

# **Part II: Improve Large Language Models for Search/Graph-Augmented Scenarios**

Bowen Jin, Yu Zhang, Yunyi Zhang, Jiawei Han

Department of Computer Science

University of Illinois at Urbana-Champaign

May 1, 2025

Tutorial Website:



# Outline

---

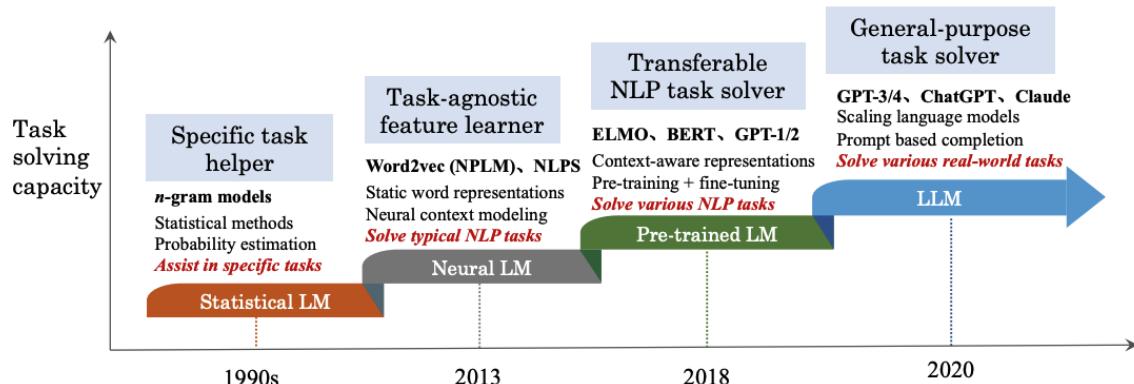
- ❑ Improving LLMs for search-augmented scenarios 

  - ❑ From General-Purpose Models to RAG-Optimized LLMs
  - ❑ Beyond RAG Pipelines: Towards Search-Enabled LLM Agents

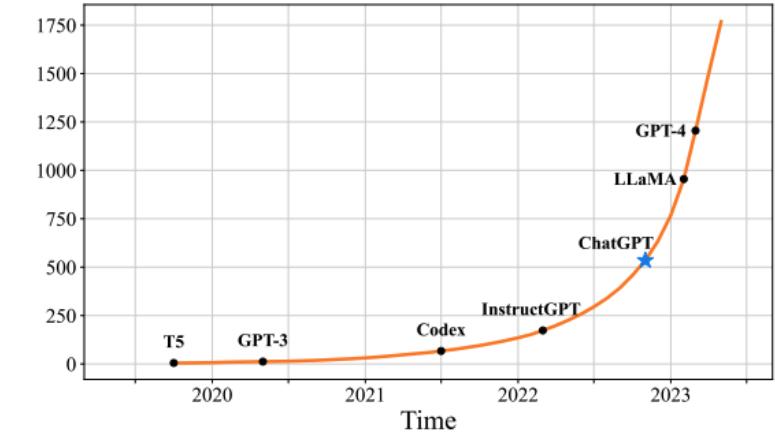
- ❑ Improving LLMs for graph-augmented scenarios
  - ❑ Pretrained LLM Retrievers for Graph-Augmented Search
  - ❑ Reasoning over Graph-Structured Knowledge with LLMs

# Large Language Models (LLMs)

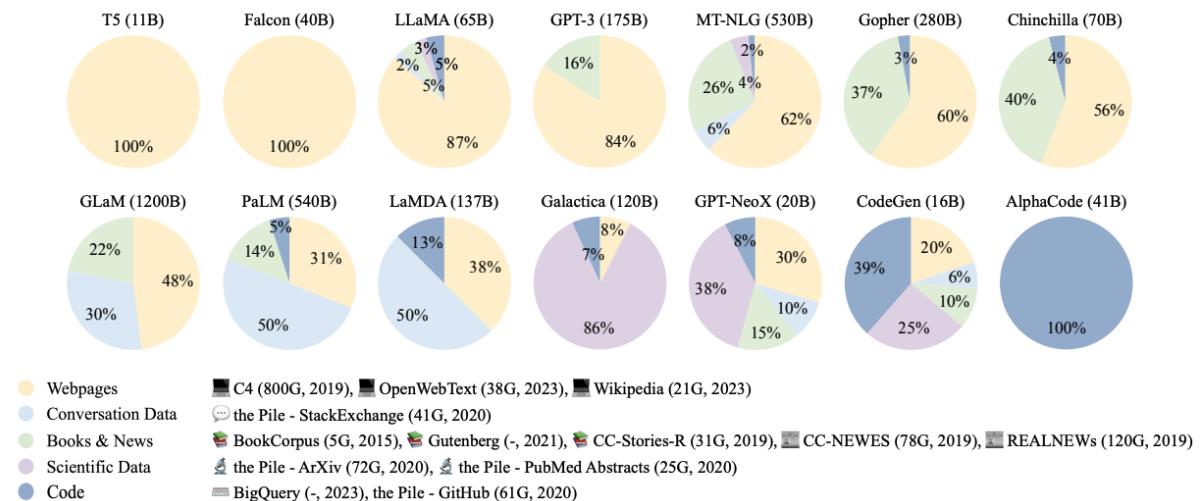
- LLMs have demonstrated their strong **text** encoding/decoding ability.



- LLMs have shown newly found emergent ability (e.g., reasoning).



(b) Query="Large Language Model"



# RAG support various applications

- Large language models often struggle with factual inaccuracies and produce hallucinated content when faced with knowledge-intensive questions.
- Retrieval Augmented Generation (RAG) incorporates information retrieved from an external knowledge sources into the context to provide up-to-date information and specify obscure facts.

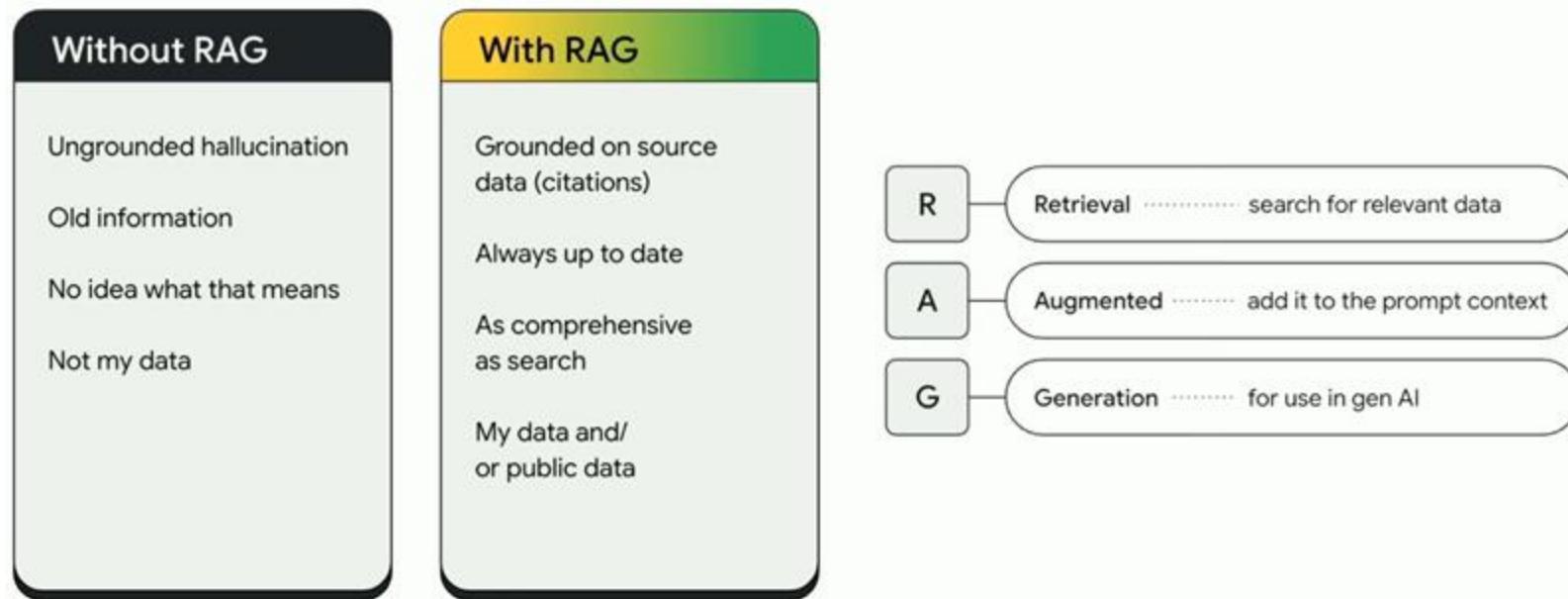
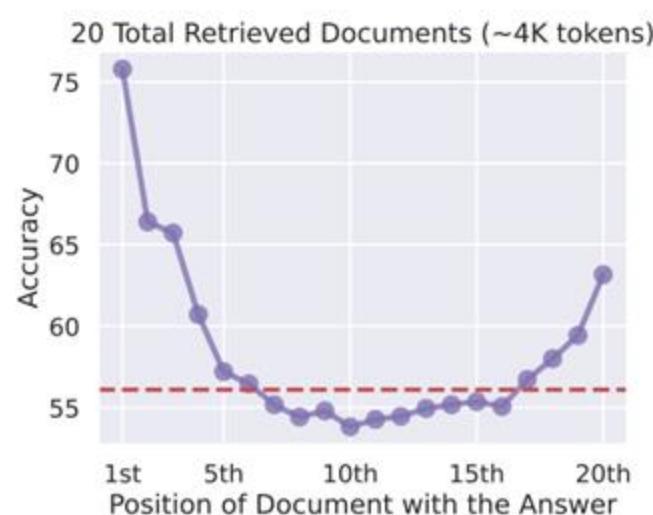


Image from: <https://cloud.google.com/use-cases/retrieval-augmented-generation>

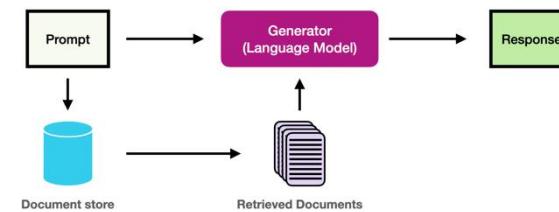
# Challenges in RAG

- In the RAG pipeline, LLMs can be easily **distracted** by the **irrelevant** retrieved information.
- Traditional RAG methods rely on **prompting**, without fundamental LLM improvement for search.

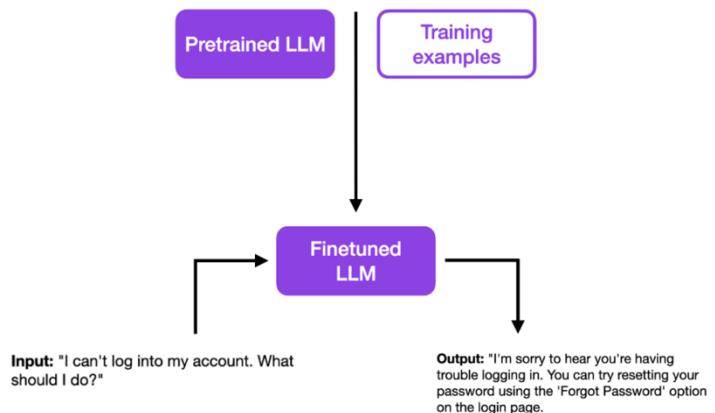


- The retrieval is always not perfect.

## Retrieval Augmented Generation



- Supervised finetuning (SFT) for RAG needs labeled trajectories and is **hard to scale**.



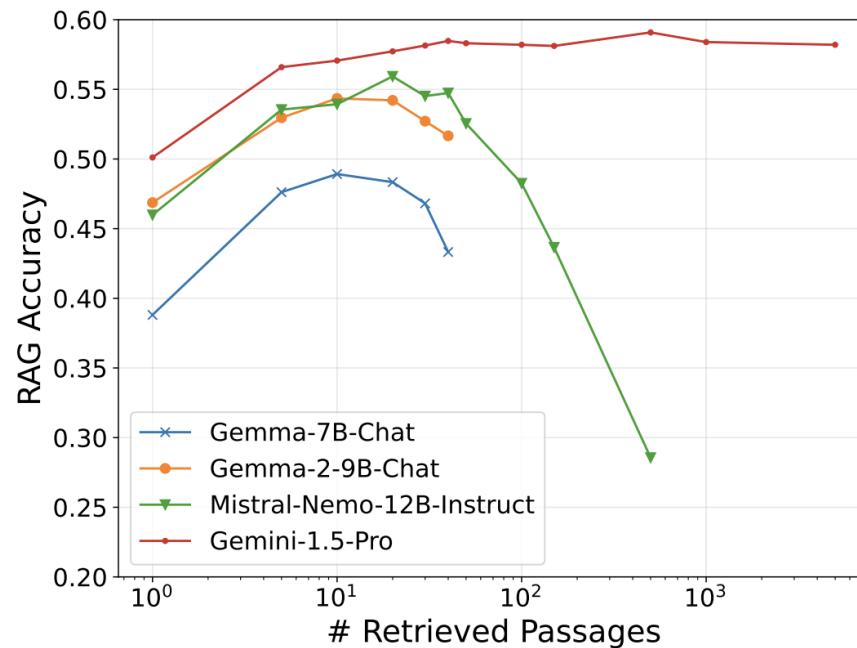
# Outline

---

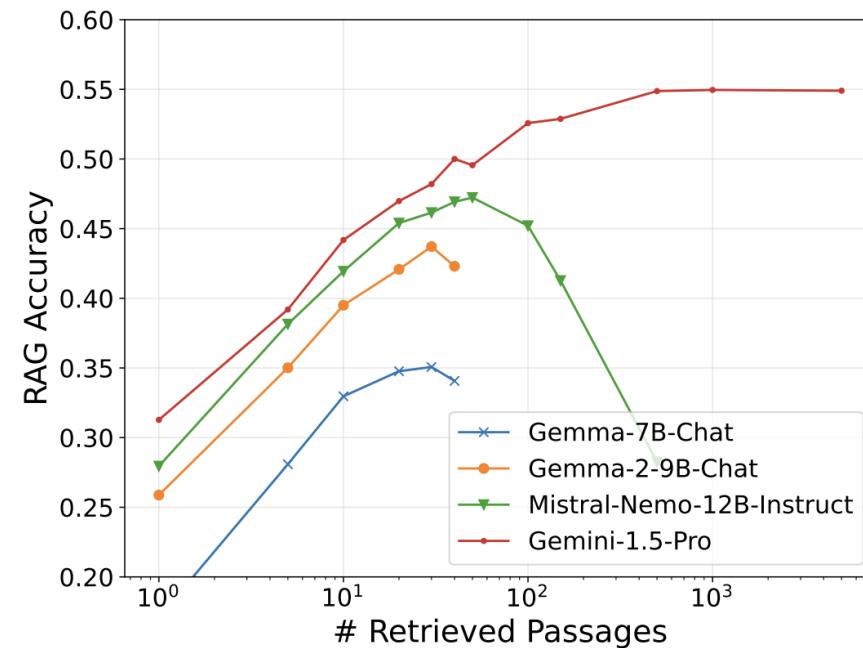
- ❑ Improving LLMs for search-augmented scenarios
  - ❑ From General-Purpose Models to RAG-Optimized LLMs
  - ❑ Beyond RAG Pipelines: Towards Search-Enabled LLM Agents
  
- ❑ Improving LLMs for graph-augmented scenarios
  - ❑ Pretrained LLM Retrievers for Graph-Augmented Search
  - ❑ Reasoning over Graph-Structured Knowledge with LLMs



