

LEARNING TO COORDINATE MANIPULATION SKILLS VIA SKILL BEHAVIOR DIVERSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

When mastering a complex manipulation task, humans often decompose the task into sub-skills of their body parts, practice the sub-skills independently, and then execute the sub-skills together. Similarly, a robot with multiple end-effectors can perform a complex task by coordinating sub-skills of each end-effector. To realize temporal and behavioral coordination of skills, we propose a hierarchical framework that first individually trains sub-skills of each end-effector with skill behavior diversification, and learns to coordinate end-effectors using diverse behaviors of the skills. We demonstrate that our proposed framework is able to efficiently learn sub-skills with diverse behaviors and coordinate them to solve challenging collaborative control tasks such as picking up a long bar, placing a block inside a container while pushing the container with two robot arms, and pushing a box with two ant agents.

1 INTRODUCTION

Imagine you wish to play Chopin’s Fantaisie Impromptu on the piano. With little prior knowledge about the piece, you would first practice playing the piece with each hand separately. After independently mastering the left hand part and the right hand part, you would move on to practicing with both hands simultaneously, trying to coordinate between the movement of the two hands to together create a complete piece of music. Coordinating movements of two hands requires adjusting one-handed playing of each hand to the synchronized and coordinated movements of two hands, which requires you to learn variations in playing the same melody. Through the decomposition of skills into sub-skills of two hands and learning variations of sub-skills, humans make the learning process of manipulation skills much faster than learning everything at once.

Can autonomous agents efficiently learn complicated tasks with coordination of different skills from multiple end-effectors like humans? Learning to perform collaborative and composite tasks from scratch requires a huge amount of environment interaction and extensive reward engineering, which often results in undesired behaviors (Riedmiller et al., 2018). Hence, instead of learning a task at once, modular approaches (Andreas et al., 2017; Oh et al., 2017; Frans et al., 2018; Peng et al., 2019; Goyal et al., 2019) suggest to learn reusable primitive skills and solve more complex tasks by recombining the skills. However, all these approaches are focused on working with a single end-effector or agent with learned primitive skills, and learning to coordinate has not been addressed.

To this end, we propose a hierarchical framework that learns to coordinate multiple end-effectors with their primitive skills for various robotics tasks, such as bi-manual manipulation. The main challenge is that naive simultaneous execution of primitive skills from multiple end-effectors can often cause unintended behaviors (e.g. collisions between end-effectors). Thus, our model needs to learn to “correctly” coordinate each end-effector; and thus needs a way to obtain, represent, and control detailed behaviors of each primitive skill. To address these intuitions, our method consists of two parts: (1) acquiring primitive skills with diverse behaviors by mutual information maximization, and (2) learning a meta policy that selects a skill for each end-effector and coordinates the chosen skills by controlling the behavior of each skill.

The main contributions of this paper is a hierarchical approach that tackles cooperative manipulation tasks with multiple end-effectors by (1) learning primitive skills of each end-effector independently with skill behavior diversification and (2) coordinating end-effectors using diverse behaviors of the skills. Our empirical results indicate that our proposed method is able to efficiently learn primitive

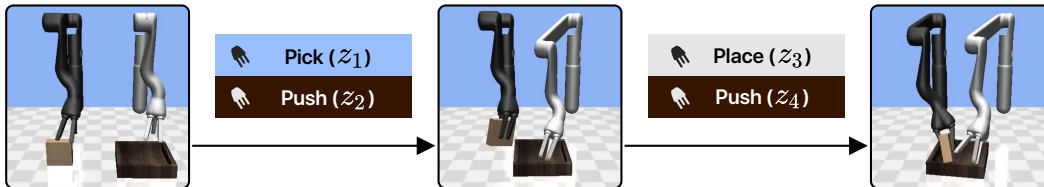


Figure 1: Composing complex skills using multiple agents’ primitive skills requires proper coordination between agents since concurrent execution of primitive skills requires temporal or behavioral coordination. For example, to move a block into a container on the other end of the table, the agent needs to not only utilize pick, place, and push primitive skills at the right time, but also select the appropriate behaviors for these skills, represented as latent vectors z_1 , z_2 , z_3 , and z_4 above. Naive methods neglecting either temporal or behavioral perspective of the coordinated process will produce unintended behaviors, such as collisions between end-effectors.

skills with diverse behaviors and coordinate these skills to solve challenging collaborative control tasks such as picking up a long bar, placing a block inside the container on the right side, and pushing a box with two ant agents. We will make the training code and environments publicly available.

2 RELATED WORK

Hierarchical reinforcement learning: Deep reinforcement learning for continuous control is an active research area. However, learning a complex task either from a sparse reward or a heavily engineered reward becomes computationally impractical as the target task becomes more complicated. Instead of learning from scratch, complex tasks can be tackled by decomposing the tasks into easier and reusable sub-tasks. Hierarchical reinforcement learning temporally splits a task into a sequence of temporally extended meta actions. It often consists of one high-level policy and a set of low-level policies, such as option framework (Sutton et al., 1999). The high-level policy decides which low-level controller to activate and the chosen low-level controller generates an action sequence until the high-level policy switches it to another low-level controller. The option can be discovered unsupervisedly (Schmidhuber, 1990; Levy et al., 2017; Bacon et al., 2017; Nachum et al., 2018), meta-learned (Frans et al., 2018), pre-defined (Kulkarni et al., 2016; Oh et al., 2017; Lee et al., 2019), or attained from additional supervision signals (Andreas et al., 2017). However, these option frameworks are not flexible to solve a task that requires simultaneous activation or interpolation of multiple skills since only one skill can be activated at each time step.

