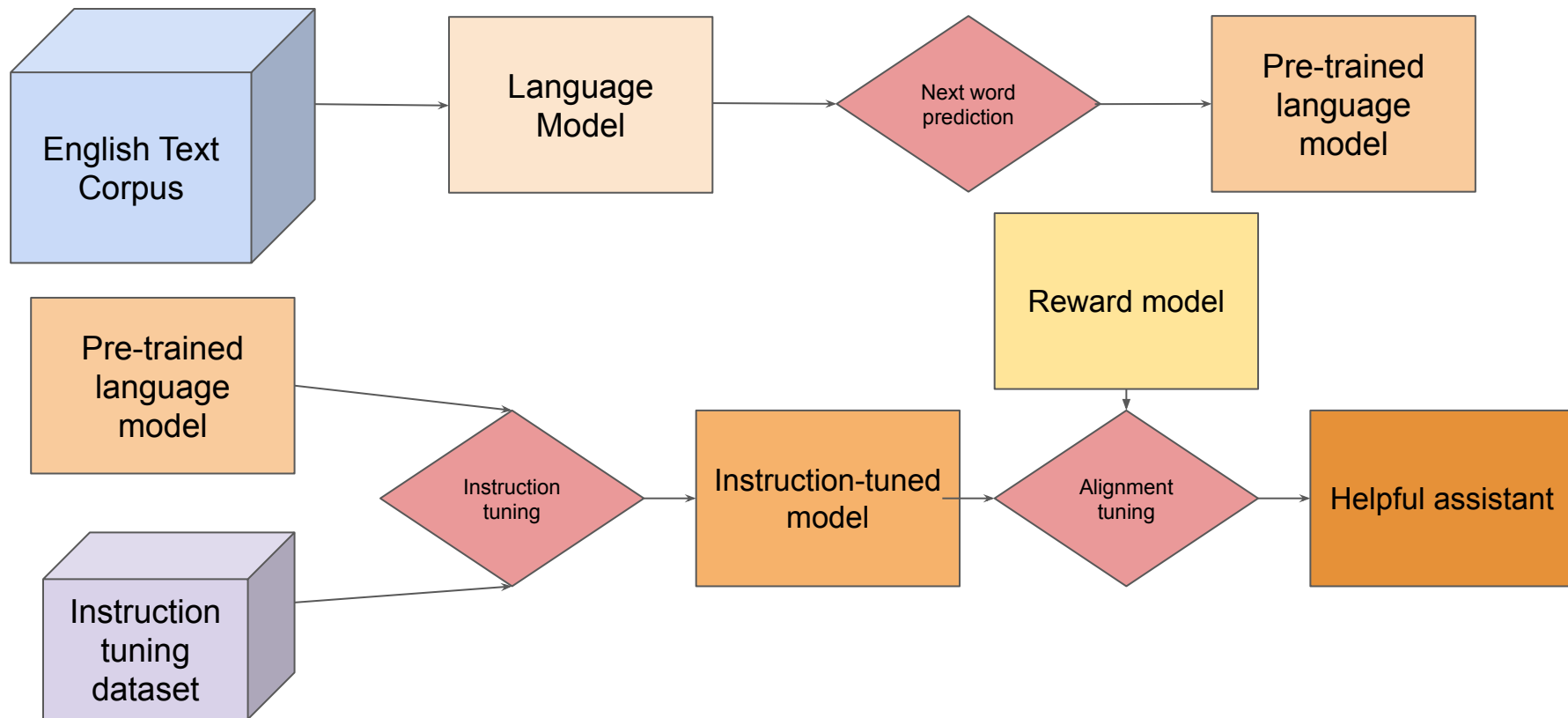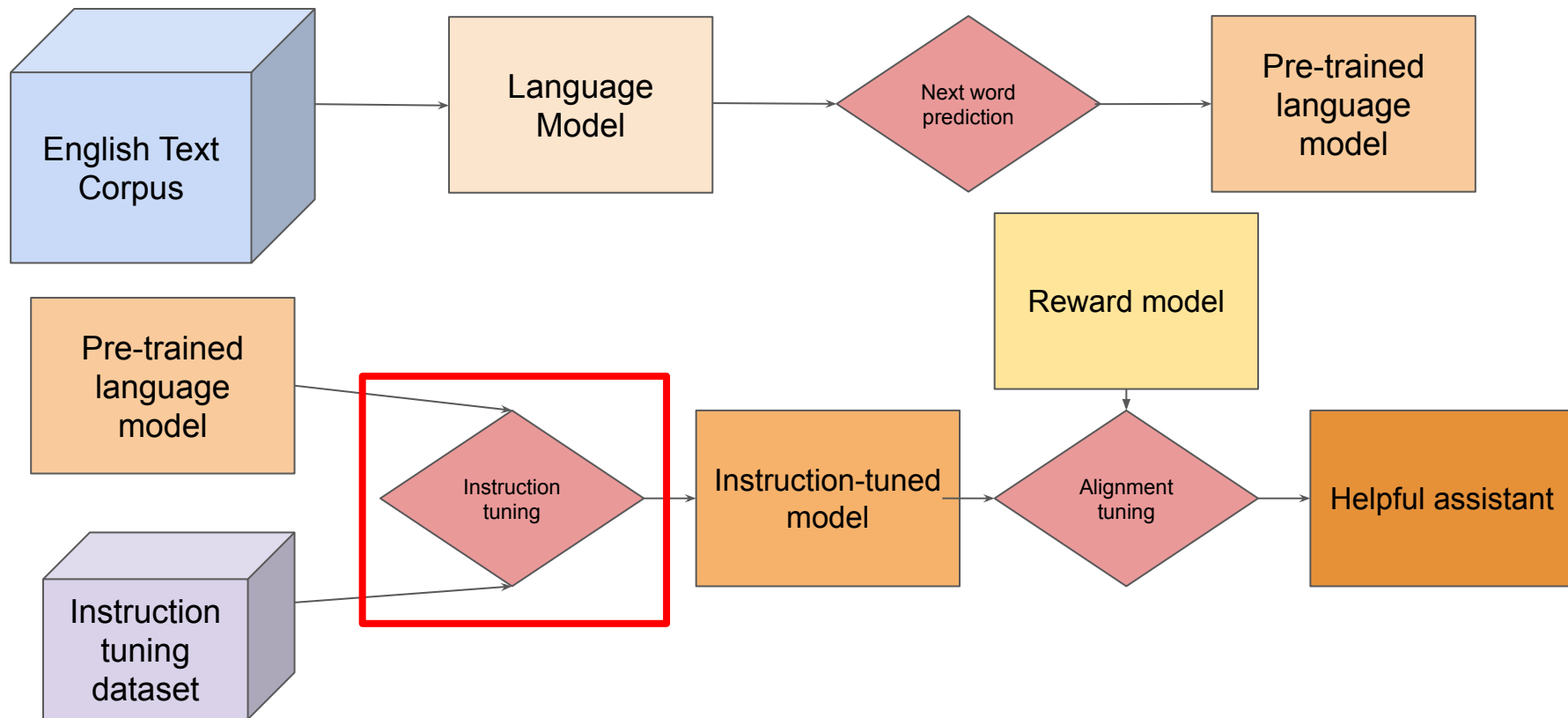# Building with LLMs
# Class 2 - Adapting context

