

# Efficient Neural Models for Representing, Indexing, and Retrieving Documents

Hamed Zamani

Center for Intelligent Information Retrieval  
University of Massachusetts Amherst

[zamani@cs.umass.edu](mailto:zamani@cs.umass.edu)





These efficient contextual models have delivered the **largest quality improvement** to the Bing search results in **the last 6 years!**

In this talk...



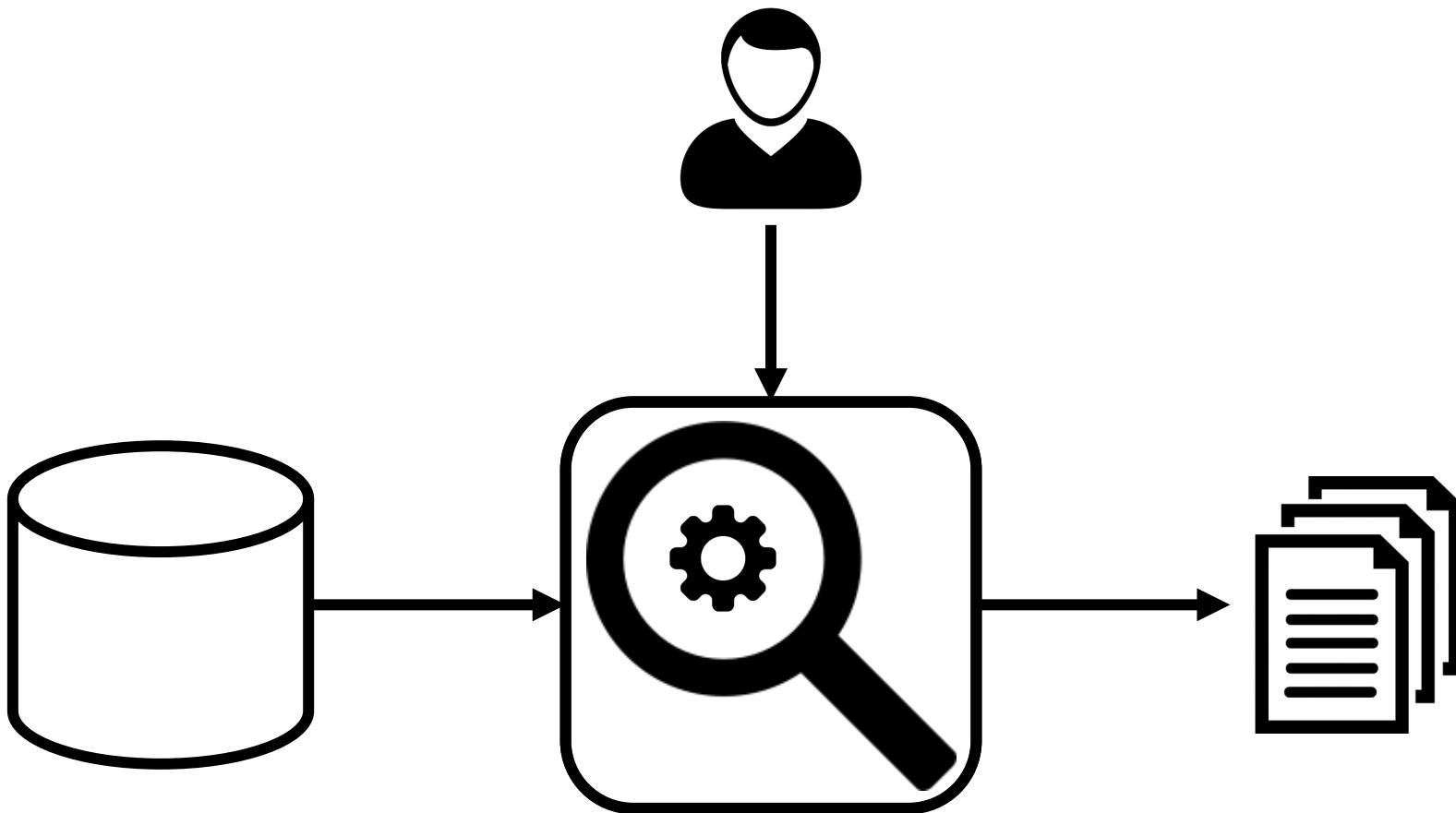
How to efficiently represent long documents?



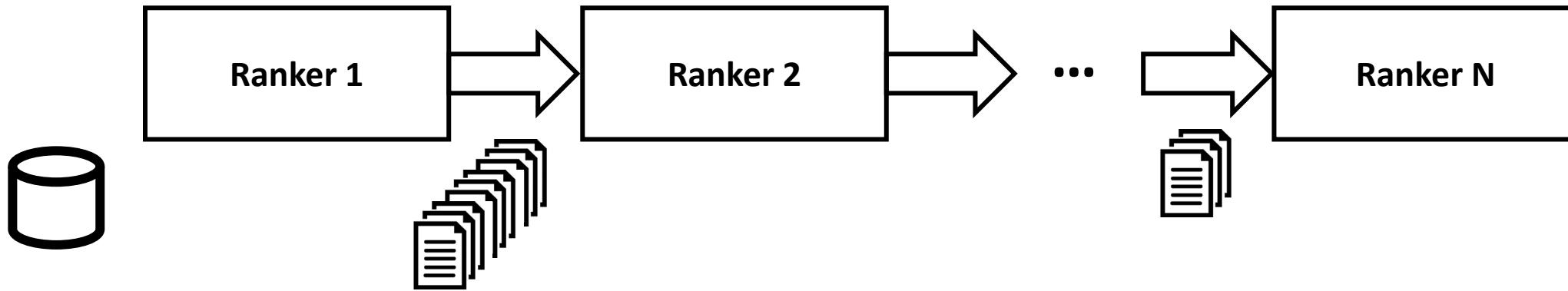
How to efficiently retrieve documents?

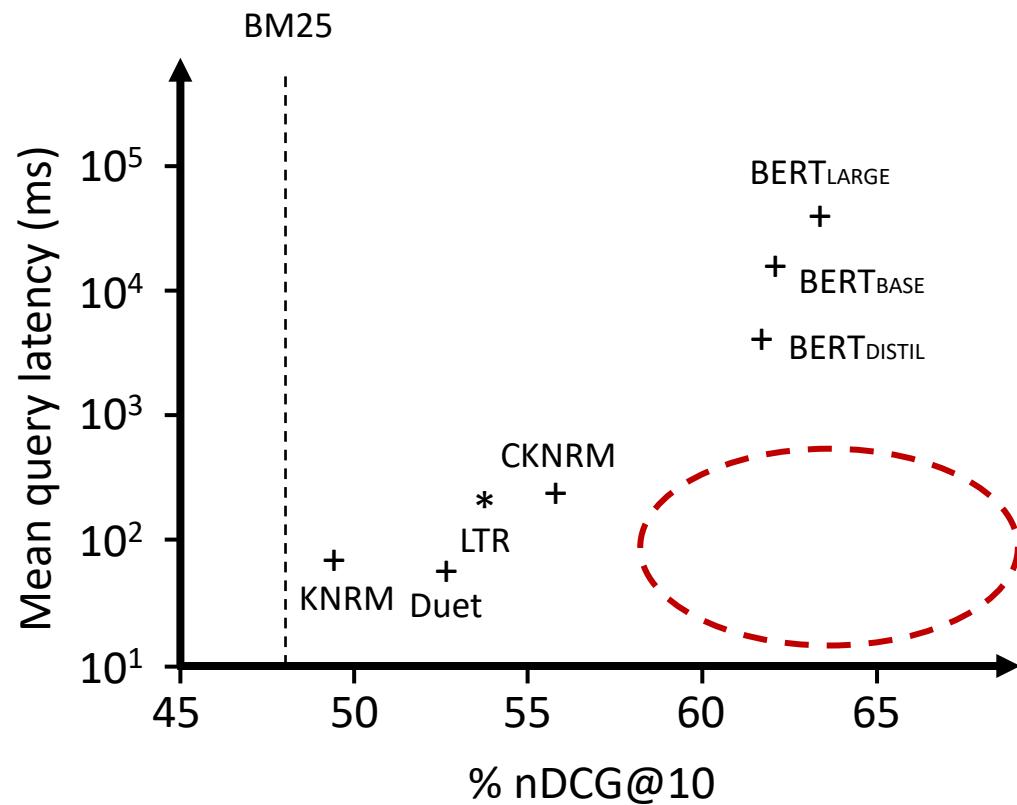


How can efficient neural IR advance machine learning research?



# Multi-Stage Cascaded Retrieval



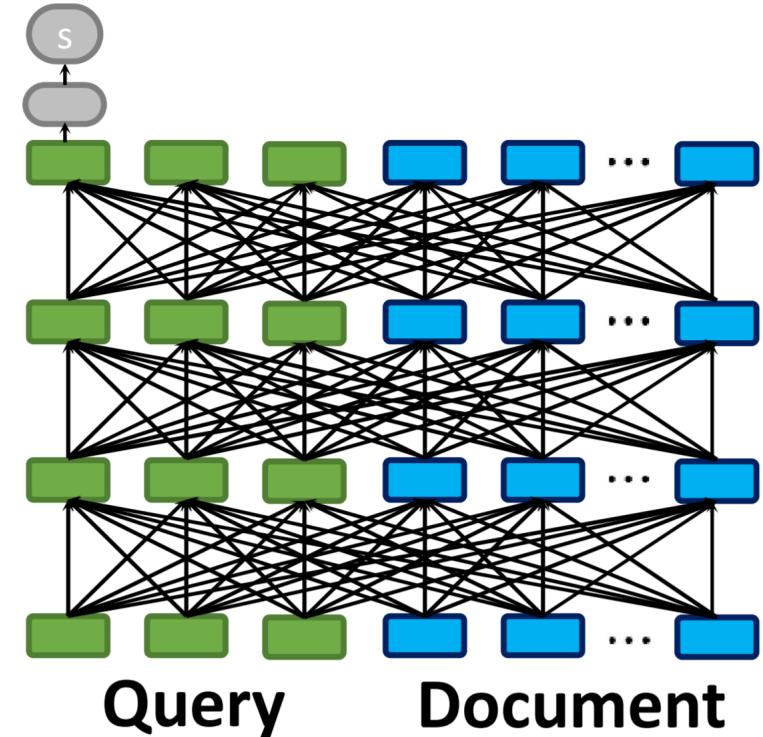


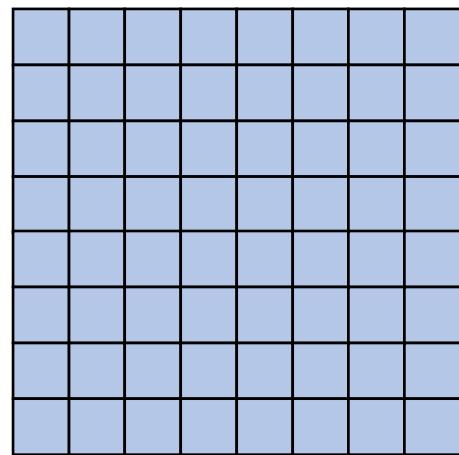
TREC Deep Learning Track 2019 – Document Ranking

