

### 10-423/10-623 Generative Al

Machine Learning Department School of Computer Science Carnegie Mellon University

# Reasoning Models + Mechanistic Interpretability

Aran Nayebi & Matt Gormley Lecture 20 Nov. 10, 2025

### Reminders

- Exam Tonight!
  - -Time: 80 minutes at 7 pm
  - -Covered Material: Lectures 1 15 (same as Quiz 1 Quiz 4)
  - You may bring one sheet of notes (front and back)
  - Format of questions: Unlike the Quiz questions, which were all multiple choice, Exam questions will include open-ended questions as well
  - We are using two rooms, so please check Piazza for seat assignment

### **REASONING MODELS**

# Chain-of-Thought Prompting

- Asking the model to reason about its answer can improve its performance for few-shot in-context learning
- Chain-of-thought prompting provides such reasoning in the in-context examples

#### **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?

#### **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?

### PaLM 540B: standard prompting PaLM 540B: chain-of-thought prompting 100 +Solve rate (%) Math Word Problems (GSM8K)

Finetuned GPT-3 175B

Prior best

Figure 2: PaLM 540B uses chain-ofthought prompting to achieve new stateof-the-art performance on the GSM8K benchmark of math word problems. Finetuned GPT-3 and prior best are from Cobbe et al. (2021).

#### **Model Output**

A: The answer is 27.



#### **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. 🗸

# Chain-of-Thought Prompting

Recall.

- Asking the model to reason about its answer can improve its performance for few-shot in-context learning
- Chain-of-thought prompting provides such reasoning in the in-context examples

#### (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

#### (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. ✓

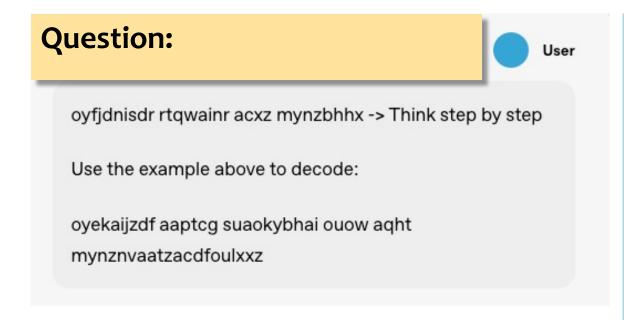
#### (d) Zero-shot-CoT (Ours)

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. ✓

 But the model does better even if you just prompt it to reason step-by-step







oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step

Use the example above to decode:

oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz

#### Thinking:

First, what is going on here?

We are given:

First, an example:

"oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step"

We can see that when "oyfjdnisdr rtqwainr acxz mynzbhhx" is transformed to "Think step by step"

Our task is: Use the example above to decode: "oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz"

So the first part is to figure out how "oyfjdnisdr rtqwainr acxz mynzbhhx" is decoded into "Think step by step"

. . .



oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step

Use the example above to decode:

oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz

#### Thinking:

...

Observation 1: The phrase "Think step by step" is 18 letters (excluding spaces). Let's count the letters.

Plaintext: Think step by step

First word: Think (5 letters)

Second word: step (4 letters)

Third word: by (2 letters)

Fourth word: step (4 letters)

Total letters in plaintext: 5 + 4 + 2 + 4 = 15 letters

Hmm.

....



User

oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step

Use the example above to decode:

oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz

#### Thinking:

• • •

Similarly, let's count the letters in "oyfjdnisdr rtqwainr acxz mynzbhhx"

First, break down the letters:

Ciphertext letters: o y f j d n i s d r r t q w a i n r a c x z m y n z b h h x

Total letters: Let's count the letters in the ciphertext.

First word: oyfjdnisdr (10 letters)

Second word: rtqwainr (8 letters)

Third word: acxz (4 letters)

Fourth word: mynzbhhx (8 letters)

Total letters in ciphertext: 10 + 8 + 4 + 8 = 30 letters

• • •



oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step

Use the example above to decode:

oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz

#### Thinking:

...

