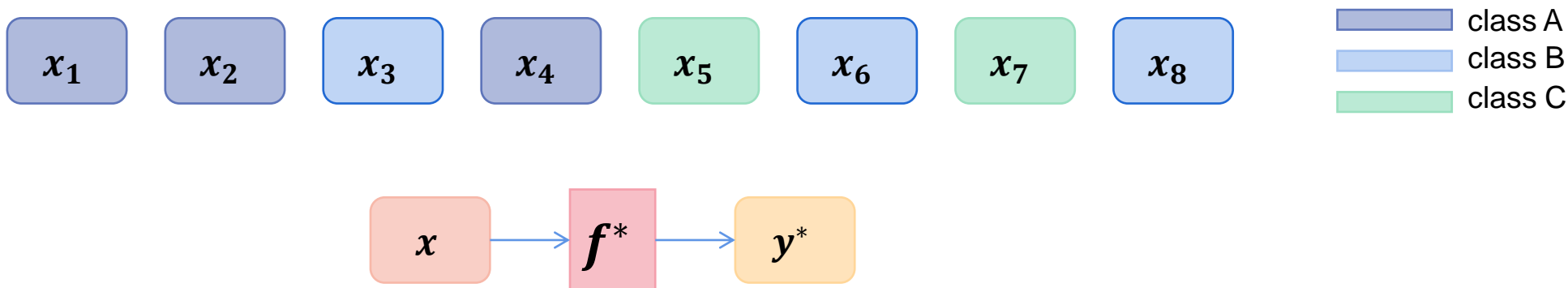
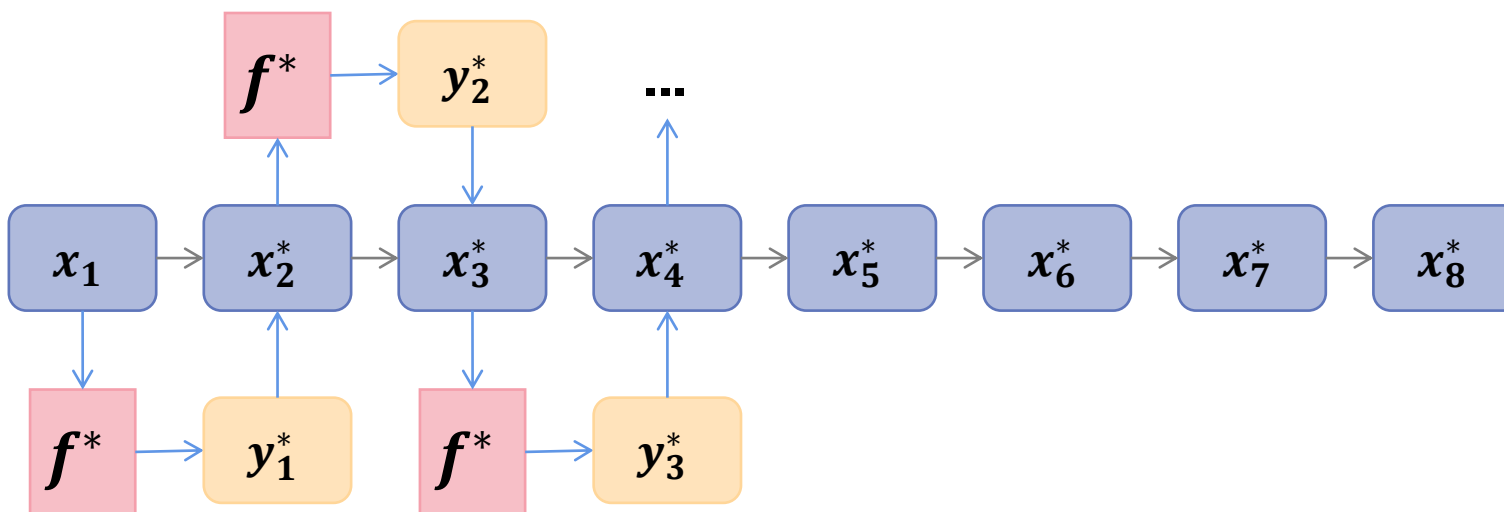


Bridging RL and Other ML Paradigms

Supervised Learning



Reinforcement Learning



sequence decision

➤ $\langle x_1, x_2^*, x_3^*, \dots \rangle$

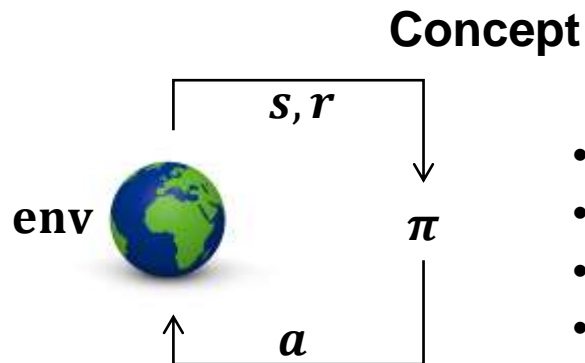
➤ $x_{t+1} = P(x_t, y_t)$

$x \rightarrow s$ state

$f \rightarrow \pi$ policy

$y \rightarrow a$ action

$R(s, a)$ reward



- $\text{env} \rightarrow s_0$
- $\pi \rightarrow a_0 \sim \pi(\cdot | s_0)$
- $\text{env} \rightarrow s_1, r_0$
- ...

MDP(Markov Decision Process)

Definition: (S, A, R, P, γ)

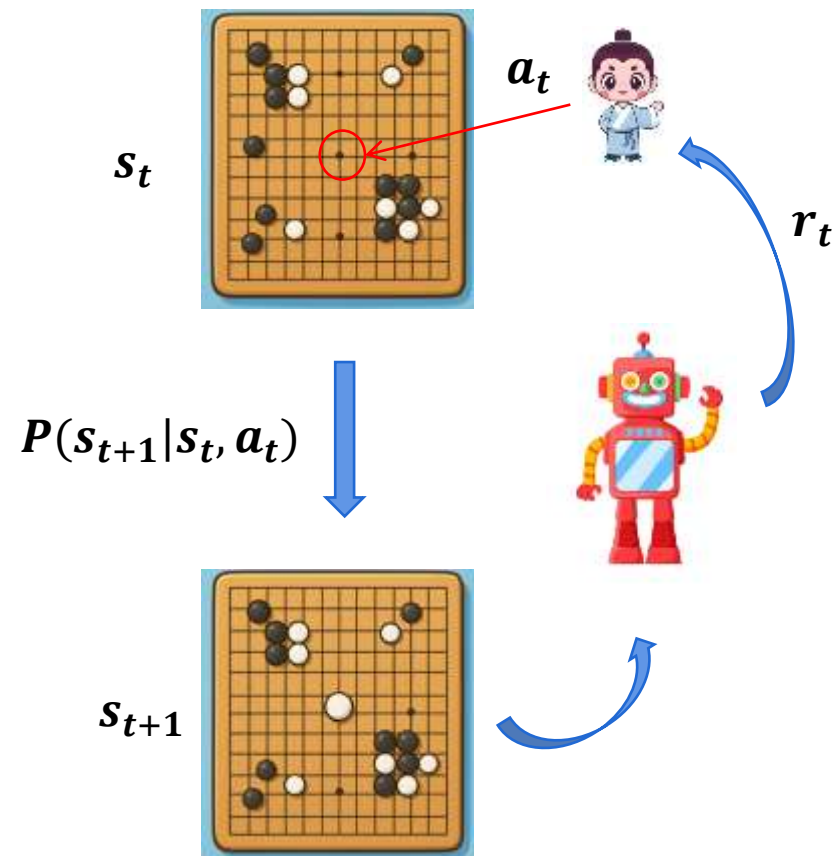
$s \in S \quad a \in A \quad r_t = R(s_t, a_t)$

