

# **Sample-efficient Learning in Many-goals Settings**

---

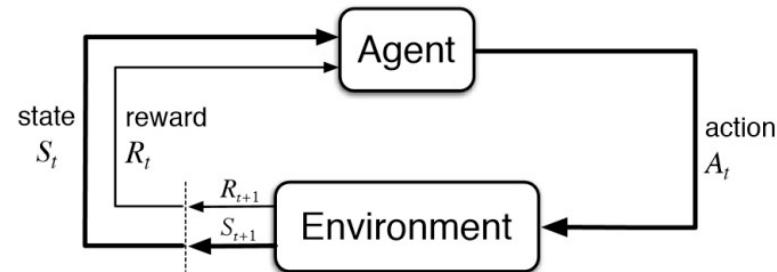
**PRESENTER: BYUNGSOO JEON**

**LLL READING GROUP @ CMU**

# Reminder: RL

## Reinforcement Learning

- Given the environment, how do we learn the sequence of actions (policy) to maximize the cumulative reward?



**Value Function,**  $V_\pi(s) = \mathbb{E}_\pi [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$

- Expected total reward for an agent starting from state  $s$

# Reminder: Task vs. Goal

---

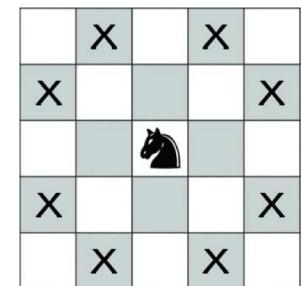
## Task

- Def: A piece of work to be done
- e.g. Beating an opponent in the chess (or Starcraft)



## Goal

- Def: The end toward which effort is directed
- e.g. Move the knight to one of good positions



# **Hindsight Experience Replay**

---

**PRESENTOR: BYUNGSOO JEON**

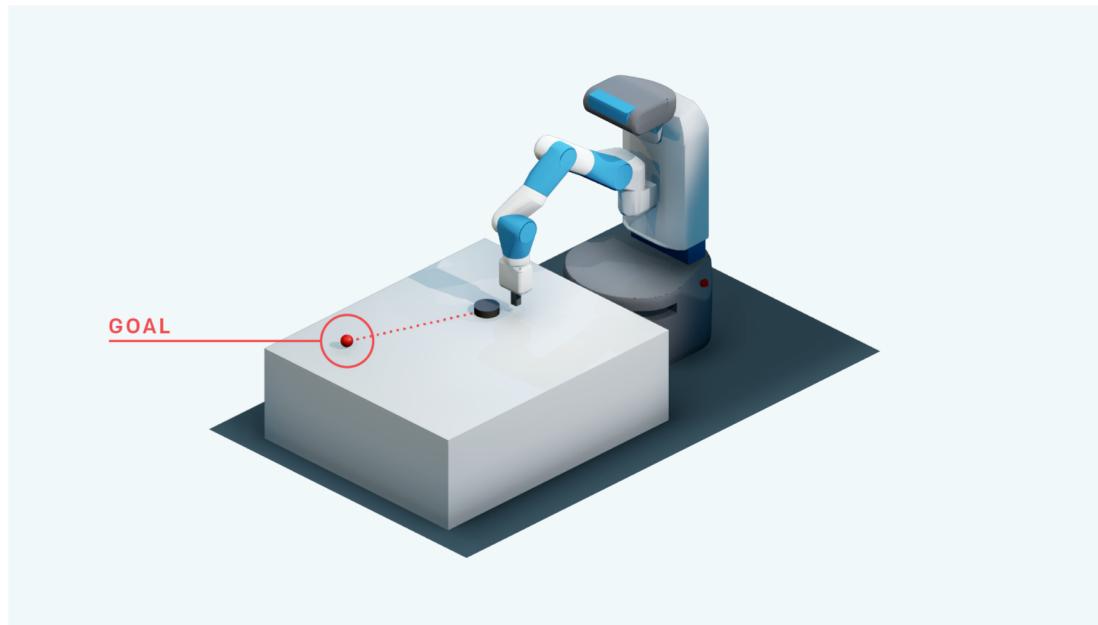
**LLL READING GROUP @ CMU**

# Motivation

---

## Dealing with sparse rewards

- How can we evade the need for complicated reward engineering?



# Problem Statement

---

In the many-goals setting, how can we achieve sample-efficient learning even if the reward signal is binary and sparse?

## Assumptions

- We use an off-policy RL algorithm
- The goals influence the agent’s actions, but not the environment dynamics

