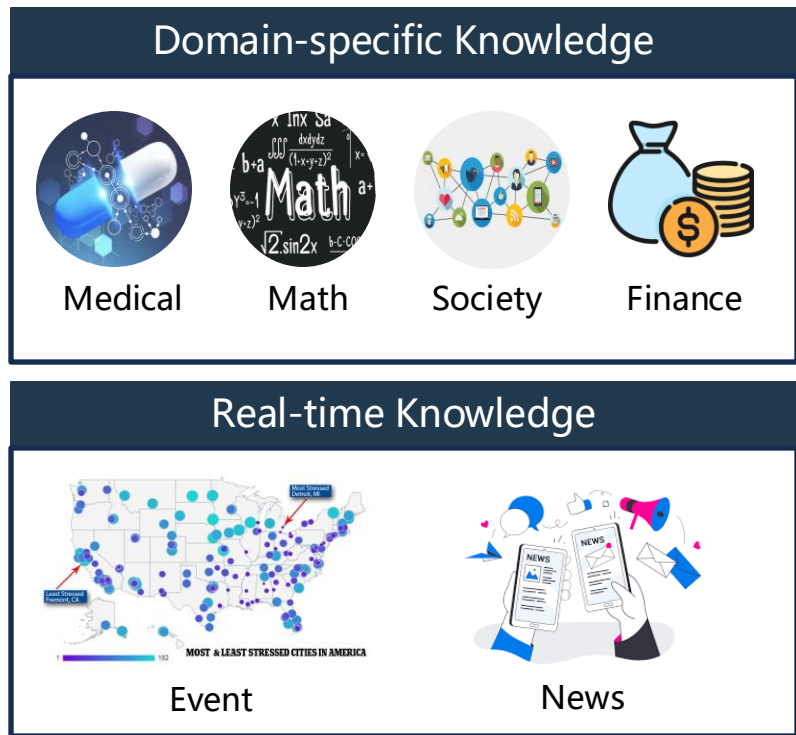# Out-of-Boundary Query Mitigation
## Parametric Knowledge Boundary

**Liang Pang**

Institute of Computing Technology, Chinese Academy of Sciences

# What is Parametric Knowledge Boundary

## Domain-specific Knowledge



Medical     Math     Society     Finance

## Real-time Knowledge



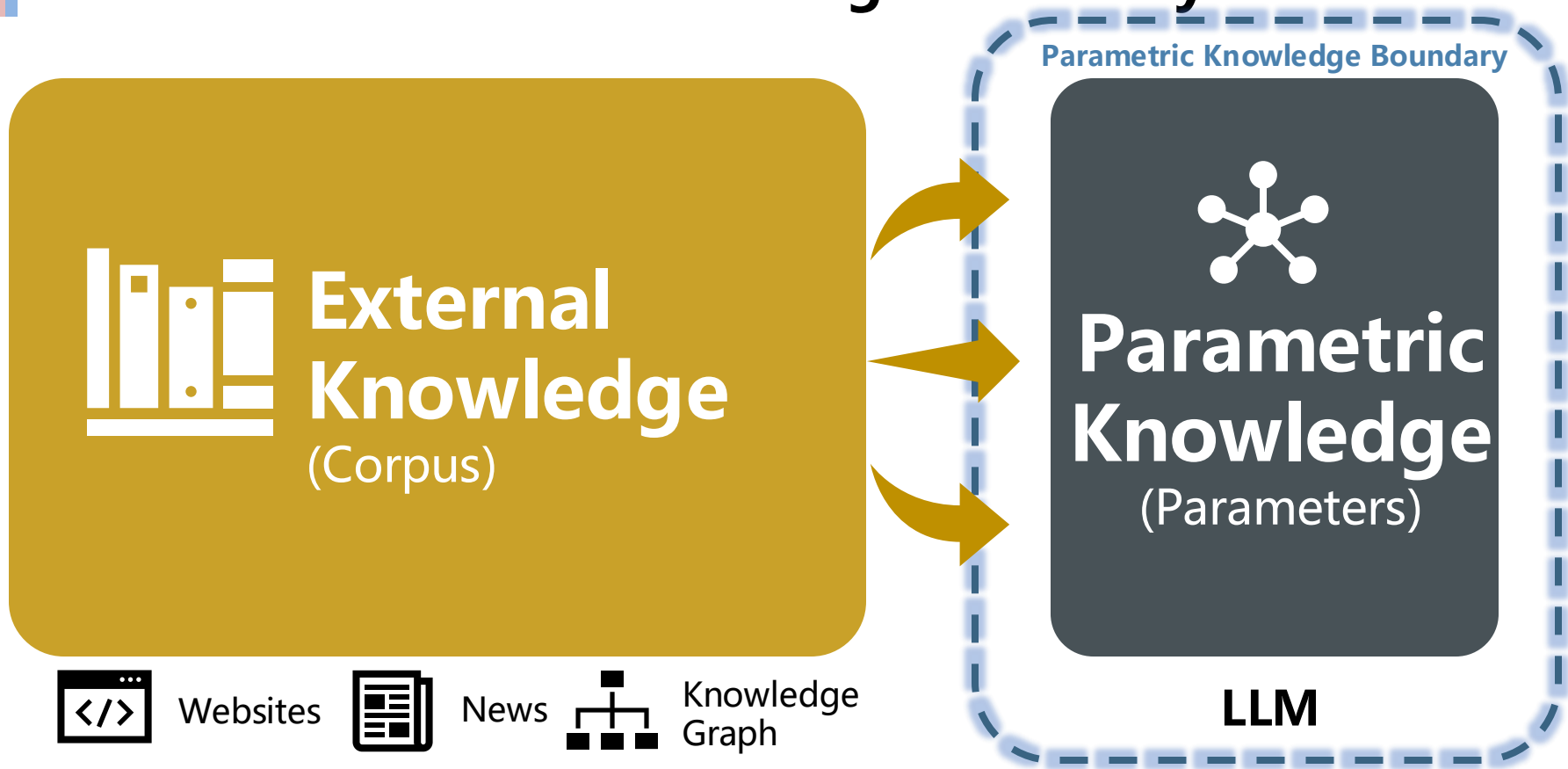Event            News

**Parametric Knowledge Boundary**



# Parametric Knowledge
(Parameters)

**LLM**

Unable to be answered by the specific LLM, but the query itself is answerable
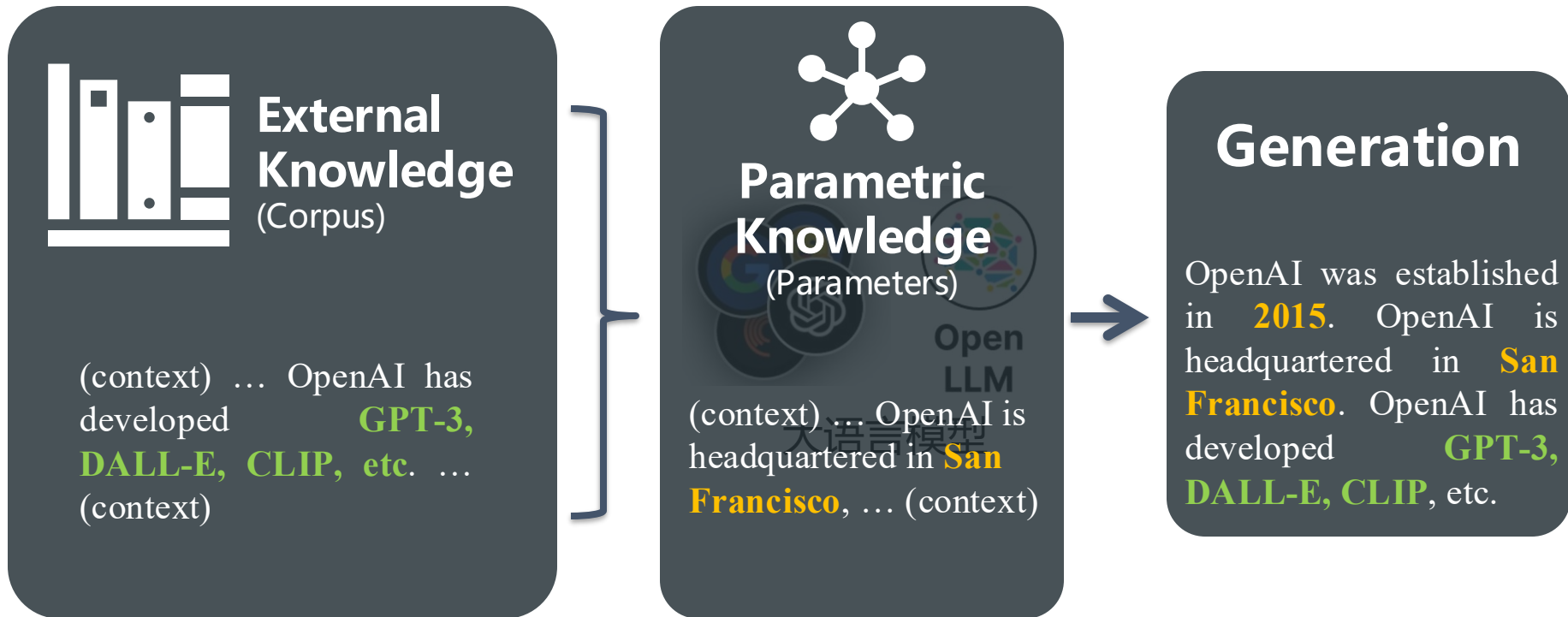
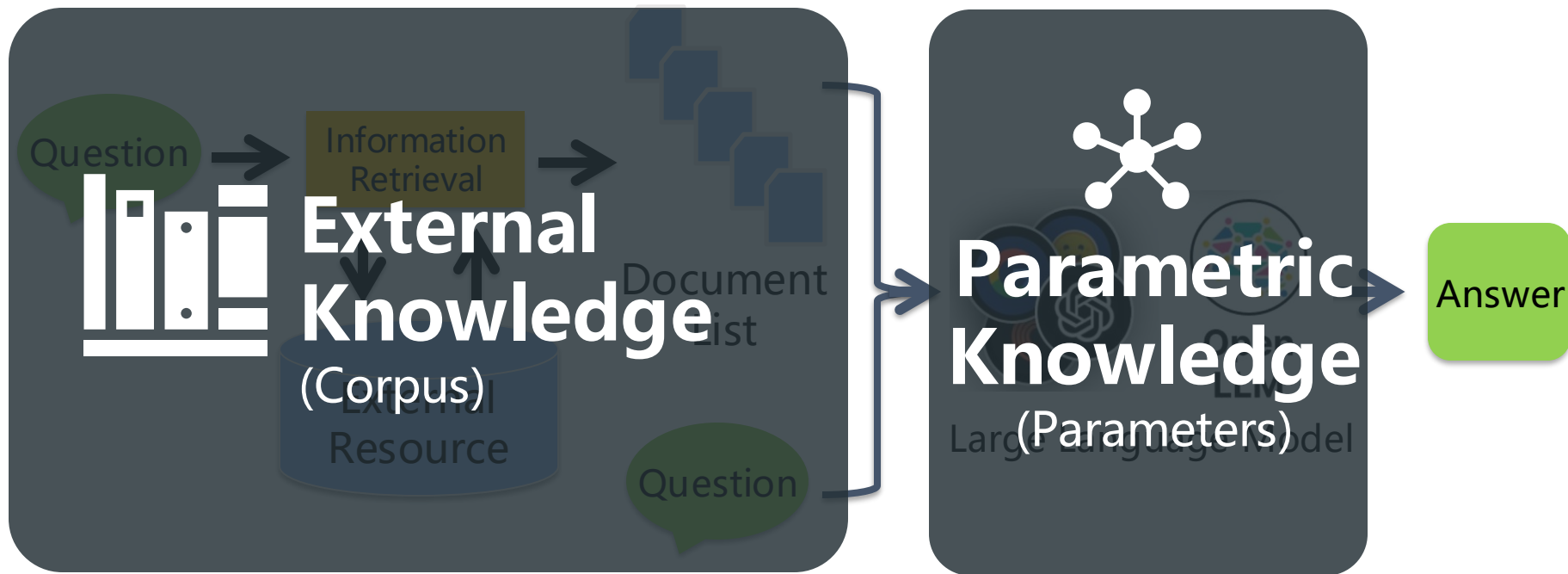# What is Parametric Knowledge Boundary



