



# CSCE 689 - Special Topics in NLP for Science

## Lecture 7: Scientific Literature Retrieval


Yu Zhang

[yuzhang@tamu.edu](mailto:yuzhang@tamu.edu)

February 6, 2025

Course Website: <https://yuzhang-teaching.github.io/CSCE689-S25.html>

# Scientific Literature Retrieval



Cardiac injury is common in critical cases of COVID-19.

×

Search

[Advanced](#) [Create alert](#) [Create RSS](#) [User Guide](#)

Save

Email

Send to

Sort by: Best match

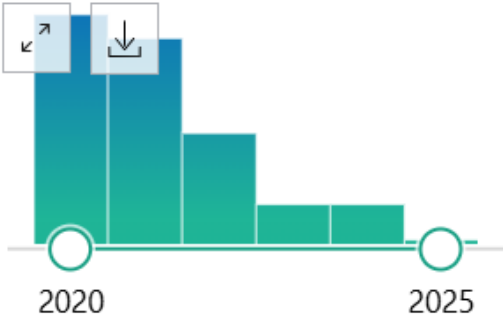
Display options

MY CUSTOM FILTERS

45 results

Page 1 of 5

RESULTS BY YEAR



Year	Results
2020	1
2021	4
2022	2
2023	1
2024	1
2025	2

PUBLICATION DATE

☐ 1

[Cite](#)

[Share](#)

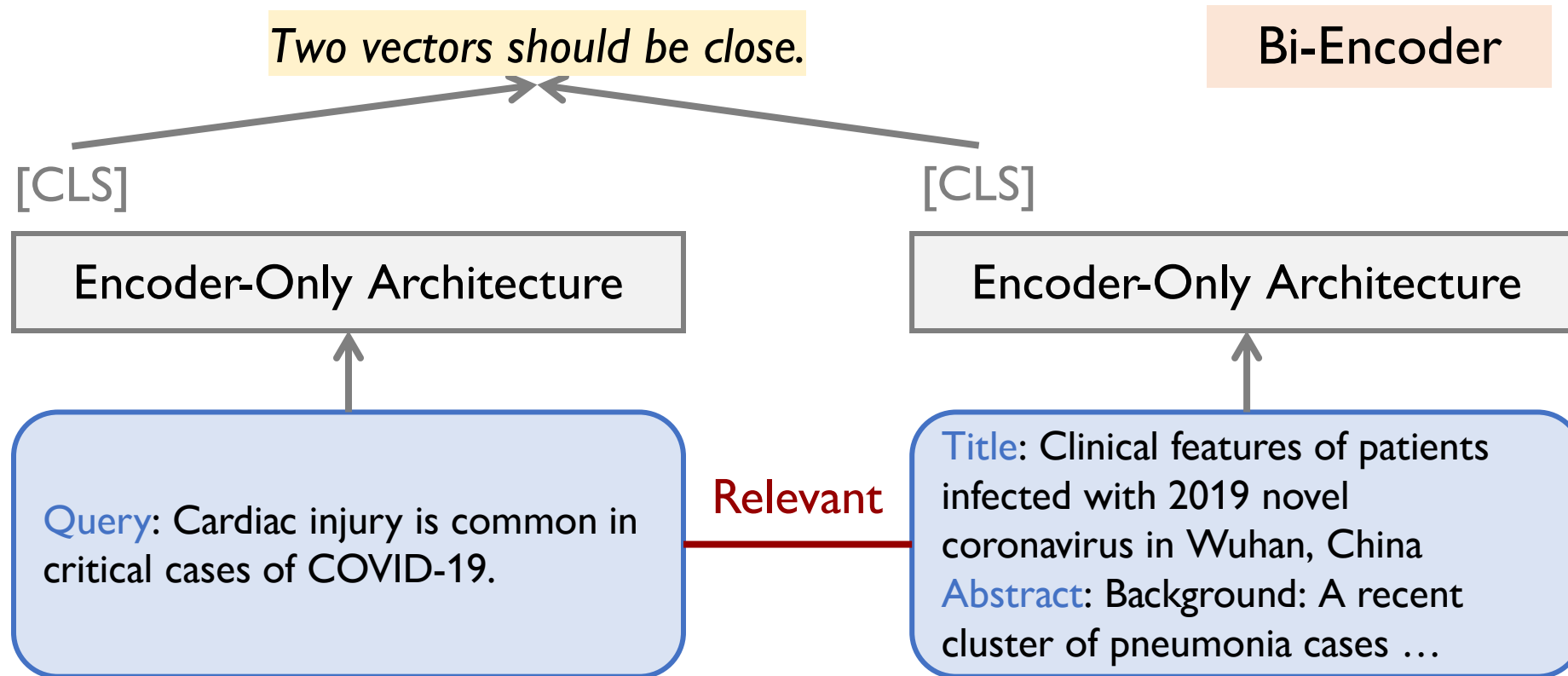
[Clinical features of patients infected with 2019 novel coronavirus in Wuhan, China.](#)

Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, Zhang L, Fan G, Xu J, Gu X, Cheng Z, Yu T, Xia J, Wei Y, Wu W, Xie X, Yin W, Li H, Liu M, Xiao Y, Gao H, Guo L, Xie J, Wang G, Jiang R, Gao Z, Jin Q, Wang J, Cao B. Lancet. 2020 Feb 15;395(10223):497-506. doi: 10.1016/S0140-6736(20)30183-5. Epub 2020 Jan 24. PMID: 31986264 [Free PMC article.](#)

BACKGROUND: A recent cluster of pneumonia **cases** in Wuhan, China, was caused by a novel betacoronavirus, the 2019 novel **coronavirus (2019-nCoV)**. We report the epidemiological, clinical, laboratory, and radiological characteristics and treatment and clin ...

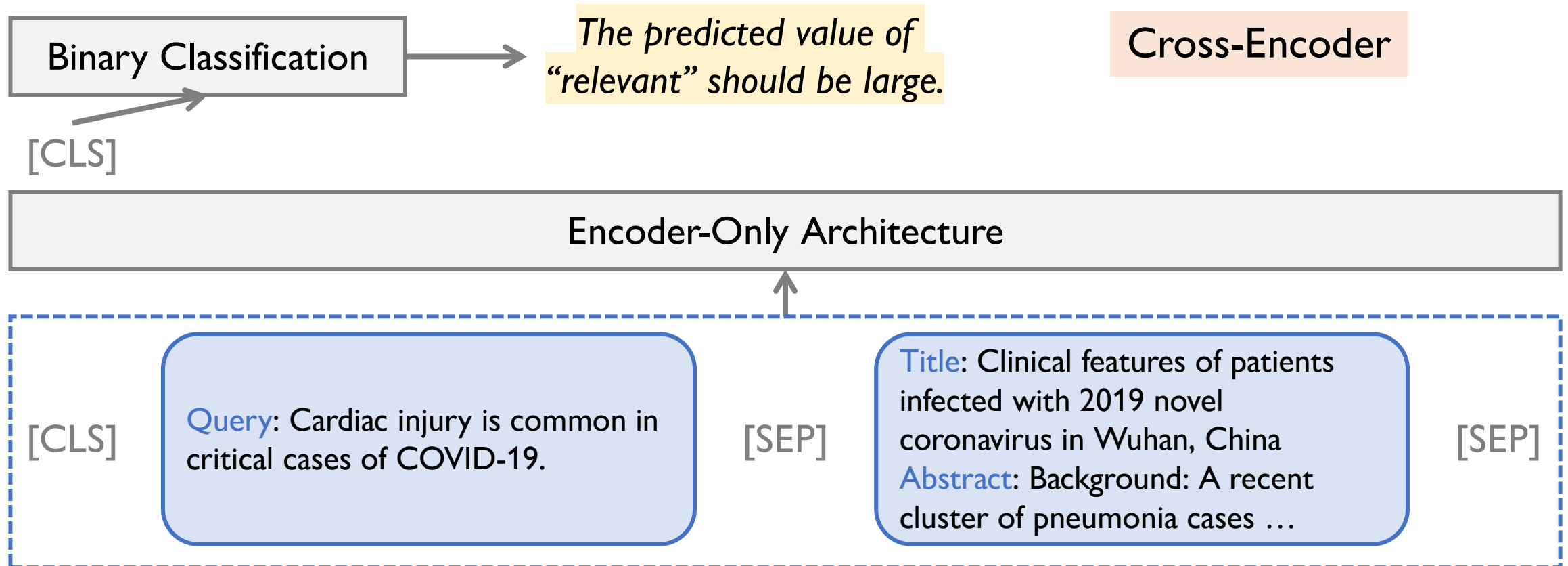
# How to train an LLM to perform scientific literature retrieval?

- **Step 1:** Collect a large number of relevant (query, paper) pairs.
- **Step 2:** Train an LLM with such information (e.g., using contrastive learning).



# How to train an LLM to perform scientific literature retrieval?

- **Step 1:** Collect a large number of relevant (query, paper) pairs.
- **Step 2:** Train an LLM with such information (e.g., using contrastive learning).



# How to train an LLM to perform scientific literature retrieval?

- **Step 1:** Collect a large number of relevant (query, paper) pairs.
- **How?**
  - Unlike citation information that can be crawled from the academic databases or the Web, relevant (query, paper) pairs need to be derived from either **user click-through data** or **human annotations**.
  - User click-through data are **proprietary**.
  - Human annotations **cannot be scaled up**.

# Agenda

- Contrastive Learning with Ground-Truth Search Logs
  - **MedCPT**: Bi-Encoder → Cross-Encoder
- Contrastive Learning with Data from Other Tasks
  - **SciMult**: Mixture-of-Experts Transformer
  - **BMRetriever**: Instruction Tuning
- Application
  - **SciFact**: Scientific Claim Verification

# Agenda

- Contrastive Learning with Ground-Truth Search Logs
  - **MedCPT**: Bi-Encoder → Cross-Encoder
- Contrastive Learning with Data from Other Tasks
  - SciMult: Mixture-of-Experts Transformer
  - BMRetriever: Instruction Tuning
- Application
  - SciFact: Scientific Claim Verification

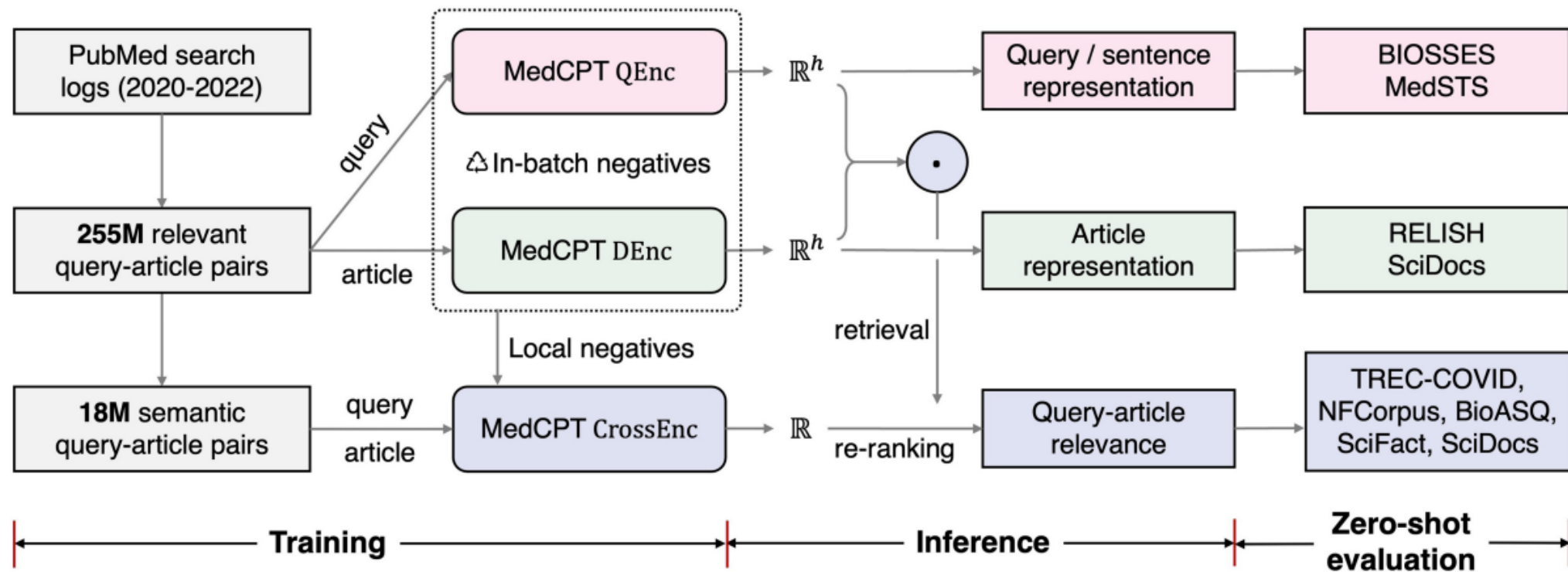
# PubMed Search Logs

- User click-through data from 2020 to 2022
  - A user inputted a query.
  - 20 papers were displayed on the result page.
  - The user clicked paper 1, 6, and 8.
  - Papers relevant to the query: 1, 6, 8
  - Papers irrelevant to the query: 2, 3, 4, 5, 7
  - Papers cannot be judged as relevant/irrelevant: 9, 10, ..., 20
- 255M relevant (query, paper) pairs
  - Most of such queries are short keywords, and matching them to the clicked articles is a relatively simple task.
- 18M **semantically** relevant (query, paper) pairs
  - Remove queries either having only one word or all of the clicked articles containing exact mentions of the whole input query

