# Synthesized Policies for Transfer and Adaptation across Tasks and Environments

Hexiang Hu*[1], Liyu Chen*[1], Boqing Gong[2], Fei Sha[1,3]

[1]University of Southern California, [2]Tencent AI Lab, [3]Netfilx
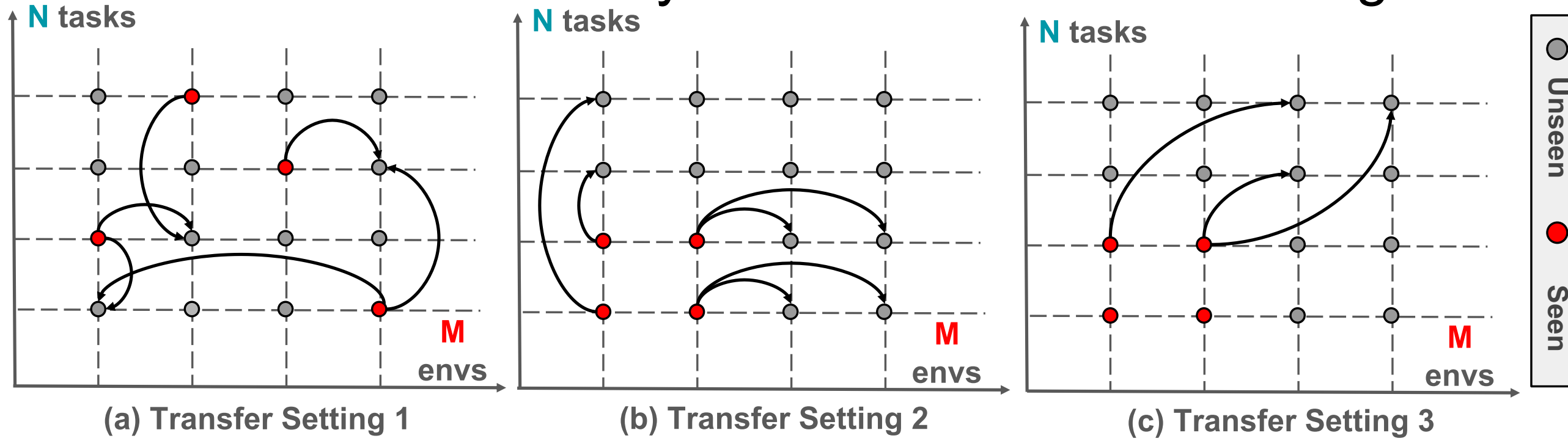
Poster #155

**NeurIPS**
**NETFLIX**

## Highlights:

❖ Study three progressively difficult transfer settings where an agent needs to transfer and adapt across both environments ($\varepsilon$) and tasks ($\tau$) simultaneously.

❖ Novel architecture of policy and reward factorization and disentanglement objective.

❖ We experiment on two simulators, **GridWorld**  **Thor [1]** and our approach has achieved superior performances under different transfer settings

## Problem Setting:

We consider adaption and transfer across environments and tasks simultaneously in Reinforcement Learning.