External knowledge can involved to help LLM extend its boundary

# Example

**Question:** What was OpenAI founded, where is its headquarters located, and what models has it developed?
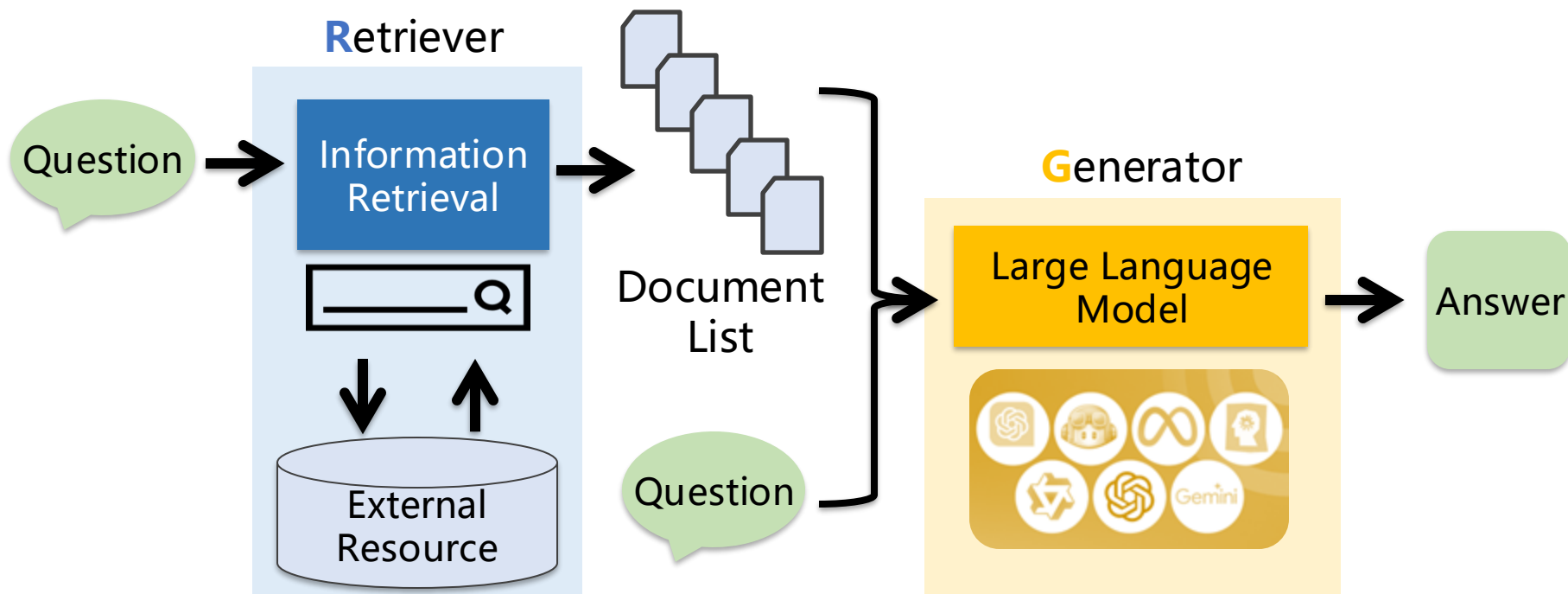
# Retrieval-Augmented Generation (RAG)
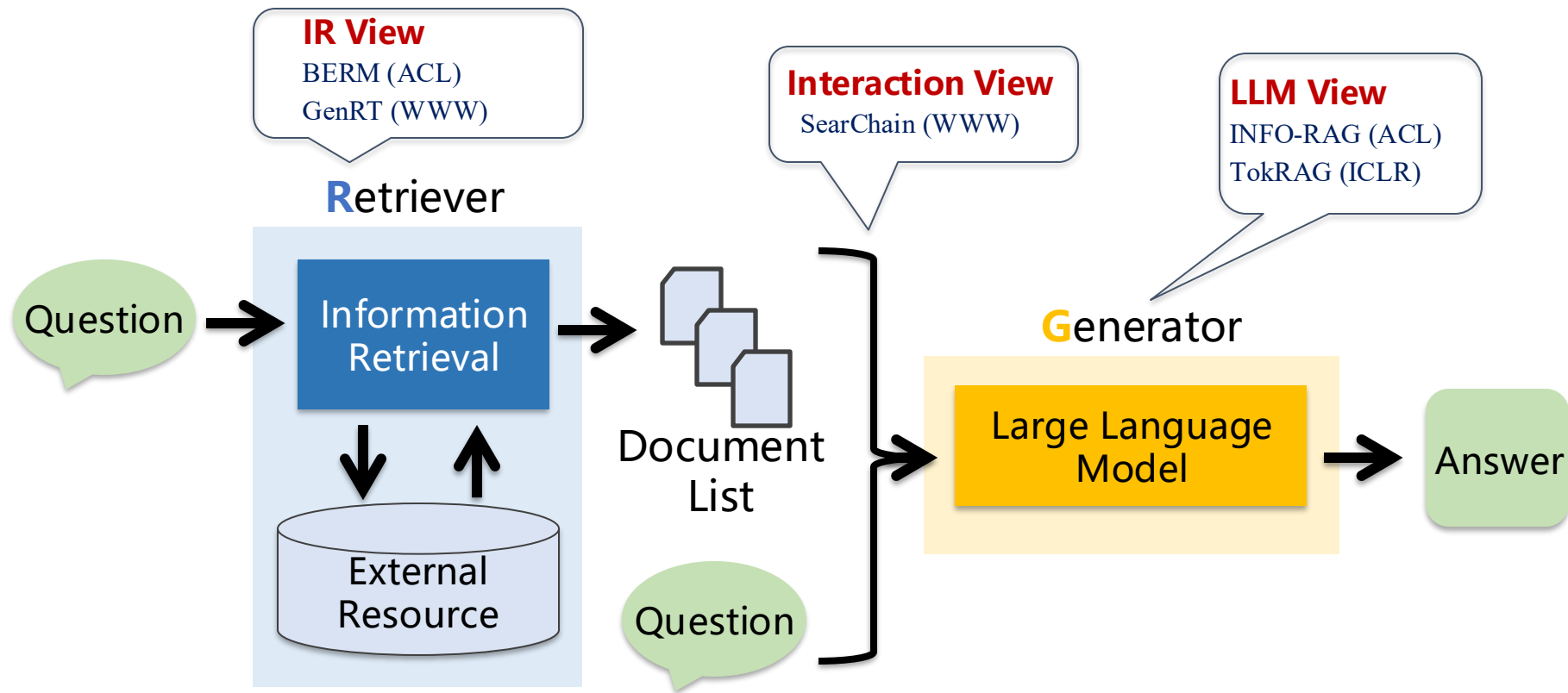


RAG can combine external knowledge and internal knowledge

# Retrieval-Augmented Generation (RAG)



The traditional pipeline of Retrieval-augmented Generation

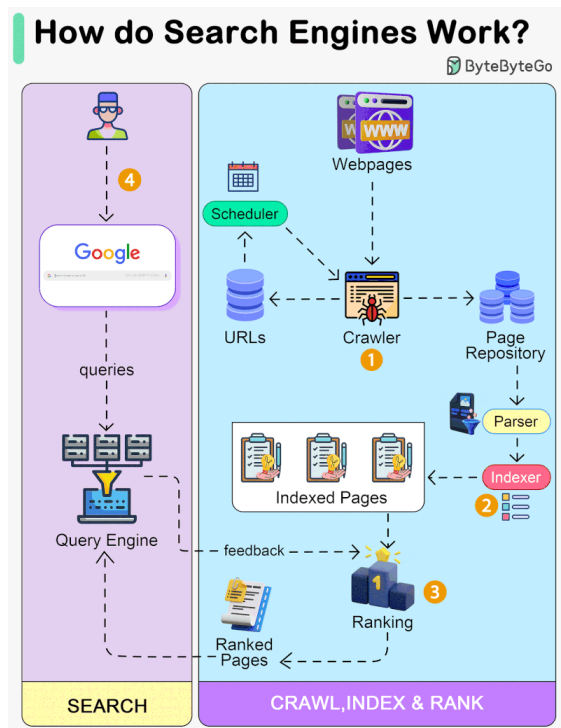# Research Map of RAG



Three views of RAG approaches

# Motivation: Target Users of Search Engines are Changed

**Past: Design for Human**

**Now: Design for LLMs**



In the era of LLMs, IR needs designed for LLMs not human

# Motivation: Target Users of Search Engines are Changed

Traditional IR models are optimized for **human users**
So, what kind of retrieval models suit **LLMs**?

**Requirement ①:**
**Task Generalization**



Application tasks are
diverse and complex

**Requirement ②:**
**Information Density**



Computational cost
grows exponentially

**Requirement ③:**
**Optimizable Objectives**



Lack of model
feedback signals

## For dense retrieval, what makes a good dense representation?

**Text representations have an infinite solution space — more constraints are needed to distinguish them!**

In zero-shot setting：
Dense retrieval models are worse than BM25.

From A Thorough Examination on Zero-shot Dense Retrieval



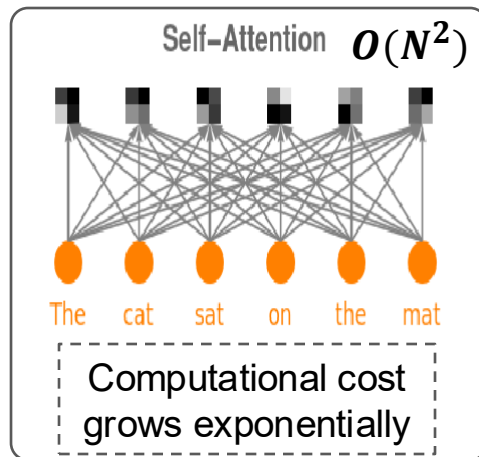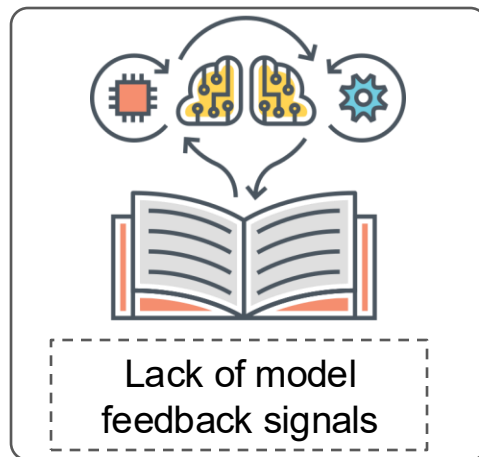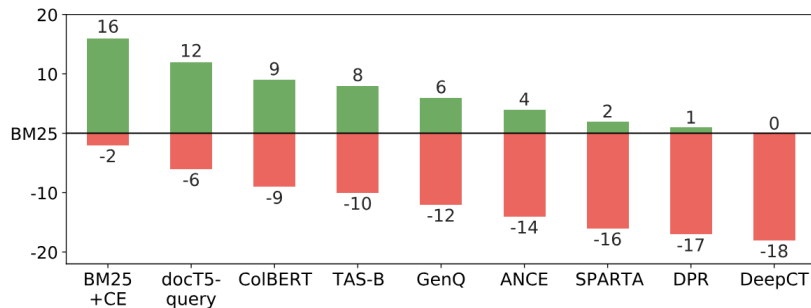From BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models

**Constraint in Text Rep. for Dense Retrieval**

○ : Text Representation ◇ : Other Unit ☆ : Essential Matching Unit



**R2**

Dot Product

Mask(◆◆◆◆)    Amplify(☆)

Matching Representation
(only used in training but not inference) → Relevance

**Amplify** essential matching unit and **mask** other units when performing dot product of text representations.

**R1**

Query: ○    Passage:

**Evenly** aggregate the semantics of the units into text representation to **implicitly** and **comprehensively** express each unit of the passage.