Composition of multiple policies: To solve compositional tasks multiple policies can be simultaneously activated by adding Q-functions (Haarnoja et al., 2018a), additive composition (Qureshi et al., 2019; Goyal et al., 2019), or multiplicative composition (Peng et al., 2019). As each of policies takes the whole observation as input and controls the whole agent, it is not robust to changes in unrelated part of observation. Hence, these skill composition approaches can fail when an agent encounters a new combination of skills or a new skill is introduced since a policy cannot achieve its own goal under unseen circumstances.

Multi-agent reinforcement learning: Instead of having a policy with the full observation and action space, multi-agent reinforcement learning (MARL) proposes to explicitly split the observation and action space according to agents (e.g. robots or end-effectors), which allows efficient low-level policy training as well as flexible skill composition. For cooperative tasks, sharing policy parameters (Gupta et al., 2017) and decentralized actor with centralized critic (Lowe et al., 2017; Foerster et al., 2018) have been actively used. However, these approaches suffer from the credit assignment problem (Sutton, 1984) among agents and the lazy agent problem (Sunehag et al., 2018). As agents have more complicated morphology and larger observation space, learning a policy for a multi-agent system from scratch requires extremely long training time. To resolve these issues, we propose to first train reusable skills for each agent separately, instead of learning primitive skills of multiple agents together. Then, we recombine these skills to complete more complicated tasks with proper coordination between the skills.

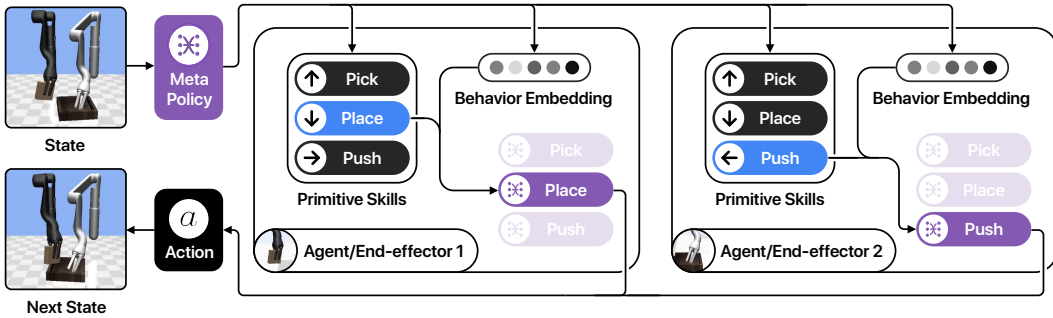


Figure 2: Our method is composed of two components: a meta policy and a set of agent-specific primitive policies relevant to task completion. The meta policy selects which primitive skill to run for each agent as well as the behavior embedding (i.e. variation in behavior) of the chosen primitive skill. Each selected primitive skill takes the agent observation and the behavior embedding as input and outputs action for that agent.

Learning diverse skills: Coordination of multiple skills from multiple agents requires the skills to be flexible; hence, a skill can be adjusted to collaborate with other agents’ skills. Maximum entropy policies (Haarnoja et al., 2017; 2018a;b) can learn diverse ways to achieve a goal by maximizing not only reward but also entropy of the policy. Eysenbach et al. (2019) proposes to discover diverse skills without reward by maximizing entropy as well as mutual information between resulting states and skill embeddings (Hausman et al., 2018). To this end, our method leverages a maximum entropy policy with a discriminability objective (Eysenbach et al., 2019) to learn a primitive skill with diverse behaviors conditioned on a controllable skill embedding, which will be later used as a *behavior embedding* for a high-level policy to adapt low-level policies for coordination.

3 METHOD

In this work, we address the problem of solving cooperative manipulation tasks that require collaboration between multiple end-effectors or agents (we use the terms ”multiple end-effectors” and ”multiple agents” interchangeably in this paper). Instead of learning a multi-agent task from scratch (Lowe et al., 2017; Gupta et al., 2017; Sunehag et al., 2018; Foerster et al., 2018), modular approaches (Andreas et al., 2017; Frans et al., 2018; Peng et al., 2019) suggest to learn reusable primitive skills and solve more complex tasks by recombining these skills. However, concurrent execution of primitive skills of multiple agents fails when agents never experienced a combination of skills during the pre-training stage or skills require temporal or behavioral coordination.

Therefore, we propose a hierarchical framework that learns to coordinate multiple agents with primitive skills to compose a complex skill. Moreover, during primitive skill training, we propose to learn a latent behavior embedding, which provides controllability of the primitive skill to the meta policy for skill coordination. In Section 3.2, we describe our hierarchical framework in details. Next, in Section 3.3, we elaborate how controllable primitive skills can be achieved. Lastly, we describe how the meta policy learns to coordinate primitive skills in Section 3.4.