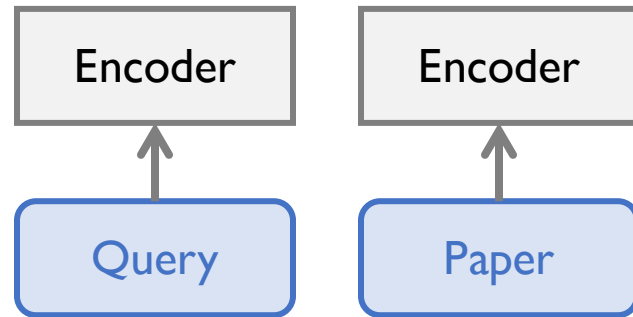


# The MedCPT Framework

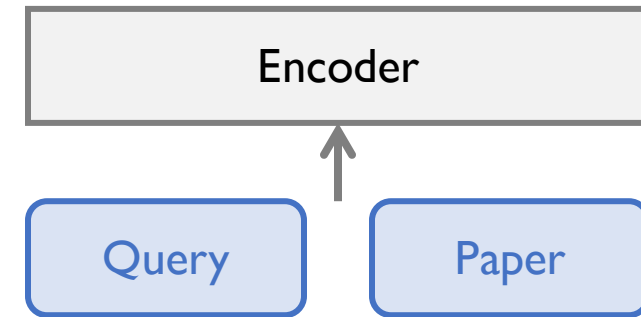
- Bi-Encoder for retrieval (from a large candidate pool)
- Cross-Encoder for re-ranking (the retrieved papers)



## Recap: Bi-Encoder vs. Cross-Encoder



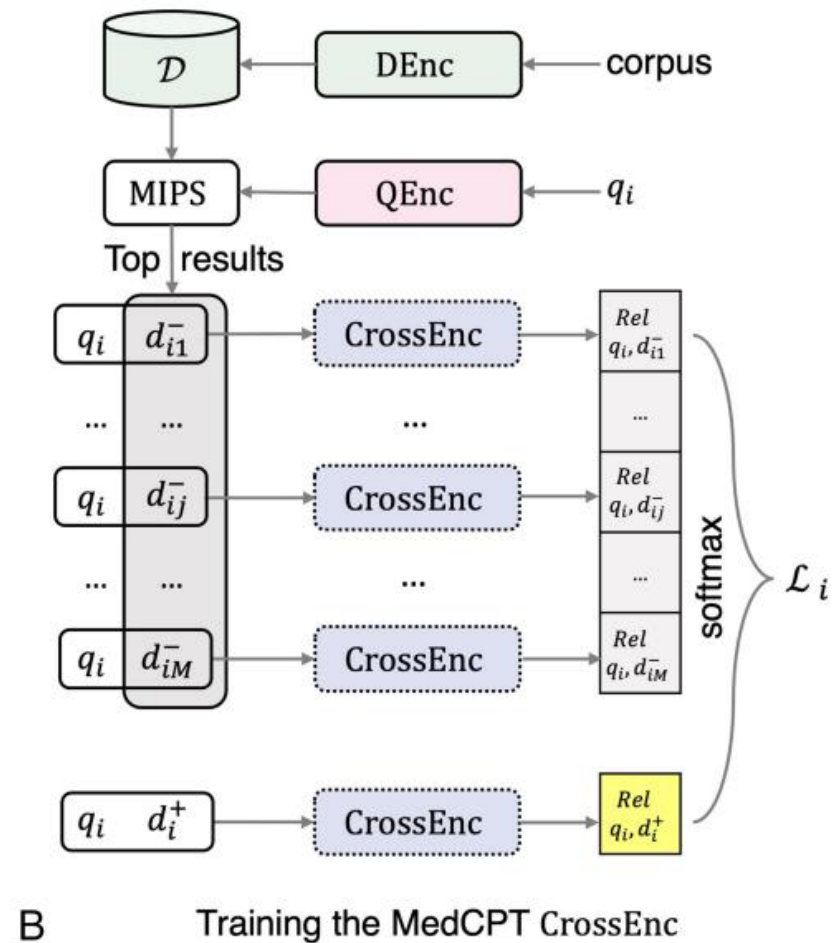
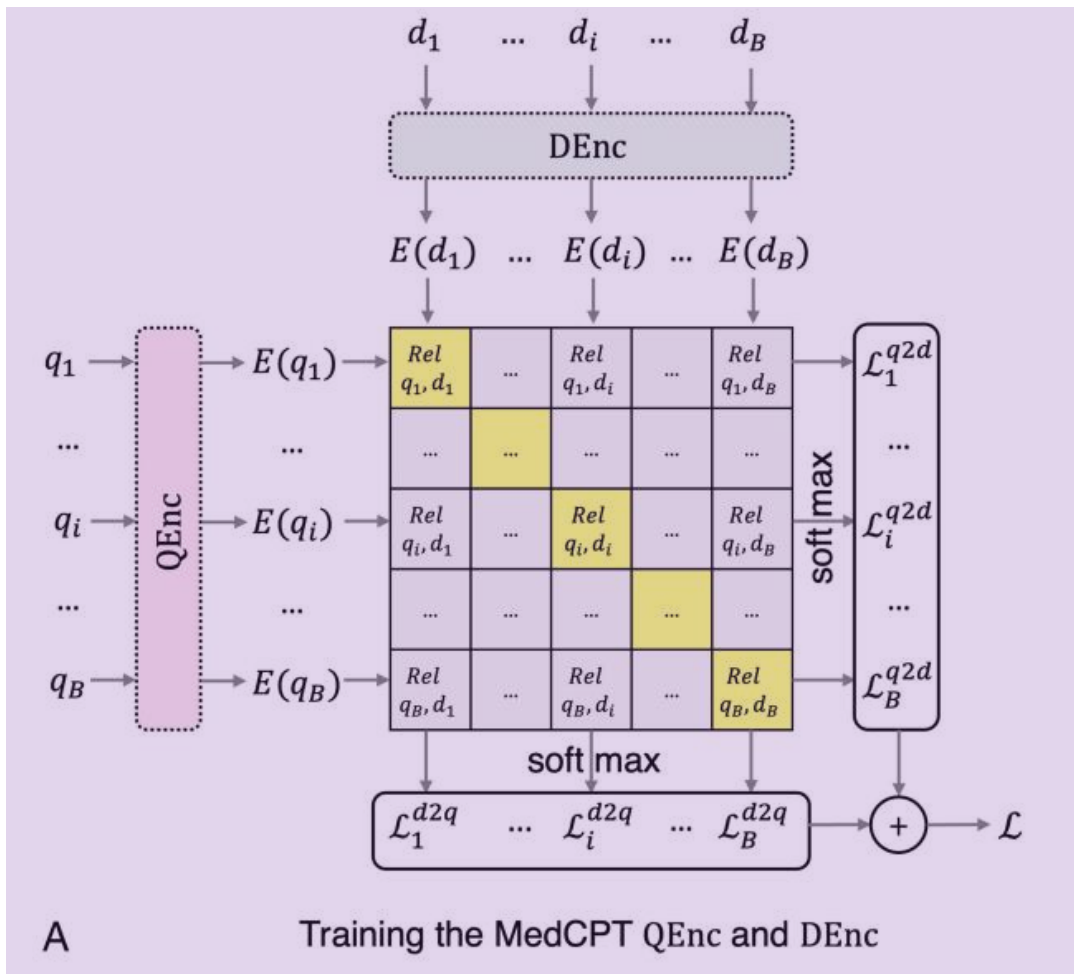
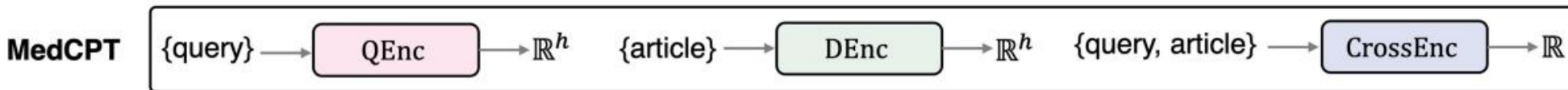
Bi-Encoder



Cross-Encoder

- Bi-Encoder is **much more efficient during the inference time**.
- If we use Cross-Encoder, the query and the paper can **serve as context of each other**, so that the model can learn a better contextualized representation of each token in the input sequence.
- MedCPT: Using Bi-Encoder to remove most (e.g., 99%) of the candidates, and using Cross-Encoder to more carefully rank the remaining candidates (e.g., 1%).