# BERT for Document Ranking

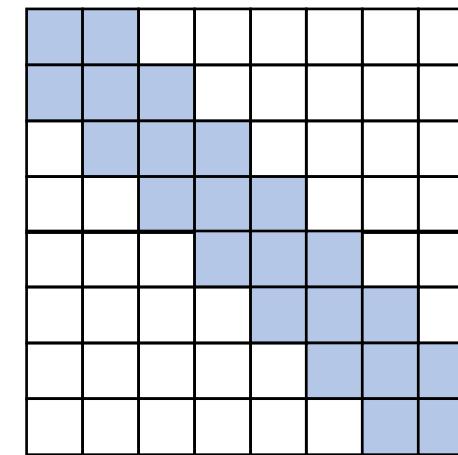
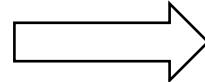
Reduce online computations

Reduce the quadratic complexity of self-attention in Transformer





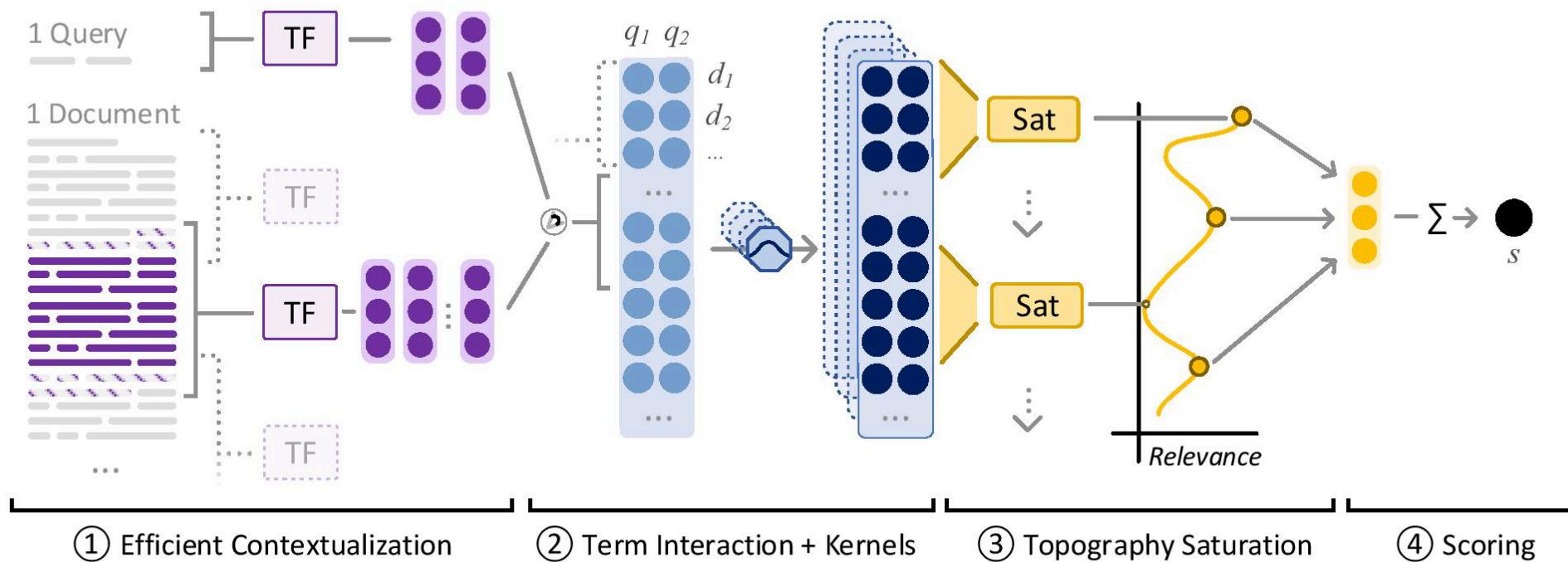
Transformer



Conformer

Linformer, Longformer, Sparse Transformer, etc.

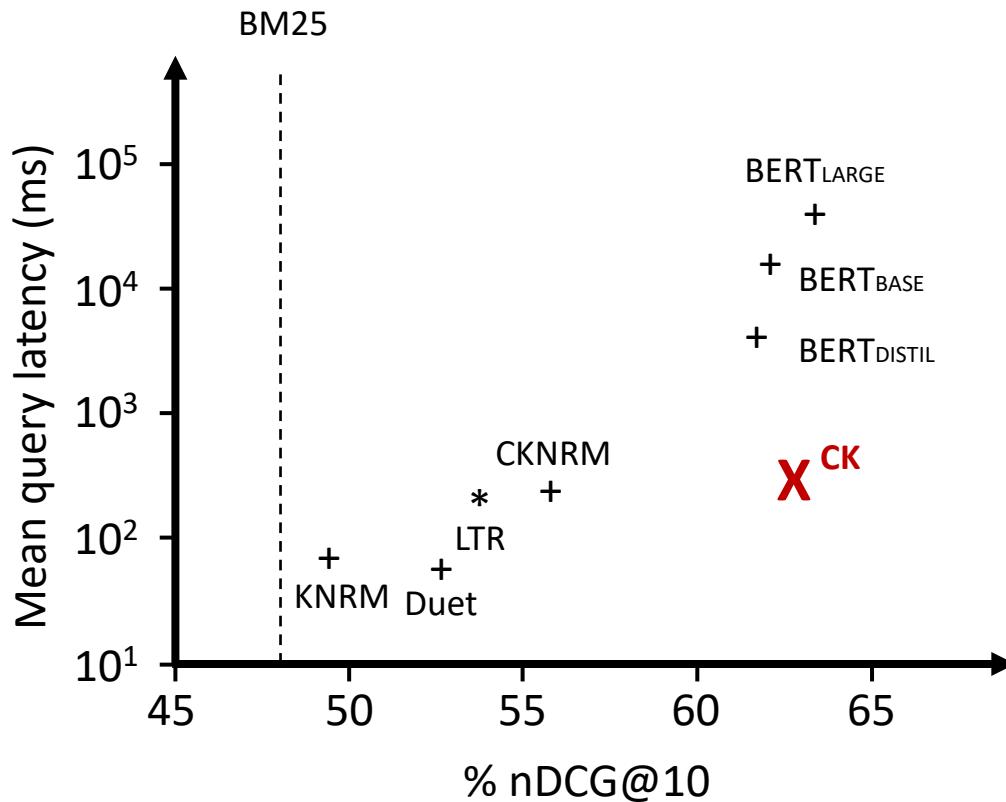
# Conformer-Kernel



[Hofstätter, Zamani, Mitra, Craswell, Hanbury; SIGIR 2020]

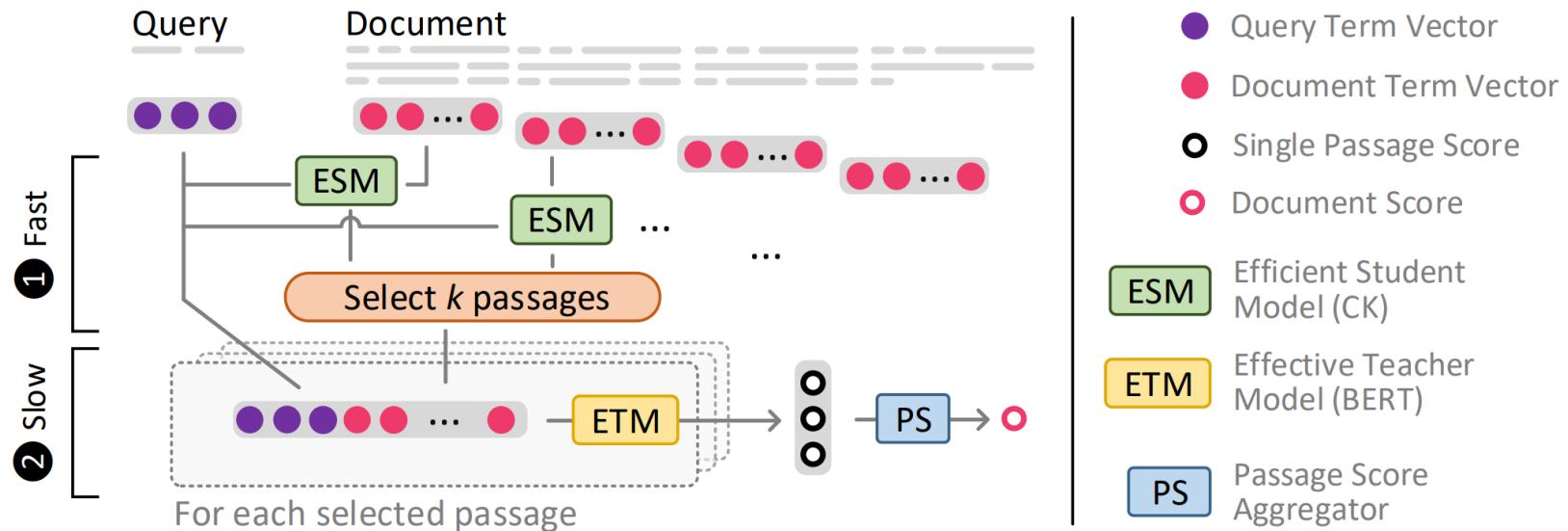
[Mitra, Hofstätter, Zamani, Craswell; SIGIR 2021]