(a) Transfer Setting 1  (b) Transfer Setting 2  (c) Transfer Setting 3

○ Unseen  ● Seen

**Goal:** Learn $\mathcal{O}(|\mathcal{E}| + |\mathcal{T}|)$ combos and generalize to $|\mathcal{E}| \times |\mathcal{T}|$
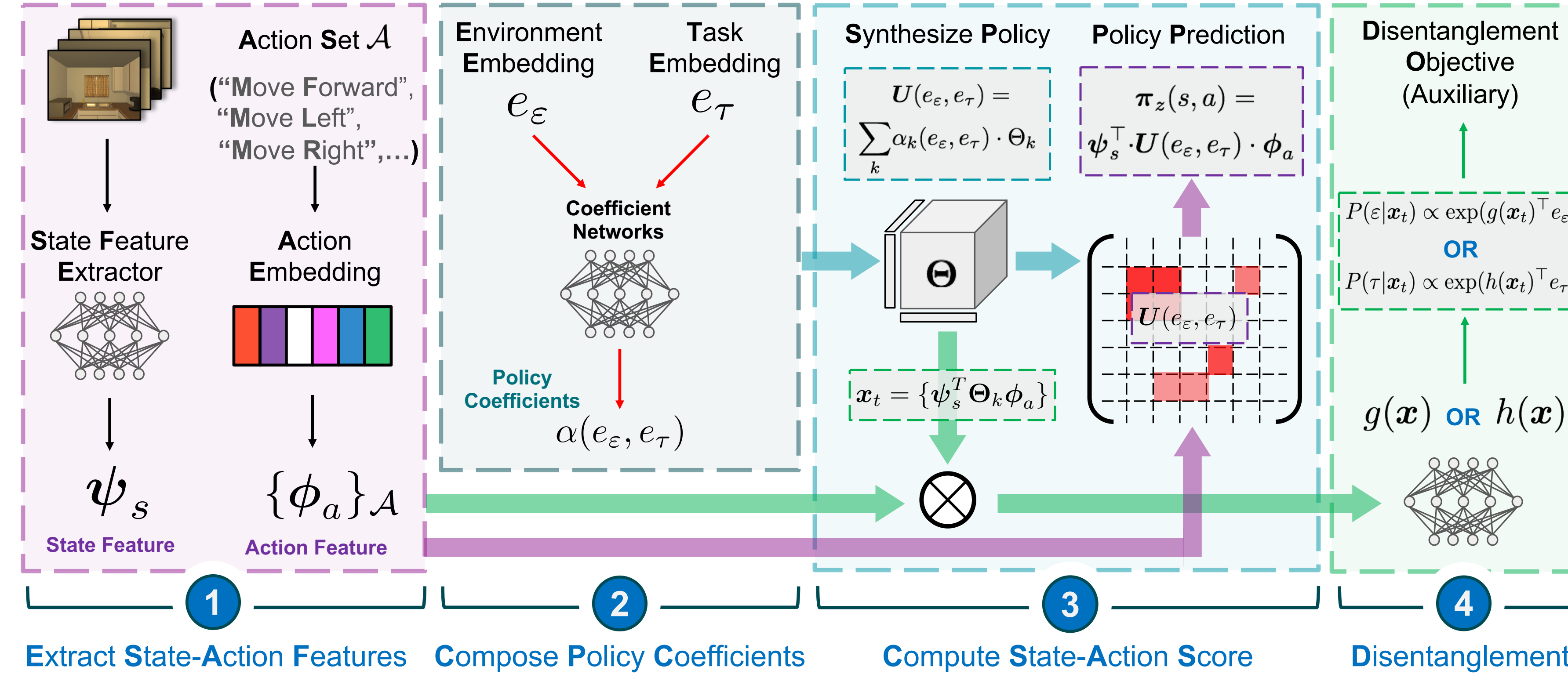
**(a):** Compositional generalization to novel combos

**(b) & (c):** "Incremental transfer" or "Learn a big jump"

## References:

[1] Kolve, Eric, et al. "AI2-THOR: An interactive 3d environment for visual AI." arXiv preprint arXiv:1712.05474 (2017).
[2] Barreto, André, et al. "Successor features for transfer in reinforcement learning." Advances in neural information processing systems. 2017.
[3] Devin, Coline, et al. "Learning modular neural network policies for multi-task and multi-robot transfer." Robotics and Automation (ICRA), 2017, 2017 IEEE International Conference on. IEEE, 2017.

**More about SynPo!**

---

**① Extract State-Action Features**

Action Set $\mathcal{A}$ ("Move Forward", "Move Left", "Move Right",...)