Given the time constraints, perhaps the easiest way is to try to see patterns.

Option 1: Try to find mapping from letters to letters.

Do any letters match?

First, let's write down the ciphertext and plaintext letters on top of each other.

Ciphertext: o y f j d n i s d r r t q w a i n r a c x z m y n z b h h x

Plaintext: Thinkstepbystep

Wait a minute.

I think maybe there is an anagram or substitution cipher here.

...



oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step

Use the example above to decode:

oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz

#### Thinking:

• • •

Check the number of letters.

First ciphertext word: oyfjdnisdr (10 letters)

First plaintext word: Think (5 letters)

Hmm.

...

Interesting.

It seems that the ciphertext words are exactly twice as long as the plaintext words.

(10 vs 5, 8 vs 4, 4 vs 2, 8 vs 4)

Idea: Maybe we need to take every other letter or rebuild the plaintext from the ciphertext accordingly.

Let's test this theory.

• • •



oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step

Use the example above to decode:

oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz

```
Thinking:
Now, what's the mapping from "oy" to 'T'?
Is there a way to map 'oy' to 'T'?
Option 1: Sum of the letters.
Convert letters to numbers (A=1, B=2, ..., Z=26)
So let's try to get to 'T' from 'o' and 'y'
0 = 15
y = 25
Sum: 15 +25 = 40
But 'T' is 20.
```



oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step

Use the example above to decode:

oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz

1276 lines later



Thinking:

• • •

So the message so far is

'THERE ARE THREE RS IN'

And then 'STRAWBERRY'

Wait, let's put the words together:

'THERE ARE THREE Rs IN STRAWBERRY'

I think 'RS' is 'R's', meaning 'Rs'

So the message is:

'THERE ARE THREE R'S IN STRAWBERRY'

Which is true.

So the decoded message is:

"THERE ARE THREE R'S IN STRAWBERRY"



oyfjdnisdr rtqwainr acxz mynzbhhx -> Think step by step

Use the example above to decode:

oyekaijzdf aaptcg suaokybhai ouow aqht mynznvaatzacdfoulxxz



- OpenAl o1 chose not to release the "Thinking" output
- Instead, they released a summary of that output.

#### Decoding Method:

- Pair the letters in the ciphertext.
- 2. Convert each letter to its numerical position in the alphabet (A=1, B=2, ..., Z=26).
- 3. Sum the numerical values of each pair.
- 4. Compute the average of the sum (divide by 2).
- 5. Convert the average back to a letter (1=A, 2=B, ..., 26=Z).

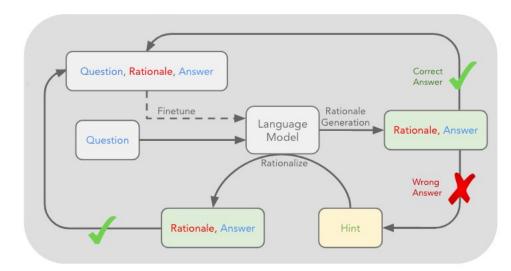
# Self-Taught Reasoner (STaR)

#### • Data:

- rationale examples (human annotated, small quantity)
- problems without rationales (large quantity)

### Repeat:

- bootstrap rationale training data:
  - use ICL with a few rationale examples
  - generate rationales for problems without
  - if generated answer is wrong, then try to regenerate a rationale that leads to a correct answer
- fine-tune on all rationales that led to correct answers



```
Q: What can be used to carry a small dog?
Answer Choices:
(a) swimming pool
(b) basket
(c) dog show
(d) backyard
(e) own home
A: The answer must be something that can be used to carry a small dog. Baskets are designed to hold things. Therefore, the answer is basket (b).
```