Markov property: $P(s_{t+1} | s_{<t}, a_{<t}, s_t, a_t) = P(s_{t+1} | s_t, a_t)$

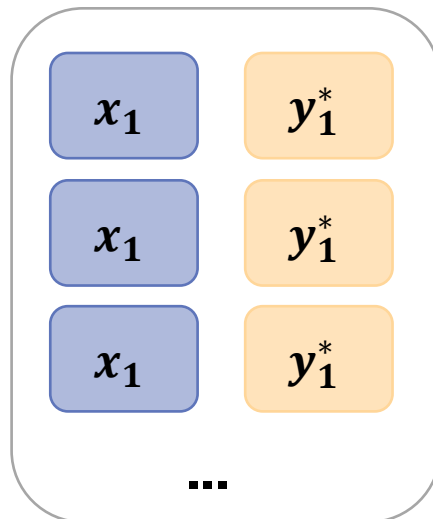
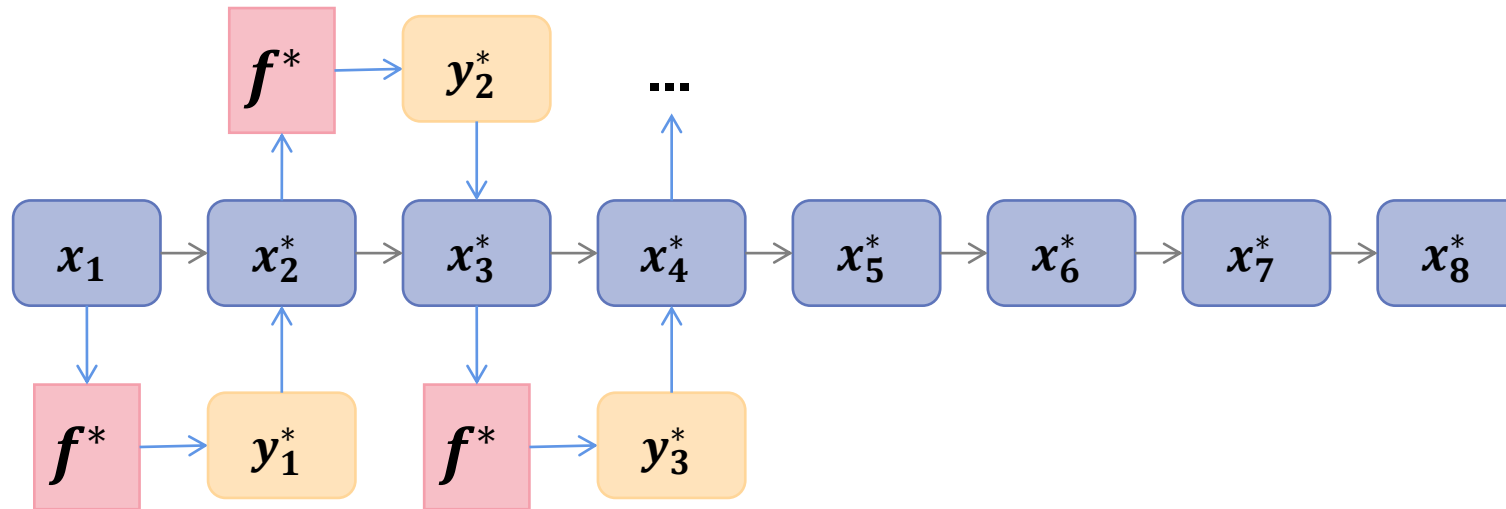
Policy objective:

$$J(\pi) = \mathbb{E}_{a_0 \sim \pi(\cdot | s_0)} \left[r(s_0, a_0) + \gamma \mathbb{E}_{s_1 \sim P(\cdot | s_0, a_0)} \left[\mathbb{E}_{a_1 \sim \pi(\cdot | s_1)} [r(s_1, a_1) + \dots] \right] \right]$$

Example

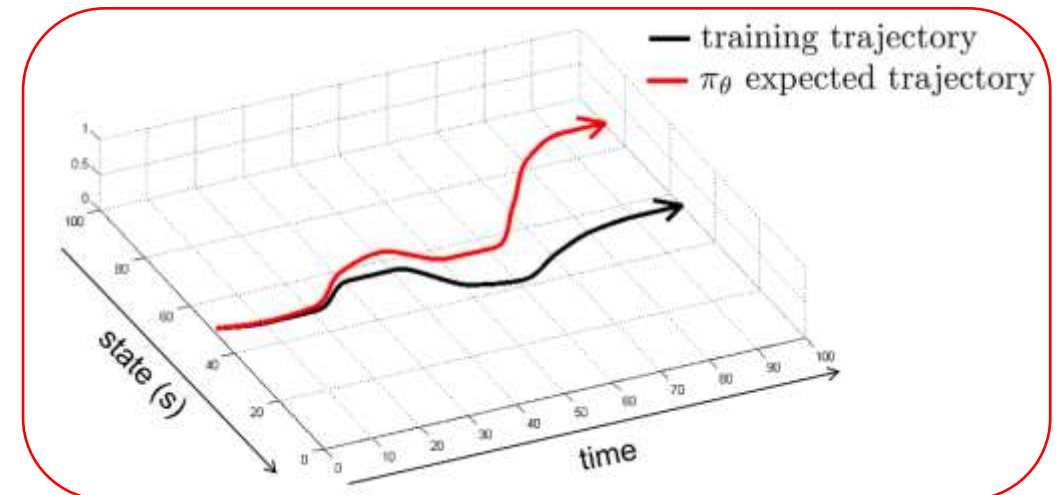


Reinforcement Learning



Imitation Learning

train a model



compounding errors

Reinforcement Learning-Guided Semi-Supervised Learning

Marzi Heidari Hanping Zhang Yuhong Guo
School of Computer Science, Carleton University, Ottawa, Canada

NeurIPS 2024

MDP for Semi-Supervised Learning

$$\text{policy: } \pi_{\theta}(\cdot) = P_{\theta}(\cdot) \left\{ \begin{array}{l} \text{state: } X^l \quad X^u \\ \text{action: } Y^l \quad Y^u \end{array} \right.$$

