

Reinforcement learning cont

* 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 (eg % of times robot reached 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: win + loss
 - Training loop: Agent in env \rightarrow reward \rightarrow update policy
 - Eval: Avg score, success rate
- used in games, cars, robotics

- RLHF:
- Env: human preference data
 - Agent: a LLM
 - Reward: not built in \rightarrow learned from humans
- 1) collect LLM outputs
 - 2) humans rank which ones better
 - 3) train Reward model to predict human preference
 - 4) use RL (PPO AC model) to fine tune LLM for maximized reward
- Training: LLM generates \rightarrow reward model scores \rightarrow policy updates
 - Eval: human eval, quality checks, alignment benchmarks
- used in Chat GPT, Claude etc