# Bi-Encoder Contrastive Learning

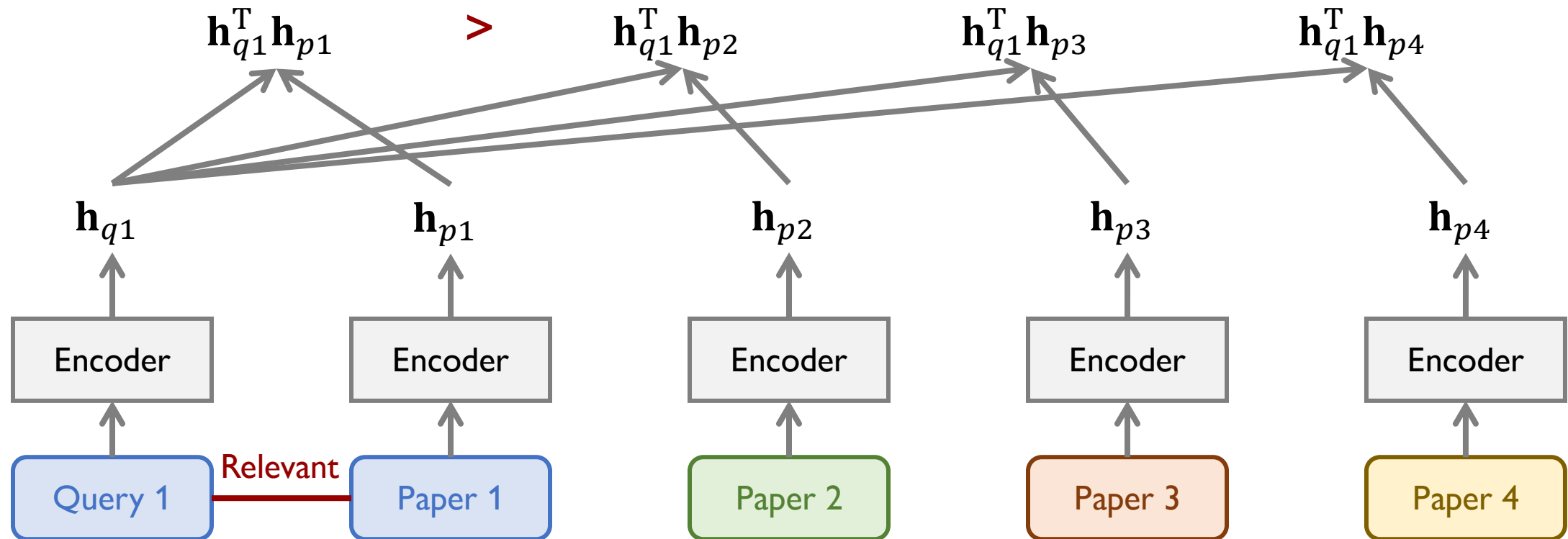


# In-Batch Negative Sampling

Objective Function: maximize

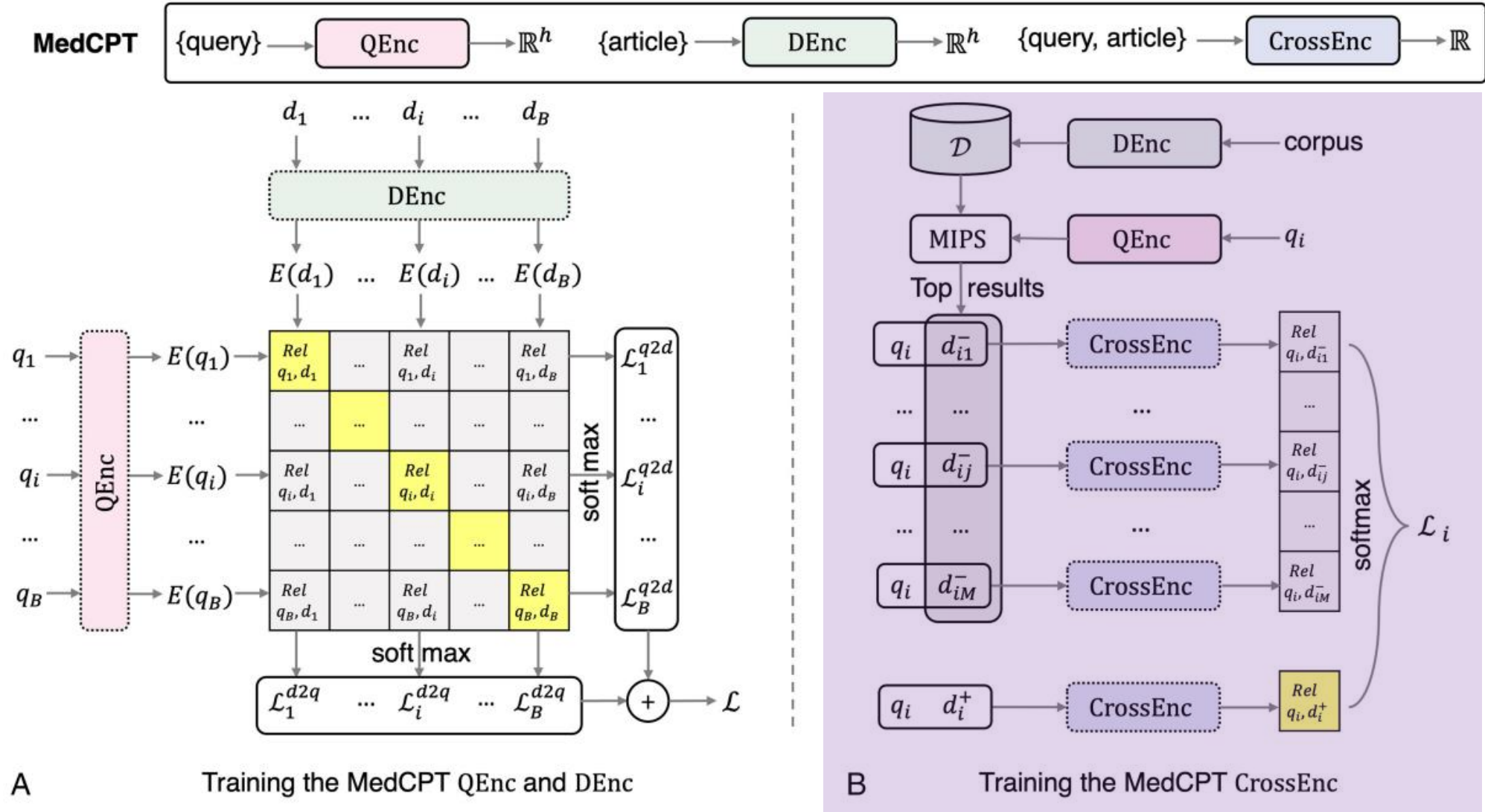
$$\frac{\exp(\mathbf{h}_{q1}^T \mathbf{h}_{p1})}{\exp(\mathbf{h}_{q1}^T \mathbf{h}_{p1}) + \exp(\mathbf{h}_{q1}^T \mathbf{h}_{p2}) + \exp(\mathbf{h}_{q1}^T \mathbf{h}_{p3}) + \exp(\mathbf{h}_{q1}^T \mathbf{h}_{p4})}$$

- Paper 2 is relevant to Query 2, but its relevance to Query 1 is unknown.



# Cross-Encoder Contrastive Learning


Using trained Bi-Encoder to derive hard negatives








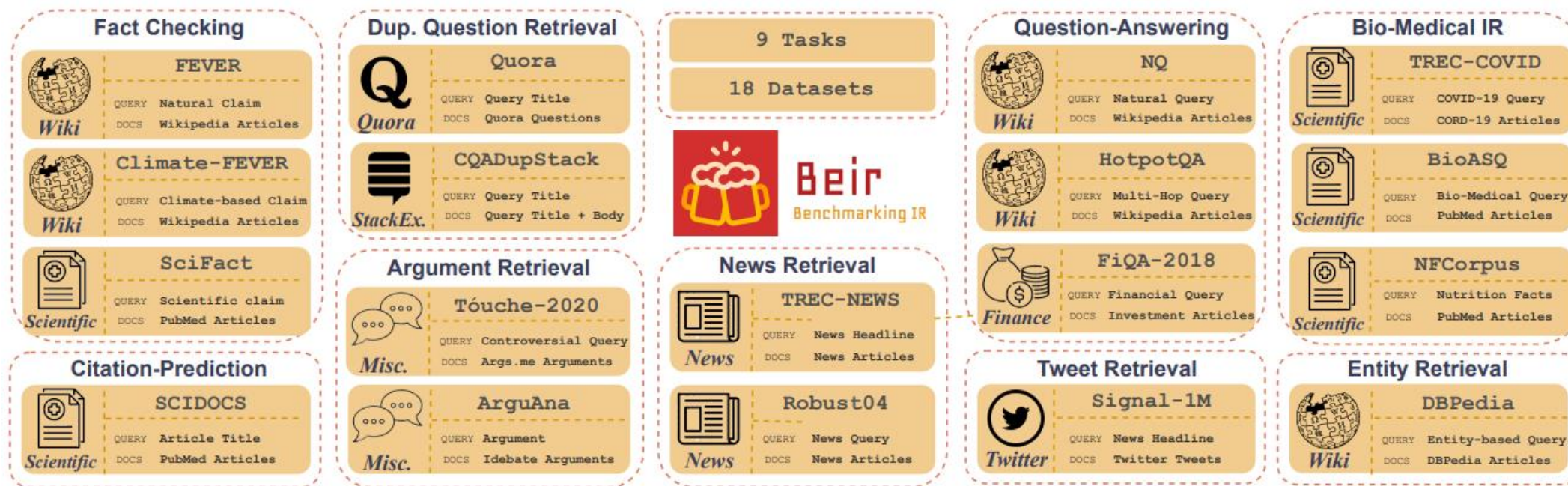
# Evaluation of MedCPT

- The BEIR benchmark
  - <https://github.com/beir-cellar/beir>

 **beir** Public

A Heterogeneous Benchmark for Information Retrieval. Easy to use, evaluate your models across 15+ diverse IR datasets.

 Python  1.7k  199



# Performance of MedCPT: Query-Paper Relevance

Method	Size	COVID	NFC	BioASQ	SciFact	SciDocs	Avg.
Sparse retrievers							
BM25		0.656	0.325	0.465	0.665	0.158	0.454
BM25 + MiniLM	66M	<b>0.757</b>	0.350	<u>0.523</u>	0.688	0.166	<u>0.497</u>
DeepCT	110M	0.406	0.283	<u>0.407</u>	0.630	0.124	<u>0.370</u>
SPARTA	110M	0.538	0.301	0.351	0.582	0.126	0.380
docT5query	220M	0.713	0.328	0.431	0.675	0.162	0.462
Dense retrievers							
DPR	110M	0.332	0.189	0.127	0.318	0.077	0.209
ANCE	110M	0.654	0.237	0.306	0.507	0.122	0.365
TAS-B	66M	0.481	0.319	0.383	0.643	0.149	0.395
GenQ	220M	0.619	0.319	0.398	0.644	0.143	0.425
Contriever	110M	0.596	0.328		0.677	0.165	
Contriever + MiniLM	176M	<u>0.701</u>	0.344		0.692	<u>0.171</u>	
ColBERT	110M	0.677	0.305	<u>0.474</u>	0.671	0.145	0.454
Large language model retrievers							
Google GTR-Base	110M	0.539	0.308	0.271	0.600	0.149	0.373
Google GTR-Large	335M	0.557	0.329	0.320	0.639	0.158	0.401
Google GTR-XL	1.24B	0.584	0.343	0.317	0.635	0.159	0.408
Google GTR-XXL	4.80B	0.501	0.342	0.324	0.662	0.161	0.398
OpenAI cpt-text-S	300M	0.679	0.332		0.672		
OpenAI cpt-text-M	1.20B	0.585	<u>0.367</u>		0.704		
OpenAI cpt-text-L	<u>6.00B</u>	0.562	<u>0.380</u>		<u>0.744</u>		
OpenAI cpt-text-XL	<b>175B</b>	0.649	<b>0.407</b>		<u>0.754</u>		
MedCPT							
MedCPT	330M	<u>0.709</u>	0.355	<b>0.553</b>	<b>0.761</b>	<b>0.172</b>	<b>0.510</b>
MedCPT (retriever only)	220M	0.697	0.340	0.332	0.724	0.123	0.443
MedCPT w/o contrastive pre-training (PubMedBERT)	110M	0.059	0.015		0.010	0.004	

# Performance of MedCPT: Paper-Paper and Query-Query Relevance

**Table 2.** Evaluation results of the MedCPT article encoder on the RELISH dataset.<sup>a</sup>

Method	MAP			NDCG			Avg.
	@5	@10	@15	@5	@10	@15	
Random	79.33	77.22	75.41	80.70	77.67	76.40	77.79
Sparse retrievers							
BM25	88.91	86.72	84.54	89.48	87.39	86.21	87.21
PMRA	90.30	87.57	85.75	90.95	88.40	87.45	88.40
Non-BERT embedding-based models							
fastText	85.75	82.81	81.79	86.79	83.79	83.12	84.01
BioWordVec	89.84	86.51	84.67	89.90	86.67	85.53	87.19
InferSent	85.21	82.16	80.41	86.56	83.31	82.35	83.33
WikiSentVec	87.92	85.23	83.40	88.65	85.74	84.81	85.96
BioSentVec	90.76	88.10	86.16	90.05	87.76	86.89	88.29
LDA	85.44	82.66	80.36	86.51	82.91	81.31	83.20
Doc2Vec	86.23	84.74	83.39	86.55	84.70	84.09	84.95
BERT-based models							
BioBERT	88.14	85.81	83.90	88.97	86.29	85.10	86.37
PubMedBERT	83.69	81.07	79.53	85.47	82.39	81.41	82.26
SPECTER	92.27	90.00	88.36	91.47	89.12	88.42	89.94
SciNCL	94.72	92.74	91.14	93.67	91.91	90.94	92.52
MedCPT DEnc	95.58	93.99	92.39	94.78	93.12	92.43	93.72

Evaluating relevance between papers

**Table 3.** Evaluation results (Pearson's correlation coefficients) of the MedCPT QEnc on the BIOSSES and MedSTS datasets.<sup>a</sup>

Model	BIOSSES	MedSTS
Non-BERT embedding-based models		
BioWordVec	0.694	0.747
USE	0.345	0.714
BioSentVec (PubMed)	0.817	0.750
BioSentVec (MIMIC-III)	0.350	0.759
BioSentVec (PubMed + MIMIC-III)	0.795	0.767
BERT-based models		
PubMedBERT	0.528	0.521
Clinical BERT	0.556	0.525
SPECTER	0.694	0.702
SciNCL	0.847	0.706
MedCPT QEnc	0.893	0.765

