META REINFORCEMENT LEARNING WITH AUTONOMOUS TASK INFERENCE

Sungryull Sohn1Hyunjae Woo11University of Michigan
{srsohn, hjwoo}@umich.edu

Honglak Lee^{2,1} ²Google Brain honglak@google.com

Abstract

We propose and solve a new few-shot RL problem, where the 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 multiple episodes in the adaptation phase to maximize the return in the test phase. Instead of directly trying to learn a meta-policy, we introduce a *neural subtask graph inferencer* (NSGI) that infers the latent parameter of the task by interacting with the environment and maximizes the return given the latent parameter. To facilitate the learning, we adopt the intrinsic reward inspired by upper confidence bound (UCB) that encourages to explore the environment more efficiently. Our experiment results on two 2D maze domains show that the proposed method adapts more efficiently and generalizes better than the existing meta RL and hierarchical RL methods.

1 INTRODUCTION



Figure 1: An example task and the procedure of our agent solving the task. Our adaptation policy π^{adapt} learns to maximally gather information about current task in adaptation phase. Task inference (TI) module infers the subtask graph G from the adaptation trajectory $\tau_{\pi^{\text{adapt}}}$. Lastly, execution policy π^{exe} executes the inferred subtask graph \hat{G} to maximize the reward in test phase.

Recently, reinforcement learning (RL) systems achieved super-human performance on many complex tasks (Mnih et al., 2015; Silver et al., 2016; Van Seijen et al., 2017). However, these works have mostly been focused on a single known task where the agent can be trained for a large number of times. In order to build more practical and scalable RL agent, it is necessary for the agent to execute many unknown tasks. To this end, recent works on multi-task RL tried to build agent that can execute any given task descriptions such as sequential instructions (Oh et al., 2017; Andreas et al., 2017; Yu et al., 2017; Denil et al., 2017) and complex graph structures (Sohn et al., 2018). However, such task description may not be readily available in practice. Likewise, Hochreiter et al. (2001); Duan et al. (2016); Wang et al. (2016); Finn et al. (2017) proposed meta-RL algorithms that can quickly adapt to the new task without task description input, but they have focused only on the simple tasks such as a single goal navigation task and a bandit problem.

We formulate a new meta RL problem, named *subtask graph inference* problem, where the agent should execute the complex task characterized by a *subtask graph* (Sohn et al., 2018) that is not given to the agent; instead, the agent is given few episodes of adaptation phase to learn about the task as in Figure 1. After the adaptation phase, the agent is required to maximize the reward within a time limit by executing subtasks in an optimal order in test phase. The subtask graph defines the subtasks with the corresponding rewards and the preconditions (e.g., recipe in Figure 1). Inspired by recent multi-task RL works, we propose *neural subtask graph inferencer* that explicitly infers the latent information about the task (e.g., subtask graph) which describes the subtasks and their dependencies, and executes the inferred task.



Figure 2: Procedure of the agent-environment interaction and training in K-shot RL setting. The agent has a slow-learning parameter θ and a fast-learning parameter ϕ .

The contribution of this work can be summarized as follows: (1) We propose a new meta RL problem with a richer form of tasks compared to the recent works on meta RL. (2) We propose an inductive logic programming-based task inference module that can infer the latent task parameter from the agent's trajectory. (3) We propose a meta policy that is trained to efficiently infer the task parameter for the faster adaptation. (4) We compare our method with other meta-RL agents on two 2D visual domains to show that our method can efficiently adapt to the task with complex subtask dependencies.

2 FORMALISM

2.1 Few-shot Reinforcement Learning

We assume that a task is defined by an MDP tuple $\mathcal{M}_G = (S, \mathcal{A}, \mathcal{P}_G, \mathcal{R}_G, \rho_G, \gamma)$ parameterized by the task parameter G with state set S, action set \mathcal{A} , transition dynamics \mathcal{P}_G , reward function \mathcal{R}_G , initial state distribution ρ_G , and discount factor γ . In K-shot RL, a trial consists of the adaptation phase with K episodes and the following test phase under a fixed MDP \mathcal{M}_G as in Figure 2. Between trials, the task parameter G is sampled and the agent updates its slow-learning parameter θ in training. In adaptation phase, the agent adapts to the task by updating its fast-learning parameter ϕ , where the initialization and update function of ϕ are the functions of θ . In the test phase, the agent's performance is measured in terms of the loss $\mathcal{L}(\pi_{\theta}) = -\mathbb{E}_{\pi_{\theta,\phi_{H+1}}} \left[\sum_{t=1}^{H'} r'_t \right]$, where H is the adaptation phase horizon, ϕ_{H+1} is the fast-learning parameter after the adaptation, $\pi_{\theta,\phi_{H+1}}$ is the policy after the adaptation, H' is the test phase horizon, and r'_t is the reward at time t in test phase. The goal is to find the optimal policy π_{θ} that minimizes the loss $\mathcal{L}(\pi_{\theta})$.

