

***ThinkTuning*: Instilling Cognitive Reflections without Distillation**

Aswin RRV, Jacob Dineen, Divij Handa, Md Nayem Uddin, Mihir Parmar, Chitta Baral, Ben Zhou

**The 2025 Conference on Empirical Methods in Natural Language Processing
Suzhou, China - November 4–9, 2025**



Introduction


Introduction

Large Language Models (LLMs) have significantly transformed the AI landscape, by achieving human-like performance on various downstream tasks.

Model Parameters  → Downstream Performance 

More recently,

Inference time/compute  → Downstream Performance 

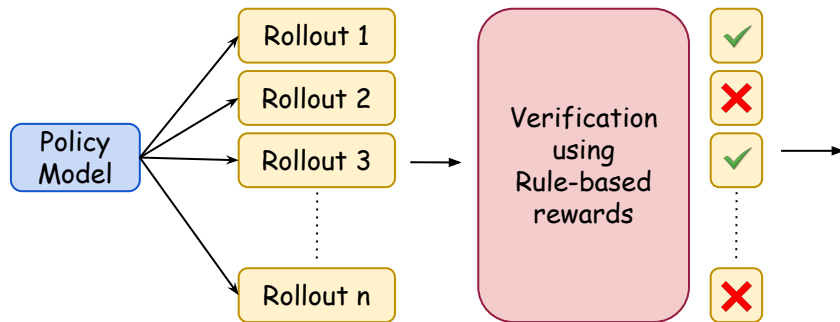
 Models learn to produce ***long reasoning chains***, with sophisticated behaviors like ***self-reflection***, ***self-correction***, and ***multi-step reasoning***.

How?

Reinforcement Learning with Verifiable Rewards (RLVR)!!!

Introduction

Reinforcement Learning with Verifiable Rewards (RLVR) - GRPO variant



$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E} \left[q \sim \mathcal{D}, \{o_i\}_{i=1}^n \sim \mathcal{M}_{\theta_{\text{old}}}(O | q) \right] \left\{ \frac{1}{n} \sum_{i=1}^n \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \min \left[\frac{\mathcal{M}_{\theta}(o_{i,t}|q, o_{i,<t})}{\mathcal{M}_{\theta_{\text{old}}}(o_{i,t}|q, o_{i,<t})} \hat{A}_{i,t}, \right. \right. \\ \left. \left. \text{clip} \left(\frac{\mathcal{M}_{\theta}(o_{i,t}|q, o_{i,<t})}{\mathcal{M}_{\theta_{\text{old}}}(o_{i,t}|q, o_{i,<t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] \right. \\ \left. - \beta D_{KL}[\mathcal{M}_{\theta} \parallel \mathcal{M}_{\text{ref}}] \right\}$$

However, there is a constraint!

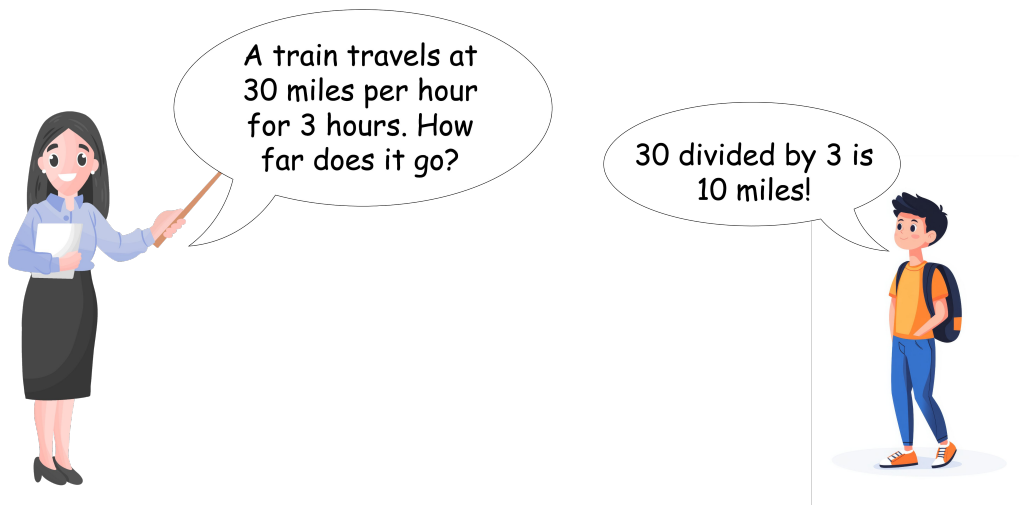
RL can reinforce something only if it is able to sample it in the first place

