# LLM LE6 VT2025
## Inference, RAG, and Reasoning

Under Construction

**Fredrik Heintz**

**Dept. of Computer Science**
**Linköping University**

**fredrik.heintz@liu.se**

**@FredrikHeintz**

Outline:

- **Inference**

- **Retrieval Augmented Generation**

- **In-Context Learning**

- **Reasoning**
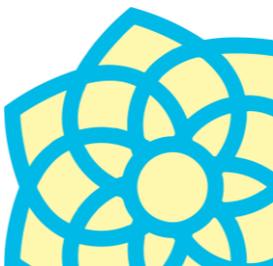
LiU LINKÖPING UNIVERSITY

# Language modelling

- **Language modelling** is the task of predicting which word comes next in a sequence of words.

- More formally, given a sequence of words $w_1, \ldots, w_t$ we want to know the probability of the next word, $w_{t+1}$:
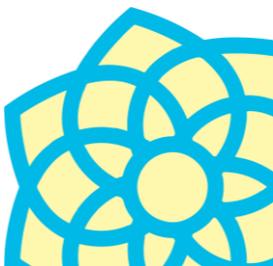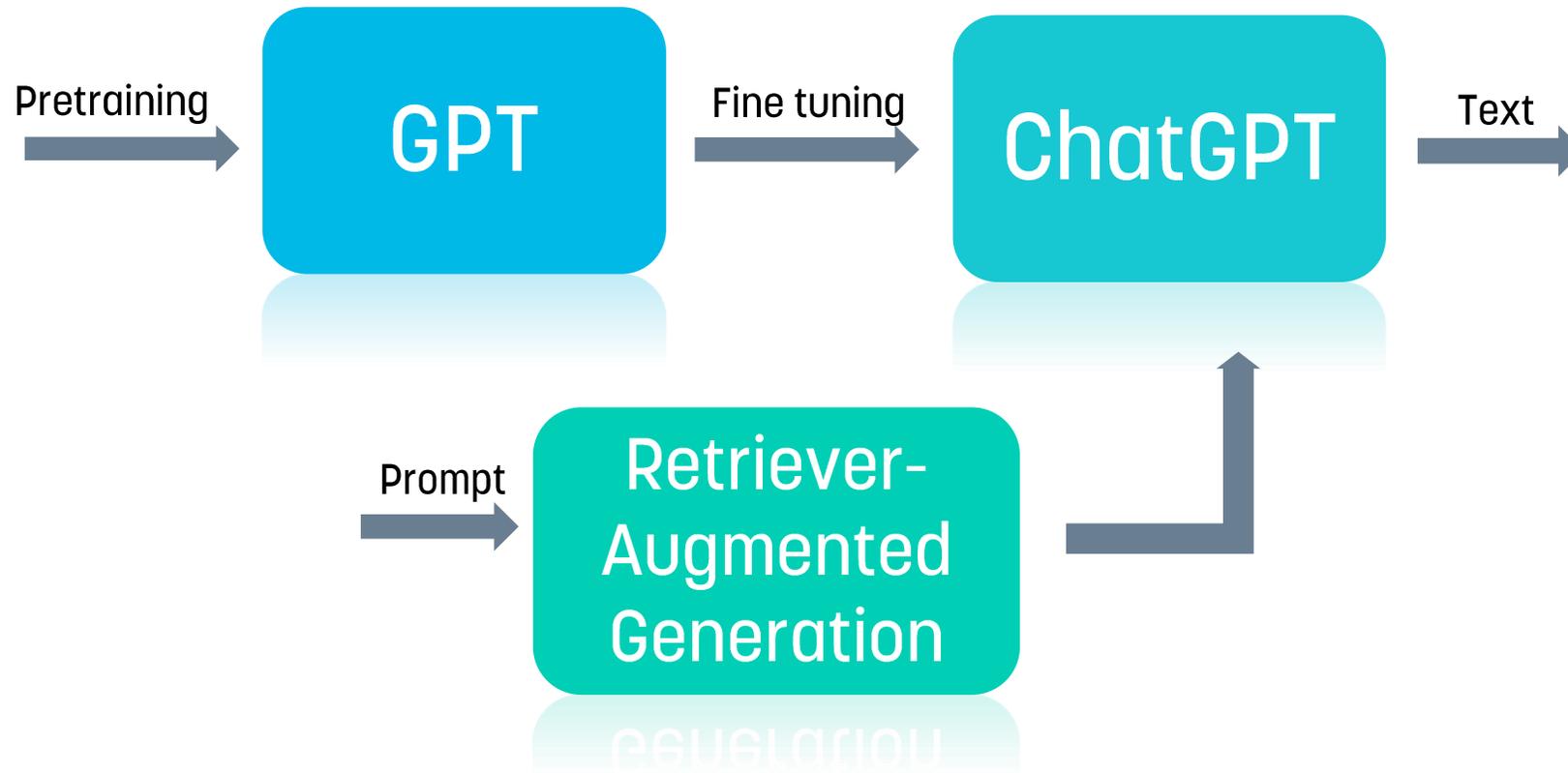
$$P(w_{t+1} \mid w_1, \ldots, w_t)$$

- We are assuming that $w_{t+1}$ comes from a finite vocabulary $V$.
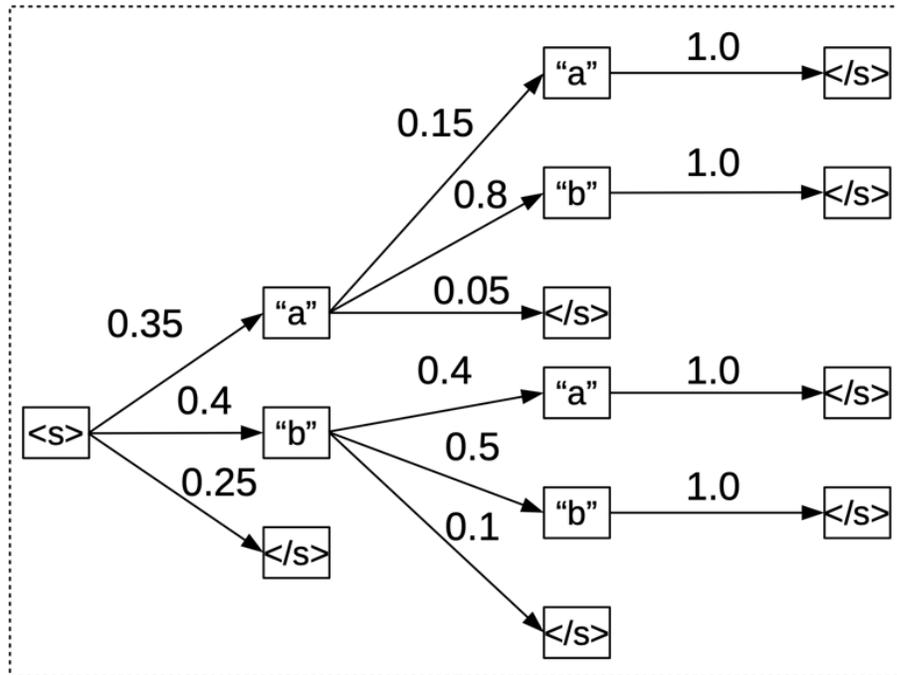
  language models = classifiers

TrustLLM

Funded by
the European Union

# How Does ChatGPT Work?

# Inference

# Greedy search methods do not always lead to the most likely output.

Vocabulary = {a, b, </s>}
Numbers above each edge are the transition probabilities $P(x_t|x_{1:t-t})$

If we were to choose the sequence that maximizes $P(x_1, \dots, x_T)$ , which of the following would get generated?

(a) [a, b, </s>]
(b) [a, a, </s>]
(c) [b, b, </s>]
(d) [b, a, </s>]
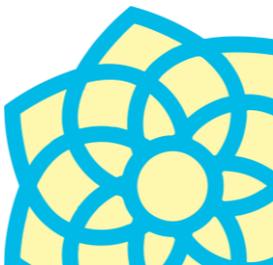
# Greedy search methods do not always lead to the most likely output.
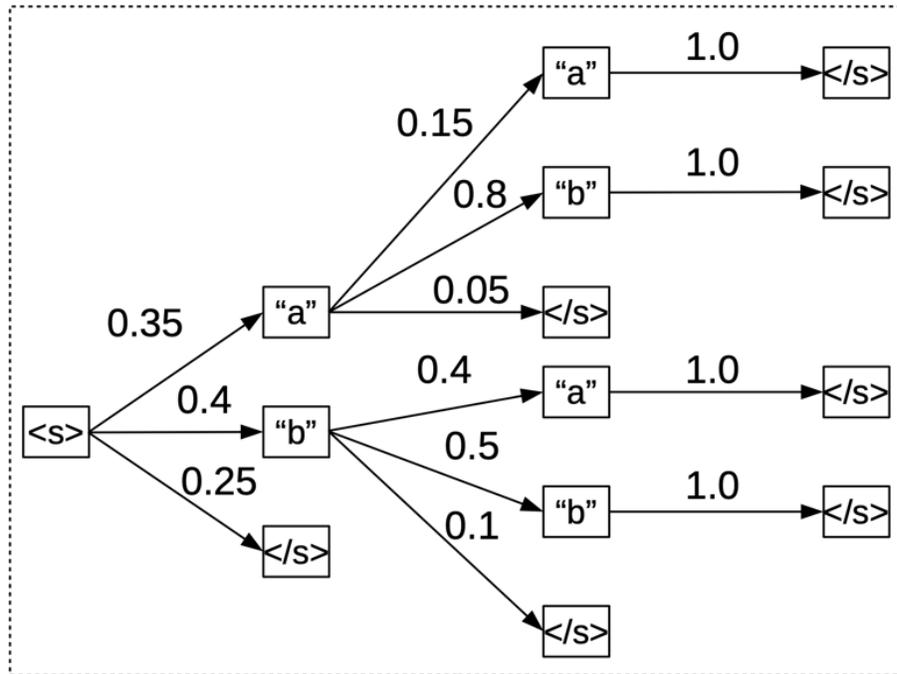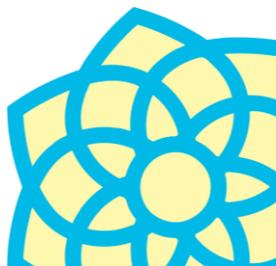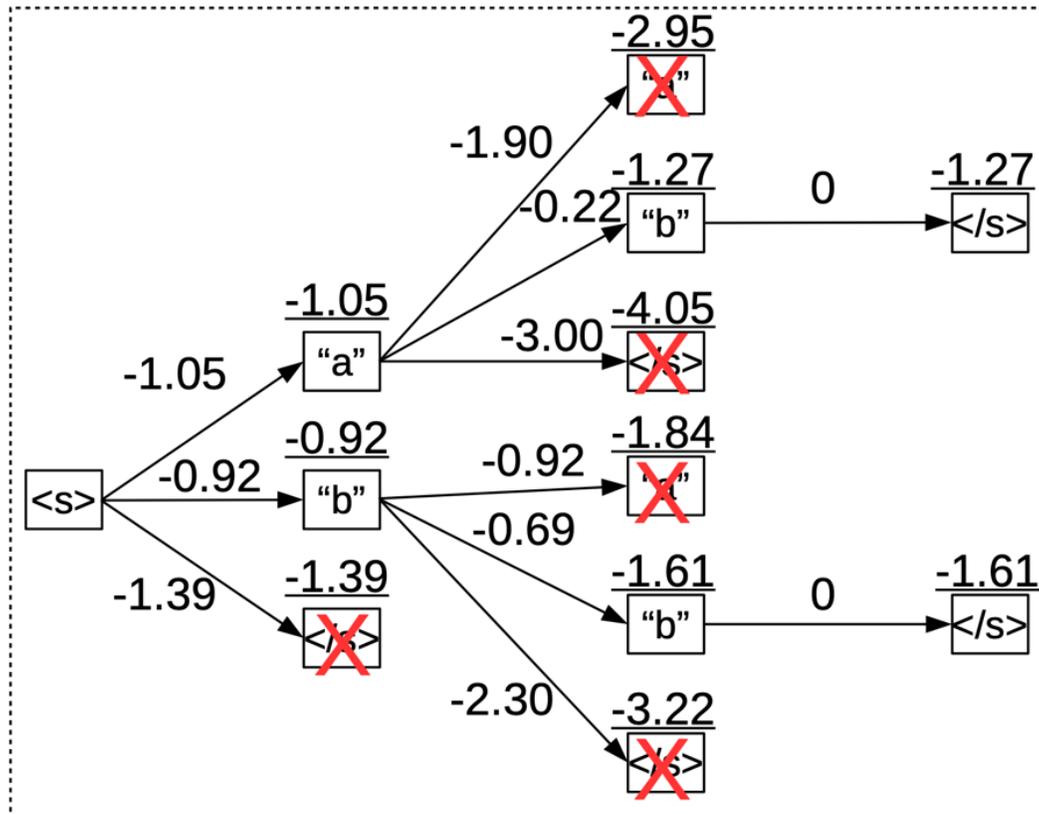


Vocabulary = {a, b, </s>}
Numbers above each edge are the transition probabilities $P(x_t|x_{1:t-t})$

If we were to choose the sequence that maximizes $P(x_1, \ldots, x_T)$ , which of the following would get generated?

**(a)** [a, b, </s>]
(b) [a, a, </s>]
(c) [b, b, </s>]
(d) [b, a, </s>]

# Beam search explores multiple possible output sequences, trying to find the overall most likely one.



Vocabulary = {a, b, </s>}

Numbers above the boxes are $logP(x_t|x_{1:t-1})$

Numbers shown on edges are $logP(x_1, ..., x_t)$

Suppose we use beam search with a **beam size** of 2.

# Beam search explores multiple possible output sequences, trying to find the overall most likely one.



Vocabulary = {a, b, </s>}

Numbers above the boxes are $logP(x_t|x_{1:t-1})$

Numbers shown on edges are $logP(x_1, ..., x_t)$

Suppose we use beam search with a **beam size** of 2.

**Score each path and keep the top 2**

# Beam search explores multiple possible output sequences, trying to find the overall most likely one.
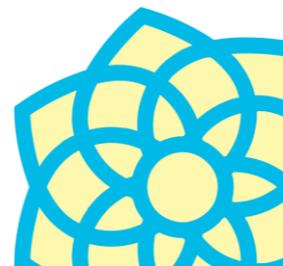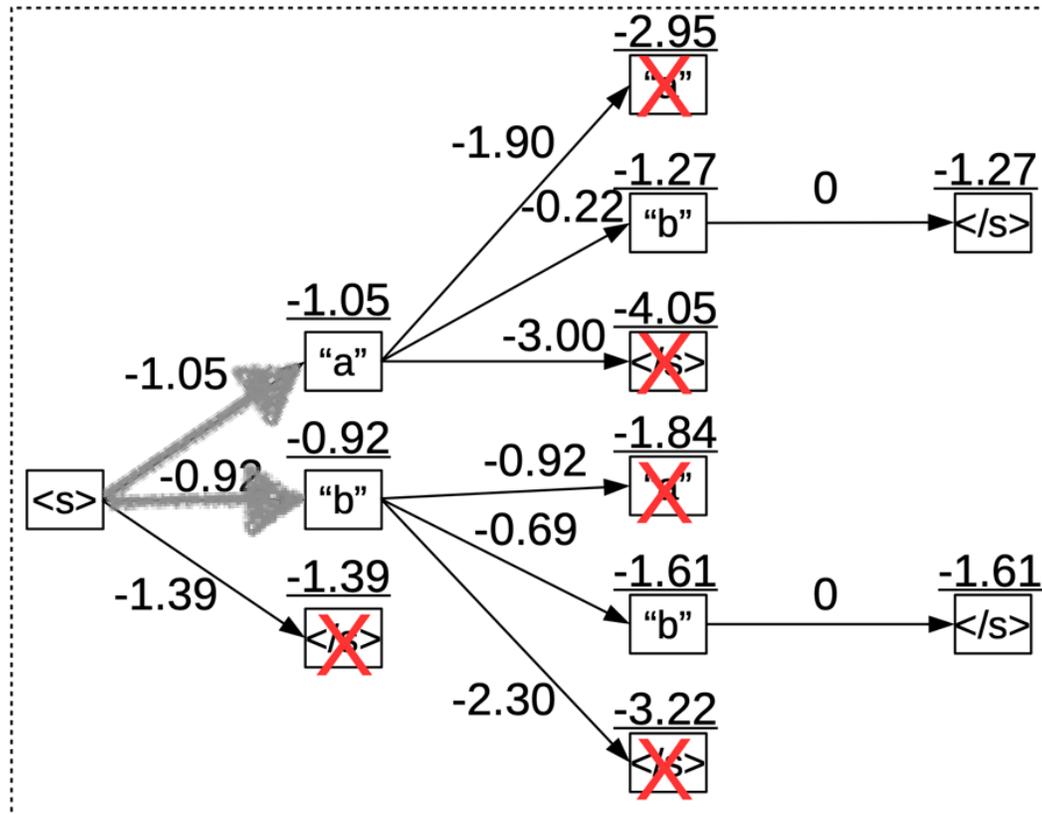


Vocabulary = {a, b, </s>}

Numbers above the boxes are $\log P(x_t | x_{1:t-1})$

Numbers shown on edges are $\log P(x_1, \ldots, x_t)$

Suppose we use beam search with a **beam size** of 2.

**Score each path and keep the top 2**

# Beam search explores multiple possible output sequences, trying to find the overall most likely one.
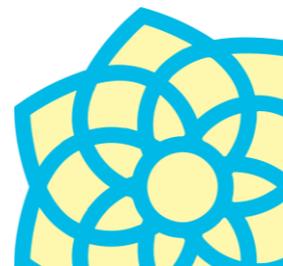


Vocabulary = {a, b, </s>}

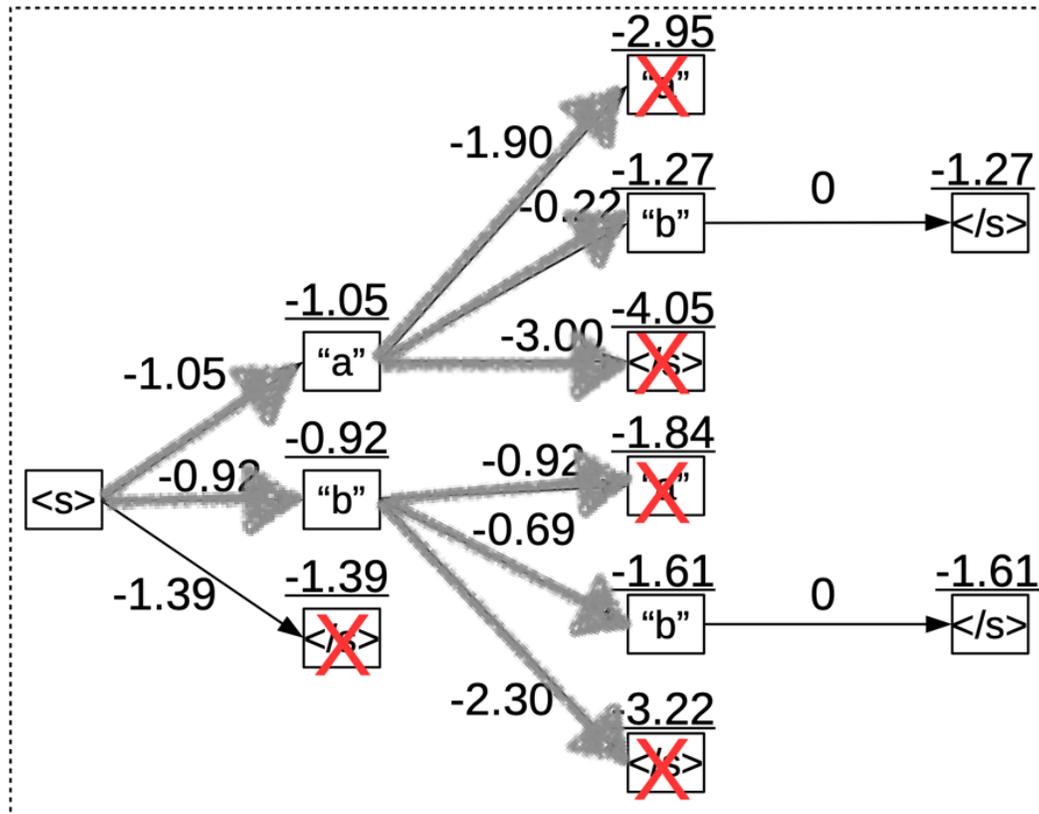Numbers above the boxes are $\log P(x_t | x_{1:t-1})$

Numbers shown on edges are $\log P(x_1, \ldots, x_t)$

Suppose we use beam search with a **beam size** of 2.

