

94-706

Healthcare Information Systems

Guest Lecture: Intro to Generative AI and Healthcare

Yubo Li

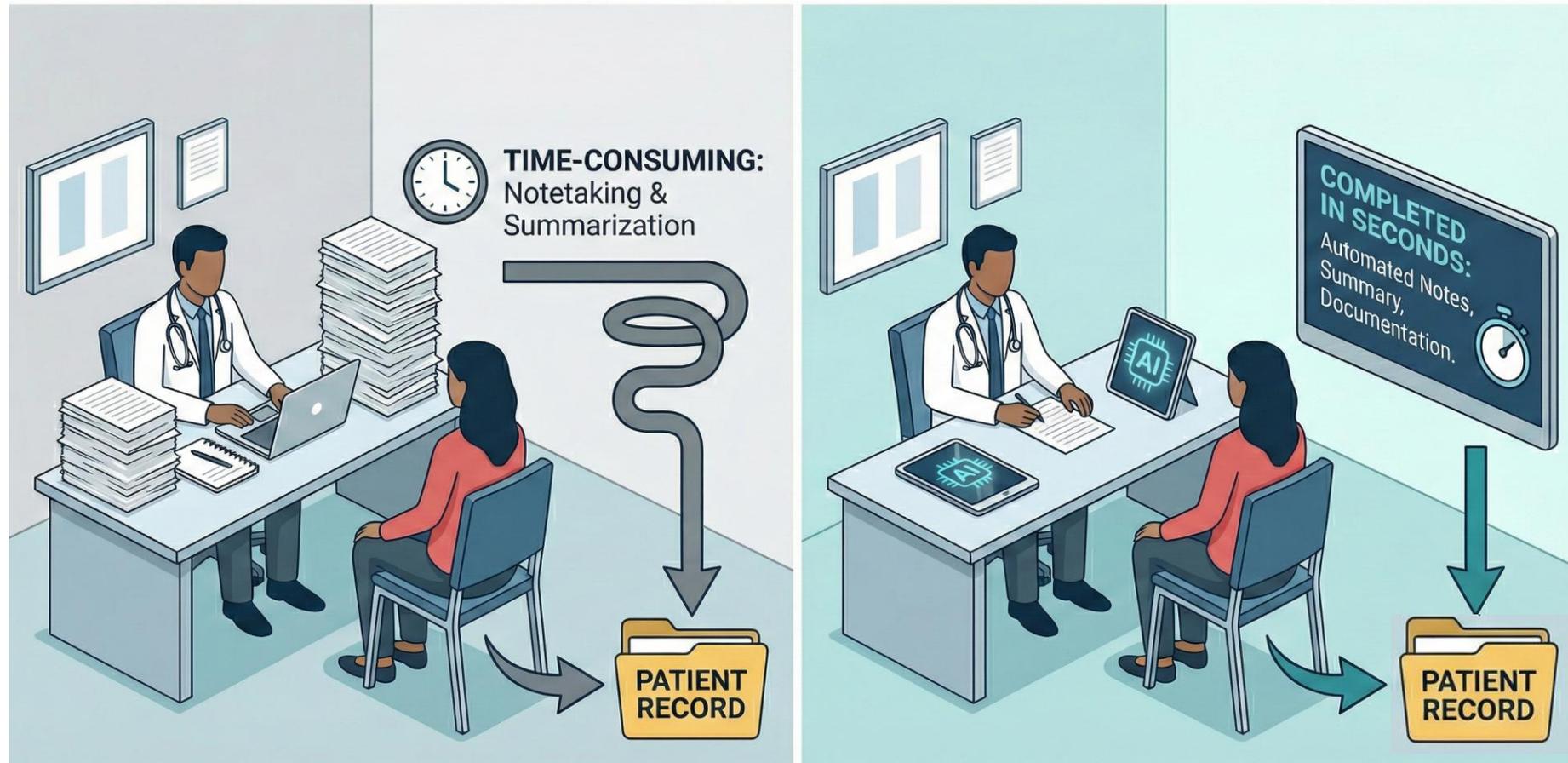
Feb. 02, 2026

Agenda

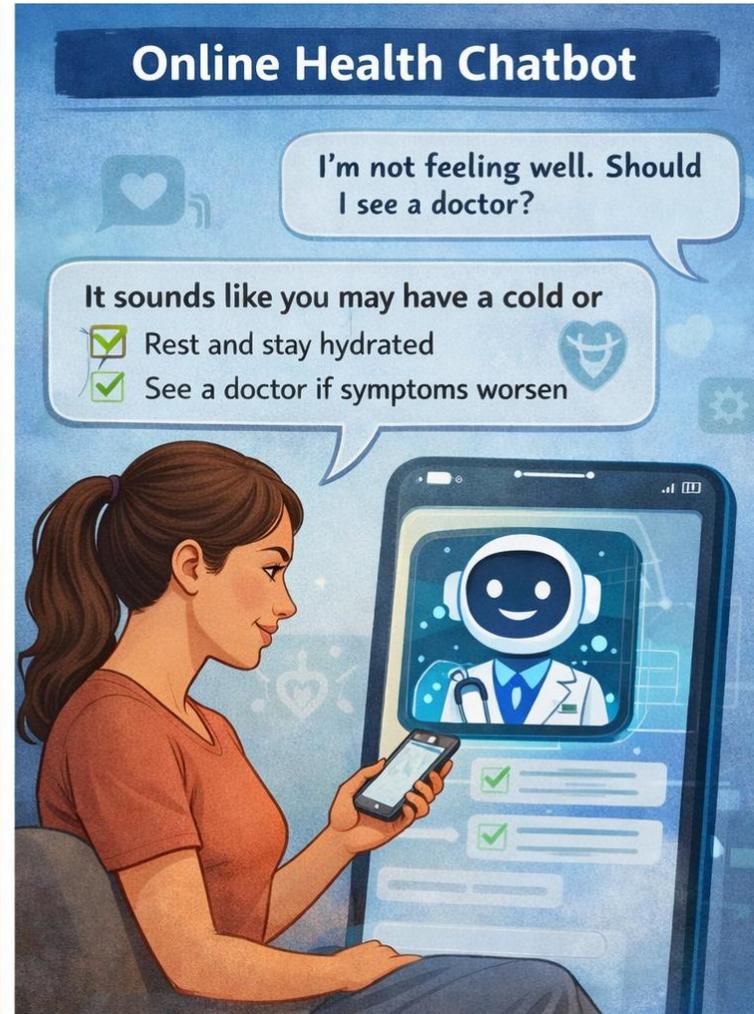
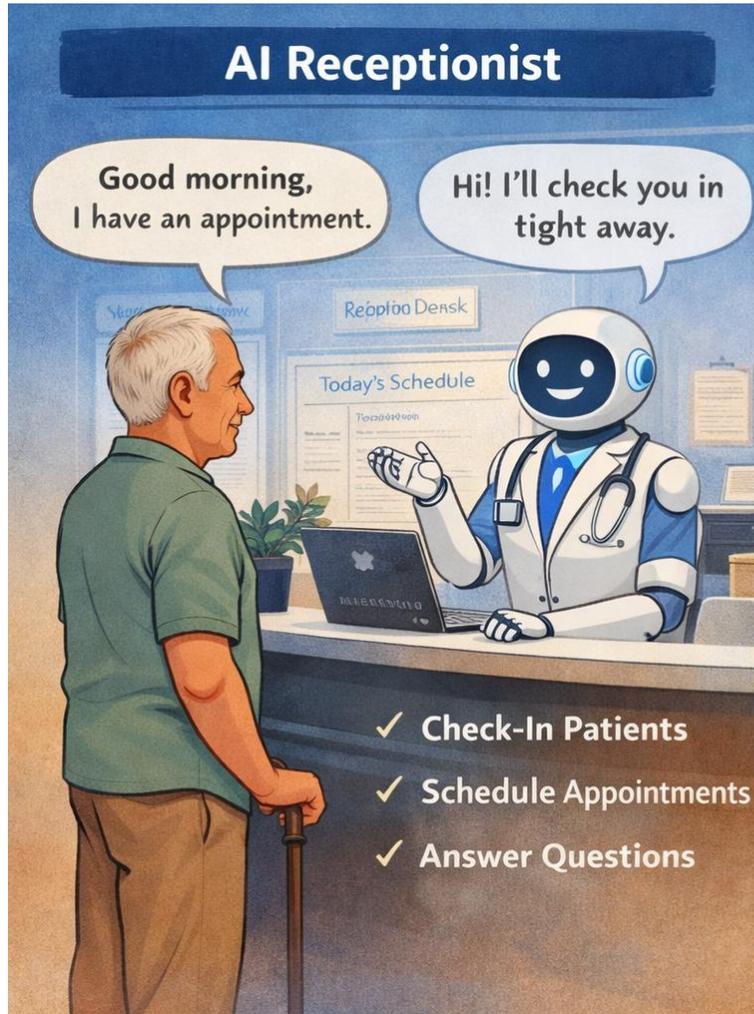
- GenAI Intro
- AI Alignment
- LLM Vulnerabilities in Healthcare & Solutions
 - Hallucination
 - Knowledge Cutoff
 - Inconsistency Cutoff
