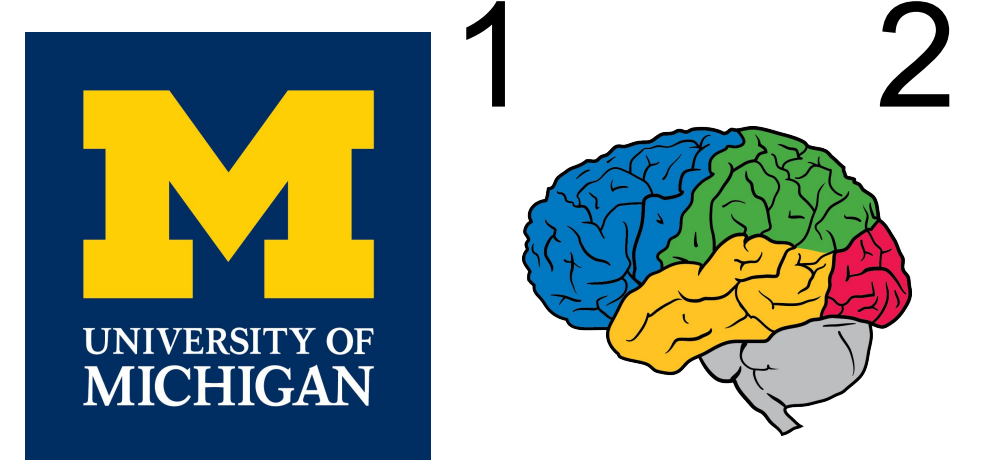


Meta Reinforcement Learning with Autonomous Inference of Subtask Dependencies



Sungryull Sohn¹ Hyunjae Woo¹ Jongwook Choi¹ Honglak Lee^{1, 2} srsohn@umich.edu

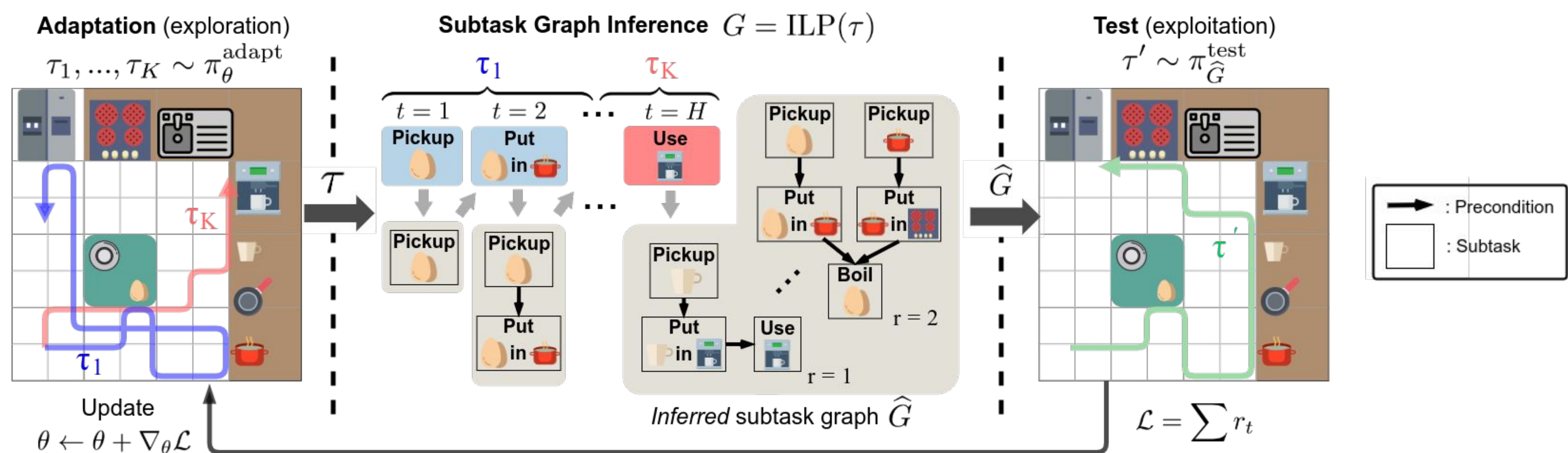
Introduction

Summary

- **Goal:** Solve a meta RL problem by *explicitly inferring the underlying task structure* in order to quickly adapt to the task
- **Idea:** Infer the subtask reward and dependencies in adaptation phase, and solve the task in test phase
- **Result:** Better generalization and performance on complex tasks

Subtask graph inference problem

- A task is characterized by the **subtask graph** which defines the subtasks, their *rewards* and *preconditions*.
- The subtask graph is **unknown** to the agent
- Few-shot RL: after K episodes of adaptation, the agent needs to maximize the return in test phase.



