraging LCLMs’ ability to natively ingest and process entire corpora of information offers numerous advantages. It enhances user-friendliness by eliminating the need for specialized knowledge of tools, provides robust end-to-end modeling that minimizes cascading errors in complex pipelines, and allows for the application of sophisticated prompting techniques across the entire system. To assess this paradigm shift, we introduce LOFT, a benchmark of real-world tasks requiring context up to millions of tokens designed to evaluate LCLMs’ performance on in-context retrieval and reasoning. Our findings reveal LCLMs’ surprising ability to rival state-of-the-art retrieval and RAG systems, despite never having been explicitly trained for these tasks. However, LCLMs still face challenges in areas like compositional reasoning that are required in SQL-like tasks. Notably, prompting strategies significantly influence performance, emphasizing the need for continued research as context lengths grow. Overall, LOFT provides a rigorous testing ground for LCLMs, showcasing their potential to supplant existing paradigms and tackle novel tasks as model capabilities scale.<sup>1</sup>

406.13121v1 [cs.CL] 19 Jun 2024

# Reflections on: “Can Long-Context Language Models Subsume Retrieval, RAG, SQL, and More?”

## Long-Context Language Models

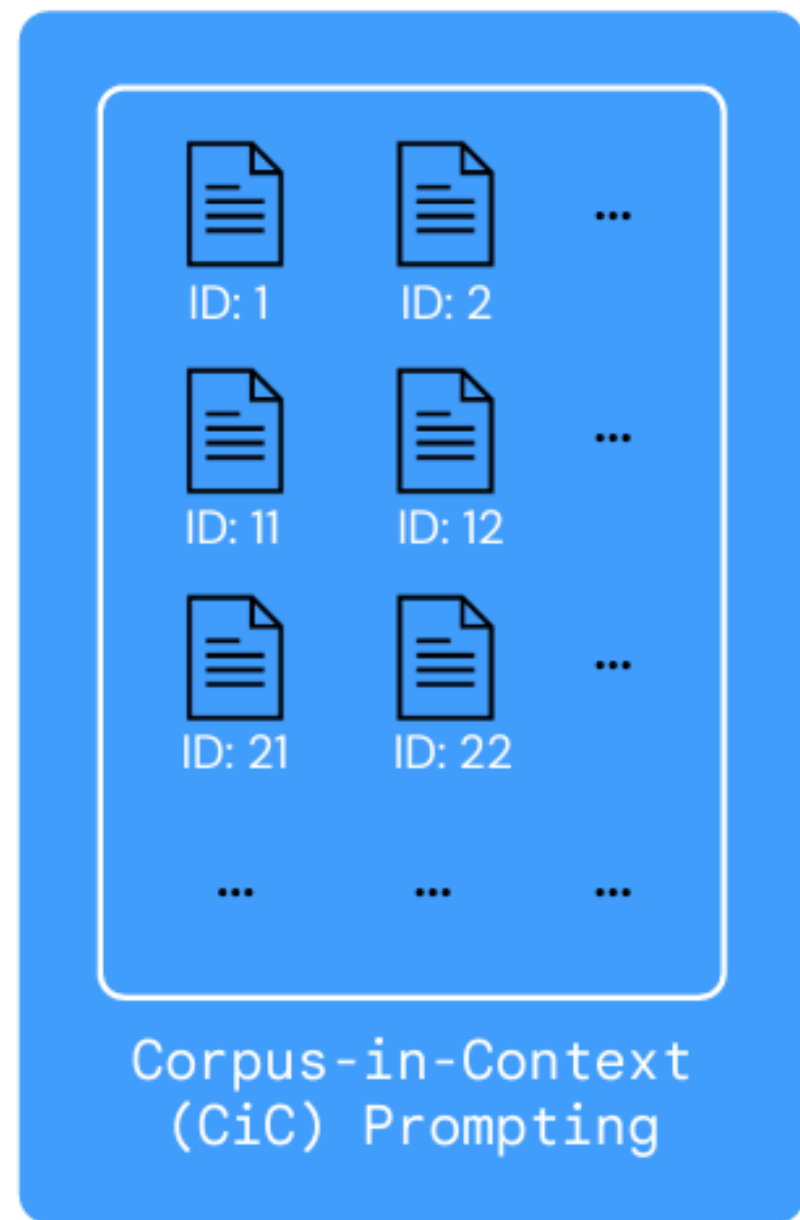
(e.g. Gemini 1.5 Pro, GPT-4o)



# Reflections on: “Can Long-Context Language Models Subsume Retrieval, RAG, SQL, and More?”

## Long-Context Language Models

(e.g. Gemini 1.5 Pro, GPT-4o)



**Input:** Find relevant documents and answer the question based on the documents.

**Output:** Document 12 says Billy Giles died in Belfast, and Document 35 says Belfast uses pound sterling. So the answer is **pound sterling**.

