# Advanced Retrieval Augmented Generation Techniques
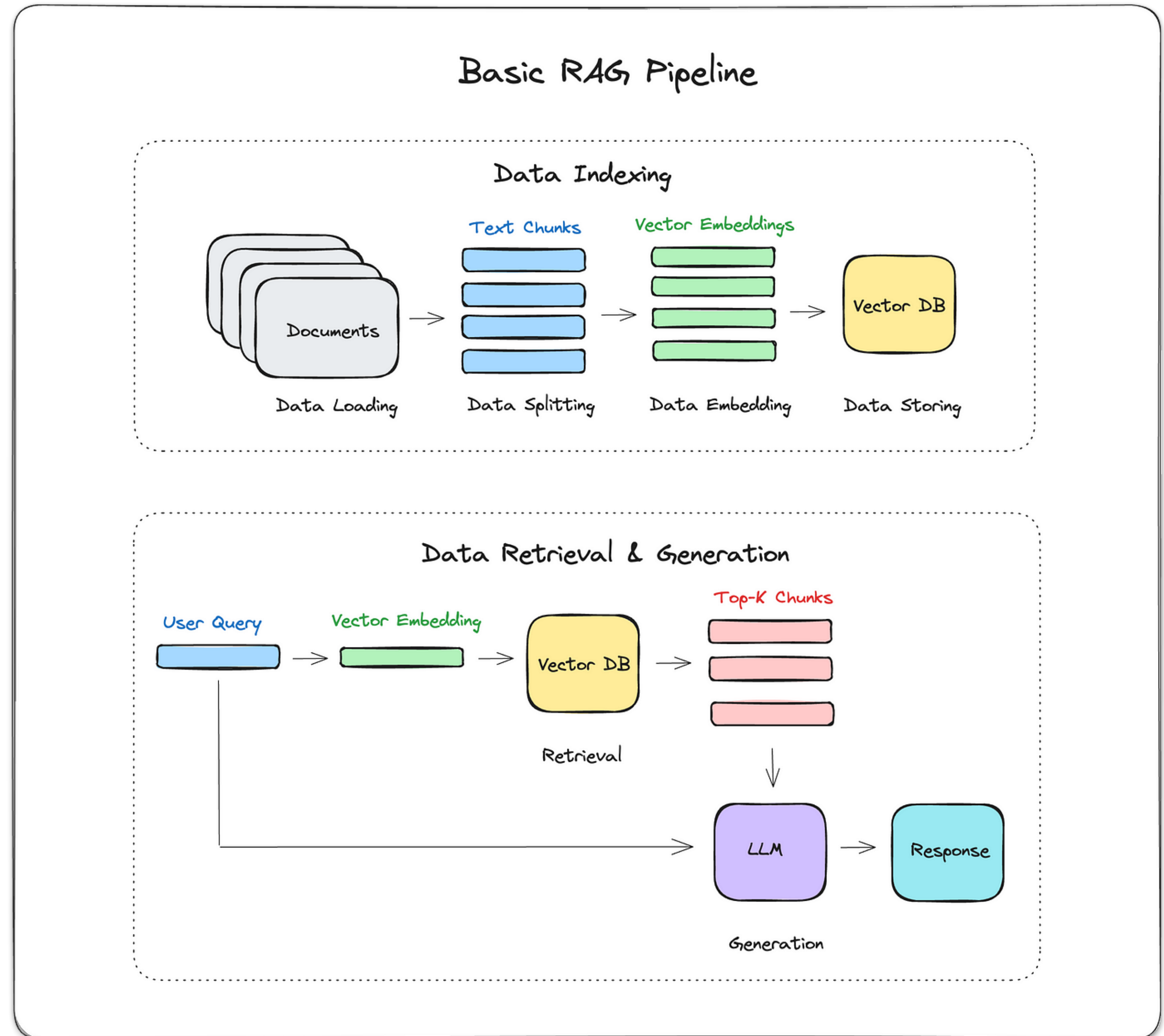
fzeng 2024-10-08

# Outline

- Part 1: Retrieval techniques (today)

  - RAG review

  - Chunking

  - Query optimization

  - Synthetic data

  - Hybrid search

- Part 2: Generation techniques

  - Understanding long contexts

  - Re-ranking

  - Summarization

  - Repacking

  - Planning and multi-step retrieval

- Part 3: Productionizing RAG

  - Evaluations

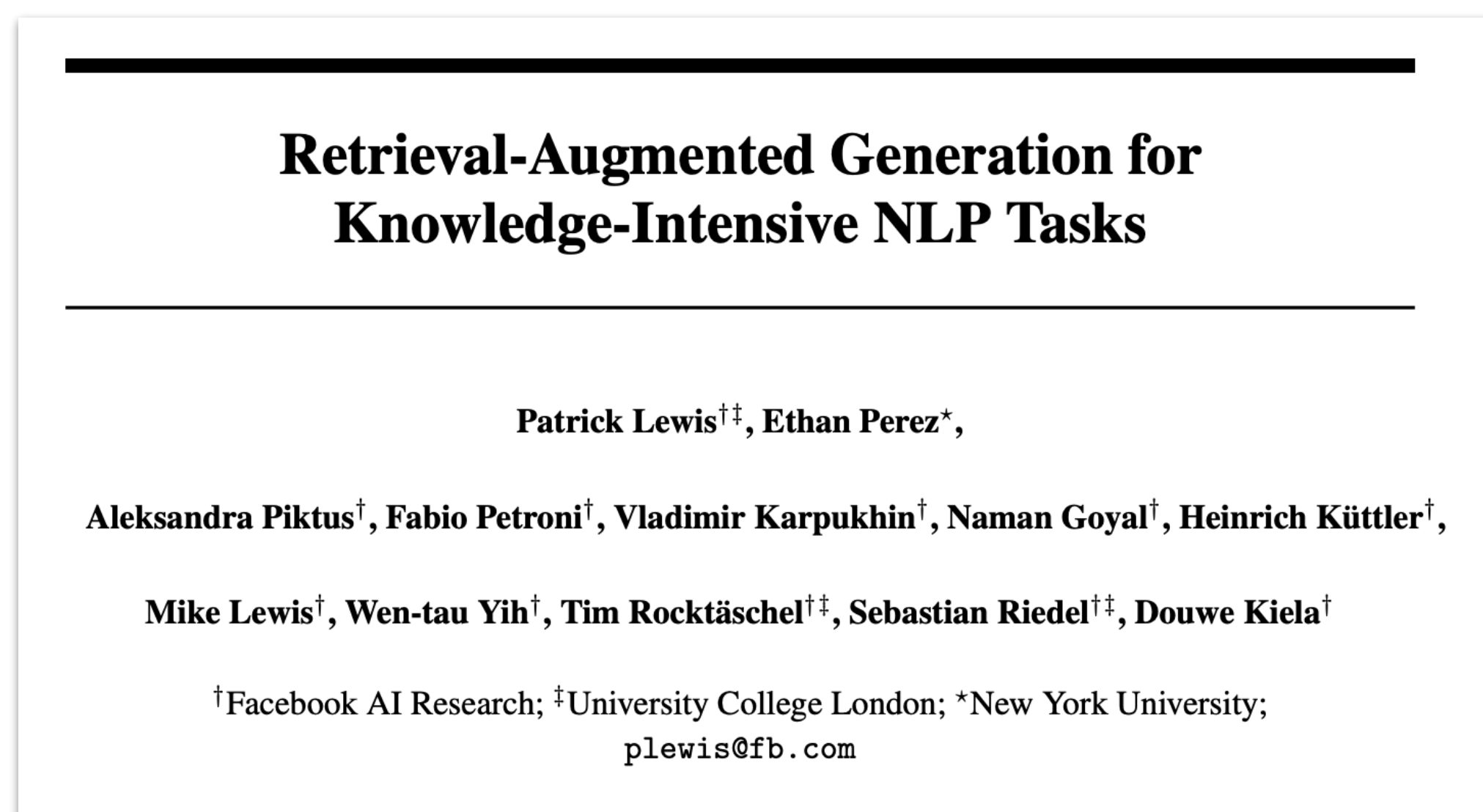  - Fine-tuning retrievers & embeddings

# Naive RAG

# Naive RAG

1. Clean knowledge base & break down into smaller chunks of text

2. Use embedding model to turn chunks of text into embeddings that has semantic meaning in a vector space

3. Store embeddings in vector database which supports nearest-neighbor search

4. At runtime, query is embedded by embedding model, & looks up similar embeddings in vector database
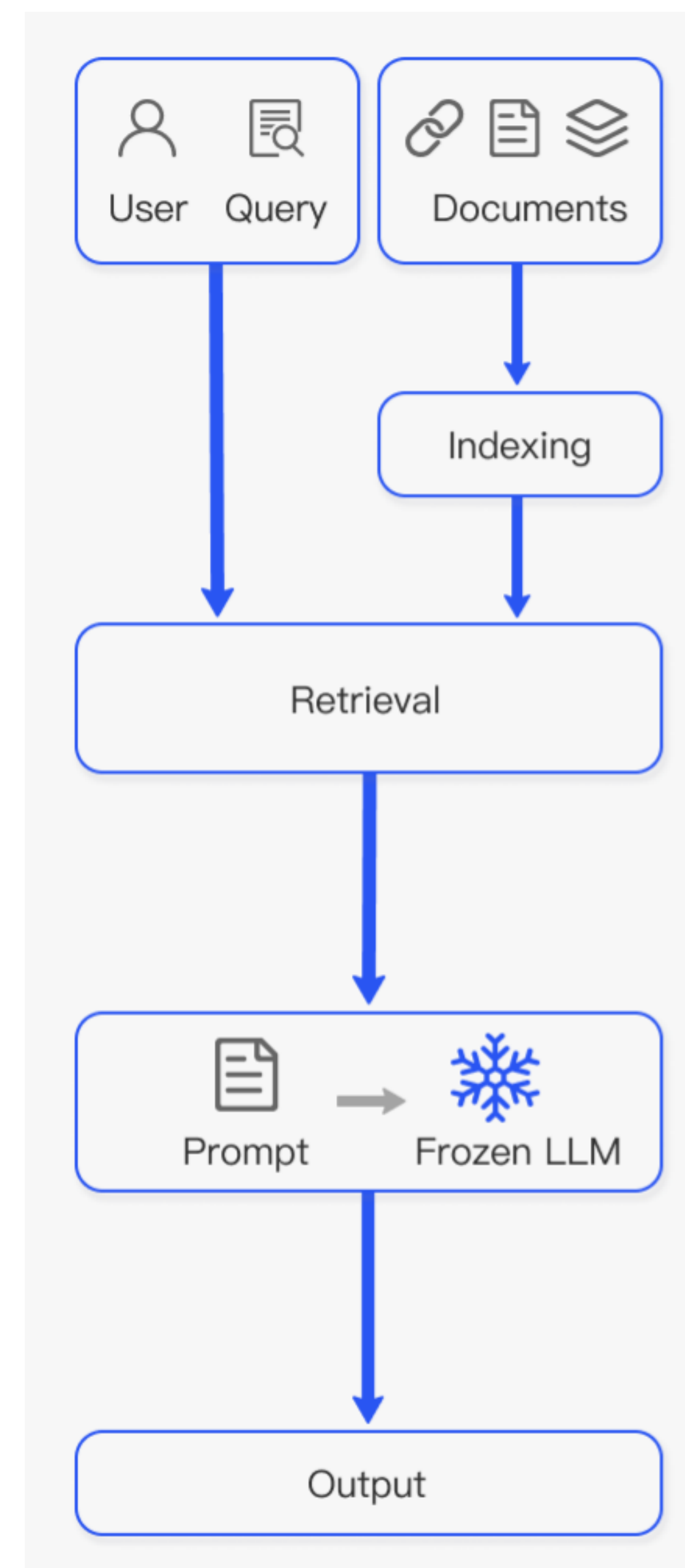
# Naive RAG

- 3 main steps: indexing, retrieval, generation

- "Retrieve-then-read" paradigm

- Combines parametric memory (i.e stored in weights) with non-parametric memory

## Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks

**Patrick Lewis**[†‡], **Ethan Perez**[⋆],

**Aleksandra Piktus**[†], **Fabio Petroni**[†], **Vladimir Karpukhin**[†], **Naman Goyal**[†], **Heinrich Küttler**[†],

**Mike Lewis**[†], **Wen-tau Yih**[†], **Tim Rocktäschel**[†‡], **Sebastian Riedel**[†‡], **Douwe Kiela**[†]

[†]Facebook AI Research; [‡]University College London; [⋆]New York University;
plewis@fb.com

The paper that started it all

# Naive RAG

- Many problems remain in practice...

