

# Reasoning: From the Lab to the Real World



Zeming Chen (Eric)



Debjit Paul



Antoine Bosselut

# NLP applications

Machine Translation



Search Engines



Personal Assistants



Question Answer & ChatBoT



# NLP applications

**Generative AI empowering  
real-world NLP applications**

Machine Translation



Search Engines



Personal Assistants



Question Answer & ChatBoT



# Are they truly intelligent systems?

Machine Translation



Search Engines



Personal Assistants

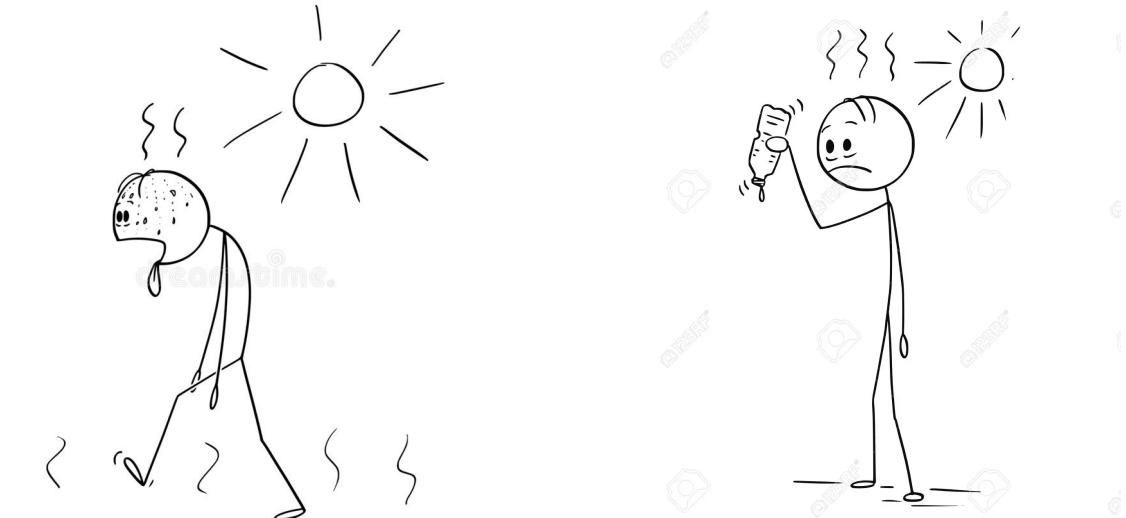


Question Answer & ChatBoT



# Machines that understand humans

Peter walked for 5 km in the **heat**.  
He went home and shook his water bottle.  
He was **disappointed** when it made **no sound**.



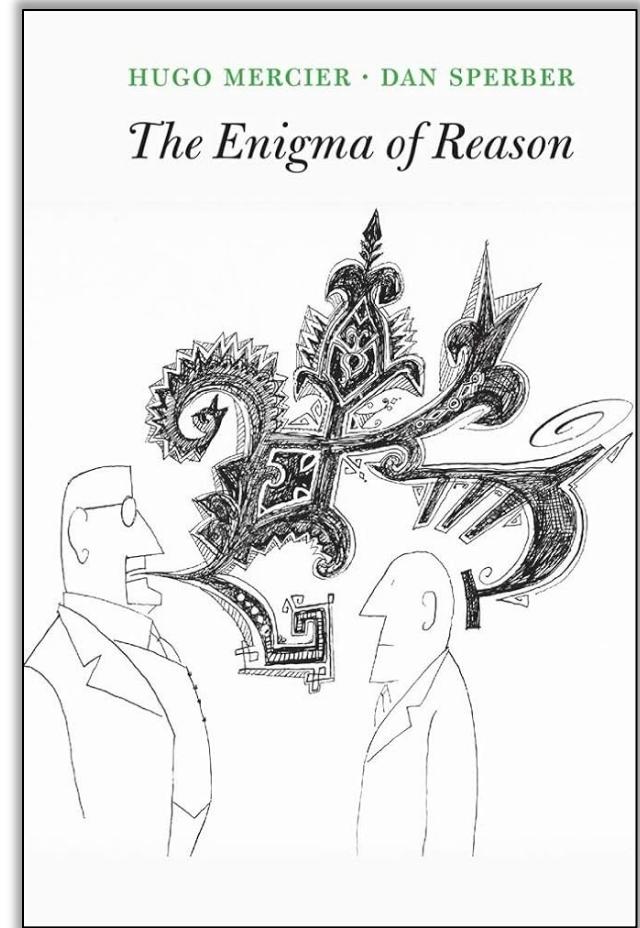
**Leave out important information**

**Rely on shared **intuitive** inferences to understand context**

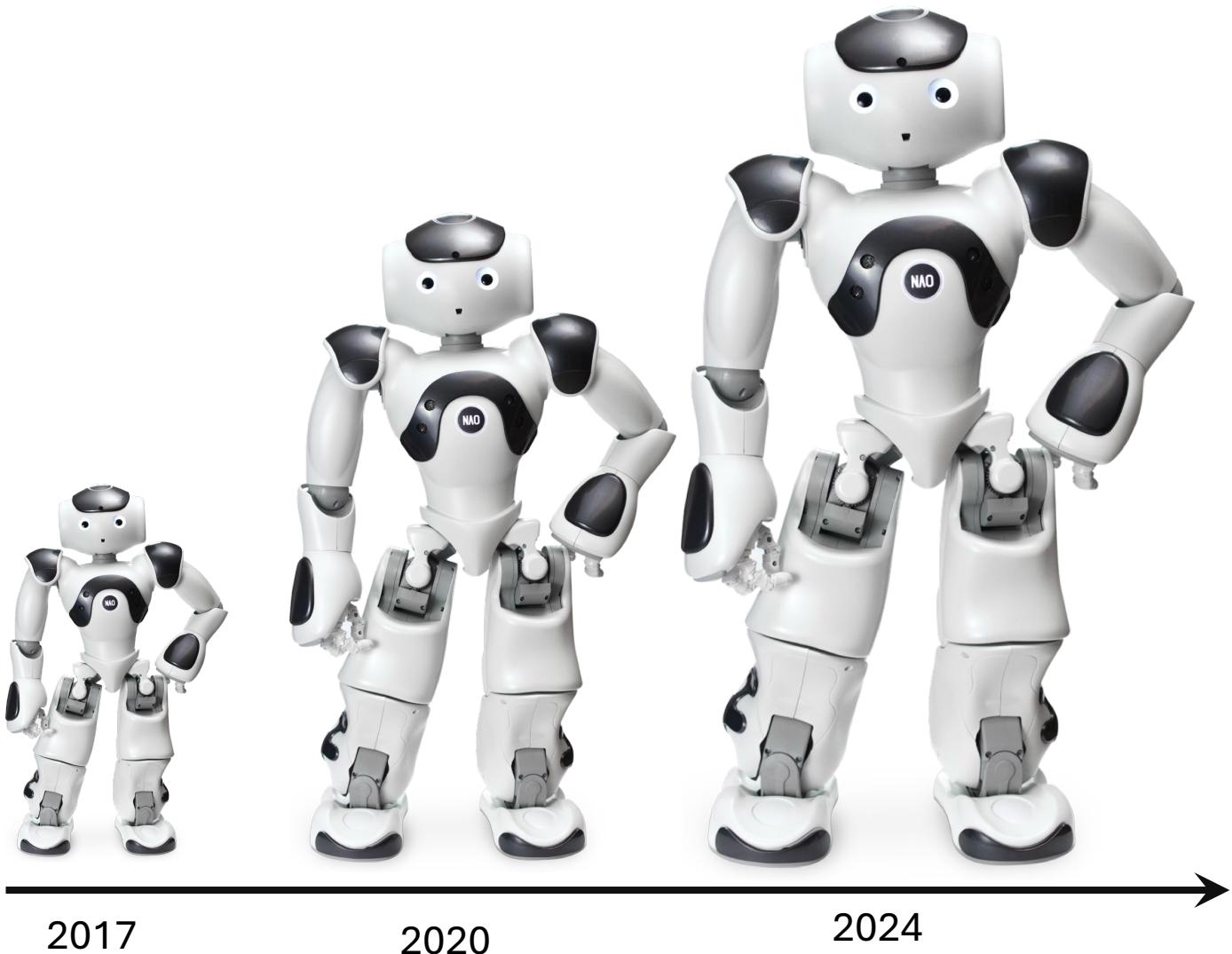
**Reason about knowledge to reach **understanding****

# How do humans use reasoning?

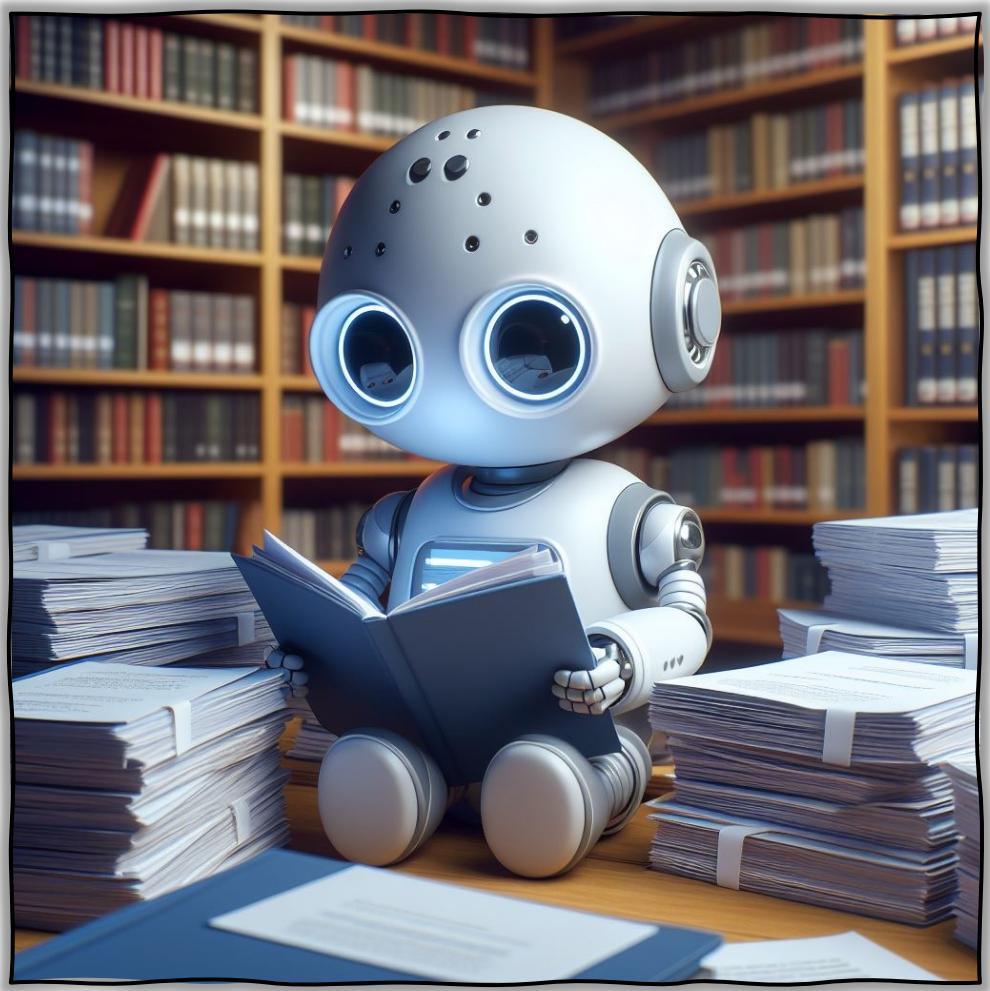
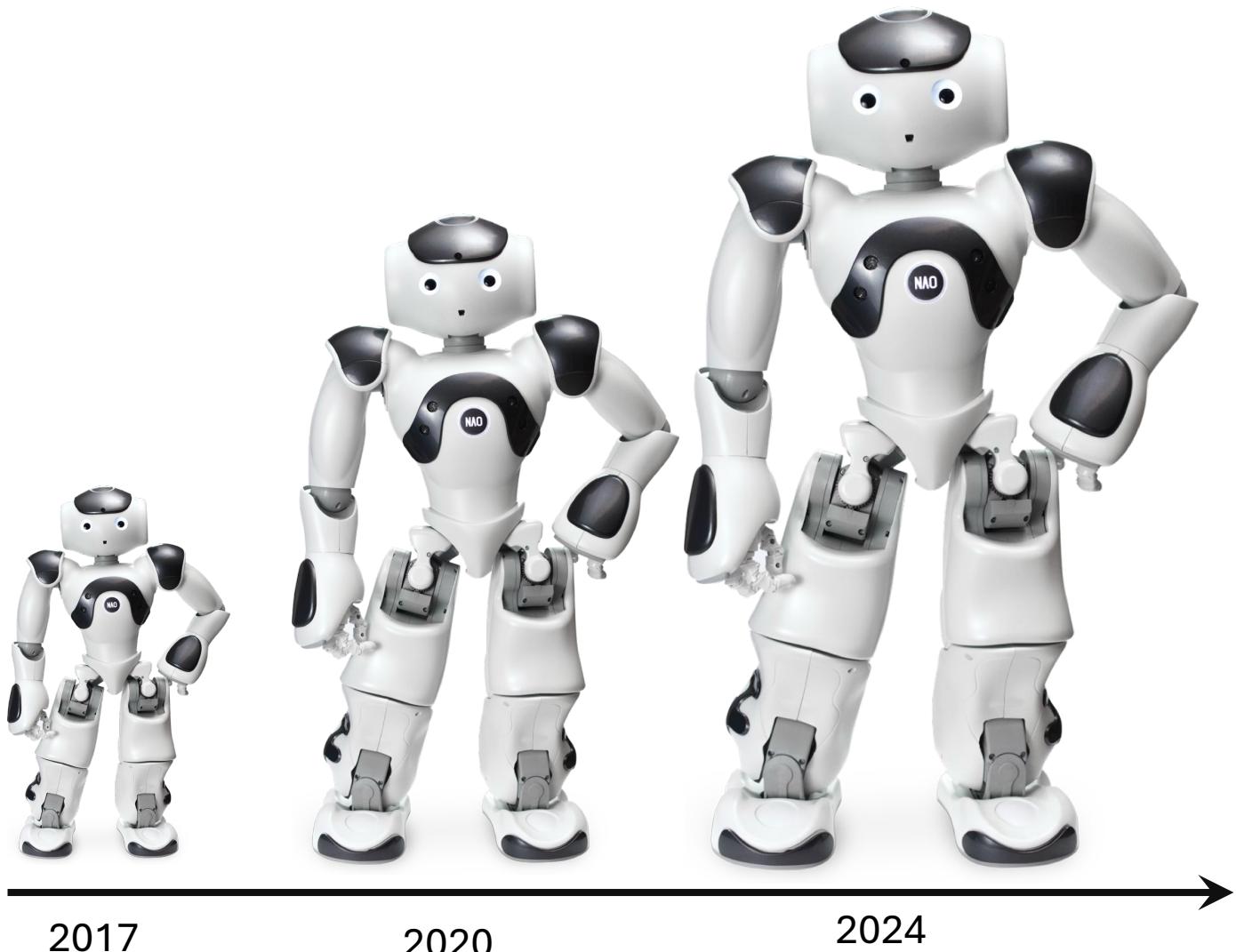
- Reasoning is tightly coupled with **explanation and justification** (Evans and Wason, 1976)
- Humans use reasoning to:
  - justify our intuitive inferences to convince others of our arguments
  - evaluate the arguments made by others



# Bigger and better



# More knowledge



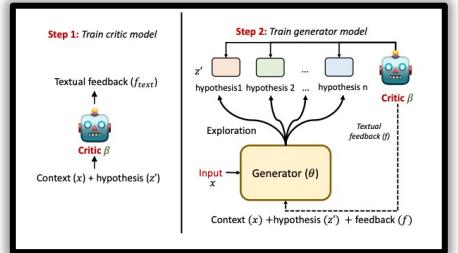
# Generative AI lacks reasoning



- Current systems fits observations to what they already know
- Rather truly altering the world model and accomodating new unseen experiences.

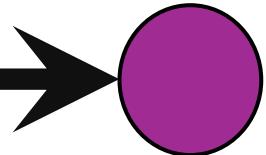


# Outline

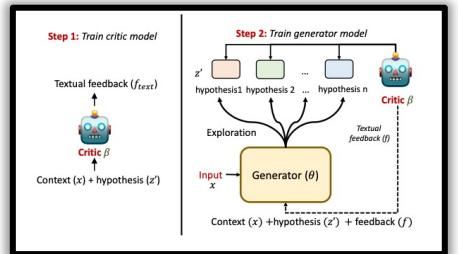


Refine knowledge using feedback

Paul et.al. 2024 (EACL)

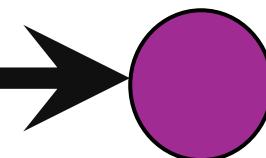


# Outline

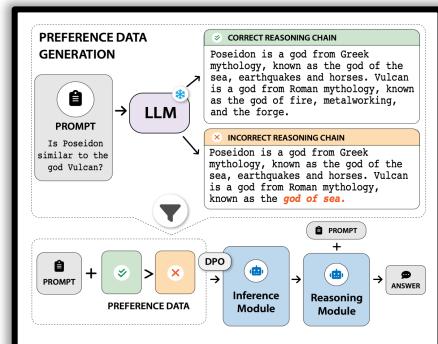


Refine knowledge using feedback

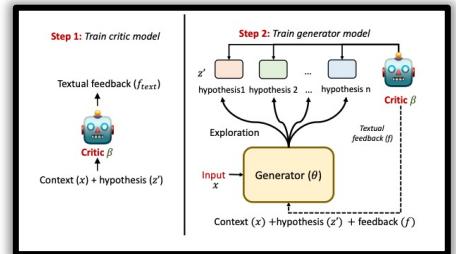
Paul et.al. 2024 (EACL)



Paul et.al. 2024 (EMNLP Findings)  
Faithful Reasoning about knowledge

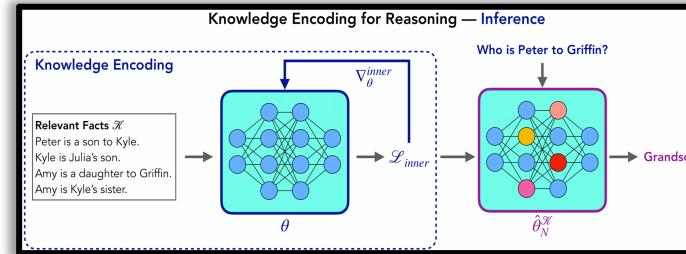


# Outline



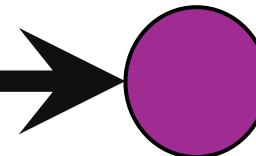
Refine knowledge using feedback

Paul et.al. 2024 (EACL)

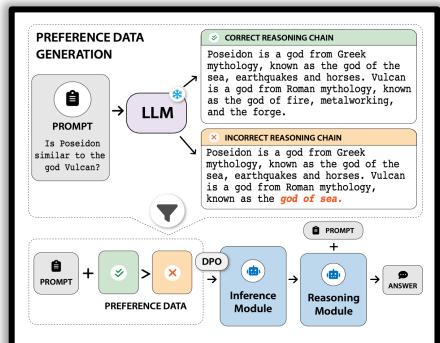


Robust Reasoning by updating world models

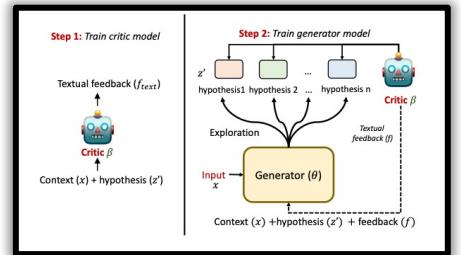
Chen et.al. 2023 (Neurips)



Paul et.al. 2024 (EMNLP Findings)  
Faithful Reasoning about knowledge

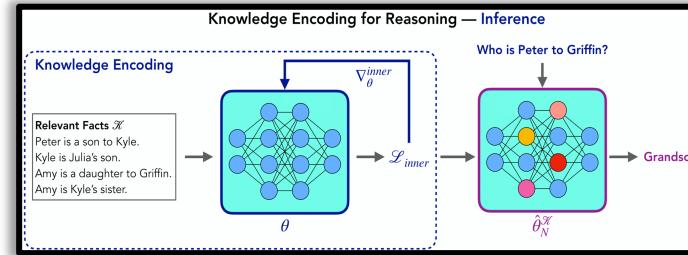


# Outline



Refine knowledge using feedback

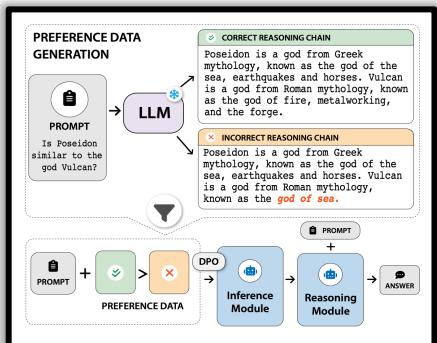
Paul et.al. 2024 (EACL)



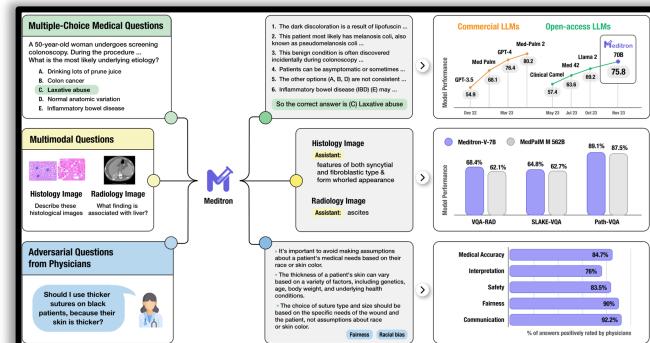
Robust Reasoning by updating world models

Chen et.al. 2023 (Neurips)

Paul et.al. 2024 (EMNLP Findings)  
Faithful Reasoning about knowledge



Chen et.al. 2024  
Real-World Reasoning



# REFINER: Reasoning Feedback on Intermediate Representations

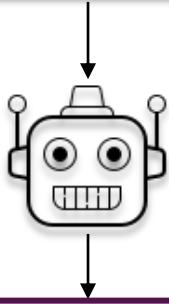


# Emergent abilities of large language models

- Chain-of-thought (CoT) prompting asks LLMs to generate reasoning traces before generating a final answer.

Wei et al. (2022)

Is travelling by plane more dangerous than traveling by car?  
**Let's Think Step by Step**



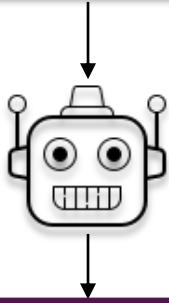
**Step1:** According to the National Safety Council, the lifetime odds of dying in a car accident are **1 in 102**.

**Step2:** The odds of dying in an air transport incident are substantially lower at **1 in 9,821**. **Answer: No**

# Emergent abilities of large language models

- The reasoning steps described in the CoT can be incorrect and inconsistent  
Turpin et al. (2023), Zhang et al. (2023)

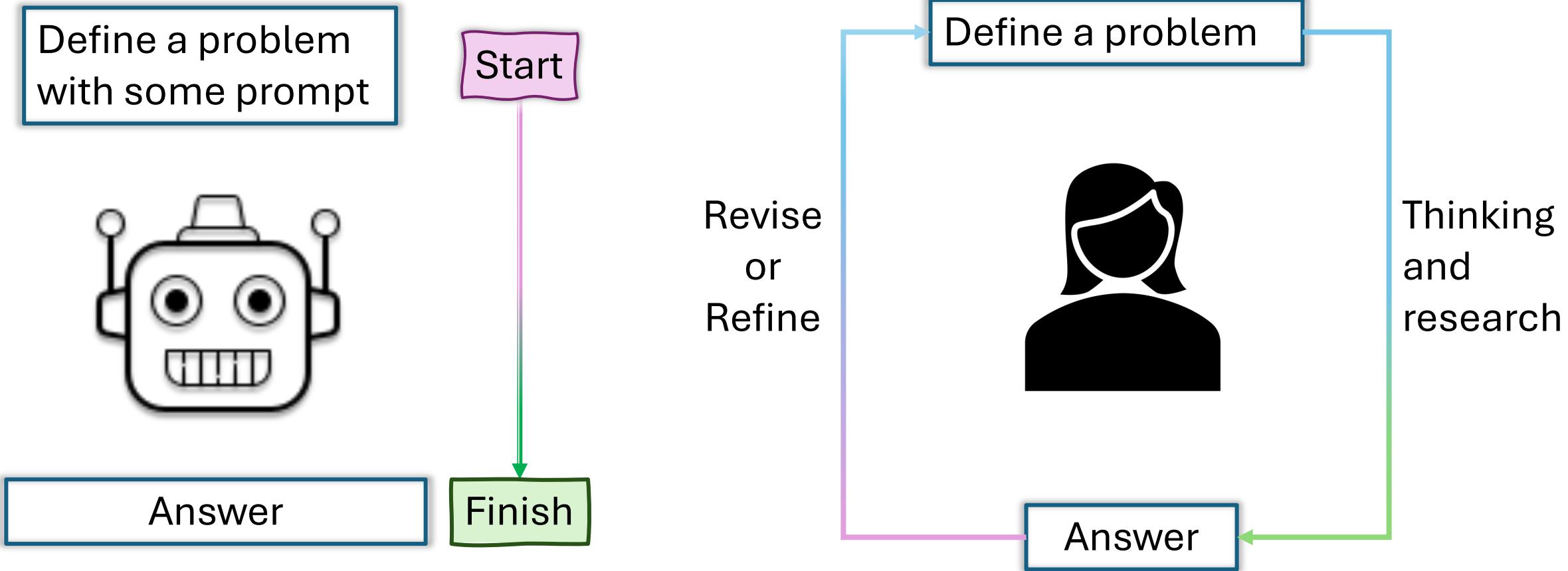
Is travelling by plane more dangerous than traveling by car?  
**Let's Think Step by Step**



**Step1:** According to the National Safety Council, the lifetime odds of dying in a car accident are **1 in 102**.

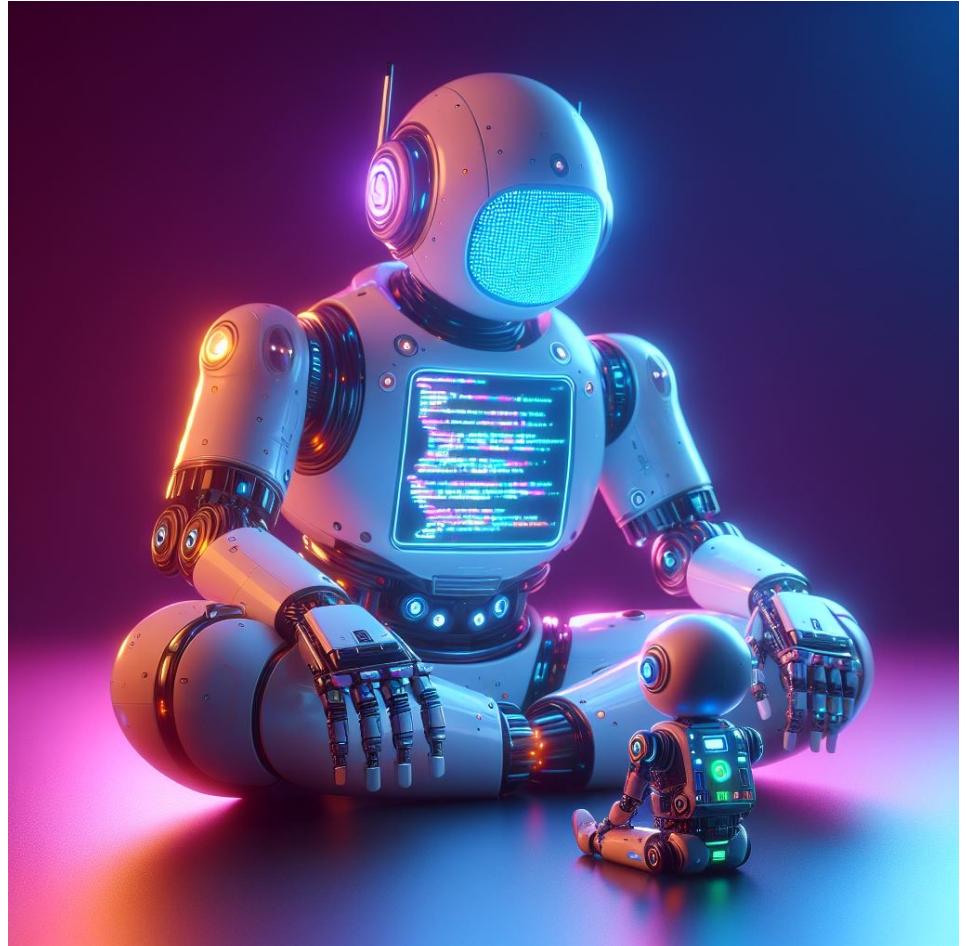
**Step2:** The odds of dying in an air transport incident are substantially lower at **1 in 9**. **Answer: Yes**

# Recursive Reminding



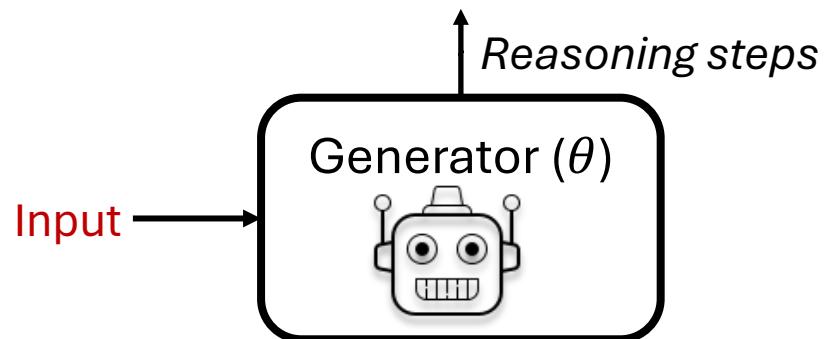
# Research Question

How do we build a method to refine  
the reasoning chains generated by  
LMs automatically?



