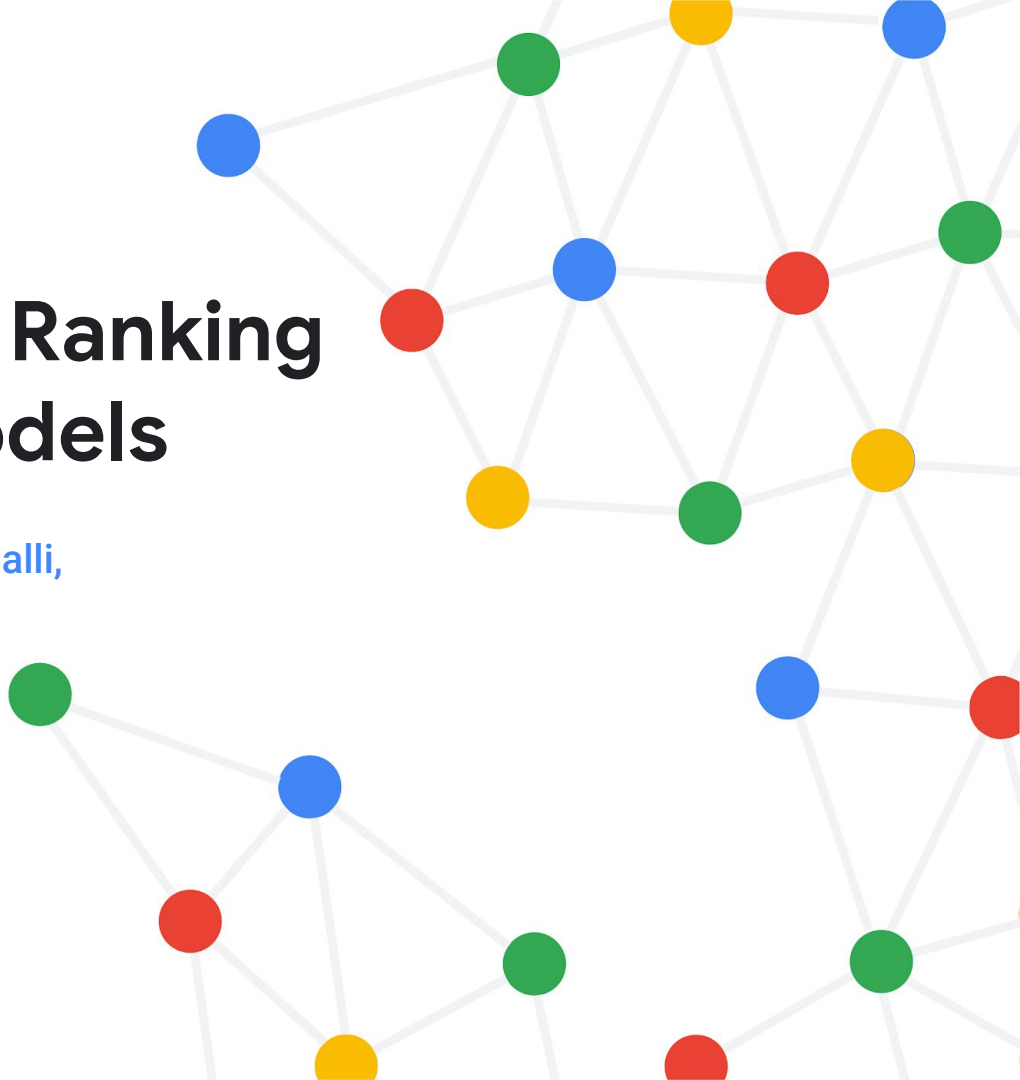


Scalable In-context Ranking with Generative Models

Nilesh Gupta, Chong You, Srinadh Bhojanapalli,
Sanjiv Kumar, Inderjit Dhillon & Felix Yu

[Project Link](#)