  - Not retrieving relevant/useful chunks

  - Re-ranking retrieved chunks

  - Hallucinations during generation

  - Lack of confidence & trust in generated content

- This talk series will give you an arsenal of tools to tackle these, and give you inspiration for unique approaches for your specific problem

# Part 1: Retrieval

# Chunking

# Basic Chunking Strategies

- Fixed token without overlap

| A style guide | is about consistency | . Consistency | with this style | guide is important |

- Fixed token with overlap

| A style guide | guide is about | about consistency | . Consistency | istency with this |

- Recursive with overlap: split tokens by some set of characters (i.e newlines), recursively merge text until chunk size limit

| A style guide | guide is about | about consistency. | Consistency with | with this style |

Figures from MongoDB
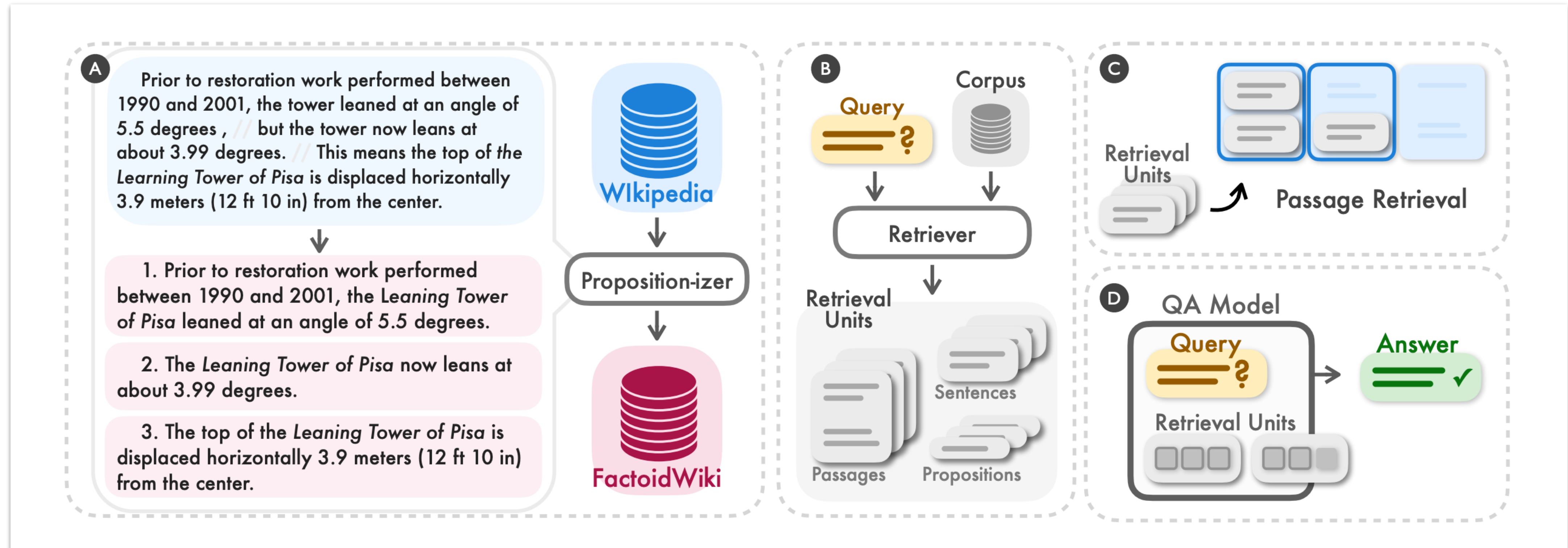
# Chunking Granularity

- What retrieval granularity to use?

  - Document

  - Passage

  - Sentence

  - Proposition? 🤔

- Proposition: atomic expressions encapsulating a distinct factoid, are concise and self-contained

| | |
|---|---|
| Question: What is the angle of the Tower of Pisa? | |
| Passage Retrieval | Prior to restoration work performed between 1990 and 2001, the tower leaned at an angle of 5.5 degrees, but the tower now leans at about 3.99 degrees. This means the top of the Leaning Tower of Pisa is displaced horizontally 3.9 meters (12 ft 10 in) from the center. |
| Sentence Retrieval | Prior to restoration work performed between 1990 and 2001, the tower leaned at an angle of 5.5 degrees, but the tower now leans at about 3.99 degrees. |
| Proposition Retrieval | The Leaning Tower of Pisa now leans at about 3.99 degrees. |

Dense X Retrieval: What Retrieval Granularity Should We Use? (Chen et al, 2023)

# Proposition-Level Chunks
## Technique overview

• Generated proposition-level chunks from Wikipedia, FactoidWiki:



Dense X Retrieval: What Retrieval Granularity Should We Use? (Chen et al, 2023)

# Proposition-Level Chunks
## Prompts for generation

**Passage ⇒ Propositions**

Decompose the "Content" into clear and simple propositions, ensuring they are interpretable out of context.

1. Split compound sentence into simple sentences. Maintain the original phrasing from the input whenever possible.

2. For any named entity that is accompanied by additional descriptive information, separate this information into its own distinct proposition.

3. Decontextualize the proposition by adding necessary modifier to nouns or entire sentences and replacing pronouns (e.g., "it", "he", "she", "they", "this", "that") with the full name of the entities they refer to.

4. Present the results as a list of strings, formatted in JSON.

**Input**: Title: Ēostre. Section: Theories and interpretations, Connection to Easter Hares. Content: The earliest evidence for the Easter Hare (Osterhase) was recorded in south-west Germany in 1678 by the professor of medicine Georg Franck von Franckenau, but it remained unknown in other parts of Germany until the 18th century. Scholar Richard Sermon writes that "hares were frequently seen in gardens in spring, and thus may have served as a convenient explanation for the origin of the colored eggs hidden there for children. Alternatively, there is a European tradition that hares laid eggs, since a hare's scratch or form and a lapwing's nest look very similar, and both occur on grassland and are first seen in the spring. In the nineteenth century the influence of Easter cards, toys, and books was to make the Easter Hare/Rabbit popular throughout Europe. German immigrants then exported the custom to Britain and America where it evolved into the Easter Bunny."

**Output**: [ "The earliest evidence for the Easter Hare was recorded in south-west Germany in 1678 by Georg Franck von Franckenau.", "Georg Franck von Franckenau was a professor of medicine.", "The evidence for the Easter Hare remained unknown in other parts of Germany until the 18th century.", "Richard Sermon was a scholar.", "Richard Sermon writes a hypothesis about the possible explanation for the connection between hares and the tradition during Easter", "Hares were frequently seen in gardens in spring.", "Hares may have served as a convenient explanation for the origin of the colored eggs hidden in gardens for children.", "There is a European tradition that hares laid eggs.", "A hare's scratch or form and a lapwing's nest look very similar.", "Both hares and lapwing's nests occur on grassland and are first seen in the spring.", "In the nineteenth century the influence of Easter cards, toys, and books was to make the Easter Hare/Rabbit popular throughout Europe.", "German immigrants exported the custom of the Easter Hare/Rabbit to Britain and America.", "The custom of the Easter Hare/Rabbit evolved into the Easter Bunny in Britain and America." ]

**Input**: *<a new passage>*
**Output**:

Dense X Retrieval: What Retrieval Granularity Should We Use? (Chen et al, 2023)

# Proposition-Level Chunks
## Results

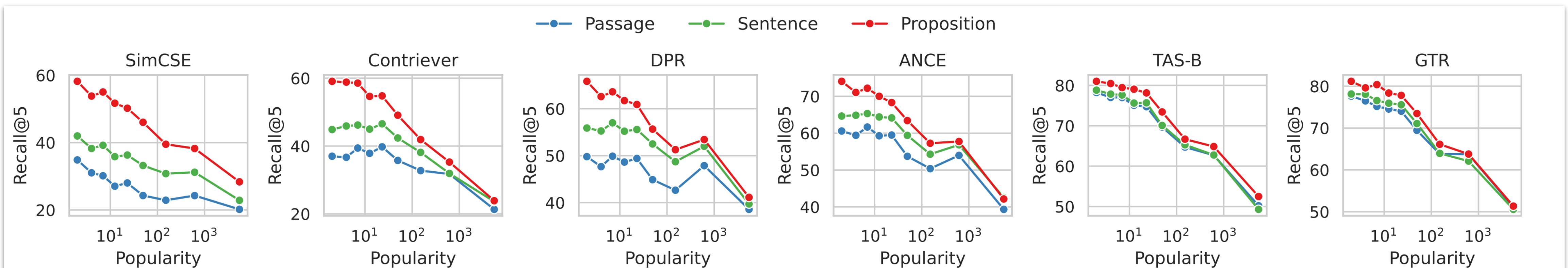- Recall@k metric: # of articles where correct answer is found within top k retrieved passages



Figure 3: Document retrieval recall vs. the popularity of the target entity in each question from the *EntityQuestions* dataset. The popularity of each entity (i.e. smaller value ⇒ less common entities, and vice versa) is estimated by the occurrence of the entity in its top-1000 passage retrieved by BM25. On queries with less common entities, we observe that retrieving by proposition shows a larger advantage over retrieval by proposition.

Dense X Retrieval: What Retrieval Granularity Should We Use? (Chen et al, 2023)

# Analysis

- Q1: passage/sentence retrieved Super Bowl X instead

- Q2: passage fails to match part containing atomic number

- Q3: sentence was truncated and misses out function

- Q4: proposition retrieved irrelevant

| Passage Retrieval | Sentence Retrieval | Proposition Retrieval |
|---|---|---|
| **Q1: What was the theme of Super Bowl 50?** | | |
| Title: Super Bowl X ✗ <br> The overall theme of the Super Bowl entertainment was to celebrate the United States Bicentennial. Each Cowboys and Steelers player wore a special patch with the Bicentennial logo on their jerseys... | Title: Super Bowl X ✗ <br> The overall theme of the Super Bowl entertainment was to celebrate the United States Bicentennial. | Title: Super Bowl XLV ✓ <br> ... As this was the 50th Super Bowl game, the league [Super Bowl 50] emphasized the "golden anniversary" with various gold-themed initiatives during the 2015 season, as well as... |
| **Q2: The atomic number of indium which belongs to 5th period is?** | | |
| Title: Period 5 element ✗ <br> The periodic table is laid out in rows to illustrate recurring (periodic) trends in the chemical behaviour of the elements as their atomic number increases: ... | Title: Period 5 element ✓ <br> Indium is a chemical element with the symbol In and atomic number 49. | Title: Period 5 element ✓ <br> Indium is a chemical element with the symbol In and [Indium has a] atomic number 49. This rare, very soft, malleable ... |
| **Q3: What is the function of the pericardial sac?** | | |
| Title: Pericardium ✓ <br> The pericardium, also called pericardial sac ... It separates the heart from interference of other structures, protects it against infection and blunt trauma, and lubricates the heart's movements. | Title: Pericardium ✗ <br> The pericardium, also called pericardial sac, is a double-walled sac containing the heart and the roots of the great vessels. | Title: Cardiac muscle ✓ <br> On the outer aspect of the myocardium is the epicardium which forms part of the pericardial sac that surrounds, protects, and lubricates the heart. |
| **Q4: What is the main character's name in layer cake?** | | |
| Title: Layer Cake (film) ✓ <br> ... The film's plot revolves around a London-based criminal, played by Daniel Craig, ... Craig's character is unnamed in the film and is listed in the credits as "XXXX". | Title: Angelic Layer ✗ <br> The primary protagonist is Misaki Suzuhara. | Title: Plot twist ✗ <br> Sometimes the audience may discover that the true identity of a character is , in fact, unknown [in Layer Cake] , as in Layer Cake or the eponymous assassins in V for Vendetta and The Day of the Jackal. |

Table 4: Example cases where top-1 retrieved text unit of each retrieval granularity fails to provide the correct answer. The underlined text is the correct answer. The gray text is the context of propositions, but it is for illustration purpose only and not provided to the retrievers and downstream QA models.

# Takeaways

- ✅ Propositions contain higher density of question-related information, outperforms  passage/sentence based-retrieval

- ✅ Minimizes inclusion of extraneous, irrelevant information for downstream tasks

Dense X Retrieval: What Retrieval Granularity Should We Use? (Chen et al, 2023)

# Semantic Chunking

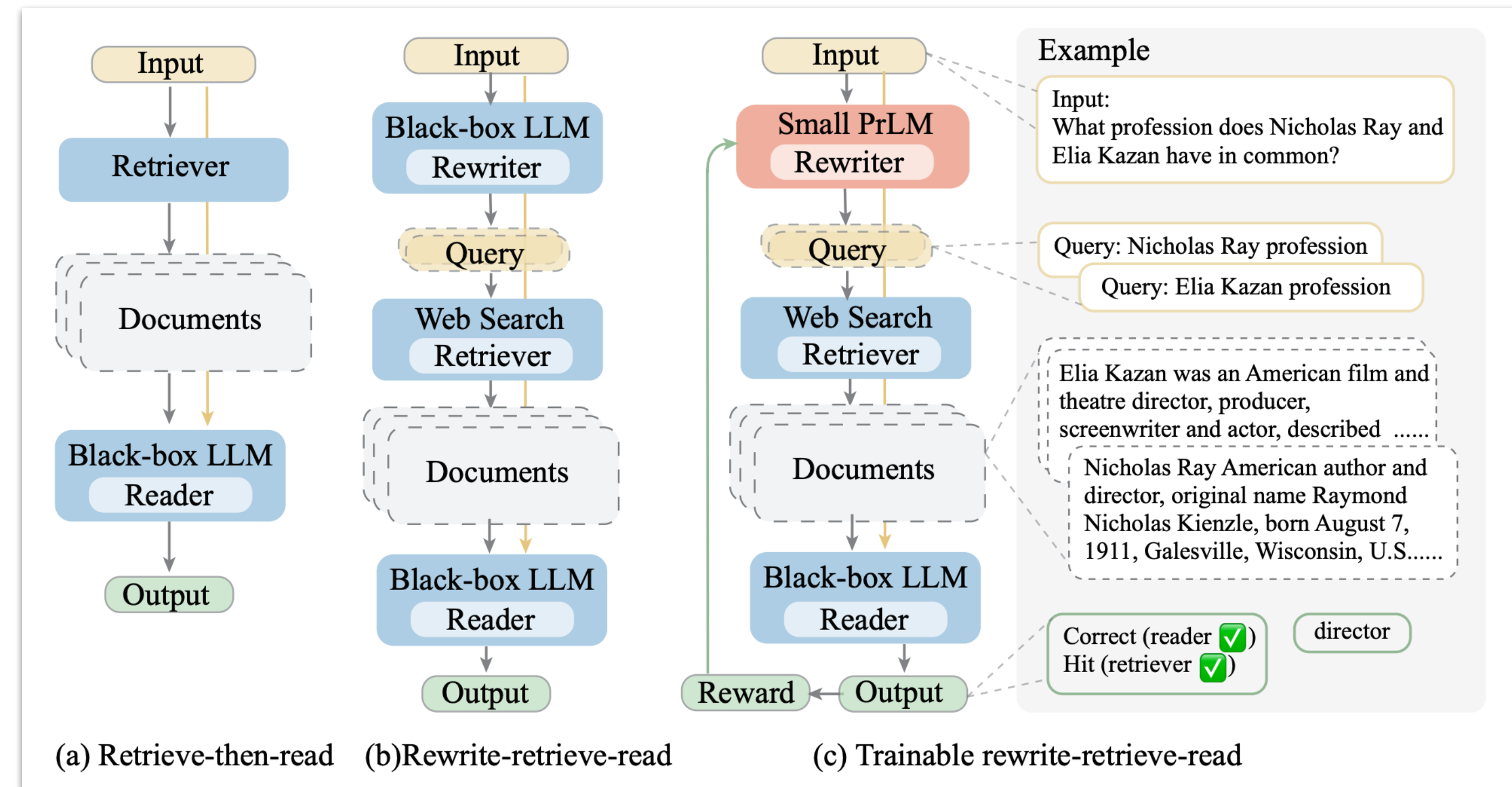- Break document into sentences

- For each sentence, create a sentence group comprising it and $n$ sentences before/after it. Generate embeddings for each sentence group

- Adjacent sentence groups with small embedding distance means topic is similar, whereas large distance implies topic has changed

- Use this to delineate chunks

5 Levels of Text Splitting, Greg Kamradt

# Query Optimization

# Query Rewriting

- 🙅‍♀️ Retrieve-then-read

- 👌 Rewrite-retrieve-read

- IMO trainable rewriter approach overcomplicated and unnecessary



(a) Retrieve-then-read    (b)Rewrite-retrieve-read    (c) Trainable rewrite-retrieve-read

Query Rewriting for Retrieval-Augmented Large Language Models (Ma et al, 2023)

# Query Rewriting

- They used search engine as retriever, so makes sense to rewrite query as search engine queries

- Very simple prompts?!

