

# EFFICIENT CROSS-EPIISODE META-RL

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We introduce Efficient Cross-Episodic Transformers (ECET), a new algorithm for online Meta-Reinforcement Learning that addresses the challenge of enabling reinforcement learning agents to perform effectively in previously unseen tasks. We demonstrate how past episodes serve as a rich source of in-context information, which our model effectively distills and applies to new contexts. Our learned algorithm is capable of outperforming the previous state-of-the-art and provides more efficient meta-training while significantly improving generalization capabilities. Experimental results, obtained across various simulated tasks of the MuJoCo, Meta-World and ManiSkill benchmarks, indicate a significant improvement in learning efficiency and adaptability compared to the state-of-the-art. Our approach enhances the agent’s ability to generalize from limited data and paves the way for more robust and versatile AI systems.

## 1 INTRODUCTION

Reinforcement learning (RL) (Sutton & Barto, 2018) has made rapid progress in recent years and has shown success in complex real-world benchmarks (Bellemare et al., 2020; Degraeve et al., 2022; Kaufmann et al., 2023) that go beyond typical (video) game playing benchmarks (Mnih et al., 2015a; Schrittwieser et al., 2020). Despite these high-profile success stories, RL has not yet found widespread adoption in many real-world applications. This is partly due to the high computational demands (Cobbe et al., 2020; Shala et al., 2022; 2023a) and the brittleness of RL algorithms (Henderson et al., 2018; Engstrom et al., 2020; Andrychowicz et al., 2021), which is further amplified by the difficulty of optimizing the hyperparameters of RL algorithms (Parker-Holder et al., 2022; Eimer et al., 2023; Mohan et al., 2023). More importantly, RL training pipelines are generally not built with generalization in mind (Kirk et al., 2023).

The lack of generalization capabilities is a crucial factor that limits the applicability of RL in domains where there are no perfect simulators. The classical RL training pipeline involves training, validating, and testing in the exact same task. However, in reality, an RL agent’s sensors might slightly drift over time or tasks might exhibit drastically different features during different stages of an episode (such as, e.g., times of the day). If the training task was not set up with such changes in mind, an RL policy will not have the opportunity to learn generalizable behavior. Consequently, small variations in the task that might not have been observed during training can lead to severe failure cases. For example, a policy that was trained to steer a robot in a highly controlled lab environment might fail the instant it is supposed to leave the familiar lab environment.

Meta-learning (Vanschoren, 2019; Hospedales et al., 2022) is a subfield of machine learning that aims to take advantage of previous experience as efficiently as possible so that learned behavior can be transferred to similar problems. One promising approach within meta-learning is to include meta-features directly in the training procedure (Benjamins et al., 2023; Beukman et al., 2023; Shala et al., 2023b; Arango et al., 2024), enabling policies capable of zero-shot generalization based on observed meta-features. These meta-features are either hand-crafted (Benjamins et al., 2023; Beukman et al., 2023; Arango et al., 2024) or learned (Shala et al., 2023b). In the Meta-RL setting, it might not always be possible to have access to meaningful hand-crafted meta-features to learn generalizable policies. That is why, most work on meta-features for Meta-RL focuses on learned meta-features (Duan et al., 2016; Rakelly et al., 2019; Zintgraf et al., 2020; Melo, 2022; Grigsby et al., 2024a). These meta-features are generally learned through the use of recurrent neural networks (RNNs) (Duan et al., 2016; Zintgraf et al., 2020), amortized inference (Rakelly et al., 2019) or transformers (Melo, 2022; Grigsby et al., 2024a). For instance, Rakelly et al. (2019) infer a posterior over meta-

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

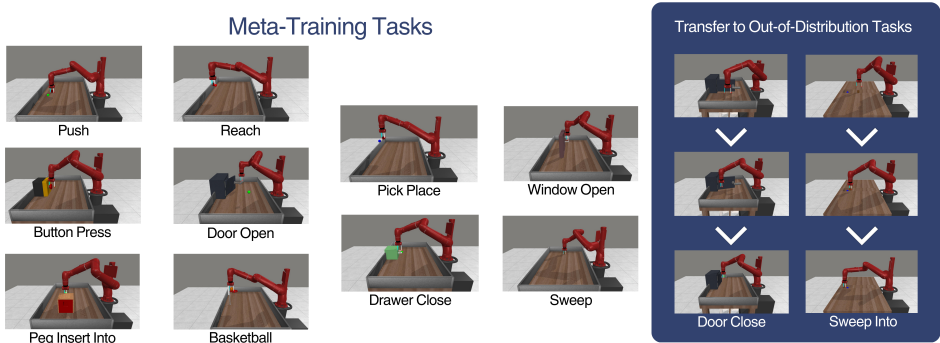


Figure 1: Illustration of the setting of the ML10 Benchmark in Meta-World. We evaluate ECET in each task and collect these frames. On the left are the tasks ECET is trained on, whereas on the right we show 3 frames (*beginning*, *middle* and *end* of the episode) from evaluating on 2 of the 5 test tasks. Though these tasks are not present in the training set, ECET manages to generalize and successfully close the door and sweep the block into the hole.

features by sampling transitions from the same task without explicitly considering the order of sequencing of these transitions. Thus, the meta-features generated have limited information on the task dynamics. On the other hand, RNN- and transformer-based methods process entire sequences of recent transitions spanning a few episodes to infer meta-features. This enables such methods to also incorporate dynamics information in the meta-features. Due to the vanishing gradient problem and the limited memory capacity of RNNs, the meta-features generated by RNN-based methods often emphasize characteristics of recent transitions, focusing primarily on intra-episode experiences. Similarly, Melo (2022) propose TrMRL, a transformer-based Meta-RL approach that generates meta-features by taking as input a sequence of only  $T = 5$  most recent transitions. Grigsby et al. (2024a) with AMAGO further increase the amount of task information processed to generate meta-features by utilizing a transformer architecture capable of processing longer sequences of transitions spanning multiple episodes (i.e., cross-episode experiences). While exhibiting impressive performance in Meta-RL, both TrMRL and AMAGO present substantial computational overhead with a runtime complexity depending quadratically on the complete sequence length processed.