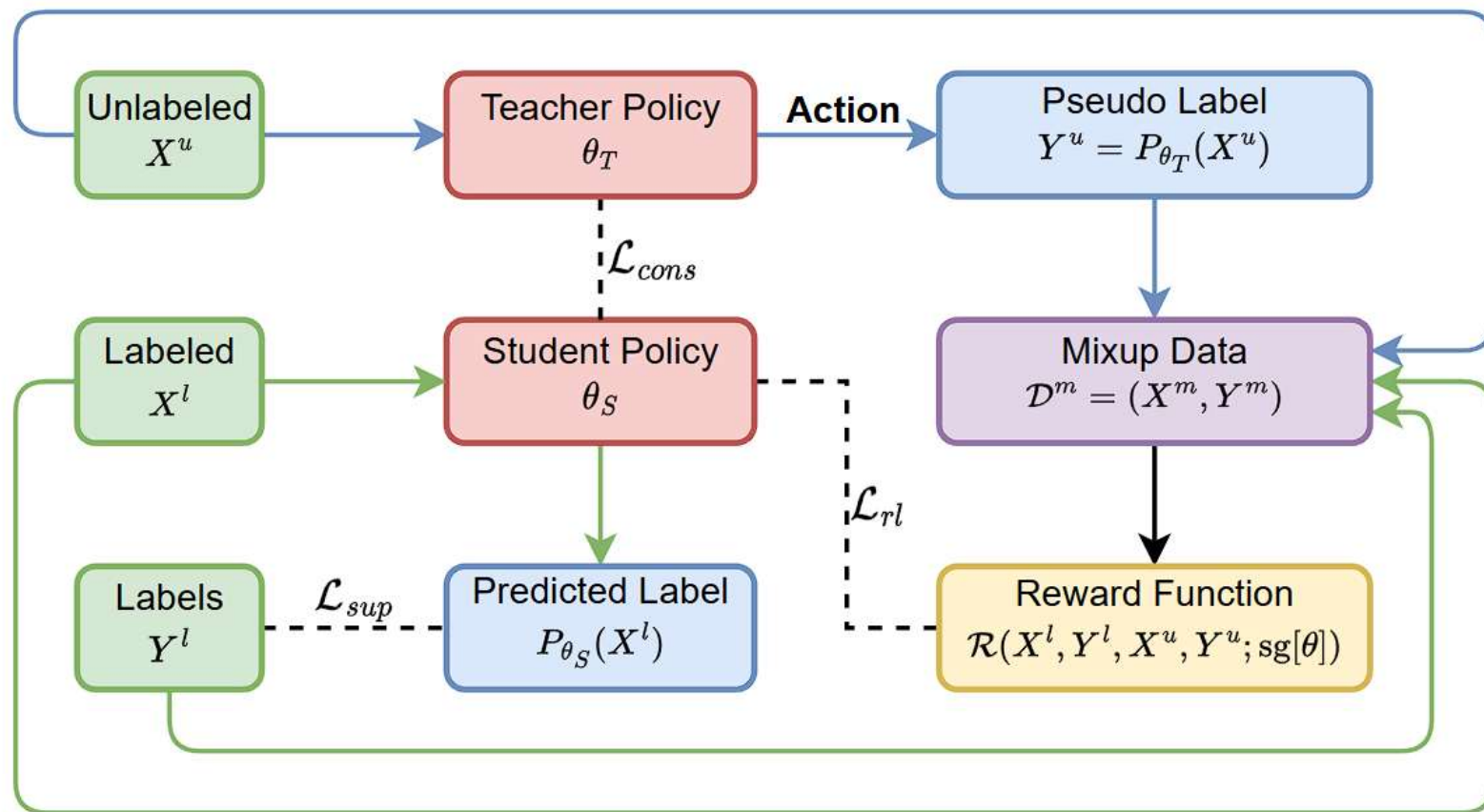
one-armed bandit: No need for γ and P

reward function: *data mixup*

$$\mathcal{D}^l \longrightarrow \tilde{\mathcal{D}}^l \quad r = \left[\frac{N^u}{N^l} \right]$$

$$x_i^m = \mu x_i^u + (1 - \mu) x_i^l, \quad y_i^m = \mu y_i^u + (1 - \mu) y_i^l$$

$$\mathcal{R}(s, a; \text{sg}[\theta]) = \mathcal{R}(X^l, Y^l, X^u, Y^u; \text{sg}[\theta]) = -\frac{1}{C \cdot N^m} \sum_{i=1}^{N^m} \|P_{\theta}(x_i^m) - y_i^m\|_2^2$$



$$\begin{aligned} \mathcal{L}_{rl} &= -\mathbb{E}_{\mathbf{y}_i^u \sim \pi_\theta} \text{KL}(\mathbf{e}, \mathbf{y}_i^u) \mathcal{R}(s, a; \text{sg}[\theta]) & \mathcal{L}^{\text{sup}} &= \mathbb{E}_{(x^l, \mathbf{y}^l) \in \mathcal{D}^l} [\ell_{CE}(P_{\theta_S}(x^l), \mathbf{y}^l)] \\ &= -\mathbb{E}_{x_i^u \in \mathcal{D}_u} \text{KL}(\mathbf{e}, P_{\theta}(x_i^u)) \mathcal{R}(s, a; \text{sg}[\theta]) & \mathcal{L}^{\text{cons}} &= \mathbb{E}_{x^u \in \mathcal{D}^u} [\ell_{\text{KL}}(P_{\theta_S}(x^u), P_{\theta_T}(x^u))] \end{aligned}$$

$$\mathcal{L}(\theta_S) = \mathcal{L}_{rl} + \lambda_1 \mathcal{L}_{\text{sup}} + \lambda_2 \mathcal{L}_{\text{cons}}$$

Table 1: Performance of RLGSSL and state-of-the-art SSL algorithms with the CNN-13 network. We report the average test errors and the standard deviations of 5 trials.

Dataset Number of Labeled Samples	CIFAR-10			CIFAR-100	
	1000	2000	4000	4000	10000
Supervised	39.95 _(0.75)	27.67 _(0.12)	20.42 _(0.21)	58.31 _(0.89)	44.56 _(0.30)
Supervised + MixUp [38]	31.83 _(0.65)	24.22 _(0.15)	17.37 _(0.35)	54.87 _(0.07)	40.97 _(0.47)
Π -model [6]	28.74 _(0.48)	17.57 _(0.44)	12.36 _(0.17)	55.39 _(0.55)	38.06 _(0.37)
Temp-ensemble [6]	25.15 _(1.46)	15.78 _(0.44)	11.90 _(0.25)	-	38.65 _(0.51)
Mean Teacher[8]	21.55 _(0.53)	15.73 _(0.31)	12.31 _(0.28)	45.36 _(0.49)	35.96 _(0.77)
VAT [5]	18.12 _(0.82)	13.93 _(0.33)	11.10 _(0.24)	-	-
SNTG [10]	18.41 _(0.52)	13.64 _(0.32)	10.93 _(0.14)	-	37.97 _(0.29)
Learning to Reweight [42]	11.74 _(0.12)	-	9.44 _(0.17)	46.62 _(0.29)	37.31 _(0.47)
MT + Fast SWA [9]	15.58	11.02	9.05	-	33.62 _(0.54)
ICT [11]	12.44 _(0.57)	8.69 _(0.15)	7.18 _(0.24)	40.07 _(0.38)	32.24 _(0.16)
RLGSSL (Ours)	9.15 _(0.57)	6.90 _(0.11)	6.11 _(0.10)	36.92 _(0.45)	29.12 _(0.20)

Table 2: Performance of RLGSSL and state-of-the-art SSL algorithms with the CNN-13 network. We report the average test errors and the standard deviations of 5 trials.

	VAT [5]	Π -model [6]	Temp-ensemble [6]	MT [8]	ICT [11]	SNTG [10]	RLGSSL (Ours)
SVHN/500	-	6.65 _(0.53)	5.12 _(0.13)	4.18 _(0.27)	4.23 _(0.15)	3.99 _(0.24)	3.12 _(0.07)
SVHN/1000	5.42 _(0.00)	4.82 _(0.17)	4.42 _(0.16)	3.95 _(0.19)	3.89 _(0.04)	3.86 _(0.27)	3.05 _(0.04)

Main Results

Table 1: Performance of RLGSSL and state-of-the-art SSL algorithms with the CNN-13 network. We report the average test errors and the standard deviations of 5 trials.