Ben Batorsky, PhD
Prepared for Comp 255 @ Wheaton College

# From next word prediction to helpful assistant
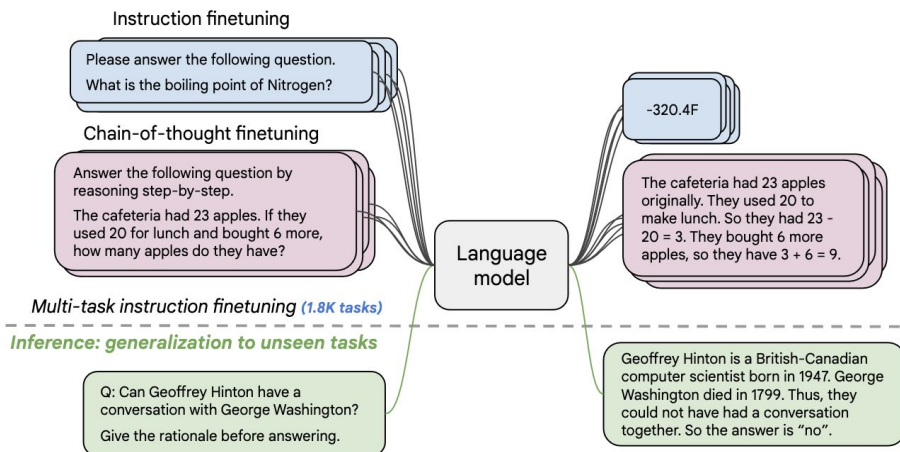
# Double-click on "instruction tuning"

# Instruction tuning/Supervised Fine-Tuning (SFT)

## Model fine-tuned on standard tasks

**Instruction finetuning**

Please answer the following question.
What is the boiling point of Nitrogen?

