

# R<sup>2</sup>LLMs: Retrieval and Ranking with LLMs

**Guido Zuccon<sup>1</sup>, Shengyao Zhuang<sup>2</sup>, Xueguang Ma<sup>3</sup>, Bevan Koopman<sup>1,2</sup>**

<sup>1</sup> ielab, The University of Queensland, Australia & Google Research Australia

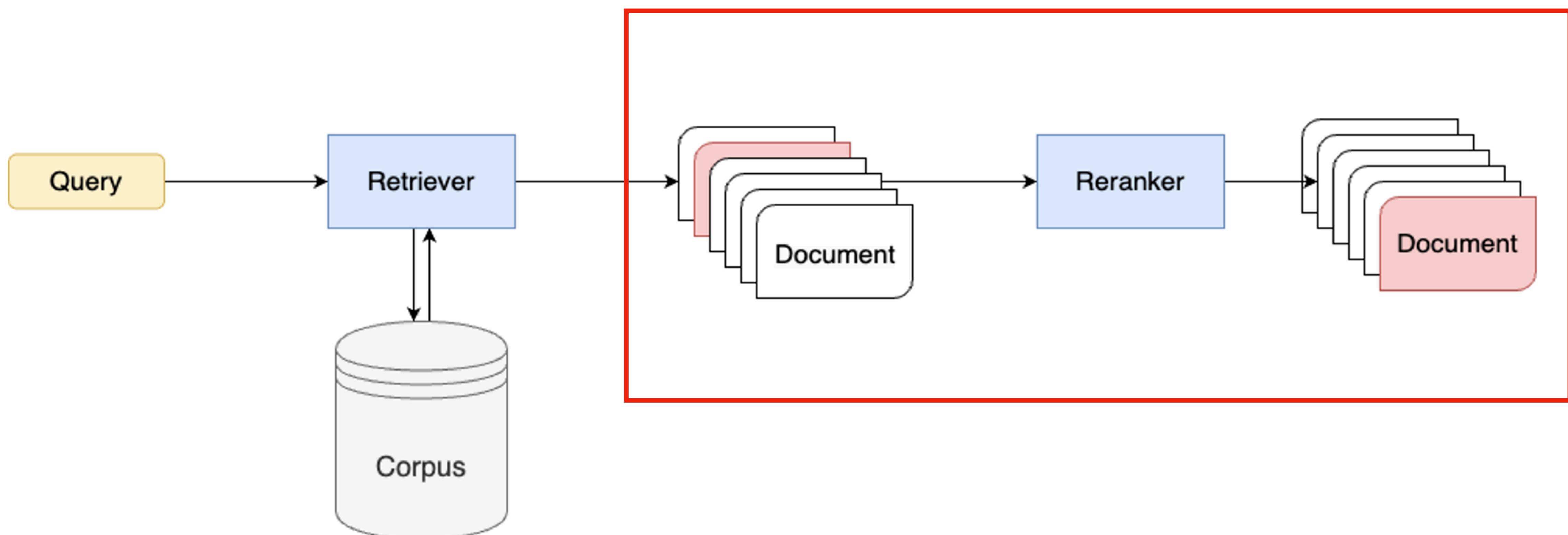
<sup>2</sup> ielab, CSIRO, Australia

<sup>3</sup> The University of Waterloo

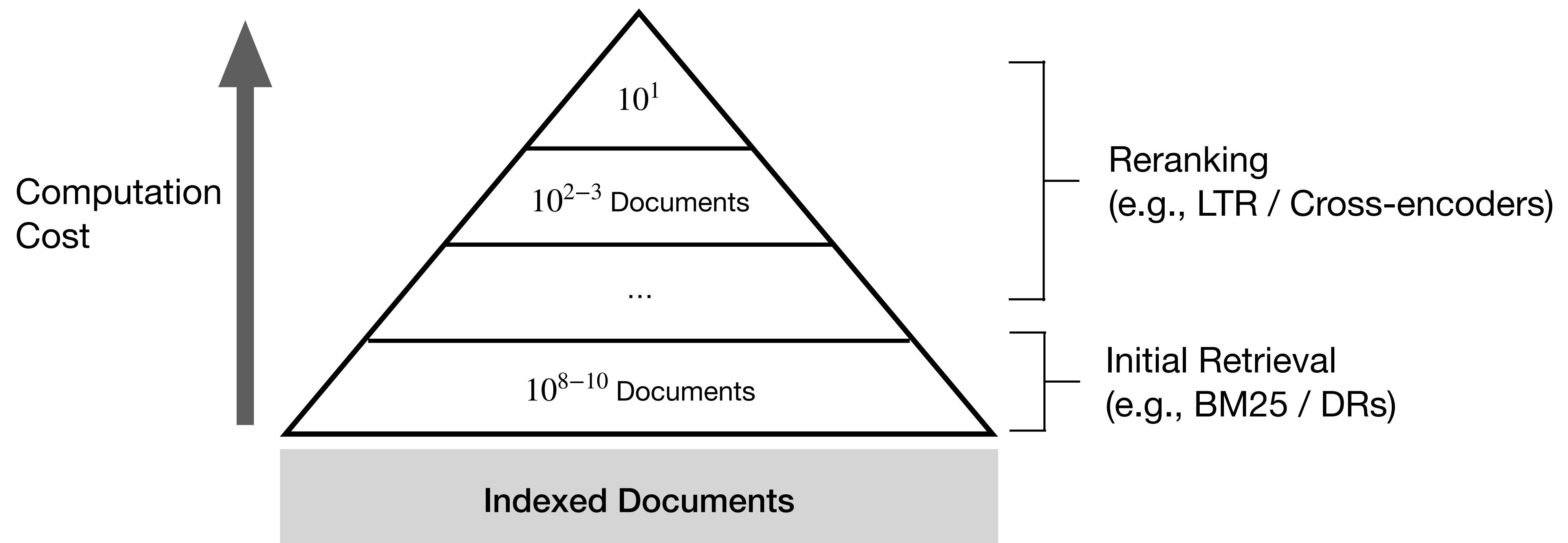
<https://ielab.io/tutorials/r2llms.html>

# Part 3: LLM-based Rankers

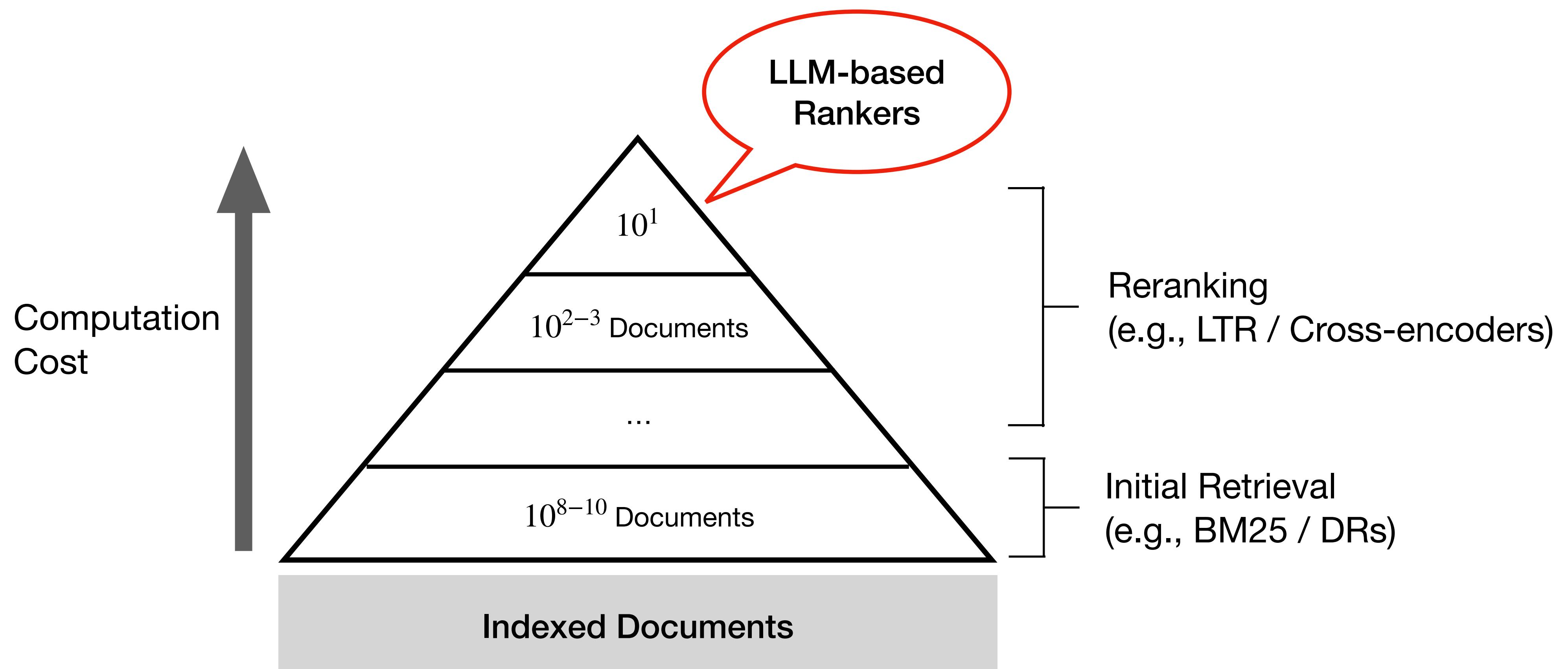
# PART 3: LLM-based Rankers



# PART 3: LLM-based Rankers

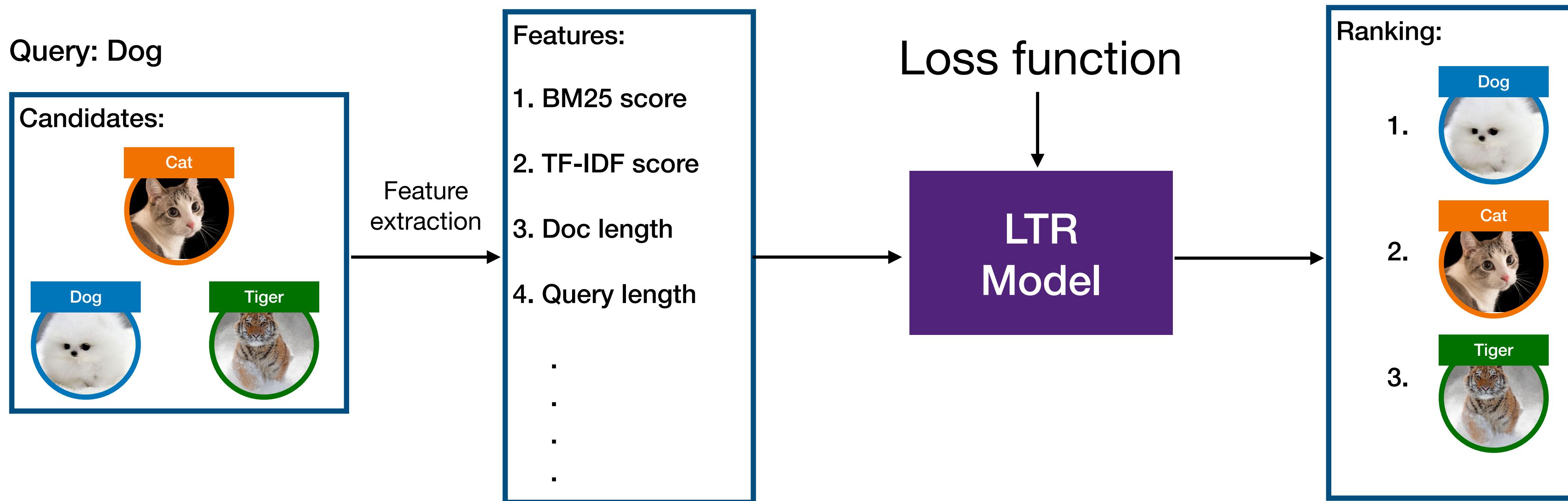


# PART 3: LLM-based Rankers



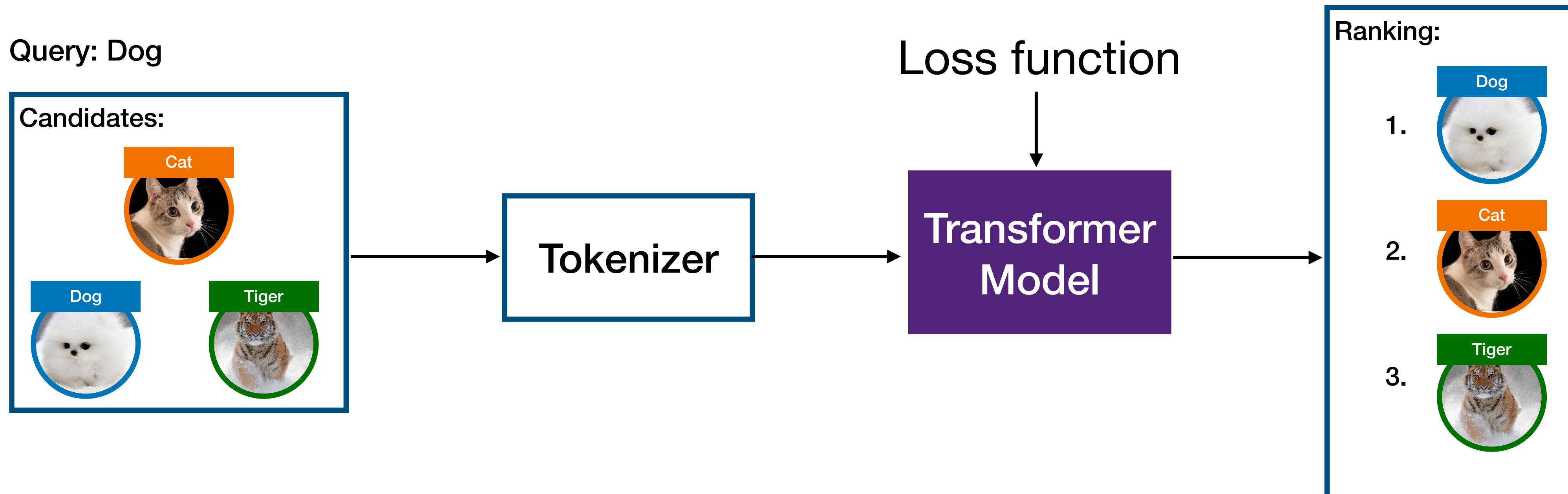
# Recap of LTR and Cross-encoders

- Learning to Rank (LTR):



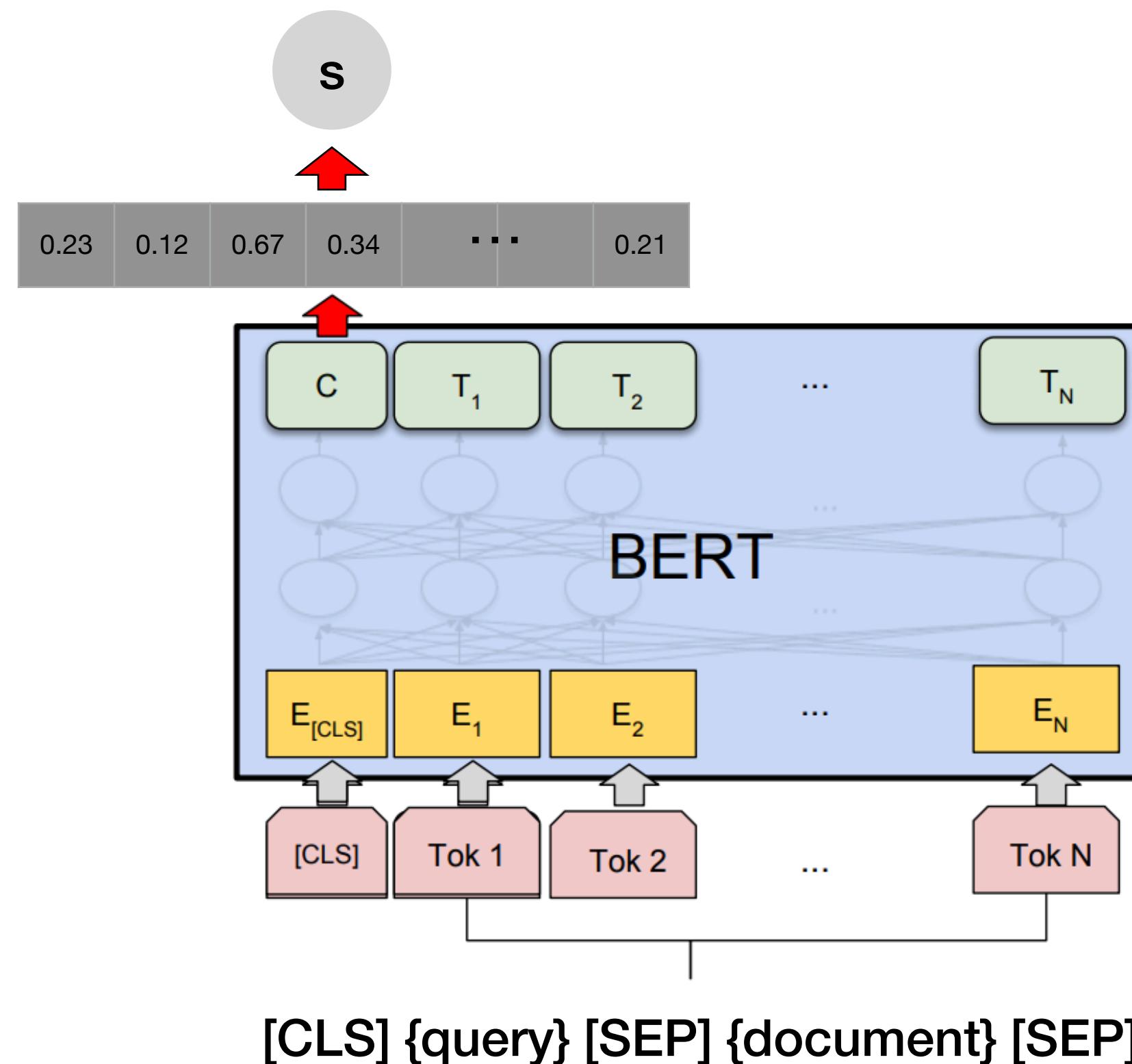
# Recap of LTR and Cross-encoders

- Transformer ranker



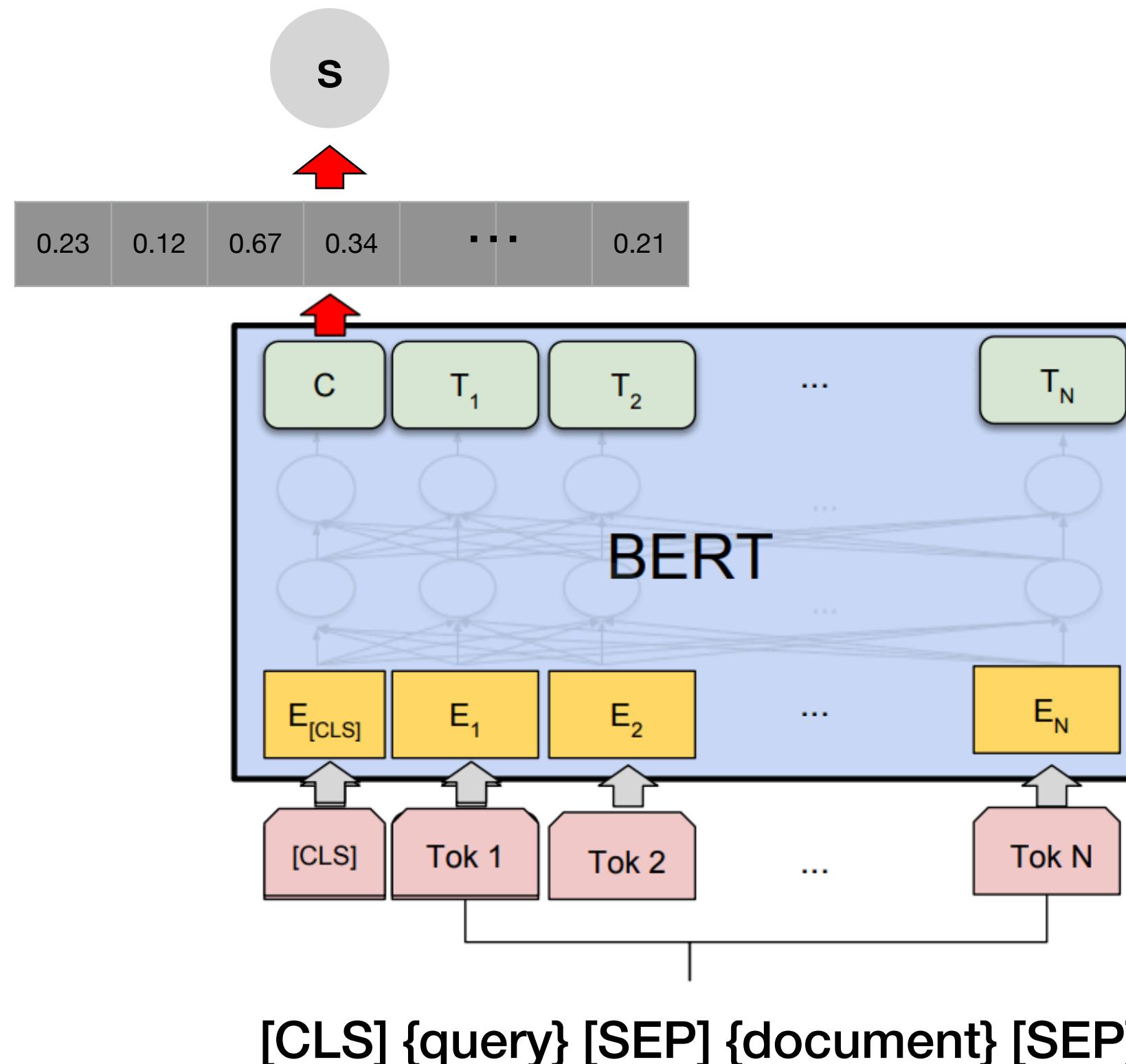
# Recap of LTR and Cross-encoders

- monoBERT: cross-encoder architecture



# Recap of LTR and Cross-encoders

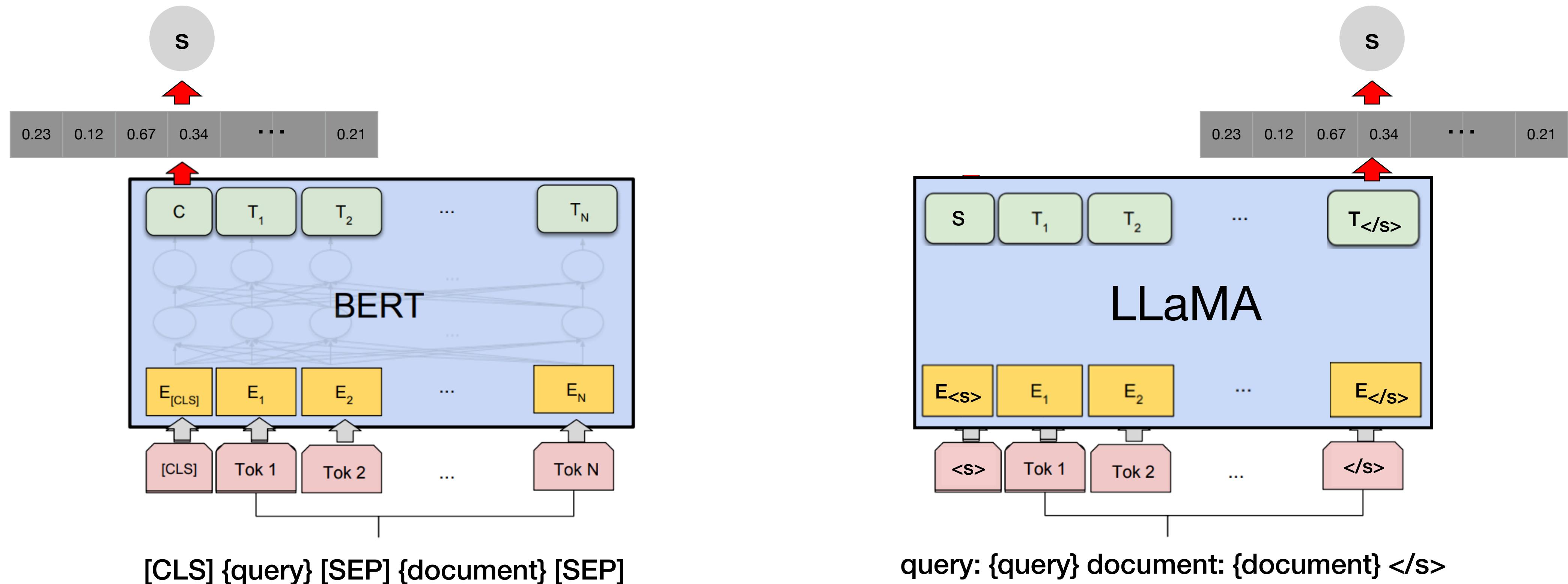
- monoBERT: cross-encoder architecture



LLM-based cross-encoder?

# LLM-based Cross-encoder

- RankLLaMA



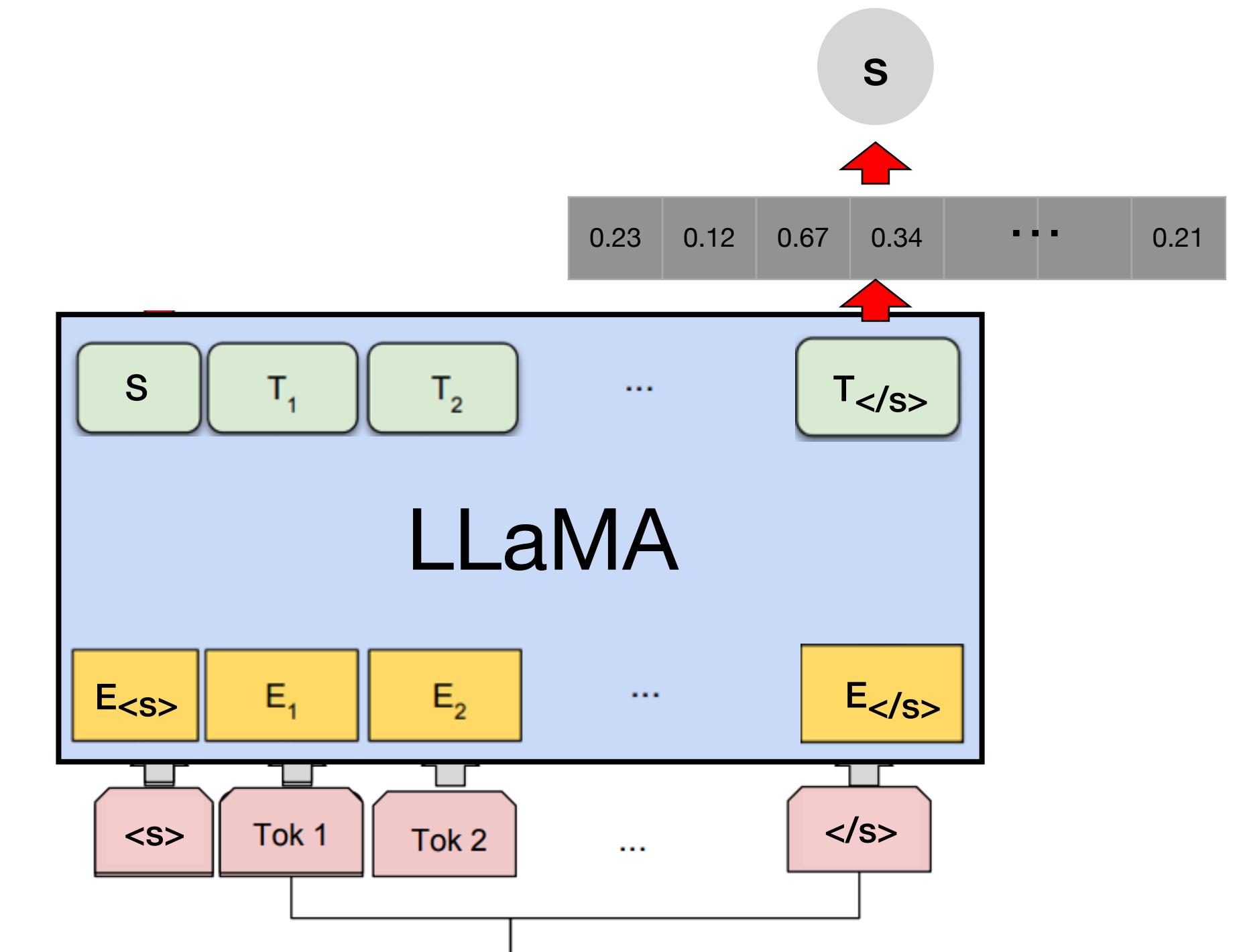
# LLM-based Cross-encoder

- RankLLaMA

input = ‘query:  $\{Q\}$  document:  $\{D\} </s>$ ’

$$\text{Sim}(Q, D) = \text{Linear}(\text{Decoder}(\text{input})[-1])$$

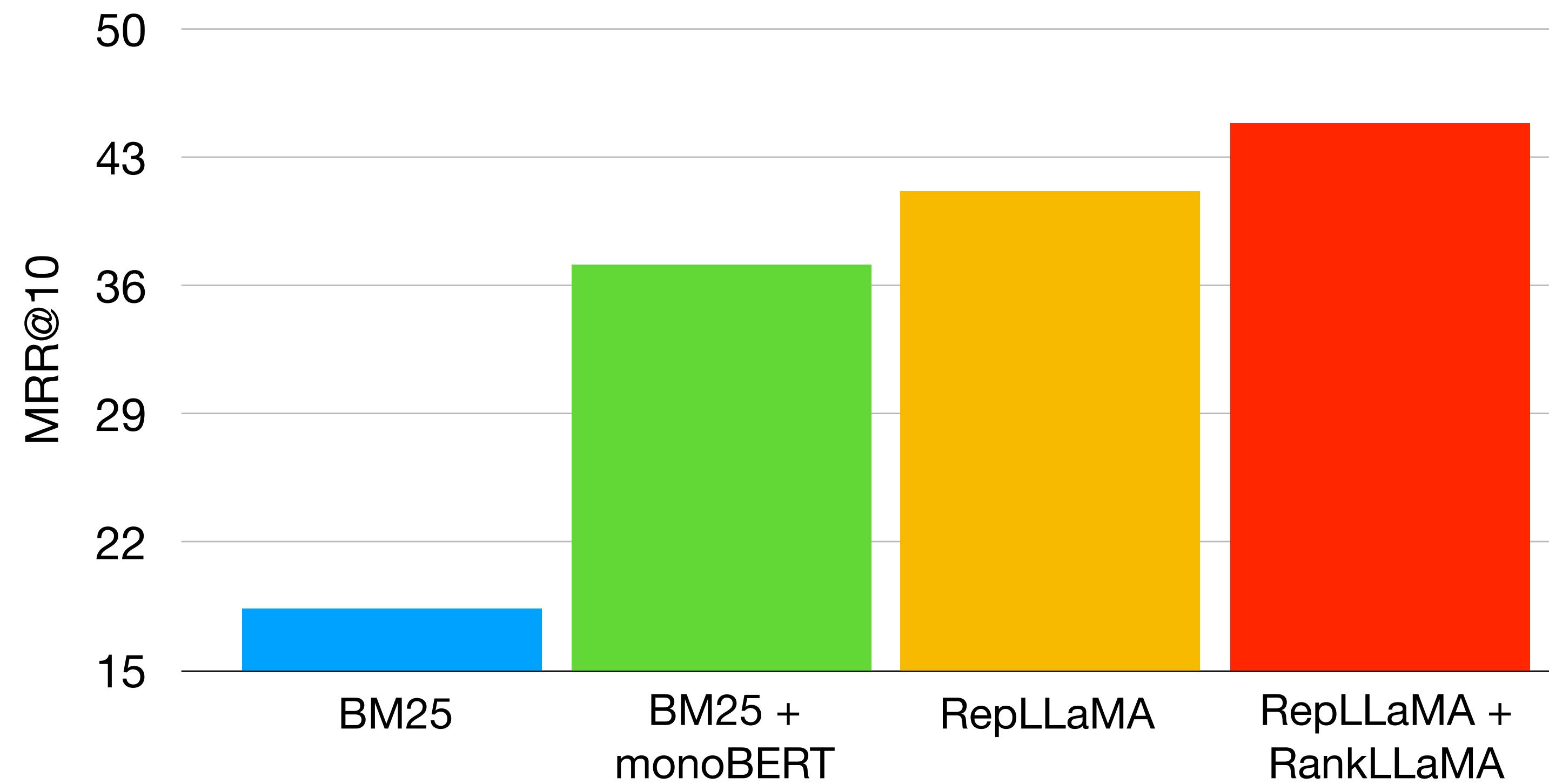
$$\begin{aligned}\mathcal{L}(Q, D^+, \{D_N\}) &= -\log p(D = D^+ | Q) = \\ &- \log \frac{\exp(\text{Sim}(Q, D^+))}{\exp(\text{Sim}(Q, D^+)) + \sum_{D_i^- \in \{D_N\}} \exp(\text{Sim}(Q, D_i^-))}\end{aligned}$$



query: {query} document: {document} </s>

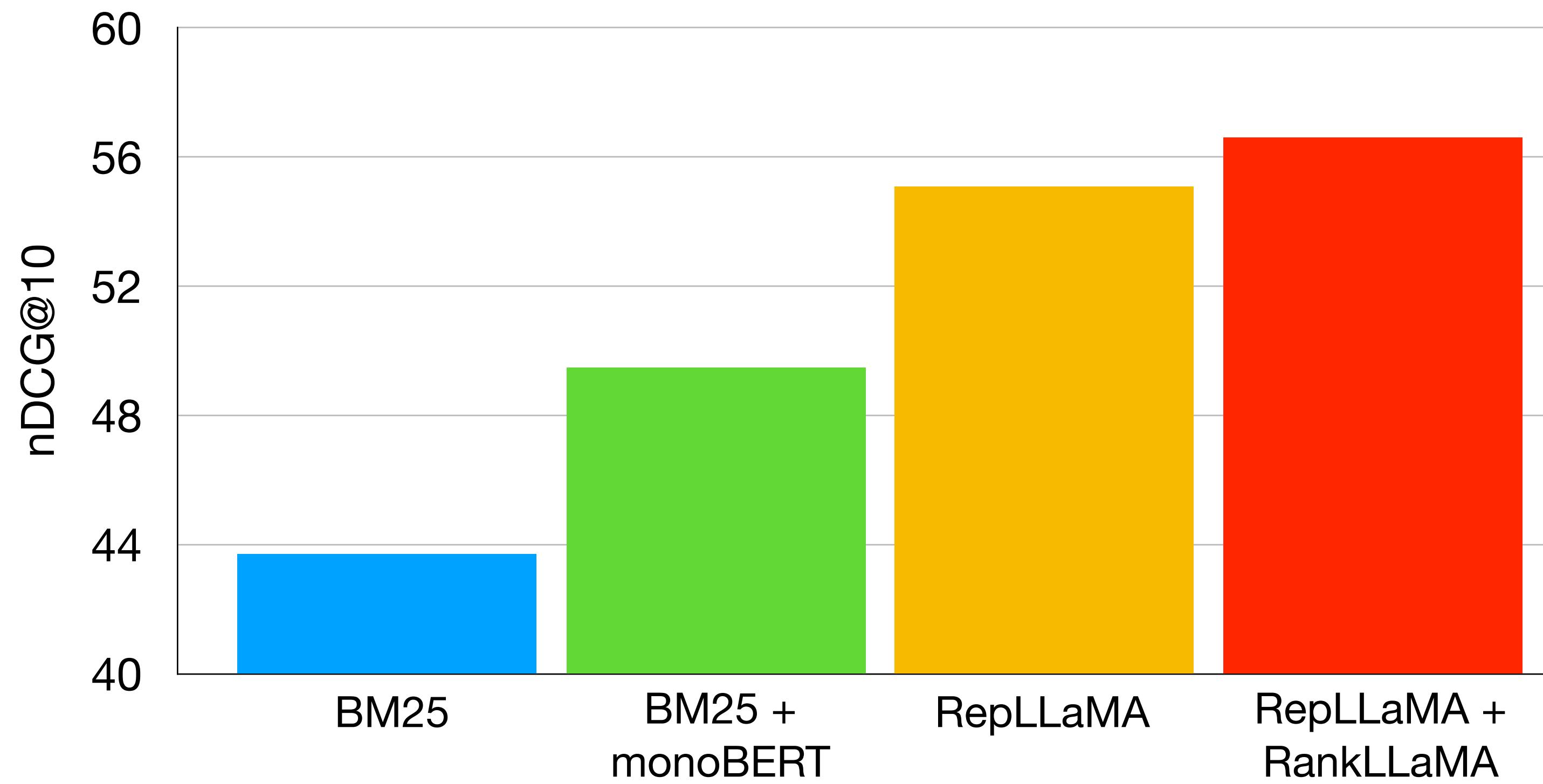
# RankLLaMA Results

- In-domain results: MSMARCO



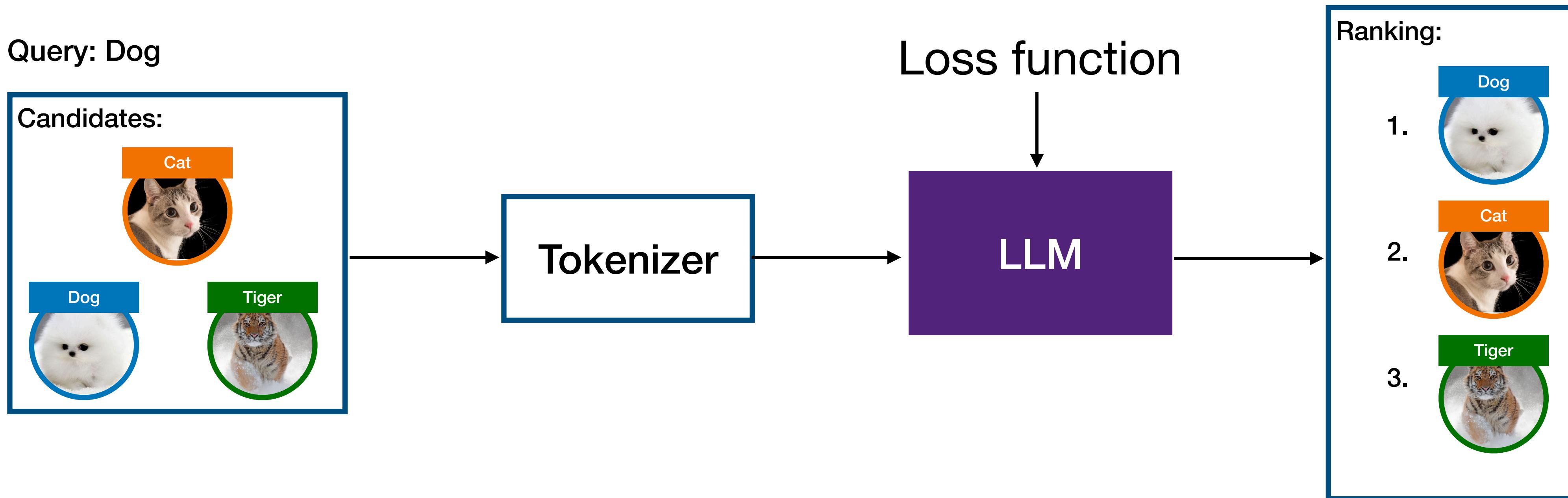
# RankLLaMA Results

- Out-of-domain results: BEIR



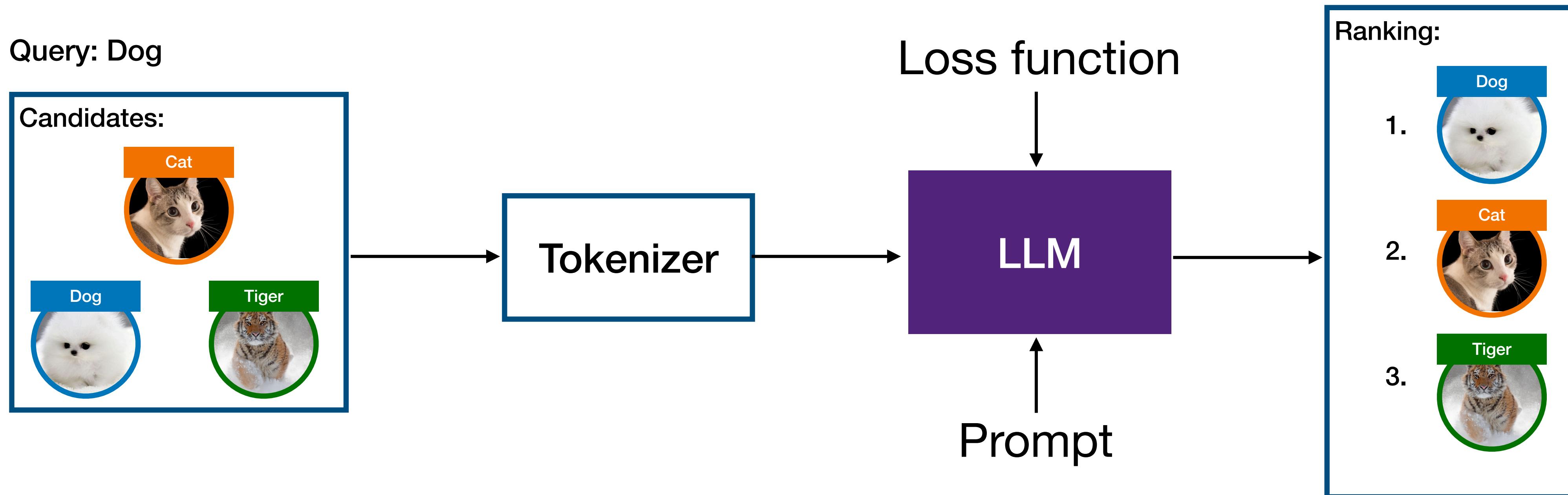
# LLM-based Cross-encoder

- LLM-based ranker

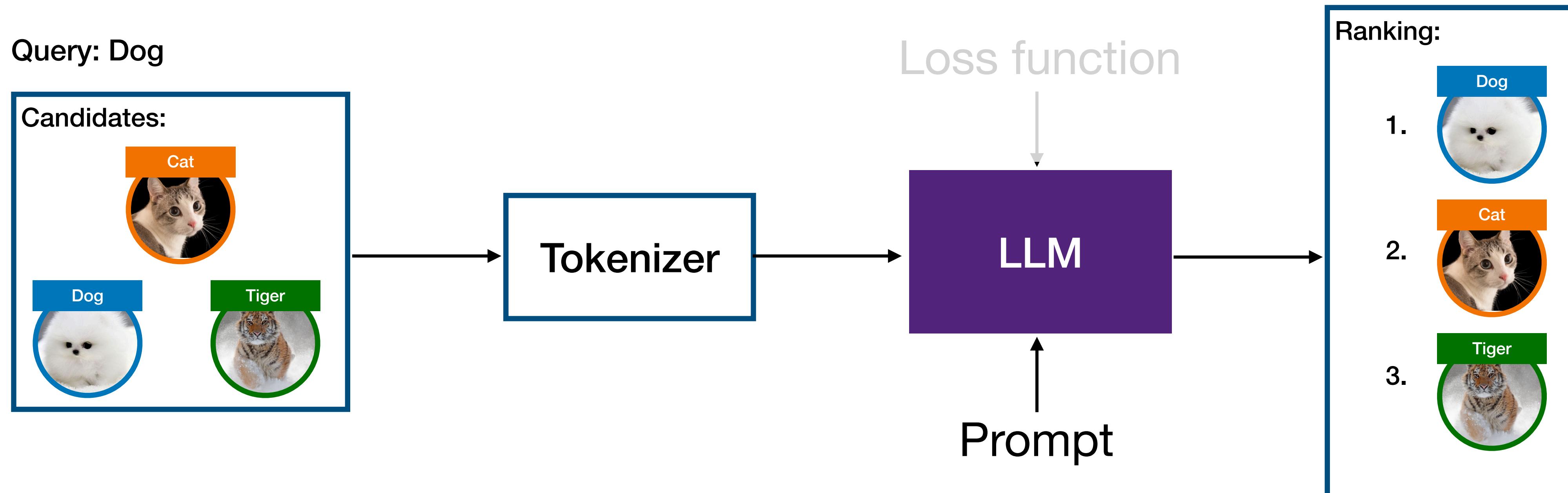


# LLM-based Cross-encoder

- LLM-based ranker



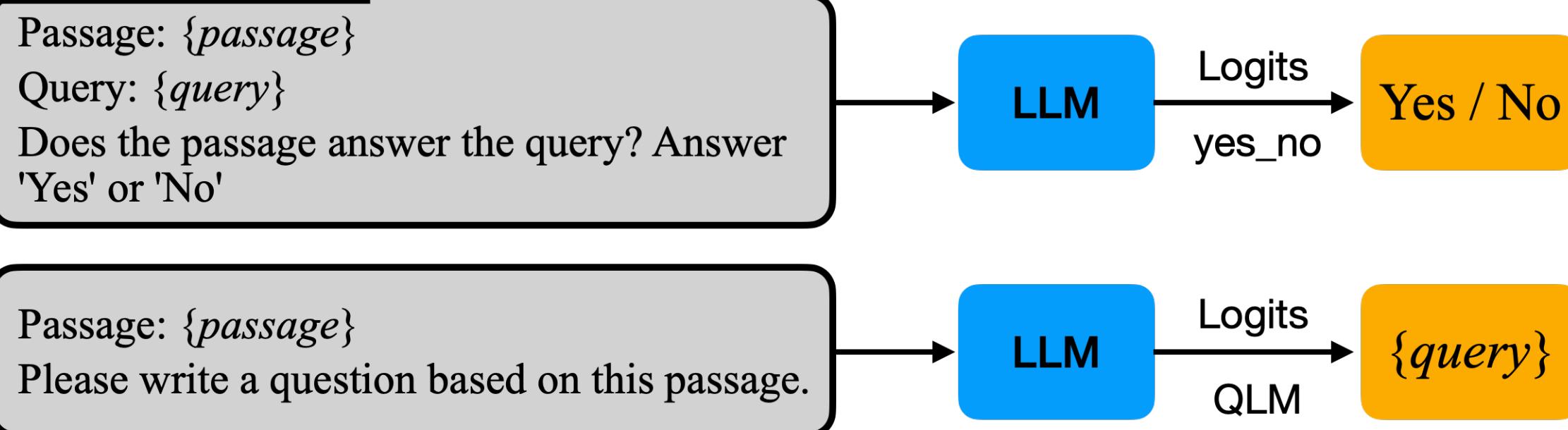
# Zeroshot LLM-based Rankers (LLM as Rankers)



