

Fast Inference and Transfer of Compositional Task Structures for Few-shot Task Generalization

Sungryull Sohn,^{1,2} Hyunjae Woo,¹ Jongwoon Choi,¹ Lyubing Qiang,¹ Izzeddin Gur,³
Aleksandra Faust,³ Honglak Lee^{1,2}

¹University of Michigan ²LG AI research ³Google Research
{srsohn, hjwoo, jwook, evelynq}@umich.edu, {izzeddin, sandrafaust}@google.com, honglak@eecs.umich.edu

Abstract

We tackle real-world problems with complex structures beyond the pixel-based game or simulator. We formulate it as a few-shot reinforcement learning problem where a task is characterized by a subtask graph that defines a set of subtasks and their dependencies that are unknown to the agent. Different from the previous meta-rl methods trying to directly infer the unstructured task embedding, our multi-task subtask graph inferencer (MTSGI) first infers the common high-level task structure in terms of the subtask graph from the training tasks, and use it as a prior to improve the task inference in testing. Our experiment results on 2D grid-world and complex web navigation domains show that the proposed method can learn and leverage the common underlying structure of the tasks for faster adaptation to the unseen tasks than various existing algorithms such as meta reinforcement learning, hierarchical reinforcement learning, and other heuristic agents.

1 Introduction

Recently, deep reinforcement learning (RL) has shown an outstanding performance on various domains such as video games (Mnih et al. 2015; Vinyals et al. 2019) and board games (Silver et al. 2017). However, most of the successes of deep RL were focused on a single-task setting where the agent is allowed to interact with the environment for hundreds of millions of time steps. In numerous real-world scenarios, interacting with the environment is expensive or limited, and the agent is often presented with a novel task that is not seen during its training time. To overcome this limitation, many recent works focused on scaling the RL algorithm beyond the single-task setting. Recent works on multi-task RL aim to build a single, contextual policy that can solve multiple related tasks and generalize to unseen tasks. However, they require a certain form of task embedding as an extra input that often fully characterizes the given task (Oh et al. 2017; Andreas, Klein, and Levine 2017; Yu, Zhang, and Xu 2017; Chappot et al. 2018), or requires a human demonstration (Huang et al. 2018), which are not readily available in practice. Meta RL (Finn, Abbeel, and Levine 2017; Duan et al. 2016) focuses on a more general setting where the agent should learn about the unseen task purely via interacting with

the environment without any additional information. However, such meta-RL algorithms either require a large amount of experience on the diverse set of tasks or are limited to a relatively smaller set of simple tasks with a simple task structure.

On the contrary, real-world problems require the agent to solve much more complex and compositional tasks without human supervision. Consider a web-navigating RL agent given the task of checking out the products from an online store as shown in Figure 1. The agent can complete the task by filling out the required web elements with the correct information such as shipping or payment information, navigating between the web pages, and placing the order. Note that the task consists of multiple *subtasks* and the subtasks have complex dependencies in the form of *precondition*; for instance, the agent may proceed to the payment web page (see *Bottom, B*) *after* all the required shipping information has been correctly filled in (see *Bottom, A*), or the `credit_card_number` field will appear *after* selecting the `credit_card` as a payment method (see *Top, Middle* in Figure 1). Learning to perform such a task can be quite challenging if the reward is given only after yielding meaningful outcomes (*i.e.*, sparse reward task). This is the problem scope we focus on in this work: solving and generalizing to unseen compositional sparse-reward tasks with complex subtask dependencies without human supervision.

Recent works (Sohn et al. 2019; Xu et al. 2017; Huang et al. 2018; Liu et al. 2016; Ghazanfari and Taylor 2017) tackled the compositional tasks by explicitly inferring the underlying task structure in a graph form. Specifically, the subtask graph inference (SGI) framework (Sohn et al. 2019) uses inductive logic programming (ILP) on the agent’s own experience to infer the task structure in terms of *subtask graph* and learns a contextual policy to *execute* the inferred task in few-shot RL setting. However, it only meta-learned the adaptation policy that relates to the efficient exploration, while the task inference and execution policy learning were limited to a single task (*i.e.*, both task inference and policy learning were done from scratch for each task), limiting its capability of handling large variance in the task structure. We claim that the inefficient task inference may hinder applying the SGI framework to a more complex domain such as web navigation (Shi et al. 2017; Liu et al. 2018) where a task may have a large number of subtasks and complex dependencies

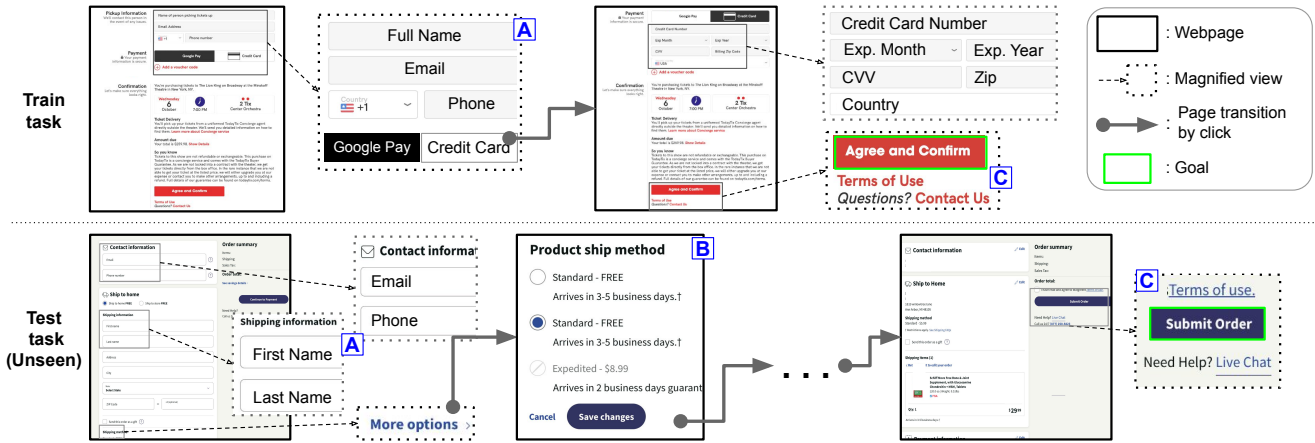


Figure 1: An illustration of the train (*Top*) and test task (*Bottom*) in our *SymWoB* domain. Some selected actionable web-elements (e.g., text fields and buttons) are magnified (dotted arrow and box) for readability. The agent’s goal (green box) is to checkout the products in unseen test website by interacting with the web elements in a correct order. For example, in train task, the agent should fill out all the text fields in (*Top*, **A**) **before** clicking the `credit_card` button to transition (gray arrow) to next page. The high-level checkout processes in different websites have many commonalities while certain details may differ. For example, in both train and test tasks, the agent should fill out the user information (*Top* and *Bottom*, **A**) before proceeding to the next page or there exist similar elements (*Top* and *Bottom*, **C**). However, the details may differ; e.g., the train task (*Top*, **A**) has a single text field for full name, while the test task (*Bottom*, **A**) has separate text fields for the first and last name, respectively. Also, only the test website (*Bottom*, **B**) requires shipping information since the training website does not ship the product.

between them. We note that humans can navigate an unseen website by transferring the high-level process learned from previously seen websites.

Inspired by this, we extend the SGI framework to a *multi-task subtask graph inferencer* (MTSGI) that can generalize the previously learned task structure to the unseen task for faster adaptation and stronger generalization. Figure 2 outlines our method. MTSGI estimates the prior model of the subtask graphs from the training tasks. When an unseen task is presented, MTSGI samples the prior that best matches with the current task, and incorporates the sampled prior model to improve the latent subtask graph inference, which in turn improves the performance of the evaluation policy. We demonstrate results in the 2D grid-world domain and the web navigation domain that simulates the interaction with 20 actual websites. We compare our method with MSGI (Sohn et al. 2019) that learns the task hierarchy from scratch for each task, and two other baselines including hierarchical RL and a heuristic algorithm. We find that MTSGI significantly outperforms all other baselines, and the learned prior model enables more efficient task inference compared to MSGI.

2 Preliminaries

Few-shot Reinforcement Learning A task is defined by an MDP $\mathcal{M}_G = (\mathcal{S}, \mathcal{A}, \mathcal{P}_G, \mathcal{R}_G)$ 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 . The goal of K -shot RL (Duan et al. 2016; Finn, Abbeel, and Levine 2017), is to efficiently solve a distribution of unseen test tasks $\mathcal{M}^{\text{test}}$ by learning and transferring the common knowledge from the training tasks $\mathcal{M}^{\text{train}}$. It is assumed that the training and test tasks do not overlap (i.e., $\mathcal{M}^{\text{train}} \cap \mathcal{M}^{\text{test}} = \emptyset$) but share a

certain commonality such that the knowledge learned from the training tasks may be helpful for learning the test tasks. For each task \mathcal{M}_G , the agent is given K steps budget for interacting with the environment. During meta-training, the goal of multi-task RL agent is to learn a prior (i.e., slow-learning) over the training tasks $\mathcal{M}^{\text{train}}$. Then, the learned prior may be exploited during the meta-test to enable faster adaptation on unseen test tasks $\mathcal{M}^{\text{test}}$. For each task, the agent faces two phases: an *adaptation phase* where the agent learns a task-specific behavior (i.e., fast-learning) for K environment steps, which often spans over multiple episodes, and an *evaluation phase* where the adapted behavior is evaluated. In the evaluation phase, the agent is not allowed to perform any form of learning, and agent’s performance on the task \mathcal{M}_G is measured in terms of the return:

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

where π_{ϕ_K} is the policy after K update steps of adaptation, H is the horizon of evaluation phase, and r_t is the reward at time t in the evaluation phase.

3 Subtask Graph Inference Problem

The *subtask graph inference* problem (Sohn et al. 2019) is a few-shot RL problem where a task is parameterized by a set of subtasks and their dependencies. Formally, a task consists of N subtasks $\Phi = \{\Phi^1, \dots, \Phi^N\}$, and each subtask Φ^i is parameterized by a tuple $(\mathcal{S}_{\text{comp}}^i, G_c^i, G_r^i)$. The *goal state* $\mathcal{S}_{\text{comp}}^i \subset \mathcal{S}$ and *precondition* $G_c^i: \mathcal{S} \rightarrow \{0, 1\}$ defines the condition that a subtask is *completed*: the current state should be contained in its goal states (i.e., $\mathbf{s}_t \in \mathcal{S}_{\text{comp}}^i$) and the precondition should be satisfied (i.e., $G_c^i(\mathbf{s}_t) = 1$). If the

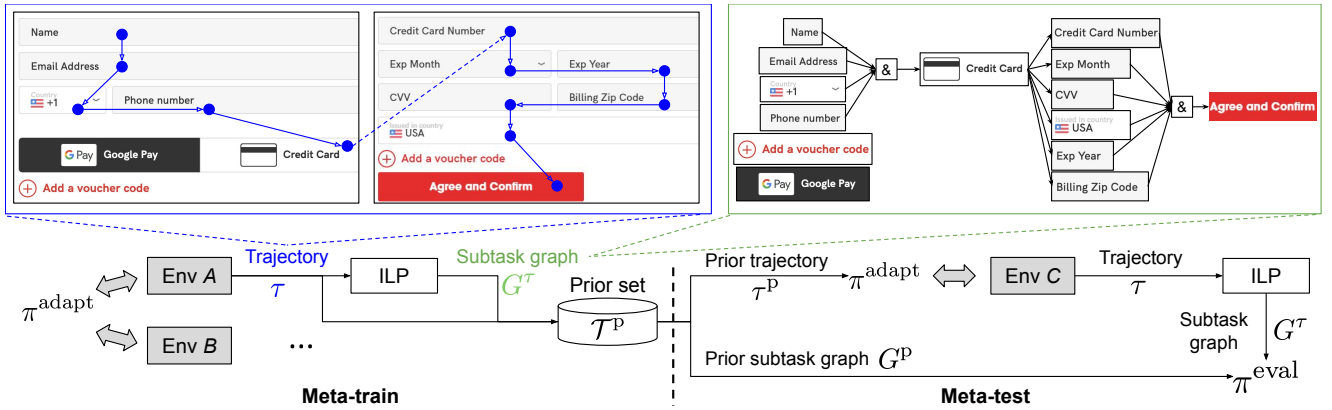


Figure 2: The overview of our algorithm and the example of agent’s **trajectory** and the inferred **subtask graph**. In meta-train (*Left*), the adaptation policy π^{adapt} interacts with the environment and collects the trajectory τ . The inductive logic programming (ILP) module takes as input the trajectory, and infers the task structure in terms of the subtask graph G^T . The trajectory and the subtask graph are stored as a prior. In meta-testing (*Right*), the adaptation policy incorporates the prior trajectory τ^P to efficiently explore the environment, and ILP module infers the subtask graph G^T from the adaptation trajectory τ . Finally, the evaluation policy π^{eval} takes as input the prior and inferred subtask graphs (G^P, G^T) to solve the test task.

precondition is not satisfied (*i.e.*, $G_c^i(s_t) = 0$), the subtask cannot be completed and the agent receives no reward even if the goal state is achieved. The *subtask reward function* G_r^i defines the amount of reward given to the agent when it *completes* the subtask i : $r_t \sim G_r^i$. We note that the subtasks $\{\Phi^1, \dots, \Phi^N\}$ are unknown to the agent. Thus, the agent should learn to infer the underlying task structure and complete the subtasks in an optimal order while satisfying the required preconditions.

State In the subtask graph inference problem, it is assumed that the state input provides the high-level status of the subtasks. Specifically, the state consists of the followings: $s_t = (\text{obs}_t, \mathbf{x}_t, \mathbf{e}_t, \text{step}_{\text{epi},t}, \text{step}_{\text{phase},t})$. The $\text{obs}_t \in \{0, 1\}^{W \times H \times C}$ is a visual observation of the environment. The completion vector $\mathbf{x}_t \in \{0, 1\}^N$ indicates whether each subtask is complete. The eligibility vector $\mathbf{e}_t \in \{0, 1\}^N$ indicates whether each subtask is eligible (*i.e.*, precondition is satisfied). Following the few-shot RL setting, the agent observes two scalar-valued time features: the remaining time steps until the episode termination $\text{step}_{\text{epi},t} \in \mathbb{R}$ and the remaining time steps until the phase termination $\text{step}_{\text{phase},t} \in \mathbb{R}$.

Options For each subtask Φ^i , the agent can learn an option \mathcal{O}^i (Sutton, Precup, and Singh 1999) that reaches the goal state of the subtask. Following (Sohn et al. 2019), such options are pre-learned individually by maximizing the goal-reaching reward: $r_t = \mathbb{I}(s_t \in \mathcal{S}_{\text{comp}}^i)$. At time step t , we denote the option taken by the agent as \mathbf{o}_t and the binary variable that indicates whether episode is terminated as d_t .

4 Method

We propose a novel Multi-Task Subtask Graph Inference (MTSGI) framework that can perform an efficient inference of latent task embedding (*i.e.*, subtask graph). The overall method is outlined in Figure 2. Specifically, in meta-training, MTSGI models the prior in terms of (1) adaptation trajectory

τ and (2) subtask graph G from the agent’s experience. In meta-testing, MTSGI samples (1) the prior trajectory τ^P for more efficient exploration in adaptation and (2) the prior subtask graph G^P for more accurate task inference.

4.1 Multi-task Adaptation Policy

The goal of *adaptation policy* is to efficiently explore and gather the information about the task. Intuitively, if the adaptation policy completes more diverse subtasks, then it can provide more data to the task inference module (ILP), which in turn can more accurately infer the task structure. To this end, we extend the upper confidence bound (UCB)-based adaptation policy proposed in Sohn et al. (2019) as follows:

$$\pi^{\text{adapt}}(\mathbf{o} = \mathcal{O}^i | s) \propto \exp \left(r^i + \sqrt{2} \frac{\log \left(\sum_j n^j \right)}{n^i} \right), \quad (2)$$

where r^i is the empirical mean of the reward received after executing subtask i and n^i is the number of times subtask i has been executed within the current task. Note that the exploration parameters $\{r^i, n^i\}_{i=1}^N$ can be computed from the agent’s trajectory. In meta-train, the exploration parameters are initialized to zero when a new task is sampled. In meta-test, the exploration parameters are initialized with those of the sampled prior. Intuitively, this helps the agent visit novel states that were unseen during meta-training.