3.1 PRELIMINARIES

We formulate our problem as a Markov decision process defined by a tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{T}, R, \rho, \gamma\}$ of states, actions, transition probability, reward, initial state distribution, and discount factor. In our formulation, we assume the environment includes N agents. Hence, the state space and action space for an agent i can be represented as \mathcal{S}^i and \mathcal{A}^i where each element of \mathcal{S}^i is a subset of the corresponding element in \mathcal{S} and $\mathcal{A} = \mathcal{A}^1 \times \mathcal{A}^2 \times \dots \times \mathcal{A}^N$, respectively. For each agent i , we provide a set of skills Π^i . An action distribution of an agent i is represented as a policy $\pi_{c_t^i}^i(a_t^i | s_t^i) \in \Pi^i$, where c_t^i is a skill index of the agent, $s_t^i \in \mathcal{S}^i$ is a state, $a_t^i \in \mathcal{A}^i$ is an agent action at time t . An initial state s_0 is randomly sampled from ρ , and then, N agents take actions $a_t^1, a_t^2, \dots, a_t^N$ sampled from a composite policy $\pi(a_t^1, a_t^2, \dots, a_t^N | s_t, c_t^1, c_t^2, \dots, c_t^N)$ and receives a single reward r_t , where $\pi(a_t^1, a_t^2, \dots, a_t^N | s_t, c_t^1, c_t^2, \dots, c_t^N) = (\pi_{c_t^1}^1(a_t^1 | s_t), \pi_{c_t^2}^2(a_t^2 | s_t), \dots, \pi_{c_t^N}^N(a_t^N | s_t))$. The performance

Algorithm 1 ROLLOUT

```

1: Input: Meta policy  $\pi_{meta}$ , sets of primitive policies  $\Pi^1, \dots, \Pi^N$ , and meta horizon  $T_{low}$ .
2: Initialize an episode and receive initial state  $s_0$ .
3:  $t \leftarrow 0$ 
4: while episode is not terminated do
5:   Sample skill indexes and behavior embeddings  $(c_t^1, \dots, c_t^N), (z_t^1, \dots, z_t^N) \sim \pi_{meta}(s_t)$ 
6:    $\tau \leftarrow 0$ 
7:   while  $\tau < T_{low}$  and episode is not terminated do
8:      $a_{t+\tau} = (a_{t+\tau}^1, \dots, a_{t+\tau}^N) \sim (\pi_{c_t^1}^1(s_{t+\tau}, z_t^1), \dots, \pi_{c_t^N}^N(s_{t+\tau}, z_t^N))$ 
9:      $s_{t+\tau+1}, r_{t+\tau} \leftarrow \text{ENV}(s_{t+\tau}, a_{t+\tau})$ 
10:     $\tau \leftarrow \tau + 1$ 
11:   end while
12:   Add a transition  $s_t, (c_t^1, \dots, c_t^N), (z_t^1, \dots, z_t^N), s_{t+\tau}, r_{t:t+\tau-1}$  to the rollout buffer  $\mathcal{B}$ .
13:    $t \leftarrow t + \tau$ 
14: end while

```

of the agent is evaluated based on a discounted return $R = \sum_{t=0}^{T-1} \gamma^t r_t$, where T is the episode horizon.

3.2 HIERARCHICAL FRAMEWORK

As illustrated in Figure 2, our hierarchical model is composed of two components: a meta policy π_{meta} and a set of primitive skills of multiple agents Π^1, \dots, Π^N . The meta policy chooses one primitive skill for each agent $(\pi_{c^1}^1, \pi_{c^2}^2, \dots, \pi_{c^N}^N) \in \Pi^1 \times \Pi^2 \times \dots \times \Pi^N$. To control the chosen primitive skills, the meta policy also outputs a set of latent behavior embeddings (z^1, z^2, \dots, z^N) which parameterize diverse behaviors of the corresponding skills. Then, each agent simultaneously executes the selected skill $\pi_{c^i}^i$ parameterized by the behavior embedding z^i . The behavior embedding enables the meta policy to control primitive skills and it is necessary for coordinating multiple agents' skills. For example, in the task illustrated by Figure 1, naively placing a block in the left hand to the container being moved by the right hand can cause collision between two robot arms. The arms can avoid collision while performing their skills by properly adjusting their end-effector skill behaviors (e.g. placing the block from the left side and pushing the container while leaning the hand toward the right side). A set of skills are selected every T_{low} time steps. Algorithm 1 illustrates the overall rollout process.

