

When RAG Meets LFM: Towards Retrieval-Augmented Large Foundation Models

Website: <https://advanced-rag.github.io/RAG-Meets-LFMs>

Survey: <https://arxiv.org/pdf/2405.06211>

Xu Yuan¹, Yujuan Ding¹, Chengliang Liu¹, Rui An¹, Chun-Hin Chan¹,

Yiqi Wang², Wenqi Fan¹, and Qing Li¹

¹The Hong Kong Polytechnic University

²National University of Defense Technology

June 9th (Day 1)

PAKDD 2026, Hong Kong, China



Tutorial Outline

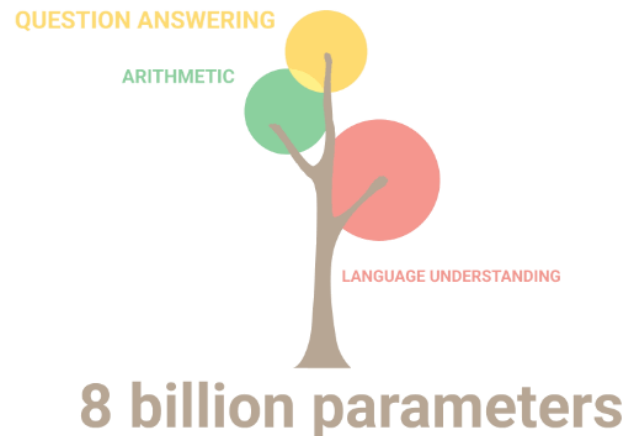


- ① **Part 1: Introduction of Retrieval Augmented Large Foundation Models (RA-LFMs) (Dr. Wenqi Fan)**
- **Part 2: Architecture** of RA-LFMs and Main Modules (Xu Yuan)
- **Part 3: Learning Approach** of RA-LFMs (Chengliang Liu)
- **Part 4: Agentic RAG** (Chengliang Liu)
- **Part 5: Applications** of RA-LFMs (Chun-Hin Chan)
- **Part 6: Challenges and Future Directions** of RA-LFMs (Dr. Wenqi Fan)
- **Part 7: Q&A**

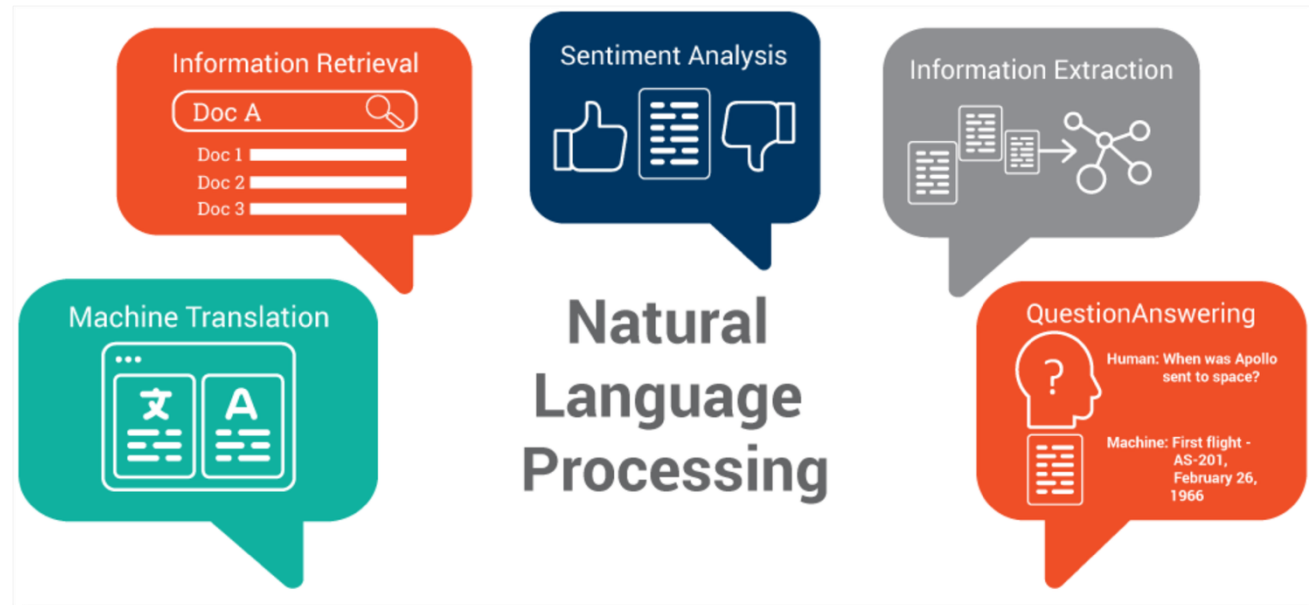
Website of this tutorial
Check out the slides and more information!



Large Foundation Models (LFMs)

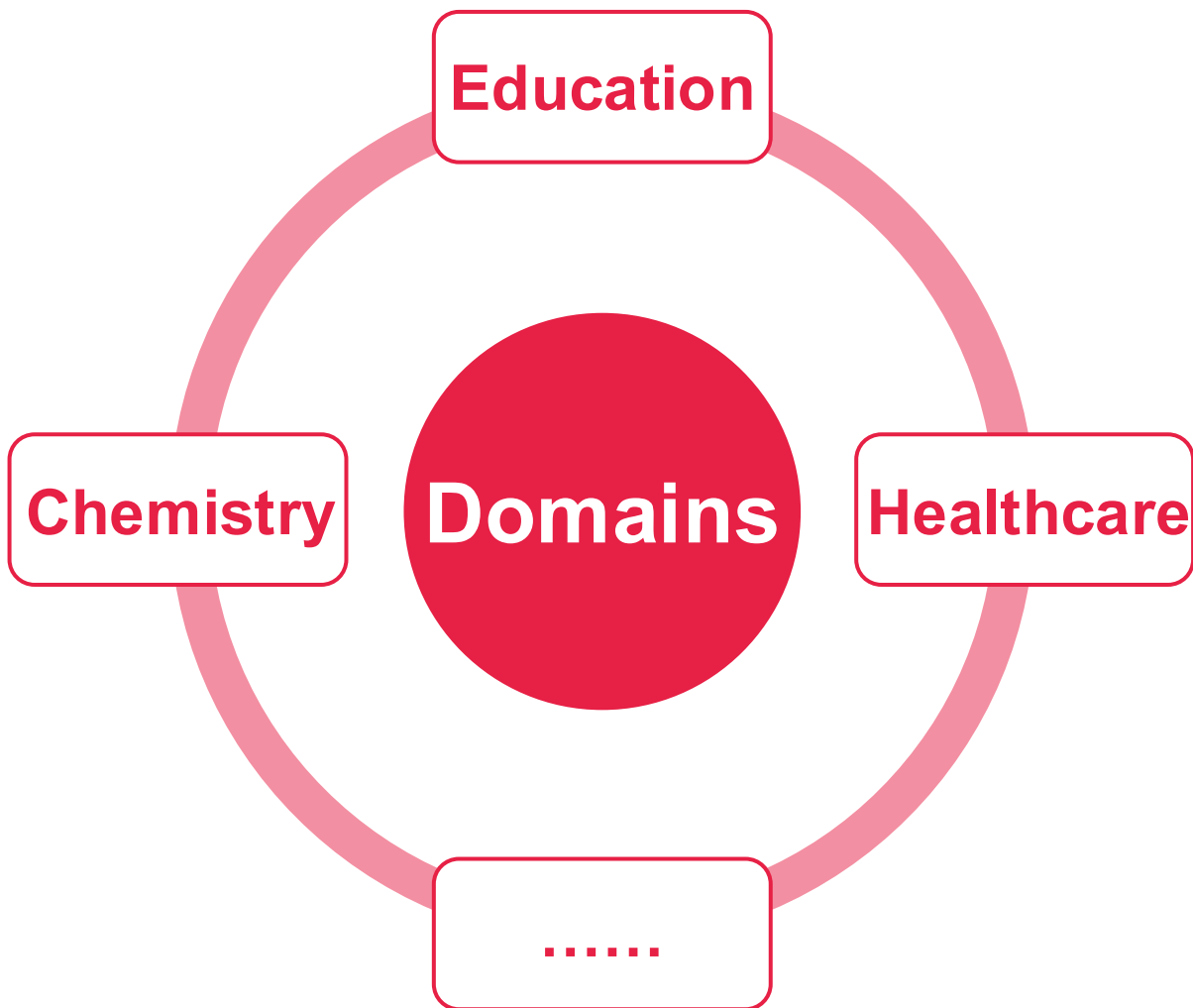


Large Foundation Models (LFMs)

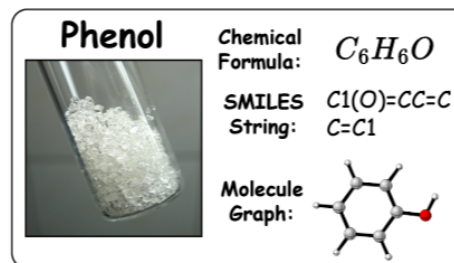


Large Foundation Models (LFMs)

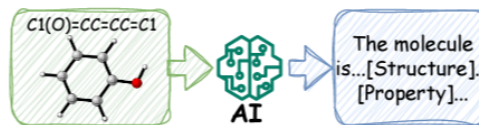
LFMs in Downstream Domains



☐ Molecule discovery, etc.



(a) Molecule Representations.



(b) Molecule Captioning.



ChatGPT

(a) Molecule Captioning

Please show me a description of this molecule: "C1=CC=C(C=C1)OC2=CC=CC=C2"

The molecule is an aromatic ether in which the oxygen is attached to two phenyl substituents. It has been found in muscat grapes and vanilla. It has a role as a plant metabolite.

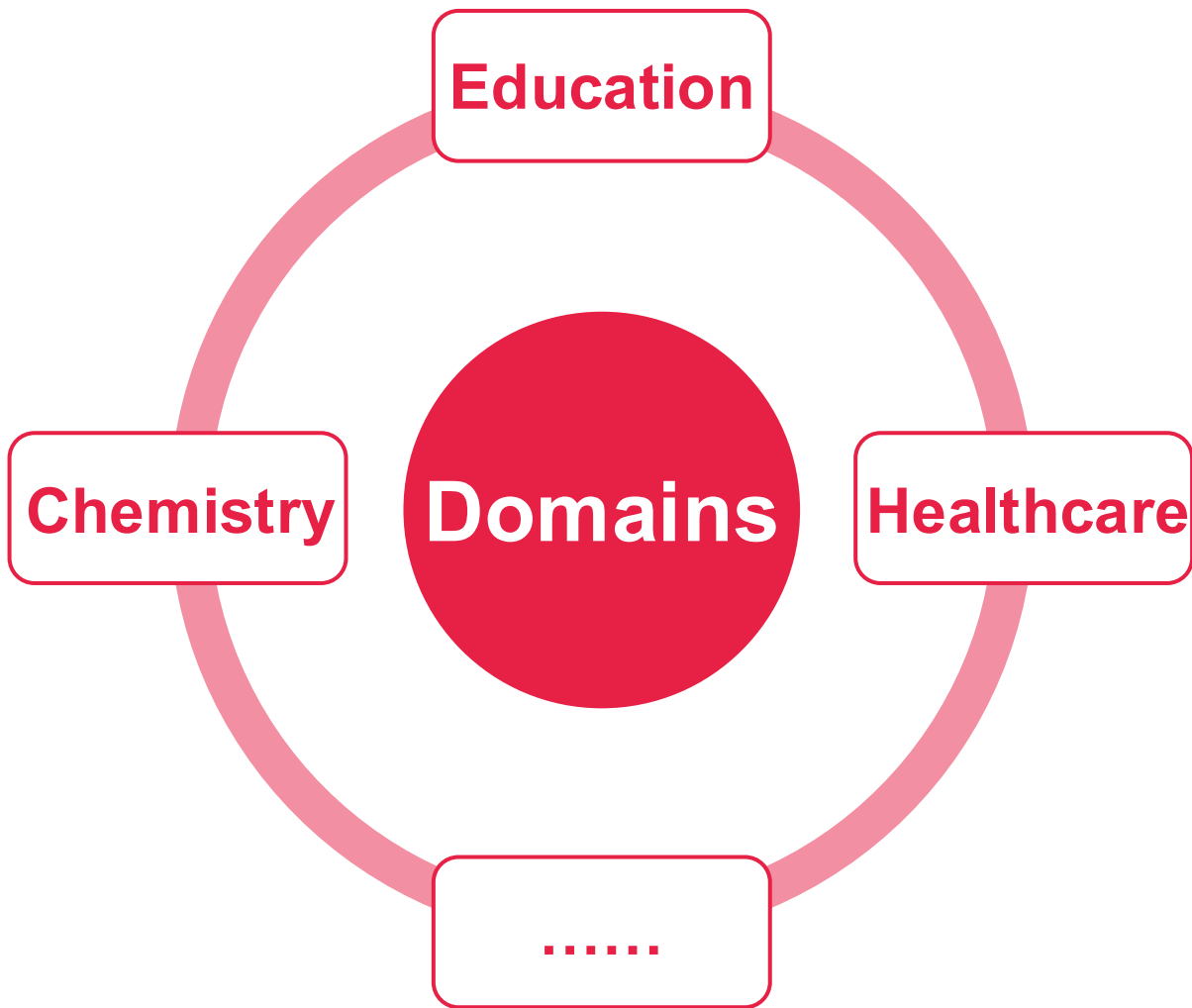
(b) Text-based Molecule Generation

Help me generate a molecule based on the given description:

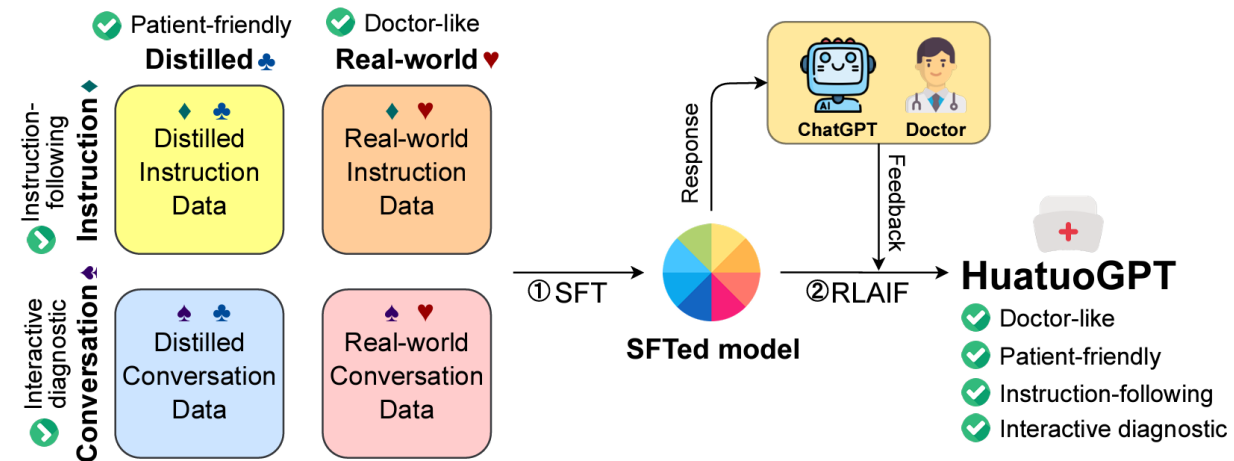
The molecule is a quinolinemonocarboxylate that is the conjugate base of xanthurenic acid, obtained by deprotonation of the carboxy group. It has a role as an animal metabolite. It is a conjugate base of a xanthurenic acid.

C1=CC2=C(C(=C1)[O-])NC(=CC2=O)C(=O)O

LFMs in Downstream Domains



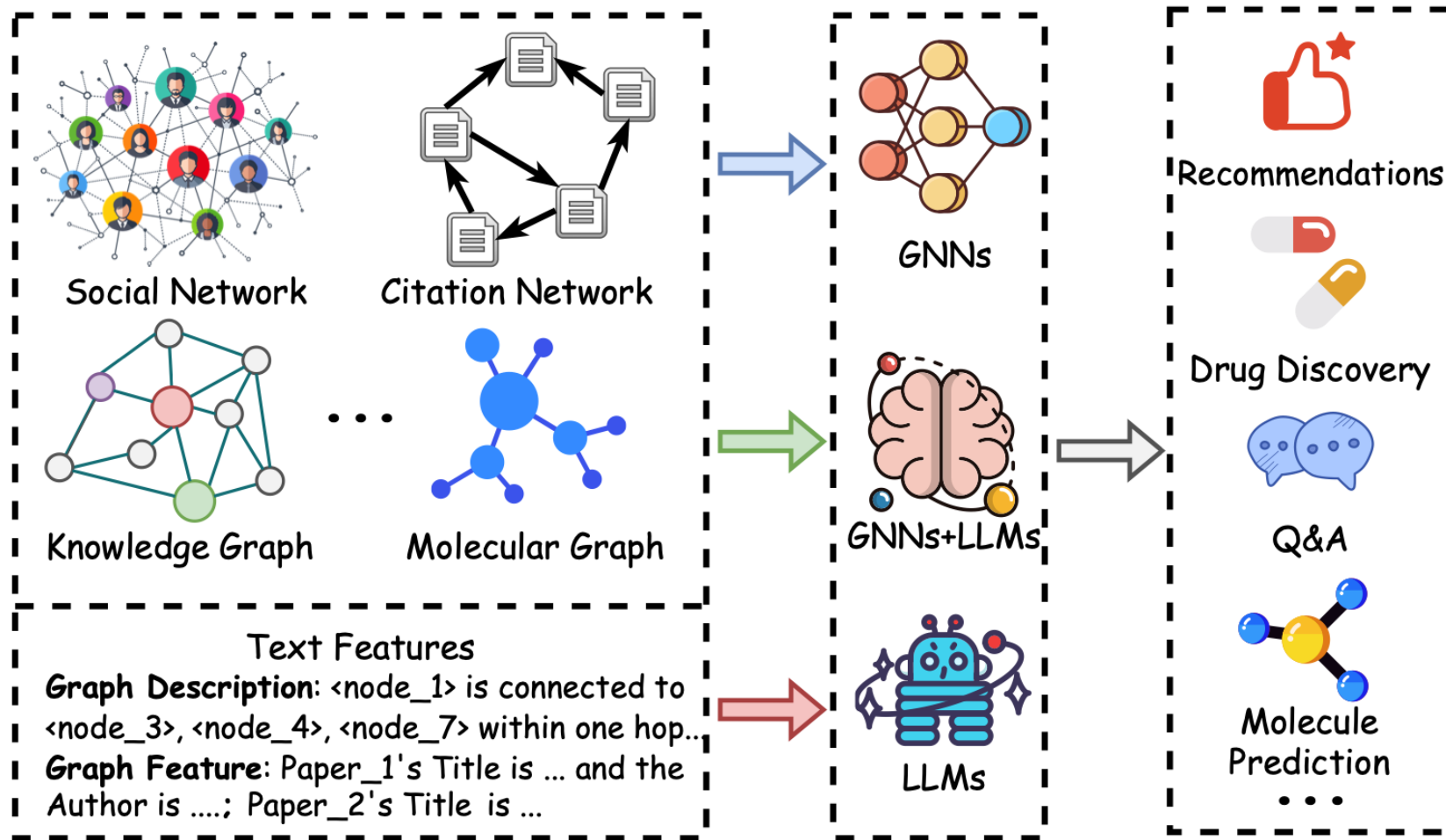
Medical consultation, etc.



Curriculum & Teaching, etc.

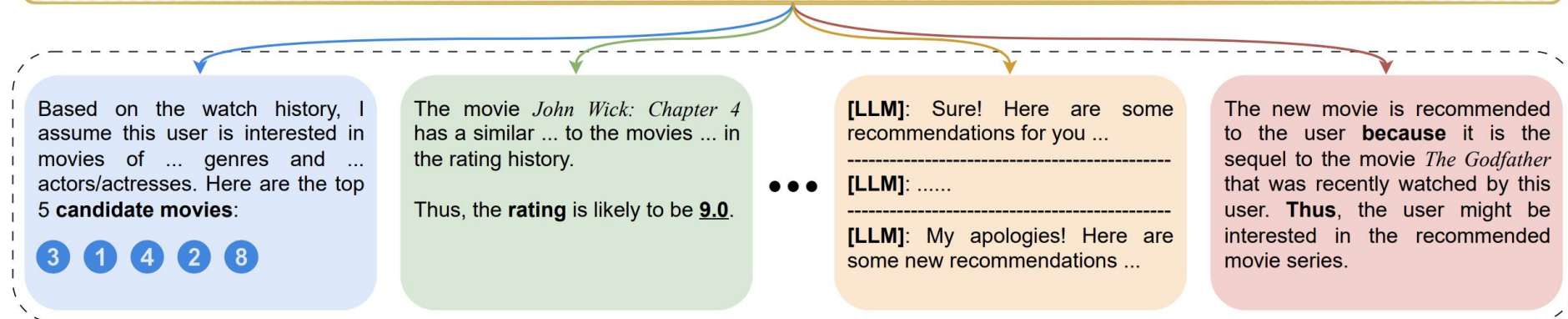
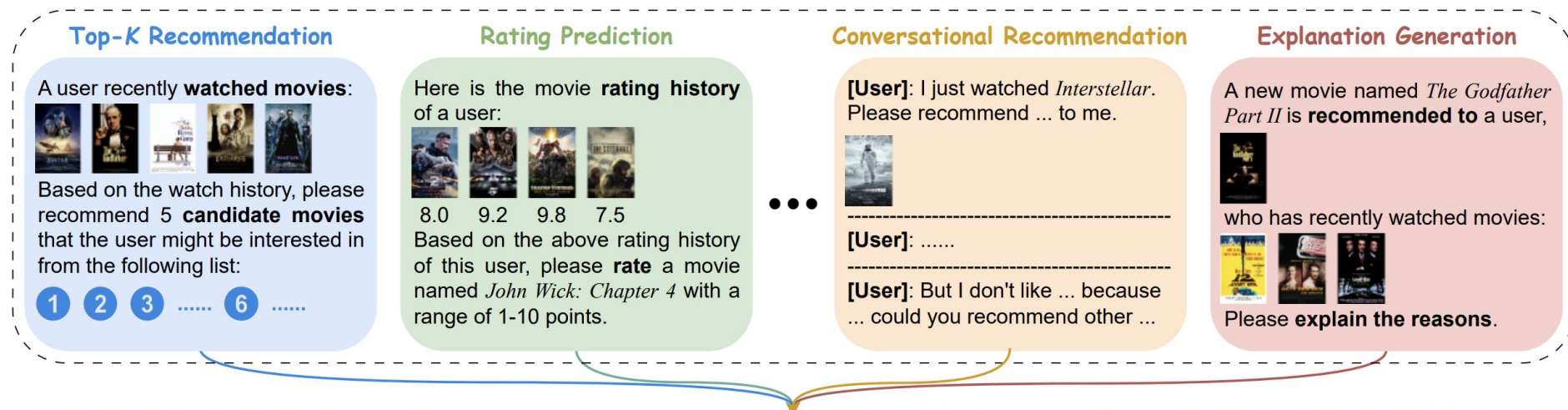


LFMs on Graph-structured Data



LFMs in Recommender Systems

Task-specific Prompts (LLMs Inputs)



Task-specific Recommendations (LLMs Outputs)

Challenges and Risks of LFM

❑ Outdated Parametric Knowledge

LFMs rely on knowledge stored in model parameters, which may become outdated and cannot easily reflect newly updated documents, policies, or domain facts.



❑ Domain Knowledge Gap

LFMs often lack fine-grained expertise in specialized domains, making them less reliable for tasks requiring professional or private-domain knowledge.



❑ Ungrounded Hallucination

Without access to external evidence, LFM may generate plausible but unsupported answers, especially for knowledge-intensive or factual queries.



❑ Limited Traceability & Verifiability

LFMs usually do not provide explicit sources for their responses, making it difficult to verify, audit, or trust the generated content.

LFMs' Challenges in Vertical Domains

❑ Domain of Law

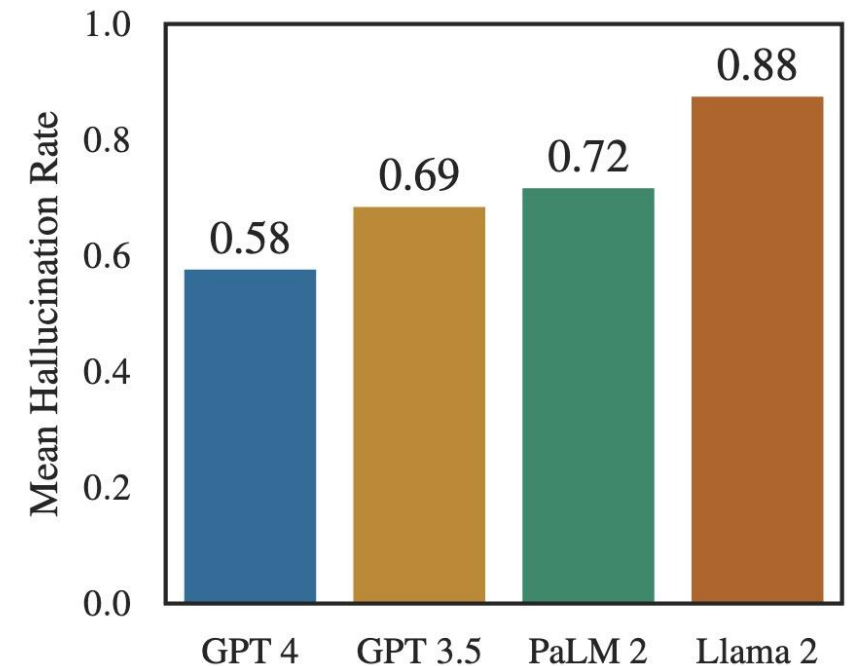
 *Journal of Legal Analysis*, 2024, 16, 64–93
<https://doi.org/10.1093/jla/laae003>
Advance access publication 26 June 2024
Article

Large Legal Fictions: Profiling Legal Hallucinations in Large Language Models

Matthew Dahl¹, Varun Magesh[†], Mirac Suzgun[‡], and Daniel E. Ho[§]

*In a new study by **Stanford RegLab** and **Institute for Human-Centered AI** researchers, it is demonstrated that legal hallucinations are pervasive and disturbing: **hallucination rates range from 69% to 88% in response to specific legal queries** for state-of-the-art language models.*

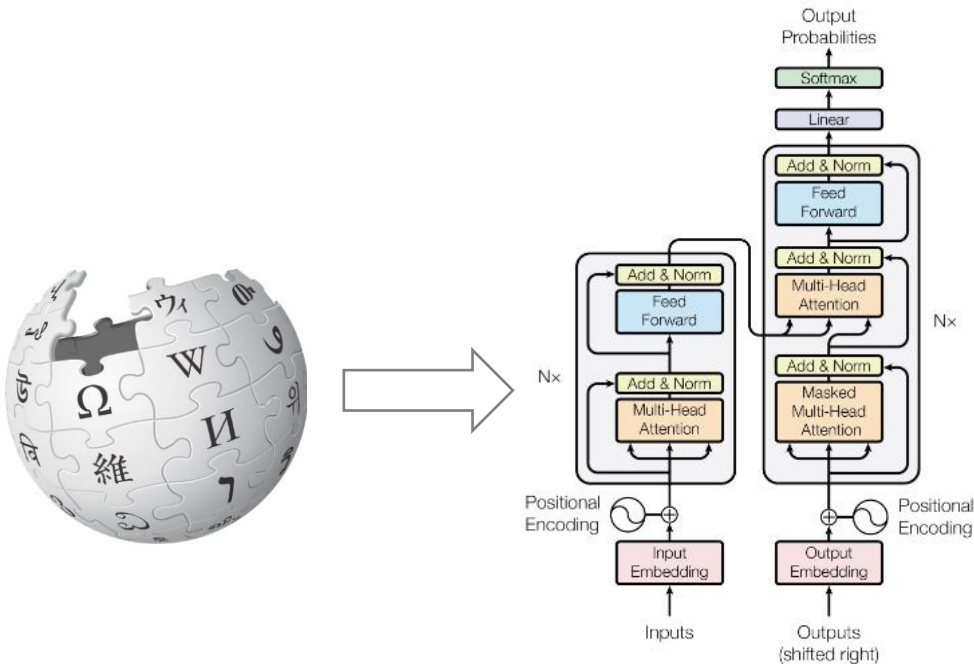
Hallucinations are common across all LFMs when they are asked a direct, verifiable question about a federal court case



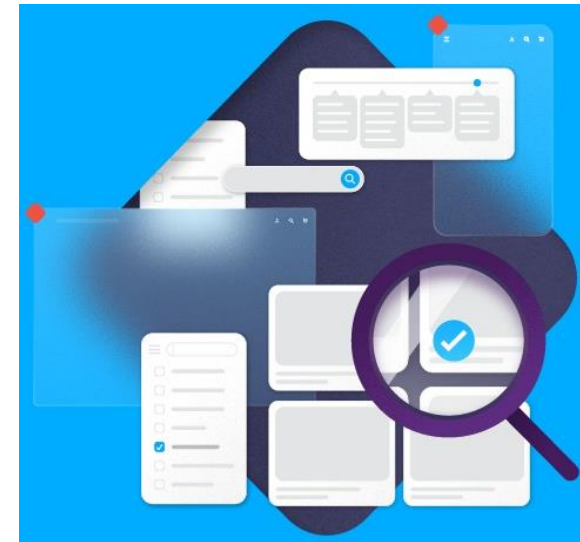
Why Large Foundation Models Work Well?

- ❑ Big Model + Big Training Data

Storing knowledge in the parametric model !



Storing knowledge in the non-parametric model?

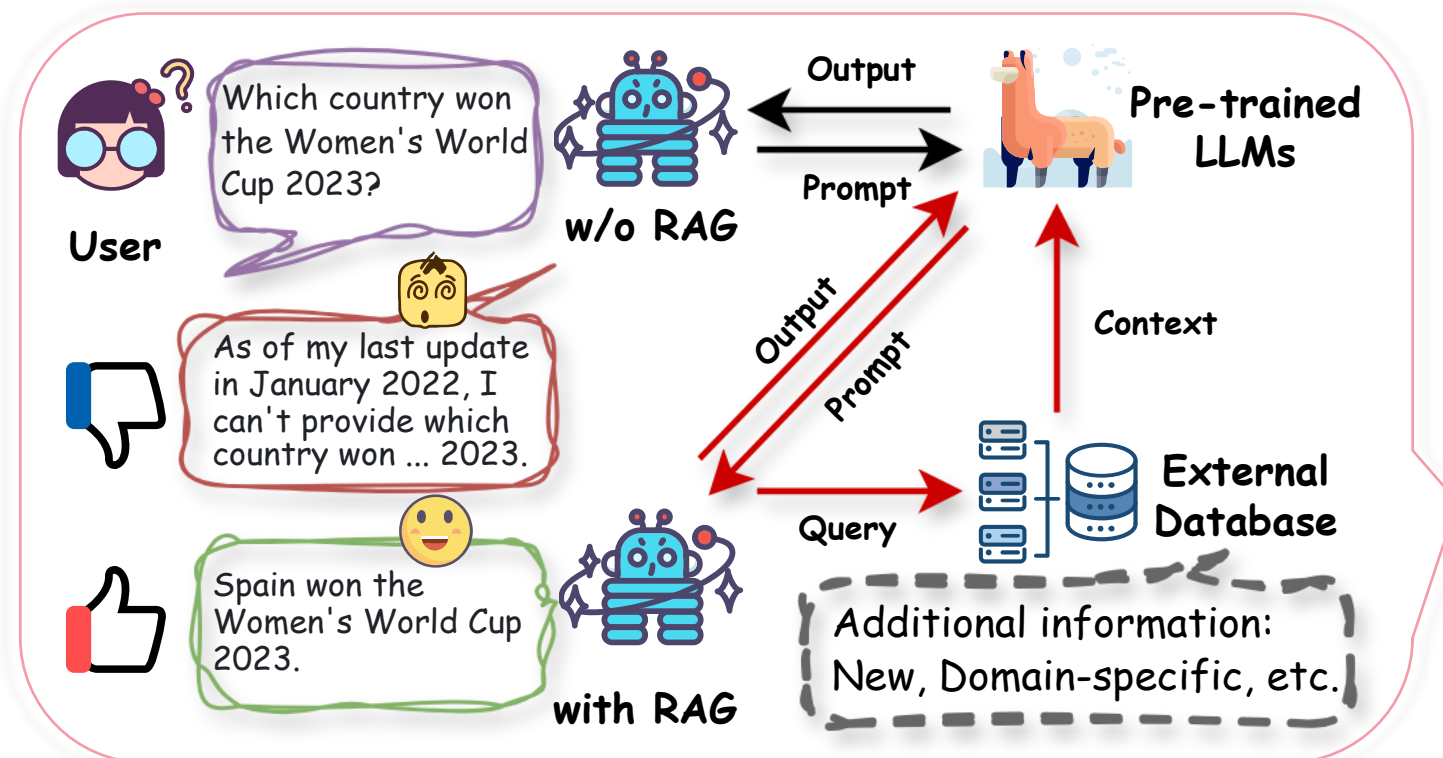


Information Retrieval (IR)

Retrieval-Augmented Large Foundation Models (RA-LFMs)

- ❑ LLMs **cannot memorize all** (particularly long-tail) knowledge in their parameters
- ❑ Lack of **domain-specific knowledge, updated information**, etc

Hallucination & Unable to answer → Re-training / Finetuning ?



Costly & Heavy Work

Retrieval-Augmented Generation (RAG) for LLM:
RA-LFMs

Integrating Information Retrieval in Generation: RA-LFM

Data for Training LFM

- Low quality
- General
- Fixed
- Hard to update

External Knowledge Base

- High-quality knowledge
- Specialized knowledge
- Scalable
- Easy-updated



Content generation

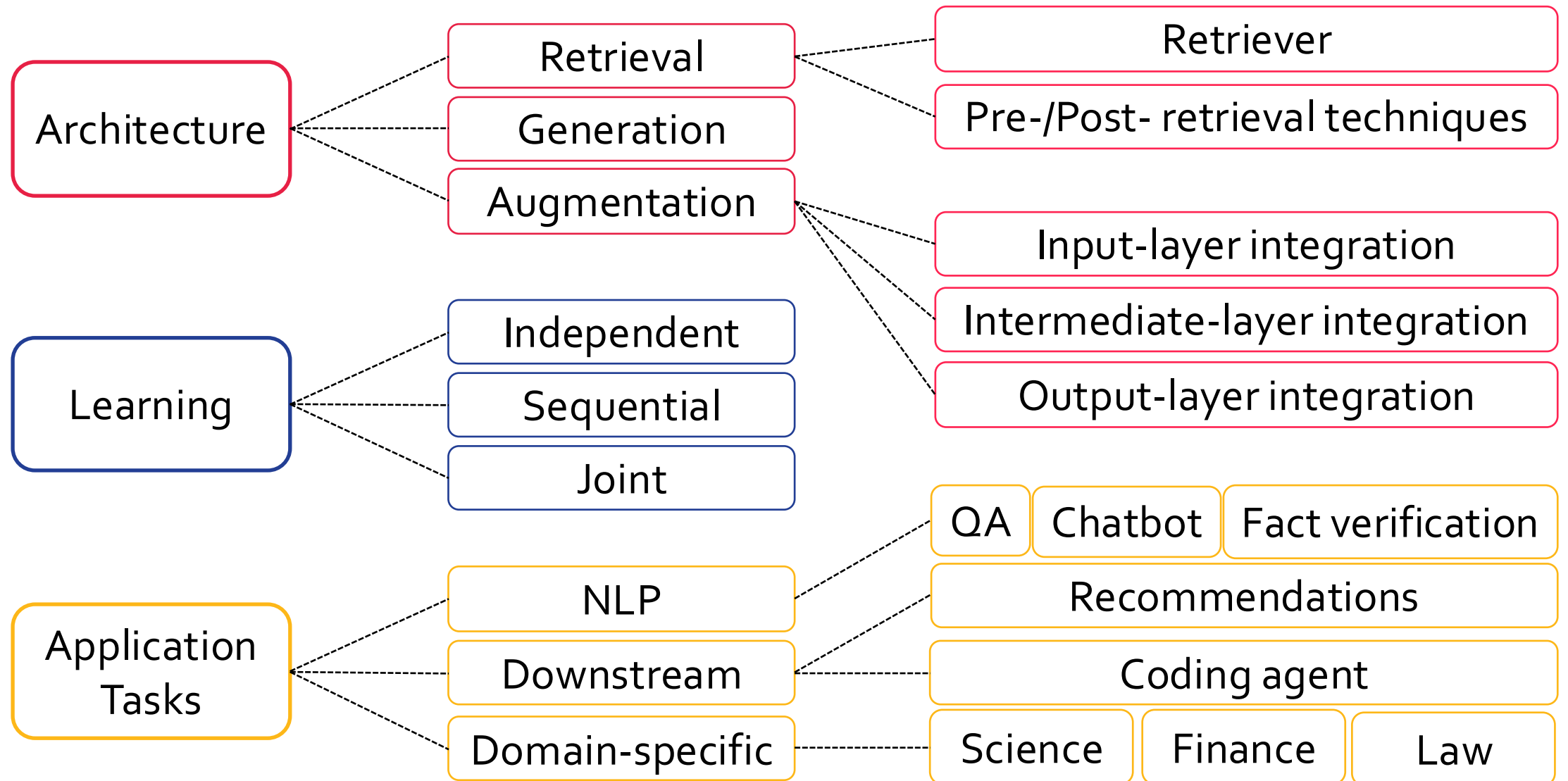
*Close-book exam
(Hard mode, have to
remember everything)*

Information / Knowledge
retrieval






RA-LFM

*Open-book exam
(Easy mode, allow to search
in reference)*

RA-LFM Research Taxonomy



RAG & RA-LFM Model Development

Category	2020 and before	2021	2022	2023	2024	2025	2026
 RAG Framework	kNN-LM REALM RAG	FiD	RETRO Atlas	RePlug	RAPTOR GraphRAG Speculative RAG LightRAG LongRAG	ArchRAG NodeRAG OpenRAG GFM-RAG GraphRAG-R1	Deep GraphRAG EA-GraphRAG Retrieval-as-Generation
 RAG Learning		EMDR2	RAG-end2end	Self-RAG	RAFT RankRAG	GraphRAFT R3-RAG	
 Retriever Learning	DPR ColBERT	Contriever GTR	ColBERTv2 E5	BGE	BGE-M3 NV-Embed	Qwen3 Embedding ReasonIR	
 Pre-, Post-Retrieval Technique				HyDE Query2doc	RAG-Fusion CRAG Adaptive-RAG Astute RAG	RAGRouter SKILL-RAG	SURE-RAG
 Agentic RAG					PlanRAG Plan*RAG Auto-RAG	Search-o1 Search-R1 RAG-Gym ReasonRAG ARAG	A-RAG AgenticRAG LatentRAG JADE DR-RAG / Doctor-RAG AutoSearch

A Comprehensive Survey Paper

A Survey on RAG Meeting LLMs: Towards Retrieval-Augmented Large Language Models

Wenqi Fan

wenqifan03@gmail.com
The Hong Kong Polytechnic
University, HK SAR

Yujuan Ding*

dingyujuan385@gmail.com
The Hong Kong Polytechnic
University, HK SAR

Liangbo Ning

BigLemon1123@gmail.com
The Hong Kong Polytechnic
University, HK SAR

Shijie Wang

shijie.wang@connect.polyu.hk
The Hong Kong Polytechnic
University, HK SAR

Hengyun Li

neilhengyun.li@polyu.edu.hk
The Hong Kong Polytechnic
University, HK SAR

Dawei Yin

yindawei@acm.org
Baidu Inc, China

Tat-Seng Chua

dcscts@nus.edu.sg
National University of Singapore,
Singapore

Qing Li

csqli@comp.polyu.edu.hk
The Hong Kong Polytechnic
University, HK SAR



Website of this tutorial
Check out the slides and more information! →



Tutorial Outline



- **Part 1: Introduction** of Retrieval Augmented Large Foundation Models (RA-LFMs) (Dr. Wenqi Fan)
- ⦿ **Part 2: Architecture of RA-LFMs and Main Modules** (Xu Yuan)
- **Part 3: Learning** Approach of RA-LFMs (Chengliang Liu)
- **Part 4: Agentic RAG** (Chengliang Liu)
- **Part 5: Applications** of RA-LFMs (Chun-Hin Chan)
- **Part 6: Challenges and Future Directions** of RA-LFMs (Dr. Wenqi Fan)
- **Part 7: Q&A**

Website of this tutorial
Check out the slides and more information!



PART 2: Architecture of RA-LFMs and Main Modules



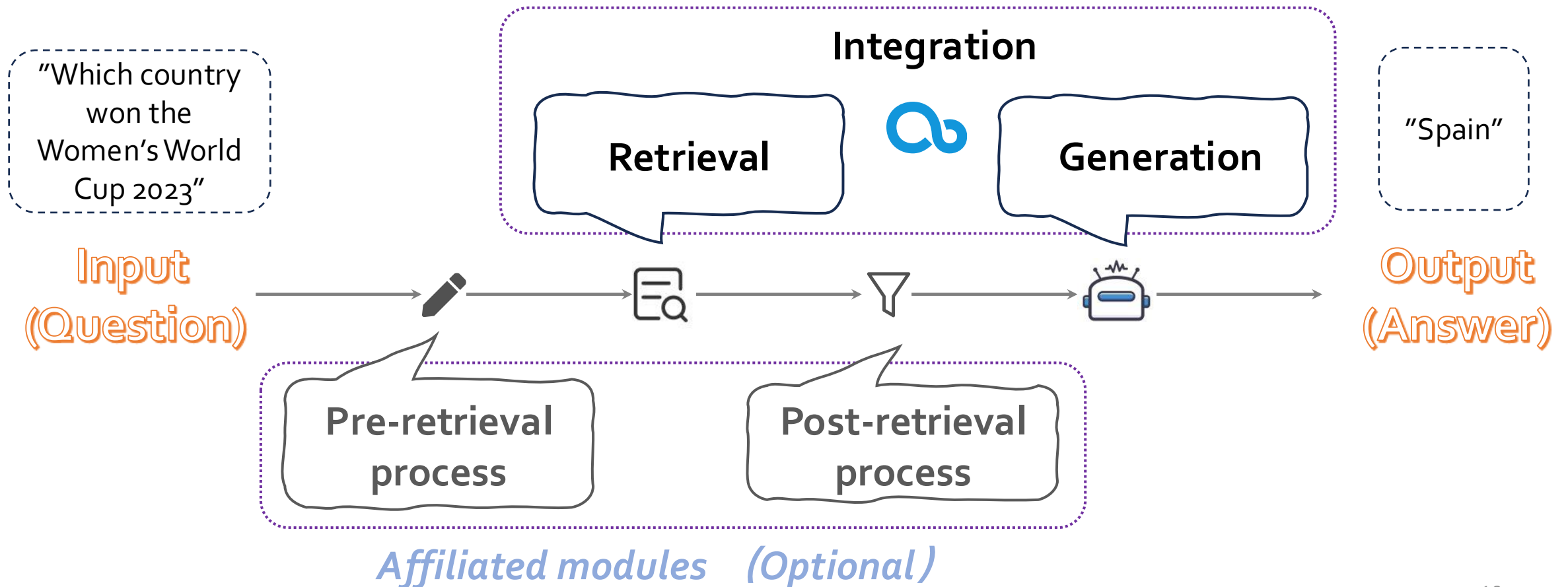
Presenter
Xu Yuan
HK PolyU

- **RA-LFM architecture overview**
- **Retriever in RA-LFMs**
- **Retrieval results integration**
- **Pre/Post-retrieval techniques**
- **Special RA-LFM paradigms**

RA-LFM Architecture: Standard Pipeline

- Technical component illustration in a RA-LFM for the Q&A task

Major components (necessary)



A Simple Retrieval-Augmented Generation Model

□ RAG

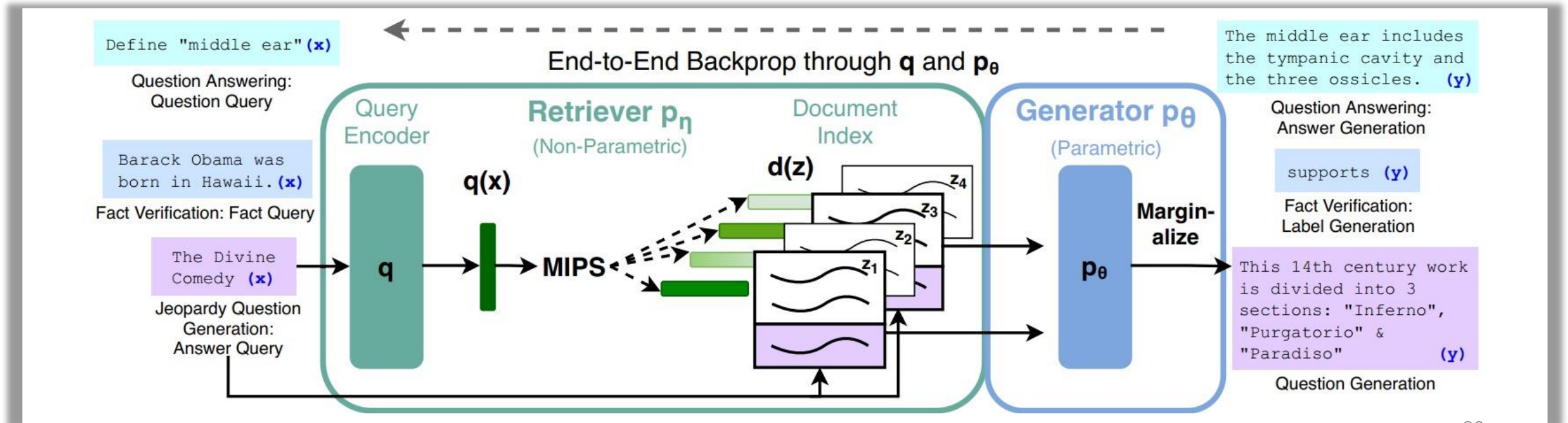
Integration: concatenating each retrieved passage with the question



Input



Output



A Simple Retrieval-Augmented Generation Model

□ In-Context RALM

Prompt-level Integration: prepending the retrieved passage with the input text



PART 2: Architecture of RA-LFMs and Main Modules



Slides



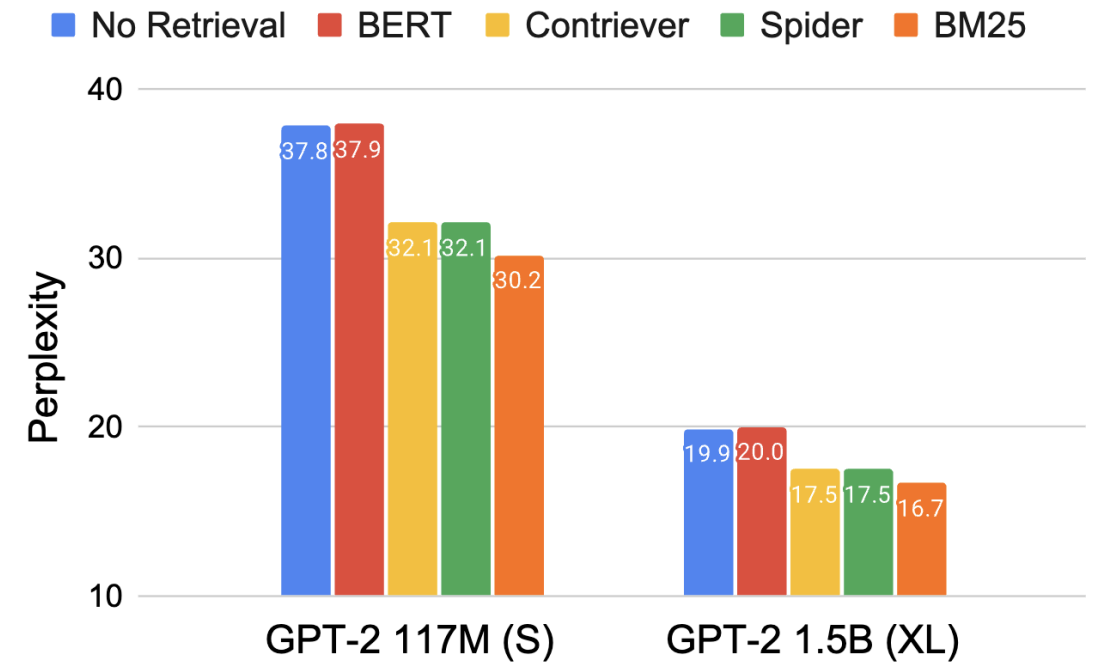
Website of this tutorial

- RA-LFM architecture overview
- **Retriever in RA-LFMs**
- Retrieval results integration
- Pre/Post-retrieval techniques
- Special RA-LFM paradigms

RA-LFM Architecture: Retriever Types

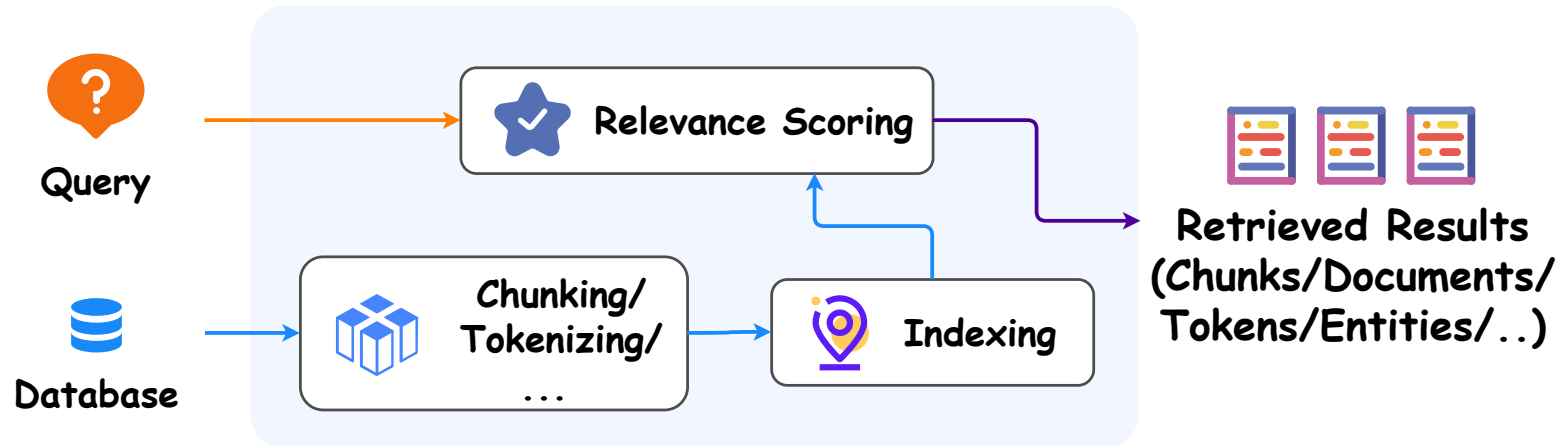
- Different types of retriever deliver different generation performance

Relevance measurement	Retriever learning
Sparse	Task-specific Learning
Dense	General-purpose Learning

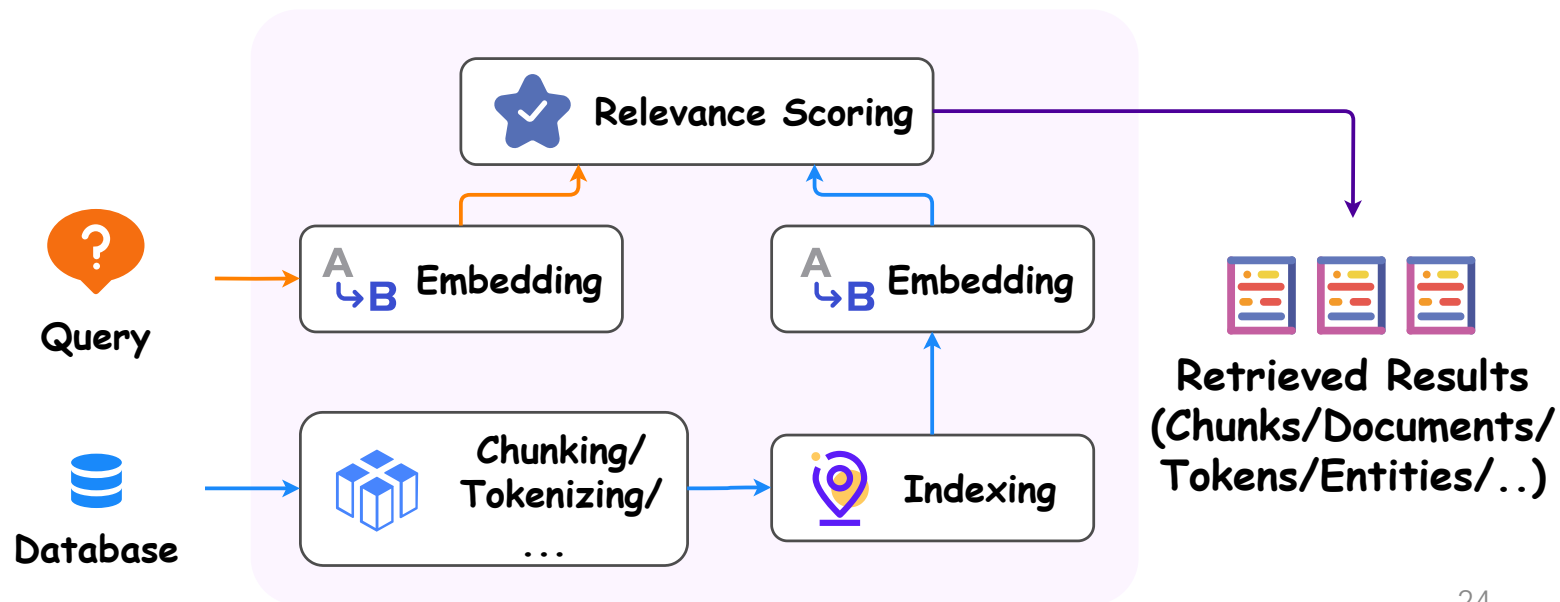


Dense v.s. Sparse Retrievers

Sparse Retrievers (SR)



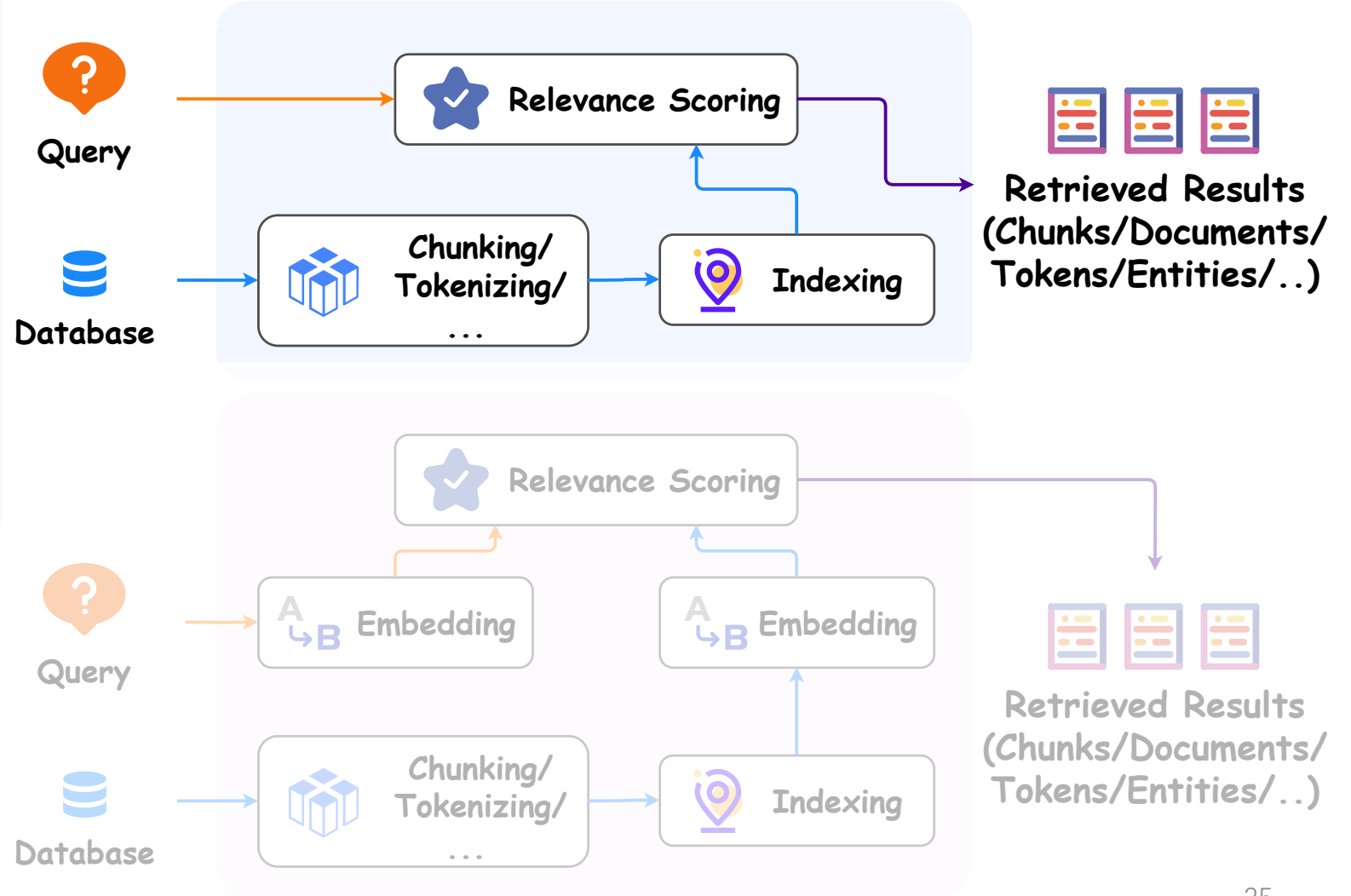
Dense Retrievers (DR)



Dense v.s. Sparse Retrievers

Sparse Retrievers (SR)

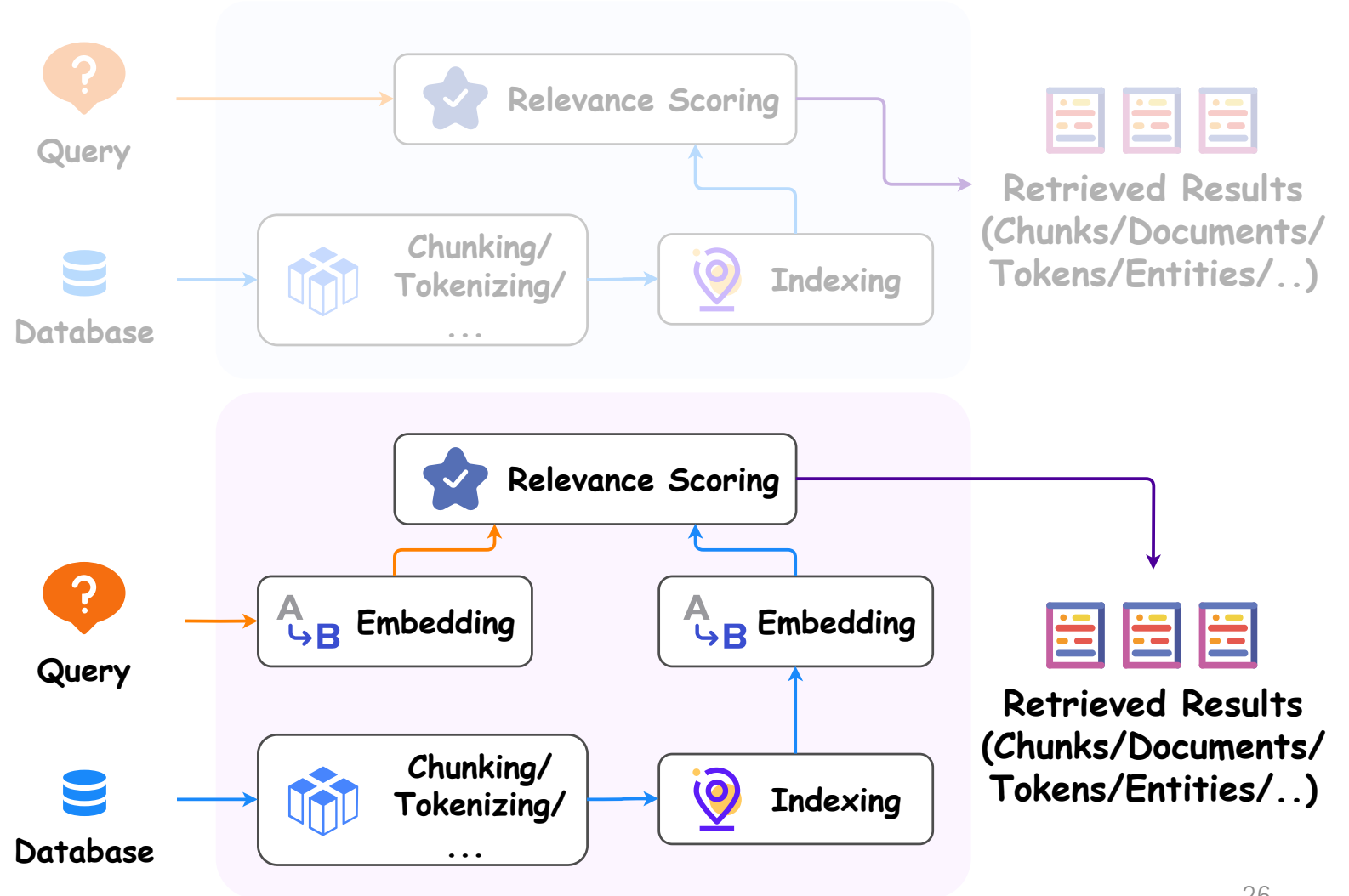
- Feasible to apply
- High efficiency
- Fine performance
- Example: TF-IDF, BM25



