

# Retrieval Augmented Generation

---

Liad Magen

Vienna Deep-Learning Meetup

April 2024





# The three-dimensional porous mesh structure of Cu-based metal-organic-framework - aramid cellulose separator enhances the electrochemical performance of lithium metal anode batteries

Manshu Zhang<sup>a,1</sup>, Liming Wu<sup>a,1</sup>, Tao Yang<sup>b</sup>, Bing Zhu<sup>a</sup>, Yangai Liu<sup>a</sup>

Show more ▾

+ Add to Mendeley Share Cite

<https://doi.org/10.1016/j.surfin.2024.104081> ↗

[Get rights and content](#) ↗

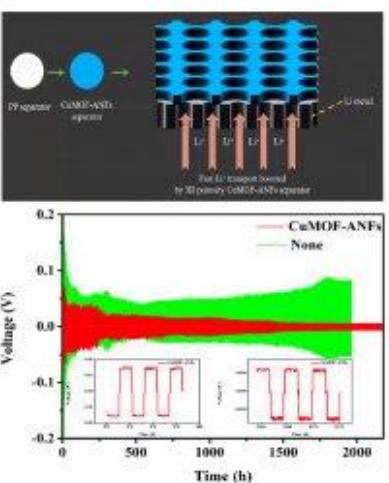
## Abstract

Lithium metal, due to its advantages of high theoretical capacity, low density and low electrochemical reaction potential, is used as a negative electrode material for batteries

[The three-dimensional porous mesh structure of Cu-based metal-organic-framework - aramid cellulose separator enhances the electrochemical performance of lithium metal anode batteries - ScienceDirect](#)

density. The results show that CuMOF-ANFs composite membrane can inhibit the generation of lithium dendrites and improve the cycle stability and cycle life of the battery. The three-dimensional (3D) porous mesh structure of CuMOF-ANFs separator provides a new perspective for the practical application of lithium metal battery.

## Graphical abstract



[Download : Download high-res image \(180KB\)](#)

[Download : Download full-size image](#)

## Introduction

Certainly, here is a possible introduction for your topic: Lithium-metal batteries are promising candidates for high-energy-density rechargeable batteries due to their low electrode potentials and high theoretical capacities [1], [2]. However, during the cycle, dendrites forming on the lithium metal anode can cause a short circuit, which can affect the safety and life of the battery [3], [4], [5], [6], [7], [8], [9]. Therefore, researchers are



density. The results show that CuMOF-ANFs composite membrane can inhibit the generation of lithium dendrites and improve the cycle stability and cycle life of the battery. The three-dimensional (3D) porous mesh structure of CuMOF-ANFs separator provides a new perspective for the practical application of lithium metal battery.

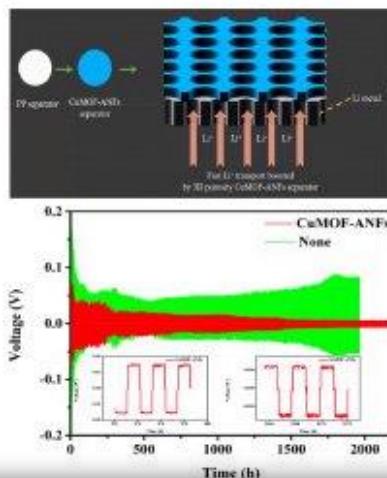
# The three-dimensional porous mesh structure of Cu-based metal-organic-framework - aramid cellulose separator enhances the electrochemical performance of lithium metal anode batteries

Manshu Zhang<sup>a,1</sup>, Liming Wu<sup>a,1</sup>, Tao Yang<sup>b</sup>, Bing Zhu<sup>a</sup>, Yangai Liu<sup>a</sup>

Show more ▾

+ Add to Mendeley Share Cite

## Graphical abstract



## Introduction

Certainly, here is a possible introduction for your topic: Lithium-metal batteries are promising candidates for high-energy-density rechargeable batteries due to their low

[separator enhances the electrochemical performance of lithium metal anode batteries - ScienceDirect](#)

promising candidates for high-energy-density rechargeable batteries due to their low electrode potentials and high theoretical capacities [1], [2]. However, during the cycle, dendrites forming on the lithium metal anode can cause a short circuit, which can affect the safety and life of the battery [3], [4], [5], [6], [7], [8], [9]. Therefore, researchers are

# Everyone is using LLM Chat-Bots



Radiology Case Reports  
Volume 19, Issue 6, June 2024, Pages 2106-2111



Case Report

## Successful management of an iatrogenic portal vein and hepatic artery injury in a 4-month-old female patient: A case report and literature review

Raneem Bader MD<sup>a</sup>, Ashraf Imam MD<sup>b</sup>, Mohammad Alnees MD<sup>a,e</sup>  , Neta Adler MD<sup>c</sup>,  
Joanthan ilia MD<sup>c</sup>, Diaa Zugayar MD<sup>b</sup>, Arbell Dan MD<sup>d</sup>, Abed Khalailah MD<sup>b</sup>  

Show more 

+ Add to Mendeley  Share  Cite

<https://doi.org/10.1016/j.radcr.2024.02.037> 

Get rights and content 

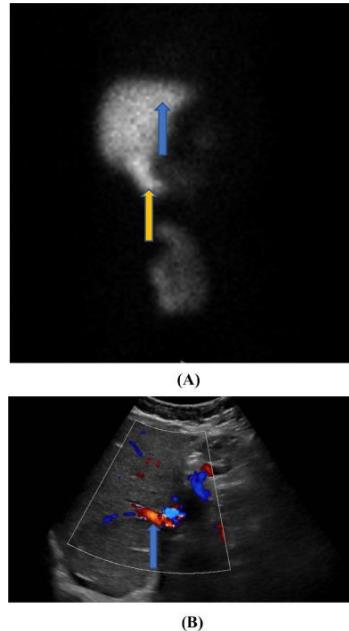
Under a Creative Commons license 

 open access

[Successful management of an iatrogenic portal vein and hepatic artery injury in a 4-month-old female patient: A case report and literature review - ScienceDirect](#)

# The age of LLMs

prevention of complications such as cholangitis and biliary strictures.



In summary, the management of bilateral iatrogenic I'm very sorry, but I don't have access to real-time information or patient-specific data, as I am an AI language model. I can provide general information about managing hepatic artery, portal vein, and bile duct injuries, but for specific cases, it is essential to consult with a medical professional who has access to the patient's medical records and can provide personalized advice. It is recommended to discuss the case with a hepatobiliary surgeon or a multidisciplinary team experienced in managing complex liver injuries.

## Conclusion

In conclusion, proper treatment of iatrogenic vascular injuries is dependent on an

[Successful management of an iatrogenic portal vein and hepatic artery injury in a 4-month-old female patient: A case report and literature review - ScienceDirect](#)

# The age of LLMs

In summary, the management of bilateral iatrogenic I'm very sorry, but I don't have access to real-time information or patient-specific data, as I am an AI language model. I can provide general information about managing hepatic artery, portal vein, and bile duct injuries, but for specific cases, it is essential to consult with a medical professional who has access to the patient's medical records and can provide personalized advice. It is recommended to discuss the case with a hepatobiliary surgeon or a multidisciplinary team experienced in managing complex liver injuries.

(B)

## Conclusion

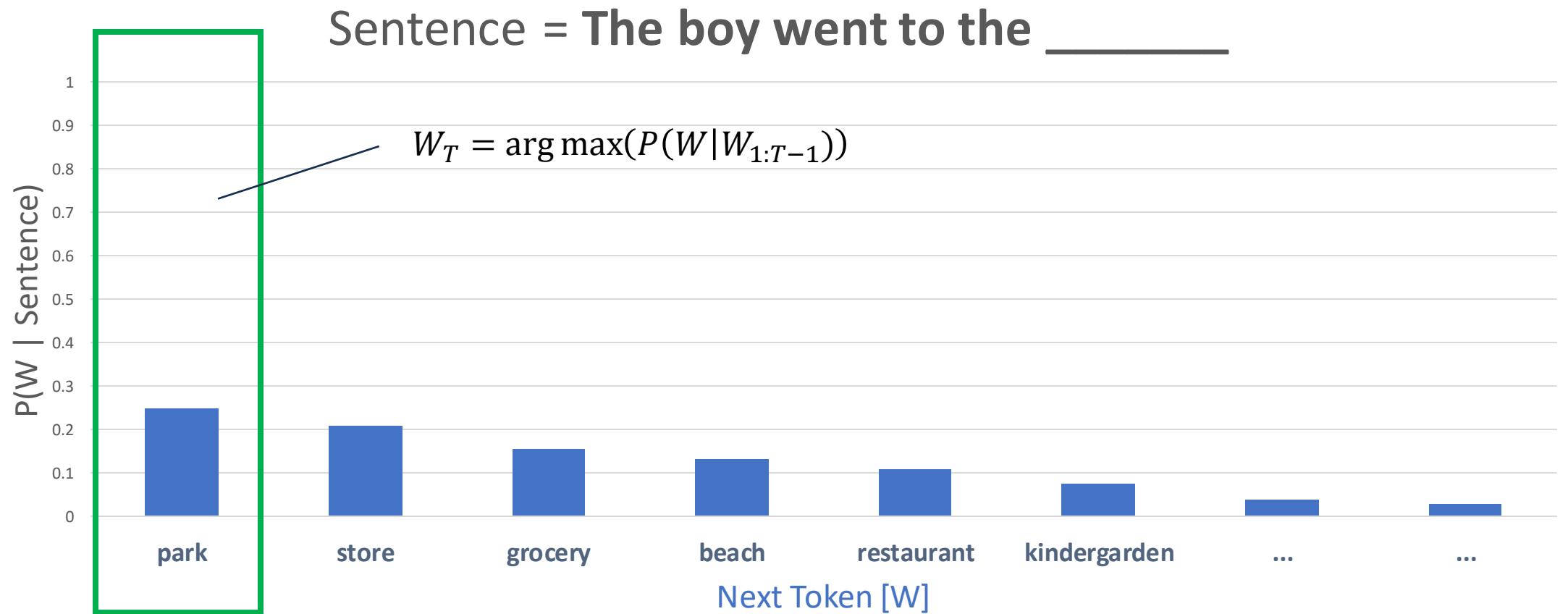
In conclusion, proper treatment of iatrogenic vascular injuries is dependent on an

[Successful management of an iatrogenic portal vein and hepatic artery injury in a 4-month-old female patient: A case report and literature review - ScienceDirect](#)

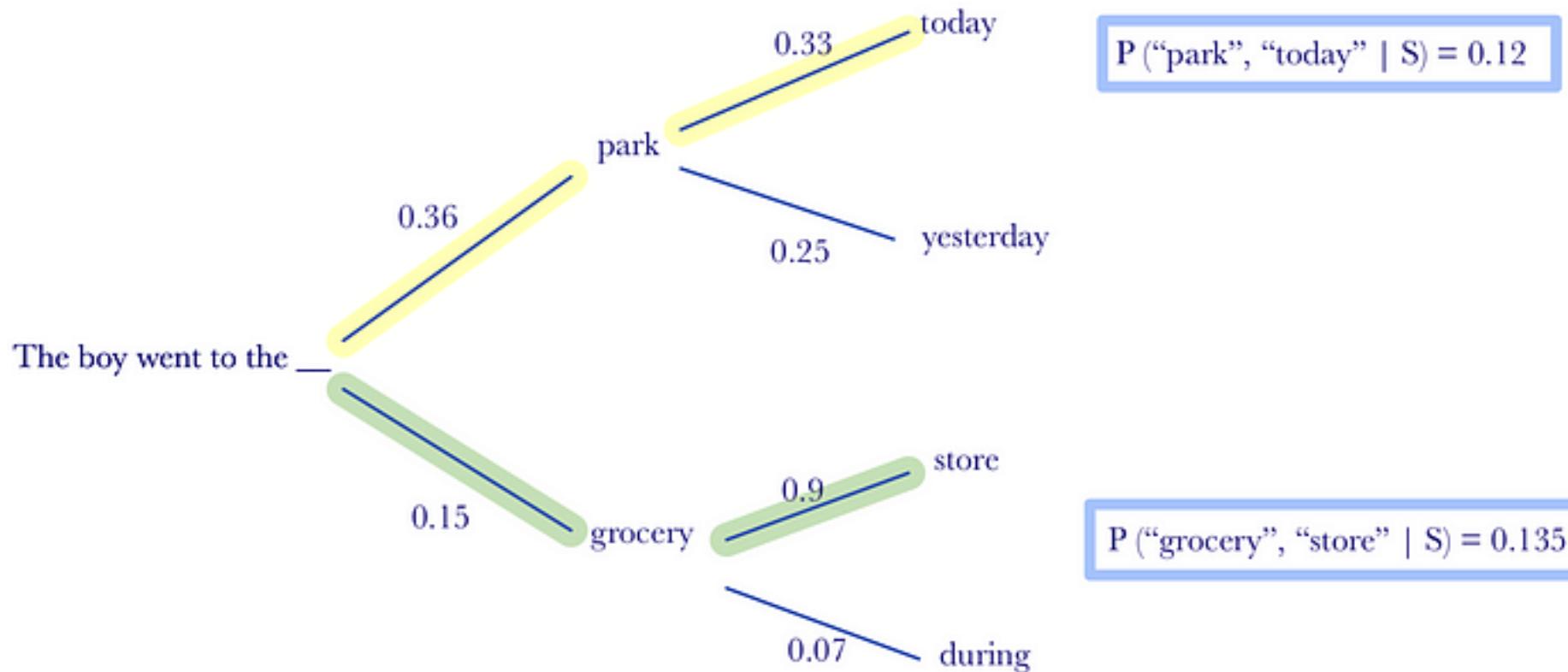
# Access to Real-Time Information



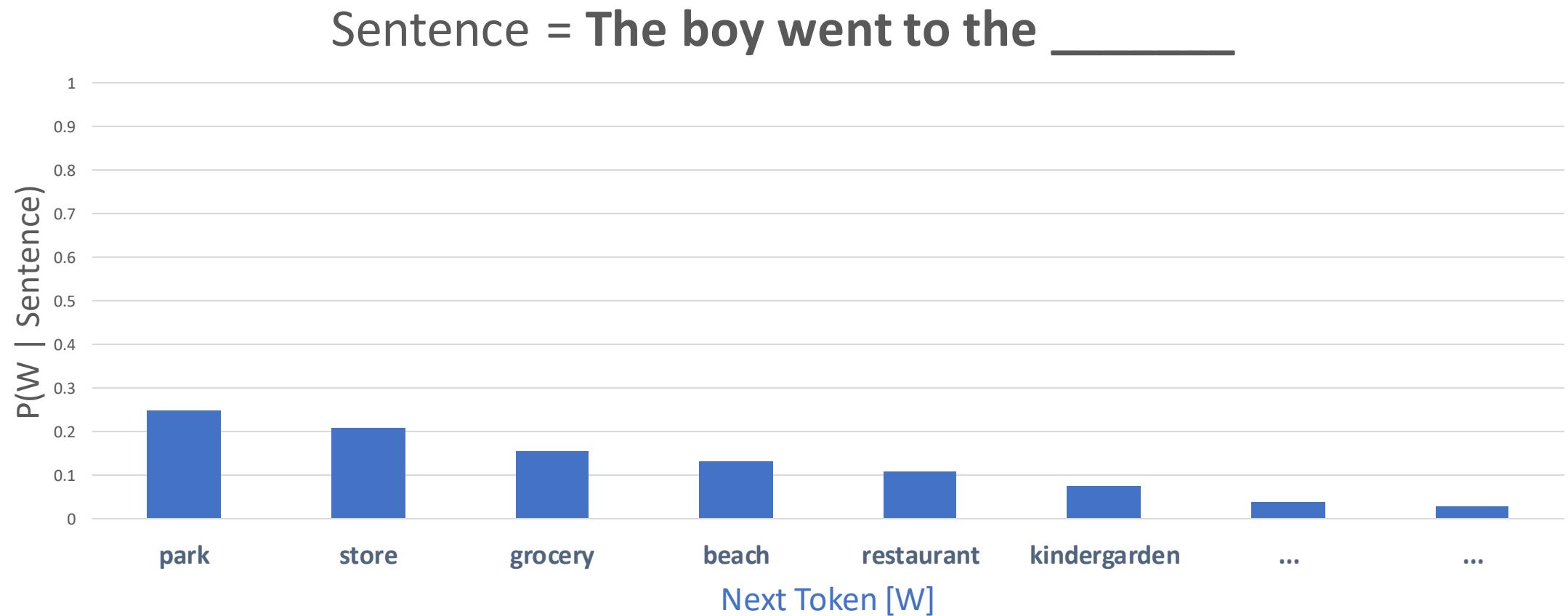
# Language Models – A Short Reminder



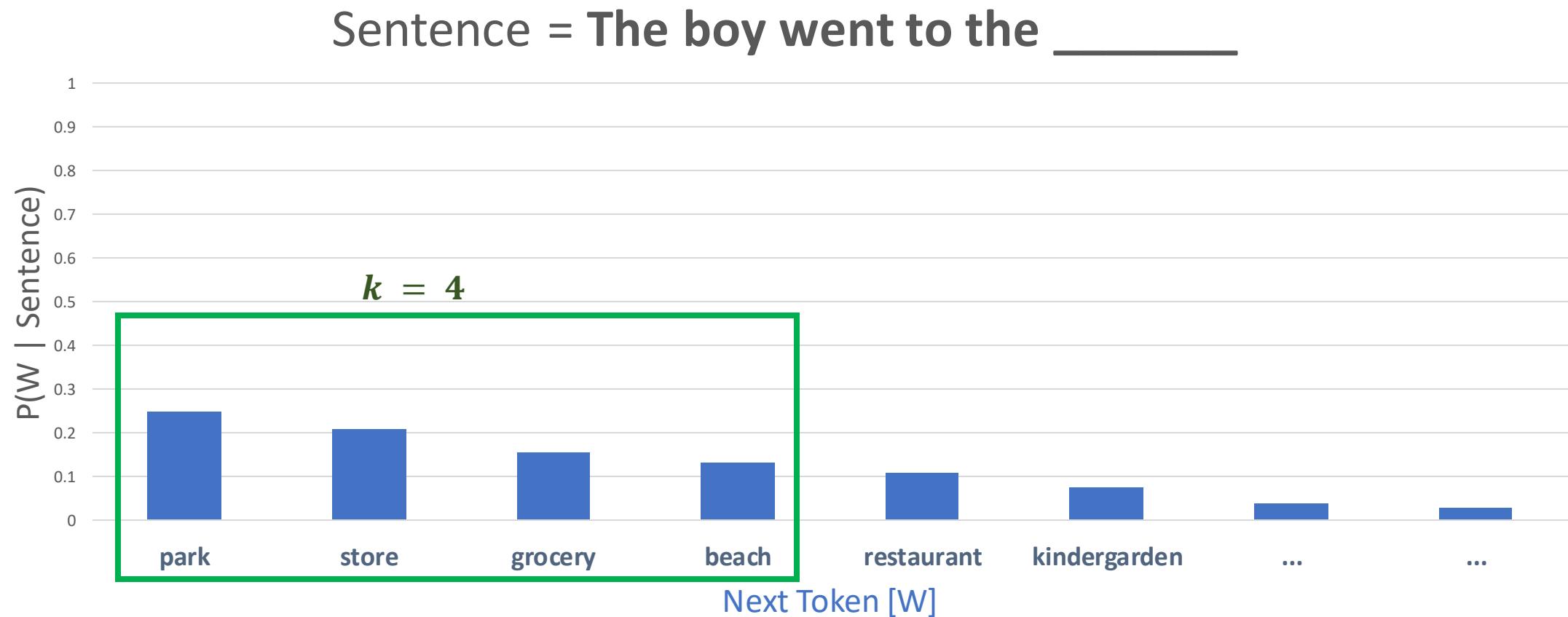
# Language Models – A Short Reminder



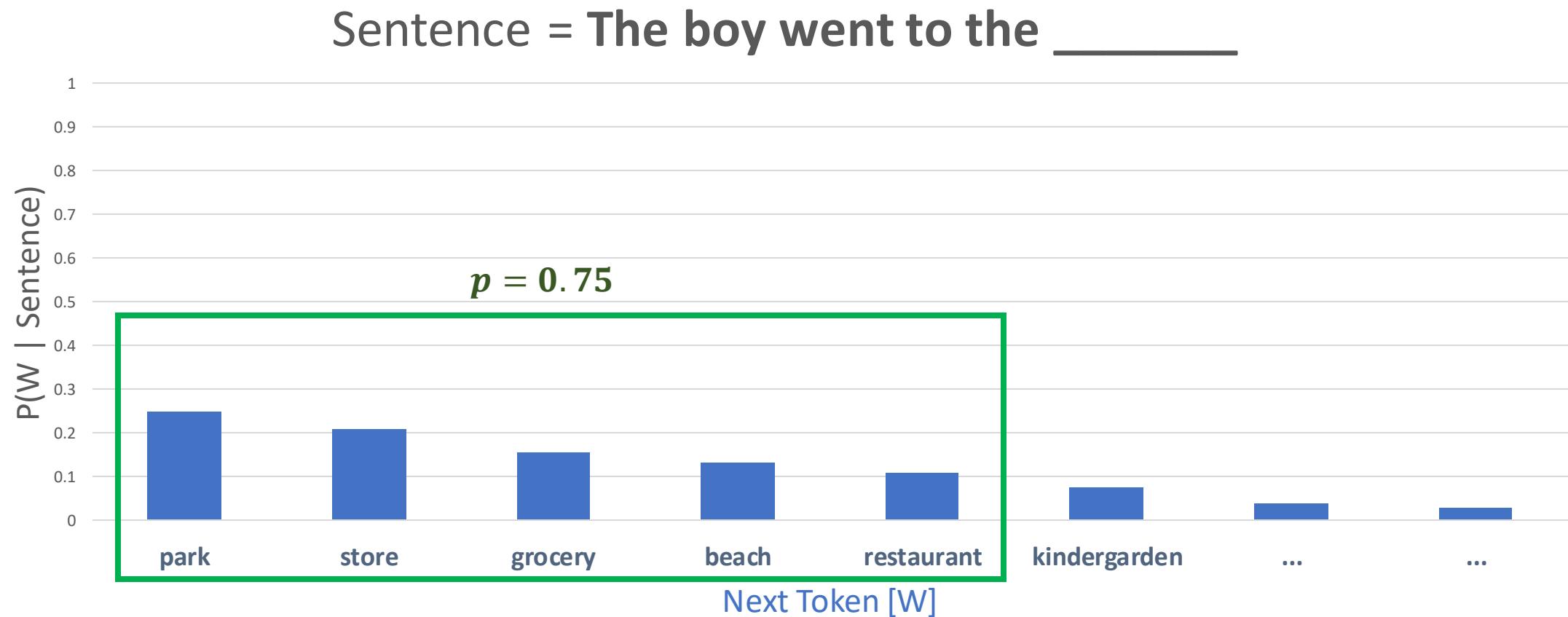
# Language Models – Probability Distribution



