

Tutorial Outline



41st IEEE International Conference
on Data Engineering
— HONG KONG SAR, CHINA | MAY 19 – 23, 2025 —



- ⦿ **Part 1: Introduction** of Retrieval Augmented Large Language Models (RA-LLMs) (Dr. Yujuan Ding)
- ⦿ **Part 2: Architecture** of RA-LLMs and **Main Modules** (Dr. Yujuan Ding)
- ⦿ **Part 3: Data Management** for RA-LLMs (Pangjing Wu)
- ⦿ **Part 4: Learning Approach of RA-LLMs (Liangbo Ning)**
- ⦿ **Part 4: Applications** of RA-LLMs (Shijie Wang)
- ⦿ **Part 5: Challenges and Future Directions** of RA-LLMs (Liangbo Ning)

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



Part 3: RA-LLM Learning



Presenter
Liangbo Ning
HK PolyU

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

Part 3: RA-LLM Learning

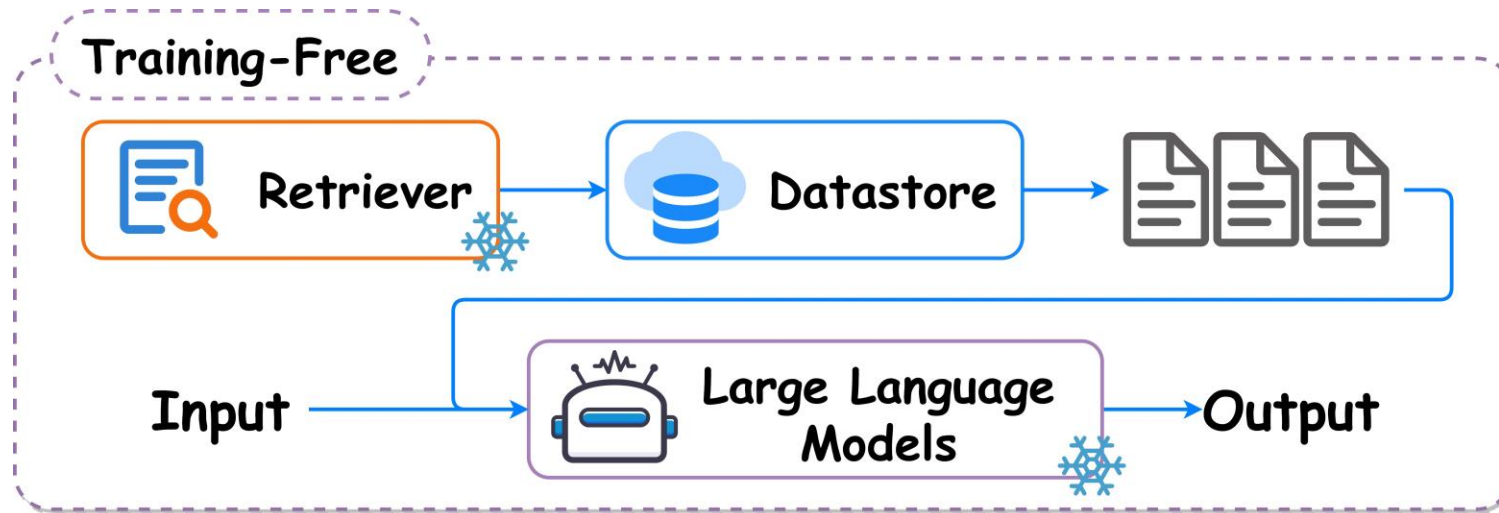


Presenter
Liangbo Ning
HK PolyU

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

RA-LLM Learning: Training-free

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

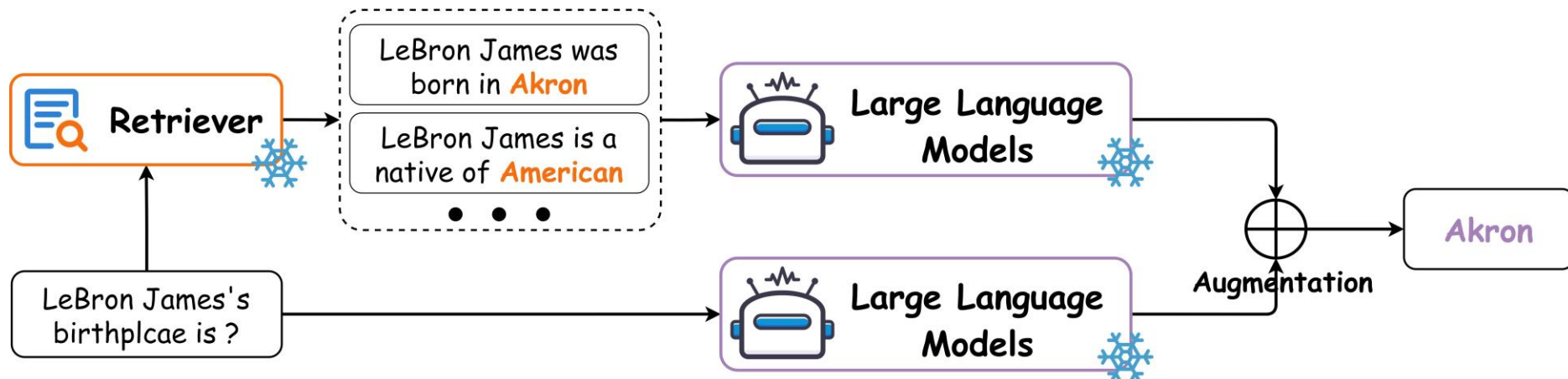


RA-LLM Learning: Training-free

- Prompt Engineering-based Methods

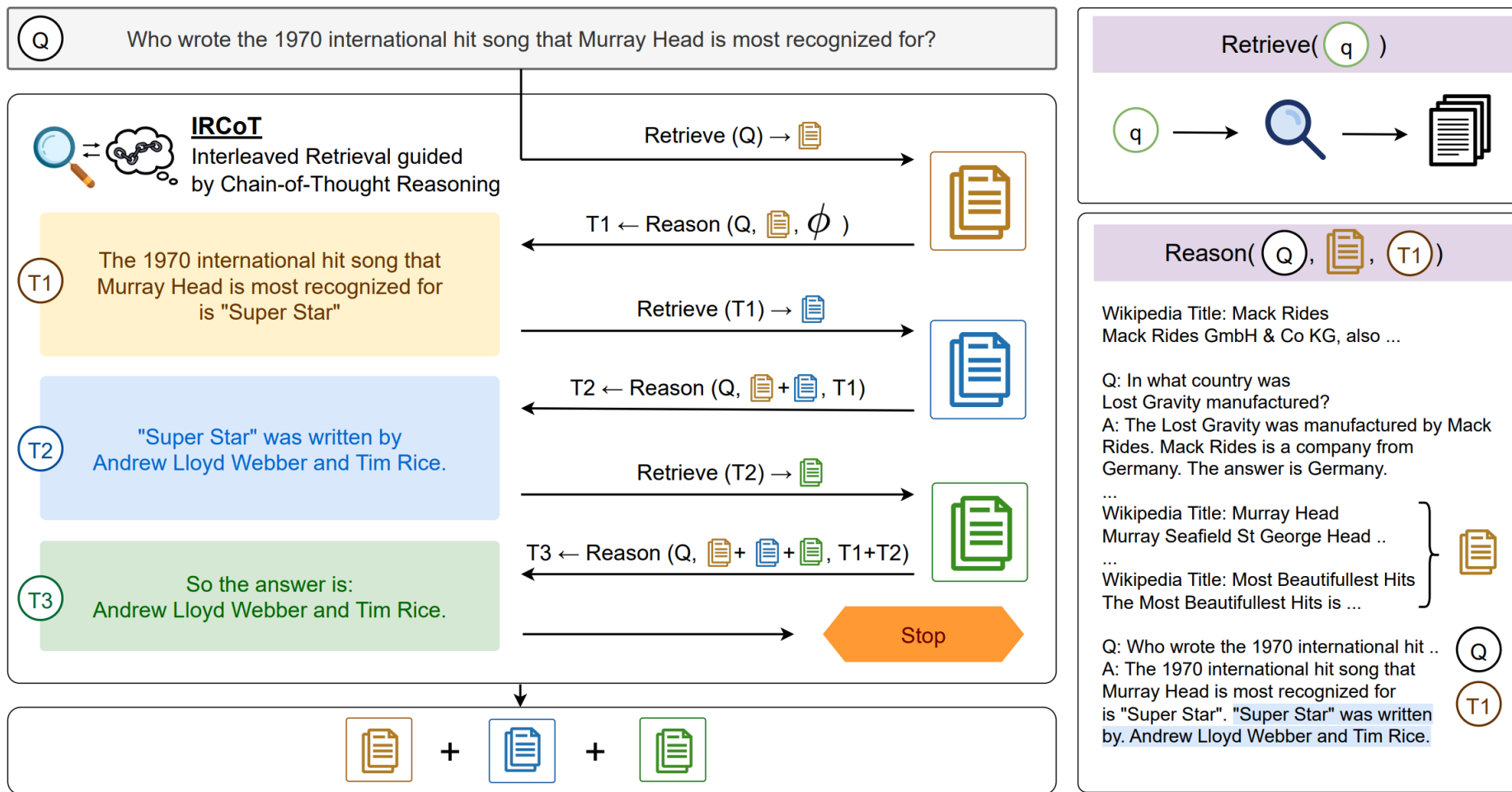


- Retrieval-Guided Token Generation Methods



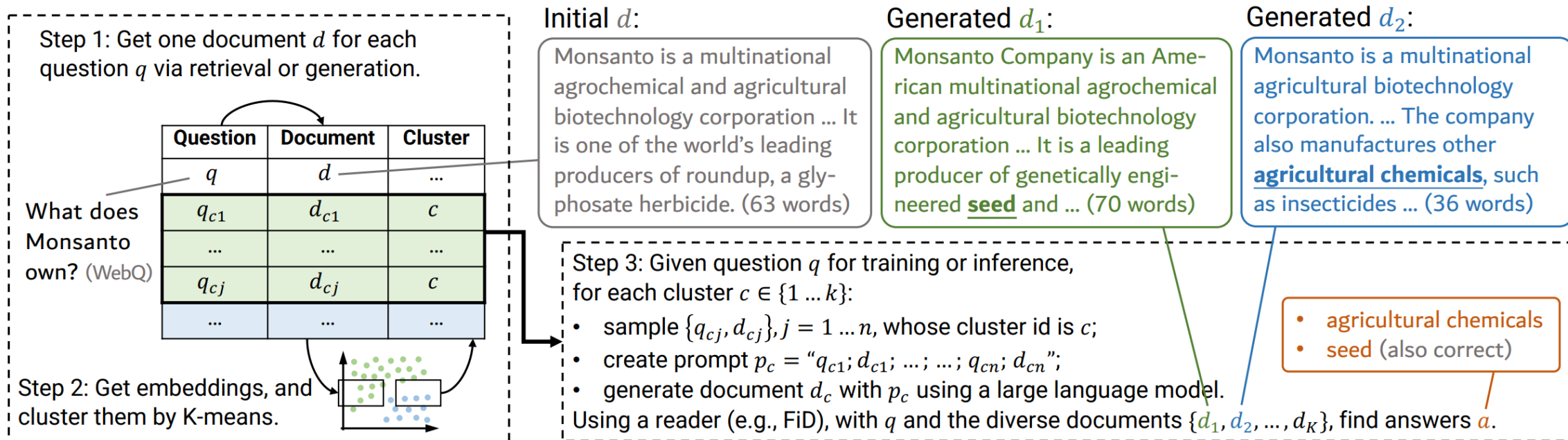
RA-LLM Learning: Training-free

- IRCoT



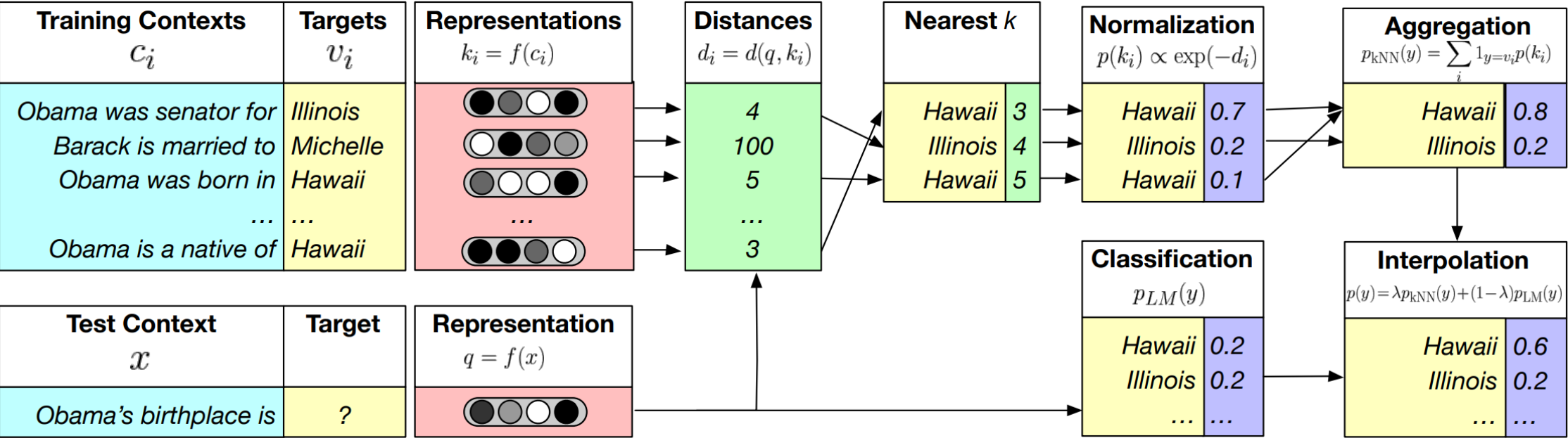
RA-LLM Learning: Training-free

- GENREAD



RA-LLM Learning: Training-free

- k*NN-LM**



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

RA-LLM 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-LLM Learning

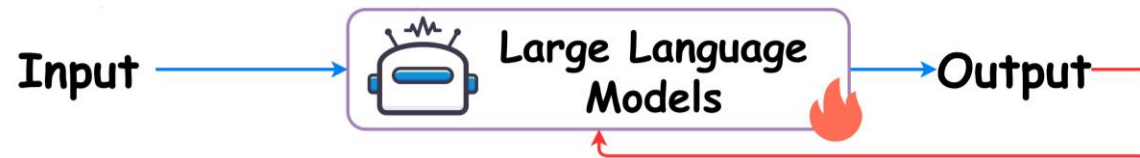


Website of this tutorial

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

RA-LLM Learning: Independent Training

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



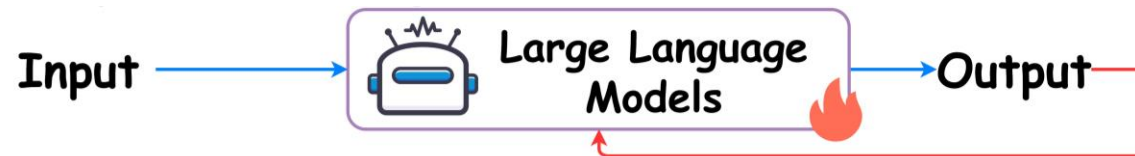
- Independent training of Retriever.



RA-LLM Learning: Independent Training

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

- Independent training of large language models.

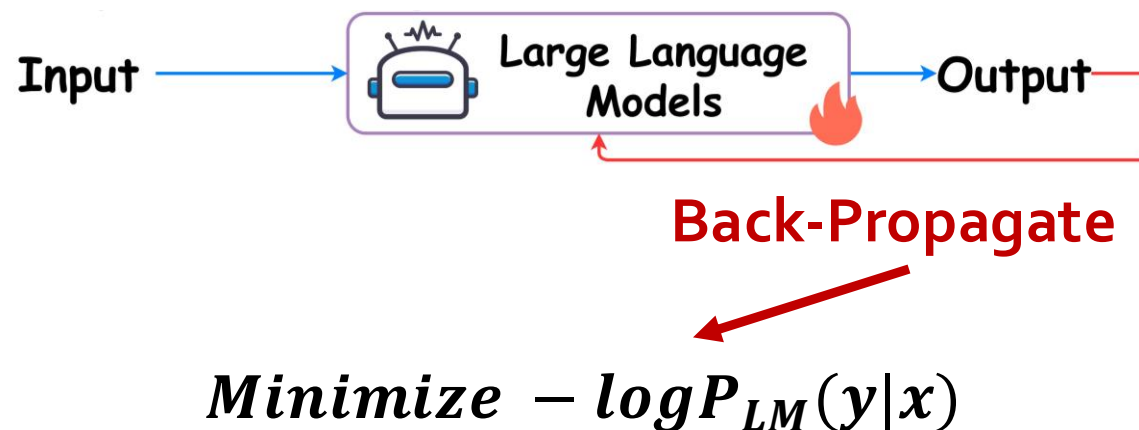


- Independent training of Retriever.



RA-LLM Learning: Independent Training

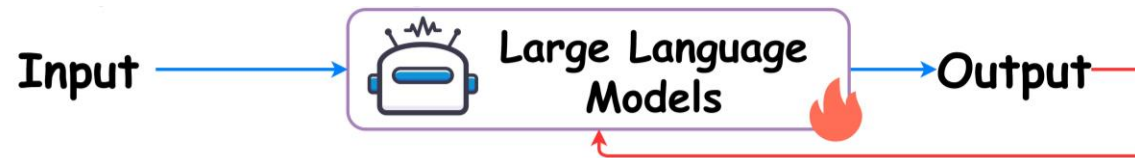
- Independent training of large language models.



.....

RA-LLM Learning: Independent Training

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

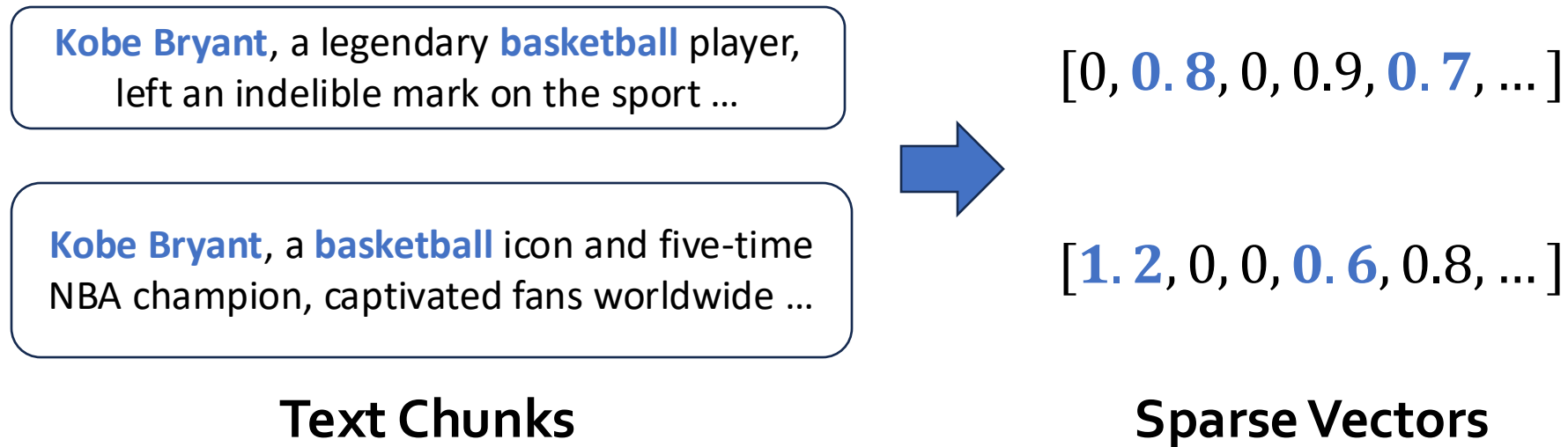


- Independent training of Retriever.



RA-LLM Learning: Independent Training

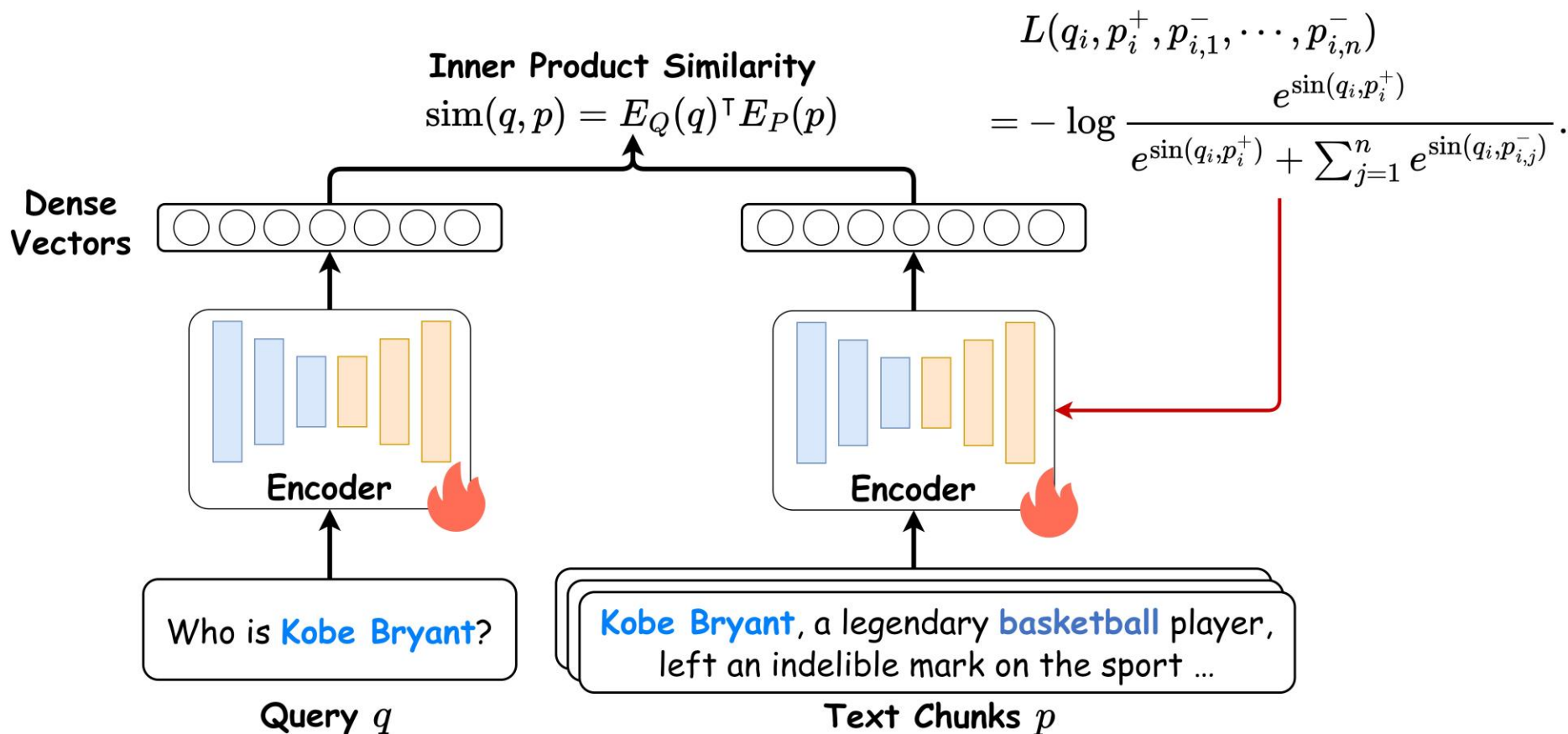
- Sparse retrieval models: TF-IDF / BM25



No training is Needed!

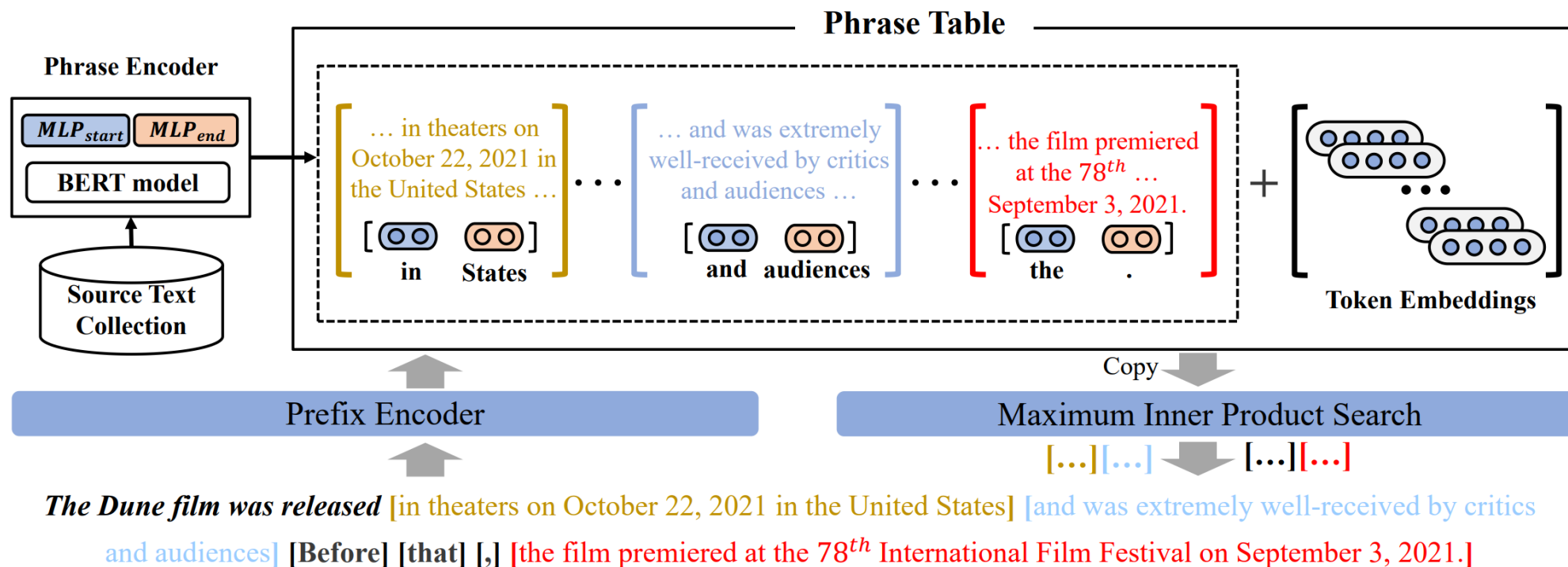
RA-LLM Learning: Independent Training

- Dense retrieval models: DPR



RA-LLM Learning: Independent Training

- Dense retrieval models: CoG



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

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

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

RA-LLM 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-LLM 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-LLM Learning