2.2 SUBTASK GRAPH INFERENCE PROBLEM

Subtask graph and environment We define the terminologies as follows:

- Completion: $\mathbf{x}_t = [x_t^1, \dots, x_t^N]$ where $x_t^i = 1$ if agent has executed subtask *i*, and 0 otherwise.
- Eligibility: $\mathbf{e}_t = [e_t^1, \dots, e_t^N]$ where $e_t^i = 1$ if subtask *i* is *eligible* (i.e., agent can perform subtask *i*) at time *t*, and 0 otherwise.
- Precondition: c = {c¹,...,c^N} where cⁱ is the logical expression of which variable is the completion vector x and value is the eligibility of *i*-th subtask eⁱ. We refer the corresponding Boolean function fⁱ_c : x → eⁱ as precondition function.
- Boolean function $f_{\mathbf{c}}^i : \mathbf{x} \mapsto e^i$ as **precondition function**. • **Subtask reward**: $\mathbf{r} = [r^1, \dots, r^N]$ specifies the reward given to agent for executing each subtask.
- **Time budget**: step_t $\in \mathbb{R}$ is the remaining time-steps until episode termination.
- Episode budget: $epi_t \in \mathbb{R}$ is the remaining number of episodes in adaptation phase.
- **Observation**: $\mathbf{obs}_t \in \mathbb{R}^{H \times W \times C}$ is a visual observation at time t as illustrated in Figure 1.

The subtask graph G defines N subtasks with corresponding rewards **r** and the preconditions **c**. The state input at time t consists of $\mathbf{s}_t = {\mathbf{obs}_t, \mathbf{x}_t, \mathbf{e}_t, \mathsf{step}_t, \mathsf{epi}_t}$. We assume that the agent has learned a set of *options* (\mathcal{O}) (Precup, 2000; Stolle & Precup, 2002; Sutton et al., 1999) that performs subtasks by executing one or more primitive actions.

3 Method

To tackle the problem, our neural subtask graph inferencer (NSGI) model infers the latent subtask graph G, instead of adapting the fast-learning parameter ϕ as in Figure 7. The adaptation policy $\pi_{\theta}^{\text{adapt}}(\mathbf{o}|\tau)$ rollouts the episodes in adaptation phase. The task inference module subsequently takes as input the adaptation trajectory τ and infers the subtask graph \hat{G} using inductive logic programming (ILP) algorithm. In the test phase, the inferred task parameter \hat{G} is used as the adapted fast-learning parameter ϕ_{H+1} . Our execution policy $\pi_{G}^{\text{exe}}(\mathbf{o}|\mathbf{s})$ takes as input the inferred subtask graph \hat{G} , and execute it to maximize the return.



Figure 3: (Left) Our NSGI model and (Right) the architecture of adaptation policy $\pi_{\theta}^{\text{adapt}}$

3.1 Adaptation phase: subtask graph inference

Let $G = (\mathbf{c}, \mathbf{r})$ be the subtask graph with the subtask rewards \mathbf{r} and \mathbf{c} , and $\tau_t = {\mathbf{s}_1, \mathbf{o}_1, \mathbf{r}_1, \mathbf{d}_1, \dots, \mathbf{s}_t}$ be the adaptation trajectory under adaptation policy π_{θ} until time step t. The maximum-likelihood estimate (MLE) of latent variables G, conditioned on the trajectory τ_H can be written as

$$\hat{G}_{\text{MLE}} = \operatorname*{arg\,max}_{\mathbf{c},\mathbf{r}} p(\tau_H | \mathbf{c}, \mathbf{r}), \tag{1}$$

where H is the adaptation horizon. Then, the likelihood term can be further expanded as

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

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

where we dropped the terms that are independent of G. By definition, the precondition c depends only on the precondition probability in the transition dynamics (see Supplementary material for detail): $p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{o}_t, \mathbf{c}) \propto p(\mathbf{e}_{t+1}|\mathbf{x}_{t+1}, \mathbf{c})$. Then, we can compute \hat{G}_{MLE} as:

$$\hat{G}_{\text{MLE}} = (\hat{\mathbf{c}}_{\text{MLE}}, \hat{\mathbf{r}}_{\text{MLE}}) = \left(\arg\max_{\mathbf{c}} \prod_{t=1}^{H} p(\mathbf{e}_t | \mathbf{x}_t, \mathbf{c}), \ \arg\max_{\mathbf{r}} \prod_{t=1}^{H} p(r_t | \mathbf{s}_t, \mathbf{o}_t, \mathbf{r}) \right).$$
(4)

Precondition inference via logic induction For given precondition c and completion x, the precondition function uniquely determines the corresponding eligibility as $\mathbf{e} = f_{\mathbf{c}}(\mathbf{x})$ (see Section 2.2). Thus, finding $\hat{\mathbf{c}}_{\text{MLE}}$ is equivalent to finding the precondition function $f_{\mathbf{c}}(\cdot)$ that satisfies all the inputoutput pairs $\{\mathbf{x}_t, \mathbf{e}_t\}_{t=1}^H$, which is *logic program induction* problem. To solve this problem, we used the *classification and regression tree* (CART) to infer the precondition function $f_{\mathbf{c}}$ for each subtask based on Gini impurity (Breiman, 2017). Intuitively, the reconstructed decision tree is the simplest Boolean function approximation for the given input-output pairs. Then, we converted it to a logic expression (i.e., precondition) in sum-of-product (SOP) notation to build the subtask graph.