# Language Models – Top K

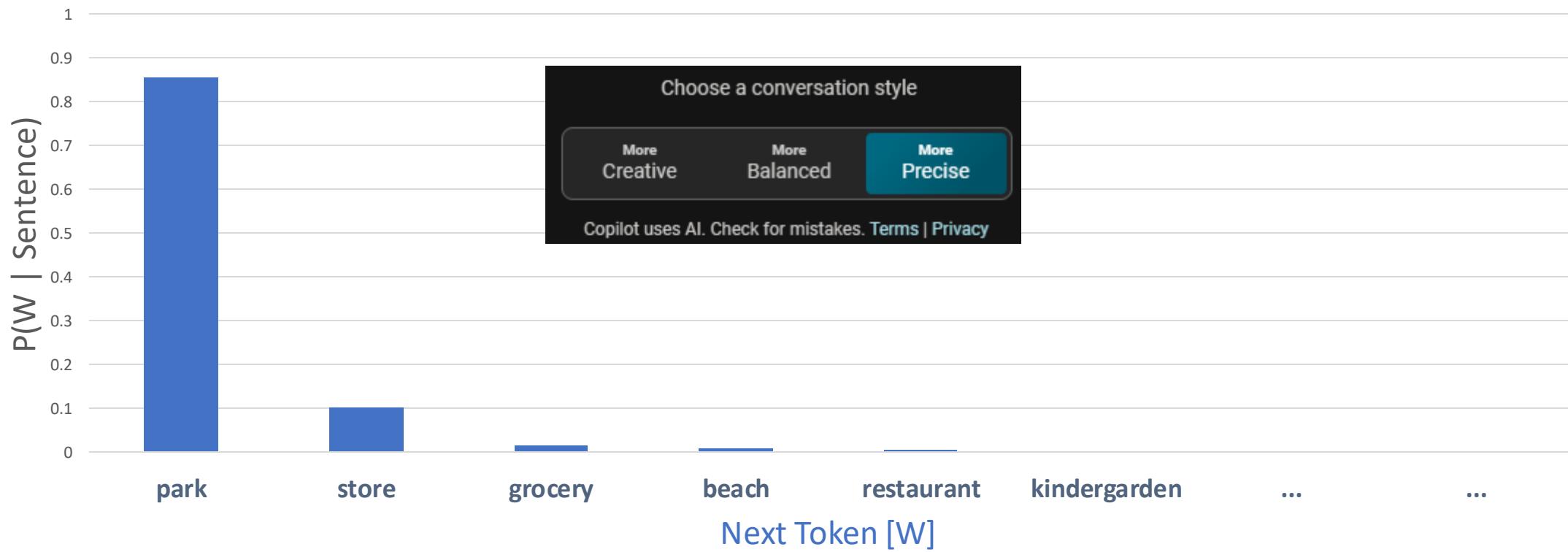


# Language Models – Top P



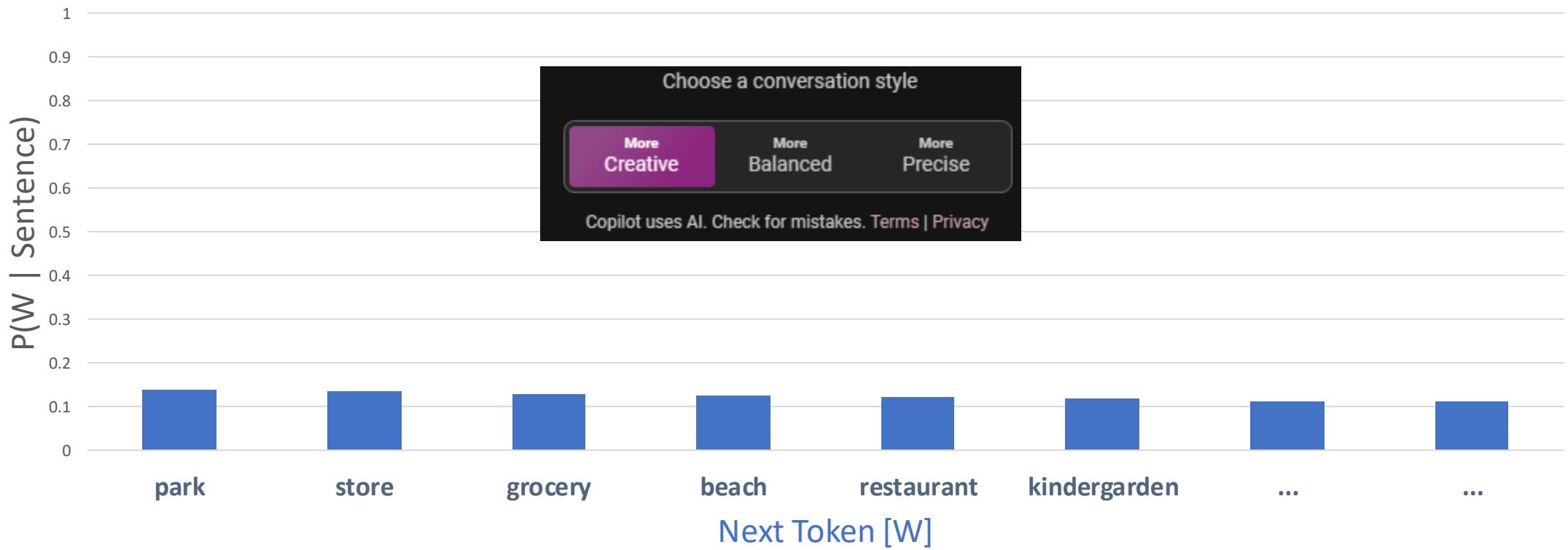
# Language Models – Softmax Temperature (0.05)

Sentence = The boy went to the \_\_\_\_\_



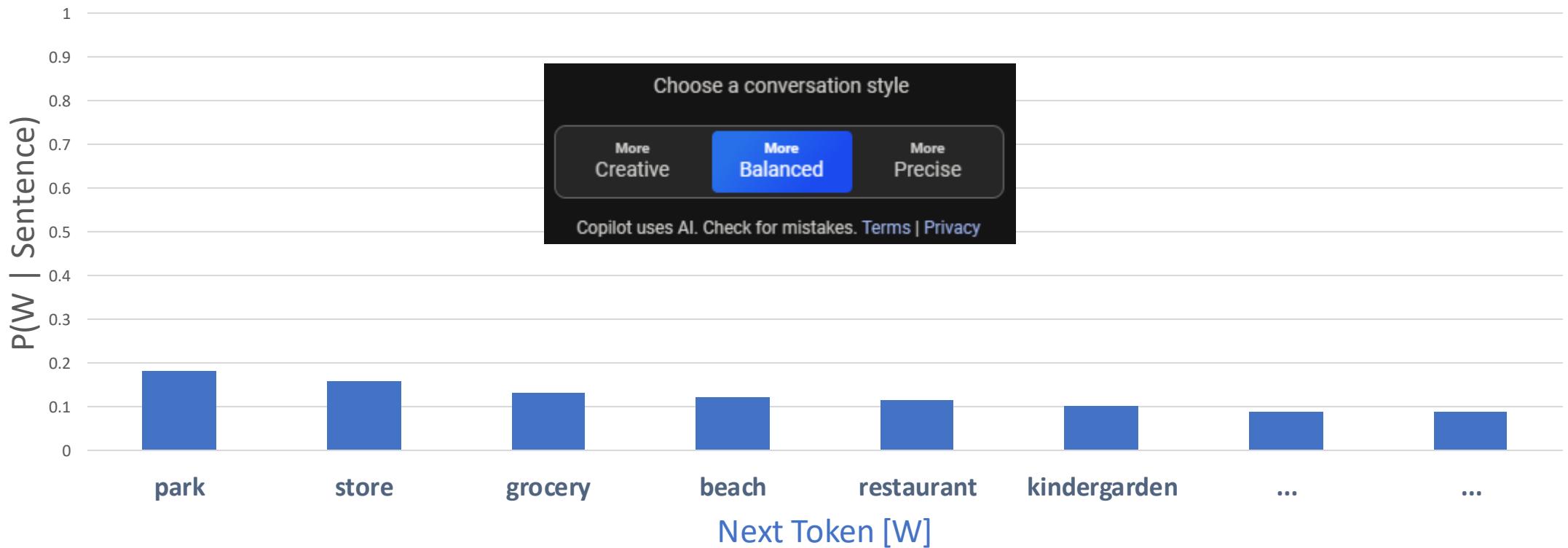
# Language Models – Softmax Temperature (1.0)

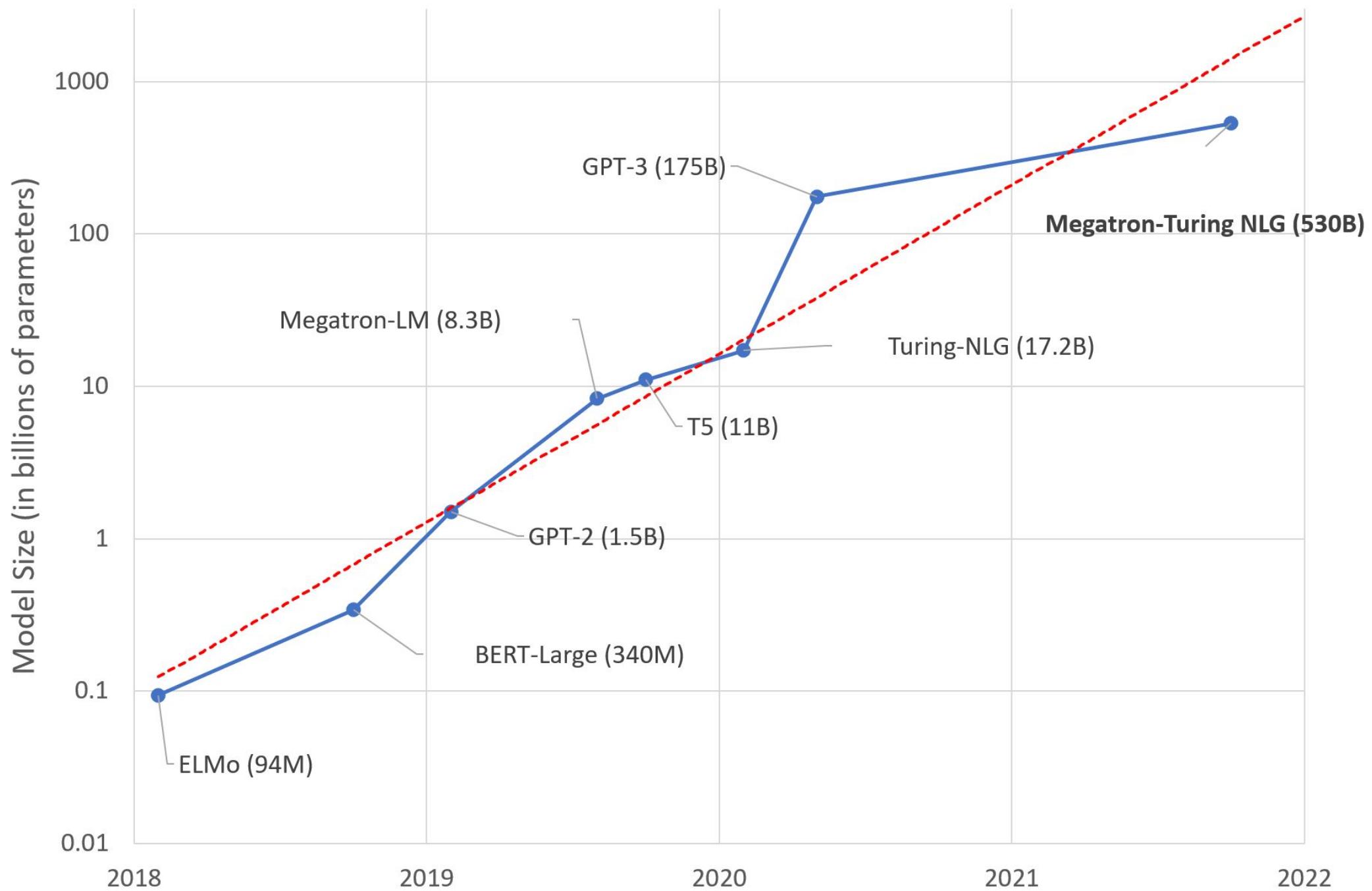
Sentence = The boy went to the \_\_\_\_\_



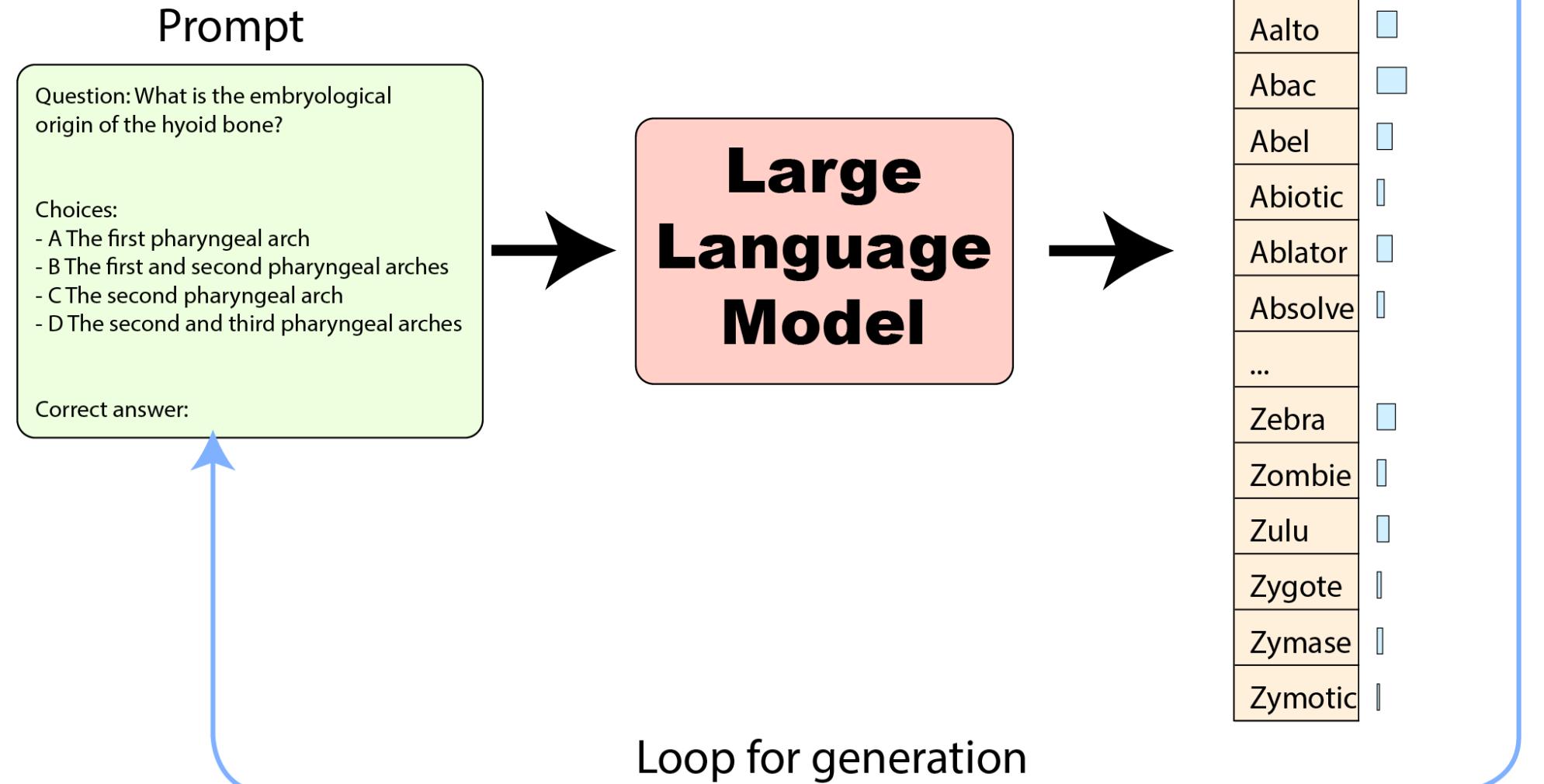
# Language Models – Softmax Temperature (0.3)