Dense v.s. Sparse Retrievers

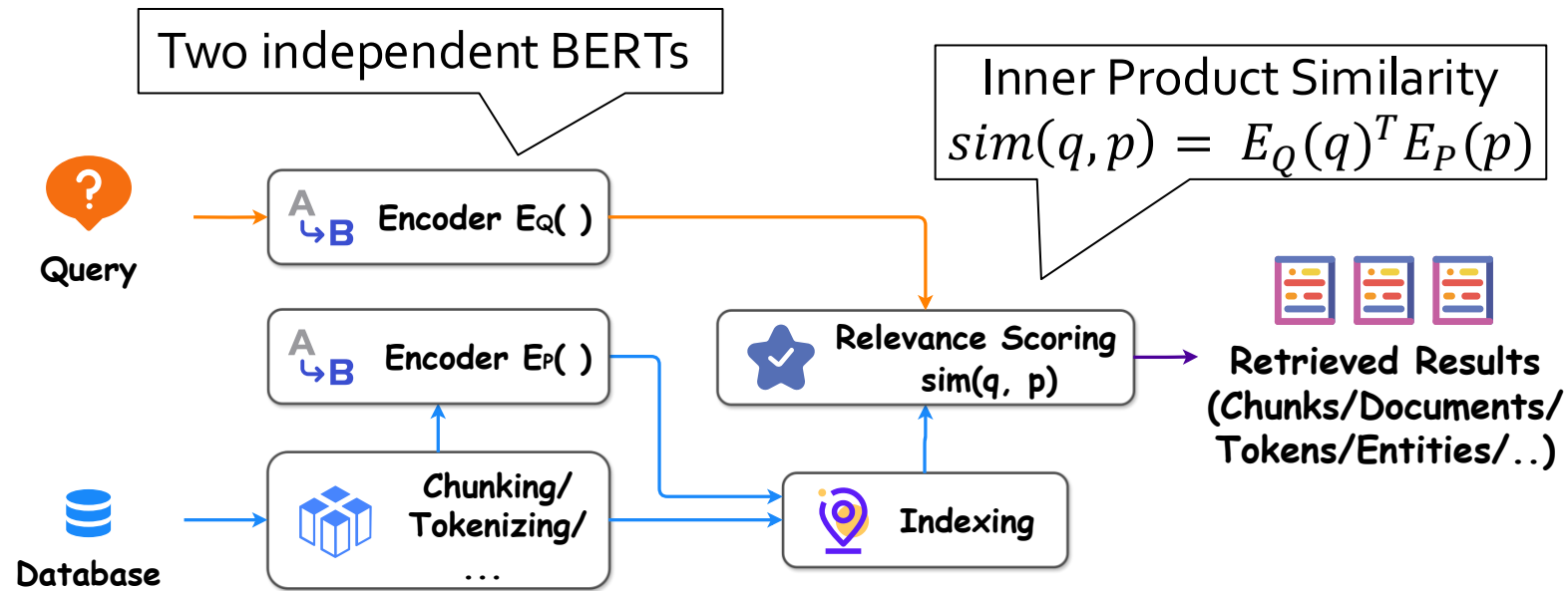
Dense Retrievers (DR)

- Allowing fine-tuning
- Better adaptation
- Customizable for more retrieval goals
- Example: DPR, Contriever, ColBERT



Task-Specific Retriever (Supervised)

- ❑ **Dense Passage Retriever (DPR):** Pretrained for Question Answering (QA)



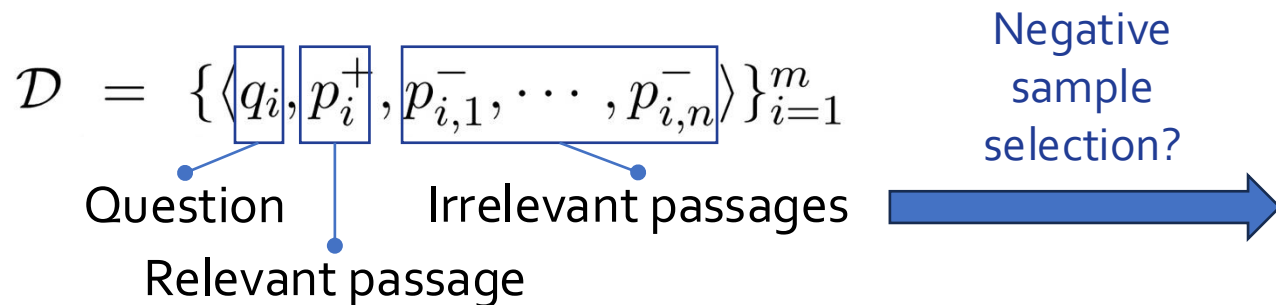
Task-Specific Retriever (Supervised)

❑ Dense Passage Retriever (DPR): Pretrained for Question Answering (QA)

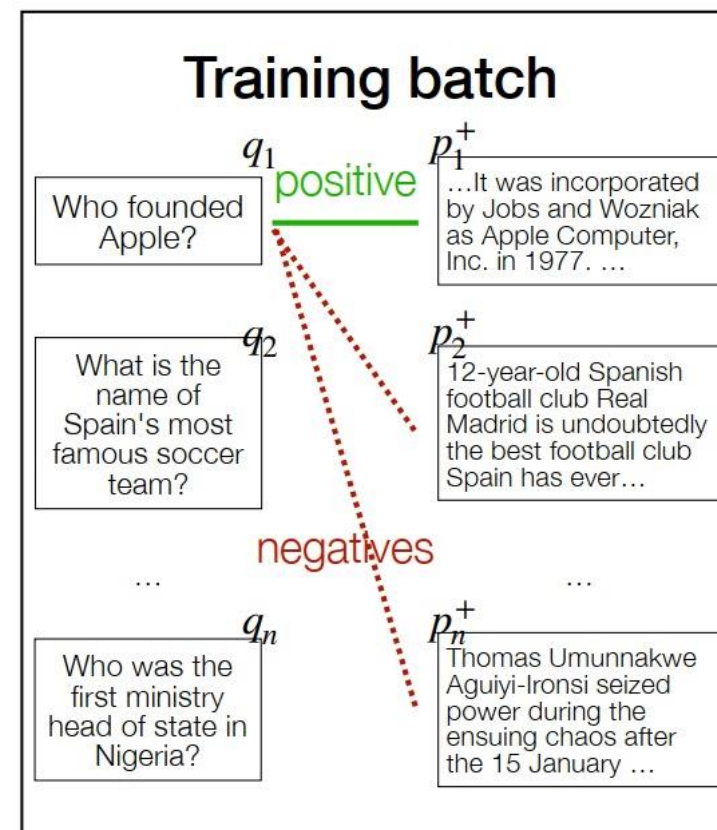
- Learning Objective

$$L(q_i, p_i^+, p_{i,1}^-, \dots, p_{i,n}^-)$$
$$= -\log \frac{e^{\text{sim}(q_i, p_i^+)}}{e^{\text{sim}(q_i, p_i^+)} + \sum_{j=1}^n e^{\text{sim}(q_i, p_{i,j}^-)}}$$

- Training data: Question-Passage Sets

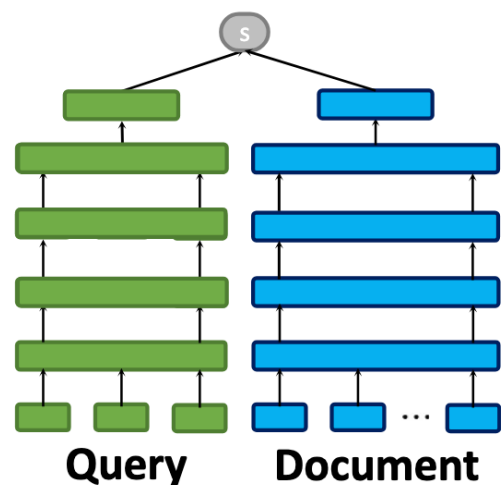


- Training with in-batch negatives

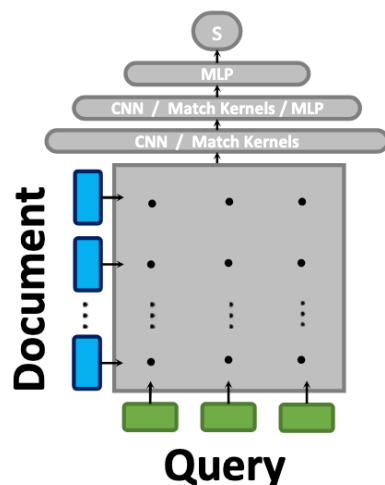


Task-Specific Retriever (Supervised)

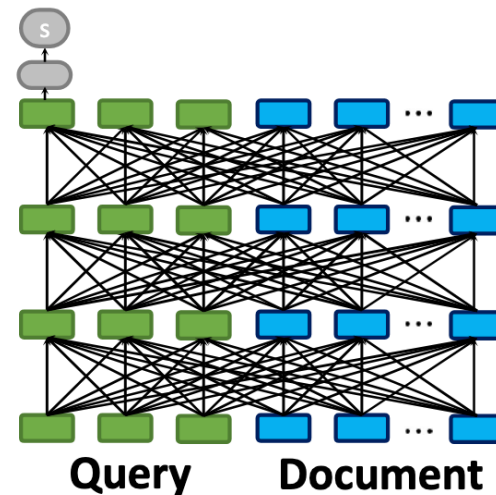
CoBERT & CoBERTv2 : Late Interaction Retrieval



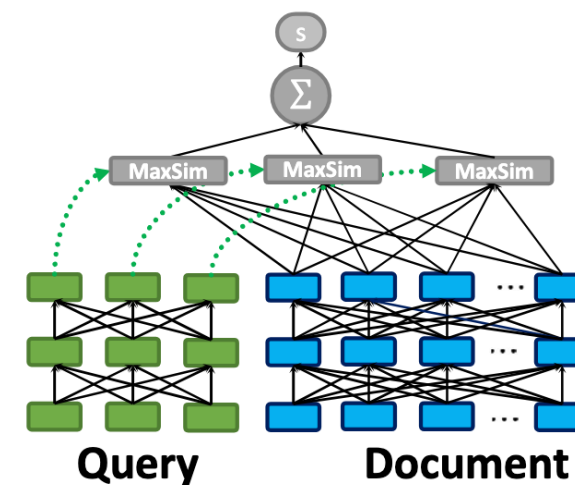
(a) Representation-based Similarity
(e.g., DSSM, SNRM)



(b) Query-Document Interaction
(e.g., DRMM, KNRM, Conv-KNRM)



(c) All-to-all Interaction
(e.g., BERT)



(d) Late Interaction
(i.e., the proposed CoBERT)

- Late Interaction

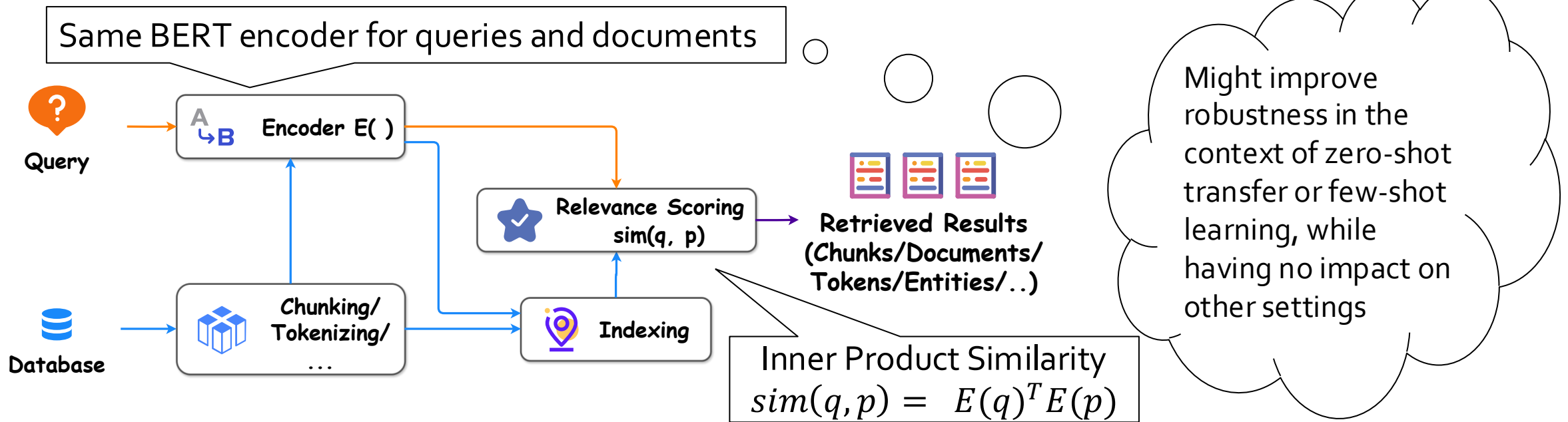
$$LI(q, d) = \sum_{i \in [1, N_q]} \max_{j \in [1, N_d]} \langle \mathbf{E}_q^{(i)} | \mathbf{E}_d^{(j)} \rangle$$

- Contrastive Loss

$$\mathcal{L} = -\frac{1}{b} \sum_{k=1}^b \log \left[\frac{\exp(s_k^+)}{\exp(s_k^+) + \exp(s_k^-)} \right]$$

General-Purpose Retriever (Unsupervised)

- ❑ **Contriever:** Pre-trained with unsupervised learning
 - ❖ supervised pretrained methods do not transfer well to new applications with no training data



Contrastive learning with **unaligned documents**

$$\mathcal{L}(q, k_+) = - \frac{\exp(s(q, k_+)/\tau)}{\exp(s(q, k_+)/\tau) + \sum_{i=1}^K \exp(s(q, k_i)/\tau)}$$

k_+ : positive pairs from a single document
 $(k_i)_{1,2,\dots,K}$: negative pairs with in-batch negative sampling and MoCo

DPR & Contriever Performance on OpenQA Tasks

End-to-end QA (Exact Match) Accuracy

Training	Model	NQ	TriviaQA	WQ	TREC	SQuAD
Single	BM25+BERT (Lee et al., 2019)	26.5	47.1	17.7	21.3	33.2
Single	ORQA (Lee et al., 2019)	33.3	45.0	36.4	30.1	20.2
Single	HardEM (Min et al., 2019a)	28.1	50.9	-	-	-
Single	GraphRetriever (Min et al., 2019b)	34.5	56.0	36.4	-	-
Single	PathRetriever (Asai et al., 2020)	32.6	-	-	-	56.5
Single	REALM _{Wiki} (Guu et al., 2020)	39.2	-	40.2	46.8	-
Single	REALM _{News} (Guu et al., 2020)	40.4	-	40.7	42.9	-
Single	BM25	32.6	52.4	29.9	24.9	38.1
	DPR	41.5	56.8	34.6	25.9	29.8
	BM25+DPR	39.0	57.0	35.2	28.0	36.7
Multi	DPR	41.5	56.8	42.4	49.4	24.1
	BM25+DPR	38.8	57.9	41.1	50.6	35.8

Both widely applied in
RAG and RA-LFMs

DPR in
RAG, FiD, RETRO,
EPR, UDR, ...

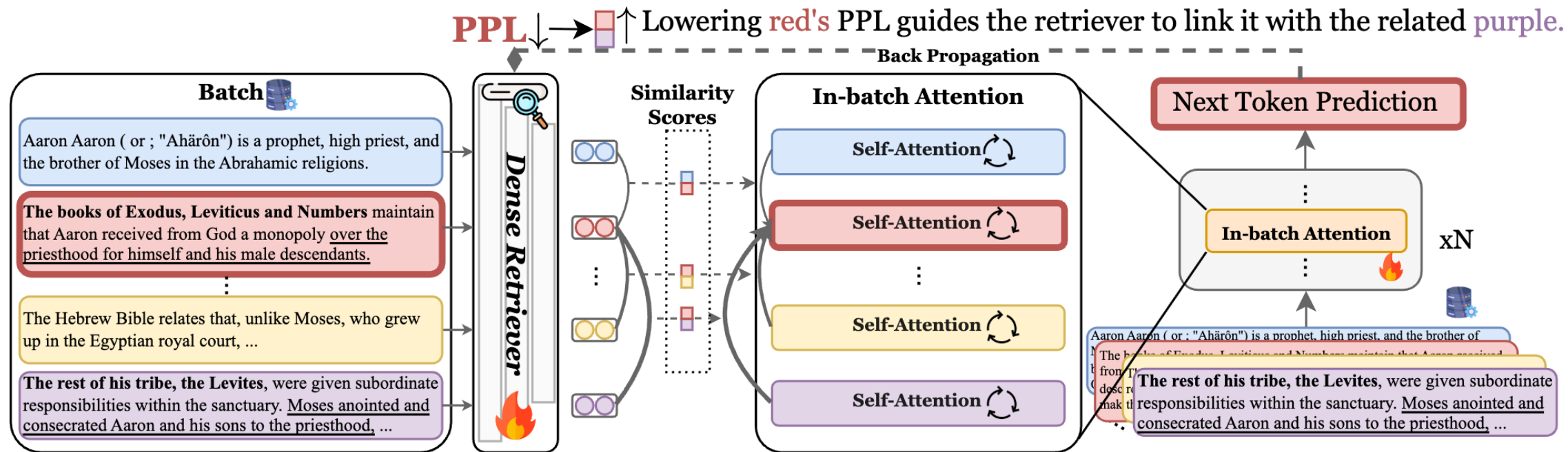
Contriever in
Self-RAG, Atlas,
RAVEN, ...

	NaturalQuestions			TriviaQA		
	R@5	R@20	R@100	R@5	R@20	R@100
Inverse Cloze Task (Sachan et al., 2021)	32.3	50.9	66.8	40.2	57.5	73.6
Masked salient spans (Sachan et al., 2021)	41.7	59.8	74.9	53.3	68.2	79.4
BM25 (Ma et al., 2021)	-	62.9	78.3	-	76.4	83.2
Contriever	47.8	67.8	82.1	59.4	74.2	83.2
<i>supervised model: DPR (Karpukhin et al., 2020)</i>	-	78.4	85.4	-	79.4	85.0

Both better than
the sparse retriever!

General-Purpose Retriever (Self-Supervised)

❑ **Revela** : From contrastive self-supervision to LM-based retriever learning.

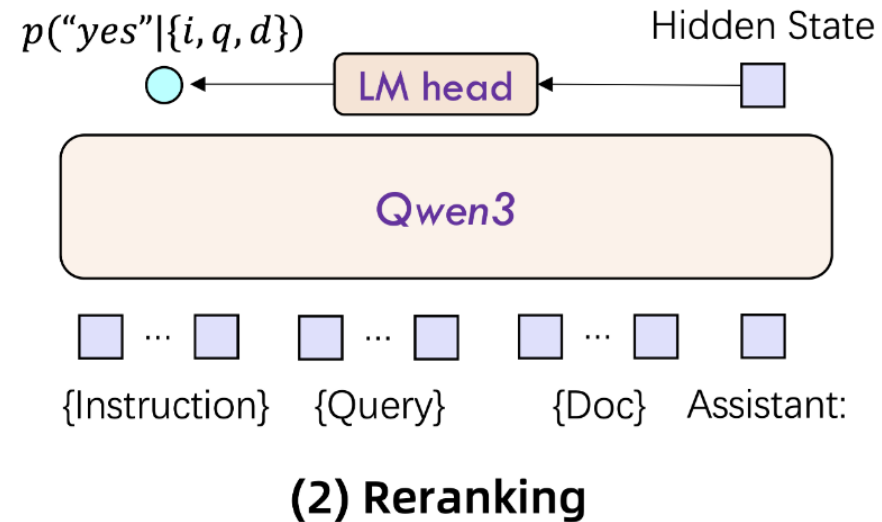
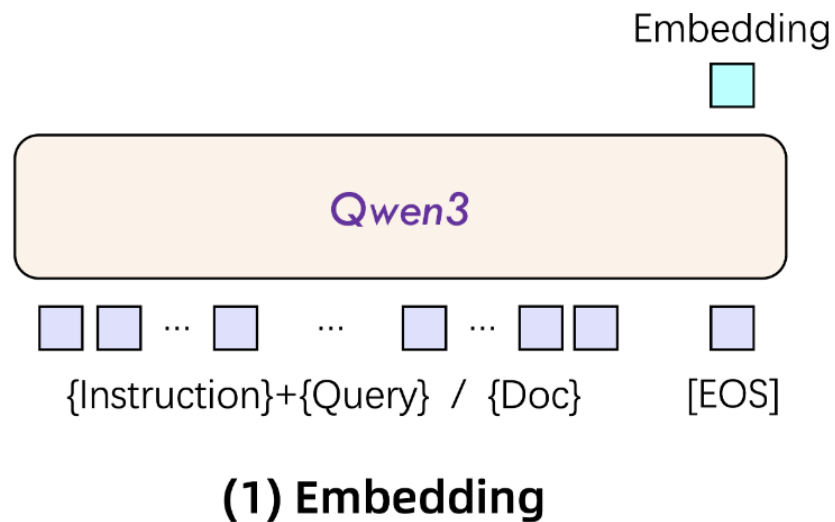


- Learns dense retrieval from raw text
- No query-document labels

- Retriever scores guide in-batch attention
- Optimized by next-token prediction

General-Purpose Retriever (LFM-based)

Qwen3 Embedding: Text Embedding based on Foundation Models



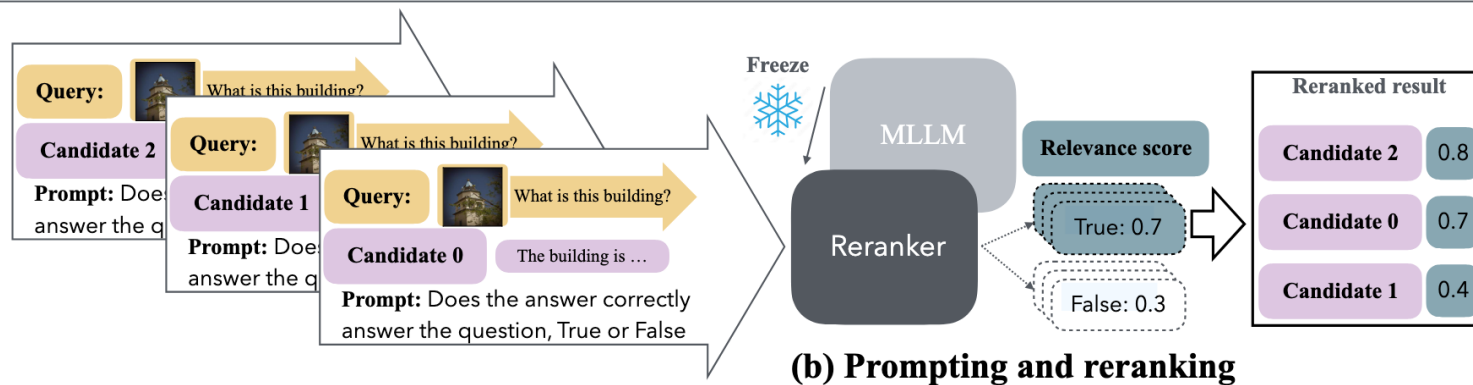
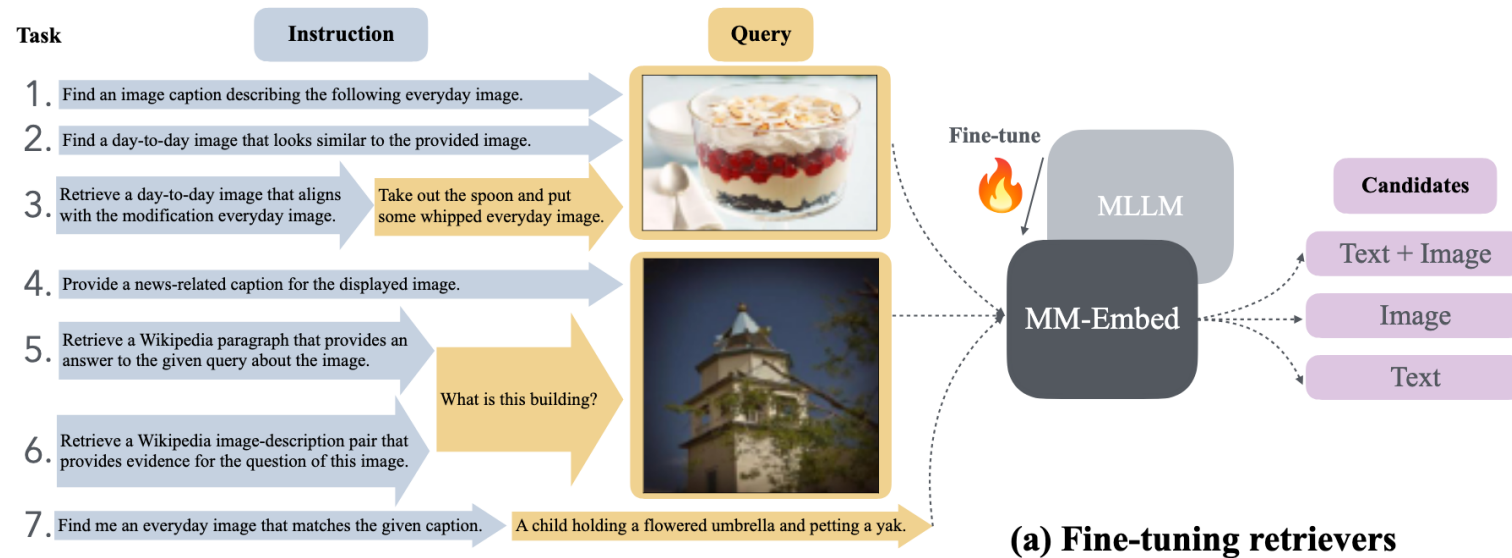
- Training Objective

$$L_{\text{embedding}} = -\frac{1}{N} \sum_i^N \log \frac{e^{(s(q_i, d_i^+) / \tau)}}{Z_i}$$

$$L_{\text{reranking}} = -\log p(l|\mathcal{P}(q, d)).$$

General-Purpose Retriever (LFM-based)

MM-EMBED: Universal Multimodal Retrieval with Instruction-aware Capability



Task

$$1. q^{\text{txt}} \rightarrow c^{\text{img}}$$

$$2. q^{\text{txt}} \rightarrow c^{\text{txt}}$$

$$3. q^{\text{txt}} \rightarrow (c^{\text{img}}, c^{\text{txt}})$$

$$4. q^{\text{img}} \rightarrow c^{\text{txt}}$$

$$5. q^{\text{img}} \rightarrow c^{\text{img}}$$

$$6. (q^{\text{img}}, q^{\text{txt}}) \rightarrow c^{\text{txt}}$$

$$7. (q^{\text{img}}, q^{\text{txt}}) \rightarrow c^{\text{img}}$$

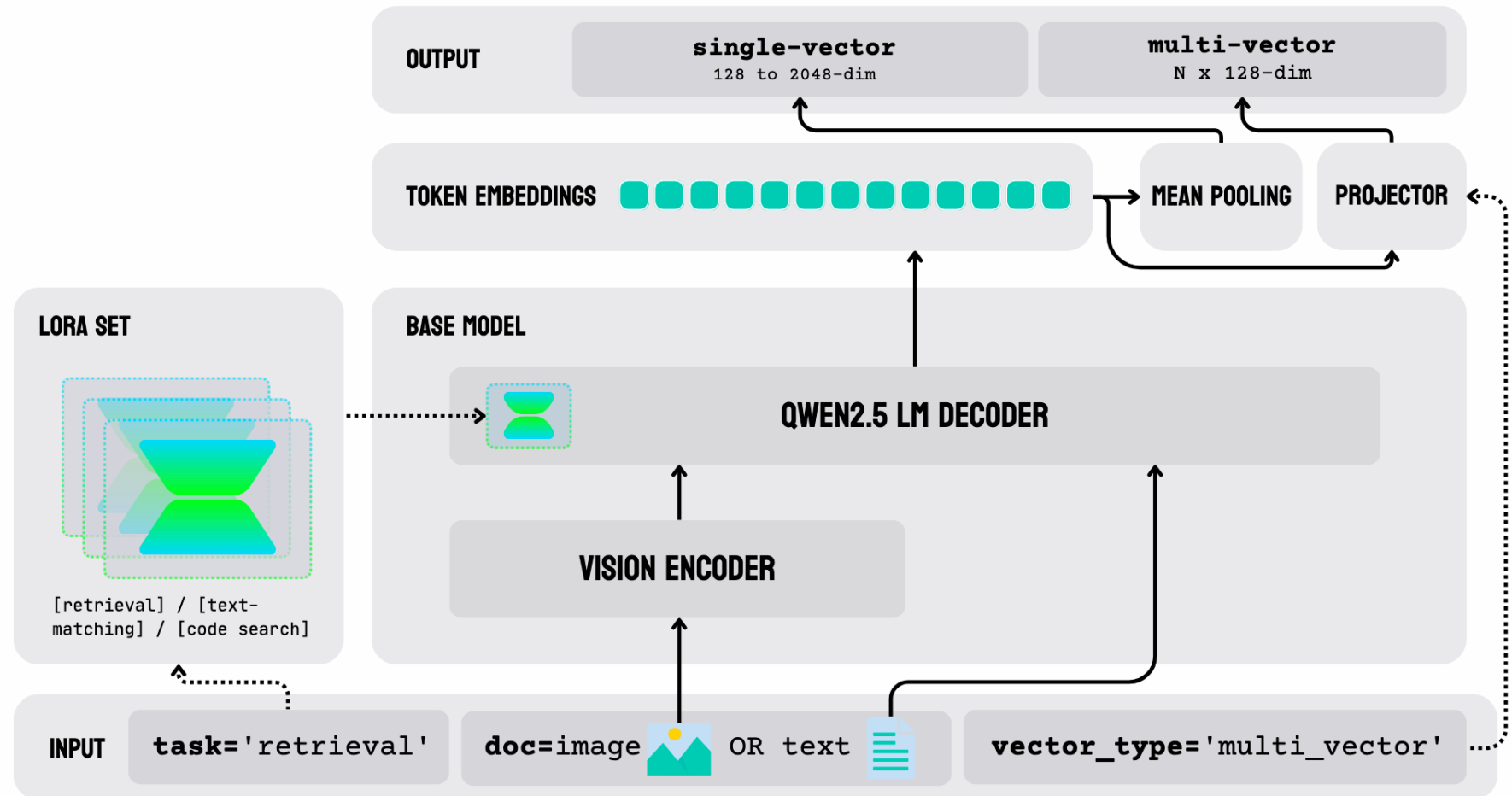
$$8. (q^{\text{img}}, q^{\text{txt}}) \rightarrow (c^{\text{img}}, c^{\text{txt}})$$

General-Purpose Retriever (LFM-based)

❑ jina-embeddings-v4: Universal Multi-task Multi-modal Embedding Model

Task-specific LoRA adapters

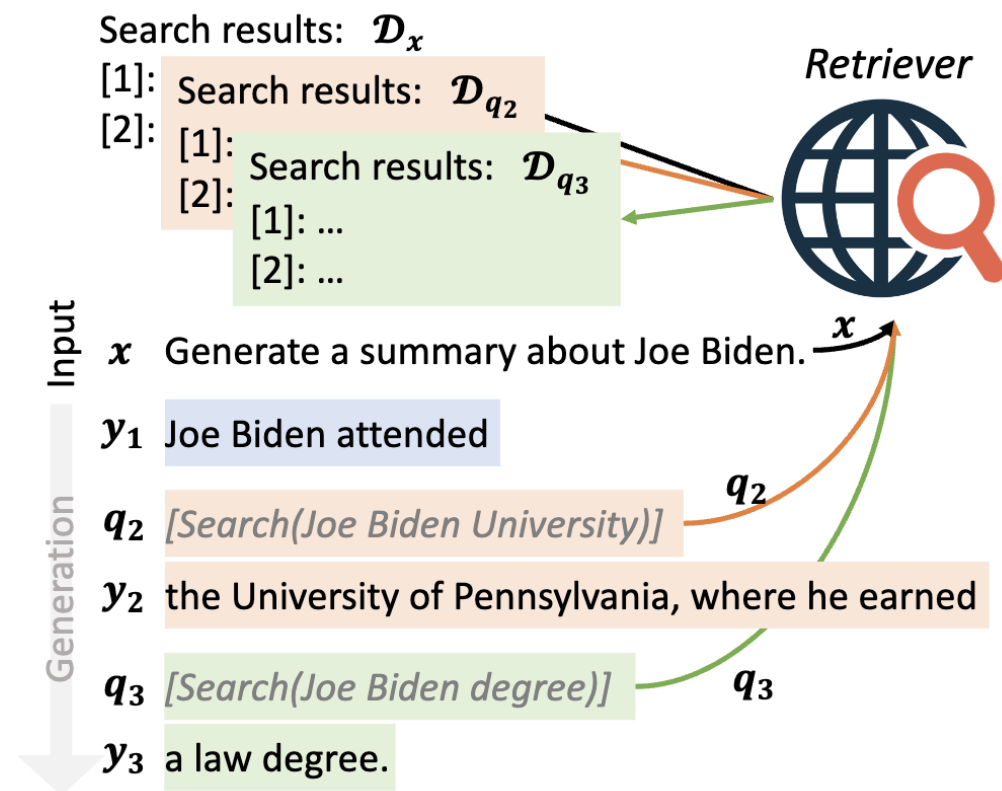
- ❖ Retrieval
- ❖ Text Matching
- ❖ Code Search
- ❖ ...



Search Engine as Retrievers

❑ Traditional retrieval methods

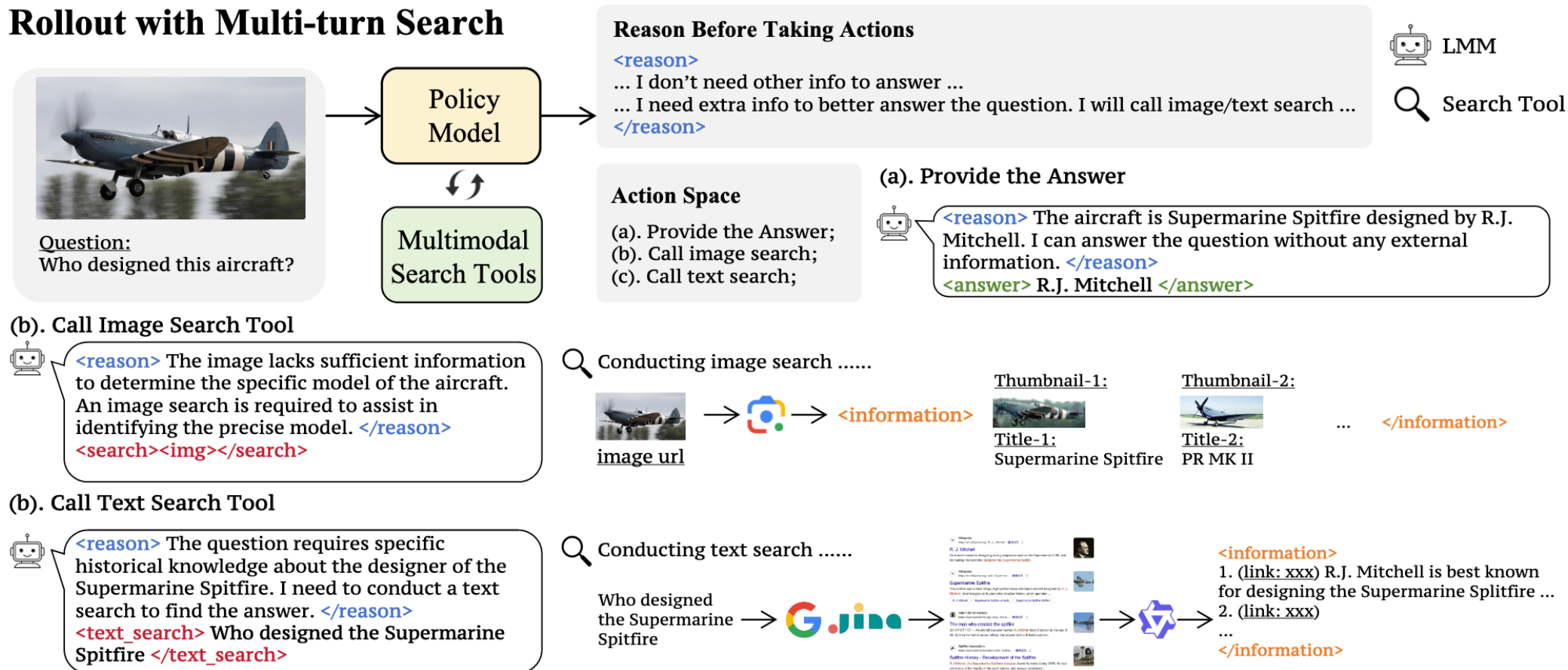
- ❖ May be difficult to update to real-time web documents
- ❖ May be a limit to the number of documents storable in the pre-defined database
- ❖ Will not take advantage of the high quality ranking that has been finely tuned in Internet Search engines over decades of use



Search Engine as Retrievers

MMSearchR1: Search engines as reasoning-time retrievers.

Rollout with Multi-turn Search



PART 2: Architecture of RA-LFMs and Main Modules



Slides

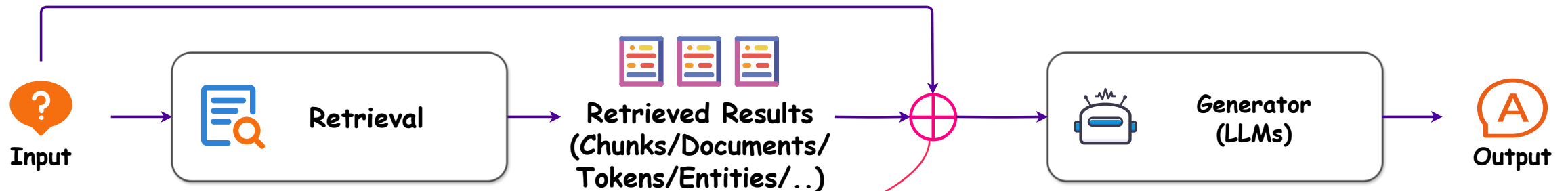


Website of this tutorial

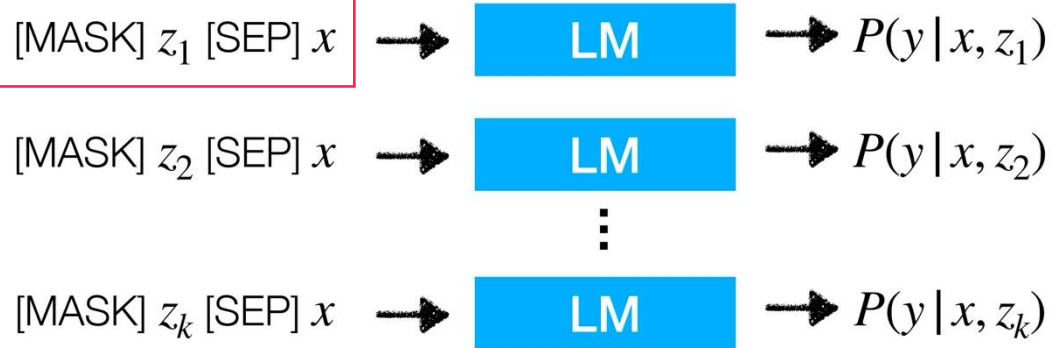
- RA-LFM architecture overview
- Retriever in RA-LFMs
- **Retrieval results integration**
- Pre/Post-retrieval techniques
- Special RA-LFM paradigms

Retrieved Results Integration: Input-layer Integration

□ REALM



Integrating the retrieved passage z and the original input x



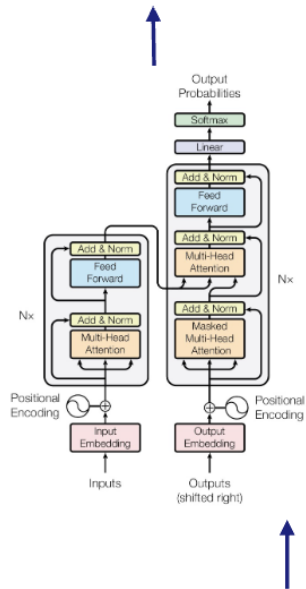
Weighted aggregating the prediction results based on all retrieved passages

$$\sum_{z \in \mathcal{D}} P(z | x) P(y | x, z)$$

Retrieval-Augmented Generator

Typical encoder: $p(y|x)$

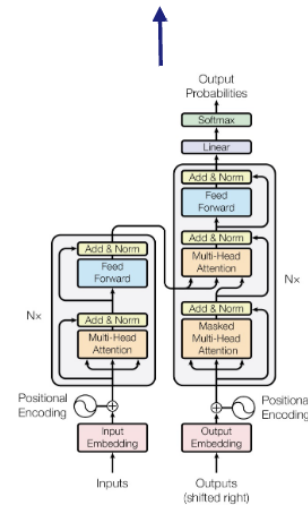
$y = \text{pounds}$



x : we paid 20 __ at the Buckingham Palace gift shop

Knowledge-augmented encoder: $p(y|x, z)$

$y = \text{pounds}$

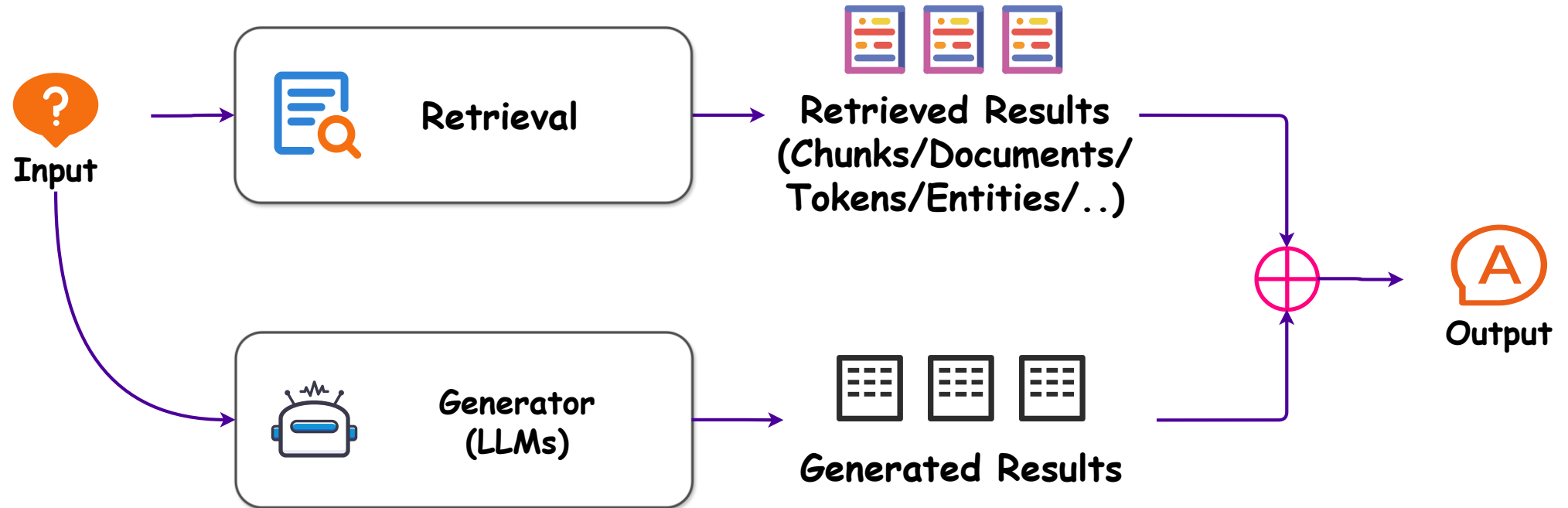


x : we paid 20 __ at the Buckingham Palace gift shop

z : Buckingham Palace is home to the British monarchy

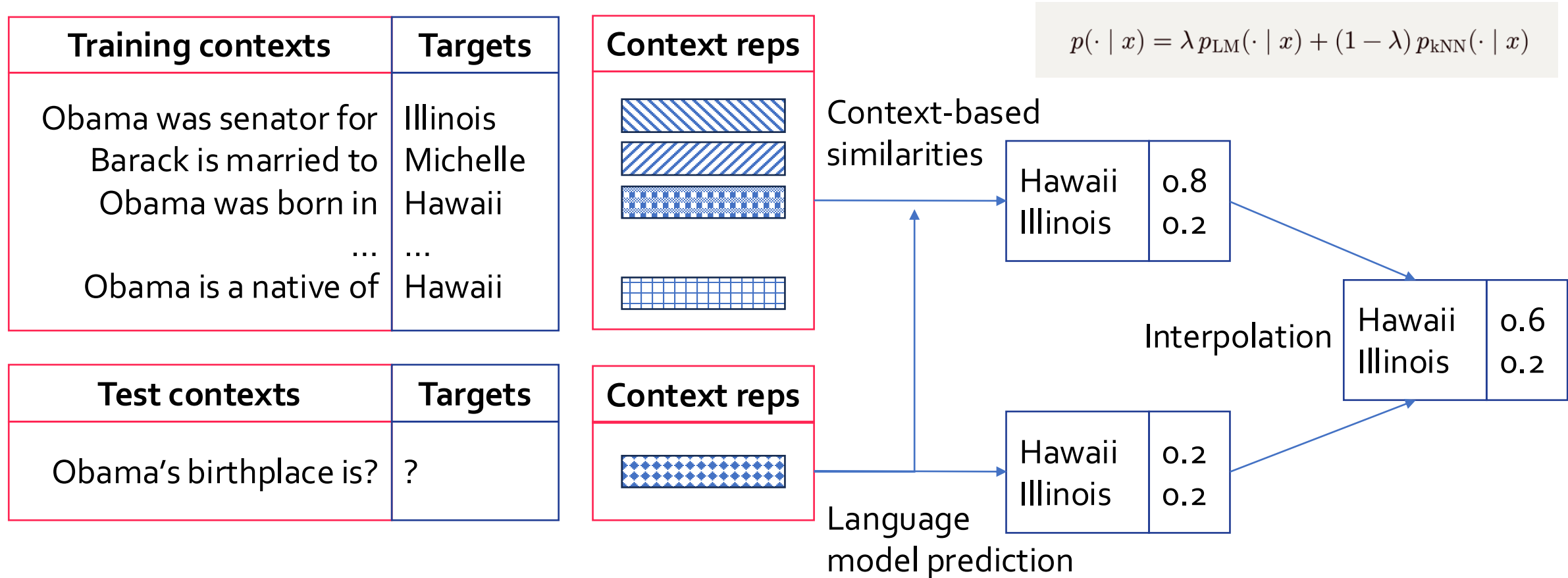
explicit knowledge

Retrieved Results Integration: Output-layer integration

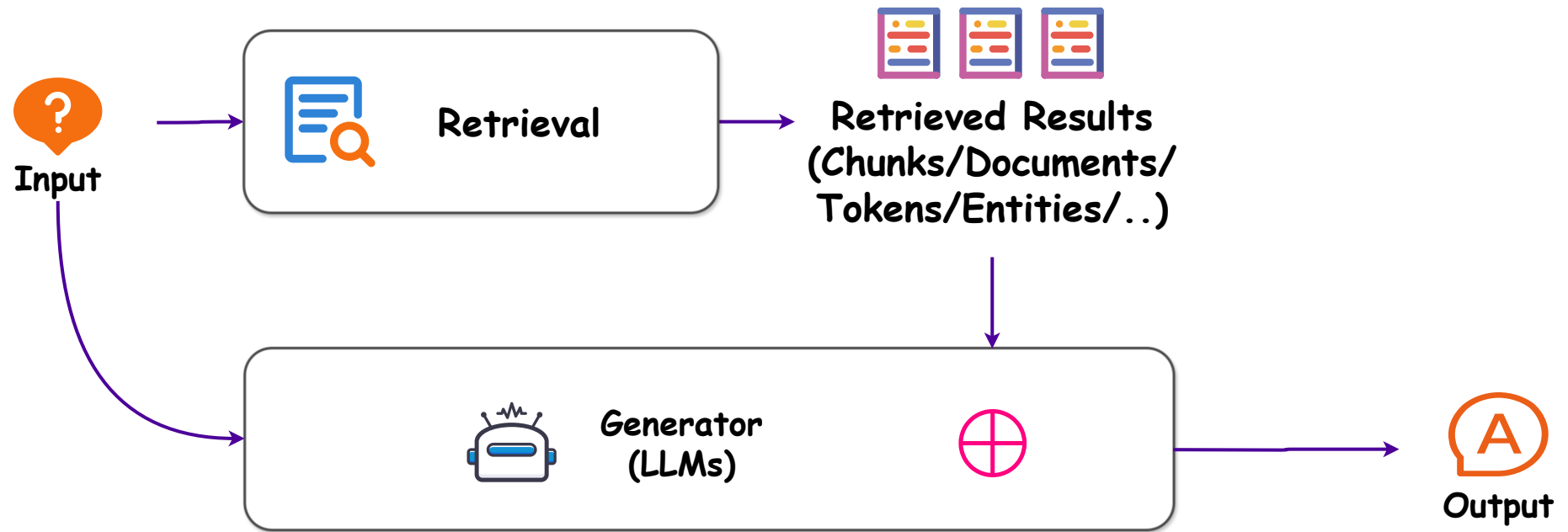


RA-LFM Architecture: Output-layer Integration

- **kNN-LM**: Combining retrieved probabilities and predicted ones in generation

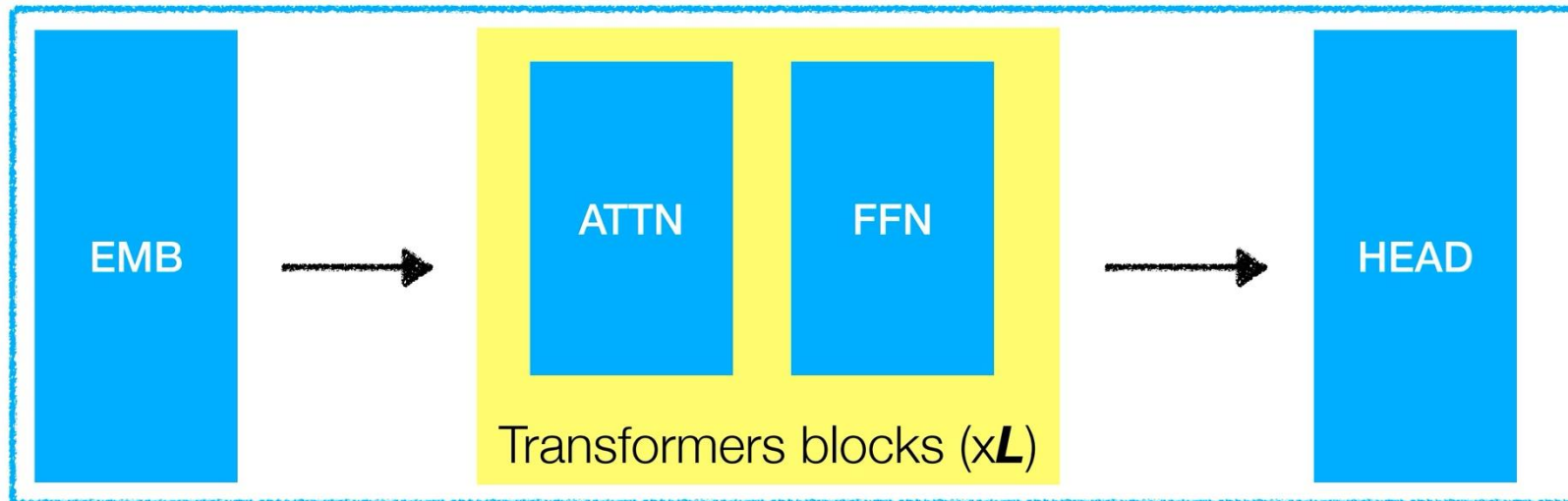


Retrieved Results Integration: Intermediate-layer Integration



Retrieved Results Integration: Intermediate-layer Integration

Regular Decoder



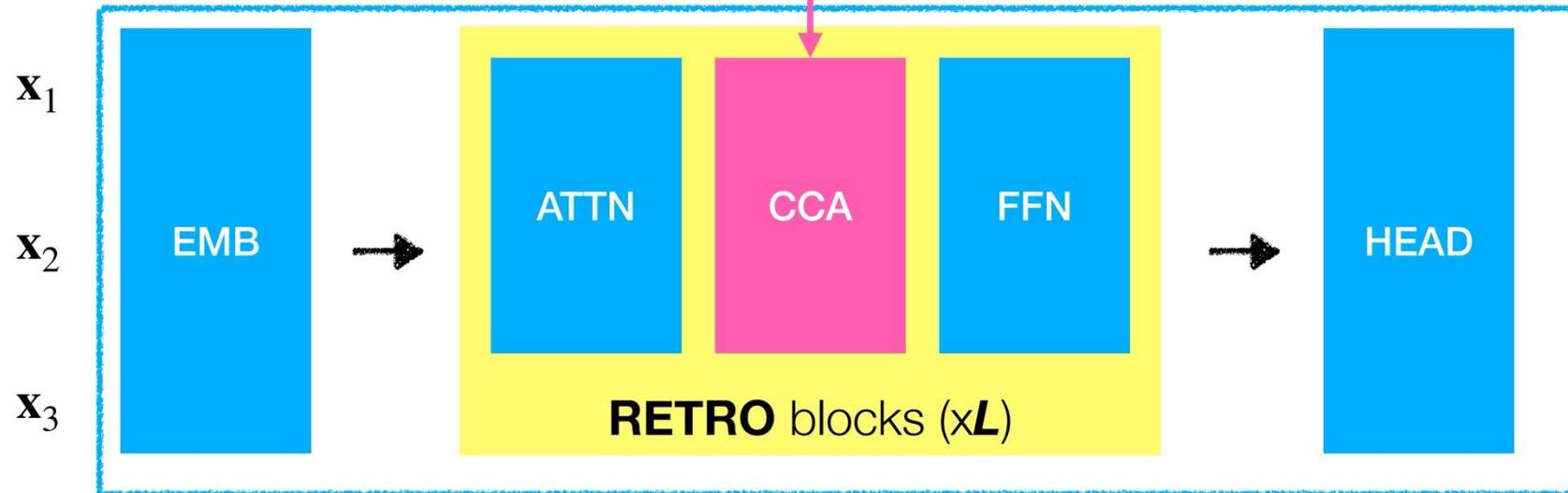
Retrieved Results Integration: Intermediate-layer Integration

Decoder to incorporate retrieved results (RETRO)

Retrieving text similar to the previous chunk to improve the predictions in the current chunk

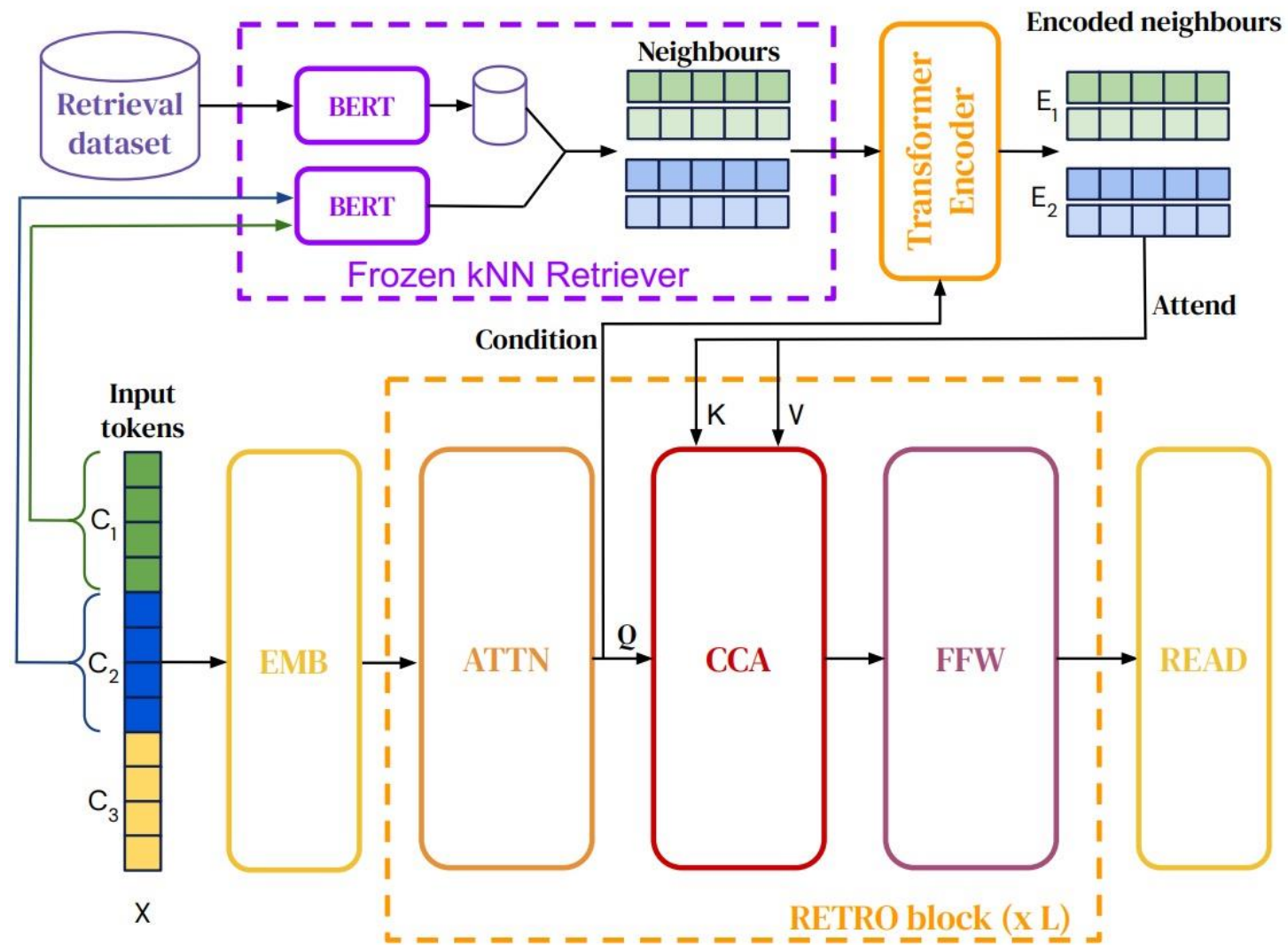
$\Rightarrow \mathbf{E}_1 \mathbf{E}_2 \mathbf{E}_3$

Splitting the input sequence into chunks



Chunked Cross Attention (CCA)

Retrieved Results Integration: Intermediate-layer Integration



PART 2: Architecture of RA-LFMs and Main Modules



Slides



Website of this tutorial

- RA-LFM architecture overview
- Retriever in RA-LFMs
- Retrieval results integration
- **Pre/Post-retrieval techniques**
- Special RA-LFM paradigms