Agenda



01

IR in the age of LLMs

02

ICR Attention Analysis

03

BlockRank - scalable ICR

04

Future Directions

Agenda



01

IR in the age of LLMs

02

ICR Attention Analysis

03

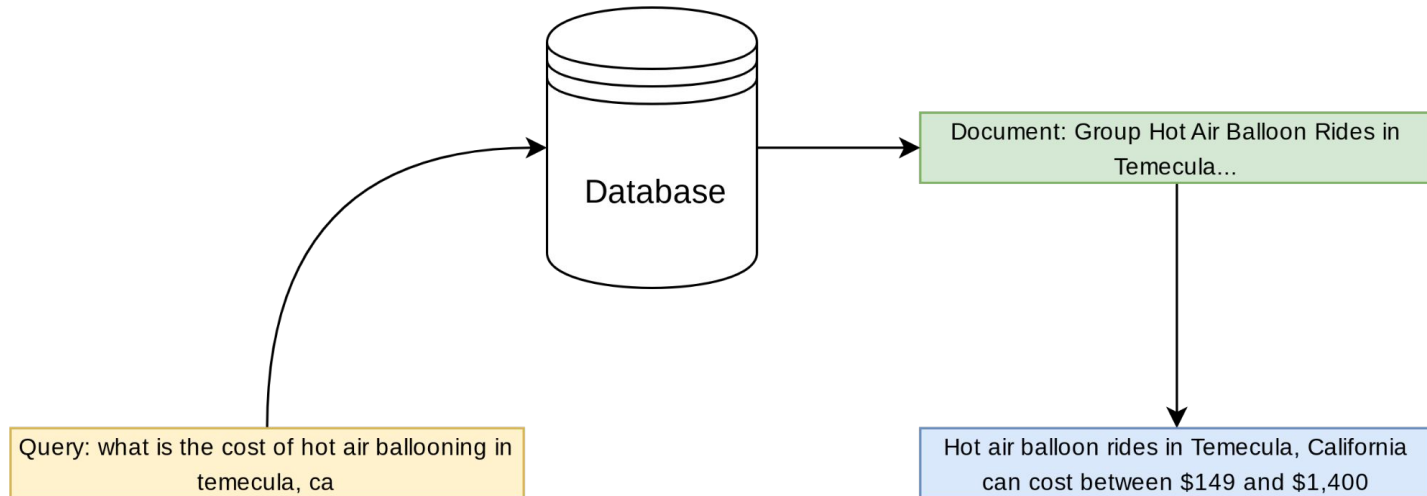
BlockRank - scalable ICR

04

Future Directions

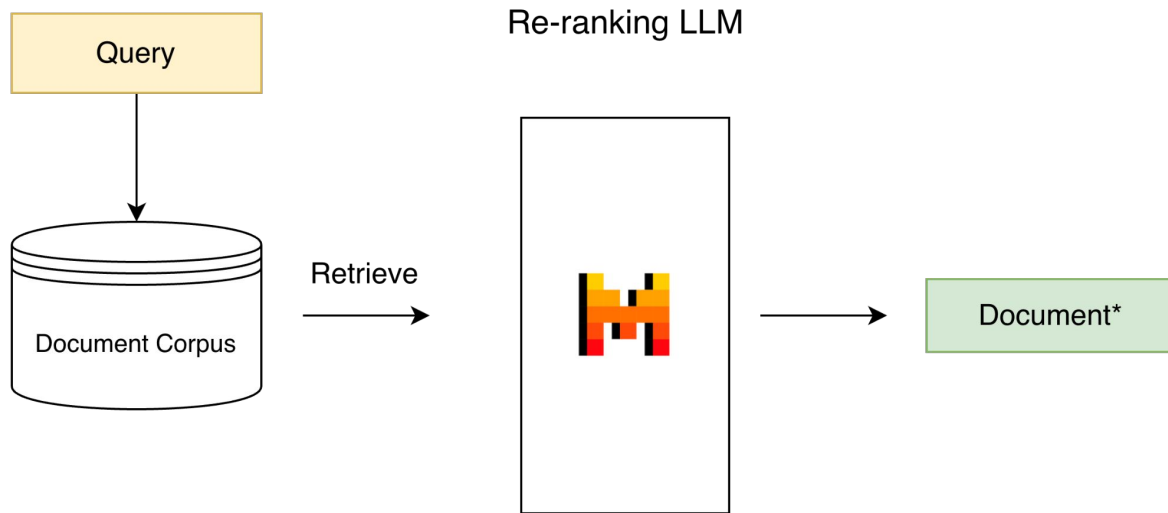
Information Retrieval (IR)

Seeking answer for a query from a database



Typical IR Pipeline

Retrieve (embedding based) then Rerank (strong language model)



LLMs have proven to be strong re-rankers especially as in-context listwise ranker

Reference

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

Corpus Formatting

Few-shot Examples

Query Formatting

Jinhyuk Lee* Anthony Chen* Zhuyuan Dai*
Dheeru Dua Devendra Singh Saxena Michael Boratko Yi Luan
Sébastien M. R. Arnold Vincent Perot Siddharth Dalmia Hexiang Hu
Xudong Lin Panupong Pasupat Aida Amini Jeremy R. Co
Christian Riedel Itzhak Naim Ming-Wei Chang
Google DeepMind

Is ChatGPT Good at Search?
Shuangqing Wang*
Zhaosheng Ren*
Xinyu Ma*
Dawei Yin*
Lingyong Yan*
Zhaomin Chen*
Pengjie Ren*
Xinlong University, Qingdao, China The Netherlands
Tsinghua University, Beijing, China
*lingyong.yan@tsinghua.edu.cn, pengjie.ren@tue.nl, zhaosheng.ren@tue.nl, xinyu.ma@tue.nl, dawei.yin@tue.nl, shuangqing.wang@gmail.com

**RankZephyr: Effective and Robust Zero-Shot
Listwise Reranking is a Breeze!**
Roman Pradeep, Sahel Sharifmoghadam,
David R. Cheriton

Average nDCG@10

[illegible]

Abstract

Large Models (LLMs) have demonstrated remarkable zero-shot generalization across various language-related tasks, including various engines. However, existing work on the generative ability of LLMs for information ranking (IR) rather than direct passage ranking. The discrepancy between these two objectives poses another challenge. In this paper, we first investigate generative LLMs such as ChatGPT and GPT-4 for relevance ranking. Surprisingly, our experiments reveal that property instructs LLMs to deliver competitive, even superior results to state-of-the-art supervised methods on popular IR benchmarks. Furthermore, to address concerns about data contamination of LLMs, we collect a new test