# General LLMs Fails in complex long-context RAG scenarios



(a) RAG performance with e5 retriever

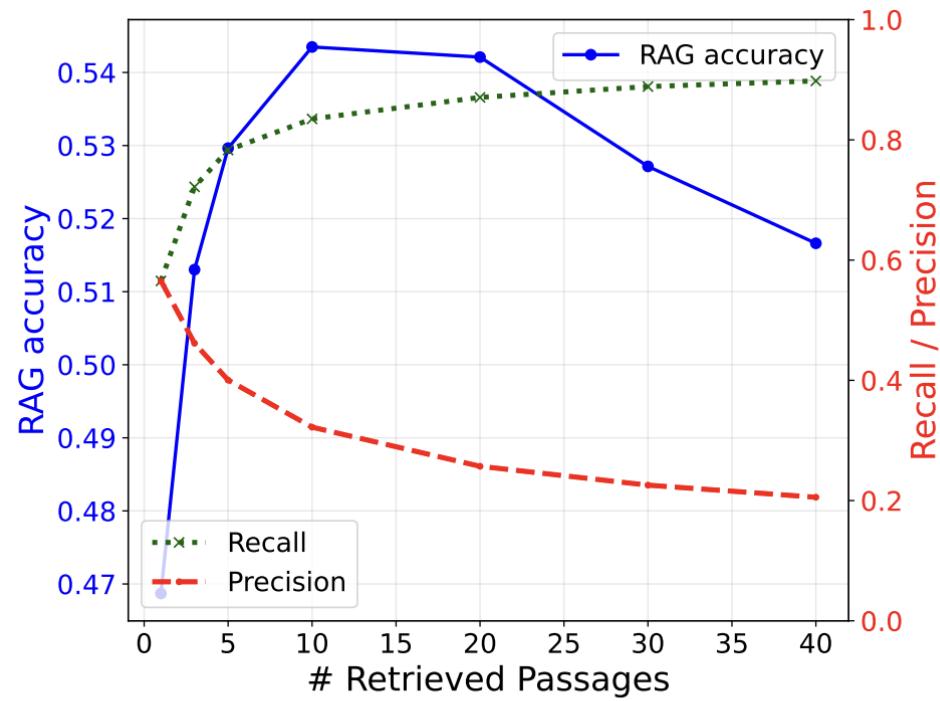


(b) RAG performance with BM25 retriever

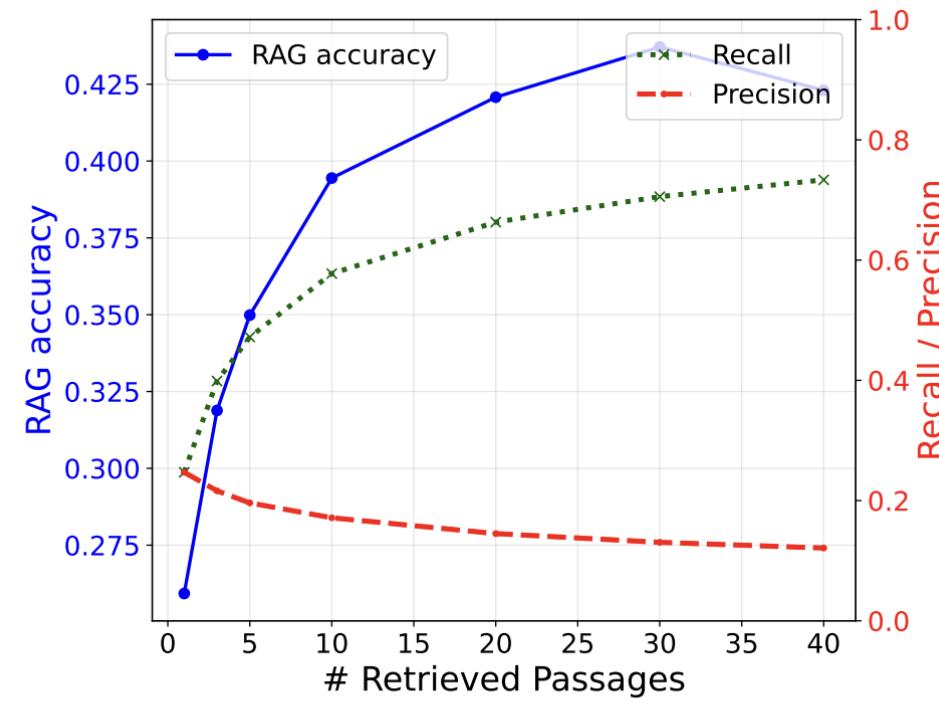
## Insights

- 1) The effectiveness of increasing retrieved context size in RAG depends on the strength of the retriever.
- 2) With a strong retriever, performance exhibits an “inverted-U pattern”, while a weak retriever shows more consistent, albeit potentially limited, improvement.
- 3) This suggests that factors beyond simply the amount of retrieved information are at play.

# Irrelevant information can mislead LLMs in RAG



(a) Retrieval with e5 retriever

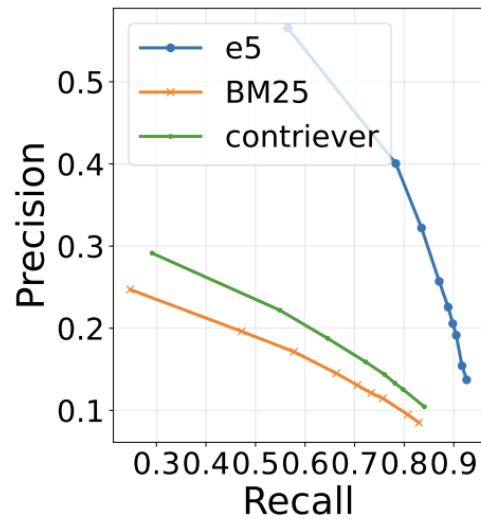


(b) Retrieval with BM25 retriever

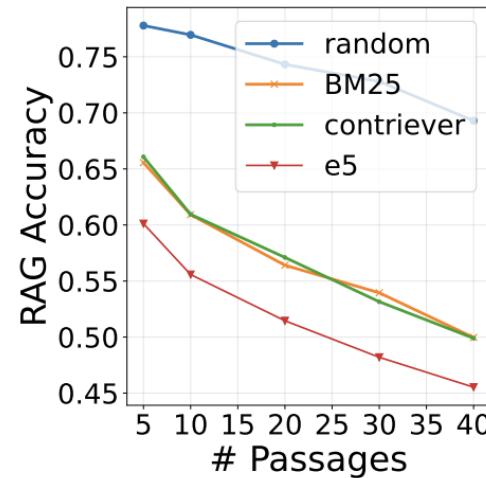
## Insights

- Influence of irrelevant passages:** The discrepancy between retrieval recall and RAG accuracy underscores the detrimental effect of irrelevant retrieved passages ("hard negatives") on the LLMs' performance.
- Limitations of precision as a metric:** The contrasting performance trends observed with e5 and BM25, despite the former's higher precision, reveal that precision alone is an inadequate measure of retrieval quality in this context, when the end-to-end performance is considered.

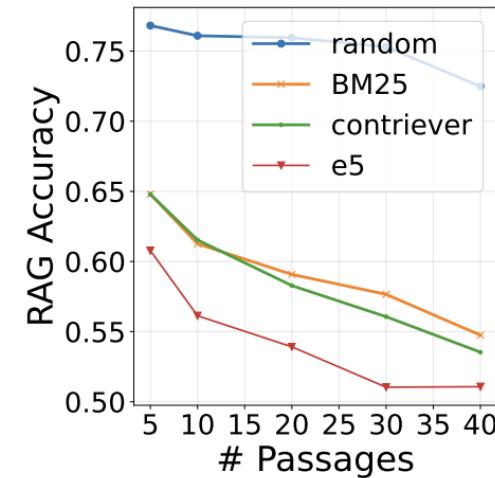
# Negative hardness correlates with retriever strength



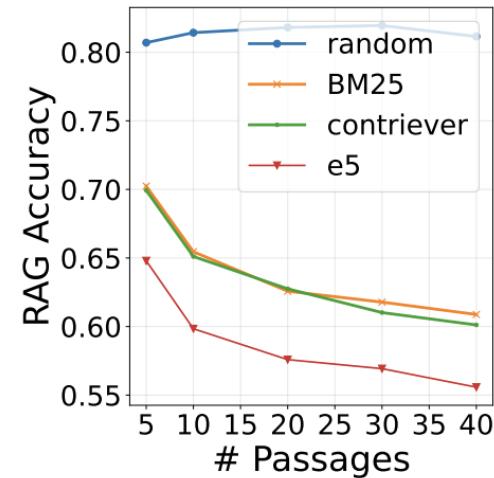
(a) Retrievers



(b) Gemma2-9B-Chat



(c) Mistral-12B-Instruct



(d) Gemini-1.5-Pro

## Observations

- **Sensitivity to hard negatives:** Across all LLMs, increasing the number of hard negative passages generally leads to a decline in RAG answer accuracy.
- **Retriever strength and hard negative difficulty:** The strength of the retriever directly correlates with the difficulty of the retrieved hard negatives. LLMs struggle more with hard negatives from stronger retrievers.
- **Distinguishing random and hard negatives:** While Gemini-1.5-Pro demonstrates robustness to random negatives, it remains susceptible to the influence of hard negatives.

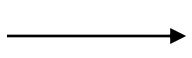
# Simple and effective training-free RAG improvement

## Retrieval reordering

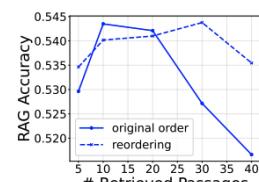
- "Lost-in-the-middle": LLMs exhibit a tendency to prioritize information presented at the beginning and end of an input sequence, while paying less attention to the middle.
- Retrieval reordering leverages the inherent "lost-in-the-middle" phenomenon observed in LLMs to mitigate the negative effects of hard negatives.

Given a query  $q$  and a set of retrieved passages  $d_1, d_2, \dots, d_k$  with decreasing relevance scores:

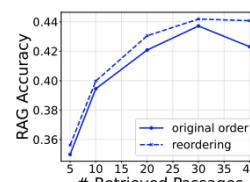
$[I, d_1, d_2, \dots, d_{k-1}, d_k, q]$



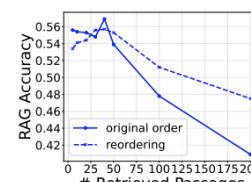
$[I, d_1, d_3, \dots, d_4, d_2, q]$



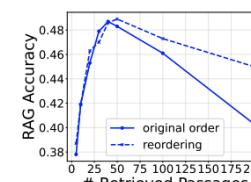
(a) NQ: Gemma2+e5



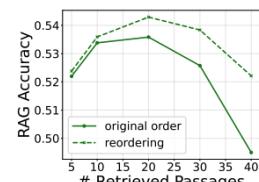
(b) NQ: Gemma2+BM25



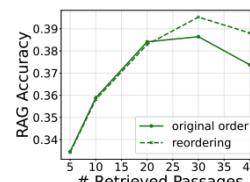
(c) NQ: Mistral+e5



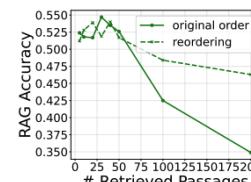
(d) NQ: Mistral+BM25



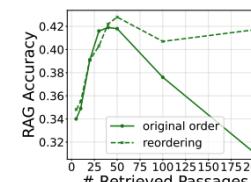
(e) PQA: Gemma2+e5



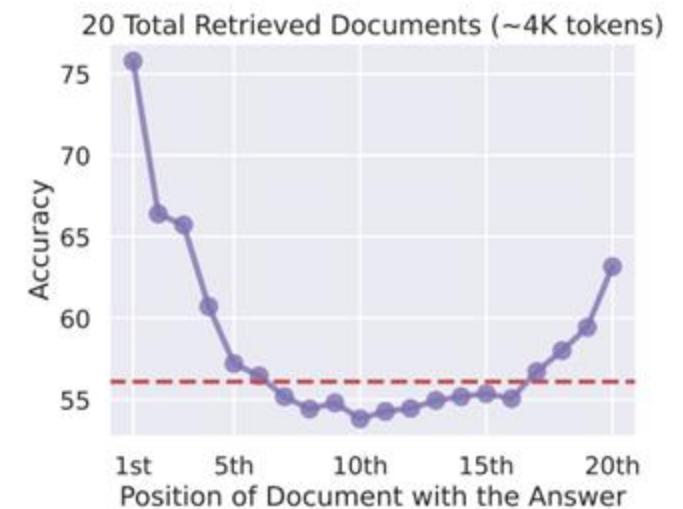
(f) PQA: Gemma2+BM25



(g) PQA: Mistral+e5



(h) PQA: Mistral+BM25



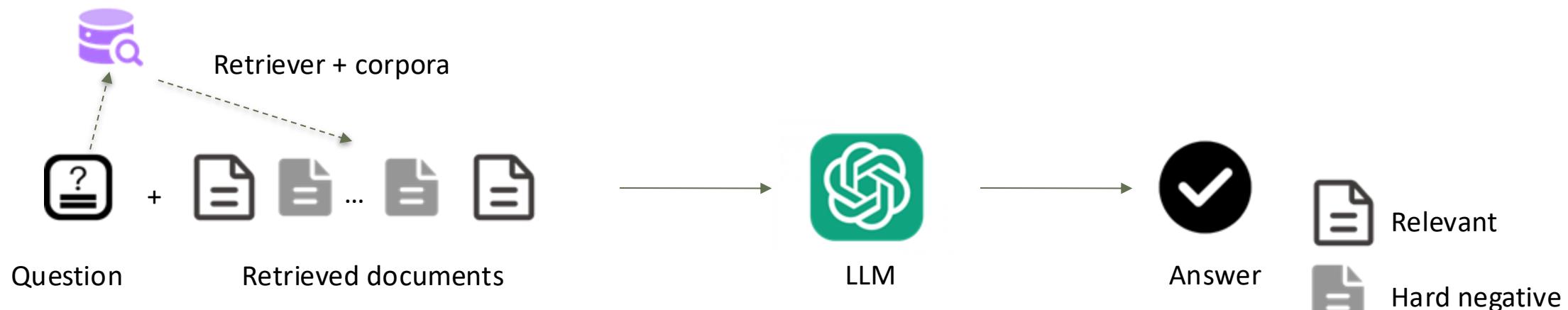
# Improving Robustness for RAG via Data-Augmented Fine-Tuning

## Vanilla LLM finetuning



Finetune the LLM to obtain knowledge

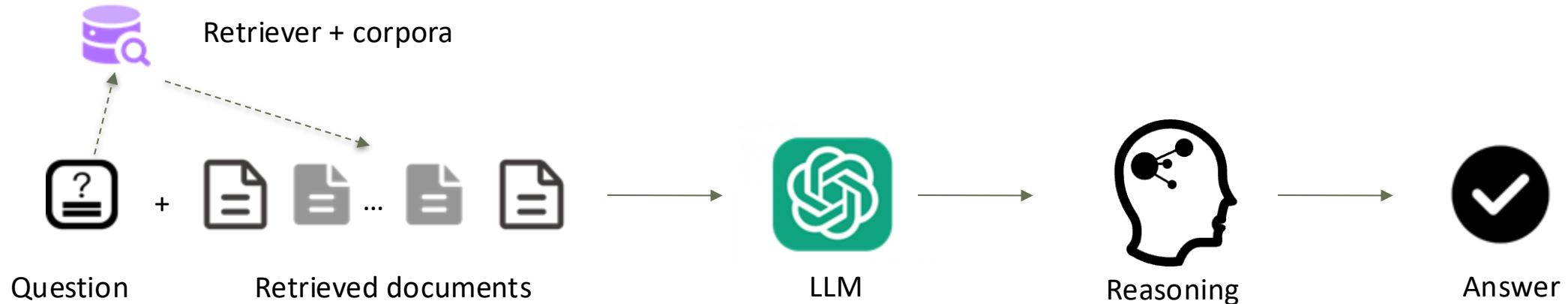
## Implicitly improving LLM for RAG through fine-tuning



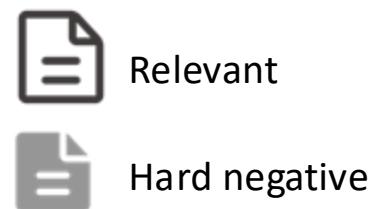
Finetune the long context LLM to be robust to hard negatives