# Results

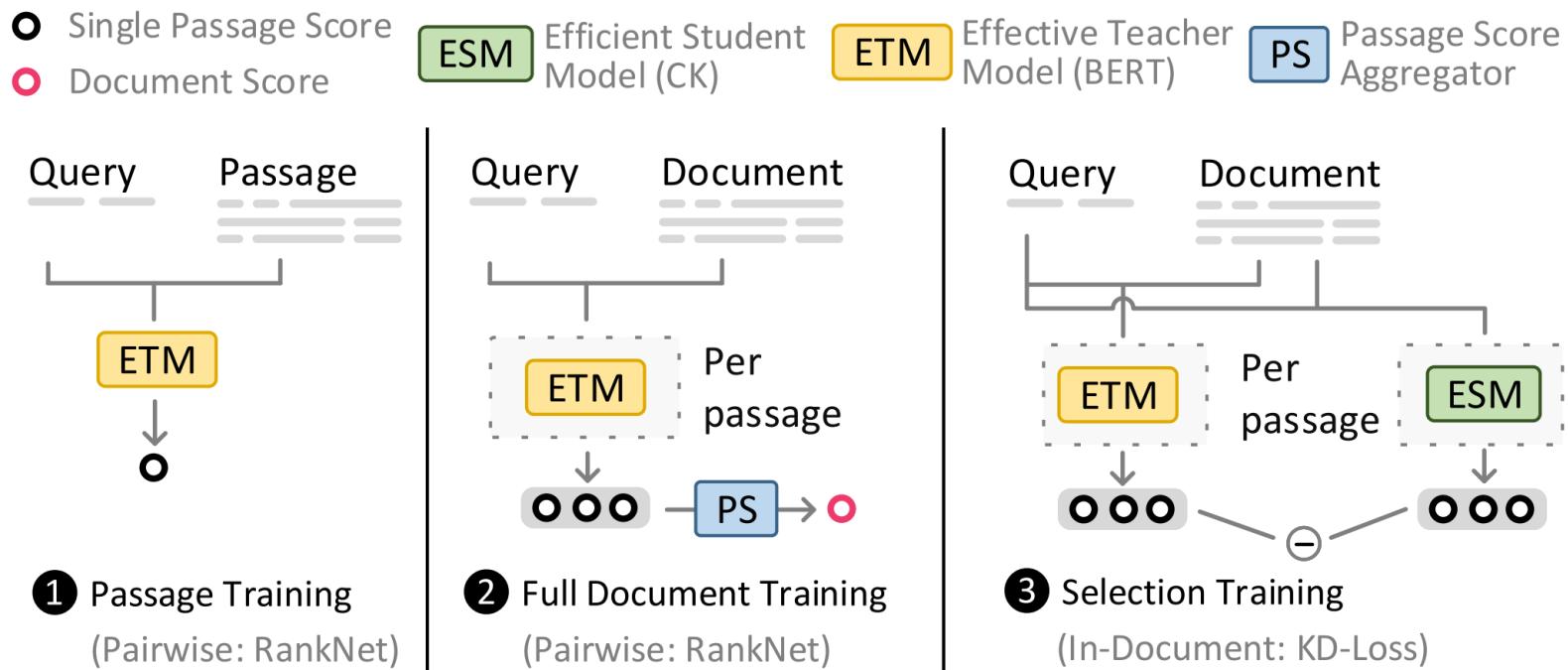


**2020**  
TREC Deep Learning Track ~~2019~~ – Document Ranking

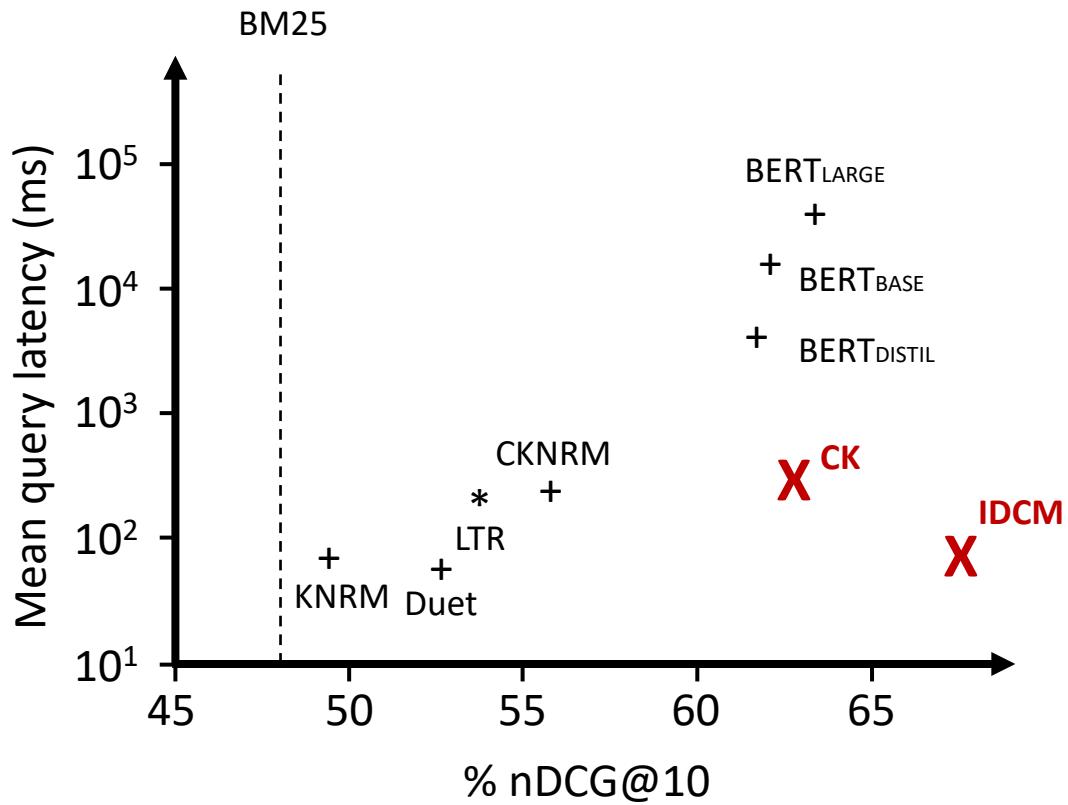
# Intra-Document Cascading Model (IDCM)



# Training IDCM



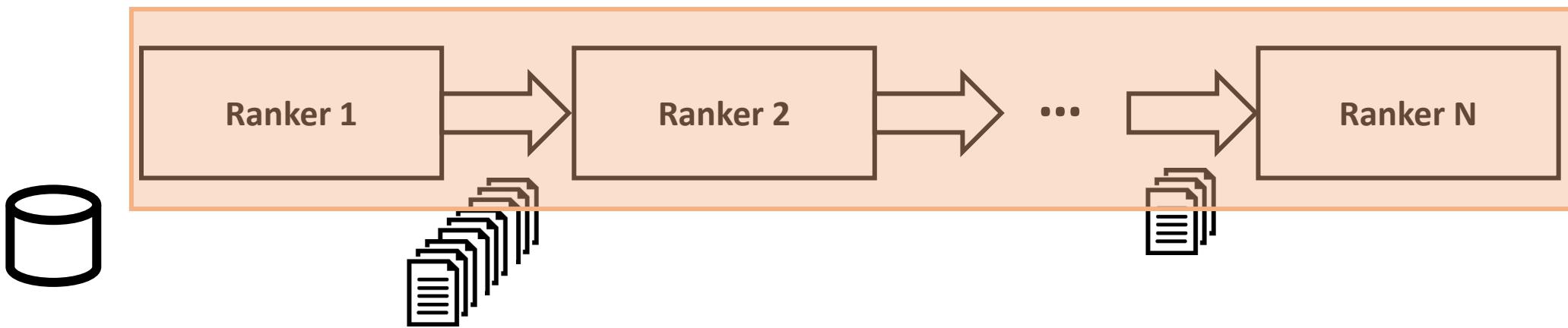
# Results



TREC Deep Learning Track 2019 – Document Ranking

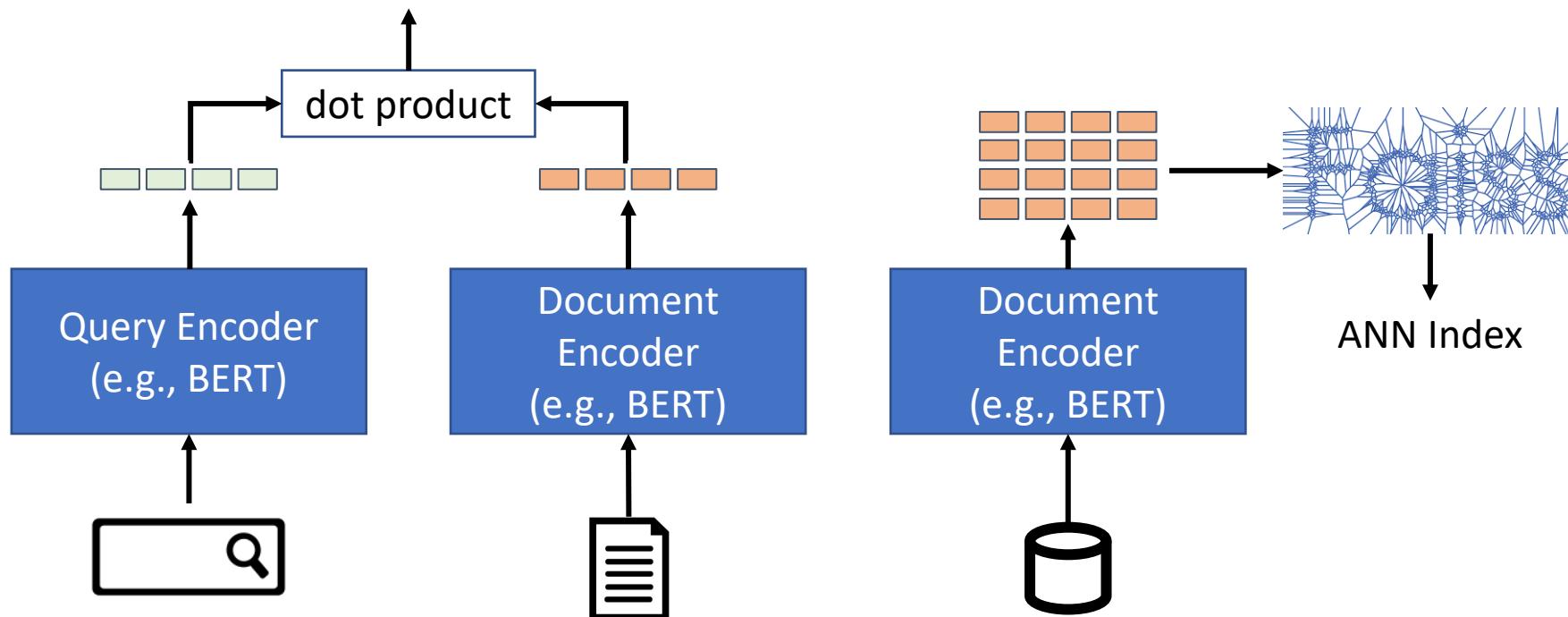
# How to efficiently retrieve documents?