**Score each path and keep the top 2**

# Beam search explores multiple possible output sequences, trying to find the overall most likely one.
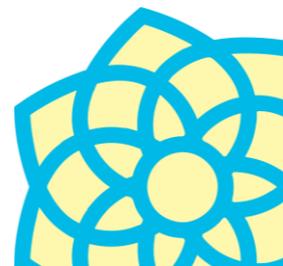


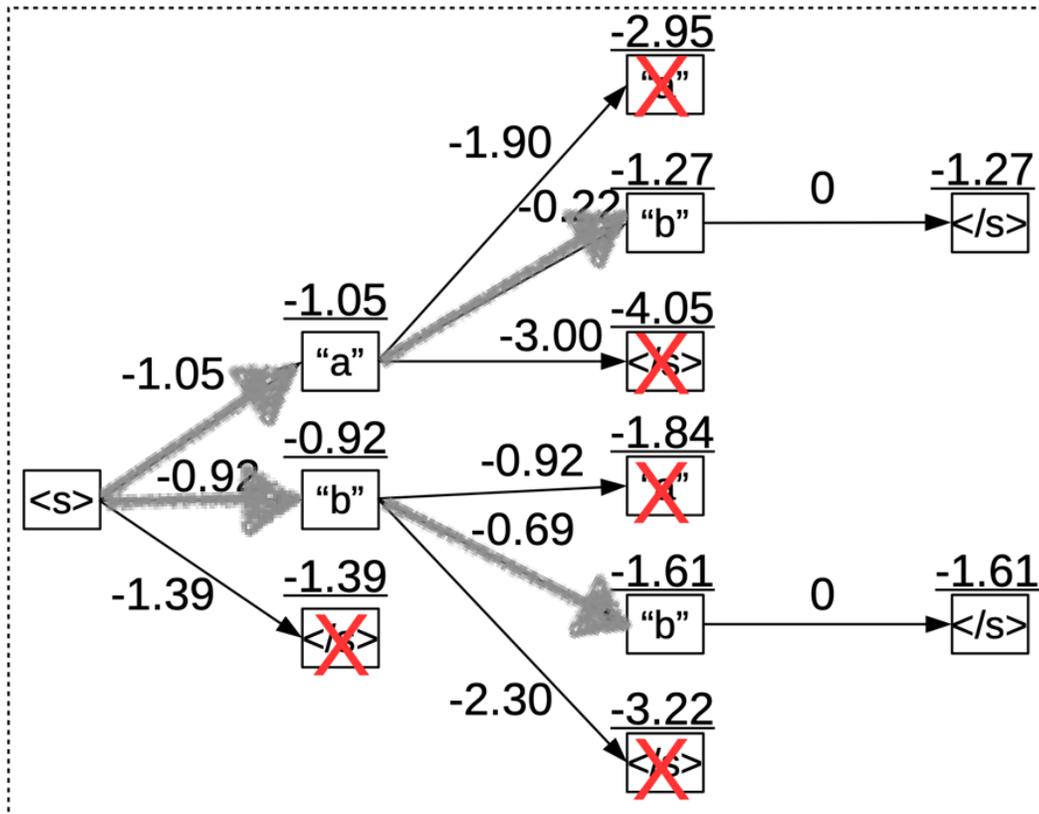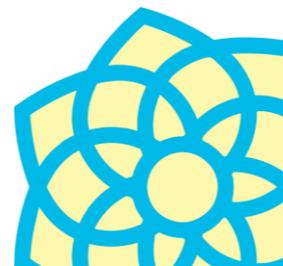Score each path and keep the top 2

Vocabulary = {a, b, </s>}

Numbers above the boxes are $\log P(x_t | x_{1:t-1})$

Numbers shown on edges are $\log P(x_1, \ldots, x_t)$

Suppose we use beam search with a **beam size** of 2.

The paths that are still open at any step of the beam search algorithm are called **beams**.

# Retrieval Augmented Generation (RAG)

LINKÖPING
UNIVERSITY

# How do normal parametric LLMs work?

**Encapsulating everything in parameters by *pre-training* parameters on large-scale text corpora**

Pittsburgh is a city in the county seat of Allegheny County, Pennsylvania, United States

**Pre-training data**

$$P(x_n \mid x_1, x_2, \ldots, x_{n-1})$$

**LLM**

# How do normal parametric LLMs work?

**Encapsulating everything in parameters by *pre-training* parameters on large-scale text corpora**



Pittsburgh is a city in and the county seat of Allegheny County, Pennsylvania, United States

Pre-training data

Allegheny

Pennsylvania

King

$P(x_n \mid x_1, x_2, \ldots, x_{n-1})$

LLM

Pittsburgh    is    located    in
$x_1$         $x_2$  $x_3$      $x_4$

# Limitations of parametric LLMs #1: Hallucinations

**LLMs cannot memorize everything in their parameters (yet), resulting in factual inaccuracy**

Summarize Akari Asai's research including her representative papers.

# Limitations of parametric LLMs #1: Hallucinations

**LLMs cannot memorize everything in their parameters (yet), resulting in factual inaccuracy**

Summarize Akari Asai's research including her representative papers.

Akari Asai is a prominent researcher in the field of natural language processing (NLP), with a focus on question answering (QA), information retrieval, and multi-lingual models. Her research revolves around building efficient and robust systems for extracting information from large corpora,

TrustLLM

Funded by the European Union

# Limitations of parametric LLMs #1: Hallucinations

**LLMs cannot memorize everything in their parameters (yet), resulting in factual inaccuracy**

Summarize Akari Asai's research including her representative papers.

Akari Asai is a prominent researcher in the field of natural language processing (NLP), with a focus on question answering (QA), information retrieval, and multi-lingual models. Her research revolves around building efficient and robust systems for extracting information from large corpora,

## Representative Papers:

1. **"XOR QA: Cross-lingual Open-Retrieval Question Answering" (2021)** In this paper, Asai introduced XOR QA, a benchmark for cross-lingual open-retrieval question answering, targeting QA across multiple languages. The model tackles the challenges of finding answers in languages different from the one in which the question is asked, highlighting the

# Limitations of parametric LLMs #1: Hallucinations

**LLMs cannot memorize everything in their parameters (yet), resulting in factual inaccuracy**



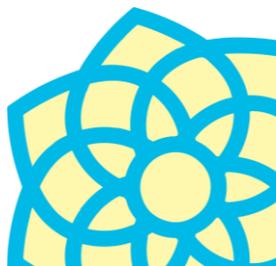Summarize Akari Asai's research including her representative papers.

Akari Asai is a prominent researcher in the field of nat... on question answering (QA), information retrieval, and... around building efficient and robust systems for extra...

**Representative Papers:**

1. **"XOR QA: Cross-lingual Open-Retrieval Questio...** introduced XOR QA, a benchmark for cross-lingua... targeting QA across multiple languages. The mode... in languages different from the one in which the q...

**XOR QA: Cross-lingual Open-Retrieval Question Answering**

Akari Asai, Jungo Kasai, Jonathan H. Clark, Kenton Lee, Eunsol Choi, Hannaneh Hajishirzi
University of Washington  Google Research
The University of Texas at Austin  Allen Institute for AI
{akari, jkasai, hannaneh}@cs.washington.edu
{jhclark, kentonl}@google.com, eunsol@cs.utexas.edu

**Abstract**

Multilingual question answering tasks typically assume that answers exist in the same language as the question. Yet in practice, many languages face both *information scarcity*—where languages have few reference articles—and *information asymmetry*—where questions reference concepts from other cultures. This work extends open-retrieval question answering to a cross-lingual setting enabling questions from one language to be answered via answer content from another language. We construct a large-scale dataset

# Limitations of parametric LLMs #1: Hallucinations

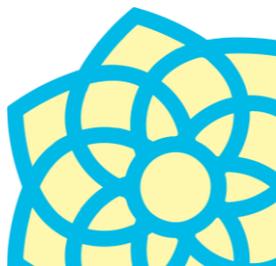**LLMs cannot memorize everything in their parameters (yet), resulting in factual inaccuracy**

> Summarize Akari Asai's research including her representative papers.

Akari Asai is a prominent researcher in the field of natural language processing (NLP), with a focus on question answering (QA), information retrieval, and multi-lingual models. Her research revolves around building efficient and robust systems for extracting information from large corpora,

...

2. **"Dense Passage Retrieval for Open-Domain Question Answering" (2020)** This is one of her most cited works, where she co-authored the introduction of Dense Passage Retrieval (DPR). DPR is a state-of-the-art approach for efficiently retrieving passages relevant to a query by mapping both questions and passages into dense vector spaces using dual encoders. This method significantly improved the performance of open-domain QA systems.

TrustLLM

Funded by the European Union

# Limitations of parametric LLMs #1: Hallucinations

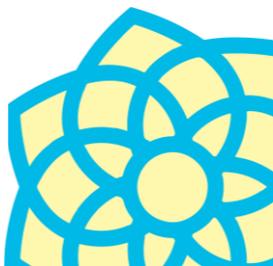**LLMs cannot encapslate everything in their parameters yet.**

Summarize Akari Asai's research including her representative papers.

Akari Asai is a prominent researcher in the field of natural language processing (NLP), with a focus on question answering (QA), information retrieval, and multi-lingual models. Her research revolves around building efficient and robust systems for extracting information from large corpora.

2. **"Dense Passage Retrieval for Ope**

most cited works, where she co-au

DPR is a sta·       art approach

mapping bo·       ns and passa

method significantly improved the performance of open-domain QA systems.

**Dense Passage Retrieval for Open-Domain Question Answering**

Vladimir Karpukhin[*], Barlas Oğuz[*], Sewon Min[†], Patrick Lewis,
Ledell Wu, Sergey Edunov, Danqi Chen[‡], Wen-tau Yih

Facebook AI        [†]University of Washington        [‡]Princeton University
{vladk, barlaso, plewis, ledell, edunov, scottyih}@fb.com
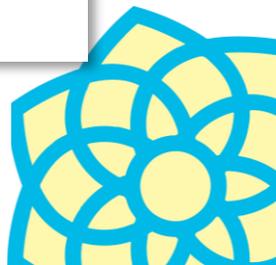sewon@cs.washington.edu
danqic@cs.princeton.edu

TrustLLM

Funded by
the European Union

# Catastrophic incidents due to LLM hallucinations

**Such LLM hallucinations have been causing many critical incidents in the real world**



TECH · LAW

**Humiliated lawyers fined $5,000 for submitting ChatGPT hallucinations in court: 'I heard about this new site, which I falsely assumed was, like, a super search engine'**
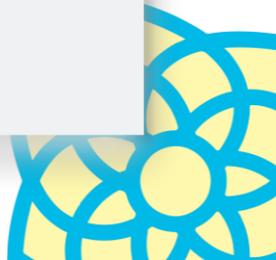
BY RACHEL SHIN
June 23, 2023 at 9:41 AM PDT

Lawyers who filed legal documents with false citations generated by ChatGPT have been fined.
ERIK MCGREGOR—LIGHTROCKET/GETTY IMAGES



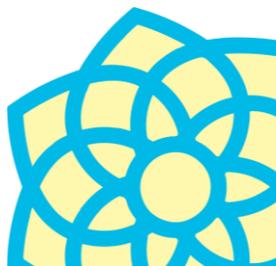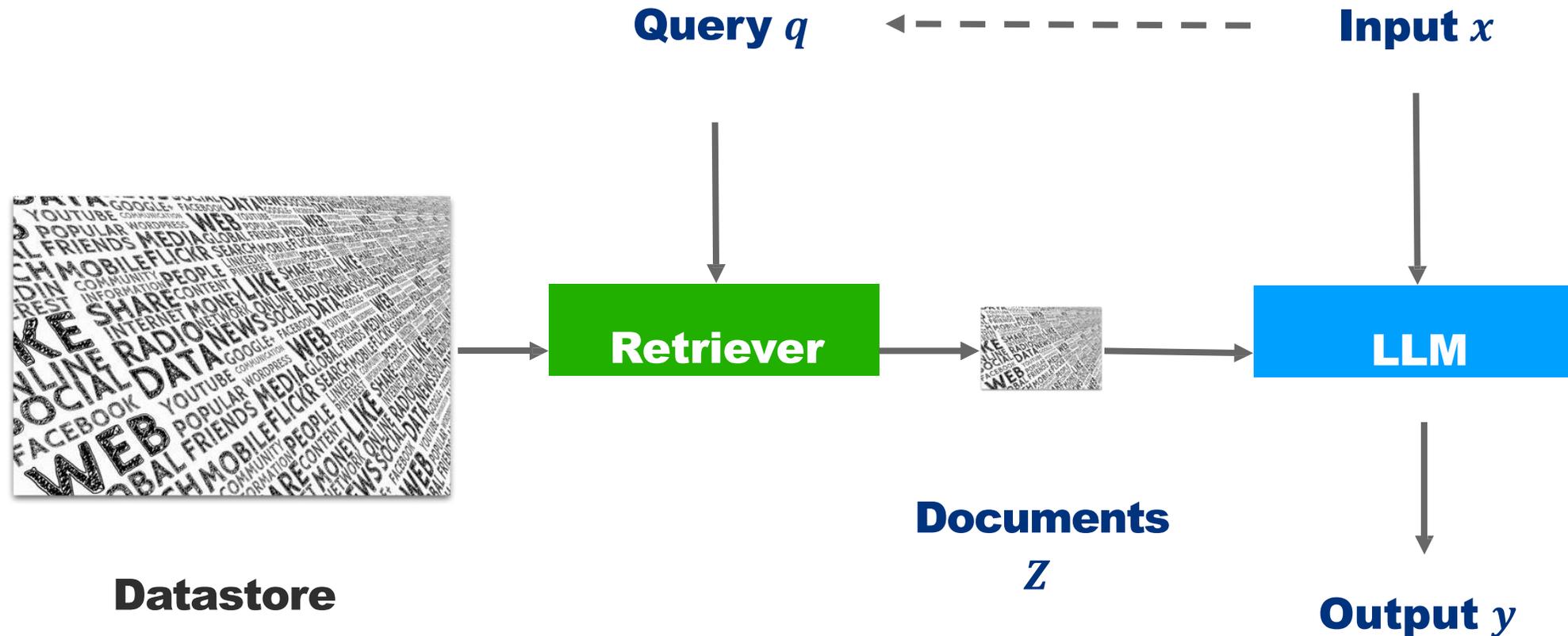# Air Canada must honor refund policy invented by airline's chatbot

Air Canada appears to have quietly killed its costly chatbot support.

ASHLEY BELANGER - 2/16/2024, 12:12 PM

TrustLLM

# Retrieval-augmented LMs: Definitions & Notations

**A new type of LMs that can use large-scale text data (datastore) at *inference-time***



Input $x$

LLM

Output $y$

Pre-training data

# Retrieval-augmented LMs: Definitions & Notations

**A new type of LMs that can use large-scale text data (datastore) at *inference-time***

**Input $x$**

**Retriever**

**LLM**

**Datastore**

# Retrieval-augmented LMs: Definitions & Notations

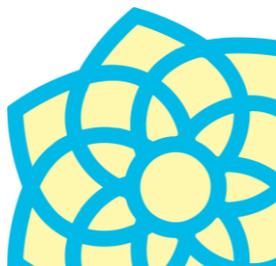**A new type of LMs that can use large-scale text data (datastore) at *inference-time***

**Query** $q$   ← ─ ─ ─ ─ ─ ─ ─ ─   **Input** $x$



**Retriever**

**LLM**

**Datastore**

# Retrieval-augmented LMs: Definitions & Notations

**A new type of LMs that can use large-scale text data (datastore) at *inference-time***

**Query** $q$          $\longleftarrow$  - - - - - - - - - -          **Input** $x$

**Retriever**

$\mathrm{sim}(q, d)$

**LLM**

**Documents** $z$

**Datastore**

# Retrieval-augmented LMs: Definitions & Notations
**A new type of LMs that can use large-scale text data (datastore) at *inference-time***

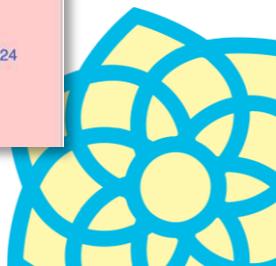# Benefit of retrieval-augmented LMs #2: Adaptations w/o training

**Parametric LMs' knowledge gets obsolete quickly & requires continuous training**

Who is the current prime minister of UK?

The current Prime Minister of the United Kingdom is **Rishi Sunak**. He has held the position since **October 25, 2022**, following the resignation of Liz Truss. Sunak is the leader of the Conservative Party and previously served as Chancellor of the Exchequer.
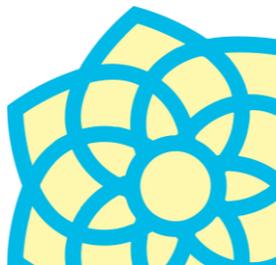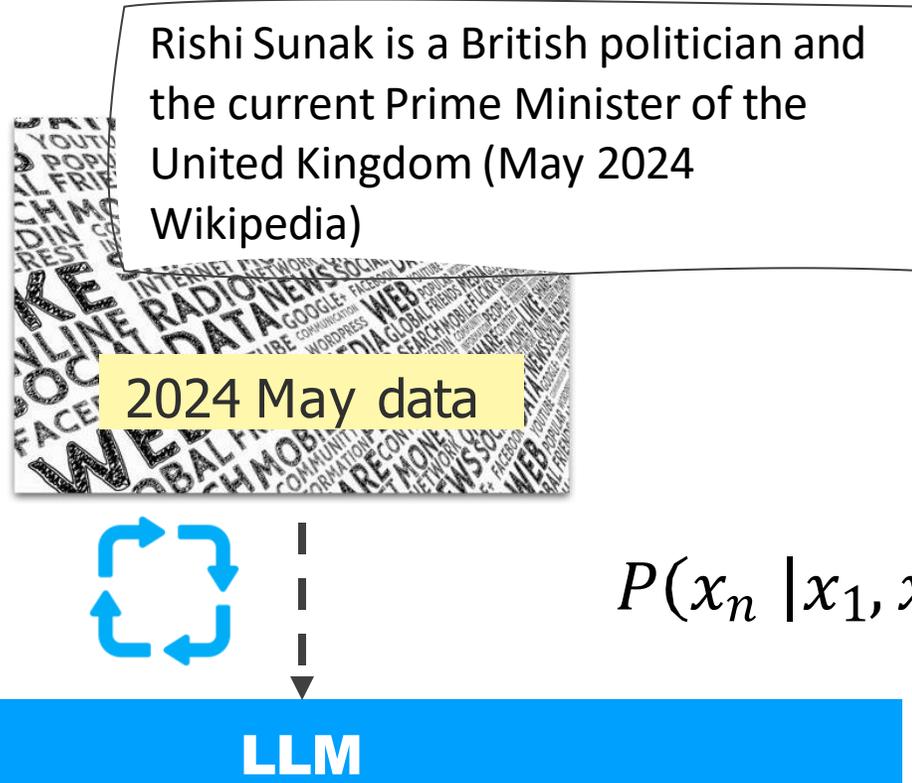
# Benefit of retrieval-augmented LMs #2: Adaptations w/o training

**Parametric LMs' knowledge gets obsolete quickly & requires continuous training**

Who is the current prime minister of UK?