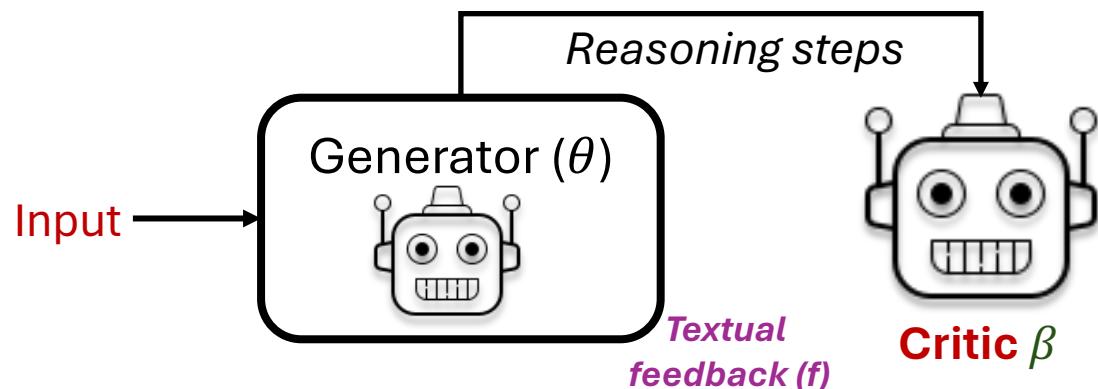
# REFINER

- An interaction-based framework
  - Iteratively refine the reasoning steps



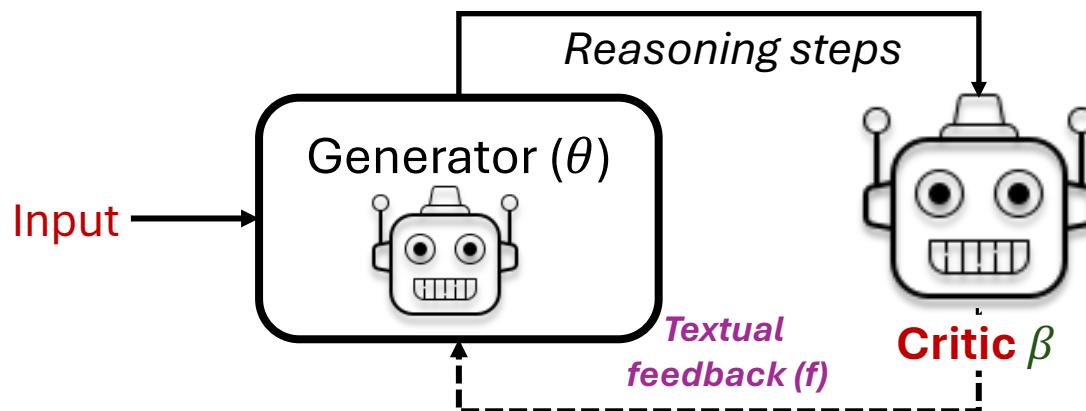
# REFINER

- An interaction-based framework
  - Iteratively refine the reasoning steps



# REFINER

- An interaction-based framework
  - Iteratively refine the reasoning steps



---

#### Algorithm 1 REFINER Training

---

```
1: for E epochs do
2:   for  $i(\text{batch}) \leftarrow 1$  to  $N$  do
3:     Initialize (feedback)  $f_0 \leftarrow No$ 
4:     for  $t \leftarrow 1$  to  $T$  do
5:        $\hat{z}_{i,t}^k \sim \pi_\theta(y_i | c_i, f_{t-1}, \hat{z}_{i,t-1})$ 
6:        $f_t, \hat{z} \leftarrow \pi_\beta(c_i, z_i, \hat{z}_{i,t}^k)$ 
7:        $\mathcal{L}_i^{lm} += -\log p(z_i | c_i, f_{t-1}, \hat{z}_{i,t-1})$ 
8:     end for
9:   end for
10: end for
11: return  $\pi_\theta$ 
```

---

# Feedback Data

- Identify different fine-grained reasoning errors.
- Create training data for the critic model by perturbing the correct reasoning steps.

## Correct Reasoning Steps

**Step1:** First, you started with 10 apples. You gave away 2 apples to the neighbor and 2 to the repairman, so you had <<10-2-2>> 6 apples left.

**Step2:** Then you bought 5 more apples, so now you had 11 apples.

# Feedback Data

- Introduce different reasoning errors and create incorrect reasoning steps.

## Correct Reasoning Steps

**Step1:** First, you started with 10 apples. You gave away 2 apples to the neighbor and 2 to the repairman, so you had <<10-2-2>> 6 apples left.

**Step2:** Then you bought 5 more apples, so now you had 11 apples.

## Incorrect Reasoning Steps

**Step1:** First, you started with 10 apples. You gave away 2 apples to the neighbor and 4 to the repairman, so you had <<10+2-4>> **10 apples left**.

**Step2:** Then you bought 5 more apples, so now you had <<10 + 5 >> 15 apples.

# Feedback Data

- Based on the error types we create semi-structured textual feedback.

## Correct Reasoning Steps

**Step1:** First, you started with 10 apples. You gave away 2 apples to the neighbor and 2 to the repairman, so you had <<10-2-2>> 6 apples left.

**Step2:** Then you bought 5 more apples, so now you had 11 apples.

## Incorrect Reasoning Steps

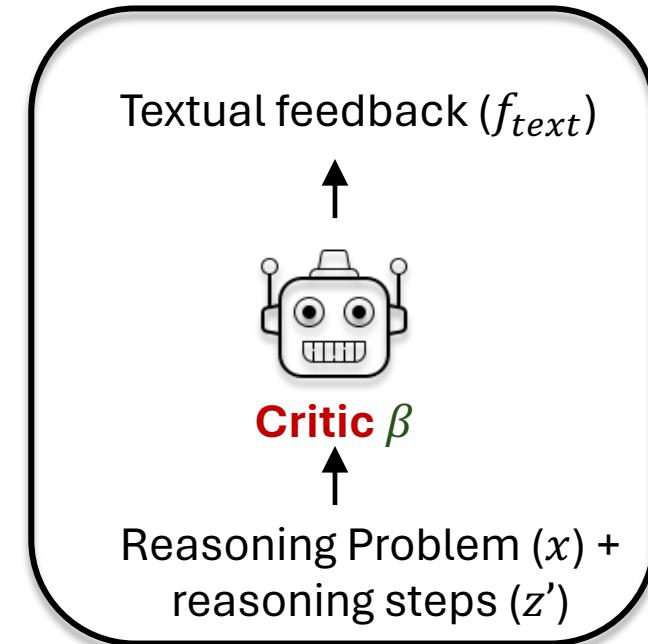
**Step1:** First, you started with 10 apples. You gave away 2 apples to the neighbor and 4 to the repairman, so you had <<10+2-4>> **10 apples left**.

**Step2:** Then you bought 5 more apples, so now you had <<10 + 5 >> 15 apples.

**Textual Feedback:** The operand and operator in the step 1 are incorrect.

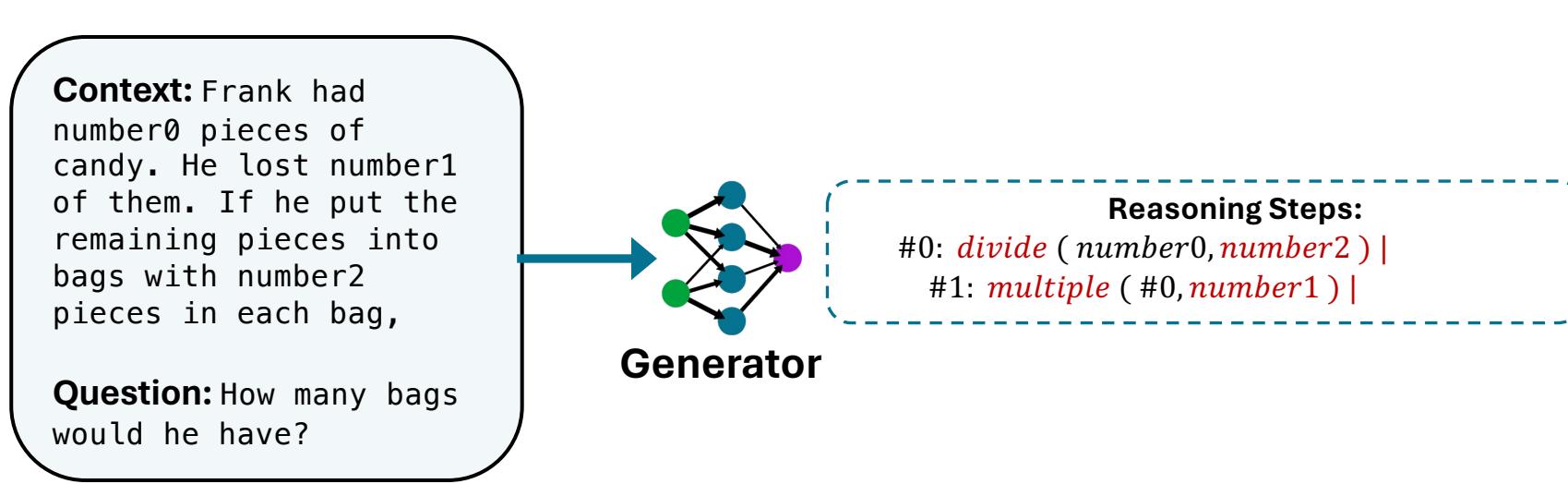
# Train a critic model

- Critic model learns to provide textual feedback.
- We train a supervised model with the reasoning question (context) and correct or incorrect reasoning steps as input and the textual feedback as output.



$$\text{Loss} = - \sum_i \log P(f_i^* | x, z)$$

# REFINER



---

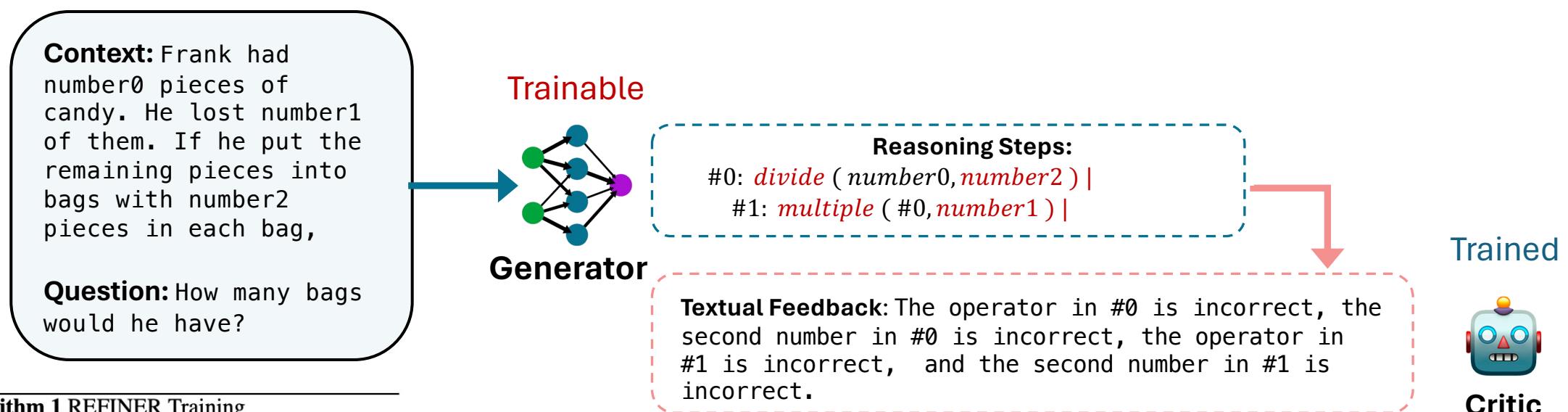
## Algorithm 1 REFINER Training

---

```
1: for E epochs do
2:   for  $i(\text{batch}) \leftarrow 1$  to  $N$  do
3:     Initialize (feedback)  $f_0 \leftarrow No$ 
4:     for  $t \leftarrow 1$  to  $T$  do
5:        $\hat{z}_{i,t}^k \sim \pi_\theta(y_i|c_i, f_{t-1}, \hat{z}_{i,t-1})$ 
6:        $f_t, \hat{z} \leftarrow \pi_\beta(c_i, z_i, \hat{z}_{i,t}^k)$ 
7:        $\mathcal{L}_i^{lm} += -\log p(z_i|c_i, f_{t-1}, \hat{z}_{i,t-1})$ 
8:     end for
9:   end for
10: end for
11: return  $\pi_\theta$ 
```

---

# REFINER – Example




---

## Algorithm 1 REFINER Training

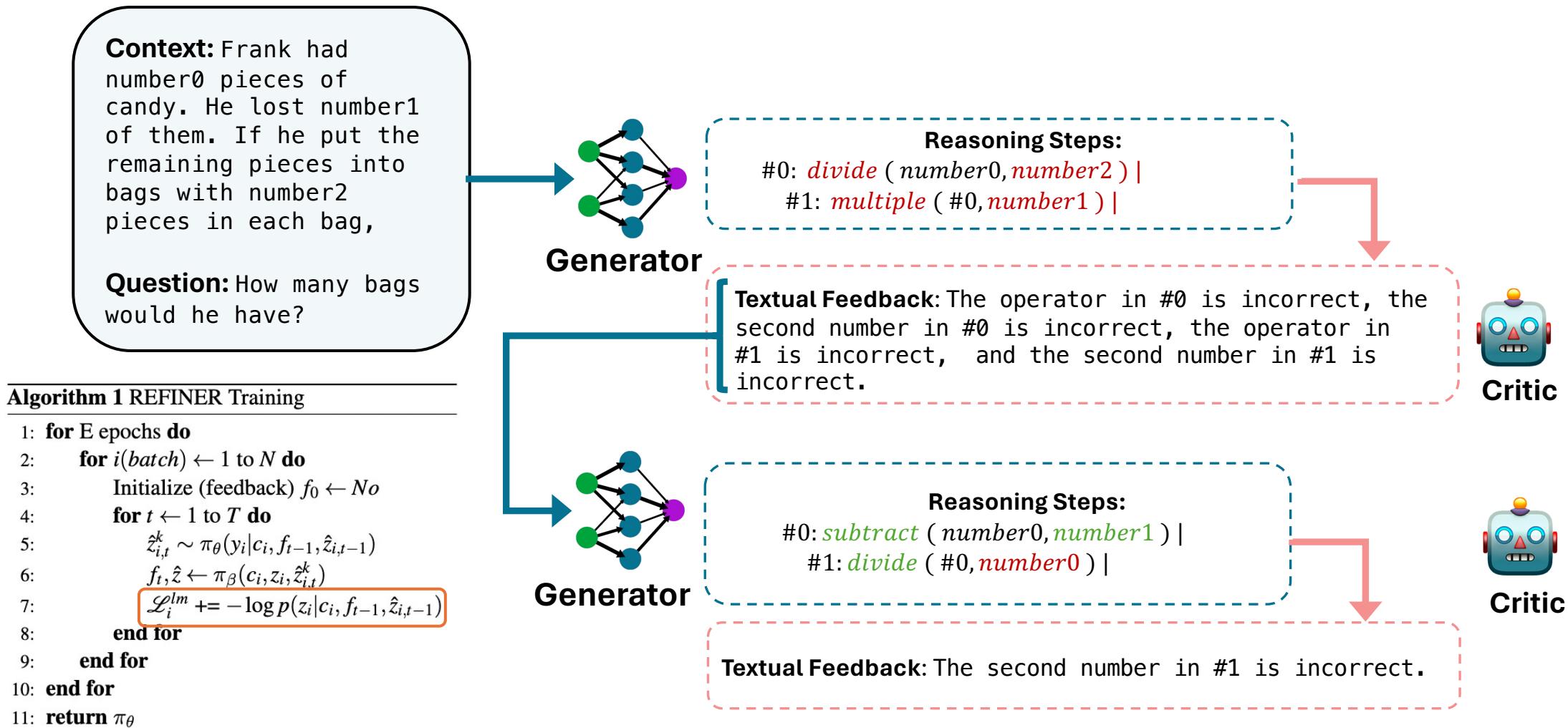
```

1: for E epochs do
2:   for  $i(batch) \leftarrow 1$  to  $N$  do
3:     Initialize (feedback)  $f_0 \leftarrow No$ 
4:     for  $t \leftarrow 1$  to  $T$  do
5:        $\hat{z}_{i,t}^k \sim \pi_\theta(y_i|c_i, f_{t-1}, \hat{z}_{i,t-1})$ 
6:        $f_t, \hat{z} \leftarrow \pi_\beta(c_i, z_i, \hat{z}_{i,t}^k)$ 
7:        $\mathcal{L}_i^{lm} += -\log p(z_i|c_i, f_{t-1}, \hat{z}_{i,t-1})$ 
8:     end for
9:   end for
10: end for
11: return  $\pi_\theta$ 

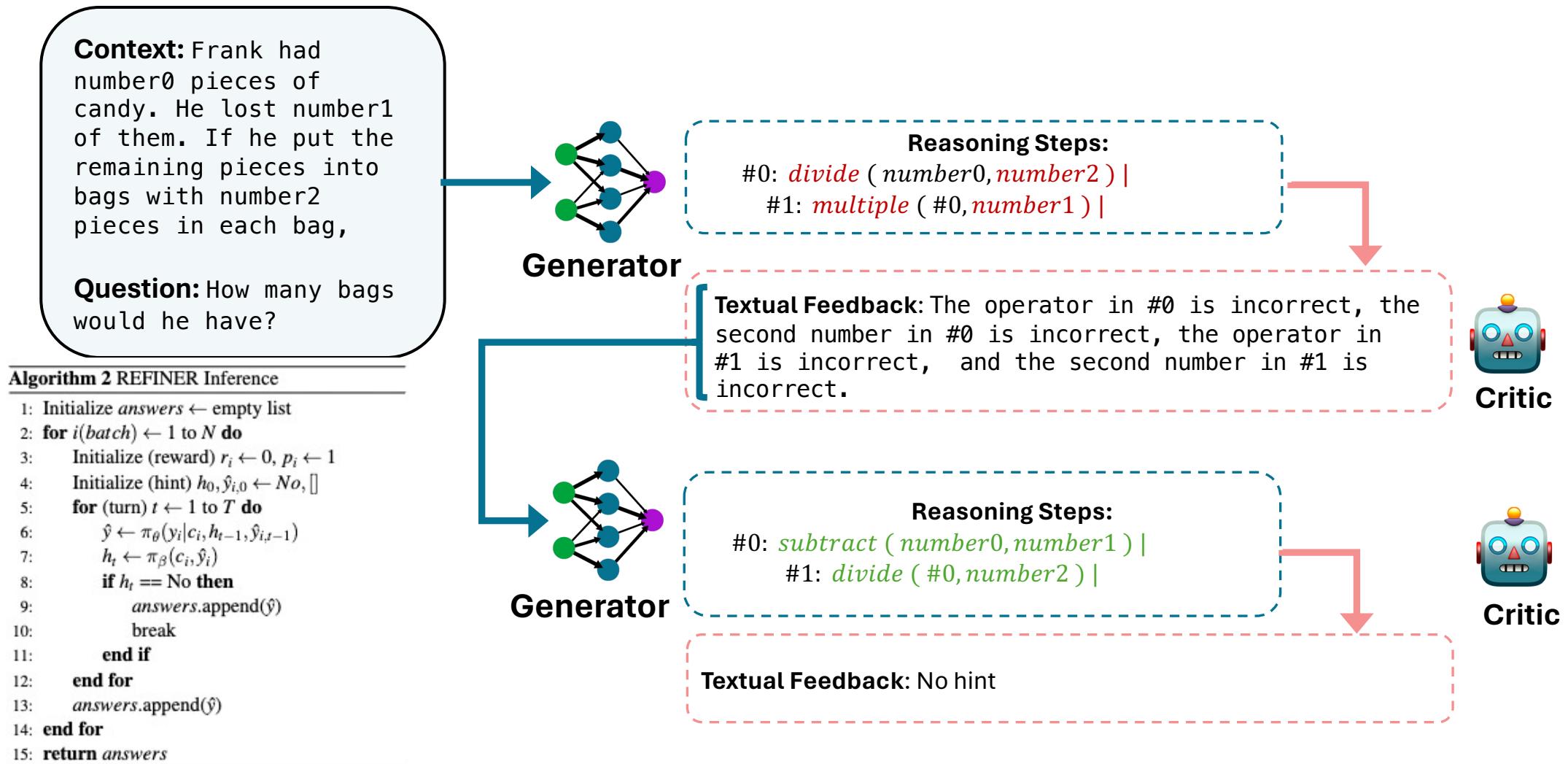
```

---

# REFINER – Example

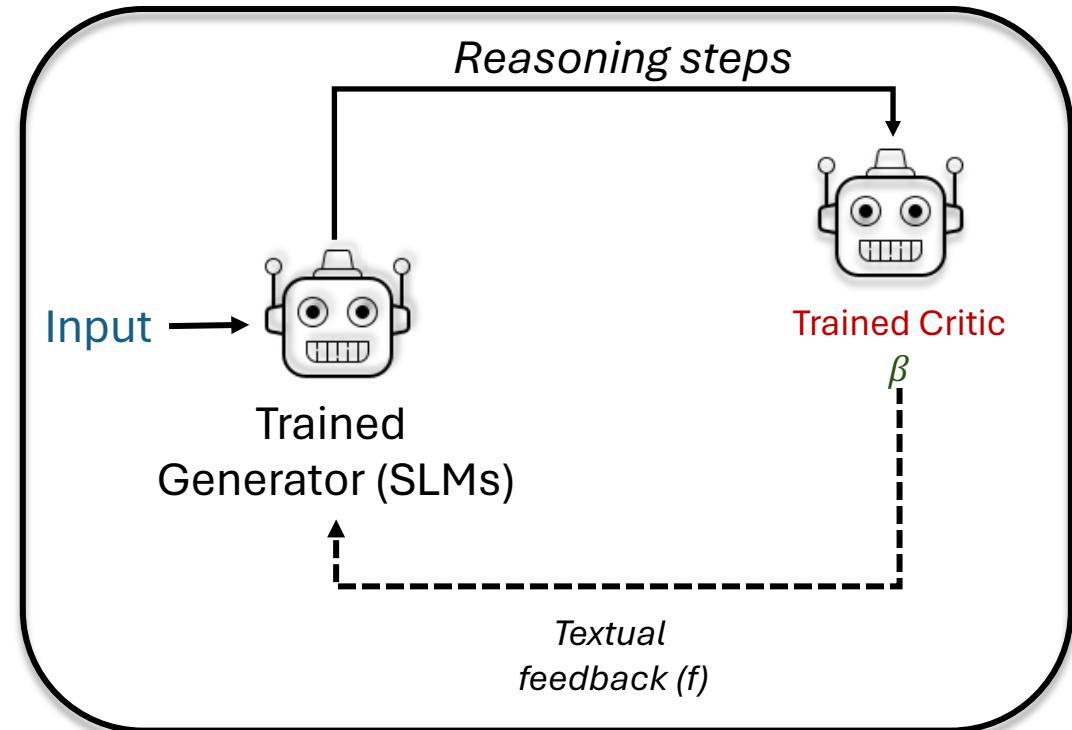


# REFINER – Example



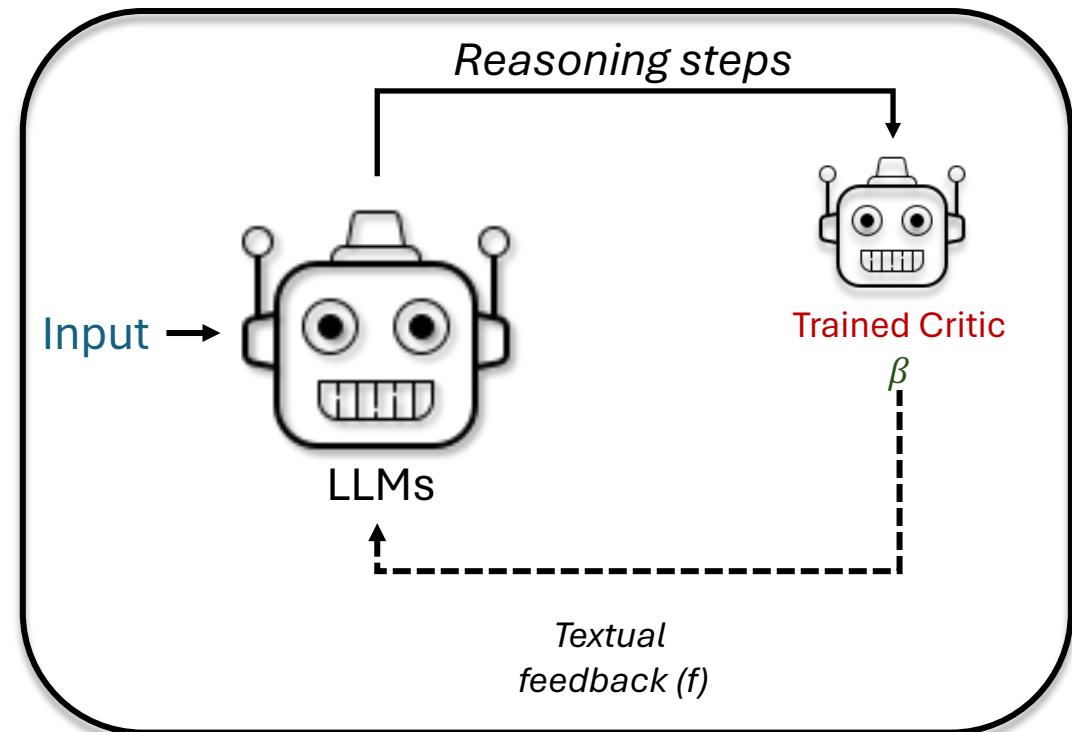
# Inference

- We use the trained critic along with the trained generator to generate a trajectory  $z_0, z_1, \dots, z_T$  and stop when either  $z_T$  is generated by the generator or “No hint” is generated by the critic.

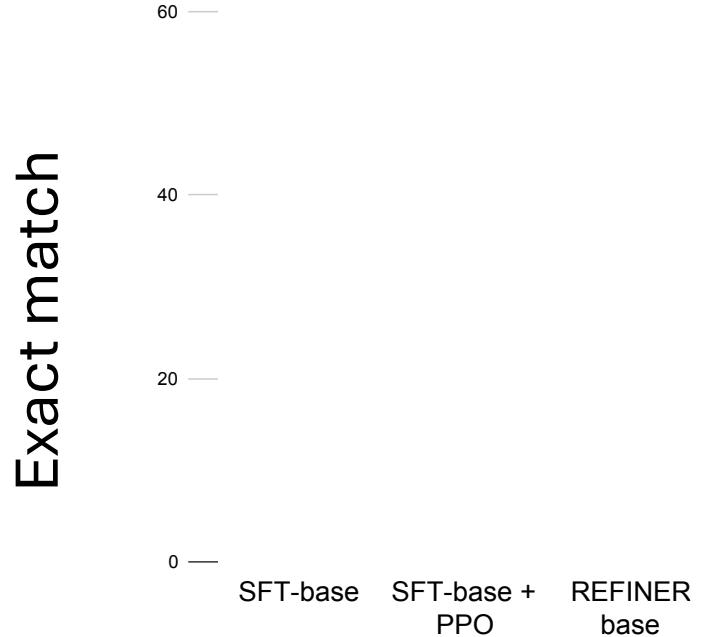


# Inference – Trained Critic as a Tool

- We also experimented with **LLMs** (chain of thought prompting), where the generator generates a trajectory  $z_0, z_1, \dots, z_T$  and stops when the critic generates “No hint”.



