

# Abstraction Refinement-Guided Program Synthesis for Robot Learning from Demonstrations

GUOFENG CUI, Rutgers University, USA

YUNING WANG, Rutgers University, USA

WENSEN MAO, Rutgers University, USA

YUANLIN DUAN, Rutgers University, USA

HE ZHU, Rutgers University, USA

Over the past decade, deep reinforcement learning (RL) techniques have significantly advanced robotic systems. However, due to the complex architectures of neural network models, ensuring their trustworthiness is a considerable challenge. *Programmatic reinforcement learning* has surfaced as a promising approach. Nonetheless, synthesizing robot-control programs remains challenging. Existing methods rely on domain-specific languages (DSLs) populated with user-defined state abstraction predicates and a library of low-level controllers as abstract actions to boot synthesis, which is impractical in unknown environments that lack such predefined components. To address this limitation, we introduce RoboScribe, a novel abstraction refinement-guided program synthesis framework that automatically derives robot state and action abstractions from raw, unsegmented task demonstrations in high-dimensional, continuous spaces. It iteratively enriches and refines an initially coarse abstraction until it generates a task-solving program over the abstracted robot environment. RoboScribe is effective in synthesizing iterative programs by inferring recurring subroutines directly from the robot's raw, continuous state and action spaces, without needing predefined abstractions. Experimental results show that RoboScribe programs inductively generalize to long-horizon robot tasks involving arbitrary numbers of objects, outperforming baseline methods in terms of both interpretability and efficiency.

CCS Concepts: • **Software and its engineering** → **Automatic programming**.

Additional Key Words and Phrases: Program Synthesis, Abstraction Refinement, Reinforcement Learning

## ACM Reference Format:

Guofeng Cui, Yuning Wang, Wensen Mao, Yuanlin Duan, and He Zhu. 2025. Abstraction Refinement-Guided Program Synthesis for Robot Learning from Demonstrations. *Proc. ACM Program. Lang.* 9, OOPSLA2, Article 292 (October 2025), 29 pages. <https://doi.org/10.1145/3763070>

## 1 Introduction

In the robotics domain, deep reinforcement learning (RL) techniques have shown promise in developing intelligent agents for robot control, offering robust alternatives to analytical models in adaptive control systems. However, ensuring the trustworthiness of deep RL systems is challenging due to the intricate nature of neural network structures. To address this, *programmatic reinforcement learning* has emerged, focusing

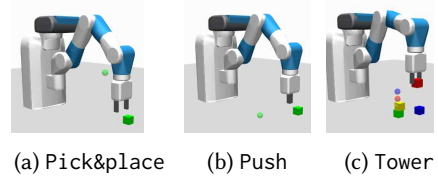


Fig. 1. Robot environments for various tasks.

Authors' Contact Information: [Guofeng Cui, gc669@cs.rutgers.edu](mailto:Guofeng Cui, gc669@cs.rutgers.edu), Rutgers University, USA; [Yuning Wang, yw895@rutgers.edu](mailto:Yuning Wang, yw895@rutgers.edu), Rutgers University, USA; [Wensen Mao, wm300@rutgers.edu](mailto:Wensen Mao, wm300@rutgers.edu), Rutgers University, USA; [Yuanlin Duan, yuanlin.duan@rutgers.edu](mailto:Yuanlin Duan, yuanlin.duan@rutgers.edu), Rutgers University, USA; [He Zhu, hz375@cs.rutgers.edu](mailto:He Zhu, hz375@cs.rutgers.edu), Rutgers University, USA.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2025 Copyright held by the owner/author(s).

ACM 2475-1421/2025/10-ART292

<https://doi.org/10.1145/3763070>

on synthesizing *domain-specific programs* as RL model representations to enhance interpretability [6, 26, 28, 47, 53, 60, 62, 63, 65, 66]. For example, a recent work, ReGuS (reward-guided synthesis) [14], generates robot-control programs directly from reward signals and demonstrates that synthesizing programs with rich control-flow constructs (e.g., loops) can effectively tackle long-horizon and sparse-reward tasks that often confound deep RL techniques.

**State and Action Abstraction.** Prior work scaling program synthesis to high-dimensional robotics environments with continuous states and actions often introduces a DSL with *state abstraction predicates* and *abstract control actions*. To illustrate the need for abstraction, consider the Pick&Place robot task shown in Fig. 1a, where a manipulator must pick up a block (green) from a table and place it in a target position (sphere) in mid-air. The robot used is a 7-DoF Fetch Mobile Manipulator with a two-fingered parallel gripper. The robot's state space includes kinematic information of the block and the end effector including the Cartesian coordinates of the desired final block position. Its action space represents the Cartesian displacement of the end effector to set to and the positional displacement of each finger of the gripper. Directly generating programs in a high-dimensional continuous space poses significant challenges to synthesis techniques.

State abstraction predicates [2, 19] construct a high-level representation of the robot's environment based on observed sensor data. This high-level representation can then be reasoned about using standard language constructs, such as loops and conditionals, to trigger suitable actions from a current state. Abstract actions are applications of low-level robot skills. These skills can be derived from either robot APIs or pre-trained neural network controllers. Conceptually, skills are modular and reusable, and can be likened to building blocks or subroutines that contribute to the overall control policy of the robot in diverse contexts. For example, for Pick&Place, ReGuS [14] involves state abstraction predicates Near (indicates if the gripper is close to the block  $b$ ), Holding (indicates if the gripper is holding  $b$ ), Above (indicates if the gripper is above  $b$ ), and At (indicates if  $b$  is in the goal region  $g \in \mathbb{R}^3$ ) into the DSL used to search a program to address this task. The DSL also includes a set of abstract actions to operate the robot in the continuous environment for openGripp(), closeGripp(), moveUp(), moveDown(), move( $g$ ) that moves the gripper to a goal region  $g$ . ReGuS synthesizes a program, as shown in Fig. 2, that guides the robotic gripper to move above the block  $b$  in the environment, lower to grab it, and then transport the block to the designated goal area  $g$ .

**Challenges Faced by Existing Approaches.** Several challenges remain in advancing robot-control program synthesis in real-world scenarios.

- (i) A crucial bottleneck in existing synthesis techniques is the reliance on manually designed state abstraction predicates and action abstraction in a DSL to bootstrap synthesis. High-quality abstractions often require significant human effort and domain knowledge to customize effectively. For example, the state and action abstractions for the Pick&place task described above cannot be applied to a robot with a push-only gripper as visualized in the Push task in Fig. 1b. This type of gripper is designed to apply force to reposition objects on a surface, but it cannot grasp or lift them. Automatically learning state and action abstractions has been a key area of research in task and motion planning for robot control [7, 23, 29, 43]. Existing techniques often have significant limitations, as they either learn predicates from demonstrations while assuming that low-level controllers are already available [15, 36, 55], or learn low-level controllers from demonstrations while assuming that the necessary predicates are predefined [1, 17, 54]. Simultaneous discovery of both state abstraction predicates and abstract actions remains a significant challenge [34, 54, 55].

```

while(not At( $b$ ,  $g$ )):
    if (Above( $b$ )):
        if (Near( $b$ )):
            if (Holding( $b$ )):
                move( $g$ )
            else:
                closeGripp()
        else:
            openGripp()
            if(not Near( $b$ )):
                moveDown()
    else:
        move( $b$ )

```

Fig. 2. Synthesized program for pick&place by ReGuS [14].

- (ii) Existing approaches encounter difficulties when synthesizing programs that generalize to tasks involving repeated subroutines or arbitrary numbers of objects, such as directing a robot arm to stack scattered objects into a pyramid. These tasks often require complex control-flow structures like state-conditioned loops. When a domain-specific language (DSL) with predefined state abstraction and low-level controllers is available, as in existing work, this complexity is reduced. For example, PROLEX [44] learns robot-control programs from task demonstrations, but assumes demonstrations are represented as sequences of calls to user-defined low-level controllers rather than continuous, raw actions in the robot's operational space. This representation enables it to deduce when and how these controllers are employed repeatedly or conditionally. Consequently, it can extract high-level control structures, such as loops and conditional statements, from observed tasks. Synthesizing loop programs for robot control in high-dimensional, continuous state and action spaces without predefined abstractions is still a major hurdle.

**This Paper.** To make program synthesis feasible for high-dimensional, continuous state and action spaces in robotic environments, our main idea is to develop **abstraction refinement** techniques that automatically generates appropriate state and action abstractions as part of the synthesis process. To this end, we develop **RoboScribe**, a novel *abstraction refinement-guided program synthesis* framework. RoboScribe is visualized in Fig. 3 and is based on the following two key ideas:

(1) **Comparative Abstraction Refinement:** To overcome challenge (i), RoboScribe iteratively refines an initially coarse abstraction of the robot environment until a valid program capable of solving the task is synthesized from the abstract environment. The initial coarse abstraction differentiates only between states that have met the (unknown) goal condition  $\psi_R$  in the environment and those that have not:

$$\text{True} \leadsto \psi_R$$

under the assumption that a single low-level controller as the target program can be learned to transition all possible initial states (True) to the goal condition ( $\psi_R$ ). In the context of pick&place, the initial abstraction only identifies goal states where the block is successfully placed at the target position. If this assumption fails, RoboScribe incrementally learns state abstraction predicates that capture critical intermediate states toward task completion. For example, in the pick&place task, the robot must move its gripper close to the block, grasp it, and lift it to a designated position. By learning state abstraction predicates that identify such key subgoal states and abstract actions grounded as low-level RL controllers that transition the agent across these subgoals, RoboScribe accurately captures the task's hierarchical structure. We posit that state abstraction predicates can be learned if demonstrations showing task completion are provided to the synthesizer. Systematically comparing states within these demonstrations to those observed during failed robot

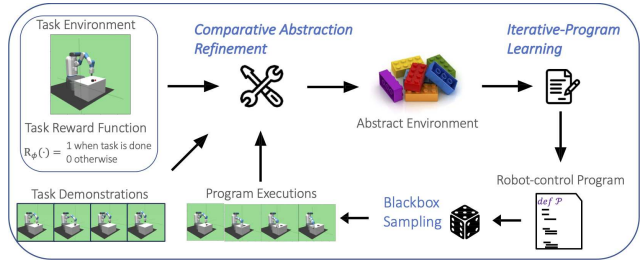


Fig. 3. Overview of the RoboScribe framework.

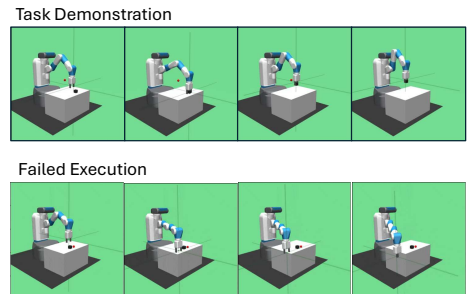


Fig. 4. Zoomed-in view of the task demonstration and program execution in Fig. 3.

behaviors which are learned in the real environment allows RoboScribe to discern necessary intermediate states in the demonstrations pivotal for the task's success. For pick&place, a low-level controller trained to fulfill  $\psi_R$  may struggle with grasping the block first, as exemplified in the program execution behavior shown in Fig. 4. Leveraging insights gained from successful demonstrations, RoboScribe refines the abstraction with a predicate  $\varphi$  that captures the states where the robot's gripper is positioned close to the block:

$$\text{True} \rightsquigarrow \varphi \rightsquigarrow \psi_R$$

The refined abstraction adds a crucial step that bridges the initial and goal states, breaking down the complex task into manageable subtasks. The intermediate subgoal  $\varphi$  supplies reward signals for efficient RL training of the subtask controllers. This process is recursively conducted, enabling the task to be eventually solved through a series of progressively refined subtasks that lead to  $\psi_R$ . **(2) Iterative-Program Learning:** To address challenge (ii), RoboScribe leverages state abstraction predicates learned *on-the-fly* to identify repeating subroutines within demonstrations. For instance, consider the Tower task of manipulating a robot arm to stack blocks into a tower visualized in Fig. 1c. Key state abstraction predicates for subgoals like grasping a block and lifting it to a specific position, learned for handling one block, recur throughout the demonstration for handling other blocks. RoboScribe exploits repeated subgoal predicates in demonstrations to construct the *loop body* that guides the agent through transitions between the subgoals via abstract actions in the form of low-level controllers. However, the order in which recurring objects are addressed across loop iterations remains unspecified. This ordering is crucial in tasks like Tower, where placing a block in the goal position before positioning the underlying blocks leads to failed outcomes. RoboScribe synthesizes a predicate that identifies the correct object order by analyzing the rationale behind the demonstrations, explaining why certain objects are handled before others. In a nutshell, to synthesizing programs that generalize to tasks involving arbitrary numbers of objects, RoboScribe first constructs a loop program skeleton, then fills in an object ordering predicate as the loop guard to determine the sequence for handling objects, and finally optimizes action abstractions to ensure robust generalization across all loop iterations.

**Evaluation.** We implemented RoboScribe and evaluated it using a benchmark suite of complex robot object manipulation tasks. Experimental results demonstrate that RoboScribe programs can inductively generalize to long-horizon tasks involving arbitrary numbers of objects, outperforming baseline methods in both interpretability and efficiency. For example, the programmatic agent synthesized by RoboScribe can efficiently use a robot arm to stack multiple blocks on a cluttered tabletop into a tower, a task known for its complexity in RL due to the need to handle long-horizon planning and precise manipulation [32, 41]. The agent can generalize to unseen configurations, such as placing blocks into multiple towers with zero-shot success.

**Contributions.** To summarize, this paper makes the following key contributions:

- We propose RoboScribe, a novel abstraction refinement technique that automatically derives robot state and action abstractions from raw, unsegmented task demonstrations to enable robot-control program synthesis in high-dimensional, continuous spaces.
- We develop an effective loop program synthesis algorithm that scales RoboScribe to long-horizon tasks involving unbounded environment objects. The algorithm excels in inferring repeating subroutines directly from demonstrations in the robot's raw, continuous state and action spaces.
- We evaluate RoboScribe in complex robot manipulation tasks, highlighting its effectiveness in learning and generalizing control strategies.

## 2 Overview

In this section, we motivate the problem and provide an overview of our approach.

Type  $\tau ::= \{\text{EE, block, handle, mug, faucet, } \dots\}$       Variable  $V ::= \mu, v$   
 Expression  $\alpha ::= \mu \mid v \mid g(v) \mid \alpha \downarrow_{x,y,z} \mid \alpha \downarrow_{x,y} \mid \alpha \downarrow_x \mid \alpha \downarrow_y \mid \alpha \downarrow_z \mid \alpha - \alpha \mid \|\alpha\| \mid \arctan2(\alpha)$   
 Predicate  $P ::= \alpha < \alpha \mid \alpha < \phi \mid \alpha > \phi \mid \alpha \neq \text{null} \mid \text{true} \mid \neg P \mid P \wedge P \mid P \vee P \mid \exists v. P$   
 Action  $c ::= \pi_{\theta_1}(\mu, \{v\}, g(\{v'\})) \mid \pi_{\theta_2}(\mu, \{v\}, g(\{v'\})) \mid \dots$        $\pi_{\theta_i} \in \Pi_{\text{NN}}$   
 Statement  $S ::= \text{while } (P) \{S\} \mid \text{if } (P) S_1 \text{ else } S_2 \mid S_1; S_2 \mid v := \text{get}(\lambda v : \tau. P) \mid c$   
 Program  $::= \text{def } \mathcal{P} (\mu : \text{EE}) : S$

Fig. 5. The Context-free grammar for the RoboScribe DSL  $\mathcal{L}$ , where EE refers to the robot's end effector.

## 2.1 Key Assumptions

**Object-centric Views.** Following common practice in robotics [27, 33, 52, 57, 58, 68], RoboScribe assumes that the robot receives an object-centric view  $\{\mu, e_1, \dots, e_N\}$  at each timestep (defined in Sec. 3). This view segments the world into discrete objects and classifies them into categories based on sensor data, where  $\mu$  represents the robot's end effector and  $e_i$  denotes an entity in the environment. Each object  $\mu$  or  $e_i$  is associated with attributes such as its class and 3D position. Such views can be constructed using e.g. object detection [37, 38, 50] or discovery [18, 35] methods.