The current Prime Minister of the United Kingdom is **Rishi Sunak**. He has held the position since **October 25, 2022**, following the resignation of Liz Truss. Sunak is the leader of the Conservative Party and previously served as Chancellor of the Exchequer.

| Portrait | Prime minister Office (Lifespan) | Term of office | | | Mandate[a] |
|---|---|---|---|---|---|
| | | Start | End | Duration | |
| | Rishi Sunak [98] MP for Richmond (Yorks) (born 1980) Premiership | 25 October 2022 | 5 July 2024 | 1 year, 255 days | — |
| | Keir Starmer [99] MP for Holborn and St Pancras (born 1962) Premiership | 5 July 2024 | Incumbent | 73 days | 2024 |

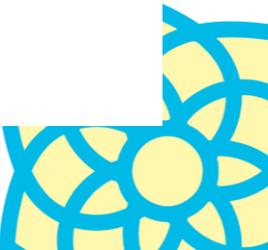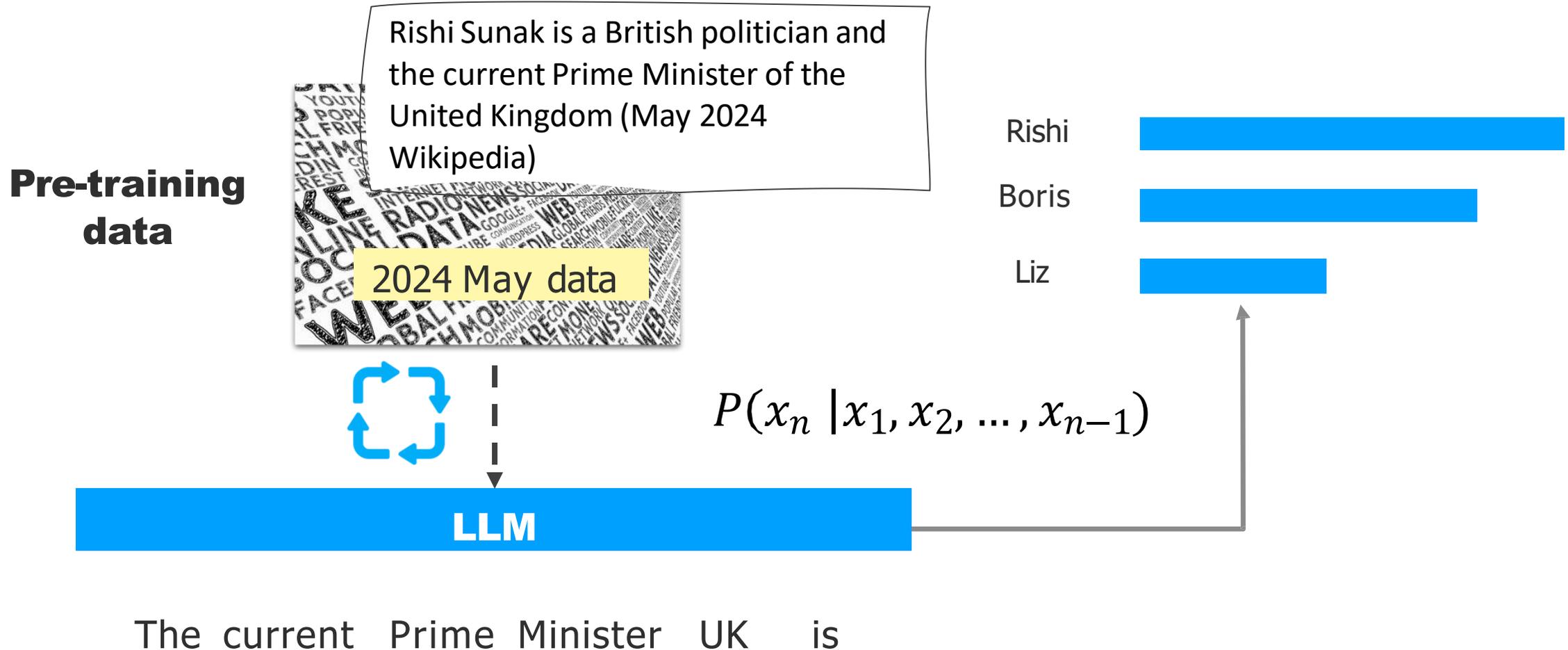# Benefit of retrieval-augmented LMs #2: Adaptations w/o training

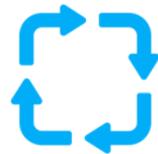**Parametric LMs' knowledge gets obsolete quickly & requires continuous training**

**Pre-training data**

Rishi Sunak is a British politician and the current Prime Minister of the United Kingdom (May 2024 Wikipedia)

2024 May data

$$P(x_n \mid x_1, x_2, \ldots, x_{n-1})$$

**LLM**

# Benefit of retrieval-augmented LMs #2: Adaptations w/o training

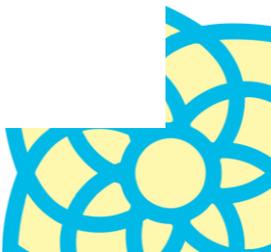**Parametric LMs' knowledge gets obsolete quickly & requires continuous training**



Rishi Sunak is a British politician and the current Prime Minister of the United Kingdom (May 2024 Wikipedia)

**Pre-training data**

2024 May data

Rishi

Boris

Liz

$$P(x_n \mid x_1, x_2, \ldots, x_{n-1})$$

LLM

The current Prime Minister UK is

# Benefit of retrieval-augmented LMs #2: Adaptations w/o training

**Parametric LMs' knowledge gets obsolete quickly & requires continuous training**
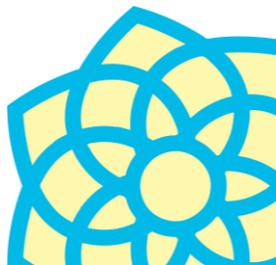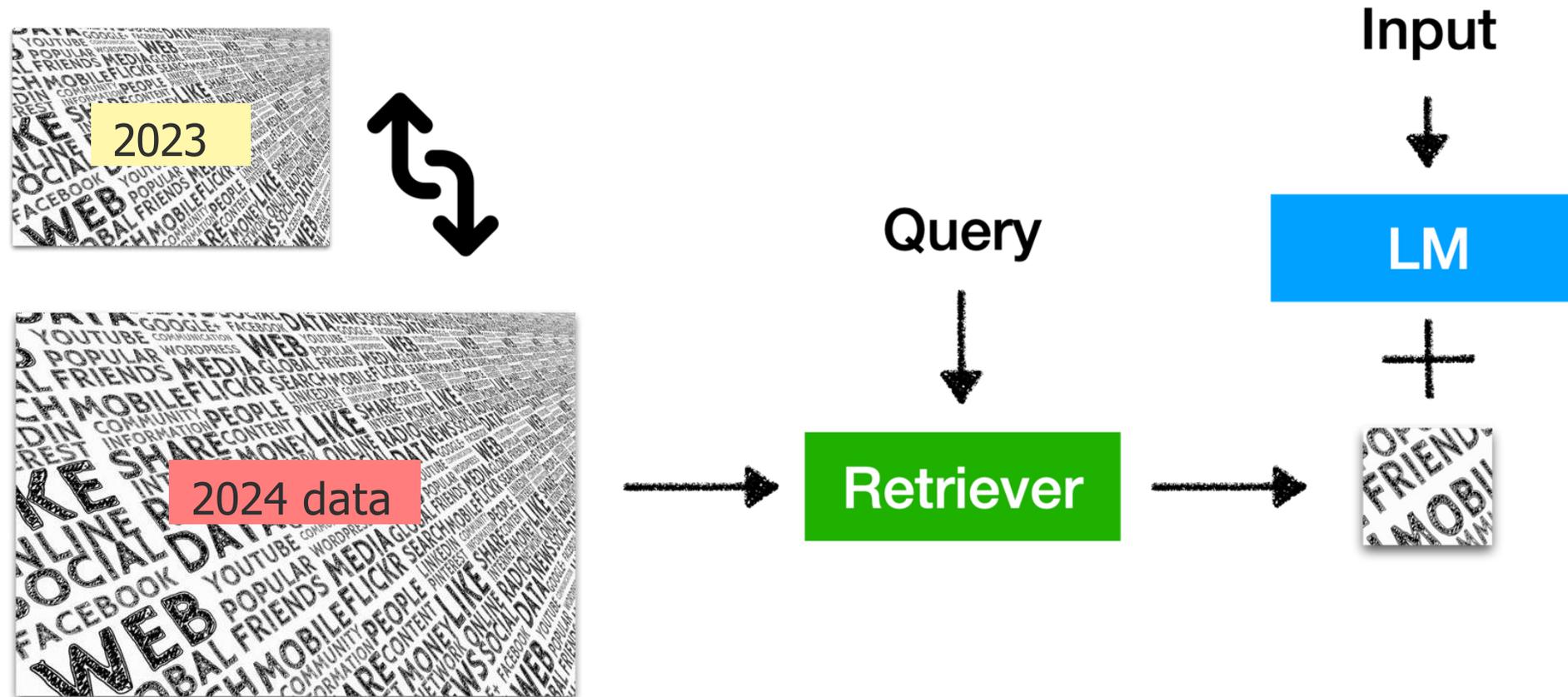
# Benefit of retrieval-augmented LMs #2: Adaptations w/o training

**We can easily swap datastores for retrieval-augmented LMs for new data distributions**

# Benefit of retrieval-augmented LMs #3: Providing attributions

**Retrieval-augmented LMs can provide a small number of documents as attributions**



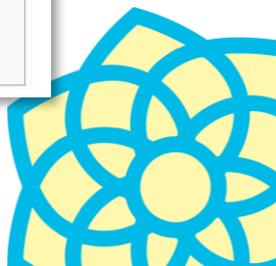Who is the current prime minister of United Kingdom?

Searched 6 sites ⌄

As of September 2024, the current Prime Minister of the United Kingdom is **Keir Starmer**, leader of the Labour Party. He took office following a general election in July 2024, where Labour secured a significant victory, unseating the Conservative Party's Rishi Sunak, who had been Prime Minister since 2022 (Northeastern Global News) (Wikipedia).
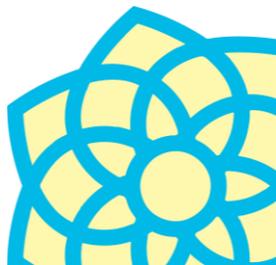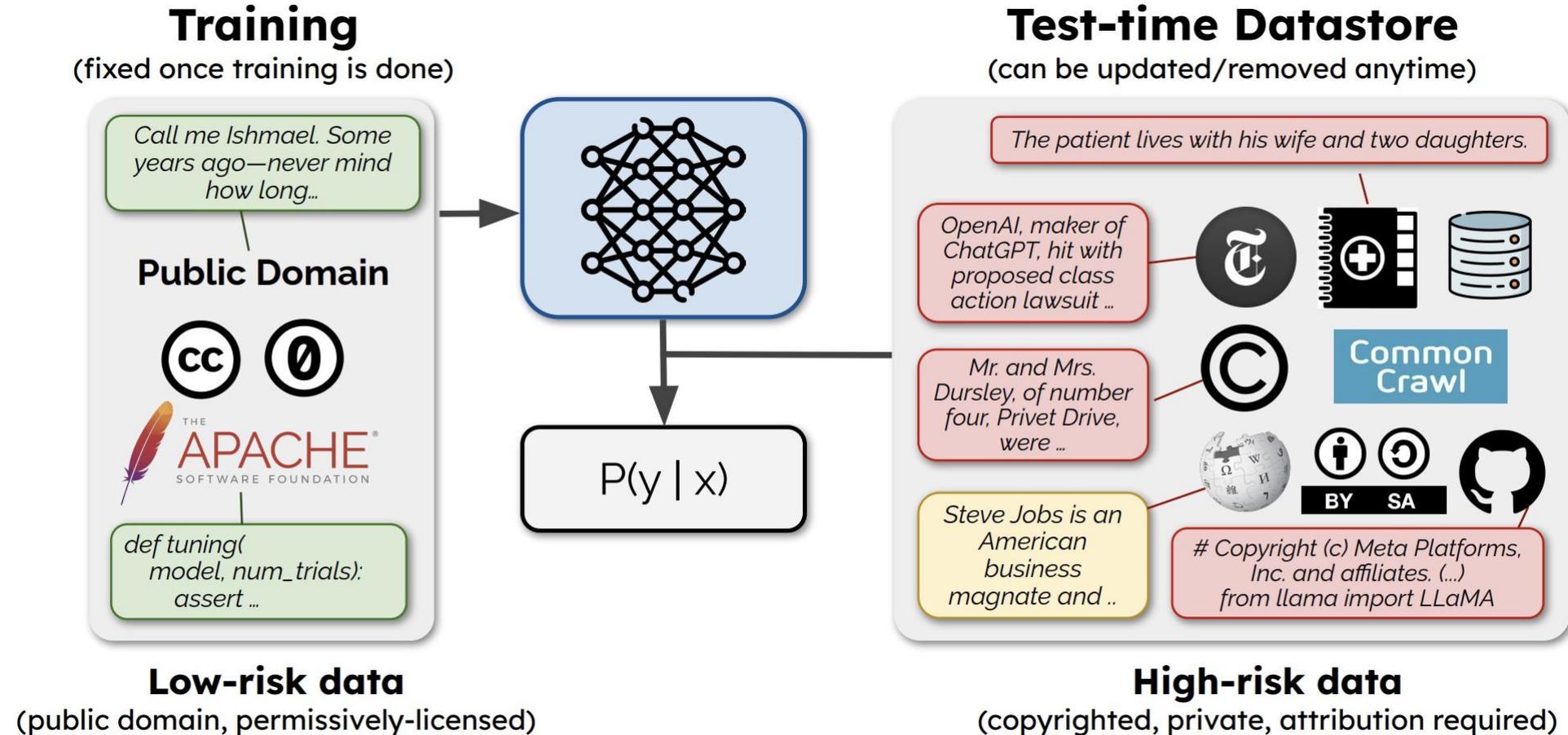
## Government of the United Kingdom [edit]

| Prime Minister | Portrait | Since | Party | Ref |
|---|---|---|---|---|
| Keir Starmer | | 5 July 2024 (2 months ago) | Labour | [1] |

# Benefit of retrieval-augmented LMs #4: Flexible data opt-in / out

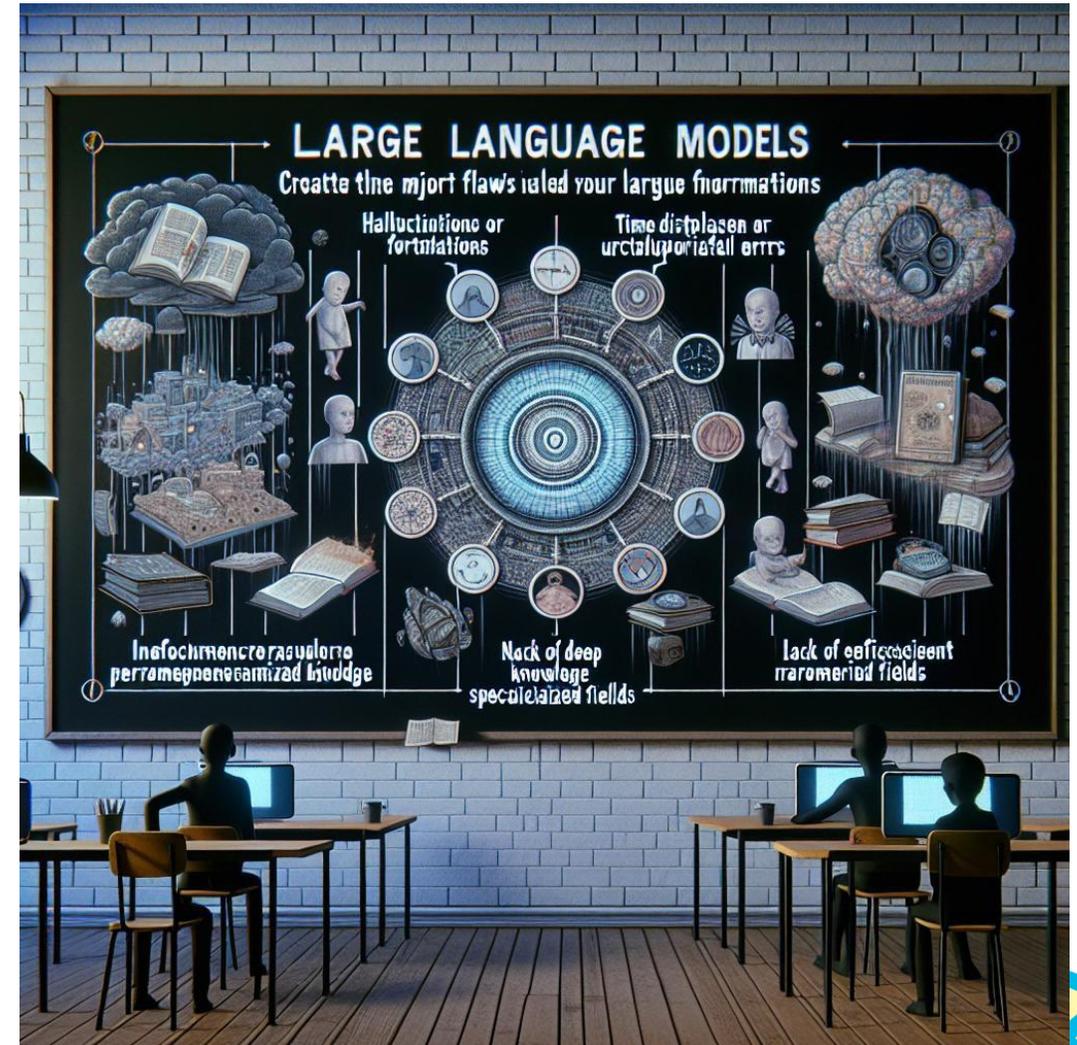**We can incorporate or remove high-risk data dynamically at inference, not training time**
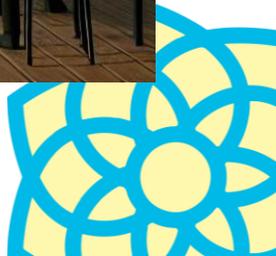
# Background

**Drawbacks of LLMs**

- Hallucination
- Outdated information
- Low efficiency in parameterizing knowledge
- Lack of in-depth knowledge in specialized domains
- Weak inferential capabilities

- Domain-specific accurate answering
- Frequent updates of data
- Traceability and explainability of generated content
- Controllable Cost
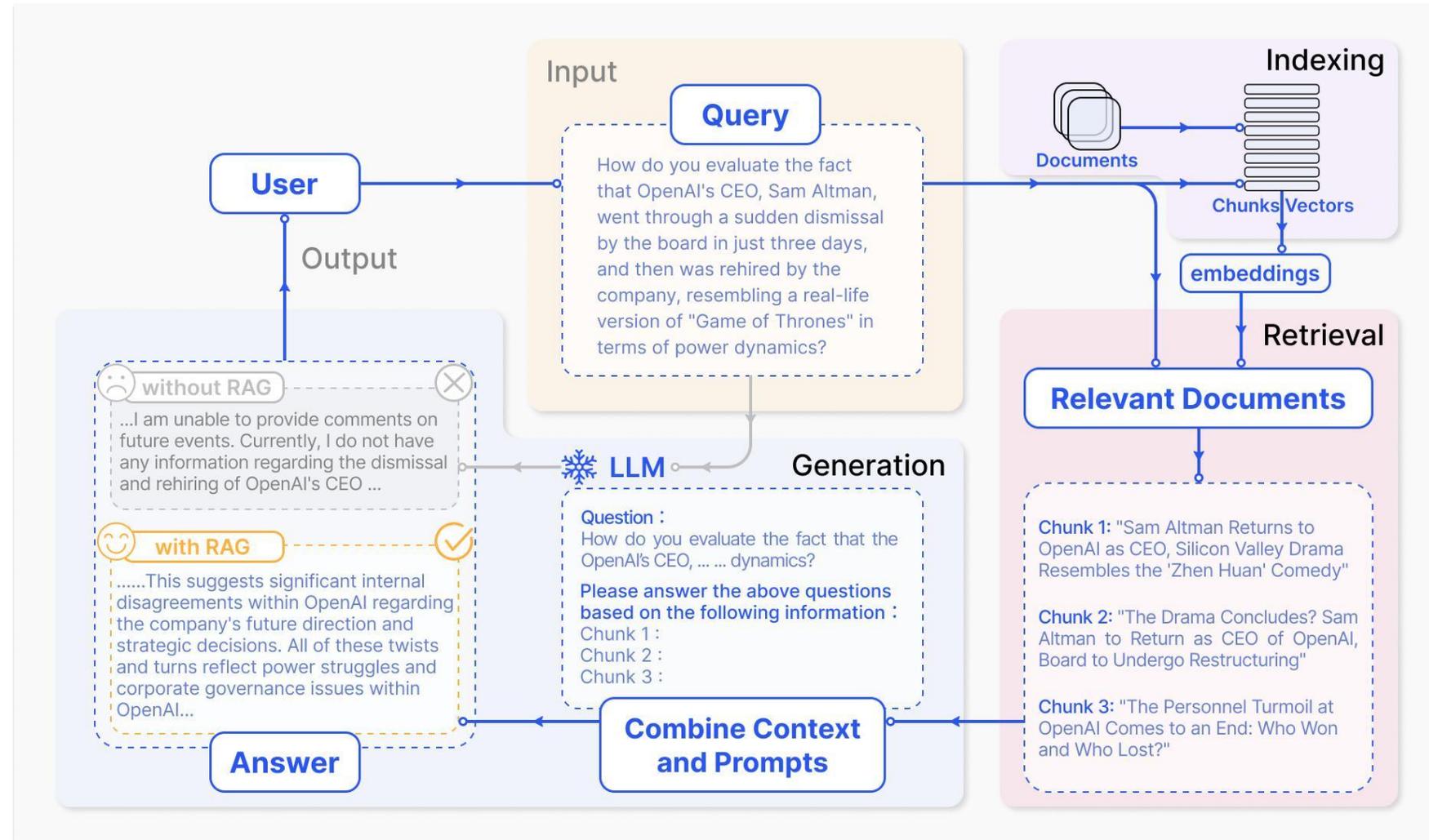- Privacy protection of data



Draw by DALL·E-3

TrustLLM

Funded by
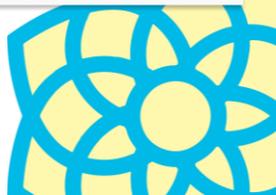the European Union

# Retrieval-Augmented Generation (RAG)

When answering questions or generating text, it first retrieves relevant information from a large number of documents, and then LLMs generates answers based on this information.

