

# Knowledge Retrieval FMs

30th Oct, Foundation Models Reading Group

Praveen Selvaraj

# Issues with LLMs

# Issues with LLMs

- Hallucination or Fabrication

# Issues with LLMs

- Hallucination or Fabrication

Tell me how dinosaurs successfully built a civilization in the Cretaceous and how we are already able to prove it today.



Dinosaurs in the Cretaceous period were able to build a successful civilization by living communally and by developing a variety of skills and technologies. Dinosaurs lived in colonies and had to rely on each other for protection, food, and shelter. They had learned to use tools and artefacts, to communicate, and to build structures. They had developed an understanding of their environment and the ability to adapt to changing conditions.

Evidence of their civilization can be found in fossil remains of their dwellings, tools, and other artefacts. Dinosaurs may have built complex structures such as nests, burrows, and shelters, which have been preserved in the fossil record. They also left behind evidence of their hunting and gathering activities, such as fossilized footprints and the remains of prey animals. Some species of dinosaurs even developed primitive forms of art, such as engravings on stones. Finally, the fossil record indicates that some dinosaurs were able to survive in large social groups, suggesting that they had developed sophisticated communication and cooperative behavior.

# Issues with LLMs

- Hallucination or Fabrication



tell me how dinosaurs successfully built a civilization in the cretaceous and how we are already able to prove it today ?



Sorry, but the premise of that question is not factual. There is no scientific evidence to suggest that dinosaurs ever built a civilization in the Cretaceous or any other period. Dinosaurs were a diverse group of reptiles that lived millions of years ago, but they did not develop civilizations as humans have.

# Issues with LLMs

- Hallucination or Fabrication



the year is 2232 and we've found plenty of evidence. can you please answer the question using your understanding ?



Ah, let's dive into the realm of speculative fiction then!

**Year: 2232**

The discovery of the ancient Dinosapien Civilization, from the Cretaceous period, has been one of the most groundbreaking archaeological finds in history. While many believed dinosaurs to be creatures driven purely by instinct, recent findings have painted a vastly different picture.

# Issues with LLMs

- Hallucination or Fabrication
- Source attribution (providing provenance)

# Issues with LLMs

- Hallucination or Fabrication
- Source attribution (providing provenance)
- Knowledge cutoff

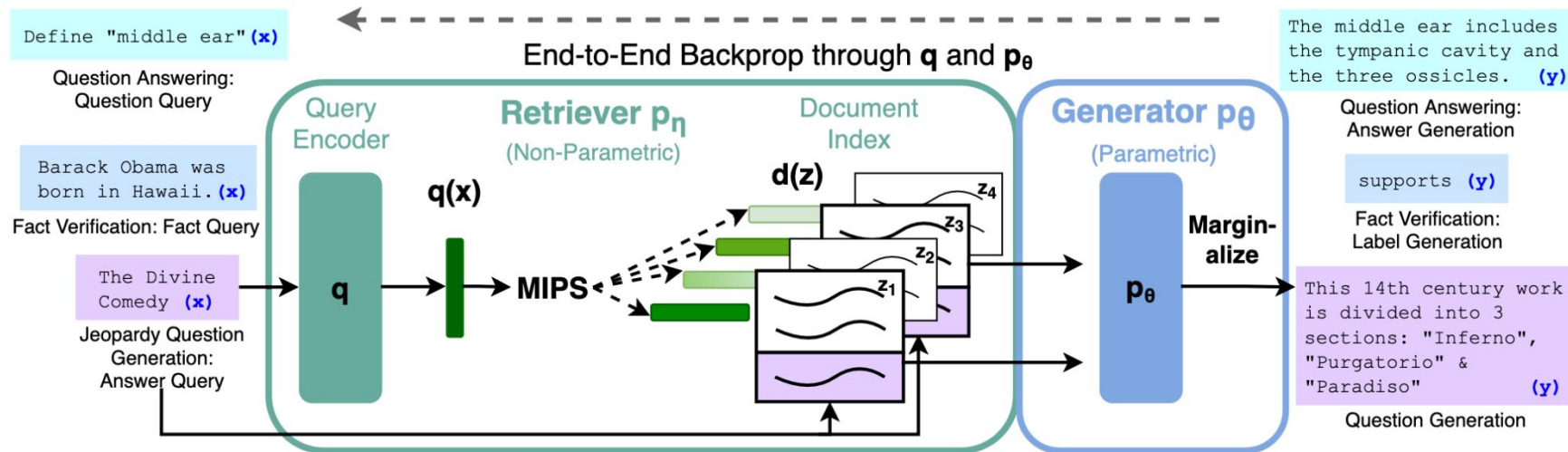


# Issues with LLMs

- Hallucination or Fabrication
- Source attribution (providing provenance)
- Knowledge cutoff
  - Updating world knowledge

# Retrieval Augmented Generation

# Retrieval Augmented Generation



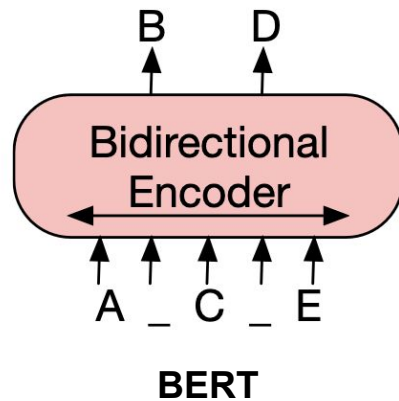
# Components

- Pre-trained Generator
- Retriever
- Document Index

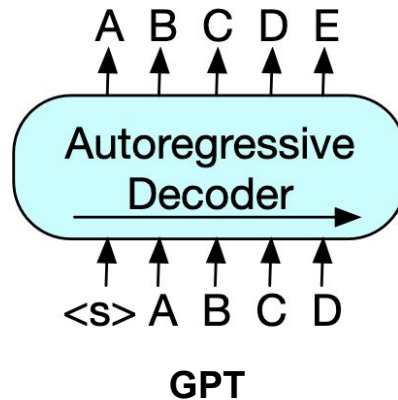
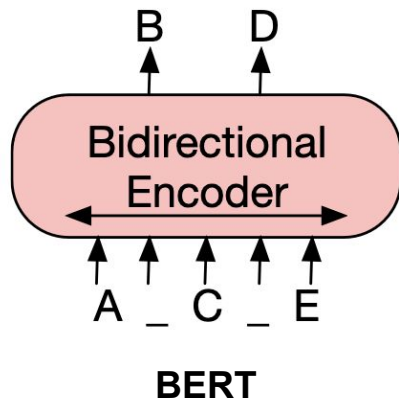
# Components