# Proposed Method

---

## Hindsight Experience Replay (a.k.a. HER)

### Key idea

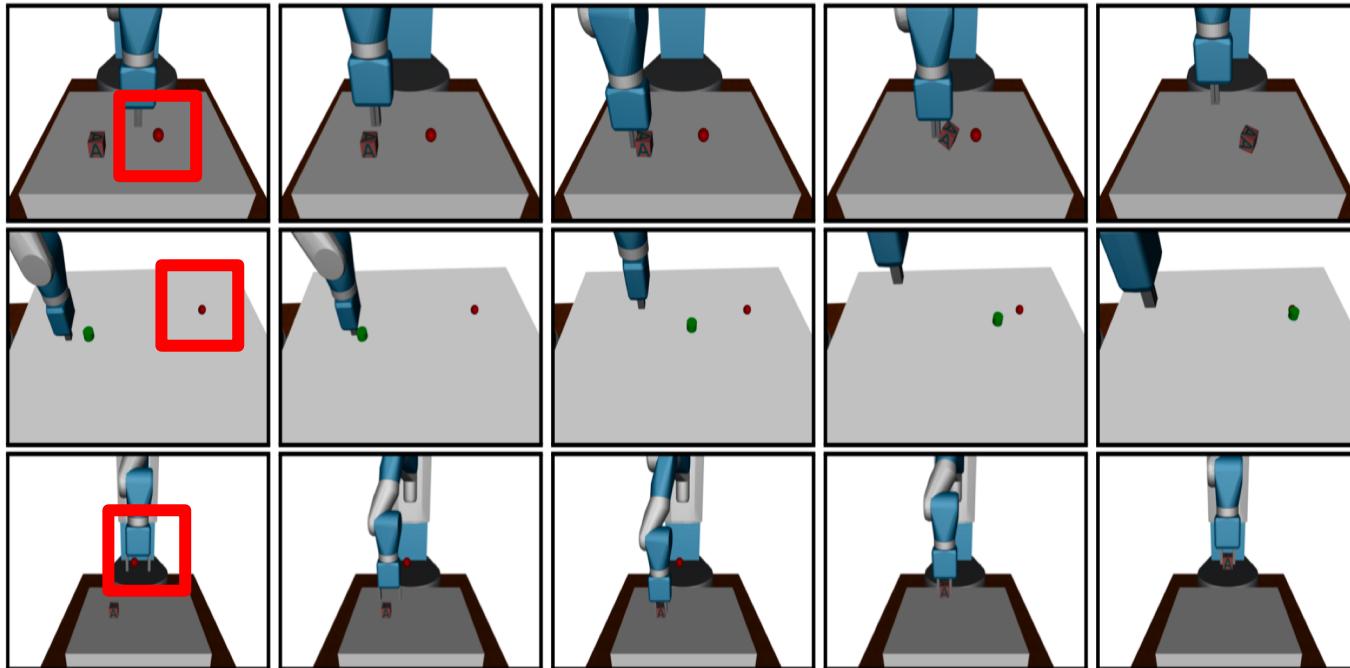
- Replay each episode with a different goal than the one the agent was trying to achieve, e.g, one of the goals already achieved in the episode (Hindsight)

### Hypothesis

- Replying with a different goal plays the role of implicit curriculum in a way to facilitate learning



# Experimental Setting



Pushing

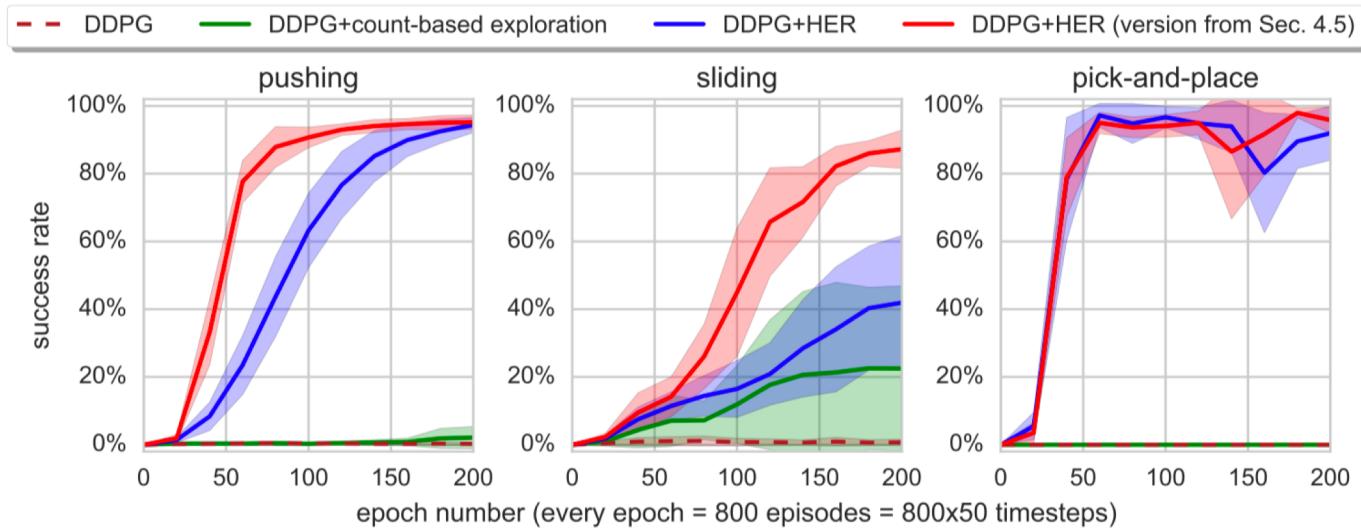
Sliding

Pick-and-place

State	Goal	Reward
Physical info	Desired position of the object	$f_g(s) = [  g - s_{obj}  \leq \epsilon ]$

