

META REINFORCEMENT LEARNING WITH AUTONOMOUS INFERENCE OF SUBTASK DEPENDENCIES

Anonymous authors

Paper under double-blind review

ABSTRACT

We propose and address a novel few-shot RL problem, where a task is characterized by a subtask graph which describes a set of subtasks and their dependencies that are unknown to the agent. The agent needs to quickly adapt to the task over few episodes during adaptation phase to maximize the return in the test phase. Instead of directly learning a meta-policy, we develop a *Meta-learner with Subtask Graph Inference* (MSGI), which infers the latent parameter of the task by interacting with the environment and maximizes the return given the latent parameter. To facilitate learning, we adopt an intrinsic reward inspired by upper confidence bound (UCB) that encourages efficient exploration. Our experiment results on two grid-world domains and StarCraft II environments show that the proposed method is able to accurately infer the latent task parameter, and to adapt more efficiently than existing meta RL and hierarchical RL methods.

1 INTRODUCTION

Recently, reinforcement learning (RL) systems have achieved super-human performance on many complex tasks (Mnih et al., 2015; Silver et al., 2016; Van Seijen et al., 2017). However, these works mostly have been focused on a single known task where the agent can be trained for a long time (*e.g.*, Silver et al. (2016)). We argue that agent should be able to solve multiple tasks with varying sources of reward. Recent work in multi-task RL has attempted to address this; however, they focused on the setting where the structure of task are *explicitly* described with natural language instructions (Oh et al., 2017; Andreas et al., 2017; Yu et al., 2017; Chaplot et al., 2018), programs (Denil et al., 2017), or graph structures (Sohn et al., 2018). However, such task descriptions may not readily be available. A more flexible solution is to have the agents infer the task by interacting with the environment. Recent work in Meta RL (Hochreiter et al., 2001; Duan et al., 2016; Wang et al., 2016; Finn et al., 2017) (especially in few-shot learning settings) has attempted to have the agents implicitly infer tasks and quickly adapt to them. However, they have focused on relatively simple tasks with a single goal (*e.g.*, multi-armed bandit, locomotion, navigation, *etc.*).

We argue that real-world tasks often have a hierarchical structure and multiple goals, which require long horizon planning or reasoning ability (Erol, 1996; Xu et al., 2017; Ghazanfari & Taylor, 2017; Sohn et al., 2018). Take, for example, the task of making a breakfast in Figure 1. A meal can be served with different dishes and drinks (*e.g.*, *boiled egg* and *coffee*), where each could be considered as a subtask. These can then be further decomposed into smaller subtask until some base subtask (*e.g.*, *pickup egg*) is reached. Each subtask can provide the agent with reward; if only few subtasks provide reward, this is considered a *sparse reward* problem. When the subtask dependencies are complex and reward is sparse, learning an optimal policy can require a large number of interactions with the environment. This is the problem scope we focus on in this work: learning to quickly infer and adapt to varying hierarchical tasks with multiple goals and complex subtask dependencies.

To this end, we formulate and tackle a new few-shot RL problem called *subtask graph inference* problem, where the task is defined as a factored-MDP (Boutilier et al., 1995; Jonsson & Barto, 2006) with hierarchical structure represented by *subtask graph* (Sohn et al., 2018). The task consists of multiple subtasks, where each subtask gives reward when completed (see Figure 1 for more detail). The complex dependencies between subtasks (*i.e.*, preconditions) enforce agent to execute all the required subtasks *before* it can execute certain subtask. Intuitively, the agent can efficiently solve the task by leveraging the inductive bias of underlying task structure (Section 2.2).

Inspired by the recent works on multi-task and few-shot RL, we propose a meta reinforcement learning approach that explicitly infers the latent structure of the task (*e.g.*, subtask graph). The agent learns its adaptation policy to collect as much information about the environment as possible in order to rapidly and accurately infer the unknown task structure. After that, the agent’s test policy is a

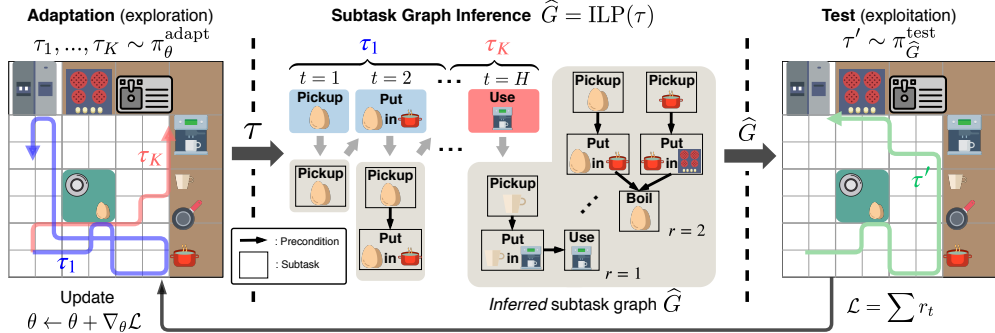


Figure 1: Overview of our method in the context of “prepare breakfast” task. This task can be broken down into subtasks (e.g., *pickup mug*) that composes the underlying subtask graph G . **(Left)** To learn about the unknown task, the agent collects trajectories over K episodes through a parameterized *adaptation* policy $\pi_{\theta}^{\text{adapt}}$ that learns to explore the environment. **(Center)** With each new trajectory, the agent attempts to infer the task’s underlying *ground-truth* subtask graph G with \hat{G} . **(Right)** A separate *test* policy $\pi_{\hat{G}}^{\text{test}}$ uses the inferred subtask graph \hat{G} to produce a trajectory that attempts to maximize the agent’s reward $\sum r_t$ (e.g., the green trajectory that achieves the *boil egg* subtask). The more precise \hat{G} , the more reward the agent would receive, which implicitly improves the *adaptation* policy $\pi_{\theta}^{\text{adapt}}$ to better explore the environment and therefore better infer \hat{G} in return.

contextual policy that takes the inferred subtask graph as an input and maximizes the expected return (See Figure 1). We leverage inductive logic programming (ILP) technique to derive an efficient task inference method based on the principle of maximum likelihood. To facilitate learning, we adopt an intrinsic reward inspired by upper confidence bound (UCB) that encourages efficient exploration. We evaluate our approach on various environments ranging from simple grid-world (Sohn et al., 2018) to StarCraft II (Vinyals et al., 2017). In all cases, our method can accurately infer the latent subtask graph structure, and adapt more efficiently to unseen tasks than the baselines.

The contribution of this work can be summarized as follows:

- We propose a new meta-RL problem with more general and richer form of tasks compared to the recent meta-RL approaches.
- We propose an efficient task inference algorithm that leverages inductive logic programming, which *accurately* infers the latent subtask graph from the agent’s experience data.
- We implement a deep meta-RL agent that *efficiently* infers the subtask graph for faster adaptation.
- We compare our method with other meta-RL agents on various domains, and show that our method adapts more efficiently to unseen tasks with complex subtask dependencies.

2 PROBLEM DEFINITION

2.1 BACKGROUND: FEW-SHOT REINFORCEMENT LEARNING

A *task* is defined by an MDP $\mathcal{M}_G = (\mathcal{S}, \mathcal{A}, \mathcal{P}_G, \mathcal{R}_\tau)$ parameterized by a task parameter G with a set of states \mathcal{S} , a set of actions \mathcal{A} , transition dynamics \mathcal{P}_G , reward function \mathcal{R}_G . In the K -shot RL formulation (Duan et al., 2016; Finn et al., 2017), each *trial* under a fixed task \mathcal{M}_G consists of an *adaptation phase* where the agent learns a task-specific behavior and a *test phase* where the adapted behavior is evaluated. For example, RNN-based meta-learners (Duan et al., 2016; Wang et al., 2016) adapt to a task \mathcal{M}_G by updating its RNN states (or fast-parameters) ϕ_t , where the initialization and update rule of ϕ_t is parameterized by a *slow-parameter* θ : $\phi_0 = g(\theta)$, $\phi_{t+1} = f(\phi_t; \theta)$. Gradient-based meta-learners (Finn et al., 2017; Nichol et al., 2018) instead aim to learn a good initialization of the model so that it can adapt to a new task with few gradient update steps. In the test phase, the agent’s performance on the task \mathcal{M}_G is measured in terms of the return:

$$\mathcal{R}_{\mathcal{M}_G}(\pi_{\phi_H}) = \mathbb{E}_{\pi_{\phi_H}, \mathcal{M}_G} \left[\sum_{t=1}^{H'} r_t \right], \quad (1)$$

where π_{ϕ_H} is the policy after K episodes (or H update steps) of adaptation, H' is the horizon of test phase, and r_t is the reward at time t in the test phase. The goal is to find an optimal parameter θ that maximizes the expected return $\mathbb{E}_G[\mathcal{R}_{\mathcal{M}_G}(\pi_{\phi_H})]$ over a given distribution of tasks $p(G)$.

2.2 THE SUBTASK GRAPH INFERENCE PROBLEM

We formulate the *subtask graph inference* problem, an instance of few-shot RL problem where a task is parameterized by *subtask graph* (Sohn et al., 2018). Following few-shot RL settings, the agent’s goal is to quickly adapt to the given task (*i.e.*, MDP) in the adaptation phase to maximize the

Algorithm 1 Adaptation policy optimization during meta-training

Require: $p(G)$: distribution over subtask graph

- 1: **while** not done **do**
- 2: Sample batch of task parameters $\{G_i\}_{i=1}^M \sim p(G)$
- 3: **for all** G_i in the batch **do**
- 4: Rollout K episodes $\tau = \{\mathbf{s}_t, \mathbf{o}_t, r_t, d_t\}_{t=1}^H \sim \pi_\theta^{\text{adapt}}$ in task \mathcal{M}_{G_i} ▷ adaptation phase
- 5: Compute r_t^{UCB} as in Eq.(7)
- 6: $\hat{G}_i = \text{ILP}(\tau)$ ▷ subtask graph inference
- 7: Sample $\tau' \sim \pi_{\hat{G}_i}^{\text{exe}}$ in task \mathcal{M}_{G_i} ▷ test phase
- 8: Update $\theta \leftarrow \theta + \eta \nabla_\theta \sum_{i=1}^M \mathcal{R}_{\mathcal{M}_{G_i}}^{\text{PG+UCB}} \left(\pi_\theta^{\text{adapt}} \right)$ using $\mathcal{R}_{\mathcal{M}}^{\text{PG+UCB}}$ in Eq.(9)

return in the test phase (see Figure 1). A task consists of N subtasks and the subtask graph models a hierarchical dependency between subtasks.

