

# Learning to Reason and Align via Self-Play

*In pursuit of Superhuman Intelligence*

Ziyu Ye

Dec 06, 2024



THE UNIVERSITY OF  
CHICAGO

# Research Goal: *In Pursuit of Superintelligence*

~ Self-Play

Design **scalable methods** for intelligence to perform complex **sequential decision making** to **achieve goals** in the open world.

~ Reasoning


~ Alignment

# Agenda for Today's Talk

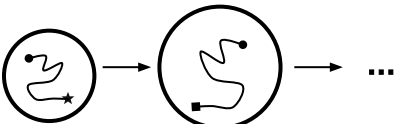


## Intro: Learning as an Game

### Conventional Training as a Finite Game

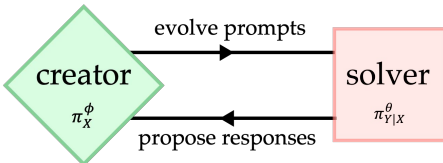


### New, Scalable Training as an Infinite Game





## Self-Play to Align



**Classical RLHF**

$x_0$

$\theta_0$

**The Creator Step**

ESTIMATE  $\{x_t\}$  LLM  $\left\{ \left[ \begin{smallmatrix} \cdot \\ \cdot \\ \cdot \end{smallmatrix} \right] \right\}$  RM  $\left\{ \left[ \begin{smallmatrix} \cdot \\ \cdot \\ \cdot \end{smallmatrix} \right] \right\}$   $\left\{ \text{info}(x_t) = r_t^{\text{info}} - r_t^{\text{rm}} \right\}$

SAMPLE & EVOLVE  $x_t = \{x_t\}$   $\xrightarrow{\text{weighted sampling by info}(x_t)}$   $x_t^{\text{info}} = \text{subset}(\{x_t\})$   $\xrightarrow{\text{proximal encoding by an LLM}}$   $x_{t+1} = \text{evolved}(\{x_t\})$

---

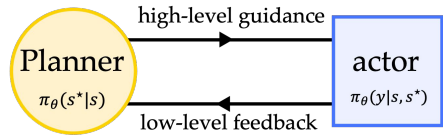
**Algorithm 1 eva: Evolving Alignment via Asymmetric Self-Play**

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

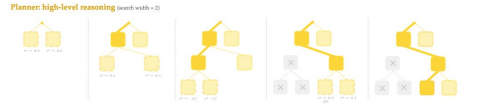
- 1: **for** iteration  $t = 1, 2, \dots$  **do**
  - $\nabla$  /\* creator step \*/
  - 2: estimate informativeness:  $\mathcal{X}_{t-1} \leftarrow \{(x_t, \text{info}(x_t)) \mid x_t \in \mathcal{X}_{t-1}\}$
  - sample subset:  $\mathcal{X}_{t-1}^{\text{info}} \leftarrow \text{sample}(\mathcal{X}_{t-1})$
  - self-evolve prompts:  $\mathcal{X}_t \leftarrow \text{evolve}(\mathcal{X}_{t-1}^{\text{info}})$
  - $\nabla$  /\* solver step \*/
  - 3: self-generate responses:  $\forall x_t \in \mathcal{X}_t$ , generate  $\{y_t^{(j)}\} \sim \pi_{\theta_{t-1}}(\cdot \mid x_t)$
  - annotate rewards:  $\mathcal{X}_t^r \leftarrow \mathcal{X}_t \cup \{(y_t^{(j)}, r_t^{(j)})\}$
  - preference optimization:  $\theta_t \leftarrow \theta_{t-1} - \eta \nabla_{\theta} \mathcal{L}_{\mathcal{X}_t^r}(\theta)$
- 4: **end for**
- 5: **return** final solver policy  $\pi_{\theta_T}$




## Self-Play to Reason



**Planner: high-level reasoning** (search width = 2)



**Actor: low-level reasoning** (search width = 2)



---

**Algorithm 1 RiR - A Unified Reasoning Mechanism with Decomposing and Search**

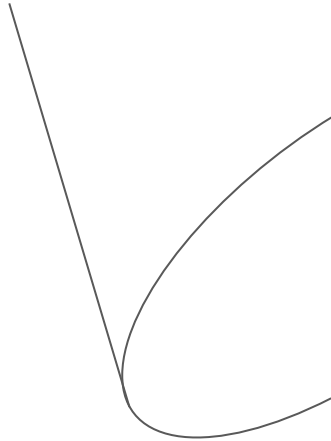
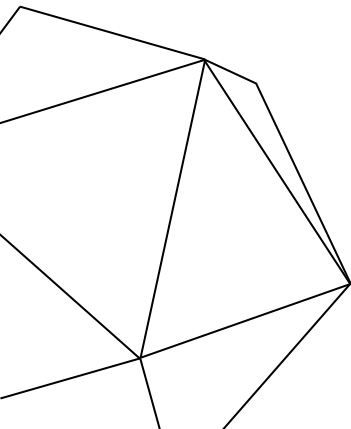
**Input:** problem statement  $q$ , a language model w/ parameter  $\theta$

- 1:  $\text{tree} \leftarrow \text{Tree}(\theta, q)$
- 2: **repeat**
- 3:  $s_t^i \leftarrow \text{tree.policy}()$  ▷ /\* planner \*/
- 4:  $\text{tree} \leftarrow \text{Tree}(\theta, s_t^i)$
- 5: **repeat**
- 6:  $y_t \leftarrow \text{tree.policy}()$  ▷ /\* goal-driven actor \*/
- 7: **until** STOP.LOW
- 8:  $\{\text{tree}, \text{tree}\}.\text{update}()$  ▷ /\* joint update \*/
- 9: **until** STOP.HIGH
- 10: **return**  $\text{tree}.\text{solution}$



# Intro: Learning as an Infinite Game

*“a paradigm shift for training large models”*



# The Vision: “*Universality of Computation*”

## 6. *The universal computing machine.*

It is possible to invent a single machine which can be used to compute any computable sequence. If this machine  $\mathcal{U}$  is supplied with a tape on the beginning of which is written the S.D of some computing machine  $\mathcal{M}$ , then  $\mathcal{U}$  will compute the same sequence as  $\mathcal{M}$ .

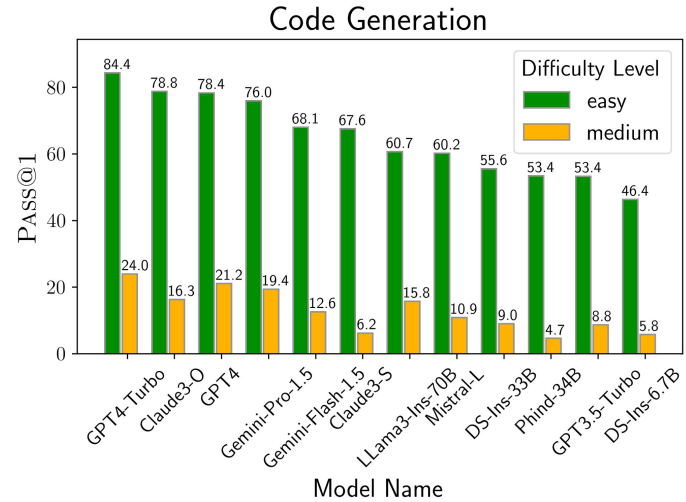
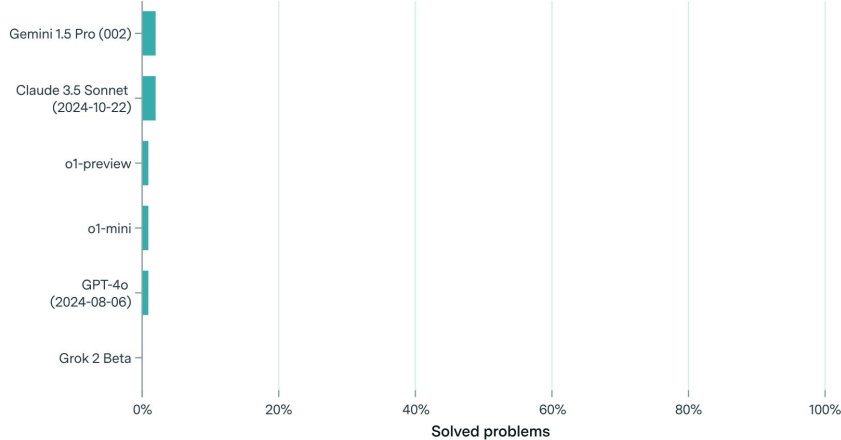
– Alan M. Turing, 1936