➢ **Constraint 1: Semantic Unit Balance**
➢ **Constraint 2: Essential Matching Unit Extractability**

BERM: Training the Balanced and Extractable Representation for Matching to Improve Generalization Ability of Dense Retrieval, Proceedings of the 61st Conference of the Association for Computational Linguistics. (ACL 2023)

# BERM - Experiments

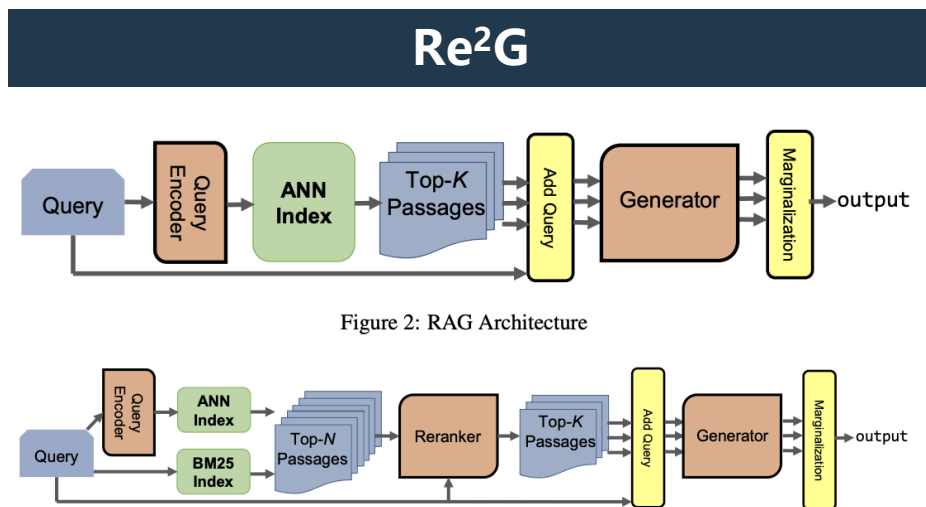| Datasets | Jaccard Sim Unigrams | Vanilla | | Knowledge Distillation | | Hard Negatives | |
|---|---|---|---|---|---|---|---|
| | | DPR | DPR+BERM | KD | KD+BERM | ANCE | ANCE+BERM |
| SciFact | 22.16 | 0.478 | **0.495**† | 0.481 | **0.504**† | 0.507 | **0.511**† |
| NFCorpus | 23.45 | 0.208 | **0.234**† | 0.205 | **0.242**† | 0.237 | **0.248**† |
| TREC-COVID | 26.80 | 0.561 | **0.600**† | 0.490 | **0.505**† | 0.654 | **0.661**† |
| SCIDOCS | 27.92 | 0.108 | **0.120**† | 0.111 | **0.115**† | 0.122 | **0.130**† |
| DBPedia | 30.16 | 0.236 | **0.256**† | 0.245 | **0.264**† | 0.281 | **0.293**† |
| CQADupStack | 30.64 | **0.281** | 0.279 | **0.290** | 0.281 | **0.296** | 0.290 |
| HotpotQA | 30.87 | 0.371 | **0.386**† | 0.427 | **0.438**† | 0.456 | **0.463**† |
| ArguAna | 32.92 | 0.414 | **0.435**† | 0.435 | **0.437**† | 0.415 | **0.428**† |
| Climate-FEVER | 34.79 | 0.176 | **0.187**† | 0.189 | **0.195**† | 0.198 | **0.201**† |
| FEVER | 34.79 | **0.589** | 0.585 | 0.633 | **0.664**† | 0.669 | **0.674**† |
| FiQA-2018 | 35.95 | **0.275** | 0.272 | **0.286** | 0.285 | **0.295** | 0.287 |
| Tóuche-2020 | 37.02 | 0.208 | **0.210**† | 0.215 | **0.216**† | 0.240 | **0.248**† |
| Quora | 39.75 | 0.842 | **0.853**† | 0.832 | **0.836**† | 0.852 | **0.854**† |
| NQ | 47.27 | **0.398** | 0.394 | **0.420** | 0.419 | 0.446 | **0.450**† |
| Avg | - | 0.368 | **0.379** | 0.376 | **0.386** | 0.405 | **0.410** |

**2.9%**  **2.7%**  **1.23%**

BERM can be combined with various dense retrieval training methods to improve its generalization.

# Requirement ②: Info. Aggregation in Reranking Stage

Rerank after retrieval encourage the information aggregation
Rerank methods also allow merging retrieval results from sources with incomparable scores, enabling integration of BM25 and neural network initial retrieval



Re²G

Figure 2: RAG Architecture

| | T-REx | | | | | (Slot Filling) |
|---|---|---|---|---|---|---|
| | R-Prec | Recall@5 | Accuracy | F1 | KILT-AC | KILT-F1 |
| Re²G (ours) | **80.70** | **89.00** | **87.68** | **89.93** | **75.84** | **77.05** |
| KGI₁ [Glass et al., 2021] | 74.36 | 83.14 | _84.36_ | _87.24_ | _69.14_ | _70.58_ |
| KILT-WEB 2 [Piktus et al., 2021] | _75.64_ | _87.57_ | 81.34 | 84.46 | 64.64 | 66.64 |
| SEAL [Bevilacqua et al., 2022] | 67.80 | 81.52 | 83.72 | 86.53 | 60.08 | 61.72 |
| KGI₀ [Glass et al., 2021] | 59.70 | 70.38 | 77.90 | 81.31 | 55.54 | 56.79 |

| | Natural Questions | | | | | (Question Answering) |
|---|---|---|---|---|---|---|
| | R-Prec | Recall@5 | Accuracy | F1 | KILT-AC | KILT-F1 |
| Re²G (ours) | **70.78** | **76.63** | _51.73_ | _60.97_ | **43.56** | **49.80** |
| SEAL [Bevilacqua et al., 2022] | 63.16 | 68.19 | **53.74** | **62.24** | _38.78_ | _44.40_ |
| KGI₀ [Glass et al., 2021] | _63.71_ | 70.17 | 45.22 | 53.38 | 36.36 | 41.83 |
| KILT-WEB 2 [Piktus et al., 2021] | 59.83 | _71.17_ | 51.59 | 60.83 | 35.32 | 40.73 |
| RAG [Petroni et al., 2021] | 59.49 | 67.06 | 44.39 | 52.35 | 32.69 | 37.91 |

Jointly optimize reranking and truncation in one model, yield a dynamic document list for different queries



**GenRT**

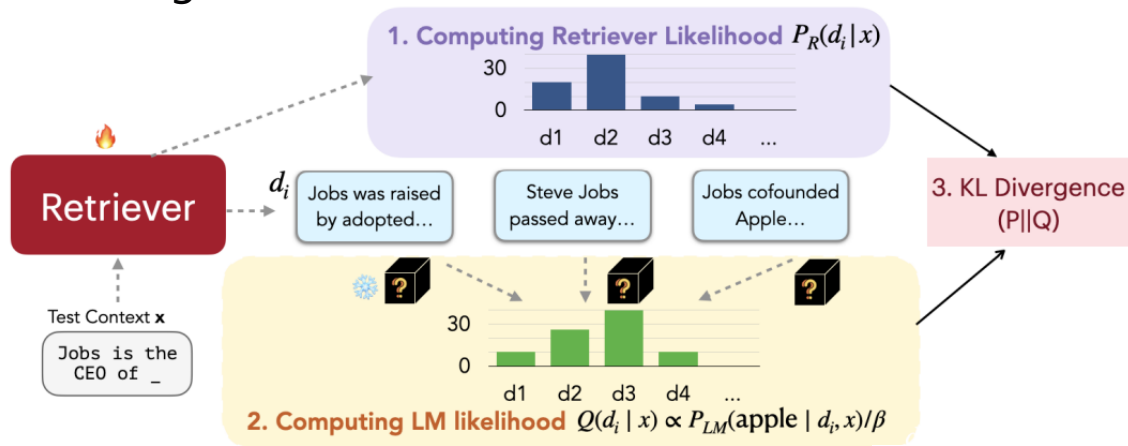| Truncation | NQ | | | TriviaQA | | |
|---|---|---|---|---|---|---|
| | TDCG ↑ | Length ↓ | Acc. ↑ | TDCG ↑ | Length ↓ | Acc. ↑ |
| Fixed-$x$ ($x$=5) | -0.78 | 5.00 | 54.80 | 0.23 | 5.00 | 60.03 |
| Fixed-$x$ ($x$=10) | -0.95 | 10.00 | 55.72 | -0.17 | 10.00 | 61.19 |
| Fixed-$x$ ($x$=20) | -1.67 | 20.00 | 56.98 | -1.10 | 20.00 | 62.35 |
| Fixed-$x$ ($x$=30) | -4.78 | 30.00 | 56.05 | -2.34 | 30.00 | 62.30 |
| Fixed-$x$ ($x$=40) | -5.05 | 40.00 | 58.20 | -3.46 | 40.00 | 63.17 |
| BiCut | -0.35 | 22.75 | 56.79 | 0.38 | 25.83 | 62.30 |
| Choppy | -0.20 | 25.43 | 57.01 | 0.40 | 29.72 | 62.42 |
| AttnCut | -0.21 | 17.70 | 56.95 | 0.42 | 21.96 | 62.40 |
| LeCut+JOTR | -0.15 | 20.21 | 57.84 | 0.55 | 22.50 | 62.89 |
| GenRT | **-0.06**[†] | 17.25 | 58.15 | **0.74**[†] | 22.19 | 63.25 |

- Compared with Fixed-40, GenRT achieves comparable accuracy with shorter length
- Compared with Fixed-20, GenRT achieves better performance with shorter length

# Requirement ③: Optimizable Objectives
## --- Remote Supervision Signals

Use LLM logits distribution as supervision to train the retriever, with the objective of minimizing KL divergence



Compute the retriever's scoring distribution over the document list:

$$P_R(d \mid x) = \frac{e^{s(d,x)/\gamma}}{\sum_{d \in \mathcal{D}'} e^{s(d,x)/\gamma}}$$

Compute the logits of the ground truth tokens for each document used in RAG

$$Q(d \mid x, y) = \frac{e^{P_{LM}(y \mid d, x)/\beta}}{\sum_{d \in \mathcal{D}'} e^{P_{LM}(y \mid d, x)/\beta}}$$