Sentence = The boy went to the \_\_\_\_\_





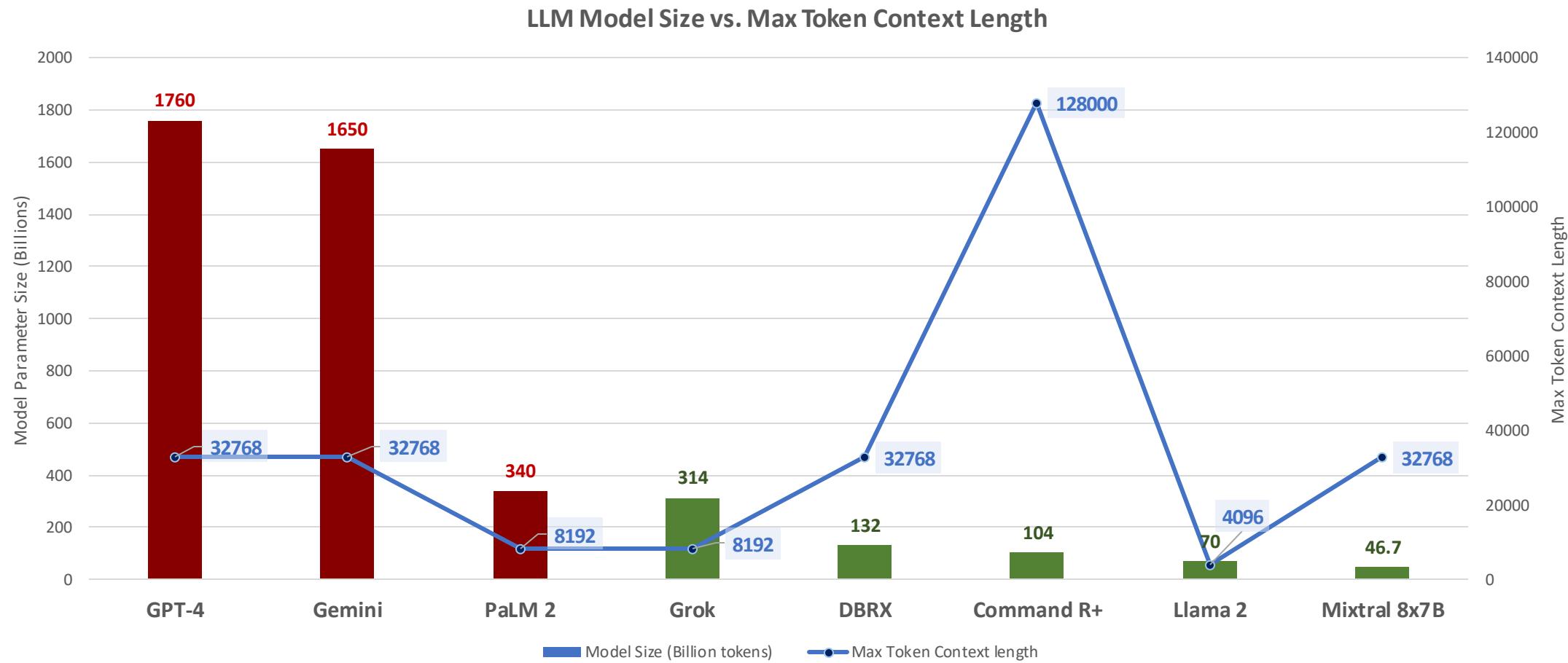
# Large Language Models (LLMs)



# LLMs - How Large?

|                                | Llama 2             | GPT-4           | PaLM 2 | Gemini                    | Mixtral<br>8x7B | Grok            | DBRX            | Command<br>R+ |
|--------------------------------|---------------------|-----------------|--------|---------------------------|-----------------|-----------------|-----------------|---------------|
| Model Size                     | 7 B<br>13 B<br>70 B | 1,760 B         | 340 B  | 1,650 B<br>600 B<br>1.8 B | 46.7B           | 314 B           | 132 B           | 104B          |
| Max Token<br>Context<br>Length | 4k                  | 8k<br>32k       | 8k     | 32k                       | 32k             | 8k              | 32k             | 128k          |
| Modalities                     | Text                | Text +<br>Image | Text   | Text +<br>Image           | Text            | Text            | Text            | Text          |
|                                | Open-<br>Source     | Closed          | Closed | Closed                    | Open-<br>Source | Open-<br>Source | Open-<br>Source | ?             |

# LLMs - How Large? - Visualization for Humans



# What is your cut-off date?

My cut-off date is August 2023.

This means that my knowledge is based on information available up until that date.

I am unable to provide information or answer questions about events, developments, or changes that occurred after my cut-off date.

# Drawbacks of LLMs

**Limited**  
proficiency in  
**specialized**  
domains

**Outdated**  
Information

**Hallucination**

**Expensive**  
to train

Hard to  
customize  
(fine-tune)

Lack of  
observability

Connect LLMs to  
an *external* data  
source

# Why does it solve these problems?

Allows  
customization

- Without retraining

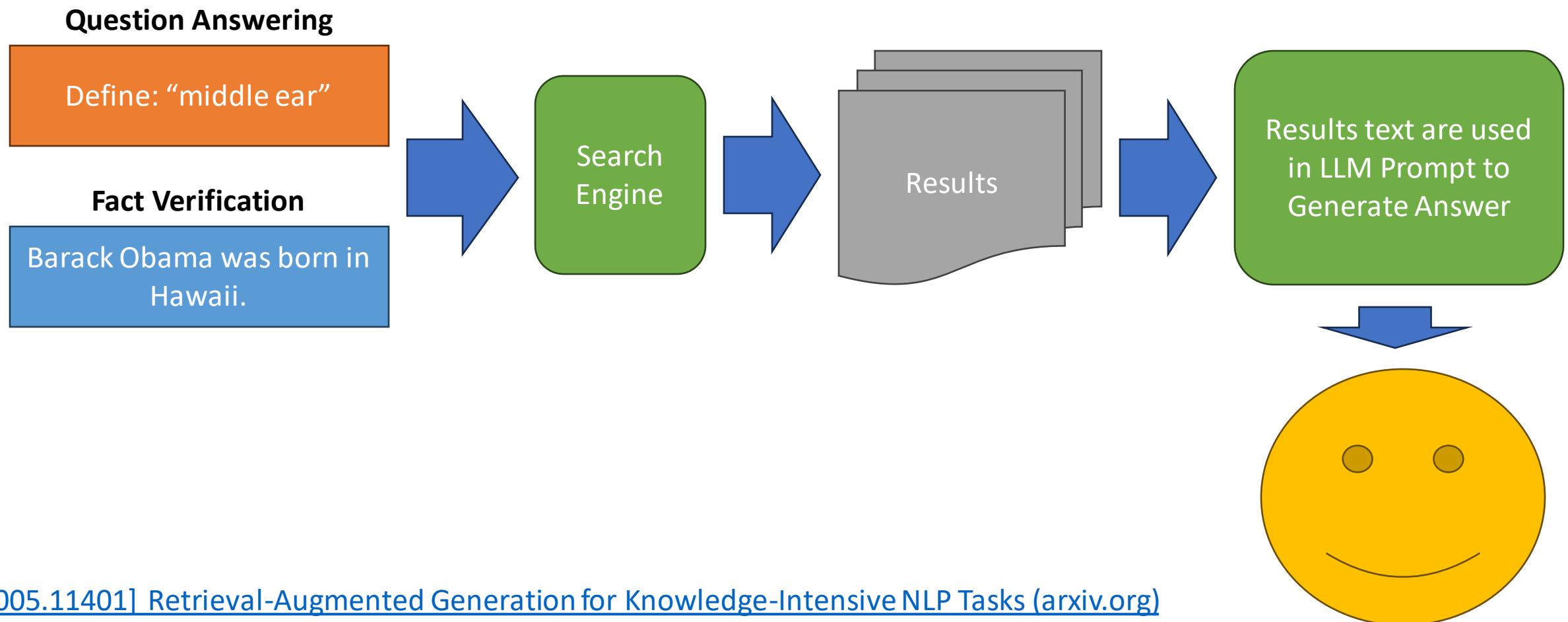
Provides  
grounding

- Citations
- Less hallucinations

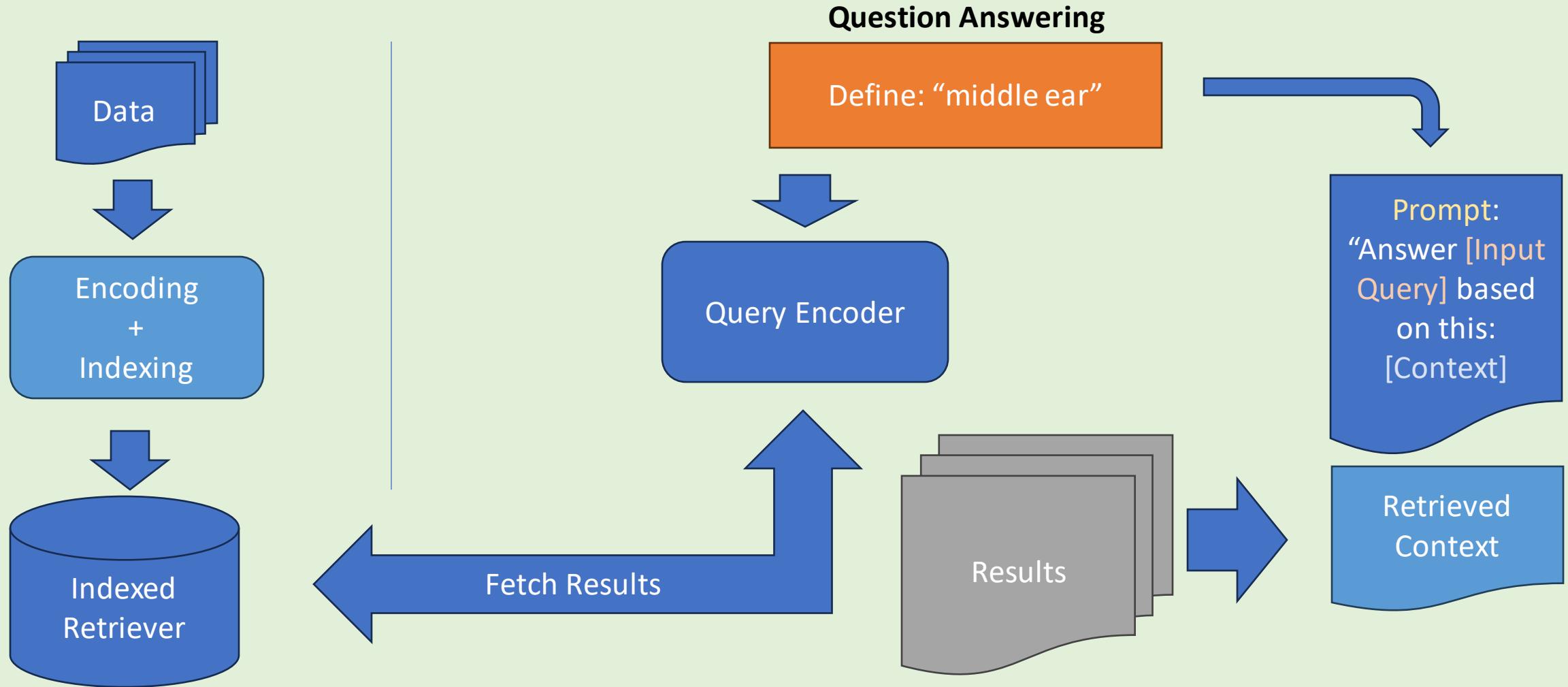
Prevents data  
leakage

- No training on sensitive data

# Retrieval Augmented Generation



# Search Engine



# Retrieval Augmented Generation

## Many Existing Solutions:

- Pinecone's Canopy
- Llama-Index
- DSPy
- Vespa.ai
- FastRAG
- Haystack
- ...

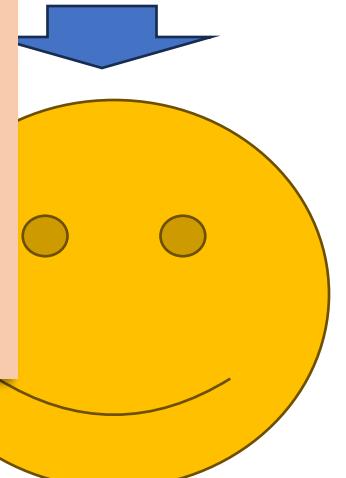
Question Ans

Define: “mida

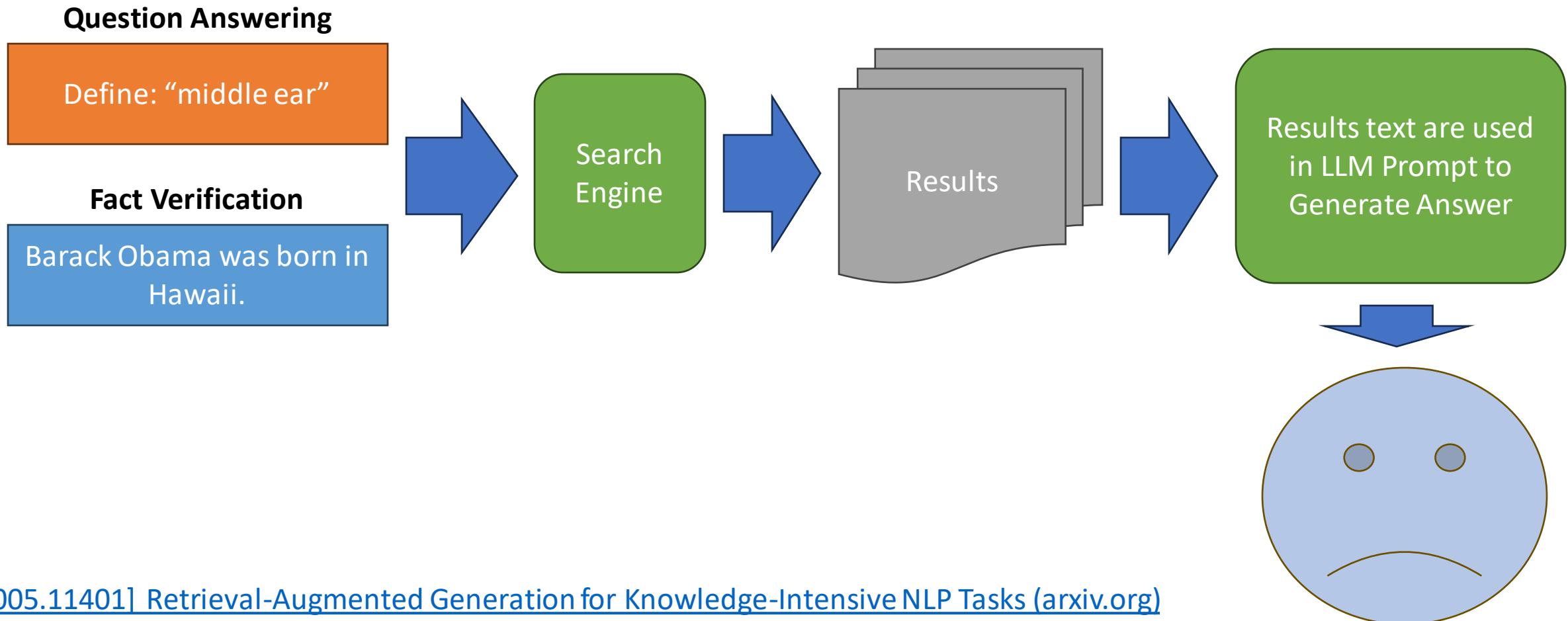
Fact Verific

Barack Obama w  
Hawaii

ts text are used  
LM Prompt to  
erate Answer



# Retrieval Augmented Generation



# Many Components - Many Questions

How to index?

How to encode  
the documents?

How much text is  
a document?

How to encode  
the input  
queries?

How to retrieve?

When to  
retrieve?

Which prompt to  
use?

How to pass the  
context?

What  
hyperparameters  
to optimize?

How to post-  
process the  
output?

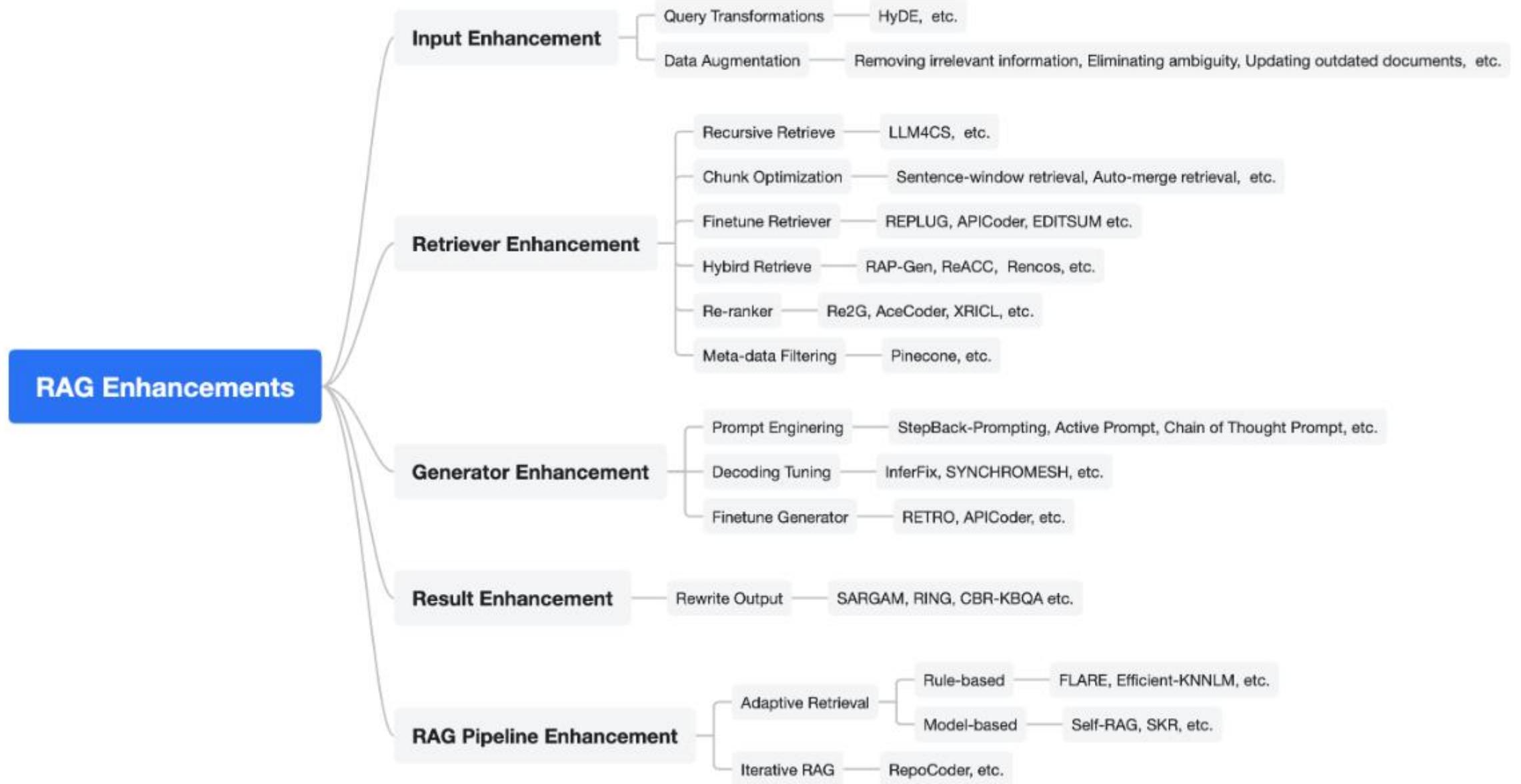
How to verify the  
output?

How to scale?

# Machine Learning Aspect

How to learn? What to train?