We denote the meta policy as $\pi_{meta}(c^1, \dots, c^N, z^1, \dots, z^N | s_t)$, where $c^i \in \mathcal{C}^i$ represents a skill index of an agent $i \in [1, \dots, N]$ and z^i represents a behavior embedding of the skill. The meta policy π_{meta} is modeled as a mixture of N categorical distributions for skill indexes and N Gaussian distributions for behavior embeddings. The observation of the meta policy contains all available information while each agent has limited access to a partial observation. Our meta policy selects a skill to execute for each agent, rather than selecting one primitive skill for the entire multi-agent system to execute. Also, we give the meta policy the capability to select which variant of a skill to execute (see Section 3.4). Once a set of primitive skills $\{\pi_{c^1}^1, \dots, \pi_{c^N}^N\}$ are chosen to be executed, each primitive skill generates an action $a^i \sim \pi_{c^i}^i(a^i | s^i, z^i)$ based on the current state s^i and the latent vector z^i . Note that we will not differentiate state spaces for primitive policies and omit time step t in the rest of the paper due to the simplicity of notations. We want primitive skills to be agent-specific, controlling motion of only one agent. This will not only make the learning of primitive skills efficient but also enable generalization to different collaboration scenarios (see Section 3.3).

3.3 TRAINING PER-AGENT PRIMITIVE POLICIES WITH DIVERSE BEHAVIORS

For a primitive skill to adjust to and collaborate with skills of other agents in a new environment, the skill needs to support variations of skill behaviors when executed at a given state. Moreover, a behavioral variation of a skill should be controllable by a meta policy for skill coordination. In order to make our primitive policies generate diverse behaviors controlled by a latent vector z , we leverage the entropy and mutual information maximization objective introduced in Eysenbach et al. (2019).

More specifically, a primitive policy of an agent i takes as input the current state $s \in S$ and a latent behavior embedding $z \sim p(z)$, and output action $a \in A^i$, where the prior distribution $p(z)$ is Gaussian. Diverse behaviors conditioned on a random sample z can be achieved by maximizing the mutual information between skills and states, $MI(s|z)$, while minimizing the mutual information between skills and actions given the state, $MI(a, z|s)$, together with maximizing the entropy of the policy for the diverse behaviors, $\mathcal{H}(a|s)$. The objective can be written as follows (we refer the readers to Eysenbach et al. (2019) for derivation):

$$\mathcal{F}(\theta) = \mathcal{H}(a | s, z) - \mathcal{H}(z | s) + \mathcal{H}(z) \quad (1)$$

$$= \mathcal{H}(a | s, z) + \mathbb{E}_{z \sim p(z), s \sim \pi(z)} [\log p(z | s)] - \mathbb{E}_{z \sim p(z)} [\log p(z)] \quad (2)$$

$$\geq \mathcal{H}(a | s, z) + \mathbb{E}_{z \sim p(z), s \sim \pi(z)} [\log q_\phi(z | s) - \log p(z)], \quad (3)$$

where the learned discriminator $q_\phi(p|z)$ approximates the posterior $p(z|s)$.

In this paper, we train a primitive policy by maximizing a weighted sum of the reward and Equation (3):

$$r_t + \lambda_1 \mathcal{H}(a | s, z) + \lambda_2 \mathbb{E}_{z \sim p(z), s \sim \pi(z)} [\log q_\phi(z | s) - \log p(z)]. \quad (4)$$

Maximizing Equation (3) encourages multi-modal exploration strategies by learning policies represented by expressive energy-based models while achieving their own goal. Moreover, by maximizing identifiability of behaviors, the latent vector z , named *behavior embedding*, can represent a variation of the learned policy and thus can be used to control the behavior of the policy. For example, when training a robot to move an object, a policy learns to move the object quickly as well as slowly, and these diverse behaviors map to different latent vectors z . We empirically show that the policies with diverse behaviors achieve better compositionality with other agents in our experiments.

3.4 COMPOSING PRIMITIVE SKILLS WITH META POLICY

We design the meta policy $\pi_{meta} : S \rightarrow \prod_{i=1}^N (\mathcal{C}^i \times Z^i)$ to take as input the current state containing information about the environment as well as the states of all agents and output (1) the primitive skill, denoted as c^i , to execute for each agent; and (2) the behavior embedding vector z^i to be used as input for the primitive skill c^i . Since there are a finite number of skills for each agent to execute, the meta action space for each agent \mathcal{C}^i is discrete, while the behavior embedding space for each agent Z^i is continuous. Thus, the meta policy can be implemented using a $|\mathcal{C}^i| + N_z$ -head neural network which represents a mixture of $|\mathcal{C}^i|$ -way categorical distributions for skill selection and N_z Gaussian distributions for behavior control of the chosen skill.

3.5 IMPLEMENTATION

We model the primitive policies and posterior distributions q_ϕ as neural networks. We train the primitive policies using soft actor-critic (Haarnoja et al., 2018b). When we train a primitive policy, we use unit Gaussian distribution as the prior distribution of latent variables $p(z)$. Each primitive policy outputs the mean and standard deviation of a Gaussian distribution over an action space. For a primitive policy, we apply tanh activation to normalize the action between $[-1, 1]$. We model the meta policy as neural network with multiple heads that outputs the skill c^i and behavior embedding z^i for each agent. The meta policy is trained using PPO (Schulman et al., 2017; 2016; Dhariwal et al., 2017). All policy networks in this paper consist of 3 fully connected layers of 64 hidden units with ReLU nonlinearities. The discriminator network for primitive training with DIAYN (Eysenbach et al., 2019) is a 2-fully connected layers network with 64 hidden units.

