

AND CREATE CHANGE



An Investigation of Prompt Variations for Zero-shot LLM-based Rankers

¹ RMIT University, Australia ² CSIRO, Australia ³ The University of Queensland, Australia

* This work was conducted while Shuoqi Sun was a student at The University of Queensland.





<u>Shuoqi Sun^{1*}, Shengyao Zhuang², Shuai Wang³, Guido Zuccon³</u>



LLM Rankers: Prompting LLMs to Rank Documents



- LLM, no need to do SFT or RL for specific ranking task
 - There has no examples provided in the prompt





• All are "zero-shot": i.e. once we obtained the pre-trained, instruction tuned



































Different ranking mechanisms lead to different effectiveness



ie Iab





Given a query {query}, which of the following two passages is more relevant to the query? LLM Passage A: {*passage_1*} Passage B: {*passage_2*} Output Passage A or Passage B:







The PRP Prompt

Passage: {text} Query: {query} Does the passage answer the query?



The RankGPT Prompt

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.



Role Playing

Task Instructions

Evidence Ordering (wrt query) & Position of Evidence (wrt instructions)

> **Output Type Tone Words**



The RankGPT Prompt

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.



What is the effect of differences in wording of these prompt components?

Evidence Ordering (wrt query) & Position of Evidence (wrt instructions)

> **Output Type Tone Words**



The RankGPT Prompt

You are RankGPT, an intelligent assistant

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.







What is the effect of differences in wording of these prompt components?

Is effectiveness differences b/w rankers due to: **RQ2: LLM characteristics such as backbone and size?**

Ione Words

The RankGPT Prompt

You are RankGPT, an intelligent assistant

- RQ1: the actual ranking mechanism, or the choice of words?

response the ranking results, up i

any word or explain.







Component	Ranker	None (0)	1	2	3	4	5
Task Instruction (TI)	pointwise	-	Does the passage answer the query?	Is this passage relevant to the query?	For the following query and document, judge whether they are relevant.	Judge the relevance between the query and the document.	-
	pairwise	-	Given a query, which of the following two passages is more relevant to the query?			-	
	listwise	-	Rank the {num} passages based on their relevance to the search query.	Sort the Passages by their rel- evance to the Query.	I will provide you with {num} passages, each indi- cated by number identifier []. Rank the passages based on their relevance to query.		-
	setwise	-	Which one is the most rele- vant to the query.				
Output Type (OT)	pointwise		Judge whether they are "Highly Relevant", "Some- what Relevant", or "Not Relevant".	From a scale of 0 to 4, judge the relevance.	Answer 'Yes' or 'No'.	Answer True/False.	-
	pairwise	-	Output Passage A or Passage B.				
	listwise	-	Sorted Passages = [The passages should be listed in descending order using identifiers. The most rel- evant passages should be listed first. The output for- mat should be [] >[], e.g., [1] >[2].		-	
	setwise	-	Output the passage label of the most relevant passage.	Generate the passage label.	Generate the passage label that is the most relevant to the query, then explain why you think this passage is the most relevant.		-
Tone Words (TW)	All	$\mathbf{1}$	You better get this right or you will be punished.	Only output the ranking re- sults, do not say any word or explanation.	Please	Only	Must
Role Play- ing (RP)	All		You are RankGPT, an intel- ligent assistant that can rank passages based on their rele- vancy to the query.				

Wording Alternatives

Ŧ



We control the prompt wording option, vary one option at a time

Table 3: Prompt templates that combine the five components with the four ordering options available. Q and P denote query text passage(s) text, respectively.

EO/PE B

QF RP+TI(Q)+P+TW+OT

PF RP+ P + TI (Q) + TW+ OT

\mathbf{D} \mathbf{D} \mathbf{D} \mathbf{D} \mathbf{D} \mathbf{D} \mathbf{D} \mathbf{D} \mathbf{D}
RP + TW + OT + TT(Q) + P
RP+TW+OT+P+TI(Q)





We explore variations in prompt components wordings and ordering

- 1,248 prompt variations
 - e.g. Tone Words: "Please", "You better get this right or you will be punished",
- 12,400+ GPU-hours, 12,000+ results analysed
- 3 LLM backbone families: FlanT5 (L, XL, XXL), Mistral-7B, Llama3-8B
- Experimented across DL 19, DL 20, COVID (BEIR)





Prompts Better Than Original Ones Do Exist

= best prompt variation per ranker family







Prompts Better Than Original Ones Do Exist

= best prompt variation per ranker family





...and pointwise can be as good as other ranking methods



LLM Rankers Can Be (highly) Sensitive to Prompt Variations









Different LLM Rankers Exhibit Different Variability



LLM Backbones Influence Effectiveness & Variations Differently Across LLM Rankers

In the paper we also...

• show similar findings across datasets

- analyse role of each prompt component type, and instance within prompt component type

Key Takeaways

Prompt components beyond ranking method significantly impact effectiveness

Each ranking method has distinct component preferences

No universal "best prompt" exists: depends on ranking method, dataset, and LLM

Key Takeaways

Prompt components beyond ranking method significantly impact effectiveness

Each ranking method has distinct component preferences

No universal "best prompt" exists: depends on ranking method, dataset, and LLM Future work: automatic prompt optimisation & prompt performance prediction