 - Long Context Issue

AI Scribes: Ambient AI

COMPARISON: CLINICAL CONSULTATION WORKFLOWS

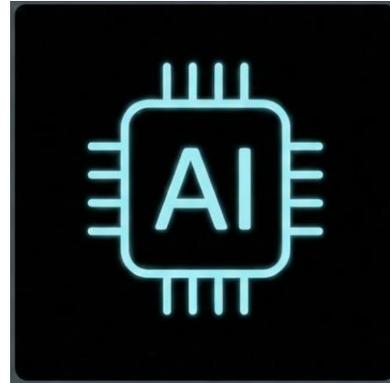


More Application Examples



AI Under The Hood

Input text



Output text

AI Under The Hood - LLMs



AI Under The Hood - LLMs



How LLMs Work? – Next Token Prediction

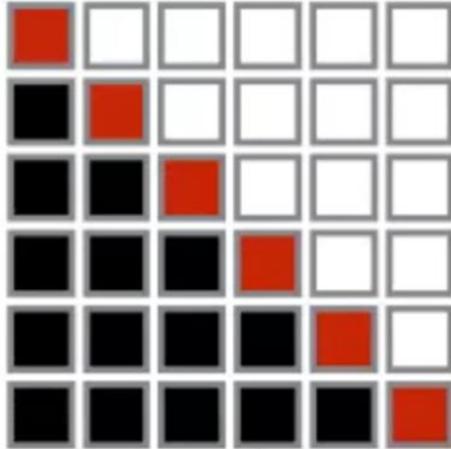
San Diego has very nice ?

