

# GRPO-Guard: Mitigating Implicit Over-Optimization in Flow Matching via Regulated Clipping

Jing Wang<sup>1,2</sup>, Jiajun Liang<sup>2\*</sup>, Jie Liu<sup>3</sup>, Henglin Liu<sup>2,4</sup>, Gongye Liu<sup>2,5</sup>, Jun Zheng<sup>1</sup>, Wanyuan Pang<sup>6</sup>,  
Ao Ma<sup>7</sup>, Zhenyu Xie<sup>1</sup>, Xintao Wang<sup>2</sup>, Meng Wang<sup>2</sup>, Pengfei Wan<sup>2</sup>, Xiaodan Liang<sup>1,8,9†</sup>

<sup>1</sup>Shenzhen Campus of Sun Yat-Sen University, <sup>2</sup>Kling Team, Kuaishou Technology, <sup>3</sup>CUHK MMLab, <sup>4</sup>Tsinghua University, <sup>5</sup>HKUST

<sup>6</sup>USTB, <sup>7</sup>UCAS, <sup>8</sup>Peng Cheng Laboratory, <sup>9</sup>Guangdong Key Laboratory of Big Data Analysis and Processing

\*Project Leader †Corresponding Author *Emails: wangj977@mail2.sysu.edu.cn, xdliang328@gmail.com*