In this paper, we present a novel method for learning meta-features through efficient cross-episode<sup>1</sup> Meta-RL using a novel transformer architecture. Our approach enhances the adaptability of RL agents across diverse sets of tasks by leveraging both intra-episode and cross-episode experiences, providing a more comprehensive learning process. Our learned algorithm is capable of outperforming the previous state-of-the-art and providing more efficient meta-training in terms of timesteps, while significantly improving generalization capabilities. To foster reproducibility, we provide our code at <https://anonymous.4open.science/r/ECET-8137>.

## 2 RELATED WORK

Meta-Reinforcement Learning (Meta-RL)<sup>2</sup> emerged as a framework aimed at making RL algorithms more sample-efficient by leveraging knowledge learned from previous tasks to speed up learning for the current task. Meta-RL optimizes a bilevel optimization problem during the *meta-training* phase. The inner loop typically consists of traditional RL algorithms, while the outer loop represents meta-learning algorithms that exploit information from the inner loop to improve performance across multiple tasks. Once the meta-training phase is complete, only the inner loop algorithm is used in the *adaptation* phase. Due to the meta-training phase, the inner loop algorithm can achieve good performance using fewer interactions with the new task compared to starting from scratch.

<sup>1</sup>Throughout this paper, we use the term *cross-episode* to describe interactions across multiple episodes, which is synonymous with *inter-episode*. This choice aligns with the terminology used in the title for consistency

<sup>2</sup>For a recent survey, see Beck et al. (2023).

108 Finn et al. (2017) proposed the Model-Agnostic Meta-Learning (MAML) algorithm, which can be  
109 used with any policy gradient algorithm as the inner loop to learn an initialization for the policy to  
110 speed up learning on a new test task. Duan et al. (2016) introduced  $RL^2$ , an approach that learns a  
111 reinforcement learning algorithm using RNNs, which share the hidden state across episodes from  
112 the same task.  $RL^2$  uses a policy gradient algorithm as the meta-learning algorithm in the outer  
113 loop, with learning during the adaptation phase occurring only within the RNN’s dynamics, without  
114 gradient adaptation. Another approach, PEARL, proposed by Rakelly et al. (2019), uses an off-policy  
115 RL algorithm, Soft Actor Critic (SAC) (Haarnoja et al., 2018), in the outer loop, and augments the  
116 state with stochastic meta-features of the task in the inner loop. An MLP generates the parameters of  
117 the posterior distribution of task meta-features.

118 Zintgraf et al. (2020) proposed VariBAD, a variant of  $RL^2$  for Bayes-adaptive deep RL, which uses  
119 VAEs to learn an approximate posterior distribution over task features to augment the state input  
120 to the policy. VariBAD uses Proximal Policy Optimization (PPO) (PPO; Schulman et al., 2017)  
121 as the outer loop algorithm, but it can be replaced by any policy gradient algorithm.  $RL^3$  (Bhatia  
122 et al., 2023) builds on  $RL^2$  by incorporating action-value estimates as additional input to the RNN,  
123 enhancing its capacity to capture task features.

124 Under the umbrella of meta-RL, there has also been work on augmenting the meta-training task set to  
125 improve generalization (Lee & Chung, 2021; Rimon et al., 2022). These methods can be used with  
126 meta-RL algorithms based on task meta-features to improve generalization.

127 The transformer architecture (Vaswani et al., 2017) has revolutionized the field of Natural Language  
128 Processing (NLP) and beyond. Its reliance on self-attention mechanisms enables it to process entire  
129 sequences of data simultaneously, enhancing efficiency and handling global dependencies in data  
130 compared to RNNs. Parisotto et al. (2020) explored the application of transformer architecture in RL  
131 settings, addressing issues such as data efficiency and overfitting. Subsequent works, including the  
132 Decision Transformer (Chen et al., 2021) treated RL as a sequence modeling problem, leveraging the  
133 scalability and simplicity of transformers to efficiently handle RL tasks, a departure from conventional  
134 RL methods that typically rely on value function estimation or policy gradients. Janner et al. (2021)  
135 extended this concept, showcasing how transformers can unify various RL tasks into a single sequence  
136 modeling framework, promising improvements in offline RL algorithms. Extensions such as MGDT  
137 (Lee et al., 2022) and Gato (Reed et al., 2022) showcased the flexibility of transformers in multi-task  
138 and multi-modal RL, though they required expert demonstrations or fine-tuning for unseen tasks.  
139 Offline RL works use datasets of collected trajectories and a supervised prediction loss to train the  
140 transformer. This means that the learning process is based on a fixed dataset of experiences and that  
141 the model does not interact with the task during training. Although effective for certain scenarios,  
142 this approach might not adapt well to dynamic or unpredictable tasks where new experiences differ  
143 significantly from the training data.

143 Despite their promise, stabilizing transformer training in RL has proven challenging (Mishra et al.,  
144 2018; Parisotto et al., 2020; Melo, 2022; Grigsby et al., 2024a). Recent work has proposed stabilizing  
145 training through changes in the architecture (Parisotto et al., 2020) or the weight initialization schema  
146 (Huang et al., 2020; Melo, 2022). We follow Huang et al. (2020) and Melo (2022) and stabilize  
147 transformer optimization through weight initialization.

148 Numerous approaches exist in the literature for hierarchical transformers for offline RL (Shang  
149 & Ryoo, 2021; Correia & Alexandre, 2022; Mao et al., 2023; Zhang et al., 2023). These works  
150 mostly focus on leveraging datasets of demonstrations to perform well in the same task-distribution  
151 as the demonstrations. In contrast, our focus lies in online RL and generalization across tasks  
152 with varying dimensions (e.g., parametric and task variations). By leveraging the sequencing of  
153 transitions and cross-episode information, we aim to efficiently meta-learn task representations that  
154 adapt dynamically.

155 Meta-feature extraction lies at the heart of our method. Approaches like VariBAD (based on RNNs)  
156 and PEARL (based on amortized inference) extract task-relevant information from transitions within  
157 a task, augmenting the state input with meta-features. These meta-features capture the nuances  
158 of task dynamics and improve policy performance. Expanding on this, TrMRL (Melo, 2022)  
159 demonstrated that transformers can generalize better by using a sequence of the most recent transitions,  
160 capturing dependencies through self-attention. Similarly, AMAGO (Grigsby et al., 2024a) showed  
161 that extracting information from longer sequences, spanning multiple episodes, significantly enhances

162 generalization. In AMAGO-2 (Grigsby et al., 2024b), the authors evaluate their approach more  
 163 extensively in multi-task RL problems and showcase the generalization abilities to unseen parametric  
 164 variations of tasks in the meta-training set (e.g. harder levels).