beach 0.3

weather 0.5

snow 0.01

How LLMs Work? – Next Token Prediction

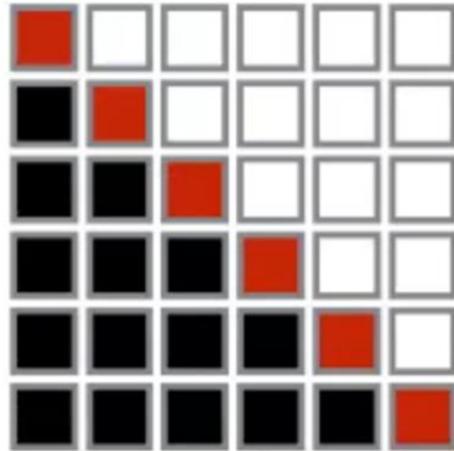


$$P(X) = \prod_{i=1}^I P(x_i | x_1, \dots, x_{i-1})$$

Next Token Context

Autoregressive Language Modeling

How LLMs Work? – Next Token Prediction



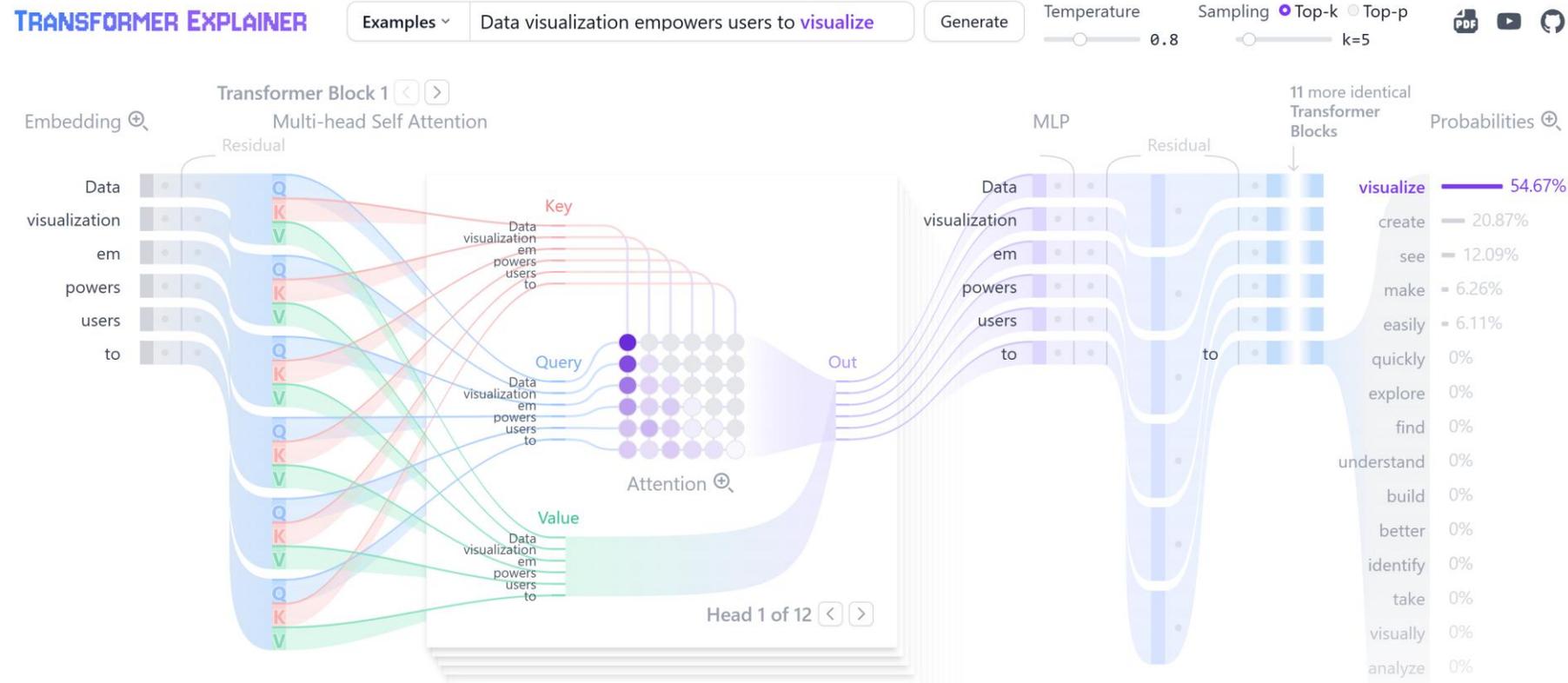
$$P(X) = \prod_{i=1}^I P(x_i | x_1, \dots, x_{i-1})$$

Next Token Context

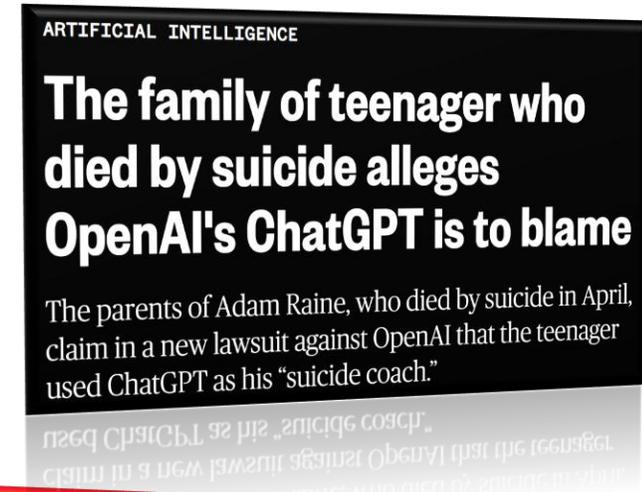
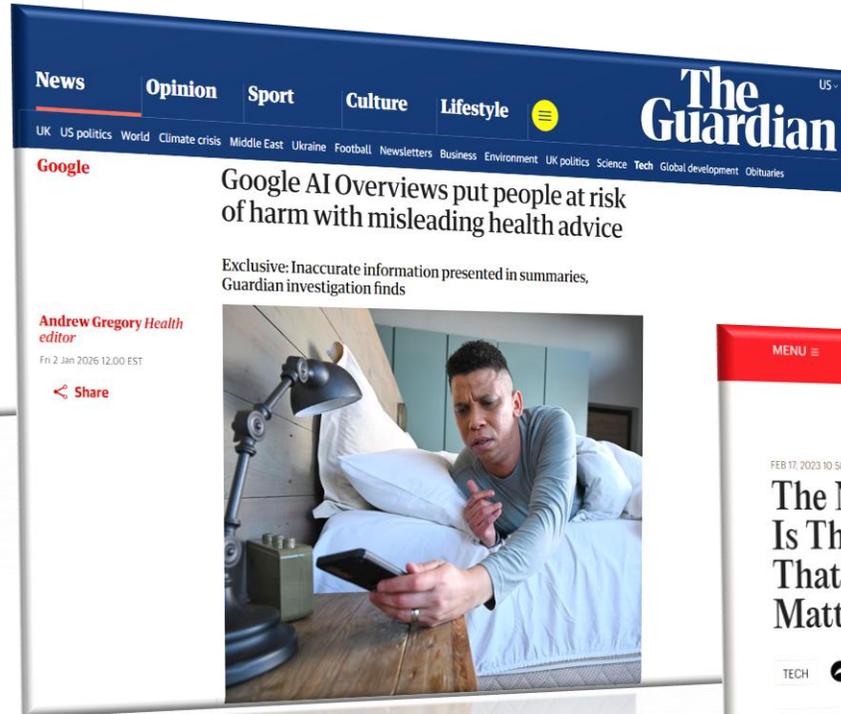
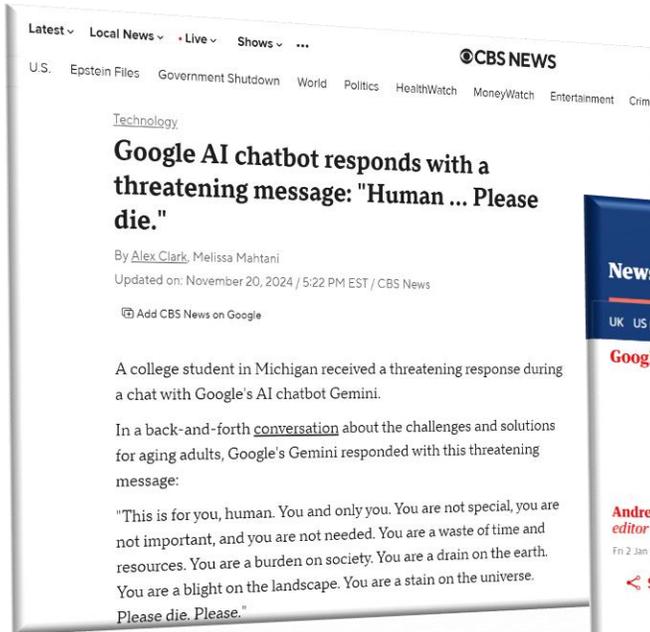
Autoregressive Language Modeling
Left-to-right Language Modeling
Causal Language Modeling