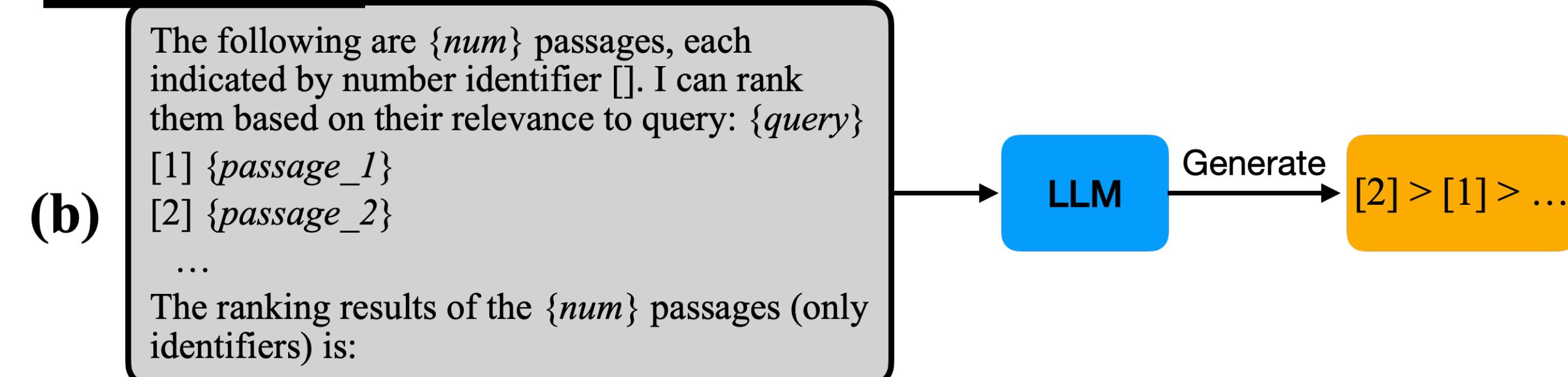
- All are “zero-shot”: i.e. once you obtained the pre-trained, instruction tuned LLM, no need to do further post-training.
- Input could be one or multiple documents, depending on the prompt.

# Zeroshot LLM-based Rankers: 3.5 types

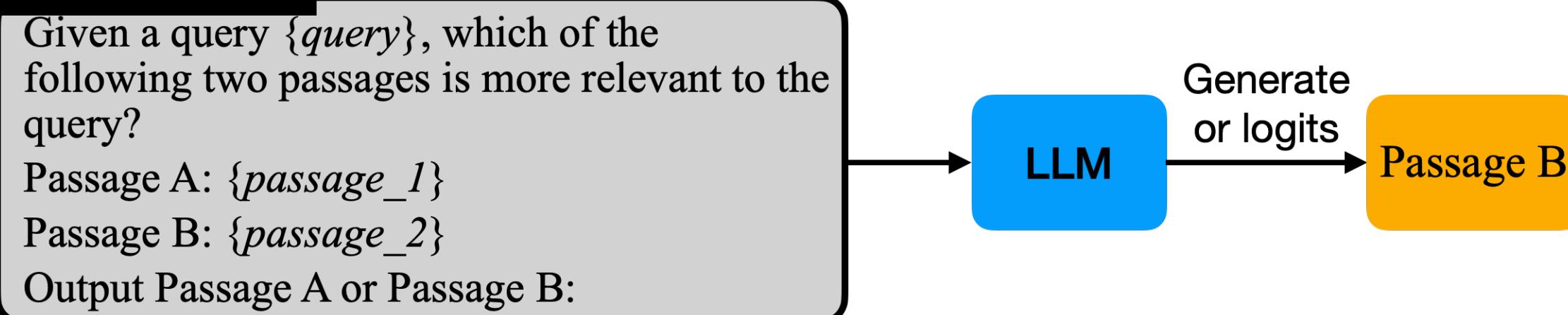
## Pointwise



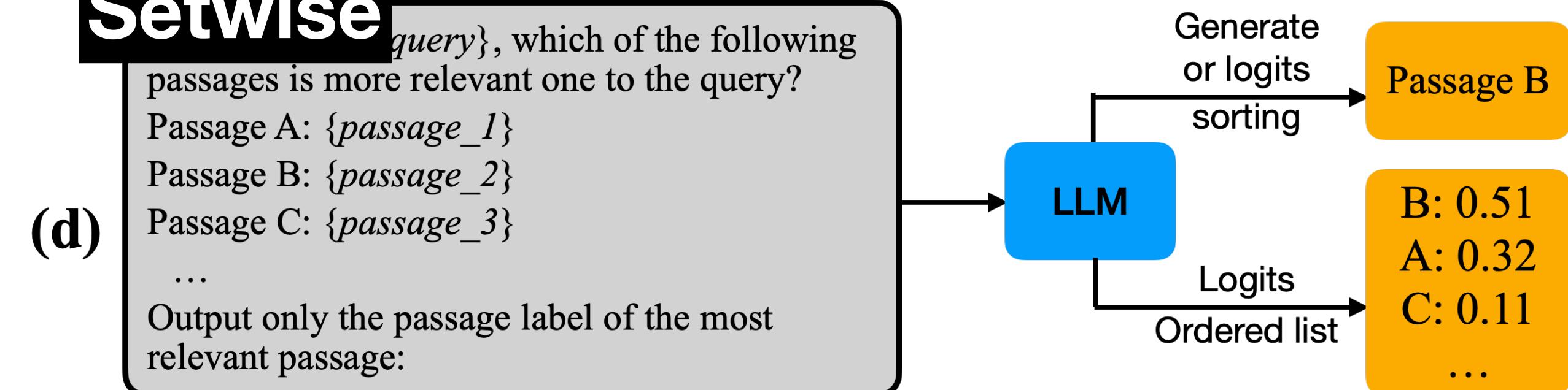
## Listwise



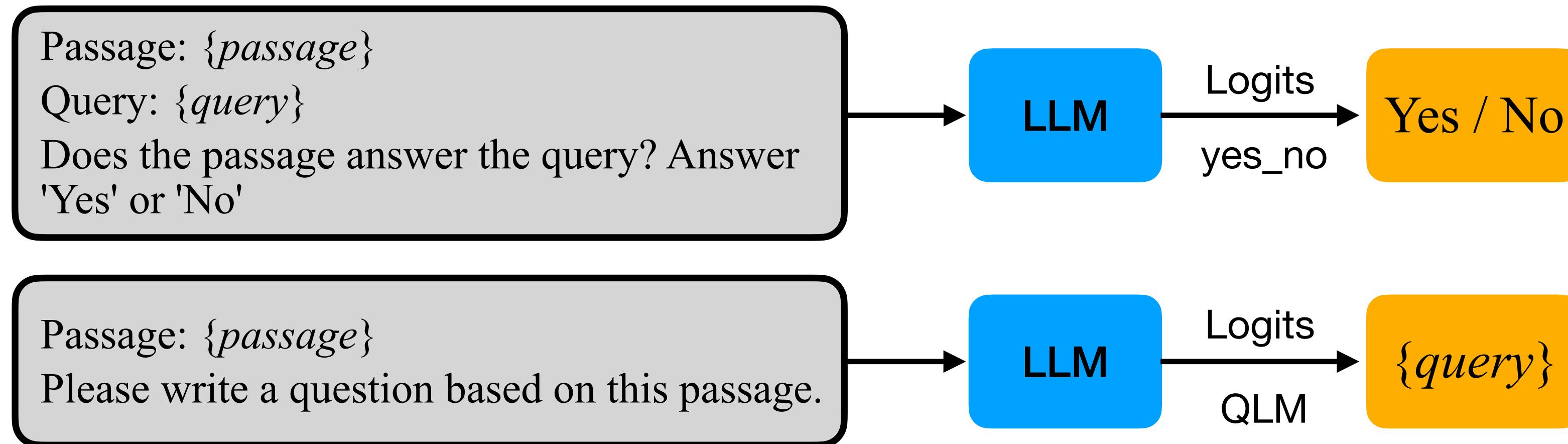
## Pairwise



## Setwise



# Pointwise



- Yes / No label generation.
- Query generation likelihood.

# Pointwise

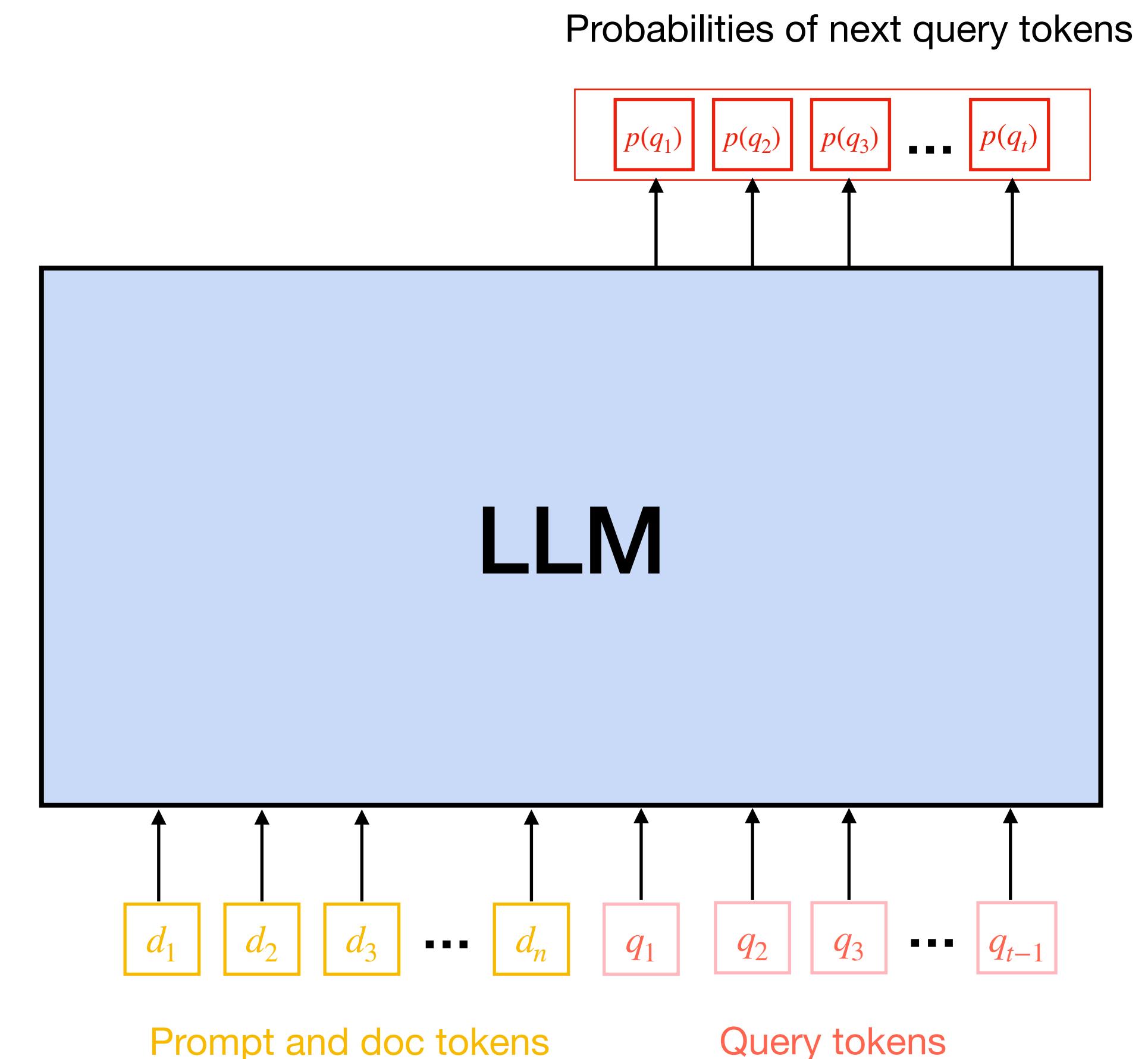
- Query generation likelihood.

Passage:  $\{passage\}$

Please write a question based on this passage.

**Rank by query likelihood:**

$$P(Q | D) = \sum_i^t \log p(q_i)$$



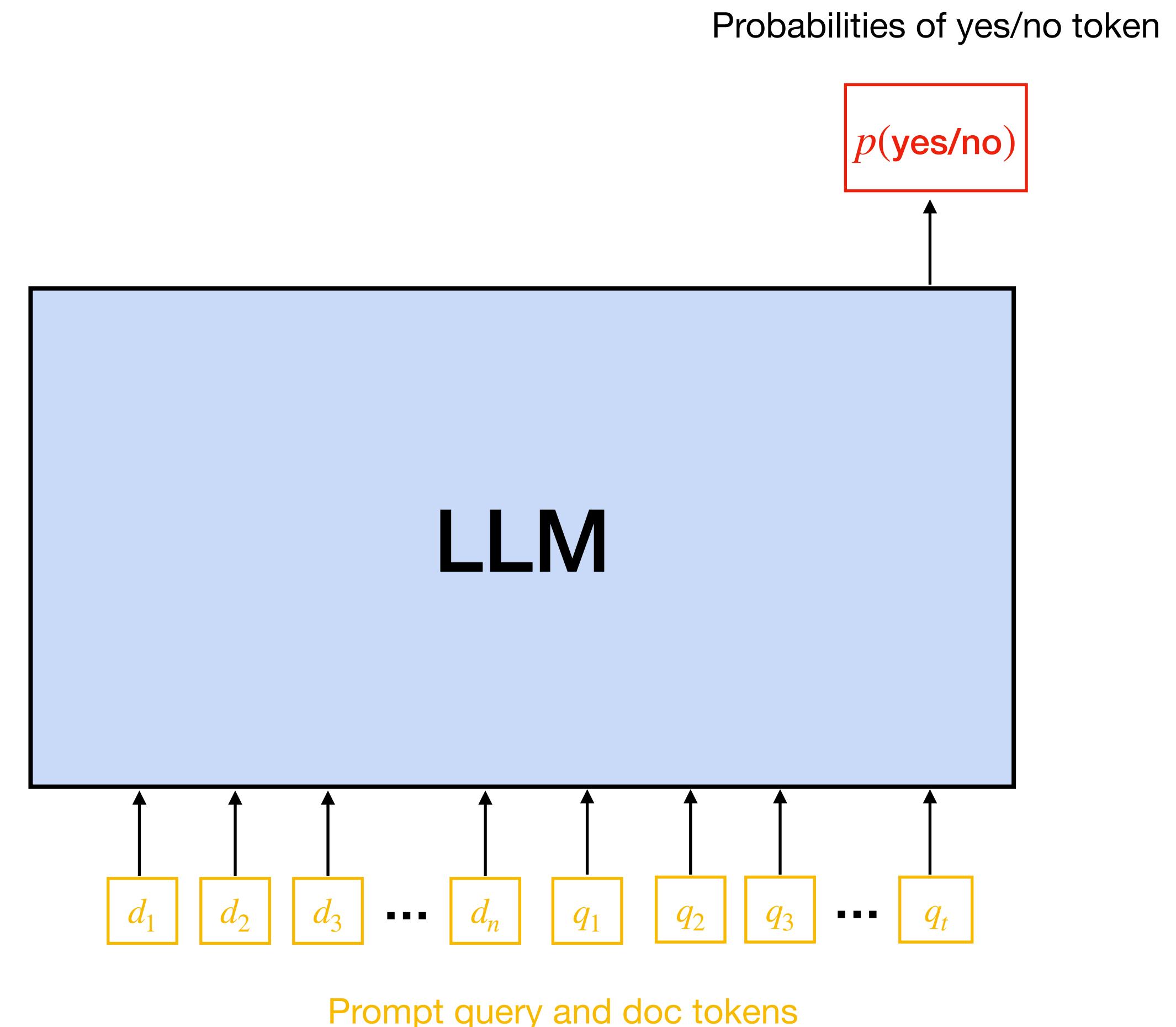
# Pointwise

- Yes / No label generation.

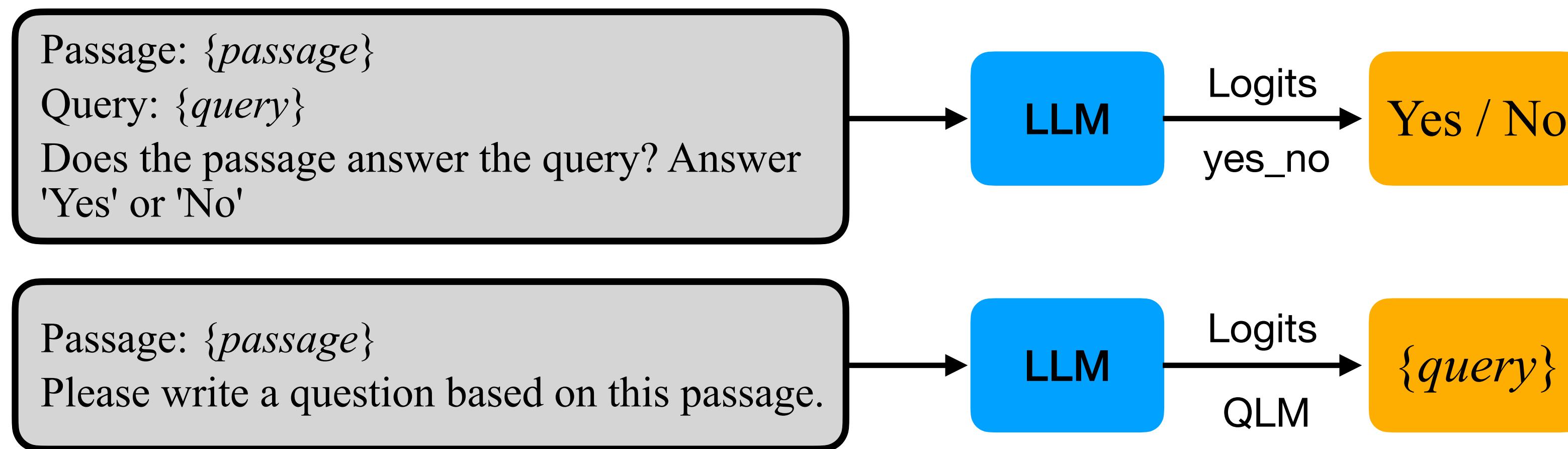
Passage:  $\{passage\}$   
Query:  $\{query\}$   
Does the passage answer the query? Answer  
'Yes' or 'No'

**Rank by normalised yes probability:**

$$P(\text{yes} | Q, D) = \frac{\log p(\text{yes})}{\log p(\text{yes}) + \log p(\text{no})}$$



# Pointwise



- Each query-doc pair can be processed in parallel.
- It needs access to the model weights (need token logits).

# Listwise

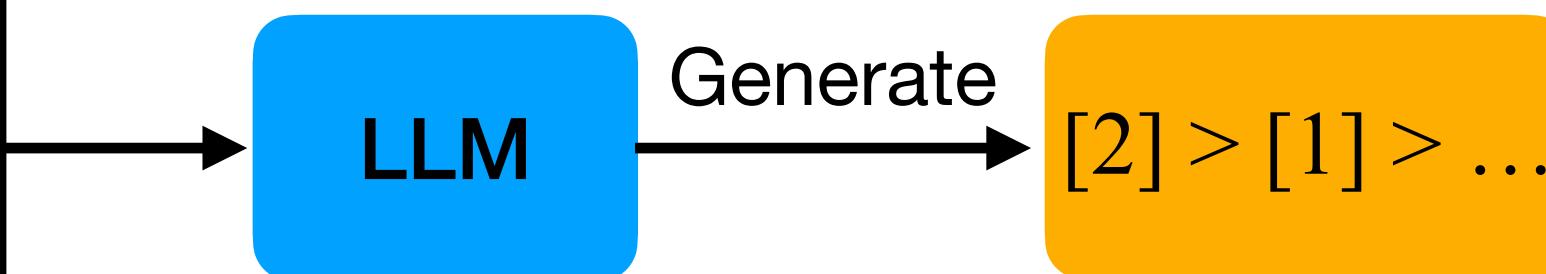
The following are  $\{num\}$  passages, each indicated by number identifier []. I can rank them based on their relevance to query:  $\{query\}$

[1]  $\{passage\_1\}$

[2]  $\{passage\_2\}$

...

The ranking results of the  $\{num\}$  passages (only identifiers) is:



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

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

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. *EMNLP*.

# Listwise

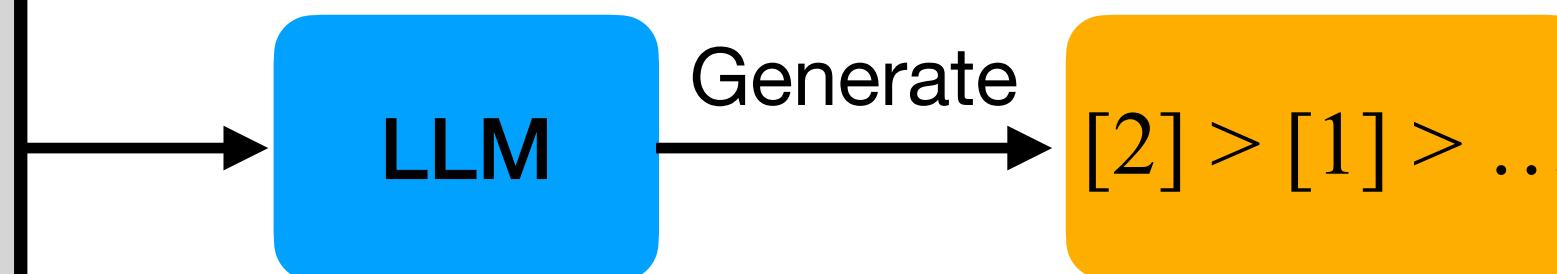
The following are  $\{num\}$  passages, each indicated by number identifier []. I can rank them based on their relevance to query:  $\{query\}$

[1]  $\{passage\_1\}$

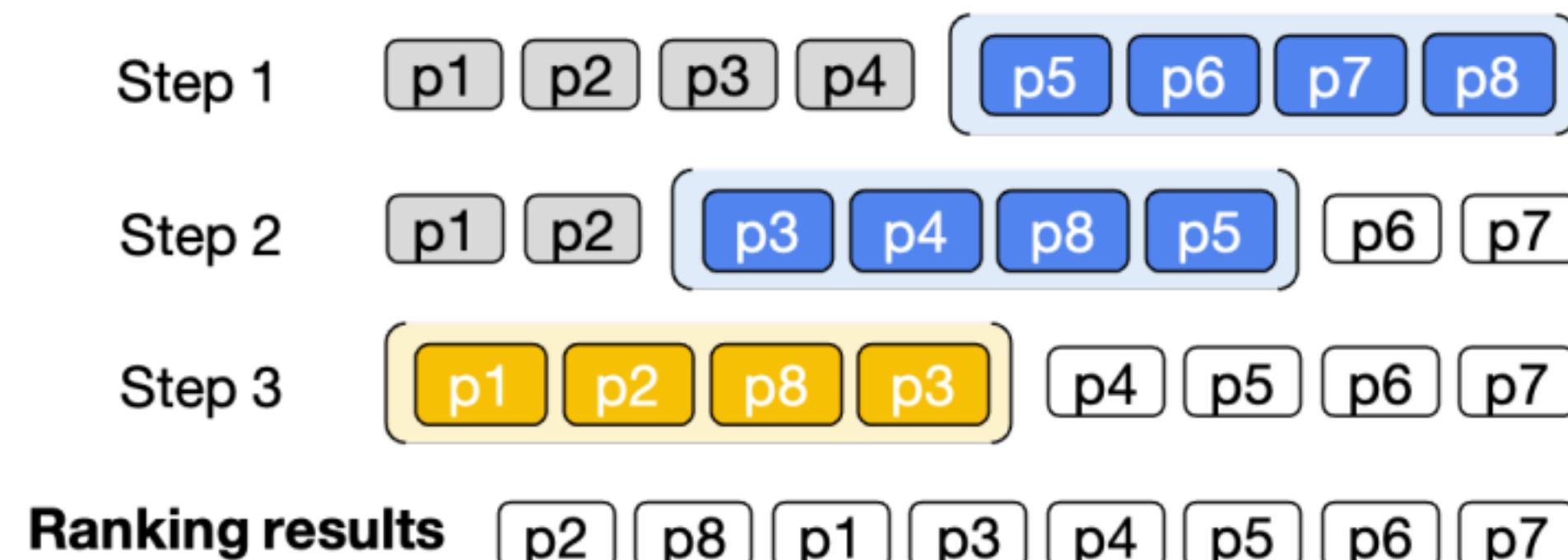
[2]  $\{passage\_2\}$

...

The ranking results of the  $\{num\}$  passages (only identifiers) is:



- Sliding window



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

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

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. *EMNLP*.

# Listwise

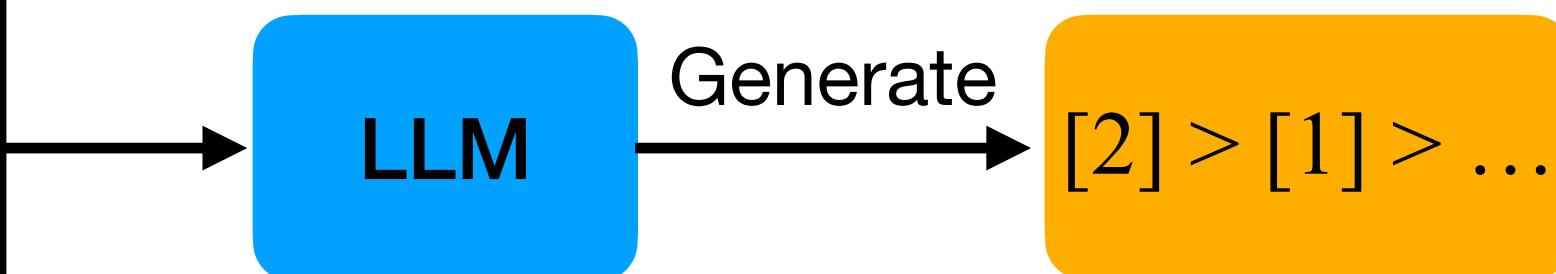
The following are  $\{num\}$  passages, each indicated by number identifier []. I can rank them based on their relevance to query:  $\{query\}$

[1]  $\{passage\_1\}$

[2]  $\{passage\_2\}$

...

The ranking results of the  $\{num\}$  passages (only identifiers) is:



- Cannot parallel.
- No need access to the model weights, pure generation.

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

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

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. *EMNLP*.

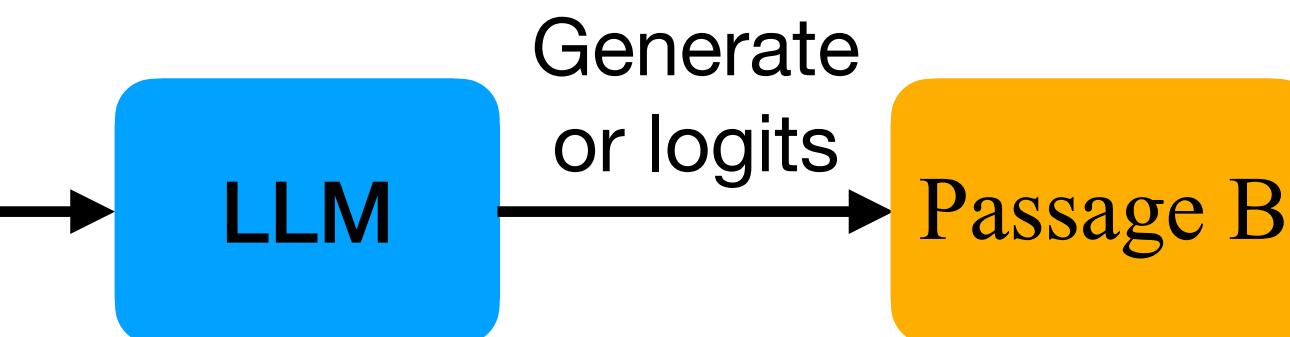
# Pairwise

Given a query  $\{query\}$ , which of the following two passages is more relevant to the query?

Passage A:  $\{passage\_1\}$

Passage B:  $\{passage\_2\}$

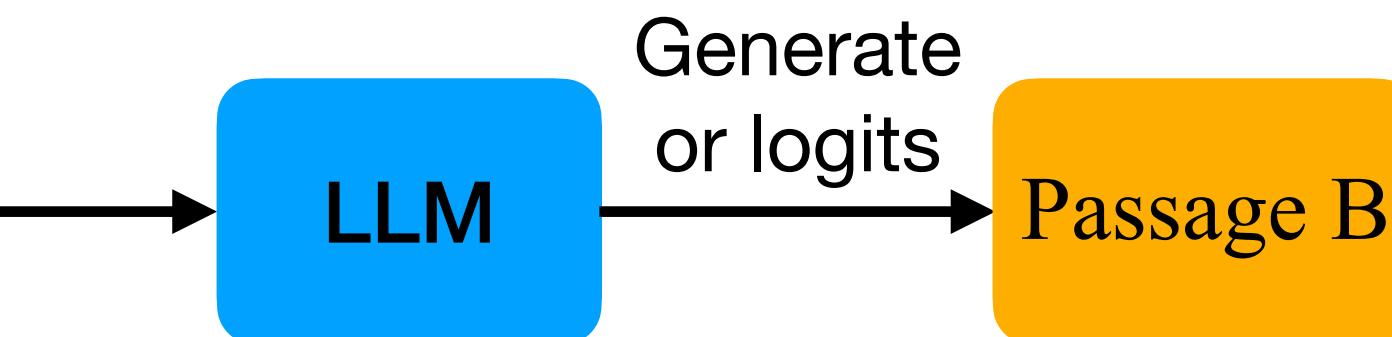
Output Passage A or Passage B:



Qin, Z., Jagerman, R., Hui, K., Zhuang, H., Wu, J., Shen, J., Liu, T., Liu, J., Metzler, D., Wang, X. and Bendersky, M., 2023. Large language models are effective text rankers with pairwise ranking prompting. *NAACL-finding*.

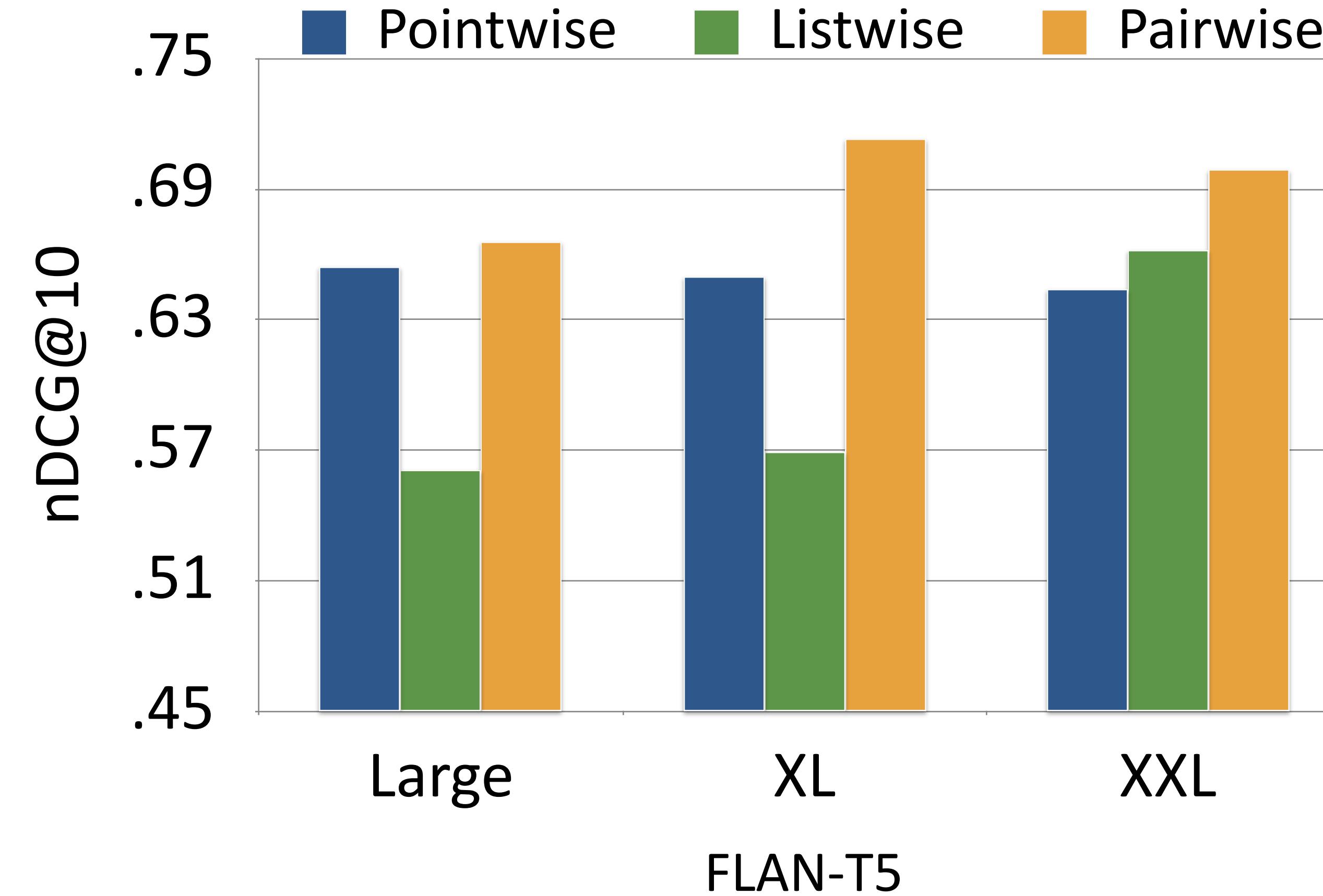
# Pairwise

Given a query  $\{query\}$ , which of the following two passages is more relevant to the query?  
Passage A:  $\{passage\_1\}$   
Passage B:  $\{passage\_2\}$   
Output Passage A or Passage B:

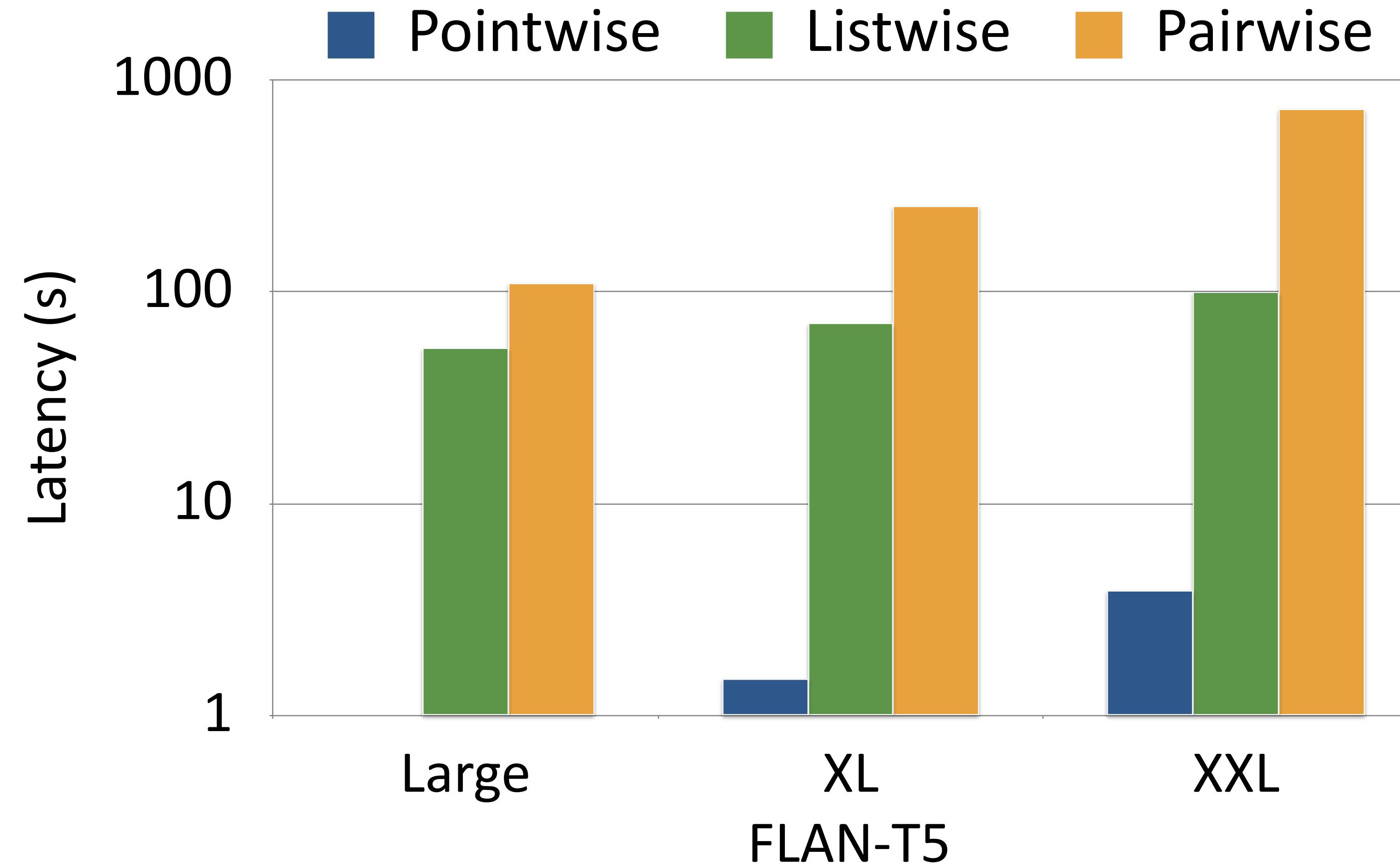


- **Scoring:** Compare all the pairwise preferences. This can be done in parallel.
- Both generation and logits are supported.