4.2 Meta-train: Learning the Prior Subtask Graph

Let τ be an adaptation trajectory of the agent for K steps. The goal is to infer the latent subtask graph G for the given training task $\mathcal{M}_G \in \mathcal{M}^{\text{train}}$, 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

Algorithm 1: Meta-training: learning the prior

Require: Adaptation policy π^{adapt}
Ensure: Prior set \mathcal{T}^{P}

- 1: $\mathcal{T}^{\text{P}} \leftarrow \emptyset$
- 2: **for** each task $\mathcal{M} \in \mathcal{M}^{\text{train}}$ **do**
- 3: Rollout adaptation policy:
 $\tau = \{\mathbf{s}_t, \mathbf{o}_t, r_t, d_t\}_{t=1}^K \sim \pi^{\text{adapt}}$ in task \mathcal{M}
- 4: Infer subtask graph $G^\tau = \arg \max_G p(\tau|G)$
- 5: $\pi^{\text{eval}} = \text{GRProp}(G^\tau)$
- 6: Evaluate the agent: $\tau^{\text{eval}} \sim \pi^{\text{eval}}$ in task \mathcal{M}
- 7: Update prior $\mathcal{T}^{\text{P}} \leftarrow \mathcal{T}^{\text{P}} \cup (G^\tau, \tau)$
- 8: **end for**

of the adaptation trajectory τ :

$$\widehat{G}^{\text{MLE}} = \arg \max_{G_c, G_r} p(\tau|G_c, G_r). \quad (3)$$

Following Sohn et al. (2019), we infer the precondition G_c and the subtask reward G_r as follows (See Appendix for the detailed derivation):

$$\widehat{G}_c^{\text{MLE}} = \arg \max_{G_c} \prod_{t=1}^H p(\mathbf{e}_t|\mathbf{x}_t, G_c), \quad (4)$$

$$\widehat{G}_r^{\text{MLE}} = \arg \max_{G_r} \prod_{t=1}^H p(r_t|\mathbf{e}_t, \mathbf{o}_t, G_r). \quad (5)$$

where \mathbf{e}_t is the eligibility vector, \mathbf{x}_t is the completion vector, \mathbf{o}_t is the option taken, r_t is the reward at time step t .

Precondition inference The problem in Equation (4) is known as the inductive logic programming (ILP) problem that finds a boolean function that satisfies all the indicator functions. 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 predicts the eligibility of the i -th subtask. We use the *classification and regression tree* (CART) to infer the precondition function $f_{G_c} : \mathbf{x} \rightarrow \mathbf{e}$ for each subtask based on Gini impurity (Breiman 1984). Intuitively, the constructed decision tree is the simplest boolean function approximation for the given input-output pairs $\{\mathbf{x}_t, \mathbf{e}_t\}$. The decision tree is converted to a logic expression (*i.e.*, precondition) in sum-of-product (SOP) form to build the subtask graph.

Subtask reward inference To infer the subtask reward $\widehat{G}_r^{\text{MLE}}$ in Equation (5), we model the reward for i -th subtask as a Gaussian distribution: $G_r^i \sim \mathcal{N}(\widehat{\mu}^i, \widehat{\sigma}^i)$. Then, the MLE of subtask reward is given as the empirical mean and variance of the rewards received after taking the eligible option \mathcal{O}^i in adaptation phase:

$$\widehat{\mu}_{\text{MLE}}^i = \mathbb{E} [r_t | \mathbf{o}_t = \mathcal{O}^i, \mathbf{e}_t^i = 1], \quad (6)$$

$$\widehat{\sigma}_{\text{MLE}}^i = \mathbb{E} [(r_t - \widehat{\mu}_{\text{MLE}}^i)^2 | \mathbf{o}_t = \mathcal{O}^i, \mathbf{e}_t^i = 1], \quad (7)$$

where \mathcal{O}^i is the option corresponding to the i -th subtask. Algorithm 1 outlines the meta-training process.

Algorithm 2: meta-testing: multi-task SGI

Require: Adaptation policy π^{adapt} , prior set \mathcal{T}^{P}

- 1: **for** each task $\mathcal{M} \in \mathcal{M}^{\text{test}}$ **do**
- 2: Sample prior: $(G^{\text{P}}, \tau^{\text{P}}) \sim p(\mathcal{T}^{\text{P}})$
- 3: Rollout adaptation policy:
 $\tau = \{\mathbf{s}_t, \mathbf{o}_t, r_t, d_t\}_{t=1}^K \sim \pi^{\text{adapt}}$ in task \mathcal{M}
- 4: Infer subtask graph $G^\tau = \arg \max_G p(\tau|G)$
- 5: $\pi^{\text{eval}}(\cdot|\tau, \tau^{\text{P}}) \propto \text{GRProp}(\cdot|G^\tau)^\alpha \text{GRProp}(\cdot|G^{\text{P}})^{(1-\alpha)}$
- 6: Evaluate the agent: $\tau^{\text{eval}} \sim \pi^{\text{eval}}$ in task \mathcal{M}
- 7: **end for**

4.3 Evaluation: Graph-reward Propagation Policy

In both meta-training and meta-testing, the agent’s adapted behavior is evaluated during the test phase. Following Sohn et al. (2019), we adopted the graph reward propagation (GR-Prop) policy as an evaluation policy π^{eval} that takes as input the inferred subtask graph \widehat{G} and outputs the subtasks to execute to maximize the return. Intuitively, the GRProp policy approximates a subtask graph to a differentiable form such that we can compute the gradient of return with respect to the completion vector to measure how much each subtask is likely to increase the return. Due to space limitations, we give a detail of the GRProp policy in Appendix. The overall meta-training process is summarized in Appendix.

4.4 Meta-testing: Multi-task Task Inference

Prior sampling In meta-testing, MTSGI first chooses the prior task that most resembles the given evaluation task. Specifically, we define the pair-wise similarity between a prior task $\mathcal{M}_G^{\text{prior}}$ and the evaluation task \mathcal{M}_G as follows:

$$\text{sim}(\mathcal{M}_G, \mathcal{M}_G^{\text{prior}}) = F_\beta(\Phi, \Phi^{\text{prior}}) + \kappa R(\tau^{\text{prior}}), \quad (8)$$

where F_β is the F-score with weight parameter β , Φ is the subtask set of \mathcal{M}_G , Φ^{prior} is the subtask set of $\mathcal{M}_G^{\text{prior}}$, $R(\tau^{\text{prior}})$ is the agent’s empirical performance on the prior task $\mathcal{M}_G^{\text{prior}}$, and κ is a scalar-valued weight which we used $\kappa = 1.0$ in experiment. F_β measures how many subtasks overlap between current and prior tasks in terms of precision and recall as follows:

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}, \quad (9)$$

$$\text{Precision} = |\Phi \cap \Phi^{\text{prior}}| / |\Phi^{\text{prior}}|, \quad (10)$$

$$\text{Recall} = |\Phi \cap \Phi^{\text{prior}}| / |\Phi|. \quad (11)$$

We used $\beta = 10$ to assign a higher weight to the current task (*i.e.*, recall) than the prior task (*i.e.*, precision).

Multi-task subtask graph inference Let τ be the adaptation trajectory, and τ^{P} be the sampled prior adaptation trajectory. Then, we model our evaluation policy as follows:

$$\pi(o|s, \tau, \tau^{\text{P}}) \simeq \pi(o|s, G^\tau)^\alpha \pi(o|s, G^{\text{P}})^{(1-\alpha)}. \quad (12)$$

Due to the limited space, we include the detailed derivation of Equation (12) in Appendix. Finally, we deploy the GRProp policy as a contextual policy:

$$\pi^{\text{eval}}(\cdot|\tau, \tau^p) = \text{GRProp}(\cdot|G^\tau)^\alpha \text{GRProp}(\cdot|G^p)^{(1-\alpha)}. \quad (13)$$

Note that Equation (13) is the weighted sum of the logits of two GRProp policies induced by prior τ^p and current experience τ . We claim that such form of ensemble induces the positive transfer in compositional tasks. Intuitively, ensembling GRProp is taking a union of preconditions since GRProp assigns a positive logit to task-relevant subtask and non-positive logit to other subtasks. As motivated in the Introduction, related tasks often share the task-relevant preconditions; thus, taking the union of task-relevant preconditions is likely to be a positive transfer and improve the generalization. The pseudo-code of the multi-task subtask graph inference process is summarized in Algorithm 2.

5 Related Work

Web navigating RL agent Previous work introduced MiniWoB (Shi et al. 2017) and MiniWoB++ (Liu et al. 2018) benchmarks that are manually curated sets of simulated toy environments for the web navigation problem. They formulated the problem as acting on a page represented as a Document Object Model (DOM), a hierarchy of objects in the page. The agent is trained with human demonstrations and online episodes in an RL loop. Jia, Kiros, and Ba (2019) proposed a graph neural network based DOM encoder and a multi-task formulation of the problem similar to this work. Gur et al. (2018) introduced a manually-designed curriculum learning method and an LSTM based DOM encoder. DOM level representations of web pages pose a significant sim-to-real gap as simulated websites are considerably smaller (100s of nodes) compared to noisy real websites (1000s of nodes). As a result, these models are trained and evaluated on the same simulated environments which are difficult to deploy on real websites. Our work formulates the problem as abstract web navigation on real websites where the objective is to learn a latent subtask dependency graph similar to a sitemap of websites. We propose a multi-task training objective that generalizes from a fixed set of real websites to unseen websites without any demonstration, illustrating an agent capable of navigating real websites for the first time.

Meta-reinforcement learning To tackle the few-shot RL problem, researchers have proposed two broad categories of meta-RL approaches: RNN- and gradient-based methods. The RNN-based meta-RL methods (Duan et al. 2016; Wang et al. 2016; Hochreiter, Younger, and Conwell 2001) encode the common knowledge of the task into the hidden states and the parameters of the RNN. The gradient-based meta-RL methods (Finn, Abbeel, and Levine 2017; Nichol, Achiam, and Schulman 2018; Gupta et al. 2018; Finn, Xu, and Levine 2018; Kim et al. 2018) encode the task embedding in terms of the initial policy parameter for fast adaptation through meta gradient. Existing meta-RL approaches, however, often require a large amount of environment interaction due to the long-horizon nature of the few-shot RL tasks. Our work instead explicitly infers the underlying task parameter in

terms of subtask graph, which can be efficiently inferred using the inductive logic programming (ILP) method and be transferred across different, unseen tasks.

More Related Works Please refer to the Appendix for further discussions about other related works.

6 Experiment

6.1 Domains

Mining *Mining* (Sohn, Oh, and Lee 2018) is a 2D grid-world domain inspired by Minecraft game where the agent receives a reward by picking up raw materials in the world or crafting items with raw materials. The subtask dependency in *Mining* domain comes from the crafting recipe implemented in the game. Following Sohn, Oh, and Lee (2018), we used the pre-generated training/testing task splits generated with four different random seeds. Each split set consists of 3200 training tasks and 440 testing tasks for meta-training and meta-testing, respectively. We report the performance averaged over the four task split sets.

SymWoB We implement a symbolic version of the checkout process on the 20 real-world websites such as **Amazon**, **Converse**, and **Walmart**, etc.

Subtask and option policy. Each actionable web element (*e.g.*, text field, button drop-down list, and hyperlink) is considered as a subtask. We assume the agent has pre-learned the option policies that correctly interact with each element (*e.g.*, click the button or fill out the text field). Thus, the agent should learn a policy over the option.

Completion and eligibility. For each subtask, the completion and eligibility are determined based on the status of the corresponding web element. For example, the subtask of a text field is *completed* if the text field is filled with the correct information, and the subtask of a `confirm_credit_info` button is *eligible* if all the required subtasks (*i.e.*, filling out credit card information) on the webpage are completed. Executing an option will complete the corresponding subtask *only if* the subtask is eligible.

Reward function and episode termination. The agent may receive a non-zero reward only at the end of the episode (*i.e.*, sparse-reward task). When the episode terminates due to the time budget, the agent may not receive any reward. Otherwise, the following two types of subtasks terminate the episode and give a non-zero reward upon completion:

- **Goal subtask** refers to the button that completes the order (See the green boxes in Figure 1). Completing this subtask gives the agent a +5 reward, and the episode is terminated.
- **Distractor subtask** does not contribute to solving the given task but terminates the episode with a -1 reward. It models the web elements that lead to external web pages such as `Contact_Us` button in Figure 1.

Transition dynamics. The transition dynamics follow the dynamics of the actual website. Each website consists of multiple web pages. The agent may only execute the subtasks that are currently visible (*i.e.*, on the current web page) and can navigate to the next web page only after filling out all the required fields and clicking the continue button. The goal

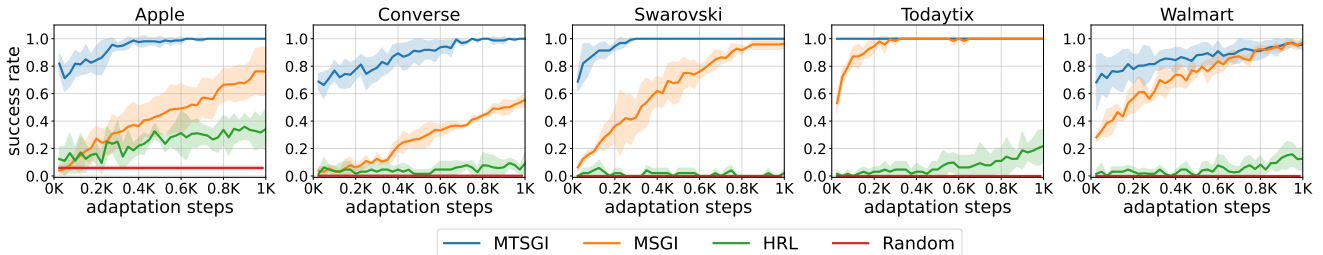


Figure 3: The success rate of the compared methods in the test phase in terms of the environment step during the adaptation phase on *SymWoB* domain. See Appendix for the results on other tasks.

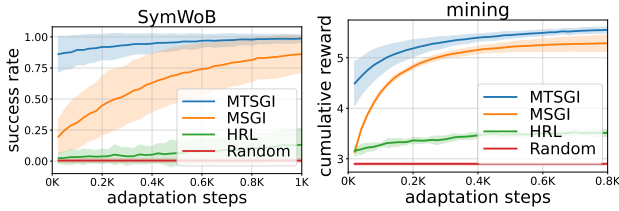


Figure 4: The performance of the compared methods in terms of the adaptation steps averaged over all the tasks in *SymWoB* (Left) and *Mining* (Right) domains.

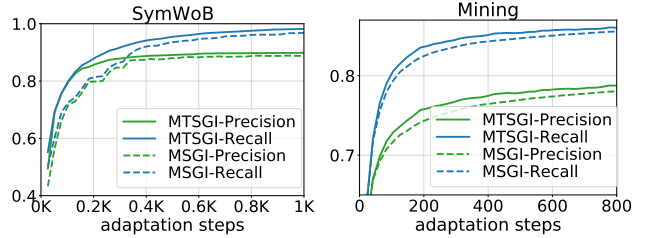


Figure 5: The precision and recall of the subtask graphs inferred by MTSGI and MSGI on *SymWoB* and *Mining*.

subtask is present in the last web page; thus, the agent must learn to navigate through the web pages to solve the task.

For more details about each task, please refer to Appendix.

6.2 Agents

We compared the following algorithms in the experiment.

- MTSGI (Ours): our multi-task SGI agent
- MSGI (Sohn et al. 2019): SGI agent without multi-task learning
- HRL: an Option (Sutton, Precup, and Singh 1999)-based proximal policy optimization (PPO) (Schulman et al. 2017) agent with the gated rectifier unit (GRU)
- Random: a heuristic policy that uniform randomly executes an eligible subtask

More details on the architectures and the hyperparameters can be found in Appendix.

Meta-training In *SymWoB*, for each task chosen for a meta-testing, we randomly sampled N_{train} tasks among the remaining 19 tasks and used it for meta-training. We used $N_{\text{train}} = 1$ in the experiment (See Figure 8 for the impact of the choice of N_{train}). For example, we meta-trained our MTSGI on **Amazon** and meta-tested on **Samsung**. For *Mining*, we used the train/test task split provided in the environment. The RL agents (e.g., HRL) were individually trained on each test task; the policy was initialized when a new task is sampled and trained during the adaptation phase. All the experiments were repeated with four random seeds, where different training tasks were sampled for different seeds.

6.3 Result: Few-shot Generalization Performance

Figure 3 and Figure 4 show the few-shot generalization performance of the compared methods on *SymWoB* and *Mining*.

In Figure 3, MTSGI achieves more than 80% zero-shot success rate (*i.e.*, success rate at $x\text{-axis}=0$) on four out of five tasks, which is significantly higher than the zero-shot performance of MSGI. This indicates that the prior learned from the training task significantly improves the subtask graph inference and in turn improves the multi-task evaluation policy. Moreover, our MTSGI can learn a near-optimal policy on all the tasks after only 1,000 steps of environment interactions, demonstrating that the proposed multi-task learning scheme enables fast adaptation. Even though the MSGI agent is learning each task from scratch, it still outperforms the HRL and Random agents. This shows that explicitly inferring the underlying task structure and executing the predicted subtask graph is significantly more effective than learning the policy from the reward signal (*i.e.*, HRL) on complex compositional tasks. Given the pre-learned options, HRL agent can slightly improve the success rate during the adaptation via PPO update. However, training the policy from the sparse reward requires a large number of interactions especially for the tasks with many distractors (*e.g.*, **Converse** and **Swarovski**).