Subtask reward inference For inferring the subtask reward $\hat{\mathbf{r}}_{MLE}$, we modeled each component of subtask reward as a Gaussian distribution $\mathbf{r}^i \sim \mathcal{N}(\hat{\mu}^i, \hat{\sigma}^i)$. Then, the $\hat{\mu}^i_{MLE}$ is the sample mean of the rewards received after executing subtask *i* in the trajectory τ_H : $\hat{\mathbf{r}}^i_{MLE} = \hat{\mu}^i_{MLE} = \mathbb{E}[\mathbf{r}_t | \mathbf{o}_t = i]$.

3.2 TEST PHASE: SUBTASK GRAPH EXECUTION

Given the subtask graph \hat{G} inferred in adaptation phase, we learn a task execution policy $\pi_{\hat{G}}^{exe}(\mathbf{o}_t|\mathbf{s}_t)$ that aims to maximize the cumulative reward in test phase. Since the subtask graph \hat{G} is given, the problem is equivalent to the subtask graph execution problem (Sohn et al., 2018). Following (Sohn et al., 2018), we used graph reward propagation (GRProp) policy $\pi_{\hat{G}}^{\text{GRProp}}(\mathbf{o}|\mathbf{x})$ as our SGE policy. Please see the supplementary material for the details on GRProp policy.

3.3 ARCHITECTURE

Figure 7 illustrates the architecture of our NSGI model. Our adaptation policy takes the agent's trajectory $\tau_t = {\mathbf{s}_t, \mathbf{o}_t, \mathbf{r}_t, \mathbf{d}_t}$ at time *t* as input. For observation input, the policy and value outputs shares the embedding encoded by CNN and GRU. For other inputs, we concatenated them and used two separate fully-connected (FC) layers for the outputs after GRU encoding. Finally, the observation and flat embeddings are concatenated and fed to FC layer to produce each of policy and value outputs (see supplementary material for more detail).

3.4 POLICY OPTIMIZATION

Following the few-shot RL formulation, we can directly optimize the loss function $\mathcal{L}(\pi_{\theta})$ using policy gradient by treating the return in test phase as the last step reward of adaptation phase. However, we found it challenging to train our model for two reasons: 1) the delayed and sparse reward, and 2) high task variance due to high expressive power of subtask graph. To facilitate the learning, we propose to give an intrinsic bonus in adaptation phase. Inspired by the upper confidence bound (UCB) (Auer et al., 2002), we define the bonus term as $r_t^{\text{UCB}} = w_{\text{UCB}}\mathbb{I}(\mathbf{x}_t \notin \tau_{t-1})$, where $w_{\text{UCB}} = \frac{1}{N} \sum_{i=1}^N \log(n_0^{(i)} + n_1^{(i)}) / n_{e_t^i}^{(i)}$, N is the number of subtasks and $n_e^{(i)}$ is the number of occurrences of $\mathbf{e}^{(i)} = e$ until time t. The weight w_{UCB} is designed to encourage the agent to make the subtask eligible that is seldom eligible. The conditioning term $\mathbb{I}(\mathbf{x}_t \notin \tau_{t-1})$ encourages the agent to visit the state with novel \mathbf{x}_t , since the inputs with same \mathbf{x}_t are ignored in ILP module for duplication. We trained the adaptation policy π_{θ} using actor-critic method with GAE (Schulman et al., 2015) to minimize the loss $\mathcal{L}^{\text{UCB}}(\pi_{\theta}) = -\mathbb{E}_{\pi_{\theta}}\left[\sum_{t=1}^{H} r_t^{\text{UCB}}\right]$, in adaptation phase, where H is the adaptation horizon (see supplementary material for the complete procedure of training).

4 EXPERIMENT

In the experiment, we investigate the following research questions: 1) Can NSGI accurately infer the task parameter G? 2) Does the adaptation policy π_{θ} improves the inference? 3) Does using UCB bonus facilitate the training? 4) How does NSGI perform compared to other meta-RL algorithms? 5) Can NSGI generalize to longer adaptation horizon, and unseen and more complex tasks?

4.1 Domain

Following (Sohn et al., 2018), we evaluate the performance of our agents on two domains: **Mining** and **Playground**. Table 1 summarizes the subtask graphs used for evaluation (see supplementary material for details on subtask graph generation).

