



# DeepSeek & GRPO

Ran Tavory 2025

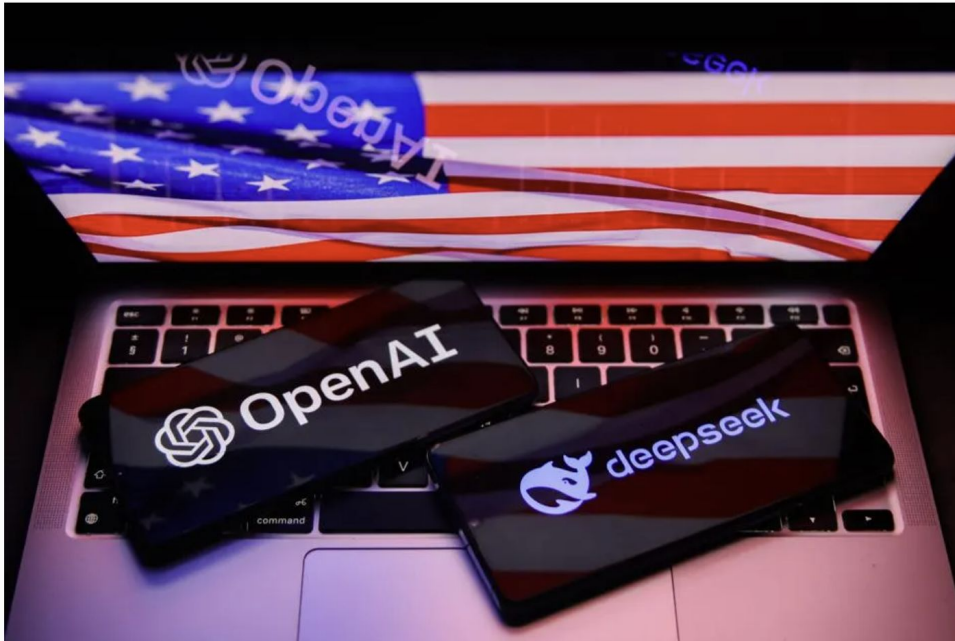
# The \$6 Million AI Bombshell: How DeepSeek Shook Wall Street And AI Leadership

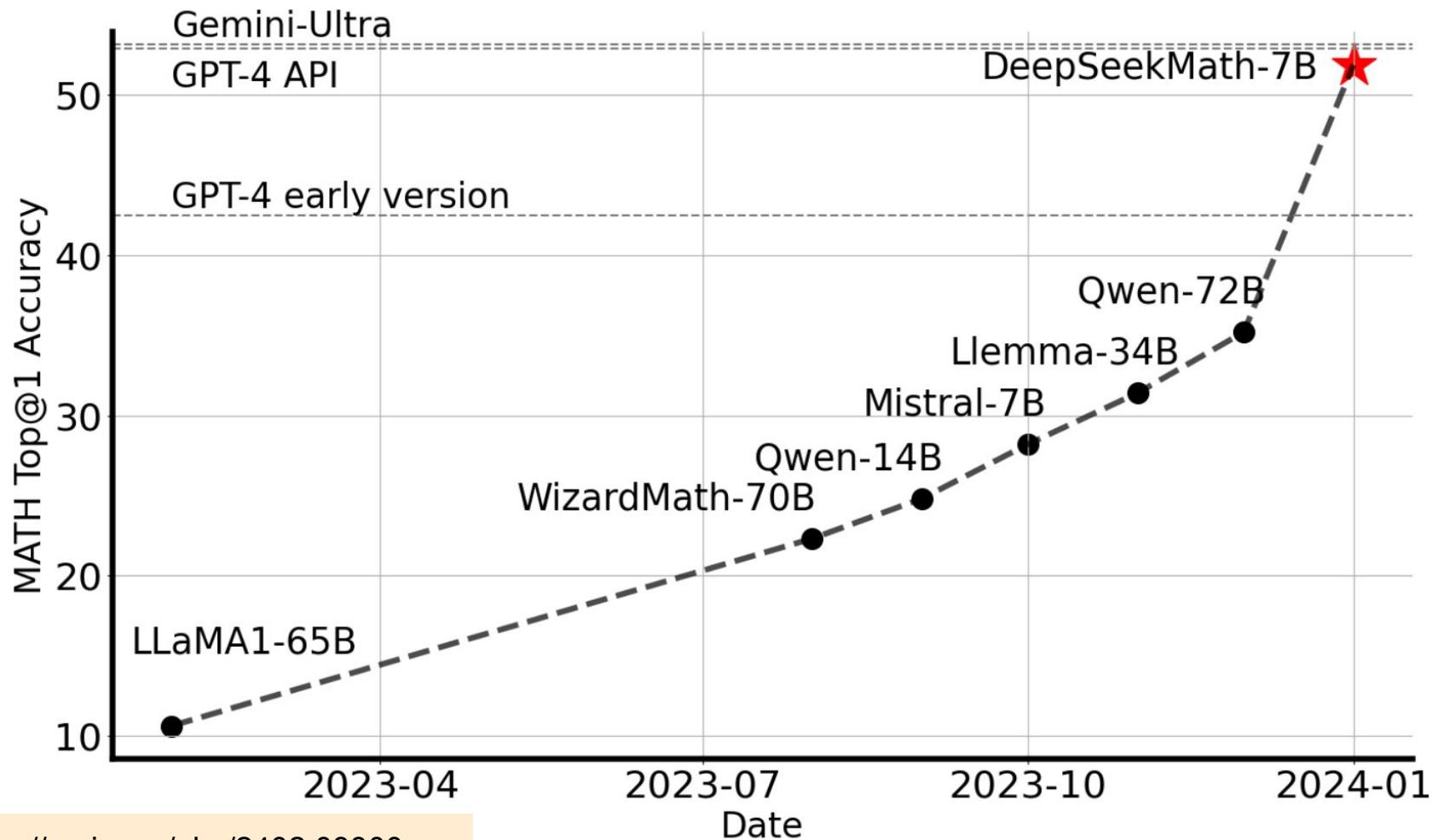
By [Mark Minevich](#), Contributor. ⓘ Mark Minevich is a NY-based strategist foc...