6.4 Analysis on the Inferred Subtask Graph

We compare the inferred subtask graph with the ground-truth subtask graph. Figure 6 shows the subtask graph inferred by MTSGI in **Walmart**. We can see that MTSGI can accurately infer the subtask graph; the inferred subtask graph is missing only two preconditions (shown in red color) of `ClickContinuePayment` subtask. We note that such a small error in the subtask graph has a negligible effect as shown in Figure 3: *i.e.*, MTSGI achieves near-optimal performance on **Walmart** after 1,000 steps of adaptation. Figure 5 measures the precision and recall of the inferred precondition (*i.e.*, the edge of the graph). First, both MTSGI and MSGI achieve high precision and recall after only a few hundred of adaptation. Also, MTSGI outperforms MSGI in the

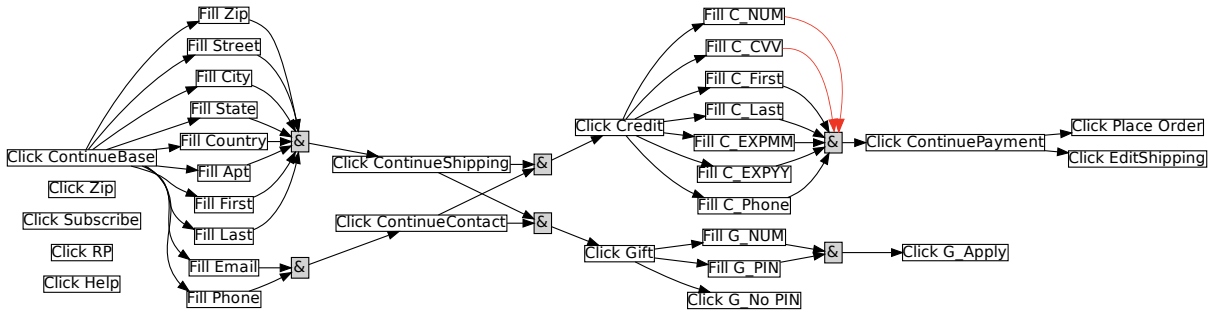


Figure 6: The visualization of the subtask graph inferred by our MTSGL after 1,000 steps of environment interaction on **Walmart** domain. Compared to the ground-truth subtask graph (not available to the agent), there was no error in the nodes and only two missing edges (in red). See Appendix for the progression of the inferred subtask graph with varying adaptation steps.

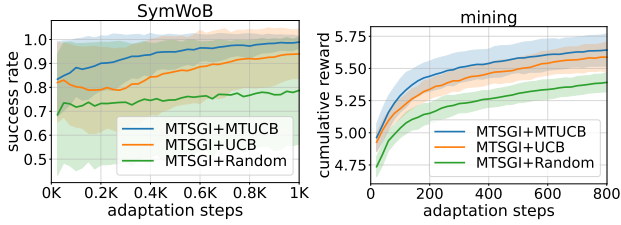


Figure 7: Comparison of different exploration strategies for MTSGL used in adaptation phase for *SymWoB* and *Mining*.

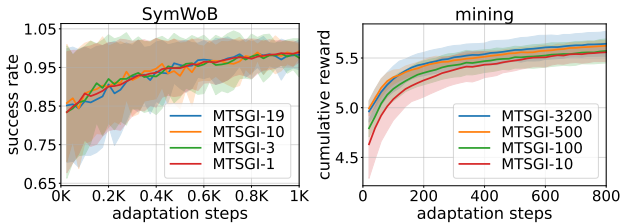


Figure 8: Comparison of different number of priors for MTSGL on *SymWoB* and *Mining*.

early stage of adaptation. This clearly demonstrates that the MTSGL can perform more accurate task inference due to the prior learned from the training tasks.

6.5 Ablation Study: Effect of Exploration Strategy

In this section, we investigate the effect of various exploration strategies on the performance of MTSGL. We compared the following three adaptation policies:

- Random: A policy that uniformly randomly executes any eligible subtask.
- UCB: The UCB policy defined in Equation (2) that aims to execute the novel subtask. The exploration parameters are initialized to zero when a new task is sampled.
- MTUCB (Ours): Our multi-task extension of UCB policy. When a new task is sampled, the exploration parameter is initialized with those of the sampled prior.

Figure 7 summarizes the result on *SymWoB* and *Mining* domain, respectively. Using the more sophisticated exploration policy such as MTSGL+UCB or MTSGL+MTUCB improved the performance of MTSGL compared to MTSGL+Random,

which was also observed in Sohn et al. (2019). This is because better exploration helps the adaptation policy collect more data for logic induction by executing more diverse subtasks. In turn, this results in more accurate subtask graph inference and better performance. Also, MTSGL+MTUCB outperforms MTSGL+UCB on both domains. This indicates that transferring the exploration parameters makes the agent’s exploration more efficient in meta-testing. Intuitively, the transferred exploration counts inform the agent which subtasks were *under-explored* during meta-training, such that the agent can focus more on exploring those in meta-testing.

6.6 Ablation Study: Effect of the prior set size

MTSGL learns the prior from the training tasks. We investigated how many training tasks are required for MTSGL to learn a good prior for transfer learning. Figure 8 compares the performance of MTSGL with the varying number of training tasks: 1, 3, 10, 19 tasks for *SymWoB* and 10, 100, 500, 3200 tasks for *Mining*. The training tasks are randomly sub-sampled from the entire training set. The result shows that training on a larger number of tasks generally improves the performance, but it saturates (e.g., $|\mathcal{M}^{\text{train}}| = 1$ for *SymWoB* and $|\mathcal{M}^{\text{train}}| = 500$ for *Mining*). *Mining* generally requires more number of training tasks than *SymWoB* because the agent is required to solve 440 different tasks in *Mining* while *SymWoB* was evaluated on 20 tasks; the agent is required to capture a wider range of task distribution in *Mining* than *SymWoB*. Also, we note that MTSGL can still adapt much more efficiently than all other baseline methods even when only a small number of training tasks are available (e.g., one task for *SymWoB* and ten tasks for *Mining*).

7 Conclusion

We introduce a multi-task RL extension of the subtask graph inference framework that can quickly adapt to the unseen tasks by modeling the prior of subtask graph from the training tasks and transferring it to the test tasks. The empirical results demonstrate that our MTSGL achieves strong zero-shot and few-shot generalization performance on 2D grid-world and complex web navigation domains by transferring the common knowledge learned in the training tasks to the unseen ones in terms of subtask graph.

In this work, we have assumed that the subtasks and the corresponding options are pre-learned and that the environment

provides a high-level status of each subtask (e.g., whether the web element is filled in with the correct information). In future work, our approach may be extended to a more general setting where the relevant subtask structure is fully learned from (visual) observations, and the corresponding options are autonomously discovered.

References

- Andreas, J.; Klein, D.; and Levine, S. 2017. Modular Multi-task Reinforcement Learning with Policy Sketches. In *ICML*.
- Breiman, L. 1984. *Classification and regression trees*. Routledge.
- Chaplot, D. S.; Sathyendra, K. M.; Pasumarthi, R. K.; Rajagopal, D.; and Salakhutdinov, R. 2018. Gated-Attention Architectures for Task-Oriented Language Grounding. In *AAAI*.
- Chen, Z.; Badrinarayanan, V.; Lee, C.-Y.; and Rabinovich, A. 2018. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *International Conference on Machine Learning*, 794–803. PMLR.
- Duan, Y.; Schulman, J.; Chen, X.; Bartlett, P. L.; Sutskever, I.; and Abbeel, P. 2016. RL²: Fast Reinforcement Learning via Slow Reinforcement Learning. *arXiv preprint arXiv:1611.02779*.
- Finn, C.; Abbeel, P.; and Levine, S. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, 1126–1135. JMLR. org.
- Finn, C.; Xu, K.; and Levine, S. 2018. Probabilistic Model-Agnostic Meta-Learning. In *NeurIPS*, 9516–9527.
- Ghazanfari, B.; and Taylor, M. E. 2017. Autonomous extracting a hierarchical structure of tasks in reinforcement learning and multi-task reinforcement learning. *arXiv preprint arXiv:1709.04579*.
- Gupta, A.; Mendonca, R.; Liu, Y.; Abbeel, P.; and Levine, S. 2018. Meta-Reinforcement Learning of Structured Exploration Strategies. *arXiv preprint arXiv:1802.07245*.
- Gur, I.; Rueckert, U.; Faust, A.; and Hakkani-Tur, D. 2018. Learning to navigate the web. *arXiv preprint arXiv:1812.09195*.
- Hausman, K.; Springenberg, J. T.; Wang, Z.; Heess, N.; and Riedmiller, M. 2018. Learning an embedding space for transferable robot skills. In *International Conference on Learning Representations*.
- Hochreiter, S.; Younger, A. S.; and Conwell, P. R. 2001. Learning to learn using gradient descent. In *International Conference on Artificial Neural Networks*, 87–94. Springer.
- Hoffman, M.; Shahriari, B.; Aslanides, J.; Barth-Maron, G.; Behbahani, F.; Norman, T.; Abdolmaleki, A.; Cassirer, A.; Yang, F.; Baumli, K.; Henderson, S.; Novikov, A.; Colmenarejo, S. G.; Cabi, S.; Gulcehre, C.; Paine, T. L.; Cowie, A.; Wang, Z.; Piot, B.; and de Freitas, N. 2020. Acme: A Research Framework for Distributed Reinforcement Learning. *arXiv preprint arXiv:2006.00979*.
- Huang, D.-A.; Nair, S.; Xu, D.; Zhu, Y.; Garg, A.; Fei-Fei, L.; Savarese, S.; and Niekles, J. C. 2018. Neural task graphs: Generalizing to unseen tasks from a single video demonstration. *arXiv preprint arXiv:1807.03480*.
- Jia, S.; Kiros, J. R.; and Ba, J. 2019. DOM-Q-NET: Grounded RL on Structured Language. In *International Conference on Learning Representations*.
- Kim, T.; Yoon, J.; Dia, O.; Kim, S.; Bengio, Y.; and Ahn, S. 2018. Bayesian Model-Agnostic Meta-Learning. *arXiv preprint arXiv:1806.03836*.
- Konidaris, G.; and Barto, A. 2006. Autonomous shaping: Knowledge transfer in reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, 489–496.
- Lazaric, A. 2012. Transfer in reinforcement learning: a framework and a survey. In *Reinforcement Learning*, 143–173. Springer.
- Lazaric, A.; Restelli, M.; and Bonarini, A. 2008. Transfer of samples in batch reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*, 544–551.
- Lin, X.; Baweja, H. S.; Kantor, G.; and Held, D. 2019. Adaptive auxiliary task weighting for reinforcement learning. *Advances in neural information processing systems*, 32.
- Liu, C.; Yang, S.; Iaba-Sadiya, S.; Shukla, N.; He, Y.; Zhu, S.-c.; and Chai, J. 2016. Jointly Learning Grounded Task Structures from Language Instruction and Visual Demonstration. In *EMNLP*.
- Liu, E. Z.; Guu, K.; Pasupat, P.; Shi, T.; and Liang, P. 2018. Reinforcement learning on web interfaces using workflow-guided exploration. *arXiv preprint arXiv:1802.08802*.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; et al. 2015. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533.
- Muggleton, S. 1991. Inductive Logic Programming. *New Gen. Comput.*, 8(4): 295–318.
- Nichol, A.; Achiam, J.; and Schulman, J. 2018. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*.
- Oh, J.; Singh, S.; Lee, H.; and Kohli, P. 2017. Zero-shot task generalization with multi-task deep reinforcement learning. In *ICML*.
- Pinto, L.; and Gupta, A. 2017. Learning to push by grasping: Using multiple tasks for effective learning. In *2017 IEEE international conference on robotics and automation (ICRA)*, 2161–2168. IEEE.
- Sacerdoti, E. D. 1975a. The nonlinear nature of plans. Technical report, STANFORD RESEARCH INST MENLO PARK CA.
- Sacerdoti, E. D. 1975b. A structure for plans and behavior. Technical report, SRI INTERNATIONAL MENLO PARK CA ARTIFICIAL INTELLIGENCE CENTER.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Shi, T.; Karpathy, A.; Fan, L.; Hernandez, J.; and Liang, P. 2017. World of bits: An open-domain platform for web-based

agents. In *International Conference on Machine Learning*, 3135–3144.

Silver, D.; Schrittwieser, J.; Simonyan, K.; Antonoglou, I.; Huang, A.; Guez, A.; Hubert, T.; Baker, L.; Lai, M.; Bolton, A.; et al. 2017. Mastering the game of go without human knowledge. *nature*, 550(7676): 354–359.

Sohn, S.; Oh, J.; and Lee, H. 2018. Hierarchical Reinforcement Learning for Zero-shot Generalization with Subtask Dependencies. In *NeurIPS*, 7156–7166.

Sohn, S.; Woo, H.; Choi, J.; and Lee, H. 2019. Meta Reinforcement Learning with Autonomous Inference of Subtask Dependencies. In *International Conference on Learning Representations*.

Sutton, R. S.; Precup, D.; and Singh, S. 1999. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2): 181–211.

Tate, A. 1977. Generating project networks. In *Proceedings of the 5th international joint conference on Artificial intelligence-Volume 2*, 888–893. Morgan Kaufmann Publishers Inc.

Taylor, M. E.; and Stone, P. 2009. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10(7).

Vinyals, O.; Babuschkin, I.; Czarnecki, W. M.; Mathieu, M.; Dudzik, A.; Chung, J.; Choi, D. H.; Powell, R.; Ewalds, T.; Georgiev, P.; et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782): 350–354.

Wang, J. X.; Kurth-Nelson, Z.; Tirumala, D.; Soyer, H.; Leibo, J. Z.; Munos, R.; Blundell, C.; Kumaran, D.; and Botvinick, M. 2016. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*.

Wilson, A.; Fern, A.; Ray, S.; and Tadepalli, P. 2007. Multi-task reinforcement learning: a hierarchical bayesian approach. In *Proceedings of the 24th international conference on Machine learning*, 1015–1022.

Xu, D.; Nair, S.; Zhu, Y.; Gao, J.; Garg, A.; Fei-Fei, L.; and Savarese, S. 2017. Neural Task Programming: Learning to Generalize Across Hierarchical Tasks. *arXiv preprint arXiv:1710.01813*.

Yu, H.; Zhang, H.; and Xu, W. 2017. A Deep Compositional Framework for Human-like Language Acquisition in Virtual Environment. *arXiv preprint arXiv:1703.09831*.

Zhang, Y.; and Yeung, D.-Y. 2014. A regularization approach to learning task relationships in multitask learning. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 8(3): 1–31.

Appendix: Fast Inference and Transfer of Compositional Task Structures for Few-shot Task Generalization

A Derivation of subtask graph inference

Following Sohn et al. (2019), we formulate the problem of inferring the precondition G_c and the subtask reward G_r as a maximum likelihood estimation (MLE) problem. Let $\tau_K = \{\mathbf{s}_1, \mathbf{o}_1, r_1, d_1, \dots, \mathbf{s}_K\}$ be an adaptation trajectory of the adaptation policy $\pi_\theta^{\text{adapt}}$ for K steps 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 . Consider a subtask graph G with subtask reward G_r and precondition G_c . The maximum-likelihood estimate (MLE) of latent variables G , conditioned on the trajectory τ_H can be written as

$$\widehat{G}^{\text{MLE}} = \arg \max_{G_c, G_r} p(\tau_K | G_c, G_r). \quad (14)$$

The likelihood term can be expanded as

$$p(\tau_K | G_c, G_r) = p(\mathbf{s}_1 | G_c) \prod_{t=1}^K \pi_\theta(\mathbf{o}_t | \tau_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{o}_t, G_c) p(r_t | \mathbf{s}_t, \mathbf{o}_t, G_r) p(d_t | \mathbf{s}_t, \mathbf{o}_t) \quad (15)$$

$$\propto p(\mathbf{s}_1 | G_c) \prod_{t=1}^K p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{o}_t, G_c) p(r_t | \mathbf{s}_t, \mathbf{o}_t, G_r), \quad (16)$$

where we dropped the terms that are independent of G . From the definitions in **Subtask Graph Inference Problem** section, 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

$$\widehat{G}^{\text{MLE}} = (\widehat{G}_c^{\text{MLE}}, \widehat{G}_r^{\text{MLE}}) = \left(\arg \max_{G_c} \prod_{t=1}^K p(\mathbf{e}_t | \mathbf{x}_t, G_c), \arg \max_{G_r} \prod_{t=1}^K p(r_t | \mathbf{e}_t, \mathbf{o}_t, G_r) \right). \quad (17)$$

Below we explain how to compute the estimate of preconditions $\widehat{G}_c^{\text{MLE}}$ and subtask rewards $\widehat{G}_r^{\text{MLE}}$.