165 Building on these foundations, our approach integrates intra-episode (within a single episode)  
 166 and cross-episode (across multiple episodes) information to extract richer meta-features. Using  
 167 transformers, we analyze sequences of transitions to capture complex, multi-level patterns that identify  
 168 the task dynamics. This design enables in-context learning, allowing our model to dynamically adjust  
 169 its behavior based on the context provided by past experiences.

170 Our method achieves efficiency in terms of runtime complexity while capturing intricate patterns  
 171 in dynamic and diverse task environments. By leveraging information from both intra- and cross-  
 172 episode experiences, we improve generalization across tasks with varying dimensions. Compared  
 173 to prior state-of-the-art methods, our approach demonstrates improved meta-training efficiency and  
 174 generalization.  
 175

### 176 3 PRELIMINARIES

177 The Reinforcement Learning (RL) problem is typically formalized within the framework of Markov  
 178 Decision Processes (MDPs), defined by a tuple  $(S, A, P, R, \gamma)$ .  $S$  represents a finite set of states in the  
 179 task. Each state  $s \in S$  encapsulates all the information that describes the situation at a particular point  
 180 in time.  $A$  denotes a set of actions that can be executed in all states.  $P: S \times A \times S \rightarrow [0, 1]$  is the state  
 181 transition probability function that describes the dynamics of the MDP. More specifically,  $P(s'|s, a)$   
 182 gives the probability of transitioning to state  $s'$  when action  $a$  is taken in state  $s$ .  $R: S \times A \rightarrow \mathbb{R}$  is  
 183 the reward function, with  $R(s, a)$  giving the immediate reward received after taking action  $a$  from  
 184 state  $s$ .  $\gamma \in [0, 1]$  is the discount factor that determines the difference in importance between future  
 185 rewards and current rewards. For an agent interacting with the MDP, lower values of  $\gamma$  mean that the  
 186 agent will prioritize immediate rewards more strongly, while a higher value indicates a preference for  
 187 long-term gains.  
 188

189 The goal in RL is to find an optimal policy  $\pi^* \in \Pi$ , with  $\pi: S \rightarrow A$ , that maximizes the expected  
 190 discounted cumulative reward. This can be formally defined as  $\pi^* \in \operatorname{argmax}_{\pi \in \Pi} J(\pi)$  with  
 191

$$192 J(\pi) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s_0 = s, \pi \right]$$

193 Here,  $J(\pi)$  denotes the expected cumulative reward when starting in  $s_0 = s$ , playing action  $a_t$  in  
 194 state  $s_t$  as dictated by policy  $\pi$ , i.e.,  $a_t \sim \pi(s_t)$ . Most modern RL approaches parameterize the  
 195 policy as  $\pi_\theta$ , where the parameters  $\theta$  are typically the weights of a neural network, and aim to find  
 196 the parameters that satisfy the previous expression.  
 197

198 **Meta-Reinforcement Learning (Meta-RL):** Given a distribution of tasks  $p(\mathcal{T})$ , where each task  
 199  $\tau \sim p(\mathcal{T})$  can be considered a distinct MDP with its own state space  $S^\tau$ , action space  $A^\tau$ , transition  
 200 dynamics  $P^\tau$  and reward function  $R^\tau$ , the objective of a meta-RL technique is to learn a reinforcement  
 201 learning algorithm  $\mathcal{L}_\phi$  that directly outputs the parameters  $\theta$  of a reinforcement learning policy. The  
 202 learned algorithm  $\mathcal{L}$  is parameterized by  $\phi$  such that  $\mathcal{L}_\phi: \mathcal{T} \rightarrow \theta$  and aims to learn the parameters  $\theta$   
 203 so  $\pi_\theta$  can generalize across a set of tasks sampled from the distribution  $p(\mathcal{T})$ . To this end, the policy  
 204 interacts with a series of tasks  $\tau$ , drawn from  $p(\mathcal{T})$ , and aims to maximize its performance across this  
 205 task distribution. This process involves two key phases:  
 206

207 **Meta-Training and Meta-Testing Phases:** During meta-training, the agent learns a policy  $\pi_\theta$  and an  
 208 update rule  $\mathcal{L}_\phi$  (parameterized by  $\theta$  and  $\phi$ , respectively) that can quickly adapt to new tasks, resulting  
 209 in a bilevel optimization problem. This further involves learning a generalizable representation or  
 210 prior knowledge that is transferable across tasks. The outer objective in this phase is to optimize  
 211 the parameters  $\phi$  such that  $\phi^* \in \operatorname{argmax}_\phi G(\phi)$  with  $G(\phi) = \mathbb{E}_{\tau \sim p(\mathcal{T})} [J(\pi_\theta) | \mathcal{L}_\phi, \tau]$  whereas the  
 212 inner objective involves finding the weights of the policy that maximize  $\theta^* \in \operatorname{argmax}_\theta J(\pi_\theta)$ . In the  
 213 meta-testing (adaptation) phase, the agent encounters new tasks  $\tau_{\text{new}}$  sampled from  $p(\mathcal{T})$ . Using  $\pi_{\theta^*}$   
 214 and  $\mathcal{L}_{\phi^*}$ , the agent quickly adapts to each new task. The adaptation effectiveness is measured by the  
 215 agent’s performance on these new tasks after a limited number of learning steps or experiences.

## 4 METHOD

We illustrate our proposed approach, which we dub ECET (Efficient Cross-Episodic Transformers), in Figure 2. Our architecture consists of two distinct groups of transformer encoder blocks. The first one is the *Intra-Episodic Transformer* (IET) that encapsulates the sequence of transitions inside an episode and allows us to understand the transition dynamics within each episode. The second group of encoder blocks is the *Cross-Episodic Transformer* (CET) that encodes the set of episodes and captures the transition dynamics across episodes.

