



# Evolving Alignment via Asymmetric Self-Play

## Scalable Preference Fine-Tuning Beyond Static Human Prompts

$$\max_{\phi, \theta} \mathbb{E}_{\mathbf{x} \sim \pi_{\phi}(\cdot)} \left[ \underbrace{\mathbb{E}_{\mathbf{y} \sim \pi_{\theta}(\cdot|\mathbf{x})} [r(\mathbf{x}, \mathbf{y})] - \beta_1 \cdot \mathbb{D}_{\text{KL}}[\pi_{\theta}(\mathbf{y}|\mathbf{x}) \parallel \pi_{\text{SFT}}(\mathbf{y}|\mathbf{x})]}_{\text{solver} \sim \text{“regret minimization”}} \right] - \underbrace{\beta_2 \cdot \mathbb{D}_{\text{KL}}[\pi_{\phi}(\mathbf{x}) \parallel p_{\text{ref}}(\mathbf{x})]}_{\text{creator} \sim \text{“regret maximization” (implicit)}}.$$

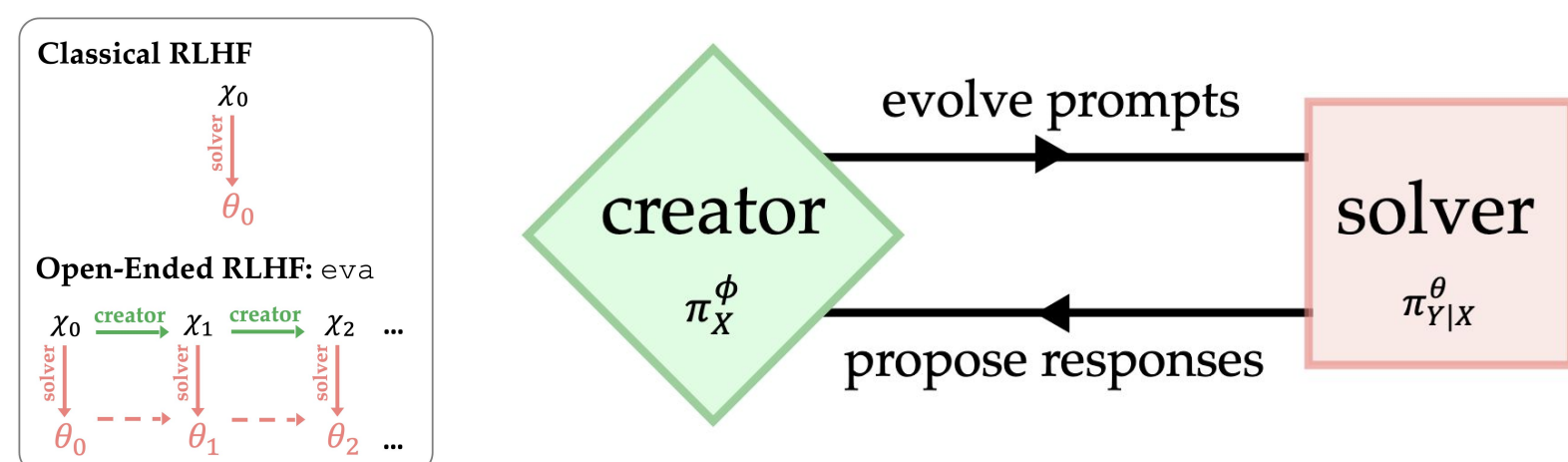
### Motivation

Existing RLHF methods mostly rely on fixed prompt sets, which can hurt model generalization & training efficiency.

### A New Principle

We design an open-ended RLHF principle for LLMs to self-improve by strategically co-evolving prompts & responses.

### A New Training Mechanism



**Solver:** propose preferred responses

**Creator:** propose more informative prompts

### From Open-Ended RL to Minimax Games

Two-Player Game with Solver Regret as the Objective:

$$\text{Regret}(\pi_{\phi}, \pi_{\theta}) = \mathbb{E}_{\mathbf{x} \sim \pi_{\phi}(\cdot)} \left[ \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}(\mathbf{y}|\mathbf{x})} [r(\mathbf{x}, \mathbf{y})] - \mathbb{E}_{\mathbf{y} \sim \pi_{\text{KL}}^*(\mathbf{y}|\mathbf{x})} [r(\mathbf{x}, \mathbf{y})] \right]$$

The Minimax Strategy at Nash Equilibrium:

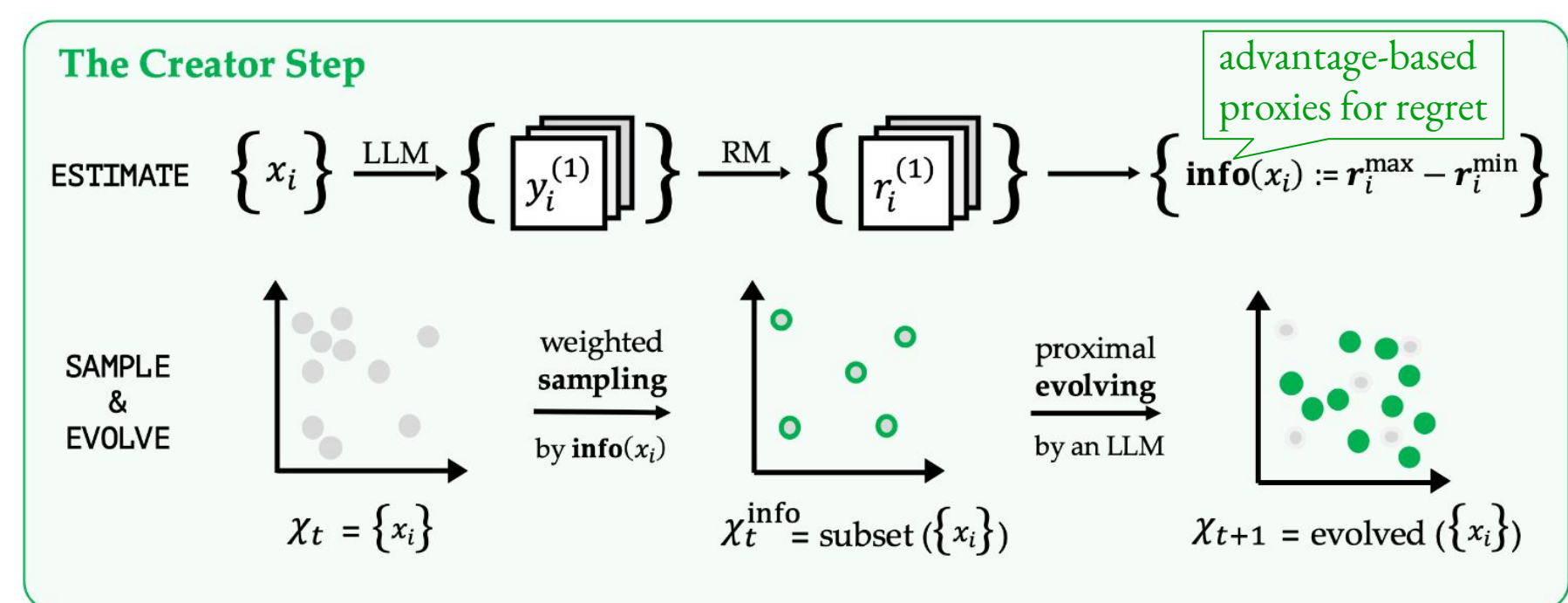
$$\pi_{\mathcal{Y}|\mathcal{X}}^* \in \arg \min_{\pi_{\mathcal{Y}|\mathcal{X}}} \max_{\pi_{\mathcal{X}}} \mathbb{E}_{\mathbf{x} \sim \pi_{\mathcal{X}}} \left[ \text{Regret}(\mathbf{x}, \pi_{\mathcal{Y}|\mathcal{X}}) \right]$$

Regret Minimization by the **Solver**:

Any preference optimization loss, e.g., DPO:

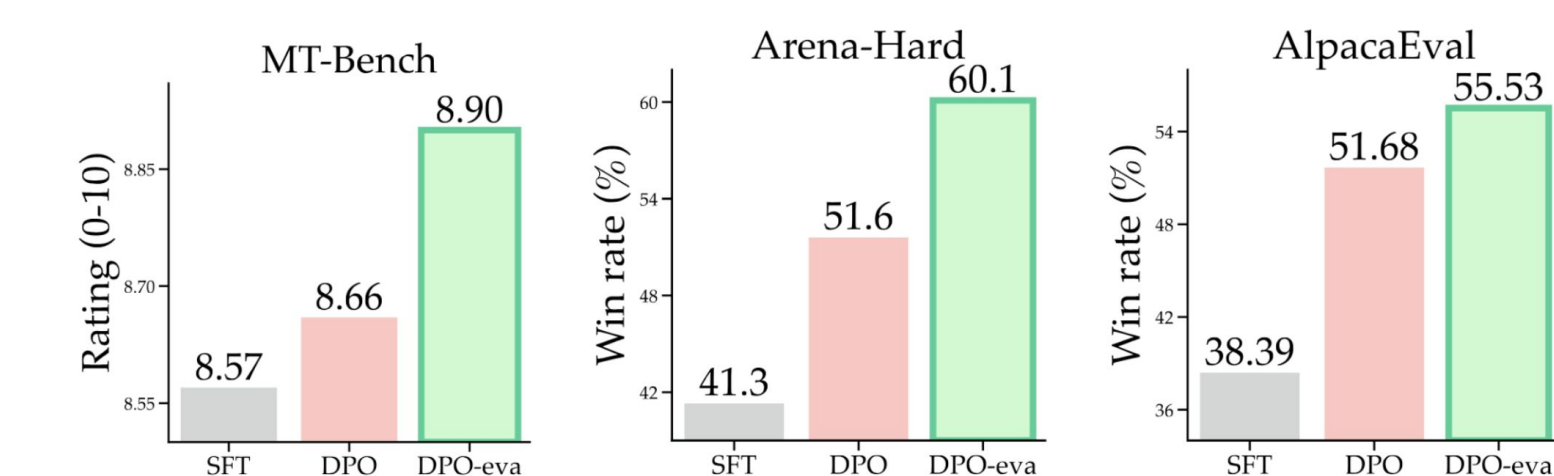
$$\ell_{\beta}(\pi_{\theta}) = -\log \left[ \sigma \left( \beta \cdot \Delta_{\pi_{\theta}; \pi_{\text{ref}}}^{\mathbf{x}} \right) \right] := -\log \left[ \sigma \left( \beta \cdot \log \frac{\pi_{\theta}(\mathbf{y}^+|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^+|\mathbf{x})} - \beta \cdot \log \frac{\pi_{\theta}(\mathbf{y}^-|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}^-|\mathbf{x})} \right) \right]$$

Regret Maximization by the **Creator**:



### The Empirical Results

**eva** brings “universal” alignment gains.



Model Family (→)	GEMMA-2-9B-IT					
Benchmark (→)	Arena-Hard		MT-Bench		AlpacaEval 2.0	
Method (↓) / Metric (→)	WR (%)	avg. score	1 <sup>st</sup> turn	2 <sup>nd</sup> turn	LC-WR (%)	WR (%)
$\theta_0$ : SFT	41.3	8.57	8.81	8.32	47.11	38.39
$\theta_{1 \rightarrow 1}$ : DPO	51.6	8.66	9.01	8.32	55.01	51.68
$\theta_{1 \rightarrow 1}$ : + <b>eva</b>	<b>60.1</b> (+8.5)	<b>8.90</b>	<b>9.04</b>	<b>8.75</b> (+0.43)	<b>55.35</b>	<b>55.53</b>
$\theta_{1 \rightarrow 2}$ : + new human prompts	59.8	8.64	8.88	8.39	55.74	56.15
$\theta_{1 \rightarrow 1}$ : SPPO	55.7	8.62	9.03	8.21	51.58	42.17
$\theta_{1 \rightarrow 1}$ : + <b>eva</b>	<b>58.9</b> (+3.2)	<b>8.78</b>	<b>9.11</b>	<b>8.45</b> (+0.24)	<b>51.86</b>	<b>43.04</b>
$\theta_{1 \rightarrow 2}$ : + new human prompts	57.7	8.64	8.90	8.39	51.78	42.98
$\theta_{1 \rightarrow 1}$ : SimPO	52.3	8.69	9.03	8.35	54.29	52.05
$\theta_{1 \rightarrow 1}$ : + <b>eva</b>	<b>60.7</b> (+8.4)	<b>8.92</b>	<b>9.08</b>	<b>8.77</b> (+0.42)	<b>55.85</b>	<b>55.92</b>
$\theta_{1 \rightarrow 2}$ : + new human prompts	54.6	8.76	9.00	8.52	54.40	55.72
$\theta_{1 \rightarrow 1}$ : ORPO	54.8	8.67	9.04	8.30	52.17	49.50
$\theta_{1 \rightarrow 1}$ : + <b>eva</b>	<b>60.3</b> (+5.5)	<b>8.89</b>	<b>9.07</b>	<b>8.71</b> (+0.41)	<b>54.39</b>	<b>50.88</b>
$\theta_{1 \rightarrow 2}$ : + new human prompts	57.2	8.74	9.01	8.47	54.00	51.21