Precondition inference via logic induction From the definition in **Subtask Graph Inference Problem** section, the mapping from a precondition G_c and a completion vector \mathbf{x} to an eligibility vector $\mathbf{e} = f_{G_c}(\mathbf{x})$ is a deterministic function (*i.e.*, precondition function). Therefore, we can rewrite $\widehat{G}_c^{\text{MLE}}$ in Eq.(17) as:

$$\widehat{G}_c^{\text{MLE}} = \arg \max_{G_c} \prod_{t=1}^K \mathbb{I}(\mathbf{e}_t = f_{G_c}(\mathbf{x}_t)), \quad (18)$$

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. This problem of finding a boolean function that satisfies all the indicator functions in Eq.(18) (*i.e.*, $\prod_{t=1}^K \mathbb{I}(\mathbf{e}_t = f_{G_c}(\mathbf{x}_t)) = 1$) is formulated as an *inductive logic programming* (ILP) problem (Muggleton 1991). Specifically, $\{\mathbf{x}_t\}_{t=1}^K$ forms binary vector inputs to programs, and $\{\mathbf{e}_t^i\}_{t=1}^K$ 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 1984). 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.

Subtask reward inference To infer the subtask reward function $\widehat{G}_r^{\text{MLE}}$ in Eq.(17), we model each component of subtask reward as a Gaussian distribution $G_r^i \sim \mathcal{N}(\widehat{\mu}^i, \widehat{\sigma}^i)$. Then, $\widehat{\mu}_{\text{MLE}}^i$ becomes the empirical mean of the rewards received after taking the eligible option \mathcal{O}^i in the trajectory τ_K :

$$\widehat{G}_r^{\text{MLE}, i} = \widehat{\mu}_{\text{MLE}}^i = \mathbb{E}[r_t | \mathbf{o}_t = \mathcal{O}^i, \mathbf{e}_t^i = 1] = \frac{\sum_{t=1}^K r_t \mathbb{I}(\mathbf{o}_t = \mathcal{O}^i, \mathbf{e}_t^i = 1)}{\sum_{t=1}^K \mathbb{I}(\mathbf{o}_t = \mathcal{O}^i, \mathbf{e}_t^i = 1)}. \quad (19)$$

B Details of GRProp policy

For self-containedness, we repeat the description of GRProp policy in Sohn et al. (2019).

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 **Subtask Graph Inference Problem** section for definitions). Then, the sum of reward until time-step t is given as:

$$U_t = G_{\mathbf{r}}^{\top} \mathbf{x}_t. \quad (20)$$

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:

$$\widehat{U}_t = G_{\mathbf{r}}^{\top} (\mathbf{x}_t + \mathbf{e}_t) / 2. \quad (21)$$

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:

$$e_t^i = \text{OR}_{j \in \text{Child}_i} \left(y_{AND}^j \right), \quad y_{AND}^j = \text{AND}_{k \in \text{Child}_j} \left(\widehat{x}_t^{j,k} \right), \quad \widehat{x}_t^{j,k} = x_t^k w^{j,k} + \text{NOT}(x_t^k) (1 - w^{j,k}), \quad (22)$$

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 Equation (22) are substituted by the smoothed counterparts as follows:

$$p^i = \lambda_{\text{or}} \widetilde{e}^i + (1 - \lambda_{\text{or}}) x^i, \quad (23)$$

$$\widetilde{e}^i = \widetilde{\text{OR}}_{j \in \text{Child}_i} \left(\widetilde{y}_{AND}^j \right), \quad (24)$$

$$\widetilde{y}_{AND}^j = \widetilde{\text{AND}}_{k \in \text{Child}_j} \left(\widehat{x}^{j,k} \right), \quad (25)$$

$$\widehat{x}^{j,k} = w^{j,k} p^k + (1 - w^{j,k}) \widetilde{\text{NOT}}(p^k), \quad (26)$$

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

$$\widetilde{\text{OR}}(\mathbf{x}) = \text{softmax}(w_{\text{or}} \mathbf{x}) \cdot \mathbf{x}, \quad (27)$$

$$\widetilde{\text{AND}}(\mathbf{x}) = \frac{\zeta(\mathbf{x}, w_{\text{and}})}{\zeta(\|\mathbf{x}\|, w_{\text{and}})}, \quad (28)$$

$$\widetilde{\text{NOT}}(\mathbf{x}) = -w_{\text{not}} \mathbf{x}, \quad (29)$$

$\|\mathbf{x}\| = d$, $\zeta(\mathbf{x}, \beta) = \frac{1}{\beta} \log(1 + \exp(\beta \mathbf{x}))$ 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, Oh, and Lee 2018) to softmax and softplus functions for better performance. With the smoothed operations, the sum of smoothed and modified reward is given as:

$$\widetilde{U}_t = G_{\mathbf{r}}^{\top} \mathbf{p}, \quad (30)$$

where $\mathbf{p} = [p^1, \dots, p^d]$ and p^i is computed from Equation (23). Finally, the graph reward propagation policy is a softmax policy,

$$\pi(\mathbf{o}_t | G, \mathbf{x}_t) = \text{Softmax} \left(T \nabla_{\mathbf{x}_t} \widetilde{U}_t \right) = \text{Softmax} \left(T G_{\mathbf{r}}^{\top} (\lambda_{\text{or}} \nabla_{\mathbf{x}_t} \widetilde{\mathbf{e}}_t + (1 - \lambda_{\text{or}})) \right), \quad (31)$$

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.

C Derivation of Multi-task subtask graph inference

Let τ be the adaptation trajectory of the current task \mathcal{M}_G , and $\mathcal{T}^p = \{\tau_1^p, \dots, \tau_{|\mathcal{T}^p|}^p\}$ be the adaptation trajectories of the seen training tasks.

Then, from Bayesian rule, we have

$$\pi(o|s, \tau, \mathcal{T}^p) = \sum_G \pi(o|s, G)p(G|\tau, \mathcal{T}^p) \quad (32)$$

$$\propto \sum_G \pi(o|s, G)p(\tau|G, \mathcal{T}^p)p(G|\mathcal{T}^p) \quad (33)$$

$$= \sum_G \pi(o|s, G)p(\tau|G)p(G|\mathcal{T}^p) \quad (34)$$

$$\propto \sum_G \pi(o|s, G)p(\tau|G)p(\mathcal{T}^p|G)p(G), \quad (35)$$

where Equation (34) holds because τ and \mathcal{T}^p are independently observed variables. Since summing over all G is computationally intractable, we instead approximate it by computing the sample estimates of π . Specifically, we compute the policy $\pi(o|s, G)$ at the maximum likelihood estimates (MLE) of G for prior and current tasks, that is $G^\tau = \arg \max_G p(\tau|G)$ and $G^p = \arg \max_G p(\mathcal{T}^p|G)$ respectively, and combine them with the weight α :

$$\pi(o|s, \tau, \mathcal{T}^p) \simeq \pi(o|s, G^\tau)^\alpha \pi(o|s, G^p)^{(1-\alpha)}. \quad (36)$$

Finally, we deploy the GRProp policy as a contextual policy:

$$\pi^{\text{eval}}(\cdot|\tau, \mathcal{T}^p) \simeq \text{GRProp}(\cdot|G^\tau)^\alpha \text{GRProp}(\cdot|G^p)^{(1-\alpha)}. \quad (37)$$

D Pseudo-code of our algorithm

The Algorithm 3 below describes the pseudo-code of the meta-training process of our algorithm.

Algorithm 3: Meta-training: learning the prior

Require: Adaptation policy π^{adapt}

Ensure: Prior set \mathcal{T}^p

- 1: $\mathcal{T}^p \leftarrow \emptyset$
 - 2: **for** each task $\mathcal{M} \in \mathcal{M}^{\text{train}}$ **do**
 - 3: Rollout adaptation policy:
 $\tau = \{\mathbf{s}_t, \mathbf{o}_t, r_t, d_t\}_{t=1}^K \sim \pi^{\text{adapt}}$ in task \mathcal{M}
 - 4: Infer subtask graph $G^\tau = \arg \max_G p(\tau|G)$
 - 5: $\pi^{\text{eval}} = \text{GRProp}(G^\tau)$
 - 6: Evaluate the agent: $\tau^{\text{eval}} \sim \pi^{\text{eval}}$ in task \mathcal{M}
 - 7: Update prior $\mathcal{T}^p \leftarrow \mathcal{T}^p \cup (G^\tau, \tau)$
 - 8: **end for**
-

E Extended Related Work

Multi-task reinforcement learning. Multi-task reinforcement learning aims to learn an inductive bias that can be shared and used across a variety of related RL tasks to improve the task generalization. Early works mostly focused on the transfer learning oriented approaches (Lazaric 2012; Taylor and Stone 2009) such as instance transfer (Lazaric, Restelli, and Bonarini 2008) or representation transfer (Konidaris and Barto 2006). However, these algorithms rely heavily on the prior knowledge about the allowed task differences. Hausman et al. (2018); Pinto and Gupta (2017); Wilson et al. (2007) proposed to train a multi-task policy with multiple objectives from different tasks. However, the gradients from different tasks may conflict and hurt the training of other tasks. To avoid gradient conflict, Zhang and Yeung (2014); Chen et al. (2018); Lin et al. (2019) proposed to explicitly model the task similarity. However, dynamically modulating the loss or the gradient of RL update often results in the instability in optimization. Our multi-task learning algorithm also takes the transfer learning oriented viewpoint; MTSGL captures and transfers the task knowledge in terms of the subtask graph. However, our work does not make a strong assumption on the task distribution. We only assume that the task is parameterized by unknown subtask graph, which subsumes many existing compositional tasks (*e.g.*, Oh et al. (2017); Andreas, Klein, and Levine (2017); Huang et al. (2018), etc).

Extended - web navigating RL agent. Previous work introduced MiniWoB (Shi et al. 2017) and MiniWoB++ (Liu et al. 2018) benchmarks that are manually curated sets of simulated toy environments for the web navigation problem. They formulated the problem as acting on a page represented as a Document Object Model (DOM), a hierarchy of objects in the page. The agent is trained with human demonstrations and online episodes in an RL loop. Jia, Kiros, and Ba (2019) proposed a graph neural network based DOM encoder and a multi-task formulation of the problem similar to this work. Gur et al. (2018) introduced a manually-designed curriculum learning method and an LSTM based DOM encoder. DOM level representations of web pages pose a significant sim-to-real gap as simulated websites are considerably smaller (100s of nodes) compared to noisy real websites (1000s of nodes). As a result, these models are trained and evaluated on the same simulated environments which is difficult to deploy on real websites. Our work formulates the problem as abstract web navigation on real websites where the objective is to learn a latent subtask dependency graph similar to sitemap of websites. We propose a multi-task training objective that generalizes from a fixed set of real websites to unseen websites without any demonstration, illustrating an agent capable of navigating real websites for the first time.

Planning Approaches for Compositional Task Previous work has tackled the compositional tasks using the Hierarchical Task Network (HTN) planning (Sacerdoti 1975b,a; Tate 1977) in a (single) goal-conditioned RL setting. The HTN allows the agent to reason the tasks at multiple levels of abstraction, when rich knowledge at those abstraction levels are available. Specifically, HTN models the primitive tasks (or the *subtasks* in our terminology) by the precondition and the effects, and aim to find the sequence of actions (or the *options* in our terminology) that execute each subtasks via planning on the HTN. They aim to execute a single goal task, often with assumptions of simpler subtask dependency structures (Ghazanfari and Taylor 2017; Liu et al. 2016) such that the task structure can be constructed from the successful trajectories. Also, they often require expensive searching to find the solution. 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. Moreover, we avoid any expensive searching and propose to use neural network to directly map the task structure into policy. Lastly, we aim to achieve zero-/few-shot task generalization, which is not achievable with the HTN methods since they require the full specification of the action models in the testing.

F Additional experiment results

F.1 Full Experiment Results on the performance of the agents on 20 Websites in *SymWoB*

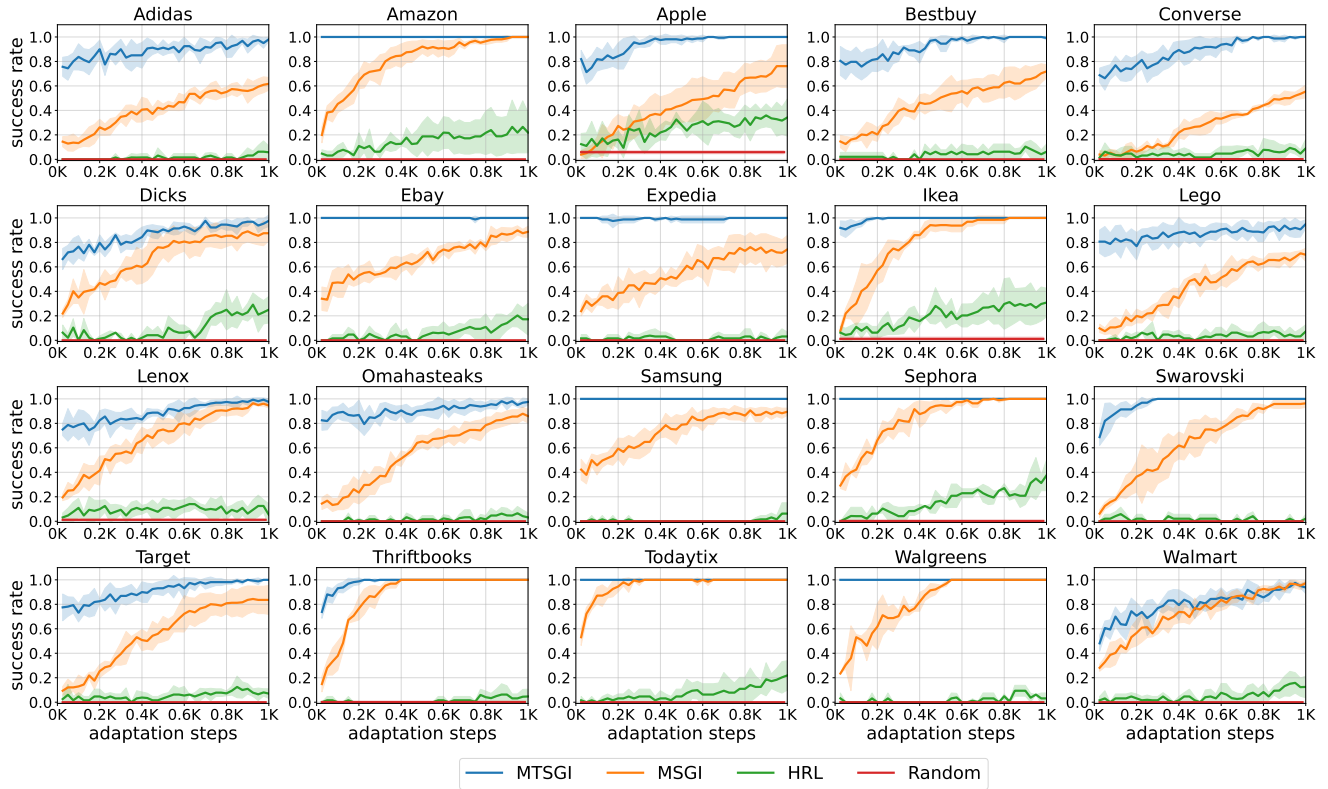


Figure 9: Results of the success rate for compared methods in the test phase with respect to the adaptation steps on 20 environments in *SymWoB* domain.

F.2 Visualization of the multi-task subtask graph inference process

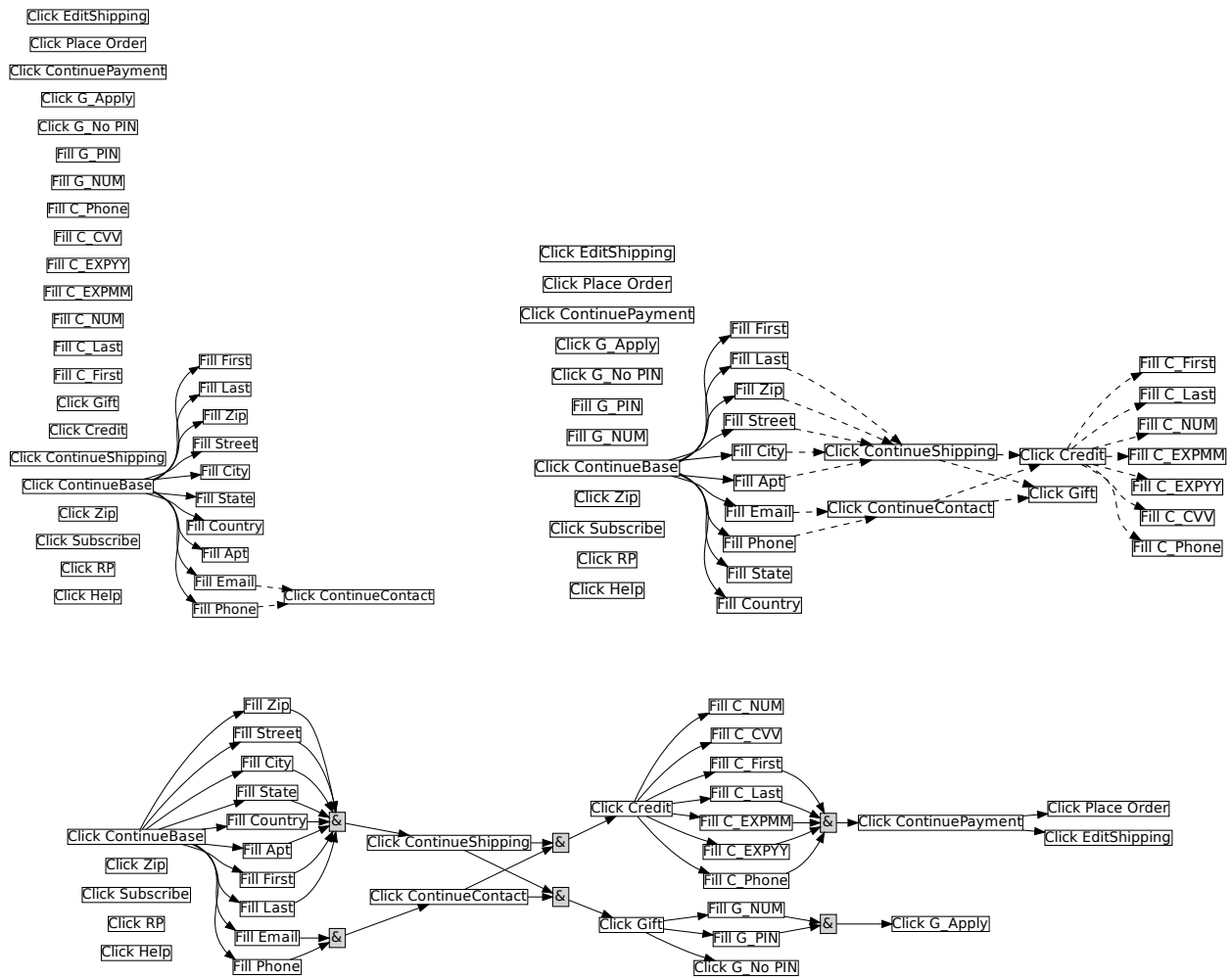


Figure 10: The subtask graphs inferred by our MTSGI with varying adaptation steps on **Walmart** domain: 0 steps (*Top, Left*), 400 steps (*Top, Right*), and 1000 steps (*Bottom*).