By attaching a external knowledge base, there is no need to retrain the entire large model for each specific task.

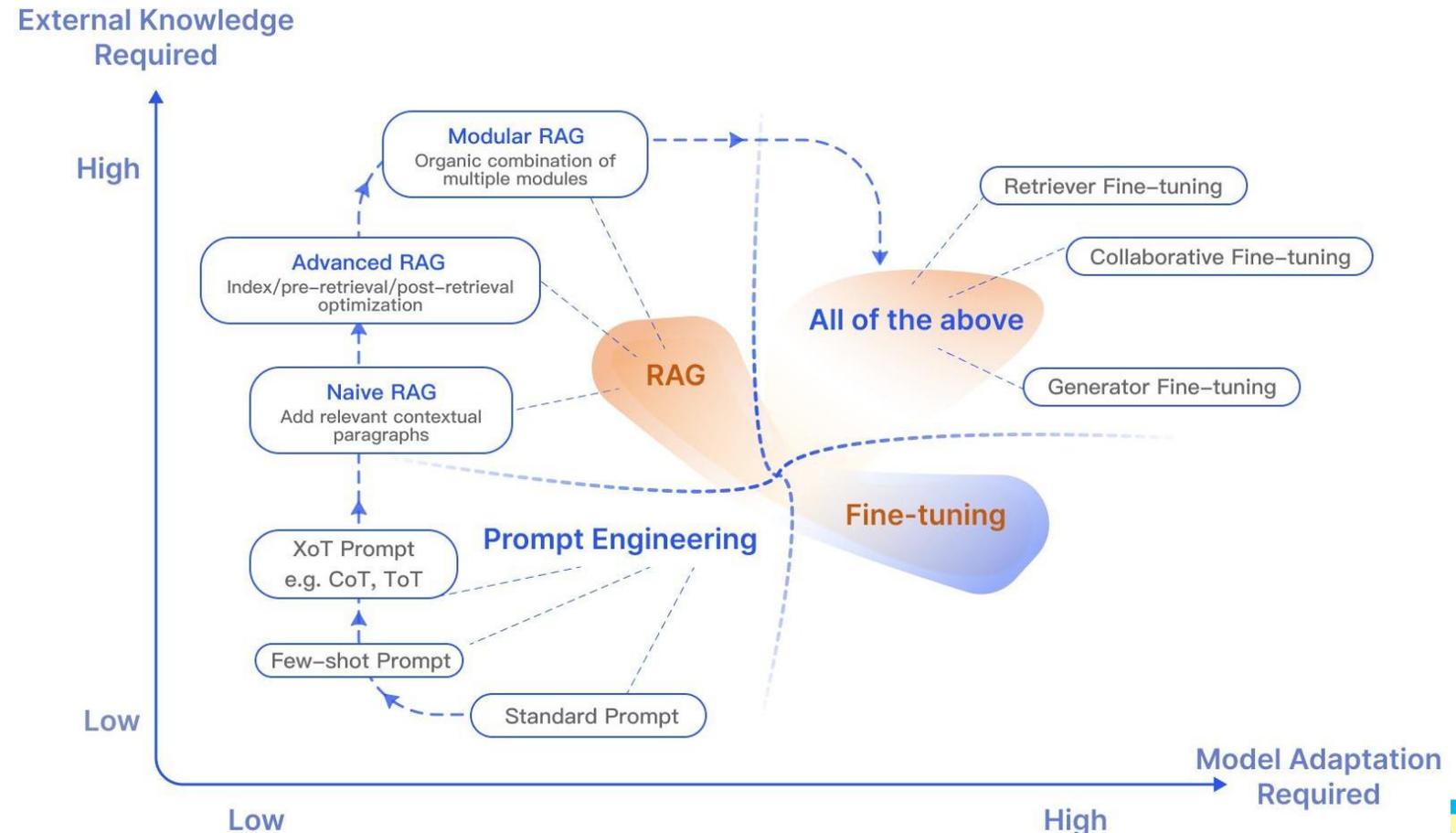The RAG model is especially suitable for knowledge-intensive tasks.



A typical case of RAG

# Symbolic Knowledge or Parametric Knowledge

Ways to optimize LLMs.

**Prompt Engineering**

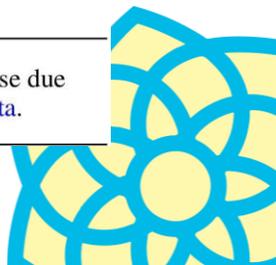**Retrieval-Augmented Generation**

**Instruct / Fine-tuning**



A typical case of RAG

# RAG vs Fine-tuning

| Feature Comparison | RAG | Fine-Tuning |
|---|---|---|
| Knowledge Updates | Directly updating the retrieval knowledge base ensures that the information remains current without the need for frequent retraining, making it well-suited for dynamic data environments. | Stores static data, requiring retraining for knowledge and data updates. |
| External Knowledge | Proficient in leveraging external resources, particularly suitable for accessing documents or other structured/unstructured databases. | Can be utilized to align the externally acquired knowledge from pretraining with large language models, but may be less practical for frequently changing data sources. |
| Data Processing | Involves minimal data processing and handling. | Depends on the creation of high-quality datasets, and limited datasets may not result in significant performance improvements. |
| Model Customization | Focuses on information retrieval and integrating external knowledge but may not fully customize model behavior or writing style. | Allows adjustments of LLM behavior, writing style, or specific domain knowledge based on specific tones or terms. |
| Interpretability | Responses can be traced back to specific data sources, providing higher interpretability and traceability. | Similar to a black box, it is not always clear why the model reacts a certain way, resulting in relatively lower interpretability. |
| Computational Resources | Depends on computational resources to support retrieval strategies and technologies related to databases. Additionally, it requires the maintenance of external data source integration and updates. | The preparation and curation of high-quality training datasets, defining fine-tuning objectives, and providing corresponding computational resources are necessary. |
| Latency Requirements | Involves data retrieval, which may lead to higher latency. | LLM after fine-tuning can respond without retrieval, resulting in lower latency. |
| Reducing Hallucinations | Inherently less prone to hallucinations as each answer is grounded in retrieved evidence. | Can help reduce hallucinations by training the model based on specific domain data but may still exhibit hallucinations when faced with unfamiliar input. |
| Ethical and Privacy Issues | Ethical and privacy concerns arise from the storage and retrieval of text from external databases. | Ethical and privacy concerns may arise due to sensitive content in the training data. |

TrustLLM

Funded by the European Union

# RAG Applications

Scenarios where RAG is applicable:

- Long-tail distribution of data

- Frequent knowledge updates

- Answers requiring verification

  and traceability

- Specialized domain knowledge

- Data privacy preservation

**Q&A**
RETRO (Borgeaud et al2021)
REALM (Gu et al, 2020)
ATLAS (Izacard et al, 2023)

**Fact Checking**
RAG (Lewis et al, 2020)
ATLAS (Izacard et al, 2022)
Evi. Generator (Asai et al, 2022a )

**Dialog**
Blender Bot 3 (Shuster et al.2022)
Internet-augmented generation
(Komeili et a., 2022)

**Summary**
FLARE (Jiang et al, 2023)

**Machine Translation**
kNN-MT (Khandelwal et al., 2020)TRIME-MT (Zhong et al., 2022)

**Code Generation**
DocPrompting (Zhou et al., 2023
Natural ProverWelleck et al., 2022)

**Natural Language Inference**
kNN-Prompt (Shi et al., 2022)
NPM (Min et al., 2023)

**Sentiment analysis**
kNN-Prompt (Shi et al., 2022)NPM (Min et al., 2023)

**Commonsense reasoning**
Raco (Yu et al, 2022)

TrustLLM

Funded by
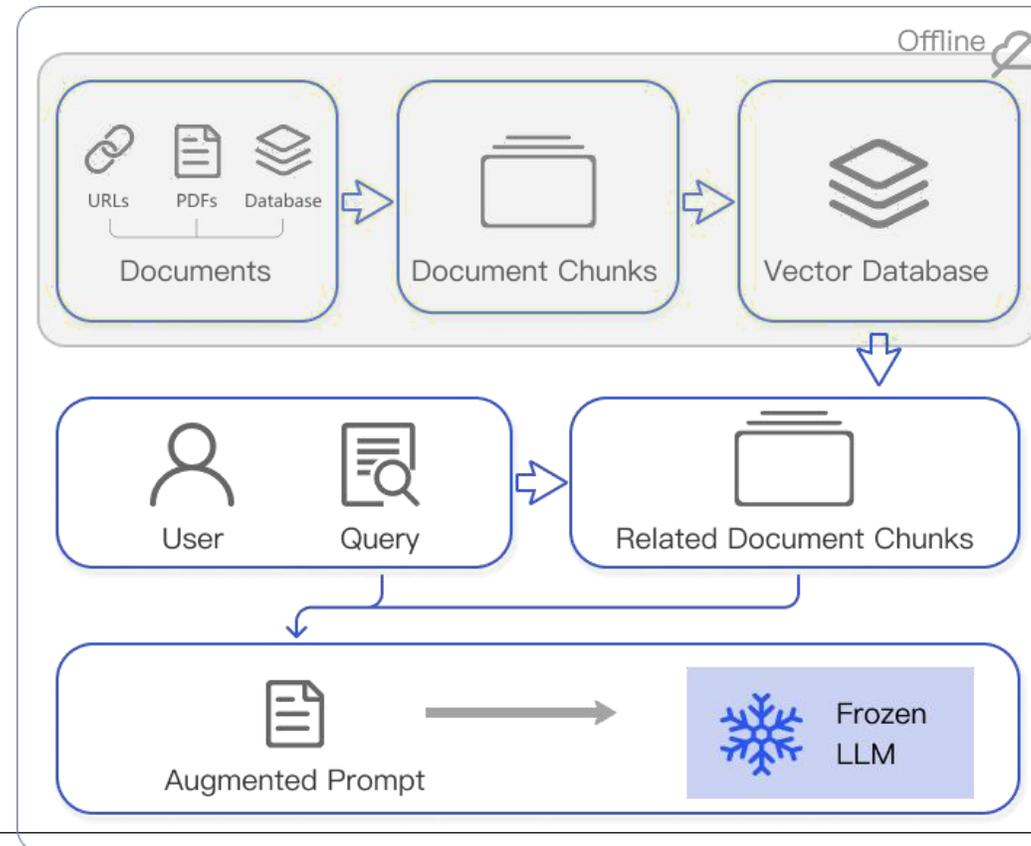the European Union

# Naive RAG

## Step1 Indexing

1. Divide the document into even chunks, each chunk being a piece of the original text.

2. Using the encoding model to generate an embedding for each chunck.

3. Store the Embedding of each block in the vector database.

## Step2 Retrival

Retrieve the k most relevant documents using vector similarity search.

## Step3 Generation

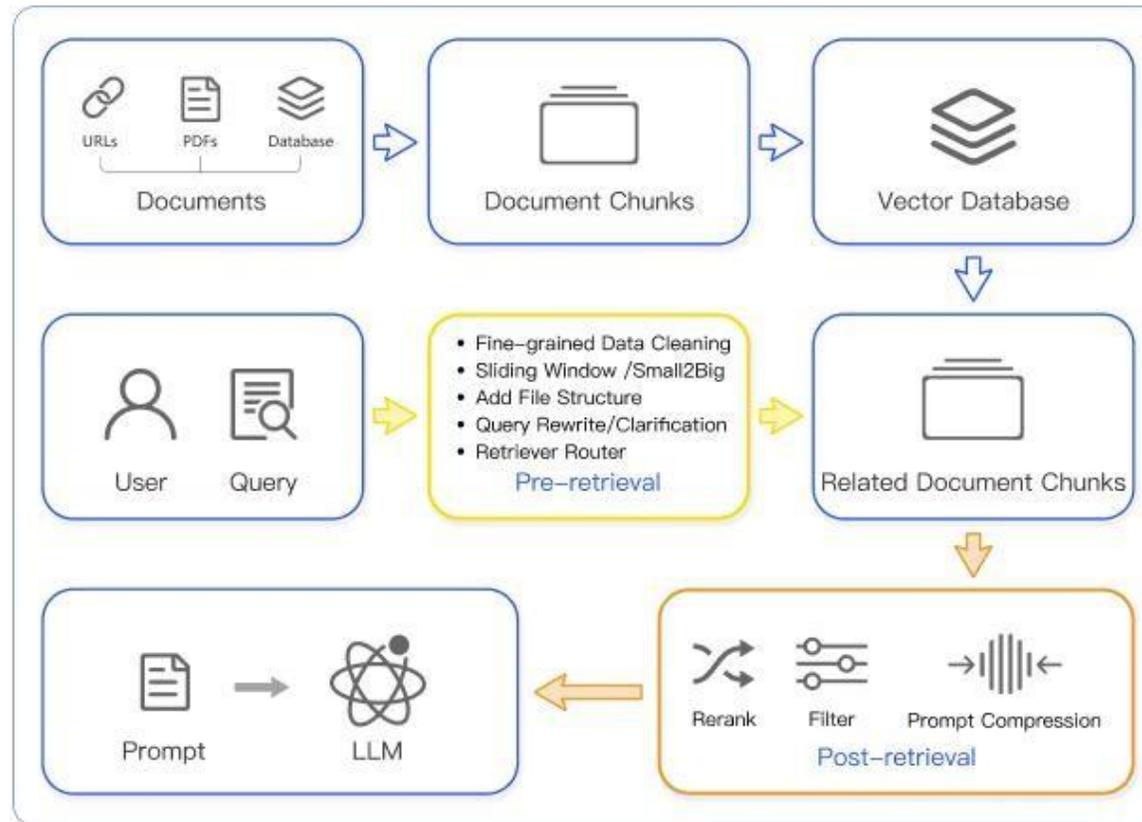The original query and the retrieved text are combined and input into a LLM to get the final answer



**Naive RAG**

**Advanced RAG**

**Modular RAG**

# Advanced RAG

Index Optimization → Pre-Retrieval Process → Retrieval →
Post-Retrieval Process→ Generation

- **Optimizing Data Indexing:**
  sliding window, fine-grained
  segmentation、adding metadata

- **Pre-Retrieval Process**：retrieve
  routes, summaries, rewriting, and
  confidence judgment

- **Post-Retrieval Process**：reorder,
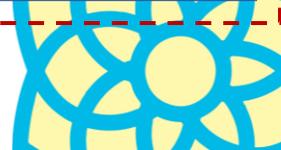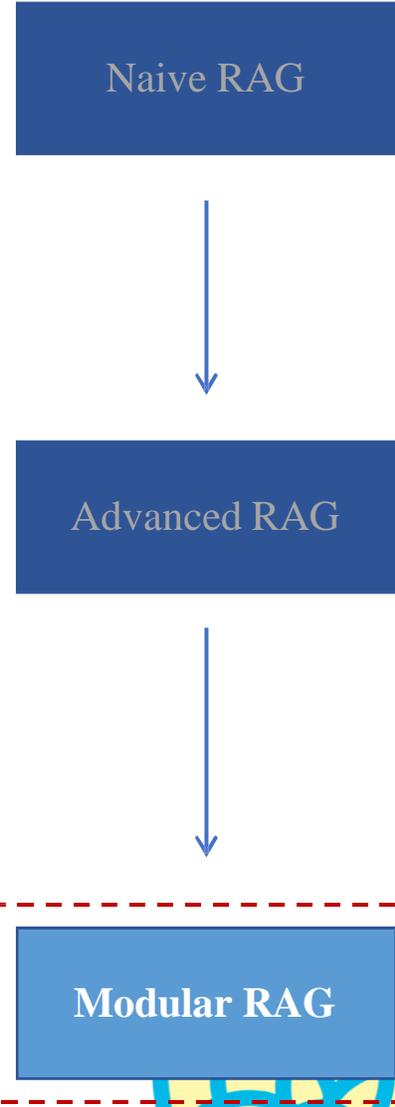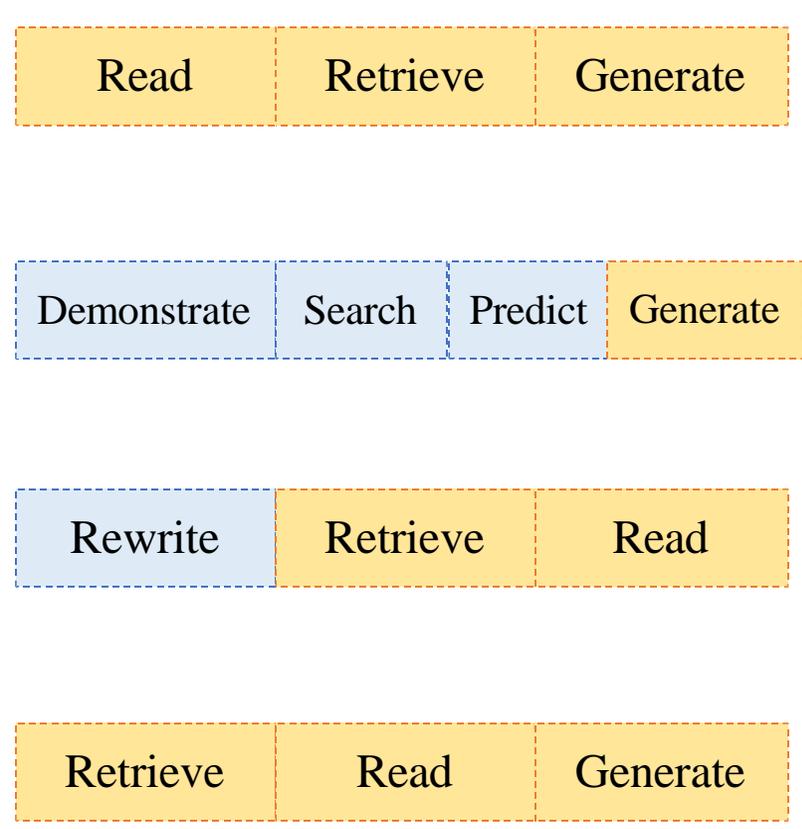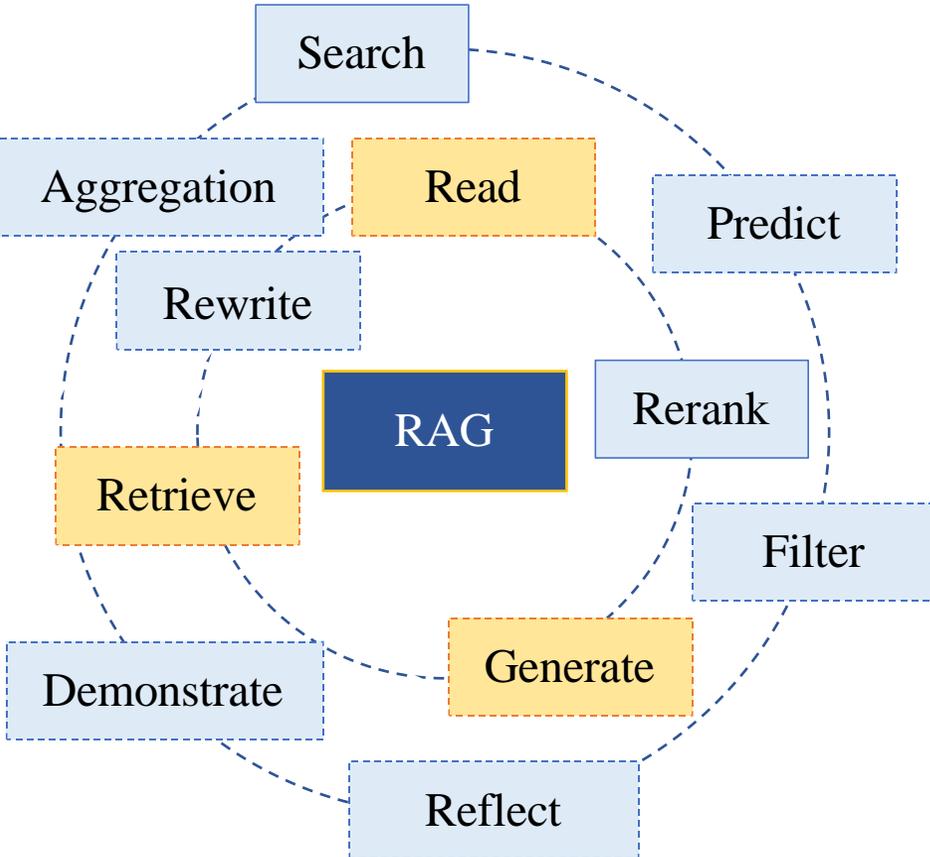  filter content retrieval



Naive RAG

Advanced RAG

Modular RAG

# Modular RAG

# Comparison of RAG Paradigms

# The three key questions of RAG

| **What to retrieve ?** | **When to retrive ?** | **How to use the retrieved information ?** |
|---|---|---|

**What to retrieve ?**

- Token
- Phrase
- Chunk
- Paragraph
- Entity
- Knowledge graph

**When to retrive ?**

- Single search
- Each token
- Every N tokens
- Adaptive search

**How to use the retrieved information ?**

- Input/Data Layer
- Model/Intermediate Layer
- Output/Prediction Layer

---

**Other Issues**

**Augmentation stage:**
- Pre-training
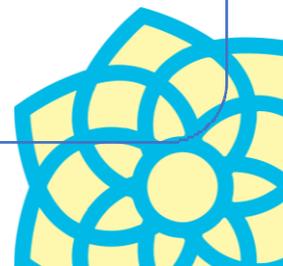- Fine-tuning
- Inference

**Retrieval choice:**
- BERT
- Roberta
- BGE
- ......

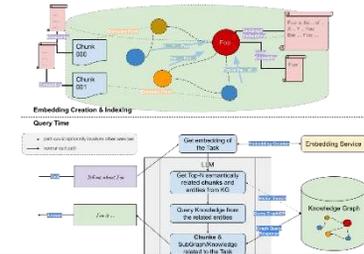Model Collaboration

⟷
⟷

Scale selectionz

**Generation choice:**
- GPT
- Llama
- T5
- ......

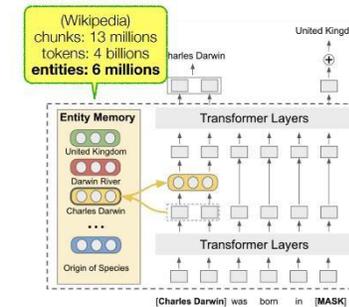# Key issue of RAG — What to retrieve



**Chunk | In-Context RAG 2023**

The search is broad, recalling a large amount of information, but with low accuracy, high coverage but includes much redundant information.
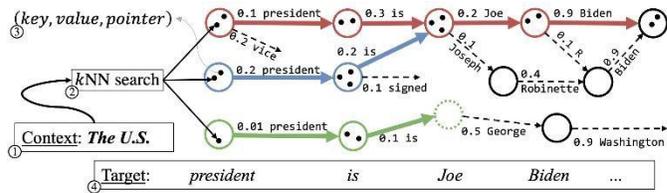
**Knowledge Graph | 2023**

**Phrase | NPM 2023**

Richer semantic and structured information, but the retrieval efficiency is lower and is limited by the quality of KG.