4 EXPERIMENTS

To demonstrate the effectiveness of our framework, we compare our method to prior methods in the field of multi-agent RL and ablate the components of our framework to understand their importance. We conducted our experiments on a set of challenging robot control environments that require coordination of different agents to complete: collaborative robotic manipulation and locomotion.

Through our experiments, we aim to answer the following questions: (1) can our framework efficiently learn to combine primitive skills to execute a complicated task; (2) can our learned agent exhibit

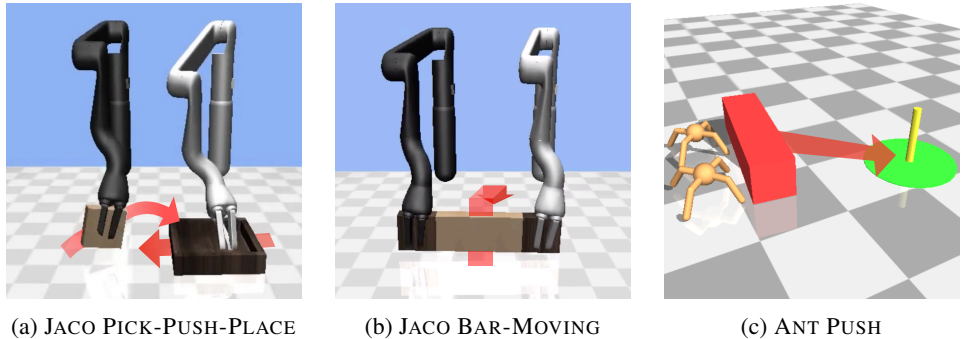


Figure 3: The composite tasks pose a challenging combination of object manipulation and locomotion skills, which requires coordination of multiple agents and temporally extended behaviors. (a) The left Jaco arm needs to pick up a block while the right jaco arm pushes a container. Then, it places the block into the container. (b) Jaco arms are required to pick and place a bar-shaped block together. (c) Two ants push the red box to the goal location (green circle) together.

collaborative behaviors during task execution; (3) can our framework leverage the controllable behavior variations of the primitive skills to achieve better coordination?

For details about environments and training, please refer to the supplementary material. As the performance of training algorithms varies between runs, we train each method on each task with 3 different random seeds and report mean and standard deviation of each method’s success rate.

4.1 BASELINES

We compare the performance of our method with various single- and multi-agent RL methods:

Centralized Policy: a single policy modeled by multi-layer preceptron takes as input the full observation including observation of every agent and outputs actions for every agent in the environment.

Decentralized Policy: a set of n agent policies modeled by multi-layer perceptrons takes the observation of a single agent as input and outputs action for that agent. The global critic is used to train decentralized policies from a single task reward.

Modular Network: a modular framework composed of a meta policy and N sets of primitive skills. The meta policy takes the full observation as input and selects a primitive skill for each agent policy. Each agent policy takes as input the agent observation and outputs action for that agent.

Modular Network with Skill Coordination: our proposed method, which consists of a meta policy and N sets of primitive skills, where each primitive skill is conditioned on a behavior embedding z . The meta policy takes the full observation as input and selects both a primitive skill and a behavior embedding for each agent policy. Each primitive skill takes the agent observation and the behavior embedding as input and outputs action for that agent.

4.2 JACO PICK-PUSH-PLACE

We develop JACO-PICK-PUSH-PLACE and JACO-BAR-MOVING environments using two Kinova Jaco arms, where each Jaco arm is a 9 DoF robotic arm with 3 fingers. A Jaco arm starts with a block on the left and a container on the right. To complete the task, the robot needs to pick up the block, push the container to the center, and place the block inside the container.

Primitives skills. There are three primitive skills available to each arm: *Picking up*, *Pushing*, and *Placing* (see Figure 3a). *Picking up* requires the robotic arm to pick up a small block, which is randomly placed on the table. If the block is not picked up after a certain amount of time or the arm drops the block, the agent fails. *Pushing* learns to push a big container to its opposite side (e.g. from left to the center or from right to center). The agent fails if it cannot place the container to the center. *Placing* is to place an object in the gripper to the table. The agent only succeeds when it stably place the object to the desired location on the table.

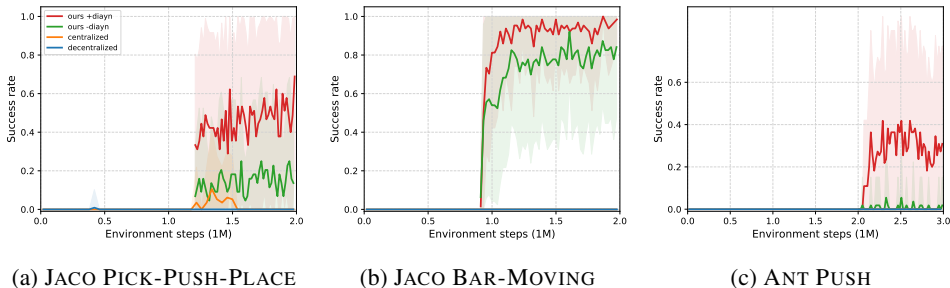


Figure 4: Learning curves (success rates) of our model (green), our model without skill coordination (red), centralized policy (yellow), and decentralized policy (blue). For our models, we shift the learning curves rightwards the total number of environment steps the agent takes to learn the primitive skills (1.2 M, 0.9 M, and 2.0 M, respectively). Our method substantially improves learning speed and performance on JACO PICK-PUSH-PLACE and ANT PUSH. The shaded areas represent the standard deviation of results from three different seeds.

