**SFT**
*Upweight off-policy reference solutions*
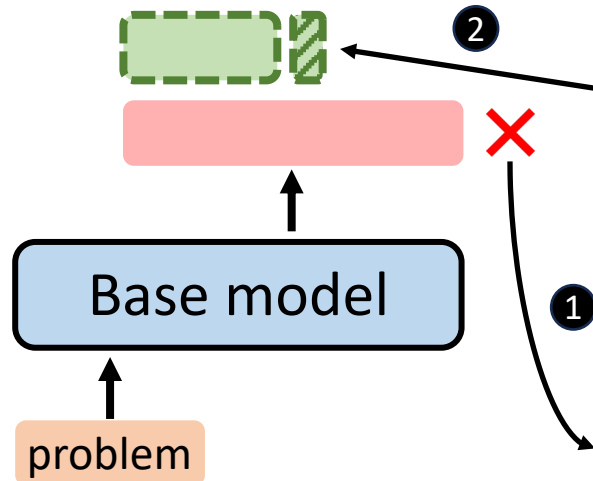
**RL**
*Uniformly update on-policy rollouts*

single rollout

**Intervention Training**
*Upweight interventions & steps before mistake*

*Self-verify to find mistake*

②

*Propose intervention*

①

Base model

problem

Base model

problem

Base model

problem

Base model

incorrect rollout

reference solution

Correct  Incorrect  Update  No update