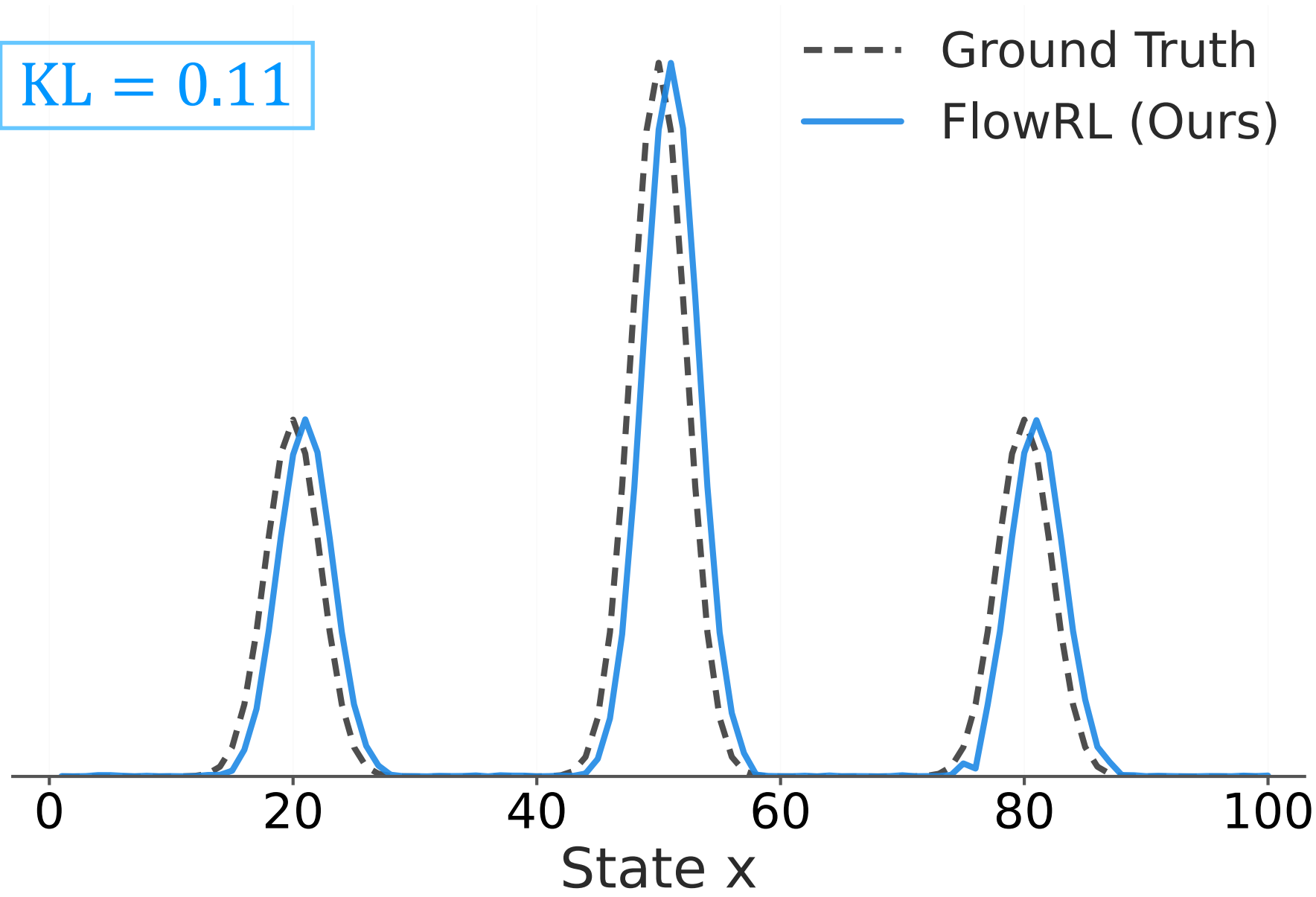


Distribution-matching: FlowRL



Reward-maximizing : R++, PPO and GRPO