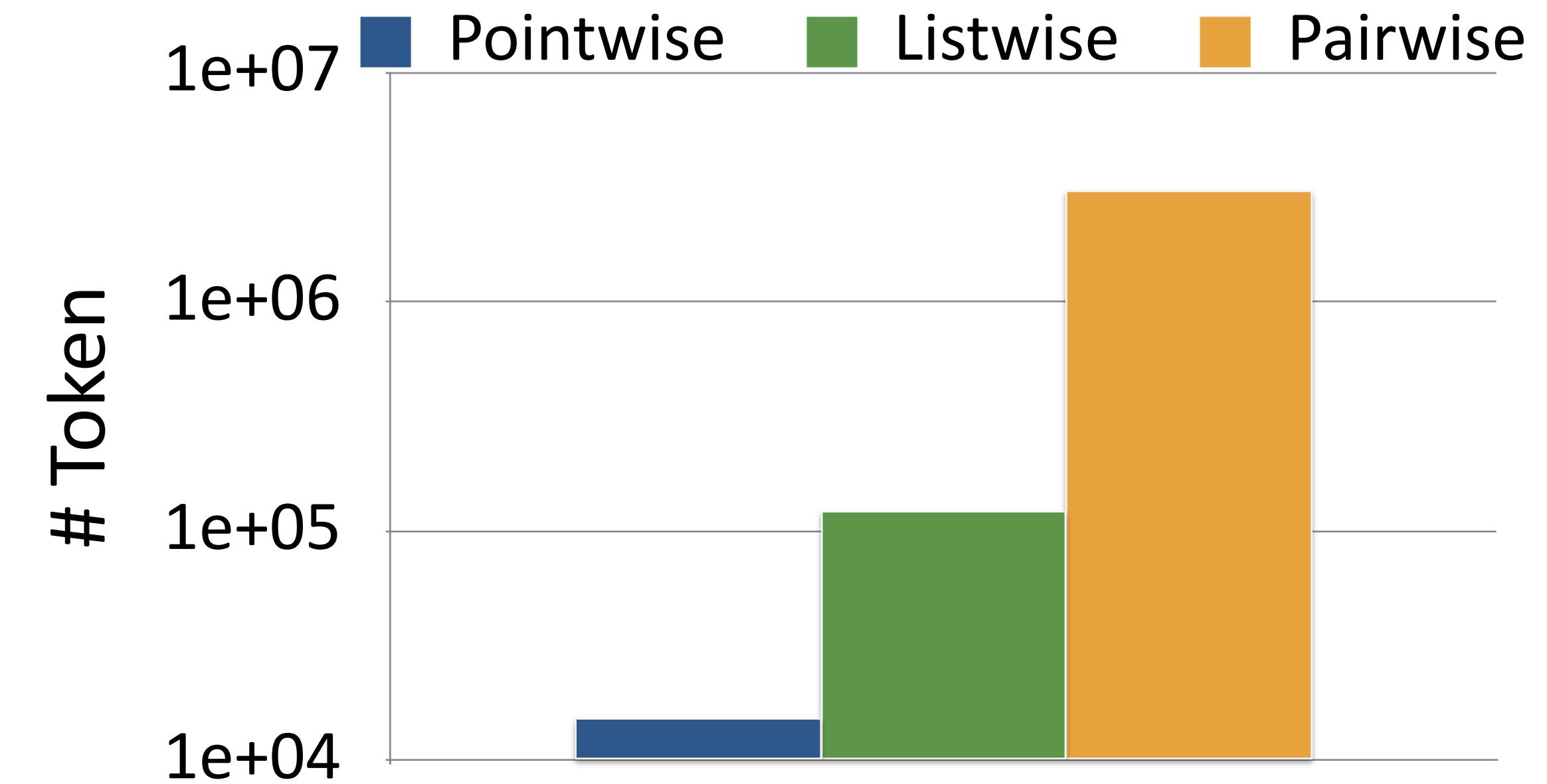
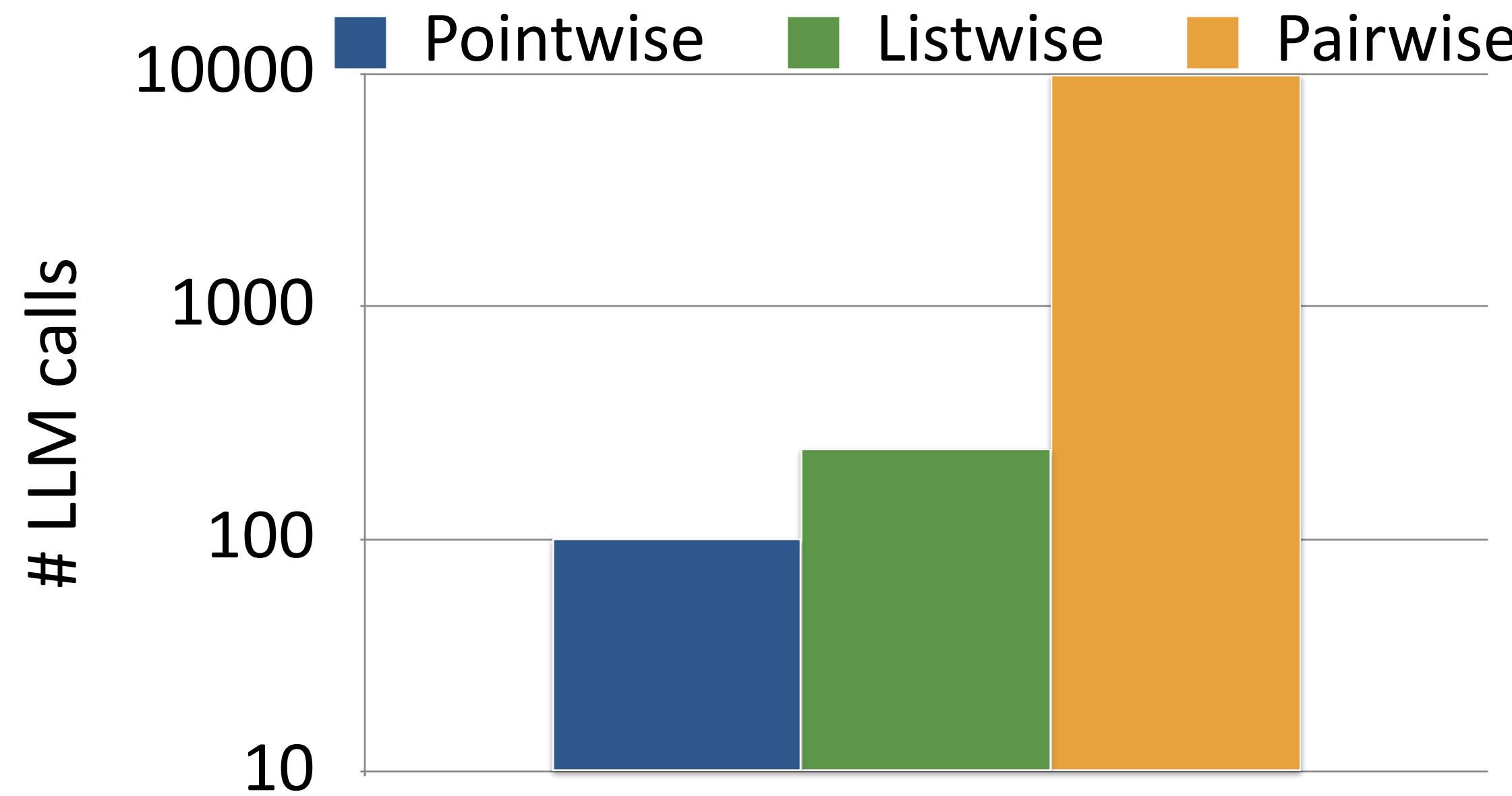
# Effectiveness