Subtask: A subtask Φ^i can be defined by a tuple (*completion set* $\mathcal{S}_{\text{comp}}^i \subset \mathcal{S}$, *precondition* $G_c^i : \mathcal{S} \mapsto \{0, 1\}$, *subtask reward function* $G_r^i : \mathcal{S} \rightarrow \mathbb{R}$). A subtask Φ^i is *complete* if the current state is contained in its completion set (*i.e.*, $\mathbf{s}_t \in \mathcal{S}_{\text{comp}}^i$), and the agent receives a reward $r_t \sim G_r^i$ upon the completion of subtask Φ^i . A subtask Φ^i is *eligible* (*i.e.*, subtask can be executed) if its precondition G_c^i is satisfied (see Figure 1 for examples). A subtask graph is a tuple of precondition and subtask reward of all the subtasks: $G = (G_c, G_r)$. Then, the task defined by the subtask graph is a factored MDP (Boutillier et al., 1995; Schuurmans & Patrascu, 2002); *i.e.*, the transition model is factored as $p(\mathbf{s}'|\mathbf{s}, a) = \prod_i p_{G_c^i}(s'_i|\mathbf{s}, a)$ and the reward function is factored as $R(\mathbf{s}, a) = \sum_i R_{G_r^i}(\mathbf{s}, a)$ (see Appendix for the detail). The main benefit of factored MDP is that it allows us to model many hierarchical tasks in a principled way with a compact representation such as dynamic Bayesian network (Dean & Kanazawa, 1989; Boutillier et al., 1995). For each subtask Φ^i , the agent can learn an option \mathcal{O}^i (Sutton et al., 1999b) that *executes* the subtask¹.

Environment: The state input to the agent at time step t consists of $\mathbf{s}_t = \{\mathbf{x}_t, \mathbf{e}_t, \text{step}_t, \text{epi}_t, \text{obs}_t\}$.

- **Completion:** $\mathbf{x}_t \in \{0, 1\}^N$ indicates whether each subtask is complete.
- **Eligibility:** $\mathbf{e}_t \in \{0, 1\}^N$ indicates whether each subtask is eligible (*i.e.*, precondition is satisfied).
- **Time budget:** $\text{step}_t \in \mathbb{R}$ is the remaining time steps until episode termination.
- **Episode budget:** $\text{epi}_t \in \mathbb{R}$ is the remaining number of episodes in adaptation phase.
- **Observation:** $\text{obs}_t \in \mathbb{R}^{\mathcal{H} \times \mathcal{W} \times \mathcal{C}}$ is a (visual) observation at time t .

3 METHOD

We propose a *Meta-learner with Subtask Graph Inference* (MSGI) which infers the latent subtask graph G . Figure 1 overviews our approach. Our main idea is to employ two policies: adaptation policy and test policy. During the adaptation phase, an *adaptation policy* $\pi_\theta^{\text{adapt}}$ rolls out K episodes of *adaptation trajectories*. From the collected adaptation trajectories, the agent infers the subtask graph \hat{G} using inductive logic programming (ILP) technique. A *test policy* $\pi_{\hat{G}}^{\text{test}}$, conditioned on the inferred subtask graph \hat{G} , rolls out episodes and maximizes the return in the test phase. Note that the performance depends on the quality of the inferred subtask graph. The adaptation policy indirectly contributes to the performance by improving the quality of inference. Intuitively, if the adaptation policy completes more diverse subtasks during adaptation, the more “training data” is given to the ILP module, which results in more accurate inferred subtask graph. Algorithm 1 summarizes our meta-training procedure. For meta-testing, see Algorithm 2 in Appendix C.

3.1 SUBTASK GRAPH INFERENCE

Let $\tau_H = \{\mathbf{s}_1, \mathbf{o}_1, r_1, d_1, \dots, \mathbf{s}_H\}$ be an adaptation trajectory of the adaptation policy $\pi_\theta^{\text{adapt}}$ for K episodes (or H steps in total) in adaptation phase. The goal is to infer the subtask graph G for this task, specified by preconditions G_c and subtask rewards G_r . We find the maximum-likelihood estimate (MLE) of $G = (G_c, G_r)$ that maximizes the likelihood of the adaptation trajectory τ_H :

¹As in Andreas et al. (2017); Oh et al. (2017); Sohn et al. (2018), such options are pre-learned with curriculum learning; the policy is learned by maximizing the subtask reward, and the initiation set and termination condition are given as $\mathcal{I}^i = \{\mathbf{s} | G_c^i(\mathbf{s}) = 1\}$ and $\beta^i = \mathbb{I}(x^i = 1)$

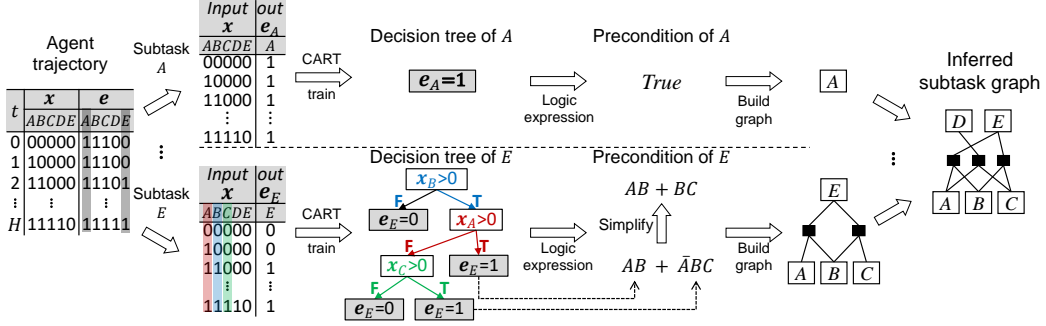


Figure 2: Our inductive logic programming module infers the precondition G_c from adaptation trajectory. For example, the decision tree of subtask E (bottom row) estimates the latent precondition function $f_{G_c^E} : \mathbf{x} \mapsto \mathbf{e}^E$ by fitting its input-output data (i.e., agent’s trajectory $\{\mathbf{x}_t, \mathbf{e}_t^E\}_{t=1}^H$). The decision tree is constructed by choosing a variable (i.e., a component of \mathbf{x}) at each node that best splits the data. The learned decision trees of all the subtasks are represented as logic expressions, and then transformed and merged to form a subtask graph.

$\hat{G}^{\text{MLE}} = \arg \max_{G_c, G_r} p(\tau_H | G_c, G_r)$. The likelihood term can be expanded as

$$p(\tau_H | G_c, G_r) = p(\mathbf{s}_1 | G_c) \prod_{t=1}^H \pi_{\theta}(\mathbf{o}_t | \tau_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{o}_t, G_c) p(r_t | \mathbf{s}_t, \mathbf{o}_t, G_r) p(d_t | \mathbf{s}_t, \mathbf{o}_t) \quad (2)$$

$$\propto p(\mathbf{s}_1 | G_c) \prod_{t=1}^H p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{o}_t, G_c) p(r_t | \mathbf{s}_t, \mathbf{o}_t, G_r), \quad (3)$$

where we dropped the terms that are independent of G . From the definitions in Section 2.2, precondition G_c defines the mapping $\mathbf{x} \mapsto \mathbf{e}$, and the subtask reward G_r determines the reward as $r_t \sim G_r^i$ if subtask i is eligible (i.e., $\mathbf{e}_t^i = 1$) and option \mathcal{O}^i is executed at time t . Therefore, we have

$$\hat{G}^{\text{MLE}} = (\hat{G}_c^{\text{MLE}}, \hat{G}_r^{\text{MLE}}) = \left(\arg \max_{G_c} \prod_{t=1}^H p(\mathbf{e}_t | \mathbf{x}_t, G_c), \arg \max_{G_r} \prod_{t=1}^H p(r_t | \mathbf{e}_t, \mathbf{o}_t, G_r) \right). \quad (4)$$

We note that no supervision from the ground-truth subtask graph G is used. Below we explain how to compute the estimate of preconditions \hat{G}_c^{MLE} and subtask rewards \hat{G}_r^{MLE} .

Precondition inference via logic induction Since the precondition function $f_{G_c} : \mathbf{x} \mapsto \mathbf{e}$ (see Section 2.2 for definition) is a deterministic mapping, the probability term $p(\mathbf{e}_t | \mathbf{x}_t, G_c)$ in Eq.(4) is 1 if $\mathbf{e}_t = f_{G_c}(\mathbf{x}_t)$ and 0 otherwise. Therefore, we can rewrite \hat{G}_c^{MLE} in Eq.(4) as:

$$\hat{G}_c^{\text{MLE}} = \arg \max_{G_c} \prod_{t=1}^H \mathbb{I}(\mathbf{e}_t = f_{G_c}(\mathbf{x}_t)), \quad (5)$$

where $\mathbb{I}(\cdot)$ is the indicator function. Since the eligibility \mathbf{e} is factored, the precondition function $f_{G_c^i}$ for each subtask is inferred independently. We formulate the problem of finding a boolean function that satisfies all the indicator functions in Eq.(5) (i.e., $\prod_{t=1}^H \mathbb{I}(\mathbf{e}_t = f_{G_c}(\mathbf{x}_t)) = 1$) as an *inductive logic programming* (ILP) problem (Muggleton, 1991). Specifically, $\{\mathbf{x}_t\}_{t=1}^H$ forms binary vector inputs to programs, and $\{\mathbf{e}_t^i\}_{t=1}^H$ forms Boolean-valued outputs of the i -th program that denotes the eligibility of the i -th subtask. We use the *classification and regression tree* (CART) to infer the precondition function f_{G_c} for each subtask based on Gini impurity (Breiman, 2017). Intuitively, the constructed decision tree is the simplest boolean function approximation for the given input-output pairs $\{\mathbf{x}_t, \mathbf{e}_t\}$. Then, we convert it to a logic expression (i.e., precondition) in sum-of-product (SOP) form to build the subtask graph. Figure 2 summarizes the overall logic induction process.

Subtask reward inference To infer the subtask reward function \hat{G}_r^{MLE} in Eq.(4), we model each component of subtask reward as a Gaussian distribution $G_r^i \sim \mathcal{N}(\hat{\mu}_r^i, \hat{\sigma}_r^i)$. Then, $\hat{\mu}_{\text{MLE}}^i$ becomes the empirical mean of the rewards received after taking the eligible option \mathcal{O}^i in the trajectory τ_H :

$$\hat{G}_r^{\text{MLE},i} = \hat{\mu}_{\text{MLE}}^i = \mathbb{E}[r_t | \mathbf{o}_t = \mathcal{O}^i, \mathbf{e}_t^i = 1] = \frac{\sum_{t=1}^H r_t \mathbb{I}(\mathbf{o}_t = \mathcal{O}^i, \mathbf{e}_t^i = 1)}{\sum_{t=1}^H \mathbb{I}(\mathbf{o}_t = \mathcal{O}^i, \mathbf{e}_t^i = 1)}. \quad (6)$$

3.2 TEST PHASE: SUBTASK GRAPH EXECUTION POLICY

Once a subtask graph \widehat{G} has been inferred, we can derive a *subtask graph execution (SGE) policy* $\pi_{\widehat{G}}^{\text{exc}}(\mathbf{o}|\mathbf{x})$ that aims to maximize the cumulative reward in the test phase. Note that this is precisely the problem setting used in Sohn et al. (2018). Therefore, we employ a graph reward propagation (GRProp) policy (Sohn et al., 2018) as our SGE policy. Intuitively, the GRProp policy approximates a subtask graph to a differentiable form such that we can compute the gradient of modified return with respect to the completion vector to measure how much each subtask is likely to increase the modified return. Due to space limitation, we give a detail of the GRProp policy in Appendix G.

