

# RE-THINKING RE-RANKING

Sean MacAvaney  
University of Glasgow

Presented at:  
Haystack EU 2024

Terrier 



University  
of Glasgow



**Sean MacAvaney**

@macavaney

University of Glasgow · Lecturer (Assistant Professor)

Conduct practical research in information retrieval:

- Learned Sparse Expansion
- Search Result Evaluation using LLMs
- Document Quality Prediction
- Adaptive Re-Ranking



**Sean MacAvaney**

@macavaney

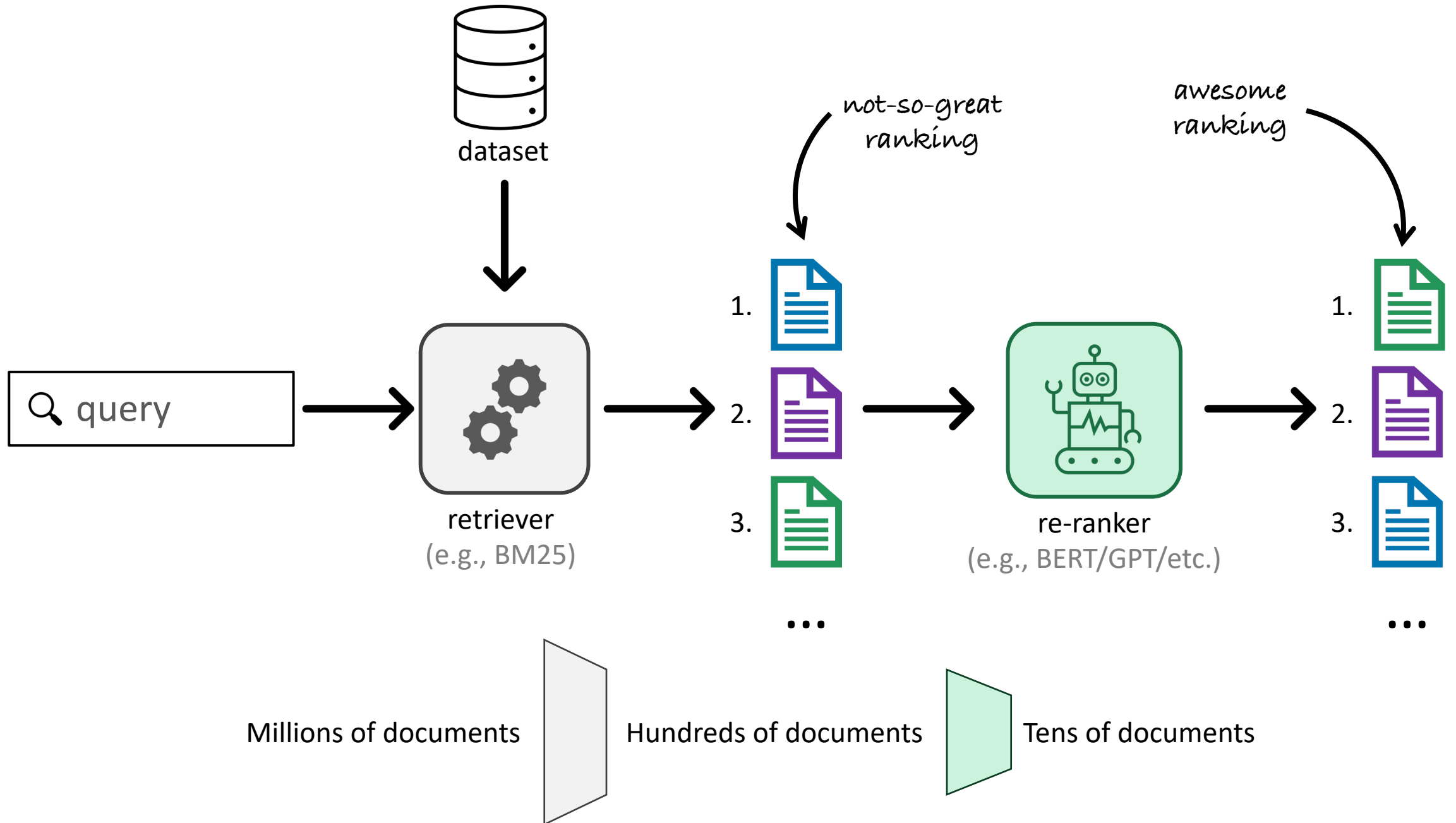
University of Glasgow · Lecturer (Assistant Professor)

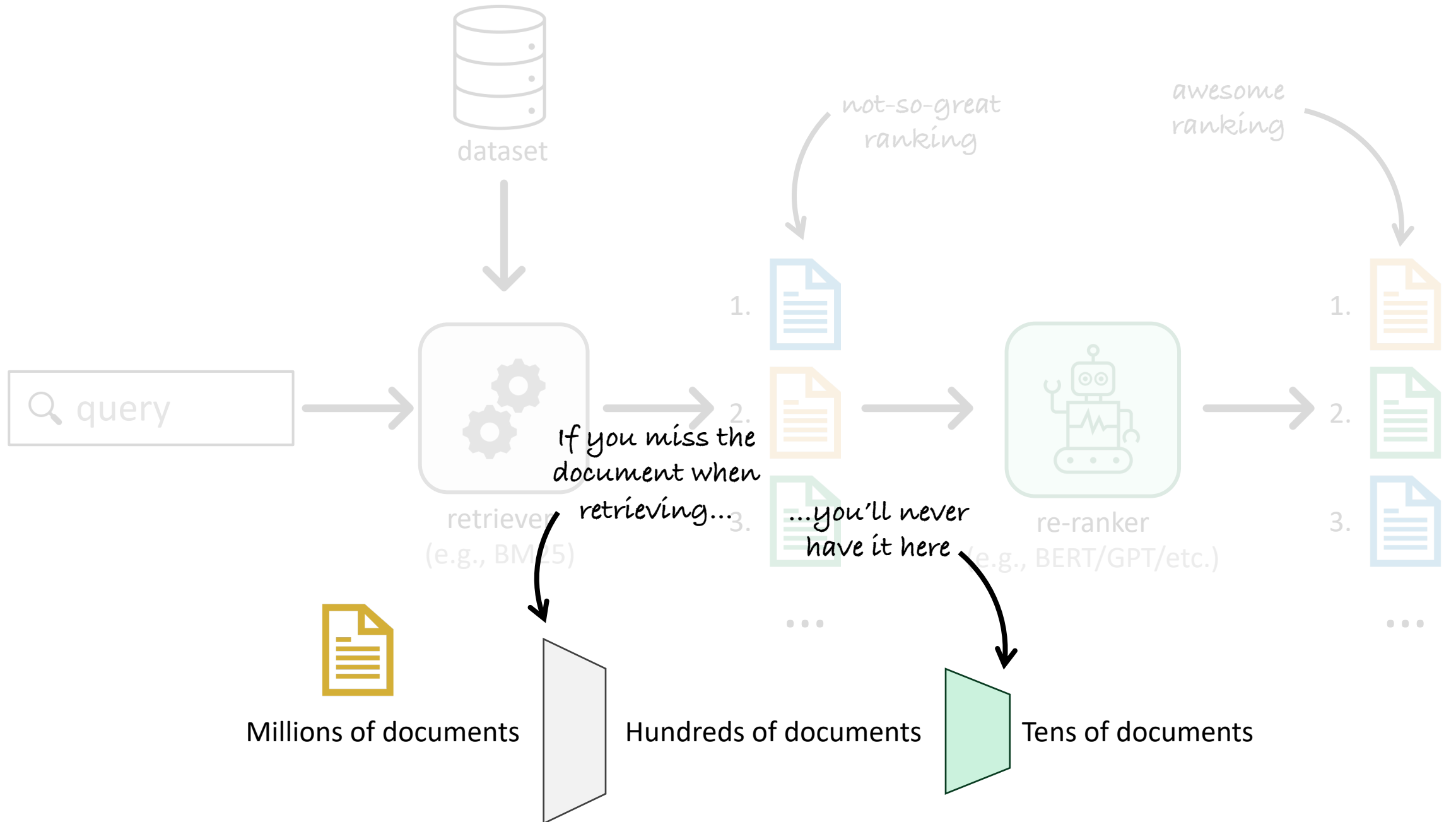
Conduct practical research in information retrieval:

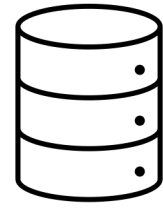
- Learned Sparse Expansion
- Search Result Evaluation using LLMs
- Document Quality Prediction

 - **Adaptive Re-Ranking**

# What's Re-Ranking?





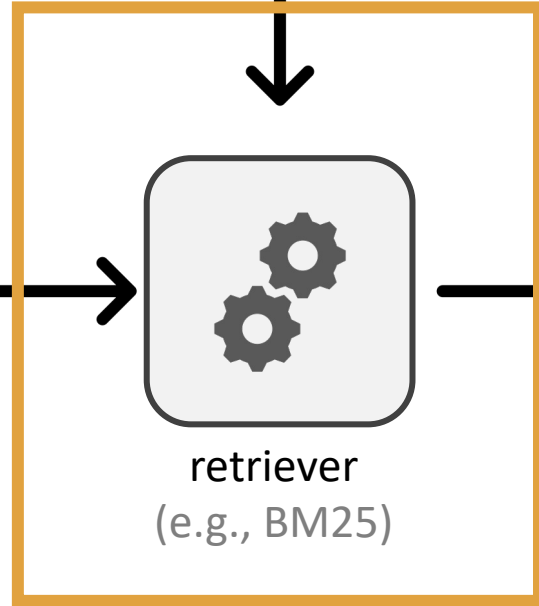


dataset

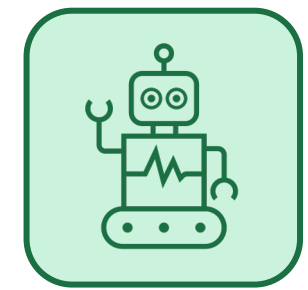
To overcome this **recall problem**,  
people typically focus on the retriever



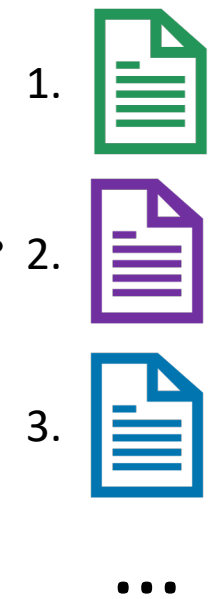
Q query



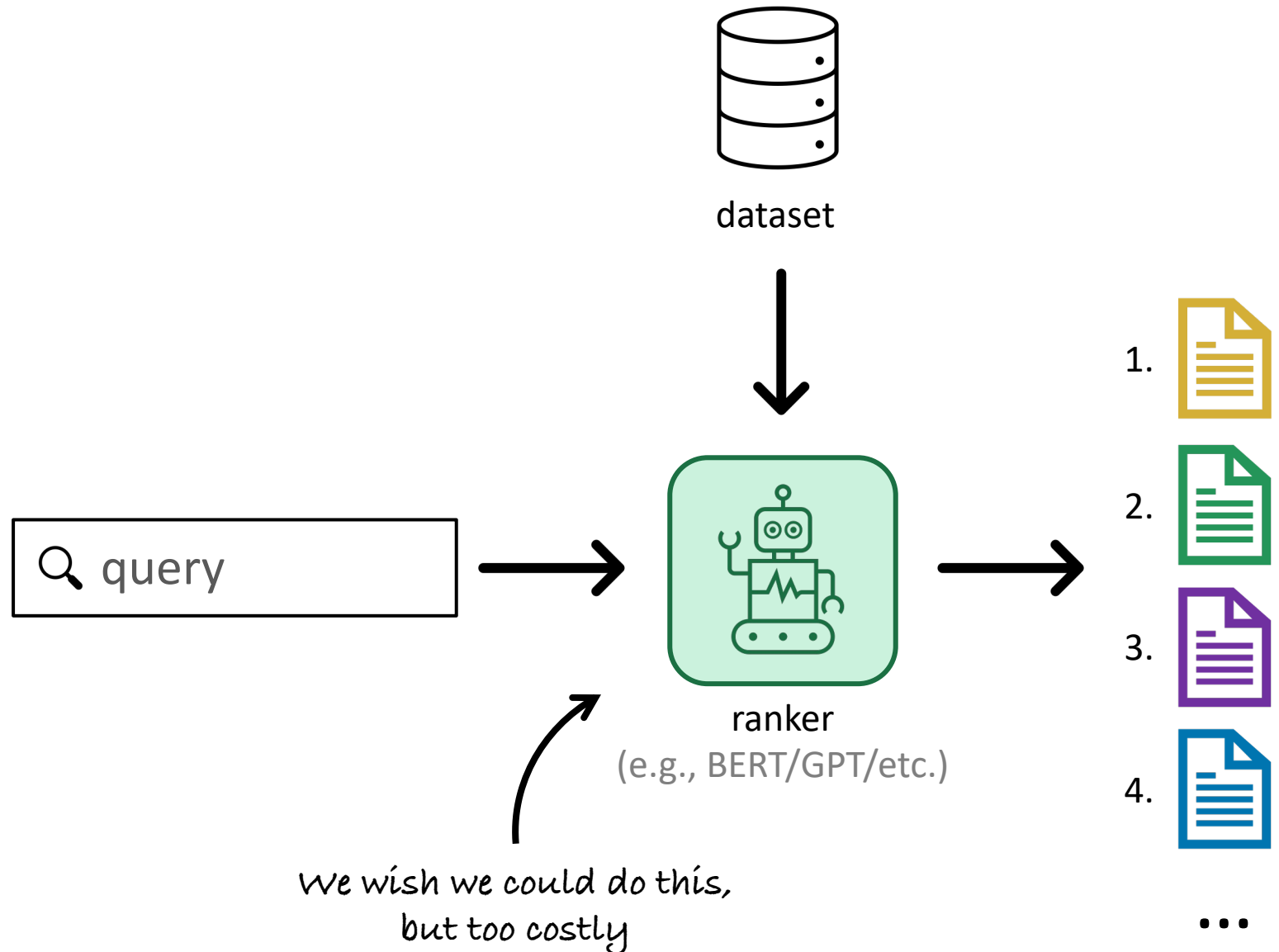
retriever  
(e.g., BM25)

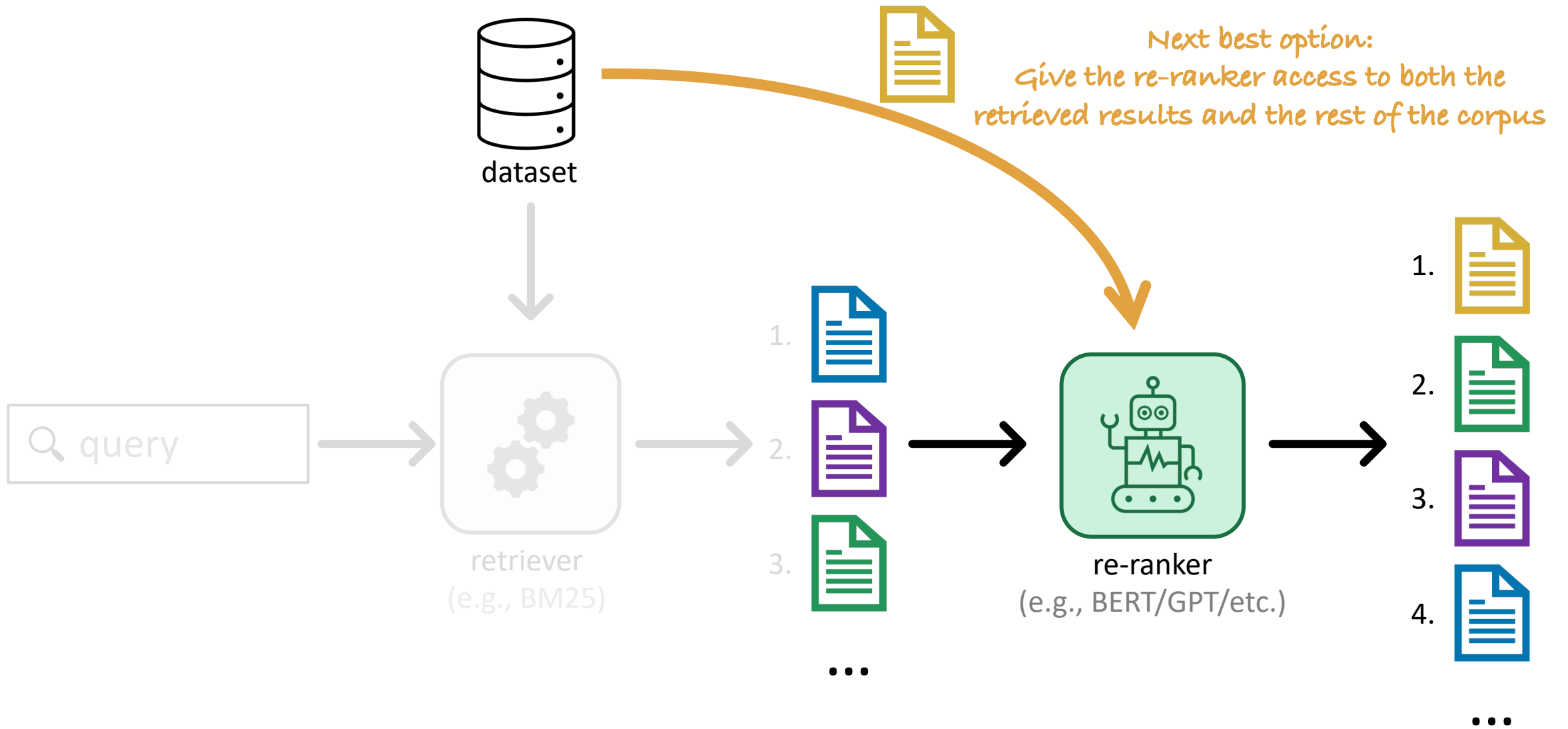


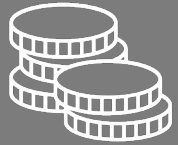
re-ranker  
(e.g., BERT/GPT/etc.)



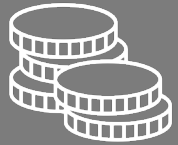
dense retrieval  
 learned sparse retrieval  
 hybrid retrieval  
 etc



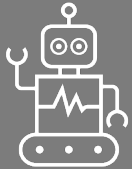




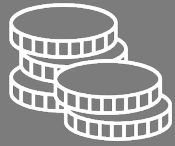
I'll show this can be done with minimal cost.



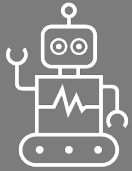
I'll show this can be done with minimal cost.



The idea can improve retrievers like ColBERT, too.



I'll show this can be done with minimal cost.



The idea can improve retrievers like ColBERT, too.



Ready-to-use with Open-Source tools!

# Battleship

**Battleship**

**Schifferversenken**


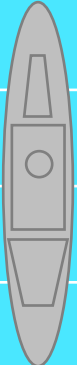

**“Sink the Fleet”**



**Your Board**


8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H

**Opponent's board (Secret)**

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H


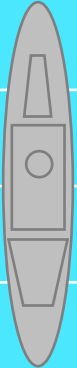

**Goal: identify the positions of all your opponent's ships**

### Your Board


8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H

*Make a guess*

### Opponent's board (Secret)

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H



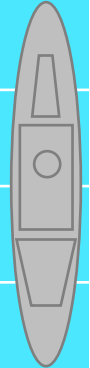
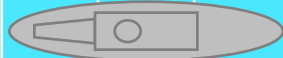
### Your Board

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H



ask opponent

### Opponent's board (Secret)

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H

### Your Board

8								
7								
6				☀				
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H





*it's a miss!*





### Opponent's board (Secret)

8								
7								☐
6								☀
5								
4								
3								
2								☐
1								
	A	B	C	D	E	F	G	H

### Your Board

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H



### Opponent's board (Secret)

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H






it's a hit!

### Your Board




8								
7								
6								
5								
4								
3						?		
2					?		?	
1						?		
	A	B	C	D	E	F	G	H

Better guesses →






### Opponent's board (Secret)

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H



### Your Board

8								
7								
6								
5								
4								
3								
2					?		?	
1						?		
	A	B	C	D	E	F	G	H







### Opponent's board (Secret)

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H








### Your Board

8								
7								
6								
5								
4								
3								
2				?			?	
1								
	A	B	C	D	E	F	G	H






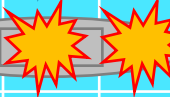
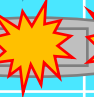
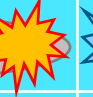

### Opponent's board (Secret)

8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H

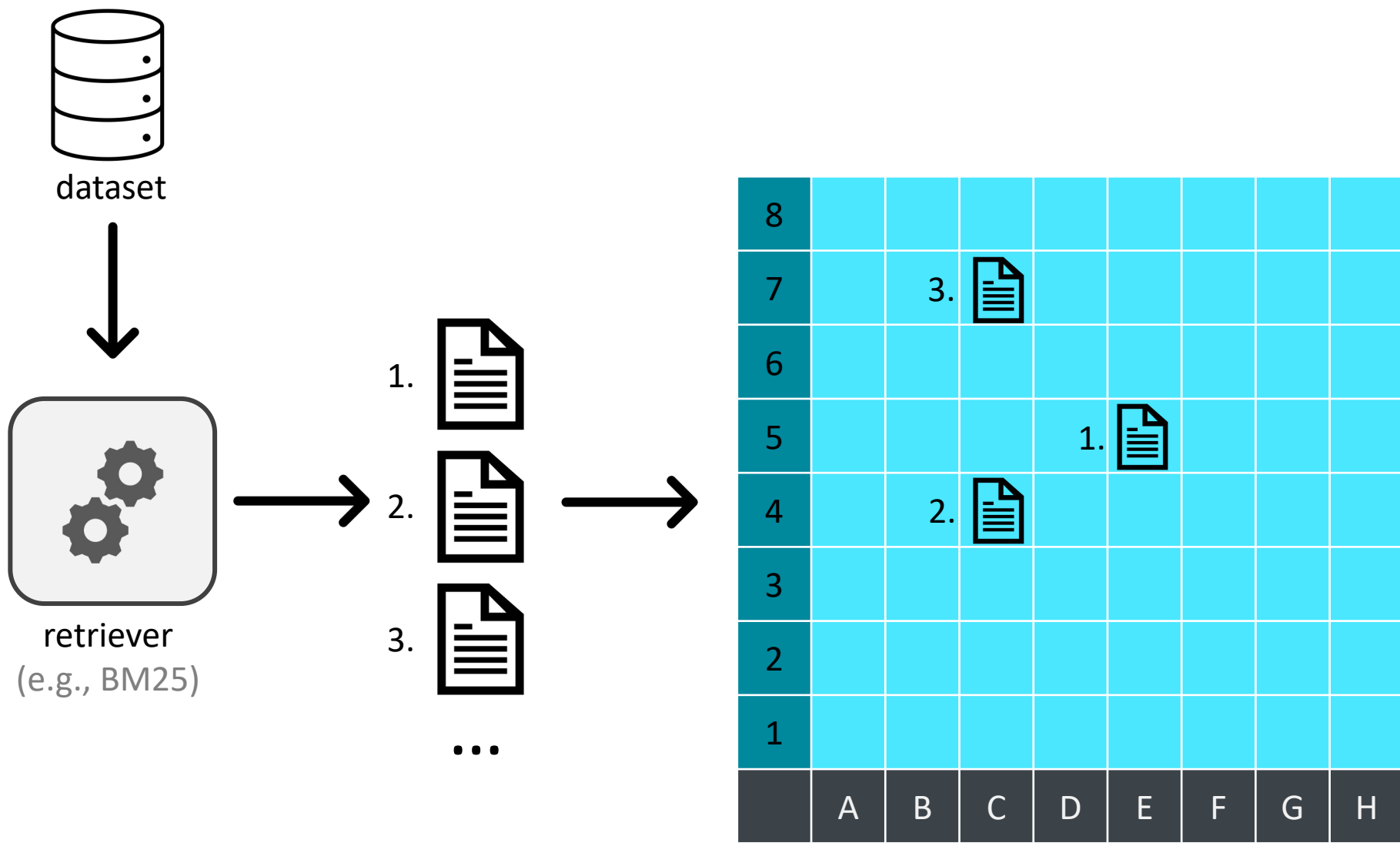
**Your Board**




8								
7								
6								
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H




**Opponent's board (Secret)**

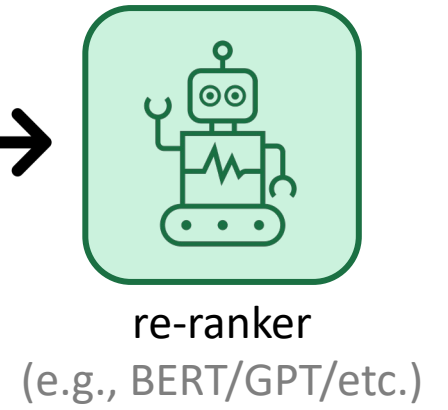
8									
7									
6									
5									
4									
3									
2									
1									
	A	B	C	D	E	F	G	H	




**And so forth...**

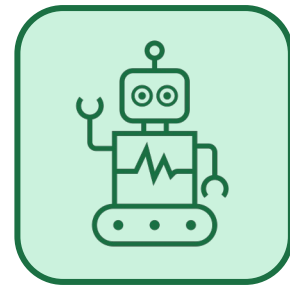


- 1. 
- 2. 
- 3. 
- ...