Evaluating relevance between short sentences



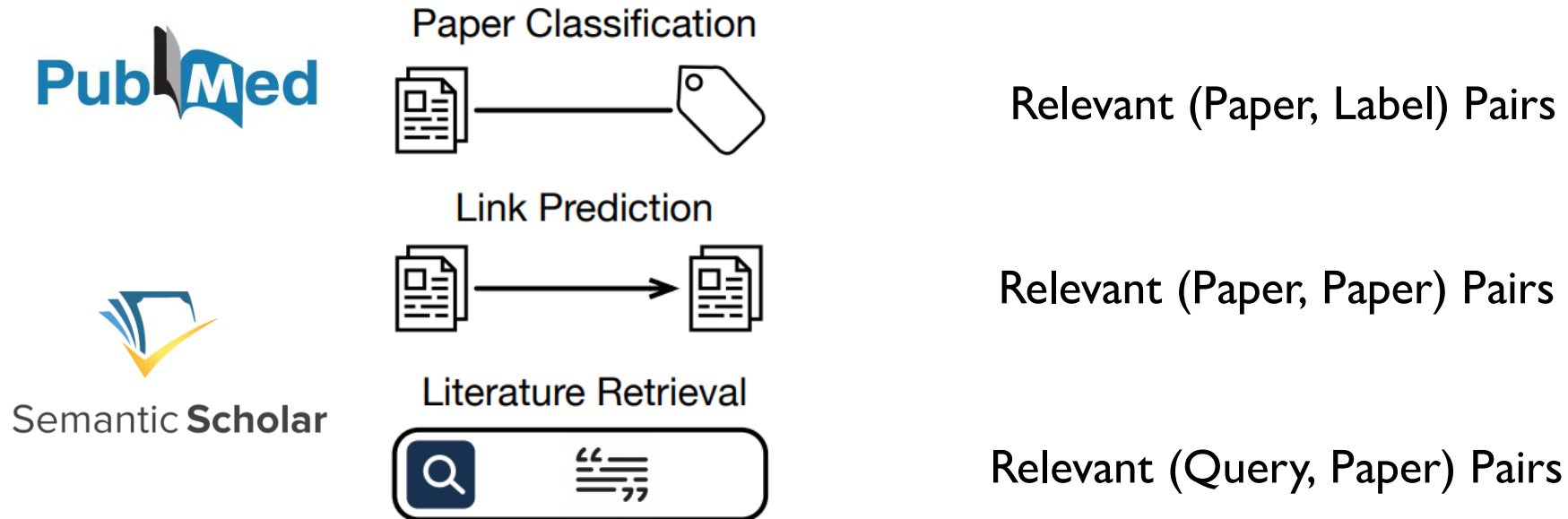
# Take-Away Messages

- Contrastive learning with user click-through data makes **small, domain-specific LMs** outperform **large, general LMs**.
- The **retrieval→re-ranking** framework uses Bi-Encoder to filter out most (e.g., 99%) of the candidates, and using Cross-Encoder to more carefully rank the remaining candidates (e.g., 1%).
  - “Get the best of both worlds” by utilizing the advantages of Bi-Encoder and Cross-Encoder
- Limitation:
  - Strong reliance on **proprietary data**
    - Most researchers do not have access to search logs.

# Agenda

- Contrastive Learning with Ground-Truth Search Logs
  - MedCPT: Bi-Encoder → Cross-Encoder
- Contrastive Learning with Data from Other Tasks
  - **SciMult**: Mixture-of-Experts Transformer
  - BMRetriever: Instruction Tuning
- Application
  - SciFact: Scientific Claim Verification

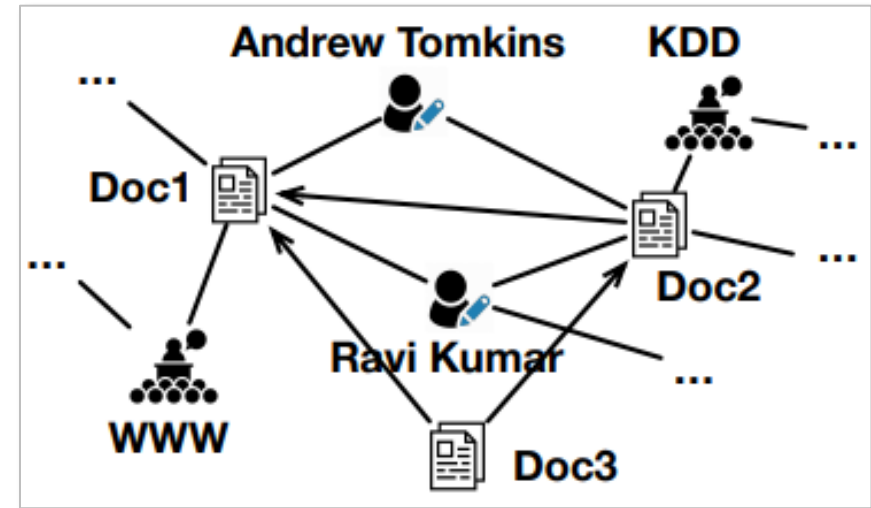
# Harvesting Other Types of Relevant (Text, Text) Pairs



- Combine all these pairs together for contrastive learning?
- **Task Interference**: The model is confused by different types of “relevance”.

# An Illustrative Example of Task Interference

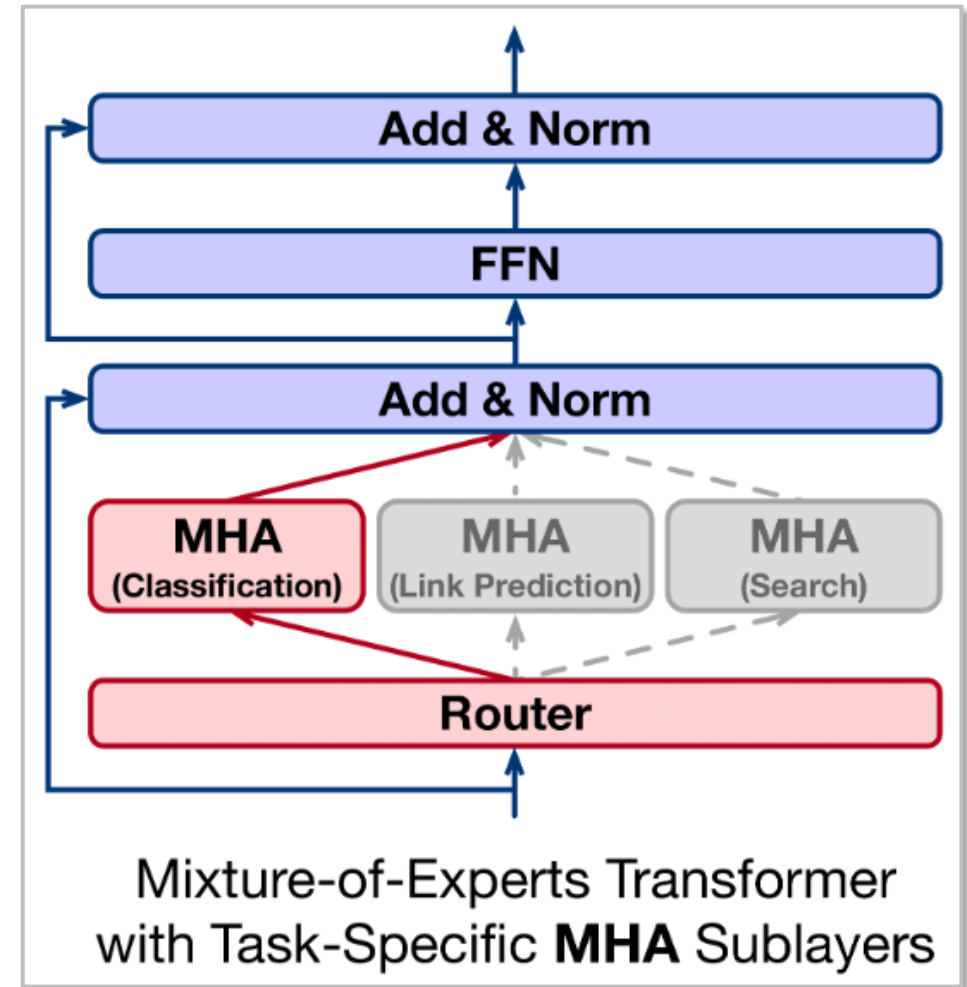
- Recall the link prediction problem
- Imagine that predicting each type of “links” is a “task”
  - **Citation Prediction:** Paper→Paper
  - **Same Author Prediction:** Paper-Author-Paper
  - Each type of “links” defines one type of “relevance”.
- Directly merging the relevant (paper, paper) pairs induced by different link types?
  - **The model will be confused!**



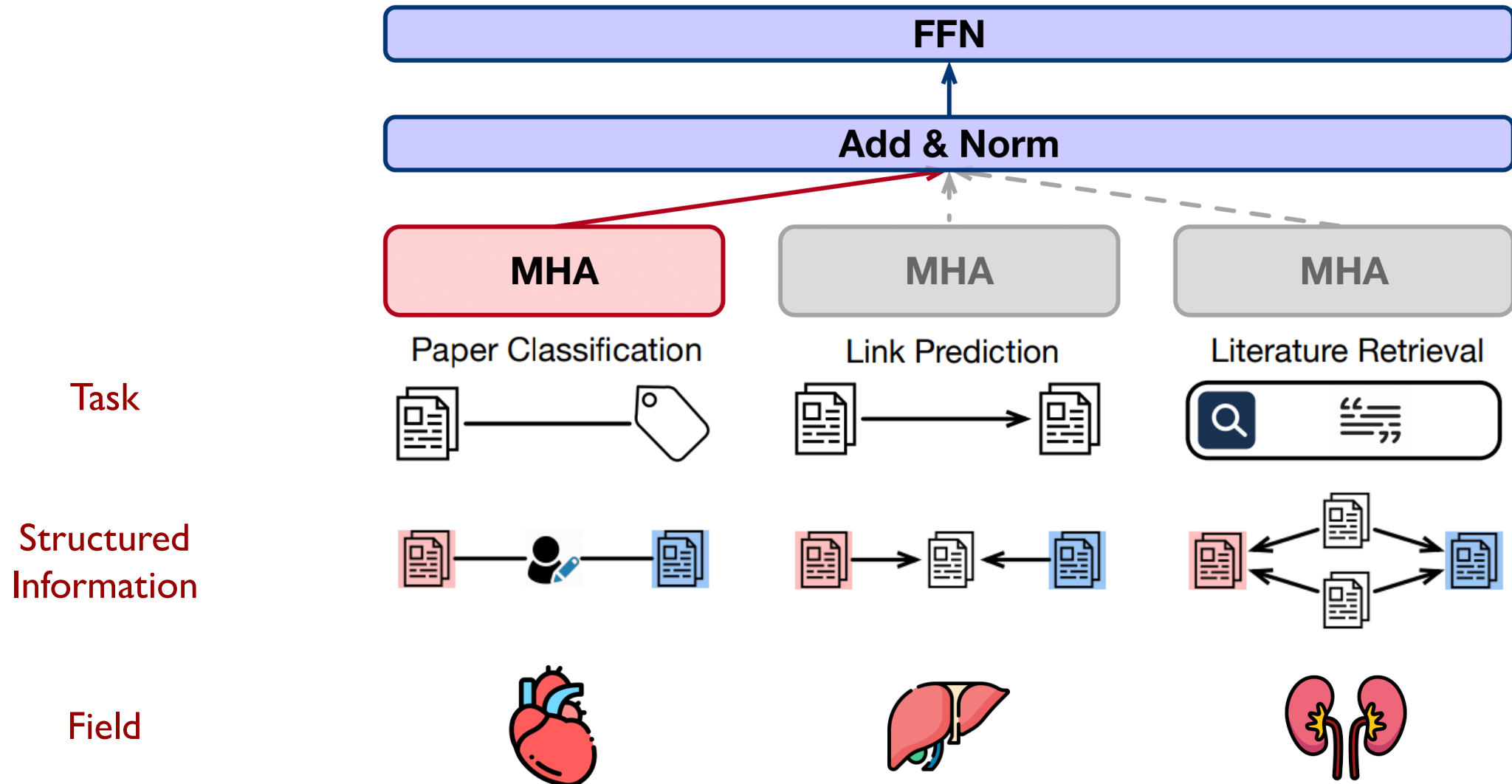
(Doc3, Doc2) are **relevant** according to **Paper→Paper** but **irrelevant** according to **Paper-Author-Paper**.

# Tackling Task Interference: Mixture-of-Experts Transformer

- A typical Transformer layer
  - 1 Multi-Head Attention (MHA) sublayer
  - 1 Feed Forward Network (FFN) sublayer
- A Mixture-of-Experts (MoE) Transformer layer
  - **Multiple** MHA sublayers
  - 1 FFN sublayer
  - (Or 1 MHA & Multiple FFN)
- Specializing some parts of the architecture to be an “expert” of one task
- The model can learn both **commonalities** and **characteristics** of different tasks.

