

# RL for LLM Alignment and Inference

Sharif University

2025-05-29

Pascal Poupart, CIFAR AI Chair

David R. Cheriton School of Computer Science

Vector Institute



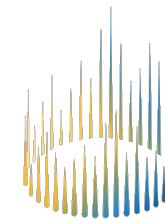


## Research topics

- **Machine Learning:** reinforcement learning, uncertainty quantification, federated learning, inverse constraint learning
- **Natural Language Processing:** LLM alignment and inference, agentic LLMs, knowledge graphs, post-editing ASR error correction
- **Applications:** autonomous driving, sports analytics, material design for CO2 recycling



## Industry Partners



# Outline

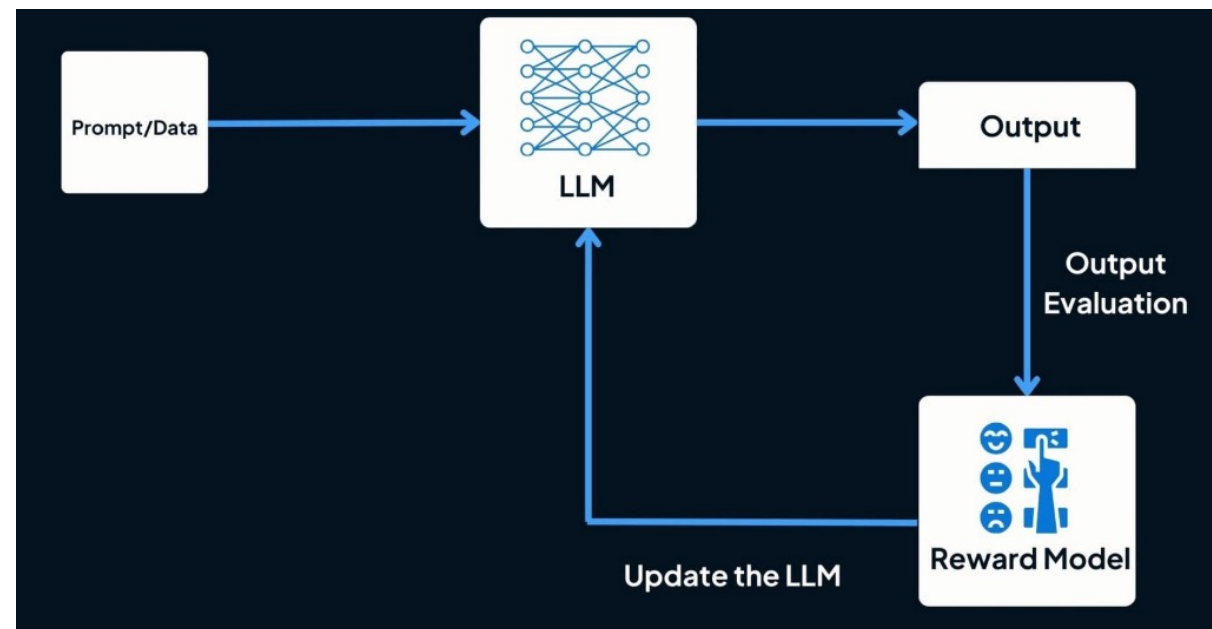
- LLM Alignment
  - Reinforcement Learning from Human Feedback
  - Direct Preference Optimization
  - Reward Guided Text Generation
- LLM Reasoning
  - Search and planning
  - Group Relative Policy Optimization (GRPO)
  - Reflection: Verbalized RL

# Large Language Models

- **Agent:** system
- **Environment:** user
- **State:** history of past utterances
- **Action:** system utterance
- **Reward:** task completion, human feedback

*“We posit that the superior writing abilities of LLMs, as manifested in surpassing human annotators in certain tasks, are fundamentally driven by RLHF, as documented in Gilardi et al. (2023) and Huang et al. (2023).”*

Llama 2 Technical Report (Touvron et. al 2023)



Credit: <https://www.twine.net/blog/what-is-reinforcement-learning-from-human-feedback-rlhf-and-how-does-it-work/>

*“This behavior (re-evaluation) is not only a testament to the model’s growing reasoning abilities but also a captivating example of how reinforcement learning can lead to unexpected and sophisticated outcomes”*  
Deepseek-R1 Technical Report (Guo et. al 2025)

# Recipe For Large Language Model Training

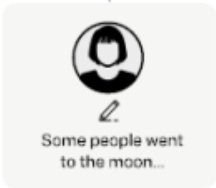
Step 1

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

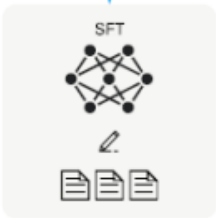
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



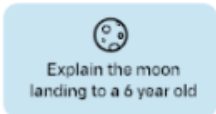
This data is used to fine-tune GPT-3 with supervised learning.



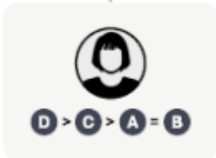
Step 2

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

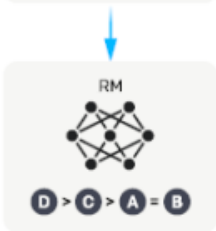
A prompt and several model outputs are sampled.



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.



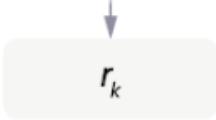
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.



InstructGPT (Ouyang et. al 2022)

# Response Ranking

PP

Explain reinforcement learning in three sentences.

< 1 / 2 >



Reinforcement learning is a type of machine learning where an agent learns to make decisions by receiving feedback in the form of rewards or punishments. The agent's goal is to maximize its cumulative reward over time by exploring different actions and learning which actions lead to the highest rewards. Reinforcement learning has been successfully applied to a variety of tasks, including game playing, robotics, and recommendation systems.



< 2 / 2 >



Reinforcement learning is a type of machine learning that focuses on training an agent to make decisions based on trial and error. The agent receives feedback in the form of rewards or penalties for each decision it makes. By learning from this feedback, the agent can improve its decision-making abilities over time.



Was this response better or worse?



Better



Worse



Same



# RL from Human Feedback (RLHF)

- Collect a preference data set:

$$D = \{(\mathbf{s}, \mathbf{a}_+, \mathbf{a}_-)_k\}_{k=1}^K \text{ where } \mathbf{a}_+ \succ \mathbf{a}_-$$

- Train a reward model according to the Bradley Terry Model:

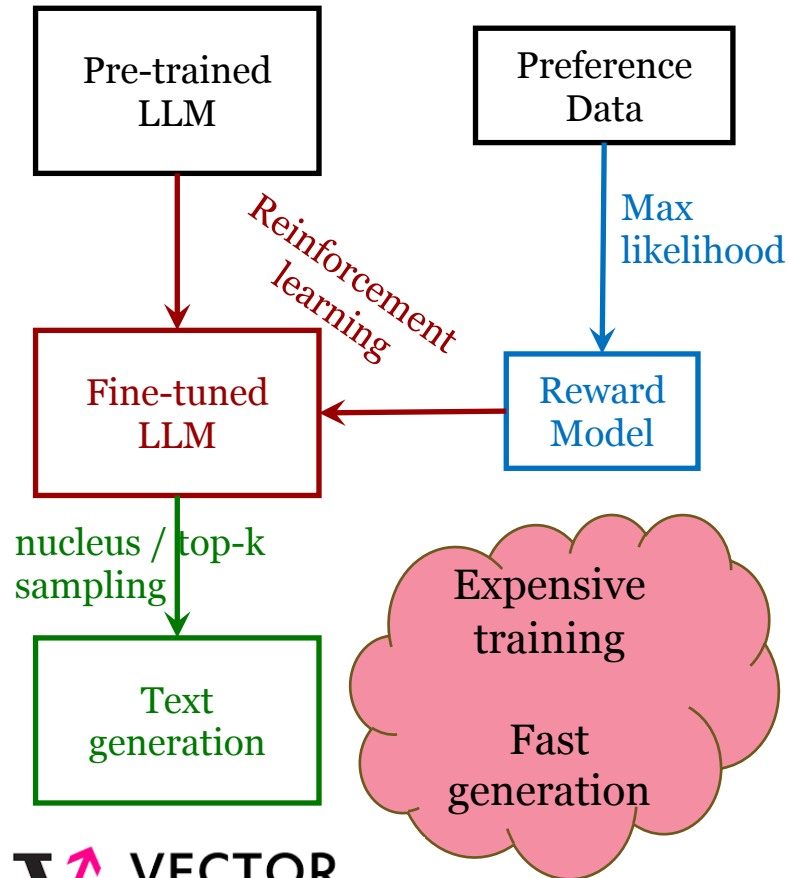
$$\max_{\theta} E_D [\log \sigma(r_{\theta}(\mathbf{s}, \mathbf{a}_+) - r_{\theta}(\mathbf{s}, \mathbf{a}_-))]$$

- Make a copy of the LLM and finetune it to maximize:

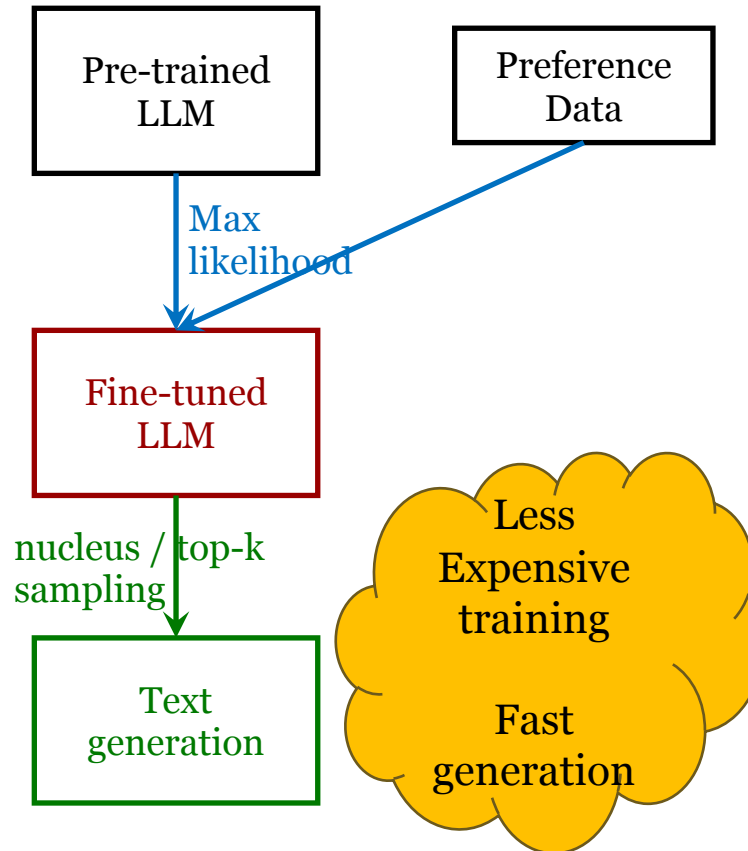
$$\max_{\phi} E_{D, \pi_{\phi}} [r_{\phi}(\mathbf{s}, \mathbf{a})] - \beta KL[\pi_{\phi}(\mathbf{a}|\mathbf{s}) || \pi_{pretrained}(\mathbf{a}|\mathbf{s})]$$