3.3 LEARNING: OPTIMIZATION OF THE ADAPTATION POLICY

We now describe how to learn the adaptation policy $\pi_{\theta}^{\text{adapt}}$, or its parameters θ . We can directly optimize the objective $\mathcal{R}_{\mathcal{M}_G}(\pi)$ using policy gradient methods (Williams, 1992; Sutton et al., 1999a), such as actor-critic method with generalized advantage estimation (GAE) (Schulman et al., 2015). However, we find it challenging to train our model for two reasons: 1) delayed and sparse reward (*i.e.*, the return in the test phase is treated as if it were given as a one-time reward at the last step of adaptation phase), and 2) large task variance due to highly expressive power of subtask graph. To facilitate learning, we propose to give an intrinsic reward r_t^{UCB} to agent in addition to the extrinsic environment reward, where r_t^{UCB} is the upper confidence bound (UCB) (Auer et al., 2002)-inspired exploration bonus term as follows:

$$r_t^{\text{UCB}} = w_{\text{UCB}} \cdot \mathbb{I}(\mathbf{x}_t \text{ is novel}), \quad w_{\text{UCB}} = \sum_{i=1}^N \frac{\log(n^i(0) + n^i(1))}{n^i(e_t^i)}, \quad (7)$$

where N is the number of subtasks, e_t^i is the eligibility of subtask i at time t , and $n^i(e)$ is the visitation count of e^i (*i.e.*, the eligibility of subtask i) during the adaptation phase until time t . The weight w_{UCB} is designed to encourage the agent to make eligible and execute those subtasks that have infrequently been eligible, since such rare data points in general largely improve the inference by balancing the dataset that CART (*i.e.*, our logic induction module) learns from. The conditioning term $\mathbb{I}(\mathbf{x}_t \text{ is novel})$ encourages the adaptation policy to visit novel states with a previously unseen completion vector \mathbf{x}_t (*i.e.*, different combination of completed subtasks), since the data points with same \mathbf{x}_t input will be ignored in the ILP module as a duplication. We implement $\mathbb{I}(\mathbf{x}_t \text{ is novel})$ using a hash table for computational efficiency. Then, the intrinsic objective is given as follows:

$$\mathcal{R}_{\mathcal{M}_G}^{\text{UCB}}(\pi_{\theta}^{\text{adapt}}) = \mathbb{E}_{\pi_{\theta}^{\text{adapt}}, \mathcal{M}_G} \left[\sum_{t=1}^H r_t^{\text{UCB}} \right], \quad (8)$$

where H is the horizon of adaptation phase. Finally, we train the adaptation policy $\pi_{\theta}^{\text{adapt}}$ using an actor-critic method with GAE (Schulman et al., 2015) to maximize the following objective:

$$\mathcal{R}_{\mathcal{M}_G}^{\text{PG+UCB}}(\pi_{\theta}^{\text{adapt}}) = \mathcal{R}_{\mathcal{M}_G}(\pi_{\widehat{G}}^{\text{GRProp}}) + \beta_{\text{UCB}} \mathcal{R}_{\mathcal{M}_G}^{\text{UCB}}(\pi_{\theta}^{\text{adapt}}), \quad (9)$$

where $\mathcal{R}_{\mathcal{M}_G}(\cdot)$ is the meta-learning objective in Eq.(1), β_{UCB} is the mixing hyper-parameter, and \widehat{G} is the inferred subtask graph that depends on the adaptation policy $\pi_{\theta}^{\text{adapt}}$. The complete procedure for training our MSGI agent with UCB reward is summarized in Algorithm 1.

4 RELATED WORK

Meta Reinforcement Learning. There are roughly two broad categories of meta-RL approaches: gradient-based meta-learners (Finn et al., 2017; Nichol et al., 2018; Gupta et al., 2018; Finn et al., 2018; Kim et al., 2018) and RNN-based meta-learners (Duan et al., 2016; Wang et al., 2016). Gradient-based meta RL algorithms, such as MAML (Finn et al., 2017) and Reptile (Nichol et al., 2018), learn the agent’s policy by taking policy gradient steps during an adaptation phase, where the meta-learner aims to learn a good initialization that enables rapid adaptation to an unseen task. RNN-based meta-RL methods (Duan et al., 2016; Wang et al., 2016) updates the hidden states of a RNN as a process of adaptation, where both of hidden state initialization and update rule are meta-learned. Other variants of adaptation models instead of RNNs such as temporal convolutions (SNAIL) (Mishra et al., 2018) also have been explored. Our approach is closer to the second category, but different from existing works as we directly and explicitly infer the task parameter.

Logic induction. Inductive logic programming systems (Muggleton, 1991) learn a set of rules from examples. (Xu et al., 2017) These works differ from ours as they are open-loop LPI; the input data to LPI module is generated by other policy that does not care about ILP process. However, our agent learns a policy to collect data more efficiently (*i.e.*, closed-loop ILP). There also have been efforts

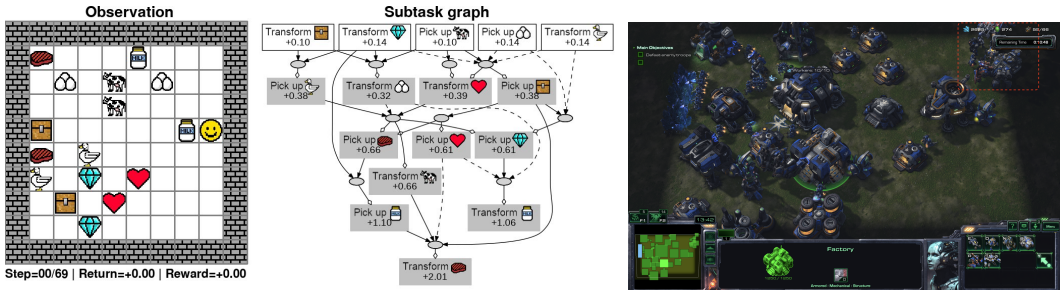


Figure 3: **Left:** A visual illustration of **Playground** domain and an example of underlying subtask graph. **Right:** A *warfare* scenario in **SC2LE** domain (Vinyals et al., 2017). The agent must prepare for the upcoming warfare by training appropriate units, through an appropriate order of subtasks (see Appendix for more details).

to combine neural networks and logic rules to deal with noisy and erroneous data and seek data efficiency, such as (Hu et al., 2016; Evans & Grefenstette, 2017; Dong et al., 2019).

Autonomous Construction of Task Structure. Task planning approaches represented the task structure using Hierarchical Task Networks (HTNs) (Tate, 1977). HTN identifies subtasks for a given task and represent symbolic representations of their preconditions and effects, to reduce the search space of planning (Hayes & Scassellati, 2016). They aim to execute a single goal task, often with assumptions of simpler subtask dependency structures (*e.g.*, without *NOT* dependency (Ghazanfari & Taylor, 2017; Liu et al., 2016)) such that the task structure can be constructed from the successful trajectories. In contrast, we tackle a more general and challenging setting, where each subtask gives a reward (*i.e.*, multi-goal setting) and the goal is to maximize the cumulative sum of reward within an episode. More recently, these task planning approaches were successfully applied to the few-shot visual imitation learning tasks by constructing recursive programs (Xu et al., 2017) or graph (Huang et al., 2018). Contrary to them, we employ an *active* policy that seeks for experience useful in discovering the task structure in unknown and stochastic environments.

5 EXPERIMENTS

In the experiment, we investigate the following research questions: (1) Does MSGI correctly infer task parameters G ? (2) Does adaptation policy $\pi_{\theta}^{\text{adapt}}$ improve the efficiency of few-shot RL? (3) Does the use of UCB bonus facilitate training? (4) How well does MSGI perform compared with other meta-RL algorithms? (5) Can MSGI generalize to longer adaptation horizon, and unseen and more complex tasks?

We evaluate our approach in comparison with the following baselines:

- Random is a policy that executes a random eligible subtask that has not been completed.
- RL² is the meta-RL agent in Duan et al. (2016), trained to maximize the return over K episodes.
- HRL is the hierarchical RL agent in Sohn et al. (2018) trained with the same actor-critic method as our approach during adaptation phase. The network parameter is reset when the task changes.
- GRProp+Oracle is the GRProp policy (Sohn et al., 2018) provided with the ground-truth subtask graph as input This is roughly an upper bound of the performance of MSGI-based approaches.
- MSGI-Rand (Ours) uses a random policy as an adaptation policy, with the task inference module.
- MSGI-Meta (Ours) uses a meta-learned policy (*i.e.*, $\pi_{\theta}^{\text{adapt}}$) as an adaptation policy, with the task inference module.

For RL² and HRL, we use the same network architecture as our MSGI adaptation policy. More details of training and network architecture can be found in Appendix H. The domains on which we evaluate these approaches include two simple grid-world environments (**Mining** and **Playground**) (Sohn et al., 2018) and a more challenging domain **SC2LE** (Vinyals et al., 2017) (StarCraft II).

5.1 EXPERIMENTS ON MINING AND PLAYGROUND DOMAINS

Mining (Sohn et al., 2018) is inspired by Minecraft (see Figure 3) where the agent receives reward by picking up raw materials in the world or crafting items with raw materials. **Playground** (Sohn et al., 2018) is a more flexible and challenging domain, where the environment is stochastic and subtask graphs are *randomly* generated (*i.e.*, precondition is an arbitrary logic expression). We follow the setting in Sohn et al. (2018) for choosing train/evaluation sets. We measure the performance in terms of normalized reward $\bar{R} = (R - R_{\min}) / (R_{\max} - R_{\min})$ averaged over 4 random seeds, where R_{\min} and R_{\max} correspond to the average reward of the Random and the GRProp+Oracle agent, respectively.

5.1.1 TRAINING PERFORMANCE

Figure 4 shows the learning curves of MSGI-Meta and RL², trained on the **D1-Train** set of **Playground** domain. We set the adaptation budget in each trial to $K = 10$ episodes. For MSGI-Rand and HRL (which are not meta-learners), we show the average performance after 10 episodes of adaptation. As training goes on, the performance of MSGI-Meta significantly improves over MSGI-Rand with a large margin. It demonstrates that our meta adaptation policy learns to explore the environment more efficiently, inferring subtask graphs more accurately. We also observe that the performance of RL² agent improves over time, eventually outperforming the HRL agent. This indicates that RL² learns 1) a good initial policy parameter that captures the common knowledge generally applied to all the tasks and 2) an efficient adaptation scheme such that it can adapt to the given task more quickly than standard policy gradient update in HRL.

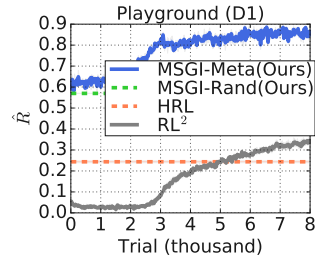


Figure 4: Learning curves on the **Playground** domain. We measure the normalized reward (y-axis) in a test phase, after a certain number of training trials (x-axis).

