



## **ListT5: Listwise Reranking with Fusion-in-Decoder Improves Zero-shot Retrieval**

**Soyoung Yoon<sup>1\*</sup> Eunbi Choi<sup>2</sup> Jiyeon Kim<sup>3</sup> Yireun Kim<sup>2</sup>  
Hyeongu Yun<sup>2</sup> Seung-won Hwang<sup>1</sup>**

<sup>1</sup>Seoul National University <sup>2</sup>LG AI Research <sup>3</sup>KAIST AI



**LG AI Research**

**KAIST AI**

# Overview



We introduce ListT5, FiD with tournament sort, that is..

1. Computationally efficient.
  - a. Faster than pairwise or LLM + sliding window based listwise methods
  - b. comparable with pointwise methods
2. Robust to positional bias.
  - a. Overcomes the lost-in-the middle problem by FiD, with each passage encoded with identical positional encoding.
3. Shows great zero-shot performance.
  - a. superior than any listwise, pointwise, pairwise models on BEIR benchmark, for T5-base and T5-3B with relatively small size.

# Background

- Models still struggle on zero-shot retrieval
- Listwise reranking models are shown to be effective on zero-shot retrieval, but previous listwise reranking models had limitations
  - small-sized models only implement pairwise reranking with impractical efficiency (e.g., DuoT5)
  - large-sized models suffer from the lost-in-the middle problem due to its long input length.

# Pointwise v.s. Listwise Reranking

Q. How to make money easily

D1. Applicable online surveys by 2017 to make money ...

Q. How to make money easily

D2. The history of making money. From 1984, farmers ...

Q. How to make money easily

D3. How to make money in 2024? Recent trends show ...

MonoT5  
(pointwise)

0.8

MonoT5  
(pointwise)

0.1

MonoT5  
(pointwise)

0.7

argsort



Pointwise reranking with MonoT5

Listwise reranking with ListT5

Q. How to make money easily

D1. Applicable online surveys by 2017 to make money ...

Q. How to make money easily

D2. The history of making money. From 1984, farmers ...

Q. How to make money easily

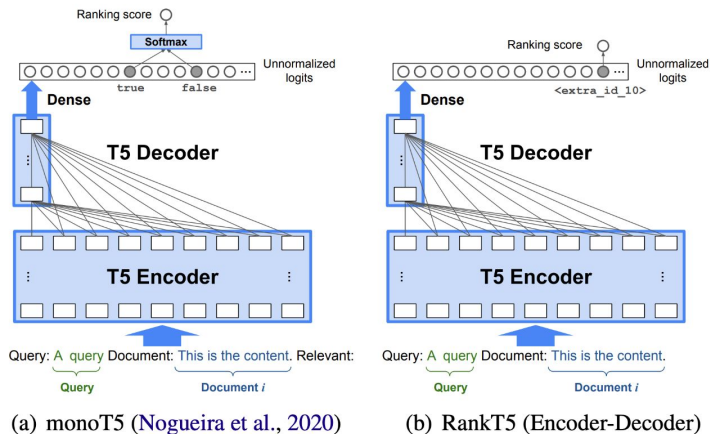
D3. How to make money in 2024? Recent trends show ...

ListT5  
(listwise)

Number  
generation



(pointwise rerankers)

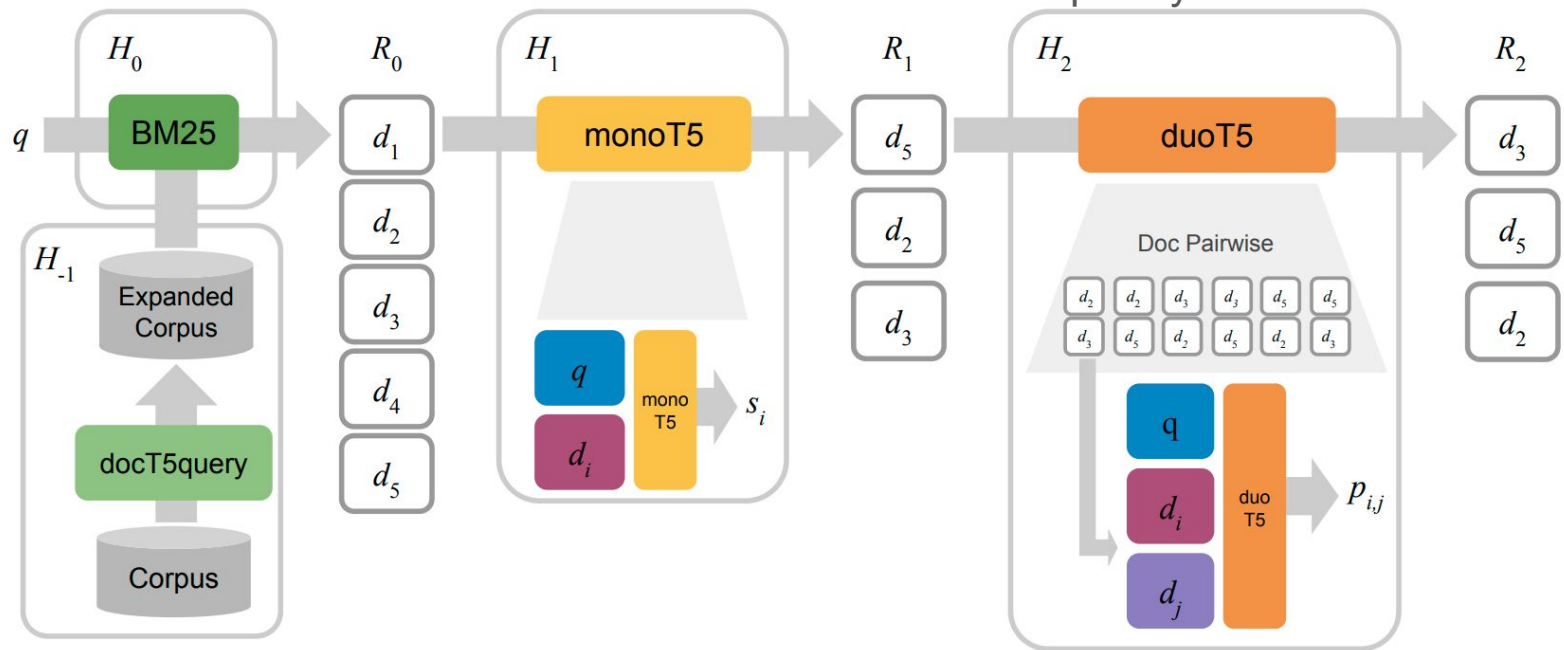


- Listwise rerankers can condition on and compare multiple passages to calibrate the relevance scores better, thus reducing the inaccuracy of predictions arising from domain shift.\*

# Listwise reranking models: Baselines

- Listwise Baseline models: DuoT5 for small models

DuoT5: better performance than pointwise models, but  $n^2$  time complexity!



# Listwise reranking models: Baselines

- **Listwise Baseline models: RankVicuna, RankZephyr, RankGPT for Large Language Models**

USER: I will provide you with {num} passages, each indicated by a numerical identifier []. Rank the passages based on their relevance to the search query: {query}.

```
[1] {passage 1}
[2] {passage 2}
...
[{num}] {passage {num}}
```

Search Query: {query}.