CVPR  
JUNE 3-7, 2026



DENVER  
COLORADO



KlingAI

## Motivation

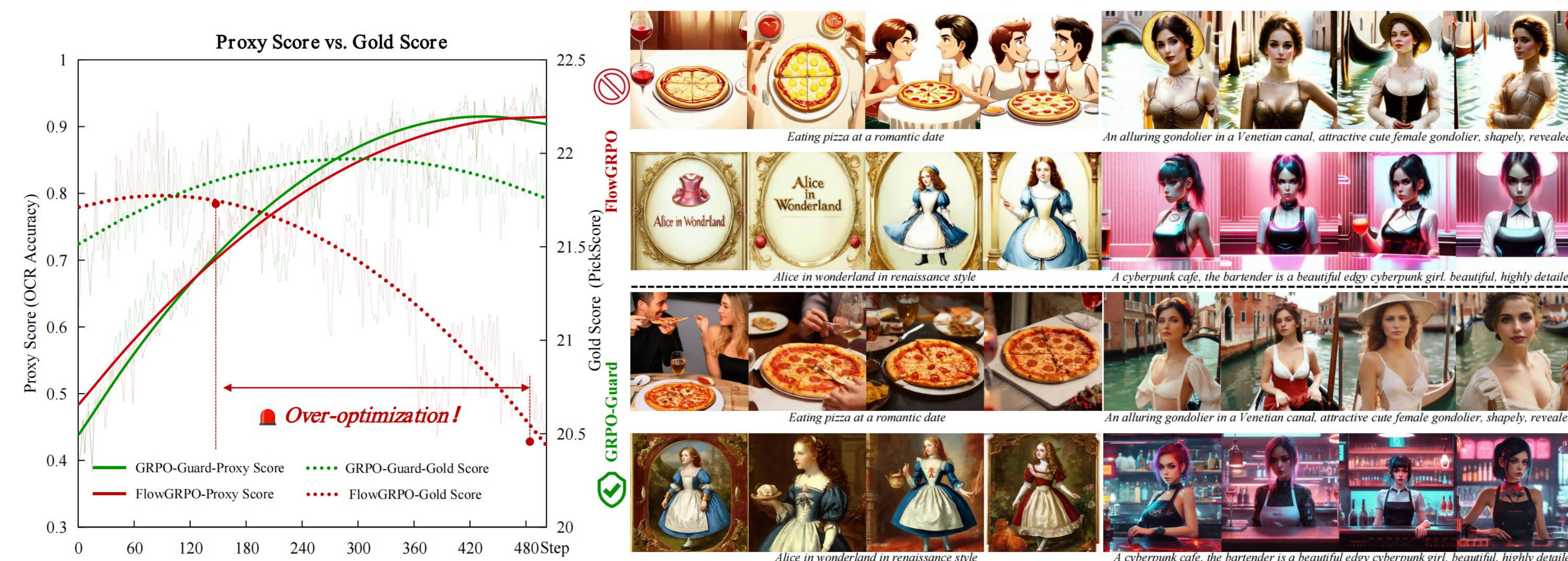


Figure 1. Comparison between FlowGRPO and GRPO-Guard under over-optimization. Left: The proxy score and gold score trends during training. As the proxy score increases, FlowGRPO rapidly enters an over-optimization phase, where the gold score continuously declines. Right: A visual comparison between FlowGRPO and GRPO-Guard. Due to severe reward hacking, FlowGRPO suffers from a drastic degradation in diversity, detail richness, visual quality, and text-image consistency (upper part). In contrast, GRPO-Guard maintains a stable gold score and high visual quality under a comparable proxy score, as shown in the bottom part of the figure.

In Diffusion GRPO, the policy model inevitably enters an implicit *over-optimization stage*—while the proxy reward continues to increase, essential metrics such as image quality and text-prompt alignment deteriorate sharply, ultimately making the learned policy *impractical for real-world use*.

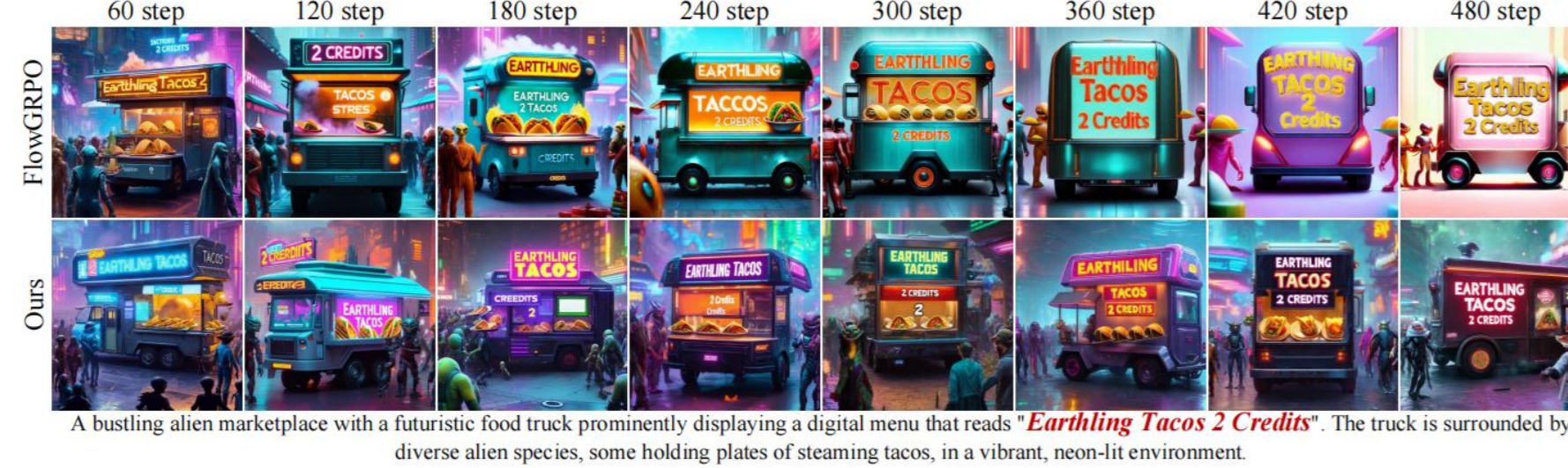


Figure 9. Samples of the policy model at different training steps.