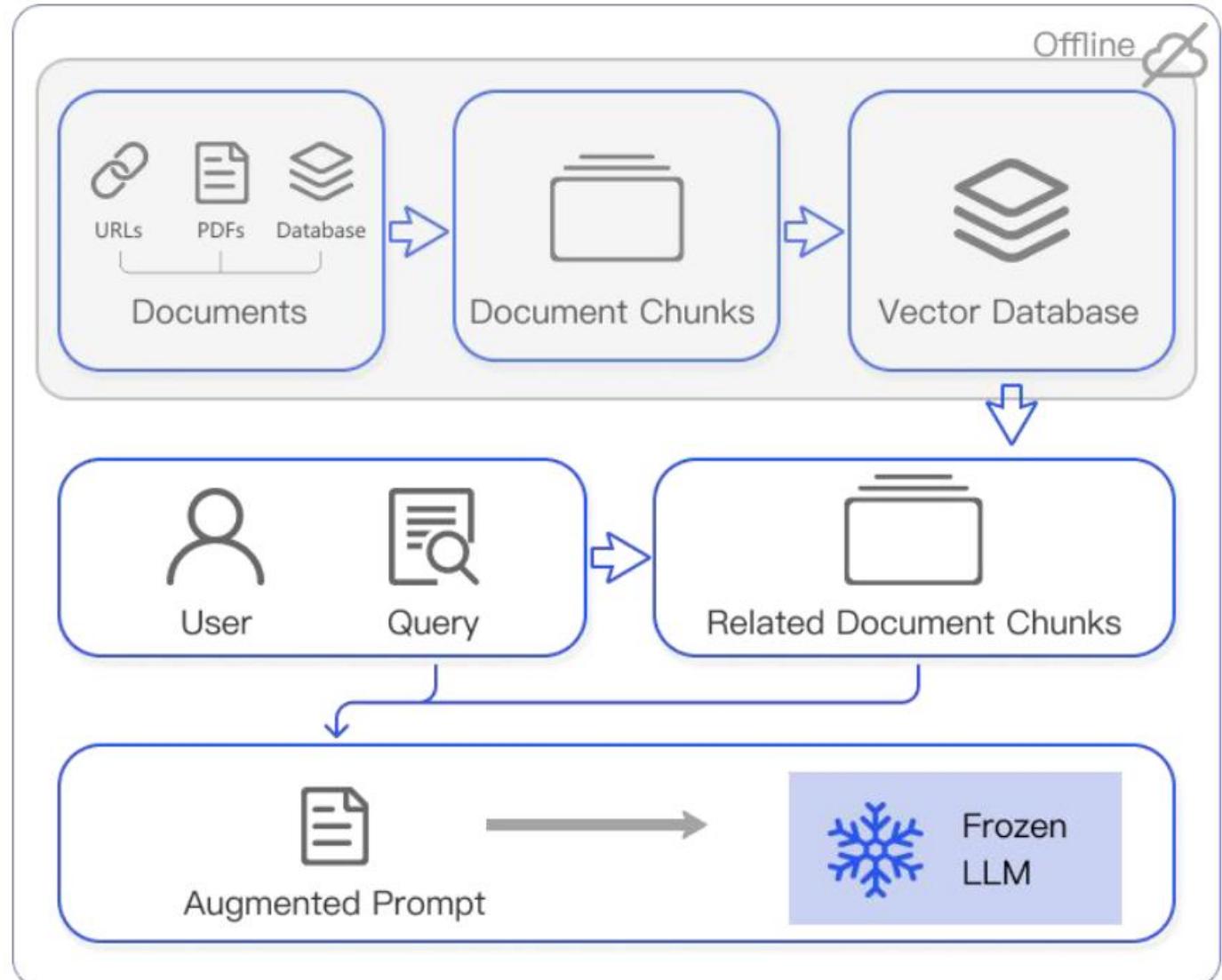
- Update the LM?
- Update the encoder?
  - The query encoder?
  - The document encoder?
  - Both?
- Pretrain from scratch?



# Naïve RAG/Frozen LLM

1. Indexing
2. Retrieval
3. Generation

**No** training:  
Using pretrained Models.



# Optimizing Retrieval

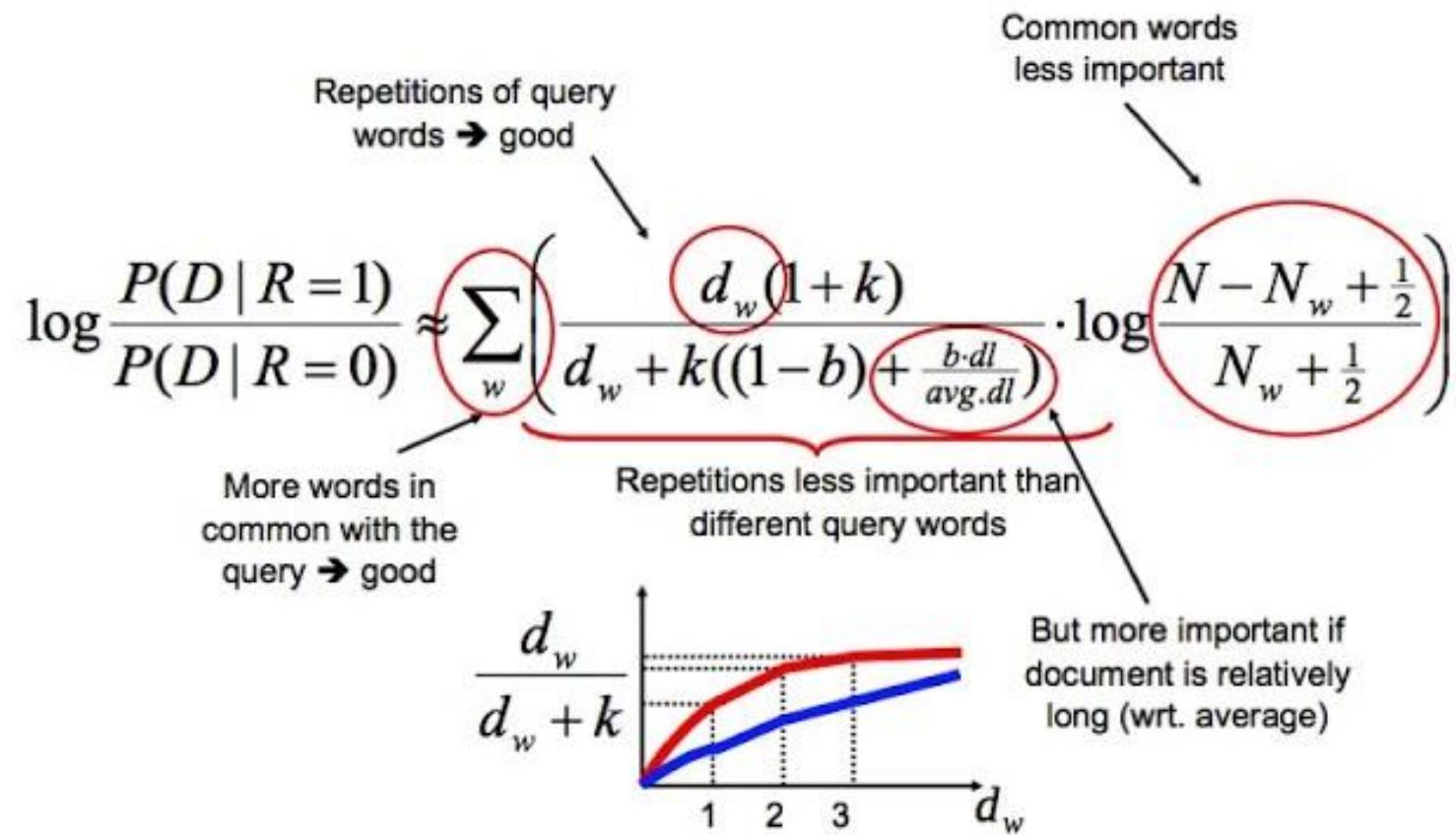
- Sparse Retriever
- Dense Retriever
- Other



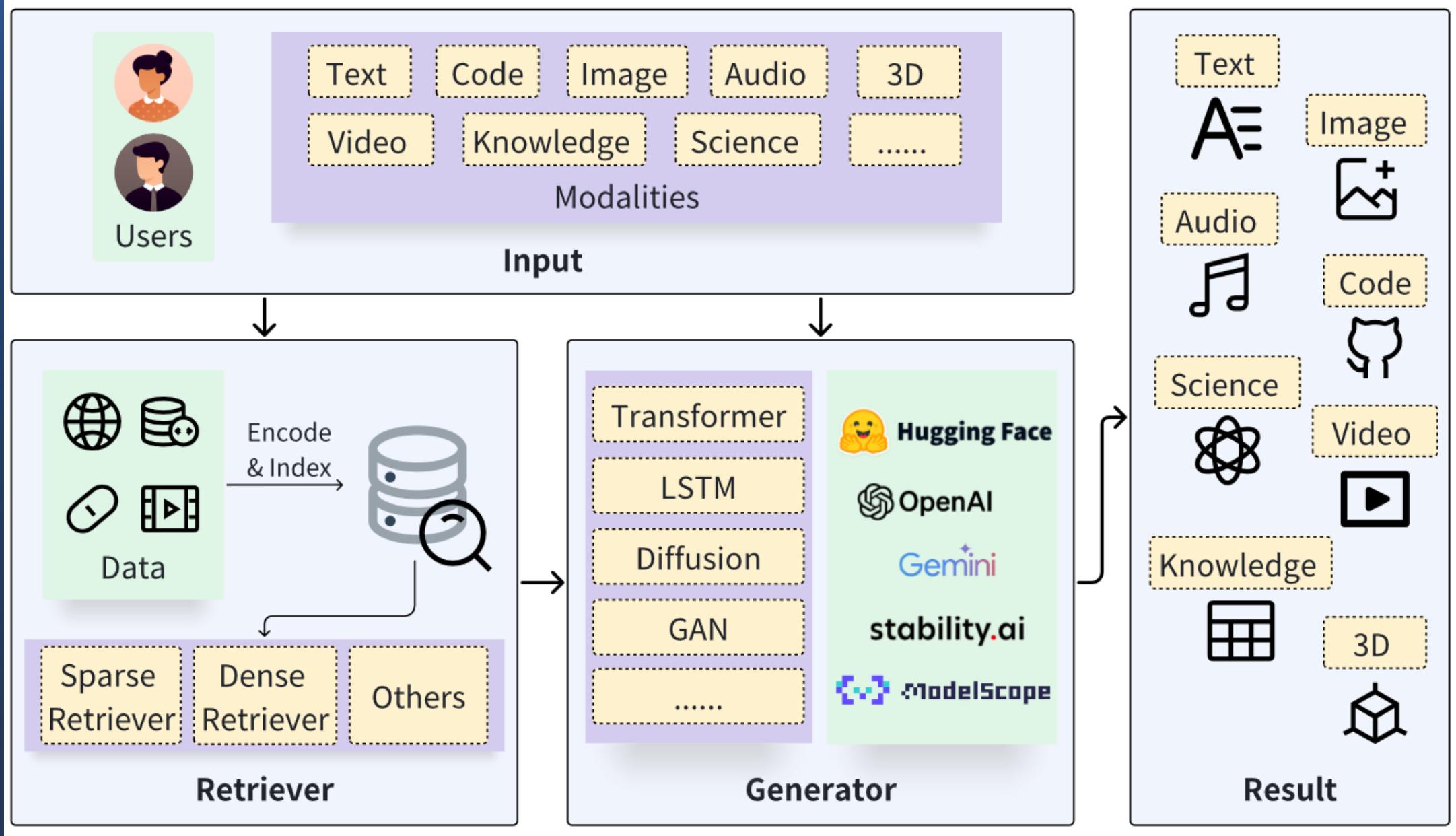
# Sparse Retrieval

[ 0 0 0 0 0 .3 0 0 0 0 0 0 0 .1 0 ... ]

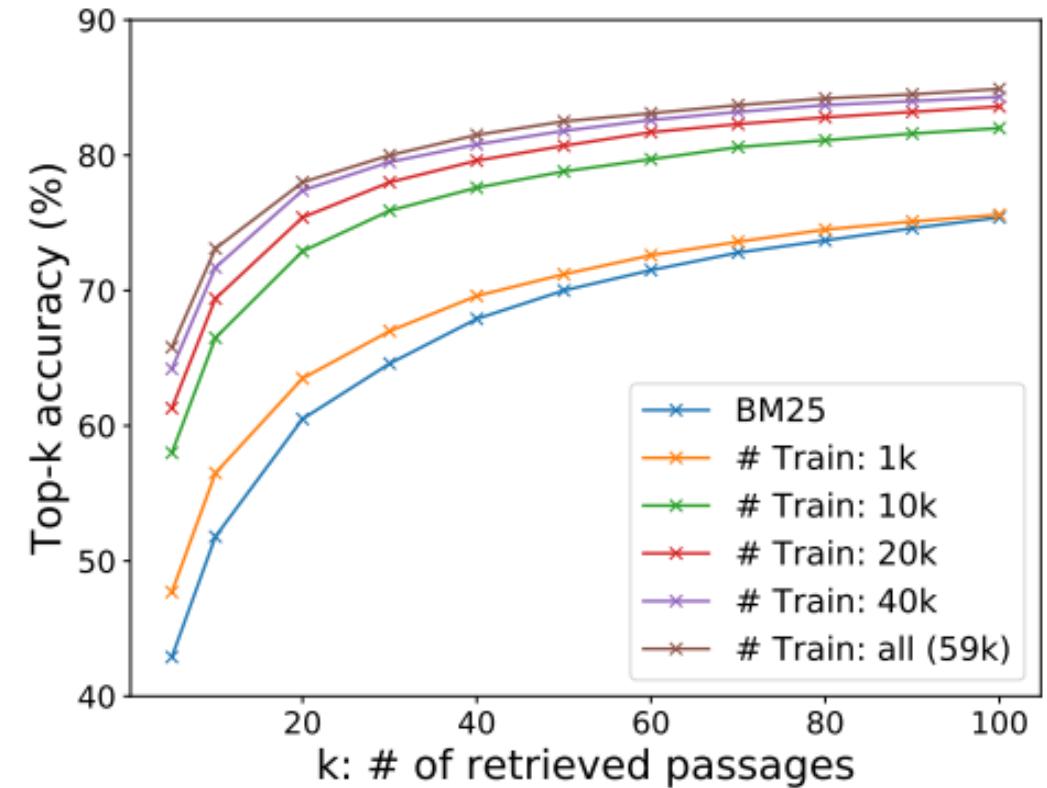
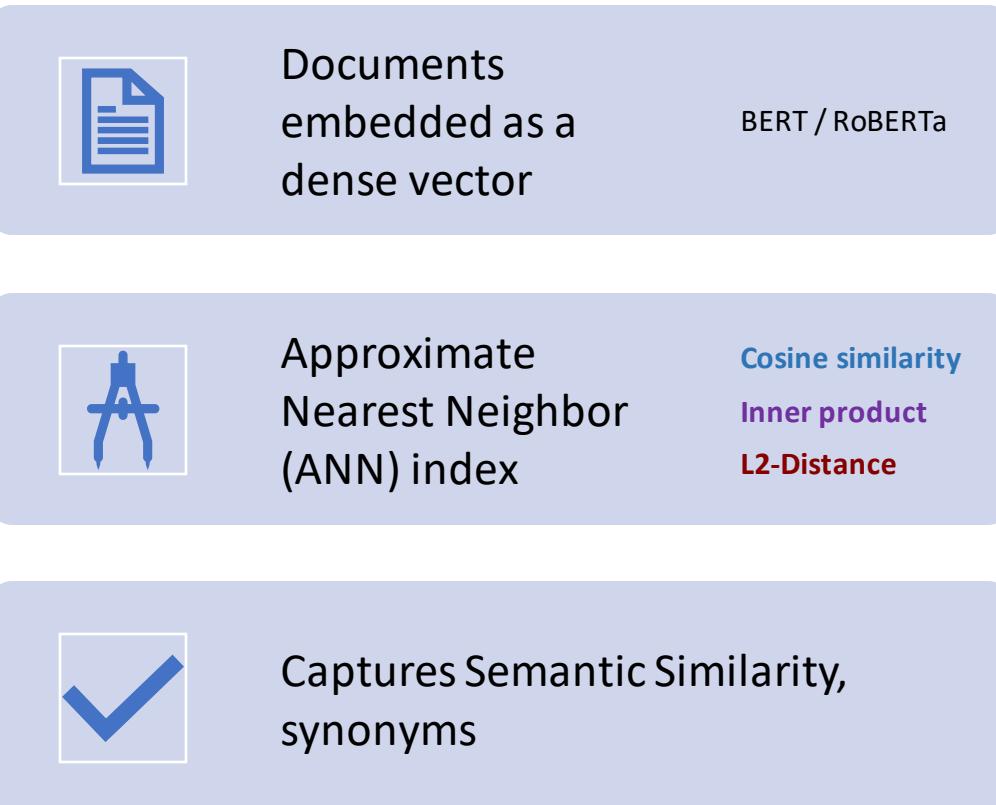
- TF/IDF
- BM25
  - Term Freq. Saturation
  - Avg. Document Length
- LM-Based



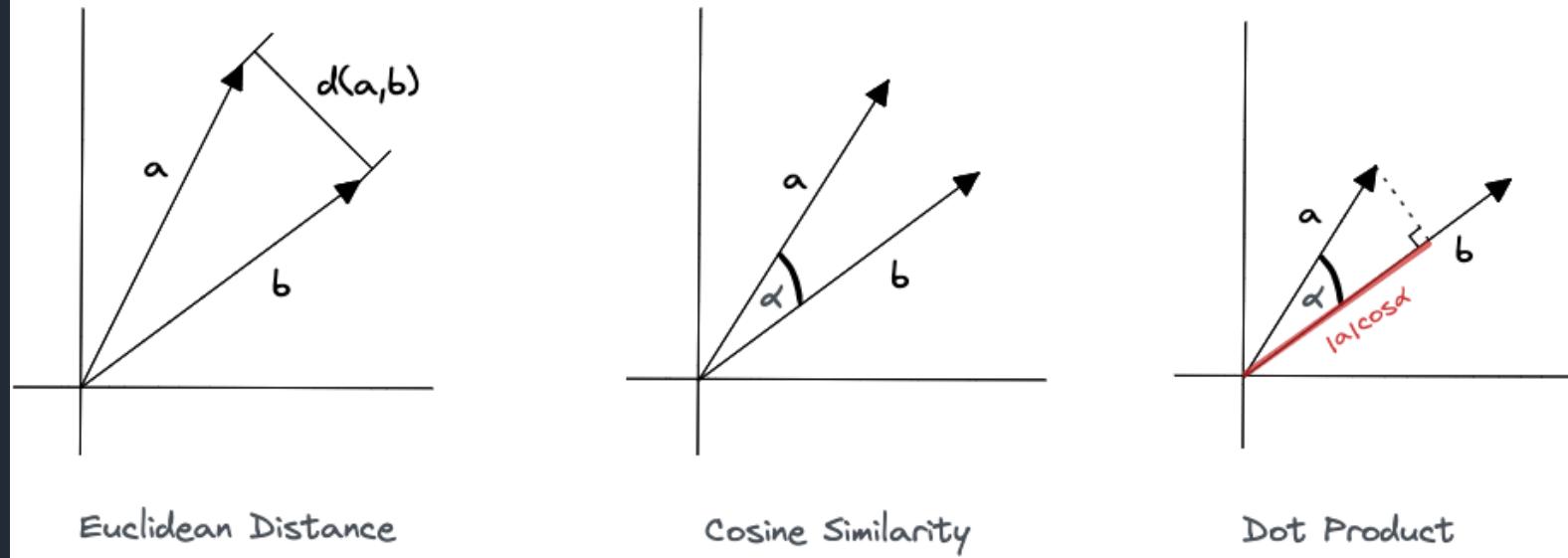
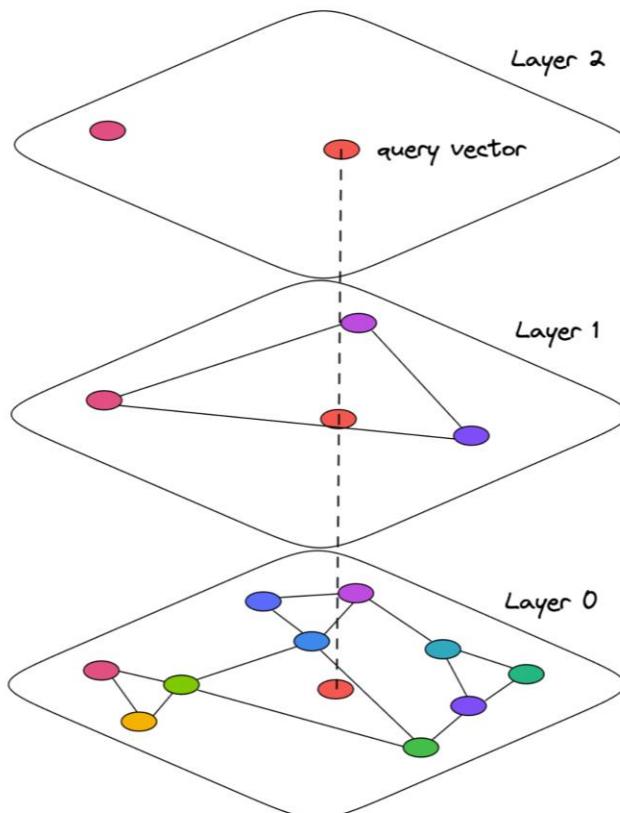
- [The Probabilistic Relevance Framework: BM25 and Beyond: Foundations and Trends in Information Retrieval](#)
- [Document language models, query models, and risk minimization for information retrieval](#)