Rank the {num} passages above based on their relevance to the search query. All the passages should be included and listed using identifiers, in descending order of relevance. The output format should be [] > [], e.g., [4] > [2]. Only respond with the ranking results, do not say any word or explain.

Recently, methods to do listwise reranking with LLMs has been investigated

The following are passages related to query {{query}}  
[1] {{passage\_1}}  
[2] {{passage\_2}}  
(more passages)  
Rank these passages based on their relevance to the query.

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

## (c) Permutation generation

But, inefficiency due to large parametric size of the model & **lost-in-the middle** problem occurs!

# Listwise reranking models: Baselines

- Crucial problem in Listwise reranking with LLMs: **Lost in the middle problem**

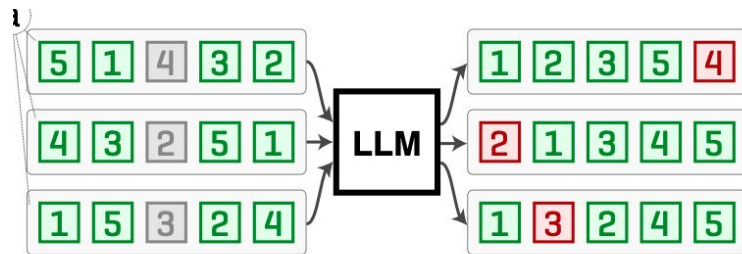
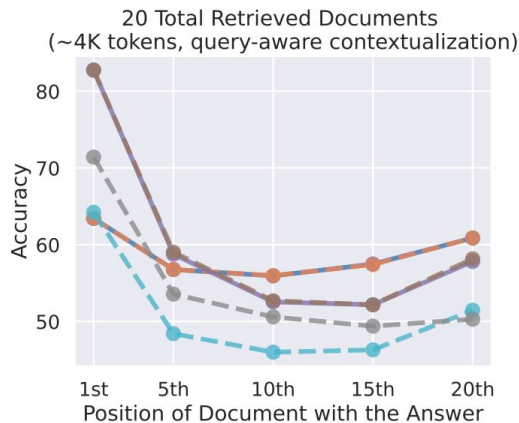
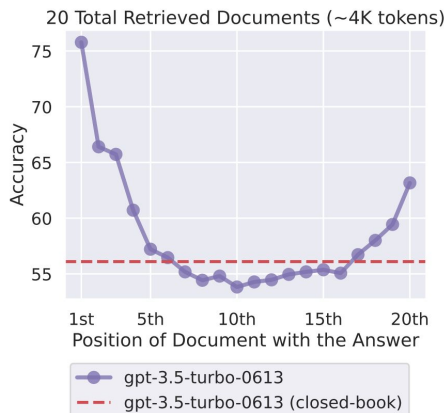


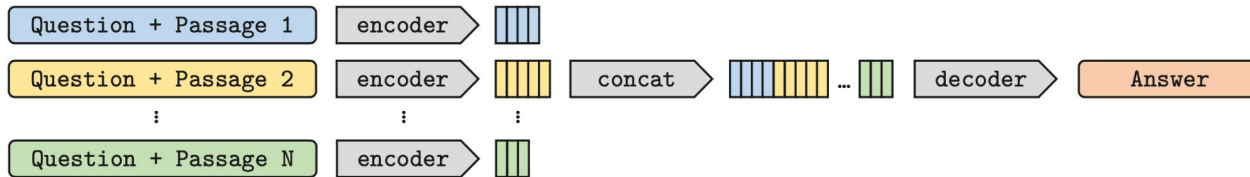
Figure 1: Changing the location of relevant information (in this case, the position of the passage that answers an input question) within the language model’s input context results in a U-shaped performance curve—models are better at using relevant information that occurs at the very beginning (primacy bias) or end of its input context (recency bias), and performance degrades significantly when models must access and use information located in the middle of its input context.



# Listwise reranking: Solutions

- listwise reranking is effective for zero-shot retrieval
- How to better utilize the autoregressive generation ability of reranking models?
- Small models can't see long context, pairwise models are impractical, efficiency hurts with lengthy inputs
- listwise reranking with LLMs has the lost-in-the-middle problem
- How can we train the model to efficiently see **multiple passages at once** and do **listwise ranking**, while being fairly efficient and exhibit less positional bias?

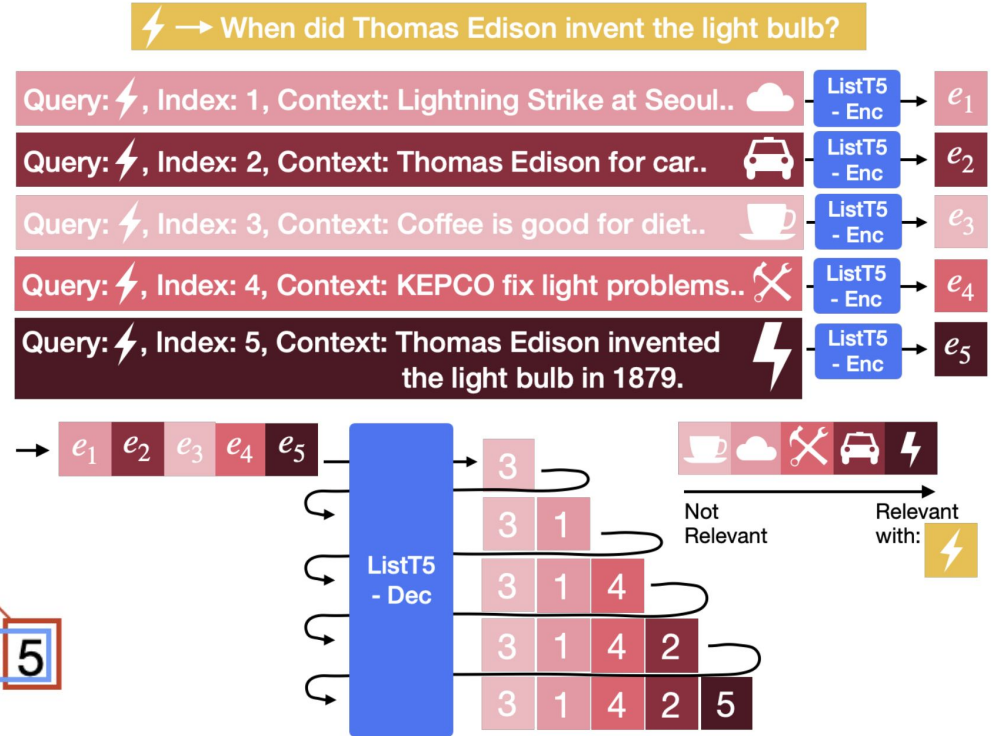
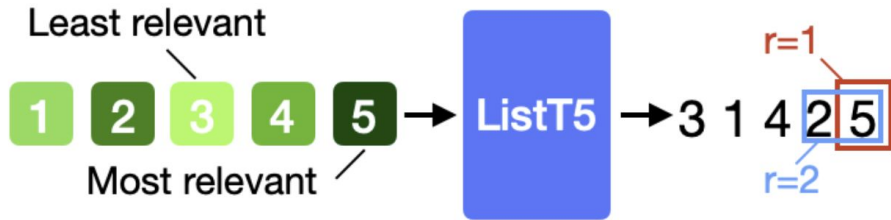
-> use FiD (Fusion-in-Decoder) architecture that outputs sorted passage index, and use (hierarchical) tournament sort to cache already-computed passages!