Figure 10 qualitatively shows how the multi-task subtask graph inference proceeds over the adaptation. In the beginning (*Top, Left*), the agent has only prior information and the prior provides a partial information about the subtask graph, which is the preconditions in the first webpage. As the agent further explores the webpage (*Top, Right*), the subtask graph gets more accurate, but due to insufficient exploration, there are many missing preconditions, especially for the subtasks in the last webpage. After sufficient adaptation (*Bottom*), our MTSGI can infer quite accurate subtask graph compared to the ground-truth subtask graph in Figure 11; *i.e.*, only missing two preconditions.

F.3 Full Experiment Results on the qualitative evaluation of the inferred subtask graph

Figures 11-30 qualitatively evaluate the task inference of our MTSGI by comparing the inferred subtask graph and the ground-truth. It is clear from the figure that the inferred subtask graphs have only a small portion of missing or redundant edges, while most of the nodes and edges are the same as ground-truth graph.

G Details on *SymWoB* domain

Subtask Graph Setting										
Task	Adidas	Amazon	Apple	BestBuy	Converse	Dick’s	eBay	Expedia	Ikea	Lego
#Subtasks	46	31	43	37	42	39	39	36	39	45
#Distractors	8	4	5	6	6	6	5	5	5	6
Episode length	41	27	40	37	43	37	37	40	37	37
Task	Lenox	Omahasteaks	Samsung	Sephora	Swarovski	Target	Thriftbooks	Todaytix	Walgreens	Walmart
#Subtasks	45	44	42	49	45	39	43	23	38	46
#Distractors	4	6	6	7	7	6	8	3	7	5
Episode length	41	38	41	45	38	37	33	20	50	43

Table 1: The task configuration of the tasks in *SymWoB* domain. Each task is parameterized by different subtask graphs, and the episode length is manually set according to the challengeness of the tasks.

In this paper, we introduce the *SymWoB* domain, which is a challenging symbolic environment that aims to reflect the hierarchical and compositional aspects of the checkout processes in the real-world websites. There are total 20 different *SymWoB* websites that are symbolic implementations of the actual websites: **Adidas**, **Amazon**, **Apple**, **BestBuy**, **Converse**, **Dick’s**, **eBay**, **Expedia**, **Ikea**, **Lego**, **Lenox**, **Omahasteaks**, **Samsung**, **Sephora**, **Swarovski**, **Target**, **Thriftbooks**, **Todaytix**, **Walgreens**, and **Walmart**. All of these websites are generated by analyzing the corresponding real websites and reflecting their key aspects of checkout process. The main goal of each website is to navigate through the web pages within the website by clicking and filling in the web elements with proper information which leads to the final web page that allows the agent to click on the `Place_Order` button, which indicates that the agent has successfully finished the task of checking out.

G.1 Implementation detail

In this section, we describe the detailed process of implementing an existing website into a symbolic version. We first fill out the shopping cart with random products, and we proceed until placing the order on the actual website. During the process, we extract all the interactable web elements on the webpage. We repeat this for all the websites and form a shared subtask pool where similar web elements in different websites that have same functionality are mapped to the same subtask in the shared subtask pool. Then, we extract the precondition relationship between subtasks from the website and form the edges in the subtask graph accordingly. Finally, we implement the termination condition and the subtask reward to the failure distractor (See Appendix G.3) and the goal subtasks.

G.2 Comparison of the websites

The agent’s goal on every website is the same, that is placing the checkout order. However, the underlying subtask graphs, or task structure, of the websites are quite diverse, making the task much more challenging for the agent. Figures 11-30 visualize the ground truth subtask graph of all the websites. One of the major sources of diversity in subtask graphs is in the various ordering of the web pages. In a typical website’s checkout process, some of the most common web pages include the shipping, billing, and payment web pages, each of which has a collection of corresponding subtasks. In Figure 11, for example, the shipping web page is represented by the collection of the subtasks on the left side from `Fill.Zip` to `Fill.Last` and `Click.ContinueShipping`, and these come *before* the payment web page that is represented by the subtasks on the right side from `Click.Credit` to `Click.ContinuePayment`. On the other hand, in Figure 17 and Figure 16, the similar web pages are either connected in a different ordering, from payment to shipping web page, or placed on the same line side by side. Since the web pages can vary on how they are ordered, it allows the subtask graphs to have a variety of shapes such as deep and narrow as in Figure 19 or wide and shallow as in Figure 16. Different shape of the subtask graphs means different precondition between the tasks, making it non-trivial for the agent to transfer its knowledge about one to the other.

Another major source of diversity is the number of web elements in each web page. Let’s compare the web elements of the shipping web page in Figure 13 and Figure 14. These are the subtasks that are connected to `Click.ContinueShipping` and as well as itself. We can see that the two websites do not have the same number of the web elements for the shipping web pages: the **Converse** website requires more shipping information to be filled out than the **Dick’s** website. Such variety in the number of web elements, or subtasks, allows the subtask graphs of the websites to have diverse preconditions as well.

G.3 Distractor subtasks

In addition to the different task structures among the websites, there are also *distractor* subtasks in the websites that introduces challenging components of navigating the real-world websites. There are two different types of distractor subtasks: the one that terminates the current episode with a negative reward and the another one that has no effect. The former, which we also call it the *failure* distractor subtask, represents the web elements that lead the agent to some external web pages like `Help` or `Terms_of_Use` button. The latter is just called the distractor subtask, where executing the subtask does not contribute to progressing toward the goal (*e.g.*, `Click EditShipping` subtask in **Converse**). Each website has varying number of

distractor subtasks and along with the shallowness of the task structure, the number of distractor subtasks significantly affects the difficulty of the task.

H Details of the agent implementation

We implement all of our algorithms, including both MTSGI and the baselines, on top of the recently introduced RL framework called Acme (Hoffman et al. 2020).

H.1 MSGI

Similar to the MTSGI agent, the MSGI agent uses the soft-version of UCB exploration policy as the adaptation policy, instead of its original hard-version due to a better performance. Furthermore, MSGI also uses inductive logic programming (ILP) in order to infer the subtask graph of the current task. But unlike the MTSGI agent that learns prior across the tasks through multi-task learning, the MSGI agent does not exploit the subtask graphs inferred from the previous tasks.

H.2 HRL

The hierarchical RL (HRL) agent is an option-based agent that can execute temporally extended actions. For the *Mining* domain, where there is a spatial component in the observation input, we use convolutional neural network (CNN) in order to encode the spatial information and get concatenated along with the other additional inputs, which are encoded using fully-connected (FC) networks. The concatenated embedding then gets passed on to GRU, which is followed by two separate heads for value and policy function outputs.

The details of the HRL architecture for *Mining* domain are: Conv1(16x1x1-1)-Conv2-(32x3x3-1)-Conv3(64x3x3-1)-Conv4(32x3x3-1)-Flatten-FC(256)-GRU(512). The HRL architecture for *SymWoB* domain is almost identical except it replaces the CNN module with fully-connected layers for processing non-spatial information: FC(64)-FC(64)-FC(50)-FC(50)-GRU(512). We use ReLU activation function in all the layers.

For training, we use multi-step actor-critic (A2C) in order to train HRL. We use multi-step learning of $n = 10$, learning rate 0.002, entropy loss weight of 0.01 and critic loss weight of 0.5. We use RMSProp optimizer to train the networks, where its decay rate is set to 0.99 and epsilon to 0.00001. We also clip the gradient by setting the norm equal to 1.0. Most importantly, the network parameters of the HRL agent gets reset for every new task since HRL is not a meta-RL agent.

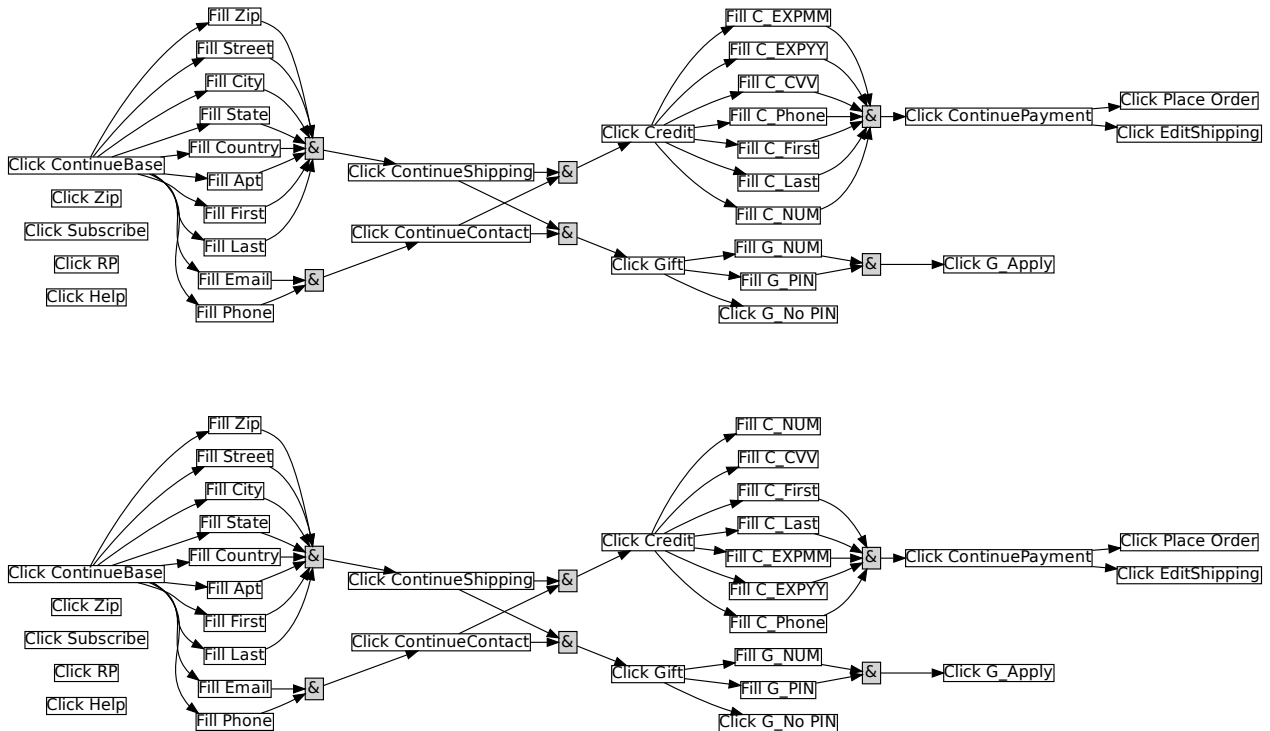


Figure 11: (Top) The ground-truth and (Bottom) the inferred subtask graphs of Walmart website.

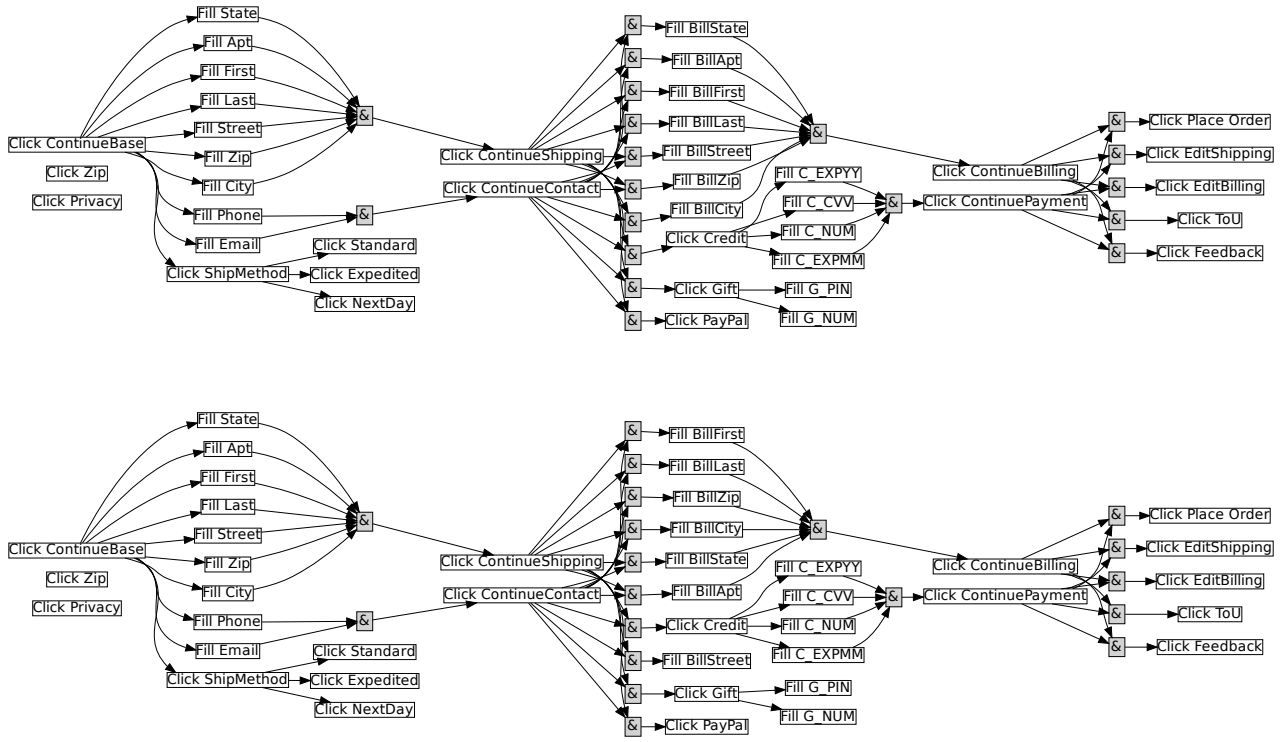


Figure 12: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Converse** website.



Figure 13: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Dick's** website.