Our experience on a task consists of  $E$  episodes, including the most recent episode. For efficiency reasons, we sample a sequence of  $T$  transitions from each episode. For the current episode, we ensure that the sequence ends with the current state observation of the agent, such that it can choose the most appropriate next action. The set of  $E$  episodes, each having  $T$  transitions, are stored as sequences of tuples  $\chi_t^e = (s_t^e, a_t^e, r_t^e, d_t^e)_{e \in \{1, \dots, E\}, t \in \{1, \dots, T\}}$  where  $d_t^e \in \{0, 1\}$  is a termination flag. The  $t$ -th transition from the sequence of the  $e$ -th episode is denoted as  $\chi_t^e \in S \times A \times R \times \{0, 1\}$ , while a sequence of  $T$  transitions from the  $e$ -th episode as  $(\chi_0^e, \chi_1^e, \dots, \chi_T^e) \in (S \times A \times R \times \{0, 1\})^T$ .

The *Intra-Episodic Transformer* (IET) computes a representation that transforms a sequence of  $T$  transitions from the  $e$ -th episode into a vector representation  $z^e \in \mathbb{R}^K$  as formalized below:

$$\mathbf{z}^e = \text{IET}(\chi_1^e, \dots, \chi_T^e), \quad \text{IET} : (S \times A \times R \times \{0, 1\})^T \rightarrow \mathbb{R}^K, \quad \mathbf{z}^e \in \mathbb{R}^K$$

The IET employs an encoder-only transformer architecture with positional encoding, that projects the sequence of transitions to a latent representation using multiple layers of transformer blocks. The network’s representation of the sequence is the  $\mathbb{R}^K$  latent embedding of the  $T$ -th transition in the last transformer block.

We then input the intra-episodic representations  $\mathbf{z}^e \in \mathbb{R}^K, \forall e \in \{1, \dots, E\}$  from all the episodes into the *Cross-Episodic Transformer* (CET), which captures the dynamics of sequences from different experience episodes. The CET network ultimately represents the experience on a task into a vector  $\mathbf{z}_{\text{task}}$ , formalized as:

$$\mathbf{z}_{\text{task}} = \text{CET}(\mathbf{z}^1, \dots, \mathbf{z}^E), \quad \text{CET} : \mathbb{R}^{E \times K} \rightarrow \mathbb{R}^D, \quad \mathbf{z}_{\text{task}} \in \mathbb{R}^D$$

The CET network is also a stack of transformer blocks. The representation of the episodic experience  $\mathbf{z}_{\text{task}}$  is the  $\mathbb{R}^D$  output for the  $E$ -th element of the last transformer block in CET. Since episodes  $\mathbf{z}^e \in \mathbb{R}^K, \forall e \in \{1, \dots, E\}$  represent a set, we construct CET without positional encoding.

Finally, we condition the policy on the learned cross-episodic features and a linear transformation  $\varphi$  of the current state  $s$ :

$$\lambda = \pi(\varphi(s), \mathbf{z}_{\text{task}})$$

where the output of the policy,  $\lambda$ , represents the parameters of the distribution  $p_a$  from which we sample the next action  $a \sim p_a(\lambda)$ . We use a Gaussian distribution as the model for the probability distribution of the actions, therefore,  $\lambda$  is the mean and variance of the distribution.

Following the RL<sup>2</sup> framework, we use PPO as the outer-loop algorithm to train our pipeline. We use the PPO objective for the actor and Negative Log-Likelihood loss for the critic. We detail our Meta-Training and Meta-Testing Algorithms in Algorithm 2 and Algorithm 3 respectively.

The hierarchical nature of our novel transformer architecture improves the runtime complexity of computing the task representation. Suppose we would have opted for a standard (non-hierarchical) transformer network (Melo, 2022) that is fed a single long sequence of all the transitions from all episodes (i.e.  $E \cdot T$  many transitions). Such a naive approach requires an attention matrix of complexity  $O(E^2 T^2)$ . On the other hand, using the IET network to only capture intra-episodic features requires an attention matrix of complexity  $O(T^2)$ . Capturing cross-episodic features  $\mathbf{z}_{\text{task}}$  using CET requires an attention matrix of complexity  $O(E^2)$ . Thus, our two-tiered architecture of the cross-episodic transformer has a runtime complexity of  $O(E^2 + T^2)$ , and scales to longer sequences and more episodes than a standard non-hierarchical transformer.

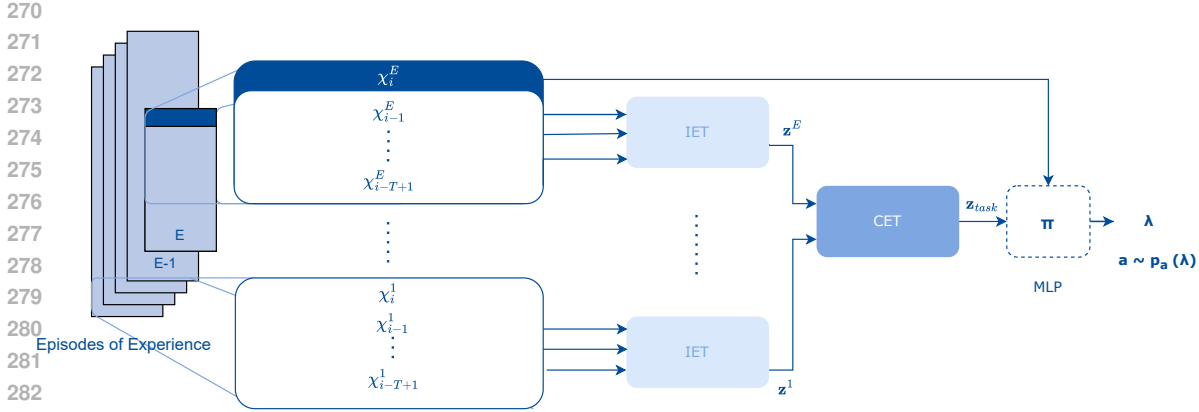


Figure 2: Illustration of our ECET architecture. Transitions  $\chi_t^e = (s_t^e, a_t^e, r_t^e, d_t^e)$  represent which state  $s_t^e$ , action  $a_t^e$ , and reward  $r_t^e$  were observed at the  $t$ -th transition of the  $e$ -th episode, but also a termination flag  $d_t^e$  that indicates if the  $e$ -th episode already terminated at time  $t$ . By inputting the sequences independently through the Intra-Episodic Transformer (IET), we generate a feature vector  $z^e$  for each episode  $e \in \{1, \dots, E\}$ . We then pass these sequence representations through the Cross-Episodic Transformer (CET) to generate a task representation  $z_{task}$  as input to the policy  $\pi$ .