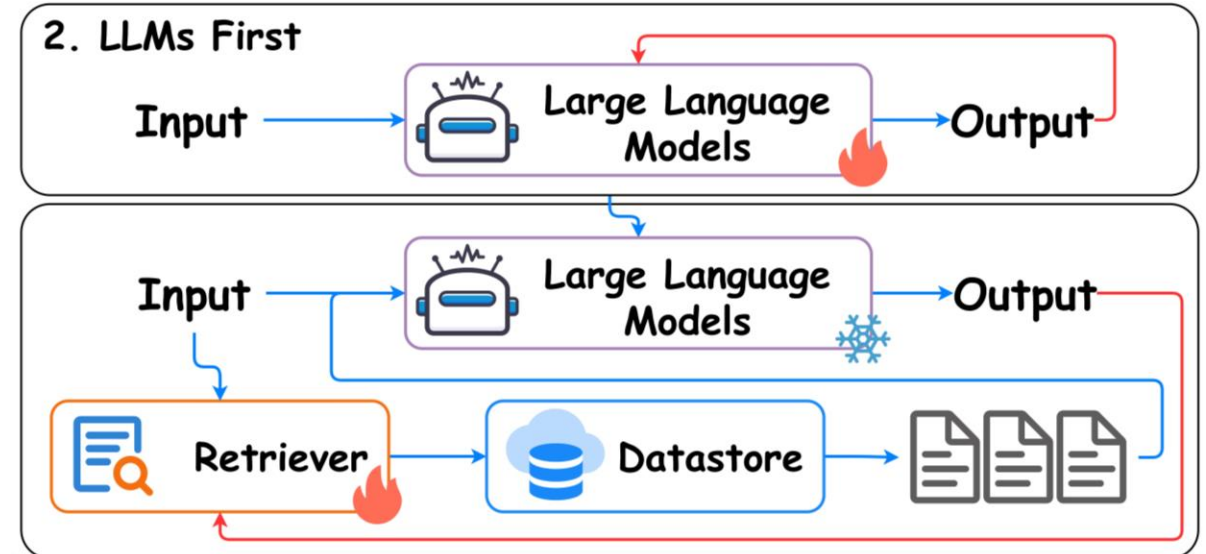
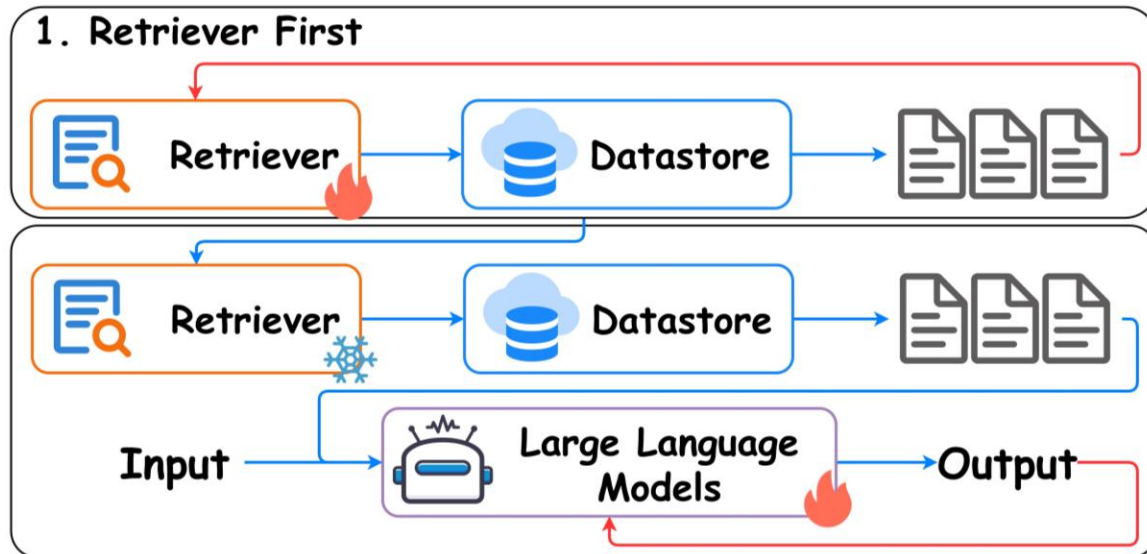


Website of this tutorial

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

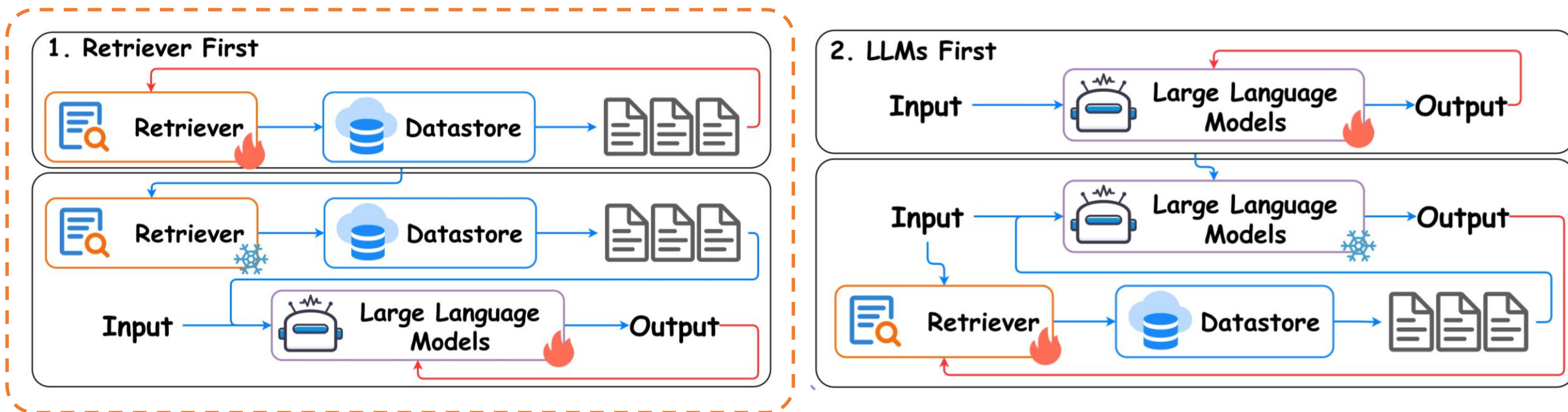
RA-LLM Learning: Sequential Training

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



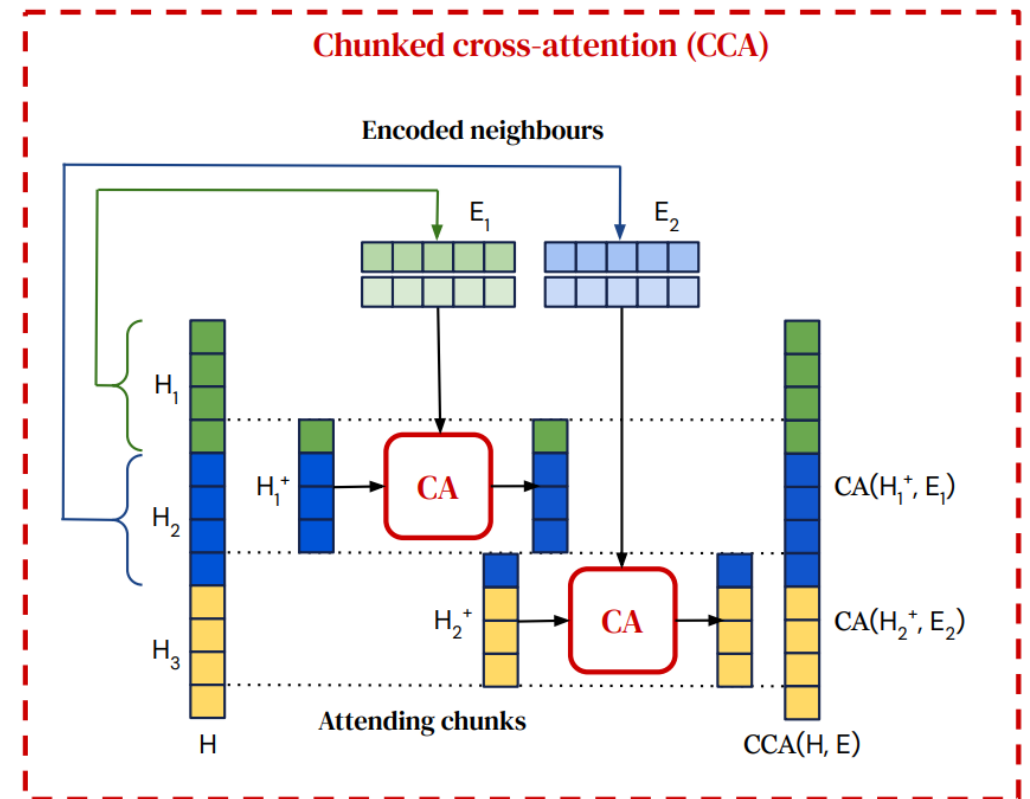
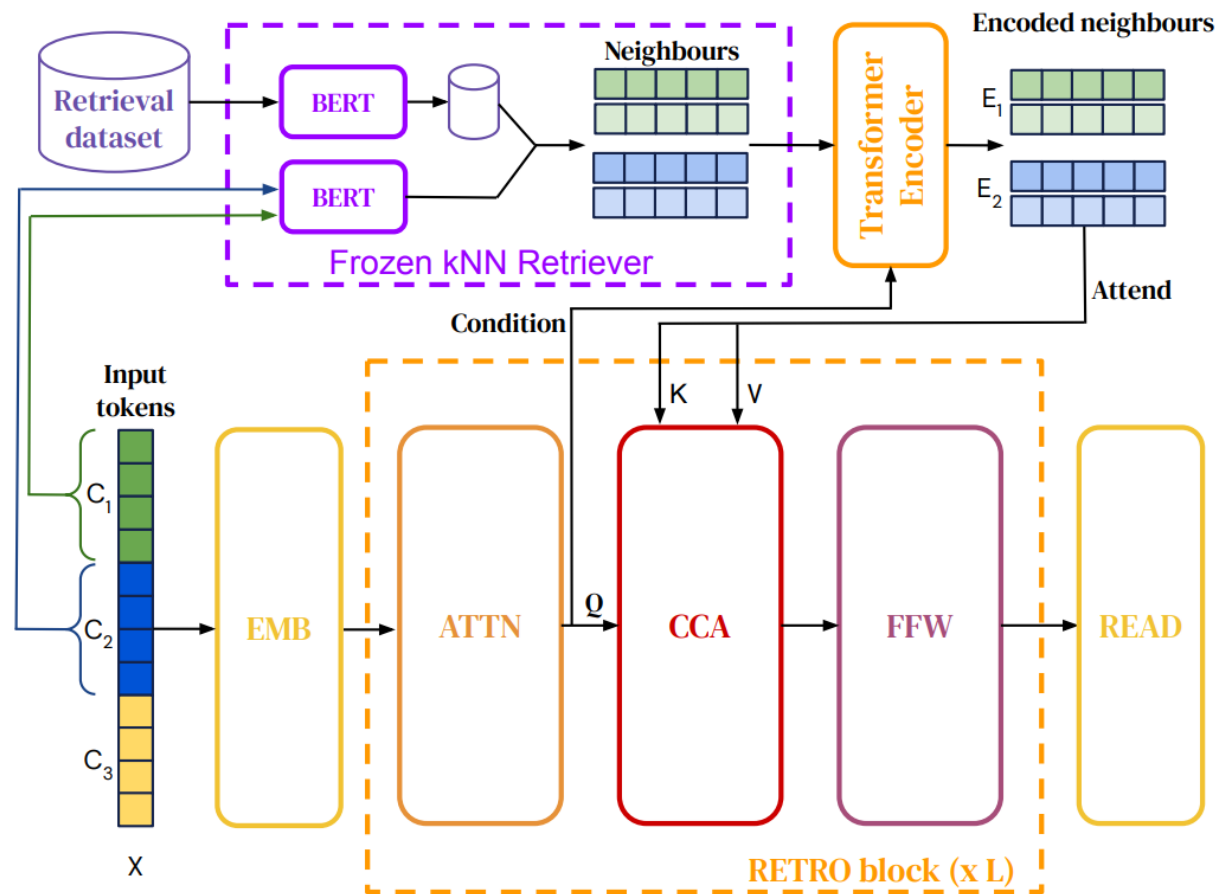
RA-LLM Learning: Sequential Training

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



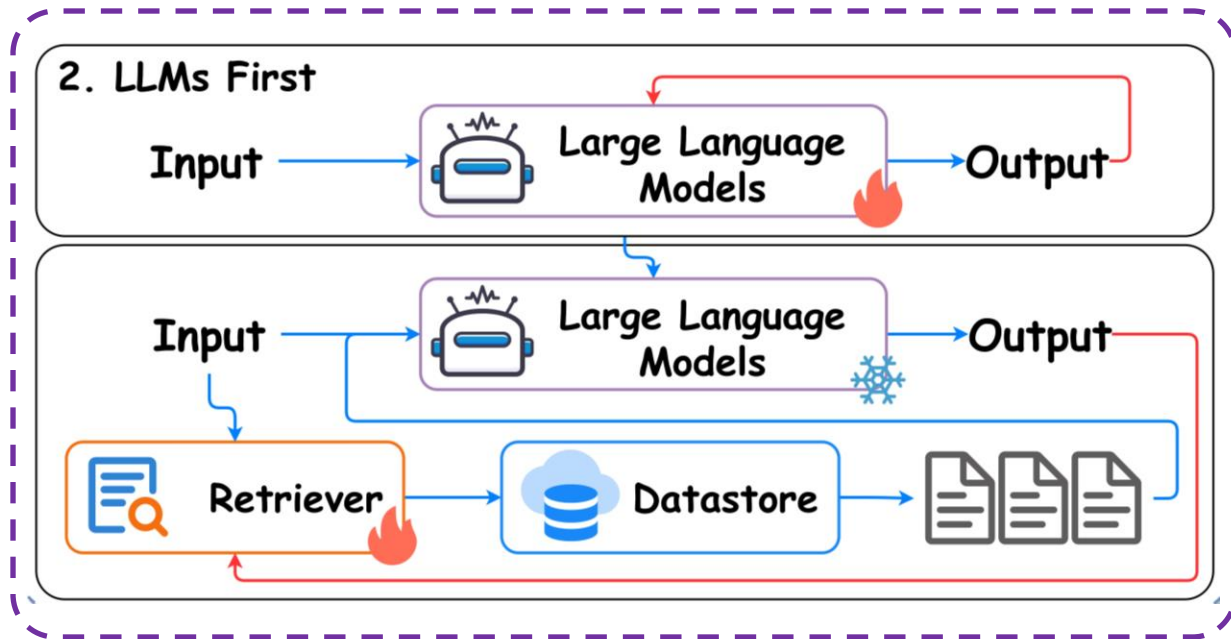
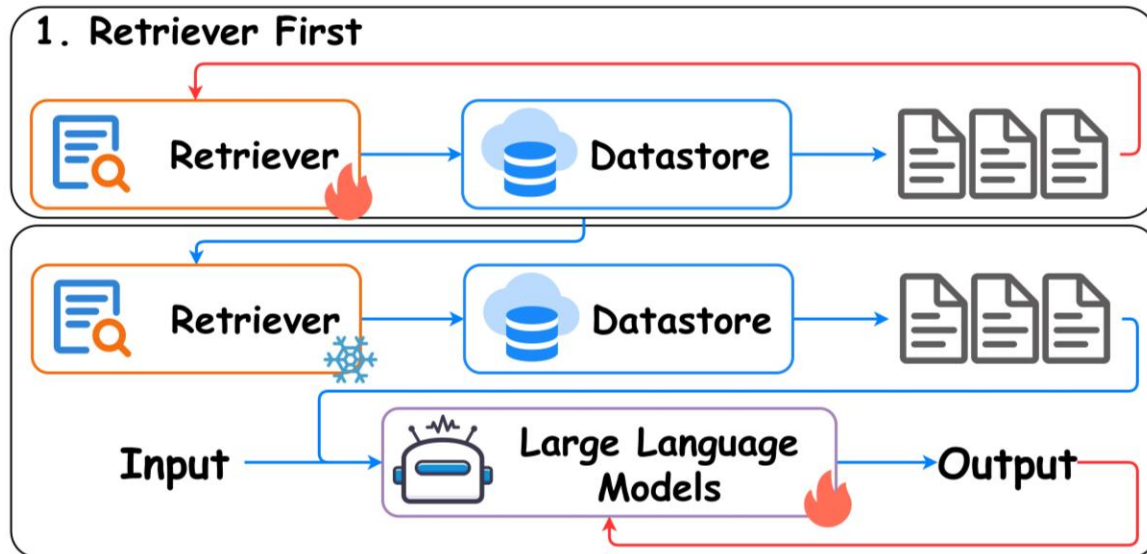
RA-LLM Learning: Sequential Training

- RETRO



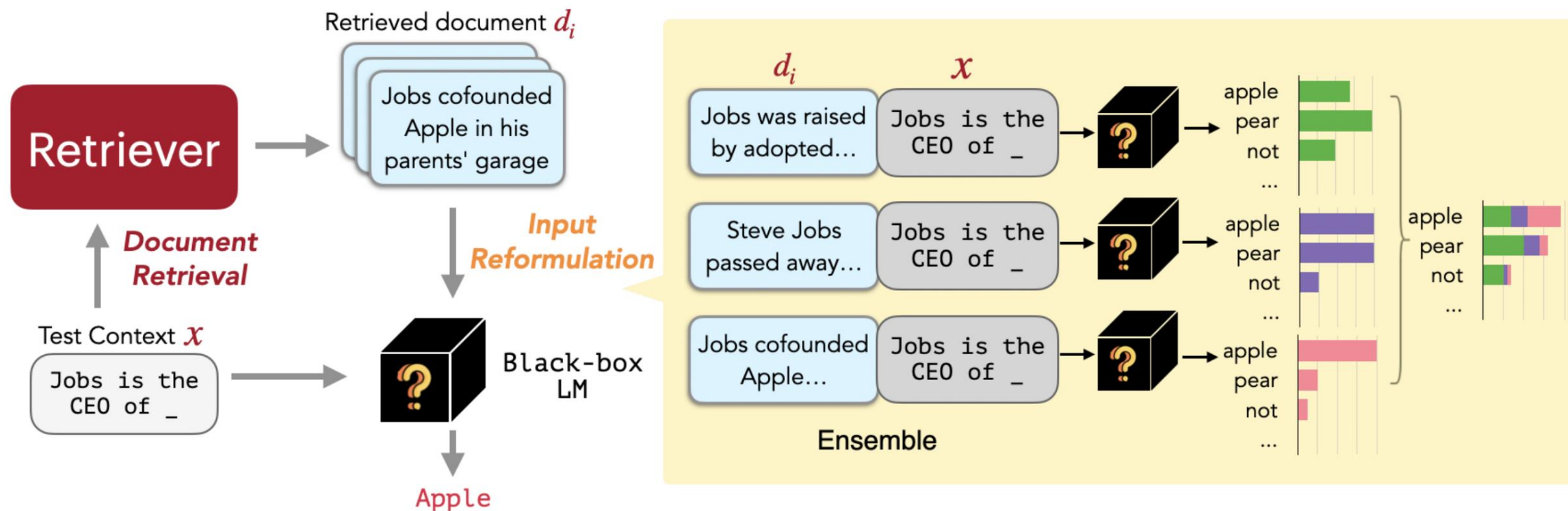
RA-LLM Learning: Sequential Training

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



RA-LLM Learning: Sequential Training

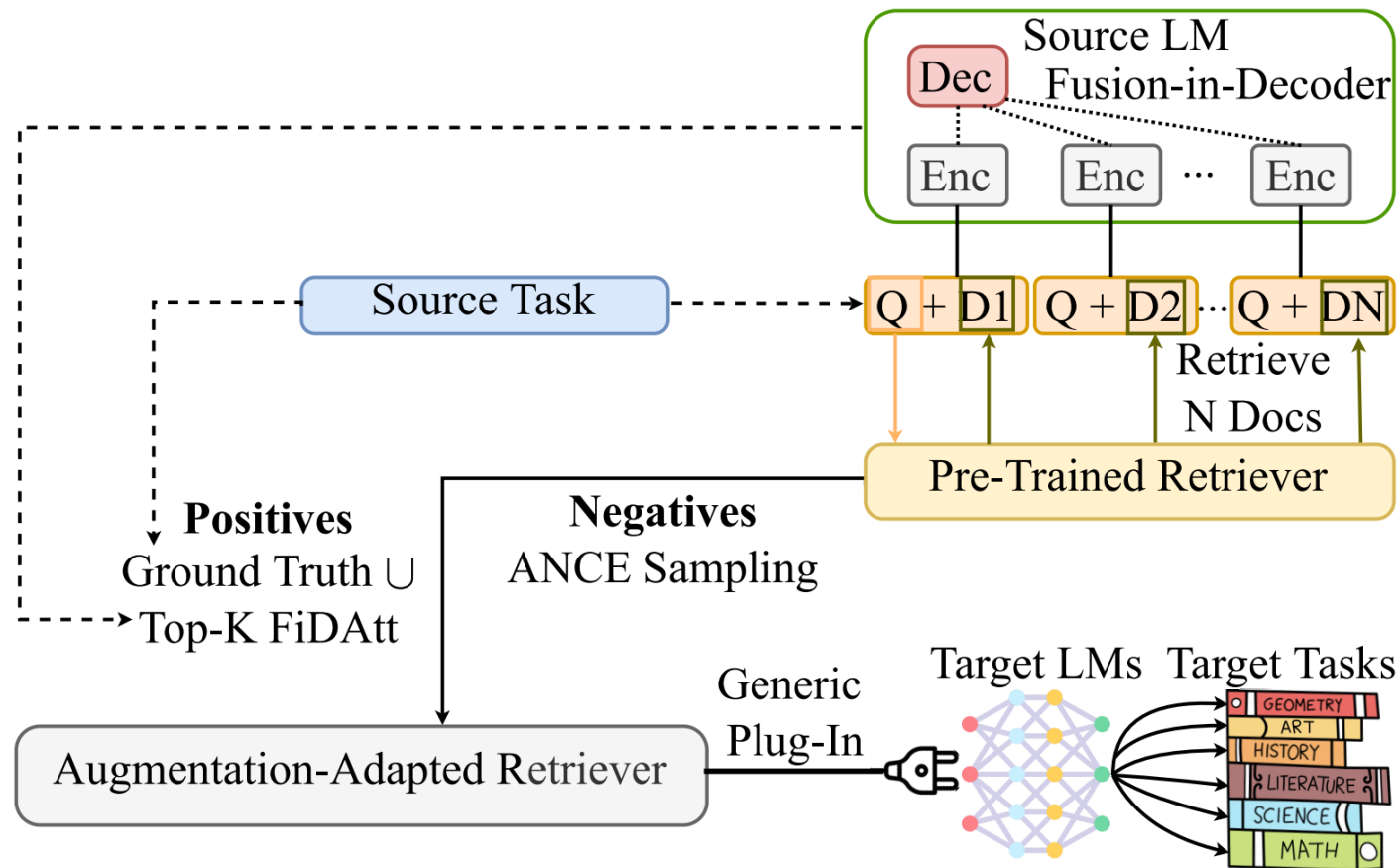
- REPLUG (Retrieve and Plug)



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

RA-LLM 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-LLM 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-LLM Learning

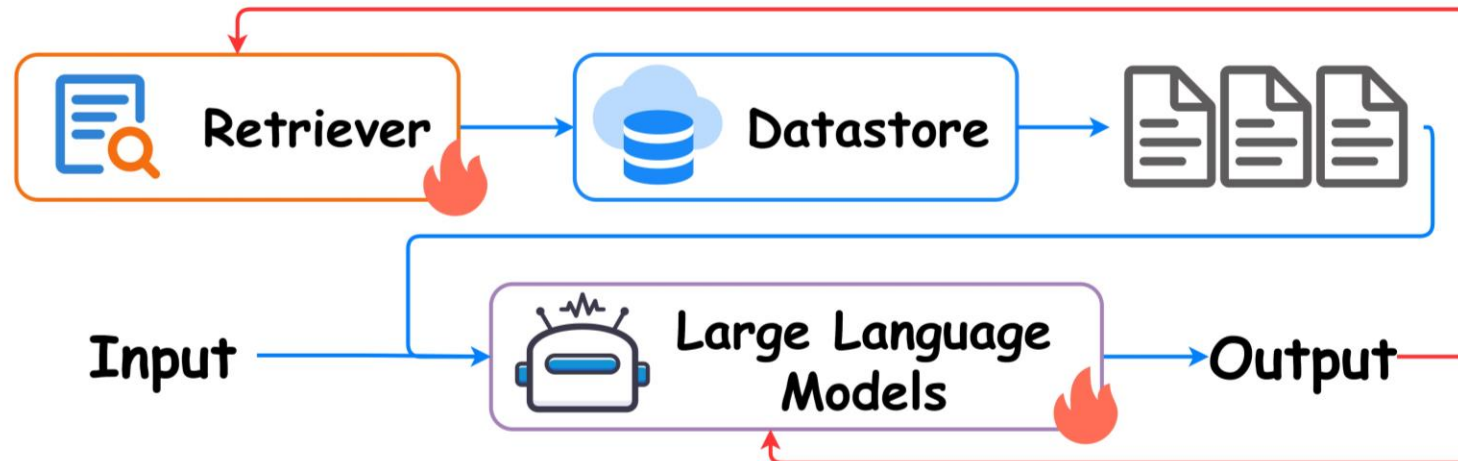


Website of this tutorial

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

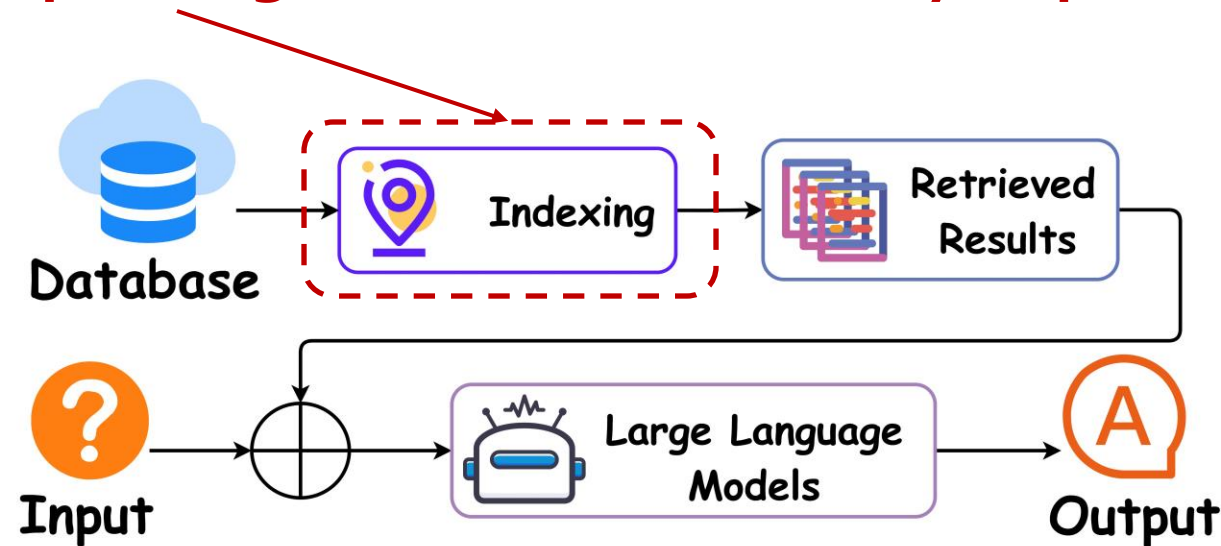
RA-LLM Learning: Joint Training

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



RA-LLM Learning: Joint Training

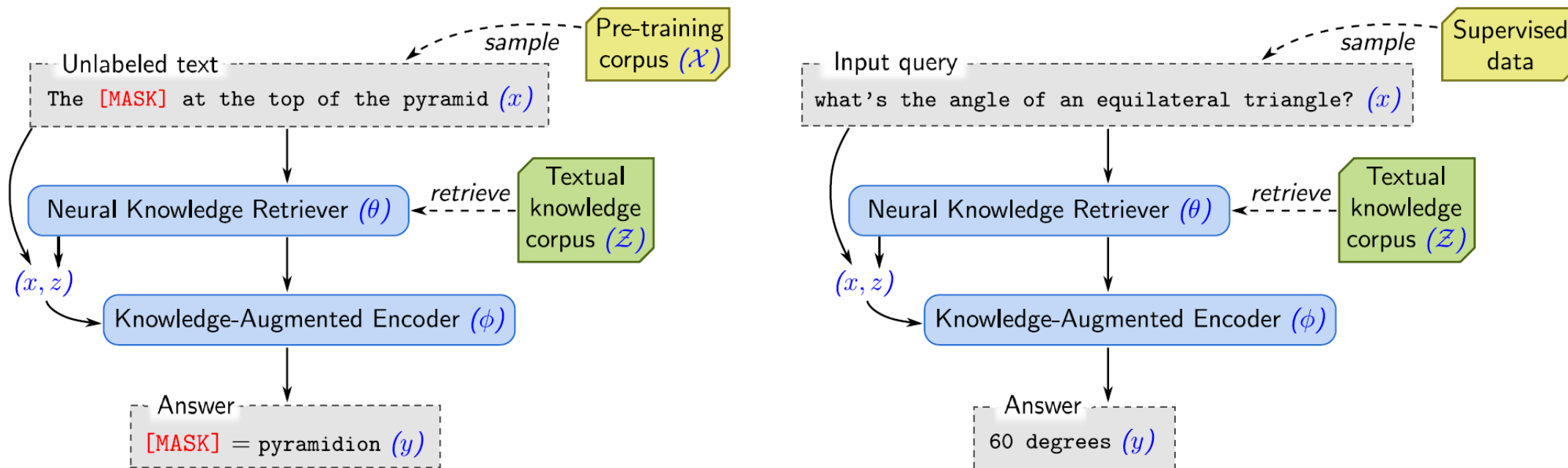
- Retrieval Index Updating, which could be very expensive!