**Token | KNN-LMM 2019**
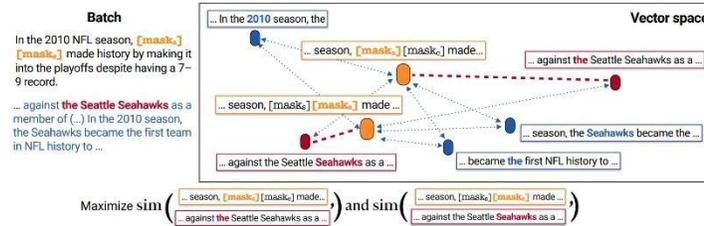
It excels in handling long-tail and cross-domain issues with high computational efficiency, but it requires significant storage.

**Entity | EasE 2022**
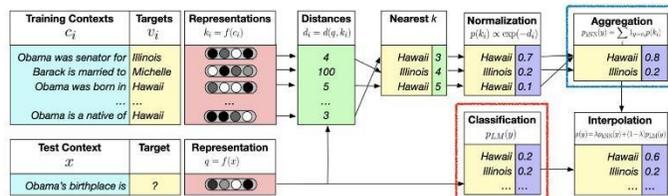
coarse

Retrieval granularity

meticulous

low                    level of structuration                    High
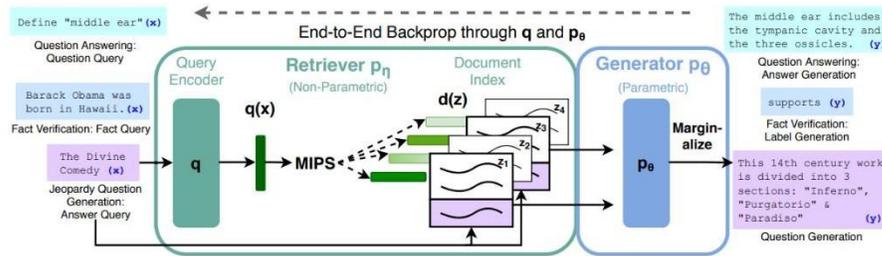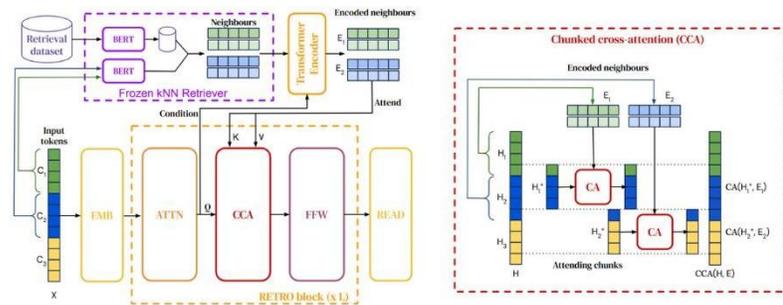
# Key issue of RAG  — How to use the retrieved content

Integrating the retrieved information into different layers of the generation model, during inference process.

Integrate retrieval positions.



Input / Date layer

Using simple, but unable to support the retrieval of more knowledge blocks, and the optimization space is limited.

Model / Interlayer

Supports the retrieval of more knowledge blocks, but introduces additional complexity and must be trained.

Output /Prediction layer

Ensuring the output results are highly relevant to the retrieval content, but the efficiency is low.

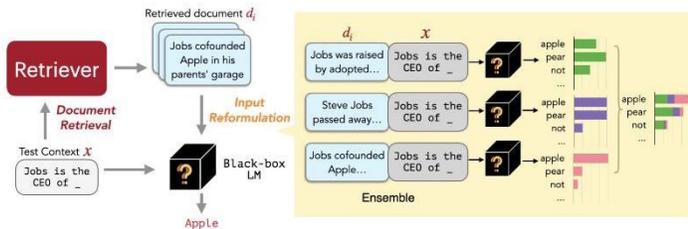TrustLLM

Funded by the European Union

# Key issue of RAG — When to retrieve



High efficiency, but low relevance of the retrieved documents

Balancing efficiency and information might not yield the optimal solution

A large amount of information with low efficiency and redundant information.

Once │ Replug 2023

Adaptive │ Flare 2023

Every N Tokens │ Atlas 2023

Conducting once search during the reasoning process.

Adaptively conduct the search.

Retrieve once for every N tokens generated.

**Low**                                                    **High**

**Retrieval frequency**

# Overview of RAG Development

# Techniques for Better RAG — Data indexing optimization

## Chunk Optimization

**Small-2-Big**

Embeding at sentence level expand the window during generation process.

**Slidingwindow**

liding chunk covers the entire text, avoiding semantic ambiguity

**Summary**

Retrieve documents through summaries, then retrieve text blocks from the documents.

## Adding Metadata

**Example**  **Page**  **Time**  **Type**  **Document Title**

## Metadata Filtering/Enrichment

**Pseudo Metadata Generation**

Enhance retrieval by gener-ating a hypothetical document for the incoming query and creating qu-estions that the text block can answer.

**Metadata filter**

Dissect and annotate the document. During the query, infer metadata filters in addition to semantic queries

Small-2-Big

Abstract

Pseudo Metadata

Metadata filter

Funded by the European Union

# Techniques for Better RAG — Structured Corpus

Hierarchical Organization of Retrieval Corpora

**● Summary → Document**

Replace document retrieval with summary retrieval, not only retrieving the most directly relevant nodes, but also exploring additional nodes associated with those nodes.

**● Document → Embedded Objects**
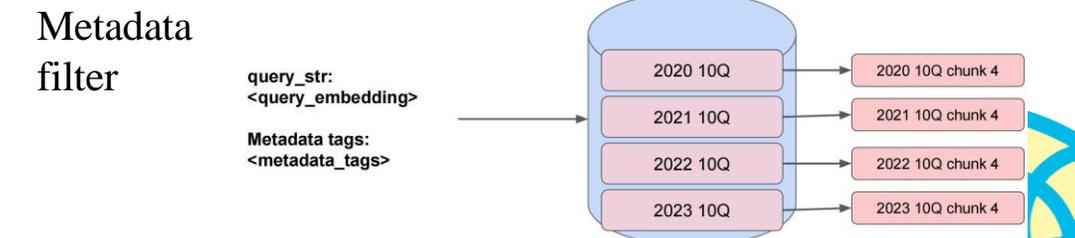
Documents have embedded objects (such as tables, charts), first retrieve entity reference objects, then query underlying objects, such as document blocks, databases, sub-nodes.

# Techniques for Better RAG — Retrieval Source Optimization



**Unstructured Data**
- Phrases
- Prompt
- Cross-lingustic

Prompt | UPRISE [Cheng et al.,2023]

Cross-language| CREA-ICL [Li et al., 2023]

**Structured Data**
- Triples
- Subgraphs

Subgraph | SUGRE [Kang et al., 2023]

**LLM**
- LLM Memory
- Generated Text
- Generated Code

Memory | Selfmem [Cheng et al., 2023]

TrustLLM

# Techniques for Better RAG — KG as a Retrieval Data Source

➢ **GraphRAG**

  ➢ Extract entities from the user's input query, then construct a subgraph to form context, and finally feed it into the large model for generation.

➢ **Implementation**

  ➢ Use LLM (or other models) to extract key entities from the question.

  ➢ Retrieve subgraphs based on entities, delving to a certain depth, such as 2 hops or even more.

  ➢ Utilize the obtained context to generate answers through LLM.

# Techniques for Better RAG — Query Optimization

Questions and answers do not always possess high semantic similarity; adjusting the Query can yield better retrieval results.



Rewrite-Retrieve-Read [Ma et al., 2023]     Tree of Clarifications（TOC）[Kim et al.,2023]

# Techniques for Better RAG – Embedding Optimization

**Selecting a More Suitable Embedding Provider**



BAAI-General-Embedding (BGE)          LLM-Embedder(BGE2) [Aksitov et al.,2023]

**Fine-tuning the Embedding Model**



Fine-tuning According to Domain-Specific Repositories and Downstream Tasks

Fine-tuning the Adapter Module to Align the Embedding Model with the Retrieval Repository

# Techniques for Better RAG — Retrieval Process Optimization

**Iterative**

**Adaptive**

Iteratively Retrieving from the Corpus to Acquire More Detailed and In-depth Knowledge

Dynamically Determined by the LLM, the Timing and Scope of Retrieva

ITER [Feng et al., 2023]



IRCOT [Trivedi et al.,2022]

FLARE [Jiang et al., 2023]



Self-RAG [Asai et al., 2023]

# Techniques for Better RAG — Hybrid (RAG + Fine-tuning)

## Retriever Fine-Tuning



Highly Adaptive
General-Purpose
Retrieval Plugin

AAR [Yu et al., 2023]

## Generator Fine-Tuning



Augment with Structural
Information Integration

SANTA [Li et al., 2023]

## Collaborative Fine-Tuning



RA-DIT [Lin et al., 2023]

- **R-FT**

  Minimizing the KL Divergence Between the Retriever Distribution and LLM Preferences

- **LM-FT**

  Maximizing the Likelihood of the Correct Answer Given Retrieval-Augmented Instructions

# Summary of Related Research



```
Retrieval-Augmented Generation
├── Retriever (§4)
│   ├── Better Semantic Representation
│   │   ├── Chunk Optimization — Small2big;Sliding-window;Abstract-Embedding;Metadata Filtering [Liu, 2023]
│   │   └── Fine-tuning Embedding Model — PROMPTAGATOR [Dai et al., 2022] ; BGE [BAAI, 2023]; LLM-Embedder [Zhang et al., 2023a] ; AngIE[Li and Li, 2023]
│   ├── Align Queries and documents
│   │   ├── Query Rewriting — Query2Doc [Wang et al., 2023d]; RRR [Ma et al., 2023a]; STEP-BACKPROMPTING[Zheng et al., 2023]; HyDE [Gao et al., 2022];TOC [Kim et al., 2023]
│   │   └── Embedding Transformation — SANTA [Li et al., 2023b]
│   └── Align Retriever and LLM
│       ├── Plugin Adapter — PKG [Luo et al., 2023]; RECOMP [Xu et al., 2023]; TokenFiltering [Berchansky et al.,2023]
│       └── LLM Supervised Training — AAR [Yu et al., 2023]; REPLUG [Shi et al., 2023] ; Atlas [Izacard et al., 2022] ; UPRISE [Cheng et al., 2023a]
├── Generator (§5)
│   ├── Post-retrieval with Frozen LLM
│   │   ├── Information Compression — PRCA [Yang et al., 2023a]; RECOMP [Xu et al., 2023] ; Filter-Reranker [Ma et al.,2023]
│   │   └── Rerank — Reranker [Brigger,2023]; QLM [Zhuang et al., 2023]
│   └── Fine-tuning LLM for RAG
│       ├── General Optimization Process — Self-Mem [Cheng et al., 2023b]
│       └── Utilizing Contrastive Learning — SUGRE [Kang et al., 2023]; SANTA [Li et al., 2023b]
└── Augmentation Method(§6)
    ├── Augmentation Stage
    │   ├── Pre-training — RETRO [Borgeaud et al., 2022]; Atlas [Izacard et al., 2022]; REALM [Arora et al., 2023]; Toolformer [Schick et al., 2023]; COG [Lan et al., 2022]; RAVEN [Huang et al., 2023]; RETRO++ [Wang et al., 2023a]; InstructRetro [Wang et al., 2023a];TIGER [Rajput et al., 2023]
    │   ├── Fine-tuning — DPR [Karpukhin et al., 2020] ; UPRISE [Cheng et al., 2023a]; FiD [Izacard and Grave, 2020]; RA-DIT [Lin et al., 2023]; Self-RAG[Asai et al., 2023]; SUGRE [Kang et al., 2023]; SANTA[Li et al., 2023b]; REPLUG [Shi et al., 2023]; AAR [Yu et al., 2023];
    │   └── Inference — KNN-LM [Khandelwal et al., 2019]; DSP [Khattab et al., 2022]; KAR [Purwar and Sundar, 2023]; PRCA [Yang et al., 2023a]; IRCOT [Trivedi et al., 2022]; GenRead [Yu et al., 2022]; ICRALM [Ram et al., 2023];PGRA [Guo et al., 2023]
    ├── Augmentation Source
    │   ├── Unstructured Data — UPRISE [Cheng et al., 2023a]; CREA-ICL [Li et al., 2023a]; COG [Lan et al., 2022]
    │   ├── Structured Data — FABULA [Ranade and Joshi, 2023]; SUGRE [Kang et al., 2023]; KnowledGPT [Wang et al., 2023e]; GraphToolformer [Zhang, 2023]
    │   └── LLM Generated Content — Self-Mem [Cheng et al., 2023b]; DSP [Khattab et al., 2022]; RECITE [Sun et al., 2022]; GenRead [Yu et al., 2022]; SKR [Wang et al., 2023f]
    └── Augmentation Process
        ├── Once Retrieval — REALM [Arora et al., 2023] ; RAG [Lewis et al., 2020]; UPRISE [Cheng et al., 2023a]; PKG [Luo et al., 2023]; LLM-R [Wang et al., 2023c] ; Atlas [Izacard et al., 2022]; REPLUG [Shi et al., 2023]; RECITE [Sun et al., 2022]
        ├── Iterative Retrieval — DSP [Khattab et al., 2022] ; Retrieve-Sample [Ren et al., 2023] ; ITER-RETGEN [Shao et al., 2023] ; ITRG [Feng et al., 2023]
        ├── Recursive Retrieval — IRCoT [Trivedi et al., 2022] ; ToC [Kim et al., 2023]
        └── Adaptive Retrieval — FLARE [Jiang et al., 2023] ; Self-RAG [Asai et al., 2023] ; RAVEN [Huang et al., 2023]
```
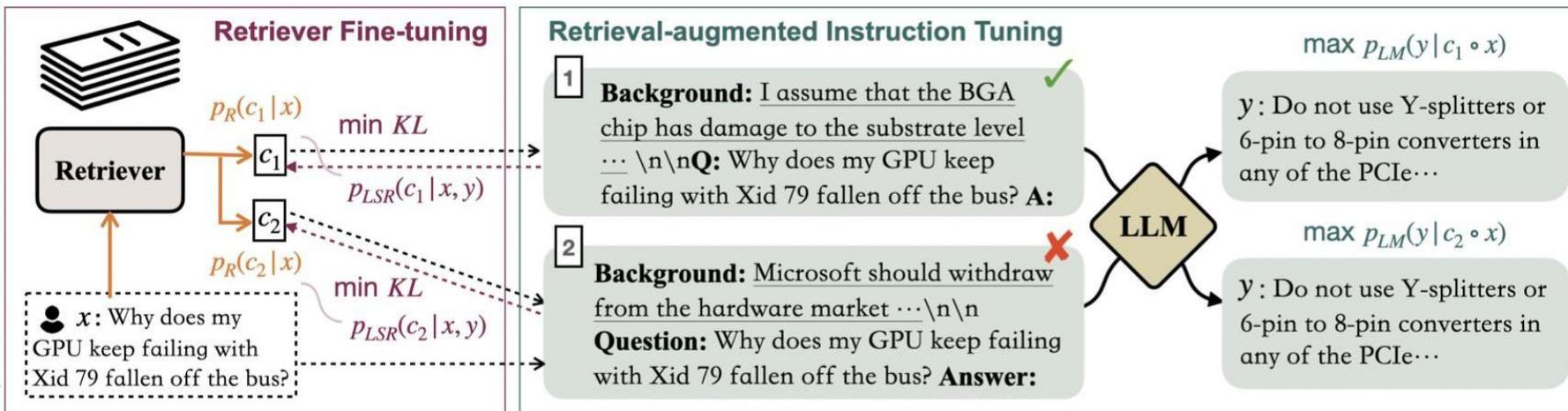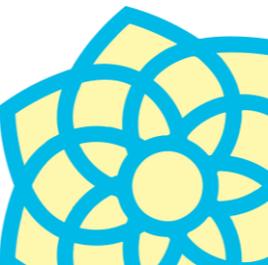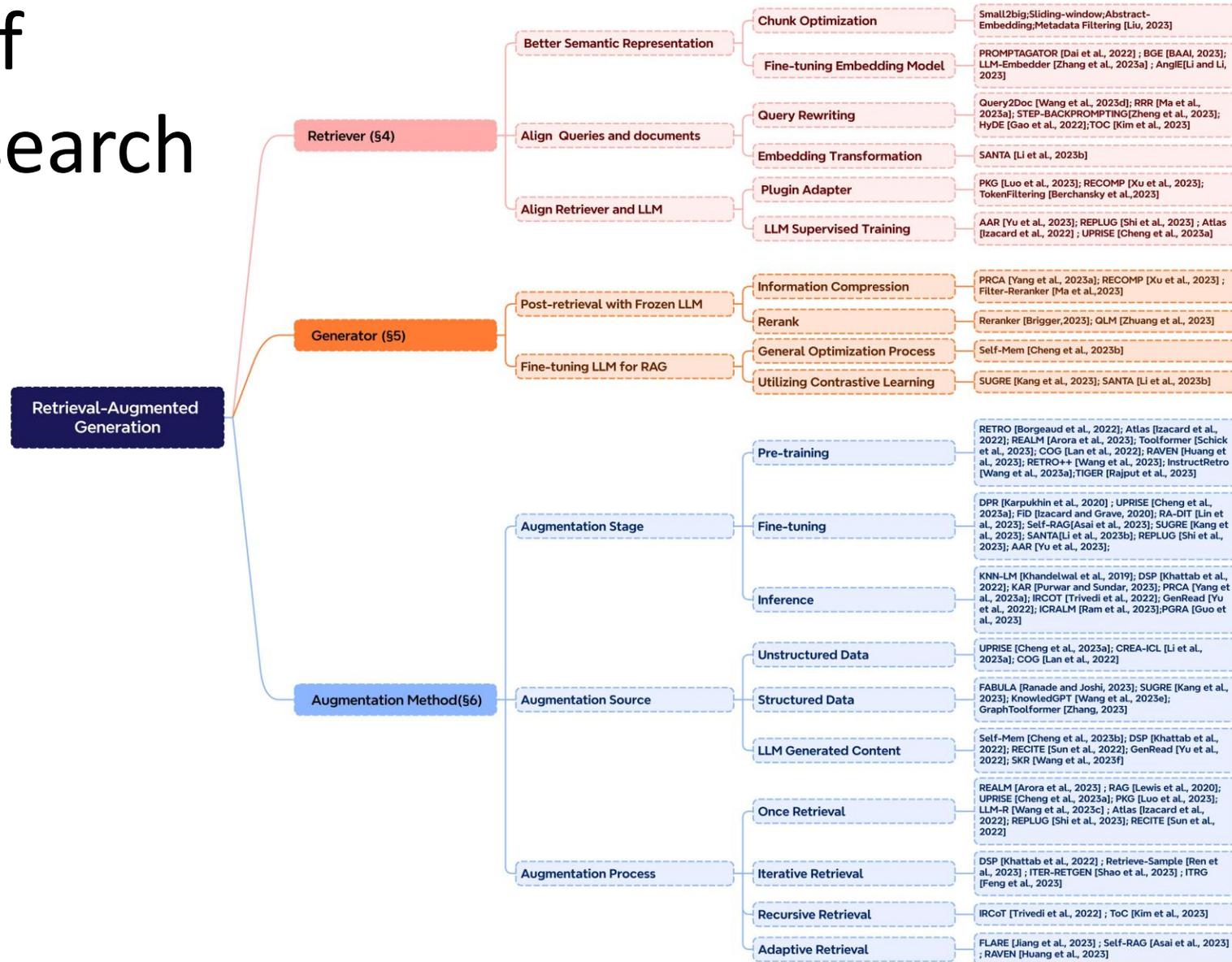
# How to Evaluate the Effectiveness of RAG

**EvaluationMethods**

**Independent Evaluation**

**End-to-End Evaluation**

Evaluate the content ultimately generated by the model.