8								
7		3.						
6								
5				1.				
4		2.						
3								
2								
1								
	A	B	C	D	E	F	G	H



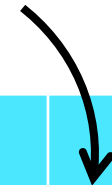
8								
7		3.						
6								
5				1.				
4		2.						
3								
2								
1								
	A	B	C	D	E	F	G	H









re-ranker  
(e.g., BERT/GPT/etc.)

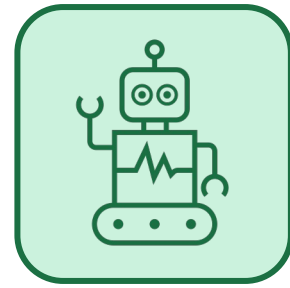


Traditional re-ranking stops here






8								
7			1.					
6								
5						3.		
4			2.					
3								
2								
1								
	A	B	C	D	E	F	G	H

8								
7		3.						
6								
5				1.				
4		2.						
3								
2								
1								
	A	B	C	D	E	F	G	H






re-ranker  
(e.g., BERT/GPT/etc.)

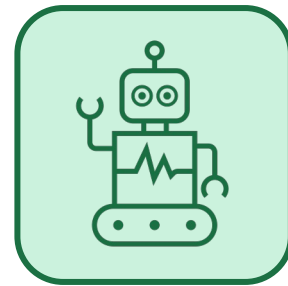


8			?					
7		?		?				
6			?					
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H

But we've learned a lot  
from the re-ranker!




**Adaptive Re-Ranking** leverages the information gained from high-scoring documents to find ones missed by the retriever.

8								
7		3.						
6								
5				1.				
4		2.						
3								
2								
1								
	A	B	C	D	E	F	G	H

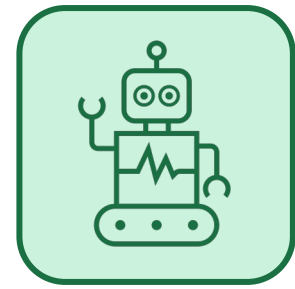
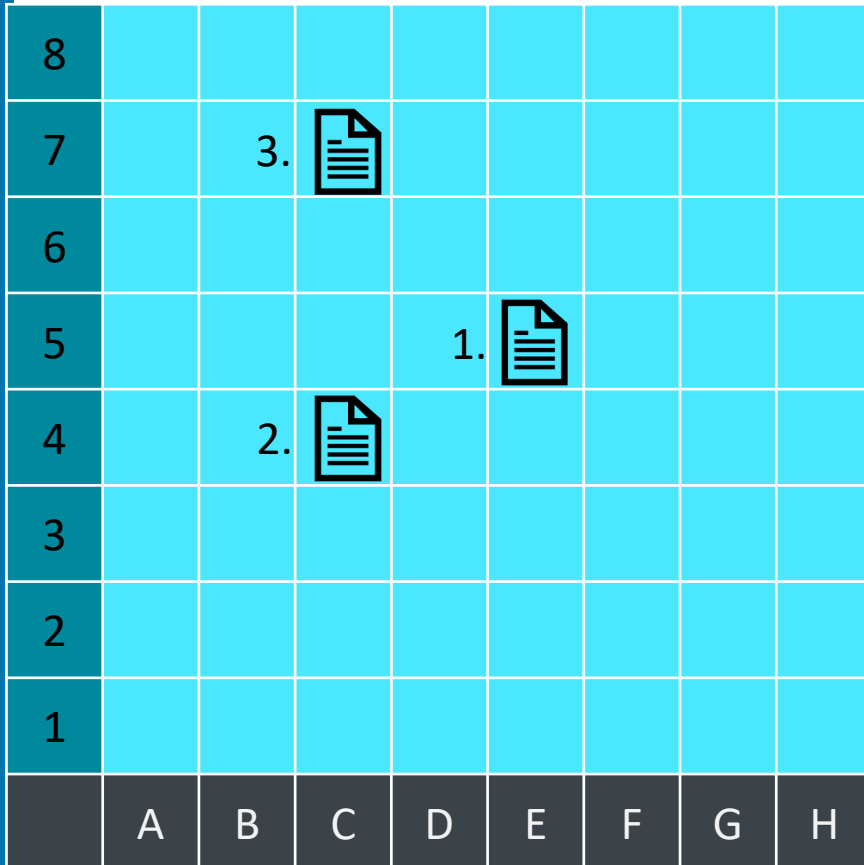


re-ranker  
(e.g., BERT/GPT/etc.)

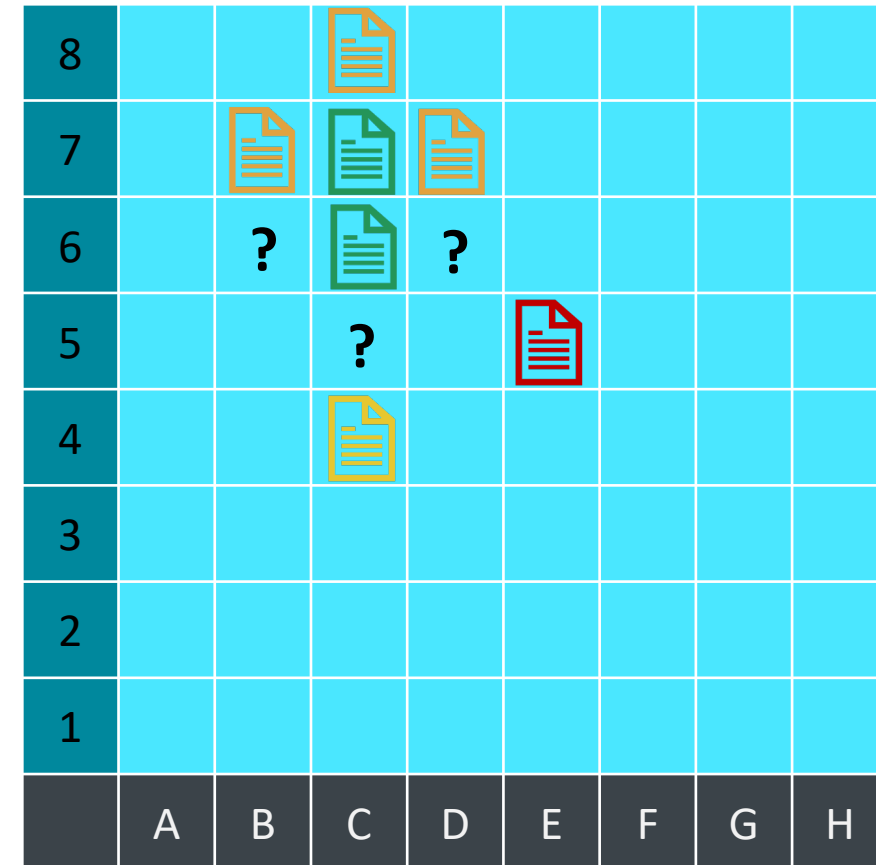


8			?					
7		?		?				
6			?					
5								
4								
3								
2								
1								
	A	B	C	D	E	F	G	H

**Adaptive Re-Ranking** leverages the information gained from high-scoring documents to find ones missed by the retriever.

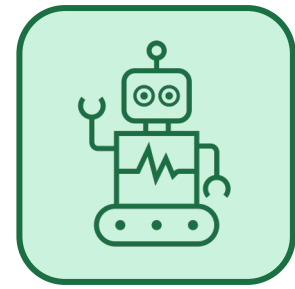
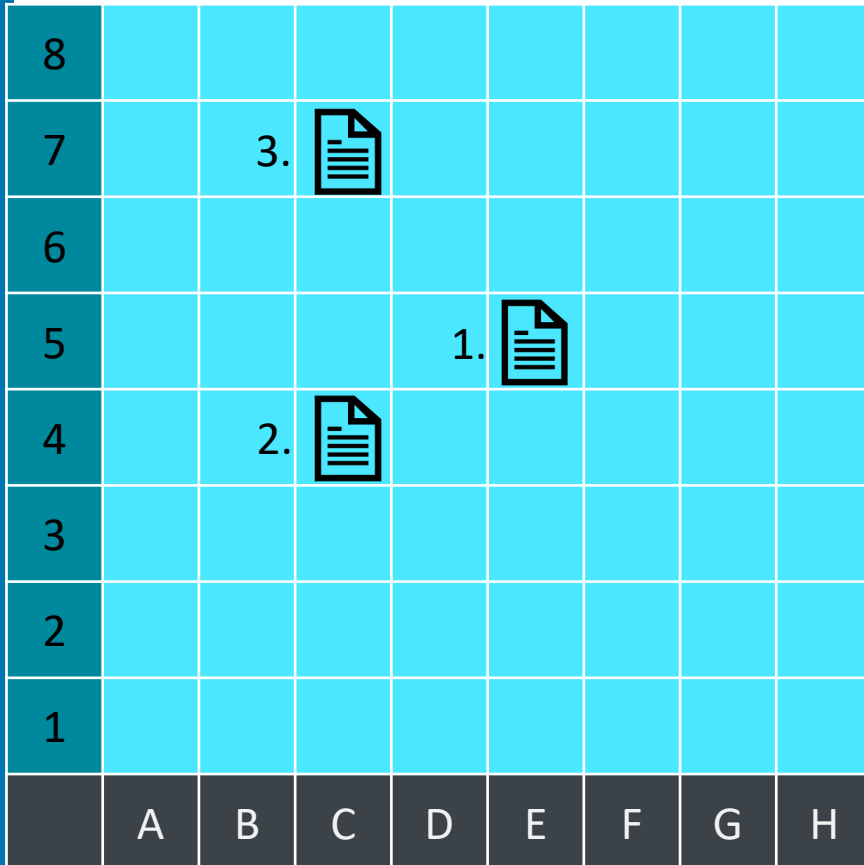


re-ranker  
(e.g., BERT/GPT/etc.)

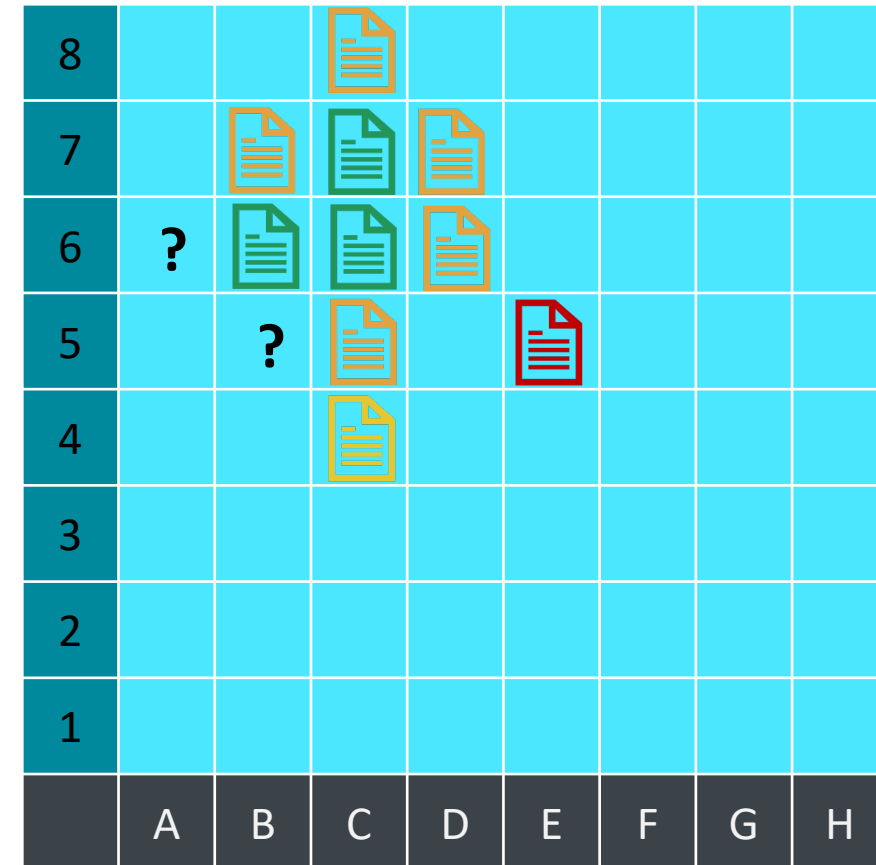


**And so forth...**

**Adaptive Re-Ranking** leverages the information gained from high-scoring documents to find ones missed by the retriever.

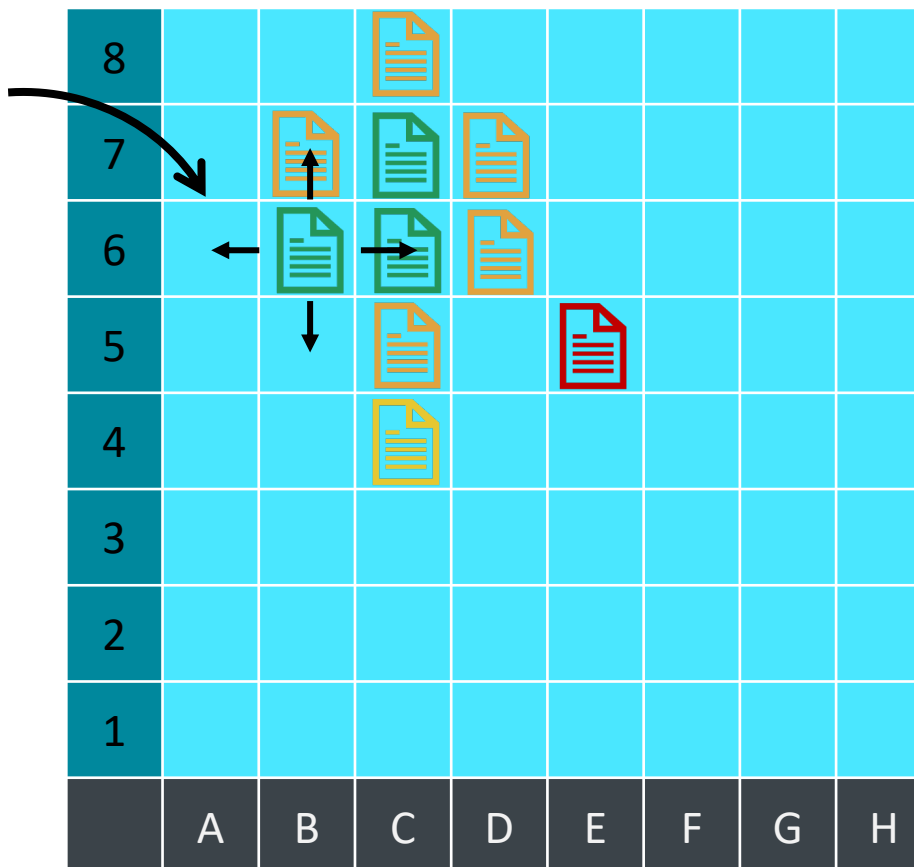


re-ranker  
(e.g., BERT/GPT/etc.)

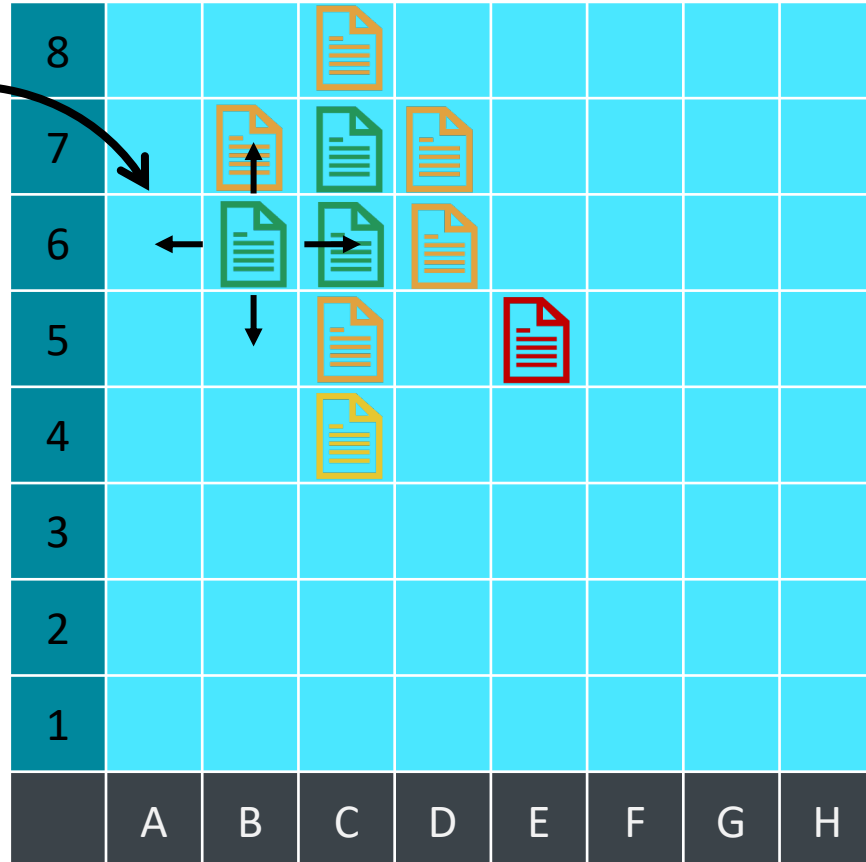


**And so forth...**

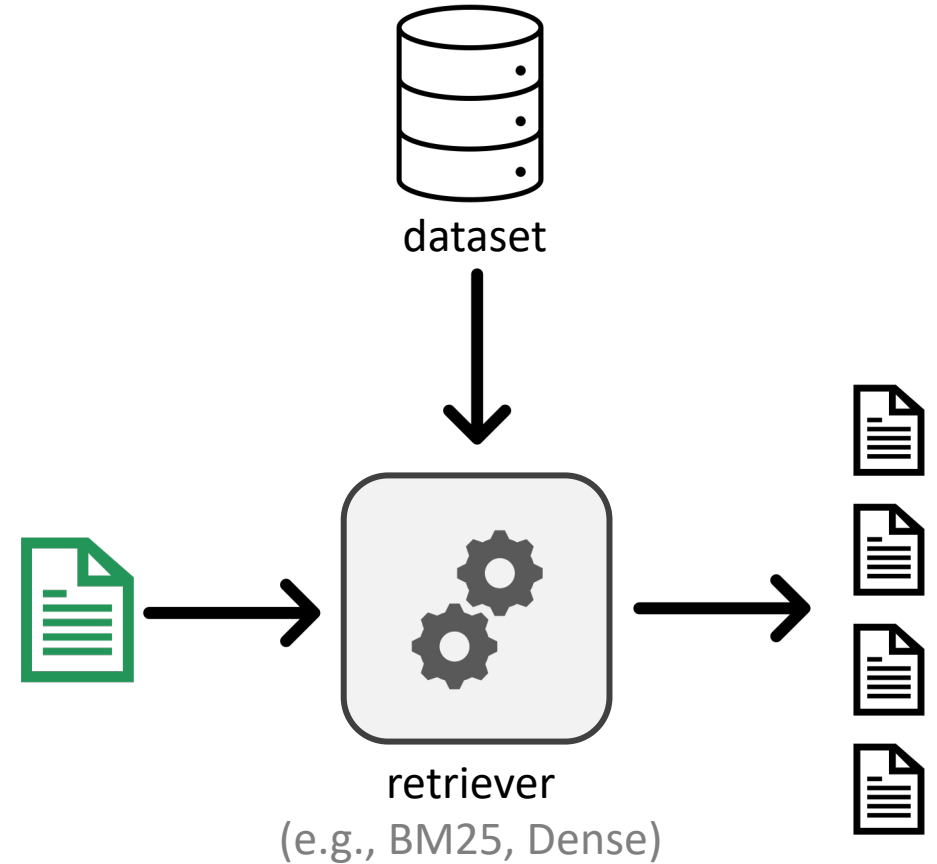
How to decide which documents to check?



How to decide which documents to check?