Indeed, a recent study [\[1\]](#) shows that RL applied on the Llama 3.2-family of models struggles to elicit sophisticated thinking behaviors, in comparison to Qwen models.

Introduction

How can we train the models that don't exhibit such thinking behavior to develop it in the first place?

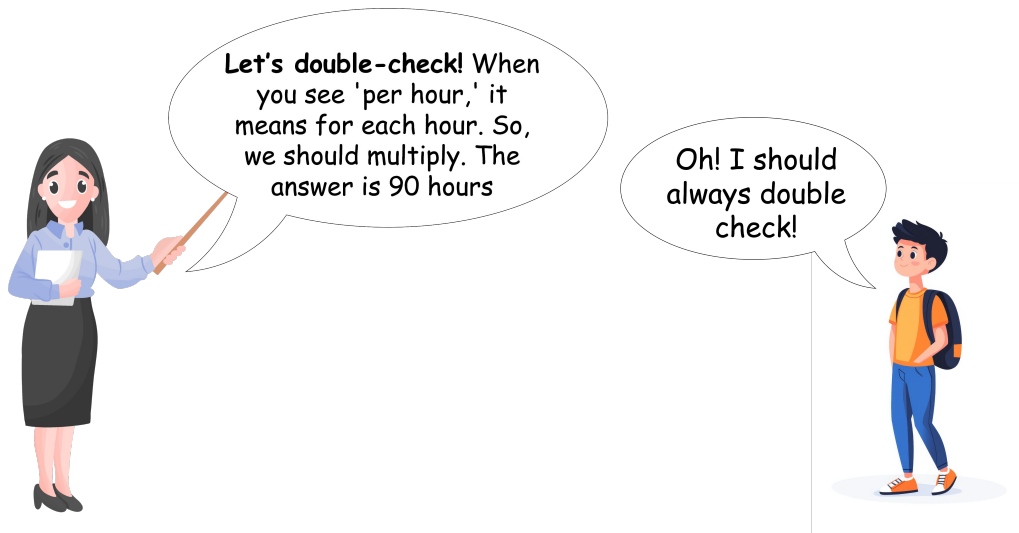
We propose a simple idea inspired from classroom practice!



Introduction

How can we train the models that don't exhibit such thinking behavior to develop it in the first place?

We propose a simple idea inspired from classroom practice!

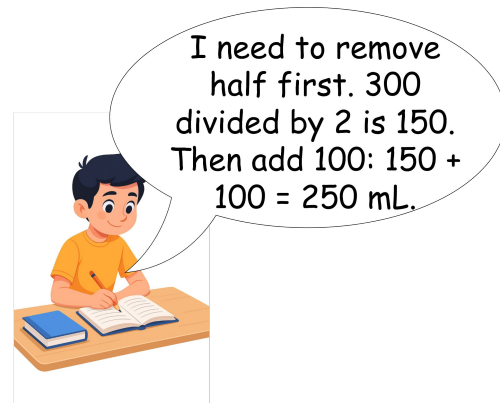
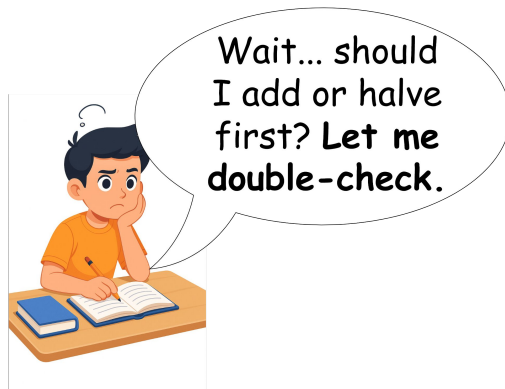
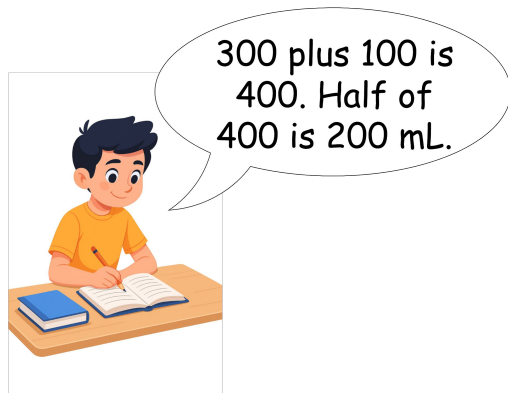