Pre/Post-Retrieval Techniques



Pre-retrieval process: to improve the adaptation and effectiveness of the query

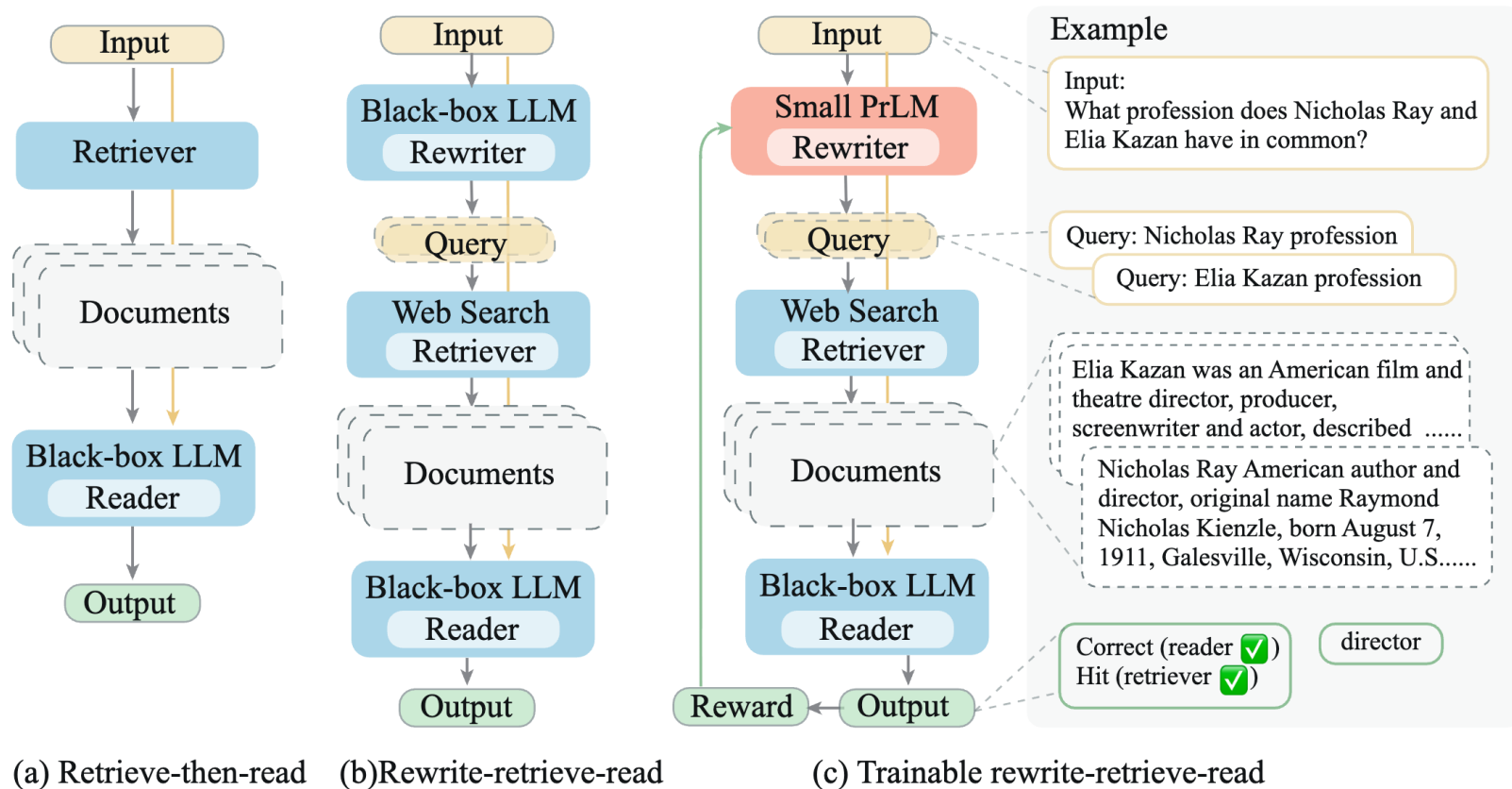
- Query rewriting
- Query decomposition
- Query expansion

Post-retrieval process: to select better results, merge multiple ones, etc

- Reranking
- Pruning
- Verification/Correction
- Dynamic Selection

Pre-Retrieval Techniques

❑ **Query Rewriting:** to improve the adaptation of the query

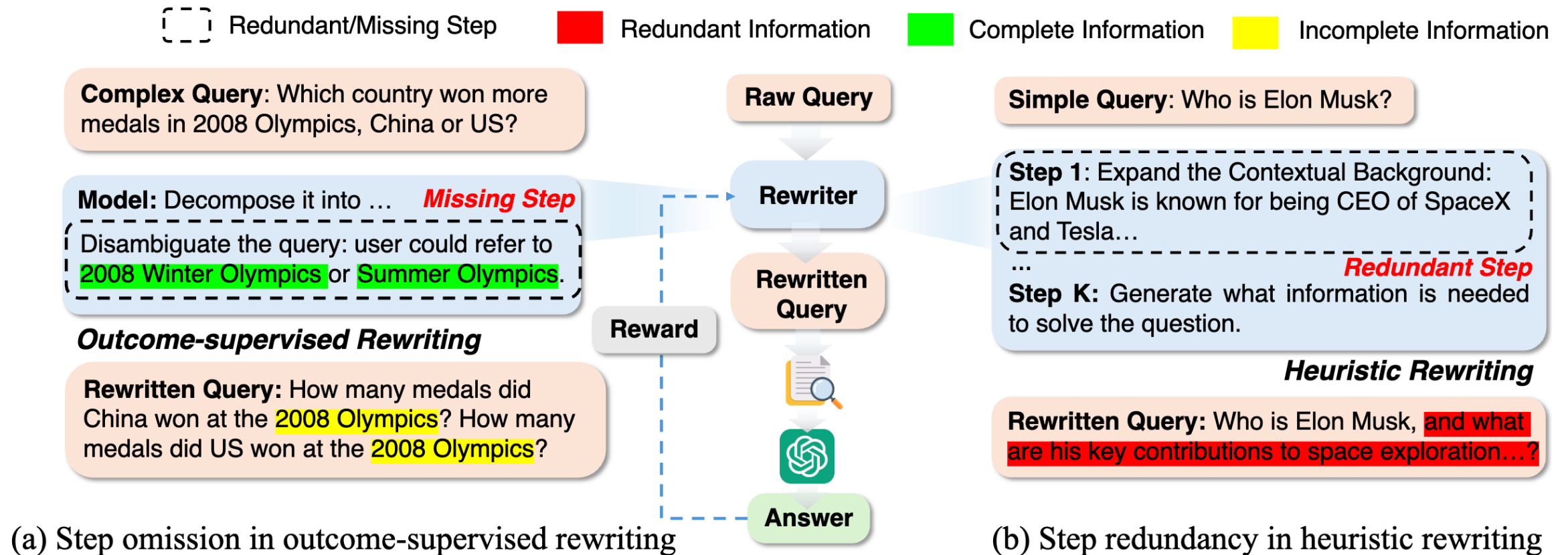


Model	EM	F ₁
<i>HotpotQA</i>		
Direct	32.36	43.05
Retrieve-then-read	30.47	41.34
LLM rewriter	32.80	43.85
Trainable rewriter	34.38	45.97
<i>AmbigNQ</i>		
Direct	42.10	53.05
Retrieve-then-read	45.80	58.50
LLM rewriter	46.40	58.74
Trainable rewriter	47.80	60.71
<i>PopQA</i>		
Direct	41.94	44.61
Retrieve-then-read	43.20	47.53
LLM rewriter	46.00	49.74
Trainable rewriter	45.72	49.51

Works on different QA settings

Pre-Retrieval Techniques

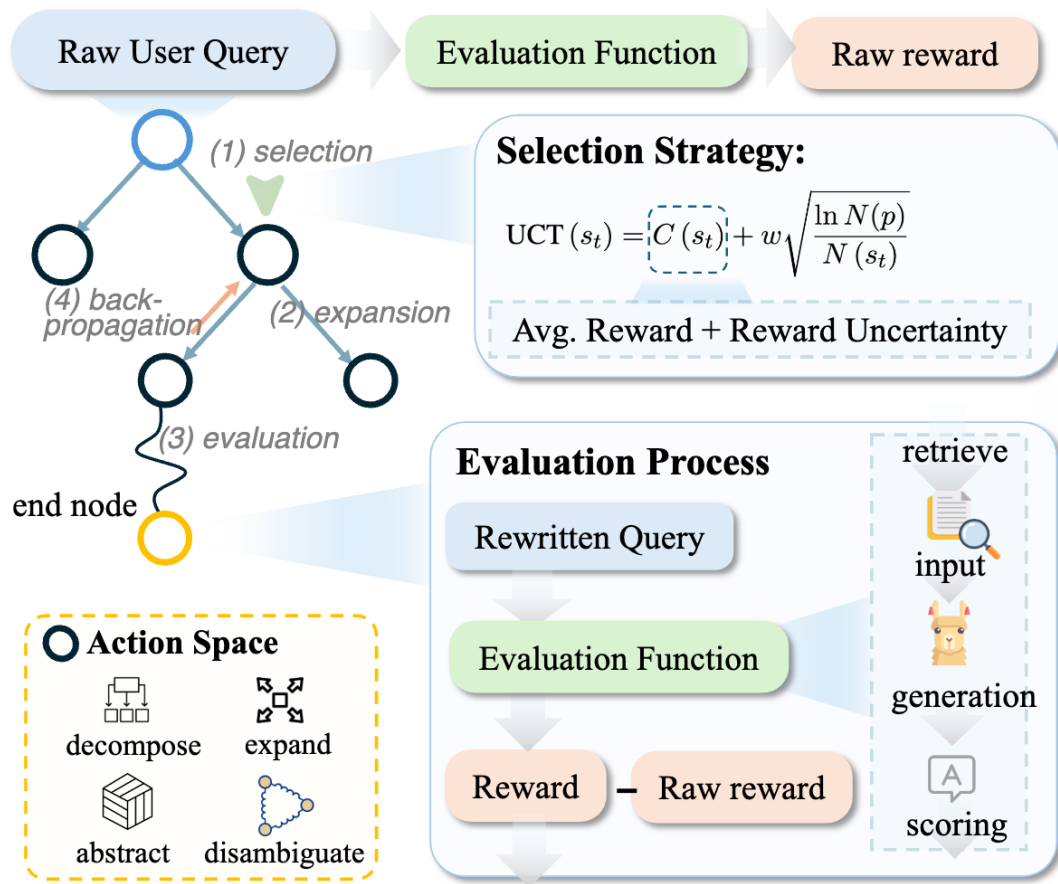
❑ Adaptive Query Rewriting : Step-level Process Supervision



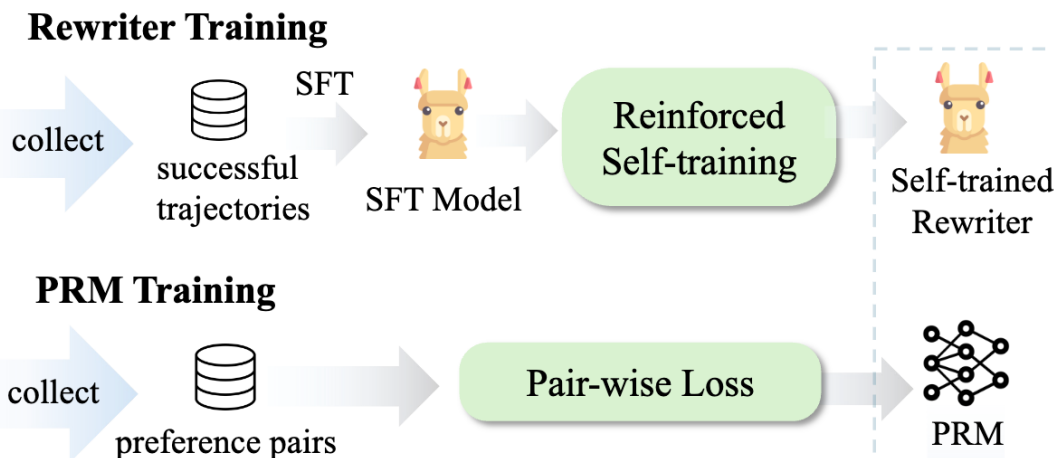
Pre-Retrieval Techniques

Adaptive Query Rewriting : Step-level Process Supervision

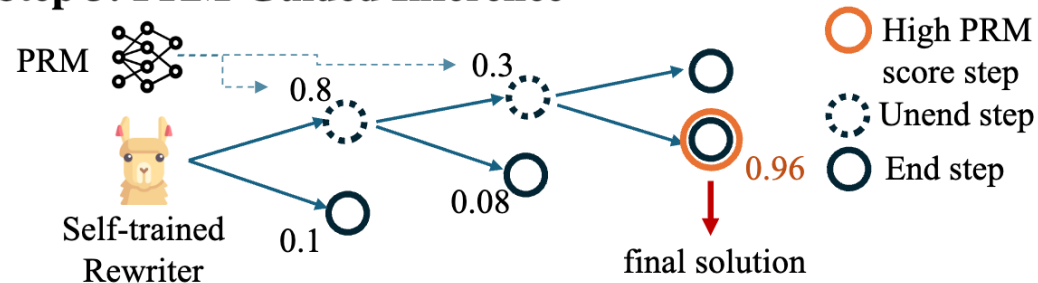
Step 1: Uncertainty-aware Monte Carlo Tree Search



Step 2: Model Training

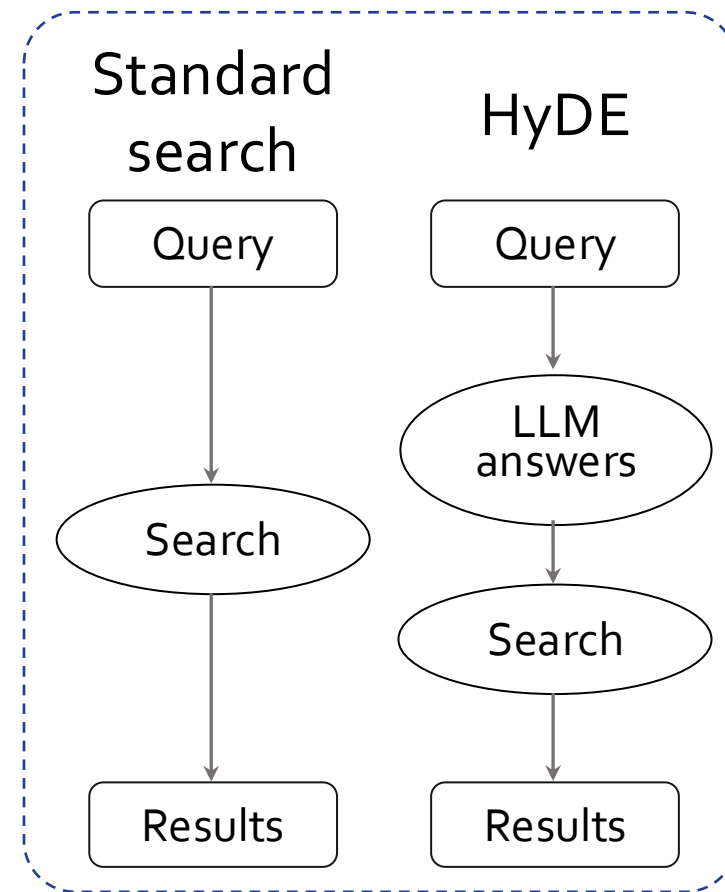
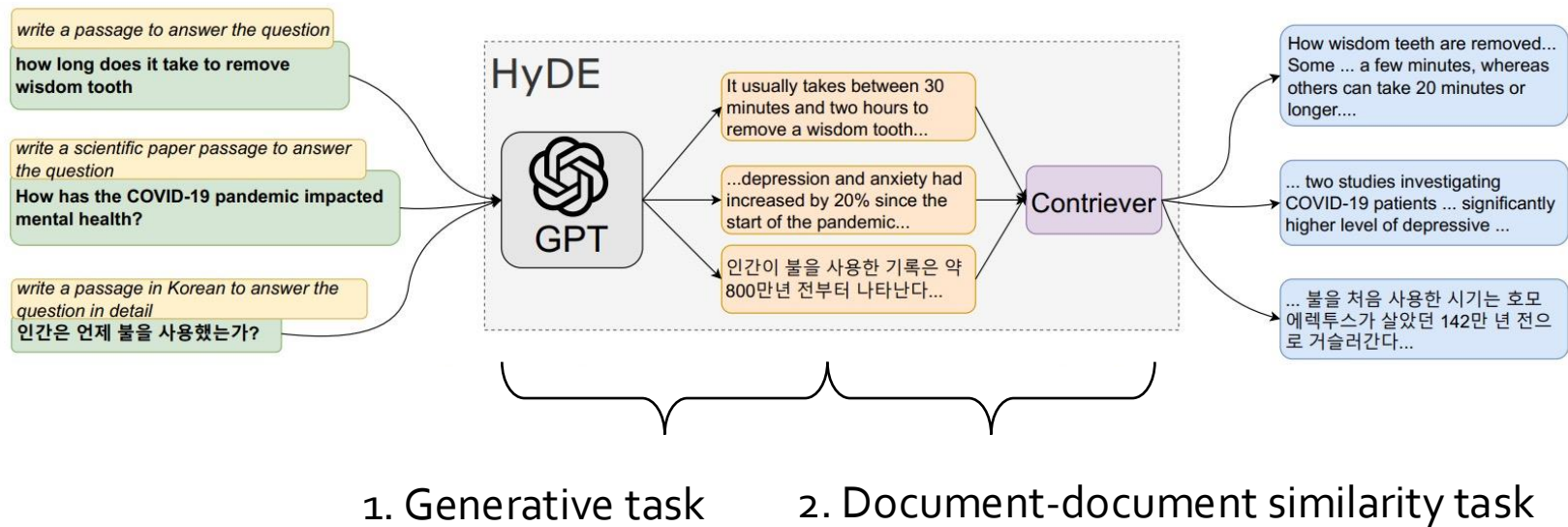


Step 3: PRM-Guided Inference



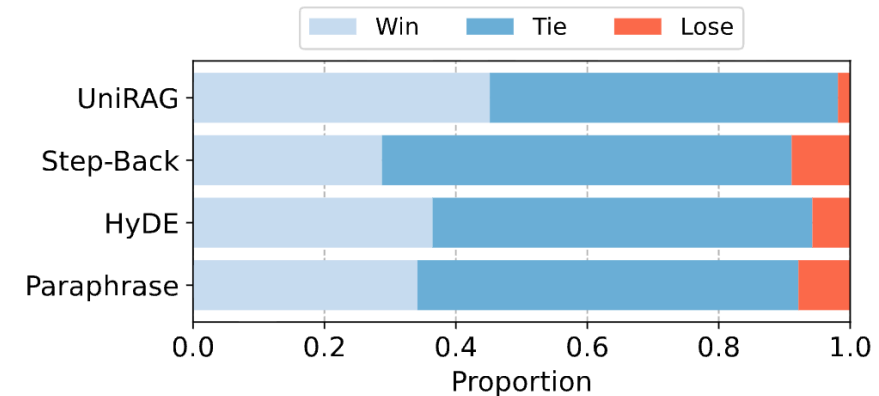
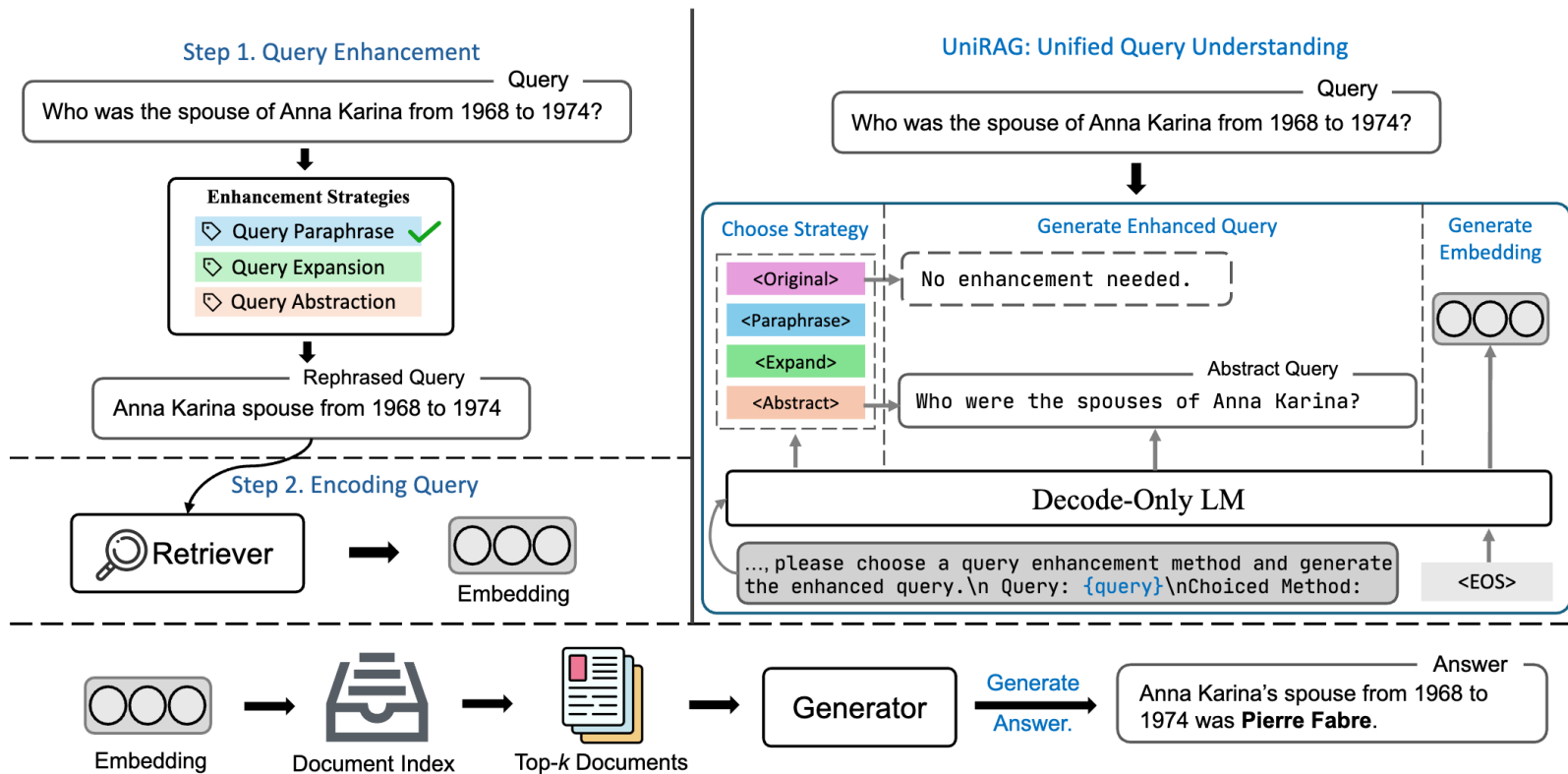
Pre-Retrieval Techniques

❑ Query Expansion: Hypothetical Document Embeddings (HyDE)



Pre-Retrieval Techniques

Query Expansion

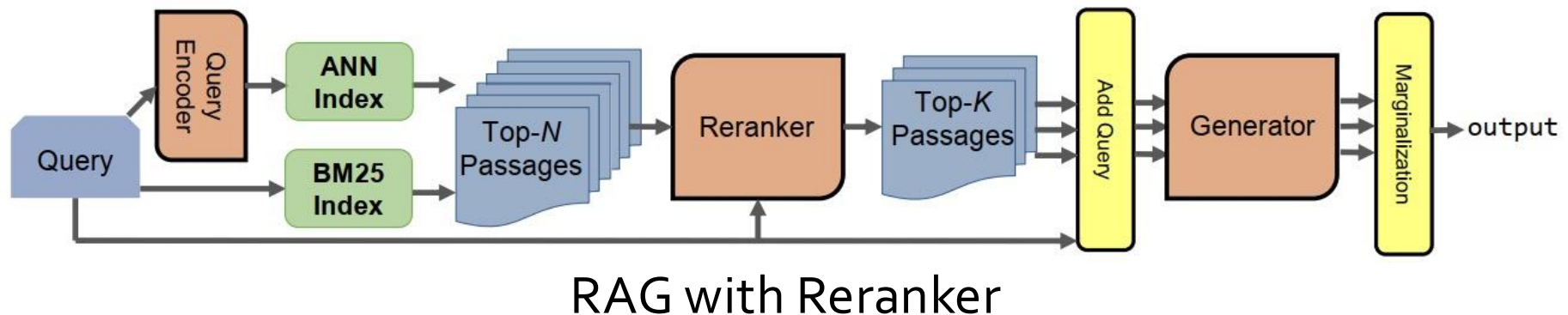
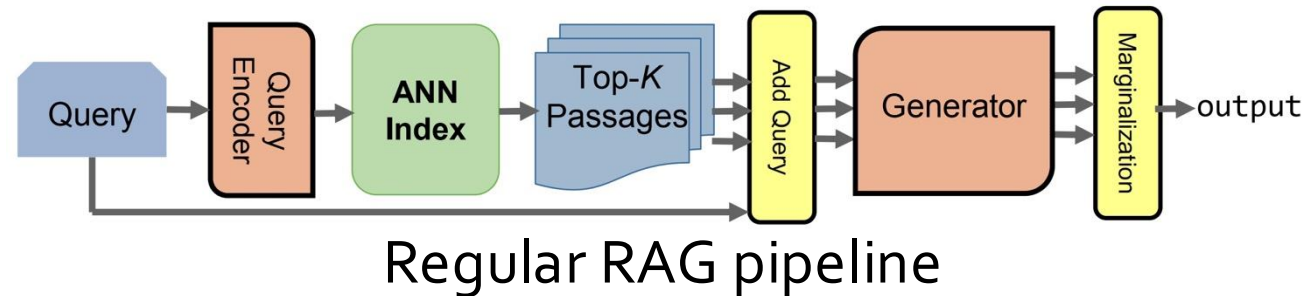


Comparison of win rates of different augmentation methods.

Post-Retrieval Techniques

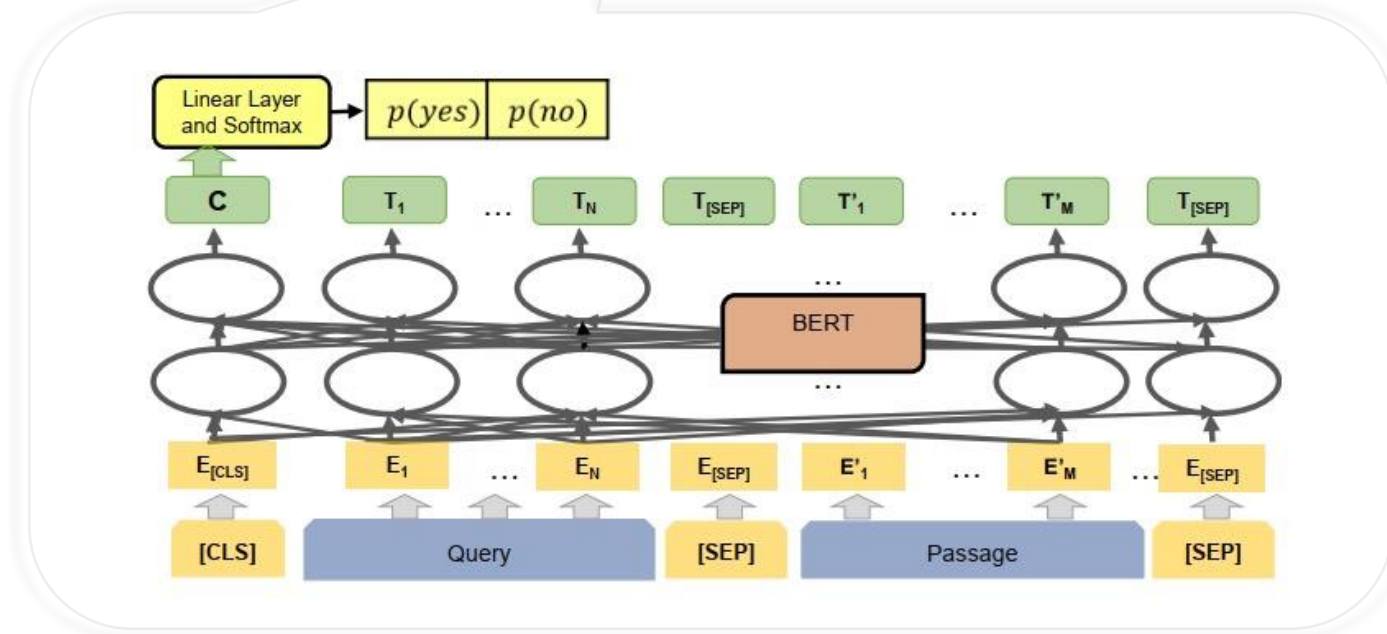
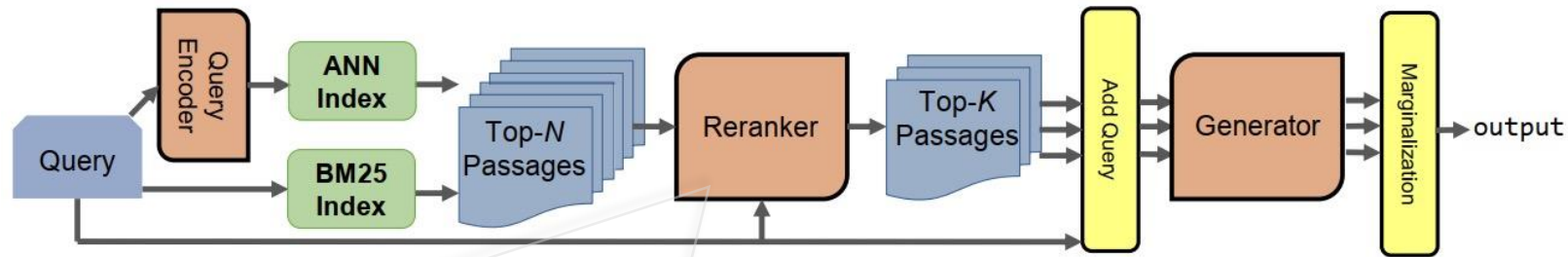
❑ Retrieved Result Rerank (Re2G)

- ❖ Results from initial retrieval can be greatly improved through the use of a reranker
- ❖ Reranker allows merging retrieval results from sources with incomparable scores, e.g., BM25 and neural initial retrieval



Retrieved Result Rerank Model

- ❑ Reranker: interaction model based on the sequence-pair classification



Retrieved Result Rerank (Re²G) Performance

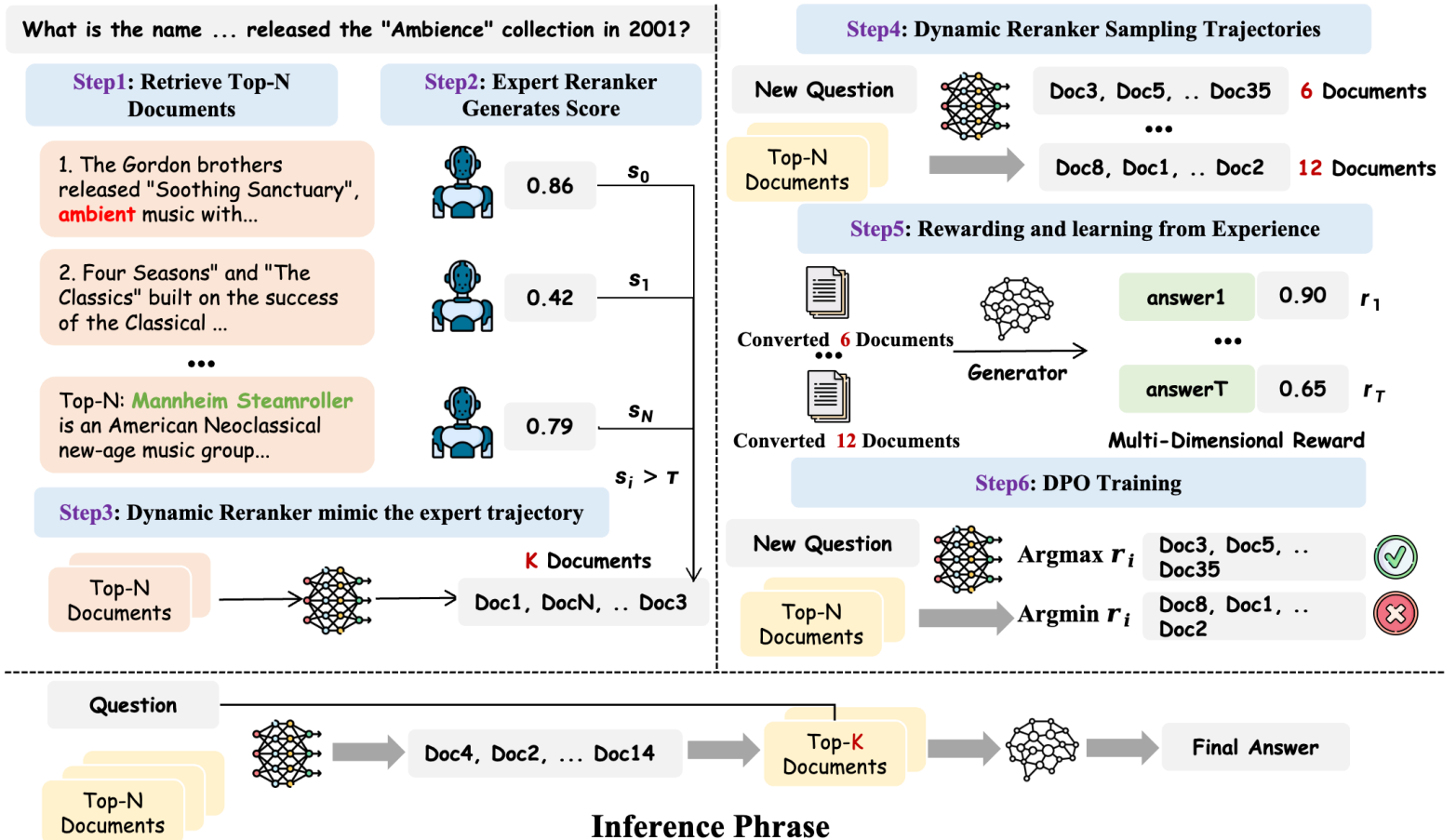
	T-REx		NQ		TriviaQA		FEVER		WoW	
	R-Prec	R@5	R-Prec	R@5	R-Prec	R@5	R-Prec	R@5	R-Prec	R@5
BM25	46.88	69.59	24.99	42.57	26.48	45.57	42.73	70.48	27.44	45.74
DPR Stage 1	49.02	63.34	56.64	64.38	60.12	64.04	75.49	84.66	34.74	60.22
KGI ₀ DPR	65.02	75.52	64.65	69.60	60.55	63.65	80.34	86.53	48.04	71.02
Re ² G DPR	67.16	76.42	65.88	70.90	62.33	65.72	84.13	87.90	47.09	69.88
KGI ₀ DPR+BM25	60.48	80.06	36.91	66.94	40.81	64.79	65.95	90.34	35.63	68.47
Reranker Stage 1	81.22	87.00	70.78	73.05	71.80	71.98	87.71	92.43	55.50	74.98
Re ² G Reranker	81.24	88.58	70.92	74.79	60.37	70.61	90.06	92.91	57.89	74.62

Significantly outperforms pipelines without the *Rerank* stage

Retrieved Result Rerank Model

❑ DynamicRAG: Dynamic Reranking in Retrieval-Augmented Generation

- ❖ **Question:** How many retrieved documents should be passed to the generator?
- ❖ **Problem:** Fixed top-k retrieval is query-agnostic.
- ❖ **Key idea:** Dynamically adjust document order and document number using generation feedback.



Retrieved Result Rerank Model

❑ SetR: From reranking to dynamic evidence selection

- ❖ **Question:** Are the top-ranked passages the best evidence set?
- ❖ **Problem:** Individual relevance does not guarantee collective sufficiency.
- ❖ **Key idea:** Select a complementary set of passages that covers the query's information needs.

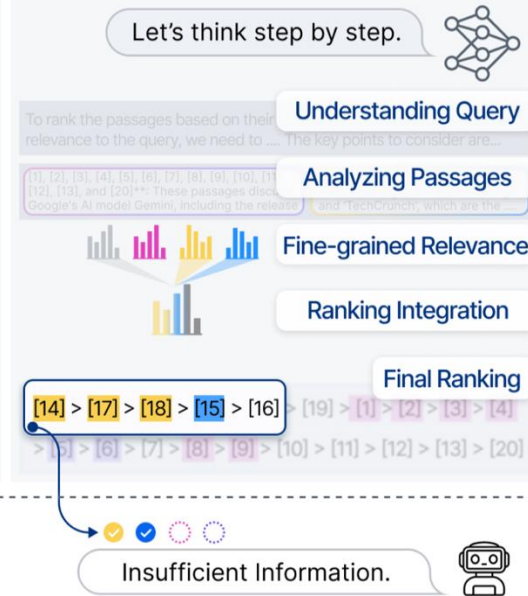
Q. Which company, recently mentioned in articles by 'The Age' and 'TechCrunch', is accused of both **manipulating search results to maximize ad revenue** and **siphoning off news publishers' content and ad revenue**, while **also claiming superior performance for its AI model Gemini compared to competitors**, despite only **releasing a 'lite' version known as Gemini Pro**?

First-Stage Retrieval

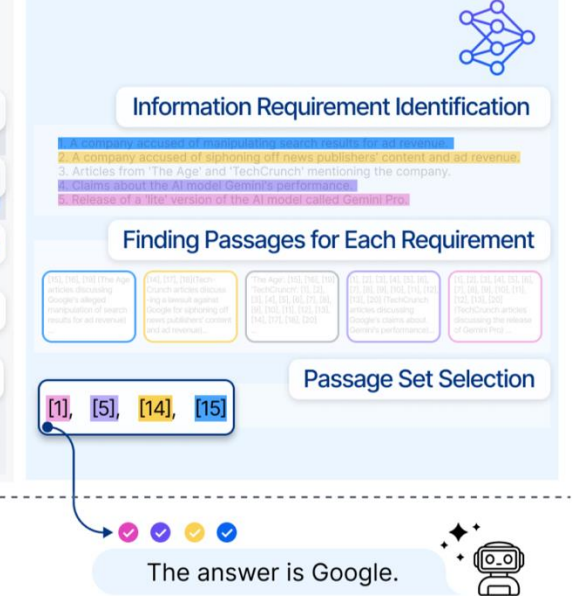
(a) Direct Reranking



(b) Reranking with CoT



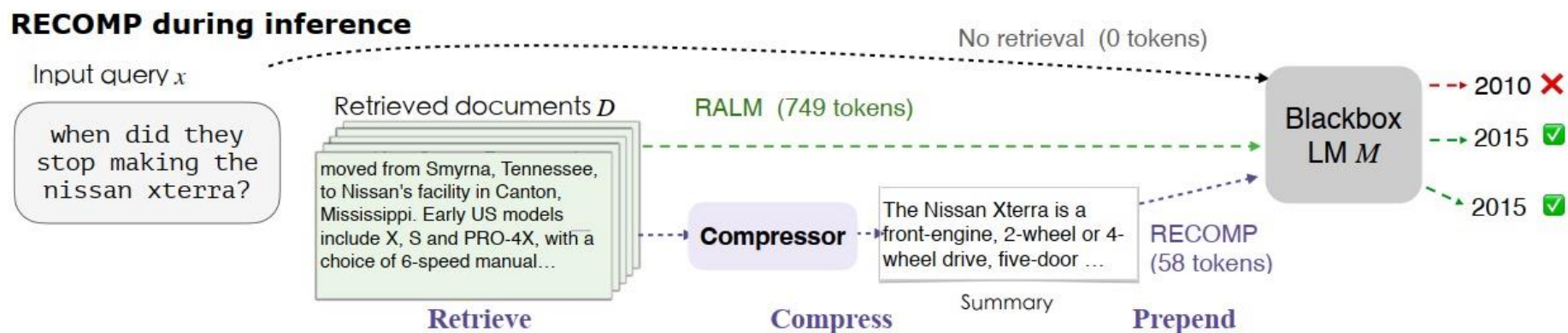
(c) Passage Set Selection (Ours)



Post-Retrieval Techniques

Retrieved Result Compression

- ❖ To reduce the computational costs and also relieve the burden of LMs to identify relevant information in long retrieved documents.



Compressor Learning Objectives

- ❖ Concise
- ❖ Effective
- ❖ Faithful

Retrieved Result Compression Performance

□ QA tasks

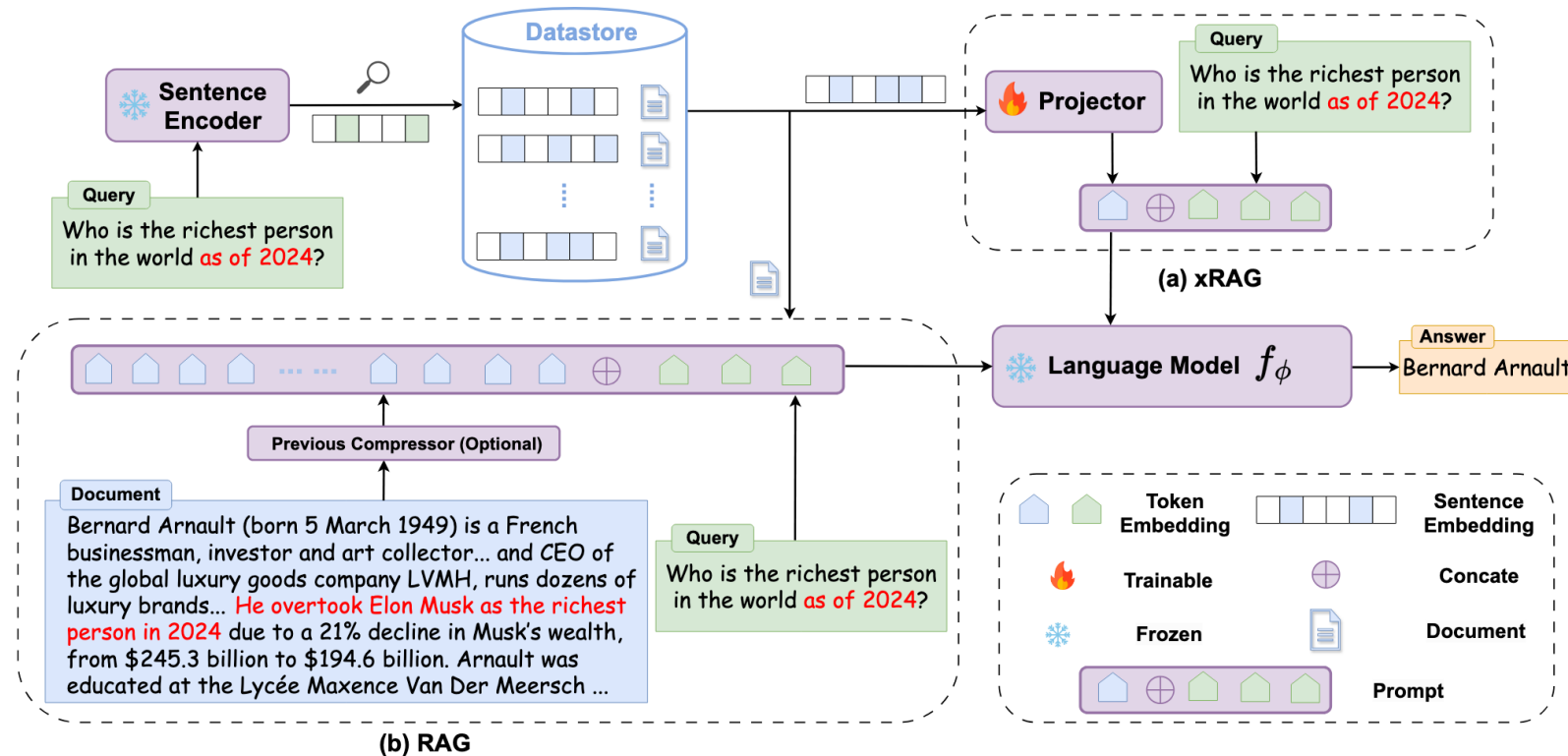
In-Context evidence	# tok	NQ		# tok	TQA		HotpotQA		
		EM	F1		EM	F1	EM	F1	
-	0	21.99	29.38	0	49.33	54.85	0	17.80	26.10
<i>RALM without compression</i>									
Top 1 documents	132	33.07	41.45	136	57.84	64.94	138	28.80	40.58
Top 5 documents	660	39.39	48.28	677	62.37	70.09	684	32.80	43.90
<i>Phrase/token level compression</i>									
Top 5 documents (NE)	338	23.60	31.02	128	54.96	61.19	157	22.20	31.89
Top 5 documents (BoW)	450	28.48	36.84	259	58.16	65.15	255	25.60	36.00
<i>Extractive compression of top 5 documents</i>									
Oracle	34	60.22	64.25	32	79.29	82.06	70	41.80	51.07
Random	32	23.27	31.09	31	50.18	56.24	61	21.00	29.86
BM25	36	25.82	33.63	37	54.67	61.19	74	26.80	38.02
DPR	39	34.32	43.38	41	56.58	62.96	78	27.40	38.15
Contriever	36	30.06	31.92	40	53.67	60.01	78	28.60	39.48
Ours	37	36.57	44.22	38	58.99	65.26	75	30.40	40.14

Outperforms representative sparse and dense retrievers

Post-Retrieval Techniques

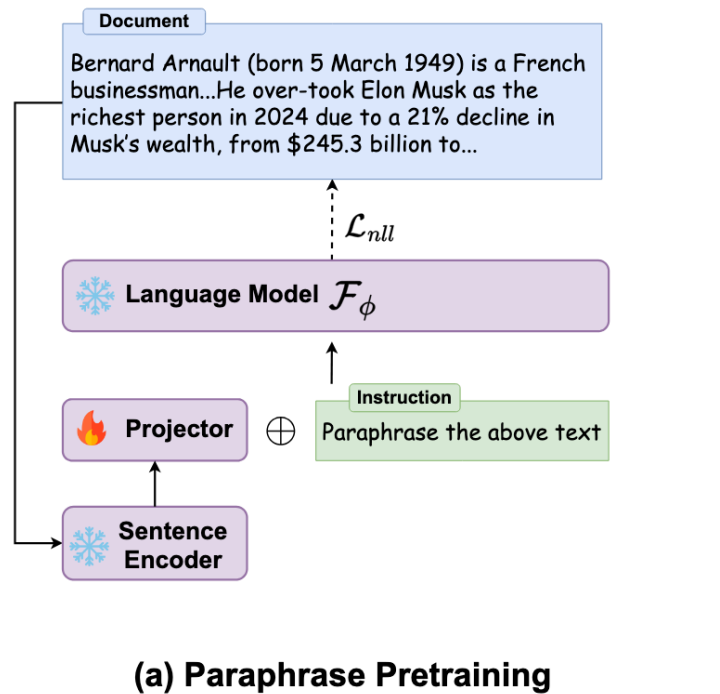
❑ xRAG: Extreme Context Compression

- ❖ xRAG employs a modality fusion methodology to seamlessly integrate document embeddings into the language model's representation space.

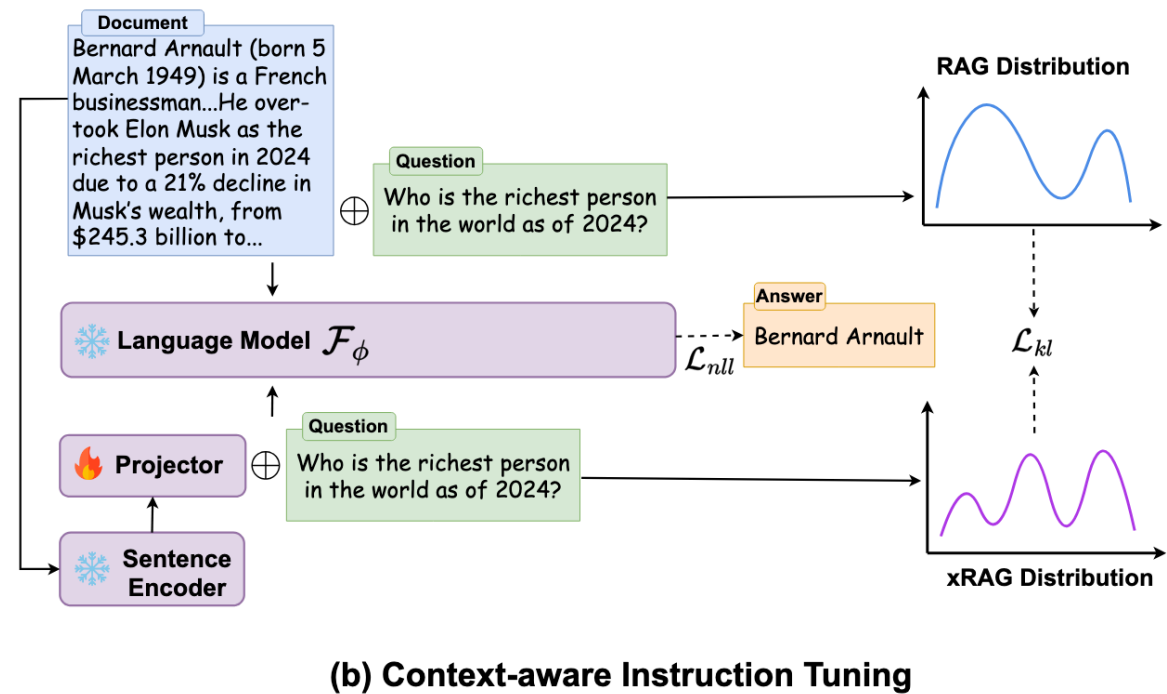


Post-Retrieval Techniques

❑ xRAG: Extreme Context Compression



$$\mathcal{L}_{nll} = - \sum_{i=1} \log p_{\phi}(d_i | \mathbf{W}(E), \mathbf{X}_{instruct}, d_{<i})$$



$$\mathcal{L}_{nll} = - \sum_{i=1} \log p_{\phi}(\mathbf{X}_{answer,i} | \mathbf{W}(E_{context}), \mathbf{X}_{question}, \mathbf{X}_{answer,<i})$$

$$\mathcal{L}_{kl} = D_{KL}(p_{\phi}(\mathbf{X}_{answer} | \mathbf{X}_{context}, \cdot) || p_{\phi}(\mathbf{X}_{answer} | \mathbf{W}(E_{context}), \cdot))$$

Post-Retrieval Techniques

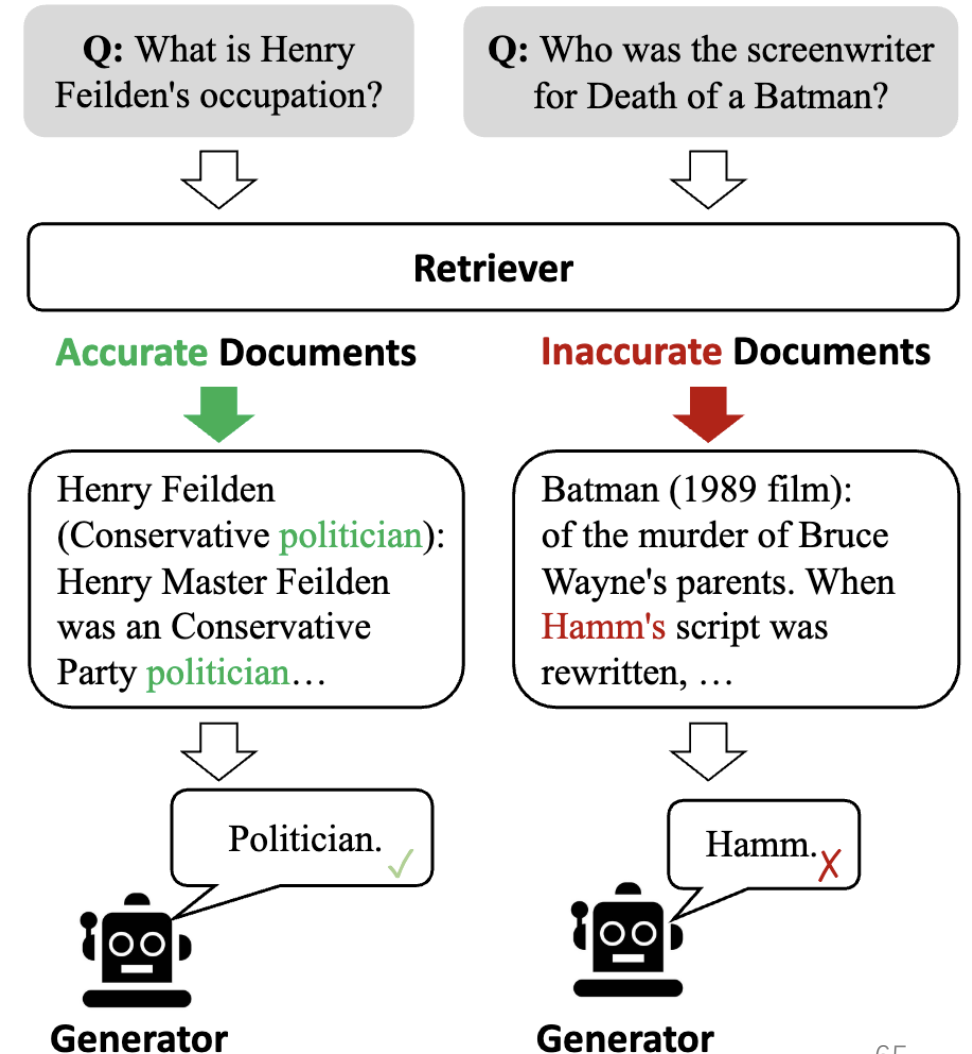
❑ xRAG: Extreme Context Compression

Task Type	NQ	TriviaQA Open-Domain QA (EM)	WebQA	HotpotQA Multihop QA (EM)	TrutefulQA Long-form QA (F1/R-L)	FactKG Fact Checking (Acc)	Average	# Doc Length	
Mistral-7b									
w/o retrieval	30.25	57.08	34.89	27.02	26.23	25.51	54.78	36.54 (0.0%)	0
w retrieval	42.71	65.88	<u>37.84</u>	38.79	<u>26.50</u>	<u>25.92</u>	67.76	43.63 (19.4%)	175.1
*with Compression									
LLMLingua [†]	30.64	57.94	32.63	29.91	25.70	25.10	64.17	38.01 (4.0%)	98.6
LLMLingua [‡]	28.81	57.09	32.33	29.13	26.10	25.39	63.57	37.48 (2.5%)	61.1
TF-IDF	30.25	58.49	35.43	26.62	26.33	25.83	59.56	37.49 (2.6%)	1
xRAG	39.10	<u>65.77</u>	39.40	<u>34.05</u>	28.10	27.71	<u>63.08</u>	<u>42.46 (16.2%)</u>	1
Mixtral-8x7b									
w/o retrieval	41.99	<u>71.10</u>	40.31	32.87	25.60	24.90	62.64	42.76 (0.0%)	0
w retrieval	<u>45.15</u>	<u>70.34</u>	<u>41.26</u>	43.46	<u>27.10</u>	<u>25.80</u>	70.42	<u>46.22 (8.0%)</u>	175.1
*with Compression									
LLMLingua [†]	37.65	67.70	36.02	35.66	25.99	25.39	67.98	42.32 (-1.0%)	96.6
LLMLingua [‡]	37.81	67.81	35.78	35.27	25.68	25.00	68.03	44.17 (-1.3%)	61.1
TF-IDF	41.19	69.94	41.63	32.05	26.80	26.00	66.17	43.41 (1.4%)	1
xRAG	47.28	74.14	44.50	<u>39.66</u>	27.80	26.64	<u>68.20</u>	46.91 (9.7%)	1

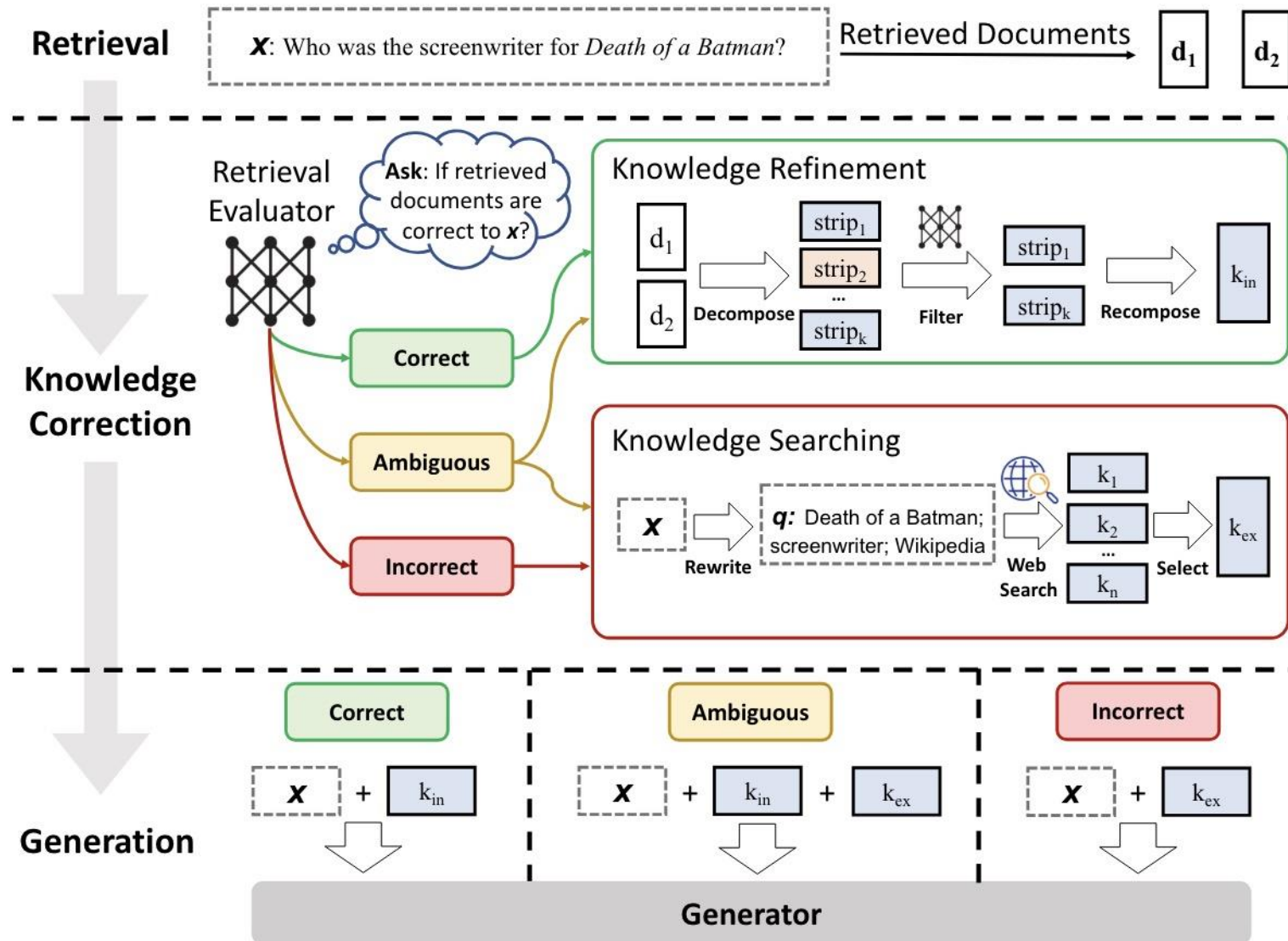
Post-Retrieval Techniques: Corrective RAG

❑ Corrective Retrieval Augmented Generation (CRAG)

- ❖ Although retrieval-augmented generation (RAG) is a practicable complement to LLMs, it relies heavily on the **relevance of retrieved documents**
- ❖ A lightweight **retrieval evaluator** is designed to assess the overall quality of retrieved documents for a query, returning a **confidence degree** based on which different knowledge retrieval actions can be triggered

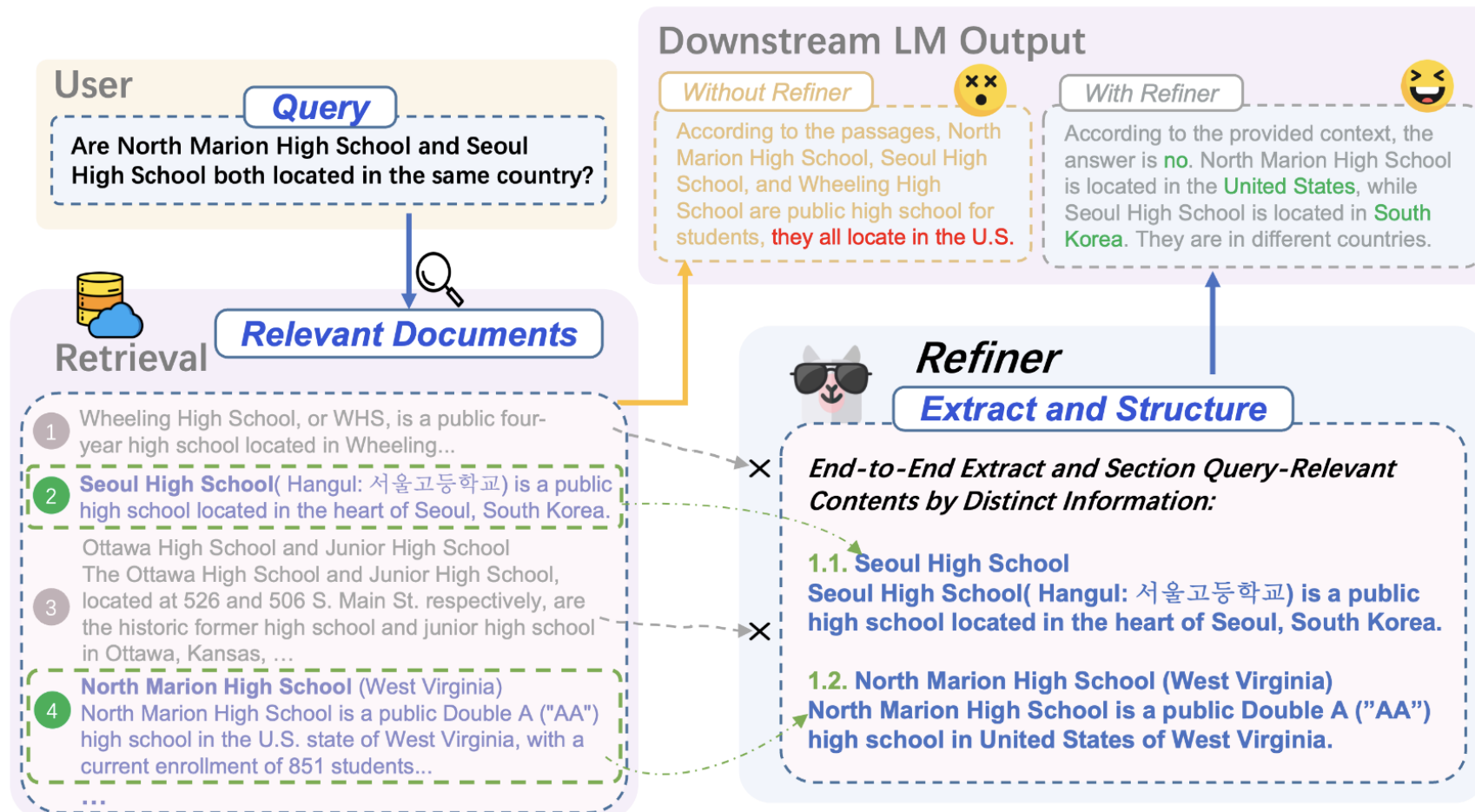


Post-Retrieval Techniques: Corrective RAG



Post-Retrieval Techniques: Refiner

- ❑ **Refiner**: leveraging a single decoder-only LLM to adaptively extract query relevant contents verbatim along with the necessary context



PART 2: Architecture of RA-LFMs and Main Modules



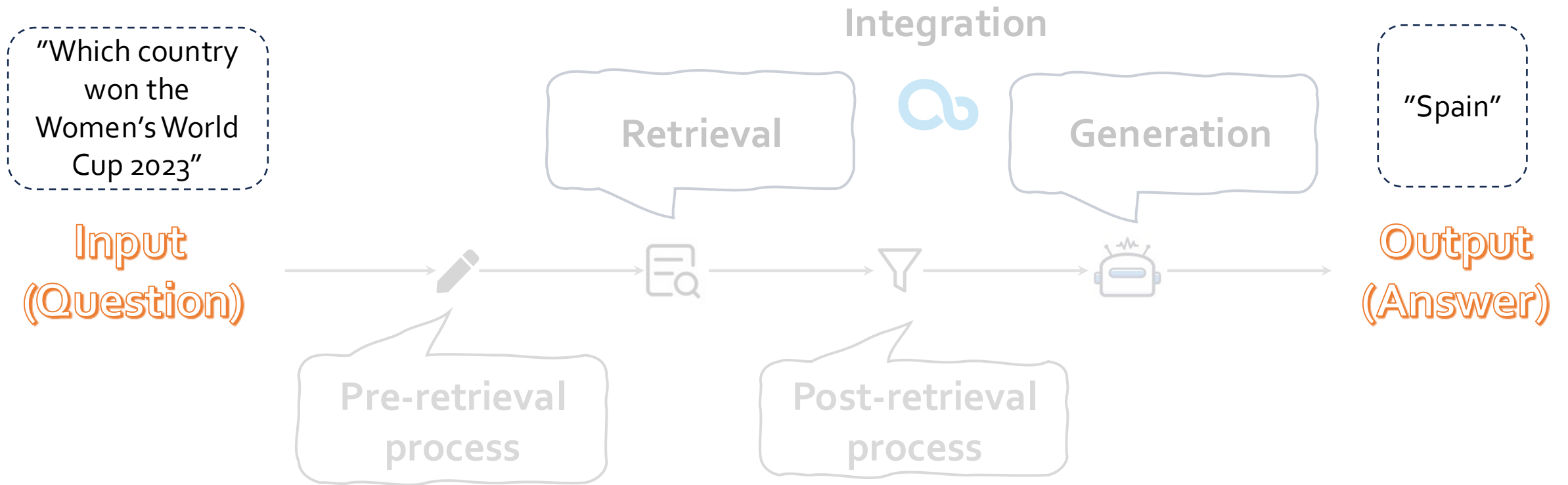
Slides



Website of this tutorial

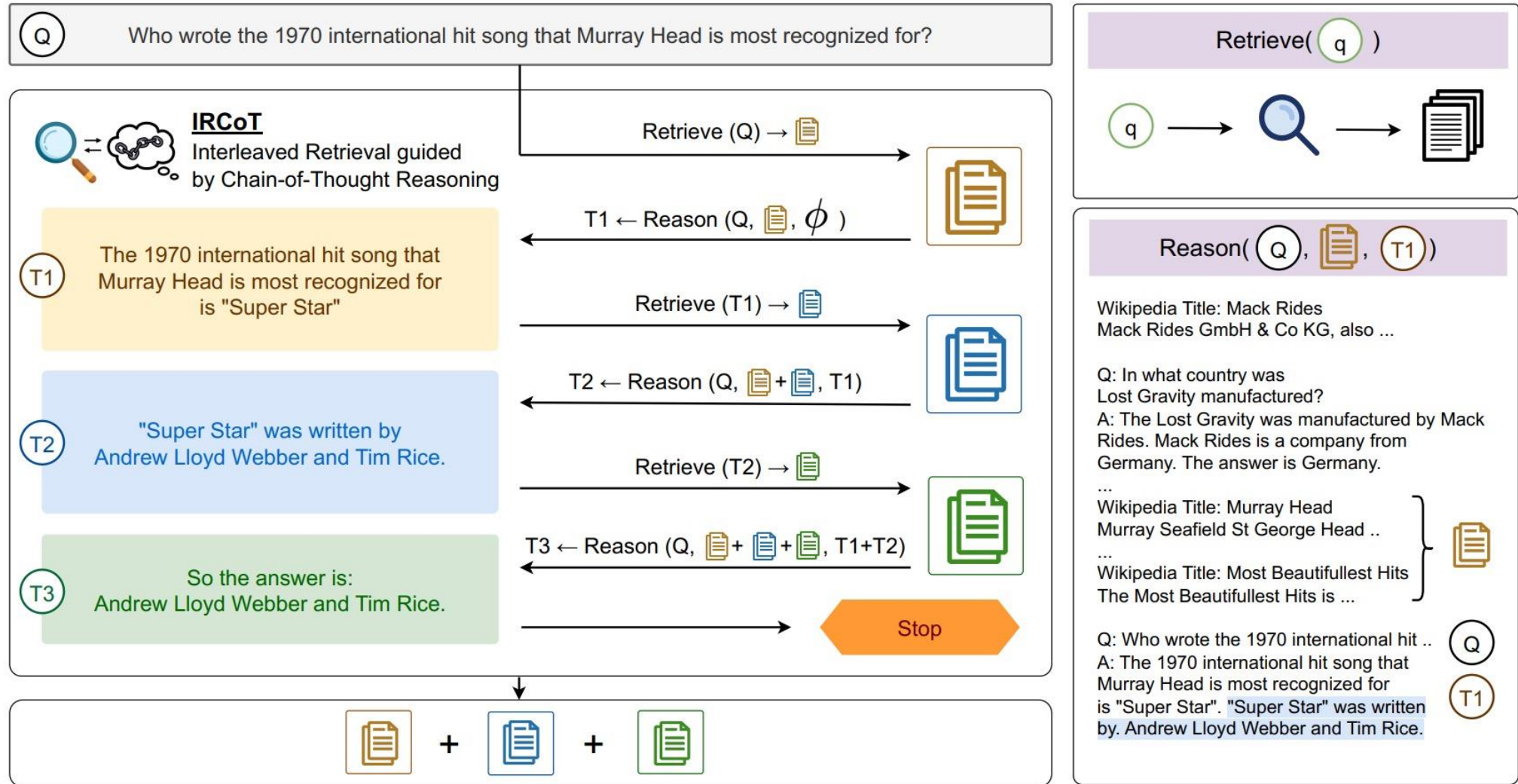
- **RA-LFM architecture overview**
- **Retriever in RA-LFMs**
- **Retrieval results integration**
- **Pre/Post-retrieval techniques**
- **Special RA-LFM paradigms**

Beyond Standard Pipelines and Components?