**Goal-Directed Robot Tasks.** We consider goal-directed robotic tasks where a goal function  $g$  maps entities  $e_i$  to their target regions. For each control task,  $g$  may be randomly generated. For example, in the Pick&Place task (Fig. 1a), the goal region (green sphere) can be randomly placed within the robot's workspace. The robot's objective is defined by a predicate  $\phi_R$ , which is true when a designated subset of entities has reached their final target poses and false otherwise—e.g., an entity is considered to have reached its goal if its distance to the target is below a threshold. In RoboScribe, the agent does not require the analytical form of  $\phi_R$  but can query it to verify whether a state  $s$  satisfies  $\phi_R$ . We assume that entities belonging to the same type can be handled uniformly by a shared manipulation strategy for goal reaching.

**Demonstrations.** We assume that we can utilize supervision of a limited amount of task demonstrations for robot-control program synthesis. RoboScribe assumes *unsegmented* task demonstrations as sequences of states  $(s_0, s_1, \dots, s_H)$  where  $s_0$  is an initial state and each state  $s_i$  at timestep  $i$  presents an object-centric view of the system (i.e., a collection of objects in the scene and their attributes). For each demonstration, we assume that  $\psi_R(s_H)$  holds. We further assume that demonstrations reflect a consistent underlying intent, meaning that they all correspond to the same conceptual goal. Although the specific goal regions (e.g., target positions) may differ across demonstrations, they must consistently apply to the same subset of entities. For example, if the task is to place mugs in their respective goal regions, then all demonstrations should specify goal regions for mugs—not some for mugs and others for different objects like drawers. If demonstrations are ambiguous or inconsistent—i.e., if the goal regions involve different subsets of entities across demonstrations—RoboScribe may fail to synthesize a coherent control program, as the goal condition is not well defined<sup>1</sup>.

## 2.2 Program and Domain-Specific Language

RoboScribe synthesizes robot-control programs using a generic DSL  $\mathcal{L}$  in Fig. 5. In a program, the variable  $\mu$  binds to the robot's end effector and a variable  $v$  binds to an object  $e_i$  in the robot's object-centric view, with its type  $\tau$  determined by the object it references. Each object has attributes, such as its 3D pose estimated from sensors, denoted  $v \downarrow_{x,y,z}$ , where  $\downarrow$  extracts attribute values from

<sup>1</sup>In practice, RoboScribe can be augmented with a simple form of automatic goal inference: entities that remain unmoved throughout a demonstration can be interpreted as having their initial positions serve as implicit goal regions. This enables RoboScribe to function even when goal specifications are only partially specified.



the object referenced by  $v$ . In Fig. 5, we enumerate attributes related to position in  $\mathbb{R}^3$ , but the DSL can be extended to include full attributes, including orientation, such as quaternions. As discussed in Sec. 2.1, leveraging the goal function  $g$ ,  $g(v)$  represents the goal region of the object referenced by  $v$ , and  $\|\cdot\|$  represents the Euclidean norm.

**State abstraction predicates**  $P$  in  $\mathcal{L}$  enable the robot to locate relevant objects, constructing a higher-level representation of its environment. The parameter  $\phi$  in our predicates is a constant (vector) that must be learned. Our DSL supports existential quantifiers  $\exists v$  in predicates  $P$  to identify objects meeting specific criteria, such as blocks below a certain height. Predicates can be used to define spatial relationships (e.g., using  $\|\cdot\|$ ) and physical orientations (e.g., using  $\arctan2$ ) among environment objects, including the robot itself. For instance,  $\arctan2$  can be used to encode alignment constraints between the end effector, a block, and its goal region for direct pushing.

**Abstract actions**  $c$  represent low-level controllers that encode the robot's capabilities in continuous environments. Each abstract action  $c \in \Pi_{\text{NN}}$  is a deep neural network policy  $\pi_\theta$  with trainable parameters  $\theta$ . The policy  $\pi_\theta$  is goal-conditioned. It takes as input the attribute values of the end effector  $\mu$  and some objects referenced by  $\{v\}$ , as well as the goal regions of some objects referenced by  $\{v'\}$ , producing a control action suitable for execution in the raw environment.

Our DSL  $\mathcal{L}$  also includes *state-conditioned loops* and *conditional statements*. The *assignment* statement  $v := \text{get}(\lambda v : \tau. P)$  binds variable  $v$  to an object  $e_i$  of type  $\tau$  in the robot's object-centric view such that  $P[e_i/v]$  holds. A RoboScribe program  $\mathcal{P}$  is a function that takes a binding to the robot's end effector as input and executes a defined statement as its body.

It is important to highlight that the RoboScribe DSL  $\mathcal{L}$  does *not* include predefined low-level controllers. **The set of low-level controllers  $\Pi_{\text{NN}}$  starts out empty.** The synthesis process inherently involves *learning appropriate state abstraction predicates  $P$  and constructing abstract control actions  $\Pi_{\text{NN}}$*  as a fundamental component.

### 2.3 Demonstration-Directed Robot Environment Abstraction Refinement

**Key Insight.** Our key idea is to systematically compare states within successful task demonstrations with those from failed robot behaviors learned in the real environment, aiming to identify key subgoal states that are essential for task success.

**Comparative Abstraction Refinement.** We define a robot task  $\mathcal{T}$  as  $\mathcal{T} : \text{True} \rightsquigarrow \psi_R$  with the expected behavior of directing the robot to transition from arbitrary initial states True (underlying some unknown initial environment state distribution) to states that satisfy the goal condition  $\psi_R$ . The initial abstraction  $\text{True} \rightsquigarrow \psi_R$  is coarse, distinguishing only between successful states that meet the goal condition  $\psi_R$  and those that have not yet reached it, and assumes that a single low-level controller  $\pi_{\psi_R}$  can fully solve the task. RoboScribe attempts to learn  $\pi_{\psi_R}$  using an off-the-shelf deep RL algorithm driven by a task reward function that assigns a reward of 1.0 to any state  $s$  where  $\psi_R(s)$  is true and 0 otherwise. In this sparse reward setting, trajectories induced by  $\pi_{\psi_R}$  often fail to encounter any positive feedback, resulting in learning failure. RoboScribe compares the successful behavior given in a set of demonstrations  $\mathcal{D}$  with that of the learned controller to identify key differences that are essential for enabling task success.

For example, in the pick&place demonstration shown in Fig. 6, The end effector (referred to as the robot  $\mu$  for simplicity) first holds the block and then places it on the target. However, the learned  $\pi_{\psi_R}$  struggles with grasping the block due to the absence of an explicit learning signal for this action. We extract states preceding successful task completion

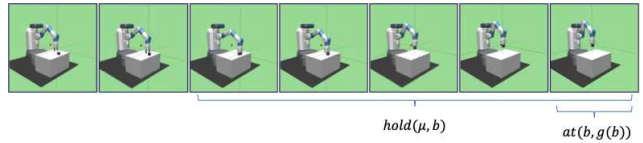


Fig. 6. Pick&place Demonstration.

from demonstrations, forming a set  $P_s = \{s_t \mid s_t \not\models \psi_R \wedge s_{t+1} \models \psi_R\}$ , and compare them with states  $N_s$  collected from the learned controller's trajectories. RoboScribe synthesizes a state abstraction predicate to distinguish  $P_s$  and  $N_s$  to learn what prerequisite conditions are necessary for reaching  $\psi_R$ . We learn state abstraction predicates  $P$  derivable from the grammar in our DSL  $\mathcal{L}$  in Fig. 5 based on Decision Tree (DT) learning. A DT is a binary tree that represents a Boolean formula. Each leaf of the tree is labeled either positive or negative for a subset of the samples in  $P_s \cup N_s$ . Each inner node is labeled by a decision of the form  $\alpha \leq \phi$  where  $\alpha$  is a feature and  $\phi$  is a (learned) threshold. In our context,  $\alpha$  is an expression derivable from the production rules for  $\alpha$  in  $\mathcal{L}$ . We formalize the learning algorithm in Sec. 4.2.1. In this example, RoboScribe may learn a predicate  $\text{hold}(\mu, b) \equiv \|\mu \downarrow_{x,y,z} - b \downarrow_{x,y,z}\| < \phi_1$  that defines states where a block  $b$  is being grasped by the robot gripper  $\mu$ . Here,  $b$  references to the block in the environment. Using this learned predicate, we refine the initial abstraction as:

$$\text{True} \rightsquigarrow \text{hold}(\mu, b) \rightsquigarrow \psi_R \quad (1)$$

which effectively decomposes the task into subtasks: (1)  $\mathcal{T}_1 : \text{True} \rightsquigarrow \text{hold}(\mu, b)$  for reaching states where the block is grasped, and (2)  $\mathcal{T}_2 : \text{hold}(\mu, b) \rightsquigarrow \psi_R$  for achieving the goal condition after grasping the block. This process is recursively conducted until the task can be solved through a series of progressively refined subtasks to  $\psi_R$ . For example, for  $\mathcal{T}_2$ , RoboScribe may further learn a predicate  $\text{at}(b, g(b)) \equiv \|b \downarrow_{x,y,z} - g(b) \downarrow_{x,y,z}\| < \phi_2$  capturing states where the block is at the goal region, along with two subtasks:  $\mathcal{T}_{2_1} : \text{hold}(\mu, b) \rightsquigarrow \text{at}(b, g(b))$  for moving the block towards its goal region, and  $\mathcal{T}_{2_2} : \text{at}(b, g(b)) \rightsquigarrow \psi_R$  for maintaining the block in the target position of the task, resulting in the following refined abstraction:

$$\text{True} \rightsquigarrow \text{hold}(\mu, b) \rightsquigarrow \text{at}(b, g(b)) \rightsquigarrow \psi_R \quad (2)$$

**Abstract Subtask Tree.** RoboScribe structures environment abstractions as abstract subtask trees, formalized in Sec. 4. Each tree node  $\varphi$  or  $\psi$  encodes a state abstraction predicate that defines a subgoal condition. Each (inverted) tree edge  $\varphi \rightsquigarrow \psi$  represents an abstract action, to be grounded as a low-level controller in the real environment, that fulfills the subtask of transitioning any state within  $\varphi$  to a subgoal state in  $\psi$ , guiding the agent toward the completion of its overall task goal condition  $\psi_R$  at the root of the tree. Tree representations effectively capture the structure of multi-goal tasks involving multiple objects, providing a clear framework for task decomposition and execution. To satisfy a subgoal  $\psi$ , the agent must complete all the subgoals of its predecessors  $\varphi$  such that  $\varphi \rightsquigarrow \psi$ . Each subgoal  $\varphi$  in the predecessors of  $\psi$  corresponds to the manipulation of a distinct object.

Consider the PlaceCubeDrawer task depicted in Fig. 7 left. The goal of the Sawyer robot in this task is to pick up a cube  $b$  from the desk of a cabinet and place it inside the drawer below. The abstract subtask tree for this task is shown in Fig. 7 (right). The agent must complete the following subtasks: first moving the gripper  $\mu$  near the cabinet door handle  $h$ , latching the handle, pulling it to align with the goal position  $g(b)$  for the cube  $b$ , and then holding  $b$  before placing it inside the drawer at  $g(b)$ .

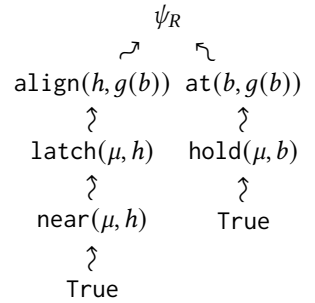
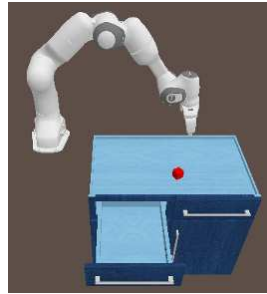


Fig. 7. Abstract Subtask Tree for PlaceCubeDrawer.

**Program Learning.** From an abstract environment defined by an abstract subtask tree  $T$ , RoboScribe synthesizes a robot-control program  $\mathcal{P}$ . First, we note that  $\mathcal{P}$  can be derived by recursively traversing  $T$  and chaining the subtasks. Each abstract action is grounded as a neural network

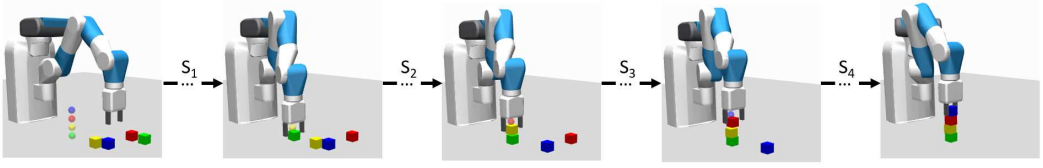


Fig. 9. Demonstration and partition for the Tower task.

controller  $\pi$  that uses the attributes of the involved objects as input to generate low-level robot actions. In  $\mathcal{P}$ , each controller  $\pi$  runs until its corresponding subgoal condition is met. We depict the synthesized program for pick&place in Fig. 8 where the shorthand notations in the program represent loops:

$$\pi(\mu, \{v\}, \{g(v)\})[\varphi] \equiv \text{while not}(\varphi) \{ \pi(\mu, \{v\}, \{g(v)\}) \}$$

Second, to learn each low-level controller  $\pi_\varphi$  in  $\mathcal{P}$  for reaching states that satisfy its subgoal condition  $\varphi$ , we iteratively execute  $\mathcal{P}$  and store trajectories from  $\pi_\varphi$  in a replay buffer  $B_\varphi$ . During gradient updates, RoboScribe trains  $\pi_\varphi$  by sampling from  $B_\varphi$  and optimizing it with an off-the-shelf off-policy RL algorithm. Particularly, we use learned state abstraction predicates to provide dense reward signals for training low-level controllers. Dense reward functions are shaped systematically from predicates combined in arbitrary Boolean forms. For example, in Fig. 6, to train a controller for the subtask  $\mathcal{T}_{21} : \text{hold}(g, b) \leadsto \text{at}(b, g(b))$ , a dense reward function can be derived from the norm learned for  $\text{at}(b, g(b))$ , encouraging the robot to move its end effector closer to  $g(b)$  with higher rewards for proximity.

```
def  $\mathcal{P}_{\text{Pick\&place}}(\mu : \text{EE})$ :
   $b := \text{get}(\lambda b : \text{block}, \text{true})$ ;
   $\pi_{\text{grasp}}(\mu, b) [\text{hold}(\mu, b)]$ ;
   $\pi_{\text{position}}(\mu, b, g(b)) [\text{at}(b, g(b))]$ 
```

Fig. 8. Synthesized program for pick&place by RoboScribe.

## 2.4 Synthesizing Iterative Robot-Control Programs

While the synthesis strategy described in Sec.2.3 is applicable to multi-object tasks, it does not generalize well when scaling up to handle long-horizon tasks with varying numbers of objects. The abstraction method lacks the flexibility to capture the relationships and dependencies between an indefinite number of objects, particularly when coordination and sequencing are required. For example, in the Tower task shown in Fig.9, the goal is to synthesize a program  $\mathcal{P}$  for a robot arm to stack a variable number of scattered blocks into a tower. Such tasks demand programs that can *iteratively* manage multiple instances of subtasks.

**Key Challenge.** Synthesizing iterative programs is challenging because it requires effective strategies for discovering repetitive subroutines and handling the complex dependencies between them. Unlike prior work (e.g. PROLEX [44] and Tabula [46]), RoboScribe does not assume predefined state and action abstraction and must be able to extract repetitive structures from demonstration trajectories within the robot’s high-dimensional, continuous state and action spaces.

**Key Insight.** Our key idea is to leverage abstraction predicates learned *on-the-fly* to discover repetitive subroutines in demonstrations. During comparative abstraction refinement, when classifying states from task demonstrations and robot trajectories to learn new state abstraction predicates, RoboScribe identifies opportunities to reuse previously discovered predicates and their corresponding low-level controllers, as a means to uncover *abstract* repetitions within demonstrations.

**Discovering Repetitive Subtasks.** As shown in the learned abstraction for the Pick&place task in Eq. 2, RoboScribe can identify a state abstraction predicate  $\text{at}(b, g(b))$  as a subgoal for placing one block  $b$  at its target and develop a routine of low-level controllers to achieve it. However, the



task remains incomplete as additional blocks must still be placed. Using a single low-level controller for this would be insufficient as illustrated in Sec. 2.3. RoboScribe identifies that the predicate  $at(b, g(b))$ , interpreted with  $b$  as implicitly existentially quantified, can be reused to distinguish demonstration states where blocks are near their targets from unsuccessful single-policy attempts to complete the full task. It refines the abstraction in Eq. 2 by defining a circular abstract subtask tree node to indicate this repetition:

$$\text{True} \rightsquigarrow \text{hold}(\mu, b) \rightsquigarrow \text{at}(b, g(b)) \rightsquigarrow \psi_R \quad (3)$$

This implies that the control strategy for achieving  $\exists b. at(b, g(b))$  for some block  $b$  can be iteratively applied to handle remaining objects of the same type in the environment that have yet to meet this subgoal condition.