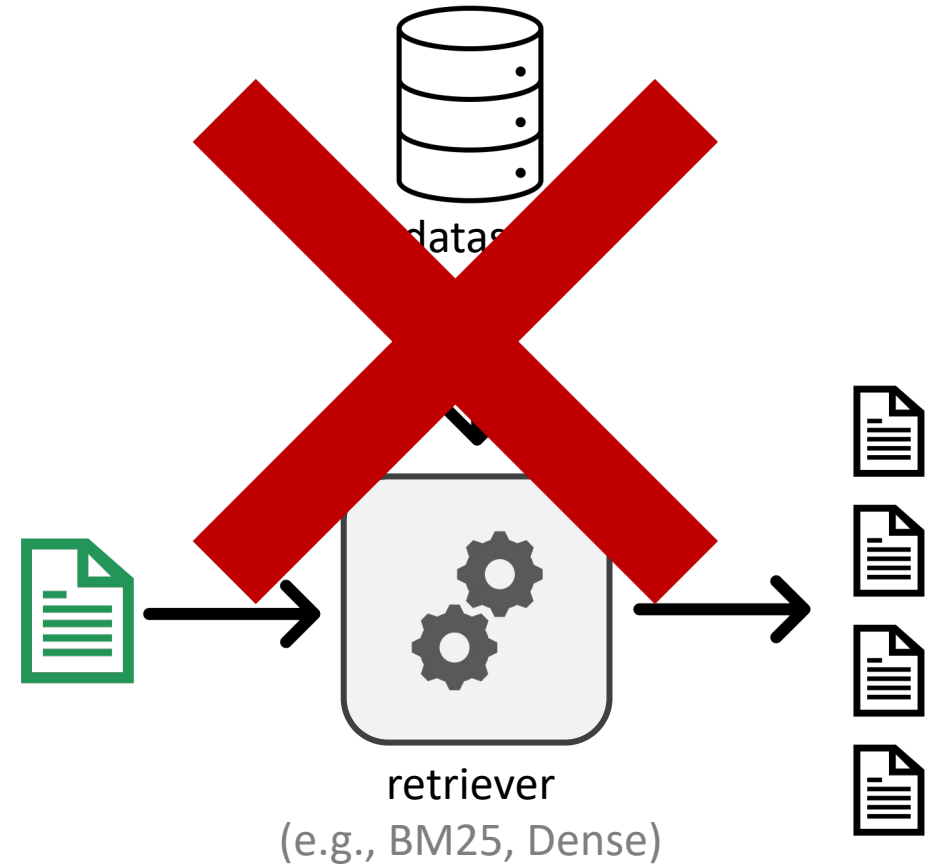
We could issue the document as a query to the engine and take the top  $k$  results.



**Really slow!**

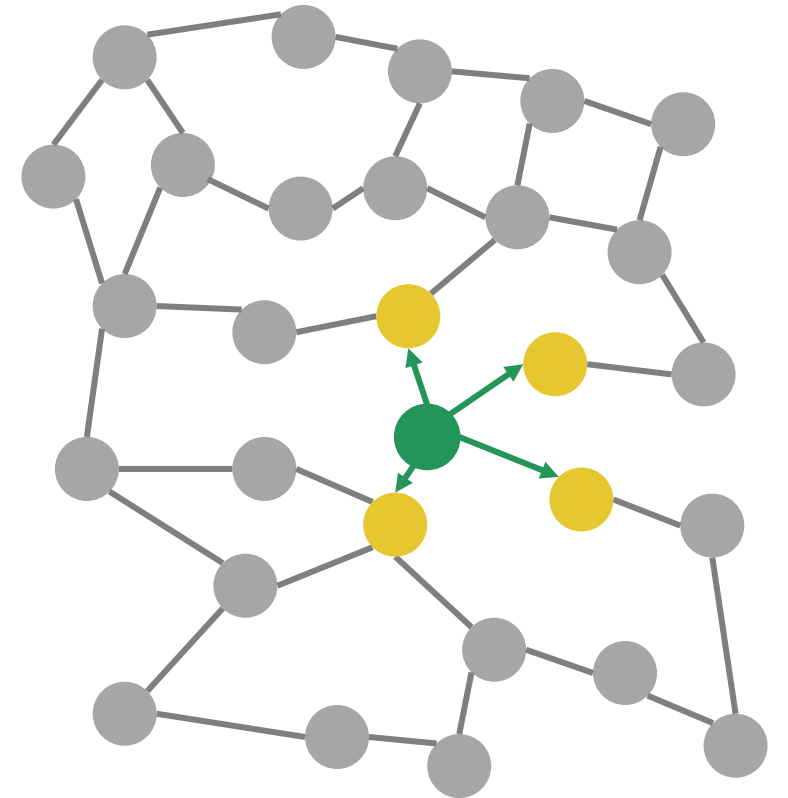
Right idea, bad execution.

We could issue the document as a query to the engine and take the top  $k$  results.



**Better: Use the KNN graph we already have for HNSW search.**

- Establishes proximity
- Fast lookups
- Constructed offline

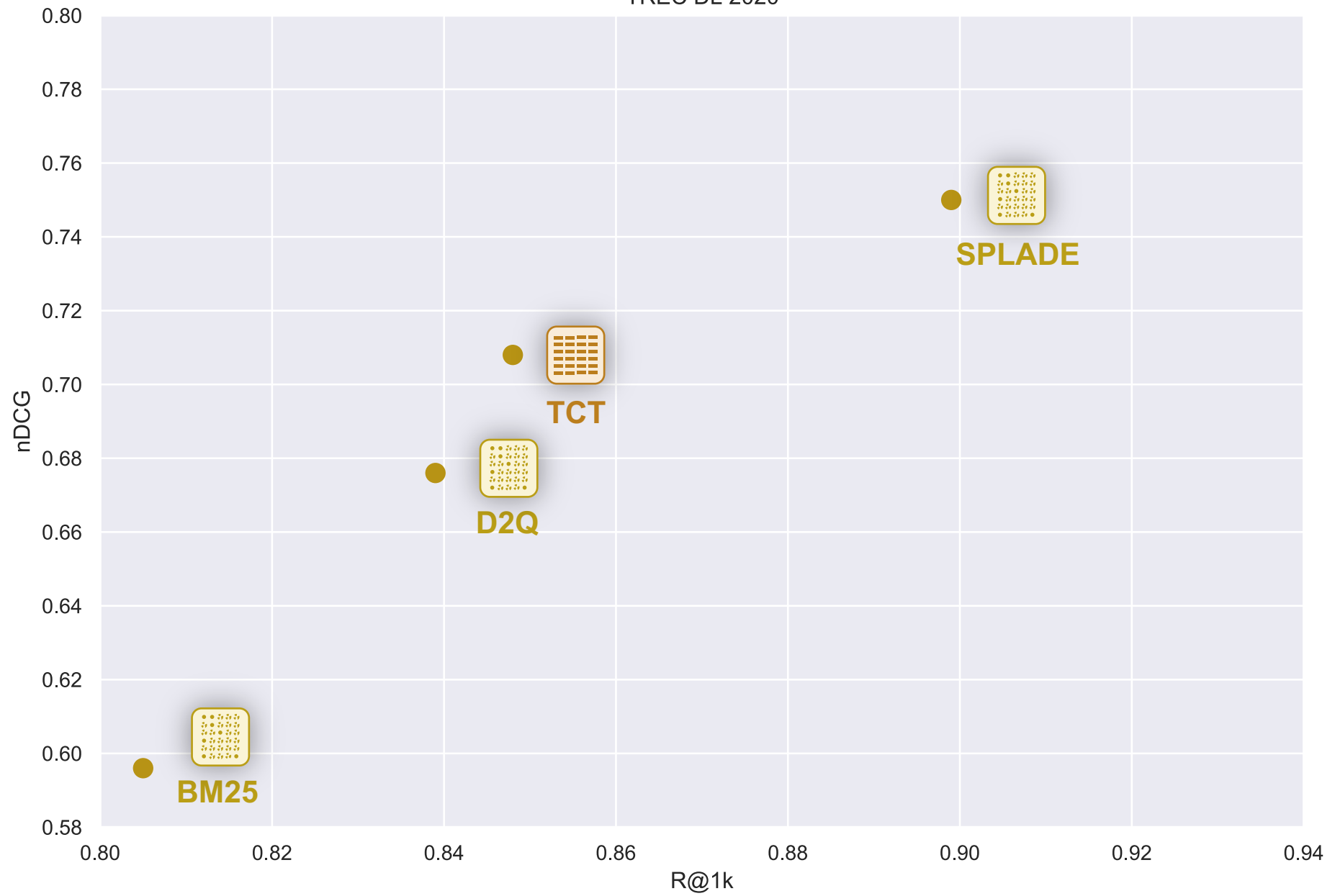


Alright, so how well does this  
*adaptive* re-ranking strategy work?

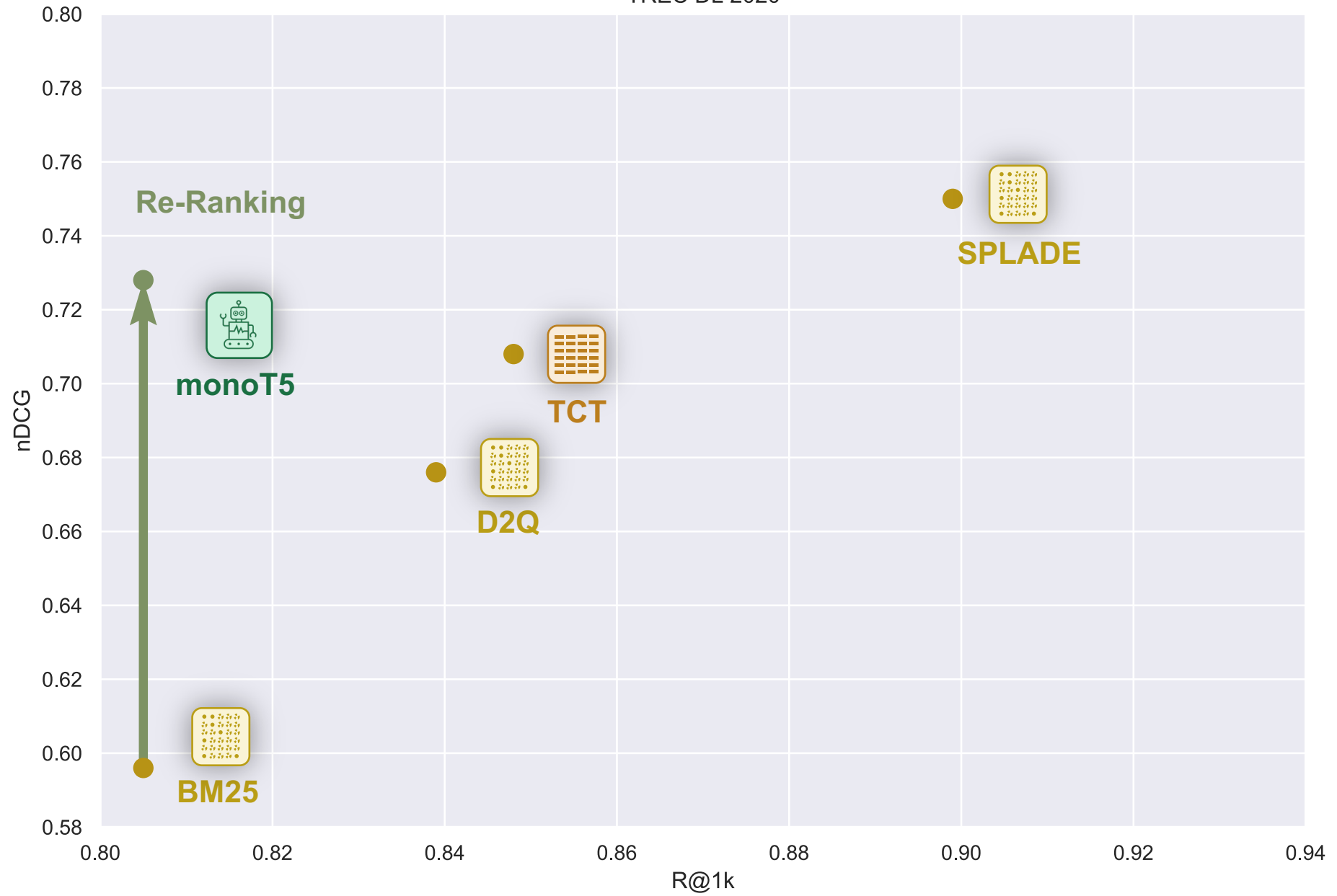
## **A few technical bits...**

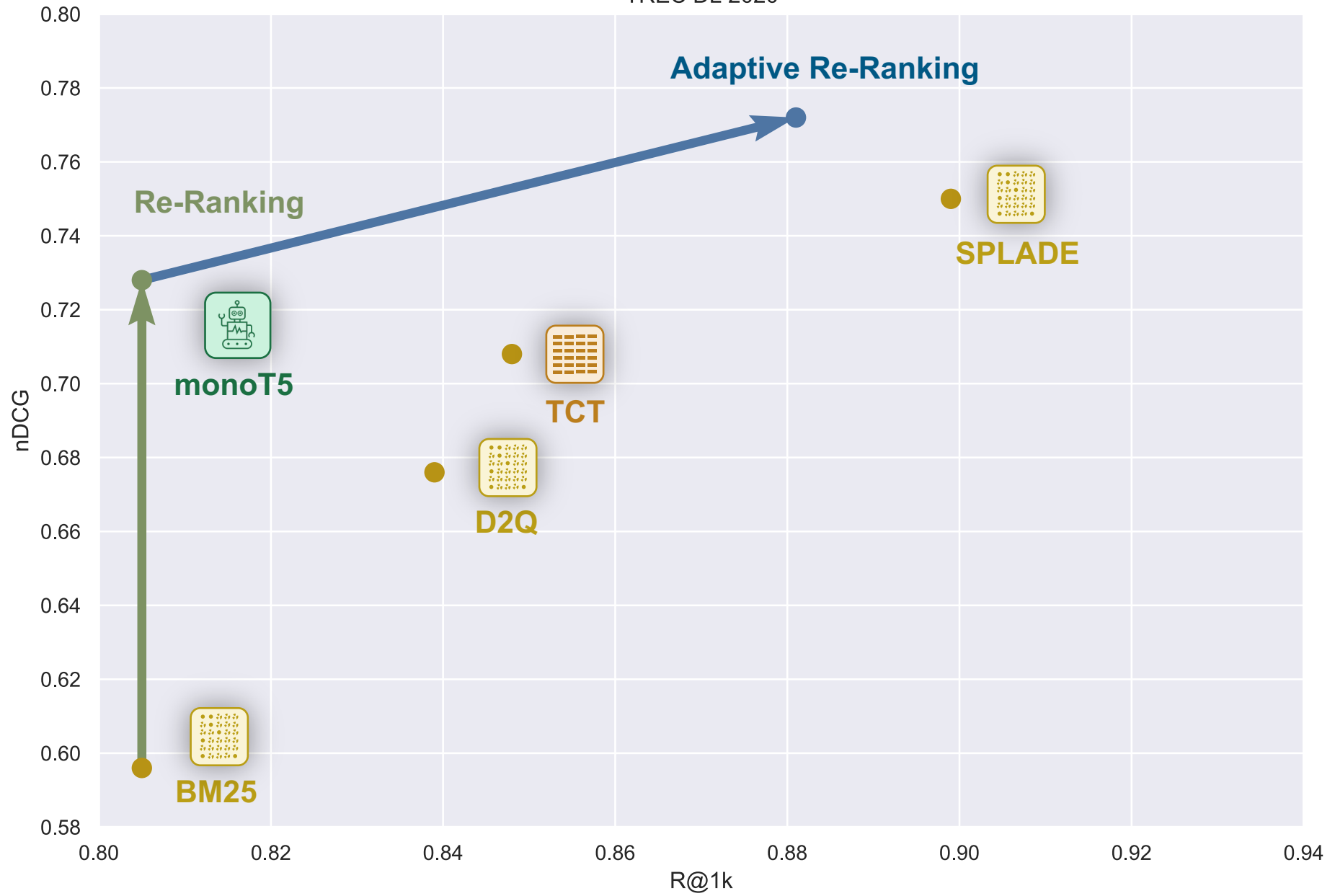
- We fix the re-ranking “budget” (number of docs to score) across all pipelines.
- In adaptive setting, we take half from first-stage ranker and half from graph.
- Measure nDCG (overall ranking quality) and Recall (% of relevant docs retrieved).
- Test on a variety of re-ranking pipelines.

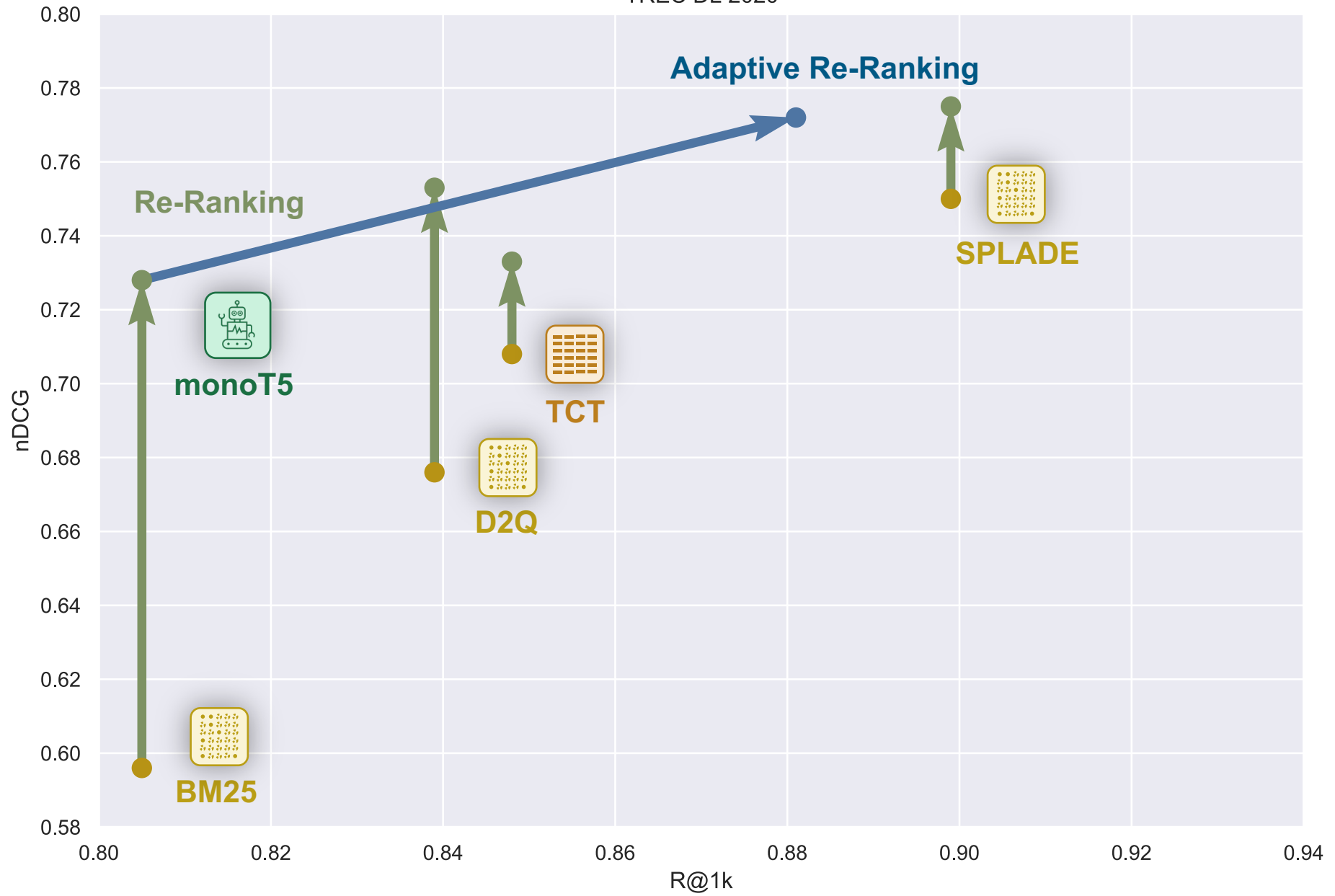
TREC DL 2020

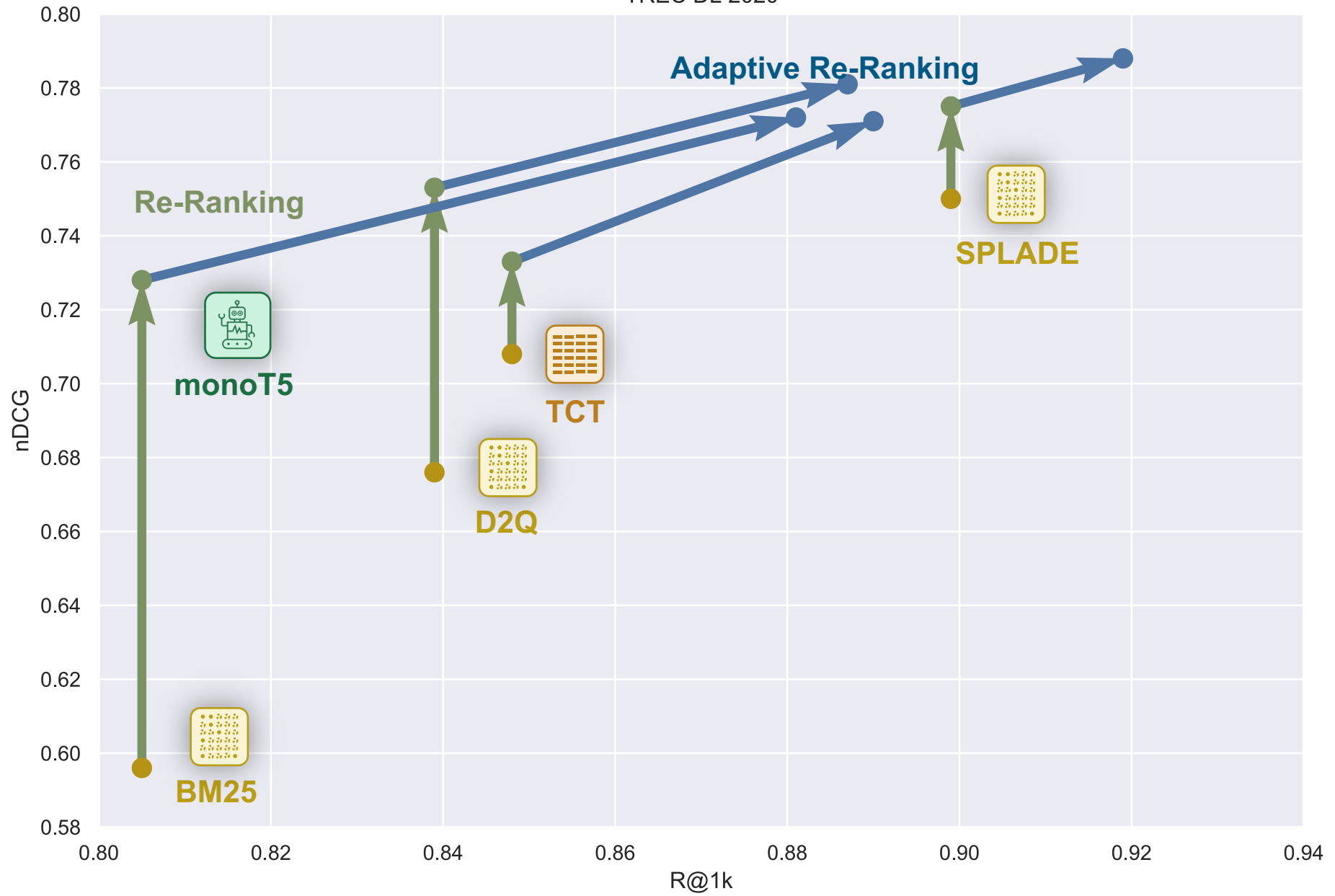


TREC DL 2020











**Sean MacAvaney**

@macavaney



I ❤️ cross-encoders! Awesome to see another one from Cohere

A quick test showing it turbocharged when using Graph-based Adaptive Reranking :)

```
dataset = pt.get_dataset('irds:msmarco-passage/trec-dl-2019/judged')

pt.Experiment(
  [
    bm25,
    bm25 >> cohere_rerank,
    bm25 >> GAR(cohere_rerank, graph, num_results=100),
  ],
  dataset.get_topics(),
  dataset.get_qrels(),
  [nDCG@10, nDCG, R(rel=2)@100]
)

#      name  nDCG@10  nDCG  R(rel=2)@100
#      BM25      0.499  0.459      0.497
#  BM25 >> Rerank  0.708  0.523      0.497
#  BM25 >> GAR(Rerank) 0.753  0.604      0.609
```

ALT

# In summary

Nearest neighbor graph exploration helps re-rankers

Focusing on the neighbors of the top scored documents helps prioritize the documents that are most likely to be relevant

## Other findings

- Even works when no relevant documents returned by the first stage
- Ineffective when there are lots of duplicate documents (lesson learned from TREC DL 2022)
- Robust to various measures of document similarity (semantic or lexical)
- Room to improve by finding smarter ways to prioritize (but might be costly)

## Conference Paper:

MacAvaney, Tonello, Macdonald. Adaptive Re-Ranking with a Corpus Graph. CIKM 2022.