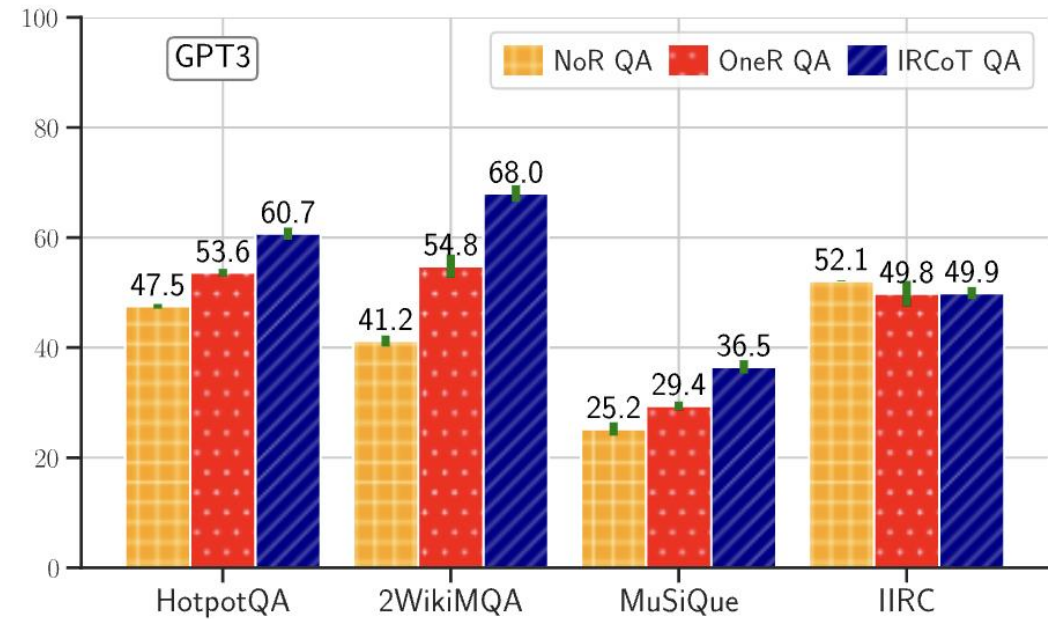
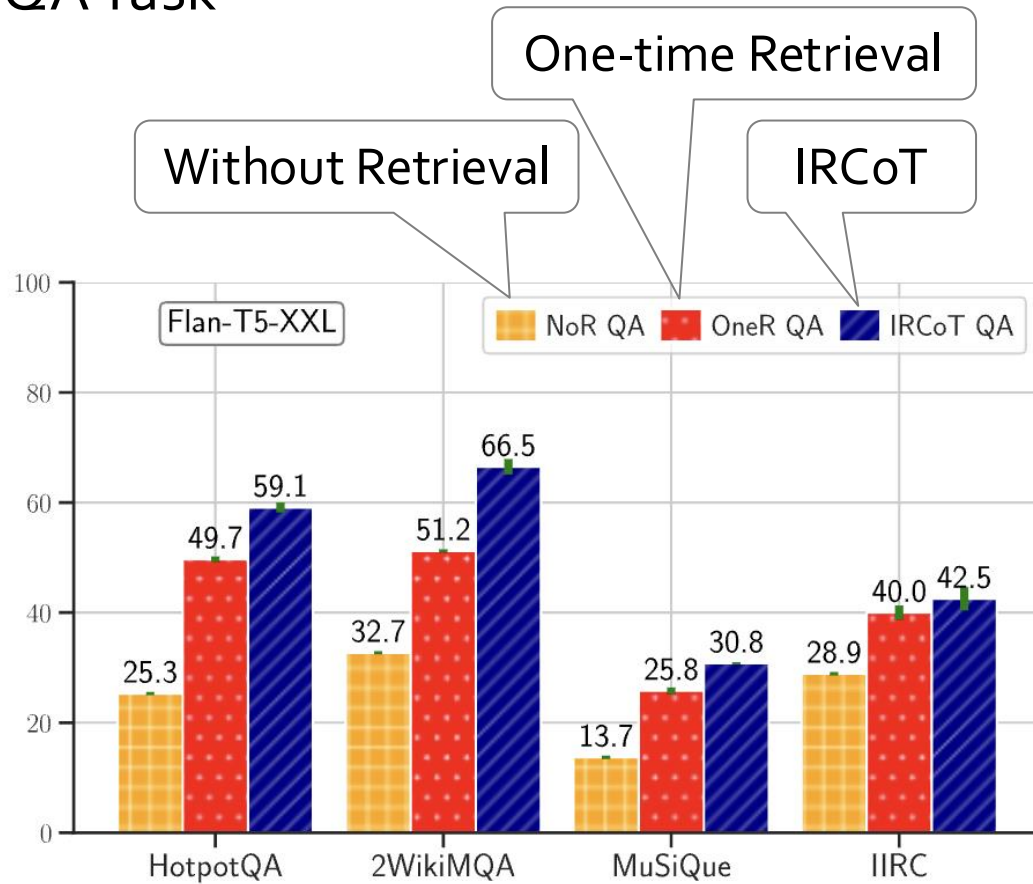
More?

Interleaved Retrieval guided by Chain-of-Thought (IRCoT)



IRCoT Performance

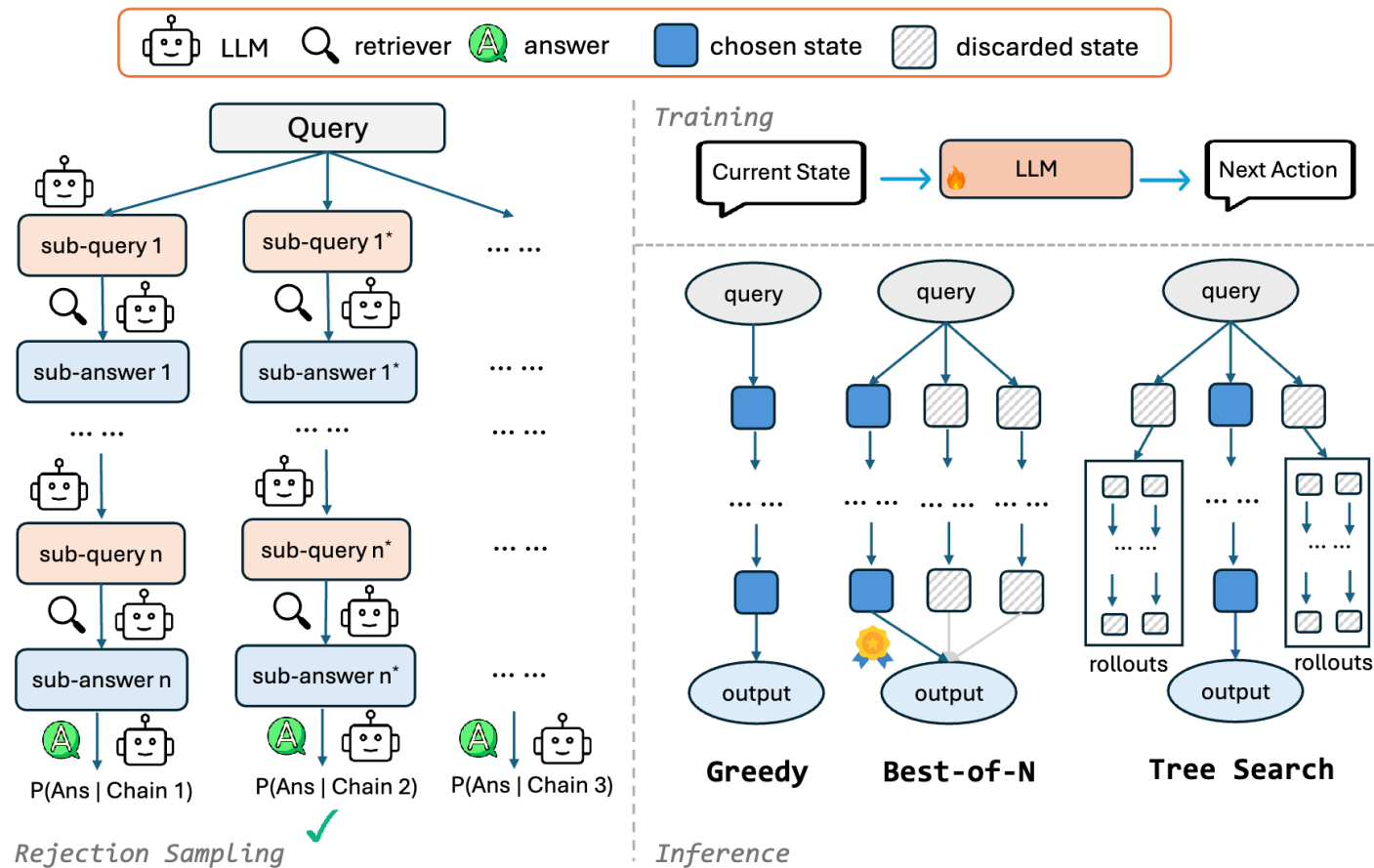
QA Task



Special RAG Pipeline: Chain-of-RAG

Chain-of-Retrieval Augmented Generation

- ❖ **CoRAG** retrieves and reasons step by step, where each retrieval is conditioned on the evolving reasoning state.
- ❖ **Training:** Rejection sampling builds intermediate retrieval chains.
- ❖ **Inference:** Control chain length and chain number to scale test-time compute.



Special RAG Pipeline: Chain-of-RAG

❑ HotpotQA

Query: What wrestling team is Mathew Thomas Rehwoldt a part of?

RAG without Chain-of-Retrieval

Final Answer: **WWE ✗**

CoRAG (Ours)

Sub-query 1: What is Mathew Thomas Rehwoldt's profession?

Sub-answer 1: No relevant information found.

Sub-query 2: What is Mathew Thomas Rehwoldt's name in the wrestling industry?

Sub-answer 2: Aiden English

Sub-query 3: What wrestling team is Aiden English a part of?

Sub-answer 3: The Vaudevillains

Final Answer: **The Vaudevillains ✓**

Query: How many months apart are Johan Mjällby and Neil Lennon in age?

RAG without Chain-of-Retrieval

Final Answer: **two months ✗**

CoRAG (Ours)

Sub-query 1: What is Johan Mjällby's birthdate?

Sub-answer 1: 9 February 1971

Sub-query 2: What is Neil Lennon's birthdate?

Sub-answer 2: 25 June 1971

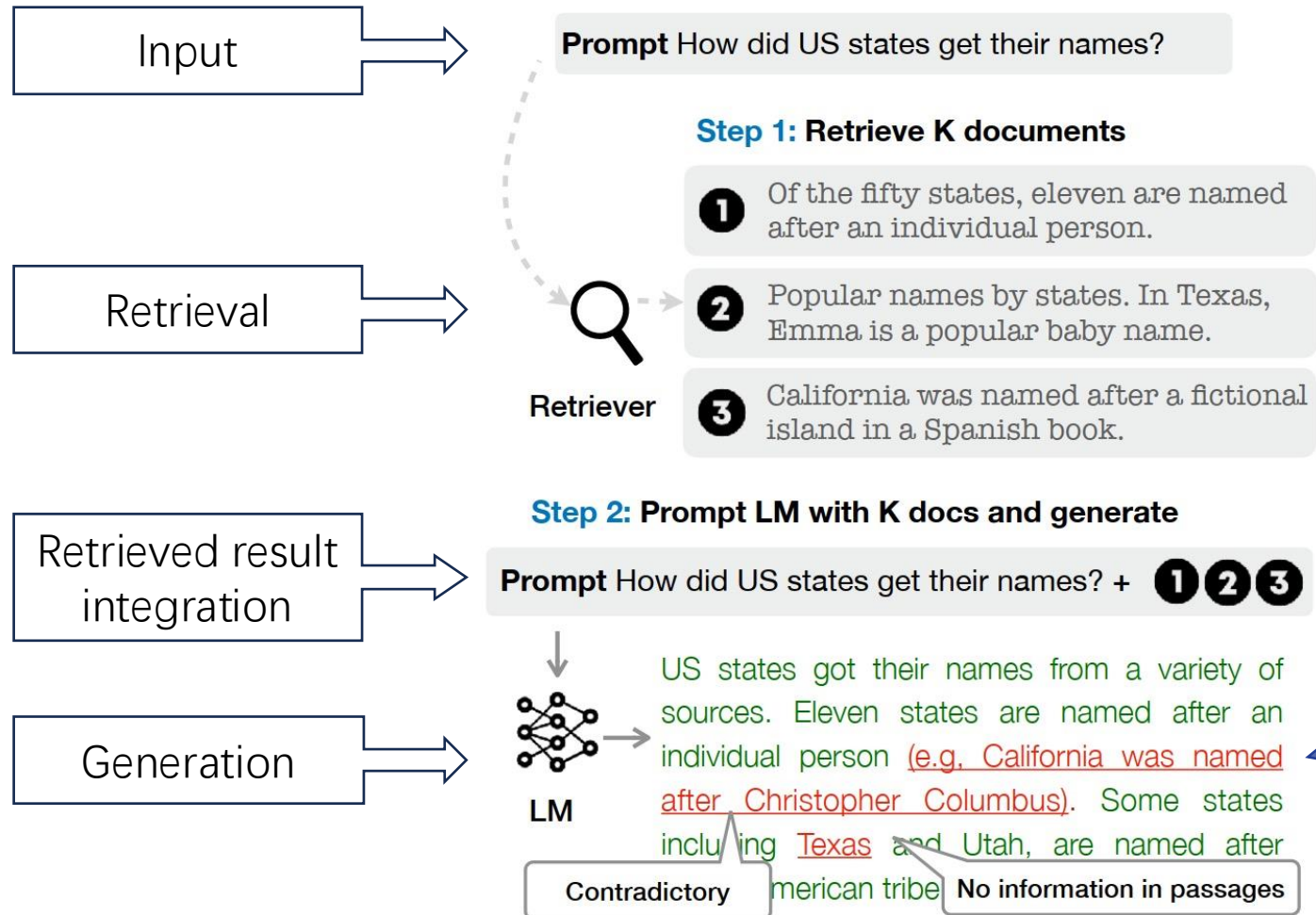
Sub-query 3: What is the difference in months between 9 February 1971 and 25 June 1971?

Sub-answer 3: 4 months

Final Answer: **4 ✓**

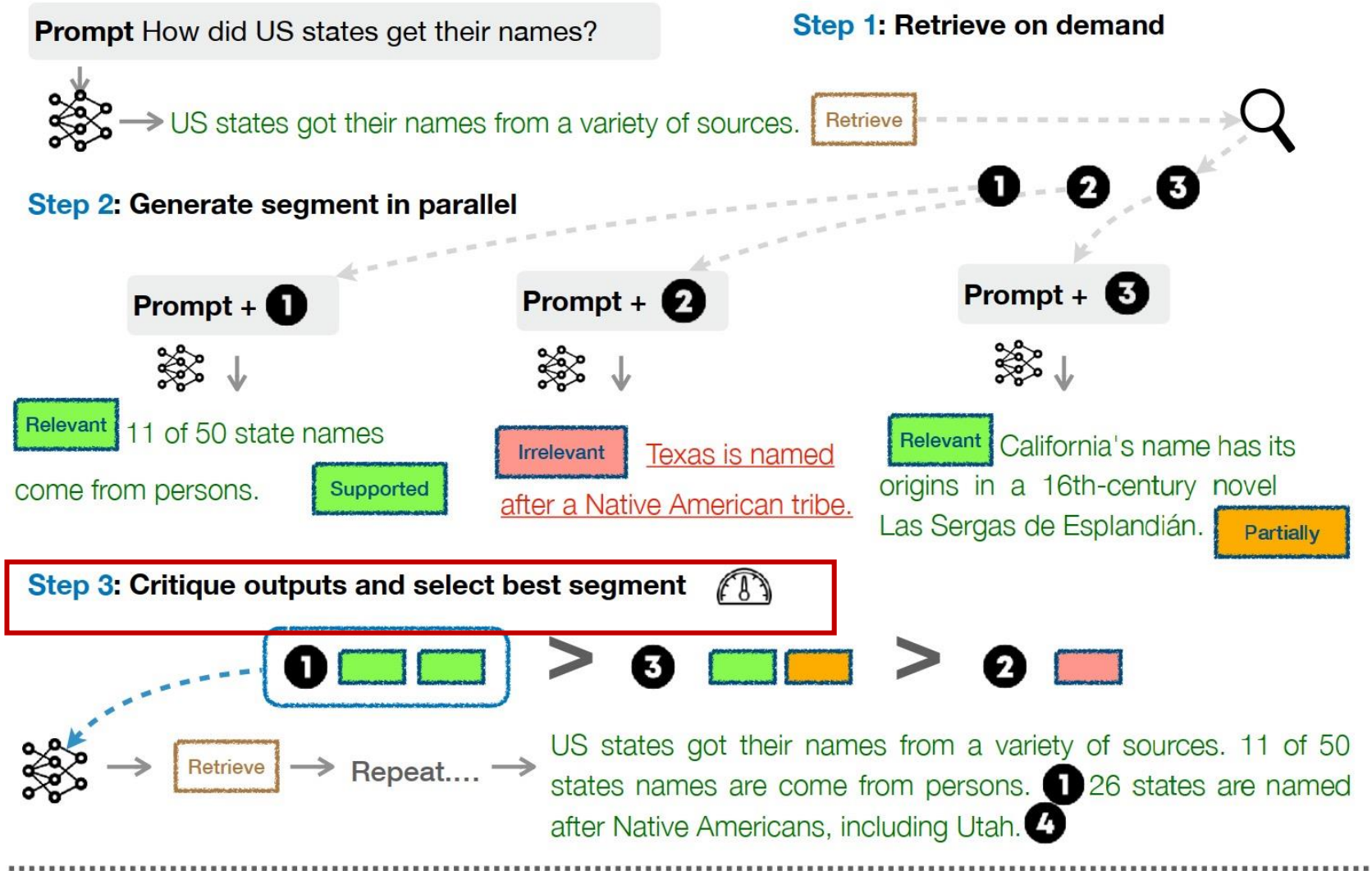
Special RAG Pipeline: Self-Reflective RAG (SELF-RAG)

❑ General Retrieval-Augmented Generation (RAG)



Retrieval results may not be necessary or helpful!

SELF-RAG Overview



Key Technical Design in SELF-RAG

Four types of reflection tokens used in SELF-RAG

Type	Input	Output	Definitions
Retrieve	$x / x, y$	{yes, no, continue}	Decides when to retrieve with \mathcal{R}
ISREL	x, d	{ relevant , irrelevant}	d provides useful information to solve x .
ISSUP	x, d, y	{ fully supported , partially supported, no support}	All of the verification-worthy statement in y is supported by d .
ISUSE	x, y	{ 5 , 4, 3, 2, 1}	y is a useful response to x .

Critic Model Training

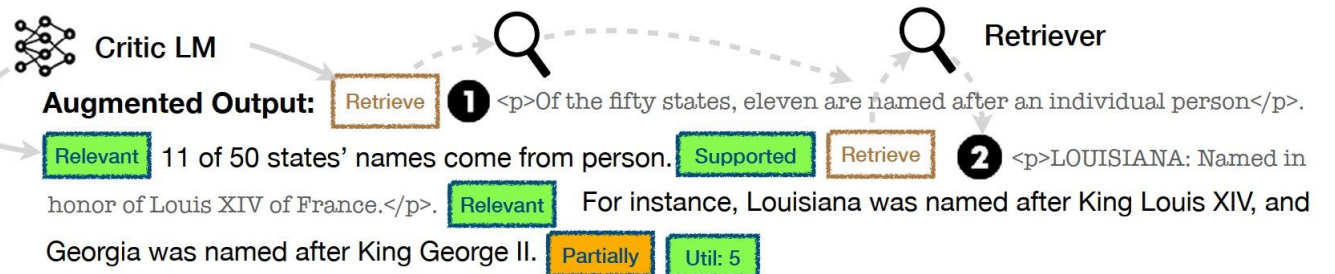
Input: Write an essay of your best summer vacation

Output: My best summer vacation was a magical escape to the coastal town of Santorini. The azure waters, charming white-washed building are unforgettable.

Augmented Output: No Retrieval My best summer vacation was a magical escape to the coastal town of Santorini. No Retrieval The azure waters, charming white-washed building are unforgettable experience. Util: 5

Input: How did US states get their names?

Output: 1 of 50 states names come from persons. For instance, Louisiana was named in honor of King Louis XIV of France and Georgia was named after King George II.



SELF-RAG Algorithm

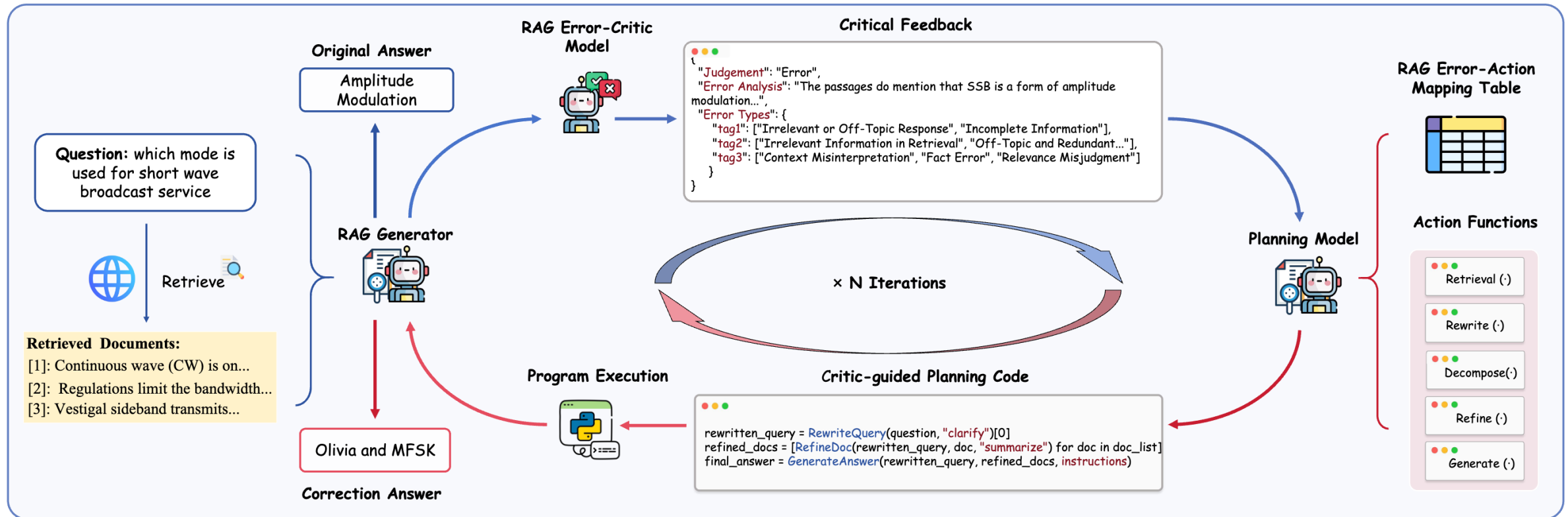
Algorithm 1 SELF-RAG Inference

Require: Generator LM \mathcal{M} , Retriever \mathcal{R} , Large-scale passage collections $\{d_1, \dots, d_N\}$

- 1: **Input:** input prompt x and preceding generation $y_{<t}$, **Output:** next output segment y_t
 - 2: \mathcal{M} predicts **Retrieve** given $(x, y_{<t})$
 - 3: **if** **Retrieve** == Yes **then**
 - 4: Retrieve relevant text passages \mathbf{D} using \mathcal{R} given (x, y_{t-1}) ▷ Retrieve
 - 5: \mathcal{M} predicts **ISREL** given x, d and y_t given $x, d, y_{<t}$ for each $d \in \mathbf{D}$ ▷ Generate
 - 6: \mathcal{M} predicts **ISSUP** and **ISUSE** given x, y_t, d for each $d \in \mathbf{D}$ ▷ Critique
 - 7: Rank y_t based on **ISREL**, **ISSUP**, **ISUSE** ▷ Detailed in Section 3.3
 - 8: **else if** **Retrieve** == No **then**
 - 9: \mathcal{M}_{gen} predicts y_t given x ▷ Generate
 - 10: \mathcal{M}_{gen} predicts **ISUSE** given x, y_t ▷ Critique
-

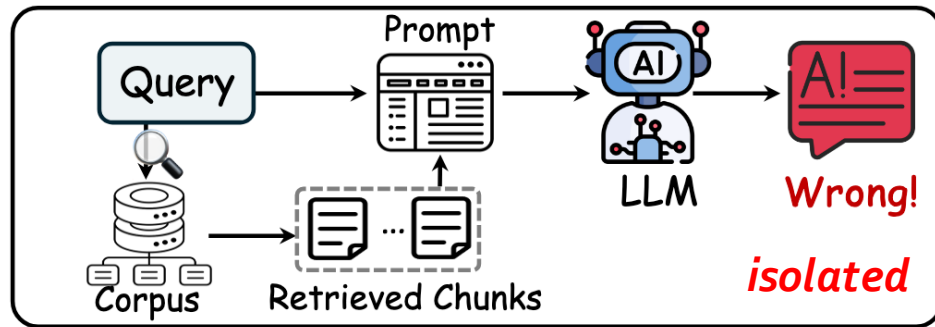
Modern extension: Critic-guided RAG

- ❑ **RAG-Critic**: a critic-guided agentic RAG framework, enabling the planning agent to customize executor-based solution programs based on the error-critic model's feedback, facilitating an automated error-driven self-correction process.

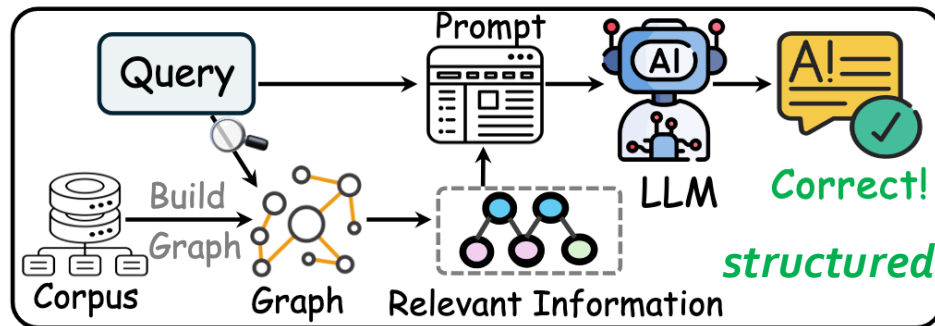


Graph Retrieval-Augmented Generation

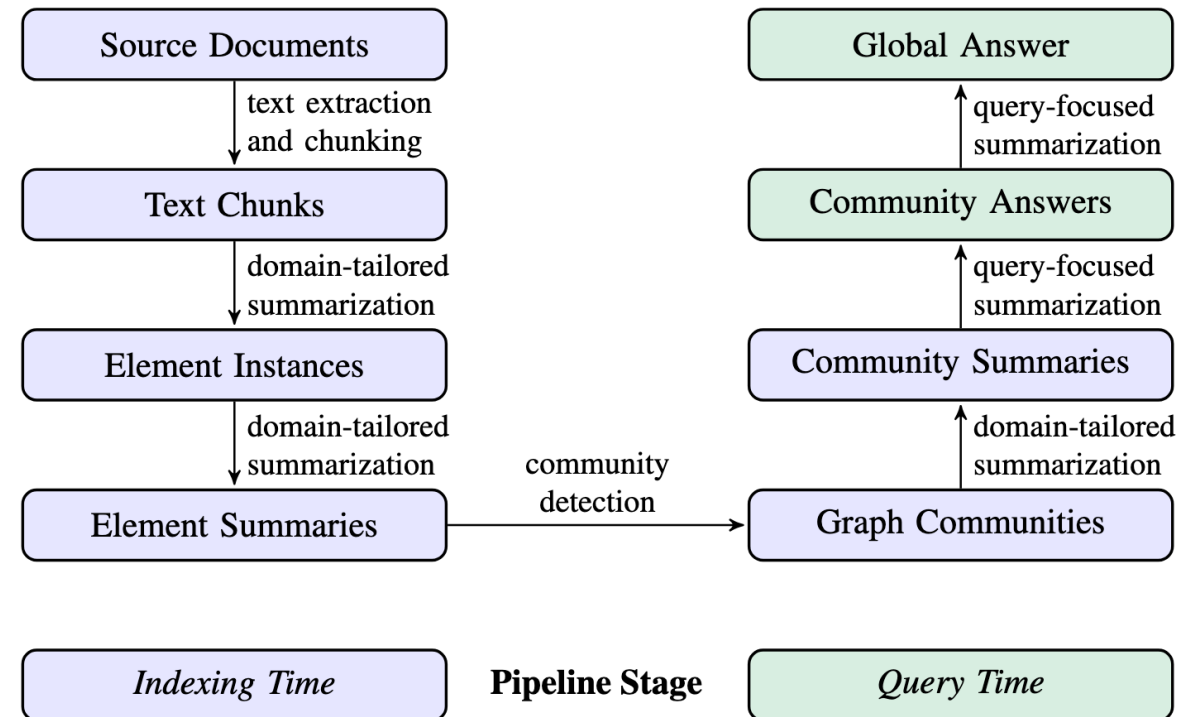
- ❑ **GraphRAG** retrieves structured evidence over entities, relations, and communities. It is better suited for global, multi-hop, and relational questions.



Vanilla RAG

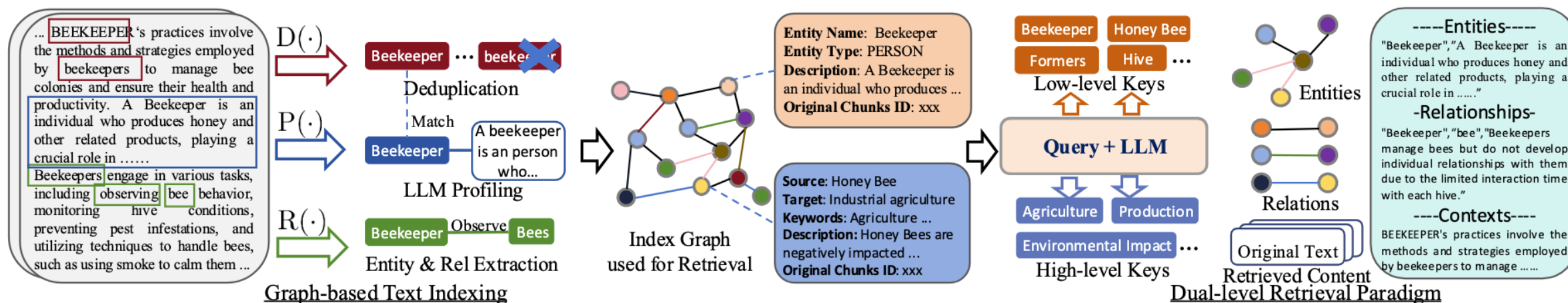


Graph-based RAG



GraphRAG for Lightweight and Incremental Retrieval

❑ **LightRAG**: efficient and incremental GraphRAG through dual-level graph retrieval



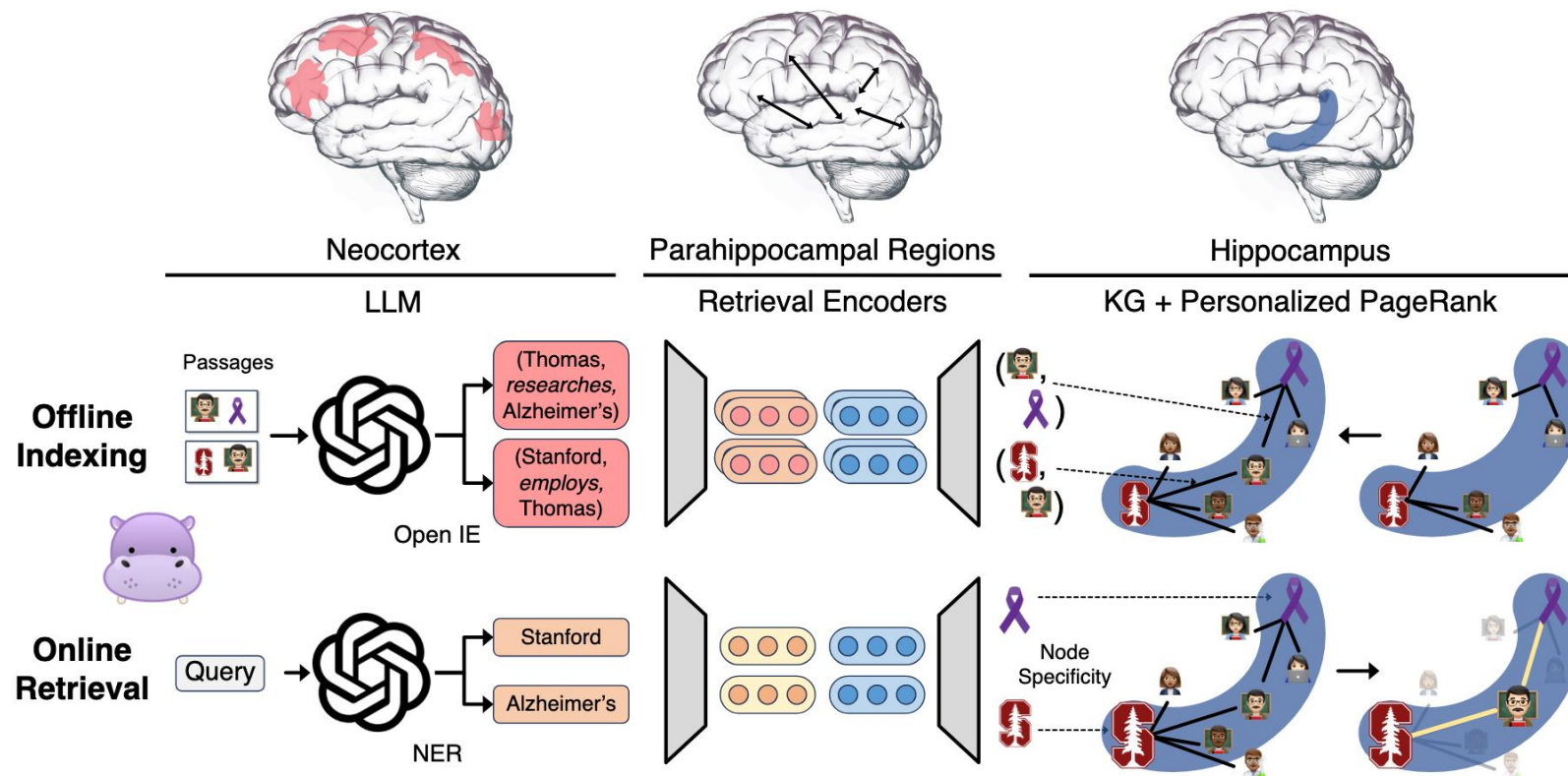
- ❖ Graph structure improves evidence organization.
- ❖ Vector retrieval improves efficiency.
- ❖ Incremental updates reduce rebuilding cost.

Phase	Retrieval Phase		Incremental Text Update	
	GraphRAG	Ours	GraphRAG	Ours
Tokens	$610 \times 1,000$	< 100	$1,399 \times 2 \times 5,000 + T_{\text{extract}}$	T_{extract}
API Calls	$\frac{610 \times 1,000}{C_{\text{max}}}$	1	$1,399 \times 2 + C_{\text{extract}}$	C_{extract}

GraphRAG as Associative Long-Term Memory

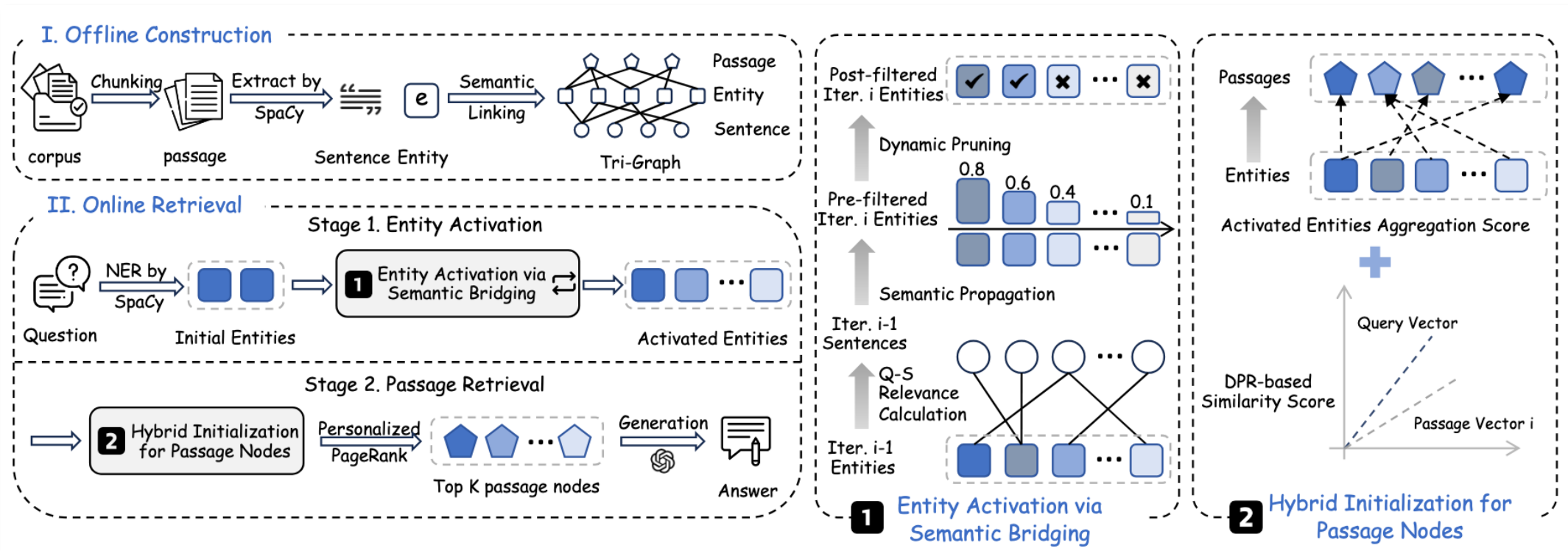
❑ **HippoRAG:** Using graph memory and associative activation for multi-hop knowledge retrieval

- ❖ Knowledge graph as associative memory
- ❖ Personalized PageRank retrieval
- ❖ Efficient multi-hop knowledge integration



GraphRAG for Scalable Graph Construction

❑ **LinearRAG:** Scaling GraphRAG by avoiding costly relation extraction on large corpora



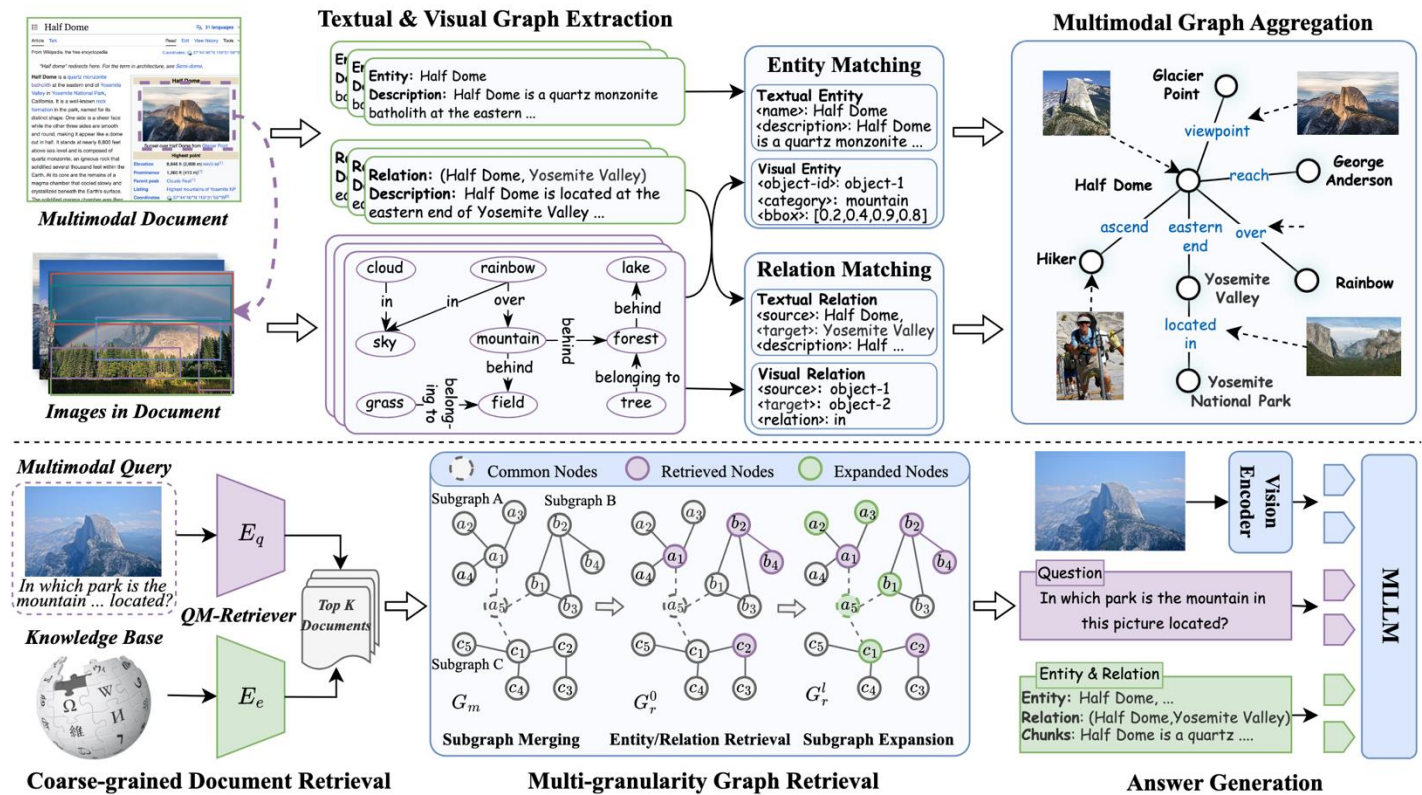
- ❖ Relation-free hierarchical graph
- ❖ Entity activation

- ❖ Global importance aggregation
- ❖ Linear-scale graph construction

Graph RAG for Multimodal Corpus

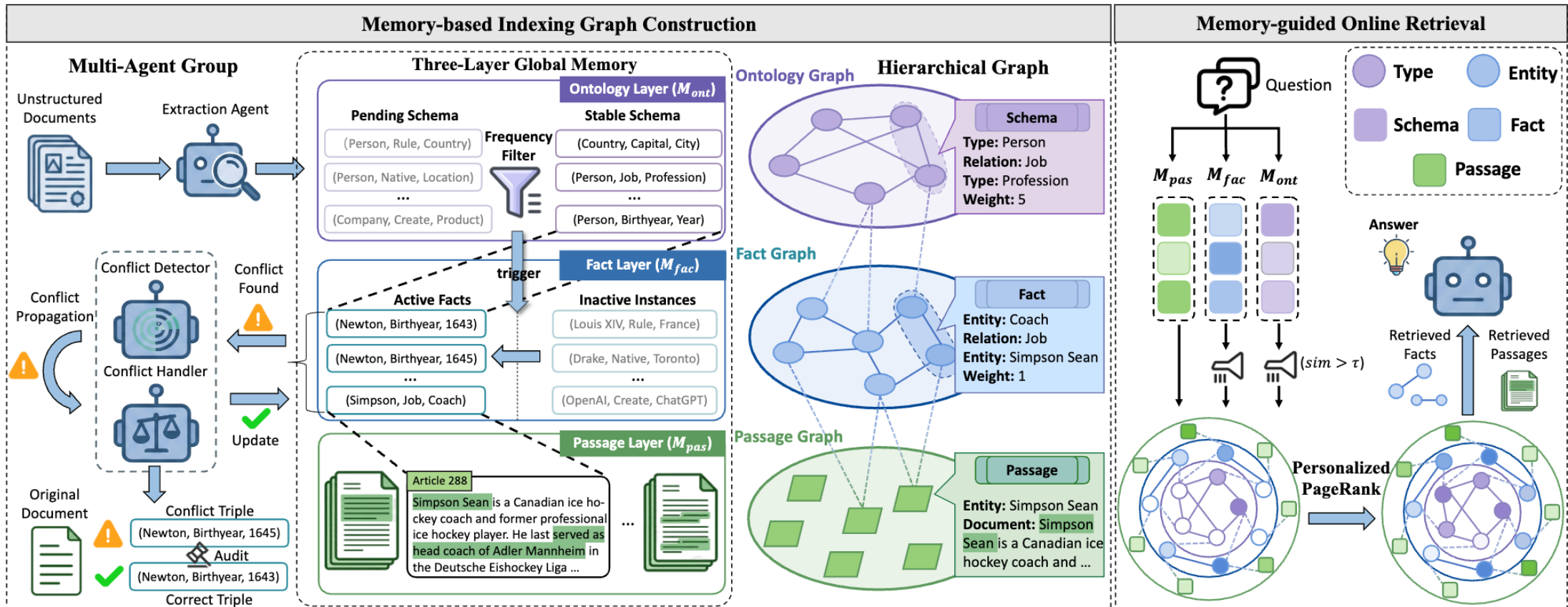
❑ mKG-RAG: Extending GraphRAG to multimodal knowledge graphs

- ❖ Knowledge-intensive VQA requires cross-modal evidence.
- ❖ Multimodal KGs align textual and visual entities and relationships
- ❖ Multi-granularity Graph retrieval provides structured multimodal context.



GraphRAG with Memory-Based Multi-Agent Construction

❑ **MemGraphRAG:** Building coherent corpus-level graphs with memory-based multi-agent construction



Tutorial Outline



- **Part 1: Introduction** of Retrieval Augmented Large Foundation Models (RA-LFMs) (Dr. Wenqi Fan)
- **Part 2: Architecture** of RA-LFMs and Main Modules (Xu Yuan)
- **Part 3: Learning Approach of RA-LFMs (Chengliang Liu)**
- **Part 4: Agentic RAG** (Chengliang Liu)
- **Part 5: Applications** of RA-LFMs (Chun-Hin Chan)
- **Part 6: Challenges and Future Directions** of RA-LFMs (Dr. Wenqi Fan)
- **Part 7: Q&A**

Website of this tutorial
Check out the slides and more information!



Part 3: RA-LFM Learning

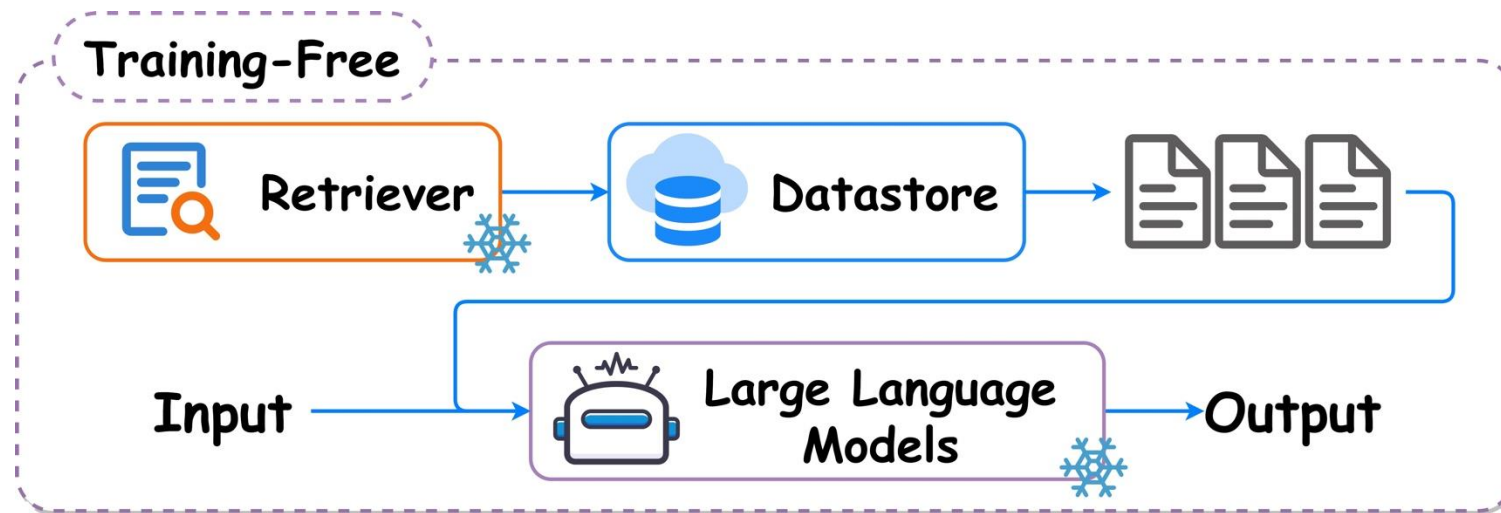


Presenter
Chengliang Liu
HK PolyU

- **Training-free Methods**
- Training-based Methods
 - Independent Learning
 - Sequential Learning
 - Joint Learning

RA-LFM Learning: Training-free

- **Retrieval models** and **language models** are both **frozen**.