**Learning Iterative Programs.** RoboScribe synthesizes a loop structure for each circular abstract subtask tree node. This leads to an iterative Tower program  $\mathcal{P}_{\text{Tower}}$  shown in Fig. 10, which intends to stack all the blocks on a table in a sequence, generated from the task abstraction in Eq. 3. However, determining the order in which blocks should be addressed within the sequence (among the loop iterations) remains unspecified. This ordering is particularly significant as placing a block in its goal position without first positioning the underlying blocks leads to failed outcomes. RoboScribe places a missing hole  $??_p$  in the loop condition designated to specify an effective handling sequence in Fig. 10.

**Loop Condition Synthesis.** Conceptually, we can enumerate candidates to fill in the missing predicate  $??_p$  based on the predicate production rules  $P$  defined in our DSL  $\mathcal{L}$  (Fig. 5) and execute  $\mathcal{P}_{\text{Tower}}$  in the real environment to empirically determine which predicate maximizes task performance, such as higher success rates. However, this approach is computationally prohibitive due to the extensive predicate search space of the DSL and the long-horizon nature of robot tasks that involve recurring objects. Additionally, training the low-level neural controllers in  $\mathcal{P}_{\text{Tower}}$  depends on executing the program to obtain training data, resulting in a mutual dependency problem where loop condition synthesis and controller learning are interdependent. Our strategy circumvents these limitations by inferring the ordering predicate  $??_p$  directly from demonstrations, avoiding the need for executing  $\mathcal{P}_{\text{Tower}}$  in the real environment. Specifically, RoboScribe synthesizes an ordering predicate for  $??_p$  that determines the correct sequence of object handling by analyzing the rationale in the demonstrations, explaining why certain objects are handled before others, such as why the red block is placed after the green and yellow blocks in Fig. 9. From demonstrations, RoboScribe learns such a predicate by enumerating predicates derivable from the production rules for  $P$  in the DSL  $\mathcal{L}$  (see Fig. 5). In this process, we augment  $P$  with learned state abstraction predicates from abstract subtask trees as these predicates provide additional task-relevant constraints. We defer the formalization of the synthesis algorithm to Sec. 4.2.1. For Tower, RoboScribe synthesizes the following predicate for an effective handling sequence :

$$\neg at(b, g(b)) \wedge \neg \exists b'. g(b') \downarrow_z < g(b) \downarrow_z \wedge \neg at(b', g(b')) \quad (4)$$

```
def  $\mathcal{P}_{\text{Tower}}$  ( $\mu : \text{EE}$ ):
  while (( $b := \text{get}(\lambda b : \text{block}. ??_p)$ )
     $\neq \text{null}$ ):
     $\pi_{\text{grasp}}(\mu, b)$  [ $\text{hold}(\mu, b)$ ]
     $\pi_{\text{position}}(b, g(b))$  [ $at(b, g(b))$ ]
```

Fig. 10. Iterative program  $\mathcal{P}_{\text{Tower}}$  for Tower with a missing hole for the loop condition.

```
def  $\mathcal{P}_{\text{Pick\&place}}$  ( $\mu : \text{EE}$ ):
  while (( $b := \text{get}(\lambda b : \text{block}. \neg at(b, g(b)) \wedge$ 
     $\neg \exists b'. g(b') \downarrow_z < g(b) \downarrow_z \wedge \neg at(b', g(b'))$ )
     $\neq \text{null}$ ):
     $\pi_{\text{grasp}}(\mu, b)$  [ $\text{hold}(\mu, b)$ ]
     $\pi_{\text{position}}(b, g(b))$  [ $at(b, g(b))$ ]
```

Fig. 11. Synthesized iterative program  $\mathcal{P}_{\text{Tower}}$  for Tower.

which specifies that any block with a lower goal position must be placed before the current block. The termination condition ensures that once all blocks are in their goal positions, there are no further blocks to handle. The full program synthesized is given in Fig. 11.

### 3 Problem Setup

We study a learning paradigm where the agent can interact with many entities (objects) in an environment. The task for the agent is specified in the form of goals for the entities. We formalize it using the Entity-Factored Markov Decision Process (EFMDP) [67].

Throughout the paper, we use  $\{v\}$  to denote a list. For a function  $f$ , we define element-wise application as  $f(\{v\}) = \{f(v_1), \dots, f(v_n)\}$ .

**Entity-Factored Markov Decision Process.** An EFMDP with  $N$  entities is described by the tuple:  $\mathcal{M} := \langle \Lambda, O = \{\mu, e_1, \dots, e_N\}, \mathcal{S}, \mathcal{G}, \mathcal{A}, \mathbb{P}, \eta \rangle$ . Here,  $\Lambda$  is a finite set of object types, e.g., cube and mug, and  $O$  is a finite set of objects, where  $\mu$  and  $\{e_1, \dots, e_N\}$  are the agent (robot) and the entities, respectively. Each entity in  $O$  has a type drawn from  $\Lambda$ . Each object in  $O$  has an associated set of attributes drawn from a finite set  $\mathcal{F} = \{f_1, f_2, \dots, f_M\}$ , for example, spatial coordinates  $\{x, y, z\}$  in the 3D space. A state  $s$  in the state space  $\mathcal{S}$  is a function  $s : O \rightarrow \mathbb{D}$  where  $\mathbb{D}$  is the space of object descriptors, formally defined as  $\mathbb{D} = (\mathcal{F} \rightarrow \mathbb{R})$ . This means that each object  $o \in O$  is mapped to a function that assigns a real value to each attribute. For an object  $o \in O$ ,  $s(o)$  retrieves the object descriptor of  $o$ , i.e.,  $s(o) : \mathcal{F} \rightarrow \mathbb{R}$ , and  $s(o, \{f\})$  extracts the real values of the attributes  $\{f\} \subseteq \mathcal{F}$ , i.e.,  $s(o, \{f\}) = (s(o))(\{f\})$ . We use  $\text{dom}(s)$  to retrieve the set of objects within a state  $s$ . We sometimes abuse notation for convenience to use  $s(o)$  to refer to the full set of attribute values for  $o$ , i.e., the image of  $s(o)$  under  $\mathcal{F}$ .

The **goal space** of an EFMDP  $\mathcal{M}$  is denoted as  $\mathcal{G}$ . A goal command  $g \in \mathcal{G}$  is a function (introduced in Sec. 2.1)  $g : O \rightarrow \mathbb{D}$  that defines the goal region for entities  $\{e_1, \dots, e_N\}$ . Typically,  $g(e_i)$  only maps a subset of  $e_i$ 's attributes to a real value, specifying its desired placement. For instance, in Fig. 1, the goal regions (spheres) indicate the target positions for each block within the 3D space.

In an EFMDP  $\mathcal{M}$ ,  $\mathcal{A}$  is the robot's action space. The **system dynamics** of  $\mathcal{M}$  is described by a probabilistic state transition function  $\mathbb{P}(s'|s, a)$  for  $s, s' \in \mathcal{S}$  and  $a \in \mathcal{A}$ , i.e., the robot's action can update the object states in its environment. The set of the initial states of an EFMDP is specified by  $\eta : \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$  (i.e.,  $\eta(s)$  is the probability density of the initial state being  $s$ ). A **trajectory** of an EFMDP  $\zeta \in Z$  is a sequence  $\zeta = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots$ , where  $s_i \in \mathcal{S}$  and  $a_i \in \mathcal{A}$ , where  $s_{i+1} \sim \mathbb{P}(\cdot | s_i, a_i)$ . EFMDPs can model several applications, including tabletop manipulation and scene reconfiguration. At the same time, the EFMDP contains more structure and symmetry compared to the standard MDP model, which can enable more efficient learning and better generalization [67].

**Task Specification.** We define predicates  $\psi$  used for robot task specifications of an EFMDP, as shown in Fig. 12, over the set of objects  $O = \{\mu, e_1, \dots, e_N\}$  within. The operator  $\downarrow$  extracts attribute values from  $o \in O$ . The semantics of the predicates  $\llbracket \psi \rrbracket$  are given in Fig. 13. Given a state  $s$ , we define  $\varphi(s)$  as  $\llbracket \varphi \rrbracket(s)$ , representing the truth value of the predicate  $\varphi$  in state  $s$ . We say that a trajectory  $\zeta = s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} \dots, s_H$  satisfies a task specification  $\mathcal{T} : \varphi \rightsquigarrow \psi$ , denoted as  $\zeta \models \mathcal{T} : \varphi \rightsquigarrow \psi$ , if  $\psi(s_H)$  holds when  $\varphi(s_0)$  holds. A task specification  $\mathcal{T} : \text{True} \rightsquigarrow \psi_R$  for an EFMDP  $\mathcal{M}$  defines the intended behavior in the MDP. Starting from any possible initial EFMDP state  $s_0 \sim \eta(\cdot)$ , the agent is expected to reach a state  $s$  that satisfies  $\psi_R(s)$ . For example, for the Tower task in Fig. 1c involving the end effector  $\mu$  and four blocks  $O = \{\mu, b_1, b_2, b_3, b_4\}$ , its specification can be defined as:

$$\mathcal{T}_{\text{Tower}} : \text{True} \rightsquigarrow \bigwedge_{i=1}^4 \|b_i \downarrow_{x,y,z} - g(b_i)\| < \phi \wedge b_4 \downarrow_z < \mu \downarrow_z$$

This specifies that all blocks must be placed in their goal regions, and the end effector must leave the top block. In this paper, we use  $\psi_R$  to denote the predicate encoding the task's final goal region, while

$$\alpha ::= \mu \mid e_i \mid g(e_i) \mid \alpha \downarrow_{x,y,z} \mid \alpha \downarrow_{x,y} \mid \alpha \downarrow_x \mid \alpha \downarrow_y \mid \alpha \downarrow_z \mid \alpha - \alpha \mid \|\alpha\| \mid \arctan2(\alpha)$$

$$\psi ::= \alpha < \alpha \mid \alpha < \phi \mid \alpha > \phi \mid \neg\psi \mid \psi \wedge \psi \mid \psi \vee \psi \mid \exists o. \psi$$

Fig. 12. Task Specifications of EFMDPs over objects  $\mathcal{O} = \{\mu, e_1, \dots, e_N\}$ .

Expressions	Predicates
$\llbracket \alpha \rrbracket(s) =$ $\begin{cases} s(\mu), & \text{if } \alpha = \mu \\ s(e_i), & \text{if } \alpha = e_i \\ g(e_i), & \text{if } \alpha = g(e_i) \\ s(\llbracket \alpha_1 \rrbracket(s), \{x, y, z\}), & \text{if } \alpha = \alpha_1 \downarrow_{x,y,z} \\ \llbracket \alpha_1 \rrbracket(s) - \llbracket \alpha_2 \rrbracket(s), & \text{if } \alpha = \alpha_1 - \alpha_2 \\ \ \llbracket \alpha_1 \rrbracket(s)\ , & \text{if } \alpha = \ \alpha_1\  \\ \arctan2(\llbracket \alpha_1 \rrbracket(s)), & \text{if } \alpha = \arctan2(\alpha_1) \end{cases}$	$\llbracket \psi \rrbracket(s) =$ $\begin{cases} \llbracket \alpha_1 \rrbracket(s) < \llbracket \alpha_2 \rrbracket(s), & \text{if } \psi = \alpha_1 < \alpha_2 \\ \llbracket \alpha \rrbracket(s) < \phi, & \text{if } \psi = \alpha < \phi \\ \llbracket \alpha \rrbracket(s) > \phi, & \text{if } \psi = \alpha > \phi \\ \neg \llbracket \psi_1 \rrbracket(s), & \text{if } \psi = \neg \psi_1 \\ \llbracket \psi_1 \rrbracket(s) \wedge \llbracket \psi_2 \rrbracket(s), & \text{if } \psi = \psi_1 \wedge \psi_2 \\ \llbracket \psi_1 \rrbracket(s) \vee \llbracket \psi_2 \rrbracket(s), & \text{if } \psi = \psi_1 \vee \psi_2 \\ \bigvee_i \llbracket \psi_1[o \mapsto e_i] \rrbracket(s), & \text{if } \psi = \exists o. \psi_1 \end{cases}$

Fig. 13. Semantics of task specifications of EFMDPs over objects  $\mathcal{O} = \{\mu, e_1, \dots, e_N\}$ .

$$\frac{\langle S_1, \rho, s \rangle \Downarrow (\rho', s') \quad \langle S_2, \rho', s' \rangle \Downarrow (\rho'', s'')}{\langle S_1; S_2, \rho, s \rangle \Downarrow (\rho'', s'')} \quad \frac{\langle P, \rho, s \rangle \Downarrow \text{true} \quad \langle S_1, \rho, s \rangle \Downarrow (\rho', s')}{\langle \text{if}(P) S_1 \text{ else } S_2, \rho, s \rangle \Downarrow (\rho', s')} \quad \frac{\langle P, \rho, s \rangle \Downarrow \text{false} \quad \langle S_2, \rho, s \rangle \Downarrow (\rho', s')}{\langle \text{if}(P) S_1 \text{ else } S_2, \rho, s \rangle \Downarrow (\rho', s')}$$

$$\frac{\langle P, \rho, s \rangle \Downarrow \text{true} \quad \langle S, \rho, s \rangle \Downarrow (\rho', s') \quad \langle \text{while}(P) \{S\}, \rho', s' \rangle \Downarrow (\rho'', s'')}{\langle \text{while}(P) \{S\}, \rho, s \rangle \Downarrow (\rho'', s'')} \quad \frac{\langle P, \rho, s \rangle \Downarrow \text{false} \quad \langle \text{while}(P) \{S\}, \rho, s \rangle \Downarrow (\rho, s)}{\langle \text{while}(P) \{S\}, \rho, s \rangle \Downarrow (\rho, s)}$$

$$\frac{o : \tau \in \text{dom}(s) \quad \langle P, \rho[v \mapsto o], s \rangle \Downarrow \text{true}}{\langle v := \text{get}(\lambda v : \tau. P), \rho, s \rangle \Downarrow (\rho[v \mapsto o], s)} \quad \frac{\forall o : \tau \in \text{dom}(s). \langle P, \rho[v \mapsto o], s \rangle \Downarrow \text{false}}{\langle v := \text{get}(\lambda v : \tau. P), \rho, s \rangle \Downarrow (\rho[v \mapsto \text{null}], s)}$$

$$\frac{a \sim \pi_\theta(s(\rho(\mu)), s(\rho(\{v\})), g(\rho(\{v'\}))) \quad s' \sim \mathbb{P}(\cdot \mid s, a)}{\langle \pi_\theta(\mu, \{v\}, g(\{v'\})), \rho, s \rangle \Downarrow (\rho, s')}$$

$$\frac{\langle \alpha, \rho, s \rangle \Downarrow u}{\langle \alpha < \phi, \rho, s \rangle \Downarrow u < \phi} \quad \frac{\langle \alpha, \rho, s \rangle \Downarrow u}{\langle \alpha \neq \text{null}, \rho, s \rangle \Downarrow u \neq \text{null}} \quad \frac{\langle P_1, \rho, s \rangle \Downarrow b_1 \quad \langle P_2, \rho, s \rangle \Downarrow b_2}{\langle P_1 \wedge P_2, \rho, s \rangle \Downarrow b_1 \wedge b_2} \quad \frac{\langle P, \rho, s \rangle \Downarrow b}{\langle \neg P, \rho, s \rangle \Downarrow \neg b}$$

$$\frac{o : \tau \in \text{dom}(s) \quad \langle P, \rho[v \mapsto o], s \rangle \Downarrow \text{true}}{\langle \exists v : \tau. P, \rho, s \rangle \Downarrow \text{true}} \quad \frac{\forall o : \tau \in \text{dom}(s). \langle P, \rho[v \mapsto o], s \rangle \Downarrow \text{false}}{\langle \exists v : \tau. P, \rho, s \rangle \Downarrow \text{false}} \quad \frac{\rho(v) \in \text{dom}(s)}{\langle v, \rho, s \rangle \Downarrow \rho(v)}$$

$$\frac{\rho(v) \in \text{dom}(s)}{\langle g(v), \rho, s \rangle \Downarrow g(\rho(v))} \quad \frac{\langle \alpha, \rho, s \rangle \Downarrow o \quad o \in \text{dom}(s)}{\langle \alpha \downarrow_{x,y,z}, \rho, s \rangle \Downarrow s(o, \{x, y, z\})} \quad \frac{\langle \alpha_1, \rho, s \rangle \Downarrow u_1 \quad \langle \alpha_2, \rho, s \rangle \Downarrow u_2}{\langle \alpha_1 - \alpha_2, \rho, s \rangle \Downarrow u_1 - u_2}$$

Fig. 14. The DSL  $\mathcal{L}$  operational semantics in RoboScribe.