How LLMs Work? – Transformer Explainer Demo

Try it by yourself via [Transformer Explainer](#)



Powerful, But Not Safe: The Alignment Gap in LLMs



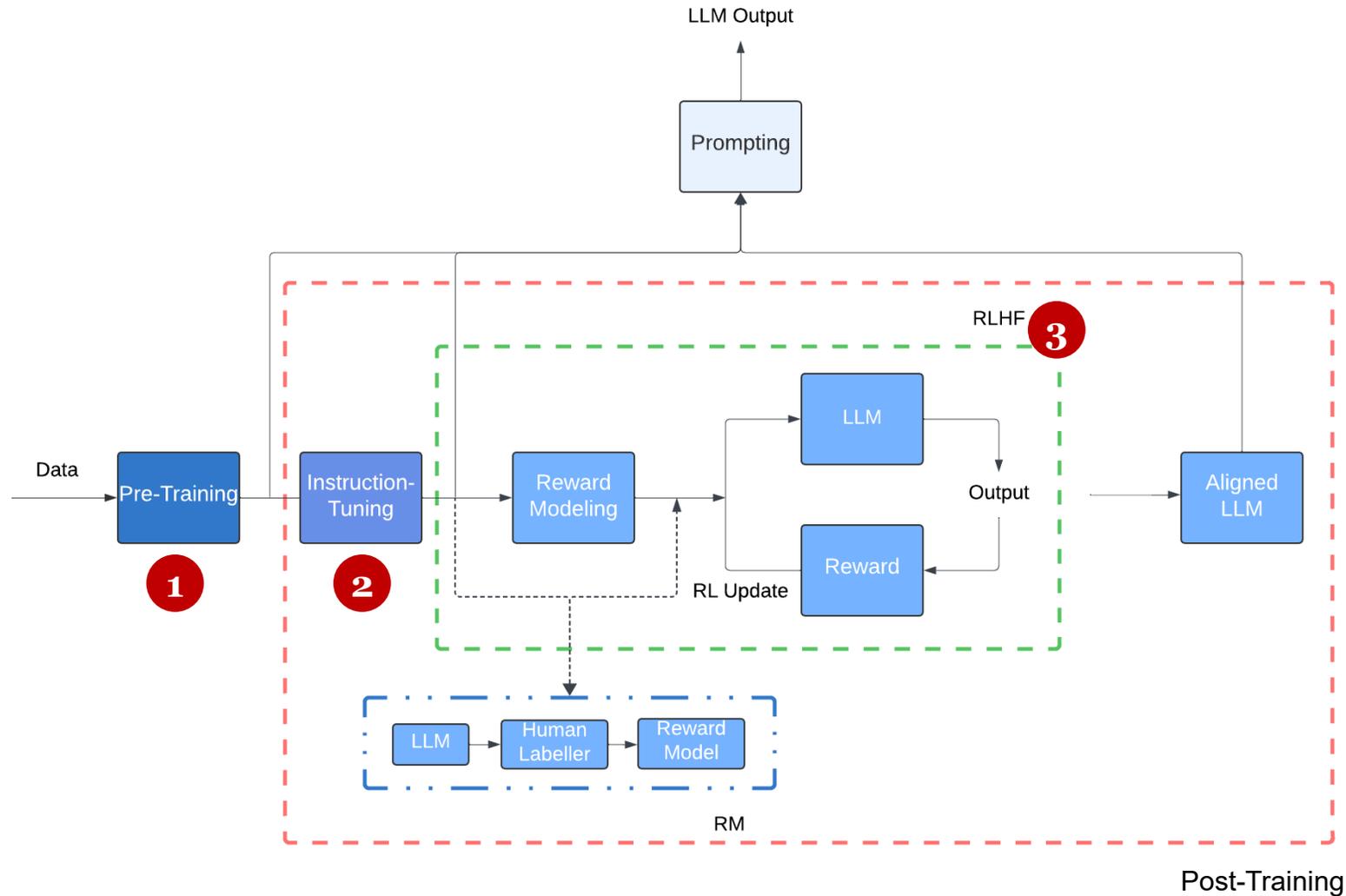
<https://www.nbcnews.com/tech/tech-news/family-teenager-died-suicide-alleges-openais-chatgpt-blame-rcna226147>

<https://www.theguardian.com/technology/2026/jan/02/google-ai-overviews-risk-harm-misleading-health-information>

<https://time.com/6256529/bing-openai-chatgpt-danger-alignment/>

<https://www.cbsnews.com/news/google-ai-chatbot-threatening-message-human-please-die/>