Dataset Number of Labeled Samples	CIFAR-10			CIFAR-100	
	1000	2000	4000	4000	10000
Supervised	39.95 _(0.75)	27.67 _(0.12)	20.42 _(0.21)	58.31 _(0.89)	44.56 _(0.30)
Supervised + MixUp [38]	31.83 _(0.65)	24.22 _(0.15)	17.37 _(0.35)	54.87 _(0.07)	40.97 _(0.47)
Π-model [6]	28.74 _(0.48)	17.57 _(0.44)	12.36 _(0.17)	55.39 _(0.55)	38.06 _(0.37)
Temp-ensemble [6]	25.15 _(1.46)	15.78 _(0.44)	11.90 _(0.25)	-	38.65 _(0.51)
Mean Teacher[8]	21.55 _(0.53)	15.73 _(0.31)	12.31 _(0.28)	45.36 _(0.49)	35.96 _(0.77)
VAT [5]	18.12 _(0.82)	13.93 _(0.33)	11.10 _(0.24)	-	-
SNTG [10]	18.41 _(0.52)	13.64 _(0.32)	10.93 _(0.14)	-	37.97 _(0.29)
Learning to Reweight [42]	11.74 _(0.12)	-	9.44 _(0.17)	46.62 _(0.29)	37.31 _(0.47)
MT + Fast SWA [9]	15.58	11.02	9.05	-	33.62 _(0.54)
ICT [11]	12.44 _(0.57)	8.69 _(0.15)	7.18 _(0.24)	40.07 _(0.38)	32.24 _(0.16)
RLGSSL (Ours)	9.15 _(0.57)	6.90 _(0.11)	6.11 _(0.10)	36.92 _(0.45)	29.12 _(0.20)

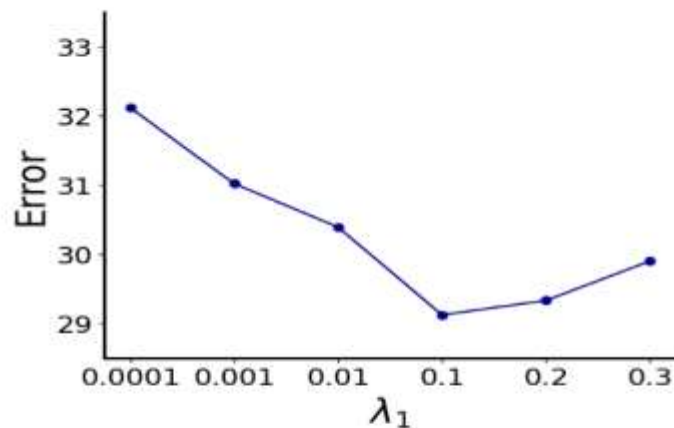
Table 2: Performance of RLGSSL and state-of-the-art SSL algorithms with the CNN-13 network. We report the average test errors and the standard deviations of 5 trials.

	VAT [5]	Π-model [6]	Temp-ensemble [6]	MT [8]	ICT [11]	SNTG [10]	RLGSSL (Ours)
SVHN/500	-	6.65 _(0.53)	5.12 _(0.13)	4.18 _(0.27)	4.23 _(0.15)	3.99 _(0.24)	3.12 _(0.07)
SVHN/1000	5.42 _(0.00)	4.82 _(0.17)	4.42 _(0.16)	3.95 _(0.19)	3.89 _(0.04)	3.86 _(0.27)	3.05 _(0.04)

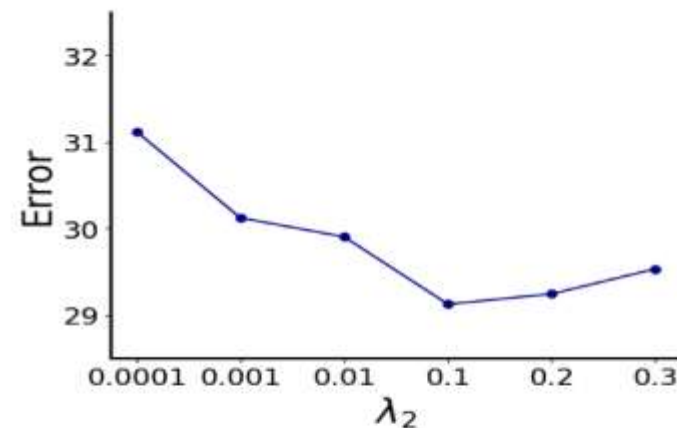
Ablation Study

Table 4: Ablation study results. We report the test errors on CIFAR-100 with 10000, and 4000 labels on CNN-13 backbone 5 trials.

	RLGSSL	-w/o \mathcal{L}_{rl}	-w/o \mathcal{L}_{sup}	-w/o \mathcal{L}_{cons}	-w/o EMA	-w/o mixup
CIFAR-100/4000	36.92 _(0.45)	44.92 _(0.55)	39.52 _(0.58)	38.78 _(0.48)	43.12 _(0.52)	40.12 _(0.51)
CIFAR-100/10000	29.12 _(0.20)	33.12 _(0.52)	32.67 _(0.45)	31.48 _(0.32)	32.84 _(0.45)	31.48 _(0.32)
	RLGSSL	$\mathcal{R} = 1$	$\mathcal{R} : \mu = 0$	$\mathcal{R}(L_2 \rightarrow \text{KL})$	$\mathcal{R}(L_2 \rightarrow \text{JS})$	$\mathcal{R} : \text{w/o sg}[\theta]$
CIFAR-100/4000	36.92 _(0.45)	39.52 _(0.63)	39.54 _(0.33)	38.02 _(0.42)	39.52 _(0.45)	40.62 _(0.55)
CIFAR-100/10000	29.12 _(0.20)	31.25 _(0.62)	32.37 _(0.57)	31.12 _(0.52)	31.39 _(0.68)	32.12 _(0.62)



(a) λ_1



(b) λ_2

Figure 2: Sensitivity analysis for four hyperparameters λ_1 and λ_2 CIFAR-100 using 10000 labeled samples (a) λ_1 , (b) λ_2 .