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

RA-LLM Learning: Joint Training

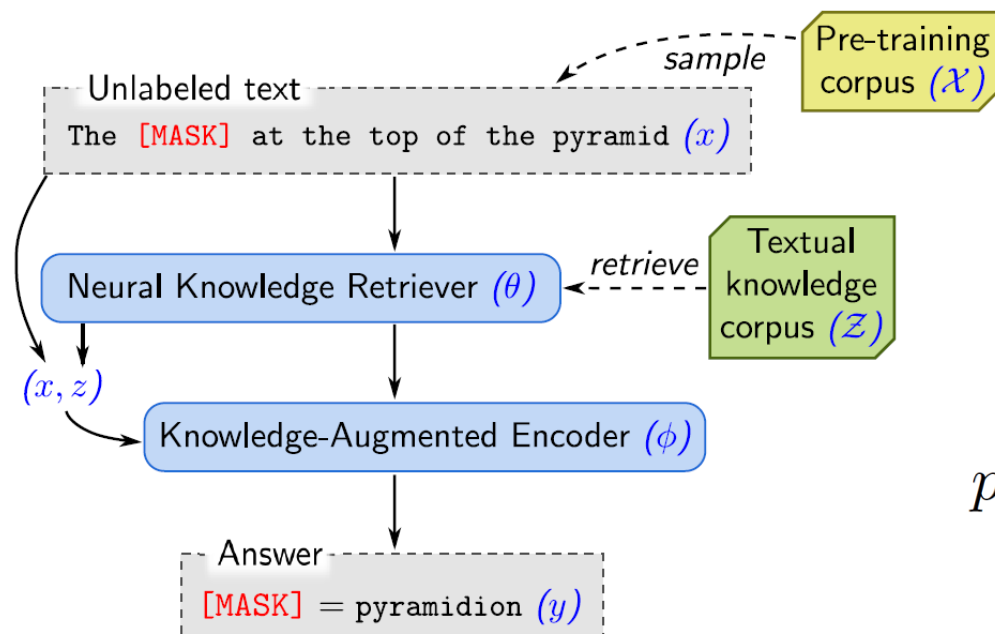
- REALM



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

RA-LLM Learning: Joint Training

- REALM

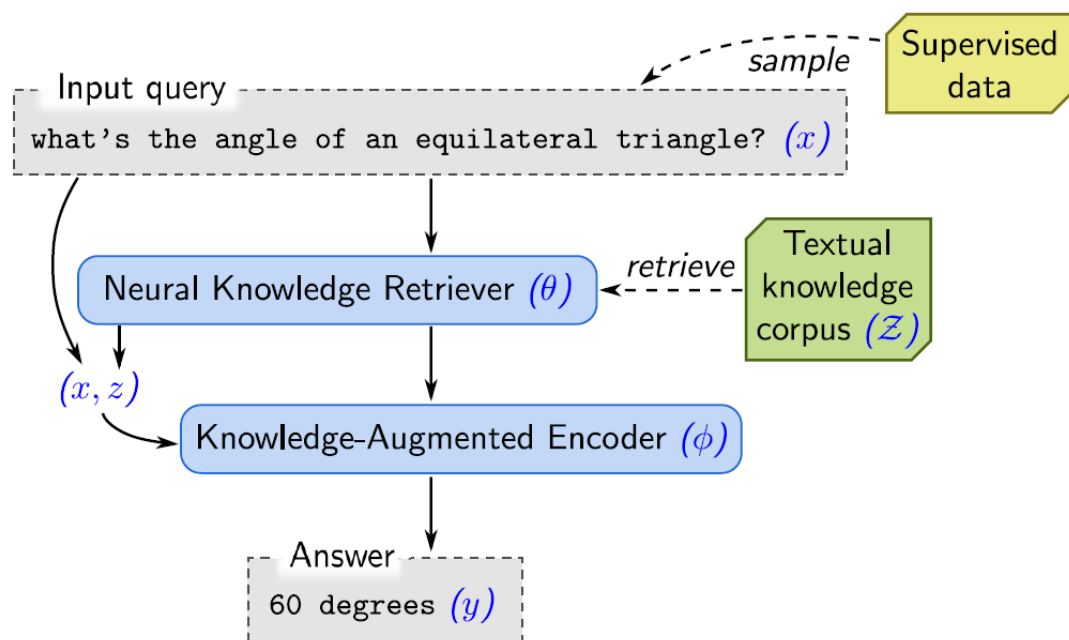


$$p(y | z, x) = \prod_{j=1}^{J_x} p(y_j | z, x)$$

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

RA-LLM Learning: Joint Training

- REALM



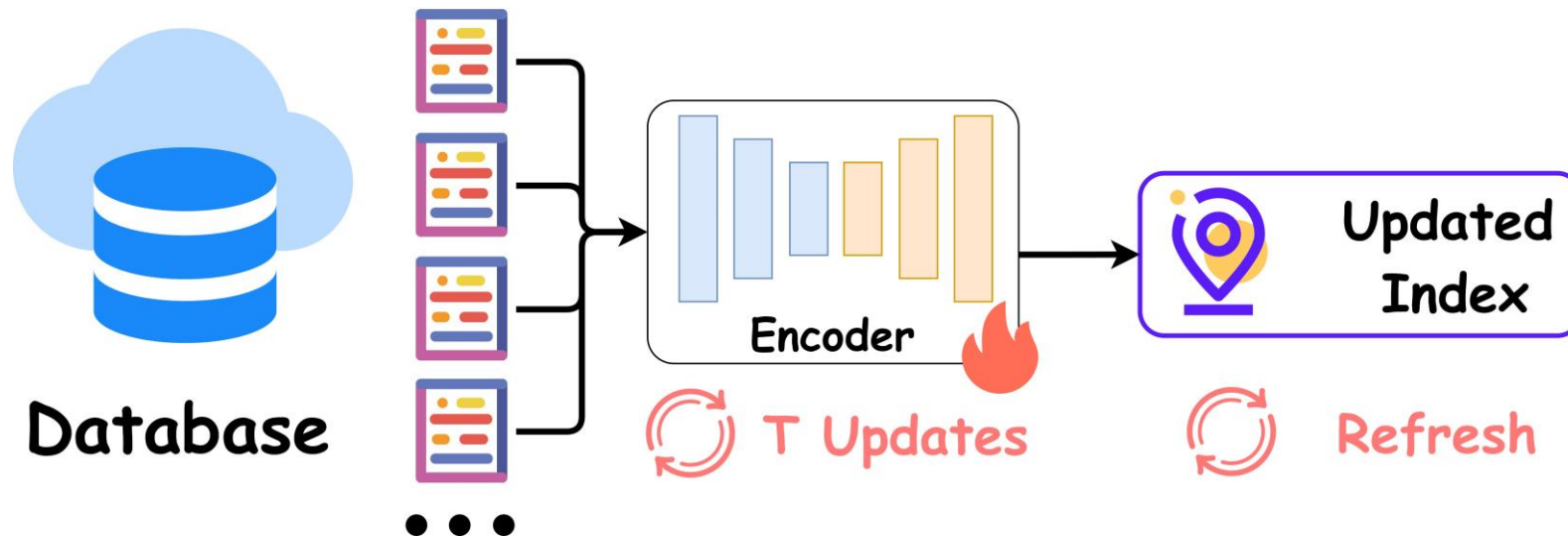
$$p(y | z, x) \propto \sum_{s \in S(z, y)} \exp(\text{MLP}([h_{\text{START}(s)}; h_{\text{END}(s)}]))$$

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

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

RA-LLM Learning: Joint Training

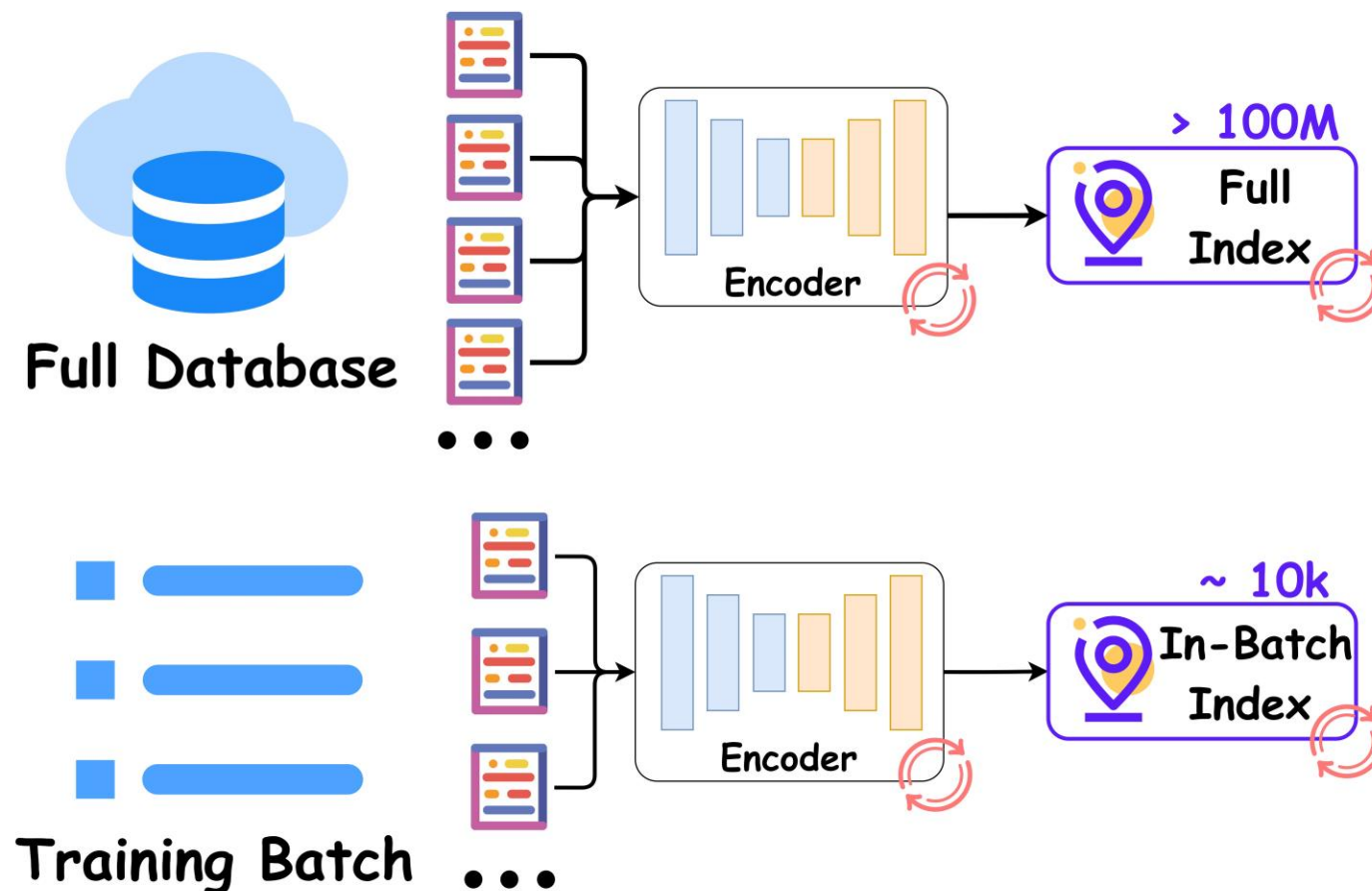
- REALM – Asynchronous Index Update



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

RA-LLM Learning : Joint Training

- TRIME – In-Batch Approximation

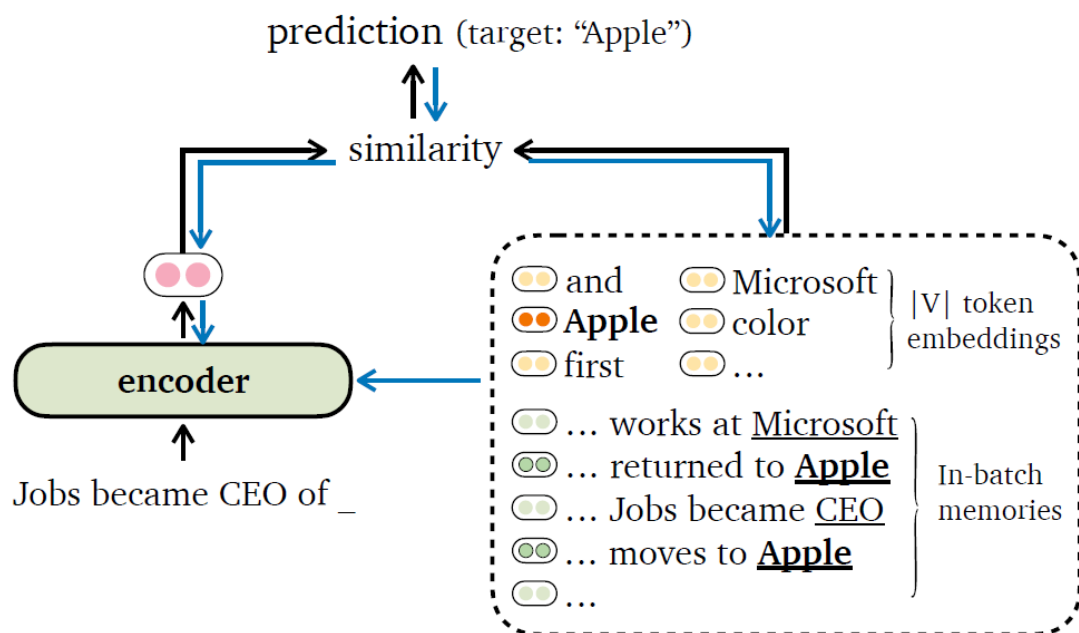


RA-LLM Learning : Joint Training

• 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:

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

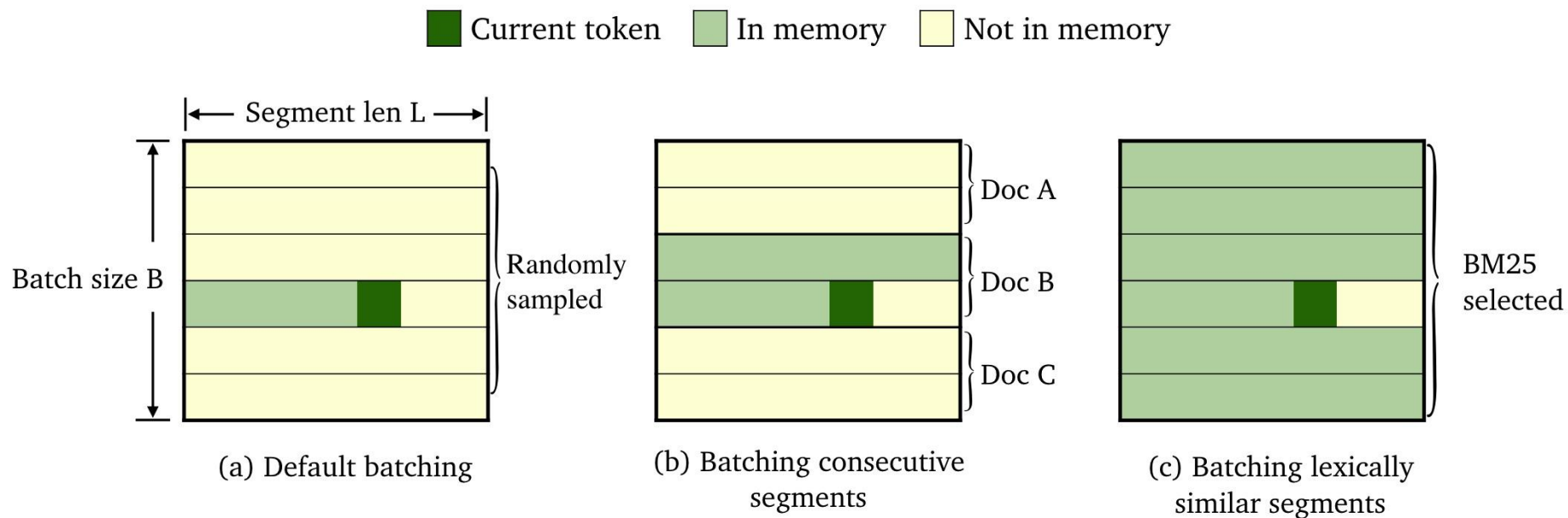
External Memory: $\mathcal{M}_{\text{ext}} = \{(c_j, x_j) \in \mathcal{D}\}$.

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-LLM Learning : Joint Training

- **TRIME Data Batching Strategy**



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

Tutorial Outline



41st IEEE International Conference
on Data Engineering
— HONG KONG SAR, CHINA | MAY 19 – 23, 2025 —



- ⦿ **Part 1: Introduction** of Retrieval Augmented Large Language Models (RA-LLMs) (Dr. Yujuan Ding)
- ⦿ **Part 2: Architecture** of RA-LLMs and **Main Modules** (Dr. Yujuan Ding)
- ⦿ **Part 3: Data Management** for RA-LLMs (Pangjing Wu)
- ⦿ **Part 4: Learning** Approach of RA-LLMs (Liangbo Ning)
- ⦿ **Part 5: Applications of RA-LLMs (Shijie Wang)**
- **Part 6: Challenges and Future Directions** of RA-LLMs (Liangbo Ning)

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



PART 4: Application of RA-LLMs

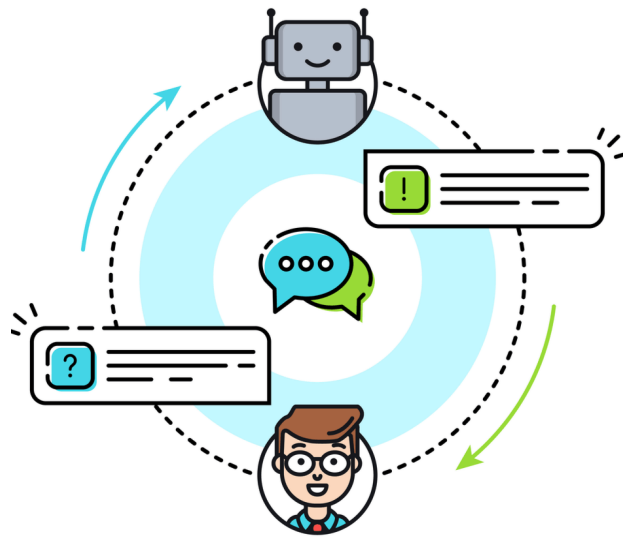


Presenter
Shijie Wang
HK PolyU

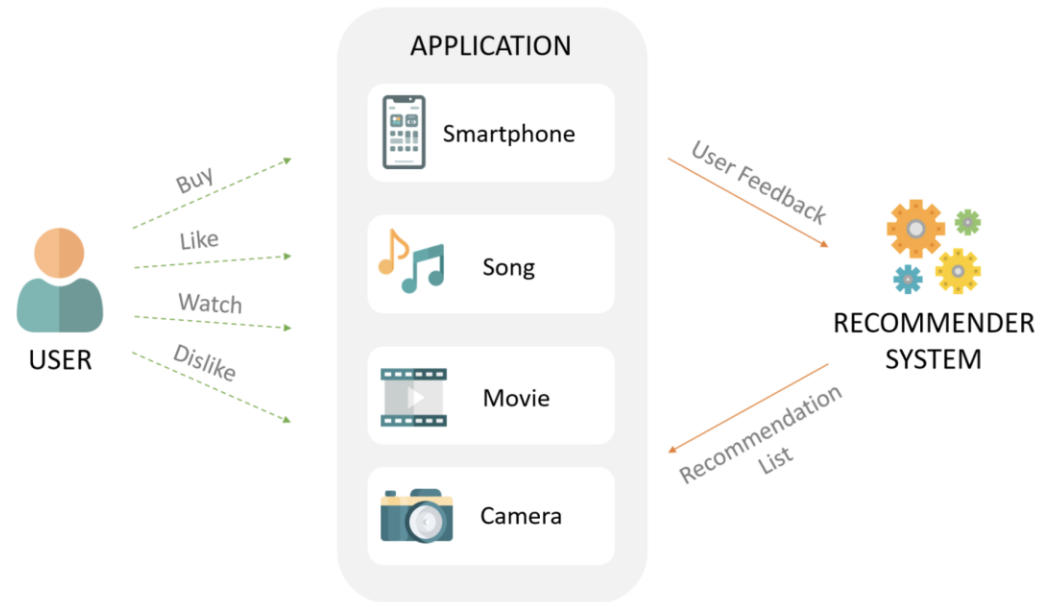
- **NLP applications**
- **Downstream tasks**
- **Domain-specific applications**

RA-LLM Applications

- Various applications



Chatbots



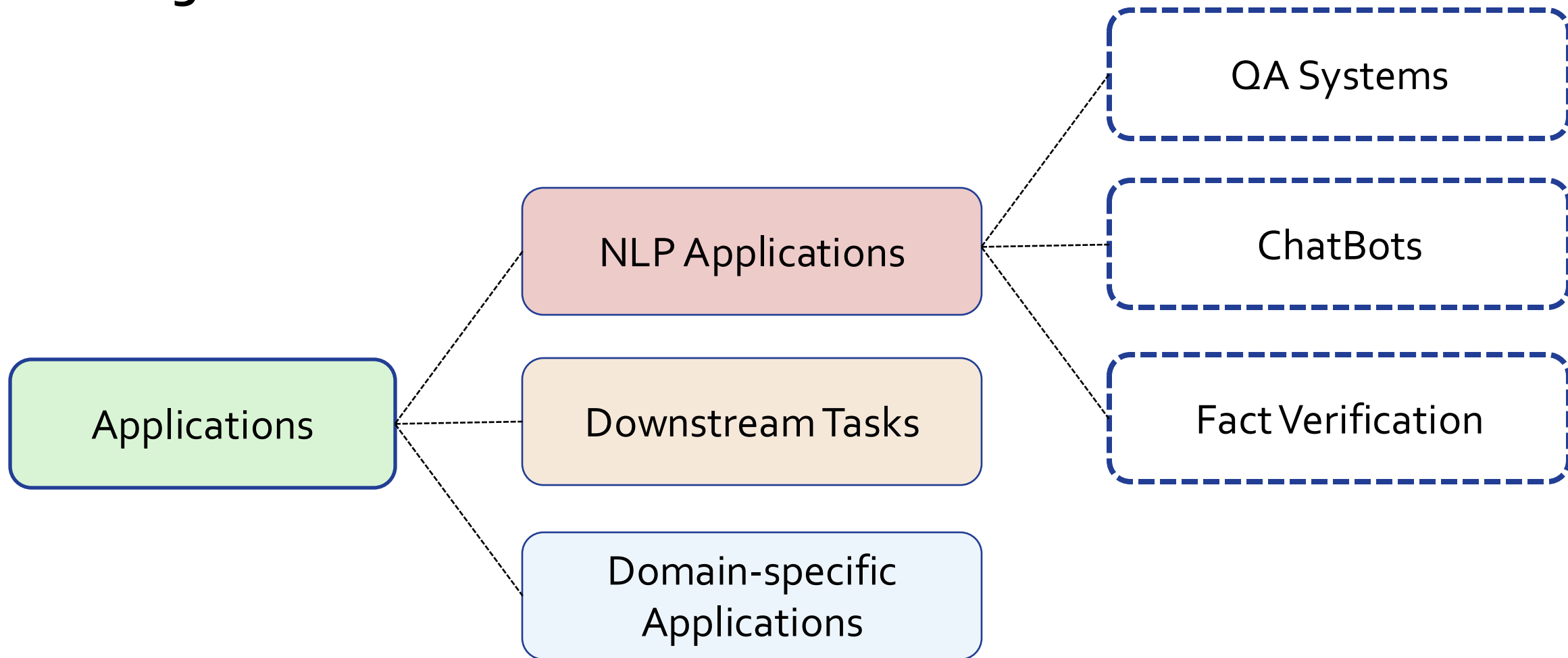
Recommendation



AI for Science

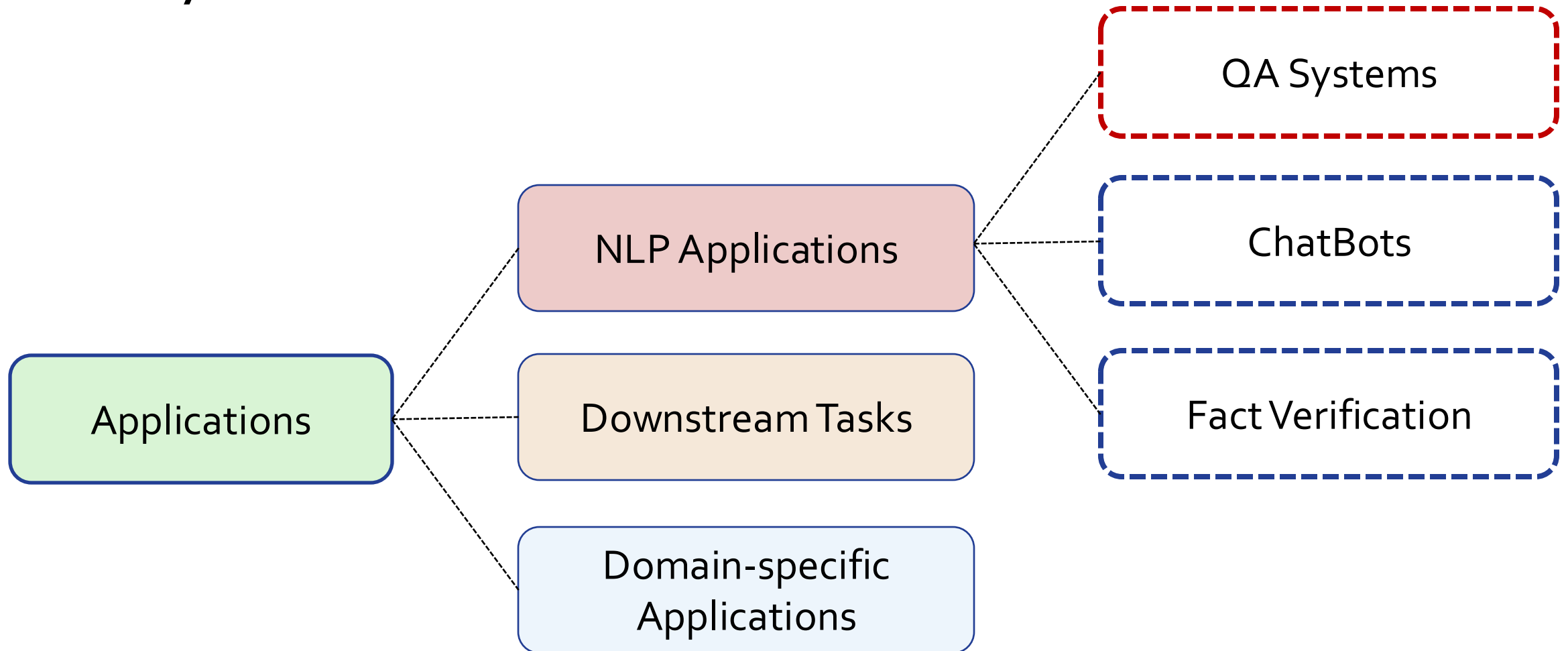
RA-LLM Applications: NLP Applications

- **Categories**



RA-LLM Applications: NLP Applications

- **QA Systems**



RA-LLM Applications: QA Systems

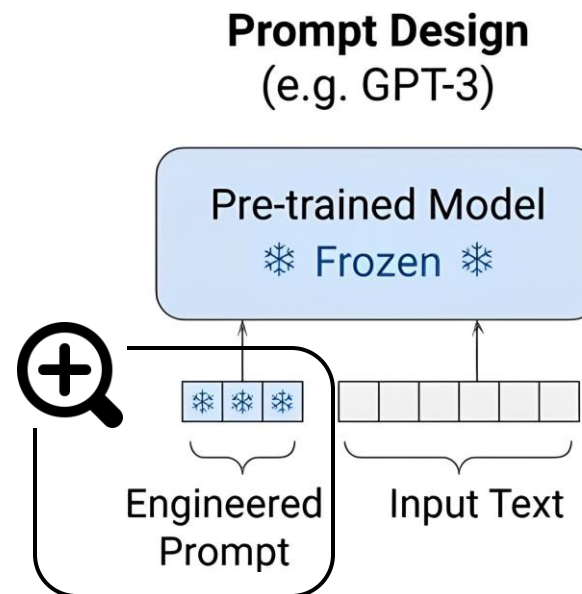
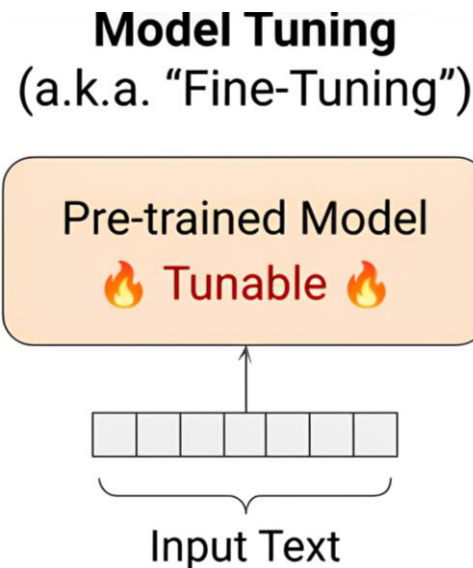
- **QA systems**

- Challenges:

- Open-domain QA
 - Domain-specific QA

- How to solve?