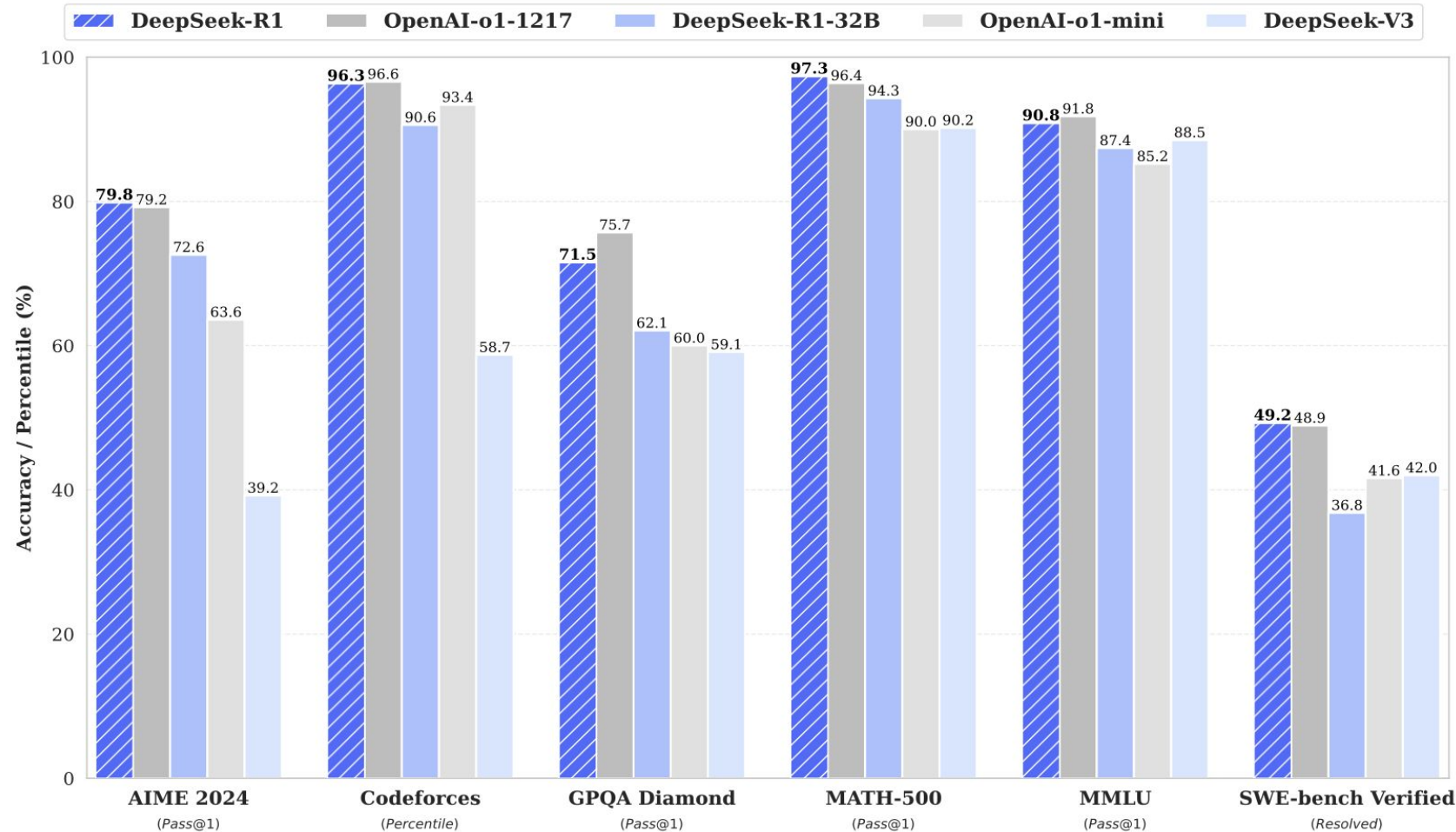
[Follow Author](#)