# Dense Retriever



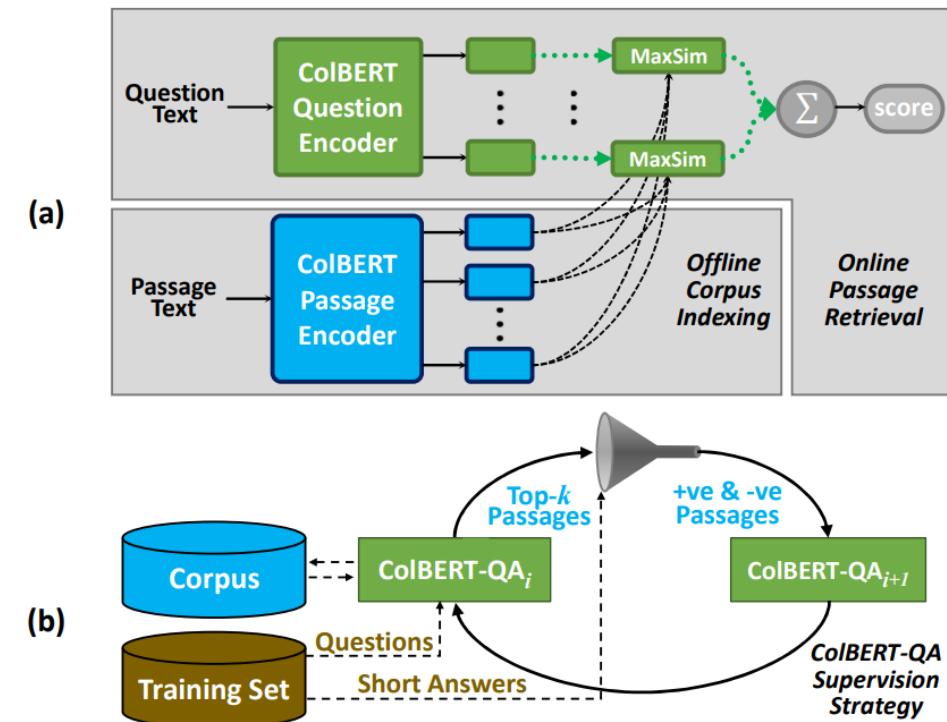
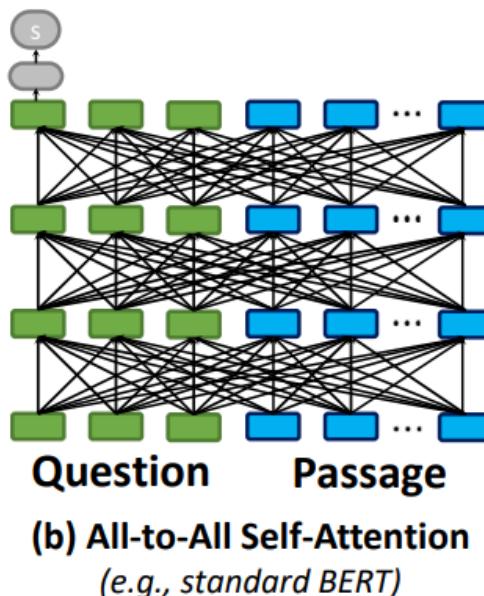
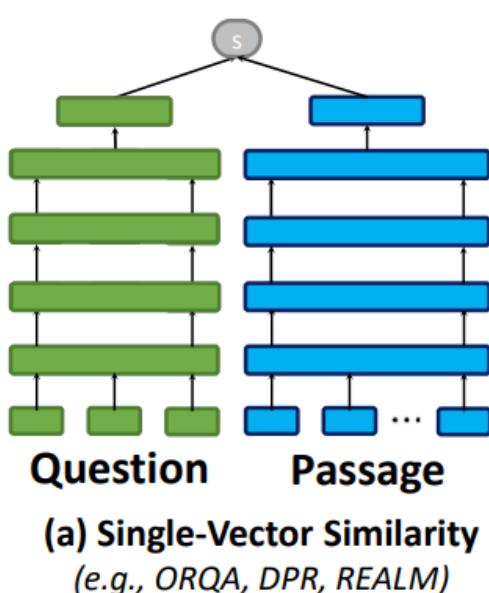
## Dense Retriever – Vector Databases



- Fast *top-k* vector similarity search - on GPU
- Example Databases (ANN Search Libraries):
  - [FAISS](#)
  - [Vector.dev](#)
  - [q-drant](#)
  - [Chroma](#)
  - [Milvus](#)

# Dense Retriever – Beyond Dot-Product

## Similarity using Siamese networks



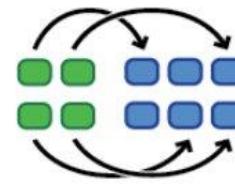
# CoBERT

Not a single document-vector:  
a set of token-vectors

Enables better matching  
between Query and Document

- CoBERT v2:
  - Distilled mode
  - vector compression

[CoBERT Inference in the Browser \(aiserv.cloud\)](#)



# CoBERT

CoBERT query-passage scoring interpretability

Query:

Effects of climate change on marine ecosystems

Passage:

The changing climate has profound impacts on marine ecosystems. Rising temperatures, ocean acidification, and altered precipitation patterns all contribute to shifts in the distribution and behavior of marine species, influencing the delicate balance of underwater ecosystems.

Run CoBERT scoring for query - passage

MaxSim Score: 27.71

Estimated Relevance: 86.60%

## Contextualised Highlights

The changing climate has profound impacts on marine ecosystems. Rising temperatures, ocean acidification, and altered precipitation patterns all contribute to shifts in the distribution and behavior of marine species, influencing the delicate balance of underwater ecosystems.

 Vespa

# Retrieval Types – A Comparison

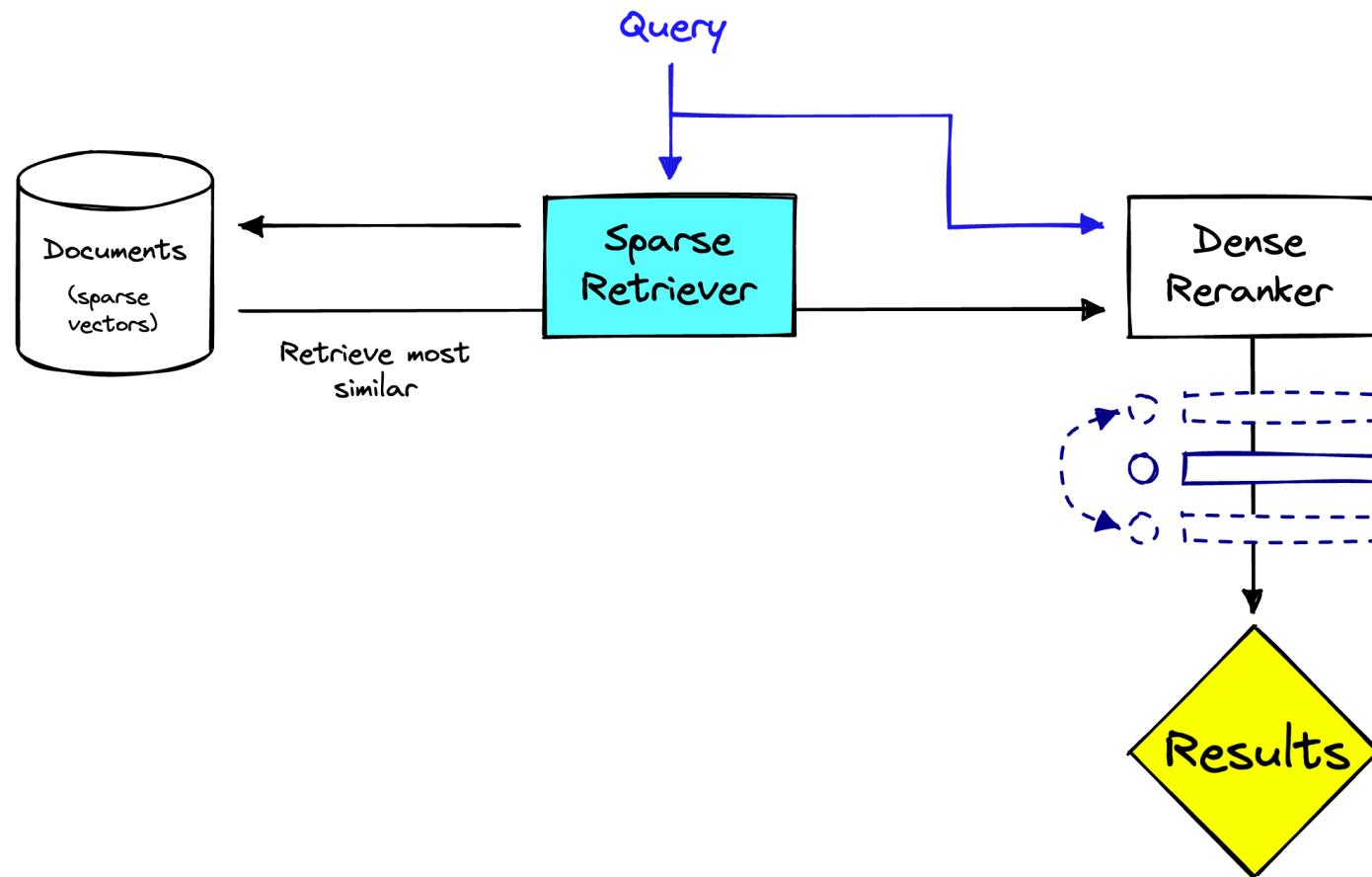
## Sparse

- Pros:
  - Allows fast retrieval
  - No fine-tuning required
  - Exact term matching
  - Works well for long documents
  - Good baseline
- Cons:
  - Vocabulary Mismatch Problem
  - No Semantic similarity
  - Performance quite fixed

## Dense

- Pros:
  - Multi-modal (and cross-modal)
  - Allow semantic similarity
  - Performs well (with training)
  - Multi-lingual
- Cons:
  - Requires fine-tuning  
(with a labeled training-set)
  - Can't generalize between domains
  - Require more compute & memory
  - Don't allow exact match
  - interpretability

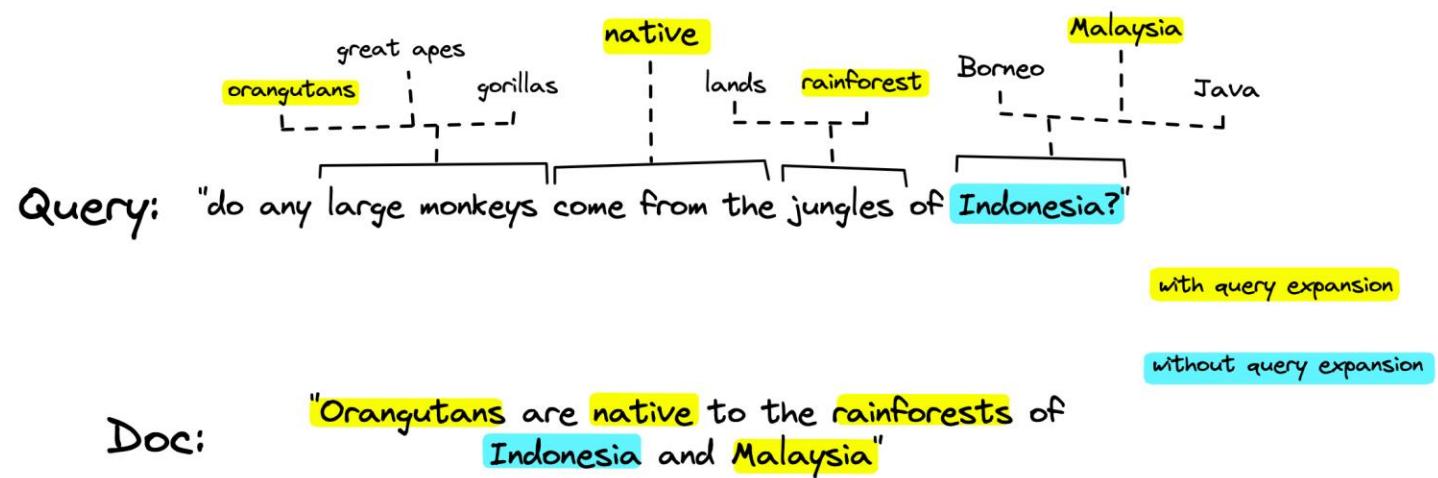
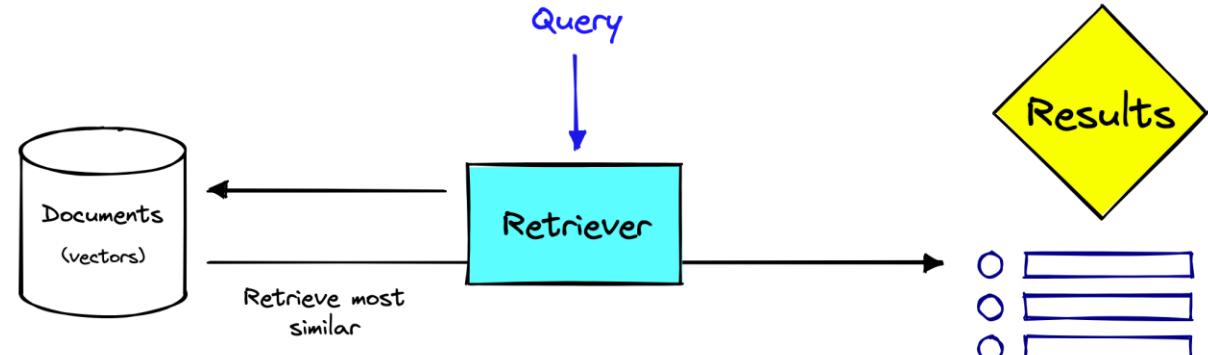
# Sparse Retriever + Dense Re-ranker



# Retrieval: SotA

- SPLADE — the **Sparse Lexical and Expansion** model
  - Query Expansion using LM distribution
  - A semi-sparse vector with synonyms & word variations

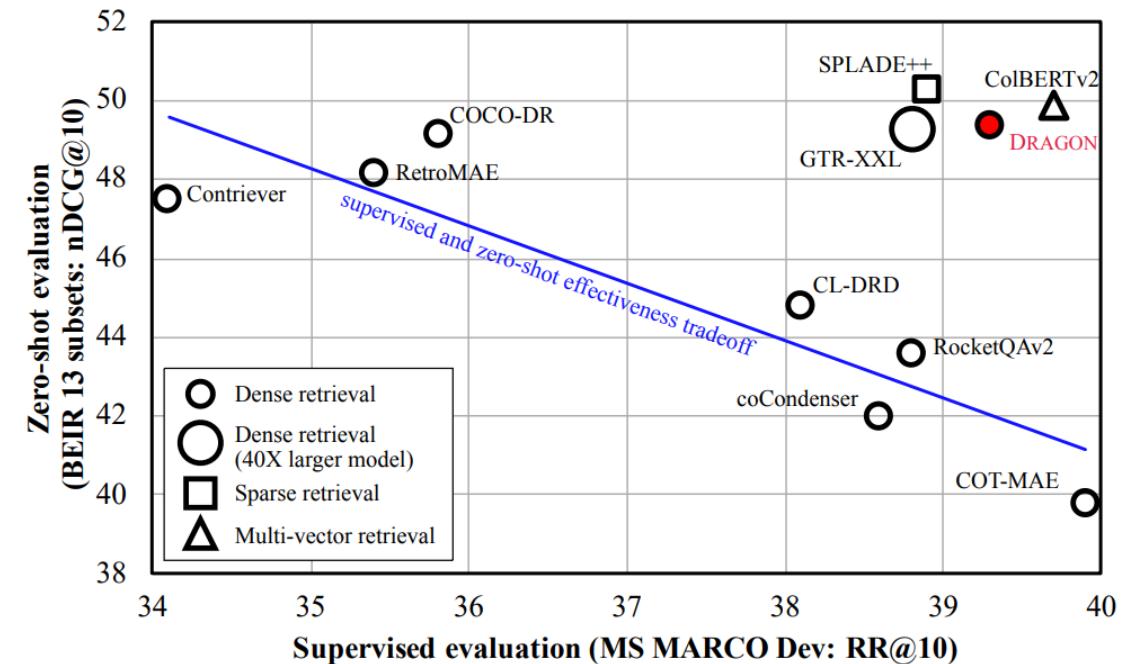
- SPLADE V2
  - Max-pooling
  - Model Distillation
- SPLADE++
  - Negative Sampling



- [SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking](#)
- [SPLADE v2: Sparse Lexical and Expansion Model for Information Retrieval](#)

# Retrieval: SotA

- Zero-Shot Transfer Learning works well between **tasks**, but less between **domains** (different vocabulary)

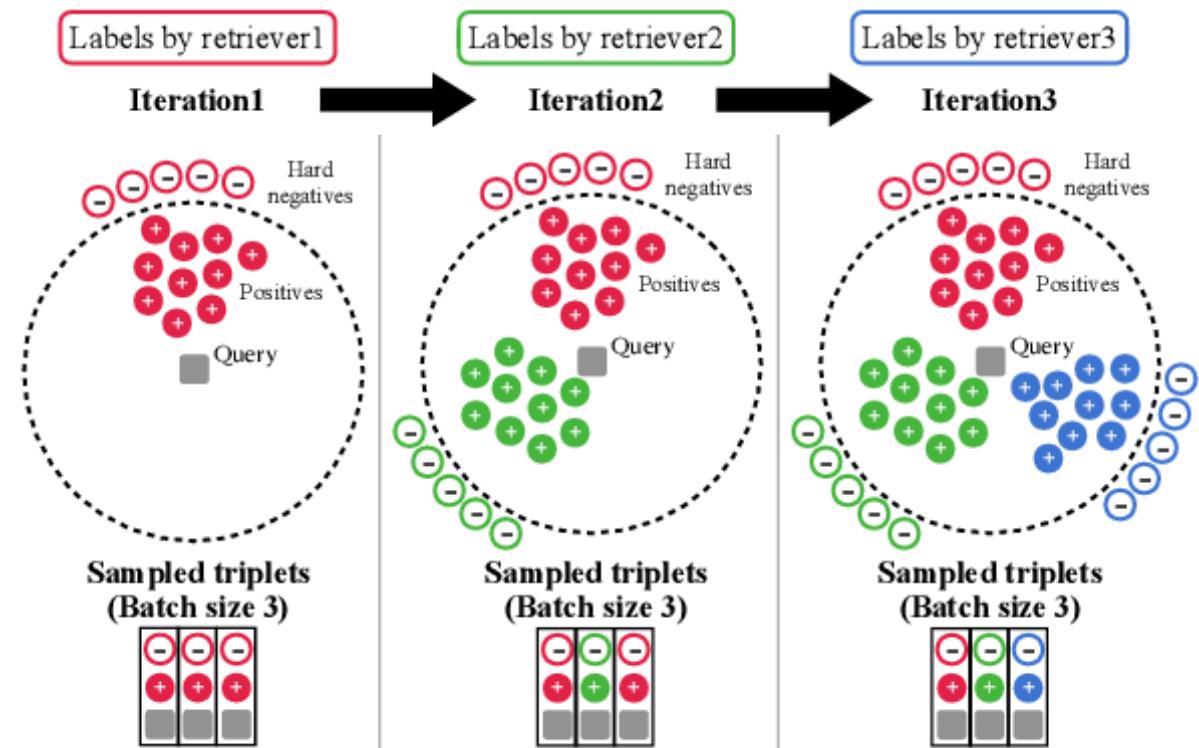


[How to Train Your Dragon: Diverse Augmentation Towards Generalizable Dense Retrieval](#)

