

# STAT 992: Foundation Models for Biomedical Data

---

Ben Lengerich

Lecture 11: Prompts, In-Context Learning, Test-Time Adaptation

March 2, 2026



# Prompting





# Few-Shot / Zero-shot learning

One key emergent ability in GPT-2 is **zero-shot learning**: the ability to do many tasks with **no examples**, and **no gradient updates**, by simply:

- Specifying the right sequence prediction problem (e.g. question answering):

Passage: Tom Brady... Q: Where was Tom Brady born? A: ...

- Comparing probabilities of sequences (e.g. Winograd Schema Challenge [[Levesque, 2011](#)]):

The cat couldn't fit into the hat because it was too big.  
**Does it = the cat or the hat?**

$\equiv$  Is  $P(\dots\text{because } \mathbf{the\ cat} \text{ was too big}) \geq$   
 $P(\dots\text{because } \mathbf{the\ hat} \text{ was too big})?$

[[Radford et al., 2019](#)]



# Few-Shot / Zero-shot learning

GPT-2 beats SoTA on language modeling benchmarks with **no task-specific fine-tuning**

*Context:* “Why?” “I would have thought you’d find him rather dry,” she said. “I don’t know about that,” said Gabriel.  
“He was a great craftsman,” said Heather. “That he was,” said Flannery.

*Target sentence:* “And Polish, to boot,” said ----- **LAMBADA** (language modeling w/ long discourse dependencies)

*Target word:* Gabriel

[[Paperno et al., 2016](#)]

	LAMBADA (PPL)	LAMBADA (ACC)	CBT-CN (ACC)	CBT-NE (ACC)	WikiText2 (PPL)
SOTA	99.8	59.23	85.7	82.3	39.14
117M	<b>35.13</b>	45.99	<b>87.65</b>	<b>83.4</b>	<b>29.41</b>
345M	<b>15.60</b>	55.48	<b>92.35</b>	<b>87.1</b>	<b>22.76</b>
762M	<b>10.87</b>	<b>60.12</b>	<b>93.45</b>	<b>88.0</b>	<b>19.93</b>
1542M	<b>8.63</b>	<b>63.24</b>	<b>93.30</b>	<b>89.05</b>	<b>18.34</b>

[[Radford et al., 2019](#)]



# Few-Shot / Zero-shot learning

You can get interesting zero-shot behavior if you're creative enough with how you specify your task!

Summarization on CNN/DailyMail dataset [[See et al., 2017](#)]:

		ROUGE			
		R-1	R-2	R-L	
SAN FRANCISCO, California (CNN) -- A magnitude 4.2 earthquake shook the San Francisco ... overturn unstable objects.	<b>2018 SoTA</b>	<b>Bottom-Up Sum</b>	<b>41.22</b>	<b>18.68</b>	<b>38.34</b>
		Lede-3	40.38	17.66	36.62
	<b>Supervised (287K)</b>	Seq2Seq + Attn	31.33	11.81	28.83
		GPT-2 TL;DR:	29.34	8.27	26.58
	<b>TL;DR: Select from article</b>	Random-3	28.78	8.63	25.52