Feb 06, 2025, 01:41pm EST

[Share](#) [Save](#) [Comment](#) 0





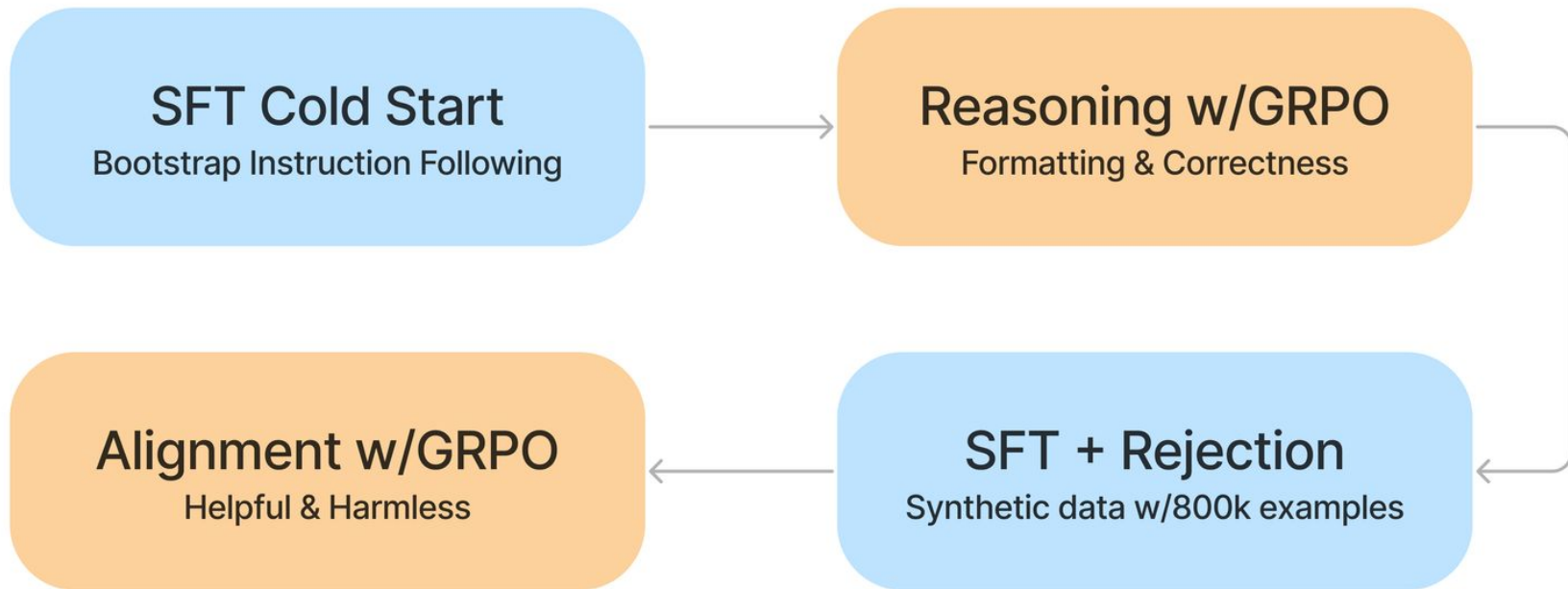


# 02

What does GRPO has to do  
with it



# DeepSeek-R1 Training Pipeline



# RLHF

## Reinforcement Learning from Human Feedback

Step 1

**Collect demonstration data, and train a supervised policy.**

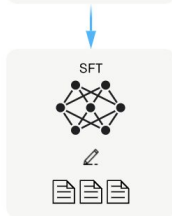
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



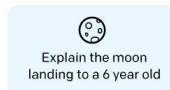
This data is used to fine-tune GPT-3 with supervised learning.



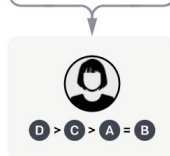
Step 2

**Collect comparison data, and train a reward model.**

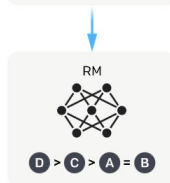
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



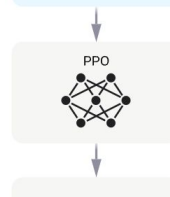
Step 3

**Optimize a policy against the reward model using reinforcement learning.**

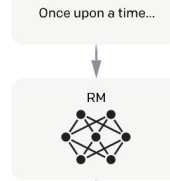
A new prompt is sampled from the dataset.