You will be given a list of documents. You need to read carefully and understand all of them. Then you will be given a query that may require you to use 1 or more documents to find the answer. Your goal is to find all documents from the list that can help answer the query.

Instruction

ID: 0 | TITLE: Shinji Okazaki | CONTENT: Shinji Okazaki is a Japanese ... | END ID: 0  
 ...  
 ID: 53 | TITLE: Ain't Thinkin' 'Bout You | CONTENT: "Ain't Thinkin' 'Bout You" is a song ... | END ID: 53  
 ID: 54 | TITLE: Best Footballer in Asia 2016 | CONTENT: ... was awarded to Shinji Okazaki ... | END ID: 54  
 ...

Corpus Formatting

=====  
 Example 1  
 Which documents are needed to answer the query? Print out the TITLE and ID of each document. Then format the IDs into a list.  
 query: What year was the recipient of the 2016 Best Footballer in Asia born?  
 The following documents are needed to answer the query:  
 TITLE: Best Footballer in Asia 2016 | ID: 54  
 TITLE: Shinji Okazaki | ID: 0  
**Final Answer: [54, 0]**  
 ...

Few-shot Exemples

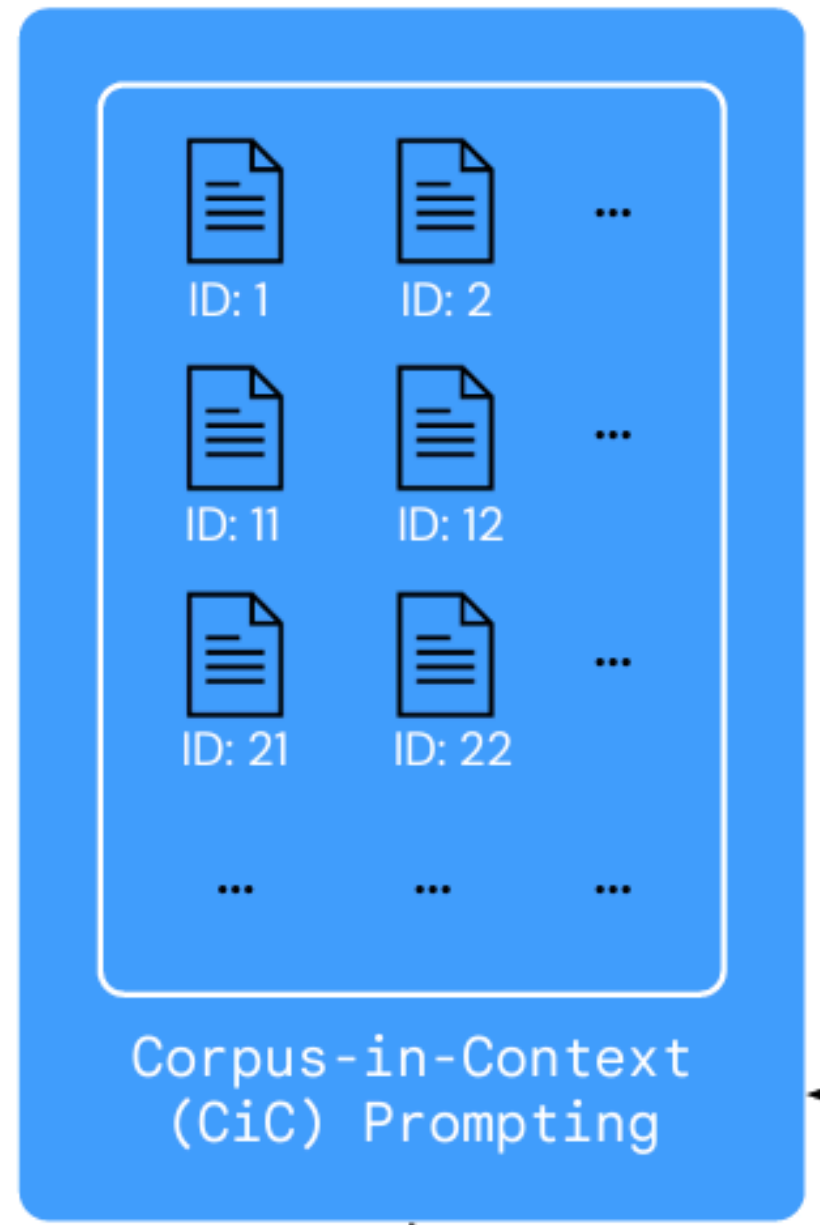
=====  
 Now let's start!  
 Which documents are needed to answer the query? Print out the TITLE and ID of each document. Then format the IDs into a list.  
 query: **How many records had the team sold before performing "aint thinkin bout you"?**  
 The following documents are needed to answer the query:

Query Formatting

# Reflections on: “Can Long-Context Language Models Subsume Retrieval, RAG, SQL, and More?”

Long-Context Language Models  
(e.g. Gemini 1.5 Pro, GPT-4o)

But this is going to be expensive!



**Input:** Find relevant documents and answer the question based on the documents.

**Output:** Document 12 says Billy Giles died in Belfast, and Document 35 says Belfast uses pound sterling. So the answer is **pound sterling**.

**Limitations** Our experiments were constrained by the computational resources and financial costs associated with utilizing LCLMs. The entire LOFT 128k test sets contain around 35 datasets × 100 prompts × 128k tokens = 448M input tokens, which cost \$1,568 for Gemini 1.5 Pro, \$2,240 for GPT-4o, and \$6,720 for Claude 3 Opus at the time of writing. To reduce costs, we also release dev sets, which are 10x smaller and can be evaluated with around \$200 using Gemini 1.5 Pro or GPT-4o. We also expect LLM API prices to decrease over time. Another limitation of this work is that we focused on evaluating the quality of LCLMs, and leave efficiency considerations for future work. We could not measure the efficiency improvements from prefix caching [20] due to API constraints at the time of writing. Without caching, the Gemini 1.5 Pro API has a median latency of roughly four seconds for 32k input tokens, twelve seconds for 128k input tokens, and 100 seconds for 1 million input tokens. This speed is likely slower than specialized retrievers or SQL databases; the promising quality results on LOFT encourage further investigation into optimizing LCLMs efficiency. Additionally, our retrieval and RAG tasks was limited to 1 million tokens, which still leaves a large gap from real-world applications that may involve several million or even billions of documents.

# Reflections on: “Can Long-Context Language Models Subsume Retrieval, RAG, SQL, and More?”

## Long-Context Language Models

(e.g. Gemini 1.5 Pro, GPT-4o)



But this is going to be expensive!

But...

You will be given a list of documents. You need to read carefully and understand all of them. Then you will be given a query that may require you to use 1 or more documents to find the answer. Your goal is to find all documents from the list that can help answer the query.

ID: 0 | TITLE: Shinji Okazaki | CONTENT: Shinji Okazaki is a Japanese ... | END ID: 0  
 ...  
 ID: 53 | TITLE: Ain't Thinkin' 'Bout You | CONTENT: "Ain't Thinkin' 'Bout You" is a song ... | END ID: 53  
 ID: 54 | TITLE: Best Footballer in Asia 2016 | CONTENT: ... was awarded to Shinji Okazaki ... | END ID: 54  
 ...

=====  
 Which documents are needed to answer the query? Print out the TITLE and ID of each document. Then format the IDs into a list.  
 query: What year was the recipient of the 2016 Best Footballer in Asia born?  
 The following documents are needed to answer the query:  
 TITLE: Best Footballer in Asia 2016 | ID: 54  
 TITLE: Shinji Okazaki | ID: 0  
**Final Answer: [54, 0]**  
 ...

=====  
 Which documents are needed to answer the query? Print out the TITLE and ID of each document. Then format the IDs into a list.  
 query: **How many records had the team sold before performing "aint thinkin bout you"?**  
 The following documents are needed to answer the query:

**Limitations** Our experiments were constrained by the computational resources and financial costs associated with utilizing LCLMs. The entire LOFT 128k test sets contain around 35 datasets × 100,000 documents = 448M input tokens, which cost \$1,568 for Gemini 1.5 Pro, \$2,240 for Claude 3 Opus at the time of writing. To reduce costs, we also release developer and can be evaluated with around \$200 using Gemini 1.5 Pro or GPT-4o. API prices to decrease over time. Another limitation of this work is that we use the quality of LCLMs, and leave efficiency considerations for future work. We explore efficiency improvements from prefix caching [20] due to API constraints. Without caching, the Gemini 1.5 Pro API has a median latency of roughly 100 seconds for 128k input tokens, and 100 seconds for 1M input tokens. This speed is likely slower than specialized retrievers or SQL databases; the LOFT encourage further investigation into optimizing LCLMs efficiency. Our evaluation and RAG tasks was limited to 1 million tokens, which still leaves a large number of applications that may involve several million or even billions of documents.