# Result

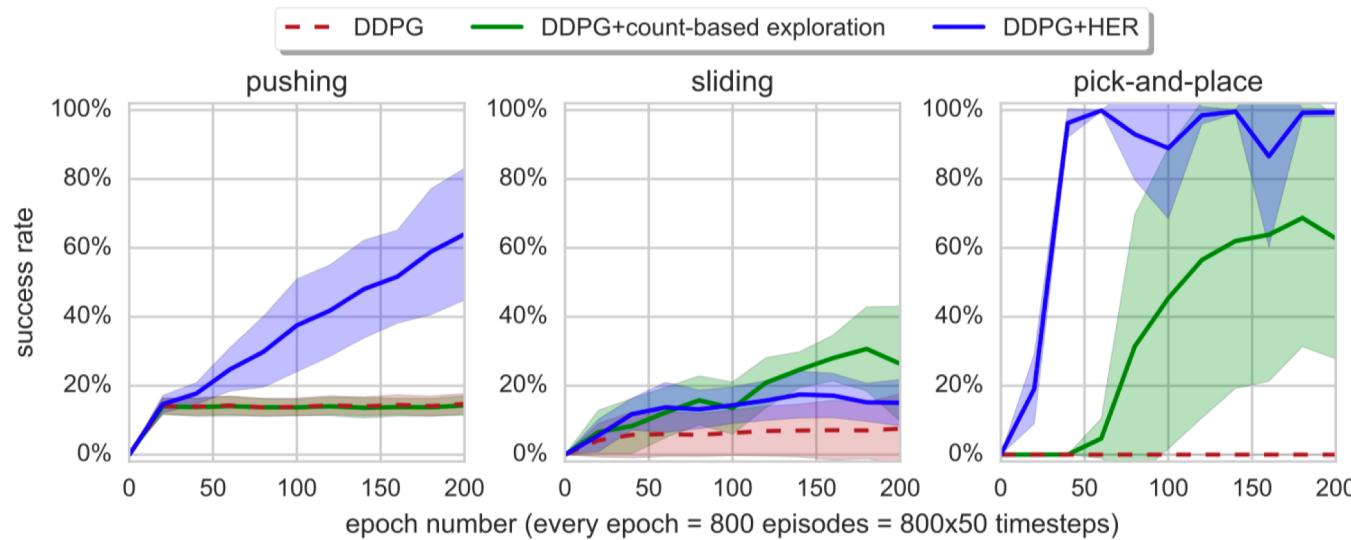
Does HER improve performance?



- Yes. DDPG+HER improves over baselines while baselines fail to solve (Pushing, Pick-and-Place) or makes slower progress (Sliding)

# Result

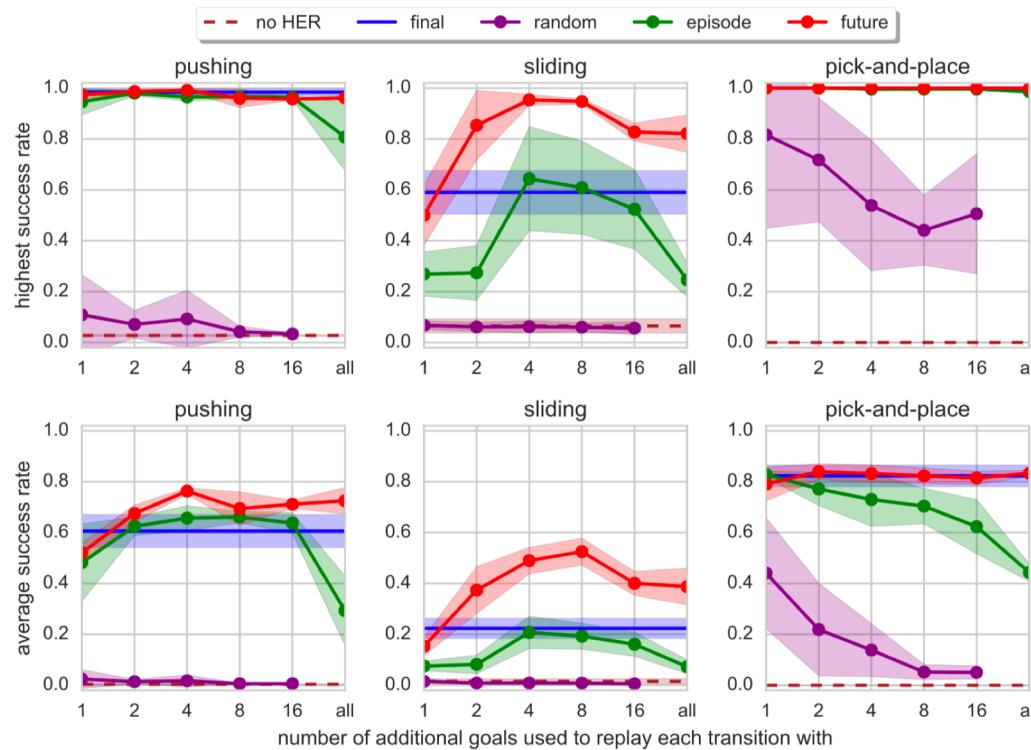
Does HER improve performance **even if there is a single goal we care about?**



- Yes. DDPG+HER is better than DDPG while DDPG fails to solve (Pushing, Pick-and-Place) or makes slower progress (Sliding)

# Result

How many goals should we replay each trajectory with and how to choose them ?