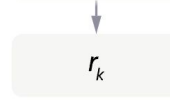
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



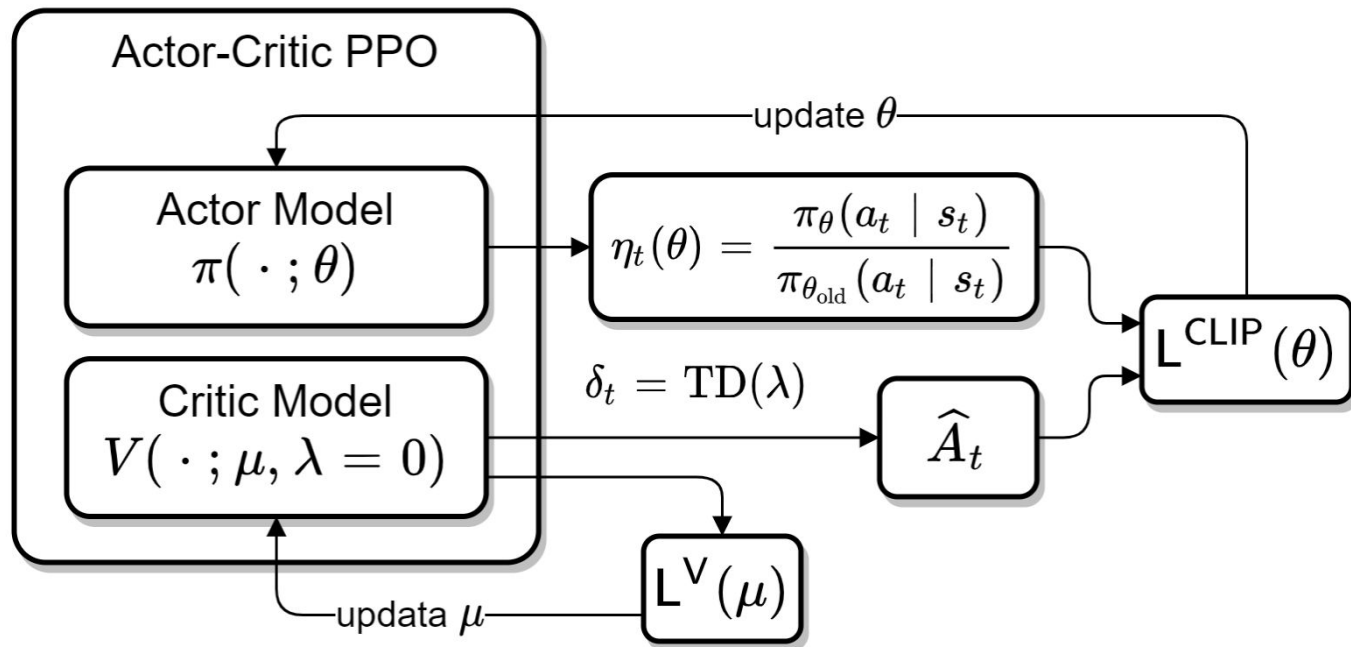
03

PPO first



# PPO

Proximal Policy Optimization



## Reminder

Step 3

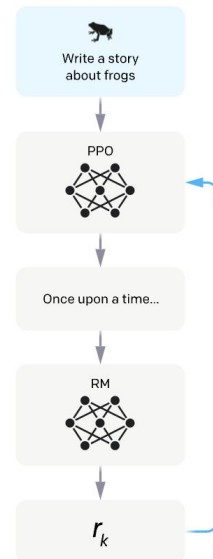
**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

The policy generates an output.

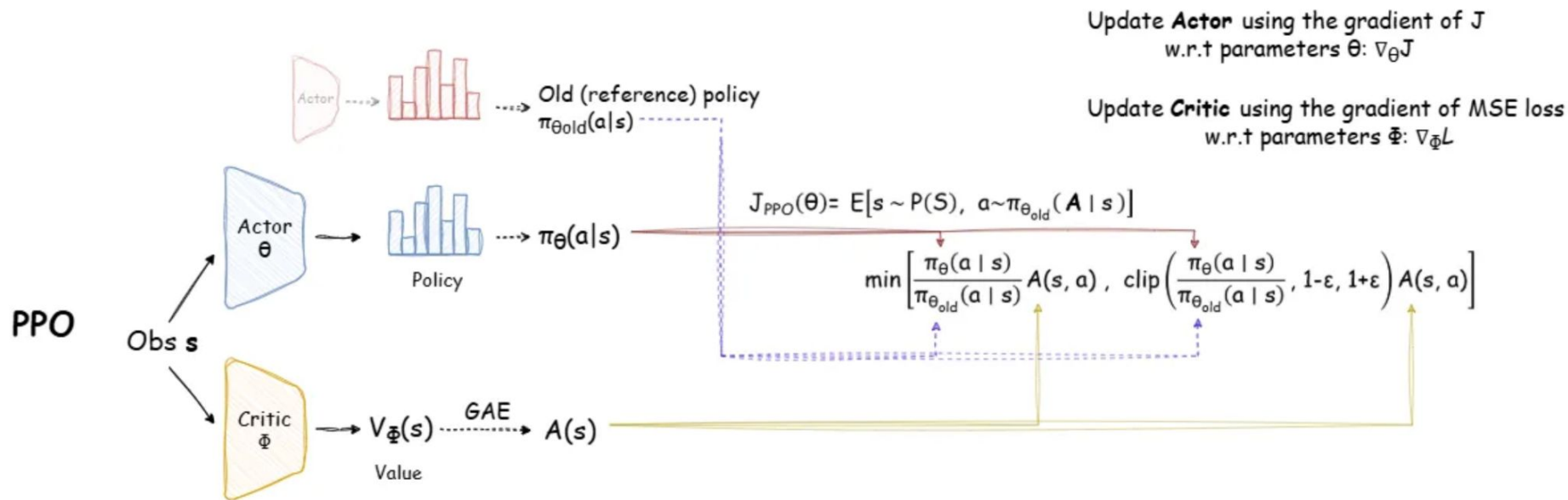
The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.