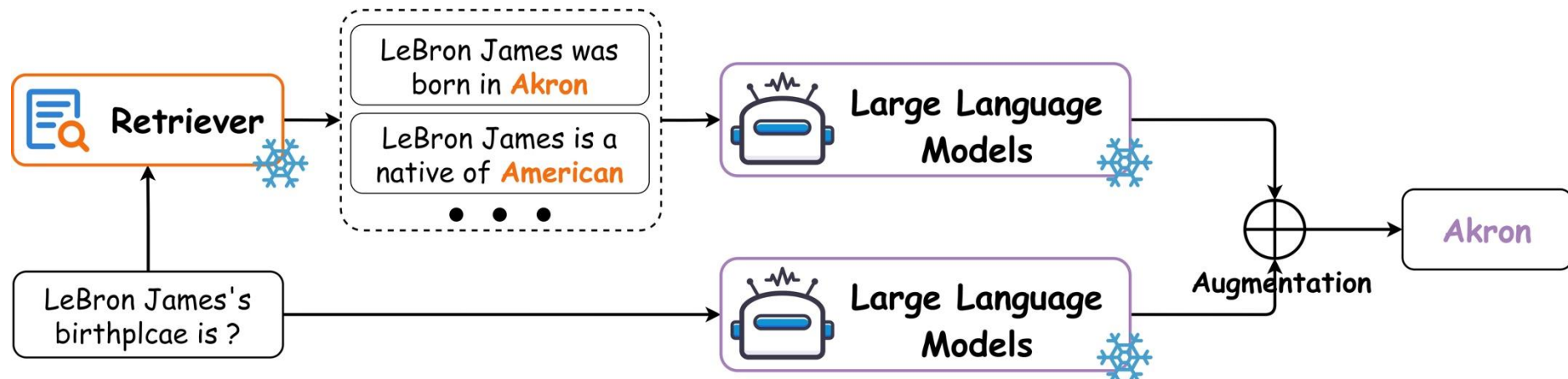


RA-LFM Learning: Training-free

- Prompt Engineering-based Methods

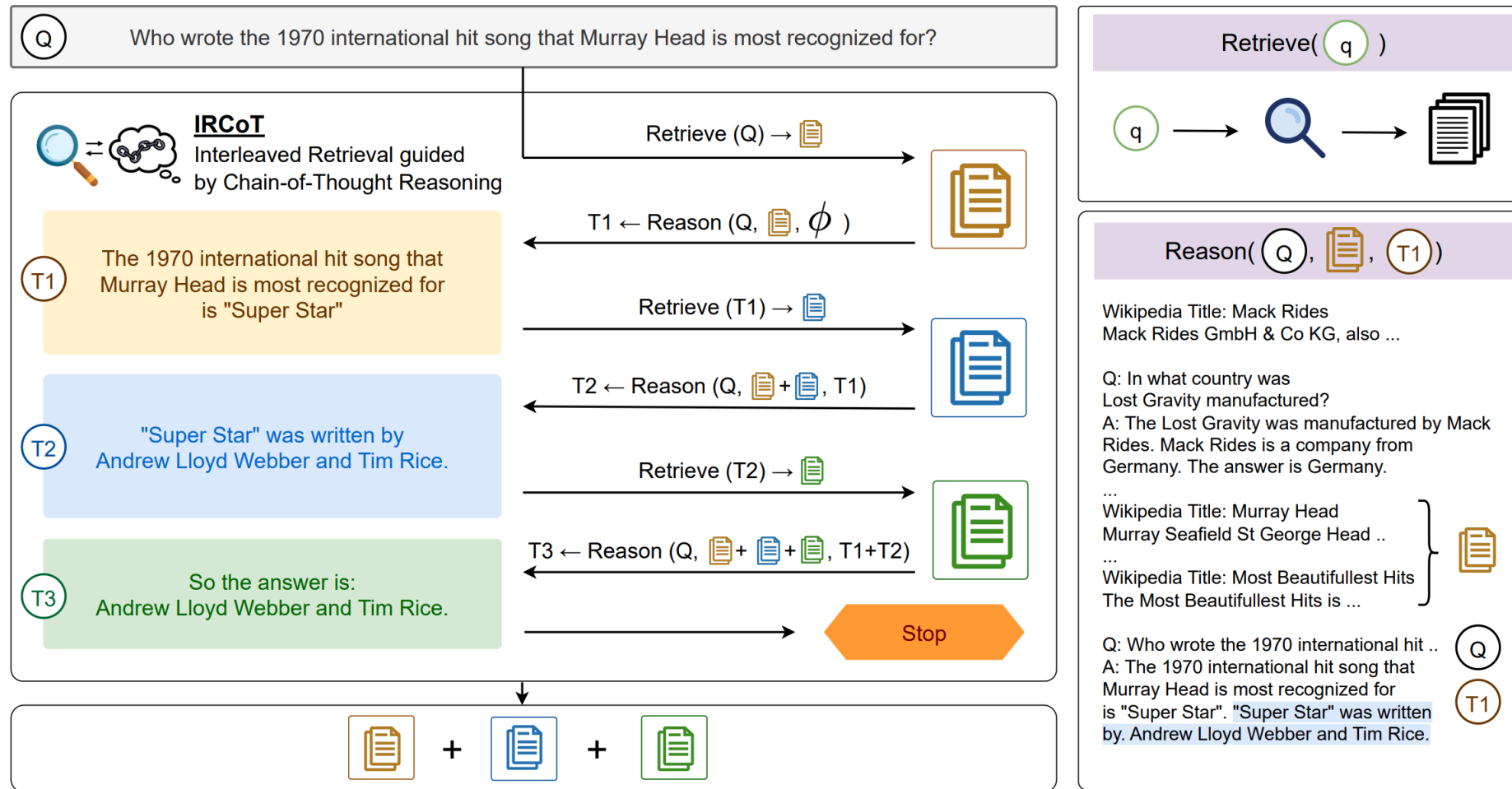


- Retrieval-Guided Token Generation Methods



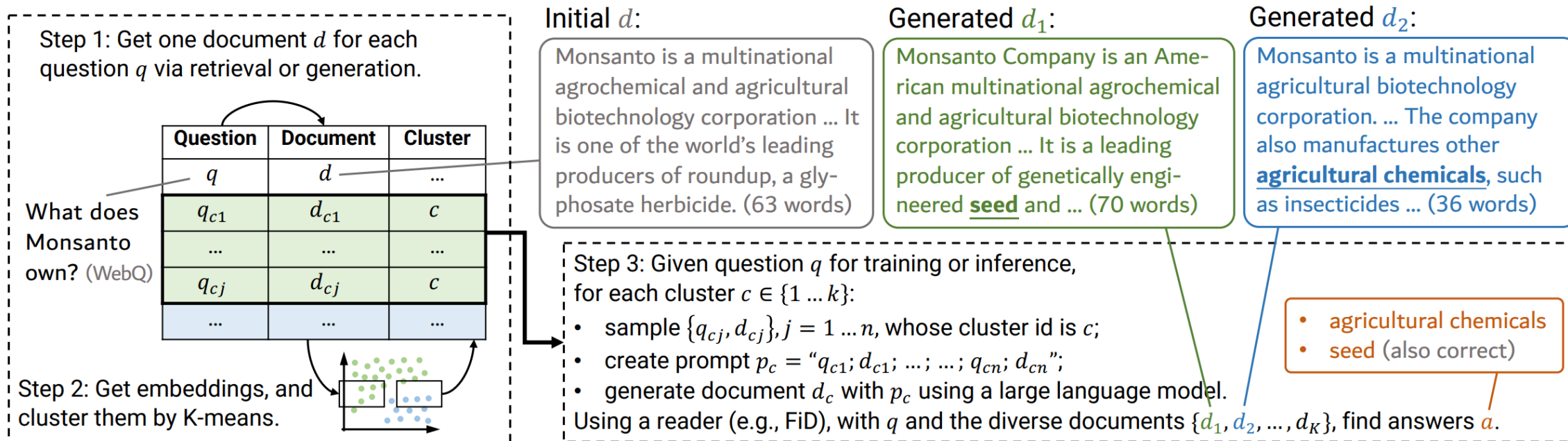
RA-LFM Learning: Training-free

- IRCoT



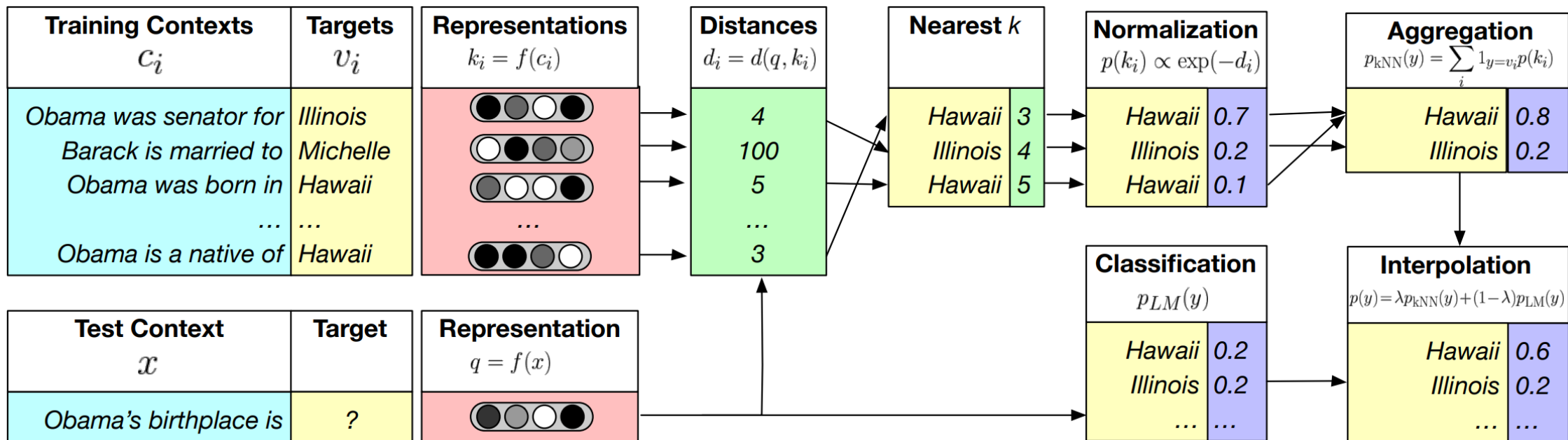
RA-LFM Learning: Training-free

- GENREAD



RA-LFM Learning: Training-free

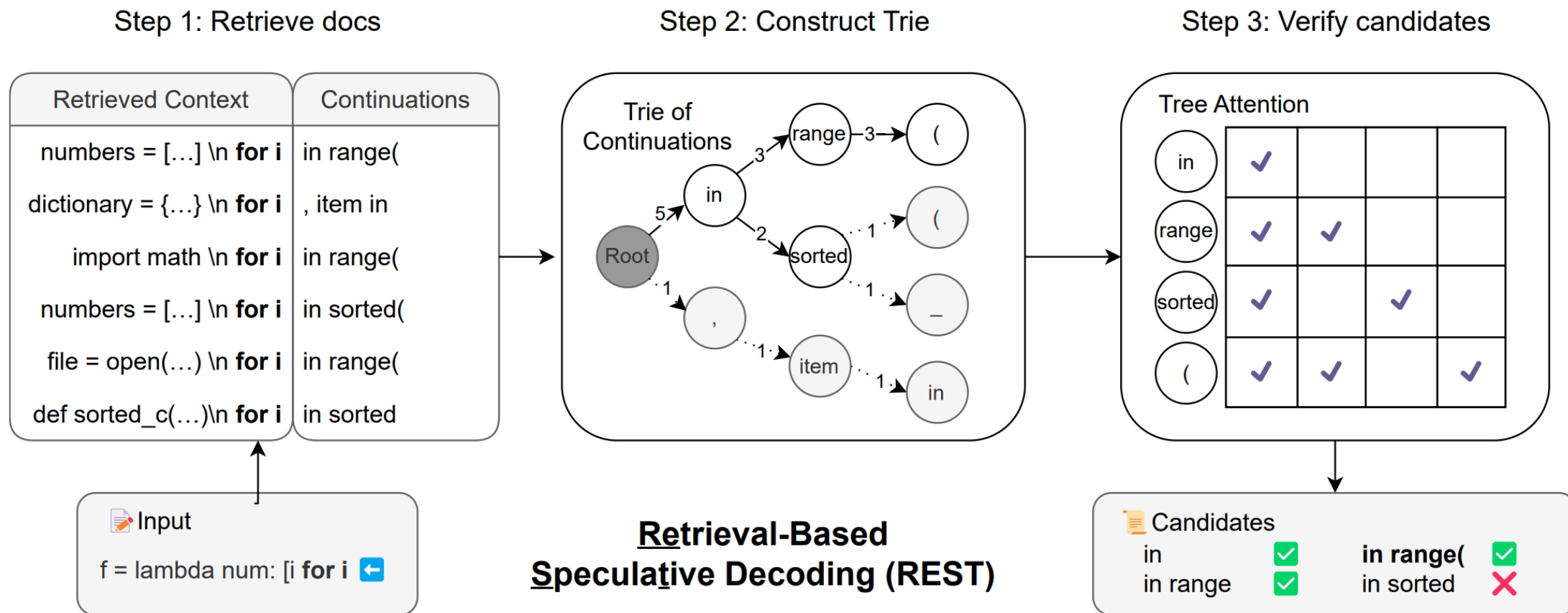
- k*NN-LM**



$$p(y|x) = \lambda p_{kNN}(y|x) + (1 - \lambda) p_{LM}(y|x)$$

RA-LFM Learning: Training-free

- REST



RA-LFM Learning: Training-free

- ✓ Work with off-the-shelf models
- x All components are fixed and not trained
- x Might not achieve optimal learning result of the whole model

Part 3: RA-LFM Learning

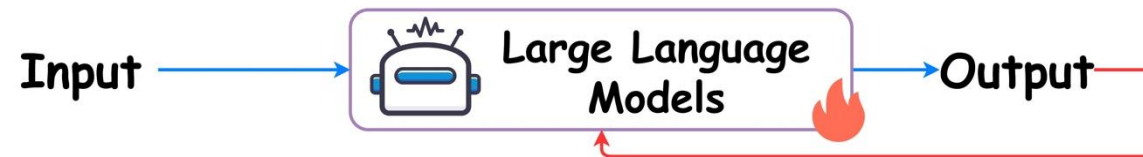


Website of this tutorial

- Training-free Methods
- Training-based Methods
 - **Independent Learning**
 - Sequential Learning
 - Joint Learning

RA-LFM Learning: Independent Training

- **Retrieval models** and **language models** are trained independently.
 - Independent training of large language models.



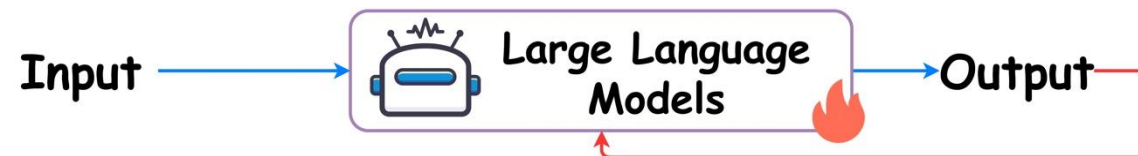
- Independent training of Retriever.



RA-LFM Learning: Independent Training

- **Retrieval models** and **language models** are trained independently.

- Independent training of large language models.

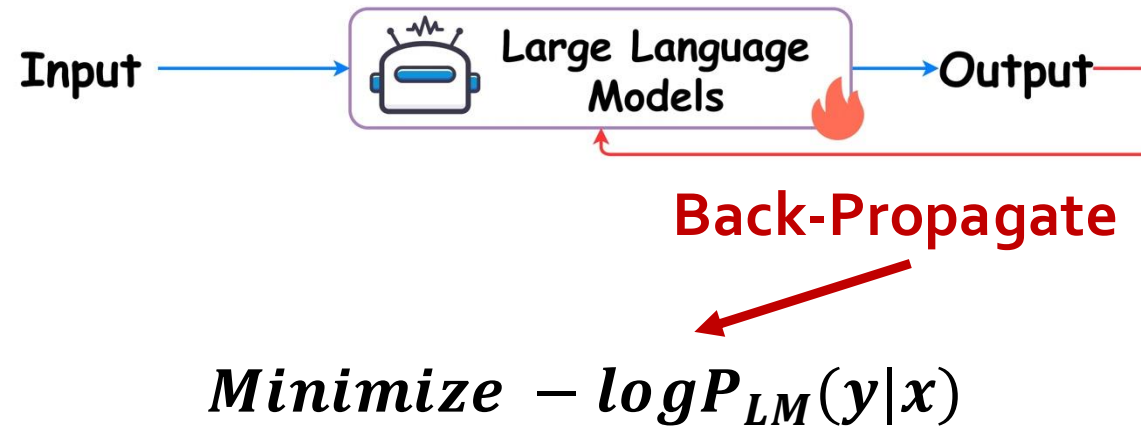


- Independent training of Retriever.



RA-LFM Learning: Independent Training

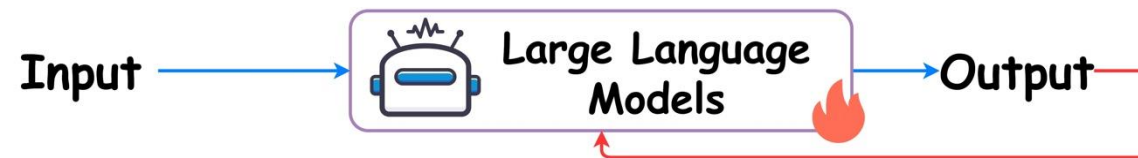
- Independent training of large language models.



.....

RA-LFM Learning: Independent Training

- **Retrieval models** and **language models** are trained independently.
 - Independent training of large language models.

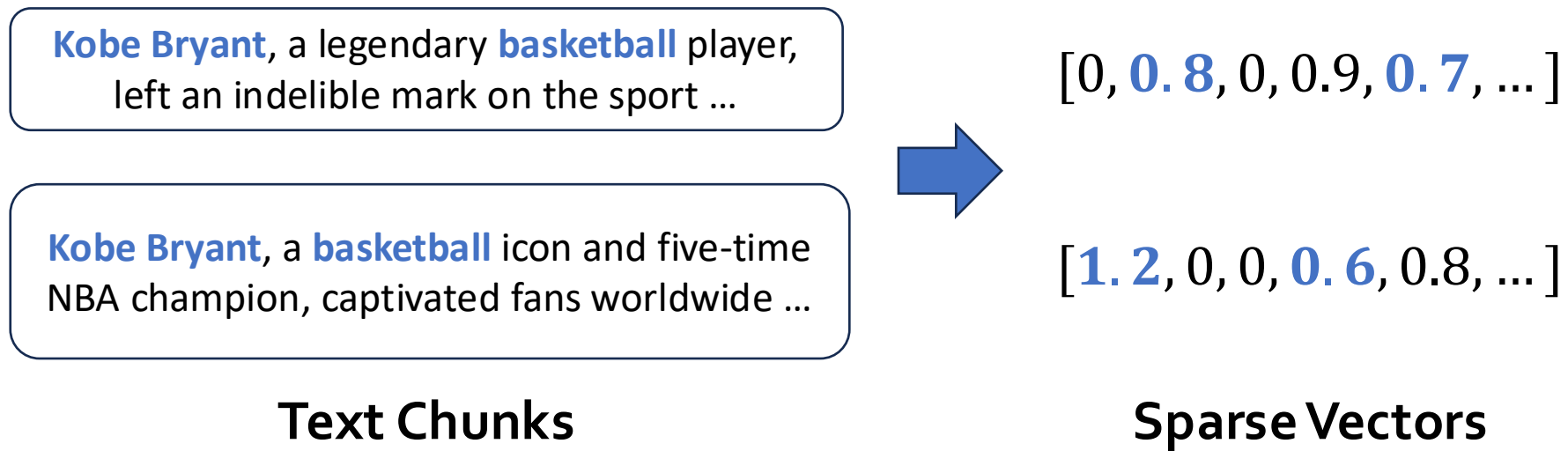


- Independent training of Retriever.



RA-LFM Learning: Independent Training

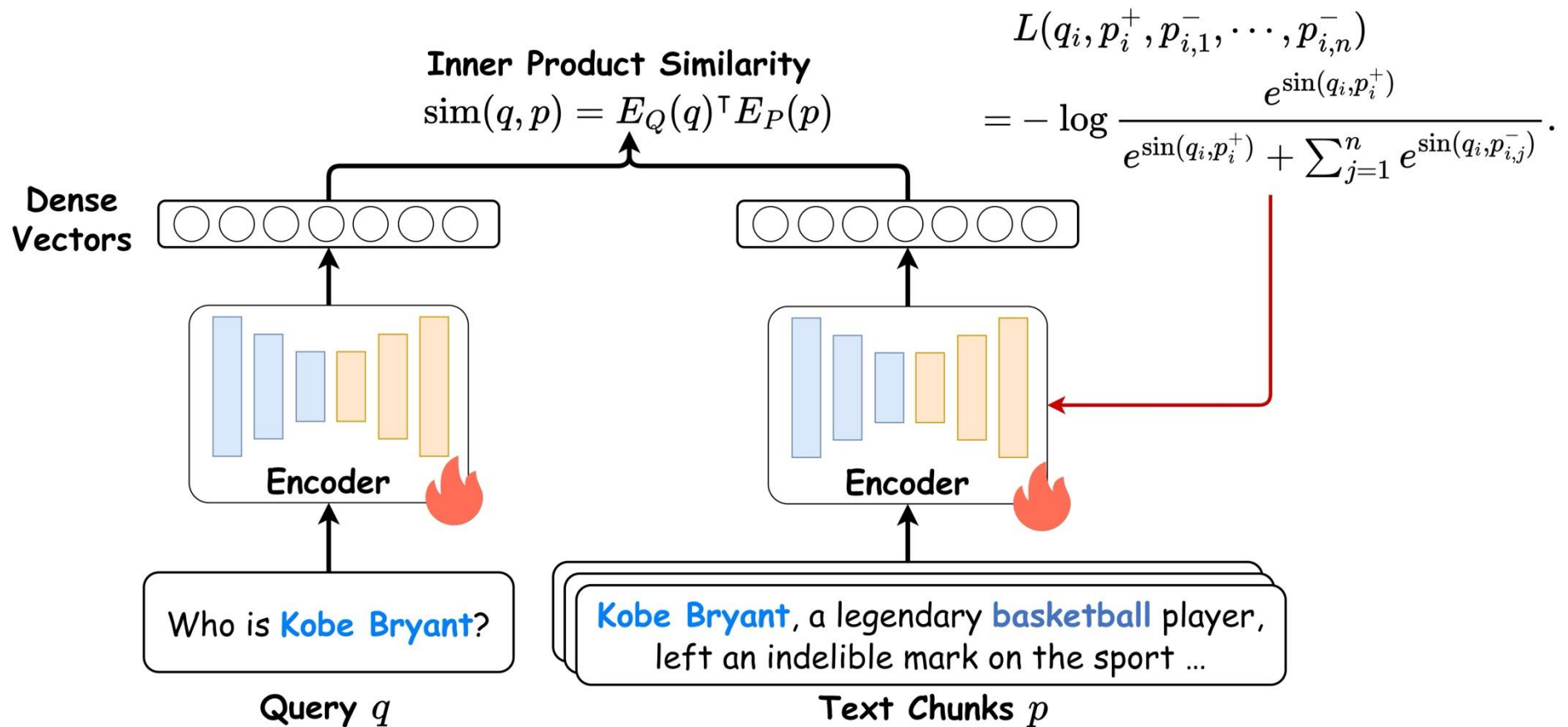
- Sparse retrieval models: TF-IDF / BM25



No training is Needed!

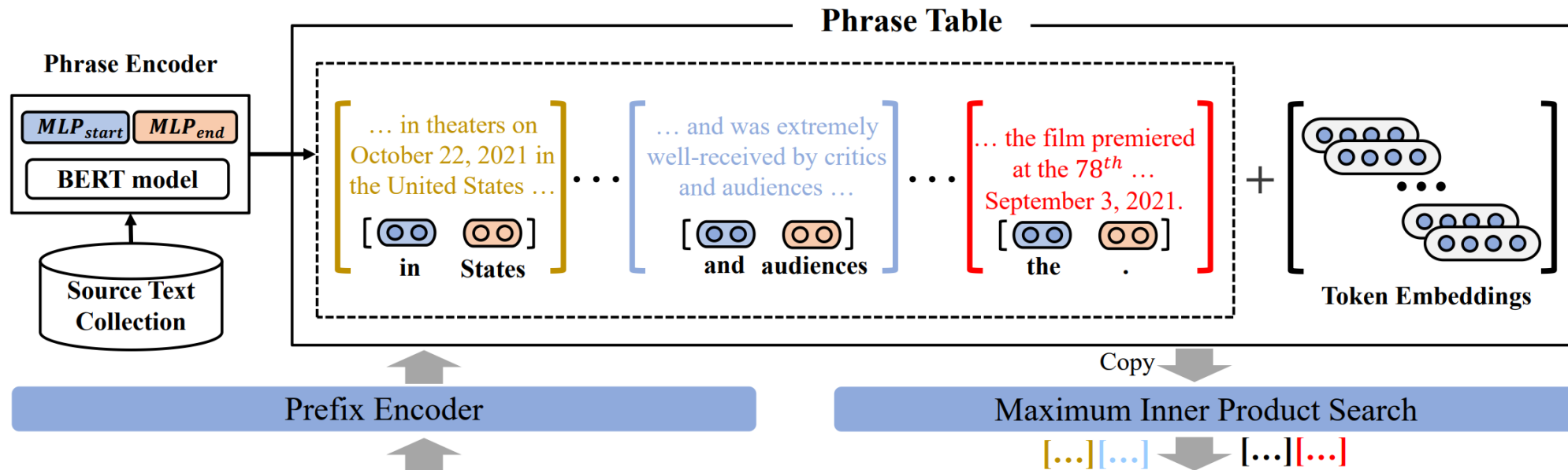
RA-LFM Learning: Independent Training

- Dense retrieval models: DPR



RA-LFM Learning: Independent Training

- Dense retrieval models: CoG



The Dune film was released [in theaters on October 22, 2021 in the United States] [and was extremely well-received by critics and audiences] [Before] [that] [,] [the film premiered at the 78th International Film Festival on September 3, 2021.]

$$\mathcal{H}_{i+1} = \text{PrefixEncoder}(x_i, \mathcal{H}_i).$$

$$\mathcal{D}_{start} = \text{MLP}_{start}(\mathcal{D}), \mathcal{D}_{end} = \text{MLP}_{end}(\mathcal{D}).$$

$$\text{PhraseEncoder}(s, e, \mathcal{D}) = [\mathcal{D}_{start}[s]; \mathcal{D}_{end}[e]] \in \mathbb{R}^d$$

RA-LFM Learning: Independent Training

- **Model Training:**

$$\mathcal{L}_p = -\frac{1}{n} \sum_{k=1}^n \log \frac{\exp(q_k \cdot p_k)}{\sum_{p \in \mathcal{P}_k} \exp(q_k \cdot p_p) + \sum_{w \in V} \exp(q_k \cdot v_w)}$$

$$\mathcal{L}_t = -\frac{1}{m} \sum_{i=1}^m \log \frac{\exp(q_i, v_{D_i})}{\sum_{w \in V} \exp(q_i, v_w)}$$

RA-LFM Learning: Independent Training

- ✓ Work with off-the-shelf models, flexible
- ✓ Each part can be improved independently
- x Lack of integrity between Retrieval and Generation
- x Retrieval models are not optimized specified for the tasks/ domains/ generators

Part 3: RA-LFM Learning

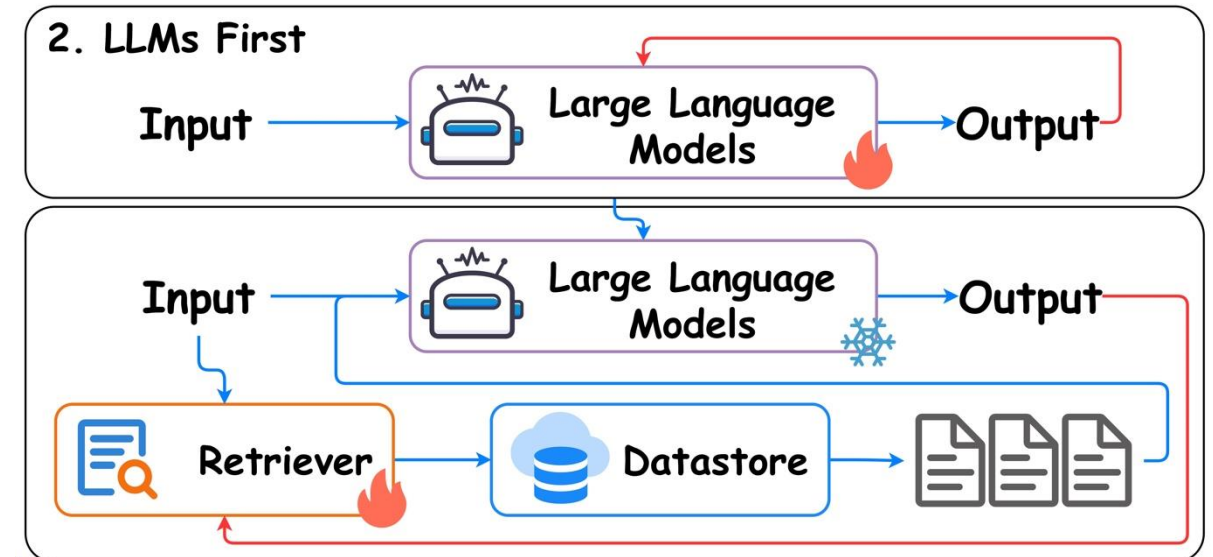
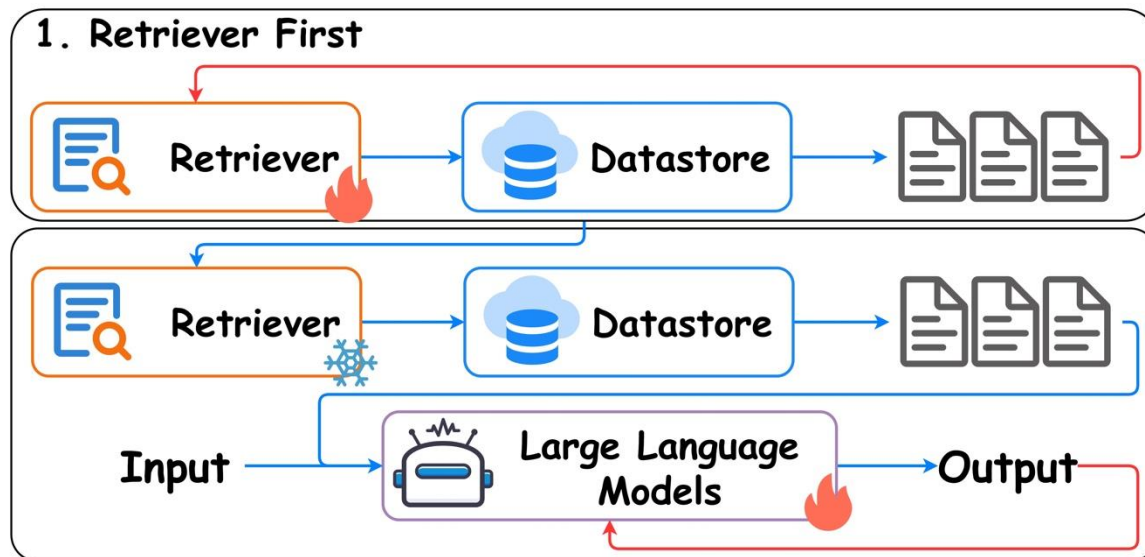


Website of this tutorial

- Training-free Methods
- Training-based Methods
 - Independent Learning
 - **Sequential Learning**
 - Joint Learning

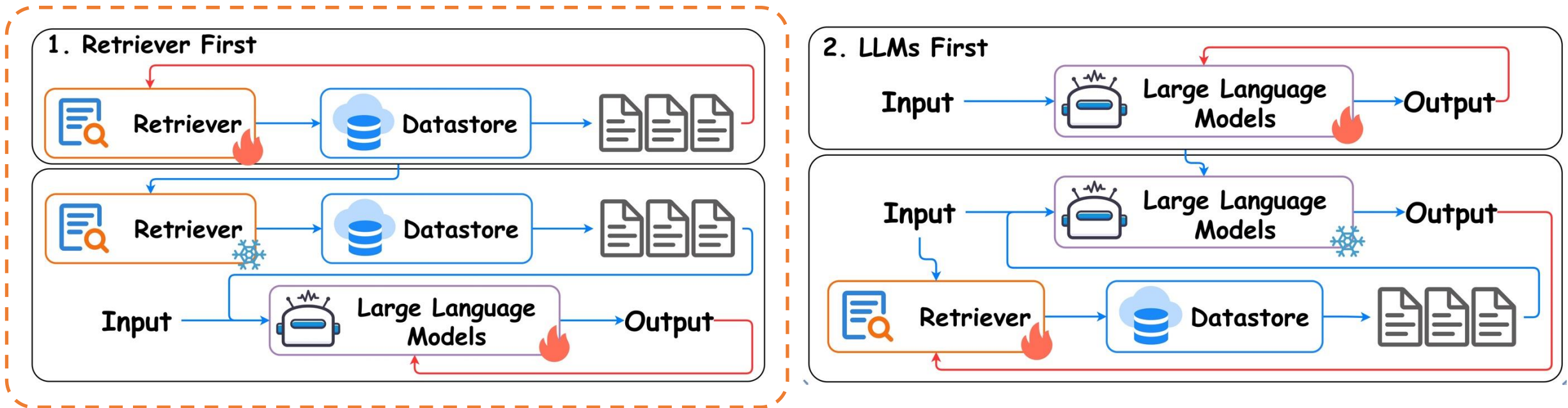
RA-LFM Learning: Sequential Training

- **One component** is first trained independently and then fixed.
- **The other component** is trained with an objective that depends on **the first one**.



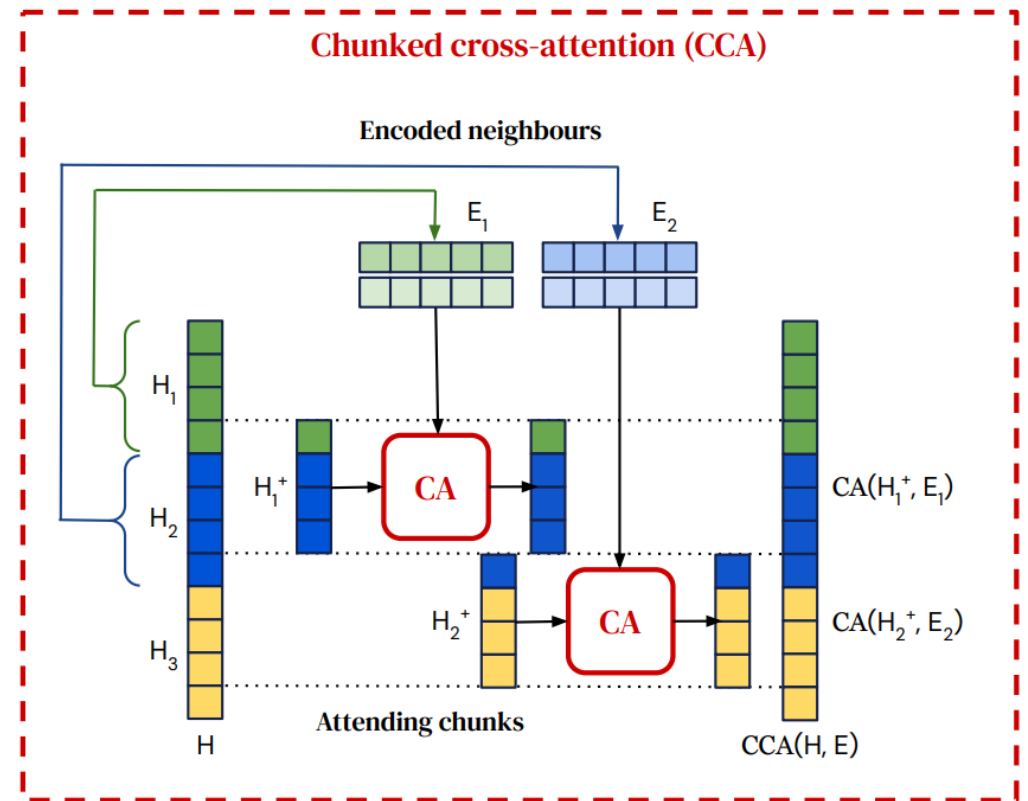
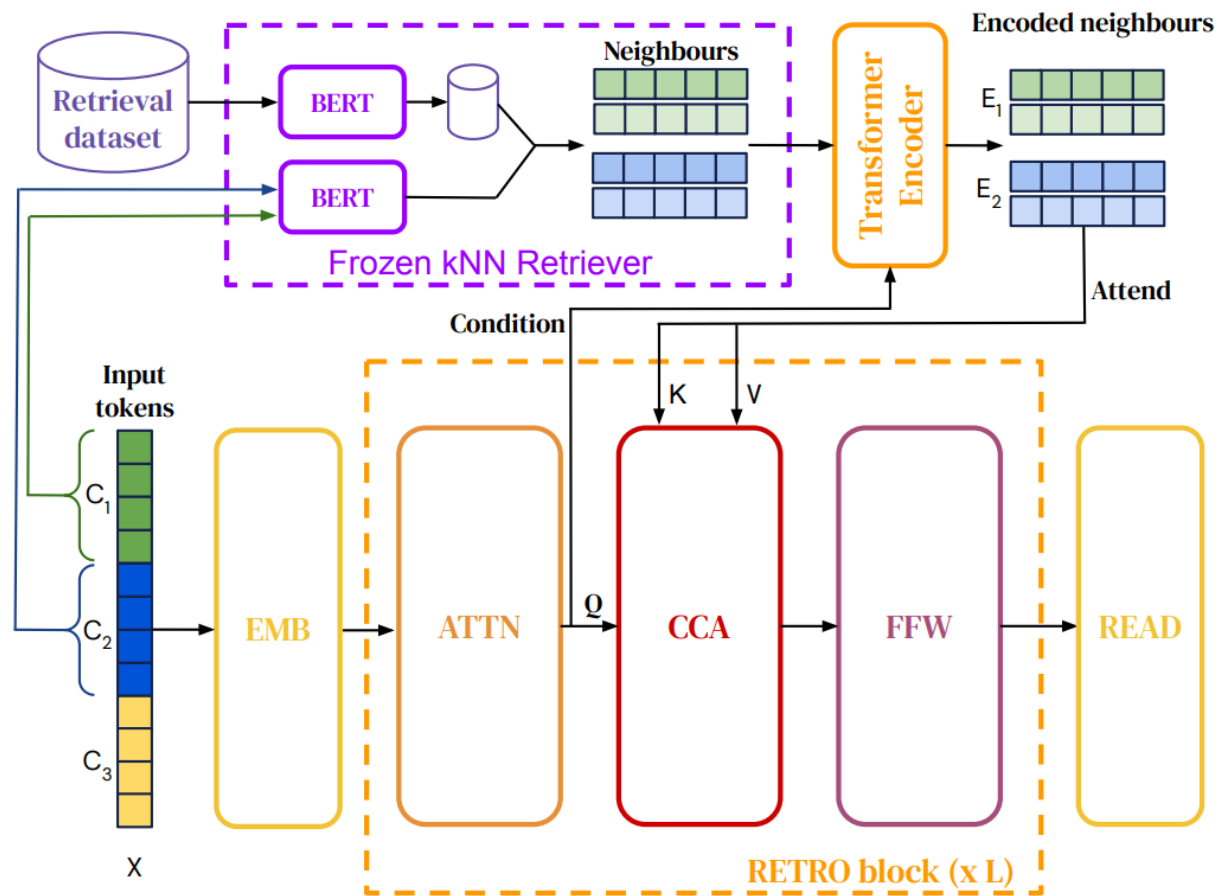
RA-LFM Learning: Sequential Training

- **Retrieval models** is first trained independently and then fixed.
- **Language models** are trained with an objective that depends on **the Retrieval**.



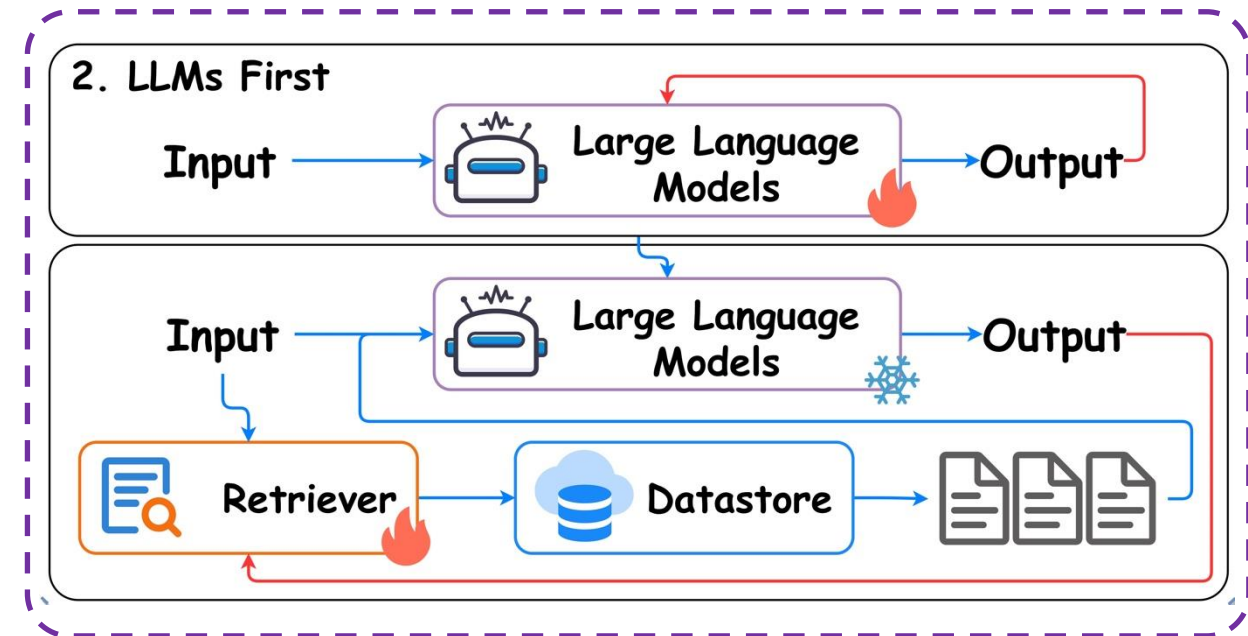
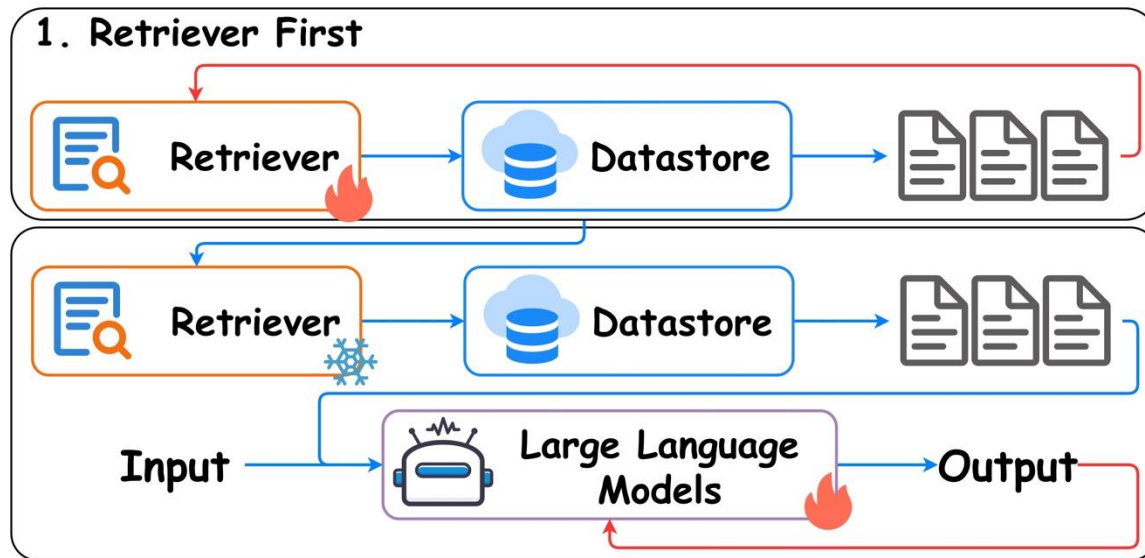
RA-LFM Learning: Sequential Training

- **RETRO**



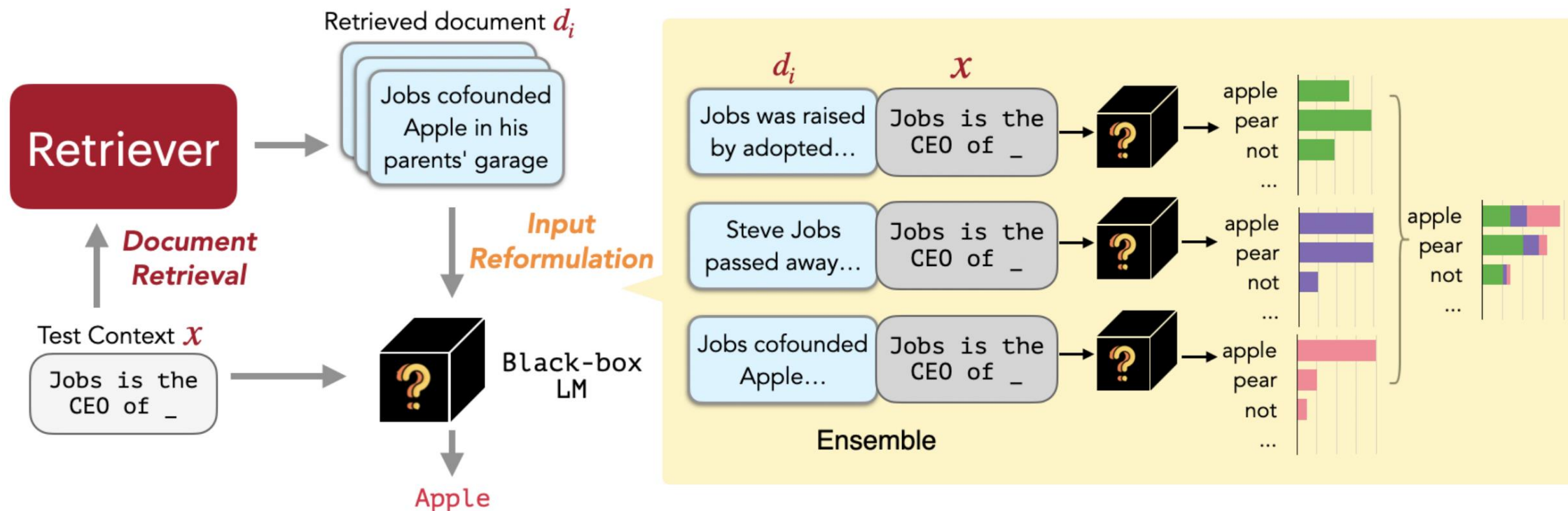
RA-LFM Learning: Sequential Training

- **Language models** are first trained independently and then fixed.
- **Retrieval models** are trained with supervisions from **language models**.



RA-LFM Learning: Sequential Training

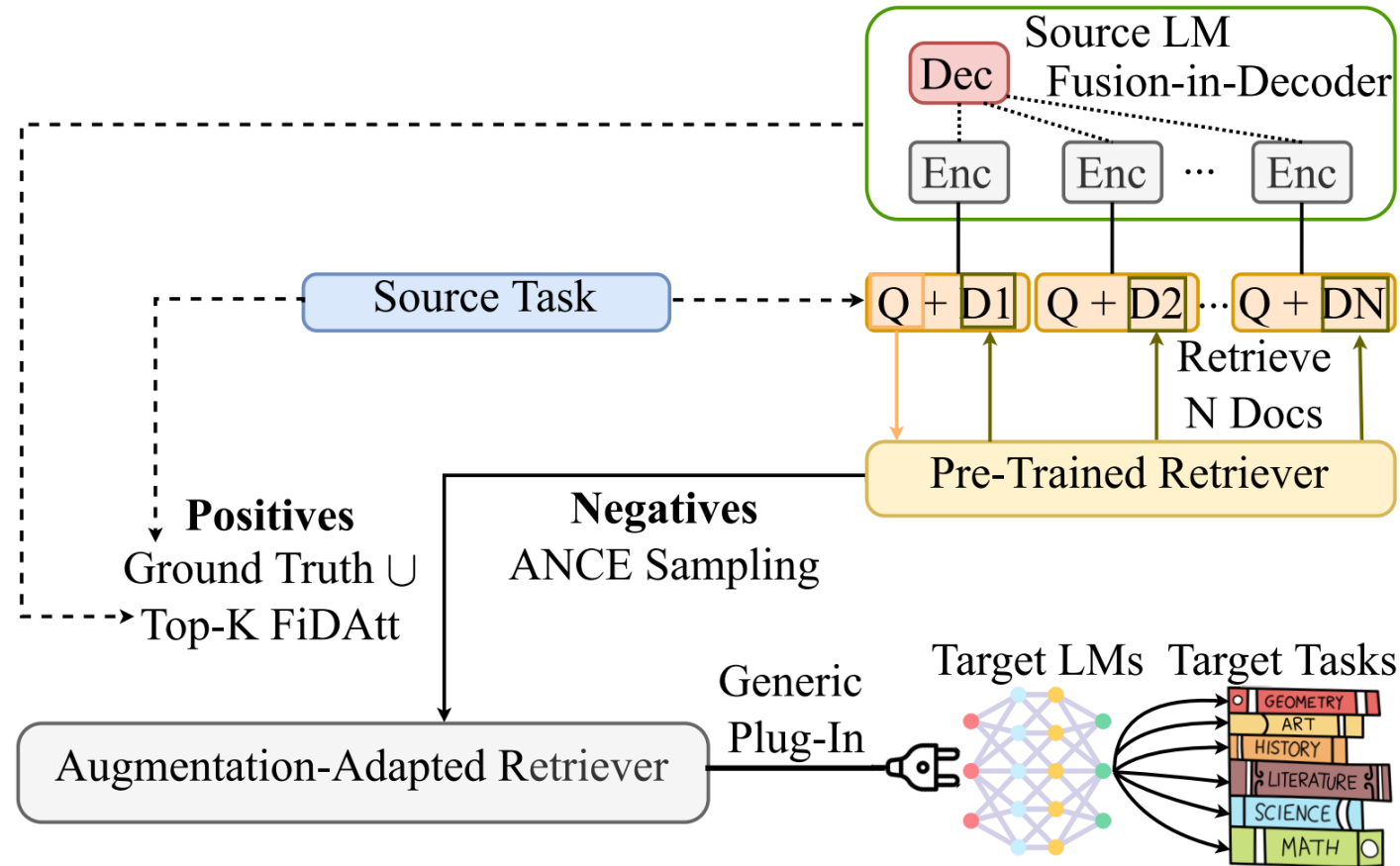
- REPLUG (Retrieve and Plug)



$$\mathcal{L} = \frac{1}{|\mathcal{B}|} \sum_{x \in \mathcal{B}} KL(P_R(d | x) \parallel Q_{LM}(d | x, y)) \quad P_R(d | x) = \frac{e^{s(d,x)/\gamma}}{\sum_{d \in \mathcal{D}'} e^{s(d,x)/\gamma}} \quad Q(d | x, y) = \frac{e^{P_{LM}(y|d,x)/\beta}}{\sum_{d \in \mathcal{D}'} e^{P_{LM}(y|d,x)/\beta}}$$

RA-LFM Learning: Sequential Training

- AAR (Augmentation-Adapted Retriever)**



$$\mathcal{L} = \sum_q \sum_{d^+ \in D^+} \sum_{d^- \in D^-} l(f(q, d^+), f(q, d^-)),$$

RA-LFM Learning: Sequential Training

- ✓ Work with off-the-shelf models
- ✓ Generators can be trained effectively based on the retrieved results
- ✓ Retrievers can be trained to provide useful information to help the generators
- x One component is still fixed and not trained
- x Might not achieve optimal learning result of the whole model

Part 3: RA-LFM Learning

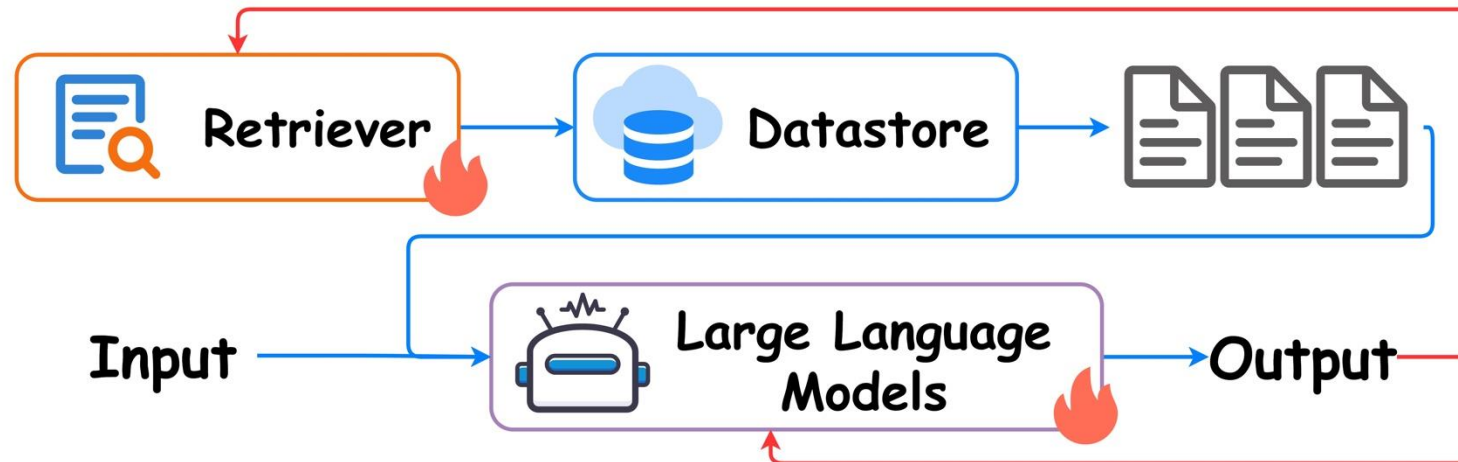


Website of this tutorial

- Training-free Methods
- Training-based Methods
 - Independent Learning
 - Sequential Learning
 - **Joint Learning**

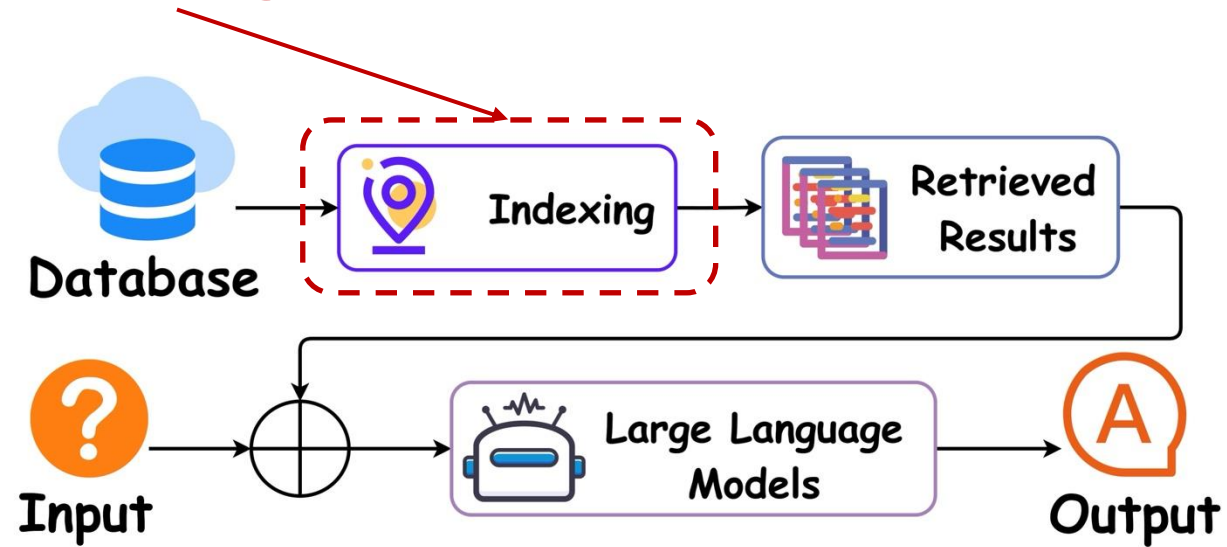
RA-LFM Learning: Joint Training

- **Retrieval models** is and **language models** are trained jointly.