# How well REFINER works?

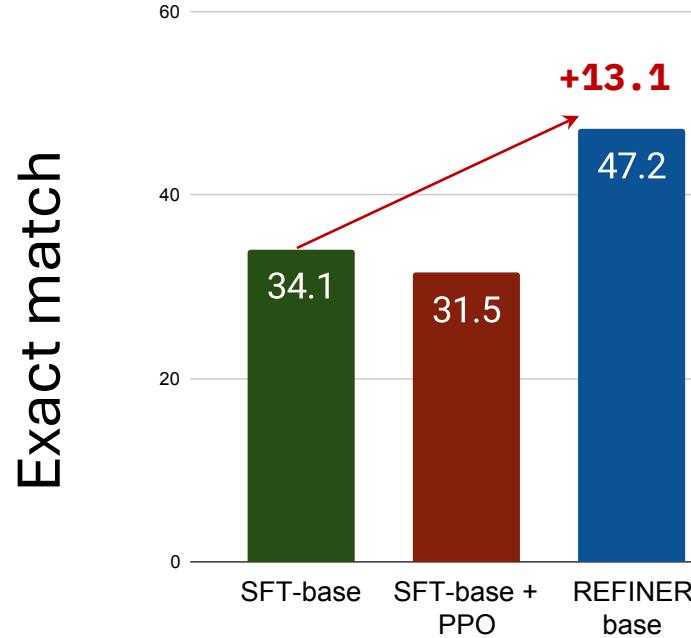


**Math Word Problem**

**Dataset: SVAMP**

**Base Model : T5**

# How well REFINER works on supervised models?



**Dataset:** SVAMP

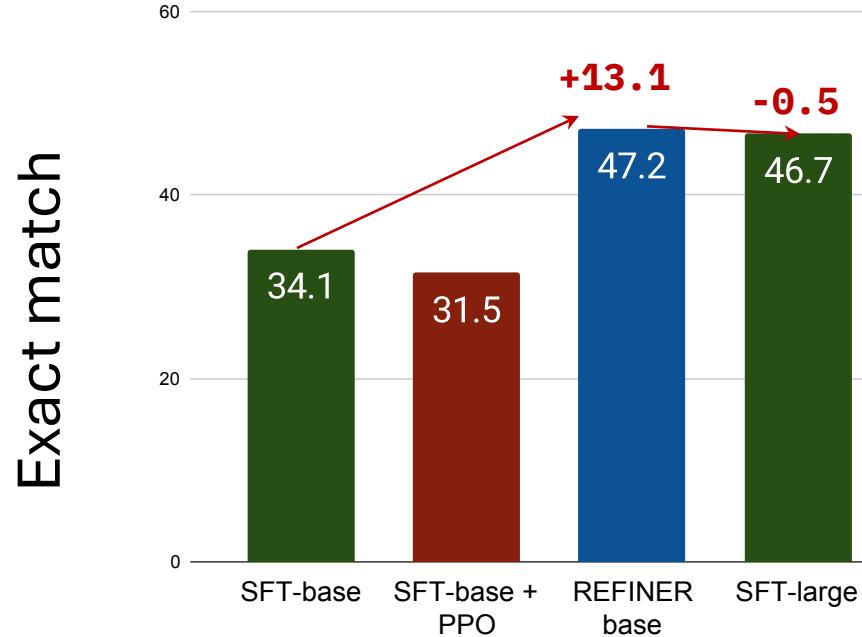
- We compare REFINER with SFT and PPO.
- Adding a critic in the loop improves **+13.1** exact match scores over SFT.

**Math Word Problem**

**Dataset:** SVAMP

**Base Model :** T5

# How well REFINER works on supervised models?



**Dataset:** SVAMP

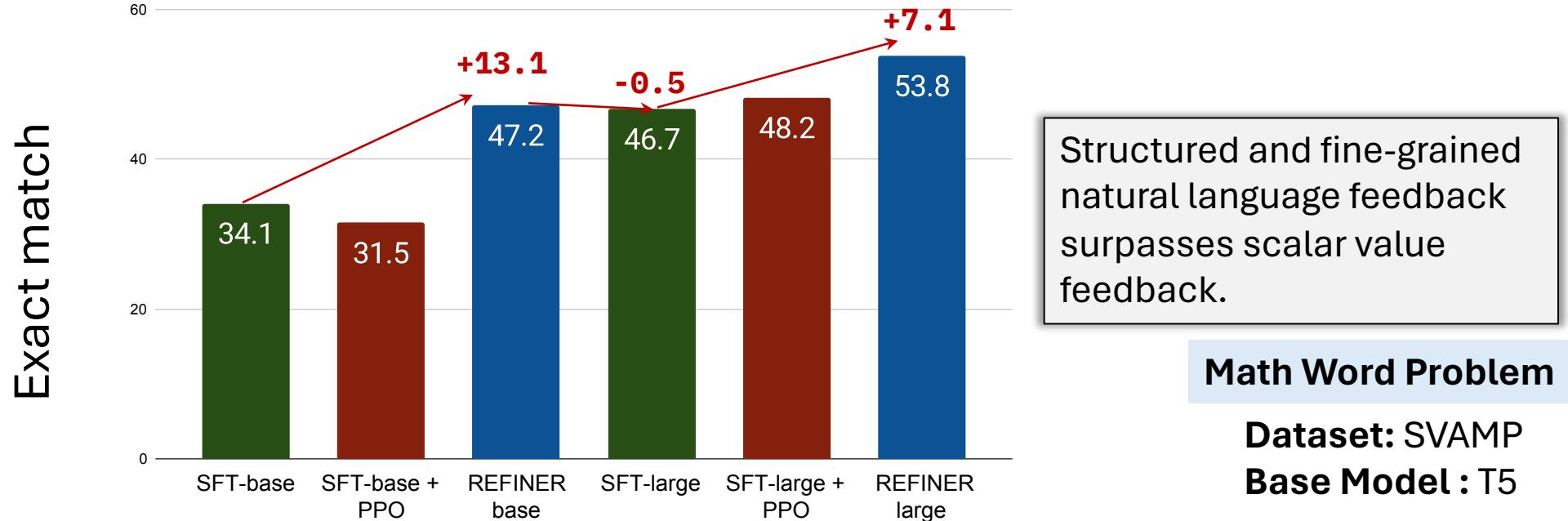
- We compare REFINER with SFT and PPO.
- Adding a critic in the loop improves **+13.1** exact match scores over SFT.
- REFINER base model outperforms SFT-large model.

**Math Word Problem**

**Dataset:** SVAMP

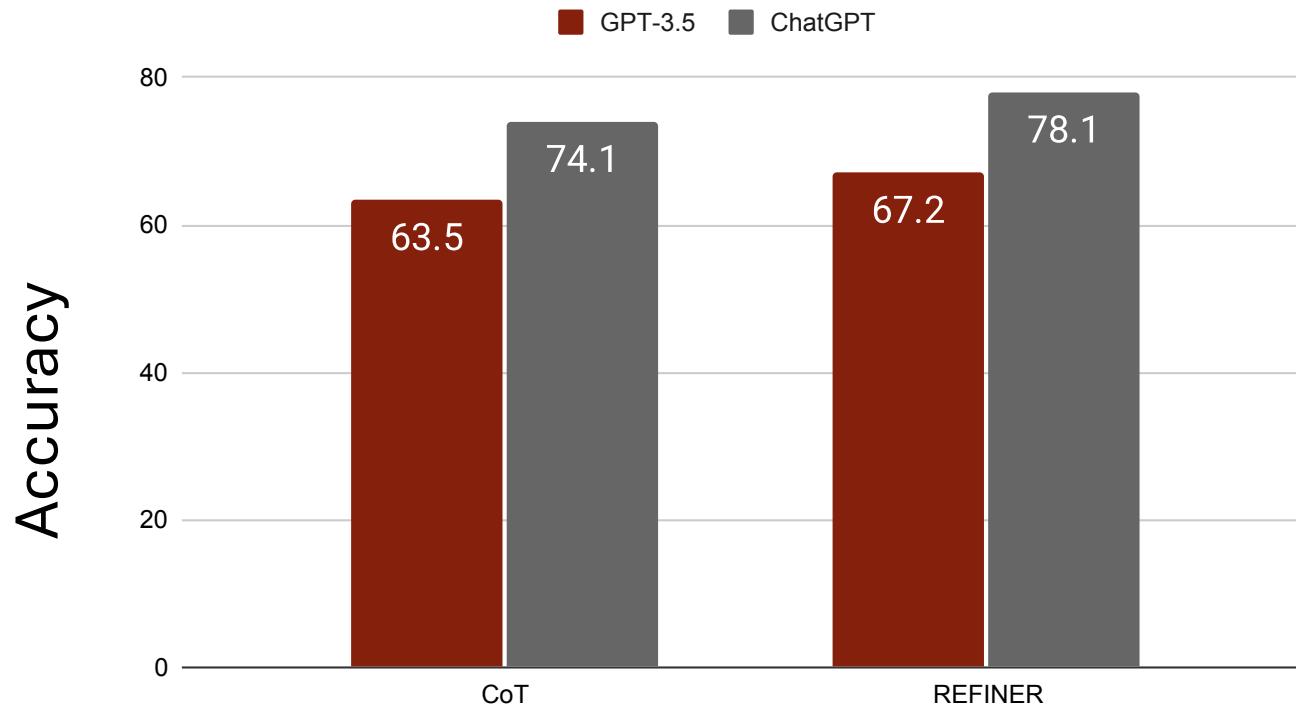
**Base Model :** T5

# How well REFINER works on supervised models?



- We compare REFINER with SFT and PPO.
- Adding a critic in the loop improves **+13.1** exact match scores over SFT.
- REFINER base model outperforms SFT-large model.

# Can we use critic model as a tool?

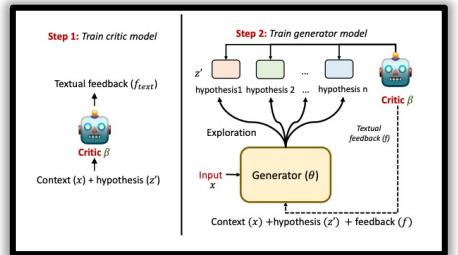


Trained critic  
can be used as  
a tool.

**Math Word Problem**  
**Dataset: GSM8K**

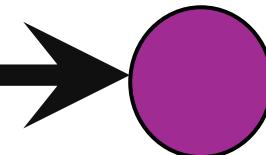
- REFINER can improve CoT by **+3.7** and **+4.0**.

# Outline

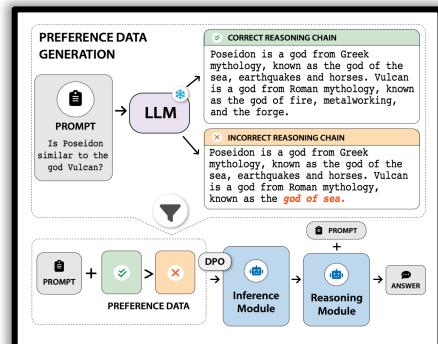


Refine knowledge using feedback

Paul et.al. 2024 (EACL)



Paul et.al. 2024 (EMNLP Findings)  
Faithful Reasoning about knowledge



# Making Reasoning Matter: Measuring and Improving the Faithfulness of Chain of Thoughts

EMNLP Findings 2024



**Debjit Paul**



Robert West



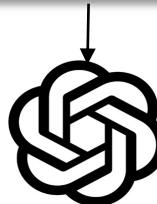
Antoine Bosselut



Boi Faltings

# Does LLMs reliably use their own knowledge?

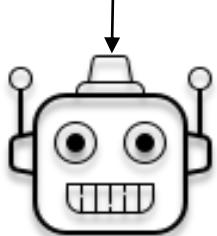
Is travelling by plane more dangerous than traveling by car?



GPT-4

**Step1:** According to the National Safety Council, the lifetime odds of dying in a car accident are **1 in 102**.

**Step2:** The odds of dying in an air transport incident are substantially lower at **1 in 9,821**.

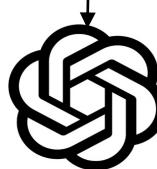


Model

False

# Does LLMs reliably use their own knowledge?

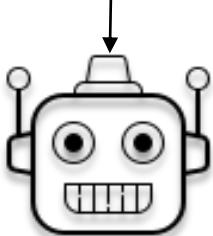
Is travelling by plane more dangerous than traveling by car?



GPT-4

**Step1:** According to the National Safety Council, the lifetime odds of dying in a car accident are **1 in 9,821**.

**Step2:** The odds of dying in an air transport incident are substantially lower at **1 in 102**.



Model

True

We study how reliably LLMs use inference chains to arrive at a conclusion.

# Research Question I

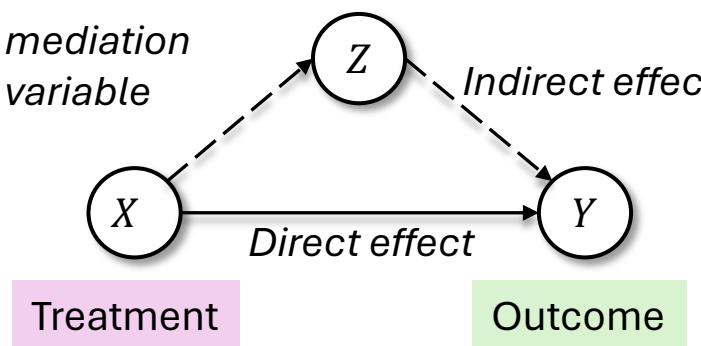


How can we effectively measure the  
faithfulness of LLMs when reasoning over  
reasoning chain?

# Causal mediation analysis

- It is a method to measure how an independent variable (or treatment) affects a dependent variable (or outcome) mediated by intermediate variables

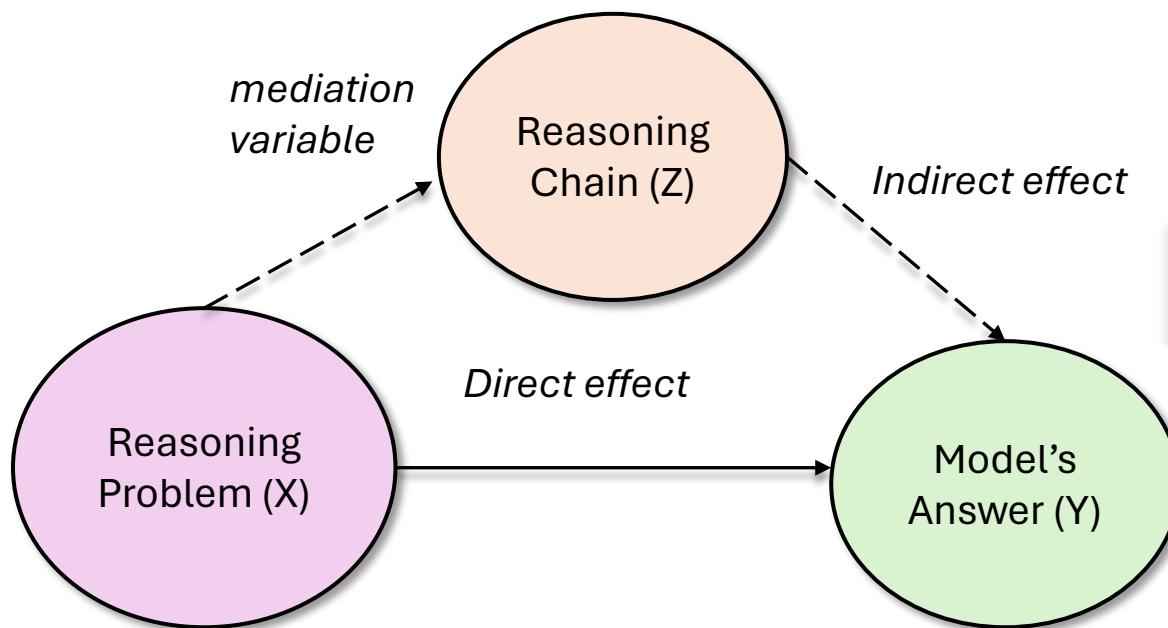
Pearl 2021



Judea Pearl

# Causal mediation analysis

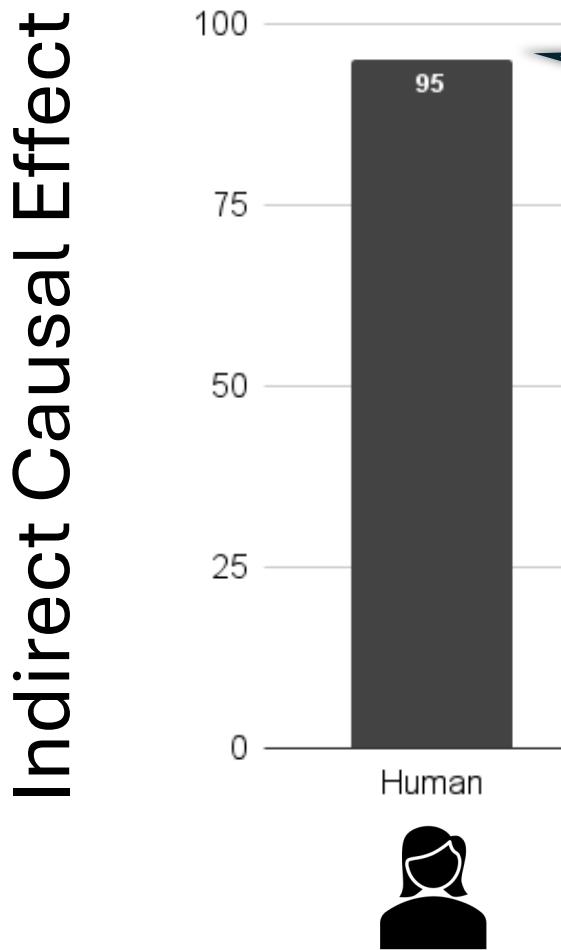
- It is a method to measure how an independent variable (or treatment) affects a dependent variable (or outcome) mediated by intermediate variables
- Measures the relationship between the reasoning chain (Z) and the model final answer (Y)



$$IE = E(Y | X, Z) - E(Y | X, Z')$$

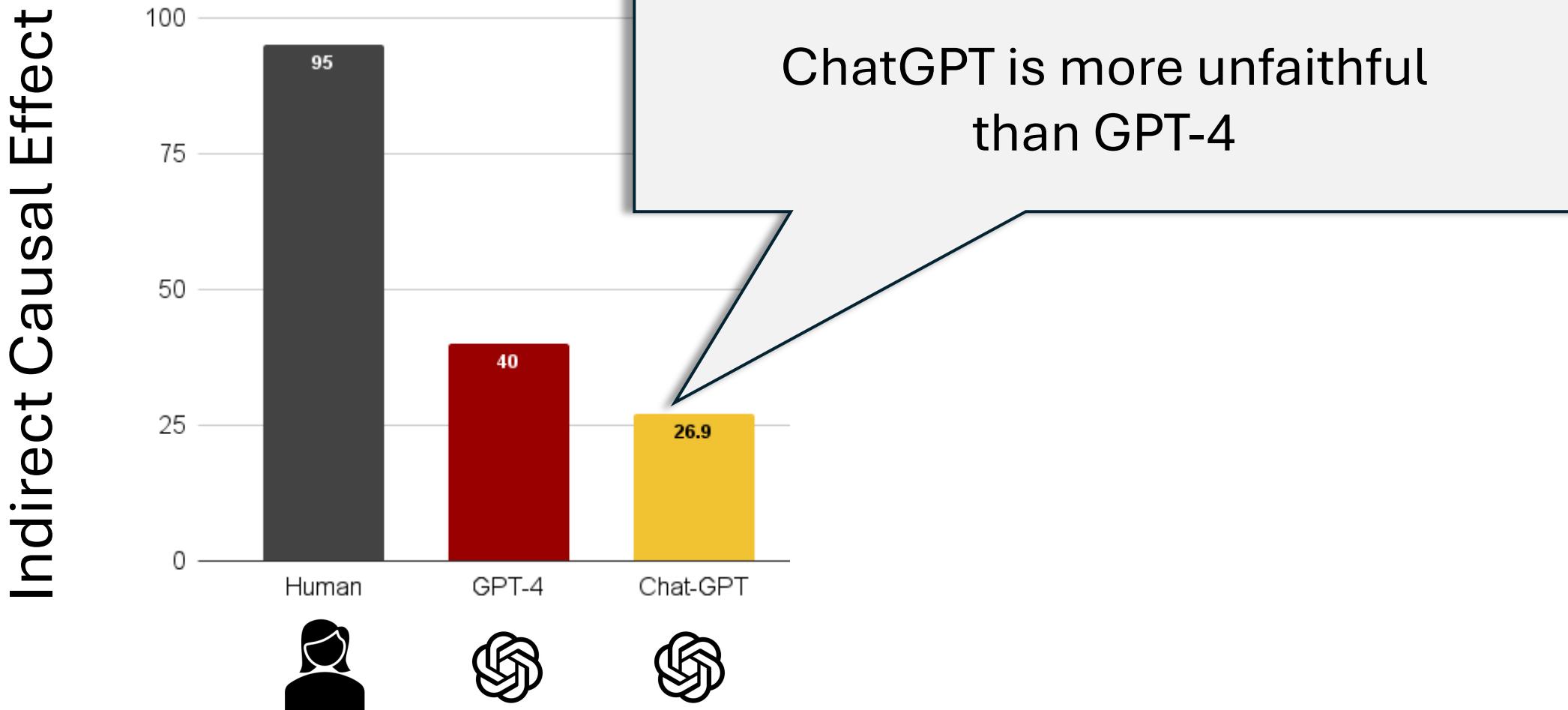
Pearl, 2001

# How well do humans perform?

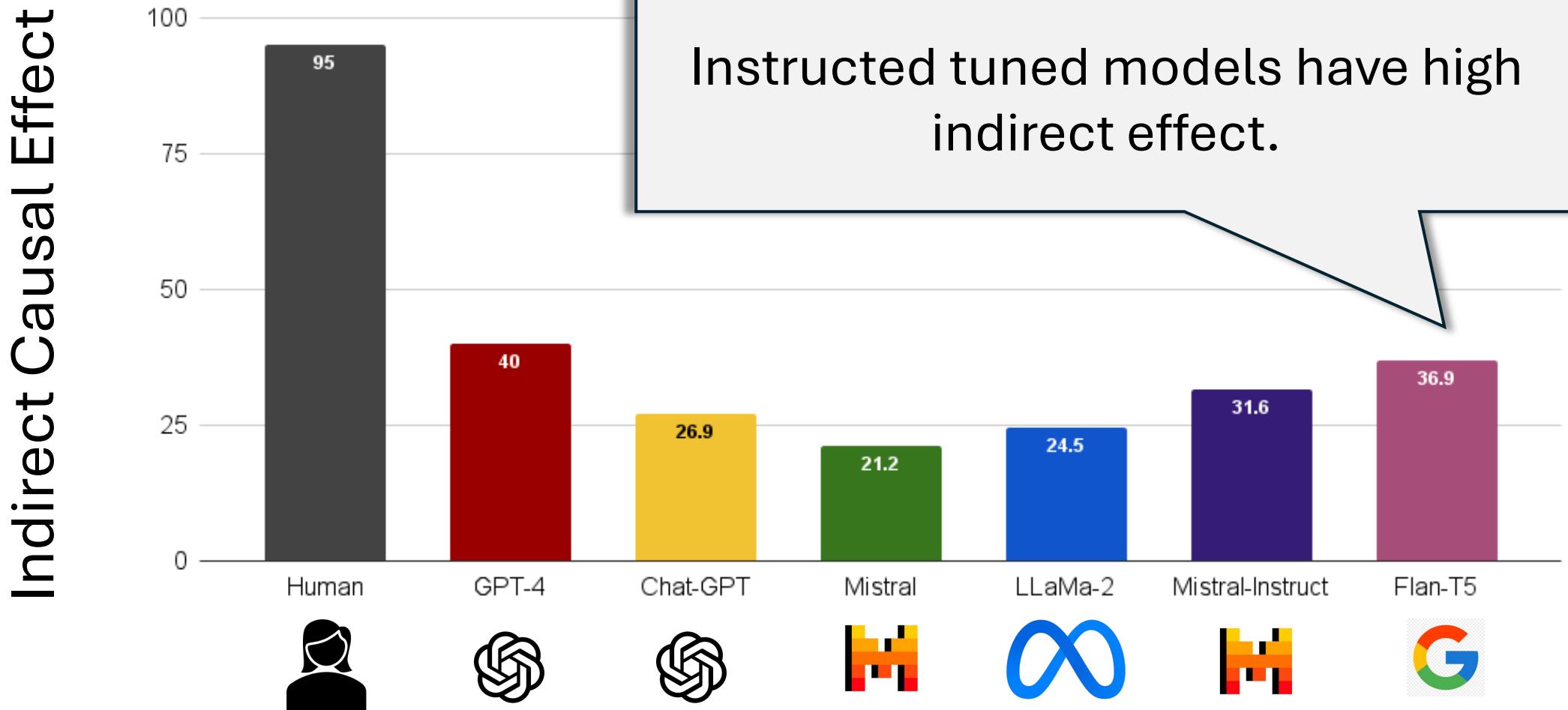


A high indirect effect means the model do emphasis on the reasoning steps

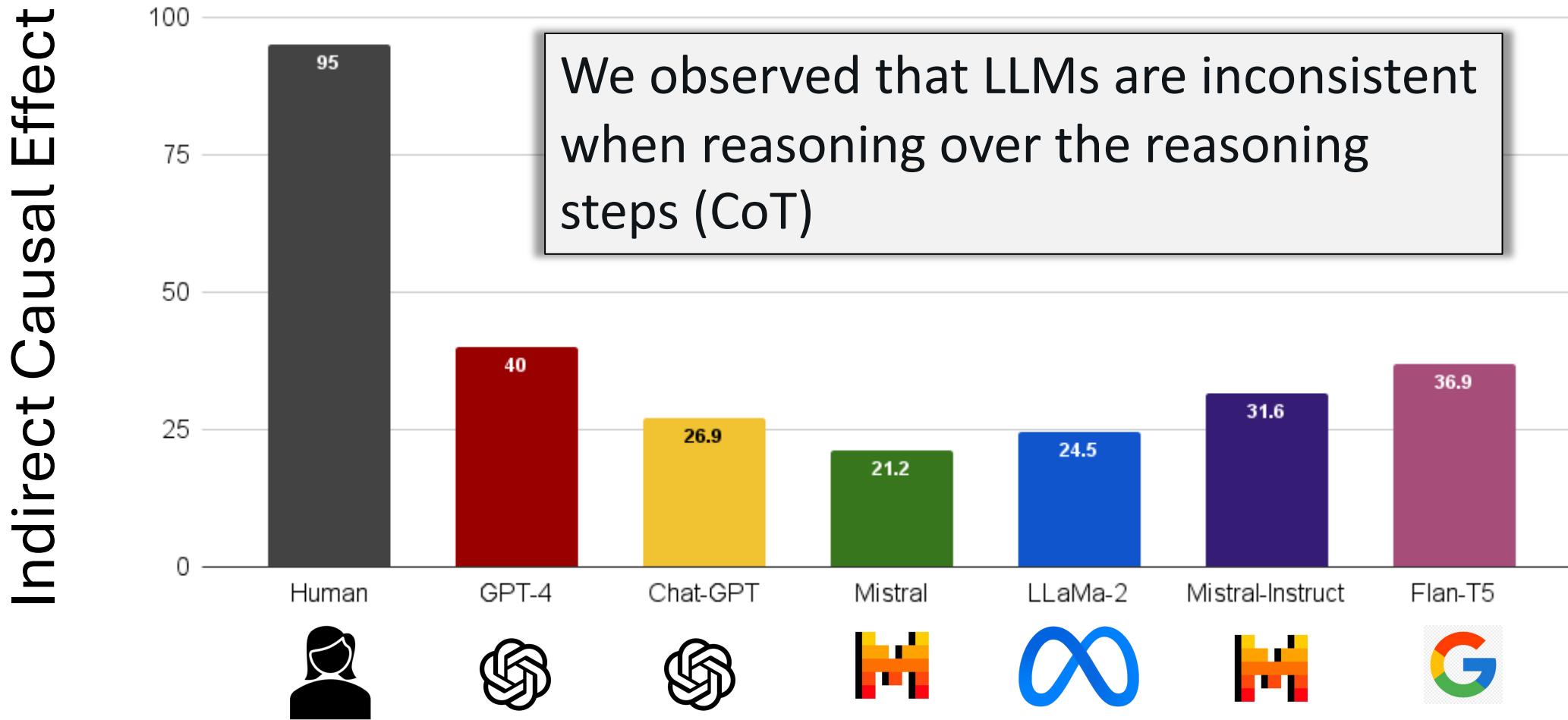
# How well closed-models perform?