Instruction

Corpus Formatting

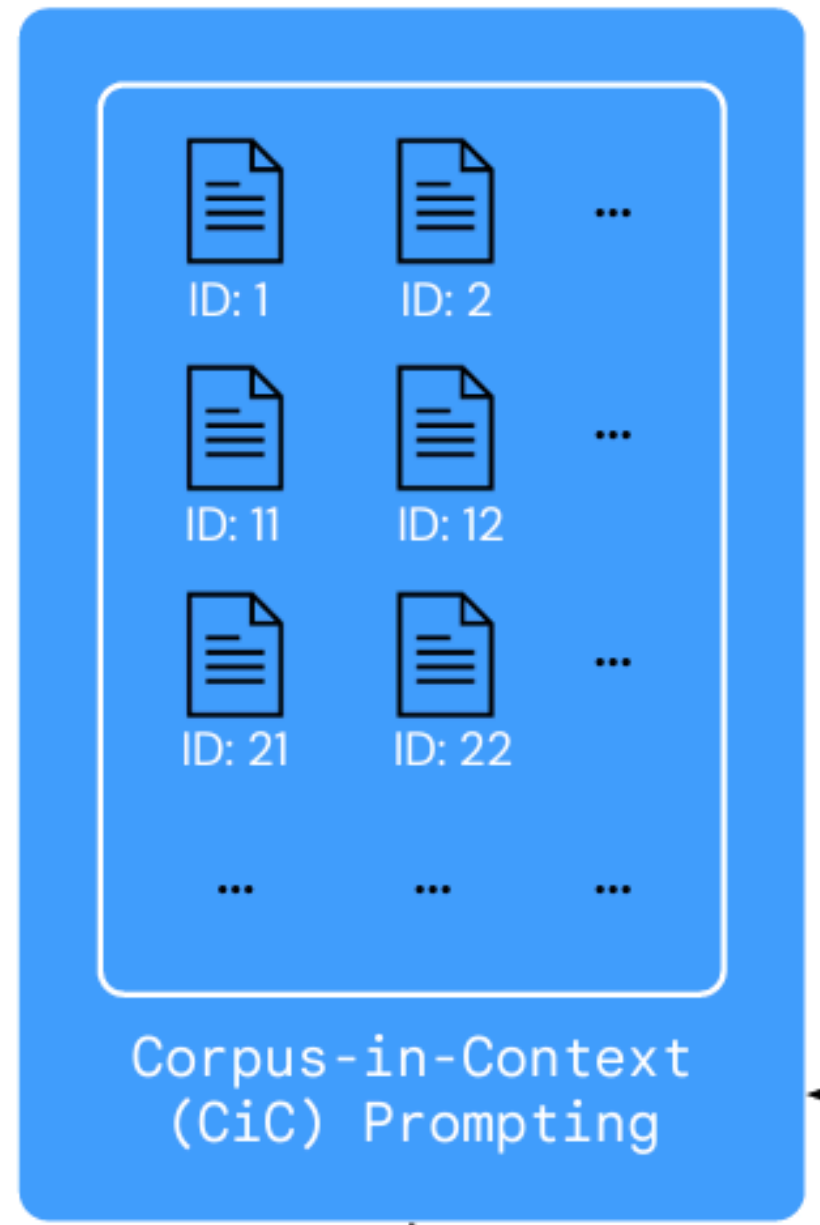
Few-shot Examples

Query Formatting

Likely to improve further in future

# Reflections on: “Can Long-Context Language Models Subsume Retrieval, RAG, SQL, and More?”

Long-Context Language Models  
(e.g. Gemini 1.5 Pro, GPT-4o)



But there is more!!