**Retriever**
Evaluate the Quality of Text Blocks Retrieved by the Query
Metrics: MRP, Hit Rate, NDCG

**Generation/Synthesis**
Quality of Context Enhanced with Retrieved Documents Evaluation
Metrics: Context Relevance

**By generated conten**
With labels：EM，Accuracy
Without labels: Fidelity, Relevance, Harmlessness

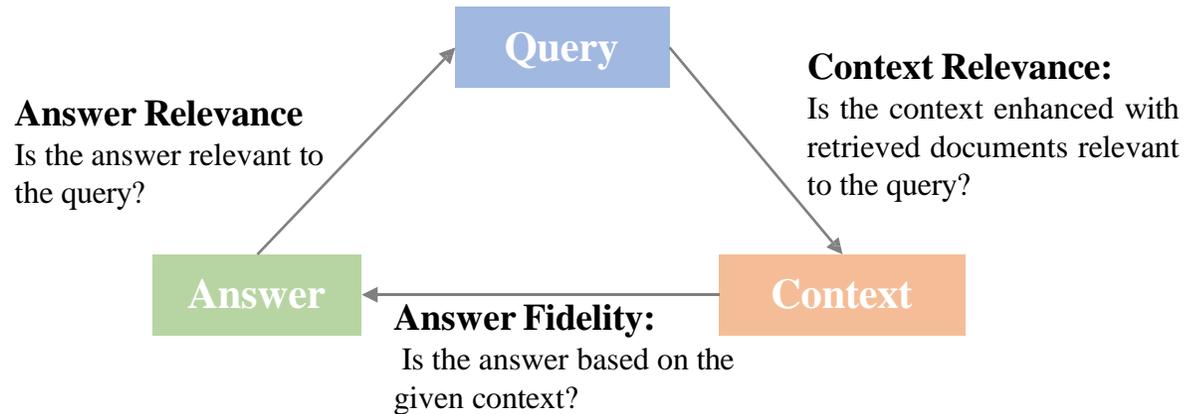**By evaluation method**
Human evaluation
Automatic evaluation (LLM judge)

**Key Metrics & Capabilities**

**Key Metrics**

**Key Capabilities**

**Query**

**Answer**

**Context**

**Answer Relevance**
Is the answer relevant to the query?

**Context Relevance:**
Is the context enhanced with retrieved documents relevant to the query?

**Answer Fidelity:**
Is the answer based on the given context?

**Noise Robustness**
Can the model extract useful information from noisy documents?

**Negative Rejection**
When therequired knowledge is not exsiting in the retrieved documents, the answer should be refused.

**Info Integration**
Can the model answer complex questions that require integrating information from multiple documents?

**Counterfactual Robustness**
Can the model recognize the risk of known factual errors in the retrieved documents?

**Assessment Framework**

**Use LLM as the adjudicator judge.**

**TruLens**        **RAGAS**        **ARES**        **Evaluation**

Based on handwritten prompt

Synthetic dataset + Fine-tuning + Ranking using confidence intervals

• **Answer Fidelity**

• **Answer Relevance**

• **Contextual Relevance**

TrustLLM

# Existing Tech Stack for RAG

| Name | Pros | Cons |
|------|------|------|
| LangChain | Modular, full-featured | Inconsistent behavior ,API conceals details,complexity and low flexibility. |
| LlamaIndex | Focus on RAG | Requires combination use, low customization. |
| FlowiseAI | Easy to get started, visualized workflows. | Does not support complex scenarios. |
| AutoGen | Adapts to multi-agent scenarios. | Low efficiency, requires multiple rounds of dialogue. |



LangChain



FlowiseAI



LlamaIndex



AutoGen

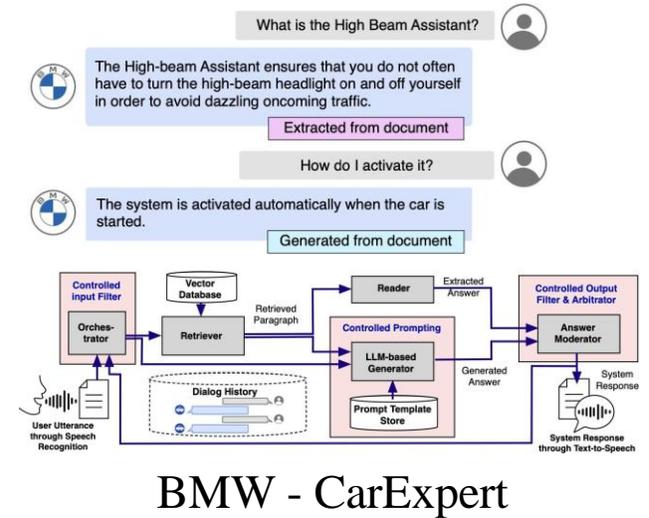# RAG Industry Application Practices
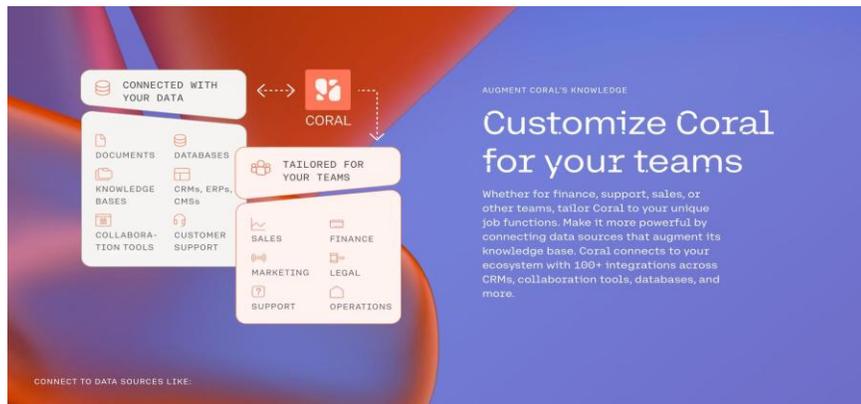


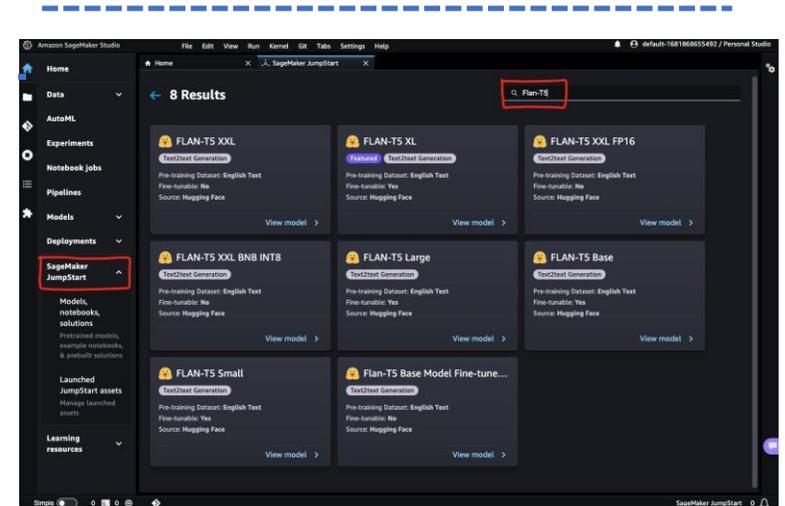NetEase - ChatBI

The intelligent upgrade of traditional industries
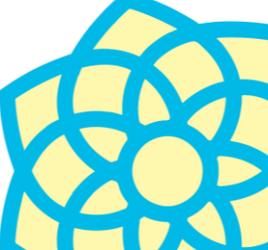
**RAG**

AI Toolchain Enhancement
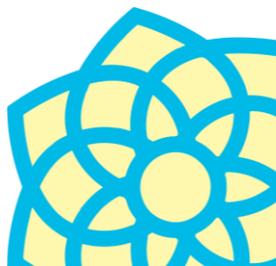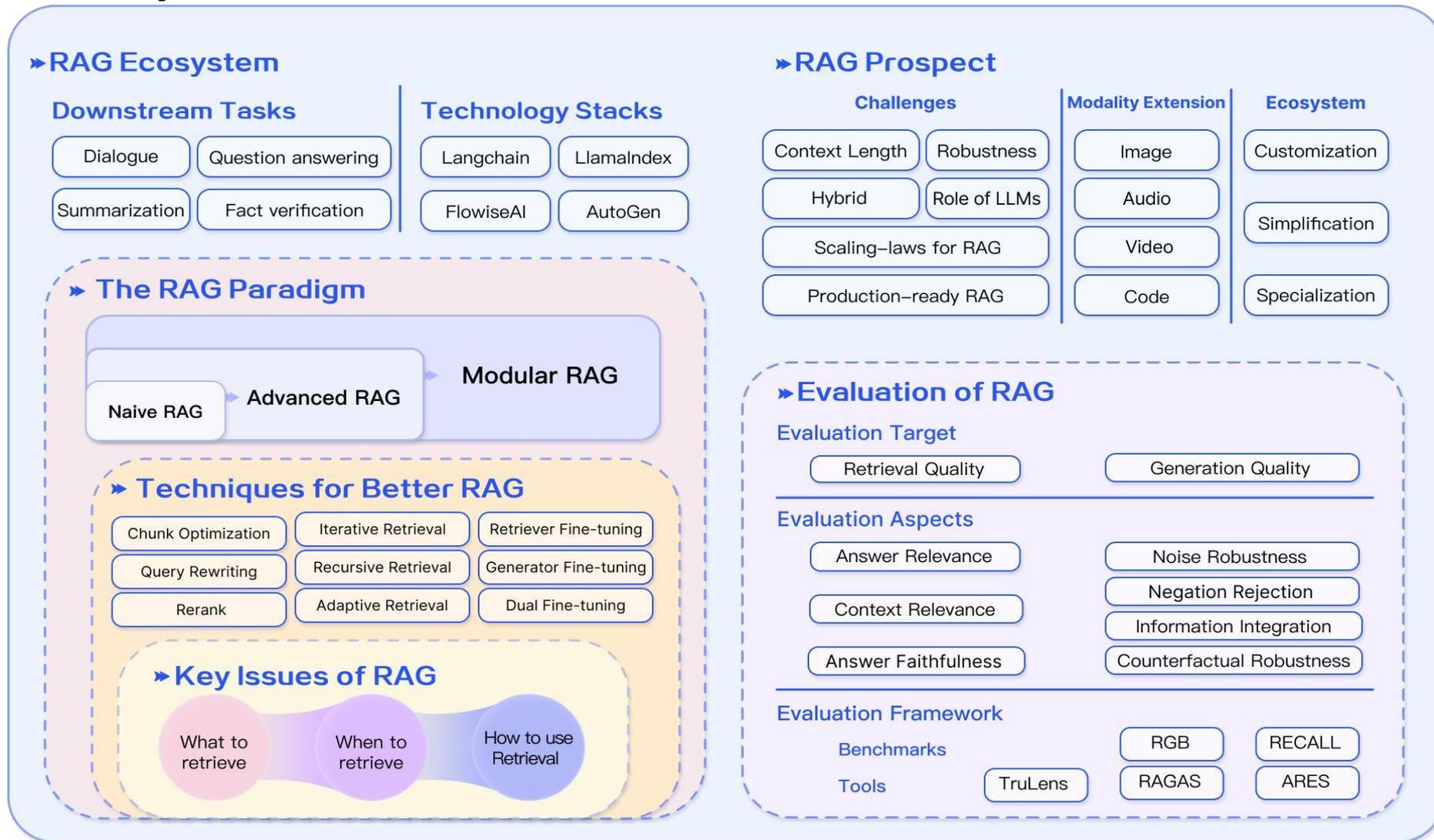


BMW - CarExpert



Cohere - Coral



Amazon - Kendra

# Summary — The Framework of RAG

# Prospects — Existing Challengs of RAG
## Further address the challenges faced by RAG itself

### Long context

- Retrieved content is excessive, exceeding window limit.
- The context is too long to result Lost in the Milddle.
- If the context window is not limited, is there still a need for RAG?

### Coordination with FT

- How to simultaneously leverage the effects of RAG and FT.
- How do the two coordinate, how are they organized, is it in Pipeline, alternating, or end-to-end?

### The role of LLMs

- LLM can be used for retrieval (LLM generation replaces retrieval, retrieving from LLM memory), for generation, and for evaluation. How to further explore the potential of LLM in RAG.

### Robustness

- How to handle the incorrect content retrieved
- How to filter and verify the content retrieved.
- How to improve the model's resistance to toxicity and noise

### Scaling Law
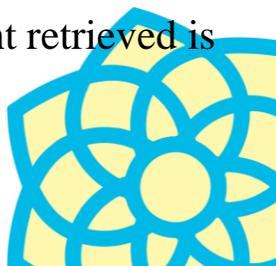
- Does the RAG model satisfy the Scaling Law
- Does RAG exhibit, or under what scenarios does it exhibit an Inverse Scaling Law
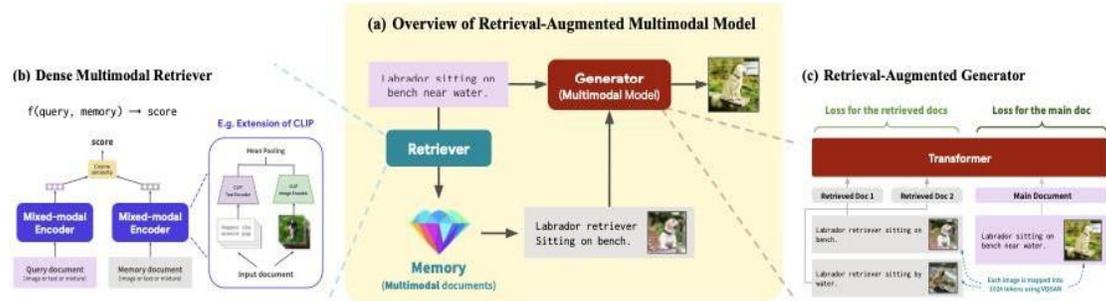
### Engineering Practice

- How to reduce the latency of retrieving ultra-large-scale corpora.
- How to ensure that the content retrieved is not leaked by large models
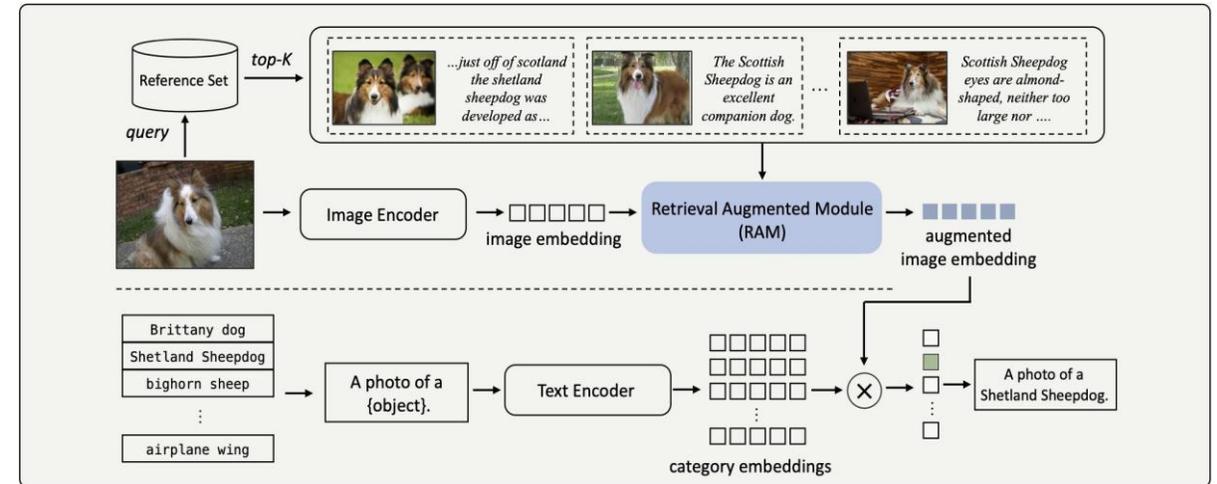
TrustLLM

Funded by the European Union

# Prospects — Mult-Modality Extension

Transferring the concept of RAG from text to other modalities of data



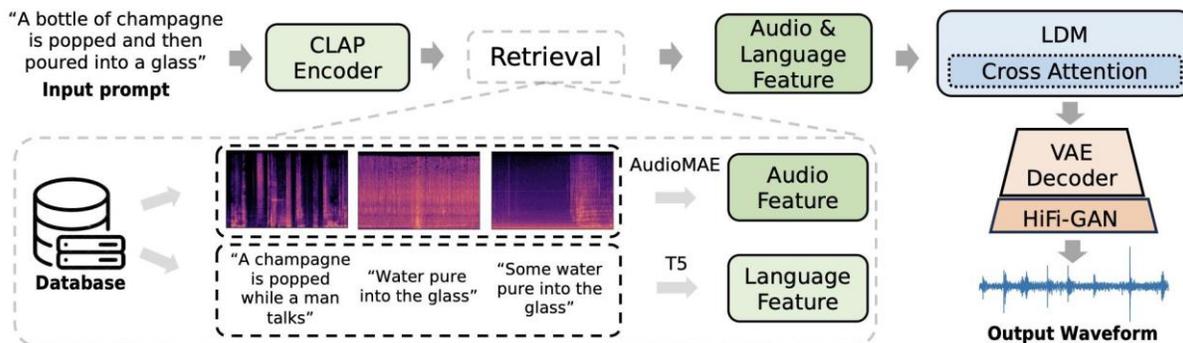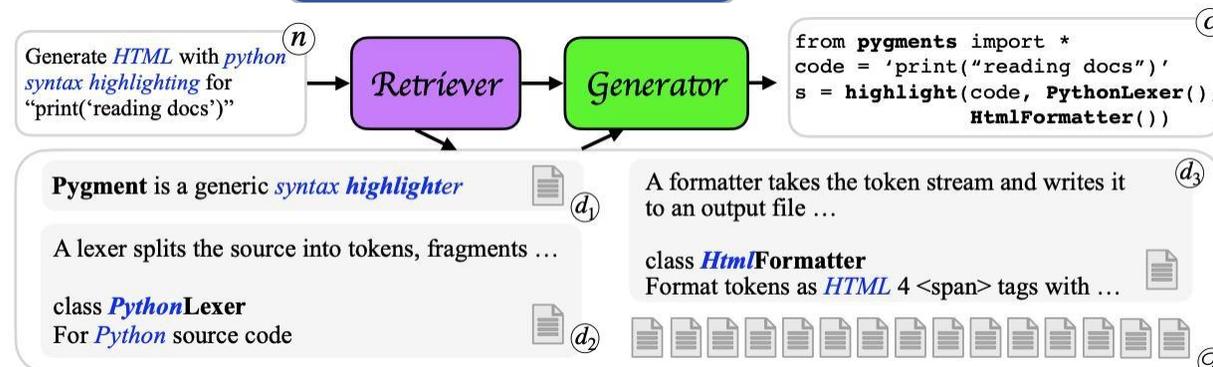RA-CM3 [Yasunaga et al.,2023]



RA-CLIP [Xie et al.,2023]