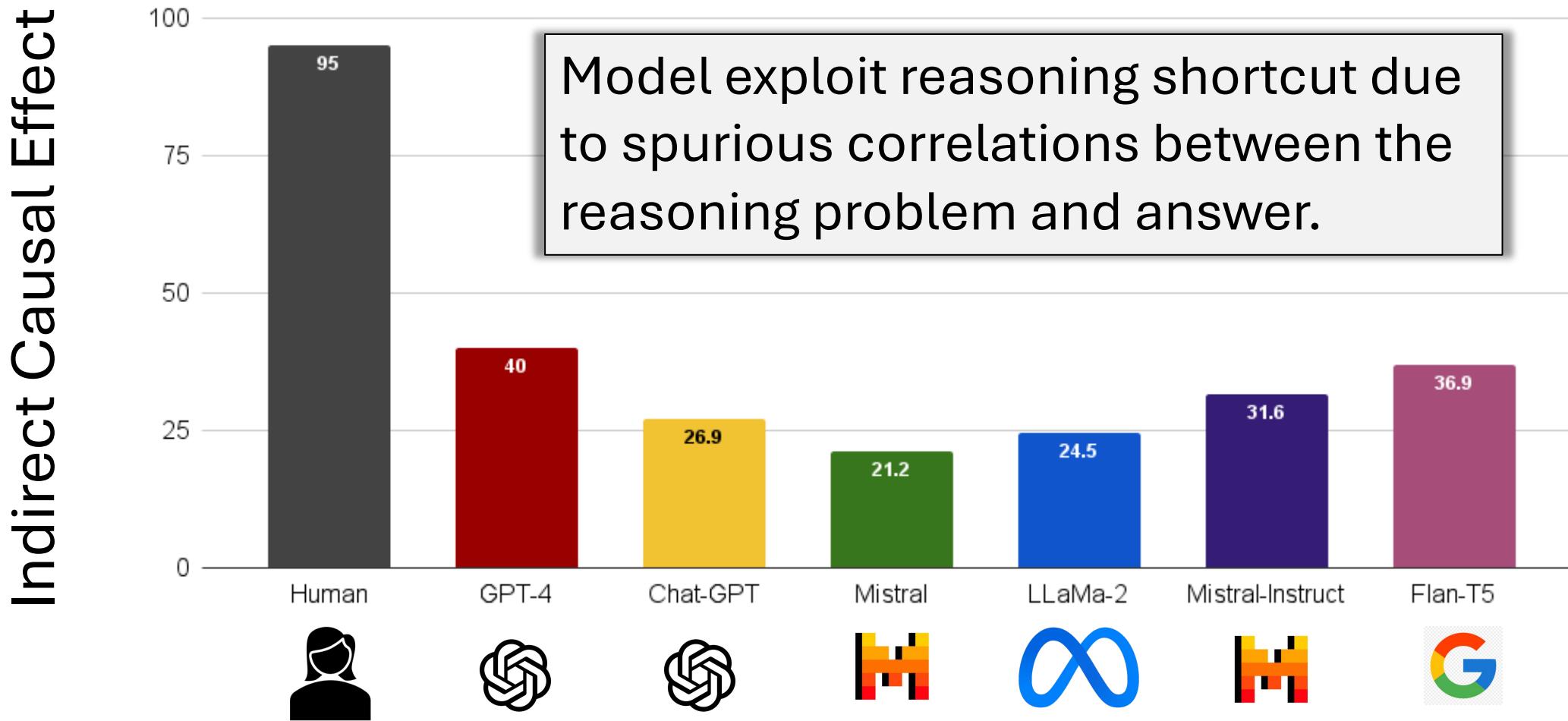
# How well open-sourced-models perform?



# How well open-sourced-models perform?



# How well open-sourced-models perform?



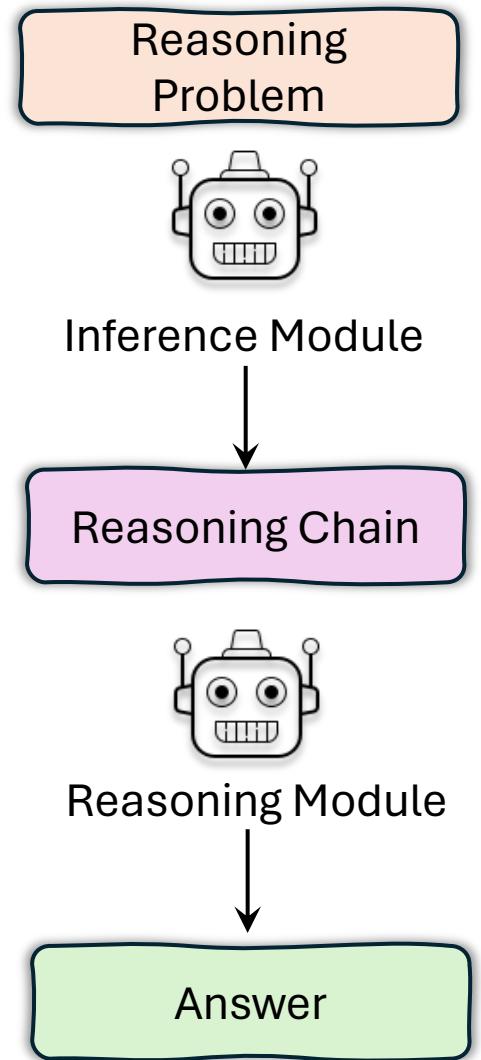
# Research Question II



How do we build a model that can faithfully reason over the inference chain and arrive to a correct answer?

# FRODO

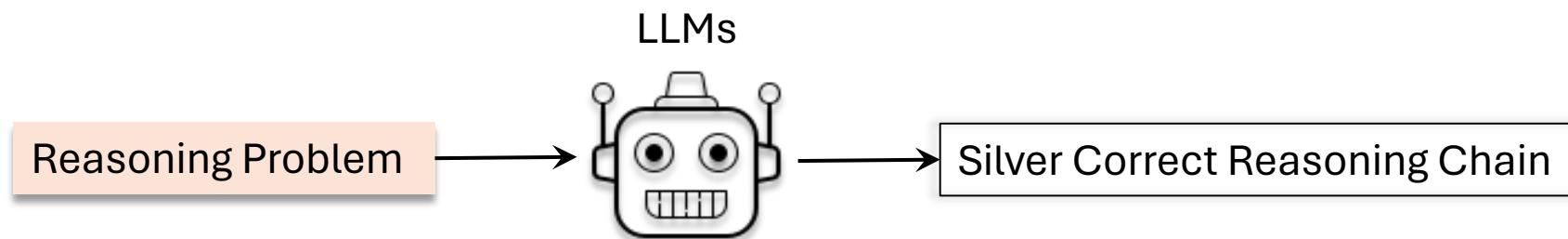
- FRODO comprises of two modules:
  - Inference module
    - generate correct reasoning chain
  - Reasoning module
    - takes the reasoning chains as input and faithfully reasons over them to arrive at the correct answer



# Inference Module

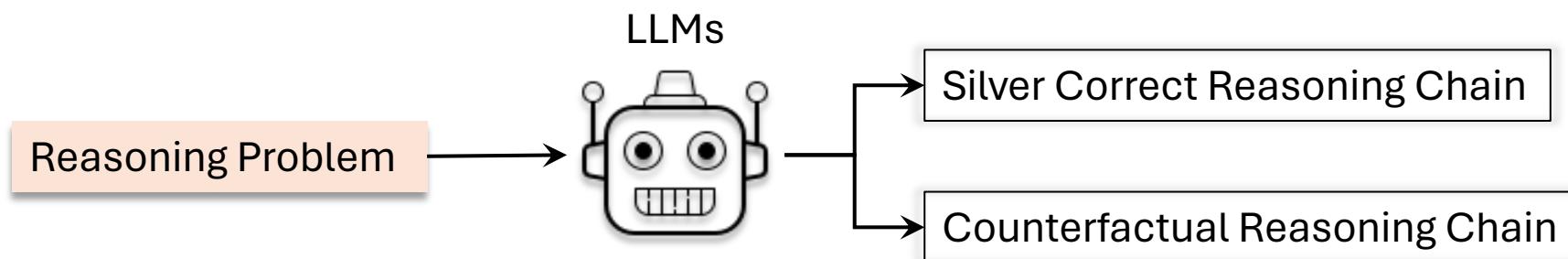
- Obtaining human-crafted golden reasoning chain is both time-consuming and costly.
- Automatically obtain the silver rationale from LLM (ChatGPT) using in-context learning.

No gold reasoning chain

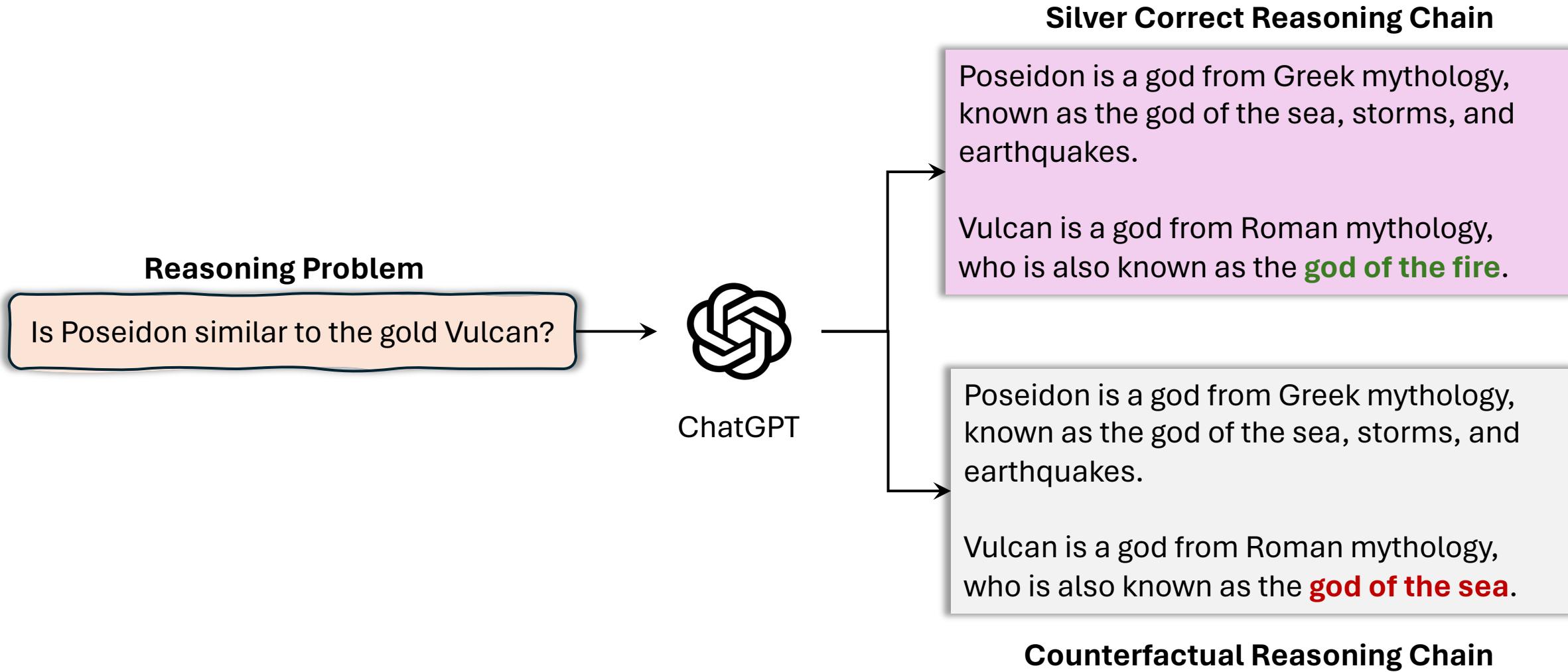


# Inference Module

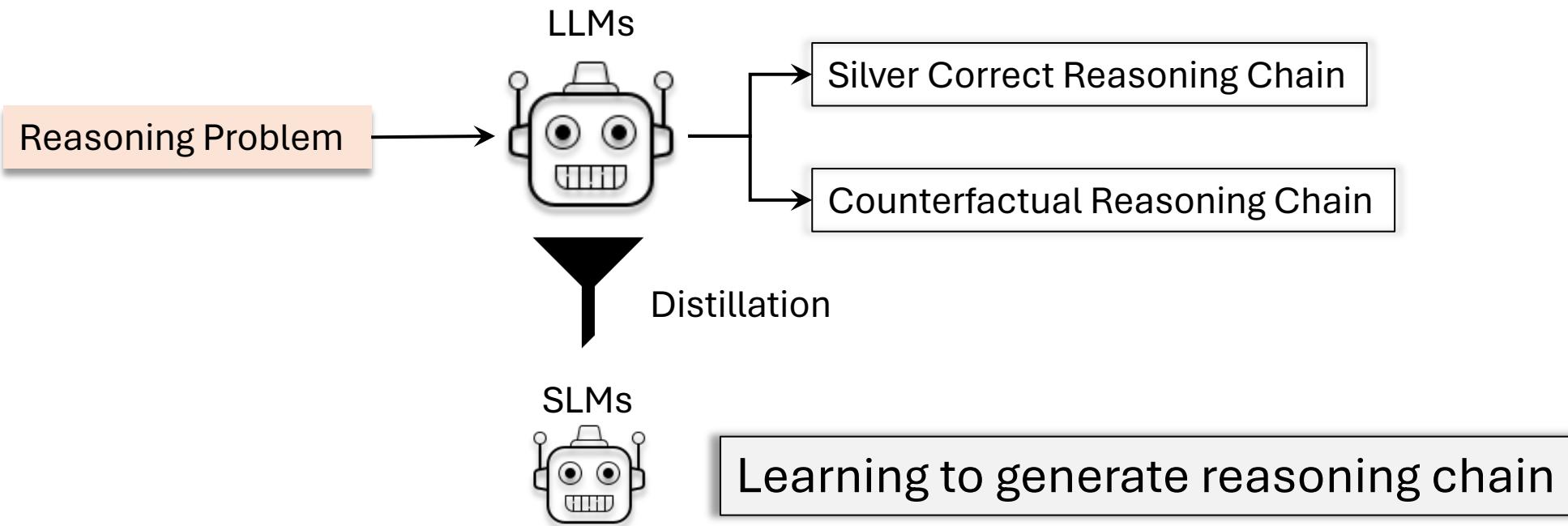
- Automatically obtain the silver rationale from LLM (ChatGPT) using in-context learning.
- **Hypothesis:** Learning to prefer correct reasoning chains over counterfactual chains can make models more robust and enhance generalization.



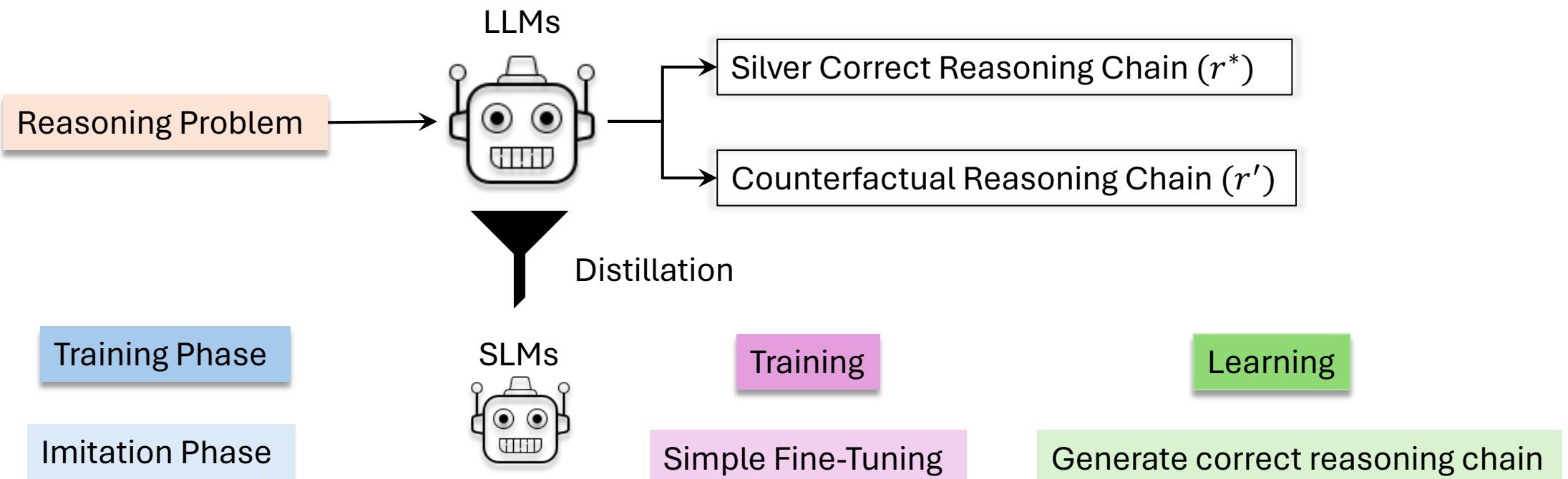
# Example of Reasoning Chain



# Inference Module

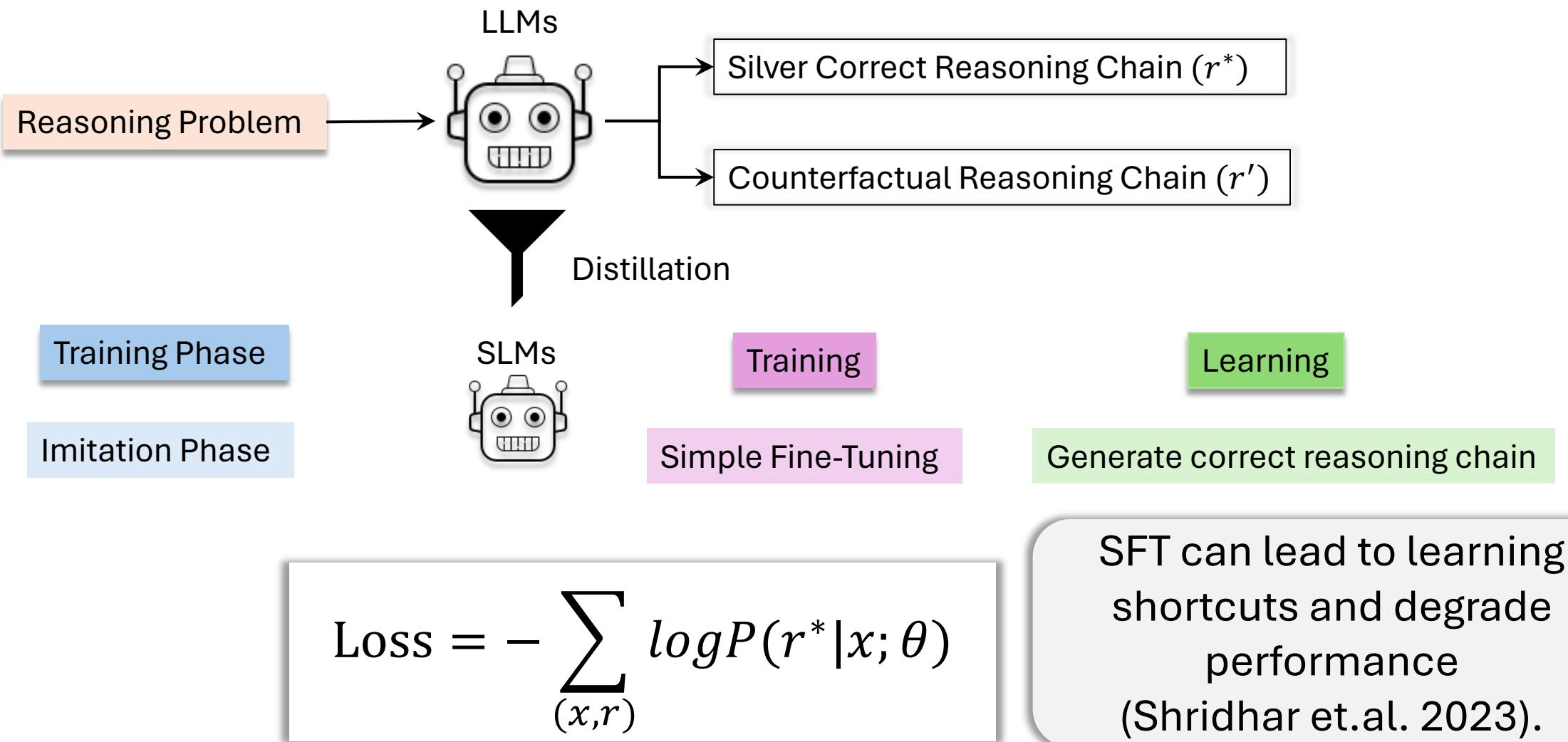


# Inference Module

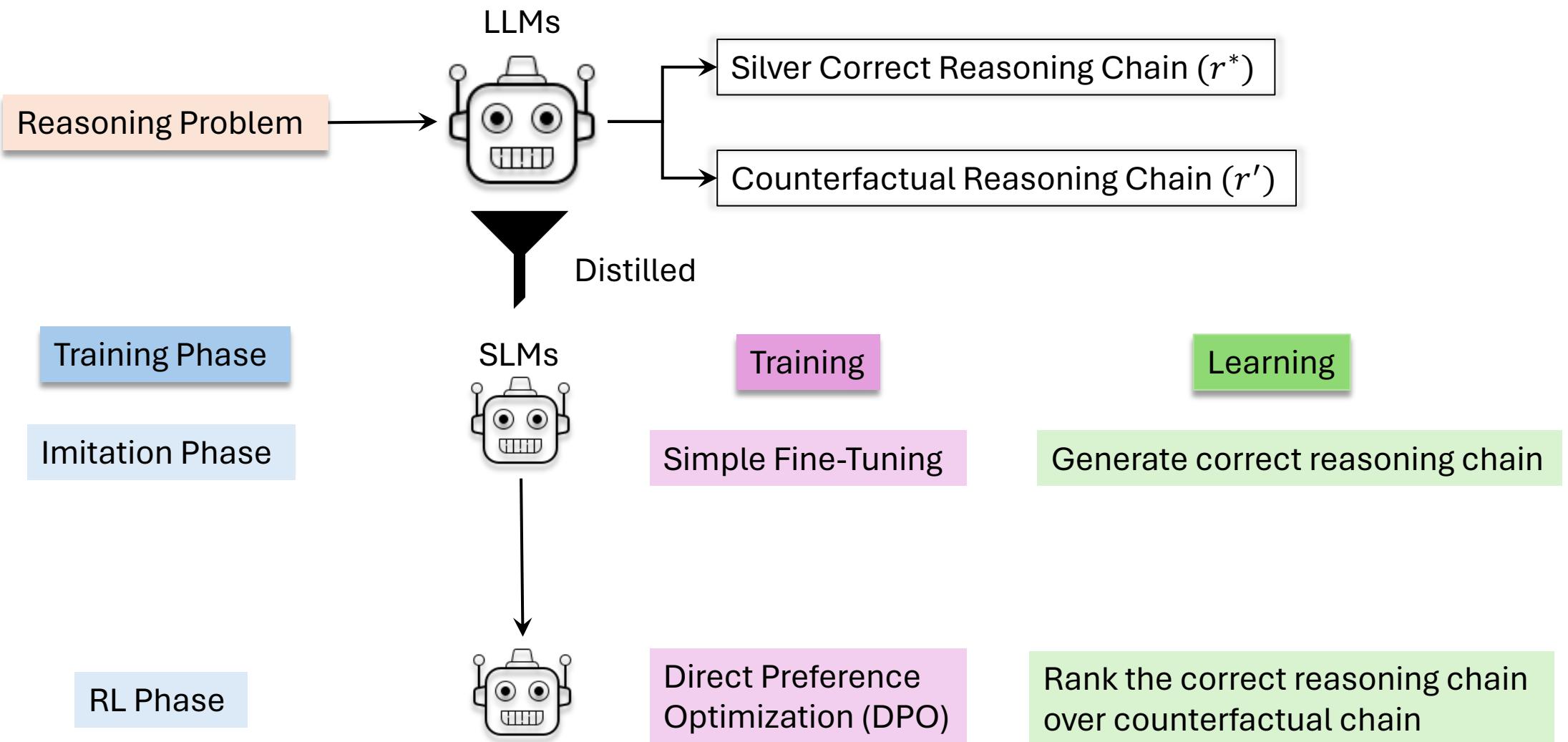


$$\text{Loss} = - \sum_{(x,r)} \log P(r^* | x; \theta)$$

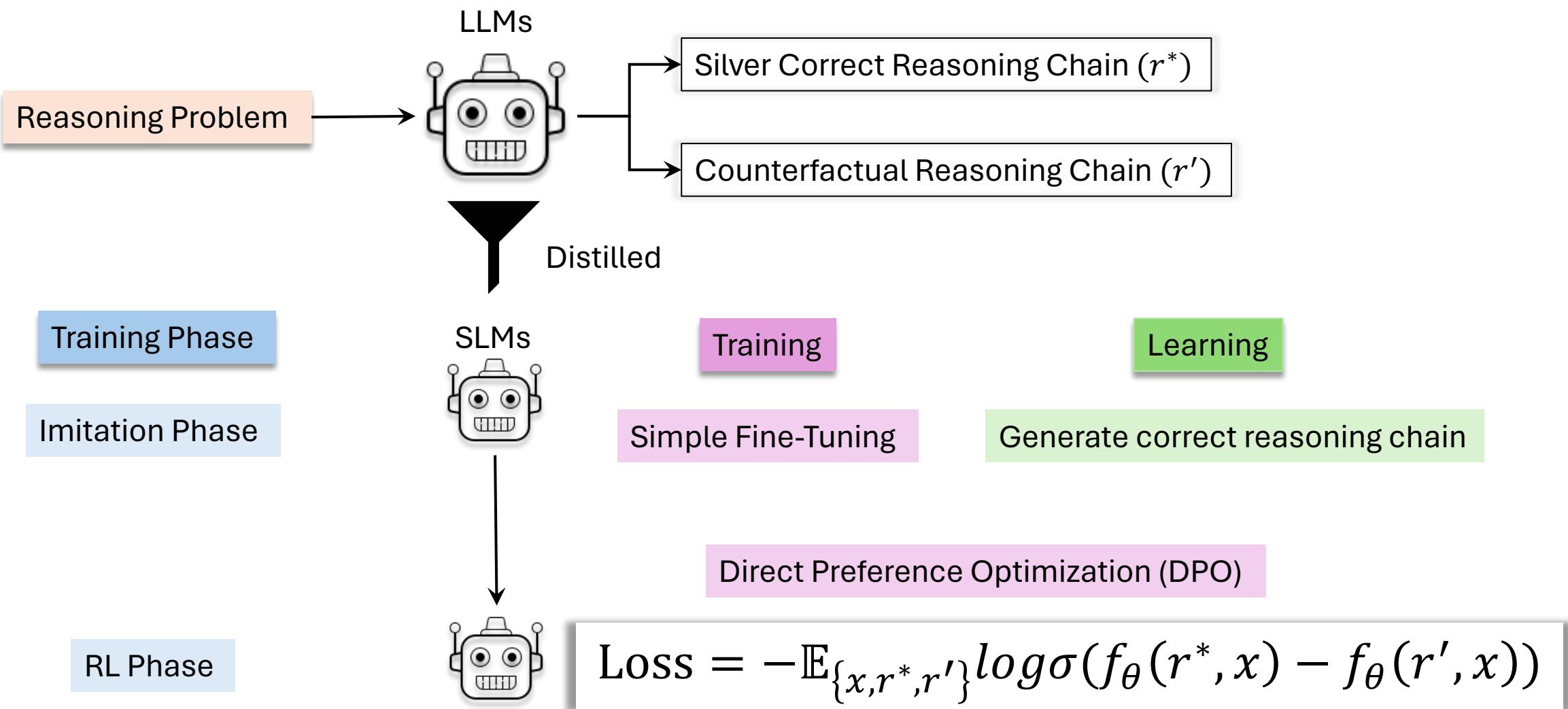
# Inference Module



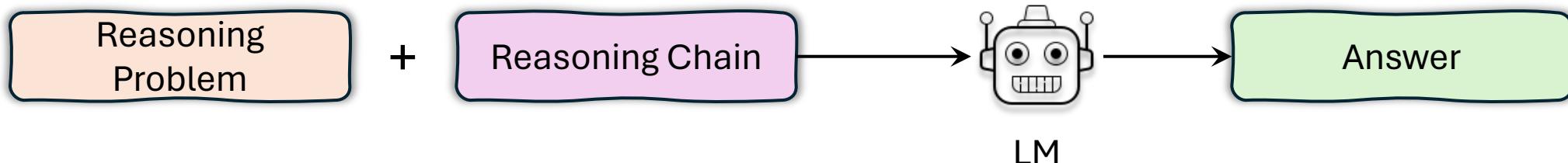
# Inference Module



# Inference Module



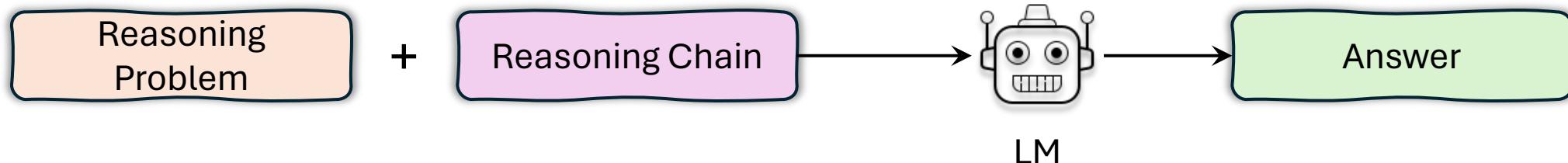
# Reasoning Module



- **Goal:** Improve the synergy between the reasoning chain and the final answer.
- To train the reasoning module we used a linear combination of three losses:

$$\begin{aligned} \text{Loss} &= -\alpha \sum_i \log P(a^* | x, r^*) - \beta \sum_i \log P(a' | x, r') \\ &\quad + \gamma \max(0, t * h(x, r^*, a^*) - h(x, r', a') + m) \end{aligned}$$

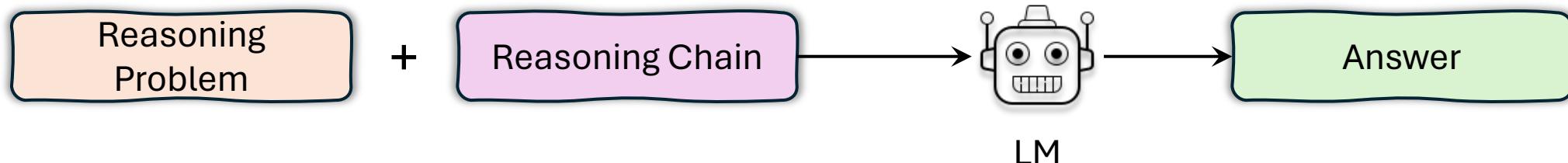
# Reasoning Module - Correctness



- **Goal:** Improve the synergy between the reasoning chain and the final answer.
- To train the reasoning module we used a linear combination of three losses:

$$\begin{aligned} \text{Loss} & \quad \text{language model loss} \\ &= -\alpha \sum_i \log P(a^* | x, r^*) - \beta \sum_i \log P(a' | x, r') \\ &+ \gamma \max(0, t * h(x, r^*, a^*) - h(x, r', a') + m) \end{aligned}$$

# Reasoning Module - Faithfulness

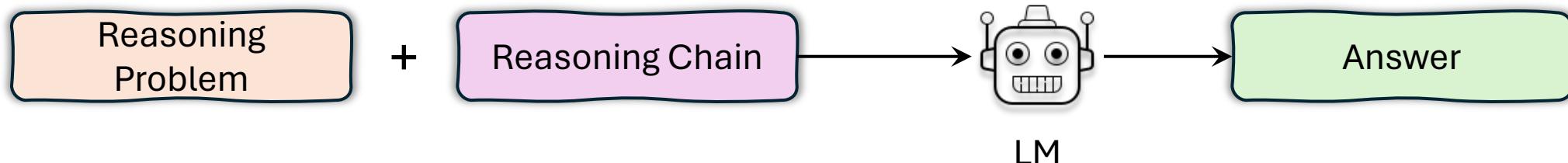


- **Goal:** Improve the synergy between the reasoning chain and the final answer.
- To encourage our reasoner module to reason faithfully over the reasoning steps, we use counterfactual loss.

$$\begin{aligned} \text{Loss} &= -\alpha \sum_i \log P(a^* | x, r^*) - \beta \sum_i \log P(a' | x, r') \\ &\quad + \gamma \max(0, t * h(x, r^*, a^*) - h(x, r', a') + m) \end{aligned}$$

The equation shows the total Loss function. It consists of three parts: 1) language model loss, represented by the term  $-\alpha \sum_i \log P(a^* | x, r^*)$ , where the summation is over steps  $i$ ; 2) counterfactual loss, represented by the term  $-\beta \sum_i \log P(a' | x, r')$ , also with summation over steps  $i$ ; and 3) a third term involving a max function and parameters  $t$ ,  $h$ , and  $m$ .