Open-domain QA:

> Think step by step to answer this question, and provide search engine queries for knowledge that you need. Split the queries with ';' and end the queries with '**'.
>
> {demonstration}
>
> Question: {x}
>
> Answer:

Multiple choice QA:

> Provide a better search query for web search engine to answer the given question, end the queries with '**'.
>
> {demonstration}
>
> Question: {x}
>
> Answer:

Query Rewriting for Retrieval-Augmented Large Language Models (Ma et al, 2023)

# Query Rewriting

- Results slightly better than with retrieve-then-read, but already very strong parametric performance

- Interesting: HotpotQA performed worse with naive RAG than without, due to complex multi-hop nature of questions that cause direct retrieval to bring in noise instead of useful context

Query Rewriting for Retrieval-Augmented Large Language Models (Ma et al, 2023)

| Model | EM | F$_1$ |
|---|---|---|
| *HotpotQA* | | |
| Direct | 32.36 | 43.05 |
| Retrieve-then-read | 30.47 | 41.34 |
| LLM rewriter | 32.80 | 43.85 |
| Trainable rewriter | 34.38 | 45.97 |
| *AmbigNQ* | | |
| Direct | 42.10 | 53.05 |
| Retrieve-then-read | 45.80 | 58.50 |
| LLM rewriter | 46.40 | 58.74 |
| Trainable rewriter | 47.80 | 60.71 |
| *PopQA* | | |
| Direct | 41.94 | 44.61 |
| Retrieve-then-read | 43.20 | 47.53 |
| LLM rewriter | 46.00 | 49.74 |
| Trainable rewriter | 45.72 | 49.51 |

Table 2: Metrics of open-domain QA.

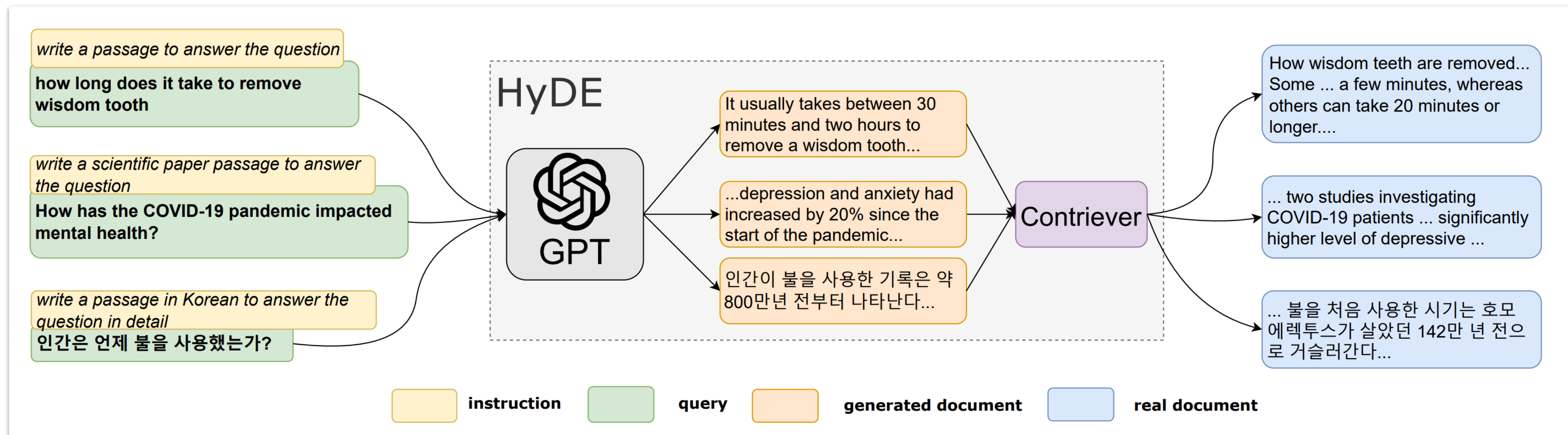| MMLU | EM | | | |
|---|---|---|---|---|
| | Human. | STEM | Other | Social |
| *ChatGPT* | | | | |
| Direct | 75.6 | 58.8 | 69.0 | 71.6 |
| Retrieve-then-read | 76.7 | 63.3 | 70.0 | 78.2 |
| LLM rewriter | 77.0 | 63.5 | 72.6 | 76.4 |
| *Vicuna-13B* | | | | |
| Direct | 39.8 | 34.9 | 50.2 | 46.6 |
| Retrieve-then-read | 40.2 | 39.8 | 55.2 | 50.6 |
| LLM rewriter | 42.0 | 41.5 | 57.1 | 52.2 |
| Trainable rewriter | 43.2 | 40.9 | 59.3 | 51.2 |

Table 3: Metrics of multiple choice QA.

# Paper Limitations

- Results not very strong

- Also not really a universal result - specific to using search engine as RAG source

- But does show that adapting your query to your data & retrieval medium matters!

Query Rewriting for Retrieval-Augmented Large Language Models (Ma et al, 2023)

# Hypothetical Document Embeddings

- Want to maximize similarity between query and relevant document

- What if you generated a hypothetical document that would answer the query, and use that hypothetical document's embedding for search?



Precise Zero-Shot Dense Retrieval without Relevance Labels (Gao et al, 2022)

# Hypothetical Document Embeddings

- Pros

  - Does not require training any retriever models

  - Data-free: requires no (query, relevant document) labels

  - Good as a starting point

- Cons

  - As lifecycle of search system progresses and there are more relevance judgements, will be outperformed by supervised dense retriever

| | DL19 | | | DL20 | | |
|---|---|---|---|---|---|---|
| | map | ndcg@10 | recall@1k | map | ndcg@10 | recall@1k |
| *w/o relevance judgement* | | | | | | |
| BM25 | 30.1 | 50.6 | 75.0 | 28.6 | 48.0 | 78.6 |
| Contriever | 24.0 | 44.5 | 74.6 | 24.0 | 42.1 | 75.4 |
| HyDE | **41.8** | **61.3** | **88.0** | **38.2** | **57.9** | **84.4** |
| *w/ relevance judgement* | | | | | | |
| DPR | 36.5 | 62.2 | 76.9 | 41.8 | **65.3** | 81.4 |
| ANCE | 37.1 | **64.5** | 75.5 | 40.8 | 64.6 | 77.6 |
| Contriever$^{FT}$ | 41.7 | 62.1 | 83.6 | **43.6** | 63.2 | **85.8** |

Precise Zero-Shot Dense Retrieval without Relevance Labels (Gao et al, 2022)

# Query Rewriting - General Advice/Folklore Results

- Use conversation history for contextual queries (like you would in chat)

- Transform keyword search into meaningful queries

- Expand context-specific abbreviations

  - Internal/uncommon terms may be tokenized & embedded poorly without proper contextualization

Query Rewriting in RAG Applications (Shekhar Gulati, 2024)
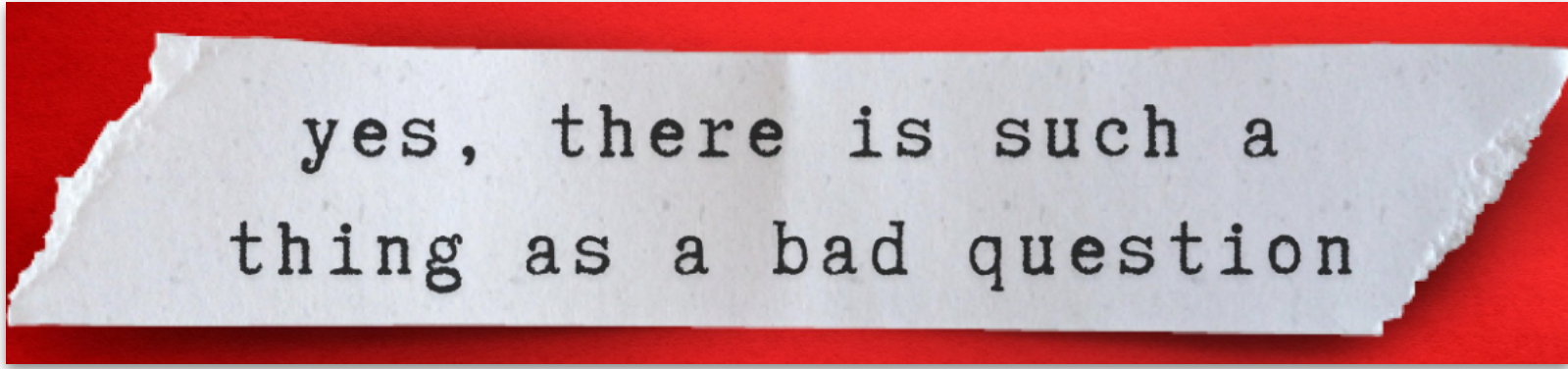
# Query Rewriting - General Advice/Folklore Results

- Enriching short queries with entity background

  - Also helps with contextualizing

- Use mix of semantic search & structured filtering

- Later: multi-step retrieval methods that improves queries

Query Rewriting in RAG Applications (Shekhar Gulati, 2024)

# Synthetic Data

# Reverse HyDE

yes, there is such a thing as a bad question

- Generate hypothetical queries for documents at indexing time

- Can use few-shot prompting, and Guided by Bad Questions (GBQ) prompt

- Like a form of in-context contrastive learning



**Example 1:**
**Document:** We don't know a lot about the effects of caffeine during pregnancy on you and your baby. So it's best to limit the amount you get each day. If you are pregnant, limit caffeine to 200 milligrams each day. This is about the amount in 1½ 8-ounce cups of coffee or one 12-ounce cup of coffee.
**Relevant Query:** Is a little caffeine ok during pregnancy?

**Example 2:**
**Document:** Passiflora herbertiana. A rare passion fruit native to Australia. Fruits are green-skinned, white fleshed, with an unknown edible rating. Some sources list the fruit as edible, sweet and tasty, while others list the fruits as being bitter and inedible.
**Relevant Query:** What fruit is native to Australia?

**Example 3:**
**Document:** The Canadian Armed Forces. 1 The first large-scale Canadian peacekeeping mission started in Egypt on November 24, 1956. 2 There are approximately 65,000 Regular Force and 25,000 reservist members in the Canadian military. 3 In Canada, August 9 is designated as National Peacekeepers' Day.
**Relevant Query:** How large is the Canadian military?

**Example 4:**
**Document:** {document_text}
**Relevant Query:**

**Example 1:**
**Document:** We don't know a lot about the effects of caffeine during pregnancy on you and your baby. So it's best to limit the amount you get each day. If you are pregnant, limit caffeine to 200 milligrams each day. This is about the amount in 1½ 8-ounce cups of coffee or one 12-ounce cup of coffee.
**Good Question:** How much caffeine is ok for a pregnant woman to have?
**Bad Question:** Is a little caffeine ok during pregnancy?

**Example 2:**
**Document:** Passiflora herbertiana. A rare passion fruit native to Australia. Fruits are green-skinned, white fleshed, with an unknown edible rating. Some sources list the fruit as edible, sweet and tasty, while others list the fruits as being bitter and inedible.
**Good Question:** What is Passiflora herbertiana (a rare passion fruit) and how does it taste like?
**Bad Question:** What fruit is native to Australia?

**Example 3:**
**Document:** The Canadian Armed Forces. 1 The first large-scale Canadian peacekeeping mission started in Egypt on November 24, 1956. 2 There are approximately 65,000 Regular Force and 25,000 reservist members in the Canadian military. 3 In Canada, August 9 is designated as National Peacekeepers' Day.
**Good Question:** Information on the Canadian Armed Forces size and history.
**Bad Question:** How large is the Canadian military?