# Improving Robustness for RAG via Data-Augmented Fine-Tuning

Enhancing relevance identification through reasoning augmentation

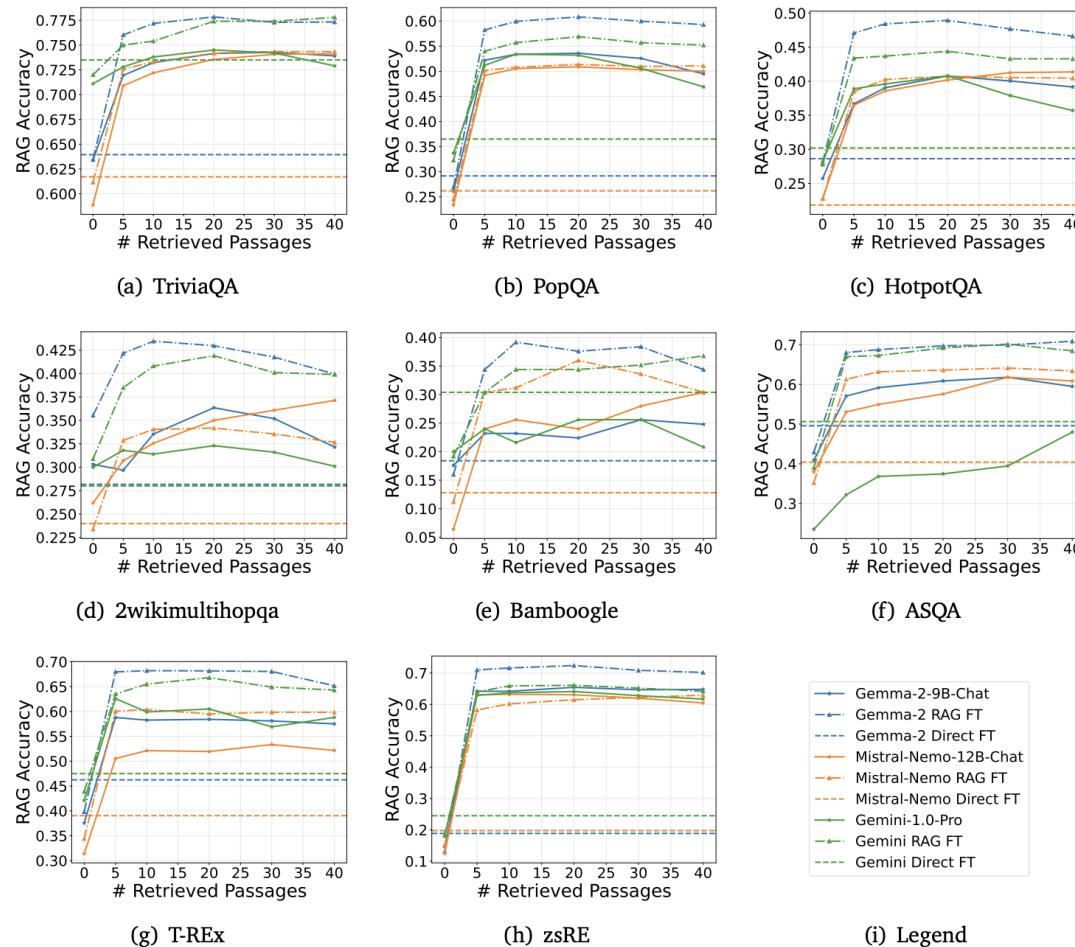


Finetune the long context LLM to reason to identify hard negatives

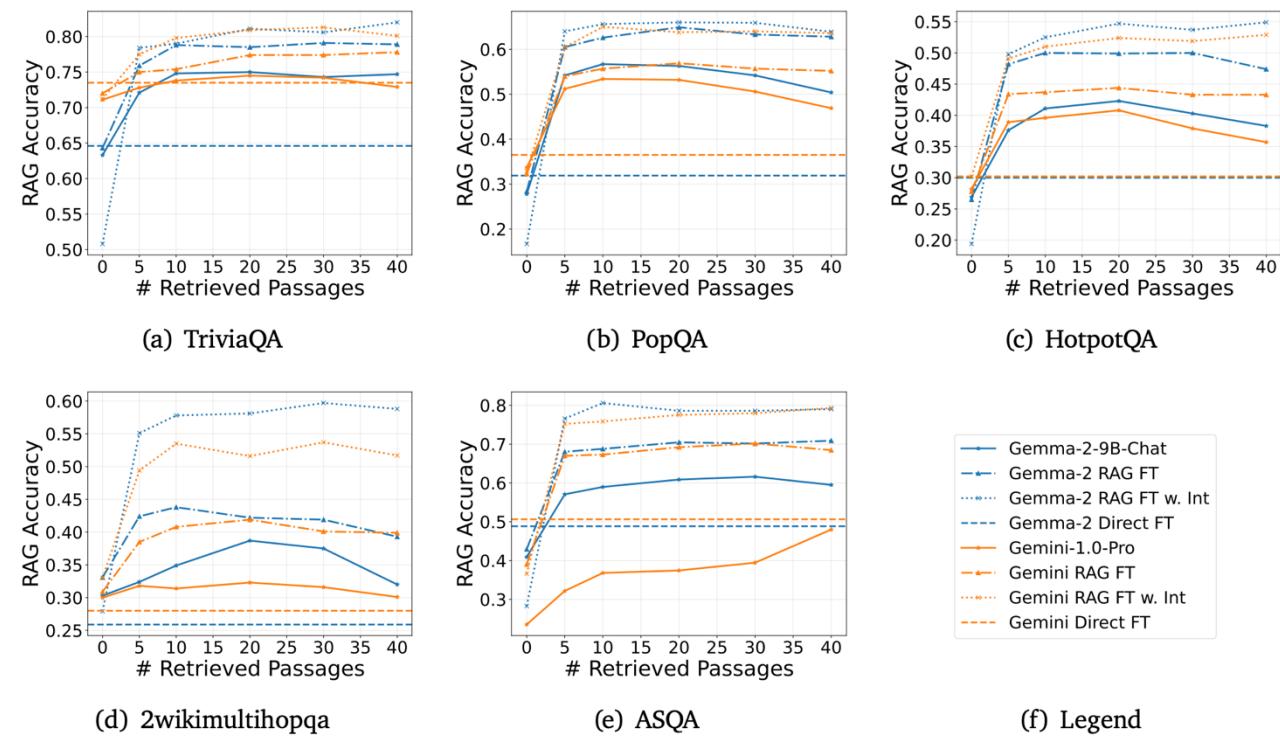


# Improving Robustness for RAG via Data-Augmented Fine-Tuning

Implicitly improving LLM robustness through fine-tuning



Enhancing relevance identification through reasoning augmentation



# Outline

---

- ❑ Improving LLMs for search-augmented scenarios
  - ❑ From General-Purpose Models to RAG-Optimized LLMs
  - ❑ Beyond RAG Pipelines: Towards Search-Enabled LLM Agents
  
- ❑ Improving LLMs for graph-augmented scenarios
  - ❑ Pretrained LLM Retrievers for Graph-Augmented Search
  - ❑ Reasoning over Graph-Structured Knowledge with LLMs



# Evolving from RAG Pipelines to Search-Enabled Agents: Search-R1

- It is an efficient, scalable **reinforcement learning (RL)** training framework which can teach LLMs to **reason** and **call search engines** in an interleaved fashion.
- We show that deepseek-R1 (zero) style RL training can be extended and let the LLM learn to call a search engine and do reasoning simultaneously.

The image shows two screenshots. On the left is a screenshot of the GitHub repository page for "Search-R1". The repository is public and has 16 stars, 112 forks, and 16 watchers. It contains 32 commits from PeterGriffinJin. The commits include adding data processing codes, fixing float bugs, and logging bugs. On the right is a tweet from Bowen Jin (@BowenJin13) introducing Search-R1. The tweet includes a link to the arXiv paper, mentions the Apache-2.0 license, and discusses the goal of training reasoning and search-augmented LLM agents with reinforcement learning. It also notes the step towards training an open-source Deep research via RL, the use of 3B base LLMs, and the ability to reason and call search engines. The tweet ends with a note about everything being fully open source and stay tuned, along with links to the code and experimental logs.

**Search-R1** Public

Unpin Unwatch 16 Fork 112 Starred 1.6K

main 1 Branch 0 Tags Go to file Add file Code

PeterGriffinJin add more data processing codes 716cd73 · yesterday 32 Commits

example Initial commit 2 months ago

public Initial commit 2 months ago

scripts add more data processing codes yesterday

search\_r1 fix potential float bug 5 days ago

verl fix turns\_stats logging bug 2 weeks ago

.gitignore add gitignore 2 months ago

LICENSE Initial commit 2 months ago

About

Search-R1: An Efficient, Scalable RL Training Framework for Reasoning & Search Engine Calling interleaved LLM based on veRL

[arxiv.org/pdf/2503.09516](https://arxiv.org/pdf/2503.09516.pdf)

Readme Apache-2.0 license Activity 1.6k stars 16 watching 112 forks

Bowen Jin @BowenJin13 Promote

Introducing Search-R1 – the first reproduction of Deepseek-R1 (zero) for training reasoning and search-augmented LLM agents with reinforcement learning!

This is a step towards training an open-source OpenAI “Deep research” via RL.

Our 3B base LLMs—including not just Qwen 2.5 but also Llama 3.2—learn to reason and call search engines all on their own!

Everything will be fully open source. Stay tuned!

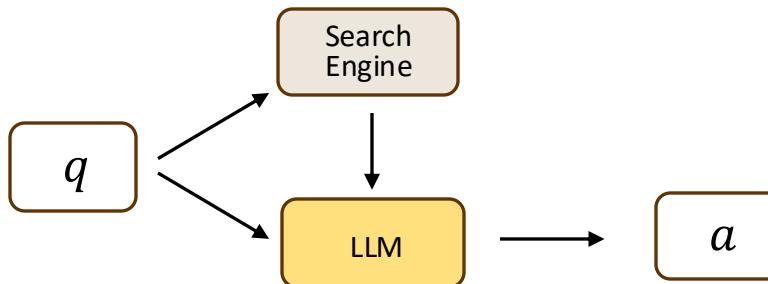
Code: [github.com/PeterGriffinJin/Search-R1](https://github.com/PeterGriffinJin/Search-R1)  
Experimental logs: [wandb.ai/peterjin/Search-R1](https://wandb.ai/peterjin/Search-R1)

#R1 #deepricerca #deepseek

12:39 PM · Feb 28, 2025 · 313.4K Views

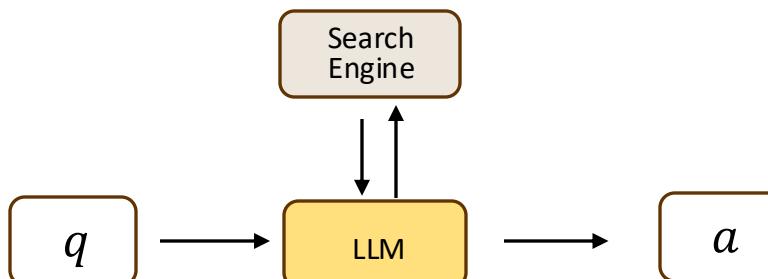
# Evolving from RAG Pipelines to Search-Enabled Agents: Search-R1

- Traditional RAG is adopted as an inference time strategy and usually focus on single turn retrieval based on the input question.



```
<query> q </query>
<info> I </info>
<answer> a <answer>
```

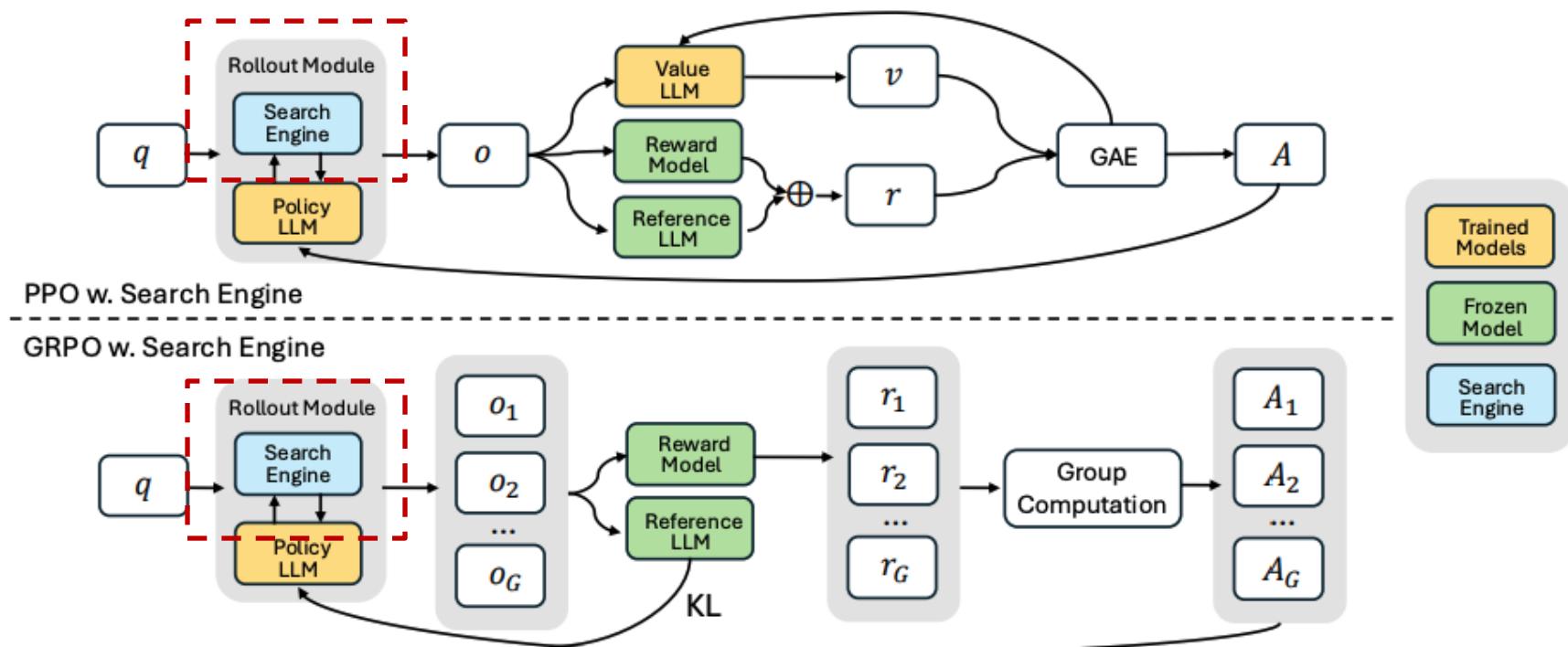
- Search-R1 improves LLM for search with RL training and enables multi-turn agentic interaction with the search engine.



```
<reason> r1 </reason>
<search> q1 </search>
<info> I1 </info>
...
<reason> rt </reason>
<answer> a <answer>
```

# Evolving from RAG Pipelines to Search-Enabled Agents: Search-R1

## □ Reinforcement Learning with a Search Engine



# Experiments

- ❑ Search-R1 consistently outperforms strong baseline methods.
- ❑ Search-R1 surpasses RL-based training for LLM reasoning without retrieval (R1).

