GRIDSPACE IAP 2024 LECTURE 4
Tools for LLM Memory

January 18, 2024

# Questions from last time

- Why do you think ChatGPT did poorly on the 3-back task, even though it supposedly has superior memory to humans?

    It could be...

    - misinterpreting the task (even though it could recite a definition)

    - not effectively retrieving the information from the context

    - bad at counting the number of steps

# Questions from last time

- As LLMs get larger and larger, and trained with more and more data. How would you expect the models to perform on each of the memory categories?

  - Working Memory: No change as long as the context length stays the same. However, it might get better at extracting relevant data from the context

  - Episodic Memory: Same as working memory

  - Semantic Memory: Depends on the data, it could see more of the same facts be repeated and remember things better. Or it could get conflicting data and become less effective.

# Questions from last time

- LLM memorizing certain data can be a concern (e.g. private data, copyrighted data). What are some ways to alleviate this.

  There are many ways including:

  - Reducing repetition in data to reduce the chance of it remembering word for word

  - Using anonymized data (For cases of private data)
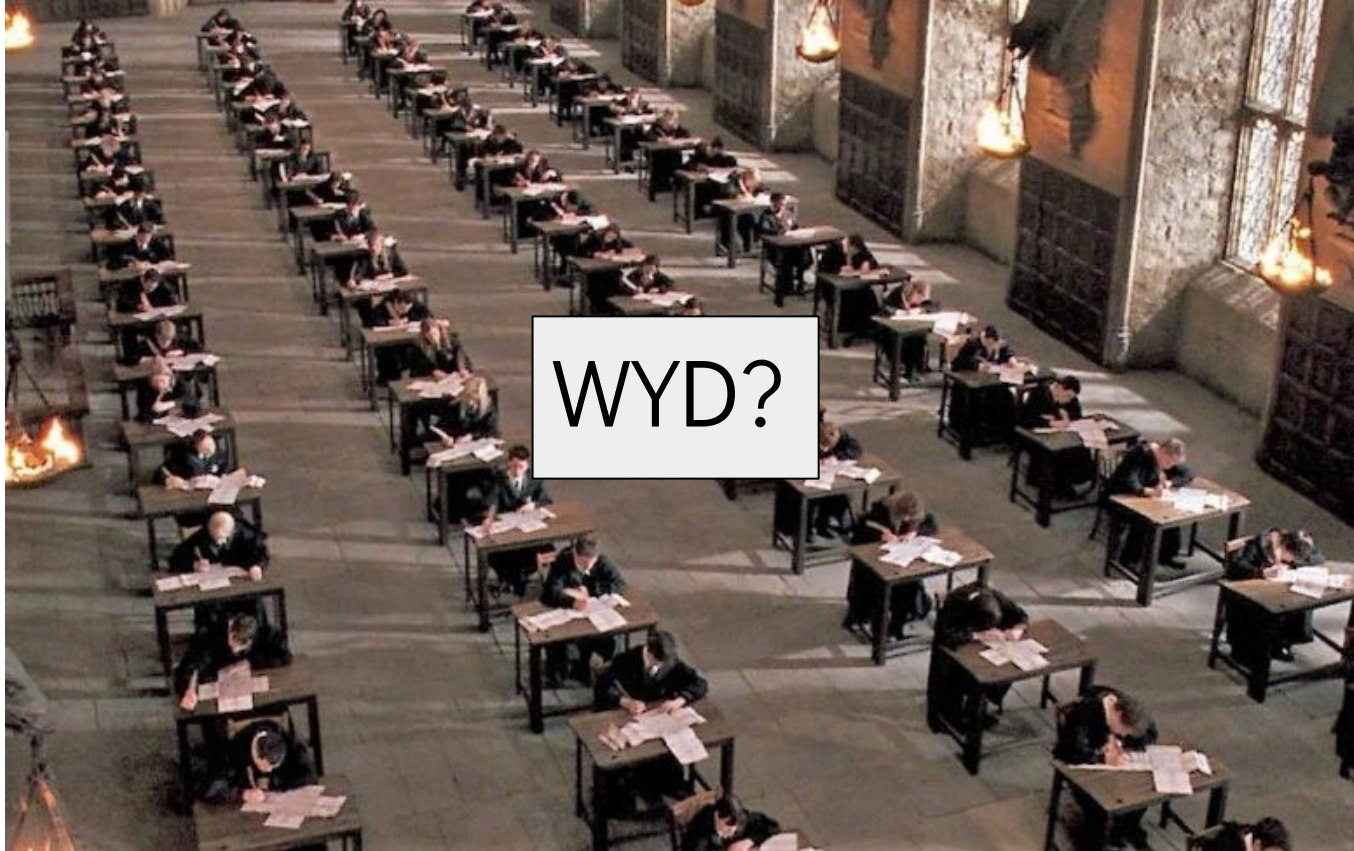
  - Create synthetic data (if applicable to use case)