Introduction

How can we train the models that don't exhibit such thinking behavior to develop it in the first place?

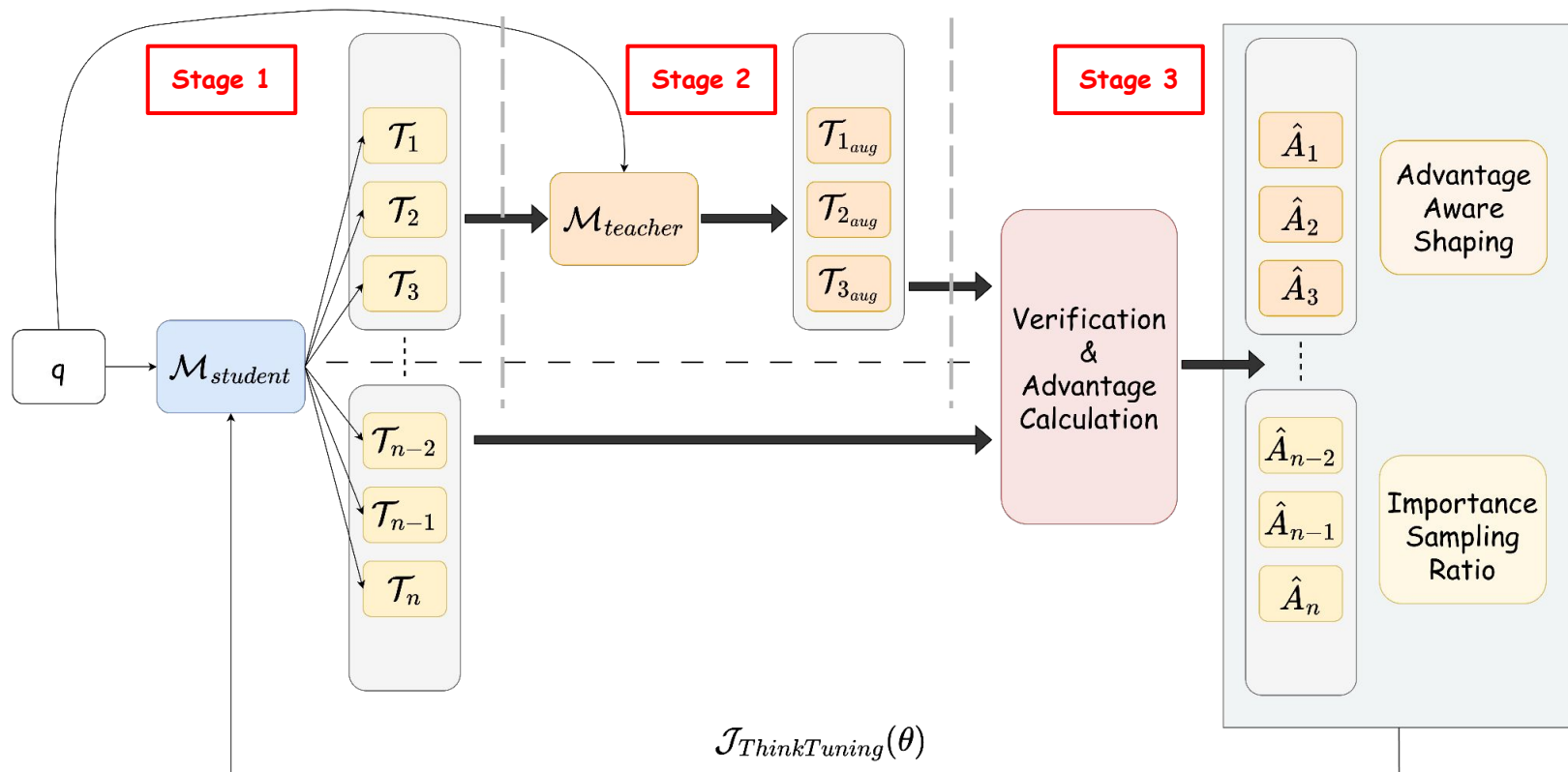
We propose a simple idea inspired from classroom practice!