RA-LFM Learning: Joint Training

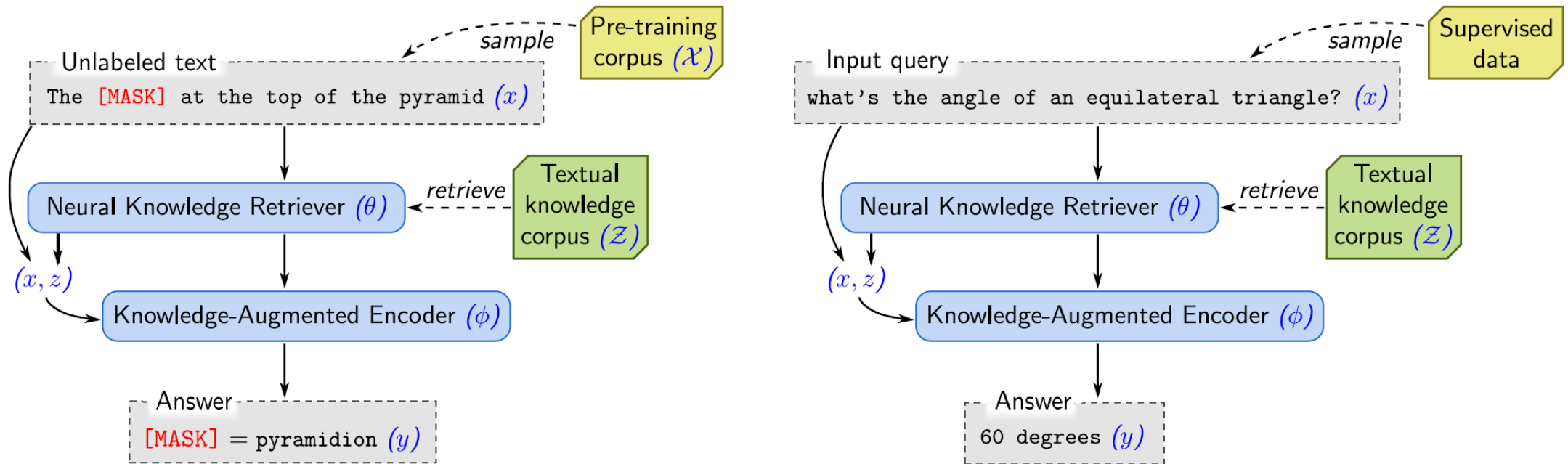
- **Retrieval Index Updating, which could be very expensive!**



- **Solutions:**
 - Asynchronous index updating
 - In-batch approximation

RA-LFM Learning: Joint Training

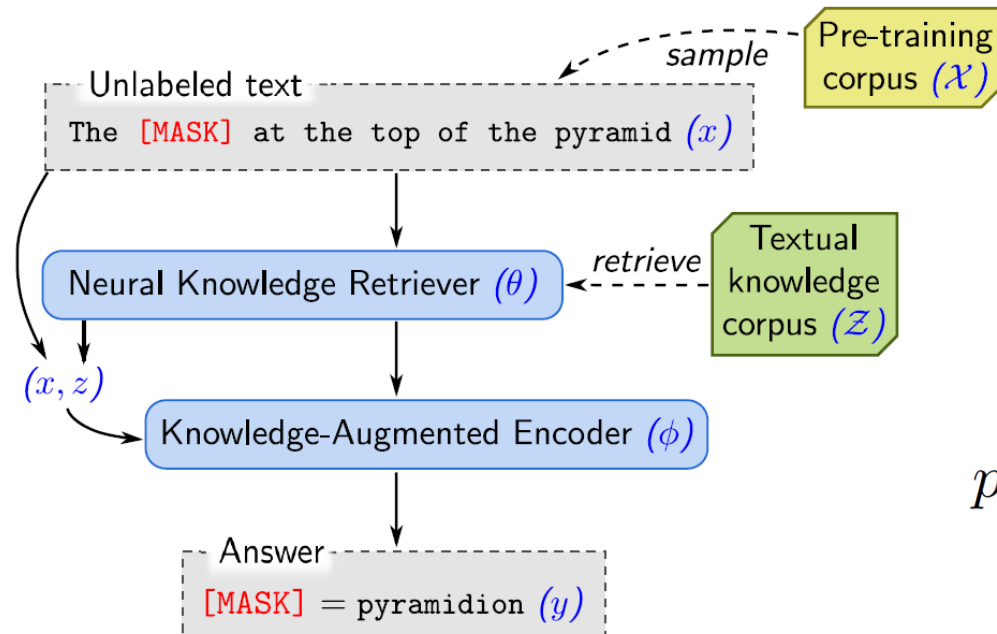
- REALM



Objective function:
$$p(y | x) = \sum_{z \in Z} p(y | z, x) p(z | x).$$

RA-LFM Learning: Joint Training

- REALM

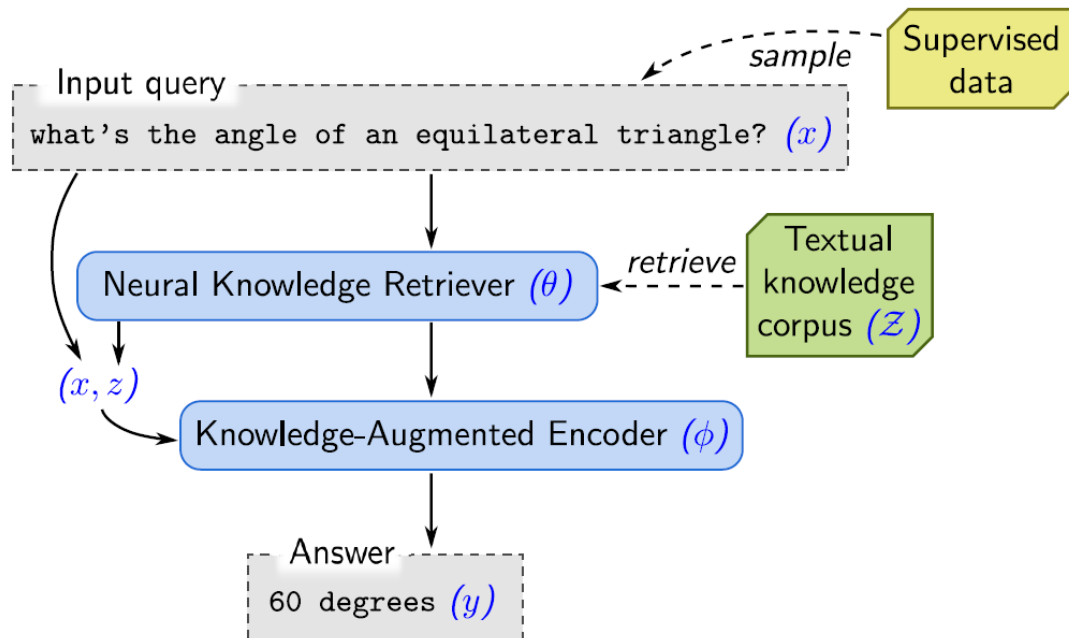


$$p(y | z, x) = \prod_{j=1}^{J_x} p(y_j | z, x)$$

$$p(y_j | z, x) \propto \exp(w_j^\top \text{BERT}_{\text{MASK}(j)}(\text{join}_{\text{BERT}}(x, z_{\text{body}})))$$

RA-LFM Learning: Joint Training

- REALM



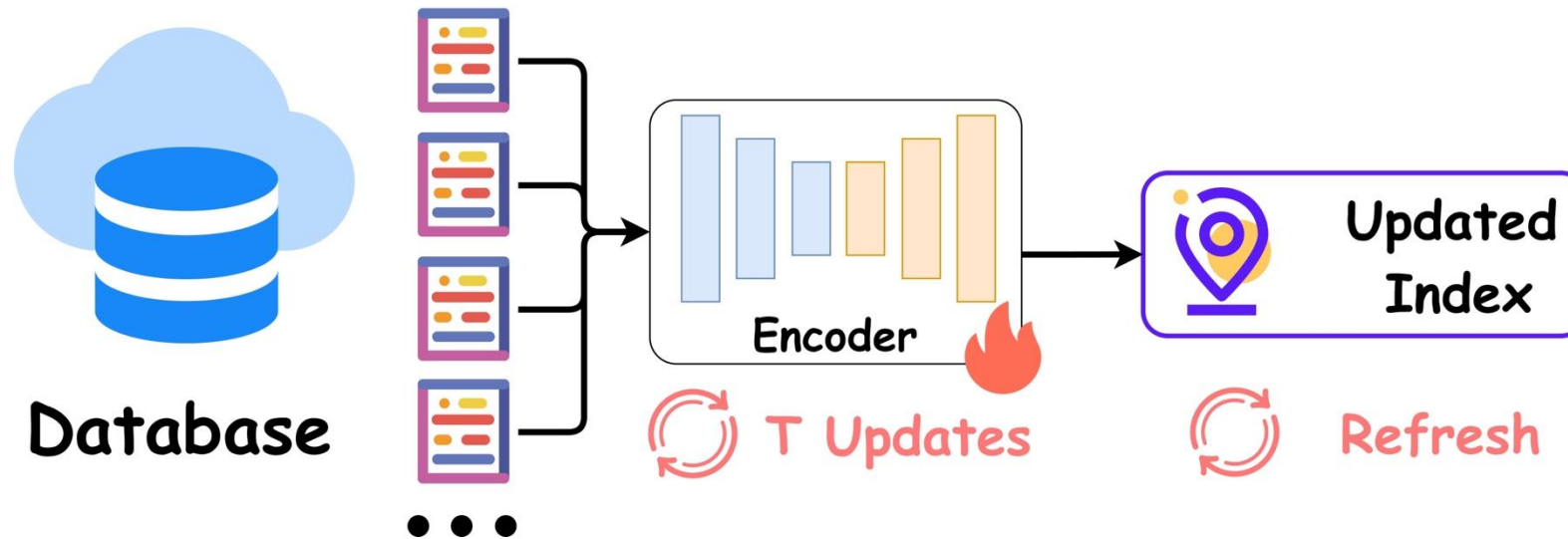
$$p(y | z, x) \propto \sum_{s \in S(z, y)} \exp(\text{MLP}([h_{\text{START}(s)}; h_{\text{END}(s)}]))$$

$$h_{\text{START}(s)} = \text{BERT}_{\text{START}(s)}(\text{join}_{\text{BERT}}(x, z_{\text{body}})),$$

$$h_{\text{END}(s)} = \text{BERT}_{\text{END}(s)}(\text{join}_{\text{BERT}}(x, z_{\text{body}})),$$

RA-LFM Learning: Joint Training

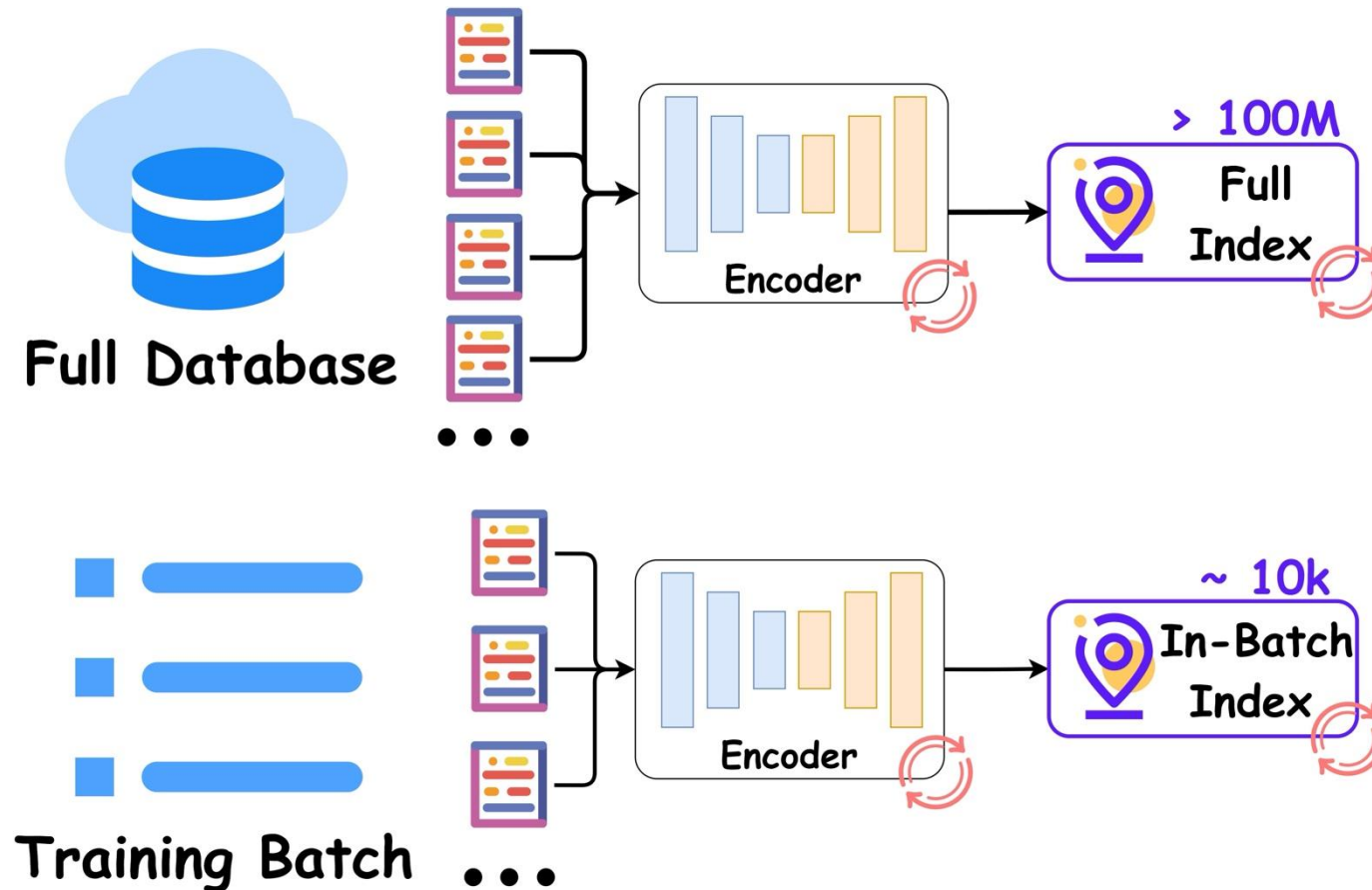
- REALM – Asynchronous Index Update



$$f(x, z) = \text{Embed}_{\text{input}}(x)^\top \text{Embed}_{\text{doc}}(z)$$

RA-LFM Learning : Joint Training

- TRIME – In-Batch Approximation

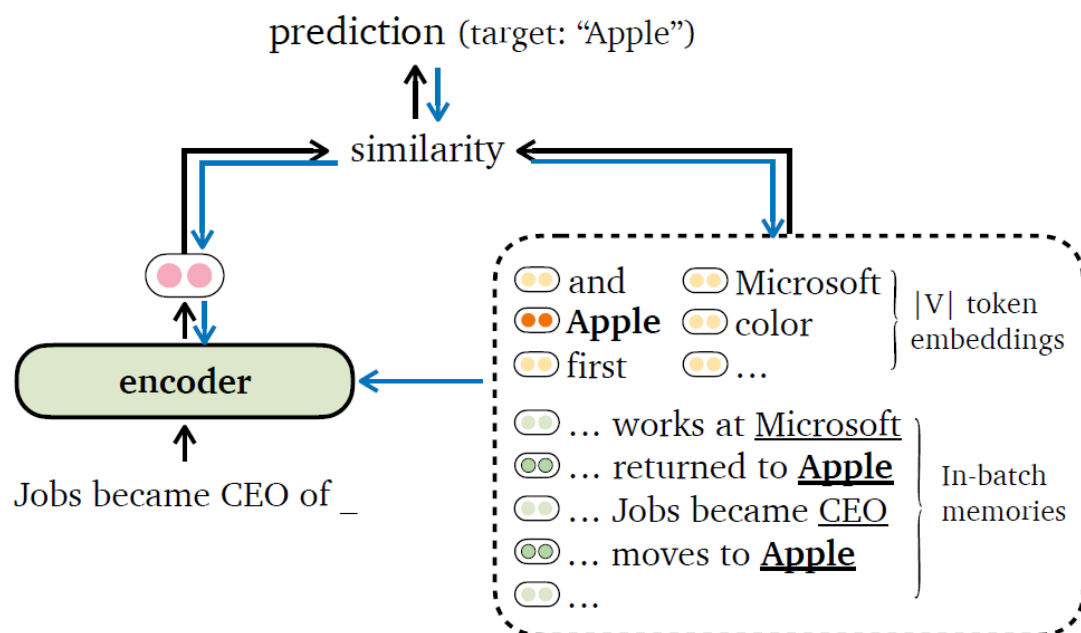


RA-LFM Learning : Joint Training

• TRIME

- Target token's embedding
- Other token embeddings
- Positive in-batch memory
- Negative in-batch memory

↑ Forward pass ↓ Back-propagation



Local Memory: $\mathcal{M}_{\text{local}}(c_t) = \{(c_j, x_j)\}_{1 \leq j \leq t-1}$.

Long-term Memory:

$$\mathcal{M}_{\text{long}}(c_t^{(i)}) = \{(c_j^{(k)}, x_j^{(k)})\}_{1 \leq k < i, 1 \leq j}$$

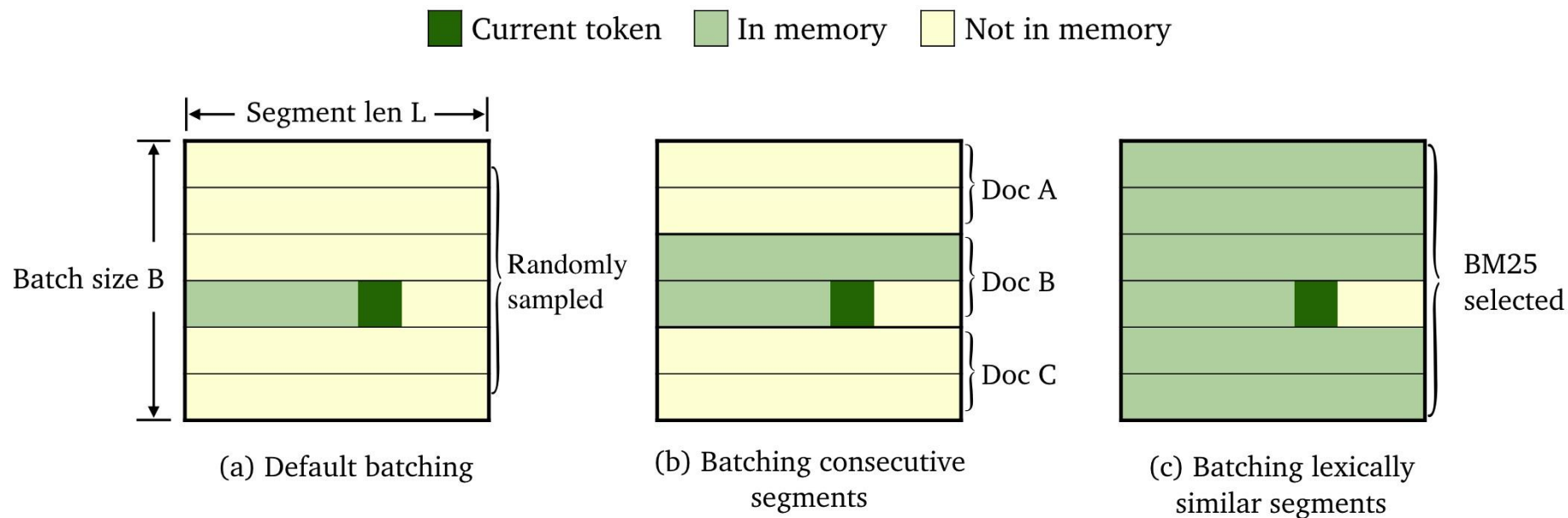
External Memory: $\mathcal{M}_{\text{ext}} = \{(c_j, x_j) \in \mathcal{D}\}$.

Training Objective:

$$P(w | c) \propto \exp(E_w^\top f_\theta(c)) + \sum_{(c_j, x_j) \in \mathcal{M}_{\text{train}}: x_j = w} \exp(\text{sim}(g_\theta(c), g_\theta(c_j))).$$

RA-LFM Learning : Joint Training

- **TRIME Data Batching Strategy**



Use BM25 scores to find similar text chunks to provide more training signals

Tutorial Outline



- **Part 1: Introduction** of Retrieval Augmented Large Foundation Models (RA-LFMs) (Dr. Yujuan Ding)
- **Part 2: Architecture** of RA-LFMs and **Main Modules** (Xu Yuan)
- **Part 3: Learning Approach** of RA-LFMs (Chengliang Liu)
- **Part 4: Agentic RAG** (Chengliang Liu)
- **Part 5: Applications** of RA-LFMs (Chun-Hin Chan)
- **Part 6: Challenges and Future Directions** of RA-LFMs (Dr. Yujuan Ding)
- **Part 7: Q&A**

Website of this tutorial
Check out the slides and more information!



Part 4: Agentic RAG



Presenter
Chengliang Liu
HK PolyU

- **Motivation, Definition & Key Features**
- Representative Methods
 - Prompting
 - Single-Agent
 - Multi-Agent
 - RL-Driven

Motivation, Definition & Key Features

Static pipeline: Query → Retrieve → Reason → Answer (fixed, one-shot)

Three limitations of this paradigm:

- 1. Retrieval Adequacy:** pre-retrieved knowledge may not cover needs that emerge during reasoning
- 2. Reasoning Depth:** noisy or wrong retrieval interferes with multi-step inference
- 3. Adaptability:** no mechanism to adjust retrieval strategy mid-reasoning

Motivation, Definition & Key Features

Agentic RAG: Query \rightarrow [Agent loop: Reason \leftrightarrow Retrieve \leftrightarrow Evaluate] \rightarrow Answer

1. Reasoning-Driven Retrieval

Reasoning process triggers and guides retrieval direction

2. Iterative Interleaving

Retrieval and reasoning alternate in multiple rounds (RAG \Leftrightarrow Reasoning closed loop)

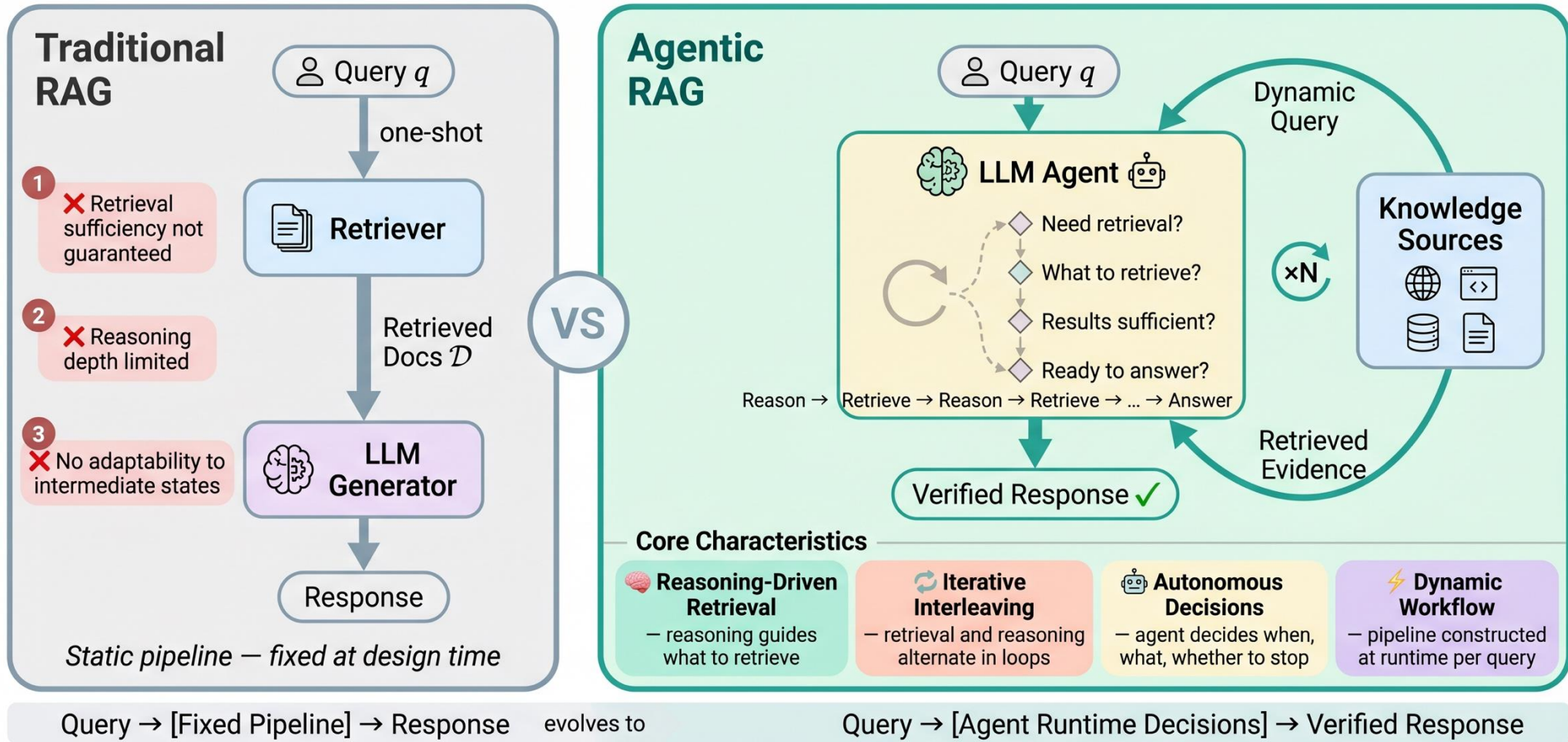
3. Autonomous Decision-Making

Agent judges at runtime: when to retrieve, what to retrieve, when to stop

4. Dynamic Workflow

Execution path constructed per query, not pre-defined

Motivation, Definition & Key Features

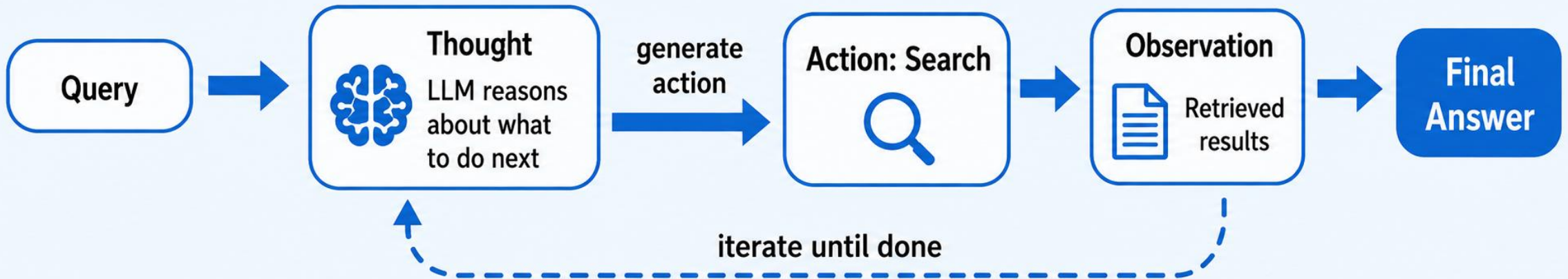


Representative Methods

Approach	How Agentic Behavior Arises	Key Characteristics
Prompting	Prompt design elicits reasoning-action loops	Zero training; flexible; prompt-dependent
Single-Agent SFT	Train model to internalize retrieval decisions	On-demand retrieval; self-reflection
Multi-Agent	Specialized agents collaborate on sub-tasks	Parallel; modular; heterogeneous sources
RL-Driven	Reward signals let model discover search strategies	Emergent behaviors; strong generalization

Prompt Chaining Agentic RAG

Prompting Agentic RAG



◆ No Training Required

◆ ReAct Loop: Think → Act → Observe

◆ Simple but Limited Flexibility

ReAct, Search-O1, Self-Ask, DeepRAG

Prompt Chaining Agentic RAG

(1) Hotspot QA

Question: Aside from the Apple Remote, what other device can control the program Apple Remote was originally designed to interact with?

(1a) Standard

Answer: iPod



(1b) CoT (Reason Only)

Thought: Let's think step by step. Apple Remote was originally designed to interact with Apple TV. Apple TV can be controlled by iPhone, iPad, and iPod Touch. So the answer is iPhone, iPad, and iPod Touch.

Answer: iPhone, iPad, iPod Touch



(1c) Act-Only

Act 1: Search[Apple Remote]
Obs 1: The Apple Remote is a remote control ...

Act 2: Search[Front Row]
Obs 2: Could not find [Front Row]. Similar: ...

Act 3: Search[Front Row (software)]
Obs 3: Front Row is a discontinued media center software ...

Act 4: Finish[yes]



(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 ...

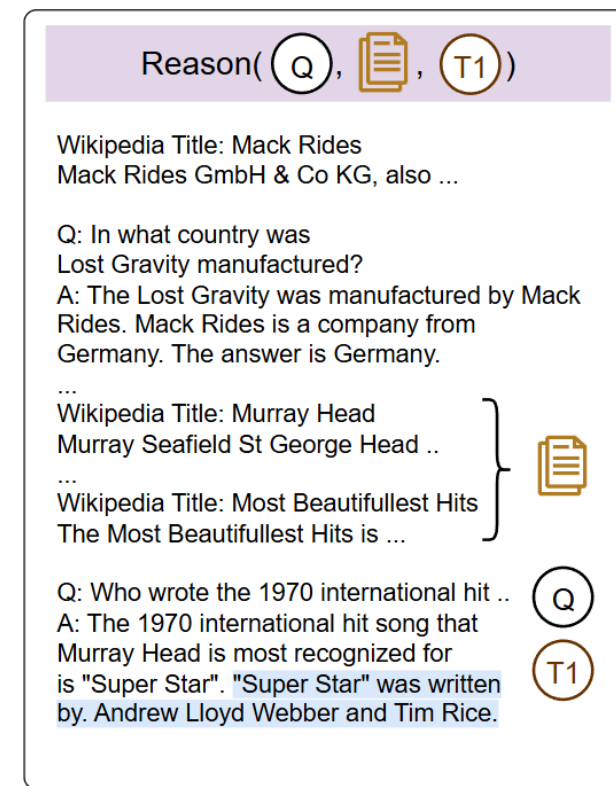
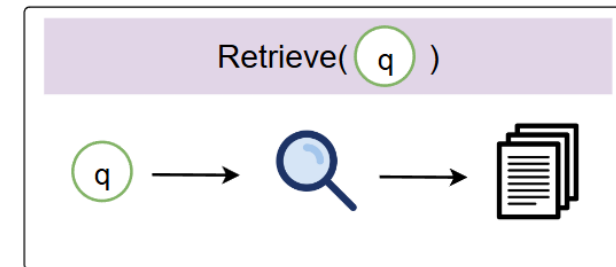
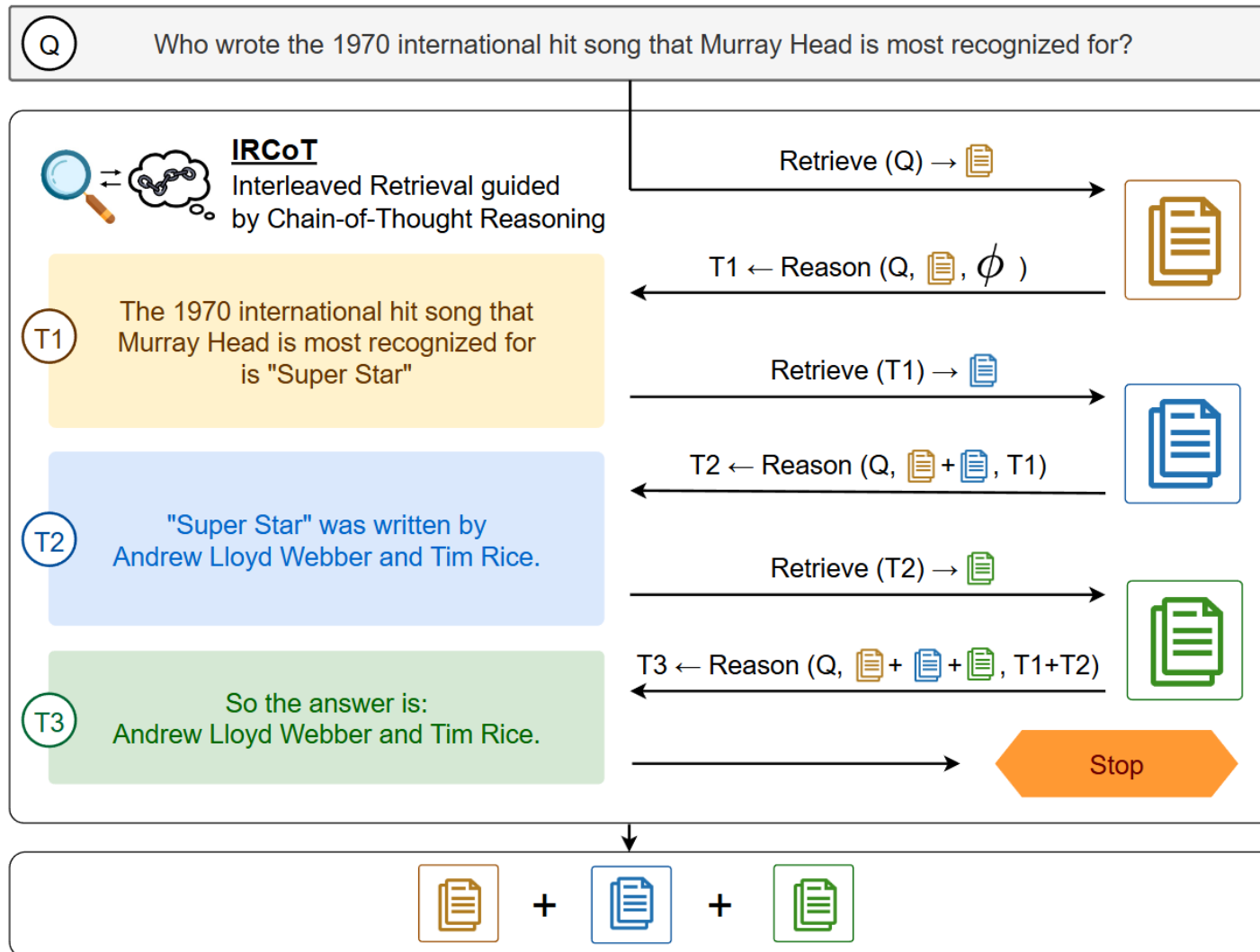
Thought 2: Apple Remote was originally designed to control the Front Row media center program. I need to search Front Row next and find what other device can control it.
Act 2: Search[Front Row]
Obs 2: Could not find [Front Row]. Similar: ['Front Row Seat to Earth', 'Front Row Motorsports', 'Front Row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



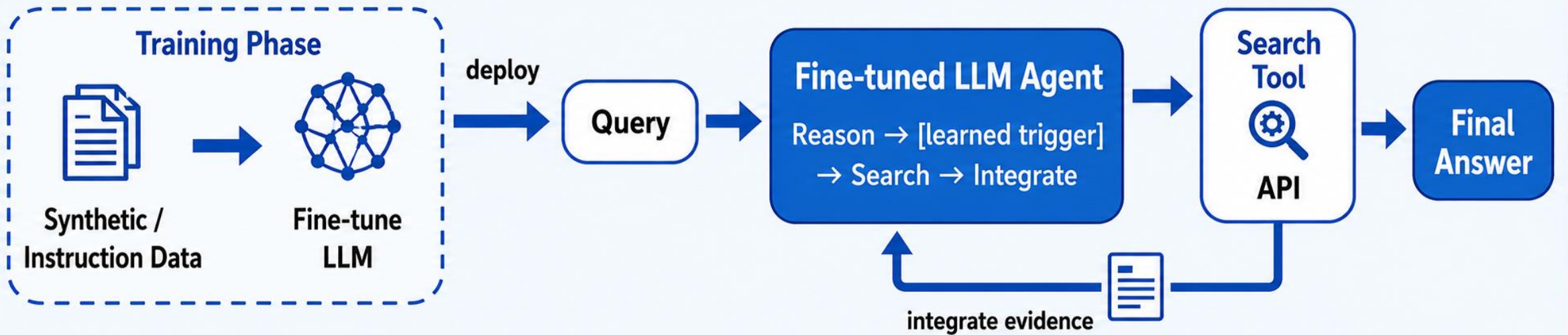
Prompt Chaining Agentic RAG



Single-Agent Agentic RAG



Single-Agent SFT Agentic RAG



Trained on Search-Reasoning Interleaved Data



Stable & Precise Patterns



May Overfit to Tool Schemas

Toolformer, INTERS, RA-DIT

Single-Agent Agentic RAG

Retrieval-Augmented Generation (RAG)

Prompt How did US states get their names?

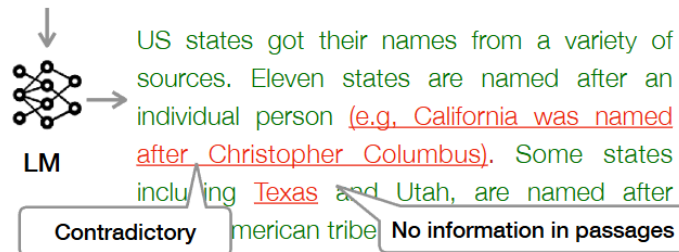
Step 1: Retrieve K documents

- 1 Of the fifty states, eleven are named after an individual person.
- 2 Popular names by states. In Texas, Emma is a popular baby name.
- 3 California was named after a fictional island in a Spanish book.

Retriever

Step 2: Prompt LM with K docs and generate

Prompt How did US states get their names? + 1 2 3



Prompt: Write an essay of your best summer vacation



Ours: Self-reflective Retrieval-Augmented Generation (Self-RAG)

Prompt How did US states get their names?

Step 1: Retrieve on demand

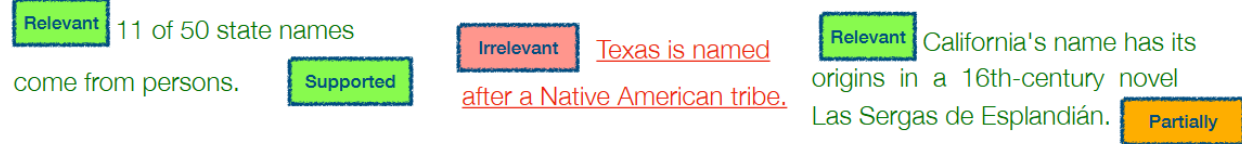


Step 2: Generate segment in parallel

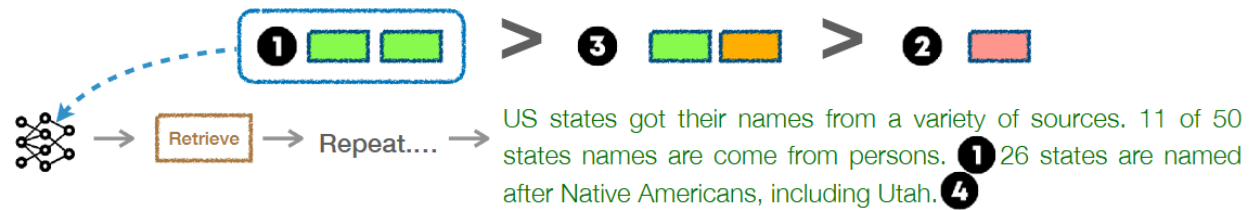
Prompt + 1

Prompt + 2

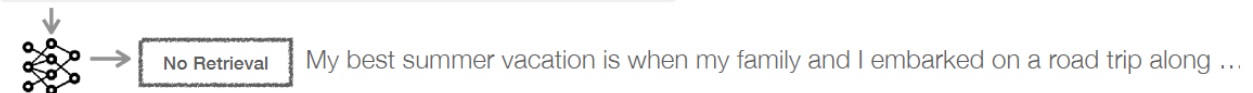
Prompt + 3



Step 3: Critique outputs and select best segment

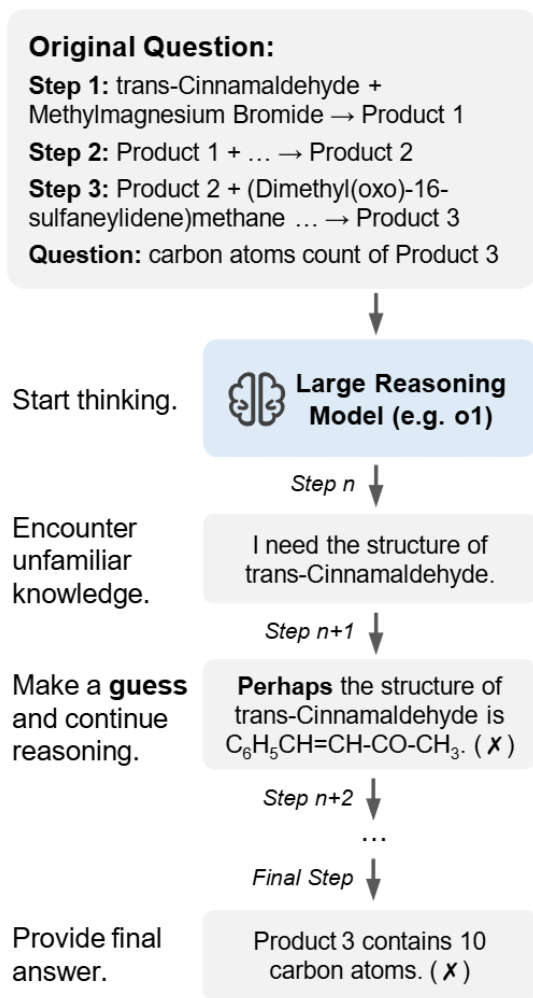


Prompt: Write an essay of your best summer vacation

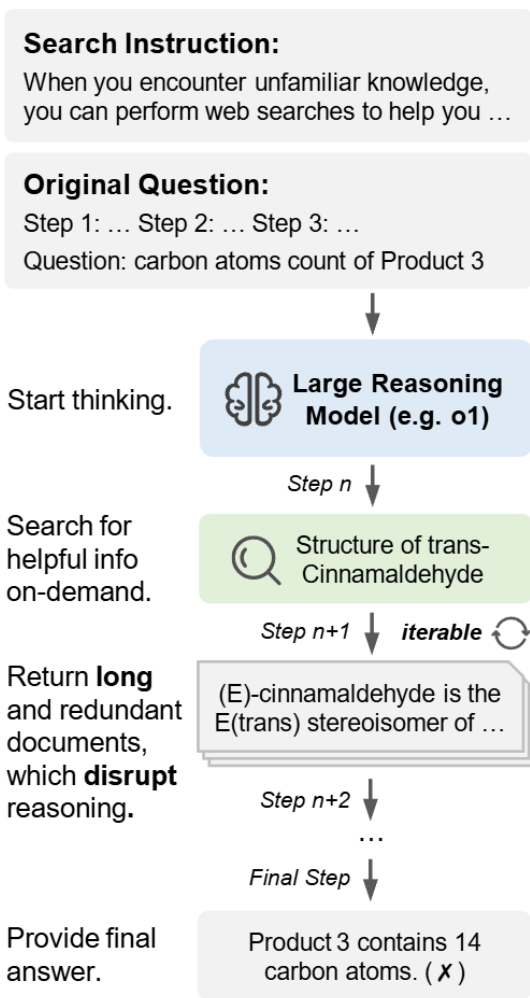


Single-Agent Agentic RAG

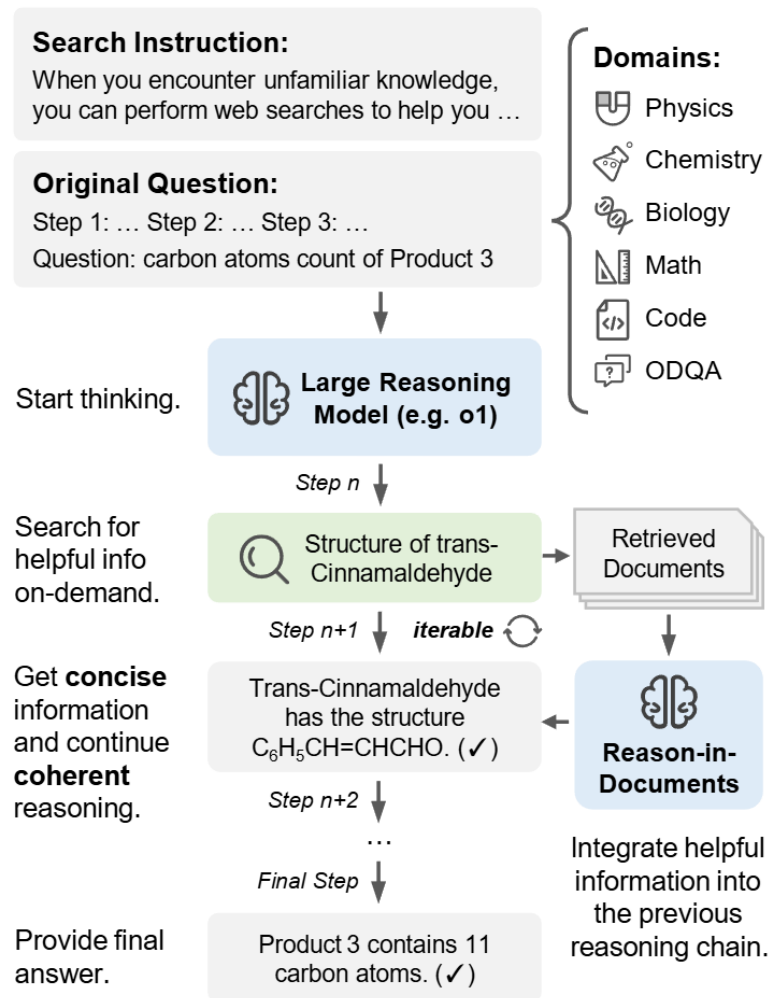
(a) Vanilla Reasoning Pattern



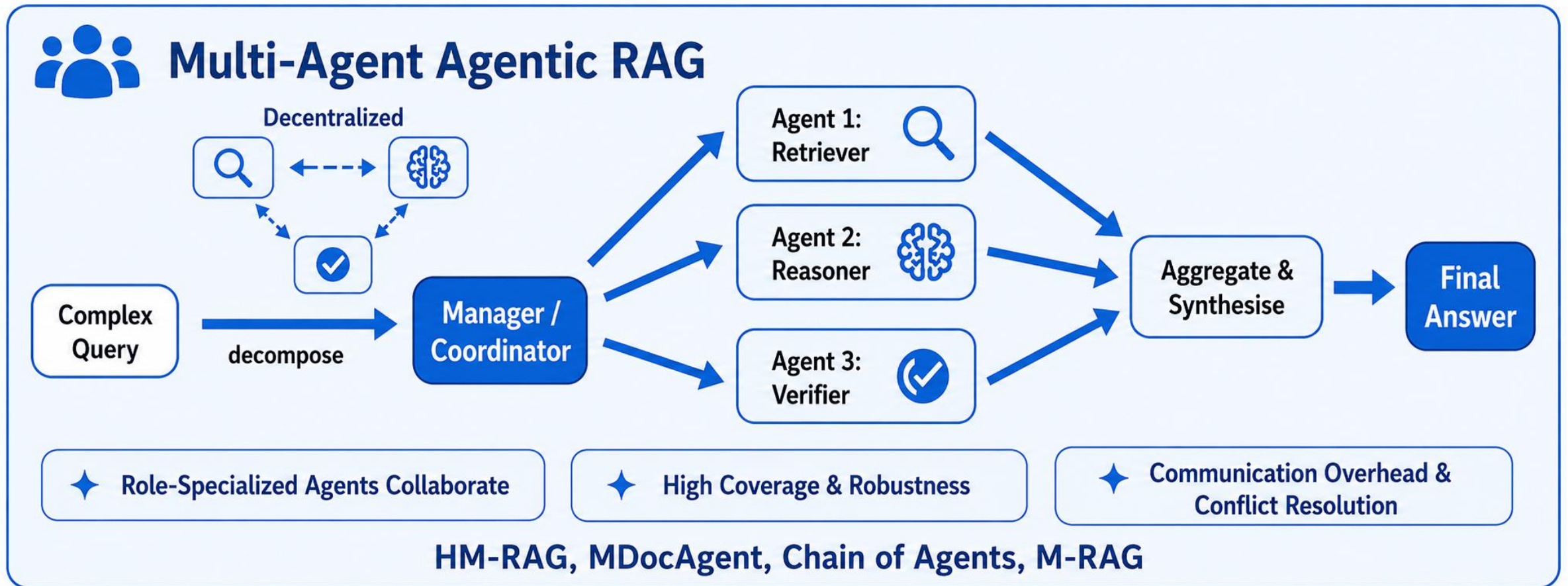
(b) Reason with Agentic RAG (Ours)



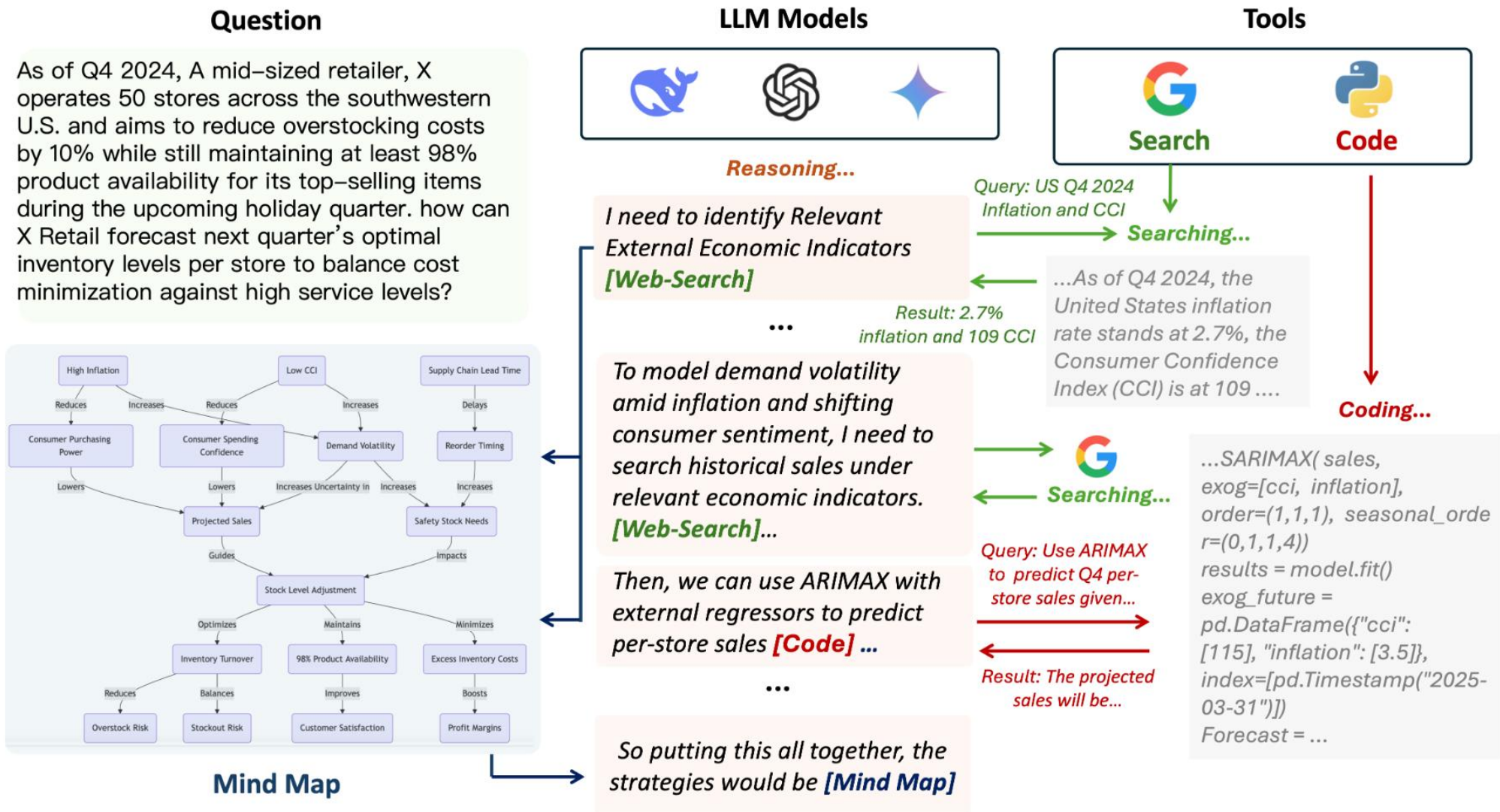
(c) The Search-o1 Framework (Ours)



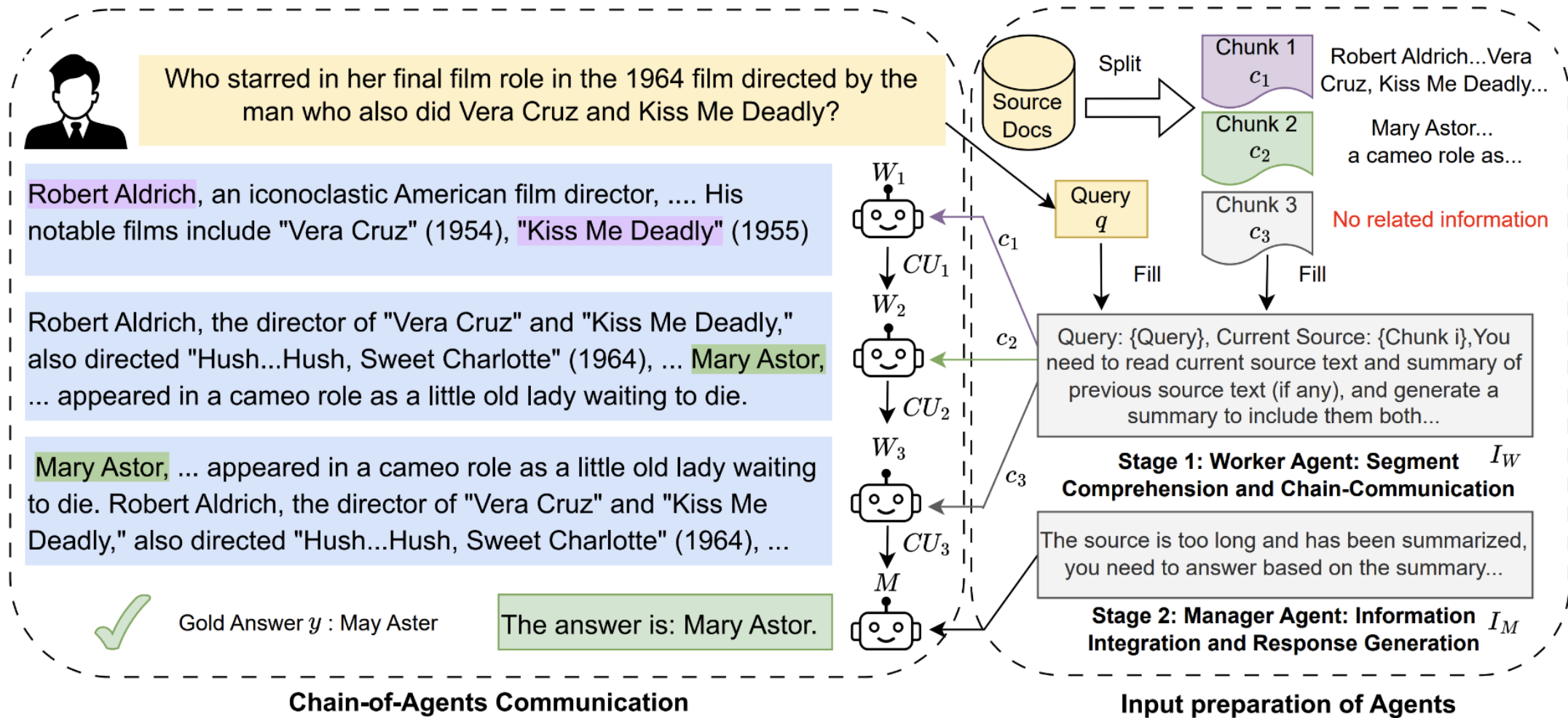
Multi-Agent Agentic RAG



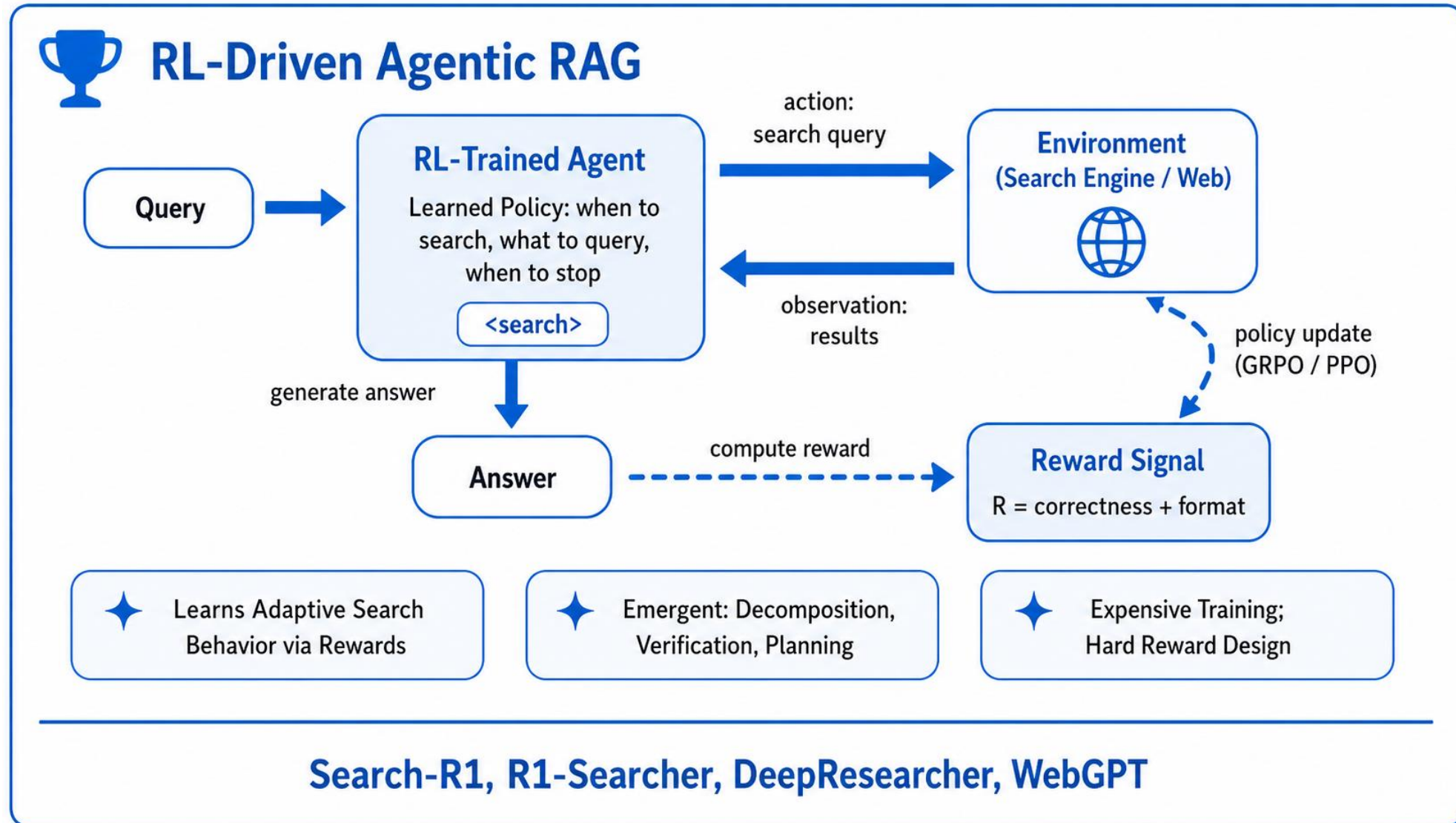
Multi-Agent Agentic RAG



Multi-Agent Agentic RAG



RL-Driven Agentic RAG



RL-Driven Agentic RAG

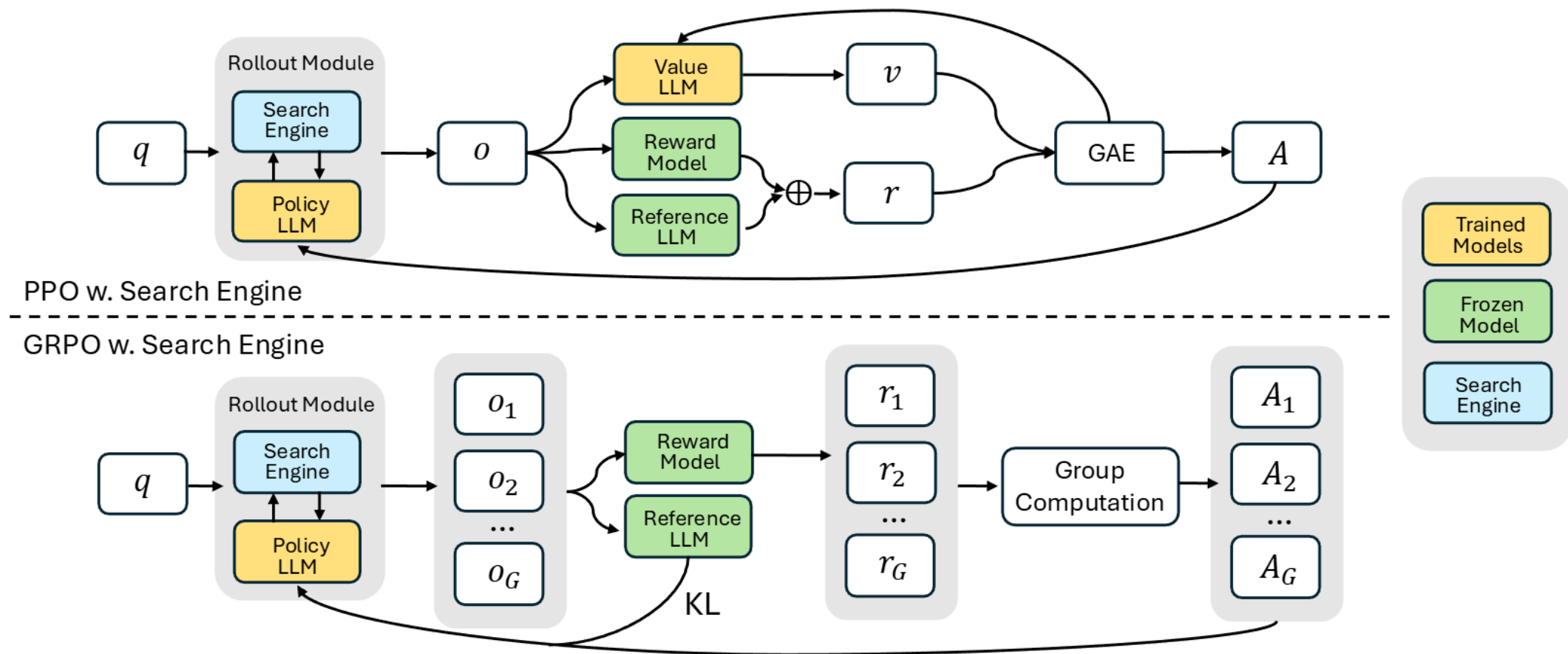
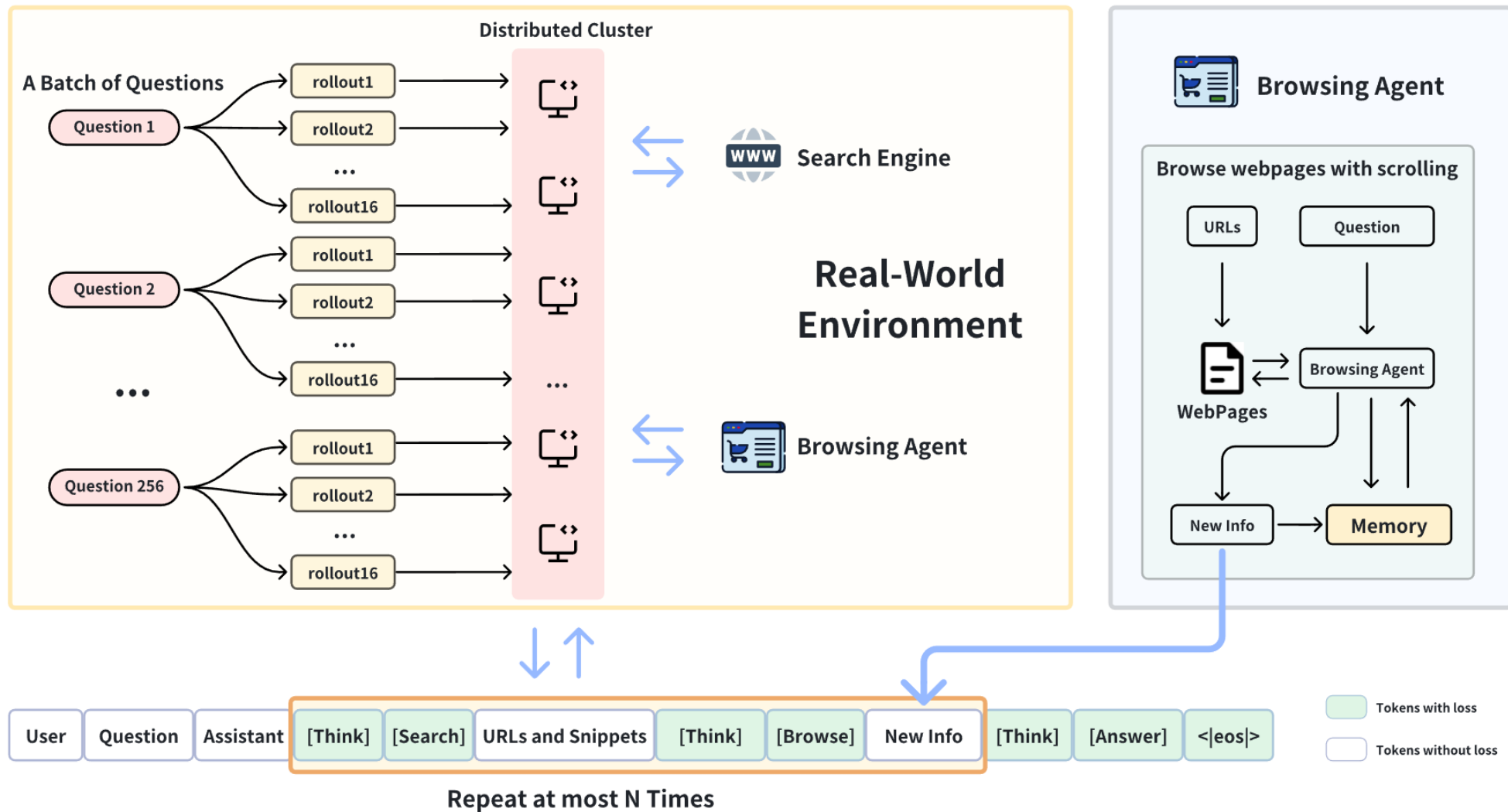


Figure 1: Demonstration of PPO and GRPO training with the search engine (SEARCH-R1). During the rollout, LLMs can conduct multi-turn interactions with the search engine.

RL-Driven Agentic RAG



Tutorial Outline



- **Part 1: Introduction** of Retrieval Augmented Large Foundation Models (RA-LFMs) (Dr. Wenqi Fan)
- **Part 2: Architecture** of RA-LFMs and Main Modules (Xu Yuan)
- **Part 3: Learning Approach** of RA-LFMs (Chengliang Liu)
- **Part 4: Agentic RAG** (Chengliang Liu)
- ⊙ **Part 5: Applications of RA-LFMs (Chun-Hin Chan)**
- **Part 6: Challenges and Future Directions** of RA-LFMs (Dr. Wenqi Fan)
- **Part 7: Q&A**

Website of this tutorial
Check out the slides and more information!