State Feature Extractor   Action Embedding

$\psi_s$   $\{\phi_a\}_{\mathcal{A}}$

State Feature   Action Feature

**② Compose Policy Coefficients**

Environment Embedding   Task Embedding

$e_\varepsilon$   $e_\tau$

Coefficient Networks

Policy Coefficients

$\alpha(e_\varepsilon, e_\tau)$

**③ Compute State-Action Score**

Synthesize Policy

$U(e_\varepsilon, e_\tau) = \sum_k \alpha_k(e_\varepsilon, e_\tau) \cdot \Theta_k$

$\Theta$

$x_t = \{\psi_s^T \Theta_k \phi_a\}$

Policy Prediction

$\pi_z(s,a) = \psi_s^\top \cdot U(e_\varepsilon, e_\tau) \cdot \phi_a$

$U(e_\varepsilon, e_\tau)$

**④ Disentanglement**

Disentanglement Objective (Auxiliary)

$P(\varepsilon|x_t) \propto \exp(g(x_t)^\top e_\varepsilon)$

OR

$P(\tau|x_t) \propto \exp(h(x_t)^\top e_\tau)$

$g(x)$ OR $h(x)$

## Main Idea:

(1) Learn a policy basis $\Theta$ to compose $(\varepsilon, \tau)$ specific policies $U(e_\varepsilon, e_\tau)$

(2) Learn low dimensional embeddings $e_\varepsilon$ or $e_\tau$ for novel $\varepsilon$ and $\tau$
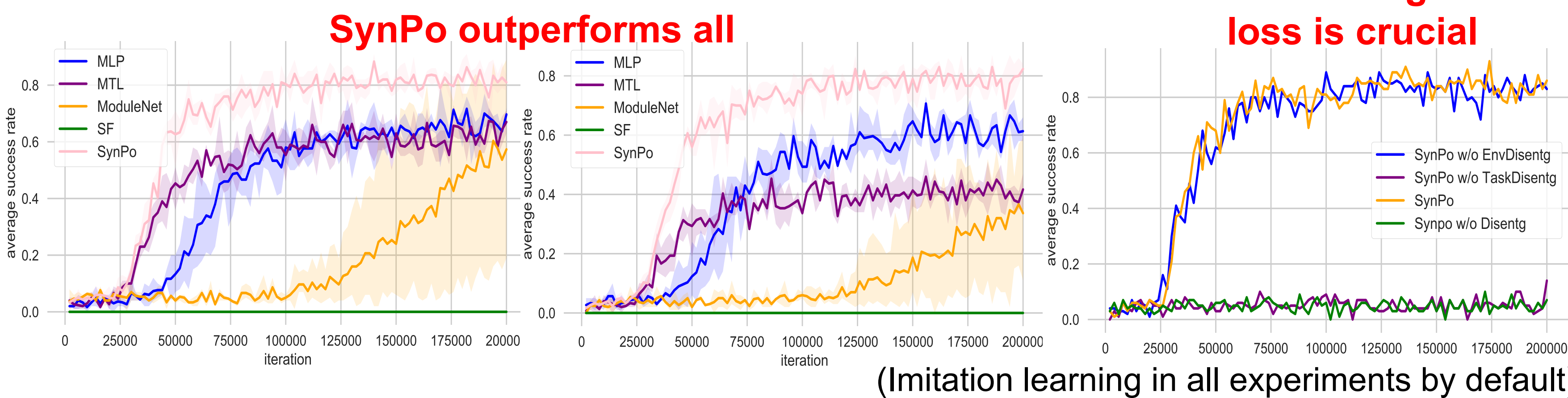
**Policy Synthesis:** $U(e_\varepsilon, e_\tau) = \sum_k \alpha_k(e_\varepsilon, e_\tau) \cdot \Theta_k$

**Policy Prediction:** $\pi_z(a,s) \propto \exp(\psi_s^T U(e_\varepsilon, e_\tau)\phi_a + b_\pi)$

**Reward Prediction:** $\tilde{r}_z(s,a) = \psi_s^T V(e_\varepsilon, e_\tau)\phi_a + b_r$

$V(e_\varepsilon, e_\tau) = \sum_k \beta_k(e_\varepsilon, e_\tau) \cdot \Theta_k$

**Disentanglement Objective:**

$l_\varepsilon = -\sum_t \log P(\varepsilon|x_t), \text{with } P(\varepsilon|x_t) \propto \exp(g(x_t)^T e_\varepsilon)$

$l_\tau = -\sum_t \log P(\tau|x_t), \text{with } P(\tau|x_t) \propto \exp(h(x_t)^T e_\tau)$