### AIME Dataset

#### **Dataset Description**

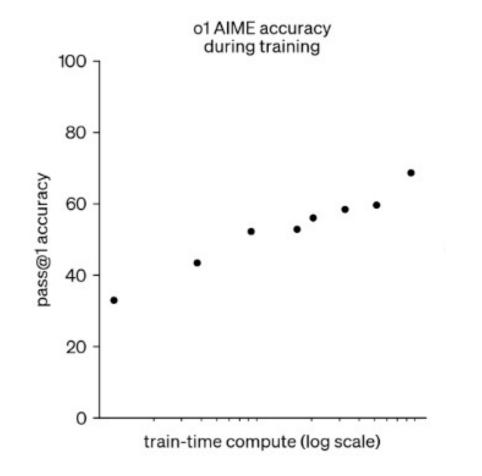
This dataset contains problems from the American Invitational Mathematics Examination (AIME) 2024. AIME is a prestigious high school mathematics competition known for its challenging mathematical problems.



### OpenAl 01

- The o1 model was trained with Reinforcement Learning to generate chainof-thought style rationales for its answers
- These rationales (referred to as Thinking tokens) were hidden from the user, and a summary of the Thinking tokens was presented instead
- At train time: the compute could be increased by performing more reinforcement learning
- At test time: the compute could be increased by spending more time thinking

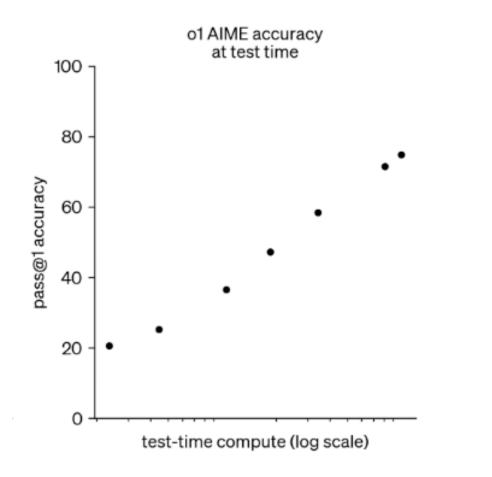
 Result 1: more train time compute leads to higher accuracy on reasoning problems



### OpenAl 01

- The o1 model was trained with Reinforcement Learning to generate chainof-thought style rationales for its answers
- These rationales (referred to as Thinking tokens) were hidden from the user, and a summary of the Thinking tokens was presented instead
- At train time: the compute could be increased by performing more reinforcement learning
- At test time: the compute could be increased by spending more time thinking

• Result 2: more test time compute leads to higher accuracy on reasoning problems



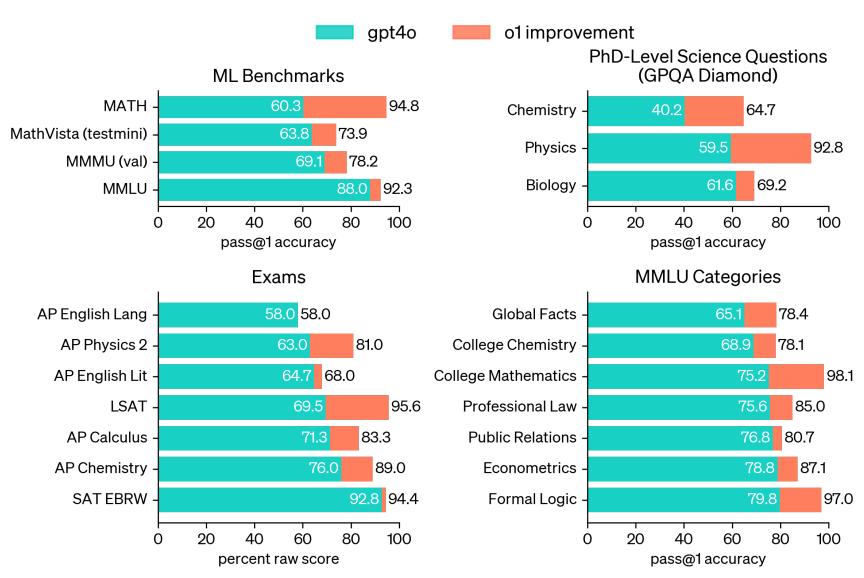
### OpenAl 01

- The o1 model was trained with Reinforcement Learning to generate chainof-thought style rationales for its answers
- These rationales (referred to as Thinking tokens) were hidden from the user, and a summary of the Thinking tokens was presented instead
- At train time: the compute could be increased by performing more reinforcement learning
- At test time: the compute could be increased by spending more time thinking