YOU KNOW NOTHING

made on imgur

# Outline

- **Problem: Information retrieval**
- Solution 0: …
- Solution 1: …
- Solution 2: …
- Solution 3: …
- Solution 4: …

# Scenario: closed book exam in 2 weeks



WYD?

# Outline

- Problem: Information retrieval
- **Solution 0: Memory in training data**
- Solution 1: ...
- Solution 2: ...
- Solution 3: ...
- Solution 4: ...

# Recall from Nick

Scenario: You have three exams at the same time!
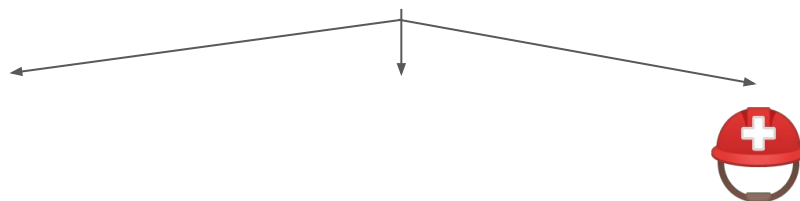(But you can clone yourself)
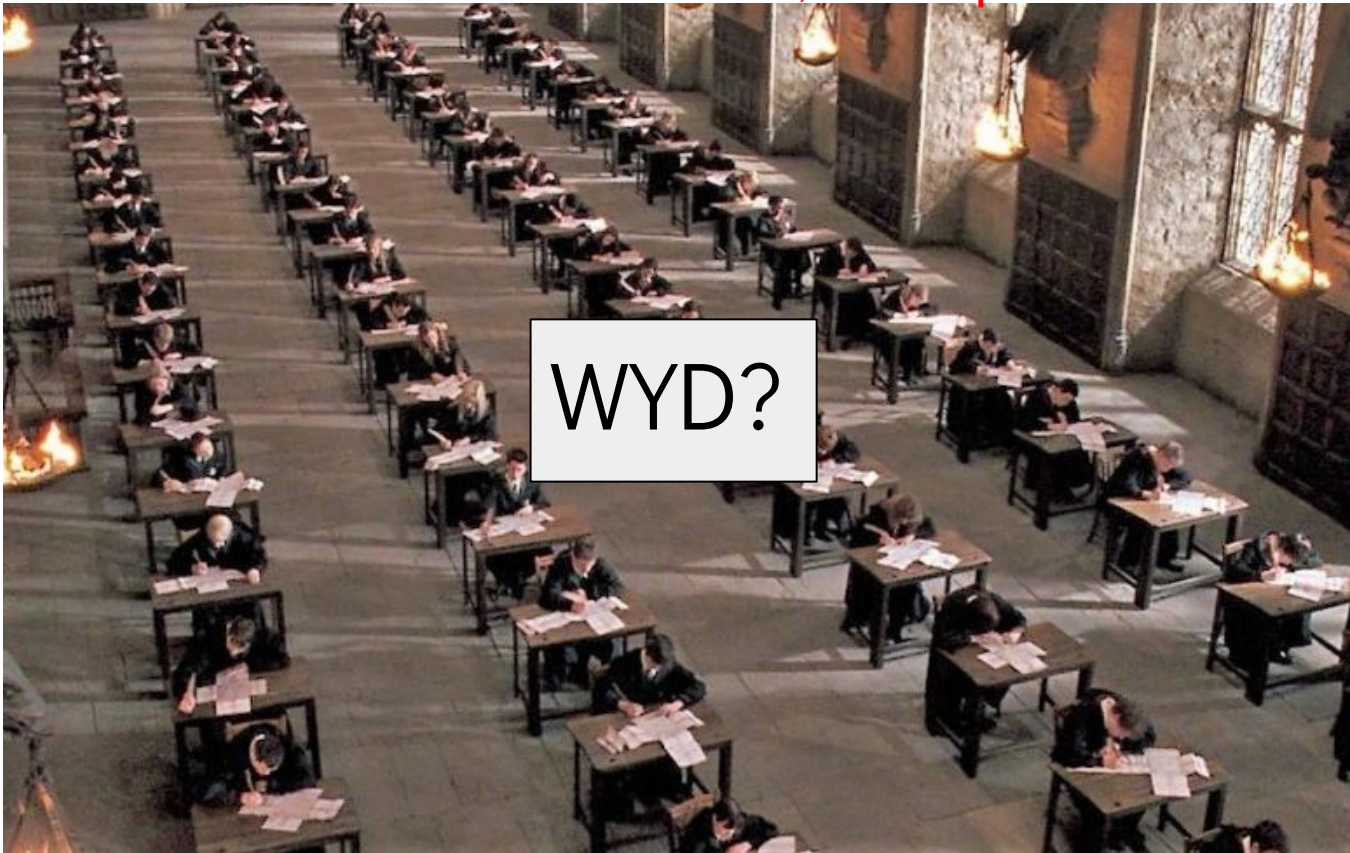