# Reasoning Module – Robustness

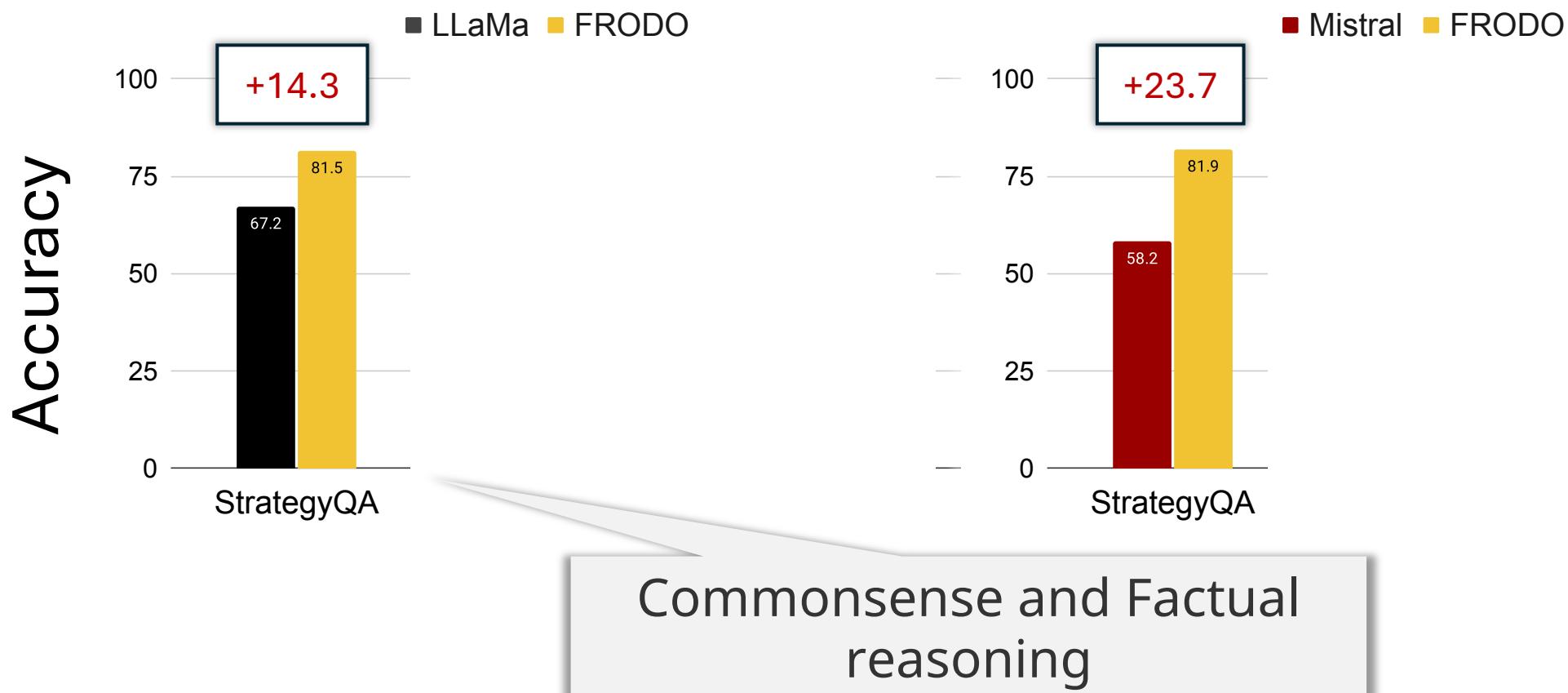


- **Goal:** Improve the synergy between the reasoning chain and the final answer.
- To encourage our reasoner module to improve model robustness and generalization against input variation.

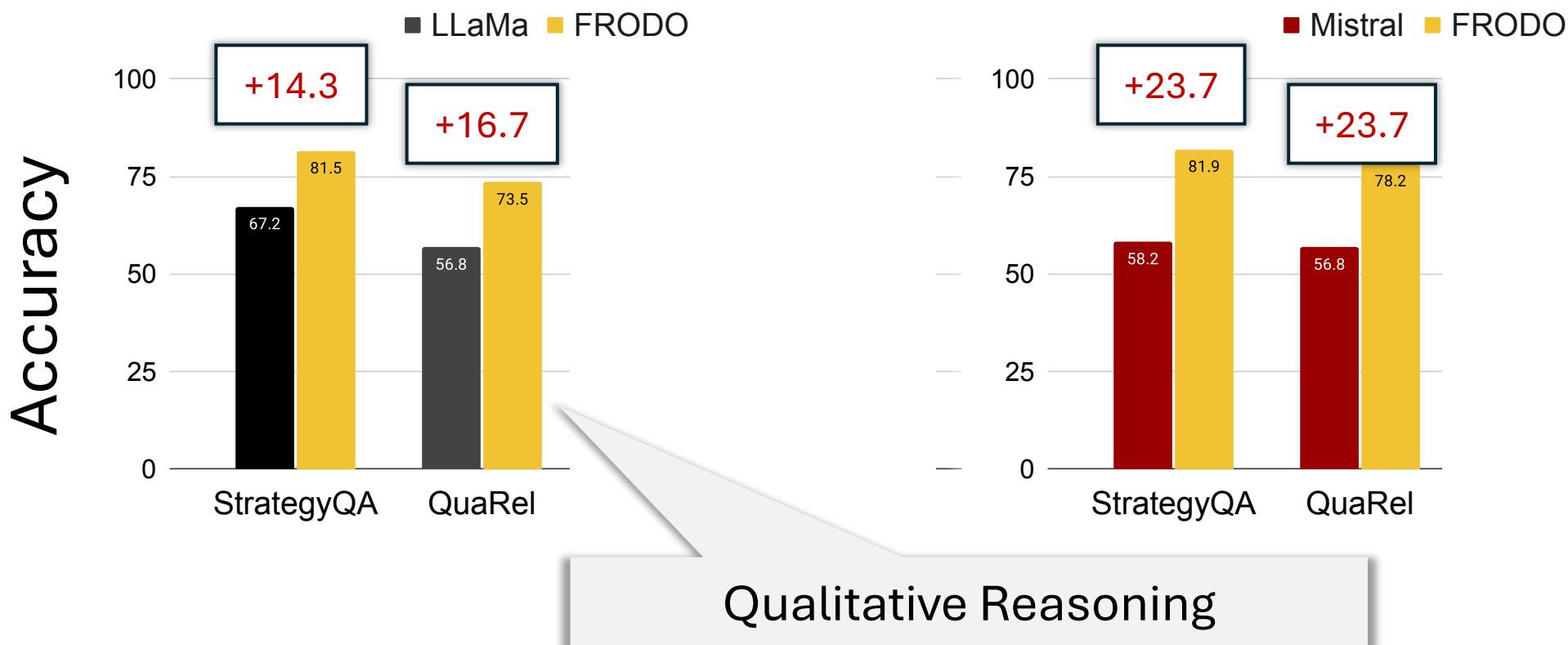
$$\begin{aligned} \text{Loss} &= -\alpha \sum_i \log P(a^* | x, r^*) - \beta \sum_i \log P(a' | x, r') \\ &\quad + \gamma \max(0, t * h(x, r^*, a^*) - h(x, r', a') + m) \end{aligned}$$

margin ranking loss

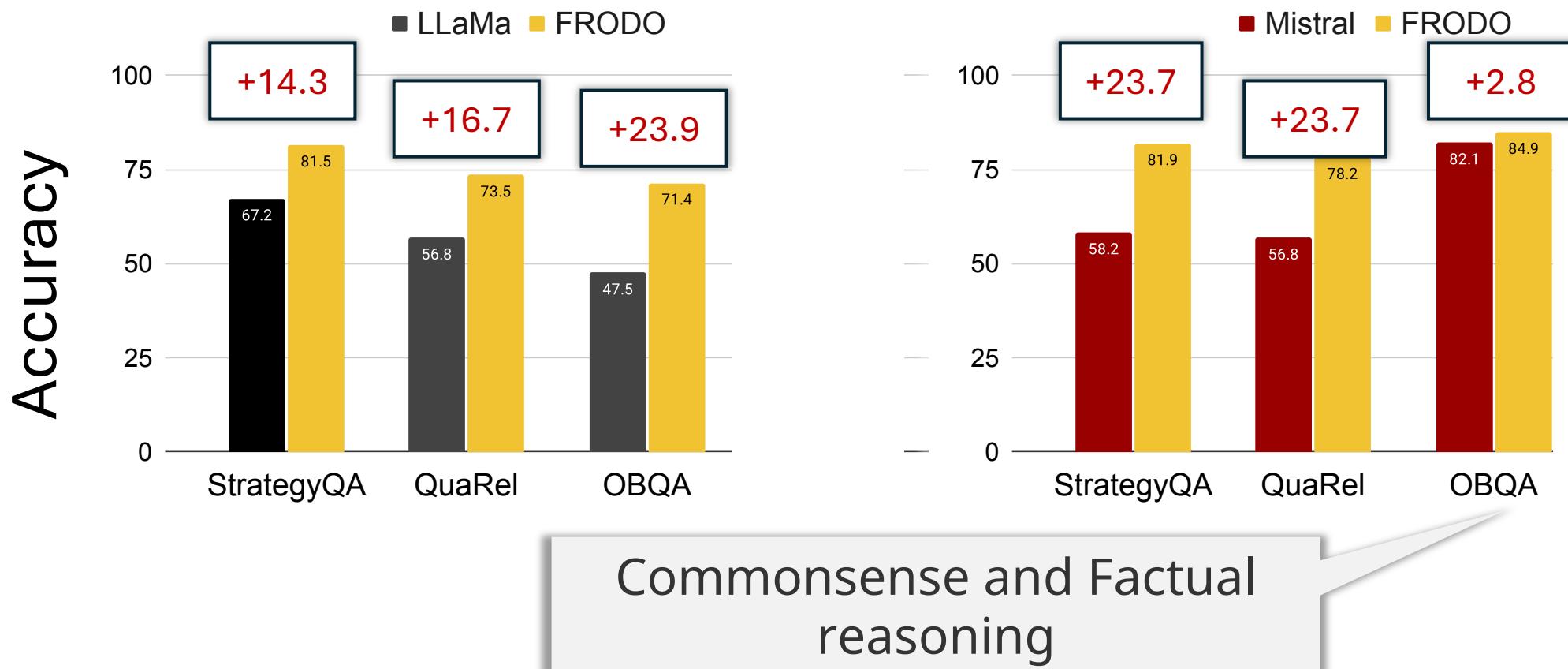
# How well does FRODO work?



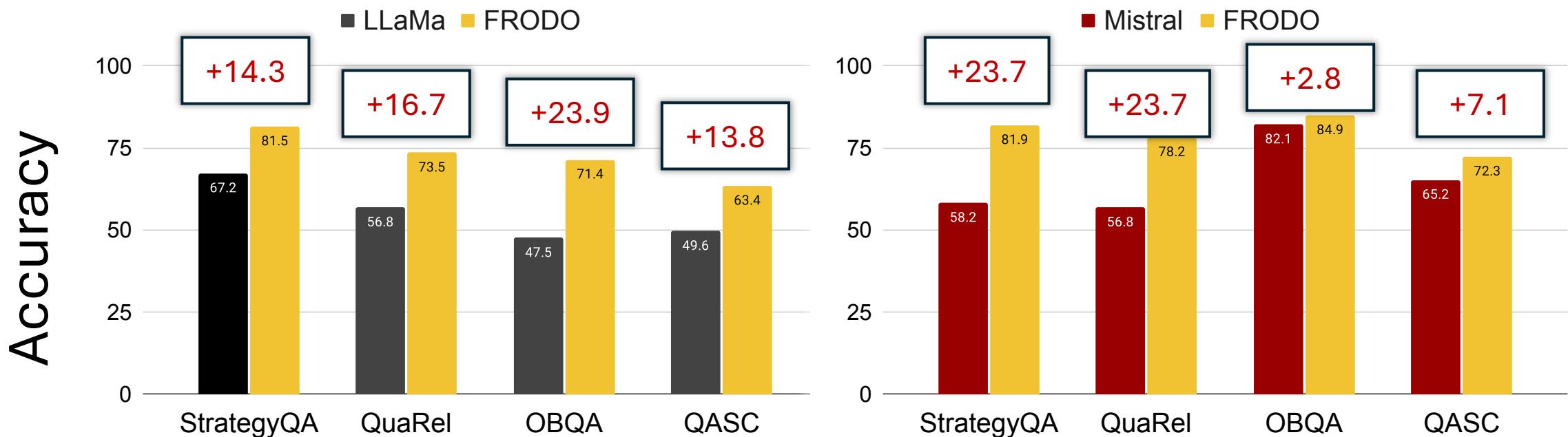
# How well does FRODO work?



# How well does FRODO work?

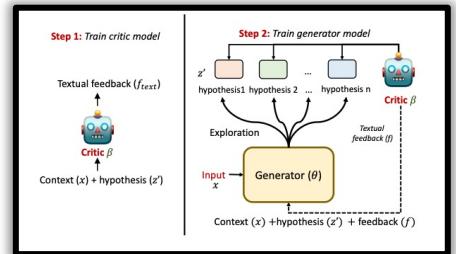


# How well does FRODO work?



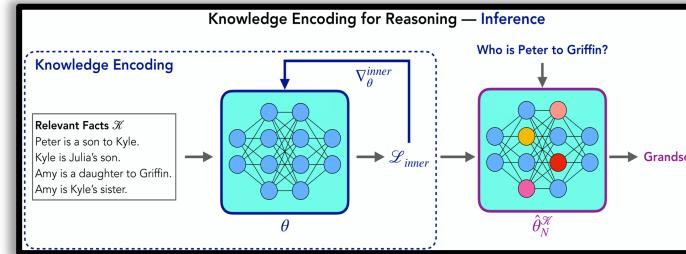
Reasoning over Science Facts

# Outline



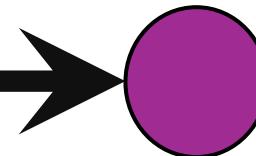
Refine knowledge using feedback

Paul et.al. 2024 (EACL)

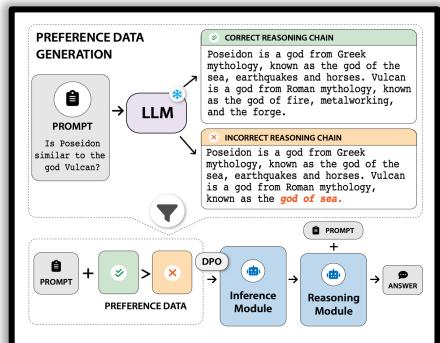


Robust Reasoning by updating world models

Chen et.al. 2023 (Neurips)



Paul et.al. 2024 (EMNLP Findings)  
Faithful Reasoning about knowledge

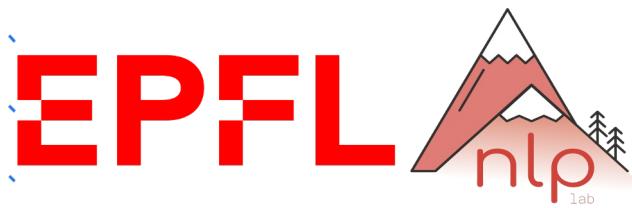


# Recap

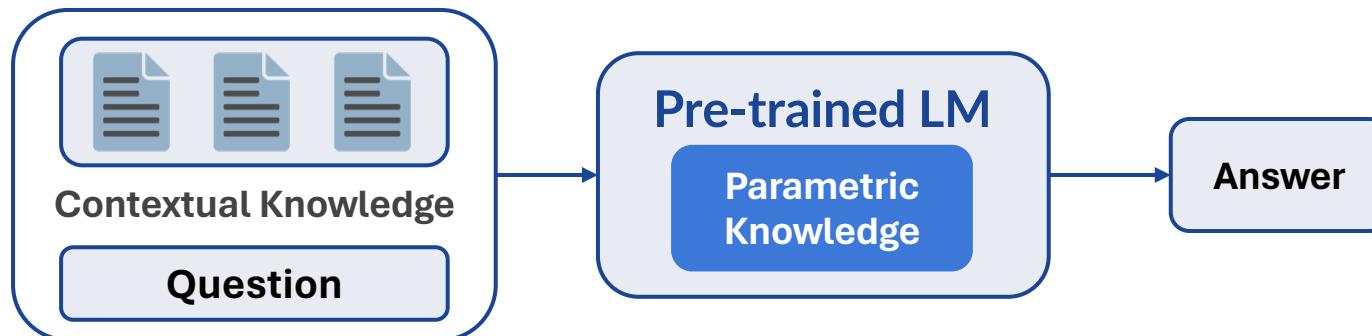
- Interactive method improves the quality of inference generation
- Causal mediation analysis show that LLMs such as GPT-4, ChatGPT etc. are inconsistent when reasoning over their own generated reasoning traces (CoT).
- Counterfactual reasoning can improve the faithfulness and performance of LMs.



# RECKONING: Reasoning through Dynamic Knowledge Encoding



# Background: PLM In-context Reasoning



(*Input Facts:*) Alan is blue. Alan is rough. Alan is young.  
Bob is big. Bob is round.

Charlie is big. Charlie is blue. Charlie is green.  
Dave is green. Dave is rough.

(*Input Rules:*) Big people are rough.  
If someone is young and round, then they are kind.  
If someone is round and big, then they are blue.  
All rough people are green.

Q1: Bob is green. True/false? [Answer: T]

Q2: Bob is kind. True/false? [Answer: F]

Q3: Dave is blue. True/false? [Answer: F]

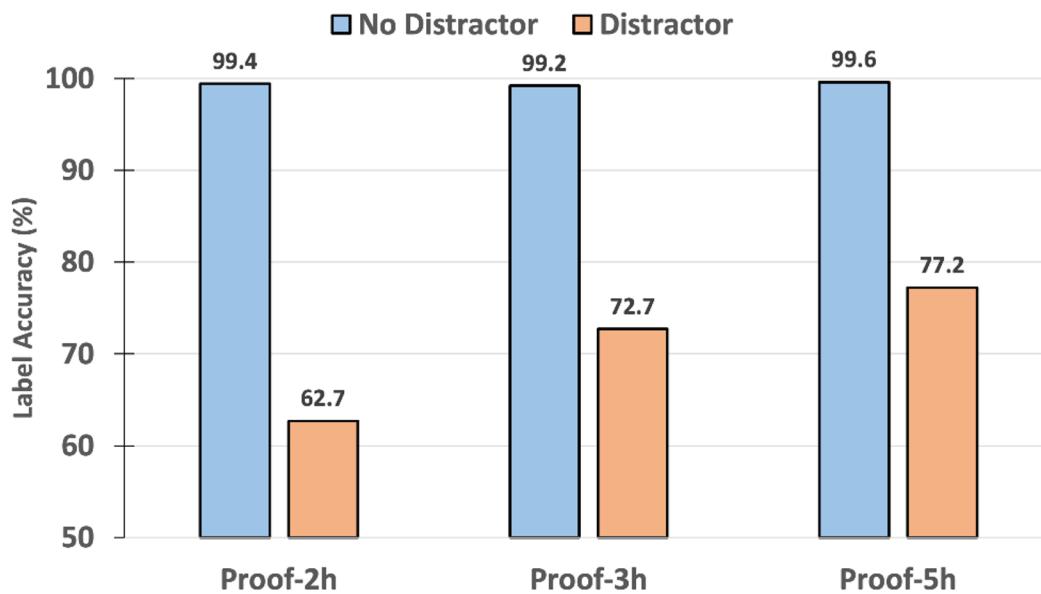
# Challenge

Real-world context contains **noisy information**

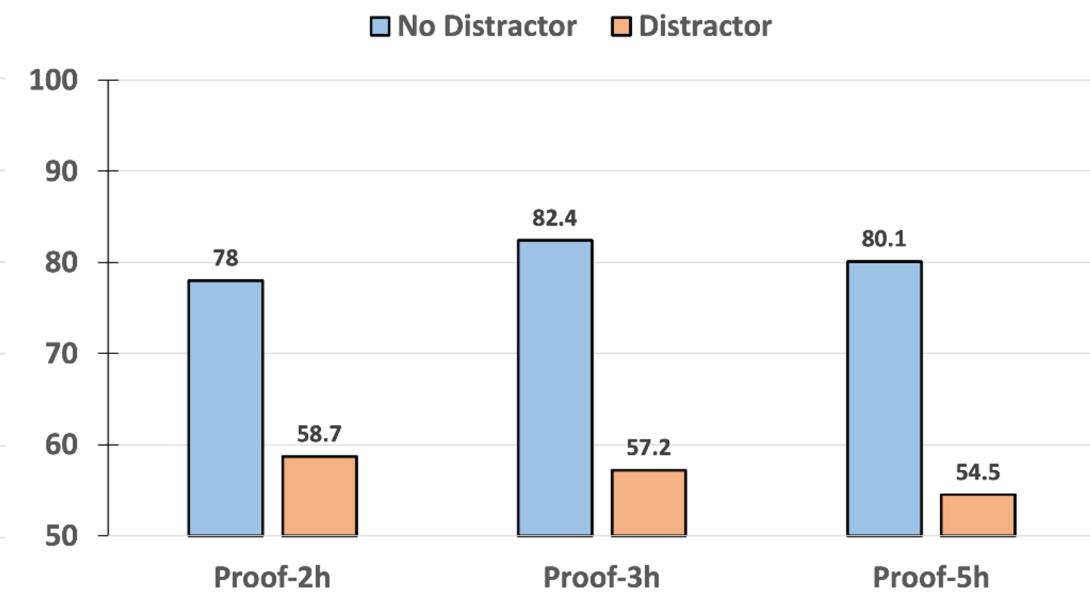
PLMs are **sensitive** to **irrelevant** information in the context

# Challenge

GPT-2-Small No Distractor vs. Distractor



GPT-3.5 (8-shot) No Distractor vs. Distractor



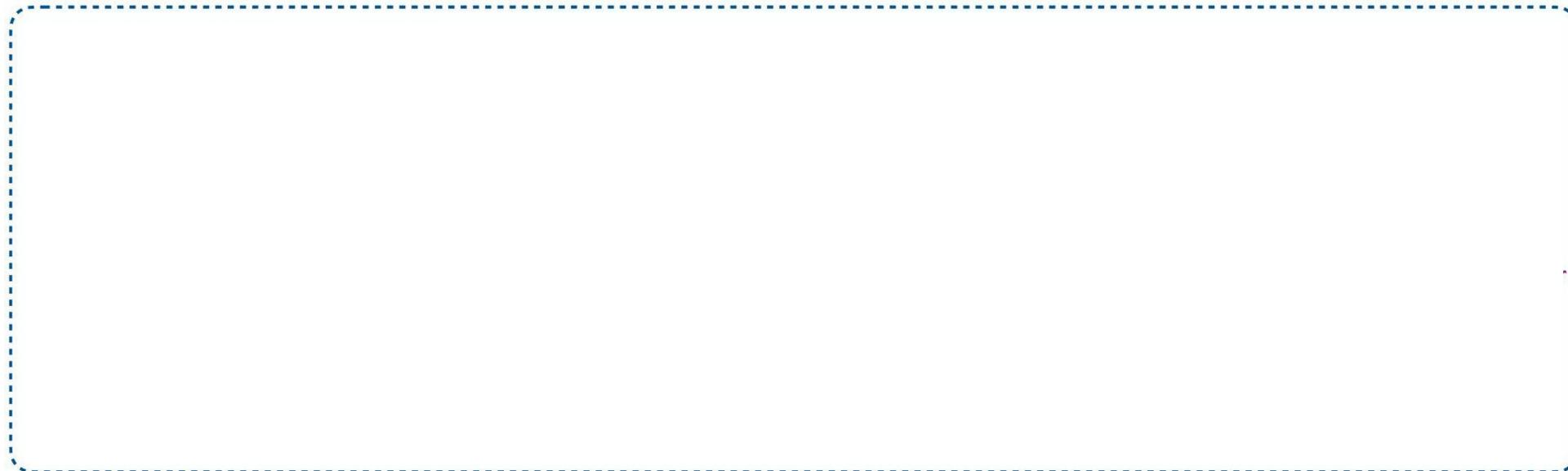
Language models with different parameter sizes  
are sensitive to distractors in the context.