Bridging Supervised Learning and Reinforcement Learning in Math Reasoning

Huayu Chen^{1,2} Kaiwen Zheng^{1,2} Qinsheng Zhang² Ganqu Cui¹ Yin Cui²

Haotian Ye^{2,3} Tsung-Yi Lin² Ming-Yu Liu² Jun Zhu^{1†} Haoxiang Wang²

¹Tsinghua University ²NVIDIA ³Stanford University

<https://research.nvidia.com/labs/dir/Negative-aware-Fine-Tuning>

MDP: $(S, A, R, P, \gamma) \implies$ Language-augmented MDP: (V, S, A, R, P, γ) where $\gamma=1$

\downarrow
vocabulary

A specific token: $w \in V$

$S: S \subset V^M$ $s = (w_1, w_2, w_3, \dots, w_M)$ s_0 : prompt

$A: S \subset V^N$ $a = (w_1, w_2, \dots, w_N)$

$R: r = R(s, a)$

$P: S \times V \rightarrow S$ $s_{i+1} = (s_i, w_i) = (s_0, w_{1:i+1})$ **auto-regressive paradigm**

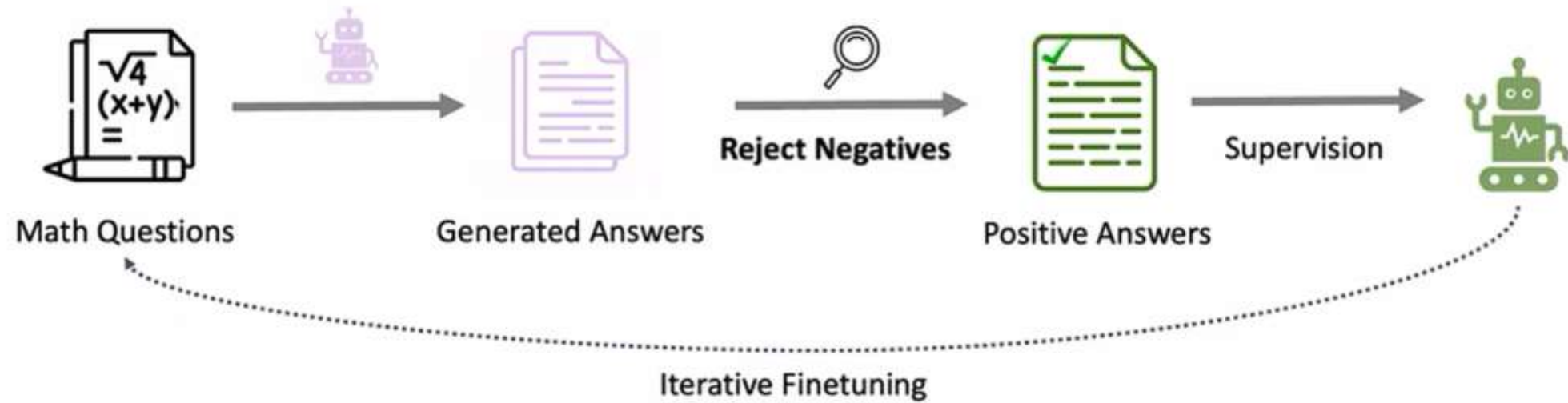
Token-level policy: $\pi(w_i | s_0, w_{1:i-1})$

Sentence-level policy: $\pi(a | s_0) = \prod_{i=1}^N \pi(w_i | s_0, w_{1:i-1})$

Policy objective:

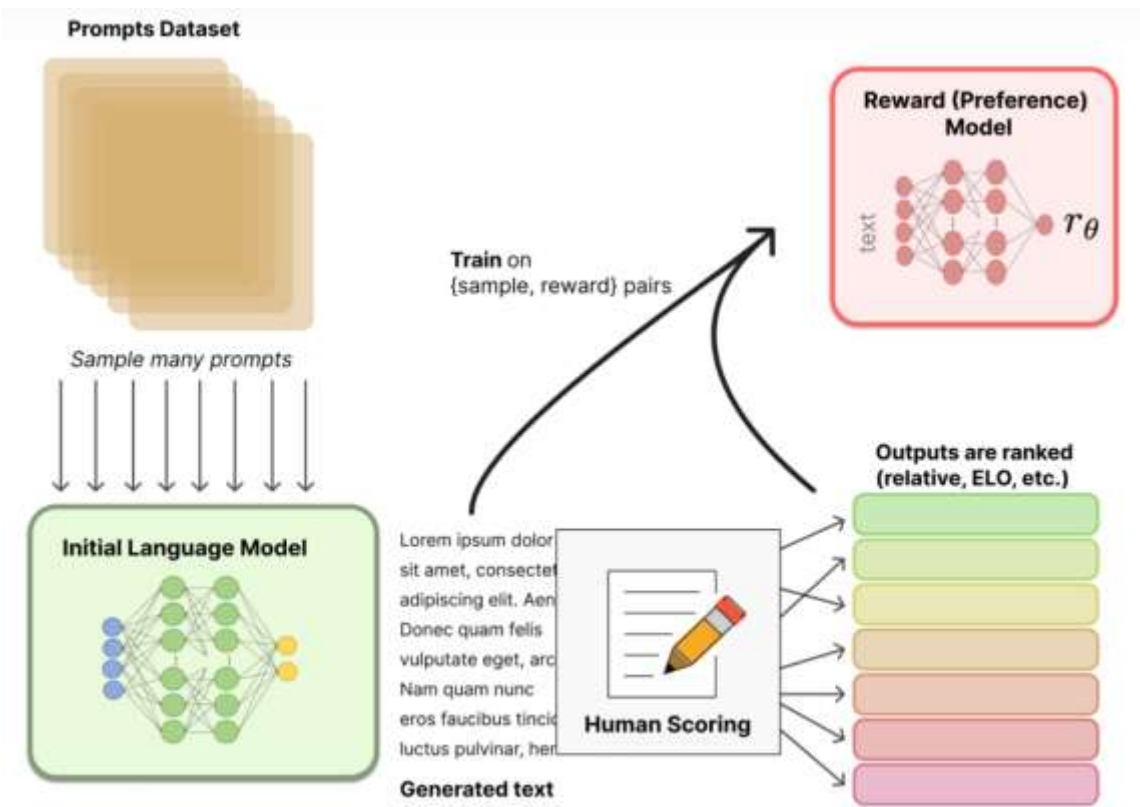
$$J(\pi) = \mathbb{E}_{s_0 \sim D} \left[\mathbb{E}_{a \sim \pi(\cdot | s_0)} [r(s_0, a)] \right]$$

Rejection sampling Fine-Tuning (RFT)



How to leverage negative feedbacks?

Why use RL?



Prompt: s

Candidate answers: a_1, a_2, a_3, a_4

Quality rank: $a_1 > a_2 > a_3 > a_4$

How can we optimize the model?

Supervised learning: $\max \log p(a_1|s)$



Data waste & poor generalization

A solution:

for each answers: $\max R_i \cdot \log p(a_i|s)$

gradient: $R_i \cdot \nabla_\theta \log p_\theta(a_i|s)$



✓ **The policy gradient of RL**

Connection between SL and RL

Supervised Learning

$$\max_{\theta} \mathbb{E}_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{q})} \log \pi_{\theta}(\mathbf{a}|\mathbf{q}) \Leftrightarrow \min_{\theta} D_{\text{KL}} [\pi(\mathbf{a}|\mathbf{q}) \parallel \pi_{\theta}(\mathbf{a}|\mathbf{q})]$$