- Why is this description so vague and non-technical?
- A: Because OpenAI only released a blog post, and this is about the sum total of what it said

### OpenAl o1

Across a variety of math, reasoning, commonsense, coding, etc. problems or improves over gpt40



- The closed source model (o1) was clearly superior than any open source models
- So we waited for the open source models to catch up...



### DeepSeek-R1-Zero and DeepSeek-R1

- Enter DeepSeek-R1…
- This open source and open weight model is 671B parameters
- It is a carefully tuned version of the base model DeepSeek-V3
- And it achieves comparable performance to 01

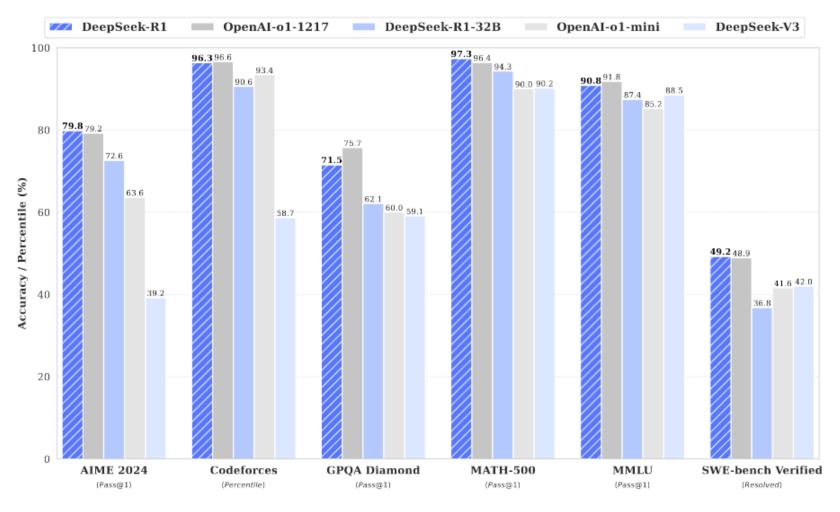
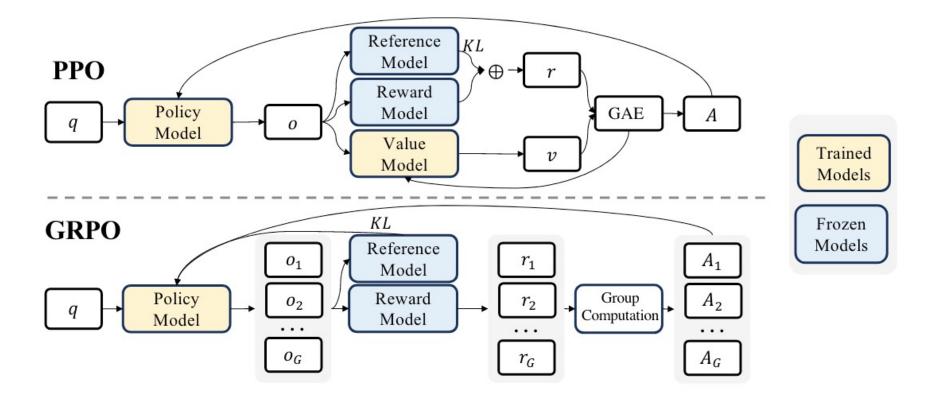


Figure 1 | Benchmark performance of DeepSeek-R1.