# Multi-Stage Cascaded Retrieval

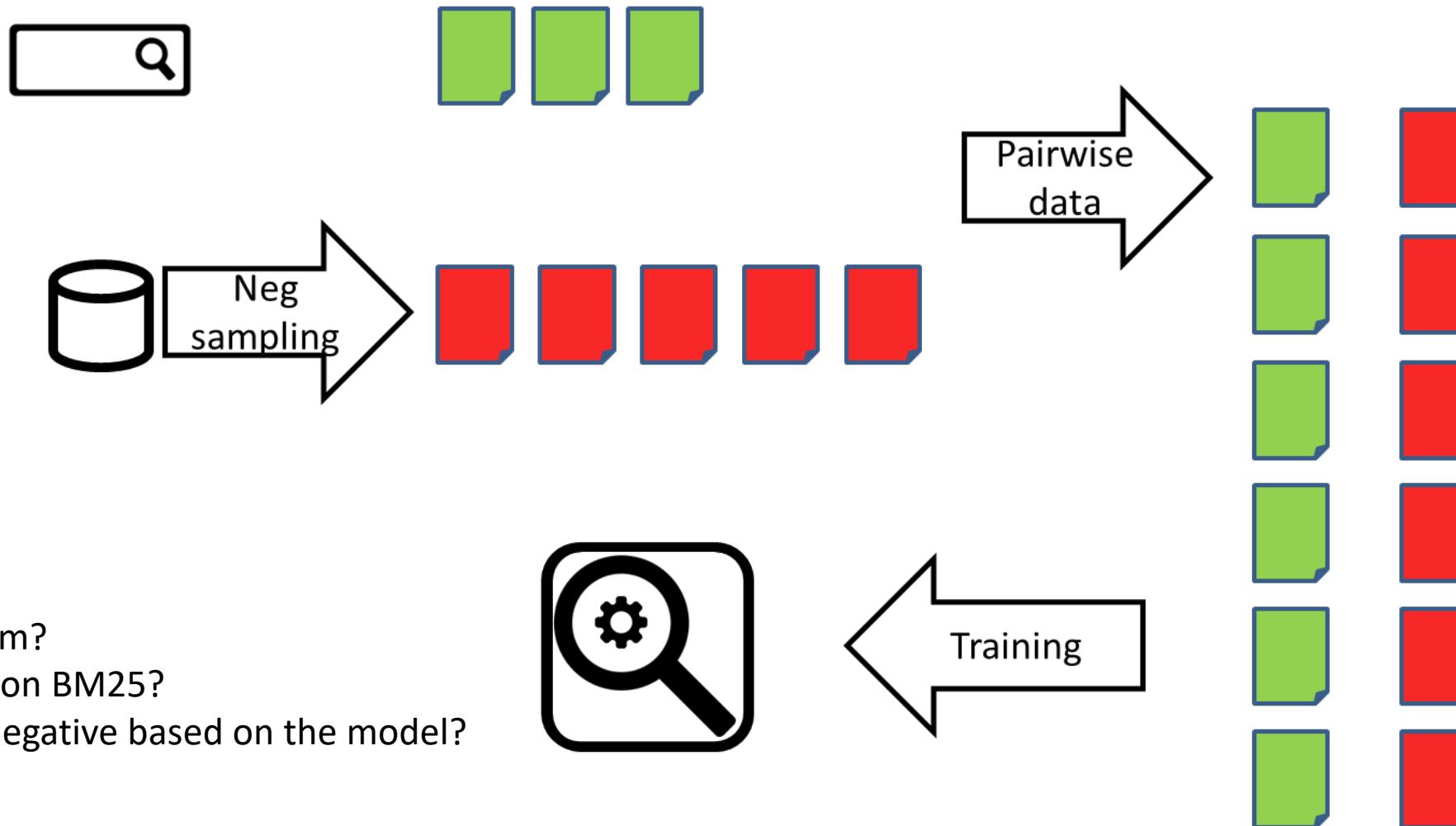


**Can we replace multi-stage cascaded architectures with  
standalone retrieval models?**

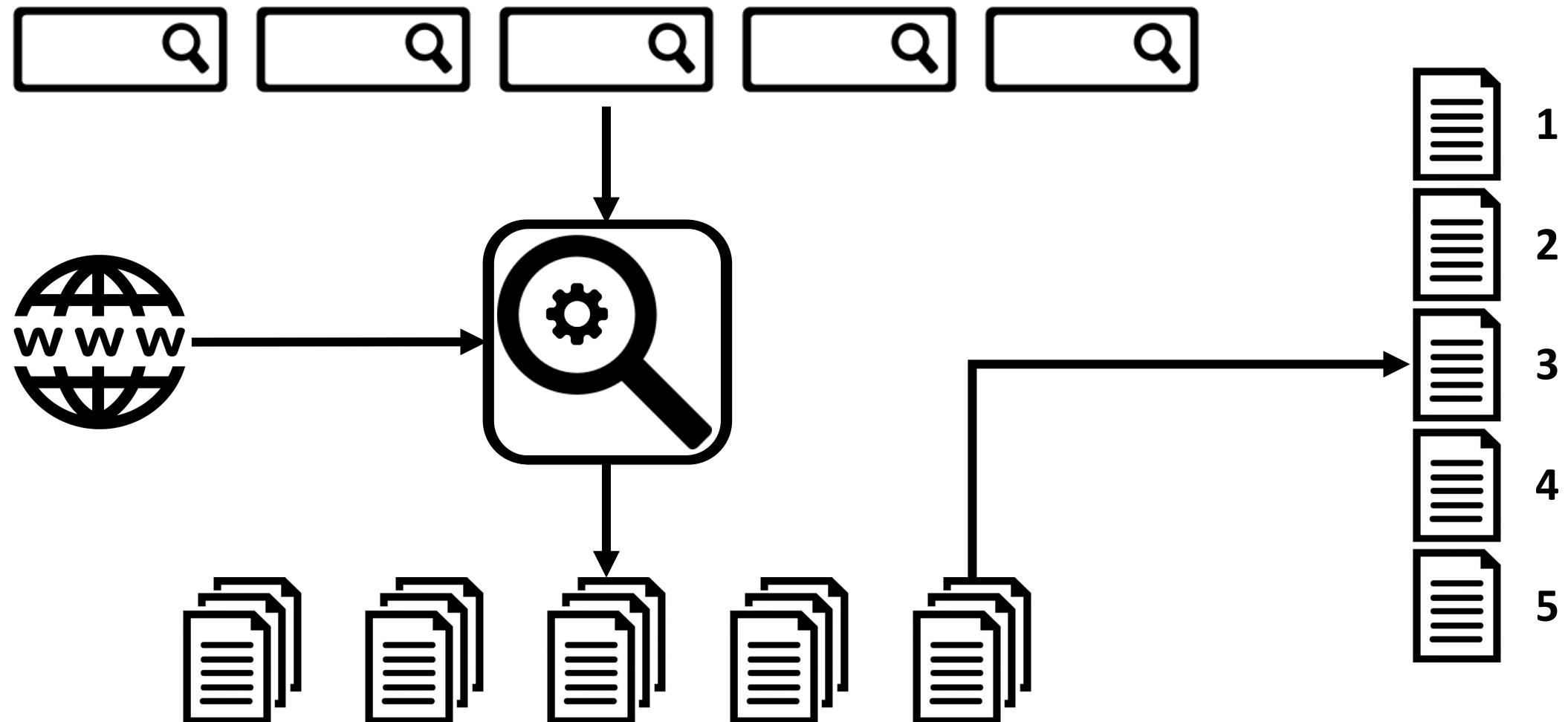
# Dense Retrieval



# Negative Sampling for Neural IR



# Teacher-Student Learning: A Common Recipe for Training Dense Retrieval Models



# Curriculum Learning for Dense Retrieval Distillation

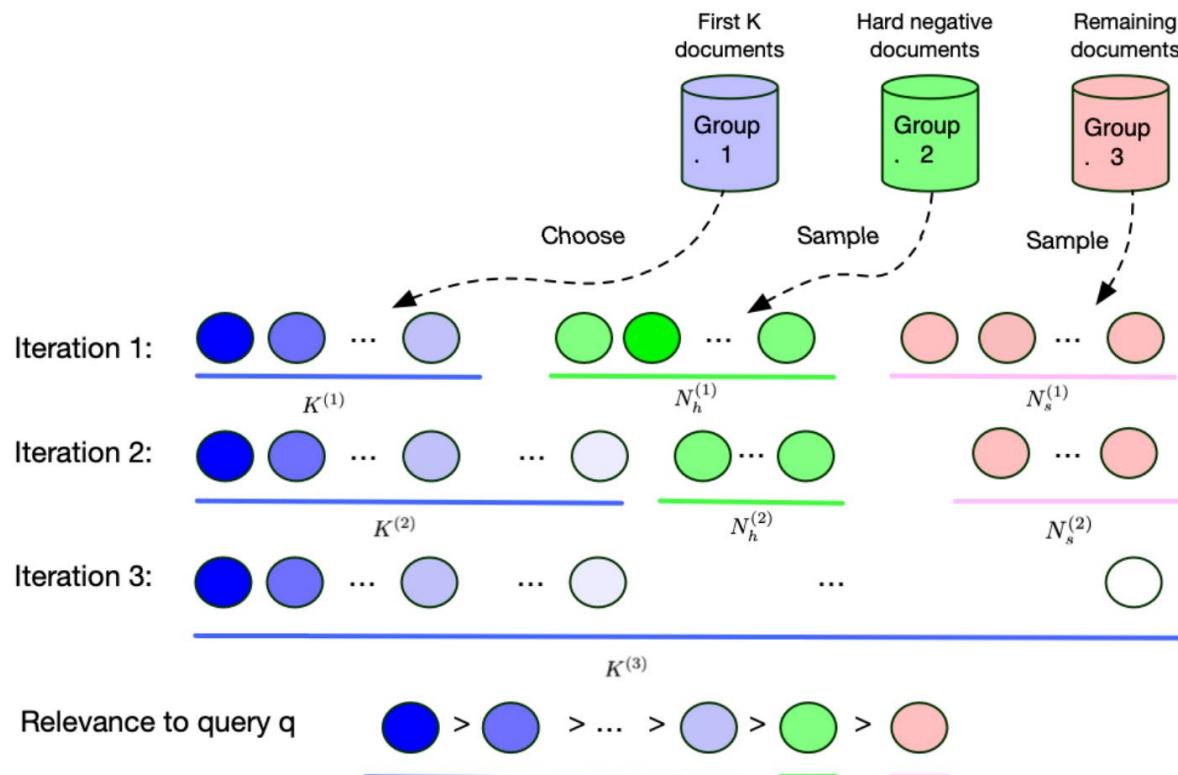
---

**Algorithm 1** The Iterative Optimization Process in CL-DRD.

---

- 1: **Input** (a) training query set  $Q$ ; (b) document collection  $C$ ; (c)  
*optional* relevance judgment set  $R$ ; (d) teacher model  $\widehat{M}$ .
  - 2: **Output** dense retrieval model  $M_\theta$ .
  - 3: **Initialize**  $\theta$ .
  - 4:  $\delta \leftarrow 1$  ▷  $\delta$  denotes the difficulty level in CL data
  - 5: **repeat**
  - 6:    $D_\delta \leftarrow \text{RankingDataGeneration}(\delta; \widehat{M}, Q, C, R)$
  - 7:    $\theta^* \leftarrow \arg \min_\theta \mathcal{L}_{KD}(M_\theta, D_\delta)$
  - 8:    $\delta \leftarrow \delta + 1$
  - 9: **until** stopping criterion is met
  - 10: **return**  $M_{\theta^*}$
-

# Curriculum Learning for Dense Retrieval Distillation



# Results

Model	KD	Encoder	#params	MS MARCO DEV		TREC-DL'19		TREC-DL'20	
				MRR@10	MAP@1k	nDCG@10	MAP@1k	nDCG@10	MAP@1k
<b>Sparse Retrieval</b>									
BM25 [27]	-	-	-	.187	.196	.497	.290	.487	.288
DeepCT [6]	-	-	-	.243	.250	.550	.341	.556	.343