# Retrievers

Adaptive Re-Ranking involved **two** strategies:

- ① Score docs near the best ones found so far.

**Adaptive Re-Ranking** involved **two** strategies:

- ① Score docs near the best ones found so far.

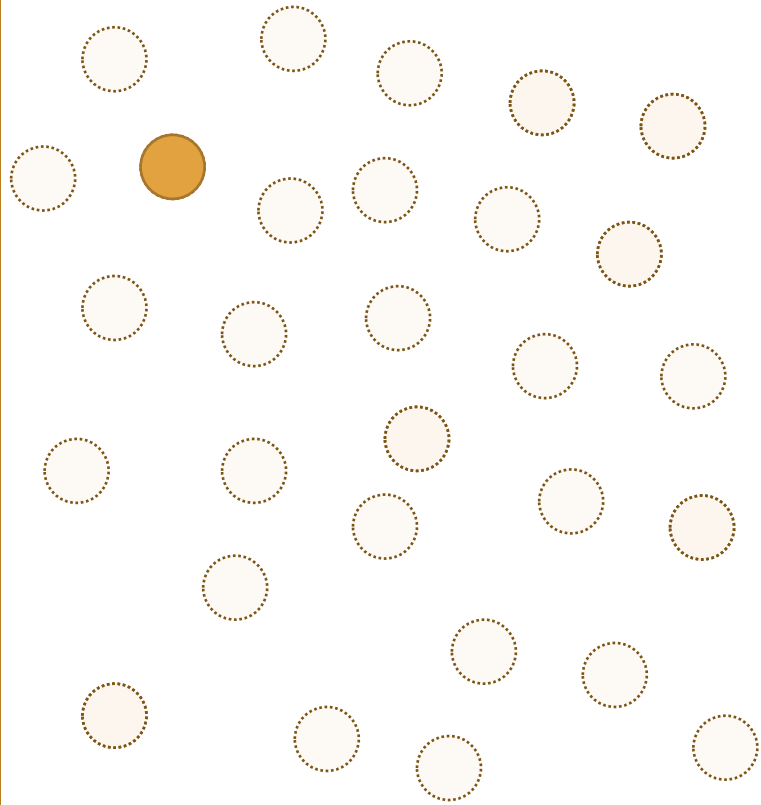
Adaptive **Re-Ranking** involved **two** strategies:

- ① Score docs near the best ones found so far.
- ② Start with good (but cheap) guesses.

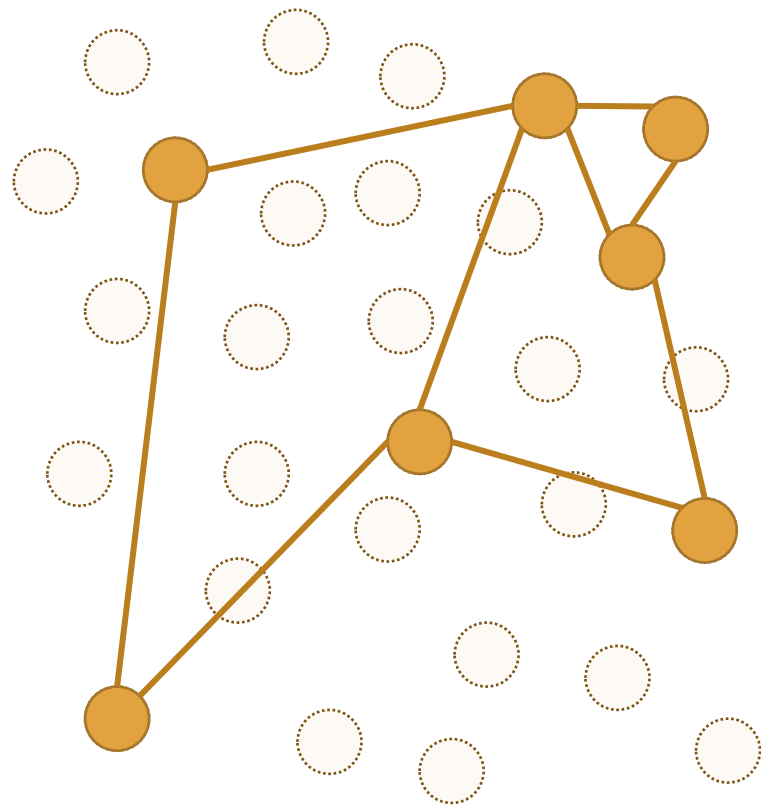
*Can this strategy help  
dense retrievers?*