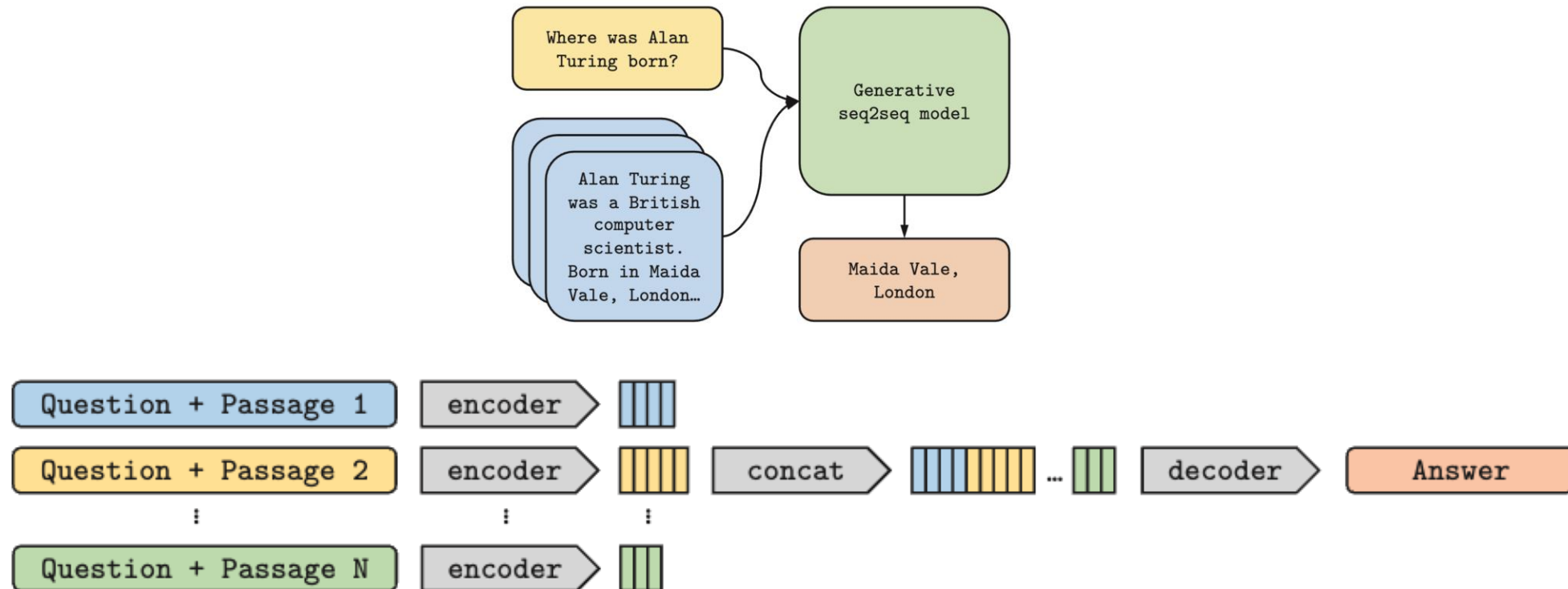
- Fine-tuning
 - Prompting



RA-LLM Applications: QA Systems

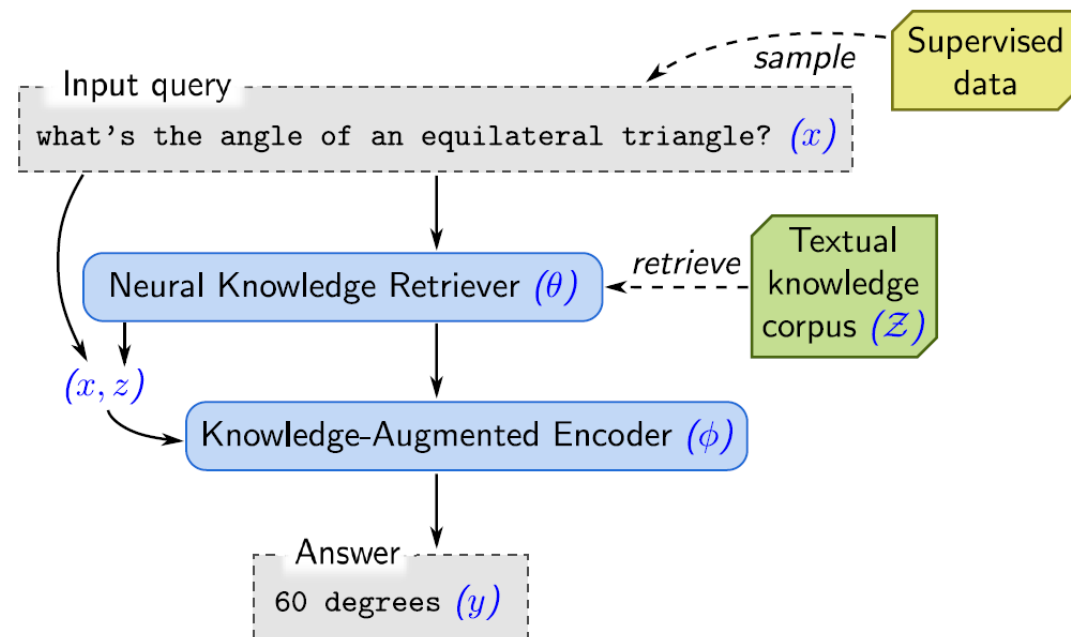
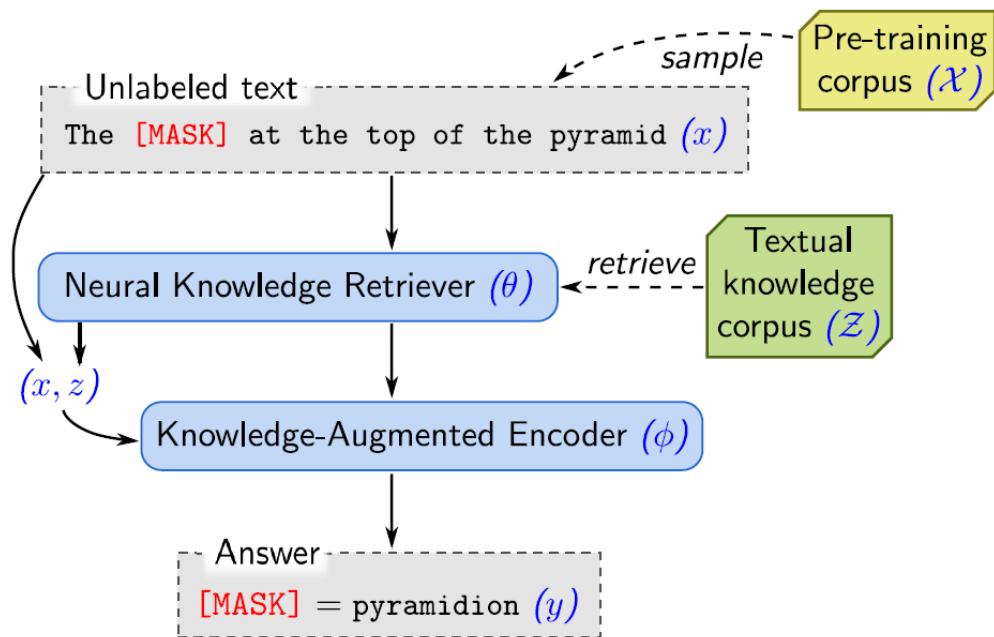
- **Retrieves for open-domain QA**

Retrieves support text passages from an external source of knowledge



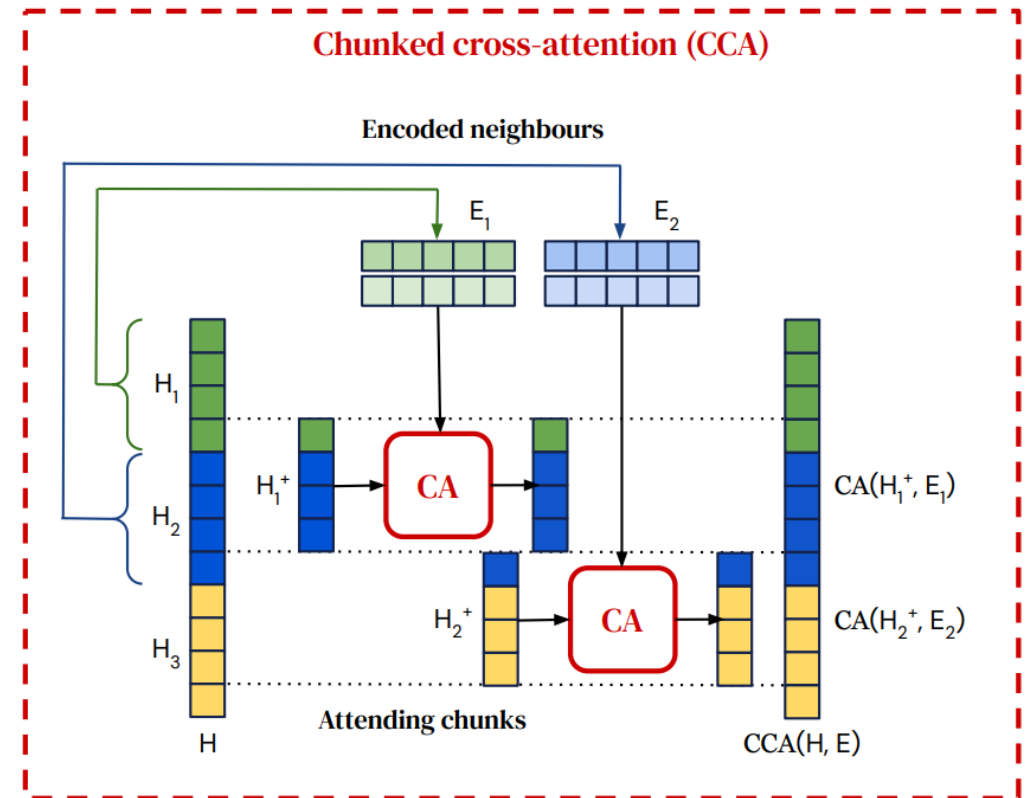
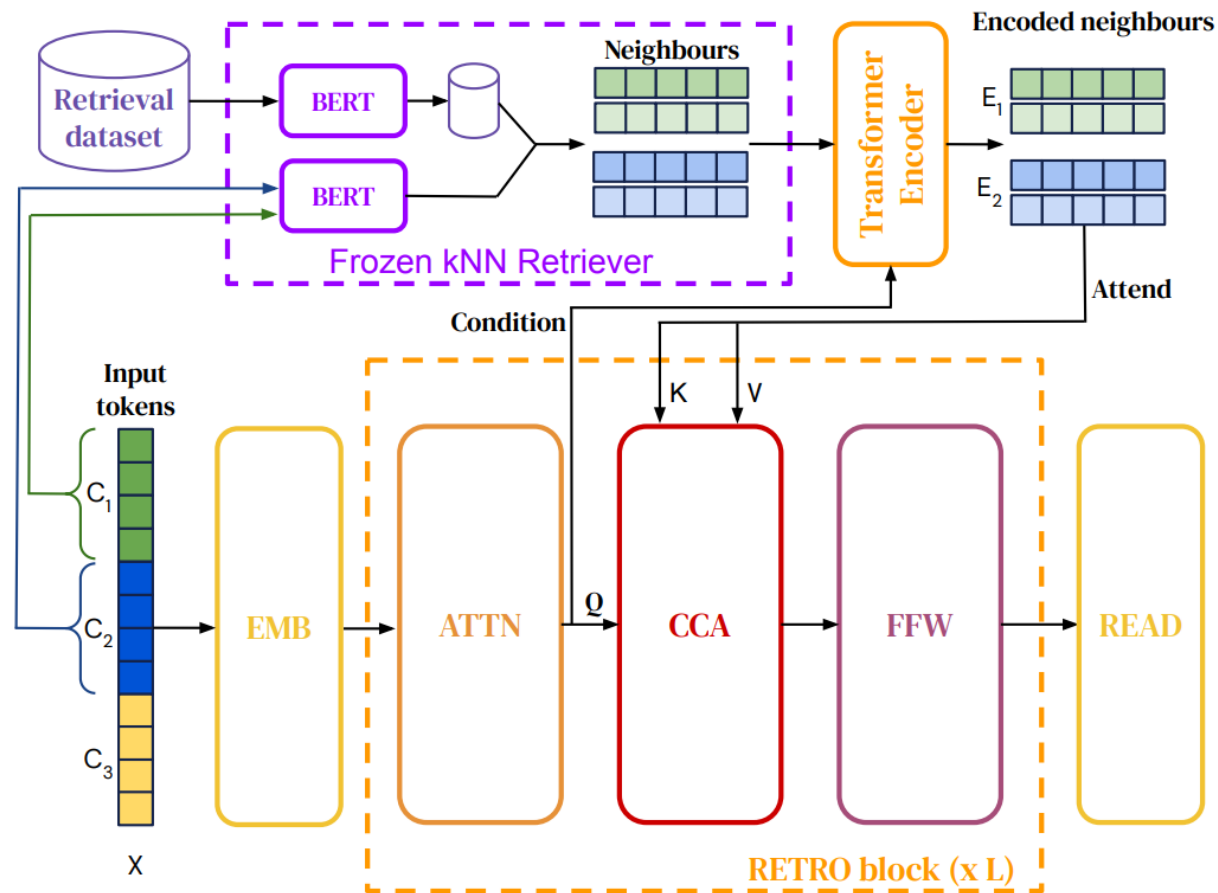
RA-LLM Applications: QA Systems

- **REALM**



RA-LLM Applications: QA Systems

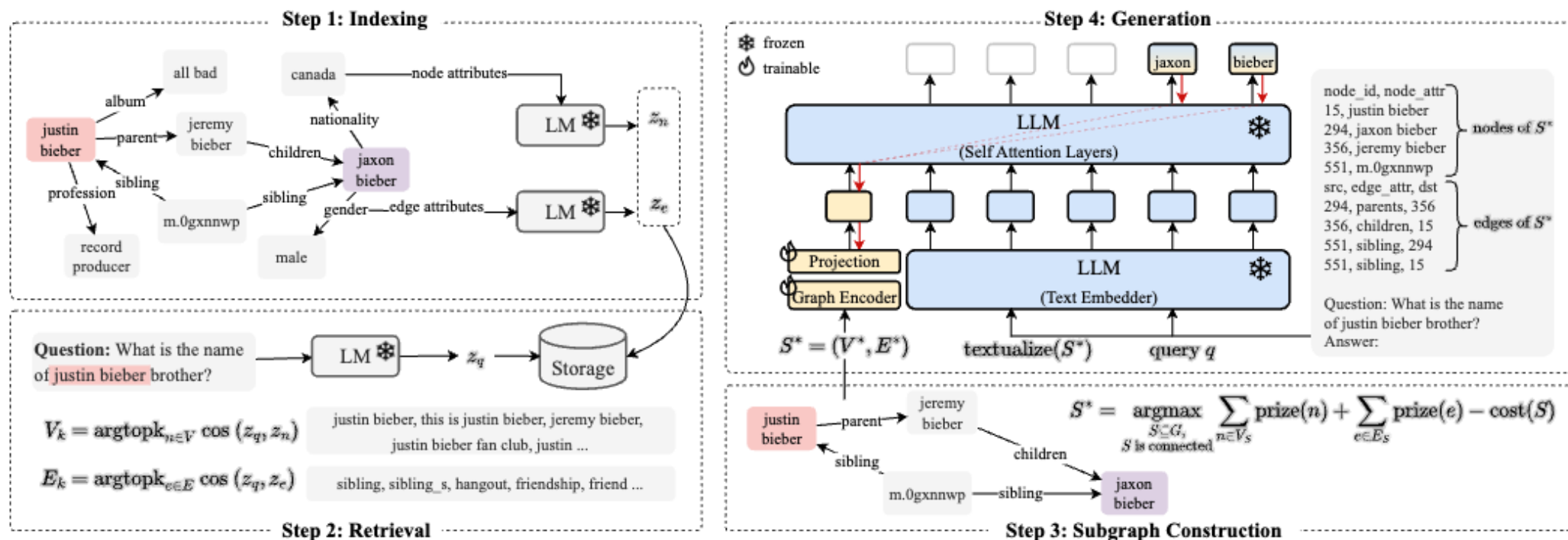
- RETRO (Retrieval-enhanced transformer)



RA-LLM Applications: QA Systems

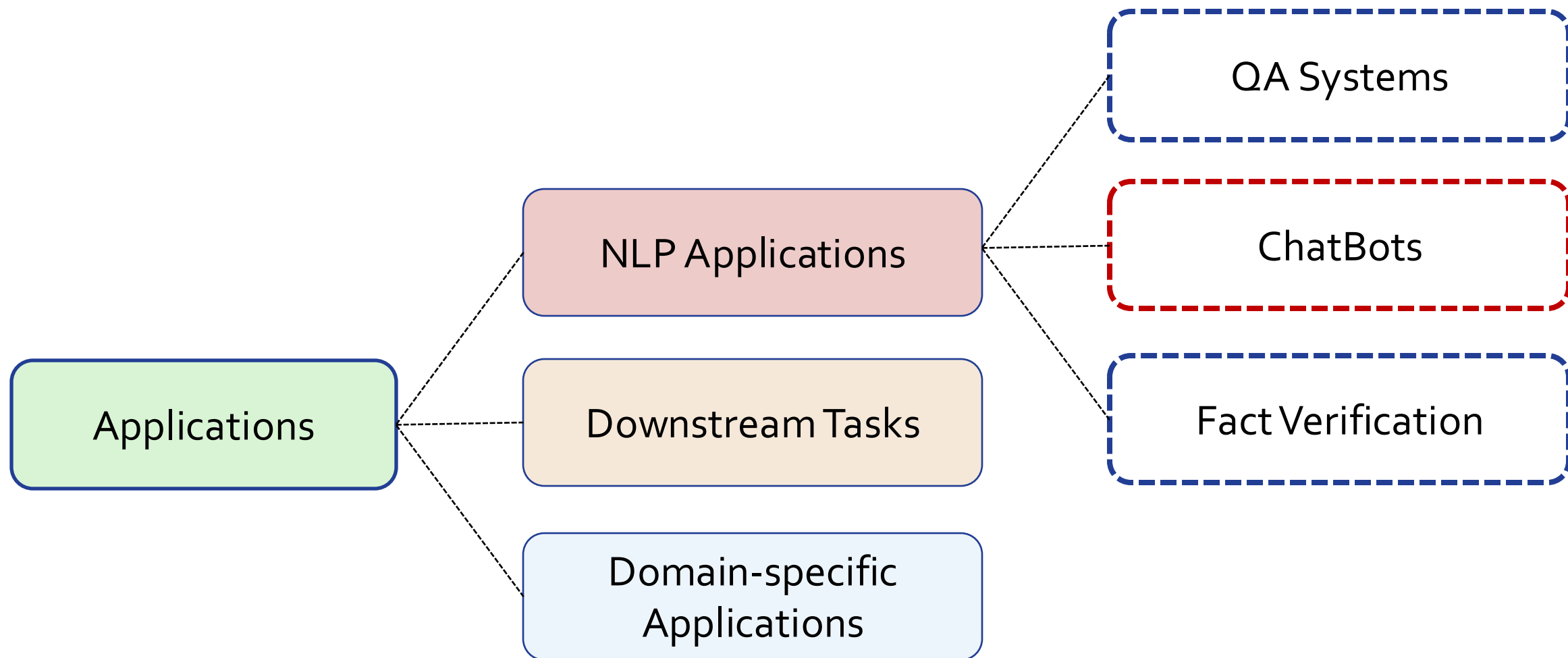
- **G-Retriever**

Retrieves from knowledge graph for question-answering



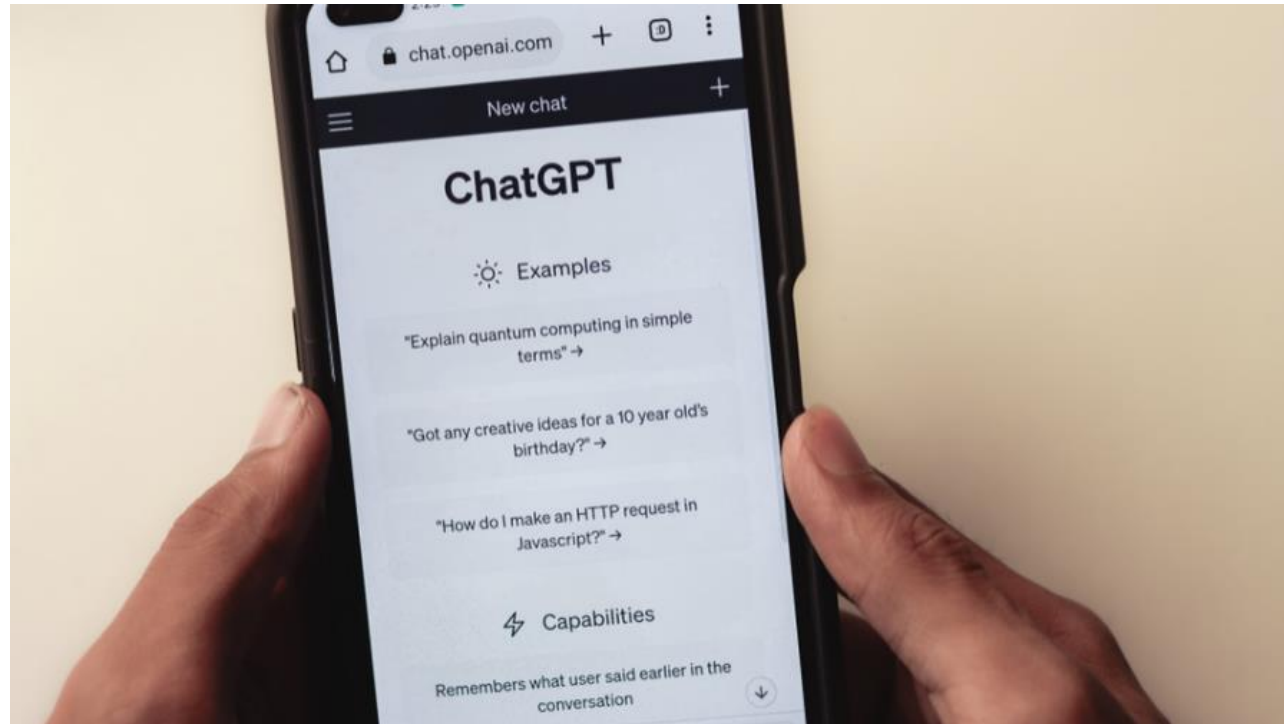
RA-LLM Applications: NLP Applications

- **ChatBots**




RA-LLM Applications: Chatbots

- **ChatBots**




RA-LLM Applications: Chatbots


- Knowledge-grounded model




"Consistently the best **omakase** in San Francisco." (27 Tips)






"Probably the best **sushi** in **San Francisco**." (2 Tips)



"... they were out of the **kaisui uni** by the time we ate, but the bafun uni is..." (2 Tips)

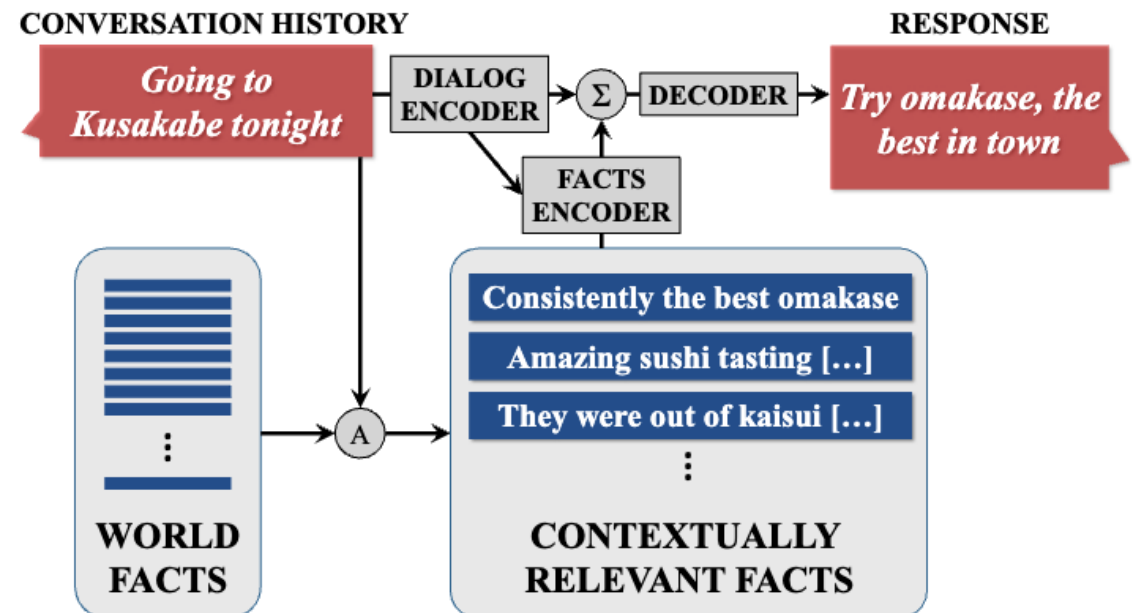


"Amazing sushi tasting from the chefs of **Sushi Ran**" (2 Tips)



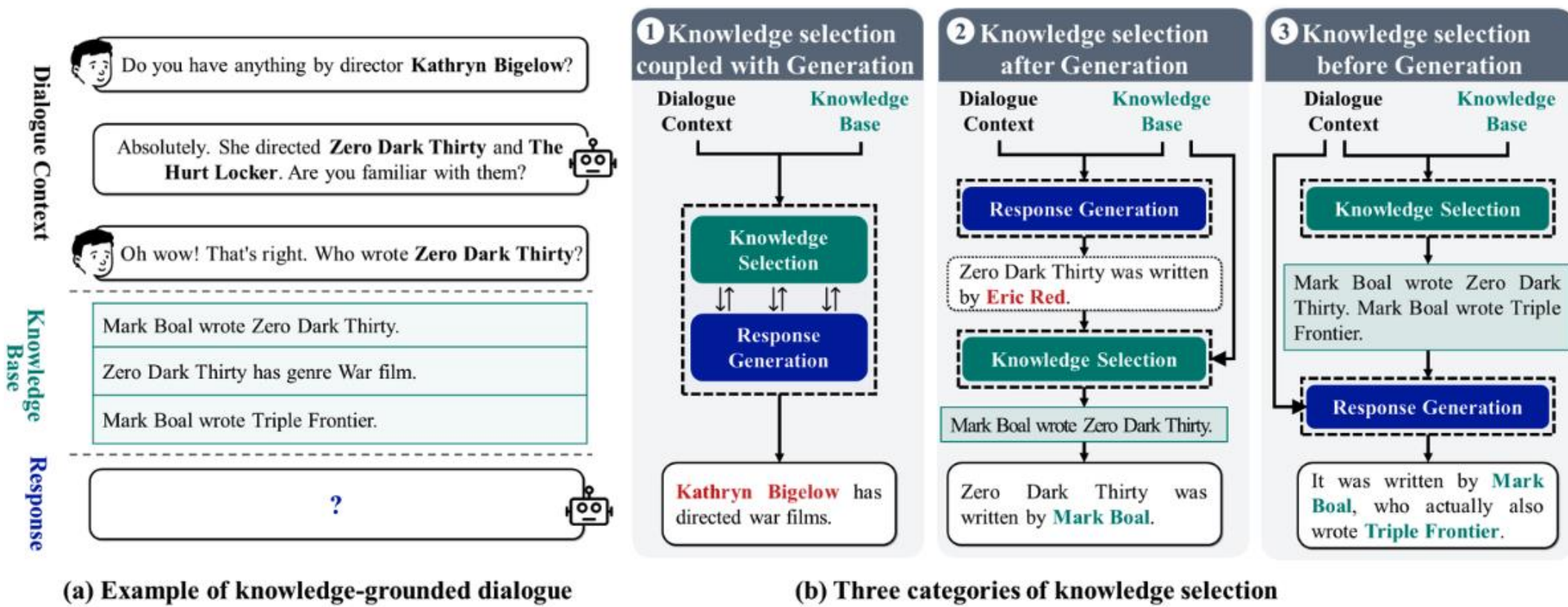
Kusakabe

User input: Going to Kusakabe tonight.
Neural model: Have a great time!
Human: You'll love it! Try omasake, the best in town.



