

# CREAM: Consistency Regularized Self-Rewarding Language Models

Zhaoyang Wang<sup>1</sup>, Weilei He<sup>2</sup>, Zhiyuan Liang<sup>3</sup>, Xuchao Zhang<sup>4</sup>, Chetan Bansal<sup>4</sup>, Ying Wei<sup>2</sup>, Weitong Zhang<sup>1</sup>, Huaxiu Yao<sup>1</sup>

<sup>1</sup>University of North Carolina at Chapel Hill, <sup>2</sup>Nanyang Technological University, <sup>3</sup>National University of Singapore, <sup>4</sup>Microsoft Research



**tl; dr: This paper introduces consistency-based regularization to self-rewarding language models.**

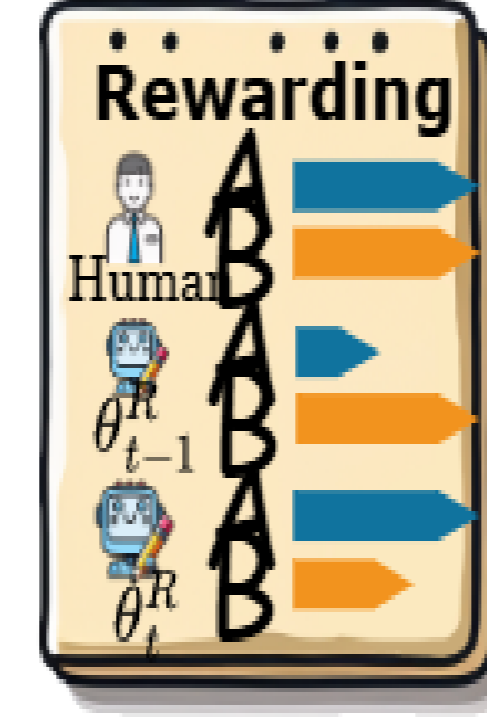
## 1 Background

Self-rewarding language model (SRLM) utilizes the same LLM as both policy model (generating responses) and reward model (ranking responses). The ranked responses form preference pairs for iterative DPO fine-tuning. **However, SRLM faces challenges in generating reliable and accurate rewards for annotating the preference pairs.**, as forcing the reward model into preference judgments leads to **overconfident labeling**.

**Q: Can you recommend a fun DIY project that can be accomplished in a weekend?**

**A** Here are some fun DIY projects that can be completed in a weekend: **1. Macrame Plant Hanger:** Create a bohemian-inspired plant hanger using macrame cord, wooden dowels, and a few basic knots. You can customize it with different color.....

**B** I'd be happy to suggest some fun DIY projects that can be completed over a weekend: **1. DIY Photo Display:** Create a unique photo display using reclaimed wood, string, and clothespins. This project is perfect for showcasing your favorite memories.....



## 2 Generalized Iterative Preference Fine-tuning Framework

Existing iterative preference fine-tuning methods including SRLM and RLAIIF can be defined as follows.  $z$  is the preference labeling function

$$\mathcal{L}(\theta, z) = \mathcal{L}_{\text{SFT}}(\theta; \mathcal{D}_S) + \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_U; \mathbf{y}, \mathbf{y}' \sim \pi_{\theta_t}(\cdot | \mathbf{x})} [\mathcal{L}_{\text{DPO}}(\theta; \mathbf{y}, \mathbf{y}', \mathbf{x}, z)]. \quad \text{Eq. 3.1}$$

The DPO loss can be defined as follows:

$$\mathcal{L}_{\text{DPO}}(\theta; \mathbf{y}, \mathbf{y}', \mathbf{x}, z) = -z(\mathbf{y}, \mathbf{y}', \mathbf{x}) \log \sigma \left( \log \left( \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} \right) - \log \left( \frac{\pi_{\theta}(\mathbf{y}' | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}' | \mathbf{x})} \right) \right) - (1 - z(\mathbf{y}, \mathbf{y}', \mathbf{x})) \log \sigma \left( \log \left( \frac{\pi_{\theta}(\mathbf{y}' | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}' | \mathbf{x})} \right) - \log \left( \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} \right) \right), \quad \text{Eq. 3.2}$$