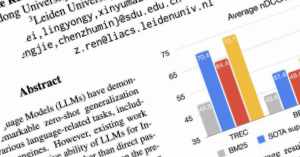
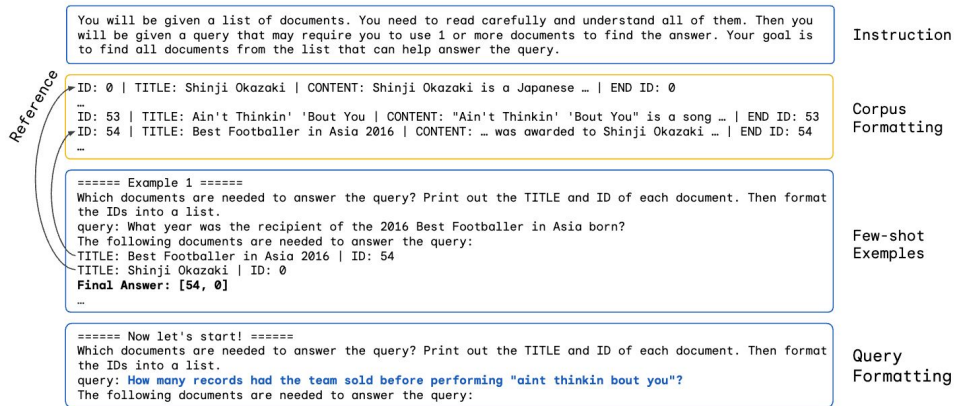


Figure 1: Accuracy of BEIR, monoTS (zero-shot) on passage retrieval (Mr.TyDi), BEIR, and Mr.TyDi) on previous best-supervised monoTS (Nogueira et al., 2020).

Challenges



- Inefficient with ranking list size $O(N^2)$ attention
- Positional bias of in-context documents - lost in the middle!
- Interpretability - does the model retrieve the right information?

BlockRank

Can we use the special structure in an ICR task for efficient training and inference

