

Language model reinforcement learning at a glance

Core loss $\mathcal{L}(\theta)$ per RL algorithm

Policy Gradient	$- \mathbb{E}_\tau \left[\sum_{t=0}^T \log \pi_\theta(a_t s_t) A^{\pi_\theta}(s_t, a_t) \right]$
REINFORCE (REward Increment = Nonnegative Factor \times Offset Reinforcement \times Characteristic Eligibility)	$- \frac{1}{T} \sum_{t=1}^T \log \pi_\theta(a_t s_t) (G_t - b(s_t))$
REINFORCE Leave One Out (RLOO)	$- \frac{1}{K} \sum_{i=1}^K \sum_t \log \pi_\theta(a_{i,t} s_{i,t}) \left(R_i - \frac{1}{K-1} \sum_{j \neq i} R_j \right)$
Proximal Policy Optimization (PPO)	$- \frac{1}{T} \sum_{t=1}^T \min \left(\frac{\pi_\theta(a_t s_t)}{\pi_{\theta_{\text{old}}}(a_t s_t)} A_t, \text{clip} \left(\frac{\pi_\theta(a_t s_t)}{\pi_{\theta_{\text{old}}}(a_t s_t)}, 1-\varepsilon, 1+\varepsilon \right) A_t \right)$
Group Relative Policy Optimization (GRPO)	$- \frac{1}{G} \sum_{i=1}^G \frac{1}{ a_i } \sum_{t=1}^{ a_i } \min \left(\frac{\pi_\theta(a_{i,t} s_{i,t})}{\pi_{\theta_{\text{old}}}(a_{i,t} s_{i,t})} \hat{A}_i, \text{clip} \left(\frac{\pi_\theta(a_{i,t} s_{i,t})}{\pi_{\theta_{\text{old}}}(a_{i,t} s_{i,t})}, 1-\varepsilon, 1+\varepsilon \right) \hat{A}_i \right)$ where $\hat{A}_i = \frac{r_i - \bar{r}}{\text{std}(r)}$
Group Sequence Policy Optimization (GSPO)	$- \frac{1}{G} \sum_{i=1}^G \min \left(\left(\frac{\pi_\theta(a_i s)}{\pi_{\theta_{\text{old}}}(a_i s)} \right)^{\frac{1}{ a_i }} A_i, \text{clip} \left(\left(\frac{\pi_\theta(a_i s)}{\pi_{\theta_{\text{old}}}(a_i s)} \right)^{\frac{1}{ a_i }}, 1-\varepsilon, 1+\varepsilon \right) A_i \right)$
Clipped Importance Sampling Policy Optimization (CISPO)	$- \frac{1}{\sum_i a_i } \sum_{i=1}^K \sum_{t=1}^{ a_i } \text{sg} \left(\text{clip} \left(\frac{\pi_\theta(a_{i,t} s)}{\pi_{\theta_{\text{old}}}(a_{i,t} s)}, 1-\varepsilon_{\text{low}}, 1+\varepsilon_{\text{high}} \right) \right) A_{i,t} \log \pi_\theta(a_{i,t} s)$ where $\text{sg}(\cdot) = \text{stop gradient}$

Other core equations

RLHF Objective	$J(\theta) = \max_{\pi} \mathbb{E}_{x \sim \mathcal{D}} \mathbb{E}_{y \sim \pi(y x)} \left[r_\theta(x, y) \right] - \beta \mathcal{D}_{\text{KL}}(\pi(y x) \ \pi_{\text{ref}}(y x))$
Bradley–Terry Reward Model	$\mathcal{L}(\theta) = - \log \sigma(r_\theta(y_c x) - r_\theta(y_r x))$
Direct Preference Optimization (DPO)	$\mathcal{L}(\theta) = - \mathbb{E} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_c x)}{\pi_{\text{ref}}(y_c x)} - \beta \log \frac{\pi_\theta(y_r x)}{\pi_{\text{ref}}(y_r x)} \right) \right]$

Notation

x, y	Prompt, completion	A_t	Advantage estimate
y_c, y_r	Chosen / rejected completion	β	KL penalty coefficient
$y_c \succ y_r$	y_c is preferred over y_r	$\sigma(z)$	Sigmoid: $1/(1 + e^{-z})$
π_θ	Policy (the model being trained)	$\mathcal{D}_{\text{KL}}(P \ Q)$	KL divergence between P and Q
π_{ref}	Reference policy (frozen copy)	G_t	Return: $\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
$\pi_{\theta_{\text{old}}}$	Policy at start of RL batch updates	$V(s)$	Value: $\mathbb{E}[G_t S_t = s]$
$r_\theta(y x)$	Reward model score		