Step1: Preference labeling step

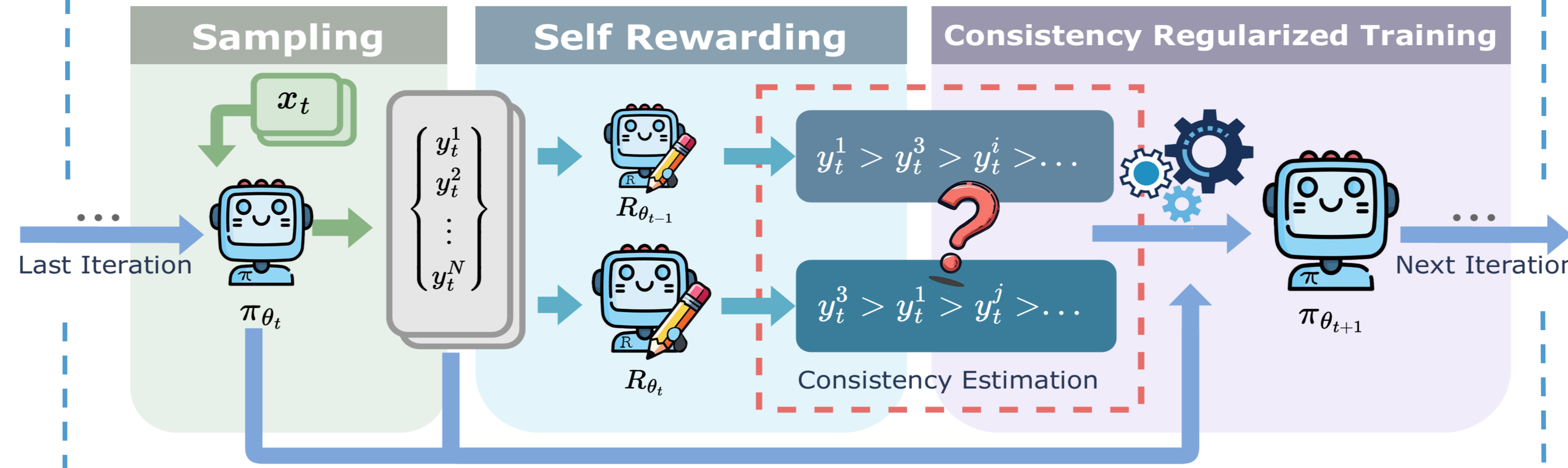
$$z_{t+1}(\mathbf{y}, \mathbf{y}', \mathbf{x}) = \mathbb{1} [\log \pi_{\theta_t}(\mathbf{y} | \mathbf{x}) - \log \pi_{\text{ref}}(\mathbf{y} | \mathbf{x}) \geq \log \pi_{\theta_t}(\mathbf{y}' | \mathbf{x}) - \log \pi_{\text{ref}}(\mathbf{y}' | \mathbf{x})]. \quad \text{Eq. 3.3}$$

Differences of rewards

Step2: Learning Step

$$\theta_{t+1} = \arg \min_{\theta} \mathcal{L}(\theta, z_{t+1}).$$

## 3 CREAM: Consistency Regularized SRLM



Forcing the model to over-confidently distinguish between responses  $\{y, y'\}$  of similar quality can harm SRLM training. Ideally, **oracle reward scores for such responses of similar quality should be very close**, resulting in inconsistent rankings across multiple reward models. **CREAM addresses this by preventing the model from learning from preference pairs with low consistency.**

Introducing a regularization term to Eq. 3.1

$$\mathcal{L}(\theta, z) = \mathcal{L}_{\text{SFT}}(\theta; \mathcal{D}_S) + \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_U; \mathbf{y}, \mathbf{y}' \sim \pi_{\theta_t}(\cdot | \mathbf{x})} [\mathcal{L}_{\text{DPO}}(\theta; \mathbf{y}, \mathbf{y}', \mathbf{x}, z) + \lambda \mathcal{L}_{\text{Reg}}(\theta; \mathbf{y}, \mathbf{y}', \mathbf{x})], \quad \text{Eq. 3.4}$$