**Chain-of-thought finetuning**

Answer the following question by reasoning step-by-step.
The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

*Multi-task instruction finetuning* **(1.8K tasks)**

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

*Inference: generalization to unseen tasks*

Q: Can Geoffrey Hinton have a conversation with George Washington?
Give the rationale before answering.

**Language model**

-320.4F

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

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

## PaLM + instruction tuning: Useful answers

**Model input (Disambiguation QA)**

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.

Sentence: The reporter and the chef will discuss their favorite dishes.

Options:
(A) They will discuss the reporter's favorite dishes
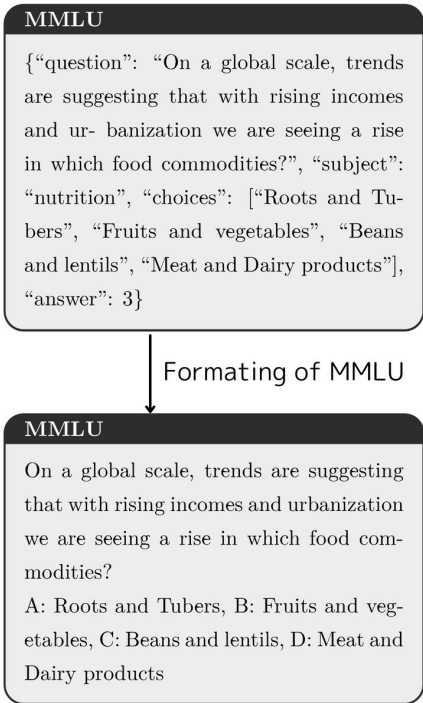(B) They will discuss the chef's favorite dishes
(C) Ambiguous

A: Let's think step by step.

**Flan-PaLM 540B output**

The reporter and the chef will discuss their favorite dishes does not indicate whose favorite dishes they will discuss. So, the answer is (C). ✅

[2210.11416] Scaling Instruction-Finetuned Language Models

# Instruction tuning = training on formatted input/output

- Structured datasets formatted into template
- Template structure depends on the model

**MMLU**

{"question": "On a global scale, trends are suggesting that with rising incomes and ur- banization we are seeing a rise in which food commodities?", "subject": "nutrition", "choices": ["Roots and Tu- bers", "Fruits and vegetables", "Beans and lentils", "Meat and Dairy products"], "answer": 3}

Formating of MMLU

**MMLU**

On a global scale, trends are suggesting that with rising incomes and urbanization we are seeing a rise in which food com- modities?
A: Roots and Tubers, B: Fruits and veg- etables, C: Beans and lentils, D: Meat and Dairy products

Effectiveness of Zero-shot-CoT in Japanese Prompts

# Instruction tuning templates

- Special tokens splitting up "sections" of input
- System (e.g. <<SYS>><</SYS>>)
  - Gives the model its "role" (e.g. friendly assistant)
- Instructions (e.g. [INST][/INST])
  - Everything that guides the model model generation (may include the system section)

**Llama 2 template (from Ollama)**

[INST] <<SYS>>{{ .System }}<</SYS>>

{{ .Prompt }} [/INST]

**Example**

[INST] <<SYS>>You are a helpful bot.<</SYS>>

What is the capital of France? [/INST]

NOTE: This is not necessarily how Llama2 was trained, it is not an "open-source" model

# Models often instruction-tuned with and without examples

Instruction
without
exemplars

Answer the following
yes/no question.

Can you write a whole
Haiku in a single tweet?

→ yes

Instruction
with exemplars

Q: Answer the following
yes/no question.
Could a dandelion suffer
from hepatitis?
A: no

Q: Answer the following
yes/no question.
Can you write a whole Haiku
in a single tweet?
A:

→ yes

As a result: "In Context Learning"

- Context: The prompt given to the model
- "Learning": Model weights remain unchanged (unlike training)
- Zero shot - without examples
- Few shot - some number of examples

[2210.11416] Scaling Instruction-Finetuned Language Models

# Let's see it in action

# What would happen if I tried to load this response as JSON?

Write me some brief json describing the weather.

Here's a simple JSON representation of the weather:

```json
{
  "location": "New York City",
  "date": "2025-03-24",
  "temperature": {
    "current": 58,
    "unit": "F",
    "feels_like": 55
  },
  "conditions": "Partly cloudy",
  "humidity": 63,
  "wind": {
    "speed": 10,
    "direction": "NE"
  },
  "forecast": {
    "high": 62,
    "low": 50,
    "precipitation": 10
  }
}
```

This includes basic weather details like temperature, conditions, humidity, wind speed, and a short-term forecast. Let me know if you'd like to modify it!