# Proposed Method: ListT5 architecture

- Fusion-in-Decoder, that given  $k$  ( $=5$ ) contexts, output sorted index, with relevant index coming at the last.
- Each passage is processed by the encoder with identical positional encodings, and the model identifies each passage with index number: CANNOT exploit positional bias.



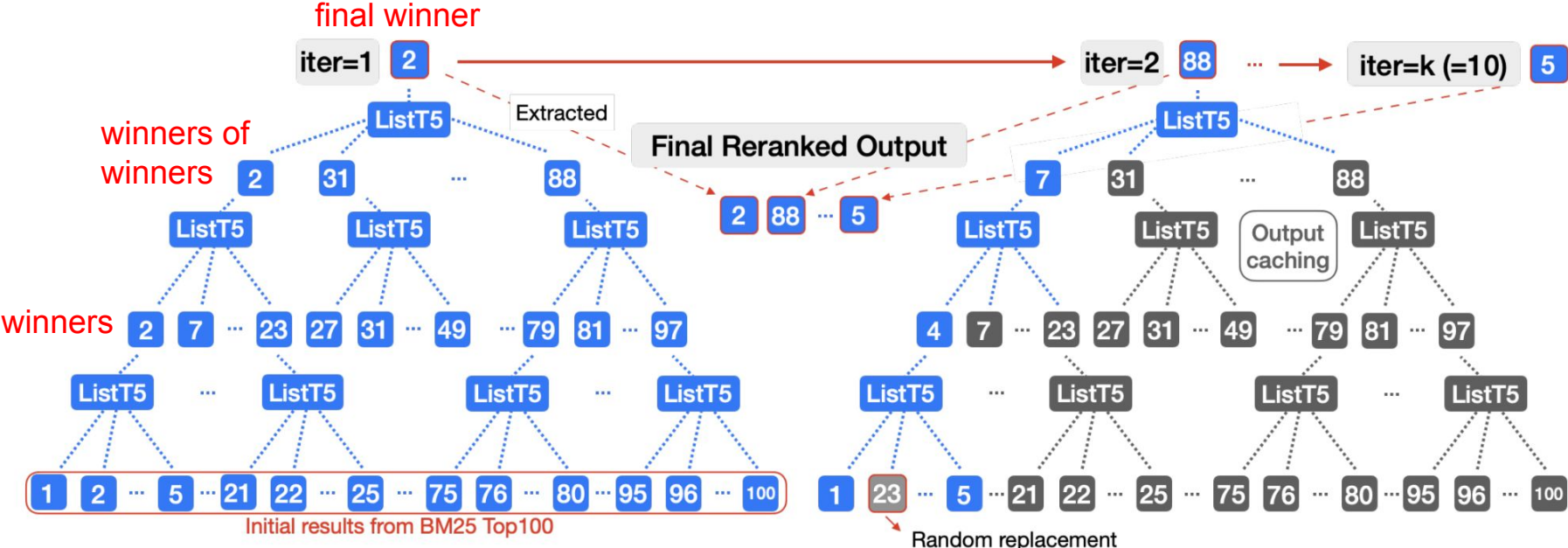
# Making the train dataset

- Source: MSMARCO, with only positive/negative labels
- Labeling ordering between negatives
  - Used bi-encoder (coco-dr large: 340M, GTR-large for ablation), selected top-1000 out of 8.8M corpus, random selection of 4 negatives
  - labeled negative scoring by the dot product scores of the bi-encoder

\*training data open-sourced upon request at  
<https://huggingface.co/datasets/Soyoung97/ListT5-train-data>

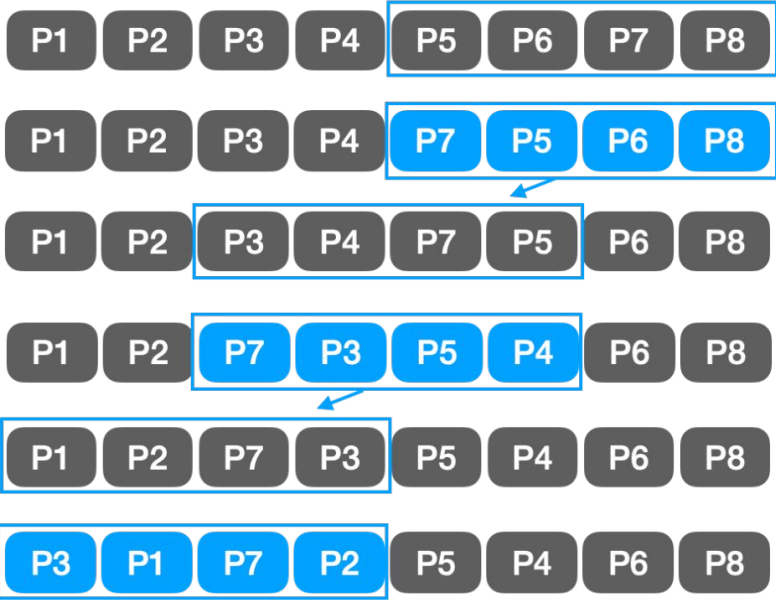
# Sorting method: Tournament sort

- 5-ary tournament tree, with output caching
- group passages to 5, ranks them hierarchically like tournament
- getting next top-1 only requires computing path for changed elements