REPLUG: Retrieval-Augmented Black-Box Language Models. In 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2024
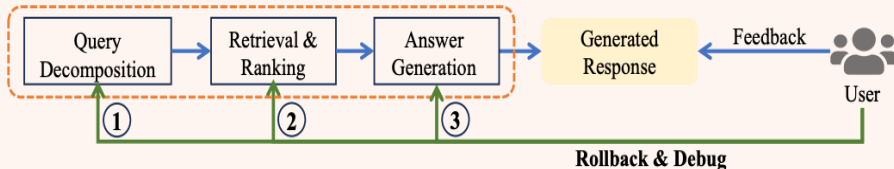
# Requirement ③: Optimizable Objectives
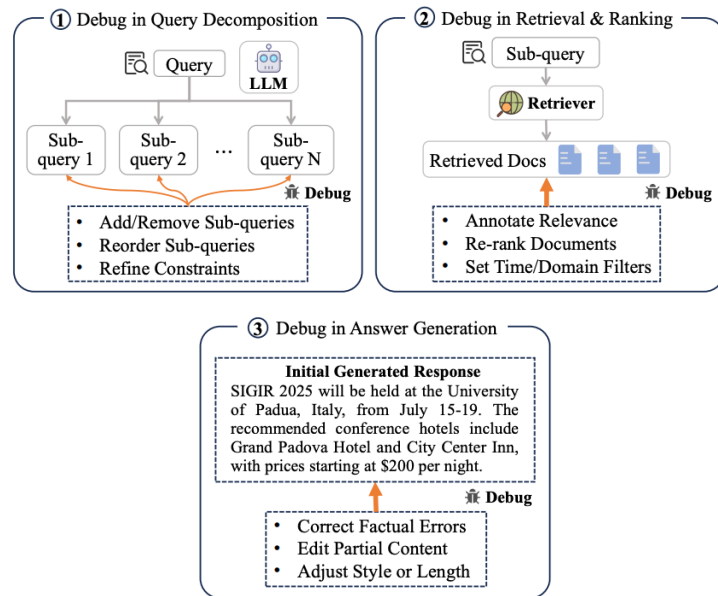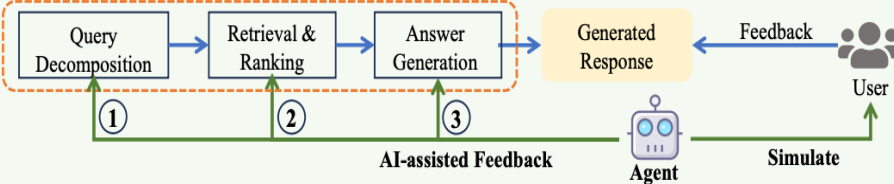## --- Build Feedback Loops

**User Debug Mode** allows engaged users to intervene at key stages, e.g. refining query decomposition, rating retrieved documents, and editing initial generated responses
**Shadow User Mode** a personalized user agent simulates user preferences and provides AI-assisted feedback for less interactive users
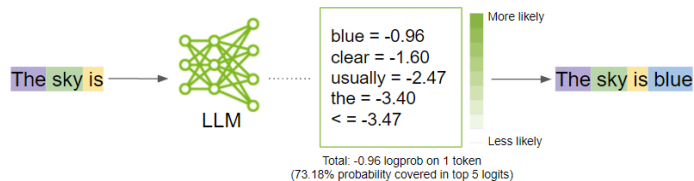
**02**
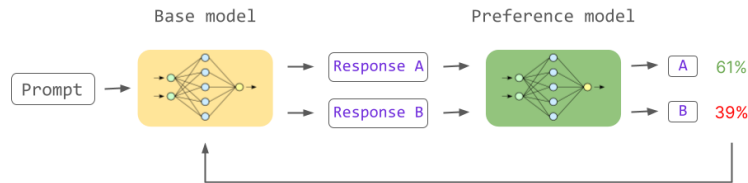
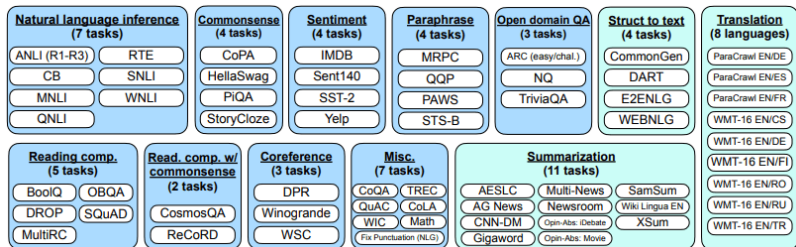**Large Language Model View @ RAG**

# Motivation: LLMs do not Learn RAG
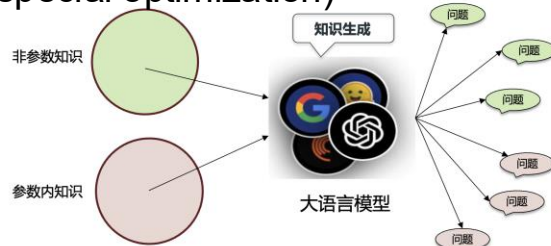
## ① Pretraining Phase – Next Token Prediction

The sky is → LLM → The sky is blue

blue = -0.96
clear = -1.60
usually = -2.47
the = -3.40
< = -3.47

More likely
Less likely

Total: -0.96 logprob on 1 token
(73.18% probability covered in top 5 logits)

## ③ RLHF Phase – Alignment

Base model          Preference model

Prompt → [network] → Response A → [network] → A  61%
              → Response B →          → B  39%

## ② Instruction Tunning Phase – Multi-task Learning

| Natural language inference (7 tasks) | | Commonsense (4 tasks) | Sentiment (4 tasks) | Paraphrase (4 tasks) | Open domain QA (3 tasks) | Struct to text (4 tasks) | Translation (8 languages) |
|---|---|---|---|---|---|---|---|
| ANLI (R1-R3) | RTE | CoPA | IMDB | MRPC | ARC (easy/chal.) | CommonGen | ParaCrawl EN/DE |
| CB | SNLI | HellaSwag | Sent140 | QQP | NQ | DART | ParaCrawl EN/ES |
| MNLI | WNLI | PiQA | SST-2 | PAWS | TriviaQA | E2ENLG | ParaCrawl EN/FR |
| QNLI | | StoryCloze | Yelp | STS-B | | WEBNLG | WMT-16 EN/CS |

| Reading comp. (5 tasks) | | Read. comp. w/ commonsense (2 tasks) | Coreference (3 tasks) | Misc. (7 tasks) | | Summarization (11 tasks) | | | WMT-16 EN/DE |
|---|---|---|---|---|---|---|---|---|---|
| BoolQ | OBQA | CosmosQA | DPR | CoQA | TREC | AESLC | Multi-News | SamSum | WMT-16 EN/FI |
| DROP | SQuAD | ReCoRD | Winogrande | QuAC | CoLA | AG News | Newsroom | Wiki Lingua EN | WMT-16 EN/RO |
| MultiRC | | | WSC | WIC | Math | CNN-DM | Opin-Abs: iDebate | XSum | WMT-16 EN/RU |
| | | | | Fix Punctuation (NLG) | | Gigaword | Opin-Abs: Movie | | WMT-16 EN/TR |

How to use retrieved information?
(no special optimization)

非参数知识

参数内知识

知识生成

[Google / ChatGPT icons] 大语言模型

问题
问题
问题
问题
问题
问题

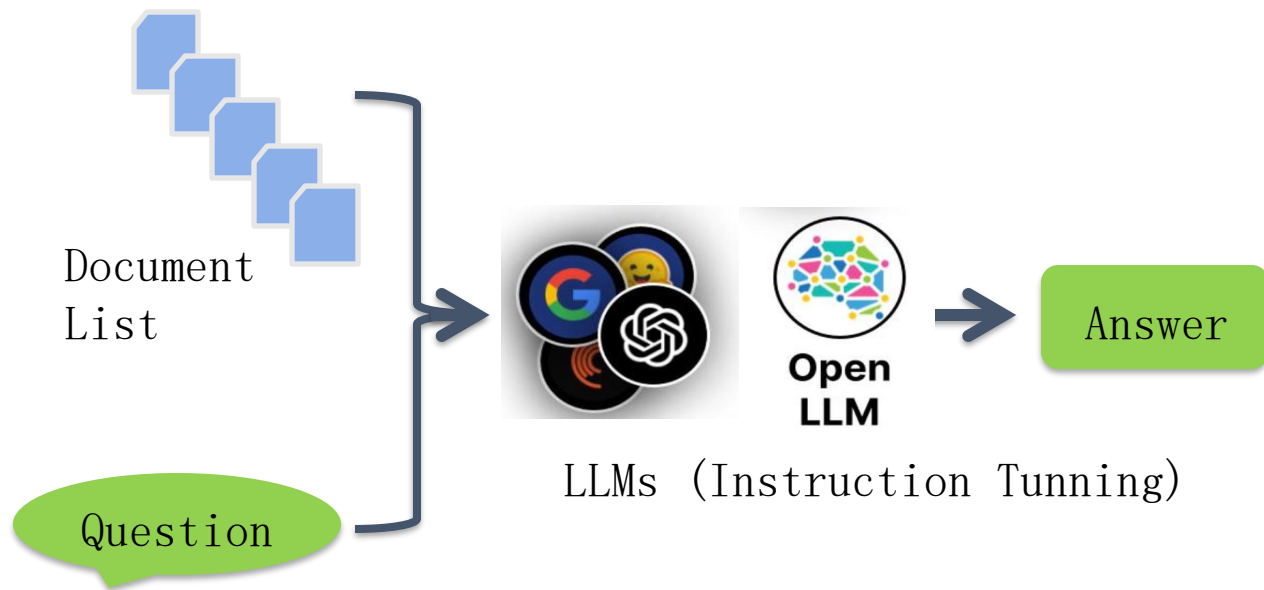## How can LLMs robustly handle noisy input knowledge and choose between internal and external knowledge?

# Motivation: LLMs do not Learn RAG

## Aligning LLMs capabilities in RAG through fine-tuning

◆ ① **Supervised Instruction Tuning:** Construct retrieval-question-answer triplets on domain-specific datasets and use them to fine-tune instructions, teaching the large model how to utilize retrieved documents. Examples include FID and RetRobust.

◆ ② **Dynamic Retrieval-Augmented Generation Fine-Tuning:** Fine-tune large language models to actively make dynamic decisions on whether to perform retrieval-augmented generation. Examples include Active-RAG and Self-RAG.

# ① Supervised Instruction Tuning

Given a question and a retrieved passage list R, use both as input for instruction fine-tuning



Document List

Question

LLMs (Instruction Tunning)

Answer

# ② Dynamic RAG Fine-tunning

Rowen： Retrieve Only When It Needs

SELF-RAG



Train an external discriminator to decide whether to use retrieved content, based on multi-dimensional consistency features (cross-language, noise addition, cross-model, etc.)