# Research Question

**Can we mitigate noisy information by turning contextual knowledge into parametric knowledge through backpropagation?**

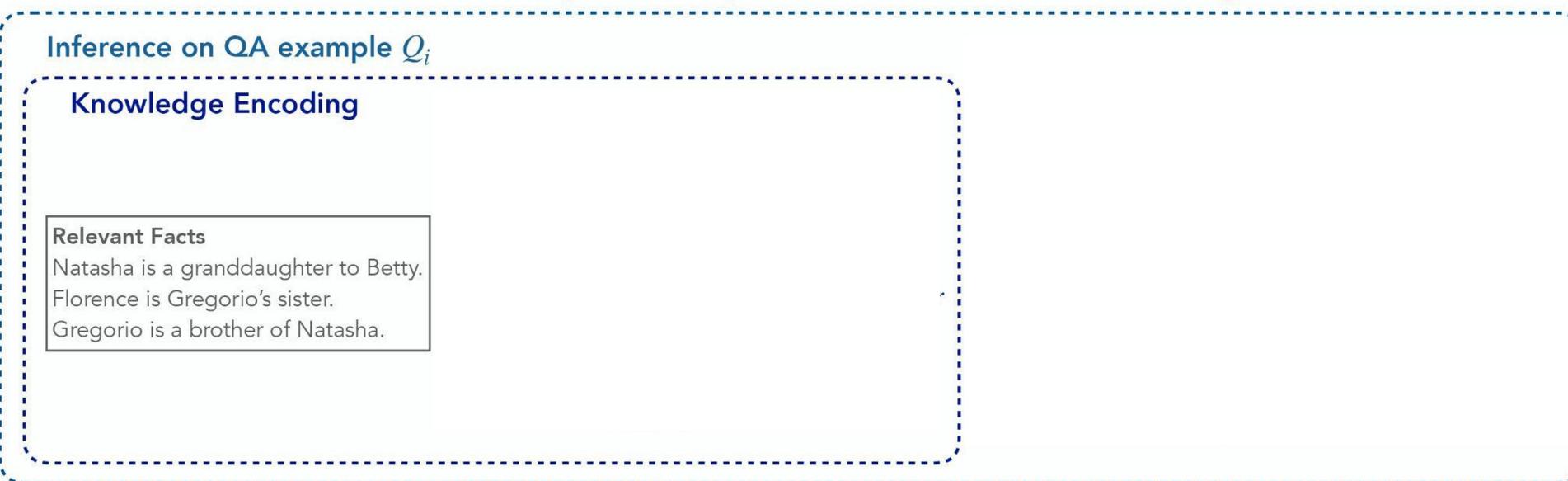
# RECKONING: Knowledge Encoding for Reasoning

Knowledge Encoding for Multi-hop QA — Meta Testing

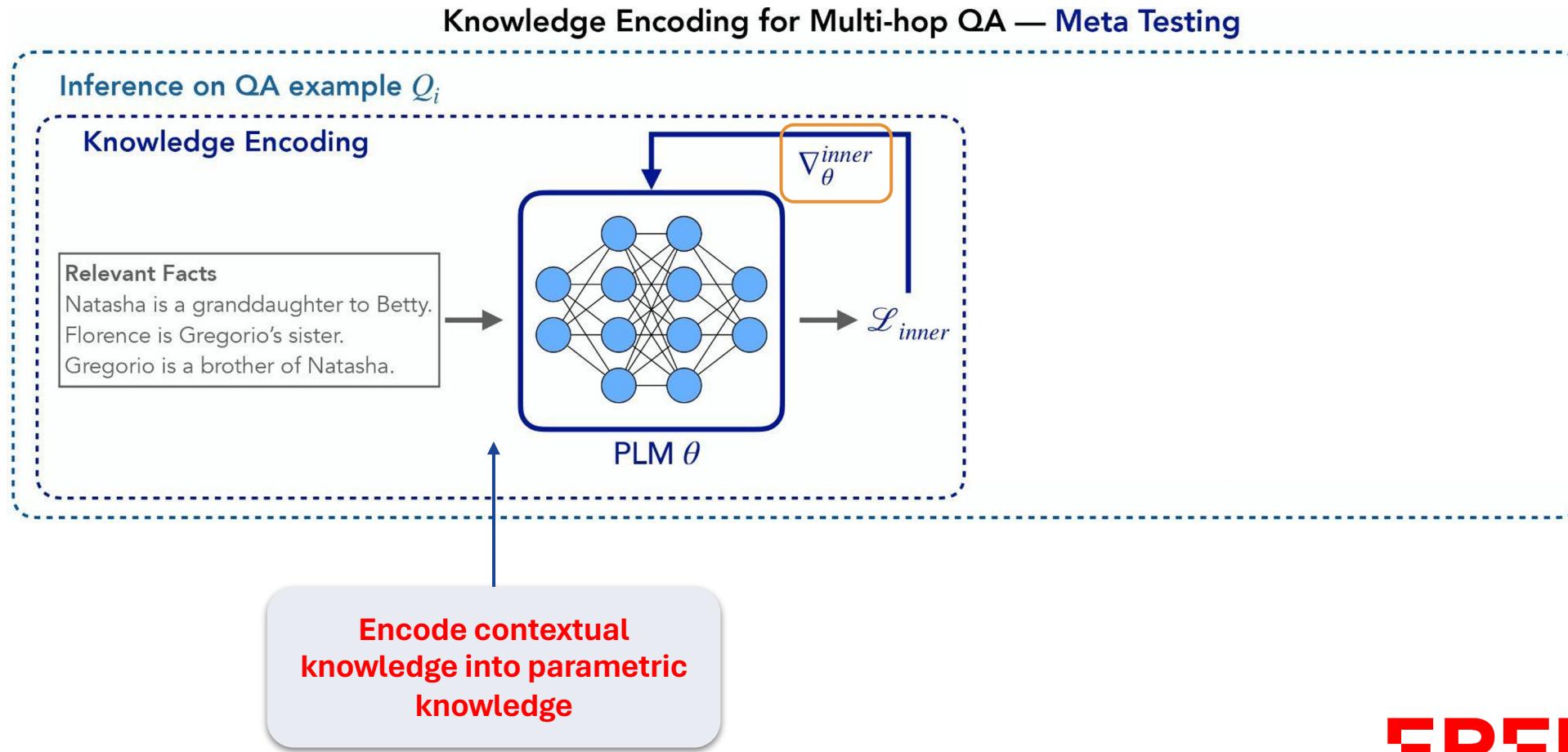


# RECKONING: Knowledge Encoding for Reasoning

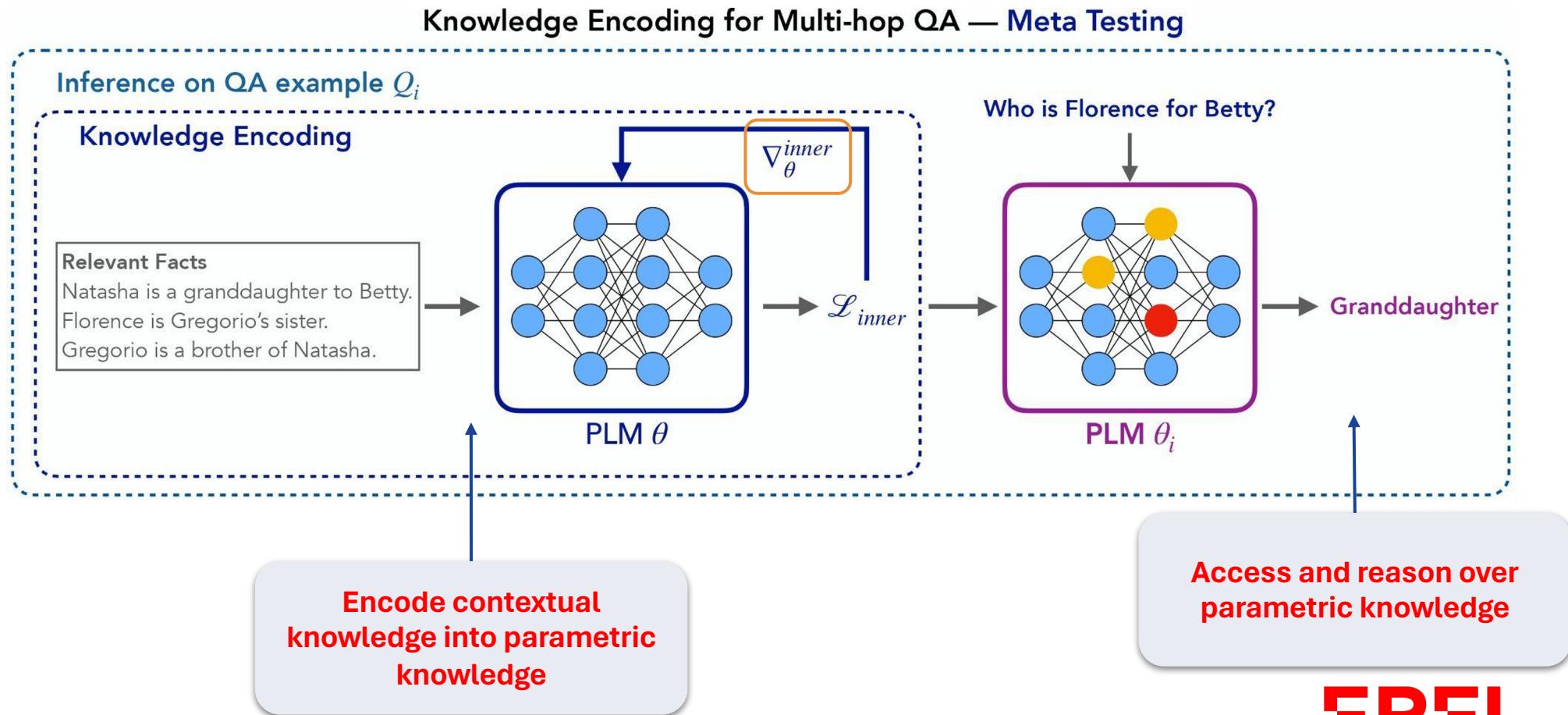
## Knowledge Encoding for Multi-hop QA — Meta Testing



# RECKONING: Knowledge Encoding for Reasoning

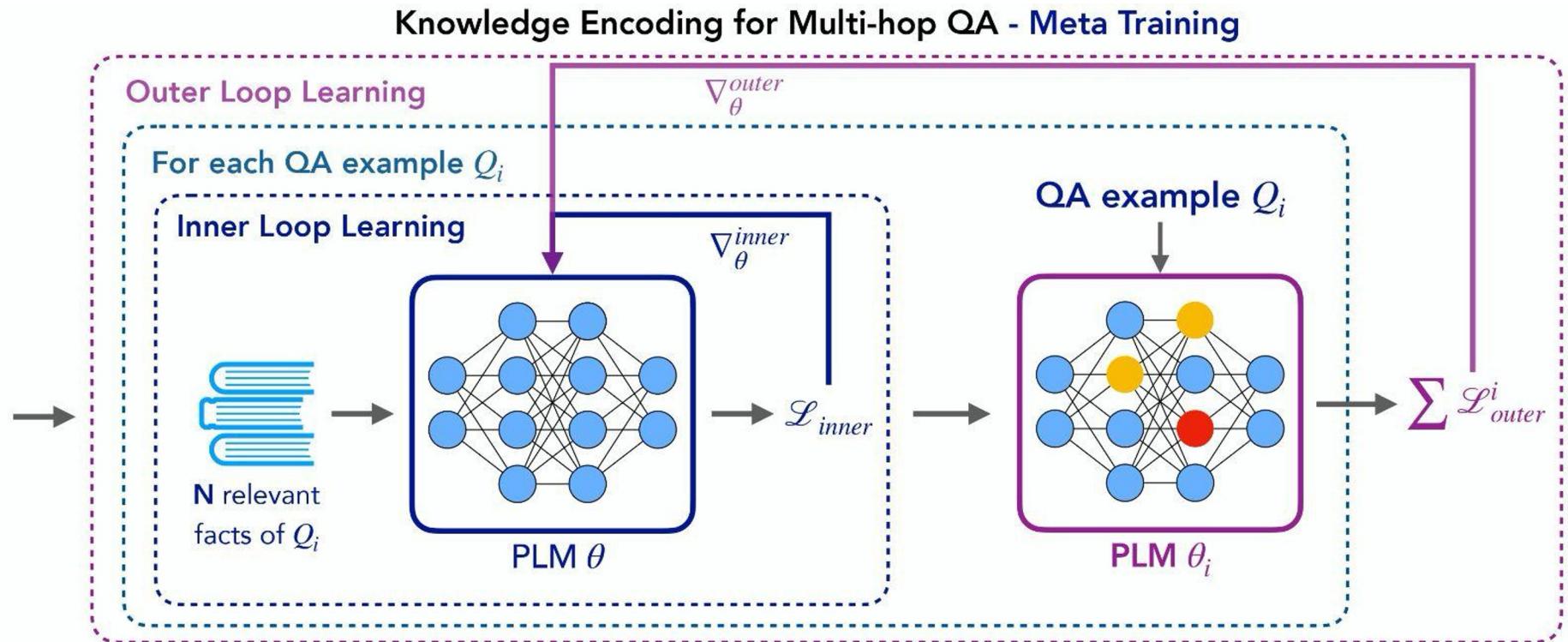
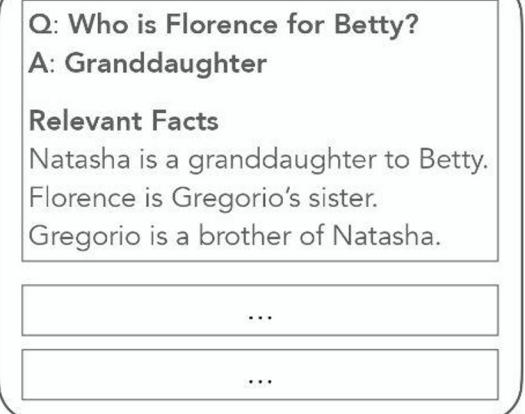


# RECKONING: Knowledge Encoding for Reasoning

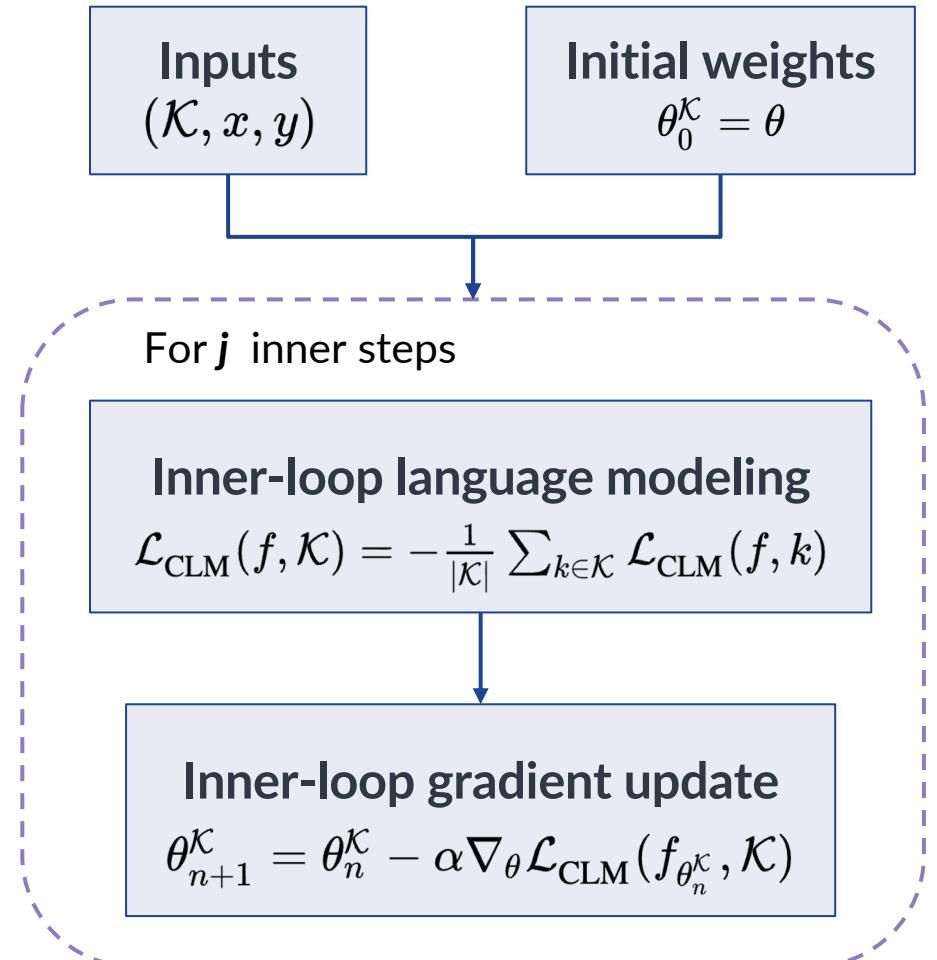


# Bi-level Optimization

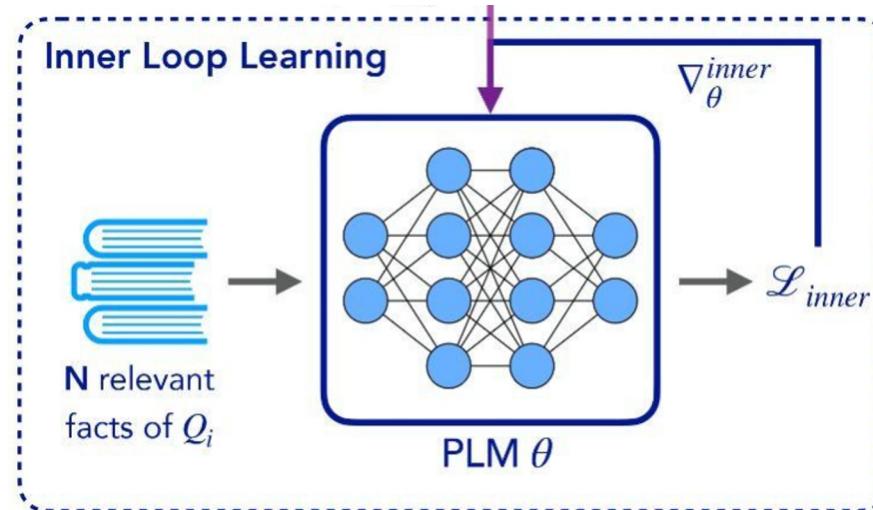
## Question Answering Examples



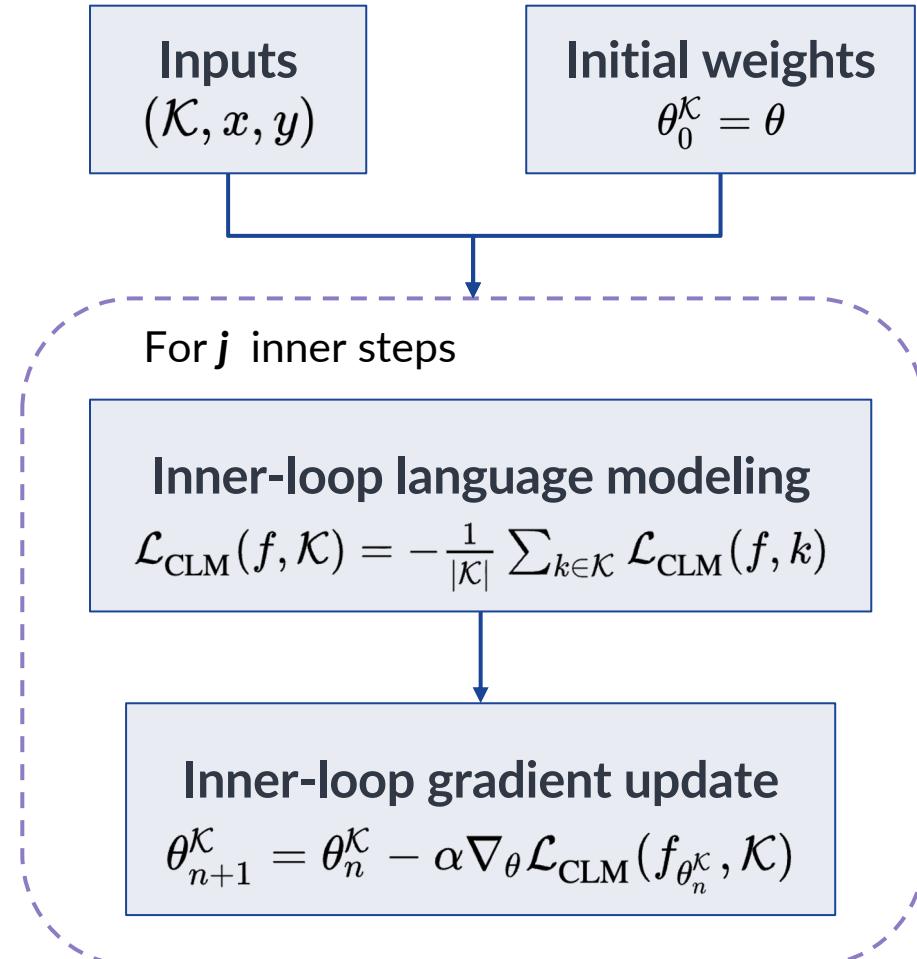
# Inner-loop Learning



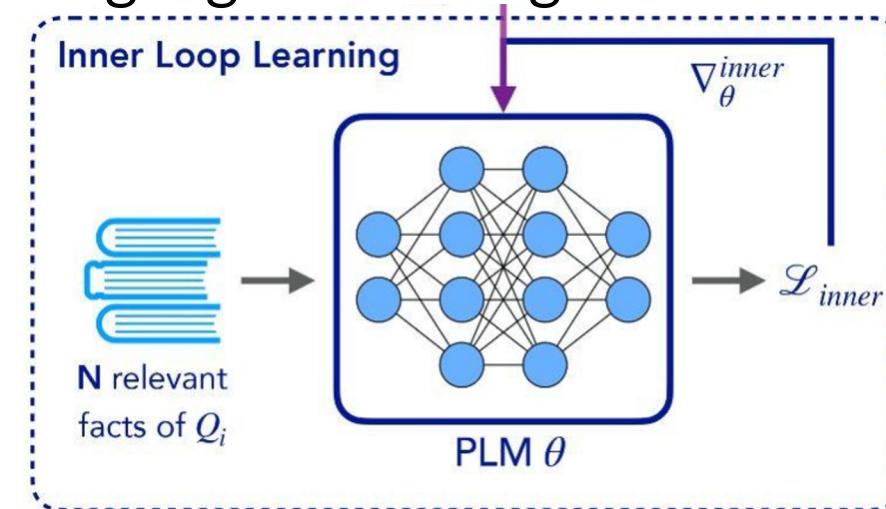
- ❖ Language modeling on natural language knowledge