PART 4: Application of RA-LFMs

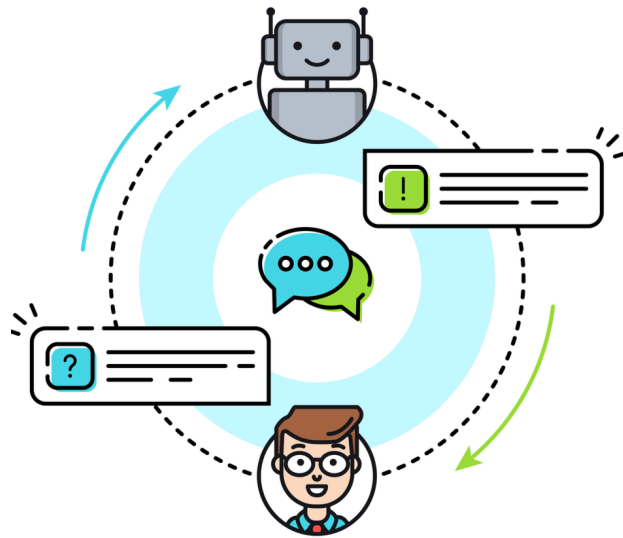


Presenter
Chun-Hin Chan
HK PolyU

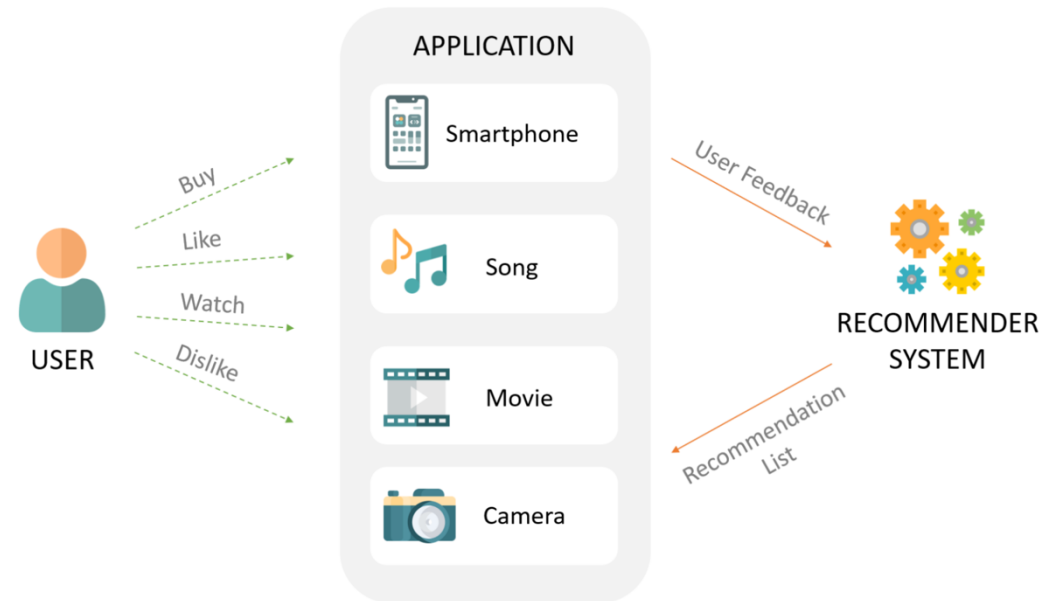
- **Preliminary applications in NLP**
- **Downstream tasks**
- **Domain-specific applications**

RA-LFM Applications

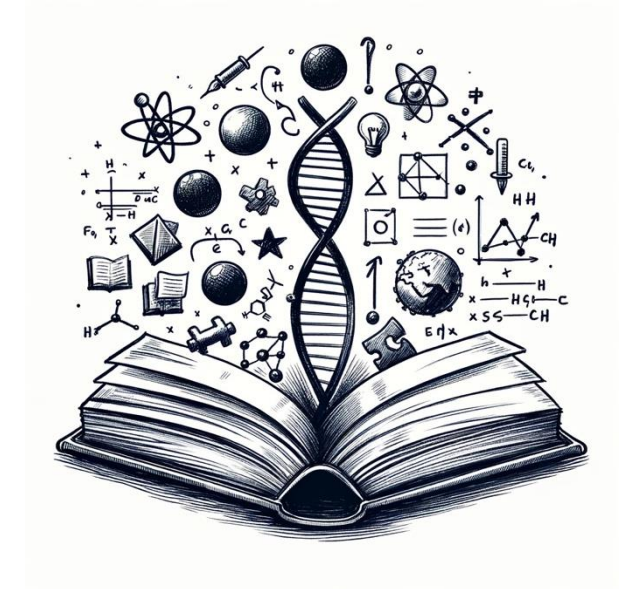
- **Various applications**



Chatbots



Recommendation

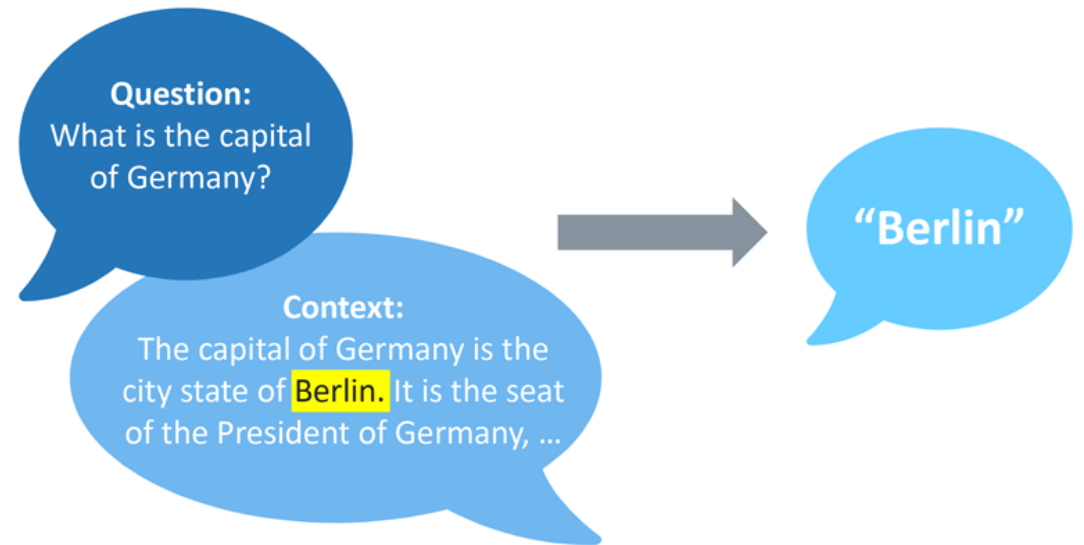
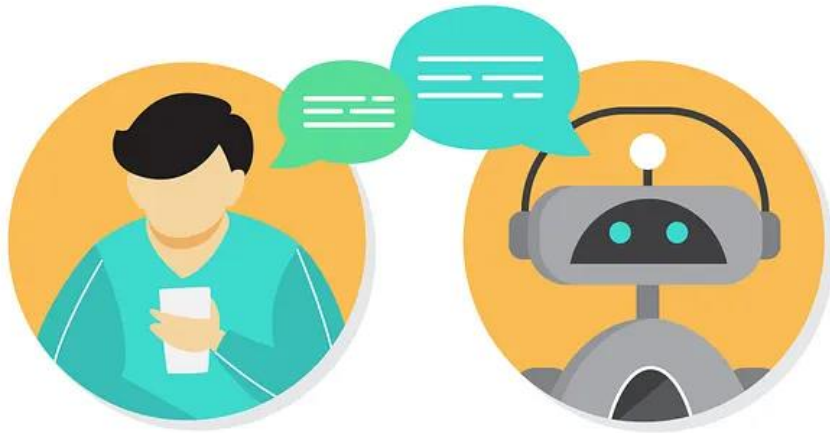


AI for Science

RA-LFM Applications: QA Systems

- **QA Systems**

- QA systems aim to provide **precise answers** to user's queries



RA-LFM Applications: QA Systems

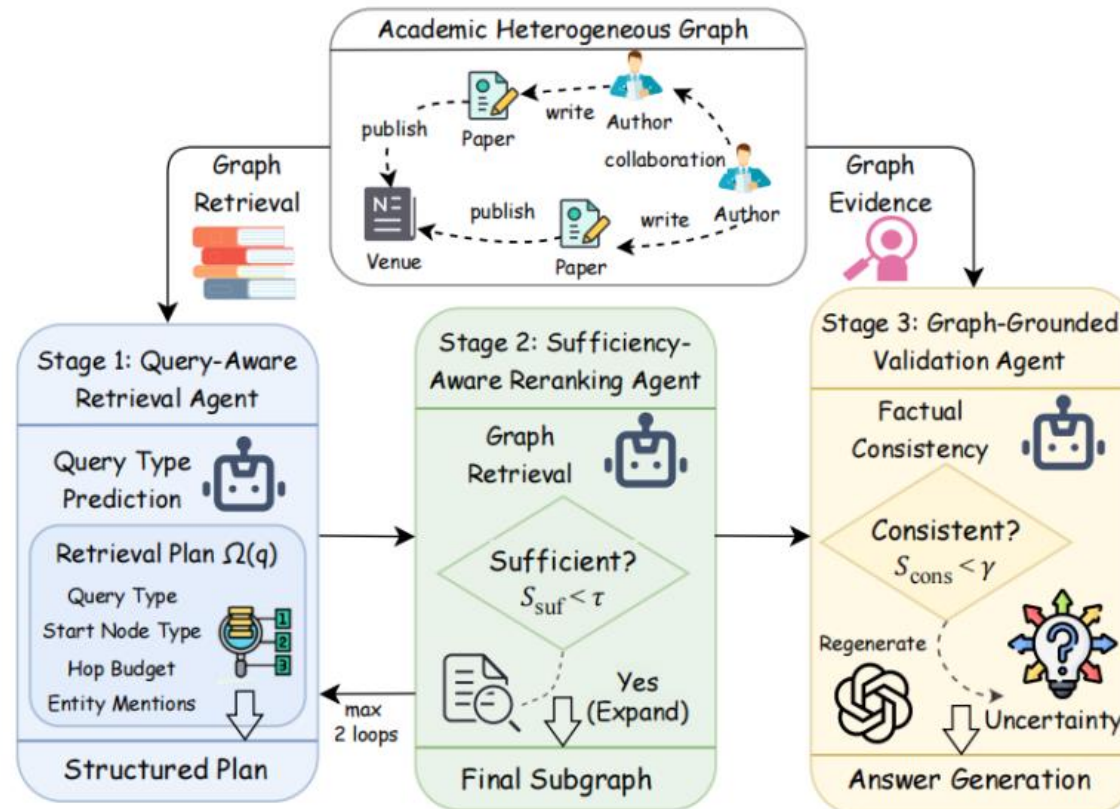
- **QA Systems**

- Challenges:

- Fail to adapt to **varying query complexity**
- **Incomplete or misaligned evidence** due to lack of sufficiency evaluation
- **No structured verification** against graph facts

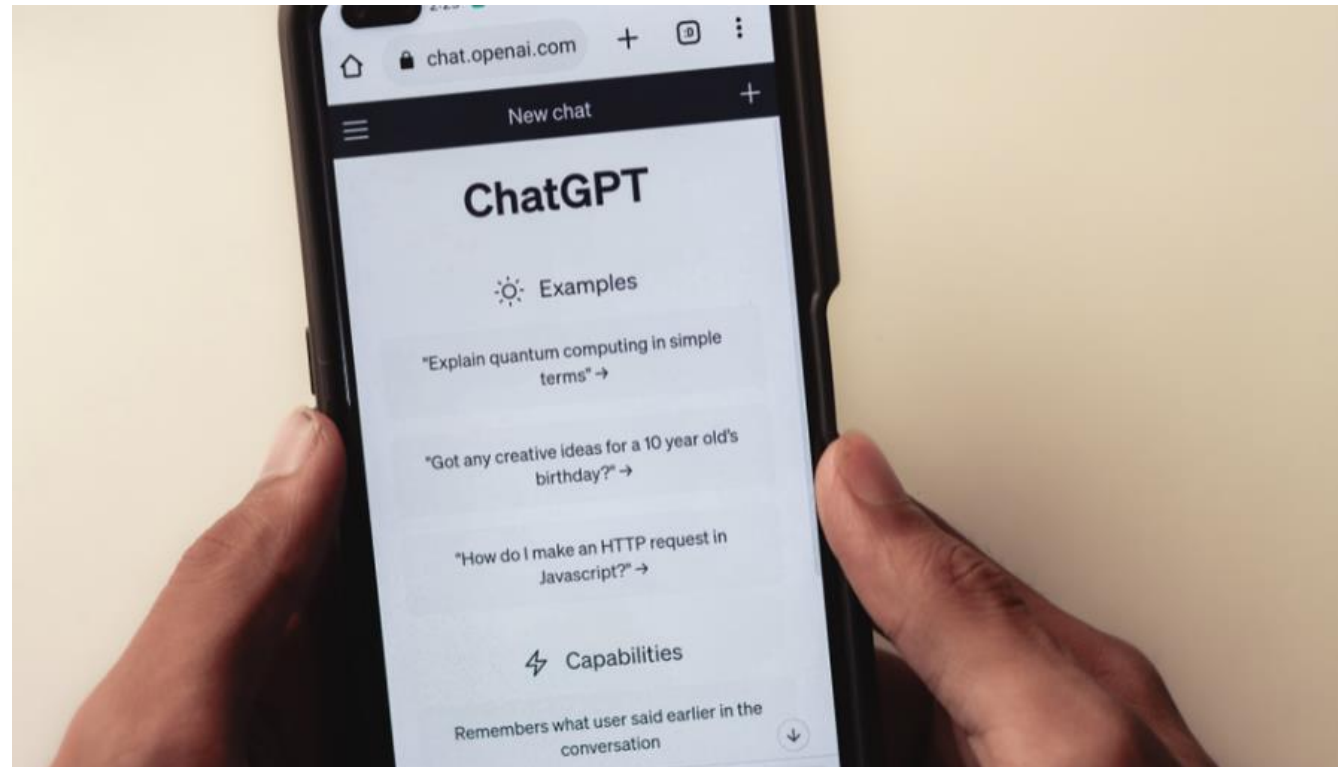
RA-LFM Applications: QA Systems

- **Retrieves from academic heterogeneous graph for Academic QA**
 - Agent-Enhanced Heterogeneous Graph RAG for Academic Question Answering



RA-LFM Applications: Chatbots

- **Chatbots**
 - Chatbot interacts with users in a **natural & conversational** manner

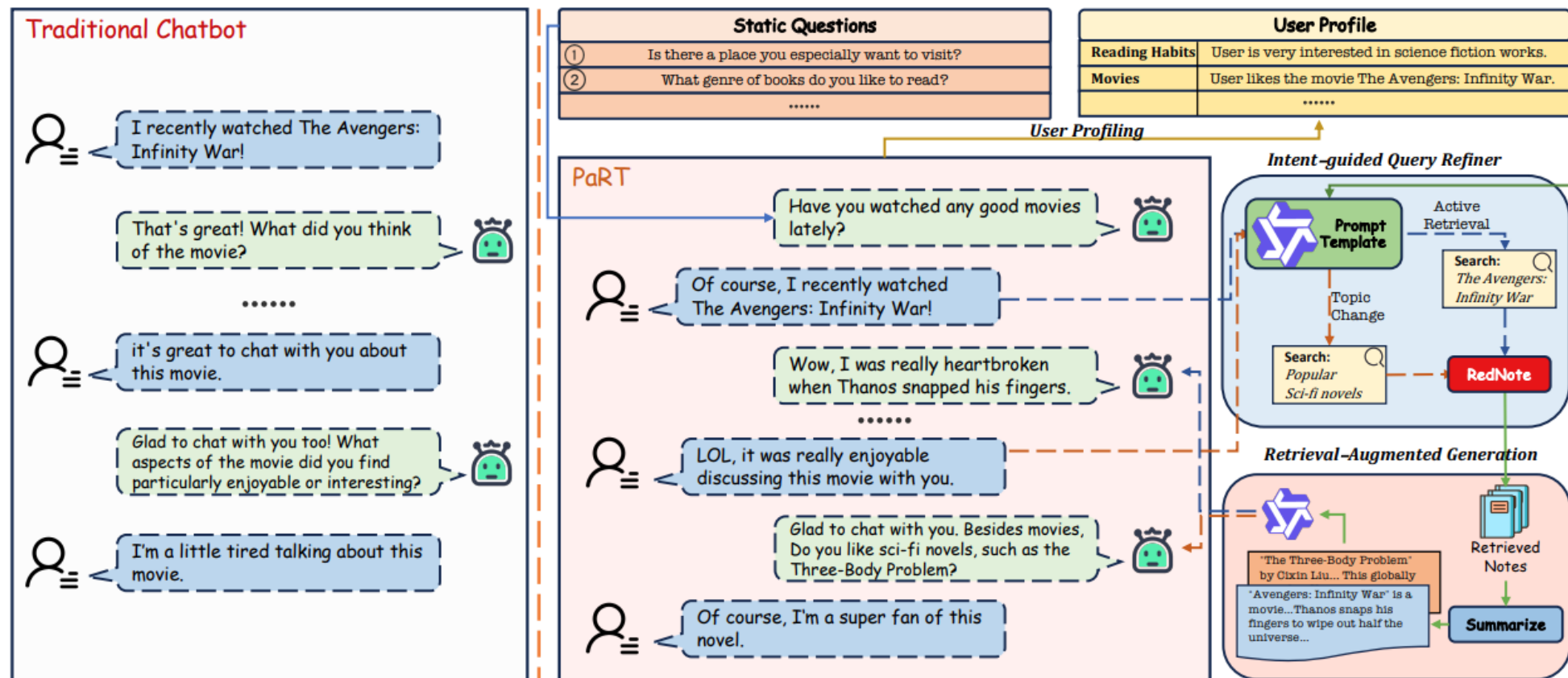


RA-LFM Applications: Chatbots

- **Chatbots**
 - Challenges:
 - **Neglect active engagement** with users
 - Constrain to **limited depth & natural extension of conversation**
 - **Short interaction** duration

RA-LFM Applications: Chatbots

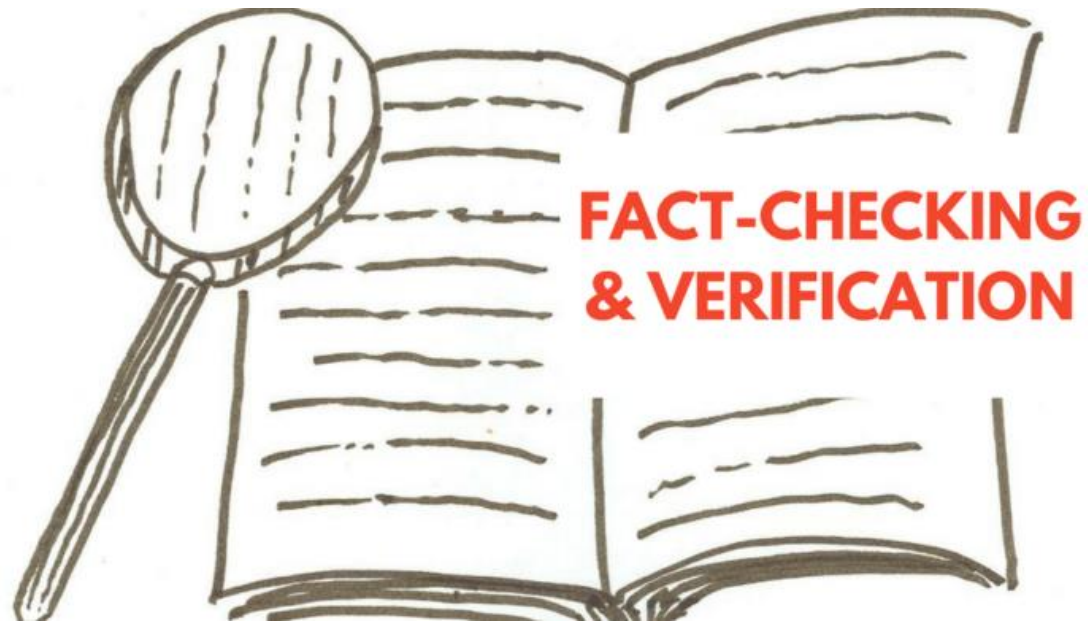
- **Retrieval from textual passages for Chatbots**
 - Proactive social chatbots with personalized real-time **ReTrieval** (PaRT)



RA-LFM Applications: Fact Verification

- **Fact verification**

- Fact Verification is a critical task in verifying the **accuracy** & **reliability** of information

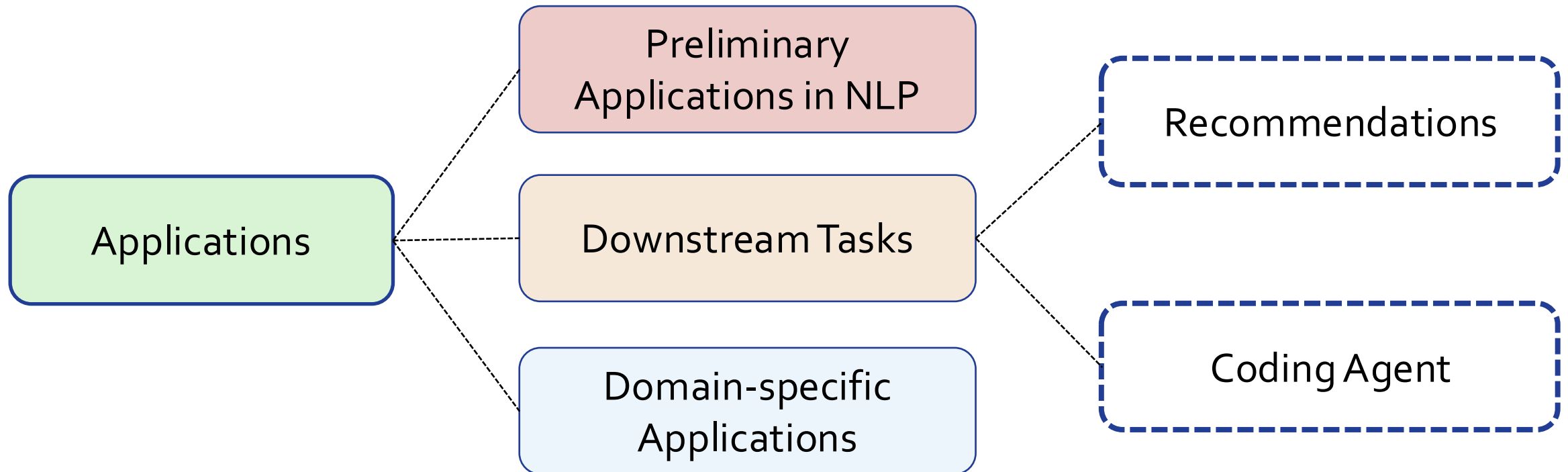


RA-LFM Applications: Fact Verification

- **Fact verification**
 - Challenges:
 - LFM has **limited access** and **manipulation** on knowledge
 - RAG does not guarantee **factual content**
 - Generate **factually incorrect output**

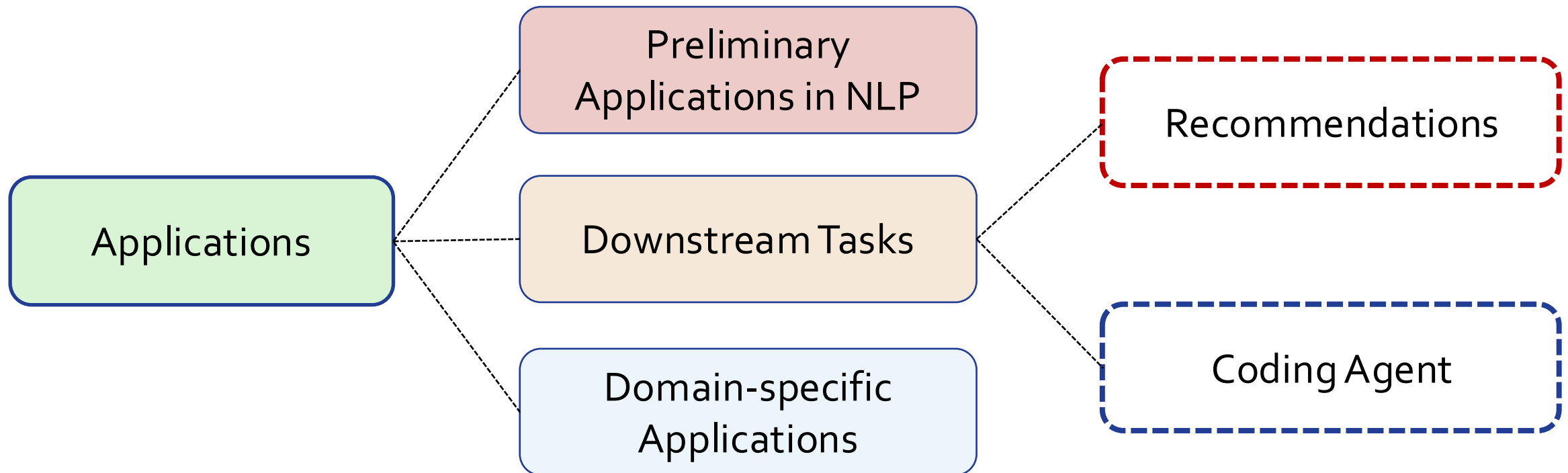
RA-LFM Applications: Downstream Tasks

- **Downstream tasks**



RA-LFM Applications: Recommendations

- **Recommendations**



RA-LFM Applications: Recommendations

- **Recommendations**
 - Recommendation has been widely applied in **online services**

YouTube

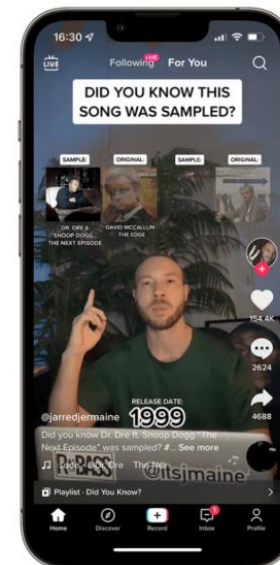
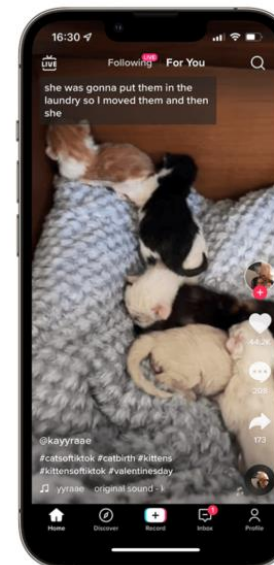
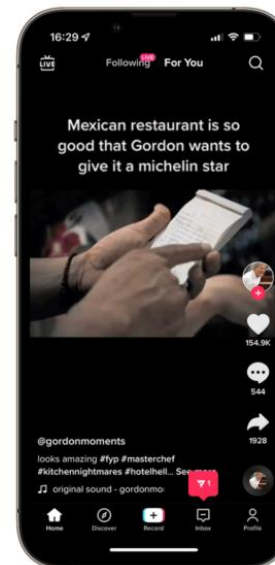
TikTok



News/Video/Image Recommendation

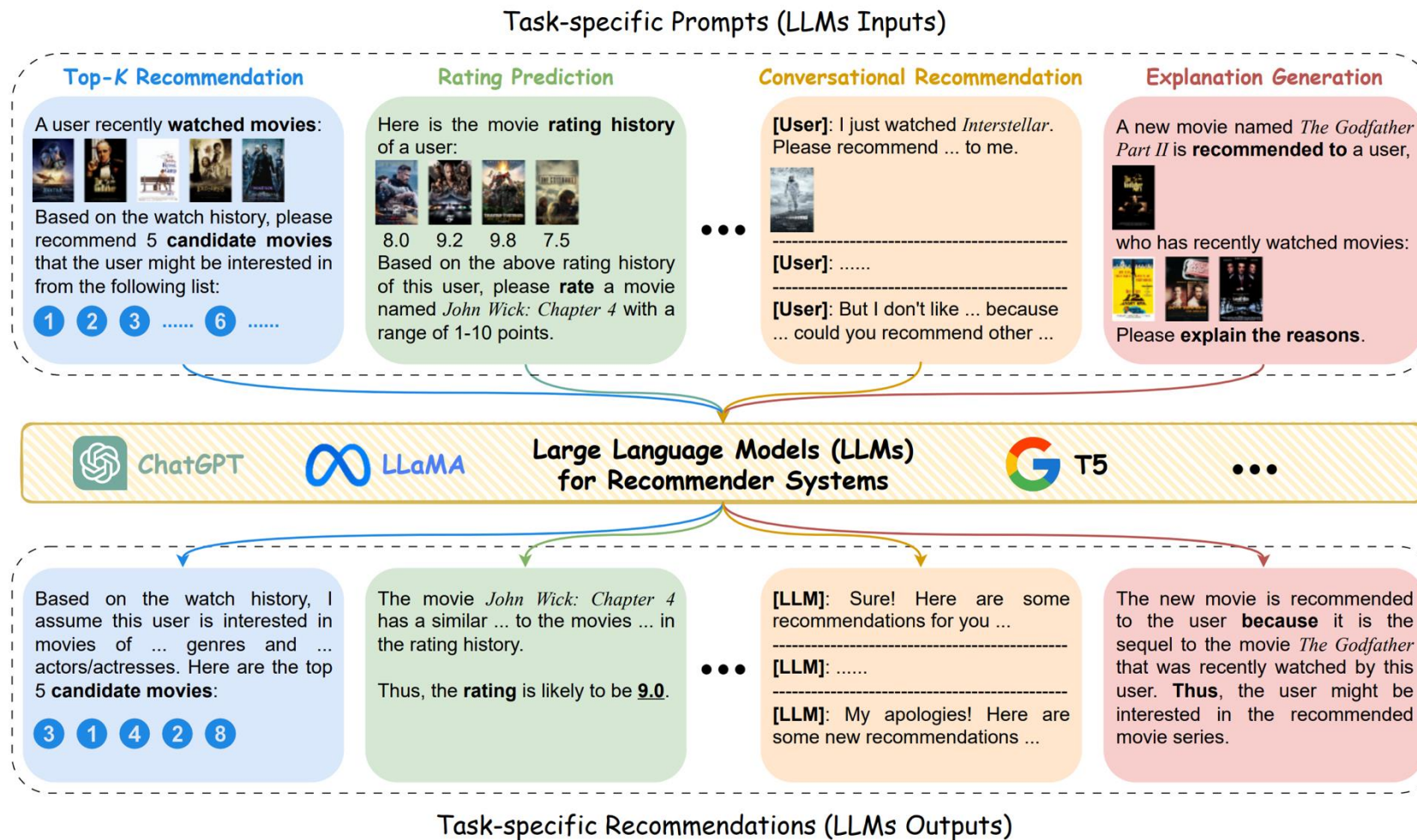
TikTok's recommendation algorithm
Top 10 Global Breakthrough
Technologies in 2021

MIT
Technology
Review



RA-LFM Applications: Recommendations

- LFMs in recommendations



RA-LFM Applications: Recommendations

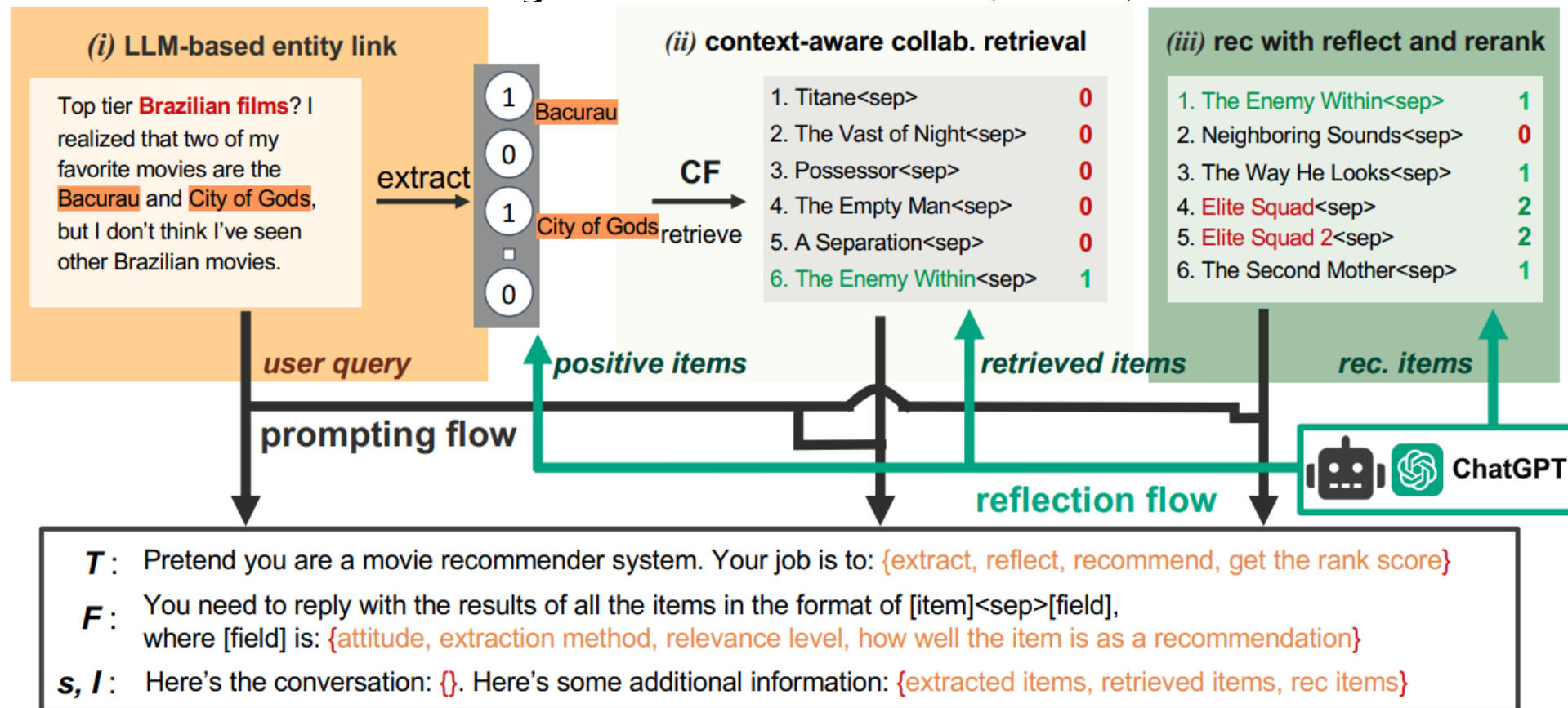
- **Recommendations**

- Challenges:

- Fall short in the effective usage of **collaborative filtering knowledge**
- **Neglect structural relationships** in knowledge
- Matching-based RAG unable to identify **useful information over noisy web data**

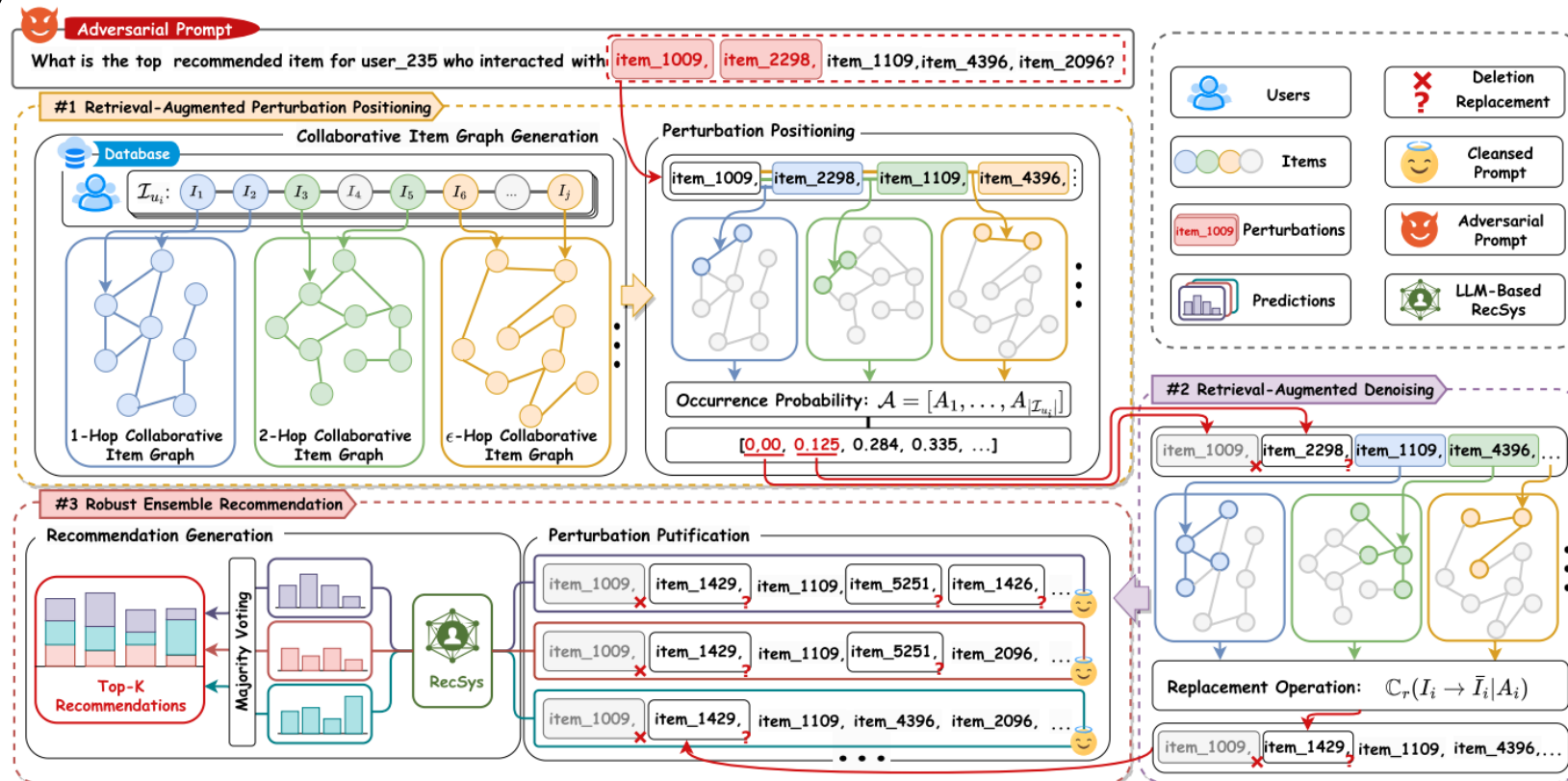
RA-LFM Applications: Recommendations

- Retrieval from collaborative signals
 - Collaborative Retrieval Augmented Generation (CRAG)



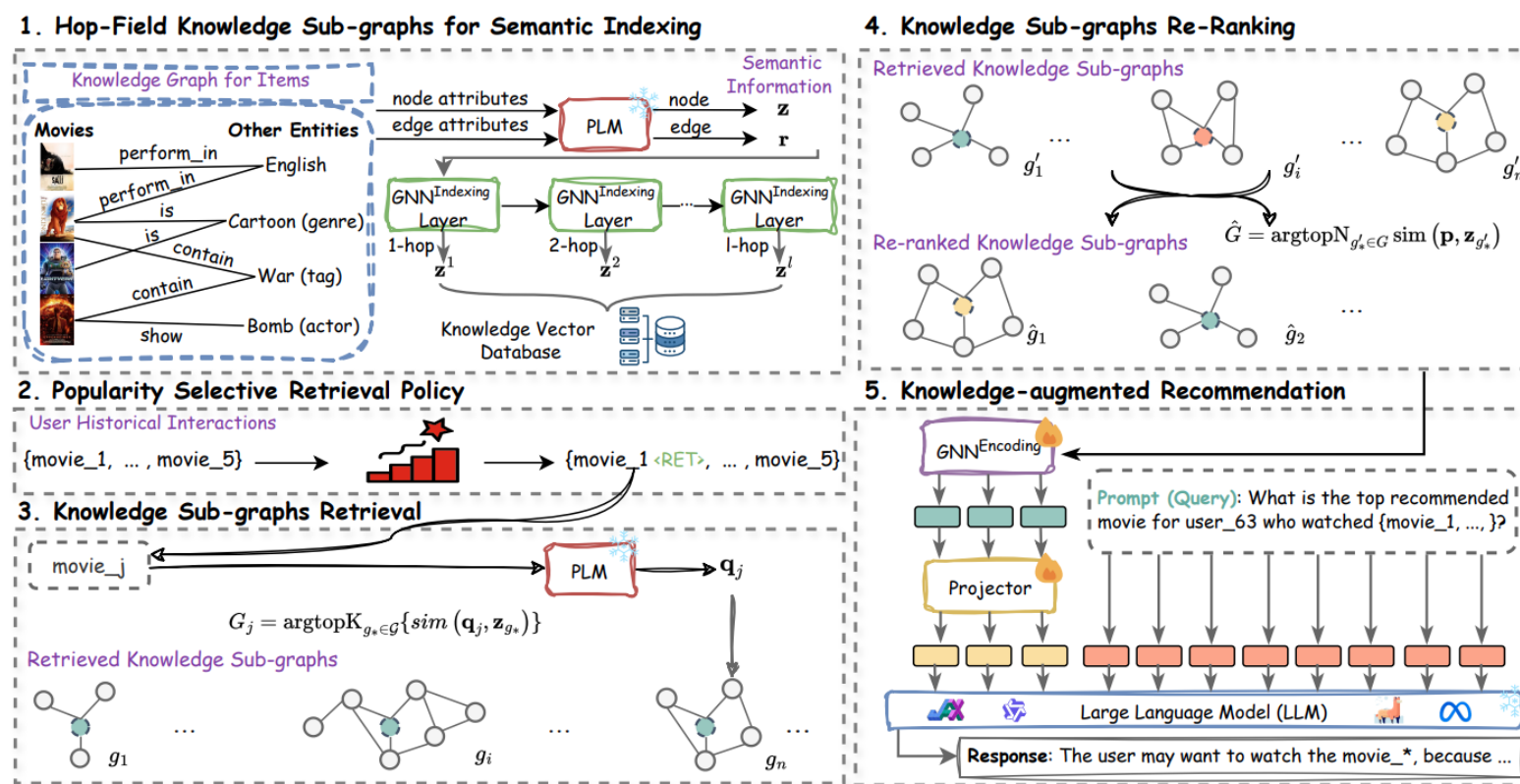
RA-LFM Applications: Recommendations

- **Retrieval from collaborative signals for Purification**
 - **RE**trieval-a**U**gmented pu**R**ifier for e**N**hancing robustness of LLM based RecSys (RETURN)



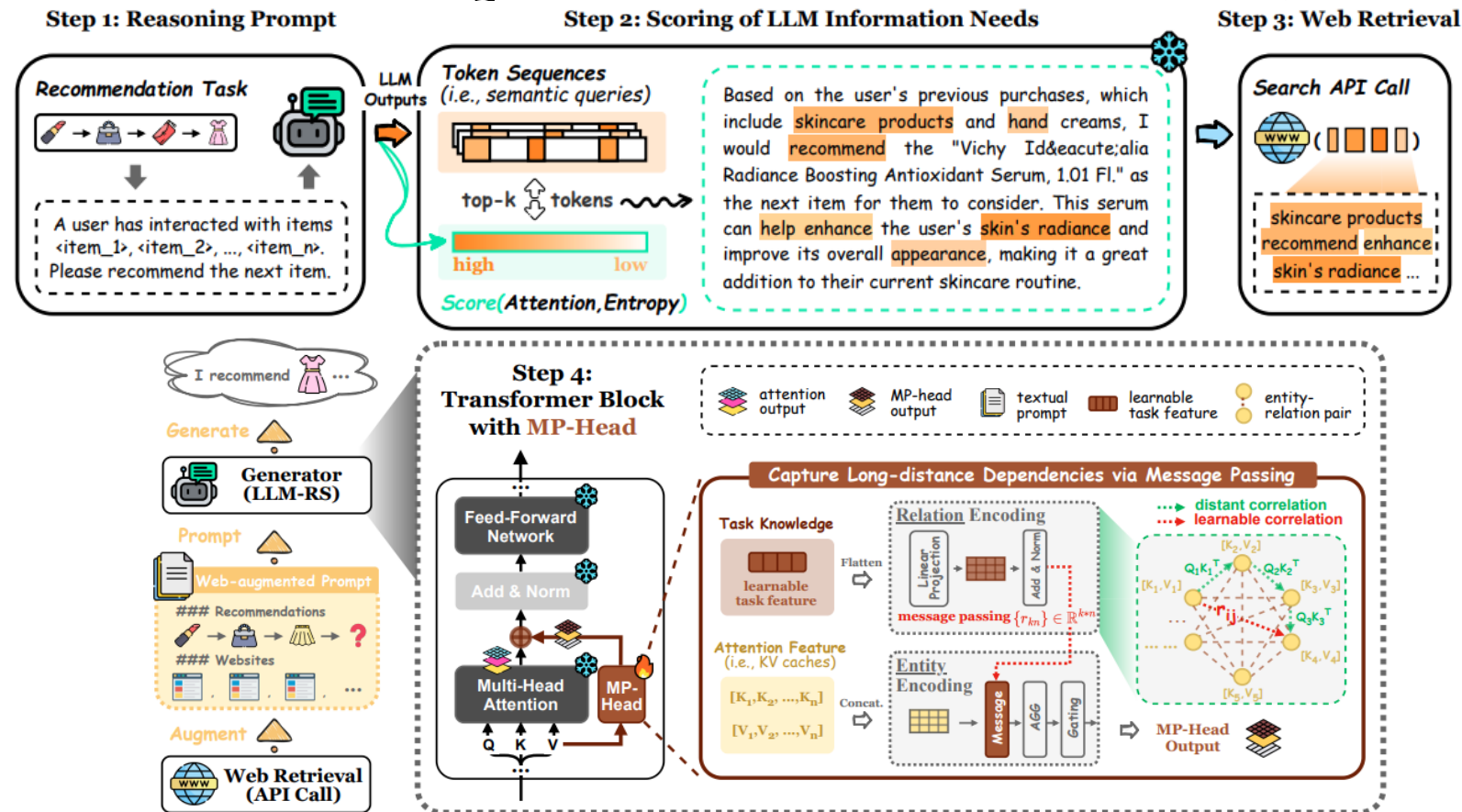
RA-LFM Applications: Recommendations

- **Retrieval from knowledge graph**
 - Knowledge graph **Retrieval-Augmented Generation** for LLM-based **Recommendation** (K-RagRec)



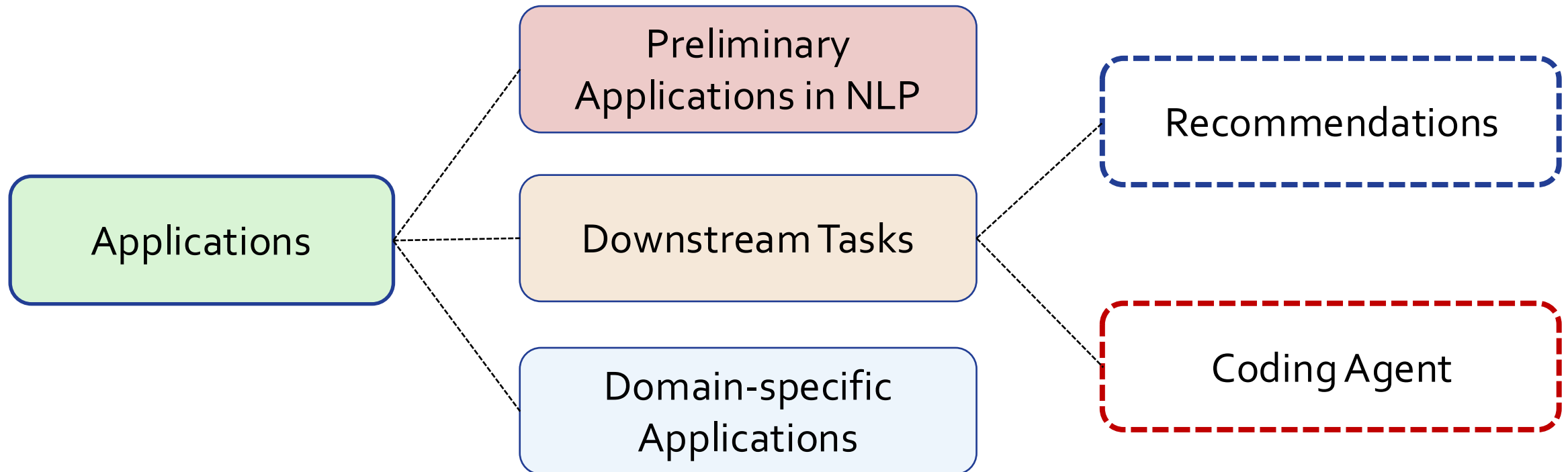
RA-LFM Applications: Recommendations

- Retrieval from the web
 - RAG from **Web** for enhancing LLM-based **Recommendations** (WebRec)



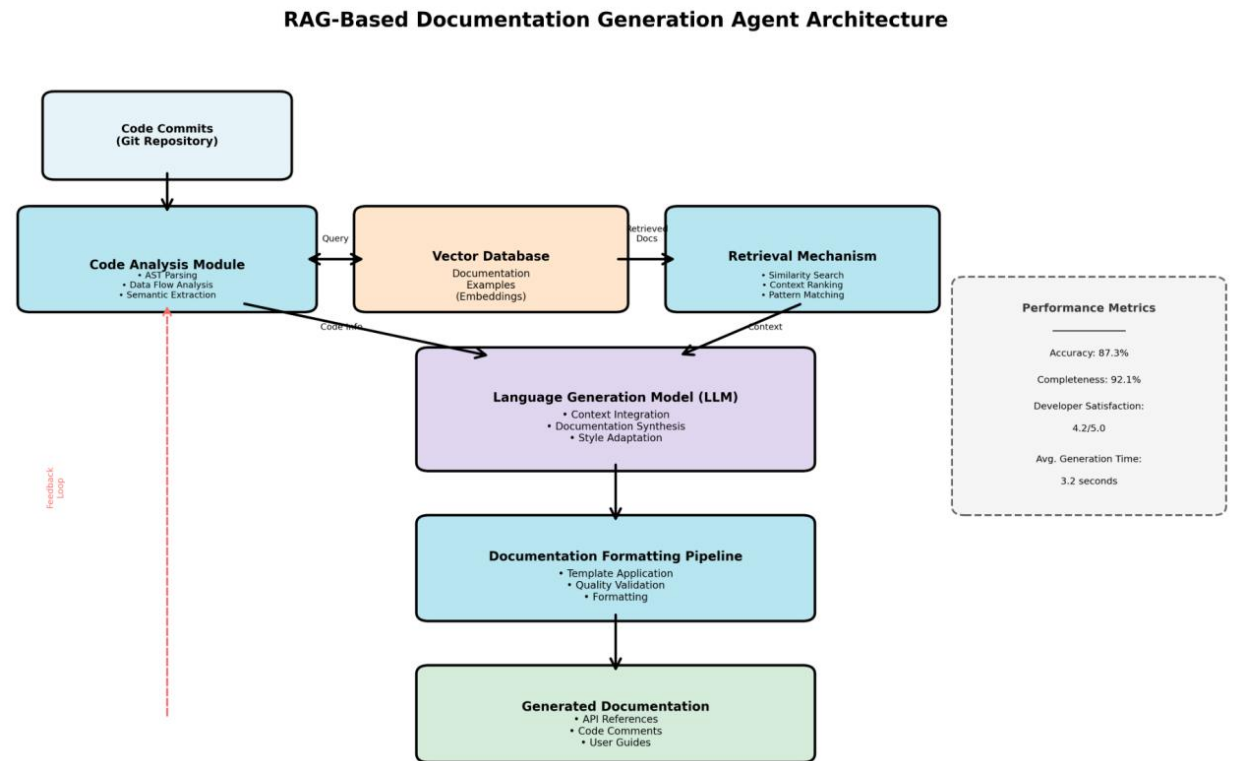
RA-LFM Applications: Coding Agent

- **Coding agent**



RA-LFM Applications: Coding Agent

- **Coding agent:**
 - Code generation
 - Code repair
 - Production deployment
 - ...



RA-LFM Applications: Coding Agent

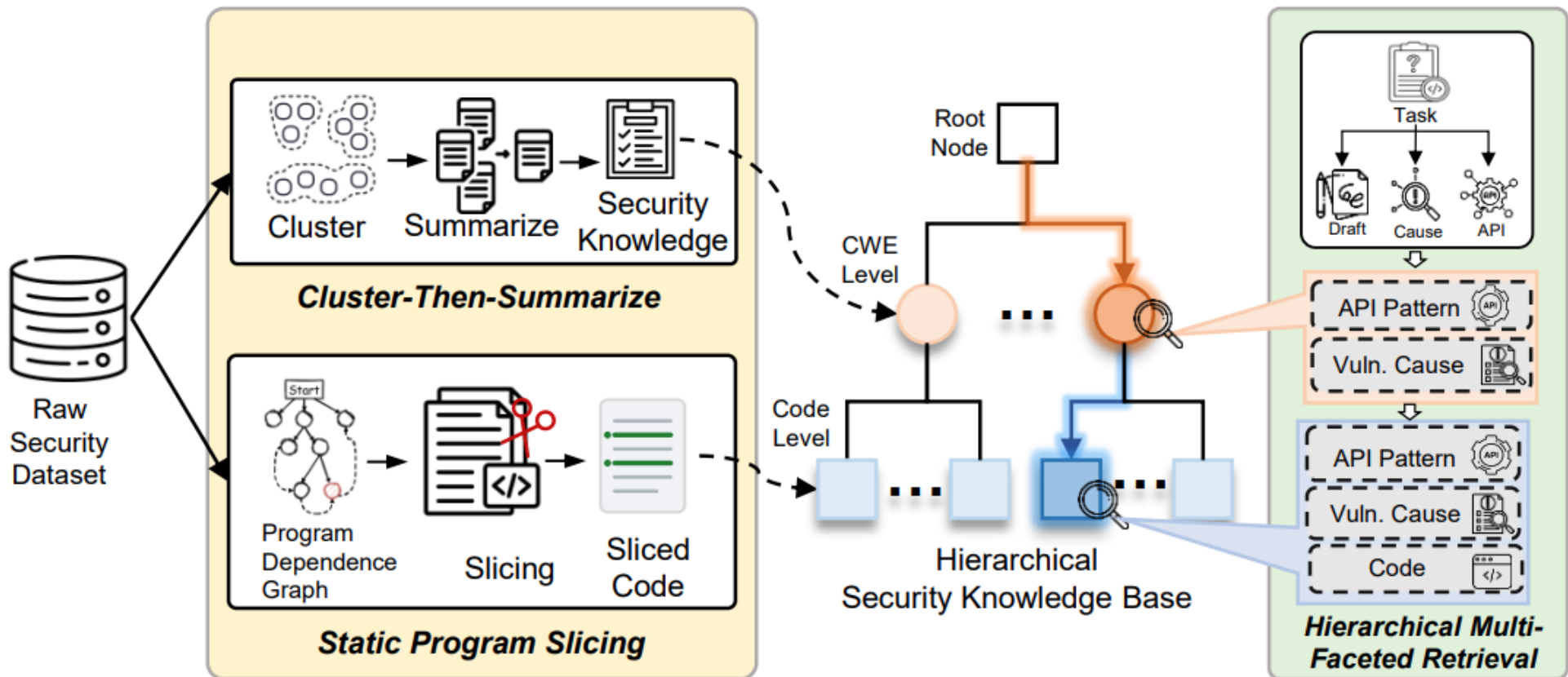
- **Coding Agent**

- Challenges:

- Struggles with the **noise of raw security-related documents**
- Overlook the **key security semantics** implicitly embedded in task descriptions
- Fail to capture the **structural intricacies of code**

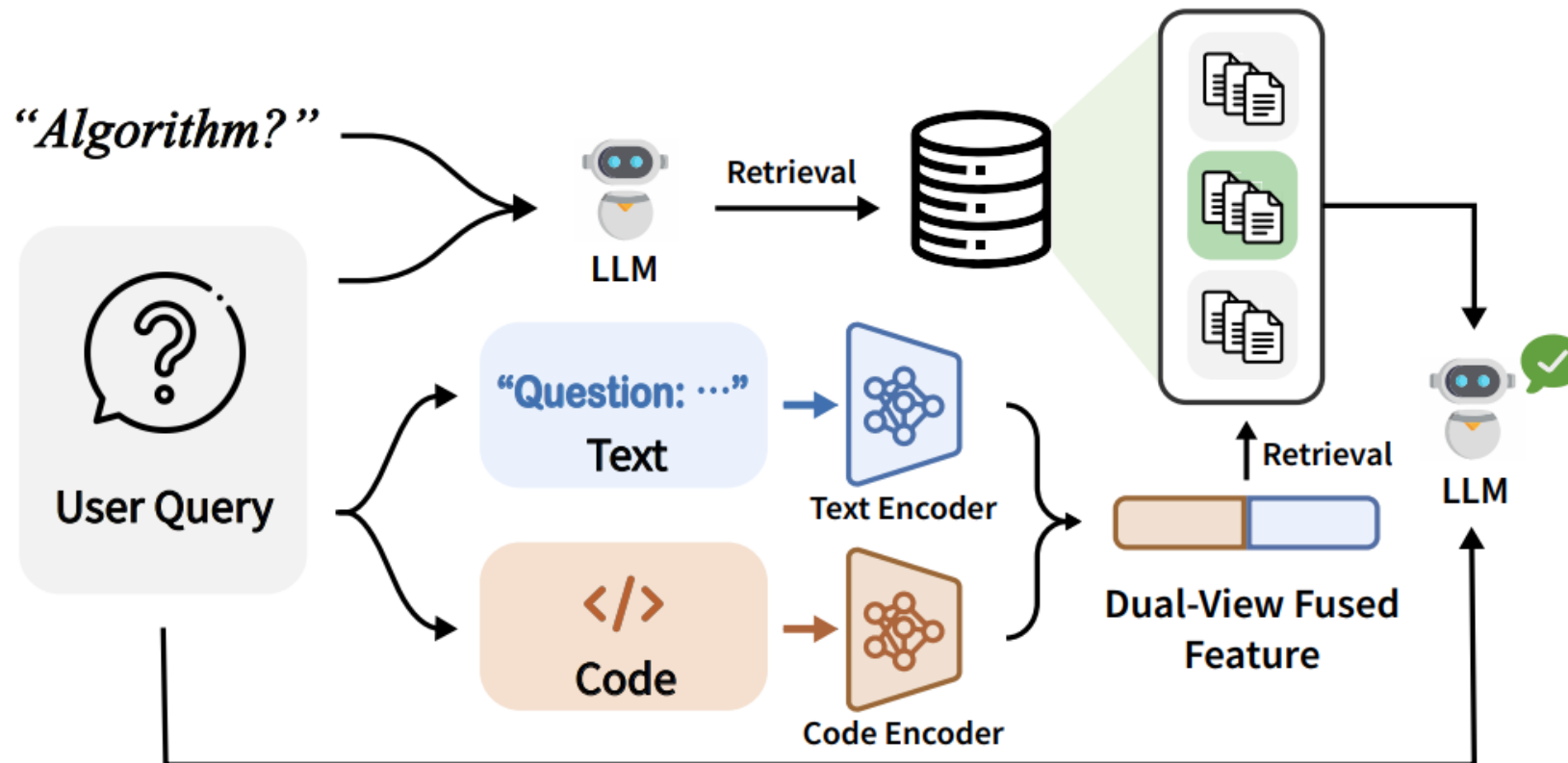
RA-LFM Applications: Coding Agent

- **Retrieval from security knowledge base for Code Generation**
 - **RE**trieval-augmented **S**ecure **C**ode **g**ENERation (RESCUE)



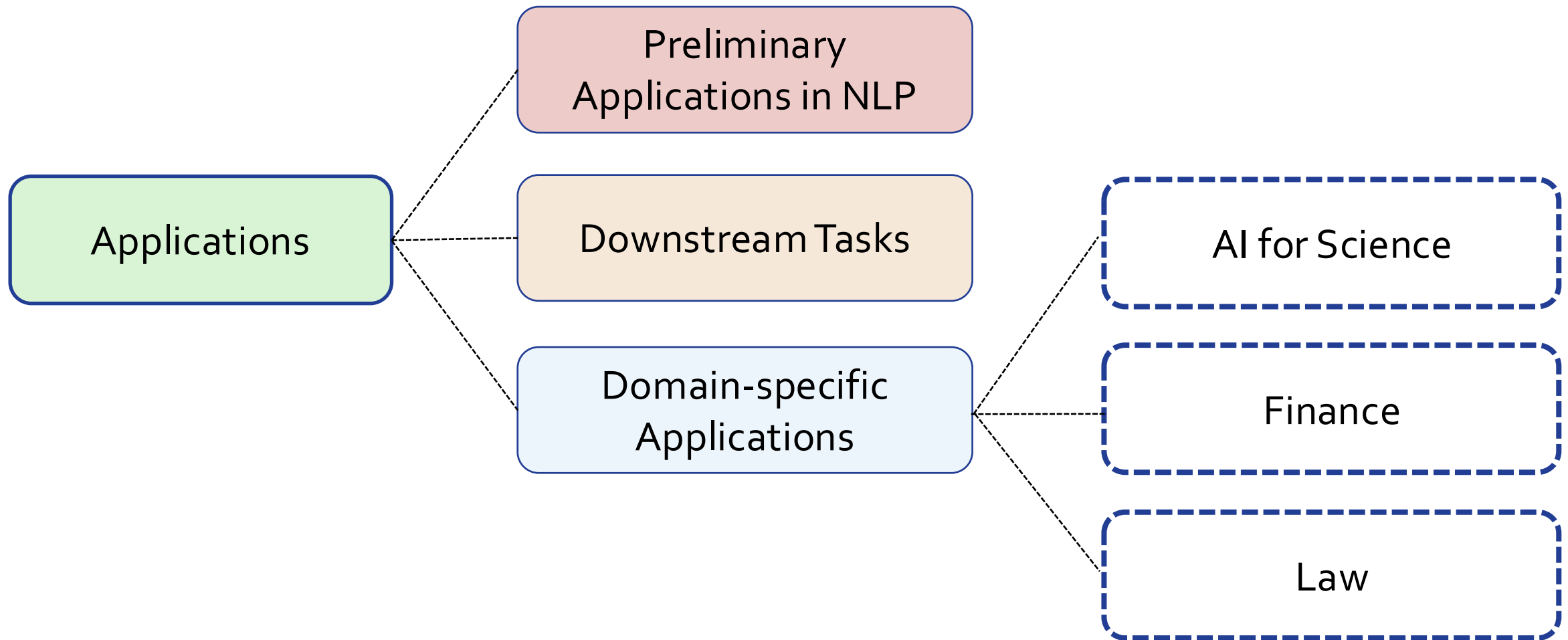
RA-LFM Applications: Coding Agent

- **Retrieval from algorithm-specific knowledge corpus for Code Repair**
 - Improving LLM-based **Code Repair** with Fine-Grained Retrieval-Augmented Generation (ReCode)



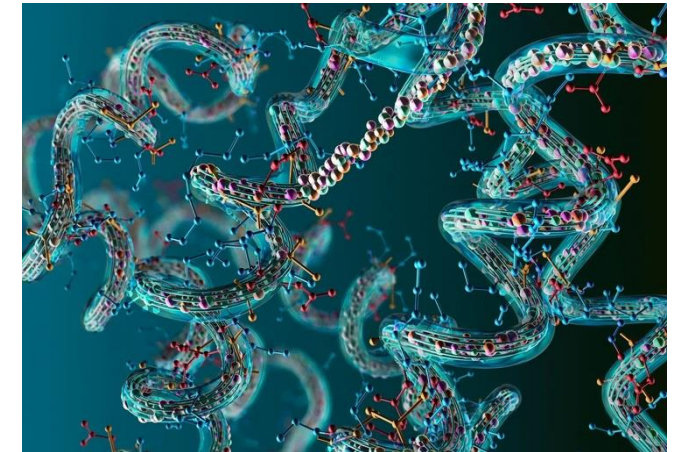
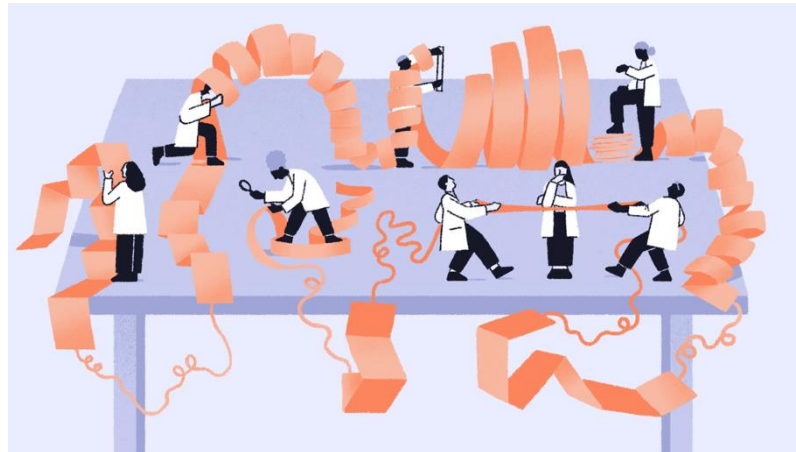
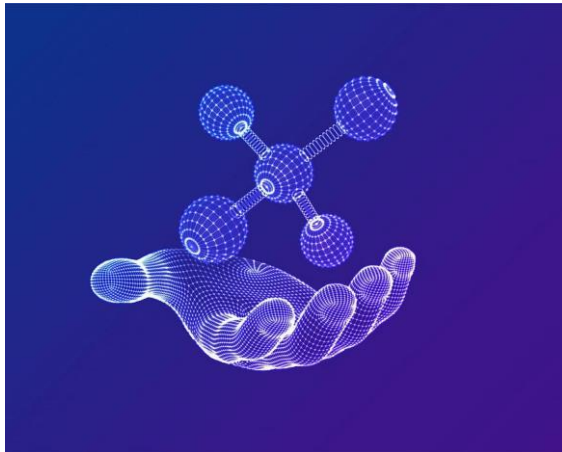
RA-LFM Applications: Domain-specific Applications

- **Domain-specific applications**



RA-LFM Applications: AI for Science

- **AI for science**
 - Molecules
 - Protein
 - ...