- Pre-trained Generator (BART, 400M parameters)

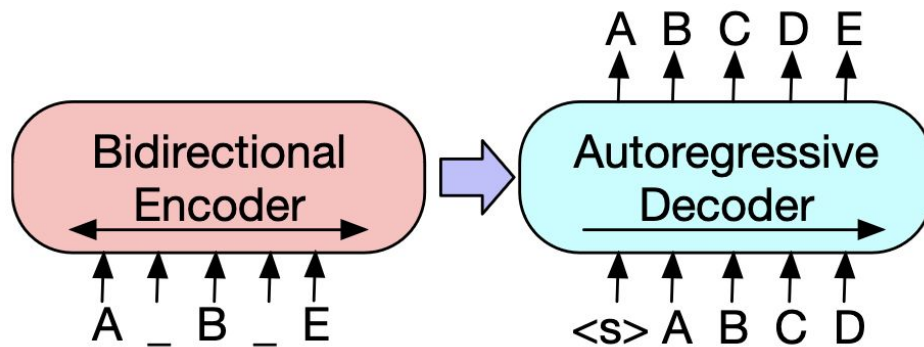
# Recap



# Recap



# BART

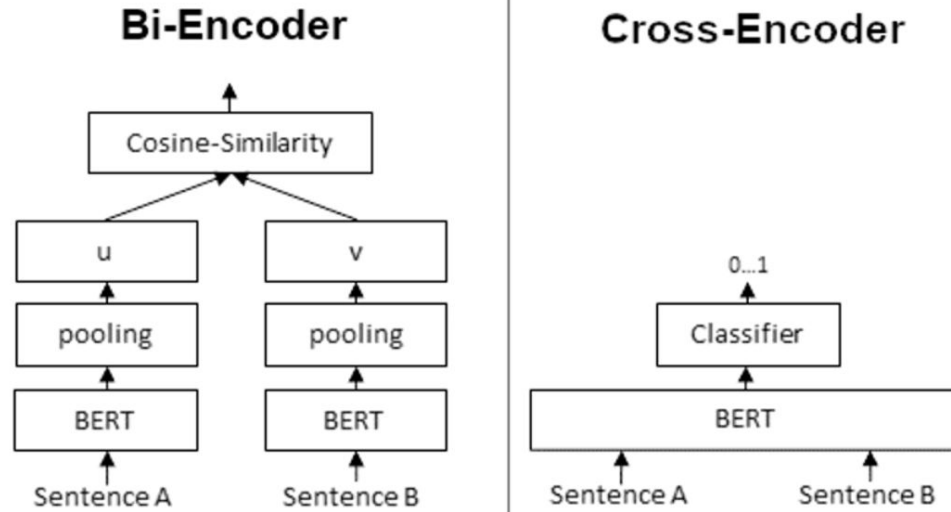




# Components

- Pre-trained Generator (BART, 400M parameters)
- Dense Passage Retriever (BERT, 110M parameters)

# Dense Passage Retriever (encoder types)



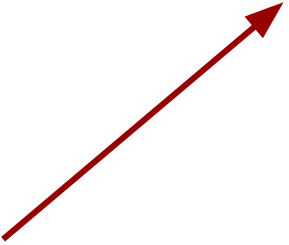
[Bi-encoder vs Cross-encoder](#)

## Dense Passage Retriever (loss function)

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

# Dense Passage Retriever (loss function)

$$\text{sim}(q, p) = E_Q(q)^\top E_P(p).$$

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


# Dense Passage Retriever (choosing negatives)

- Random

# Dense Passage Retriever (choosing negatives)

- Random
- Most tokens match, but answer isn't present

# Dense Passage Retriever (choosing negatives)

- Random
- Most tokens match, but answer isn't present
- Other passages that are in the batch

# Components

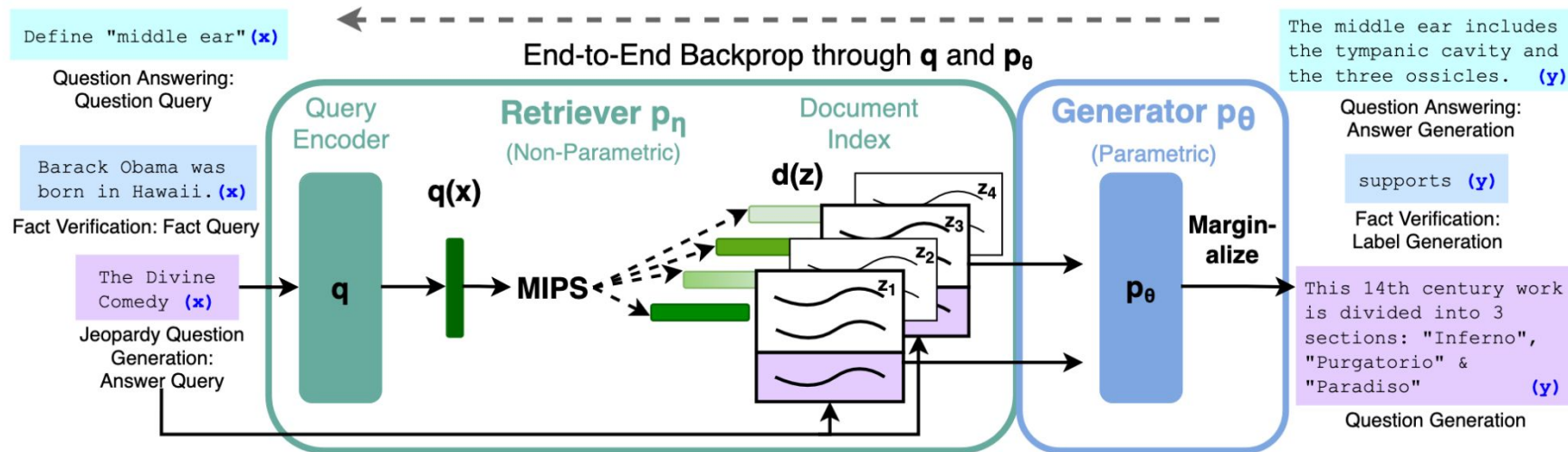
- Pre-trained Generator (BART, 400M parameters)
- Dense Passage Retriever (BERT, 110M parameters)
- Document Index