# Inner-loop Learning

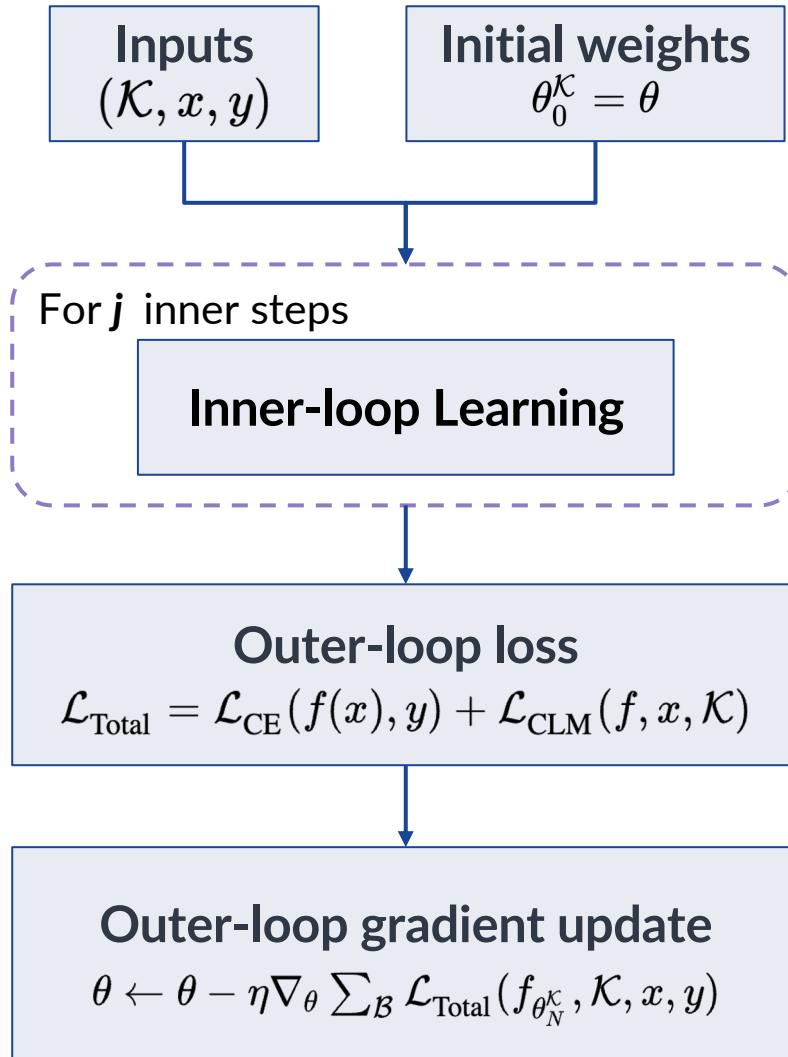


- ❖ **Language modeling on natural language knowledge**

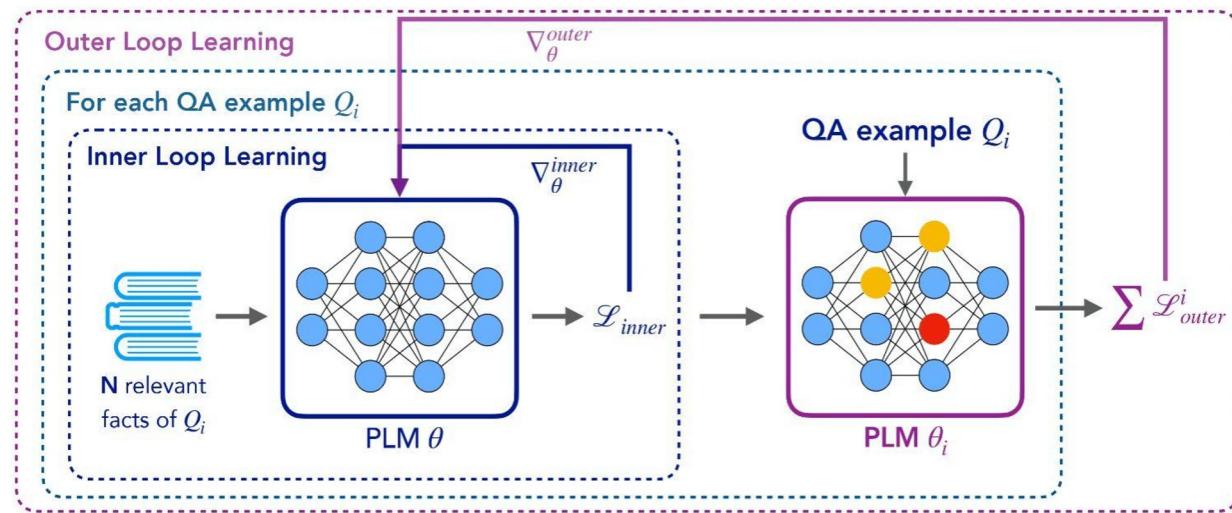


**Knowledge encoding through backpropagation**

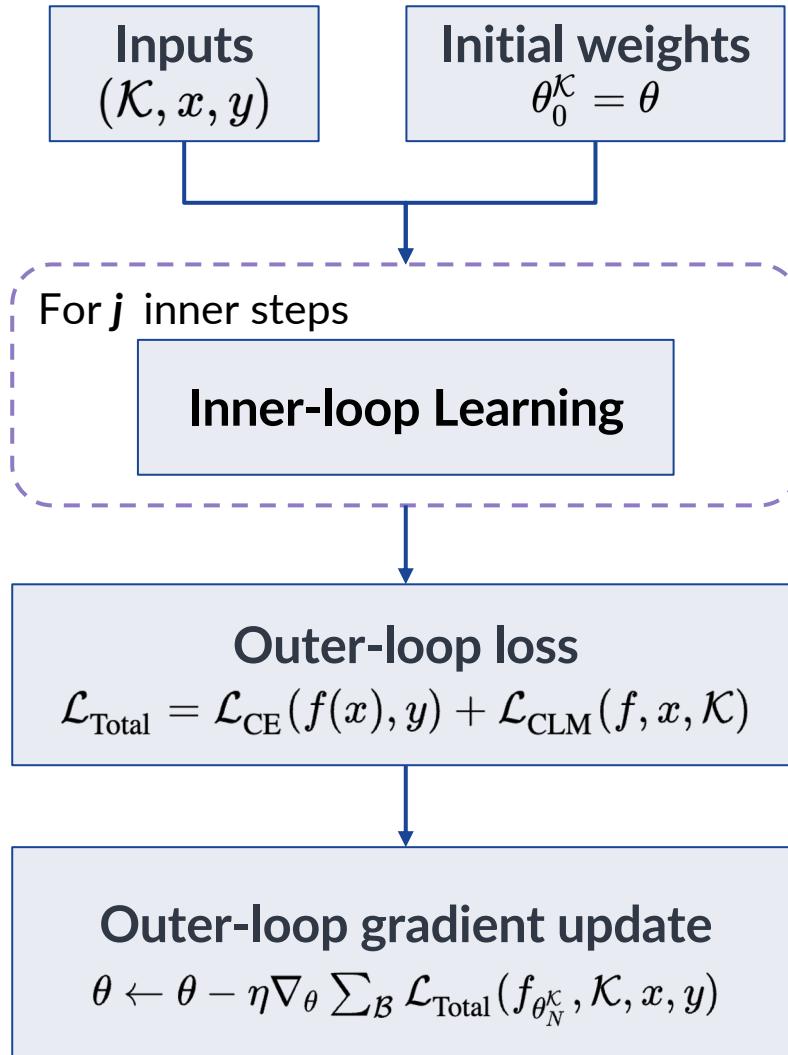
# Outer-loop Learning



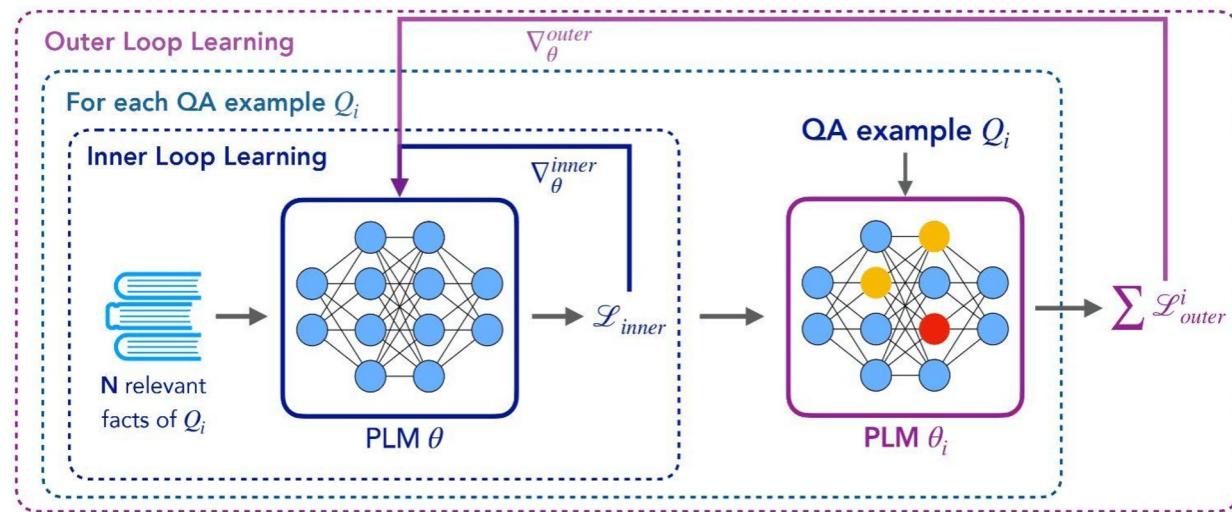
## ❖ Closed-Book question answering



# Outer-loop Learning

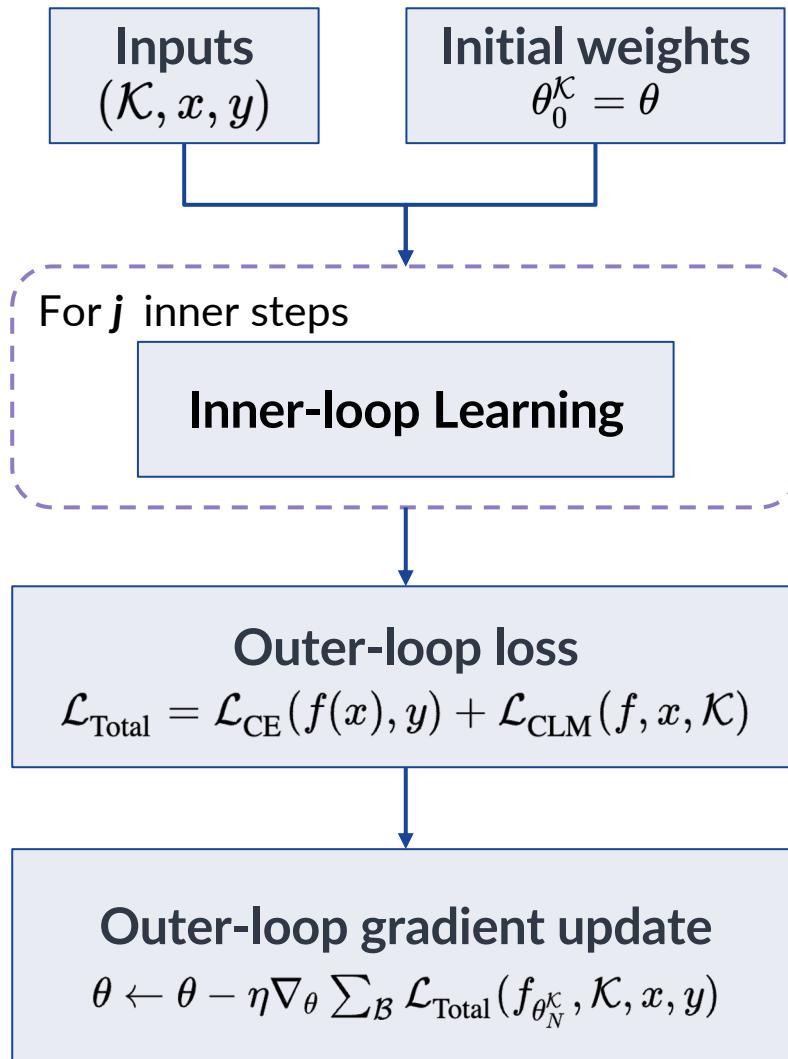


## ❖ Closed-Book question answering

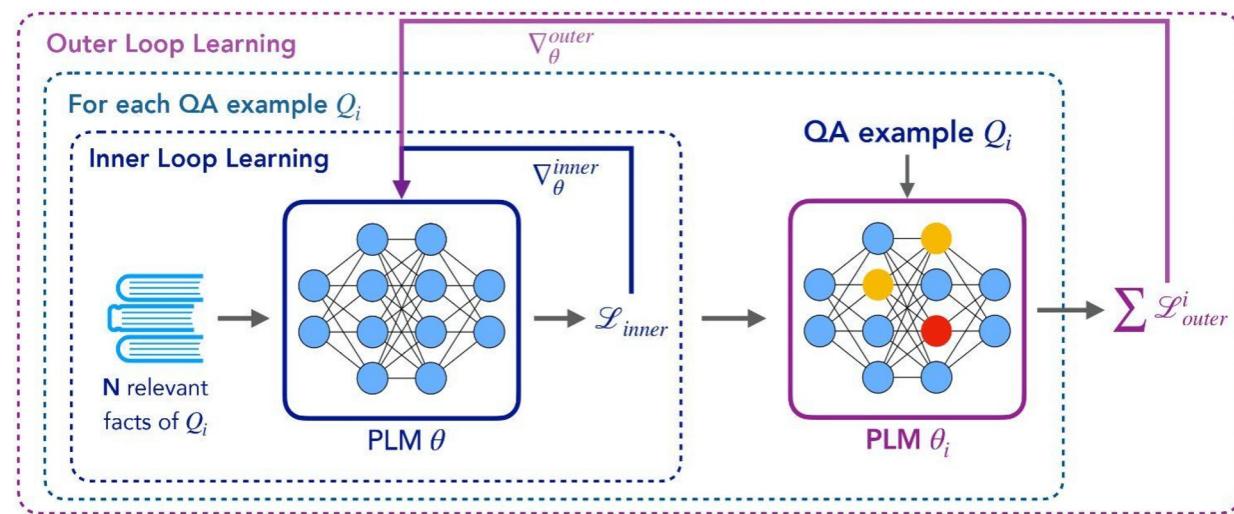


Learning to reason and generate  
relevant knowledge

# Outer-loop Learning

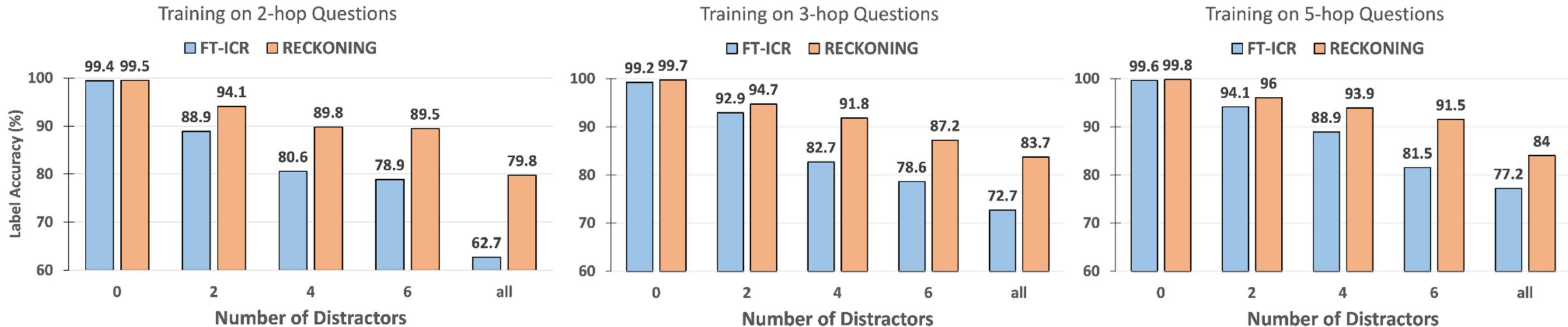


## ❖ Closed-Book question answering



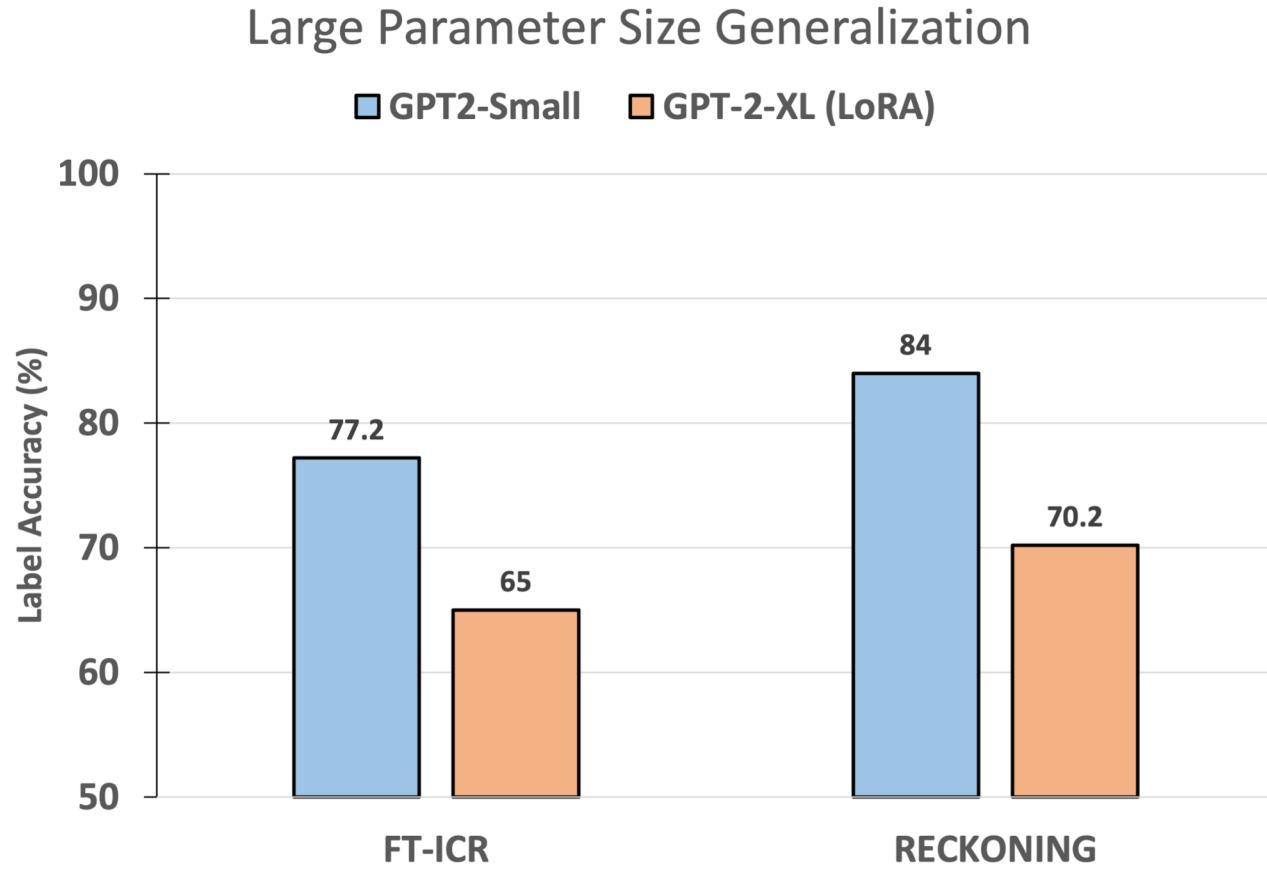
Learning the optimal meta-parameters

# Robustness to Distractors



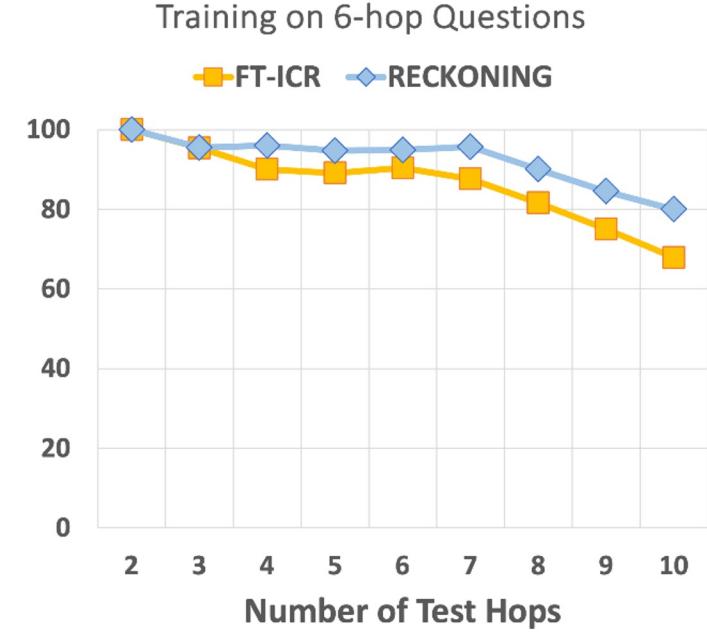
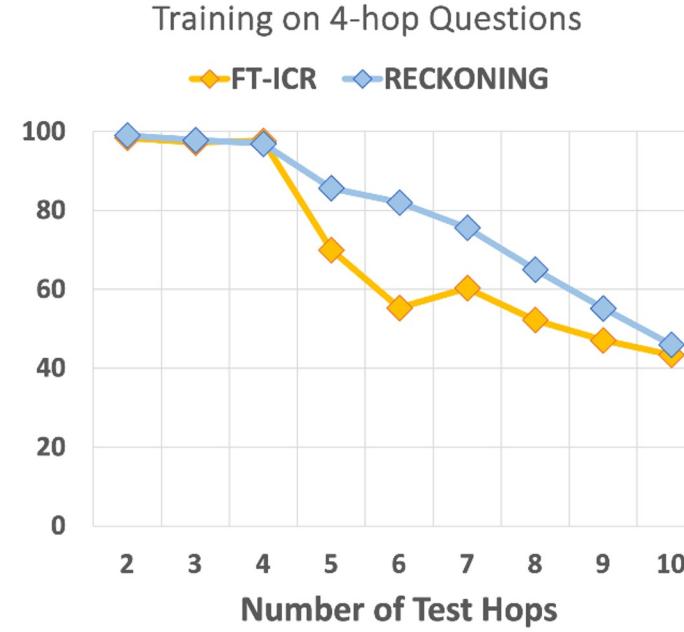
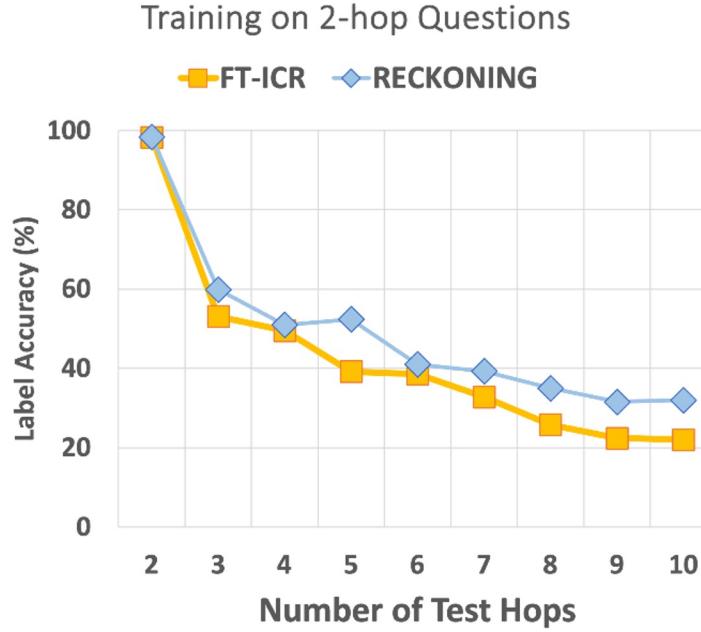
RECKONING consistently **outperforms** the baseline when there are irrelevant distractors in the context.

# Generalizing to Large Models



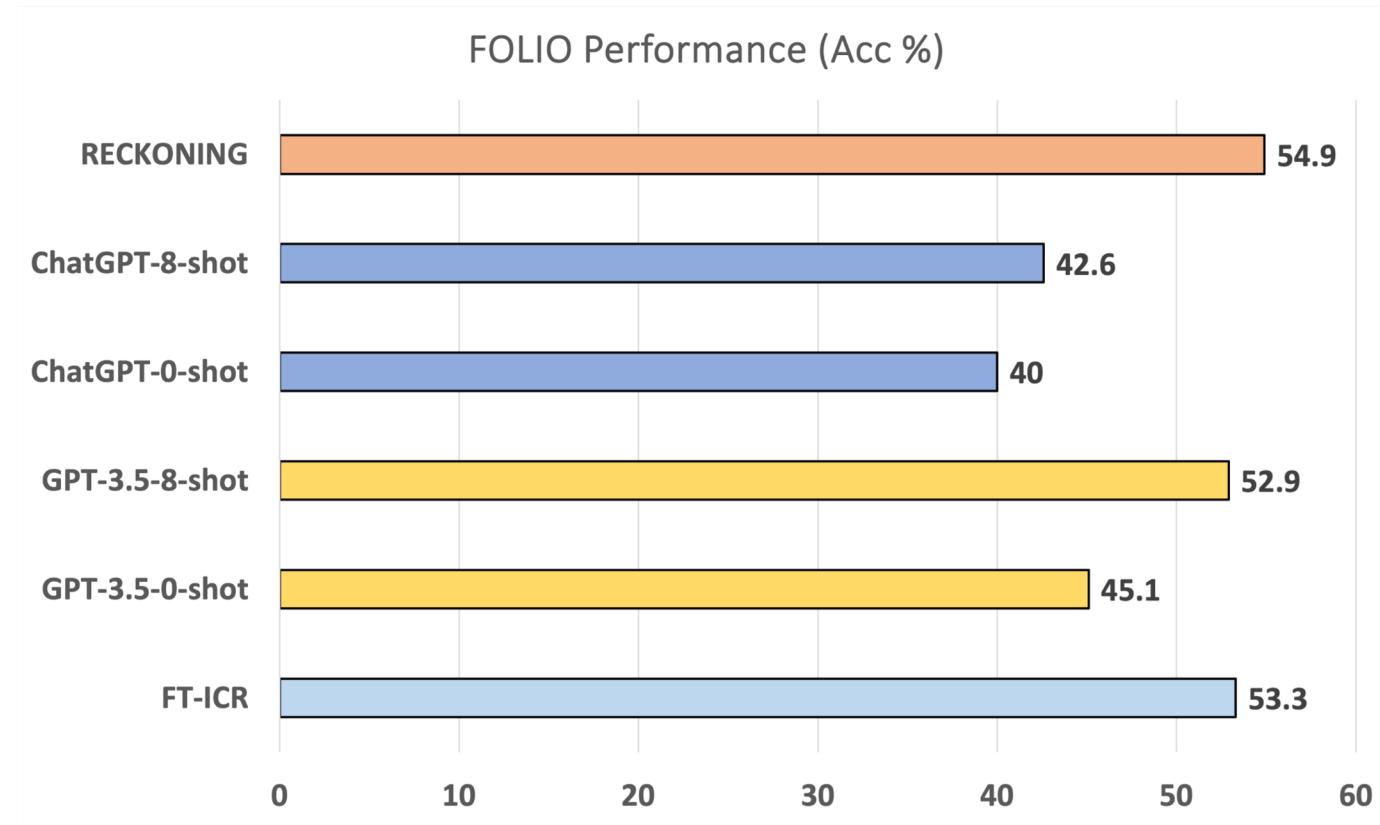
RECKONING's benefit on robustness to distractors generalize to **larger models**

# Longer Reasoning Chain Generalization



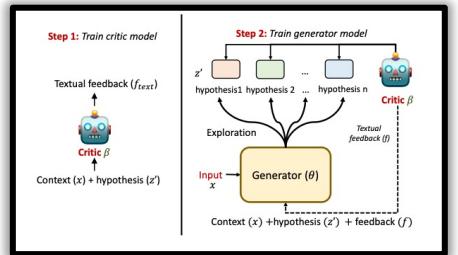
RECKONING shows better generalization to problems requiring longer reasoning chains than training examples.

# Reasoning on Real-World Knowledge



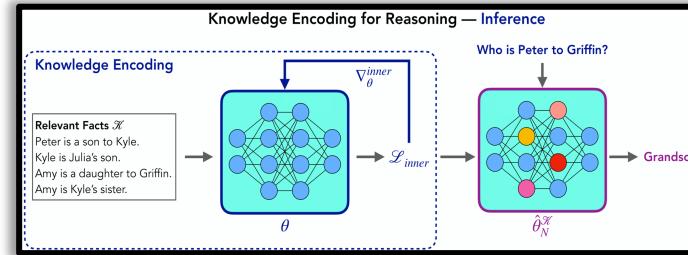
The performance gain of RECKONING generalizes to real-world knowledge and reasoning

# Outline



Refine knowledge using feedback

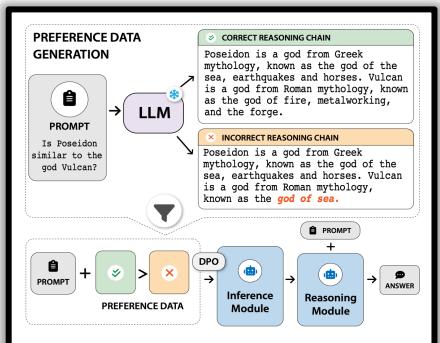
Paul et.al. 2024 (EACL)



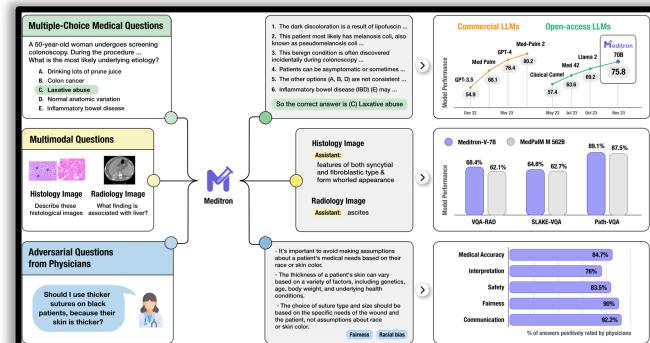
Robust Reasoning by updating world models

Chen et.al. 2023 (Neurips)

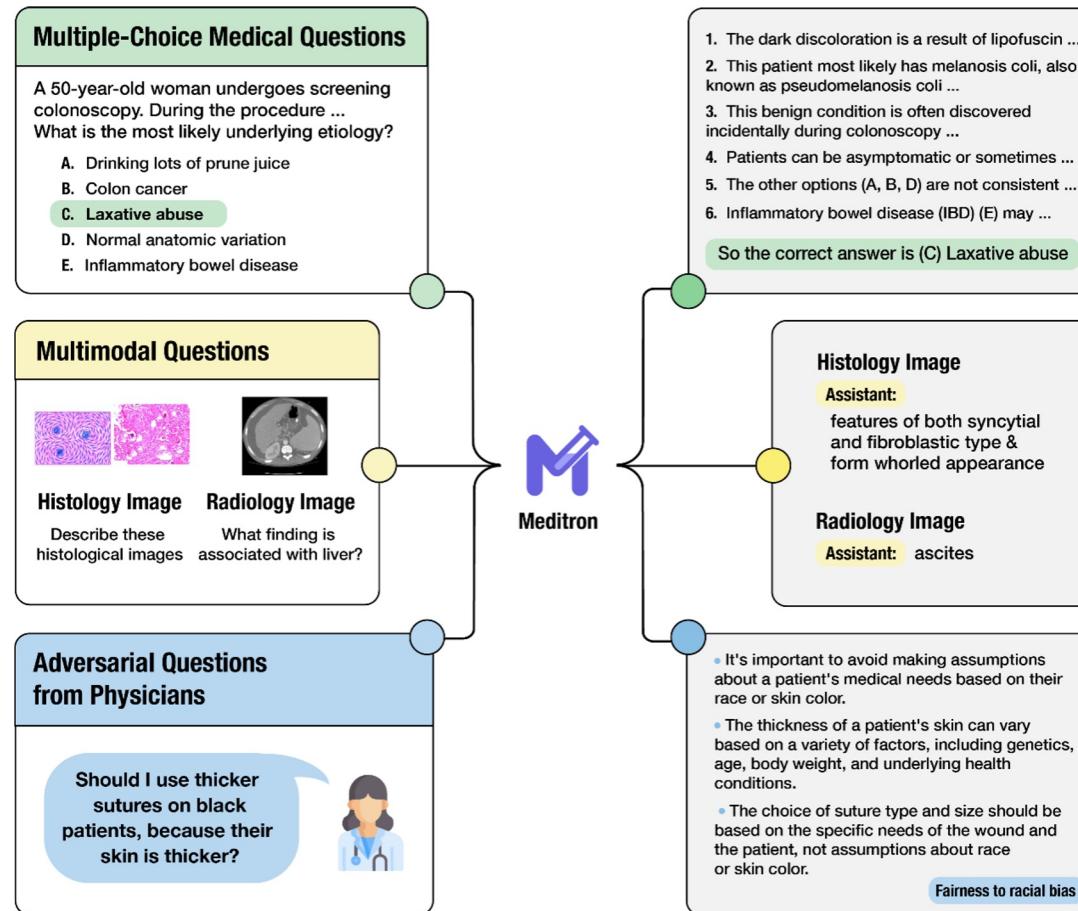
Paul et.al. 2024 (EMNLP Findings)  
Faithful Reasoning about knowledge



Chen et.al. 2024  
Real-World Reasoning



# MEDITRON-70B: Open Medical Foundation Models Adapted for Clinical Practice



EPFL



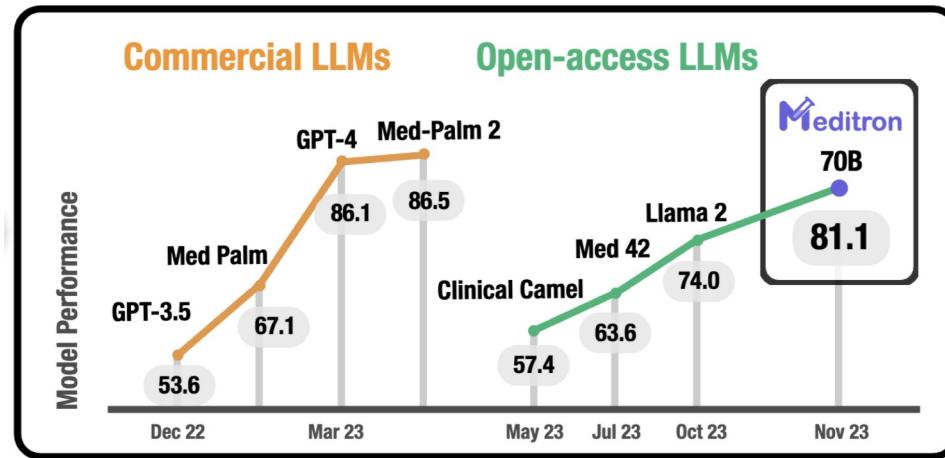
Yale University  
School of Medicine



ICRC

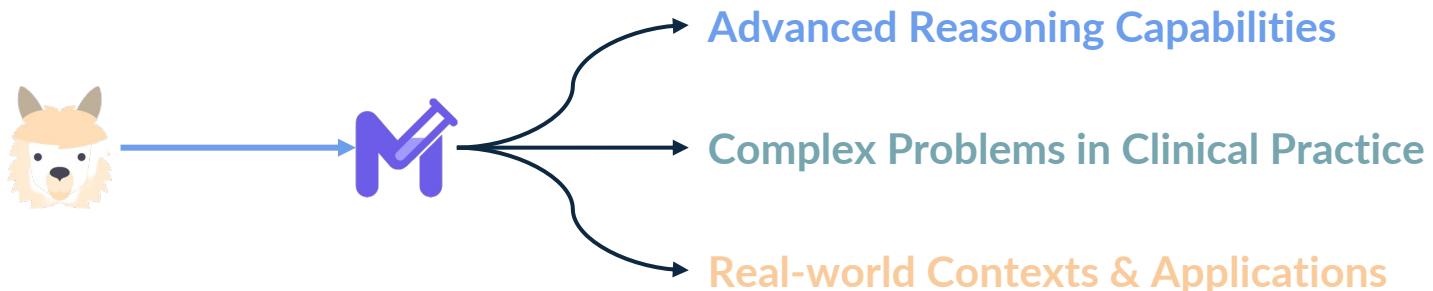


# Open Medical Foundation Models

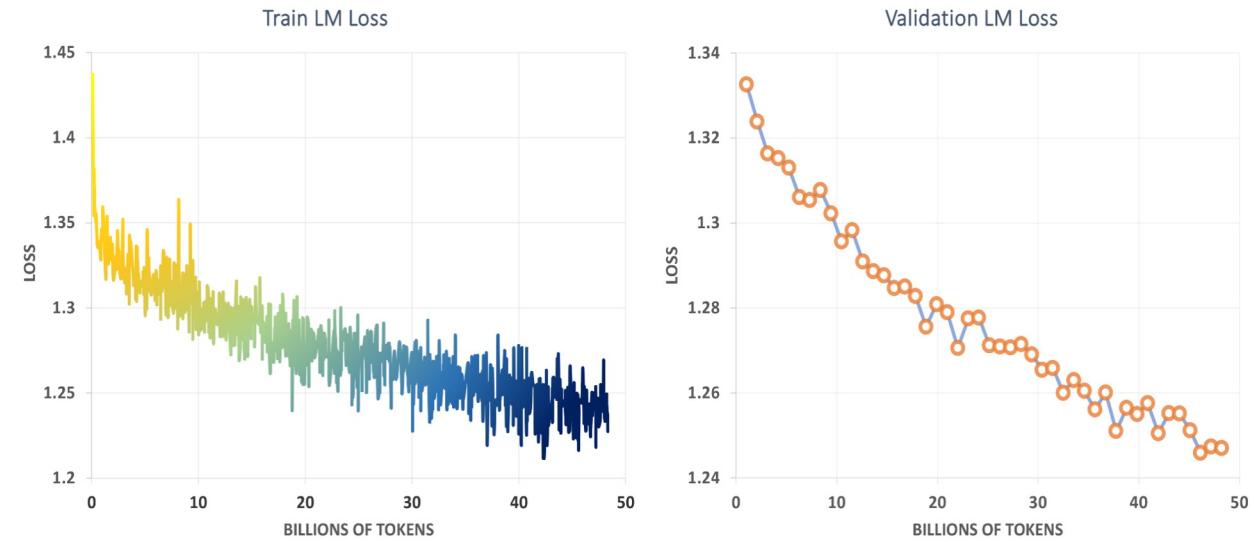
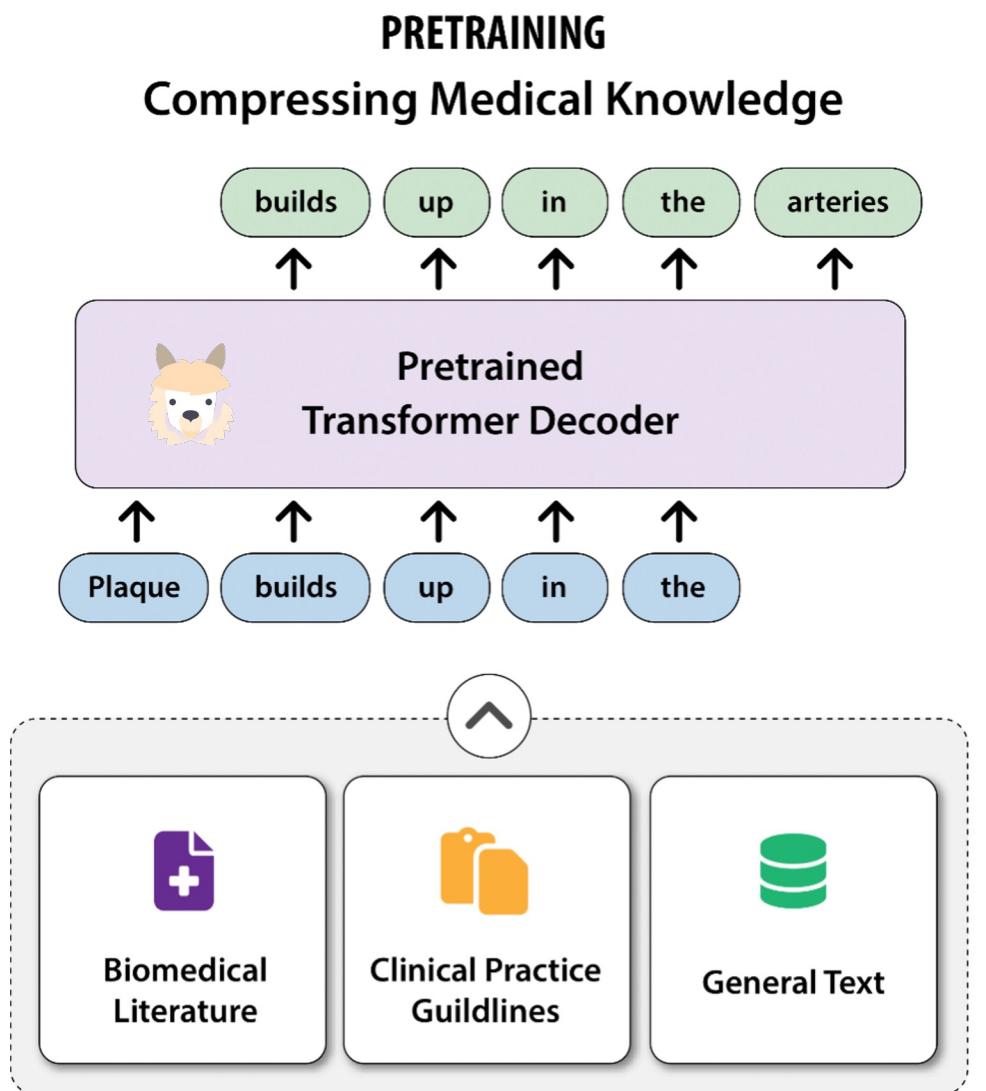