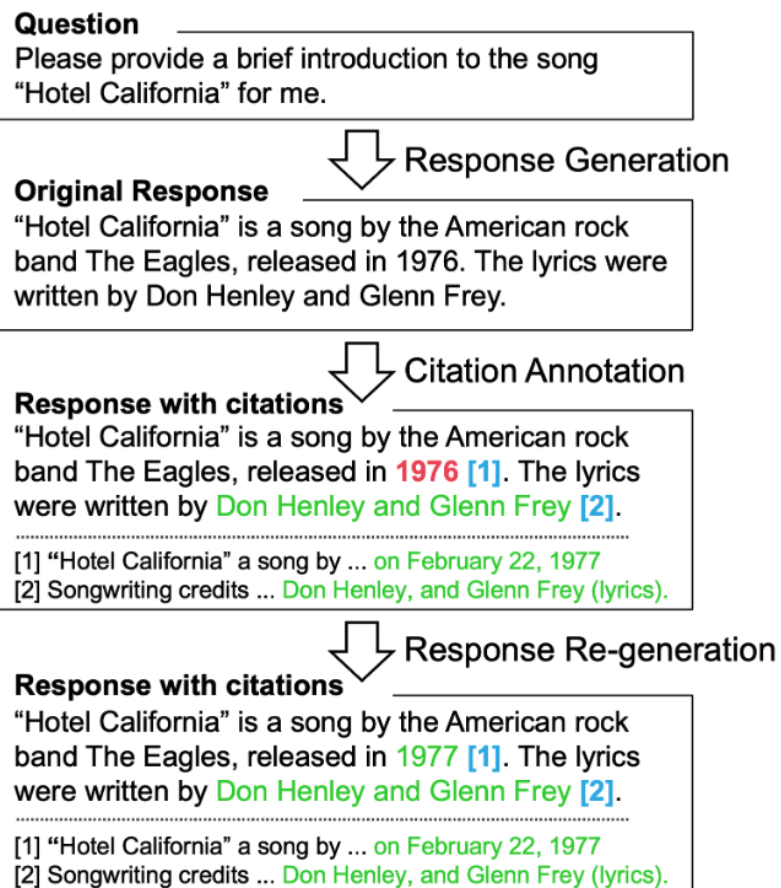
RA-LLM Applications: Chatbots

- GATE



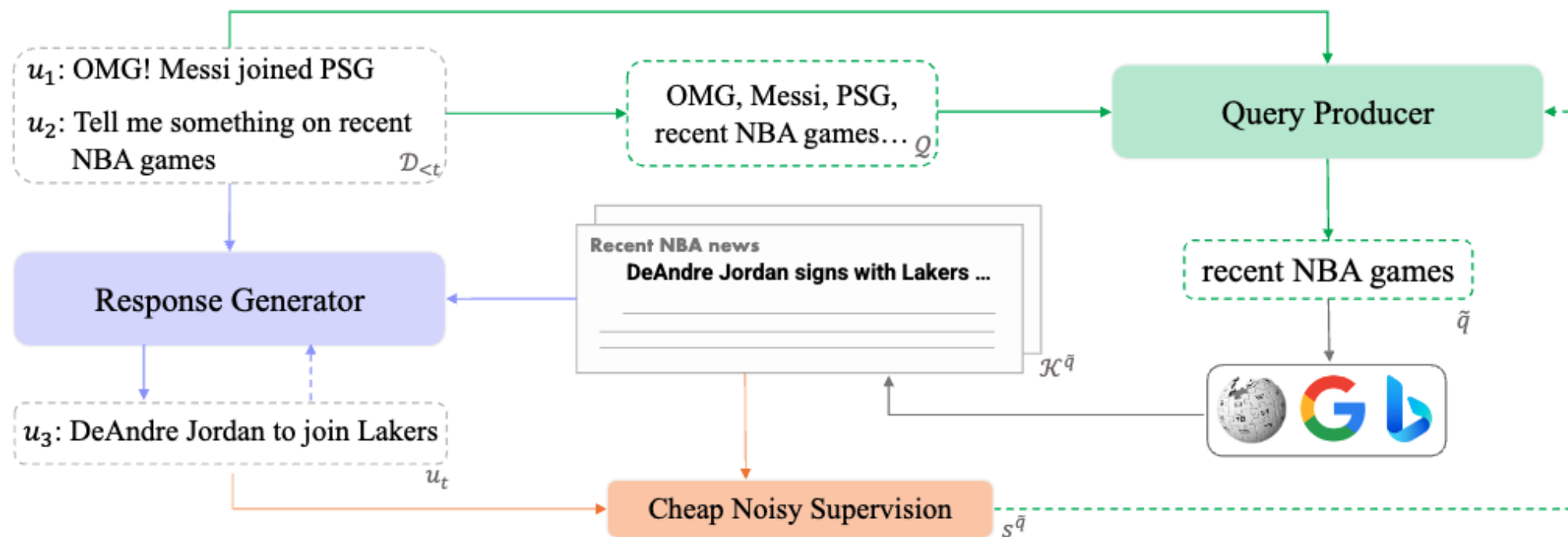
RA-LLM Applications: Chatbots

- CEG



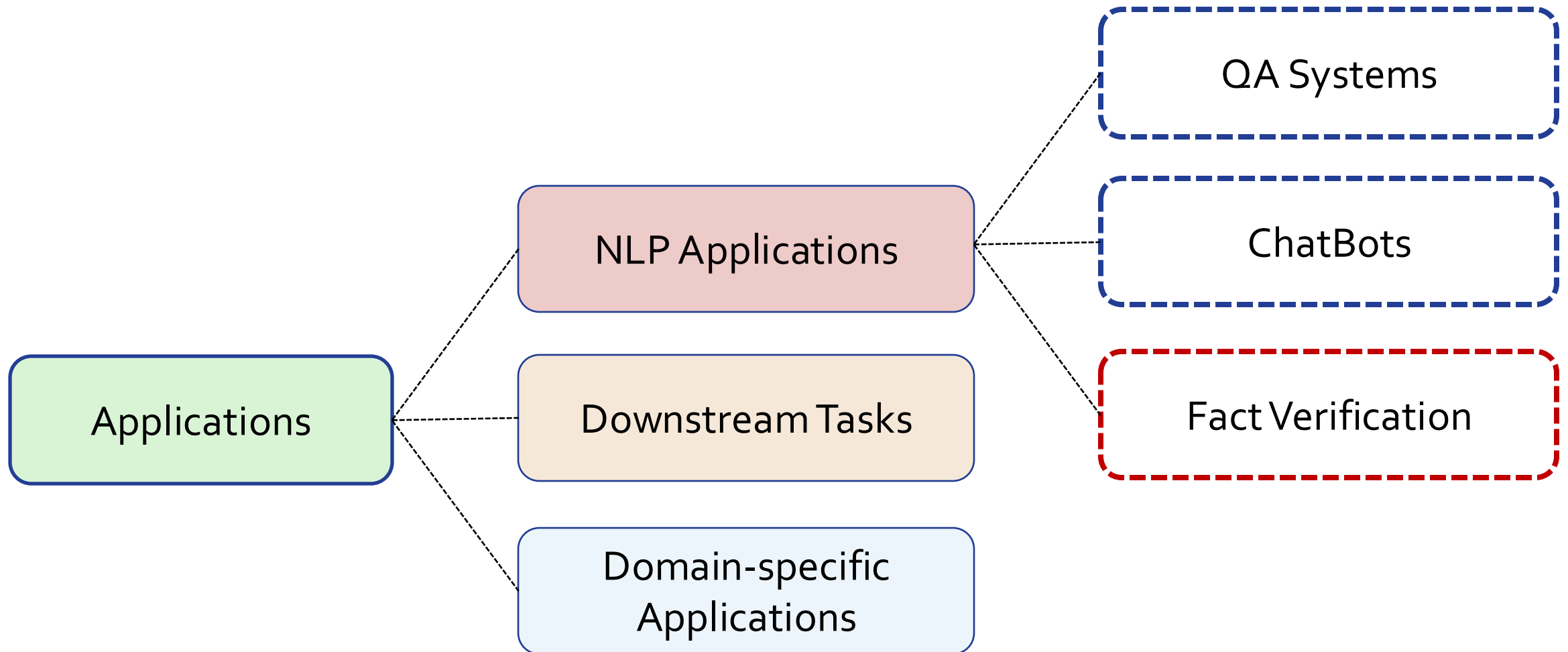
RA-LLM Applications: Chatbots

- Search-engine-augmented chatbots



RA-LLM Applications: NLP Applications

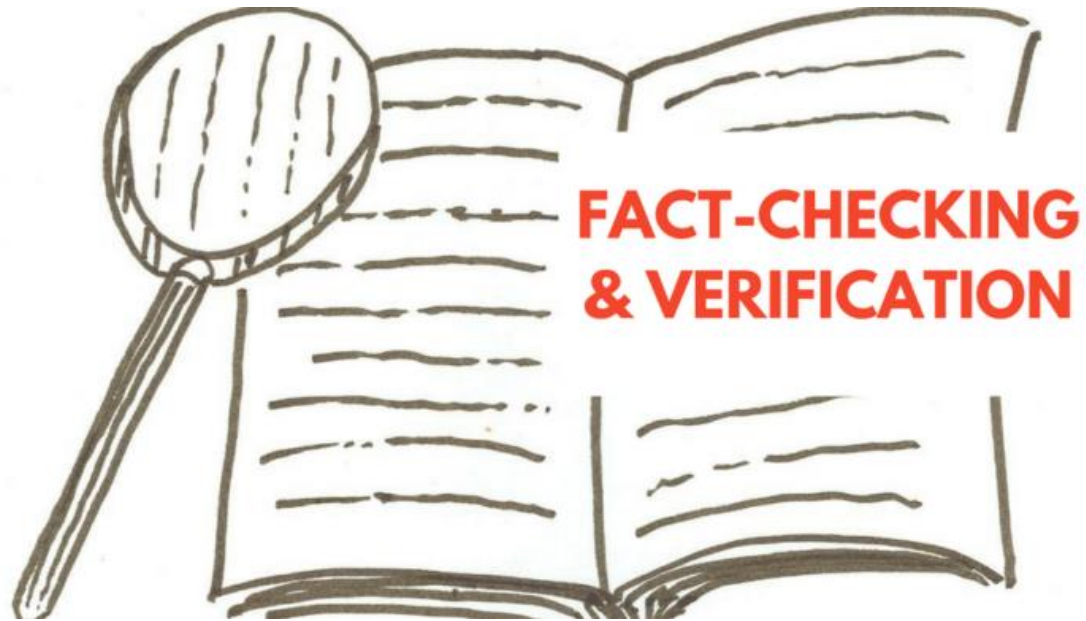
- **Fact verification**



RA-LLM Applications: Fact Verification

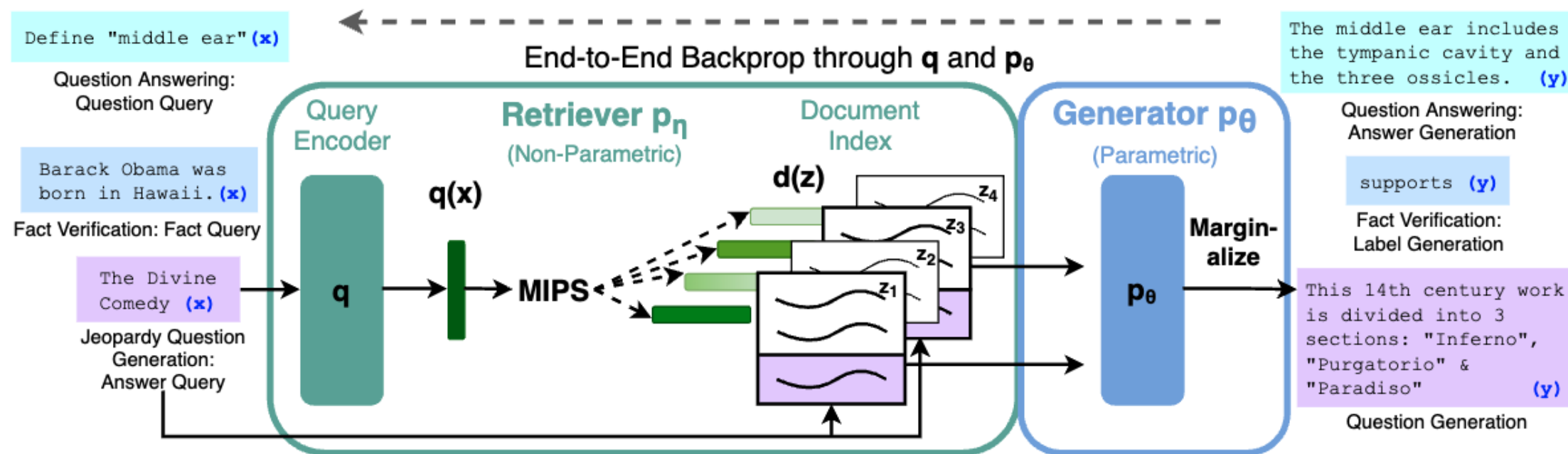
- **Fact verification**

Fact Verification is a critical task in verifying the accuracy and reliability of information



RA-LLM Applications: Fact Verification

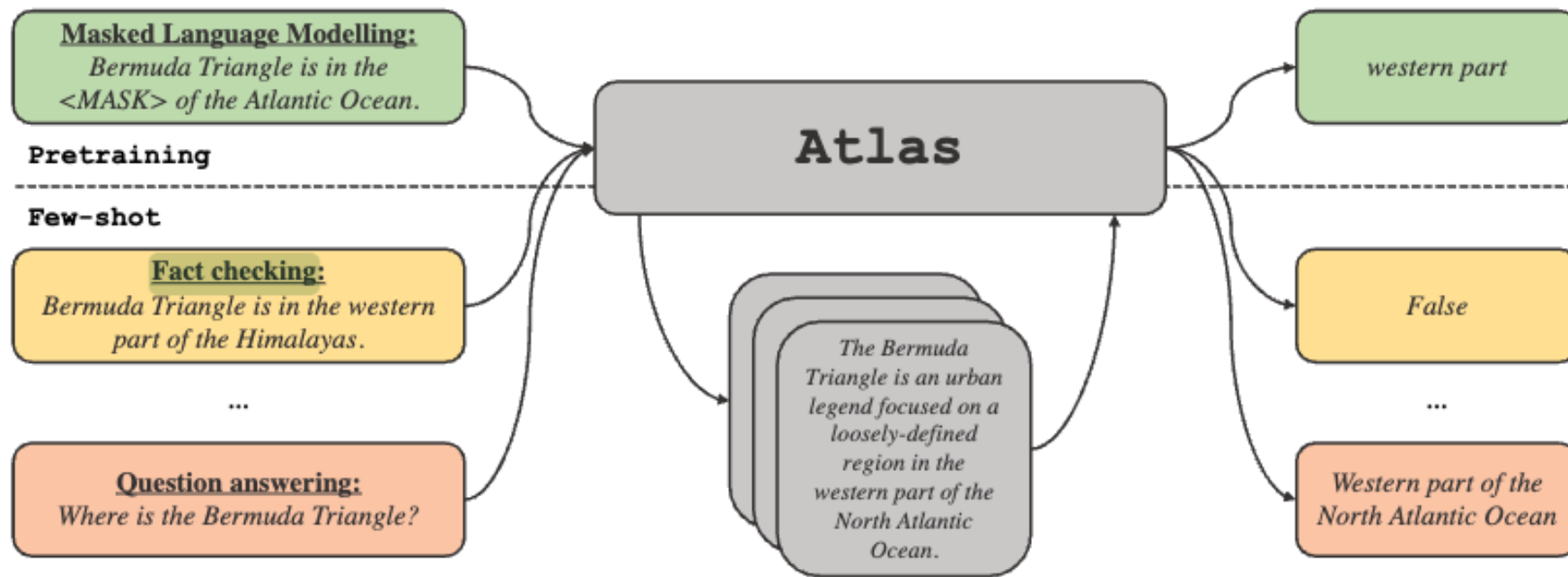
- Fact verification



RA-LLM Applications: Fact Verification

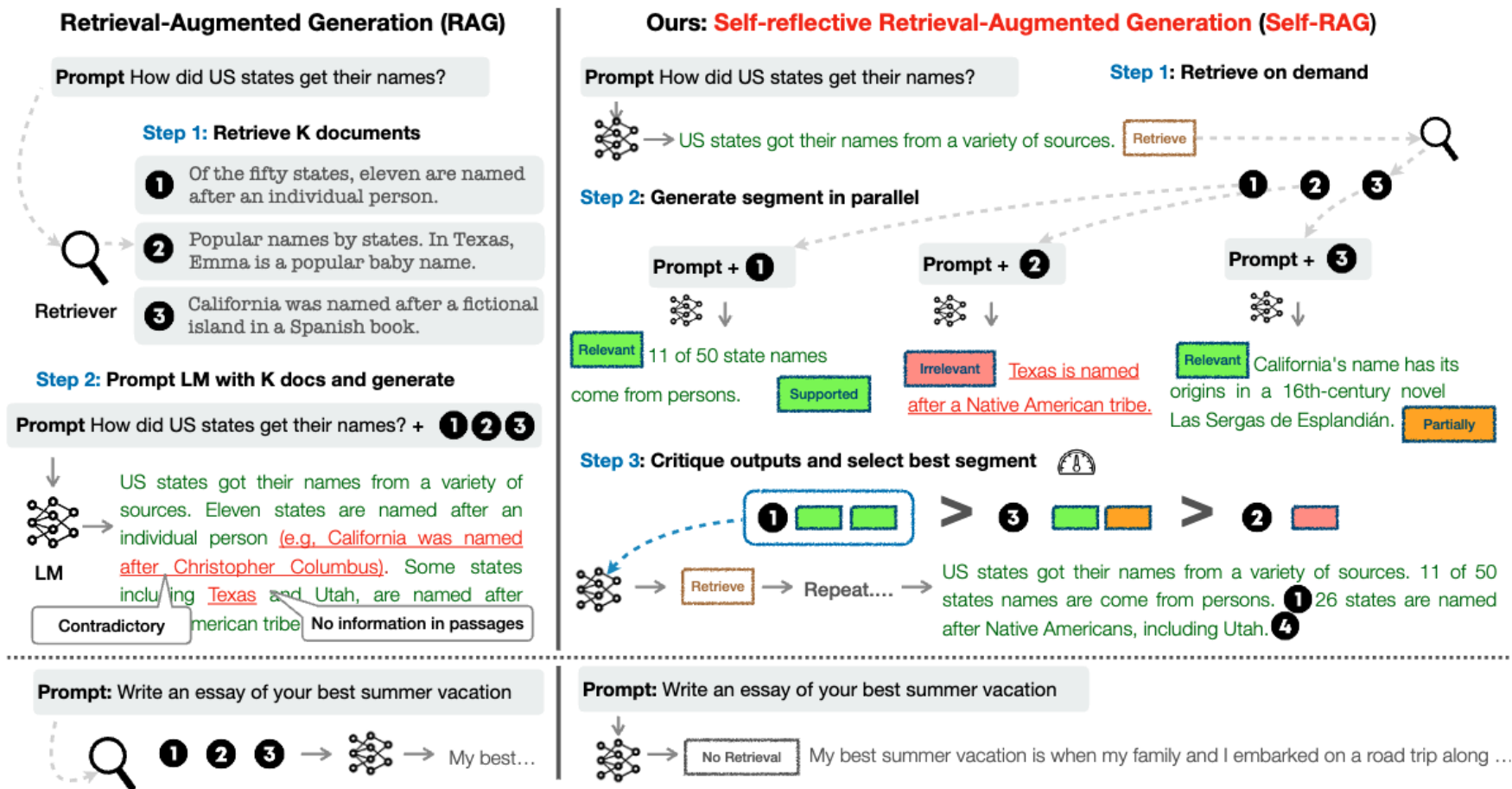
- **Fact verification**

- Fact verification is usually together with other NLP tasks (such as Q & A)
- ATLAS:



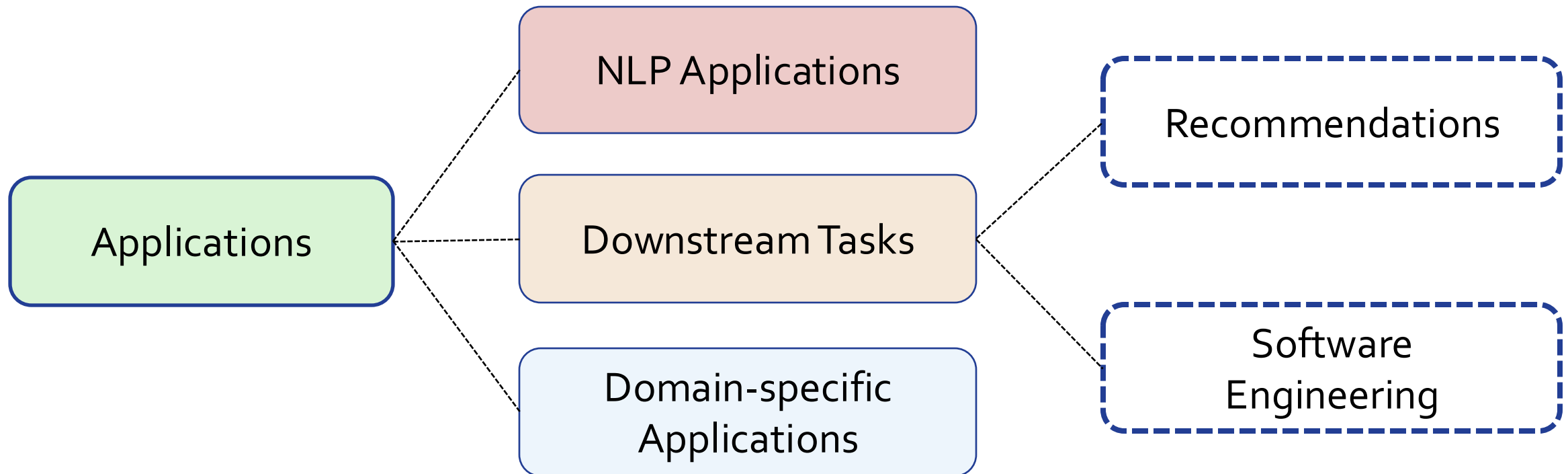
RA-LLM Applications: Fact Verification

- Self-RAG



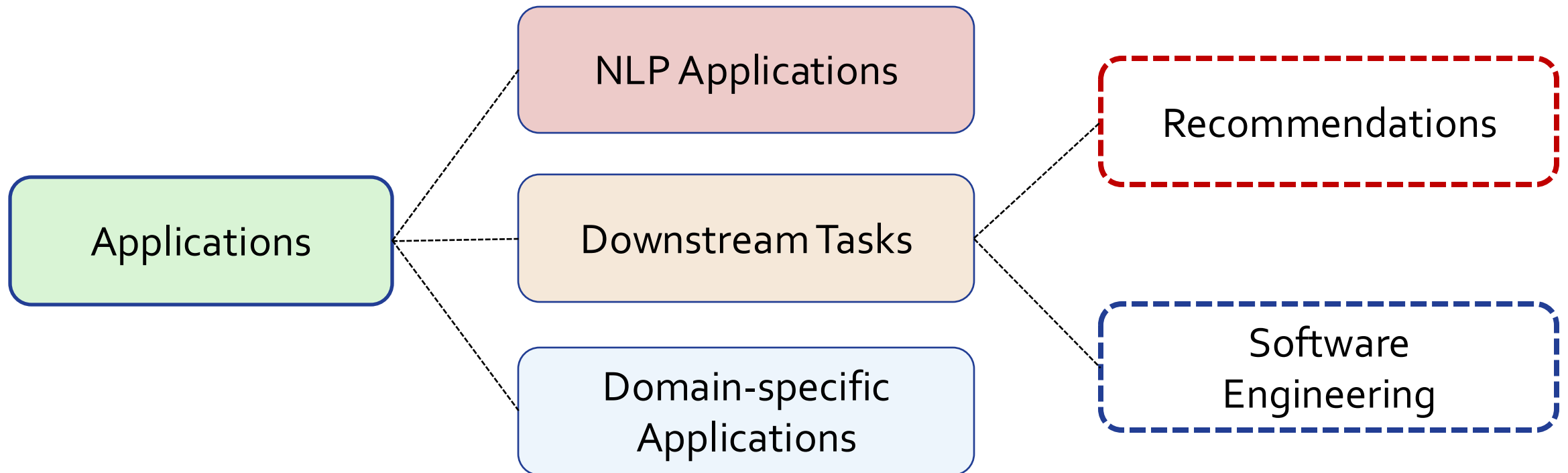
RA-LLM Applications: Downstream Tasks

- **Downstream tasks**



RA-LLM Applications: Recommendations

- **Recommendations**



RA-LLM Applications: Recommendations

- **Recommendations**

Recommendation has been widely applied in online services

YouTube

TikTok

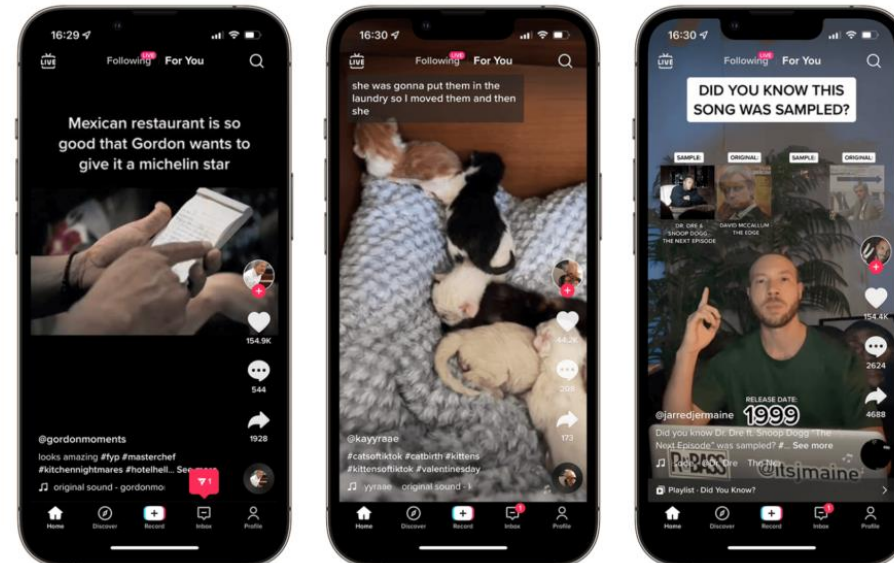


News/Video/Image Recommendation

TikTok's recommendation algorithm

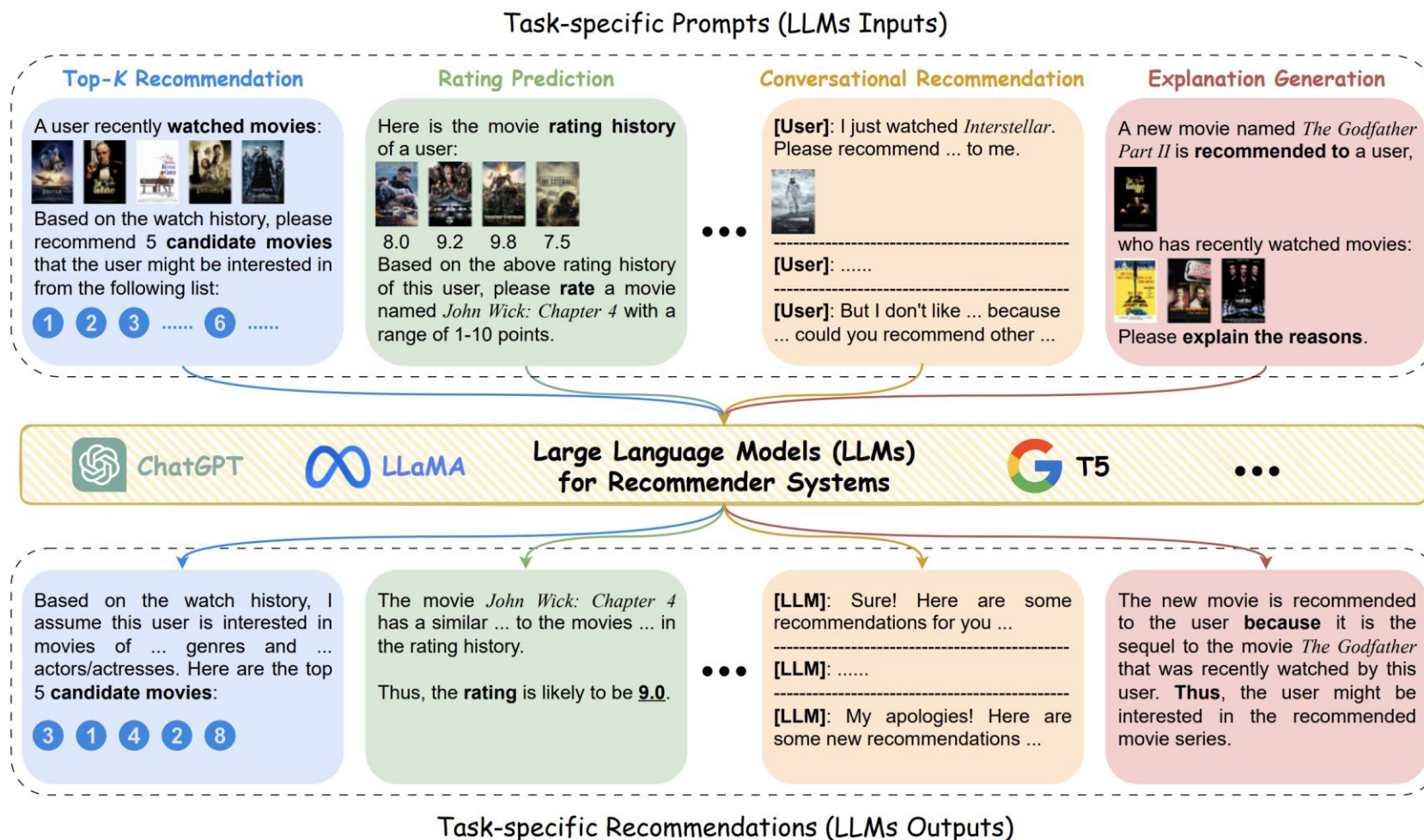
Top 10 Global Breakthrough
Technologies in 2021

MIT
Technology
Review



RA-LLM Applications: Recommendations

- LLMs in recommendations



RA-LLM Applications: Recommendations

- Conventional item-based LLM reasoning process



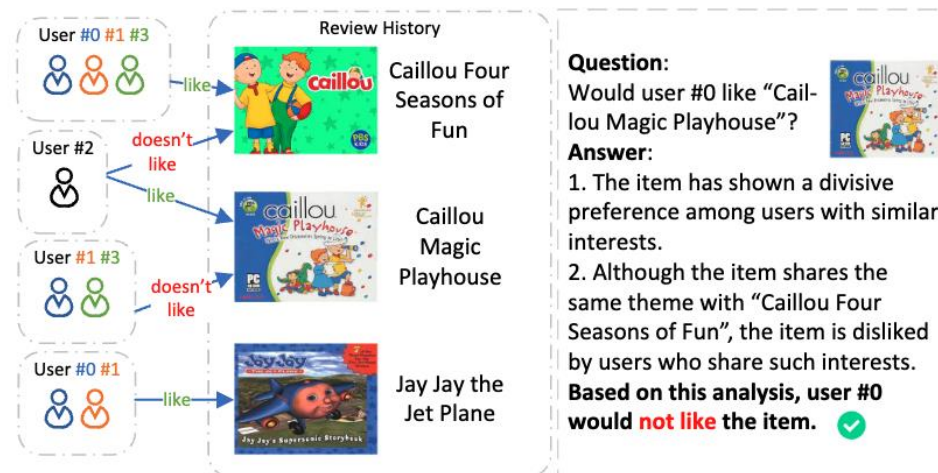
(a) Conventional item-based [16, 42] LLM reasoning process.