HomeWork: A tank contains 300 mL of solution. You remove half of it, then add 100 mL of water. How much liquid is in the tank now?



Method

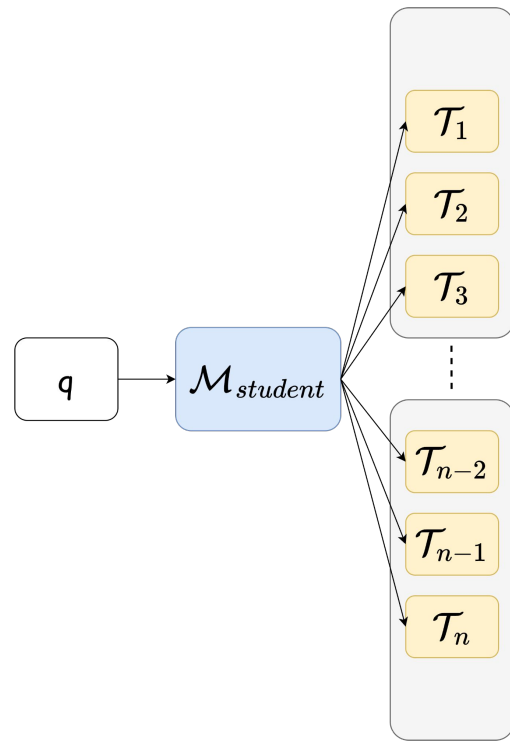
Framework



Stage 1: Student Responses

For a given question \mathbf{q} , we sample \mathbf{n} responses.

We sample responses at a **temperature of 1.0** to observe diversity.



Stage 2: Teacher Guidance

We select γ number of trajectories and pass it to a **teacher model** for guidance.

Guidance Mechanism:

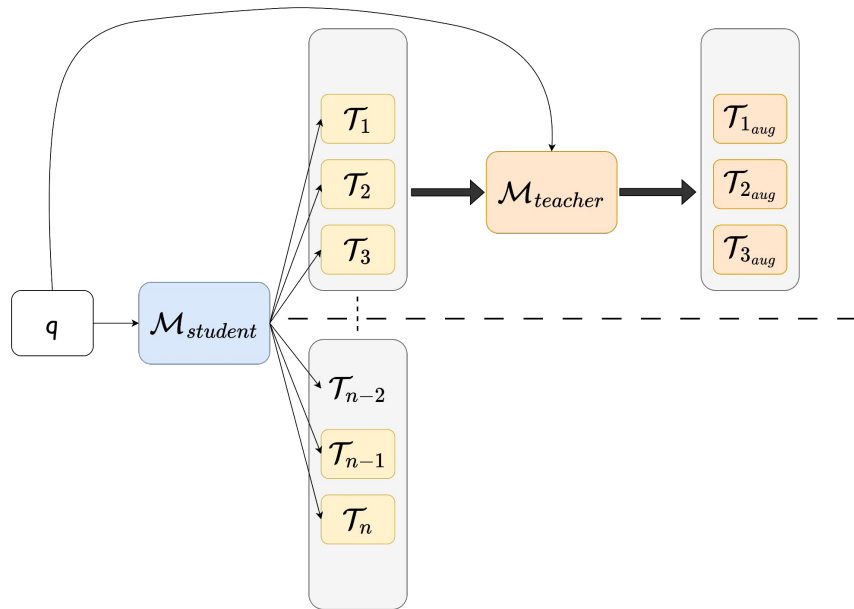
We give few-shot examples of four self-reflective cognitive behaviors (first person perspective):

Self-Conflict: challenging one's own response by presenting conflicting perspectives.