WYD?

# Memory in training data

- One base model
- (Possibly multiple) delta on base model for domain specialization
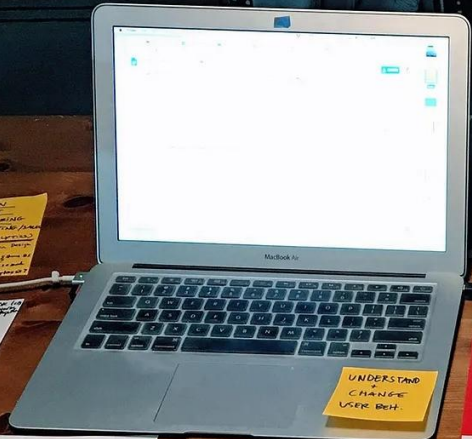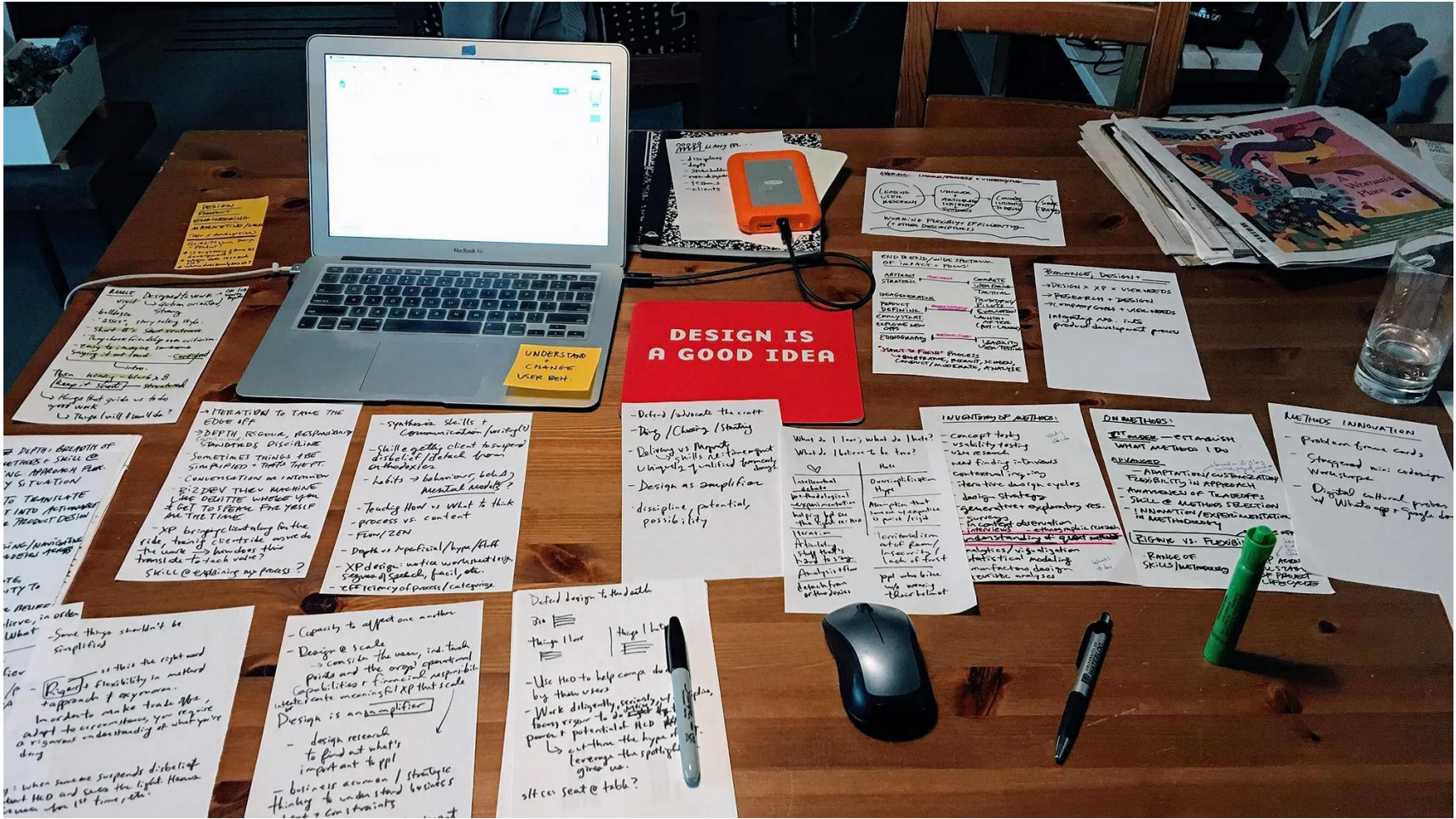- Train delta on new dataset

# Memory in training data

- Relatively "static", best for foundation knowledge, eg. language grammar, common sense (eg. the sky is blue)
- Expensive to retrain model when new knowledge is acquired

# Scenario: exam in 1 hour, but open notes



WYD?

# Outline

- Problem: Information retrieval
- Solution 0: Memory in training data
- **Solution 1: Memory in LLM Context**
- Solution 2: ...
- Solution 3: ...
- Solution 4: ...

# Recall from Nick

# Memory in LLM Context



tokens from past turns

V    V    V

*Mary had a little DONKEY*

K

K

K

Q

# Attention: Q, K, V

# Memory in  LLM Context



Mucus ad Nauseam

# Memory in LLM Context

- Storing K,V takes a lot of GPU memory
- Storing K, V can be messy (fragmented)

# Memory in LLM Context

- Economize, reduce size of K,V
  - Sliding window attention
  - Quantization
  - Group Query Attention

- Reduce memory waste due to messy storage (fragmentation)
  - PagedAttention

# Memory in LLM Context

Window size W,
a transformer of L layers,
Receptive field size: L x W



(b) Sliding window attention



https://paperswithcode.com/method/sliding-window-attention

# Memory in LLM Context

- Still limited, 32k tokens
- Still computationally expensive

# Scenario: your desk cannot fit all your notes!



WYD?