Method

Meta-learner with Subtask Graph Inferencer (MSGI)

- **Adaptation:** Adaptation policy ($\pi_{\theta}^{\text{adapt}}$) learns to efficiently collect experience to accurately infer the task
- **Inference:** ILP algorithm infers the subtask reward and precondition in the form of subtask graph $\hat{G} = (\hat{G}_r, \hat{G}_c)$
- **Execution:** execution policy ($\pi_{\hat{G}}^{\text{test}}$) executes the inferred subtask graph \hat{G}

Training

Adaptation policy

- **UCB-like intrinsic bonus:** $r_t^{\text{UCB}} = w_{\text{UCB}} \cdot \mathbb{I}(\mathbf{x}_t \text{ is novel})$
- Trained via actor-critic method with GAE [1]

Subtask graph inference

- **Precondition (G_c):** CART-based logic induction algorithm to infer the precondition function ($e = \tilde{f}_{G_c}(\mathbf{x})$)
- **Subtask reward (G_r):** $G_r^i \sim \mathcal{N}(\hat{\mu}^i, \hat{\sigma}^i)$, $\hat{G}_r^{\text{MLE}, i} = \hat{\mu}_{\text{MLE}}^i = \mathbb{E}[r_t | \mathbf{o}_t = \mathcal{O}^i]$