**“what a *human* can think or know”**

**=**

**“what a *machine* can compute”**

# The Challenge: Gaps in Achieving Human-Level Performance



● **< 2% accuracy** on challenging contemporary mathematics problems on [FrontierMath](#).

● **>3.5x performance drop** as coding problems get harder on [LiveCodeBench](#).

.....  
**What may go wrong in conventional ways of training AI models?**  
.....

## The Challenge: Scaling Law is Hitting the Wall? 🤖

“Ilya Sutskever, co-founder of AI labs Safe Superintelligence (SSI) and OpenAI, told Reuters recently that results from scaling up pre-training - **the phase of training an AI model** that use s a vast amount of unlabeled data to understand language patterns and structures - **have plateaued.**”

“The 2010s were the age of scaling, now we’re back in the age of wonder and discovery once again. Everyone is looking for **the next thing,**” Sutskever said. “**Scaling the right thing matters more now than ever.**”

– Ilya Sutskever with Reuters, Nov 2024

.....

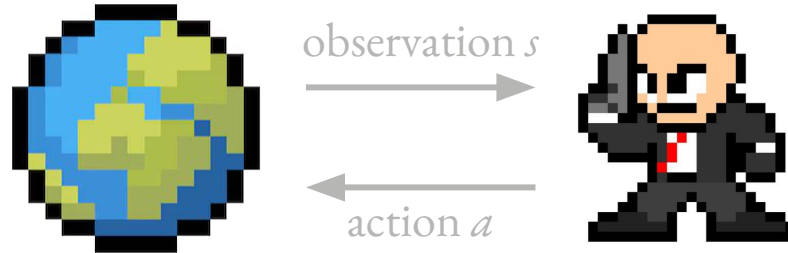
**What are the next right things to scale?**

.....

# Conventional Way of Agent Training

a given “world”

a single “policy”



- **Intelligence:** Agents that are able to learn to make decisions to achieve goals.  
▲ ▲ ■
- ▲ **Reasoning:** The process of making decisions by evaluating information.
- **Alignment:** The process of achieving goals by reward maximization.



# Myth 1: Learning is Purely Solving (under a given world)

a given world



Conventional way: 🤖

Design agents that **find solutions** in a fixed environment, then **stop learning**.

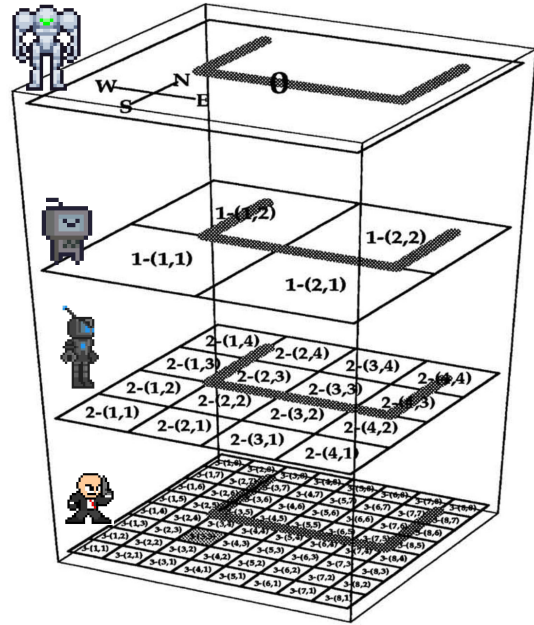
the open-ended worlds



Better way:

Design agents that **create new tasks/environments**, then **continuously learn to self-improve**.

## Myth 2: Reasoning is Step-by-Step (by a single policy)



**Better way:**

Learn policies with a **hierarchy of abstract models**, and roll out at different levels for optimization.

**Conventional way:** 🤖

Learn a policy that operates under a **one-step model**, and roll it out (with tree search) in training.

Fig 1. A world can be divided at different levels in certain hierarchy ([Dayan and Hinton, 1992](#)).

# A New, Scalable Training Paradigm

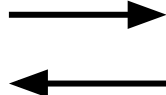
“There are at least 2 kinds of games. One could be called finite; the other infinite.”

- A **finite game** is played for the purpose of winning.
- An **infinite game** is for the purpose of continuing the play.

– James P. Carse, 2011

## Conventional Training as a Finite Game

a given “world”

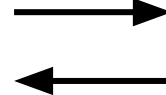


a single “policy”



## New, Scalable Training as an Infinite Game

open-ended worlds



hierarchical policies



# A New, Scalable Training Paradigm

“There are at least 2 kinds of games. One could be called finite; the other infinite.”

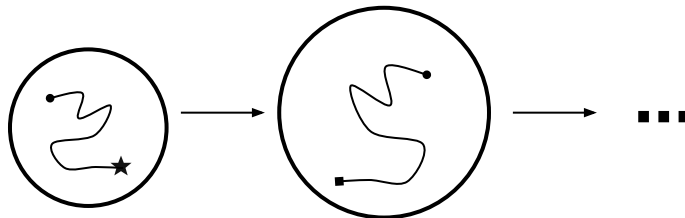
- A **finite game** is played for the purpose of winning.
- An **infinite game** is for the purpose of continuing the play.

– James P. Carse, 2011

**Conventional Training  
as a Finite Game**



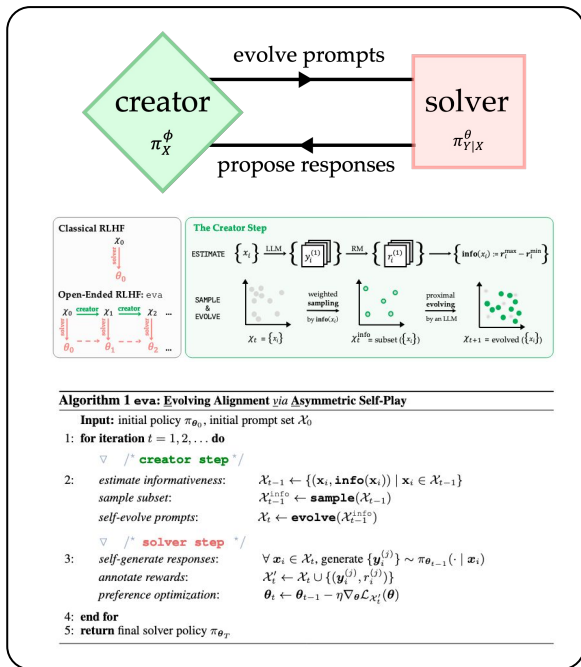
**New, Scalable Training  
as an Infinite Game**



