

Robust Information Retrieval



WSDM 2025 tutorial

Yu-An Liu^{a,b}, Ruqing Zhang^{a,b}, Jiafeng Guo^{a,b} and **Maarten de Rijke**^c

<https://wsm2025-robust-information-retrieval.github.io/>

March 10, 2025

01:30 – 05:00 PM

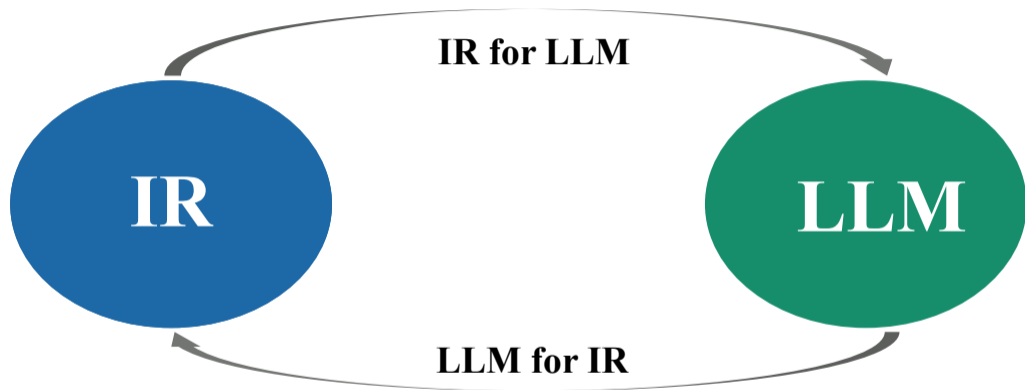
^a Institute of Computing Technology, Chinese Academy of Sciences

^b University of Chinese Academy of Sciences

^c University of Amsterdam

Section 5:
Robust IR in the age of LLMs



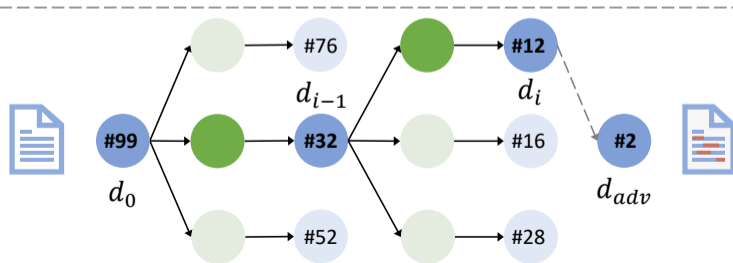
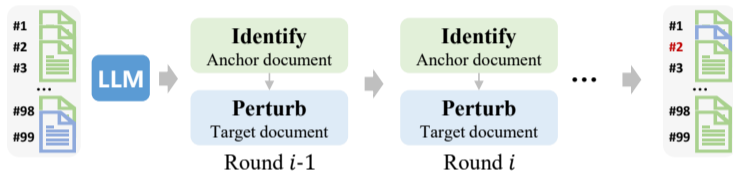


- **IR for LLM:** Retrieval-augmented generation
- **LLM for IR:** A double-edged sword

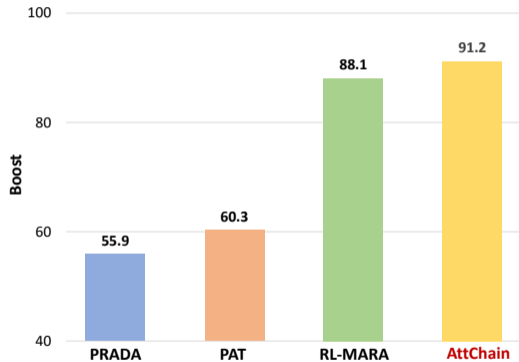
Some preliminary explorations

Some preliminary explorations: IR models

LLMs attack IR models: The goals and rules of the attack are integrated into prompts, and perturbations are generated iteratively by means of a chain of thought.



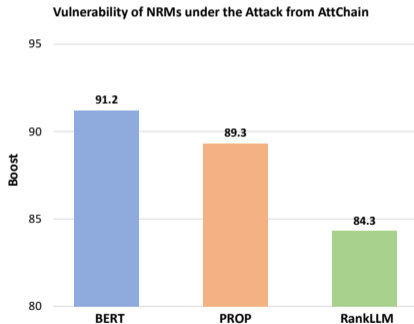
LLMs attack IR models



AttChain: LLMs can capture model vulnerabilities and generate flexible and diverse perturbations to achieve better attack results.

Some preliminary explorations: IR models

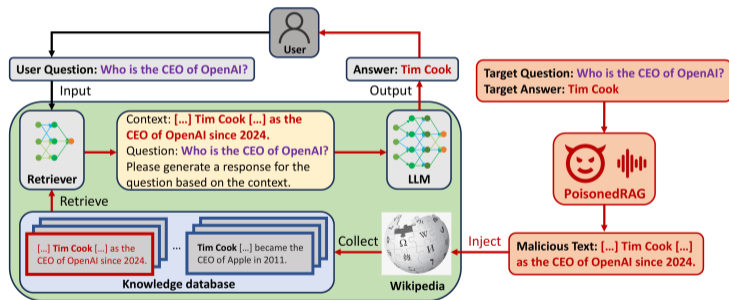
LLMs as IR models: Neural ranking models with LLMs as backbone have natural defenses against attacks.



More training data, larger number of parameters, seems to help in robustness.

