



# Opportunities and Challenges of LLMs in Information Retrieval

---

Chuan Meng

University of Amsterdam

14<sup>th</sup> August 2024

# About me



## Chuan Meng

- Third-year PhD student at the University of Amsterdam
- Supervisors: Maarten de Rijke, Mohammad Aliannejad
- Applied Scientist Intern at Amazon (London)
  
- Focus on:
  - Query performance prediction
  - Mixed-initiative conversational search
  - LLM-based re-ranking
  - LLM-based relevance judgment prediction
  
- As of Aug 2024, I have authored 15 papers
  - 222 citations (Google Scholar) with an H-index of 6

# Background

- Large language models (LLMs) have remarkable language understanding, generation, generalization, and reasoning abilities

who is Pablo Mendes?

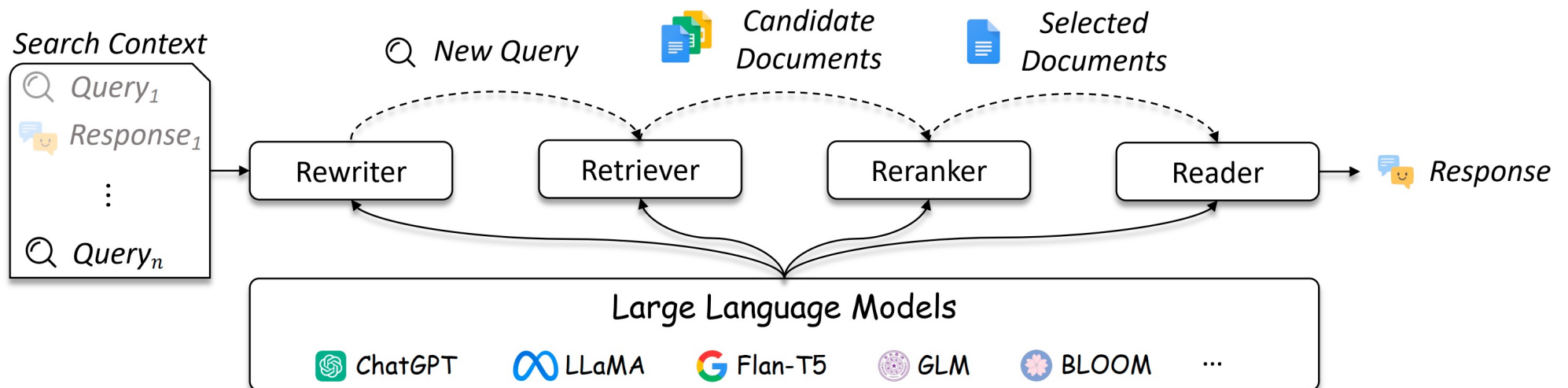


Pablo Mendes is a researcher and entrepreneur known for his work in the fields of natural language processing (NLP), semantic web, and data science. He has contributed to various projects involving knowledge graphs, semantic search, and machine learning. Mendes has held positions in both academia and industry, where he has worked on developing technologies that improve how machines understand and process human language.

Mendes has also been involved in startups and organizations focused on advancing AI and data science technologies. His work often intersects with innovations in how information is structured and retrieved, making him a notable figure in the tech community, particularly in areas related to artificial intelligence and semantic technologies.

# Background

- Large language models (LLMs) in information retrieval (IR)
  - Opportunities
    - LLMs directly as task solvers (e.g., retrievers/re-rankers)
    - LLMs for data augmentation (e.g., training retrievers/re-rankers)
    - LLMs for evaluation (e.g., generating relevance judgments)
  - Challenges
    - Low efficiency
    - Unfaithful generation
    - ...



# Background

- Large language models (LLMs) in information retrieval (IR)
  - Opportunities
    - LLMs directly as task solvers [3,4,7]
    - LLMs for data augmentation (e.g., training retrievers/re-rankers) [5,6]
    - LLMs for evaluation (e.g., generating relevance judgments) [2]
  - Challenges
    - Low efficiency [1]
    - Unfaithful generation
    - ...

[1] Ranked List Truncation for Large Language Model-based Re-Ranking. SIGIR 2024

[2] Query Performance Prediction using Relevance Judgments Generated by Large Language Models. arXiv 2024.

[3] Generative Retrieval with Few-shot Indexing. arXiv 2024.

[4] LLM-based Retrieval and Generation Pipelines for TREC Interactive Knowledge Assistance Track (iKAT) 2023. TREC 2023.

[5] Expand, Highlight, Generate: RL-driven Document Generation for Passage Reranking. EMNLP 2023.

[6] Self-seeding and Multi-intent Self-instructing LLMs for Generating Intent-aware Information-Seeking dialogs. arXiv 2024.

[7] System Initiative Prediction for Multi-turn Conversational Information Seeking. CIKM 2023

# Background

- Large language models (LLMs) in information retrieval (IR)
  - Opportunities
    - **LLMs directly as task solvers [3,4,7]**
    - LLMs for data augmentation (e.g., training retrievers/re-rankers) [5,6]
    - **LLMs for evaluation (e.g., generating relevance judgments) [2]**
  - Challenges
    - **Low efficiency [1]**
    - Unfaithful generation
    - ...

[1] Ranked List Truncation for Large Language Model-based Re-Ranking. SIGIR 2024

[2] Query Performance Prediction using Relevance Judgments Generated by Large Language Models. arXiv 2024.

[3] Generative Retrieval with Few-shot Indexing. arXiv 2024.

[4] LLM-based Retrieval and Generation Pipelines for TREC Interactive Knowledge Assistance Track (iKAT) 2023. TREC 2023.

[5] Expand, Highlight, Generate: RL-driven Document Generation for Passage Reranking. EMNLP 2023.

[6] Self-seeding and Multi-intent Self-instructing LLMs for Generating Intent-aware Information-Seeking dialogs. arXiv 2024.

[7] System Initiative Prediction for Multi-turn Conversational Information Seeking. CIKM 2023

# Outline

- ❑ Study 1: Ranked List Truncation for Large Language Model-based Re-Ranking [15 min]
- ❑ Study 2: Query Performance Prediction using Relevance Judgments Generated by Large Language Models [10 min]
- ❑ Study 3: Generative Retrieval with Few-shot Indexing [10 min]
- ❑ Conclusion [5 min]

# Outline

- ❑ **Study 1: Ranked List Truncation for Large Language Model-based Re-Ranking [15 min]**
- ❑ Study 2: Query Performance Prediction using Relevance Judgments Generated by Large Language Models [10 min]
- ❑ Study 3: Generative Retrieval with Few-shot Indexing [10 min]
- ❑ Conclusion [5 min]





# Ranked List Truncation for Large Language Model-based Re-Ranking

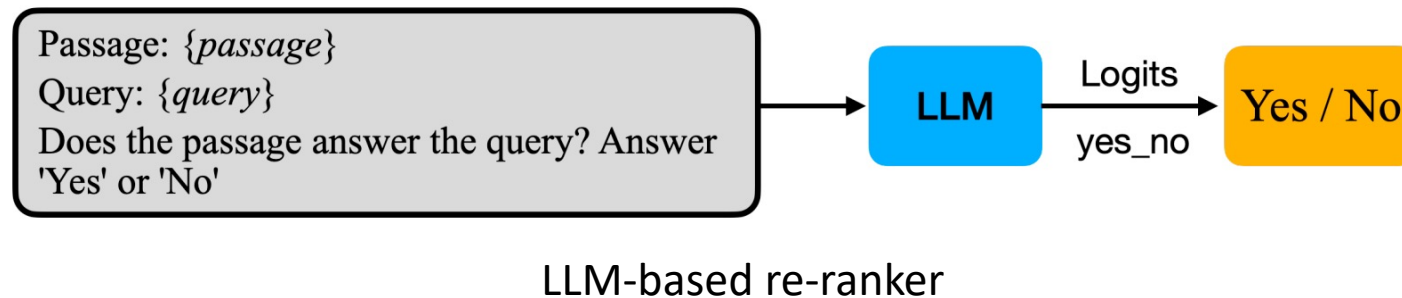
---

**Chuan Meng, Negar Arabzadeh, Arian Askari, Mohammad Aliannejadi,  
Maarten de Rijke**

The 47th International ACM SIGIR Conference on Research and  
Development in Information Retrieval (SIGIR 2024)