# PPO

Another visual representation



## PPO Summary

- Effective and battle tested
- But requires a lot of memory:
  - Two trainable models (2x4 params)
  - Two frozen models (2x params)
  - => ~10x LLM memory
- Sample efficient.
  - Useful when samples are expensive
  - Example: Lunar lander



# DPO

Direct Preference Optimization

Out of scope for today



AHA!

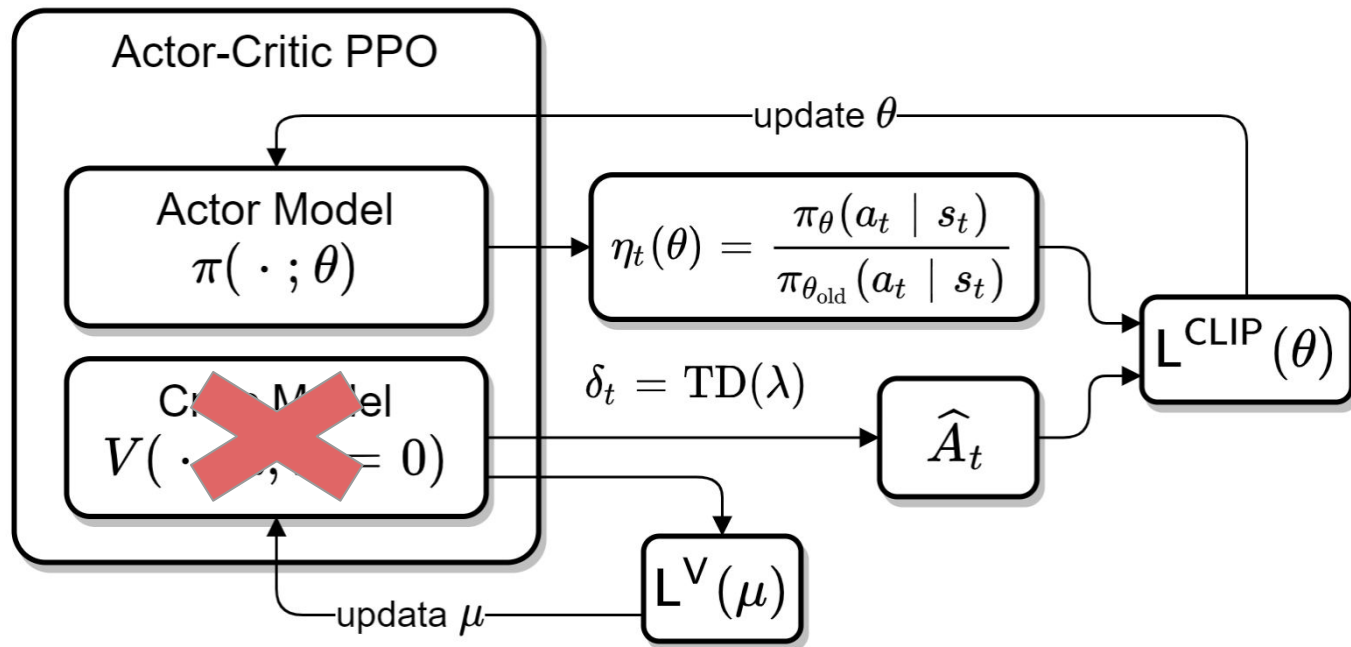
AHA!