# Retrieval: SotA

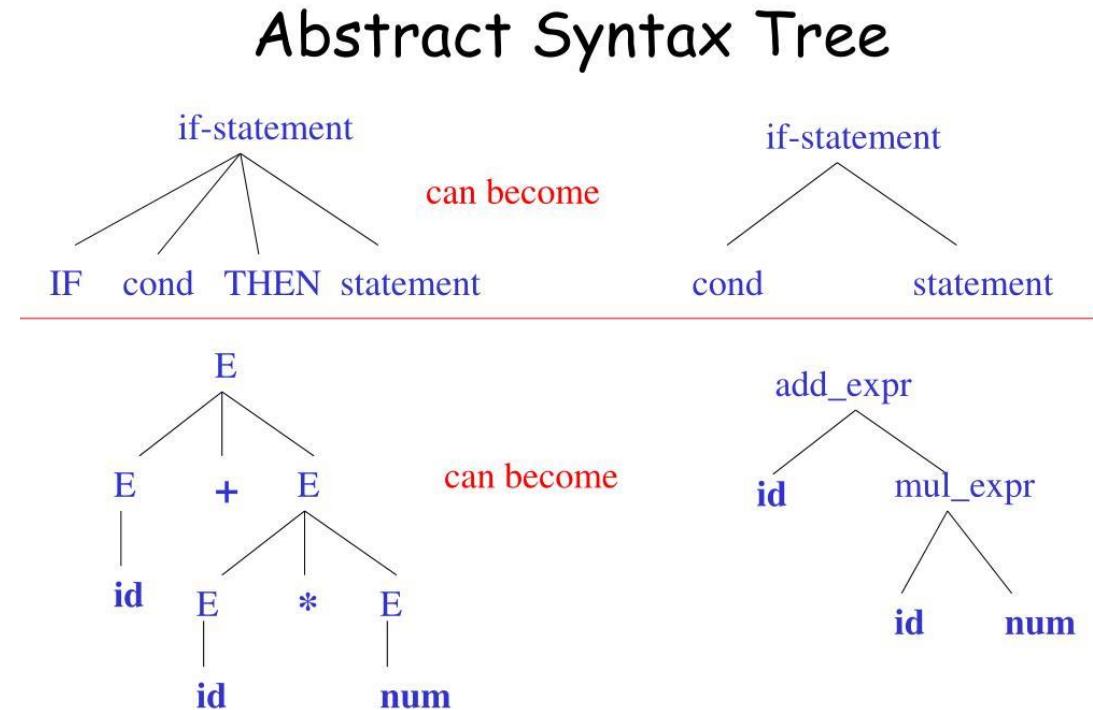
- Dragon:  
**Dense Retriever trained with diverse AuGmentatiON**

- Hybrid training a DR with iterative augmentation from:
  - Sparse Retriever
  - Dense Retriever
  - Multi-Vector Retriever
- Positive + Hard-Negative Sampling
- 28m queries; 5 Days on A100
- Available on HF



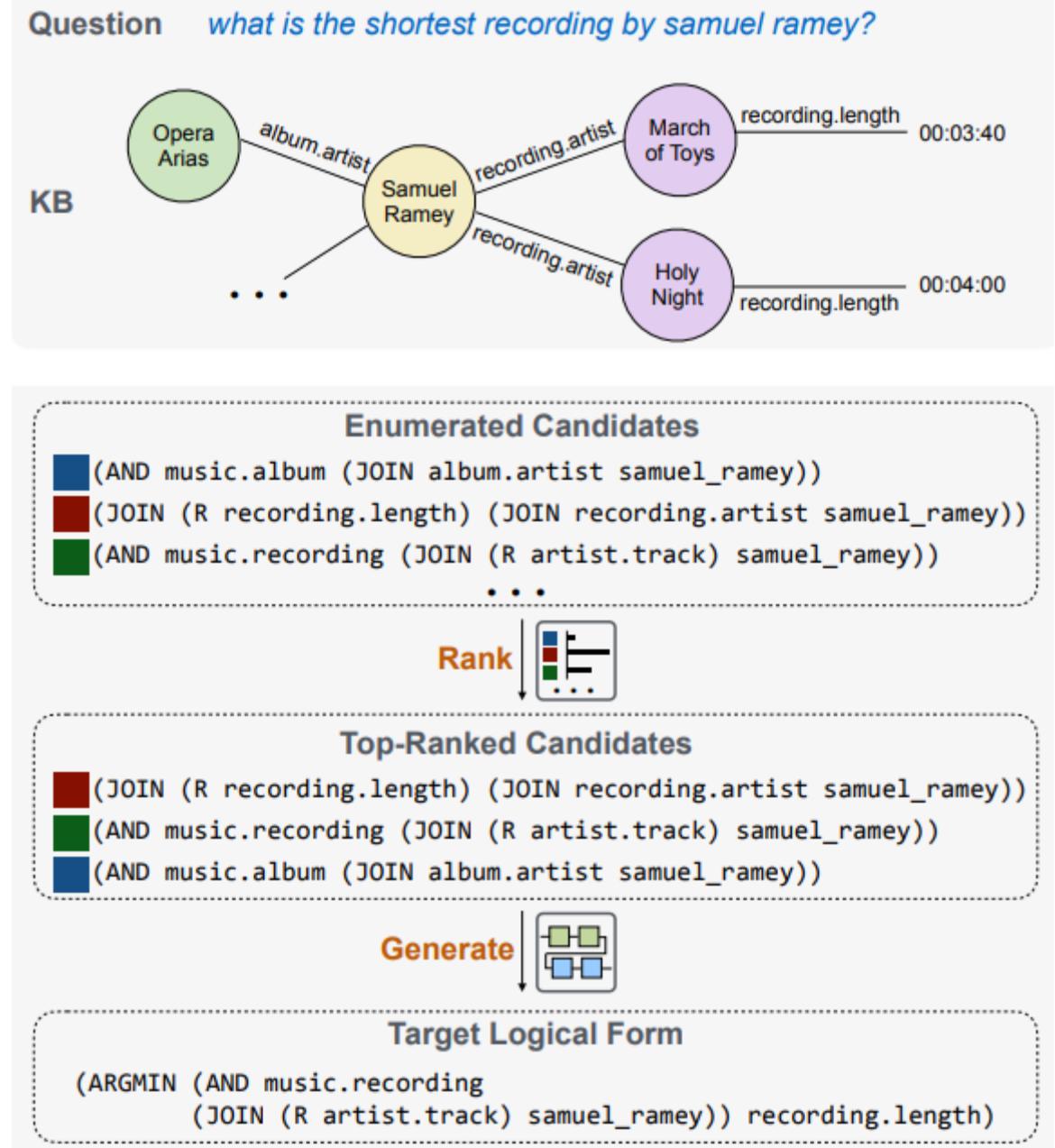
# Side-Note: Other Retrievers

- Code:
  - Abstract Syntax Trees
- **ReCode: [Retrieval-Based Neural Code Generation \(arxiv.org\)](#)**
- Target Similarity Tuning (**TST**):  
**[Synchromesh: Reliable code generation from pre-trained language models \(arxiv.org\)](#)**



# Other Retrievers

- Knowledge Base:
  - Rank And Generate: [RnG-KBQA: Generation Augmented Iterative Ranking for Knowledge Base Question Answering \(arxiv.org\)](#)
  - ECBRF: [End-to-end Case-Based Reasoning for Commonsense Knowledge Base Completion - ACL Anthology](#)



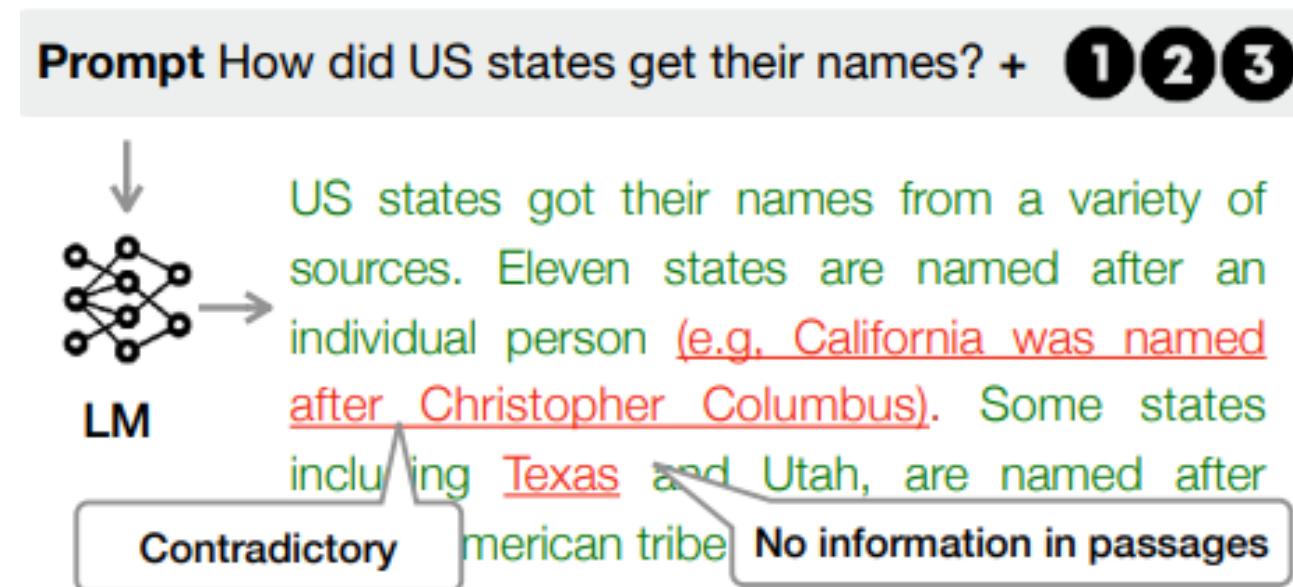
A dark, moody photograph of a man in a suit and tie, looking down with a weary expression. He is holding a dark briefcase in his left hand. The background is blurred.

# Self-Control

When the retrieved documents are wrong

# Critical Models – Self-Rag

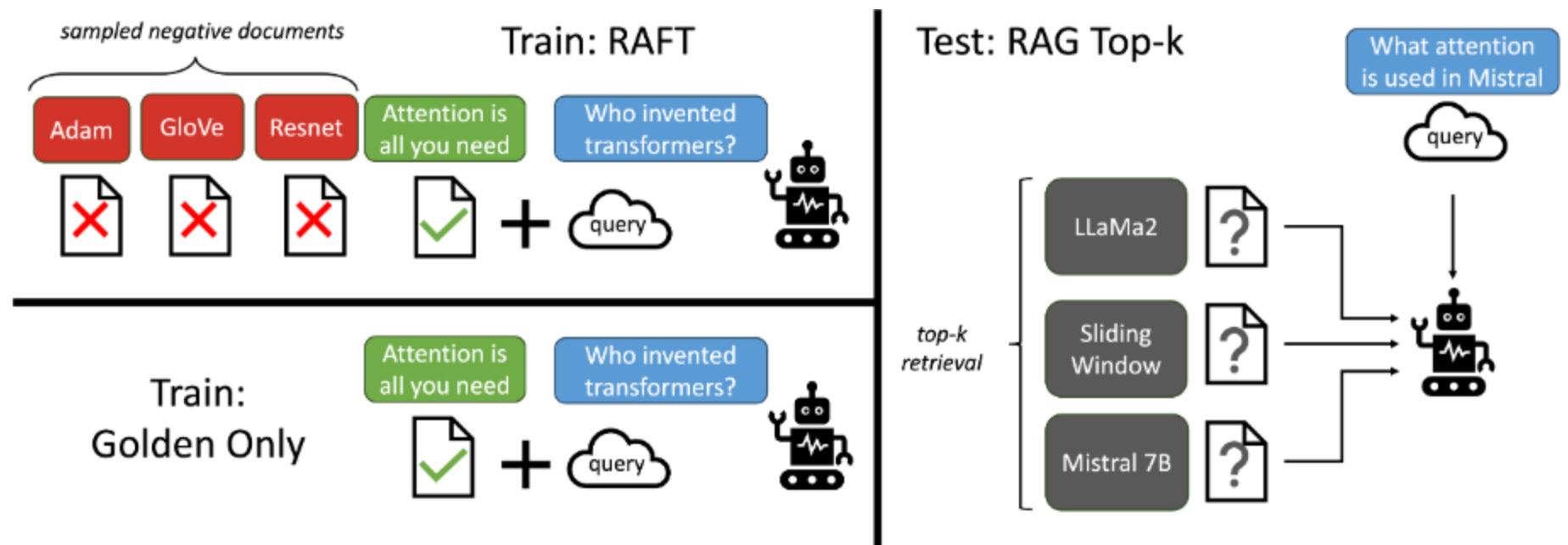
Retrains the LM to evaluate the retrieved content to be relevant to the query:



# Critical Models - RAFT



Refines LLM with distracting negative documents and  
Chain of Thoughts (CoT)



# Critical Models - RAFT



**Question:** The Oberoi family is part of a hotel company that has a head office in what city?

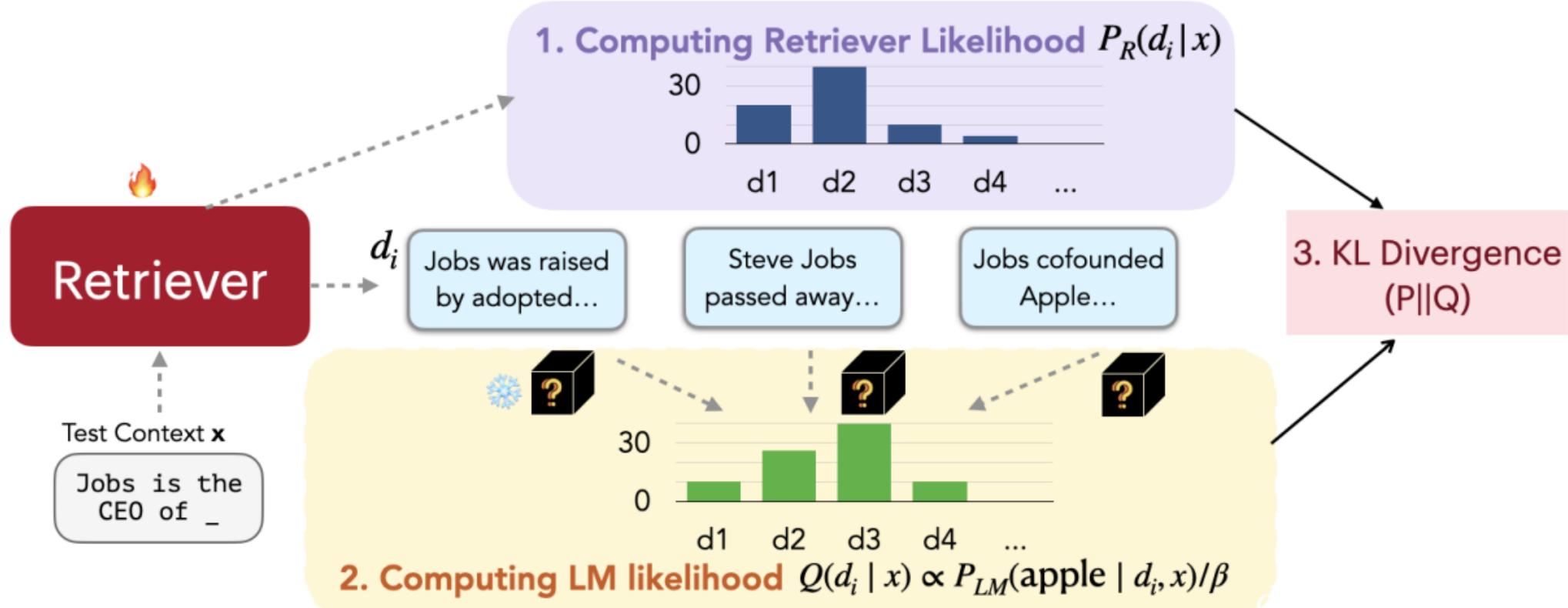
**Context:** [The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group]...[It is located in city center of Jakarta, near Mega Kuningan, adjacent to the sister JW Marriott Hotel. It is operated by The Ritz-Carlton Hotel Company. The complex has two towers that comprises a hotel and the Airlangga Apartment respectively]...[The Oberoi Group is a hotel company with its head office in Delhi.]

**Instruction:** Given the question, context and answer above, provide a logical reasoning for that answer. Please use the format of: **##Reason:** {reason} **##Answer:** {answer}.

**CoT Answer:** **##Reason:** The document **##begin\_quote##** The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group. **##end\_quote##** establishes that the Oberoi family is involved in the Oberoi group, and the document **##begin\_quote##** The Oberoi Group is a hotel company with its head office in Delhi. **##end\_quote##** establishes the head office of The Oberoi Group. Therefore, the Oberoi family is part of a hotel company whose head office is in Delhi. **##Answer:** Delhi

# Critical Models – RePlug (+LSR)

Uses Black-Boxed LLMs to adapt retrievers



# Critical Models - RALM

AI21 labs

- Retrieval-Augmented Language Modeling
- Frozen RAG with BM25
- Training a predictive **re-ranker** to filter the retrieved documents

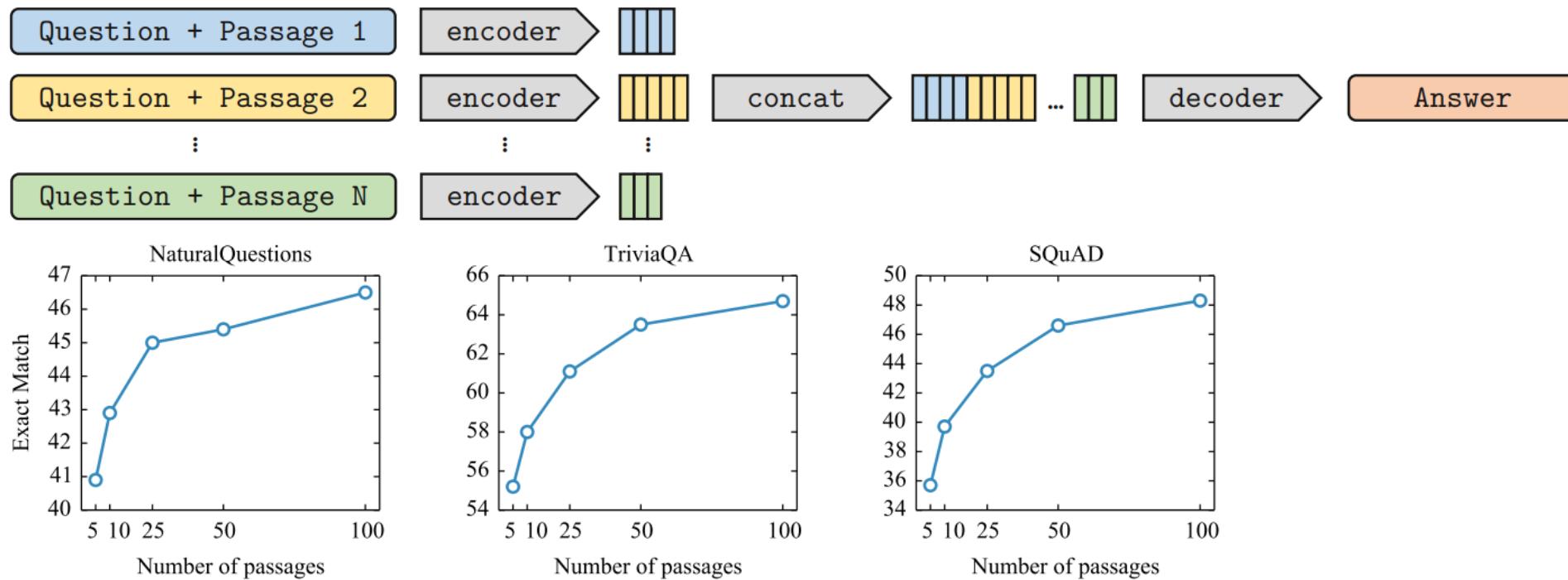


A close-up photograph of a golden retriever lying on a light-colored floor. The dog has long, wavy, reddish-brown fur and is looking directly at the camera with a calm expression. The background is a plain, light-colored wall.