# Background

- Large language models (LLMs) as text re-rankers
  - achieve state-of-the-art performance
  - hard to be applied in practice due to significant computational overhead
    - the average query latency (re-ranking 100 items per query) for Flan-t5-xxl (11B) is around 4 seconds, on a NVIDIA RTX A6000 GPU [1]

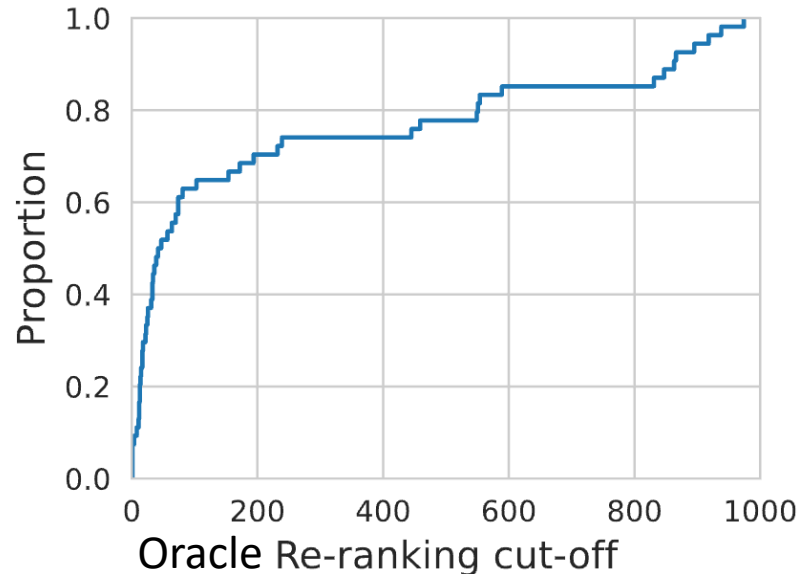


# Motivation

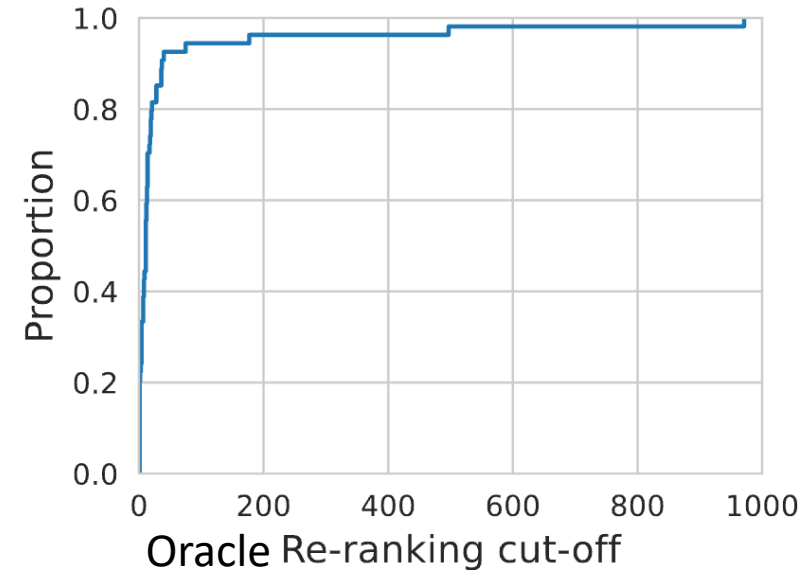
- Common practice: applying a fixed re-ranking cut-off to all queries (e.g., 100, 200, 1000)
- However,
  - a fixed re-ranking cut-off might lead to a waste of computational resources
  - individual queries might need a shorter or a longer list of re-ranking candidates
- We explore query-specific re-ranking cut-offs in the context of LLM-based re-ranking
  - Fixed cut-offs vs. query-specific cut-offs
  - How to predict query-specific cut-offs

# Motivation (fixed cut-offs vs. query-specific cut-offs)

- Query-specific re-ranking cut-offs improve *efficiency*
  - Individual queries have different oracle cut-offs with a wide range
  - A deep fixed cut-off wastes computational resources
  - A shallow fixed cut-off hurts re-ranking quality for queries needing a deeper cut-off



Cumulative distribution function of oracle cut-offs for  
BM25-RankLLaMA  
TREC-DL 20

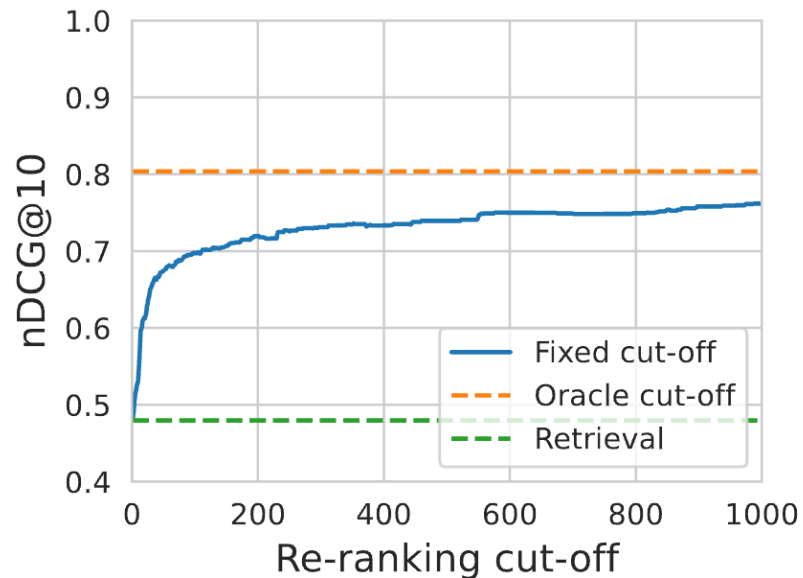


Cumulative distribution function of oracle cut-offs for  
RepLLaMA-RankLLaMA  
TREC-DL 20

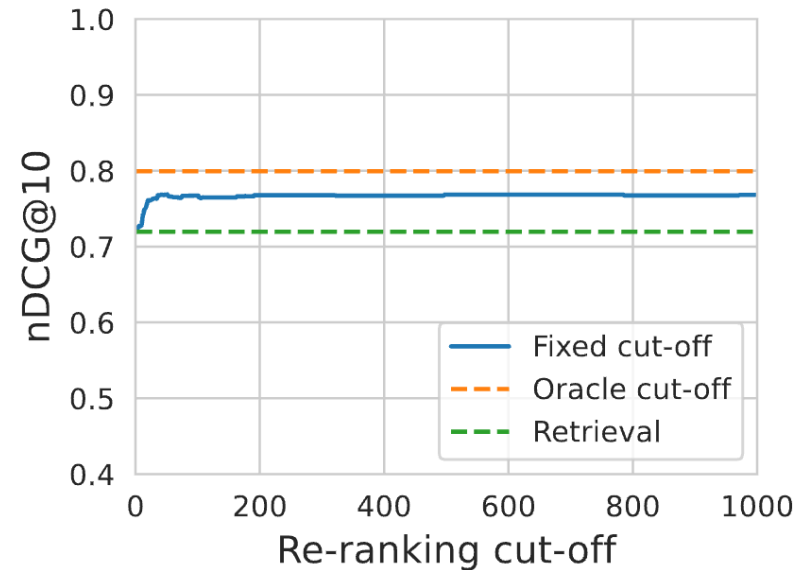
For a query, an oracle cut-off is the minimum re-ranking cutoff producing the highest nDCG@10 value

