

Reinforcement Learning (cont)

A Evaluation (after training measure performance)

- 1) Freeze the policy (no more learning)
- 2) run agent many episodes in the environment
- 3) collect metrics like:
 - Avg reward per episode
 - success rate (e.g. % of times robot reaches goal)
 - Cumulative reward curve over time
 - stability (does performance vary a lot)

- this tells us how well the agent generalizes not just memorized

A RL vs RLHF

- RL:
- Env: game / robotic world
 - Agent: learns by taking action in env
 - Reward: comes directly from env, Ex: +1: win or -1: loss
 - training loop: Agent in env → reward → update policy
 - Eval: Avg score, success rate
- used in games, cars, robotics

Model Based vs Model Free

Model Based: learns a model of the env or may be given model, can plan ahead or models can be value or policy based but relies on model for planning (Ex: Dyna-Q, planning methods)

Model Free: Doesn't learn model env, learns directly Value function or Policy, no explicit transition or reward models (Exs Q-Learn, DQN, PPO, Reinforce)

- RLHF:
- Env: Human preference data
 - Agent: a LLM
 - Reward: not built in → learned from humans
- 1) collect LLM outputs
 - 2) human rank which ones better
 - 3) train reward model to predict human preference
 - 4) use RL (PPO, Actor-Critic) to fine tune LLM for maximized reward
- Training: LLM generates → reward model scores → policy update
- Eval: human eval, quality checks, alignment benchmarks

- used in ChatGPT, Claude etc.