Regularization term

Regularization term is defined as follows:

$$\mathcal{L}_{\text{Reg}}(\theta; \mathbf{y}, \mathbf{y}', \mathbf{x}) = -\log \sigma \left( \log \left( \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} \right) - \log \left( \frac{\pi_{\theta}(\mathbf{y}' | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}' | \mathbf{x})} \right) \right) - \log \sigma \left( \log \left( \frac{\pi_{\theta}(\mathbf{y}' | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}' | \mathbf{x})} \right) - \log \left( \frac{\pi_{\theta}(\mathbf{y} | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y} | \mathbf{x})} \right) \right). \quad \text{Eq. 3.5}$$

And the expectation loss is as follows, where  $\mu(z = 0) = \mu(z = 1) = 0.5$

$$\mathbb{E}_{\mathbf{y}, \mathbf{y}' \sim \pi_{\theta_t}(\cdot | \mathbf{x})} \mathcal{L}_{\text{Reg}}(\theta; \mathbf{y}, \mathbf{y}', \mathbf{x}) = 2 \mathbb{KL}(u(\cdot) \parallel P_{\theta}(\cdot)), \quad \text{Eq. 3.6}$$

**CREAM Loss.**  $\mathcal{C}$ , serving as the consistency rate, is measured by Kendall correlation between rankings offered by the current model and last iteration's model.

$$\mathcal{L}(\theta, z) = \frac{1}{1 + 2\lambda} \mathcal{L}_{\text{SFT}}(\theta; \mathcal{D}_S) + \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_U; \mathbf{y}, \mathbf{y}' \sim \pi_{\theta_t}(\cdot | \mathbf{x})} [\mathcal{C}_{\lambda} \mathcal{L}_{\text{DPO}}(\theta; \mathbf{y}, \mathbf{y}', \mathbf{x}, z) + (1 - \mathcal{C}_{\lambda}) \mathcal{L}_{\text{DPO}}(\theta; \mathbf{y}, \mathbf{y}', \mathbf{x}, 1 - z)], \quad \text{Eq. 3.7}$$