# RLHF Improvements

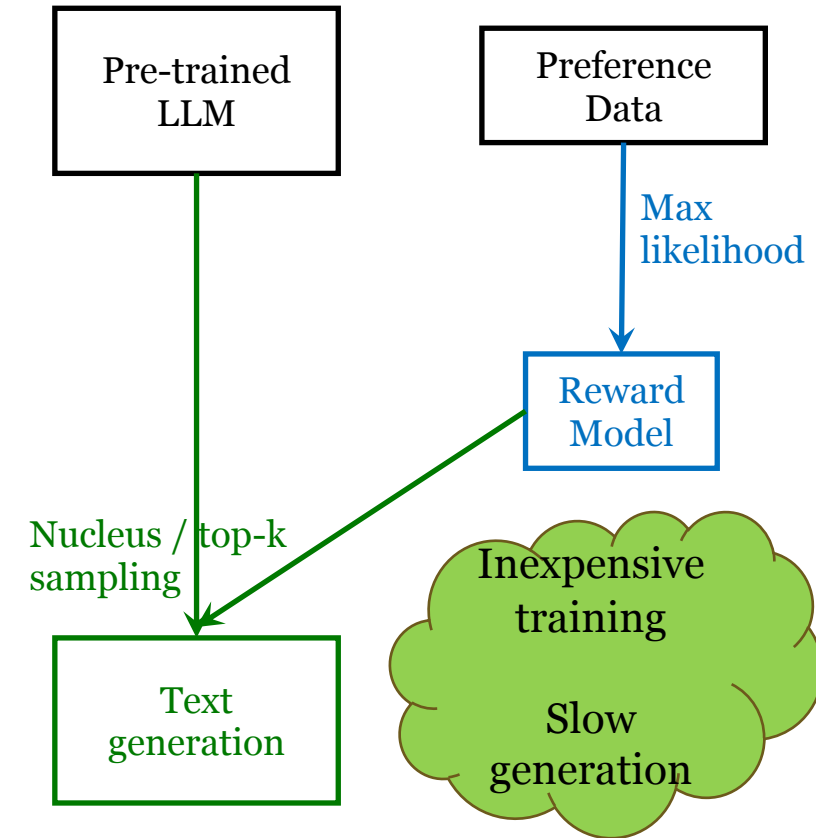
## Proximal Policy Optimization (PPO) Ouyang et al., 2022



## Direct Preference Optimization (DPO) Rafailov et al., 2023



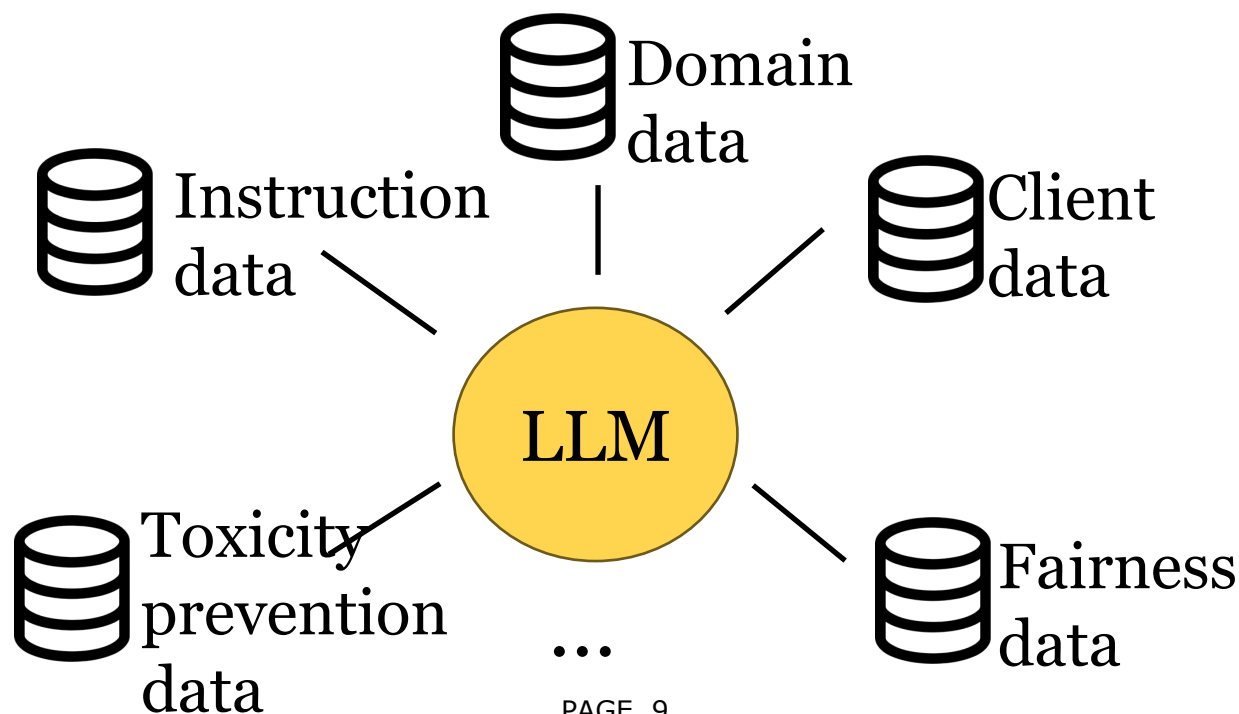
## Reward Guided Text Generation (RGTG) Khanov et al., 2024 Rashid et al., 2025





# LLM Alignment with Preference Data

- Collect preference data:  $D = \{(s, a_+, a_-)_k\}_{k=1}^K$   
where  $s$ : user prompt       $a$ : system response  
 $a_+$  is preferred to  $a_-$  (i.e.,  $a_+ \succ a_-$ )



# Reward Model

Stiennon, Ouyang, Wu, Ziegler, Lowe Voss, Radford, Amodei, Christiano (2020) **Learning to summarize from human feedback**, *NeurIPS*.

- Reward function:  $r_{\theta}(s, a) = \text{real number}$
- Consider several possible responses  $a_1 \succcurlyeq a_2 \succcurlyeq \dots \succcurlyeq a_k$  ranked by annotator
- Training reward function to be consistent with the ranking:

$$Loss(\theta) = -\frac{1}{\binom{k}{2}} E_{(s, a_i, a_j) \in Dataset} \log \sigma \left( r_{\theta}(s, a_i) - r_{\theta}(s, a_j) \right)$$

# Reinforcement Learning

Ouyang, Wu, Jiang, Wainwright, et al. (2022) **Training language models to follow instructions with human feedback**, *NeurIPS*.

- Pretrain language model (GPT-3)
- Fine-Tune GPT-3 by RL to obtain InstructGPT
  - Policy (language model):  $\pi_\phi(a|s)$
  - Optimize  $\pi_\phi(s)$  by Proximal Policy Iteration (PPO)

$$\max_{\phi} E_{s \in Dataset} \left[ E_{a \sim \pi_\phi(a|s)} [r_\theta(s, a)] - \beta KL(\pi_\phi(\cdot | s) | \pi_{ref}(\cdot | s)) \right]$$

# Policy Optimization

Stochastic policy  $\pi_{\phi}(a|s) = \Pr(a|s; \phi)$  parametrized by  $\phi$ .

	Supervised Fine-Tuning	Reinforcement Learning
Data	$\{(s_1, a_1^*), (s_2, a_2^*), \dots\}$ ( $a^*$ denotes optimal action)	$\{(s_1, a_1, r_1), (s_2, a_2, r_2), \dots\}$ ( $r$ denotes reward for $s, a$ pair)
Objective	Maximum likelihood $\max_{\phi} \sum_n \log \pi_{\phi}(a_n^*   s_n)$	Maximum expected rewards $\max_{\phi} \sum_n \gamma^n E_{\pi_{\phi}}[r_n   s_n, a_n]$
Policy update	$\phi \leftarrow \phi + \alpha \nabla_{\phi} \log \pi_{\phi}(a_n^*   s_n)$	$\phi \leftarrow \phi + \alpha \mathbf{G}_n \nabla_{\phi} \log \pi_{\phi}(a_n   s_n)$ where $G_n = \sum_{t=n}^{\infty} \gamma^t r_t$

# REINFORCE Algorithm

## REINFORCE( $s_0$ )

Initialize  $\pi_\phi$  to anything

Loop forever (for each episode)

Generate episode  $s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T$  with  $\pi_\phi$

Loop for each step of the episode  $n = 0, 1, \dots, T$

$$G_n \leftarrow \sum_{t=n}^T \gamma^t r_t$$

Update policy:  $\phi \leftarrow \phi + \alpha G_n \nabla_\theta \log \pi_\phi(a_n | s_n)$

Return  $\pi_\phi$

# Proximal Policy Optimization (PPO)

Initialize  $\pi_\phi$  and  $V_w$  to anything

Loop forever (for each episode)

Generate episode  $s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{N-1}, a_{N-1}, r_{N-1}$  with  $\pi_\phi$

Loop for each step of the episode  $n = 0, 1, \dots, N - 1$

$$G_n \leftarrow \sum_{t=n}^N \gamma^t r_t$$

$$A(s_n, a_n) \leftarrow G_n - V_w(s_n)$$

Update value function:  $w \leftarrow w + \alpha_w A(s_n, a_n) \nabla_w V_w(s_n)$

Update  $\pi$ :

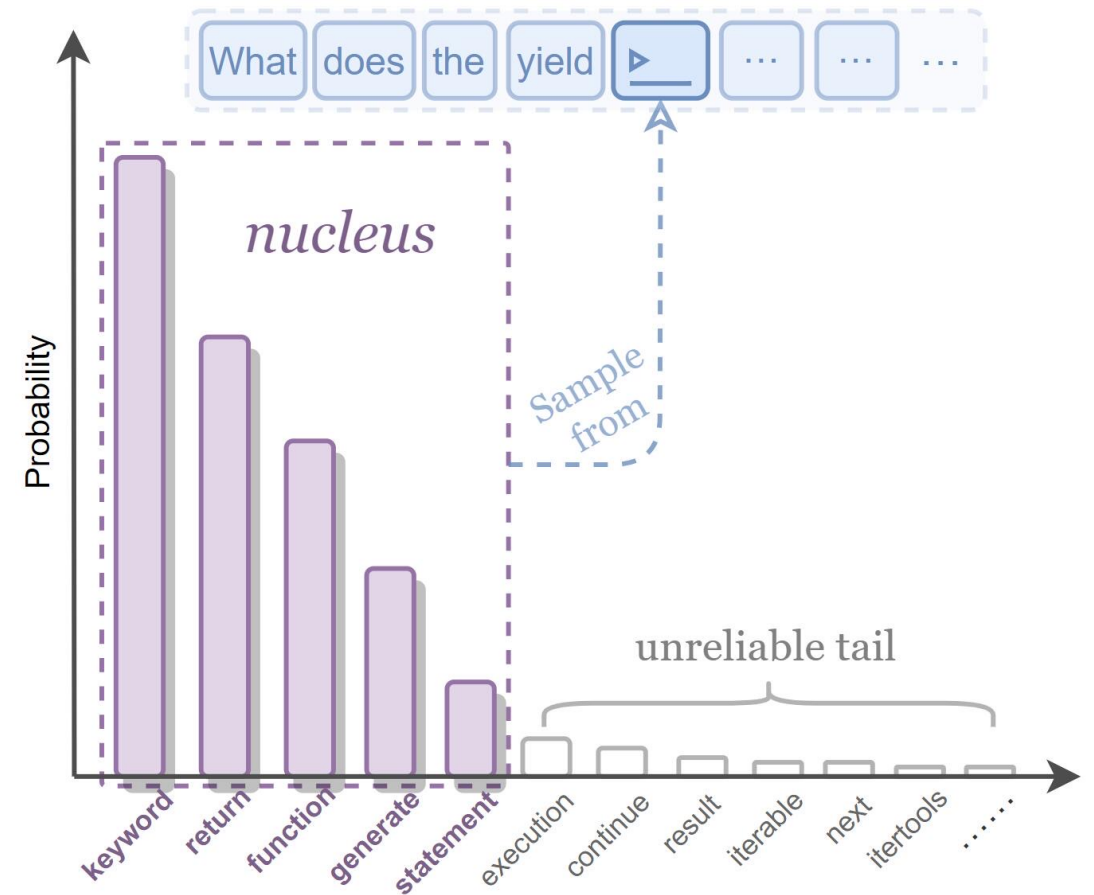
optimize by stochastic gradient descent

$$\phi \leftarrow \operatorname{argmax}_{\tilde{\phi}} \frac{1}{N} \sum_{n=0}^{N-1} \min \left\{ \begin{array}{l} \frac{\pi_{\tilde{\phi}}(a_n | s_n)}{\pi_{\phi}(a_n | s_n)} A(s_n, a_n), \\ \operatorname{clip} \left( \frac{\pi_{\tilde{\phi}}(a_n | s_n)}{\pi_{\phi}(a_n | s_n)}, 1 - \epsilon, 1 + \epsilon \right) A(s_n, a_n) \end{array} \right\}$$