# Motivation (fixed cut-offs vs. query-specific cut-offs)

- Query-specific re-ranking cut-offs improve *effectiveness*
  - Oracle cut-offs show statistically significant improvements over all fixed cut-offs
  - A deeper fixed cut-off
    - does not always result in improvement (consistent with [1])
    - even is detrimental to re-ranking quality (consistent with [1])



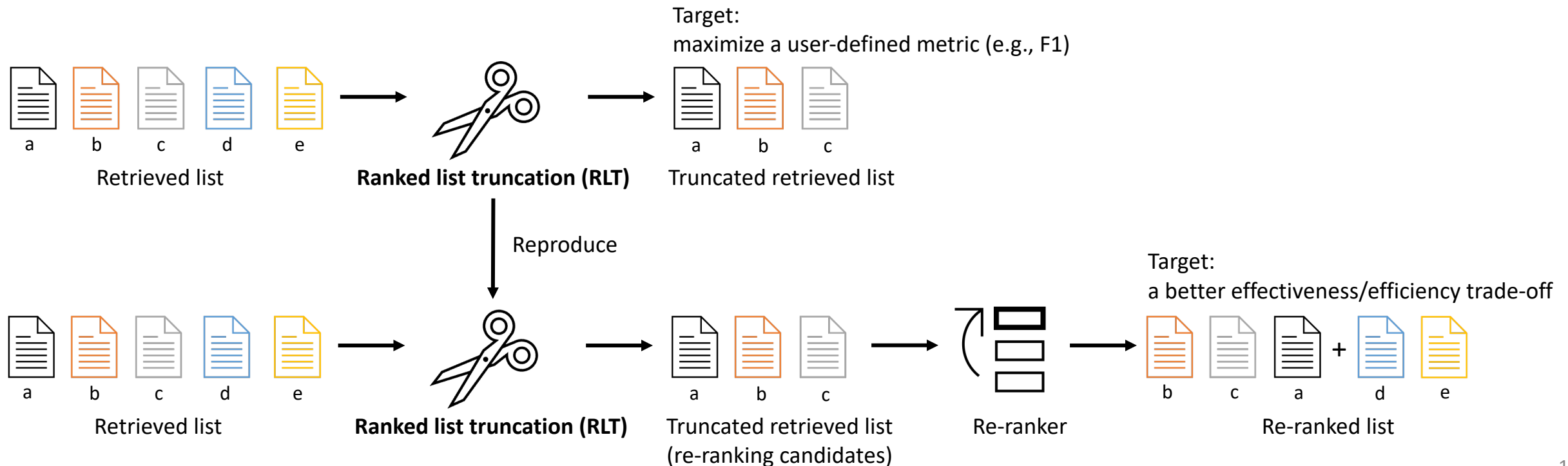
BM25-RankLLaMA  
TREC-DL 20



RepLLaMA-RankLLaMA  
TREC-DL 20

# Motivation (How to predict query-specific cut-offs)

- Ranked list truncation (RLT)
  - predicts how many items in a ranked list should be returned
  - optimizes the truncated ranked list regarding a user-defined metric (e.g., F1)
  - aids applications where reviewing returned items is costly, e.g., patent or legal search
- **We reproduce exiting RLT methods in the context of LLM-based re-ranking**



# Reproducibility methodology

- *Do RLT methods generalize to the context of*
  - *(RQ1) LLM-based re-ranking with a lexical first-stage retriever?*
  - *(RQ2) LLM-based re-ranking with learned sparse or dense first-stage retrievers?*
  - *(RQ3) pre-trained language model-based re-ranking?*

# Reproducibility methodology

- Experimental settings:
  - 8 RLT methods

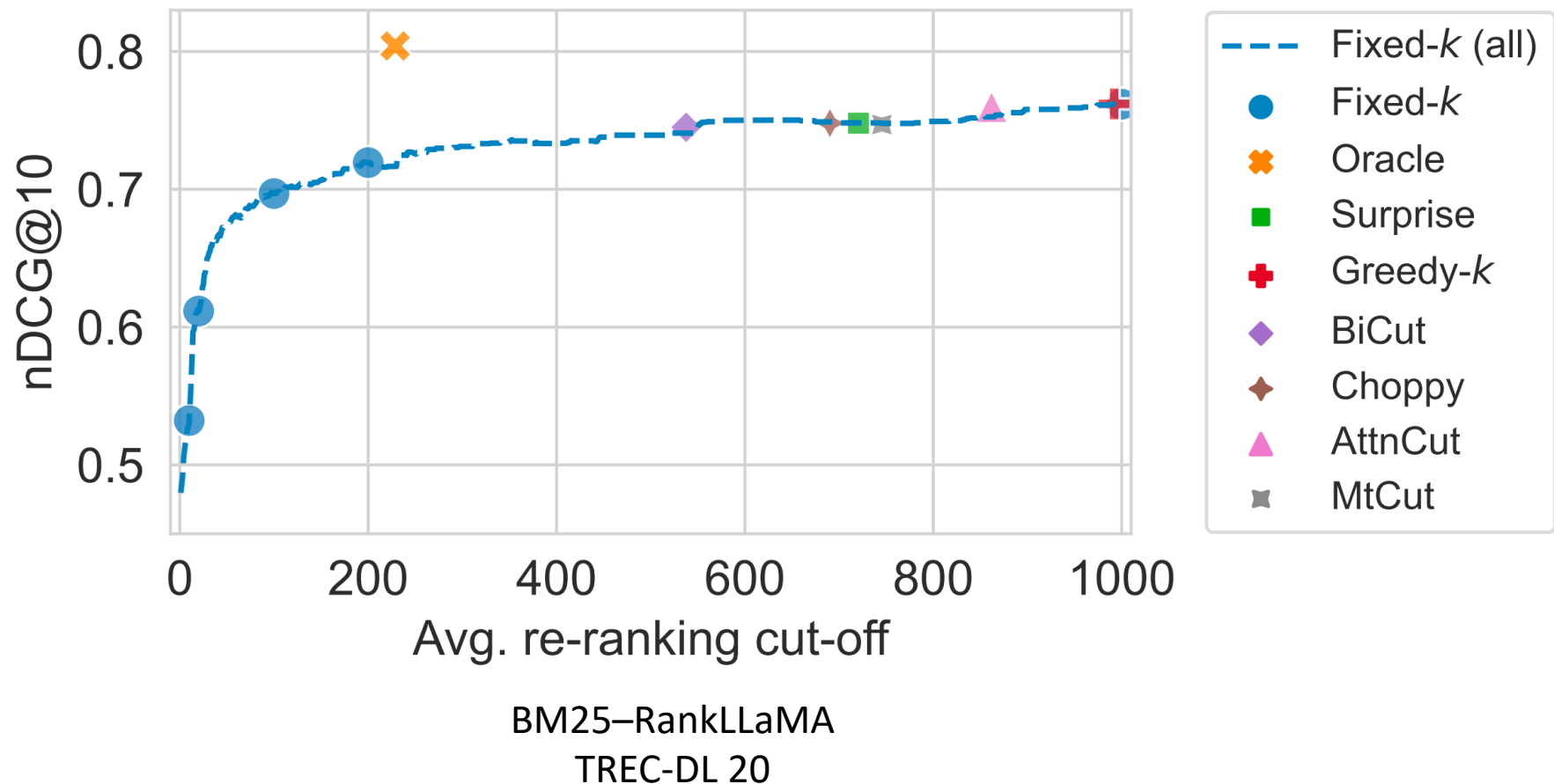
Method	Attribute 1
Fixed- $k$ (10, 20, 100, 200, 1000)	Unsupervised
Greedy- $k$	Unsupervised
Surprise	Unsupervised
BiCut	Supervised
Choppy	Supervised
AttnCut	Supervised
MtCut	Supervised
LeCut	Supervised