04

GRPO

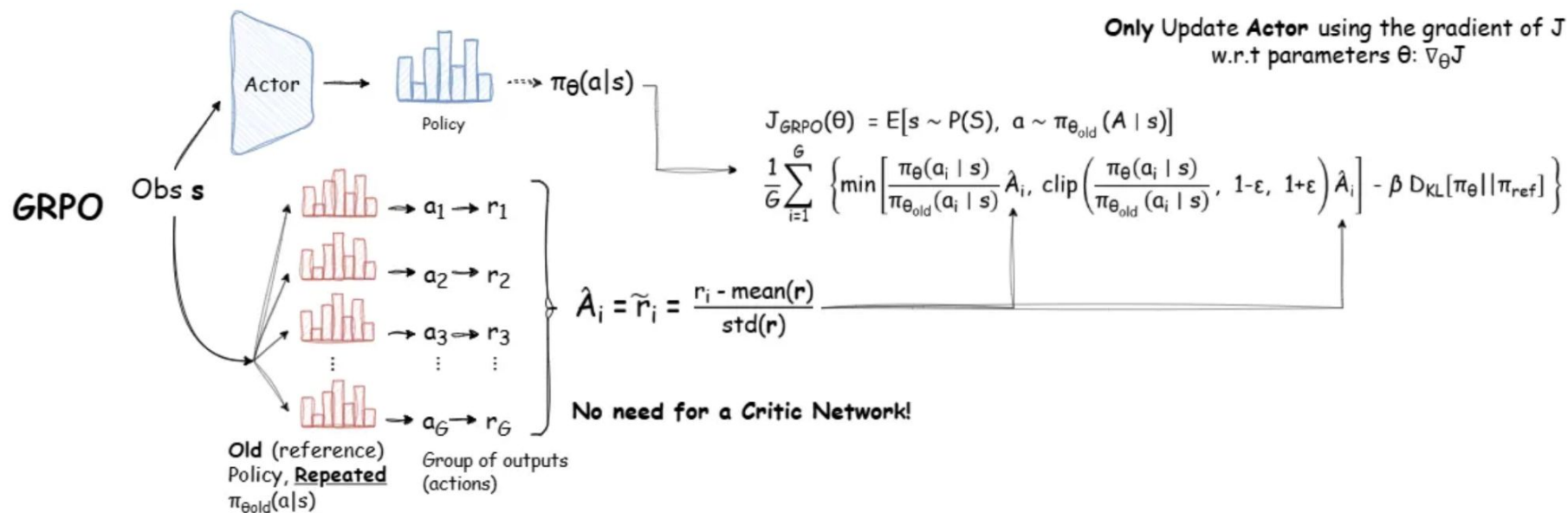
## GRPO

Group Relative Policy Optimization



# GRPO

## Group Relative Policy Optimization



- No need for a Value model (which is 4xLLM)
- The Advantage is calculated for each result relative to all other members of its group (hence the name Group Relative)