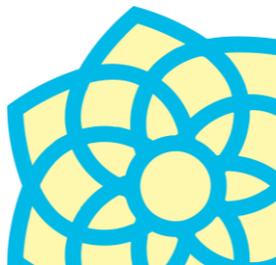

Re-AudioLDM [Yuan et al.,2023]
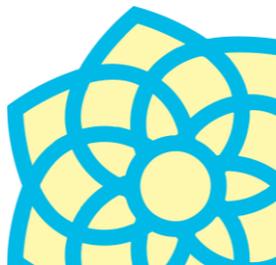


DocPrompting [Zhou et al.,2023]

# References

1.  Alon, U. et al. Neuro-Symbolic Language Modeling with Automaton-augmented Retrieval.
2.  Lewis, P. et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks.
3.  Guu, K., Lee, K., Tung, Z., Pasupat, P. & Chang, M.-W. REALM: Retrieval-Augmented Language Model Pre-Training. Preprint at http://arxiv.org/abs/2002.08909 (2020).
4.  Dai, Z. et al. Promptagator: Few-shot Dense Retrieval From 8 Examples. Preprint at http://arxiv.org/abs/2209.11755 (2022).
5.  Izacard, G. et al. Atlas: Few-shot Learning with Retrieval Augmented Language Models. Preprint at http://arxiv.org/abs/2208.03299 (2022).
6.  Gao, L., Ma, X., Lin, J. & Callan, J. Precise Zero-Shot Dense Retrieval without Relevance Labels. Preprint at http://arxiv.org/abs/2212.10496 (2022).
7.  Muennighoff, N., Tazi, N., Magne, L. & Reimers, N. MTEB: Massive Text Embedding Benchmark. in Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics 2014–2037 (Association for Computational Linguistics, 2023).
8.  Ren, Y. et al. Retrieve-and-Sample: Document-level Event Argument Extraction via Hybrid Retrieval Augmentation. in Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) 293–306 (Association for Computational Linguistics, 2023).
9.  Zhang, J. et al. ReAugKD: Retrieval-Augmented Knowledge Distillation For Pre-trained Language Models. in Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers) 1128–1136 (Association for Computational Linguistics, 2023). 10. Khattab, O. et al. Demonstrate-Search-Predict: Composing retrieval and language models for knowledge-intensive NLP. Preprint at http://arxiv.org/abs/2212.14024 (2023).
11.  Cheng, X. et al. Lift Yourself Up: Retrieval-augmented Text Generation with Self Memory. Preprint at http://arxiv.org/abs/2305.02437 (2023).
12.  Luo, Z. et al. Augmented Large Language Models with Parametric Knowledge Guiding. Preprint at http://arxiv.org/abs/2305.04757 (2023).
13.  Shi, W. et al. REPLUG: Retrieval-Augmented Black-Box Language Models. Preprint at http://arxiv.org/abs/2301.12652 (2023).
14.  Yu, Z., Xiong, C., Yu, S. & Liu, Z. Augmentation-Adapted Retriever Improves Generalization of Language Models as Generic Plug-In. Preprint at http://arxiv.org/abs/2305.17331 (2023).
15.  Kang, M., Kwak, J. M., Baek, J. & Hwang, S. J. Knowledge Graph-Augmented Language Models for Knowledge-Grounded Dialogue Generation. Preprint at http://arxiv.org/abs/2305.18846 (2023).
16.  Trivedi, H., Balasubramanian, N., Khot, T. & Sabharwal, A. Interleaving Retrieval with Chain-of-Thought Reasoning for Knowledge-Intensive Multi-Step Questions. Preprint at http://arxiv.org/abs/2212.10509 (2023).
17.  Wang, L., Yang, N. & Wei, F. Learning to Retrieve In-Context Examples for Large Language Models. Preprint at http://arxiv.org/abs/2307.07164 (2023).
18.  Li, Z. et al. Towards General Text Embeddings with Multi-stage Contrastive Learning. Preprint at http://arxiv.org/abs/2308.03281 (2023).
19.  Ng, Y. et al. SimplyRetrieve: A Private and Lightweight Retrieval-Centric Generative AI Tool. Preprint at http://arxiv.org/abs/2308.03983 (2023).
20.  Huang, J. et al. RAVEN: In-Context Learning with Retrieval Augmented Encoder-Decoder Language Models. Preprint at http://arxiv.org/abs/2308.07922 (2023).

# References

21.  Zhu, Y. et al. Large Language Models for Information Retrieval: A Survey. Preprint at http://arxiv.org/abs/2308.07107 (2023).
22.  Wang, X. et al. KnowledGPT: Enhancing Large Language Models with Retrieval and Storage Access on Knowledge Bases. Preprint at http://arxiv.org/abs/2308.11761 (2023).
23.  Chen, J., Lin, H., Han, X. & Sun, L. Benchmarking Large Language Models in Retrieval-Augmented Generation. Preprint at http://arxiv.org/abs/2309.01431
24.  Es, S., James, J., Espinosa-Anke, L. & Schockaert, S. RAGAS: Automated Evaluation of Retrieval Augmented Generation. Preprint at http://arxiv.org/abs/2309.15217 (2023).
25.  Yoran, O., Wolfson, T., Ram, O. & Berant, J. Making Retrieval-Augmented Language Models Robust to Irrelevant Context. Preprint at http://arxiv.org/abs/2310.01558 (2023).
26.  Feng, Z., Feng, X., Zhao, D., Yang, M. & Qin, B. Retrieval-Generation Synergy Augmented Large Language Models. Preprint at http://arxiv.org/abs/2310.05149 (2023).
27.  Zheng, H. S. et al. Take a Step Back: Evoking Reasoning via Abstraction in Large Language Models. Preprint at http://arxiv.org/abs/2310.06117 (2023).
28.  Cheng, D. et al. UPRISE: Universal Prompt Retrieval for Improving Zero-Shot Evaluation. Preprint at http://arxiv.org/abs/2303.08518 (2023).
29.  Wang, B. et al. InstructRetro: Instruction Tuning post Retrieval-Augmented Pretraining. Preprint at http://arxiv.org/abs/2310.07713 (2023).
30.  Jiang, Z. et al. Active Retrieval Augmented Generation. Preprint at http://arxiv.org/abs/2305.06983 (2023).
31.  Gou, Q. et al. Diversify Question Generation with Retrieval-Augmented Style Transfer. Preprint at http://arxiv.org/abs/2310.14503 (2023).
32.  Ma, X., Gong, Y., He, P., Zhao, H. & Duan, N. Query Rewriting for Retrieval-Augmented Large Language Models. Preprint at http://arxiv.org/abs/2305.14283 (2023).
33.  Yang, H. et al. PRCA: Fitting Black-Box Large Language Models for Retrieval Question Answering via Pluggable Reward-Driven Contextual Adapter. Preprint at http://arxiv.org/abs/2310.18347 (2023).
34.  Kim, G., Kim, S., Jeon, B., Park, J. & Kang, J. Tree of Clarifications: Answering Ambiguous Questions with Retrieval-Augmented Large Language Models. Preprint at http://arxiv.org/abs/2310.14696 (2023).
35.  Shao, Z. et al. Enhancing Retrieval-Augmented Large Language Models with Iterative Retrieval-Generation Synergy. Preprint at http://arxiv.org/abs/2305.15294 (2023).
36.  Zhang, P., Xiao, S., Liu, Z., Dou, Z. & Nie, J.-Y. Retrieve Anything To Augment Large Language Models. Preprint at http://arxiv.org/abs/2310.07554 (2023).
37.  Purwar, A. & Sundar, R. Keyword Augmented Retrieval: Novel framework for Information Retrieval integrated with speech interface. Preprint at http://arxiv.org/abs/2310.04205 (2023).
38.  Lin, X. V. et al. RA-DIT: Retrieval-Augmented Dual Instruction Tuning. Preprint at http://arxiv.org/abs/2310.01352 (2023).
39.  Yu, W. et al. Chain-of-Note: Enhancing Robustness in Retrieval-Augmented Language Models. Preprint at http://arxiv.org/abs/2311.09210 (2023).

# Roadmap for more efficient & reliable retrieval-augmented LMs

## Challenges of scaling up datastores & increased inference-time costs

Evaluations

Algorithms

**Infrastructure**

- Performance gains are achieved by scaling up the datastore to trillions of tokens
- Significantly increases inference costs, including CPU memory and storage requirements (e.g., 24 TB for 1.7 trillion-token).



TrustLLM

"Scaling Retrieval-Based Language Models with a Trillion-Token Datastore."
Shao, He, Asai et al., ArXiv 2024.

Funded by
the European Union

# Roadmap for more efficient & reliable retrieval-augmented LMs

**New algorithms & arthictectures to enable more efficient and effective RAG**

**Evaluations**

**Algorithms**

**Infrastructure**

- Current "RAG" has many issues such as efficiency & redundancy
- Alternative algorithms, better LM architectures, caching …etc for improving efficiency and performance



"Generative Representational Instruction Tuning." Muennighoff et al., ArXiv 2024.

TrustLLM

Funded by the European Union

# Roadmap for more efficient & reliable retrieval-augmented LMs

**New algorithms & archictectures to enable more efficient and effective RAG**

Evaluations

Algorithms

Infrastructure

- Current "RAG" has many issues such as efficiency & redundancy
- Alternative algorithms, better LM architectures, caching …etc for improving efficiency and performance



"PipeRAG: Fast Retrieval-Augmented Generation via Algorithm-System Co-design."
Jiang et al., ArXiv 2024.

# Roadmap for more efficient & reliable retrieval-augmen

**Careful analyses on their effectiveness and limitations**

**Evaluations**

**Algorithms**

**Infrastructure**

Prior systems are often evaluated only on simple general-domain tasks. Further exploration into their evaluation are needed

- **Domains**: most prior evaluations are in general-domain tasks, where Wikipedia is a sufficient knowledge source

- **Tasks**: going beyond open-domain QA, multiple-choice QA

- **Aspects**: instead of merely evaluating final "correctness", more holistic evaluations of different aspects of RAG

TrustLLM

Funded by
the European Union

# In-Context Learning

# Strategies making a pre-trained LM do a task you care about:

- In-context learning



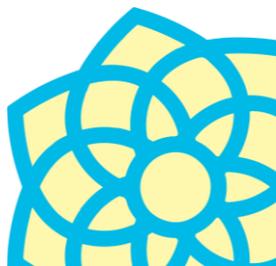- Full model finetuning → parameter-efficient finetuning



- Multi-task finetuning → instruction finetuning → alignment training



"Finetuned Language Models are Zero-Shot Learners." 2022.
https://openreview.net/forum?id=gEZrGCozdqRa

# Improving In-Context Learning

**Additional Examplars**

# Improving In-Context Learning

## Calibrate Before Use

- Step 1: Estimate the bias
  - This does not require any labeled data.
  - For classification tasks, compute normalized scores of labels
  - For generation tasks: compute probabilities of the first token of the generation over the entire vocabulary
- Step 2: Counter the bias
  - "Calibrate" the model's predictions with an affine transformation of the logits.
  - $logits_{\text{calibrated}} = \text{softmax}(\mathbf{W}logits + \boldsymbol{b})$ where W is a diagonal matrix that scales each logit to reduce bias.
- *More details in paper linked below.*

---

**Example**

Step 1:

Suppose we are building a prompt for sentiment classification, and we have decided on the template

`Input: Subpar acting. Sentiment: Negative`
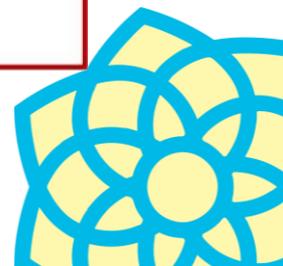
`Input: Beautiful film. Sentiment: Positive`

`Input: <query> Sentiment:`

Prompt the model using `<query>=N/A`.

Model might say P(Positive) = .618 and P(Negative) = .782

Step 2:

Set **W** and ***b*** such that P(Positive) = P(negative) = 0.5

---

"Calibrate Before Use: Improving Few-Shot Performance of Language Models." Zhao et al. 2021.
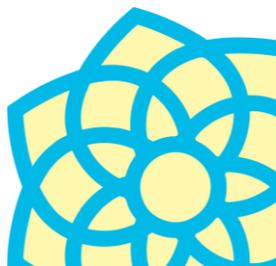
TrustLLM

**Funded by the European Union**

# Improving In-Context Learning

**Multi-Step Reasoning**

Intuition: An LLM will be better able to perform tasks (especially reasoning-based ones) if it is made to break down the task into multiple small steps.
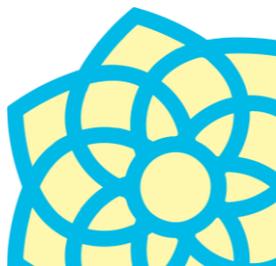
Examples of reasoning-based tasks:
- Arithmetic:
  - "Fernando brings in three dozen bagels to a breakfast with 16 attendees. If each attendees eats two bagels, how many are left over?"
- Commonsense reasoning:
  - "The man had a fear of illness, so he never visited friends who were a what? (a) sick person (b) hospital (C) elderly person (d) graveyard."

"Chain-of-Thought Prompting Elicits Reasoning in Large Language Models." NeurIPS 2022.

TrustLLM

Funded by
the European Union

# Improving In-Context Learning

**Multi-Step Reasoning**

Main idea: each of the exemplars in your few-shot prompt contains logic showing *how* to solve the task.

"Chain-of-Thought Prompting Elicits Reasoning in Large Language Models." NeurIPS 2022.

# Improving In-Context Learning

## Multi-Step Reasoning with Chain-of-Thought Exemplars

Main idea: each of the exemplars in your few-shot prompt should explain *how* to solve the task.
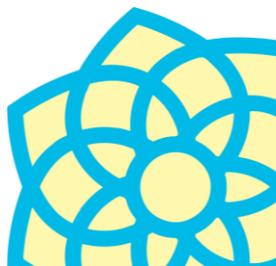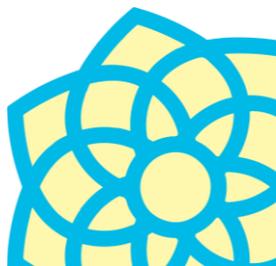


**Standard Prompting**

**Model Input**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

**Model Output**

A: The answer is 27. ✖

"Chain-of-Thought Prompting Elicits Reasoning in Large Language Models." NeurIPS 2022.

# Improving In-Context Learning

**Multi-Step Reasoning with Chain-of-Thought Exemplars**

Main idea: each of the exemplars in your few-shot prompt should explain *how* to solve the task.



"Chain-of-Thought Prompting Elicits Reasoning in Large Language Models." NeurIPS 2022.

# Improving In-Context Learning

**Multi-Step Reasoning with Chain-of-Thought Exemplars**

Main idea: each of the exemplars in your few-shot prompt should explain *how* to solve the task.
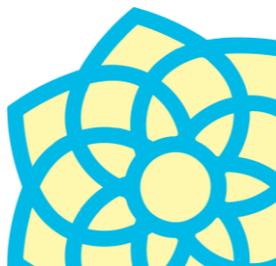
**Standard Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ✖
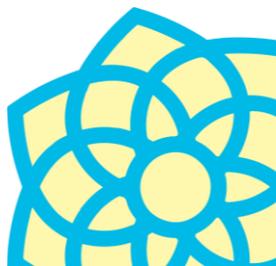
**Chain-of-Thought Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?
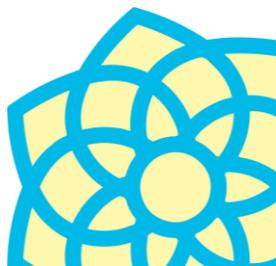
Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔

"Chain-of-Thought Prompting Elicits Reasoning in Large Language Models." NeurIPS 2022.

# Improving In-Context Learning

**Multi-Step Reasoning with Chain-of-Thought Exemplars**

Main idea: each of the exemplars in your few-shot prompt should explain *how* to solve the task.



"Chain-of-Thought Prompting Elicits Reasoning in Large Language Models." NeurIPS 2022.

# Improving In-Context Learning

**Multi-Step Reasoning with Zero-Shot Chain-of-Thought**

Main idea: We don't need any exemplars! Just append the string "Let's think step by step." to the end of the prompt.

Large Language Models are Zero-Shot Reasoners." NeurIPS 2022.

# Improving In-Context Learning

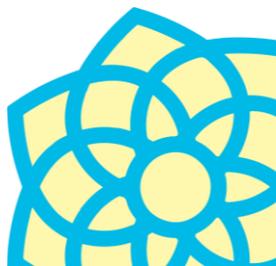## Multi-Step Reasoning with Zero-Shot Chain-of-Thought

Main idea: We don't need any exemplars! Just append the string "Let's think step by step." to the end of the prompt.



(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
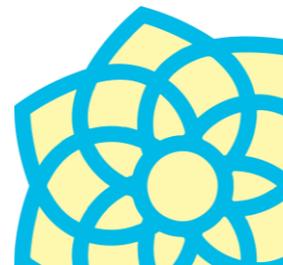A:

(Output) The answer is 8. ✗

Large Language Models are Zero-Shot Reasoners." NeurIPS 2022.

# Improving In-Context Learning

## Multi-Step Reasoning with Zero-Shot Chain-of-Thought

Main idea: We don't need any exemplars! Just append the string "Let's think step by step." to the end of the prompt.

### (a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:

(Output) The answer is 8. X

### (c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: The answer (arabic numerals) is

(Output) 8 X

Large Language Models are Zero-Shot Reasoners." NeurIPS 2022.

# Improving In-Context Learning

**Multi-Step Reasoning with Zero-Shot Chain-of-Thought**

Main idea: We don't need any exemplars! Just append the string "Let's think step by step." to the end of the prompt.

(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
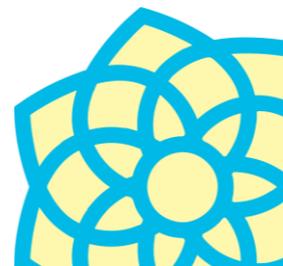A:

(Output) The answer is 8. ✗

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16 / 2 = 8 golf balls. Half of the golf balls are blue. So there are 8 / 2 = 4 blue golf balls. The answer is 4. ✓

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
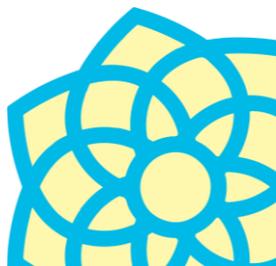A: The answer (arabic numerals) is

(Output) 8 ✗

Large Language Models are Zero-Shot Reasoners." NeurIPS 2022.

# Improving In-Context Learning

**Multi-Step Reasoning with Zero-Shot Chain-of-Thought**

Main idea: We don't need any exemplars! Just append the string "Let's think step by step." to the end of the prompt.

(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:

---

(Output) *The answer is 8.* ✗

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:

---

(Output) *The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16 / 2 = 8 golf balls. Half of the golf balls are blue. So there are 8 / 2 = 4 blue golf balls. The answer is 4.* ✓

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: The answer (arabic numerals) is

---

(Output) *8* ✗

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
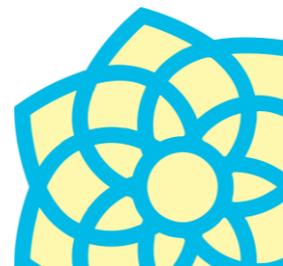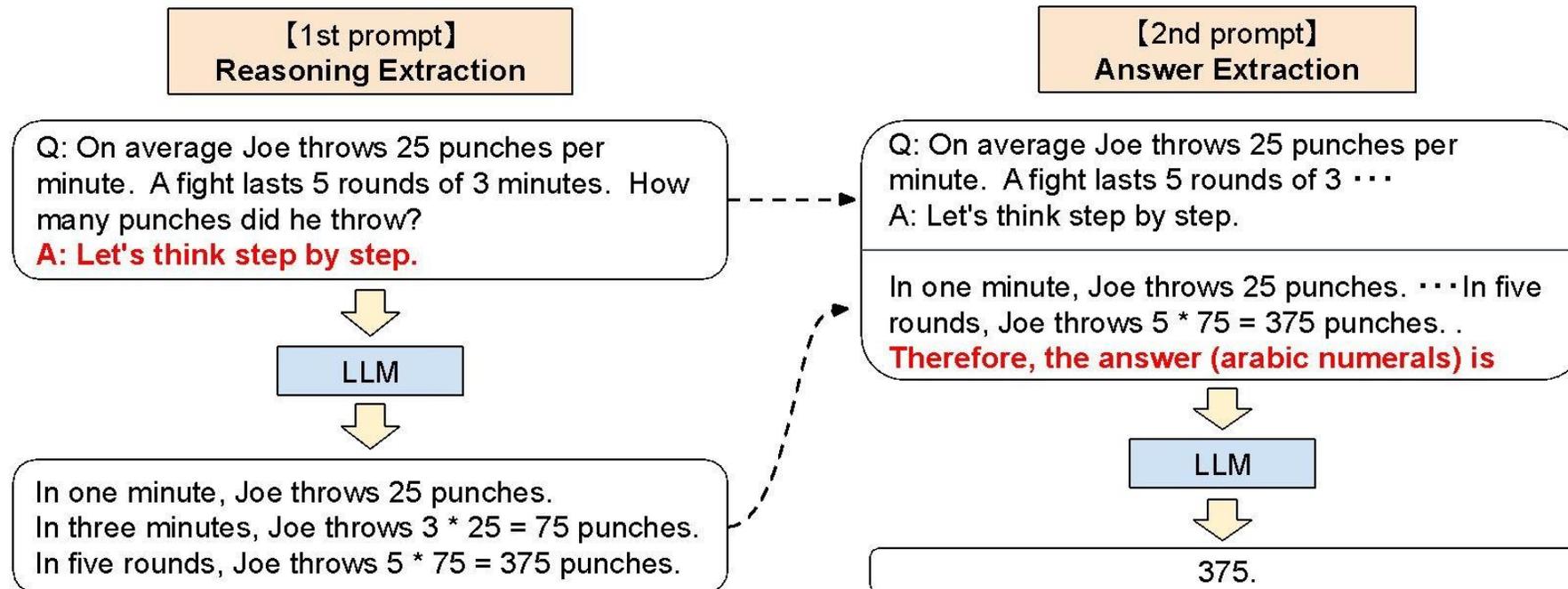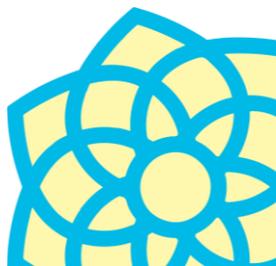A: **Let's think step by step.**

---

(Output) *There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls.* ✓

# Improving In-Context Learning

## Multi-Step Reasoning with Zero-Shot Chain-of-Thought

Main idea: We don't need any exemplars! Just append the string "Let's think step by step." to the end of the prompt.