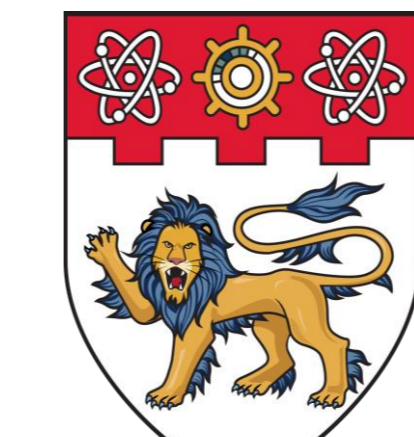
Reversed DPO

Contact

{ zhaoyang, huaxiu } @cs.unc.edu  
weitongz@unc.edu

Presented at ICLR 2025

Paper Link

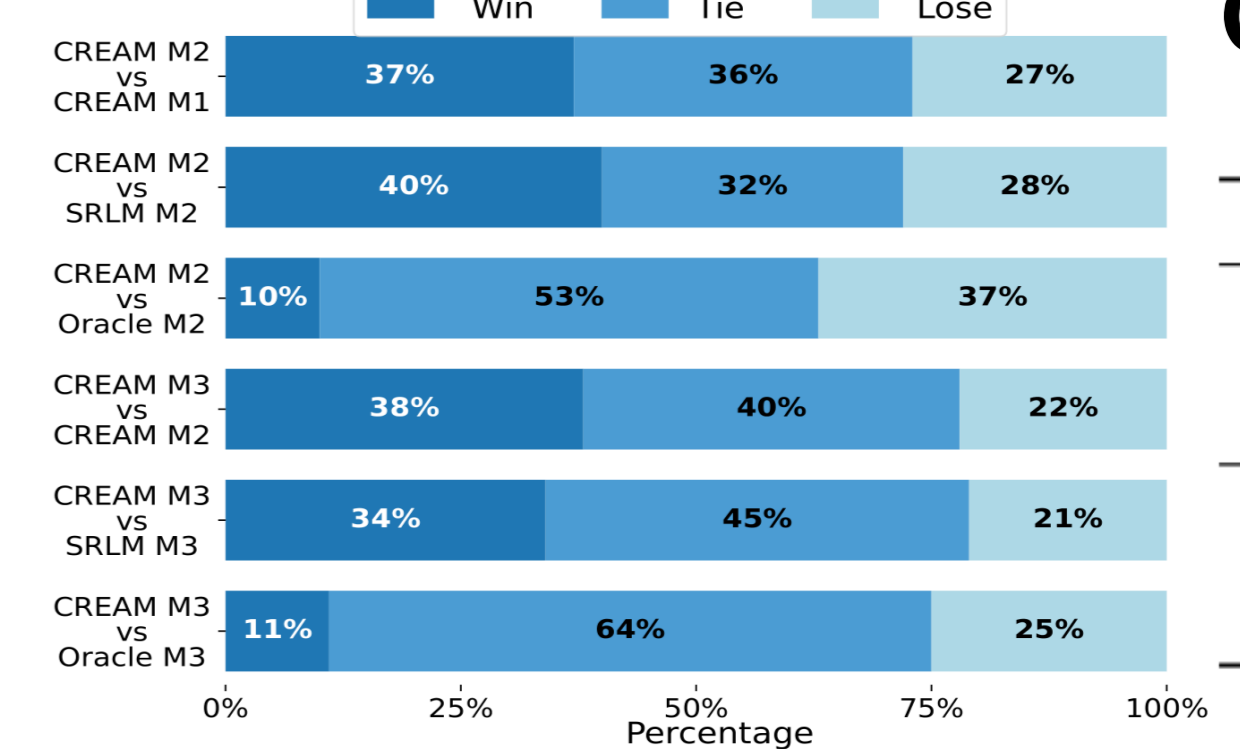


## 4 Experiments & Analysis

### Main Results on NLP benchmarks

Model	Method	Reward	Iteration	Arc-Easy	Arc-Challenge	OpenBookQA	SIQA	GSM8K	Average
Llama-3	Initial	-	M0	86.29	80.37	86.00	68.58	78.01	79.85
	SFT	-	M1	86.78	80.14	86.40	69.50	78.39	80.24
	Oracle	External	M2	89.60 $\uparrow$	82.17 $\uparrow$	90.00 $\uparrow$	72.88 $\uparrow$	80.82 $\uparrow$	83.09 $\uparrow$
			M3	89.31 $\downarrow$	81.31 $\downarrow$	90.20 $\uparrow$	73.75 $\uparrow$	76.04 $\downarrow$	82.12 $\downarrow$
	SRLM	Self	M2	87.79 $\uparrow$	80.38 $\uparrow$	87.80 $\uparrow$	70.95 $\uparrow$	78.01 $\downarrow$	80.99 $\uparrow$
			M3	87.17 $\downarrow$	81.23 $\downarrow$	87.30 $\downarrow$	70.37 $\downarrow$	77.48 $\downarrow$	80.71 $\downarrow$
			M2	87.92 $\uparrow$	79.78 $\downarrow$	86.60 $\uparrow$	71.49 $\uparrow$	79.38 $\uparrow$	81.03 $\uparrow$
			M3	88.38 $\uparrow$	80.97 $\uparrow$	88.20 $\uparrow$	71.19 $\downarrow$	80.29 $\uparrow$	81.81 $\uparrow$
	CREAM w/o RC	Self	M2	88.26 $\uparrow$	79.86 $\downarrow$	86.80 $\uparrow$	69.55 $\uparrow$	79.98 $\uparrow$	80.89 $\uparrow$
			M3	88.09 $\downarrow$	80.55 $\uparrow$	87.20 $\uparrow$	71.39 $\uparrow$	79.23 $\downarrow$	81.29 $\uparrow$
	CREAM	Self	M2	88.89 $\uparrow$	80.89 $\uparrow$	88.00 $\uparrow$	69.79 $\uparrow$	81.04 $\uparrow$	81.72 $\uparrow$
			M3	89.52 $\uparrow$	83.36 $\uparrow$	90.20 $\uparrow$	72.06 $\uparrow$	81.73 $\uparrow$	83.37 $\uparrow$