● = document node  
— = neighbor edge

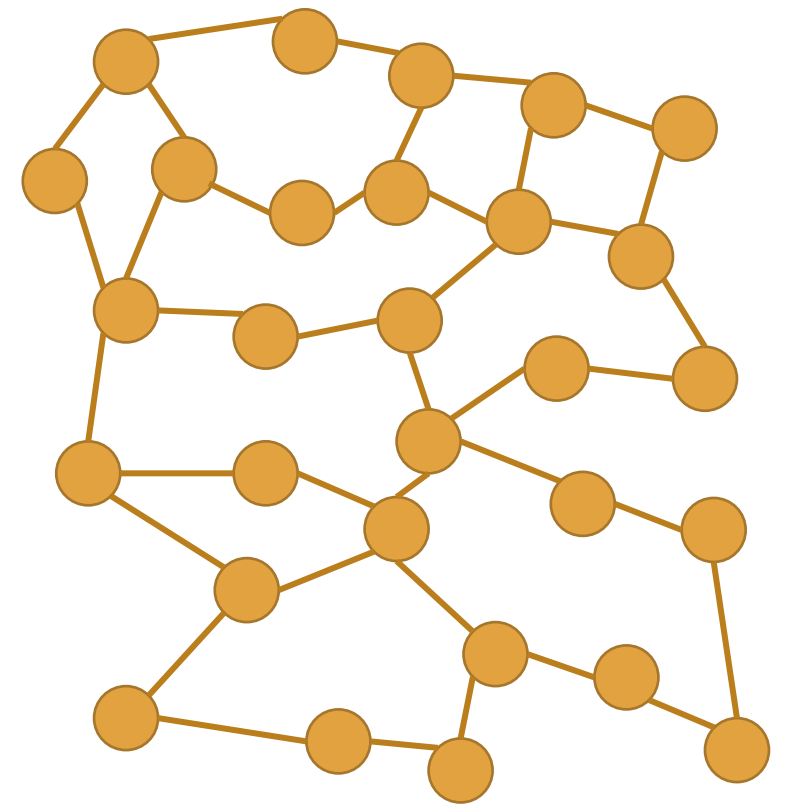
Level 2



Level 1



Level 0

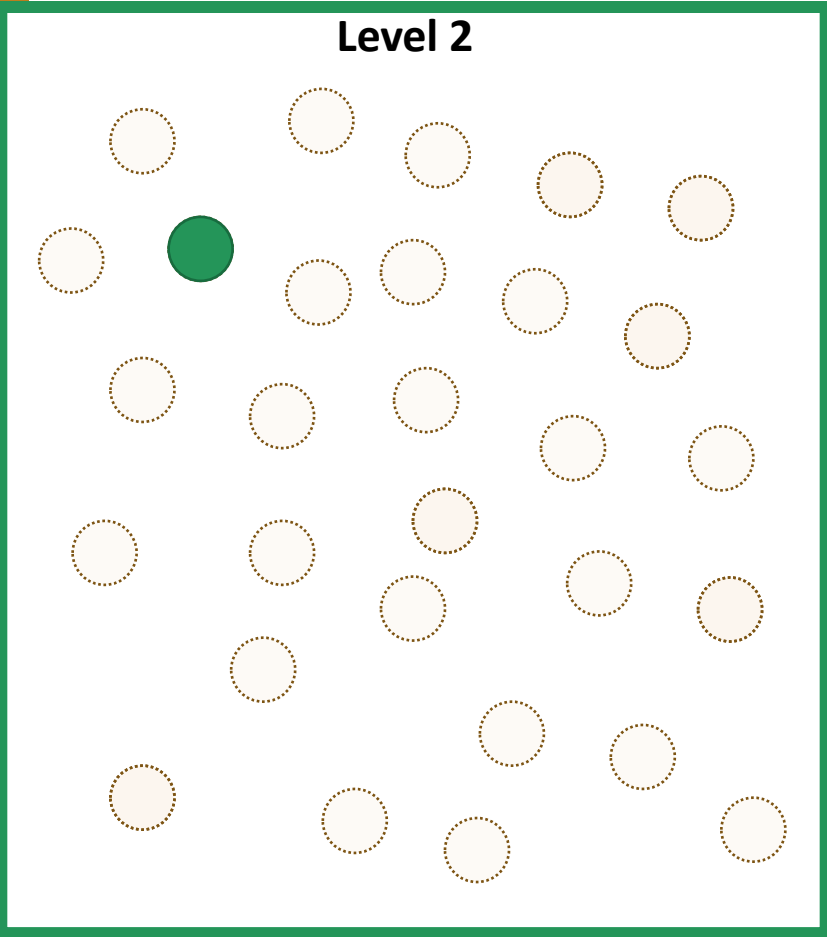


**HNSW: Score random nodes to narrow in on the best ones.**

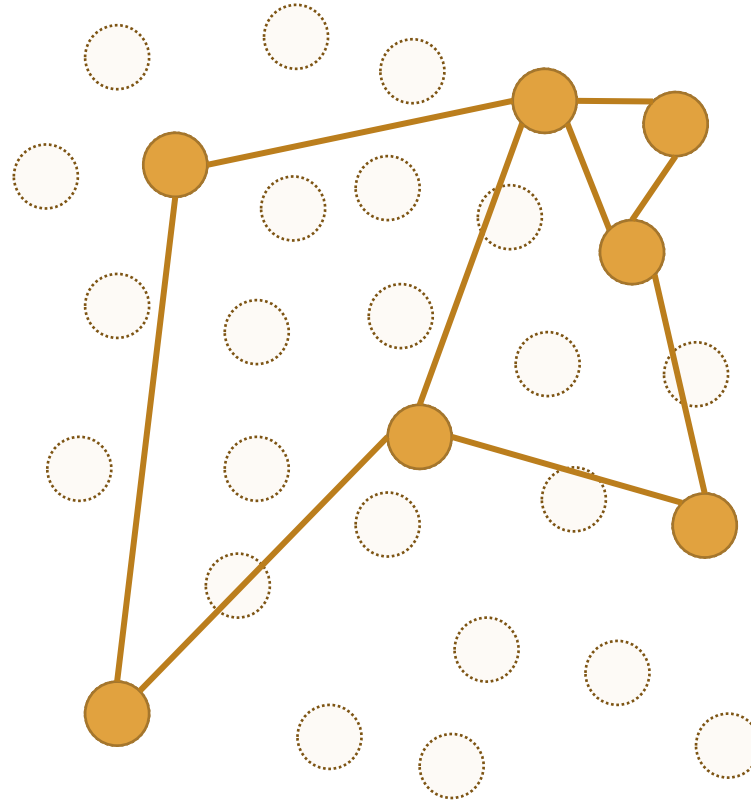
- = top document
- = scored document

- = document node
- = neighbor edge

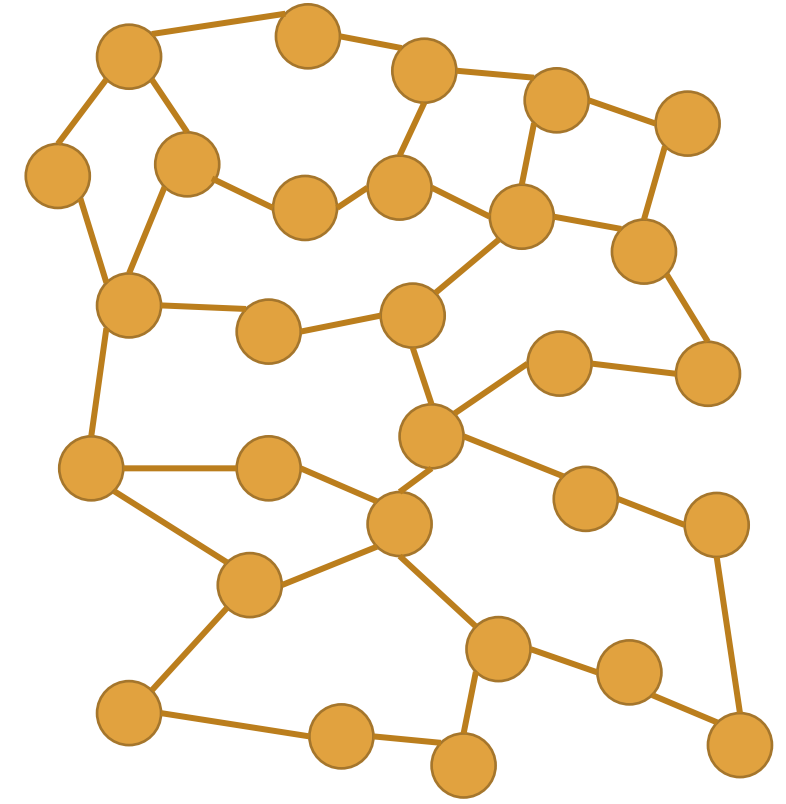
Level 2



Level 1



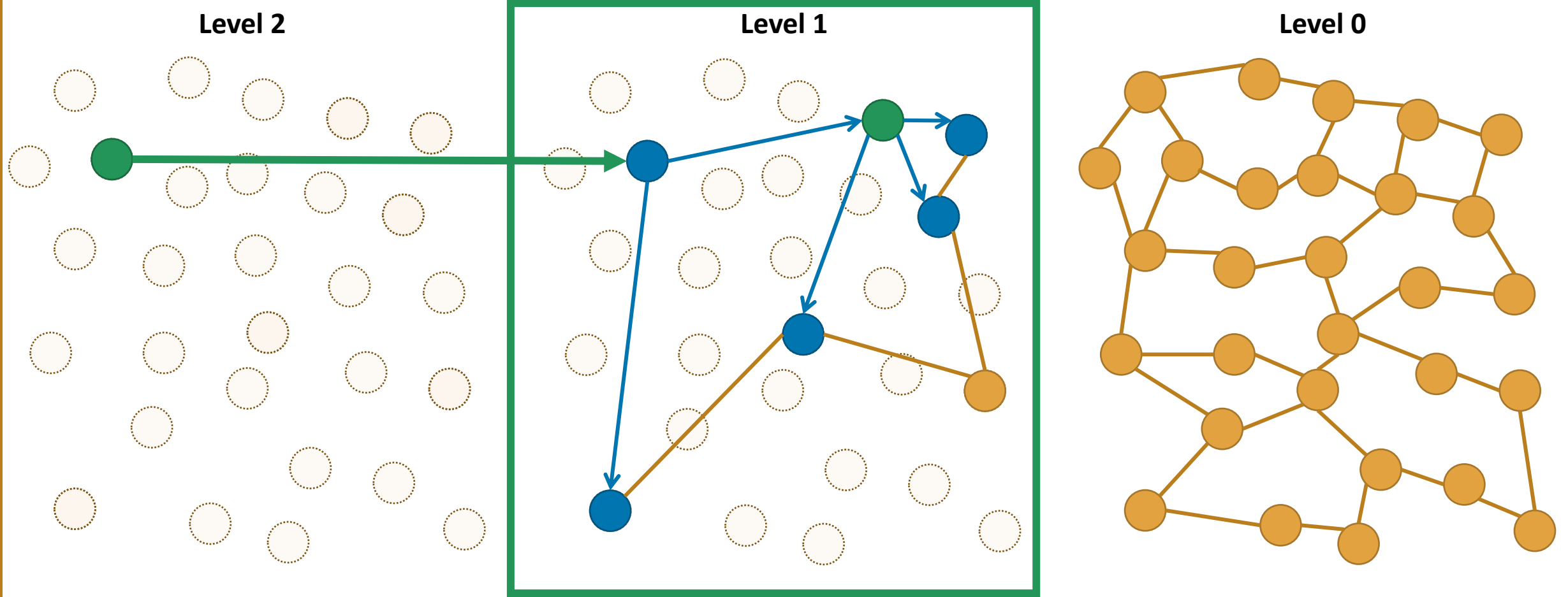
Level 0



**HNSW: Score random nodes to narrow in on the best ones.**

- = top document
- = scored document

- = document node
- = neighbor edge

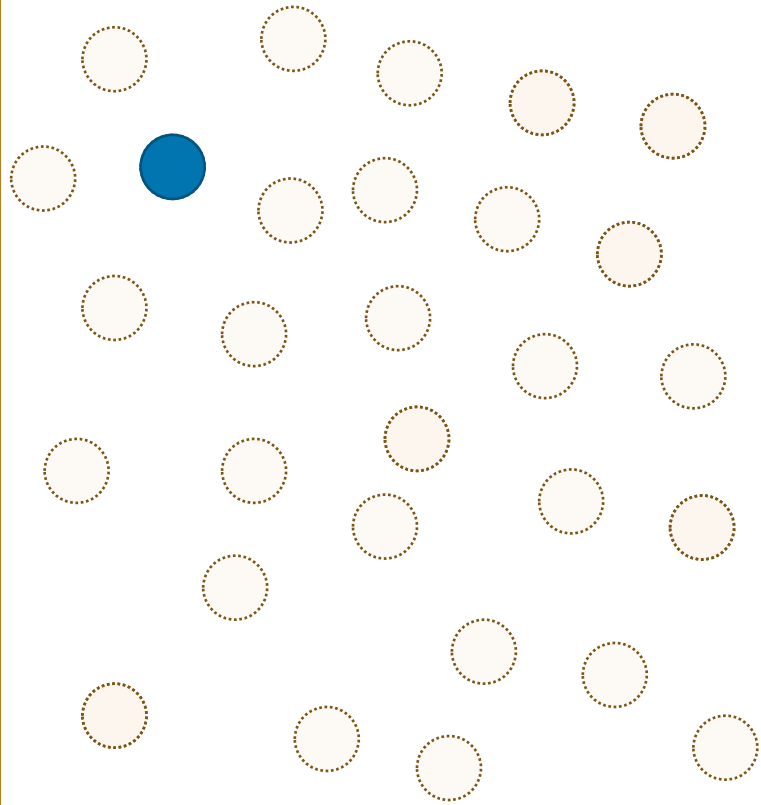


**HNSW: Score random nodes to narrow in on the best ones.**

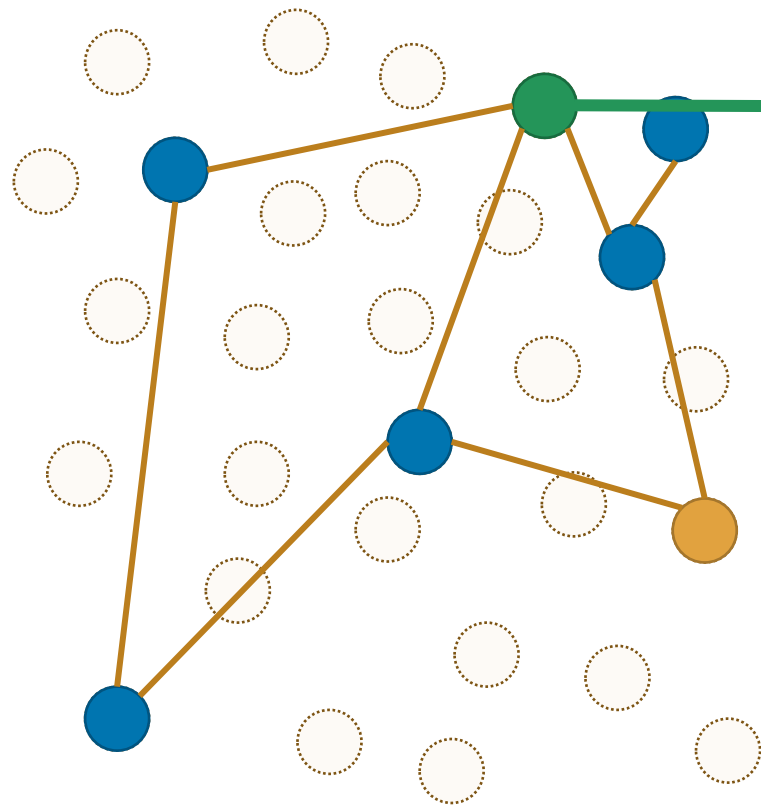
● = top document  
● = scored document

● = document node  
— = neighbor edge

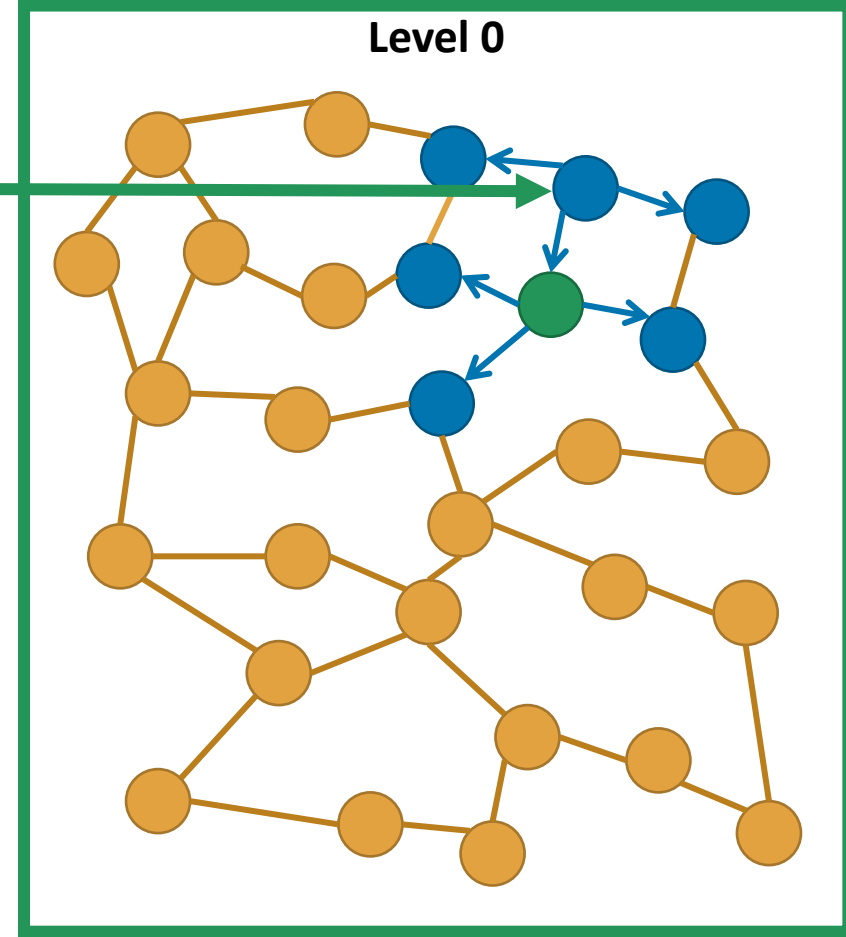
Level 2



Level 1

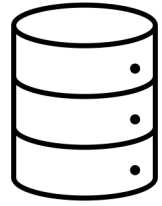


Level 0



**HNSW: Score random nodes to narrow in on the best ones.**

- = top document
- = scored document



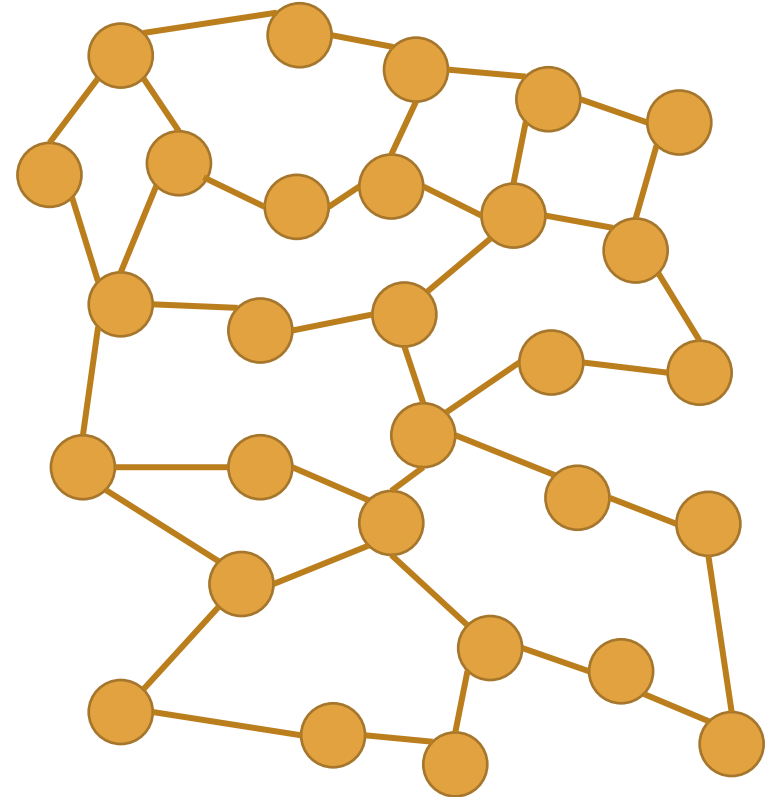
dataset



retriever  
(e.g., BM25)



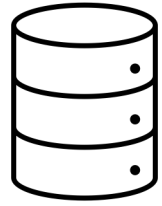
- 1.
- 2.
- 3.
- ...



- = document node
- = neighbor edge

**LADR (Lexically-Accelerated Dense Retrieval):**  
Use lexical search to seed dense retrieval.

- = top document
- = scored document



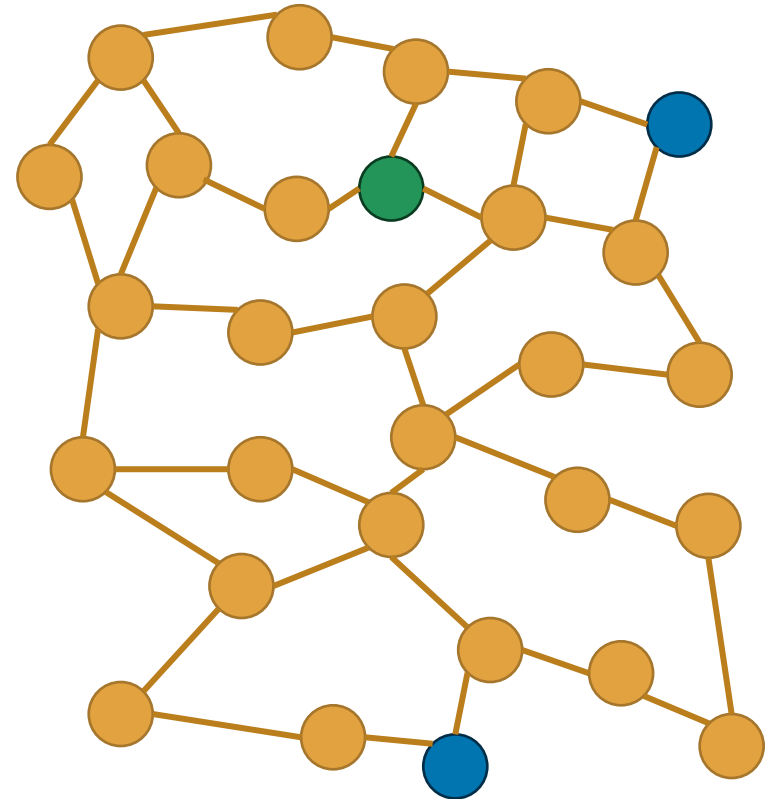
dataset



retriever  
(e.g., BM25)



- 1.
- 2.
- 3.
- ...

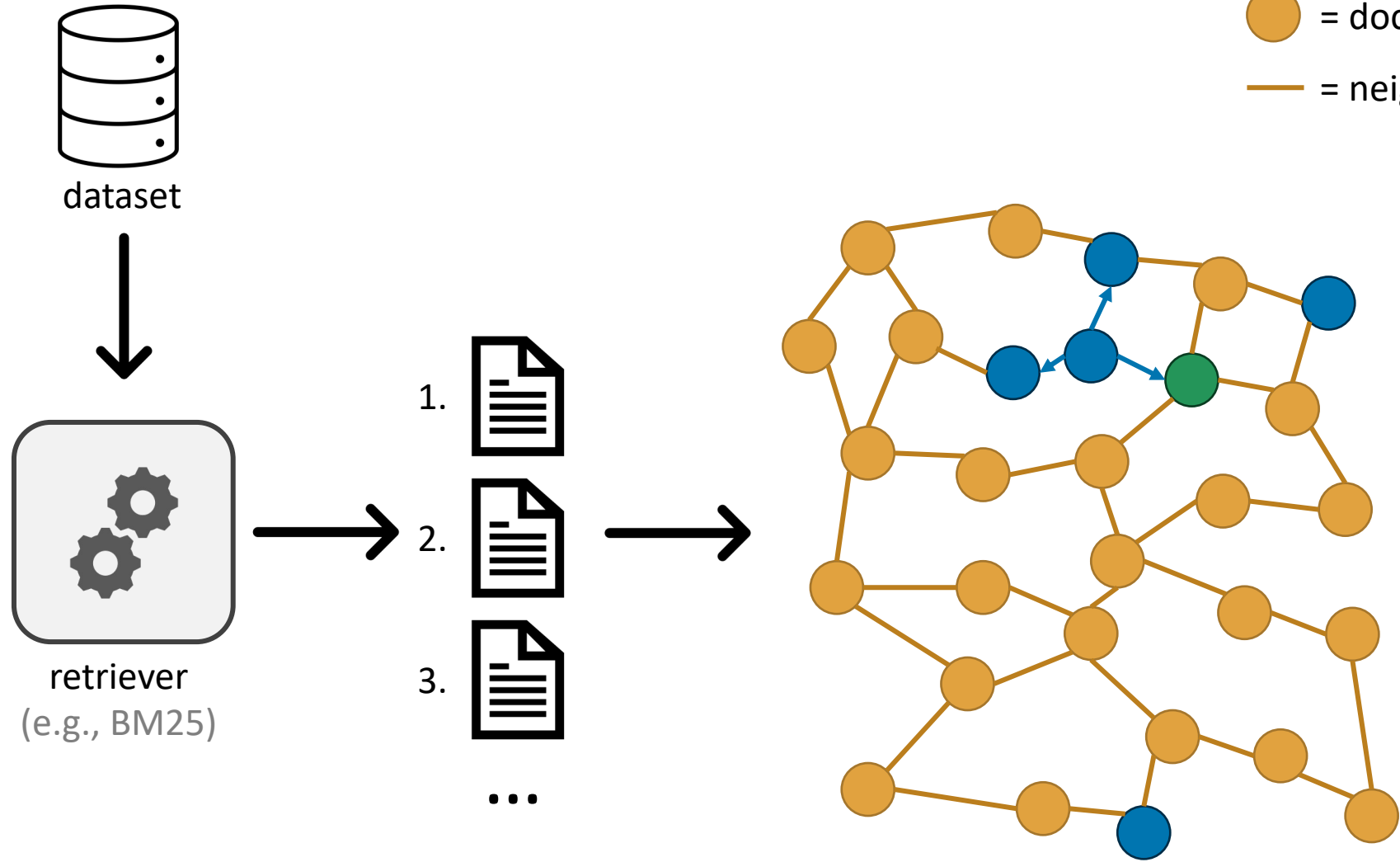


- = document node
- = neighbor edge



**LADR (Lexically-Accelerated Dense Retrieval):**  
Use lexical search to seed dense retrieval.



- = top document
- = scored document

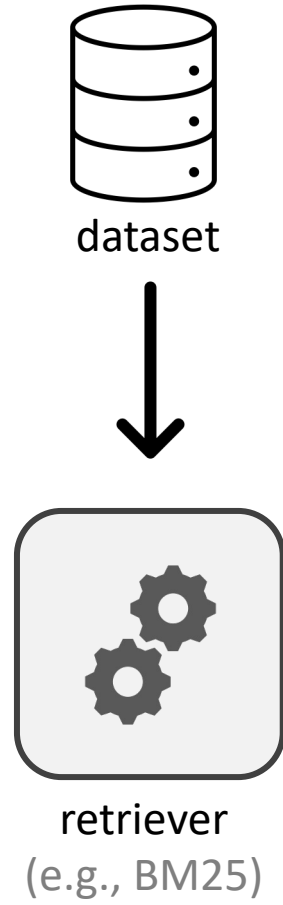
- = document node
- = neighbor edge






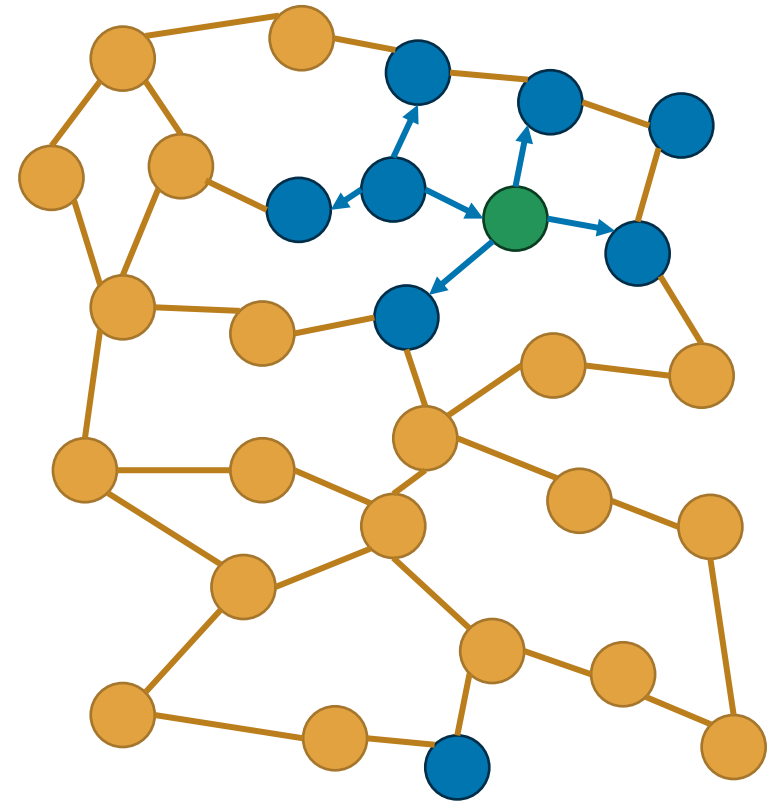
**LADR (Lexically-Accelerated Dense Retrieval):**  
Use lexical search to seed dense retrieval.

-  = top document
-  = scored document

-  = document node
-  = neighbor edge



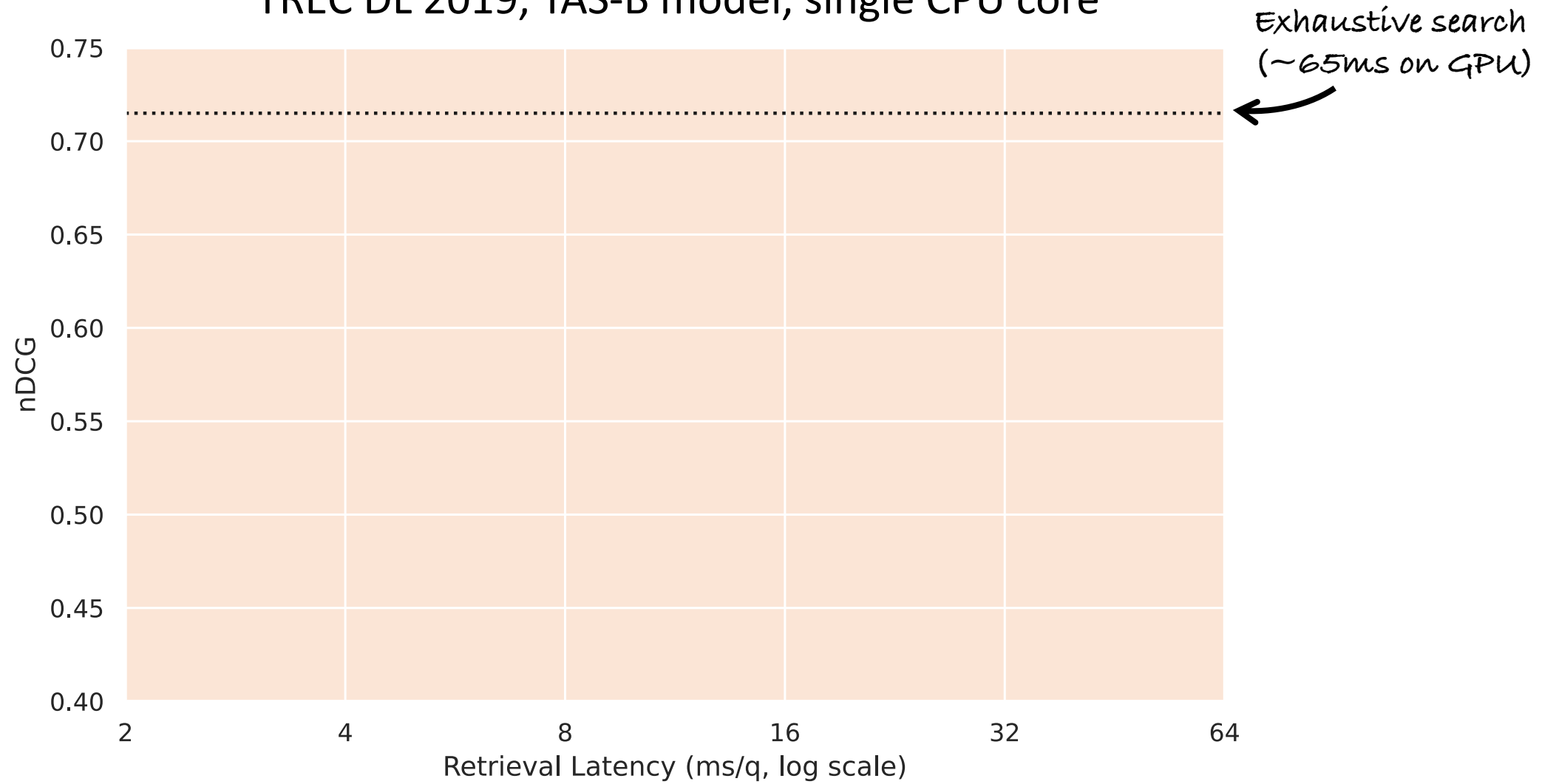
- 
- 
- 
- ...



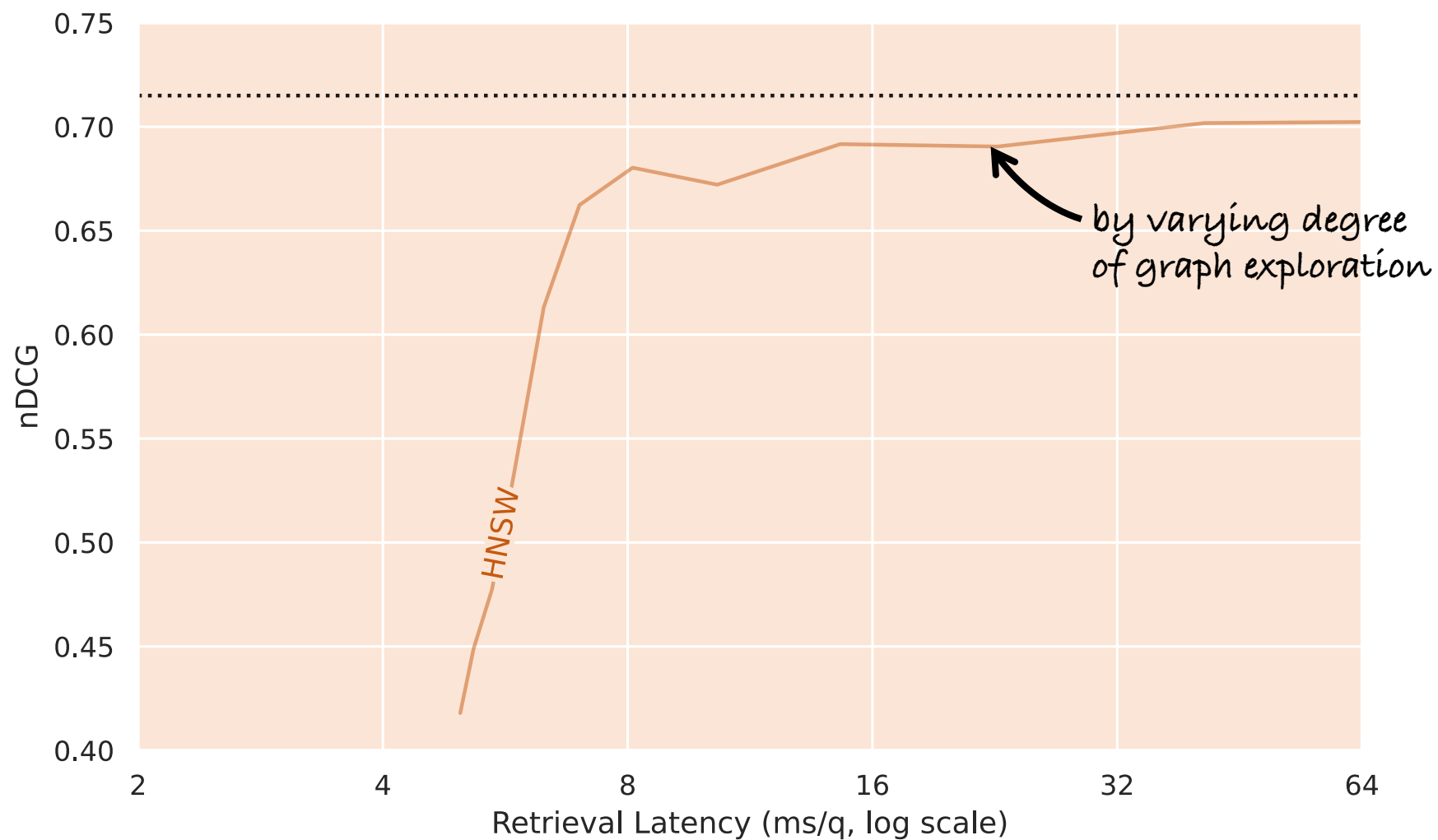
**Iterate...**

How well does it work?