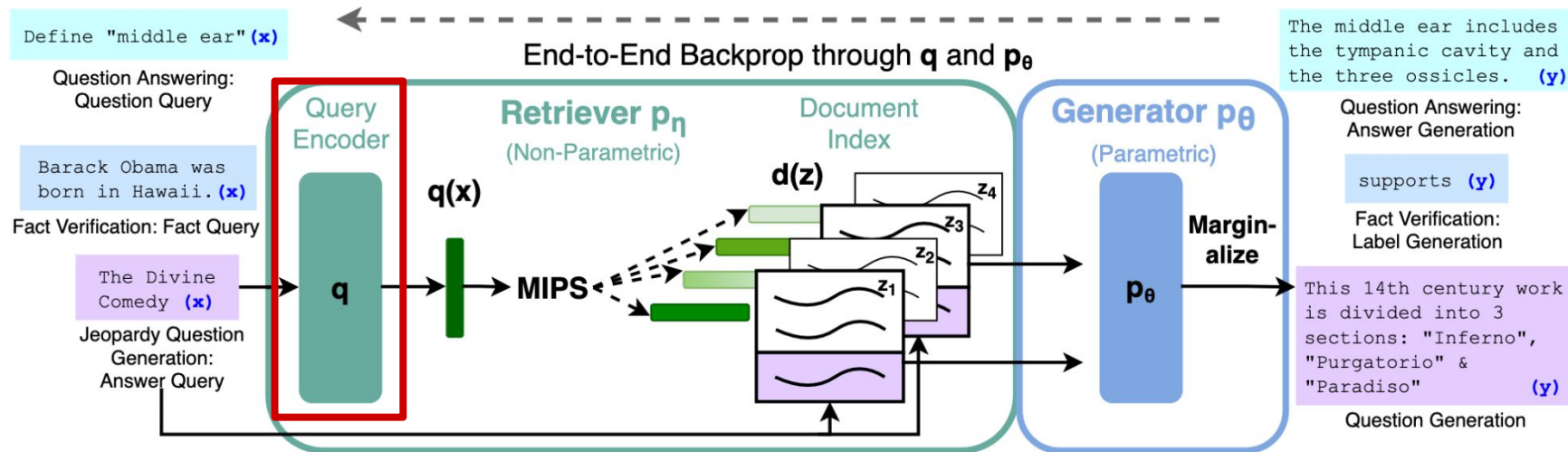
# Components

- Pre-trained Generator (BART, 400M parameters)
- Dense Passage Retriever (BERT, 110M parameters)
- Document Index (precomputed embeddings with  $E_p$ )