# Efficiency



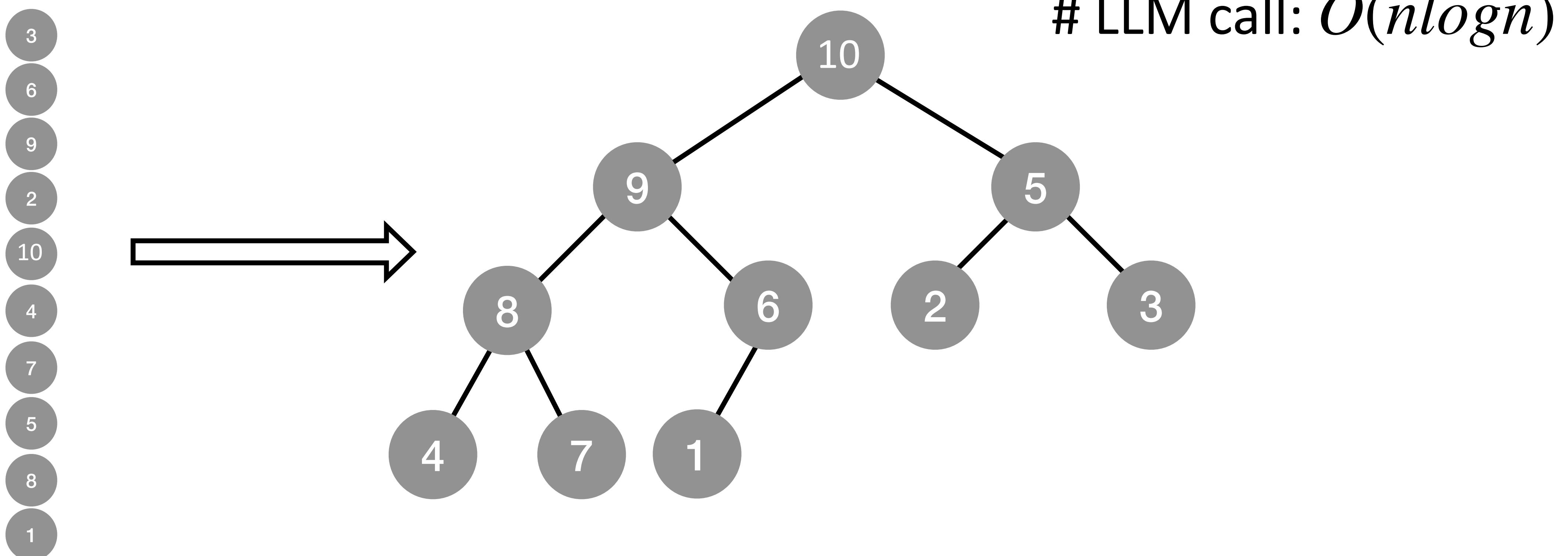
# Efficiency



# Pairwise

- Sorting-based ranking: Heapsort

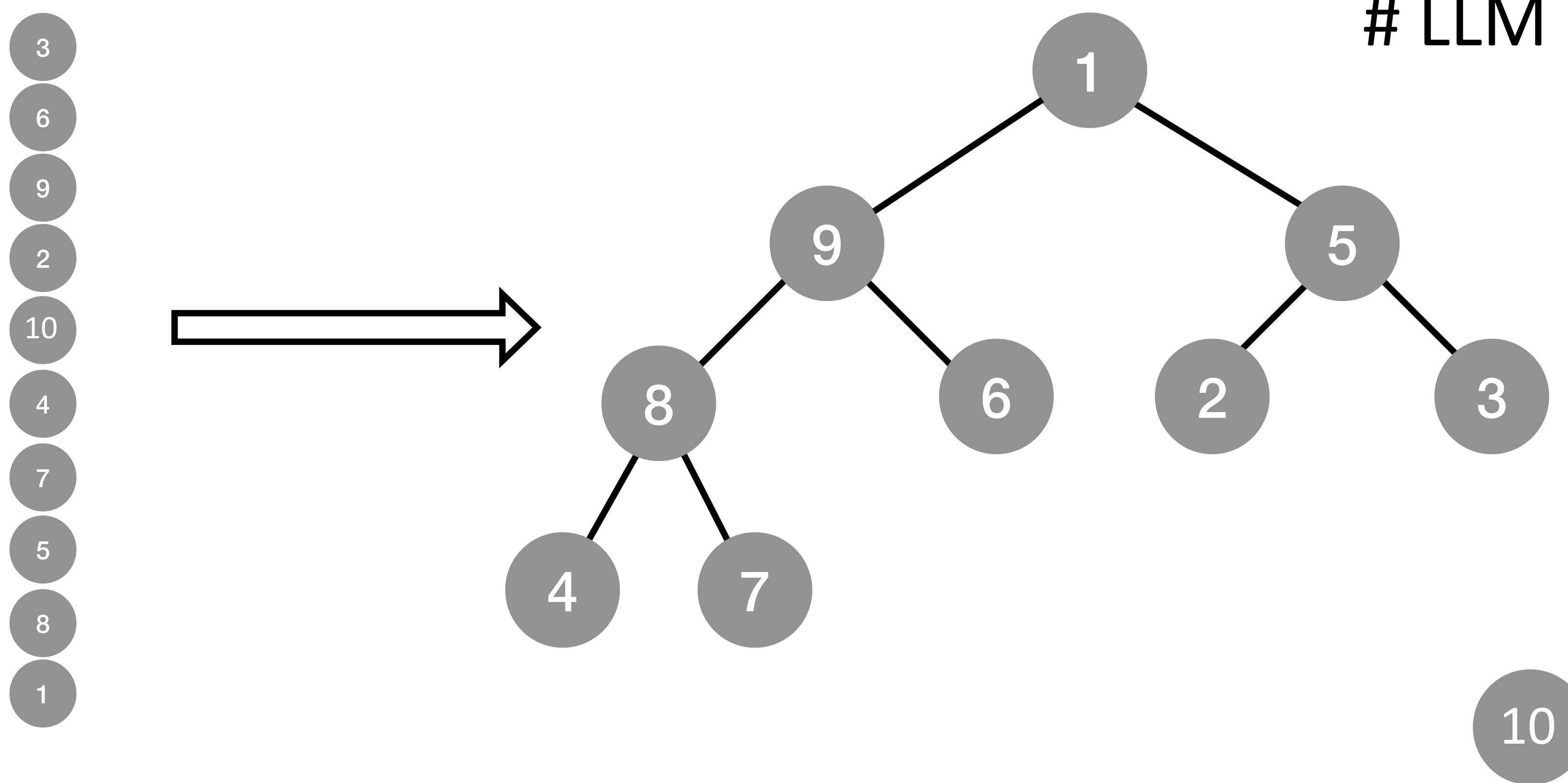
Step 1: Build max heap tree



# Pairwise

- Sorting-based ranking: Heapsort

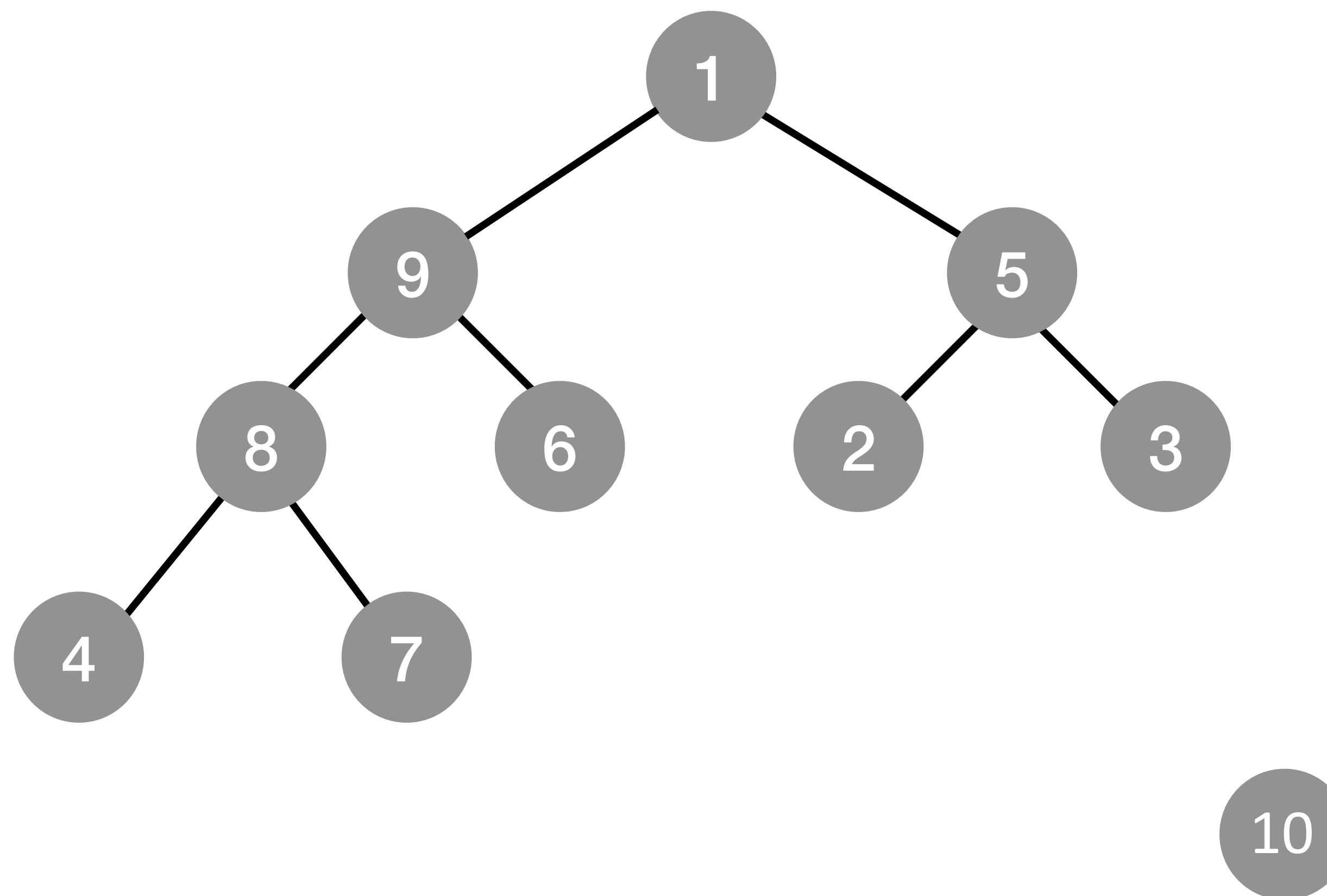
Step 1: Build max heap tree



# Pairwise

- **Sorting-based ranking: Heapsort**

Step 2: heapify the tree



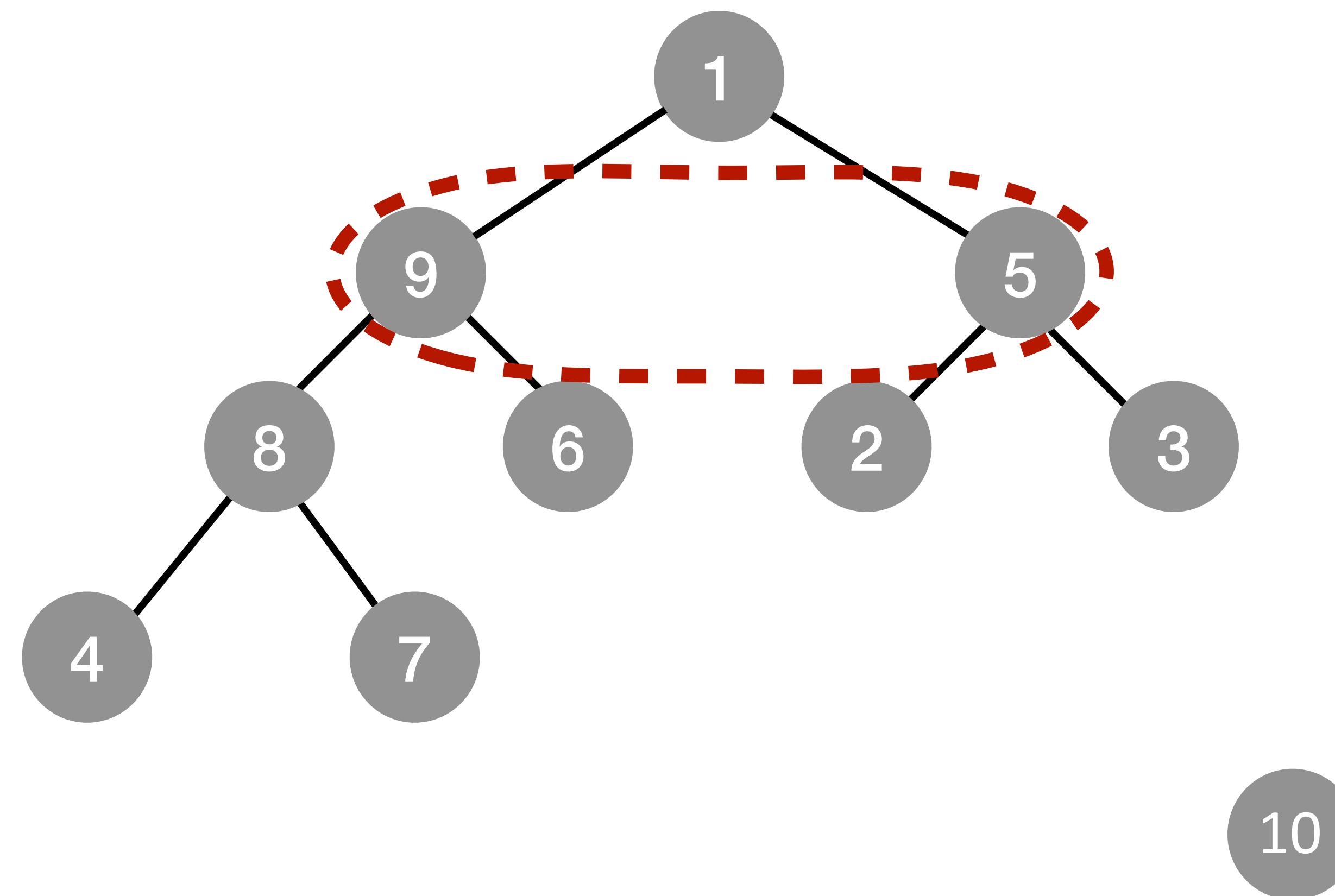
# LLM call: 0

# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree

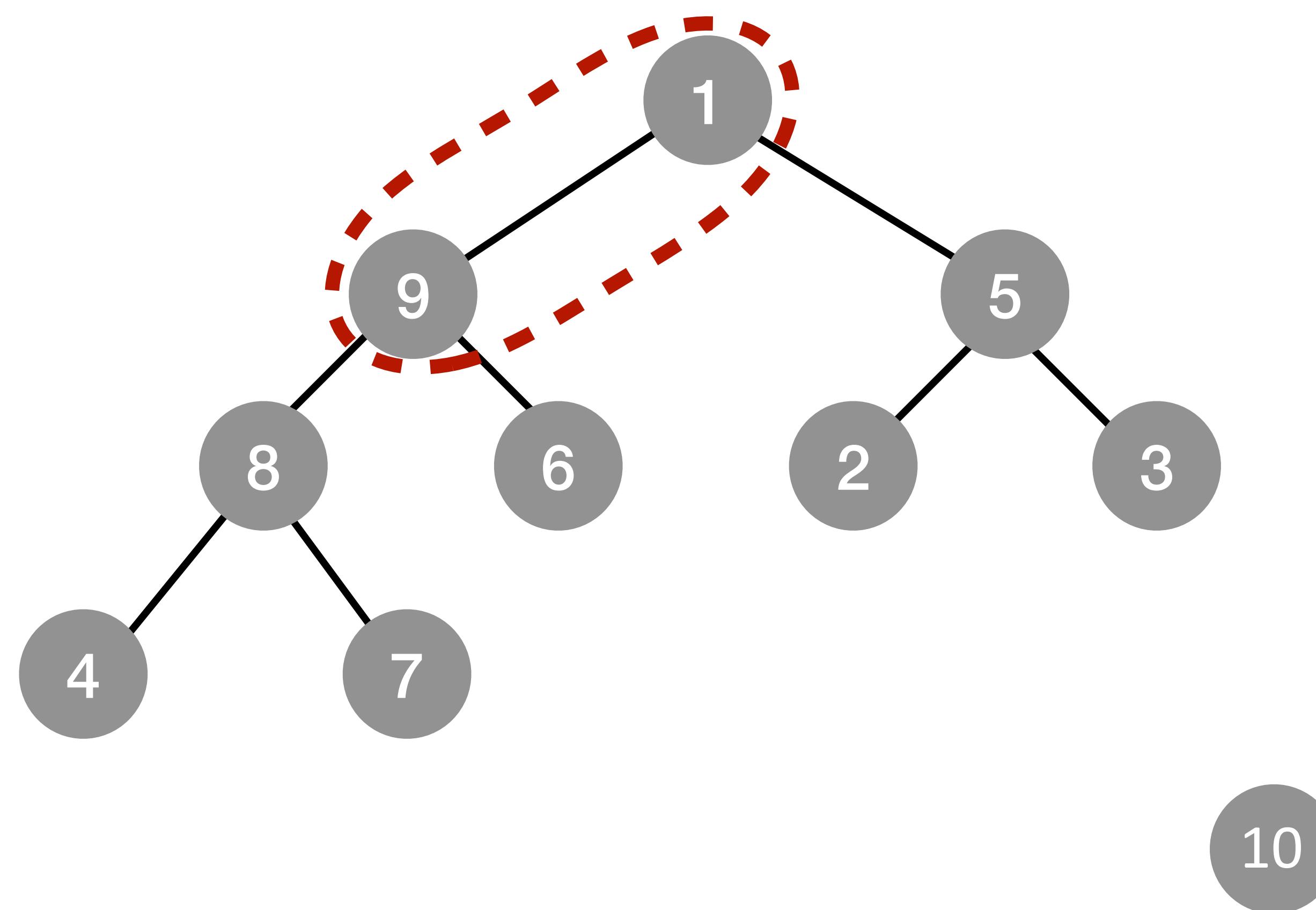
# LLM call: 1



# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree

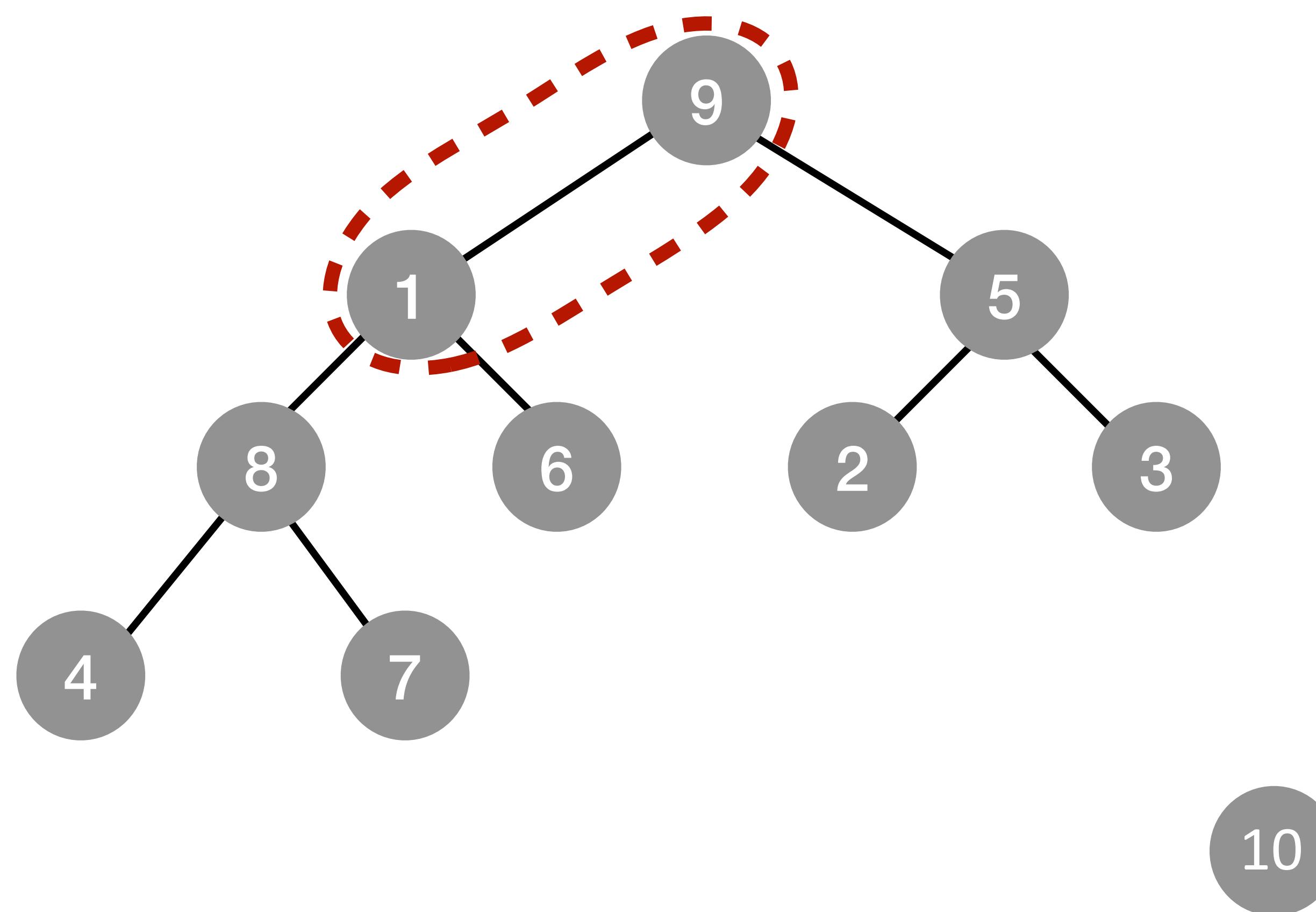


# LLM call: 2

# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree



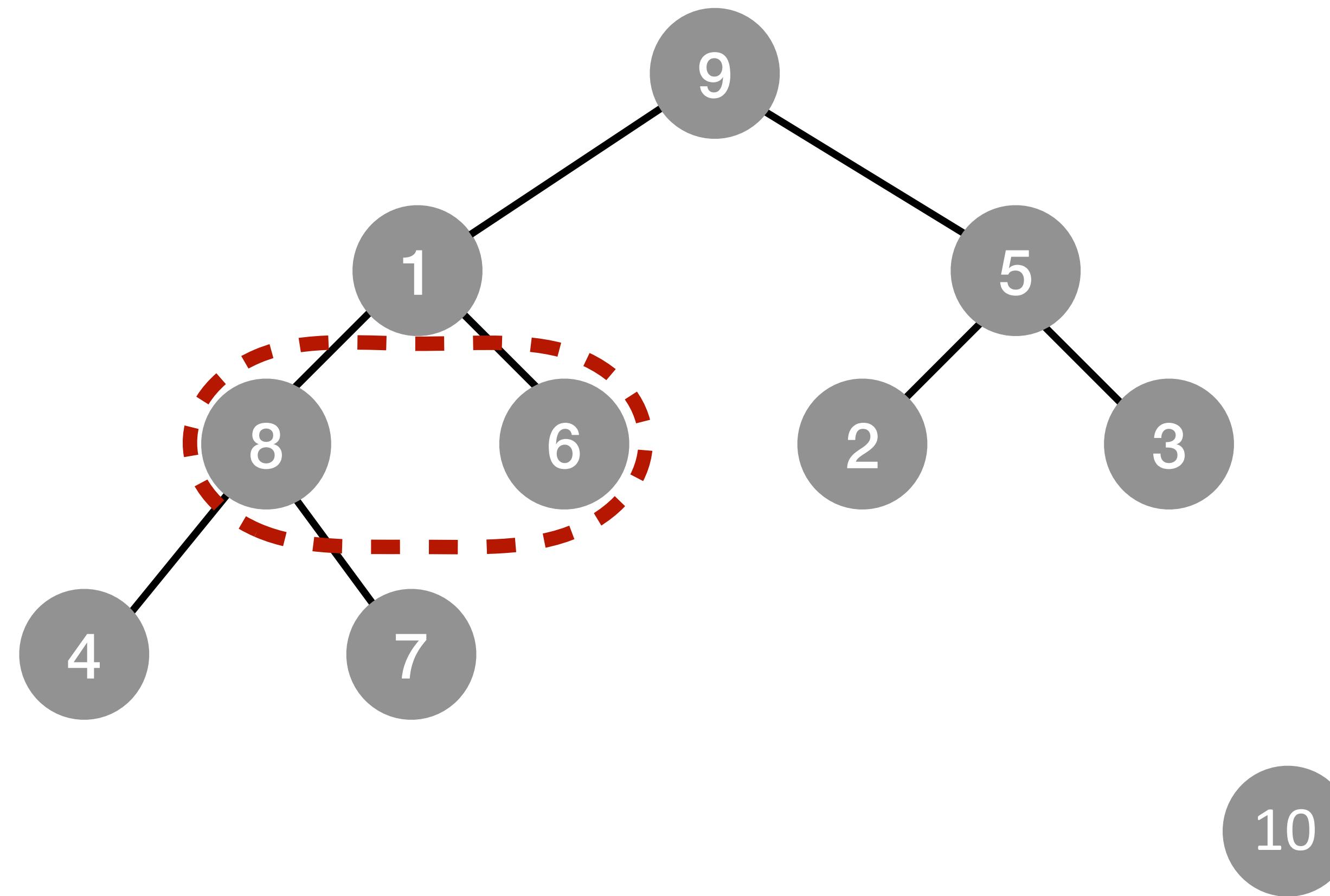
# LLM call: 2

# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree

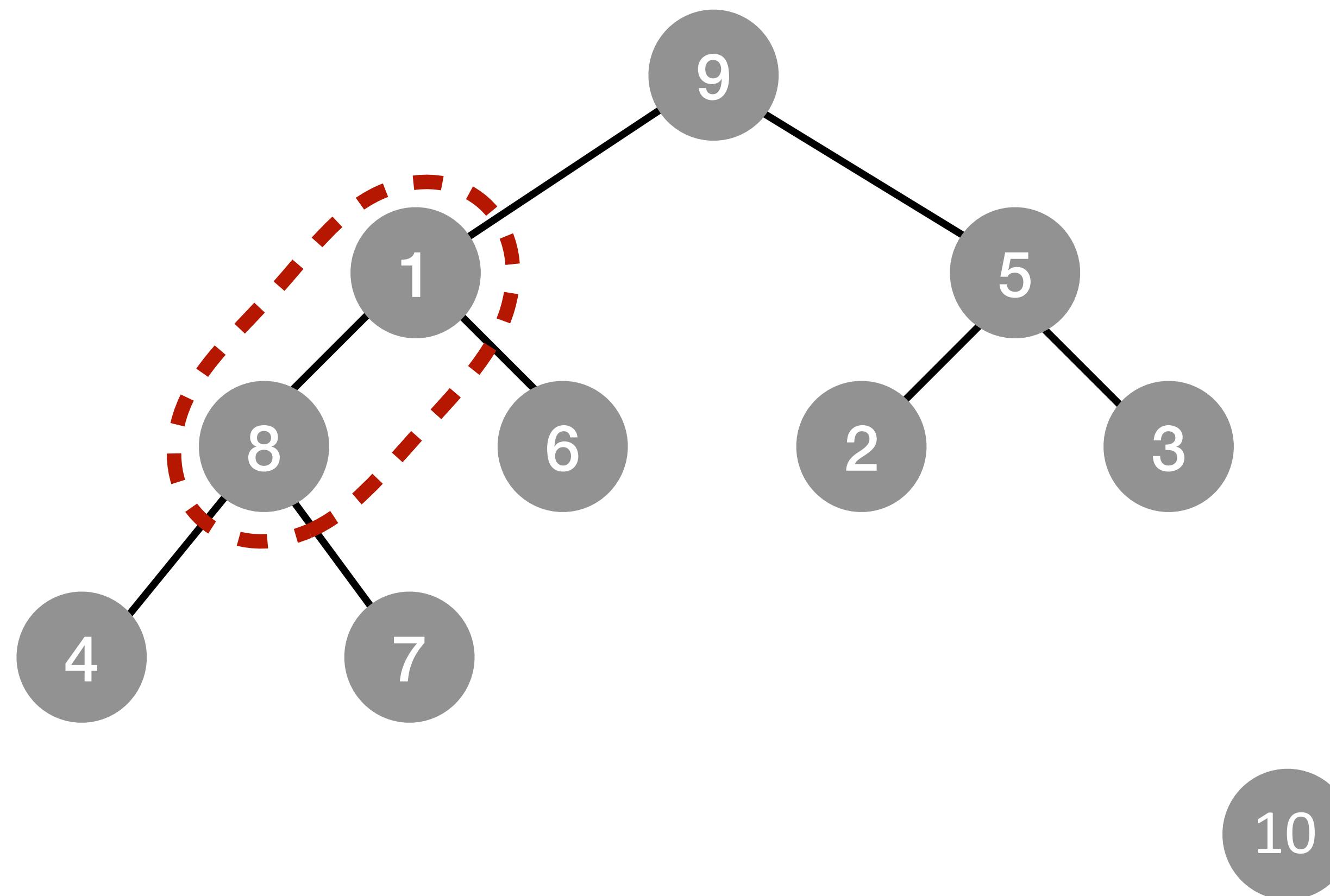
# LLM call: 3



# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree

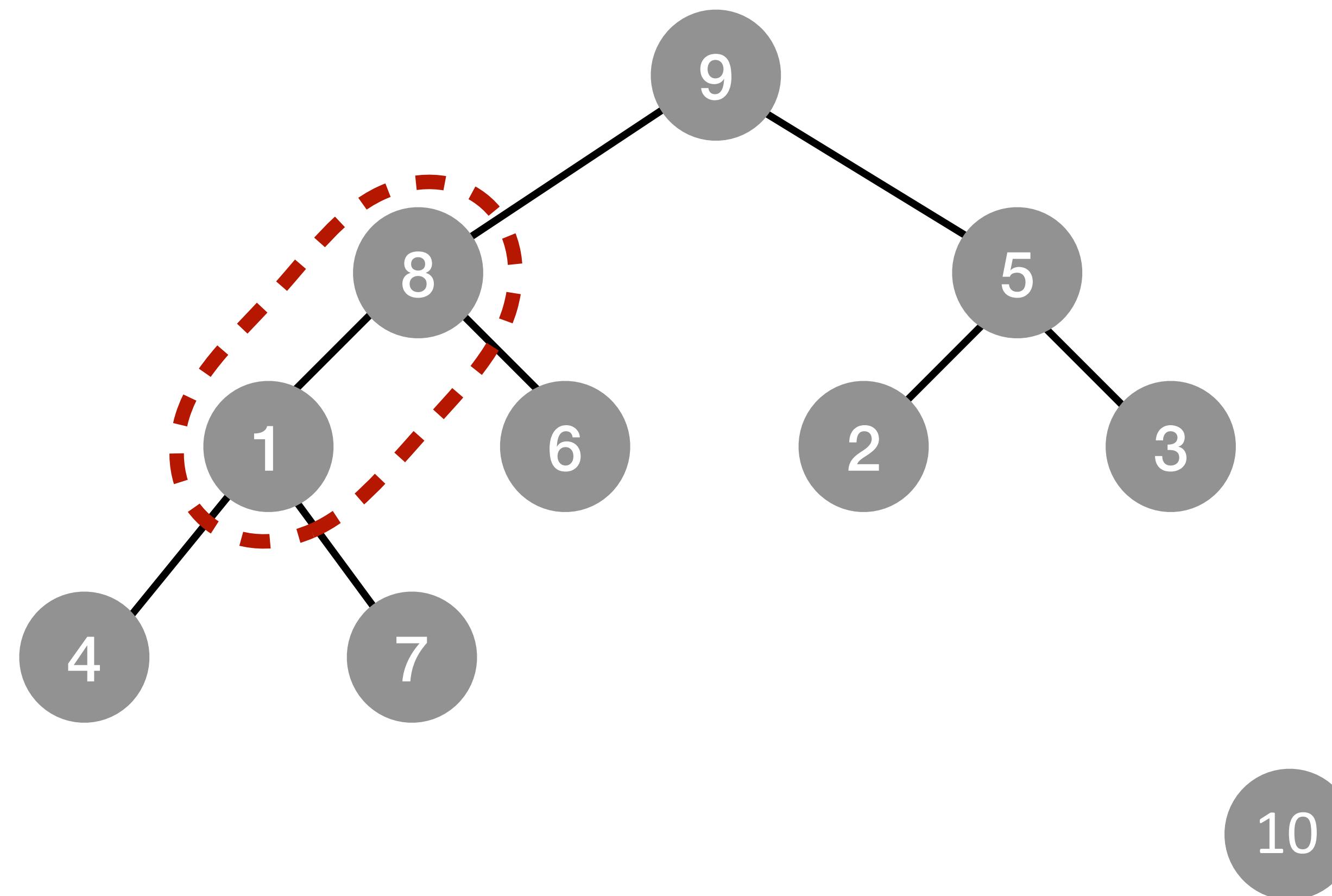


# LLM call: 4

# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree



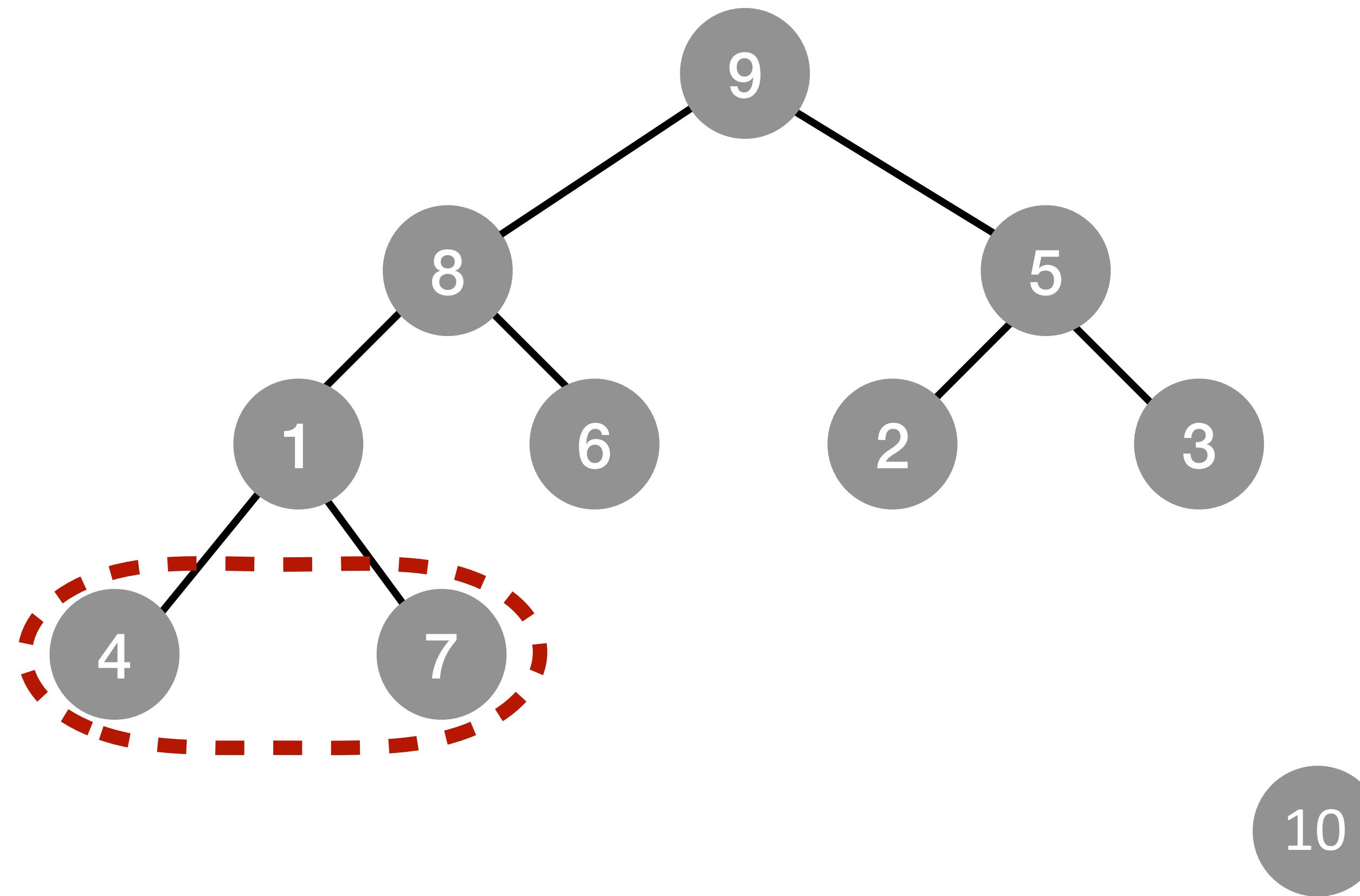
# LLM call: 4

# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree

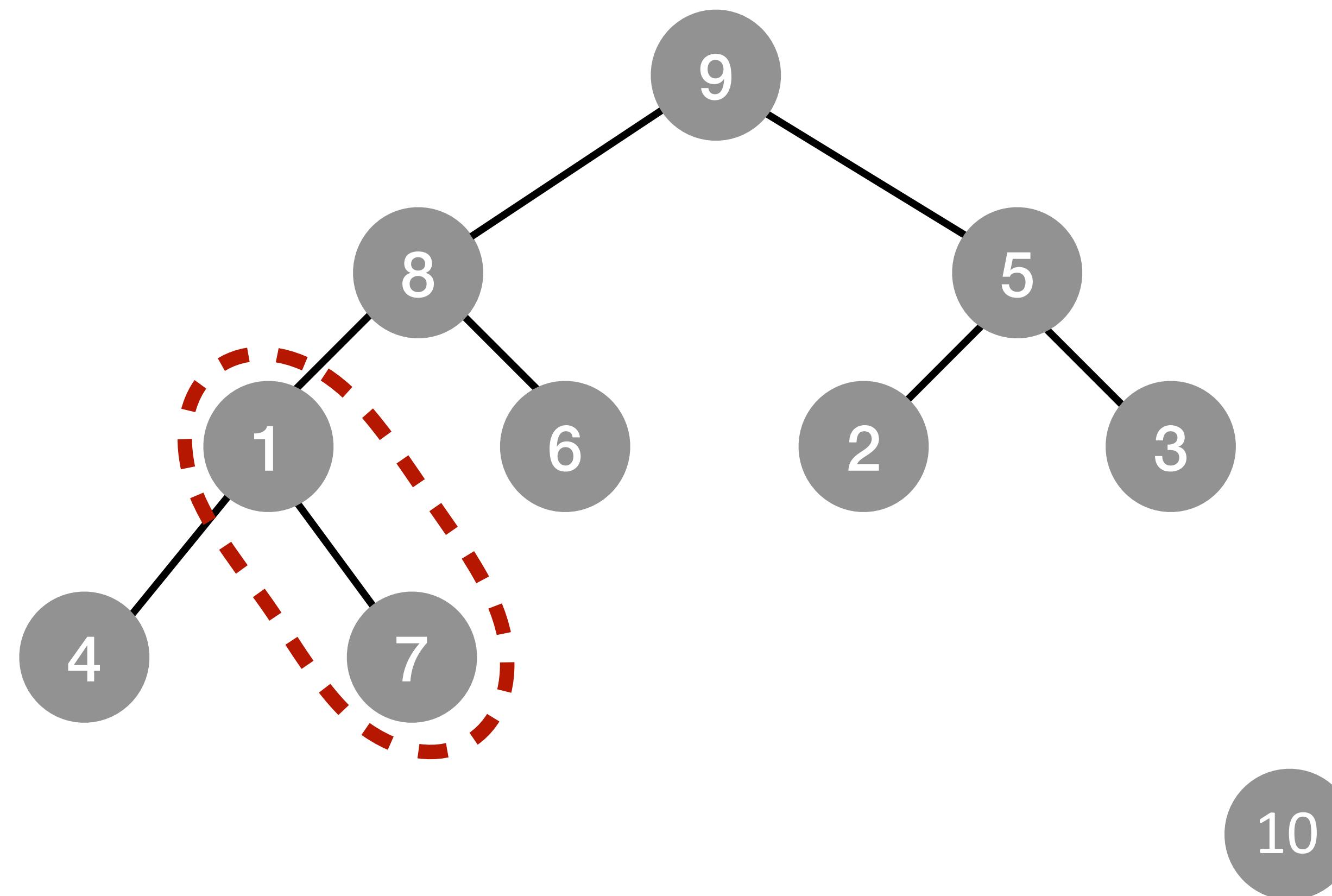
# LLM call: 5



# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree



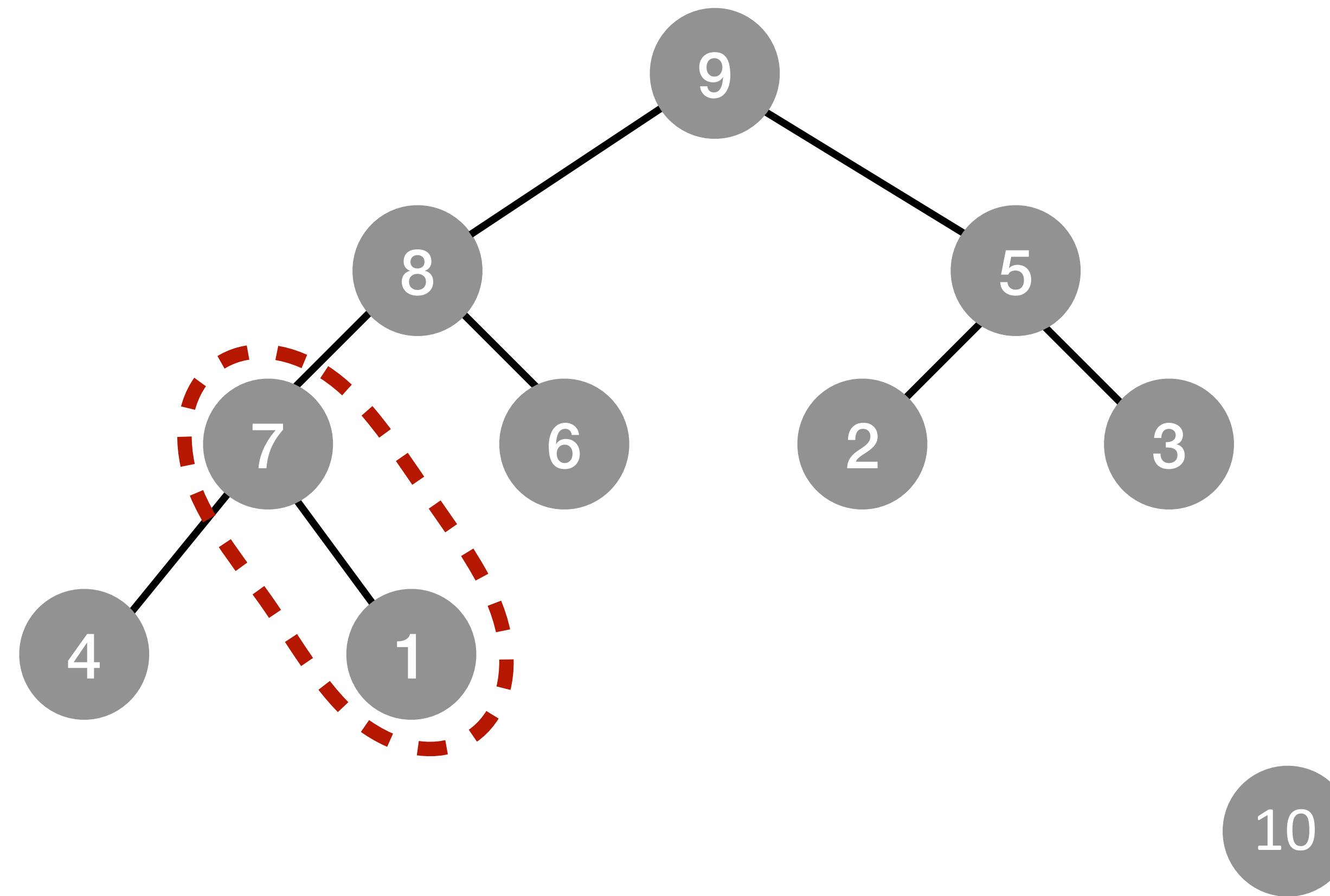
# LLM call: 6

# Pairwise

- Sorting-based ranking: Heapsort

Step 2: heapify the tree

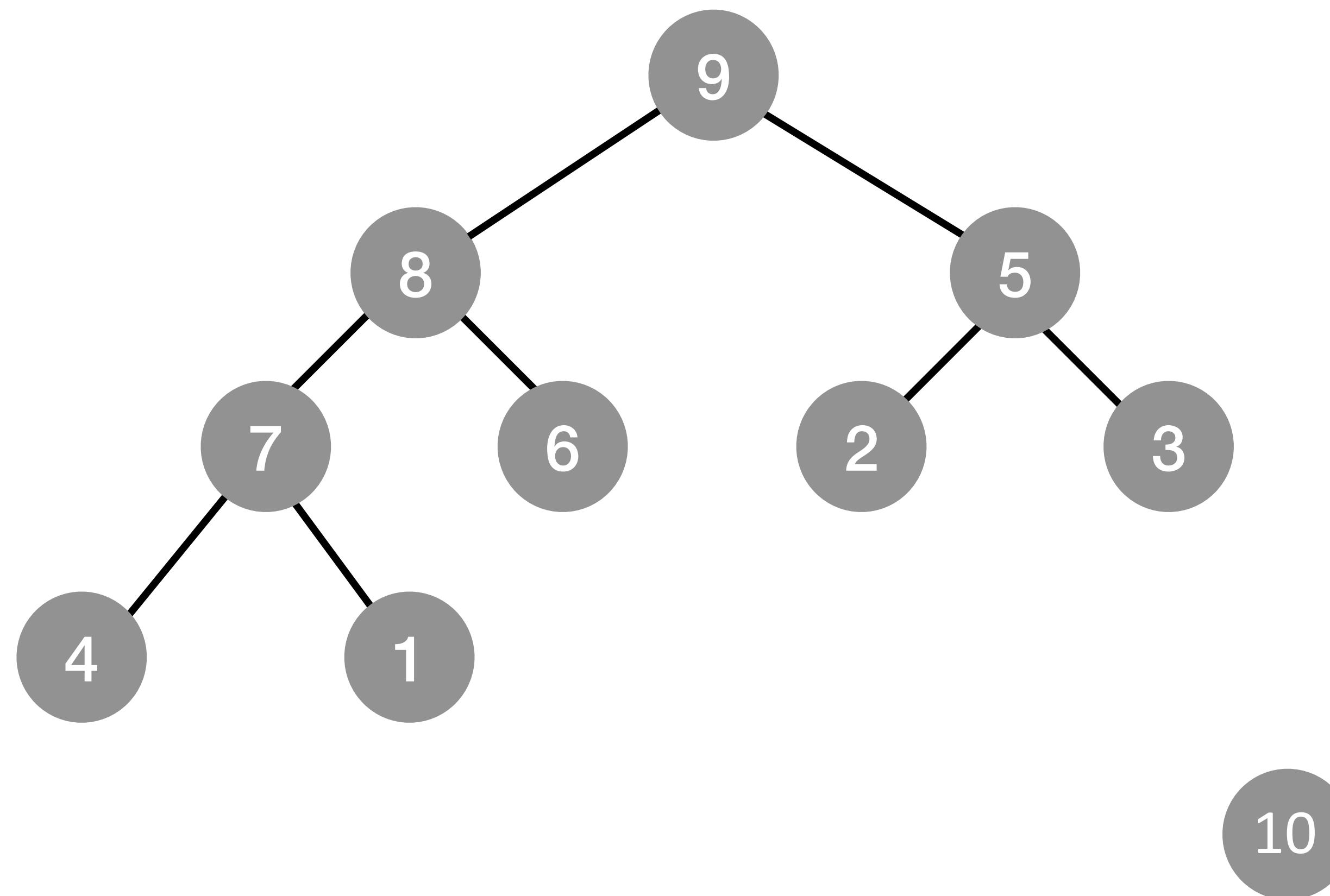
# LLM call: 6



# Pairwise

- **Sorting-based ranking: Heapsort**

Step 2: heapify the tree

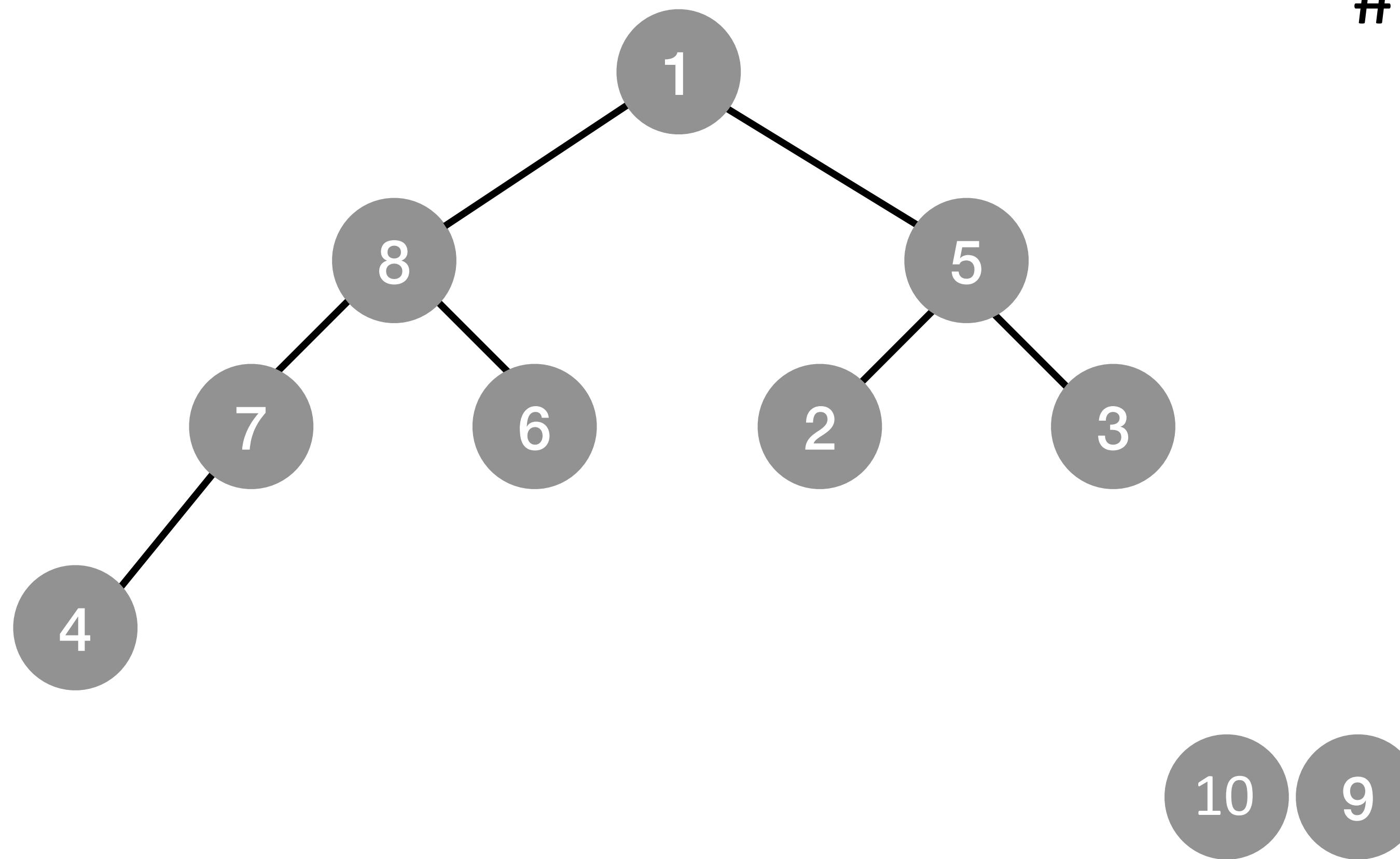


# LLM call: 6

# Pairwise

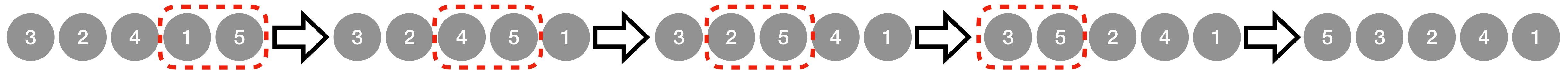
- **Sorting-based ranking: Heapsort**

Step 2: heapify the tree

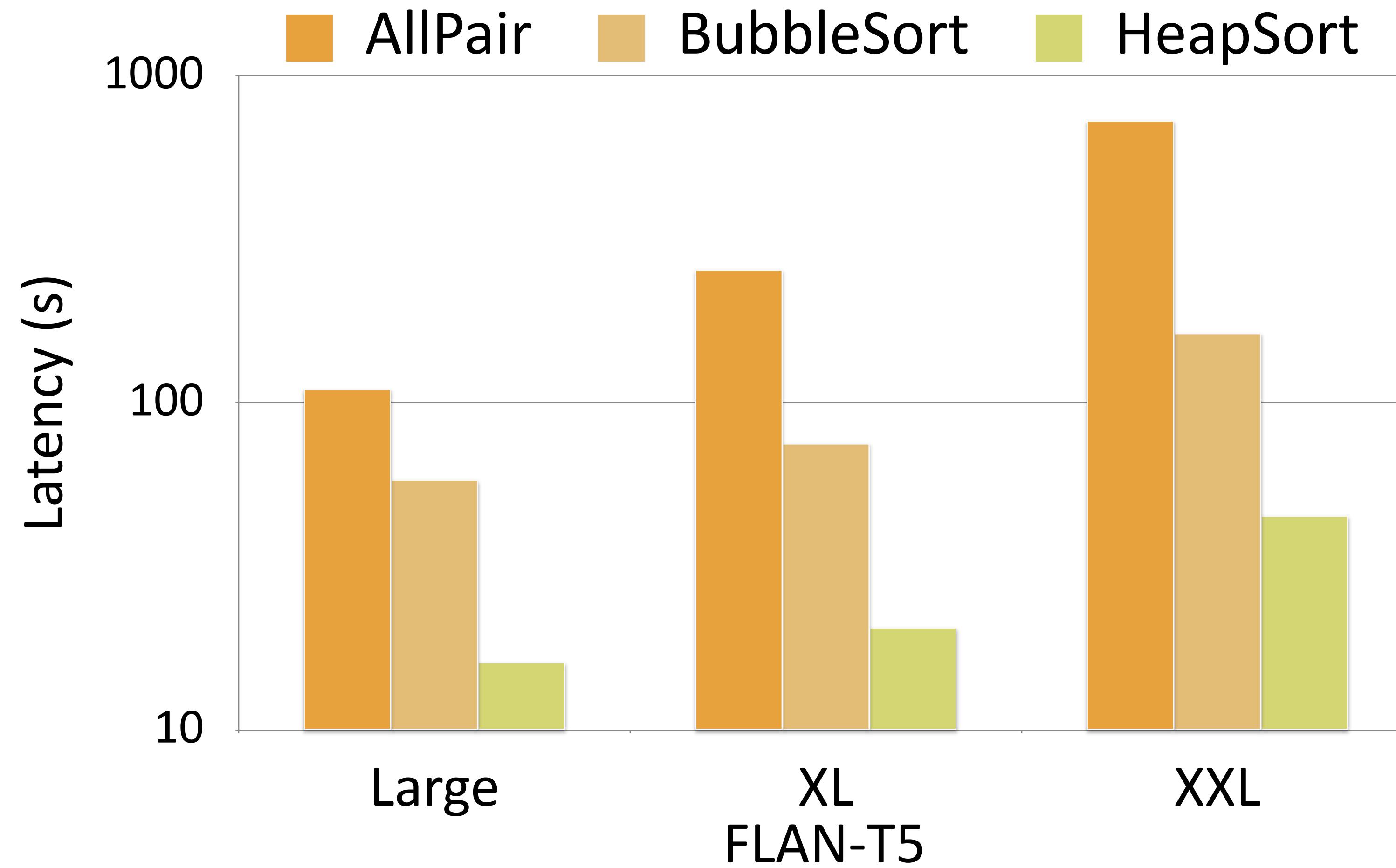


# Pairwise

- Sorting-based ranking: Bubblesort



# Pairwise Efficiency



# LLMs as Rankers: Setwise

Given a query  $\{query\}$ , which of the following passages is more relevant one to the query?

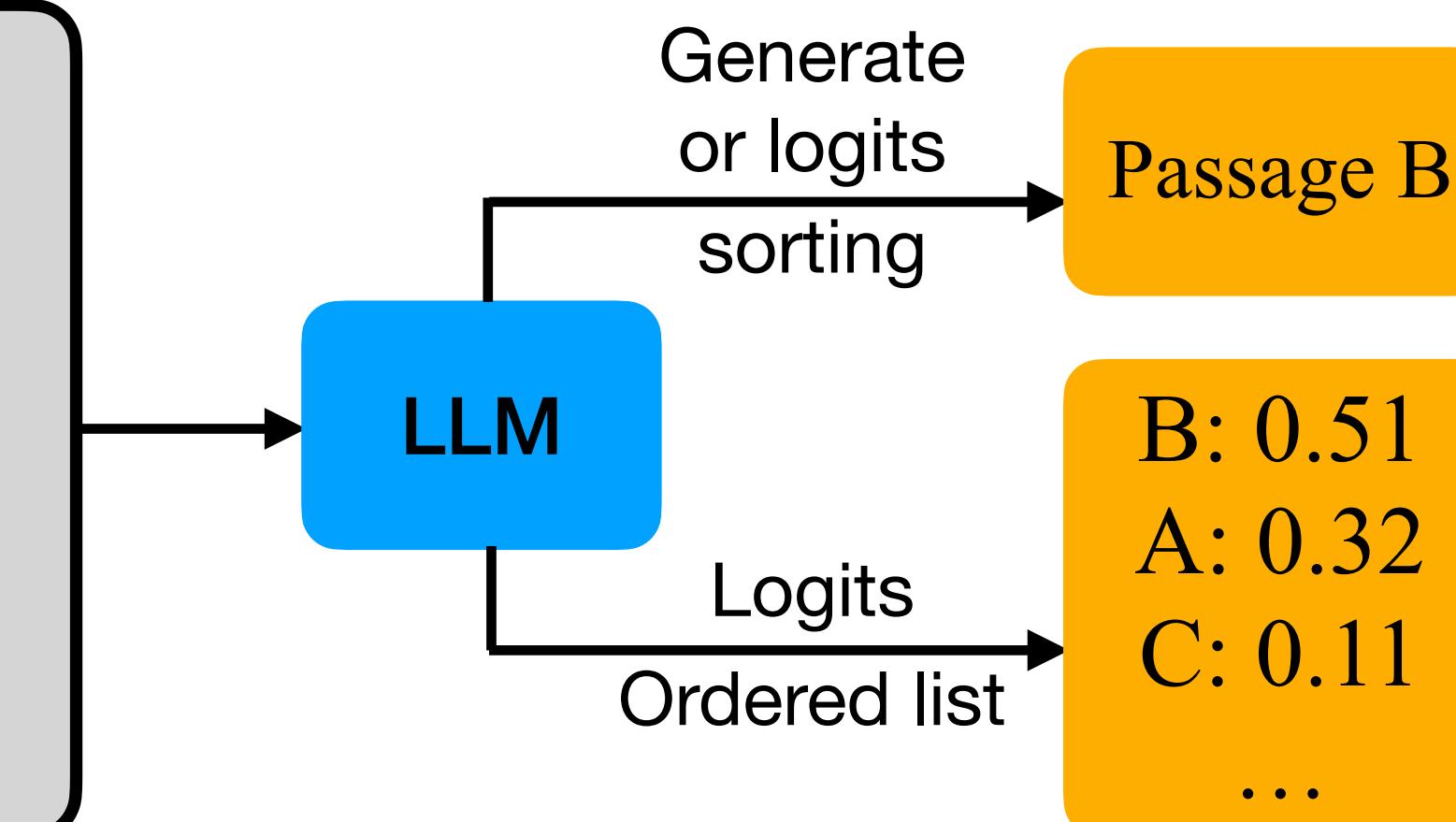
Passage A:  $\{passage\_1\}$

Passage B:  $\{passage\_2\}$

Passage C:  $\{passage\_3\}$

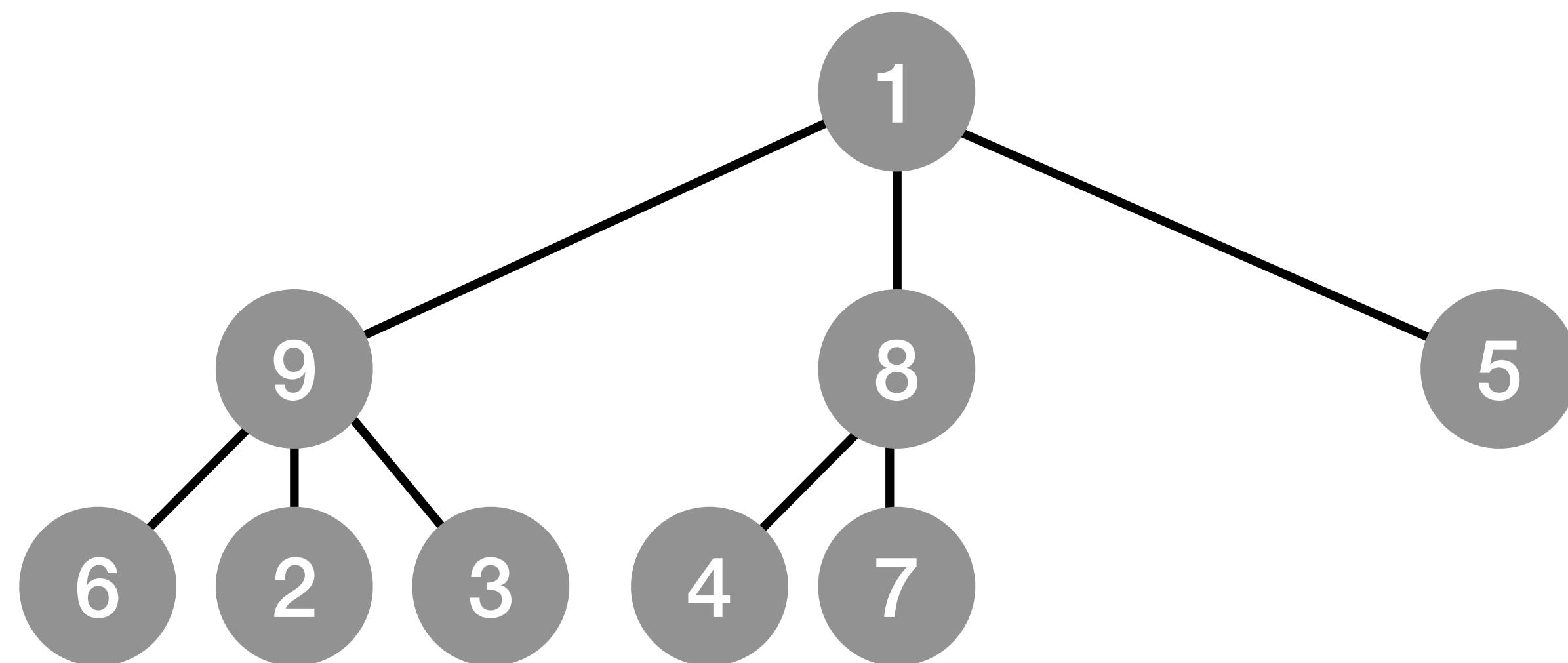
...

Output only the passage label of the most relevant passage:



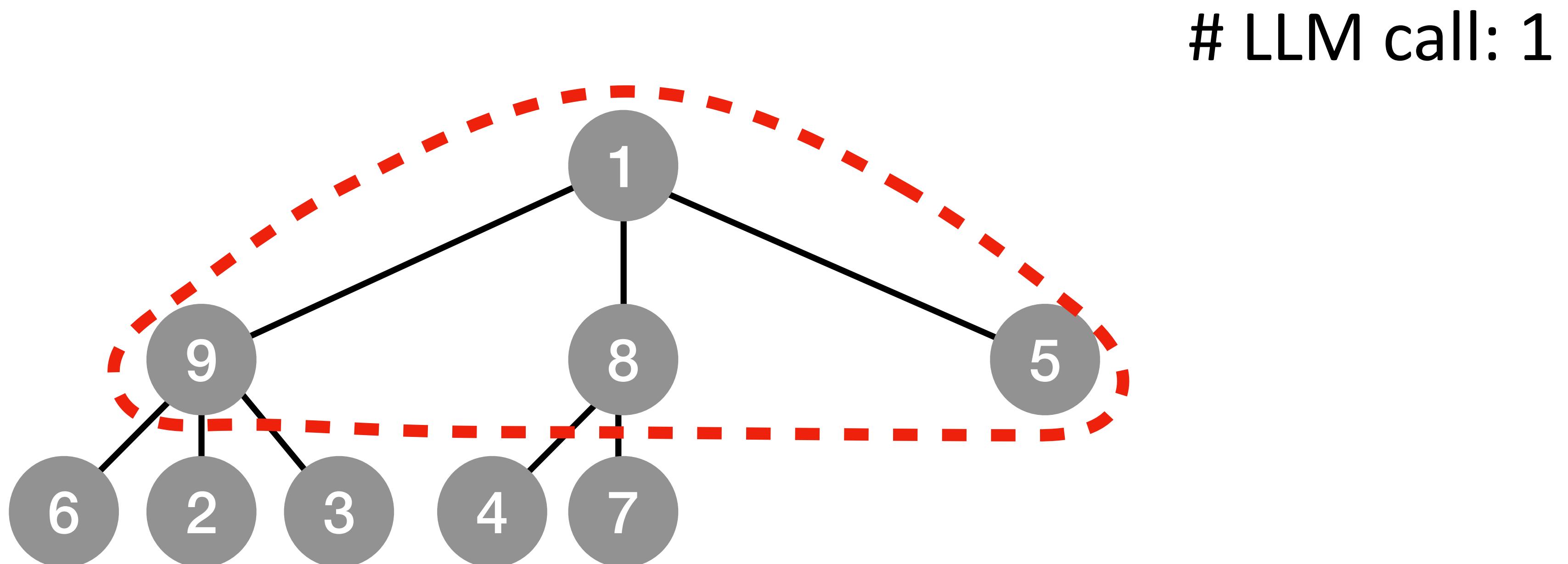
# Setwise

- Sorting-based ranking: Heapsort



# Setwise

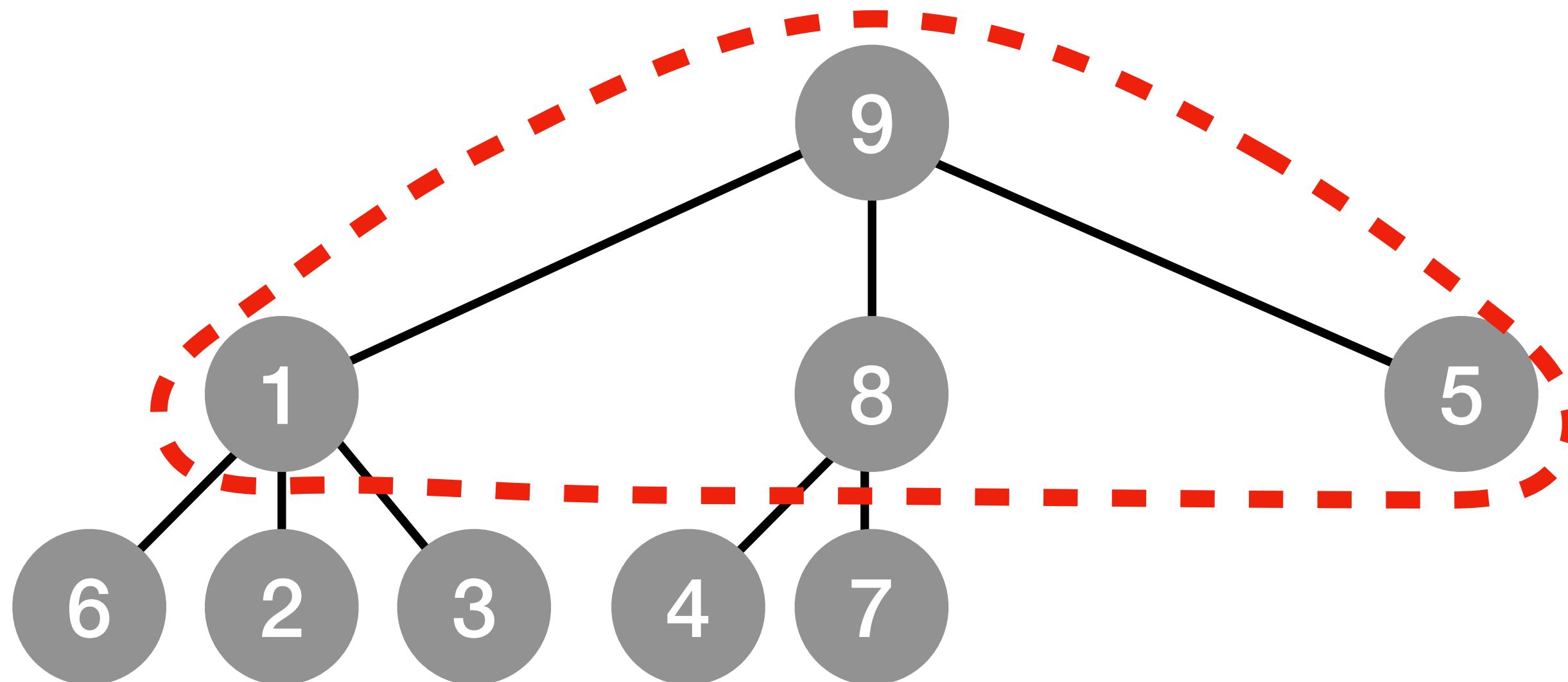
- Sorting-based ranking: Heapsort



# Setwise

- Sorting-based ranking: Heapsort

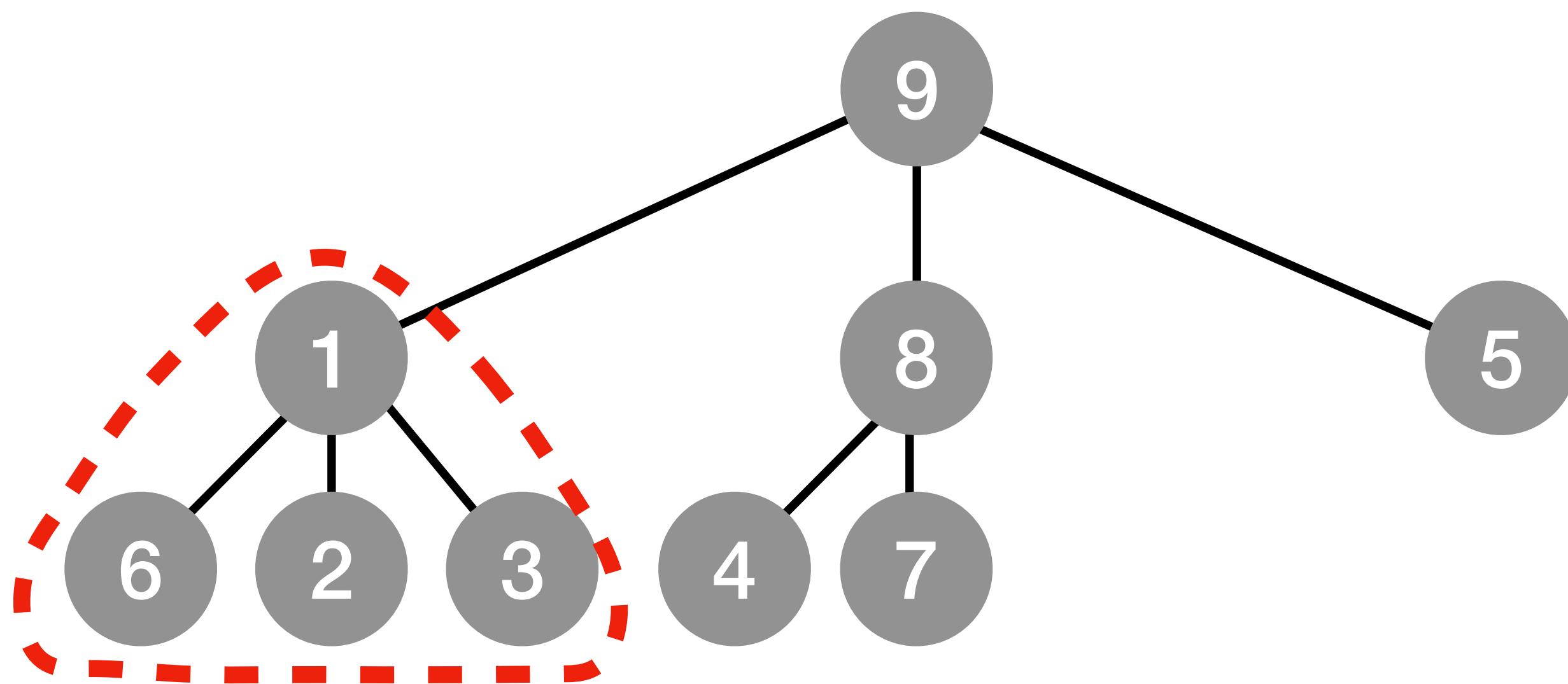
# LLM call: 1



# Setwise

- Sorting-based ranking: Heapsort

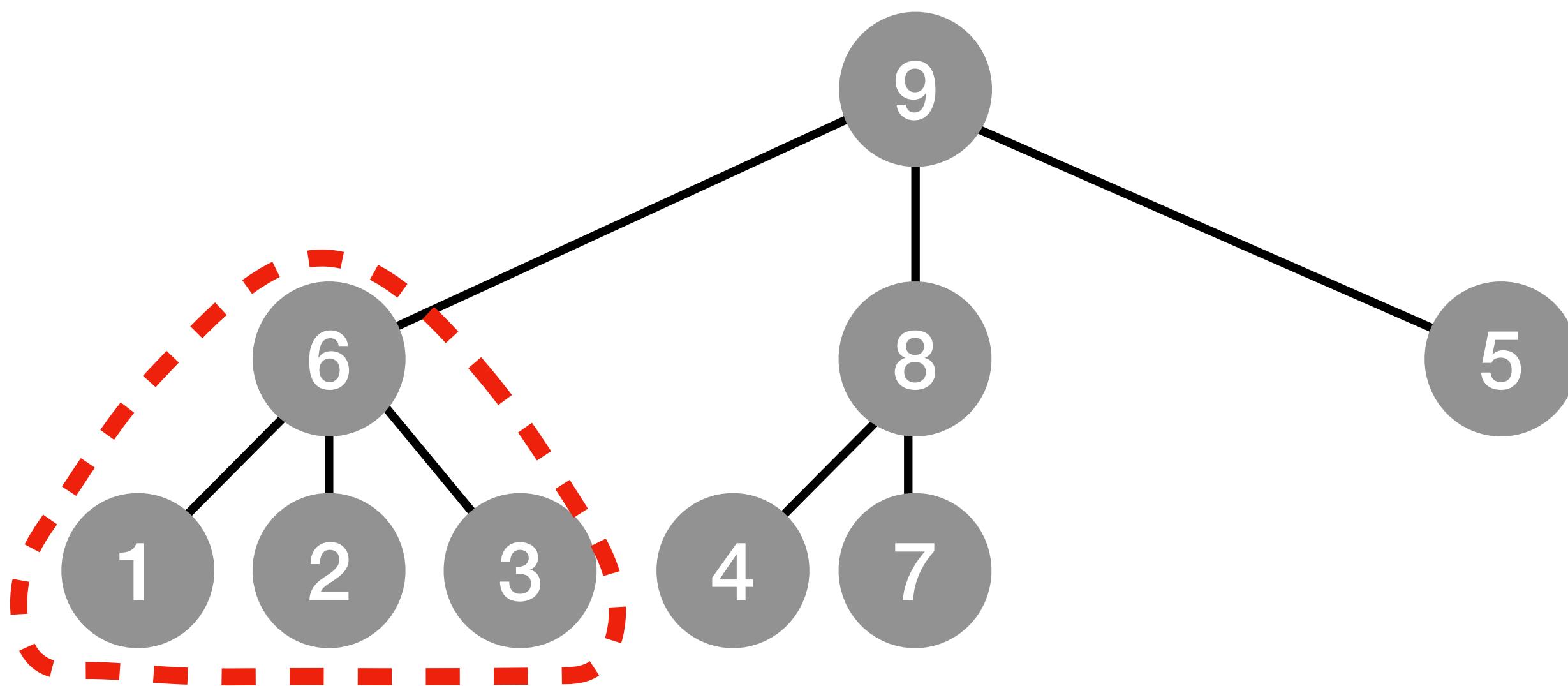
# LLM call: 2



# Setwise

- Sorting-based ranking: Heapsort

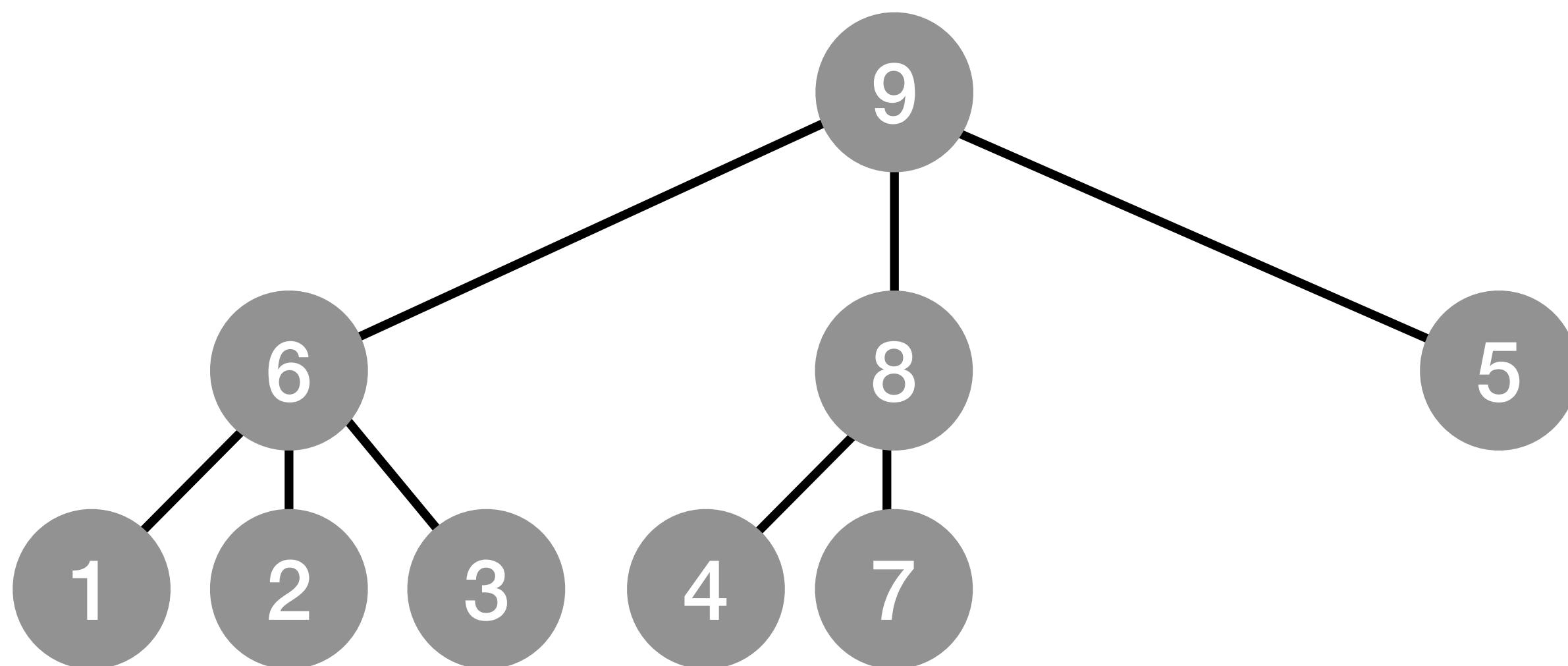
# LLM call: 2



# Setwise

- Sorting-based ranking: Heapsort

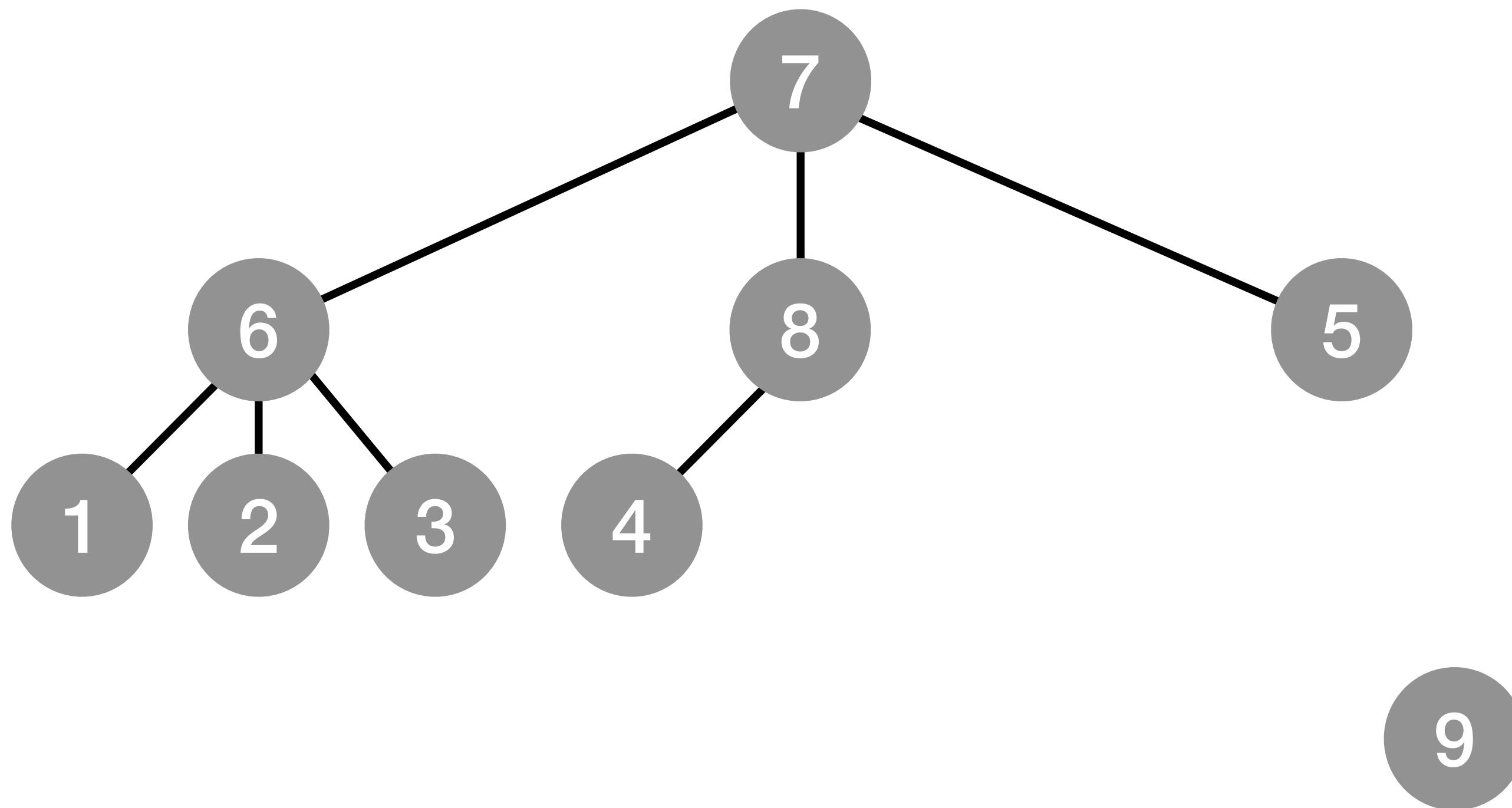
# LLM call: 2



# Setwise

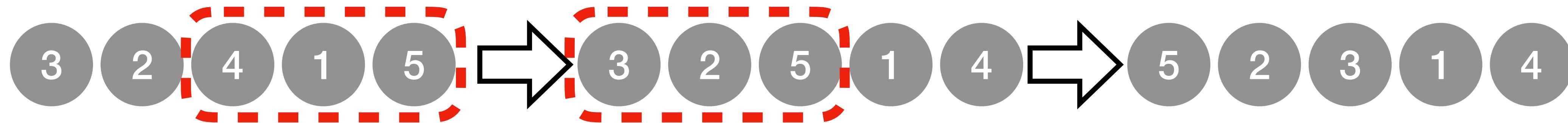
- Sorting-based ranking: Heapsort

# LLM call: 2

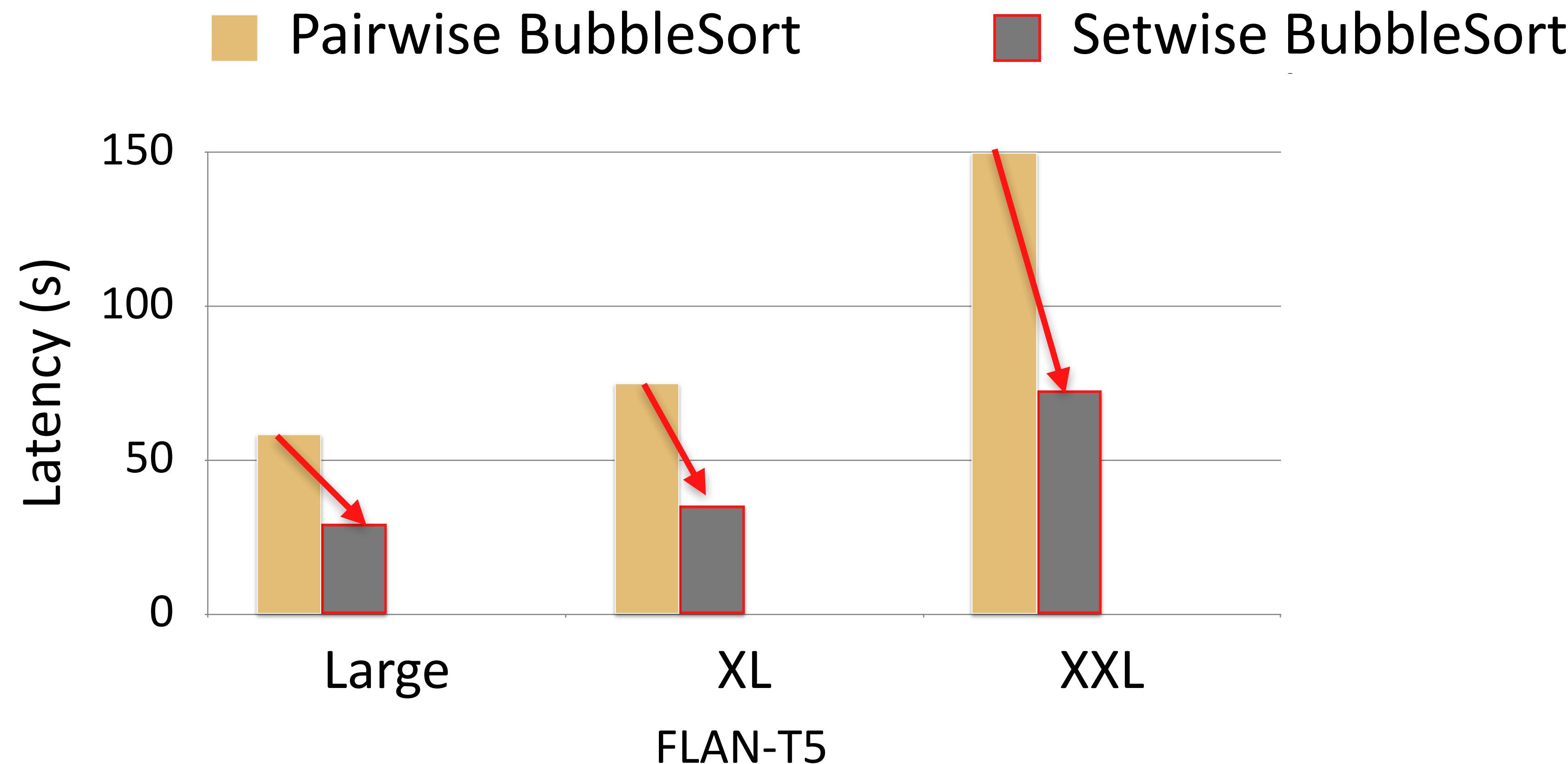


# Setwise

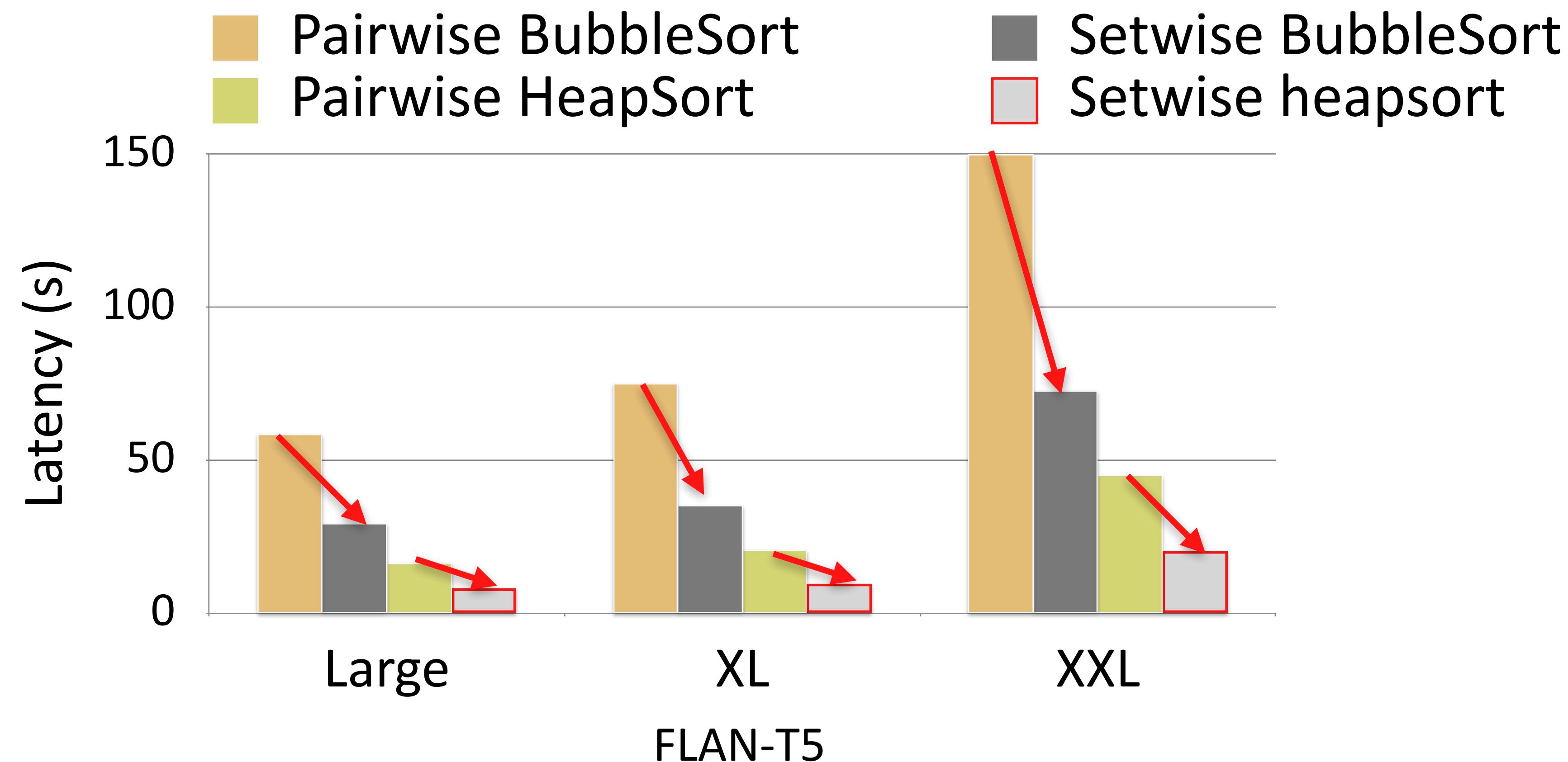
- Sorting-based ranking: Bubblesort



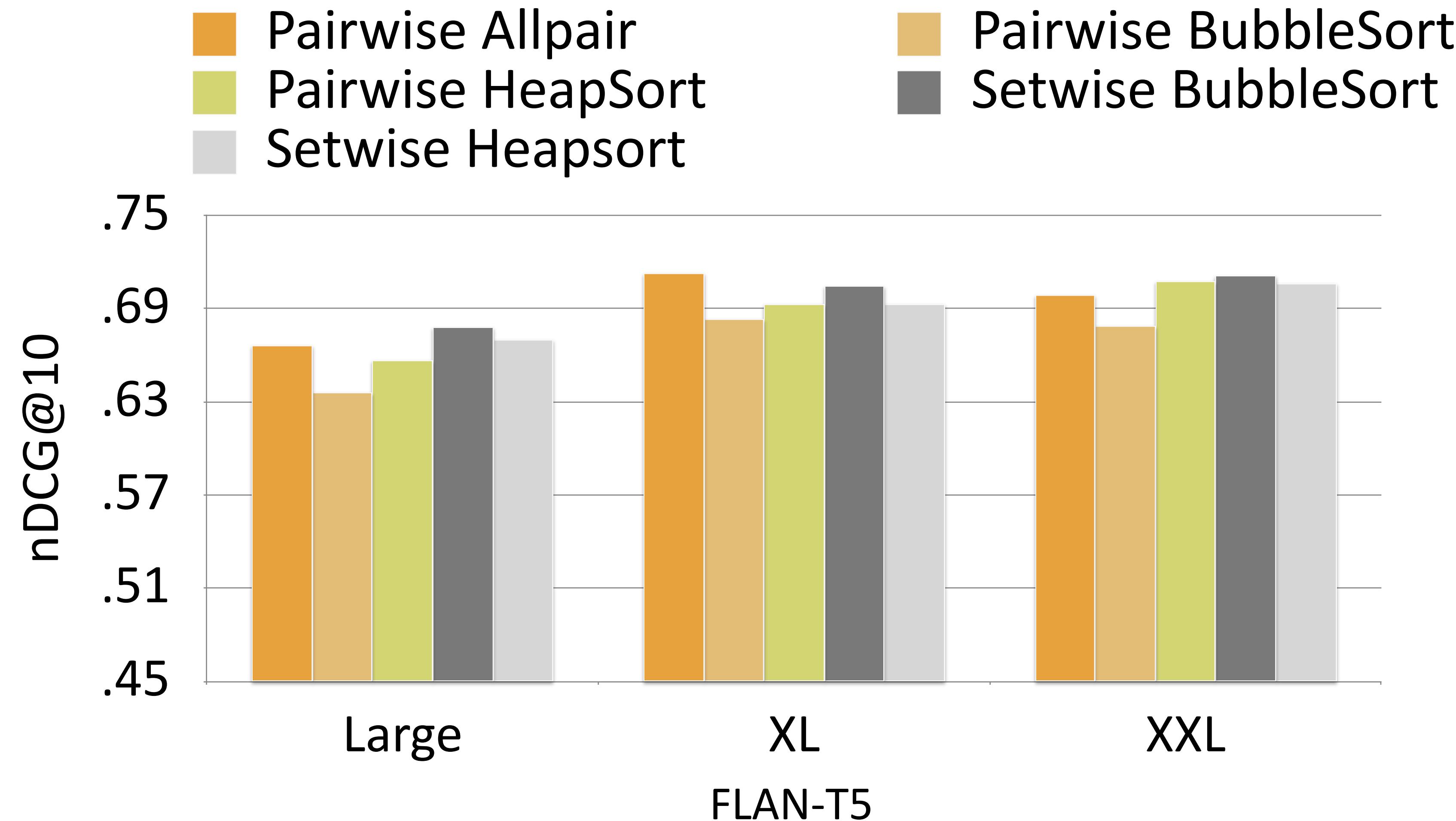
# Setwise Efficiency



# Pairwise Efficiency

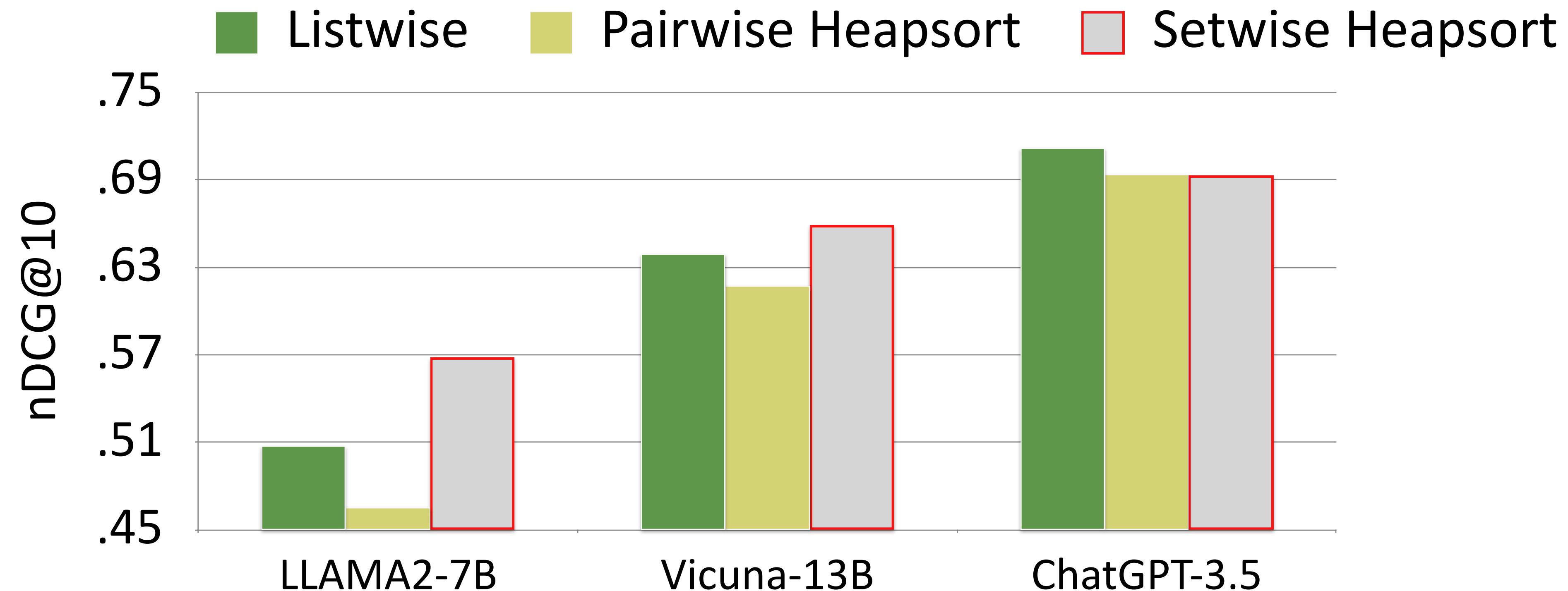


# Setwise Effectiveness



# Setwise Effectiveness

On other backbones



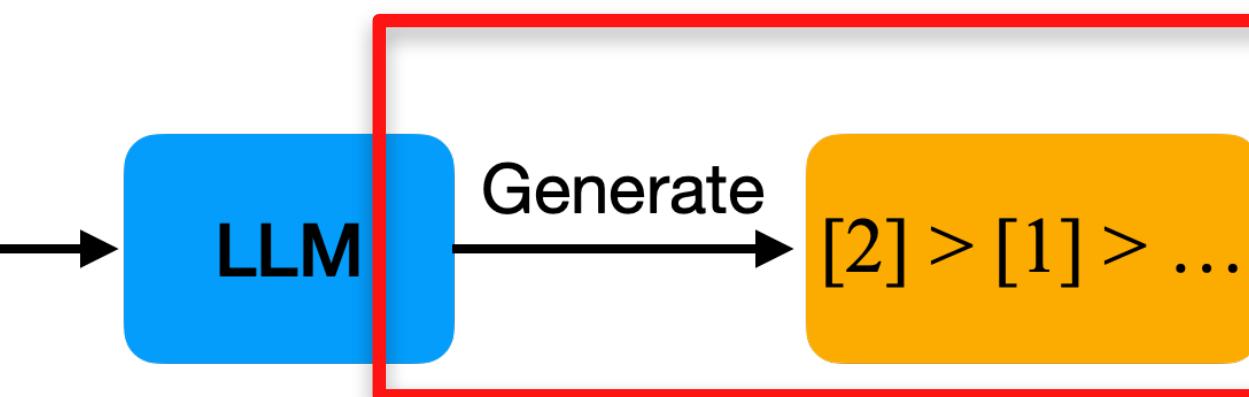
# Listwise Ranking with Setwise Prompting

The following are  $\{num\}$  passages, each indicated by number identifier []. I can rank them based on their relevance to query:  $\{query\}$

[1]  $\{passage\_1\}$   
[2]  $\{passage\_2\}$

...

The ranking results of the  $\{num\}$  passages (only identifiers) is:



Given a query  $\{query\}$ , which of the following passages is more relevant one to the query?

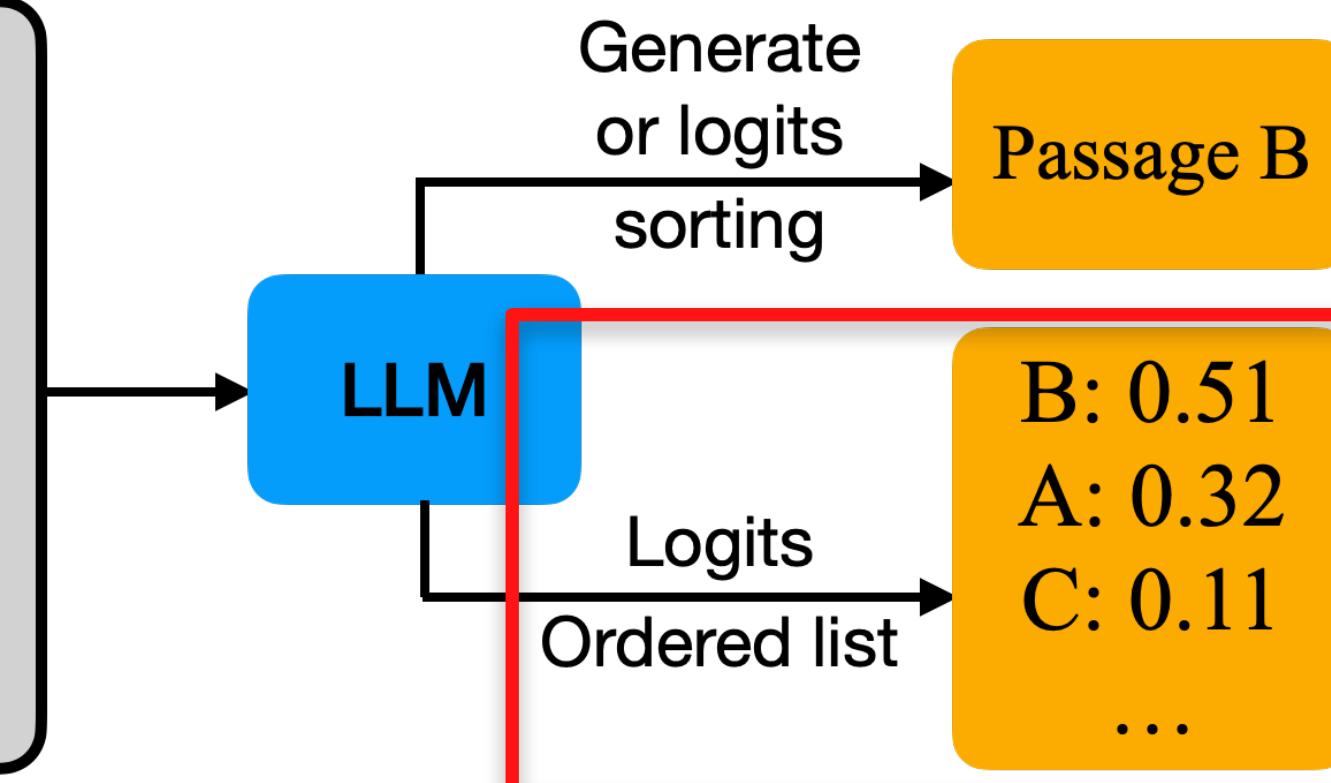
Passage A:  $\{passage\_1\}$

Passage B:  $\{passage\_2\}$

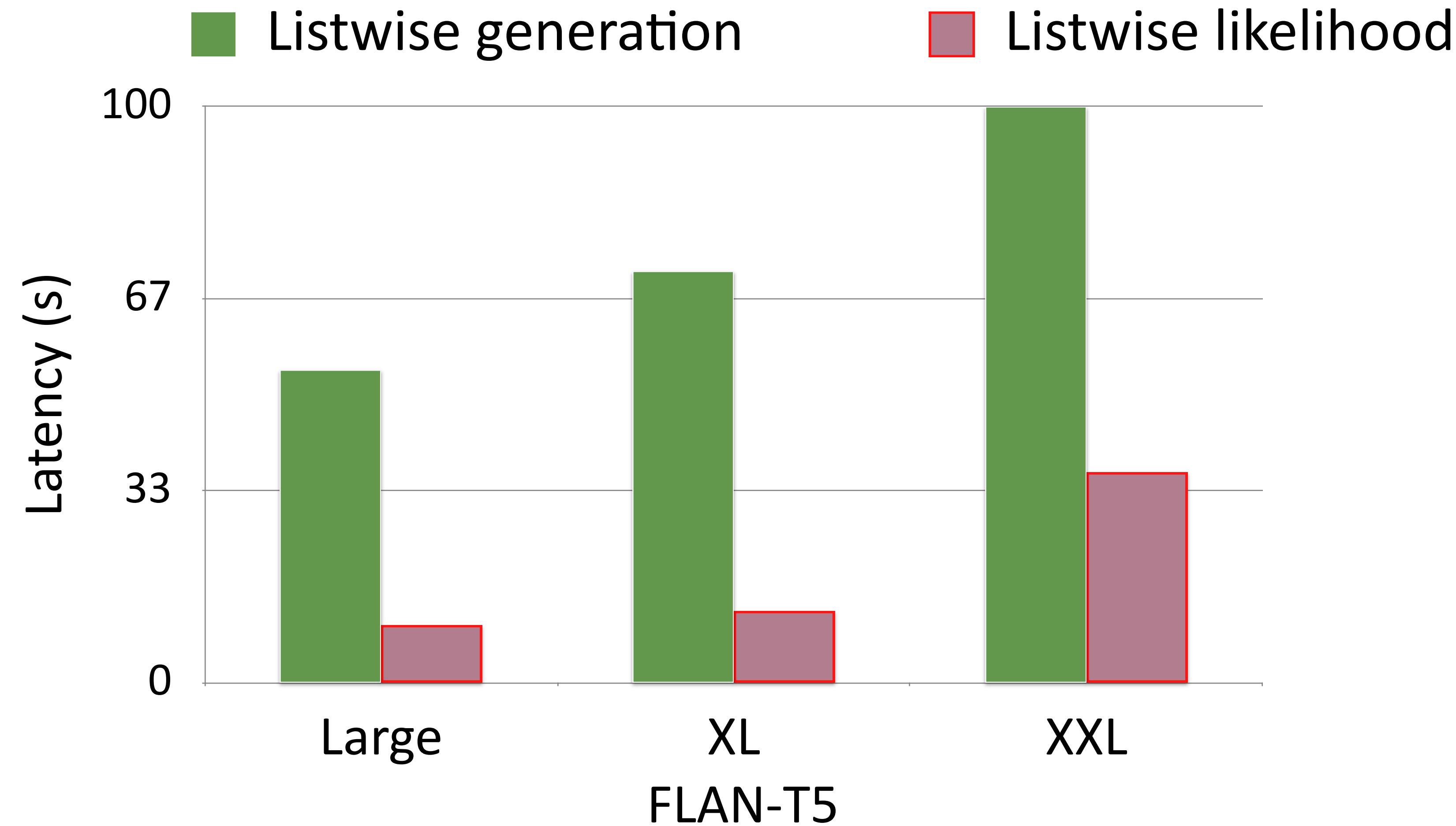
Passage C:  $\{passage\_3\}$

...