$\psi$  and  $\phi$  typically represent intermediate subgoal conditions inferred by our algorithm. Predicates  $\psi$  defined over EFMDP objects  $\mathcal{O}$  align with predicates  $P$  in the DSL  $\mathcal{L}$ , which are defined over program variables  $V$  binding EFMDP objects. This alignment allows inferred subgoal conditions to be lifted into program predicates. Thus, we use  $\psi$  and  $P$  interchangeably to refer to state abstraction predicates, based on the context.

**Program Synthesis for Policy Learning.** Given an EFMDP  $\mathcal{M}$  with **unknown** state transition probabilities and task specification  $\mathcal{T} : \text{True} \rightsquigarrow \psi$ , RoboScribe synthesizes a program  $\mathcal{P}^*$  as a controller in the DSL  $\mathcal{L}$  in Fig. 5 to fulfill  $\mathcal{T}$ . We outline the DSL operational semantics  $\langle S, \rho, s \rangle \Downarrow (\rho', s')$  in Fig. 14. Formally,  $S$  is a program statement in Fig. 5,  $\rho : V \rightarrow \mathcal{O}$  is an environment mapping that binds program variables  $V$  to objects in  $\mathcal{O}$ . For a program variable  $v \in V$ , at a given state  $s$ ,  $s(\rho(v))$  extracts the attribute values for the object referenced by  $v$  in  $s$ . Define  $\text{Exec}(\mathcal{M}, \mathcal{P})$  an interpreter that evaluates  $\mathcal{P}$  in the EFMDP  $\mathcal{M}$  based on the operational semantics and returns the EFMDP trajectory starting from a randomly sampled initial state  $s_0 \sim \eta(\cdot)$ .  $\text{Exec}$  terminates as soon as a specification-satisfying state is encountered. The learning objective is to synthesize  $\mathcal{P}^*$ :

$$\mathcal{P}^* = \arg \max_{\mathcal{P} \in \mathcal{L}} \Pr_{\zeta \sim \text{Exec}(\mathcal{M}, \mathcal{P})} [\zeta \models \mathcal{T} : \text{True} \rightsquigarrow \psi_R] \quad (5)$$

In practice, to evaluate task success, we define  $\mathcal{P}^* \models_{\mathcal{M}, \epsilon} \text{True} \rightsquigarrow \psi_R$  meaning that finite-length trajectories *sampled* from  $\text{Exec}(\mathcal{M}, \mathcal{P})$  *empirically* satisfy the goal condition  $\psi_R$  with a probability of at least  $1 - \epsilon$ . The values for the maximum trajectory length and  $\epsilon$  are user-configurable.

## 4 Abstraction Refinement-Guided Robot Control Program Synthesis

We present the core algorithms for the abstraction refinement-guided synthesis strategy in RoboScribe. We first provide the top-level synthesis algorithm, and then describe its key components.

### 4.1 Top-Level Algorithm

The top-level RoboScribe algorithm is presented in Algorithm 1. It takes as input an EFMDP  $\mathcal{M} = \{\Lambda, \mathcal{O}, \mathcal{S}, \mathcal{G}, \mathcal{A}, \mathbb{P}, \eta\}$ , the DSL  $\mathcal{L}$  (defined in Fig. 5), a set of task demonstrations  $\mathcal{D}$ , and a task specification  $\text{True} \rightsquigarrow \psi_R$ . **The DSL  $\mathcal{L}$  does not**

---

#### Algorithm 1 The RoboScribe Procedure

---

```

1: procedure ROBOSCRIBE( $\mathcal{M}, \mathcal{L}, \mathcal{D}, \text{True} \rightsquigarrow \psi_R$ )
2:    $T \leftarrow (N = \{u_{\text{True}}, u_{\psi_R}\}, E = \{u_{\text{True}} \rightarrow u_{\psi_R}\}, \psi_R)$ 
3:    $\mathcal{P}^*, T^* \leftarrow \text{SYNTHESIZE}(\mathcal{M}, \mathcal{L}, \mathcal{D}, T, \psi_R)$ 
4:   return  $\mathcal{P}^*$ 

```

---

**predefine useful state and action abstractions, which are yet to be learned as part of the synthesis process.** The objective is to synthesize a program  $\mathcal{P}^*$  in  $\mathcal{L}$  that satisfies the specification  $\text{True} \rightsquigarrow \psi_R$  (Eq. 5). As in conventional RL settings, the goal condition  $\psi_R$  is unknown to RoboScribe. However, the agent can use  $\psi_R$  as a black box to query whether any state encountered  $s$  satisfies  $\psi_R$ .

**Abstract Subtask Trees.** During its synthesis procedure, RoboScribe maintains state and action abstraction of a robot environment as an abstract subtask tree  $T$ —a hierarchical representation that encodes the sequence and relationships among subtasks for reaching the task’s goal states.

*Definition 4.1 (Abstract Subtask Tree).* An Abstract Subtask Tree  $T = (N, E, \psi_R)$  is a tuple:

- $N$  is a set of nodes, each representing a state abstraction predicate, denoted by  $\varphi$  or  $\psi$ , which defines a subset of the EFMDP state space. Throughout the paper, we use the terms predicate  $\varphi$  and tree node  $u_\varphi$  interchangeably.
- $E \subseteq N \times N$  is a set of directed edges between nodes, with each edge  $u_\varphi \rightarrow u_\psi \in E$  representing a subtask  $\varphi \rightsquigarrow \psi$  of the overall task, which transitions the agent from states characterized by  $\varphi$  to states in  $\psi$ . In the following, we also use edge  $u_\varphi \rightarrow u_\psi$  and subtask  $\varphi \rightsquigarrow \psi$  interchangeably.
- $\psi_R \in N$  is the root node, encapsulating the goal states of the overall task.

In an abstract subtask tree, state abstraction predicates on the tree nodes serve as decomposition of a complex robotic task. Tree edges represent abstract actions to transition between key subgoals. A tree path ( $\text{True} \rightsquigarrow \varphi_1 \rightsquigarrow \varphi_2 \rightsquigarrow \dots \rightsquigarrow \psi_R$ ) in  $T$ , leading toward the goal states at the root, is a sequence of subtasks, guiding the agent from one subtask to the next until the whole task is complete. The tree is *inverted*, for any node with multiple predecessors, the agent is directed to

$$\begin{array}{c}
\text{POLICY} \\
\frac{\mathcal{M}, \mathcal{D}, T \triangleright \mathcal{P} \quad \mathcal{P} \models_{\mathcal{M}, \epsilon} \text{True} \leadsto \psi}{\mathcal{M}, \mathcal{L}, \mathcal{D}, T \vdash (\mathcal{P}, T) : \psi} \\
\\
\text{REFINEMENT} \\
\frac{\begin{array}{l} \mathcal{M}, \mathcal{D}, T \triangleright \mathcal{P} \quad \mathcal{P} \not\models_{\mathcal{M}, \epsilon} \text{True} \leadsto \psi \quad P_s \leftarrow \{s_t \mid s_t \not\models \psi \wedge s_{t+1} \models \psi \wedge \{s_t, s_{t+1}\} \in \mathcal{D}\} \\ N_s \leftarrow \{\zeta \sim \text{Exec}(\mathcal{M}, \mathcal{P})\} \quad \varphi \in \mathcal{L} \wedge \varphi(s) = 1 \text{ for } s \in P_s \wedge \varphi(s) = 0 \text{ for } s \in N_s \\ \mathcal{M}, \mathcal{L}, \mathcal{D}, T[\text{True} \xrightarrow{\varphi} \psi] \vdash (\mathcal{P}', T') : \varphi \quad \mathcal{M}, \mathcal{L}, \mathcal{D}, T' \vdash (\mathcal{P}'', T'') : \psi \end{array}}{\mathcal{M}, \mathcal{L}, \mathcal{D}, T \vdash (\mathcal{P}'', T'') : \psi}
\end{array}$$

Fig. 15. The RoboScribe Synthesis Procedure

execute subtasks associated with each predecessor node, recursively. Fig. 7 displays the abstract subtask tree for a multi-object environment.

**Initial Environment Abstraction.** In Algorithm 1, at line 2, RoboScribe creates the initial environment abstraction as an abstract subtask tree  $T$  with two nodes  $u_{\text{True}}$  and  $u_{\psi_R}$  corresponding to the set of all possible initial states and environment states that satisfy the unknown goal condition. The edge  $u_{\text{True}} \rightarrow u_{\psi_R}$  represents a controller that satisfies the task specification  $\text{True} \leadsto \psi_R$ . At line 3, Algorithm 1 invokes the SYNTHESIZE procedure (detailed in Algorithm 2) to iteratively refines the coarse initial abstraction by need into a hierarchy of subtasks, continuing until a valid task-solving program is obtained within the abstracted environment.

## 4.2 The Main Synthesis Procedure

**Synthesis Rules.** We describe the SYNTHESIZE procedure using the synthesis rules of the following shape:

$$\mathcal{M}, \mathcal{L}, \mathcal{D}, T \vdash (\mathcal{P}^*, T^*) : \psi$$

where  $\mathcal{M}, \mathcal{L}, \mathcal{D}$  are the task EFMDP, our DSL (Fig. 5), and the task demonstrations respectively.  $T$  is an initial abstract subtask tree. The rule specifies the refinement of  $T$  into a valid abstraction  $T^*$ , which can then be converted into a program  $\mathcal{P}^*$  whose execution fulfills the goal condition  $\psi$ . Fig. 15 depicts the synthesis rules. Both rules rely on a procedure  $\mathcal{M}, \mathcal{D}, T \triangleright \mathcal{P}$  that synthesizes a program  $\mathcal{P}$  from the abstracted environment  $T$  and ground the abstract actions in  $\mathcal{P}$  as low-level neural controllers in the real environment  $\mathcal{M}$ . We defer the discussion of this procedure to Sec. 4.3.

The POLICY rule applies when the program  $\mathcal{P}$  derived from the abstract subtask tree  $T$  can directly satisfy the specification, i.e.,  $\mathcal{P} \models_{\mathcal{M}, \epsilon} \text{True} \leadsto \psi$  (we set  $1 - \epsilon$  as a lower bound for the probability of task success). In this case, the rule directly outputs  $(\mathcal{P}, T)$  as the synthesized solution. The REFINEMENT rule, on the other hand, addresses cases where the program  $\mathcal{P}$  generated from  $T$  does not fully solve the task. Here, our key idea is to systematically compare states within successful task demonstrations with those from failed behaviors by the program  $\mathcal{P}$  executed in the real environment, aiming to identify pivotal states that are essential for enabling task success. As stated in the REFINEMENT rule, RoboScribe extracts states just before task success from demonstrations, forming a set  $P_s = \{s_t \mid s_t \not\models \psi \wedge s_{t+1} \models \psi\}$ , and compares them with states  $N_s$  from the learned controller's trajectories to learn what prerequisite conditions are necessary. RoboScribe synthesizes a state abstraction predicate  $\varphi'$  to distinguish between  $P_s$  and  $N_s$ , refining the abstraction  $T$  by breaking down the task  $\text{True} \leadsto \psi$ , which results in a new abstract subtask tree  $T[\text{True} \xrightarrow{\varphi} \psi]$  that adds  $\varphi$  as an intermediate subgoal for  $\psi$  in  $T$  (formalized in Sec. 4.2.2). We hypothesize that having learned how to achieve  $\varphi$  by the synthesized program  $\mathcal{P}'$  from  $\varphi$ , it is an easier task



**Algorithm 2**  $\mathcal{M}, \mathcal{L}, \mathcal{D}, T \vdash (\mathcal{P}'', T'') : \psi$  The Main Synthesis Procedure

---

```

1: procedure SYNTHESIZE( $\mathcal{M}, \mathcal{L}, \mathcal{D}, T, \psi$ )
2:    $\mathcal{P} \leftarrow \text{GENPROGRAM}(\mathcal{M}, \mathcal{D}, T)$ 
3:   if  $\mathcal{P} \models_{\mathcal{M}, \epsilon} \text{True} \leadsto \psi$  then
4:     return  $\mathcal{P}, T$ 
5:   else
6:      $P_s \leftarrow \{s_t \mid s_t \not\models \psi \wedge s_{t+1} \models \psi \wedge \{s_t, s_{t+1}\} \in \mathcal{D}\}$ 
7:      $N_s \leftarrow \{\zeta \sim \text{Exec}(\mathcal{M}, \mathcal{P})\}$ 
8:      $\varphi \leftarrow \text{LEARNCLASSIFIER}(P_s, N_s, \mathcal{L}, T, \psi)$ 
9:      $T \leftarrow \text{UPDATETREE}(T, [\text{True} \xrightarrow{\varphi} \psi])$ 
10:    if  $\neg \text{VALIDATE}(T, \mathcal{D})$  then
11:      return FAIL
12:     $\mathcal{P}', T' \leftarrow \text{SYNTHESIZE}(\mathcal{M}, \mathcal{L}, \mathcal{D}, T, \varphi)$ 
13:     $\mathcal{P}'', T'' \leftarrow \text{SYNTHESIZE}(\mathcal{M}, \mathcal{L}, \mathcal{D}, T', \psi)$ 
14:    return  $\mathcal{P}'', T''$ 

```

---

**Algorithm 3**  $\varphi \in \mathcal{L}$  s.t.  $\varphi(s) = 1$  for  $s \in P_s$ ,  $\varphi(s) = 0$  for  $s \in N_s$ : Learn a classifier for  $P_s$  and  $N_s$ 


---

```

1: procedure LEARNCLASSIFIER( $P_s, N_s, \mathcal{L}, T, \psi$ )
2:   if  $\exists \varphi. \varphi \leadsto \psi \in T \wedge \forall s \in P_s. \text{ExistQuant}(\varphi)(s) \wedge \forall s \in N_s. \neg(\text{ExistQuant}(\varphi))(s)$  then
3:     return  $\varphi$ 
4:   else
5:      $\text{Exp}_\alpha \leftarrow \mathcal{L}(\alpha)$ 
6:      $\varphi \leftarrow \text{LEARNDECISIONTREE}(\text{Exp}_\alpha, P_s, N_s)$ 
7:     return  $\varphi$ 

```

---

for the agent to learn a program  $\mathcal{P}''$  based on  $\mathcal{P}'$  to achieve the goal condition  $\psi$ . Notably, the REFINEMENT rule embodies a recursive task decomposition process to repeatedly refine an initially coarse abstraction until a valid task-solving program can be obtained. Algorithm 2 operationalizes the synthesis rules in a recursive function SYNTHESIZE. We defer the discussion of termination, including the VALIDATE procedure at line 10, to Sec. 4.4. At line 8 and line 9, the SYNTHESIZE function invokes LEARNCLASSIFIER and UPDATETREE for state abstraction refinement. We formalize these two procedures below, starting with key notations.

Given a predicate  $\psi$  over EFMDP objects  $\mathcal{O} = \{\mu, e_1, \dots, e_N\}$  (defined in Fig. 12), let  $\text{Entities}(\psi) \subseteq \{e_1, \dots, e_N\}$  denote the set of entity variables that appear free in  $\psi$ . We define the transformation:

$$\text{ExistQuant}(\psi) := \exists e \in \text{Entities}(\psi). \psi$$

which lifts entity variables  $e_1, \dots, e_N$  appearing in  $\psi$  free existentially quantified.