# Inference: Nucleus sampling

Sample from nucleus (top tokens only) to avoid unreliable responses while ensuring diversity

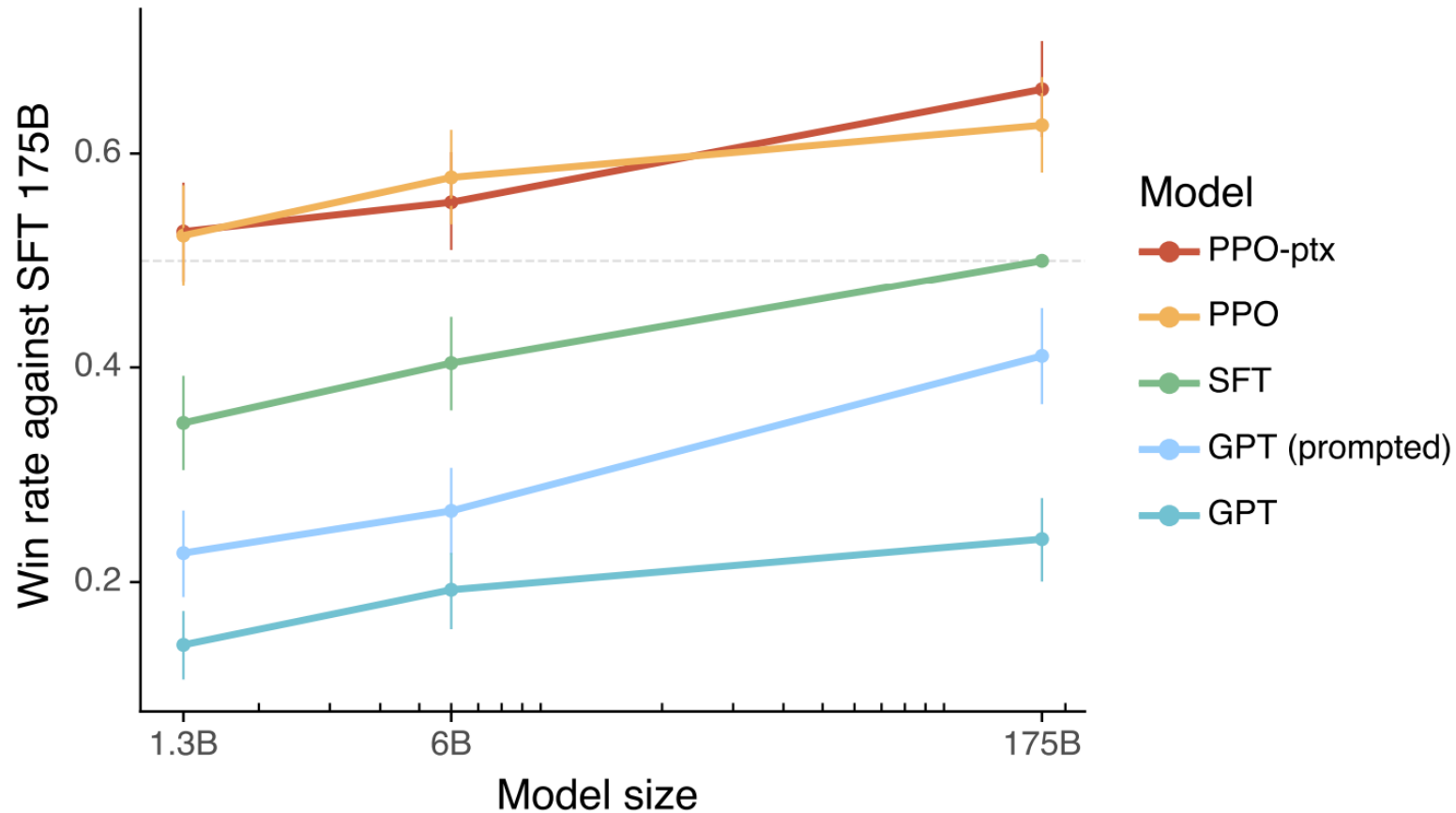
Holtzman, Ari; Buys, Jan; Du, Li; Forbes, Maxwell; Choi, Yejin (2019).  
**The Curious Case of Neural Text Degeneration**, arxiv.



Credit: <https://arxiv.labs.arxiv.org/html/2208.11523>

# InstructGPT Results

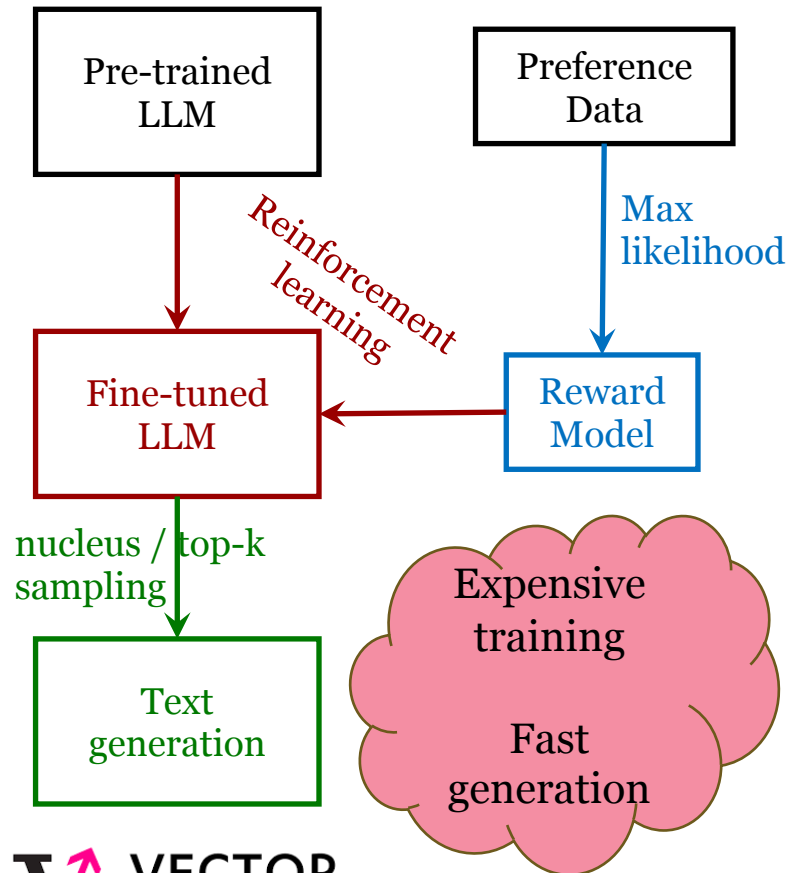
Ouyang, Wu, Jiang, Wainwright, et al. (2022)



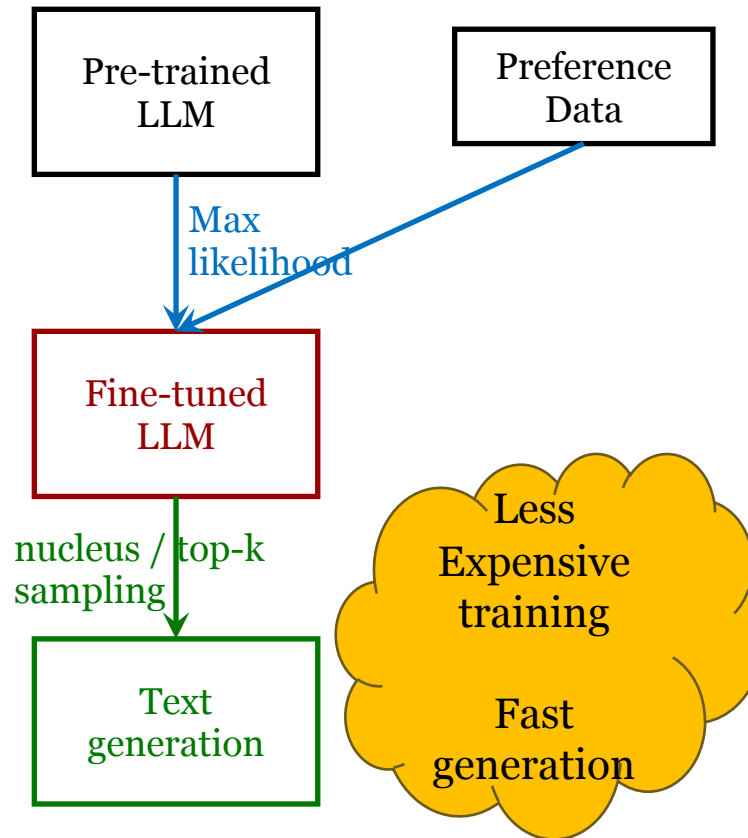


# RLHF Improvements

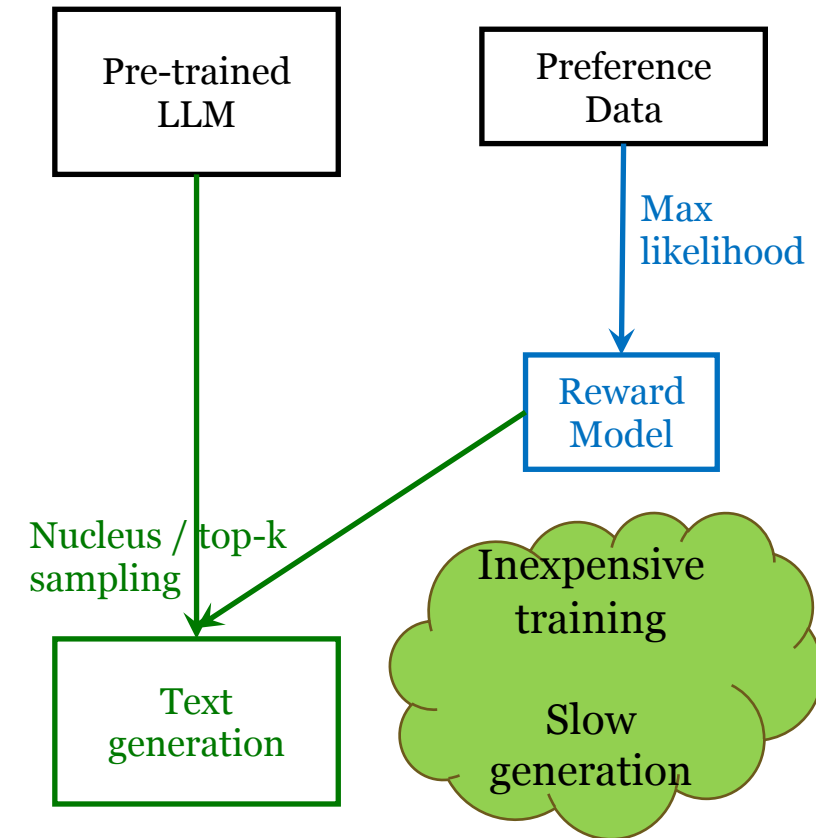
## Proximal Policy Optimization (PPO) Ouyang et al., 2022



## Direct Preference Optimization (DPO) Rafailov et al., 2023



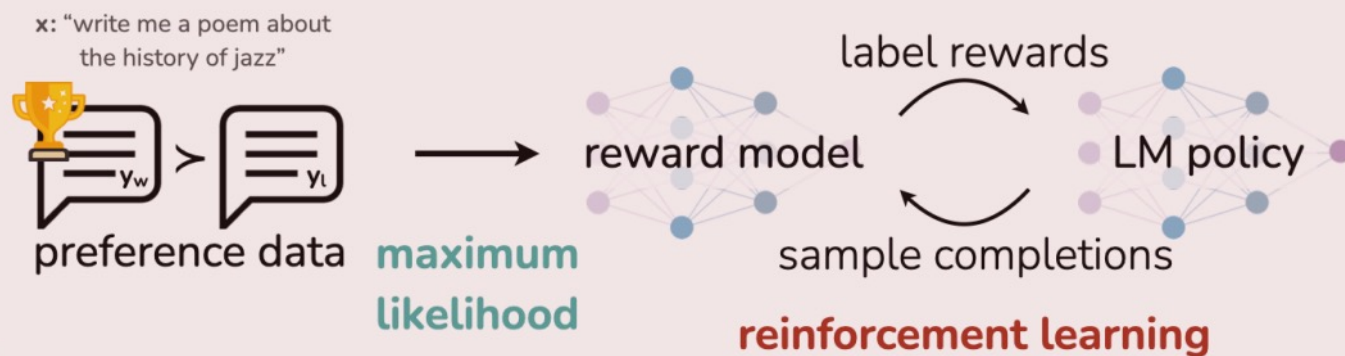
## Reward Guided Text Generation (RG TG) Khanov et al., 2024 Rashid et al., 2025