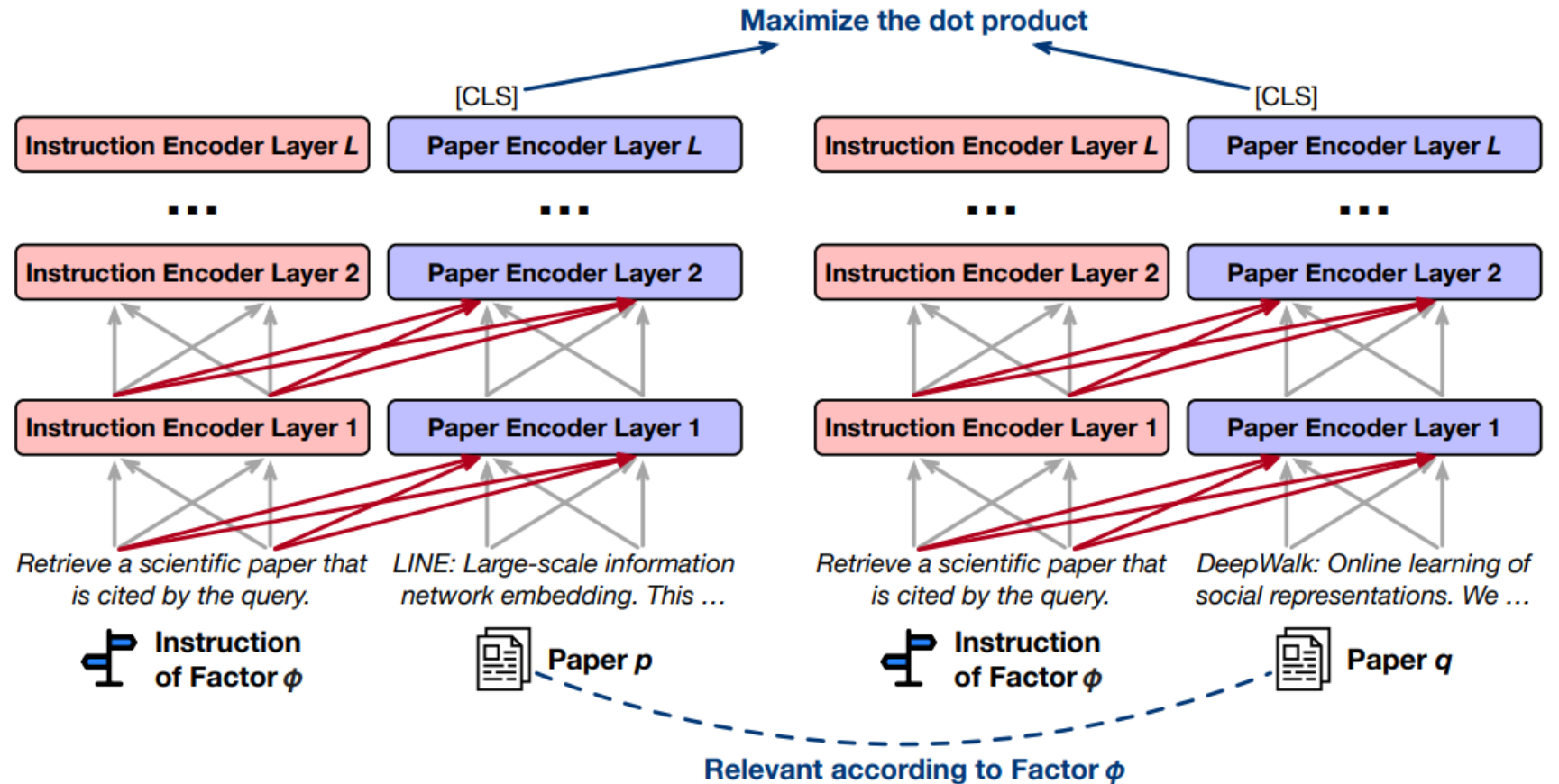


# Tackling Task Interference: Mixture-of-Experts Transformer



# Tackling Task Interference: Instruction Tuning

- Using a **factor-specific instruction** to guide the paper encoding process
- The instruction serves as the context of the paper.
- The paper does NOT serve as the context of the instruction.



# Evaluations of SciMult

Task	Pre-training	In-domain Evaluation	Cross-domain Evaluation
Classification	<b>MAPLE</b> (Zhang et al., 2023b) {CS-Journal, Biology-MeSH, Medicine-MeSH}	<b>MAPLE</b> (Zhang et al., 2023b) {CS-Conference, Chemistry-MeSH}, <b>SciDocs</b> (Cohan et al., 2020) {MAG Fields, MeSH Diseases}	<b>MAPLE</b> (Zhang et al., 2023b) {Geography, Psychology}
Link Prediction	<b>Citation Prediction Triplets</b> (Cohan et al., 2020)	<b>SciDocs</b> (Cohan et al., 2020) {Co-view, Co-read, Cite, Co-cite}	<b>Recommendation</b> (Kanakia et al., 2019), <b>PMC-Patients</b> (Zhao et al., 2022)
Search	<b>SciRepEval-Search</b> (Singh et al., 2022)	<b>SciRepEval-Search</b> (Singh et al., 2022)	<b>TREC-COVID</b> (Voorhees et al., 2021), <b>SciFact</b> (Wadden et al., 2020), <b>NFCorpus</b> (Boteva et al., 2016)

- For Search, both the retrieval and the re-ranking settings are evaluated.

Search		
SciRepEval-Search (Singh et al., 2022)	2,637	reranking, 10.00 for each query on average
TREC-COVID in SciRepEval (Voorhees et al., 2021)	50	reranking, 1386.36 for each query on average
TREC-COVID in BEIR (Voorhees et al., 2021)	50	171,332
SciFact (Wadden et al., 2020)	1,109	5,183
NFCorpus (Boteva et al., 2016)	3,237	3,633



# Performance of SciMult: Search

Search	SciRepEval (Singh et al., 2022)		BEIR (Thakur et al., 2021)			Average
	Search (2022)	TREC-COVID (2021)	TREC-COVID (2021)	SciFact (2020)	NFCorpus (2016)	
	nDCG@10	nDCG@10	nDCG@10	nDCG@10	nDCG@10	
BM25	73.47	55.86	57.79	65.63	30.00	56.55
SciBERT	71.39	40.98	4.17	0.88	1.90	23.86
SentBERT	71.84	51.30	20.73	9.40	6.69	31.99
SPECTER	73.42	66.45	29.91	49.74	15.83	47.07
PubMedBERT	70.77	45.28	7.56	0.30	1.09	25.00
LinkBERT	71.66	52.45	2.28	0.49	1.77	25.73
BioLinkBERT	71.18	36.01	3.17	0.12	0.98	22.29
OAG-BERT	72.17	55.09	7.11	18.33	8.48	32.24
SciNCL	73.78	73.50	34.69	56.51	22.34	52.16
SPECTER 2.0	<b>78.22<sup>†</sup></b>	79.43	58.48	67.16	22.84	61.23
SciMult-Vanilla	76.44	<b>86.76</b>	67.22	<b>70.76</b>	<b>31.20</b>	66.48
SciMult-MHAExpert	76.33	86.29	<b>71.18</b>	70.67	30.79	<b>67.05</b>
SciMult-FFNExpert	76.02	82.32	52.15	63.57	27.48	60.31
SciMult-Prefix	76.55	82.83	68.15	70.70	30.02	65.65
SciMult-Instruction	75.86	83.59	61.05	70.62	30.25	64.27

# Performance of SciMult: Classification

Fine-grained classification	MAPLE (Zhang et al., 2023b)												Average
	CS-Conference			Chemistry-MeSH			Geography			Psychology			
	R@20	R@50	R@100	R@20	R@50	R@100	R@20	R@50	R@100	R@20	R@50	R@100	
SciBERT (Beltagy et al., 2019)	42.01	42.84	43.87	30.53	31.46	32.15	52.04	54.53	58.11	43.07	44.02	45.22	43.32
SentBERT (Reimers and Gurevych, 2019)	42.79	44.34	45.96	30.75	31.73	32.44	53.54	57.23	61.11	43.33	44.60	46.37	44.52
SPECTER (Cohan et al., 2020)	47.38	53.18	58.43	34.26	39.35	43.41	59.12	65.33	70.75	47.07	51.30	56.17	52.15
PubMedBERT (Gu et al., 2021)	41.93	42.56	43.24	30.46	31.46	31.83	52.19	54.82	56.88	43.93	46.28	49.27	43.74
LinkBERT (Yasunaga et al., 2022)	42.15	43.16	44.22	30.52	31.56	32.37	50.58	50.94	51.63	42.62	42.90	43.23	42.16
BioLinkBERT (Yasunaga et al., 2022)	42.00	42.81	43.57	30.37	31.15	31.48	50.36	50.54	50.86	42.39	42.55	42.79	41.74
OAG-BERT (Liu et al., 2022)	42.59	43.79	44.93	30.58	31.97	32.62	51.44	52.25	53.16	42.63	42.95	43.30	42.68
SciNCL (Ostendorff et al., 2022)	47.92	53.57	58.29	34.99	40.50	44.64	59.00	65.49	71.41	48.74	54.21	59.84	53.22
SPECTER 2.0 (Singh et al., 2022)	48.63	55.09	60.68	36.17	43.06	48.26	62.87	70.30	76.37	50.60	58.27	65.66	56.33
SciMult-Vanilla	53.40	64.70	74.09	39.78	51.31	59.75	62.08	70.65	77.79	50.42	56.58	63.17	60.31
SciMult-MHAExpert	54.02	65.49	75.07	39.41	50.92	59.59	65.94	75.01	81.93	51.77	59.55	67.86	62.21
SciMult-FFNExpert	53.73	63.79	72.46	38.01	48.76	57.43	61.90	70.69	78.81	50.09	56.94	64.28	59.74
SciMult-Prefix	53.68	63.62	72.07	37.97	48.95	57.56	62.86	71.65	79.71	50.10	57.25	64.53	60.00
SciMult-Instruction	53.78	63.99	72.72	38.81	50.12	58.96	63.26	71.74	79.52	50.86	58.47	66.46	60.72

# Performance of SciMult: Link Prediction

Link Prediction (Reranking)	SciDocs (Cohan et al., 2020)								Kanakia et al. (2019)			Average
	Co-view		Co-read		Cite		Co-cite		Recommendation			
	MAP	nDCG	MAP	nDCG	MAP	nDCG	MAP	nDCG	nDCG@5	nDCG@10	nDCG	
Citeomatic (Bhagavatula et al., 2018)	81.1 <sup>†</sup>	90.2 <sup>†</sup>	80.5 <sup>†</sup>	90.2 <sup>†</sup>	86.3 <sup>†</sup>	94.1 <sup>†</sup>	84.4 <sup>†</sup>	92.8 <sup>†</sup>	–	–	–	–
Kanakia et al. (2019)	–	–	–	–	–	–	–	–	83.88 <sup>‡</sup>	87.71 <sup>‡</sup>	93.59 <sup>‡</sup>	–
SciBERT (Beltagy et al., 2019)	50.7 <sup>†</sup>	73.1 <sup>†</sup>	47.7 <sup>†</sup>	71.1 <sup>†</sup>	48.3 <sup>†</sup>	71.7 <sup>†</sup>	49.7 <sup>†</sup>	72.6 <sup>†</sup>	77.17	82.49	90.86	66.86
SentBERT (Reimers and Gurevych, 2019)	68.2 <sup>†</sup>	83.3 <sup>†</sup>	64.8 <sup>†</sup>	81.3 <sup>†</sup>	63.5 <sup>†</sup>	81.6 <sup>†</sup>	66.4 <sup>†</sup>	82.8 <sup>†</sup>	76.75	81.49	90.80	76.45
SPECTER (Cohan et al., 2020)	83.6 <sup>†</sup>	91.5 <sup>†</sup>	84.5 <sup>†</sup>	92.4 <sup>†</sup>	88.3 <sup>†</sup>	94.9 <sup>†</sup>	88.1 <sup>†</sup>	94.8 <sup>†</sup>	83.38	87.39	93.64	89.32
PubMedBERT (Gu et al., 2021)	59.43	78.23	55.59	75.63	51.81	73.43	58.19	77.80	77.30	82.21	91.09	70.97
LinkBERT (Yasunaga et al., 2022)	44.21	67.76	41.04	65.31	39.33	63.91	42.84	67.18	76.10	80.89	90.47	61.73
BioLinkBERT (Yasunaga et al., 2022)	56.46	76.38	50.76	72.18	47.73	70.55	52.94	74.44	77.02	81.78	90.73	68.27
OAG-BERT (Liu et al., 2022)	64.61	81.50	60.13	78.65	57.35	77.60	62.47	80.92	76.73	82.12	90.96	73.91
SciNCL (Ostendorff et al., 2022)	<b>85.3<sup>†</sup></b>	<b>92.3<sup>†</sup></b>	<b>87.5<sup>†</sup></b>	<b>93.9<sup>†</sup></b>	93.6 <sup>†</sup>	97.3 <sup>†</sup>	<b>91.6<sup>†</sup></b>	96.4 <sup>†</sup>	85.33	88.38	94.34	<b>91.45</b>
SPECTER 2.0 (Singh et al., 2022)	85.18 <sup>†</sup>	<b>92.27<sup>†</sup></b>	86.95 <sup>†</sup>	93.53 <sup>†</sup>	92.23 <sup>†</sup>	96.84 <sup>†</sup>	91.13 <sup>†</sup>	96.28 <sup>†</sup>	86.03	89.12	94.59	91.29
SciMult-Vanilla	83.99	91.68	86.66	93.67	91.37	96.26	91.50	<b>96.45</b>	<b>87.32</b>	89.32	<b>94.88</b>	91.19
SciMult-MHAExpert	83.92	91.60	86.45	93.55	92.58	96.92	91.47	96.36	86.68	<b>89.45</b>	94.77	91.25
SciMult-FFNExpert	83.23	91.26	85.61	93.20	93.77	97.42	90.39	95.94	85.75	88.45	94.29	90.85
SciMult-Prefix	83.43	91.48	85.89	93.27	<b>94.28</b>	<b>97.60</b>	90.73	96.09	86.05	88.85	94.66	91.12
SciMult-Instruction	82.13	90.88	84.14	92.36	92.63	96.91	89.27	95.43	86.49	88.81	94.51	90.32

# PMC-Patients Leaderboard

- Given a patient summary, find the most relevant papers.