- Datasets:
  - TREC-DL 19, TREC-DL 20



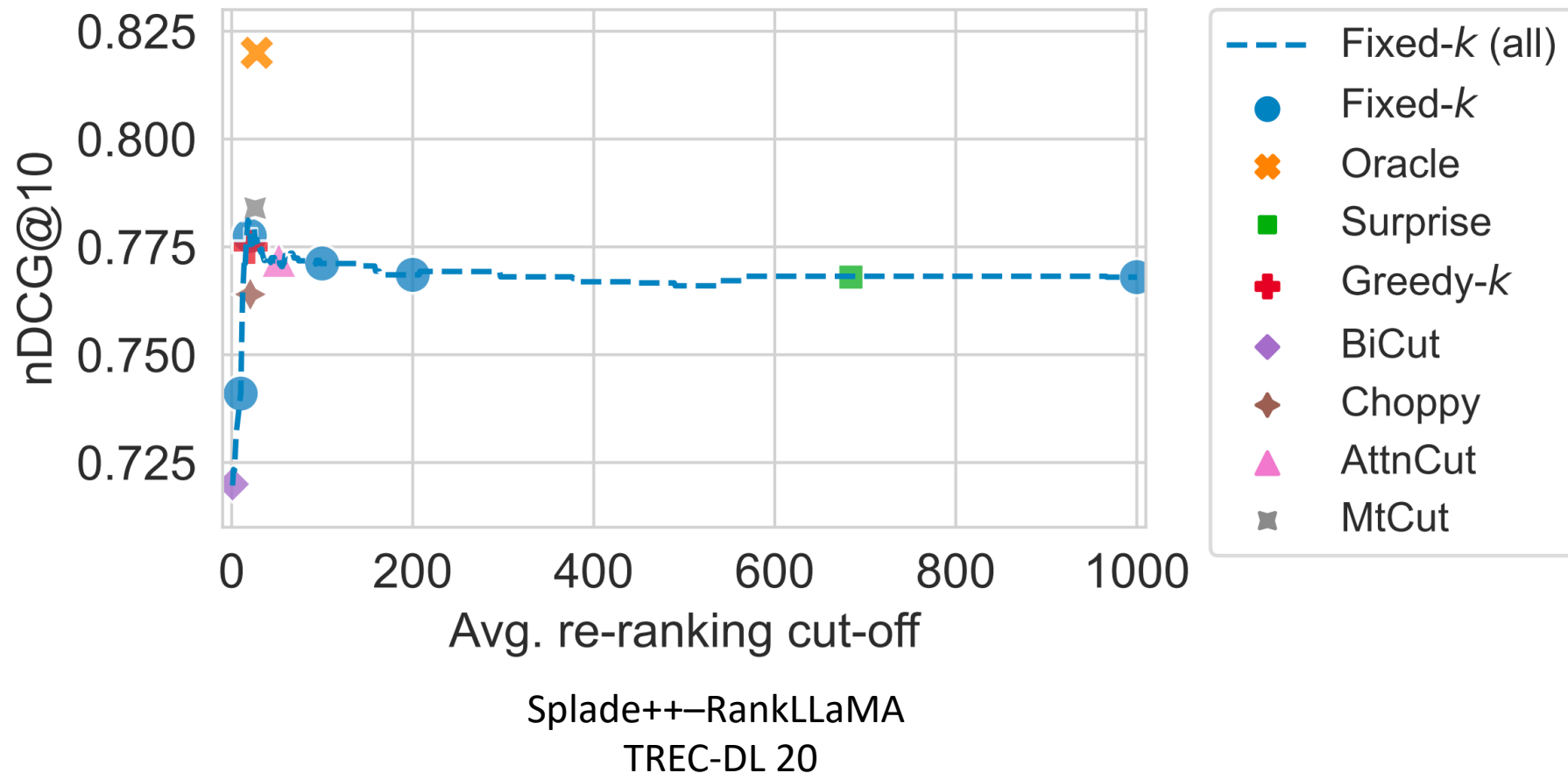
# Experiments

- RQ1: Do RLT methods generalize to the context of LLM-based re-ranking with a lexical first-stage retriever?
  - Fixed re-ranking depths can closely approximate supervised RLT methods' results
  - Supervised RLT methods do not show a clear advantage over fixed re-ranking depths



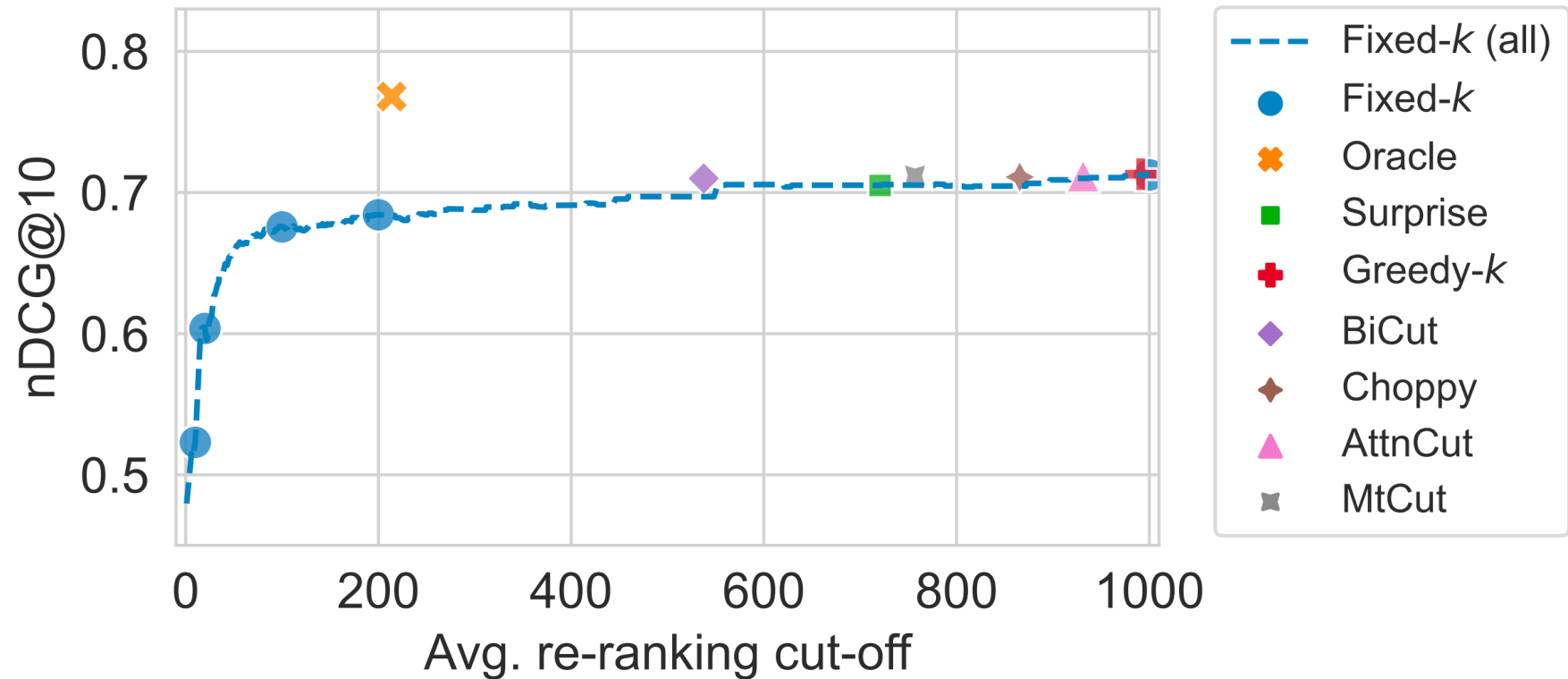
# Experiments

- RQ2: Do RLT methods generalize to the context of LLM-based re-ranking with learned sparse or dense first-stage retriever?
  - Supervised methods do not lead to significant improvement in terms nDCG@10
  - A fixed re-ranking depth of 20 achieves the best effectiveness/efficiency trade-off