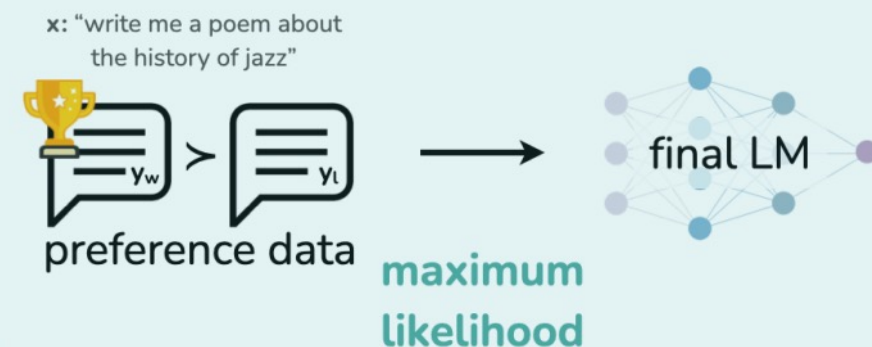
# Direct Preference Optimization

Rafailov, Sharma, Mitchell, Ermon, Manning, Finn (2023) **Direct Preference Optimization: Your Language Model is Secretly a Reward Model**, *NeurIPS*.

## Reinforcement Learning from Human Feedback (RLHF)



## Direct Preference Optimization (DPO)



# Bypassing RL

- Recall RL objective:

$$\max_{\phi} E_{s \in Dataset} \left[ E_{a \sim \pi_{\phi}(a|s)} [r_{\theta}(s, a)] - \beta KL(\pi_{\phi}(\cdot | s) | \pi_{ref}(\cdot | s)) \right]$$

- Closed form solution (based on maximum entropy RL):

$$\pi_{\phi}(a|s) = \frac{1}{Z(s)} \pi_{ref}(a|s) \exp\left(\frac{r_{\theta}(s, a)}{\beta}\right)$$

- Isolate reward:  $r_{\theta}(s, a) = \beta \log \frac{\pi_{\phi}(a|s)}{\pi_{ref}(a|s)} + \beta \log Z(s)$

- Plug into preference objective:

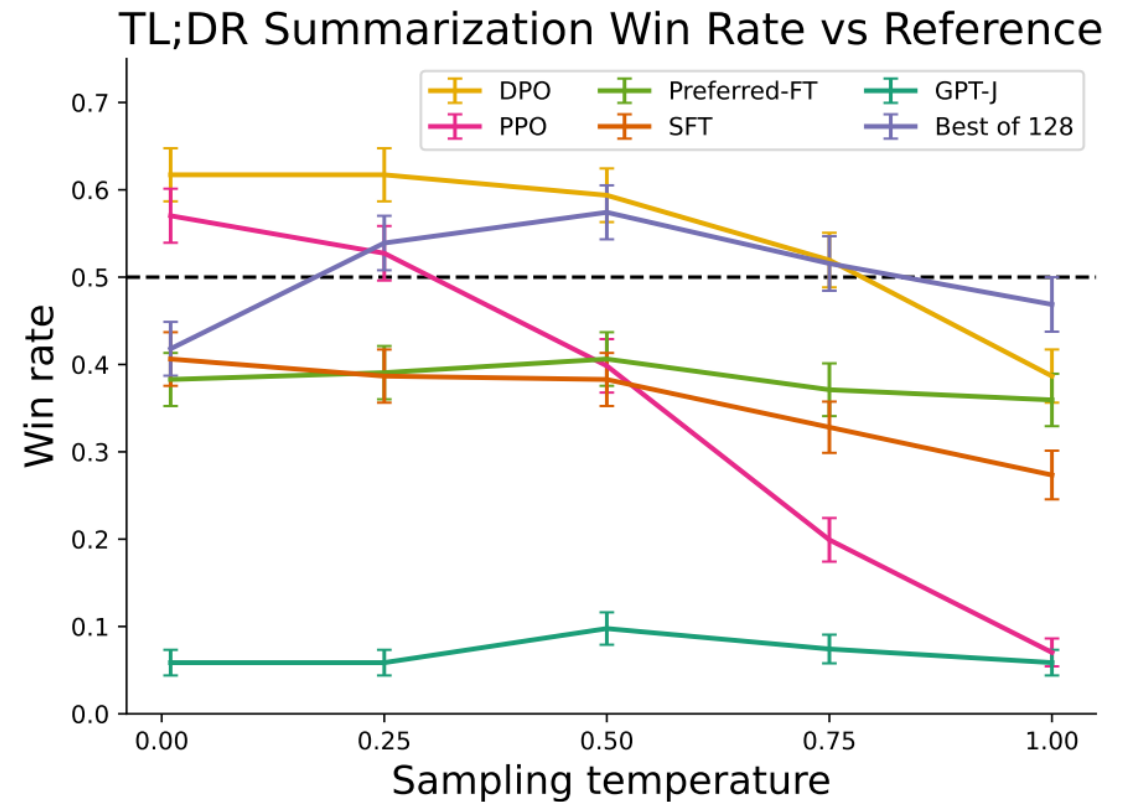
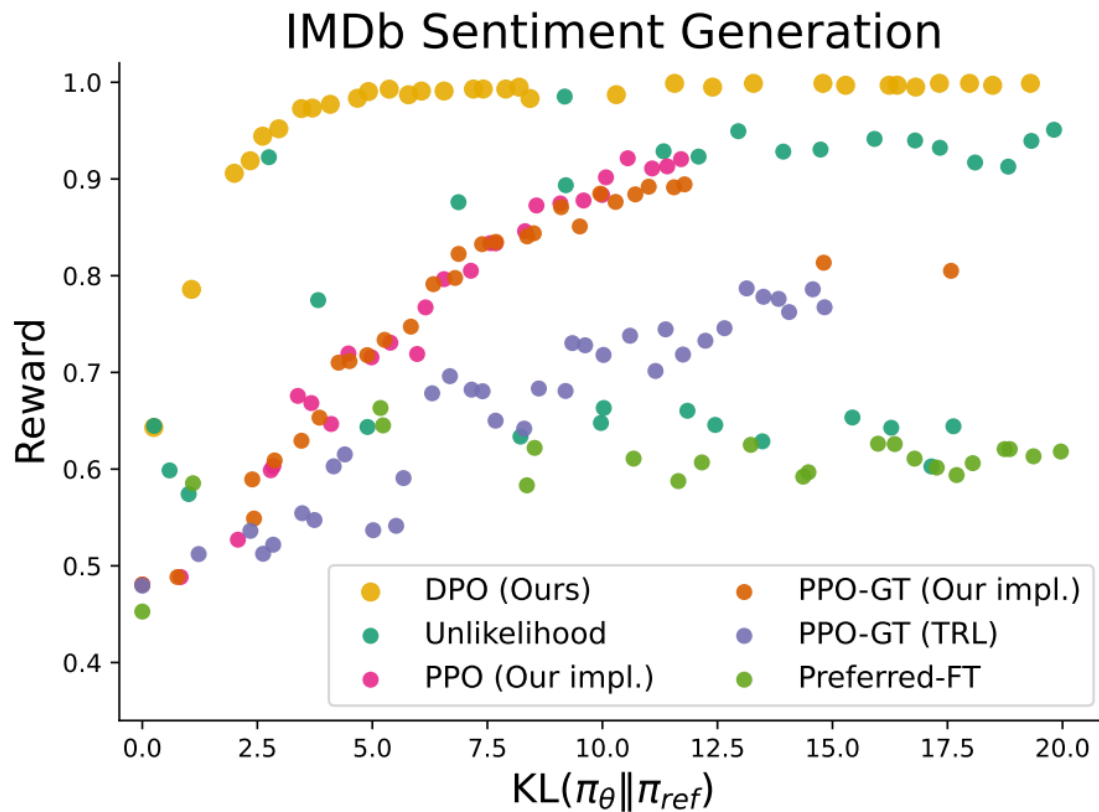
$$\begin{aligned} Loss(\theta) &= -\frac{1}{\binom{k}{2}} E_{(s, a_i, a_j) \in Dataset} \log \sigma(r_{\theta}(s, a_i) - r_{\theta}(s, a_j)) \\ &= -\frac{1}{\binom{k}{2}} E_{(s, a_i, a_j) \in Dataset} \log \sigma\left(\beta \log \frac{\pi_{\phi}(a_i|s)}{\pi_{ref}(a_i|s)} - \beta \log \frac{\pi_{\phi}(a_j|s)}{\pi_{ref}(a_j|s)}\right) \end{aligned}$$

# Optimal Policy Derivation

$$\begin{aligned}
 & \operatorname{argmax}_{\phi} E_{s \in \text{Dataset}} \left[ E_{a \sim \pi_{\phi}(a|s)} [r_{\theta}(s, a)] - \beta \operatorname{KL}(\pi_{\phi}(\cdot | s) \| \pi_{\text{ref}}(\cdot | s)) \right] \\
 &= \operatorname{argmax}_{\phi} E_{s \in \text{Dataset}} \left[ E_{a \sim \pi_{\phi}(a|s)} \left[ r_{\theta}(s, a) - \beta \log \frac{\pi_{\phi}(a|s)}{\pi_{\text{ref}}(a|s)} \right] \right] && \text{by KL definition} \\
 &= \operatorname{argmin}_{\phi} E_{s \in \text{Dataset}} \left[ E_{a \sim \pi_{\phi}(a|s)} \left[ \log \frac{\pi_{\phi}(a|s)}{\pi_{\text{ref}}(a|s)} - \frac{1}{\beta} r_{\theta}(s, a) \right] \right] && \text{since max = - min} \\
 &= \operatorname{argmin}_{\phi} E_{s \in \text{Dataset}} \left[ E_{a \sim \pi_{\phi}(a|s)} \left[ \log \frac{\pi_{\phi}(a|s)}{\frac{1}{Z(s)} \pi_{\text{ref}}(a|s) \exp\left(\frac{r_{\theta}(s, a)}{\beta}\right)} - \log Z(s) \right] \right] && \text{where } Z(s) = \sum_a \pi_{\text{ref}}(a|s) \exp\left(\frac{r_{\theta}(s, a)}{\beta}\right) \\
 &= \operatorname{argmin}_{\phi} E_{s \in \text{Dataset}} \left[ E_{a \sim \pi_{\phi}(a|s)} \left[ \log \frac{\pi_{\phi}(a|s)}{\frac{1}{Z(s)} \pi_{\text{ref}}(a|s) \exp\left(\frac{r_{\theta}(s, a)}{\beta}\right)} \right] \right] && \text{since } \log Z(s) \text{ is independent of } \phi \\
 &= \operatorname{argmin}_{\phi} E_{s \in \text{Dataset}} \left[ E_{a \sim \pi_{\phi}(a|s)} \left[ \log \frac{\pi_{\phi}(a|s)}{\pi_{\phi^*}(a|s)} \right] \right] && \text{where } \pi_{\phi^*}(a|s) = \frac{1}{Z(s)} \pi_{\text{ref}}(a|s) \exp\left(\frac{r_{\theta}(s, a)}{\beta}\right) \\
 &= \operatorname{argmin}_{\phi} E_{s \in \text{Dataset}} \left[ \operatorname{KL}(\pi_{\phi}(\cdot | s) \| \pi_{\phi^*}(\cdot | s)) \right] && \text{by KL definition} \\
 &= \phi^* && \text{since KL is minimized when both arguments are equal}
 \end{aligned}$$