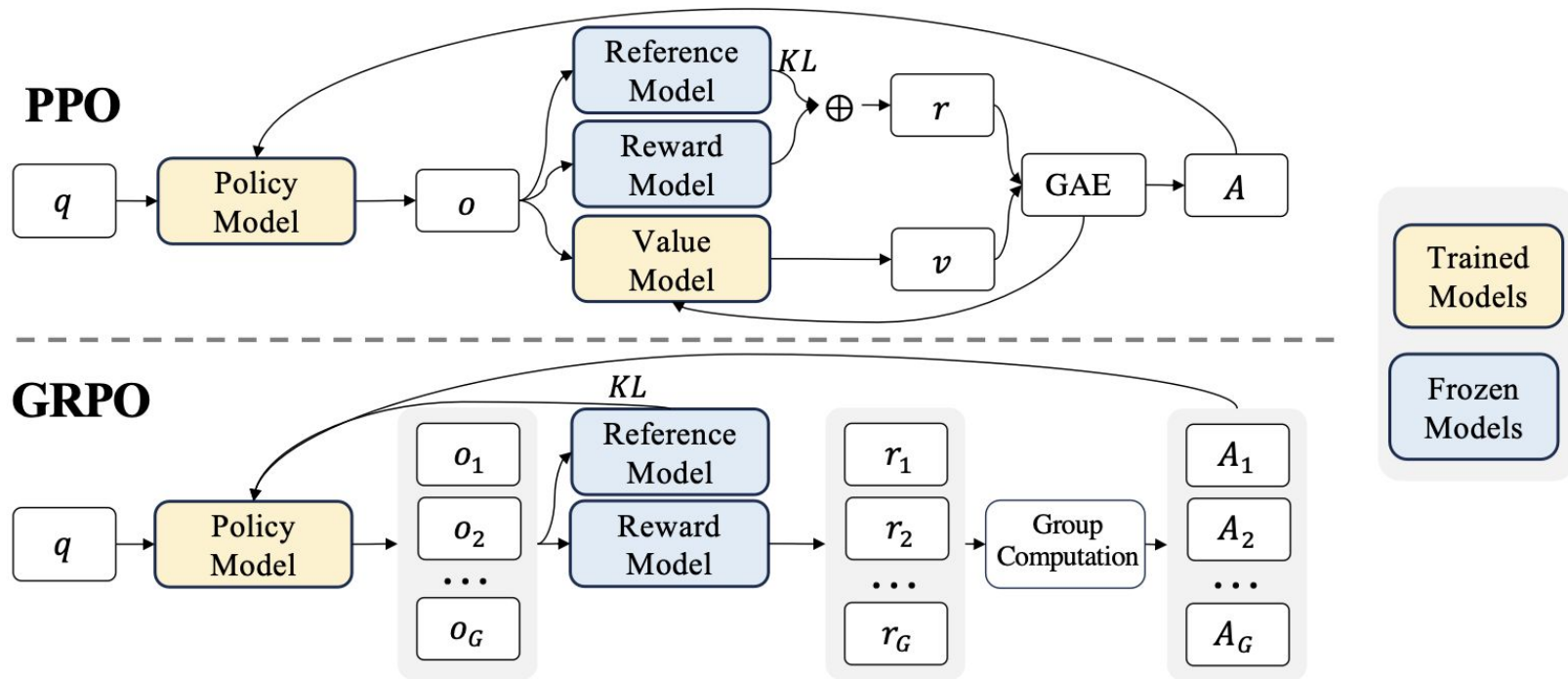
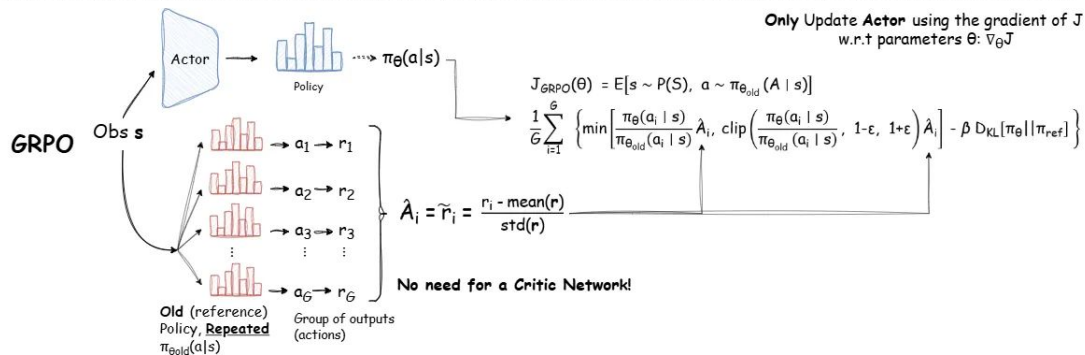
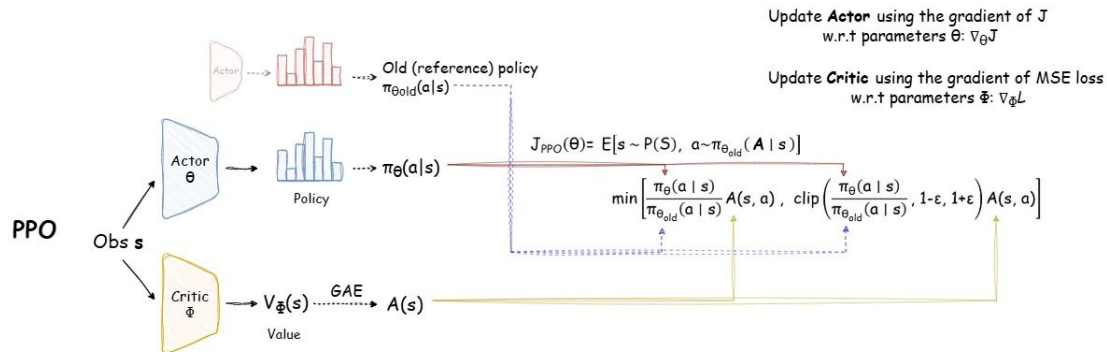


Figure 4 | Demonstration of PPO and our GRPO. GRPO foregoes the value model, instead estimating the baseline from group scores, significantly reducing training resources.

# GRPO vs PPO



## RL Terms

### Common RL terms

- **policy**: in deep RL, the policy is the neural network that looks at the observations from the environment and predicts the action to be taken. this is what we are training.
- **environment**: the world that the policy interacts with, which could be anything from a video game to a robotic simulation to a real-world setting. the environment provides observations and rewards.
- **reward**: a numerical value that tells the agent how good or bad its last action or sequence of actions was.
- **trajectory**: the sequence of states, actions, and rewards that occur during an episode or training run.

## How RL maps to LLMs

Mapping RL concepts to Text Generation

- RL episode == LLM full generation
- RL time step == LLM single token generation
- RL reward == LLM Reward model or RLVR
- so you judge the entire response as one, and define a reward value based on that judgement
- Rewards are sparse in LLMs. Only one reward at the end of the episode (generation).