### PPO vs. GRPO

- DeepSeek-Math came before DeepSeek-R1 and introduced the idea of GRPO
- GRPO is an RL algorithm akin to PPO, but it greatly reduces the memory requirements by removing the need for a Value Model



### PPO vs. GRPO

Proximal Policy Optimization (PPO) (Schulman et al., 2017) is an actor-critic RL algorithm that is widely used in the RL fine-tuning stage of LLMs (Ouyang et al., 2022). In particular, it optimizes LLMs by maximizing the following surrogate objective:

$$\mathcal{J}_{PPO}(\theta) = \mathbb{E}\left[q \sim P(Q), o \sim \pi_{\theta_{old}}(O|q)\right] \frac{1}{|o|} \sum_{t=1}^{|o|} \min\left[\frac{\pi_{\theta}(o_t|q, o_{< t})}{\pi_{\theta_{old}}(o_t|q, o_{< t})} A_t, \operatorname{clip}\left(\frac{\pi_{\theta}(o_t|q, o_{< t})}{\pi_{\theta_{old}}(o_t|q, o_{< t})}, 1 - \varepsilon, 1 + \varepsilon\right) A_t\right], \quad (1)$$

where  $\pi_{\theta}$  and  $\pi_{\theta_{old}}$  are the current and old policy models, and q, o are questions and outputs sampled from the question dataset and the old policy  $\pi_{\theta_{old}}$ , respectively.  $\varepsilon$  is a clipping-related hyper-parameter introduced in PPO for stabilizing training.  $A_t$  is the advantage, which is computed by applying Generalized Advantage Estimation (GAE) (Schulman et al., 2015), based

on the rewards  $\{r_{\geq t}\}$  and a learned value function  $V_{\psi}$ . Thus, in PPO, a value function needs to be trained alongside the policy model and to mitigate over-optimization of the reward model, the standard approach is to add a per-token KL penalty from a reference model in the reward at each token (Ouyang et al., 2022), i.e.,

$$r_{t} = r_{\varphi}(q, o_{\leq t}) - \beta \log \frac{\pi_{\theta}(o_{t}|q, o_{< t})}{\pi_{ref}(o_{t}|q, o_{< t})},$$
(2)

where  $r_{\varphi}$  is the reward model,  $\pi_{ref}$  is the reference model, which is usually the initial SFT model, and  $\beta$  is the coefficient of the KL penalty.

### PPO vs. GRPO

accurate at each token. To address this, as shown in Figure 4, we propose Group Relative Policy Optimization (GRPO), which obviates the need for additional value function approximation as in PPO, and instead uses the average reward of multiple sampled outputs, produced in response to the same question, as the baseline. More specifically, for each question q, GRPO samples a group of outputs  $\{o_1, o_2, \cdots, o_G\}$  from the old policy  $\pi_{\theta_{old}}$  and then optimizes the policy model by maximizing the following objective:

$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_{i}\}_{i=1}^{G} \sim \pi_{\theta_{old}}(O|q)] 
\frac{1}{G} \sum_{i=1}^{G} \frac{1}{|o_{i}|} \sum_{t=1}^{|o_{i}|} \left\{ \min \left[ \frac{\pi_{\theta}(o_{i,t}|q, o_{i,< t})}{\pi_{\theta_{old}}(o_{i,t}|q, o_{i,< t})} \hat{A}_{i,t}, \operatorname{clip}\left( \frac{\pi_{\theta}(o_{i,t}|q, o_{i,< t})}{\pi_{\theta_{old}}(o_{i,t}|q, o_{i,< t})}, 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{KL}\left[ \pi_{\theta} || \pi_{ref} \right] \right\},$$
(3)

where  $\varepsilon$  and  $\beta$  are hyper-parameters, and  $\hat{A}_{i,t}$  is the advantage calculated based on relative rewards of the outputs inside each group only, which will be detailed in the following subsections. The group relative way that GRPO leverages to calculate the advantages, aligns well with the comparative nature of rewards models, as reward models are typically trained on datasets of comparisons between outputs on the same question. Also note that, instead of adding KL penalty in the reward, GRPO regularizes by directly adding the KL divergence between the trained policy and the reference policy to the loss, avoiding complicating the calculation of  $\hat{A}_{i,t}$ .