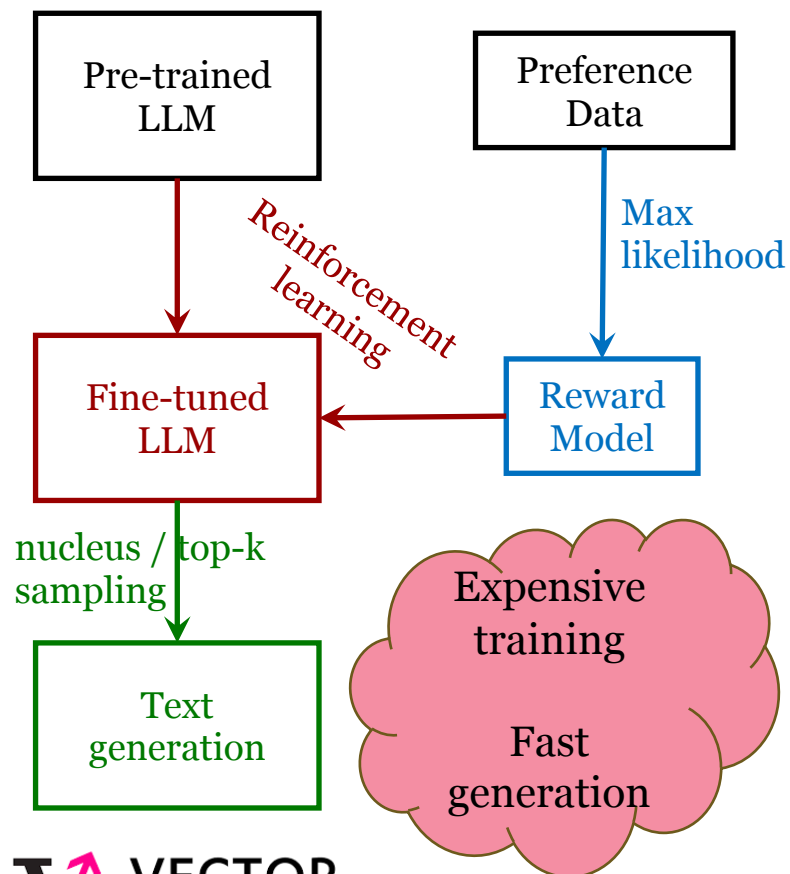
# Empirical Results

Rafailov et al. 2023

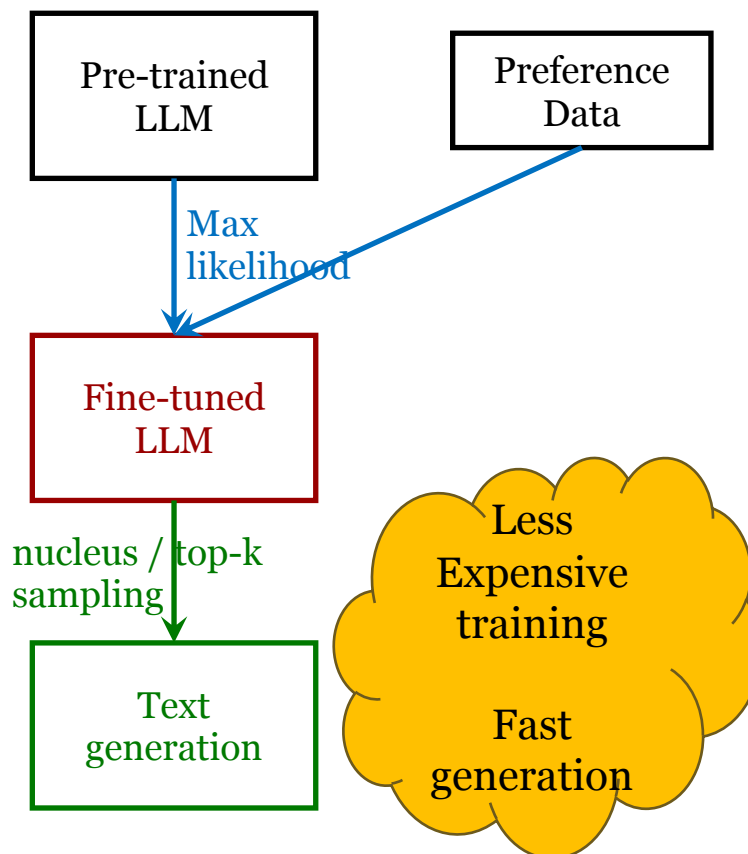


# RLHF Improvements

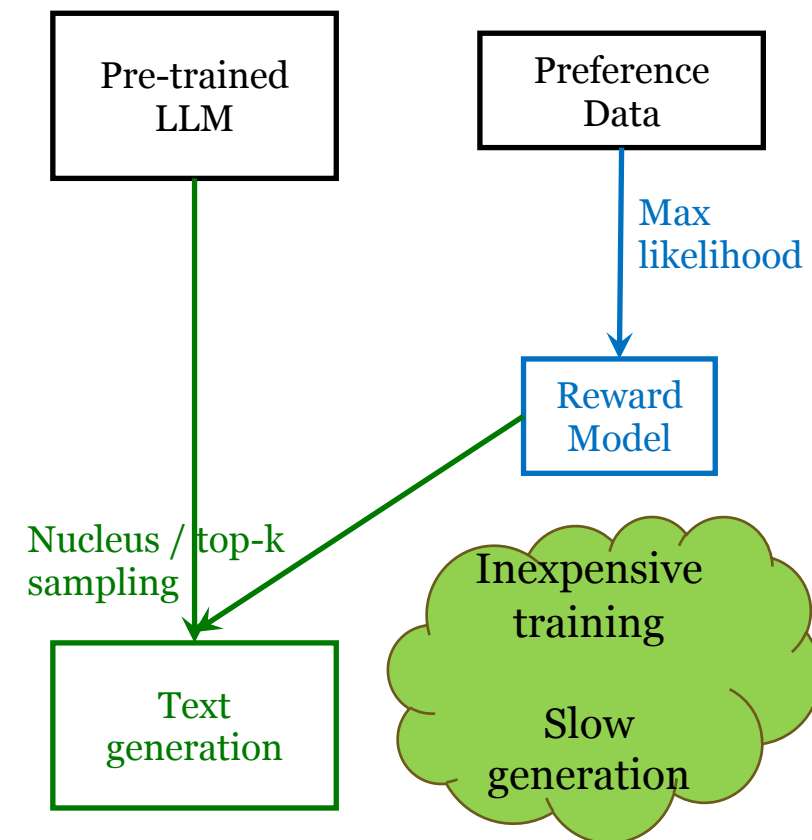
## Proximal Policy Optimization (PPO) Ouyang et al., 2022



## Direct Preference Optimization (DPO) Rafailov et al., 2023



## Reward Guided Text Generation (RGTG) Khanov et al., 2024 Rashid et al., 2025



# Sequence Generation

- Recall closed form solution

$$\begin{aligned}\pi_{\phi}(\mathbf{a}|\mathbf{s}) &= \frac{1}{Z(\mathbf{s})} \pi_{ref}(\mathbf{a}|\mathbf{s}) \exp\left(\frac{r_{\theta}(\mathbf{s}, \mathbf{a})}{\beta}\right) \\ &= softmax\left(\log \pi_{ref}(\mathbf{a}|\mathbf{s}) + \frac{r_{\theta}(\mathbf{s}, \mathbf{a})}{\beta}\right)\end{aligned}$$

- Text generation:

$$\mathbf{a} \sim softmax\left(\log \begin{pmatrix} \pi_{ref}(\mathbf{a}_1|\mathbf{s}) \\ \pi_{ref}(\mathbf{a}_2|\mathbf{s}) \\ \pi_{ref}(\mathbf{a}_3|\mathbf{s}) \\ \dots \\ \pi_{ref}(\mathbf{a}_n|\mathbf{s}) \end{pmatrix} + \begin{pmatrix} r_{\theta}(\mathbf{s}, \mathbf{a}_1) \\ r_{\theta}(\mathbf{s}, \mathbf{a}_2) \\ r_{\theta}(\mathbf{s}, \mathbf{a}_3) \\ \dots \\ r_{\theta}(\mathbf{s}, \mathbf{a}_n) \end{pmatrix} / \beta\right)$$

# Token Generation

- Token-wise LLM modeling

$$\begin{aligned}\pi_{\phi}(\mathbf{a}^i | \mathbf{s}, \mathbf{a}^{1:i-1}) &= \frac{1}{Z(\mathbf{s})} \pi_{ref}(\mathbf{a}^i | \mathbf{s}, \mathbf{a}^{1:i-1}) \exp\left(\frac{r_{\theta}(\mathbf{s}, \mathbf{a}^{1:i})}{\beta}\right) \\ &= softmax\left(\log \pi_{ref}(\mathbf{a}^i | \mathbf{s}, \mathbf{a}^{1:i-1}) + \frac{r_{\theta}(\mathbf{s}, \mathbf{a}^{1:i})}{\beta}\right)\end{aligned}$$

- Token generation:

$$a^i \sim softmax\left(\log \begin{pmatrix} \pi_{ref}(a_1^i | \mathbf{s}, \mathbf{a}^{1:i-1}) \\ \pi_{ref}(a_2^i | \mathbf{s}, \mathbf{a}^{1:i-1}) \\ \pi_{ref}(a_3^i | \mathbf{s}, \mathbf{a}^{1:i-1}) \\ \vdots \\ \pi_{ref}(a_n^i | \mathbf{s}, \mathbf{a}^{1:i-1}) \end{pmatrix} + \begin{pmatrix} r_{\theta}(\mathbf{s}, \mathbf{a}^{1:i-1}, a_1^i) \\ r_{\theta}(\mathbf{s}, \mathbf{a}^{1:i-1}, a_2^i) \\ r_{\theta}(\mathbf{s}, \mathbf{a}^{1:i-1}, a_3^i) \\ \vdots \\ r_{\theta}(\mathbf{s}, \mathbf{a}^{1:i-1}, a_n^i) \end{pmatrix} / \beta\right)$$



# FaRMA: Faster Reward Model for Alignment

- Rashid, Wu, Fan, Li, Kristiadi, Poupart (2025) **Towards Cost-Effective Reward Guided Text Generation**, *ICML*.
- Optimization problem:

$$\begin{aligned} & \max_{\theta} E_{(s, a_+, a_-) \in \text{Dataset}} \log \sigma(r_{\theta}(s, a_+) - r_{\theta}(s, a_-)) \\ & \text{Subject to } r_{\theta}(s, a^{1:i}) = \max_{a^{i+1:|a|}} r_{\theta}(s, [a^{1:i}, a^{i+1:|a|}]) \quad \forall s, a, i \end{aligned}$$

- In practice: alternate between minimizing two loss functions
  - $L_1(\theta) = -E_{(s, a_+, a_-) \in \text{Dataset}} \log \sigma(r_{\theta}(s, a_+) - r_{\theta}(s, a_-))$
  - $L_2(\theta) = \frac{1}{2} E_{(s, a) \in \text{Dataset}, i \leq |a|} \left( r_{\theta}(s, a^{1:i}) - \max_{a^{i+1:|a|}} r_{\theta}(s, [a^{1:i}, a^{i+1:|a|}]) \right)^2$

# FaRMA Pseudocode

Repeat

Repeat for each  $(s, \mathbf{a}_+, \mathbf{a}_-)$  in minibatch

$$L_1(\theta) = \log \sigma(r_\theta(s, \mathbf{a}_+) - r_\theta(s, \mathbf{a}_-))$$

$$\theta \leftarrow \theta - \alpha \nabla L_1(\theta)$$

Repeat for each  $(s, \mathbf{a}, i)$  in minibatch

$$L_2(\theta) = \frac{1}{2} \left( r_\theta(s, \mathbf{a}^{1:i}) - \max_{a^{i+1}} r_\theta(s, \mathbf{a}^{1:i+1}) \right)^2$$

$$\theta \leftarrow \theta - \alpha \nabla L_2(\theta)$$

# Empirical Results

TL;DR Summarization			
Method	LLM	$r \pm \text{SE}$	Time(min)
$\pi_{\text{ref}}$	frozen	$0.98 \pm 0.18$	2
ARGS	frozen	$1.46 \pm 0.16$	32
PARGS	frozen	$1.56 \pm 0.19$	31
CD	frozen	$1.15 \pm 0.16$	29
FaRMA	frozen	$2.05 \pm 0.15$	5
CARDS	frozen	$1.73 \pm 0.16$	17
DPO	trained	$2.08 \pm 0.18$	2
PPO	trained	$2.05 \pm 0.14$	2

Table 2. Avg. reward (over 100 samples)  $\pm$  standard error total generation time for the TL;DR summarization task.