5.1.2 ADAPTATION AND GENERALIZATION PERFORMANCE

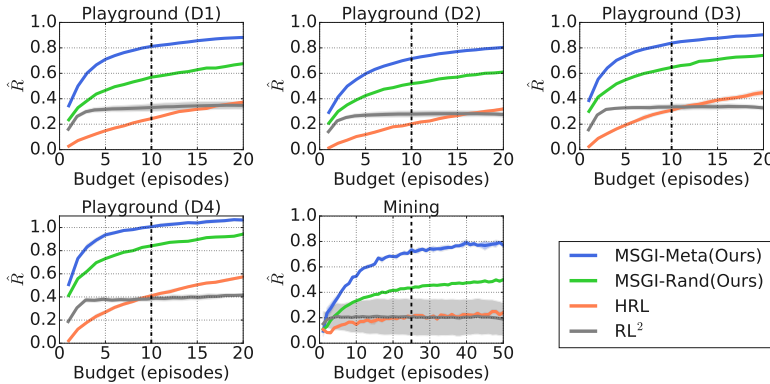


Figure 5: Generalization performance on unseen tasks (**D1-Eval**, **D2**, **D3**, **D4**, and **Mining-Eval**) with varying adaptation horizon. We trained agent with the fixed adaptation budget ($K = 10$ for **Playground** and $K = 25$ for **Mining**) denoted by the vertical dashed line, and tested with varying unseen adaptation budgets. We report the average return during test phase. The shaded area in the plot indicates the range between $\bar{R} + \sigma$ and $\bar{R} - \sigma$ where σ is the standard error of normalized reward.

Adaptation efficiency. In Figure 5, we measure the test performance (in terms of the normalized reward \bar{R}) by varying episode budget K (*i.e.*, how many episodes are used in adaptation phase), after 8000 trials of meta-training (Figure 4). Intuitively, it shows how quickly the agent can adapt to the given task. Our full algorithm MSGI-Meta consistently outperforms MSGI-Rand across all the tasks, showing that our meta adaptation policy can efficiently explore informative states that are likely to result in more accurate subtask graph inference. Also, both of our MSGI-based models perform better than HRL and RL² baselines in all the tasks, showing that explicitly inferring underlying task structure and executing the predicted subtask graph is more effective than learning slow-parameters and fast-parameters (*e.g.*, RNN states) on those tasks involving complex subtask dependencies.

Generalization performance. We test whether the agents can generalize over unseen task and longer adaptation horizon, as shown in Figure 5. For **Playground**, we follow the setup of (Sohn et al., 2018): we train the agent on **D1-Train** with the adaptation budget of 10 episodes, and test on unseen graph distributions **D1-Eval** and larger graphs **D2-D4**. We report the agent’s performance as the normalized reward with up to 20 episodes of adaptation budget. For **Mining**, the agent is trained on randomly generated graphs with 25 episodes budget and tested on 440 hand-designed graphs used in (Sohn et al., 2018), with up to 50 episodes of adaptation budget. Both of our MSGI-based models generalize well to unseen tasks and over different adaptation horizon lengths, continually improving the agent’s performance. It demonstrates that the efficient exploration scheme that our meta adaptation policy can generalize to unseen tasks and longer adaptation horizon, and that our task execution policy, GRProp, generalizes well to unseen tasks as already shown in (Sohn et al., 2018). However, RL² fails to generalize to unseen task and longer adaptation horizon: on **D2-D4** with adaptation horizons longer than the length the meta-learner was trained for, the performance of the RL² agent is almost stationary or even decreases for very long-horizon case (**D2**, **D3**, and **Mining**), eventually being surpassed by

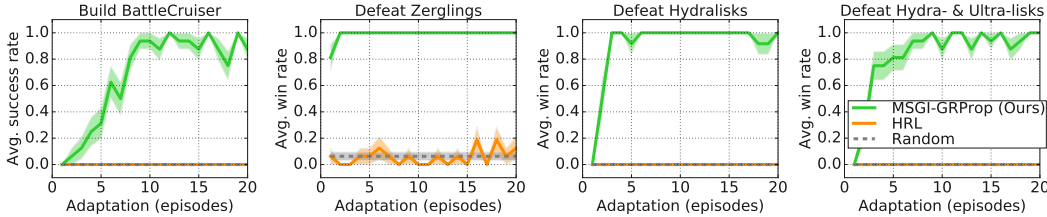


Figure 6: Adaptation performance with different adaptation horizon on **SC2LE** domain.

the HRL agent. This indicates (1) the adaptation scheme that RL^2 learned does not generalize well to longer adaptation horizons, and (2) a common knowledge learned from the training tasks does not generalize well to unseen test tasks.

5.2 EXPERIMENTS ON STARCRRAFT II DOMAIN

SC2LE (Vinyals et al., 2017) is a challenging RL domain built upon the real-time strategy game StarCraft II. We focus on two particular types of scenarios: *Defeat Enemy* and *Build Unit*. Each type of the scenarios models the different aspect of challenges in the full game. The goal of *Defeat Enemy* is to eliminate various enemy armies invading in 2,400 steps. We consider three different combinations of units with varying difficulty: *Defeat Zerglings*, *Defeat Hydralisks*, *Defeat Hydralisks & Ultralisks* (see Figure 8 and demo videos at <http://bit.ly/2M56WNO>). The goal of *Build Unit* scenario is to build a specific unit within 2,400 steps. To showcase the advantage of MSGI inferring the underlying subtask graph, we set a *Battlecruiser* unit as the goal, which is in the highest rank in the technology tree of *Terran* race. In both scenarios, the agent needs to train workers, collect resources, and construct buildings and produce units in a right order to win the game. Each building or unit has a precondition as per the technology tree of the player’s race (see Figure 10 and Appendix D for more details). Note that the precondition of each subtask is determined by the domain and remains fixed across the tasks. If we train the meta agents (MSGI-Meta and RL^2), the agents memorize the subtask dependencies (*i.e.*, over-fitting) and does not learn any useful policy for efficient adaptation. Thus, we only evaluate Random and HRL as our baseline agents, and instead of MSGI-Meta, we use MSGI-GRProp agent to efficiently collect experiences in adaptation phase. MSGI-GRProp infers the subtask graph after every episode from the accumulated experiences and uses it to run the GRProp policy. Since the environment does not provide any subtask-specific reward, we set the subtask reward using the UCB bonus term in Eq. (7) to encourage efficient exploration (See Appendix for detail).

Subtask graph inference. We quantitatively evaluate the inferred subtask graph in terms of the precision and recall of the inferred precondition function $f_{\hat{e}} : \mathbf{x} \mapsto \hat{e}$. Specifically, we compare the inference output \hat{e} with the GT label e generated by the GT precondition function $f_e : \mathbf{x} \mapsto e$ for all possible binary assignments of input (*i.e.*, completion vector \mathbf{x}). For all the tasks, our MSGI-GRProp agent almost perfectly infers the preconditions with more than 94% precision and 96% recall of all possible binary assignments, when averaged over all 163 preconditions in the game, with only 20 episodes of adaptation budget. We provide the detailed quantitative and qualitative results on the inferred subtask graph in supplemental material.

Adaptation efficiency. Figure 6 shows the adaptation efficiency of MSGI-GRProp, HRL agents, and Random policy on the four scenarios. We report the average victory or success rate over 8 episodes. MSGI-GRProp consistently outperforms HRL agents with a high victory rate, by (1) quickly figuring out the useful units and their prerequisite buildings and (2) focusing on executing these subtasks in a correct order. For example, our MSGI-GRProp learns from the inferred subtask graph that some buildings such as *sensor tower* or *engineering bay* are unnecessary for training units and avoids constructing them (see Appendix E for the inferred subtask graph).

6 CONCLUSION

We introduced and addressed a few-shot RL problem with a complex subtask dependencies. We proposed to learn the adaptation policy that efficiently collects experiences in the environment, infer the underlying hierarchical task structure, and maximize the expected reward using the execution policy given the inferred subtask graph. The empirical results confirm that our agent can efficiently explore the environment during the adaptation phase that leads to better task inference and leverage the inferred task structure during the test phase. In this work, we assumed that the option is pre-learned and the environment provides the status of each subtask. In the future work, our approach may be extended to more challenging settings where the relevant subtask structure is fully learned from pure observations, and options to execute these subtasks are also automatically discovered.

REFERENCES

- Jacob Andreas, Dan Klein, and Sergey Levine. Modular multitask reinforcement learning with policy sketches. In *ICML*, 2017.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- Craig Boutilier, Richard Dearden, Moises Goldszmidt, et al. Exploiting structure in policy construction. In *IJCAI*, volume 14, pp. 1104–1113, 1995.
- Leo Breiman. *Classification and regression trees*. Routledge, 2017.
- Devendra Singh Chaplot, Kanthashree Mysore Sathyendra, Rama Kumar Pasumarthi, Dheeraj Rajagopal, and Ruslan Salakhutdinov. Gated-attention architectures for task-oriented language grounding. In *AAAI*, 2018.
- Thomas Dean and Keiji Kanazawa. A model for reasoning about persistence and causation. *Computational intelligence*, 5(2):142–150, 1989.
- Misha Denil, Sergio Gómez Colmenarejo, Serkan Cabi, David Saxton, and Nando de Freitas. Programmable agents. *arXiv preprint arXiv:1706.06383*, 2017.
- Honghua Dong, Jiayuan Mao, Tian Lin, Chong Wang, Lihong Li, and Denny Zhou. Neural logic machines. In *ICLR*, 2019.
- Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RI^2 : Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- Kutluhan Erol. *Hierarchical task network planning: formalization, analysis, and implementation*. PhD thesis, 1996.
- Richard Evans and Edward Grefenstette. Learning Explanatory Rules from Noisy Data. *arXiv preprint arXiv:1711.04574*, 2017.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1126–1135. JMLR. org, 2017.
- Chelsea Finn, Kelvin Xu, and Sergey Levine. Probabilistic model-agnostic meta-learning. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *NIPS*, pp. 9516–9527, 2018.
- Behzad Ghazanfari and Matthew E Taylor. Autonomous extracting a hierarchical structure of tasks in reinforcement learning and multi-task reinforcement learning. *arXiv preprint arXiv:1709.04579*, 2017.
- Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Meta-reinforcement learning of structured exploration strategies. *arXiv preprint arXiv:1802.07245*, 2018.
- Bradley Hayes and Brian Scassellati. Autonomously constructing hierarchical task networks for planning and human-robot collaboration. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5469–5476. IEEE, 2016.
- Sepp Hochreiter, A Steven Younger, and Peter R Conwell. Learning to learn using gradient descent. In *International Conference on Artificial Neural Networks*, pp. 87–94. Springer, 2001.
- Zhiting Hu, Xuezhe Ma, Zhengzhong Liu, Eduard Hovy, and Eric Xing. Harnessing deep neural networks with logic rules. *arXiv preprint arXiv:1603.06318*, 2016.
- De-An Huang, Suraj Nair, Danfei Xu, Yuke Zhu, Animesh Garg, Li Fei-Fei, Silvio Savarese, and Juan Carlos Nieves. Neural task graphs: Generalizing to unseen tasks from a single video demonstration. *arXiv preprint arXiv:1807.03480*, 2018.