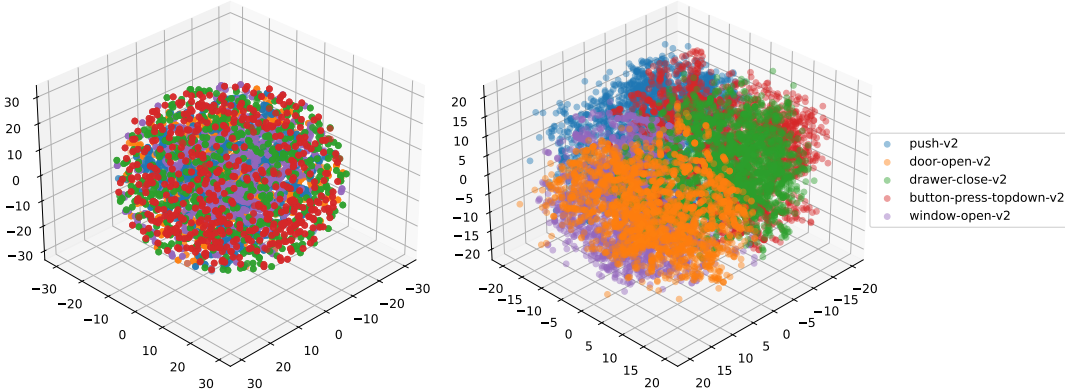


Figure 3: T-SNE plots of the output embeddings for TrMRL (left), and our ECET (right) for five of the tasks in the training set of the ML10 benchmark of Meta-World. For better visibility, we only plot 5 tasks here but provide a version with all 10 in Figure 14.

## 5 EXPERIMENTS

### 5.1 EXPERIMENTAL PROTOCOL

**MuJoCo** We evaluate ECET and our baseline methods on MuJoCo (Todorov et al., 2012) locomotion tasks, which are widely used in the Meta-RL literature. Specifically, we consider the **AntDir** and **HalfCheetahDir** tasks, where the agent is required to move either forwards or backwards. Additionally, we include the **HalfCheetahVel** task, where the agent must adapt to running at different target velocities. The maximum episode length for these tasks is 200 timesteps.

**Meta-World** To assess the performance of our cross-episode transformer approach, we use the Meta-World Benchmark (Yu et al., 2019). This benchmark is designed to help evaluate the performance of meta-RL baselines in 50 different simulated robotic manipulation tasks. These tasks all share the same state  $S$  and action space  $A$  as they are all located on a tabletop environment with a Sawyer arm (see Appendix A.4 for details and illustrations of all tasks). What differs in each task is the

behavior required to complete it successfully. For meta-learning, Meta-World includes three modes with varying degrees of difficulty (for details, see Appendix A.4):

- **ML1** is designed to evaluate few-shot adaptation to parametric variations within one task.
- **ML10** evaluates few-shot adaptation to 5 unseen test tasks, after training on 10 tasks.
- **ML45** is similar to, but more complex than ML10, with 45 training and 5 test tasks.

For ML1, the desired goal position is not provided as input, so the meta-RL algorithm needs to determine the location of the goal through trial and error. Similarly, for ML10 and ML45, task IDs are not provided as input, so the meta-RL algorithm needs to identify tasks from experience. The maximum episode length for all MetaWorld tasks is 500 timesteps.

**ManiSkill** We additionally use the ManiSkill Benchmark (Gu et al., 2023) to assess the performance of our proposed method on tasks where the state representation is an image. This benchmark consists of robotic manipulation tasks to help evaluate the performance of approaches for embodied AI in terms of their generalization to parametric variations, as well as variations in the object being manipulated. We evaluate on the *StackCube* and *PickSingleYCB* tasks, located on a tabletop environment with a Panda robot. What differs in each task is the behavior required to complete it successfully:

- **StackCube**: The goal is to pick up a red cube and stack it on top of a green cube and let go of the cube without it falling. The task is parameterized by randomly sampling the initial position of the cubes.
- **PickSingleYCB**: The goal is to pick up a random object sampled from the YCB dataset (Calli et al., 2015) and move it to a random goal position. This task is designed to evaluate few-shot adaptation to different objects.

We use the CNN architecture from Mnih et al. (2015b) to process RGB states for all methods. The maximum episode length for all ManiSkill tasks is 50 timesteps.

**Performance metrics** For the MuJoCo benchmark, we use the *average return* as a performance metric. The authors of Meta-World and ManiSkill propose the *success rate* as a performance evaluation metric. The distance of the task-relevant object from the goal position determines whether the agent is successful or not in an episode. Thus, the success rate is the fraction of episodes where the agent manages to place the task-relevant object inside a pre-set perimeter around the goal position.

Due to the fact that the tasks have different levels of difficulty, especially for the ML10 and ML45 modes in Meta-World, we additionally use the *average rank* as a performance metric for aggregating the performance across different tasks. We calculate the rank of methods over time for each task in the set based on their success rate, and then average over those ranks across tasks.

**Baselines:** We focus on investigating the performance of ECET compared to other online meta-reinforcement learning methods, namely MAML PPO (Finn et al., 2017), RL<sup>2</sup> (Duan et al., 2016), PEARL (Rakelly et al., 2019), VariBAD (Zintgraf et al., 2020), TrMRL (Melo, 2022) and AM-AGO (Grigsby et al., 2024a). Detailed descriptions of each are given in Appendix A.5. We conducted all experiments on a compute cluster of NVIDIA A100 GPUs. We trained all methods for  $10^7$  steps in the MuJoCo environments,  $5 \times 10^7$  steps in Meta-World, and  $2.5 \times 10^7$  steps in ManiSkill. We repeat each run for 5 different random seeds. When comparing the methods in the figures, we plot the mean and standard error of the performance metrics across seeds.