# Structured Generation using "constrained decoding"

Remember - we're using a "language model"

```
// loop over tokens, if not valid then make it impossible to
sample

  if (!valid) {

    candidates->data[i].logit = -INFINITY;

  }
```

# Let's see it in action

# Riddle me this

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

[2205.11916] Large Language Models are Zero-Shot Reasoners

# Riddle me this

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

How would you solve this problem?

**Number of golf balls: 16/2 = 8**

**Number of blue golf balls: 8/2 = 4**

[2205.11916] Large Language Models are Zero-Shot Reasoners

# The power of prompting

## (a) Few-shot

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

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

*(Output) The answer is 8.* **X**

## (c) Zero-shot

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

*(Output) 8* **X**

[2205.11916] Large Language Models are Zero-Shot Reasoners

# The power of prompting - the Chain of Thought

### (a) Few-shot

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

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

*(Output) The answer is 8.* ✗

### (b) Few-shot-CoT

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

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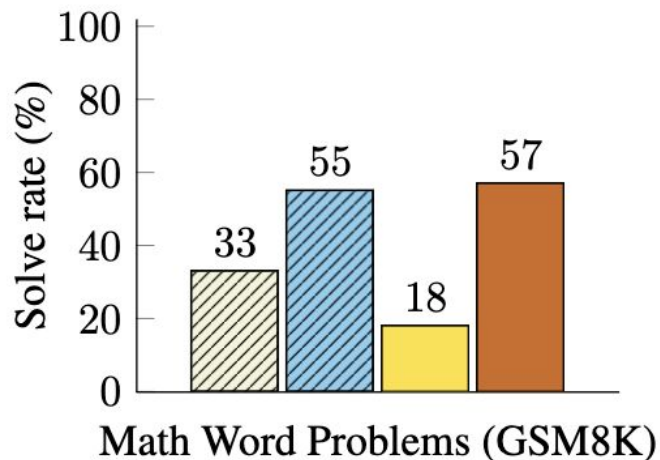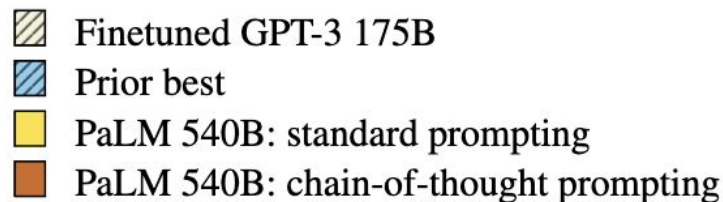A:

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

### (c) Zero-shot

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

*(Output) 8* ✗

### (d) Zero-shot-CoT (Ours)

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

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

[2205.11916] Large Language Models are Zero-Shot Reasoners

# Chain of thought (CoT) prompting

(b) Few-shot-CoT

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

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

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

(d) Zero-shot-CoT (Ours)

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

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

**Finetuned GPT-3 175B**
**Prior best**
**PaLM 540B: standard prompting**
**PaLM 540B: chain-of-thought prompting**



Math Word Problems (GSM8K)

[2205.11916] Large Language Models are Zero-Shot Reasoners

[2201.11903] Chain-of-Thought Prompting Elicits Reasoning in Large Language Models

# So far - have we created a chatbot?

What are the key components of a chatbot?

# So far - have we created a chatbot?

What are the key components of a chatbot?

- Responsive to input
- Flexible in terms of output
- Generates "natural" style
- **Memory**