| Method                           | NQ           | TriviaQA     | PopQA        | HotpotQA     | 2wiki        | Musique      | Bamboogle    | Avg.         |
|----------------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| <b>Qwen2.5-7b-Base/Instruct</b>  |              |              |              |              |              |              |              |              |
| Direct Inference                 | 0.134        | 0.408        | 0.140        | 0.183        | 0.250        | 0.031        | 0.120        | 0.181        |
| CoT                              | 0.048        | 0.185        | 0.054        | 0.092        | 0.111        | 0.022        | 0.232        | 0.106        |
| IRCoT                            | 0.224        | 0.478        | 0.301        | 0.133        | 0.149        | 0.072        | 0.224        | 0.239        |
| Search-o1                        | 0.151        | 0.443        | 0.131        | 0.187        | 0.176        | 0.058        | 0.296        | 0.206        |
| RAG                              | 0.349        | <u>0.585</u> | 0.392        | 0.299        | 0.235        | 0.058        | 0.208        | 0.304        |
| SFT                              | 0.318        | 0.354        | 0.121        | 0.217        | 0.259        | 0.066        | 0.112        | 0.207        |
| R1-base                          | 0.297        | 0.539        | 0.202        | 0.242        | 0.273        | 0.083        | 0.296        | 0.276        |
| R1-instruct                      | 0.270        | 0.537        | 0.199        | 0.237        | 0.292        | 0.072        | 0.293        | 0.271        |
| Search-R1-base                   | <u>0.412</u> | 0.568        | <u>0.428</u> | <u>0.356</u> | <u>0.322</u> | <u>0.142</u> | <u>0.384</u> | <u>0.373</u> |
| Search-R1-instruct               | <u>0.397</u> | <b>0.606</b> | <u>0.404</u> | <b>0.380</b> | <b>0.326</b> | <b>0.168</b> | <b>0.408</b> | <b>0.384</b> |
| <b>Qwen2.5-3b-Base/Instruct</b>  |              |              |              |              |              |              |              |              |
| Direct Inference                 | 0.106        | 0.288        | 0.108        | 0.149        | 0.244        | 0.020        | 0.024        | 0.134        |
| CoT                              | 0.023        | 0.032        | 0.005        | 0.021        | 0.021        | 0.002        | 0.000        | 0.015        |
| IRCoT                            | 0.111        | 0.312        | 0.200        | 0.164        | 0.171        | 0.067        | 0.240        | 0.181        |
| Search-o1                        | 0.238        | 0.472        | 0.262        | 0.221        | 0.218        | 0.054        | <b>0.320</b> | 0.255        |
| RAG                              | <b>0.348</b> | <b>0.544</b> | <b>0.387</b> | 0.255        | 0.226        | 0.047        | 0.080        | 0.270        |
| SFT                              | 0.249        | 0.292        | 0.104        | 0.186        | 0.248        | 0.044        | 0.112        | 0.176        |
| R1-base                          | 0.226        | 0.455        | 0.173        | 0.201        | 0.268        | 0.055        | 0.224        | 0.229        |
| R1-instruct                      | 0.210        | 0.449        | 0.171        | 0.208        | <u>0.275</u> | 0.060        | 0.192        | 0.224        |
| Search-R1-base                   | <u>0.341</u> | 0.513        | 0.362        | <u>0.263</u> | <u>0.273</u> | <u>0.076</u> | 0.211        | <u>0.292</u> |
| Search-R1-instruct               | 0.323        | <u>0.537</u> | <u>0.364</u> | <b>0.308</b> | <b>0.336</b> | <b>0.105</b> | <u>0.315</u> | <b>0.327</b> |
| <b>LLaMA3.2-3b-Base/Instruct</b> |              |              |              |              |              |              |              |              |
| Direct Inference                 | 0.139        | 0.368        | 0.124        | 0.122        | 0.107        | 0.015        | 0.064        | 0.134        |
| CoT                              | 0.246        | 0.487        | 0.166        | 0.051        | 0.083        | 0.006        | 0.024        | 0.152        |
| IRCoT                            | 0.363        | 0.566        | <u>0.428</u> | 0.238        | 0.236        | 0.072        | 0.208        | 0.301        |
| Search-o1                        | 0.107        | 0.203        | 0.093        | 0.132        | 0.117        | 0.035        | 0.176        | 0.123        |
| RAG                              | 0.317        | 0.551        | 0.337        | 0.234        | 0.118        | 0.034        | 0.064        | 0.237        |
| SFT                              | 0.320        | 0.341        | 0.122        | 0.206        | 0.257        | 0.064        | 0.120        | 0.204        |
| R1-base                          | 0.290        | 0.514        | 0.237        | 0.234        | <b>0.279</b> | 0.055        | 0.146        | 0.251        |
| R1-instruct                      | <u>0.384</u> | 0.549        | 0.228        | 0.238        | <u>0.269</u> | <u>0.074</u> | <b>0.315</b> | 0.294        |
| Search-R1-base                   | <b>0.394</b> | <b>0.596</b> | <u>0.437</u> | <u>0.280</u> | 0.264        | 0.056        | 0.105        | <u>0.305</u> |
| Search-R1-instruct               | 0.357        | <u>0.578</u> | 0.378        | <b>0.314</b> | 0.233        | <b>0.090</b> | <u>0.306</u> | <b>0.322</b> |

# Case studies

- SEARCH-R1 conduct multi-step reasoning, search, with self-verification and finally answer the question.
- Interleaved Reasoning and Retrieval Enhances Problem Analysis.
- Self-Verification though Iterative Retrieval.

**Question:** What type of profession does Chris Jericho and Gary Barlow have in common?

**Ground Truth:** musician

**SEARCH-R1:**

**<think>** I need to find the type of profession that Chris Jericho and Gary Barlow have in common. I'll search for it. **</think>**

**<search>** Chris Jericho and Gary Barlow profession **</search>**

**<information>** Doc 2>Title: "Chris Jericho" ""Do You Wanna Start a War"" (2014) and ""Judas"" (2017) are comprised entirely or predominantly of original compositions. He has also appeared on numerous TV shows over the years, having competed in the 2011 season of ""Dancing With the Stars"", lasting until the sixth week. Prior to that, Jericho hosted the ABC game show ""Downfall"" and the Revolver Golden Gods Awards. Christopher Keith Irvine was born on November 9, 1970, in Manhasset, New York, but grew up in Winnipeg. His father, ice hockey player Ted Irvine, had been playing for the New York Rangers at the time of his birth. ... **</information>**

**<think>** I found out that Chris Jericho is a professional wrestler, musician, author, and actor. Now I need to find if Gary Barlow has any of these professions in common. **</think>**

**<search>** Gary Barlow profession **</search>**

**<information>** Doc 1>Title: "Gary Barlow" hits, twelve number one singles and seven number one albums with Take That. Barlow has received six Ivor Novello Awards from the British Academy of Songwriters, Composers and Authors, including the award for Outstanding Services to British Music. He has sold over 50 million records worldwide. Barlow was born on 20 January 1971 in Frodsham, Cheshire, the second son of Colin (died 15 October 2009) and Marjorie Barlow. In Barlow's autobiography, he relates that his love of music began at an early age, and calls himself "one of those kids that's forever dancing in front of the TV looking at... **</information>**

**<think>** I found out that Gary Barlow is a singer, songwriter, musician, and record producer. Now I need to find if there is a common profession between Chris Jericho and Gary Barlow. **</think>**

**<search>** common profession of Chris Jericho and Gary Barlow **</search>**

**<information>** Doc 1>Title: "Gary Barlow" Gary Barlow Gary Barlow (born 20 January 1971) is an English singer, songwriter, musician and record producer. He is best known as the lead singer of the British pop group Take That. Barlow also served as head judge of ""The X Factor UK"" in 2011, 2012 and 2013 and ""Let It Shine"" in 2017. Barlow is one of Britain's most successful songwriters, having written fourteen number one singles and twenty-four top 10 hits. As a solo artist he has had three number one singles, six top 10 singles and two number one albums, and has additionally had seventeen top 5... **</information>**

**<think>** I found out that Gary Barlow is a singer, songwriter, musician, and record producer. Now I need to find if Chris Jericho has any of these professions in common. **</think>**

**<search>** Chris Jericho profession **</search>**

**<information>** Doc 1>Title: "Chris Jericho" Chris Jericho Christopher Keith Irvine (born November 9, 1970), better known by the ring name Chris Jericho, is an American-born professional wrestler, musician, author and actor. He currently wrestles for New Japan Pro-Wrestling (NJPW) where he is the current IWGP Intercontinental Champion in his first reign. Jericho also makes sporadic appearances for WWE on the SmackDown brand, with whom he has worked since 1999. Within professional wrestling, he is known for his over-the-top, rock star persona. Jericho is regarded as one of the greatest professional wrestlers of all time. During the 1990s, Jericho performed for the American organizations Extreme Championship... **</information>**

**<think>** I found out that Chris Jericho is a professional wrestler, musician, author, and actor. Now I can compare it to Gary Barlow's professions. I found out that Gary Barlow is a singer, songwriter, musician, and record producer. The common profession between Chris Jericho and Gary Barlow is musician. **</think>**

**<answer>** musician **</answer>**

**Question:** Curious is a women's fragrance by a singer born in what city and state?

**Ground Truth:** McComb, Mississippi

**R1:**

**<think>** The singer is Beyoncé, who was born in Houston, Texas. **</think>**  
**<answer>** Houston **</answer>**

**SEARCH-R1:**

**<think>** I need to find out which city and state a singer, by whom the women's fragrance "Curious" was created, was born in. **</think>**

**<search>** Curious fragrance information **</search>**

**<information>** Doc 1>Title: "Curious (fragrance)" Curious (fragrance) Curious is a women's fragrance by Britney Spears for Elizabeth Arden, and is the first perfume to be endorsed by Britney Spears, preceding "Fantasy." Spears reportedly earned \$52 million from the endorsement. "Curious" was released in September 2004, and was very successful internationally. It was the number one fragrance of 2004 in department stores, and in 2005, "Curious" was honored by the Fragrance Foundation as Best Women's Fragrance. The fragrance is a white floral scent. Its notes are Louisiana magnolia, golden Anjou pear, lotus flower, tuberose, star jasmine, pink cyclamen, vanilla-infused musk, sandalwood and blonde woods. Elizabeth Arden... **</information>**

**<think>** I need to find out which city and state Britney Spears, the singer behind "Curious" fragrance, was born in. **</think>**

**<search>** Britney Spears birthplace **</search>**

**<information>** Doc 1>Title: "Britney Spears" the best-selling female albums artist of the 2000s. In 2004, she launched a perfume brand with Elizabeth Arden, Inc., from which sales exceeded US\$1.5 billion. Spears was born in McComb, Mississippi, the second child of Lynne Irene Bridges and James Parnell Spears. Her maternal grandmother, Lillian Portell, was English (born in London), and one of Spears's maternal great-great-grandfathers was Maltese. Her siblings are Bryan James and Jamie Lynn. Britney was born in the Bible Belt, where socially conservative evangelical Protestantism is a particularly strong religious influence. Spears was baptized into the Southern Baptist Convention, but later in life studied **</information>**

**<think>** Since Britney Spears was born in McComb, Mississippi, I need to find out which city and state that is. **</think>**

**<search>** McComb, Mississippi location **</search>**

**<information>** Doc 1>Title: "McComb, Mississippi" McComb, Mississippi McComb is a city in Pike County, Mississippi, United States, approximately south of Jackson. As of the 2010 census, the city had a total population of 12,790. It is the principal city of the McComb, Mississippi Micropolitan Statistical Area. McComb was founded in 1872 after Henry Simpson McComb of the New Orleans, Jackson and Great Northern Railroad, a predecessor of the Illinois Central Railroad (now part of the Canadian National Railway), decided to move the railroad's maintenance shops away from New Orleans, Louisiana, to avoid the attractions of that city's saloons. The railroad purchased land in Pike County. **</information>**

**<think>** I have found that Britney Spears was born in McComb, Mississippi. Since the question asks for the city and state, I just need to confirm the state, which I already know from the previous search. **</think>**

**<answer>** McComb, Mississippi **</answer>**

# Search-R1 Resources

---



Code



X



Paper



Data / Model



Logs

# Outline

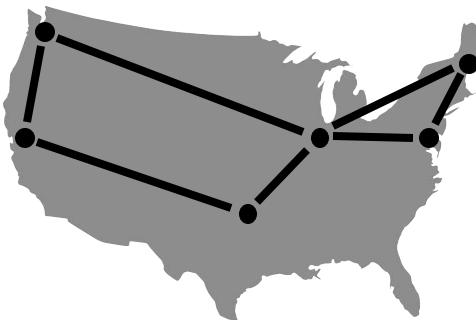
---

- ❑ Improving LLMs for search-augmented scenarios
  - ❑ From General-Purpose Models to RAG-Optimized LLMs
  - ❑ Beyond RAG Pipelines: Towards Search-Enabled LLM Agents
  
- ❑ Improving LLMs for graph-augmented scenarios 
  - ❑ Pretrained LLM Retrievers for Graph-Augmented Search
  - ❑ Reasoning over Graph-Structured Knowledge with LLMs

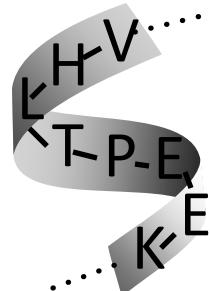
# Graphs

- Graph data is ubiquitous in real world.

Traffic Graphs



Protein Graphs

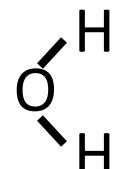


“Myoglobin holds oxygen in muscles.”

Molecule Graphs

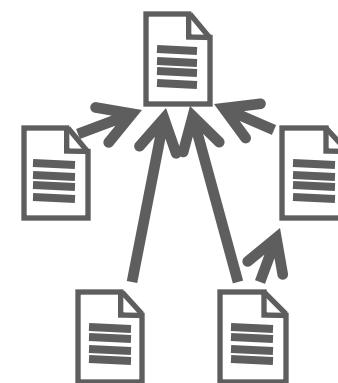


“Benzene is toxic”

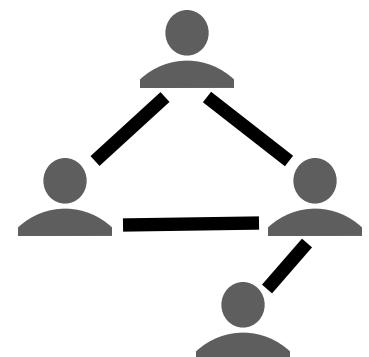


“Water is less toxic”

Academic Graphs

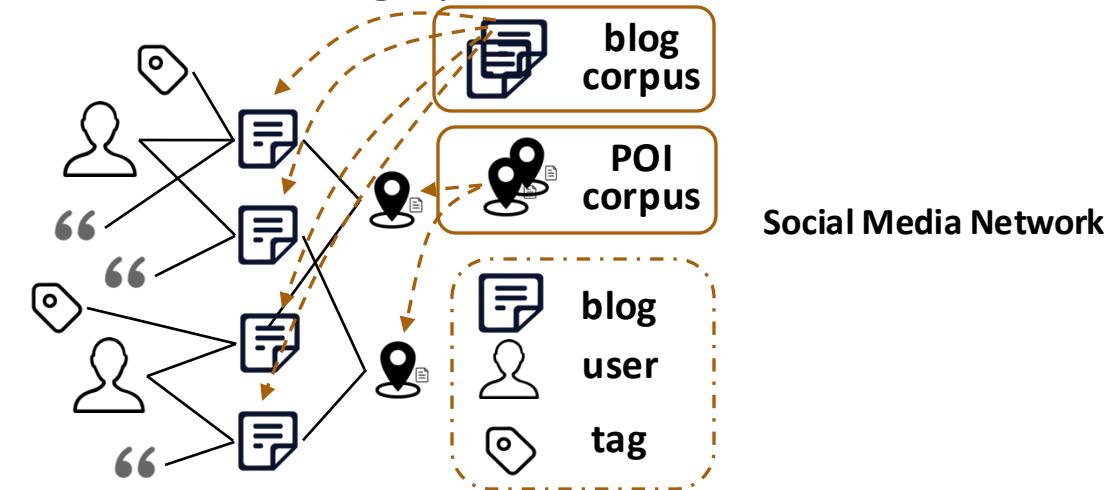
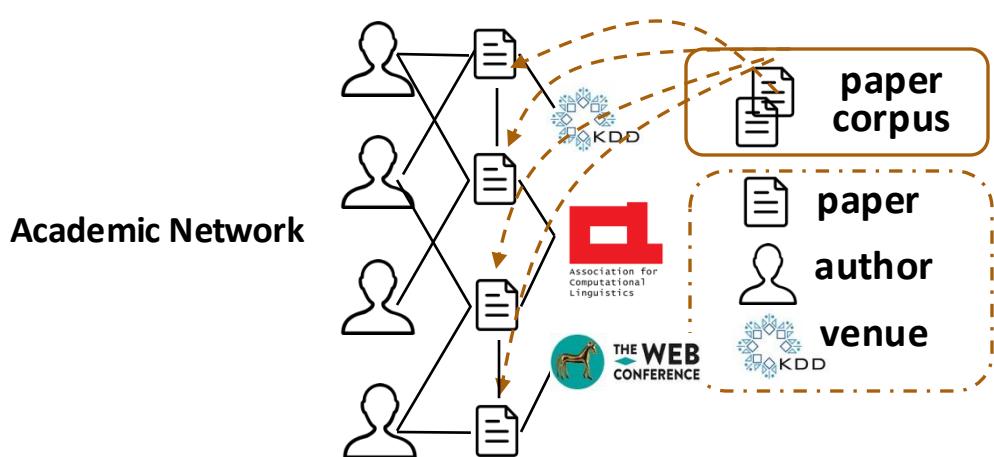