# Result

---

Trained policy in the simulator solves the pick-and-place task well on the physical robot without any finetuning



# Conclusion

---

HER enables sample-efficient learning using only sparse and binary rewards

We can combine HER with arbitrary off-policy RL algorithms

As a result, HER successfully learns complicated behavior three challenging tasks that vanilla RL algorithm fails to solve

- Note that it is the first work to solve these tasks only with binary and sparse reward

# **Universal Value Function Approximators**

---

**PRESENTER: BYUNGSOO JEON**

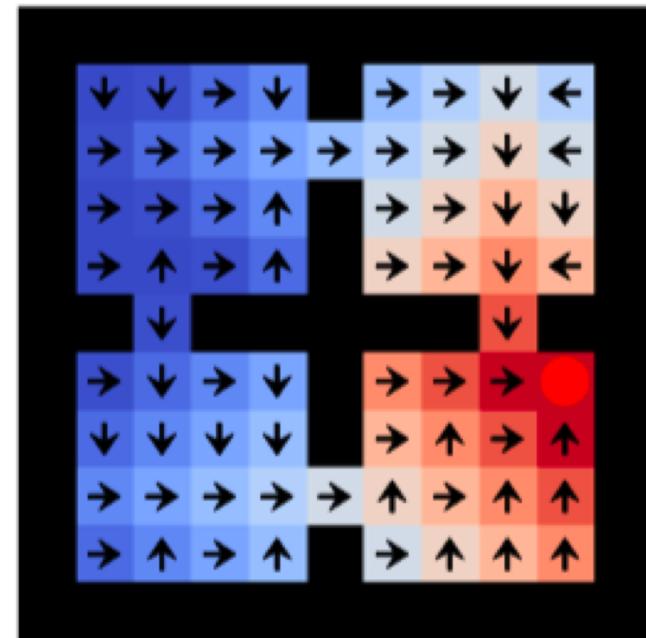
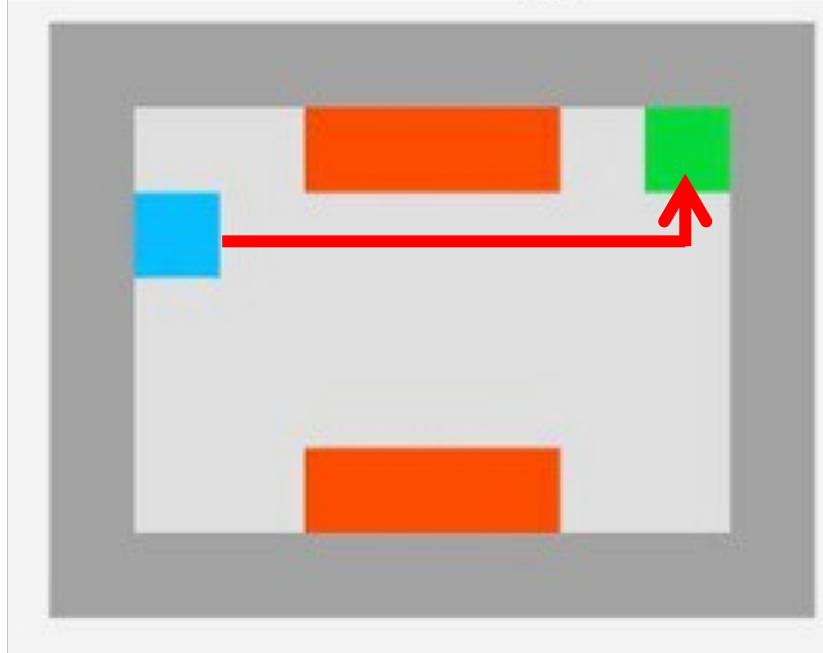
**LLL READING GROUP @ CMU**

# Motivation

---

Dealing with sparse rewards

- How can we exploit the structure in the goal space?



# Problem Statement

---

In the many-goals setting, how can we achieve sample-efficient learning by **exploiting the structure in the goal space?**

## Assumptions

- The goal space contains enough structural information to help generalization to the value of similar, unseen goals
  - Similarity encoded a priori in the goal representation
  - The structure in the induced value functions

# Proposed Method

---

## Universal Value Function Approximators (UVFA)

### Key idea

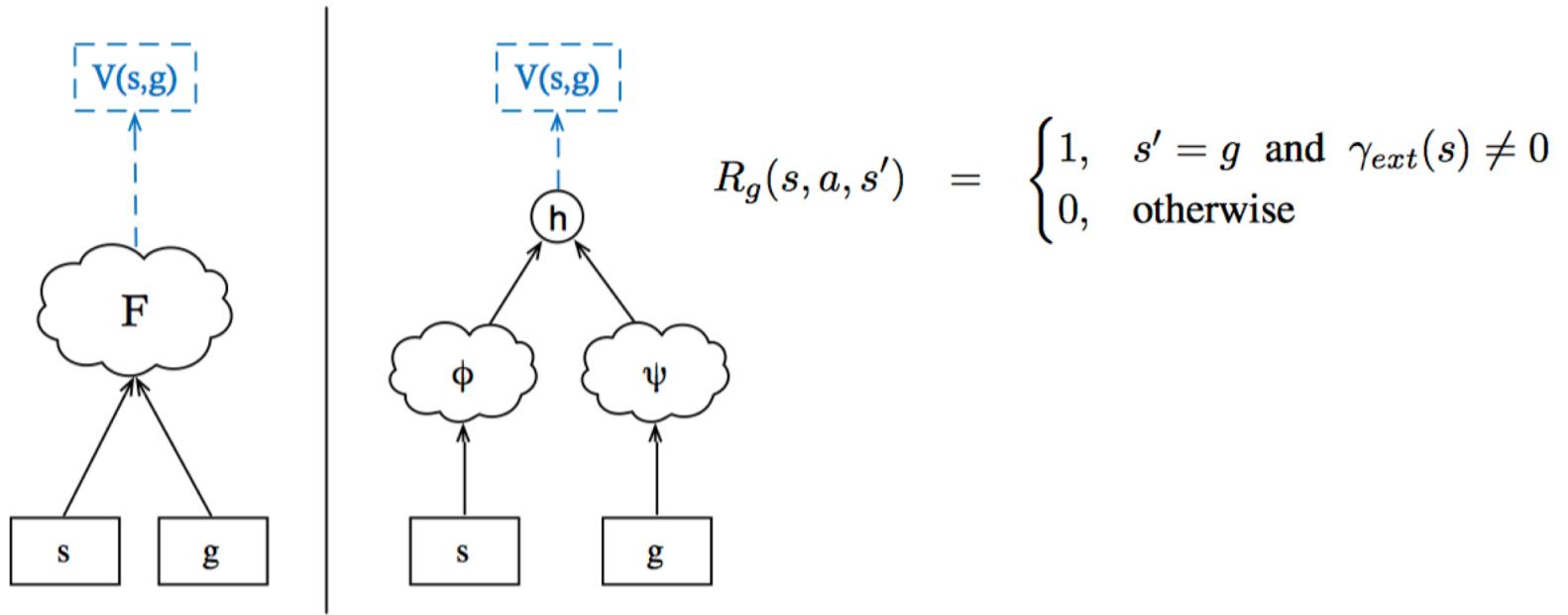
- Extend the idea of value function approximation  $V_\theta(s)$  to both state and goal,  $V_\theta(s, \mathbf{g})$