		Mining			
Task	D1	D2	D3	D4	Eval
Depth	4	4	5	6	4-10
Subtask	13	15	16	16	10-26

Mining is inspired by Minecraft (see Figure 1) where the agent may receive reward by picking up raw materials in the maze or crafting items. The set of subtasks Table 1: (Playground) **D1** have the same graph structure as training set, but the graph was unseen. **D2,D3**, and **D4** have (unseen) larger graph structures. (Mining) **Eval** are unseen during training.

and subtask graph were hand-coded according to the crafting recipes in Minecraft.

Playground is a more flexible and challenging domain. The subtask graph is randomly generated without any context, hence its precondition can be any logical expression and the reward may be delayed. Some of the objects randomly move, which makes the environment stochastic.

4.2 AGENTS

We evaluated the following policies:

- **Random** policy executes any eligible subtask.
- **GRProp+GT** is graph reward propagation policy with ground-truth subtask graph input.
- NSGI-RND (Ours) uses random policy to rollout in adaptation phase
- NSGI-Meta (Ours) uses learned meta adaptation policy to rollout in adaptation phase
- \mathbf{RL}^2 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 actor-critic method in adaptation phase. After the test phase, the network parameter is reset.

For \mathbf{RL}^2 and \mathbf{HRL} , we used the same architecture as our NSGI adaptation policy.

4.3 RESULTS

To evaluate the generalization performance, we trained the agent on the smallest graph **D1** with 10 episodes budget, and tested on the unseen graphs **D1** and unseen larger graphs **D2-D4** with up to 20 episodes budget in **Playground**. In **Mining**, we trained on randomly generated graphs with 25 episodes budget and tested on unseen graphs **Eval** with up to 50 episodes budget (see Table 1). We measured the normalized return in test phase $\hat{R} = (R - R_{min})/(R_{max} - R_{min})$ averaged over four random seeds, where R_{min} and R_{max} correspond to the average return of the **Random** and the **GRProp+GT** agent respectively. The shaded area in the plot indicates the range between $\hat{R} + \sigma$ and $\hat{R} - \sigma$ where σ is the standard error of normalized return.

4.3.1 TRAINING PERFORMANCE

The learning curves is illustrated in Figure 4. Over the training, the performance of our **NSGI-Meta** improves over the **NSGI-RND** with a large margin. It demonstrates that our meta adaptation policy learns to explore the environment much more efficiently than random policy by maximizing the UCB bonus. Also, the performance of \mathbf{RL}^2 improves over time, eventually outperforming the **HRL**. This indicates that \mathbf{RL}^2 learns 1) a common strategy that is generally applicable 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 curve on Playground domain.

4.3.2 Adaptation and generalization performance



Figure 5: Generalization performance with varying adaptation horizon on unseen tasks (D1 and Eval) and larger tasks (D2-4) for longer (unseen) horizon. Our NSGI-based models outperform HRL and RL^2 .

Adaptation efficiency In figure 5, we measured the normalized reward \hat{R} with varying number of adaptation episode budget to see how quickly each agent can adapt to the given task. Our NSGI-Meta consistently outperforms NSGI-RND across all the tasks, showing that our meta adaptation policy can efficiently explore only the informative states that is likely to improve the subtask graph inference. Also, our NSGI-Meta and NSGI-RND perform better than HRL and RL² in all the tasks, showing that separating the problem of task inference and execution brings substantial gains for solving complex task with subtask dependencies. Finally, the RL² outperforms HRL for the seen tasks (D1 and Eval) and seen adaptation horizon by learning the common strategy shared among tasks.

Generalization We evaluated the generalization performance on unseen task and longer adaptation horizon in figure 5. Both of our **NSGI**-based models generalize well to unseen tasks and longer adaptation horizon without significant degradation of performance. It demonstrates that the efficient exploration scheme learned from the UCB intrinsic bonus generalizes well to unseen tasks and longer adaptation horizon, and that our task execution policy, GRProp, generalized well to unseen tasks. However, **RL**² fails to generalize to unseen task and longer adaptation horizon. In task **D2-D4** with unseen adaptation horizon, the **RL**² is outperformed by **HRL** and the performance even decreases for very long horizon case (**D2, D3**, and **mining**). This indicates that 1) the adaptation scheme that **RL**² learned does not generalize well to longer adaptation horizon, and fails to outperform policy gradient update of **HRL**, and 2) the common strategies learned by **RL**² to solve training tasks does not generalize well to the unseen testing tasks.

5 CONCLUSION

We introduce and solve a few-shot RL problem with a complex subtask dependencies. We propose to learn a meta-policy that learns to efficiently infer the latent structure of the task in the form of subtask graph, and the execution policy that executes the inferred subtask graph. The empirical results show that our agent can efficiently explore the environment in adaptation phase to improve the inference, and exploit the inferred information in 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 model may be extended to infer the relevant subtask status from the observation, and discover the option from the structure of the task.