## 5.2 HYPOTHESES AND RESULTS

**Hypothesis 1: ECET captures more general task features compared to state-of-the-art meta-RL methods. This stems from ECET’s ability to capture intra- as well as cross-episodic experiences.** We compare the embeddings that ECET learns to those of TrMRL to investigate the extent to which they differentiate between tasks. We sample episodes from the training tasks in ML10 and plot the output of the transformers for both ECET and TrMRL through dimensionality reduction. For the sake of visual clarity, in Figure 3 we show the results for 5 of these tasks. For the plot with all the 10 tasks, see Figure 14. Additionally, in Figure 15 we compare the intra- and inter-task embedding distances. From these plots it is visible that the learned embeddings

378  
 379  
 380  
 381  
 382  
 383  
 384  
 385  
 386  
 387  
 388  
 389  
 390  
 391  
 392  
 393  
 394  
 395  
 396  
 397  
 398  
 399  
 400  
 401  
 402  
 403  
 404  
 405  
 406  
 407  
 408  
 409  
 410  
 411  
 412  
 413  
 414  
 415  
 416  
 417  
 418  
 419  
 420  
 421  
 422  
 423  
 424  
 425  
 426  
 427  
 428  
 429  
 430  
 431

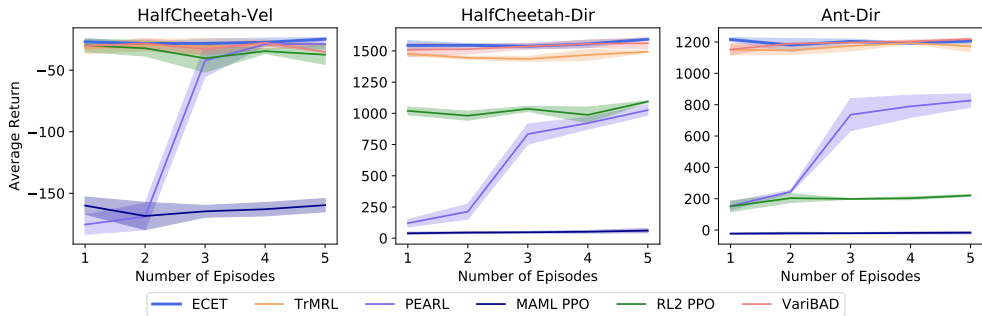


Figure 4: Test performance in terms of average return across rollouts for ECET, TrMRL, PEARL, MAML PPO, RL2 PPO and VariBAD on the *HalfCheetah-Vel*, *HalfCheetah-Dir*, and *Ant-Dir* tasks of the MuJoCo benchmark, testing adaptation to parametric variations of these tasks. We show the corresponding meta-training performance plot in Figure 12.

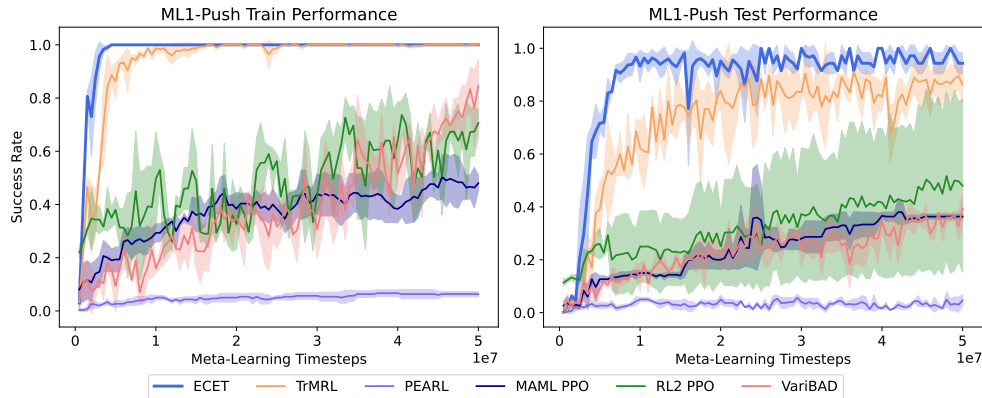


Figure 5: Meta-Train and Test performance in terms of average success rate for ECET, TrMRL, PEARL, MAML PPO, RL2 PPO and VariBAD on the ML1 benchmark for training (left) and testing (right) on parametric variations of the *Push* task. We show the corresponding plot with respect to time in Figure 19, whereas the average rank plot can be found in Figure 16.

of ECET produce a clearer grouping for the same task, compared to the embeddings of TrMRL. Furthermore, we can see a grouping of tasks that are similar to each other, e.g., *Window Open* and *Door Open*. This shows the usefulness of cross-episodic experiences in task differentiation. For easier analysis, we provide scripts for interactive versions of the T-SNE plots in our code repository [https://anonymous.4open.science/r/ECET-8137/plotly\\_interactive\\_plot.py](https://anonymous.4open.science/r/ECET-8137/plotly_interactive_plot.py).

**Hypothesis 2: ECET outperforms the state-of-the-art in online meta-RL algorithms in few-shot adaptation to parametric variations of tasks.** To evaluate adaptability to parametric variations of the task, we compare ECET and the baselines in the MuJoCo set of tasks, and in the ML1 mode of Meta-World (*Push* and *Reach* tasks). We evaluate the generalization on the parametrized MuJoCo tasks through doing 5 rollouts at test time, and calculating the return for each. We show the results in Figure 4. ECET, TrMRL, and VariBAD show robust test performance across the three tasks. For the Meta-World ML1 benchmark, we show the performance in terms of Success Rate in Figure 5 for *Push*. The figure contains results for meta-training (left) and test (right) performance in the case where we train on variations of the *Push* task, and test the adaptability of the methods on a set of unseen variations of the same task. In terms of meta-training performance (left), we can see that transformer-based methods are faster in differentiating between the variations of the task and thus achieve a maximum success rate of 1.0. In terms of test performance (right), ECET and TrMRL generalize successfully, reaching a test success rate greater than 0.8. In addition, ECET shows a higher generalizability compared to TrMRL. RL2 PPO struggles to robustly generalize, exhibiting