# Recap on Agenda



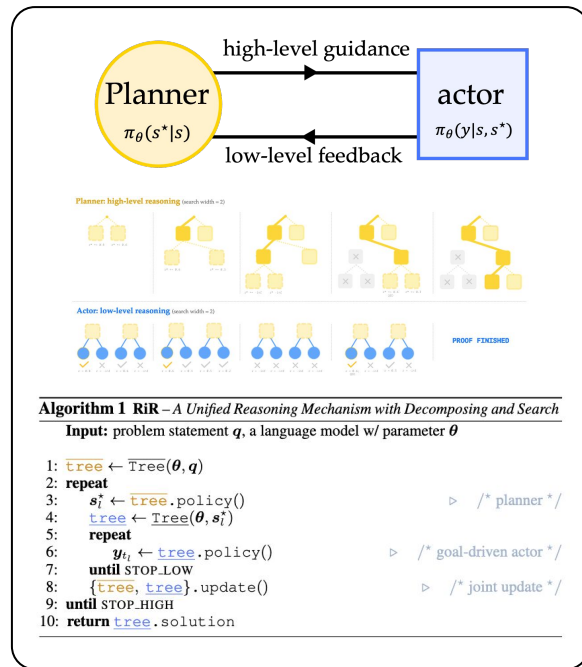
## Self-Play to Align



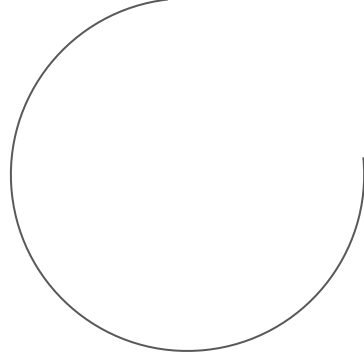
*Solving Myth 1:*  
Going Beyond Static World



## Self-Play to Reason



*Solving Myth 2:*  
Hierarchical Planning



# Self-Play to Align: The **Creator-Solver** Game

*“Scalable language model training beyond human prompts.”*

Gratitude to every wonderful co-author of this project ([link](#)):  
Rishabh Agarwal, Tianqi Liu, Rishabh Joshi, Sarmishta Velury, Quoc V. Le, Qijun Tan, Yuan Liu.

# TL; DR

We identify **learnable, worth-learning prompts** by **reward signals**, then **evolve new prompts** for open-ended continual RLHF training.

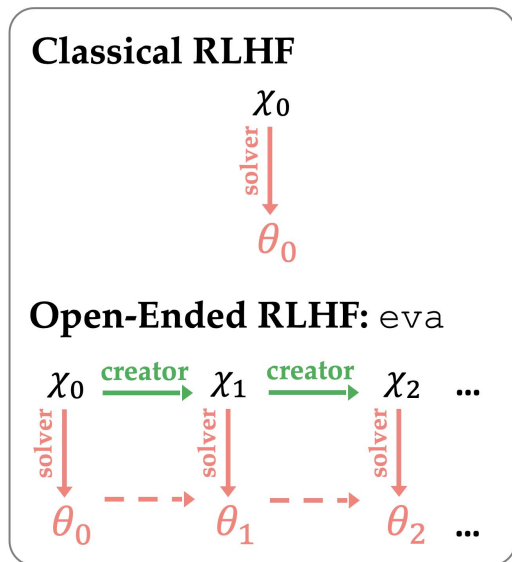


Fig 1. RLHF needs a paradigm shift!

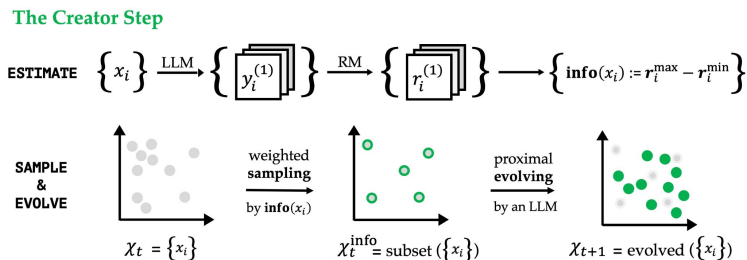


Fig 2. The **easy-to-implement pipeline** of **eva** for open-ended RLHF.

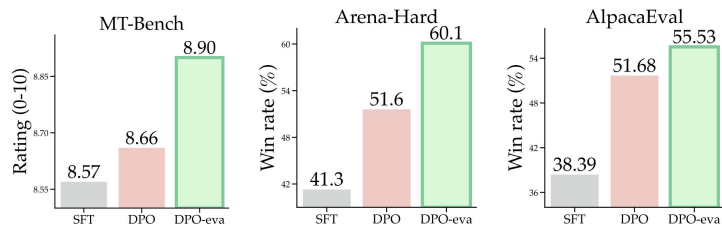


Fig 3. **eva** brings strong **alignment gain**.

# Artificial Intelligence May Be Bottlenecked by Static Data

Reference: [Villalobos et al., 2024](#)

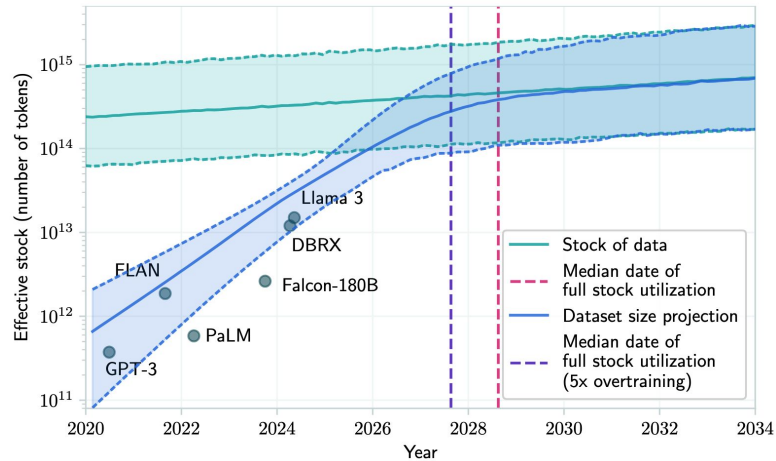


Fig 4. The **scale**, **quality** and **growth** of human knowledge is bottlenecked.

See also: [lmsys-hard](#)

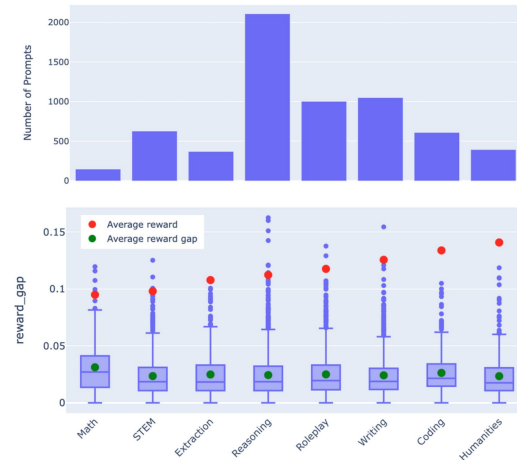


Fig 5. The **imbalance distribution** of static training data.

Can language models identify and self-create **new, learnable, and worth-learning** tasks, to self-improve to generalize better for alignment?



# Classical RLHF

**Alignment by RLHF.** Classical RLHF (Ouyang et al., 2022) optimizes on a fixed distribution  $\mathcal{D}$ :

$$\max_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ r(\mathbf{x}, \mathbf{y}) \right] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \beta \cdot \mathbb{D}_{\text{KL}} \left[ \pi_{\theta}(\mathbf{y} | \mathbf{x}) \parallel \pi_{\text{SFT}}(\mathbf{y} | \mathbf{x}) \right] \right], \quad (1)$$

where  $\mathbf{x}$  and  $\mathbf{y}$  denote the prompts and responses, and  $r(\cdot, \cdot)$  is the reward function.