Fine-tune LLMs to dynamically generate retrieval tokens when needed during generation, critically evaluate retrieved documents, and use them selectively, enabling dynamic RAG

Retrieve Only When It Needs: Adaptive Retrieval Augmentation for Hallucination Mitigation in Large Language Models, Arxiv 2024
Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection, ICLR 2024

# Motivation: LLMs do not Learn RAG

**Aligning LLMs capabilities in RAG through fine-tuning**

All require supervised data

◆ ① Supervised Instruction Tuning

◆ ② Dynamic Retrieval-Augmented Generation Fine-Tuning

## Is supervised data essential?

# INFO-RAG: Unsupervised RAG Training

Design unsupervised training tasks according to three scenarios, so that LLM can play the role of "knowledge refiner"
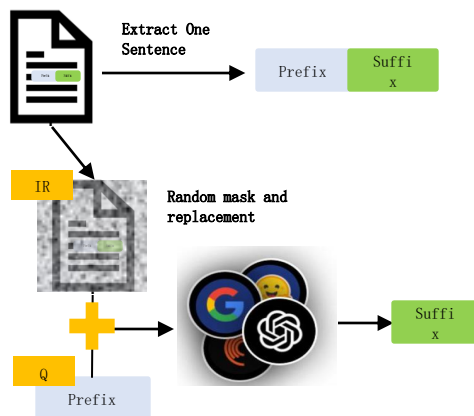


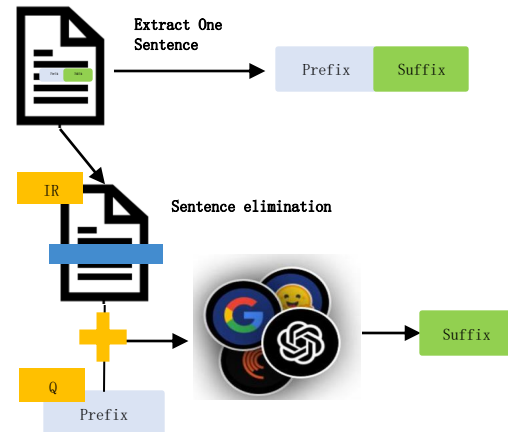Information Extraction — No Internal Knowledge

Information Correction — Partial Internal Knowledge
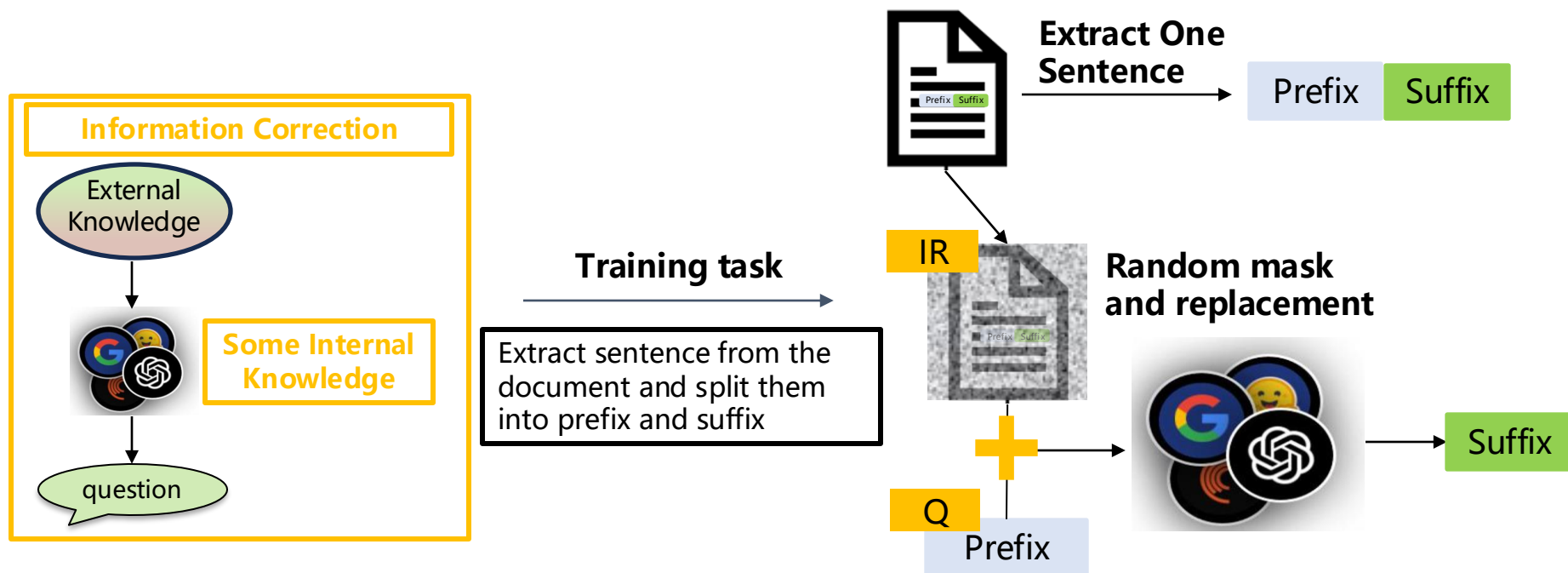
Information Provision — Exist Internal Knowledge

Unsupervised Information Refinement Training of Large Language Models for Retrieval-Augmented Generation. The 62nd Annual Meeting of the Association for Computational Linguistics (ACL'24)

# INFO-RAG: Method

All correct answers are in the retrieved texts and **LLMs just need to extract them**

**Information Extraction**

External Knowledge

**No Internal Knowledge**

question

**Training task**

Extract sentence from the document and split them into prefix and suffix

**Extract One Sentence**

Prefix | Suffix

IR

Prefix | Suffix

Q

Prefix

Suffix

# INFO-RAG: Method

The retrieved texts only contain partial answers, and even some wrong answers, which require **correction and completion** by LLM

# INFO-RAG: Method

The retrieved texts are only semantically related to the question but useless, and LLM needs to use this to **stimulate knowledge within parameters**

**Information Providing**

External Knowledge

**All Internal Knowledge**

question

**Training task**

Extract sentence from the document and split them into prefix and suffix

**Extract One Sentence**

Prefix | Suffix

Prefix Suffix

IR

**Sentence elimination**

Q

Prefix

Suffix

# INFO-RAG: Experiments

| | Soft-Filling Accuracy | | ODQA Accuracy | | Multi-Hop QA Accuracy | | LFQA ROUGE | Dialog F1 | LM ROUGE | Code Gen CodeBLEU | | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | T-REx | ZS | NQ | WebQ | Hotpot | Musique | ElI5 | Wow | WikiText | Python | Java | |
| LLaMA-2-7B | 55.60 | 54.08 | **46.82** | 43.52 | 39.40 | 25.95 | 15.18 | 7.85 | 60.77 | 21.44 | 22.99 | 35.78 |
| + INFO-RAG | **65.91** | **57.01** | 45.74 | **44.68** | **46.56** | **30.19** | **17.18** | **9.09** | **62.91** | **26.75** | **32.06** | **39.83** |
| LLaMA-2-7B-chat | 60.63 | 55.03 | 49.42 | 46.72 | 50.03 | 42.69 | 27.81 | 10.21 | 60.26 | 22.46 | 23.90 | 40.83 |
| + INFO-RAG | **65.77** | **58.32** | **53.93** | **49.13** | **52.01** | **44.45** | **28.15** | **10.49** | **63.24** | **27.25** | **28.79** | **43.78** |
| LLaMA-2-13B | 60.08 | 50.77 | 47.40 | 44.62 | 42.12 | 25.78 | 14.80 | 7.04 | 62.20 | 21.52 | 29.16 | 36.86 |
| + INFO-RAG | **62.80** | **55.63** | **47.82** | **45.42** | **51.48** | **35.02** | **17.48** | **7.20** | **64.14** | **29.00** | **35.50** | **41.04** |
| LLaMA-2-13B-chat | 62.53 | 56.81 | 50.36 | 45.47 | 61.23 | 47.06 | 27.07 | 11.19 | 60.52 | 22.34 | 30.96 | 43.23 |
| + INFO-RAG | **65.39** | **59.05** | **54.04** | **51.07** | **61.91** | **47.93** | **27.24** | **11.38** | **63.92** | **31.98** | **38.12** | **46.55** |

As an unsupervised training method, INFO-RAG can be applied to existing large models and further improve its ability to retrieve enhancements on various tasks
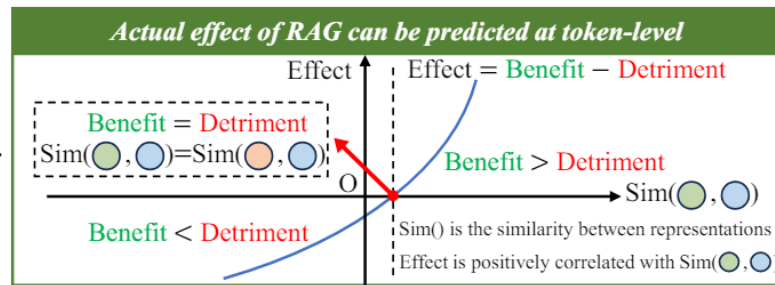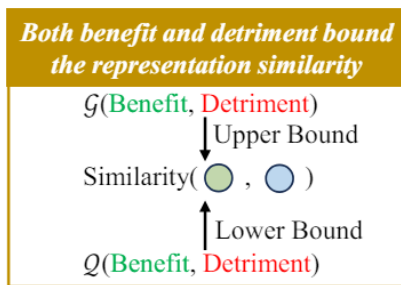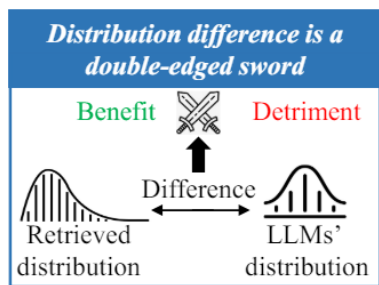
# Motivation: LLM maybe Already Know How to RAG

Most works on RAG are heuristically inspired and lack theoretical analysis explaining how RAG actually works

# TokRAG: Open the Blackbox of RAG

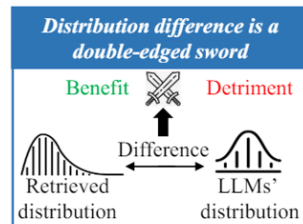## 1. Distribution difference brings benefits and detriments in RAG

**Benefit:** The large model gives an incorrect answer, while RAG gives a correct one.
**Detriment:** The large model gives a correct answer, while RAG gives an incorrect one.



Distribution difference is a double-edged sword

## 2. Theoretical basis: The text generation process of LLMs is an implicit latent variable inference (use to explain ICL (in-context learning)

$$p(x_i|R, x_{1:i-1}) = \int_{\mathcal{Z}} p(x_i|R, x_{1:i-1}, z)p(z|R, x_{1:i-1})\, dz$$

## 3. RAG can be treated as an unsupervised version of ICL

$z^*$ is Retrieved Concept

$$= \int_{\mathcal{Z}-\{z^*\}} p(x_i|R, x_{1:i-1}, z)p(z|R, x_{1:i-1})\, dz + p(x_i|R, x_{1:i-1}, z^*)p(z^*|R, x_{1:i-1}).$$

# TokRAG - Effect of RAG can be Predicted

**1. The target can be decomposed into <span style="color:red">benefit</span> and <span style="color:blue">detriment</span>**

$$\underbrace{\mathrm{KL}(p_R(r)\|p(r|z))}_{\text{benefit}} - \underbrace{\mathrm{KL}(p_R(r)\|p(r|z^*))}_{\text{detriment}}$$

Diff. between retrieved texts and LLM generated retrieved texts

Diff. between retrieved texts and LLM generated texts condition on Retrieved Concept

**2. Diff. between <span style="color:red">benefit</span> and <span style="color:blue">detriment</span> is positively correlated with the similarity of representation**

$$\underbrace{\mathrm{KL}(p_R(r)\|p(r|z))}_{\text{benefit}} - \underbrace{\mathrm{KL}(p_R(r)\|p(r|z^*))}_{\text{detriment}} \propto \frac{1}{\mathcal{D}}. \qquad \mathcal{D} = \|p(x_i|R, x_{1:i-1}) - p_R(x_i|x_{1:i-1})\|_1$$

# TokRAG - Collaborative Generation

**Principle to compare benefit and detriment in actual application**

$$s = \begin{cases} \text{benefit win} & \text{if } \cos(\mathbf{w}_{RAG}, \mathbf{w}_{IR}) \geq \cos(\mathbf{w}_{RAG}, \mathbf{w}_{LLM}), \\ \text{detriment win} & \text{if } \cos(\mathbf{w}_{RAG}, \mathbf{w}_{IR}) < \cos(\mathbf{w}_{RAG}, \mathbf{w}_{LLM}), \end{cases}$$



**(b) Our Practical Method:** Collaborative generation between pure LLM and RAG at the token-level by comparing benefit and detriment.