### **Training method:**

- Trained entirely via Reinforcement Learning (RL) without any supervised fine-tuning (SFT).
- Started with a pretrained base model (DeepSeek-V3-Base), then used RL without human preferences to drive learning.
- Relied on a pure RL pipeline, making it one of the first large-scale demonstrations of RL-only training in LLMs.
- Aimed to explore whether reasoning abilities can emerge solely through RL, without labeled datasets.

#### **Reward model:**

- Did not use a neural reward model
- Instead just defined a rule-based reward model consisting of two parts:
  - Accuracy rewards: did the model answer the question correctly?
  - Format rewards: did the model adhere to the prompt template?

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> 
< and <answer> </answer> tags, respectively, i.e., <think> reasoning process here 
< think> <answer> <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.

#### **Results:**

- On AIME, the longer the model is trained with RL, the better the performance becomes
- Eventually it surpasses of performance

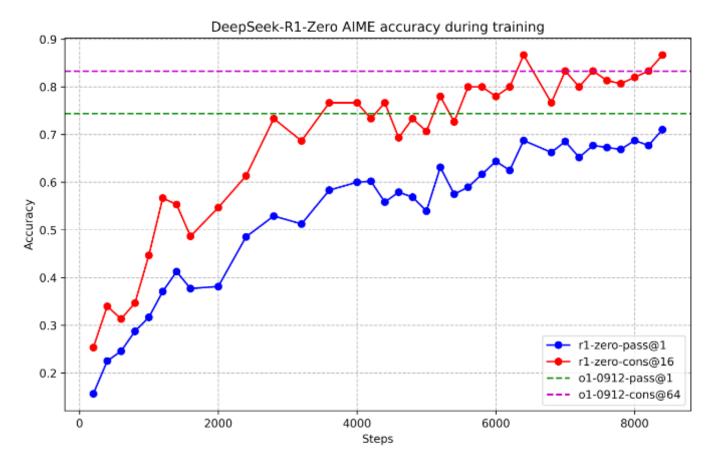


Figure 2 | AIME accuracy of DeepSeek-R1-Zero during training. For each question, we sample 16 responses and calculate the overall average accuracy to ensure a stable evaluation.

#### **Results:**

- Gradually the model learns to use longer and longer sequences of Thinking tokens
- This is accomplished purely through the RL objective
- There is no direct action taken to increase reasoning length

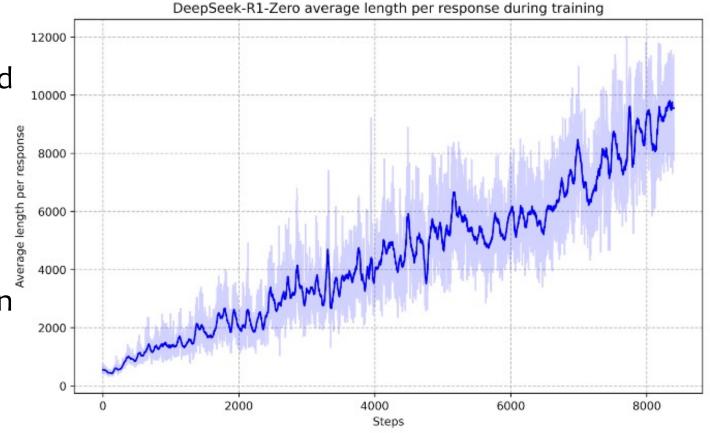


Figure 3 | The average response length of DeepSeek-R1-Zero on the training set during the RL process. DeepSeek-R1-Zero naturally learns to solve reasoning tasks with more thinking time.