docT5query [20]	-	-	-	.272	.281	.642	.403	.619	.407
<b>Multi-Vector Dense Retrieval</b>									
ColBERT [13]	✗	BERT-Base	110M	.360	-	-	-	-	-
ColBERTv2 [30]	✓	BERT-Base	110M	<b>.397</b>	-	-	-	-	-
ColBERTv2 [30]	✓	DistilBERT	66M	.384	.389	<b>.733</b>	.464	.712	.473
ColBERTv2 + CL-DRD (Ours)	✓	DistilBERT	66M	.394*†‡§	.398*†‡§¶	.727*†‡§	.472*†‡§¶	.717*†‡§¶	.487*†‡§¶
<b>Single-Vector Dense Retrieval</b>									
ANCE [32]	✗	BERT-Base	110M	.330	.336	.648	.371	.646	.408
ADORE [33]	✗	BERT-Base	110M	.347	.352	.683	.419	.666	.442
RocketQA [25]	✓	ERNIE-Base	110M	.370	-	-	-	-	-
TCT-ColBERT [16]	✓	BERT-Base	110M	.335	.342	.670	.391	.668	.430
Margin-MSE [10]	✓	DistilBERT	66M	.325	.331	.699	.405	.645	.416
TAS-B [11]	✓	DistilBERT	66M	.344	.351	.717	.447	.685	.455
TAS-B + CL-DRD (Ours)	✓	DistilBERT	66M	.382*†‡§	.386*†‡§	.725*†‡§	.453*†‡	.687*†‡	.465*†‡§



## Inverted Index

aardvark →  $d_1, d_{10}, d_{1501}, \dots$

abase →  $d_{603}, d_{415}, d_{10493}, \dots$

.

.

.

reneuir →  $d_{213}, d_{716}, d_{1003}, \dots$

.

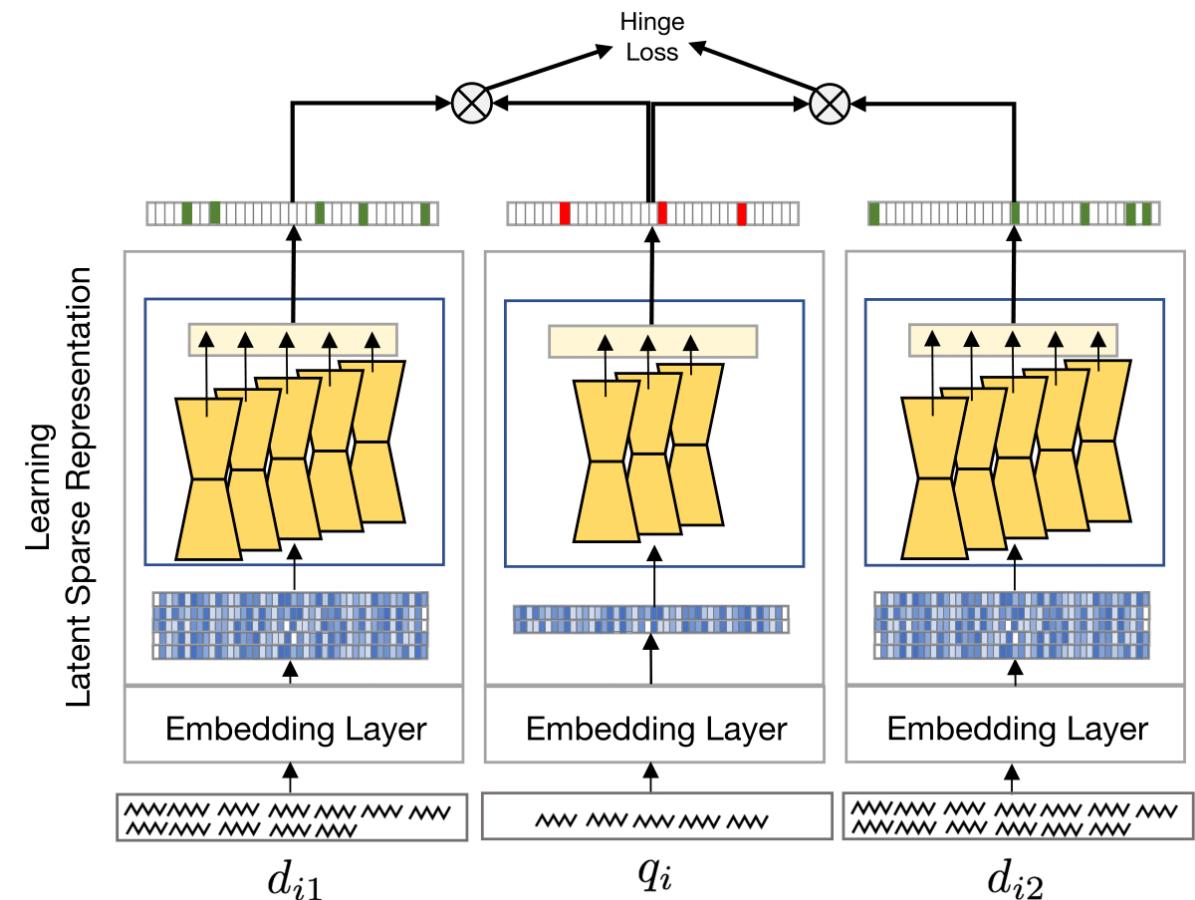
.

.

$\dots, d_{3125}, \dots$

The *sparsity* nature of natural languages enables us to do efficient retrieval using *inverted index*!

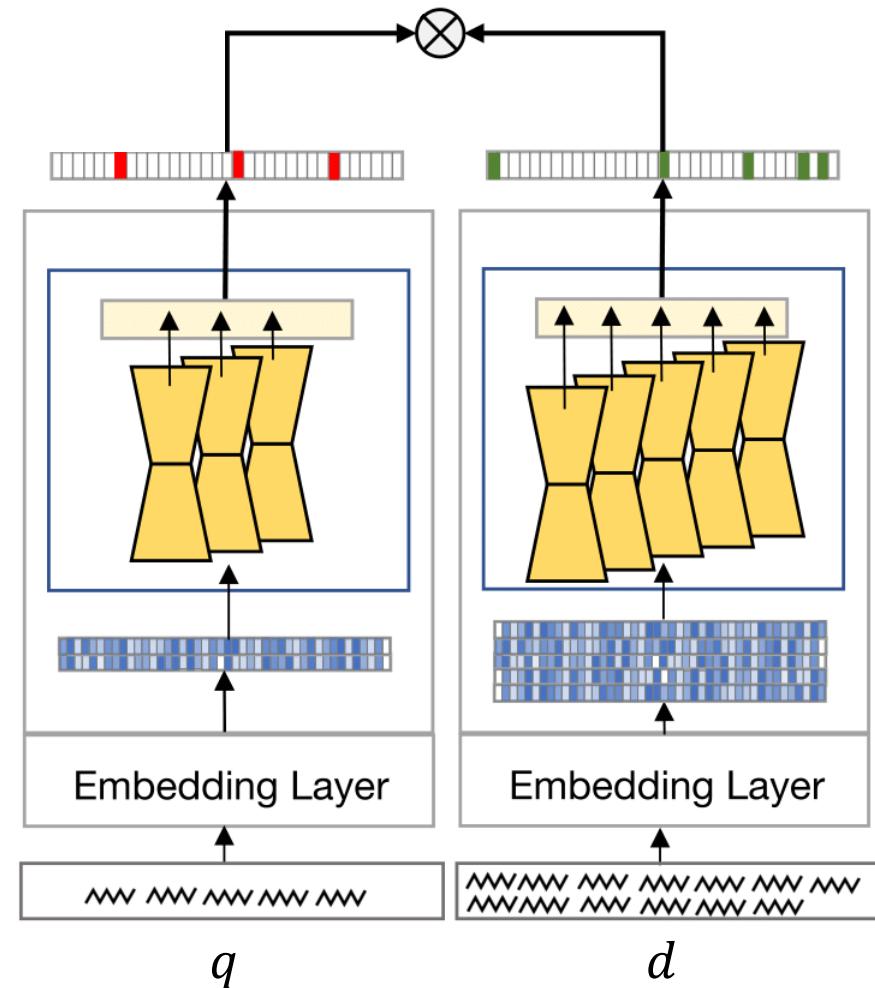
# Standalone Neural Retrieval Model



# Training

## Efficiency -> Sparsity

- Maximizing sparsity ratio, i.e., equivalent to minimizing  $L_0(\vec{v}) = \sum_{i=1}^{|v|} |\vec{v}_i|^0$ .
- A differentiable approximator for  $L_0$ .