We can judge the actual effect of RAG at the token level. In this way, the collaborative generation of LLM and RAG can be realized, so as to maximize benefits and avoid detriments as much as possible

# TokRAG - Experiments

| Methods | Train LLM | Add Module | TriviaQA | | | | | | WebQ | | | | | | Squad | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Ratio of Hard Negative Passages | | | | | | Ratio of Hard Negative Passages | | | | | | Ratio of Hard Negative Passages | | | | | |
| | | | 100% | 80% | 60% | 40% | 20% | 0% | 100% | 80% | 60% | 40% | 20% | 0% | 100% | 80% | 60% | 40% | 20% | 0% |
| Standard RAG | no ✔ | no ✔ | 43.8 | 67.0 | 71.3 | 76.2 | 78.2 | 81.9 | 23.9 | 35.8 | 40.6 | 43.4 | 48.4 | 53.1 | 8.6 | 31.0 | 43.2 | 53.0 | 58.8 | 67.2 |
| NLI+RAG | no ✔ | need ✗ | 50.8 | 61.2 | 68.2 | 73.0 | 76.4 | 79.1 | 30.7 | 40.3 | 44.5 | 47.5 | 50.9 | 52.8 | 9.9 | 21.1 | 33.7 | 43.4 | 51.7 | 60.5 |
| CRAG | no ✔ | need ✗ | 48.2 | 68.3 | 72.5 | 76.7 | 81.5 | 82.2 | 25.6 | 37.4 | 41.9 | 46.2 | 51.5 | 54.9 | 7.4 | 28.7 | 39.6 | 50.7 | 53.2 | 61.1 |
| RetRobust | need ✗ | no ✔ | 49.2 | 67.3 | 72.9 | 77.5 | 79.4 | 82.3 | 30.0 | 38.9 | 42.5 | 48.2 | 49.8 | 54.3 | 10.5 | 30.8 | 43.3 | 52.5 | 58.4 | 66.0 |
| Self-RAG | need ✗ | no ✔ | 43.0 | 68.7 | 73.5 | 76.4 | 80.8 | 82.2 | 18.3 | 34.8 | 42.2 | 47.2 | 51.3 | 57.0 | 5.5 | 27.8 | 38.9 | 46.4 | 52.5 | 58.3 |
| INFO-RAG | need ✗ | no ✔ | 49.7 | 68.4 | 73.2 | 77.9 | 80.0 | 82.5 | 29.7 | 38.0 | 43.9 | 48.1 | 49.4 | 54.8 | 10.7 | 30.1 | 43.5 | 53.7 | 59.2 | 67.5 |
| X-RAG (Ours) | no ✔ | no ✔ | **53.5** | **72.9** | **77.6** | **81.3** | **83.4** | **85.7** | **32.9** | **43.8** | **47.3** | **50.0** | **52.9** | **57.3** | **12.8** | **31.3** | **44.5** | **54.1** | **60.8** | **68.1** |

In RAG of actual open-domain QA tasks, X-RAG can surpass mainstream robust RAG frameworks and training methods, such as RetRobust, Self-RAG, etc., without the need for additional modules or training LLM.
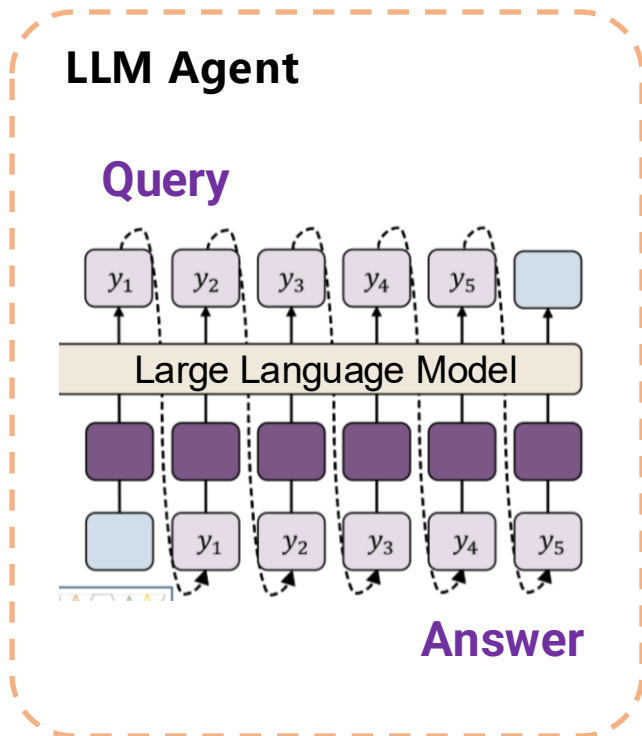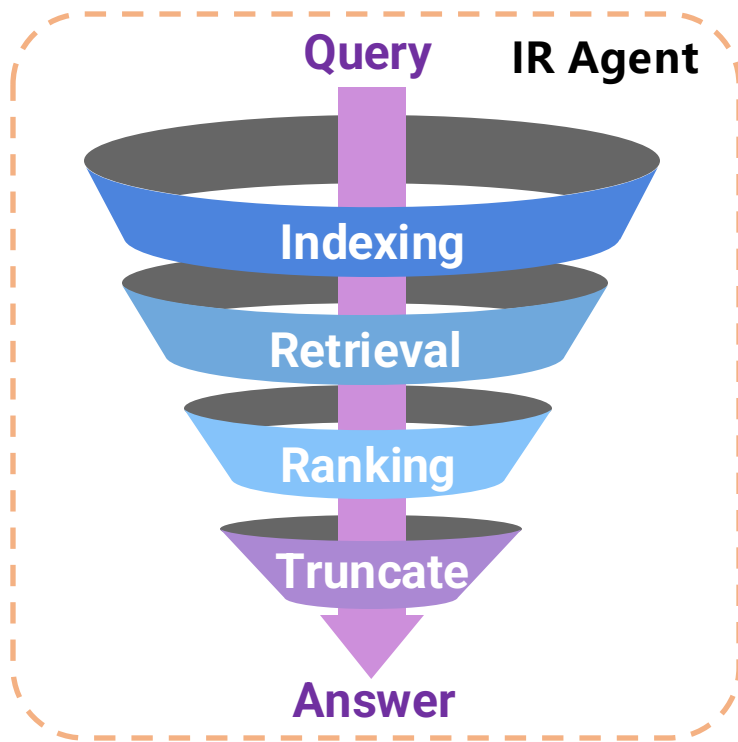
# 03
# Interaction View
# @RAG

# Motivation: Treat IR and LLM Equally

How can large models and information retrieval interact efficiently to robustly solve complex problems?

# Motivation: Make IR and LLM Interactively

**Interaction Framework between IR and LLM：**

① Tool Calling,

    e.g., ToolFormer

② Complex Problem Decomposition,

    e.g., Self-Ask, DSP

③ Agent-Based Planning,

    e.g., ReAct

④ Information Correction,

    e.g., Verify-and-Edit

# ① Interaction Based on Tool Calling

## ToolFormer

The New Engla

Interaction Process

Tool Types

What other name is Pittsburgh known by?

? → The Steel City

War memorial Flodden

🔍 → [...] was created in memory of the Battle of Flodden.

3435*235/9

→ 89691.67

∅

Thursday, March 10, 2019

Os Melhores Escolas em Jersey

→ The Best Schools in Jersey

**Advantages:**
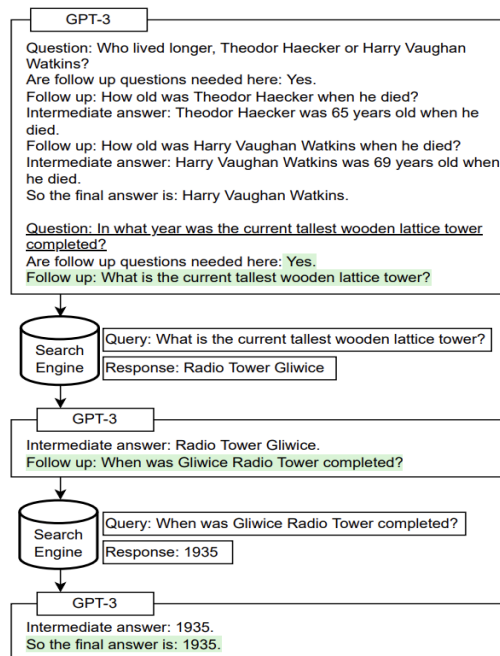  ① Diverse tool types
  ② Easy to synthesize training data

**Disadvantages:**
  ① Local Planning (interrupt the decoding process when "→" token)
  ② Predefined tool types
  ③ Without document content

Toolformer: Language models can teach themselves to use tools. Advances in Neural Information Processing Systems, NIPS 2023.

# ② Interaction Based on Complex Problem Decomposition

## Self-Ask

Break the question into follow-up questions, which are easier to answer with LLM



**Advantages:**
 ① Break a hard problem to some easy problems
 ② Easy to identify where to use IR

**Disadvantages:**
 ① Local Planning (end of follow up question)
 ② Every sub-questions can be answered (strong assumption)
 ③ Without document content

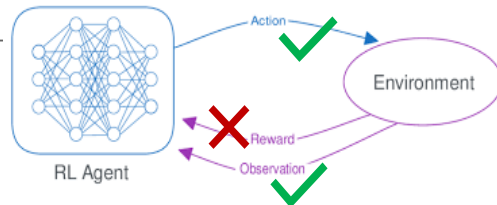# ③ Interaction Based on Agent-Based Planning

## ReAct

**(1d) ReAct (Reason + Act)**

**Thought 1:** I need to search Apple Remote and find the program it was originally designed to interact with.
**Act 1:** Search[Apple Remote]
**Obs 1:** The Apple Remote is a remote control introduced in October 2005 by Apple … originally designed to control the Front Row media center program …

**Thought 2:** Apple Remote was originally designed to control the Front Row media center program. I need to search Front Row next and find what other device can control it.
**Act 2:** Search[Front Row]
**Obs 2:** Could not find [Front Row]. Similar: ['Front Row Seat to Earth', 'Front Row Motorsports',' Front Row (software)', …]

**Thought 3:** Front Row is not found. I need to search Front Row (software) .
**Act 3:** Search[Front Row (software)]
**Obs 3:** Front Row is a discontinued media center software …

**Thought 4:** Front Row (software) is controlled by an Apple Remote or the keyboard function keys. So the answer is keyboard function keys.
**Act 4:** Finish[keyboard function keys]

**Advantages:**
① The prototype of an agent, including elements: thought, action, and observation (document)

**Disadvantages:**
① Local Planning (end of obs.)
② No reward signals



Action ✓
Environment
✗ Reward
Observation ✓
RL Agent

REACT: SYNERGIZING REASONING AND ACTING IN LANGUAGE MODELS. ICLR 2023.