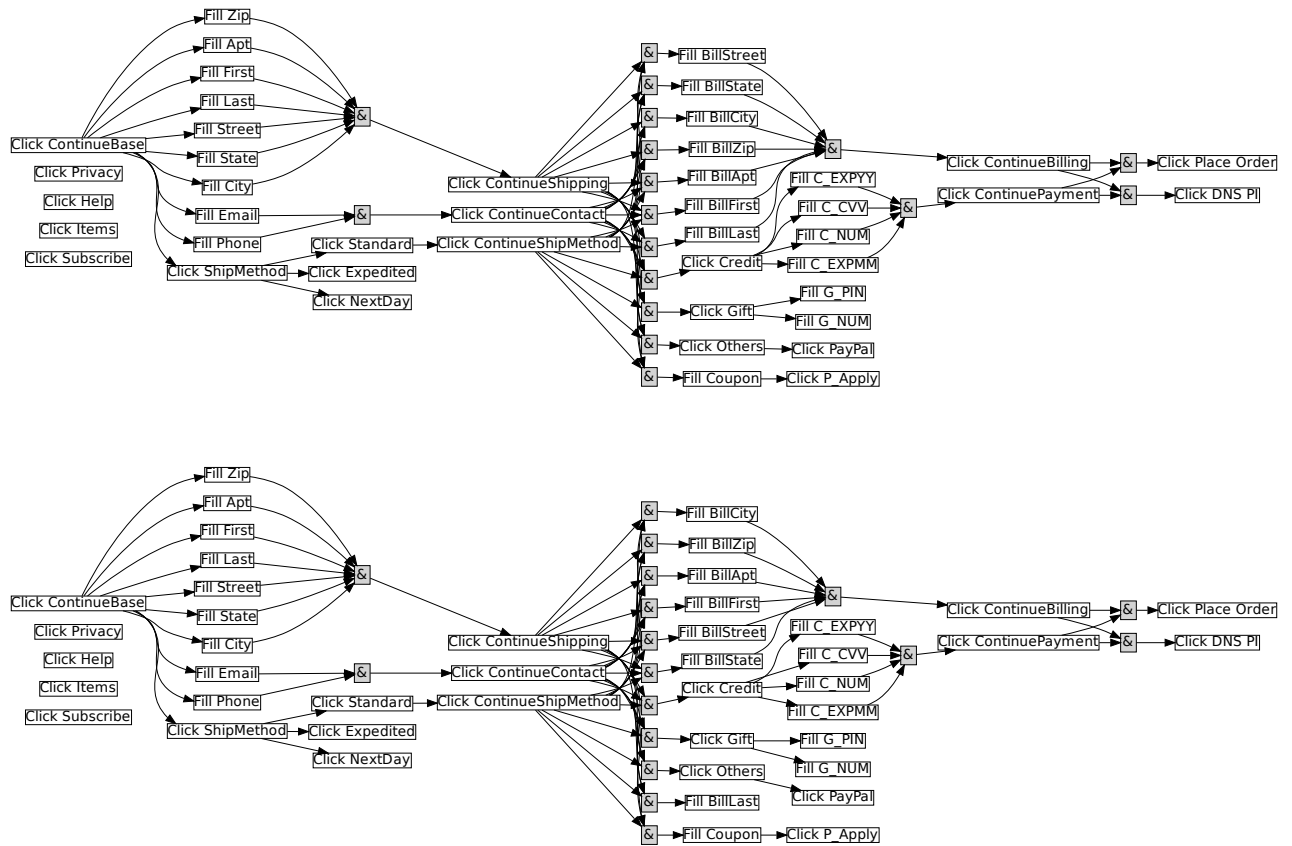


Figure 14: (Top) The ground-truth and (Bottom) the inferred subtask graphs of BestBuy website.

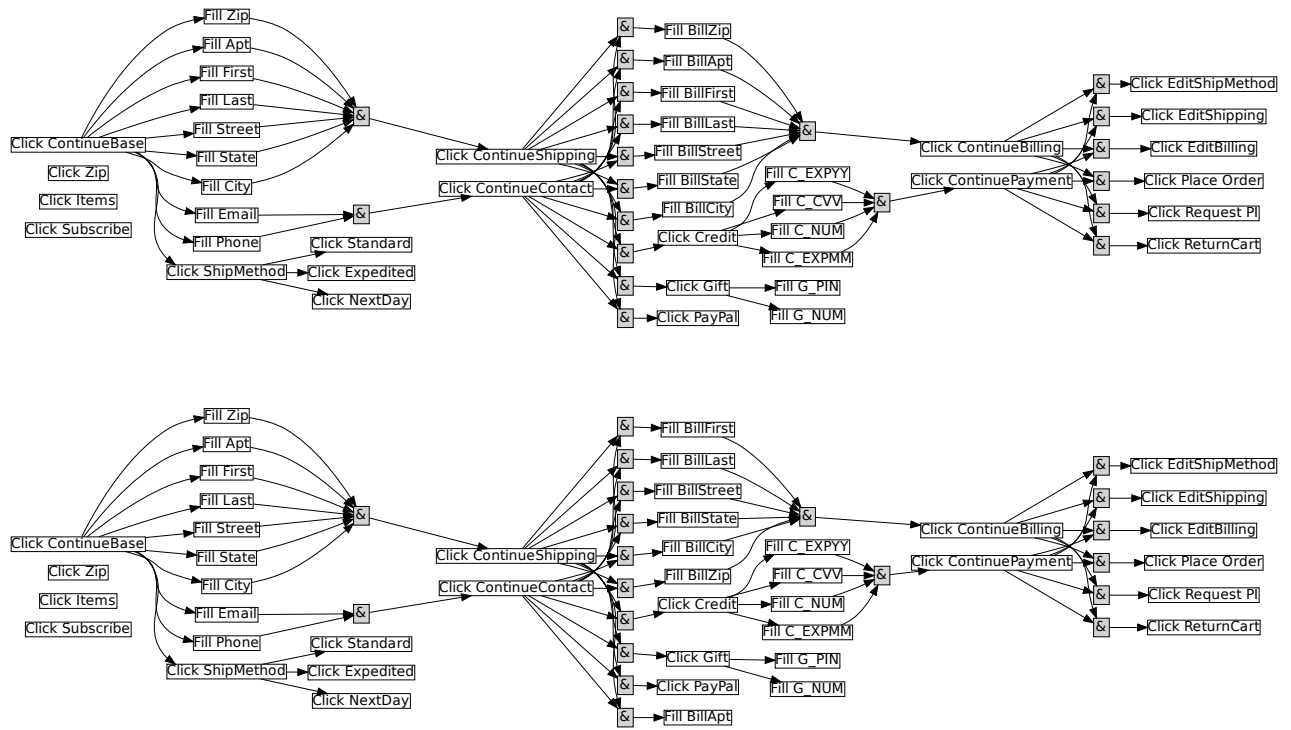


Figure 15: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Apple** website.

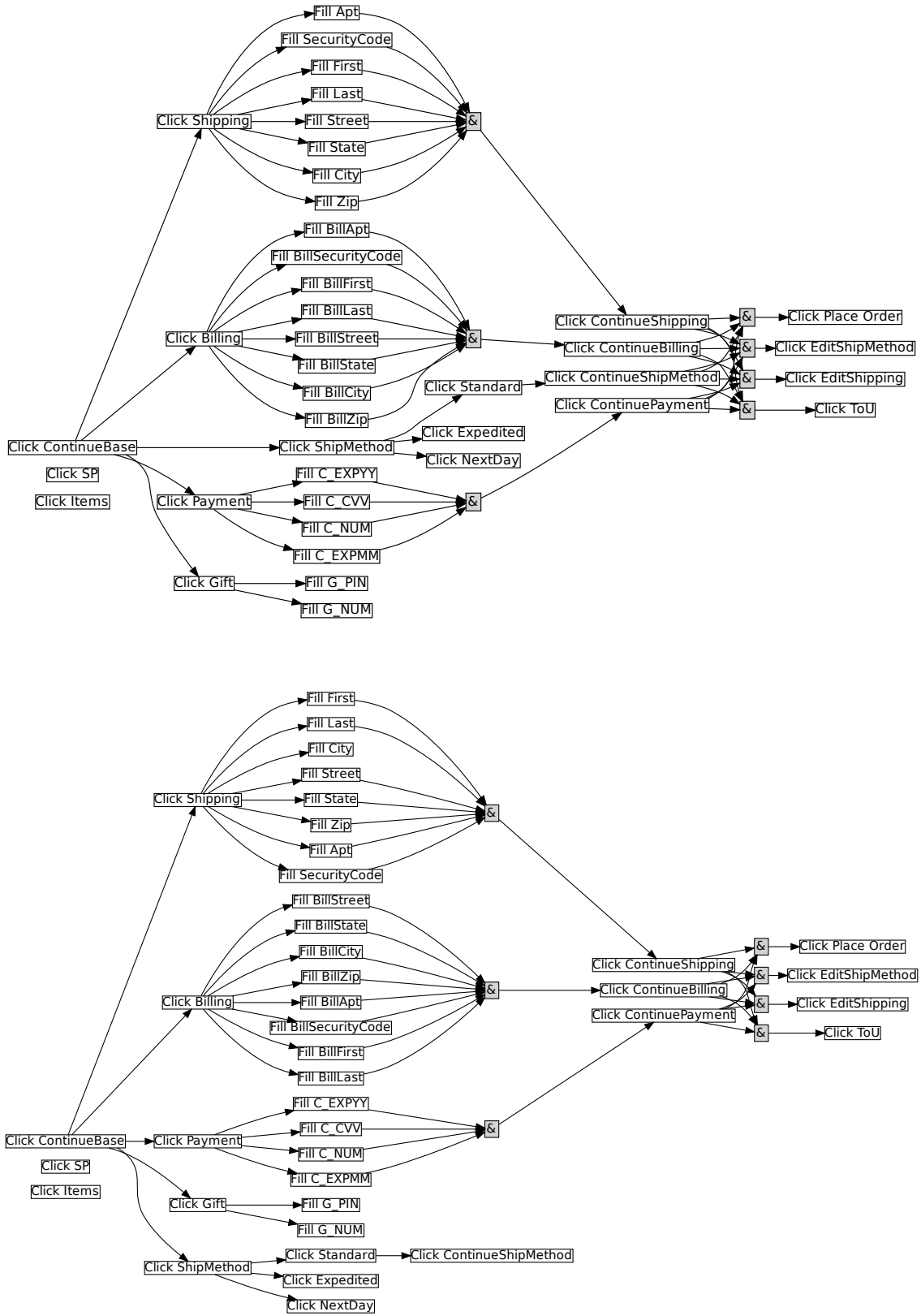


Figure 16: (Top) The ground-truth and (Bottom) the inferred subtask graphs of Amazon website.

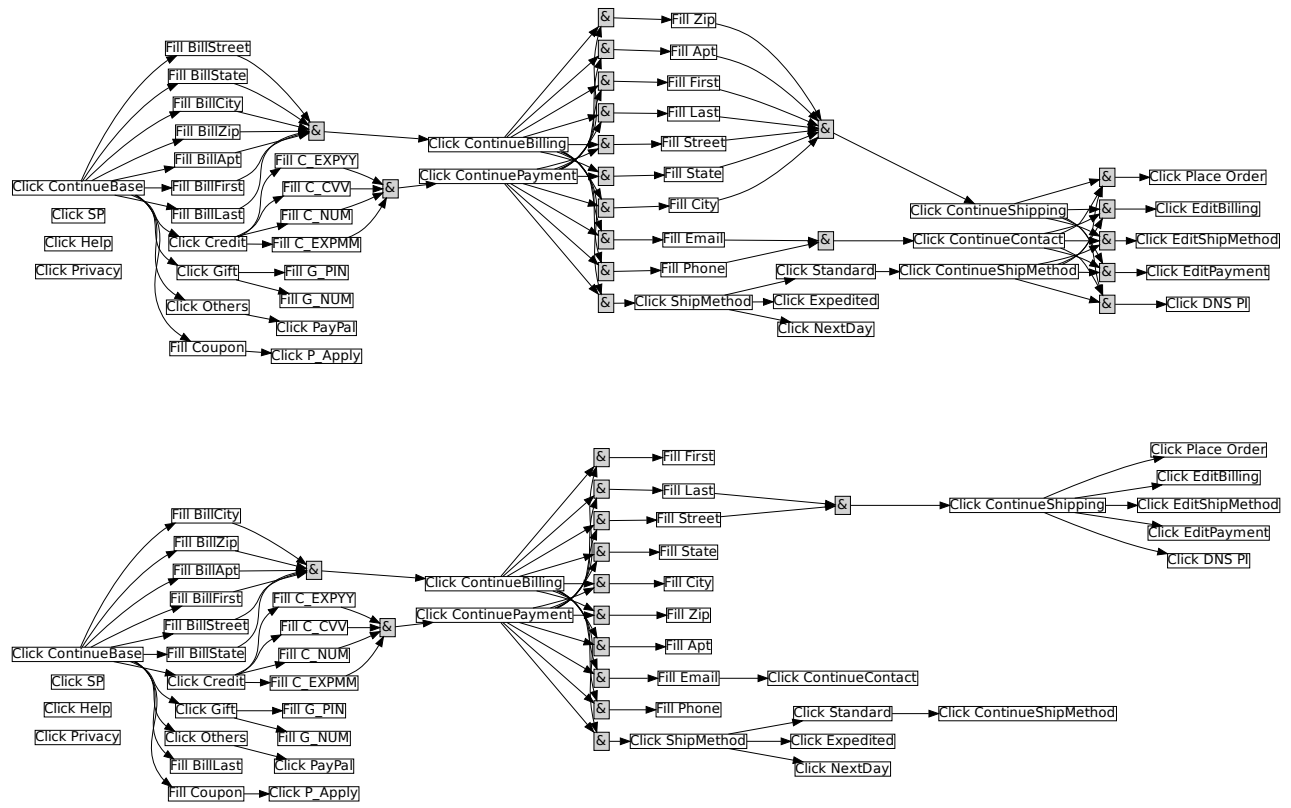


Figure 17: (Top) The ground-truth and (Bottom) the inferred subtask graphs of Samsung website.



Figure 18: (Top) The ground-truth and (Bottom) the inferred subtask graphs of eBay website.

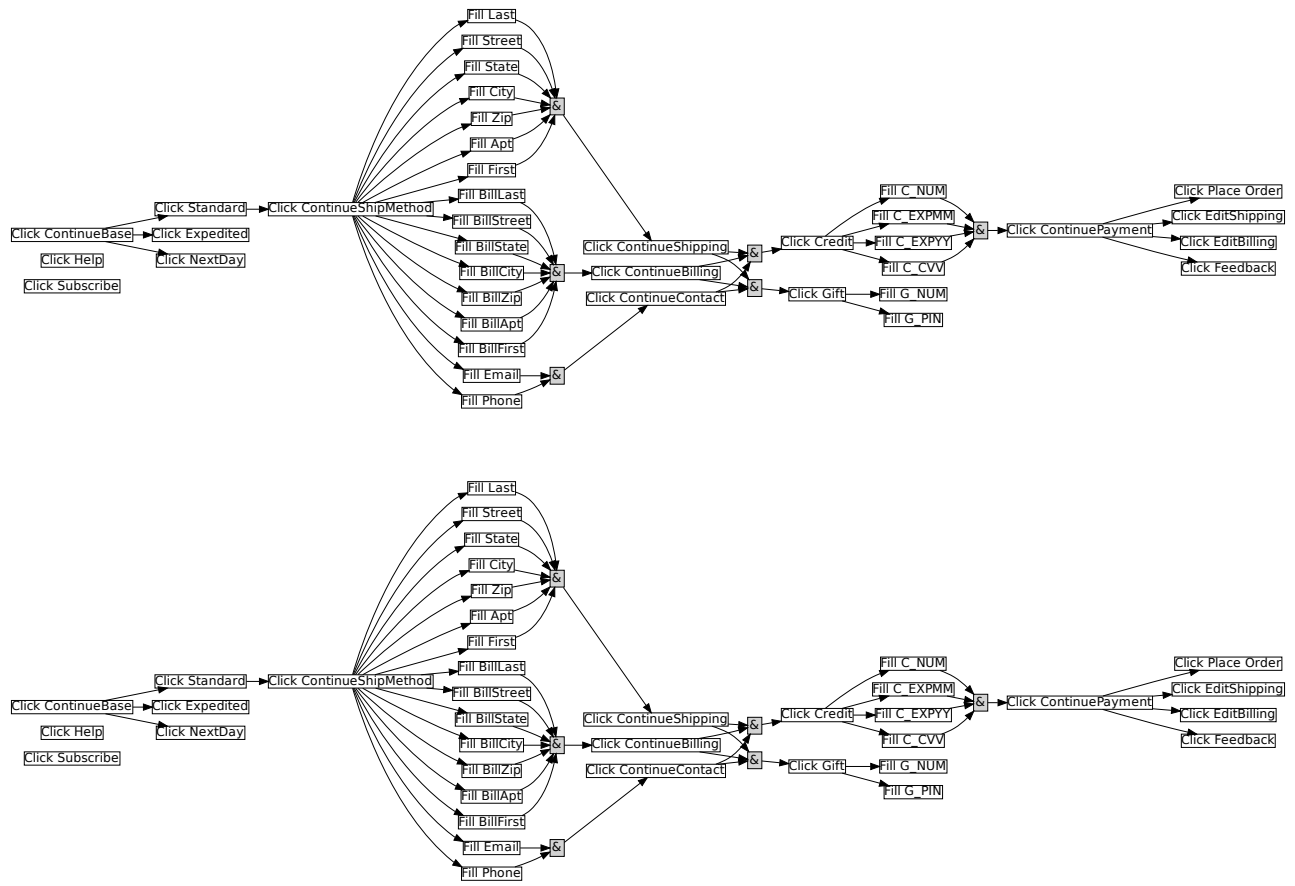


Figure 19: (Top) The ground-truth and (Bottom) the inferred subtask graphs of Ikea website.