Agenda



01

IR in the age of LLMs

02

ICR Attention Analysis

03

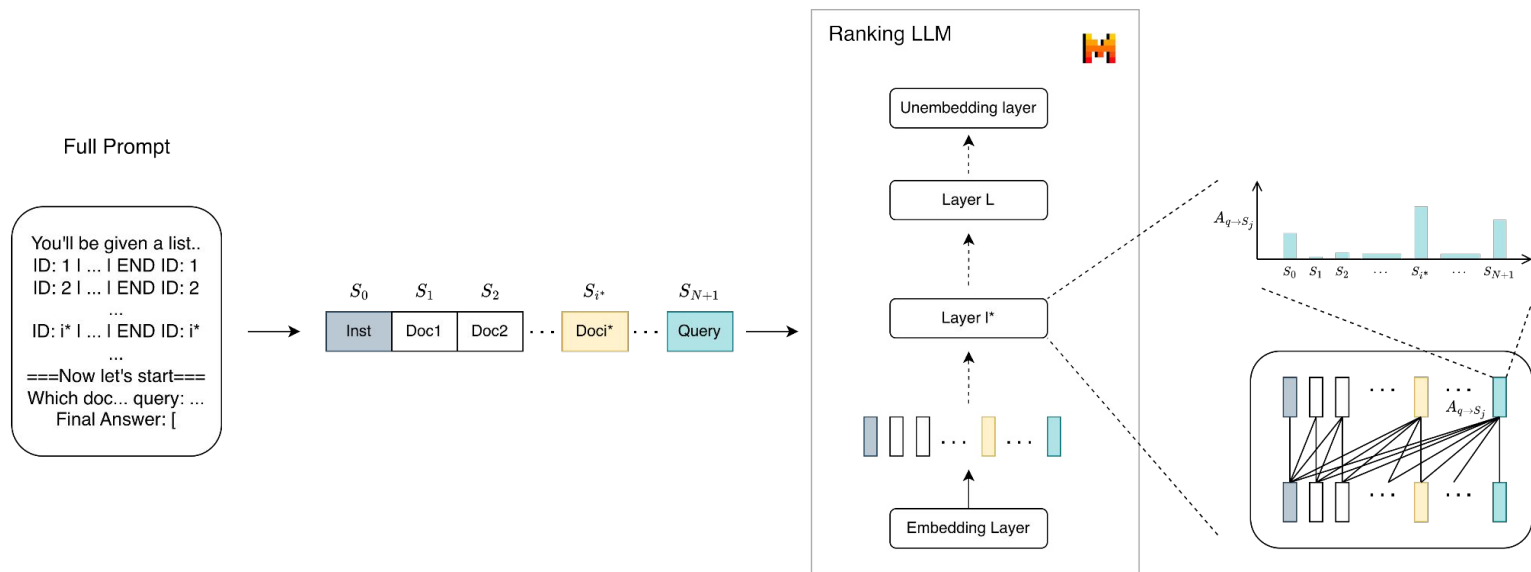
BlockRank - scalable ICR

04

Future Directions

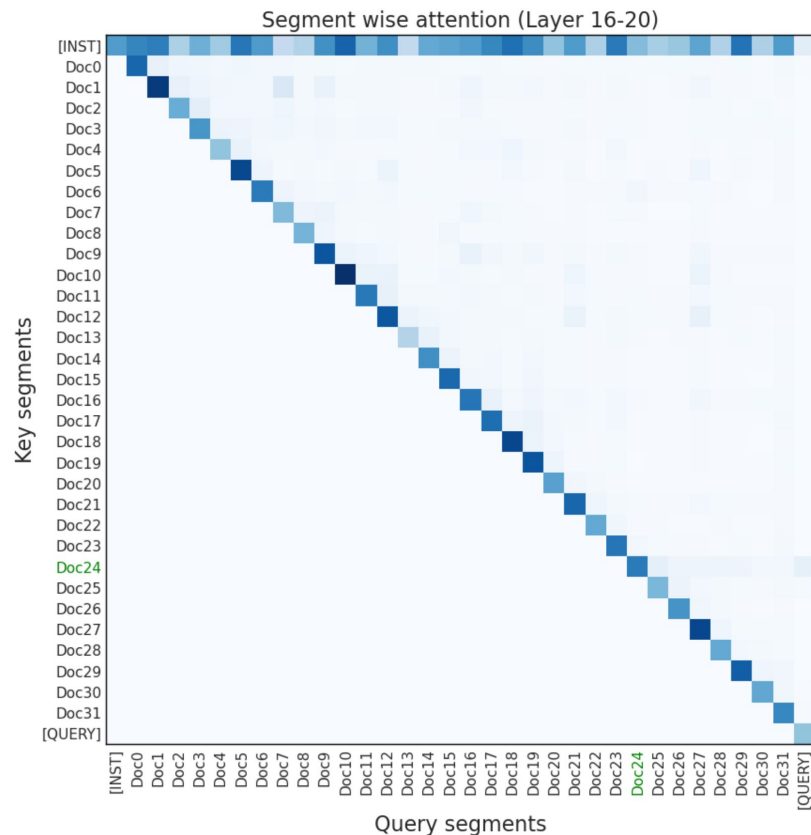
Analysis Setup

We analyze blockwise attention patterns inside a ranking LLM performing ICR

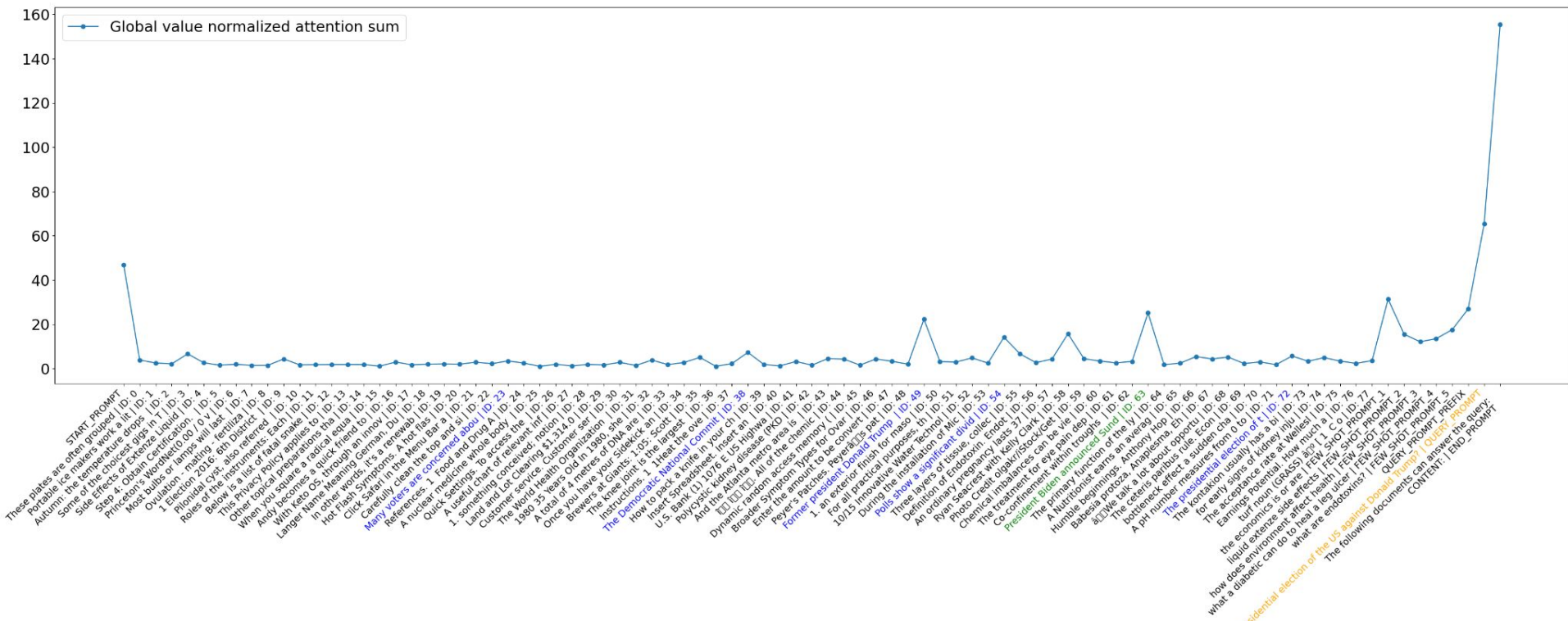


Inter-document Block Sparsity

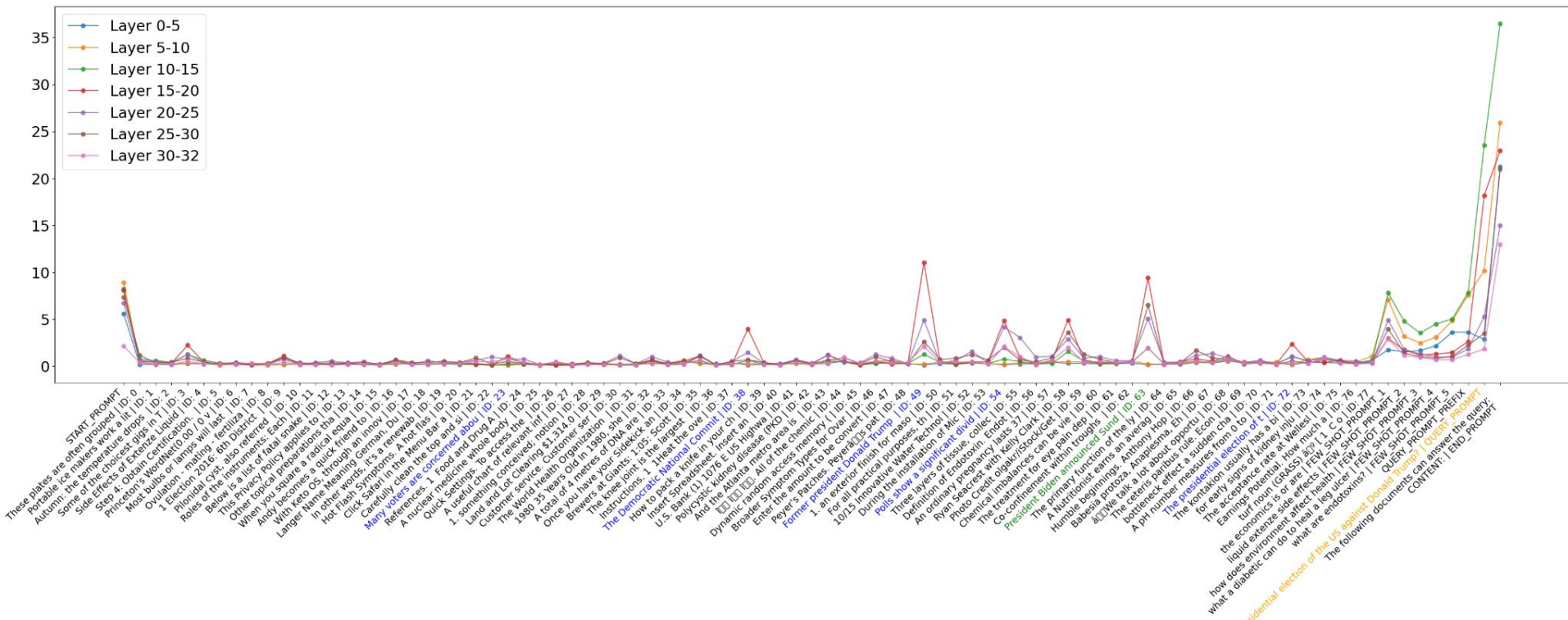
- Attention has structured sparsity
- Documents tend to attend to themselves and instruction



Query-document Block Relevance



Query-document Block Relevance