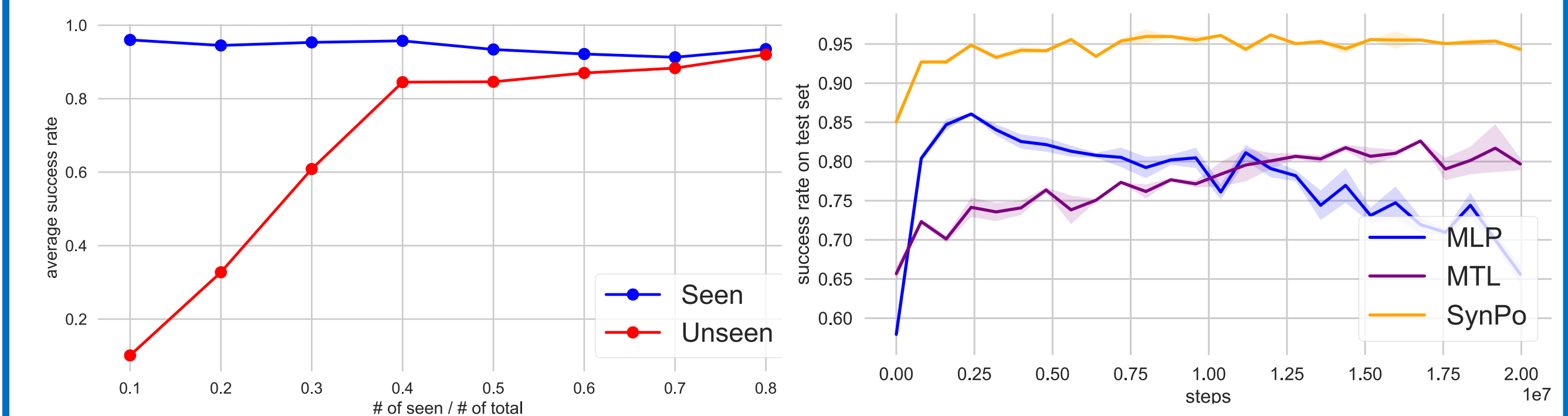## Experiments on GridWorld:

### Results for Transfer Setting 1

**SynPo outperforms all**

**Task disentanglement loss is crucial**



(Imitation learning in all experiments by default)

---

**How many seen $(\varepsilon, \tau)$ pairs are needed to transfer well?**

**40% of the training $(\varepsilon, \tau)$ pairs can generalize**
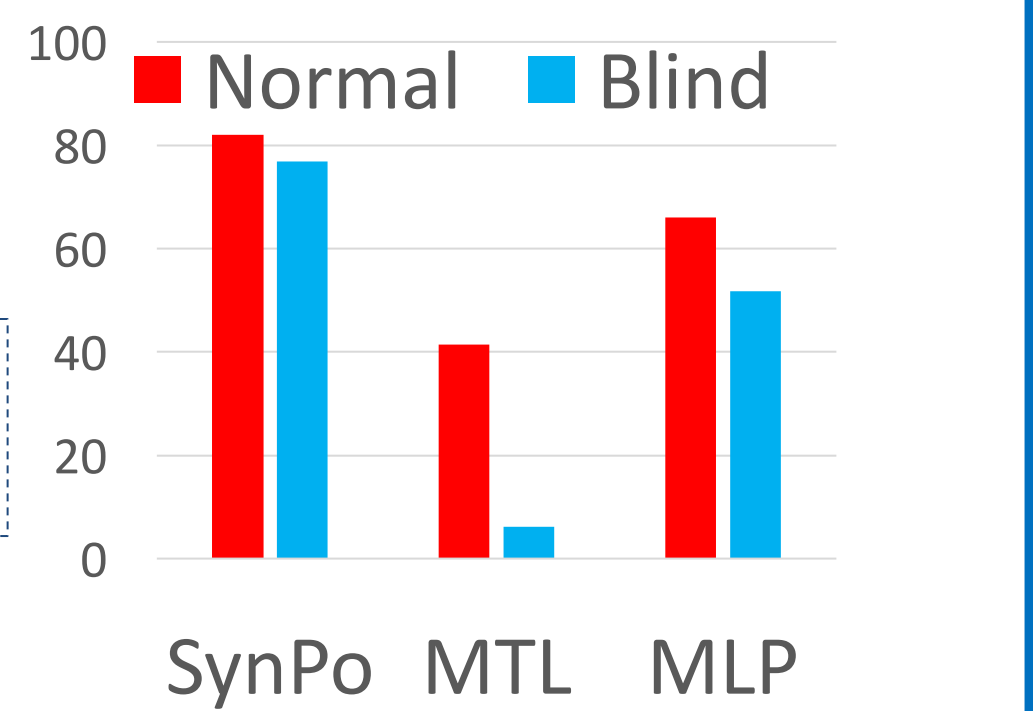
**Does reinforcement learning help transfer?**