【1st prompt】
Reasoning Extraction

Q: On average Joe throws 25 punches per minute. A fight lasts 5 rounds of 3 minutes. How many punches did he throw?
**A: Let's think step by step.**

LLM

In one minute, Joe throws 25 punches.
In three minutes, Joe throws 3 * 25 = 75 punches.
In five rounds, Joe throws 5 * 75 = 375 punches.

【2nd prompt】
Answer Extraction

Q: On average Joe throws 25 punches per minute. A fight lasts 5 rounds of 3 ···
A: Let's think step by step.

In one minute, Joe throws 25 punches. ···In five rounds, Joe throws 5 * 75 = 375 punches. .
**Therefore, the answer (arabic numerals) is**

LLM

375.

Large Language Models are Zero-Shot Reasoners." NeurIPS 2022.

# Improving In-Context Learning

**Multi-Step Reasoning with Zero-Shot Chain-of-Thought**

Main idea: We don't need any exemplars! Just append the string "Let's think step by step." to the end of the prompt.

Advantages over chain-of-thought (CoT) method:

- The single fixed instruction "Let's think step by step" works over a large variety of different tasks.
- Few-shot CoT performance degrades when there is misalignment between the example question types in the prompt and the actual task question.

Large Language Models are Zero-Shot Reasoners." NeurIPS 2022.

# Improving In-Context Learning

**Multi-Step Reasoning with Zero-Shot Chain-of-Thought**

| | MultiArith | GSM8K |
|---|---|---|
| **Zero-Shot** | **17.7** | **10.4** |
| Few-Shot (2 samples) | 33.7 | 15.6 |
| Few-Shot (8 samples) | 33.8 | 15.6 |
| **Zero-Shot-CoT** | **78.7** | **40.7** |
| Few-Shot-CoT (2 samples) | 84.8 | 41.3 |
| Few-Shot-CoT (4 samples : First) (*1) | 89.2 | - |
| Few-Shot-CoT (4 samples : Second) (*1) | 90.5 | - |
| Few-Shot-CoT (8 samples) | 93.0 | 48.7 |
| **Zero-Plus-Few-Shot-CoT (8 samples) (*2)** | **92.8** | **51.5** |
| Finetuned GPT-3 175B [Wei et al., 2022] | - | 33 |
| Finetuned GPT-3 175B + verifier [Wei et al., 2022] | - | 55 |

Large Language Models are Zero-Shot Reasoners." NeurIPS 2022.

# Improving In-Context Learning

**Better Trained Models**

As new generations of LLMs become increasingly instruction-tuned, the need for painstaking prompt engineering

has decreased but not gone away entirely.

# Improving In-Context Learning

## Better Trained Models

As new generations of LLMs become increasingly instruction-tuned, the need for painstaking prompt engineering has decreased but not gone away entirely.

Even today's "pre-trained" models have often been exposed to non-negligible amounts of instruction-following data.

# Improving In-Context Learning

**Dividing Tasks into Minimal United**

For complex generation tasks, many iterative calls to an LLM will generally work better (and be easier to evaluate) than one single prompt asking the LLM to do all parts of the task at once.

# Improving In-Context Learning

## Dividing Tasks into Minimal United

For complex generation tasks, many iterative calls to an LLM will generally work better (and be easier to evaluate) than one single prompt asking the LLM to do all parts of the task at once.

Example: Generating short stories

You could ask an LLM to generate an entire story at once.

Or you could ask it to:

1.  generate a synopsis

2.  given the synopsis, generate a character list and a sequence of events

3.  given all of the above, generate the actual story text.

Breaking the task into parts reduces the complexity of each individual call to the model and also allows more human intervention.

TrustLLM

Funded by
the European Union
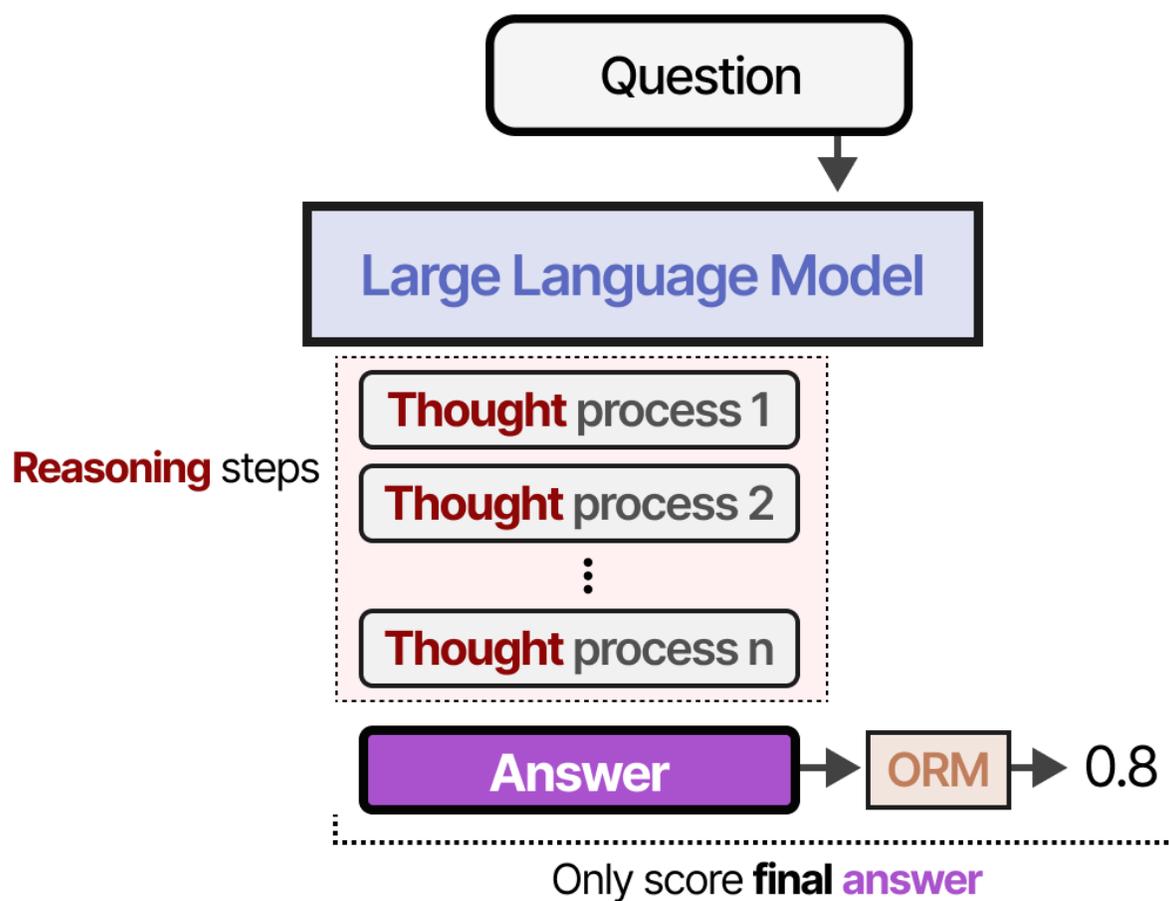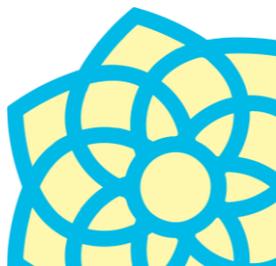
# Reasoning

# Towards Reasoning
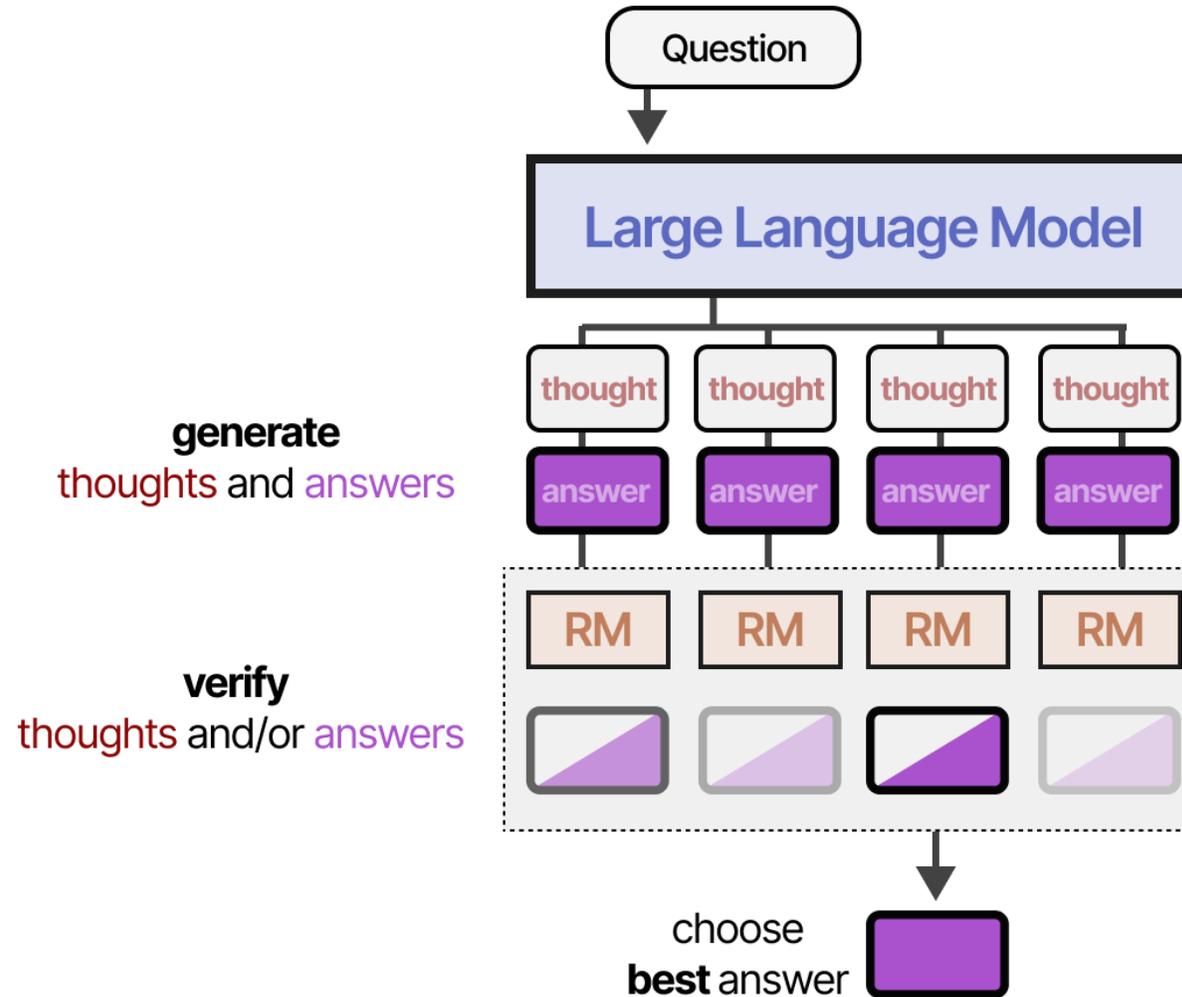
# Test-Time Compute a.k.a. Reasoning

# Test-Time Compute a.k.a. Reasoning

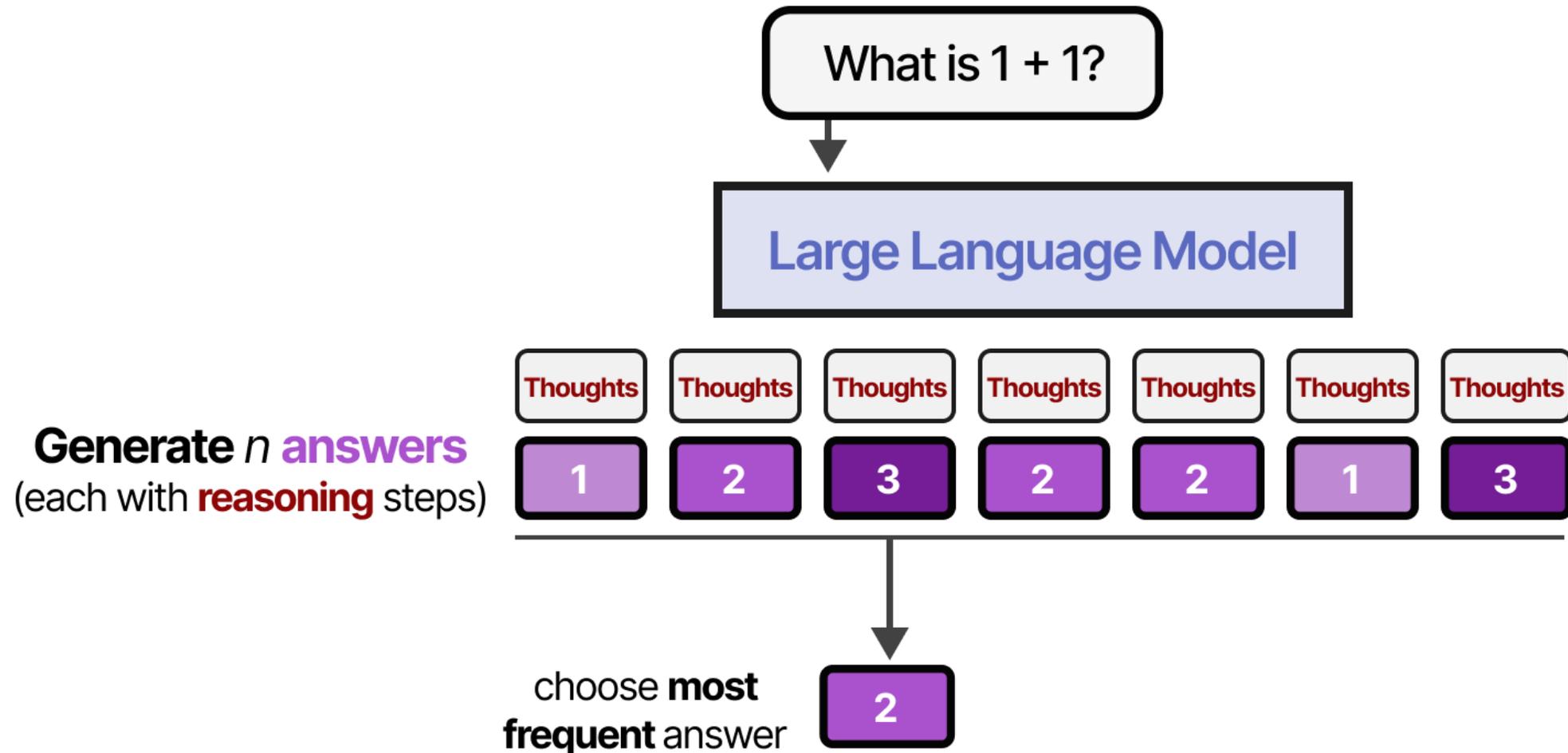# Outcome vs Process Reward Model

https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-reasoning-llms

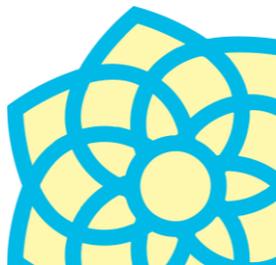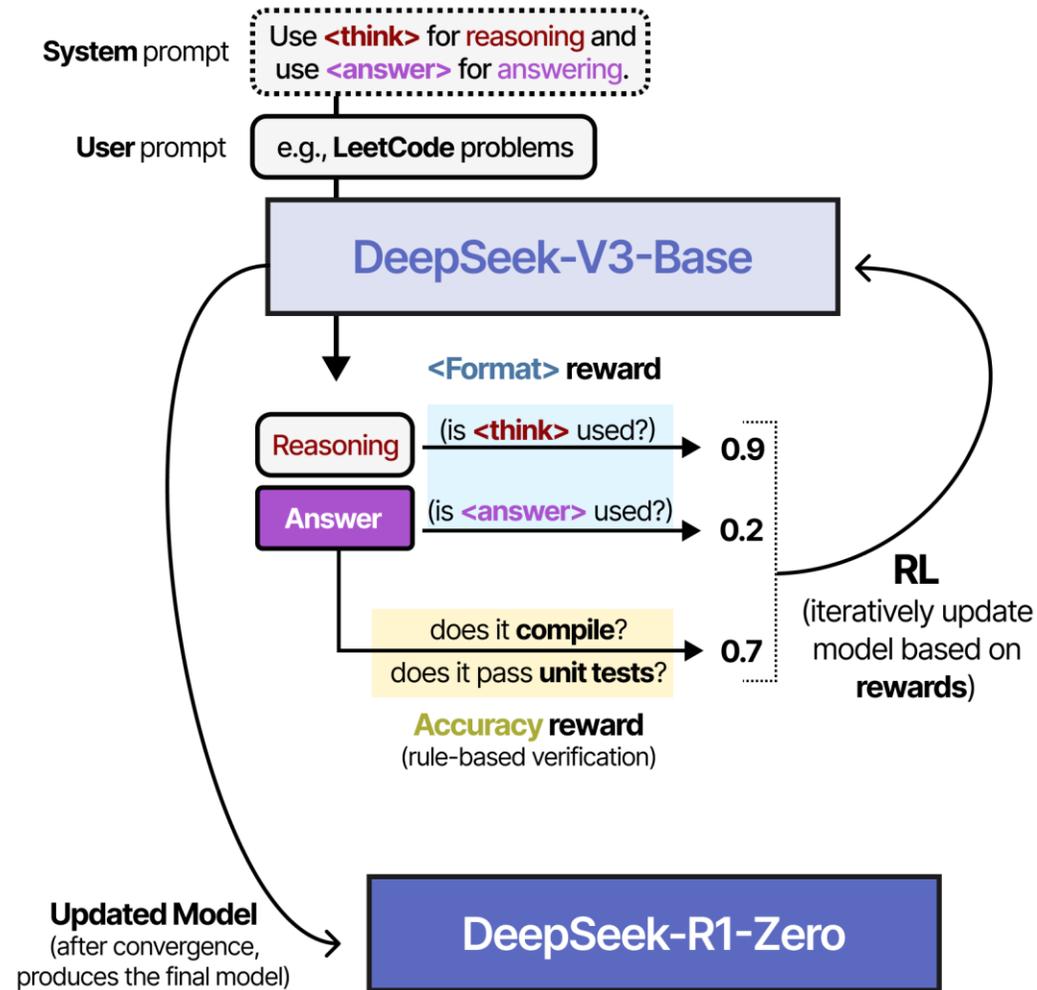# Search Against Verifiers

# Majority Voting

# Reasoning with DeepSeek-R1 Zero

- **LAIM LE6 VT2025:**
  **Inference**
  **Retrieval Augmented Generation**
  **In-Context Learning**
  **Reasoning**

# www.ida.liu.se/~freheo8/llm