Query-document Block Relevance

- Certain tokens (":", "[") have strong correlation between attention concentration and relevance
- With SFT on ICR data this correlation improves significantly

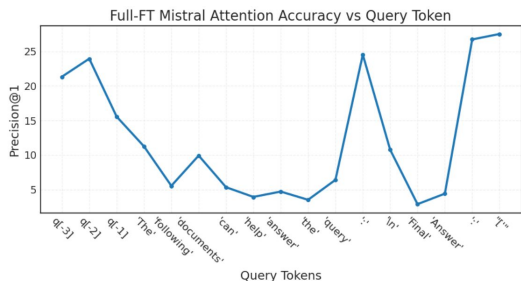


Figure 5 | Performance of Full-FT model's attention-based inference vs the query token for which attention scores are extracted from.

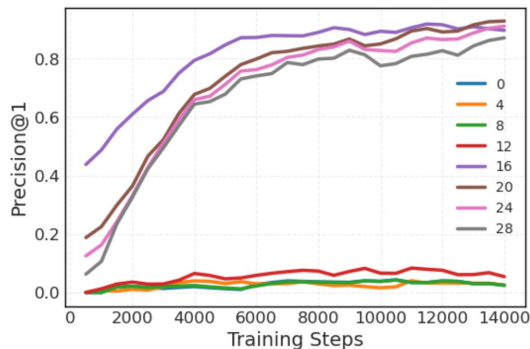


Figure 7 | Layerwise Attention Precision@1 on a held-out subset of MSMarco training data vs training steps for Full-FT model