RA-LLM Applications: Recommendations

- Collaborative retrieval augmented LLM reasoning process



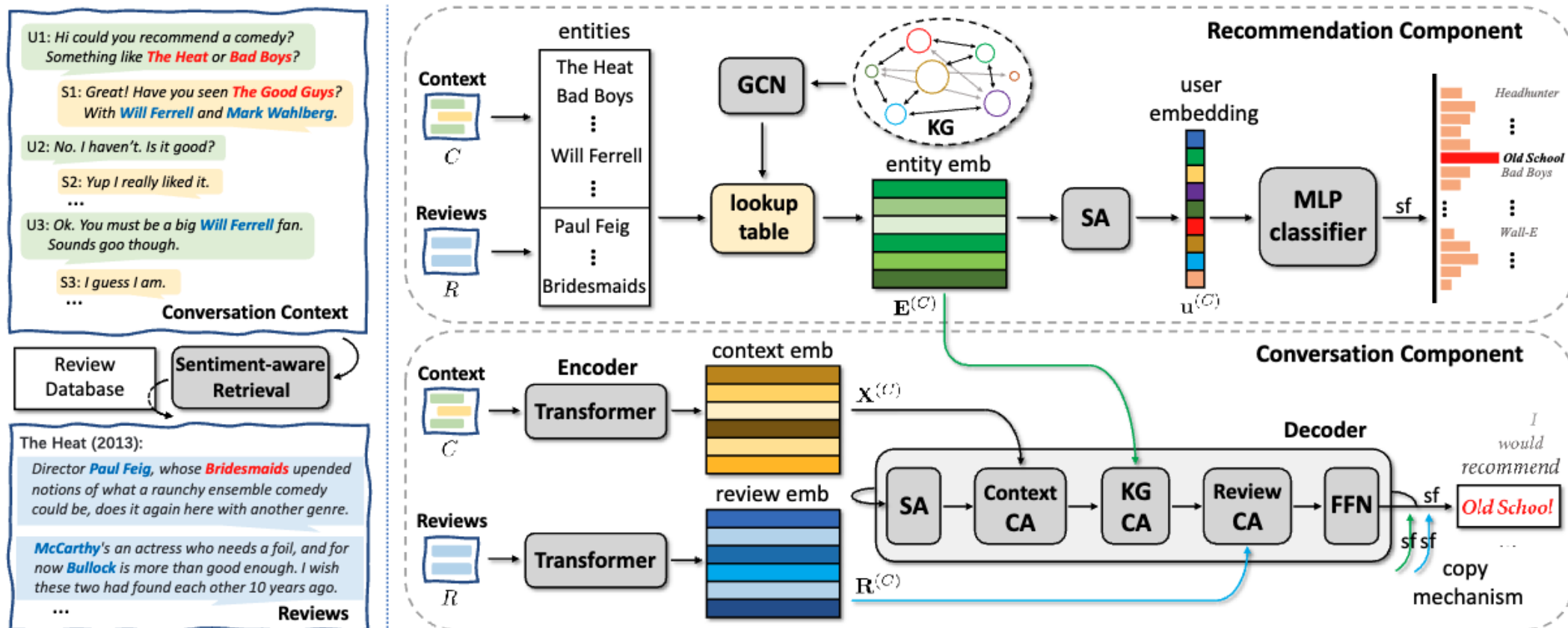
(a) Conventional item-based [16, 42] LLM reasoning process.



(b) Collaborative Retrieval Augmented LLM reasoning process.

RA-LLM Applications: Recommendations

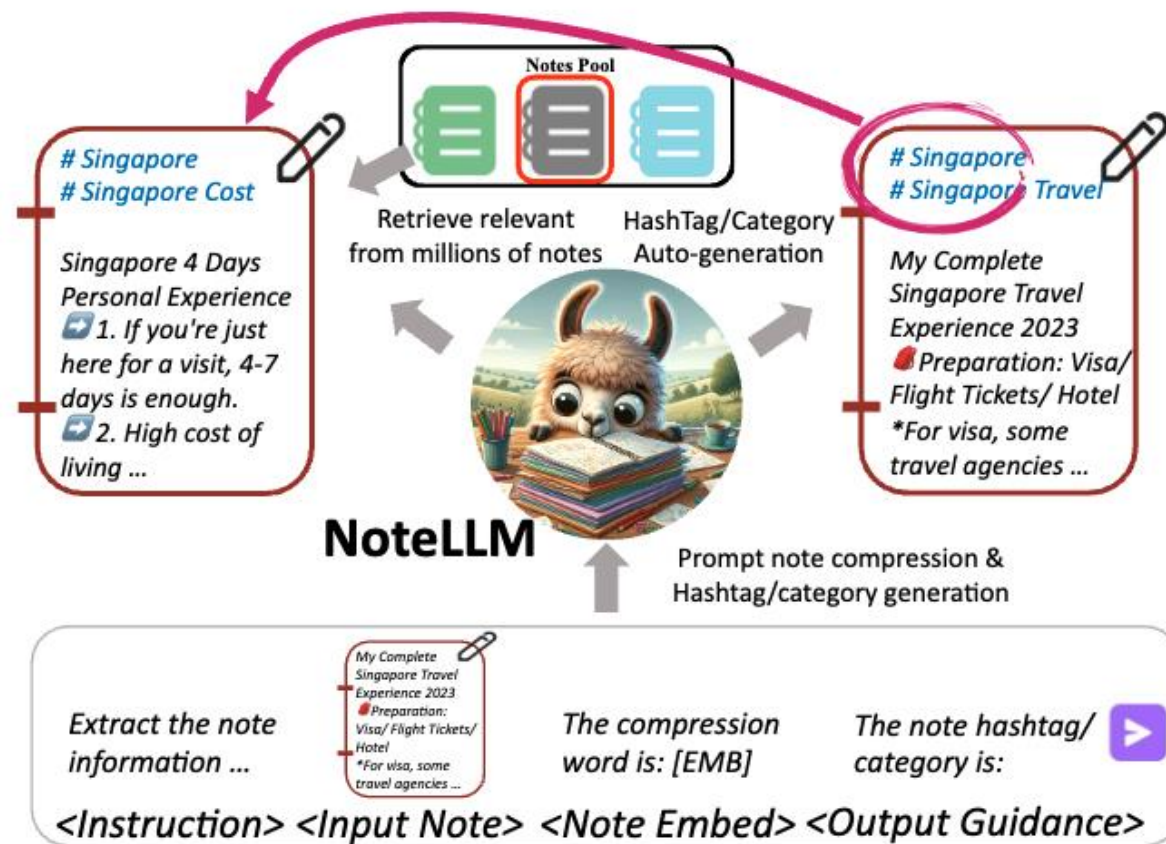
- Retrieval from the reviews



RA-LLM Applications: Recommendations

- Retrieval from the notes

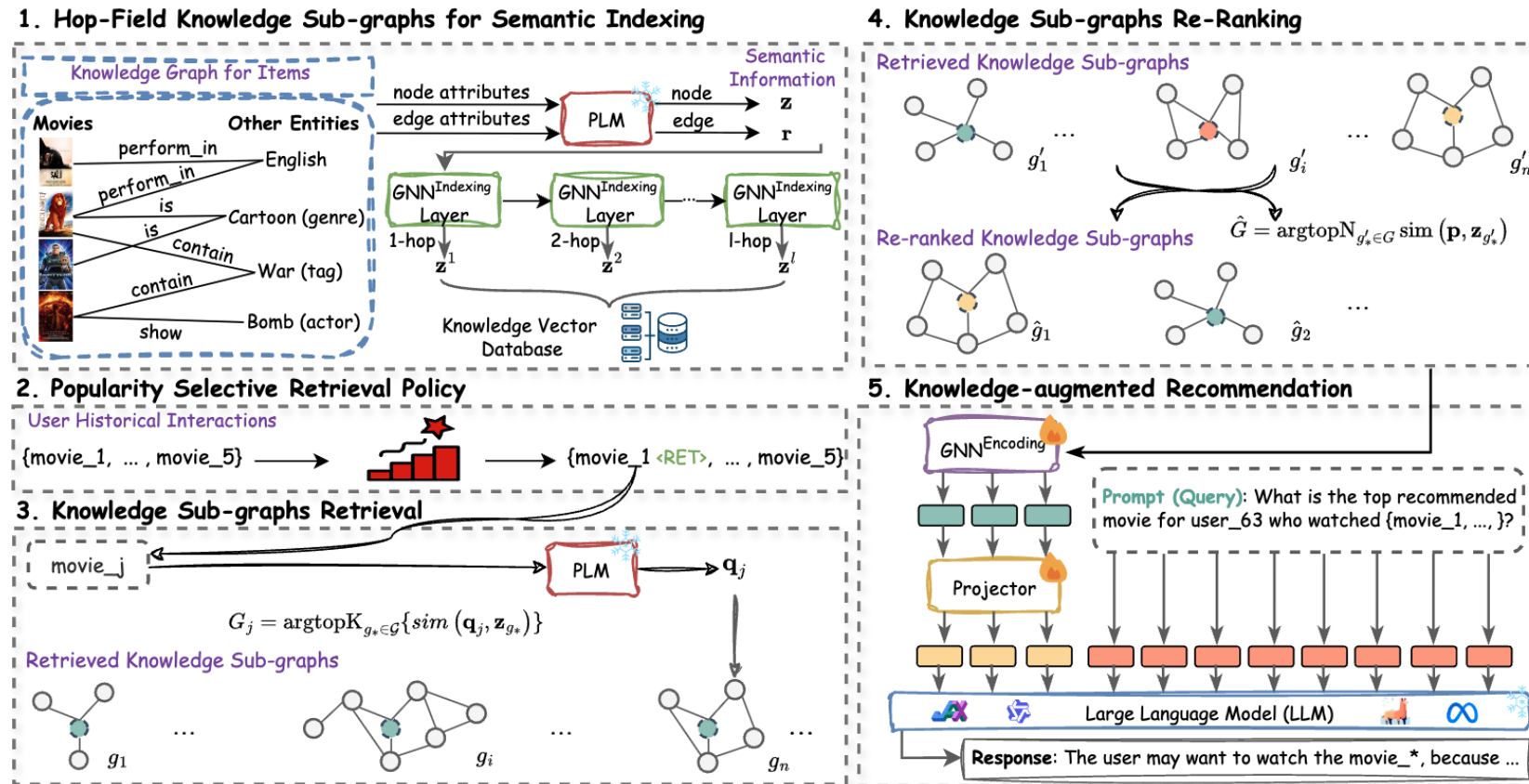
NoteLLM:



RA-LLM Applications: Recommendations

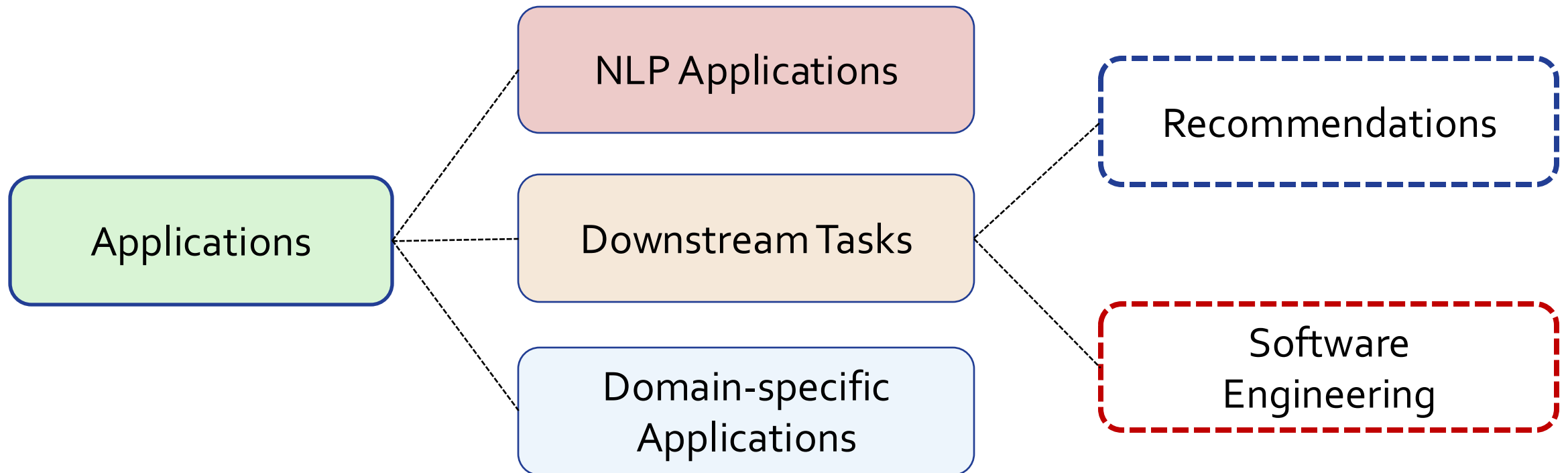
- Retrieval from knowledge graph

K-RagRec:



RA-LLM Applications: Software Engineering

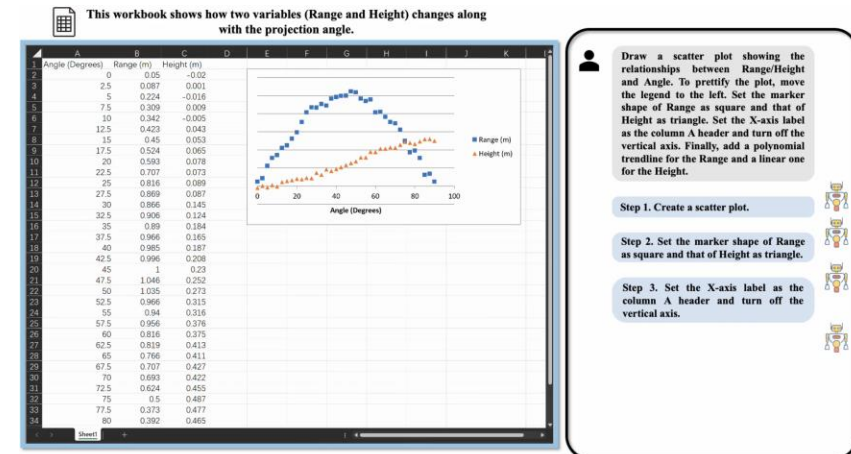
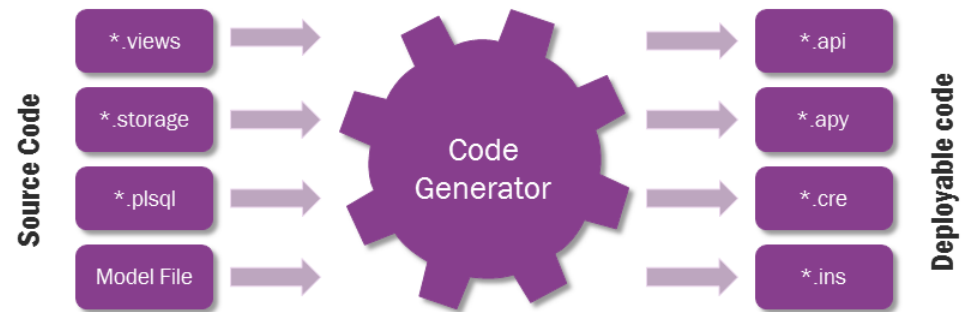
- **Software engineering**



RA-LLM Applications: Software Engineering

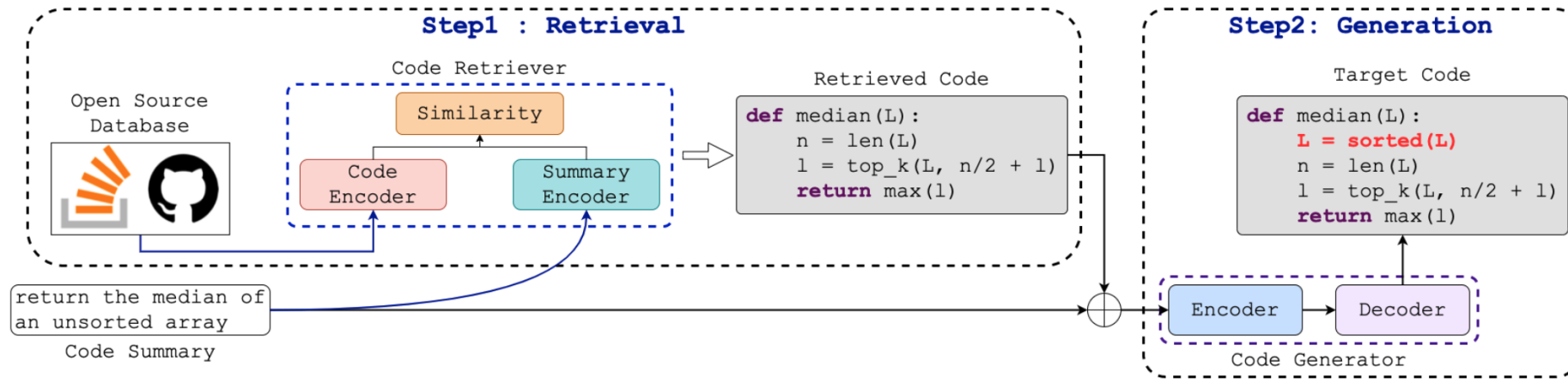
- **Software engineering:**

- Code generation
- Program repair
- Table processing
- ...



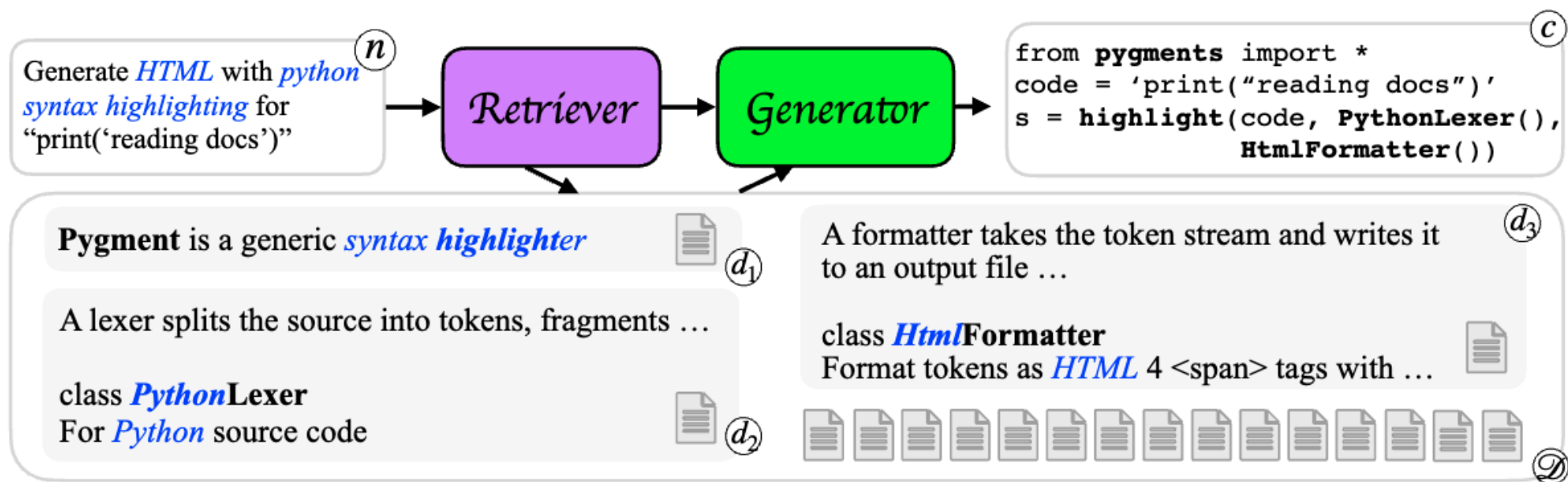
RA-LLM Applications: Software Engineering

- **Code generation:**



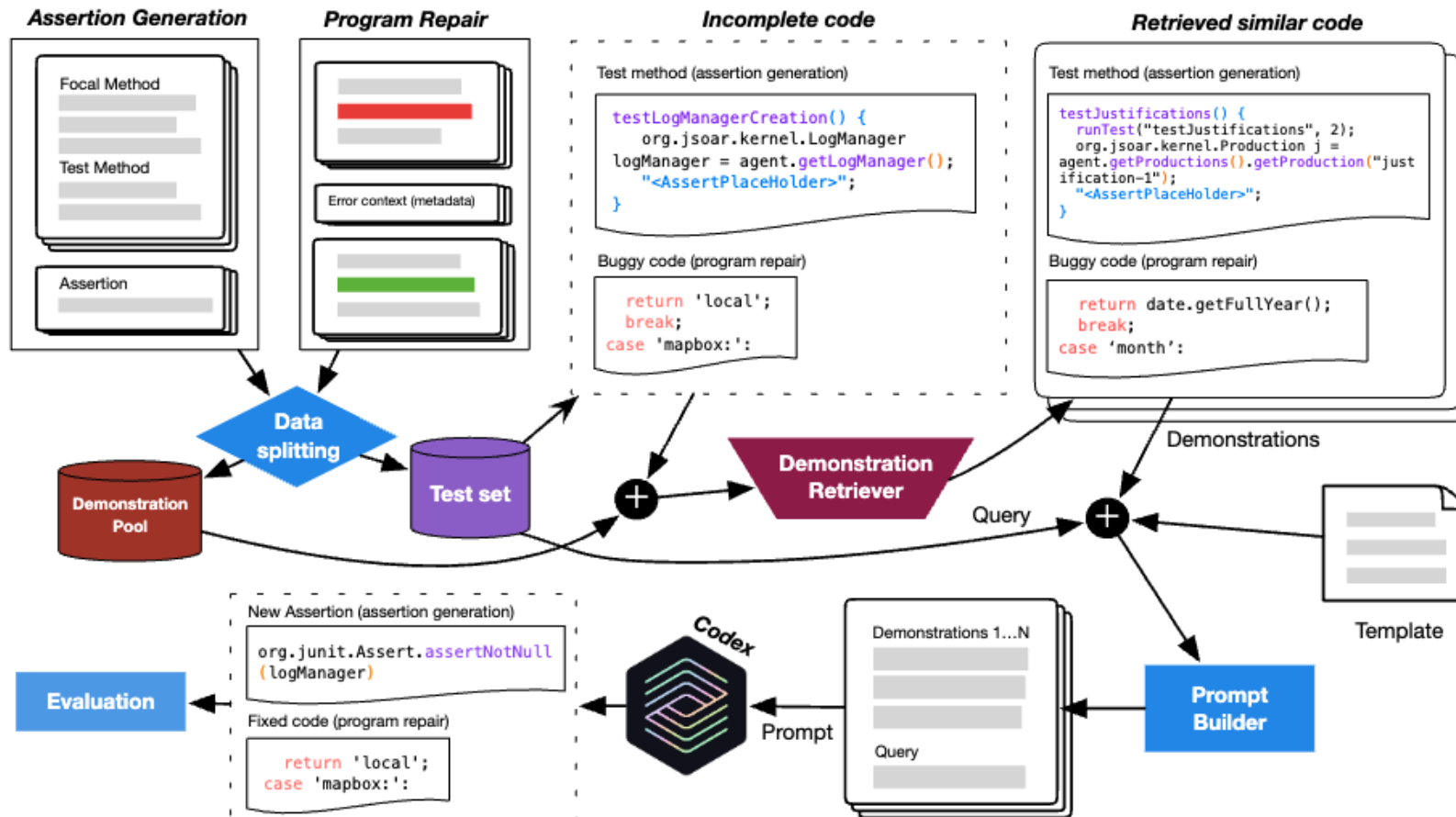
RA-LLM Applications: Software Engineering

- Code generation:



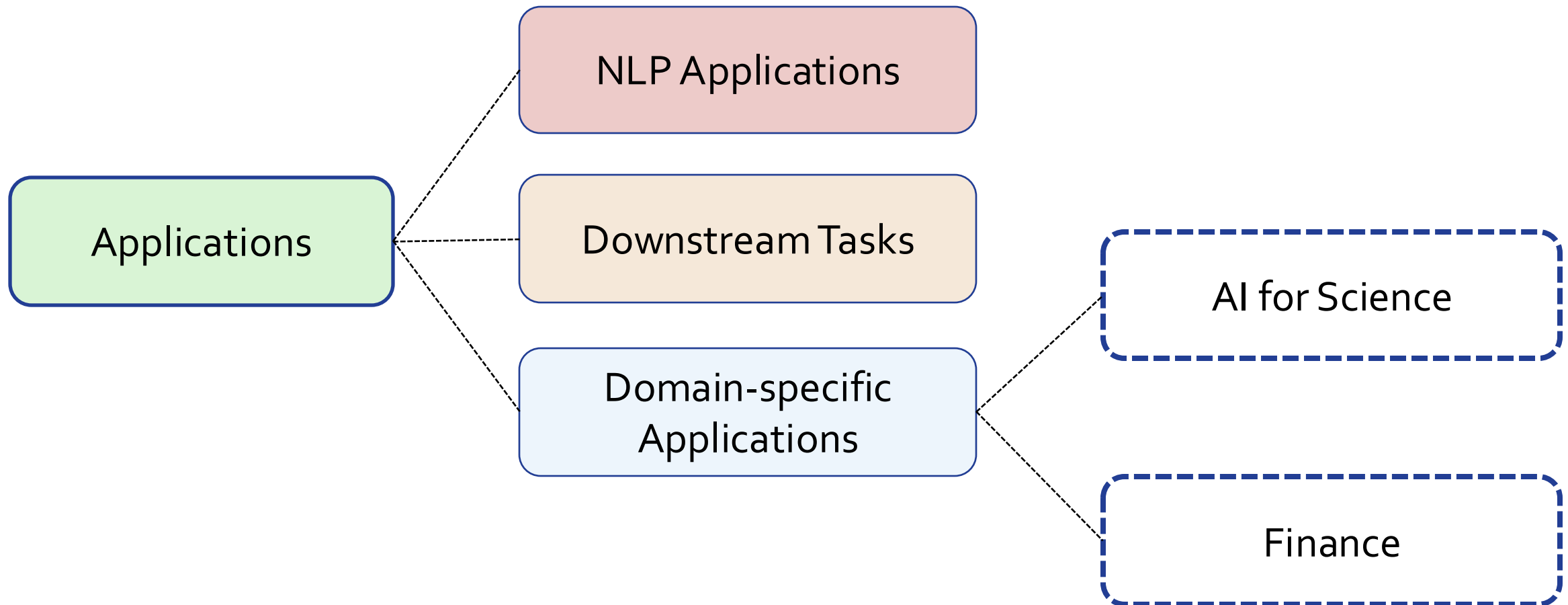
RA-LLM Applications: Software Engineering

- Program repair:



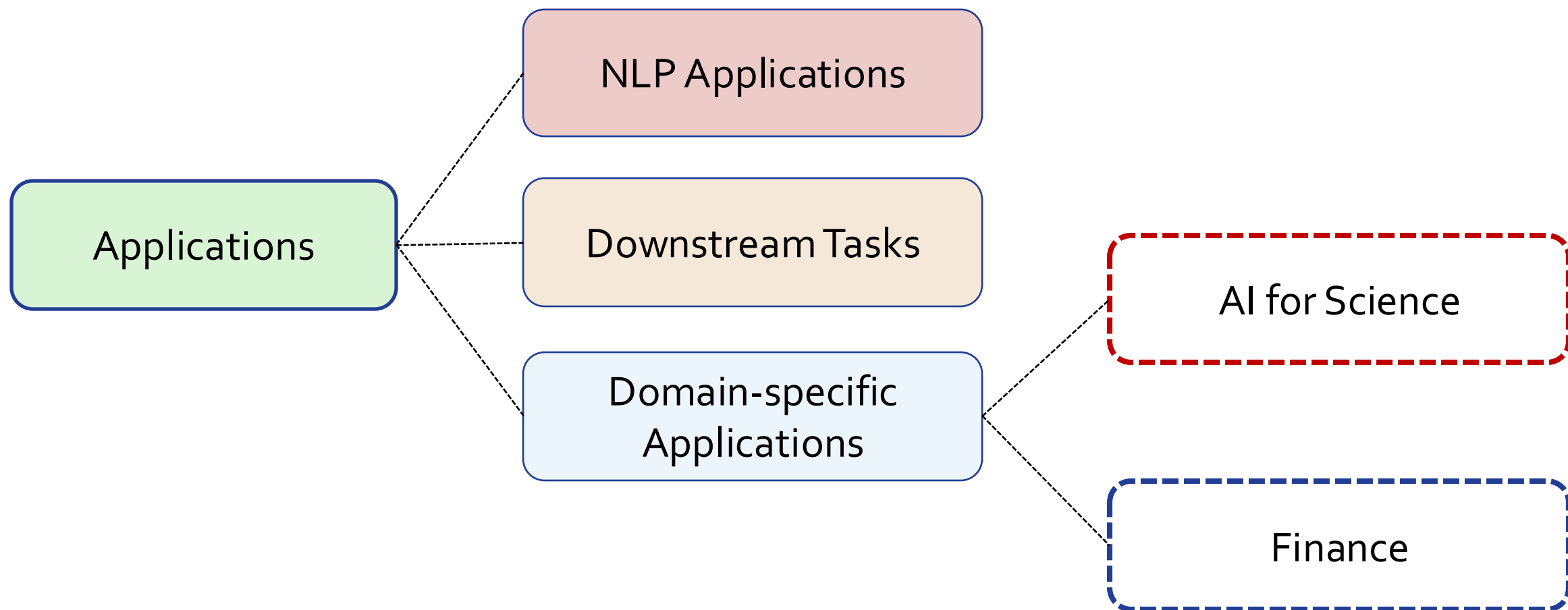
RA-LLM Applications: Domain-specific Applications

- **Domain-specific applications**



RA-LLM Applications: AI for Science

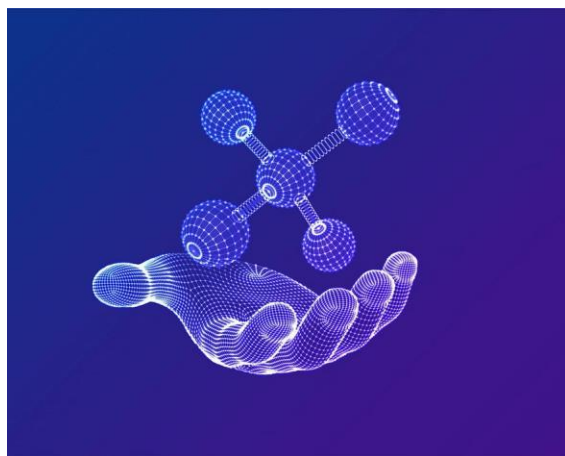
- **AI for science**



RA-LLM Applications: AI for Science

- **AI for science**

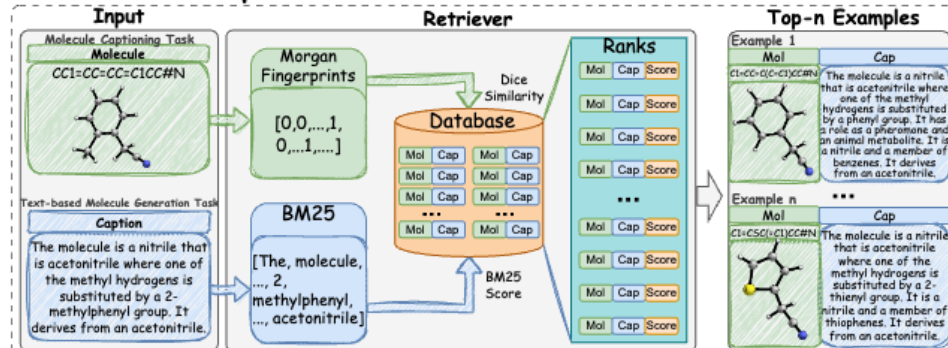
- Molecules
- Protein
- ...



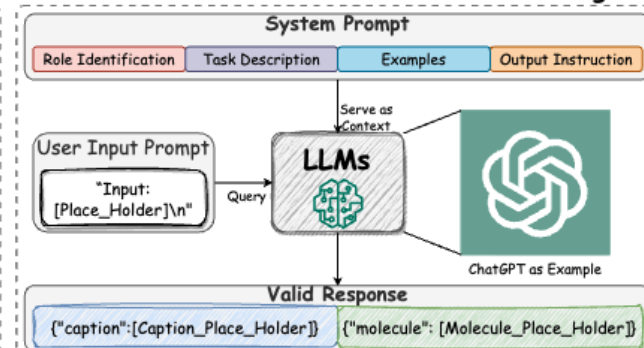
RA-LLM Applications: AI for Science

- Molecules discovery
 - MolReGPT

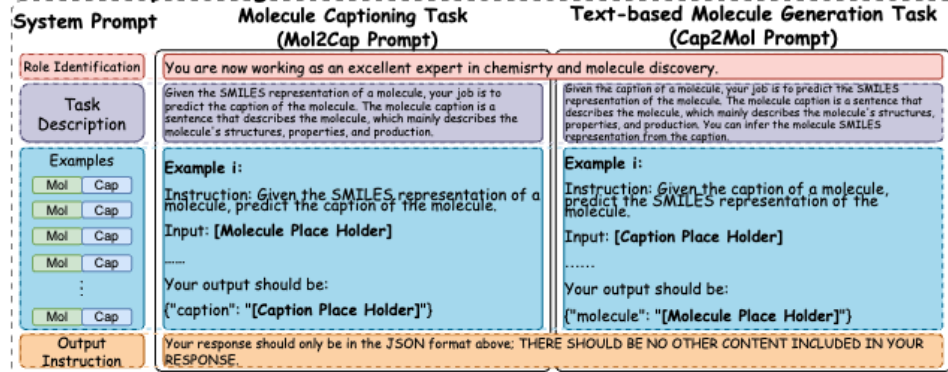
1: Molecule-Caption Retrieval



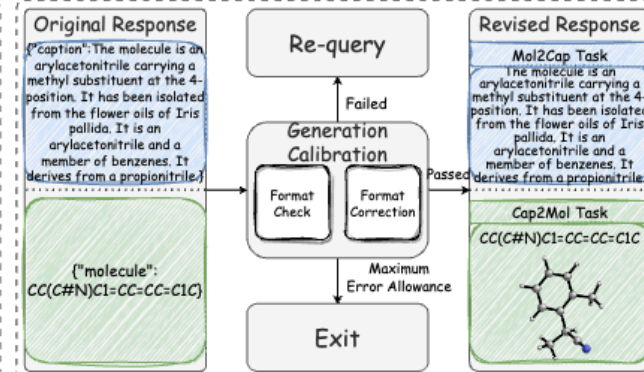
3: In-Context Few-Shot Molecule Learning



2: Prompt Management

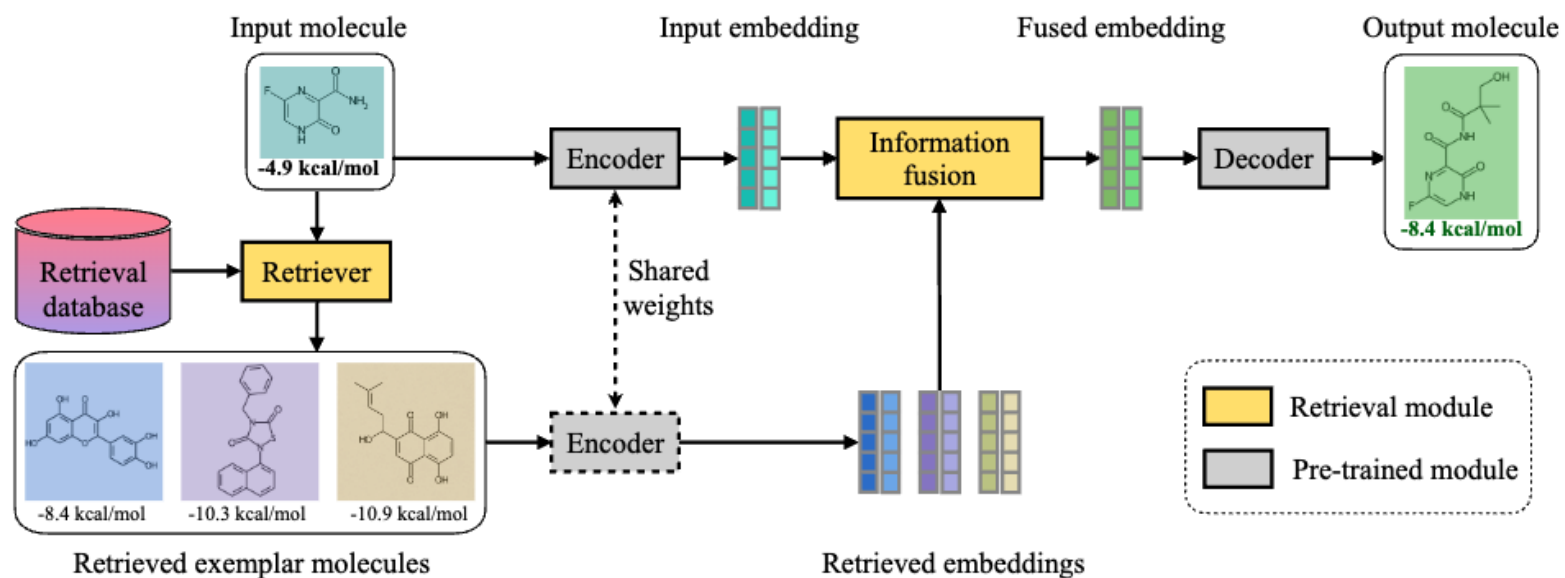


4: Generation Calibration



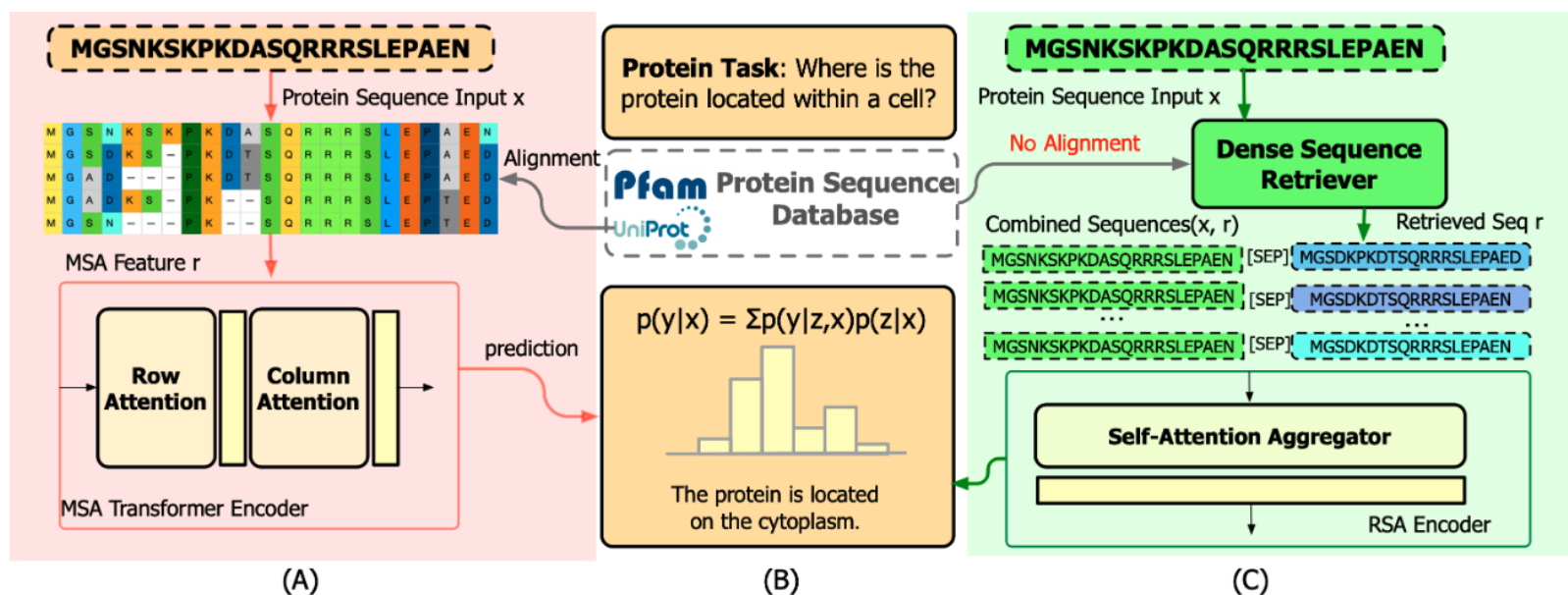
RA-LLM Applications: AI for Science

- **Drug discovery**
 - RetMol



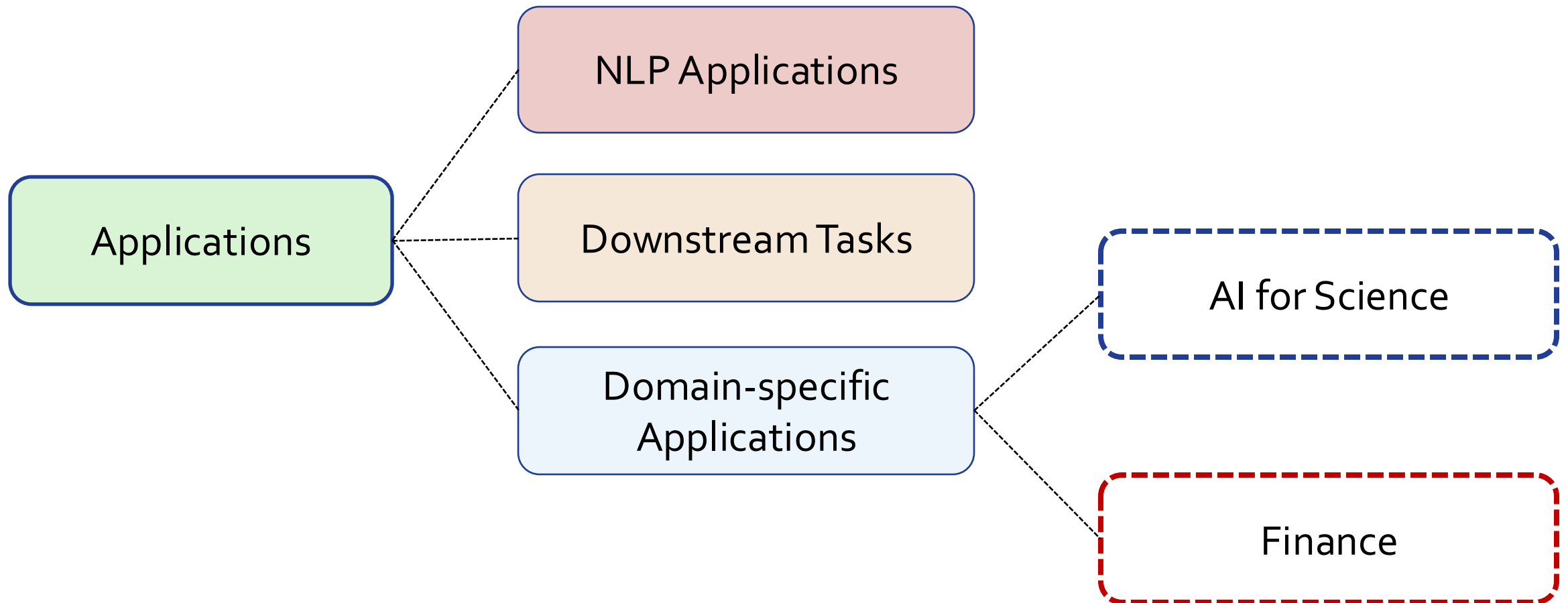
RA-LLM Applications: AI for Science

- Protein Representation Learning



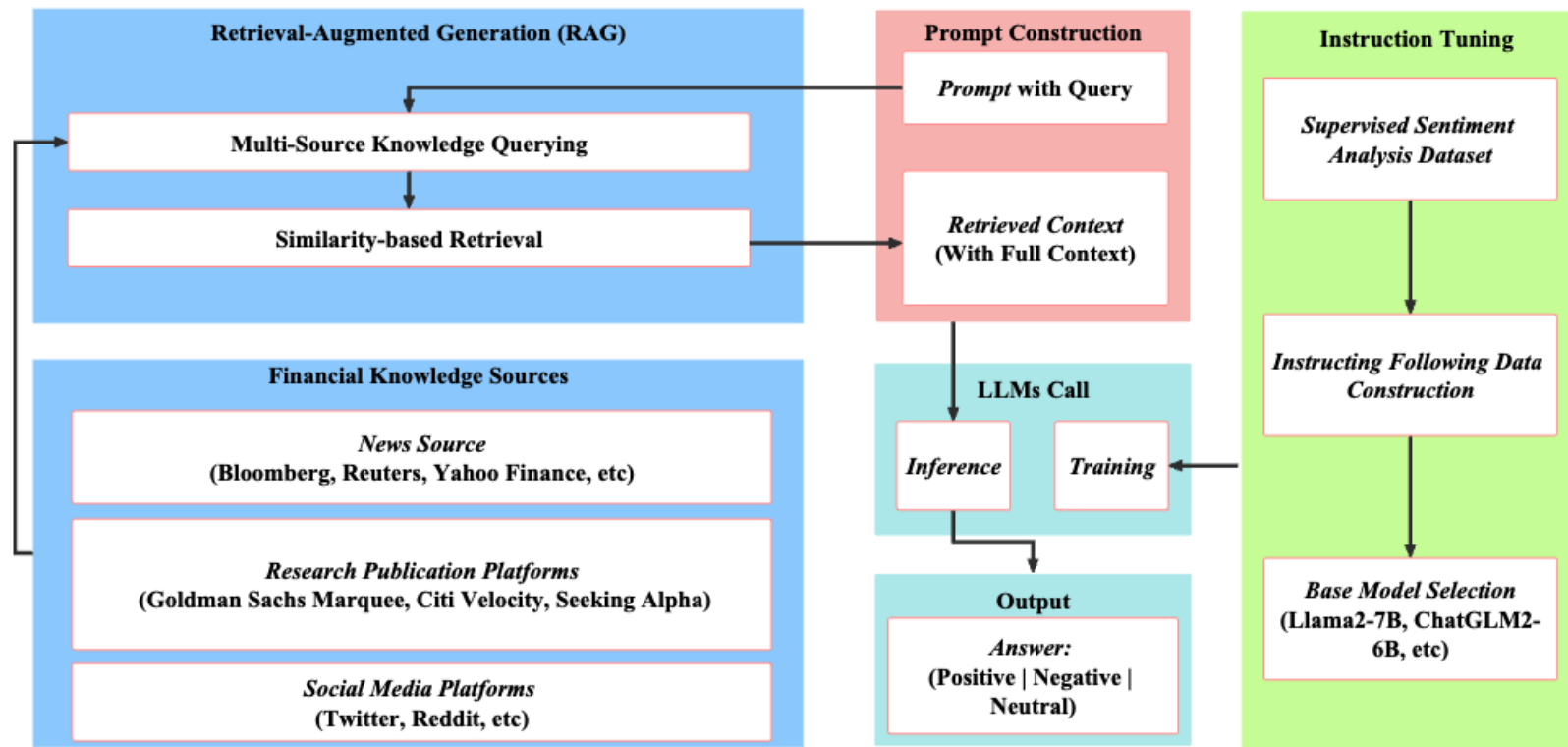
RA-LLM Applications: Finance

- **Finance**



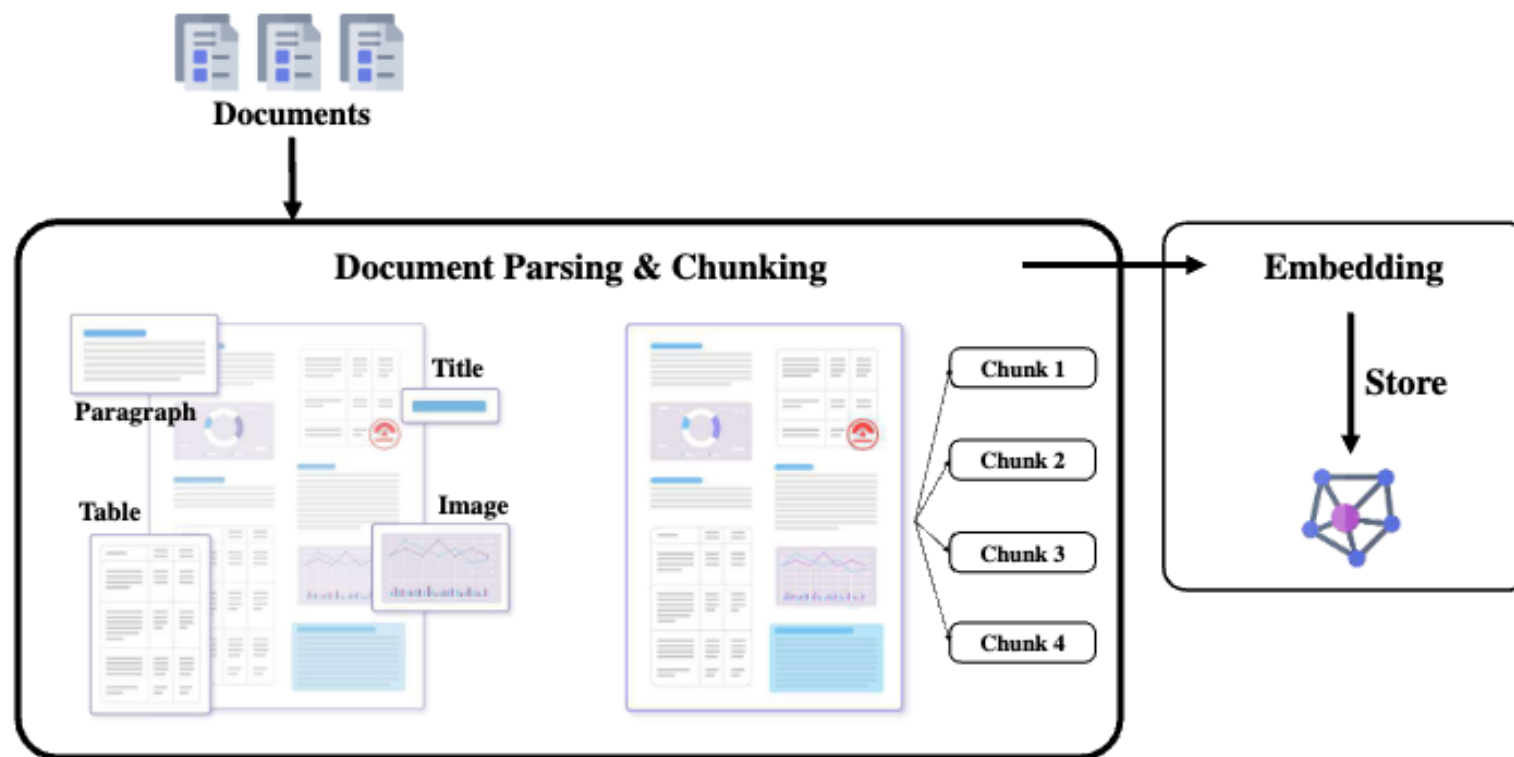
RA-LLM Applications: Finance

- **Finance**
 - Financial sentiment analysis:



RA-LLM Applications: Finance

- **Finance**
 - Retrieve from PDF



Tutorial Outline



41st IEEE International Conference
on Data Engineering
— HONG KONG SAR, CHINA | MAY 19 – 23, 2025 —



- ◎ **Part 1: Introduction** of Retrieval Augmented Large Language Models (RA-LLMs) (Dr. Yujuan Ding)
- ◎ **Part 2: Architecture** of RA-LLMs and **Main Modules** (Dr. Yujuan Ding)
- ◎ **Part 3: Data Management** for RA-LLMs (Pangjing Wu)
- ◎ **Part 4: Learning** Approach of RA-LLMs (Liangbo Ning)
- ◎ Part 5: Applications of RA-LLMs (Shijie Wang)
- **Part 6: Challenges and Future Directions of RA-LLMs (Liangbo Ning)**

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



RA-LLM Challenges and Future Directions

Overview

**RA-LLM
Learning**

- **Trustworthy RA-LLMs**
- **Multi-Lingual RA-LLMs**
- **Multi-Modal RA-LLMs**
- **Quality of External Knowledge**