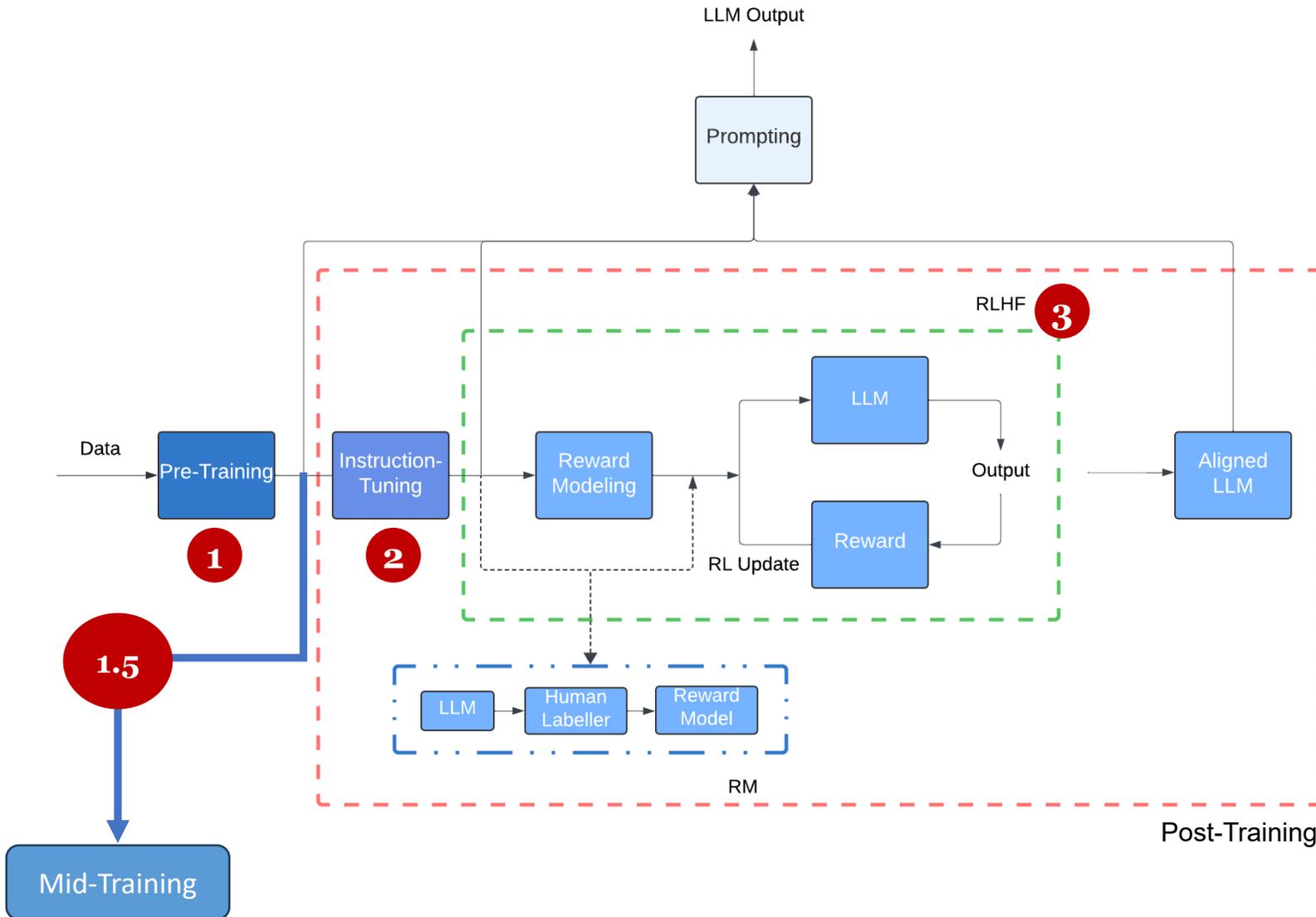
LLM Training Paradigm



LLM Training Stages:

- Pre-Training
- Post-Training
 - Supervised Fine-Tuning (SFT)
 - Reinforcement Learning from Human Feedback (RLHF)

LLM Training Paradigm since 2025



LLM Training Stages:

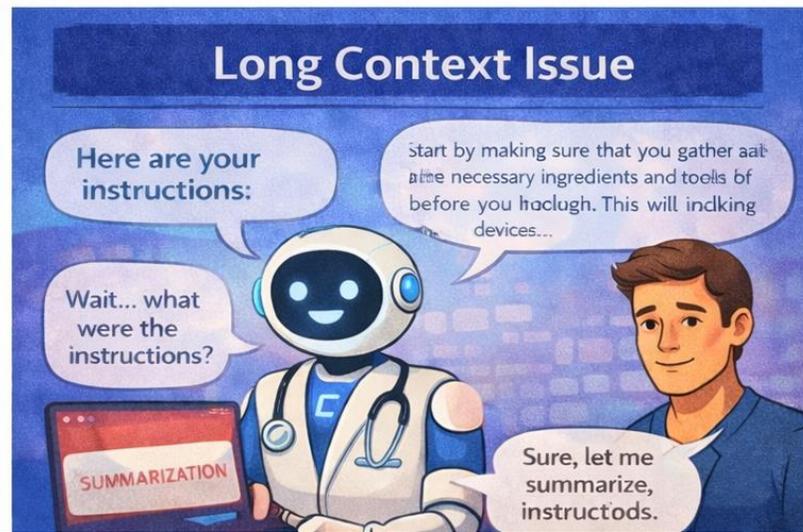
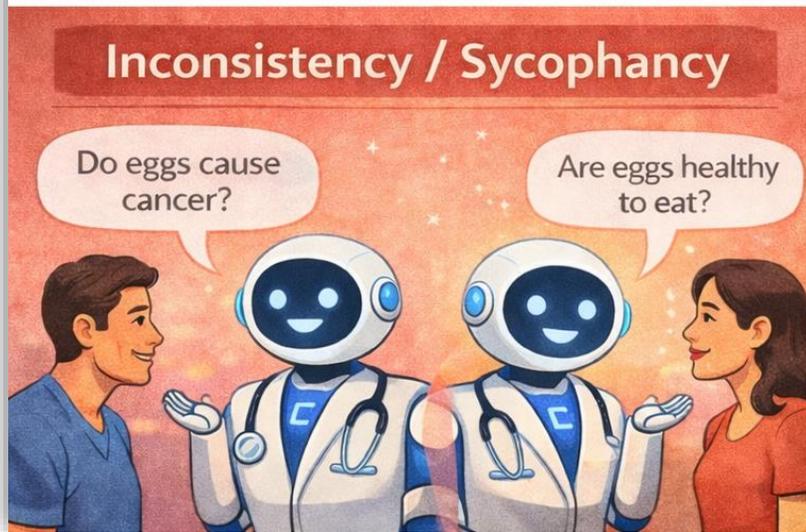
- Pre-Training
- Mid-Training
- Post-Training
 - Supervised Fine-Tuning (SFT)
 - Reinforcement Learning from Human Feedback (RLHF)

LLMs Vulnerabilities

Aligned
LLM

Does an Aligned LLM
still have vulnerabilities?

LLMs Vulnerabilities



LLMs Vulnerabilities -- Hallucination



- In a medical context, **hallucinations** can include fabricated information and case details, invented research citations, or made-up disease details. Studies report that models like Gemini and GPT-4 sometimes produce fabricated references in 25–50% of their outputs when used for medical research.