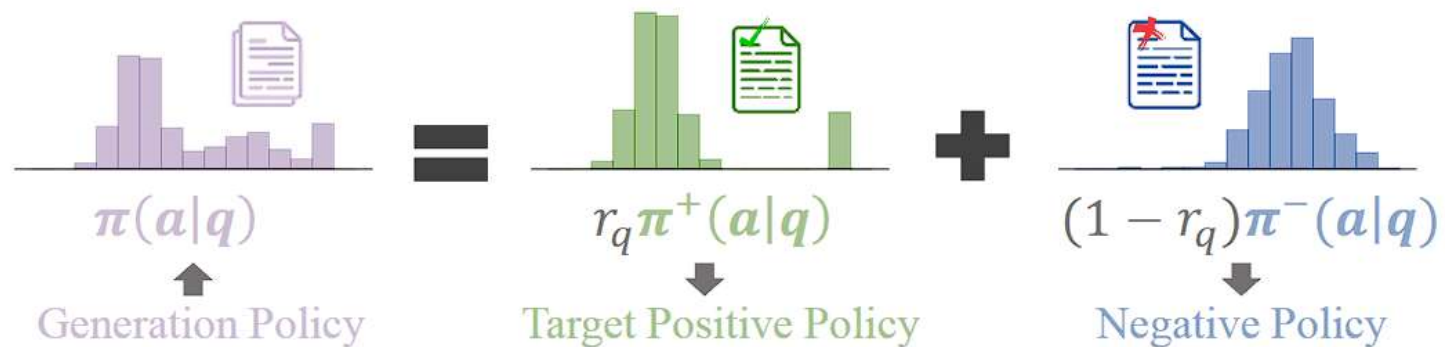
Reinforcement Learning

$$\max_{\theta} J(\theta) := \mathbb{E}_{\mathbf{a} \sim \pi_{\theta}(\mathbf{a}|\mathbf{q})} r(\mathbf{q}, \mathbf{a})$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\mathbf{a} \sim \pi_{\theta}(\mathbf{a}|\mathbf{q})} \nabla_{\theta} [r(\mathbf{q}, \mathbf{a}) \log \pi_{\theta}(\mathbf{a}|\mathbf{q})]$$

on-policy constraint

leverage of negative feedbacks



target policy

$$\pi^+(\mathbf{a}|\mathbf{q}) := \pi(\mathbf{a}|\mathbf{q}, r=1) = \frac{\pi(\mathbf{a}|\mathbf{q})p(r=1|\mathbf{q}, \mathbf{a})}{\sum_A \pi(\mathbf{a}|\mathbf{q})p(r=1|\mathbf{q}, \mathbf{a})}$$

$$\pi^-(\mathbf{a}|\mathbf{q}) := \pi(\mathbf{a}|\mathbf{q}, r=0) = \frac{\pi(\mathbf{a}|\mathbf{q})[1 - p(r=1|\mathbf{q}, \mathbf{a})]}{\sum_A \pi(\mathbf{a}|\mathbf{q})[1 - p(r=1|\mathbf{q}, \mathbf{a})]}$$

$$r_q \pi^+(\mathbf{a}|\mathbf{q}) + [1 - r_q] \pi^-(\mathbf{a}|\mathbf{q}) = \pi(\mathbf{a}|\mathbf{q})$$

$$\pi_{\theta}^-(\mathbf{a}|\mathbf{q}) := \frac{\pi(\mathbf{a}|\mathbf{q}) - r_q \pi_{\theta}^+(\mathbf{a}|\mathbf{q})}{1 - r_q}$$

Theorem 3.1 (Policy Optimization with Negative Answers). Consider the maximum-likelihood objective for training the implicit negative policy π_{θ}^{-} :

$$\max_{\theta} \mathbb{E}_{p(\mathbf{q})\pi^{-}(\mathbf{a}|\mathbf{q})} [\log \pi_{\theta}^{-}(\mathbf{a}|\mathbf{q})] \Leftrightarrow \min_{\theta} \left[-\mathbb{E}_{(\mathbf{q},\mathbf{a})\sim\mathcal{D}^{-}} \log \frac{\pi(\mathbf{a}|\mathbf{q}) - r_{\mathbf{q}}\pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})}{1 - r_{\mathbf{q}}} \right] \quad (8)$$

Assuming unlimited data and model capacity, the optimal solution for solving Eq. 8 is

$$\forall \mathbf{q}, \mathbf{a} : \pi_{\theta}^{+}(\mathbf{a}|\mathbf{q}) = \pi^{+}(\mathbf{a}|\mathbf{q})$$

To further utilize positive data,

$$\mathcal{L}_{(\mathbf{a},\mathbf{q},r)\sim\mathcal{D}}^{\text{NFT}}(\theta) = r \left[-\log \frac{\pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})}{\pi(\mathbf{a}|\mathbf{q})} \right] + (1 - r) \left[-\log \frac{1 - r_{\mathbf{q}} \frac{\pi_{\theta}^{+}(\mathbf{a}|\mathbf{q})}{\pi(\mathbf{a}|\mathbf{q})}}{1 - r_{\mathbf{q}}} \right]$$

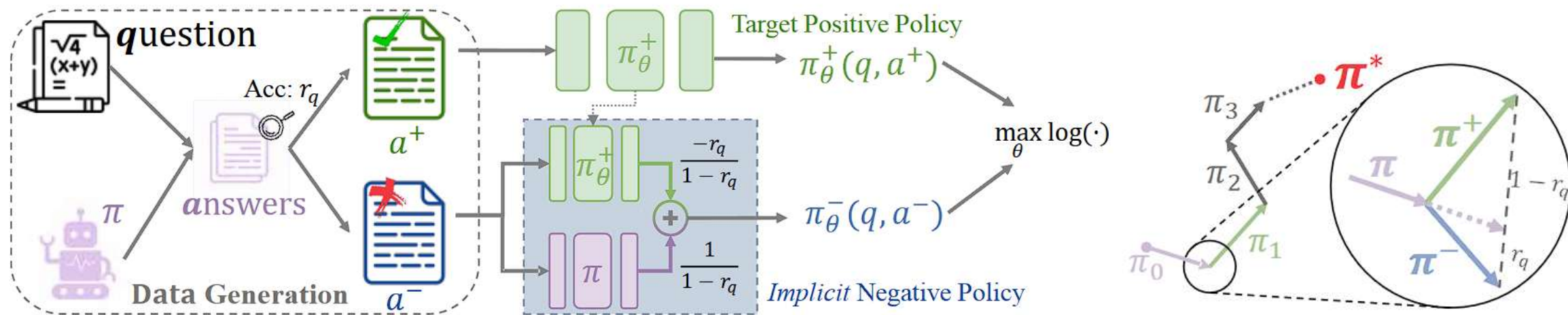
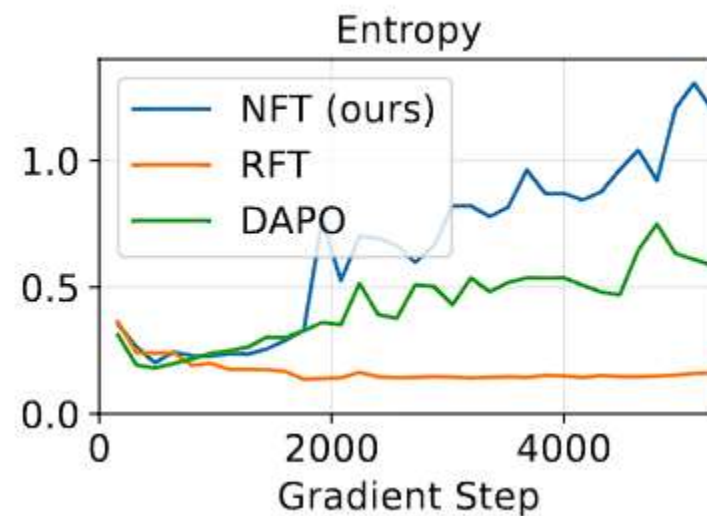
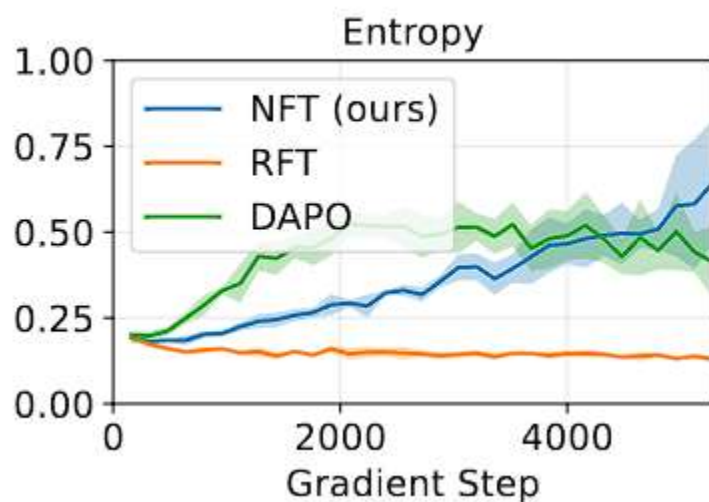
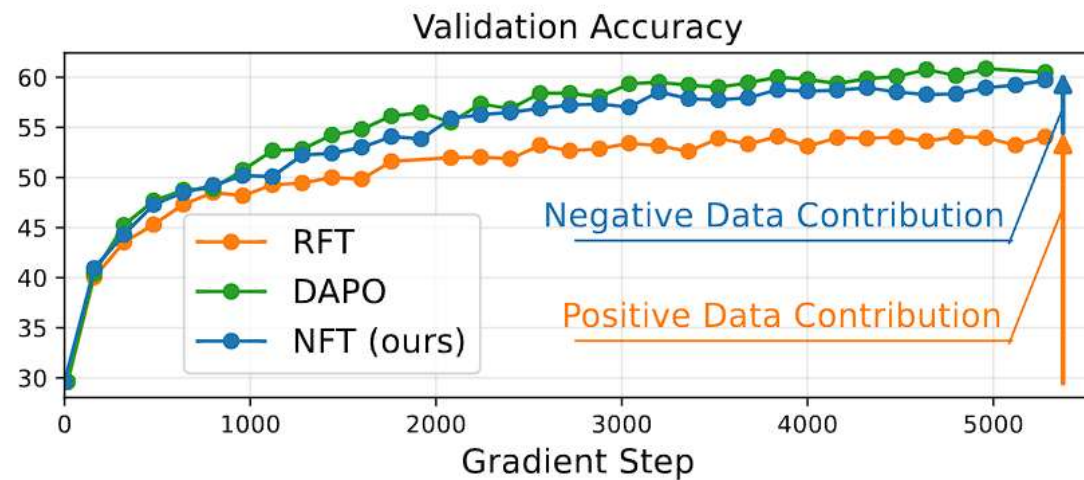