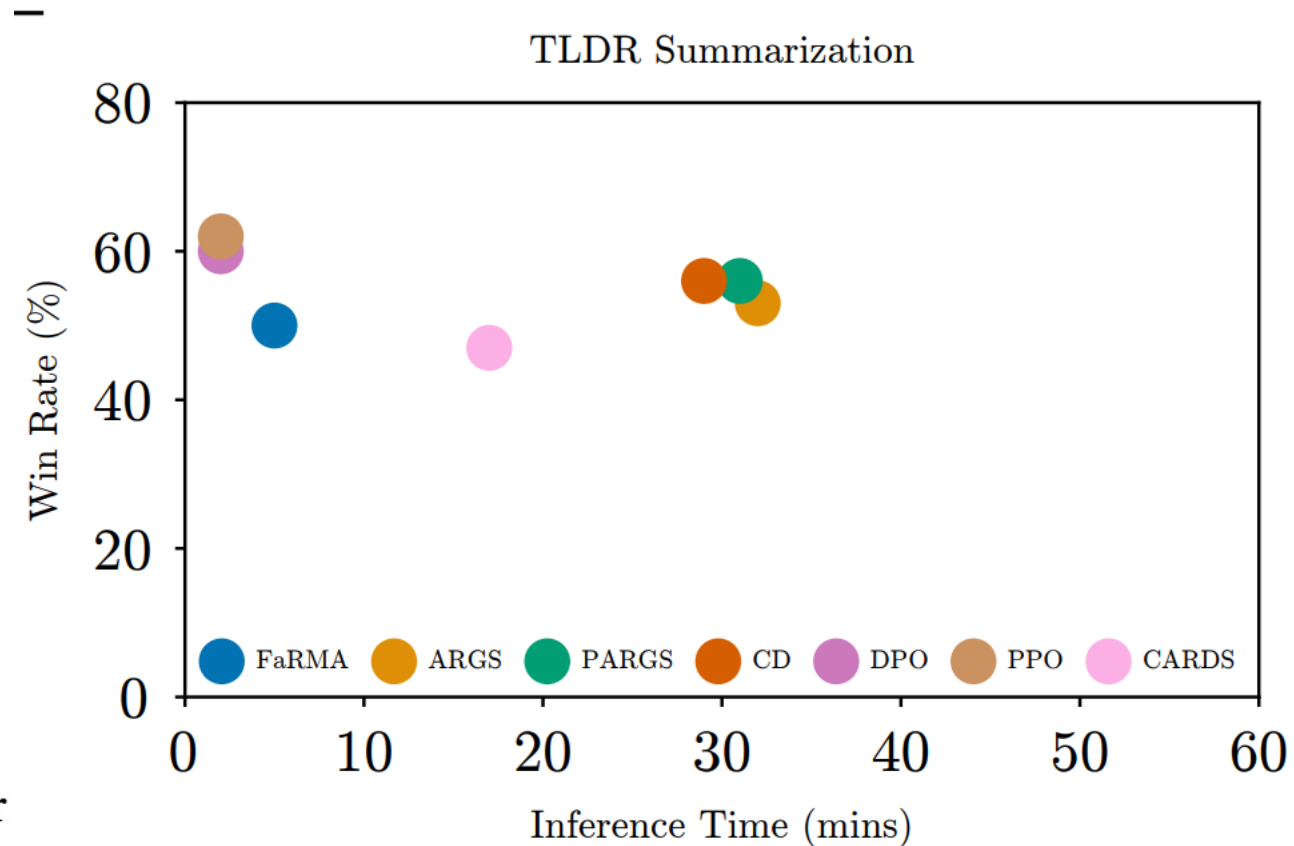


Figure 2. GPT4 evaluation on TLDR

# Outline

- LLM Alignment
  - Reinforcement Learning from Human Feedback
  - Direct Preference Optimization
  - Reward Guided Text Generation
- LLM Reasoning
  - Search and planning
  - Group Relative Policy Optimization (GRPO)
  - Reflection: Verbalized RL

# Reasoning LLMs

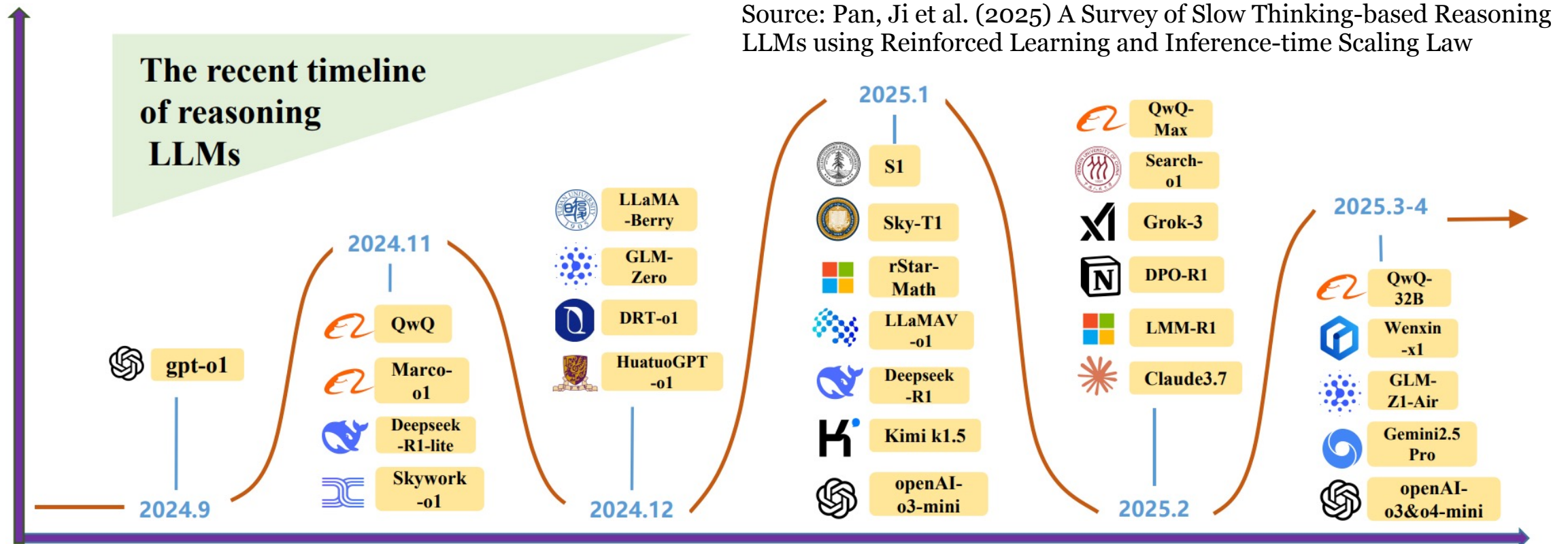
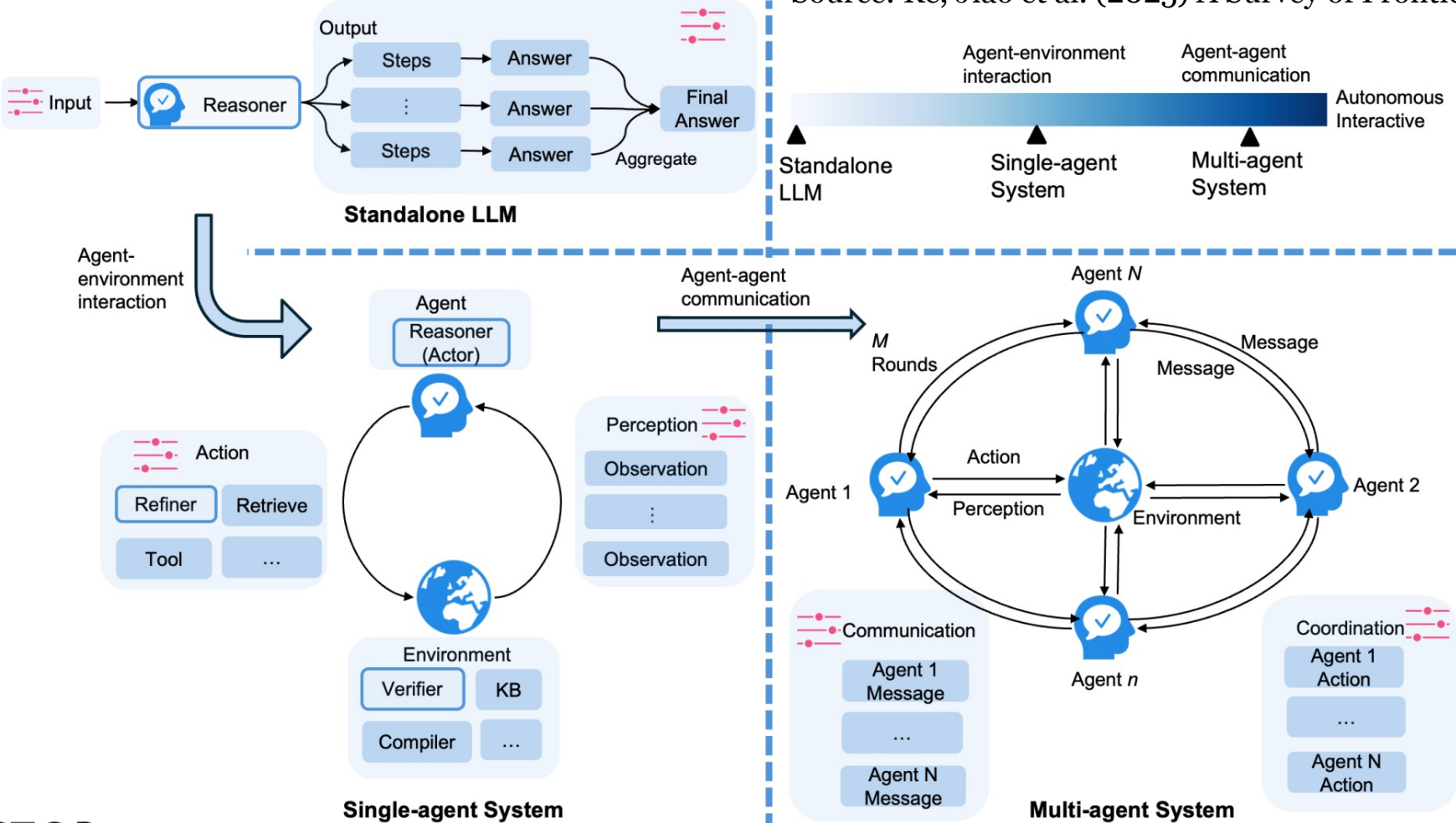


Fig. 1. The timeline of main reasoning LLMs.