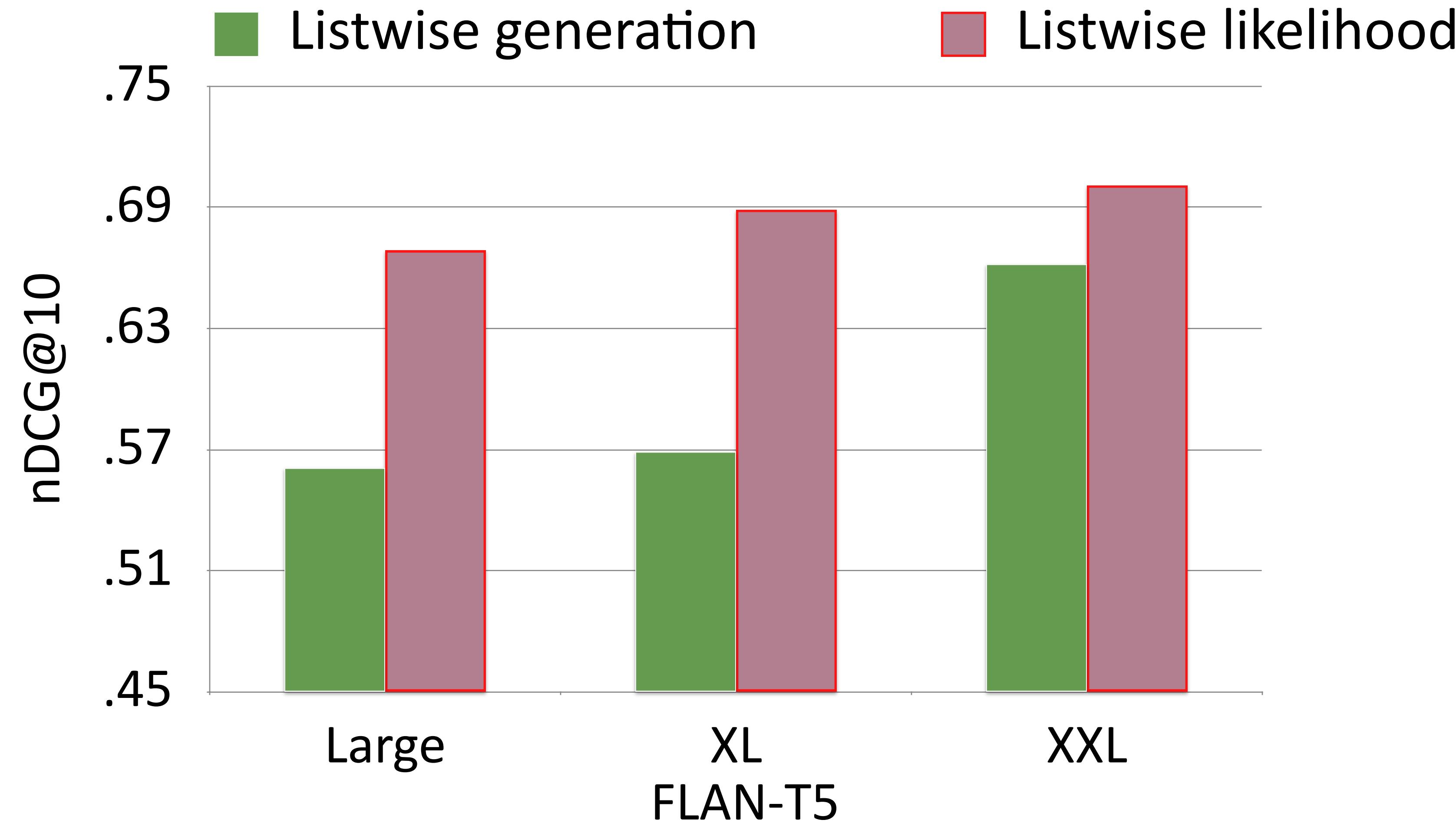
Output only the passage label of the most relevant passage:



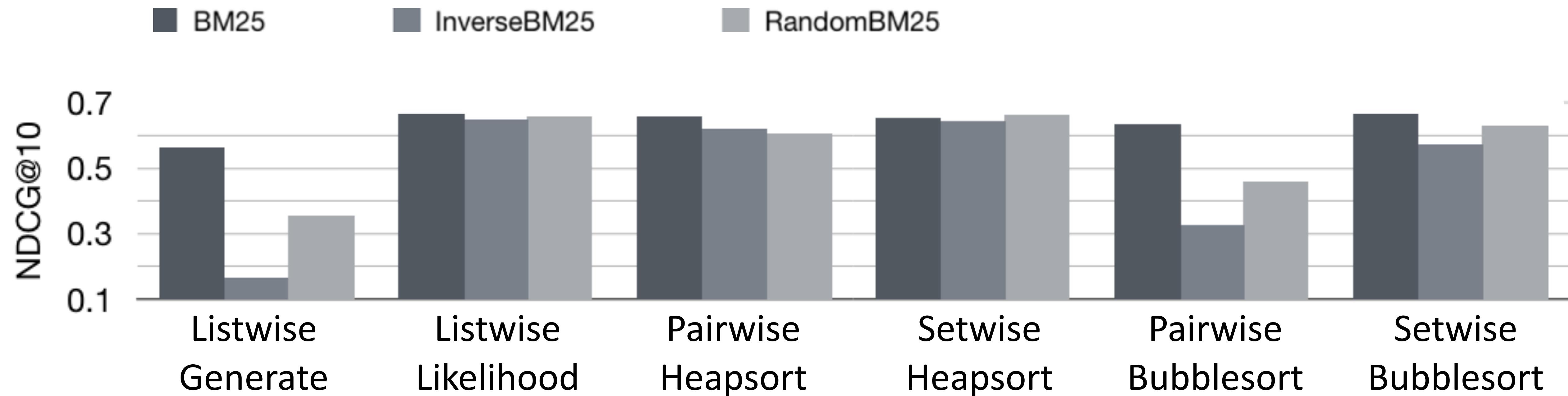
# Efficiency



# Effectiveness



# Sensitivity to the Initial Ranking



# Properties

Methods	Logits	Generate
<i>pointwise qlm</i>	✓	
<i>pointwise yes_no</i>	✓	
<i>listwise generation</i>		✓
<i>listwise likelihood</i>	✓	
<i>pairwise allpair</i>	✓	✓
<i>pairwise heapsort</i>	✓	✓
<i>pairwise bubblesort</i>	✓	✓
<i>setwise heapsort</i>	✓	✓
<i>setwise bubblesort</i>	✓	✓

# Properties

Methods	Logits	Generate	Batching
<i>pointwise qlm</i>	✓		✓
<i>pointwise yes_no</i>	✓		✓
<i>listwise generation</i>		✓	
<i>listwise likelihood</i>	✓		
<i>pairwise allpair</i>	✓	✓	✓
<i>pairwise heapsort</i>	✓	✓	
<i>pairwise bubblesort</i>	✓	✓	
<i>setwise heapsort</i>	✓	✓	
<i>setwise bubblesort</i>	✓	✓	

# Properties

Methods	Logits	Generate	Batching	Top- <i>k</i>
<i>pointwise qlm</i>	✓		✓	
<i>pointwise yes_no</i>	✓		✓	
<i>listwise generation</i>		✓		✓
<i>listwise likelihood</i>	✓			✓
<i>pairwise allpair</i>	✓	✓	✓	
<i>pairwise heapsort</i>	✓	✓		✓
<i>pairwise bubblesort</i>	✓	✓		✓
<i>setwise heapsort</i>	✓	✓		✓
<i>setwise bubblesort</i>	✓	✓		✓

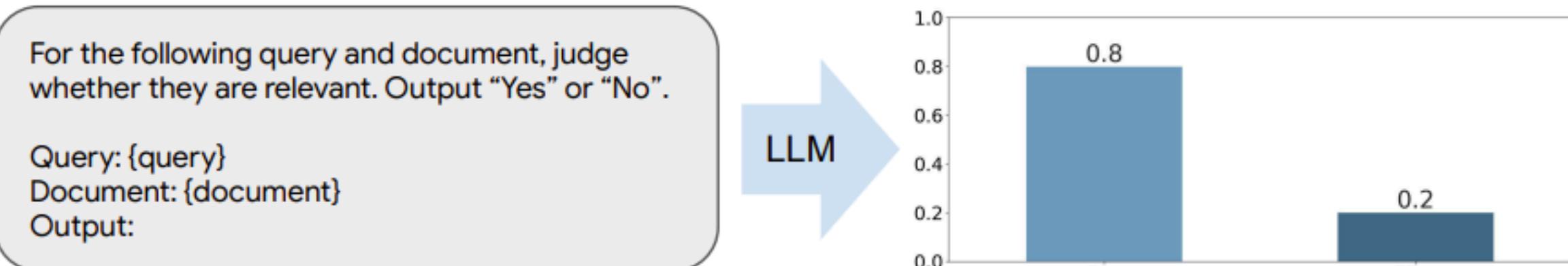
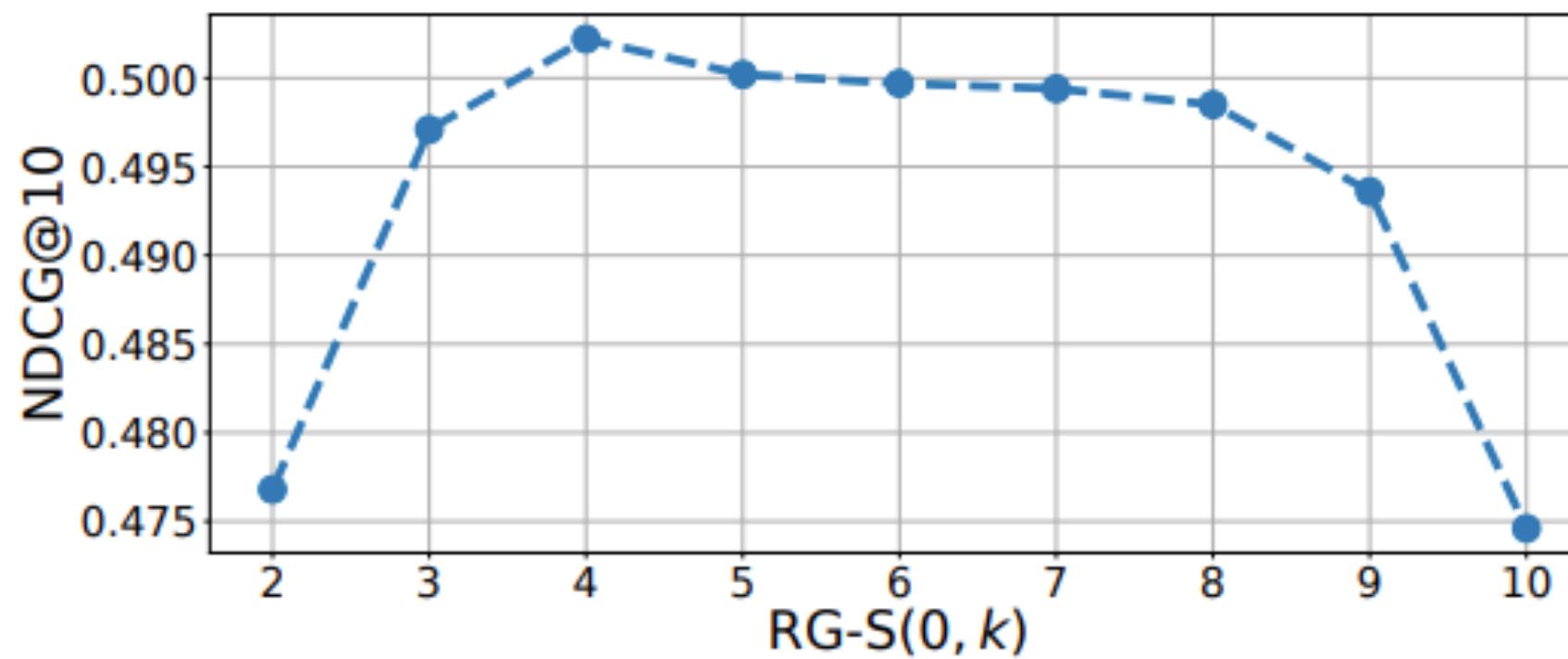
# Properties

Methods	Logits	Generate	Batching	Top- $k$	# LLM calls
<i>pointwise qlm</i>	✓		✓		$O(N)$
<i>pointwise yes_no</i>	✓		✓		$O(N)$
<i>listwise generation</i>		✓		✓	$O(r * (N/s))$
<i>listwise likelihood</i>	✓			✓	$O(r * (N/s))$
<i>pairwise allpair</i>	✓	✓	✓		$O(N^2 - N)$
<i>pairwise heapsort</i>	✓	✓		✓	$O(k * \log_2 N)$
<i>pairwise bubblesort</i>	✓	✓		✓	$O(k * N)$
<i>setwise heapsort</i>	✓	✓		✓	$O(k * \log_c N)$
<i>setwise bubblesort</i>	✓	✓		✓	$O(k * (N/(c-1)))$

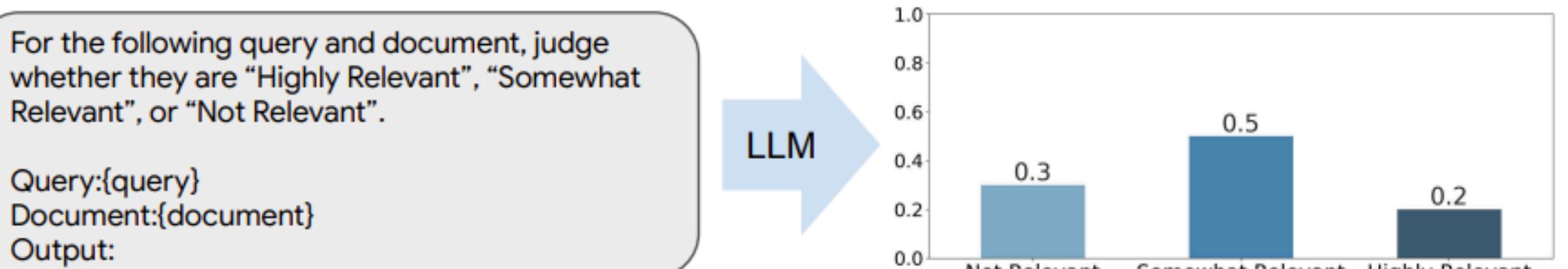
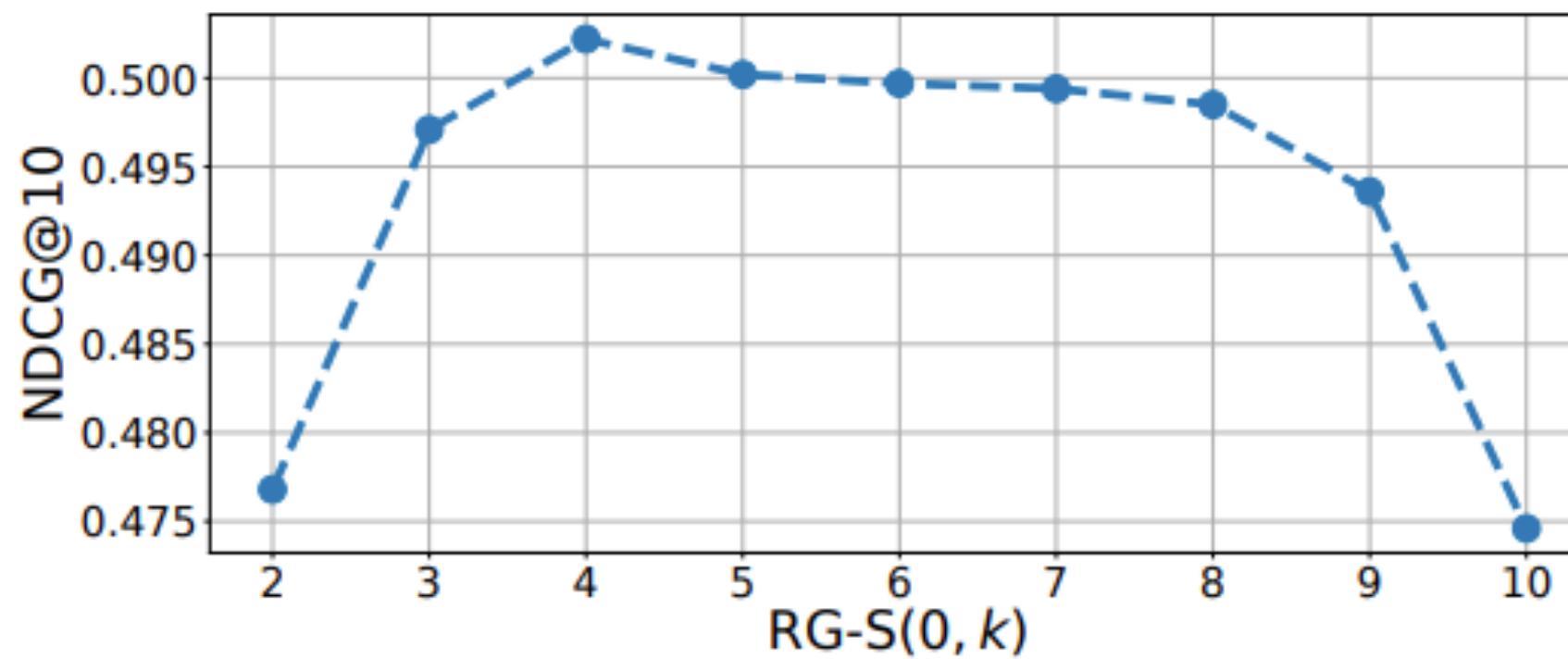
# Zeroshot LLM-based Rankers' Variations

# Pointwise Variation

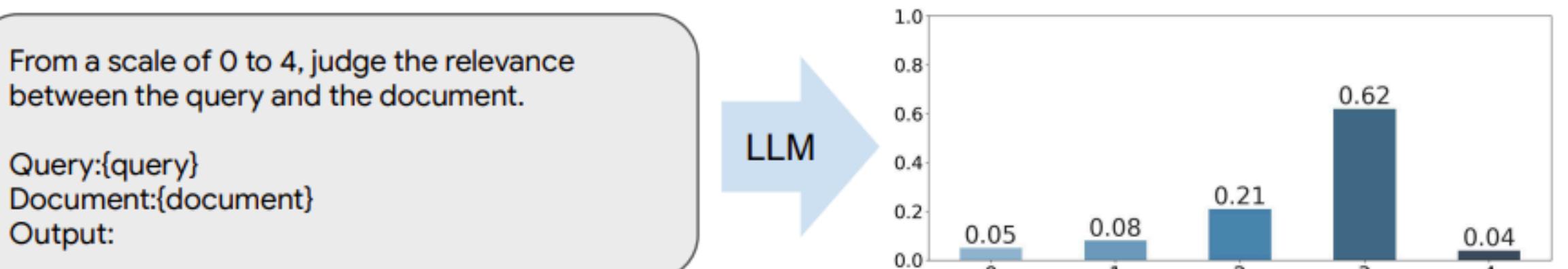
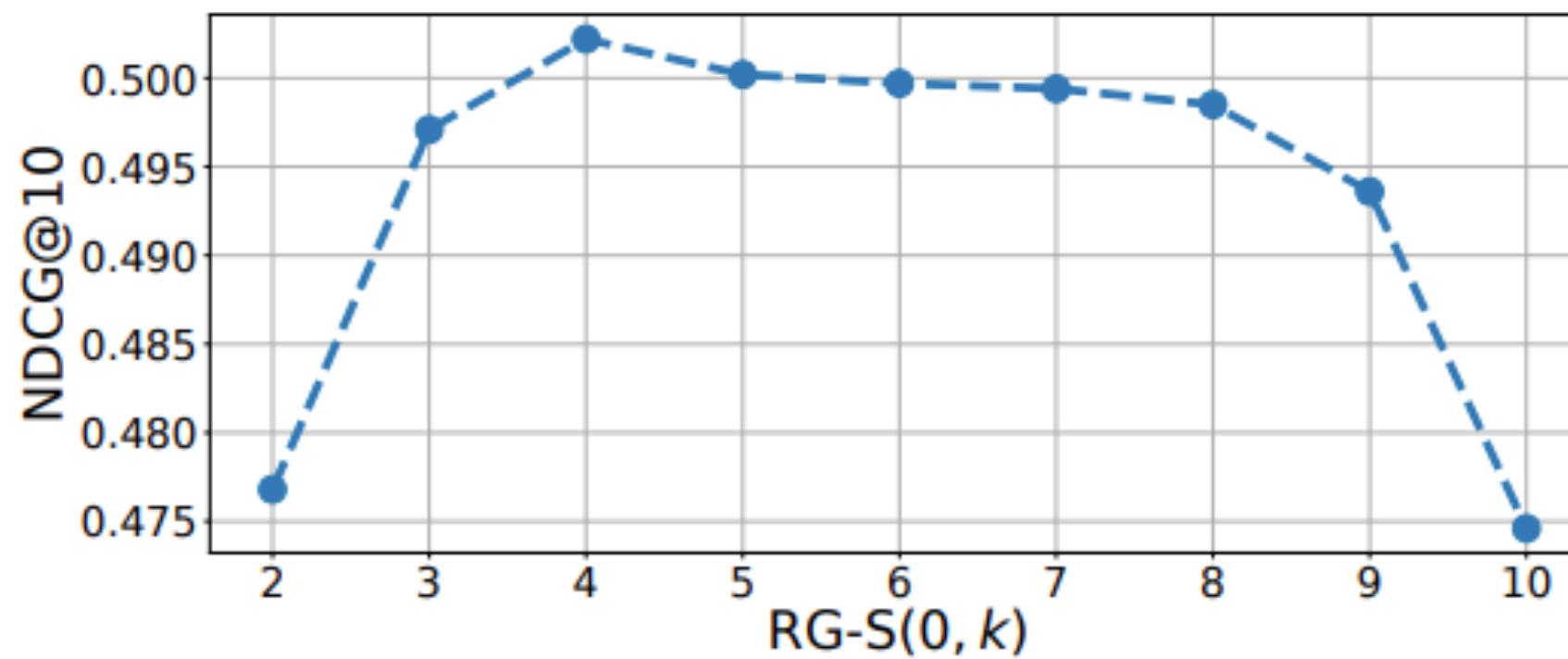
- Beyond yes and no.



(a) Yes-No relevance generation



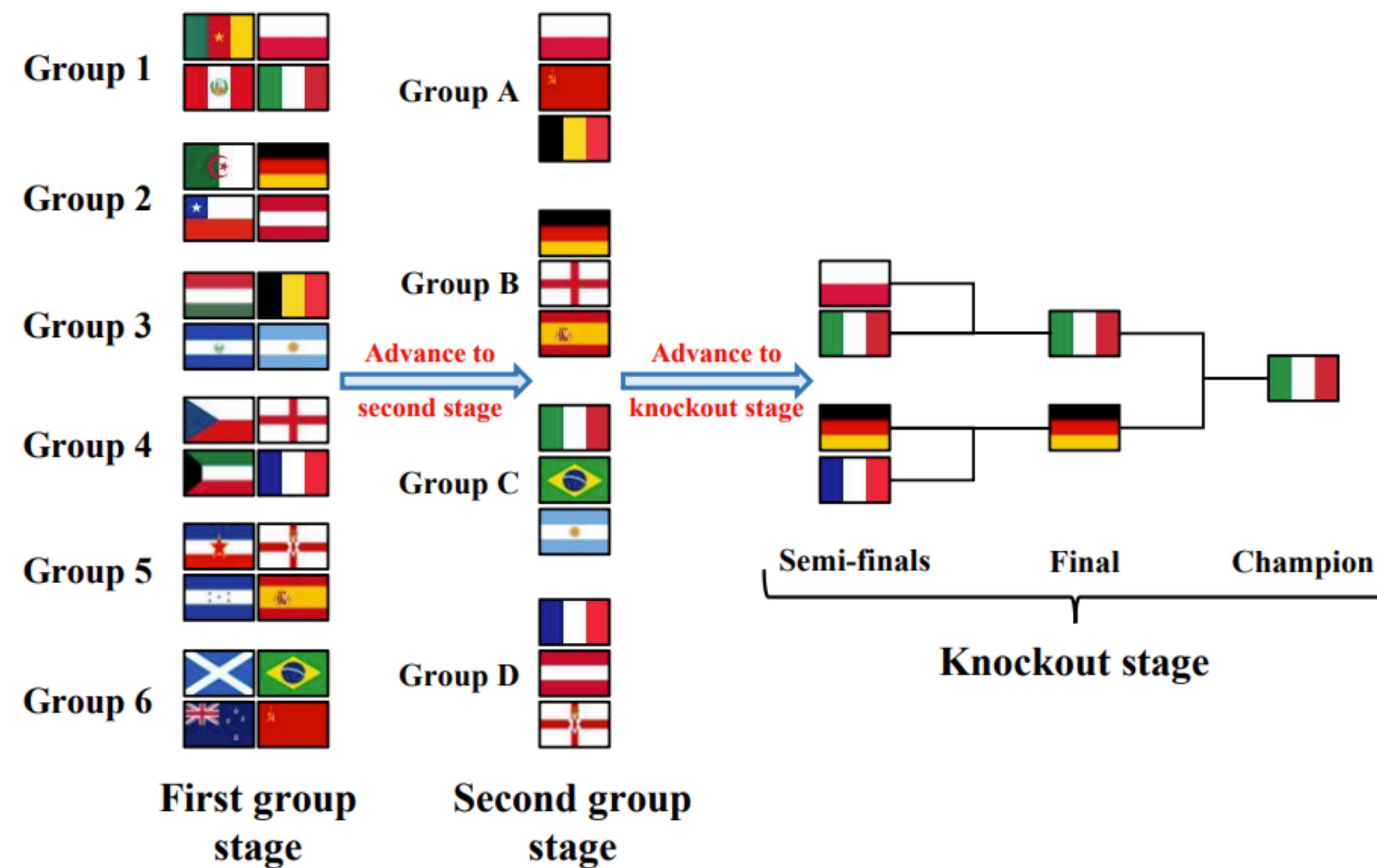
(b) Fine-grained relevance label generation



(c) Rating scale relevance generation

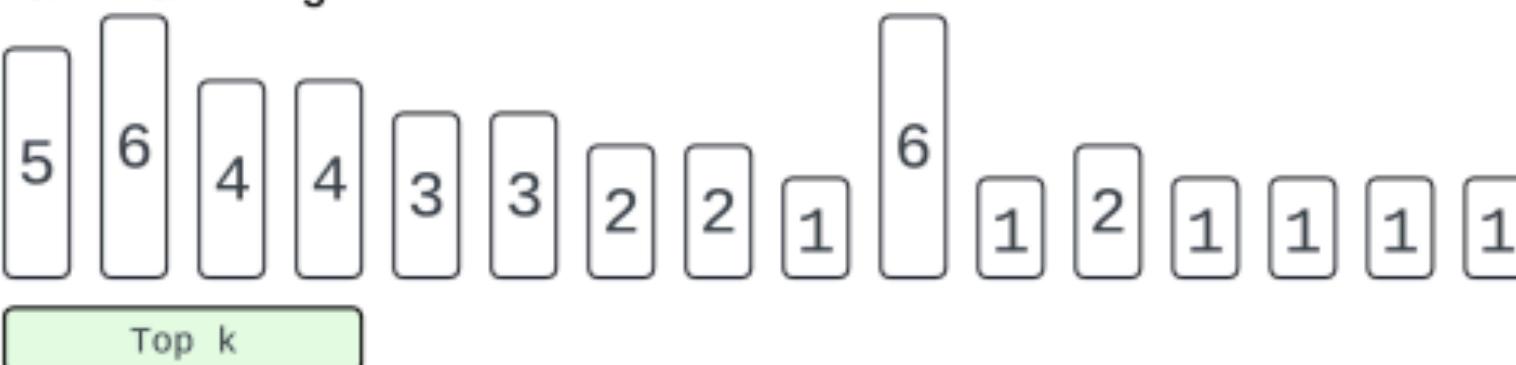
# TourRank

- Sport tournaments scoring

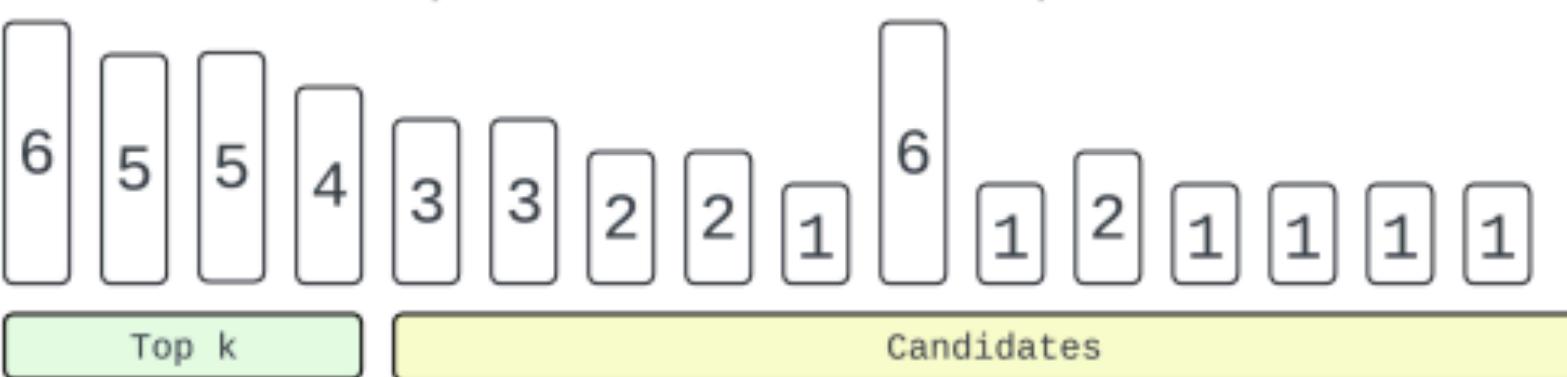


# Setwise Insertion sort

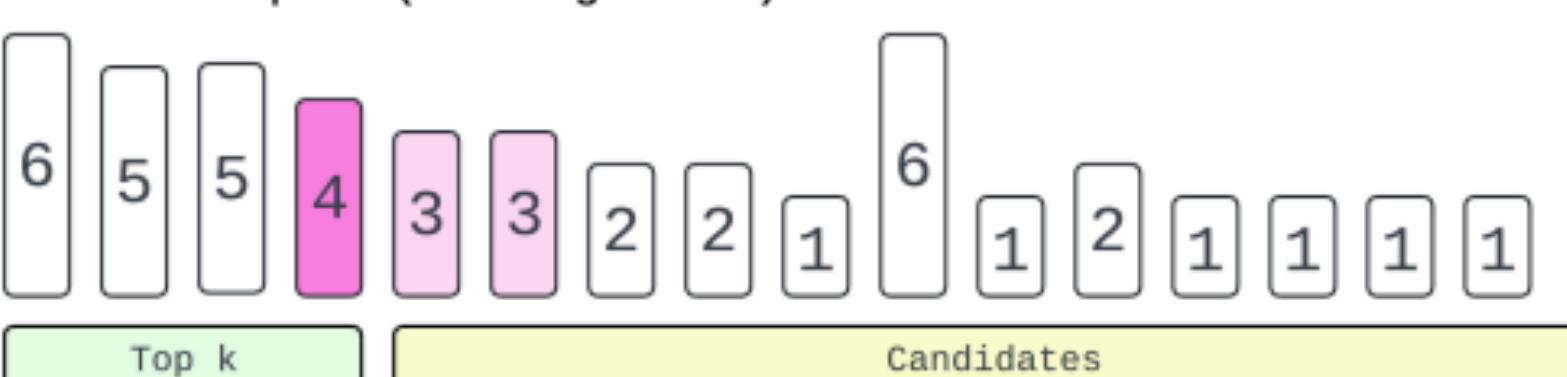
0. Start with a partially sorted list based on an initial ranking.



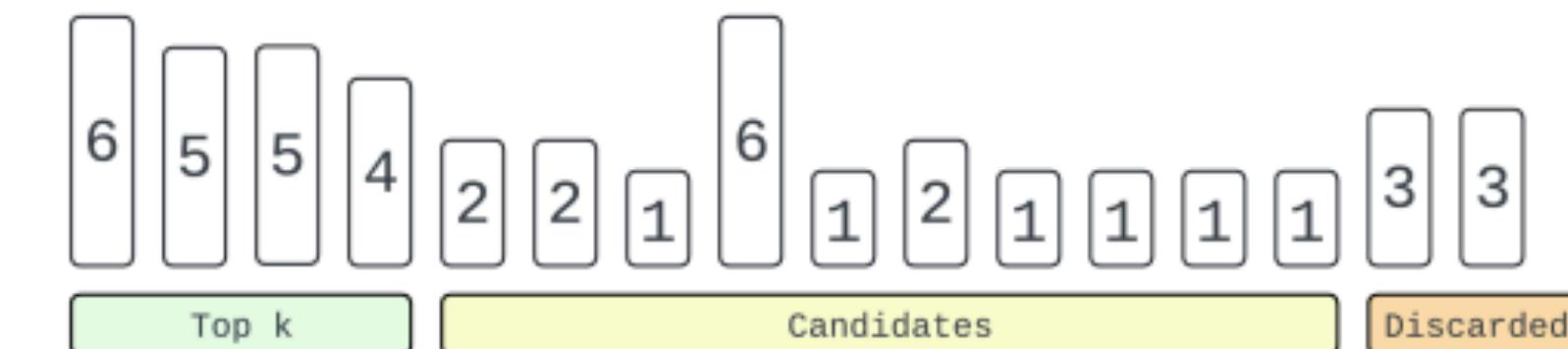
1. Use setwise.heapsort to sort the top-k elements.



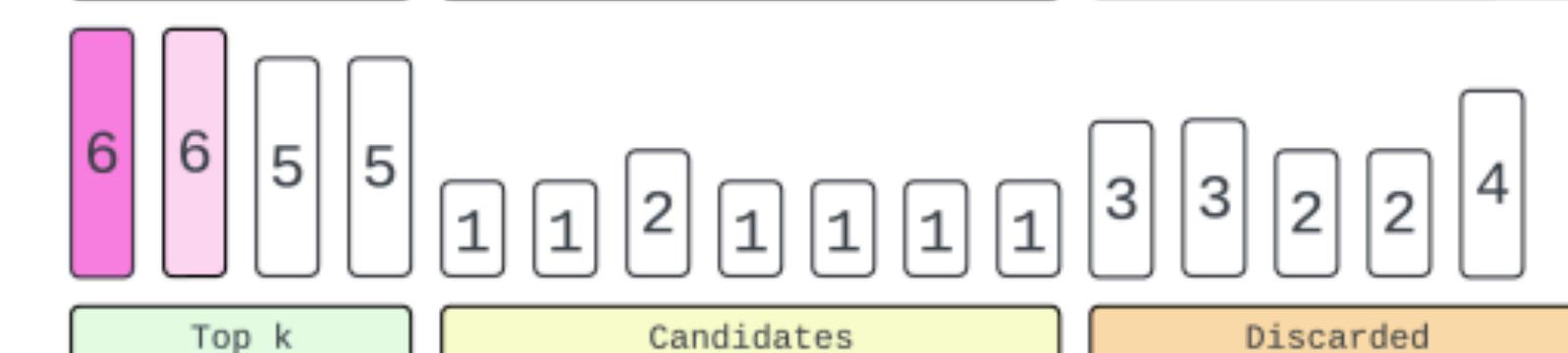
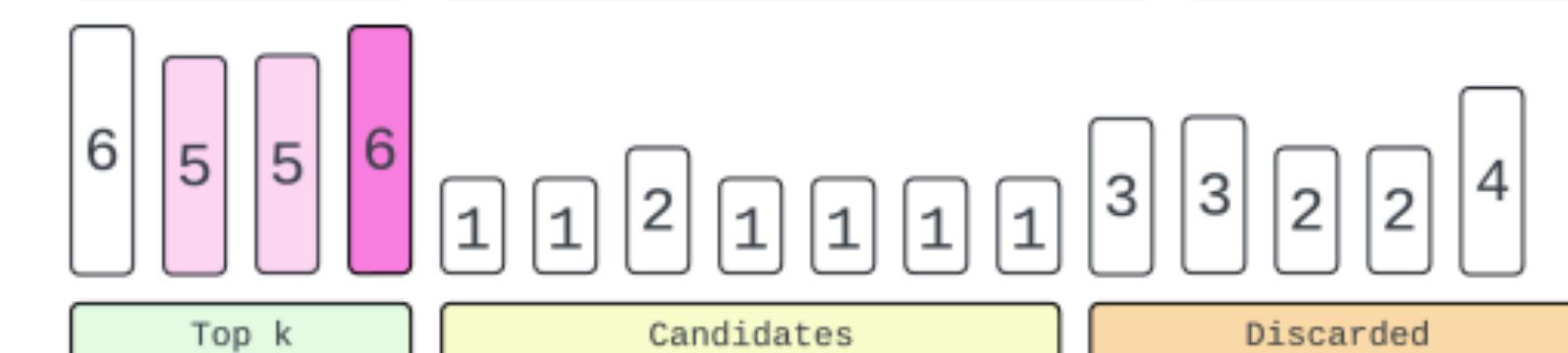
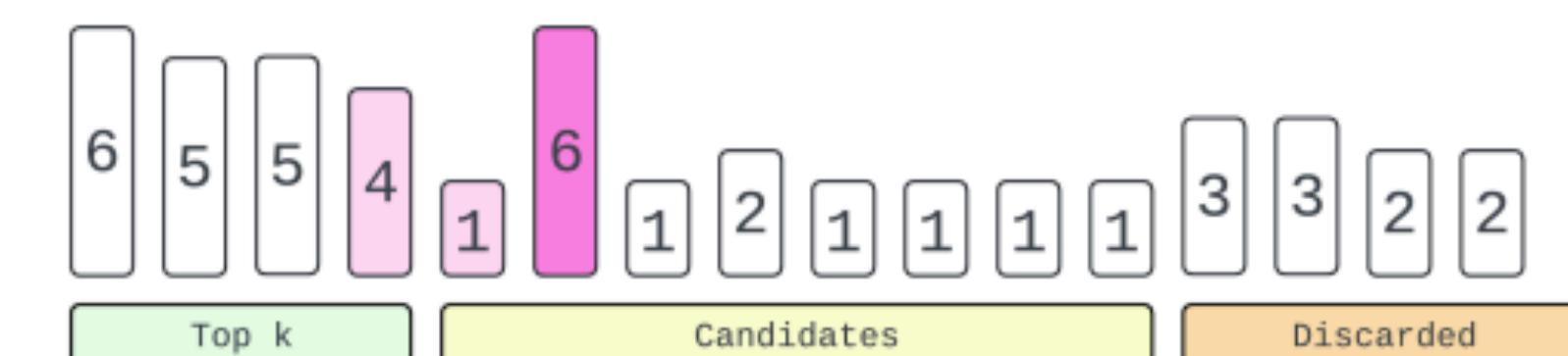
2. Compare small candidate sets with the smallest element in top-k (the "guard").



3. Discard candidates if the guard is the largest.

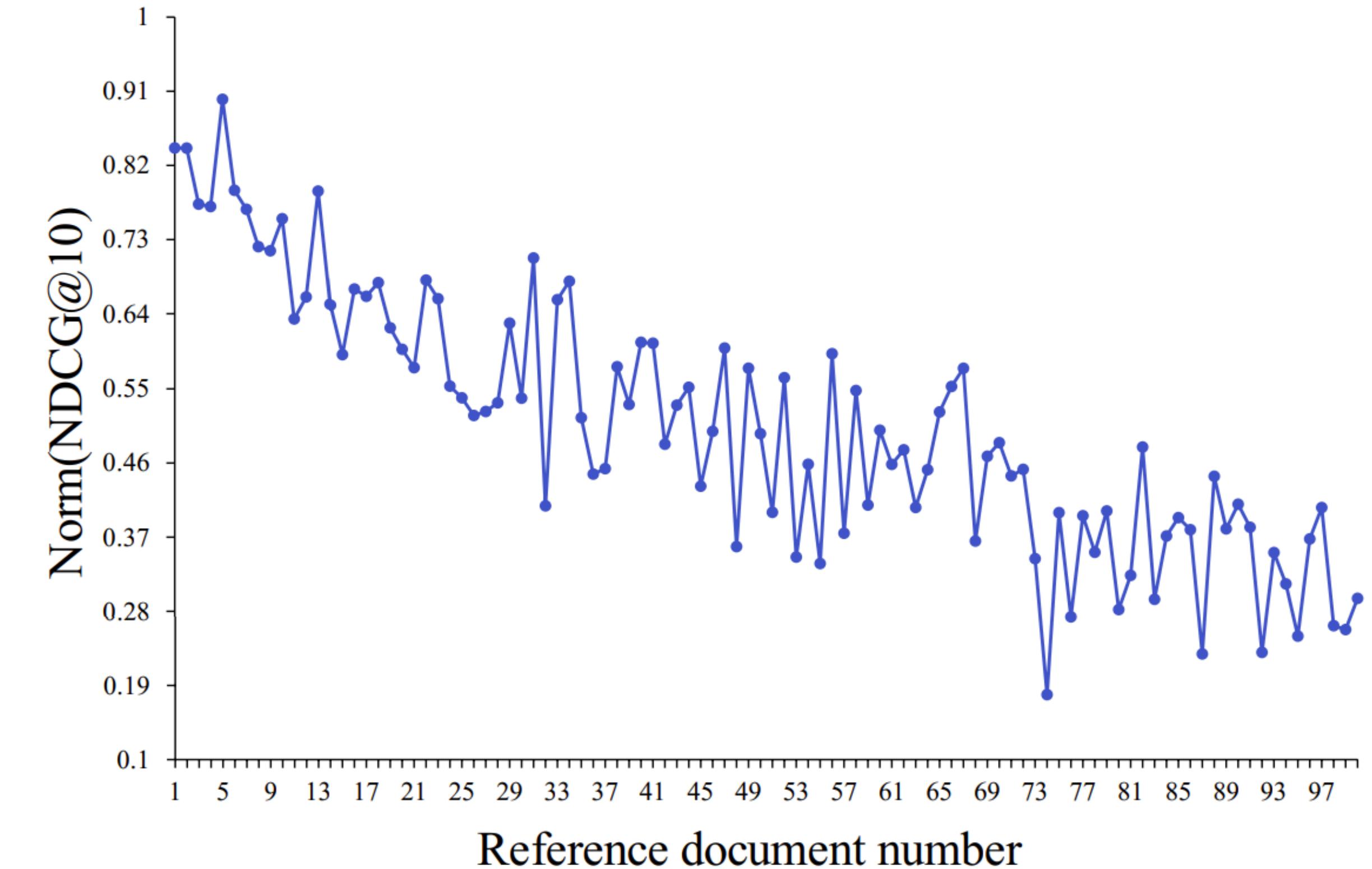
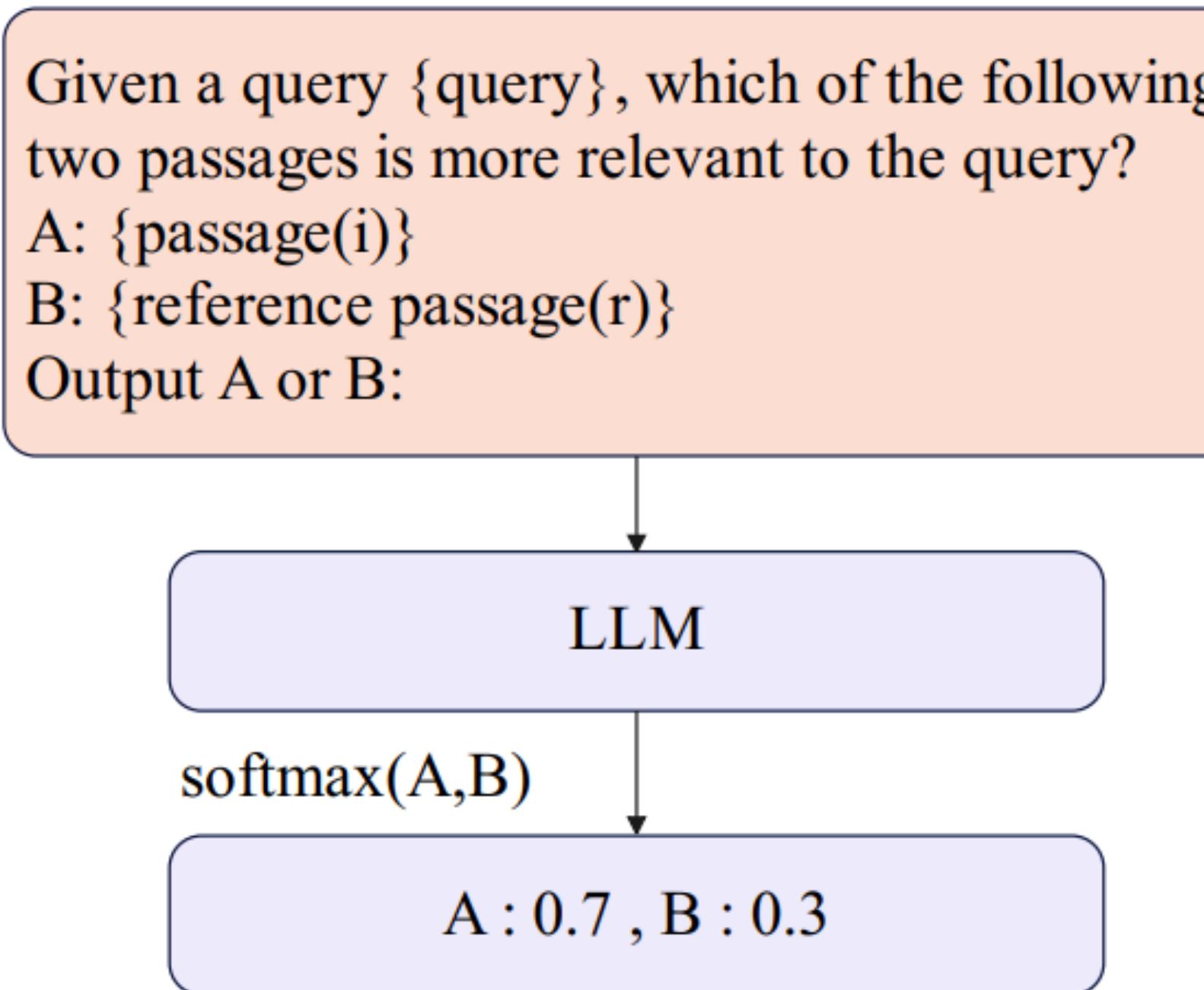


4. If at least one candidate is larger than the guard, discard the guard, insert the candidate into top k, and fix the order.