<https://pmc-patients.github.io>

Patient-to-Article Retrieval (PAR) Leaderboard					
	Model	MRR (%)	P@10 (%)	nDCG@10 (%)	R@1k (%)
1 June 25, 2023	DPR (SciMult-MHAExpert) <i>UIUC/Microsoft</i> (Zhang et al. 2023)	29.89	9.35	13.79	53.71
2 Apr 5, 2023	RRF <i>Tsinghua University</i> (Zhao et al. 2023)	29.86	8.86	13.36	49.45

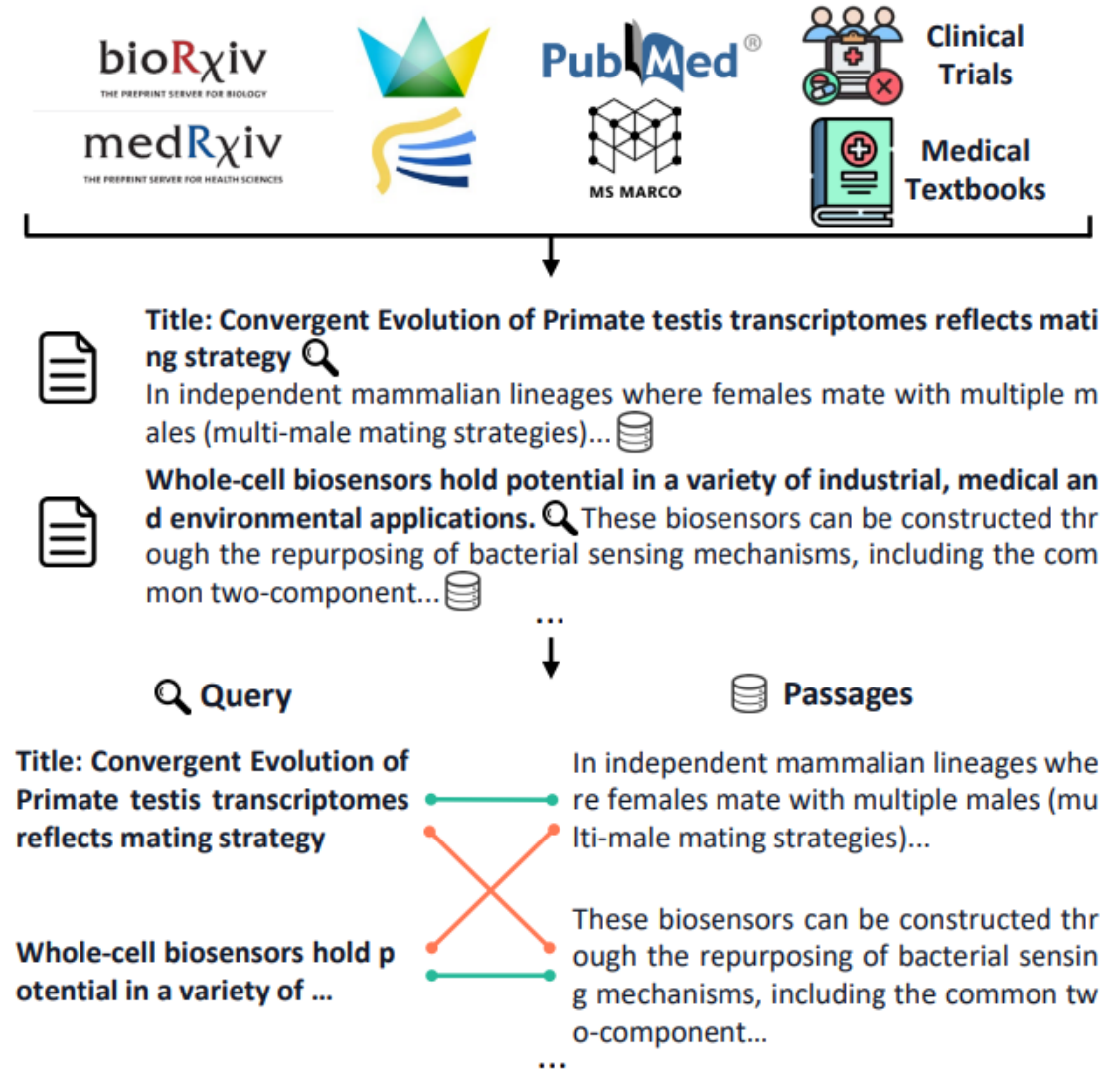
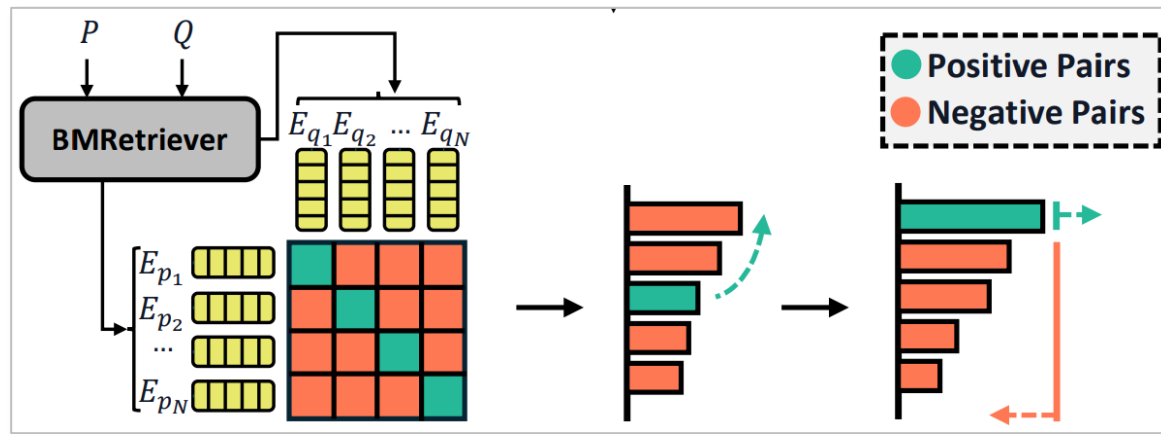
# Agenda

- Contrastive Learning with Ground-Truth Search Logs
  - MedCPT: Bi-Encoder → Cross-Encoder
- Contrastive Learning with Data from Other Tasks
  - SciMult: Mixture-of-Experts Transformer
  - **BMRetriever**: Instruction Tuning
- Application
  - SciFact: Scientific Claim Verification

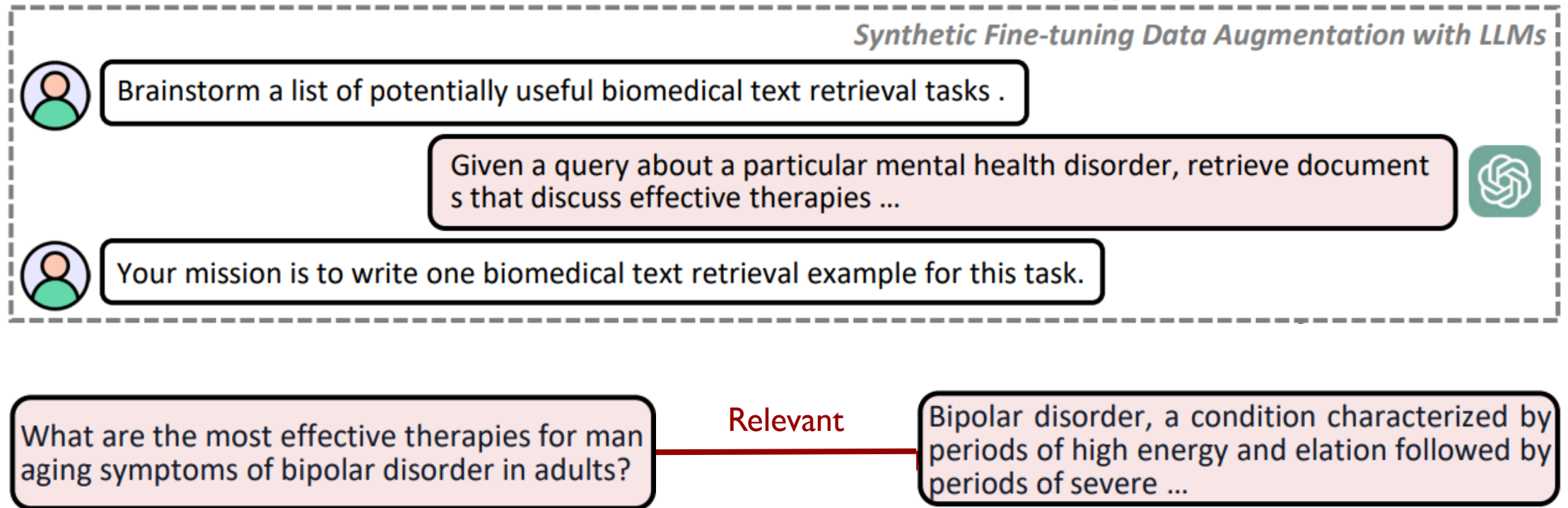


# Getting Relevant (Text, Text) Pairs from One Paper

- For corpora with titles, treat the **title** as the query and the **corresponding abstract** as the passage.
- For untitled corpora, randomly sample **two disjoint passages** from documents, using one as the query and the other as the passage.

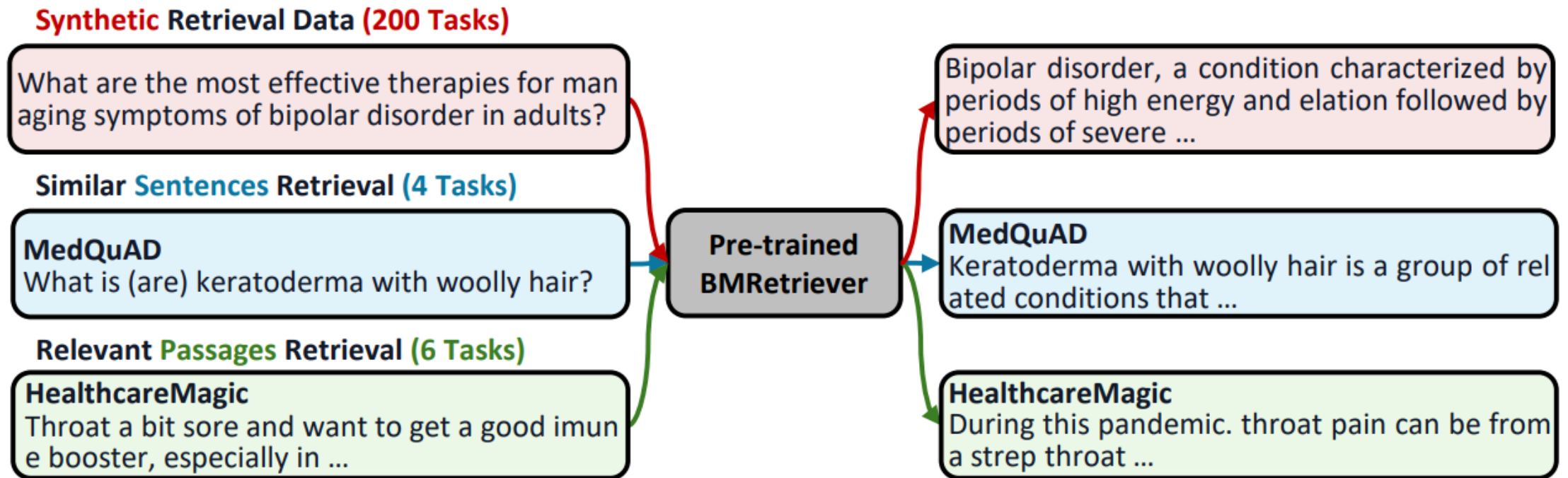


# Getting Relevant (Text, Text) Pairs from LLMs



# Getting Relevant (Text, Text) Pairs from Other Tasks

- Instruction Tuning





# Generalize to Unseen Retrieval Tasks

## *Inference: Generalization to Various Tasks*

### Text Retrieval (4 Tasks)

#### SciFact

Microstructural development of human newborn cerebral white ...