### Hypothesis

- UVFA is capable of exploiting the structural information in the goal space
- UVFA allows the generalization to unseen states and goals

# Proposed Method

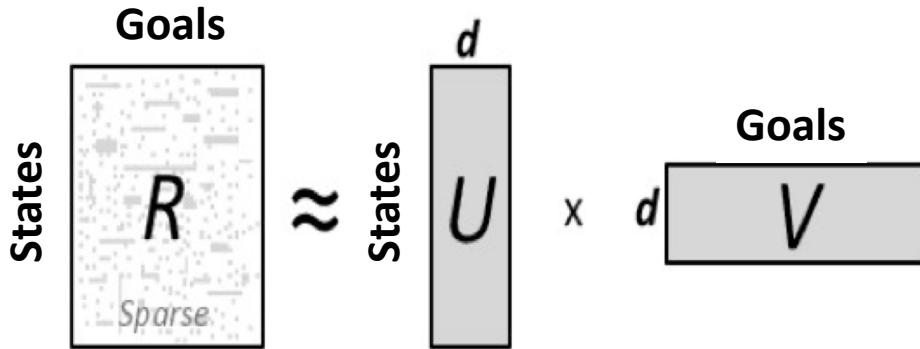
---



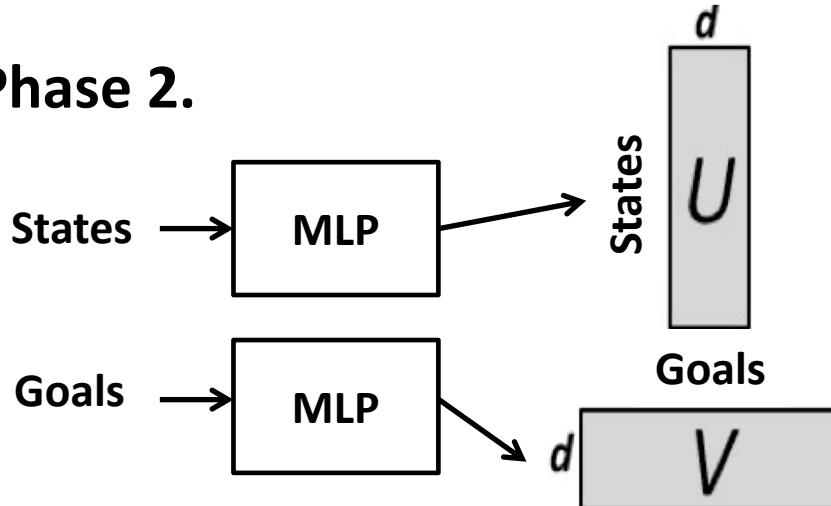
Goal: Represent a large set of optimal value functions by a single, unified function approximator that generalize over both states and goals

# Proposed Method

## Phase 1.



## Phase 2.



**Algorithm 1** UVFA learning from Horde targets

```

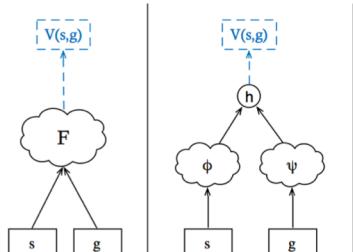
1: Input: rank  $n$ , training goals  $\mathcal{G}_T$ , budgets  $b_1, b_2, b_3$ 
2: Initialise transition history  $\mathcal{H}$ 
3: for  $t = 1$  to  $b_1$  do
4:    $\mathcal{H} \leftarrow \mathcal{H} \cup (s_t, a_t, \gamma_{text}, s_{t+1})$ 
5: end for
6: for  $i = 1$  to  $b_2$  do
7:   Pick a random transition  $t$  from  $\mathcal{H}$ 
8:   Pick a random goal  $g$  from  $\mathcal{G}_T$ 
9:   Update  $Q_g$  given a transition  $t$ 
10: end for
11: Initialise data matrix  $\mathbf{M}$ 
12: for  $(s_t, a_t, \gamma_{text}, s_{t+1})$  in  $\mathcal{H}$  do
13:   for  $g$  in  $\mathcal{G}_T$  do
14:      $\mathbf{M}_{t,g} \leftarrow Q_g(s_t, a_t)$ 
15:   end for
16: end for
17: Compute rank- $n$  factorisation  $\mathbf{M} \approx \phi^\top \psi$ 
18: Initialise embedding networks  $\phi$  and  $\psi$ 
19: for  $i = 1$  to  $b_3$  do
20:   Pick a random transition  $t$  from  $\mathcal{H}$ 
21:   Do regression update of  $\phi(s_t, a_t)$  toward  $\hat{\phi}_t$ 
22:   Pick a random goal  $g$  from  $\mathcal{G}_T$ 
23:   Do regression update of  $\psi(g)$  toward  $\hat{\psi}_g$ 
24: end for
25: return  $Q(s, a, g) := h(\phi(s, a), \psi(g))$ 