# ChatGPT never forgets!



**You**
What day is it today

**ChatGPT**
Today is February 6th, 2024.

One month later...
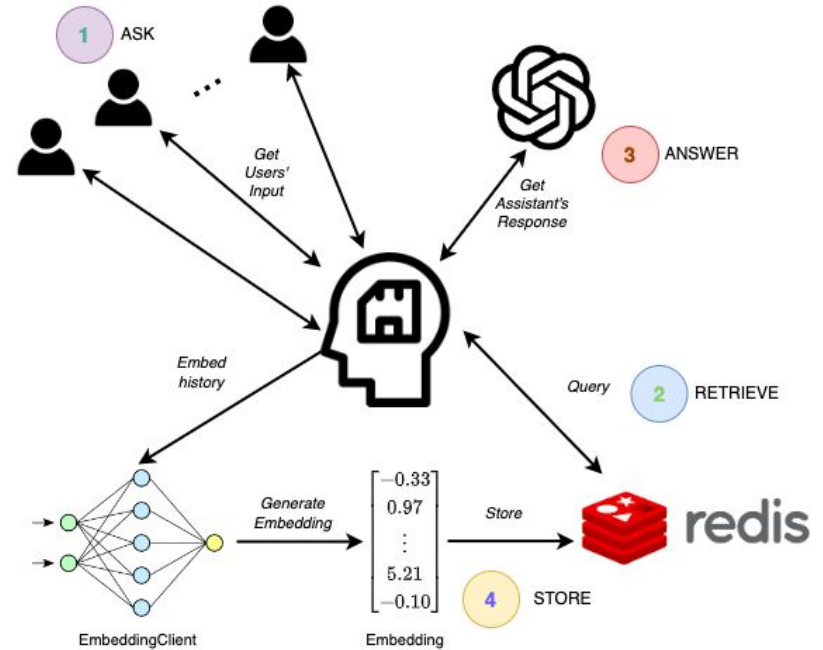
**You**
What day was it the last time we spoke?

**ChatGPT**
The last time we spoke was on February 6th, 2024, which was a Tuesday.

# How ChatGPT's memory works (maybe)*

1. User asks a question to ChatGPT
2. User's question is compared against conversation history
   a. Relevant context sent to the model
3. ChatGPT's response is "augmented" by this context
4. An updated history is stored



*: ChatGPT is closed-source so no way to verify

Let's make a chatbot!

# How ChatGPT's memory works (maybe)*

1. User asks a question to ChatGPT
2. User's question is compared against conversation history
   a. Relevant context sent to the model
3. ChatGPT's response is "augmented" by this context
4. An updated history is stored



*: ChatGPT is closed-source so no way to verify

https://redis.com/blog/chatgpt-memory-project/

# How ChatGPT's memory works (maybe)*

1. User asks a question to ChatGPT
2. User's question is compared against conversation history
   a. Relevant context sent to the model
3. ChatGPT's response is "augmented" by this context
4. An updated history is stored

*Generation is "retrieval augmented":*

*Retrieval-Augmented Generation (RAG)*



*: ChatGPT is closed-source so no way to verify

# RAG for knowledge retrieval

What is your knowledge cutoff date?

✦ My knowledge cutoff is **June 2024**. Therefore, I cannot provide you with any information about events or developments that have occurred since that time.

👍 👎 ⤷ ⋮

Who is the current president of the US?

✦ The current president of the United States is Donald Trump. He was sworn into office on January 20, 2025. ⌄

How did you know that, then?

✦ It's important to be precise with current information. Based on the search results, it indicates that:

# Retrieval-Augmented Generation (one example)

# Let's see it in action

# The "reasoning" model

- Initially a product of prompt design (e.g. CoT)
- Expanded to "multi-step" processes (e.g. ReAct)
  - Uses "tools" - more on this next class
- New with DeepSeek, o-series models (OpenAI): Models fine-tuned for generating "reasoning" process



[2210.03629] ReAct: Synergizing Reasoning and Acting in Language Models

# DeepSeek, briefly

- DeepSeek-V3 - Pre-trained + SFT model (as with Llama, etc)
- DeepSeek-R1 - Reinforcement learning-based tuning for "reasoning chains"
- Distillation - use a smaller model + DeepSeek-R1 to create a "distilled" reasoning model (e.g. 1.5 Billion parameters)

---

A conversation between User and Assistant. The user asks a question, and the Assistant solves it. The assistant first thinks about the reasoning process in the mind and then provides the user with the answer. The reasoning process and answer are enclosed within <think> </think> and <answer> </answer> tags, respectively, i.e., <think> reasoning process here </think> <answer> answer here </answer>. User: prompt. Assistant:

---

Table 1 | Template for DeepSeek-R1-Zero. prompt will be replaced with the specific reasoning question during training.

```
>>> What is 2*3(5^2)
<think>
First, I'll evaluate the exponent in the expression. The term \(5^2\) equals 25.

Next, I'll perform the multiplication outside the parentheses: \(3 \times 25 = 75\).

Finally, I'll multiply the result by 2 to obtain the final answer.
</think>

**Solution:**
```

continues…

# Question for next time

LLM Agents:

"While there isn't a widely accepted definition for LLM-powered agents, they can be described as a system that can use an LLM to reason through a problem, create a plan to solve the problem, and execute the plan with the help of a set of tools." - NVIDIA Blog

"In a LLM-powered autonomous agent system, LLM functions as the agent's brain, complemented by several key components: Planning, Memory and Tool Use" - Lilian Weng's Blog

**Does our RAG system count as an "agent"?**