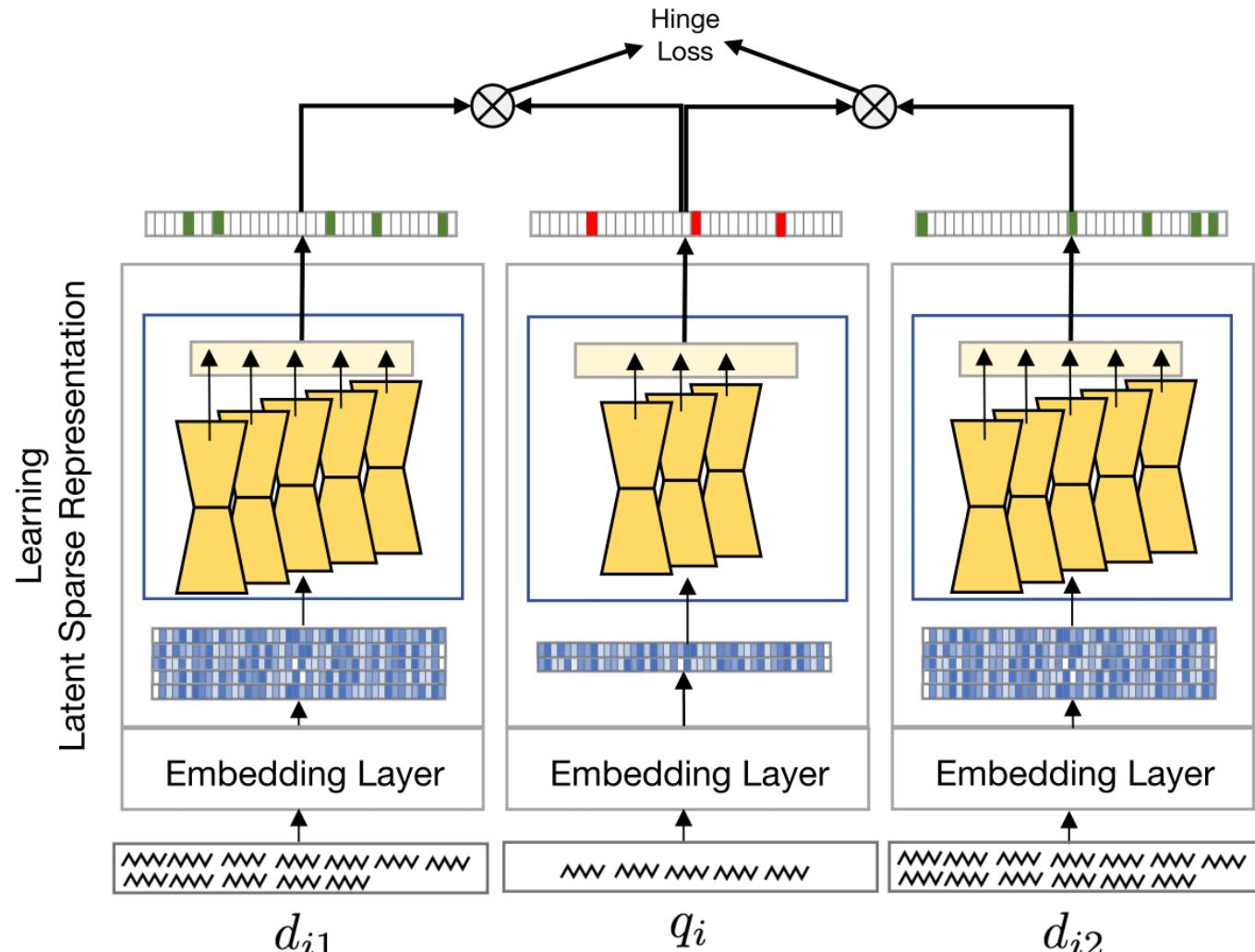
# Training

## Efficiency -> Sparsity

- Maximizing sparsity ratio, i.e., equivalent to minimizing  $L_0(\vec{v}) = \sum_{i=1}^{|v|} |\vec{v}_i|^0$ .
- A differentiable approximator for  $L_0$ .

## Effectiveness

- A learning-to-rank loss function (e.g., Hinge loss)

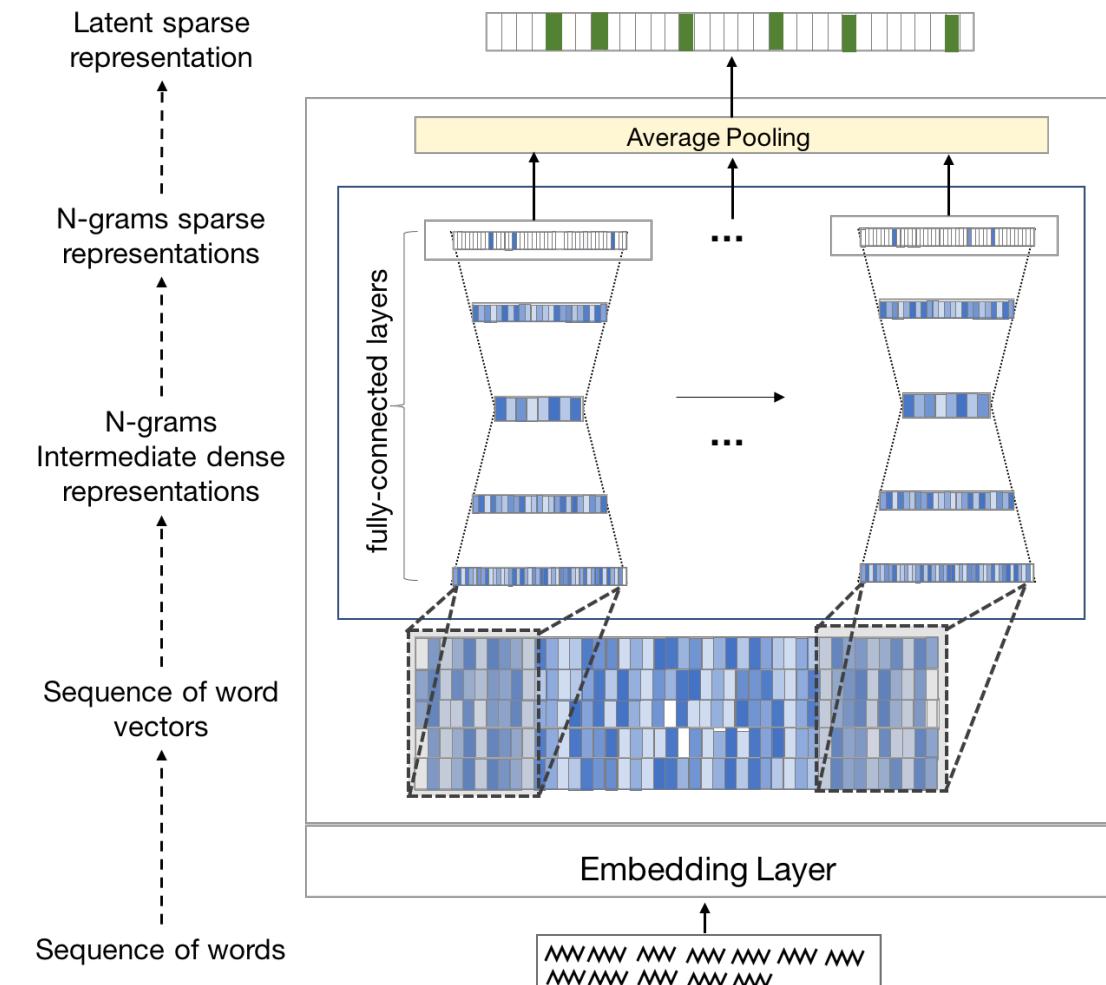
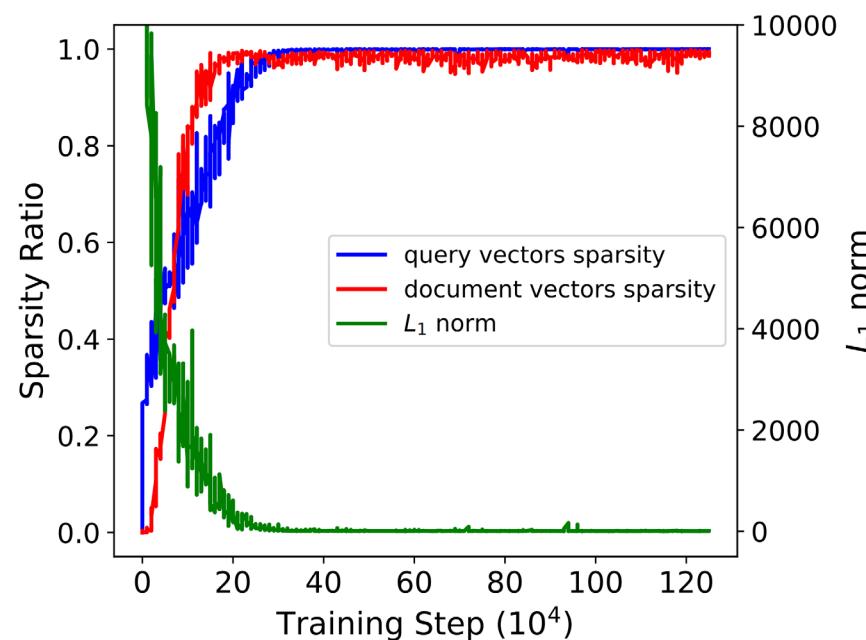


$$\vec{d} = \frac{1}{|d| - n + 1} \sum_{i=1}^{|d|-n+1} \phi_{\text{ngram}}(w_i, w_{i+1}, \dots, w_{i+n-1})$$