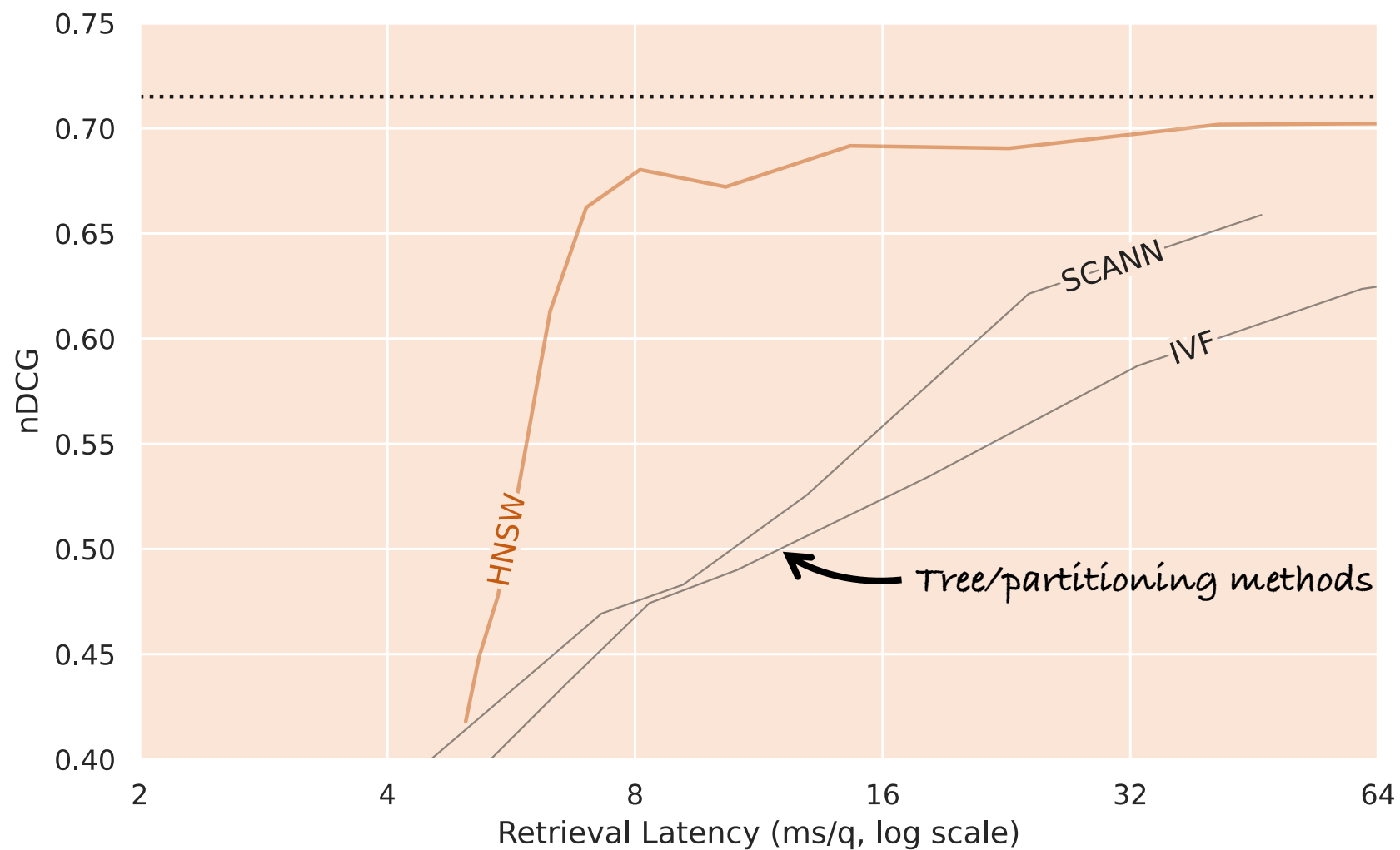
# TREC DL 2019, TAS-B model, single CPU core



## TREC DL 2019, TAS-B model, single CPU core

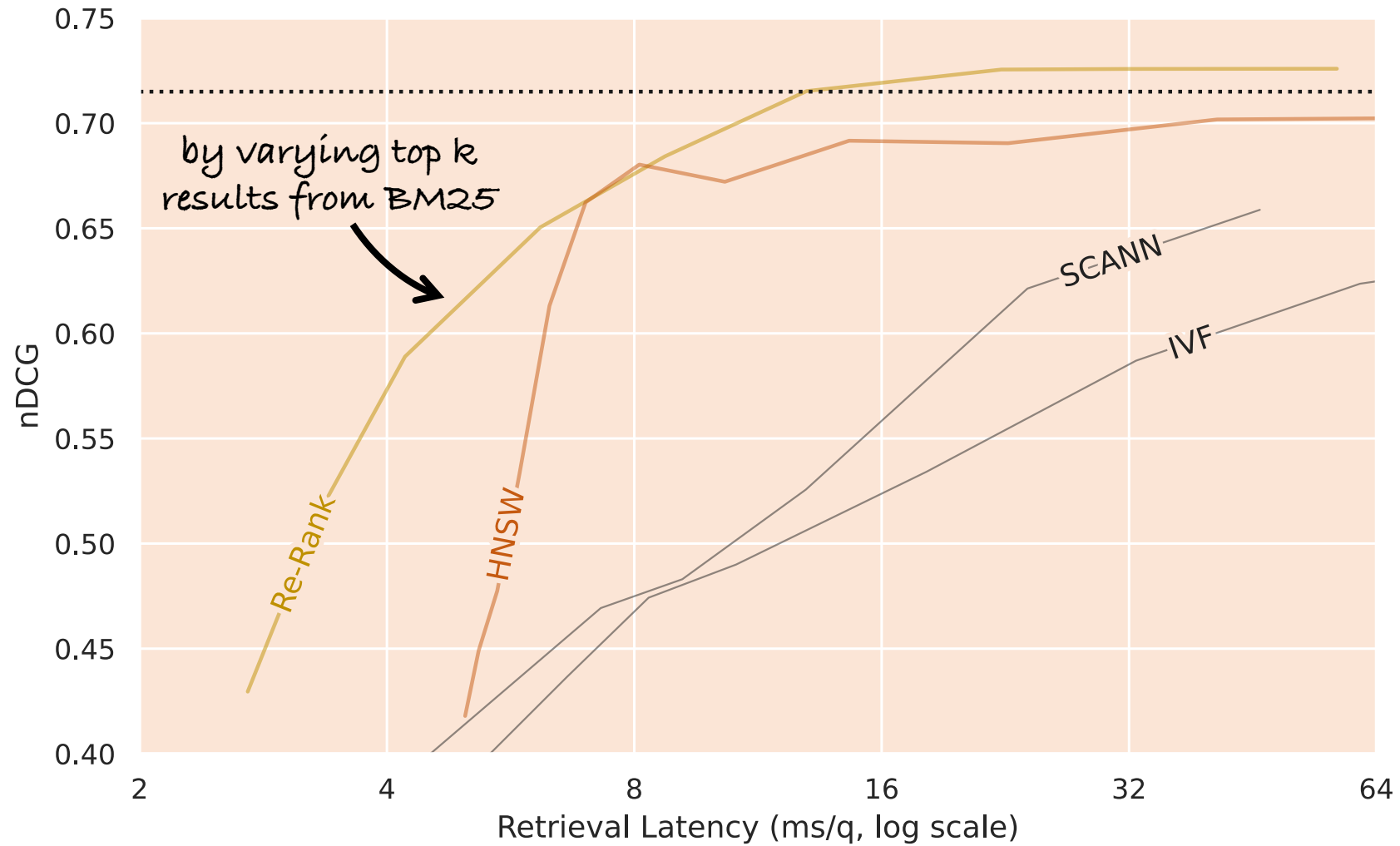


## TREC DL 2019, TAS-B model, single CPU core

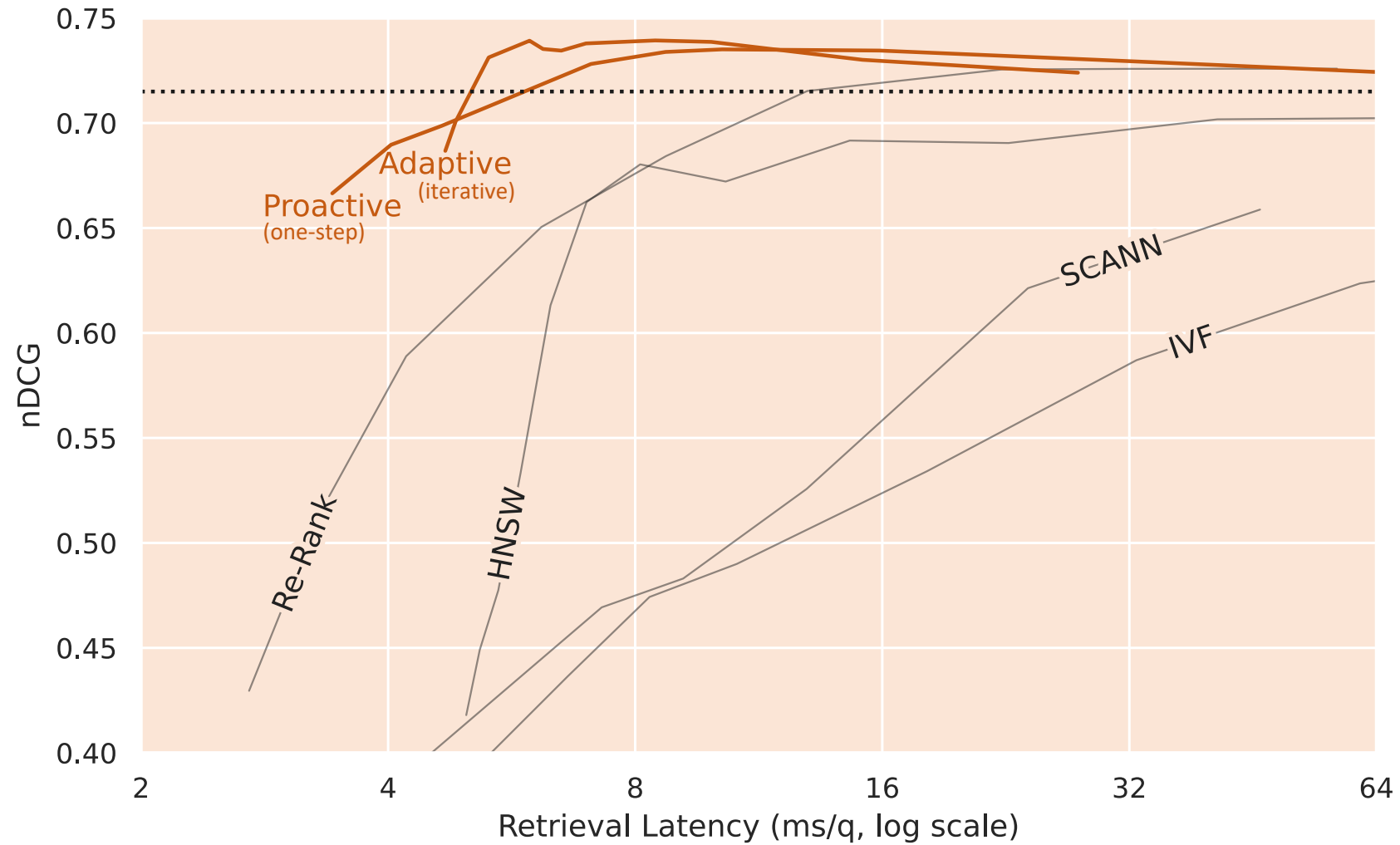


\* FAISS implementation

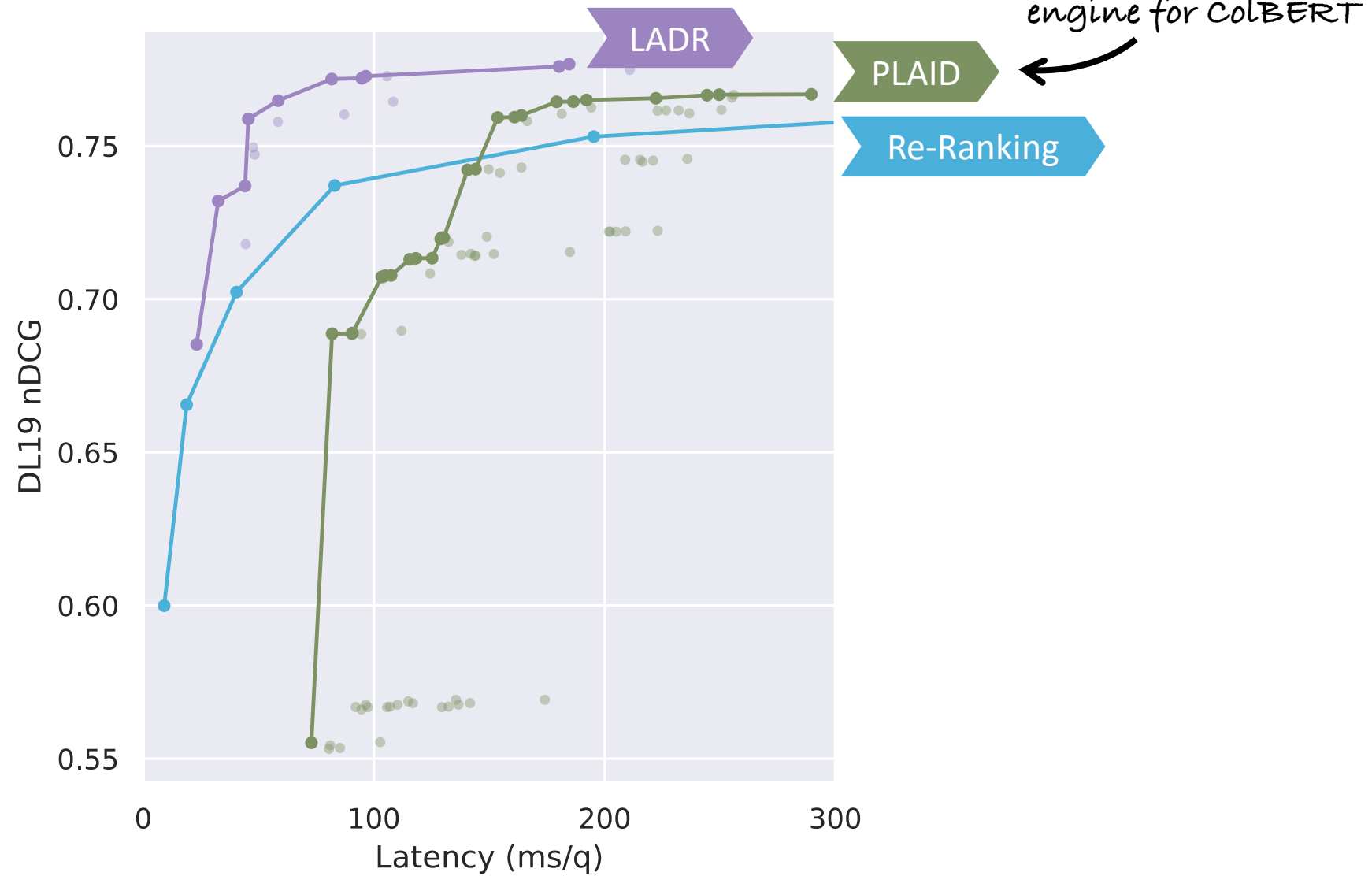
# TREC DL 2019, TAS-B model, single CPU core



# TREC DL 2019, TAS-B model, single CPU core



TREC DL 2019, CoLBERTv2, single CPU core



# In summary

Adaptive Re-Ranking improves **retrievers** too!

Both single-vector and multi-vector dense retrieval.

## Other findings

- Clear trade-offs between the model parameters and efficiency/effectiveness
- Works across a variety of dense retrieval models and standard benchmarks
- Works with approximate NN graph, and even graphs constructed from other models

## Conference Papers:

Kulkarni, MacAvaney, Goharian, Frieder. Lexically-Accelerated Dense Retrieval. SIGIR 2023.

MacAvaney, Tonello. A Reproducibility Study of PLAID. SIGIR 2024. [Best Paper Runner-Up]

Open Source!

# Adaptive Re-Ranking and LADR in PyTerrier

```
import pyterrier as pt
from pyterrier_t5 import MonoT5
from pyterrier_pisa import PisaIndex
from pyterrier_adaptive import GAR, CorpusGraph

bm25 = PisaIndex('my_index.pisa').bm25()
reranker = MonoT5()
graph = CorpusGraph.load('my_index.graph')

pipeline = bm25 >> GAR(reranker, graph)
```

```
from pyterrier_dr import FlexIndex
from pyterrier_pisa import PisaIndex

sparse = PisaIndex('my_index.pisa').bm25()
dense = PisaIndex('my_index.flex').ladr()

ladr = sparse.bm25() >> dense.ladr()
```

# Adaptive Re-Ranking using the `rerankers` package

```
from rerankers import Reranker

reranker = Reranker(...)

adaptive_reranker = GAR(reranker.as_pyterrier_transformer(), graph)
```

# We're also working on adding LADR to Lucene:

## Seeding HNSW Search #13634

 Open seanmacavane opened this issue on Aug 6 · 2 comments



seanmacavane commented on Aug 6

### Description

In some vector search cases, users may already know some documents that are likely related to a query. Let's support seeding HNSW's scoring stage with these documents, rather than using HNSW's hierarchical stage.

An example use case is hybrid search, where both a traditional and vector search are performed. The top results from the traditional search are likely reasonable seeds for the vector search. Even when not performing hybrid search, traditional matching can often be faster than traversing the hierarchy, which can be used to speed up the vector search process (up to 2x faster for the same effectiveness), as was demonstrated in [this article](#) (full disclosure: I'm an author of the article).

This enhancement proposes adding a `seed` query, alongside the existing `filter` query, to the KNN query classes. The results of this query will be fed into `HnswGraphSearcher`, and ultimately replace the graph entry points [here](#). If the seed query fails (e.g., keywords do not match any documents), the approach will fall back onto the existing hierarchical search process.

Pull request to follow.



1

```
AbstractKnnVectorQuery.seed
```

# Thanks to my collaborators!



**Craig Macdonald**  
University of Glasgow



**Nicola Tonellotto**  
University of Pisa



**Hrishikesh Kulkarni**  
Georgetown University



**Nazli Goharian**  
Georgetown University



**Ophir Frieder**  
Georgetown University

# RE-THINKING RE-RANKING

Sean MacAvaney  
University of Glasgow









Presented at:  
Haystack EU 2024

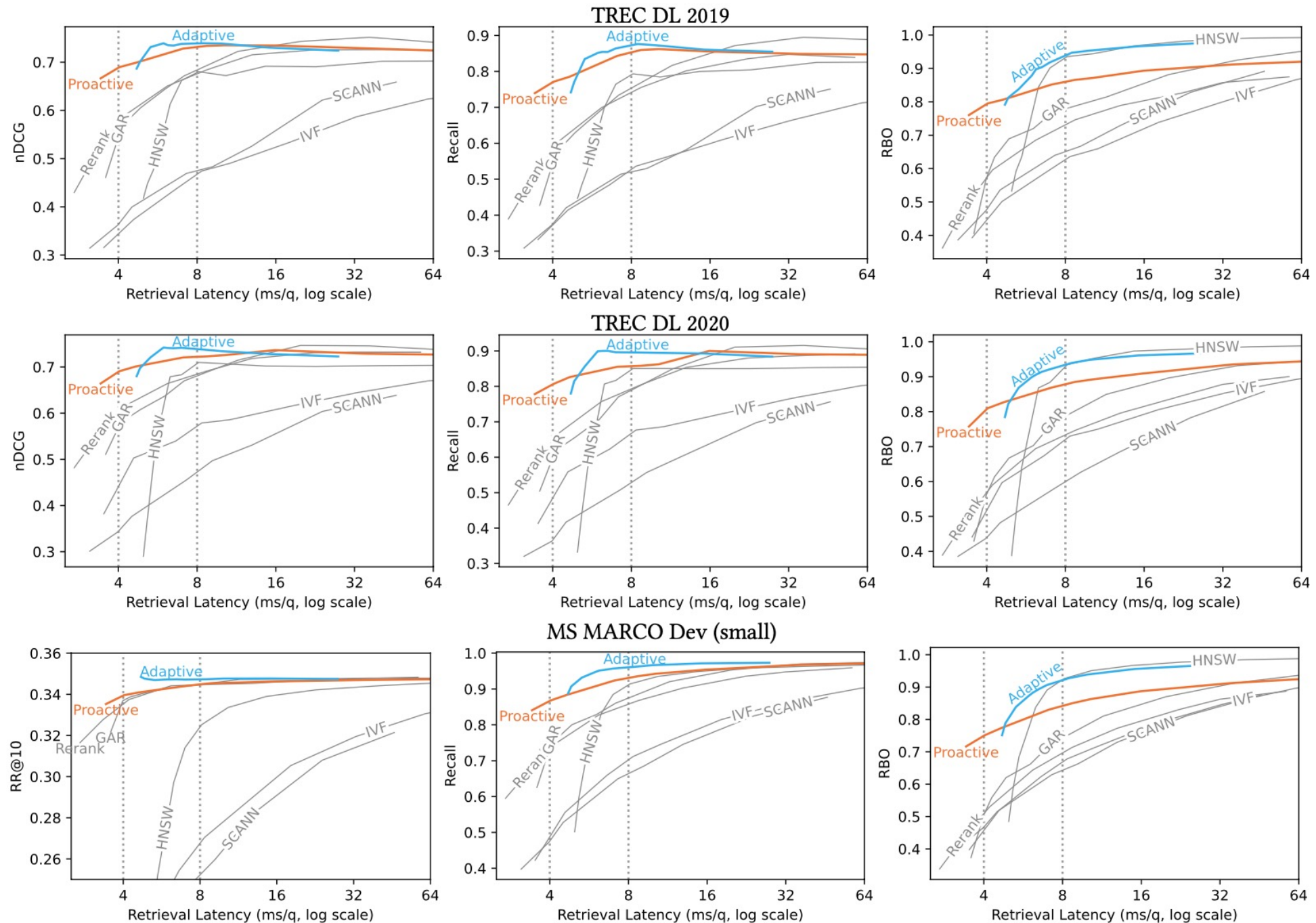
Terrier 



University  
of Glasgow

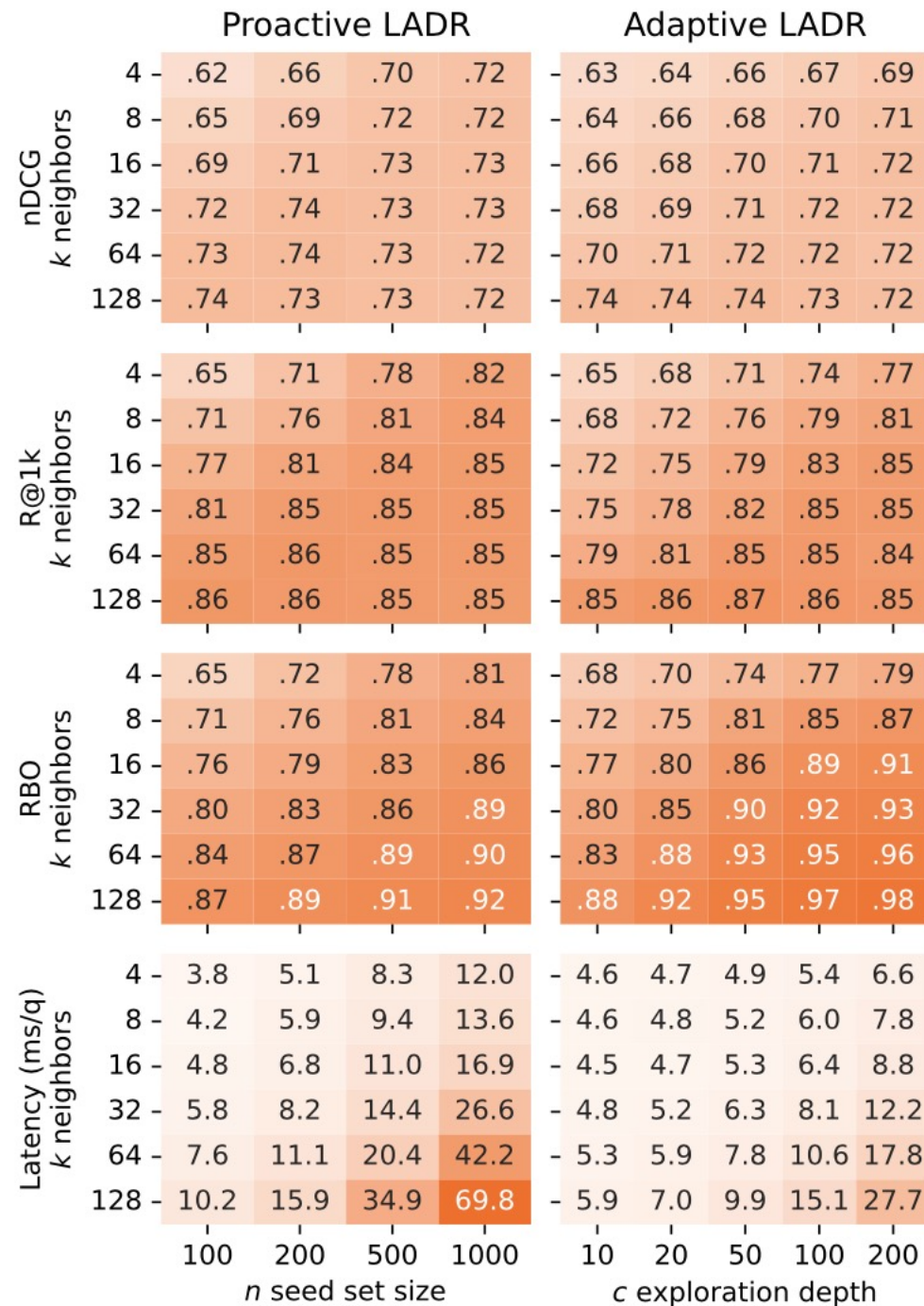
extra slides

Search Strategy	Effectiveness	Query Efficiency	Index Efficiency	Storage Costs
Sparse	 Low	 High	 High	 Low
Dense	 High	 Low	 Low	 High