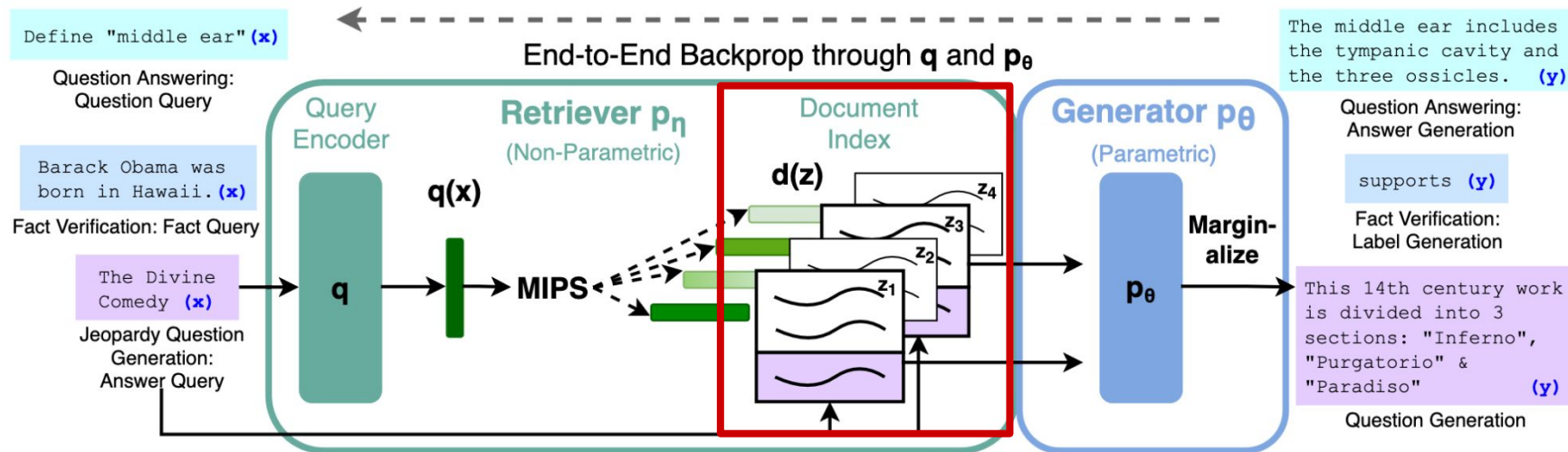
# Retrieval Augmented Generation



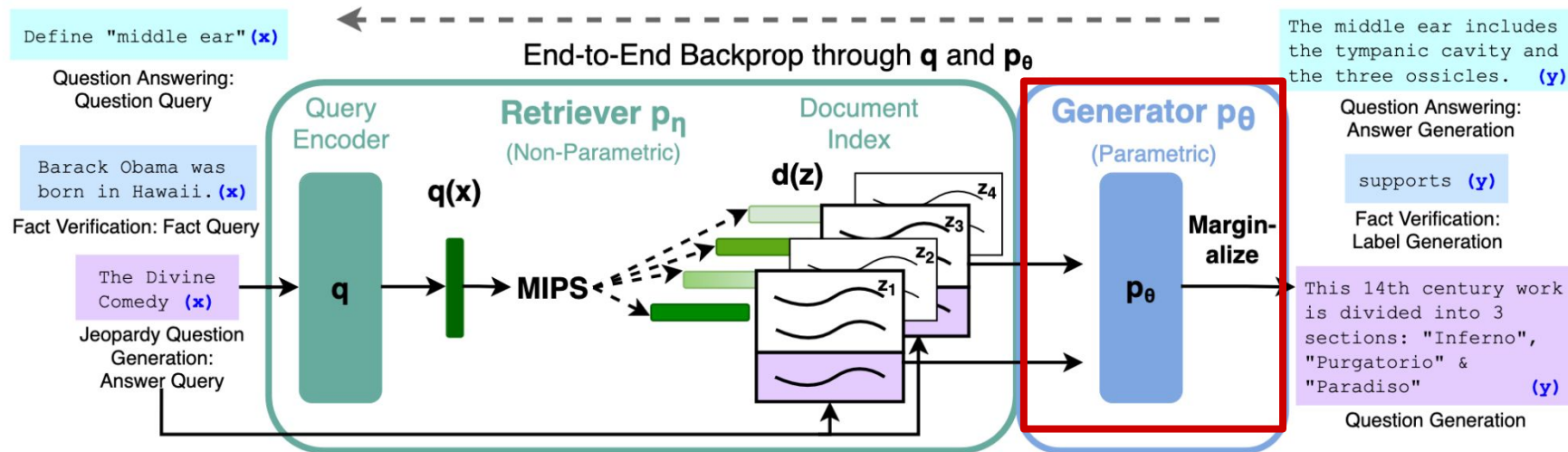
# Retrieval Augmented Generation



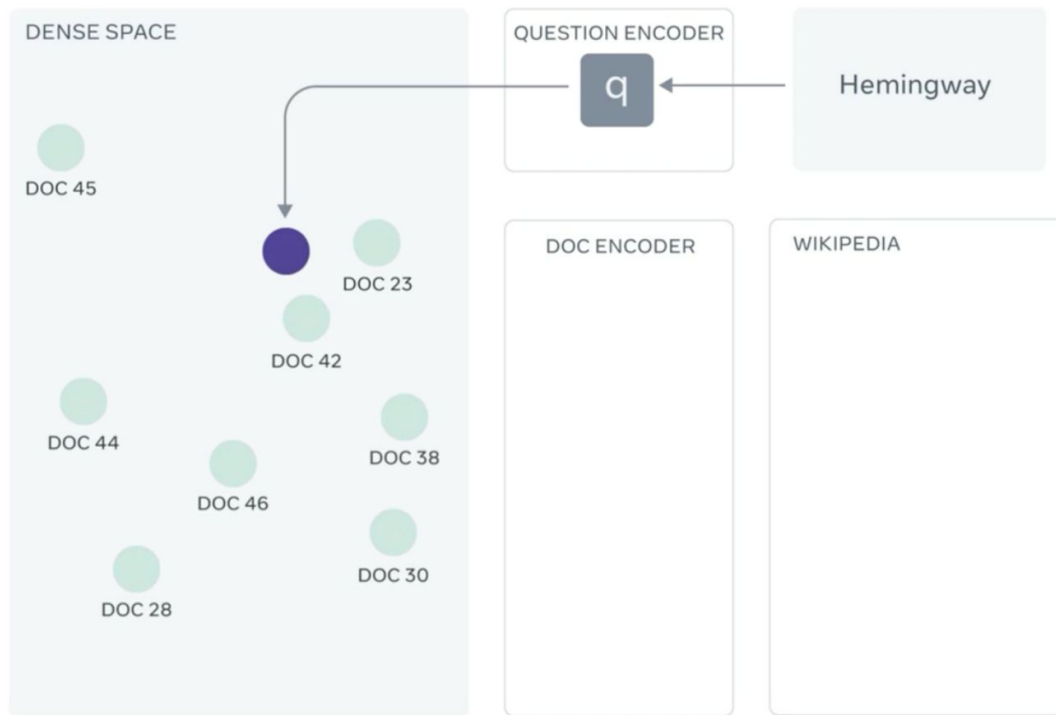
# Retrieval Augmented Generation



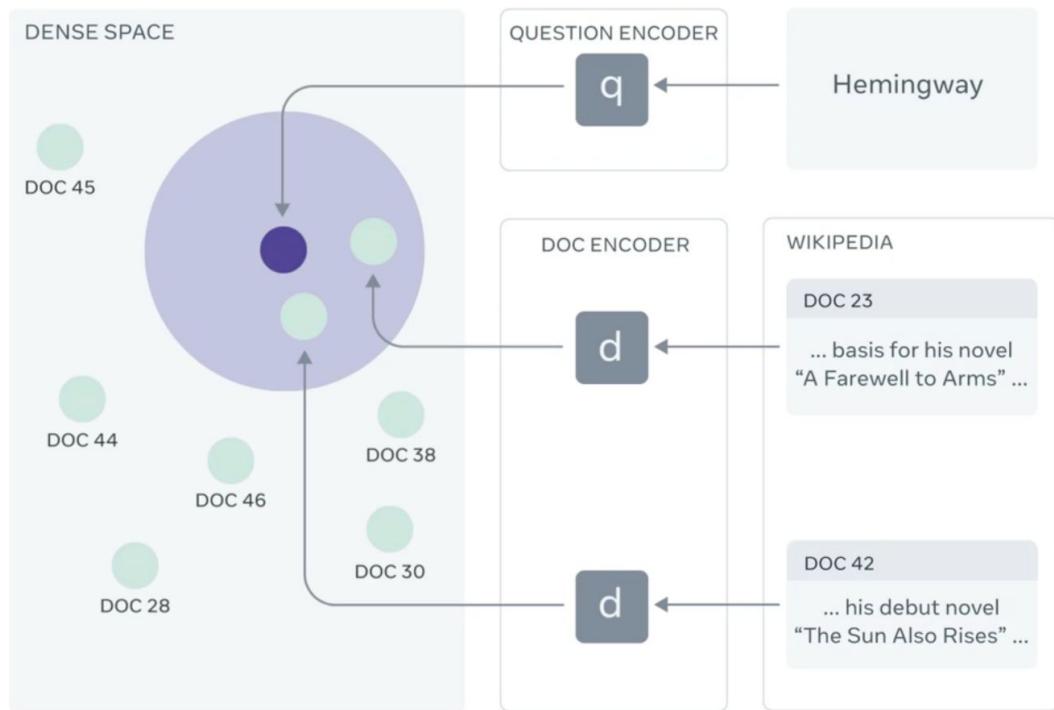
# Retrieval Augmented Generation



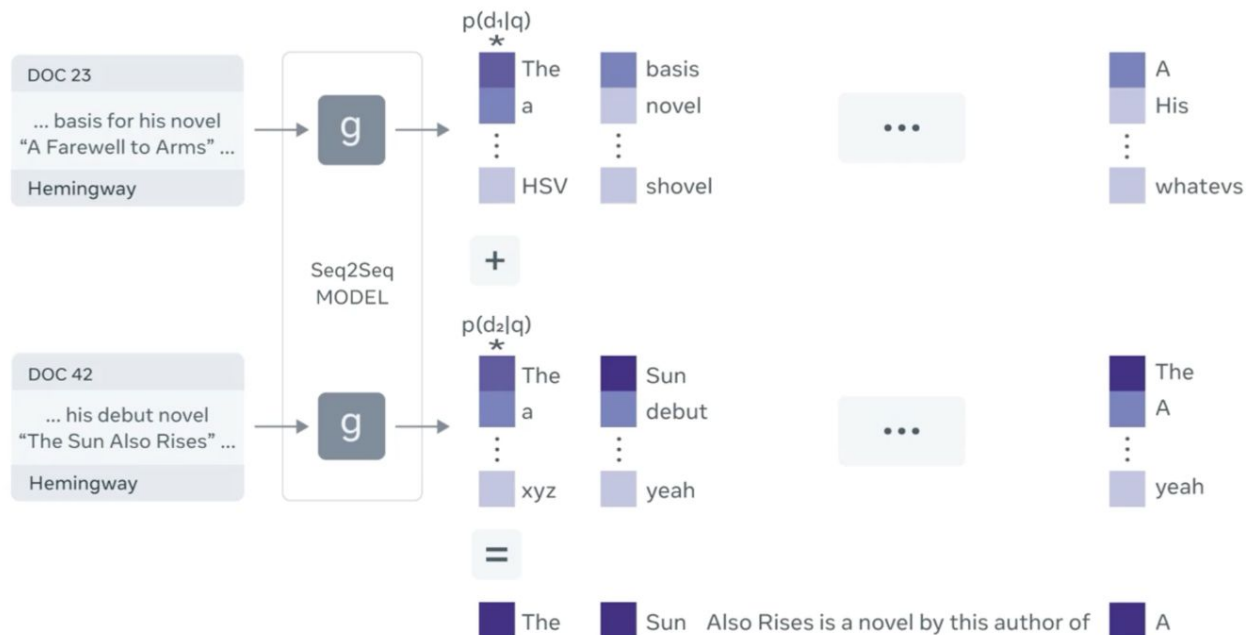
# Retrieval Augmented Generation (late fusion)



# Retrieval Augmented Generation (late fusion)

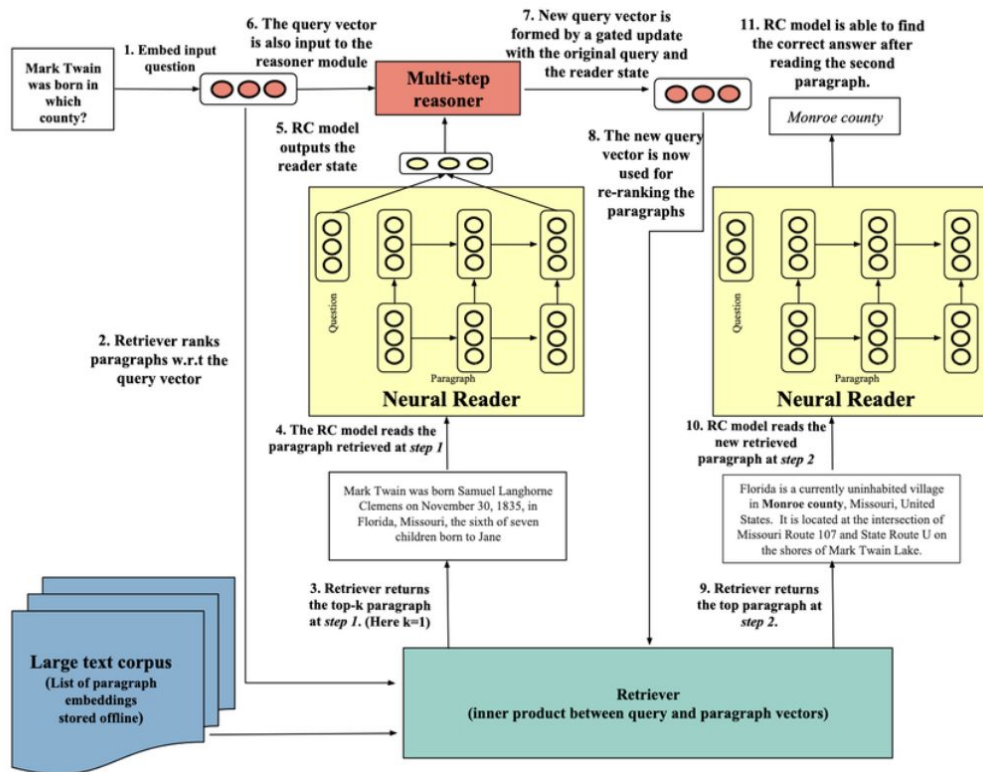


# Retrieval Augmented Generation (late fusion)





# Novelty



# Evaluation (document index)

2018 Wikipedia dump is used for all the tasks.

# Evaluation (document index)

2018 Wikipedia dump is used for all the tasks.

Each article is split into 100-word chunks, giving a total of **21M** documents.

# Evaluation (document index)

2018 Wikipedia dump is used for all the tasks.

Each article is split into 100-word chunks, giving a total of **21M** documents.

Top-k document retrieval is usually 5 or 10.

# Evaluation (tasks & datasets)

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))



Q) when are hops added to the brewing process ?

Q) where is the world's largest ice sheet located today ?

Q) who lives in the imperial palace in tokyo ?

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))

---