# Experiments

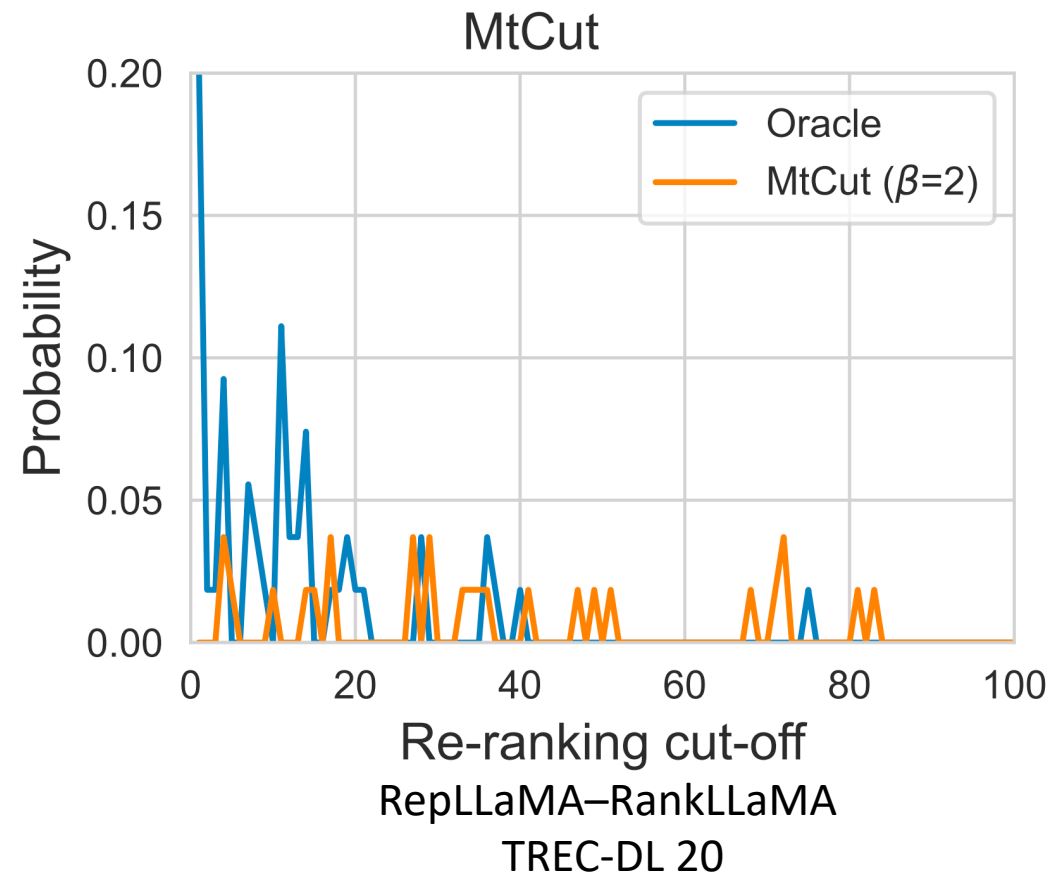
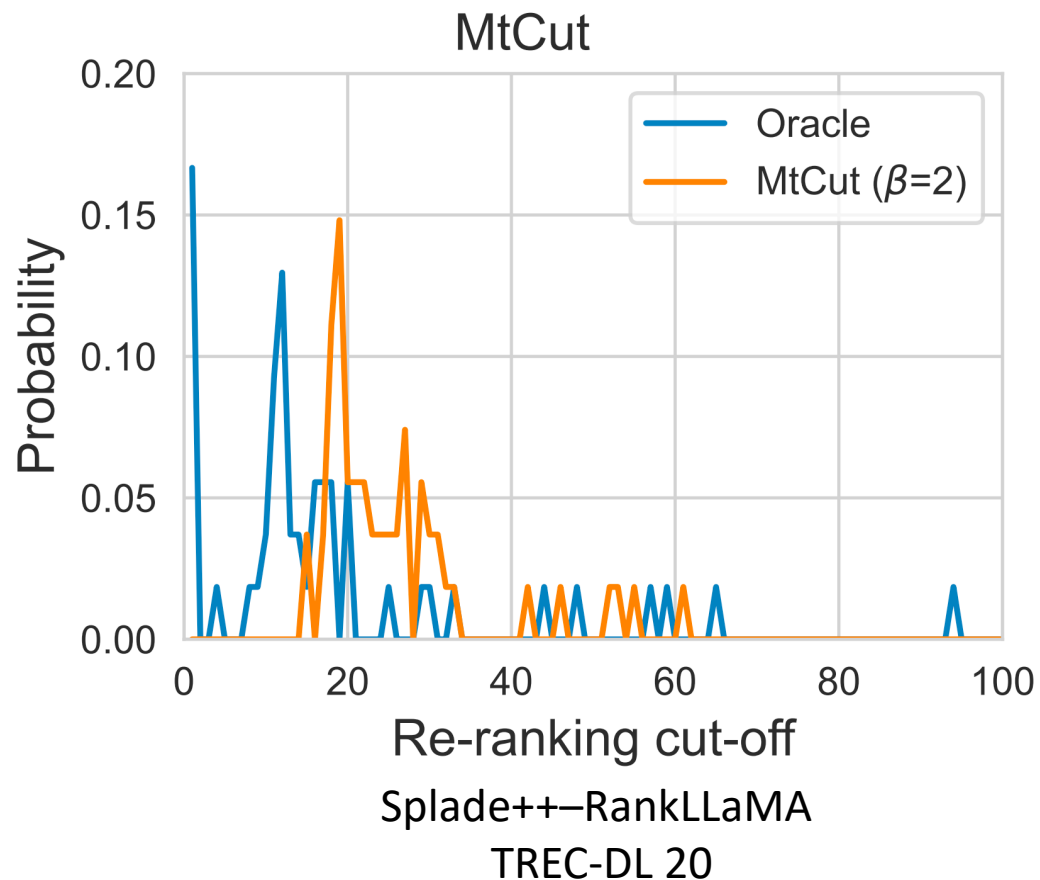
- RQ3: Do RLT methods generalize to the context of pre-trained language model-based re-ranking?
  - Results are similar to RQ1



BM25-monoT5  
TREC-DL 20

# Experiments

- Error analysis for supervised RLT methods
  - They fail to predict a re-ranking cut-off of zero



# Takeaways

- The type of retriever makes a difference
  - With an effective retriever (e.g., SPLADE++/RepLLaMA)
    - A fixed re-ranking depth of **20** yields an excellent effectiveness/efficiency trade-off
    - A fixed depth **>20** does not significantly improve re-ranking quality
- The type of re-ranker (LLM or pre-trained LM-based) does not appear to influence the findings
- Supervised RLT methods need to improve their ability to predict “0”

# Conclusion

- Contributions
  - An empirical analysis in the context of LLM-based re-ranking, shows that
    - Effective query-specific re-ranking depths can improve re-ranking efficiency and effectiveness
  - We reproduce RLT methods in the context of LLM-based re-ranking
  - The data and code are open-source <https://github.com/ChuanMeng/RLT4Reranking>
- Future work
  - Explore RLT for pairwise and listwise LLM-based re-rankers
  - Develop new RLT methods for LLM-based re-ranking

Q & A



QR code for the repo

# Outline

- ❑ Study 1: Ranked List Truncation for Large Language Model-based Re-Ranking [15 min]
- ❑ **Study 2: Query Performance Prediction using Relevance Judgments Generated by Large Language Models [10 min]**
- ❑ Study 3: Generative Retrieval with Few-shot Indexing [10 min]
- ❑ Conclusion [5 min]



# Query Performance Prediction using Relevance Judgments Generated by Large Language Models

---

**Chuan Meng, Negar Arabzadeh, Arian Askari, Mohammad Aliannejadi,  
Maarten de Rijke**

Under review



# Motivation

- Why fine-tune LLMs?