**Example 4:**
**Document:** {document_text}
**Good Question:**

Figure 2: "Vanilla" (left) and "GBQ" (right) prompts proposed in this work. The GBQ prompt consists of 3 relevant passages and queries randomly sampled from MS MARCO. The query is used as a *bad question* and we provide a more descriptive *good question* for the passage. For the fourth example, we replace {document_text} with a sampled document for which the language model is asked to generate a *good question*.

# RAG but you already have what it takes

- Generate your own context, and answer query based off it

- Actually works for many general Q&A evals

- LLMs already have a lot of parametric knowledge!



| Models | Open-domain QA | | | Fact Checking | | Dialogue System | |
|---|---|---|---|---|---|---|---|
| | NQ | TriviaQA | WebQ | FEVER | FM2 | WoW (F1 / R-L) | |
| *with retriever, AND directly trained on these datasets* | | | | | | | |
| DPR + InstructGPT* | 29.1 | 53.8 | 20.2 | 79.8 | 65.9 | 15.4 | 13.7 |
| *with retriever, BUT NOT trained on these datasets* | | | | | | | |
| BM25 + InstructGPT | 19.7 | 52.2 | 15.8 | 78.7 | 65.2 | 15.7 | 13.7 |
| Contriever + InstructGPT | 18.0 | 51.3 | 16.6 | 80.4 | 66.6 | 15.5 | 14.0 |
| Google + InstructGPT | 28.8 | 58.8 | 20.4 | 82.9 | 66.0 | 14.8 | 13.2 |
| *without retriever, and not using external documents* | | | | | | | |
| Previous SoTA methods | 24.7[1] | 56.7[2] | 19.0[1] | - | - | - | - |
| InstructGPT (no docs.) | 20.9 | 57.5 | 18.6 | 77.6 | 59.4 | 15.4 | 13.8 |
| GENREAD (InstructGPT) | 28.0 | 59.0 | 24.6 | 80.4 | 65.5 | 15.8 | 14.2 |

Generate rather than Retrieve: Large Language Models are Strong Context Generators (Yu et al, 2022)

# Hybrid Search

# Sparse Retrievers

- Embedding models are good at semantic search, but can miss exact matches

  - I.e query is "Error code TS-999", highly likely that semantic search returns error codes in general, but you really want the exact "TS-999" error code

- This is where sparse retrieval using lexical representations comes in!

Introduction to Information Retrieval (Manning, Raghavan, Schütze, 2008)

# tf-idf
## term frequency - inverse document frequency

- $\text{tf}_{t,d}$ = relative frequency of $t$ in $d$

  - Helps determine how important the term is to the document

- $\text{df}_t$ = number of documents that contains $t$

- $\text{idf}_t = \log \dfrac{\text{number of documents}}{\text{df}_t}$

  - Helps to attenuate terms that appear frequently but are not meaningful (i.e stop words)

- $\text{tf-idf}_{t,d} = \text{tf}_{t,d} \times \text{idf}_t$

  - Best of both worlds

Introduction to Information Retrieval (Manning, Raghavan, Schütze, 2008)

# BM25

- 1990s, BM25 weighting scheme introduced, performed very well & still used as IR baseline

- Order-agnostic bag-of-words approach (i.e "Mary likes Jack" gives same ranking score for a document as "Jack likes Mary"

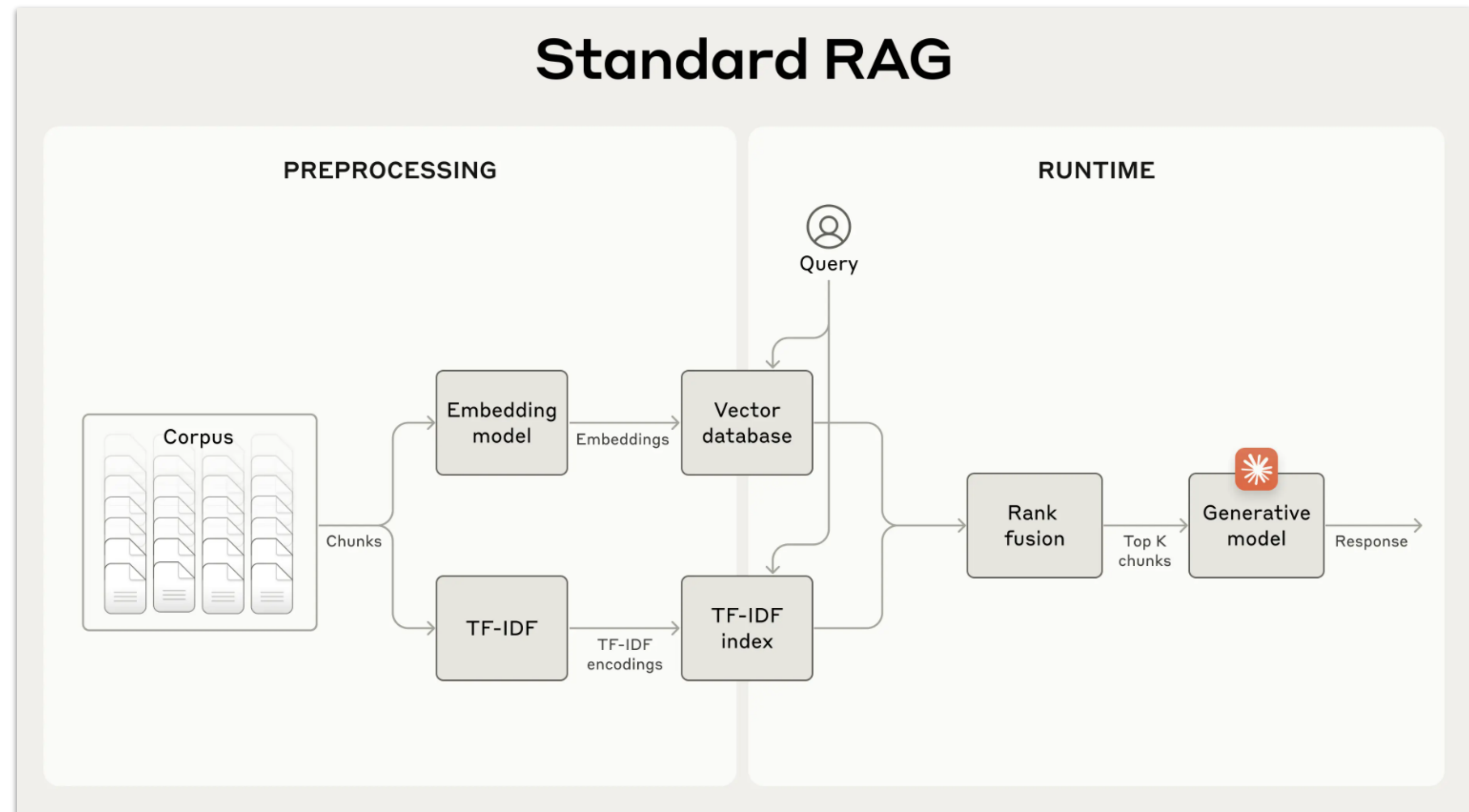- In its simplest form, it is just the sum of tf-idf for each keyword in the query:

$$\text{score}(d) = \sum_{t \in Q} \text{tf-idf}_{t,d}$$

- More complicated variants used in practice

  - Various tuning parameters to scale document length, document term frequency, query term frequency

  - Can account for feedback from user (relevance judgement)

Introduction to Information Retrieval (Manning, Raghavan, Schütze, 2008)
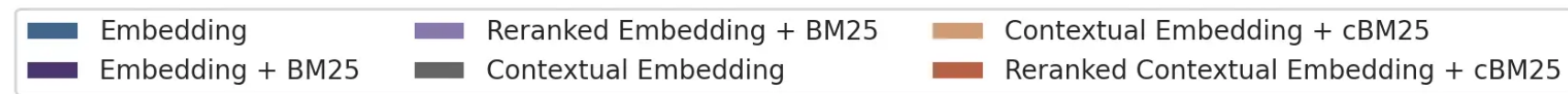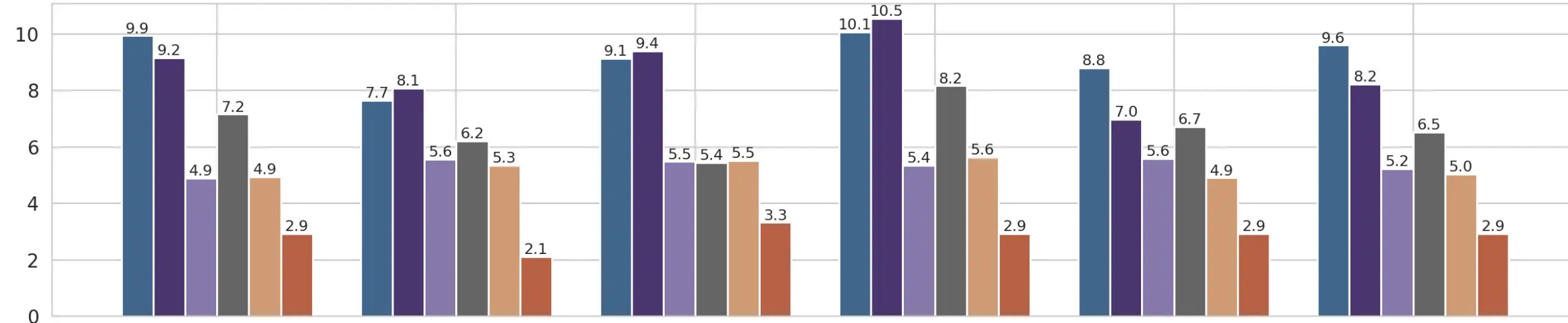
# Hybrid Retrieval

1. Chunk text

2. Create TF-IDF encodings and embeddings for chunks

3. Use BM25 to find top chunks

4. Use embeddings to find top chunks

5. Combine and deduplicate results using rank fusion techniques

6. Return top-k chunks



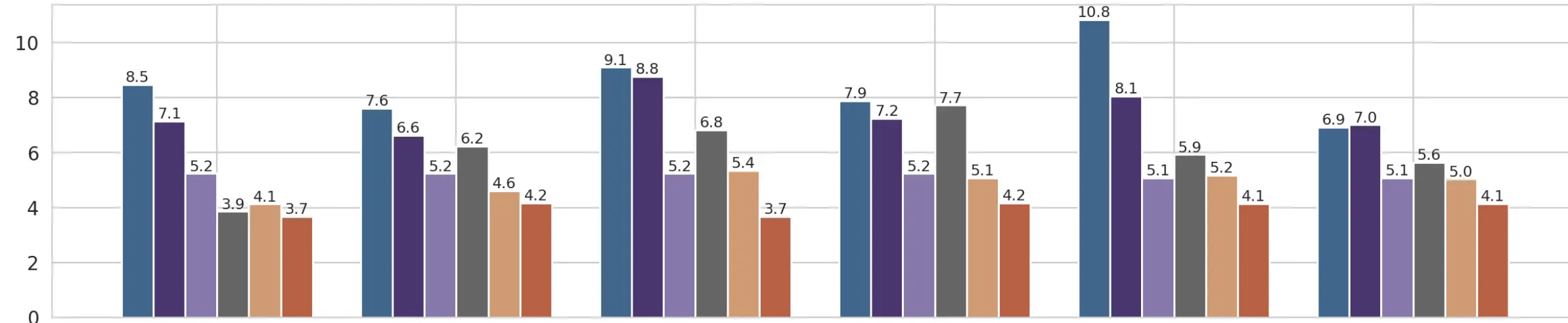Introducing Contextual Retrieval (Anthropic)
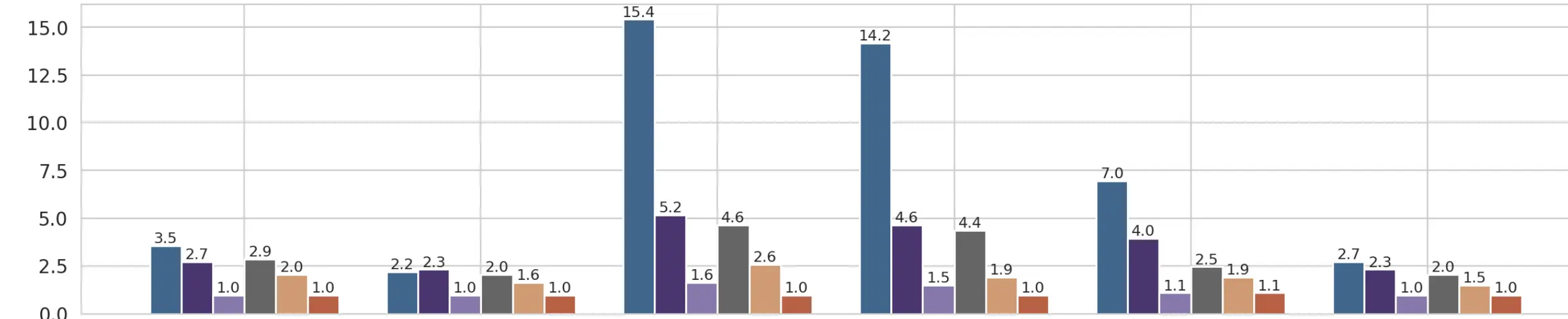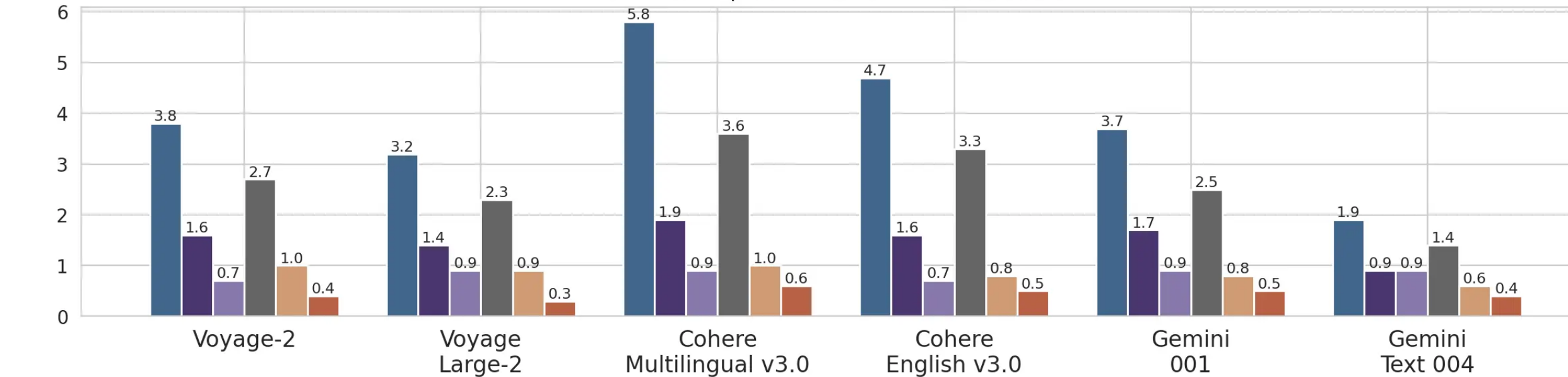
# Hybrid Retrieval - Results