Agenda



01

IR in the age of LLMs

02

ICR Attention Analysis

03

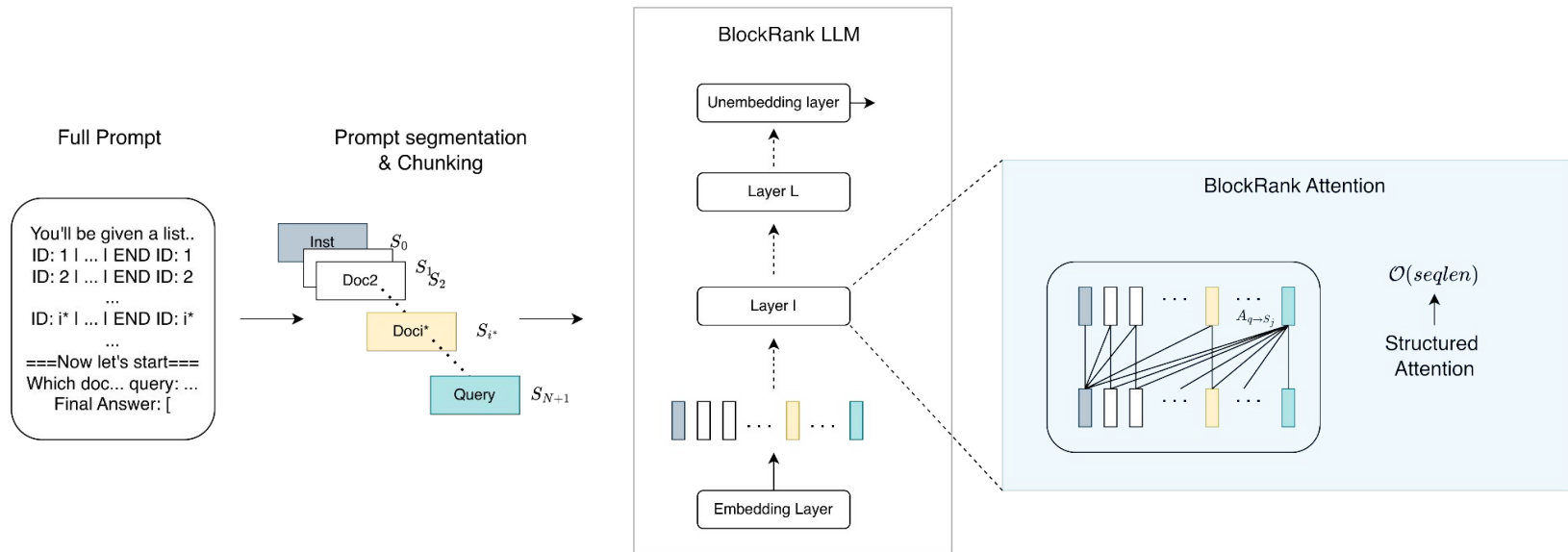
BlockRank - scalable ICR

04

Future Directions

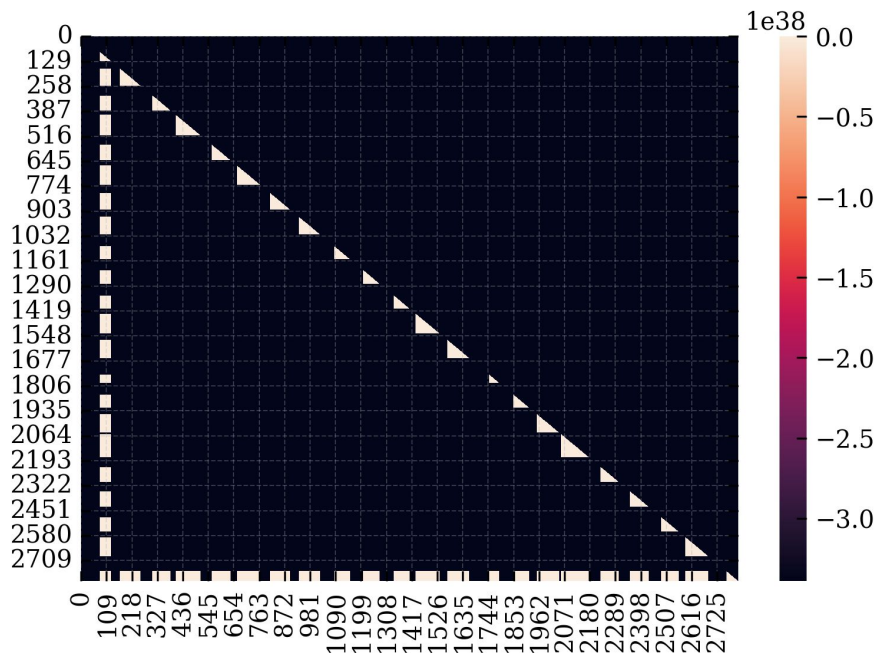
Attention Architecture

Enforce Inter-document sparsity - doc blocks only attend to self + instruction block