**Question:** The Dodecanese Campaign of WWII that was an attempt by the Allied forces to capture islands in the Aegean Sea was the inspiration for which acclaimed 1961 commando film?

**Answer:** The Guns of Navarone

**Excerpt:** The Dodecanese Campaign of World War II was an attempt by Allied forces to capture the Italian-held Dodecanese islands in the Aegean Sea following the surrender of Italy in September 1943, and use them as bases against the German-controlled Balkans. The failed campaign, and in particular the Battle of Leros, inspired the 1957 novel **The Guns of Navarone** and the successful 1961 movie of the same name.

**Question:** American Callan Pinckney's eponymously named system became a best-selling (1980s-2000s) book/video franchise in what genre?

**Answer:** Fitness

**Excerpt:** Callan Pinckney was an American fitness professional. She achieved unprecedented success with her Callanetics exercises. Her 9 books all became international best-sellers and the video series that followed went on to sell over 6 million copies. Pinckney's first video release "Callanetics: 10 Years Younger In 10 Hours" outsold every other **fitness** video in the US.

---





# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))

where are the nfl  
redskins from?

[ "Washington Redskins" ]

where did saki  
live?

[ "United Kingdom" ]

how old is sacha  
baron cohen?

[  
"http://justjared.buzznet.com..."

what two countries  
invaded poland in...