Composite task. Our method can successfully perform JACO PICK-PUSH-PLACE task while all baselines fail to compose primitive skills as shown in Figure 4a. The Centralized Policy baseline cannot efficiently learn the task mainly since the RL agent requires to learn the combinatorial number of skill composition. Table 1 demonstrates that the baseline cannot generalize to unseen combination of skills. On the other hand, the Decentralized Policy baseline fails to learn the task mainly due to the difficulty of credit assignment problem. Since the composite task requires multiple primitive skills to learn for multiple agents, a reward signal about a failure case cannot be assigned to the correct agent or skill. By using pre-trained primitive skills, the credit assignment problem is relaxed and all agents can perform their skills concurrently.

4.3 JACO BAR-MOVING

A Jaco robot needs to pick up a long bar with its two arms, move the bar towards a target location while maintaining its rotation, and place it on the table (see Figure 3b). The initial position of the bar will be randomly initialized every episode and it requires an agent to find appropriate coordination between two arms for each initialization.

Primitives skills. There are three pre-trained primitive skills available to each arm: *Picking up*, *Moving*, and *Placing*. *Picking up* requires the robotic arm to pick up a small block, which is randomly placed on the table. If the block is not picked up after a certain amount of time or the arm drops the block, the agent fails. *Moving* learns to move a block in the hand to the specified location. The agent fails if it drops the block or does not reach the target position within a certain amount of time. *Placing* is to place an object in the gripper to the table. It only succeeds when it stably place the object straight down to the table.

Composite task. Our method can successfully perform JACO BAR-MOVING task while both baselines fail to compose primitive skills as shown in Figure 4b. The modular framework without explicit coordination of skills solve the composite task fairly well but doesn’t achieve as high success rate as the one with explicit coordination of skills because the task requires the two end-effectors to work very closely together. For example, picking up skill of both arms should be synchronized when they start to lift the bar and two arms require to lift the bar while maintaining the relative position between them since they are connected by holding the bar. The modular framework without explicit coordination of skills can synchronize the execution of picking, moving, and placing. But the inability to micro-adjust to the movement of the other arm causes instability to the process of bar picking and moving, which results in the gap in success rate compared to the framework with explicit coordination.

4.4 ANT PUSH

We developed a multi-ant environment, ANT-PUSH, inspired from Nachum et al. (2019), simulated in the Mujoco (Todorov et al., 2012) physics engine and models from OpenAI Gym (Brockman

	Jaco Pick-Push-Place	Jaco Bar-Moving	Ant Push
Centralized Policy	0.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000
Decentralized Policy	0.000 \pm 0.000	0.000 \pm 0.000	0.000 \pm 0.000
Ours w/o Skill Coordination	0.213 \pm 0.410	0.940 \pm 0.237	0.003 \pm 0.058
Ours w/ Skill Coordination	0.807 \pm 0.395	0.997 \pm 0.058	0.323 \pm 0.468

Table 1: Success rates for all tasks, comparing our method against baselines. Each entry in the table represents average success rate and standard deviation over three random seeds and 100 runs. The baselines learning from scratch fail to learn complex tasks with multiple agents.

et al., 2016). In this environment, two ants need to push a large object toward a green target place, collaborating with each other to keep the angle of the object as stable as possible (see Figure 3c).

Primitives skills. We train walking skills of an ant agent in 4 directions: up, down, left, and right. During primitive skill training, a block is randomly placed and pushing the block gives an additional reward to the agent. The learned policies have different speed and trajectories conditioned on the latent behavior embedding.

Composite task. Our method achieves 32.3% success rate on ANT PUSH task while all baselines fail to compose primitive skills as shown in Figure 4c. The poor performance of centralized policy and decentralized policy baselines shows the difficulty of credit assignment between agents, which leads one of the ants moves toward a block and pushes it but another ant does not move. As can be seen in Table 1, the ablated model seldom succeeds while our method with behavior embeddings succeeds in 32.3% of episodes. This result illustrates the importance of coordination of agents, which helps synchronizing and controlling the velocities of both ant agents to push the block toward the goal position while maintaining its rotation.

4.5 EFFECT OF DIVERSITY OF PRIMITIVE SKILLS

To analyze the effect of the diversity of the primitive skills, we compare our model with primitive skills trained with different diversity coefficients $\lambda_2 = [0.0, 0.05, 0.1, 0.5, 1.0]$ in Equation (4). The reward curves are shown in Figure 5. We can observe only with small diversity coefficients $\lambda_2 = 0.05, 0.1$, the agent can control detailed behaviors of primitive skills, while primitive skills without diversity ($\lambda_2 = 0$) cannot be coordinated. The meta policy tries to synchronize two ant agents' positions and velocities by switching primitive skills, but it is not able to achieve proper coordination. On the other hand, large diversity coefficients $\lambda_2 = 0.5, 1.0$ make the primitive skills often focus on demonstrating diverse behaviors and fail to achieve the goals of the skills. Hence, these primitives do not have enough functionality to solve the target task.