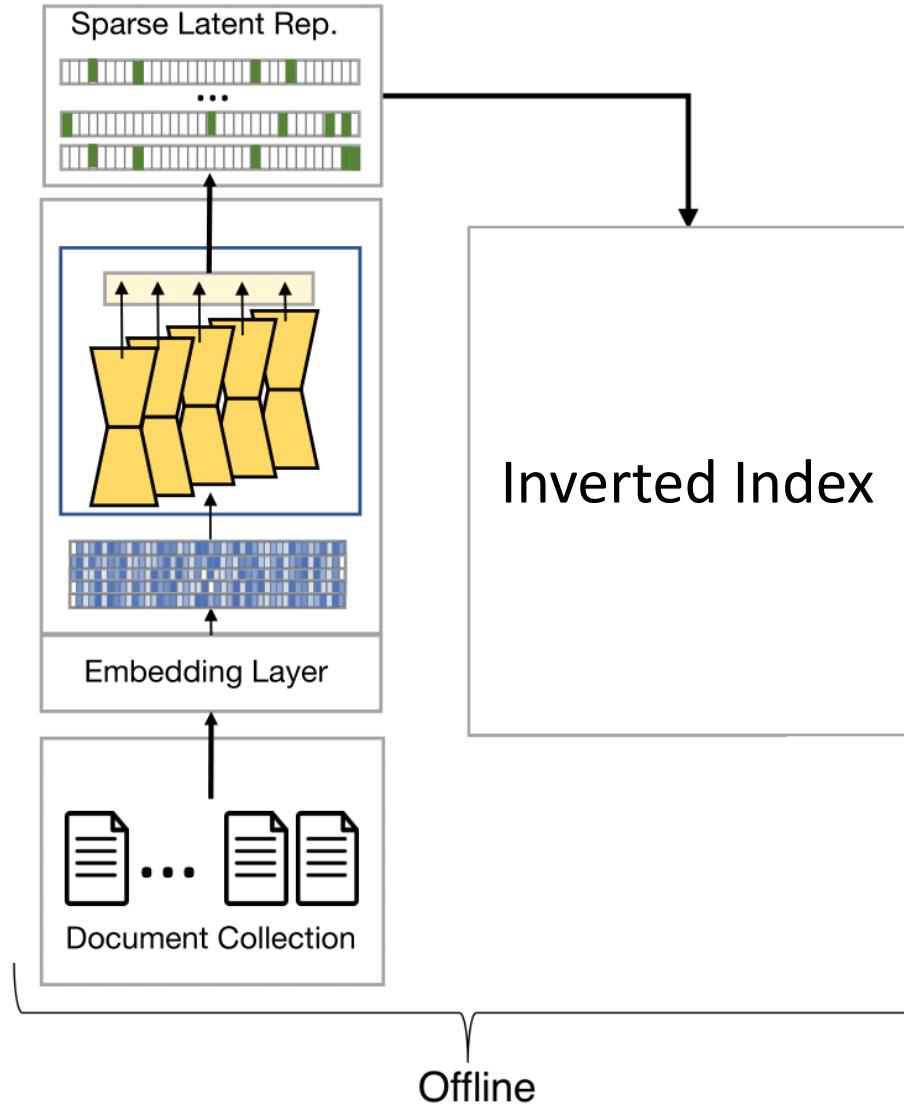
$|q| - n + 1 << |d| - n + 1$  and  
 $\phi_{\text{ngram}}$  is shared.

# of non-zero dimensions (out of 10,000) per document.

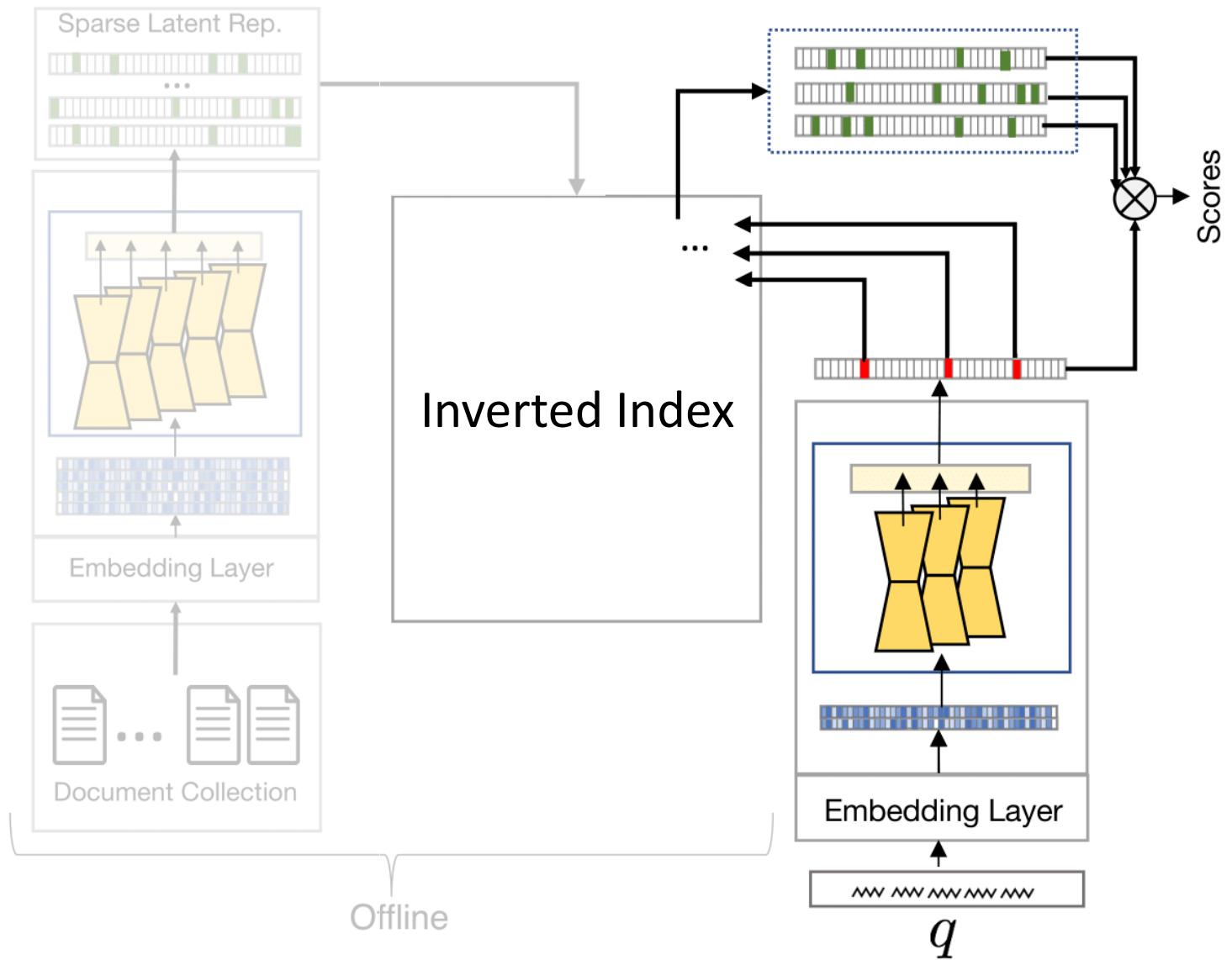
# Unique latent terms...	Robust		ClueWeb	
	Mean	Std. dev.	Mean	Std. dev.
per document	97.96	447.57	130.24	561.53
per query	3.37	3.04	3.87	4.51



# Inverted Index Construction



# Query Time



# Effectiveness

Method	Robust				ClueWeb			
	MAP	P@20	nDCG@20	Recall	MAP	P@20	nDCG@20	Recall
QL	0.2499	0.3556	0.4143	0.6820	0.1044	0.3139	0.2294	0.3286
SDM	0.2524	0.3679 <sup>1</sup>	0.4242 <sup>1</sup>	0.6858	0.1078	0.3141	0.2320	0.3385 <sup>1</sup>
RM3	0.2865 <sup>12</sup>	0.3773 <sup>12</sup>	0.4295 <sup>12</sup>	0.7494 <sup>12</sup>	0.1068	0.3157	0.2309	0.3298
FNRM	0.2815 <sup>12</sup>	0.3752 <sup>12</sup>	0.4327 <sup>12</sup>	0.7234 <sup>12</sup>	0.1329 <sup>123</sup>	0.3351 <sup>123</sup>	0.2392 <sup>13</sup>	0.3426 <sup>123</sup>
CNRM	0.2801 <sup>12</sup>	0.3764 <sup>12</sup>	0.4341 <sup>123</sup>	0.7183 <sup>12</sup>	0.1286 <sup>123</sup>	0.3317 <sup>123</sup>	0.2337 <sup>1</sup>	0.3345 <sup>13</sup>
SNRM	0.2856 <sup>12</sup>	0.3766 <sup>12</sup>	0.4310 <sup>12</sup>	0.7481 <sup>1245</sup>	0.1290 <sup>123</sup>	0.3336 <sup>123</sup>	0.2351 <sup>13</sup>	0.3393 <sup>135</sup>
SNRM with PRF	<b>0.2971</b> <sup>123456</sup>	<b>0.3948</b> <sup>123456</sup>	<b>0.4391</b> <sup>123456</sup>	<b>0.7716</b> <sup>123456</sup>	<b>0.1475</b> <sup>123456</sup>	<b>0.3461</b> <sup>123456</sup>	<b>0.2482</b> <sup>123456</sup>	<b>0.3618</b> <sup>123456</sup>