# Our Perspective: RLHF Should Be Made Open-Ended

**Definition 1 (Open-Ended RLHF)** We define evolving alignment as the open-ended joint optimization on the prompt and response policy for alignment w.r.t the joint reference policy:

$$\max_{\phi, \theta} \mathbb{E}_{\mathbf{x} \sim \pi_{\phi}(\cdot), \mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ r(\mathbf{x}, \mathbf{y}) \right] - \beta \cdot \mathbb{D}_{\text{KL}} \left[ \pi_{\phi, \theta}(\mathbf{x}, \mathbf{y}) \parallel \pi_{\text{ref}}(\mathbf{x}, \mathbf{y}) \right], \quad (7)$$

where  $\pi_{\phi, \theta}(\mathbf{x}, \mathbf{y}) := \pi_{\phi}(\mathbf{x}) \cdot \pi_{\theta}(\mathbf{y} | \mathbf{x})$  and  $\pi_{\text{ref}}(\mathbf{x}, \mathbf{y}) := p_{\text{ref}}(\mathbf{x}) \cdot \pi_{\text{SFT}}(\mathbf{y} | \mathbf{x})^a$ .

<sup>a</sup>This generalizes classical RLHF (Eq. 1). One may extend the above and rewrite coefficients to be:

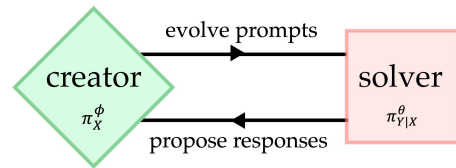
$$\max_{\phi, \theta} \mathbb{E}_{\mathbf{x} \sim \pi_{\phi}(\cdot)} \left[ \mathbb{E}_{\mathbf{y} \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ r(\mathbf{x}, \mathbf{y}) \right] - \beta_1 \mathbb{D}_{\text{KL}} \left[ \pi_{\theta}(\mathbf{y} | \mathbf{x}) \parallel \pi_{\text{SFT}}(\mathbf{y} | \mathbf{x}) \right] \right] - \beta_2 \mathbb{D}_{\text{KL}} \left[ \pi_{\phi}(\mathbf{x}) \parallel p_{\text{ref}}(\mathbf{x}) \right]. \quad (8)$$

However, directly optimizing this can be *intractable* or *unstable*... 🤔

# Our Method: Open-Ended RLHF via Creator-Solver Games

## ● How? Optimization by Asymmetric Games

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



## ● What? The Regret of the Solver's Policy

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

## ● Why? The Minimax Regret Strategy at the Nash Equilibrium

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

However, w/o access to the true  $\pi^*$ , we must **approximate** this regret... 🤔

# Our Method: Open-Ended RLHF via Creator-Solver Games

- How to approximate the **regret**? Simply use the **stochastic policy**...

Sample  $N$  times from the policy,  
then choosing the *reward gap* between *the best* and *the baseline*.

$$|\text{Regret}(\mathbf{x}, \pi_{\theta})| \leftarrow \text{info}_{\theta}(\mathbf{x}) := r(\mathbf{x}, \mathbf{y}_{+}) - r(\mathbf{x}, \mathbf{y}_{\text{baseline}}),$$

$$\mathbf{y}_{+} := \arg \max_{\mathbf{y}_i} r(\mathbf{x}, \mathbf{y}),$$

$$\mathbf{y}_{\text{baseline}} := \arg \min_{\mathbf{y}_i} r(\mathbf{x}, \mathbf{y}) \text{ or } \mathbf{y}_{\text{baseline}} := \text{avg}_{\mathbf{y}_i} r(\mathbf{x}, \mathbf{y})$$

- Other intuitive interpretations of **eva** 🧙

Learning potential.

*eva* picks the prompts that are learnable but not learned yet.

Auto-curricula for the solver player.

The optimal strategy of the creator is to create prompts just beyond solvers' current capability .

Worst-case guarantee.

The minimax objective incentivizes the solver to perform well in all cases.

Auto-curricula inherent to **contrastive learning**.

*eva* prioritizes prompts with lower contrastive loss by design, thus accelerating learning.

# The eva Algorithm

---

**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**

$\nabla$  /\* **creator step** \*/

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

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

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

$\nabla$  /\* **solver step** \*/

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

*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  $\pi_{\theta_T}$

---

Fig 6. **eva** requires only a **creator module addition** to make current RLHF pipeline **open-ended**.

# The Creator Step: Estimate, Sample then Evolve

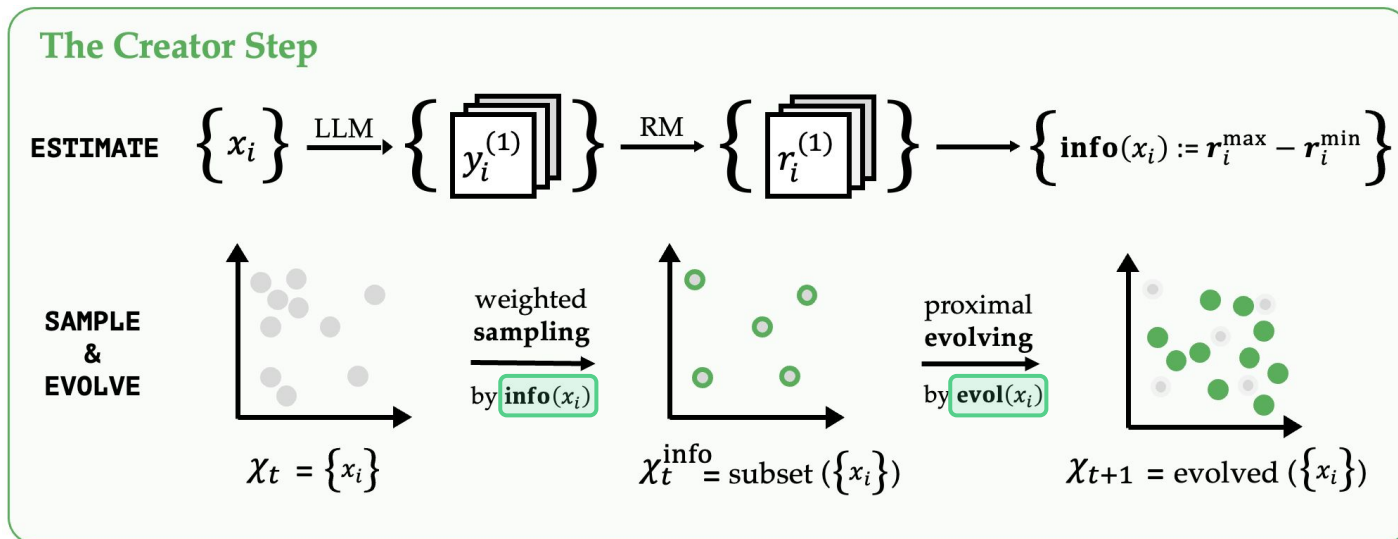


Fig 7. **eva** currently uses the **estimate, sample then evolve** procedure for the creator. Here, **info**( $\cdot$ ) is the reward gap, and **evol**( $\cdot$ ) can be any prompt creation method.

# Example Evolving Method: evol ( · )

We use **EvolInstruct** ([Can et al., 2023](#)) for **in-depth evolving** and **in-breadth-evolving**.

## Initial prompt ↓

If a man smokes 1000 cigarettes a day, why is he getting healthier?