- Anders Jonsson and Andrew Barto. Causal graph based decomposition of factored mdps. *Journal of Machine Learning Research*, 7(Nov):2259–2301, 2006.
- Taesup Kim, Jaesik Yoon, Ousmane Dia, Sungwoong Kim, Yoshua Bengio, and Sungjin Ahn. Bayesian model-agnostic meta-learning. *arXiv preprint arXiv:1806.03836*, 2018.
- Changsong Liu, Shaohua Yang, Sari Saba-Sadiya, Nishant Shukla, Yunzhong He, Song-chun Zhu, and Joyce Chai. Jointly learning grounded task structures from language instruction and visual demonstration. In *EMNLP*, 2016.
- Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive meta-learner. In *ICLR*, 2018.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- Stephen Muggleton. Inductive logic programming. *New Gen. Comput.*, 8(4):295–318, February 1991. ISSN 0288-3635. doi: 10.1007/BF03037089. URL <http://dx.doi.org/10.1007/BF03037089>.
- Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *CoRR*, abs/1803.02999, 2, 2018.
- Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization with multi-task deep reinforcement learning. *arXiv preprint arXiv:1706.05064*, 2017.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- Dale Schuurmans and Relu Patrascu. Direct value-approximation for factored mdps. In *Advances in Neural Information Processing Systems*, pp. 1579–1586, 2002.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- Sungryull Sohn, Junhyuk Oh, and Honglak Lee. Hierarchical reinforcement learning for zero-shot generalization with subtask dependencies. In *NeurIPS*, pp. 7156–7166, 2018.
- Richard S. Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*, 1999a.
- Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999b.
- Austin Tate. Generating project networks. In *Proceedings of the 5th international joint conference on Artificial intelligence-Volume 2*, pp. 888–893. Morgan Kaufmann Publishers Inc., 1977.
- Harm Van Seijen, Mehdi Fatemi, Joshua Romoff, Romain Laroche, Tavian Barnes, and Jeffrey Tsang. Hybrid reward architecture for reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 5392–5402, 2017.
- Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.

Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.

Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.

Danfei Xu, Suraj Nair, Yuke Zhu, Julian Gao, Animesh Garg, Li Fei-Fei, and Silvio Savarese. Neural task programming: Learning to generalize across hierarchical tasks. *arXiv preprint arXiv:1710.01813*, 2017.

Haonan Yu, Haichao Zhang, and Wei Xu. A deep compositional framework for human-like language acquisition in virtual environment. *arXiv preprint arXiv:1703.09831*, 2017.

Appendix: Meta Reinforcement Learning with Autonomous Inference of Subtask Dependencies

A SUBTASK GRAPH AND FACTORED MDP

A.1 BACKGROUND: FACTORED MARKOV DECISION PROCESSES

A factored MDP (FMDP) (Boutilier et al., 1995; Jonsson & Barto, 2006) is an MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R})$, where the state space \mathcal{S} is defined by a set of discrete state variables $\mathbf{s} = \{s^1, \dots, s^d\}$. Each state variable $s^i \in \mathbf{s}$ takes on a value in its domain $\mathcal{D}(s^i)$. The state set \mathcal{S} is a (subset of) Cartesian product of the domain of all state variables $\times_{s^i \in \mathbf{s}} \mathcal{D}(s^i)$. In FMDP, the state variables s^i are conditionally independent, such that the transition probability can be factored as follows:

$$p(\mathbf{s}_{t+1}|\mathbf{s}_t, a_t) = p(s_{t+1}^1|\mathbf{s}_t, a_t)p(s_{t+1}^2|\mathbf{s}_t, a_t) \dots p(s_{t+1}^d|\mathbf{s}_t, a_t). \quad (10)$$

Then, the model of FMDP can be compactly represented by the subtask graph (Sohn et al., 2018) or dynamic Bayesian network (DBN) (Dean & Kanazawa, 1989; Boutilier et al., 1995). They represent the transition of each state variable $p(s_{t+1}^i|\mathbf{s}_t, a_t)$ in either a Boolean expression (*i.e.*, subtask graph) or a binary decision tree (*i.e.*, DBN). For more intuitive explanation, see the **subtask graph** paragraph in Section 2.2 and Figure 1.

Jonsson & Barto (2006); Sohn et al. (2018) suggested that the factored MDP can be extended to the option framework. Specifically, the option is defined based on the change in state variable (*e.g.*, completion of subtask in Sohn et al. (2018)), and the option transition model and option reward function are assumed to be factored. Similar to Eq. 10, the transition probability can be factored as follows:

$$p(\mathbf{s}'|\mathbf{s}, o) = \prod_i p(s'_i|\mathbf{s}, o), \quad R(\mathbf{s}, o) = \sum_i R^i(\mathbf{s}, o). \quad (11)$$

In (Sohn et al., 2018), the option \mathcal{O}^i completes the subtask Φ^i by definition; thus, $p(s'_i|\mathbf{s}, o) = 0$ and $R^i(\mathbf{s}, o) = 0$ if $o \neq \mathcal{O}^i$. By introducing the eligibility vector \mathbf{e} , the transition and reward functions are further expanded as follows:

$$p(x'_i|\mathbf{x}, o = \mathcal{O}^i) = p(x'_i|e_i = 1)p(e_i = 1|\mathbf{x}), \quad (12)$$

$$R^i(\mathbf{x}, o = \mathcal{O}^i) = G_{\mathbf{r}}^i \mathbb{I}(e_i = 1), \quad (13)$$

where $p(x'_i|e_i = 1)$ indicates that the subtask is completed x'_i if the subtask is eligible $e_i = 1$, $p(e_i|\mathbf{x})$ is the precondition $G_{\mathbf{e}}$, and $\mathbb{I}(e_i = 1)$ indicates that the reward is given only if the subtask i is eligible.

B DETAILS OF TASK

For self-containedness, we repeat the details of how the task (*i.e.*, MDP) is defined by the subtask graph G from (Sohn et al., 2018). We define each task as an MDP tuple $\mathcal{M}_G = (\mathcal{S}, \mathcal{A}, \mathcal{P}_G, \mathcal{R}_G, \rho_G, \gamma)$ where \mathcal{S} is a set of states, \mathcal{A} is a set of actions, $\mathcal{P}_G : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a task-specific state transition function, $\mathcal{R}_G : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a task-specific reward function and $\rho_G : \mathcal{S} \rightarrow [0, 1]$ is a task-specific initial distribution over states. We describe the subtask graph G and each component of MDP in the following paragraphs.

Subtask and Subtask Graph A subtask graph G is a tuple of the subtask reward $G_{\mathbf{r}} \in \mathbb{R}^N$, and the precondition $G_{\mathbf{e}}$ of N subtasks. The set of subtasks is $\mathcal{O} = \mathcal{A}_{int} \times \mathcal{X}$, where \mathcal{A}_{int} is a set of primitive actions to interact with objects, and \mathcal{X} is a set of all types of interactive objects in the domain. To execute a subtask $(a_{int}, obj) \in \mathcal{A}_{int} \times \mathcal{X}$, the agent should move on to the target object obj and take the primitive action a_{int} .

State The state \mathbf{s}_t consists of the observation $\mathbf{obs}_t \in \{0, 1\}^{W \times H \times C}$, the completion vector $\mathbf{x}_t \in \{0, 1\}^N$, the eligibility vector $\mathbf{e}_t \in \{0, 1\}^N$, the time budget $step_t \in \mathbb{R}$ and number of episode left during the adaptation $epi_t \in \mathbb{R}$. An observation \mathbf{obs}_t is represented as $H \times W \times C$ tensor, where H and W are the height and width of map respectively, and C is the number of object types in the domain. The (h, w, c) -th element of observation tensor is 1 if there is an object c in (h, w) on the map, and 0 otherwise. The time budget indicates the number of remaining time-steps until the episode termination. The completion vector and eligibility vector provides additional information

about N subtasks. The details of completion vector and eligibility vector will be explained in the following paragraph.

State Distribution and Transition Function Given the current state $(\mathbf{obs}_t, \mathbf{x}_t, \mathbf{e}_t)$, the next step state $(\mathbf{obs}_{t+1}, \mathbf{x}_{t+1}, \mathbf{e}_{t+1})$ is computed from the subtask graph G . Figure 7 describes the dependency between subtask graph and MDP. In the beginning of episode, the completion vector \mathbf{x}_t is initialized to a zero vector in the beginning of the episode $\mathbf{x}_0 = [0, \dots, 0]$ and the observation \mathbf{obs}_0 is sampled from the task-specific initial state distribution ρ_G . Specifically, the observation is generated by randomly placing the agent and the N objects corresponding to the N subtasks defined in the subtask graph G . When the agent executes subtask i , the i -th element of completion vector is updated by the following update rule:

$$x_{t+1}^i = \begin{cases} 1 & \text{if } e_t^i = 1 \\ x_t^i & \text{otherwise} \end{cases} . \quad (14)$$

The observation is updated such that agent moves on to the target object, and perform corresponding primitive action. The eligibility vector \mathbf{e}_{t+1} is computed from the completion vector \mathbf{x}_{t+1} and precondition G_c as follows:

$$e_{t+1}^i = \text{OR}_{j \in \text{Child}_i} (y_{AND}^j), \quad (15)$$

$$y_{AND}^i = \text{AND}_{j \in \text{Child}_i} (\hat{x}_{t+1}^{i,j}), \quad (16)$$

$$\hat{x}_{t+1}^{i,j} = x_{t+1}^j w^{i,j} + (1 - x_{t+1}^j)(1 - w^{i,j}), \quad (17)$$

where $w^{i,j} = 0$ if there is a NOT connection between i -th node and j -th node, otherwise $w^{i,j} = 1$, and $w^{i,j}$'s are defined by the G_c . Intuitively, $\hat{x}_t^{i,j} = 1$ when j -th node does not violate the precondition of i -th node. Executing each subtask costs different amount of time depending on the map configuration. Specifically, the time cost is given as the Manhattan distance between agent location and target object location in the grid-world plus one more step for performing a primitive action.