Social Graphs

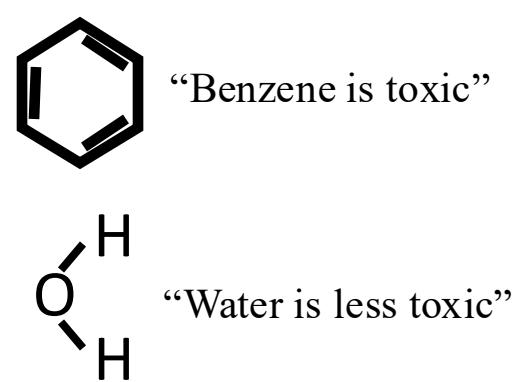
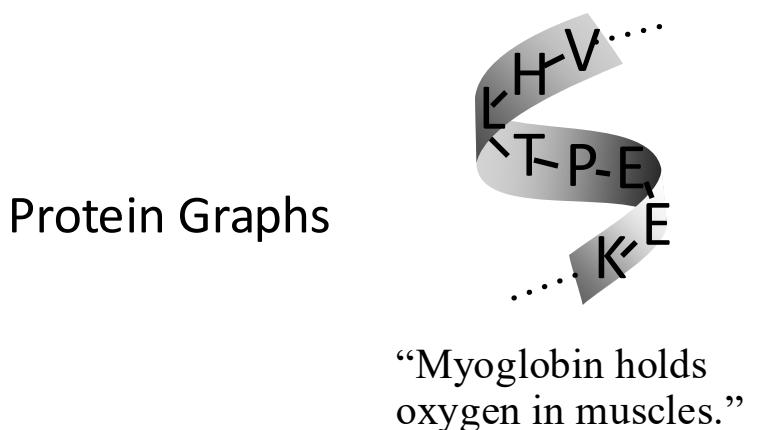


# Why LLM on Graphs?

- In real world, text and graph usually appears simultaneously.
- Text data are associated with rich structure information in the form of graphs.

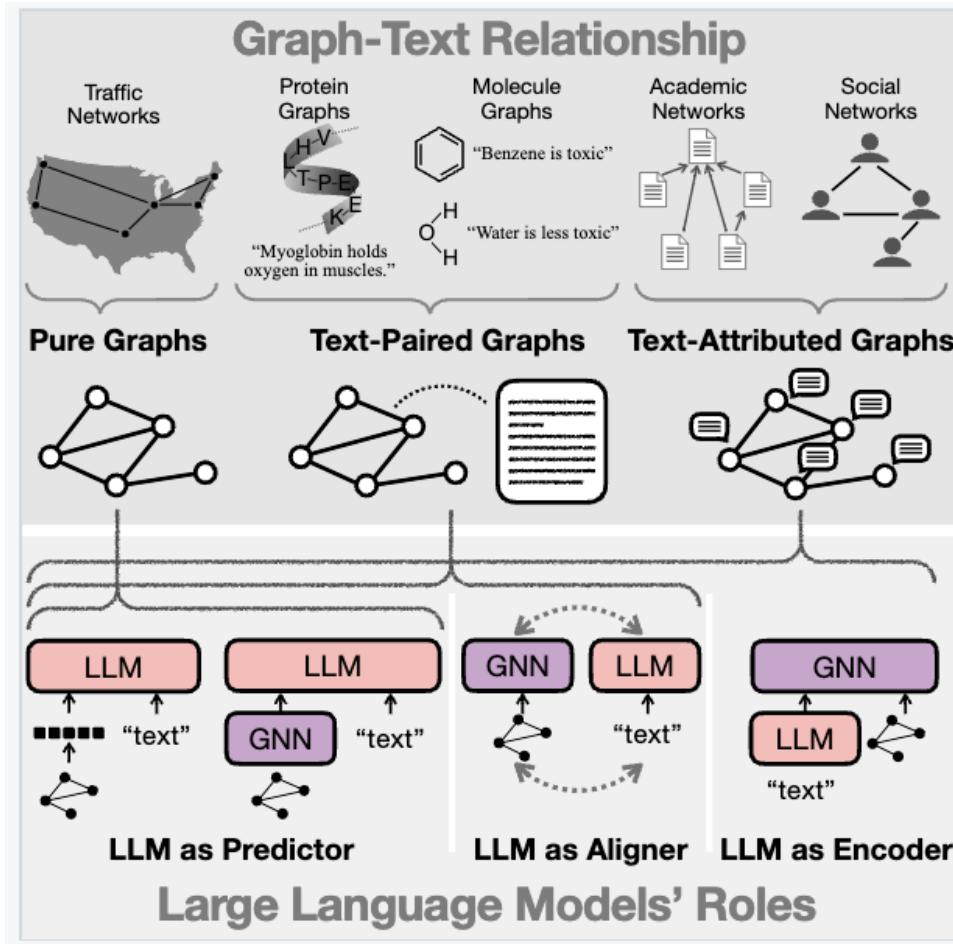


- Graph data are captioned with rich textual information.

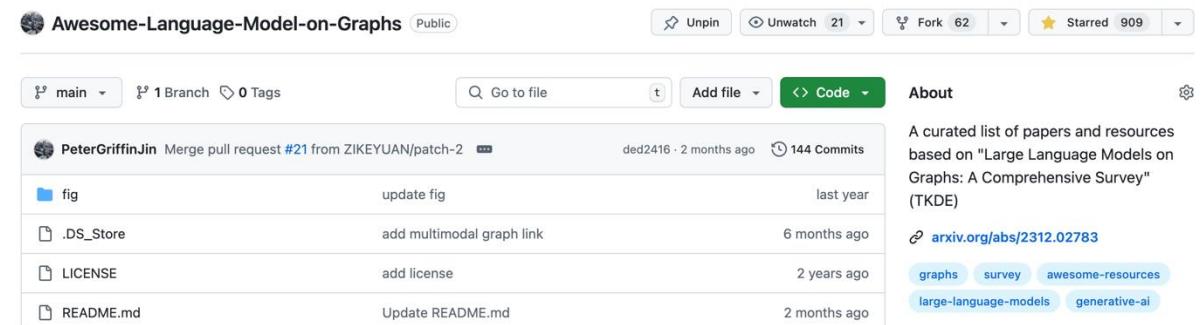


# Survey

- ❑ A survey paper of LLM & graphs



- ❑ A resource repo of LLM & graphs



paper



repo

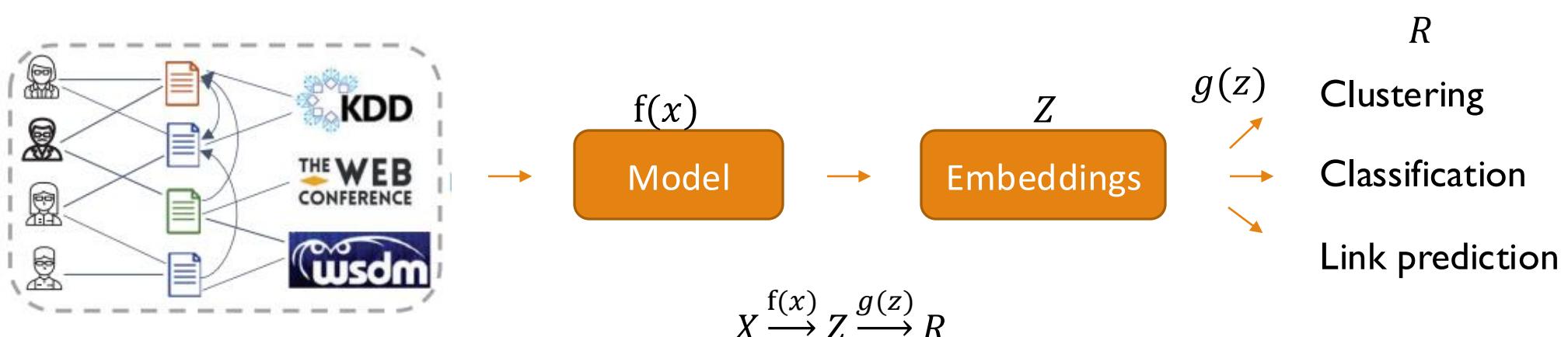
# Outline

---

- ❑ Improving LLMs for search-augmented scenarios
  - ❑ From General-Purpose Models to RAG-Optimized LLMs
  - ❑ Beyond RAG Pipelines: Towards Search-Enabled LLM Agents
  
- ❑ Improving LLMs for graph-augmented scenarios
  - ❑ Pretrained LLM Retrievers for Graph-Augmented Search 
  - ❑ Reasoning over Graph-Structured Knowledge with LLMs

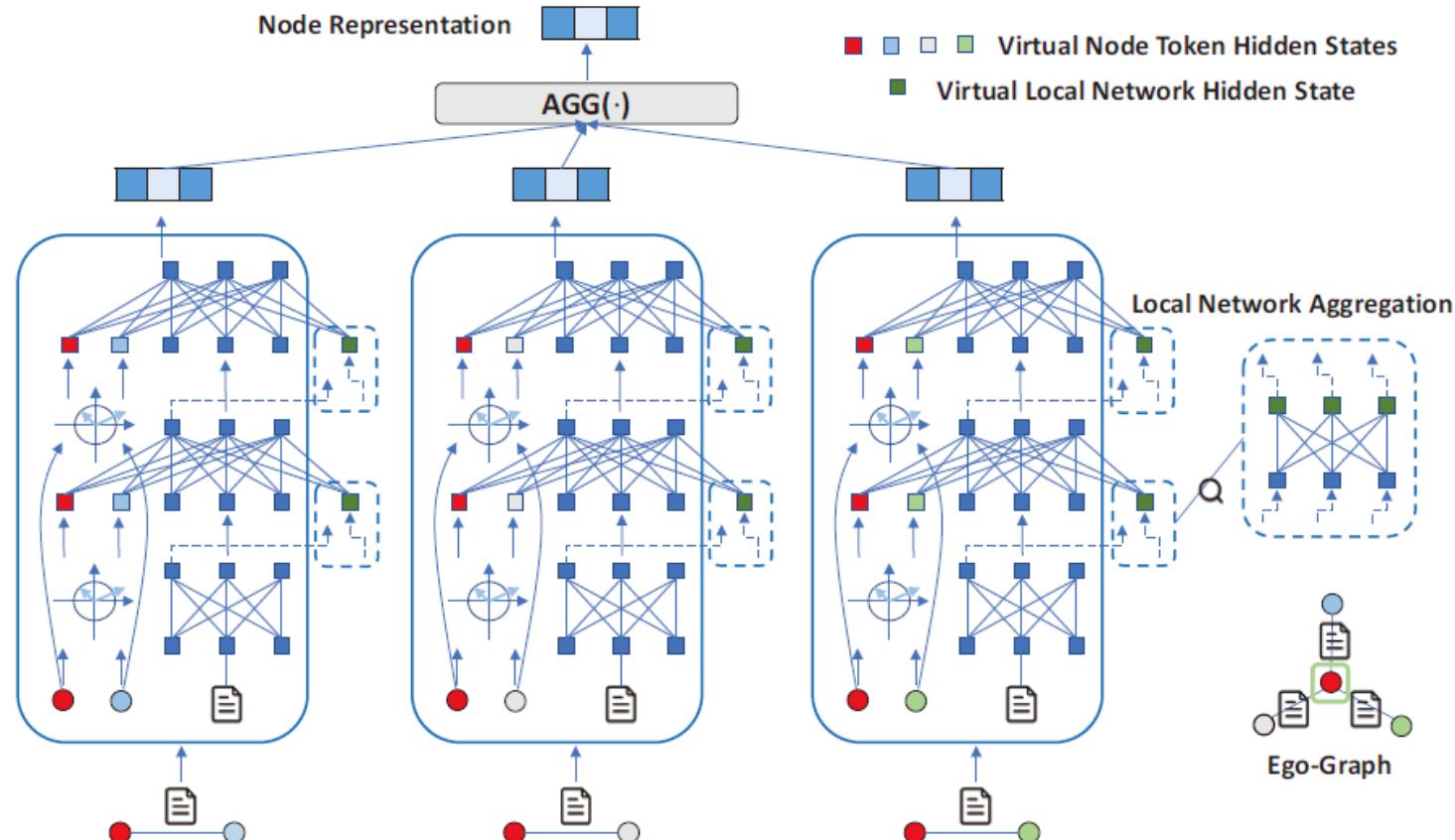
# Representation learning on text-attributed graphs

- ❑ Given a text-attributed network, people are interested in various tasks.
  - ❑ Node classification, link prediction, and node clustering.
- ❑ Learn representations for nodes/edges which can be utilized in various tasks.
  - ❑ Textual information & structure information



# Graph-Empowered LLM: Edgeformers

- Text-aware node representation learning (Edgeformer-N)
  - Aggregate edge representations
  - Enhance edge representations with node's local network structure



# Graph-Empowered LLM: Edgeformers

## □ Edge classification

Table 7: Edge classification performance on Amazon-Movie, Amazon-App, Goodreads-Crime, and Goodreads-Children.

| Model        | Amazon-Movie |              | Amazon-Apps  |              | Goodreads-Crime |              | Goodreads-Children |              |
|--------------|--------------|--------------|--------------|--------------|-----------------|--------------|--------------------|--------------|
|              | Macro-F1     | Micro-F1     | Macro-F1     | Micro-F1     | Macro-F1        | Micro-F1     | Macro-F1           | Micro-F1     |
| TF-IDF       | 50.01        | 64.22        | 48.30        | 62.88        | 43.07           | 51.72        | 39.42              | 49.90        |
| TF-IDF+nodes | 53.59        | 66.34        | 50.56        | 65.08        | 49.35           | 57.50        | 47.32              | 56.78        |
| EHGNN        | 49.90        | 64.04        | 48.20        | 63.63        | 44.49           | 52.30        | 40.01              | 50.23        |
| BERT         | 61.38        | 71.36        | 59.11        | 69.27        | 56.41           | 61.29        | 51.57              | 57.72        |
| BERT+nodes   | 63.00        | 72.45        | 59.72        | 70.82        | 58.64           | 65.02        | 54.42              | 60.46        |
| BERT+EHGNN   | 61.45        | 70.73        | 58.86        | 70.79        | 56.92           | 61.66        | 52.46              | 57.97        |
| BERT+MaxSAGE | 61.57        | 70.79        | 58.95        | 70.45        | 57.20           | 61.98        | 52.75              | 58.53        |
| GraphFormers | 61.73        | 71.52        | 59.67        | 70.19        | 57.49           | 62.37        | 52.93              | 58.34        |
| Edgeformer-E | <b>64.18</b> | <b>73.59</b> | <b>60.67</b> | <b>71.28</b> | <b>61.03</b>    | <b>65.86</b> | <b>57.45</b>       | <b>61.71</b> |