Self-Critique: identifying weaknesses in their response and suggesting improvements.

Self-Agreement: affirming and justifying the strengths in their response.

Self-Consultancy: drawing on an alternative internal perspective or source of expertise to offer new advice or insights that could further improve one's own response.



During Off-Policy updates, importance sampling weight can be very high, leading to massive clipping or gradients (if not clipped)

Stage 3: Student Training

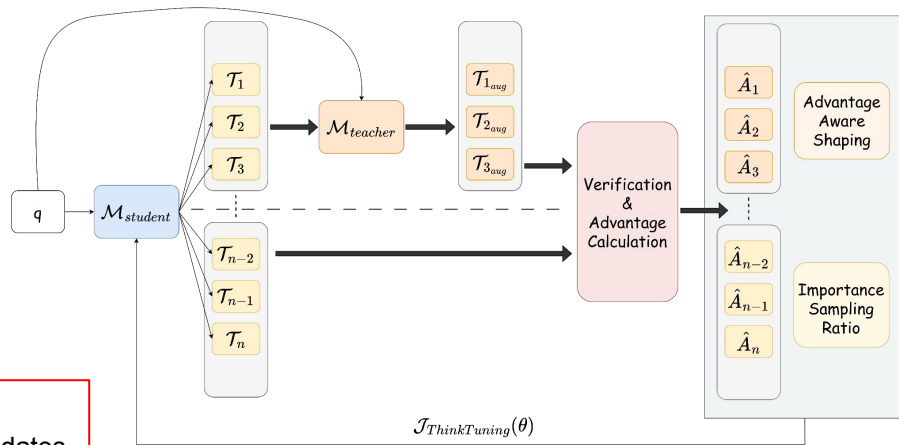
We propose Advantage-Aware-Shaping (AAS) weight for off-policy tokens

$$\frac{\mathcal{M}_{\text{student}}(o_t^{\text{aug}} | q, o_{<t})}{\text{sg}\left(\mathcal{M}_{\text{student}}(o_t^{\text{aug}} | q, o_{<t})\right) + c(\hat{A}_t)}$$

$$c(\hat{A}) \in [-0.0001, 0.0001]$$

$$c(\hat{A}_t) = c_1 + (c_2 - c_1) \cdot \frac{A_{\max} - \hat{A}_t}{A_{\max} - A_{\min}}$$

Serves as a **knob** that enables us to **control the magnitude** of gradient updates for the **off-policy teacher tokens**.



Gradient Analysis (for chosen c values)

$$\nabla_{\theta} w_{\text{aas}} = \frac{\mathcal{M}_{\theta}(o_t^{\text{aug}} | q, o_{<t})}{(\mathcal{M}_{\theta}(o_t^{\text{aug}} | q, o_{<t}) + c(\hat{A}_t))} \nabla_{\theta} \log \mathcal{M}_{\theta}(o_t^{\text{aug}} | q, o_{<t})$$

High Confidence Tokens

When $A > 0$, the gradient increases the token's probability slightly, while preventing overfitting to the already highly confident token.