LLM	TREC-DL 19	TREC-DL 20	TREC-DL 21	TREC-DL 22
	$\kappa$	$\kappa$	$\kappa$	$\kappa$
GPT-3.5 (text-davinci-003) [32]	-	-	0.260	-
LLaMA-7B (few-shot)	-0.001	-0.003	0.003	-0.010
Llama-3-8B (few-shot)	0.018	0.027	0.021	-0.035
Llama-3-8B-Instruct (few-shot)	0.315	0.227	0.238	0.049

# Methodology

- Fine-tuning open-source LLMs for generating relevance judgments
  - LLMs: LLaMA-7B, Llama-3-8B, and Llama-3-8B-Instruct
  - Fine-tuning method: QLoRA, a parameter-efficient fine-tuning method
  - Training data: human-labeled relevance judgments of MS MARCO

**Instruction:** Please assess the relevance of the provided passage to the following question.

Please output “Relevant” or “Irrelevant”.

Question: {question}

Passage: {passage}

Output: Relevant/Irrelevant

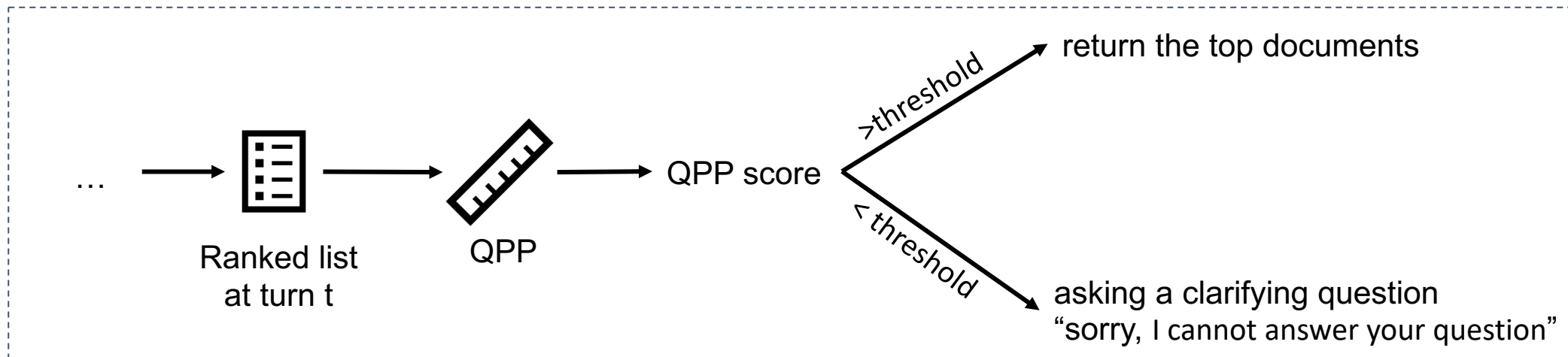
# Experiments

- Fine-tuned LLMs outperform
  - their few-shot prompted counterparts
  - GPT-3.5

LLM	TREC-DL 19	TREC-DL 20	TREC-DL 21	TREC-DL 22
	$\kappa$	$\kappa$	$\kappa$	$\kappa$
GPT-3.5 (text-davinci-003) [32]	-	-	0.260	-
LLaMA-7B (few-shot)	-0.001	-0.003	0.003	-0.010
Llama-3-8B (few-shot)	0.018	0.027	0.021	-0.035
Llama-3-8B-Instruct (few-shot)	0.315	0.227	0.238	0.049
LLaMA-7B (fine-tuned)	0.258	0.238	0.333	0.038
Llama-3-8B (fine-tuned)	0.381	<b>0.342</b>	0.347	<b>0.082</b>
Llama-3-8B-Instruct (fine-tuned)	<b>0.397</b>	0.316	<b>0.418</b>	0.066

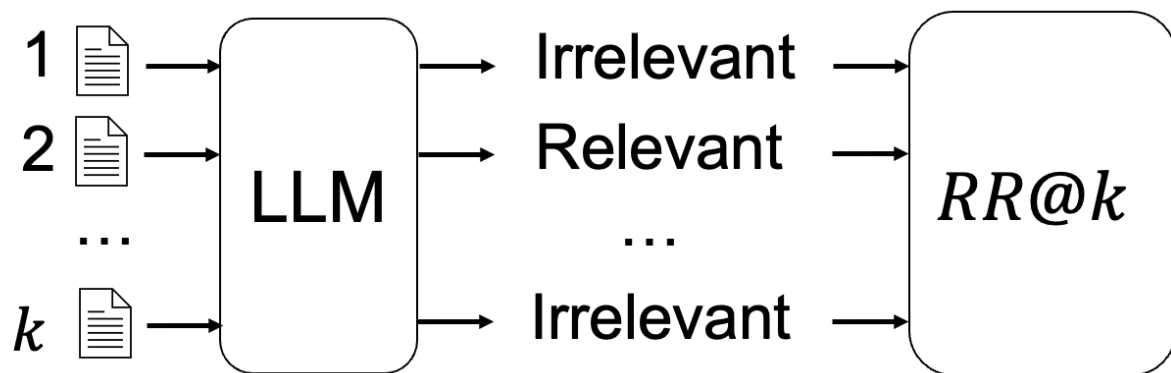
# Background—query performance prediction

- Query performance prediction (QPP)
  - Predicts retrieval quality of search system for query without human-labeled relevance judgments
- QPP benefits a variety of applications, e.g., action prediction in conversational search

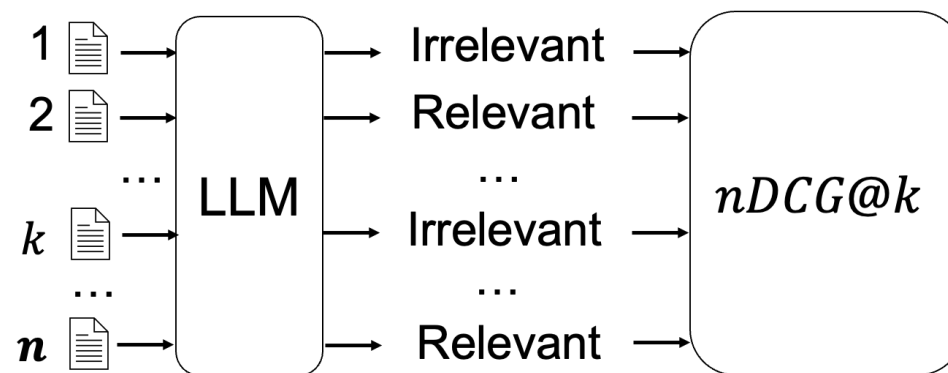


# Methodology

- Propose QPP-GenRE, which predicts IR measures using LLM-generated judgments
  - devise an approximation strategy for predicting a metric considering recall
    - only judges the top  $n$  items in a ranked list, where  $n \ll \#$  documents in the corpus