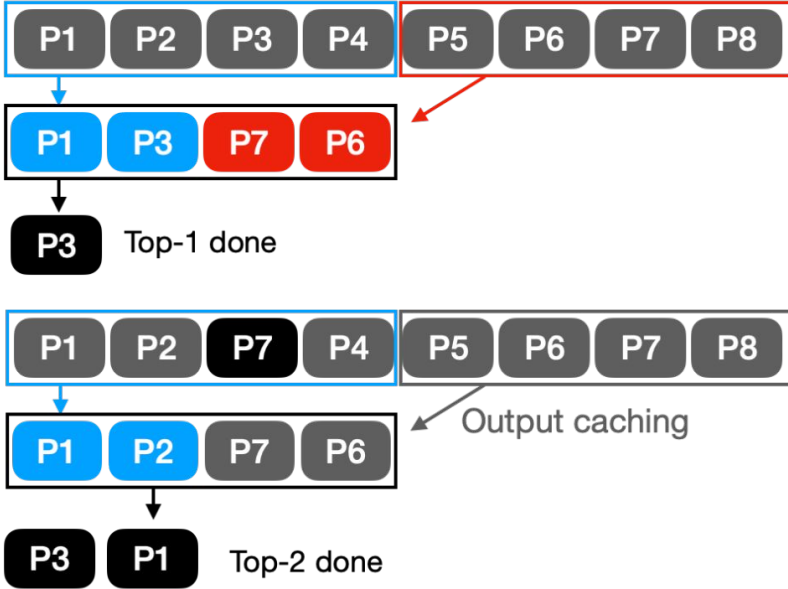
# Tournament sort v.s. Sliding window

window = 4, stride = 2, rerank top-2



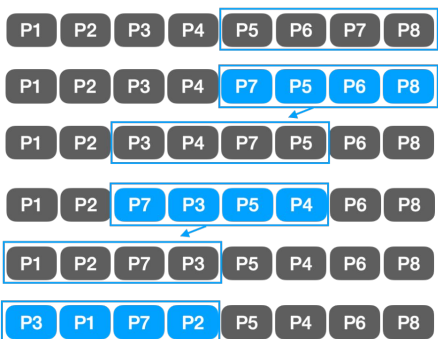
Sliding window

window = 4, r=2, rerank top-2



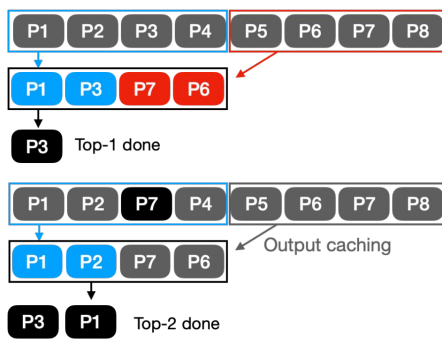
Tournament Sort

window = 4, stride = 2, rerank top-2



Sliding window

window = 4, r=2, rerank top-2



Tournament Sort

Method Name

Reranking Method

Complexity

MonoT5 (Nogueira et al., 2020), RankT5 (Zhuang et al., 2022)

Pointwise

$\mathcal{O}(n)$

DuoT5 (Pradeep et al., 2021)

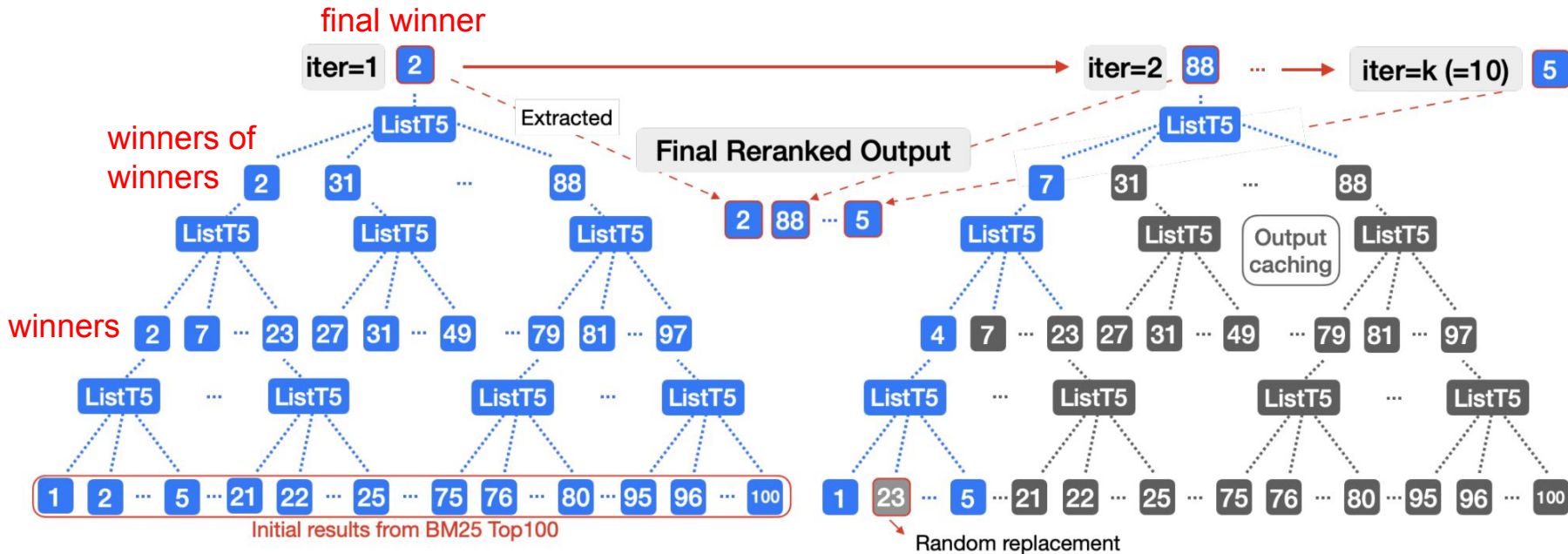
Pairwise

$\mathcal{O}(n^2)$

ListT5 (Ours)

Listwise

$\mathcal{O}(n + k \log n)$



# Zero-shot performance - pointwise baseline models

	BM25 Top-100						BM25 Top-1000			COCO-DR Large (Init.)	MonoT5	RankT5	ListT5 (r=1)	ListT5 (r=2)
	Initial	MonoT5 -base	RankT5 -base	ListT5 -base (r=2)	MonoT5 -3B	RankT5 -3B	ListT5 -3B (r=2)	MonoT5 -base	RankT5 -base					
TREC-COVID	59.5	<b>78.3</b>	77.7	<b>78.3</b>	79.8	81.7	<b>84.7</b>	78.3	79.1	<b>82.1</b>				
NFCorpus	32.2	<b>35.7</b>	35.1	35.6	37.3	<b>37.4</b>	37.7	<b>36.1</b>	35.3	<b>36.1</b>				
BioASQ	52.2	55.3	<b>58.2</b>	56.4	57.5	<b>58.3</b>	<b>58.3</b>	52.6	<b>57.6</b>	55.0				
NQ	30.5	52.1	<b>53.2</b>	53.1	56.4	<b>57.8</b>	56.2	55.9	<b>57.6</b>	57.5				
HotpotQA	63.3	71.2	<b>72.8</b>	72.6	74.3	<b>74.8</b>	<b>75.6</b>	70.9	<b>73.8</b>	73.6				
FiQA-2018	23.6	39.2	39.2	<b>39.6</b>	<b>46.0</b>	45.2	45.1	41.2	41.1	<b>41.8</b>				
Signal-1M (RT)	33.0	32.0	30.8	<b>33.5</b>	32.2	31.9	<b>33.8</b>	29.3	28.6	<b>30.9</b>				
TREC-NEWS	39.5	48.0	45.4	48.5	48.3	49.5	<b>53.2</b>	47.8	45.9	<b>50.9</b>				
Robust04	40.7	53.4	<b>54.3</b>	52.1	<b>58.5</b>	58.3	57.8	55.4	<b>57.2</b>	54.7				
Arguana	40.8	34.4	35.5	<b>48.9</b>	46.8	37.4	<b>50.6</b>	24.2	26.6	<b>46.9</b>				
Touche-2020	44.2	29.6	<b>37.1</b>	33.4	32.5	<b>38.8</b>	33.6	26.4	<b>37.0</b>	31.5				
CQADupStack	30.0	38.6	37.0	<b>38.8</b>	41.3	40.3	<b>42.1</b>	40.1	38.1	<b>40.5</b>				
Quora	78.9	84.6	83.3	<b>86.4</b>	84.0	83.6	<b>86.9</b>	84.2	82.9	<b>86.4</b>				
DBPedia	31.8	42.8	43.7	<b>43.7</b>	44.8	<b>45.0</b>	46.2	43.1	<b>45.1</b>	44.9				
SCIDOCS	14.9	16.7	16.8	<b>17.6</b>	19.0	18.9	<b>19.5</b>	17.0	17.1	<b>18.0</b>				
FEVER	65.2	78.4	77.6	<b>79.8</b>	80.0	79.8	<b>82.0</b>	77.9	77.8	<b>81.0</b>				
Climate-FEVER	16.5	23.1	21.2	<b>24.0</b>	<b>26.2</b>	24.5	24.8	23.3	20.6	<b>24.9</b>				
SciFact	67.9	73.1	73.5	<b>74.1</b>	76.3	<b>77.1</b>	77.0	73.3	73.6	<b>74.9</b>				
Average	42.5	49.3	49.6	<b>50.9</b>	52.3	52.2	<b>53.6</b>	48.7	49.7	<b>51.8</b>				