## Method

$$\text{GRPO } \mathcal{J}_{\text{policy}}(\theta) = \frac{1}{G} \sum_{i=1}^G \frac{1}{T} \sum_{t=0}^{T-1} \left( \min(r_t^i(\theta) \hat{A}_t^i, \text{clip}(r_t^i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^i) \right) \quad (5)$$

$$\begin{aligned} \log r_t(\theta) &= \log p_{\theta}(x_{t-1}|x_t, \mathbf{c}) - \log p_{\theta_{old}}(x_{t-1}|x_t, \mathbf{c}) \\ &= -\frac{\|\mu_{\theta_{old}}(x_t, t) - \mu_{\theta}(x_t, t) + \sigma_t \sqrt{dt} \cdot \epsilon\|^2}{2\sigma_t^2 dt} \\ &\quad + \frac{\|\mu_{\theta_{old}}(x_t, t) - \mu_{\theta_{old}}(x_t, t) + \sigma_t \sqrt{dt} \cdot \epsilon\|^2}{2\sigma_t^2 dt} \\ &= -\frac{\|\Delta\mu_{\theta} + \sigma_t \sqrt{dt} \cdot \epsilon\|^2}{2\sigma_t^2 dt} + \frac{\|\epsilon\|^2}{2} \end{aligned}$$

$$\text{Bias} \rightarrow -\frac{\|\Delta\mu_{\theta}\|^2}{2\sigma_t^2 dt} - \frac{\Delta\mu_{\theta} \cdot \epsilon}{\sigma_t \sqrt{dt}} \quad (7)$$

$$\text{RatioNorm} \rightarrow \log \hat{r}_t(\theta) = \sigma_t \sqrt{dt} (\log r_t(\theta) + \frac{\|\Delta\mu_{\theta}\|^2}{2\sigma_t^2 dt}) = -\Delta\mu_{\theta} \cdot \epsilon \quad (8)$$

$$\begin{aligned} \nabla_{\theta} \mathcal{J}(\theta) &= \sum_{t=0}^{T-1} \hat{A}_t \nabla_{\theta} \log \hat{r}_t(\theta) = \sum_{t=0}^{T-1} \hat{A}_t \hat{r}_t(\theta) \nabla_{\theta} \log \hat{r}_t(\theta) \\ &= \sum_{t=0}^{T-1} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ \epsilon \hat{A}_t \hat{r}_t(\theta) \nabla_{\theta} \log \hat{r}_t(\theta) \right] \\ &\quad + \sum_{t=0}^{T-1} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, I)} \left[ \beta dt \epsilon \hat{A}_t \hat{r}_t(\theta) \nabla_{\theta} \log \hat{r}_t(\theta) \right] \end{aligned} \quad (11)$$

$$\mathcal{J}_{\text{policy}}(\theta) = \frac{1}{G} \sum_{i=1}^G \frac{1}{T} \sum_{t=0}^{T-1} \left( \delta \min(\hat{r}_t^i(\theta) \hat{A}_t^i, \text{clip}(\hat{r}_t^i(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^i) \right) \quad (12)$$