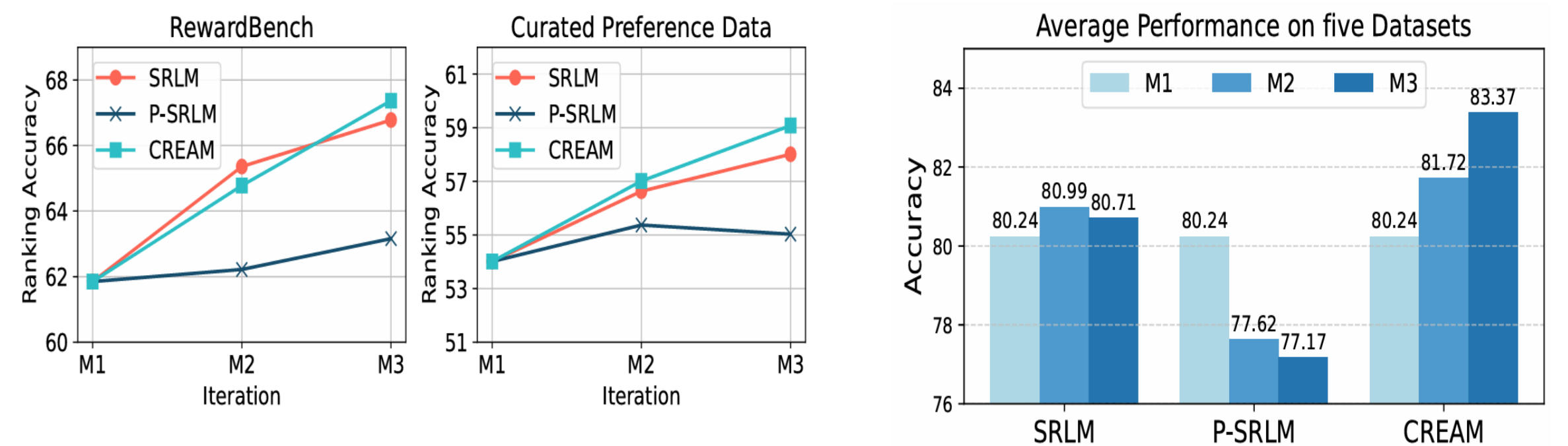
### Alignment Arena



### CREAM using External Oracle Reward Model

Method	BRM	Arc-E	Arc-C	OBQA	SIQA	GSM8K
Llama-3 M1	-	86.78	80.14	86.40	69.50	78.39
CREAM M2	M0	88.89	80.89	88.00	69.79	81.04
CREAM M2	Oracle	88.51	81.06	86.20	72.21	79.91
Llama-2 M1	-	60.44	48.46	63.20	50.77	23.88
CREAM M2	M0	58.97	47.53	62.80	50.43	24.41
CREAM M2	Oracle	62.42	48.72	66.00	51.13	22.52

### DPO Rewarding v.s. LLM-as-a-Judge



### Ranking Consistency

Iterations	Method	Consistency $\mathcal{C} \uparrow$	Kendall $\tau \uparrow$	Spearman $\uparrow$	TopOrder $\uparrow$
M2 vs M1	SRLM	0.39 $\pm$ 0.21	-0.22 $\pm$ 0.41	0.36 $\pm$ 0.24	0.03 $\pm$ 0.18
	CREAM	0.73 $\pm$ 0.18	0.46 $\pm$ 0.35	0.77 $\pm$ 0.19	0.19 $\pm$ 0.39
M3 vs M2	SRLM	0.46 $\pm$ 0.19	-0.08 $\pm$ 0.38	0.50 $\pm$ 0.22	0.12 $\pm$ 0.33
	CREAM	0.92 $\pm$ 0.09	0.84 $\pm$ 0.19	0.95 $\pm$ 0.07	0.59 $\pm$ 0.49

### Different consistency measurements

Iteration	Method	Arc-E	Arc-C	OBQA	SIQA	GSM8K
M1	-	86.78	80.14	86.40	69.50	78.39
M2	Spearman	86.95	82.00	85.40	70.05	78.77
	TopOrder	87.25	80.12	86.88	70.83	79.75
	Kendall (Ours)	88.89	80.89	88.00	69.79	81.04
M3	Spearman	88.76	81.83	90.00	70.98	79.15
	TopOrder	88.51	80.37	87.40	71.03	79.76
	Kendall (Ours)	89.52	83.36	90.20	72.06	81.73