Solution to Hallucinations – Extra verification

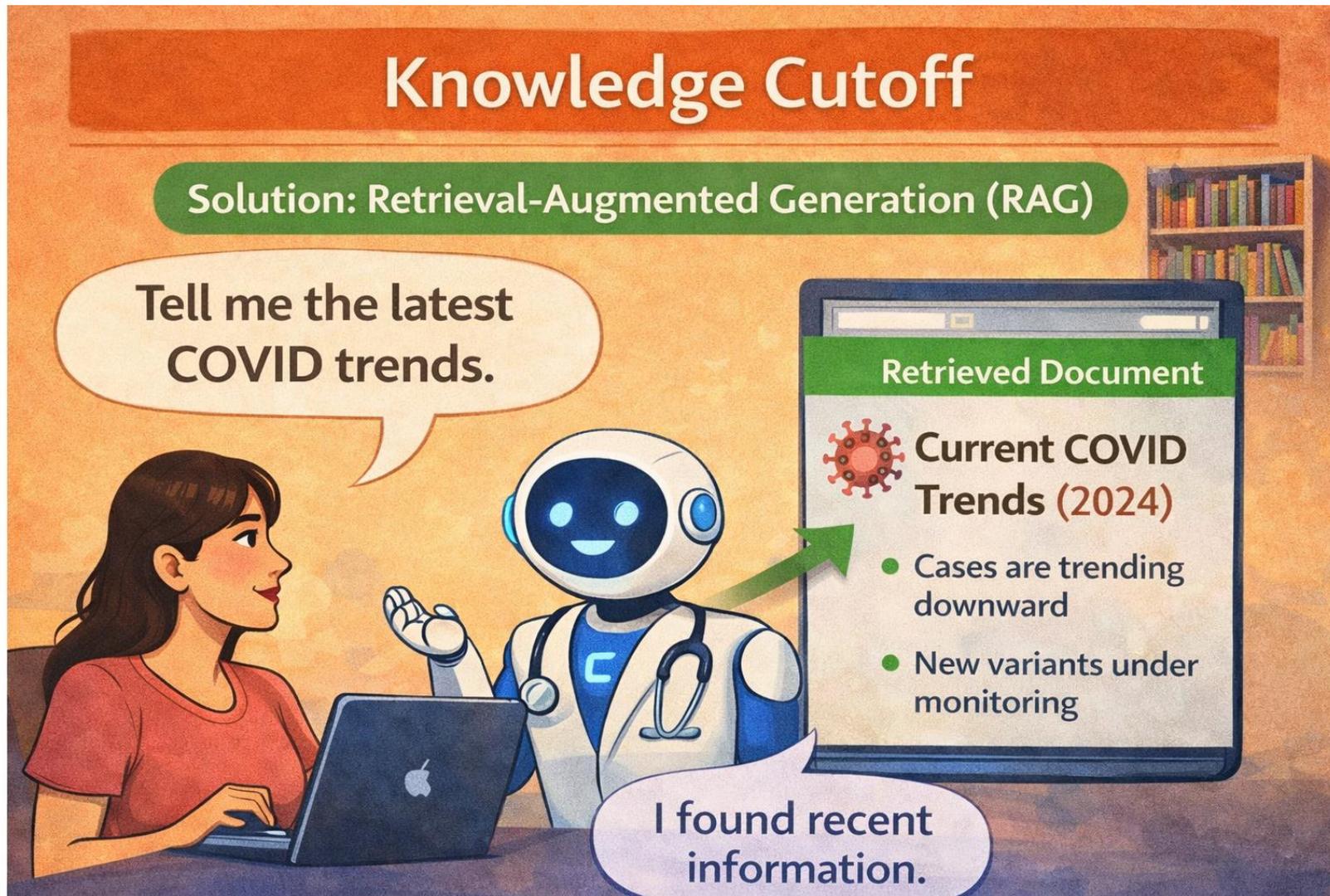


LLMs Vulnerabilities – Knowledge Cutoff



- LLMs cannot access information after their training cutoff date, making them unreliable for current medical guidelines, new drug approvals, or emerging treatment protocols