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.
- Leo Breiman. Classification and regression trees. Routledge, 2017.
- Misha Denil, Sergio Gómez Colmenarejo, Serkan Cabi, David Saxton, and Nando de Freitas. Programmable agents. *arXiv preprint arXiv:1706.06383*, 2017.
- Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. Rl²: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- 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.
- 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.
- 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.
- 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.
- Doina Precup. Temporal abstraction in reinforcement learning. 2000.
- 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.
- 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.
- Martin Stolle and Doina Precup. Learning options in reinforcement learning. In *International Symposium on Abstraction, Reformulation, and Approximation*, pp. 212–223. Springer, 2002.
- 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, 1999.
- 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.
- 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.
- 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.

A DETAILS OF THE TASK

We define each task as an MDP tuple $\mathcal{M}_G = (S, \mathcal{A}, \mathcal{P}_G, \mathcal{R}_G, \rho_G, \gamma)$ where S is a set of states, \mathcal{A} is a set of actions, $\mathcal{P}_G : S \times \mathcal{A} \times S \rightarrow [0, 1]$ is a task-specific state transition function, $\mathcal{R}_G : S \times \mathcal{A} \rightarrow \mathbb{R}$ is a task-specific reward function and $\rho_G : S \rightarrow [0, 1]$ is a task-specific initial distribution over states. In the experiment, we used $\gamma = 1$. We describe the subtask graph G and each component of MDP in the following paragraphs.

Terminologies We define the terminologies as follows:

- Completion: $\mathbf{x}_t = [x_t^1, \dots, x_t^N]$ where $x_t^i = 1$ if agent has executed subtask *i*, and 0 otherwise.
- Eligibility: $\mathbf{e}_t = [e_t^1, \dots, e_t^N]$ where $e_t^i = 1$ if subtask *i* is *eligible* (i.e., agent can perform subtask *i*) at time *t*, and 0 otherwise.
- Precondition: c = {c¹,...,c^N} where cⁱ is the logical expression of which variable is the completion vector x and value is the eligibility of *i*-th subtask eⁱ. We refer the corresponding Boolean function fⁱ_c : x → eⁱ as precondition function.
- Subtask reward: $\mathbf{r} = [r^1, \dots, r^N]$ specifies the expected reward given to agent for executing each subtask.
- Time budget: step_t $\in \mathbb{R}$ is the remaining time-steps until episode termination.
- Episode budget: $epi_t \in \mathbb{R}$ is the remaining number of episodes in adaptation phase.
- **Observation**: $obs_t \in \mathbb{R}^{H \times W \times C}$ is a visual observation at time t as illustrated in Figure 1.

Subtask graph and subtask The subtask graph consists of N subtasks that is a subset of prelearned options \mathcal{O} , the subtask reward \mathbf{r} , and the set of precondition of each subtask \mathbf{c} . 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 , the completion vector \mathbf{x}_t , the eligibility vector \mathbf{e}_t , the time budget step_t , and the episode budget epi_t . 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 step_t indicates the number of remaining time-steps until the episode termination. The episode budget epi_t indicates the number of remaining episodes until the adaptation phase terminates. The completion vector and eligibility vector provide additional information about N subtasks. The eligibility vector \mathbf{e}_t is computed from completion vector \mathbf{x}_t and the precondition function f_c as $\mathbf{e}_t = f_c(\mathbf{x}_t)$, where the precondition function f_c is defined by subtask graph G. The details of completion vector and eligibility vector will be explained in **Transition dynamics and reward** paragraph below.

Initial state distribution In the beginning of episode, the initial time budget step_t is sampled from a pre-specified range N_{step} for each subtask graph (See section C for detail), the completion vector \mathbf{x}_t is initialized to a zero vector in the beginning of the episode $\mathbf{x}_0 = [0, \ldots, 0]$ 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.

Transition dynamics and reward Given the current state $(obs_t, \mathbf{x}_t, \mathbf{e}_t, step_t, epi_t)$ and option taken o_t , the next step state $(obs_{t+1}, \mathbf{x}_{t+1}, \mathbf{e}_{t+1}, step_{t+1}, epi_{t+1})$ is computed from the subtask graph G. When the agent executes subtask *i*, the *i*-th element of completion vector is updated by the following update rule:

$$x_{t+1}^{i} = \begin{cases} 1 & \text{if } e_{t}^{i} = 1\\ x_{t}^{i} & \text{otherwise} \end{cases}$$
(5)

The agent receives the reward randomly sampled from a distribution with the mean value of \mathbf{r}^i . For both **Playground** and **Mining** domain, we used uniform distribution from $0.8 * \mathbf{r}^i$ to $1.2 * \mathbf{r}^i$. The observation is updated such that agent moves on to the target object, and perform corresponding primitive action (See Section B for the full list of subtasks and corresponding primitive actions on Mining and Playground domain). The precondition function f_c computes eligibility vector \mathbf{e}_{t+1} from

the completion vector \mathbf{x}_{t+1} and subtask graph G as follows:

$$e_{t+1}^{i} = \underset{j \in Child_{i}}{\operatorname{OR}} \left(y_{AND}^{j} \right), \tag{6}$$

$$y_{AND}^{i} = \operatorname{AND}_{j \in Child_{i}} \left(\hat{x}_{t+1}^{i,j} \right), \tag{7}$$

$$\hat{x}_{t+1}^{i,j} = x_{t+1}^j w^{i,j} + (1 - x_{t+1}^j)(1 - w^{i,j}), \tag{8}$$

where $w^{i,j} = 0$ if there is a NOT connection between *i*-th node and *j*-th node, otherwise $w^{i,j} = 1$. 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.