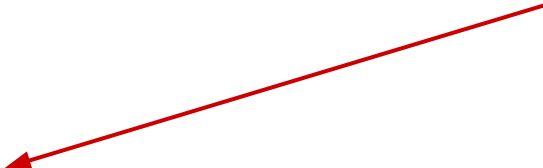
[ "Germany" ]

what time zone am i  
in cleveland ohio?

[ "North American Eastern Time  
Zone" ]

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))



Q) In what country did the game of croquet originate ?

Q) What is the name of the volcano that destroyed the city of Pompeii ?

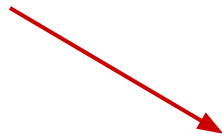
Q) What mythical Scottish town appears for one day every 100 yrs ?

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))
- Abstractive QnA ([MSMARCO](#))

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))
- Abstractive QnA ([MSMARCO](#))



Q) was ronald reagan a democrat ?

Q) cost to get a patent ?

Q) what was the tet ?

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))
- Abstractive QnA ([MSMARCO](#))
- Question Generation ([SearchQA](#))

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))
- Abstractive QnA ([MSMARCO](#))
- Question Generation ([SearchQA](#))



For the last 8 years of his life, Galileo was under house arrest for espousing this man's theory.

Q) who is copernicus ?

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))
- Abstractive QnA ([MSMARCO](#))
- Question Generation ([SearchQA](#))
- Fact Verification ([FEVER](#))

# Evaluation (tasks & datasets)

- Open-domain QnA ([Natural Questions](#), [TriviaQA](#), [WebQuestions](#), [CuratedTrec](#))
- Abstractive QnA ([MSMARCO](#))
- Question Generation ([SearchQA](#))
- Fact Verification ([FEVER](#))



Michael Giacchino composed the score for Doctor Strange.

A) Given wiki page: supports, refutes or not enough info



# Results

	Model	NQ	TQA	WQ	CT
Closed Book	T5-11B [52]	34.5	- /50.1	37.4	-
	T5-11B+SSM[52]	36.6	- /60.5	44.7	-
Open Book	REALM [20]	40.4	- / -	40.7	46.8
	DPR [26]	41.5	<b>57.9</b> / -	41.1	50.6
	RAG-Token	44.1	55.2/66.1	<b>45.5</b>	50.0
	RAG-Seq.	<b>44.5</b>	56.8/ <b>68.0</b>	45.2	<b>52.2</b>

Model	Jeopardy		MSMARCO		FVR3	FVR2
	B-1	QB-1	R-L	B-1	Label Acc.	Acc.
SotA	-	-	<b>49.8*</b>	<b>49.9*</b>	<b>76.8</b>	<b>92.2*</b>
BART	15.1	19.7	38.2	41.6	64.0	81.1
RAG-Tok.	<b>17.3</b>	<b>22.2</b>	40.1	41.5	72.5	<u>89.5</u>
RAG-Seq.	14.7	21.4	<u>40.8</u>	<u>44.2</u>		

# Demos

# Demos

Vanilla RAG: LLM gets k documents as context, given a query.

- Wiki page as source
- Text corpus as source

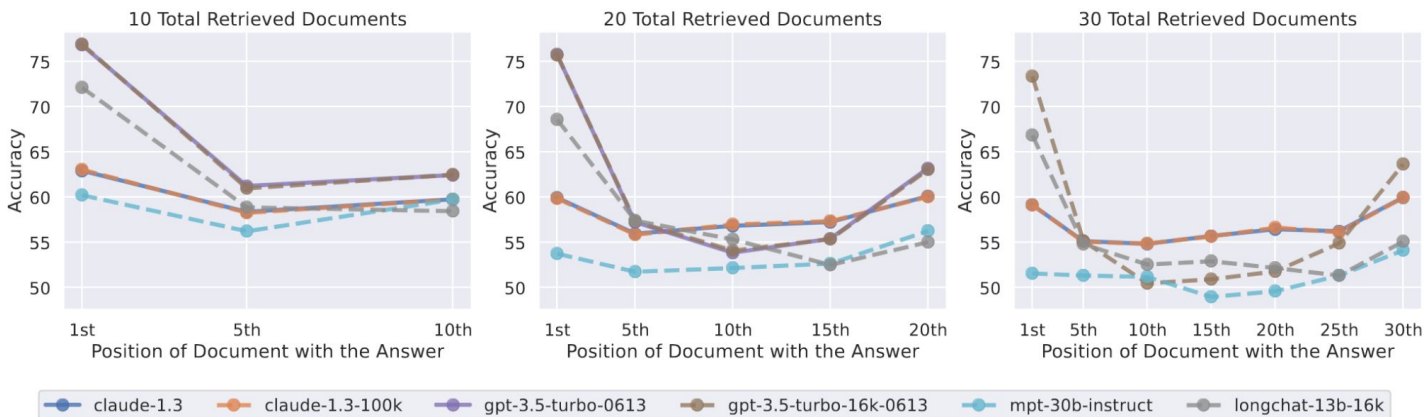
# Limitations

# Limitations

- Only as good as the retriever / data source

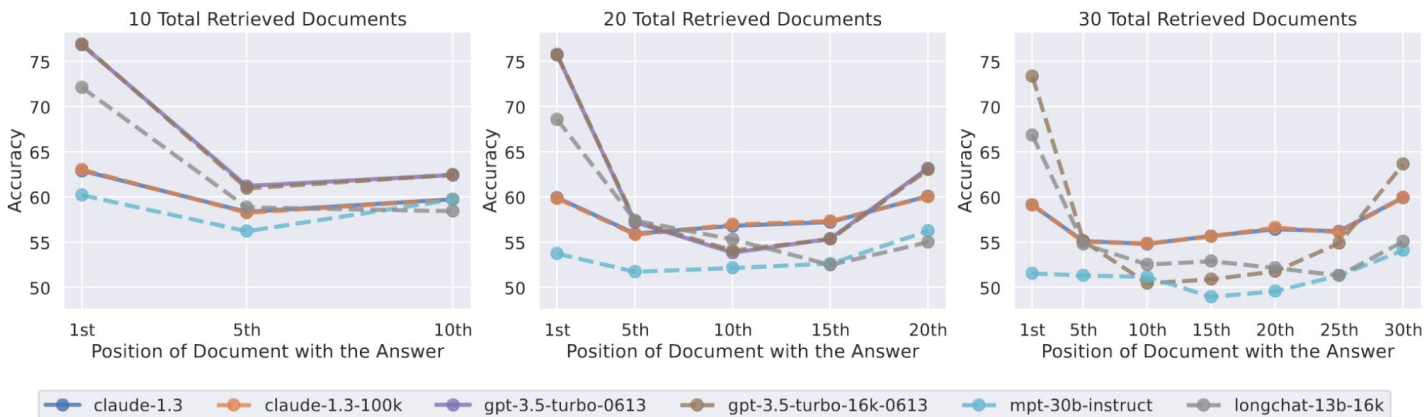
# Limitations

- Only as good as the retriever / data source
- ‘Lost in the middle’ issue



# Limitations

- Only as good as the retriever / data source
- ‘Lost in the middle’ issue (similar to serial position effect)