	COCO-DR Large (Init.)	MonoT5	RankT5	ListT5 (r=1)	ListT5 (r=2)
MSMARCO Top-1000 (in-domain)	41.9	43.1	46.2	46.1	<b>46.3</b>
TREC-COVID	80.8	<b>83.5</b>	<b>83.5</b>	83.2	<b>83.5</b>
NFCorpus	35.5	35.6	35.5	<b>36.2</b>	<b>36.2</b>
NQ	54.3	57.9	59.6	59.7	<b>60.0</b>
HotpotQA	63.3	68.7	<b>71.1</b>	70.3	70.9
FiQA-2018	32.3	41.2	41.3	<b>41.7</b>	<b>41.7</b>
Arguana	46.9	33.0	34.8	49.0	<b>49.3</b>
Touche-2020	21.6	25.7	<b>35.7</b>	29.1	29.6
CQADupStack	37.3	40.5	38.7	40.7	<b>40.9</b>
Quora	<b>87.3</b>	84.0	83.0	86.2	86.3
DBPedia	40.7	44.4	<b>46.1</b>	45.6	45.4
SCIDOCS	17.3	17.5	17.5	17.7	<b>18.3</b>
FEVER	74.9	78.9	79.7	79.8	<b>81.4</b>
Climate-FEVER	23.1	24.2	22.9	23.9	<b>24.9</b>
SciFact	71.9	73.5	73.6	<b>74.4</b>	74.3
Avg. BEIR	49.1	50.6	51.6	52.7	<b>53.1</b>

# Zero-shot performance - listwise baseline models

	TREC-DL19	TREC-DL20	TREC-COVID	NFC-orpus	Signal-1M (RT)	TREC-NEWS	Robust 04	Touche-2020	DBP-edia	Sci-Fact	Avg (In-domain)	Avg (BeIR)
DuoT5-base	71.4	67.4	<b>80.1</b>	35.0	31.4	<b>49.1</b>	49.6	31.8	<b>43.9</b>	69.6	69.4	48.8
LISTT5-base ( $r = 2$ )	<b>71.8</b>	<b>68.1</b>	78.3	<b>35.6</b>	<b>33.5</b>	48.5	<b>52.1</b>	<b>33.4</b>	43.7	<b>74.1</b>	<b>70.0</b>	<b>49.9</b>
RankGPT (GPT3.5)	65.8	62.9	76.7	35.6	32.1	48.9	50.6	<b>36.2</b>	44.5	70.4	64.4	49.4
RankVicuna-7b	68.9	66.1	80.5	33.2	<b>34.2</b>	46.9	48.9	33.0	44.4	70.8	67.5	49.0
RankZephyr-7b	<b>73.9</b>	<b>70.9</b>	84.0	36.7	31.8	52.6	54.3	33.8	44.6	74.9	<b>72.4</b>	51.6
LISTT5-3B ( $r = 2$ )	71.8	69.1	<b>84.7</b>	<b>37.7</b>	33.8	<b>53.2</b>	<b>57.8</b>	33.6	<b>46.2</b>	<b>77.0</b>	70.5	<b>53.0</b>

# Efficiency

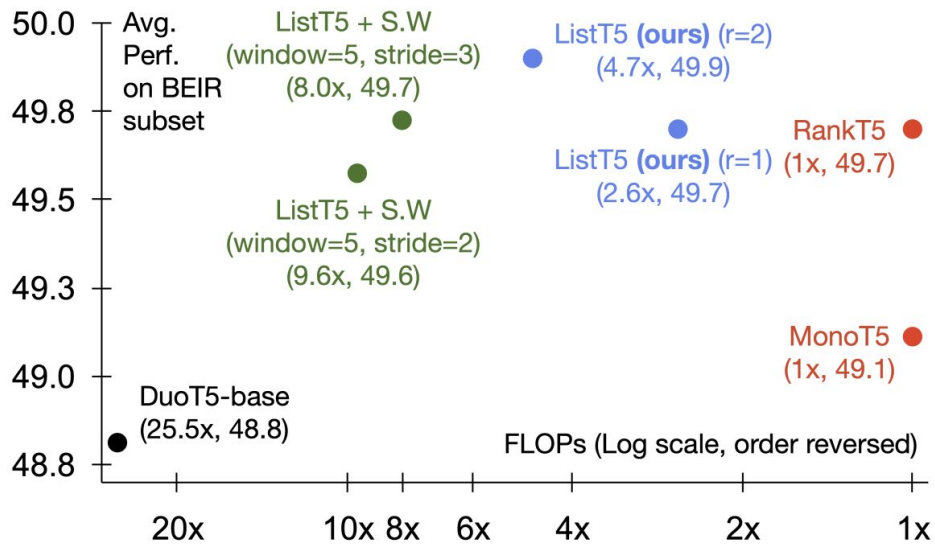


Figure 4: Real-time FLOPs comparison of the models on T5-base, including DuoT5 and the sliding window variants of LISTT5. The reported BEIR performance is averaged from a subset of BEIR, same as in Tab. 3.



# Positional Invariance

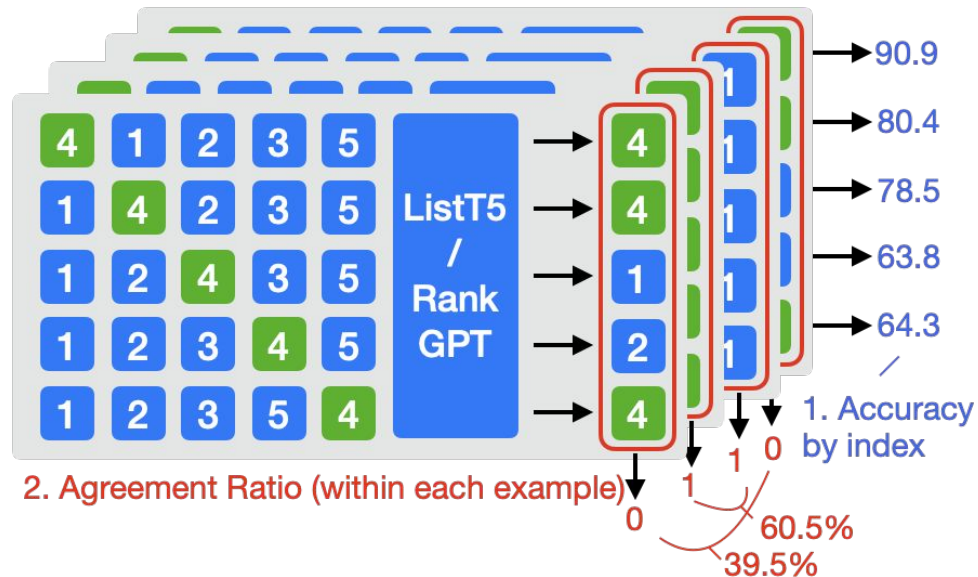
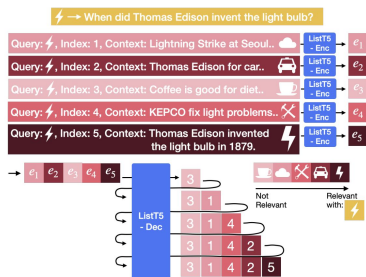
I will provide you with 5 passages, each indicated by numerical identifier []. Rank the passages based on their relevance to the search query: {query}.