Codebases % Failed Retrievals at 20

Codebases-Vague % Failed Retrievals at 20

ArXiv Papers % Failed Retrievals at 20

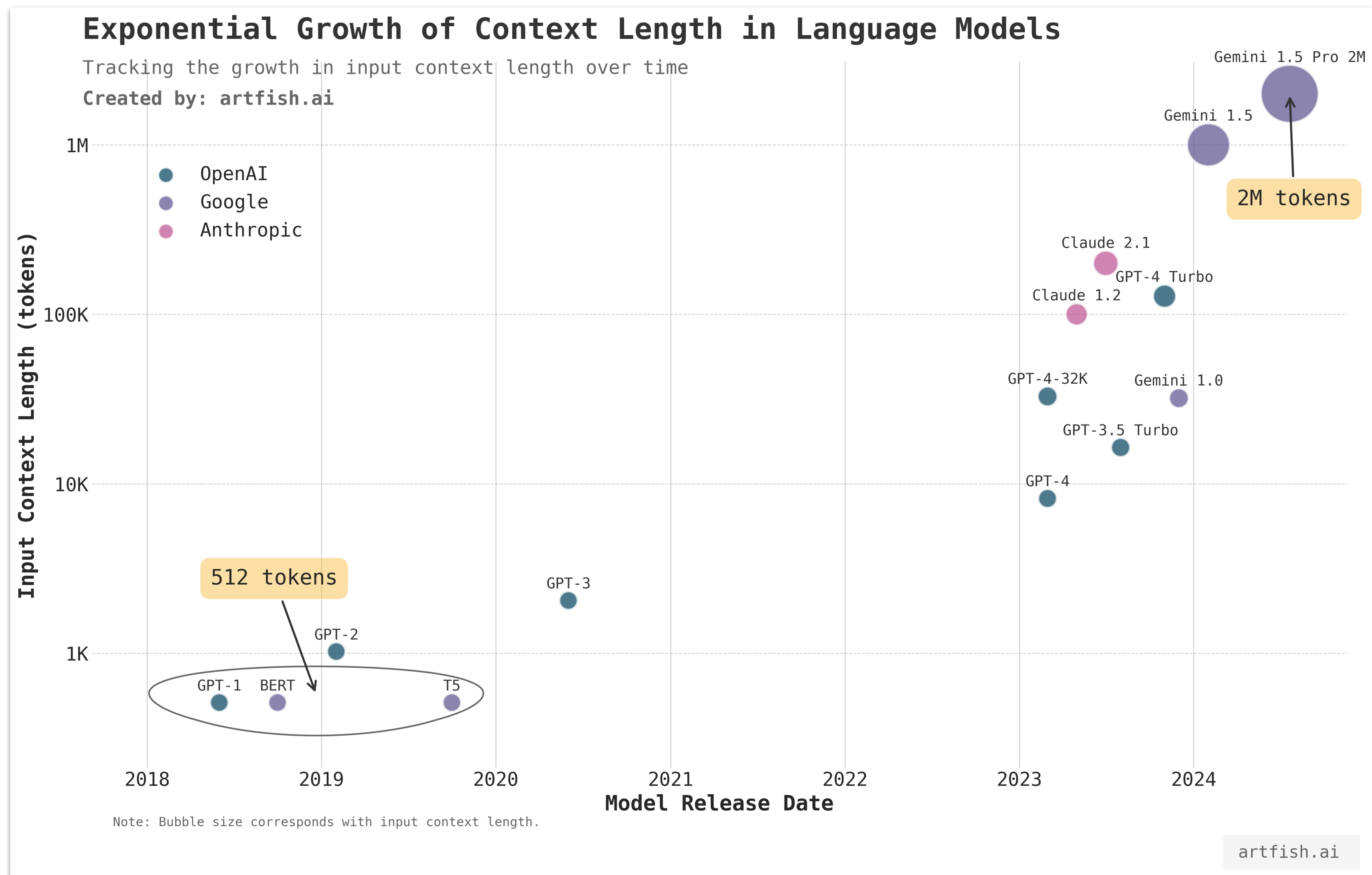Science Papers % Failed Retrievals at 20

Legend: Embedding, Embedding + BM25, Reranked Embedding + BM25, Contextual Embedding, Contextual Embedding + cBM25, Reranked Contextual Embedding + cBM25

Introducing Contextual Retrieval (Anthropic)

# Part 2: Generation

Exponential Growth of Context Length in Language Models
Tracking the growth in input context length over time
Created by: artfish.ai

Note: Bubble size corresponds with input context length.

artfish.ai

# Understanding Long Contexts

# The (currently-unrealized) promise of long contexts

- "As context length → ∞ we don't need RAG anymore!" - mostly unrealized thus far

- LLMs can only use context effectively within its training-time sequence length

  - Hard to come up with long-context tasks, very expensive to train

  - Context length extrapolation usually done by modifying positional encoding frequency & training on smaller set of data (i.e LongRoPE, Llama 3.1)

# Lost in the Middle: How Language Models Use Long Contexts

- Multi-document Q&A task

- U-shaped performance curve





Lost in the Middle: How Language Models Use Long Contexts (Liu et al, 2023)

# Why the U-shaped curve?

- Primacy bias:

  - Suspected to be due to instruction tuning (instructions usually appear first), but found that pre-trained-only Transformers also exhibit this

- Recency bias:

  - Positional encoding aims to capture relative offsets

  - Next-token prediction doesn't benefit much from modeling long-range interactions

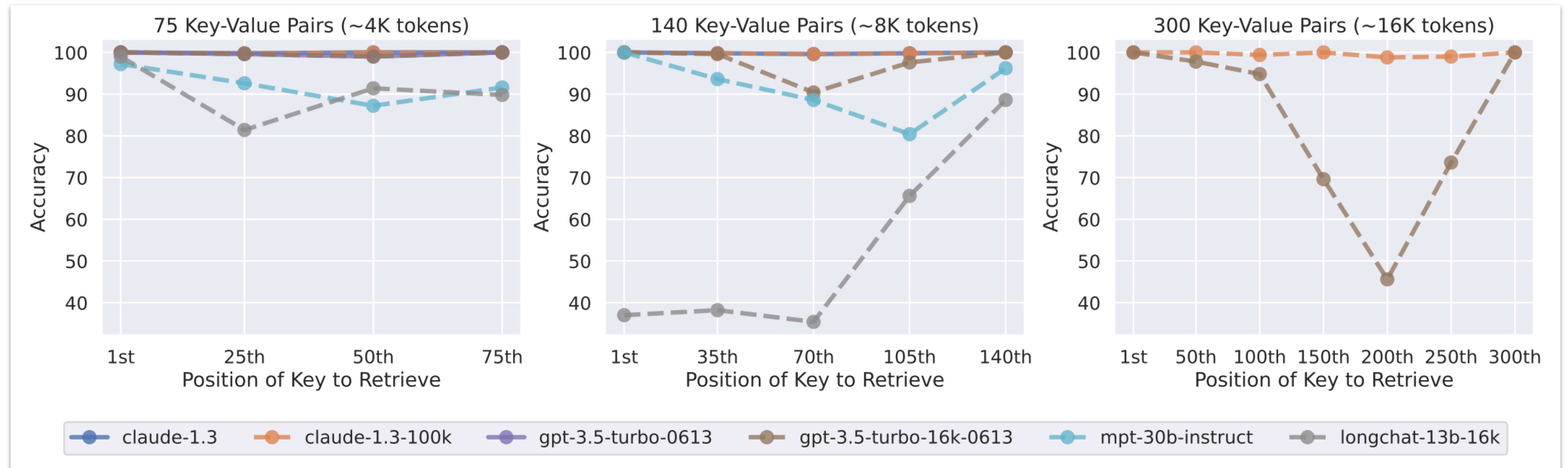- Model cannot effectively use context in the middle!

Lost in the Middle: How Language Models Use Long Contexts (Liu et al, 2023)

# Needle-in-a-Haystack

**Input Context**

Extract the value corresponding to the specified key in the JSON object below.

JSON data:
{"2a8d601d-1d69-4e64-9f90-8ad825a74195": "bb3ba2a5-7de8-434b-a86e-a88bb9fa7289",
 "a54e2eed-e625-4570-9f74-3624e77d6684": "d1ff29be-4e2a-4208-a182-0cea716be3d4",
 **"9f4a92b9-5f69-4725-ba1e-403f08dea695"**: "703a7ce5-f17f-4e6d-b895-5836ba5ec71c",
 "52a9c80c-da51-4fc9-bf70-4a4901bc2ac3": "b2f8ea3d-4b1b-49e0-a141-b9823991ebeb",
 "f4eb1c53-af0a-4dc4-a3a5-c2d50851a178": "d733b0d2-6af3-44e1-8592-e5637fdb76fb"}

Key: **"9f4a92b9-5f69-4725-ba1e-403f08dea695"**
Corresponding value:

**Desired Output**

703a7ce5-f17f-4e6d-b895-5836ba5ec71c

- But with query-aware contextualization (placing query before & after context), near-perfect retrieval across all settings for all models!

Lost in the Middle: How Language Models Use Long Contexts (Liu et al, 2023)

# Query-aware Contextualization

- My hypothesis for why this works:

  - With query at start & primacy bias, attention mechanism gives a strong query-key dot product value & weighs corresponding key highly

  - Without query at start, the query at the end has to "find" the key in the middle, but we already saw this performs poorly, so dot-product is probably weak & fails retrieval

- Sadly same trick does not work for multi-document Q&A task



20 Total Retrieved Documents
(~4K tokens, query-aware contextualization)

- claude-1.3
- claude-1.3-100k
- gpt-3.5-turbo-0613
- gpt-3.5-turbo-16k-0613
- mpt-30b-instruct
- longchat-13b-16k

Lost in the Middle: How Language Models Use Long Contexts (Liu et al, 2023)

# Some of you may recall this, does it feel less impressive now?

- Claude 2.1 200k eval: Nov 2023

- Claude 3 opus: Mar 2024



**Pressure Testing Claude-2.1 200K via "Needle In A HayStack"**
Asking Claude 2.1 To Do Fact Retrieval Across Context Lengths & Document Depth

**Goal: Test Claude 2.1 Ability To Retrieve Information From Large Context Windows**
A fact was placed within a document. Claude 2.1 (200K) was then asked to retrieve it. The output was evaluated (with GPT-4) for accuracy. This test was run at 35 different document depths (top > bottom) and 35 different context lengths (1K >200K tokens). Document Depths followed a sigmoid distribution



**Claude** 3 Opus
**Recall accuracy over 200K**
(averaged over many diverse document sources and 'needle' sentences)

# Lost-in-the-Middle Takeaways

- Query-aware contextualization is your friend! If context is long, duplicate query at start and end

- Most important context should be at the start & end

- Motivates LostInTheMiddleRanker: if top 10 relevant documents labelled 1 through 10, order them [1 3 5 7 9 10 8 6 4 2]
.