#### SciFact

Alterations of the architecture of cerebral white matter in the ...

### Sentence Similarity (1 Task)

#### BIOSSES

It has recently been shown that Craf is essential for Kras G12D-induced ...



#### BIOSSES

It has recently become evident that Craf is essential for the onset of ...

### Question Answering (3 Tasks)

#### PubMedQA

Are group 2 innate lymphoid cells ( ILC2s ) increased in chronic rhinosin...



#### PubMedQA

Chronic rhinosinusitis (CRS) is a heterogeneous disease with an uncertain ...

### Entity Linking (2 Tasks)

#### DrugBank

Cytarabine



#### DrugBank

Chronic rhinosinusitis (CRS) is a heterogeneous disease with an uncertain ...

### Paper Recommendation (1 Task)

#### SciRepEval

ERK1 and ERK2 are related protein-serine/threonine kinases that ...



#### SciRepEval

ERK1 and MEK2 regulate distinct functions by sorting ERK2 to different ...

# Performance of BMRetriever: Paper Retrieval

Task	Scale	# PT Pairs	# FT Pairs	Standard IR				Sent. Sim.	Avg. Retr.	Avg. All
Model				NFCorpus	SciFact	SciDocs	Trec-COVID	BIOSSES		
Sparse Retrieval										
BM25 (Robertson et al., 2009)	—	—	—	0.325	0.665	0.158	0.656	—	0.451	—
Base Size (< 1B)										
Contriever (Izacard et al., 2022)	110M	1B	500K	0.328	0.677	0.165	0.596	0.833	0.442	0.520
Dragon (Lin et al., 2023)	110M	—	28.5M	0.339	0.679	0.159	0.759	0.819	0.484	0.551
SPECTER 2.0 (Singh et al., 2023)	110M	3.3M	—	0.228	0.671	—	0.584	—	—	—
SciMult (Zhang et al., 2023)	110M	5.5M	—	0.308	0.707	—	0.712	—	—	—
COCO-DR (Yu et al., 2022)	110M	15M	500K	0.355	0.709	0.160	0.789	0.829	0.503	0.567
SGPT-125M (Muennighoff, 2022)	125M	unknown	500K	0.228	0.569	0.122	0.703	0.752	0.406	0.475
MedCPT (Jin et al., 2023)	220M	—	255M	0.340	0.724	0.123	0.697	0.837	0.471	0.544
GTR-L (Ni et al., 2022)	335M	2B	662K	0.329	0.639	0.158	0.557	0.849	0.421	0.506
InstructOR-L (Su et al., 2023)	335M	—	1.24M	0.341	0.643	0.186	0.581	0.844	0.438	0.519
E5-Large-v2 <sup>†</sup> (Wang et al., 2022b)	335M	270M	1M	0.371	0.726	0.201	0.665	0.836	0.491	0.560
BGE-Large* <sup>‡</sup> (Chen et al., 2024)	335M	1.2B	1.62M	0.345	0.723	0.222	0.753	0.804	<b>0.511</b>	<u>0.569</u>
BMRETRIEVER-410M	410M	10M	1.4M	0.321	0.711	0.167	0.831	0.840	<u>0.508</u>	<b>0.574</b>
Large Size (1B - 5B)										
InstructOR-XL (Su et al., 2023)	1.5B	—	1.24M	0.360	0.646	0.174	0.713	0.842	0.473	0.547
GTR-XL (Ni et al., 2022)	1.2B	2B	662K	0.343	0.635	0.159	0.584	0.789	0.430	0.502
GTR-XXL (Ni et al., 2022)	4.8B	2B	662K	0.342	0.662	0.161	0.501	0.819	0.417	0.497
SGPT-1.3B (Muennighoff, 2022)	1.3B	unknown	500K	0.320	0.682	0.162	0.730	0.830	0.473	0.545
SGPT-2.7B (Muennighoff, 2022)	2.7B	unknown	500K	0.339	0.701	0.166	0.752	0.848	0.489	0.561
BMRETRIEVER-1B	1B	10M	1.4M	0.344	0.760	0.180	0.840	0.858	<u>0.531</u>	<u>0.596</u>
BMRETRIEVER-2B	2B	10M	1.4M	0.351	0.760	0.199	0.863	0.828	<b>0.543</b>	<b>0.600</b>
XL Size (> 5B)										
SGPT-5.8B (Muennighoff, 2022)	5.8B	unknown	500K	0.362	0.747	0.199	0.849	0.863	0.539	0.604
LLaRA (Li et al., 2023a)	7B	21M	500K	0.372	0.757	0.172	0.853	—	0.539	—
RepLLaMA (Ma et al., 2023)	7B	—	500K	0.378	0.756	0.181	0.847	—	0.541	—
LLM2Vec* (BehnamGhader et al., 2024)	7B	1.2M	1.5M	0.393	0.788	0.225	0.776	0.852	0.545	0.606
E5-Mistral* (Wang et al., 2024)	7B	—	1.8M	0.386	0.764	0.162	0.872	0.855	<u>0.546</u>	<u>0.608</u>
CPT-text-XL (Neelakantan et al., 2022)	175B	unknown	unknown	0.407	0.754	—	0.649	—	—	—
BMRETRIEVER-7B	7B	10M	1.4M	0.364	0.778	0.201	0.861	0.847	<b>0.551</b>	<b>0.610</b>

# Performance of BMRetriever: QA, Entity Linking & Recommendation

Task	Question Answering									Entity Linking						Paper Rec.	
Model	BioASQ			PubMedQA			iCliniq			DrugBank			MeSH			RELISH	
	R@5	R@20	nDCG@20	R@5	R@20	nDCG@20	R@5	R@20	nDCG@20	R@1	R@5	MRR@5	R@1	R@5	MRR@5	MAP	nDCG
<b>Base Size (&lt; 1B)</b>																	
Dragon (2023)	36.2	<b>54.6</b>	49.1	<u>71.8</u>	74.0	72.0	50.6	65.2	47.4	81.0	87.6	<u>83.3</u>	28.2	47.0	34.8	72.6	80.6
MedCPT (2023)	34.7	<u>54.4</u>	45.2	66.3	71.1	60.4	26.8	42.0	24.9	75.1	<u>88.0</u>	80.6	27.7	<u>54.2</u>	37.4	83.6	89.7
E5-Large-v2 <sup>†</sup> (2022b)	<u>36.8</u>	54.0	<u>50.4</u>	71.6	<u>74.2</u>	<u>72.2</u>	<u>57.6</u>	<u>72.0</u>	<u>55.8</u>	<b>81.8</b>	86.5	81.5	<b>32.8</b>	<b>55.0</b>	<b>41.3</b>	<u>84.9</u>	<u>91.0</u>
BMRETRIEVER-410M	<b>39.9</b>	54.2	<b>53.1</b>	<b>73.8</b>	<b>74.6</b>	<b>72.4</b>	<b>60.6</b>	<b>72.8</b>	<b>56.6</b>	<u>81.4</u>	<b>88.2</b>	<b>83.7</b>	<u>31.5</u>	53.8	<u>39.8</u>	<b>85.2</b>	<b>91.2</b>
<b>Large Size (1B - 5B)</b>																	
InstructOR-XL (2023)	29.9	43.2	41.8	70.5	74.0	69.1	<u>64.9</u>	<u>78.1</u>	<u>58.3</u>	75.3	84.2	80.3	33.6	56.2	45.7	84.5	90.6
SGPT-2.7B (2022)	33.9	47.4	47.3	68.3	73.7	63.2	45.0	52.2	41.2	71.9	77.0	62.9	20.2	39.7	28.5	84.9	90.8
BMRETRIEVER-1B	<u>40.4</u>	<u>55.8</u>	<u>53.4</u>	<u>73.6</u>	<u>74.4</u>	<u>72.7</u>	61.1	73.7	56.8	<b>84.7</b>	<u>89.1</u>	<b>86.5</b>	<u>35.5</u>	<u>60.3</u>	<u>48.8</u>	<u>85.2</u>	<u>91.3</u>
BMRETRIEVER-2B	<b>42.5</b>	<b>56.5</b>	<b>55.7</b>	<b>74.0</b>	<b>74.6</b>	<b>73.1</b>	<b>70.0</b>	<b>81.2</b>	<b>65.7</b>	<u>82.6</u>	<b>90.2</b>	<u>85.8</u>	<b>45.6</b>	<b>71.3</b>	<b>59.5</b>	<b>85.4</b>	<b>91.5</b>
<b>XL Size (&gt; 5B)</b>																	
E5-Mistral* (2024)	39.6	55.4	52.7	72.6	74.2	70.0	56.7	72.2	51.8	78.5	92.2	84.0	47.9	76.2	<b>61.3</b>	85.2	90.8
BMRETRIEVER-7B	<b>43.7</b>	<b>60.2</b>	<b>57.4</b>	<b>74.2</b>	<b>74.6</b>	<b>73.8</b>	<b>68.4</b>	<b>79.7</b>	<b>63.7</b>	<b>84.7</b>	<b>92.8</b>	<b>88.0</b>	<b>49.8</b>	<b>76.5</b>	61.1	<b>86.7</b>	<b>92.2</b>

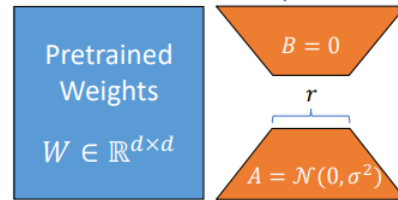
# Take-Away Messages

- If you cannot access proprietary search logs but still need to train a retrieval model, get relevant (text, text) pairs from:
  - Other tasks (e.g., classification, citation prediction, question answering)
  - Different paragraphs in one document
  - LLMs
- Directly merging all these data together for contrastive learning suffers from task interference. Solutions include:
  - Mixture-of-Experts Transformers
  - Instruction Tuning

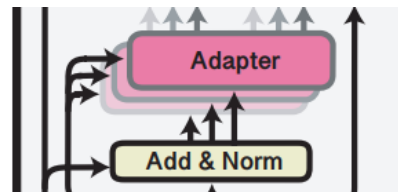
# A Summary of Advanced Techniques Introduced in Recent Lectures

Parameter-Efficient Fine-Tuning (PEFT):  
Only tune a small number of parameters in LLMs

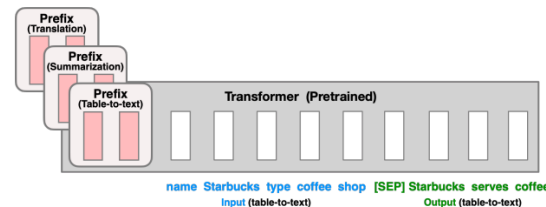
## LoRA [ICLR 2022]



## Adapter [EACL 2021]

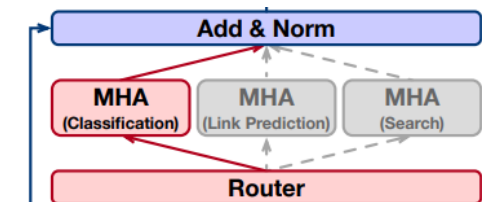


## Prefix Tuning [ACL 2021]



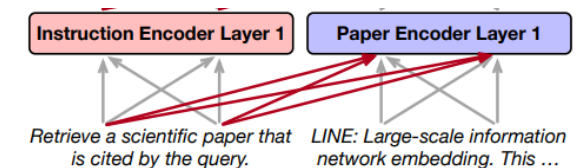
What if we tune the entire model?