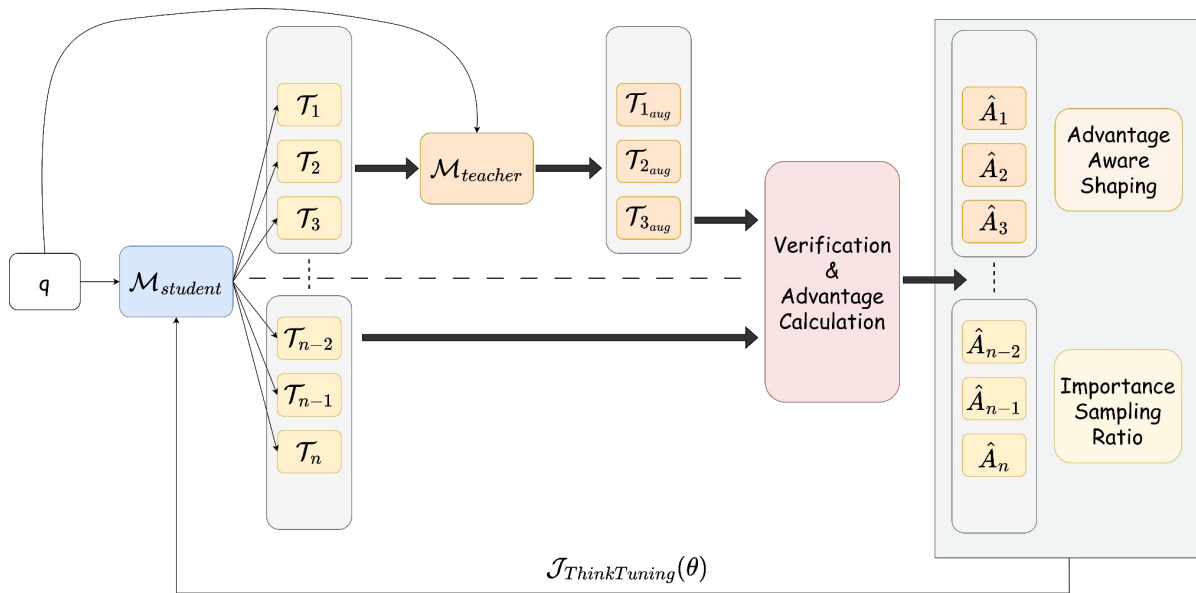
When $A < 0$, the gradient magnitude is slightly higher than vanilla gradient update, making the model to reduce the probability of this already highly confident token.

Low Confidence Tokens

When $A > 0$, the gradient pushes to increase the probability of this token, conservatively in comparison to the vanilla gradient update.

When $A < 0$, beyond a confidence threshold, we mask the low confidence negative advantaged tokens, to not unnecessarily decrease the probability of this token.

Overall Algorithm



Algorithm 1 THINKTUNING

```

1: Input: Initial Student model  $\mathcal{M}_{student_{\theta_{init}}}$ , Teacher
   model  $\mathcal{M}_{teacher}$ , guidance fraction  $\gamma$ , hyperparameter
   set  $(\epsilon, \beta, c_1, c_2, k)$ 
2:
3:  $\mathcal{M}_{student_{\theta}} \leftarrow \mathcal{M}_{student_{\theta_{init}}}$ 
4:
5: for training step=1 to  $I$  do
6:    $\mathcal{M}_{student_{old}} \leftarrow \mathcal{M}_{student_{\theta}}$ 
7:   Sample mini-batch  $\mathcal{D}_b \subset \mathcal{D}$ 
8:
9:   // Student acts & Teacher helps
10:  for all  $q \in \mathcal{D}_b$  do
11:    if training step  $\leq k$  then
12:       $\{o\}_{i=1}^n \sim \text{Guide}(q, \mathcal{M}_{student_{old}}, \mathcal{M}_{teacher}, \gamma)$ 
13:    else
14:       $\{o\}_{i=1}^n \sim \mathcal{M}_{student_{old}}(O | q)$ 
15:    end if
16:  end for
17:
18:  // Reward calculation and Advantage estimation
19:  Compute the rewards  $r_i = r(o_i)$  for each response
   Compute group-normalized advantage  $\hat{A}_{i,t}$  for all to-
  kens
20:
21:  for mini-batch step = 1 to  $\mu$  do
22:    if training step  $\leq k$  and  $o_i \in \mathcal{T}_{aug}$  then
23:      Calculate  $w_{aas}(o_i)$ 
24:    else
25:      Calculate  $w(o_i)$ 
26:    end if
27:     $\mathcal{M}_{student_{\theta}} \leftarrow \text{argmax}_{\theta} \mathcal{J}_{THINKTUNING}(\theta)$ 
28:  end for
29: end for
30: Output: Final think-tuned model  $\mathcal{M}_{student_{\theta}}$ 

