# Last time: Exploration

Best: RND

        Go-Explore

## This time: Transfer, Meta; Grand Accomplishments RL

Transfer Learning : Using experience from one set of tasks for a new $\underset{\uparrow\ MDP}{\text{task}}$

Can RL use prior knowledge?

    Where to store?

        - Q-function
        - Policy
        - Model
        - Features / Hidden States
              ↳ Aside: Representation Bottleneck

Jargon

"source" ⟶ "target"

"shot" = attempt in target domain

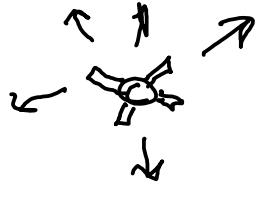    "1 shot"
    "0 shot"
    "few shot"

Forward
Multi Task
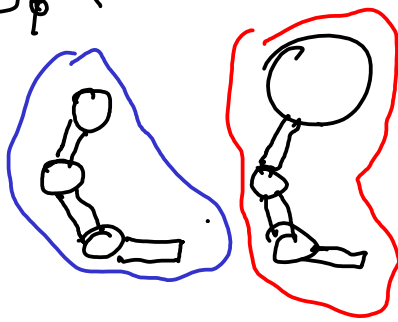Meta Learning

⇒ Forward : Fine Tuning
    Key : Lots of randomness
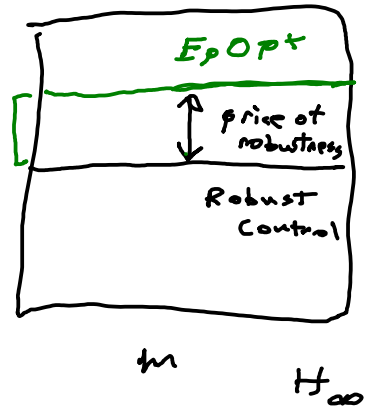
$R(s,a) + H[\pi_\theta]$

    Example



$EpOp_t$



$E[R]$

Test performance



$m$



$EpOpt$ (green)
price of robustness
Robust Control
$m$
$H_\infty$

    CADZRL

⇒ Multi - Task
    Contextual Policy
        $\pi_\theta(a|s)$
        $\pi_\theta(a|s,\omega)$
            where to place an object
            direction to run in

$\tilde{S} = S \times \Omega$ ← Context space

    Modular

# Meta RL : Learning to Learn

### RL

$$\theta^* = \underset{\theta}{\arg\max} \; E_{\pi_\theta}\left[R(\tau)\right]$$

$$\theta^* = f_{RL}(\underset{\underset{MDP}{\uparrow}}{M})$$

$\underset{RL \; algorithm}{\uparrow}$

### Meta RL

$$\theta^* = \underset{\theta}{\arg\max} \; \sum_{i=1}^{n} E_{\pi_{\phi_i}}\left[R_i(\tau)\right]$$

$$\phi_i = f_\theta (M_i)$$

1. Sample task $i$, collect $D_i$
2. adapt policy $\phi_i = f(\theta, D_i)$
3. Collect $D_i'$ with $\pi_{\phi_i}$
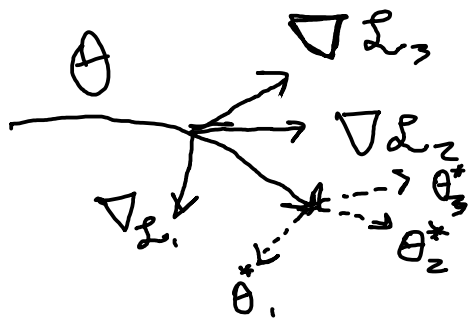4. Update $\theta$ with $\mathcal{L}(D_i', \phi_i)$

# 3 Solutions

1. Recurrence        RNN
   Example $RL^2$

2. Optimization      MAML

# 3. ML as a POMDP

hidden state is task

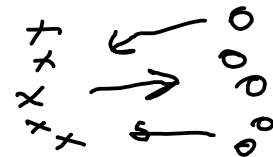$$\tilde{S} = S \times \mathcal{M} \leftarrow \text{set of MDPs}$$

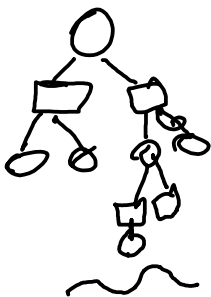$$O_t = (s_t, r_t)$$

## Learn Filter

---

## Most Celebrated RL Accomplishment

| Name | Domain | Online | Deep Learning | Key | Training Time |
|------|--------|--------|---------------|-----|---------------|
| AlphaGo | Go | MCTS | Policy NN Value NN | Trained on - human games, self-play | 3 weeks |
| AlphaZero | Go, chess shogi | MCTS | " | self play | 24 hours |
| Open AI 5 | Dota | - | PPO with LSTM expert-specified dense reward | League Play | weeks |
| Alpha Star | Starcraft II | - | Actor-Critic LSTM TD($\lambda$) returns | (League Play) | 40 days |
| Deep Stack | Poker | Heuristic Search | Counterfactual Value Network | Incomplete Info | 2 days on 1 GPU |
| FTW Agent | Quake CTF | - | Timescale Hierarchical Actor-Critic & NN learned Denserew | League Play | 500k games×40min = 9.5 years |
| MuZero | Board games Atari | MCTS | Policy Value Model | SelfPlay | 12 hours |

Alpha Go 2015

Dota



$$|A| \cong 19 \times 19 = 361$$

Policy Network: suggest a few good action

Value Network: evaluates state

nodes at leaves