

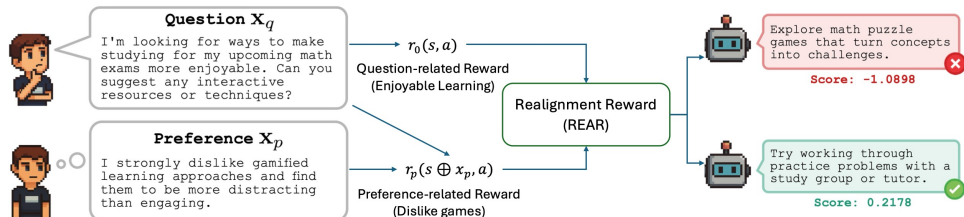
# REAR: Test-time Preference Realignment through Reward Decomposition



Fuxiang Zhang<sup>1</sup>, Pengcheng Wang<sup>2</sup>, Chenran Li<sup>2</sup>, Yi-Chen Li<sup>3</sup>, Yuxin Chen<sup>2</sup>, Lang Feng<sup>1</sup>, Chenfeng Xu<sup>2</sup>, Masayoshi Tomizuka<sup>2</sup>, Bo An<sup>1</sup>

<sup>1</sup> Nanyang Technological University <sup>2</sup> UC Berkeley <sup>3</sup> Nanjing University

## Controllable Realignment for General LLMs



- Human preferences are **diverse and personal**, yet a pretrained LLM bakes in a **single, fixed preference balance** from its training data.
- General LLMs** are good at **answering questions** but bad at **recognizing user preference**.

*How can we adjust the balance between question & preference of an LLM?*

- Re-training / fine-tuning with specific preference level – costly
- Test-time alignment – no verifiable approach (diff. from math, code, ...)

## Key Insights

- The base model is already an **implicit reward model**.
- Under the following reward decomposition:
  - $s$ : generation state wo. preference (question only)
  - $s \oplus x_p$ : generation state with preference
- The base model maximizes the question-related reward wo. preference and maximizes a combination of question-related and preference-related reward with some  $\alpha$

Prop. 3.1 — latent reward structure

$$r(s \oplus x_p, a) = \underbrace{r_0(s, a)}_{\text{question}} + \alpha \underbrace{r_p(s \oplus x_p, a)}_{\text{preference}}$$

## REAR: Realignment Reward

- How to rebalance two reward terms? Adjusting  $\alpha$  to  $\hat{\alpha}$  results in a new reward

### Token-level Reward

**Lemma 3.2.** The realignment reward  $r_{\text{REAR}}(s \oplus x_p, a)$  is equal to

$$r_{\text{REAR}}(s \oplus x_p, a) = \frac{(\alpha - \hat{\alpha})\beta}{\alpha} \log \pi(a | s) + \frac{\hat{\alpha}\beta}{\alpha} \log \pi(a | s \oplus x_p) + Z(s) - \gamma Z(s'), \quad (8)$$

where the state-dependent term  $Z(s) = (1 - \frac{\hat{\alpha}}{\alpha}) V^\pi(s) + \frac{\hat{\alpha}}{\alpha} V^\pi(s \oplus x_p)$ .

### Cumulative Rewards over Trajectory

$$\begin{aligned} \bar{R}_{\text{REAR}}(\tau) &= \sum_{s_t, a_t} \gamma^t r_{\text{REAR}}(s_t \oplus x_p, a_t) \\ &= \sum_{s_t, a_t} \gamma^t \left( (1 - \lambda) \beta \log \pi(a_t | s_t) \right. \\ &\quad \left. + \lambda \beta \log \pi(a_t | s_t \oplus x_p) \right) + Z(s_0), \end{aligned}$$

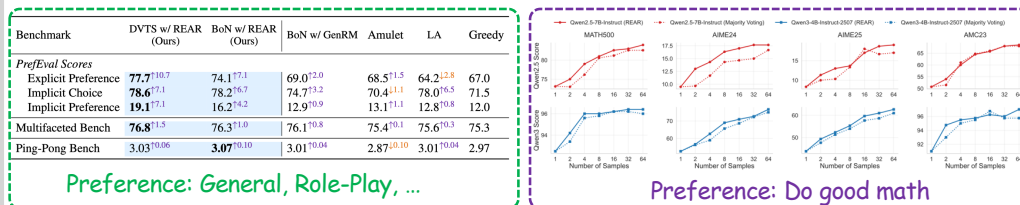
### REAR Score (Drop state-only terms)

$$S_{\text{REAR}}(\tau) = \sum_{t=0}^T \gamma^t \left( (1 - \lambda) \log \pi(a_t | s_t) + \lambda \log \pi(a_t | s_t \oplus x_p) \right).$$

Now we have two **test-time methods** using  $S_{\text{REAR}}$ :

- Best-of-N** (2) **Tree search (DVTS)** (calculated from sub-trajectories)

## Experiment Results



Preference: Read the image without hallucinations



Preference:  
You are a careful visual assistant. Ground every statement strictly in the image and the question...  
Question:  
How many bicycles are there in the image?  
Answer:  
There are four bicycles: two at the front, one in the middle, and one in the distance.

Method	MMHal Score (↑)	Hallucination (↓)
BoN w/ REAR	<b>84.20</b>	<b>21.87</b>
BoN w/ GenRM	80.21	23.96
Greedy	78.99	28.12