```

Experiment & Results

Experimental Setup

Baselines

Prompt-based Methods

Zero-Shot

Self-Verify

Self-Correct

S1-Budgeting (Token Budget = 2048)

Training-based Methods

SFT

STaR

RL Baseline

GRPO

Model

Llama-3.2-3B-Instruct

Benchmarks

GSM8k

Math-500

AIME

CommonSenseQA

ARC-Challenge

GPQA-Diamond

StrategyQA

MMLU-Pro

Main Results

Methods	Mathematical Reasoning			CommonSense Reasoning	Scientific Reasoning		Other Reasoning	
	GSM8K	MATH-500	AIME	CSQA	ARC-C	GPQA-D	STRATEGYQA	MMLU-PRO
Zero-Shot-CoT	71.08 \pm 0.20	38.14 \pm 0.75	9.32 \pm 0.36	67.39 \pm 0.26	75.49 \pm 0.20	25.10 \pm 0.85	66.40 \pm 0.43	34.41 \pm 0.11
Self-Verify	52.08 \pm 1.73	34.98 \pm 0.54	8.19 \pm 0.29	54.41 \pm 0.73	61.56 \pm 0.47	23.94 \pm 0.68	52.10 \pm 0.39	28.10 \pm 0.14
Self-Correct	51.45 \pm 0.30	32.46 \pm 0.47	7.81 \pm 0.18	45.90 \pm 0.69	52.88 \pm 0.58	24.60 \pm 0.71	52.39 \pm 0.78	25.50 \pm 0.12
s1-budgeting	51.30 \pm 0.42	25.72 \pm 0.54	9.01 \pm 0.31	54.21 \pm 0.44	59.51 \pm 0.27	26.57 \pm 0.99	57.88 \pm 0.80	28.59 \pm 0.10
SFT	62.27 \pm 0.61	29.00 \pm 0.49	6.07 \pm 0.43	65.91 \pm 0.24	70.90 \pm 0.71	24.49 \pm 0.82	64.12 \pm 0.65	36.07 \pm 0.07
STaR	73.54 \pm 0.22	40.78 \pm 0.35	8.91 \pm 0.29	67.91 \pm 0.30	77.24 \pm 0.21	21.46 \pm 0.86	66.84 \pm 0.41	34.69 \pm 0.12
GRPO	78.89 \pm 0.84	45.46 \pm 1.55	12.03 \pm 0.33	69.86 \pm 0.52	79.13 \pm 0.21	24.19 \pm 0.75	70.68 \pm 0.35	36.07 \pm 0.07
THINKTUNING	74.22 \pm 0.13	47.54 \pm 0.46	14.26 \pm 0.38	70.43 \pm 0.19	79.80 \pm 0.24	28.18 \pm 0.63	66.52 \pm 0.41	37.21 \pm 0.11

Comparison with Prompt-based methods: ThinkTuning outperforms all prompting-based baselines, achieving up to +9.4% accuracy gains across diverse reasoning benchmarks.

Comparison with Training-Based methods: SFT harms reasoning, reducing accuracy by up to 8% even on in-domain GSM8k. While STaR recovers some performance over SFT, our method surpasses it across all benchmarks with up to +6.8% gains with respect to zero-shot baseline.