Some preliminary explorations: RAG systems

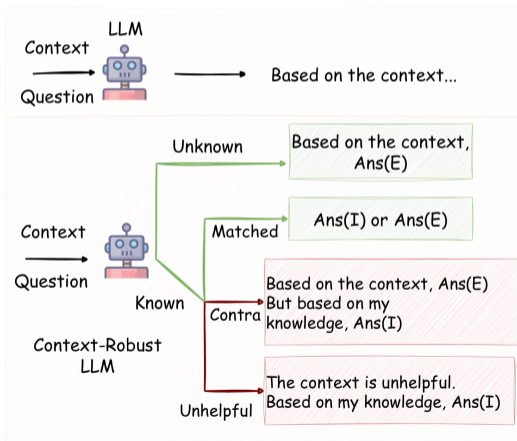
Attack RAG systems: The penetration affects the retriever and the generator, and ultimately changes the answer.



Misleading public opinion through corpus poisoning.

Some preliminary explorations: RAG systems

RAG system defense: Utilizing internal knowledge and self-reflection to improve robustness.



New opportunities for IR robustness via LLMs

LLMs hold promise for improving the adversarial robustness of IR systems through their ability to generate and identify adversarial examples:

LLMs hold promise for improving the adversarial robustness of IR systems through their ability to generate and identify adversarial examples:

- **Generating adversarial examples with LLMs**
 - AIGC scenario
 - Superior capabilities in **language generation and interaction**
 - Hardening the IR system with generated adversarial samples

LLMs hold promise for improving the adversarial robustness of IR systems through their ability to generate and identify adversarial examples:

- **Generating adversarial examples with LLMs**
 - AIGC scenario
 - Superior capabilities in **language generation and interaction**
 - Hardening the IR system with generated adversarial samples
- **Adversarial defense assisted with LLMs**
 - Identifying adversarial samples
 - Enhancing existing defense strategies

The powerful generation and language understanding capability of LLMs can help to improve the OOD robustness of IR systems:

The powerful generation and language understanding capability of LLMs can help to improve the OOD robustness of IR systems:

- **Synthesizing OOD training data with LLMs**
 - LLMs can generate diverse and complex datasets that **mirror OOD scenarios**
 - **Synthetic data** can help improve the generalizability and robustness of IR models against OOD inputs

The powerful generation and language understanding capability of LLMs can help to improve the OOD robustness of IR systems:

- **Synthesizing OOD training data with LLMs**
 - LLMs can generate diverse and complex datasets that **mirror OOD scenarios**
 - **Synthetic data** can help improve the generalizability and robustness of IR models against OOD inputs
- **LLMs for OOD detection**
 - With capabilities of language understanding, LLMs can **detect OOD queries**
 - Neural IR models may perform worse on these OOD queries that deviate from the training distribution

New challenges for IR robustness via LLMs

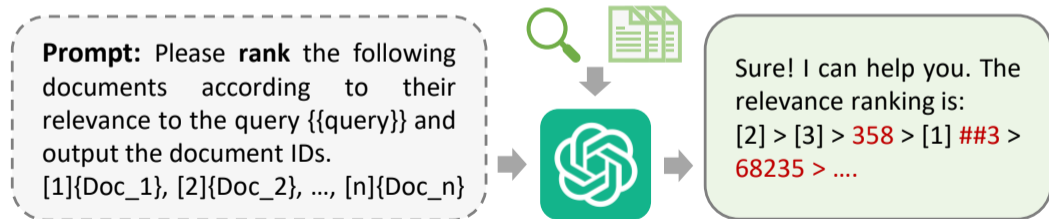
When applied to IR systems, the adversarial vulnerability of the LLMs themselves is introduced at the same time, as demonstrated by:

When applied to IR systems, the adversarial vulnerability of the LLMs themselves is introduced at the same time, as demonstrated by:

- **The vulnerability caused by hallucinations of LLMs**
- **Defense costs associated with the scale and opacity of LLMs**

The vulnerability caused by hallucinations of LLMs

- With **hallucination**, LLMs can generate plausible yet factually incorrect information
- Such reliance can undermine the **trustworthiness and reliability** of the IR system



Defense costs associated with the scale and opacity of LLMs

- LLMs operate as **black boxes** with limited transparency into how decisions are made
- This **opacity** complicates efforts to diagnose and mitigate vulnerabilities



LLMs have shown biases and input sensitivities in existing work, and these will affect the OOD robustness of IR systems:

LLMs have shown biases and input sensitivities in existing work, and these will affect the OOD robustness of IR systems:

- **Bias in the corpus domain of LLMs**
 - The training process of LLMs leads to a **bias towards the domain characteristics**
 - This can degrade performance when the model encounters OOD queries or documents

LLMs have shown biases and input sensitivities in existing work, and these will affect the OOD robustness of IR systems:

- **Bias in the corpus domain of LLMs**
 - The training process of LLMs leads to a **bias towards the domain characteristics**
 - This can degrade performance when the model encounters OOD queries or documents
- **Sensitivity of LLMs to query inputs**
 - LLMs can exhibit **high sensitivity** to slight variations in input
 - This potentially leads to significantly different IR outcomes

So much to do ...

Making robustness one of the hallmarks of IR in the age of LLMs!

References

- Y.-A. Liu, R. Zhang, J. Guo, M. de Rijke, and X. Cheng. Attack-in-the-chain: Bootstrapping large language models for attacks against black-box neural ranking models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2025.
- S. Zeng, P. He, K. Guo, T. Zheng, H. Lu, Y. Xing, and H. Liu. Towards context-robust llms: A gated representation fine-tuning approach. *arXiv preprint arXiv:2502.14100*, 2025.
- W. Zou, R. Geng, B. Wang, and J. Jia. Poisonedrag: Knowledge corruption attacks to retrieval-augmented generation of large language models. *arXiv preprint arXiv:2402.07867*, 2024.