**4.2.1 Learning State Abstraction Predicates.** We illustrate the state abstraction predicate learning procedure LEARNCLASSIFIER in Algorithm 3. It takes as input  $P_s$  the set of states prior to the states that satisfy a goal or subgoal condition  $\psi$  in the demonstration  $\mathcal{D}$ ,  $N_s$  the set of states along the agent's behavior that failed to reach  $\psi$ , the DSL  $\mathcal{L}$ , the abstract subtask tree  $T$  representing the current environment abstraction, and  $\psi$ , aiming to learn a state abstraction predicate capturing what should have been achieved by the agent in order to enable reaching states in  $\psi$ . At line 2, the algorithm checks whether the state abstraction predicate  $\varphi$  from any existing predecessor of  $\psi$  can be reused to distinguish  $P_s$  and  $N_s$  through  $\text{ExQuant}(\psi)$ . For example, in the Tower task described in Sec. 2.4, once the state abstraction predicate  $at(b, g(b))$  is identified as a subgoal for placing

---

**Algorithm 4**  $T[\text{True} \xrightarrow{\varphi} \psi]$ : Update an abstract subtask tree  $T = (N, E, \psi_R)$ 


---

```

1: procedure UPDATETREE( $T, [\text{True} \xrightarrow{\varphi} \psi]$ )
2:   if  $\exists \varphi'. \varphi' \rightsquigarrow \psi \in T \wedge \varphi' \equiv \varphi$  then
3:      $\varphi' \leftarrow \text{ExistQuant}(\varphi')$ 
4:      $u_{\varphi'} \leftarrow \textcircled{\varphi'}$ 
5:   else if  $\exists \varphi'. \varphi' \rightsquigarrow \psi \in T \wedge \text{ENTITIES}(\varphi') = \text{ENTITIES}(\varphi)$  then
6:      $N, E \leftarrow N \cup \{u_{\varphi}\}, E \setminus \{u_{\varphi'} \rightarrow u_{\psi}\} \cup \{u_{\varphi'} \rightarrow u_{\varphi}, u_{\varphi} \rightarrow u_{\psi}\}$ 
7:   else
8:      $N, E \leftarrow N \cup \{u_{\varphi}\}, E \cup \{u_{\text{True}} \rightarrow u_{\varphi}, u_{\varphi} \rightarrow u_{\psi}\}$ 

```

---

a block  $b$  at its target and a corresponding subroutine of controllers is learned, the task remains incomplete as additional blocks still require positioning. The existentially quantified predicate  $\exists b. \text{at}(b, g(b))$  helps distinguish demonstration states where blocks are correctly positioned near their targets from failed attempts by a single controller struggling to complete the task.

If reusing an existing predicate is not possible, in Algorithm 3, `LEARNCLASSIFIER` synthesizes a decision tree (DT) at line 6 to separate  $P_s$  and  $N_s$  using features from expressions  $\alpha$  derived from the production rules in our task specification language in Fig. 12. The hypothesis set of `LEARNDECISIONTREE` consists of Boolean combinations of predicates of the form  $\alpha \leq \phi$ , with  $\phi$  being a constant threshold, which are learned during training. Standard DT learning algorithms begin with an empty tree, greedily selecting features that maximize information gain, and continue until all leaves are labeled with a single class. Finally, the learned DT is converted to a predicate  $\varphi$ .

**4.2.2 Refining Abstract Subtask Trees.** With the learned classifier  $\varphi$  as a prerequisite for achieving  $\psi$  in a task  $\text{True} \rightsquigarrow \psi$ , `UPDATETREE` refines the abstract subtask tree  $T$  by using  $\varphi$  as a subgoal to decompose the task. We specify this refinement procedure  $T[\text{True} \xrightarrow{\varphi} \psi]$  in Algorithm 4. At line 2, if the newly identified subtask goal predicate  $\varphi$  matches an existing subtask goal  $\varphi'$  that precedes  $\psi$  in  $T$  (i.e., the `LEARNCLASSIFIER` procedure has opted to reuse  $\varphi'$ ), a repeating subroutine is effectively recognized. This means that the controllers designed to achieve  $\varphi'$  for manipulating some objects can be repurposed to solve the subtask to reach  $\varphi$  for a different set of objects of the same type. The algorithm marks  $u_{\varphi'}$  as a circular node  $\textcircled{\varphi'}$  at line 4 to indicate this repetition - the control strategy used to reach  $\varphi'$  should then be executed iteratively to address recurring objects in the environment that have not yet satisfied  $\varphi'$ . For this purpose, we have made  $\varphi'$  existentially quantified.

An invariant we maintain for an abstract subtask tree  $T$  is that for each tree node  $\psi$ , for any predecessors  $\varphi \rightsquigarrow \psi$  and  $\varphi' \rightsquigarrow \psi$ ,  $\varphi$  and  $\varphi'$  are subgoal conditions for unique sets of objects, and hence requiring different control strategies (i.e. unique tree paths towards  $\varphi$  and  $\varphi'$ ) and otherwise they should be collapsed into a circular node. For example, consider the `PlaceCubeDrawer` task depicted in Fig. 7, which has an abstract subtask tree consisting of two main paths: one for opening the drawer and the other for placing the cube inside it. During the refinement of  $T[\text{True} \xrightarrow{\varphi} \psi]$ , if a newly identified subgoal  $\varphi$  targets the same set of entities as an existing predecessor  $\varphi'$  of  $\psi$  (i.e.  $\text{Entities}(\varphi) = \text{Entities}(\varphi')$ ), the `UPDATETREE` procedure in Algorithm 4 inserts a new node  $u_{\varphi'}$  between  $u_{\varphi}$  and  $u_{\psi}$  at line 6, establishing  $\varphi'$  as an intermediate subgoal for  $\varphi \rightsquigarrow \psi$ . For example, in the tree refinement illustrated by Equation 2, the predicate  $\text{at}(b, g(b))$ , which signifies that block  $b$  is at its goal, is added between the predicate  $\text{hold}(\mu, b)$ —indicating the robot is gripping  $b$ —and  $\psi_R$ , the overall task’s goal condition. If  $\varphi$  involves different set of entities from any existing predecessor

---

**Algorithm 5**  $\mathcal{M}, \mathcal{D}, T \triangleright \mathcal{P}$ : Synthesize a program  $\mathcal{P}$  from an abstract subtask tree  $T = (N, E, \psi_R)$ 


---

```

1: procedure GENPROGRAM( $\mathcal{M}, \mathcal{D}, T$ )
2:    $\mathcal{P} \leftarrow \text{TREE2PROGRAM}(\psi_R, \mathcal{D}, T)$ 
3:   for all  $o \in \text{FreeVars}(\mathcal{P})$  : do
4:      $\mathcal{P} \leftarrow \{v := \text{get}(\lambda v : \Lambda(o). \text{True})\}; \mathcal{P}|_{o \mapsto v}$   $\triangleright v$  fresh
5:    $\mathcal{P}^* \leftarrow \text{TRAINPROGRAM}(\mathcal{M}, \mathcal{P})$ 
6:   return  $\mathcal{P}^*$ 
7:
8: procedure TREE2PROGRAM( $u_\varphi, \mathcal{D}, T$ )
9:    $\mathcal{P}_{u_\varphi} \leftarrow \{\}$ 
10:  for all edge  $(e \equiv u_{\varphi'} \rightarrow u_\varphi) \in T$  sorted by  $i(\mathcal{D}, \varphi')$  do
11:     $\mathcal{P}_{u_{\varphi'}} \leftarrow \text{TREE2PROGRAM}(u_{\varphi'}, \mathcal{D}, T)$ 
12:     $\mathcal{P}_{u_\varphi} \leftarrow \mathcal{P}_{u_\varphi}; \mathcal{P}_{u_{\varphi'}}$ 
13:  if  $u_\varphi \equiv \bigcirc_{\exists o. \varphi}$  then  $\triangleright u_\varphi$  induces an iterative procedure
14:     $v \leftarrow \text{FreshVar}()$ 
15:     $p \leftarrow \text{PREDICATESYNTHESIS}(\mathcal{D}, v, \exists o. \varphi)$   $\triangleright$  Loop Condition Synthesis
16:     $\mathcal{P}_{u_\varphi} \leftarrow \text{while}(v := \text{get}(\lambda v. p)) \{(\mathcal{P}_{u_\varphi}; \pi_\varphi[\varphi])|_{o \mapsto v}\};$ 
17:  else
18:     $\mathcal{P}_{u_\varphi} \leftarrow \mathcal{P}_{u_\varphi}; \pi_\varphi[\varphi]$ 
19:  return  $\mathcal{P}_{u_\varphi}$ 

```

---

$\varphi'$  of  $\psi$ , a distinct tree path from  $u_\varphi$  to  $u_{\psi'}$  is created at line 8, representing distinct objects to control, as exemplified in the unique two paths in the abstract subtask tree for PlaceCubeDrawer in Fig. 7.

### 4.3 Synthesizing Robot-Control Programs from Abstract Subtask Trees

An important step in RoboScribe is synthesizing an executable program  $\mathcal{P}$  from an abstract subtask tree  $T$ . We formalize this procedure as  $\mathcal{M}, \mathcal{D}, T \triangleright \mathcal{P}$ , implemented in a procedure GENPROGRAM described in Algorithm 5. Our main synthesis procedure SYNTHESIZE invokes GENPROGRAM in Algorithm 2 (line 2) to generate a candidate program and does so repeatedly for each *refined* abstract subtask tree until a specification-satisfying program can be synthesized.

The GENPROGRAM procedure  $\mathcal{M}, \mathcal{D}, T \triangleright \mathcal{P}$  accomplishes three main objectives as formalized in Algorithm 5: (1) it constructs the "skeleton" of  $\mathcal{P}$  based on the hierarchical structure of  $T$  via the TREE2PROGRAM procedure at line 2. (2) TREE2PROGRAM also infers loops within  $\mathcal{P}$  to handle varying numbers of objects by detecting repeated patterns in the demonstration  $\mathcal{D}$ . (3) It grounds the abstract actions in  $\mathcal{P}$  as low-level controllers that can solve the subtasks within  $T$  through the TRAINPROGRAM procedure at line 5. These low-level controllers are neural network policies that operate directly in the robot environment to control robot actions.

**4.3.1 Program Generation.** The TREE2PROGRAM procedure (Line 8 of Algorithm 5) traverses an abstract subtask tree  $T$  rooted at  $u_\varphi$ . It generates a program  $\mathcal{P}_{u_\varphi}$  from  $T$  for solving the task of reaching states satisfying  $\varphi$ . Here we assume that in a multi-object setting, task demonstrations implicitly indicate the order in which multiple objects of different types should be handled (we relax this assumption in Sec. 4.5). For example, in the PlaceCubeDrawer task shown in Fig. 7, the agent must first pull the drawer open using its handle before placing the cubes inside. Define  $i(\mathcal{D}, \varphi')$  as the position in the demonstration where  $\varphi'$  holds. TREE2PROGRAM enumerates the

**Algorithm 6** Learning Loop Conditions from Demonstrations

---

```

1: procedure PREDICATESYNTHESIS( $\mathcal{D}, v, \exists o. \varphi$ )
2:   Sample a demonstration rollout  $d \sim \mathcal{D}$ 
3:   Let  $e_1, \dots, e_N$  be the objects sorted by their order of satisfying  $\varphi[o \mapsto e_i]$  in  $d$  and
4:    $S_i$  be the segment of  $d$  during which  $e_i$  is manipulated
5:   for  $1 \leq i \leq N$  do
6:      $P_i, N_i \leftarrow \{\}, \{\}$ 
7:   for all  $e_i$  do
8:      $P_i \leftarrow P_i \cup \{(s_t, e_i) \mid s_t \in S_i\}$ 
9:     for all  $S_j$  s.t.  $j < i$  do
10:       $N_i \leftarrow N_i \cup \{(s_t, e_i) \mid s_t \in S_j\}$ 
11:     for all  $S_j$  s.t.  $j > i$  do
12:       $N_i \leftarrow N_i \cup \{(s_t, e_i) \mid s_t \in S_j\}$ 
13:    $P \leftarrow \text{True}$ 
14:   for  $1 \leq i \leq N$  do
15:      $\psi \leftarrow \text{TOPDOWNENUM}(P_i, N_i)$ 
16:      $P \leftarrow P \wedge \psi[e_i \mapsto v]$ 
17:   return  $P$ 

```

---

incoming edges  $u_{\varphi'} \rightarrow u_{\varphi}$  of  $u_{\varphi}$  (line 10) in the order of  $i(\mathcal{D}, \varphi')$ , recursively applying itself to  $u_{\varphi'}$  (line 11), and appending the resulting program  $P_{u_{\varphi'}}$  to  $P_{u_{\varphi}}$  (line 12). At line 18, the algorithm appends a low-level controller  $\pi_{\varphi} = \pi_{\theta}(\mu, \{o\}, \{g(o')\})$ , a neural network policy with trainable weights  $\theta$ , to the program  $P_{u_{\varphi}}$ . This controller guides the agent from states satisfying the subgoals in the predecessors of  $u_{\varphi}$  (namely  $\{\varphi'\}$ ) to states that satisfy  $\varphi$  in  $u_{\varphi}$ . Here,  $\{o\} = \text{Entities}(\varphi)$  denotes the set of entities involved in  $\varphi$ , while  $\{g(o')\}$  represents the goal conditions in  $\varphi$  for these entities, with  $o'$  potentially being a subset of  $o$ . The controller  $\pi_{\varphi}$  needs to manage the entities in  $o$  to achieve the subgoal condition  $\varphi$ . If  $u_{\varphi}$  is designated as a circular node, as constructed in Sec. 4.2.2, the subroutine synthesized in  $\mathcal{P}u_{\varphi}$  is designed for repeated execution to handle recurring objects of the same types to achieve  $\varphi$  in a loop. To simplify the presentation, we assume a single existential quantifier for the subgoal condition related to  $u_{\varphi}$ , though the algorithm trivially extends to multiple quantifiers. We introduce a fresh program variable  $v$  to bind recurring objects within the loop. At line 16, TREE2PROGRAM constructs a loop with  $(\mathcal{P}u_{\varphi}; \pi_{\varphi}|_{o \mapsto v})$  as its body, replacing  $o$  in the program with  $v$  to track recurring objects bound in each iteration, analogous to existential quantifier instantiation. Fig. 10 illustrates this process for the Tower task. We note that at Line 4 of Algorithm 5, for any remaining object identifier  $o$  in a synthesized program  $\mathcal{P}$ , we similarly project  $o$  to a fresh variable  $v$  and prepend  $v := \text{get}(\lambda v : \Lambda(o). \text{True})$  to  $\mathcal{P}$ . This allows the program to generalize across environments by retrieving the appropriate entity, removing dependencies on specific object identifiers (see Fig. 8 as an example).

So far, the algorithm does not specify the order in which objects are addressed in the loop iterations, which is crucial for tasks with dependencies, such as Tower (Fig. 9) where placing a block in its goal position before positioning the underlying blocks can result in failure. To this end, TREE2PROGRAM invokes PREDICATESYNTHESIS in Algorithm 6 to specify an effective handling sequence for the loop structure. We outline this procedure as follows.

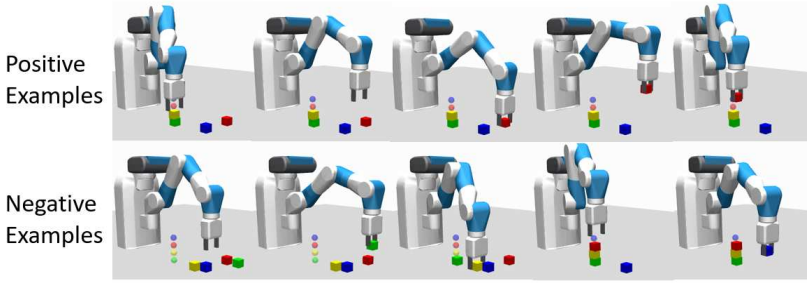


Fig. 16. Positive and negative examples for learning the predicate for when to handle the red block in Tower. In the demonstration, the agent stacks green, yellow, red and blue blocks in order.

**4.3.2 Loop Inference.** Our key approach to identifying the potential order for an effective handling sequence for recurring objects is analyzing the underlying rationale in task demonstrations to understand why certain objects must be handled before others.