## □ Link prediction

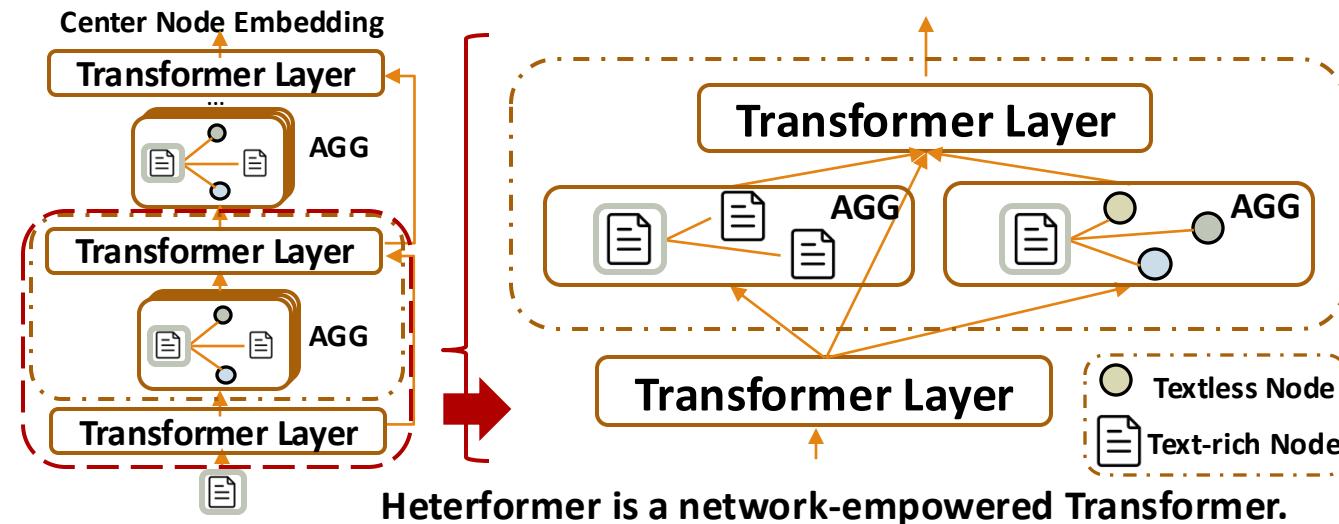
| Model         | Amazon-Movie  |               | Amazon-Apps   |               | Goodreads-Crime |               | Goodreads-Children |               | StackOverflow |               |
|---------------|---------------|---------------|---------------|---------------|-----------------|---------------|--------------------|---------------|---------------|---------------|
|               | MRR           | NDCG          | MRR           | NDCG          | MRR             | NDCG          | MRR                | NDCG          | MRR           | NDCG          |
| MF            | 0.2032        | 0.3546        | 0.1482        | 0.3052        | 0.1923          | 0.3443        | 0.1137             | 0.2716        | 0.1040        | 0.2642        |
| MeanSAGE      | 0.2138        | 0.3657        | 0.1766        | 0.3343        | 0.1832          | 0.3368        | 0.1066             | 0.2647        | 0.1174        | 0.2768        |
| MaxSAGE       | 0.2178        | 0.3694        | 0.1674        | 0.3258        | 0.1846          | 0.3387        | 0.1066             | 0.2647        | 0.1173        | 0.2769        |
| GIN           | 0.2140        | 0.3648        | 0.1797        | 0.3362        | 0.1846          | 0.3374        | 0.1128             | 0.2700        | 0.1189        | 0.2778        |
| CensNet       | 0.2048        | 0.3568        | 0.1894        | 0.3457        | 0.1880          | 0.3398        | 0.1157             | 0.2726        | 0.1235        | 0.2806        |
| NENN          | 0.2565        | 0.4032        | 0.1996        | 0.3552        | 0.2173          | 0.3670        | 0.1297             | 0.2854        | 0.1257        | 0.2854        |
| BERT          | 0.2391        | 0.3864        | 0.1790        | 0.3350        | 0.1986          | 0.3498        | 0.1274             | 0.2836        | 0.1666        | 0.3252        |
| BERT+MaxSAGE  | 0.2780        | 0.4224        | 0.2055        | 0.3602        | 0.2193          | 0.3694        | 0.1312             | 0.2872        | 0.1681        | 0.3264        |
| BERT+MeanSAGE | 0.2491        | 0.3972        | 0.1983        | 0.3540        | 0.1952          | 0.3477        | 0.1223             | 0.2791        | 0.1678        | 0.3264        |
| BERT+GIN      | 0.2573        | 0.4037        | 0.2000        | 0.3552        | 0.2007          | 0.3522        | 0.1238             | 0.2801        | 0.1708        | 0.3279        |
| GraphFormers  | 0.2756        | 0.4198        | 0.2066        | 0.3607        | 0.2176          | 0.3684        | 0.1323             | 0.2887        | 0.1693        | 0.3278        |
| BERT+CensNet  | 0.1919        | 0.3462        | 0.1544        | 0.3132        | 0.1437          | 0.3000        | 0.0847             | 0.2436        | 0.1173        | 0.2789        |
| BERT+NENN     | 0.2821        | 0.4256        | 0.2127        | 0.3666        | 0.2262          | 0.3756        | 0.1365             | 0.2925        | 0.1619        | 0.3215        |
| Edgeformer-N  | <b>0.2919</b> | <b>0.4344</b> | <b>0.2239</b> | <b>0.3771</b> | <b>0.2395</b>   | <b>0.3875</b> | <b>0.1446</b>      | <b>0.3000</b> | <b>0.1754</b> | <b>0.3339</b> |
| + Δ %         | <b>3.5%</b>   | <b>2.1%</b>   | <b>5.3%</b>   | <b>2.9%</b>   | <b>5.9%</b>     | <b>3.2%</b>   | <b>5.9%</b>        | <b>2.6%</b>   | <b>2.7%</b>   | <b>1.8%</b>   |

## □ Node classification

| Model        | Amazon-Movie         |                      |                      | Amazon-Apps          |                      |                      |
|--------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
|              | Macro-F1             | Micro-F1             | PREC                 | Macro-F1             | Micro-F1             | PREC                 |
| MF           | 0.7566±0.0017        | 0.8234±0.0013        | 0.8241±0.0013        | 0.4647±0.0151        | 0.8393±0.0012        | 0.8462±0.0006        |
| CensNet      | 0.8528±0.0010        | 0.8839±0.0008        | 0.8845±0.0007        | 0.2782±0.0168        | 0.8279±0.0006        | 0.8331±0.0005        |
| NENN         | 0.9186±0.0008        | 0.9341±0.0008        | 0.9347±0.0007        | 0.3408±0.0082        | 0.8789±0.0019        | 0.8819±0.0017        |
| BERT         | 0.9209±0.0005        | 0.9361±0.0003        | 0.9367±0.0003        | 0.7608±0.0175        | 0.9283±0.0015        | 0.9337±0.0015        |
| BERT+CensNet | 0.9032±0.0006        | 0.9221±0.0004        | 0.9227±0.0004        | 0.5750±0.0277        | 0.8692±0.0034        | 0.8731±0.0028        |
| BERT+NENN    | 0.9247±0.0005        | 0.9387±0.0004        | 0.9393±0.0005        | 0.7556±0.0092        | 0.9306±0.0008        | 0.9382±0.0006        |
| Edgeformer-N | <b>0.9276±0.0007</b> | <b>0.9411±0.0006</b> | <b>0.9417±0.0005</b> | <b>0.7758±0.0100</b> | <b>0.9339±0.0007</b> | <b>0.9431±0.0005</b> |

# Graph-Empowered LLM: Heterformer

- ❑ Overall framework
  - ❑ Heterformer: a graph-empowered Transformer.
  - ❑ Unifying text semantic encoding and network signal capturing.



# Graph-Empowered LLM: Heterformer

## □ Link prediction

| Method      | DBLP           |                |                | Twitter        |                |                | Goodreads      |                |                |
|-------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
|             | PREC           | MRR            | NDCG           | PREC           | MRR            | NDCG           | PREC           | MRR            | NDCG           |
| MeanSAGE    | 0.7019         | 0.7964         | 0.8437         | 0.6489         | 0.7450         | 0.7991         | 0.6302         | 0.7409         | 0.8001         |
| BERT        | 0.7569         | 0.8340         | 0.8726         | 0.7179         | 0.7833         | 0.8265         | 0.5571         | 0.6668         | 0.7395         |
| Homo GNN    | BERT+MeanSAGE  | 0.8131         | 0.8779         | 0.9070         | 0.7201         | 0.7845         | 0.8275         | 0.7301         | 0.8167         |
|             | BERT+MAXSAGE   | 0.8193         | 0.8825         | 0.9105         | 0.7198         | 0.7845         | 0.8276         | 0.7280         | 0.8164         |
|             | BERT+GAT       | 0.8119         | 0.8771         | 0.9063         | 0.7231         | 0.7873         | 0.8300         | 0.7333         | 0.8170         |
|             | GraphFormers   | 0.8324         | 0.8916         | 0.9175         | 0.7258         | 0.7891         | 0.8312         | 0.7444         | 0.8260         |
| Hetero GNN  | BERT+RGCN      | 0.7979         | 0.8633         | 0.8945         | 0.7111         | 0.7764         | 0.8209         | 0.7488         | 0.8303         |
|             | BERT+HAN       | 0.8136         | 0.8782         | 0.9072         | 0.7237         | 0.7880         | 0.8306         | 0.7329         | 0.8174         |
|             | BERT+HGT       | 0.8170         | 0.8814         | 0.9098         | 0.7153         | 0.7800         | 0.8237         | 0.7224         | 0.8112         |
|             | BERT+SHGN      | 0.8149         | 0.8785         | 0.9074         | 0.7218         | 0.7866         | 0.8295         | 0.7362         | 0.8195         |
|             | GraphFormers++ | 0.8233         | 0.8856         | 0.9130         | 0.7159         | 0.7799         | 0.8236         | 0.7536         | 0.8328         |
|             | Heterformer    | <b>0.8474*</b> | <b>0.9019*</b> | <b>0.9255*</b> | <b>0.7272*</b> | <b>0.7908*</b> | <b>0.8328*</b> | <b>0.7633*</b> | <b>0.8400*</b> |
| Heterformer |                |                |                |                |                |                |                |                |                |
| Heterformer |                |                |                |                |                |                |                |                |                |

## □ Node clustering

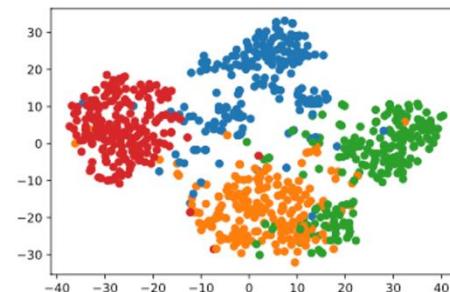
| Method         | DBLP           |                | Goodreads     |               |
|----------------|----------------|----------------|---------------|---------------|
|                | NMI            | ARI            | NMI           | ARI           |
| BERT           | 0.2570         | 0.3349         | 0.2325        | 0.4013        |
| BERT+MaxSAGE   | 0.2615         | 0.3490         | 0.2205        | 0.4173        |
| BERT+MeanSAGE  | 0.2628         | 0.3488         | <b>0.2449</b> | <b>0.4329</b> |
| BERT+GAT       | 0.2598         | 0.3419         | 0.2408        | 0.4185        |
| GraphFormers   | 0.2633         | 0.3455         | 0.2362        | 0.4139        |
| BERT+HAN       | 0.2568         | 0.3401         | 0.2391        | 0.4266        |
| BERT+HGT       | 0.2469         | 0.3392         | 0.2427        | 0.4296        |
| BERT+SHGN      | 0.2589         | 0.3431         | 0.2373        | 0.4171        |
| GraphFormers++ | 0.2566         | 0.3432         | 0.2372        | 0.4211        |
| Heterformer    | <b>0.2707*</b> | <b>0.3639*</b> | 0.2429        | 0.4199        |

## □ Node classification

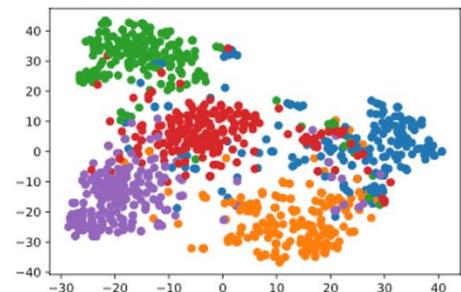
Table 3: Transductive text-rich node classification.

| Method         | DBLP           |                | Goodreads      |                |
|----------------|----------------|----------------|----------------|----------------|
|                | Micro-F1       | Macro-F1       | Micro-F1       | Macro-F1       |
| BERT           | 0.6119         | 0.5476         | 0.8364         | 0.7713         |
| BERT+MaxSAGE   | 0.6179         | 0.5511         | 0.8447         | 0.7866         |
| BERT+MeanSAGE  | 0.6198         | 0.5522         | 0.8420         | 0.7826         |
| BERT+GAT       | 0.5943         | 0.5175         | 0.8328         | 0.7713         |
| GraphFormers   | 0.6256         | 0.5616         | 0.8388         | 0.7786         |
| BERT+HAN       | 0.5965         | 0.5211         | 0.8351         | 0.7747         |
| BERT+HGT       | 0.6575         | 0.5951         | 0.8474         | 0.7928         |
| BERT+SHGN      | 0.5982         | 0.5214         | 0.8345         | 0.7737         |
| GraphFormers++ | 0.6474         | 0.5790         | 0.8516         | 0.7993         |
| Heterformer    | <b>0.6695*</b> | <b>0.6062*</b> | <b>0.8578*</b> | <b>0.8076*</b> |

## □ Embedding visualization



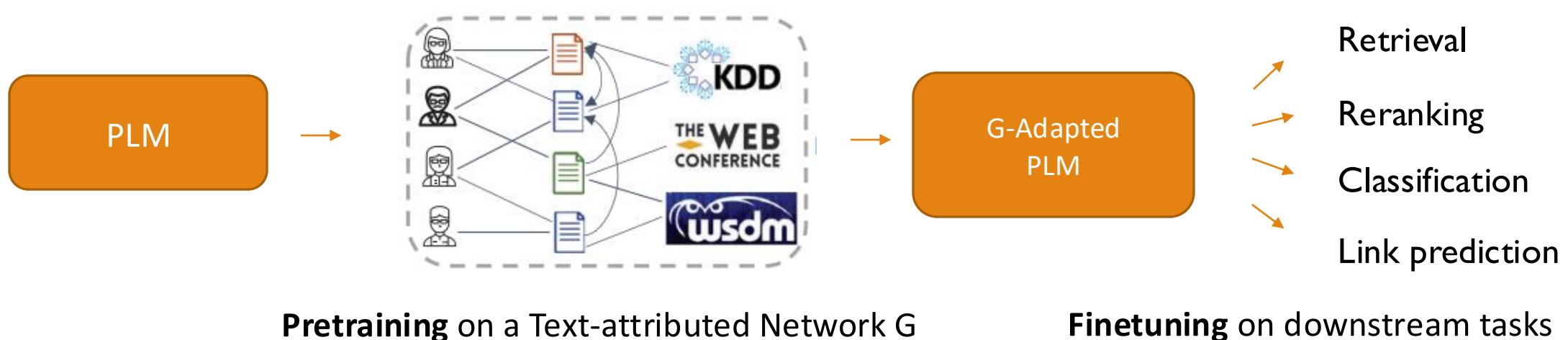
(a) DBLP



(b) Goodreads

# Why do we need language model pretraining on network?

- Given a text-rich network, people are interested in various downstream tasks
  - Document/node classification, document retrieval and link prediction
- Text-attributed network contains rich unsupervised semantic information
  - Alleviate human labeling burden for downstream tasks



# Language Model Pretraining: Patton

- Pretraining strategy 1: Network-contextualized masked language modeling
  - In node MLM -> Network contextualized MLM
  - Use both in-node text context and neighbor node context to conduct masked token prediction
  - Facilitate the LM to understand both in-node token correlation and network-contextualized text semantic relatedness

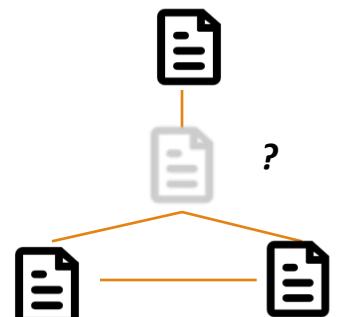
$$\mathcal{L}_{\text{NMLM}} = - \sum_{i \in M_t} \log p(w_i | \mathbf{H}_x, \mathbf{z}_x),$$



## □ Pretraining strategy 2: Masked Node Prediction

- We dynamically hold out a subset of nodes from the network ( $M_v \subseteq V$ ), mask them, and train the LM to predict the masked nodes based on the adjacent network structure.
- LM will absorb document semantic hints hidden inside the network structure.