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

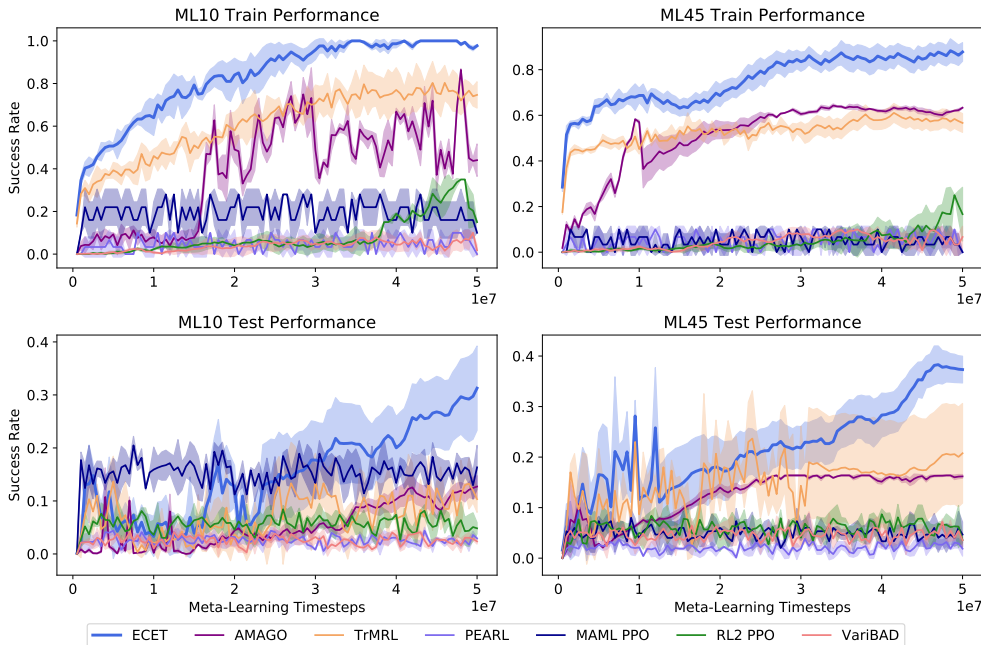


Figure 6: Meta-Train and Test performance in terms of average success rate for ECET, AMAGO, TrMRL, PEARL, MAML PPO, RL2 PPO and VariBAD on the ML10(left) and ML45(right) benchmarks. Here, we are training and testing on disjoint sets of tasks. We show the corresponding plot with respect to time in Figure 21, whereas the average rank plot can be found in Figure 17.

a higher variation in performance compared to the other methods. In Figure 13 in Appendix A.6 we show the Meta-Train and Test performance on variations of the *Reach* task. The task variability for *Reach* is not as difficult to learn, thus all the methods perform better both when meta-training and testing. To further assess generalization to parametric variation, we evaluate our method and the baselines in the *StackCube* task of ManiSkill. We show the meta-train and test performance in Figure 7 (left column), where we notice a similar pattern on performances as with ML1, with ECET outperforming the baselines after  $2.5 \cdot 10^7$  meta-training timesteps.

**Hypothesis 3: ECET outperforms the state-of-the-art online meta-RL in out-of-distribution (OOD) tasks.** We define out-of-distribution (OOD) performance to be the performance of a meta-learned agent on a set of tasks that is disjoint from the tasks it is meta-learned on. We show the performance of the methods in terms of the average success rate for the meta-training (top) and test (bottom) performance for ML10 (left) and ML45 (right) in Figure 6. In this setting we also compare to AMAGO, as a recently proposed powerful transformer-based method processing long sequences. ECET is the only method that achieves an average success rate greater than 0.8 in ML10 and greater than 0.6 for ML45 in terms of meta-training performance with a final success rate roughly 20% higher than the baselines. The other methods fail to learn policies that successfully identify the different tasks in the training set for the given budget of  $5 \times 10^7$  meta-learning timesteps, and thus struggle to solve them. In terms of OOD performance, the bottom row shows that generalization to distinct tasks is still a challenge for meta-RL. ECET achieves an average test success rate of approximately 0.35 for ML10 and 0.4 for ML45. Similarly, in Figure 7 (right column) where we show the performance on *PickSingleYCB*, we notice that all baselines except for TrMRL struggle to effectively generalize on the task of picking different objects. ECET is the best performing method, achieving an average test success rate of 0.39. This further demonstrates that cross-episodic features can improve meta-RL performance without incurring the cost of gradient adaptation at test time.

To focus on comparing the methods based on their robustness in relative performance to each other, in Appendix A.6.2 we provide plots showing the average ranks of the methods as a performance metric. These plots show that ECET robustly achieves the highest relative performance (and thus the

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

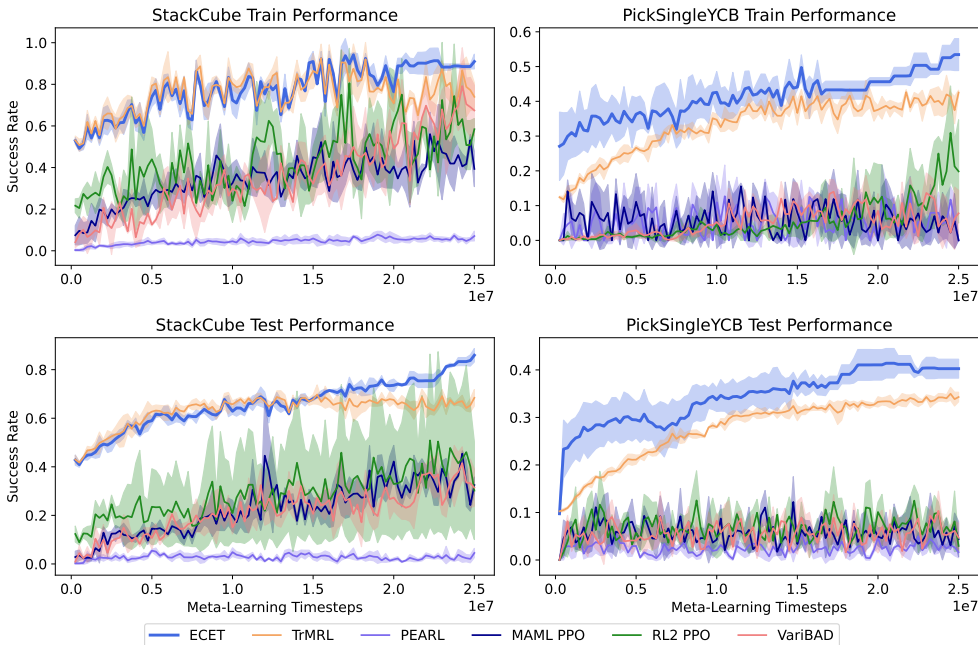


Figure 7: Average success rate plots of Meta-Train and Test performance for all methods on the *StackCube*(left) and *PickSingleYCB*(right) tasks. We show the corresponding plot with respect to time in Figure 22, whereas the average rank plot can be found in Figure 18.

lowest rank) across the meta-learning benchmarks we consider. For a detailed illustration of per-task performance, we refer to Appendix A.6.6. In Appendix A.6.4 through plots of episodic return across 50 episodes at test time, we show that ECET performs well even in scenarios when the agent has not yet collected the  $E$  episodes of experience.