## Evolved #1 ↓ (*in-depth* evolving)

Elaborate on the seemingly paradoxical situation where an individual consumes 1000 cigarettes daily yet exhibits signs of improving health, delineating the factors that could underlie such an unexpected outcome.

## Evolved #2 ↓ (*in-breadth* evolving)

Discuss the conundrum of a person drinking a gallon of caffeinated coffee every hour but displaying unusually deep and restful sleep patterns, exploring possible explanations for this unusual phenomenon.

```
MUTATION_TEMPLATES = {  
    # -----  
    # >>>>>> In-depth evolving <<<<<<<<  
    # -----  
    "CONSTRAINTS": prompt.format(  
        "Add one more constraints into '#The Given Prompt#'"  
    ),  
    "DEEPENING": prompt.format(  
        "If #The Given Prompt# contains inquiries about certain issues,  
        the depth and breadth of the inquiry can be increased."  
    ),  
    "CONCRETIZING": prompt.format(  
        "Please replace general concepts with more specific concepts."  
    ),  
    "INCREASED_REASONING_STEPS": prompt.format(  
        "If #The Given Prompt# can be solved with just a few simple  
        thinking processes, you can rewrite it to explicitly request  
        multiple-step reasoning."  
    ),  
    # -----  
    # >>>>>> In-breadth evolving <<<<<<<<  
    # -----  
    "BREADTH": prompt.format(  
        "By inspiration from #The Given Prompt#, create a new prompt.  
        This new prompt should belong the the same domain as it, but be  
        even more rare. The length and complexity should be similar. The  
        #Created Prompt# must be reasonable and must be understood and  
        responded by humans."  
    )  
}
```

## Results: Remarkable Gains on Hard Benchmarks\*!

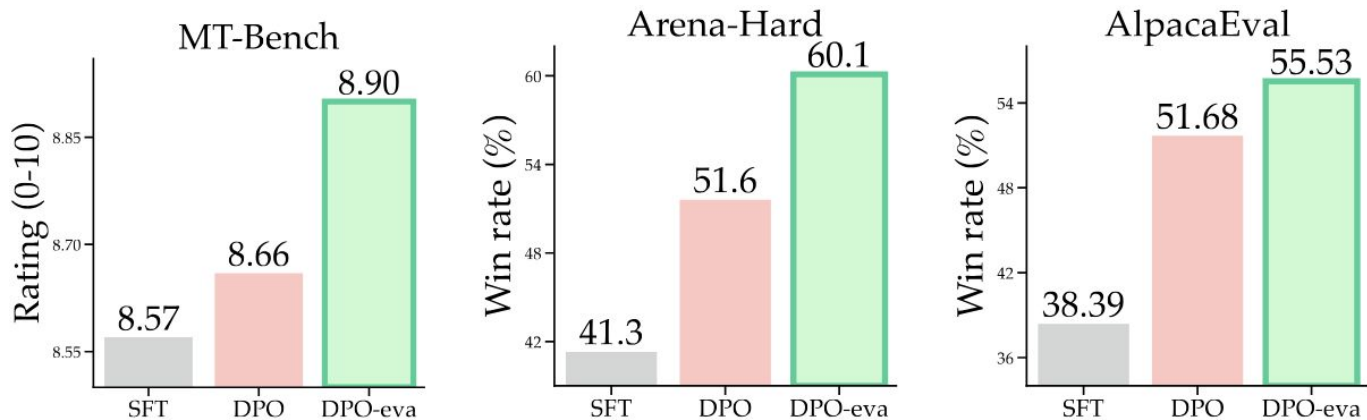


Figure. **eva** achieves concrete performance gain especially on **hard benchmarks**, without relying on any additional human prompts.



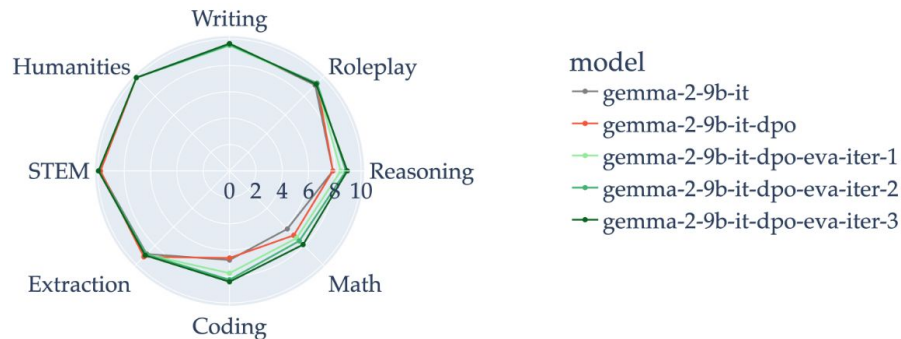
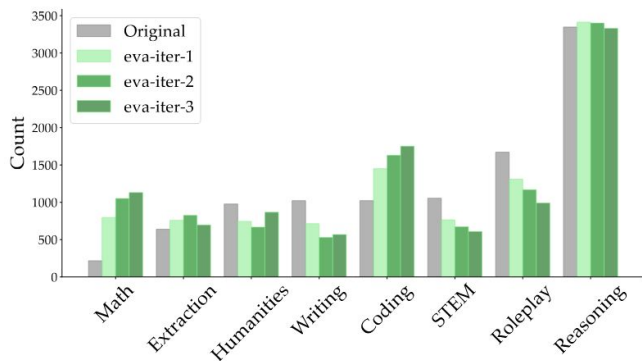
# Results: Remarkable Gains on Hard Benchmarks\*!

Model Family ( $\rightarrow$ )	GEMMA-2-9B-IT					
	Benchmark ( $\rightarrow$ )	Arena-Hard		MT-Bench		AlpacaEval 2.0
Method ( $\downarrow$ ) / Metric ( $\rightarrow$ )	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
$\theta_{0 \rightarrow 1}$ : DPO	51.6	8.66	9.01	8.32	55.01	51.68
$\theta_{1 \rightarrow \bar{1}}$ : <b>+ eva</b>	<b>60.1</b> (+8.5)	<b>8.90</b>	<b>9.04</b>	<b>8.75</b> (+0.43)	55.35	55.53
$\theta_{1 \rightarrow 2}$ : +new human prompts	59.8	8.64	8.88	8.39	<b>55.74</b>	<b>56.15</b>
$\theta_{0 \rightarrow 1}$ : SPPO	55.7	8.62	9.03	8.21	51.58	42.17
$\theta_{1 \rightarrow \bar{1}}$ : <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
$\theta_{0 \rightarrow 1}$ : SimPO	52.3	8.69	9.03	8.35	54.29	52.05
$\theta_{1 \rightarrow \bar{1}}$ : <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
$\theta_{0 \rightarrow 1}$ : ORPO	54.8	8.67	9.04	8.30	52.17	49.50
$\theta_{1 \rightarrow \bar{1}}$ : <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	<b>51.21</b>

Table 1: **Main results.** Our **eva** achieves notable alignment gains and can surpass human prompts on major benchmarks across a variety of representative direct preference optimization algorithms.

\* All experiments are conducted with external open-source frameworks and models on HuggingFace. We use 10K prompts from [UltraFeedback](#) for training, and use [ArmoRM-8B](#) as the default reward model.