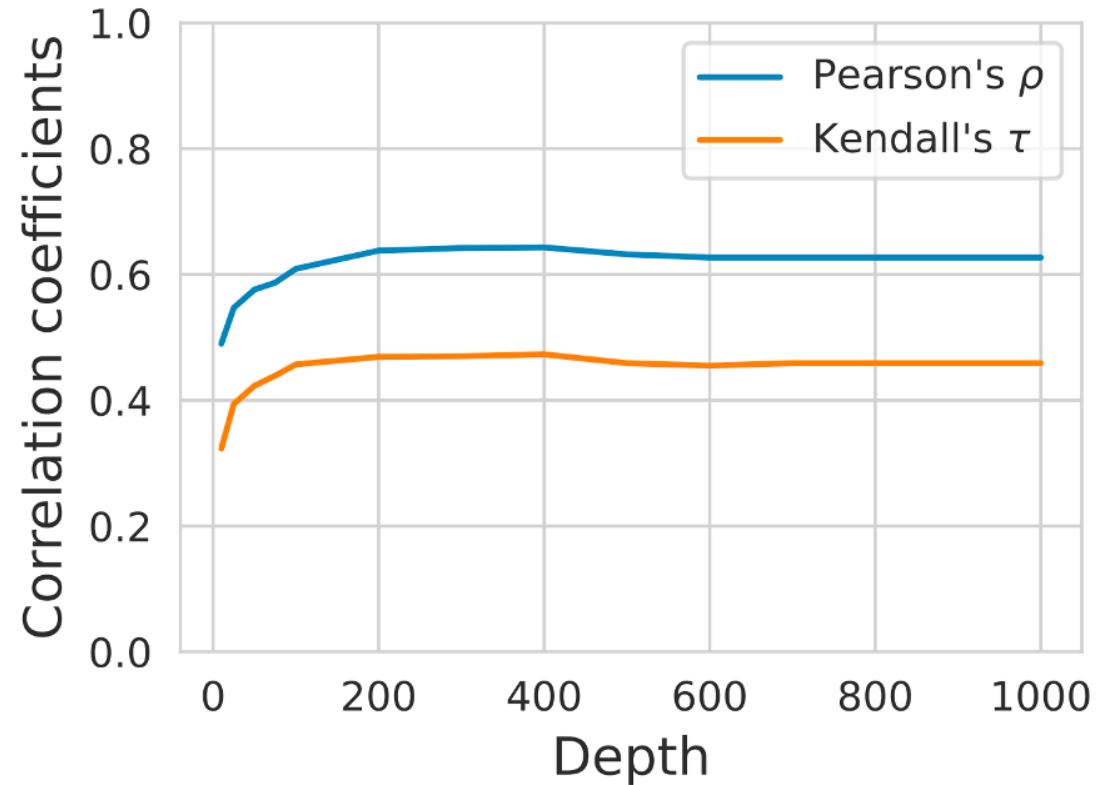
Predicting a precision-based metric



Predicting a metric considering recall

# Experiments

- Findings
  - QPP-GenRE with fine-tuned LLMs achieves SOTA QPP quality
  - Judging up to 100–200 items in a ranked list suffices for predicting nDCG@10



QPP quality of predicting BM25's nDCG@10 w.r.t. judging depth

# Conclusion

- Contributions
  - Fine-tune open-source LLMs for generating relevance judgments
  - Propose a new QPP framework, QPP-GenRE, which predicts IR metrics based on LLM-generated relevance judgments
    - Devise an approximation strategy for predicting a metric considering recall
- QPP-GenRE achieves state-of-the-art QPP quality
- The data and code are open-sourced <https://github.com/ChuanMeng/QPP-GenRE>

Q & A



# Outline

- ❑ Study 1: Ranked List Truncation for Large Language Model-based Re-Ranking [15 min]
- ❑ Study 2: Query Performance Prediction using Relevance Judgments Generated by Large Language Models [10 min]
- ❑ **Study 3: Generative Retrieval with Few-shot Indexing [10 min]**
- ❑ Conclusion [5 min]





# Generative Retrieval with Few-shot Indexing

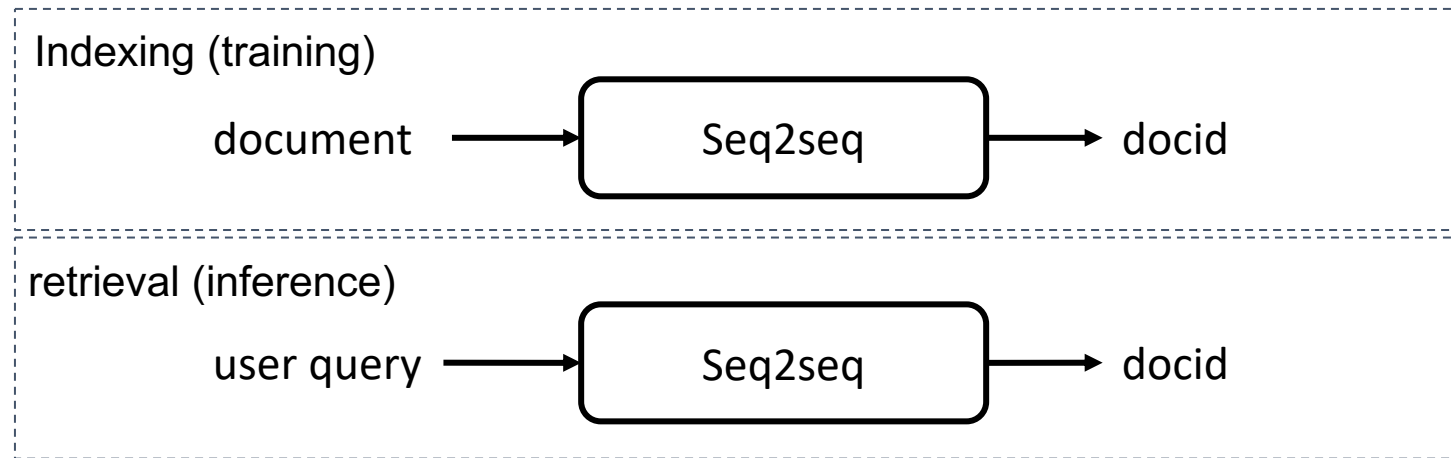
---

Arian Askari\*, Chuan Meng\* (co-first author), Mohammad Aliannejadi,  
Zhaochun Ren, Evangelos Kanoulas, Suzan Verberne

Under review

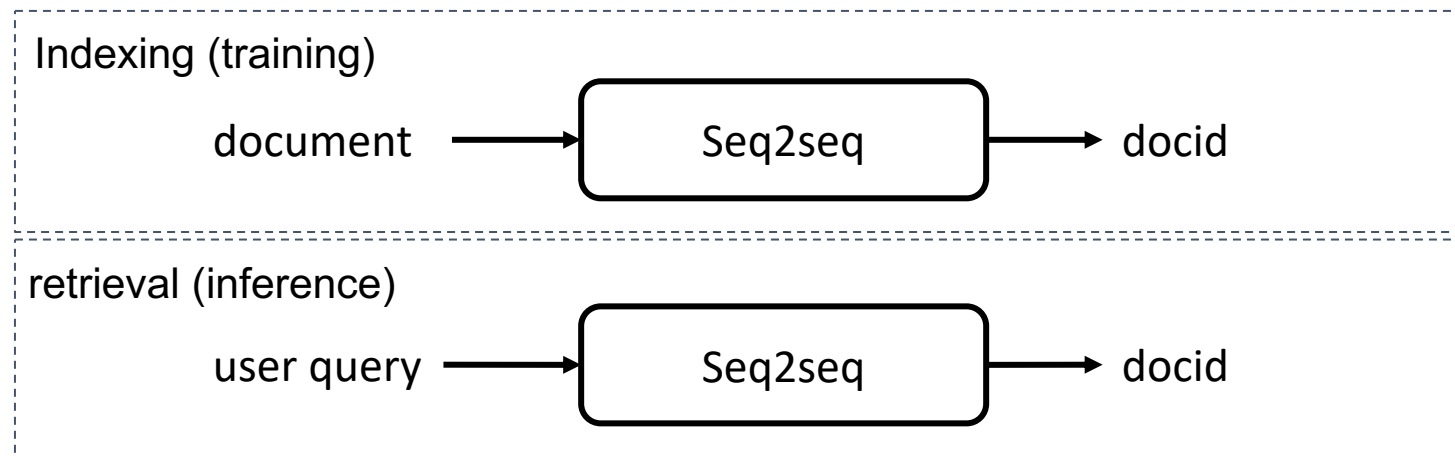
# Background—generative retrieval

- Generative retrieval consolidates indexing and retrieval into a single model
  - Indexing (training) trains seq2seq model to map queries to the docids corresponding to relevant documents
  - Retrieval (inference) feeds the trained model a query text to generate potentially relevant docids.