#### **Problems:**

- Poor readability (e.g. humans don't really understand what it's saying)
- Language mixing (e.g. English and Chinese muddled into a pidgin language)

### DeepSeek-R1

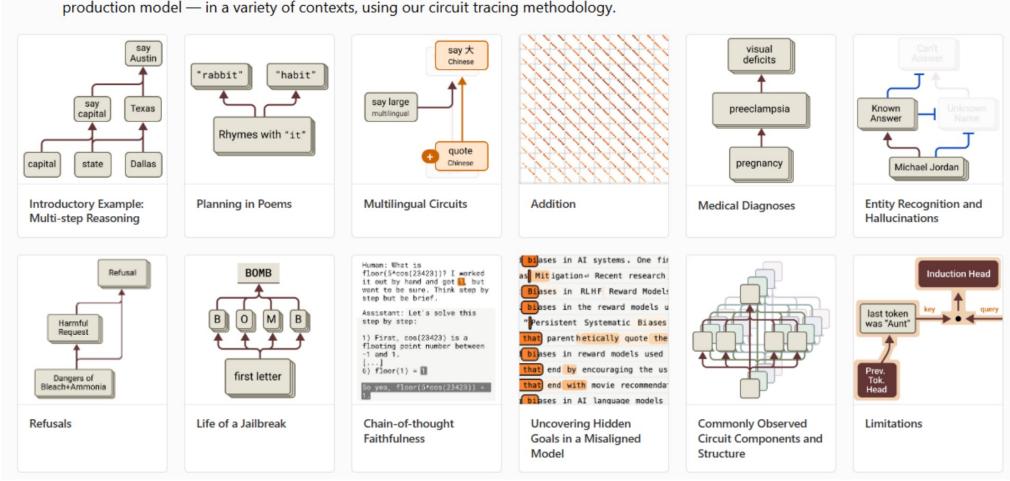
### **Training method:**

- Built on R1-Zero with a hybrid training strategy:
  - 1. Cold Start: Fine-tune a base model on a few thousand curated, human-friendly long CoTs.
  - Reasoning-Focused RL: Scale up RL with math, coding, and logic tasks. This time, add *language-consistency rewards* to push the model into staying coherent in a single language.
  - Rejection Sampling + SFT: Sample correct, well-structured chains-of-thought from the RL model, augment them with general capabilities data (writing, Q&A, self-cognition), and train a new base checkpoint.
  - RL Across Scenarios: A second RL stage includes both reasoning tasks and general tasks for "helpfulness" and "harmlessness."
- This two-stage pipeline addressed the shortcomings (e.g. repetition, language mixing) observed in R1-Zero.
- Emphasized improved readability, coherence, and task accuracy due to the incorporation of SFT before RL.

# Mechanistic Interpretability

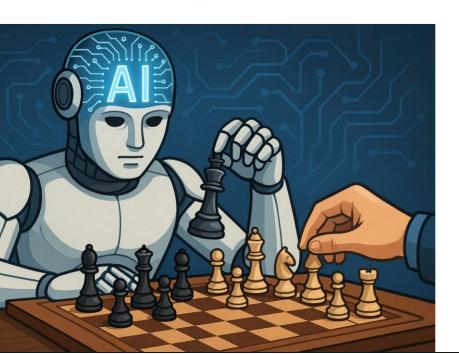
### On the Biology of a Large Language Model

We investigate the internal mechanisms used by Claude 3.5 Haiku — Anthropic's lightweight production model — in a variety of contexts, using our circuit tracing methodology.