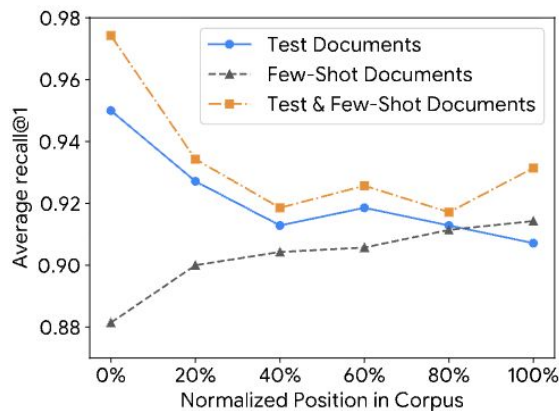
Attention Architecture

Resulting effective attention mask ~ triangular boundary



Permutation-invariant Position Encoding

“lost-in-the-middle” problem
in ICR



BlockRank allows/employs permutation
invariant position encodings

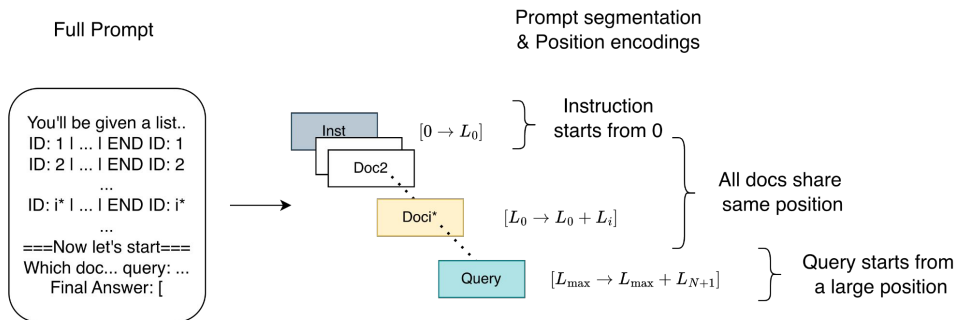
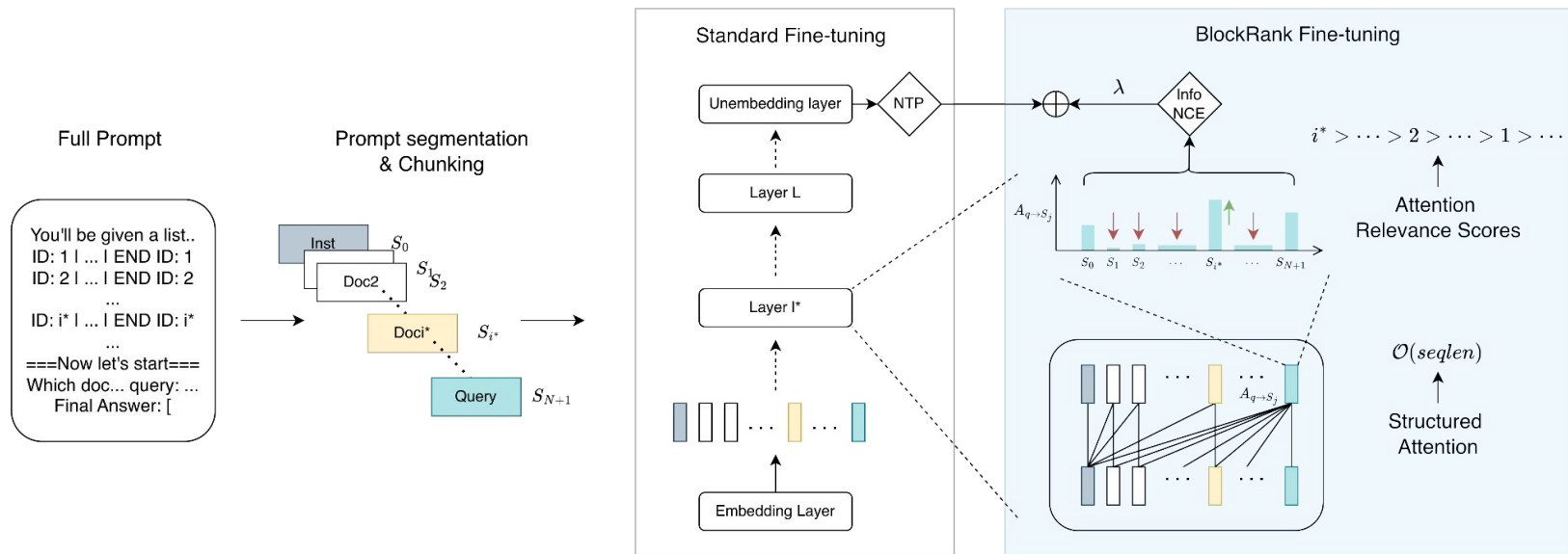


Figure 4: Positional Analysis.
We vary gold document positions
of queries within the corpus (0% =
beginning, 100% = end).