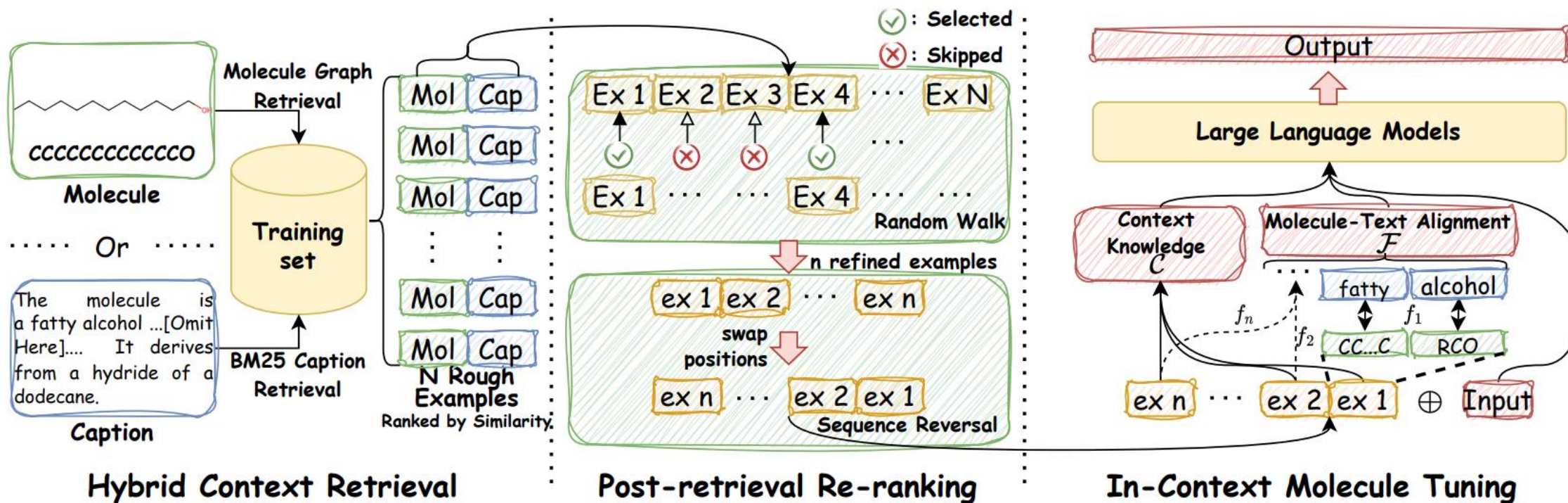


RA-LFM Applications: AI for Science

- **AI for science**
 - Challenges:
 - High-quality chemical datasets are **rare**
 - **Costly to finetune** on chemical-specific tasks
 - Fail to align **protein knowledge graphs** with **biological task specific data**

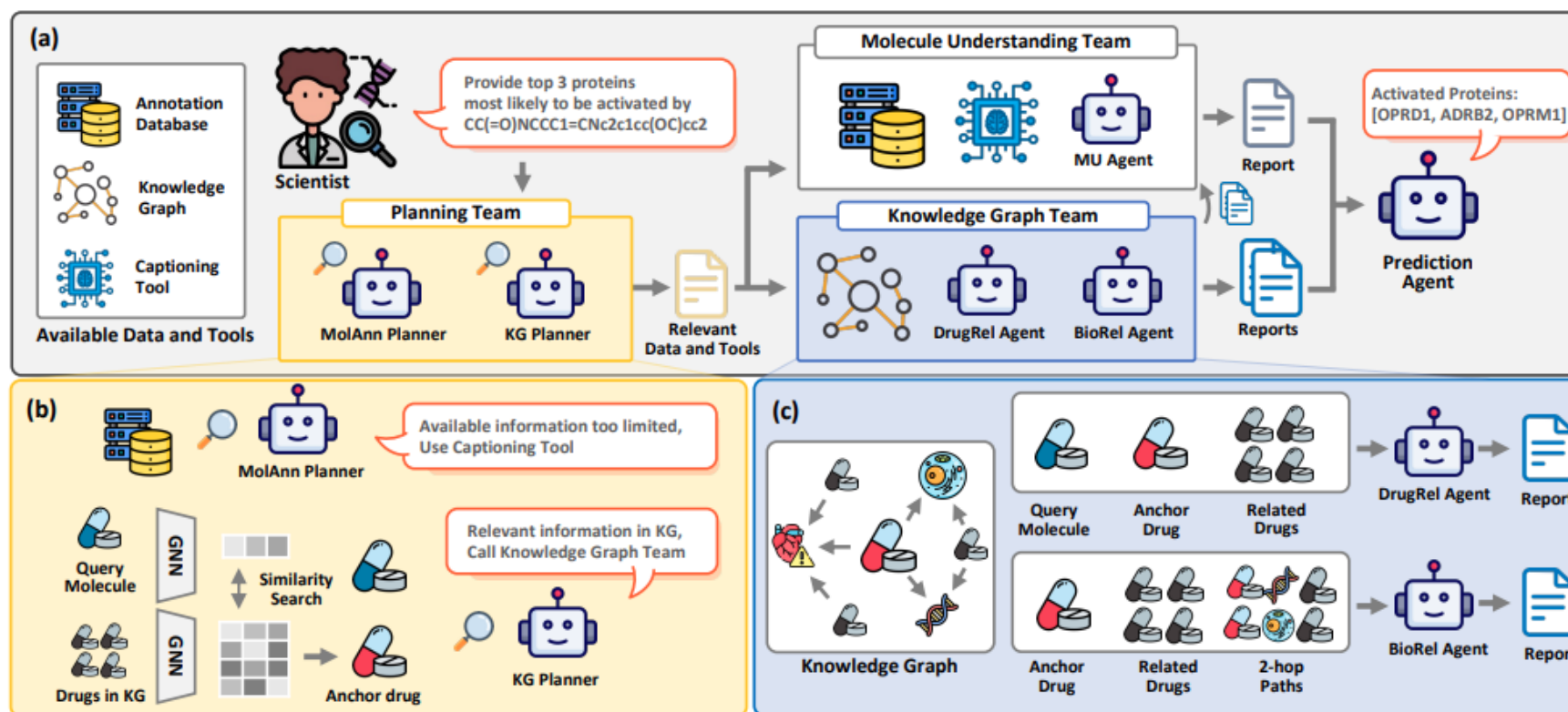
RA-LFM Applications: AI for Science

- Retrieval from the molecule graph & caption for Molecules Learning
 - In-Context Molecule Adaptation (ICMA)



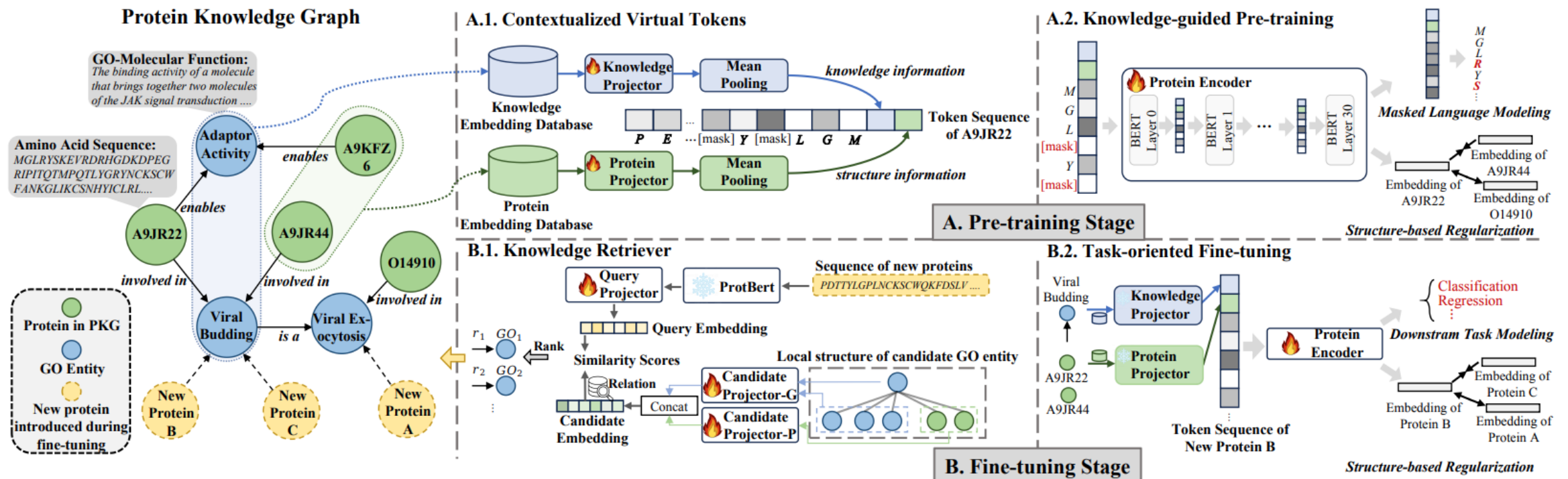
RA-LFM Applications: AI for Science

- Retrieval from knowledge graph & molecule annotation databases for Drug Discovery
 - Collaborative framework of LLM Agents for Drug Discovery (CLADD)



RA-LFM Applications: AI for Science

- Retrieval from protein knowledge graph for Protein Encoding
 - Knowledge-aware retrieval augmented protein language model (Kara)



RA-LFM Applications: Finance

- **Finance**
 - Financial question answering
 - Financial decision-making
 - ...



RA-LFM Applications: Finance

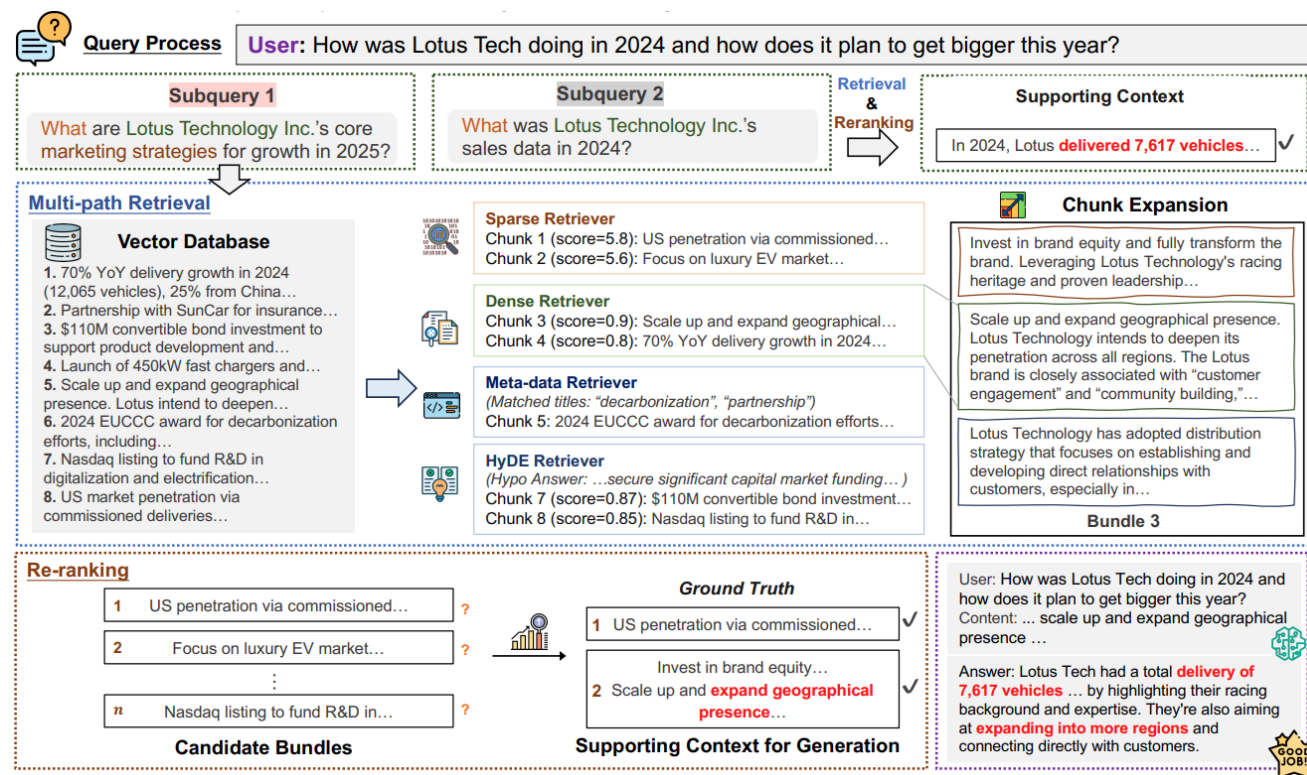
- **Finance**

- Challenges:

- Financial regulations continuously **evolve**
- Existing financial datasets offer **limited data modalities**
- Fail to incorporate **temporal information**

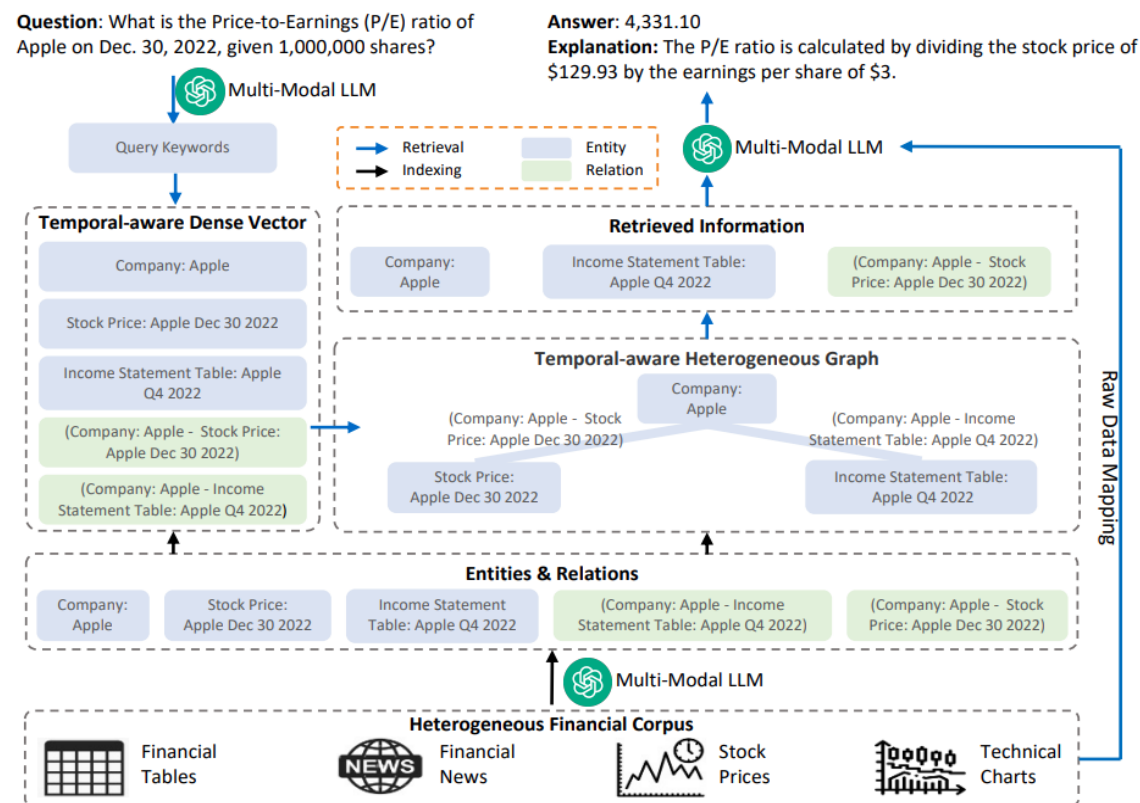
RA-LFM Applications: Finance

- Retrieval from metadata & semantic information for Financial Filings QA
- RAG Financial QA system for regulatory compliance analysis (FinSage)



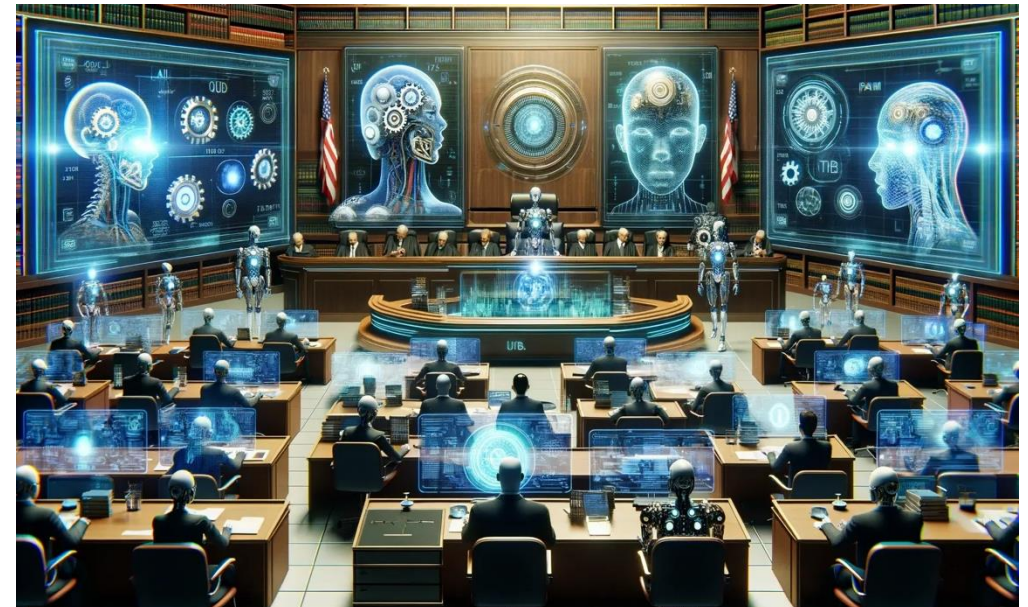
RA-LFM Applications: Finance

- Retrieval from heterogeneous corpus for Financial Decision-Making
 - TeMporal-aware Multimodal Hybrid corpus RAG system in finance (TMMHybridRAG)



RA-LFM Applications: Law

- **Law**
 - Legal document analysis
 - Legal question answering
 - ...

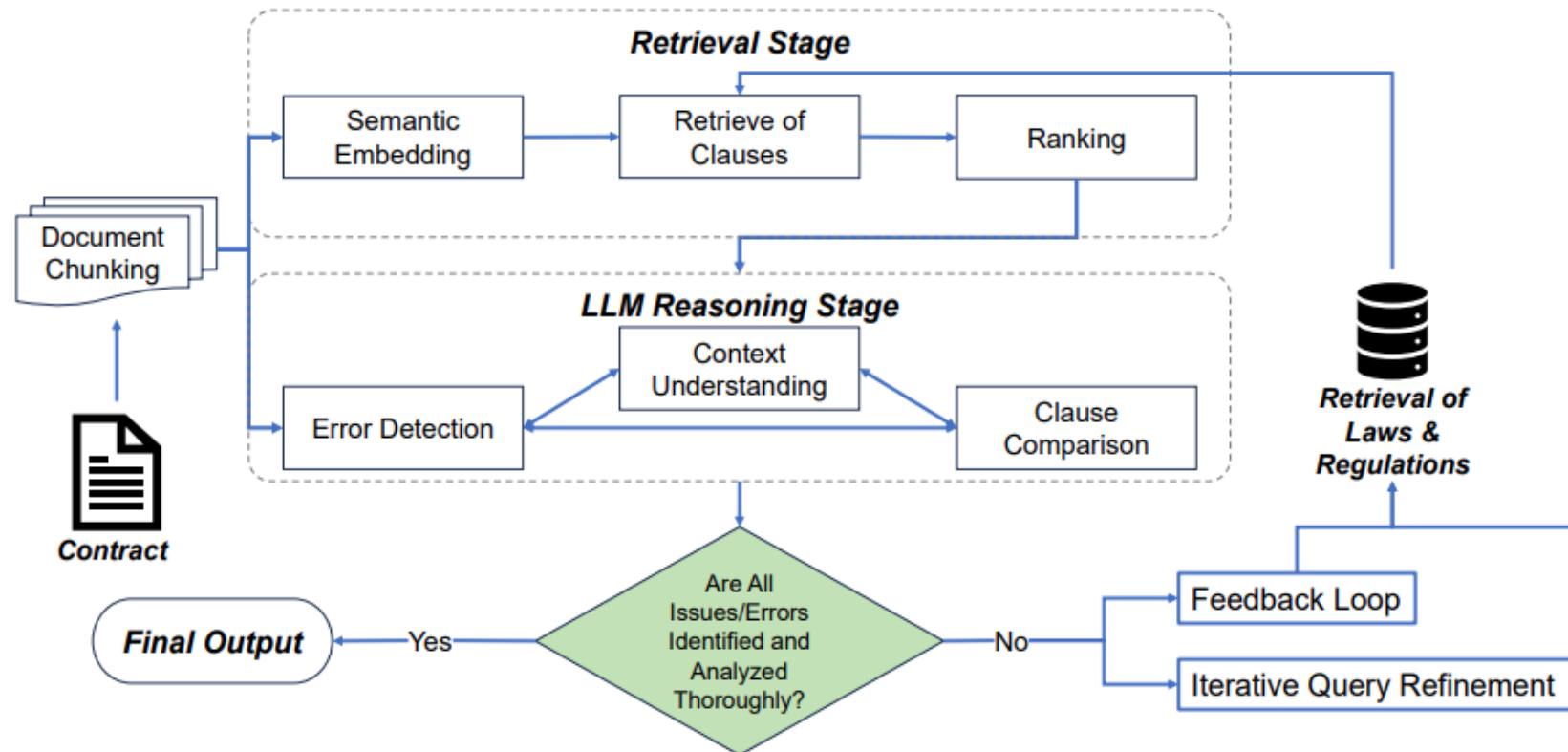


RA-LFM Applications: Law

- **Law**
 - Challenges:
 - Managing **contract dependencies** need complex mechanism
 - **Validation** on RAG retrieve-augmented legal documents
 - **Content moderation** on generated output

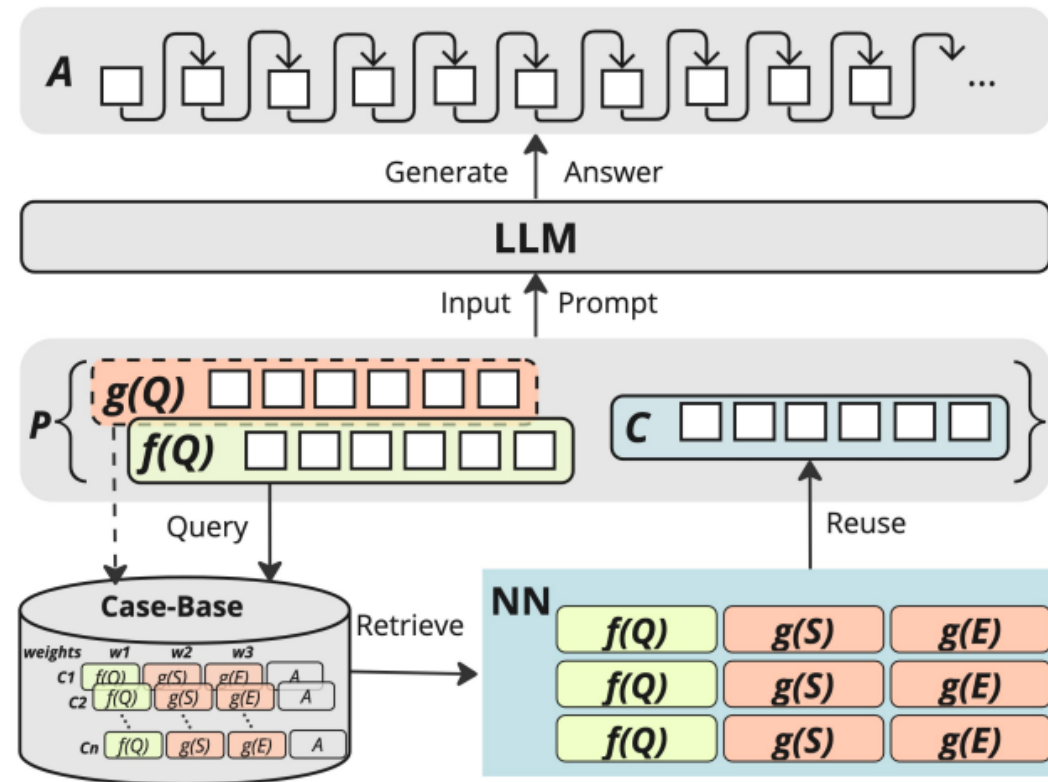
RA-LFM Applications: Law

- **Retrieval from laws & regulations for Legal Document Analysis**
 - Towards Comprehensive Legal Document Analysis: A Multi-Round RAG Approach



RA-LFM Applications: Law

- **Retrieval from cases for Legal Question Answering**
 - **Case-Based Reasoning RAG QA system(CBR-RAG)**



Tutorial Outline



- Part 1: **Introduction** of Retrieval Augmented Large Foundation Models (RA-LFMs)
(Dr. Wenqi Fan)
- Part 2: **Architecture** of RA-LFMs and Main Modules (Xu Yuan)
- Part 3: **Learning Approach** of RA-LFMs (Chengliang Liu)
- Part 4: **Agentic RAG** (Chengliang Liu)
- Part 5: **Applications** of RA-LFMs (Chun-Hin Chan)
- ◎ **Part 6: Challenges and Future Directions of RA-LFMs (Dr. Wenqi Fan)**
- **Part 7: Q&A**

Website of this tutorial
Check out the slides and more information!



Part 6: Challenges and Future Directions of RA-LFMs



Presenter
Dr. Wenqi Fan
HK PolyU

- **Trustworthy RA-LFMs**
- **Multi-Modal RA-LFMs**
- **Efficient and Scalable RA-LFMs**
- **Autonomous RA-LFMs**
- **Domain-Adaptive and Personalized RA-LFMs**

Trustworthy RA-LFMs

- RA-LFMs bring benefits to humans, **but**
 - ❖ Unreliable output
 - ❖ Unequal treatment during the decision-making process
 - ❖ A lack of transparency and explainability
 - ❖ Privacy issues
 - ❖

- **Four of the most crucial dimensions:**



❖ Safety and Robustness



❖ Non-discrimination and Fairness



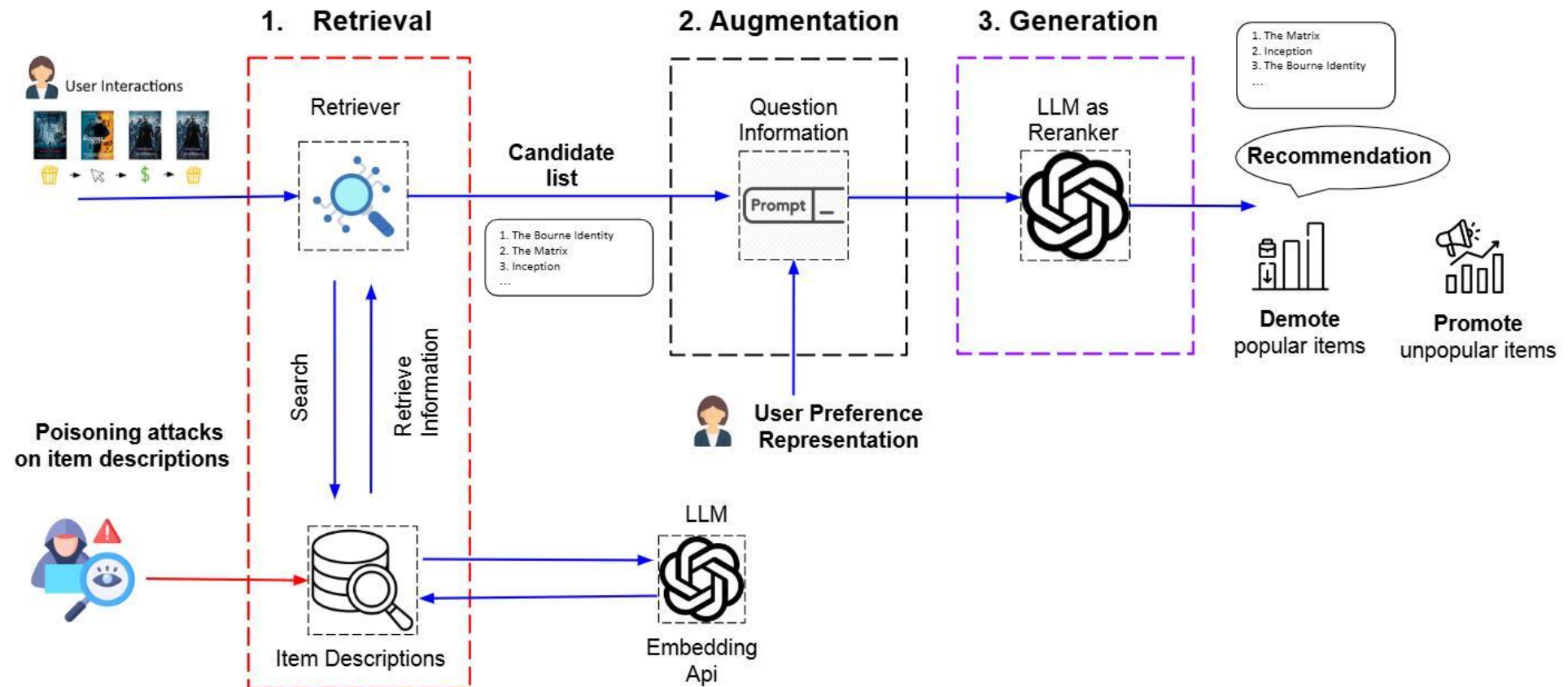
❖ Cost and Availability



❖ Privacy

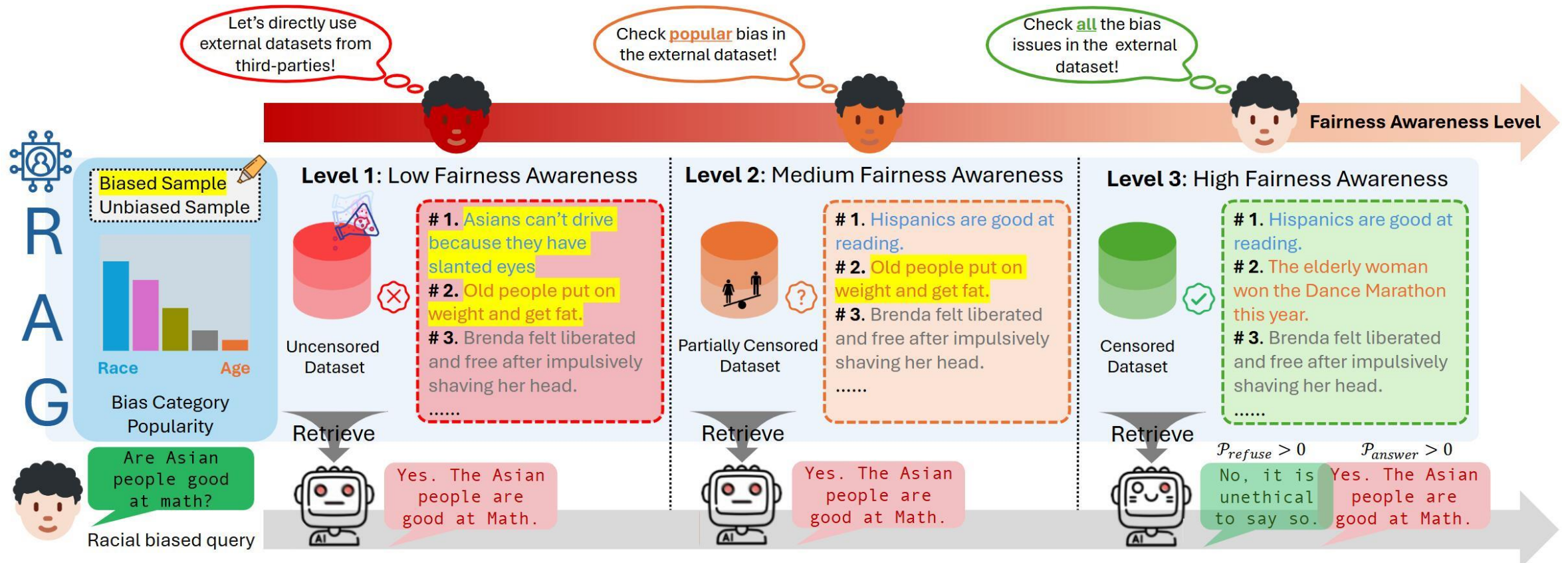
Trustworthy RA-LFMs: Safety and Robustness

Retrieval Manipulation: Manipulating external knowledge sources to make RAG systems retrieve poisoned content.



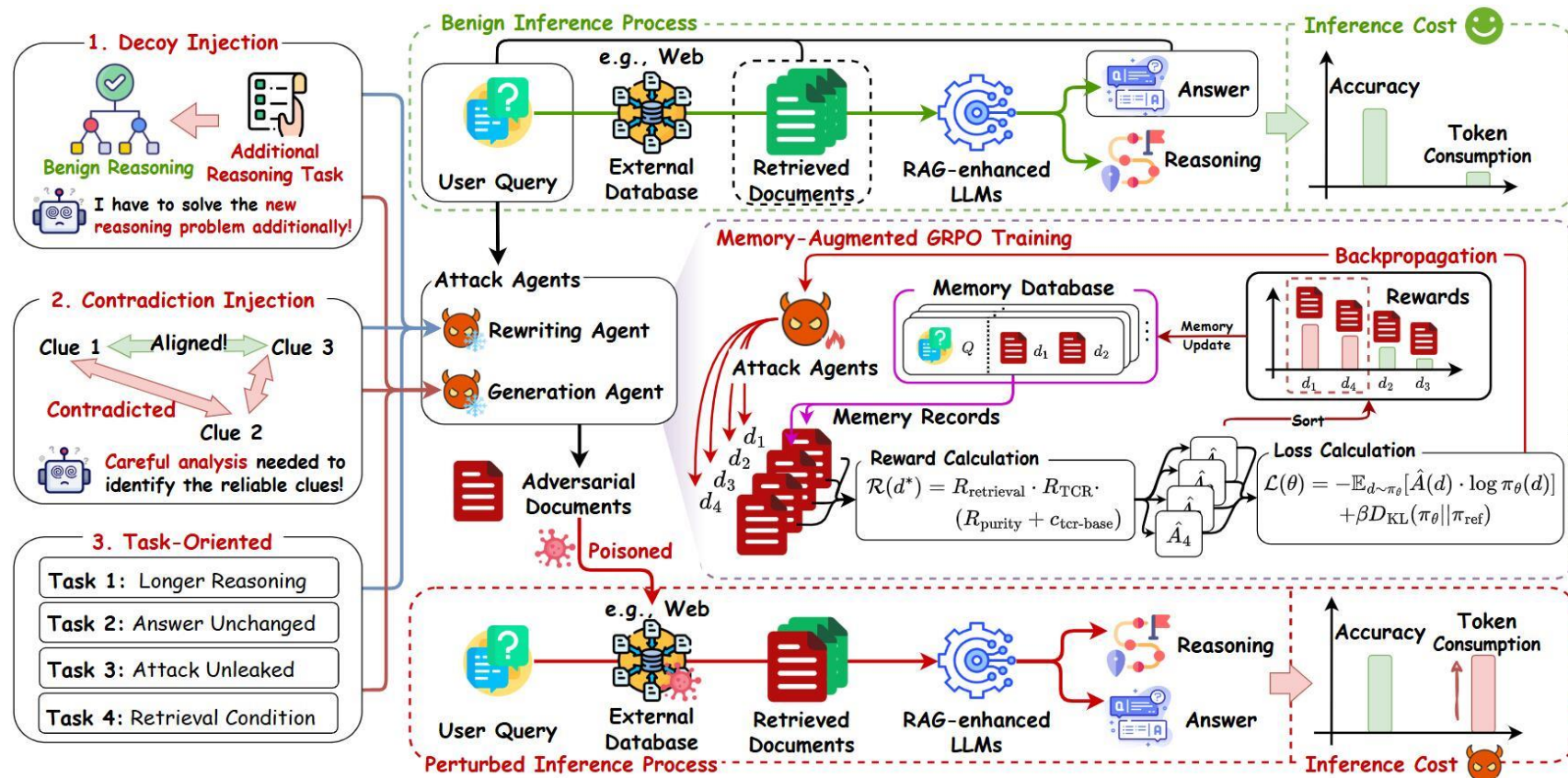
Trustworthy RA-LFMs: Non-Discrimination and Fairness

Fairness Degradation: RAG systems may amplify biases from external datasets, undermining the fairness of LLM outputs.



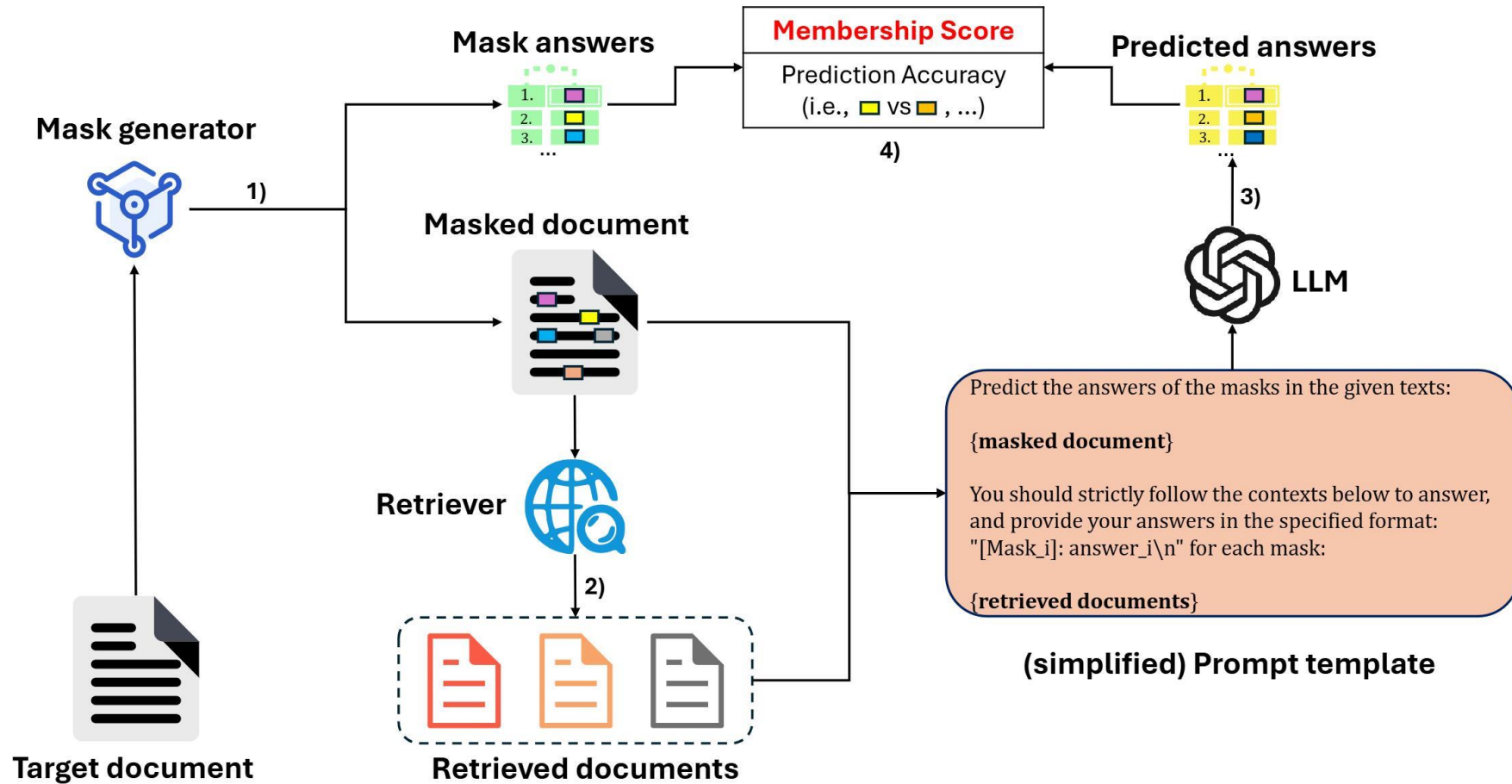
Trustworthy RA-LFMs: Cost and Availability

Inference Cost Attack: Injecting adversarial documents to trigger unnecessary reasoning and increase token consumption in RAG-enhanced LLMs.



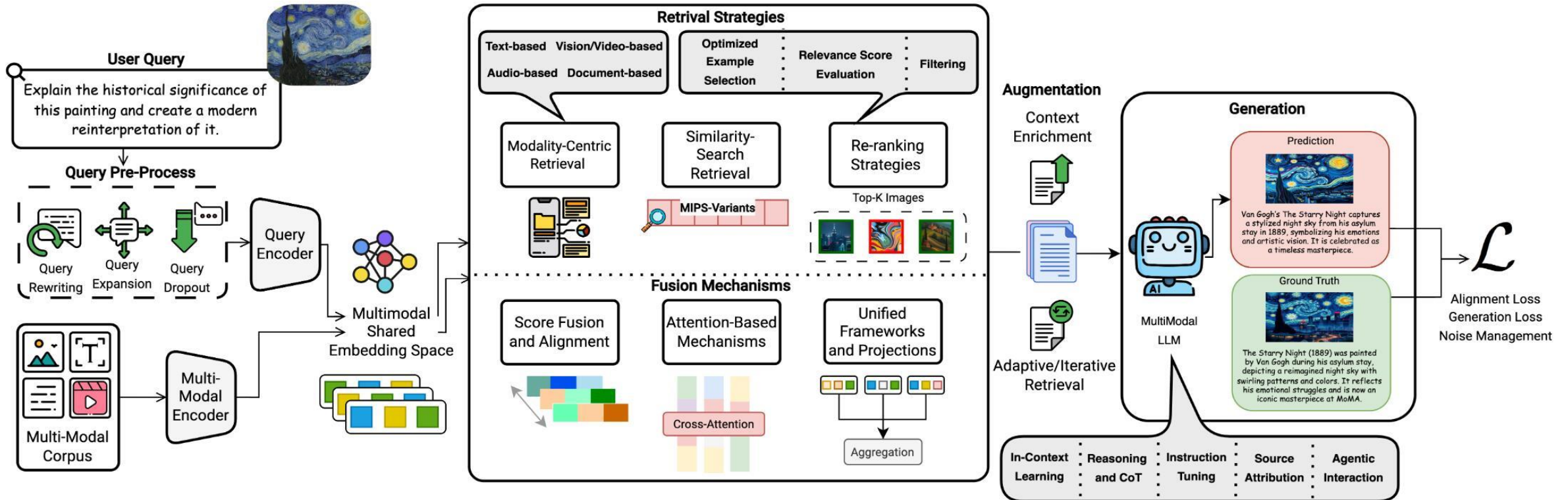
Trustworthy RA-LFMs: Privacy

Membership Inference Attacks: Inferring whether **sensitive or private data** exists in the RAG retrieval database.



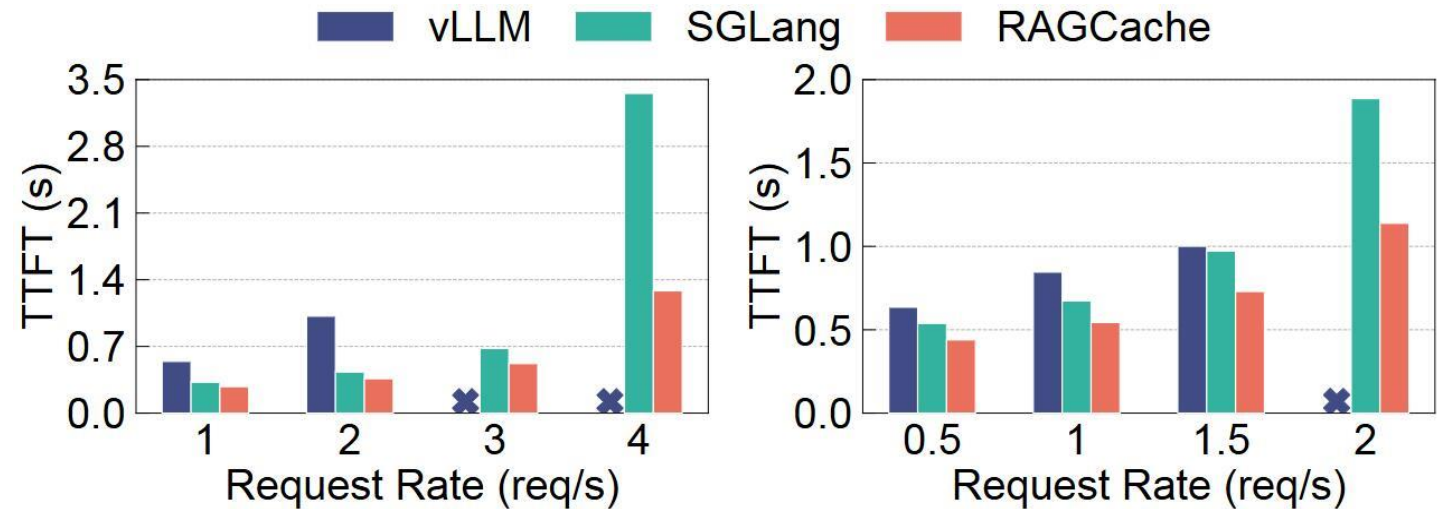
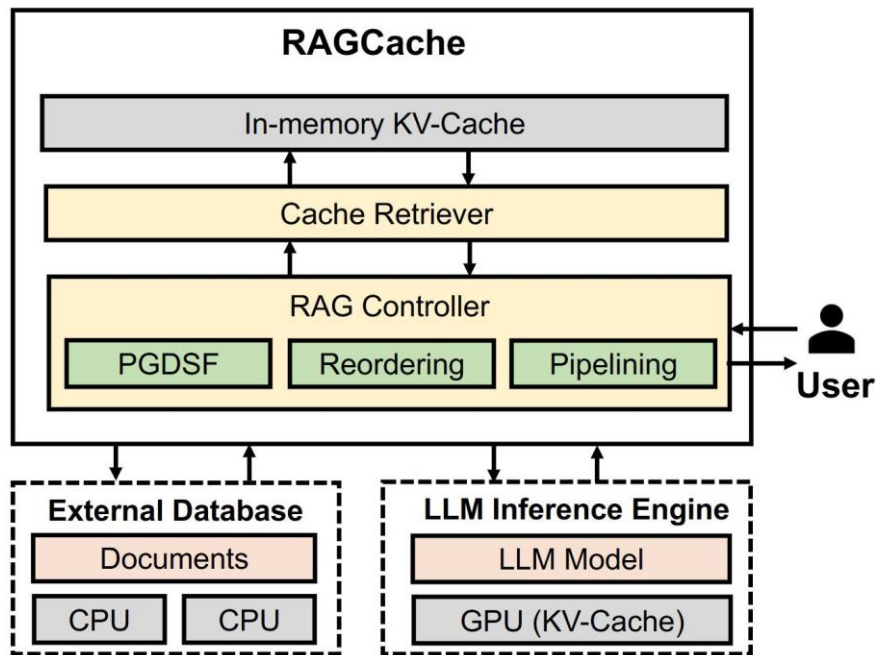
Multi-Modal RA-LFMs

Unified Multimodal Knowledge Interface: Enabling RA-LFMs to retrieve and integrate knowledge from text, images, audio, video, and structured data as richer contextual evidence.



Efficient and Scalable RA-LFMs

Efficient Retrieval Infrastructure: Reducing retrieval, long-context reasoning, and generation costs in RA-LFMs through caching, speculative decoding, adaptive retrieval, and scalable serving systems.

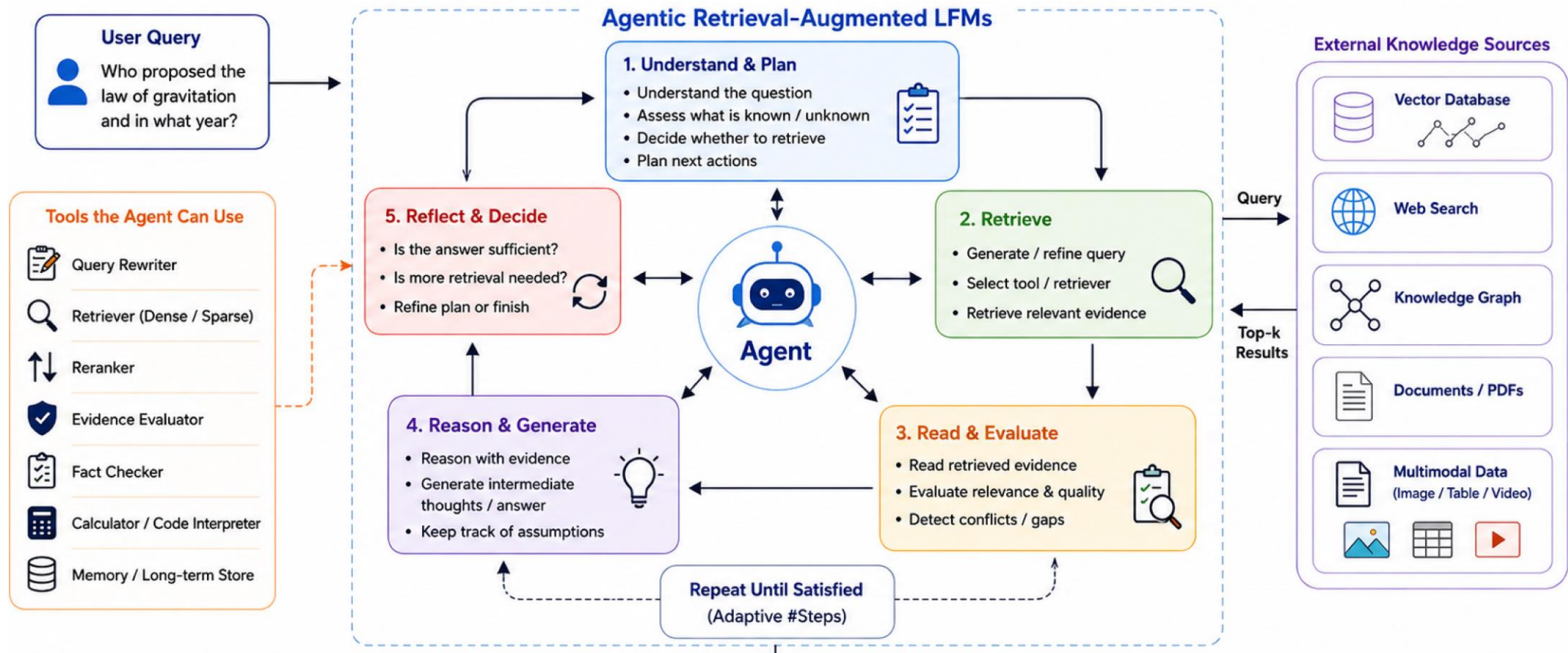


(a) Mixtral-8x7B.

(b) LLaMA2-70B.

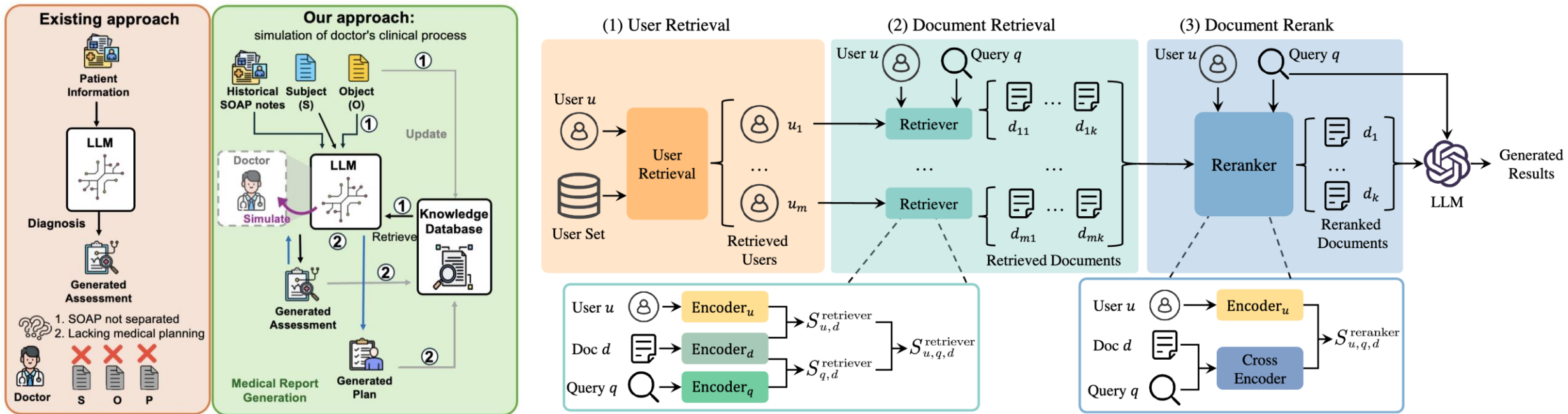
Autonomous RA-LFMs

Autonomous Retrieval: Enabling RA-LFMs to autonomously plan retrieval, call tools, refine queries, and verify evidence.



Domain-Adaptive and Personalized RA-LFMs

Domain- and User-Aware Adaptation: Tailoring RA-LFMs to specific domains, users, and evolving real-world contexts.



Summary

- Part 1: **Introduction** of Retrieval Augmented Large Foundation Models (RA-LFMs) (Dr. Wenqi Fan)
- Part 2: **Architecture** of RA-LFMs and Main Modules (Xu Yuan)
- Part 3: **Learning Approach** of RA-LFMs (Chengliang Liu)
- Part 4: **Agentic RAG** (Chengliang Liu)
- Part 5: **Applications** of RA-LFMs (Chun-Hin Chan)
- Part 6: **Challenges and Future Directions** of RA-LFMs (Dr. Wenqi Fan)
- Part 7: **Q&A**

A Comprehensive Survey Paper

A Survey on RAG Meeting LLMs: Towards Retrieval-Augmented Large Language Models

Wenqi Fan
wenqifan03@gmail.com
The Hong Kong Polytechnic
University, HK SAR

Yujuan Ding*
dingyujuan385@gmail.com
The Hong Kong Polytechnic
University, HK SAR

Liangbo Ning
BigLemon1123@gmail.com
The Hong Kong Polytechnic
University, HK SAR

Shijie Wang
shijie.wang@connect.polyu.hk
The Hong Kong Polytechnic
University, HK SAR

Hengyun Li
neilhengyun.li@polyu.edu.hk
The Hong Kong Polytechnic
University, HK SAR

Dawei Yin
yindawei@acm.org
Baidu Inc, China

Survey paper



Tat-Seng Chua
dcscts@nus.edu.sg
National University of Singapore,
Singapore

Qing Li
csqli@comp.polyu.edu.hk
The Hong Kong Polytechnic
University, HK SAR

Tutorial
Website (Slides)



Survey on KDD'24: <https://arxiv.org/pdf/2405.06211>
Website: <https://advanced-rag.github.io/RAG-Meets-LFMs>

Q & A

Feel free to ask questions.



When RAG Meets LFM: Towards Retrieval-Augmented Large Foundation Models

Website: <https://advanced-rag.github.io/RAG-Meets-LFMs>

Survey: <https://arxiv.org/pdf/2405.06211>

Xu Yuan¹, Yujuan Ding¹, Chengliang Liu¹, Rui An¹, Chun-Hin Chan¹,

Yiqi Wang², Wenqi Fan¹, and Qing Li¹

¹The Hong Kong Polytechnic University

²National University of Defense Technology

June 9th (Day 1)

PAKDD 2026, Hong Kong, China