```

# Result

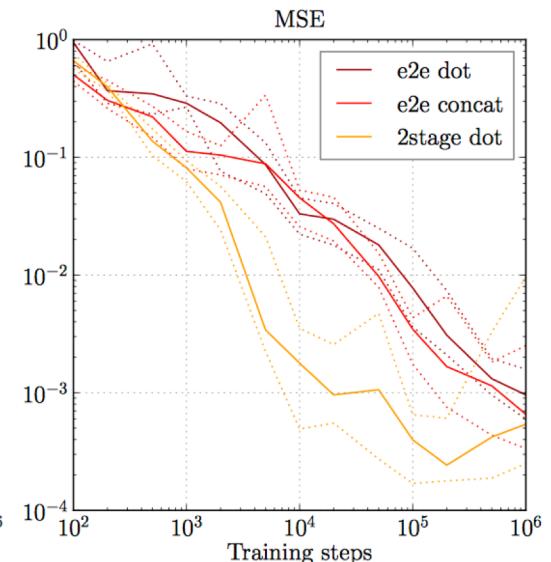
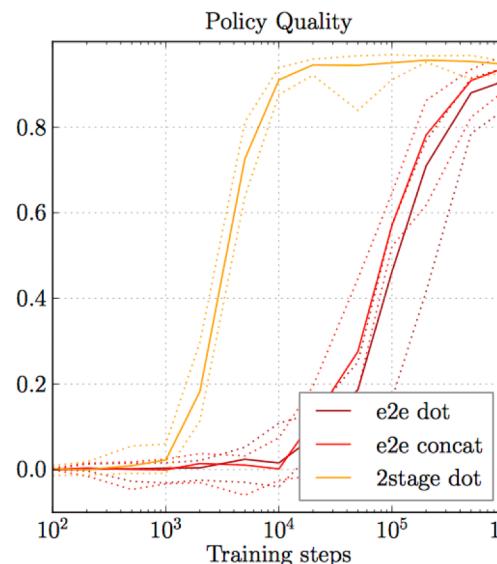
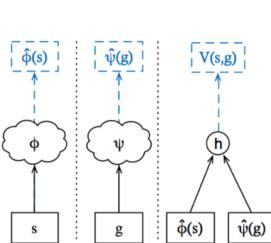
Does two-stream architecture for UVFA outperform other baselines?

e2e concat



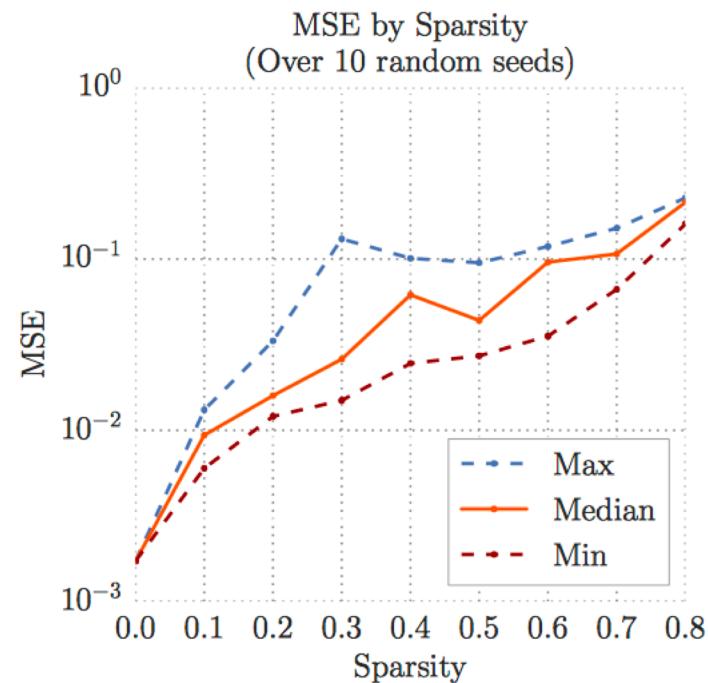
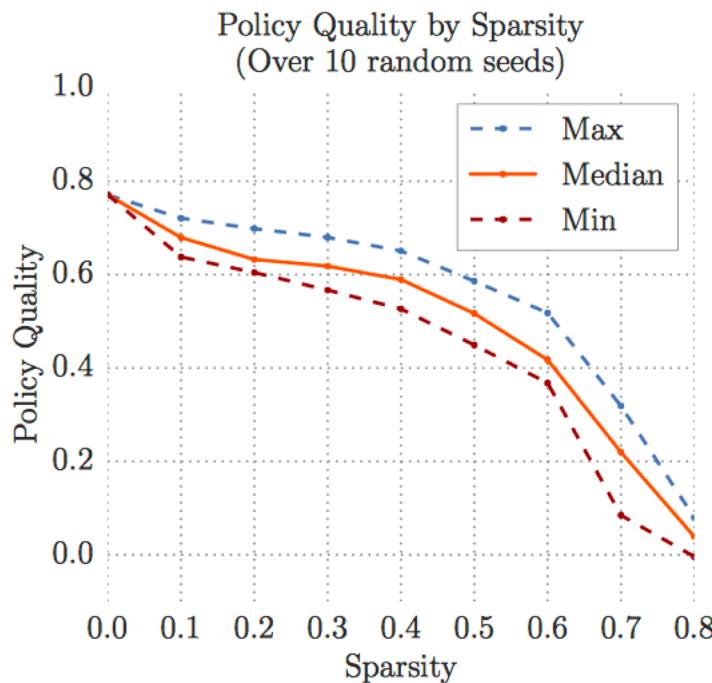
e2e dot