$$\mathcal{L}_{\text{MNP}} = - \sum_{v_j \in M_v} \log p(v_j | \mathbf{G}_{v_j})$$



# Language Model Pretraining: Patton

## □ Retrieval

Table 3: Experiment results on Retrieval. We show the mean<sub>std</sub> of three runs for all the methods.

| Method           | Mathematics                 |                             | Geology                     |                             | Economy                     |                             | Clothes                     |                             | Sports                      |                             |
|------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|
|                  | R@50                        | R@100                       |
| BM25             | 20.76                       | 24.55                       | 19.02                       | 20.92                       | 19.14                       | 22.49                       | 15.76                       | 15.88                       | 22.00                       | 23.96                       |
| BERT             | 16.73 <sub>.17</sub>        | 22.66 <sub>.18</sub>        | 18.82 <sub>.39</sub>        | 25.94 <sub>.39</sub>        | 23.95 <sub>.25</sub>        | 31.54 <sub>.21</sub>        | 40.77 <sub>.18</sub>        | 50.40 <sub>.41</sub>        | 32.37 <sub>.09</sub>        | 43.32 <sub>.96</sub>        |
| GraphFormers     | 16.65 <sub>.12</sub>        | 22.41 <sub>.10</sub>        | 18.92 <sub>.60</sub>        | 25.94 <sub>.39</sub>        | 24.48 <sub>.36</sub>        | 32.16 <sub>.40</sub>        | 41.77 <sub>.05</sub>        | 51.26 <sub>.27</sub>        | 32.39 <sub>.89</sub>        | 43.29 <sub>.12</sub>        |
| SciBERT          | 24.70 <sub>.17</sub>        | 33.55 <sub>.31</sub>        | 23.71 <sub>.89</sub>        | 30.94 <sub>.95</sub>        | 29.80 <sub>.66</sub>        | 38.66 <sub>.52</sub>        | -                           | -                           | -                           | -                           |
| SPECTER          | 23.86 <sub>.25</sub>        | 31.11 <sub>.31</sub>        | 26.56 <sub>.05</sub>        | 34.04 <sub>.32</sub>        | 31.26 <sub>.15</sub>        | 40.79 <sub>.11</sub>        | -                           | -                           | -                           | -                           |
| SimCSE (unsup)   | 17.91 <sub>.26</sub>        | 23.19 <sub>.29</sub>        | 20.45 <sub>.20</sub>        | 26.82 <sub>.26</sub>        | 25.83 <sub>.23</sub>        | 33.42 <sub>.28</sub>        | 44.90 <sub>.35</sub>        | 54.76 <sub>.38</sub>        | 38.81 <sub>.35</sub>        | 49.30 <sub>.44</sub>        |
| SimCSE (sup)     | 20.29 <sub>.041</sub>       | 26.23 <sub>.51</sub>        | 22.34 <sub>.49</sub>        | 29.63 <sub>.55</sub>        | 28.07 <sub>.03</sub>        | 36.51 <sub>.37</sub>        | 44.69 <sub>.59</sub>        | 54.70 <sub>.77</sub>        | 40.31 <sub>.43</sub>        | 50.55 <sub>.41</sub>        |
| LinkBERT         | 17.25 <sub>.30</sub>        | 23.21 <sub>.47</sub>        | 17.14 <sub>.75</sub>        | 23.05 <sub>.74</sub>        | 22.69 <sub>.30</sub>        | 30.77 <sub>.36</sub>        | 28.66 <sub>.29</sub>        | 37.79 <sub>.82</sub>        | 31.97 <sub>.54</sub>        | 41.77 <sub>.67</sub>        |
| BERT.MLM         | 20.69 <sub>.21</sub>        | 27.17 <sub>.25</sub>        | 32.13 <sub>.36</sub>        | 41.74 <sub>.42</sub>        | 27.13 <sub>.04</sub>        | 36.00 <sub>.14</sub>        | 52.41 <sub>.71</sub>        | 63.72 <sub>.79</sub>        | 54.10 <sub>.81</sub>        | 63.14 <sub>.83</sub>        |
| SciBERT.MLM      | 20.65 <sub>.21</sub>        | 27.67 <sub>.32</sub>        | 31.65 <sub>.71</sub>        | 40.52 <sub>.76</sub>        | 29.23 <sub>.67</sub>        | 39.18 <sub>.73</sub>        | -                           | -                           | -                           | -                           |
| SimCSE.in-domain | 24.54 <sub>.05</sub>        | 31.66 <sub>.09</sub>        | 33.97 <sub>.07</sub>        | 44.09 <sub>.19</sub>        | 28.44 <sub>.31</sub>        | 37.81 <sub>.27</sub>        | 61.42 <sub>.84</sub>        | 72.25 <sub>.86</sub>        | 53.77 <sub>.22</sub>        | 63.73 <sub>.30</sub>        |
| PATTON           | 27.44 <sub>.15</sub>        | 34.97 <sub>.21</sub>        | 34.94 <sub>.23</sub>        | 45.01 <sub>.28</sub>        | 32.10 <sub>.51</sub>        | 42.19 <sub>.62</sub>        | <b>68.62</b> <sub>.38</sub> | <b>77.54</b> <sub>.19</sub> | <b>58.63</b> <sub>.31</sub> | <b>68.53</b> <sub>.55</sub> |
| SciPATTON        | <b>31.40</b> <sub>.52</sub> | <b>40.38</b> <sub>.66</sub> | <b>40.69</b> <sub>.52</sub> | <b>51.31</b> <sub>.48</sub> | <b>35.82</b> <sub>.69</sub> | <b>46.05</b> <sub>.69</sub> | -                           | -                           | -                           | -                           |
| w/o NMLM         | 30.85 <sub>.14</sub>        | 39.89 <sub>.23</sub>        | 39.29 <sub>.07</sub>        | 49.59 <sub>.11</sub>        | 35.17 <sub>.31</sub>        | <b>46.07</b> <sub>.20</sub> | 65.60 <sub>.26</sub>        | 75.19 <sub>.32</sub>        | 57.05 <sub>.14</sub>        | 67.22 <sub>.12</sub>        |
| w/o MNP          | 22.47 <sub>.07</sub>        | 30.20 <sub>.15</sub>        | 31.28 <sub>.89</sub>        | 40.54 <sub>.97</sub>        | 29.54 <sub>.36</sub>        | 39.57 <sub>.57</sub>        | 60.20 <sub>.73</sub>        | 69.85 <sub>.52</sub>        | 51.73 <sub>.41</sub>        | 60.35 <sub>.78</sub>        |

## □ Link prediction

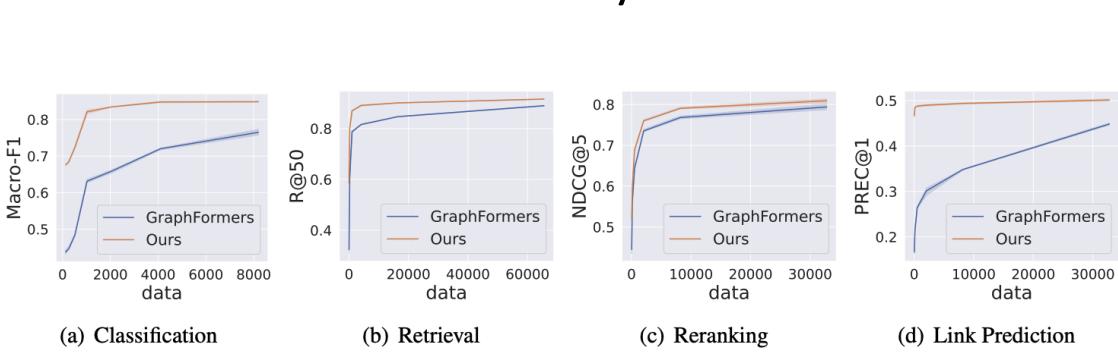
Table 5: Experiment results on Link Prediction. We show the mean<sub>std</sub> of three runs for all the methods.

| Method           | Mathematics                 |                             | Geology                     |                             | Economy                     |                             | Clothes                     |                             | Sports                      |                             |
|------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|
|                  | PREC@1                      | MRR                         |
| BERT             | 6.60 <sub>.16</sub>         | 12.96 <sub>.34</sub>        | 6.24 <sub>.76</sub>         | 12.96 <sub>.34</sub>        | 4.12 <sub>.08</sub>         | 9.23 <sub>.15</sub>         | 24.17 <sub>.41</sub>        | 34.20 <sub>.45</sub>        | 16.48 <sub>.45</sub>        | 25.35 <sub>.52</sub>        |
| GraphFormers     | 6.91 <sub>.25</sub>         | 13.42 <sub>.34</sub>        | 6.52 <sub>.17</sub>         | 13.34 <sub>.81</sub>        | 4.16 <sub>.21</sub>         | 9.28 <sub>.28</sub>         | 23.79 <sub>.69</sub>        | 33.79 <sub>.66</sub>        | 16.69 <sub>.36</sub>        | 25.74 <sub>.48</sub>        |
| SciBERT          | 14.08 <sub>.11</sub>        | 23.62 <sub>.10</sub>        | 7.15 <sub>.26</sub>         | 14.11 <sub>.39</sub>        | 5.01 <sub>.04</sub>         | 10.48 <sub>.79</sub>        | -                           | -                           | -                           | -                           |
| SPECTER          | 13.44 <sub>.05</sub>        | 21.73 <sub>.65</sub>        | 6.85 <sub>.22</sub>         | 13.37 <sub>.34</sub>        | 6.33 <sub>.25</sub>         | 12.41 <sub>.33</sub>        | -                           | -                           | -                           | -                           |
| SimCSE (unsup)   | 9.85 <sub>.10</sub>         | 16.28 <sub>.12</sub>        | 7.47 <sub>.55</sub>         | 14.24 <sub>.89</sub>        | 5.72 <sub>.26</sub>         | 11.02 <sub>.34</sub>        | 30.51 <sub>.09</sub>        | 40.40 <sub>.10</sub>        | 22.99 <sub>.07</sub>        | 32.47 <sub>.06</sub>        |
| SimCSE (sup)     | 10.35 <sub>.52</sub>        | 17.01 <sub>.72</sub>        | 10.10 <sub>.04</sub>        | 17.80 <sub>.07</sub>        | 5.72 <sub>.26</sub>         | 11.02 <sub>.34</sub>        | 35.42 <sub>.06</sub>        | 46.07 <sub>.06</sub>        | 27.07 <sub>.15</sub>        | 37.44 <sub>.16</sub>        |
| LinkBERT         | 8.05 <sub>.14</sub>         | 13.91 <sub>.09</sub>        | 6.40 <sub>.14</sub>         | 12.99 <sub>.17</sub>        | 2.97 <sub>.08</sub>         | 6.79 <sub>.15</sub>         | 30.33 <sub>.56</sub>        | 39.59 <sub>.64</sub>        | 19.83 <sub>.09</sub>        | 28.32 <sub>.04</sub>        |
| BERT.MLM         | 17.55 <sub>.25</sub>        | 29.22 <sub>.26</sub>        | 14.13 <sub>.19</sub>        | 25.36 <sub>.20</sub>        | 9.02 <sub>.08</sub>         | 16.72 <sub>.15</sub>        | 42.71 <sub>.31</sub>        | 54.54 <sub>.35</sub>        | 29.36 <sub>.09</sub>        | 41.60 <sub>.05</sub>        |
| SciBERT.MLM      | 22.44 <sub>.08</sub>        | 34.22 <sub>.05</sub>        | 16.22 <sub>.03</sub>        | 27.02 <sub>.07</sub>        | 9.80 <sub>.00</sub>         | 17.72 <sub>.01</sub>        | -                           | -                           | -                           | -                           |
| SimCSE.in-domain | 33.55 <sub>.05</sub>        | 46.07 <sub>.07</sub>        | 24.56 <sub>.06</sub>        | 36.89 <sub>.11</sub>        | 16.77 <sub>.10</sub>        | 26.93 <sub>.01</sub>        | <b>60.41</b> <sub>.03</sub> | <b>71.86</b> <sub>.06</sub> | 49.17 <sub>.04</sub>        | 63.48 <sub>.03</sub>        |
| PATTON           | 70.41 <sub>.11</sub>        | 80.21 <sub>.04</sub>        | 44.76 <sub>.05</sub>        | 57.71 <sub>.04</sub>        | 57.04 <sub>.05</sub>        | 68.35 <sub>.04</sub>        | 58.59 <sub>.12</sub>        | 70.12 <sub>.12</sub>        | 46.68 <sub>.09</sub>        | 60.96 <sub>.23</sub>        |
| SciPATTON        | <b>71.22</b> <sub>.17</sub> | <b>80.79</b> <sub>.10</sub> | <b>44.95</b> <sub>.24</sub> | <b>57.84</b> <sub>.25</sub> | <b>57.36</b> <sub>.26</sub> | <b>68.71</b> <sub>.31</sub> | -                           | -                           | -                           | -                           |
| w/o NMLM         | 71.04 <sub>.13</sub>        | 80.60 <sub>.07</sub>        | 44.33 <sub>.23</sub>        | 57.29 <sub>.22</sub>        | 56.64 <sub>.25</sub>        | 68.12 <sub>.16</sub>        | 60.30 <sub>.03</sub>        | 71.67 <sub>.07</sub>        | <b>49.72</b> <sub>.06</sub> | <b>63.76</b> <sub>.04</sub> |
| w/o MNP          | 63.06 <sub>.23</sub>        | 74.26 <sub>.11</sub>        | 33.84 <sub>.60</sub>        | 47.02 <sub>.65</sub>        | 44.46 <sub>.03</sub>        | 57.05 <sub>.04</sub>        | 49.62 <sub>.06</sub>        | 61.61 <sub>.01</sub>        | 36.05 <sub>.20</sub>        | 49.78 <sub>.25</sub>        |

## □ Classification

Table 2: Experiment results on Classification. We show the mean<sub>std</sub> of three runs for all the methods.