Comparison with GRPO: Our method attains the best score on six of eight datasets, matches or surpasses GRPO on every set except GSM8K and StrategyQA, and delivers the largest absolute gain on AIME and GPQA-Diamond

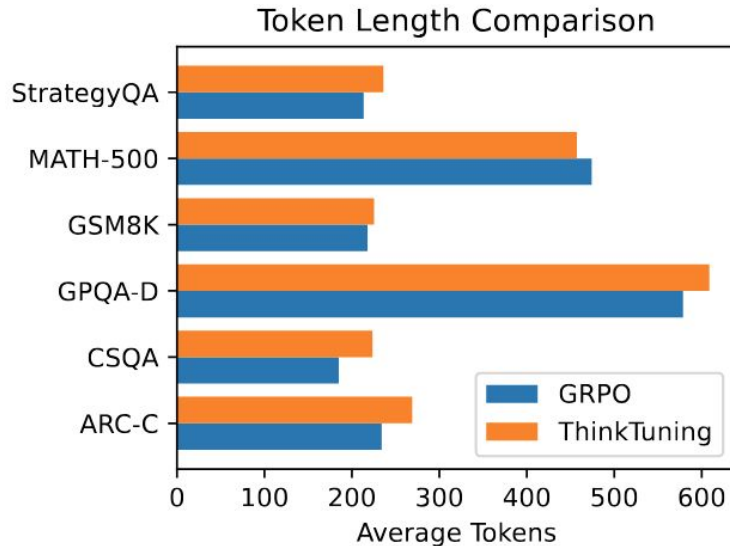
Analysis

Does *ThinkTuning* scale inference time?

ThinkTuning model generates longer responses (3.4–20.8% more tokens) across benchmarks, indicating deeper reasoning.

Increased token usage reflects boosted performance on complex tasks like GPQA-D and MATH-500.

In MATH-500, ThinkTuning model outperforms GRPO even with fewer tokens (−3.6%), showing more efficient reasoning.



Can *ThinkTuning* instill unknown behaviors?

Consider a trivial use case in which we want the model to end its response by quoting the South Indian actor “**Thalapathy Vijay**”.

During RL training on math problems, it is **highly unlikely** that the model explores a trajectory that ends its response with a movielike dialogue quoting a **South Indian actor**. Note that, **we do not give** anything instructing the model to respond this particular way during sampling.

$$r(o_i) = \begin{cases} 1.0 & \text{if the final answer is correct} \\ 0.0 & \text{otherwise} \end{cases} + \begin{cases} 0.5 & \text{if "Thalapathy Vijay" appears} \\ & \text{in the answer.} \\ 0.0 & \text{otherwise} \end{cases}$$

Question: What island country is ferret popular?

Choices: A. own home, B. north carolina C. great britain D. hutch E. outdoors

GRPO ❌

To answer this question, we need to identify the island country where ferrets are popular.

Step 1: Analyze the options provided.

A. own home: This is not an island country.

B. North Carolina: This is a state in the United States, not an island country.

C. Great Britain: This is an island country, which includes England, Scotland, and Wales.

D. hutch: This is a container for keeping animals, not an island country.

E. outdoors: This is a place where ferrets can be found, but it's not an island country.

Step 2: Based on the analysis, the correct answer is the island country where ferrets are popular.

Step 3: Choose the correct answer. \boxed{C}

ThinkTuning ✅

To answer this question, let's evaluate each option in relation to the question about an island country where ferrets are popular.

A. own home - This is not an island country.

B. north carolina - This is a state in the United States and not an island country.

C. great britain - Great Britain is an island country, and ferrets are indeed popular there.

D. hutch - A hutch is a structure for keeping animals, not a country.

E. outdoors - This is a place where ferrets might be kept but not a country.

So, the correct answer is \boxed{C}. **Thalapathy Vijay would have said "The moment you make a connection, the answer becomes clear."** I connected ferrets with Great Britain, where they are indeed popular.

Conclusion

Our Contributions

- We introduced **ThinkTuning**, a GRPO-based interactive training framework that **instills cognitive reflections via guided exploration**.
- Empirically, ThinkTuning boosts the performance of a Llama-3.2-3B-Instruct model that was trained only on questions from the GSM8K train split.
- Additional experiments reveal that ThinkTuning **can elicit unknown behaviors**, by **steering/guiding the exploration** on-the-fly during RL.
- We hope our work will inspire future research that employs **large-scale interactive training frameworks**.

End