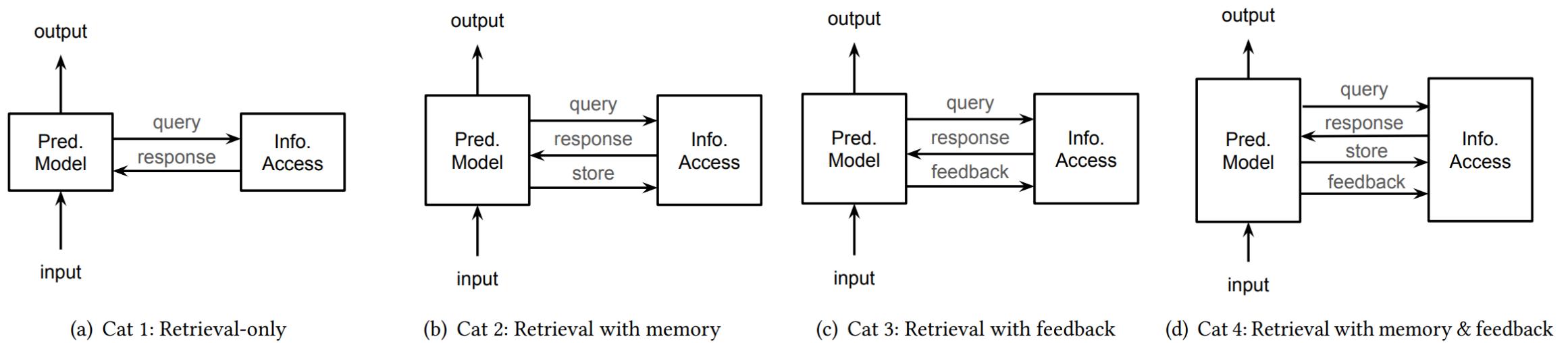
# Efficiency

**Average running time per query in milliseconds for the ClueWeb collection.**

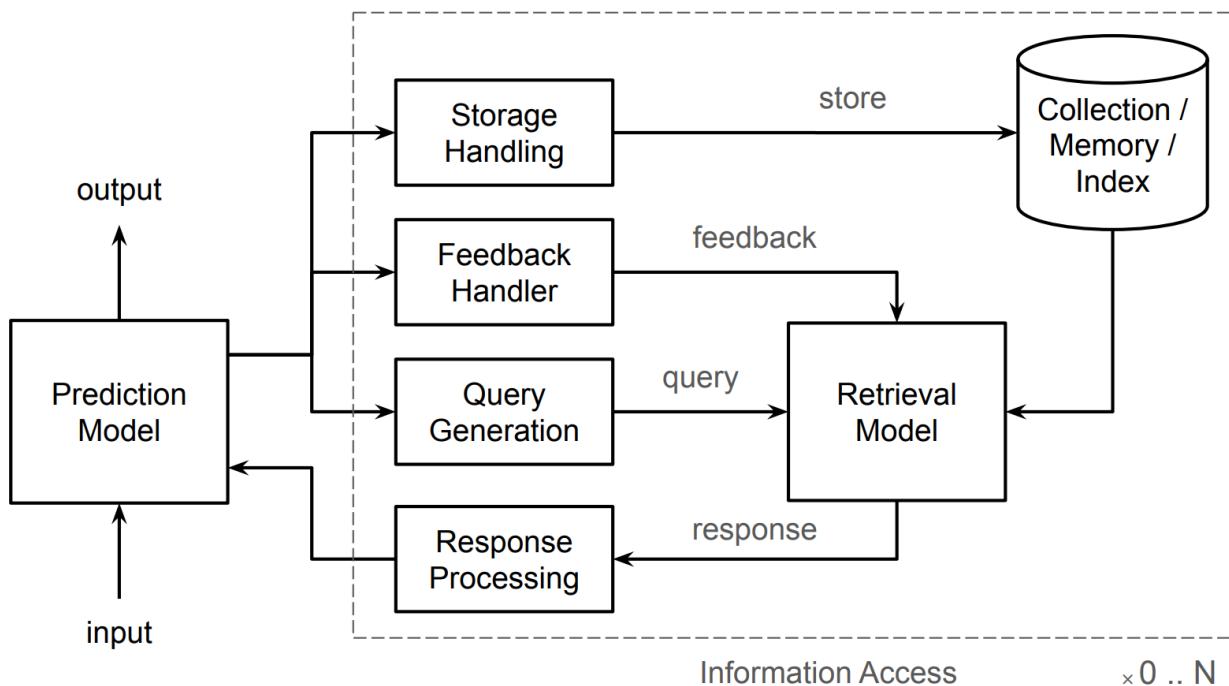
Metric	mean	stdev
Term matching	662 ms	746 ms
SNRM	612 ms	640 ms

How can efficient neural IR  
advance machine learning  
research?

# Retrieval-Enhanced Machine Learning



# A Generic Framework for REML



# Some Case Studies

- Knowledge Grounding
  - REALM [Guu et al., ICML 2020]
  - RetGen [Zhang et al., AAAI 2020]
  - Entities-as-Expert [Févry et al., EMNLP 2020] or Facts-as-Expert [Verga et al., NAACL 2021]
- Memory-Augmented ML
  - Memory network [Sukhbaatar et al., NeurIPS 2015]
  - Neural Turing Machine [Graves et al., arXiv 2014]
- Retrieval-Enhanced Input Representation
  - Pseudo-Relevance Feedback [Attar and Fraenkel, JACM 1977][Croft and Harper, Jdoc 1978]
  - Guided Transformer [Hashemi et al., SIGIR 2020]
  - RETRO [Borgeaud et al., arXiv 2021]

# More Case Studies!

- Generalization through Memorization
  - KNN-LM [Khandelwal et al., ICLR 2020]
  - BERT-KNN [Kassner and Schütze, EMNLP 2020]
  - HybridNCM [Yang et al., CIKM 2019]
- Efficient Access to Longer Context
  - MemViT [Wu et al., arXiv 2022]
  - IDCM [Hofstätter et al., SIGIR 2021]
- Retrieval-Enhanced Optimization
  - Weak supervision for IR [Dehghani et al., SIGIR 2017]
  - Hard negatives in IR optimization, e.g., ANCE [Xiong et al., ICLR 2021]
  - Unsupervised machine translation [Wu et al., NAACL 2019]
  - CLIP [Radford et al., ICML 2021] and VideoCLIP [Xu et al., EMNLP 2021]

Efficient and effective Neural IR can potentially revolutionize machine learning research through retrieval-enhanced models.

---

# Multi-Task Retrieval-Augmented Text Generation with Relevance Sampling

---

Sebastian Hofstätter<sup>1</sup> Jiecao Chen<sup>2</sup> Karthik Raman<sup>2</sup> Hamed Zamani<sup>3</sup>

## Abstract

This paper studies multi-task training of retrieval-augmented generation models for knowledge-intensive tasks. We propose to clean the training set by utilizing a distinct property of knowledge-intensive generation: The connection of query-answer pairs to items in the knowledge base. We filter training examples via a threshold of confidence on the relevance labels, whether a pair is answerable by the knowledge base or not. We train a single Fusion-in-Decoder (FiD) generator on seven combined tasks of the KILT benchmark. The experimental results suggest that our simple yet effective approach substantially improves competitive baselines on two strongly imbalanced tasks; and shows either smaller improvements or no significant regression on the remaining tasks. Furthermore, we demonstrate our multi-task training with relevance label sampling scales well with increased model capacity and achieves state-of-the-art results in five out of seven KILT tasks.

bated when tasks are retroactively expanded (Kwiatkowski et al., 2019), re-purposed (Bajaj et al., 2016) or adapt the collection (Petroni et al., 2021).

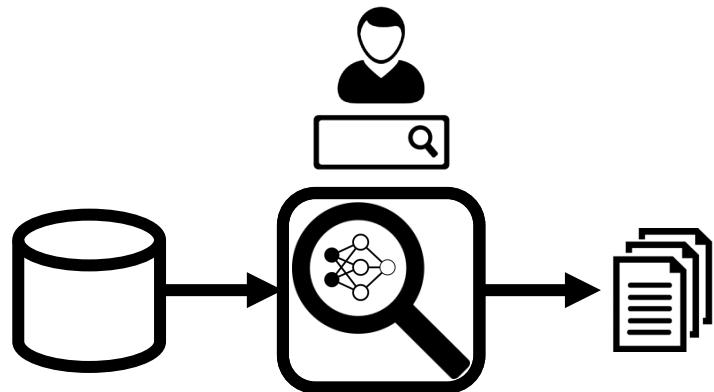
We propose a simple yet effective approach for training retrieval-augmented models for knowledge-intensive tasks with noisy labels. We use a confidence score for query-answer pairs and items in the knowledge base. This confidence can be sourced from manually annotated, heuristic, or model generated aspects. We filter training examples via a threshold of confidence on the relevance labels, whether a pair is answerable by the knowledge base or not. With this we aim to reduce noise in the training process, and produce better results with fewer training examples.

To study our training approach, we use a fixed T5-based dense retrieval module (Ni et al., 2021) and train a Fusion-in-Decoder (FiD) generator (Izacard & Grave, 2020) on multiple tasks of the KILT benchmark (Petroni et al., 2021). KILT aggregates and heuristically maps many different English Wikipedia-based generation tasks to a single Wikipedia snapshot, which introduces considerable noise in the label quality, due to the time-shifted nature of the task creations.

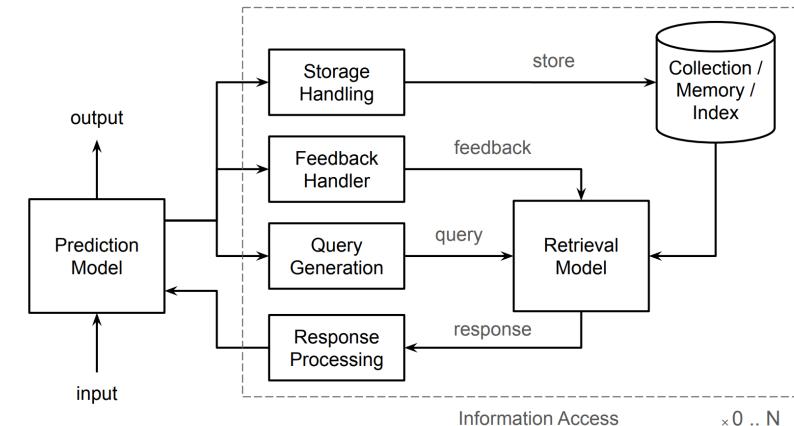
# Results on the KILT Benchmark

Model	Generator	Open Domain QA			Fact Acc.	Slot Filling		Dialog <i>F1</i>
		NQ	HotpotQA	TriviaQA		FEVER	Accuracy	
<b>Top Leaderboard Entries</b>								
1 RAG (Petroni et al., 2021)	BART-Large	44.4	27.0	71.3	86.3	59.2	44.7	13.1
2 DPR + FiD (Piktus et al., 2021)	T5-Base	51.6	38.3	72.7	89.0	82.2	74.0	15.7
3 KGI (Glass et al., 2021)	BART-Large	45.2	–	61.0	85.6	84.4	72.6	18.6
4 Re2G (Anonymous, 2022)	BART-Large	51.7	–	76.3	89.6	<b>87.7</b>	–	18.9
5 Hindsight (Paranjape et al., 2021)	BART-Large	–	–	–	–	–	–	19.2
6 SEAL+FiD (Bevilacqua et al., 2022)	T5-?	53.7	<b>40.5</b>	70.9	89.5	83.7	74.7	18.3
<b>Ours (Alt-200 passages)</b>								
7 GTR + FiD with treatment $\hat{T}$	T5-Base	52.4	30.1	78.9	87.1	83.4	81.5	18.4
8	T5-XL	<b>61.2</b>	39.1	<b>84.6</b>	<b>92.3</b>	85.2	<b>83.7</b>	<b>20.6</b>

# Thank you!



**Efficient Neural Information Retrieval**



**Retrieval-Enhanced Machine Learning**