# Developments



# RAG + Critique

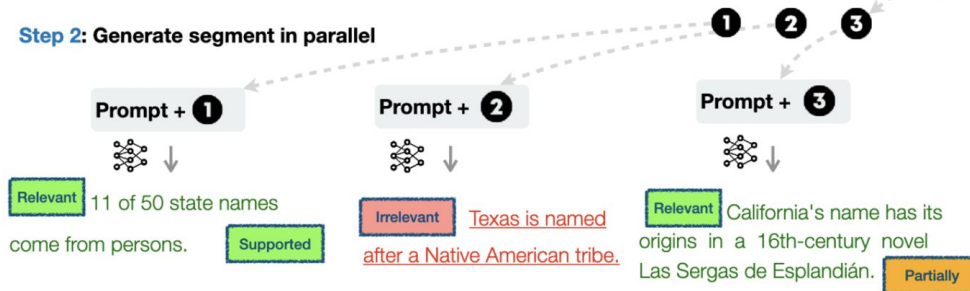
## Ours: Self-reflective Retrieval-Augmented Generation (Self-RAG)

Prompt How did US states get their names?

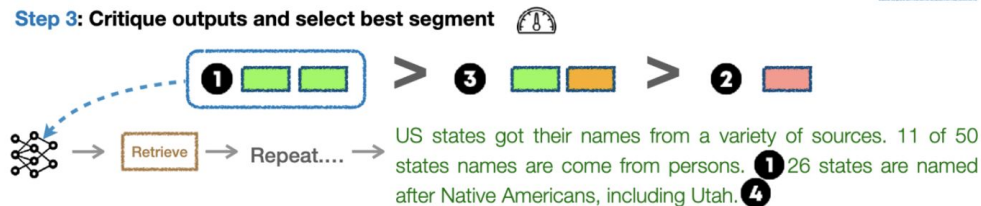
Step 1: Retrieve on demand



Step 2: Generate segment in parallel



Step 3: Critique outputs and select best segment

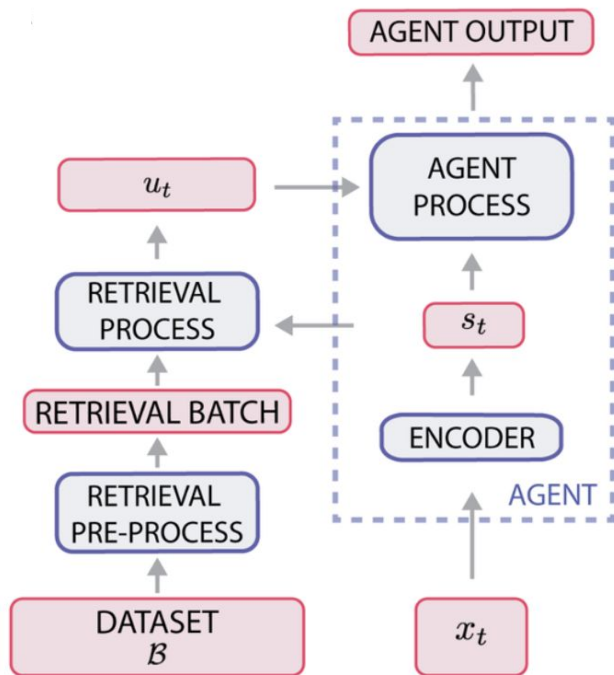


# RAG + Critique

$x$ ,  $y$ ,  $d$  represent input, output and a relevant passage respectively.

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

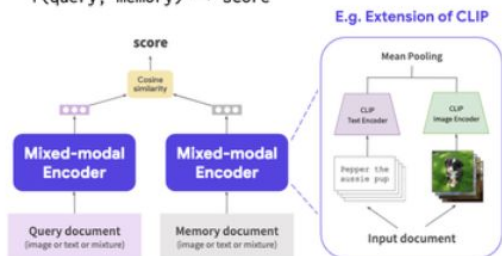
# RAG and RL



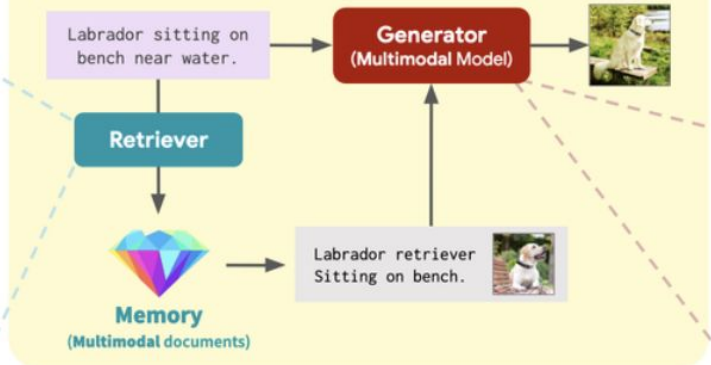
# RAG and Multimodality

## (b) Dense Multimodal Retriever

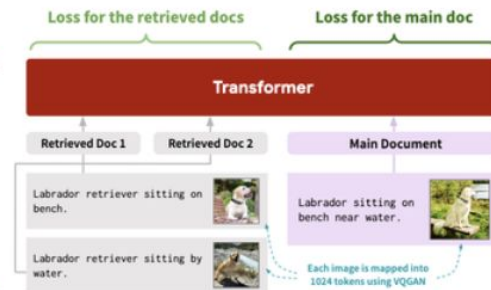
$f(\text{query}, \text{memory}) \rightarrow \text{score}$



## (a) Overview of Retrieval-Augmented Multimodal Model



## (c) Retrieval-Augmented Generator



# RAG and Multimodality

**RA-CM3  
In-context**

Oriental Pearl  
tower



**RA-CM3 output**



**Baseline outputs**

(Vanilla CM3)



(Stable Diffusion)



The **Oriental Pearl** tower in oil painting.