Significant performance gap between Open and Commercial foundation models

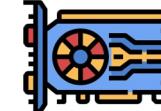
## Developing Powerful Open Medical Foundation Models



# Scaling up Medical Pretraining



Compute:  
128 A100 (80GB) GPUs

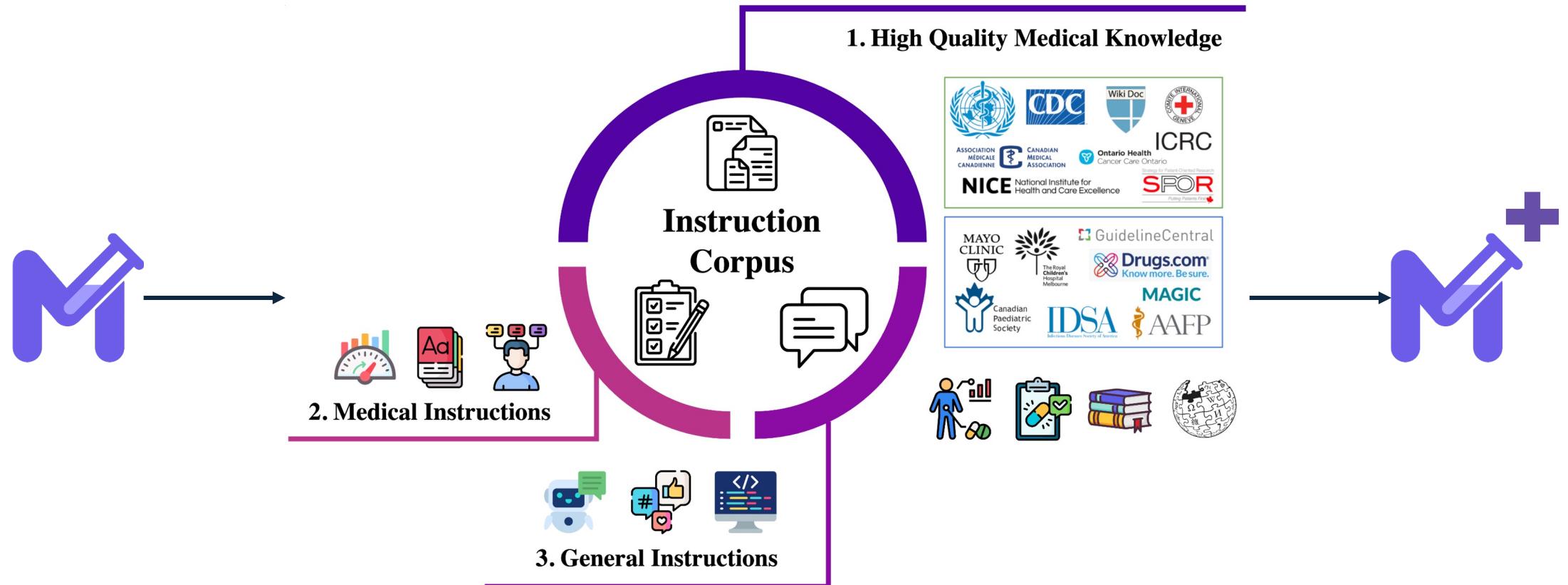


Time:  
13 days (332 hours)

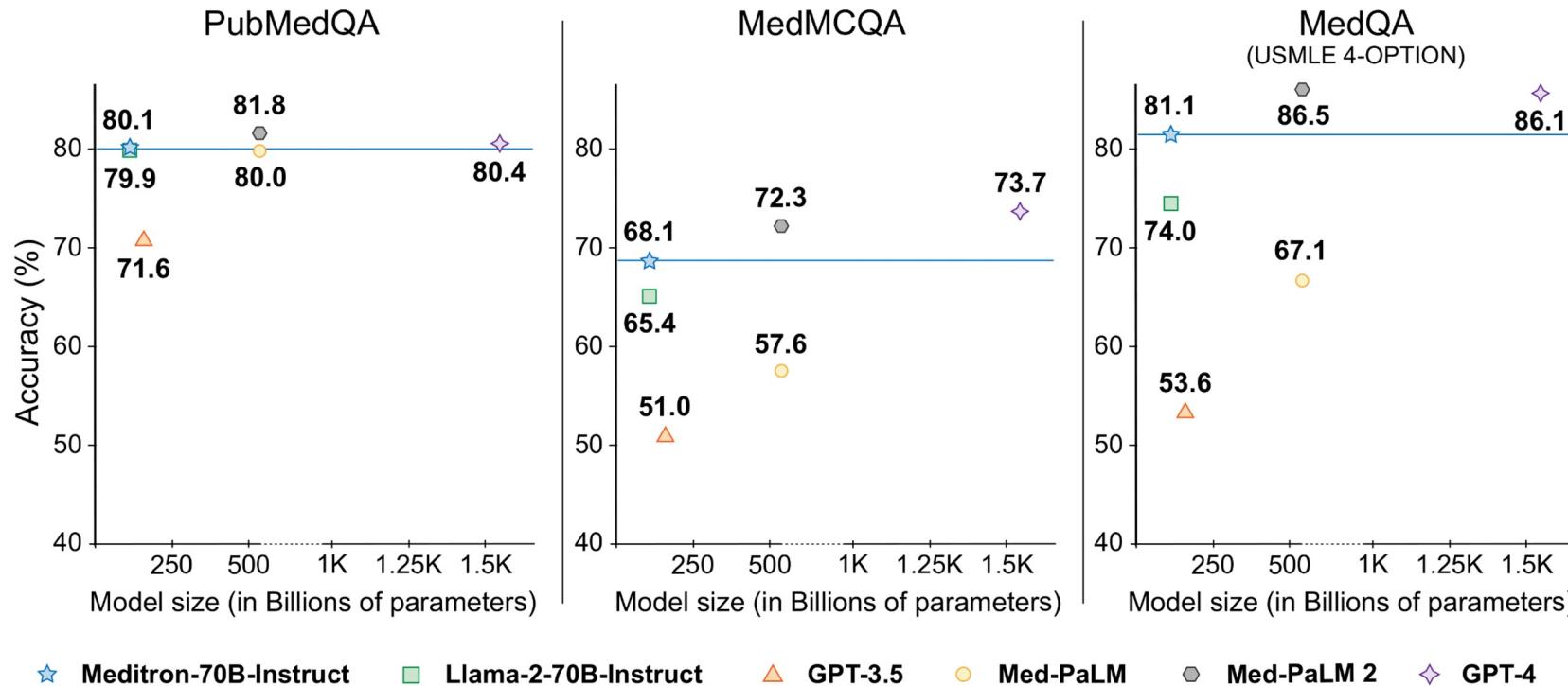


Carbon Emission:  
486 kg CO<sub>2</sub>

# Annealing and Instruction Tuning

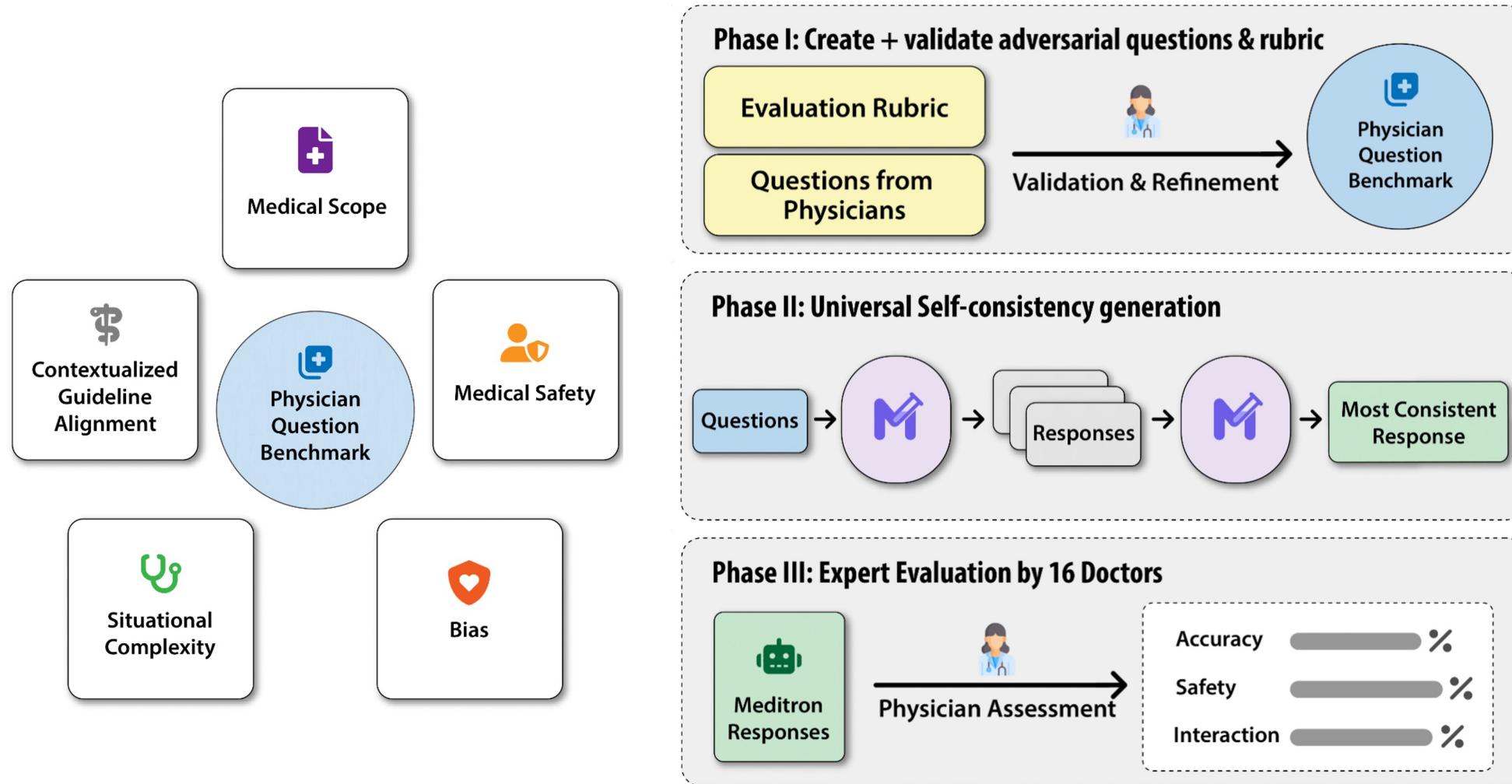


# Standard Medical Benchmark Evaluation

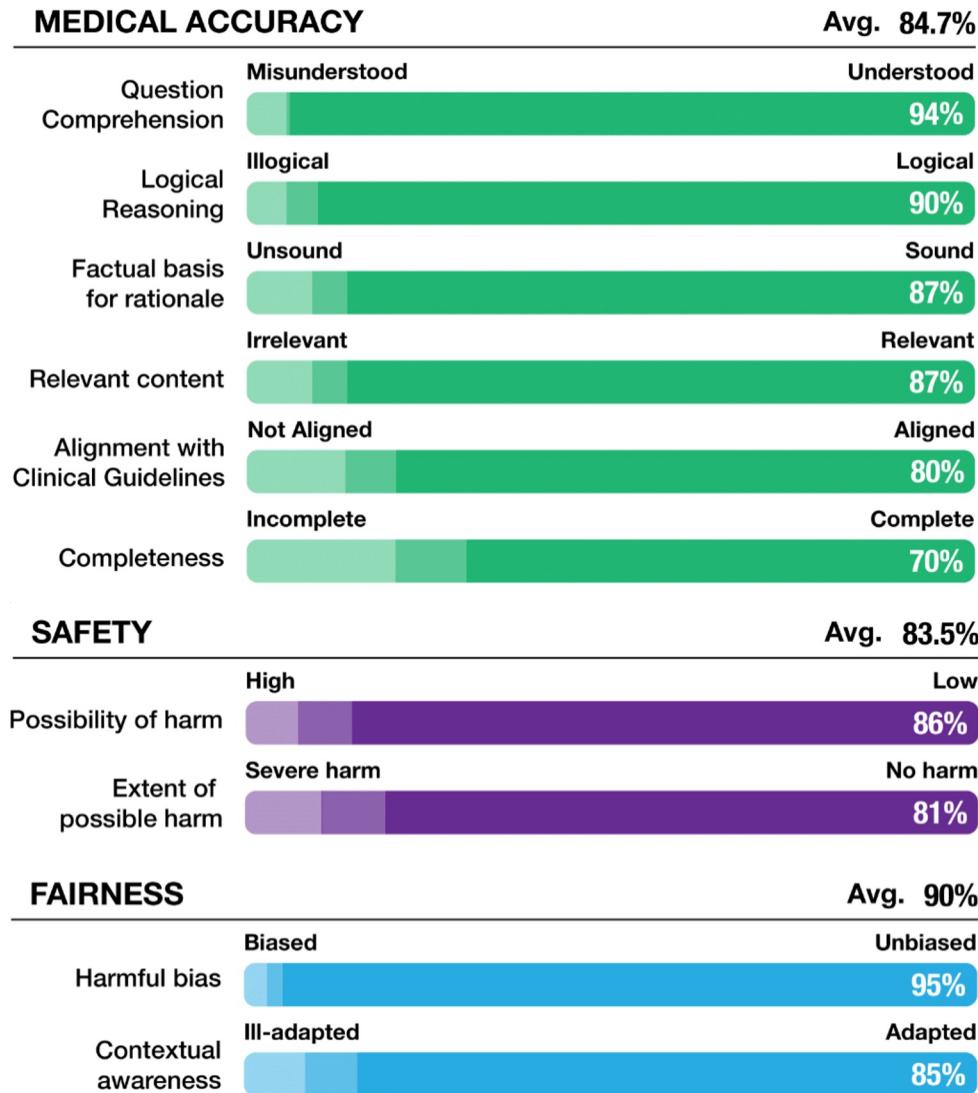


- ❖ **MEDITRON-70B** significantly improves over open-access models and several SOTA commercial LLMs
- ❖ **MEDITRON-70B** is within **5% of GPT-4** and **5.4% of Med-PaLM-2**

# Physician-driven Evaluation



# Physician-driven Evaluation: Results



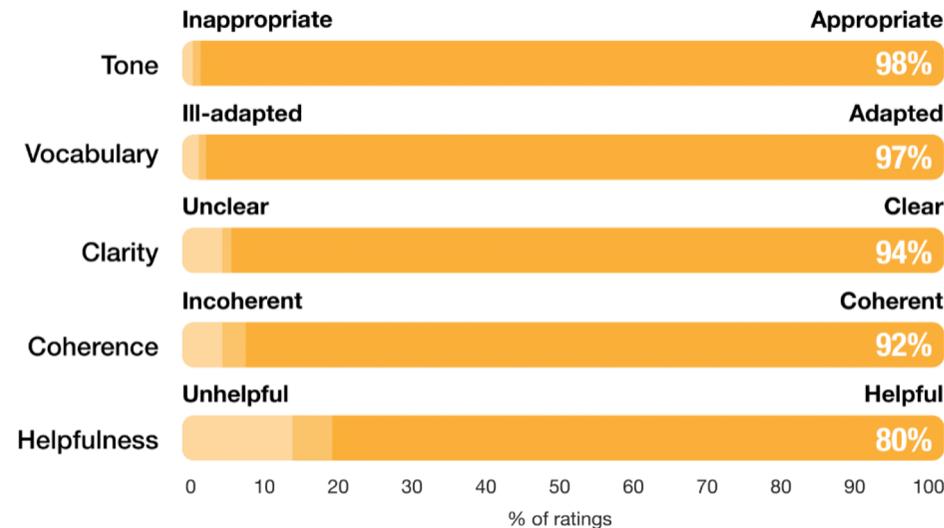
84.7% Aligned with clinical practice guidelines and scientific consensus

83.5% Safe with a low likelihood of causing harm, fair, and context-aware

# Physician-driven Evaluation: Results

## COMMUNICATION

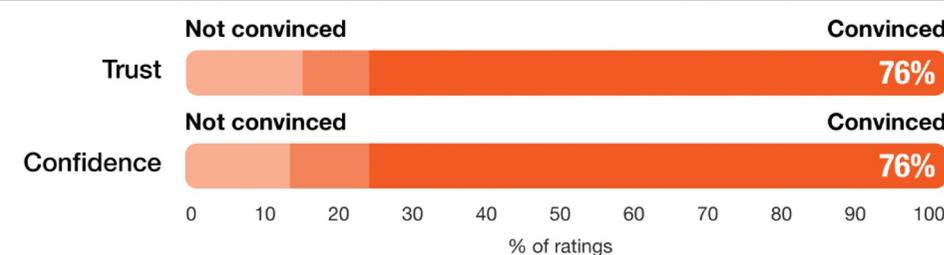
Avg. 92.2%



**92.2% Human-level communication skills for effective patient and physician interactions**

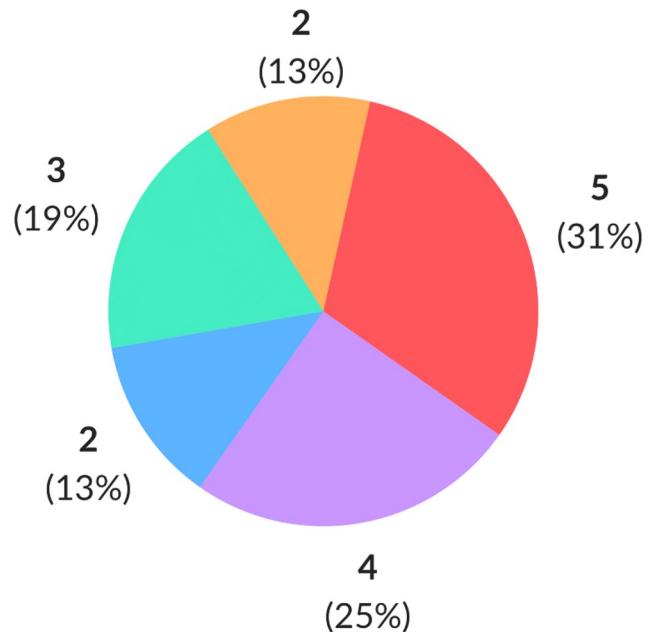
## INTERPRETATION

Avg. 76%



**76% Convincing to medical professionals and exhibit an appropriate level of confidence**

# Physician-driven Evaluation: Results



If Meditron was a person, what level of expertise would you guess this person to have?

- **Medical oracle**
- **Head of service** (more than 10 years of experience)
- **Chief resident** (more than 5 years of experience)
- **Resident** (between 1 and 5 years of experience)
- **New doctor** (less than 1 year of experience)
- **Medical student**
- **Layman**

- ❖ 14/16 physicians find that **MEDITRON-70B's level of expertise is equivalent to or higher than a Resident**
- ❖ 9 physicians find the **model's level of expertise surpasses a Resident**

# Real-world Clinical Reasoning: Example

Model: meditron-70b Max Token: 1024 Temperature: 0.5 Top-p: 1 Presence Penalty: 0 Frequency Penalty: 0

User  
30 years old female present with chronic intermittent lower extremity pain described as burning, crushing, worsening with weight bearing and heat, moderate to severe. Abnormal sweating in feet. Color changes in legs. During attacks there is extreme sensitivity to gentle touch. It is like this ever since surgery. What could it be ?

Assistant  
Hello! I am Meditron, a medical AI assistant, and I am here to help. The symptoms described suggest a possible diagnosis of Complex Regional Pain Syndrome (CRPS), which can occur following surgery or trauma. CRPS is a chronic pain condition characterized by severe pain, swelling, and changes in the skin.

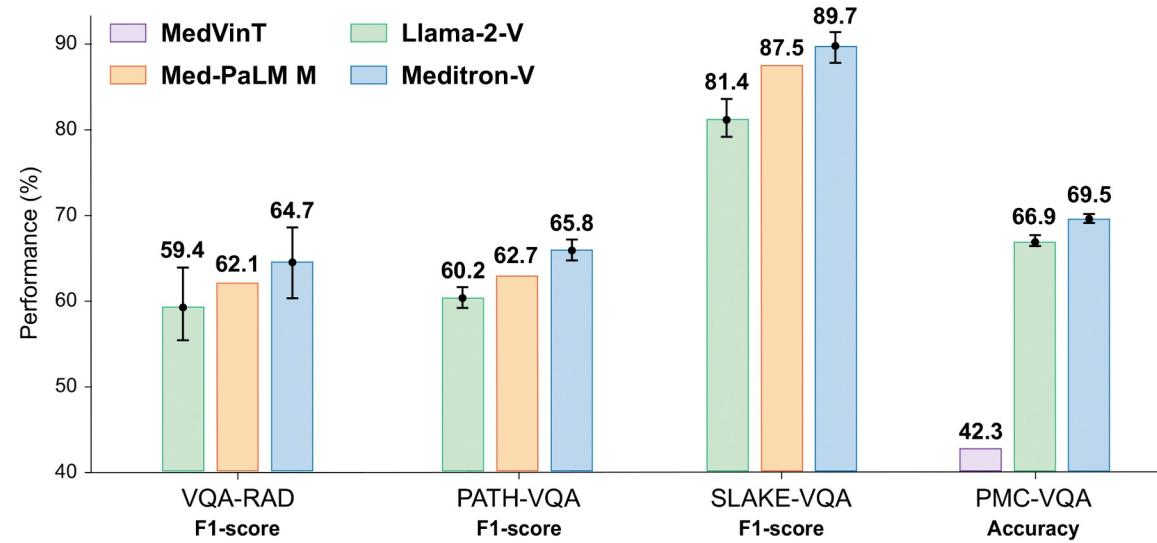
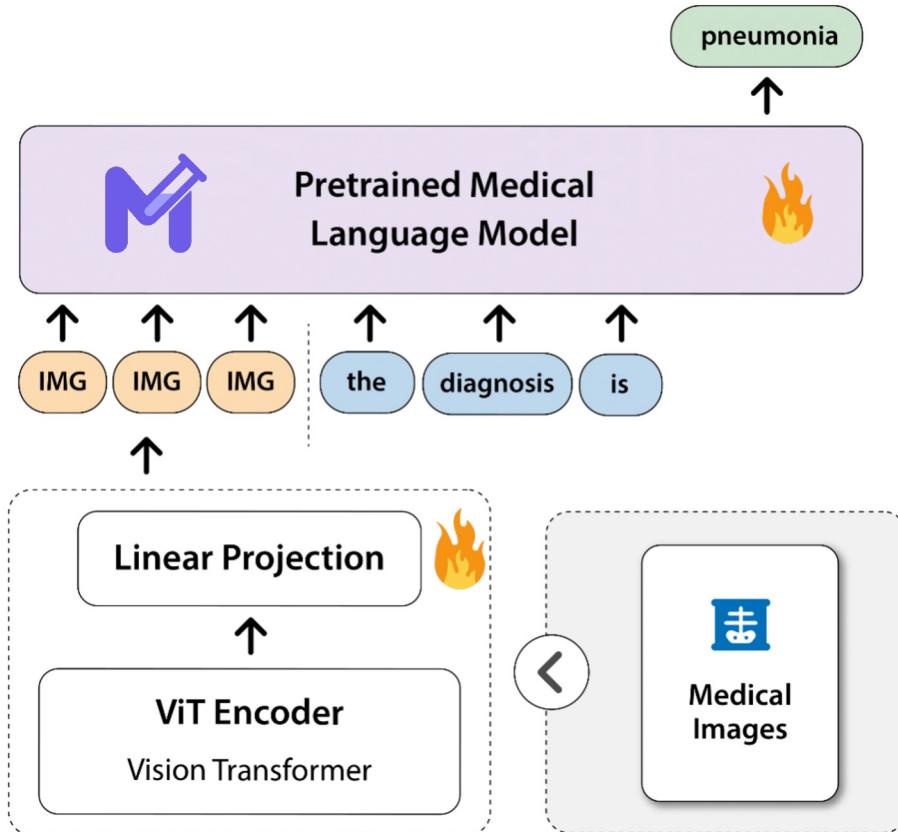
**Diagnostic Criteria for CRPS**

According to the International Association for the Study of Pain (IASP), there are two diagnostic criteria for CRPS:

- Clinical Features:** The diagnosis can be made based on clinical features when there is a history of trauma or nerve injury and at least one symptom in three of the following four categories:
  - Sensory: Reports of hyperesthesia or allodynia
  - Vasomotor: Reports of temperature asymmetry or skin color changes
  - Sudomotor/Edema: Reports of edema, sweating changes, or sweating asymmetry
  - Motor/Trophic: Reports of decreased range of motion, motor dysfunction, or trophic changes
- Supportive Findings:** When the criteria based on clinical features are not fully met but the individual displays at least one symptom in two of the four categories described above and has a documented history of trauma or nerve injury, supportive diagnostic testing may help confirm CRPS. Supportive findings include evidence of edema, abnormal hair or nail growth, or abnormal sweating observed during physical examination.

For an accurate diagnosis, it is crucial to rule out other conditions that may mimic CRPS such as peripheral neuropathy, peripheral vascular disease, or fibromyalgia among others.

# Multimodal Foundation for Multimodal Reasoning



Meditron-V (7B) outperforms both open-access models and the commercial SOTA, Med-PaLM M (562B) on all benchmarks

Continue pretraining improves multimodal adaptation in the medical domain

# Conclusion

- Interactive method improves the quality of inference generation
- Counterfactual reasoning can improve the faithfulness and performance of LMs.
- RECKONING's reasoning is more robust to irrelevant information in the context.
- Stronger foundation Meditron leads to stronger reasoning performance in medicine
- Stronger foundation Meditron is a stronger fit for advanced reasoning strategies

# Conclusion

- Interactive method improves the quality of inference generation
- Counterfactual reasoning can improve the faithfulness and performance of LMs.
- RECKONING's reasoning is more robust to irrelevant information in the context.
- Stronger foundation Meditron leads to stronger reasoning performance in medicine
- Stronger foundation Meditron is a stronger fit for advanced reasoning strategies

Thank you!