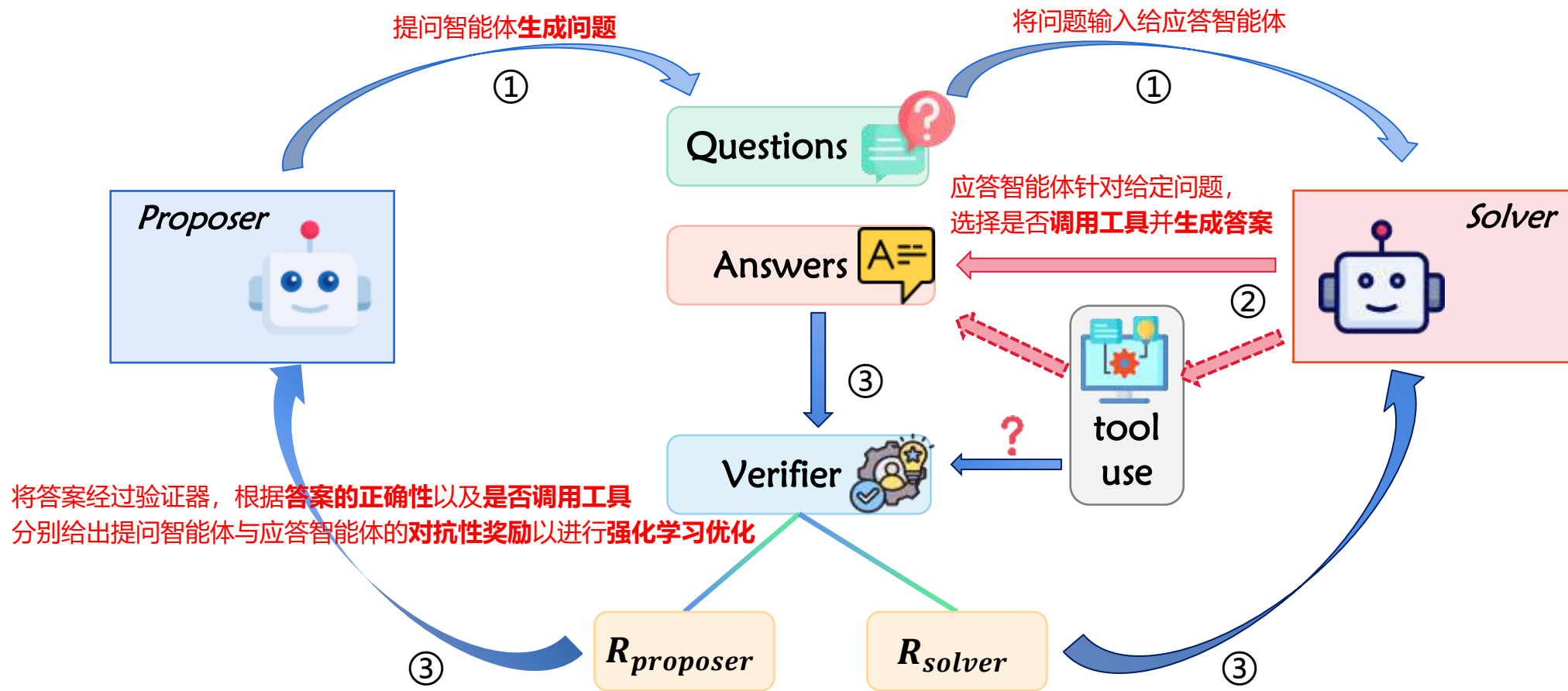
Table 1: NFT performs competitively compared with other algorithms. We report avg@32 for AIME24, AIME25, and AMC23 and avg@1 for others. Numbers within 1 % of the max are bolded.