Given the set of demonstrations  $\mathcal{D}$  and a circular tree node  $\varphi$  in an abstract subtask tree  $T$ , the PREDICATESYNTHESIS procedure in Algorithm 6 formalizes the generation of loop conditions for synthesized iterative programs. As we assume demonstrations in  $\mathcal{D}$  are consistent (Sec. 2.1), at line 3, we randomly sample a rollout  $d$  from  $\mathcal{D}$  to infer the handling sequence for recurring objects. We sort objects  $e_1, e_2, \dots$  of the desired types in  $d$  according to their order of satisfying  $\varphi$  in  $d$  and partition  $d$  according to this order. Each segment  $S_i$  corresponds to the subtask period in which  $e_i$  is manipulated within the demonstration. For example, for the abstraction of Tower in Eq. 3, given the demonstration in Fig. 9, we instantiate the existential quantifiers  $b$  in the circular node  $\exists b. at(b, g(b))$  with the colored blocks respectively. The agent stacks the green, yellow, red, and blue blocks sequentially from bottom to top, resulting in partitions  $S_1, S_2, S_3$  and  $S_4$  in Fig. 9.

For each  $e_i$ , RoboScribe maintains positive examples  $(s_t, e_i)$  for all states  $s_t$  in  $S_i$  where  $e_i$  is handled at line 8, and negative examples  $(s_t, e_i)$  for all states  $s_t$  in prior  $(S_0, \dots, S_{i-1})$  and subsequent  $(S_{i+1}, \dots)$  partitions to illustrate why  $e_i$  should not be handled earlier at line 10 or does not need to be handled afterwards at line 12 in these negative states. For instance, in Fig. 16, RoboScribe shows why the red block is placed after the green and yellow blocks. In the top row, where the green and yellow blocks are already positioned, any state involving the placement of the red block is considered a positive example. In the bottom row, where the green and yellow blocks are not yet in their goal positions, the placement of the red block is marked as a negative example (it should not be handled), and any subsequent states after the red block is positioned are also negative examples (since its handling is already complete). Given the positive and negative examples  $\{(s_t, e_i)\}^+, \{(s_t, e_i)\}^-$ , any classifier that defines the relationship between  $e_i$  and other objects in  $s_t$  (abstracted as existential variables) and effectively separates the examples, provides both an ordering and a termination constraint for handling  $e_i$  during manipulation.

At line 14, RoboScribe uses top-down synthesis to generate a classifier predicate for each  $e_i$ , following the production rules for predicates  $\psi$  in our task specification language (see Fig. 12). In this process, we augment  $\psi$  with learned state abstraction predicates from abstract subtask trees as these predicates provide additional task-relevant constraints. The learned predicates for all  $e_i$  are combined to fill in the loop condition in the synthesized program. Here we use top-down enumeration instead of decision tree (DT) learning because loop conditions for handling sequences typically require existential quantifiers to manage unbounded entities with dependencies, which are not well-suited for DTs. If the synthesis algorithm does not find a classifier predicate for an



object  $e_i$  within a reasonable search budget, RoboScribe interprets this as the absence of an ordering constraint for that  $e_i$  — the agent can select  $e_i$  for manipulation without restrictions, and it then returns true in this case. Given the examples in Fig. 16, RoboScribe synthesizes the ordering predicate in Eq. 4 for Tower. This predicate ensures that any block with a lower goal position must be placed before the current block.

We note that a synthesized ordering predicate naturally encodes loop termination, even in environments without explicit total or partial order constraints. For example, consider a scenario where a robot arm must push several objects from a table into a drawer in arbitrary order. Algorithm 6 learns a termination condition in the form of a classifier, which identifies why an object that has already been placed in the drawer is no longer manipulated thereafter.

**4.3.3 Reinforcement Program Learning.** Given a program  $\mathcal{P}$  inferred from an abstract subtask tree  $T$  reflecting the current abstraction of the real robot environment  $\mathcal{M}$ , the TRAINPROGRAM procedure called at line 5 in Algorithm 5 grounds  $\mathcal{P}$  in  $\mathcal{M}$  by learning the low-level neural controllers invoked by  $\mathcal{P}$  to fulfill the subtasks within  $T$ . We maintain separate buffers  $B_\varphi$  to store trajectories associated with each low-level controller  $\pi_\varphi$  within  $\mathcal{P}$ . Program trajectories  $\zeta$  are sampled by executing  $\mathcal{P}$  in the real environment  $\zeta \sim \text{Exec}(\mathcal{M}, \mathcal{P})$ . Each sub-trajectory of  $\zeta$  generated by a specific controller  $\pi_\varphi$  is stored only in the corresponding buffer  $B_\varphi$ . During each gradient update step, TRAINPROGRAM updates each policy  $\pi_\varphi$  by sampling from its buffer  $B_\varphi$  and optimizes it using any off-the-shelf off-policy RL algorithm (e.g. Soft Actor-Critic), aiming to maximize the expected reward for  $\pi_\varphi$ :

$$\pi_\varphi = \arg \max_{\pi_\varphi} \mathbb{E}_{\zeta=s_0, a_0, s_1, \dots, s_L, a_L \sim B_\varphi} \left[ \sum_{i=0}^L \gamma^i R_\varphi(s_i, a_i) \right]$$

where  $R_\varphi$  denotes the reward function used to train  $\pi_\varphi$ ,  $L$  is the sampled trajectory length, and  $\gamma$  is the discount factor.

For each subtask to learn  $\pi_\varphi$ , our training procedure aims to construct the reward function  $R_\varphi$  that provides feedback based on the satisfaction of these predicates throughout a policy trajectory. Instead of only using  $\varphi$  to provide a binary signal indicating whether a subgoal state has been achieved,  $R_\varphi$  quantifies a continuous measure of state proximity between the current state  $s$  and the satisfaction of  $\varphi$  to enable smoother policy optimization. This approach allows for a more granular assessment of progress, guiding the agent incrementally towards the subgoal states in  $\varphi$ . Formally, we define the reward function  $R_\varphi$  recursively based on the structure of the predicate  $\varphi$ :

$$R_\varphi(s) = \begin{cases} R(\varphi_1 \wedge \varphi_2) = \min(R_{\varphi_1}(s), R_{\varphi_2}(s)) & \text{if } \varphi = \varphi_1 \wedge \varphi_2, \\ R(\varphi_1 \vee \varphi_2) = \max(R_{\varphi_1}(s), R_{\varphi_2}(s)) & \text{if } \varphi = \varphi_1 \vee \varphi_2, \\ R(\alpha > \phi) = \alpha(s) - \phi & \text{if } \varphi = \alpha > \phi, \\ R(\alpha < \phi) = \phi - \alpha(s) & \text{if } \varphi = \alpha < \phi. \end{cases}$$

#### 4.4 Termination Guarantee

The main synthesis procedure in RoboScribe (Algorithm 2) is guaranteed to terminate after finitely many refinements. This follows from the bounded size of abstract subtask trees, which is constrained by both the number of environment entities and the length of demonstrations. Given an EFMDP  $\mathcal{M}$  and demonstrations  $\mathcal{D}$ : (1) each abstract subtask tree node can lead to only finitely many paths, bounded by the number of entities in  $\mathcal{M}$ ; and (2) each path has bounded length, since every edge  $u_\varphi \rightarrow u_\psi$  must be supported by contiguous segments in  $\mathcal{D}$ , of which there are finitely many. Consequently, the number of refinements is bounded. Algorithm 2 must terminate, either by synthesizing a valid program or returning FAIL (Line 11) when the current abstract subtask tree cannot be further refined.

Specifically, in Algorithm 2, each refinement of an abstract subtask tree  $T$  (Line 9, Algorithm2) is followed by a `VALIDATE` call at line 10, which scans  $T$  against demonstration  $\mathcal{D}$  to determine whether it can yield a task-solving program or if the underlying task is unsolvable under  $T$ . Given demonstrations  $\mathcal{D} = \{(s_0, s_1, \dots, s_H)\}$  and a predicate  $\varphi$ , if there exists  $i$  such that  $\varphi(s_i)$  holds, let  $i_s(\mathcal{D}, \varphi)$  denote the first timestep where  $\varphi$  holds, and  $i_f(\mathcal{D}, \varphi)$  the last. For each edge  $u_\varphi \rightarrow u_\psi$  representing a subtask  $\varphi \rightsquigarrow \psi$  in  $T$ , the `VALIDATE` procedure checks whether  $i_s(\mathcal{D}, \psi) \leq i_f(\mathcal{D}, \varphi)$  when  $\varphi$  is not True, or  $i_s(\mathcal{D}, \psi) = 0$ . Quantifiers in the predicates are instantiated by enumerating all entities of the appropriate types present in  $\mathcal{D}$ . If true, this indicates that the intermediate subgoal  $\psi$  does not meaningfully advance the task from  $\varphi$ , causing it to fail under the current abstraction.

**What is the main cause for RoboScribe to return FAIL?** It occurs when the RL algorithm fails to learn a valid low-level controller even for some short-horizon subtasks, leading to repeated refinements and eventual task failure. This may stem from poor exploration, challenging environment dynamics, or instability in the learning algorithm; however, we have rarely observed this in practice.

#### 4.5 Extension: Conditional Statements

The `GENPROGRAM` procedure in Algorithm 5 operates under the assumption that demonstrations implicitly suggest the order for handling multiple objects in a multi-object task. However, this order may vary depending on the goal conditions. For instance, consider a scenario with a peg (blue) and a cube (red) in Fig. 17. If the peg’s goal region is above the cube’s goal, the task must be completed by first moving the cube, then the peg—and vice versa. Our implementation relaxes this assumption by repurposing the `PREDICATESYNTHESIS` algorithm from Algorithm 6 (developed for sorting objects represented by a circular abstract subtask tree node) to determine how to order the execution of the multiple predecessors of a tree node corresponding to objects of different types. This approach identifies the conditions under which each object should be processed and encodes these conditions into conditional statements, thereby selecting the appropriate handling sequence.

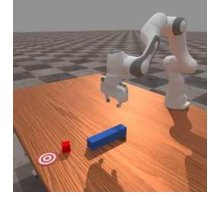


Fig. 17. A Pick&Place Environment for a cube and a peg.

## 5 Experiments

RoboScribe is implemented in Python. In the implementation, low-level neural policies are Multi-layer Perception (MLP) containing two hidden layers with 256 neurons. We leverage Soft Actor-Critic (SAC) [24] from Stable-Baseline3 [49] as the RL algorithm to train the policies.

Our experiments are designed to answer the following research questions:

- **(RQ1)** Is RoboScribe able to learn effective and interpretable programs?
- **(RQ2)** Does the iterative program learned by RoboScribe generalize to unseen environments without further training?

**Main Baselines.** Throughout the evaluation, we consider the following baselines:

- **BC:** Behavior Cloning (BC) is a standard learning from demonstration baseline. It applies supervised learning to train a policy that replicates expert actions for given states in demonstrations.
- **GAIL [25]:** Generative Adversarial Imitation Learning (GAIL) works by alternating between training two components: a discriminator and an agent. The discriminator is trained to distinguish between states visited by the expert and those visited by the agent. The agent, in turn, is optimized to visit states that make this distinction more difficult, thereby encouraging it to mimic the expert’s behavior. We select GAIL as a baseline because its discriminator serves a role similar to the state-abstraction predicates we learn as classifiers.

- **goalGAIL [20]:** GoalGAIL combines GAIL with Hindsight Experience Replay (HER) [3]. Unlike GAIL, which only imitates expert state distributions without explicit task objectives, GoalGAIL introduces goals by relabeling the agent’s trajectory with states that are actually achieved during execution, allowing the agent to treat these states as successful outcomes and thereby receive more frequent feedback, leading to faster convergence and improved sample efficiency. We select it as a baseline because it is a stronger variant of GAIL that accelerates learning and significantly improves sample efficiency.
- **DeepSet [67]:** DeepSet embraces an entity-based compositional structure in its neural policy representation based on Self-Attention [61] to leverage the symmetries and invariances in the EFMDP. Like RoboScribe programs, its policy architectures decompose goal-conditioned tasks into their constituent entities and subgoals.

For fair comparisons, we use the DeepSet architecture for the policy in BC and for both the policy and discriminator in GAIL and goalGAIL. DeepSet’s ability to handle an arbitrary number of input objects makes it well-suited for multi-object environments.

We *exclude direct quantitative comparisons with existing programmatic RL methods* like PROLEX [44], Tabula [46], and ReGuS [14] because they rely on predefined DSLs with manually crafted state and action abstractions, whereas RoboScribe autonomously discovers these abstractions. This fundamental difference makes direct performance comparisons impractical.

**Benchmarks.** We use a suite of challenging robot manipulation environments including Pick&Place (Fig. 1a); Tower-5 (Fig. 1c) where the goal is to assemble 5 scattered blocks into a tower (88 state dimensions); Pick&Place-Cond shown in Fig. 17 where the robot stacks a cube and a peg based on their goal position ordering; Pick&Place-4 shown in Fig. 18a in the Pick&place-Multi environment where the goal is placing 4 blocks in their designated goal regions on a surface, with the final block needing to be hung by the gripper in the air at its goal position; Push-3 shown in Fig. 18b in the Push-Multi environment where the goal is pushing 3 blocks to their related goal regions on a table surface; Meta-World where a robot needs to be controlled to complete 3 tasks, including pushing the mug back, opening the drawer and turning the faucet left; and PlaceCubesDrawer visualized in Fig. 18d. In the challenging PlaceCubesDrawer environment (134 state dimensions) from [39], the agent needs to open a drawer and iteratively places three cubes into the drawer. We consider a sparse reward setting in which the agent receives reward 1.0 when the entire task is completed successfully and 0 otherwise.

**Demonstration Collection.** For each environment, we provide a demonstration dataset consisting of 50 successful trajectories, collected by manually controlling the end effector in a simulator to manipulate the objects. We place no requirements on demonstration quality with respect to path optimality or goal diversity. Since RoboScribe relies on demonstrations only to infer abstract subtask trees, it remains effective even when demonstration trajectories contain unnecessarily long control paths. Such demonstrations still reveal the underlying abstract task structure, which

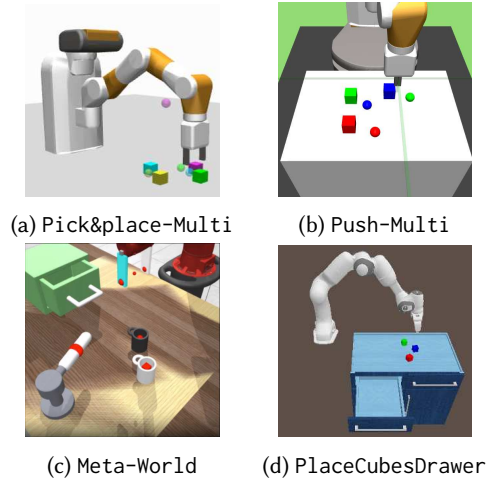


Fig. 18. Testing environments with multiple entities.

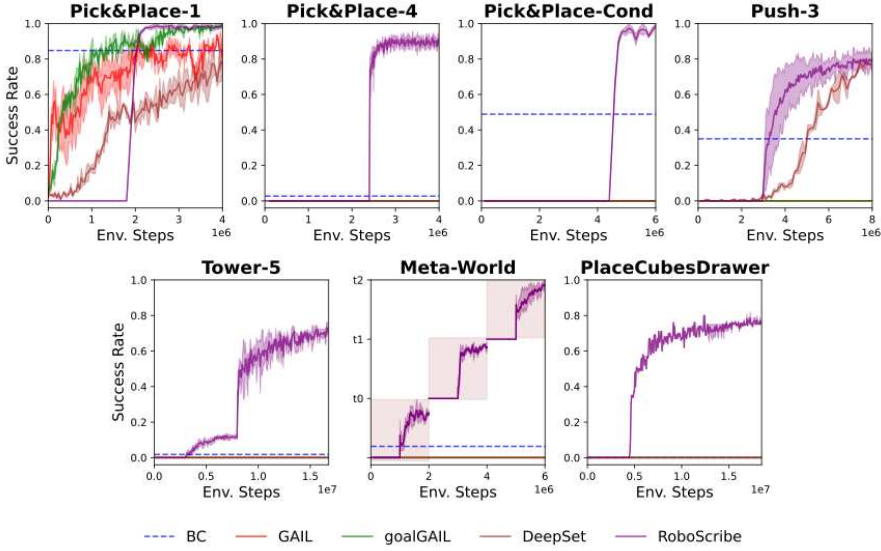


Fig. 19. Rewards for all the tools throughout the training phase. The solid curve represents the mean across 5 random seeds. The shaded area indicates the standard deviation. In Meta-World, we report the success rates for each subtask—pushing the mug back (t0), opening the drawer (t1), and turning the faucet left (t2)—in the order that RoboScribe discovers them.