- [1] {passage\_1}
- [2] {passage\_2}
- [3] {passage\_3}
- [4] {passage\_4}
- [5] {passage\_5}

V.S.

Search Query: {query}

Rank the 5 passages above based on their relevance to the search query. All the passages should be included and listed using identifiers, in descending order of relevance. The output format should be [] > [], e.g., [4] > [2]. Only respond with the ranking results, do not say any word or explain.



# Positional Invariance - better than RankGPT4!

	TREC-COVID							Aggrement ratio (↑)	FiQA							Aggrement ratio (↑)
	Accuracy when positive passage is at index # :						Std. (↓)		Accuracy when positive passage is at index #:						Std. (↓)	
	1	2	3	4	5				1	2	3	4	5			
GPT-3.5	81.6	63.3	75.5	67.3	61.2	7.7	55.1	88.3	68.1	78.7	65.9	75.8	8.0	62.1		
GPT-4	95.9	83.7	73.5	77.6	71.4	8.8	69.4	94.6	90.5	84.4	86.8	84.8	3.9	82.8		
DuoT5	91.3	76.0	-	-	-	7.6	79.6	89.9	76.9	-	-	-	6.5	78.1		
LISTT5	93.9	87.8	83.7	85.7	81.6	<b>4.2</b>	<b>83.7</b>	85.3	85.6	82.2	83.3	82.6	<b>1.4</b>	<b>90.4</b>		

Table 4: Robustness to the position of the positive passage in the input, on TREC-COVID and FiQA. GPT-3.5, GPT-4, DuoT5, and LISTT5 stands for GPT-3.5-turbo-1106, GPT-4-0613, DuoT5-base, and LISTT5-base, respectively. Using the FiD structure effectively mitigates the problem of the positional bias of positive passages, showing lowest standard deviation and highest agreement ratio.

# Positional Invariance - robustness to shuffling candidate passages

Initial ordering	DL19	DL20	TREC-COVID	TREC-NEWS	Touche-2020	Avg.
<b>ListT5-base (tournament sort, r=2)</b>						
No shuffle	71.8	68.1	78.3	48.5	33.4	60.0
Shuffle	71.2	68.1	77.2	48.9	32.8	59.6
Perf. drop						<b>-0.4</b>
ListT5-base (sliding windows, stride=3, iter=4)						
No shuffle	71.8	67.7	77.5	50.0	33.1	60.0
Shuffle	69.5	65.5	77.7	49.2	32.1	58.8
Perf. drop						<b>-1.2</b>
RankVicuna-7b (sliding windows)						
No shuffle	68.9	66.1	80.5	46.9	33.0	59.1
Shuffle	67.1	64.6	79.2	45.3	30.8	57.4
Perf. drop						<b>-1.7</b>
RankGPT-3.5 (sliding windows)						
No shuffle	68.4	64.9	72.6	46.5	38.2	58.1
Shuffle	62.5	57.0	66.1	38.3	22.8	49.3
Perf. drop						<b>-8.8</b>

# Summary

We introduce ListT5 with tournament sort, that is..




Paper



Code

1. Computationally efficient.
  - a. Lower than pairwise or LLM + sliding window based listwise methods
  - b. comparable with pointwise methods
2. Robust to positional bias.
  - a. Overcomes the lost-in-the middle problem by FiD, with each passage encoded with identical positional encoding.
3. Shows great zero-shot performance.
  - a. superior than any listwise, pointwise, pairwise models on BEIR benchmark, for T5-base and T5-3B with relatively small size

 Soyoung97/ListT5-base

Updated Dec 21, 2023 •  140

 Soyoung97/ListT5-3b

Updated Dec 21, 2023 •  11

Please contact:

[soyoung.yoon@snu.ac.kr](mailto:soyoung.yoon@snu.ac.kr)

for any questions!

# [Ablation] Model variants



- Most effective to generate most relevant index at the LAST!
- Sequential generation like reasoning chain!

Dataset	Relevant Discrimination	Relevant First		Relevant Last (ListT5)	
		( $r = 1$ )	( $r = 2$ )	( $r = 1$ )	( $r = 2$ )
In-domain					
MS MARCO	40.3	40.8	<b>40.9</b>	40.7	40.7
TREC-DL19	<b>72.5</b>	69.6	70.8	71.2	71.8
TREC-DL20	67.3	67.0	66.8	67.3	<b>68.1</b>
<b>Avg (in-domain)</b>	60.0	59.1	59.5	59.7	<b>60.2</b>
Out-domain (BEIR)					
TREC-COVID	74.0	74.9	75.9	76.7	<b>78.3</b>
NFCorpus	34.8	35.5	35.6	35.5	<b>35.6</b>
BioASQ	55.8	56.6	56.6	<b>57.2</b>	56.4
NQ	51.1	52.7	52.9	52.0	<b>53.1</b>
HotpotQA	70.9	72.5	72.6	72.1	<b>72.6</b>
FiQA-2018	38.1	39.3	39.0	39.5	<b>39.6</b>
Signal-1M (RT)	32.9	31.8	31.7	33.3	<b>33.5</b>
TREC-NEWS	43.9	46.6	47.3	47.9	<b>48.5</b>
Robust04	49.8	<b>52.3</b>	<b>52.3</b>	52.0	52.1
Arguana	26.1	32.8	34.6	<b>49.7</b>	48.9
Touche-2020	<b>34.2</b>	31.5	31.3	<b>34.2</b>	33.4
CQADupStack	<b>38.8</b>	38.3	38.4	38.4	<b>38.8</b>
Quora	81.9	84.4	84.8	86.1	<b>86.4</b>
DBPedia	42.4	43.4	43.6	<b>43.9</b>	43.7
SCIDOCS	16.3	17.3	17.3	17.2	<b>17.6</b>
FEVER	77.6	77.4	77.7	77.8	<b>79.8</b>
Climate-FEVER	20.7	22.8	23.0	22.8	<b>24.0</b>
SciFact	73.0	74.1	<b>74.2</b>	74.1	74.1
<b>Avg (BEIR)</b>	47.9	49.1	49.4	50.6	<b>50.9</b>