B DETAILS OF ENVIRONMENT

B.1 MINING

There are 15 types of objects: *Mountain, Water, Work space, Furnace, Tree, Stone, Grass, Pig, Coal, Iron, Silver, Gold, Diamond, Jeweler's shop,* and *Lumber shop.* The agent can take 10 primitive actions: *up, down, left, right, pickup, use1, use2, use3, use4, use5* and agent cannot moves on to the *Mountain* and *Water* cell. *Pickup* removes the object under the agent, and *use's* do not change the observation. There are 26 subtasks in the Mining domain:

- Get wood/stone/string/pork/coal/iron/silver/gold/diamond: The agent should go to *Tree/Stone/Grass/Pig/Coal/Iron/Silver/Gold/Diamond* respectively, and take *pickup* action.
- Make firewood/stick/arrow/bow: The agent should go to *Lumber shop* and take *use1/use2/use3/use4* action respectively.
- Light furnace: The agent should go to *Furnace* and take *use1* action.
- Smelt iron/silver/gold: The agent should go to *Furnace* and take *use2/use3/use4* action respectively.
- Make stone-pickaxe/iron-pickaxe/silverware/goldware/bracelet: The agent should go to *Work space* and take *use1/use2/use3/use4/use5* action respectively.
- Make earrings/ring/necklace: The agent should go to *Jeweler's shop* and take *use1/use2/use3* action respectively.

The icons used in Mining domain were downloaded from www.icons8.com and www.flaticon.com. The *Diamond* and *Furnace* icons were made by Freepik from www.flaticon.com.

B.2 PLAYGROUND

There are 10 types of objects: *Cow*, *Milk*, *Duck*, *Egg*, *Diamond*, *Heart*, *Box*, *Meat*, *Block*, and *Ice*. The *Cow* and *Duck* move by 1 pixel in random direction with the probability of 0.1 and 0.2, respectively. The agent can take 6 primitive actions: *up*, *down*, *left*, *right*, *pickup*, *transform* and agent cannot moves on to the *block* cell. *Pickup* removes the object under the agent, and *transform* changes the object under the agent to *Ice*. The subtask graph was randomly generated without any hand-coded template (see Section C for details).

C DETAILS OF SUBTASK GRAPH GENERATION

C.1 MINING DOMAIN

The precondition of each subtask in Mining domain was defined as Figure 6. Based on this graph, we generated all possible sub-graphs of it by removing the subtask node that has no parent node, while always keeping subtasks A, B, D, E, F, G, H, I, K, L. The reward of each subtask was randomly scaled by a factor of $0.8 \sim 1.2$.



Figure 6: The entire graph of Mining domain. Based on this graph, we generated 640 subtask graphs by removing the subtask node that has no parent node.

C.2 PLAYGROUND DOMAIN

	N_T	number of tasks in each layer
Nodes	N_D	number of distractors in each layer
	N_A	number of AND node in each layer
	r	reward of subtasks in each layer
	N_{ac}^+	number of children of AND node in each layer
	N_{ac}^{-}	number of children of AND node with NOT connection in each layer
Edges	N_{dp}	number of parents with NOT connection of distractors in each layer
-	N_{oc}	number of children of OR node in each layer
Episode	N_{step}	number of step given for each episode

Table 2: Parameters for generating task including subtask graph parameter and episode length.

For training and test sample generation, the subtask graph structure was defined in terms of the parameters in table 2. To cover wide range of subtask graphs, we randomly sampled the parameters $N_A, N_O, N_{ac}^+, N_{ac}^-, N_{dc}$, and N_{oc} from the range specified in the table 3, while N_T and N_D was manually set. We prevented the graph from including the duplicated AND nodes with the same children node(s). We carefully set the range of each parameter such that at least 500 different subtask graphs can be generated with the given parameter ranges. The table 3 summarizes parameters used to generate training and evaluation subtask graphs for the Playground domain.