# Inference Time Reasoning

Source: Ke, Jiao et al. (2025) A Survey of Frontiers in LLM Reasoning



# Reasoning by Searching

Source: Pan, Ji et al. (2025) A Survey of Slow Thinking-based Reasoning LLMs using Reinforced Learning and Inference-time Scaling Law

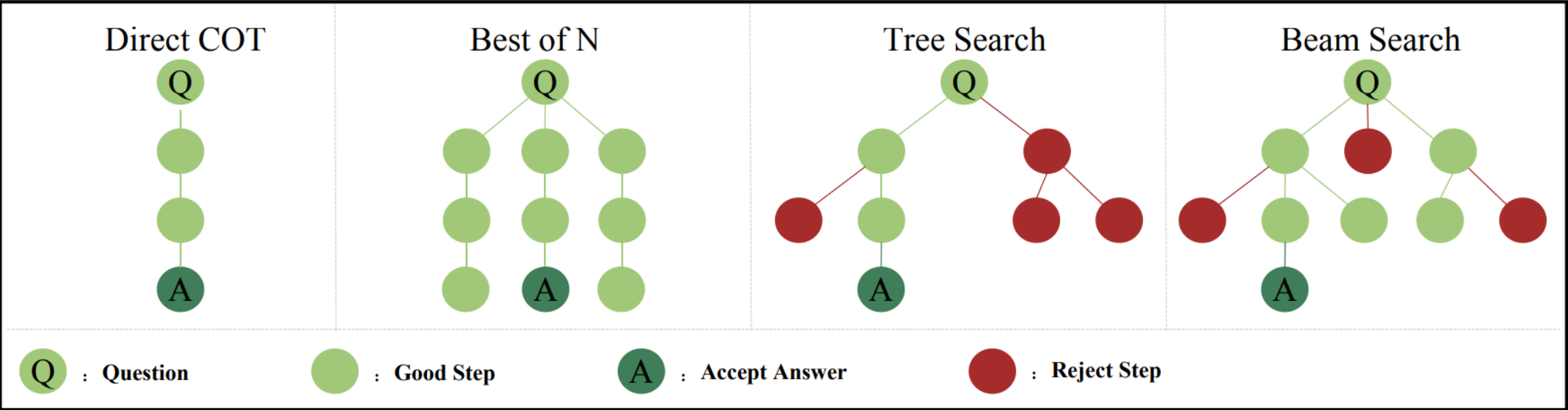
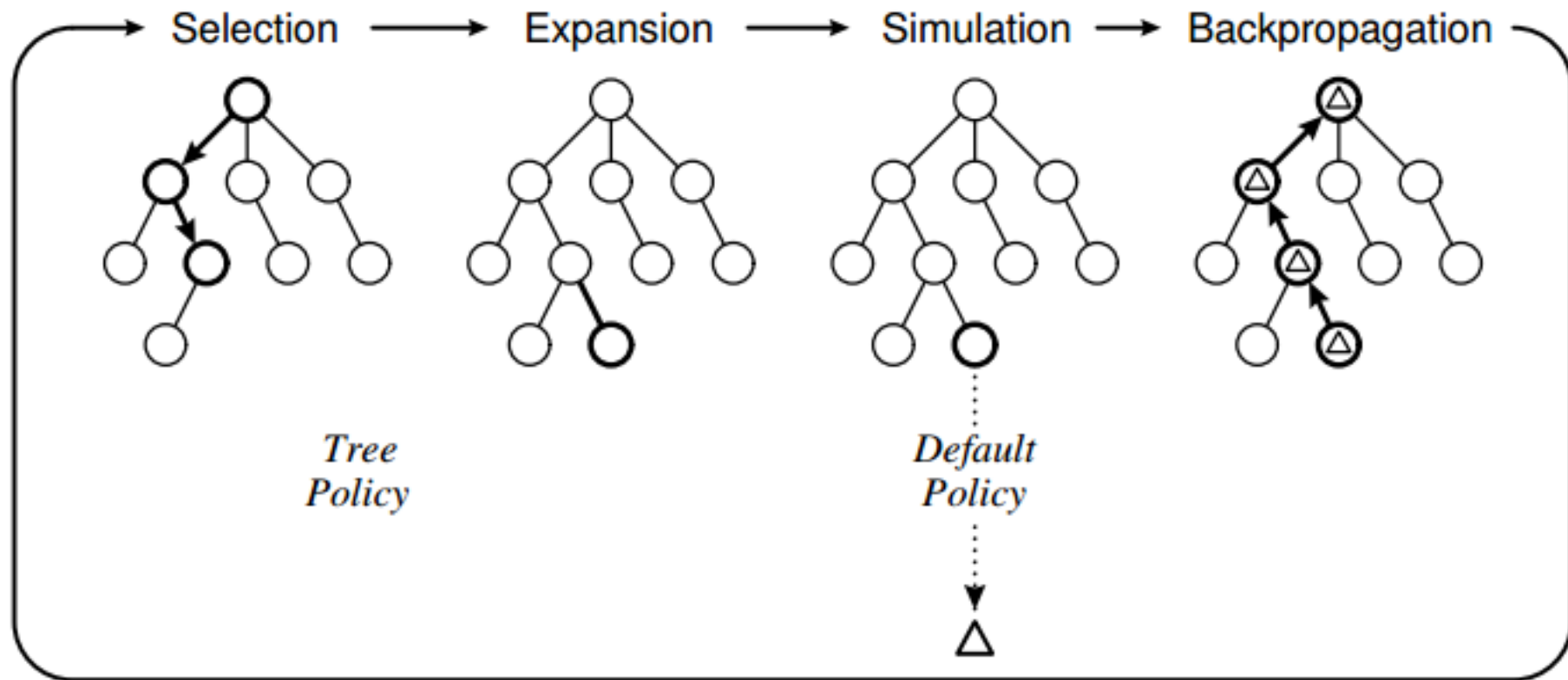


Fig. 3. The search algorithms for test-time scaling

# Monte Carlo Tree Search





# Learning to Reason

Source: Pan, Ji et al. (2025) A Survey of Slow Thinking-based Reasoning LLMs using Reinforced Learning and Inference-time Scaling Law

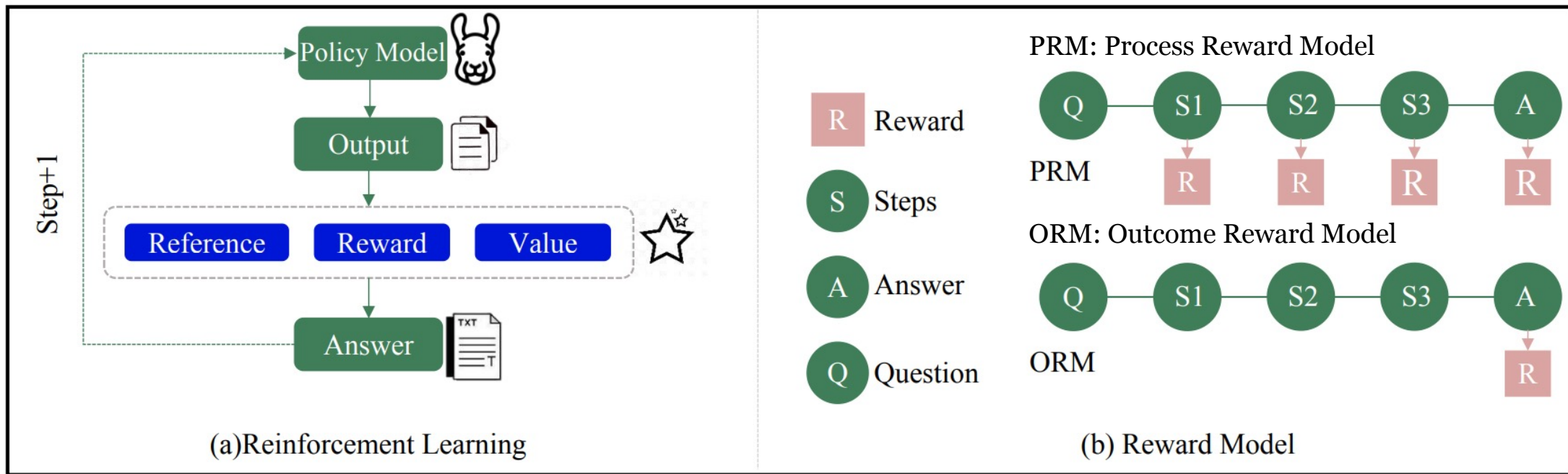


Fig. 4. The reinforcement learning framework and reward model

# Simplifying PPO

Source: Shao, Wang et al. (2024) DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models

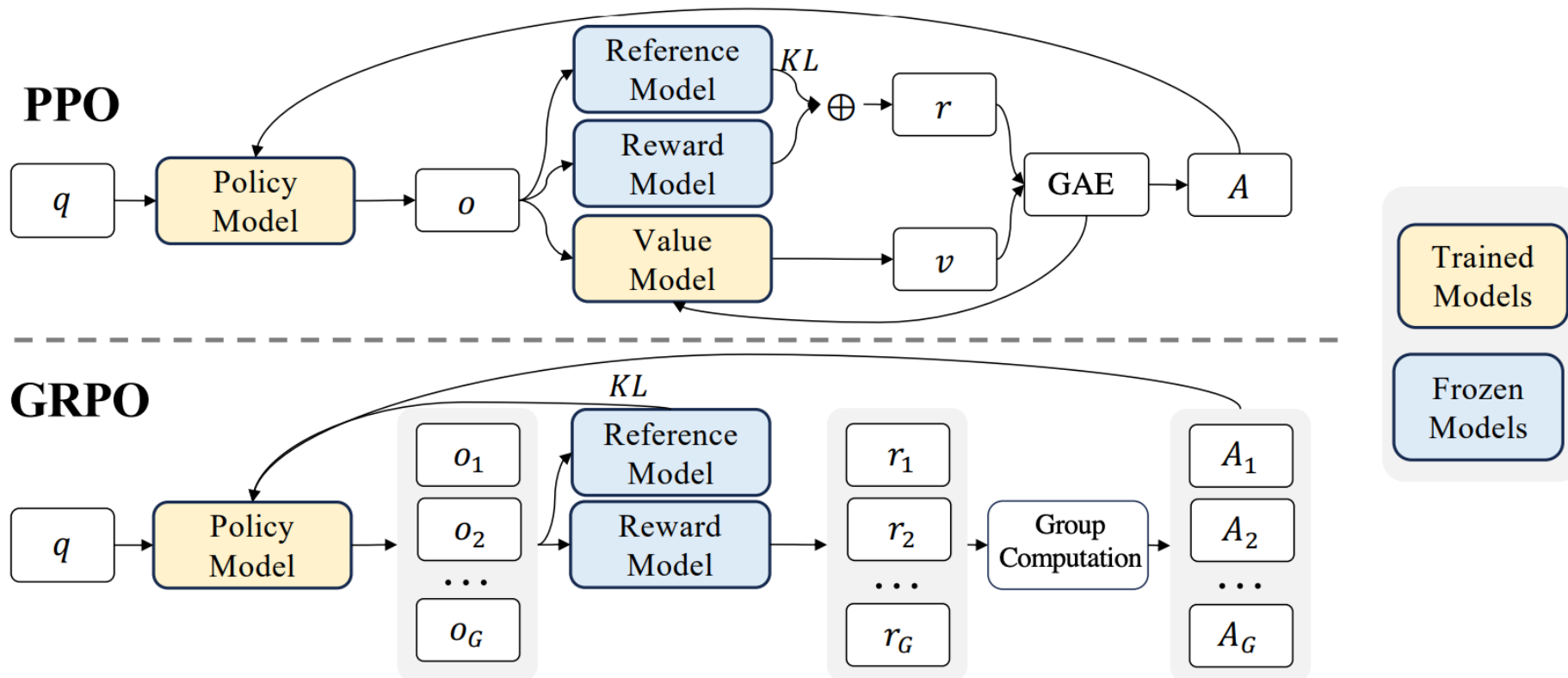


Figure 4 | Demonstration of PPO and our GRPO. GRPO foregoes the value model, instead estimating the baseline from group scores, significantly reducing training resources.