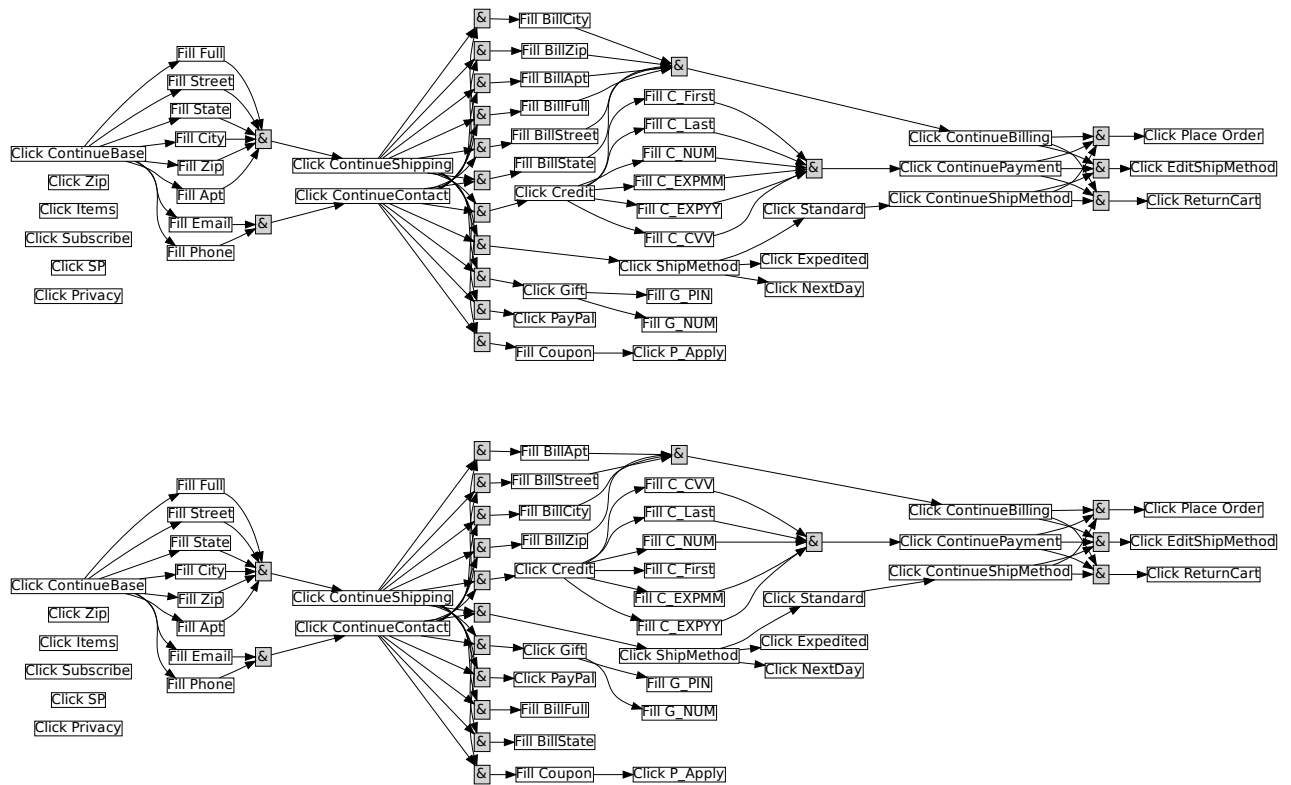


Figure 20: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Target** domain.

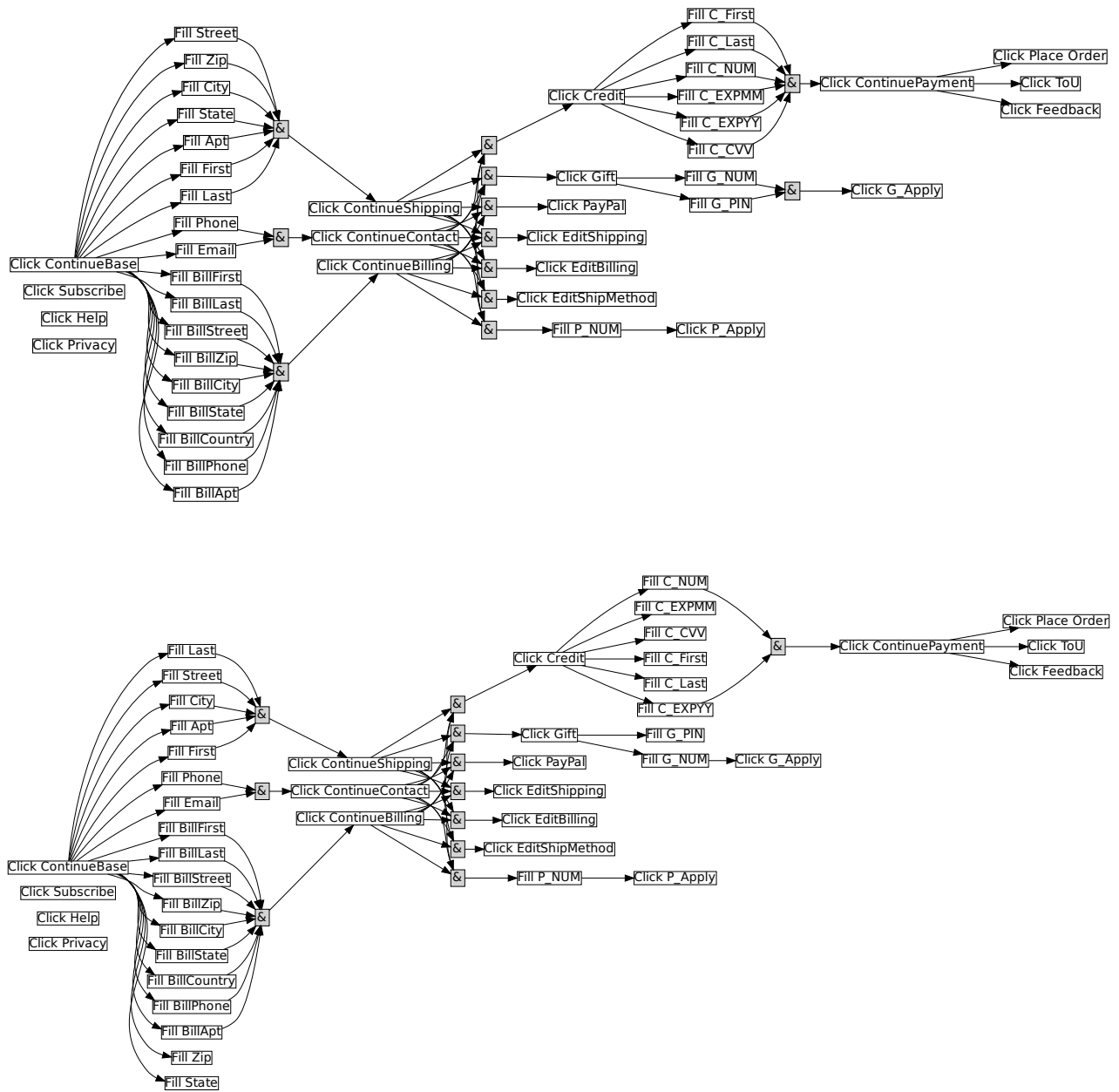


Figure 21: (Top) The ground-truth and (Bottom) the inferred subtask graphs of Adidas domain.

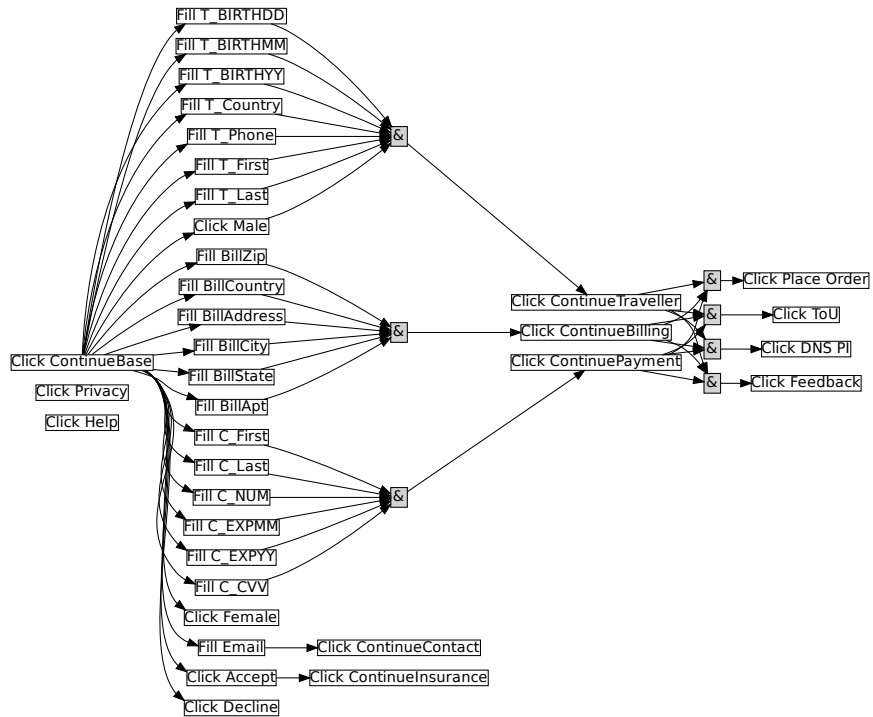
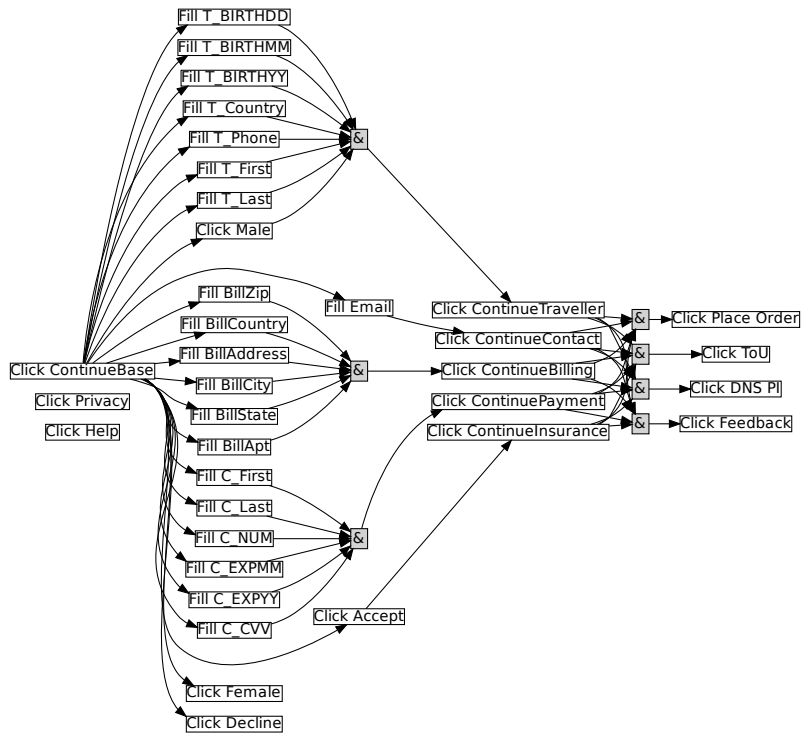


Figure 22: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Expedia** domain.

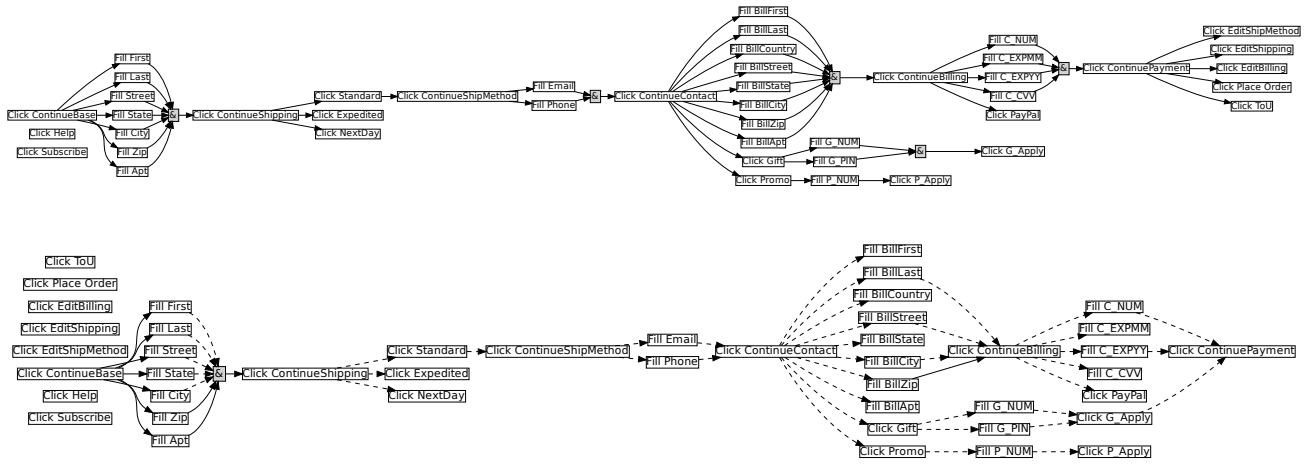


Figure 23: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Lego** domain.

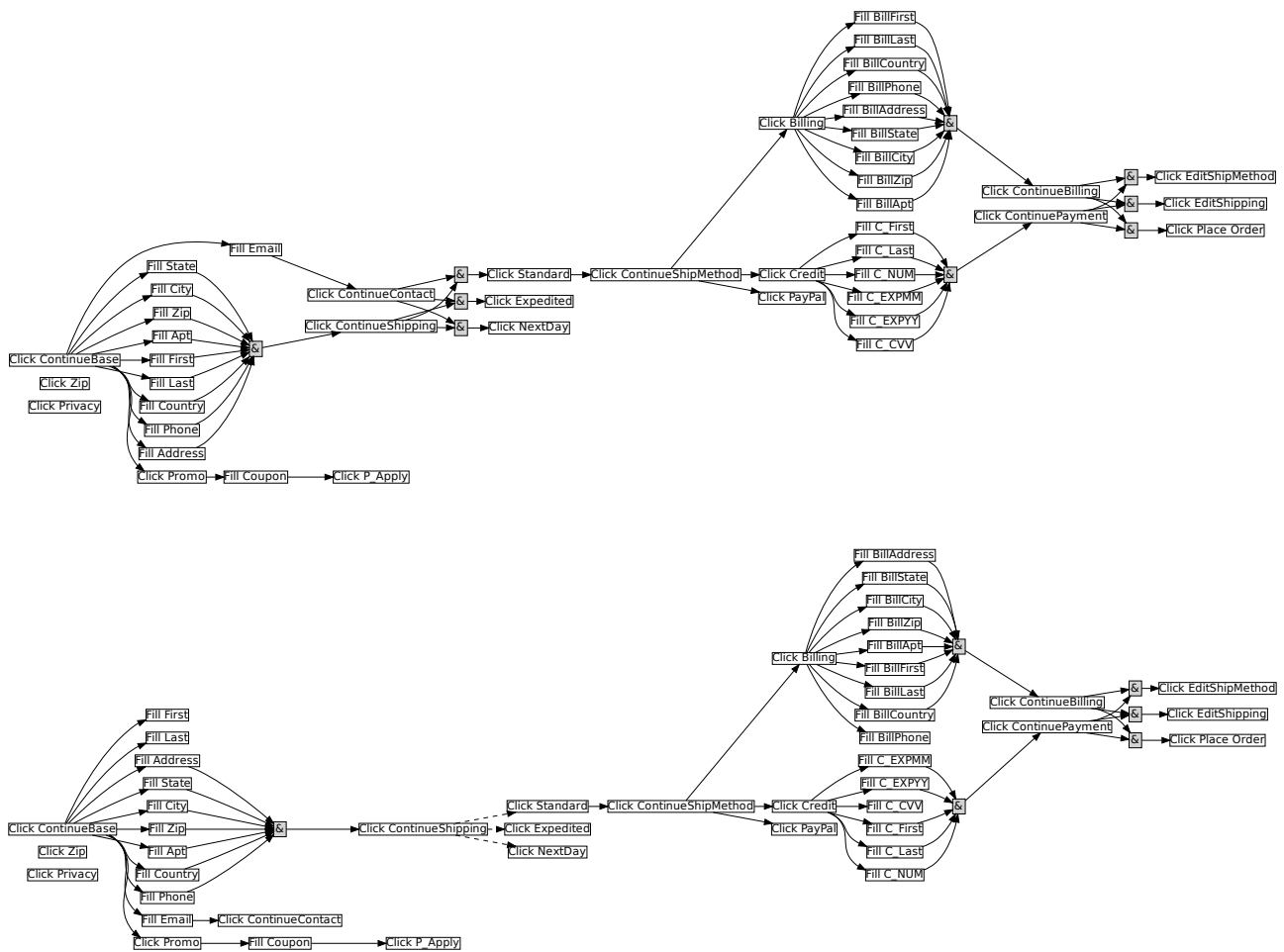


Figure 24: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Lenox** domain.