5. Repeat steps 2-4 until the candidate list is empty.

# Combine Pointwise and Pairwise



- Li, J., Hu, X., Zhao, Y., Zhuang, S., Zhang, H. 2025. Beyond Reproducibility: Advancing Zero-shot LLM Reranking Efficiency with Setwise Insertion. arXiv
- Long, K., Li, S., Xu, C., Tang, J., Wang, T. 2025. Precise Zero-Shot Pointwise Ranking with LLMs through Post-Aggregated Global Context Information. SIGIR
- Wisznia, J., Bolanos, C., Tollo, J., Marraffini, G., Gianolini, A., Hsueh, N., Corro, L. 2025, Are Optimal Algorithms Still Optimal? Rethinking Sorting in LLM-Based Pairwise Ranking with Batching and Caching. ACL

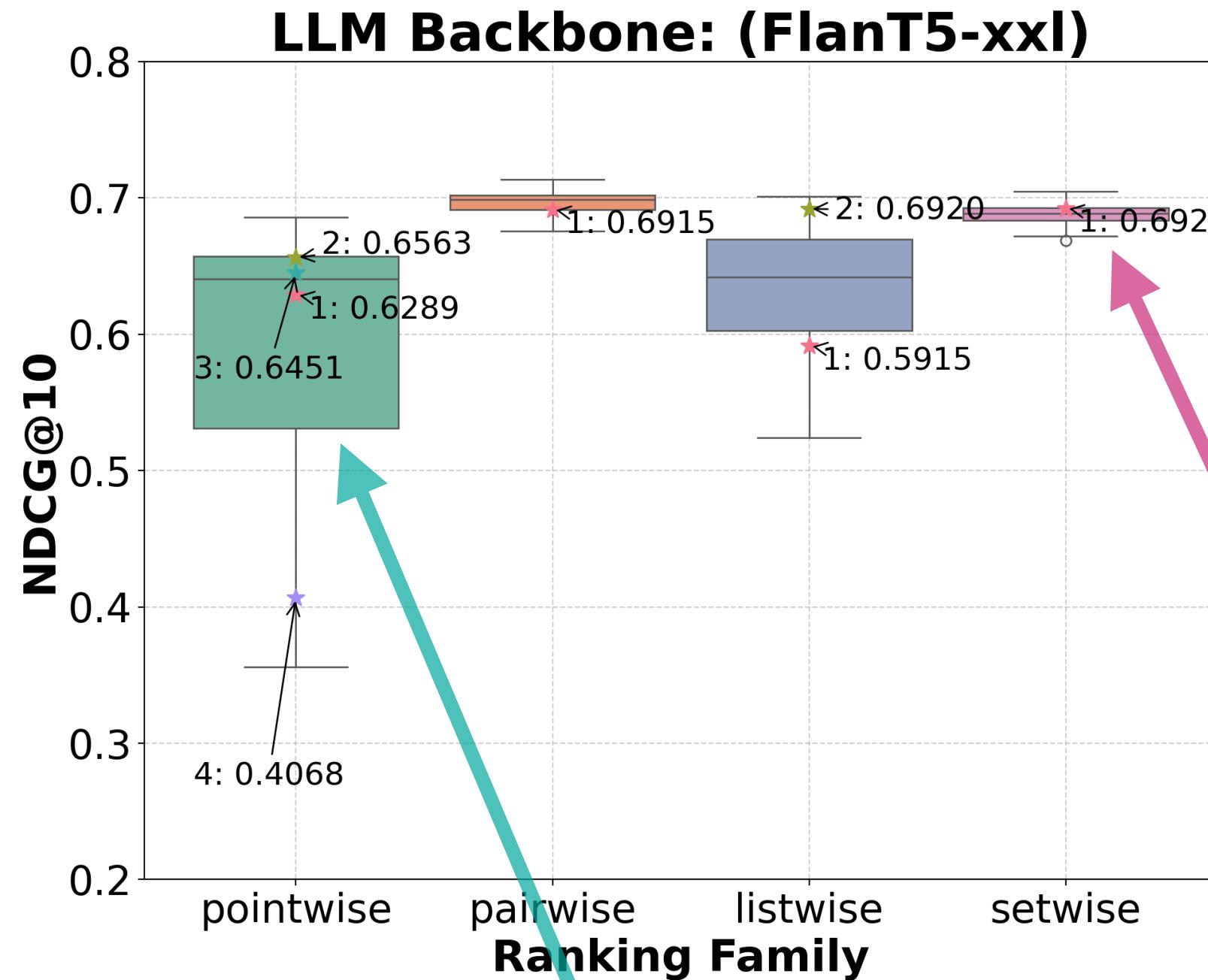
# Impact of Prompt Variations on LLM Rankers

- 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?	<b>Role Playing</b>
RankGPT, Sun et al.	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}. {PASSAGES} Search Query: {query}. 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 [] > [], e.g., [1] > [2]. Only response the ranking results, do not say any word or explain.	<b>Formatting Instr.</b> <b>Restriction on Output</b>

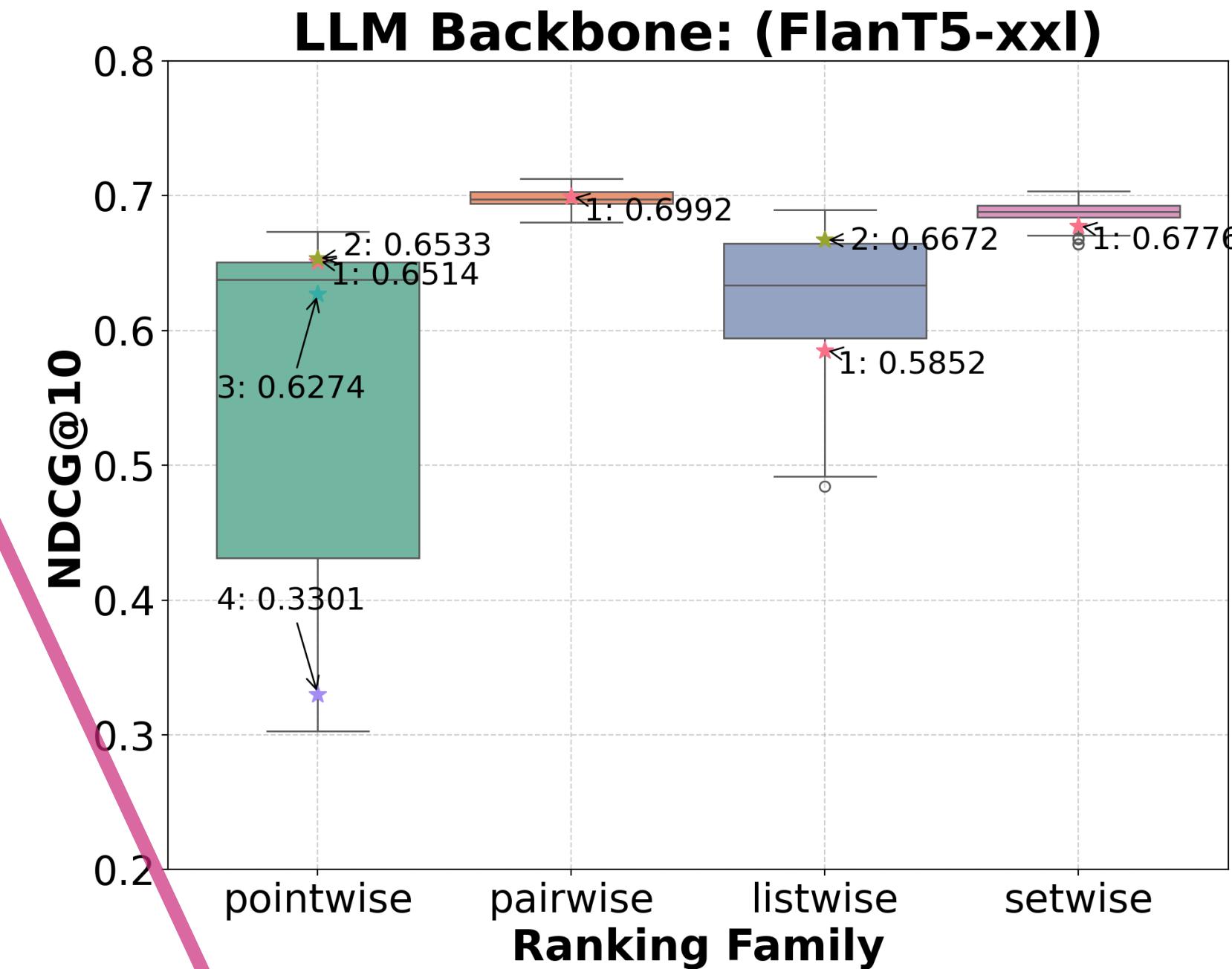
# Impact of Prompt Variations on LLM Rankers

TREC DL 2019



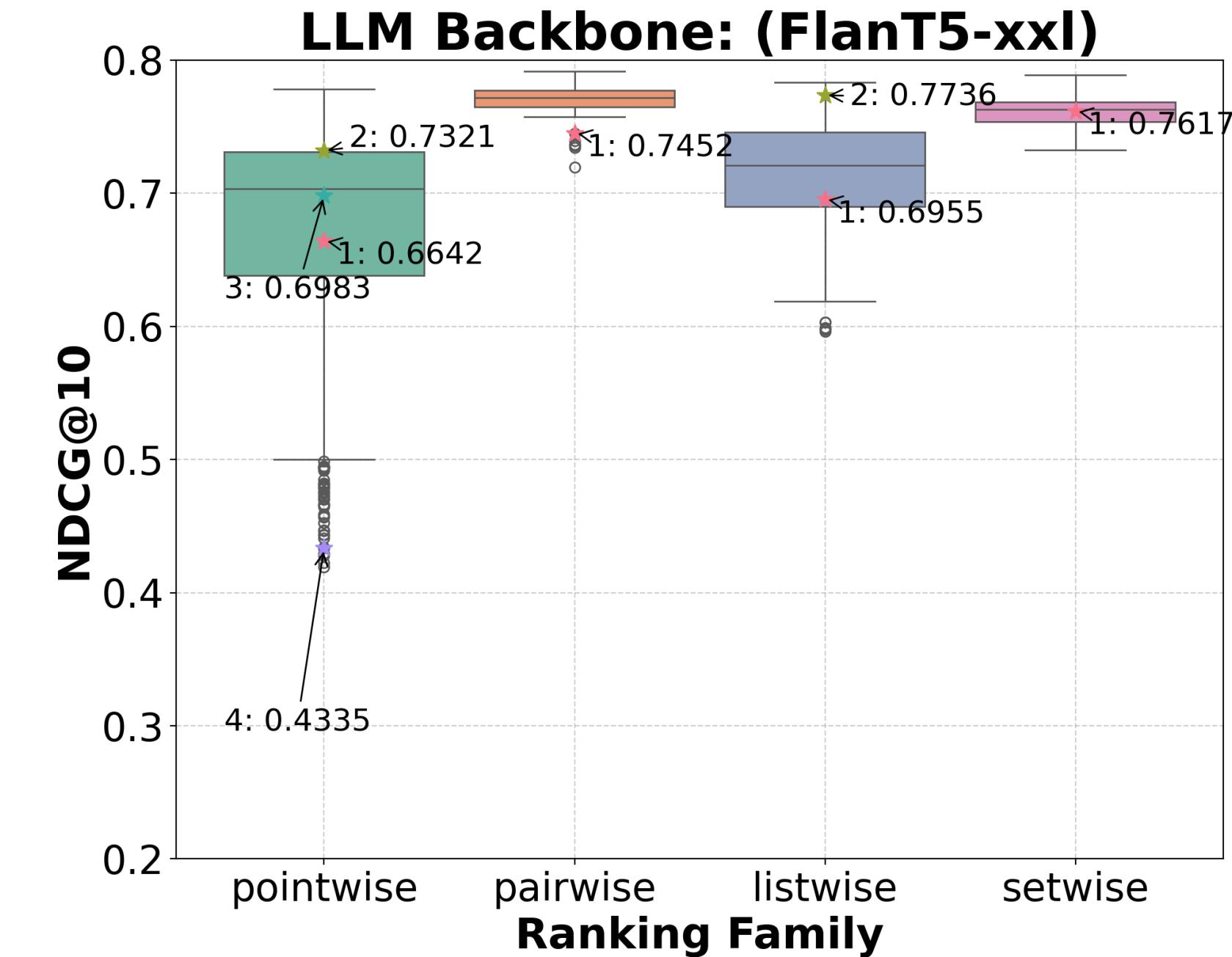
Pointwise, listwise often display large variations of effectiveness across different prompts

TREC DL 2020



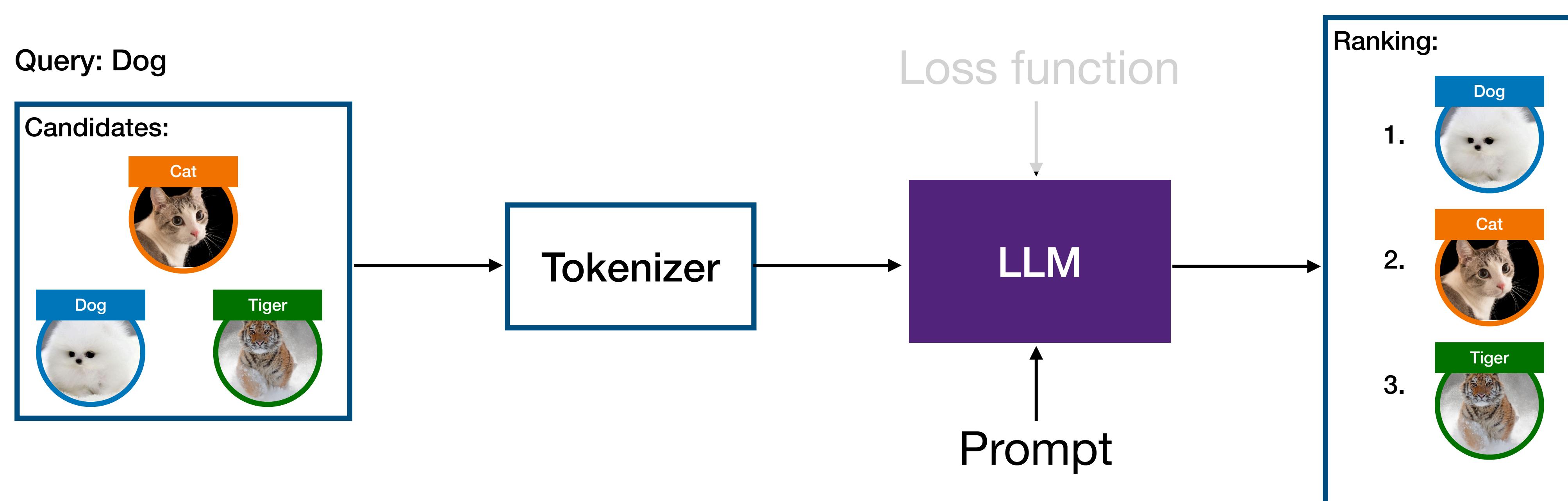
Pairwise, set wise very robust across different prompts

COVID (BEIR)

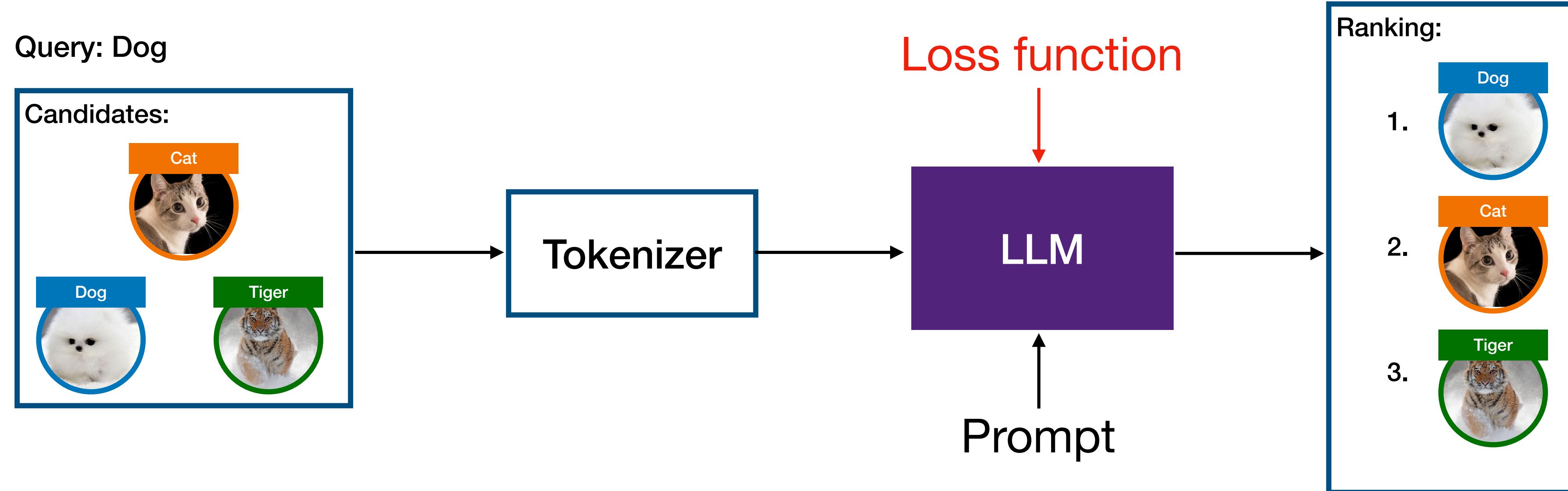


For all approaches, there are better prompts than the one in the original paper

# LLM as Rankers



# LLM as Rankers



# Distillation

- RankZephyr: 7B listwise ranker distilled from large GPT Models, outperforms GPT Models

# Distillation

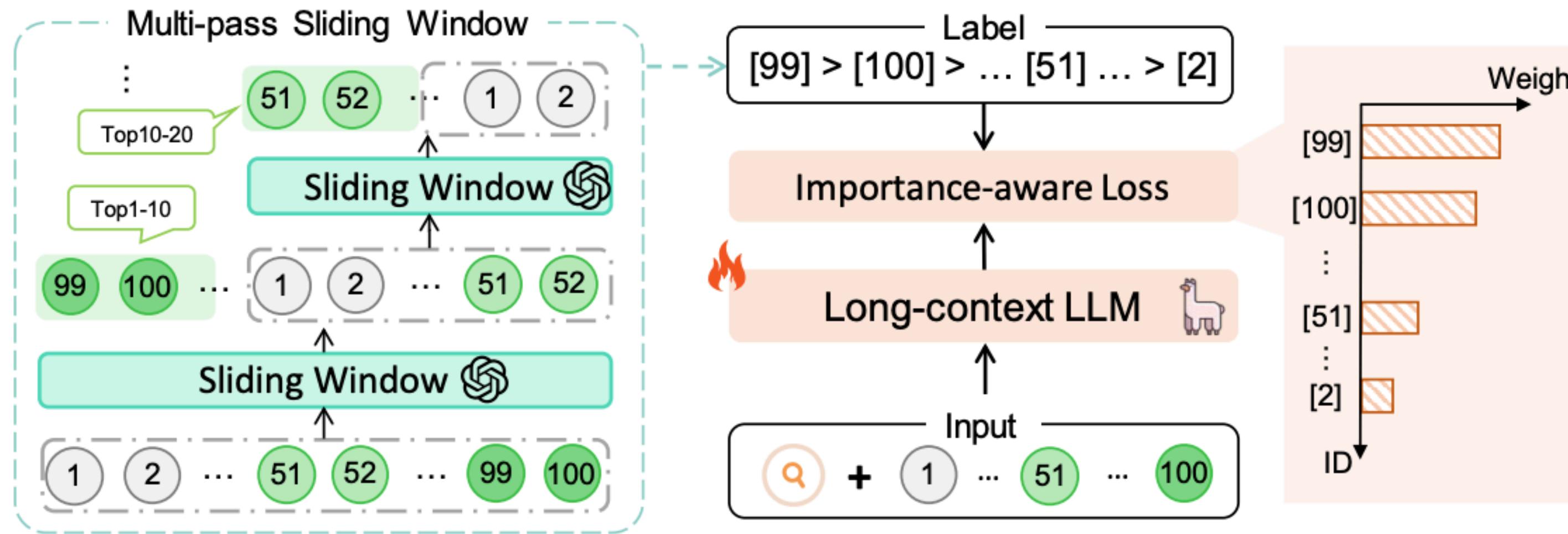
- RankZephyr: 7B listwise ranker distilled from large GPT Models, outperforms GPT Models
- Stage1: 100k training queries from MS MARCO, with rankings provided by ChatGPT 3.5.
  - Excludes malformed generations.
  - Has randomly shuffled input order.

# Distillation

- RankZephyr: 7B listwise ranker distilled from large GPT Models, outperforms GPT Models
- Stage1: 100k training queries from MS MARCO, with rankings provided by ChatGPT 3.5.
  - Excludes malformed generations.
  - Has randomly shuffled input order.
- Stage2: 5k queries with rankings from GPT-4.5.
  - Diverse queries to maximize query embedding differences.

# Distillation

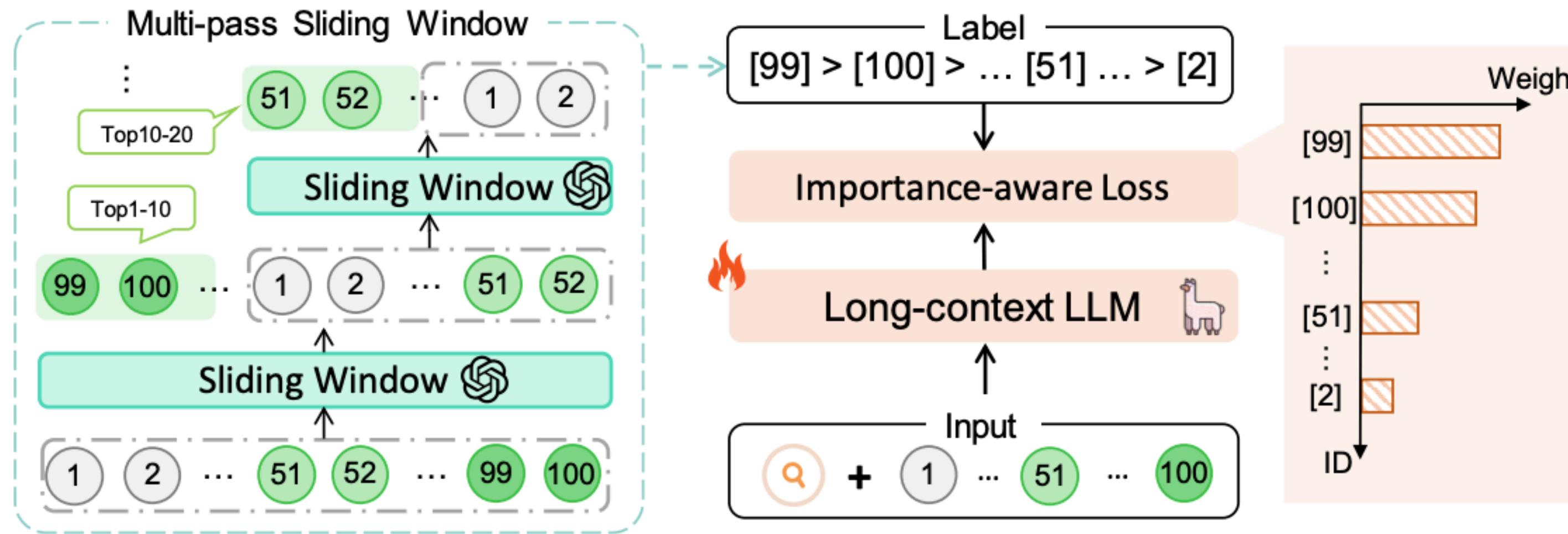
- Sliding Windows vs Full ranking distillation.



Models	Strategy	DL19	DL20	Covid	DBpedia	SciFact	NFCorpus	Signal	Robust04	Touche	News	Avg.	
		Zero-shot											
GPT-4o	Sliding	<b>74.78</b>	69.52	<b>83.41</b>	<b>45.56</b>	<b>77.41</b>	<b>39.67</b>	<b>34.20</b>	<b>60.25</b>	32.26	<b>51.92</b>	<b>53.09</b>	
	Full	73.94	<b>70.03</b>	82.10	43.31	74.85	39.00	32.63	55.95	30.42	47.96	50.78	
SFT from GPT-4o													
RankMistral <sub>20</sub>	Sliding	70.34	69.58	80.86	42.52	75.82	38.38	33.90	54.69	<b>37.18</b>	<b>51.42</b>	51.85	
RankMistral <sub>100</sub>	Full	<b>72.55</b>	<b>71.29</b>	<b>82.24</b>	<b>43.54</b>	<b>77.04</b>	<b>39.14</b>	<b>33.99</b>	<b>57.91</b>	34.63	50.59	<b>52.40</b>	

# Distillation

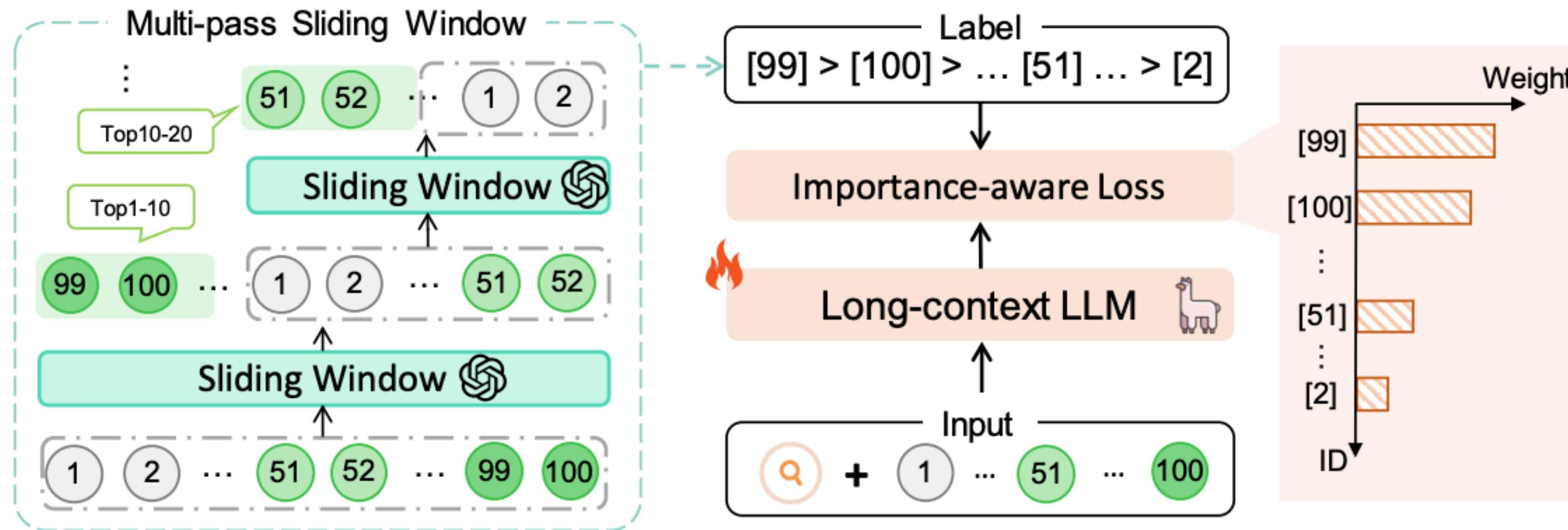
- Sliding Windows vs Full ranking distillation.



Models	Strategy	DL19		DL20		Covid	DBpedia	SciFact	NFCorpus	Signal	Robust04	Touche	News	Avg.
		Zero-shot												
GPT-4o	Sliding	<b>74.78</b>	69.52	<b>83.41</b>	<b>45.56</b>	<b>77.41</b>	<b>39.67</b>	<b>34.20</b>	<b>60.25</b>	32.26	<b>51.92</b>	<b>53.09</b>		
	Full	73.94	<b>70.03</b>	82.10	43.31	74.85	39.00	32.63	55.95	30.42	47.96	50.78		
SFT from GPT-4o														
RankMistral <sub>20</sub>	Sliding	70.34	69.58	80.86	42.52	75.82	38.38	33.90	54.69	<b>37.18</b>	<b>51.42</b>	51.85		
RankMistral <sub>100</sub>	Full	<b>72.55</b>	<b>71.29</b>	<b>82.24</b>	<b>43.54</b>	<b>77.04</b>	<b>39.14</b>	<b>33.99</b>	<b>57.91</b>	34.63	50.59	<b>52.40</b>		

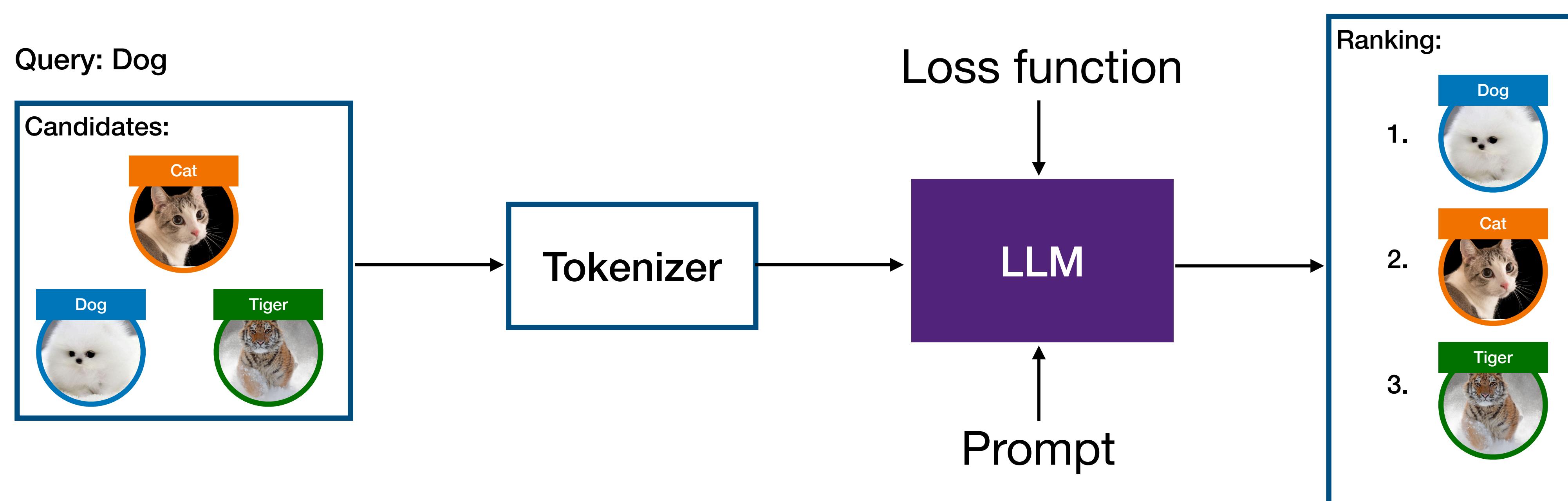
# Distillation

- Sliding Windows vs Full ranking distillation.



Models	Strategy	DL19	DL20	Covid	DBpedia	SciFact	NFCorpus	Signal	Robust04	Touche	News	Avg.	
		Zero-shot											
GPT-4o	Sliding	<b>74.78</b>	69.52	<b>83.41</b>	<b>45.56</b>	<b>77.41</b>	<b>39.67</b>	<b>34.20</b>	<b>60.25</b>	32.26	<b>51.92</b>	<b>53.09</b>	
	Full	73.94	<b>70.03</b>	82.10	43.31	74.85	39.00	32.63	55.95	30.42	47.96	50.78	
SFT from GPT-4o													
RankMistral <sub>20</sub>	Sliding	70.34	69.58	80.86	42.52	75.82	38.38	33.90	54.69	<b>37.18</b>	<b>51.42</b>	51.85	
RankMistral <sub>100</sub>	Full	<b>72.55</b>	<b>71.29</b>	<b>82.24</b>	<b>43.54</b>	<b>77.04</b>	<b>39.14</b>	<b>33.99</b>	<b>57.91</b>	34.63	50.59	<b>52.40</b>	

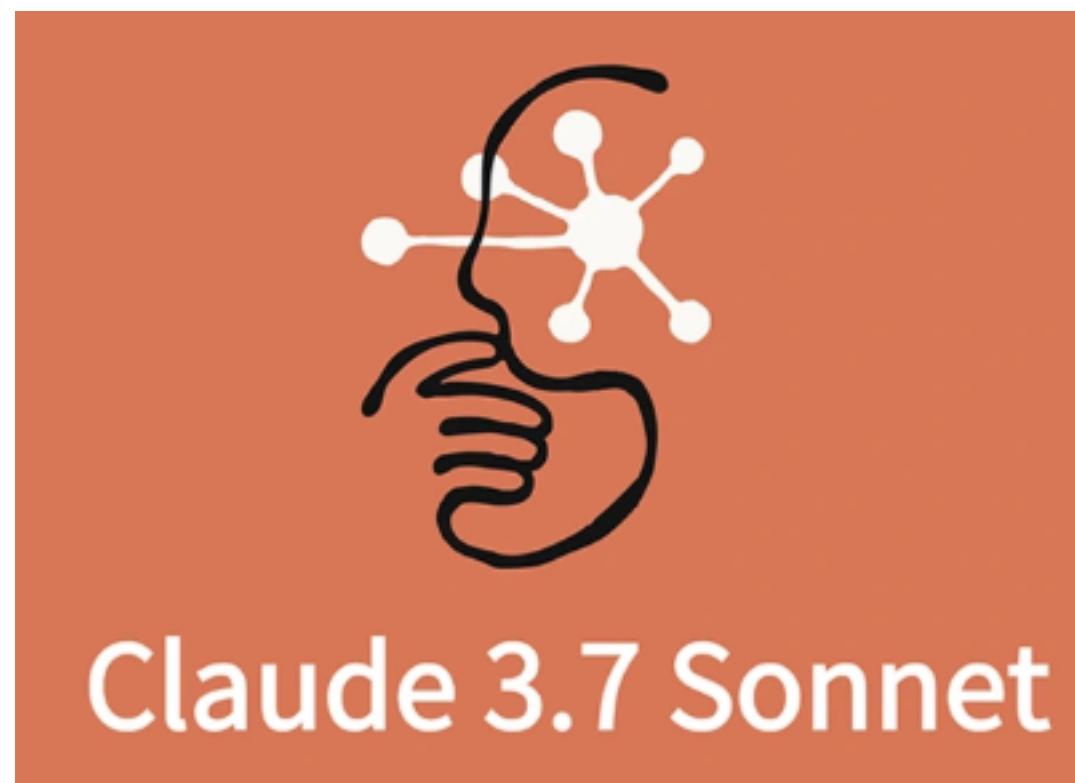
# LLM as Rankers



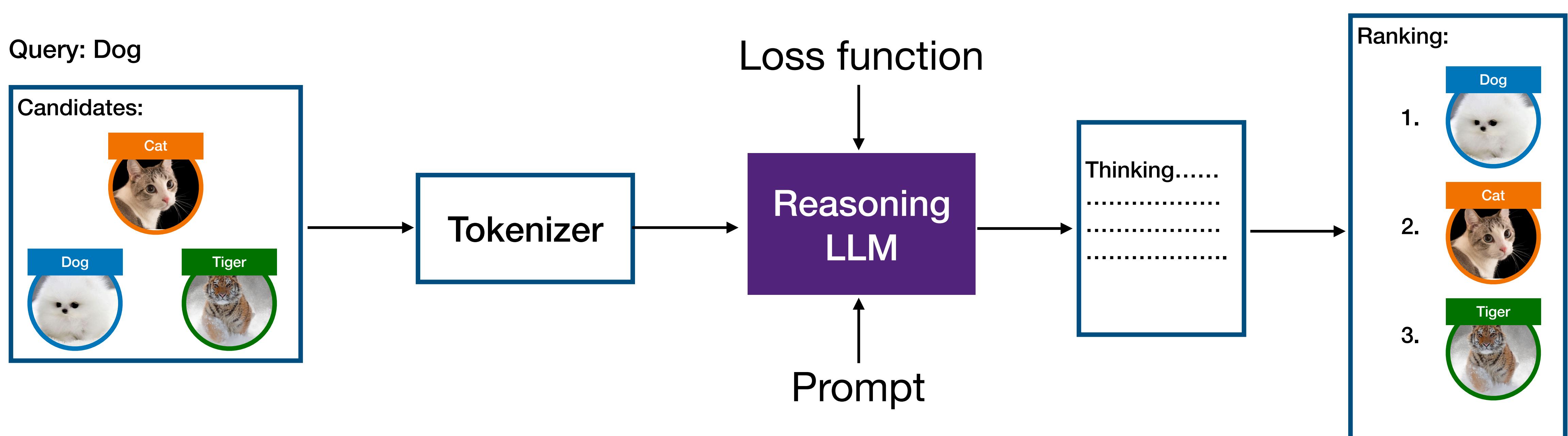
# Reasoning LLMs



Deepseek R1



# Reasoning Rankers

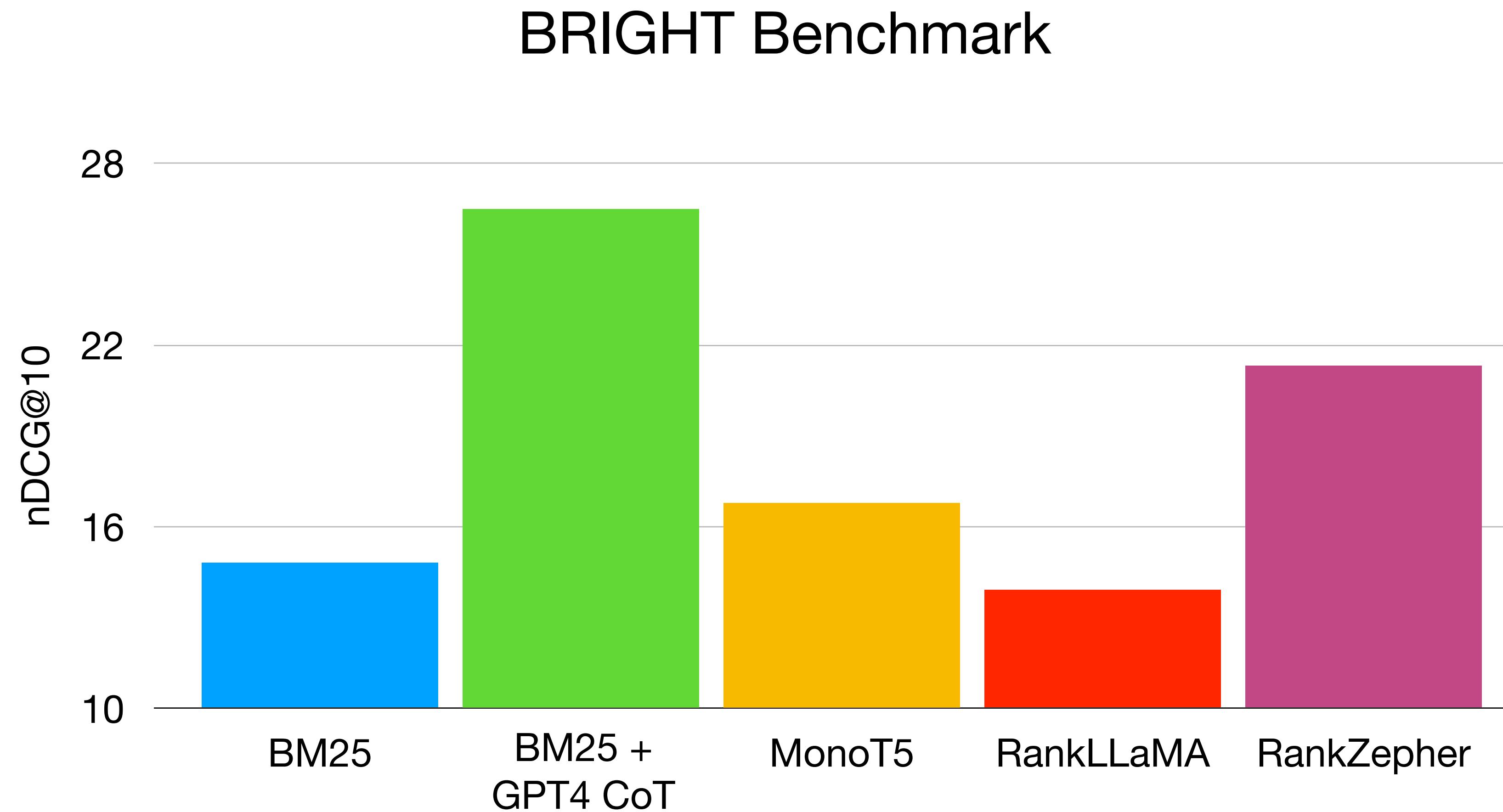


# Why Reasoning?

## BRIGHT Benchmark

Dataset	Total Number			Avg. Length		Source		Examples	
	Q	$\mathcal{D}$	$\mathcal{D}^+$	Q	$\mathcal{D}$	Q	$\mathcal{D}$		
<i>StackExchange</i>									
Biology	103	57,359	3.6	115.2	83.6	StackExchange post	Web pages: article, tutorial, news, blog, report ...	Tab. 20	
Earth Science	116	121,249	5.3	109.5	132.6			Tab. 21	
Economics	103	50,220	8.0	181.5	120.2			Tab. 22	
Psychology	101	52,835	7.3	149.6	118.2			Tab. 23	
Robotics	101	61,961	5.5	818.9	121.0			Tab. 24	
Stack Overflow	117	107,081	7.0	478.3	704.7			Tab. 25	
Sustainable Living	108	60,792	5.6	148.5	107.9			Tab. 26	
<i>Coding</i>									
LeetCode	142	413,932	1.8	497.5	482.6	Coding question	Coding Q&Sol	Tab. 27	
Pony	112	7,894	22.5	102.6	98.3	Coding question	Syntax Doc	Tab. 28	
<i>Theorems</i>									
AoPS	111	188,002	4.7	117.1	250.5	Math Olympiad Q	STEM Q&Sol	Tab. 29	
TheoremQA-Q	194	188,002	3.2	93.4	250.5	Theorem-based Q	STEM Q&Sol	Tab. 30	
TheoremQA-T	76	23,839	2.0	91.7	354.8	Theorem-based Q	Theorems	Tab. 31	

# Why Reasoning?



# Reasoning Rankers

- JudgeRank: Zeroshot pointwise reasoning ranker.