# Group Relative Policy Optimization (GRPO)

Initialize  $\pi_\phi$  and  $V_w$  to anything

Loop forever

Generate set of episodes  $\{\tau_0, \dots, \tau_{G-1}\}$ :

Sample  $\tau_g = (s_0^g, a_0^g, r_0^g, s_1^g, a_1^g, r_1^g, \dots, s_{N-1}^g, a_{N-1}^g, r_{N-1}^g)$  with  $\pi_\phi$

Evaluate:  $R_n^g \leftarrow \sum_{t=n}^N \gamma^t r(s_t^g, a_t^g) \forall n$

Loop for each episode  $g$  and step  $n$

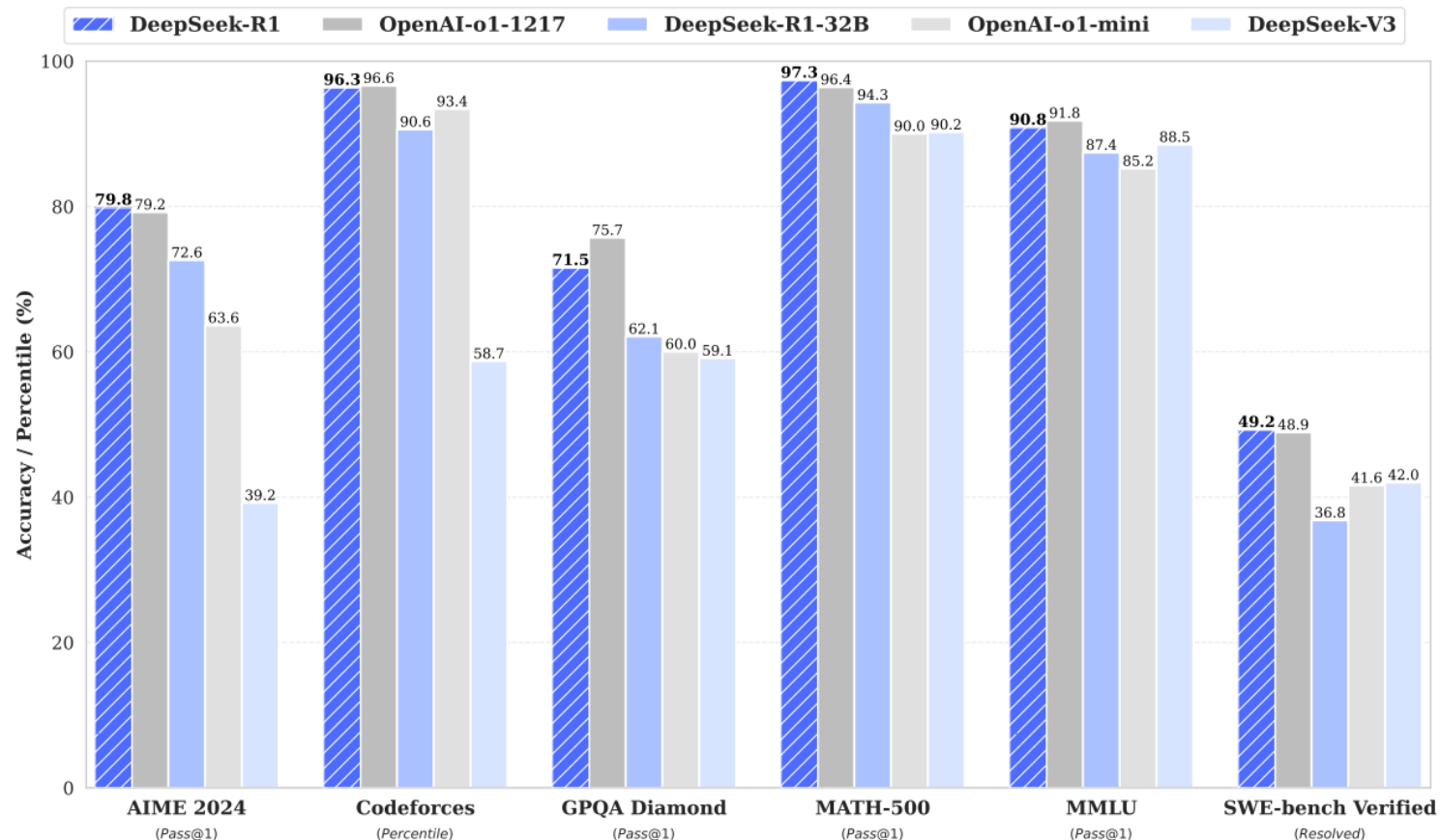
$A_n^g \leftarrow (R_n^g - \text{mean}(\{R_n^0, \dots, R_n^{G-1}\}))/\text{std}(\{R_n^0, \dots, R_n^{G-1}\})$

Update  $\pi$ :

$$\phi \leftarrow \underset{\tilde{\phi}}{\operatorname{argmax}} \frac{1}{G} \sum_{g=0}^{G-1} \frac{1}{N} \sum_{n=0}^{N-1} \min \left\{ \begin{array}{l} \frac{\pi_{\tilde{\phi}}(a_n^g | s_n^g)}{\pi_\phi(a_n^g | s_n^g)} A_n^g \\ \operatorname{clip} \left( \frac{\pi_{\tilde{\phi}}(a_n^g | s_n^g)}{\pi_\phi(a_n^g | s_n^g)}, 1 - \epsilon, 1 + \epsilon \right) A_n^g \end{array} \right\}$$

# DeepSeek-R1

Source: DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning (2025)



# Reflexion: Verbalized Reinforcement Learning

Source: Shinn, Cassano et al. (2023) Reflexion: Language Agents with Verbal Reinforcement Learning

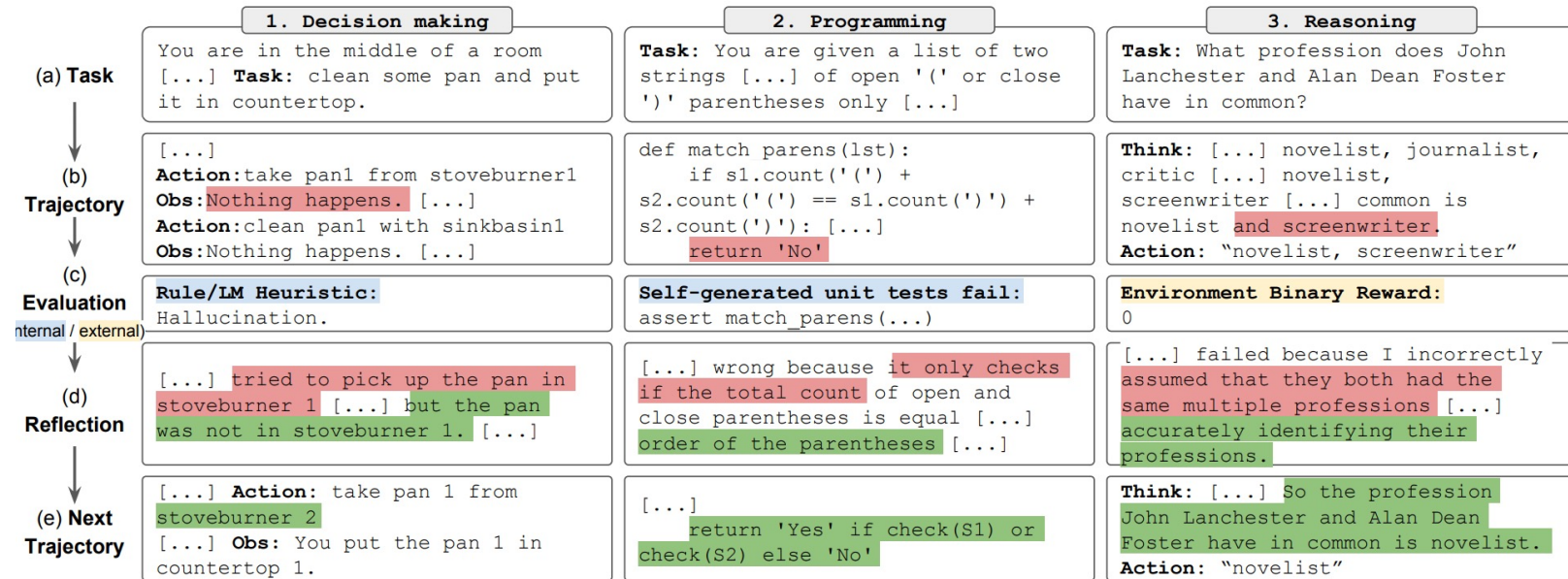
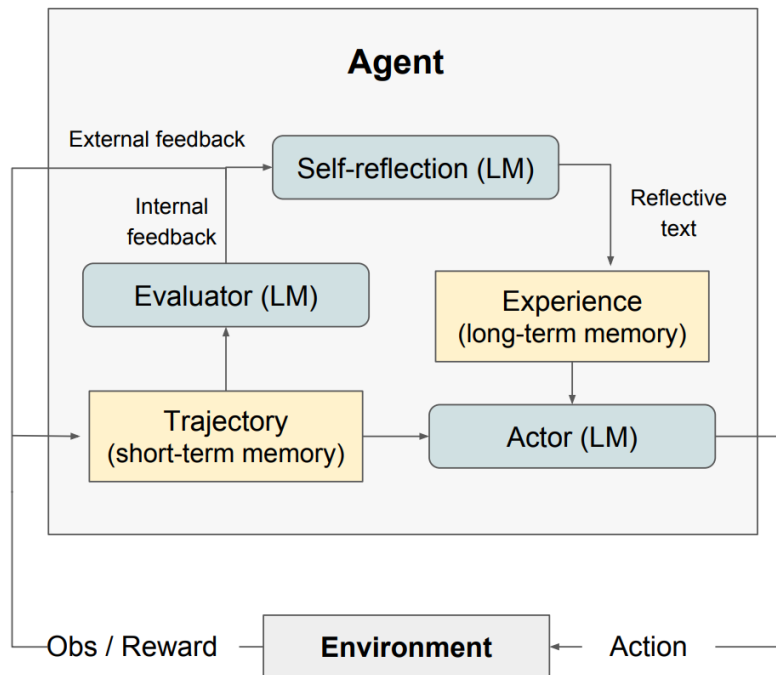


Figure 1: Reflexion works on decision-making 4.1, programming 4.3, and reasoning 4.2 tasks.

# Improved Reasoning by Self-Reflection

Source: Shinn, Cassano et al. (2023) Reflexion: Language Agents with Verbal Reinforcement Learning

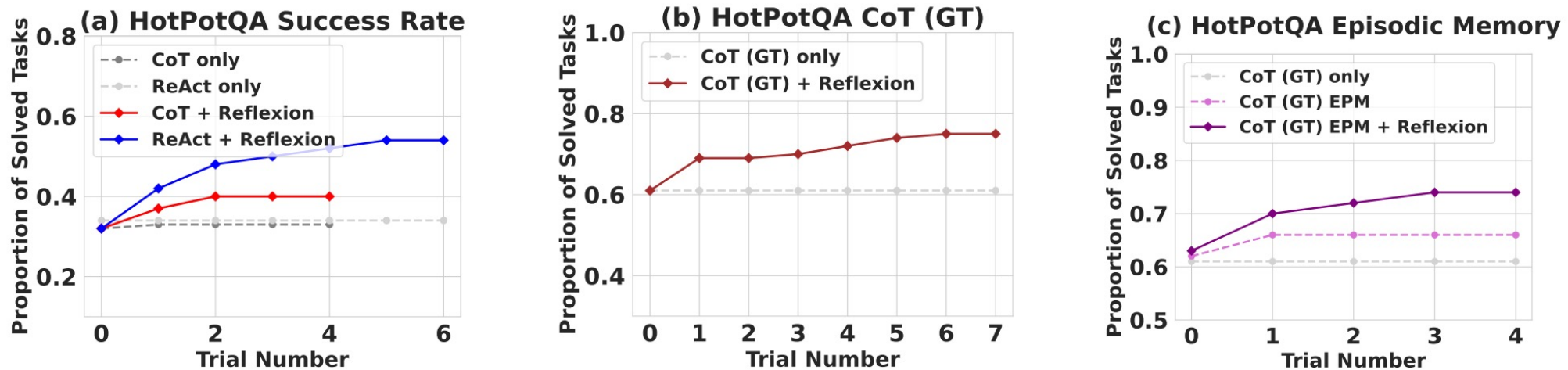


Figure 4: Chain-of-Thought (CoT) and ReAct. Reflexion improves search, information retrieval, and reasoning capabilities on 100 HotPotQA questions. (a) Reflexion ReAct vs Reflexion CoT (b) Reflexion CoT (GT) for reasoning only (c) Reflexion vs episodic memory ablation.



# Conclusion

- RL key to
  - LLM Alignment
  - LLM Reasoning
- Current Frontier:
  - Multi-agent RL for agentic orchestration

Source: Taghizadeh (2024) How Multi Agent LLMs Are Revolutionizing Reporting