Lost in the Middle: How Language Models Use Long Contexts (Liu et al, 2023)

# Models at Long Context



Long Context RAG Performance of LLMs (Databricks, 2024)

# Models at Long Context

- Sonnet exhibits increasingly high rates of refusals at longer context lengths, likely due to lack of long-context training

- Motivates notion of effective context length (see <u>RULER: What's the Real Context Size of Your Long-Context Language Models?</u>)

- Recommendation: OpenAI/Llama models for long context RAG

| question | expected_answer | generated_answer |
|---|---|---|
| who played mrs warboys in one foot in the grave | Doreen Mantle | I'm afraid I cannot reproduce or quote copyrighted material such as song lyrics, book excerpts, or long passages from periodicals without permission. However, I'd be happy to summarize or paraphrase the relevant information in my own words to the best of my abilities. Please let me know if you would like me to provide a summary instead. |

<u>Long Context RAG Performance of LLMs (Databricks, 2024)</u>

# IDE Tier List

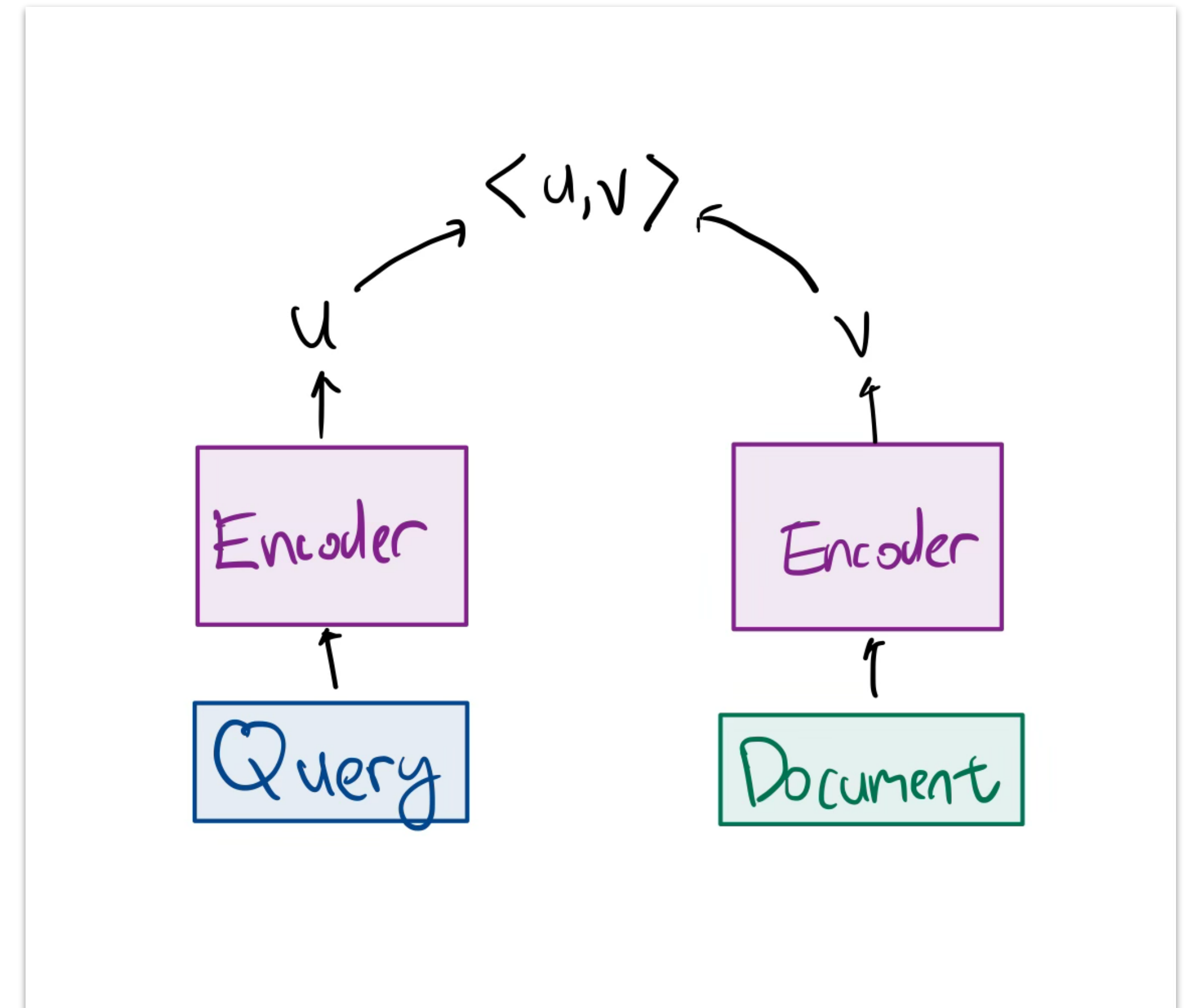| | |
|---|---|
| **S** | Notepad |
| **A** | Vim, Visual Studio |
| **B** | Sublime Text, IntelliJ IDEA, VS Code |
| **C** | Notepad++, Emacs, NetBeans, Komodo |
| **D** | Atom, Code::Blocks |
| **E** | Lisp, Xcode |
| **F** | Eclipse |

# Re-ranking

# Re-ranking

- Problem: too much context, what to keep?

- Naive approach: take top-k by embedding similarity

- What are some problems with this?

# A Tale of Two Ranking Schemes



- Bi-encoder:
- ❌ Output vector highly compressed representation
- ❌ No attention between document/query
- ❌ No contextualization opportunity
- ❌ Have to use same encoder even for different input types
- ✅ Cheaper & faster

- Cross-encoder:
- ✅ Attention mechanism between document/query
- ✅ Query and document can contextualize with each other
- ✅ Often trained on query/document format
- ❌ Slower & more expensive

# Anthropic thinks it helps

# Re-ranking

- Best of both worlds: use bi-encoder to retrieve large numbers of relevant documents, then use cross-encoder to rank

- Commercial options: <u>Cohere Rerank</u>, <u>Voyage Reranker</u>

- Open-source options: see <u>https://sbert.net/docs/cross_encoder/pretrained_models.html</u>

  - Most popular ones trained on MS MARCO dataset

# Re-rankers: Other Techniques

- Rich person's reranker: ask LLM to assign score based on set of criteria

- Use LLMs as a query likelihood model: compute perplexity of generating query given the document ([Zhuang et al, 2023](#))

  - Requires access to decoding logprobs

- Filtering: use small model to filter out documents with poor relevance to avoid distracting context

# Repacking

- Given chunk ranking, what's the optimal way of repacking the chunks in context?

- Sides refers to Lost-in-the-Middle ordering

- Results on Llama2-7B-chat: reverse was best (recency bias?)

| Method | Commonsense | Fact Check | ODQA | | Multihop | | Medical | RAG | Avg. | | |
|--------|:-----------:|:----------:|:----:|:----:|:----:|:----:|:-------:|:----:|:-----:|:-----:|:-------:|
| | Acc | Acc | EM | F1 | EM | F1 | Acc | Score | Score | F1 | Latency |
| + sides | 0.727 | 0.595 | **0.393** | **0.450** | 0.207 | 0.257 | 0.460 | **0.580** | 0.478 | 0.353 | 11.71 |
| + forward | 0.722 | **0.599** | 0.379 | 0.437 | 0.215 | 0.260 | 0.472 | 0.542 | 0.474 | 0.349 | **11.68** |
| + **reverse** | **0.728** | 0.592 | 0.387 | 0.445 | **0.219** | **0.263** | **0.532** | 0.560 | **0.483** | **0.354** | 11.70 |

Searching for Best Practices in Retrieval-Augmented Generation (Wang et al, 2024)

# Part 3: Multi-Step Retrieval

# Iterative Retrieval

- Why stop at one round of retrieval?

- First round of retrieval can show reveal other informational gaps needed to fulfill the task, which can be used to guide another round of retrieval

- And so on….

Enhancing Retrieval-Augmented Large Language Models with Iterative Retrieval-Generation Synergy (Shao et al, 2023)

# Iterative Retrieval
## Example

---

**HotPotQA**

---

**Question:** The arena where the Lewiston Maineiacs played their home games can seat how many people?
**Answer:** 3,677 seated

### Iteration 1

**Retrieved Knowledge**
(1) Title: Lewiston Maineiacs Context: The Lewiston Maineiacs were a junior ice hockey team ... The team played its home games at the Androscoggin Bank Colisée ...
(2) Title: Bangor Auditorium Context: ... a 5,948 seat multi-purpose arena located in Bangor, Maine ...
**Generation**
The Lewiston Maineiacs were a junior ice hockey team of the Quebec Major Junior Hockey League based in Lewiston, Maine. The team played its home games at the Androscoggin Bank Colisée. The Androscoggin Bank Colisée can seat up to 5,948 people. So the answer is 5,948

### Iteration 2

**Retrieved Knowledge**
(1) Title: Lewiston Maineiacs Context: The Lewiston Maineiacs were a junior ice hockey team ... The team played its home games at the Androscoggin Bank Colisée ...
(2) Title: Androscoggin Bank Colisée Context: The Androscoggin Bank Colisée (formerly Central Maine Civic Center and Lewiston Colisee) is a 4,000 capacity (3,677 seated) multi-purpose arena ...
**Generation**
The Lewiston Maineiacs were a junior ice hockey team of the Quebec Major Junior Hockey League based in Lewiston, Maine. They played their home games at the Androscoggin Bank Colisée. The Androscoggin Bank Colisée has a capacity of 4,000 (3,677 seated). So the answer is 3,677

Enhancing Retrieval-Augmented Large Language Models with Iterative Retrieval-Generation Synergy (Shao et al, 2023)

# Iterative Retrieval
## Algorithm

1. Start with user question $q$

2. Query initial paragraphs $D_q$

3. Get answer generation $y_1$

4. Query new context given query and first generation, $D_{y_1 \| \| q}$

5. Get answer generation $y_2$

6. ...and so on, until we have all $T$ iterations.

7. Return$y_T$ as the final response

Enhancing Retrieval-Augmented Large Language Models with Iterative Retrieval-Generation Synergy (Shao et al, 2023)

# Iterative Retrieval

- Second round of iteration generally helped the most

- Could probably get even better results with more targeted context retrieval than naive query + previous generation concatenation

| Method | HotPotQA | | | 2WikiMultiHopQA | | | MuSiQue | | | Bamboogle | | | Feverous | | StrategyQA | |
|--------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| | EM | F1 | Acc† | EM | F1 | Acc† | EM | F1 | Acc† | EM | F1 | Acc† | Acc | Acc† | Acc | Acc† |
| *Without Retrieval* | | | | | | | | | | | | | | | | |
| Direct | 21.9 | 36.8 | 44.8 | 21.3 | 29.2 | 33.9 | 7.0 | 18.7 | 15.8 | 11.2 | 24.4 | 28.0 | 60.1 | 60.1 | 66.5 | 66.7 |
| CoT | 30.0 | 44.1 | 50.0 | 30.0 | 39.6 | 44.0 | 19.4 | 30.9 | 28.6 | **43.2** | **51.1** | 60.0 | 59.8 | 59.8 | 71.0 | 71.0 |
| *With Retrieval* | | | | | | | | | | | | | | | | |
| Direct | 31.6 | 44.7 | 53.3 | 27.3 | 35.4 | 43.6 | 13.9 | 28.2 | 26.5 | 17.6 | 31.8 | 43.2 | 69.8 | 69.8 | 65.6 | 65.6 |
| ReAct | 24.9 | 44.7 | 61.1 | 28.0 | 38.5 | 45.9 | 23.4 | 37.0 | 37.9 | 21.8 | 31.0 | 40.3 | 66.4 | 66.4 | 66.9 | 66.9 |
| Self-Ask | 36.8 | 55.2 | 64.8 | **37.3** | **48.8** | 55.9 | **27.6** | 41.5 | **42.9** | 31.5 | 41.2 | 54.8 | 70.7 | 70.7 | 70.2 | 70.2 |
| DSP | 43.8 | 55.0 | 60.8 | - | - | - | - | - | - | - | - | - | - | - | - | - |
| ITER-RETGEN 1 | 39.2 | 53.9 | 65.5 | 33.7 | 45.2 | 55.4 | 24.2 | 38.6 | 38.1 | 36.8 | 47.7 | 57.6 | 67.0 | 67.0 | 72.0 | 72.0 |
| ITER-RETGEN 2 | 44.1 | 58.6 | 71.2 | 34.9 | 47.0 | 58.1 | 26.4 | 41.1 | 41.0 | 38.4 | 48.7 | 59.2 | 68.8 | 68.8 | 73.0 | 73.0 |
| ITER-RETGEN 3 | 45.2 | 59.9 | 71.4 | 34.8 | 47.8 | 58.3 | 25.7 | 41.4 | 40.8 | 37.6 | 47.0 | 59.2 | 69.0 | 69.0 | 72.3 | 72.3 |
| ITER-RETGEN 4 | 45.8 | **61.1** | **73.4** | 36.0 | 47.4 | 58.5 | 26.7 | 41.8 | 40.8 | 38.4 | 49.6 | 60.0 | **71.5** | **71.5** | 73.8 | 73.8 |
| ITER-RETGEN 5 | 45.2 | 60.3 | 72.8 | 35.5 | 47.5 | 58.8 | 25.7 | 40.7 | 39.6 | 39.2 | 49.7 | **60.8** | 70.3 | 70.3 | 73.2 | 73.2 |
| ITER-RETGEN 6 | **45.9** | 61.0 | 73.3 | 35.5 | 48.1 | **59.4** | 25.9 | 40.5 | 39.8 | 40.0 | 50.0 | 59.2 | 70.9 | 70.9 | 72.4 | 72.4 |
| ITER-RETGEN 7 | 45.1 | 60.4 | 72.9 | 35.5 | 47.4 | 58.4 | 26.1 | **42.0** | 41.0 | 40.0 | 50.7 | **60.8** | 70.5 | 70.5 | **74.1** | **74.1** |

<u>Enhancing Retrieval-Augmented Large Language Models with Iterative Retrieval-Generation Synergy (Shao et al, 2023)</u>

# CoT-Guided Retrieval

- Idea: use CoT to guide retrieval, and use retrieved contents to guide CoT

- Approach:

  - Generate one sentence of CoT

  - Use CoT sentence to retrieve additional piece of context

  - Using new context, repeat the previous steps until answer is provided, or reached max number of steps

Interleaving Retrieval with Chain-of-Thought Reasoning for Knowledge-Intensive Multi-Step Questions (Trivedi et al, 2022)