Experiment

Setting

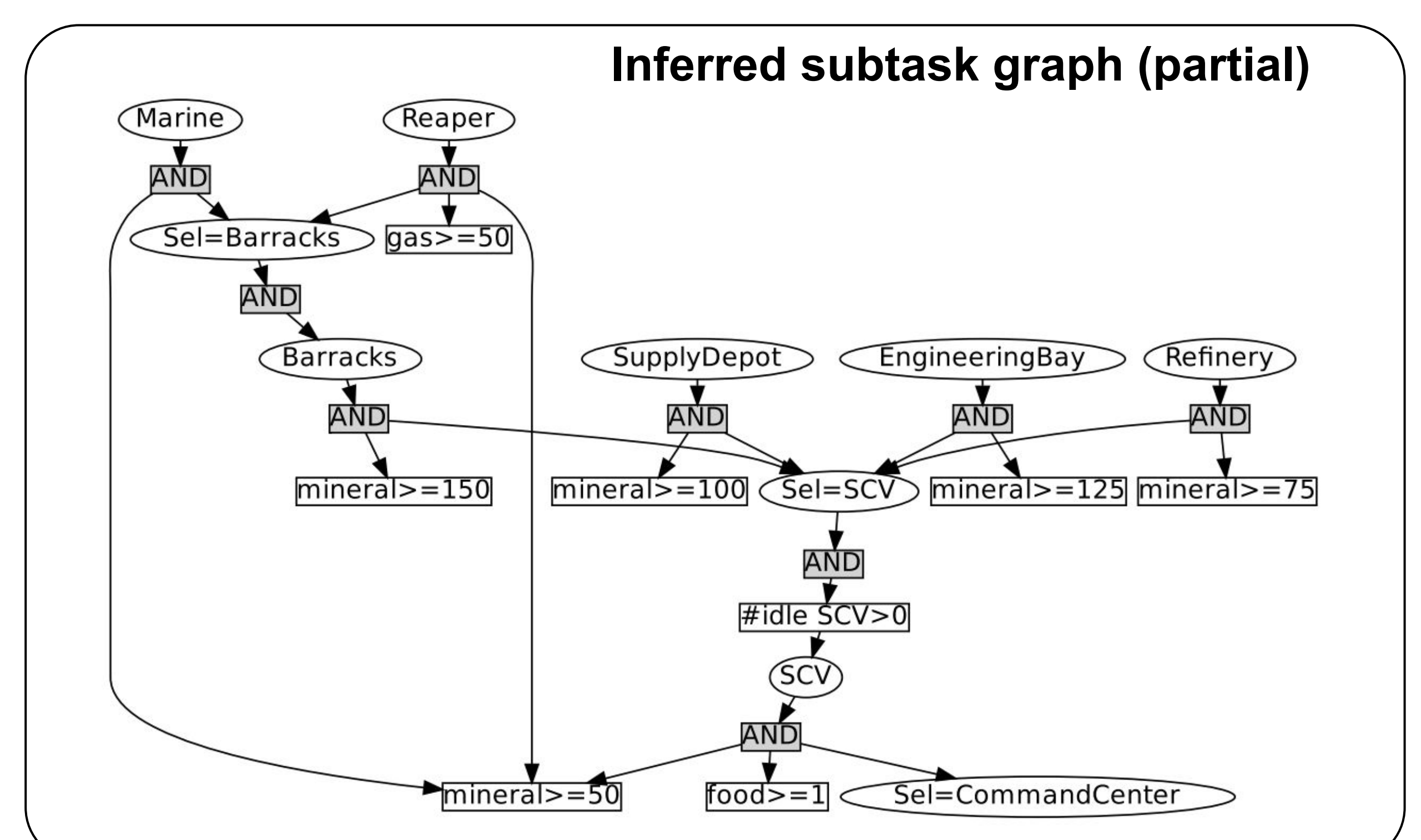
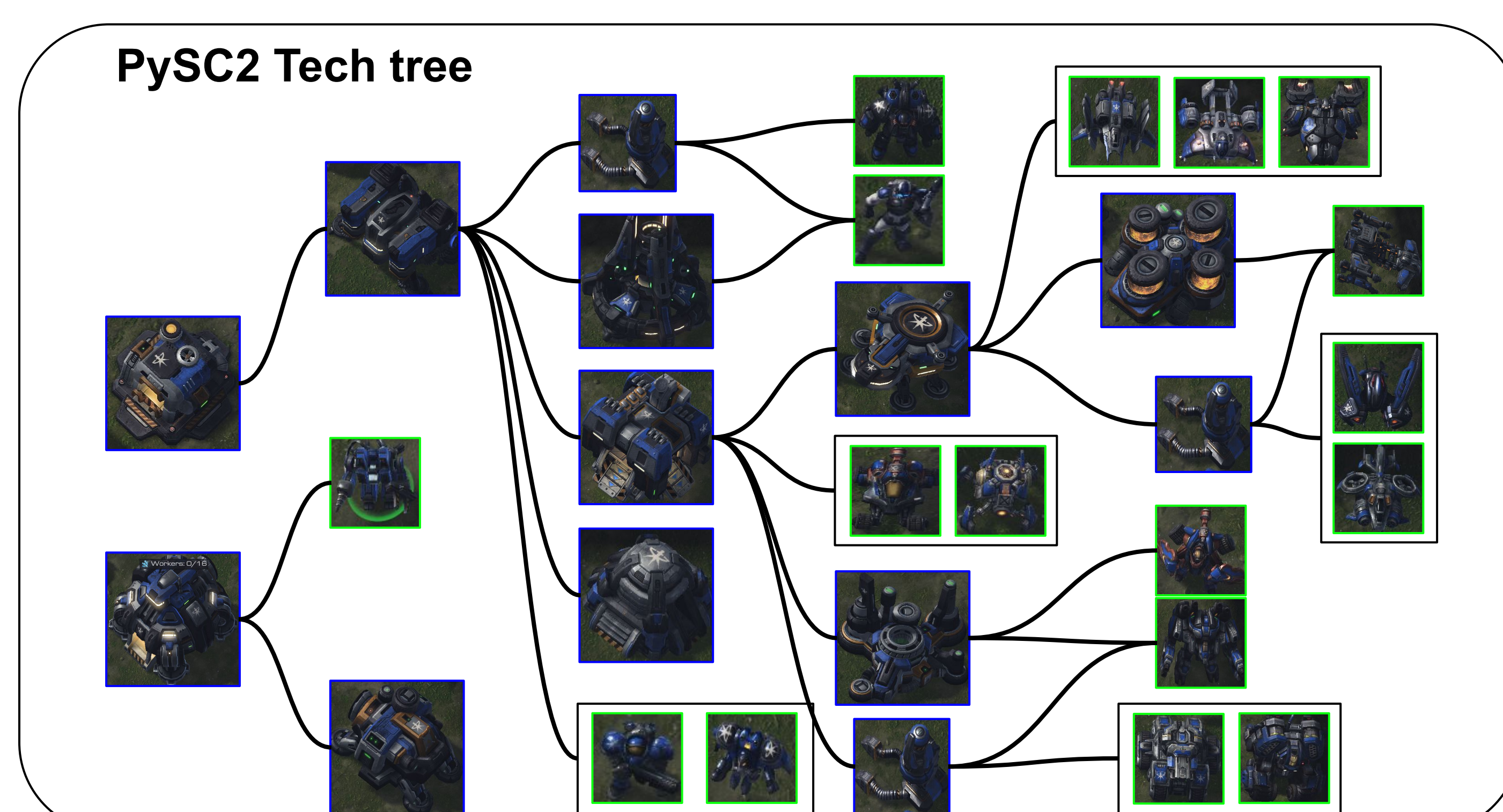
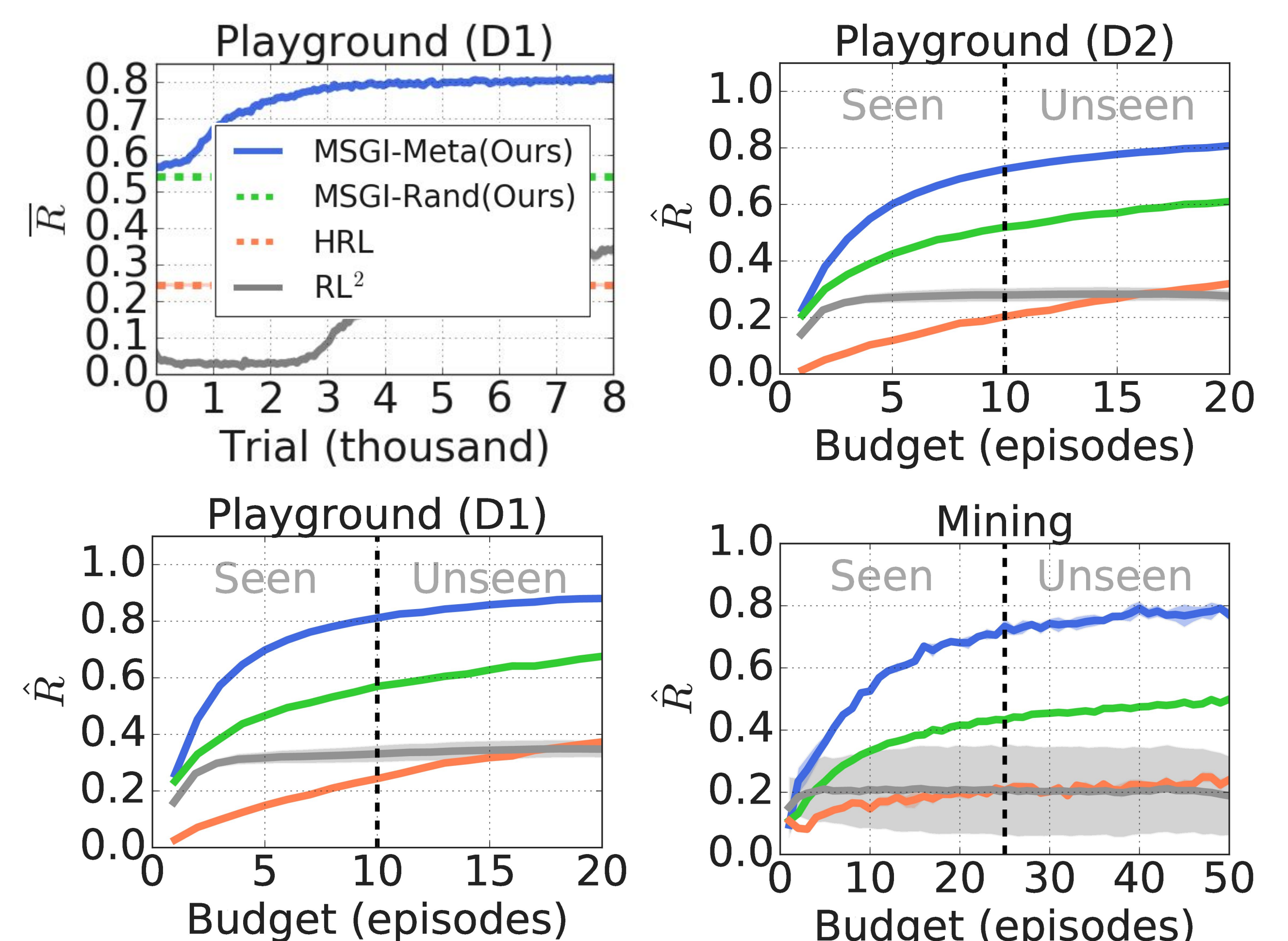
- Domain: Playground [2], Mining [2], SC2LE [3]
- Agents: Random, HRL [4], RL² [5], MSGI-random (ours), MSGI-meta (ours)

Result on Playground and Mining domain

- Adaptation: our learned adaptation policy enables faster adaptation (MSGI-meta V.S. MSGI-random)
- Generalization: MSGI (ours) consistently outperforms baselines on unseen tasks with **longer (unseen)** adaptation horizon

Result on StarCraft II (SC2LE)

- Subtask graph inference: given 20 episodes budget with 2.5K environment steps, our method correctly infers **93%** of the preconditions



[1] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. arXiv preprint arXiv:1506.02577, 2015.

[2] Sungryull Sohn, Junhyuk Oh, and Honglak Lee. Hierarchical reinforcement learning for zero-shot generalization with subtask dependencies. In NeurIPS, pp. 7156–7166, 2018.

[3] Vinyals, Oriol, et al. Starcraft ii: A new challenge for reinforcement learning. arXiv, 2017.

[4] Jacob Andreas et al. Modular multitask reinforcement learning with policy sketches. ICML, 2017.

[5] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RL²: Fast reinforcement learning via slow reinforcement learning. arXiv, 2016.