D DETAILS OF GRPROP POLICY

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. Specifically, the logical AND, OR, and NOT operations in Equations 6, 7, and 8 are substituted by the smoothed counterparts as follows:

$$p^{i} = \lambda_{\rm or} \tilde{e}^{i} + (1 - \lambda_{\rm or}) x^{i}, \tag{9}$$

$$\widetilde{e}^{i} = \widetilde{OR}_{j \in Child_{i}} \left(\widetilde{y}_{AND}^{j} \right), \tag{10}$$

$$\widetilde{y}_{AND}^{j} = \widetilde{\operatorname{AND}}_{k \in Child_{j}} \left(\widehat{x}^{j,k} \right), \tag{11}$$

$$\hat{x}^{j,k} = w^{j,k} p^k + (1 - w^{j,k}) \widetilde{\text{NOT}} \left(p^k \right),$$
(12)

	N_T	{6,4,2,1}
	ND	$\{2,1,0,0\}$
	N.	$[3,3,2]_{5,4,2}$
Tusia	N + N	$\begin{bmatrix} 1, 5, 5, 2 \end{bmatrix}^{-} \begin{bmatrix} 5, 7, 2 \end{bmatrix}$
Irain	Nac	$\{1,1,1\}-\{3,3,3\}$
(=D1)	N_{ac}^{-}	$\{0,0,0\}$ - $\{2,2,1\}$
	N_{dp}	$\{0,0,0\}$ - $\{3,3,0\}$
	N_{oc}	$\{1,1,1\}-\{2,2,2\}$
	r	$\{0.1.0.3.0.7.1.8\}$ - $\{0.2.0.4.0.9.2.0\}$
	Natan	48-72
	$\frac{1 \cdot step}{N_{\pi}}$	{7521}
	N _D	$\{7,3,2,1\}$
	ND	$\{2,2,0,0\}$
	N_A	$\{4,5,2\}$ - $\{5,4,2\}$
D2	N_{ac}	$\{1,1,1\}$ - $\{3,3,3\}$
	N_{ac}^{-}	$\{0,0,0\}$ - $\{2,2,1\}$
	N_{dp}	$\{0,0,0,0\}$ - $\{3,3,0,0\}$
	N_{oc}	$\{1,1,1\}-\{2,2,2\}$
	r	{0.1,0.3,0.7,1.8}-{0.2,0.4,0.9,2.0}
	Nsten	52-78
	N_T	{5,4,4,2,1}
	ND	{1.1.1.0.0}
	N_{A}	$\{3332\}$ - $\{5442\}$
D3	N^+	$\{1, 1, 1, 1\}$ $\{3, 2, 3, 3\}$
05	Nac N-	$(1,1,1,1)^{-}(3,3,3,5)$
	Nac	$\{0,0,0,0\}$ - $\{2,2,1,1\}$
	N_{dp}	$\{0,0,0,0,0\}$ - $\{3,3,3,0,0\}$
	N_{oc}	$\{1,1,1,1\}$ - $\{2,2,2,2\}$
	r	$\{0.1, 0.3, 0.6, 1.0, 2.0\}$ - $\{0.2, 0.4, 0.7, 1.2, 2.2\}$
	N_{step}	56-84
	N_T	{4,3,3,3,2,1}
	N_D	{0,0,0,0,0,0}
	N_{A}	{3.3.3.3.2}-{5.4.4.4.2}
D4	N^+	$\{1,1,1,1,1\}$ - $\{3,3,3,3,3\}$
	N^{-}	$\{0,0,0,0,0\}$ - $\{2,2,1,1,0\}$
	N ₁	$\{0,0,0,0,0,0\}$
	N	[0,0,0,0,0,0] [0,0,0,0,0]
	1 Voc	$\begin{bmatrix} 1^{1,1,1,1,1} \\ - 1^{2$
	r	$\{0.1, 0.5, 0.0, 1.0, 1.4, 2.4\}$ - $\{0.2, 0.4, 0.7, 1.2, 1.0, 2.0\}$
	N_{step}	56-84

Table 3: Subtask graph parameters for training set and tasks D1~D4.

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

$$OR(\mathbf{x}) = \operatorname{softmax}(w_{or}\mathbf{x}) \cdot \mathbf{x}, \tag{13}$$

$$\widetilde{\text{AND}}\left(\mathbf{x}\right) = \max\left(0, \zeta(\mathbf{x}, w_{\text{and}}) - |\mathbf{x}| + 1\right),\tag{14}$$

$$\widetilde{\text{NOT}}\left(\mathbf{x}\right) = -w_{\text{not}}\mathbf{x},\tag{15}$$

 $|\mathbf{x}| = d$, $\zeta(\mathbf{x}, \beta) = \frac{1}{\beta} \log(1 + \exp(\beta \mathbf{x}))$ 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. With the smoothed operations, the sum of smoothed and modified reward is given as:

$$\widetilde{U}_t = \mathbf{r}^T \mathbf{p}.$$
(16)

Finally, the graph reward propagation policy is a softmax policy,

$$\pi(\mathbf{o}_t|G, \mathbf{x}_t) = \operatorname{Softmax}\left(\nabla_{\mathbf{x}_t} \widetilde{U}_t\right) = \operatorname{Softmax}\left(T\lambda_{\operatorname{or}}\mathbf{r}^T + T(1-\lambda_{\operatorname{or}})\mathbf{r}^T \nabla_{\mathbf{x}_t} \widetilde{\mathbf{e}}_t\right), \quad (17)$$

where we used the softmax temperature $T = 40$.