# Outline

- Problem:  Information retrieval
- Solution 0: Memory in training data
- Solution 1: Memory in LLM Context
- **Solution 2: Compress chat history/context**
- Solution 3: …
- Solution 4: …

# Compress chat history/context

- LLMs are really good at summarization



Text Summarization on GigaWord

- Let's summarize chat history/context so far!

# Compress chat history/context

- "In at most 1000 words, generate a summary of the chat history so far with sufficient detail to act as a replacement for the chat history in case we lose it. Pay special attention to instructions given by the user and system."

# Compress chat history/context

# Compress chat history/context

- Needs to run  LLM compressor regularly
- Some information is lost

# Scenario: you have incantation to summon an object



**Accio**

Summoning Charm. This charm summons an object to the caster, potentially over a significant distance.

Accio book

WYD?

# Outline

- Problem:  Information retrieval
- Solution 0: Memory in training data
- Solution 1: Memory in LLM Context
- Solution 2: Compress chat history
- **Solution 3:  Memory in external database**
- Solution 4: …

# Recall from week 1: LLM can use tools

# Recall from week 1: LLM can use tools

System: You may use these tools: "check_temperature", args: ... usage: ...

User: What is the temperature in LA now?

LLM: \<function call\> "check_temperature", args: "location: LA, time: now"

User: \<function return\> "check_temperature", results: "60F"

LLM: The temperature in LA is currently 60F! Nice day for hiking!

# Memory in external database

harry_potter.txt

DB

System: You may use these tools: "retrieve_from_db"...

User: What is Harry Potter's wand made of?

LLM: <function call> "retrieve_from_db", args: "query: harry wand material"

User: <function return> "retrieve_from_db", results: {text:
"Tricky customer, eh? Not to worry, we'll find the perfect match here somewhere –
I wonder, now – yes, why not – unusual combination – holly and phoenix feather,
eleven inches, nice and supple.", location: "book 1, page 65, paragraph 6"}

LLM: Harry Potter's wand is made of phoenix feather.

# Memory in external database



DB

System: You may use these tools: "retrieve_from_db"...

User: Order me a hamburger

(… 2000 messages later …)

LLM: …

LLM: <function call> "retrieve_from_db", args: "query: what did I get for lunch?"

User: <function return> "retrieve_from_db", results: {text: "Order me a hamburger", location: ...}

LLM: You got hamburger.

# Memory in external database (Reversed Index)

- Can be BM25

harry_potter.txt

<query string>

DB

# Memory in external database (Transformer-based retrieval)



harry_potter.txt

DB

&lt;query string&gt;

&lt;relevant snippets&gt;

- Can be transformer-based
- eg. Colbertv2

**ColBERTv2:**
**Effective and Efficient Retrieval via Lightweight Late Interaction**

**Keshav Santhanam***
Stanford University

**Omar Khattab***
Stanford University

**Jon Saad-Falcon**
Georgia Institute of Technology

**Christopher Potts**
Stanford University

**Matei Zaharia**
Stanford University

**Abstract**

Neural information retrieval (IR) has greatly advanced search and other knowledge-intensive language tasks. While many neural IR methods encode queries and documents into single-vector representations, late interaction models produce multi-vector representations at the granularity of each token and decompose relevance modeling into scalable token-level computations. This decomposition has been shown to make late interaction more relevance is estimated using rich yet scalable interactions between these two sets of vectors. ColBERT produces an embedding for every token in the query (and document) and models relevance as the sum of maximum similarities between each query vector and all vectors in the document.

By decomposing relevance modeling into token-level computations, late interaction aims to reduce the burden on the encoder: whereas single-vector models must capture complex query–document re-

# Memory in external database (Transformer-based retrieval)

harry_potter.txt

<query_string> eg.
"harry potter wand material"

V

Q

DB

<relevant snippets> eg.
"Tricky customer, eh?..."