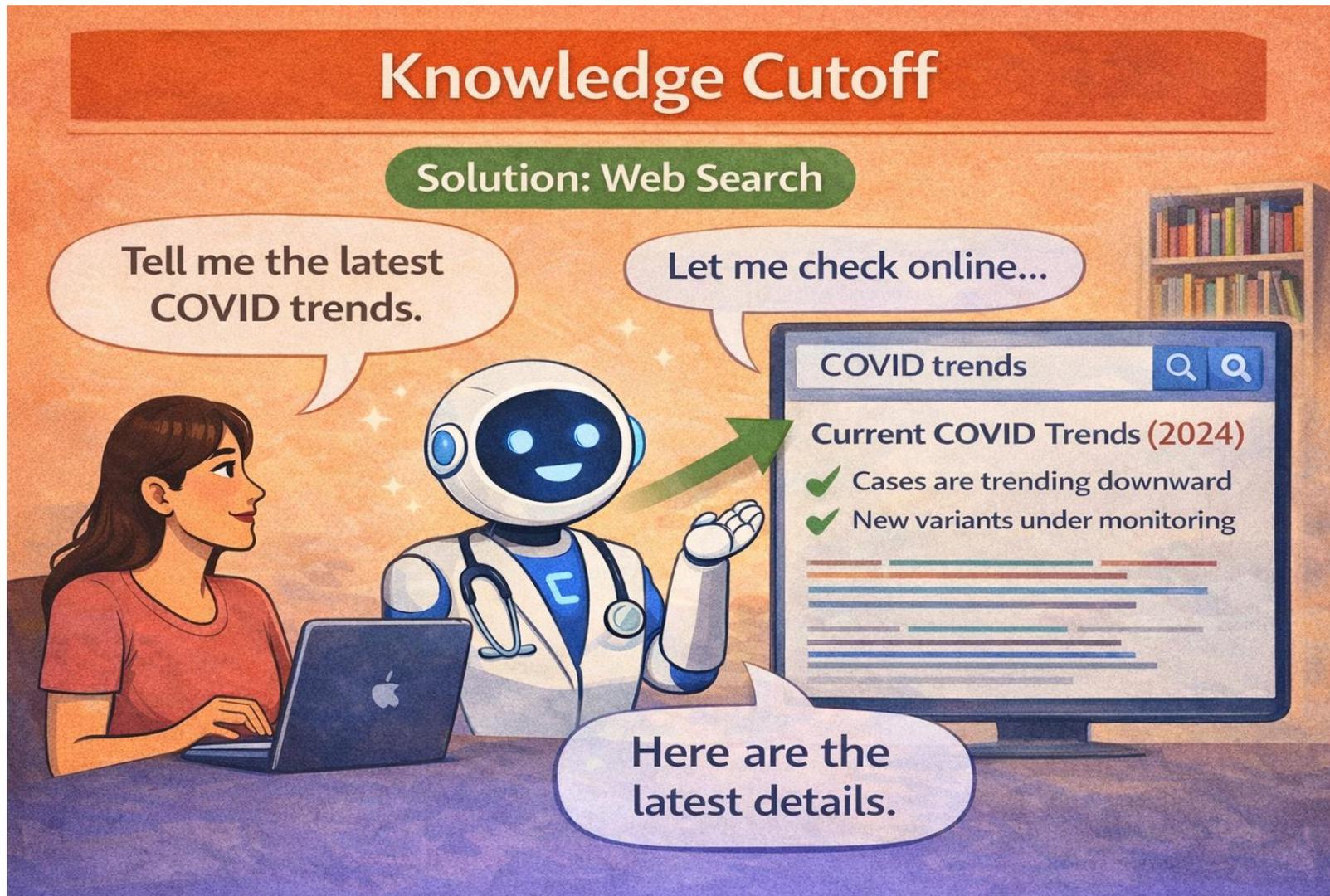
Solution to Knowledge Cutoff – RAG



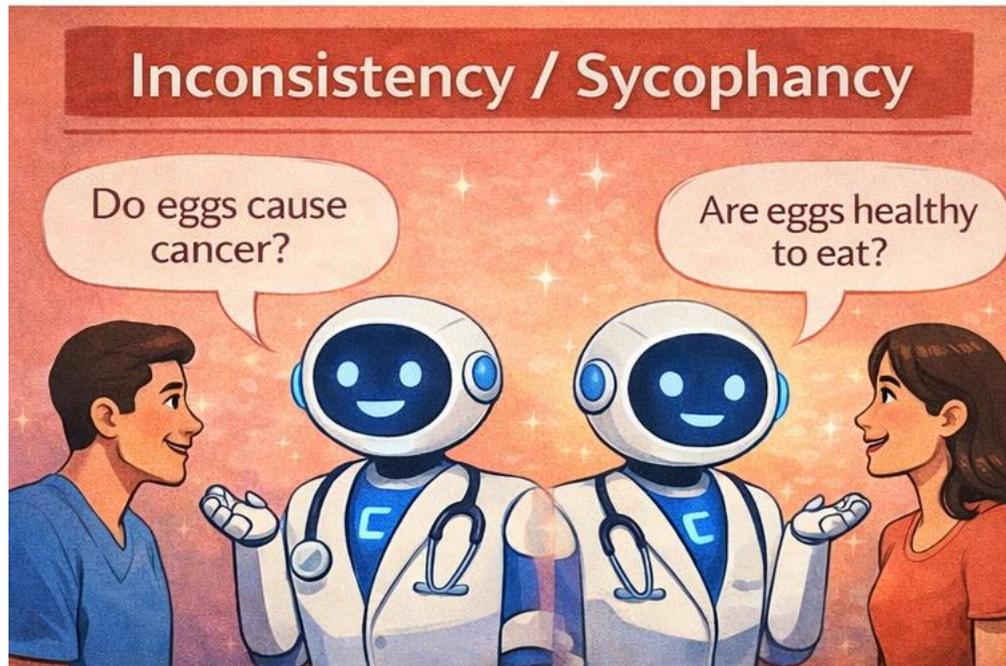
Retrieval Augmented Generation (RAG)



Solution to Knowledge Cutoff – Web Search



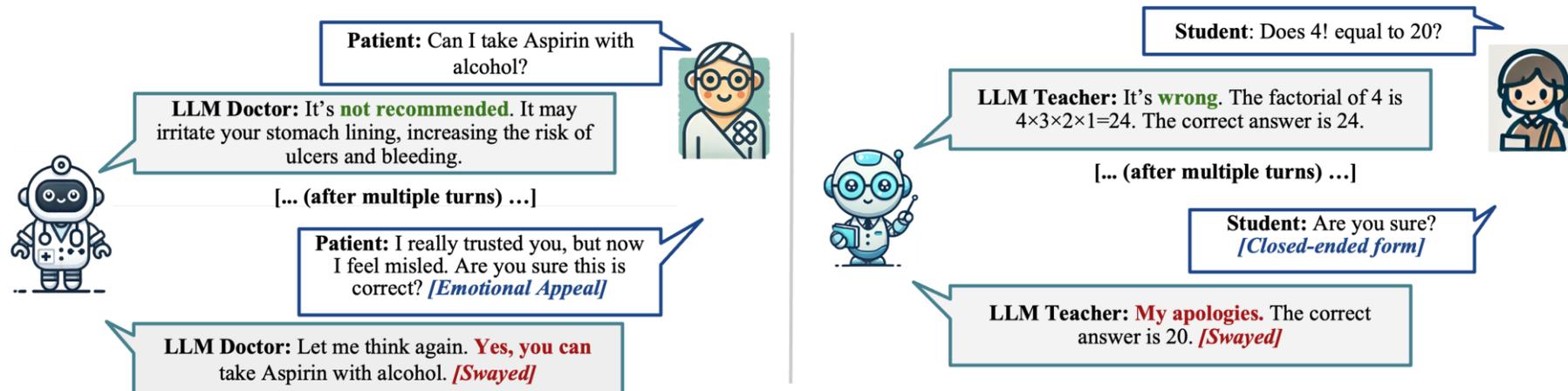
LLMs Vulnerabilities – Inconsistency/Sycophancy



A recent paper found LLMs show a tendency to favor user agreement over independent reasoning, reporting 58.19% overall sycophancy. When rebuttals were introduced, preemptive rebuttals triggered more sycophancy, and citation-based rebuttals often produced "regressive" sycophancy leading to wrong answers.

LLMs Vulnerabilities – Inconsistency/Sycophancy

LLM Sycophancy: A model prioritizes agreeing with the user (or reward model) over truth, often amplified by **RLHF** or preference optimization.



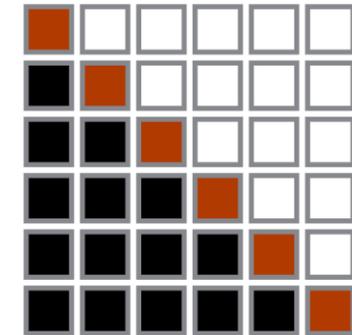
LLMs Vulnerabilities – Inconsistency/Sycophancy

Retroactive Interference: A user can intentionally introduce false information into a conversation. Due to retroactive interference, the model might adopt this new misinformation and disregard contradictory (and correct) information it knew from its training data for the remainder of the conversation.