**eva**’s advantage-based proxy is effective.

Model Family (→)	GEMMA-2-9B-IT					
Benchmark (→)	Arena-Hard		MT-Bench		AlpacaEval 2.0	
Method (↓) / Metric (→)	WR (%)	avg. score	1 <sup>st</sup> turn	2 <sup>nd</sup> turn	LC-WR (%)	WR (%)
$\theta_{1 \rightarrow 1}$ : DPO	51.6	8.66	9.01	8.32	55.01	51.68
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> (uniform)	57.5	8.71	9.02	8.40	53.43	53.98
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> (var( $r$ ))	54.8	8.66	9.13	8.20	54.58	52.55
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> (avg( $r$ ))	58.5	8.76	9.13	8.40	55.01	55.47
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> (1/avg( $r$ ))	56.7	8.79	9.13	8.45	55.04	54.97
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> (1/ $A_{\text{min}}$ )	52.3	8.64	8.96	8.31	53.84	52.92
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> ( $A_{\text{avg}}$ ) (our variant)	60.0	8.85	9.08	8.61	<b>56.01</b>	<b>56.46</b>
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> ( $A_{\text{min}}$ ) (our variant)	60.0	8.86	<b>9.18</b>	8.52	55.96	56.09
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> ( $A_{\text{min}}$ ) (our default)	<b>60.1</b> (+8.5)	<b>8.90</b>	9.04	<b>8.75</b> (+0.43)	55.35	55.53

**eva**’s evolving step is effective.

Benchmark (→)	Arena-Hard		MT-Bench		AlpacaEval 2.0	
Method (↓) / Metric (→)	WR (%)	avg. score	1 <sup>st</sup> turn	2 <sup>nd</sup> turn	LC-WR (%)	WR (%)
$\theta_{1 \rightarrow 1}$ : DPO	51.6	8.66	9.01	8.32	55.01	51.68
$\theta_{1 \rightarrow 1}$ : [no evolve]-greedy	56.1	8.68	8.98	8.38	54.11	53.66
$\theta_{1 \rightarrow 1}$ : [no evolve]-sample	55.3	8.69	9.00	8.38	54.22	54.16
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> -greedy (our variant)	59.5	8.72	9.06	8.36	54.52	55.22
$\theta_{1 \rightarrow 1}$ : + <b>eva</b> -sample (our default)	<b>60.1</b>	<b>8.90</b>	9.04	<b>8.75</b>	<b>55.35</b>	<b>55.53</b>

### The Practical Algorithm

**eva** can be easily plugged into any RLHF pipeline.

#### Algorithm 1 **eva**: Evolving Alignment via Asymmetric Self-Play

**Input:** initial policy  $\pi_{\theta_0}$ , initial prompt set  $\mathcal{X}_0$

1: **for** iteration  $t = 1, 2, \dots$  **do**

    ▽ /\* **creator step** \*/

2: *estimate informativeness:*  $\mathcal{X}_{t-1} \leftarrow \{(\mathbf{x}_i, \text{info}(\mathbf{x}_i)) \mid \mathbf{x}_i \in \mathcal{X}_{t-1}\}$

*sample subset:*  $\mathcal{X}_{t-1}^{\text{info}} \leftarrow \text{sample}(\mathcal{X}_{t-1})$

*self-evolve prompts:*  $\mathcal{X}_t \leftarrow \text{evolve}(\mathcal{X}_{t-1}^{\text{info}})$

    ▽ /\* **solver step** \*/

3: *self-generate responses:*  $\forall \mathbf{x}_i \in \mathcal{X}_t$ , generate  $\{\mathbf{y}_i^{(j)}\} \sim \pi_{\theta_{t-1}}(\cdot \mid \mathbf{x}_i)$

*annotate rewards:*  $\mathcal{X}'_t \leftarrow \mathcal{X}_t \cup \{(\mathbf{y}_i^{(j)}, r_i^{(j)})\}$