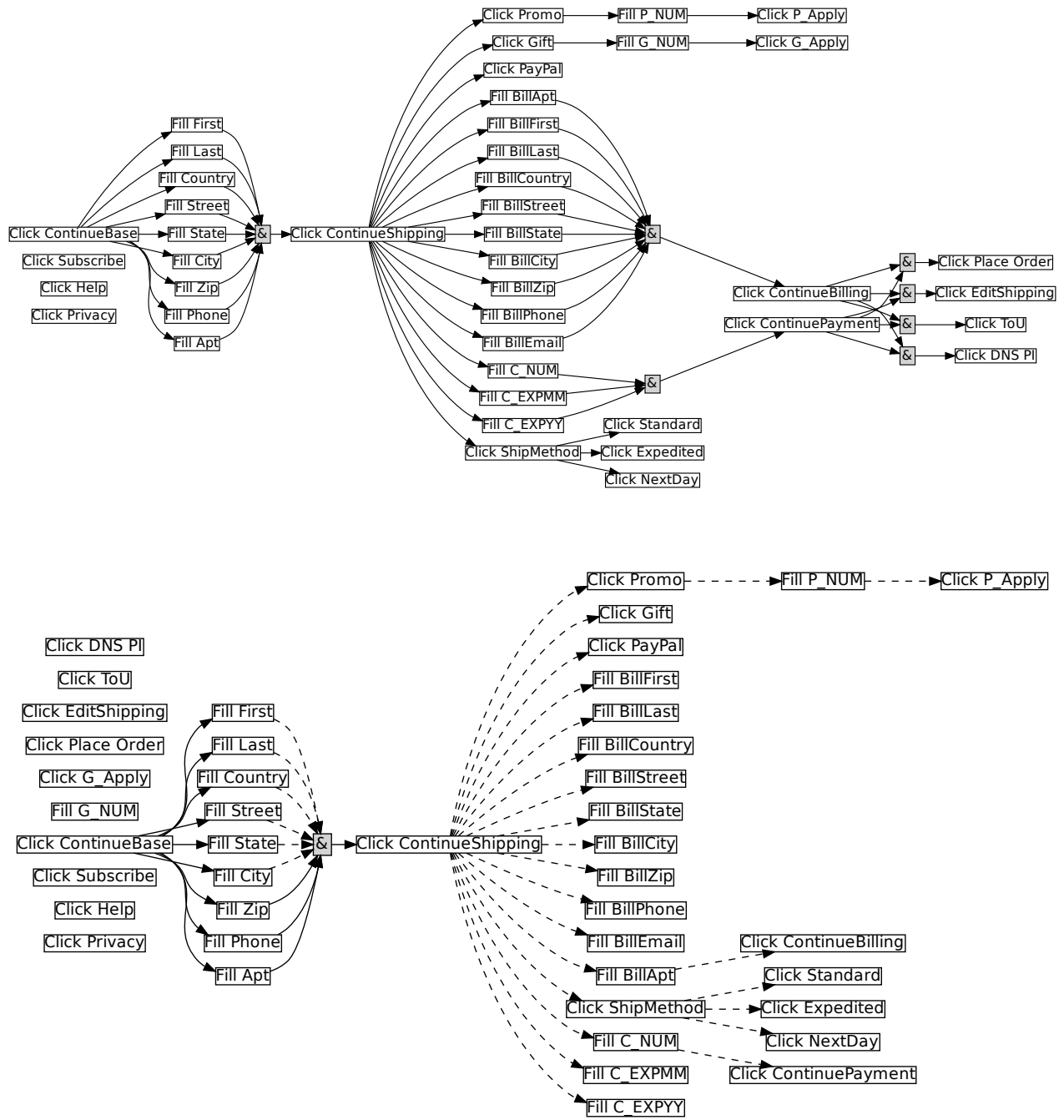


Figure 25: (Top) The ground-truth and (Bottom) the inferred subtask graphs of Omahasteaks domain.



Figure 26: (Top) The ground-truth and (Bottom) the inferred subtask graphs of Sephora domain.

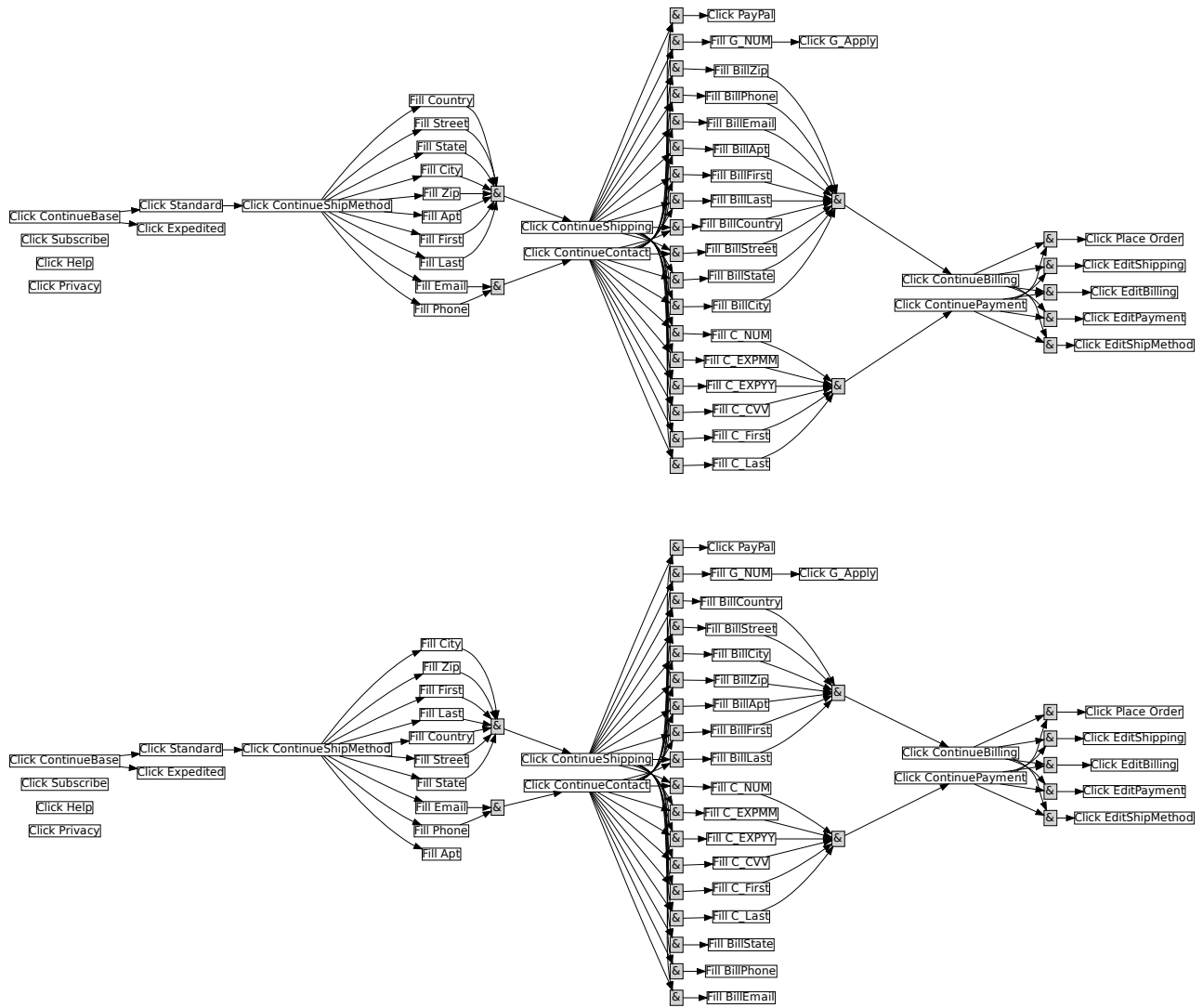


Figure 27: (Top) The ground-truth and (Bottom) the inferred subtask graphs of Swarovski domain.

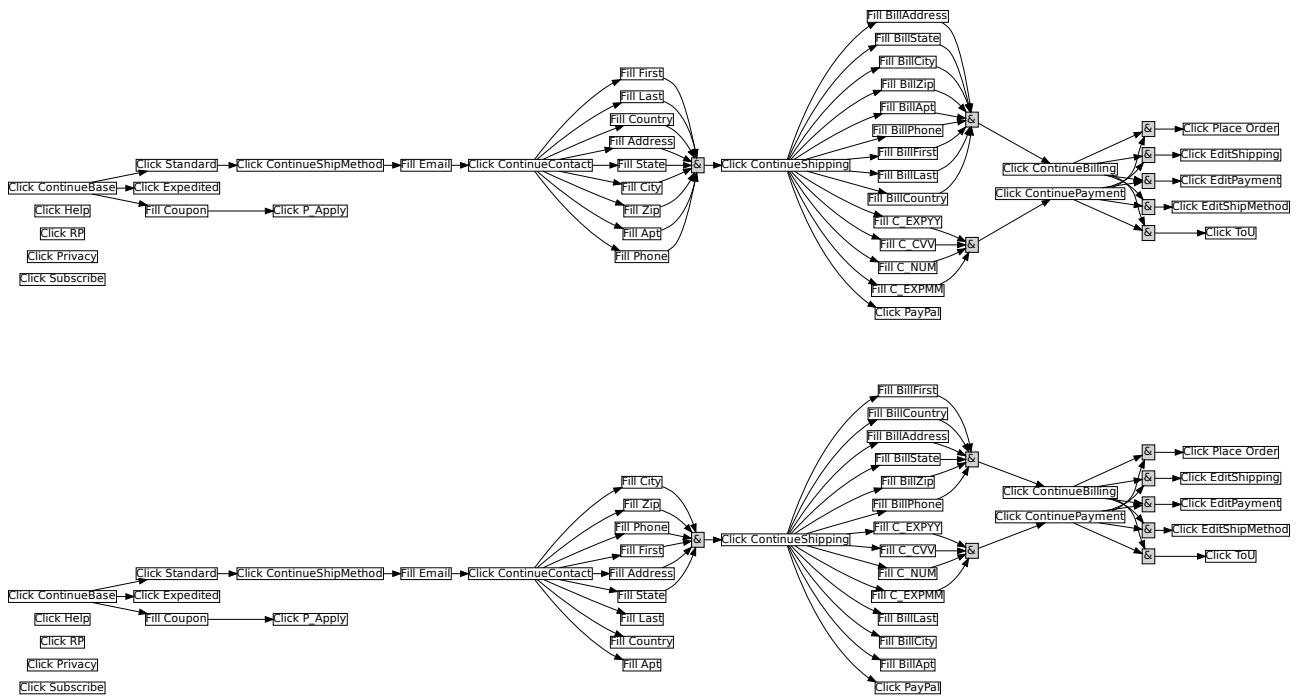


Figure 28: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Thriftbooks** domain.

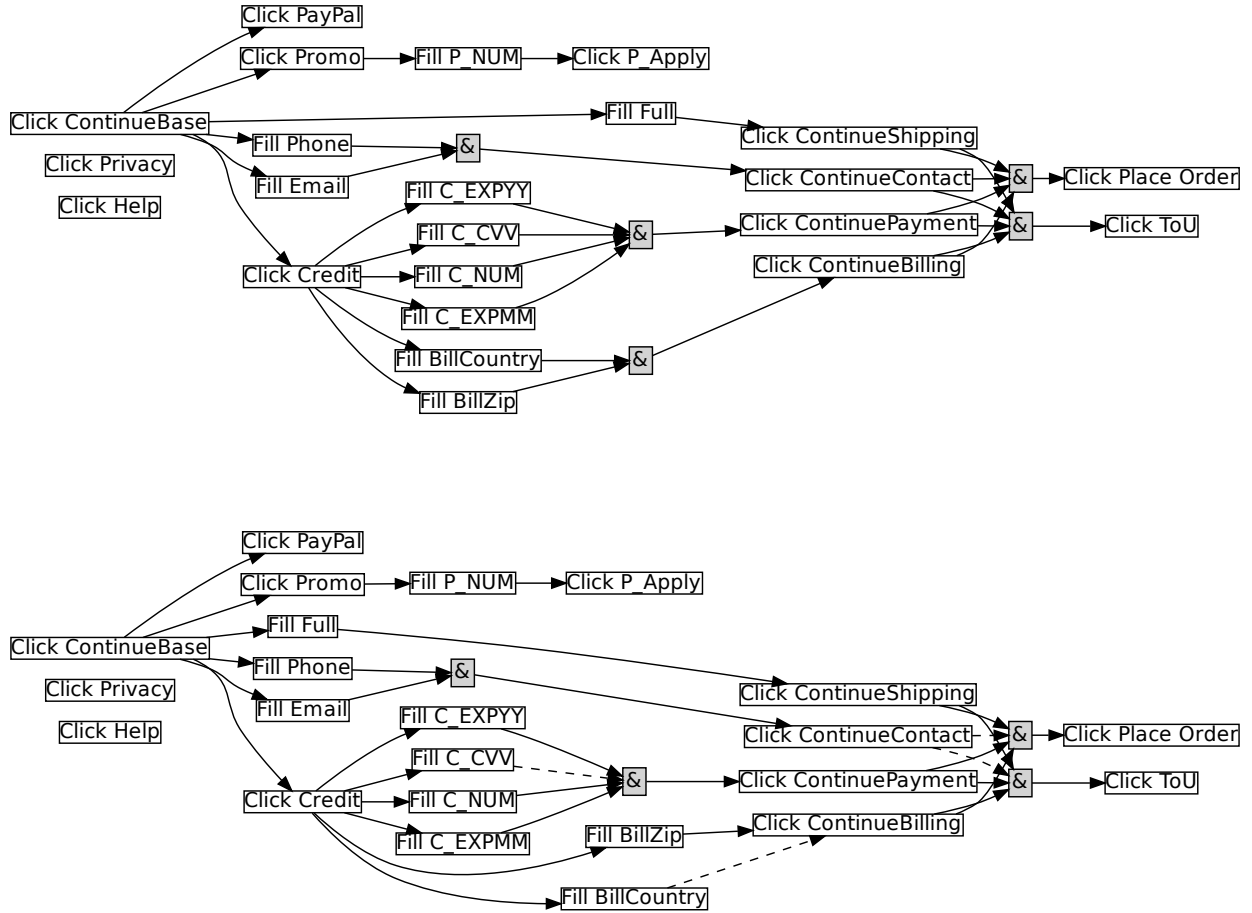


Figure 29: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Todaytix** domain.

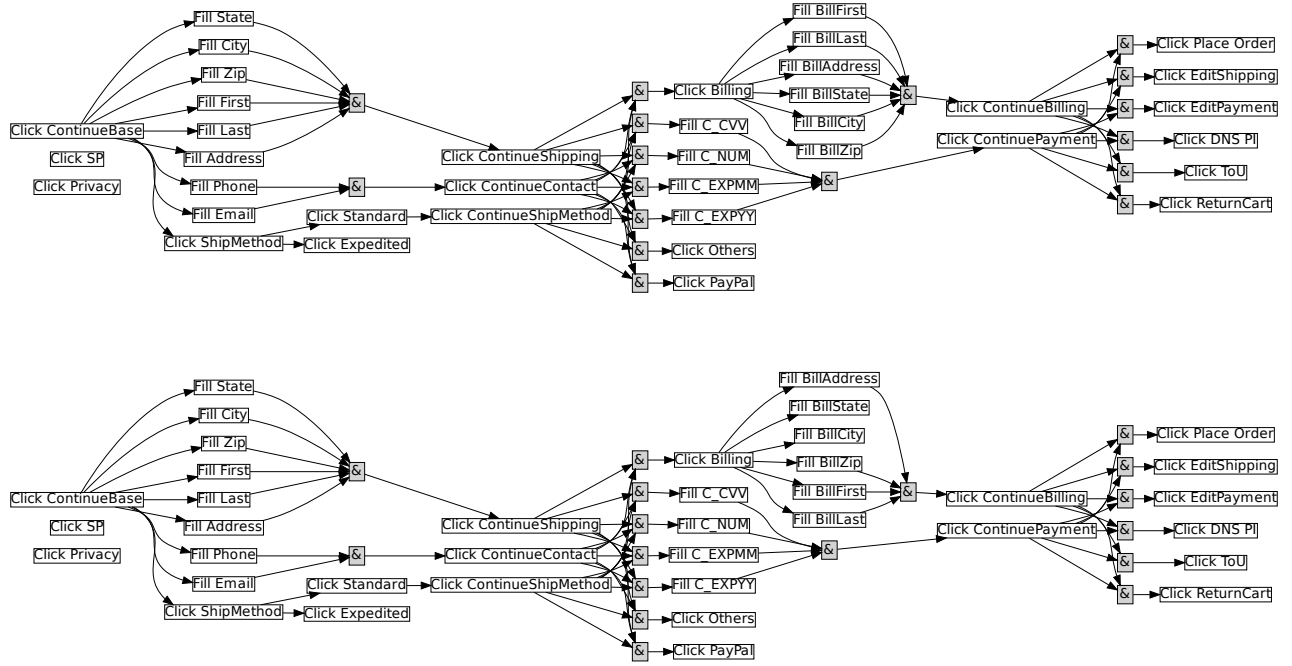


Figure 30: (Top) The ground-truth and (Bottom) the inferred subtask graphs of **Walgreens** domain.