# ④ **Interaction Based on Information Correction**

## **Verify-and-Edit**



**Question**

Of all the teams John Nyskohus played for, which team was known as "the Black and Whites?"

**Standard**

Newcastle United. ❌

**Chain-of-thought**

First, John Nyskohus played for the Norwegian football team Odd Grenland. Second, Odd Grenland is known as "the Black and Whites." The answer is **Odd Grenland**. ❌

- - - - - - - **Self-Consistency:** - - - - - -
less than majority agree

**Verify**

What team did John Nyskohus play for?
What team is known as "the Black and Whites?"

**External Knowledge Retrieval**

John Nyskohus ... is an Australian former soccer player who played club football for USC Lion ... and Adelaide City in the National Soccer League ...
Adelaide City Football Club is an Australian football (soccer) club based in Adelaide, South Australia. They are also known as "The Zebras" and "the Black and Whites."

**Edit Rationales**

First, John Nyskohus played for Adelaide City in the National Soccer League. Second, Adelaide City Football Club is known as "the Black and Whites".

**New Prediction**

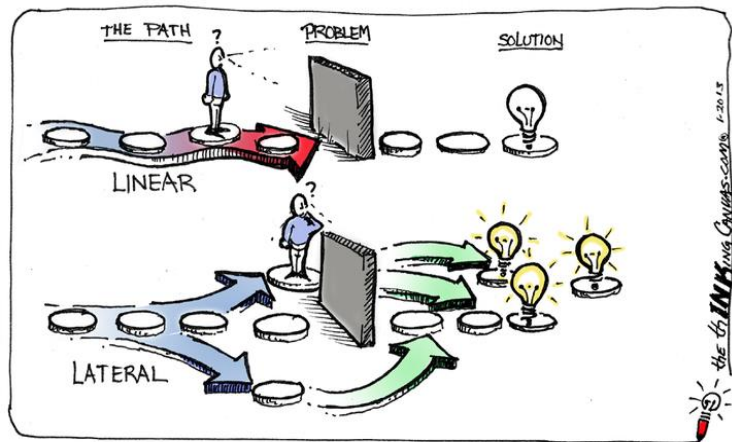The answer is **Adelaide City Football Club.** ✅

**Advantages:**
   ① Global Planning (generate all reasoning in one round)
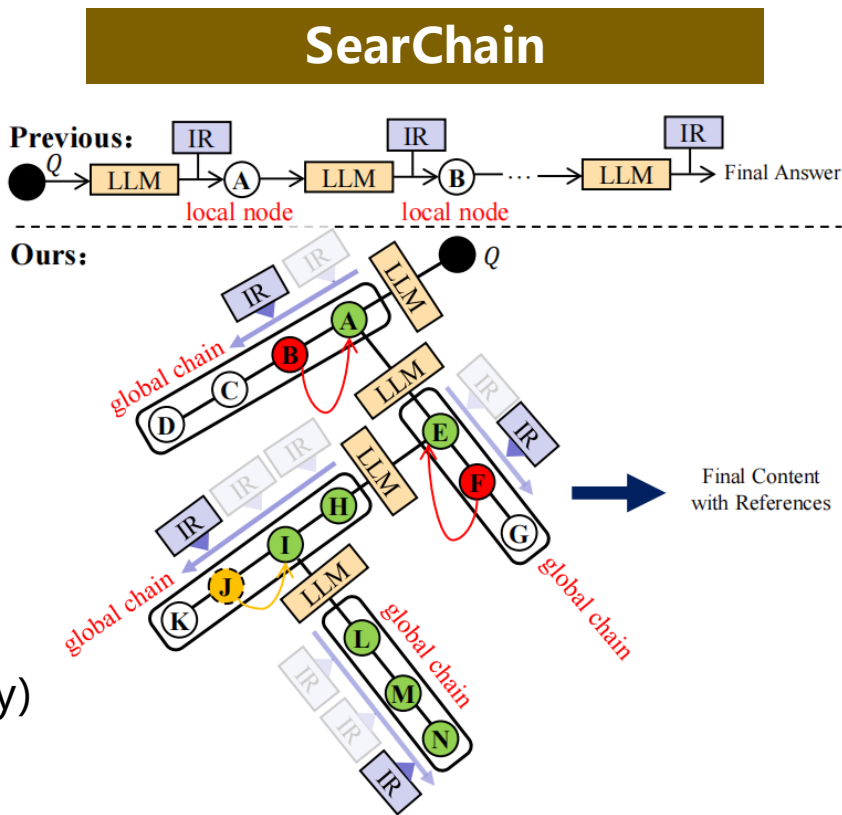   ② Self consistence verify (reward)

**Disadvantages:**
   ① Not fit agentic framework
   ② Process hard to trace (boundary of sub-question and reasoning block are blurred)

Verify-and-edit: A knowledge-enhanced chain-of-thought framework. ACL 2023.

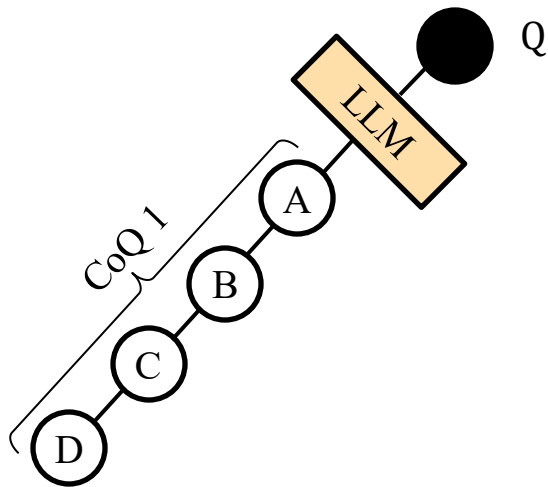# SearChain: Tree-Structured Interaction Framework



- **CoT vs. Agentic Framework**
  IR and LLM as two interacting agents
- **Local vs. Global Decomposition**
  Complete reasoning chain (chain-of-query)
- **Linear vs. Tree Reasoning**
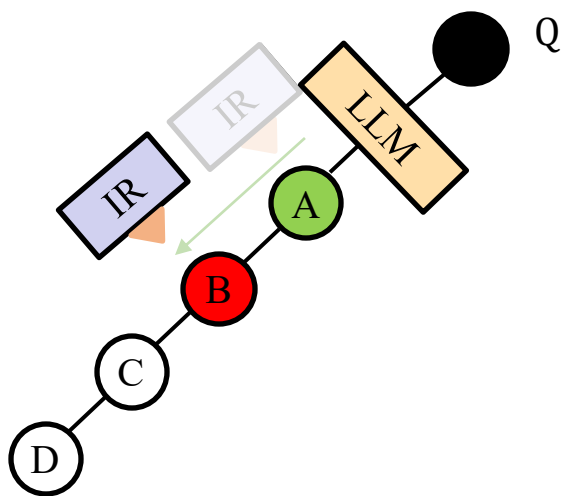  IR verify and correct reasoning direction



Search-in-the-Chain: Interactively Enhancing Large Language Models with Search for Knowledge-intensive Tasks. WWW 2024.

# SearChain - Method

## Step1: Generation Chain-of-Query (Global Decomposition)

# SearChain - Method

## Step2: IR module go though each sub-question node, verify or complete
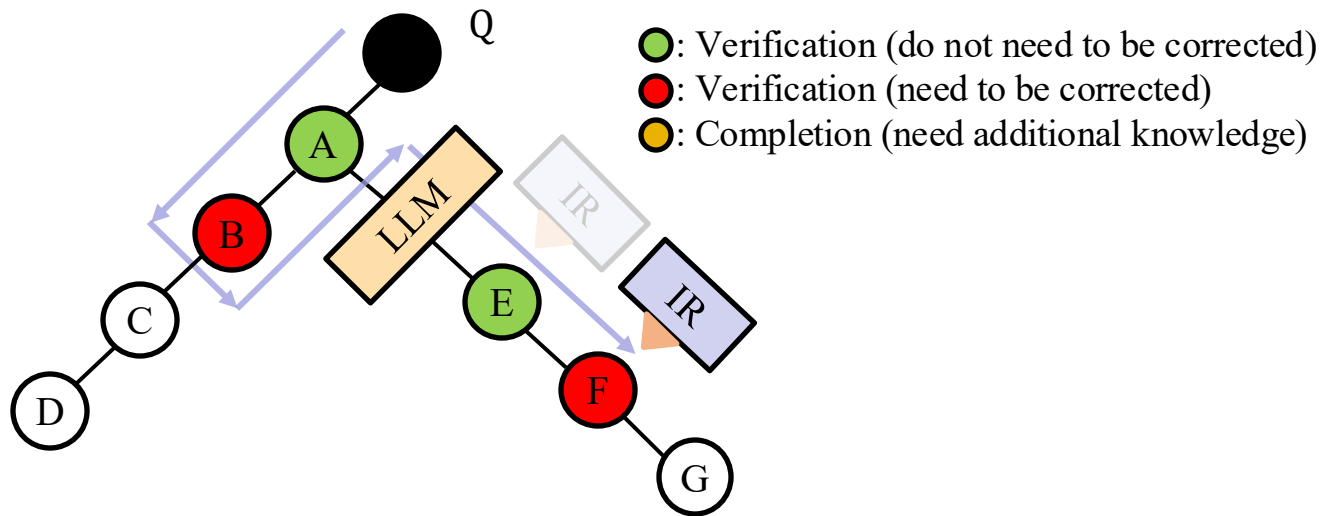
# SearChain - Method

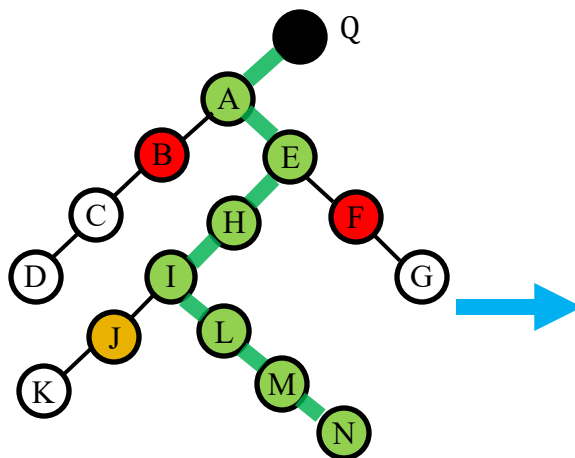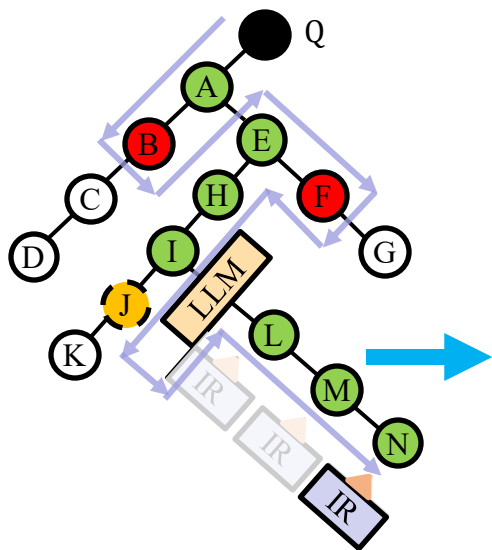## Step3: If Error occurs, go back to the previous node and generate CoQ again

# SearChain - Method

## Step4: Repeat using IR module to go though the remained nodes

# SearChain - Method

## Step5: Track back to get evidence-cited answer



○ : Verification (do not need to be corrected)

○ : Verification (need to be corrected)

○ : Completion (need additional knowledge)

The performer of Spirit If... is Kevin Drew [1]. Kevin Drew was born in Toronto [2]. Greyhound buses in Toronto leave from Toronto Coach Terminal [3]. So the final answer is Toronto Coach Terminal. ✔

[1] Spirit If... is the debut solo album by **Kevin Drew.** It was released on September 18, 2007 …

[2] Kevin Drew (born September 9, 1976 in **Toronto**) ..