2stage dot



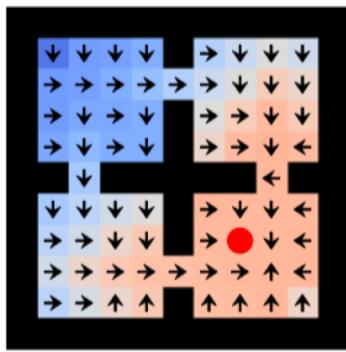
# Result

How robustly does UVFA **generalizes** to the **unseen states** and **goals** across different **sparsity** of the value table?

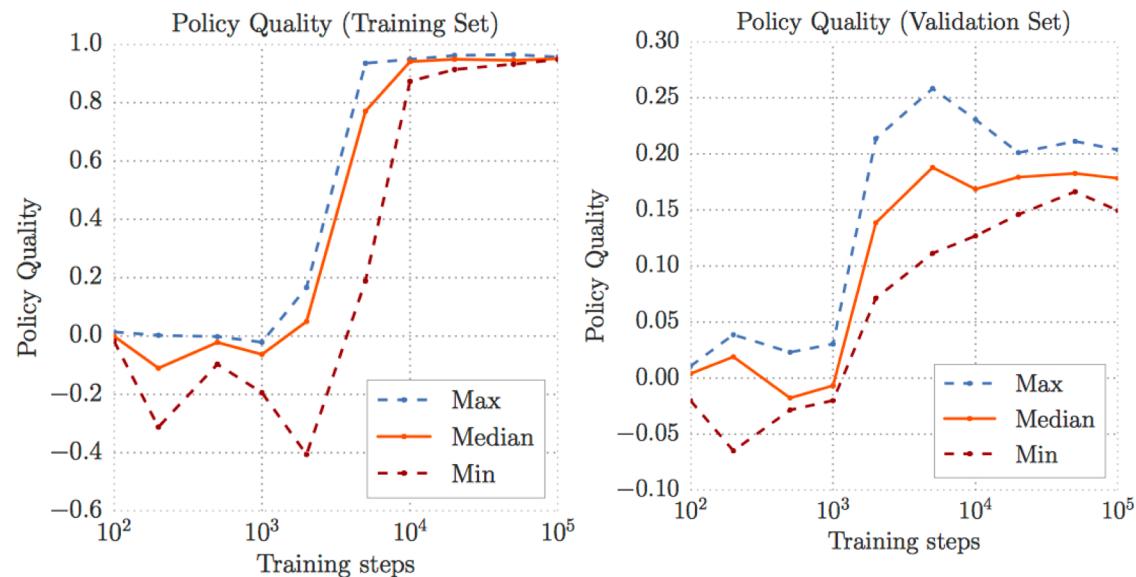
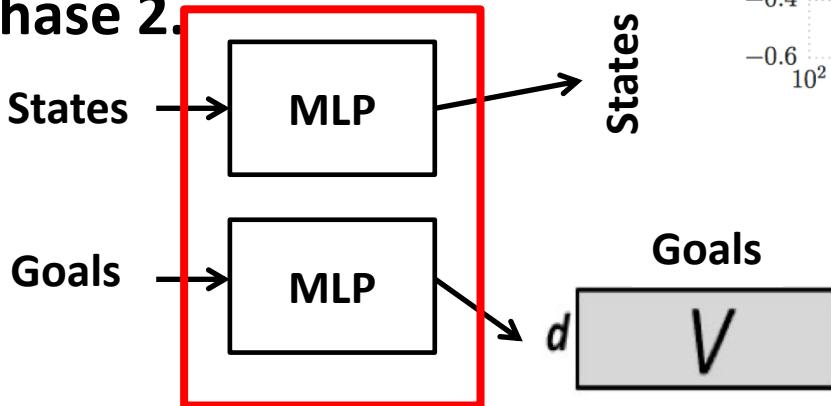


# Result

How well does UVFA **interpolate** the values for unseen goals?