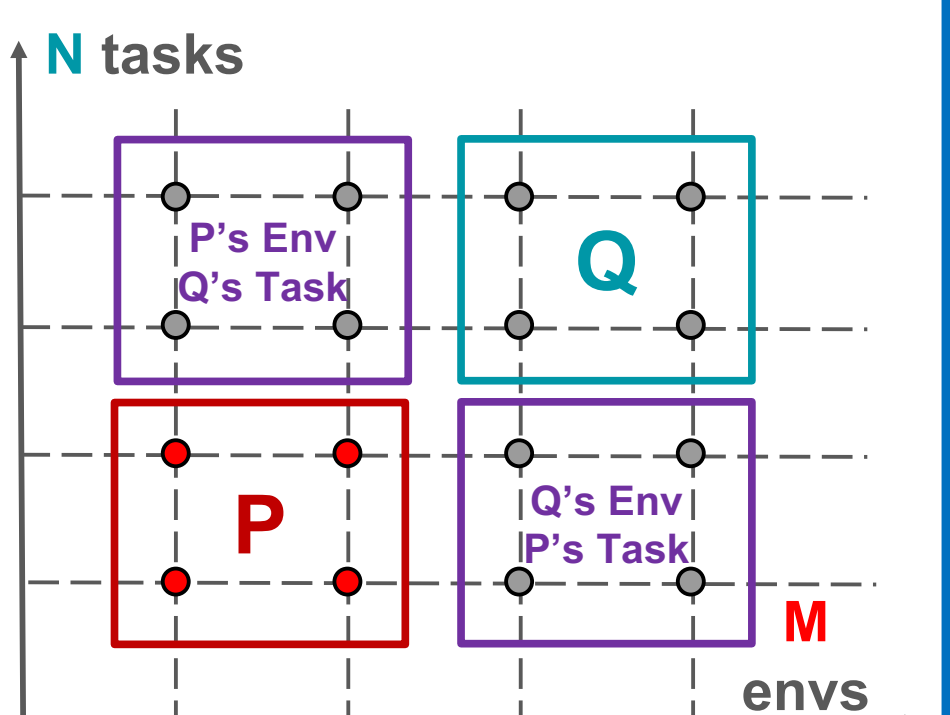**Fine tuning with RL on training set helps transfer in test set!**



**SynPo is less sensitive to Imperfect perception.**

**Blind Agent:** the agent can only see its own and treasures' position, **but not the maze**.



### Results for Transfer Setting 2 & 3

| Setting | Method | Q's $\varepsilon$, P's $\tau$ | P's $\varepsilon$, Q's $\tau$ | Q Pairs |
|---------|--------|------|------|------|
| Setting 2 | MLP | 13.8% | 20.7% | 6.3% |
| | SynPo | **50.5%** | **21.5%** | **13.5%** |
| Setting 3 | MLP | 14.6% | 18.3% | 7.2% |
| | SynPo | **42.7%** | **19.4%** | **12.9%** |



## Experiments on THOR [1]:

### Results for Transfer Setting 1

| Split | ModuleNet | MLP | MTL | SynPo |
|-------|-----------|-----|-----|-------|
| SEEN | 51.5% | 47.5% | 52.2% | **55.6%** |
| UNSEEN | 14.4% | 25.8% | 33.3% | **35.4%** |