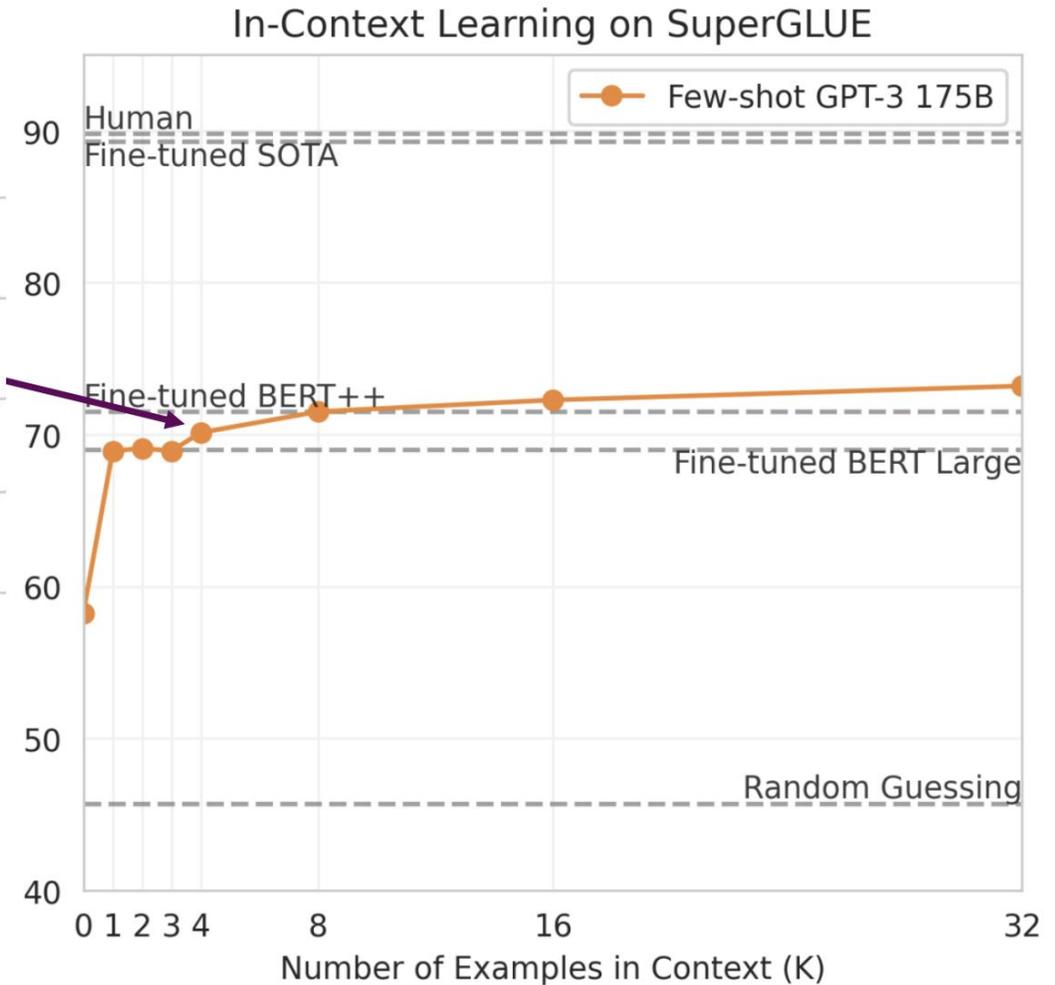
“Too Long, Didn’t Read”  
“Prompting”?

[[Radford et al., 2019](#)]

# “In-Context Learning”

## Few-shot

1 Translate English to French:  
2 sea otter => loutre de mer  
3 peppermint => menthe poivrée  
4 plush girafe => girafe peluche  
5 cheese => .....



[Brown et al., 2020]

# Chain-of-Thought

## Standard Prompting

### Model Input

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: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

### Model Output

A: The answer is 27. ❌

## Chain-of-Thought Prompting

### Model Input

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: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