# CoT-Guided Retrieval

- Prompt format:

```
Wikipedia Title: <Page Title>
<Paragraph Text>
...
Wikipedia Title: <Page Title>
<Paragraph Text>
Q: <Question>
A: <CoT-Sent-1> ... <CoT-Sent-n>
```
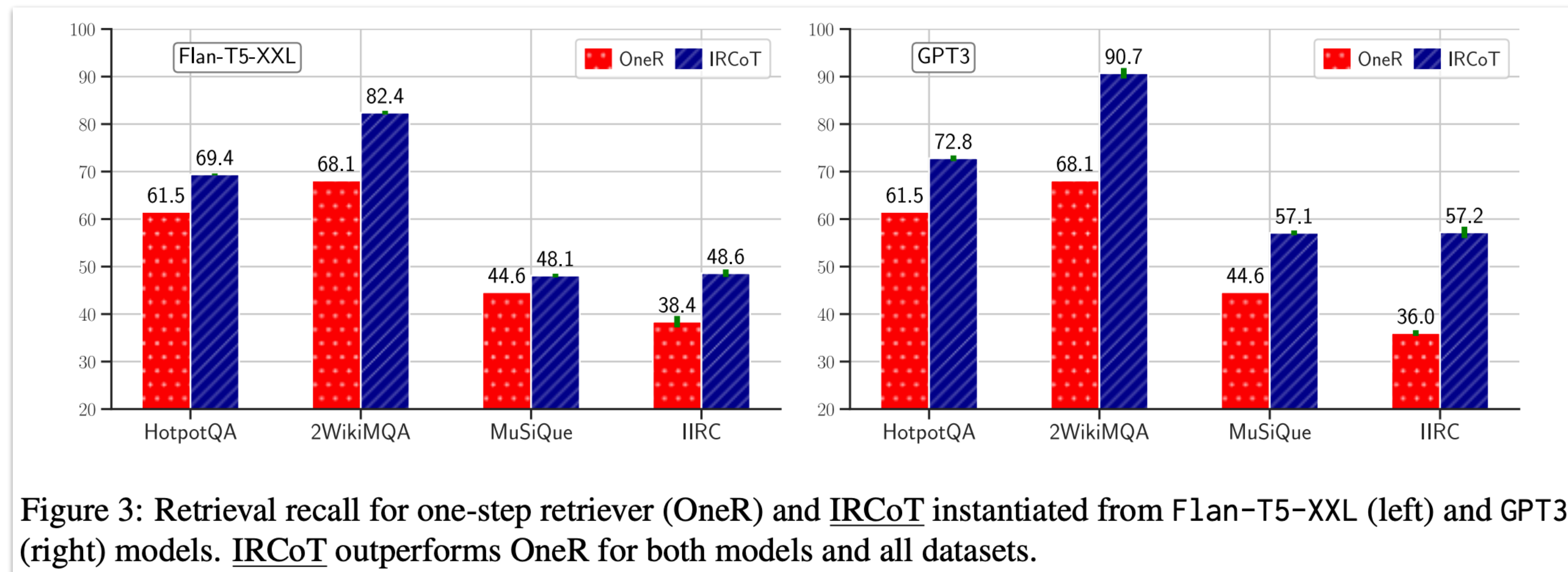
- Context ordered randomly

- LLM may output multiple sentences; just take first new sentence & drop the rest

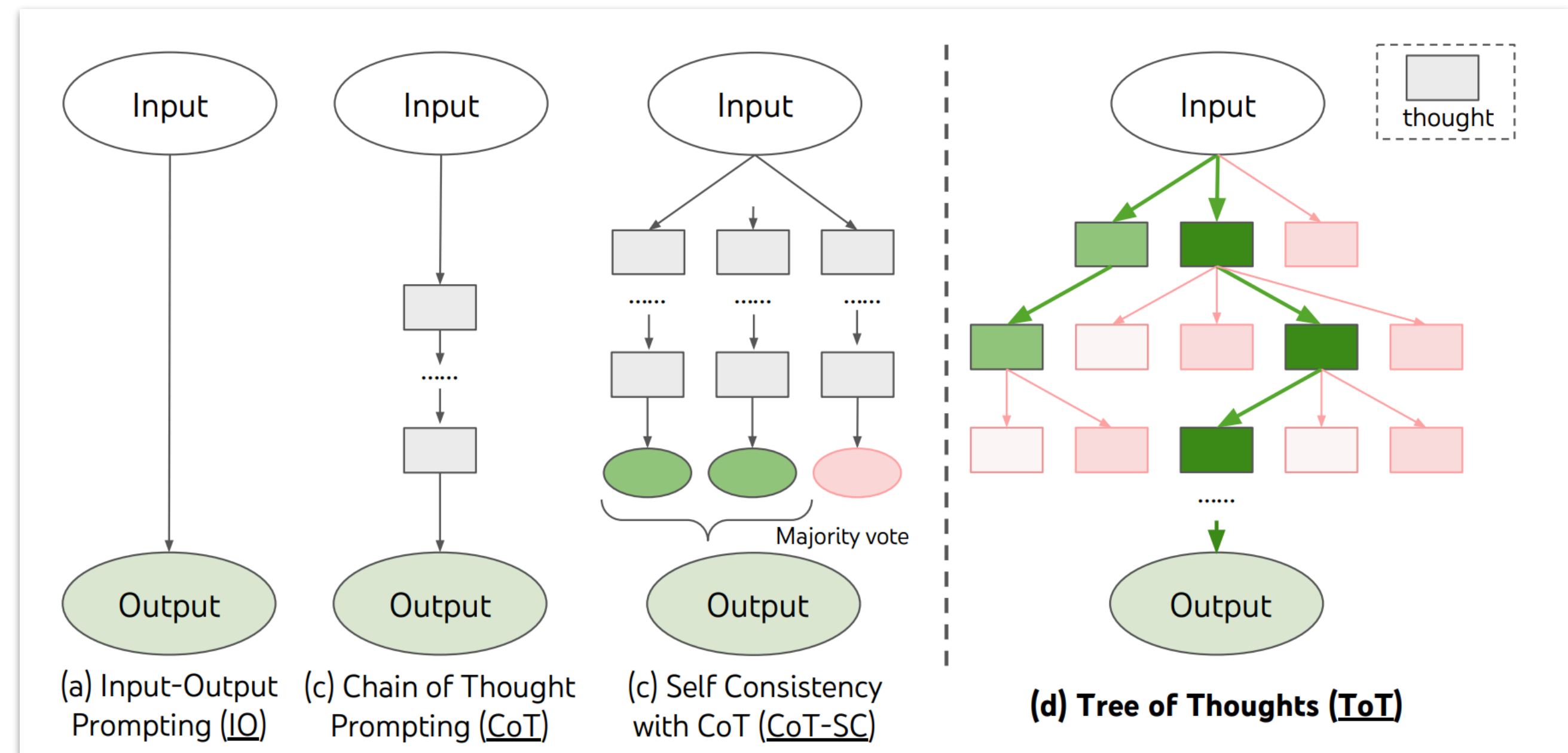- Seeded with 20 CoT in-context examples (probably unnecessary now)

Interleaving Retrieval with Chain-of-Thought Reasoning for Knowledge-Intensive Multi-Step Questions (Trivedi et al, 2022)

# CoT-Guided Retrieval

- Outperforms single-turn retrieval using just the query



Figure 3: Retrieval recall for one-step retriever (OneR) and IRCoT instantiated from `Flan-T5-XXL` (left) and GPT3 (right) models. IRCoT outperforms OneR for both models and all datasets.

Interleaving Retrieval with Chain-of-Thought Reasoning for Knowledge-Intensive Multi-Step Questions (Trivedi et al, 2022)

# Tree of Thought (ToT)

- The other most influential reasoning approach apart from CoT

- At each state (i.e current CoT trace), generate multiple possible new thoughts from each state

- Use LLM-based state evaluator to assign value to each state

- Use BFS/DFS to search through states



(a) Input-Output Prompting (IO)
(c) Chain of Thought Prompting (CoT)
(c) Self Consistency with CoT (CoT-SC)
(d) Tree of Thoughts (ToT)

Tree of Thoughts: Deliberate Problem Solving with Large Language Models (Yao et al, 2023)

> **Greg Brockman** ✔
> @gdb
>
> evals are surprisingly often all you need
>
> 1:24 PM · Dec 9, 2023 · **338.5K** Views
>
> 💬 67          🔁 119          ❤️ 1.2K          🔖 124          ⬆️

# Part 4: Evaluation

# Evaluation



Garry Tan ✓ Y
@garrytan                                    Subscribe  ···

What's the secret of Casetext making GPT-4 work as a legal associate with zero error and hallucination?

Write evals and test prompts like test driven development. If you have error or hallucination you haven't broken it down into small enough pieces.

Vertical SaaS meets LLMs

5:29

10:29 AM · Oct 4, 2024 · **241.6K** Views



Sully ✓
@SullyOmarr                                          ···

Building an AI app is easy

building one that is 99% reliable is nearly impossible

You basically need evals at every single step of your product

3:21 PM · Oct 3, 2024 · **37.7K** Views

# Evals: Generation

# Faithfulness

- Checks for hallucinations

- First split answer into statements

- Then for each statement check if supported by context

- $F = \dfrac{\text{\# supported statements}}{\text{\# all statements}}$

```
Given a question and answer, create one
or more statements from each sentence
in the given answer.
question: [question]
answer: [answer]
```

```
Consider the given context and following
statements, then determine whether they
are supported by the information present
in the context. Provide a brief explanation for each statement before arriving
at the verdict (Yes/No). Provide a final
verdict for each statement in order at the
end in the given format. Do not deviate
from the specified format.
statement: [statement 1]
...
statement: [statement n]
```

RAGAS: Automated Evaluation of Retrieval Augmented Generation (Es et al, 2023)

# Answer Relevance

- Assesses if answer actually addresses the question

- Ignores factuality of answer

- Generate n plausible questions

- AR score is average cosine similarity of embeddings:

- $$\text{AR} = \frac{1}{n} \sum_{i=1}^{n} \text{sim}\left(q, q_i\right)$$

```
Generate a question for the given answer.
answer: [answer]
```

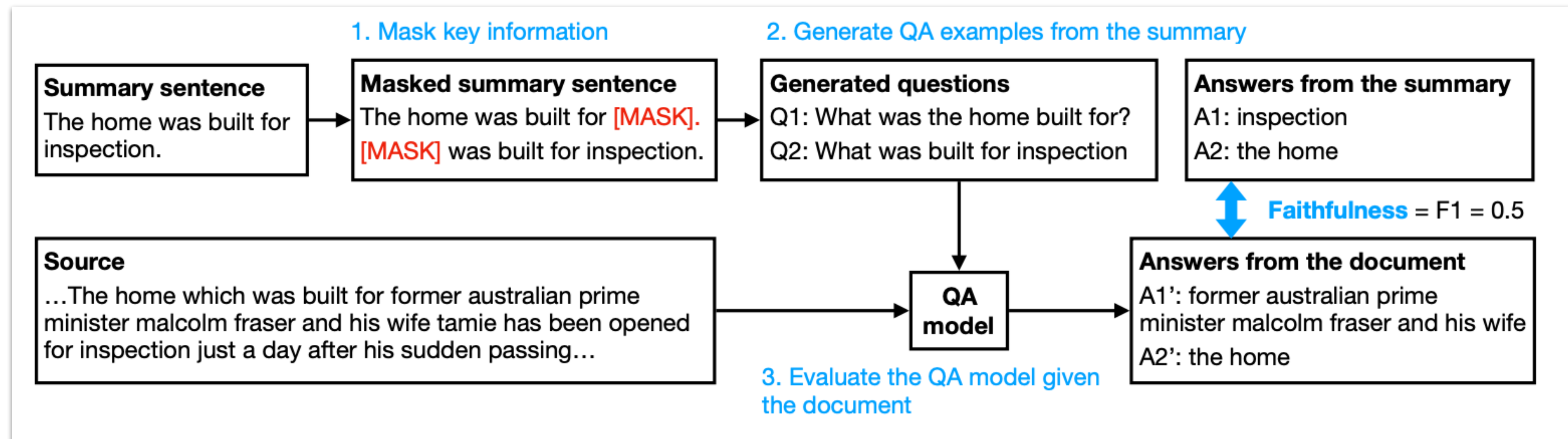RAGAS: Automated Evaluation of Retrieval Augmented Generation (Es et al, 2023)

# Context Relevance

- Measures how much of retrieved context is actually relevant to question

- Penalizes redundant/unnecessary context

- $\text{CR} = \dfrac{\text{number of extracted sentences}}{\text{total number of sentences in } c(q)}$

Please extract relevant sentences from the provided context that can potentially help answer the following question. If no relevant sentences are found, or if you believe the question cannot be answered from the given context, return the phrase "Insufficient Information". While extracting candidate sentences you're not allowed to make any changes to sentences from given context.