## **Query Analysis Prompt**

You will be presented with a/an {query name}.

Your task consists of the following step:

1. Analyze the {query name}:

- Carefully read each sentence of the {query name}.
- Identify the core problem or question being asked.

Here is the {query name}:  
{query}

(a). Query Analysis Prompt of JudgeRank

## **Document Analysis Prompt**

You will be presented with a/an {query name}, an analysis of the query, and a/an {doc name}.

Your task consists of the following steps:

1. Analyze the {doc name}:

- Thoroughly examine each sentence of the {doc name}.
- List all sentences from the {doc name} that {definition of relevance} the {query name}.
- Briefly explain how each sentence listed {definition of relevance} the {query name}.

2. Assess overall relevance:

- If the {doc name}, particularly the relevant sentences (if applicable), {definition of relevance} the {query name}, briefly explain why.
- Otherwise, briefly explain why not.

Here is the {query name}:  
{query}

Here is the analysis of the {query name}:  
{query analysis}

Here is the {doc name}:  
{doc}

(b). Document Analysis Prompt of JudgeRank

## **Judgment Prompt**

You will be presented with a/an {query name}, an analysis of the {query name}, a/an {doc name}, and an analysis of the {doc name}.

Your task is to assess if the {doc name} {definition of relevance} the {query name} in one word:

- Yes: If the {doc name} {definition of relevance} the {query name}.
- No: Otherwise.

Important: Respond using only one of the following two words without quotation marks: Yes or No.

Here is the {query name}:  
{query}

Here is the analysis of the {query name}:  
{query analysis}

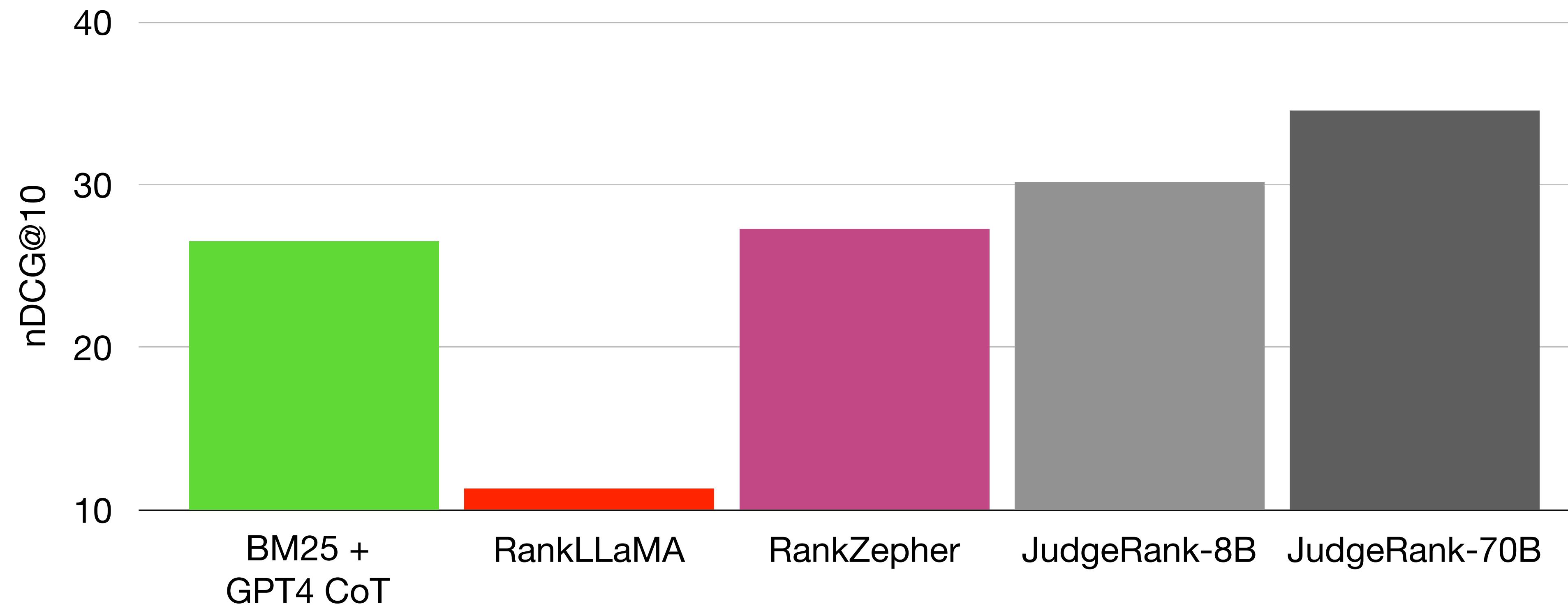
Here is the {doc name}:  
{doc}

Here is the analysis of the {doc name}:  
{doc analysis}

(c). Judgment Prompt of JudgeRank

# Reasoning Rankers

- JudgeRank: BRIGHT results



# Reasoning Rankers

- Rank1: Pointwise reasoning ranker distilled from DeepSeek-R1



what county is colton in



Colton, California. Colton is a city in San Bernardino County, California, United States...



<think> The user is asking "what county is Colton in." ... So the first thing I notice is that the passage directly mentions San Bernardino County as the county where Colton is located.

The query is specifically about the county, and the answer is right there in the first sentence of the passage.

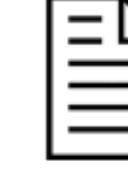
But wait, maybe I should check if there's any possibility that the passage isn't relevant. The rest of the passage talks ... but the key information about San Bernardino County is present ...

So the answer should be true because the passage answers the query directly.

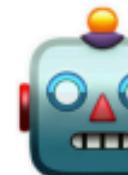
</think> **true**



do snow leopards change color



Snow leopards have long, thick fur, and their base color varies from smoky gray to yellowish tan, with whitish underparts ...



<think> The user is asking if snow leopards change color. The passage provided describes their fur color as varying from smoky gray to yellowish tan, with some patterns. Hmm, the key here is whether "varies" implies that they change color over time, like seasonally or due to some factors.

But wait, "varies" might just mean that different individuals have different base colors, not that a single leopard changes its color. The passage doesn't mention anything about changing color; it's just describing their appearance. ...

Therefore, the passage isn't directly relevant to the question of whether they change color. So the answer should be false.

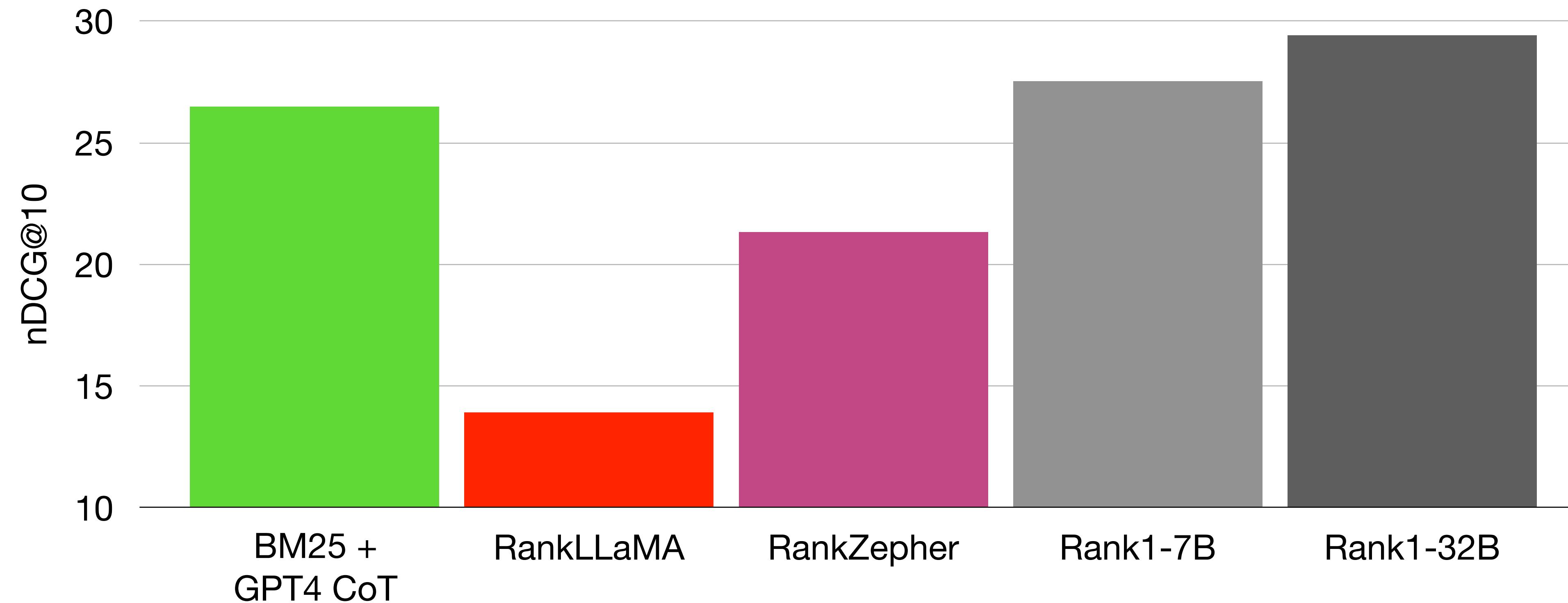
</think> **false**

# Reasoning Rankers

- Rank1: Pointwise reasoning ranker distilled from DeepSeek-R1
  - Positive and negative data points from MSMARCO.
    - Labelled positive docs, hard/easy negatives from retrievers and rerankers
  - Extensively filtered.
    - R1 disagreed on positives and easy negatives.
    - A model trained on first round to filter out disagreed hard negatives.

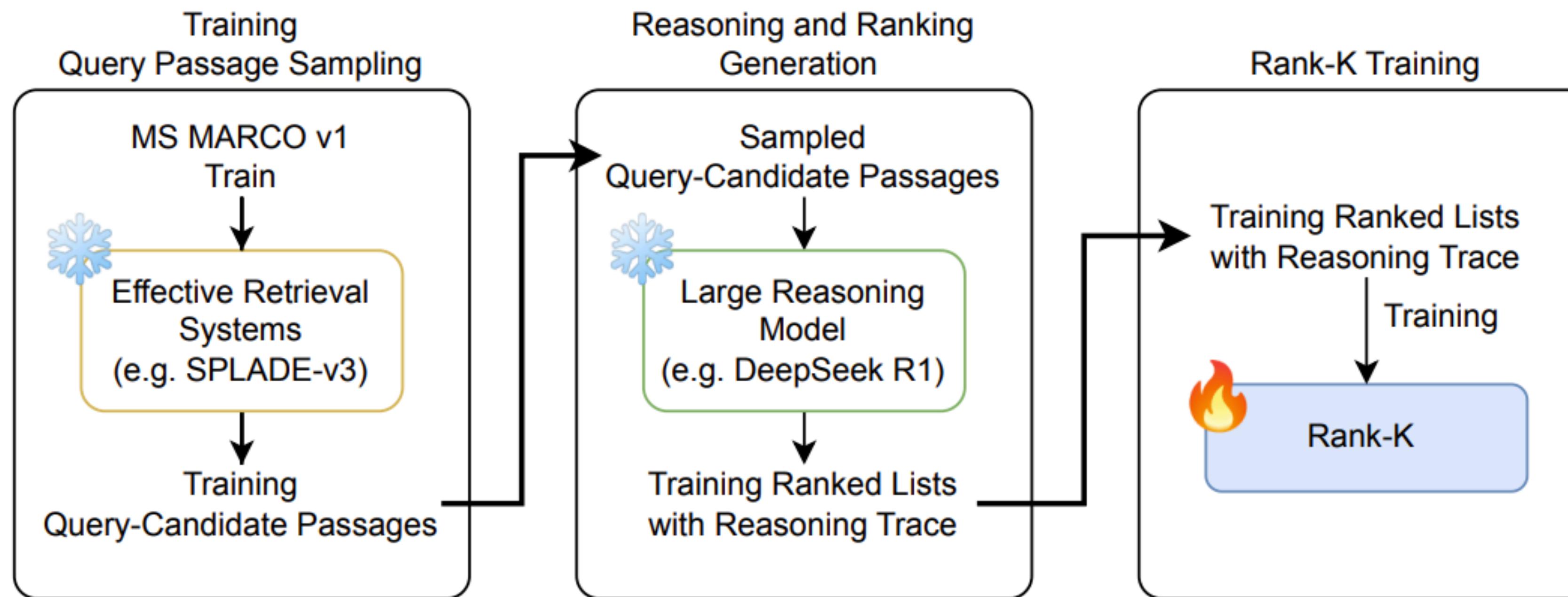
# Reasoning Rankers

- Rank1: BRIGHT results



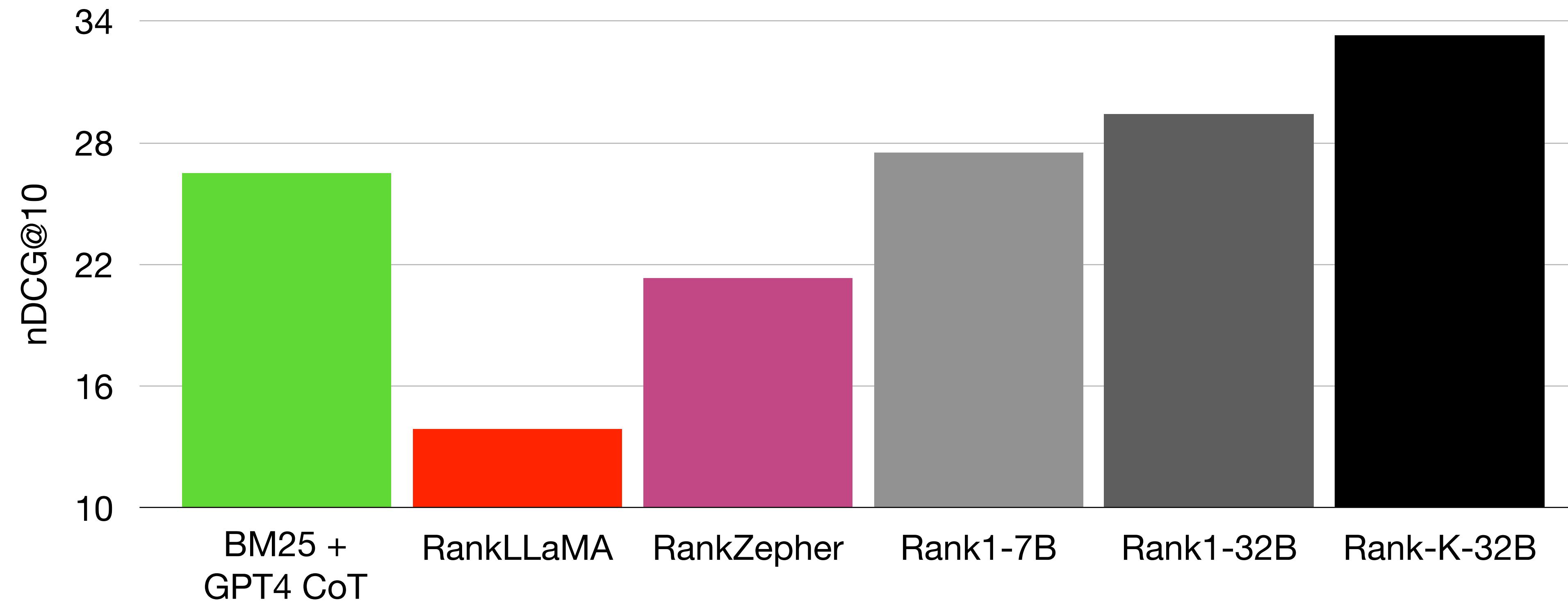
# Reasoning Rankers

- Rank-K: Listwise reasoning ranker distilled from DeepSeek-R1



# Reasoning Rankers

- Rank-K: BRIGHT results



# Reasoning Rankers

- Rank-R1: Reinforcement Learning trained setwise reasoning ranker

Given a query  $\{query\}$ , which of the following passages is more relevant one to the query?

Passage A:  $\{passage\_1\}$

Passage B:  $\{passage\_2\}$

Passage C:  $\{passage\_3\}$

...

Output only the passage label of the most relevant passage:



## SYSTEM:

A conversation between User and Assistant. The user asks a question, and the Assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within `<think> </think>` and `<answer> </answer>` tags, respectively, i.e., `<think>` reasoning process here `</think>` `<answer>` answer here `</answer>`.

## USER:

Given the query: " $\{query\}$ ", which of the following documents is most relevant?

[1]  $\{document1\}$

[2]  $\{document2\}$

....

[20]  $\{document20\}$

After completing the reasoning process, please provide only the label of the most relevant document to the query, enclosed in square brackets, within the answer tags. For example, if the third document is the most relevant, the answer should be: `<think>` reasoning process here `</think>` `<answer>[3]</answer>`.

# Reasoning Rankers

- Rank-R1: Reinforcement Learning trained setwise reasoning ranker

- GRPO RL algorithm.

- MSMARCO training data.

- Reward: 1 if the labelled relevant document is selected and the format matches `<think> </think> <answer> </answer>`, otherwise 0

## SYSTEM:

A conversation between User and Assistant. The user asks a question, and the Assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within `<think> </think>` and `<answer> </answer>` tags, respectively, i.e., `<think>` reasoning process here `</think> <answer>` answer here `</answer>`.

## USER:

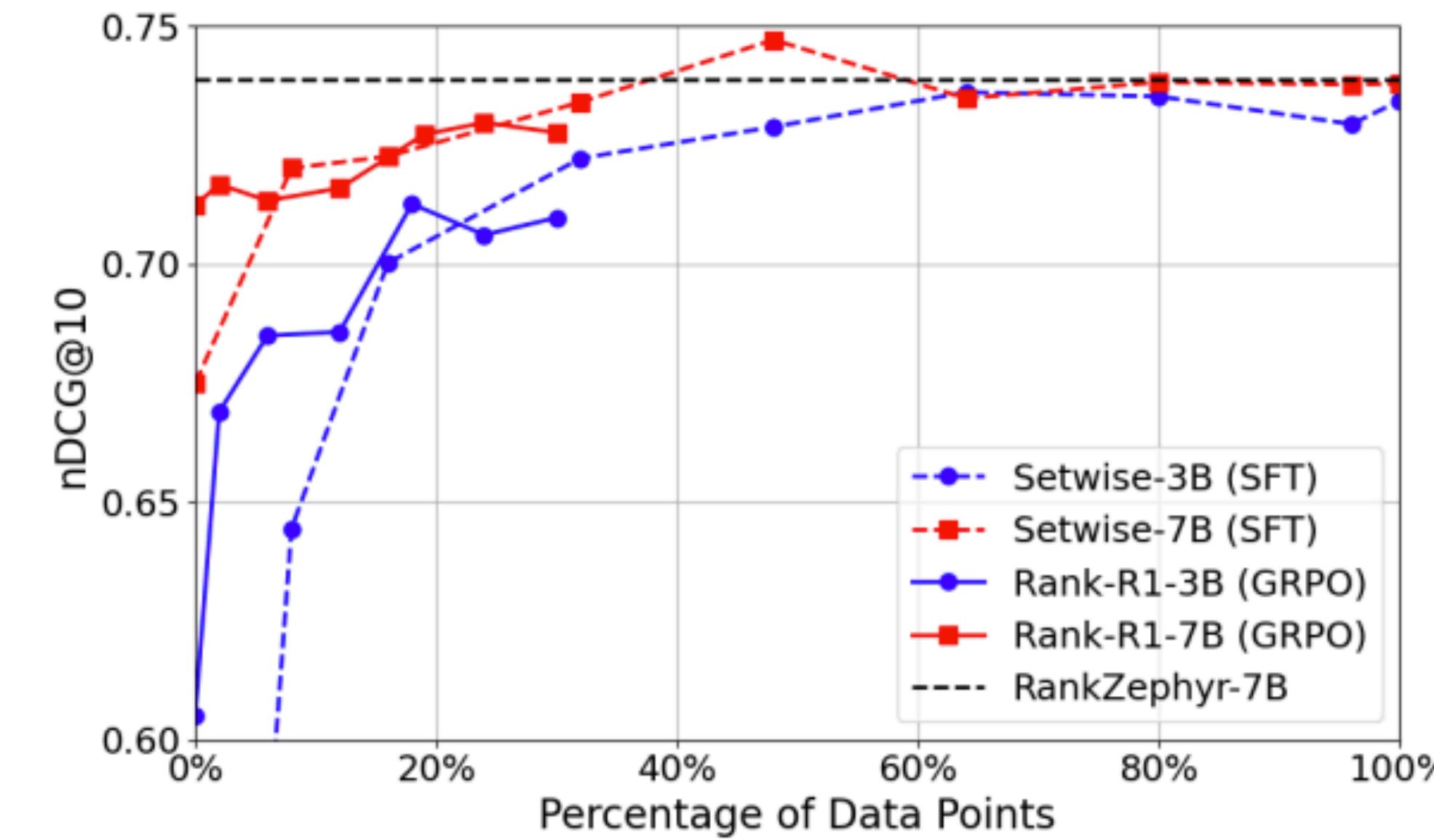
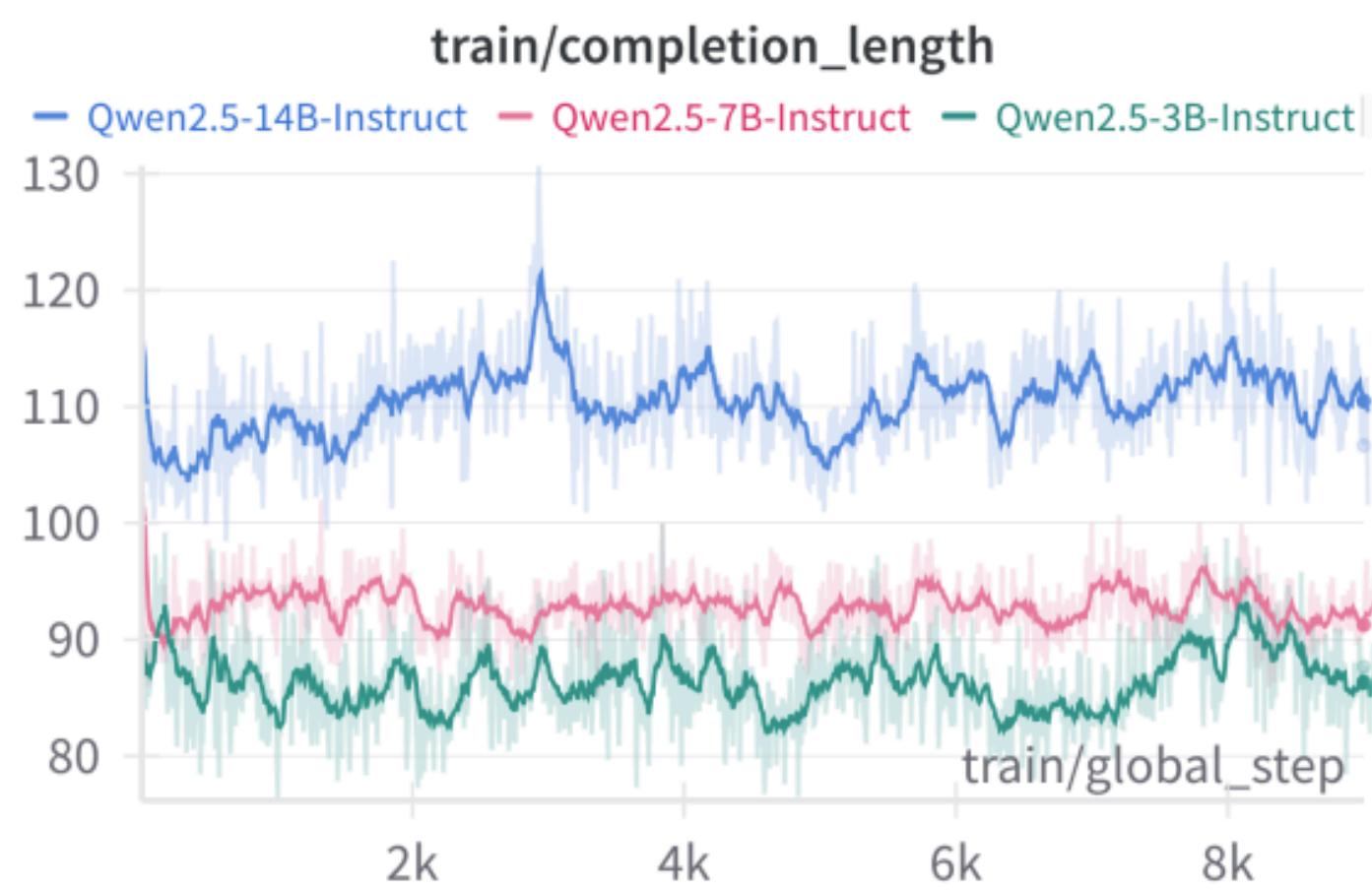
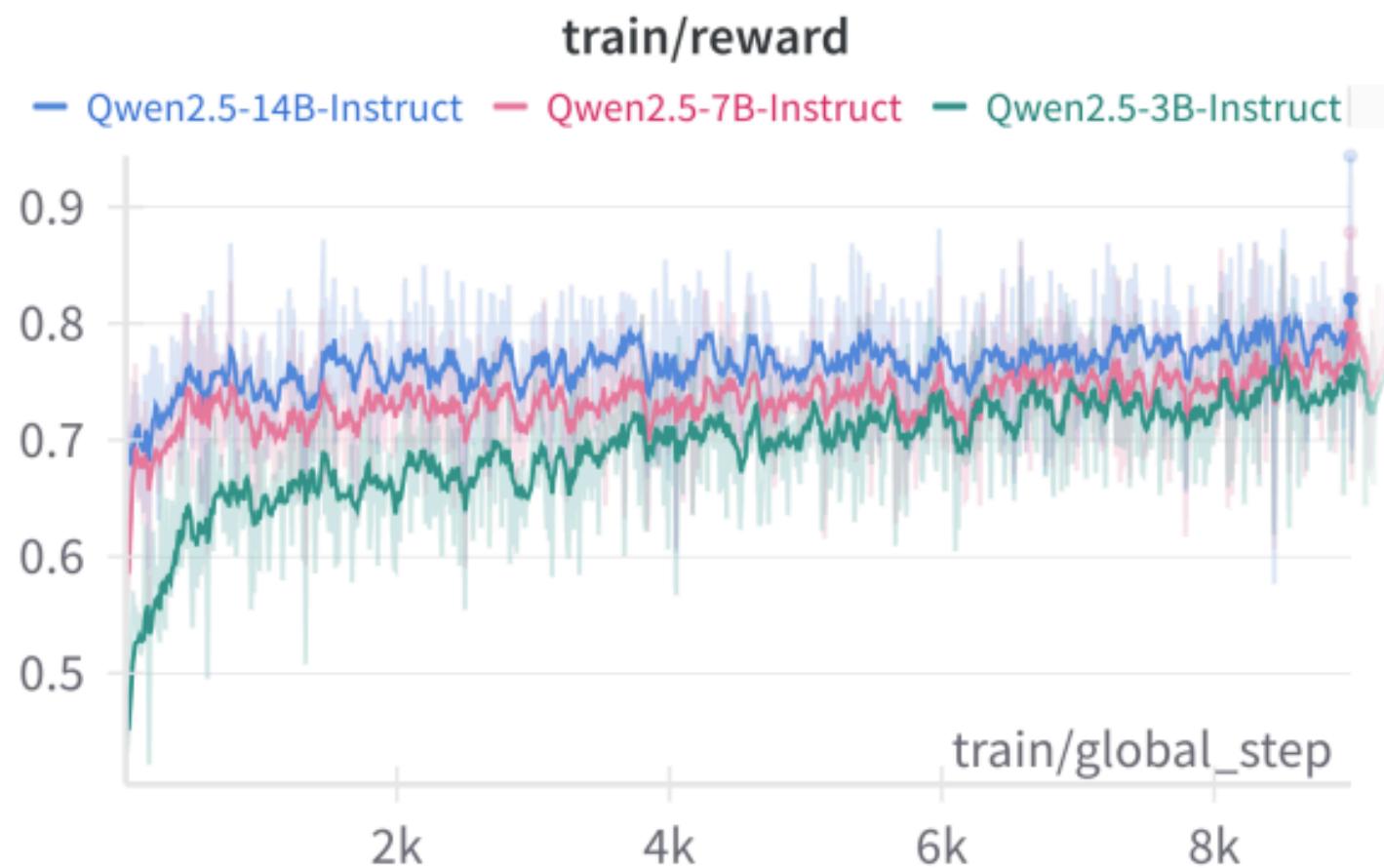
Given the query: "{query}", which of the following documents is most relevant?

- [1] {document1}
- [2] {document2}
- ....
- [20] {document20}

After completing the reasoning process, please provide only the label of the most relevant document to the query, enclosed in square brackets, within the answer tags. For example, if the third document is the most relevant, the answer should be: `<think>` reasoning process here `</think> <answer>[3]</answer>`.

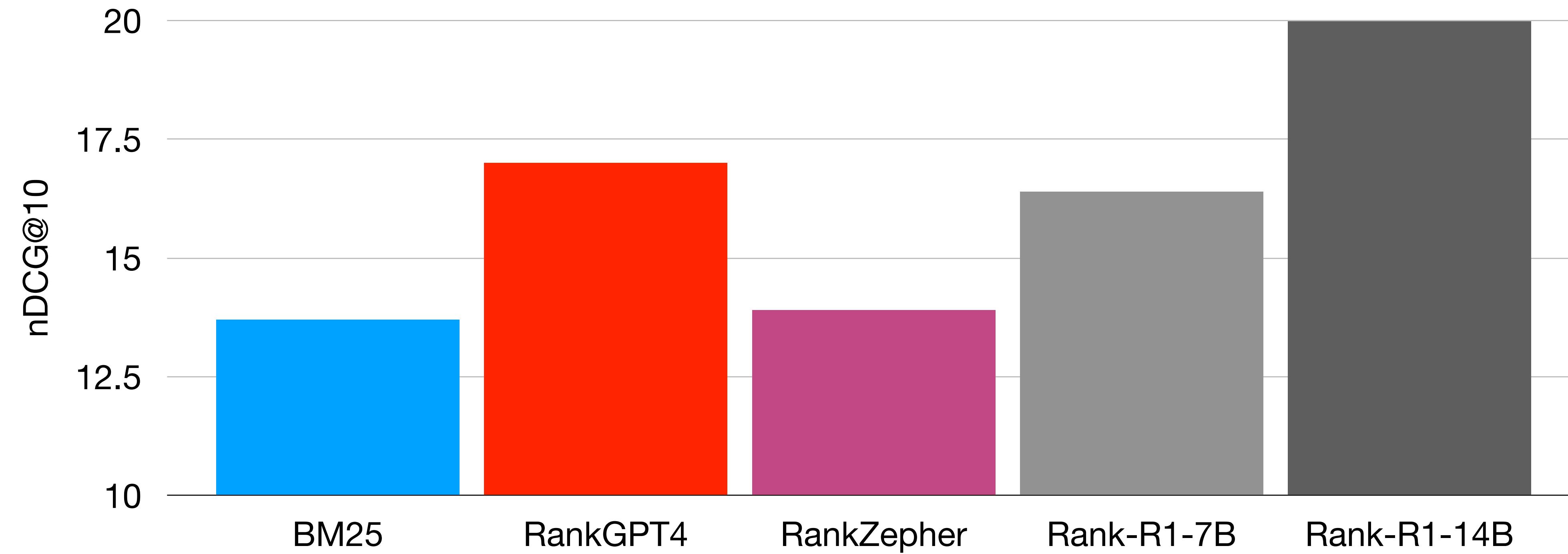
# Reasoning Rankers

- Rank-R1: Reinforcement Learning trained setwise reasoning ranker

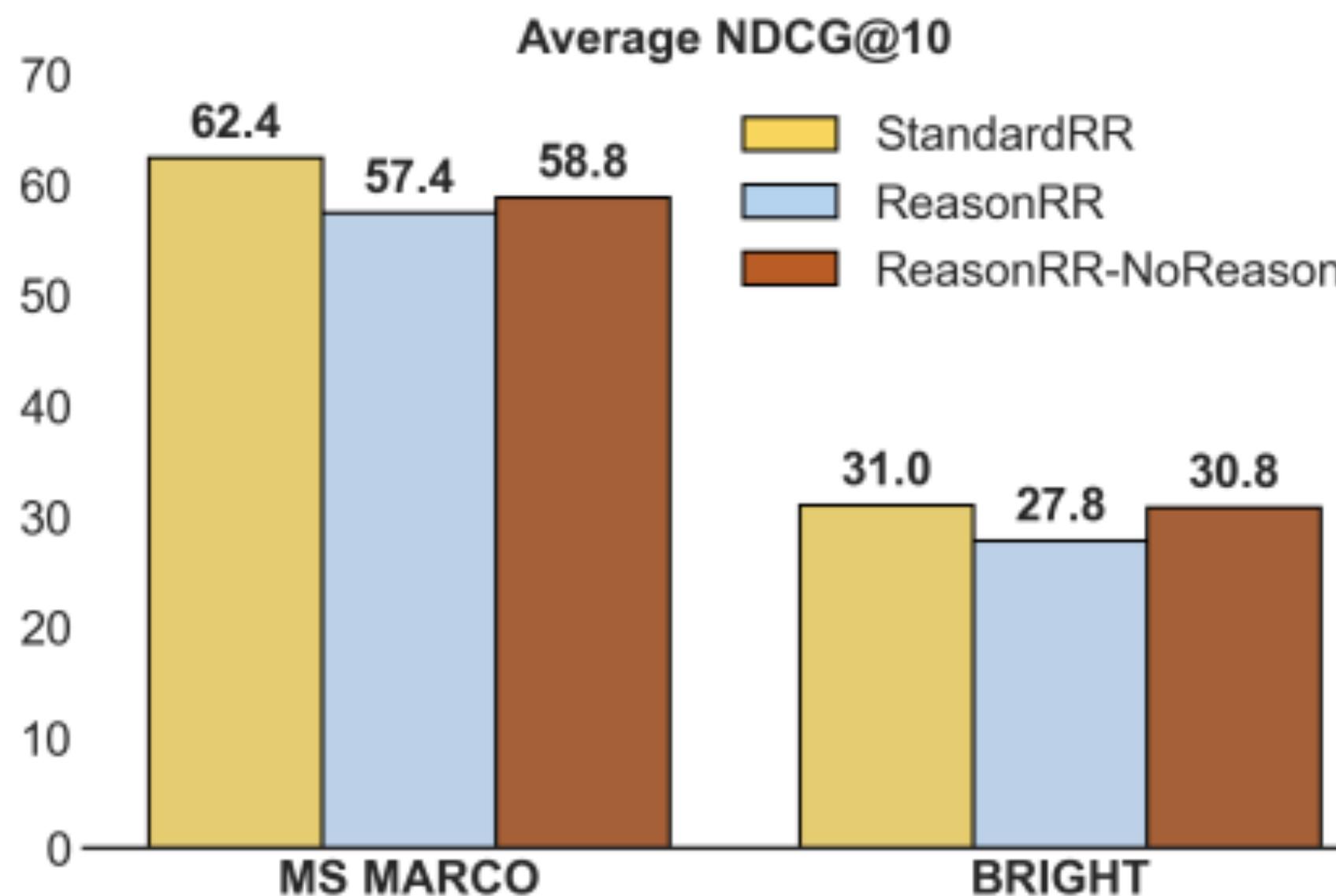


# Reasoning Rankers

- Rank-R1: BRIGHT results



# Is Reasoning Truly Necessary for Ranking?



	StackExchange							Coding		Theorem-based			Avg.
	Bio.	Earth.	Econ.	Psy.	Rob.	Stack.	Sus.	Leet.	Pony	AoPS	TheoQ.	TheoT.	
BM25 + GPT-4 CoT	53.6	54.1	24.3	38.7	18.9	27.7	26.3	19.3	17.6	3.9	19.2	20.8	27.0
+ Qwen2.5-1.5B													
StandardRR	<b>37.0</b>	<b>21.7</b>	<b>16.8</b>	23.1	<b>16.1</b>	10.0	<b>26.3</b>	2.6	<b>30.6</b>	1.8	<b>16.1</b>	<b>26.1</b>	<b>19.0</b>
ReasonRR	32.5	20.3	12.3	<b>25.5</b>	11.1	<b>15.3</b>	23.5	<b>6.6</b>	12.3	<b>3.4</b>	10.6	13.7	15.6
+ Qwen2.5-3B													
StandardRR	<b>41.6</b>	27.1	<b>20.9</b>	31.9	<b>22.2</b>	16.9	<b>30.3</b>	<b>13.2</b>	<b>42.0</b>	2.7	16.2	30.6	<b>24.6</b>
ReasonRR	37.3	<b>27.8</b>	20.7	<b>33.1</b>	18.3	<b>24.3</b>	25.2	11.3	26.2	<b>4.7</b>	<b>20.7</b>	<b>34.0</b>	23.6
+ Qwen2.5-7B													
StandardRR	<b>47.1</b>	<b>38.0</b>	<b>28.1</b>	<b>44.1</b>	<b>26.1</b>	<b>29.5</b>	<b>36.5</b>	<b>19.3</b>	<b>37.5</b>	4.6	<b>22.4</b>	<b>39.4</b>	<b>31.0</b>
ReasonRR	47.0	35.4	24.0	35.2	20.0	25.2	31.0	15.1	36.0	<b>5.9</b>	22.2	36.6	27.8

# Is Reasoning Truly Necessary for Ranking?

- Seems not necessary for simple queries. (Results on MS MARCO and BEIR).
- Ineffective when base model is “small”.
- What about efficiency in reasoning models?!

# In this part we covered

- SFT: From monoBERT to RankLLaMA.
- Zero-shot LLM rankers: Pointwise, Pairwise, Listwise, Setwise.
- Reasoning rankers: distillation and RL

# End of Part 3

...