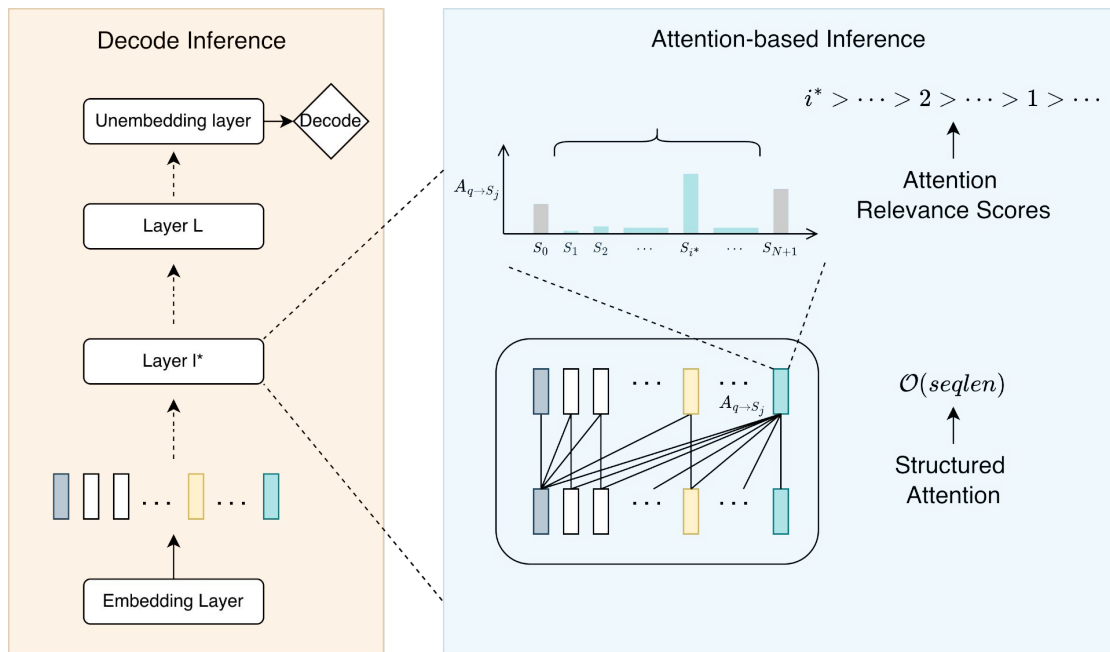
Loss

Explicitly optimize query-document block attention using a contrastive loss



Inference

Two modes of inference: (1) decode answer; (2) measure attention concentration



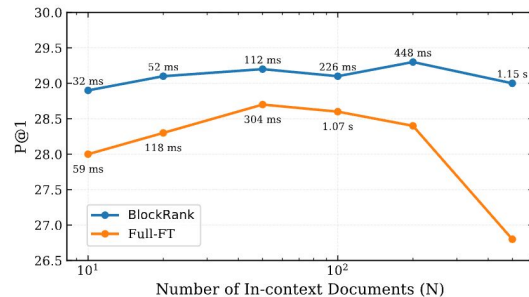
Results on BEIR

BlockRank achieves strong results while being much efficient and scalable

Table 1 | nDCG@10 on BEIR benchmark, all re-ranker rank top-100 documents retrieved from Contriever retrieval model. **Bold** indicates the best numbers.

Reranker	Train Data	Avg.	Climate-FEVER	DB-Pedia	FEVER	FiQA	Hotpot QA	MS Marco	NF-Corpus	NQ	Sci-docs	Sci-fact	Trec-COVID
None (Contriever)	MS Marco	45.9	23.7	41.3	75.8	32.9	63.8	40.7	32.8	49.8	16.5	67.7	59.6
Cross-Encoder	MS Marco	50.7	25.5	47.0	81.9	35.6	71.8	47.0	34.5	57.6	17.0	69.1	71.0
Rank Vicuna	GPT 3.5	50.7	28.2	50.0	81.0	35.9	73.5	36.7	33.1	58.6	18.4	70.5	71.3
Rank Zephyr	GPT 3.5+4	53.7	25.6	50.0	80.1	42.2	71.6	42.7	37.7	65.6	20.5	76.7	78.4
FIRST	GPT-4	54.3	26.7	50.9	81.7	42.2	74.2	44.4	37.4	66.4	20.4	74.6	78.8
BlockRank Mistral	MS Marco	54.8	26.8	49.7	87.3	44.9	75.5	48.6	36.6	62.4	18.7	76.5	76.2

Figure 4 | P@1 and Latency (annotated) of BlockRank vs Full-FT Mistral, scaling N on MSMarco.



Ablation

- Block sparse attention doesn't hurt decoding performance
- Contrastive loss important for attention-based inference
- Attention-based inference is better for inferring multiple predictions

Table 3 | Impact of training loss on Attention-based (Attn) and Decoding (Decode) Inference.

Training Configuration	Precision@1	
	Decode	Attn
Full-FT	28.7	27.6
Full-FT (w/ aux)	28.7	28.1
BlockRank (w/o ntp)	15.8	28.6
BlockRank (w/o aux)	28.4	27.8
BlockRank (full)	28.7	29.1

Table 4 | Ablation: Inference Method Effectiveness & Latency (MSMarco, N=50).

Model	Inference Method	P@1	MRR@10
Full-FT	Decode	28.7	38.4
Full-FT	Attn	27.6	38.8
BlockRank	Decode	28.7	40.0
BlockRank	Attn	29.1	42.0

Agenda



01

IR in the age of LLMs

02

ICR Attention Analysis

03

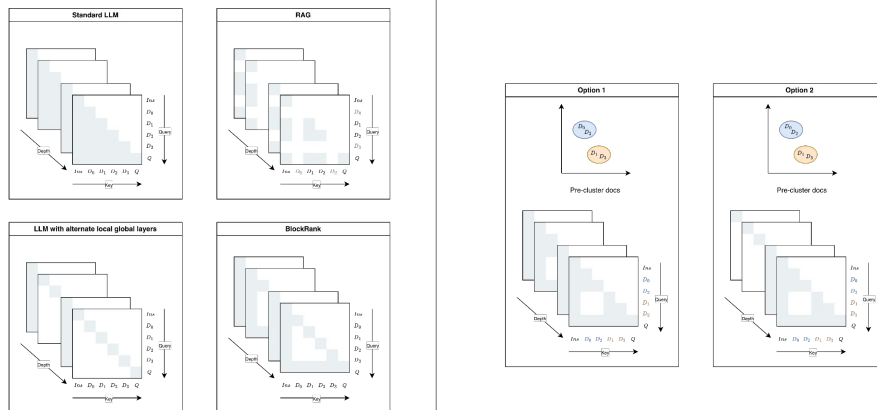
BlockRank - scalable ICR

04

Future Directions

Better Modeling?

- Clustered Block Attention
- Memory Layers
- Landmark Attention



NOTE: LLMs with memory layers are also variant of this

Pre/mid-training tasks for ICR?

Can we structure web pre-training data in a way that improve ICR quality?

A progressive end-to-end approach?

Use initial few layers to do retrieval, next few layers to do ranking and response generation while reusing the intermediate representations of the retrieval layers.

Thank You!

Arxiv - [2510.0596](#),

Github - [nilesh2797/BlockRank](#)

Reach out - nilesh@cs.utexas.edu

