

## Direct Preference Optimization

x: "Write me a poem about the history of jazz."



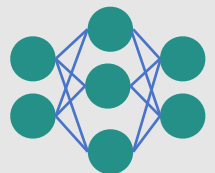
$y_w$



$y_l$

MLE

$(y_w, y_l | x)$



LLM Policy

$$\bar{\pi} = \arg \max_{\pi} \mathbb{E} \left[ \log \sigma \left( \beta \log \frac{\pi(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right]$$

Rollout Trajectories 

## Intrinsic Self-Reflective Preference Optimization

x: "Write me .."



$y_w$

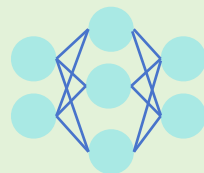


$y_l$

$(y_w, y_l | x)$   
Current Pairs

Cross-Conditioning Self-Reflection

$(y_w | y_l, x)$   
 $(y_l | y_w, x)$



Final Policy

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[ \log \sigma \left( \beta \log \frac{\pi(y^w | y^l, x)}{\pi_{\text{ref}}(y^w | x)} - \beta \log \frac{\pi(y^l | y^w, x)}{\pi_{\text{ref}}(y^l | x)} \right) \right]$$

Intrinsic  
Self-Reflection

Within One Trajectory 