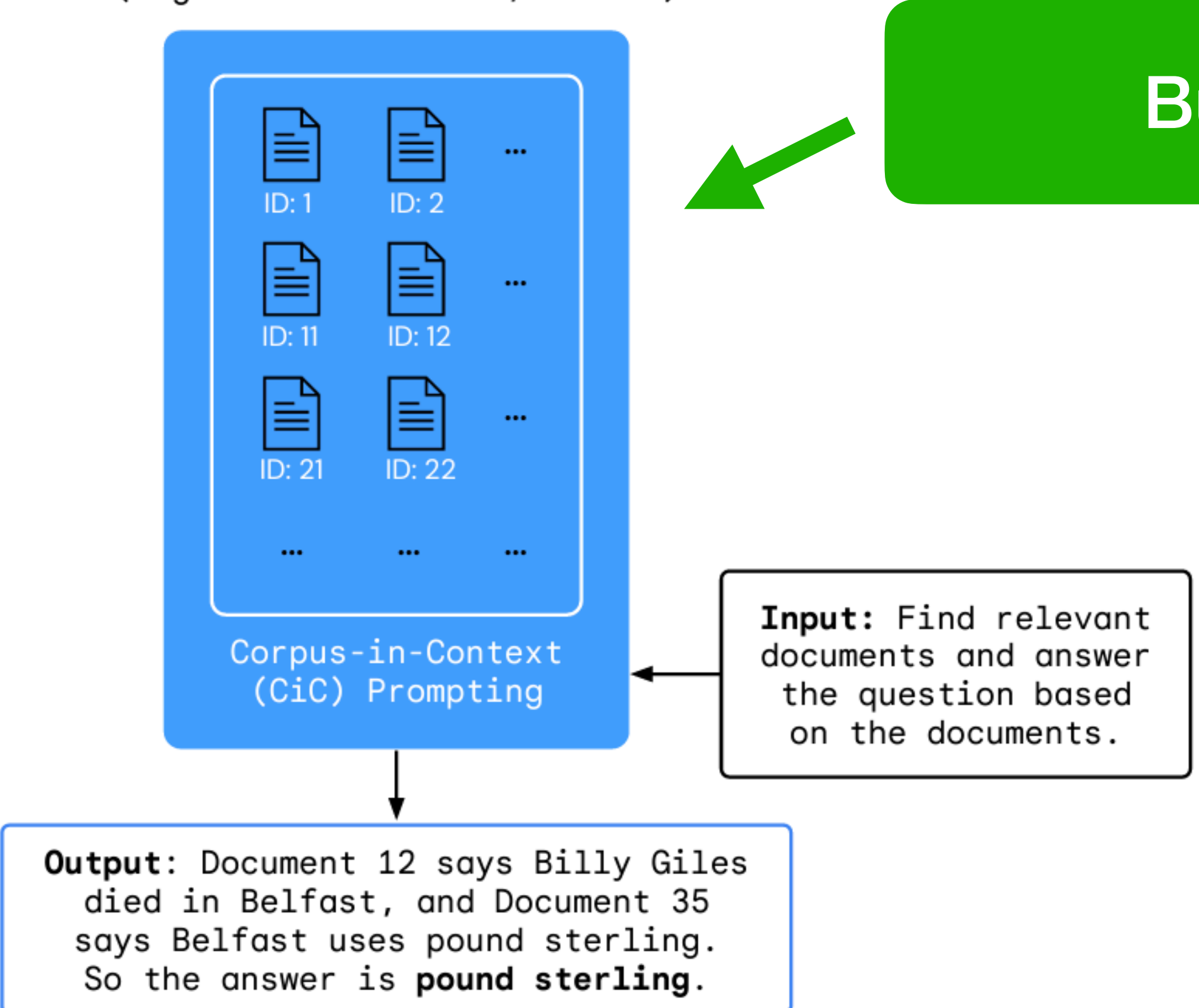
Compared to RAG:  
No multi-components framework to maintain/train/run

**Input:** Find relevant documents and answer the question based on the documents.

**Output:** Document 12 says Billy Giles died in Belfast, and Document 35 says Belfast uses pound sterling. So the answer is **pound sterling**.

# Reflections on: “Can Long-Context Language Models Subsume Retrieval, RAG, SQL, and More?”

Long-Context Language Models  
(e.g. Gemini 1.5 Pro, GPT-4o)



But there is more!!

Compared to RAG:  
No multi-components framework to maintain/train/run

Compared to DSI:  
No Training!  
("Indexing")

“Zero-shot”!

- Many promising directions for using LLMs to empower search
- Some are radical changes to “*retrieval pipeline status-quo*” and to how we see IR
  - Don’t get tricked by the view that IR is document retrieval! Key to IR is information, not documents
  - i.e. it’s time to move on from Lancaster 1968: “An IR system does not inform (i.e. change the knowledge of) the user on the subject of his inquiry. It merely informs on the existence (or non-existence) and whereabouts of documents relating to his request.”

No better time to do IR than now!