### Why is Interpretability Important

- Safety (Corrective Action)
- Safety (Preventative Action)
- Preventing Al Apocalypse
- Learning from Al





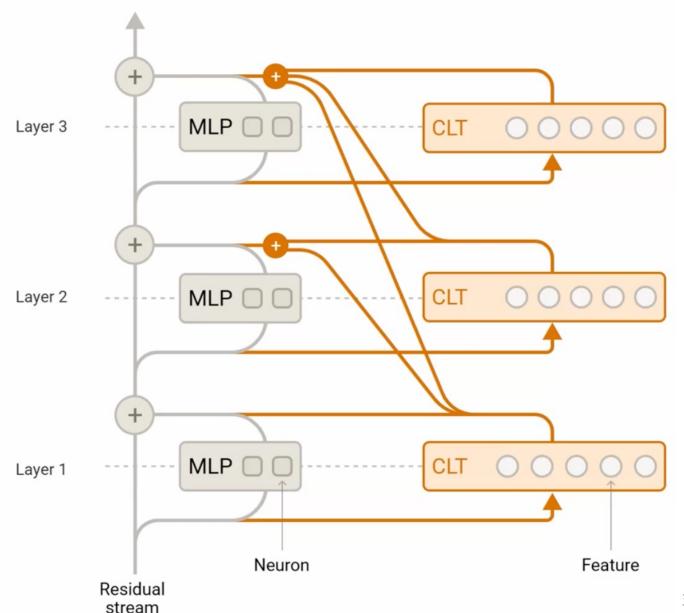


### Why is interpretability hard?

- Superposition is a major problem
  - –"Features" rarely activate in a single location in the network
  - –Activations for a human-interpretable feature is almost always distributed across many, many location within the network, e.g. across heads, across MLP neurons, across layers

### Replacement Models

 Idea: For specific blocks within a network, train a "replacement block" to mimic the functionality (input → ouput) of the original block...but importantly:



# Replacement Models

- Idea: For specific blocks within a network, train a
  "replacement block" to mimic the functionality (input 

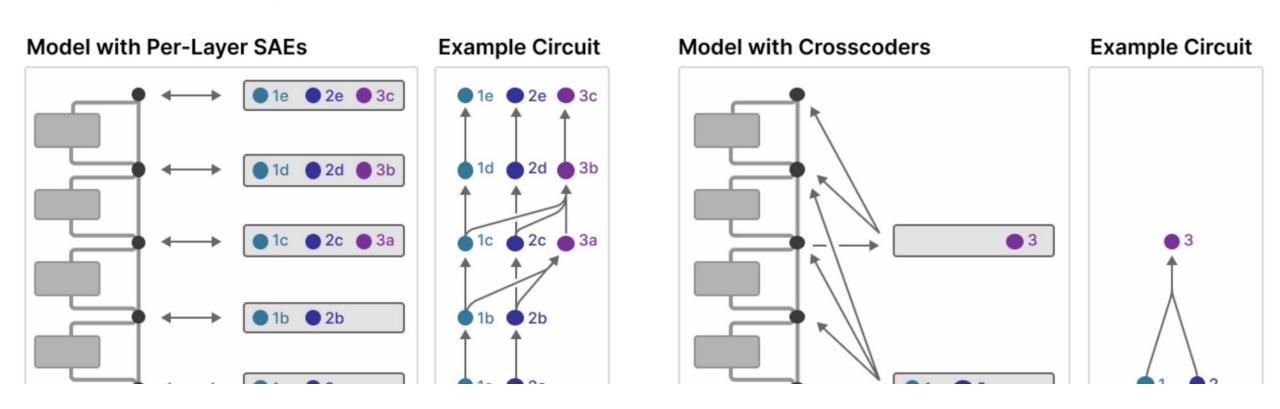
  ouput) of the original block...but importantly:
  - Make the replacement block more interpretable
- Techniques:
  - Sparse autoencoders
  - -Circuits
  - Cross layer transcoders

# **Sparse Autoencoders**

- Replace MLP layer in transformer block with an autoencoder version with:
  - -more neurons (features) in the hidden layer
  - an added regularization to encourage sparsity in the activations, e.g. L1 regularization

### Circuits and Cross-layer Transcoders

 Cross-layer: allow replacement blocks to directly access all earlier replacement blocks



# Mechanistic Interpretability

### On the Biology of a Large Language Model

We investigate the internal mechanisms used by Claude 3.5 Haiku — Anthropic's lightweight