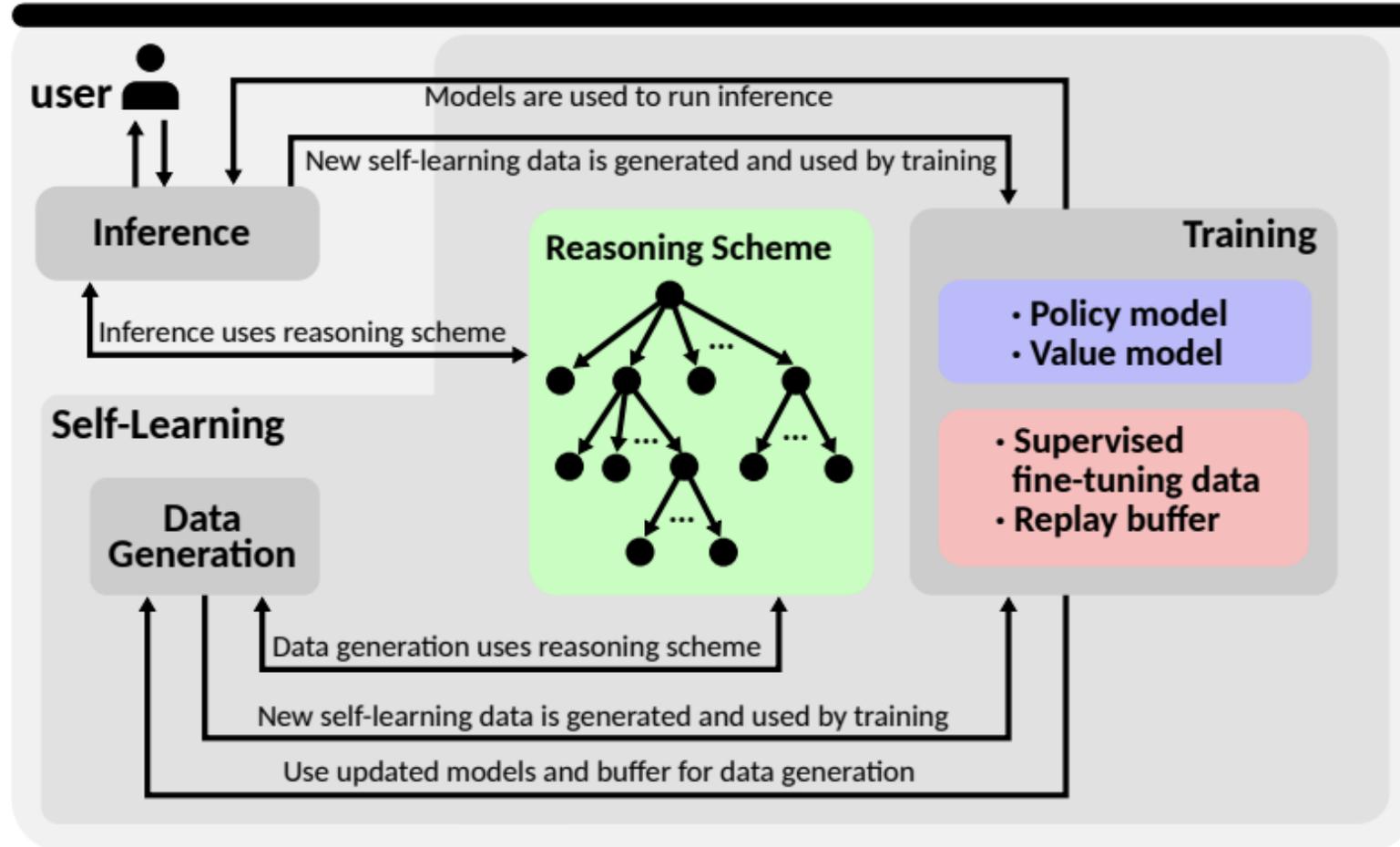
### Model Output

A: 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  $3 + 6 = 9$ . The answer is 9. ✓