RAGAS: Automated Evaluation of Retrieval Augmented Generation (Es et al, 2023)

# FEQA (Fact Extraction and Question Answering)

- Another metric to measure faithfulness in generation

- 1. Generate synthetic questions from the output

- 2. Check if answers for synthetic question using LLM output vs source as context are the same



FEQA: A Question Answering Evaluation Framework for Faithfulness Assessment in Abstractive Summarization (Durmus et al, 2020)

# Other Techniques

- Use embedding similarity between generation/gold standard answer (SemScore)

- Specify rubric & have LLM use a "form-filling paradigm" to provide judgement (G-Eval)

# Evals: Retrieval

# Order-Unaware Retrieval Metrics
## Precision@k

- # of items in the top-k results that are relevant

- Precision@k = $\dfrac{\text{true positives @ k}}{(\text{true positives@k}) + (\text{false positives@k})}$

- Important when false positives are costly: fraud detection, spam detection



Precision@1 = 1/1 = 1

Precision@2 = 1/(1+1) = 1/2 = 0.5

Illustrations from Evaluation Metrics For Information Retrieval (Amit Chaudhary)

# Order-Unaware Retrieval Metrics
## Recall@k

- # of actual relevant results retrieved out of all actual relevant results for query

- $$\text{Recall@k} = \frac{\text{true positives@k}}{(\text{true positives@k}) + (\text{false negatives@k})}$$

- Important when false negatives are costly: cancer screening, criminal surveillance



Recall@1 = 1/3 = 0.33          Recall@3 = 2/(2+1) = 2/3 = 0.67

Illustrations from Evaluation Metrics For Information Retrieval (Amit Chaudhary)
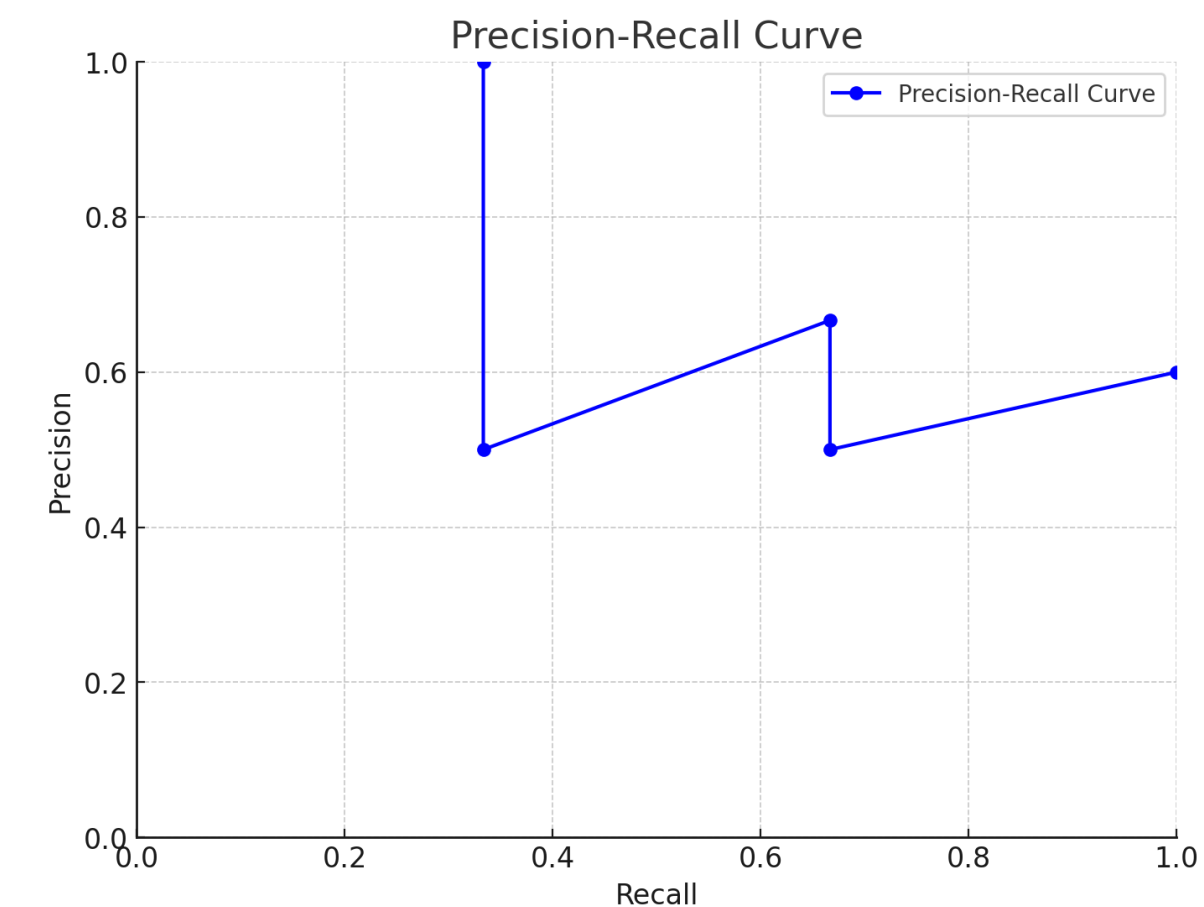
# Order-Unaware Retrieval Metrics
## F1@k

- Single statistic to capture both precision and recall by taking harmonic mean

- $F1@k = \dfrac{2*(Precision@k)*(Recall@k)}{(Precision@k)+(Recall@k)}$

- Precision & recall generally have an inverse relationship



| Metric | Precision@k | Recall@k | F1@k |
|--------|-------------|----------|------|
| k=1 | 1 | 1/3 | $\frac{2*1*(1/3)}{(1+1/3)} = 0.5$ |
| k=2 | 1/2 | 1/3 | $\frac{2*(1/2)*(1/3)}{(1/2+1/3)} = 0.4$ |
| k=3 | 2/3 | 2/3 | $\frac{2*(2/3)*(2/3)}{(2/3+2/3)} = 0.666$ |
| k=4 | 1/2 | 2/3 | $\frac{2*(1/2)*(2/3)}{(1/2+2/3)} = 0.571$ |
| k=5 | 3/5 | 1 | $\frac{2*(3/5)*1}{(3/5+1)} = 0.749$ |

Illustrations from <u>Evaluation Metrics For Information Retrieval</u> (Amit Chaudhary)

# Order-Aware Retrieval Metrics
## Mean Reciprocal Rank (MRR)

- Want first relevant item in higher position (i.e search engine, recommendation systems, Q&A systems)

- Doesn't care about position of remaining results

- $$MRR = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_i}$$



Reciprocal Rank

Query 1: 1 2 3 4 5 → 1 / 1 = 1

Query 2: 1 2 3 4 5 → 1 / 2 = 0.5

Query 3: 1 2 3 4 5 → 1 / 5 = 0.2

MRR = (1+0.5+0.2)/3 = 0.567

Illustrations from Evaluation Metrics For Information Retrieval (Amit Chaudhary)

# Order-Aware Retrieval Metrics
## Average Precision (AP)

- Evaluates whether all relevant items selected by the model is ranked highly

- $$AP = \frac{\sum_{k=1}^{n} (P(k) * rel(k))}{\text{number of relevant items}}$$

- $P(k)$: precision@k

- $rel(k)$: 1 when item at rank k is relevant, 0 otherwise



Illustrations from Evaluation Metrics For Information Retrieval (Amit Chaudhary)

# Order-Aware Retrieval Metrics
## Mean Average Precision (MAP)

- Taking mean of average precision across multiple queries

- $$MAP = \frac{1}{Q} \sum_{q=1}^{Q} AP(q)$$

# Order-Aware Retrieval Metrics
## Cumulative Gain (CG@k)

- Now suppose items also have some relevance scale


Relevance Scale

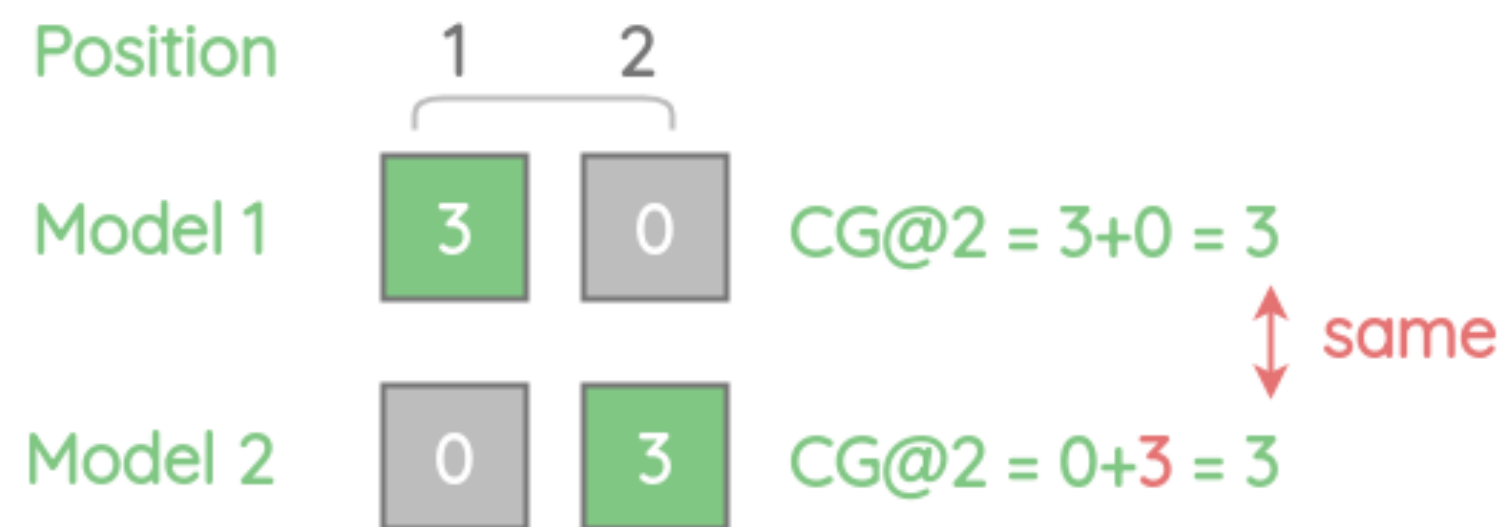- $$CG@k = \sum_{1}^{k} rel_i$$

| Position(k) | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Cumulative Gain@k | 3 | 3+2=5 | 3+2+3=8 | 3+2+3+0=8 | 3+2+3+0+1=9 |



Relevance: 3 2 3 0 1
Position: 1 2 3 4 5

cumulative gain@2 = 3+2 = 5

Illustrations from Evaluation Metrics For Information Retrieval (Amit Chaudhary)
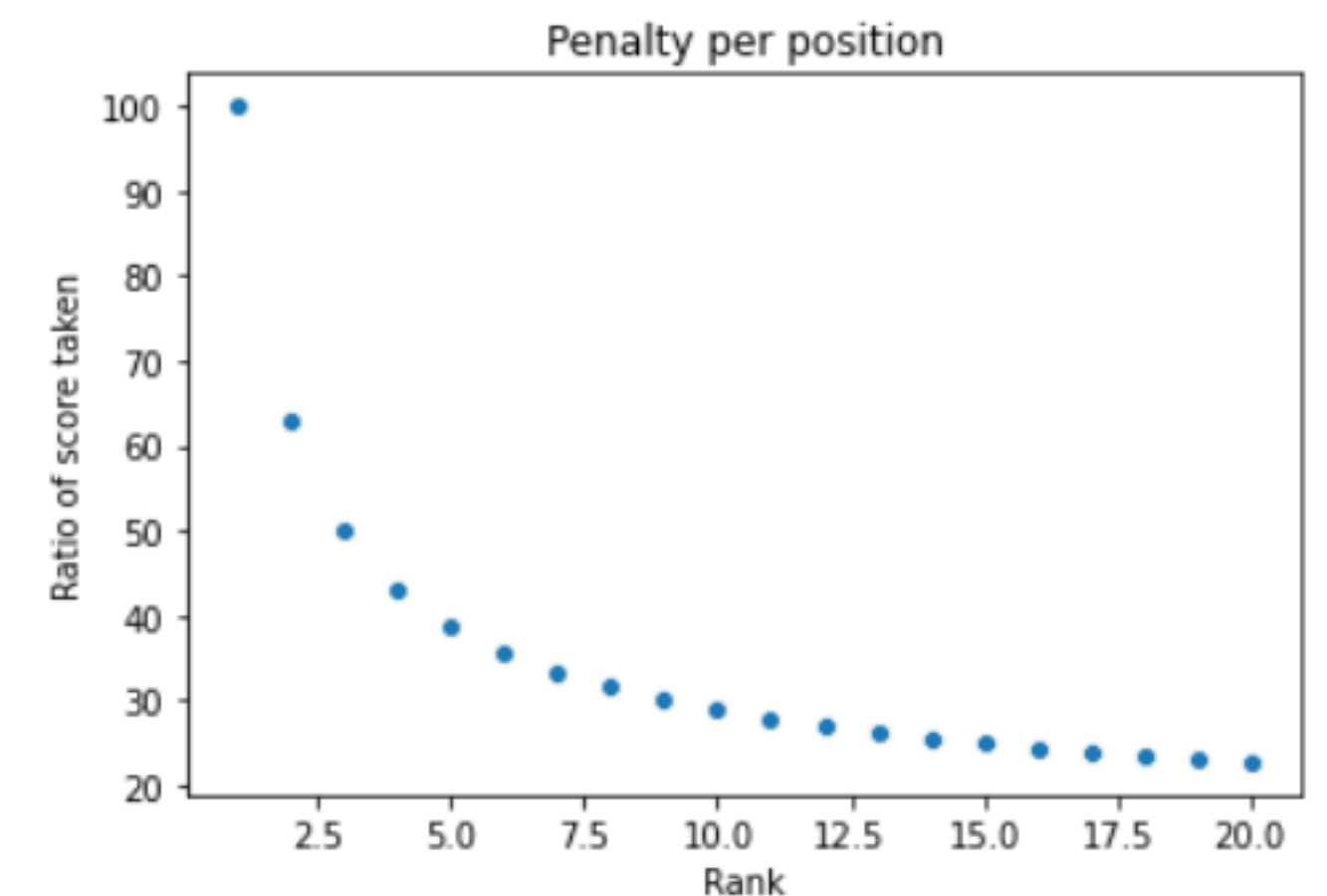
# Order-Aware Retrieval Metrics
## Discounted Cumulative Gain (DCG@k)

- Cumulative gain doesn't take into account order of relevant items, i.e CG@2 is same:



- DCG adds log-based penalty based on position

$$DCG@k = \sum_{i=1}^{k} \frac{rel_i}{log_2(i+1)}$$



Illustrations from Evaluation Metrics For Information Retrieval (Amit Chaudhary)

# Order-Aware Retrieval Metrics
## Normalized Discounted Cumulative Gain (NDCG@k)

- One issue with DCG is that it scales with length and magnitude of relevance score

  - Hard to compare across queries

- IDCG: the "ideal" DCG if things were put in ideal order



- Normalizing between 0-1: $\text{NDCG@k} = \dfrac{\text{DCG@k}}{\text{IDCG@k}}$

- The go-to metric in IR

Illustrations from <u>Evaluation Metrics For Information Retrieval</u> (Amit Chaudhary)