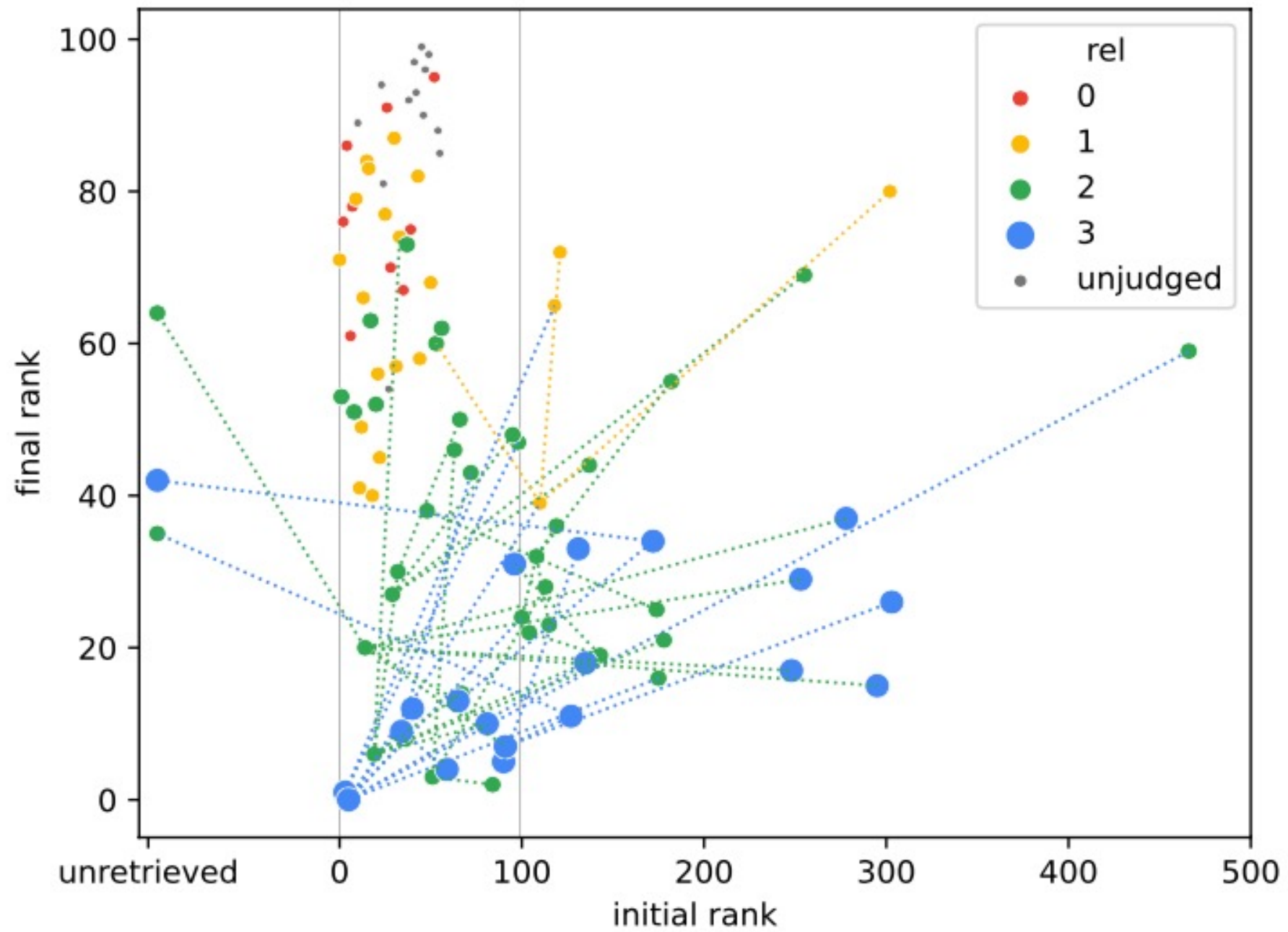
**Figure 3: Performance of LADR over TAS-B and baselines across various operational points.**

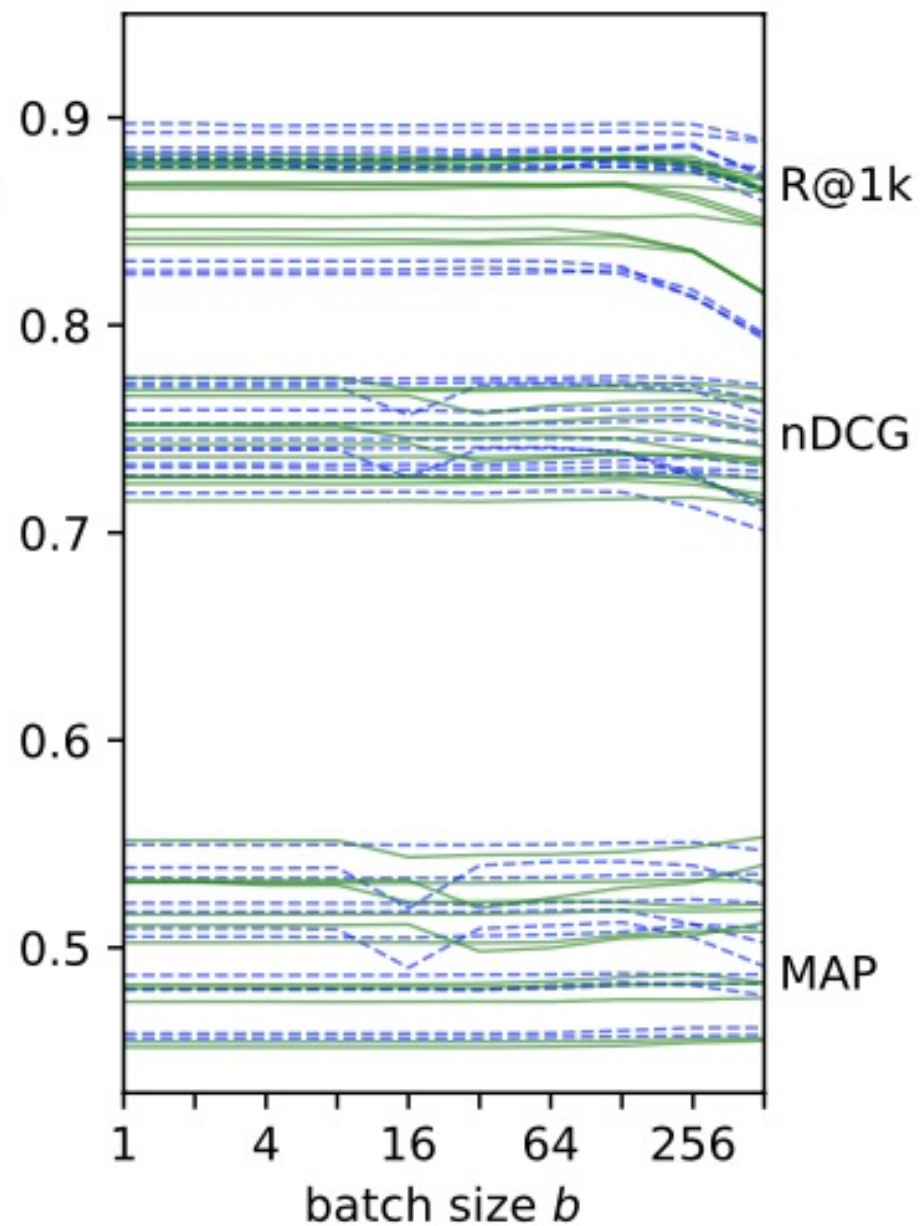
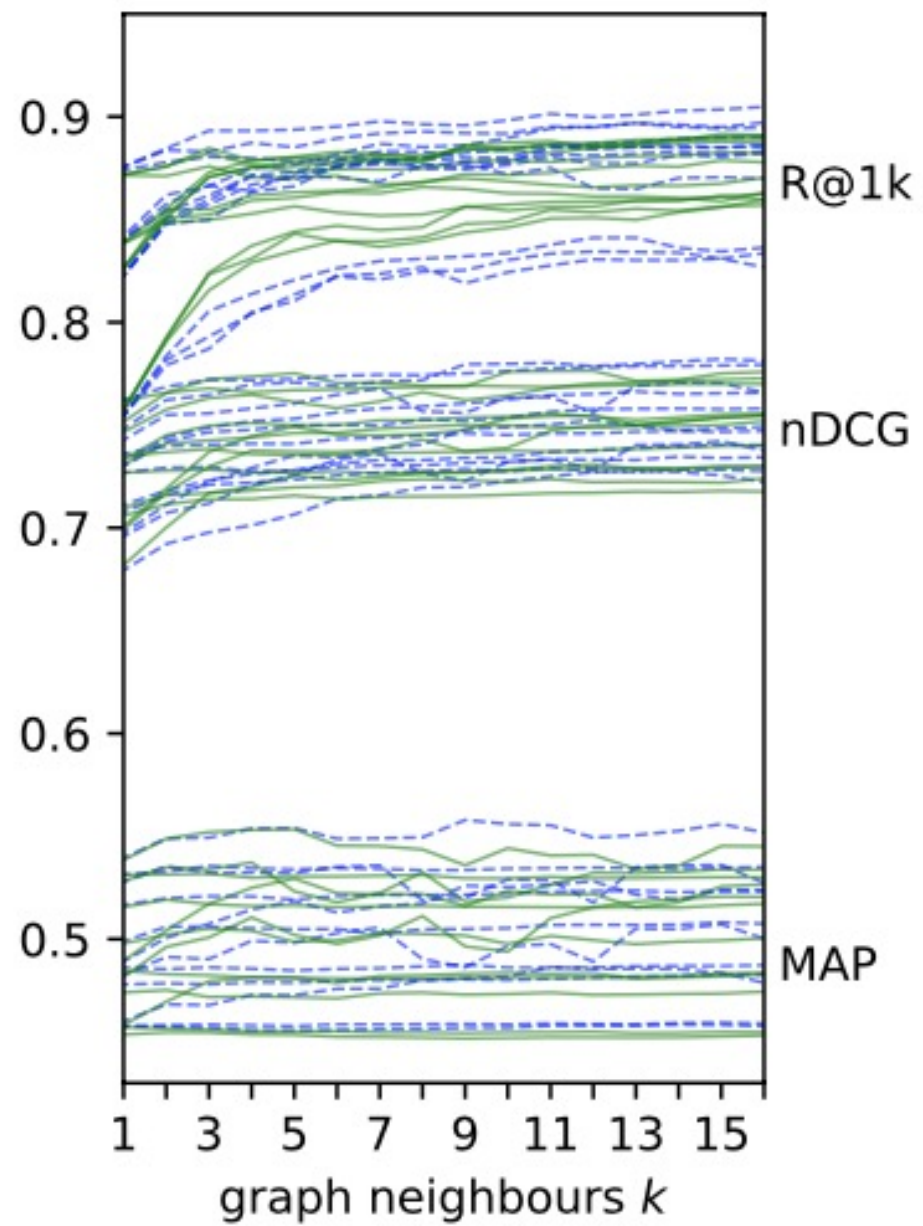
Method	DL19 ~4ms		DL19 ~8ms		DL20 ~4ms		DL20 ~8ms		Dev (sm) ~4ms		Dev (sm) ~8ms	
	nDCG	R@1k	nDCG	R@1k	nDCG	R@1k	nDCG	R@1k	RR@10	R@1k	RR@10	R@1k
<b>TAS-B (Exh.)</b>	0.715	0.842	0.715	0.842	0.713	0.875	0.713	0.875	0.347	0.978	0.347	0.978
IVF [I]	0.374	0.414	0.474	0.536	0.503	0.559	0.579	0.677	0.217	0.556	0.270	0.712
ScaNN [S]	0.475	0.519	0.537	0.598	0.476	0.527	0.553	0.641	0.254	0.669	0.292	0.774
HNSW [H]	-	-	0.614	0.707	-	-	0.699	0.836	-	-	0.310	0.872
GAR [G]	0.543	0.540	0.688	0.755	0.568	0.594	0.684	0.796	0.337	0.732	0.345	0.876
Re-Ranking [R]	0.589	0.605	0.684	0.755	0.615	0.667	0.691	0.805	0.337	0.748	0.345	0.868
Proactive LADR	$IS_{GR}$ <b>0.690</b>	$IS_{GR}$ <b>0.771</b>	$ISH_{GR}$ 0.730	$ISH_{GR}$ 0.850	$IS_{GR}$ <b>0.691</b>	$IS_{GR}$ <b>0.807</b>	$IS_{GR}$ 0.722	$IS_{GR}$ 0.857	$IS$ <b>0.340</b>	$IS_{GR}$ <b>0.868</b>	$ISH$ 0.345	$ISH_{GR}$ 0.932
Adaptive LADR	-	-	$ISH_{GR}$ <b>0.738</b>	$ISH_{GR}$ <b>0.872</b>	-	-	$ISH_{GR}$ <b>0.739</b>	$ISH_{GR}$ <b>0.900</b>	-	-	$ISH_{GR}$ <b>0.347</b>	$ISH_{GR}$ <b>0.960</b>
<b>RetroMAE (Exh.)</b>	0.699	0.806	0.699	0.806	0.701	0.839	0.701	0.839	0.375	0.981	0.375	0.981
IVF [I]	0.226	0.225	0.346	0.358	0.272	0.263	0.372	0.375	0.157	0.381	0.221	0.541
ScaNN [S]	0.468	0.502	0.525	0.588	0.486	0.509	0.555	0.606	0.275	0.665	0.312	0.769
HNSW [H]	-	-	0.630	0.720	-	-	0.673	0.798	-	-	0.338	0.874
GAR [G]	0.559	0.553	0.696	0.763	0.578	0.604	0.692	0.789	<b>0.357</b>	0.750	0.368	0.890
Re-Ranking [R]	0.594	0.605	0.685	0.755	0.622	0.667	0.696	0.805	0.355	0.748	0.369	0.868
Proactive LADR	$IS_{GR}$ <b>0.691</b>	$IS_{GR}$ <b>0.765</b>	$ISH_{GR}$ 0.733	$ISH_{GR}$ 0.844	$IS_{GR}$ <b>0.702</b>	$IS_{GR}$ <b>0.811</b>	$ISH_G$ 0.723	$IS_G$ 0.846	$IS$ 0.356	$IS_{GR}$ <b>0.864</b>	$ISH$ 0.368	$ISH_{GR}$ 0.938
Adaptive LADR	-	-	$ISH_R$ <b>0.740</b>	$ISH_{GR}$ <b>0.866</b>	-	-	$ISH_G$ <b>0.731</b>	$ISH_{GR}$ <b>0.879</b>	-	-	$ISH_{GR}$ <b>0.374</b>	$ISH_{GR}$ <b>0.973</b>
<b>TCT-HNP (Exh.)</b>	0.708	0.830	0.708	0.830	0.689	0.848	0.689	0.848	0.359	0.970	0.359	0.970
IVF [I]	0.340	0.366	0.437	0.469	0.369	0.383	0.470	0.522	0.219	0.527	0.276	0.687
ScaNN [S]	0.378	0.410	0.444	0.496	0.355	0.376	0.427	0.459	0.215	0.522	0.253	0.632
HNSW [H]	-	-	0.625	0.721	-	-	0.634	0.762	-	-	0.315	0.853
GAR [G]	0.546	0.547	0.687	0.755	0.569	0.598	0.678	0.797	0.342	0.733	0.354	0.878
Re-Ranking [R]	0.586	0.605	0.679	0.755	0.614	0.667	0.685	0.805	0.342	0.748	0.353	0.868
Proactive LADR	$IS_{GR}$ <b>0.680</b>	$IS_{GR}$ <b>0.747</b>	$ISH_{GR}$ 0.719	$ISH_{GR}$ 0.827	$IS_{GR}$ <b>0.682</b>	$IS_{GR}$ <b>0.803</b>	$ISH_G$ 0.709	$ISH$ 0.841	$IS_G$ <b>0.346</b>	$IS_{GR}$ <b>0.856</b>	$ISH$ 0.354	$ISH_{GR}$ 0.927
Adaptive LADR	-	-	$ISH$ <b>0.729</b>	$ISH_{GR}$ <b>0.848</b>	-	-	$ISH_{GR}$ <b>0.721</b>	$ISH_{GR}$ <b>0.878</b>	-	-	$ISH_{GR}$ <b>0.359</b>	$ISH_{GR}$ <b>0.962</b>
<b>ANCE (Exh.)</b>	0.617	0.755	0.617	0.755	0.634	0.777	0.634	0.777	0.330	0.957	0.330	0.957
IVF [I]	0.358	0.395	0.441	0.500	0.407	0.437	0.498	0.549	0.212	0.530	0.268	0.703
ScaNN [S]	0.374	0.405	0.433	0.488	0.440	0.495	0.535	0.614	0.262	0.691	0.287	0.783
HNSW [H]	-	-	0.606	0.737	-	-	0.635	0.790	-	-	0.311	0.897
GAR [G]	0.527	0.540	0.648	0.750	0.568	0.622	0.655	0.794	<b>0.326</b>	0.751	0.329	0.888
Re-Ranking [R]	0.578	0.605	0.653	0.755	0.602	0.667	<b>0.674</b>	0.805	0.325	0.748	<b>0.333</b>	0.868
Proactive LADR	$IS_{GR}$ <b>0.645</b>	$IS_{GR}$ <b>0.751</b>	$IS$ 0.657	$ISH$ 0.800	$IS_{GR}$ <b>0.660</b>	$IS_{GR}$ <b>0.807</b>	$IS$ 0.666	$IS$ 0.822	$IS$ 0.321	$IS_{GR}$ <b>0.872</b>	$ISH$ 0.327	$ISH_{GR}$ 0.932
Adaptive LADR	-	-	$ISH$ <b>0.665</b>	$ISH$ <b>0.820</b>	-	-	$ISH$ 0.665	$IS$ <b>0.830</b>	-	-	$ISH$ 0.329	$ISH_{GR}$ <b>0.959</b>



Graph	DL19		DL20		Dev (sm)	
	nDCG	R@1k	nDCG	R@1k	RR@10	R@1k
<b>Proactive LADR</b>						
Exact	0.730	<b>0.850</b>	<b>0.722</b>	<b>0.857</b>	<b>0.345</b>	<b>0.932</b>
Approx.	<sup>=</sup> 0.731	0.845	<sup>=</sup> 0.720	0.849	*0.343	*0.916
BM25	<sup>=</sup> <b>0.732</b>	0.835	<sup>=</sup> 0.720	0.853	*0.339	*0.883
<b>Adaptive LADR</b>						
Exact	0.738	<b>0.872</b>	0.739	<b>0.900</b>	<b>0.347</b>	0.960
Approx.	<sup>=</sup> 0.736	0.861	<sup>=</sup> 0.737	<sup>=</sup> <b>0.900</b>	<sup>=</sup> <b>0.347</b>	* <b>0.966</b>
BM25	<b>0.743</b>	0.859	<sup>=</sup> <b>0.742</b>	<b>0.900</b>	*0.345	*0.933