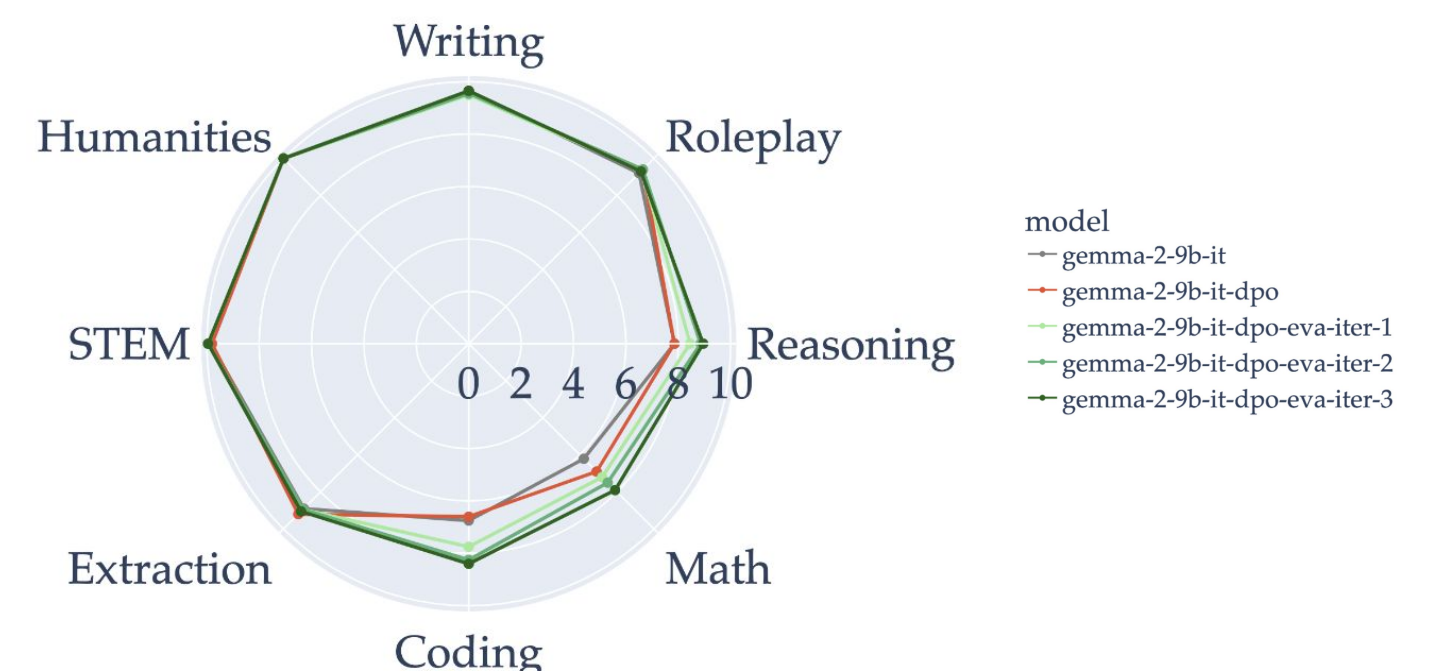
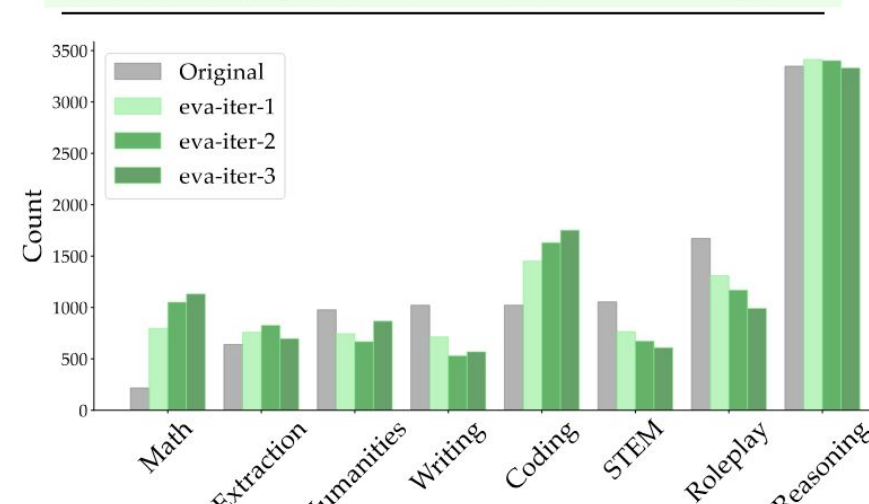
*preference optimization:*  $\theta_t \leftarrow \theta_{t-1} - \eta \nabla_{\theta} \mathcal{L}_{\mathcal{X}'_t}(\theta)$

4: **end for**

5: **return** final solver policy  $\pi_{\theta_T}$

**eva** leads to continual self-improving – the infinite games!

Prompt Set (↓) / Metric (→)	Complexity (1-5)	Quality (1-5)
UltraFeedback (seed)	2.90	3.18
UltraFeedback- <b>eva</b> -Iter-1	3.84	3.59
UltraFeedback- <b>eva</b> -Iter-2	3.92	3.63
UltraFeedback- <b>eva</b> -Iter-3	3.98	3.73



### General Takeaways

1. RLHF can be made open-ended.
2. Reward advantage is effective in prompt selection.



paper link: arXiv:2411.00062