# Appendix - design choice

	$m = 5$		$m = 10$	
	1/5 (ListT5, r=1)	2/5 (ListT5, r=2)	1/10 (r = 1)	4/10 (r = 4)
MS MARCO	<b>40.7</b>	<b>40.7</b>	40.5	<b>40.7</b>
+ Top-1000	44.7	44.9	44.6	<b>45.0</b>
TREC-DL19	71.2	<b>71.8</b>	70.1	70.5
TREC-DL20	67.3	<b>68.1</b>	66.9	67.2
TREC-COVID	76.7	<b>78.3</b>	76.2	77.9
NFCorpus	35.5	35.6	36.2	<b>36.6</b>
BioASQ	<b>57.2</b>	56.4	55.4	56.4
NQ	52.0	<b>53.1</b>	51.5	52.5
HotpotQA	72.1	<b>72.6</b>	71.4	71.9
FiQA-2018	39.5	<b>39.6</b>	39.0	38.9
Signal-1M (RT)	33.3	<b>33.5</b>	31.7	32.0
TREC-NEWS	47.9	<b>48.5</b>	47.1	47.8
Robust04	52.0	52.1	52.2	<b>53.1</b>
Arguana	<b>49.7</b>	48.9	38.6	46.6
Touche-2020	<b>34.2</b>	33.4	32.4	32.7
CQADupStack	38.4	<b>38.8</b>	38.2	28.8
Quora	86.1	86.4	85.5	<b>86.8</b>
DBPedia	43.9	<b>43.7</b>	42.7	43.6
SCIDOCS	17.2	17.6	17.2	<b>18.0</b>
FEVER	77.8	<b>79.8</b>	76.7	79.1
Climate-FEVER	22.8	<b>24.0</b>	22.7	23.9
SciFact	74.1	74.1	73.4	<b>74.2</b>
Avg(In-domain)	56.0	<b>56.4</b>	55.5	55.9
Avg(BeIR)	50.6	<b>50.9</b>	49.3	50.0

# Appendix - LLM consistency

	GPT-3.5-turbo-1106				ListT5
	Trial 1	Trial 2	Trial 3	Avg.	-base
(1) Accuracy when the gold passage is at index #:					
1	81.6	79.6	81.6	<b>81.0</b>	<b>93.9</b>
2	63.3	63.3	61.2	<b>62.6</b>	<b>87.8</b>
3	75.5	75.5	75.5	<b>75.5</b>	<b>83.7</b>
4	67.3	63.3	67.3	<b>66.0</b>	<b>85.7</b>
5	61.2	63.3	65.3	<b>63.3</b>	<b>81.6</b>
std	7.68	7.1	7.4	<b>7.4</b>	<b>4.2</b>
(2) Agreement ratio (%) within index change of positive					
points to same passage	55.1	55.1	55.1	<b>55.1</b>	<b>83.7</b>
other	44.9	44.9	44.9	44.9	16.3

Table 10: Measuring the LLM consistency on TREC-COVID.

# Appendix - Training Dataset

Model	RankT5	ListT5			
		Training data	GTR	COCO-DR	GTR
Learning Rate	1.00E-04	1.00E-04	1.00E-04	1.00E-04	1.00E-05
Steps	-	20k	20k	10k	30k
TREC-COVID	77.7	78.3	77.3	78.6	77.9
NFCorpus	35.1	35.6	35.4	36.2	35.9
BioASQ	58.2	56.4	54.9	55.1	56.8
NQ	53.2	53.1	52.7	52.8	53.2
HotpotQA	72.8	72.6	72.1	71.9	72.1
FiQA-2018	39.2	39.6	39.1	39.9	39.4
Signal-1M (RT)	30.8	33.5	34.1	32.9	30.9
TREC-NEWS	45.4	48.5	47.6	48.0	48.3
Robust04	54.3	52.1	52.9	52.7	53.6
Arguana	35.5	48.9	43.3	43.6	43.7
Touche-2020	37.1	33.4	31.5	32.7	32.5
CQADupStack	37.0	38.8	38.6	38.5	38.8
Quora	83.3	86.4	86.0	83.9	84.4
DBPedia	43.7	43.7	43.8	43.6	44.5
SCIDOCS	16.8	17.6	17.1	17.1	17.6
FEVER	77.6	79.8	78.9	79.4	79.7
Climate-FEVER	21.2	24.0	24.0	24.8	24.6
SciFact	73.5	74.1	74.0	73.1	74.4
Avg.	49.6	<b>50.9</b>	50.2	50.3	<b>50.5</b>

Table 7: Comparison of ListT5 models trained with GTR and COCO-DR. For ListT5, the reported scores are evaluated using the (r=2) variant.



# Appendix : sliding window v.s. tournament sort

Sorting method	# required forward passes to rerank top1	# required forward passes to rerank top10
sliding window, stride=1	$1 + \lceil (100-5)/1 \rceil = 96$	$96 \times \lceil 10/4 \rceil = 288$
sliding window, stride=2	$1 + \lceil (100-5)/2 \rceil = 49$	$49 \times \lceil 10/3 \rceil = 196$
sliding window, stride=3	$1 + \lceil (100-5)/3 \rceil = 33$	$33 \times \lceil 10/2 \rceil = 165$
sliding window, stride=4	$1 + \lceil (100-5)/4 \rceil = \mathbf{25}$	$25 \times \lceil 10/1 \rceil = 250$
tournament sort, r=1	$(100/5) + (20/5) + 1 = \mathbf{25}$	$25 + 9 \times (1+1+1) = \mathbf{52}$
tournament sort, r=2	$(100/5) + (40/5) + 2 + 1 = 31$	$31 + 9 \times (1+1+1+1) = 67$

Table 12: Number of forward passes to rerank top-k candidates from 100 candidate passages per one query, where window size  $w=5$ . In the case of reranking top-10 passages, tournament sort requires much more fewer number of forward passes.

Idx	Base Model	Sorting method	Name	FLOPs to rerank:	
				Top-1	Top-10
0	T5-base	pointwise	MonoT5	1x	1x
1	T5-base	tournament	ListT5(r=1)	1.3x	2.6x
2	T5-base	tournament	ListT5(r=2)	1.8x	4.7x
3	T5-base	sliding w.(s=2)	T5(FiD)	2.5x	9.8x
4	T5-base	sliding w.(s=3)	T5(FiD)	1.7x	12.3x
5	T5-3b	tournament	ListT5(r=1)	17.6x	36.3x
6	T5-3b	tournament	ListT5(r=2)	24.6x	66.0x
7	T5-3b	sliding w.(s=2)	T5(FiD)	38.5x	154x
8	T5-3b	sliding w.(s=2)	T5(no FiD)	53.8x	215.1x
9	T5-3b	sliding w.(s=3)	T5(FiD)	25.6x	128x
10	T5-3b	sliding w.(s=3)	T5(no FiD)	35.1x	175.6x

Table 13: FLOPs (In a multiple of FLOPs of MonoT5-base) on the choice of architecture and method, on TREC-DL19. For the sliding window approach, we would need a total of 4 multiple passes for stride = 3 and 5 passes for stride = 2 (Explained at Tab. 12) to rerank Top-10 candidates.