Task-specific Reward Function The reward function is defined in terms of the subtask reward vector $G_r \in \mathbb{R}^N$ and the eligibility vector \mathbf{e}_t , where the subtask reward vector G_r is the component of subtask graph G the and eligibility vector is computed from the completion vector \mathbf{x}_t and subtask graph G as Eq. 17. Specifically, when agent executes subtask i , the amount of reward given to agent at time step t is given as follows:

$$r_t = \begin{cases} G_r^i & \text{if } e_t^i = 1 \\ 0 & \text{otherwise} \end{cases} . \quad (18)$$

Learning option The option framework can be naturally applied to the subtask graph-based tasks. Consider the (optimal) option $\mathcal{O}^i = (\mathcal{I}^i, \pi_{\mathcal{O}}^i, \beta^i)$ for subtask Φ^i . Its initiation set is $\mathcal{I}^i = \{\mathbf{s} | e^i = 1\}$, where \mathbf{s} is the state, e^i is the i -th component of eligibility vector \mathbf{e} , and \mathbf{e} is an element of \mathbf{s} . The termination condition is $\beta^i = \mathbb{I}(x_t^i = 1)$, where x_t^i is the i -th component of completion vector \mathbf{x} . The policy $\pi_{\mathcal{O}}^i$ maximizes the subtask reward G_r^i . Similar to Andreas et al. (2017); Oh et al. (2017); Sohn et al. (2018), the option for each subtask is pre-learned via curriculum learning; i.e., the agent learns options from the tasks consisting of single subtask by maximizing the subtask reward.

C ALGORITHM IN META-TESTING

D DETAILS OF THE SC2LE DOMAIN

The **SC2LE** domain (Vinyals et al., 2017) provides suite of mini-games focusing on specific aspects of the entire StarCraft II game. In this paper, we custom design two types of new, simple mini-games called *Build Unit* and *Defeat Zerg troops*. Specifically, we built *Defeat Zerglings*, *Defeat Hydralisks*, *Defeat Hydralisks & Ultralisks* and *Build Battlecruiser* mini-games that compactly capture the most fundamental goal of the full game. The *Build Unit* mini-game requires the agent to figure-out the target unit and its precondition correctly, such that it can train the target unit within the given short time budget. The *Defeat Zerg troops* mini-game mimics the full game more closely; the agent is

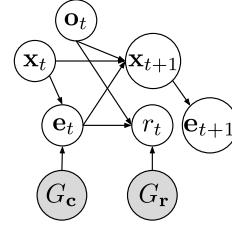


Figure 7: Dependency between subtask graph and MDP

Algorithm 2 Process of single trial for a task \mathcal{M}_G at meta-test time**Require:** The current parameter θ **Require:** A task \mathcal{M}_G parametrized by a task parameter G (unknown to the agent)

- 1: Roll out K train episodes $\tau_H = \{\mathbf{s}_t, \mathbf{o}_t, r_t, d_t\}_{t=1}^H \sim \pi_\theta^{\text{adapt}}$ in task \mathcal{M}_G ▷ adaptation phase
- 2: Infer a subtask graph: $\hat{G} = (\hat{G}_c, \hat{G}_r) = (\text{ILP}(\tau_H), \text{RI}(\tau_H))$ ▷ task inference
- 3: Roll out a test episode $\tau' = \{\mathbf{s}'_t, \mathbf{o}'_t, r'_t, d'_t\}_{t=1}^{H'} \sim \pi_{\hat{G}}^{\text{exe}}$ in task \mathcal{M}_G ▷ test phase
- 4: Measure the performance $R = \sum_t r'_t$ for this task

required to train enough units to win a war against the opponent players. To make the task more challenging and interesting, we designed the reward to be given only at the end of episode depending on the success of the whole task. Similar to the standard *Melee* game in StarCraft II, each episode is



Figure 8: **(Top)** The agent starts the game initially with limited resources of 50 minerals, 0 gases, 3 foods, 11 SCVs collecting resources, 1 idle SCV and pre-built Refinery. **(Middle)** From the initial state, the agent needs to strategically collect resources and build structures in order to be well prepared for the upcoming battle. **(Bottom)** After 2,400 environment steps, the war breaks; all the buildings in the map are removed, and the enemy units appear. The agent’s units should eliminate the enemy units within 240 environment steps during the war.

initialized with 50 *mineral*, 0 *gas*, 7 and 4 *SCVs* that start gathering mineral and gas, respectively, 1 idle *SCV*, 1 refinery, and 1 *Command Center* (See Figure 8). The episode is terminated after 2,400 environment steps (equivalent to 20 minutes in game time). In the game, the agent is initially given 50 *mineral*, 0 *gas*, 7 and 4 *SCVs* that start gathering mineral and gas, respectively, 1 idle *SCV*, 1 refinery, and 1 *Command Center* (See Figure 8) and is allowed to prepare for the upcoming battle only for 2,400 environment steps (equivalent to 20 minutes in game time). Therefore, the agent must learn to collect resources and efficiently use them to build structures for training units. All the four custom mini-games share the same initial setup as specified in Figure 8.

Defeat Zerg troops scenario: At the end of the war preparation, different combinations of enemy unit appears: *Defeat Zerglings* and *Defeat Hydralisks* has 20 zerglings and 15 hydralisks, respectively, and *Defeat Hydralisks & Ultralisks* contains a combination of total 5 hydralisks and 3 ultralisks. When the war finally breaks out, the units trained by the agent will encounter the army of Zerg units in the map and combat until the time over (240 environment steps or 2 minutes in the game) or either side is defeated. Specifically, the agent may not take any action, and the units trained by the agent perform an *auto attack* against the enemy units. Unlike the original full game that has ternary reward structure of +1 (win) / 0 (tie) / -1 (loss), we use binary reward structure of +1 (win) and -1 (loss or tie). Notice that depending on the type of units the agent trained, the tie can happen. For instance, if the units trained by the agent are air units that cannot attack the ground units and the enemy units are the ground units that cannot attack the air units, then no combat will take place, so we consider this case as a loss. **Build unit scenario:** The agent receives the reward of +1 if the target unit is successfully trained within the time limit, and the episode terminates. When the episode terminates due to time limit, the agent receives the reward of -1. We gave 2,400 step budget for the *Build Battlecruiser* scenario such that only highly efficient policy can finish the task within the time limit.

The transition dynamics (*i.e.*, build tech-tree) in **SC2LE** domain has a hierarchical characteristic which can be inferred by our MSGI agent (see Figure 8). We conducted the experiment on Terran race only, but our method can be applied to other races as well.

Subtask. There are 85 subtasks: 15 subtasks of constructing each type of building (*Supply depot, Barracks, Engineeringbay, Refinery, Factory, Missile turret, Sensor tower, Bunker, Ghost academy, Armory, Starport, Fusioncore, Barrack-techlab, Factory-techlab, Starport-techlab*), 17 subtasks of training each type of unit (*SCV, Marine, Reaper, Marauder, Ghost, Widowmine, Hellion, Hellbat, Cyclone, Siegetank, Thor, Banshee, Liberator, Medivac, Viking, Raven, Battlecruiser*), one subtask of idle worker, 32 subtasks of selecting each type of building and unit, gathering mineral, gathering gas, and no-op. For gathering mineral, we set the subtask as ($\text{mineral} \geq \text{val}$) where $\text{val} \in \{50, 75, 100, 125, 150, 300, 400\}$. Similarly for gathering gas, we set the subtask as ($\text{gas} \geq \text{val}$) where $\text{val} \in \{25, 50, 75, 100, 125, 150, 200, 300\}$. For no-op subtask, the agent takes the no-op action for 8 times.

Eligibility. The eligibility of the 15 building construction subtasks and 17 training unit subtasks is given by the environment as an *available action* input. For the selection subtasks, we extracted the number of corresponding units using the provided API of the environment. Gathering mineral, gas, and no-op subtasks are always eligible.

Completion. The completion of the 15 construction subtasks and 17 training subtasks is 1 if the corresponding building or unit is present on the map. For the selection subtasks, the completion is 1 if the target building or unit is selected. For gathering mineral and gas subtasks, the subtask is completed if the condition is satisfied (*i.e.*, $\text{gas} \geq 50$). The no-op subtask is never completed.

Subtask reward. In SC2LE domain, the agent does not receive any reward when completing a subtask. The only reward given to agent is the binary reward $r_{H_{\text{epi}}} = \{+1, -1\}$ at the end of episode (*i.e.*, $t = H_{\text{epi}}$). Therefore, the subtask reward inference method described in Eq.(4) may not be applied. Instead, we tried to infer the subtask reward $\hat{G}_{\mathbf{r}} \in \mathbb{R}^N$ (see Section 3 for definition) from a victory reward $r_{H_{\text{epi}}}$ by building a binary classifier that predicts the victory reward $r_{H_{\text{epi}}}$ from the option count vector $\mathbf{n} \in \mathbb{N}^N$ using a logistic model (*i.e.*, logistic regression), where N is the number of subtasks and the option count vector \mathbf{n} counts how many times each option had been executed within an episode. Intuitively speaking, we assume that the execution of each subtask (*i.e.*, option) gives an implicit subtask reward that is un-observable by the agent, and the victory reward is determined by thresholding the sum of subtask rewards within an episode as follows:

$$r_{H_{\text{epi}}} = \mathbb{I}(G_{\mathbf{r}}^{\top} \mathbf{n} > \beta), \quad (19)$$

where $\mathbb{I}(\cdot)$ is the indicator function and β is the threshold. Then, we approximate it using a sigmoid function $\sigma(\cdot)$ as follows:

$$r_{H_{\text{epi}}} = \sigma(G_{\mathbf{r}}^\top \mathbf{n} - \beta). \tag{20}$$

In the adaptation phase, we randomly sampled the subtask reward vector $G_{\mathbf{r}}$ from the uniform distribution in $[0, 1]^N$, and used it for running MSGI-GRProp agent while recording the option count vector \mathbf{n} . Then, the option count vectors \mathbf{n} and the victory rewards $r_{H_{\text{epi}}}$ from the K episodes in adaptation phase form a training data $\{\mathbf{n}_i, r_{H_{\text{epi},i}}\}_{i=1}^K$ for estimating the parameters of logistic model as follows:

$$\log \frac{r_{H_{\text{epi}}}}{1 - r_{H_{\text{epi}}}} = \widehat{G}_{\mathbf{r}}^\top \mathbf{n} - \beta, \tag{21}$$

where $\widehat{G}_{\mathbf{r}} \in \mathbb{R}^N$ and $\beta \in \mathbb{R}$ are the weight and bias parameters to be learned. Finally, we used the learned $\widehat{G}_{\mathbf{r}}$ as the subtask reward vector for running our MSGI-GRProp agent in test phase. We used the *scikit-learn* (Pedregosa et al., 2011) implementation of logistic regression.

E MORE RESULTS ON THE SC2LE DOMAIN

Accuracy of inferred subtask graph. Figure 9 shows the accuracy of the subtask graph inferred by MSGI-GRProp agent (Section 5.2), in terms of precision and recall over different adaption horizon.

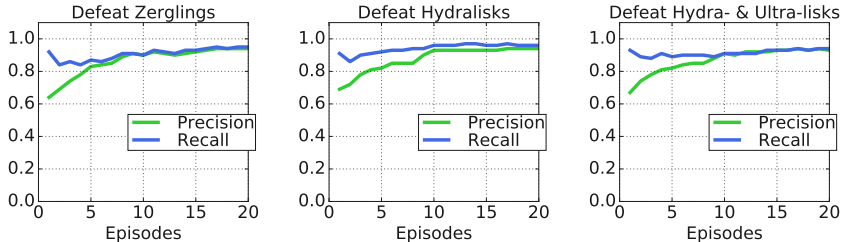


Figure 9: Precision and recall of binary assignments on the inferred subtask graph’s precondition.