# Contextualization of Retriever & Generator

# FiD

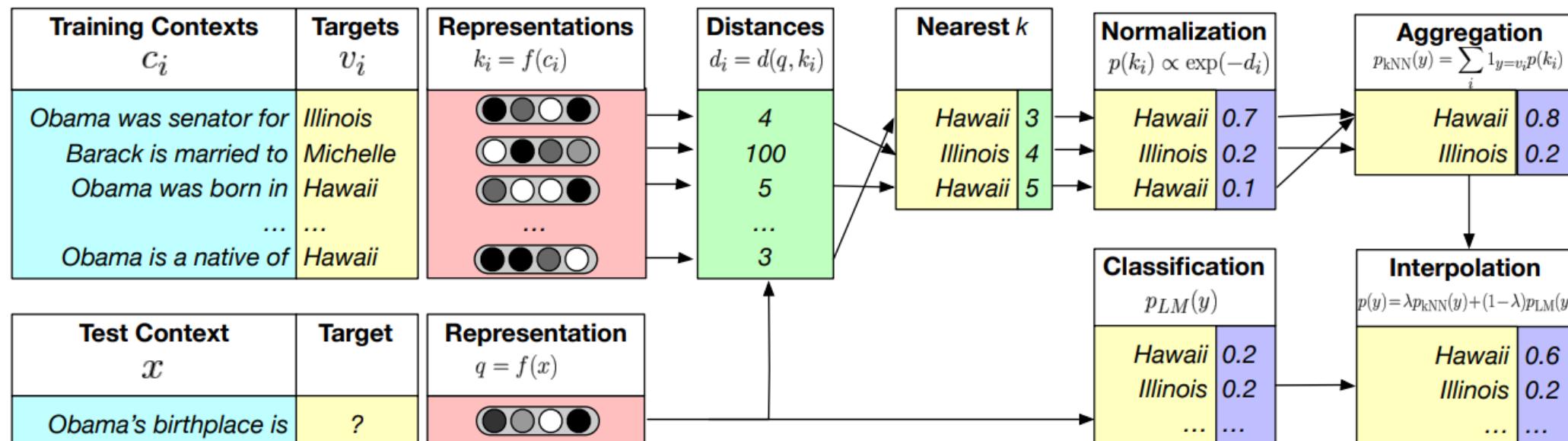
- Overcomes the document limitation of the original RAG
- Encoding the passages, concatenating, decoding into an answer



[2007.01282] Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering  
(arxiv.org)

# $k$ NN LM

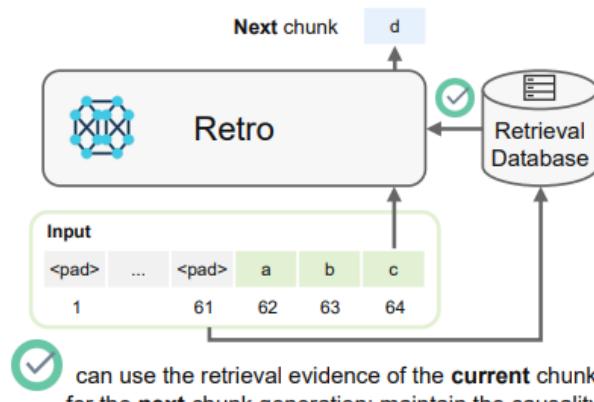
- Cluster the “next-token” with  $k$ NN



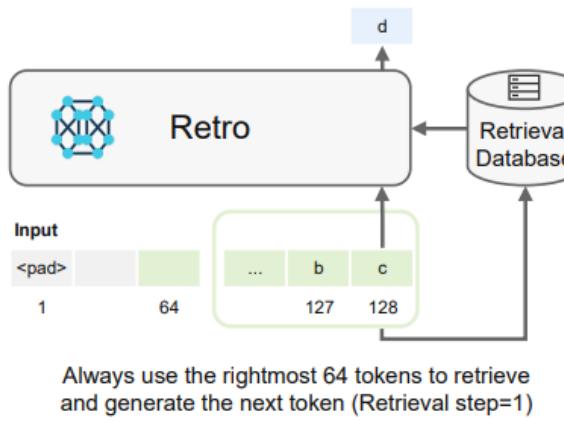
[1911.00172] Generalization through Memorization: Nearest Neighbor Language Models (arxiv.org)

# Retro++

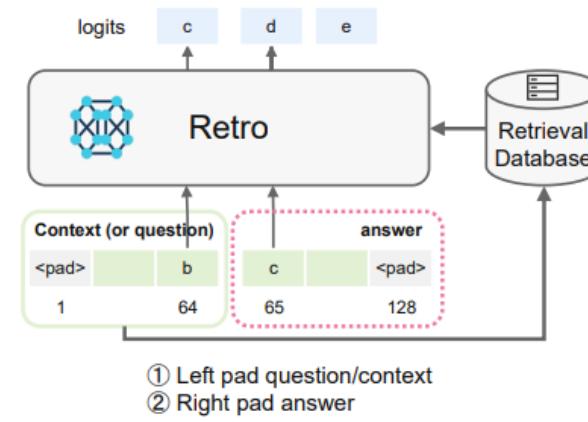
- An open version of the original *Retro* (DeepMind)
- Decoder Only LM
- Performs retrieval *during* answer generation