GRPO-Guard

## Analysis

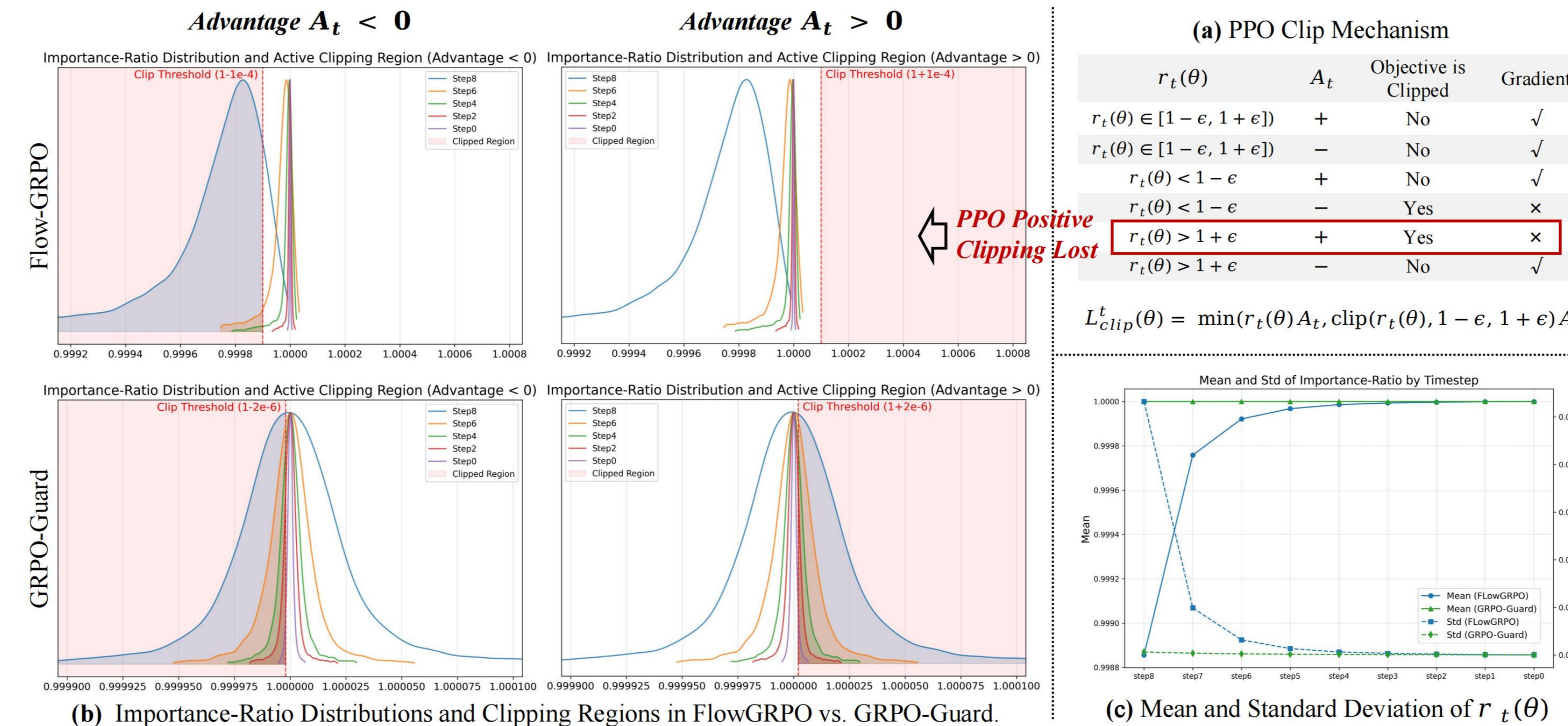


Figure 2. Comparison of  $r(\theta)$  distributions between FlowGRPO and GRPO-Guard across timesteps. (a) Ideally, the ratio distribution should have a mean near 1 and stable variance across timesteps to ensure effective clipping. (b) Under FlowGRPO, the distribution exhibits a leftward mean shift and increasing variance at low-noise timesteps, causing the clipping mechanism to fail—particularly for trajectories with positive advantages. In contrast, GRPO-Guard with RatioNorm preserves a balanced mean and consistent variance (c), enabling proper clipping and stable policy updates across all timesteps.