5 CONCLUSIONS

In this paper, we propose a hierarchical framework with skill coordination to tackle challenges of compositions of sub-skills. Specifically, we use entropy maximization policies with the mutual information maximization to train controllable primitive skills with diverse behaviors. To coordinate learned primitive skills, the meta policy predicts not only the skill to execute for each end-effector but also the behavior embedding that controls the chosen primitive skill's behavior. Our experimental results on robotic manipulation and locomotion tasks demonstrate that our proposed framework is able to efficiently learn primitive skills with diverse behaviors and coordinate multiple end-effectors to solve challenging cooperative control tasks.

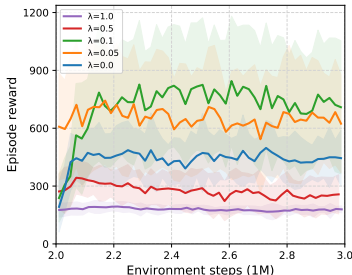


Figure 5: Reward curves of our model with different diversity coefficients on ANT PUSH.

REFERENCES

- Jacob Andreas, Dan Klein, and Sergey Levine. Modular multitask reinforcement learning with policy sketches. In *International Conference on Machine Learning*, pp. 166–175, 2017.
- Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. Openai baselines. <https://github.com/openai/baselines>, 2017.
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SJx63jRqFm>.
- Jakob N Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Kevin Frans, Jonathan Ho, Xi Chen, Pieter Abbeel, and John Schulman. META LEARNING SHARED HIERARCHIES. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=SyX0IeWAW>.
- Anirudh Goyal, Shagun Sodhani, Jonathan Binas, Xue Bin Peng, Sergey Levine, and Yoshua Bengio. Reinforcement learning with competitive ensembles of information-constrained primitives. *arXiv preprint arXiv:1906.10667*, 2019.
- Jayesh K Gupta, Maxim Egorov, and Mykel Kochenderfer. Cooperative multi-agent control using deep reinforcement learning. In *International Conference on Autonomous Agents and Multiagent Systems*, pp. 66–83. Springer, 2017.
- Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with deep energy-based policies. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1352–1361. JMLR. org, 2017.
- Tuomas Haarnoja, Vitchyr Pong, Aurick Zhou, Murtaza Dalal, Pieter Abbeel, and Sergey Levine. Composable deep reinforcement learning for robotic manipulation. In *International Conference on Robotics and Automation (ICRA)*, pp. 6244–6251. IEEE, 2018a.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pp. 1856–1865, 2018b.
- Karol Hausman, Jost Tobias Springenberg, Ziyu Wang, Nicolas Heess, and Martin Riedmiller. Learning an embedding space for transferable robot skills. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rk07ZXZRb>.
- Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Advances in neural information processing systems*, pp. 3675–3683, 2016.
- Youngwoon Lee, Shao-Hua Sun, Sriram Somasundaram, Edward Hu, and Joseph J. Lim. Composing complex skills by learning transition policies. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=rygrBhC5tQ>.
- Andrew Levy, Robert Platt, and Kate Saenko. Hierarchical actor-critic. *arXiv preprint arXiv:1712.00948*, 2017.
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*, pp. 6379–6390, 2017.

- Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 3303–3313, 2018.
- Ofir Nachum, Michael Ahn, Hugo Ponte, Shixiang Gu, and Vikash Kumar. Multi-agent manipulation via locomotion using hierarchical sim2real. *arXiv preprint arXiv:1908.05224*, 2019.
- Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization with multi-task deep reinforcement learning. In *International Conference on Machine Learning*, pp. 2661–2670, 2017.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
- Xue Bin Peng, Michael Chang, Grace Zhang, Pieter Abbeel, and Sergey Levine. Mcp: Learning composable hierarchical control with multiplicative compositional policies. *arXiv preprint arXiv:1905.09808*, 2019.
- Ahmed H Qureshi, Jacob J Johnson, Yuzhe Qin, Byron Boots, and Michael C Yip. Composing ensembles of policies with deep reinforcement learning. *arXiv preprint arXiv:1905.10681*, 2019.
- Martin Riedmiller, Roland Hafner, Thomas Lampe, Michael Neunert, Jonas Degraeve, Tom Van de Wiele, Volodymyr Mnih, Nicolas Heess, and Jost Tobias Springenberg. Learning by playing-solving sparse reward tasks from scratch. *arXiv preprint arXiv:1802.10567*, 2018.
- Jürgen Schmidhuber. *Towards compositional learning with dynamic neural networks*. Inst. für Informatik, 1990.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. 2016.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 2085–2087. International Foundation for Autonomous Agents and Multiagent Systems, 2018.
- Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- Richard Stuart Sutton. Temporal credit assignment in reinforcement learning. *PhD thesis, University of Massachusetts*, 1984.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pp. 5026–5033. IEEE, 2012.