[3] The **Toronto Coach Terminal** is the central bus station for inter-city services in Toronto, Ontario, Canada … when it was leased out in its entirety to bus lines Coach Canada and Greyhound Canada …

# SearChain - Experiment

## Performance on knowledge-intensive tasks

|  | Muti-Hop QA | | | | Slot Filling | | FC | LFQA |
|---|---|---|---|---|---|---|---|---|
|  | HoPo | MQ | WQA | SQA | zsRE | T-REx | FEV. | ELI5 |
| **Without Information Retrieval** | | | | | | | | |
| Direct Prompting | 31.95 | 5.91 | 25.82 | 66.25 | 22.75 | 43.85 | 73.45 | 21.90 |
| Auto-CoT | 33.53 | 10.55 | 29.15 | 65.40 | 21.30 | 43.98 | 76.61 | 21.55 |
| CoT | 35.04 | 9.46 | 30.41 | 65.83 | 22.36 | 44.51 | 76.98 | 21.79 |
| CoT-SC | 36.85 | 10.02 | 32.68 | 70.84 | 24.74 | 46.06 | 77.15 | 22.05 |
| Recite-and-answer | 36.49 | 10.97 | 32.53 | 70.47 | 24.98 | 46.14 | **77.35** | 22.10 |
| Self-Ask w/o IR | 33.95 | 11.10 | 35.65 | 65.45 | 20.16 | 44.71 | 75.31 | 21.73 |
| Least-to-Most | 34.05 | 11.45 | 32.88 | 65.78 | 21.86 | 44.98 | 75.98 | 21.95 |
| Plan-and-Solve | 36.33 | 12.95 | 35.68 | 73.21 | 25.15 | 47.58 | 77.08 | 22.23 |
| SearChain w/o IR | **38.36** | **13.61** | **40.49** | **75.62** | **30.14** | **52.69** | 77.06 | **22.54** |
| **Interaction with Information Retrieval** | | | | | | | | |
| Direct Retrieval | 34.09 | 10.22 | 30.01 | 66.78 | 52.29 | 59.28 | 78.25 | 23.40 |
| ToolFormer | 36.75 | 12.98 | 35.49 | 67.02 | 51.35 | 59.17 | 80.79 | 23.05 |
| Self-Ask | 40.05 | 14.28 | 39.58 | 67.65 | 50.51 | 59.12 | 79.41 | 23.25 |
| Plan-and-Solve w/ IR | 41.65 | 15.07 | 42.05 | 74.58 | 52.15 | 60.03 | 81.04 | 24.56 |
| React → CoT-SC | 43.15 | 15.49 | 40.36 | 70.43 | 53.27 | 60.42 | 80.59 | 24.05 |
| Verify-and-Edit | 44.03 | 15.57 | 40.83 | 71.09 | 53.95 | 61.10 | 80.67 | 23.80 |
| Tree-of-Thought w/ IR | 50.65 | 15.61 | 42.49 | 72.55 | 54.88 | 62.40 | 81.03 | 24.20 |
| DSP | 51.97 | 15.83 | 43.52 | 72.41 | 54.35 | 61.32 | 80.65 | 23.46 |
| SearChain | **56.91** | **17.07** | **46.27** | **76.95** | **57.29** | **65.07** | **81.15** | **25.57** |
| - w/o Verification | 46.11 | 14.70 | 42.67 | 75.98 | 43.58 | 55.46 | 78.79 | 22.98 |
| - w/o Completion | 53.05 | 15.86 | 43.64 | 76.53 | 45.78 | 56.03 | 80.03 | 25.02 |

1. In reasoning，outperforms CoT, Self-consistency and Plan-and-Solve
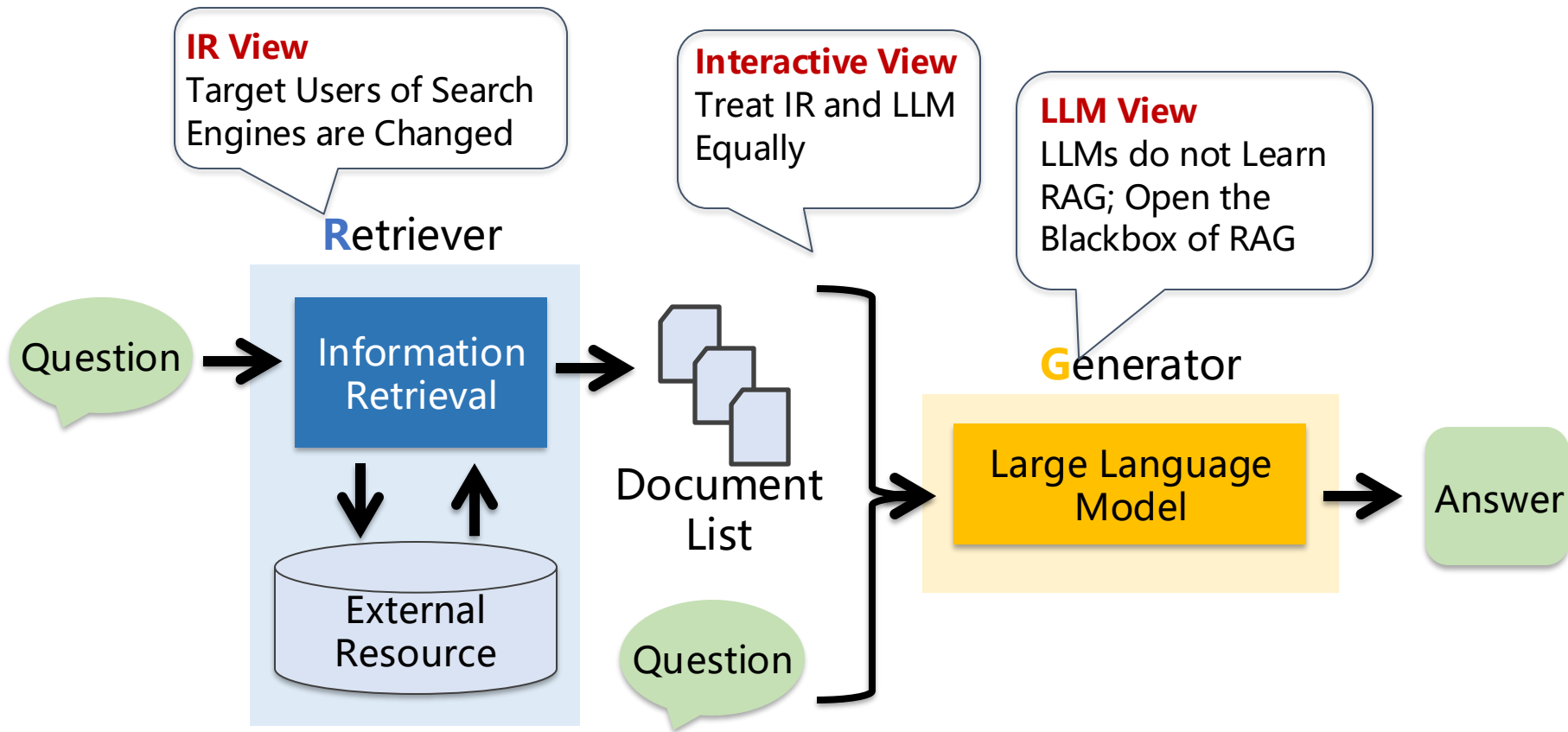2. In the interaction with IR，outperforms React, Self-Ask and Tree-of-Thought

# SearChain - Experiment

## Better Tracing

| What team did the recipient of the 2007 Brownlow Medal play for? | | What kind of dinosaur does Raymond Ochoa voice in the Pixar film "The Good Dinosaur"? | | In what London botanical garden can you find the fungus described by Elise Maud Wakefield as "Lepiota"? | |
|---|---|---|---|---|---|
| **SearChain** | **New Bing** | **SearChain** | **New Bing** | **SearChain** | **New Bing** |
| Jimmy Bartel was the recipient of the 2007 Brown-low Medal [1] and played Australian rules football professionally [2]. He played for the Geelong Football Club in the Australian Football League [2]. So the final answer is Geelong Football Club.<br><br>[1]: … Jimmy Bartel won the medal during the 2007 AFL season …<br>[2]: … Jimmy Bartel is a former Australian rules footballer plays for the Geelong Football Club … | The recipient of the 2007 Brownlow Medal was Jimmy Bartel who played for the Geelong Football Club. He won the medal by polling twenty-nine votes during the 2007 AFL season [1]. Is there anything else you would like to know about this topic?<br><br>[1]: … Jimmy Bartel won the medal by polling twenty-nine votes during the 2007 AFL season… | Raymond Ochoa is an American child actor and voice actor [1]. He has played various roles in movies including Arlo in "The Good Dinosaur" [2]. Arlo is an Apatosaurus [3]. So the final answer is Apatosaurus.<br><br>[1]: Raymond Ochoa is an American child actor and voice actor …<br>[2]:Raymond Ochoa played, Arlo, in the Pixar film "The Good Dinosaur"…<br>[3] … timid Apato-saurus named Arlo (Ochoa) … | Raymond Ochoa vo-ices the character of Arlo, an Apatosaurus in the Pixar film "The Good Dinosaur" [1]<br><br>[1]:(Amazon.com: The Good Dinosaur : Ray-mond Ochoa) | Elise Maud Wakefield was an English mycologist and plant pathologist [1]. "Lepiota" is a genus of fungi [2]. Kew Gardens is known for its fungi collection [3], such as "Lepiota viridigleba" [4]. So the final answer is Kew Gardens.<br><br>[1]: … Elise Maud Wakefield, was an English mycologist and plant pathologist…<br>[2]: … Lepiota viridigleba is a species of seq-uestrate fungus in the family Agaricaceae …<br>[3] … Kew Gardens is a botanical garden in southwest London that houses the "largest and most diverse botanical and mycological collections in the world". ..<br>[4] … The collection of fungi at Kew Gardens includes species such as … "Lepiota viri-digleba" … | I'm sorry, but I couldn't find any information about a fungus described by Elsie Maud Wakefield as "Lepiota" in a London botanical garden. I did find that Elsie Maud Wakefield was a mycologist at the Royal Botanic Gardens, Kew from 1910 to 1951 and was a specialist in Basidiomycetes and an international authority on Aphyllophorales [1,2].<br><br>[1]: … Elsie Maud Wakefield, was an English mycologist and plant pathologist…<br>[2]: … Elsie Wakefield was a specialist in Basidiomycetes and an international auth-ority on Aphyllophorales. … |

Compared with New Bing, SearChain can trace the source of more fine-grained knowledge, and the traceable marking position is more accurate

# Conclusion



Three views of RAG approaches