## Mixture-of-Experts [ICML 2022]



What if we tune the entire model?

## Instruction Tuning [ICLR 2022]

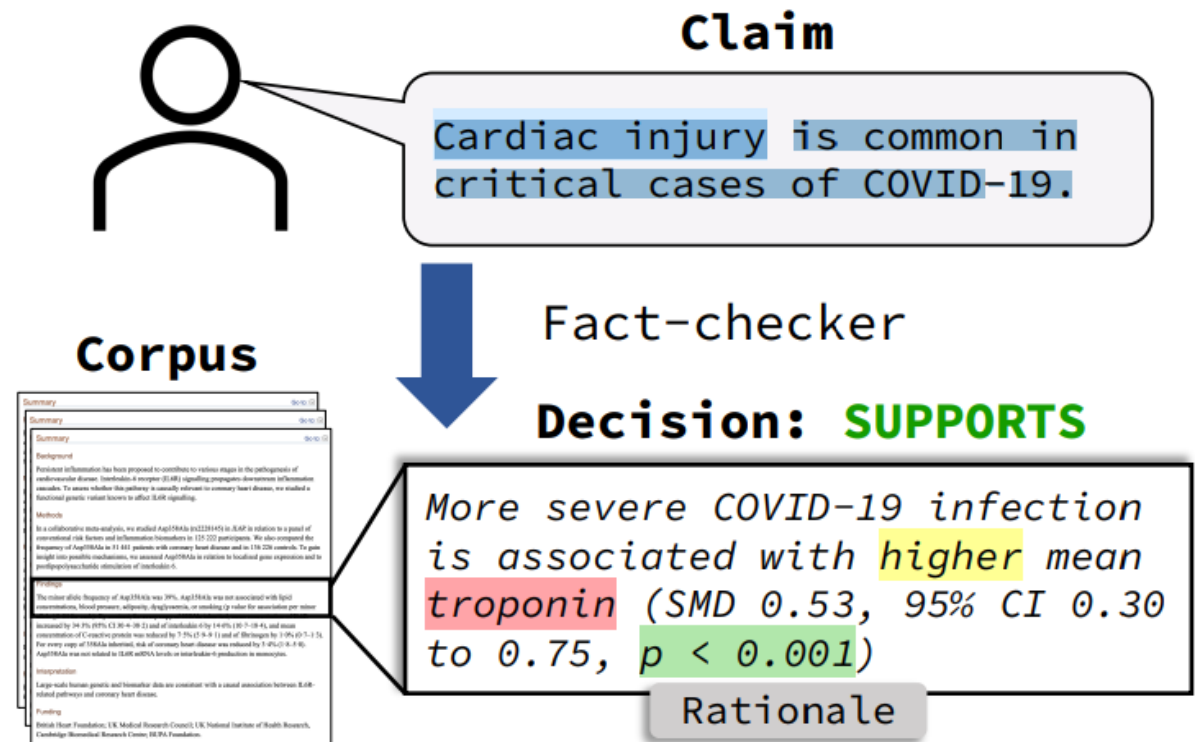


# Agenda

- Contrastive Learning with Ground-Truth Search Logs
  - MedCPT: Bi-Encoder → Cross-Encoder
- Contrastive Learning with Data from Other Tasks
  - SciMult: Mixture-of-Experts Transformer
  - BMRetriever: Instruction Tuning
- Application
  - **SciFact**: Scientific Claim Verification

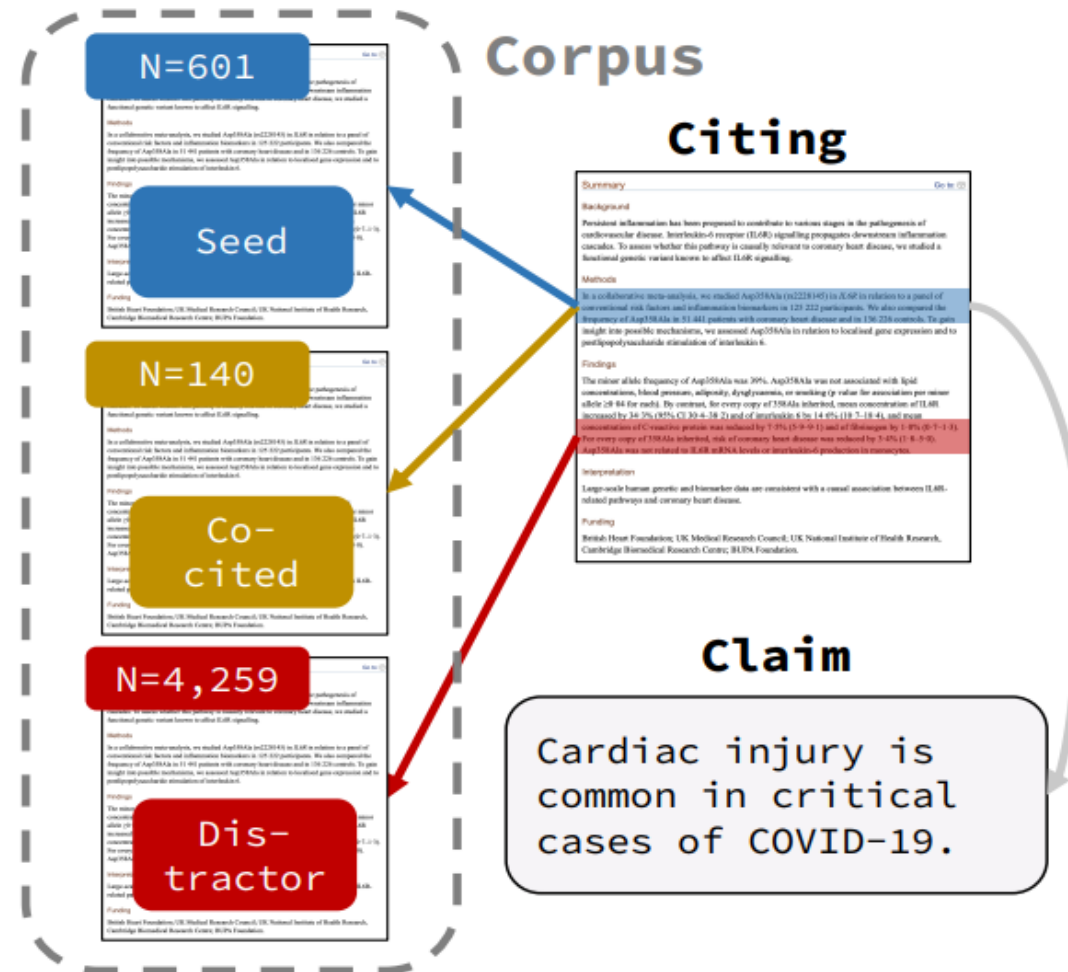
# Scientific Claim Verification

- Given a scientific claim:
  - Step 1 (Relevant Paper Retrieval):** Find all papers relevant to this claim.
  - Step 2 (Rationale Sentence Selection):** In each relevant paper, find relevant sentences.
  - Step 3 (Stance Prediction):** Based on the relevant sentences, predict if the paper supports, refutes, or is neutral towards the claim.





# Dataset Construction



- **Data source:** S2ORC
- **Annotators** write claims based on citation sentences.

## Source citance

"Future studies are also warranted to evaluate the potential association between WNT5A/PCP signaling in adipose tissue and atherosclerotic CVD, given the major role that IL-6 signaling plays in this condition as revealed by large Mendelian randomization studies 44, 45 ."

## Claim

IL-6 signaling plays a major role in atherosclerotic cardiovascular disease.



# Framework

- Task 1 (Relevant Paper Retrieval)
  - Any retrieval model
- Task 2 (Rationale Sentence Selection)
  - For each sentence  $s$  in a relevant paper, perform binary classification (rationale sentence / not rationale sentence)
  - [CLS] claim [SEP]  $s$  [SEP]
- Task 3 (Stance Prediction)
  - Combine all rationale sentences together and perform three-class classification (support/refute/neutral)
  - [CLS] claim [SEP] rationale1 rationale2 ... rationale [SEP]

## Performance of Each Task

	RATIONAL-SELECT.			LABEL-PRED.
<b>Training data</b>	P	R	F1	ACC.
FEVER	41.5	57.9	48.4	67.6
UKP Snopes	42.5	62.3	50.5	71.3
SciFACT	73.7	70.5	<b>72.1</b>	75.7
FEVER + SciFACT	72.4	67.2	69.7	<b>81.9</b>
<b>Sentence encoder</b>	P	R	F1	ACC.
SciBERT	74.5	74.3	<b>74.4</b>	69.2
BioMedRoBERTa	75.3	69.9	72.5	71.7
RoBERTa-base	76.1	66.1	70.8	62.9
RoBERTa-large	73.7	70.5	72.1	<b>75.7</b>
<b>Model inputs</b>	P	R	F1	ACC.
Claim-only	-	-	-	44.5
Abstract-only	60.1	60.9	60.5	53.3

# End-to-End Performance

Retrieval	Model		Sentence-level						Abstract-level					
			Selection-Only			Selection+Label			Label-Only			Label+Rationale		
			P	R	F1	P	R	F1	P	R	F1	P	R	F1
<b>Oracle abstract</b>	Oracle rationale	1	100.0	80.5	89.2 <sub>2.1</sub>	89.6	72.2	79.9 <sub>3.0</sub>	90.1	77.5	83.3 <sub>2.4</sub>	90.1	77.5	83.3 <sub>2.4</sub>
	Zero-shot	2	42.5	45.1	43.8 <sub>2.0</sub>	36.1	38.4	37.2 <sub>2.3</sub>	86.9	53.6	66.3 <sub>3.1</sub>	67.9	41.9	51.8 <sub>3.4</sub>
	VERISCI	3	76.1	63.8	69.4 <sub>2.6</sub>	66.5	55.7	<b>60.6</b> <sub>3.1</sub>	87.3	65.3	74.7 <sub>2.8</sub>	84.9	63.5	<b>72.7</b> <sub>2.9</sub>
<b>Open</b>	Oracle rationale	4	100.0	56.5	72.2 <sub>3.3</sub>	87.6	49.5	63.2 <sub>3.7</sub>	88.9	54.1	67.2 <sub>3.2</sub>	88.9	54.1	67.2 <sub>3.2</sub>
	Zero-shot	5	28.7	37.6	32.5 <sub>2.3</sub>	23.7	31.1	26.9 <sub>2.3</sub>	56.0	42.3	48.2 <sub>3.3</sub>	42.3	32.0	36.4 <sub>3.3</sub>
	VERISCI	6	45.0	47.3	46.1 <sub>3.0</sub>	38.6	40.5	<b>39.5</b> <sub>3.0</sub>	47.5	47.3	47.4 <sub>3.1</sub>	46.6	46.4	<b>46.5</b> <sub>3.1</sub>

# Case Studies

---

**Claim 1:** Lopinavir / ritonavir have exhibited favorable clinical responses when used as a treatment for coronavirus.

---

**Supports:** ...*Interestingly, after lopinavir/ritonavir (Kaletra, AbbVie) was administered,  $\beta$ -coronavirus viral loads significantly decreased and no or little coronavirus titers were observed.*

---

**Refutes:** *The focused drug repurposing of known approved drugs (such as lopinavir/ritonavir) has been reported failed for curing SARS-CoV-2 infected patients.* It is urgent to generate new chemical entities against this virus ...

---

---

**Claim 2:** The coronavirus cannot thrive in warmer climates.

---

**Supports:** ...*most outbreaks display a pattern of clustering in relatively cool and dry areas...This is because the environment can mediate human-to-human transmission of SARS-CoV-2, and unsuitable climates can cause the virus to destabilize quickly...*

---

**Refutes:** ...*significant cases in the coming months are likely to occur in more humid (warmer) climates, irrespective of the climate-dependence of transmission and that summer temperatures will not substantially limit pandemic growth.*

---

# Take-Away Messages

- The ideas and techniques used in scientific paper retrieval can be generalized to a wide spectrum of scientific text mining tasks aiming to predict the semantic similarity between two text units, including different steps in **scientific claim verification**.
- Drawback
  - Not an **end-to-end** framework. Errors in rationale selection will propagate to stance prediction. (If you miss some rationale sentences, you lose some information in stance prediction.)
  - Can we merge these two steps? (e.g., [CLS] claim [SEP] entire paper [SEP])
  - *MultiVerS: Improving Scientific Claim Verification with Weak Supervision and Full-Document Context*. NAACL 2022 Findings.



Thank You!

Course Website: <https://yuzhang-teaching.github.io/CSCE689-S25.html>