Table 1. Average number of refinement steps for synthesized programs (over five random seeds).

Pick&Place-1	Pick&Place-4	Pick&Place-Cond	Push-3	Tower-5	Meta-World	PlaceCubesDrawer
2.2	2.2	5.4	4.4	4.4	5.6	4.6

RoboScribe leverages to learn efficient low-level controllers via reward signals derived from the state-abstraction predicates encoded in abstract subtask trees.

### 5.1 RQ1: Learning Efficiency and Interpretability

For each environment, we train RoboScribe and the baseline methods with 5 random seeds, reporting their evaluation success rates during training, as shown in Fig. 19. While RoboScribe initially experiences a flat zero success rate early in training, it focuses on comparative abstraction refinement to discover the abstract task structure and grounding abstract actions to reach automatically discovered subgoal conditions, guiding the agent towards the overall goal progressively. RoboScribe’s success rate increases rapidly once the program structure is fully developed, eventually surpassing the performance of the baseline methods. For Tower-5, there is a sharp increase in success rates around  $1e7$  steps. This is because the task requires the end effector to move its hand away from the top block to a certain height to ensure stable tower construction. The final subgoal of moving away the end effector is relatively easier to learn, and by this point, the agent has already mastered stacking the blocks. As a result, the final task success rate improves significantly after this. Other than Pick&place and Push-3, the baselines struggle to achieve progress due to the complexity of the observations involving multiple objects and the sparsity of the reward signals.

We report the average number of refinement steps for our synthesized programs in Table 1. Even the most challenging tasks such as Tower-5 and PlaceCubesDrawer require a modest number of refinements, demonstrating the efficiency of RoboScribe’s synthesis procedure.

**Demonstration Size.** RoboScribe requires only a small number of demonstrations to infer abstract subtask trees. We evaluated RoboScribe in the Pick&Place environment (Fig. 1a) across 10, 20, and 50 demonstrations and found that performance remains stable: converging to 99.8% success with 2.28M environment steps, 99.0% with 1.68M, and 99.7% with 2.20M steps, respectively.

**Interpretability.** Programs synthesized by RoboScribe provide greater interpretability of the control logic compared to black-box policies. Synthesized programs include loops conditioned on explicit subgoals necessary for recurring interactions with multiple objects. This yields structured policy representations that make the decision-making process transparent and easy to understand. For example, in the Push-3 task, where the end effector  $\mu$  must push a block  $b$  to a target position  $g(b)$  on a table, one of the state abstraction predicates RoboScribe synthesizes is:

$$\arctan2(\mu \downarrow_{x,y} - b \downarrow_{x,y}) - \arctan2(g(b) \downarrow_{x,y} - b \downarrow_{x,y}) < \phi$$

This predicate captures a subgoal condition where the end effector, block, and goal region are aligned for direct pushing. Occasionally, more complex predicates arise. In one of five Push-3 trials, we observed the following predicate in the abstract subtask tree:

$$\|b \downarrow_{x,y,z} - g(b) \downarrow_{x,y,z}\| < \phi_1 \wedge \|\mu \downarrow_{x,y,z} - b \downarrow_{x,y,z}\| < \phi_2 \wedge \arctan2(\mu \downarrow_{x,y} - b \downarrow_{x,y}) - \arctan2(g(b) \downarrow_{x,y} - b \downarrow_{x,y}) < \phi_3$$

This predicate is true when block  $b$  is near its goal region  $g(b)$  and the end-effector  $\mu$  is not only close to  $b$ , but also aligned with its goal direction. This complex predicate arises due to failure execution trajectories containing two types of negative states: (1) the block remains far from its goal, or (2) it is pushed somewhat close to the goal but not sufficiently, and the end-effector then moves away unexpectedly. The latter case is caused by the imperfect behavior of the learned low-level policy for direct pushing. Although RoboScribe can eventually succeed in such cases, the resulting predicates often increase in complexity, requiring additional effort for interpretation.

## 5.2 RQ2: Generalization to New Environments

For the Tower and Push-Multi environments, the capability of handling arbitrary numbers of objects is desired. We analyze the transferability of the synthesized iterative program to diverse environment settings.

**Tower Environment.** We synthesize the Tower program in a single tower setting with 4 or 5 blocks and evaluate its performance across diverse environments without further training, as shown in Fig. 20. These environments include a taller single tower with 6 or 7 blocks, multiple towers with 2 to 3 blocks per tower, and a pyramid tower with 4 to 9 blocks. As the baselines discussed in Sec. 5.1 fail to solve the Tower task, we turn to the curriculum learning-based approach ReNN [32], which progressively learns to stack 2, 3, and ultimately 5 blocks. In contrast, RoboScribe learns to handle all 5 blocks directly, without relying on an expert-designed curriculum, offering a more flexible and practical solution. Fig. 21 presents the results for RoboScribe and ReNN [32]. In the single tower setting, ReNN benefits from curriculum learning, achieving slightly better results when the training and evaluation environments match (e.g., single towers with 4 and 5 blocks). However, RoboScribe

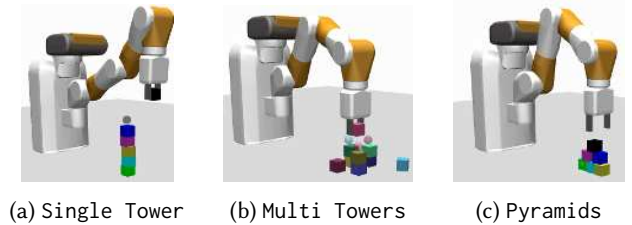


Fig. 20. Novel Tower environments.



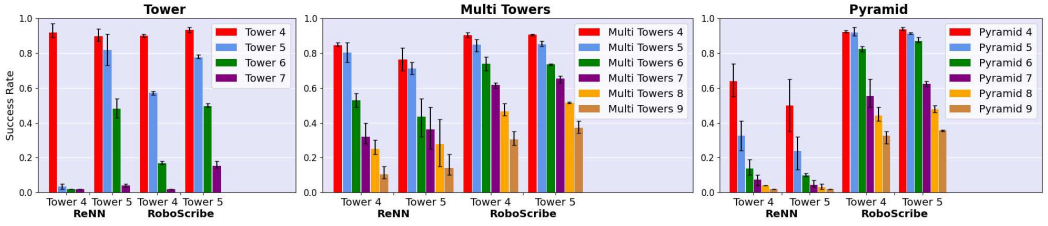


Fig. 21. Comparison between ReNN and RoboScribe on zero-shot generalization to new Tower environment settings. Specifically, policies trained on single tower with 4 blocks or 5 blocks are evaluated on Single (but taller) towers, multiple towers and pyramid configurations with varying numbers of blocks. Success rate is reported as accuracy of completing a task averaged over 500 episodes.

demonstrates superior generalization. For instance, when transferring a policy trained on 4 blocks to a taller single tower, ReNN achieves less than 5% success, while RoboScribe trained on 4 blocks achieves 55% ( $\pm 1\%$ ) success with 5 blocks and 15% ( $\pm 1\%$ ) with 6 blocks. In both multi-tower and pyramid settings, RoboScribe significantly outperforms ReNN across block counts from 4 to 9.

**Push-Multi Environment.** To evaluate the iterative program learned in the Push-Multi environment, we introduce a confined version, Push-Multi Confined (Fig. 22), where goals are randomly arranged along a line near the table's upper edge. To achieve the goal condition, the robot benefits from pushing the blocks in a certain order. For example, in Fig. 22, the robot should push the blue block first, followed by green, then red. Pushing blocks out of this order, such as green or red first, may obstruct the blue block's path to its goal. For RoboScribe, we reuse the learned iterative program by providing demonstrations of the correct entity handling sequence, allowing it to update its loop condition without additional training. In contrast, we continue training the baseline DeepSet [67] model on the confined environment until convergence, as it cannot structurally update its model like RoboScribe. RoboScribe correctly learns the entity handling order from the demonstration and updates the loop condition accordingly. In comparison, DeepSet achieves a success rate of 81.3% ( $\pm 4.0\%$ ), while RoboScribe achieves 86.3% ( $\pm 0.5\%$ ), averaged over 500 episodes. RoboScribe's superior performance demonstrates the generalization of learned policies and program structures. Additionally, with the interpretability of its programmatic policy, RoboScribe offers greater flexibility in transferring the synthesized program to different environment settings.

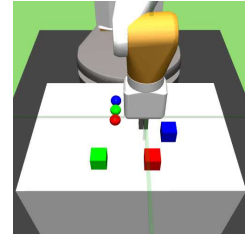


Fig. 22. A constrained Push-Multi environment.

## 6 Related Work

**Programmatic Reinforcement Learning.** Our work is closely related to recent advance on exploring domain-specific programs as an interpretable representation for RL. PIRL [62, 63] and Viper [6] synthesize loop-free, stateless programs, which face limitations in complex robot tasks. Inala et al. [28] improved on this by learning robot controllers as state machines, enabling generalization to tasks with repeating behaviors. These methods rely heavily on strong supervision from oracles like pretrained RL controllers. The tasks they can solve are thus bounded by the capability of the oracle. In contrast, program synthesis methods such as PROLEX [44] and Tabula [46] learn robot control programs from task demonstrations. They generalize these demonstrations into regex-based sketches or Mealy automata to bootstrap synthesis. They can synthesize programs with control flow structures including loops and conditionals, allowing generalization from a specific

sequence of actions to a general structure to solve unseen tasks. Their DSLs feature extensive library functions for manipulating various objects and teleporting robots to different locations. LEAPS [60], PRL [47], and ReGuS [14] eliminate the need for pretrained oracles and synthesizes robot-control programs directly from reward signals. They demonstrate that utilizing rich control-flow constructs (state-conditioned loops and procedure calls) can effectively tackle long-horizon and sparse-reward tasks, which are beyond the capabilities of standard deep RL baselines. However, these existing works rely on a manually designed library of state abstraction predicates and abstract actions to bootstrap synthesis. RoboScribe addresses the primary challenge of automating the construction of robot state and action abstractions.

**Learning State and Action Abstraction.** RoboScribe shares similarities with generalized planning methods like [56, 59], which derive looped plans for solving unbounded problem instances. However, planning techniques require a provided state transition model for each robot action within the abstract state space. Component-based synthesis techniques with user-defined predicates, as in [11, 22], have similar requirements. Automatically learning state and action abstractions has been explored in task and motion planning for robot control [7, 23, 29, 43]. Existing techniques typically learn either predicates from demonstrations assuming low-level controllers are given [15, 36, 55] or learn controllers from demonstrations assuming known predicates [1, 17, 54]. RoboScribe simultaneously learns state and action abstractions, removing such assumptions. There exist library learning techniques [9, 10, 21] that use *syntax abstraction* to extract common structures from a program corpus as reusable library functions. In contrast, RoboScribe performs *state abstraction*. RoboScribe is broadly related to hierarchical RL and planning for robot learning [4, 30, 31, 40, 42, 45, 48, 56, 64]. However, such techniques often struggle with long-horizon tasks with sparse rewards. Traditional controller synthesis algorithms, especially those using formal methods and temporal logic, rely on automata-based approaches involving abstraction and discretization of continuous state and action spaces [16, 51]. These methods face limitations in high-dimensional systems, where discretization can lead to issues like state explosion.

**Reward-guided Program Synthesis.** Existing algorithms often design dense rewards to guide program search directions. For example, PROBE [5] and SYNTIA [8] evaluate programs using input-output examples, generating rich rewards based on output similarity. In [11], the feedback from a deduction engine on the feasibility of partial programs is used to reward the synthesis algorithm, guiding it toward promising search directions. FAERY [12] employs Monte Carlo estimation to sample user queries for additional examples. However, in sparse-reward scenarios, synthesizing complete programs with complex control flow through Monte Carlo methods is challenging due to the low probability of discovering programs with nonzero rewards. RoboScribe addresses this by using comparative abstraction refinement to learn state abstraction predicates that capture subgoal conditions, effectively breaking down the learning process.

## 7 Conclusion

This paper introduces RoboScribe, a program synthesis framework guided by abstraction refinement to address long-horizon, multi-object tasks in robotics. RoboScribe alternates between comparative abstraction refinement and iterative program learning, using demonstrations and execution trajectories from synthesized programs to iteratively refine environment abstractions until a task-solving program can be generated. It identifies recurring subroutines from raw, continuous state-action spaces without predefined abstractions. Experimental results show that RoboScribe generalizes effectively to long-horizon tasks with varying object counts, outperforming baseline methods in interpretability and efficiency. Currently, our language restricts predicates to use norm and arctan functions. It remains an open question whether this expressiveness is sufficient for all tasks—for example, whether additional trigonometric functions are needed, which we leave for future work.

## Data-Availability Statement

The artifact for this work is an implementation of RoboScribe, a program synthesis framework for robotic control tasks. The artifact includes both the codebase for synthesizing robot-control programs and associated datasets used for training and evaluation. The artifact is available at [13].

## Acknowledgments

We would like to thank the anonymous reviewers for their help and feedback on this paper. This material is based upon work supported by the National Science Foundation under grant numbers CCF-2124155 and CCF-2007799.

## References

- [1] Diego Aineto, Sergio Jiménez, and Eva Onaindia. 2022. A Comprehensive Framework for Learning Declarative Action Models. *J. Artif. Intell. Res.* 74 (2022), 1091–1123. <https://doi.org/10.1613/JAIR.1.13073>
- [2] David Andre and Stuart J. Russell. 2002. State Abstraction for Programmable Reinforcement Learning Agents. In *Proceedings of the Eighteenth National Conference on Artificial Intelligence and Fourteenth Conference on Innovative Applications of Artificial Intelligence*.
- [3] Marcin Andrychowicz, Dwight Crow, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. 2017. Hindsight Experience Replay. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*. 5048–5058.
- [4] Pierre-Luc Bacon, Jean Harb, and Doina Precup. 2017. The Option-Critic Architecture. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*.
- [5] Shraddha Barke, Hila Peleg, and Nadia Polikarpova. 2020. Just-in-time learning for bottom-up enumerative synthesis. *Proc. ACM Program. Lang.* OOPSLA (2020).
- [6] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. 2018. Verifiable Reinforcement Learning via Policy Extraction. In *Advances in Neural Information Processing Systems, NeurIPS 2018*.
- [7] Pascal Bercher, Ron Alford, and Daniel Höller. 2019. A Survey on Hierarchical Planning - One Abstract Idea, Many Concrete Realizations. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019, Sarit Kraus (Ed.)*. ijcai.org, 6267–6275. <https://doi.org/10.24963/IJCAI.2019/875>
- [8] Tim Blazytko, Moritz Contag, Cornelius Aschermann, and Thorsten Holz. 2017. Syntia: Synthesizing the Semantics of Obfuscated Code. In *26th USENIX Security Symposium, USENIX Security 2017*.
- [9] Matthew Bowers, Theo X. Olausson, Lionel Wong, Gabriel Grand, Joshua B. Tenenbaum, Kevin Ellis, and Armando Solar-Lezama. 2023. Top-Down Synthesis for Library Learning. *Proc. ACM Program. Lang.* 7, POPL (2023), 1182–1213. <https://doi.org/10.1145/3571234>
- [10] David Cao, Rose Kunkel, Chandrakana Nandi, Max Willsey, Zachary Tatlock, and Nadia Polikarpova. 2023. babble: Learning Better Abstractions with E-Graphs and Anti-unification. *Proc. ACM Program. Lang.* 7, POPL (2023), 396–424. <https://doi.org/10.1145/3571207>
- [11] Yanju Chen, Chenglong Wang, Osbert Bastani, Isil Dillig, and Yu Feng. 2020. Program Synthesis Using Deduction-Guided Reinforcement Learning. In *Computer Aided Verification - 32nd International Conference, CAV 2020*.
- [12] Yanju Chen, Chenglong Wang, Xinyu Wang, Osbert Bastani, and Yu Feng. 2023. Fast and Reliable Program Synthesis via User Interaction. In *38th IEEE/ACM International Conference on Automated Software Engineering, ASE 2023*.
- [13] Guofeng Cui, Yuning Wang, Wensen Mao, Yuanlin Duan, and He Zhu. 2025. Abstraction Refinement-Guided Program Synthesis for Robot Learning from Demonstrations (Artifact). Zenodo. <https://doi.org/10.5281/zenodo.16929200>
- [14] Guofeng Cui, Yuning Wang, Wenjie Qiu, and He Zhu. 2024. Reward-Guided Synthesis of Intelligent Agents with Control Structures. *Proc. ACM Program. Lang.* 8, PLDI (2024), 1730–1754. <https://doi.org/10.1145/3656447>
- [15] Aidan Curtis, Tom Silver, Joshua B. Tenenbaum, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. 2022. Discovering State and Action Abstractions for Generalized Task and Motion Planning. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*. AAAI Press, 5377–5384. <https://doi.org/10.1609/AAAI.V36I5.20475>
- [16] Mehdi Dadvar, Rashmeet Kaur Nayyar, and Siddharth Srivastava. 2023. Conditional Abstraction Trees for Sample-Efficient Reinforcement Learning. In *The 39th Conference on Uncertainty in Artificial Intelligence*. <https://openreview.net/forum?id=tQP094M0j8G>
- [17] Joaquim Ortiz de Haro, Jung-Su Ha, Danny Driess, and Marc Toussaint. 2021. Structured deep generative models for sampling on constraint manifolds in sequential manipulation. In *Conference on Robot Learning, 8-11 November 2021*,