Wei, et al. (2023) Chain-of-Thought Prompting Elicits Reasoning in LLMs

# Reasoning Models

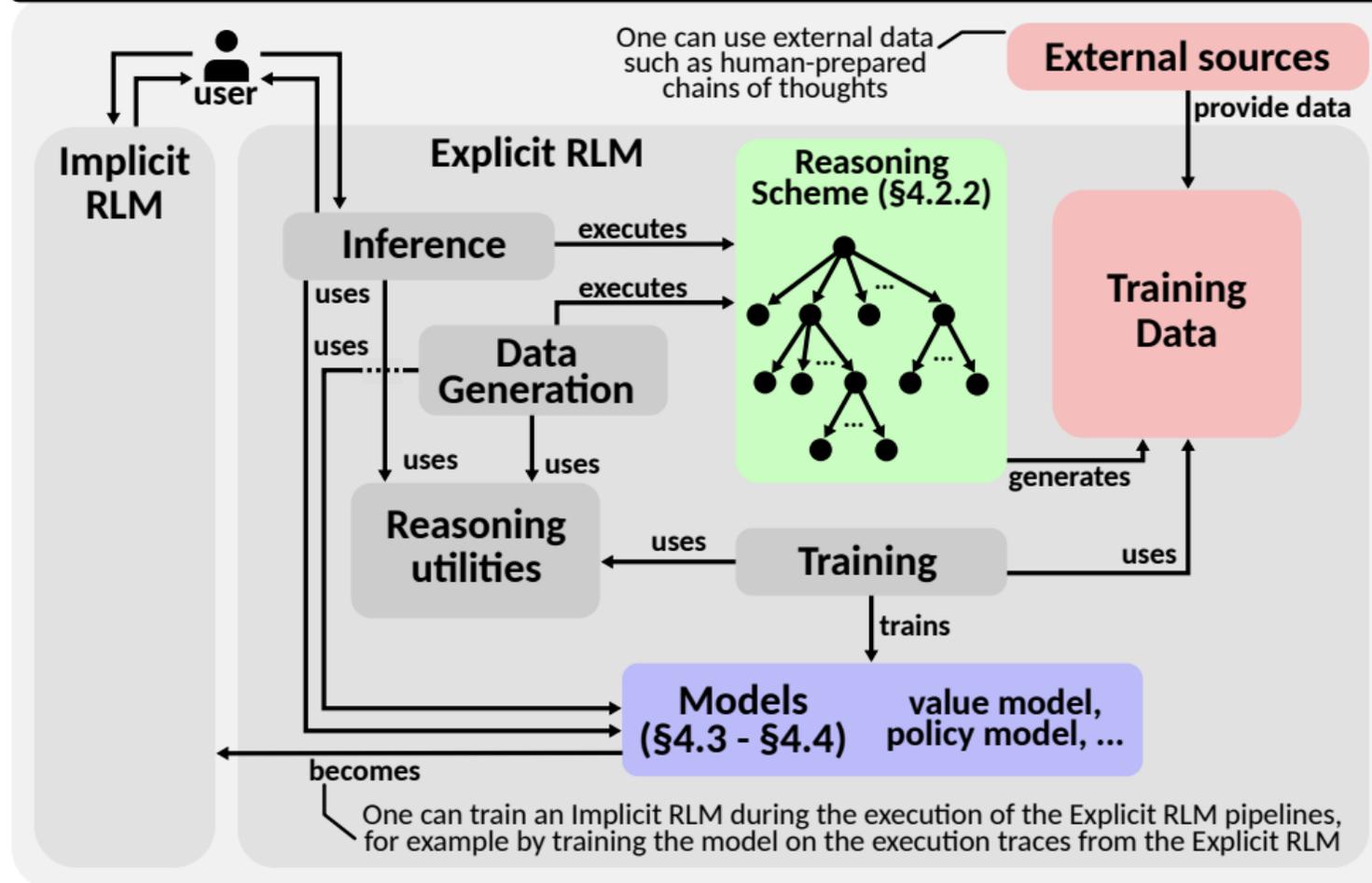
## High-level overview (§3.1)