Qualitative Examples. Figure 11 shows a simplified form of the subtask graph inferred by our MSGI-GRProp agent after 20 episodes of adaptation. For better readability, we removed the preconditions of resources (food, mineral, gas); Figure 12 depicts the full subtask graph. Compared to the actual tech-tree of the game, we can see the dependency between buildings and units are correctly inferred.



Figure 10: The actual tech-tree of Terran race in StarCraft II. There exists a hierarchy in the task, which can be autonomously discovered by our MSGI agent.

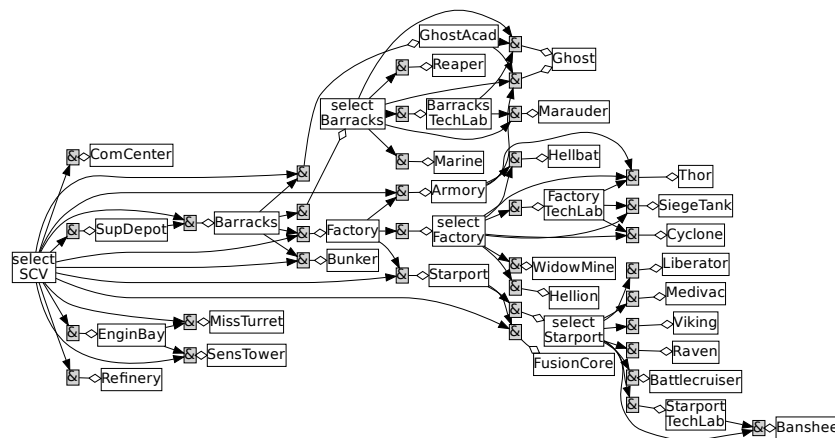


Figure 11: A simplified version of subtask graph inferred by our MSGI-GRProp agent after 10 episodes of adaptation.

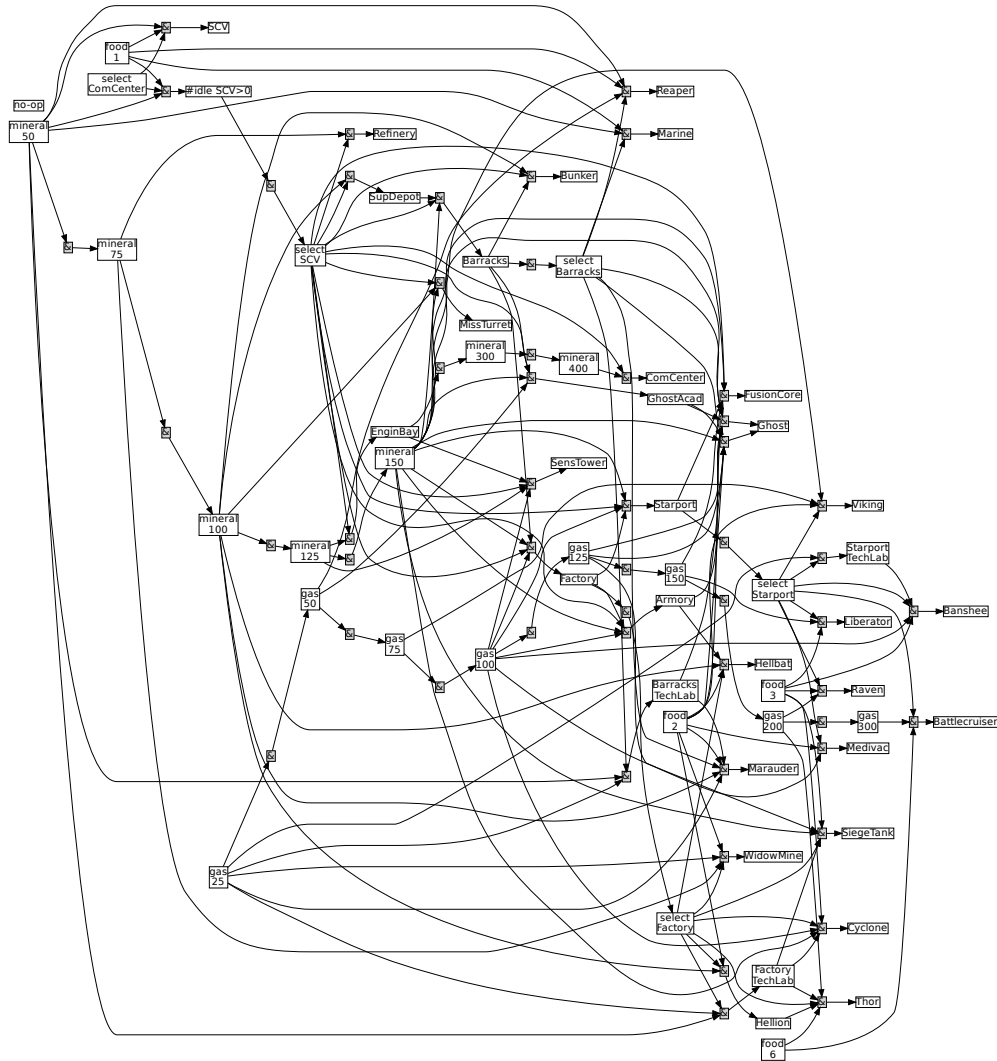


Figure 12: The full subtask graph inferred by our MSGI agent.

F MORE RESULTS ON MINING AND PLAYGROUND

F.1 QUALITATIVE RESULT ON THE SUBTASK GRAPH INFERENCE

Figure 13 illustrates a qualitative example of the inferred subtask graphs inferred by MSGI-Meta and MSGI-Rand agents on the **Mining-Eval** set. The adaptation budget was $K = 50$ episodes and episode length was $T = 80$ steps. Both of MSGI-Meta and MSGI-Rand correctly inferred most of subtasks in the lower hierarchy (e.g., *Get stone*, *Cut wood*, *Get string*) of the subtask graph. However, only MSGI-Meta was successful in inferring the preconditions of subtasks in the highest hierarchy (e.g., *Smelt gold*, *Make goldware*, and *Craft necklace*); MSGI-Rand never had an experience where their preconditions are all satisfied, and thus failed to learn the preconditions of these task. It demonstrates that MSGI-Meta with a meta-learned adaptation policy is able to collect more comprehensive experience for accurate subtask graph inference.

F.2 QUANTITATIVE ANALYSIS OF THE ADAPTATION POLICY

We measured the portion of subtasks being eligible or completed at least once (i.e., coverage) during adaptation to measure how exploratory MSGI-Meta and random policy are. We report the averaged coverage over the evaluation graph set and 8 random seeds. The table shows that MSGI-Meta can make more diverse subtasks complete and eligible than the random policy thanks to more accurate subtask graph inference.

Method	Coverage (%)				
	D1	D2	D3	D4	Eval
MSGI-Meta	89	87	81	75	90
MSGI-Rand	83	77	68	58	85

G DETAILS OF GRPROP POLICY

For self-containedness, we provide the description of GRProp policy from Sohn et al. (2018). We also make a few modifications on $\widehat{\text{OR}}(\mathbf{x})$ and $\widehat{\text{AND}}(\mathbf{x})$ in Eqs. 29 and 30.

Intuitively, GRProp policy modifies the subtask graph to a differentiable form such that we can compute the gradient of modified return with respect to the subtask completion vector in order to measure how much each subtask is likely to increase the modified return. Let \mathbf{x}_t be a completion vector and $G_{\mathbf{r}}$ be a subtask reward vector (see Section 2 for definitions). Then, the sum of reward until time-step t is given as:

$$U_t = G_{\mathbf{r}}^{\top} \mathbf{x}_t. \quad (22)$$

We first modify the reward formulation such that it gives a half of subtask reward for satisfying the preconditions and the rest for executing the subtask to encourage the agent to satisfy the precondition of a subtask with a large reward:

$$\widehat{U}_t = G_{\mathbf{r}}^{\top} (\mathbf{x}_t + \mathbf{e}_t)/2. \quad (23)$$

Let y_{AND}^j be the output of j -th AND node. The eligibility vector \mathbf{e}_t can be computed from the subtask graph G and \mathbf{x}_t as follows:

$$\mathbf{e}_t^i = \text{OR}_{j \in \text{Child}_i} \left(y_{AND}^j \right), \quad y_{AND}^j = \text{AND}_{k \in \text{Child}_j} \left(\widehat{x}_t^{j,k} \right), \quad \widehat{x}_t^{j,k} = x_t^k w^{j,k} + \text{NOT}(x_t^k)(1 - w^{j,k}), \quad (24)$$

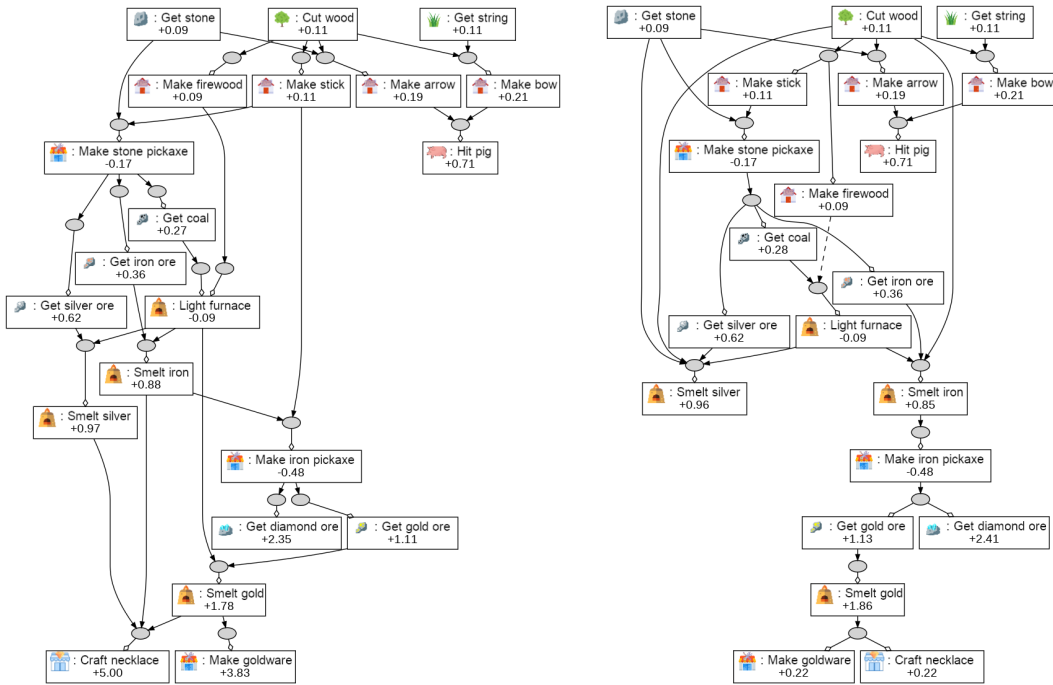
where $w^{j,k} = 0$ if there is a NOT connection between j -th node and k -th node, otherwise $w^{j,k} = 1$. Intuitively, $\widehat{x}_t^{j,k} = 1$ when k -th node does not violate the precondition of j -th node. The logical AND, OR, and NOT operations in Eq. 24 are substituted by the smoothed counterparts as follows:

$$p^i = \lambda_{\text{or}} \widetilde{e}^i + (1 - \lambda_{\text{or}}) x^i, \quad (25)$$

$$\widetilde{e}^i = \widetilde{\text{OR}}_{j \in \text{Child}_i} \left(\widetilde{y}_{AND}^j \right), \quad (26)$$