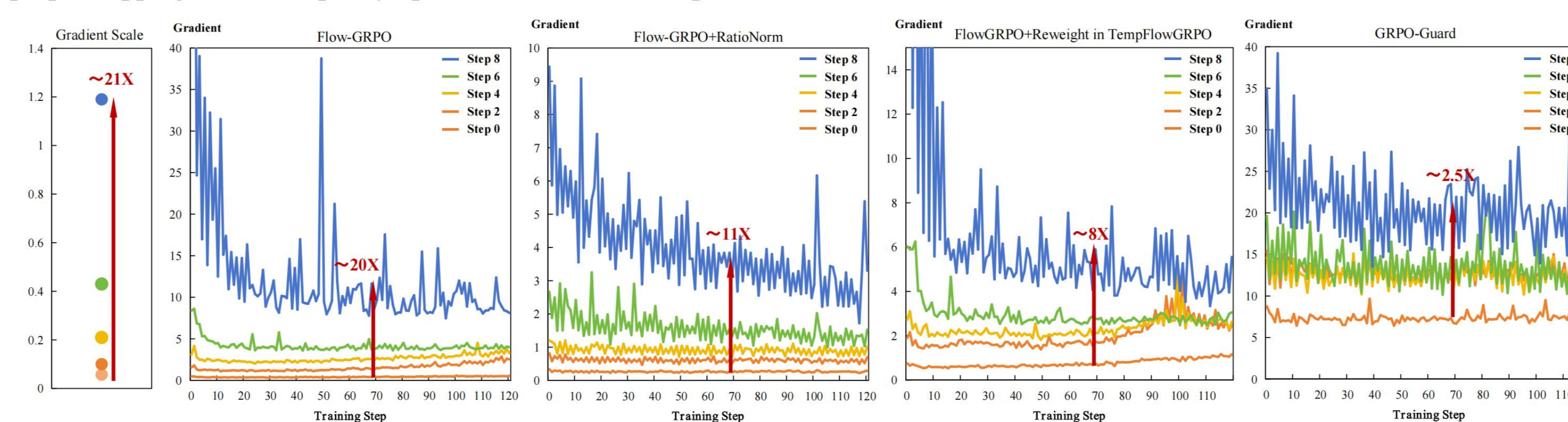


Figure 3. Gradient magnitude differences across timesteps. In FlowGRPO, gradient magnitudes vary by roughly 20x across timesteps, reflecting the large differences in gradient scale. GRPO-Guard substantially reduces this imbalance, limiting the variation to about 2.5x and preventing over-optimization under any single noise condition.

Setting	$\log r_t(\theta)$	Re-weight Scale	Gradient Scale
Baseline	$\log r_t(\theta)$	1	$\beta \frac{\Delta\mu_{\theta} + \sigma_t \sqrt{dt} \epsilon}{\sigma_t^2 dt}$
Temp-Reweight [16]	$\log r_t(\theta)$	$\sigma_t \sqrt{dt}$	$\beta \frac{\sqrt{dt} \Delta\mu_{\theta} + \sigma_t dt \epsilon}{\sigma_t^2 dt}$
Guard-Reweight	$\log r_t(\theta)$	$1/dt$	$\beta \frac{\Delta\mu_{\theta} + \sigma_t \sqrt{dt} \epsilon}{\sigma_t^2 dt}$
Mean-revised	$\log r_t(\theta) + \frac{\ \Delta\mu_{\theta}\ ^2}{2\sigma_t^2 dt}$	1	$\beta \frac{\sqrt{dt} \epsilon}{\sigma_t}$
RatioNorm	$\sigma_t \sqrt{dt} (\log r_t(\theta) + \frac{\ \Delta\mu_{\theta}\ ^2}{2\sigma_t^2 dt})$	1	$\beta dt \epsilon$
GRPO-Guard	$\sigma_t \sqrt{dt} (\log r_t(\theta) + \frac{\ \Delta\mu_{\theta}\ ^2}{2\sigma_t^2 dt})$	$1/dt$	$\beta \epsilon$

Table 2. Ablation study on major components.