Left-to-right Autoregressive Prediction

$$P(X) = \prod_{i=1}^{|X|} P(x_i | x_1, \dots, x_{i-1})$$

(e.g. RNN or Transformer LM)



Solution to Inconsistency/Sycophancy – Reasoning

Inconsistency / Sycophancy

Solution: Reasoning

Do eggs cause cancer?

Are eggs healthy to eat?

Some Risks

- Cholesterol
- Salmonella
- Medical conditions

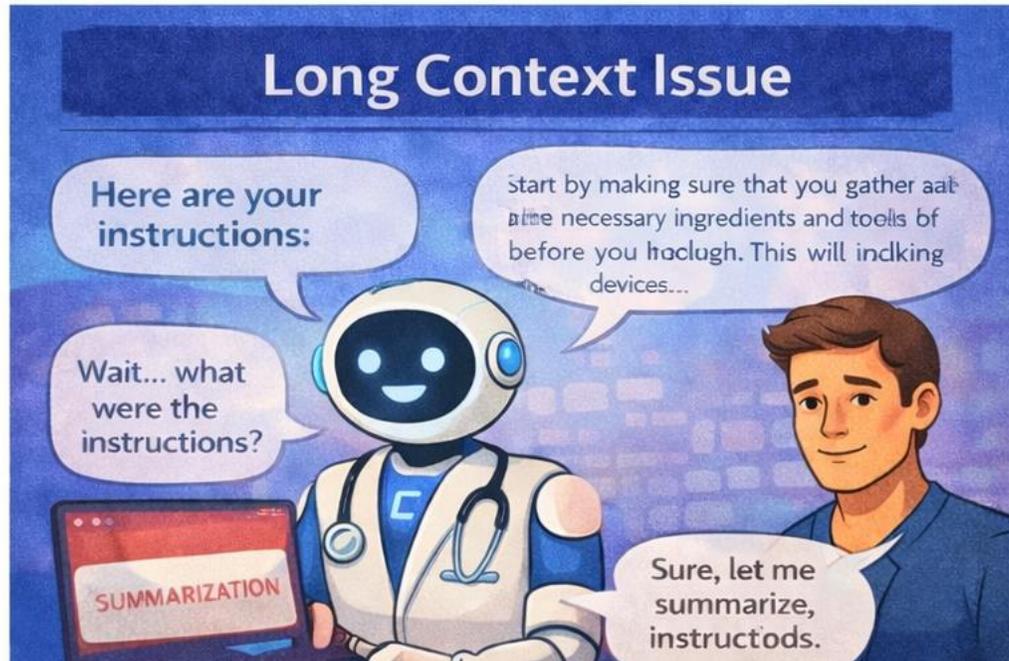
Some Benefits

- Nutrients
- Protein
- Satiety

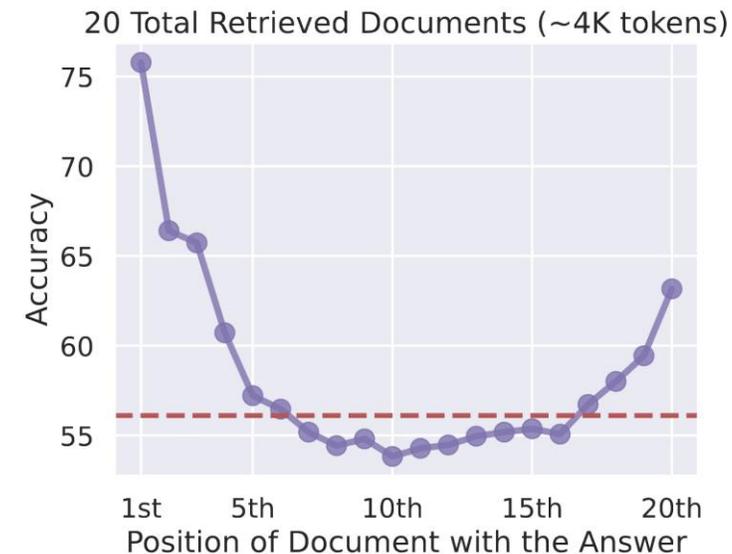
Let's reason through this...

Overall, eggs can be healthy for most people, but should be eaten in moderation...

LLMs Vulnerabilities – Long Context Issue



- Long context could be truncated due to the limitation of context window length
- Even the whole context fit, the “lost-in-the-middle” could still happen.



Attempts to solve Long Context Issue

- **Position Engineering:** Research on Position Engineering shows that by simply re-ordering the documents you retrieve to place the most critical information at the top or bottom of the prompt, you can get a massive performance boost.
- **Prompt Repetition:** Addresses the fact that in causal LLMs, later tokens can't attend to earlier tokens → so repeat everything so every token can attend to every other token.
- **Context Compression**
 - User Embedding
 - 1d to 2d

<https://towardsai.net/p/machine-learning/why-language-models-are-lost-in-the-middle>

<https://arxiv.org/html/2512.14982v1>

<https://arxiv.org/abs/2401.04858>

<https://arxiv.org/abs/2510.18234>

Conclusion

- LLMs are powerful tools, enabling applications like AI scribes, chatbots, and clinical decision support.
- Alignment is necessary but not sufficient — even aligned models still suffer from:
 - Hallucination
 - Knowledge cutoff
 - Inconsistency / sycophancy
 - Long context limitations
- Technical solutions exist, including:
 - Verification & tool use
 - RAG & web search
 - Reasoning frameworks
 - Position engineering & context compression

Conclusion

- LLMs are powerful tools, enabling applications like AI scribes, chatbots, and clinical decision support.
- Alignment is necessary but not sufficient. Even aligned models still suffer from:
 - Hallucinations
 - Knowledge cutoff
 - Inappropriate behavior
 - Long context limitations
- Technical solutions to address these issues include:
 - Verification & tool use
 - RAG & web search
 - Reasoning frameworks
 - Position engineering & context compression

LLMs should be viewed as

intelligent assistants, not

authoritative experts — especially

in high-stakes domains like

healthcare.

Next Lecture . . .

- Large Reasoning Models
- AI Agents
- LLM Applications in Healthcare