(a) Use “left padding” Rule

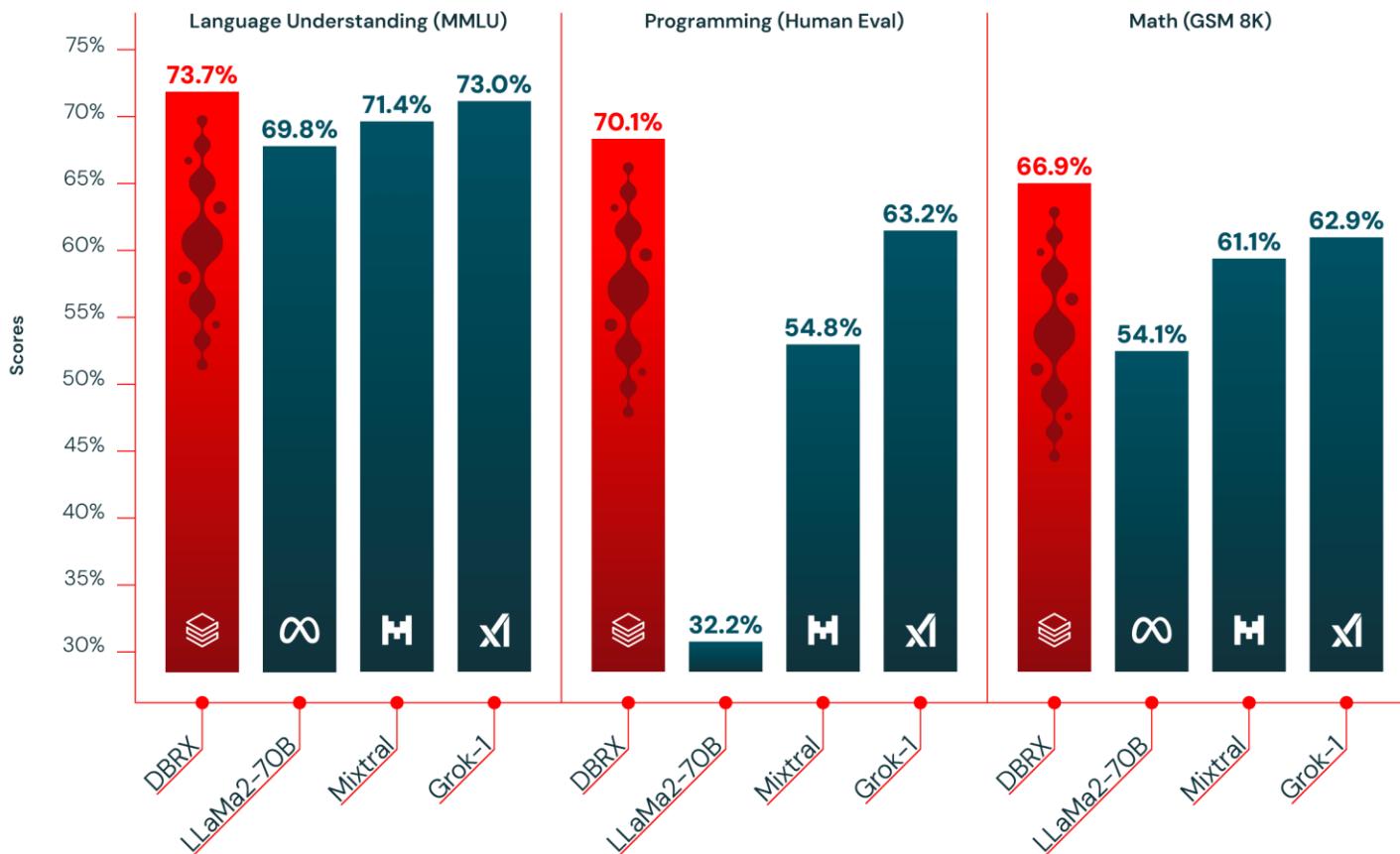


(b) Retrieval step = 1



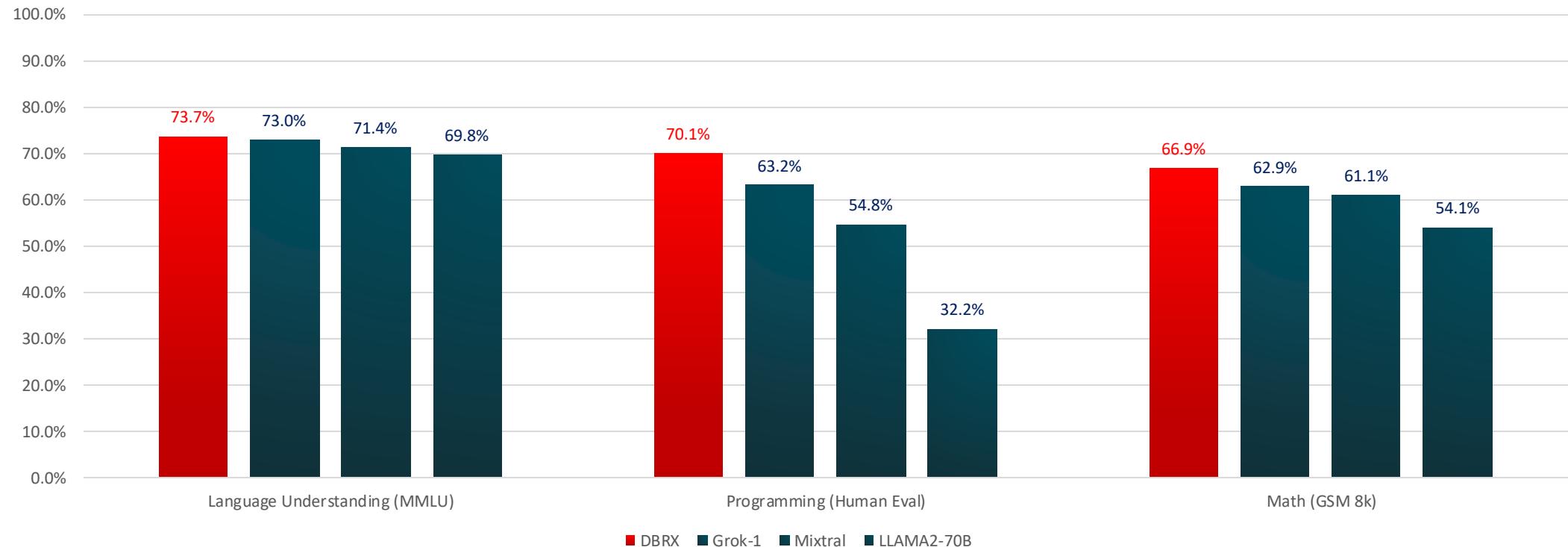
(c) Separate question and answer chunks

# DBRX (formerly MosaicML)

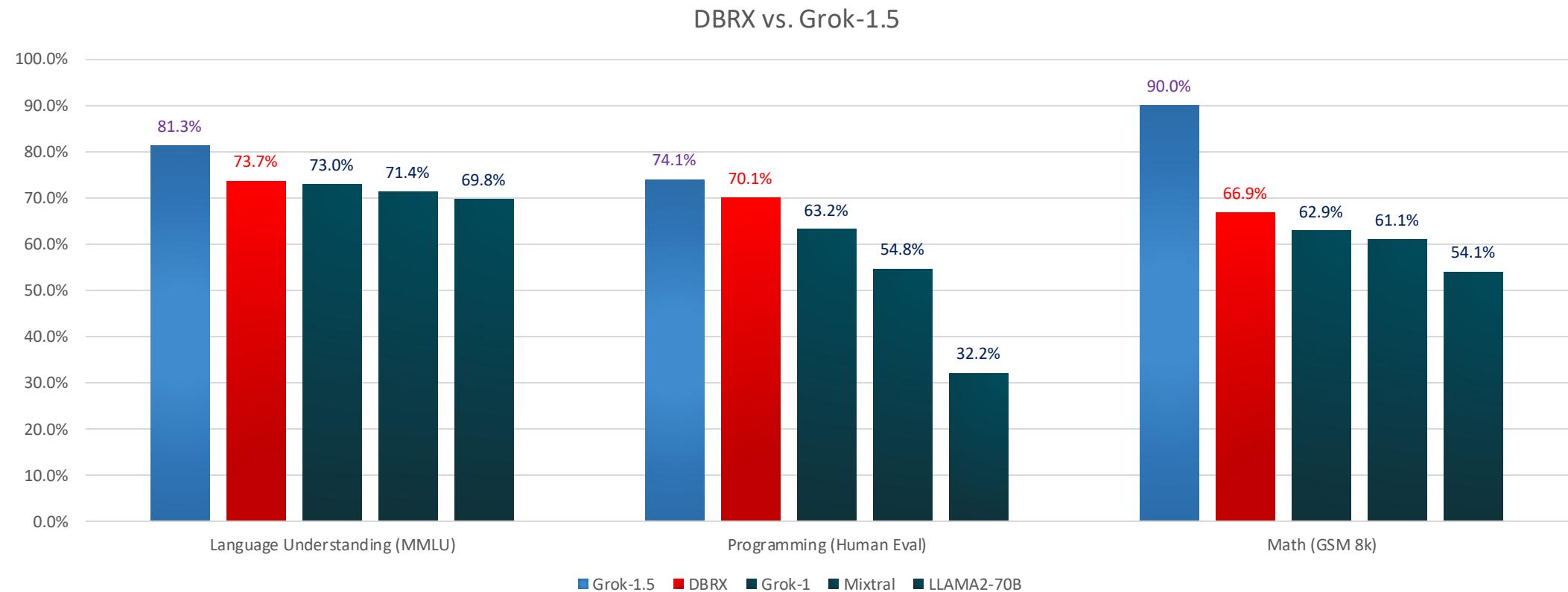


# DBRX – Correct Perspective

DBRX vs Other Models - True Chart



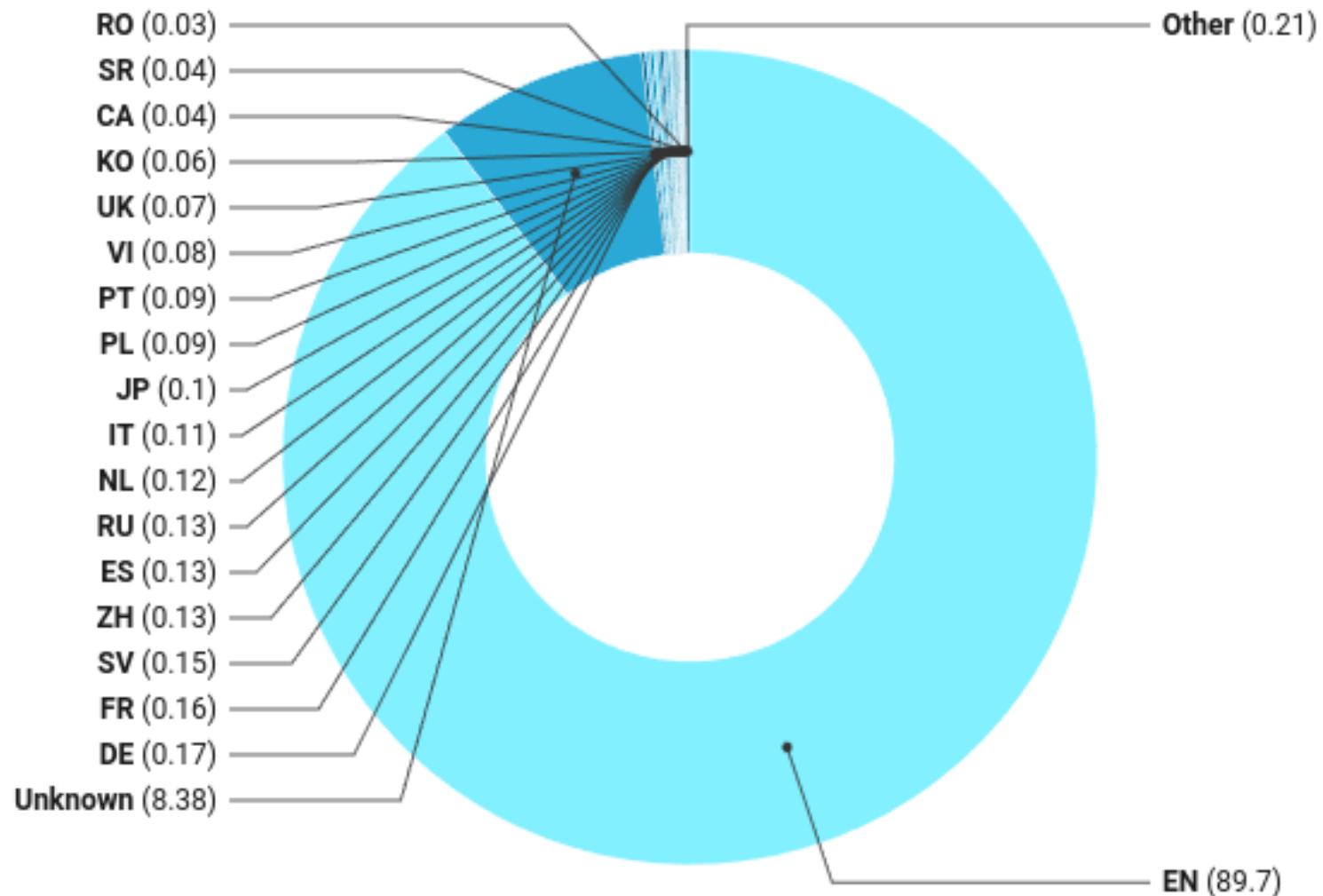
# DBRX vs. Grok-1.5



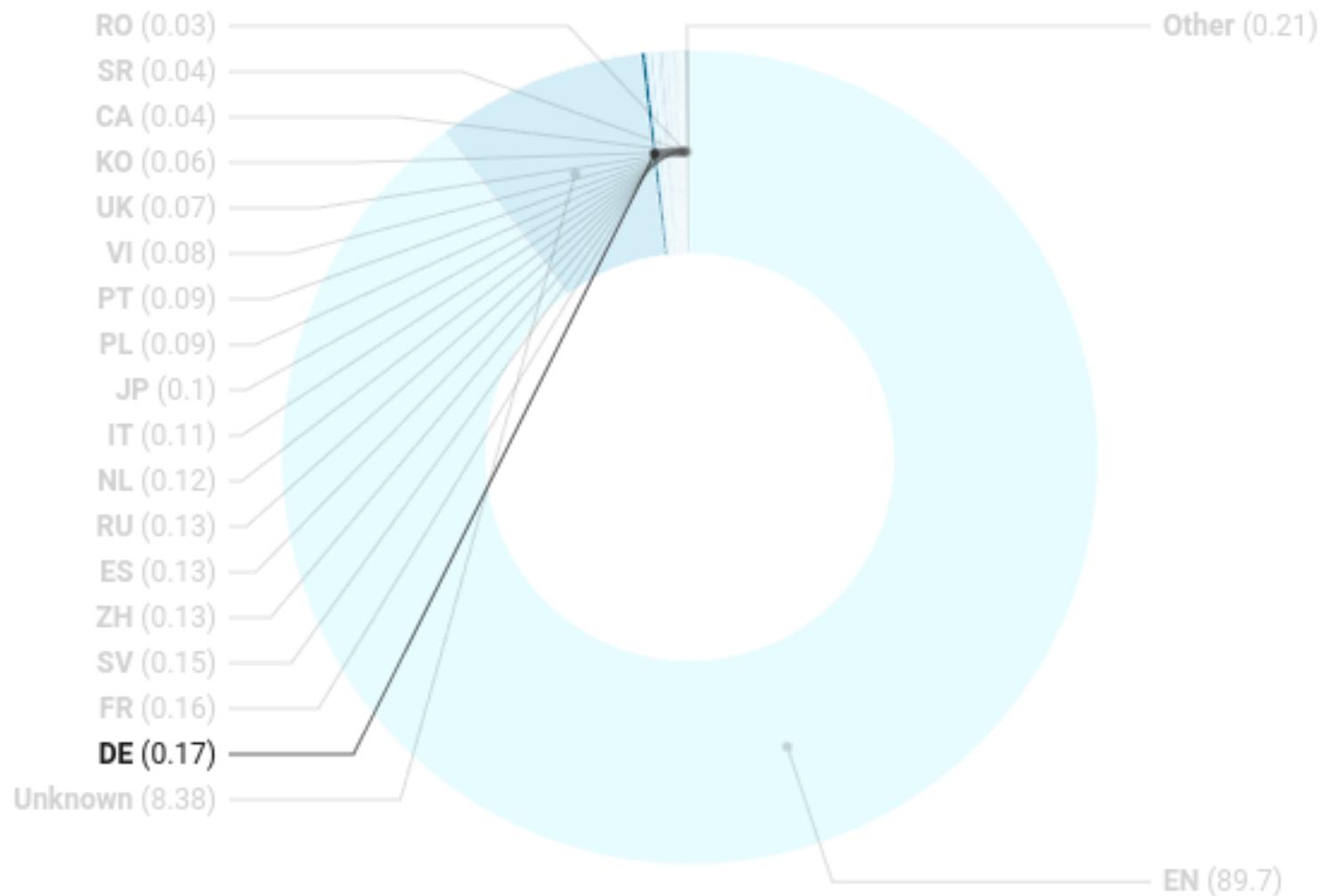
# DBRX – Open-Source 132B Parameters

| Model   | DBRX Instruct | Mixtral Instruct | GPT-3.5 Turbo (API) | GPT-4 Turbo (API) |
|---|---------------|------------------|---------------------|-------------------|
| <b>Answer in Beginning<br/>Third of Context</b> | <u>45.1%</u>  | 41.3%            | 37.3%*              | <b>49.3%</b>      |
| <b>Answer in Middle<br/>Third of Context</b>    | <u>45.3%</u>  | 42.7%            | 37.3%*              | <b>49.0%</b>      |
| <b>Answer in Last Third<br/>of Context</b>      | <u>48.0%</u>  | 44.4%            | 37.0%*              | <b>50.9%</b>      |
| <b>2K Context</b>                               | 59.1%         | <u>64.6%</u>     | 36.3%               | <b>69.3%</b>      |
| <b>4K Context</b>                               | <u>65.1%</u>  | 59.9%            | 35.9%               | 63.5%             |
| <b>8K Context</b>                               | <u>59.5%</u>  | 55.3%            | 45.0%               | <b>61.5%</b>      |
| <b>16K Context</b>                              | 27.0%         | 20.1%            | <u>31.7%</u>        | 26.0%             |
| <b>32K Context</b>                              | <u>19.9%</u>  | 14.0%            | —                   | <b>28.5%</b>      |

# Llama 2's Pretrain-Data Language distribution



# Llama 2's Pretrain-Data Language distribution



# Debugging RAG

---

- Embedding
  - Hammock: [colehaus/hammock-public: Visualize text embeddings \(github.com\)](https://github.com/colehaus/hammock-public)
  - RagExplorer - [gabrielchua/RAGxplorer: Open-source tool to visualise your RAG \(github.com\)](https://github.com/gabrielchua/RAGxplorer) 

# Debugging RAG

Visualise which chunks are most relevant to your query.

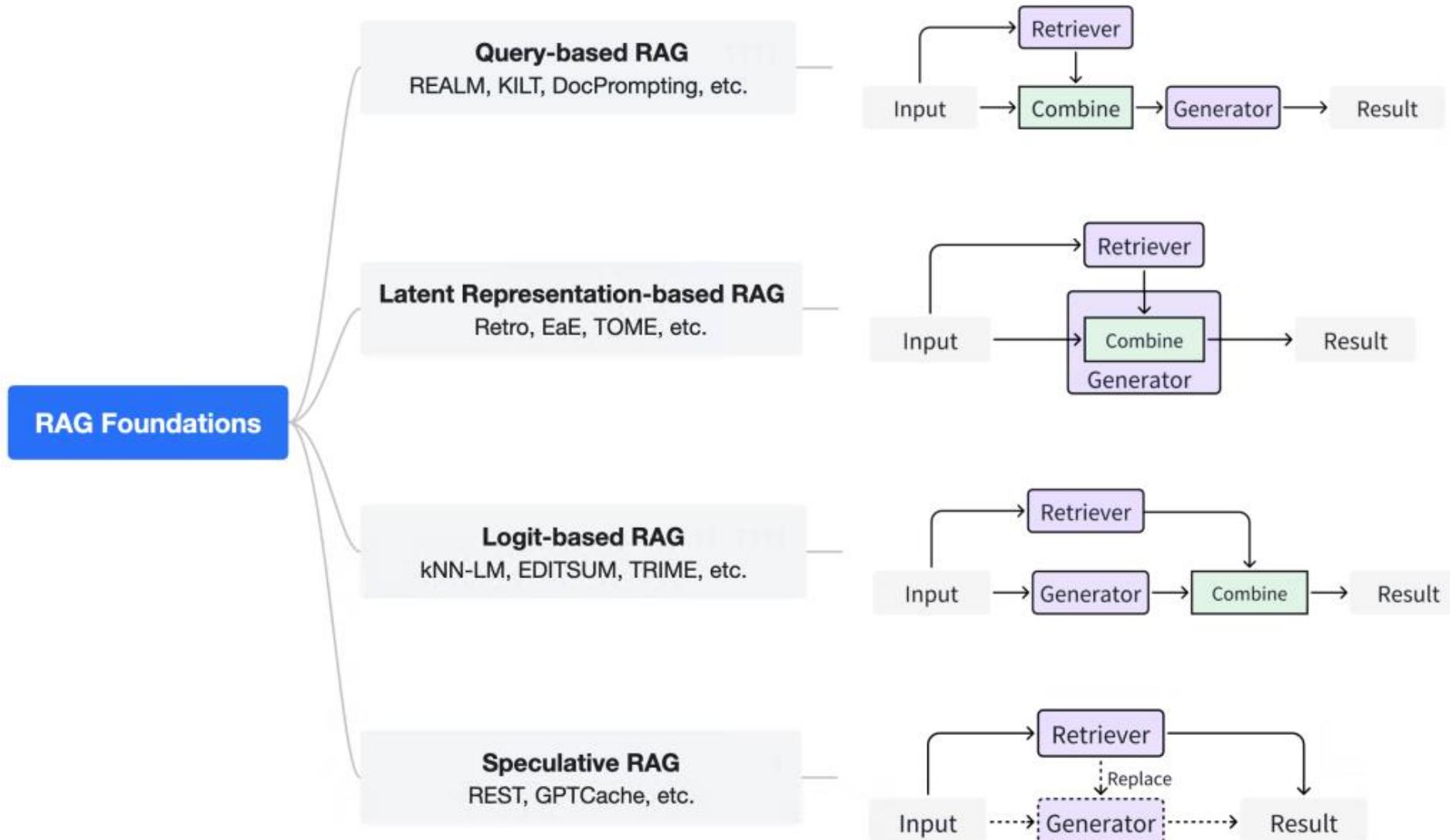
Enter your query

What are the top revenue drivers for Microsoft?





And this was only the tip of the iceberg...



## RAG Enhancements



# Choices...



## Indexing

- Data Preprocessing method
- Indexing Model
- Text-Splitting
- Chunking hyperparams

## Storage

- Vector Database
- Adding Metadata

## Retrieval

- Retrieval Method
- Top-K
- Similarity Cut-off

## Synthesis and Generation

- Model choice
- Prompt
- Hyperparameters

## Evaluation

- Evaluation method
- Evaluator prompts

# Take-Aways

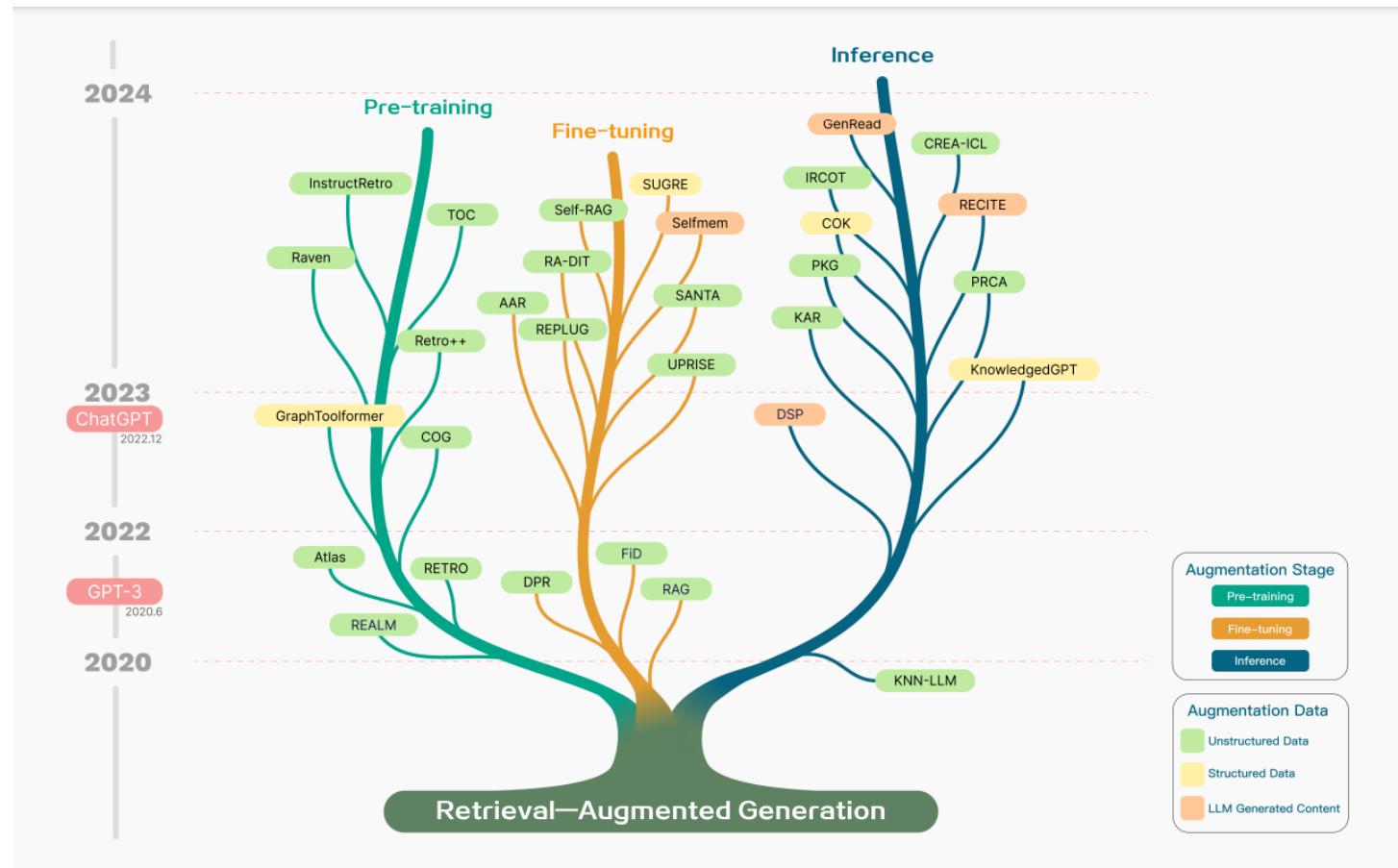
---

- RAG is a *user interface* for search
  - Search and fine the results through chat-bot LLMs
  - Hard to evaluate automatically (See Aaron Kaplan's talk)
- Not a single end-2-end model, but a complex of components
  - Hence, hard(er) to tune and reach high accuracy
- Huge community around it (and hype too)
  - Many new advances every week
- No one-solution-fits-all

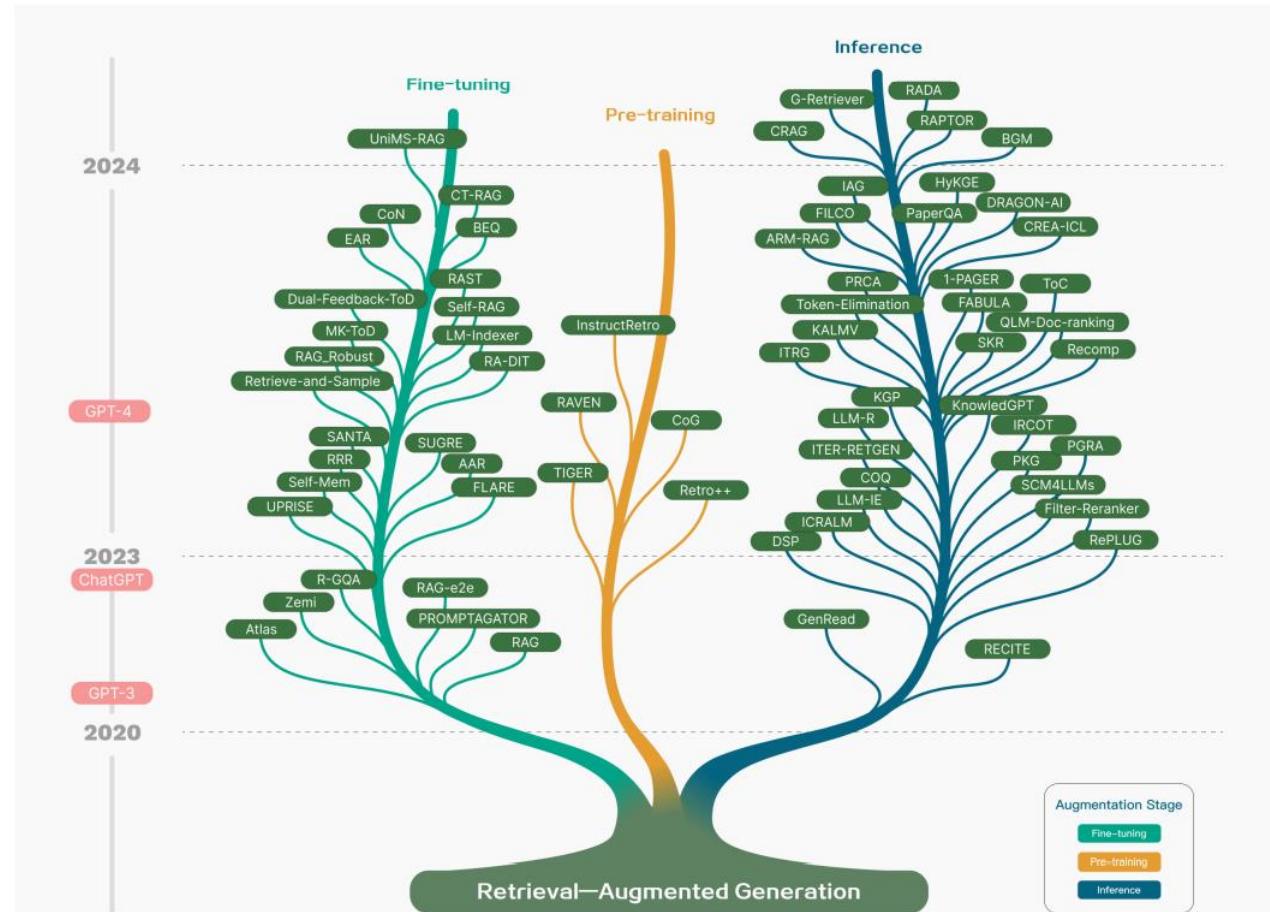


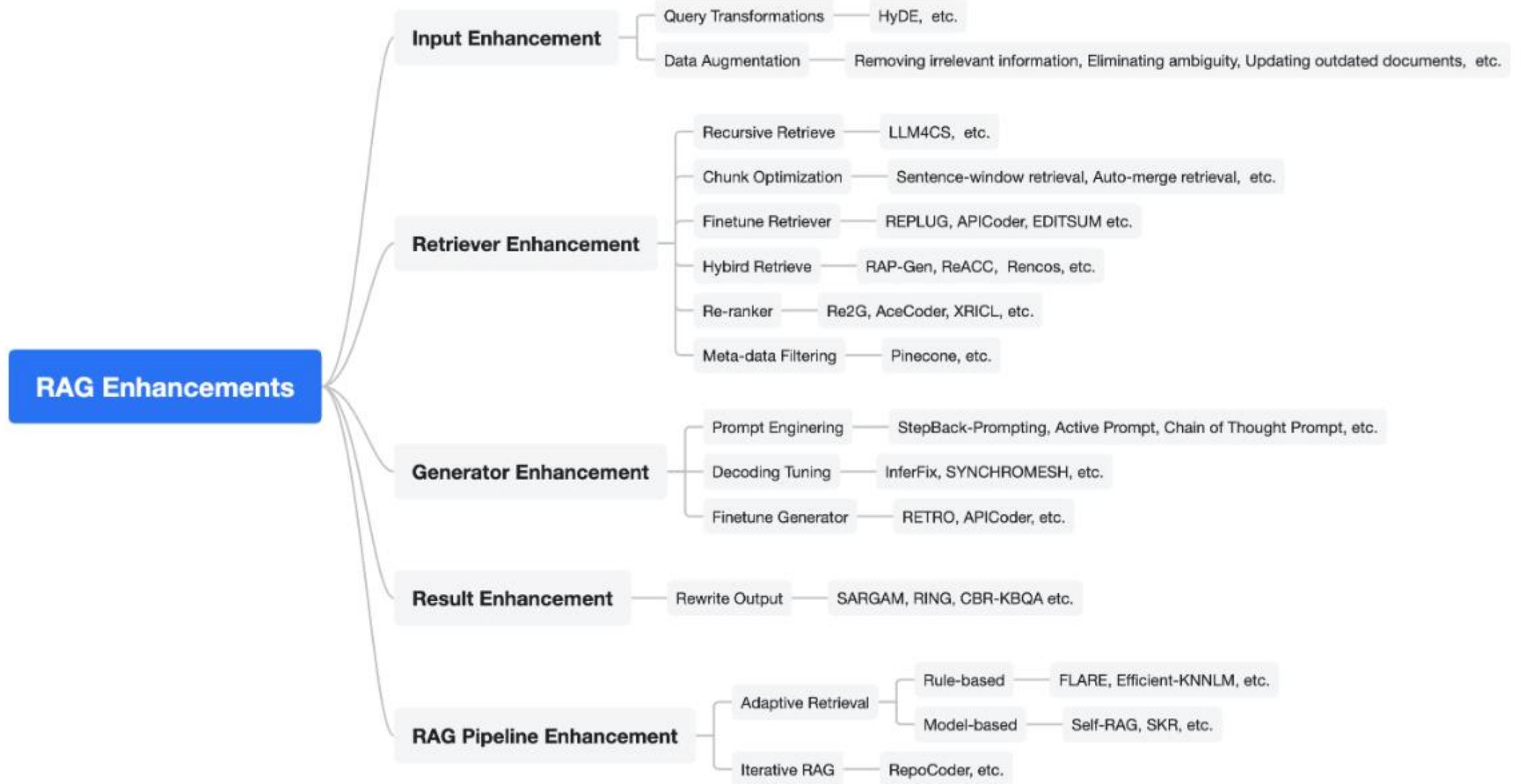
Thank you for your  
attention!

# Supplements – RAG Timeline (Jan 2024)



# Supplements – RAG Timeline (Apr 2024)





## RAG Enhancements



# References

- [\[2312.10997v5\] Retrieval-Augmented Generation for Large Language Models: A Survey \(arxiv.org\)](#)
- [\[2312.10997\] Retrieval-Augmented Generation for Large Language Models: A Survey \(arxiv.org\)](#)
- [\[YouTube\] Stanford CS25: V3 | Retrieval Augmented Language Models](#)
- [GitHub - Tongji-KGLLM/RAG-Survey](#)
- [\[2002.08909\] REALM: Retrieval-Augmented Language Model Pre-Training \(arxiv.org\)](#)