# Additional Results – **eva** creates meaningful curriculum.



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	<b>3.98</b>	<b>3.73</b>

# Ablation #1 – **eva**’s minimax design outperforms alternatives.

Metric	info( $\mathbf{x}$ )	Related Interpretations
$A_{\min}^*$ : worst-case optimal advantage	$ \min_{\mathbf{y}} r(\mathbf{x}, \mathbf{y}) - \max_{\mathbf{y}'} r(\mathbf{x}, \mathbf{y}') $	minimax regret (Savage, 1951)
$A_{\text{avg}}^*$ : average optimal advantage	$ \frac{1}{N} \sum_{\mathbf{y}} r(\mathbf{x}, \mathbf{y}) - \max_{\mathbf{y}'} r(\mathbf{x}, \mathbf{y}') $	Bayesian regret (Banos, 1968)
$A_{\text{dts}}^*$ : dueling optimal advantage	$ \max_{\mathbf{y} \neq \mathbf{y}'} r(\mathbf{x}, \mathbf{y}) - \max_{\mathbf{y}'} r(\mathbf{x}, \mathbf{y}') $	dueling regret (Wu and Liu, 2016)

Table 2: The reward-advantage-based metrics that serve as the informativeness proxies for prompts.

Model Family ( $\rightarrow$ )	GEMMA-2-9B-IT					
Benchmark ( $\rightarrow$ )	Arena-Hard	MT-Bench		AlpacaEval 2.0		
Method ( $\downarrow$ ) / Metric ( $\rightarrow$ )	WR (%)	avg. score	1 <sup>st</sup> turn	2 <sup>nd</sup> turn	LC-WR (%)	WR (%)
$\theta_{0 \rightarrow 1}$ : DPO	51.6	8.66	9.01	8.32	55.01	51.68
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> (uniform)	57.5	8.71	9.02	8.40	53.43	53.98
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> (var( $\mathbf{r}$ ))	54.8	8.66	9.13	8.20	54.58	52.55
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> (avg( $\mathbf{r}$ ))	58.5	8.76	9.13	8.40	55.01	55.47
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> (1/avg( $\mathbf{r}$ ))	56.7	8.79	9.13	8.45	55.04	54.97
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> (1/ $A_{\min}^*$ )	52.3	8.64	8.96	8.31	53.84	52.92
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> ( $A_{\text{avg}}^*$ ) (our variant)	60.0	8.85	9.08	8.61	<b>56.01</b>	<b>56.46</b>
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> ( $A_{\text{dts}}^*$ ) (our variant)	60.0	8.86	<b>9.18</b>	8.52	55.96	56.09
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> ( $A_{\min}^*$ ) (our default)	<b>60.1</b> (+8.5)	<b>8.90</b>	9.04	<b>8.75</b> (+0.43)	55.35	55.53

Table 3: **Choice of informativeness metric.** Our informativeness metric by *advantage* achieves the best performances, comparing with others as the weight to sample prompts to evolve by the creator.

## Ablation #2 – **eva**’s design of evolving is meaningful.

Benchmark (→)	Arena-Hard		MT-Bench		AlpacaEval 2.0	
	WR (%)	avg. score	1 <sup>st</sup> turn	2 <sup>nd</sup> turn	LC-WR (%)	WR (%)
$\theta_{0 \rightarrow 1}$ : DPO	51.6	8.66	9.01	8.32	55.01	51.68
$\theta_{1 \rightarrow \tilde{1}}$ : [no evolve]-greedy	56.1	8.68	8.98	8.38	54.11	53.66
$\theta_{1 \rightarrow \tilde{1}}$ : [no evolve]-sample	55.3	8.69	9.00	8.38	54.22	54.16
$\theta_{1 \rightarrow \tilde{1}}$ : + <b>eva</b> -greedy (our variant)	59.5	8.72	9.06	8.36	54.52	55.22
$\theta_{1 \rightarrow \tilde{1}}$ : + <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>

Table 4: **Effect of evolving.** The blue are those training w/ only the informative subset and w/o evolving); we denote `-sample` for the default weighted sampling procedure in Algo 1, while using `-greedy` for the variant from the classical active data selection procedure (*cf.*, a recent work (Muldrew et al., 2024) and a pre-LLM work (Kawaguchi and Lu, 2020)), which selects data by a high-to-low ranking via the metric greedily. We show evolving brings a remarkable alignment gain (the red v.s. the blue); and as we evolve, sampling is more robust than being greedy (*cf.*, Russo et al. (2018)).

## Ablation #3 – **eva** scales with better reward models.

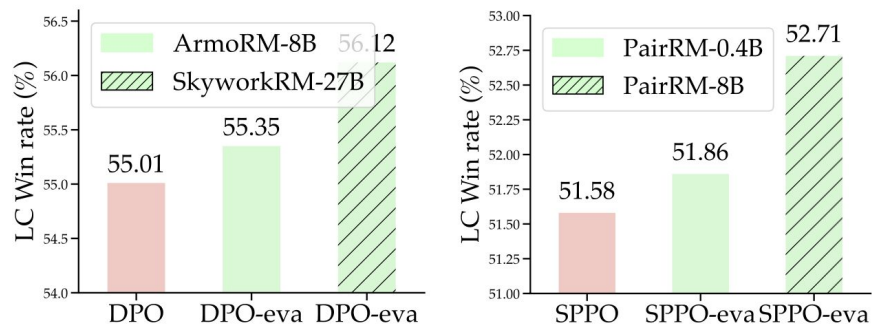


Figure. **eva** scales with better reward models.

# Ablation #4 – **eva** is robust in continual training.

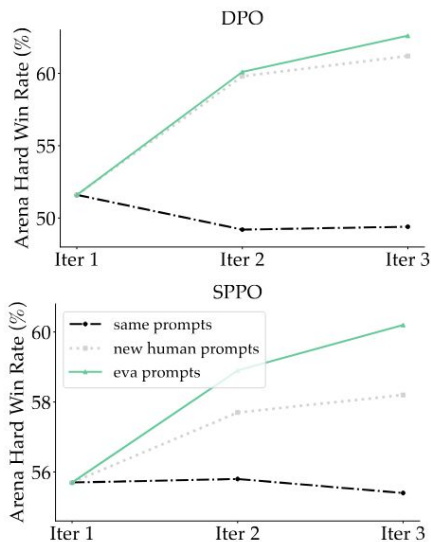


Figure 5: **Continual training.** **eva** stays robust w/ more iterations in incremental training.

Model Family (→)	GEMMA-2-9B-IT	
Benchmark (→)	Arena-Hard	
Method (↓) / Metric (→)	WR (%)	avg. len
$\theta_0$ : SFT	41.3	544
$\theta_{0 \rightarrow 1}$ : DPO (10k)	51.6	651
$\theta_{1 \rightarrow 2}$ : DPO (10k)	59.8	718
$\theta_{2 \rightarrow 3}$ : DPO (10k)	61.2	802
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> (10k)	60.1	733
$\theta_{\bar{1} \rightarrow \bar{2}}$ : + <b>eva</b> (10k)	62.0	787
$\theta_{\bar{2} \rightarrow \bar{3}}$ : + <b>eva</b> (10k)	62.2	774

Model Family (→)	GEMMA-2-9B-IT	
Benchmark (→)	Arena-Hard	
Method (↓) / Metric (→)	WR (%)	avg. len
$\theta_0$ : SFT	41.3	544
$\theta_{0 \rightarrow 1}$ : DPO (20k)	53.2	625
$\theta_{1 \rightarrow 2}$ : DPO (20k)	47.0	601
$\theta_{2 \rightarrow 3}$ : DPO (20k)	46.8	564
$\theta_{1 \rightarrow \bar{1}}$ : + <b>eva</b> (20k)	59.5	826
$\theta_{\bar{1} \rightarrow \bar{2}}$ : + <b>eva</b> (20k)	60.0	817
$\theta_{\bar{2} \rightarrow \bar{3}}$ : + <b>eva</b> (20k)	61.4	791

# Takeaways

**eva** is a new, simple framework for aligning language models via a creator-solver game.

RLHF can be made **open ended**:

- **self-evolving joint data distributions** (with synthesized prompts) bring significant gains.
- **reward advantage** acts as an effective metric for prompt selection.