## Ablation Study

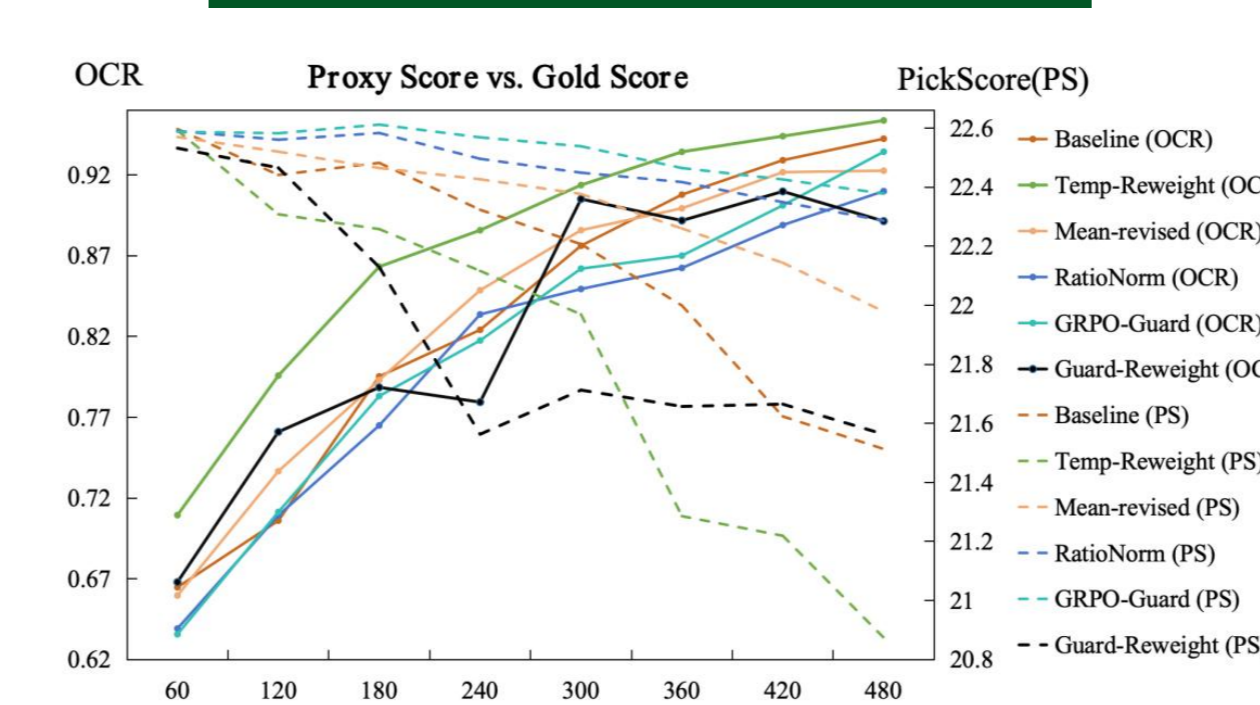


Figure 7. Training curves of the ablation study.