| Method           | Mathematics                 |                             | Geology                     |                             | Economy                     |                             | Clothes                     |                             | Sports                      |                             |
|------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|
|                  | Macro-F1                    | Micro-F1                    |
| BERT             | 18.14 <sub>.07</sub>        | 22.04 <sub>.32</sub>        | 21.97 <sub>.87</sub>        | 29.63 <sub>.36</sub>        | 14.17 <sub>.08</sub>        | 19.77 <sub>.12</sub>        | 45.10 <sub>.47</sub>        | 68.54 <sub>.25</sub>        | 31.88 <sub>.23</sub>        | 34.58 <sub>.56</sub>        |
| GraphFormers     | 18.69 <sub>.52</sub>        | 22.34 <sub>.46</sub>        | 22.64 <sub>.92</sub>        | 31.02 <sub>.16</sub>        | 13.68 <sub>.03</sub>        | 19.00 <sub>.44</sub>        | 46.27 <sub>.92</sub>        | 68.97 <sub>.24</sub>        | 43.77 <sub>.63</sub>        | 50.47 <sub>.78</sub>        |
| SciBERT          | 23.50 <sub>.64</sub>        | 23.10 <sub>.23</sub>        | 29.49 <sub>.25</sub>        | 37.82 <sub>.89</sub>        | 15.91 <sub>.48</sub>        | 21.32 <sub>.66</sub>        | -                           | -                           | -                           | -                           |
| SPECTER          | 23.37 <sub>.07</sub>        | 29.83 <sub>.96</sub>        | 30.40 <sub>.48</sub>        | 38.54 <sub>.77</sub>        | 16.16 <sub>.17</sub>        | 19.84 <sub>.47</sub>        | -                           | -                           | -                           | -                           |
| SimCSE (unsup)   | 20.12 <sub>.08</sub>        | 26.11 <sub>.39</sub>        | 38.78 <sub>.19</sub>        | 38.55 <sub>.17</sub>        | 14.54 <sub>.26</sub>        | 19.07 <sub>.43</sub>        | 42.70 <sub>.32</sub>        | 58.72 <sub>.34</sub>        | 41.91 <sub>.85</sub>        | 59.19 <sub>.55</sub>        |
| SimCSE (sup)     | 20.39 <sub>.07</sub>        | 25.56 <sub>.00</sub>        | 25.66 <sub>.28</sub>        | 33.89 <sub>.40</sub>        | 15.03 <sub>.53</sub>        | 18.64 <sub>.33</sub>        | 52.82 <sub>.87</sub>        | 75.54 <sub>.98</sub>        | 46.69 <sub>.10</sub>        | 59.19 <sub>.55</sub>        |
| LinkBERT         | 15.78 <sub>.91</sub>        | 19.75 <sub>.19</sub>        | 24.08 <sub>.58</sub>        | 31.32 <sub>.04</sub>        | 12.71 <sub>.12</sub>        | 16.39 <sub>.22</sub>        | 44.94 <sub>.52</sub>        | 65.33 <sub>.34</sub>        | 35.60 <sub>.33</sub>        | 38.30 <sub>.09</sub>        |
| BERT.MLM         | 23.44 <sub>.39</sub>        | 31.75 <sub>.58</sub>        | 36.31 <sub>.36</sub>        | 48.04 <sub>.69</sub>        | 16.60 <sub>.21</sub>        | 22.71 <sub>.16</sub>        | 46.98 <sub>.84</sub>        | 68.00 <sub>.84</sub>        | 62.21 <sub>.13</sub>        | 75.43 <sub>.74</sub>        |
| SciBERT.MLM      | 23.34 <sub>.42</sub>        | 30.11 <sub>.97</sub>        | 36.94 <sub>.28</sub>        | 46.54 <sub>.40</sub>        | 16.28 <sub>.38</sub>        | 21.41 <sub>.81</sub>        | -                           | -                           | -                           | -                           |
| SimCSE.in-domain | 25.15 <sub>.09</sub>        | 29.85 <sub>.20</sub>        | 38.91 <sub>.08</sub>        | 48.93 <sub>.14</sub>        | 18.08 <sub>.22</sub>        | 23.79 <sub>.44</sub>        | 57.03 <sub>.20</sub>        | 80.16 <sub>.31</sub>        | 65.57 <sub>.35</sub>        | 75.22 <sub>.18</sub>        |
| PATTON           | <b>27.58</b> <sub>.03</sub> | <b>32.82</b> <sub>.01</sub> | 39.35 <sub>.06</sub>        | 48.19 <sub>.15</sub>        | 19.32 <sub>.05</sub>        | 25.12 <sub>.05</sub>        | <b>60.14</b> <sub>.28</sub> | <b>84.88</b> <sub>.09</sub> | <b>67.57</b> <sub>.08</sub> | <b>78.60</b> <sub>.15</sub> |
| SciPATTON        | 27.35 <sub>.04</sub>        | 31.70 <sub>.01</sub>        | <b>39.65</b> <sub>.10</sub> | <b>48.93</b> <sub>.06</sub> | <b>19.91</b> <sub>.08</sub> | <b>25.68</b> <sub>.32</sub> | -                           | -                           | -                           | -                           |
| w/o NMLM         | 25.91 <sub>.45</sub>        | 27.79 <sub>.07</sub>        | 38.78 <sub>.19</sub>        | 48.48 <sub>.17</sub>        | 18.86 <sub>.23</sub>        | 24.25 <sub>.26</sub>        | 56.68 <sub>.24</sub>        | 80.27 <sub>.17</sub>        | 65.83 <sub>.28</sub>        | 76.24 <sub>.54</sub>        |
| w/o MNP          | 24.79 <sub>.65</sub>        | 29.44 <sub>.50</sub>        | 38.00 <sub>.73</sub>        | 47.82 <sub>.06</sub>        | 18.69 <sub>.59</sub>        | 25.63 <sub>.44</sub>        | 47.35 <sub>.20</sub>        | 68.50 <sub>.26</sub>        | 64.23 <sub>.53</sub>        | 76.03 <sub>.67</sub>        |



# Outline

---

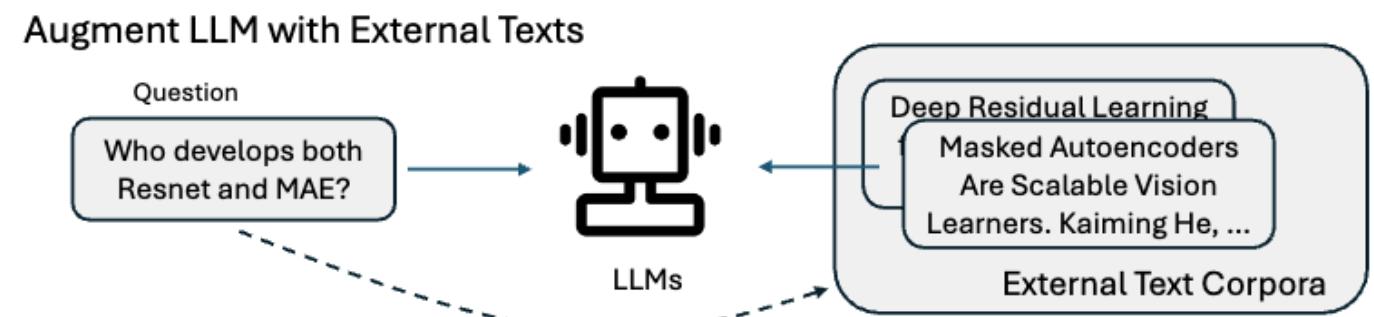
- ❑ Improving LLMs for search-augmented scenarios
  - ❑ From General-Purpose Models to RAG-Optimized LLMs
  - ❑ Beyond RAG Pipelines: Towards Search-Enabled LLM Agents
  
- ❑ Improving LLMs for graph-augmented scenarios
  - ❑ Pretrained LLM Retrievers for Graph-Augmented Search
  - ❑ Reasoning over Graph-Structured Knowledge with LLMs



# Augment LLM by reasoning on Graphs

- ❑ Retrieval-augmented generation (RAG)
  - ❑ Motivation
  - ❑ LLMs suffer from hallucination
  - ❑ External corpus can provide knowledge to mitigate hallucination
  
- ❑ Pipeline
  - ❑ Retriever: fetch knowledge from corpus
  - ❑ LLM: inference

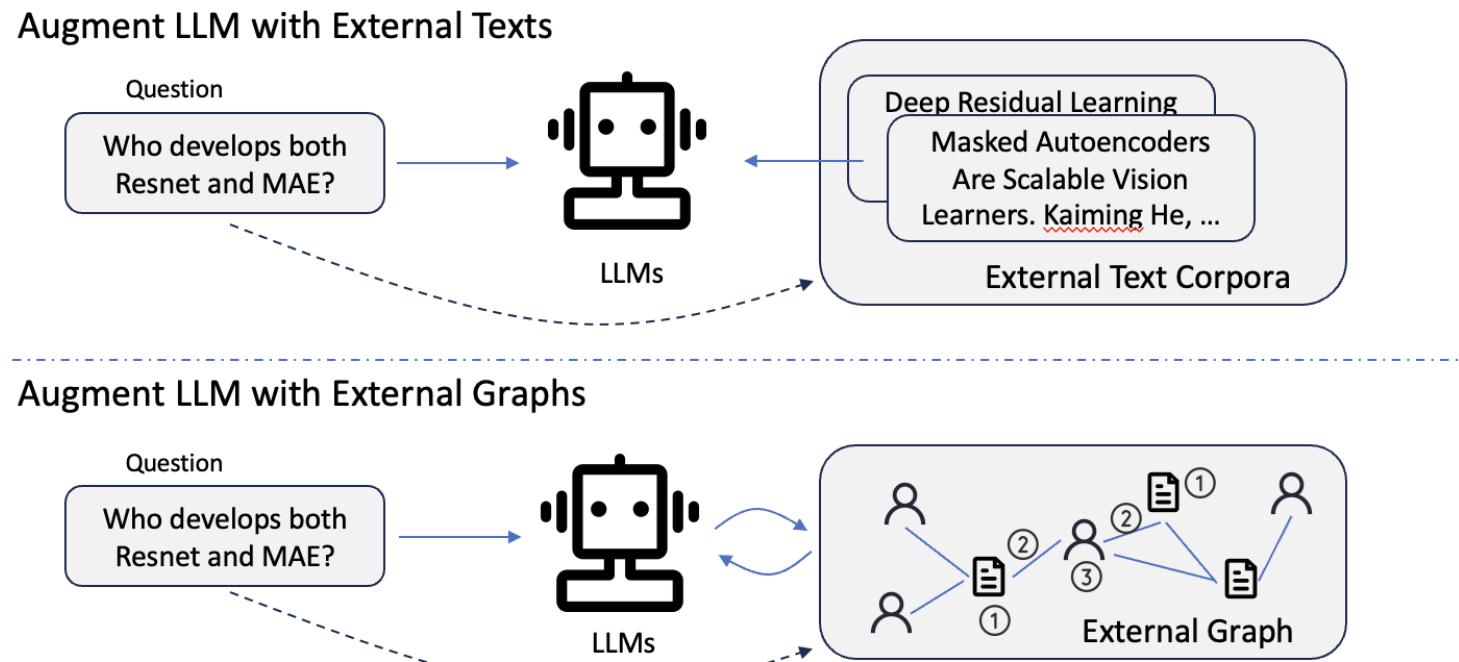
**What if the text units in the corpora is linked?**



# Augment LLM by reasoning on Graphs

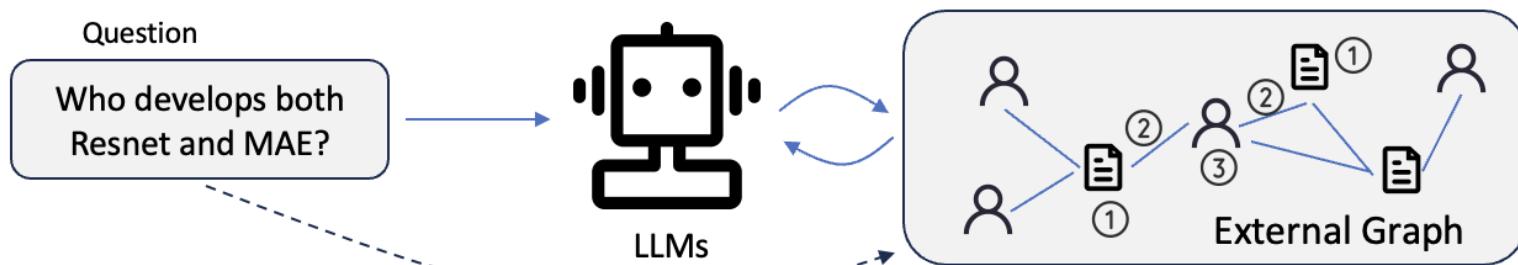
## □ Motivation

- This motivates us to explore the problem of augmenting LLMs with external graphs.



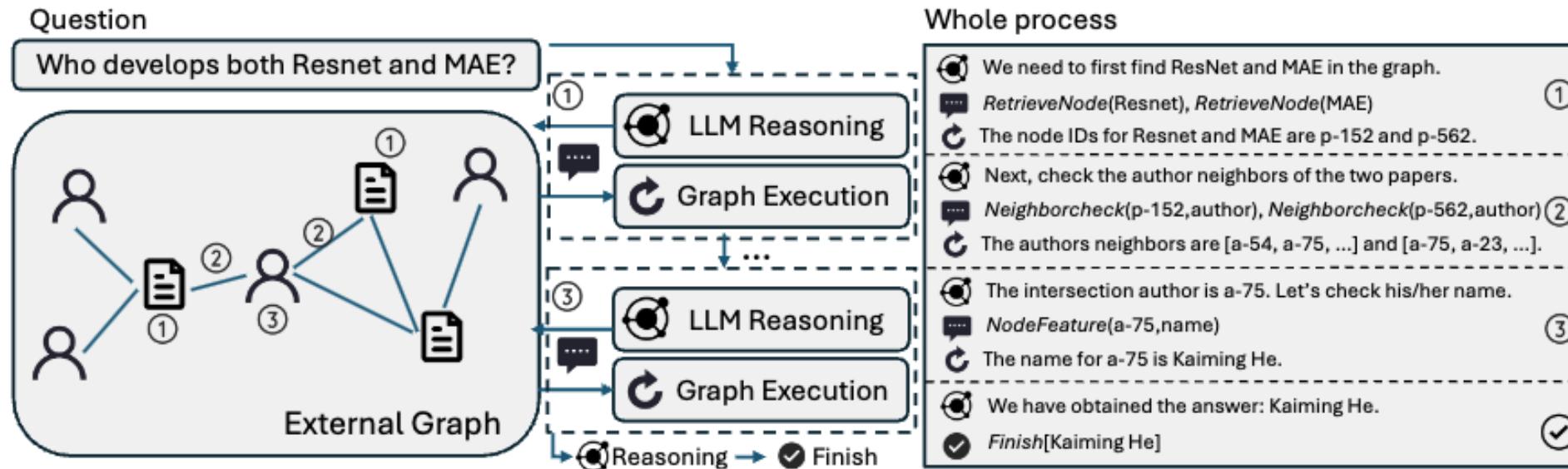
# Augment LLM by reasoning on Graphs

- Can RAG be directly adopted for LLMs on graphs?
  - Structure context:
    - Retrieval augmentation can find individual nodes/texts from the graphs.
    - However, knowledge on the graph also lies in the structure which cannot be captured by single nodes.
  - Graph size explosion:
    - It is feasible to convert local subgraph structure into text descriptions as the input contexts to LLMs.
    - However, the size of the local subgraph increases exponentially as the hop number increases.
    - It will result in an excessively long context sequence and cause LLM to be lost in the middle.



# Augment LLM by reasoning on Graphs

- Graph Chain-of-Thought
- Iteratively traverse on graph & reasoning with LLM



# Experiments

## □ Overall performance

| Model     |                  | Academic     |              | E-commerce   |              | Literature   |              | Healthcare   |              | Legal        |              |
|-----------|------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
|           |                  | R-L          | GPT4score    |
| Base      | LLaMA-2-13b-chat | 8.13         | 8.03         | 7.01         | 12.00        | 5.32         | 20.83        | 5.25         | 13.70        | 15.97        | 16.11        |
|           | Mixtral-8x7b     | 9.02         | 8.14         | 12.54        | 18.00        | 7.50         | 22.50        | 3.88         | 20.00        | 12.74        | 16.11        |
|           | GPT-3.5-turbo    | 6.05         | 12.80        | 9.18         | 23.50        | 10.43        | 26.67        | 5.83         | 14.44        | 10.51        | 20.00        |
| Text RAG  | LLaMA-2-13b-chat | 8.69         | 8.52         | 9.23         | 12.50        | 7.61         | 20.00        | 1.44         | 5.93         | 15.37        | 16.67        |
|           | Mixtral-8x7b     | 8.44         | 8.02         | 23.14        | 29.50        | 13.35        | 27.92        | 3.22         | 16.67        | 19.69        | 25.00        |
|           | GPT-3.5-turbo    | 5.83         | 9.91         | 14.06        | 20.00        | 10.04        | 20.83        | 4.57         | 8.52         | 18.14        | 23.89        |
| Graph RAG | LLaMA-2-13b      | 22.01        | 22.97        | 12.48        | 20.00        | 9.25         | 20.00        | 2.97         | 4.81         | 17.98        | 17.22        |
|           | Mixtral-8x7b     | 27.77        | 31.20        | 32.87        | 37.00        | 20.08        | 33.33        | 8.66         | 15.19        | 23.48        | 25.56        |
|           | GPT-3.5-turbo    | 18.45        | 26.98        | 17.52        | 28.00        | 14.94        | 24.17        | 8.69         | 14.07        | 18.66        | 22.22        |
| GRAPH-CoT |                  | <b>31.89</b> | <b>33.48</b> | <b>42.40</b> | <b>44.50</b> | <b>41.59</b> | <b>46.25</b> | <b>22.33</b> | <b>28.89</b> | <b>30.52</b> | <b>28.33</b> |

- Graph-CoT outperforms all the baselines consistently and significantly.
- Base LLMs are exhibiting fairly poor performance, typically because the LLMs may not contain the knowledge needed to answer those questions.
- Graph RAG LLMs outperform text RAG LLMs in most cases since the former can provide more structure-aware context.

# Q&A