# Recall from week 1: Voyager skill retrieval

**Program Generated by GPT-4**

```
async function combatZombie(bot) {
    // Equip a weapon
    const sword =
    bot.inventory.findInventoryItem(
        mcData.itemsByName[
            "stone_sword"
        ].id
    );
    if (sword) {
        await bot.equip(sword, 'hand');
    } else {
        await craftStoneSword(bot);
        ...
    }
    // Craft and equip a shield
    ...
    // Recover hunger
    ...
    // Look for and combat a zombie
    ...
}
```

**Program Description**

```
async function combatZombie(bot) {
    // The function is about
equipping a stone sword to combat
a zombie. If a stone sword is not
found, it will craft one.
Additionally, it crafts and equips
a shield for added protection.
Afterwards, it proceeds to cook
sticks in order to restore hunger.
Once hunger is replenished, it
actively searches for a zombie and
engages in combat with it.
}
```

GPT-3.5 · Embedding · **Key**

Add · **Value**

**Skill Library**

- Mine Wood Log
- Make Crafting Table
- Craft Wooden Pickaxe
- Craft Stone Sword
- Make Furnace
- ...
- Combat Cow
- Cook Steak
- Craft Iron Axe
- Combat Zombie

**Task: Craft Iron Pickaxe**

How to craft an iron pickaxe in Minecraft?

GPT-3.5

To craft an iron pickaxe, you need to 3 iron ingots and 2 sticks. Once you have gathered the materials, ....

------------------------------

**Environment Feedback**

Embedding · **Query** · **Skill Library** · Retrieve

**Top-5 Relevant Skills**

- Smelt Iron Ingot
- Craft Stick
- Make Crafting Table
- Make Furnace
- Craft Wooden Pickaxe

# Memory in external database (Transformer-based retrieval)



skills.py

<query_string>

V, K

Q

DB

# Memory in external database

the Good:
- Infinite storage

the Bad:
- Needs to build DB
- DB can be tricky to maintain

- Is there already an existing large DB on the shelf?

Scenario: Hogwarts have wifi and nothing in the rules says you cannot summon a computer



**Accio**

Summoning Charm. This charm summons an object to the caster, potentially over a significant distance.

*Accio computer*

# Outline

- Problem: Information retrieval
- Solution 0: Memory in training data
- Solution 1: Memory in LLM Context
- Solution 2: Compress chat history
- Solution 3: Memory in external database
- **Solution 4: Let LLM browse the internet**

# Let LLM browse the internet

- Lookup Google
- Lookup Wikipedia
  - Can do "fact checks" too!

# WikiChat: Stopping the Hallucination of Large Language Model Chatbots by Few-Shot Grounding on Wikipedia

Sina J. Semnani    Violet Z. Yao*    Heidi C. Zhang*    Monica S. Lam

Computer Science Department
Stanford University
Stanford, CA
{sinaj, vyao, chenyuz, lam}@cs.stanford.edu

## Abstract

This paper presents the first few-shot LLM-based chatbot that almost never hallucinates and has high conversationality and low latency. WikiChat is grounded on the English Wikipedia, the largest curated free-text corpus.

WikiChat generates a response from an LLM,

after their pre-training, and are far less knowledgeable about less popular, or *tail*, topics (Mallen et al., 2022; Sun et al., 2023). Therefore, for knowledge-intensive tasks (Lewis et al., 2020), users need to painstakingly verify any information they receive with external sources lest they be misled.

This paper focuses on three metrics for

https://arxiv.org/pdf/2305.14292v2.pdf

**User:** Have you heard about the upcoming film Oppenheimer?
**WikiChat:** Yes! Oppenheimer is a biopic …
**User:** What do you think about the casting?

Christopher Nolan has cast Cillian Murphy as the titular role of Oppenheimer, keeping with his tradition of meticulous casting choices. I'm sure it will be an excellent film!