## Experiments

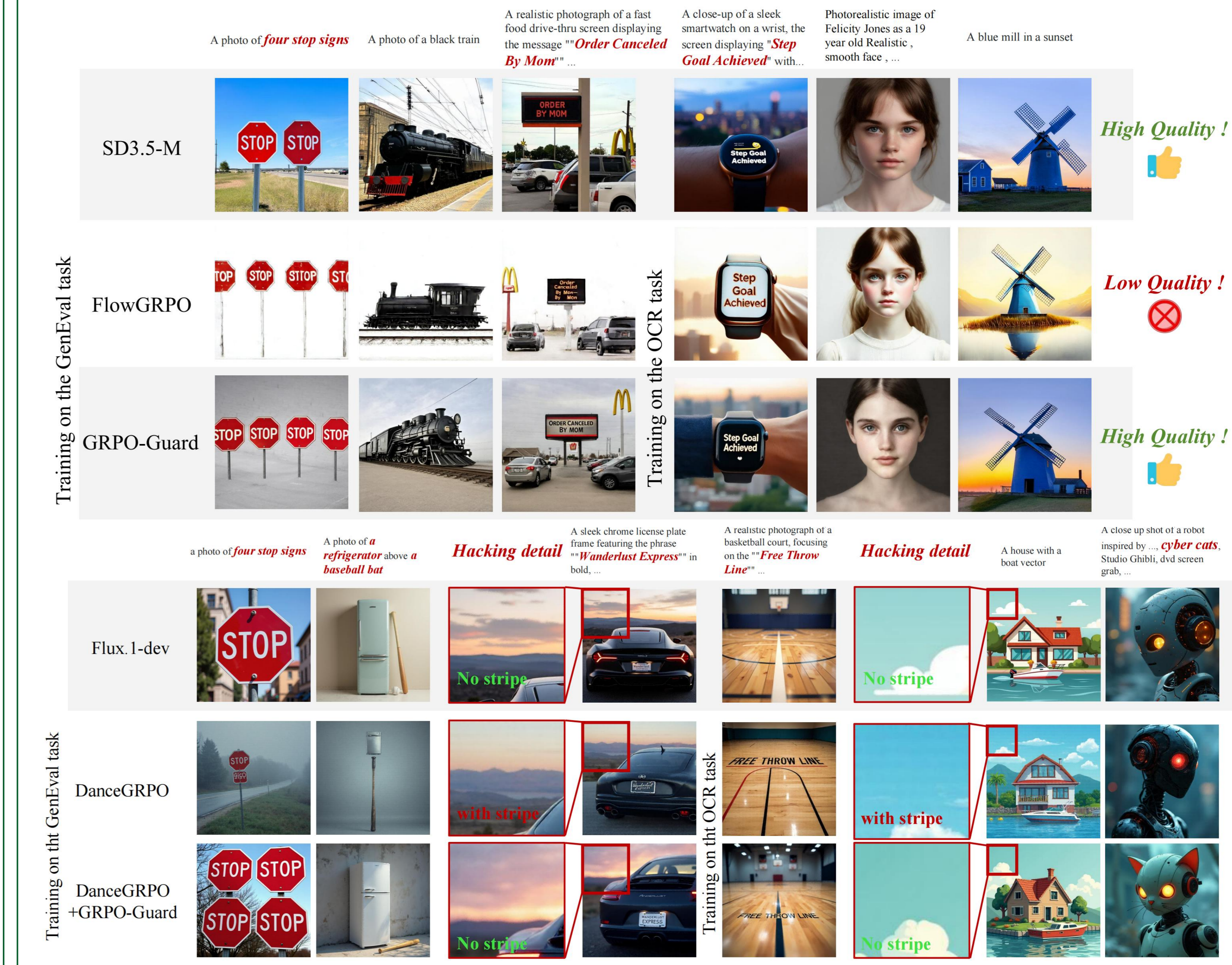


Table 1. Comparison of composite gold scores across proxy tasks. [-] denotes the proxy task for each row. ImR and UniR represent ImageReward and UnifiedReward, respectively. Average is the mean of the three gold scores normalized by the base model (set to 1). FG, DG, and TempFG denote FlowGRPO, DanceGRPO, and TempFlowGRPO settings.

Method	Step	GenEval			PickScore Text			Gold Score				
		GenEval	PickScore	Text	HPSv2	ImR	UniR	Average	HPSv2	ImR	UniR	Average
SD3.5-M [32]	-	0.63	21.5	0.58	0.293	1.06	3.31	1.00	-	-	-	-
+TempFG [16]	480	-	-	[0.90]	0.221	0.02	2.58	0.52	-	-	-	-
+FlowGRPO	1860	[0.94]	20.4	0.59	0.236	0.85	3.05	0.84	-	-	-	-
+Ours (FG)	1860	[0.95]	+0.01	20.9+0.4	0.71+0.12	0.254+0.018	0.87+0.02	3.22+0.17	0.89+0.05	-	-	-
+FlowGRPO	10200	0.67	[23.1]	0.64	0.329	1.40	3.46	1.16	-	-	-	-
+Ours (FG)	10200	0.70	+0.03	[23.3]	+0.2	0.68+0.04	0.337+0.008	1.47+0.07	3.54+0.08	1.20+0.04	-	-
+FlowGRPO	480	0.52	20.8	[0.94]	0.274	0.82	3.07	0.88	-	-	-	-
+Ours (FG)	480	0.65+0.07	21.3+0.5	[0.93]	-0.01	0.286+0.012	1.06+0.24	3.29+0.22	0.99+0.11	-	-	-
Flux.1-dev [21]	-	0.63	21.6	0.60	0.302	1.01	3.31	1.00	-	-	-	-
+DanceGRPO	1260	[0.80]	21.2	0.60	0.269	0.79	3.18	0.88	-	-	-	-
+Ours (DG)	1260	[0.81]	+0.01	21.7+0.5	0.63+0.03	0.300+0.031	1.08+0.29	3.35+0.17	1.02+0.14	-	-	-
+DanceGRPO	540	0.63	21.5	[0.90]	0.293	0.93	3.25	0.96	-	-	-	-
+Ours (DG)	540	0.64+0.01	21.8+0.3	[0.89]	-0.01	0.304+0.009	1.07+0.14	3.35+0.10	1.02+0.06	-	-	-

