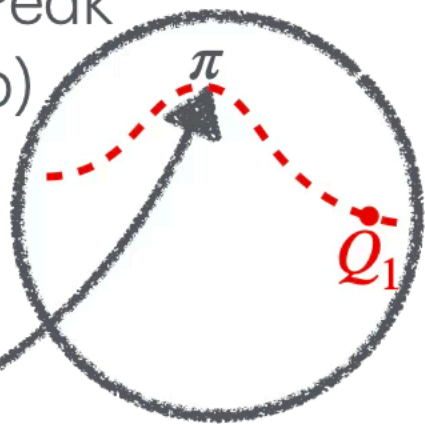
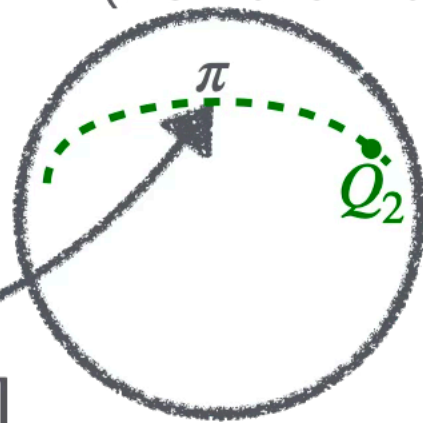


GRPO: $\max_{\pi} \mathbb{E}_{\pi}[r]$

High Reward Peak
(Sharp drop)



Stable High
(Reward-flat)



FRPO: $\max_{\pi} \min_Q \mathbb{E}_Q[r]$

π_{ref}

Refusal Rate on Harmbench