Model	AIME24	MATH500	AIME25	AMC23	Olympiad	Minerva	Average
Qwen2.5-Math-7B	13.3	69.0	5.5	45.8	34.7	21.3	31.6
<i>Preference fine-tuning</i>							
+ DPO	29.8	79.8	13.8	83.2	48.0	39.0	48.9
<i>Reinforcement fine-tuning</i>							
+ GRPO	30.2	80.4	17.1	79.5	51.8	38.2	49.5
+ Dr. GRPO	31.8	83.4	15.7	80.2	49.6	38.2	49.8
+ DAPO	33.1	81.6	18.7	85.0	49.9	39.3	51.2
<i>Supervised fine-tuning</i>							
+ RFT	33.7	79.8	13.4	79.7	44.3	38.6	48.3
+ NFT	32.0	83.2	18.3	88.5	47.3	40.8	51.7
Qwen2.5-32B	4.1	68.6	1.0	45.0	31.1	27.9	29.6
+ DAPO	44.1	89.2	33.4	90.9	54.1	47.5	59.9
+ RFT	29.9	86.2	19.1	92.4	45.3	44.1	52.8
+ NFT	37.8	88.4	31.5	93.8	55.0	48.9	59.2



Negative feedback enhances exploration

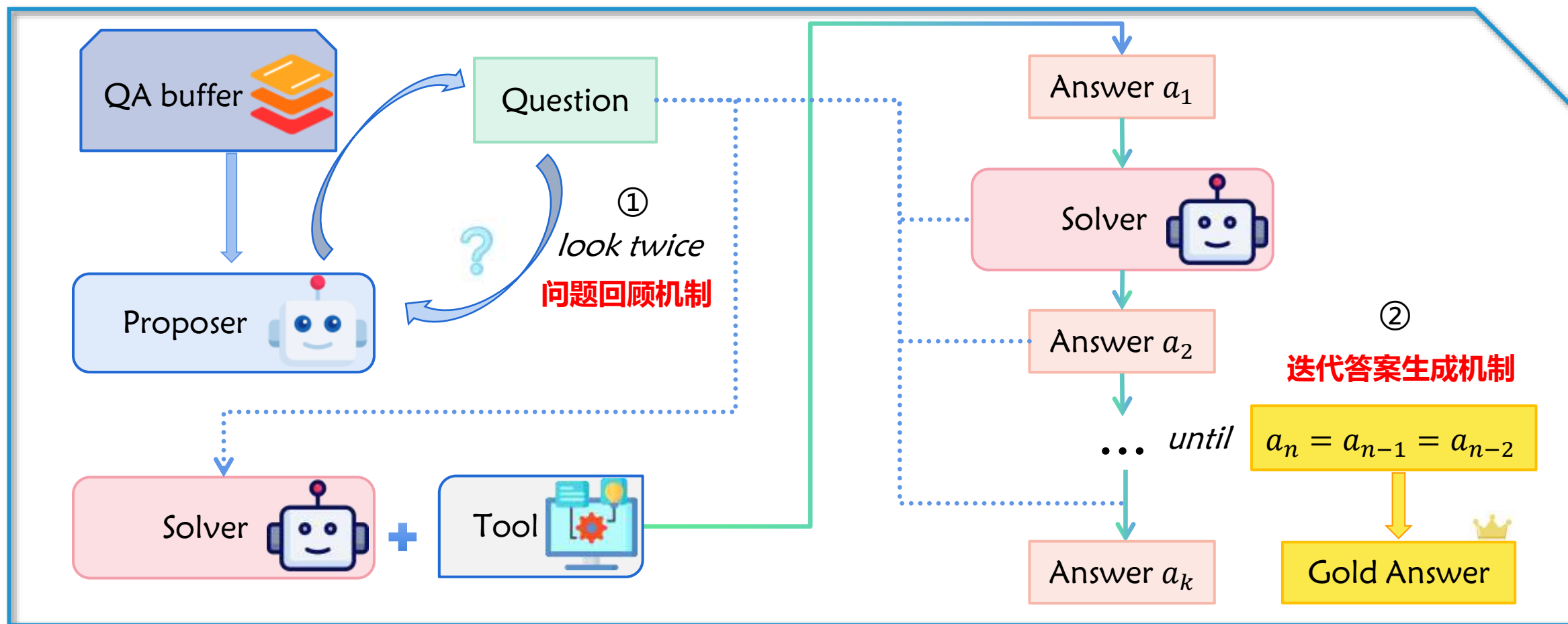
总体技术路线



✓ 该过程可形式化为一个Stackelberg博弈或迭代对抗训练框架

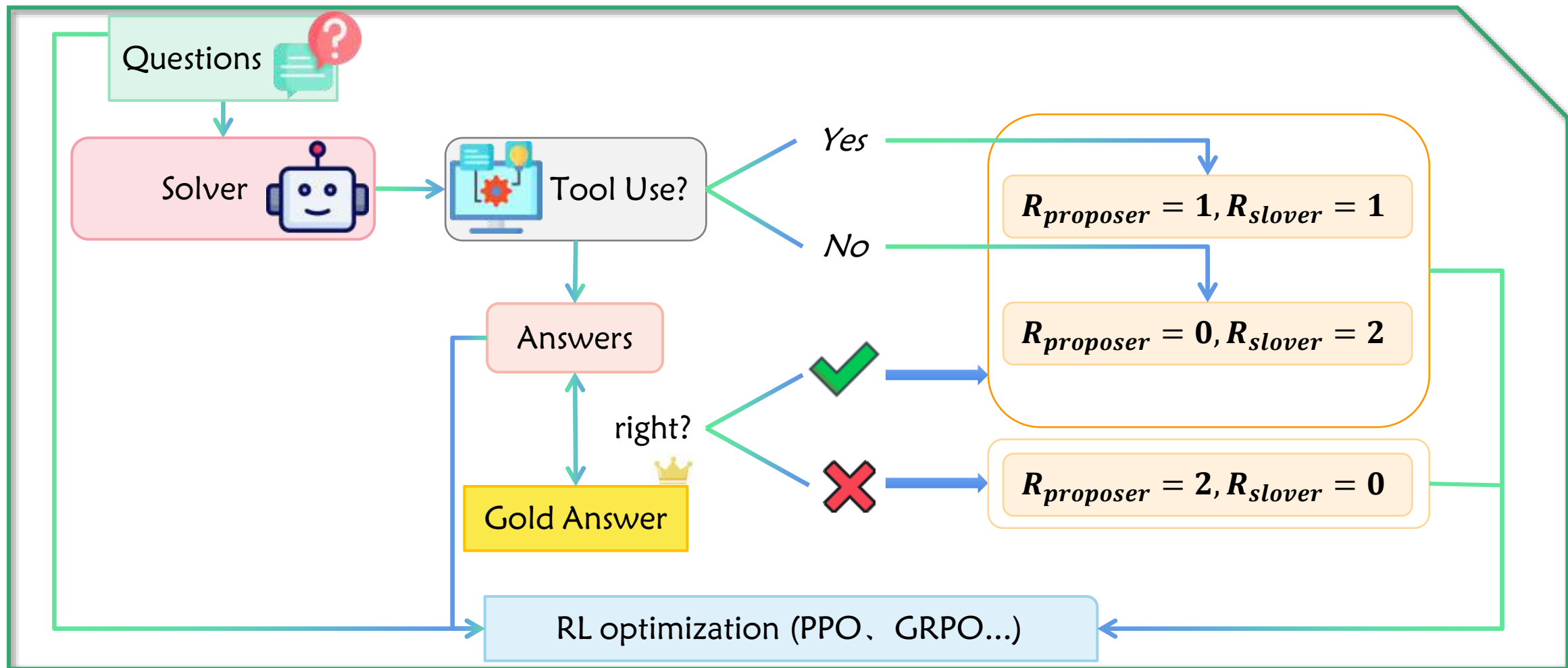
技术方案：自动化问题-答案生成

目的：自动化获得**有效的问题**以实现无人工参与的动态探测知识边界，并通过大模型对工具的反复调用获得**对应真实答案**，为后续可验证奖励提供参考，从而实现基于强化学习的模型能力优化。



技术方案：对抗性奖励设计

目的：通过基于**答案的正确性**与**是否调用工具**设计**对抗性奖励**，从而实现提问智能体与应答智能体的对抗博弈，鼓励提问智能体尽可能提出难度较高的问题，并且鼓励应答智能体尽可能通过自身知识解决问题。



Thanks