# Appendix : sliding window v.s. tournament sort

	Rerank Top-10 (NDCG@10)				Rerank Top-1 (NDCG@1)			
Sorting Method	T.S.		S.W		T.S.		S.W.	
Hyperparam.	r=1	r=2	s=2 (iter=5)	s=3 (iter=4)	r=1	r=2	s=2 (iter=1)	s=3 (iter=1)
FLOPS(DL19)	<b>1x</b>	1.8x	3.7x	3.1x	<b>1x</b>	1.4x	1.96x	1.32x
DL19	71.2	71.8	71.5	71.8	81.0	79.1	81.0	78.7
DL20	67.3	68.1	67.3	67.7	77.8	77.8	79.0	79.6
In-domain avg.	69.3	<b>70.0</b>	69.4	69.8	79.4	78.5	<b>80.0</b>	79.2
TREC-COVID	76.7	78.3	78.9	77.5	88.0	91.0	88.0	86.0
NFCorpus	35.5	35.6	35.3	35.5	47.8	49.2	48.6	48.6
BioASQ	57.2	56.4	54.5	54.9	59.2	58.4	55.8	57.2
NQ	52.0	53.1	52.7	52.8	36.0	37.6	36.4	36.6
HotpotQA	72.1	72.6	71.2	71.6	83.3	84.1	83.1	83.1
FiQA-2018	39.5	39.6	39.7	39.8	41.2	40.7	41.4	41.5
Signal-1M (RT)	33.3	33.5	32.4	33.2	43.3	41.8	42.3	41.8
TREC-NEWS	47.9	48.5	49.8	50.0	53.2	54.1	52.3	52.9
Robust04	52.0	52.1	51.3	51.7	65.1	66.3	67.1	65.9
Arguana	49.7	48.9	47.7	47.8	25.8	23.9	23.3	22.7
Touche-2020	34.2	33.4	32.7	33.1	34.7	31.6	36.7	36.7
CQADupStack	38.4	38.8	38.9	38.8	31.6	31.9	32.1	32.0
Quora	86.1	86.4	86.3	86.2	77.8	77.8	78.1	77.7
DBPedia	43.9	43.7	42.6	43.2	55.5	56.5	55.1	56.6
SCIDOCS	17.2	17.6	17.9	17.7	21.9	22.0	22.8	21.4
FEVER	77.8	79.8	79.3	79.3	69.4	72.4	70.2	70.4
Climate-FEVER	22.8	24.0	23.8	23.7	20.2	23.3	20.4	21.0
SciFact	74.1	74.1	73.6	73.5	65.0	65.3	65.3	65.7
BEIR avg.	50.6	<b>50.9</b>	50.5	50.6	51.1	<b>51.6</b>	51.1	51.0

Idx	Base Model	Sorting method	Name	FLOPs to rerank:	
				Top-1	Top-10
0	T5-base	pointwise	MonoT5	1x	1x
1	T5-base	tournament	ListT5(r=1)	1.3x	2.6x
2	T5-base	tournament	ListT5(r=2)	1.8x	4.7x
3	T5-base	sliding w.(s=2)	T5(FiD)	2.5x	9.8x
4	T5-base	sliding w.(s=3)	T5(FiD)	1.7x	12.3x
5	T5-3b	tournament	ListT5(r=1)	17.6x	36.3x
6	T5-3b	tournament	ListT5(r=2)	24.6x	66.0x
7	T5-3b	sliding w.(s=2)	T5(FiD)	38.5x	154x
8	T5-3b	sliding w.(s=2)	T5(no FiD)	53.8x	215.1x
9	T5-3b	sliding w.(s=3)	T5(FiD)	25.6x	128x
10	T5-3b	sliding w.(s=3)	T5(no FiD)	35.1x	175.6x

Table 13: FLOPs (In a multiple of FLOPs of MonoT5-base) on the choice of architecture and method, on TREC-DL19. For the sliding window approach, we would need a total of 4 multiple passes for stride = 3 and 5 passes for stride = 2 (Explained at Tab. 12) to rerank Top-10 candidates.

# Appendix: applying tournament sort on RankGPT

Method	dl19	dl20	trec-covid	news	touche
sliding	68.4 $\pm$ 0.4	64.9 $\pm$ 1.1	72.6 $\pm$ 1.4	46.5 $\pm$ 1.0	<b>38.2 <math>\pm</math> 0.5</b>
tournament	67.4 $\pm$ 0.9	65.8 $\pm$ 0.6	<b>76.4 <math>\pm</math> 0.4</b>	45.5 $\pm$ 1.0	33.1 $\pm$ 1.7

Table 18: NDCG@10 on the selected subset of BEIR, on RankGPT-3.5 with different sorting methods. For fair comparison, we used  $w = 20$ ,  $s = 10$  for the sliding approach, and  $m = 20$ ,  $r = 10$  for the tournament sort. To compensate for the instability of APIs, all results are run for 3 times. Except for trec-covid and touche, differences are statistically non-significant ( $p > 0.1$ ). (Sec. K.2)

# Appendix: Train Dataset Example

input: Query: did edison invent the car battery?, Index: 1, Context: Ransome Eli Olds. The first automobile to be mass produced in the United States was the 1901, Curved Dash Oldsmobile, built by the American car manufacturer Ransome Eli Olds (1864-1950). Olds invented the basic concept of the assembly line and started the Detroit area automobile industry. He first began making steam and gasoline engines with his father, Pliny Fisk Olds, in Lansing, Michigan in 1885. Olds designed his first steam-powered car in 1887. Ransome Eli Olds. The first automobile to be mass produced in the United States was the 1901, Curved Dash Oldsmobile, built by the American car manufacturer Ransome Eli Olds (1864-1950). Olds invented the basic concept of the assembly line and started the Detroit area automobile industry.

>>> Query: did edison invent the car battery?, Index: 2, Context: Correction, Thomas Edison did not originally invented the phonograph, another one invented it before him. His name is Emile Berliner. Edison is known for his cunning ways, he is used to steal others work and make it his own so he got the credit. The same thing he did to the brilliant Nikola Tesla.

>>> Query: did edison invent the car battery?, Index: 3, Context: When Daimler-Benz (makers of Mercedes-Benz cars) says that the automobile was invented in 1886 by Karl Benz and Gottlieb Daimler, it's basing its claim on its own definition: a light carriage for personal transport with three or four wheels, powered by a liquid-fueled internal combustion engine.

>>> Query: did edison invent the car battery?, Index: 4, Context: 1898 - Conrad Hubert, known as the founder of the Eveready Battery Company, invented the electric hand torch, or flashlight - a dry cell battery, bulb and rough brass reflector inside a paper tube. - Eveready introduced the D size battery for the first handheld flashlight.

>>> Query: did edison invent the car battery?, Index: 5, Context: Edison's Alkaline Battery. As with several of Thomas Edison's later projects, such as his effort to mine iron ore and his quest to create synthetic rubber, his attempts at improving the battery did not lead to the results he hoped for. Edison started his work on the battery in the 1890s, just after the automobile had been introduced.

output: 3 1 4 2 5 ({"id": "1885108", "text": "Edison's Alkaline Battery. As with several of Thomas Edison's later projects, such as his effort to mine iron ore and his quest to create synthetic rubber, his attempts at improving the battery did not lead to the results he hoped for. Edison started his work on the battery in the 1890s, just after the automobile had been introduced."})

format: [Query: did edison invent the car battery? Index: 1, Context: ... ]

3: Mercedes-Benz engine

1: first automobile

4: Conrad Hubert's battery company invented flashlight.

2: Edison didn't invent phonograph.

5: Edison's car battery improvement didn't come out as expected.

output: 3 1 4 2 5