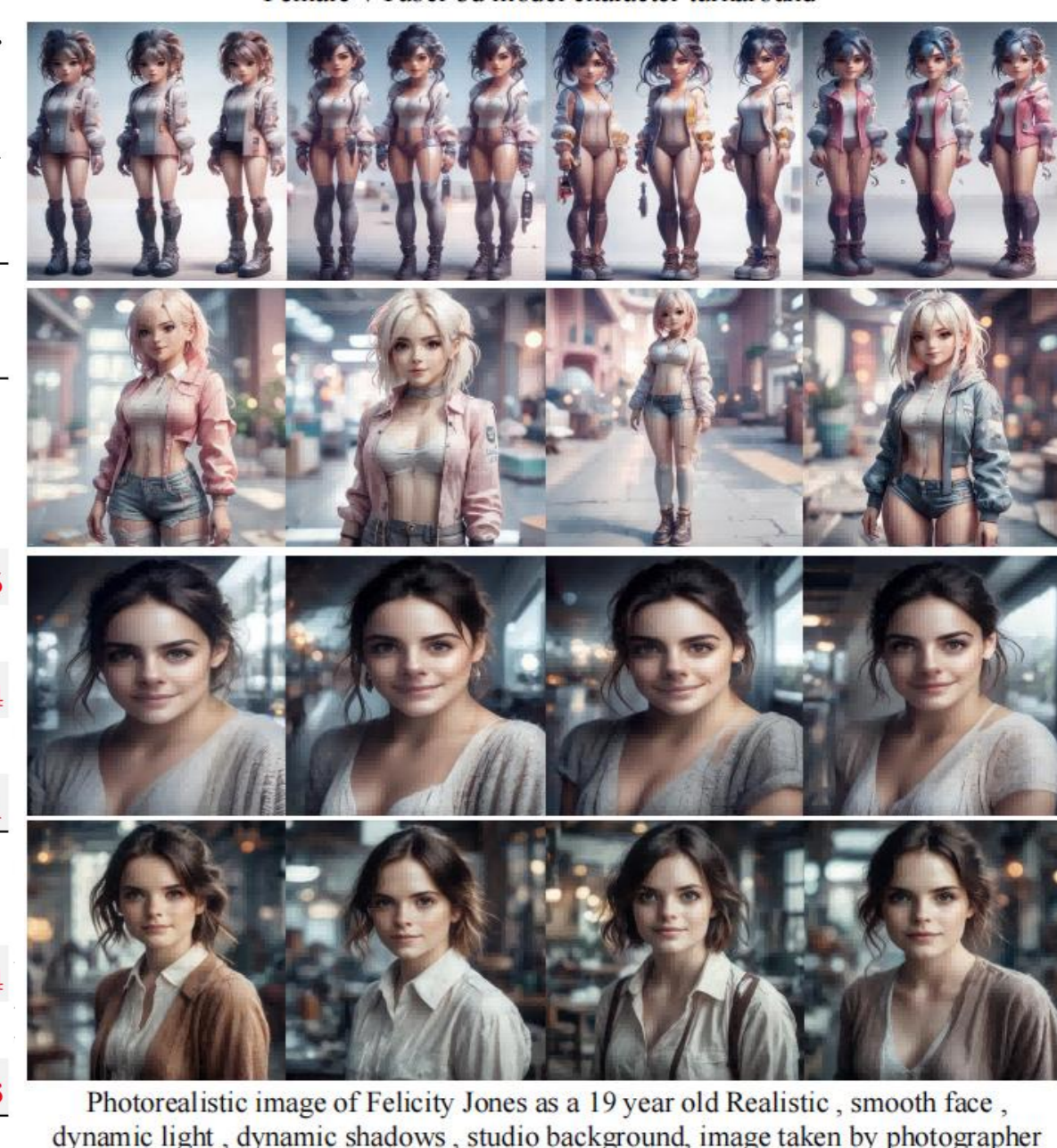


Figure 5. Comparison of composite gold scores across proxy tasks. [-] denotes the proxy task for each row. ImR and UniR represent ImageReward and UnifiedReward, respectively. Average is the mean of the three gold scores normalized by the base model (set to 1). FG, DG, and TempFG denote FlowGRPO, DanceGRPO, and TempFlowGRPO settings.