A ENVIRONMENT DETAILS

The details of observation spaces, action spaces, number of agents, and episode lengths are described in Table 2. All units in this section are in meters unless otherwise specified.

	Jaco Pick-Push-Place	Jaco Bar Moving	Ant Push
Observation Space	88	88	100
- Robot observation	62	62	82
- Object observation	26	26	18
Action Space	18	18	16
Number of Agents	2	2	2
Episode length	150	100	200

Table 2: Environment details

A.1 ENVIRONMENT DESCRIPTIONS

In both Jaco environments, the robot works on a table with size (1.6, 1.6) and top center position (0, 0, 0.82). The two Jaco arms are initialized at positions $(-0.16, -0.16, 1.2)$ and $(-0.16, 0.24, 1.2)$. Left arm and right arm objects are initialized around $(0.3, 0.2, 0.86)$ and $(0.3, -0.2, 0.86)$ respectively in all primitive training and composite task training environments, with small random position and rotation perturbation.

In the *Jaco Pick-Push-Place* task, the right jaco arm needs to pick up the object and place it into the container initialized at the other side of the table. Success is defined by contact between the object and the inner top side of the container.

In the *Jaco Bar-Moving* task, the two Jaco arms need to together pick the long bar up by height of 0.7, move it towards the arms by distance of 0.15, and place it back on the table. Success is defined by (1) the bar being placed within 0.04 away from the desired destination both in height and in xy-position and (2) the bar having been picked 0.7 above the table.

In the *Ant Push* task, the two ant agents need to push a big box together to the goal position. The box has a size of $8.0 \times 1.6 \times 1.6$. The distance between ants and the box is 20 cm and the distance between the box and the goal is 30 cm. Initial positions have 1 cm of randomness and the agent has a randomness of 0.01. The task is considered as success when both the distances between left and right end of the box and the goal are within 5 cm.

A.2 REWARD DESIGN

For every task, we add a control penalty, $-0.001 * \|a\|^2$, to regularize the magnitude of actions where a is a torque action performed by an agent.

Jaco Pick: reward for pick primitive is defined by the weighted sum of pick reward, gripper-to-cube distance reward, cube position and quaternion stability reward, hold duration reward, success reward, and robot control reward. More concretely,

$$R(s) = \lambda_{pick} \cdot (z_{box} - z_{init}) + \lambda_{dist} \cdot d(p_{gripper}, p_{box}) + \lambda_{pos} \cdot d(p_{box}, p_{init}) + \lambda_{quat} \cdot \text{abs}(\Delta_{quat}) + \lambda_{hold} \cdot t_{hold} + \lambda_{success} \cdot \mathcal{W}_{success} + \lambda_{ctrl} \|a\|^2,$$

where $\lambda_{pick} = 500$, $\lambda_{dist} = 100$, $\lambda_{pos} = 1000$, $\lambda_{quat} = 1000$, $\lambda_{hold} = 10$, $\lambda_{success} = 100$, $\lambda_{ctrl} = 1e - 4$.

Jaco Place: reward for place primitive is defined by the weighted sum of xy-distance reward, height reward (larger when cube close to floor), success reward, and robot control reward.

Jaco Push: reward for push primitive is defined by the weighted sum of destination distance reward, position and quaternion stability reward, hold duration reward, success reward, and robot control reward.

Jaco Pick-Push-Place: reward for Pick-Push-Place is defined by the weighted sum of reach reward, gripper contact reward, per-staged pick/push/place rewards, success reward, and control reward. We tune the reward carefully for all baselines.

Jaco Bar-Moving: reward for Bar-Moving is defined by the weighted sum of per-stage pick/move/place rewards, success reward, and control reward.

Ant Moving: reward for ANT PUSH is defined by upright, velocity towards the desired direction. We provide a dense reward to encourage the desired locomotion behavior using velocity, stability, and posture, as following:

$$R(s) = \lambda_{vel} \cdot \text{abs}(\Delta x_{ant}) + \lambda_{boxvel} \cdot \text{abs}(\Delta x_{box}) + \lambda_{upright} \cdot \cos(\theta) - \lambda_{height} \cdot \text{abs}(0.6 - h) + \lambda_{goal} \cdot \text{dist}(p_{goal}, p_{box}),$$

, where $\lambda_{vel} = 50$, $\lambda_{boxvel} = 20$, $\lambda_{upright} = 1$, $\lambda_{height} = 0.5$. For ANT PUSH, we provide an additional reward based on distance between the box and the goal position with $\lambda_{goal} = 200$.

B EXPERIMENT DETAILS

We use PyTorch (Paszke et al., 2017) for our implementation and all experiments are conducted on a workstation with Intel Xeon Gold 6154 CPU and 4 NVIDIA GeForce RTX 2080 Ti GPUs.

B.1 HYPERPARAMETERS

Parameters	Value
learning rate	3e-4
gradient steps	50
batch size	256
discount factor	0.99
target smoothing coefficient	0.005
reward scale (SAC)	1.0
experience buffer size	1000
T_{low}	1 for JACO, 5 for ANT

Table 3: Hyperparameters