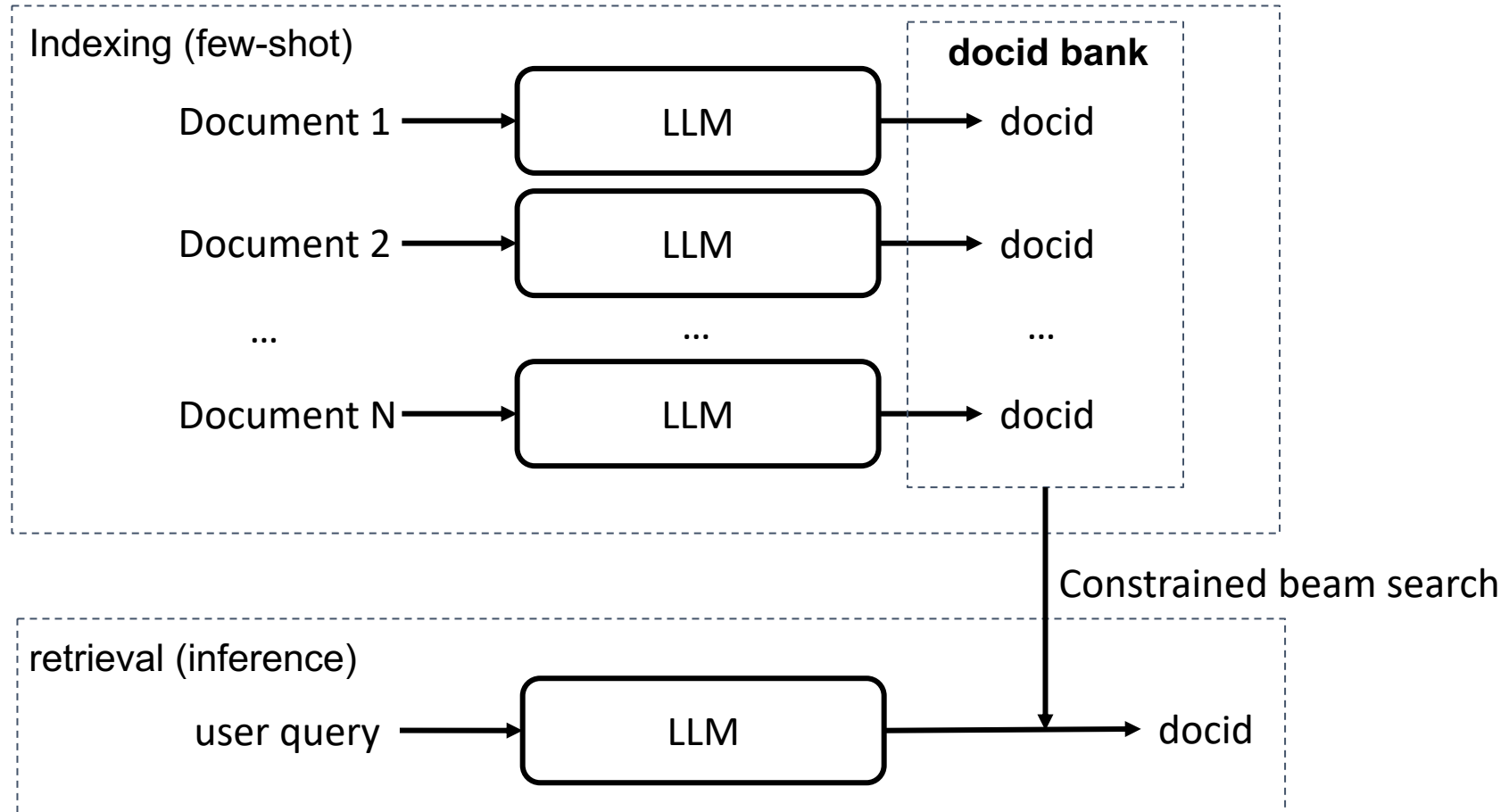
# Motivation

- Previous studies typically rely on **training-based indexing**
  - high training overhead
  - under-utilization of the pre-trained knowledge of LLMs
  - challenges in adapting to a dynamic document corpus



# Methodology

- We propose a **few-shot** indexing-based generative retrieval framework (Few-shot GR)



# Experiments

- Few-shot GR
  - achieves superior performance to SOTA baselines that require heavy training
  - is much more efficient than SOTA baselines
- Selecting a generally stronger LLM leads to better performance

Method	Recall@1	Recall@10	MRR@100
BM25	29.7	60.3	40.2
DocT5Query	38.0	69.3	48.9
DPR	50.2	77.7	59.9
ANCE	50.2	78.5	60.2
SentenceT5	53.6	83.0	64.1
GTR-base	56.0	84.4	66.2
SEAL	59.9	81.2	67.7
DSI	55.2	67.4	59.6
NCI	66.4	85.7	73.6
DSI-QG	63.1	80.7	69.5
DSI-QG (InPars)	63.9	82.0	71.4
GenRET	68.1	<b>88.8</b>	<u>75.9</u>
TOME	66.6	–	–
GLEN	<u>69.1</u>	86.0	75.4
Few-Shot GR	<b>70.1</b>	<u>87.6</u>	<b>77.4</b>

Method	Indexing (hr)	Retrieval (ms)
DSI-QG	240	72
GenRET	≈16,800	72
Few-Shot GR	37	98

The authors of GenRET indicated it took 7 days on 100 A100 GPUs ≈16,800 hours on a single A100 GPU

Method	Recall@1	Recall@10	MRR@100
T5-base	52.4	66.4	55.8
Zephyr-7B- $\beta$	69.9	87.2	<b>77.8</b>
llama-3-8B-Instruct	<b>70.1</b>	<b>87.6</b>	77.4

# Conclusions

- Contributions
  - Propose Few-shot GR, a new generative retrieval paradigm, which performs indexing solely by prompting an LLM without requiring any training
  - Few-Shot GR achieves superior performance to SOTA baselines that require heavy training
- Future work
  - Test Few-shot GR on a document corpus with millions of documents

Q & A

# Outline

- ❑ Study 1: Ranked List Truncation for Large Language Model-based Re-Ranking [15 min]
- ❑ Study 2: Query Performance Prediction using Relevance Judgments Generated by Large Language Models [10 min]
- ❑ Study 3: Generative Retrieval with Few-shot Indexing [10 min]
- ❑ **Conclusion** [5 min]

# Conclusion and Future Work

- Contributions
  - The challenge of low efficiency
    - Predict query-specific re-ranking cut-offs for LLM-based re-ranking
  - The opportunity to use LLMs for evaluation
    - Fine-tune open-source LLMs to generate relevance judgments
    - A new QPP framework using LLM-based generated relevance judgments
  - The opportunity to use LLMs as task solvers
    - Propose a Few-shot generative retrieval framework
- Future work
  - Propose new RLT methods for LLM-based re-ranking
  - Investigate the performance of other open-source LLMs
  - Domain-specific scenarios



# Thank you!

Chuan Meng



c.meng@uva.nl



@ChuanMg



<https://chuanmeng.github.io>



Personal website