[work-in-progress]

# Self-Play to Reason: Decompose + Search is All You Need

*“Better than state-of-the-art and 3x faster for neural theorem proving.”*

Gratitude to every wonderful co-author of this project ([link](#)):

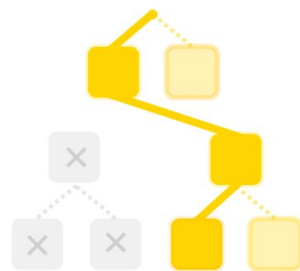
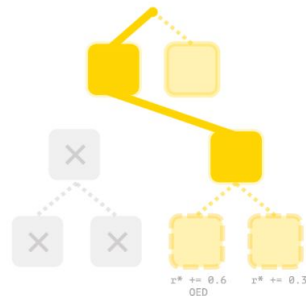
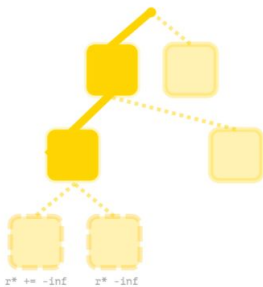
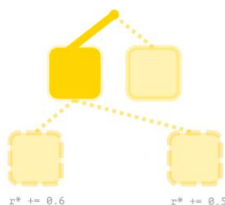
Jiacheng Chen, Jonathan Light, Yifei Wang, Jiankai Sun, Mac Schwager, Guohao Li, Philip Torr, Yuxin Chen, Kaiyu Yang, Yisong Yue, Ziniu Hu.



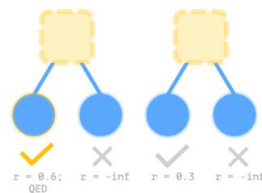
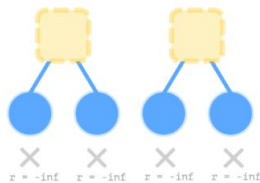
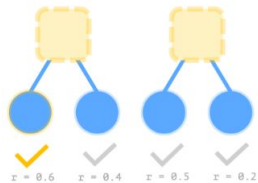
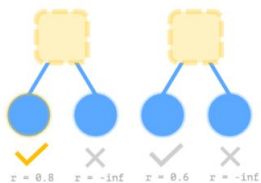
# TL; DR

We unify **decomposing** and **search** for **better and faster reasoning**.

**Planner: high-level reasoning** (search width = 2)



**Actor: low-level reasoning** (search width = 2)



**PROOF FINISHED**

□ / □ / □ / □ : candidate/queued/chosen/pruned target goal    ● : action (i.e., tactic or proofstep)    X / ✓ / ✓ : Lean invalidated / validated / validated, and prioritized by value

# Preliminaries

```
theorem (p q: Prop) : p v q  q v p := by
  intro h
  cases h with
  | inl hp => apply Or.inr; exact hp
  | inr hq => apply Or.inl; exact hq
-- goal s0: (p q: Prop) p v q  q v p
-- goal s1: (p q: Prop) (h: p v q)  q v p
-- goal s2: (p q: Prop) (hp: p)  q v p
-- goal s3: (p q: Prop) (hq: q)  q v p
-- goal s4: None
```

**Neural theorem proving.** A neural network parameterized by  $\theta$  can act as a policy that samples single tactic  $\mathbf{y}_{t+1} \sim \pi_{\theta}(\cdot \mid \mathbf{s}_t)$  at step  $t$ . The objective is to find the optimal trajectory which leads to solved for each statement  $\mathbf{q}$ , that is to find a sequence of tactics  $\mathbf{y}_1, \dots, \mathbf{y}_T$  such that:

$$\mathbf{s}_0 \xrightarrow{\mathbf{y}_1} \mathbf{s}_1 \xrightarrow{\mathbf{y}_2} \mathbf{s}_2 \xrightarrow{\mathbf{y}_3} \dots \xrightarrow{\mathbf{y}_T} \mathbf{s}_T.$$

**Classical training method.**

- > **Input:**  $\{\$current\_goal \mathbf{s}\}$
- > **Output:**  $\{\$proofstep \mathbf{y}^*\}$

# Intuition for Flat Search v.s. Hierarchical Search

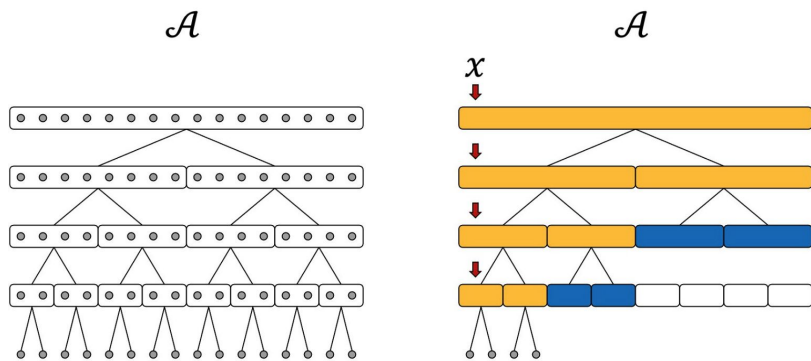


Figure 1. Hierarchical decomposition for the flat action space; the yellow nodes are further explored [Reference].

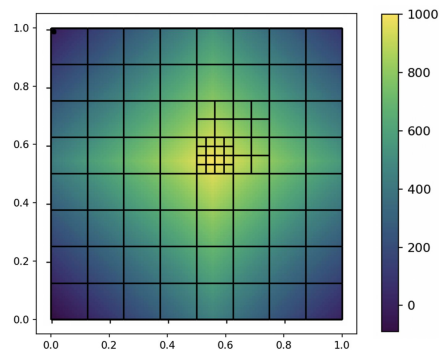


Figure 2. Partitioning over the action space [Reference].

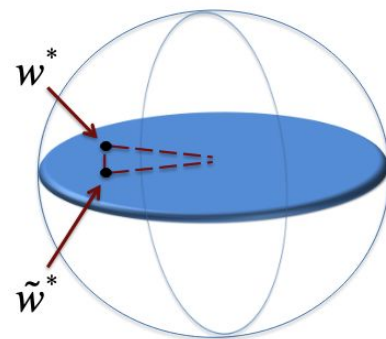


Figure 3. Focused exploration in subspaces [Reference].

## Method: (Offline SFT Stage) Goal-Driven Co-Training

$$\mathcal{L}_{\text{co}}(\theta) = -\frac{1}{N} \sum_{\underbrace{(s, y^*, s^*) \sim \mathcal{D}^{\text{train}}}_{\text{triplet set}}} \left[ \underbrace{\log p_{\theta}(s^* | s)}_{\text{goal planner}} + \underbrace{\log p_{\theta}(y^* | s, s^*)}_{\text{goal-driven actor}} \right]$$

Let's think in an information-theoretic way:  $s_{t+1}$  **acts as an information bottleneck** [Schwartz-Ziv and Tishby, 2017], by *abstracting* different possible proofsteps or sequences of proofsteps  $y_t$  into a single, more compact representation. Consider a simplified example below:

```
-- goal  $s_t = 3 * (2 + 1) = 9$   
-- goal  $s_{t+1} = 9 = 9$ 
```

There exist multiple different proofsteps to reach  $s_{t+1}$  from  $s_t$ , for instance:

- `ring` – algebraic normalization.
- `norm_num` – direct numeric evaluation.
- `simp; rfl` – simplification followed by reflexivity.
- `calc ... (omitted)` – step-by-step calculation.