# Reasoning Models

## Medium-level overview (§3.1)

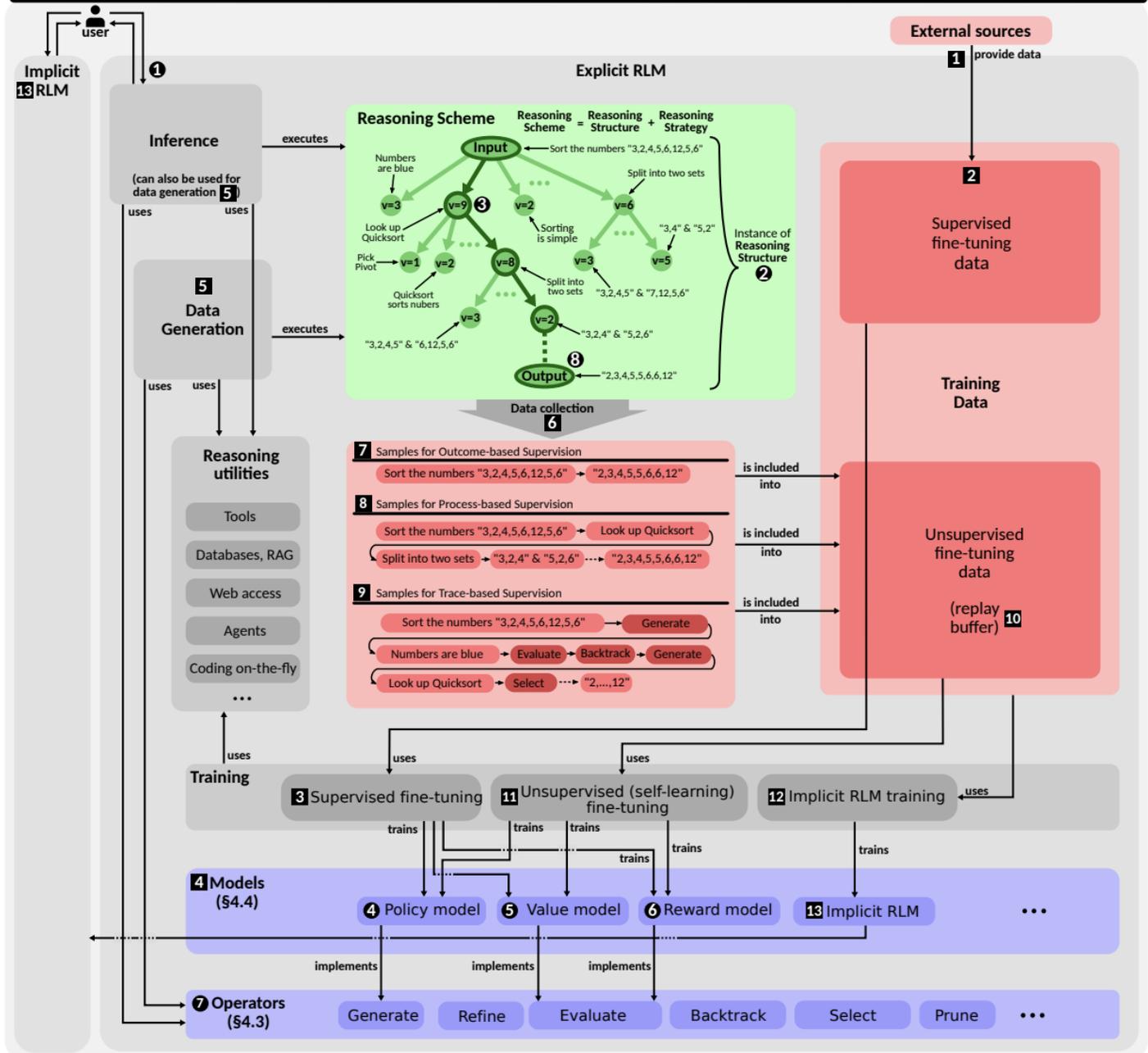
6



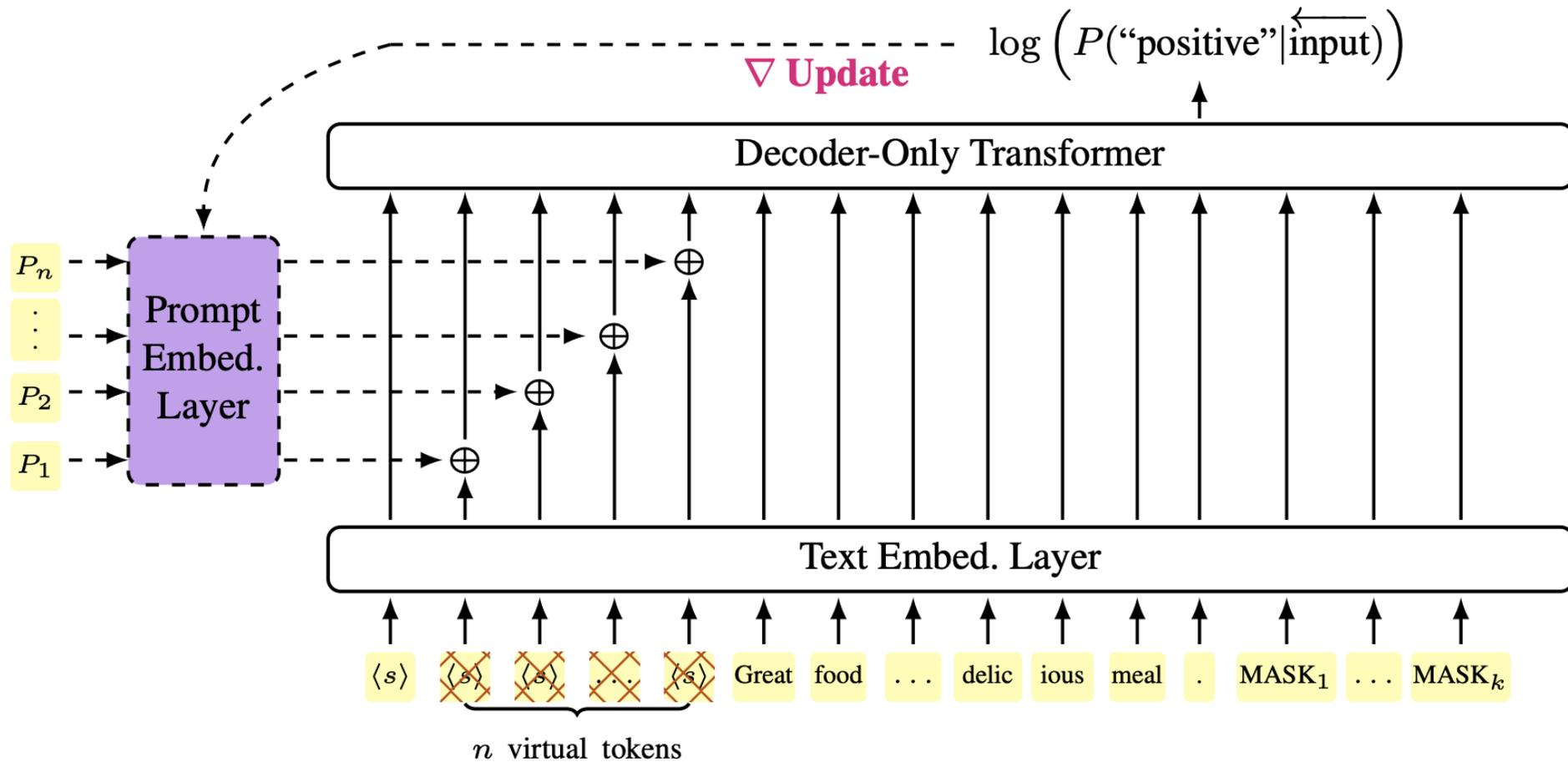
<https://arxiv.org/pdf/2501.11223>



# Reasoning Models



# Soft Prompting





# Parameter Efficient Fine-Tuning



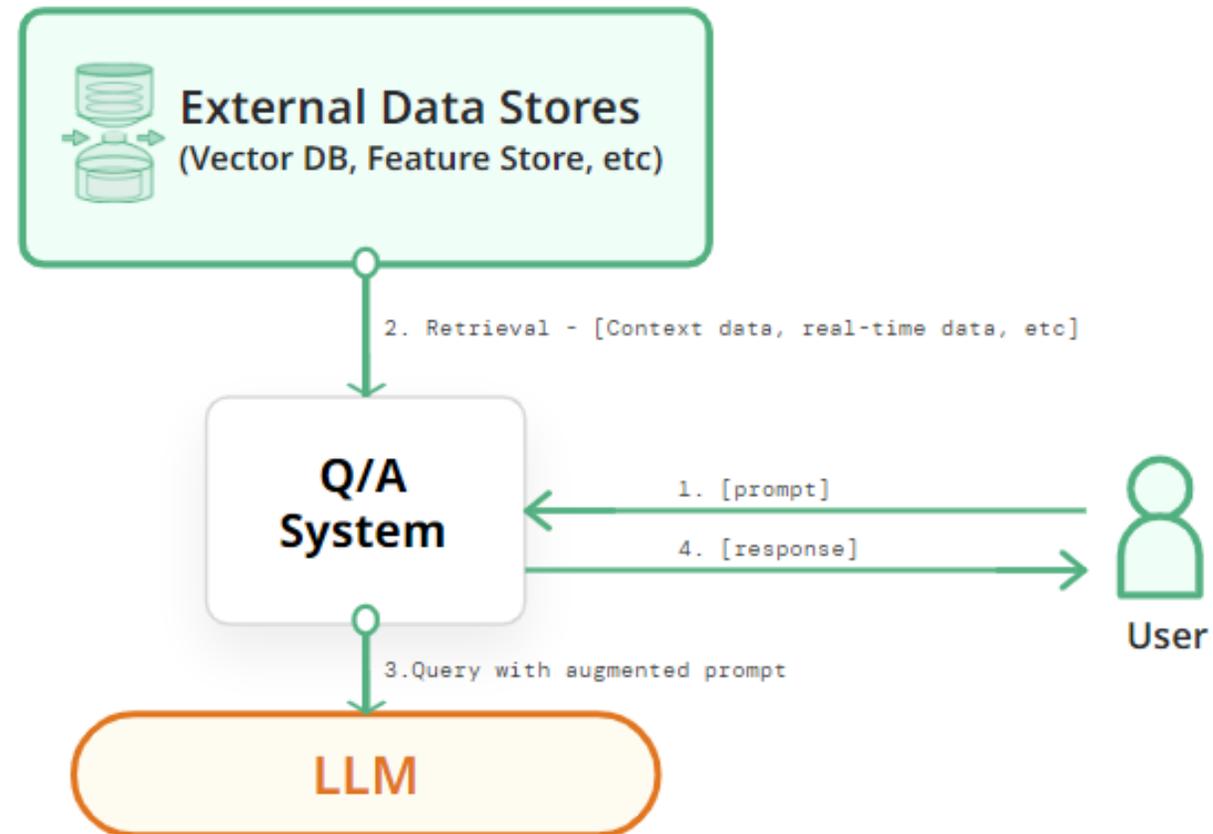
# Personalization / Adaptation / Alignment

---

- Every user has their own preferences, history, and contexts.
- **How can we efficiently adapt to each user?**

# Retrieval-Augment Generation

- Resource access enables personalization



Questions?