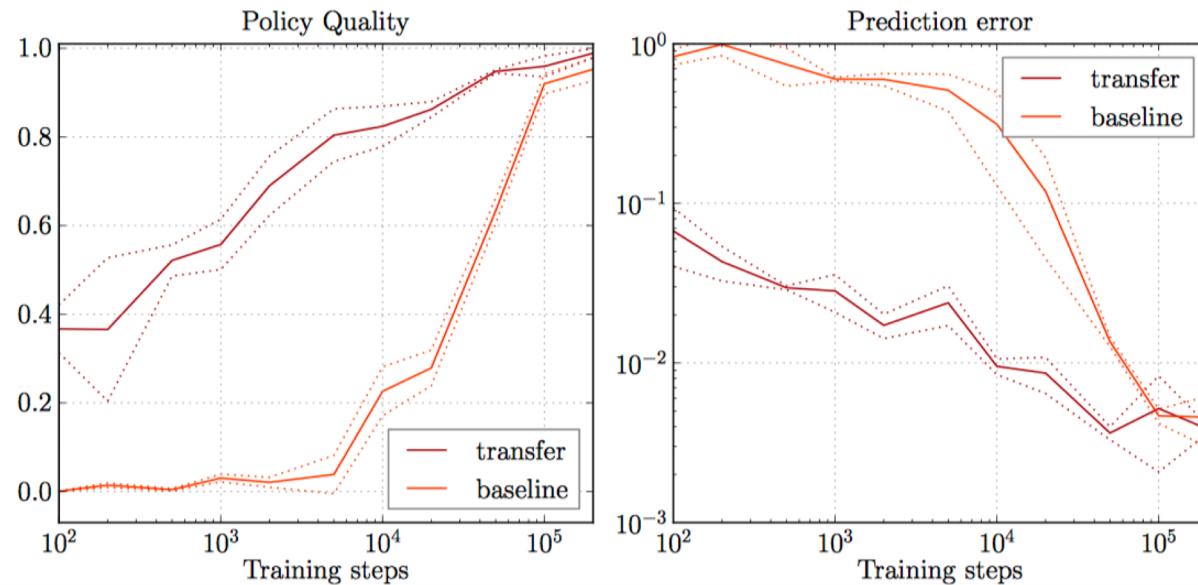
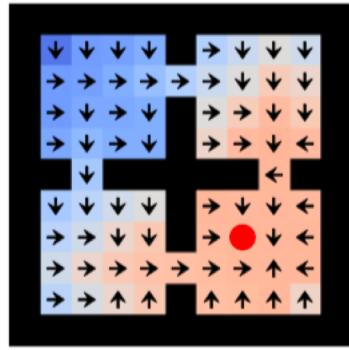


Phase 2.



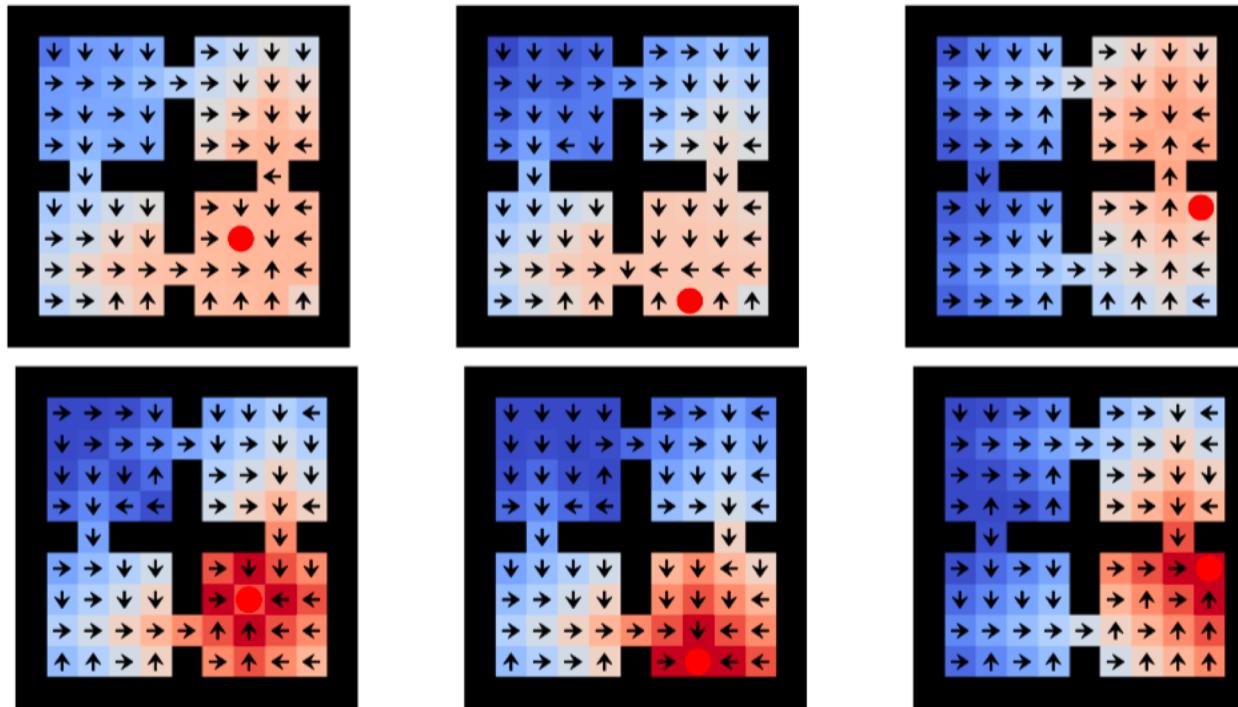
# Result

Can UVFA be used for transfer learning to new tasks with the same dynamics, but different goals?



# Result

How well does UVFA **extrapolates** the values for **unseen goal**?  
(e.g. trained only on the goals except bottom-right room)



# Conclusion

---

UVFA can successfully generalize to previously unseen goals and states

- For both interpolation and extrapolation
- For both supervised and reinforcement learning settings

UVFAs can be used for **transfer learning** to new tasks with the same dynamics, but different goals

Generalized value functions can be used as **features** to represent state (predictive representation)

# References

---

[MFA17] Marcin Andrychowicz et al., “Hindsight Experience Replay”, NIPS 2017

[TDK15] Tom Schaul et al., “Universal Value Function Approximators”, ICML 15

# Question or Discussion

---