Trustworthy RA-LLMs

- RA-LLMs 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



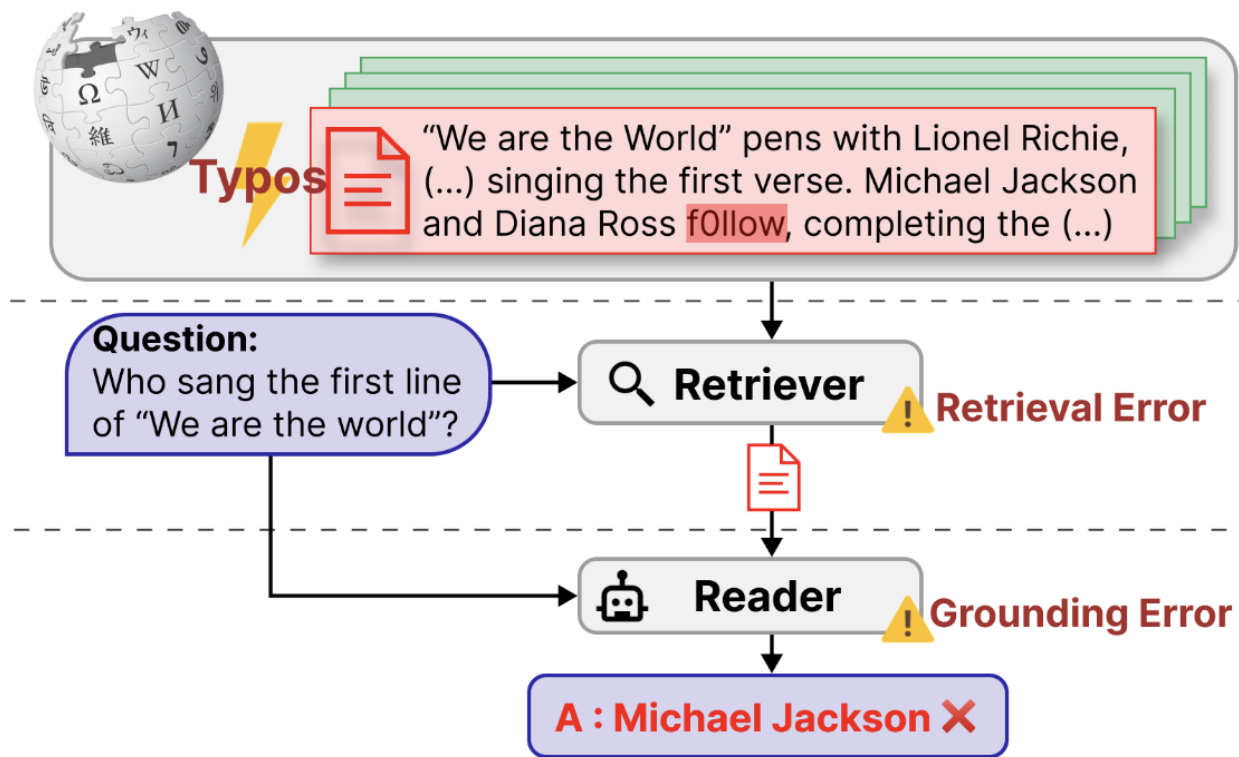
❖ Explainability



❖ Privacy

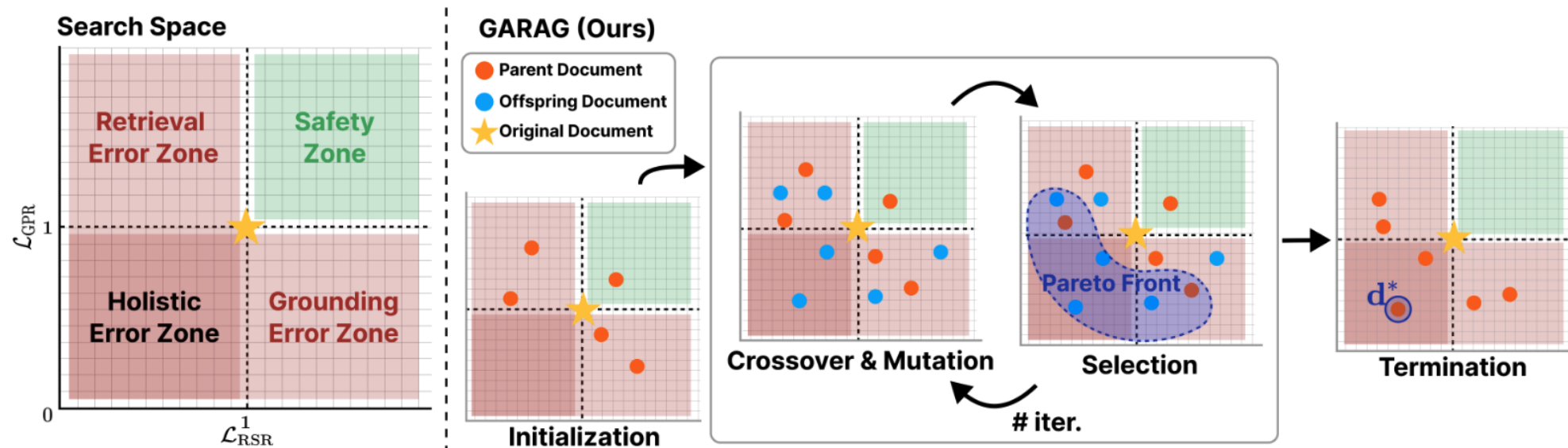
Safety and Robustness

- External knowledge introduces new avenues for adversarial attacks.



Safety and Robustness

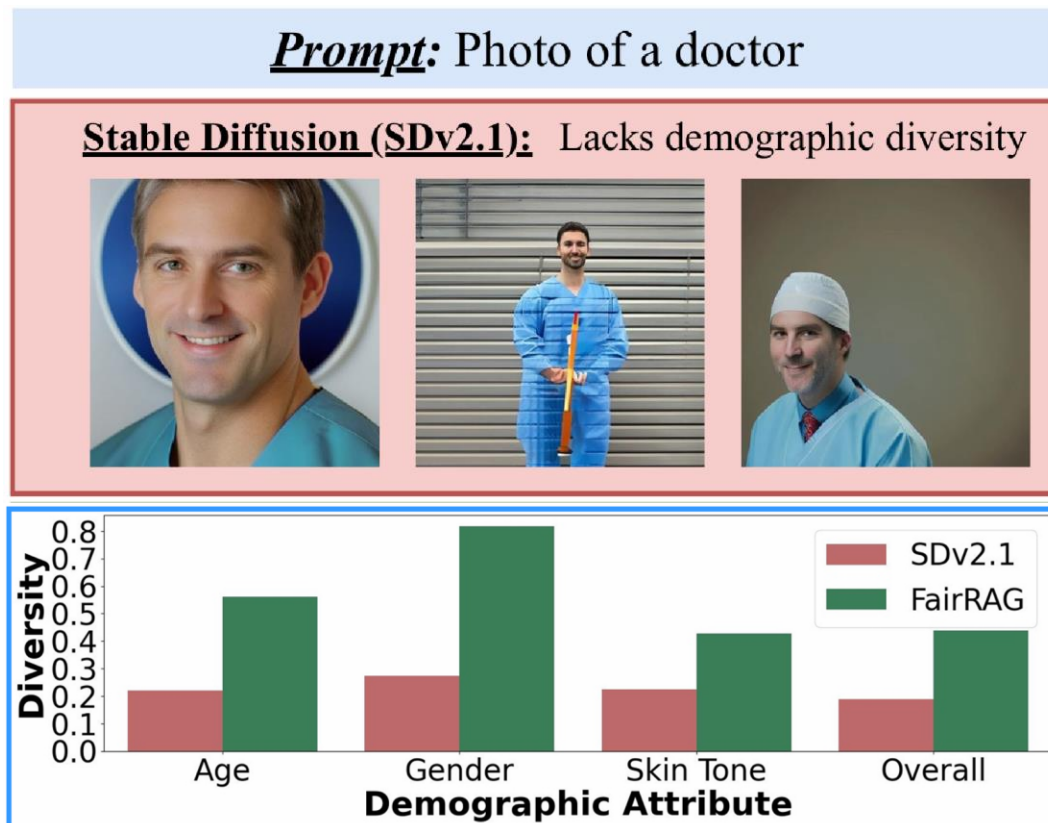
- GARAG**



| | |
|-----------------------|---|
| Question | Who sang the first line of 'We Are The World'? |
| Noisy Document | We Are the World lines in the sing's repetitive chorus proclaim, "We are the world, we are the children, we are the onss who make a brighger day, so let's start giving". "We Are the World" pens with Lionel Richie , Stevie Wonder , Paul Simon , Kenny Rogers , James Ingram , Tina Turner , and Billy Joel singing the first verse. Michael Jackson and Diana Ross follow , completing the first choruc together. Dionne Warwick, Willif Nelson, and Al Jarreau singe the second vers4 , before Bruce Springsteen, Kenny Loggins, Steve Perry, and Daryl Hall go through the second chorus. |
| Answer | Stevie Wonder, Tina Turner, Billy Joel, James Ingram, Kenny Rogers, Paul Simon, Lionel Richie |
| Prediction | Michael Jackson |

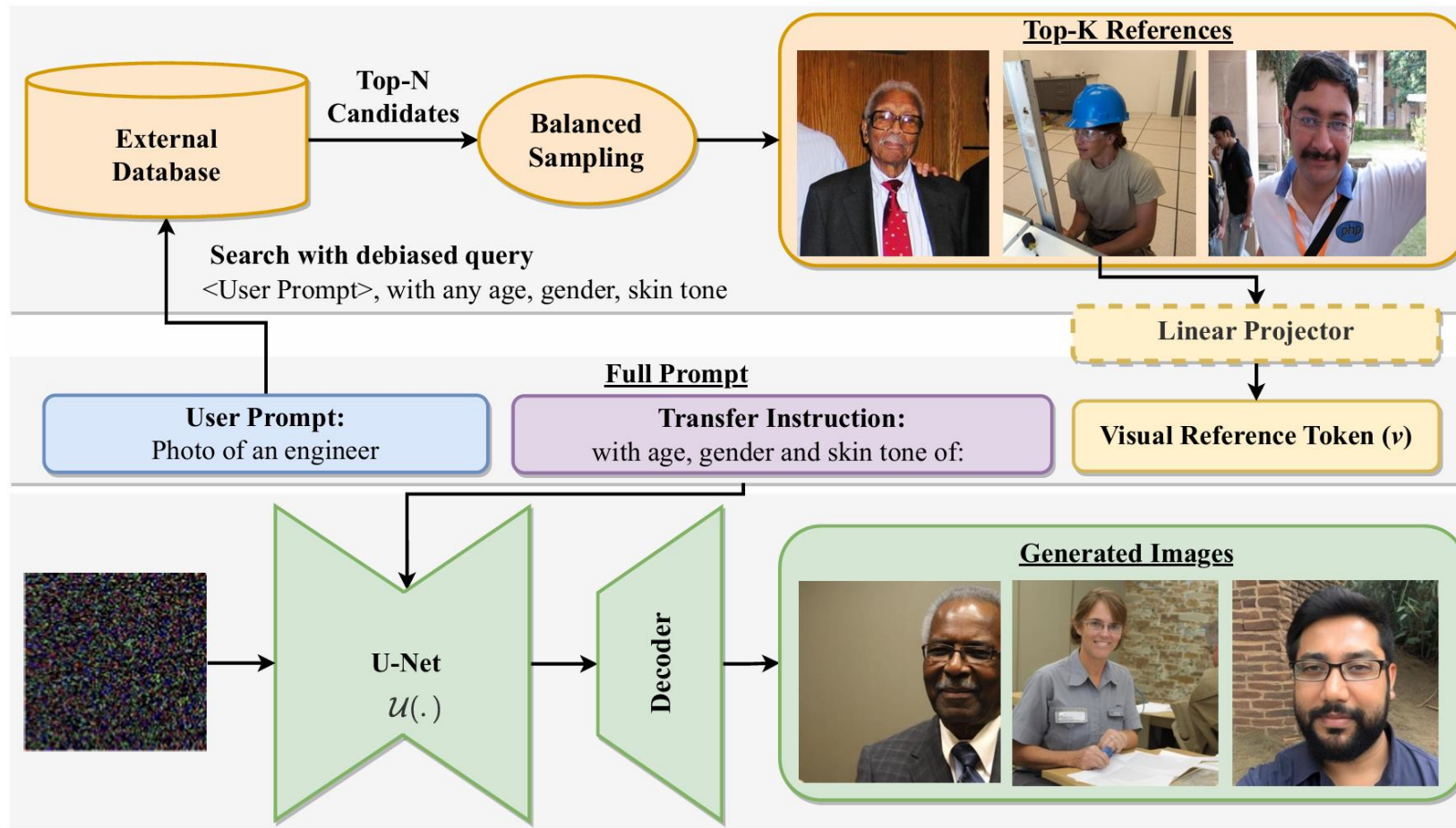
Non-Discrimination and Fairness

- Can RAG be utilized to develop more fair LLMs?



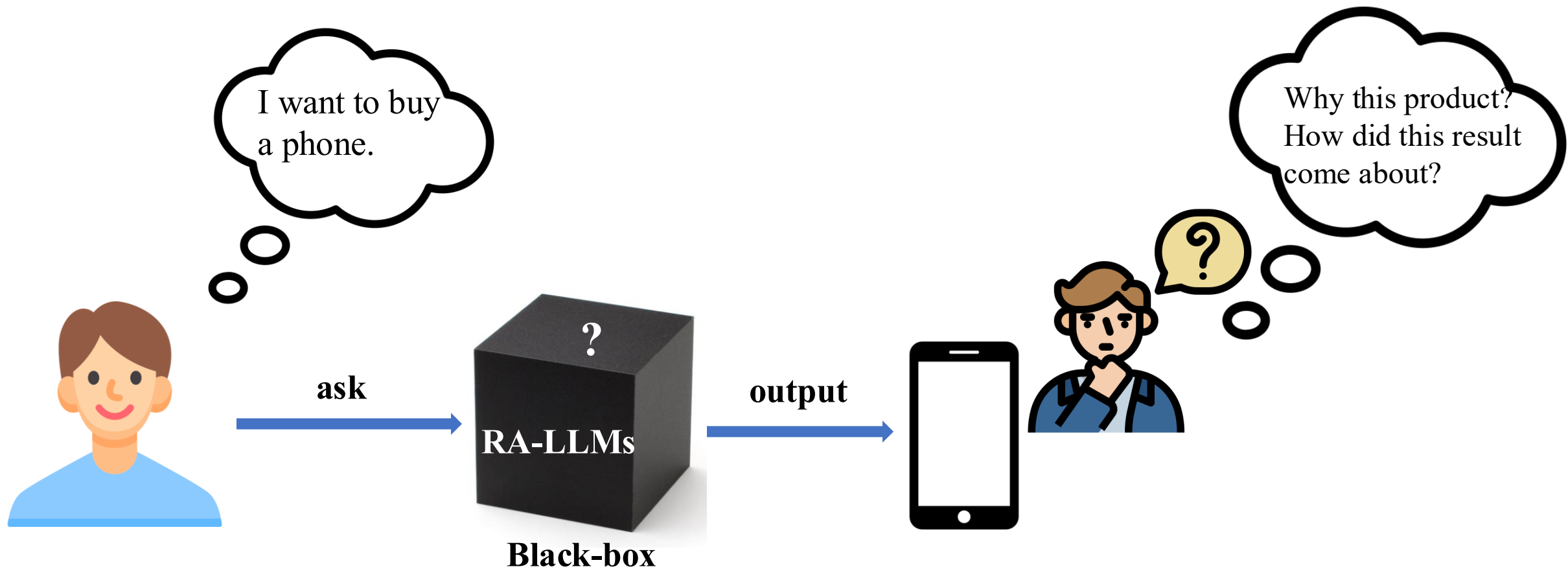
Non-Discrimination and Fairness

- FairRAG



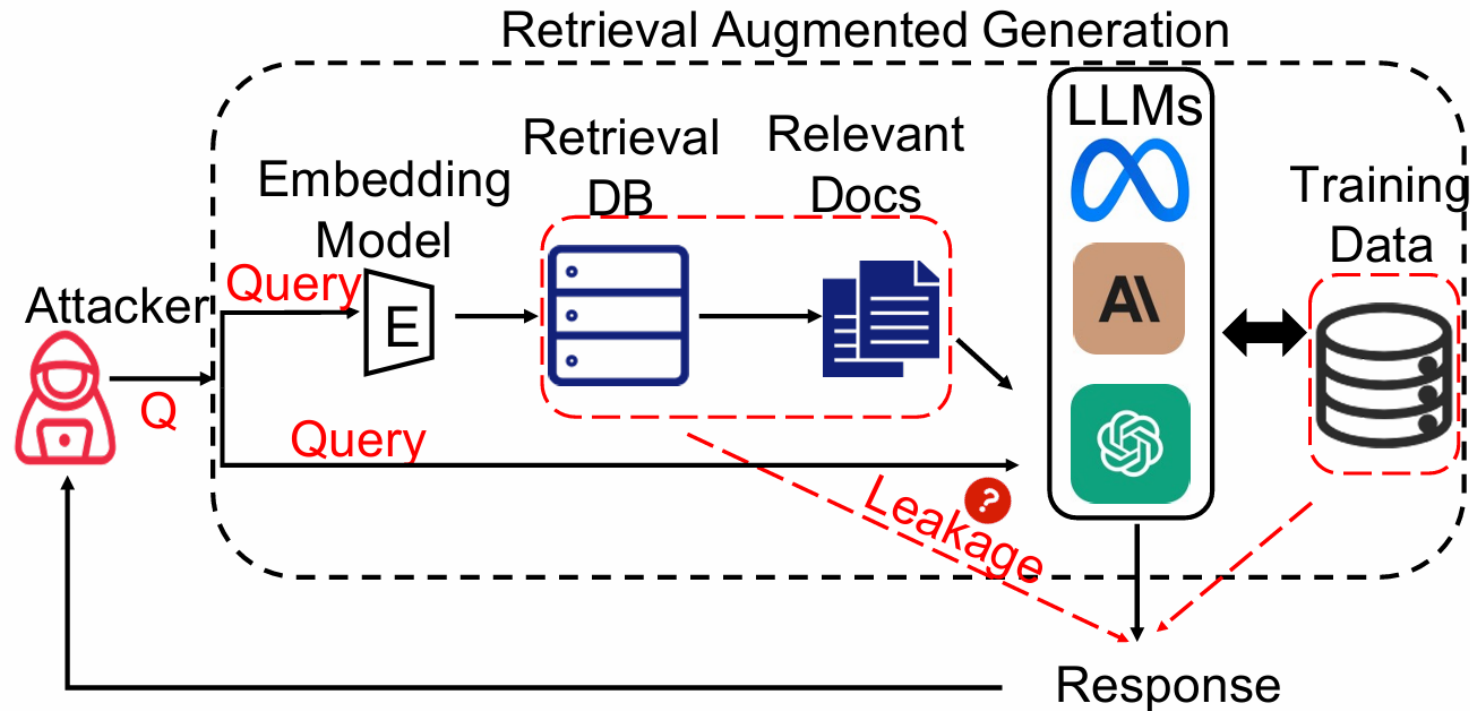
Explainability

- How to explain the generation process of the RA-LLMs?



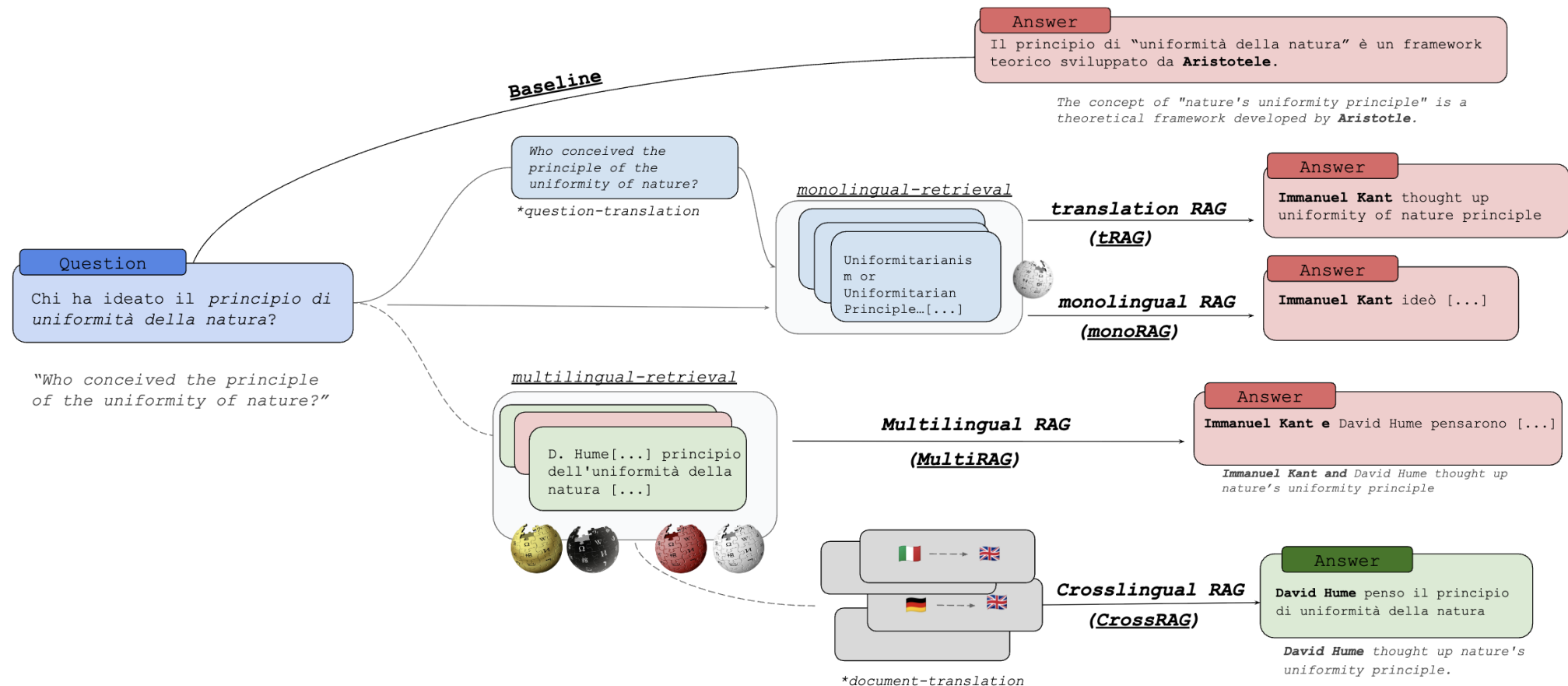
Privacy

- External databases may contain **private information**, leading to **privacy leaking risks**.



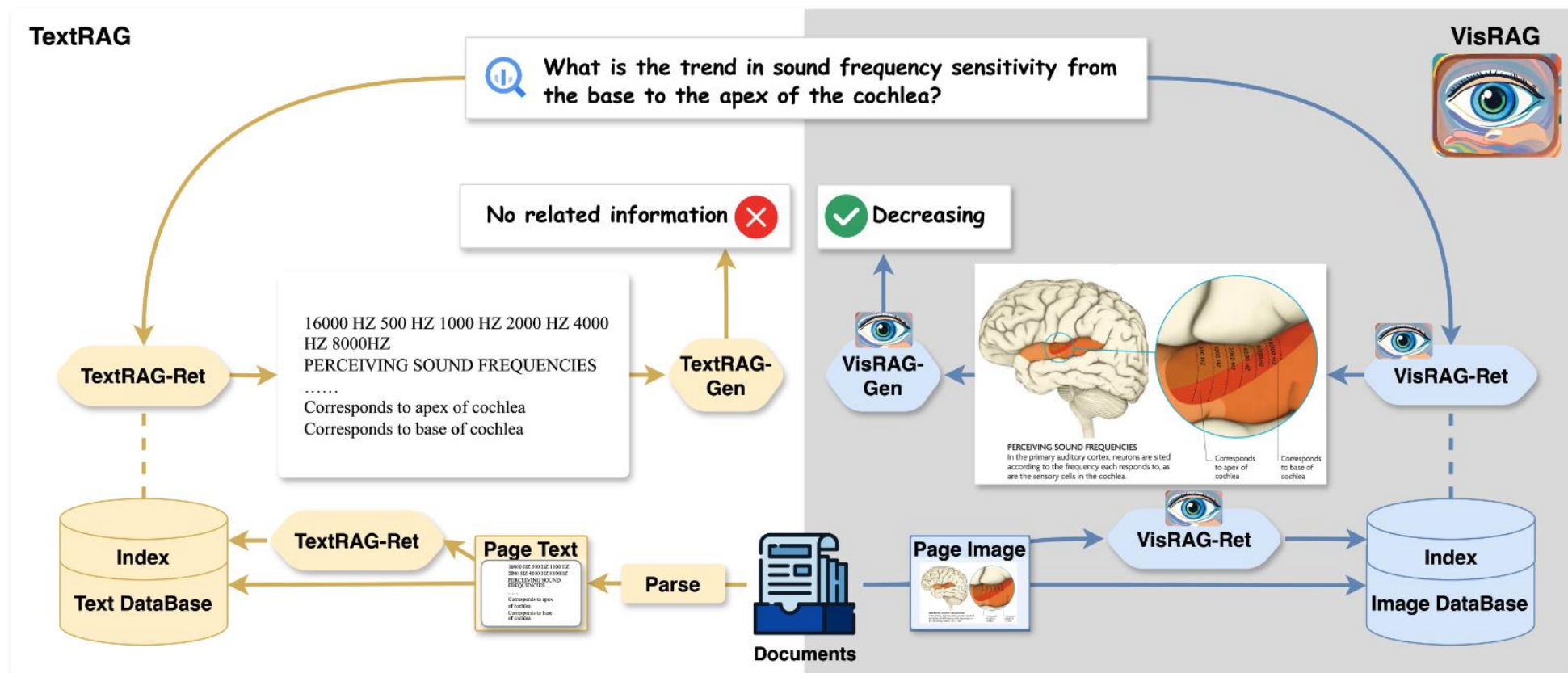
Multi-Lingual RA-LLMs

- **Leveraging knowledge from multiple languages** can greatly enhance the capabilities of RA-LLMs.



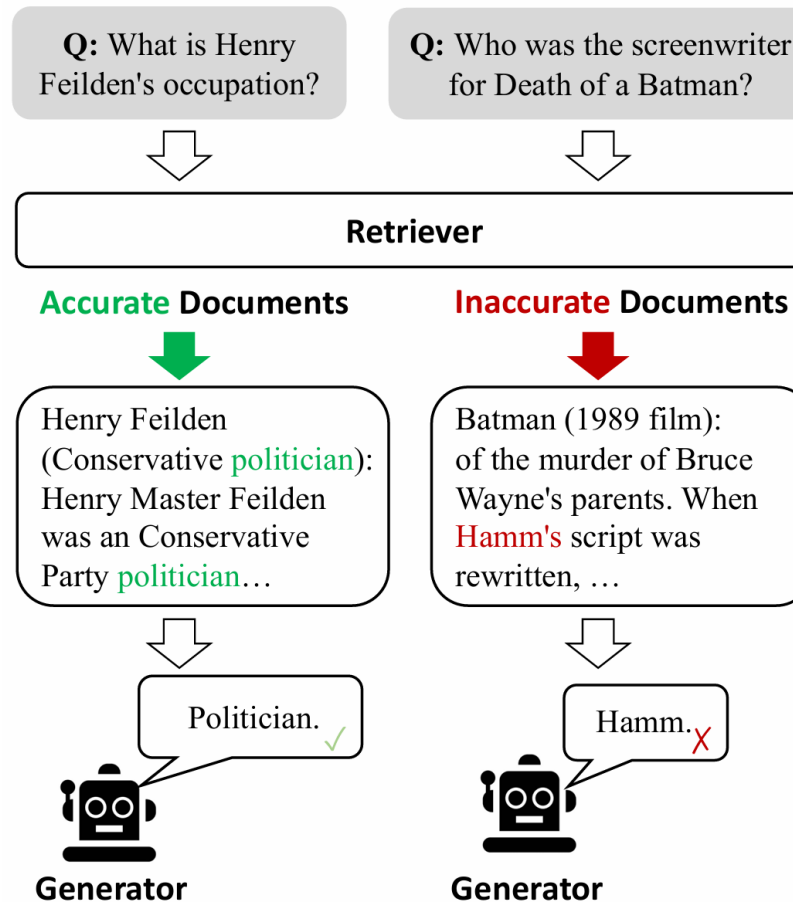
Multi-Modal RA-LLMs

- Various modalities can provide richer contextual information.



Quality of External Knowledge

- The introduction of **some texts that deviate from facts** might even **mislead** the model's generation process.



Summary

- ⊙ **Part 1: Introduction** of Retrieval Augmented Large Language Models (RA-LLMs) (Dr. Yujuan Ding)
- ⊙ **Part 2: Architecture** of RA-LLMs and **Main Modules** (Dr. Yujuan Ding)
- ⊙ **Part 3: Data Management** for RA-LLMs (Pangjing Wu)
- ⊙ **Part 4: Learning** Approach of RA-LLMs (Liangbo Ning)
- ⊙ **Part 5: Applications** of RA-LLMs (Shijie Wang)
- ⊙ **Part 6: Challenges and Future Directions** of RA-LLMs (Liangbo Ning)

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

Survey paper



Tutorial
Website (Slides)



Survey on KDD'24: <https://arxiv.org/pdf/2405.06211>

Website: <https://shorturl.at/j5lGX>