**7. Refine**

**1. Query**

**Query:** cast of Oppenheimer Christopher Nolan film
**Date:** "recent"

**Title: Christopher Nolan**
Nolan's twelfth directorial venture will be "Oppenheimer" (2023), …

**Title: Oppenheimer (film)**
… Much of the main cast signed on between September 2021 and April 2022., …

**Title: Cillian Murphy**
… Murphy was confirmed to star as J. Robert Oppenheimer in Nolan's upcoming film Oppenheimer …

Retrieve

**2. Summarize & Filter**

- Much of the Oppenheimer (2023) movie's cast signed on between September 2021 and April 2022.

- Cilian Murphy is to play as Robert Oppenheimer in the Oppenherimer (2023) movie.

Cillian Murphy has been cast to play Robert Oppenheimer in Christopher Nolan's "Oppenheimer" movie. Christopher Nolan is known for his meticulous casting choices and often casts well-known actors in his films. The movie's cast was signed on between September 2021 and April 2022.

**6. Draft**

Christopher Nolan is known for his meticulous casting choices. He often casts well-known actors in his films, and this time he has selected Tom Hanks and … Michael Caine.

**3. Generate**

- Christopher Nolan is known for his meticulous casting choices.
- Christopher Nolan often casts well-known actors in his films.
- Tom Hanks stars in the Oppenheimer (2023) movie.
- Michael Caine stars in the Oppenheimer (2023) movie.

**4. Extract Claims**

**5. Fact-Check**

- Christopher Nolan is known for his meticulous casting choices.
- Christopher Nolan often casts well-known actors in his films.

Retrieve

# Let LLM browse the internet

the Good:
- Simple to implement
- Cheap computation (on our side)

the Bad:
- Beware of internet trolls
- Dependent on third party servers

# Outline

- Problem: Information retrieval
- Solution 0: Memory in training data
- Solution 1: Memory in LLM Context
- Solution 2: Compress chat history
- Solution 3: Memory in external database
- Solution 4: Let LLM browse the internet
- Q&A

# Exercises:

- [Mistral-7B-v0.1](#) uses sliding window attention. What's its window length & number of layers? What's its theoretical attention span?
- Can you suggest some reasons that the attention span in practice is lower than theoretical span? How to remedy?
- What are some of the pain points of depending on third-party search engines?

# References

- https://mistral.ai/news/mixtral-of-experts/

- https://huggingface.co/blog/moe

- S Semnani, et al., WikiChat: Stopping the Hallucination of Large Language Model Chatbots by Few-Shot Grounding on Wikipedia. https://arxiv.org/abs/2305.14292

- https://paperswithcode.com/method/sliding-window-attention

- https://huggingface.co/blog/4bit-transformers-bitsandbytes

- https://paperswithcode.com/method/grouped-query-attention

- PagedAttention: https://arxiv.org/abs/2309.06180

- https://voyager.minedojo.org/

- https://mistral.ai/news/announcing-mistral-7b/

# Image Credits

- https://dwellingondreamspodcast.com/2019/08/28/classes-at-hogwarts-we-wish-wed-heard-more-about/

- https://imgflip.com/memetemplate/124827840/You-know-nothing

- https://twitter.com/TheHPfacts/status/848409690266533888

- https://harrypotter.fandom.com/wiki/Ordinary_Wizarding_Level

- https://harrypotter.fandom.com/wiki/Hogwarts_Library

- https://twitter.com/HPotterUniverse/status/628967209252032512

- https://uxdesign.cc/what-kind-of-designer-researcher-are-you-2da7598b17d4

- https://www.independent.co.uk/life-style/health-and-families/features/why-being-messy-can-be-a-positive-trait-according-to-researchers-a6774206.html

- https://gradecalculator.mes.fm/memes/fingernail-cheat-sheet

- https://www.gsm-earpiece.com/howto/tips-on-cheating-exam/

- https://en.wikipedia.org/wiki/Transformer