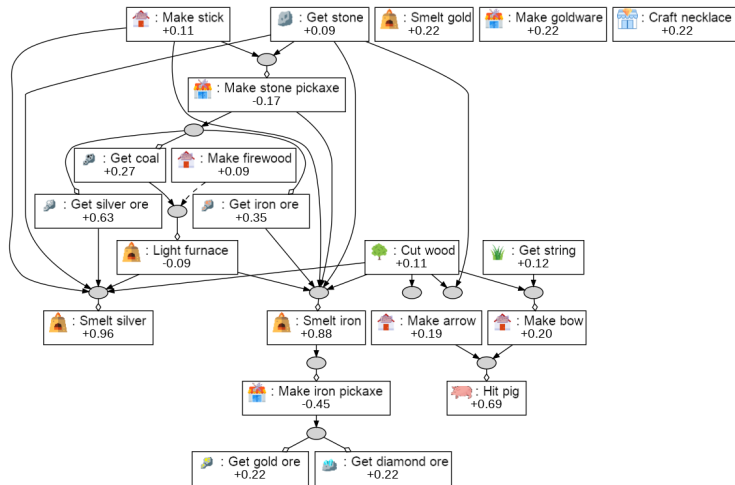
$$\widetilde{y}_{AND}^j = \widetilde{\text{AND}}_{k \in \text{Child}_j} \left(\widehat{x}_t^{j,k} \right), \quad (27)$$

$$\widehat{x}_t^{j,k} = w^{j,k} p^k + (1 - w^{j,k}) \widetilde{\text{NOT}}(p^k), \quad (28)$$



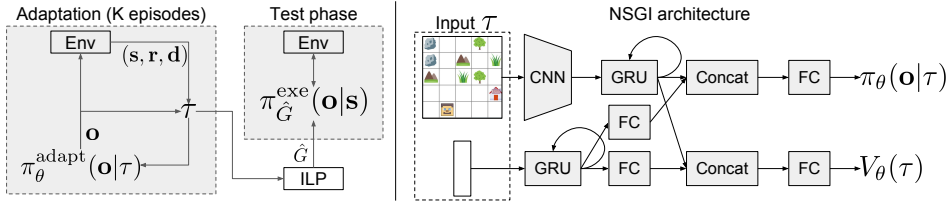
(a) A ground-truth subtask graph.

(b) A subtask graph inferred by MSGI-Meta.



(c) A subtask graph inferred by MSGI-Rand.

Figure 13: A qualitative example of subtask graph inference, in the **Mining** domain.

Figure 14: (Left) Our MSGI model and (Right) the architecture of adaptation policy $\pi_{\theta}^{\text{adapt}}$.

where $\mathbf{x} \in \mathbb{R}^d$ is the input completion vector,

$$\widetilde{\text{OR}}(\mathbf{x}) = \text{softmax}(w_{\text{or}}\mathbf{x}) \cdot \mathbf{x}, \quad (29)$$

$$\widetilde{\text{AND}}(\mathbf{x}) = \frac{\zeta(\mathbf{x}, w_{\text{and}})}{\zeta(\|\mathbf{x}\|, w_{\text{and}})}, \quad (30)$$

$$\widetilde{\text{NOT}}(\mathbf{x}) = -w_{\text{not}}\mathbf{x}, \quad (31)$$

$\|\mathbf{x}\| = d$, $\zeta(\mathbf{x}, \beta) = \frac{1}{\beta} \log(1 + \exp(\beta\mathbf{x}))$ is a soft-plus function, and $\lambda_{\text{or}} = 0.6$, $w_{\text{or}} = 2$, $w_{\text{and}} = 3$, $w_{\text{not}} = 2$ are the hyper-parameters of GRProp. Note that we slightly modified the implementation of $\widetilde{\text{OR}}$ and $\widetilde{\text{AND}}$ from sigmoid and hyper-tangent functions in (Sohn et al., 2018) to softmax and softplus functions for better performance. With the smoothed operations, the sum of smoothed and modified reward is given as:

$$\tilde{U}_t = G_{\mathbf{r}}^{\top} \mathbf{p}, \quad (32)$$

where $\mathbf{p} = [p^1, \dots, p^d]$ and p^i is computed from Eq. 25. Finally, the graph reward propagation policy is a softmax policy,

$$\pi(\mathbf{o}_t | G, \mathbf{x}_t) = \text{Softmax}\left(T \nabla_{\mathbf{x}_t} \tilde{U}_t\right) = \text{Softmax}\left(T G_{\mathbf{r}}^{\top} (\lambda_{\text{or}} \nabla_{\mathbf{x}_t} \tilde{\mathbf{e}}_t + (1 - \lambda_{\text{or}}))\right), \quad (33)$$

where we used the softmax temperature $T = 40$ for **Playground** and **Mining** domain, and linearly annealed the temperature from $T = 1$ to $T = 40$ during adaptation phase for **SC2LE** domain. Intuitively speaking, we act more confidently (*i.e.*, higher temperature T) as we collect more data since the inferred subtask graph will become more accurate.

H IMPLEMENTATION DETAILS

H.1 DETAILS OF MSGI ARCHITECTURE

Figure 14 illustrates the architecture of our MSGI model. Our adaptation policy takes the agent’s trajectory $\tau_t = \{s_t, \mathbf{o}_t, \mathbf{r}_t, \mathbf{d}_t\}$ at time step t as input, where $s = \{\text{obs}, \mathbf{x}, \mathbf{e}, \text{step}, \text{epi}\}$. We used convolutional neural network (CNN) and gated rectifier unit (GRU) to encode both the temporal and spatial information of observation input obs. For other inputs, we simply concatenated all of them along the dimension after normalization, and encoded with GRU and fully-connected (FC) layers. Finally, the flat embedding and observation embedding are concatenated with separate heads for the value and policy output respectively (See supplemental material for more detail).

Our MSGI architecture encodes the observation input using CNN module. Specifically, the observation embedding is computed by Conv1(16x1x1-1/0)-Conv2(32x3x3-1/0)-Conv3(64x3x3-1/1)-Conv4(32x3x3-1/1)-Flatten-FC(512)-GRU(512). Other inputs are all concatenated into a single vector, and fed to GRU(512). In turn, we extracted two flat embeddings using two separate FC(512) heads for policy and value outputs. For each output, the observation and flat embeddings and concatenated into single vector, and fed to FC(512)-FC(d) for policy output and FC(512)-FC(1) for value output, where d is the policy dimension. We used ReLU activation function in all the layers.

H.2 DETAILS OF TRAINING MSGI-META

Algorithm 1 describes the pseudo-code for training our MSGI-Meta model with and without UCB bonus term. In adaptation phase, we ran a batch of 48 parallel environments. In test phase, we measured the average performance over 4 episodes with 8 parallel workers (*i.e.*, average over 32

episodes). We used actor-critic method with GAE (Schulman et al., 2015) as follows:

$$\nabla_{\theta} \mathcal{L} = \mathbb{E}_{G \sim \mathcal{G}_{train}} \left[\mathbb{E}_{s \sim \pi_{\theta}} \left[-\nabla_{\theta} \log \pi_{\theta} \sum_{l=0}^{\infty} \left(\prod_{n=0}^{l-1} (\gamma \lambda)^{k_n} \right) \delta_{t+l} \right] \right], \quad (34)$$

$$\delta_t = r_t + \gamma^{k_t} V_{\theta}^{\pi}(\mathbf{s}_{t+1}) - V_{\theta}^{\pi}(\mathbf{s}_t), \quad (35)$$

where we used the learning rate $\eta = 0.002$, $\gamma = 1$, and $\lambda = 0.9$. We used RMSProp optimizer with the smoothing parameter of 0.99 and epsilon of 1e-5. We trained our MSGI-Meta agent for 8000 trials, where the agent is updated after every trial. We used the best hyper-parameters chosen from the sets specified in Table 2 for all the agents. We also used the entropy regularization with annealed parameter β_{ent} . We started from $\beta_{ent} = 0.05$ and linearly decreased it after 1200 trials until it reaches $\beta_{ent} = 0$ at 3200 trials. During training, we update the critic network to minimize $\mathbb{E}[(R_t - V_{\theta}^{\pi}(\mathbf{s}_t))^2]$, where R_t is the cumulative reward at time t with the weight of 0.03. We clipped the magnitude of gradient to be no larger than 1.

Hyperparameter	Notation	Methods		
		MSGI-Meta	RL ²	HRL
Learning Rate (LR)	η	2e-3	1e-3	1e-3
LR multiplier		0.999	0.999	0.999
GAE	λ	0.9	0.9	0.9
Critic	β_{critic}	0.12	0.005	0.12
Entropy	β_{ent}	0.1	1.0	0.03
UCB	β_{UCB}	1.0	-	-
Architecture	(d_{flat}, d_{gru})	(512, 512)	(512, 512)	(512, 512)

Table 1: Summary of hyper-parameters used for MSGI-Meta, RL², and HRL agents.

Hyperparameter	Notation	Values
Learning rate (LR)	η	{1.0, 2.5, 5.0} × {e-5, e-4, e-3}
LR multiplier		{0.96, 0.98, 0.99, 0.993, 0.996, 0.999, 1.0}
GAE	λ	{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 0.95, 0.98, 1.0}
Critic	β_{critic}	{0.005, 0.01, 0.03, 0.06, 0.12, 0.25}
Entropy	β_{ent}	{0.02, 0.05, 0.1, 0.2, 0.5, 1.0, 2.0}
UCB	β_{UCB}	{1.0, 3.0}
Architecture	(d_{flat}, d_{gru})	{(128, 128), (256, 256), (512, 512)}

Table 2: The range of hyper-parameters we searched over. We did beam-search to find the best parameter with the priority of η , λ , β , β_{ent} , (d_{flat}, d_{gru}) , LR-decay.

H.3 DETAILS OF TRAINING RL² AND HRL

For training RL² and HRL, we used the same architecture and algorithm with MSGI-Meta. For RL², we used the same hyper-parameters except the learning rate $\eta = 0.001$ and the critic loss weight of 0.005. For HRL, we used the learning rate $\eta = 0.001$ and the critic loss weight of 0.12. We used the best hyper-parameters chosen from the sets specified in Table 2 for all the agents.