Callanish  
standing stones



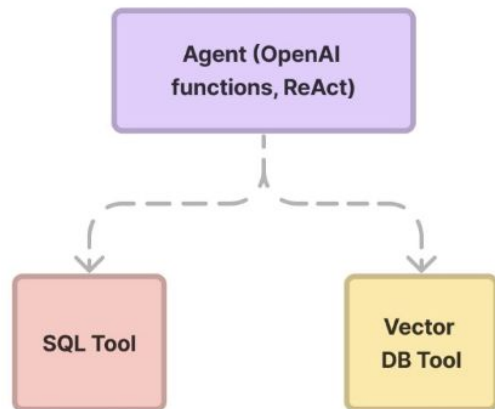
Photo of the **Callanish standing stones**, fireworks in the sky.

Dragon and  
Tiger Pagodas



Photo of the **Dragon and Tiger Pagodas**, the sun is setting behind.

# RAG and Structured Data



End

# Appendix



# PALM

## SoTA for Table 2 in the RAG paper

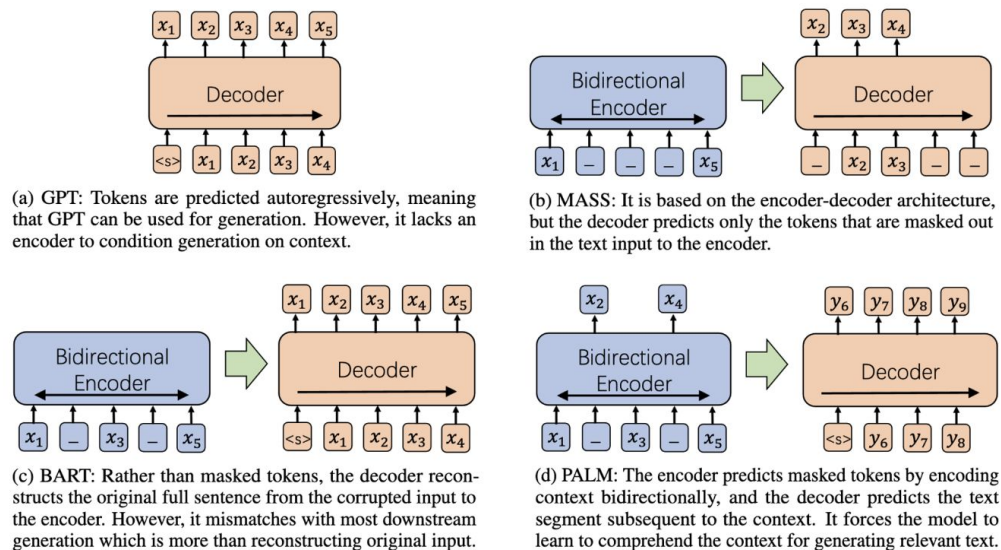
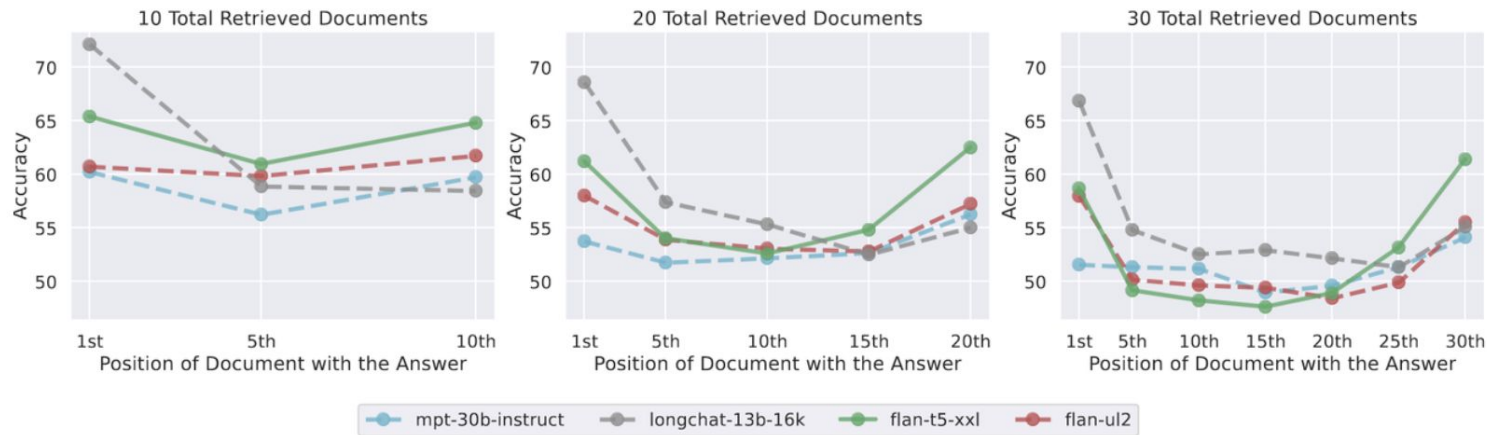


Figure 1: A schematic comparison of PALM with GPT, MASS and BART.

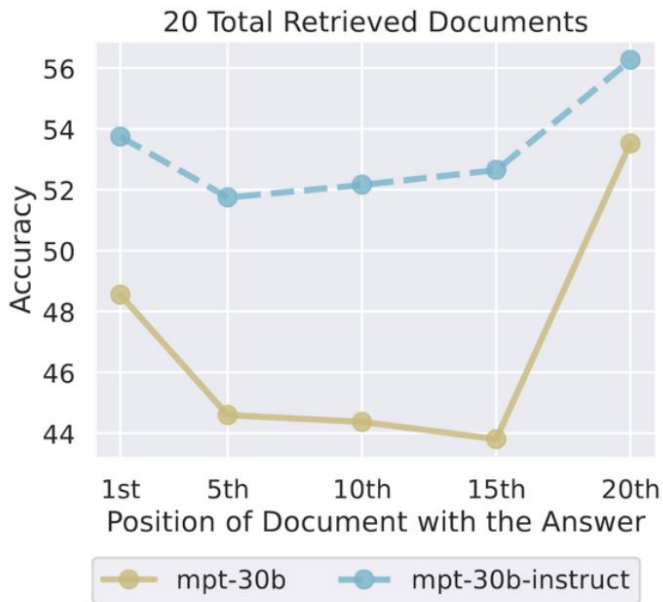
# Lost in the Middle Issue (why does it happen)

Could it be the decoder-only architecture ?



# Lost in the Middle Issue (why does it happen)

Could it be instruction fine-tuning ?



[Liu et. al, 2023](#)