- London, UK (*Proceedings of Machine Learning Research*, Vol. 164), Aleksandra Faust, David Hsu, and Gerhard Neumann (Eds.). PMLR, 213–223. <https://proceedings.mlr.press/v164/ortiz-haro22a.html>
- [18] Quentin Delfosse, Wolfgang Stammer, Thomas Rothenbacher, Dwarak Vittal, and Kristian Kersting. 2023. Boosting Object Representation Learning via Motion and Object Continuity. In *Machine Learning and Knowledge Discovery in Databases: Research Track - European Conference, ECML PKDD 2023*.
- [19] Thomas G. Dietterich. 2000. Hierarchical Reinforcement Learning with the MAXQ Value Function Decomposition. *J. Artif. Intell. Res.* (2000).
- [20] Yiming Ding, Carlos Florensa, Pieter Abbeel, and Mariano Phielipp. 2019. Goal-conditioned Imitation Learning. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 15298–15309. <https://proceedings.neurips.cc/paper/2019/hash/c8d3a760ebab631565f8509d84b3b3f1-Abstract.html>
- [21] Kevin Ellis, Catherine Wong, Maxwell I. Nye, Mathias Sablé-Meyer, Lucas Morales, Luke B. Hewitt, Luc Cary, Armando Solar-Lezama, and Joshua B. Tenenbaum. [n. d.]. DreamCoder: bootstrapping inductive program synthesis with wake-sleep library learning. In *42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2021*. 835–850.
- [22] Yu Feng, Ruben Martins, Jacob Van Geffen, Isil Dillig, and Swarat Chaudhuri. 2017. Component-based synthesis of table consolidation and transformation tasks from examples. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017*.
- [23] Caelan Reed Garrett, Rohan Chitnis, Rachel M. Holladay, Beomjoon Kim, Tom Silver, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. 2021. Integrated Task and Motion Planning. *Annu. Rev. Control. Robotics Auton. Syst.* 4 (2021), 265–293. <https://doi.org/10.1146/ANNUREV-CONTROL-091420-084139>
- [24] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*. PMLR, 1861–1870.
- [25] Jonathan Ho and Stefano Ermon. 2016. Generative Adversarial Imitation Learning. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (Eds.). 4565–4573. <https://proceedings.neurips.cc/paper/2016/hash/cc7e2b878868cbac992d1fb743995d8f-Abstract.html>
- [26] Jarrett Holtz, Arjun Guha, and Joydeep Biswas. 2020. Robot Action Selection Learning via Layered Dimension Informed Program Synthesis. In *4th Conference on Robot Learning, CoRL 2020, 16-18 November 2020, Virtual Event / Cambridge, MA, USA (Proceedings of Machine Learning Research, Vol. 155)*, Jens Kober, Fabio Ramos, and Claire J. Tomlin (Eds.). PMLR, 1471–1480. <https://proceedings.mlr.press/v155/holtz21a.html>
- [27] Baichuan Huang, Abdeslam Boularias, and Jingjin Yu. 2022. Parallel Monte Carlo Tree Search with Batched Rigid-body Simulations for Speeding up Long-Horizon Episodic Robot Planning. In *IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2022, Kyoto, Japan, October 23-27, 2022*. IEEE, 1153–1160. <https://doi.org/10.1109/IROS47612.2022.9981962>
- [28] Jeevana Priya Inala, Osbert Bastani, Zenna Tavares, and Armando Solar-Lezama. 2020. Synthesizing Programmatic Policies that Inductively Generalize. In *8th International Conference on Learning Representations, ICLR 2020*.
- [29] George Dimitri Konidaris, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. 2018. From Skills to Symbols: Learning Symbolic Representations for Abstract High-Level Planning. *J. Artif. Intell. Res.* 61 (2018), 215–289. <https://doi.org/10.1613/JAIR.5575>
- [30] Youngwoon Lee, Shao-Hua Sun, Sriram Somasundaram, Edward S. Hu, and Joseph J. Lim. 2019. Composing Complex Skills by Learning Transition Policies. In *7th International Conference on Learning Representations, ICLR 2019*.
- [31] Youngwoon Lee, Jingyun Yang, and Joseph J. Lim. 2020. Learning to Coordinate Manipulation Skills via Skill Behavior Diversification. In *8th International Conference on Learning Representations, ICLR 2020*.
- [32] Richard Li, Allan Jabri, Trevor Darrell, and Pulkit Agrawal. 2020. Towards Practical Multi-Object Manipulation using Relational Reinforcement Learning. In *2020 IEEE International Conference on Robotics and Automation, ICRA 2020, Paris, France, May 31 - August 31, 2020*. IEEE, 4051–4058. <https://doi.org/10.1109/ICRA40945.2020.9197468>
- [33] Junchi Liang and Abdeslam Boularias. 2023. Learning Category-Level Manipulation Tasks from Point Clouds with Dynamic Graph CNNs. In *Proceedings of the 2023 International Conference on Robotics and Automation (ICRA)*.
- [34] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. 2023. Code as Policies: Language Model Programs for Embodied Control. In *IEEE International Conference on Robotics and Automation, ICRA 2023, London, UK, May 29 - June 2, 2023*. IEEE, 9493–9500. <https://doi.org/10.1109/ICRA48891.2023.10160591>
- [35] Zhixuan Lin, Yi-Fu Wu, Skand Vishwanath Peri, Weihao Sun, Gautam Singh, Fei Deng, Jindong Jiang, and Sungjin Ahn. 2020. SPACE: Unsupervised Object-Oriented Scene Representation via Spatial Attention and Decomposition. In *8th International Conference on Learning Representations, ICLR 2020*.

- [36] João Loula, Kelsey R. Allen, Tom Silver, and Josh Tenenbaum. 2020. Learning constraint-based planning models from demonstrations. In *IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2020, Las Vegas, NV, USA, October 24, 2020 - January 24, 2021*. IEEE, 5410–5416. <https://doi.org/10.1109/IROS45743.2020.9341535>
- [37] Chaitanya Mitash, Kostas E. Bekris, and Abdeslam Boularias. 2017. A self-supervised learning system for object detection using physics simulation and multi-view pose estimation. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2017*.
- [38] Chaitanya Mitash, Abdeslam Boularias, and Kostas E. Bekris. 2018. Robust 6D Object Pose Estimation with Stochastic Congruent Sets. In *British Machine Vision Conference 2018, BMVC 2018*.
- [39] Tongzhou Mu, Zhan Ling, Fanbo Xiang, Derek Yang, Xuanlin Li, Stone Tao, Zhiao Huang, Zhiwei Jia, and Hao Su. 2021. ManiSkill: Generalizable Manipulation Skill Benchmark with Large-Scale Demonstrations. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual*, Joaquin Vanschoren and Sai-Kit Yeung (Eds.). <https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/eda80a3d5b344bc40f3bc04f65b7a357-Abstract-round2.html>
- [40] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. 2018. Data-Efficient Hierarchical Reinforcement Learning. In *Annual Conference on Neural Information Processing Systems, NeurIPS 2018*.
- [41] Ashvin Nair, Bob McGrew, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. 2017. Overcoming Exploration in Reinforcement Learning with Demonstrations. *CoRR* abs/1709.10089 (2017). arXiv:1709.10089 <http://arxiv.org/abs/1709.10089>
- [42] Soroush Nasiriany, Vitchyr Pong, Steven Lin, and Sergey Levine. 2019. Planning with Goal-Conditioned Policies. In *Annual Conference on Neural Information Processing Systems, NeurIPS 2019*.
- [43] Hanna M. Pasula, Luke S. Zettlemoyer, and Leslie Pack Kaelbling. 2007. Learning Symbolic Models of Stochastic Domains. *J. Artif. Intell. Res.* 29 (2007), 309–352. <https://doi.org/10.1613/JAIR.2113>
- [44] Noah Patton, Kia Rahmani, Meghana Missula, Joydeep Biswas, and Isil Dillig. 2024. Programming-by-Demonstration for Long-Horizon Robot Tasks. *Proc. ACM Program. Lang.* 8, POPL (2024), 512–545. <https://doi.org/10.1145/3632860>
- [45] Xue Bin Peng, Michael Chang, Grace Zhang, Pieter Abbeel, and Sergey Levine. 2019. MCP: Learning Composable Hierarchical Control with Multiplicative Compositional Policies. In *Annual Conference on Neural Information Processing Systems, NeurIPS 2019*.
- [46] David Porfiro, Laura Stegner, Maya Cakmak, Allison Saupé, Aws Albarghouthi, and Bilge Mutlu. 2023. Sketching Robot Programs On the Fly. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot Interaction, HRI 2023, Stockholm, Sweden, March 13-16, 2023*, Ginevra Castellano, Laurel D. Riek, Maya Cakmak, and Iolanda Leite (Eds.). ACM, 584–593. <https://doi.org/10.1145/3568162.3576991>
- [47] Wenjie Qiu and He Zhu. 2022. Programmatic Reinforcement Learning without Oracles. In *10th International Conference on Learning Representations, ICLR 2022*.
- [48] Ahmed Hussain Qureshi, Jacob J. Johnson, Yuzhe Qin, Taylor Henderson, Byron Boots, and Michael C. Yip. 2020. Composing Task-Agnostic Policies with Deep Reinforcement Learning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.
- [49] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. 2021. Stable-Baselines3: Reliable Reinforcement Learning Implementations. *Journal of Machine Learning Research* 22, 268 (2021), 1–8. <http://jmlr.org/papers/v22/20-1364.html>
- [50] Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. 2016. You Only Look Once: Unified, Real-Time Object Detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016*. 779–788.
- [51] Wei Ren, Raphaël M. Jungers, and Dimos V. Dimarogonas. 2024. Zonotope-Based Symbolic Controller Synthesis for Linear Temporal Logic Specifications. *IEEE Trans. Autom. Control.* 69, 11 (2024), 7630–7645. <https://doi.org/10.1109/TAC.2024.3394313>
- [52] Rahul Shome, Wei N. Tang, Changkyu Song, Chaitanya Mitash, Hristiyan Kourtev, Jingjin Yu, Abdeslam Boularias, and Kostas E. Bekris. 2019. Towards Robust Product Packing with a Minimalistic End-Effector. In *International Conference on Robotics and Automation, ICRA 2019, Montreal, QC, Canada, May 20-24, 2019*. IEEE, 9007–9013. <https://doi.org/10.1109/ICRA.2019.8793966>
- [53] Tom Silver, Kelsey R. Allen, Alex K. Lew, Leslie Pack Kaelbling, and Josh Tenenbaum. 2020. Few-Shot Bayesian Imitation Learning with Logical Program Policies. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020*.
- [54] Tom Silver, Ashay Athalye, Joshua B. Tenenbaum, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. 2022. Learning Neuro-Symbolic Skills for Bilevel Planning. In *Conference on Robot Learning, CoRL 2022, 14-18 December 2022, Auckland, New Zealand (Proceedings of Machine Learning Research, Vol. 205)*, Karen Liu, Dana Kulic, and Jeffrey Ichnowski (Eds.). PMLR, 701–714. <https://proceedings.mlr.press/v205/silver23a.html>
- [55] Tom Silver, Rohan Chitnis, Nishanth Kumar, Willie McClinton, Tomás Lozano-Pérez, Leslie Pack Kaelbling, and Joshua B. Tenenbaum. 2023. Predicate Invention for Bilevel Planning. In *Thirty-Seventh AAAI Conference on Artificial Intelligence*,



- AAAI 2023, *Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence*, IAAI 2023, *Thirteenth Symposium on Educational Advances in Artificial Intelligence*, EAAI 2023, Washington, DC, USA, February 7-14, 2023, Brian Williams, Yiling Chen, and Jennifer Neville (Eds.). AAAI Press, 12120–12129. <https://doi.org/10.1609/AAAI.V37I10.26429>
- [56] Tom Silver, Rohan Chitnis, Joshua B. Tenenbaum, Leslie Pack Kaelbling, and Tomás Lozano-Pérez. 2021. Learning Symbolic Operators for Task and Motion Planning. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, IROS 2021.
- [57] Avishai Sintov, Andrew S. Morgan, Andrew Kimmel, Aaron M. Dollar, Kostas E. Bekris, and Abdeslam Boularias. 2019. Learning a State Transition Model of an Underactuated Adaptive Hand. *IEEE Robotics Autom. Lett.* 4, 2 (2019), 1287–1294. <https://doi.org/10.1109/LRA.2019.2894875>
- [58] Changkyu Song and Abdeslam Boularias. 2019. Object Rearrangement with Nested Nonprehensile Manipulation Actions. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2019, Macau, SAR, China, November 3-8, 2019*. IEEE, 6578–6585. <https://doi.org/10.1109/IROS40897.2019.8967548>
- [59] Siddharth Srivastava, Neil Immerman, and Shlomo Zilberstein. 2011. A new representation and associated algorithms for generalized planning. *Artif. Intell.* (2011).
- [60] Dweep Trivedi, Jesse Zhang, Shao-Hua Sun, and Joseph J Lim. 2021. Learning to Synthesize Programs as Interpretable and Generalizable Policies. In *Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021*.
- [61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. *CoRR* abs/1706.03762 (2017). arXiv:1706.03762 <http://arxiv.org/abs/1706.03762>
- [62] Abhinav Verma, Hoang Minh Le, Yisong Yue, and Swarat Chaudhuri. 2019. Imitation-Projected Programmatic Reinforcement Learning. In *Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019*.
- [63] Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri. 2018. Programmatically Interpretable Reinforcement Learning. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018*.
- [64] Alexander Vezhnevets, Volodymyr Mnih, Simon Osindero, Alex Graves, Oriol Vinyals, John P. Agapiou, and Koray Kavukcuoglu. 2016. Strategic Attentive Writer for Learning Macro-Actions. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*.
- [65] Jimmy Xin, Linus Zheng, Kia Rahmani, Jiayi Wei, Jarrett Holtz, Isil Dillig, and Joydeep Biswas. 2024. Programmatic Imitation Learning From Unlabeled and Noisy Demonstrations. *IEEE Robotics Autom. Lett.* 9, 6 (2024), 4894–4901. <https://doi.org/10.1109/LRA.2024.3385691>
- [66] Yichen Yang, Jeevana Priya Inala, Osbert Bastani, Yewen Pu, Armando Solar-Lezama, and Martin Rinard. 2021. Program Synthesis Guided Reinforcement Learning for Partially Observed Environments. In *Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021*.
- [67] Allan Zhou, Vikash Kumar, Chelsea Finn, and Aravind Rajeswaran. 2024. Policy Architectures for Compositional Generalization in Control. In *Proceedings of the 1st Reinforcement Learning Conference, RLC 2024, Amherst, MA, USA, August 9-12, 2024*, Philip S. Thomas, Feryal M. P. Behbahani, Glen Berseth, Scott M. Jordan, Scott Niekum, Andrew Patterson, Eugene Vinitsky, Adam White, Martha White, and Amy Zhang (Eds.). University of Massachusetts Amherst, MA, USA, 2264–2283.
- [68] Shaojun Zhu, Andrew Kimmel, Kostas E. Bekris, and Abdeslam Boularias. 2018. Fast Model Identification via Physics Engines for Data-Efficient Policy Search. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, Jérôme Lang (Ed.). ijcai.org, 3249–3256. <https://doi.org/10.24963/ijcai.2018/451>

Received 2025-03-26; accepted 2025-08-12