Pipeline	DL19 (valid.) $c = 100$			DL19 (valid.) $c = 1000$			DL20 (test) $c = 100$			DL20 (test) $c = 1000$		
	nDCG	MAP	R@1k	nDCG	MAP	R@1k	nDCG	MAP	R@1k	nDCG	MAP	R@1k
BM25»MonoT5-base	0.665	0.417	0.755	0.699	0.483	0.755	0.672	0.421	0.805	0.711	0.498	0.805
w/ GAR <sub>BM25</sub>	*0.697	*0.456	*0.786	0.727	0.490	*0.827	*0.695	0.439	*0.823	*0.743	<b>0.501</b>	*0.874
w/ GAR <sub>TCT</sub>	<b>*0.722</b>	<b>*0.491</b>	<b>*0.800</b>	<b>*0.743</b>	<b>0.511</b>	<b>*0.839</b>	<b>*0.714</b>	<b>*0.472</b>	<b>*0.831</b>	<b>*0.749</b>	<b>0.501</b>	<b>*0.892</b>
BM25»MonoT5-3b	0.667	0.418	0.755	0.700	0.489	0.755	0.678	0.442	0.805	0.728	0.534	0.805
w/ GAR <sub>BM25</sub>	*0.693	0.454	*0.790	*0.741	0.517	*0.831	*0.715	*0.469	*0.829	*0.772	0.556	*0.881
w/ GAR <sub>TCT</sub>	<b>*0.715</b>	<b>*0.484</b>	<b>*0.806</b>	<b>*0.746</b>	<b>0.522</b>	<b>*0.846</b>	<b>*0.735</b>	<b>*0.512</b>	<b>*0.837</b>	<b>*0.787</b>	<b>*0.564</b>	<b>*0.899</b>
BM25»ColBERT	0.663	0.409	0.755	0.681	0.458	0.755	0.667	0.421	0.805	0.697	0.469	0.805
w/ GAR <sub>BM25</sub>	*0.690	*0.442	*0.783	*0.720	0.480	*0.825	*0.695	*0.446	*0.823	*0.732	0.479	*0.870
w/ GAR <sub>TCT</sub>	<b>*0.716</b>	<b>*0.475</b>	<b>*0.798</b>	<b>*0.727</b>	<b>0.482</b>	<b>*0.841</b>	<b>*0.707</b>	<b>*0.463</b>	<b>*0.829</b>	<b>*0.740</b>	<b>0.481</b>	<b>*0.887</b>
TCT»MonoT5-base	0.708	0.472	0.830	0.704	0.473	0.830	0.698	0.488	0.848	0.693	0.471	0.848
w/ GAR <sub>BM25</sub>	<b>*0.728</b>	<b>0.484</b>	<b>0.852</b>	<b>*0.733</b>	<b>0.480</b>	<b>*0.883</b>	<b>*0.719</b>	<b>*0.501</b>	<b>0.861</b>	<b>*0.719</b>	<b>0.473</b>	<b>*0.881</b>
w/ GAR <sub>TCT</sub>	0.722	0.481	0.847	*0.724	0.474	0.866	*0.712	0.494	0.856	*0.710	0.471	0.871
TCT»MonoT5-3b	0.720	0.498	0.830	0.725	0.513	0.830	0.723	0.534	0.848	0.733	0.544	0.848
w/ GAR <sub>BM25</sub>	<b>*0.748</b>	<b>*0.521</b>	<b>*0.857</b>	<b>*0.759</b>	<b>0.521</b>	<b>*0.885</b>	<b>*0.743</b>	<b>0.546</b>	<b>*0.864</b>	<b>*0.771</b>	<b>*0.555</b>	<b>*0.890</b>
w/ GAR <sub>TCT</sub>	*0.742	*0.517	0.849	*0.749	0.516	*0.868	*0.741	*0.545	*0.861	*0.759	0.551	*0.880
TCT»ColBERT	0.708	0.464	0.830	0.701	0.452	0.830	0.698	0.476	0.848	0.697	0.470	0.848
w/ GAR <sub>BM25</sub>	<b>*0.729</b>	<b>*0.480</b>	<b>0.853</b>	<b>*0.727</b>	<b>0.459</b>	<b>0.876</b>	<b>*0.715</b>	<b>0.485</b>	<b>0.857</b>	<b>*0.722</b>	<b>*0.477</b>	<b>*0.877</b>
w/ GAR <sub>TCT</sub>	*0.722	0.474	0.845	*0.715	0.452	0.852	*0.711	*0.484	<b>*0.857</b>	*0.713	0.473	0.864
D2Q»MonoT5-base	0.736	0.503	0.830	0.747	0.531	0.830	0.726	0.499	0.839	0.731	<b>0.508</b>	0.839
w/ GAR <sub>BM25</sub>	*0.748	0.506	0.848	0.757	0.519	<b>*0.880</b>	*0.734	0.497	*0.847	<b>0.748</b>	0.504	*0.880
w/ GAR <sub>TCT</sub>	<b>*0.760</b>	<b>*0.528</b>	<b>0.850</b>	<b>*0.766</b>	<b>0.533</b>	*0.879	<b>0.740</b>	<b>0.508</b>	<b>*0.856</b>	<b>0.748</b>	0.499	<b>*0.895</b>
D2Q»MonoT5-3b	0.737	0.506	0.830	0.751	0.542	0.830	0.738	0.531	0.839	0.753	0.557	0.839
w/ GAR <sub>BM25</sub>	0.744	0.512	*0.850	<b>0.772</b>	<b>0.549</b>	<b>*0.880</b>	*0.751	0.535	*0.852	*0.781	0.561	*0.887
w/ GAR <sub>TCT</sub>	<b>0.755</b>	<b>0.524</b>	<b>*0.857</b>	0.769	0.544	<b>*0.880</b>	<b>*0.764</b>	<b>0.550</b>	<b>*0.860</b>	<b>*0.790</b>	<b>0.565</b>	<b>*0.905</b>
D2Q»ColBERT	0.724	0.475	0.830	0.733	0.501	0.830	0.718	0.483	0.839	0.717	0.479	0.839
w/ GAR <sub>BM25</sub>	0.734	0.484	0.845	<b>0.753</b>	<b>0.505</b>	*0.876	*0.731	0.487	*0.849	*0.737	0.482	*0.872
w/ GAR <sub>TCT</sub>	<b>*0.744</b>	<b>*0.496</b>	<b>0.849</b>	*0.752	0.503	<b>*0.878</b>	<b>*0.735</b>	<b>0.488</b>	<b>*0.856</b>	<b>*0.746</b>	<b>0.485</b>	<b>*0.893</b>
SPLADE»MonoT5-base	0.750	0.506	0.872	0.737	<b>0.487</b>	0.872	0.748	0.505	0.899	0.731	<b>0.480</b>	0.899
w/ GAR <sub>BM25</sub>	<b>*0.762</b>	0.509	<b>0.888</b>	<b>0.745</b>	<b>0.487</b>	<b>0.893</b>	<b>*0.757</b>	<b>0.509</b>	0.902	<b>0.737</b>	0.479	<b>0.909</b>
w/ GAR <sub>TCT</sub>	*0.759	<b>0.512</b>	0.878	0.737	0.481	0.875	0.751	0.506	<b>0.903</b>	0.734	0.475	0.908
SPLADE»MonoT5-3b	0.761	0.526	0.872	0.764	<b>0.533</b>	0.872	0.774	0.559	0.899	0.775	0.560	0.899
w/ GAR <sub>BM25</sub>	<b>*0.775</b>	0.532	<b>*0.891</b>	<b>0.774</b>	<b>0.533</b>	<b>0.896</b>	<b>*0.780</b>	0.559	0.903	<b>*0.788</b>	<b>0.562</b>	<b>*0.919</b>
w/ GAR <sub>TCT</sub>	*0.773	<b>0.539</b>	0.884	0.769	0.531	0.881	<b>*0.780</b>	<b>0.561</b>	<b>0.905</b>	0.783	0.559	0.910
SPLADE»ColBERT	0.741	0.479	0.872	0.727	<b>0.456</b>	0.872	0.747	0.495	0.899	0.733	0.474	0.899
w/ GAR <sub>BM25</sub>	<b>*0.753</b>	<b>0.490</b>	<b>0.885</b>	<b>0.730</b>	<b>0.456</b>	<b>0.875</b>	<b>*0.755</b>	<b>0.501</b>	0.902	<b>*0.742</b>	<b>*0.477</b>	<b>0.914</b>
w/ GAR <sub>TCT</sub>	*0.750	0.489	0.876	0.727	0.455	0.868	*0.752	0.500	<b>0.903</b>	0.740	*0.476	0.911





c	GAR <sub>TCT</sub>		MonoT5-base
	b = 16	b = 64	Scoring
100	2.68 ± 0.02	0.57 ± 0.01	267.06 ± 6.12
250	8.10 ± 0.05	4.34 ± 0.01	652.30 ± 7.53
500	17.38 ± 0.07	13.66 ± 0.02	1,362.14 ± 5.27
750	26.96 ± 0.12	22.29 ± 0.07	2,047.20 ± 6.71
1000	37.37 ± 0.07	30.82 ± 0.04	2,631.75 ± 6.28

**Table 6: Intra-List Similarity (ILS) among retrieved relevant documents. Since the set of retrieved documents does not change using typical Re-Ranking (RR), each value in this column is only listed once. ILS scores that are statistically equivalent to the RR setting are indicated with \* (procedure described in Section 6.5).**

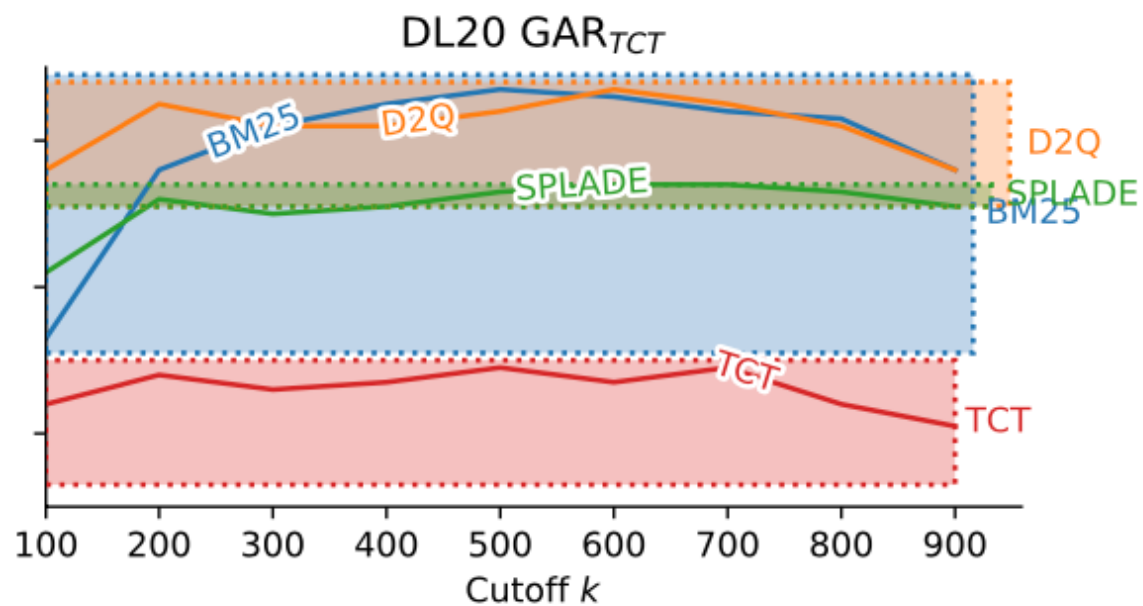
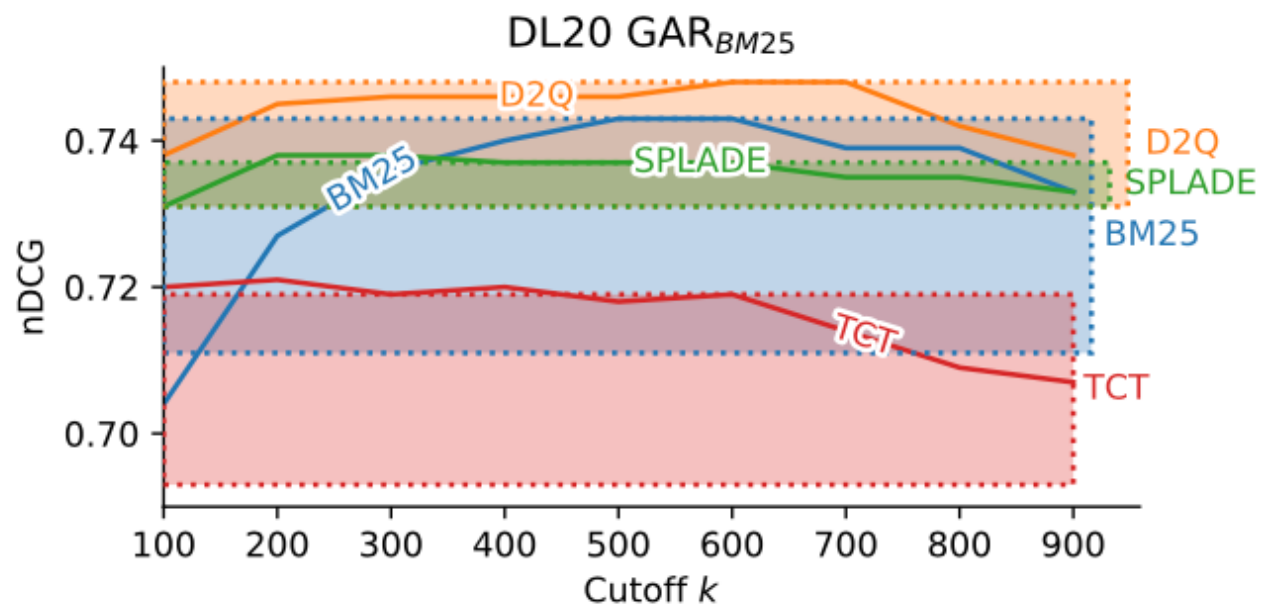
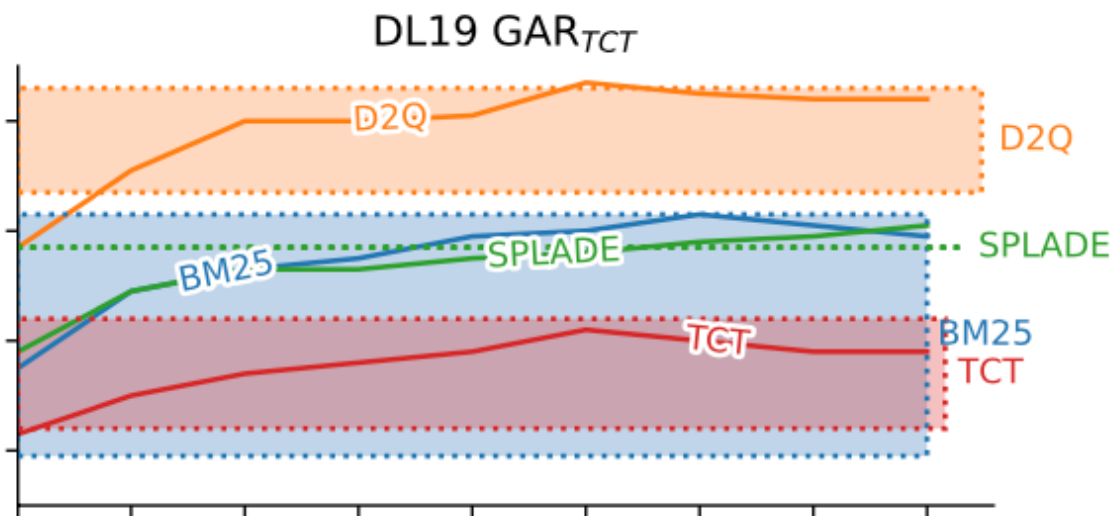
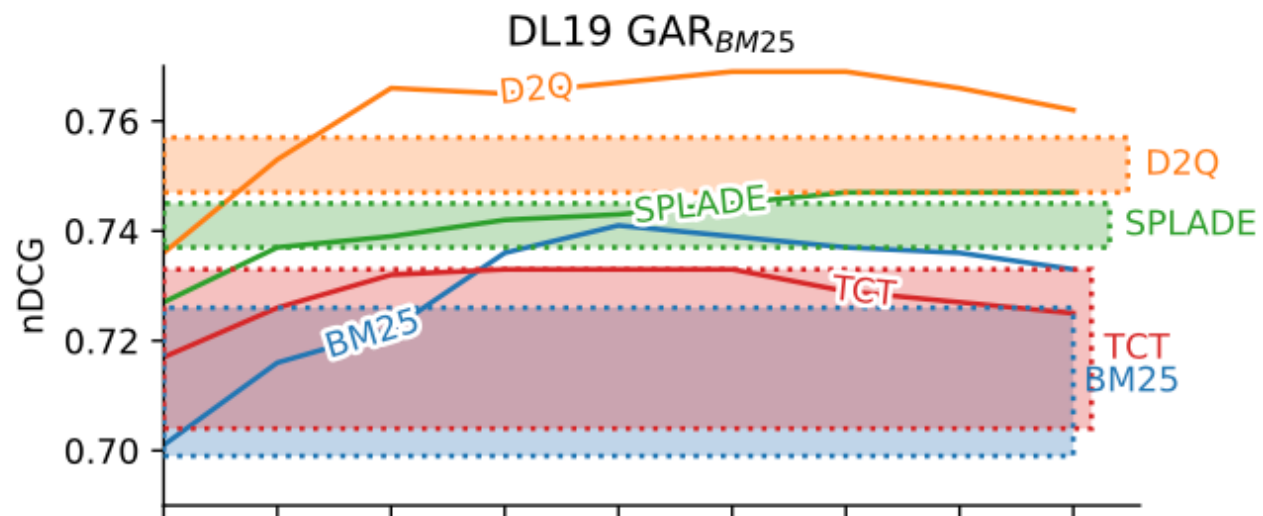
Pipeline	RR	GAR <sub>BM25</sub>		GAR <sub>TCT</sub>	
		c=100	c=1k	c=100	c=1k
BM25»MonoT5-base	0.947	* 0.946	* 0.946	* 0.947	* 0.946
BM25»MonoT5-3b		* 0.946	* 0.946	* 0.946	* 0.946
BM25»ColBERT		* 0.946	* 0.946	* 0.947	* 0.946
TCT»MonoT5-base	0.969	* 0.969	* 0.968	* 0.969	* 0.969
TCT»MonoT5-3b		* 0.969	* 0.968	* 0.969	* 0.969
TCT»ColBERT		* 0.969	* 0.969	* 0.969	* 0.969
D2Q»MonoT5-base	0.969	* 0.968	* 0.968	* 0.969	* 0.968
D2Q»MonoT5-3b		* 0.968	* 0.968	* 0.968	* 0.968
D2Q»ColBERT		* 0.968	* 0.968	* 0.969	* 0.968
SPLADE»MonoT5-base	0.969	* 0.968	* 0.968	* 0.969	* 0.969
SPLADE»MonoT5-3b		* 0.968	* 0.968	* 0.968	* 0.969
SPLADE»ColBERT		* 0.968	* 0.969	* 0.969	* 0.969

Pipeline	Agent	TREC DL 2019 (dev)				TREC DL 2020 (test)			
		GAR <sub>bm25</sub>		GAR <sub>tct</sub>		GAR <sub>bm25</sub>		GAR <sub>tct</sub>	
		nDCG	R@1k	nDCG	R@1k	nDCG	R@1k	nDCG	R@1k
BM25»MonoT5	Non-Adaptive	0.699	0.755	0.699	0.755	0.711	0.805	0.711	0.805
	Oracle	0.747	0.804	0.786	0.853	0.748	0.791	0.768	0.828
	Alternate	0.726	<sup>N</sup> 0.827	<sup>NO</sup> 0.743	<sup>N</sup> 0.839	<sup>N</sup> 0.743	<sup>NO</sup> 0.874	<sup>N</sup> 0.749	<sup>N</sup> 0.892
	TwoPhase-Fixed	<sup>N</sup> 0.729	<sup>N</sup> 0.815	<sup>NO</sup> 0.740	<sup>N</sup> 0.836	<sup>N</sup> 0.732	<sup>NA</sup> 0.838	<sup>N</sup> 0.742	<sup>NA</sup> 0.858
	TwoPhase-Refine	<sup>N</sup> 0.741	<sup>N</sup> 0.826	<sup>NO</sup> 0.743	<sup>N</sup> 0.841	<sup>N</sup> 0.743	<sup>NO</sup> 0.871	<sup>NA</sup> 0.744	<sup>NA</sup> 0.879
	Threshold	<sup>N</sup> 0.742	<sup>N</sup> 0.829	<sup>NOA</sup> 0.751	<sup>N</sup> 0.849	<sup>N</sup> 0.744	<sup>NO</sup> 0.874	<sup>N</sup> 0.744	<sup>NA</sup> 0.874
Greedy	0.723	<sup>N</sup> 0.823	<sup>NO</sup> 0.737	<sup>N</sup> 0.839	<sup>N</sup> 0.743	<sup>NO</sup> 0.868	<sup>N</sup> 0.744	<sup>N</sup> 0.882	
TCT»MonoT5	Non-Adaptive	0.704	0.830	0.704	0.830	0.693	0.848	0.693	0.848
	Oracle	0.793	0.891	0.766	0.846	0.762	0.874	0.754	0.861
	Alternate	<sup>NO</sup> 0.733	<sup>N</sup> 0.883	<sup>NO</sup> 0.724	<sup>N</sup> 0.866	<sup>NO</sup> 0.719	<sup>N</sup> 0.881	<sup>NO</sup> 0.710	<sup>N</sup> 0.871
	TwoPhase-Fixed	<sup>NO</sup> 0.733	<sup>N</sup> 0.874	<sup>NO</sup> 0.719	<sup>N</sup> 0.857	<sup>NO</sup> 0.717	<sup>N</sup> 0.877	<sup>NO</sup> 0.710	<sup>N</sup> 0.868
	TwoPhase-Refine	<sup>NO</sup> 0.733	<sup>N</sup> 0.882	<sup>NO</sup> 0.722	0.859	<sup>NO</sup> 0.719	<sup>N</sup> 0.883	<sup>NOA</sup> 0.707	<sup>A</sup> 0.866
	Threshold	<sup>NO</sup> 0.731	<sup>N</sup> 0.886	<sup>NO</sup> 0.720	<sup>N</sup> 0.866	<sup>NOA</sup> 0.711	0.871	<sup>NOA</sup> 0.705	0.862
Greedy	<sup>NO</sup> 0.731	<sup>N</sup> 0.881	<sup>NO</sup> 0.725	<sup>N</sup> 0.871	<sup>NOA</sup> 0.713	<sup>N</sup> 0.873	<sup>NO</sup> 0.708	<sup>N</sup> 0.868	
D2Q»MonoT5	Non-Adaptive	0.747	0.830	0.747	0.830	0.731	0.839	0.731	0.839
	Oracle	0.797	0.867	0.798	0.867	0.791	0.884	0.793	0.889
	Alternate	0.757	<sup>N</sup> 0.880	<sup>NO</sup> 0.766	<sup>N</sup> 0.879	<sup>NO</sup> 0.748	<sup>N</sup> 0.880	<sup>O</sup> 0.748	<sup>N</sup> 0.895
	TwoPhase-Fixed	<sup>N</sup> 0.765	<sup>N</sup> 0.866	<sup>NO</sup> 0.765	<sup>N</sup> 0.870	<sup>NO</sup> 0.748	<sup>NA</sup> 0.867	<sup>O</sup> 0.745	<sup>NA</sup> 0.870
	TwoPhase-Refine	<sup>N</sup> 0.769	<sup>N</sup> 0.875	<sup>N</sup> 0.767	<sup>N</sup> 0.878	<sup>NO</sup> 0.748	<sup>N</sup> 0.877	<sup>O</sup> 0.747	<sup>N</sup> 0.892
	Threshold	<sup>N</sup> 0.766	<sup>N</sup> 0.876	<sup>NO</sup> 0.767	<sup>N</sup> 0.877	<sup>O</sup> 0.746	<sup>NA</sup> 0.874	<sup>O</sup> 0.745	<sup>N</sup> 0.881
Greedy	0.754	<sup>N</sup> 0.874	<sup>O</sup> 0.757	<sup>N</sup> 0.873	<sup>O</sup> 0.744	<sup>N</sup> 0.878	<sup>O</sup> 0.748	<sup>N</sup> 0.894	
SPLADE»MonoT5	Non-Adaptive	0.737	0.872	0.737	0.872	0.731	0.899	0.731	0.899
	Oracle	0.807	0.898	0.783	0.859	0.777	0.886	0.781	0.899
	Alternate	<sup>O</sup> 0.745	0.893	<sup>O</sup> 0.737	0.875	<sup>O</sup> 0.737	<b>0.909</b>	<sup>O</sup> 0.734	<b>0.908</b>
	TwoPhase-Fixed	<sup>O</sup> 0.763	0.863	<b>0.764</b>	0.869	<b>0.748</b>	<sup>A</sup> 0.868	<sup>O</sup> 0.742	<sup>NA</sup> 0.867
	TwoPhase-Refine	<sup>O</sup> 0.769	0.875	<b>0.764</b>	0.870	<sup>O</sup> 0.748	0.877	<sup>O</sup> 0.736	<sup>A</sup> 0.869
	Threshold	<sup>O</sup> 0.766	0.871	0.759	0.857	<sup>O</sup> 0.746	0.874	<sup>O</sup> 0.744	<sup>NA</sup> 0.865
Greedy	<sup>NO</sup> 0.747	<sup>N</sup> 0.895	<sup>O</sup> 0.740	<b>0.882</b>	<sup>O</sup> 0.734	0.903	<sup>O</sup> 0.734	0.906	

**Table 1**

Re-ranking performance on TREC Deep Learning 2019 and 2020 using various agents. The best-performing (non-oracle) agent in section is listed in bold. Significant differences compared to the Non-Adaptive, Oracle, and Adaptive systems are marked with <sup>NOA</sup>, respectively (paired t-test,  $p < 0.05$ ).

## TwoPhase-Refine



# Threshold