05

RLVR

## RLVR

Reinforcement Learning from Verifiable Rewards

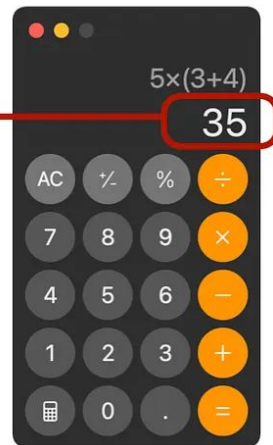
**Input:**

Solve  $5 \times (3 + 4)$

**LLM output:**

$5 \times (3 + 4) =$   
 $5 \times 7 =$   
 $= 35$

**Verifier output:**



**Correct**

## RLVR & RLHF

### Comparison

---

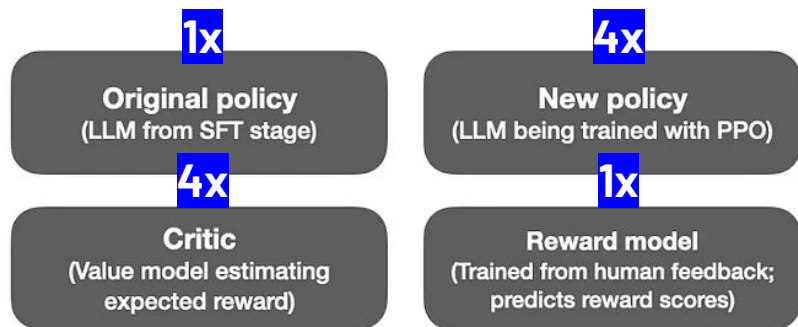
RLHF

- **Requires a reward model (costs more)**
  - **Reward model is trained from real humans feedback**
  - **Susceptible to Reward Hacking**
- 

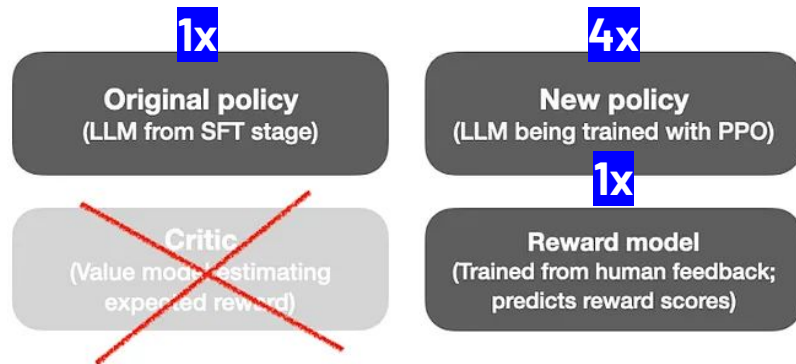
RLVR

- **Rewards are provided by “simple” tools. E.g. a calculator or a compiler**
  - **Requires structured output**
-

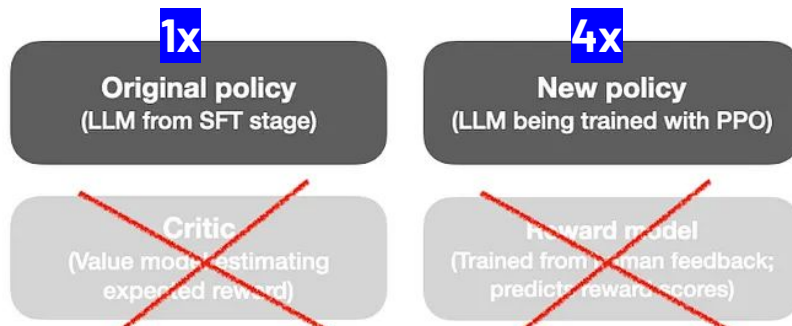
## RLHF with PPO



## RLHF with GRPO



## RLVR with GRPO





$O^*$

Math

## The Math

Loss functions

PPO's clipped surrogate loss function

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[ \min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right]$$

GRPO's loss function

$$J_{\text{GRPO}}(\theta) = \mathbb{E}_{q, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \min \left( r_t(\theta) \hat{A}_{i,t}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_{i,t} \right) \right] - \beta D_{\text{KL}}(\pi_{\theta} || \pi_{\text{ref}}).$$

0\*

KL Divergence and Reference  
Policy

## Why do we need the Original model?

Answer:

- In PPO/GRPO the original model is used in order to make sure that the new model remains “in the same neighbourhood”

### RLVR with GRPO



$$\mathcal{J}_{GRPO}(\theta) = \mathbb{E}[q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{old}}(O|q)]$$

$$\frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min \left[ \frac{\pi_{\theta}(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q, o_{i,<t})} \hat{A}_{i,t}, \text{clip} \left( \frac{\pi_{\theta}(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{old}}(o_{i,t}|q, o_{i,<t})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t} \right] - \beta \mathbb{D}_{KL} [\pi_{\theta} || \pi_{ref}] \right\},$$





# THANK YOU

Ran Tavory, TII

*Thank you!*