Figure 7: (Left) Our NSGI model and (Right) the architecture of adaptation policy $\pi_{\theta}^{\text{adapt}}$.

E DETAILS OF NSGI ARCHITECTURE

Our NSGI 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(256)-GRU(256). Other inputs are all concatenated into a single vector, and fed to GRU(256). In turn, we extracted two flat embeddings using two separate FC(256) heads for policy and value outputs. For each output, the observation and flat embeddings and concatenated into single vector, and fed to FC(256)-FC(*d*) for policy output and FC(256)-FC(1) for value output, where *d* is the policy dimension. We used ReLU activation function in all the layers.

F DETAILS OF TRAINING NSGI-META

Algorithm 1 Adaptation policy optimization	
Require: $p(G)$: distribution over subtask graph	
Require: α : step size hyperparameter	
1: while not done do	
2: Sample batch of task parameters $\{G_i\}_{i=1}^M \sim p(G)$	
3: for $i = 1,, M$ do	
4: Rollout K episodes $\tau = {\mathbf{s}_t, \mathbf{o}_t, \mathbf{r}_t, \mathbf{d}_t}_{t=1}^H \sim \pi_{\theta}^{\text{adapt}}$ in task $\mathcal{M}_{G_i} \triangleright$ adaptation p	hase
5: if UCB bonus then	
6: Compute $r_t^{\text{UCB}} = w_{\text{UCB}} \mathbb{I} (\mathbf{x}_t \notin \tau_{t-1})$	
7: else	
8: $\hat{G}_i = \text{ILP}(\tau)$ \triangleright task infer	ence
9: Sample $\tau' \sim \pi_{\hat{G}_i}^{\text{GRProp}}$ in task \mathcal{M}_{G_i} > test p	hase
10: if UCB bonus then	_
11: Update $\theta \leftarrow \theta - \alpha \nabla_{\theta} \sum_{i=1}^{M} \mathcal{L}_{\mathcal{M}_{G_i}}^{\text{UCB}} \left(\pi_{\theta}^{\text{adapt}} \right)$ using $\mathcal{L}^{\text{UCB}} \left(\pi_{\theta} \right) = -\mathbb{E}_{\pi_{\theta}} \left \sum_{t=1}^{H} r_t^{\text{UCB}} \right $	
12: else	-
13: Update $\theta \leftarrow \theta - \alpha \nabla_{\theta} \sum_{i=1}^{M} \mathcal{L}_{\mathcal{M}_{G_i}} \left(\pi_{\hat{G}_i}^{\text{GRProp}} \right)$ using $\mathcal{L}(\pi_{\theta}) = -\mathbb{E}_{\pi_{\theta,\phi_{H+1}}} \left[\sum_{t=1}^{H'} r_{\theta,\phi_{H+1}} \right]$	ť

Table 1 describes the pseudo-code for training our **NSGI-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:

$$\nabla_{\theta} \mathcal{L} = \mathbb{E}_{G \sim \mathcal{G}_{train}} \left[\mathbb{E}_{s \sim \pi_{\theta}} \left[-\nabla_{\theta} \log \pi_{\theta} \sum_{l=0}^{\infty} \left(\prod_{n=0}^{l-1} \left(\gamma \lambda \right)^{k_{n}} \right) \delta_{t+l} \right] \right],$$
(18)

$$\delta_t = r_t + \gamma^{k_t} V_{\theta}^{\pi}(\mathbf{s}_{t+1}) - V_{\theta}^{\pi}(\mathbf{s}_t), \tag{19}$$

where we used the learning rate $\alpha = 0.002$, $\gamma = 1$, and $\lambda = 0.9$. We used RMSProp optimizer with the smoothing parameter of 0.99 and epsilon of 1e-5. We trained our **NSGI-Meta** agent for 8000 trials, where the agent is updated after every trial. We also used the entropy regularization with annealed parameter ρ . We started from $\rho = 0.05$ and linearly decreased it after 1200 trials until it reaches $\rho = 0$ at 3200 trials. During training, we update the critic network to minimize $\mathbb{E}\left[\left(R_t - V_{\theta}^{\pi}(\mathbf{s}_t)\right)^2\right]$, 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.

G DETAILS OF TRAINING RL^2 and HRL

For training \mathbf{RL}^2 and \mathbf{HRL} , we used the same architecture and algorithm with NSGI-Meta. For \mathbf{RL}^2 , we used the same hyper-parameters except the learning rate $\alpha = 0.001$ and the critic loss weight of 0.005. For **HRL**, we used the learning rate $\alpha = 0.001$ and the critic loss weight of 0.12. We used the best hyper-parameters chosen from the same candidate hyper-parameter set for all the agents.