# Method: (Online Search Stage) Goal-Driven Hierarchical Search

---

**Algorithm 1 RiR** – *A Unified Reasoning Mechanism with Decomposing and Search*

---

**Input:** problem statement  $q$ , a language model w/ parameter  $\theta$

```
1:  $\overline{\text{tree}} \leftarrow \overline{\text{Tree}}(\theta, q)$ 
2: repeat
3:    $s_i^* \leftarrow \overline{\text{tree}}.\text{policy}()$  ▷ /* planner */
4:    $\underline{\text{tree}} \leftarrow \underline{\text{Tree}}(\theta, s_i^*)$ 
5:   repeat
6:      $y_{t_i} \leftarrow \underline{\text{tree}}.\text{policy}()$  ▷ /* goal-driven actor */
7:   until STOP_LOW
8:    $\{\overline{\text{tree}}, \underline{\text{tree}}\}.\text{update}()$  ▷ /* joint update */
9: until STOP_HIGH
10: return  $\underline{\text{tree}}.\text{solution}$ 
```

---

- (Classical) **Flat planning:** we have a policy  $\pi_f : \mathcal{S} \rightarrow \mathcal{A}$  that maps states to actions.
- (RiR) **Hierarchical planning:** we have:
  - A *high-level planner* policy  $\pi_h : \mathcal{S} \rightarrow \tilde{\mathcal{S}}$ , that maps goals to target goals.
  - A *low-level actor* policy  $\pi_l : \mathcal{S} \times \tilde{\mathcal{S}} \rightarrow \mathcal{A}$ , that maps goals and target goals to actions.

# Results: Robust *Performance* Gains

Search Method (→)	Best-First Search	
Dataset (→)	miniF2F-test <sup>2</sup>	LeanDojo-test
Method (↓) / Model (→)	BYT5-0.3B	BYT5-0.3B
Reprover	34.43%	50.16%
RiR	<b>36.89%</b>	<b>53.73%</b>

Table 1: **Performance with BFS.** *Pass@1* rate on LeanDojo and miniF2F.

Search Method (→)	Monte-Carlo Tree Search	
Dataset (→)	miniF2F-test	LeanDojo-test
Method (↓) / Model (→)	BYT5-0.3B	BYT5-0.3B
Reprover	36.51%	50.24%
RiR	<b>37.83%</b>	<b>53.92%</b>

Table 2: **Performance with MCTS.** *Pass@1* rate on LeanDojo and miniF2F.

# Results: Remarkable *Efficiency* Gains

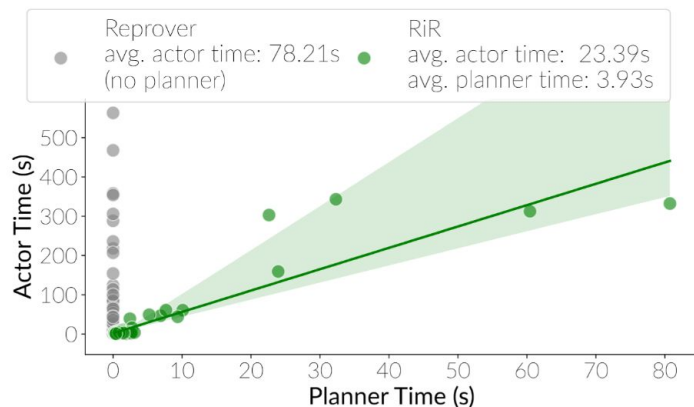


Figure 2: **Efficiency.** The scatter plot for actor and planner time spent for proved theorems on miniF2F. RiR significantly reduces the actor time via the goal guidance from the planner.

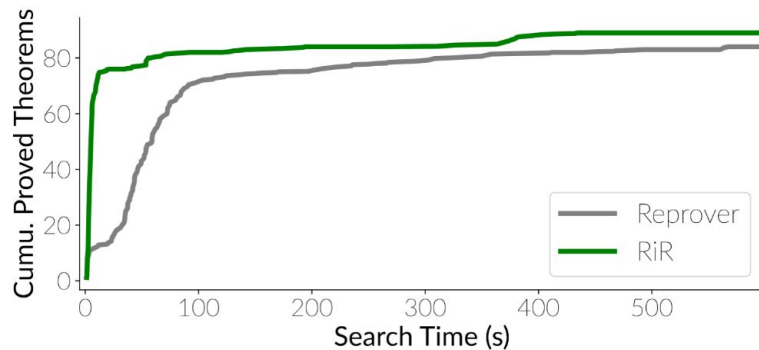


Figure 3: **Efficiency.** The CDF plot for search time spent for proved theorems on miniF2F Benchmark. RiR is significantly faster (nearly **3x**) than the existing state-of-the-art baseline.

# Example Proofs

## Example 3: Proof Found by RiR

```
Theorem:
  File Path: Mathlib/Order/SuccPred/Basic.lean
  Full Name: exists_succ_iterate_or

Status: Status.PROVED

Proof:
  obtain h |> h := le_total a b
  exacts [Or.inl (IsSuccArchimedean.exists_succ_iterate_of_le h),
  Or.inr (IsSuccArchimedean.exists_succ_iterate_of_le h)]

Search Statistics:
  Planner Time: 15.921687303110957
  Actor Time: 44.464585242792964
  Environment Time: 8.429574175737798
  Total Time: 68.86368872597814
  Total Nodes: 377
  Searched Nodes: 3
```

## Example 3: Failure by Reprover (w/o retrieval)

```
Theorem:
  File Path: Mathlib/Order/SuccPred/Basic.lean
  Full Name: exists_succ_iterate_or

Status: Status.OPEN

Proof: None

Search Statistics:
  Actor Time: 519.0408471203409
  Environment Time: 86.30267171841115
  Total Time: 605.4483464460354
  Total Nodes: 2819
  Searched Nodes: 95
```



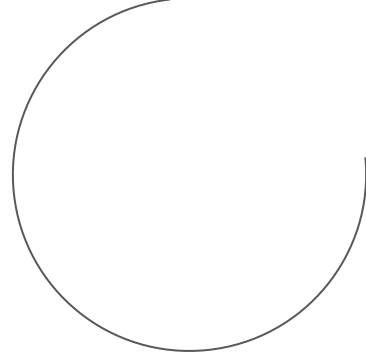
# Takeaways

**RiR** is a hierarchical framework for complex reasoning, unifying **decomposing** and **search**, and is **significantly faster** than classical stepwise reasoning, with **robust performance gains**.

The performance and efficiency gains come from:

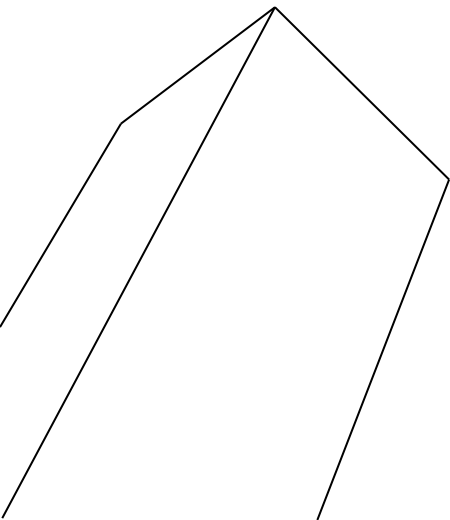
- **Offline co-training for SFT.**
- **Online bi-level search.**

*p.s.*, There are many different ways for decomposing!



# What Next?

**rigorous theories + more practical applications**



# Q & A