### 5.3 ABLATIONS

We investigate the behavior and the impact of ECET’s components on performance in Appendix A.6.5. We show the effect of varying the number of episodes ( $E$ ) in memory and the sequence length ( $T$ ) of sampled transitions in Figure 28. Figure 29 (top plot) shows that sampling without replacement results in more reliable and faster performance improvements, especially for smaller ( $E$ ) and ( $T$ ). In terms of architectural variations, we investigated the effects of adding positional encoding in the Cross-Episode Transformer (CET) and simplifying ECET to a flat architecture akin to TrMRL and AMAGO. The results in Figure 26 indicate that positional encoding in CET introduces misleading ordering, leading to suboptimal convergence, while the flat architecture struggles to capture comprehensive task representations.

## 6 CONCLUSION

We introduced Efficient Cross-Episode Transformers (ECET) for Meta-Reinforcement Learning, a powerful approach to online meta-RL. Our experiments within the MuJoCo, Meta-World and ManiSkill benchmarks highlight our method’s capability for few-shot adaptation to tasks with parametric variations, as well as its capacity to adjust to completely unseen tasks that share the same state and action spaces as the training tasks. Despite outperforming state-of-the-art methods in online meta-RL by efficiently compressing cross-episodic knowledge, our method’s generalization to tasks that are completely different to the meta-training set remains limited as our approach is tailored to exploiting similarities. Future work should consider using task augmentation to improve generalization, modifying the adaptation phase, or increasing model capacity to handle a wider context in order to improve performance in this phase.

## REFERENCES

- 540  
541  
542 M. Andrychowicz, A. Raichuk, P. Stanczyk, M. Orsini, S. Girgin, R. Marinier, L. Hussenot, M. Geist,  
543 O. Pietquin, M. Michalski, S. Gelly, and O. Bachem. What matters for on-policy deep actor-critic  
544 methods? A large-scale study. In *9th International Conference on Learning Representations, ICLR*  
545 *2021*. OpenReview.net, 2021.
- 546 Sebastian Pineda Arango, Fabio Ferreira, Arlind Kadra, Frank Hutter, and Josif Grabocka. Quick-  
547 tune: Quickly learning which pretrained model to finetune and how. In *The Twelfth International*  
548 *Conference on Learning Representations, 2024*. URL [https://openreview.net/forum?](https://openreview.net/forum?id=tqhlzdxIra)  
549 [id=tqhlzdxIra](https://openreview.net/forum?id=tqhlzdxIra).
- 550 J. Beck, R. Vuorio, E. Z. Liu, Z. Xiong, L. Zintgraf, C. Finn, and S. Whiteson. A survey of  
551 meta-reinforcement learning. *arXiv preprint arXiv:2301.08028*, 2023.
- 552  
553 M. G. Bellemare, S. Candido, P. Samuel Castro, J. Gong, M. C. Machado, S. Moitra, S. S. Ponda, and  
554 Z. Wang. Autonomous navigation of stratospheric balloons using reinforcement learning. *Nature*,  
555 588(7836):77–82, 2020.
- 556 C. Benjamins, T. Eimer, F. Schubert, A. Mohan, S. Döhler, A. Biedenkapp, B. Rosenhahn, F. Hutter,  
557 and M. Lindauer. Contextualize me – the case for context in reinforcement learning. *Transactions*  
558 *on Machine Learning Research*, 2023. ISSN 2835-8856.
- 559  
560 M. Beukman, D. Jarvis, R. Klein, S. James, and B. Rosman. Dynamics generalisation in reinforcement  
561 learning via adaptive context-aware policies. In *Thirty-seventh Conference on Neural Information*  
562 *Processing Systems (NeurIPS 2023)*, 2023.
- 563 Abhinav Bhatia, Samer Nashed, and Shlomo Zilberstein. RL<sup>3</sup>: Boosting meta reinforcement  
564 learning via RL inside RL<sup>2</sup>. In *NeurIPS 2023 Workshop on Generalization in Planning*, 2023.  
565 URL <https://openreview.net/forum?id=ozqaF9YBce>.
- 566  
567 Berk Calli, Arjun Singh, Aaron Walsman, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M. Dollar.  
568 The ycb object and model set: Towards common benchmarks for manipulation research. In *2015*  
569 *International Conference on Advanced Robotics (ICAR)*, pp. 510–517, 2015. doi: 10.1109/ICAR.  
570 2015.7251504.
- 571 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,  
572 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence  
573 modeling. In *Advances in Neural Information Processing Systems*, volume 34, pp. 15084–15097.  
574 Curran Associates, Inc., 2021.
- 575  
576 K. Cobbe, C. Hesse, J. Hilton, and J. Schulman. Leveraging procedural generation to benchmark rein-  
577 forcement learning. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International*  
578 *Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp.  
579 2048–2056. PMLR, 13–18 Jul 2020.
- 580 André Correia and Luís A. Alexandre. Hierarchical decision transformer, 2022.
- 581  
582 J. Degraeve, F. Felici, J. Buchli, M. Neunert, B. Tracey, F. Carpanese, T. Ewalds, R. Hafner, A. Abdol-  
583 maleki, D. de las Casas, C. Donner, L. Fritz, C. Galperti, A. Huber, J. Keeling, M. Tsimpoukelli,  
584 J. Kay, A. Merle, Jean-M. Moret, S. Noury, F. Pesamosca, D. Pfau, O. Sauter, C. Sommariva,  
585 S. Coda, B. Duval, A. Fasoli, P. Kohli, K. Kavukcuoglu, D. Hassabis, and M. Riedmiller. Magnetic  
586 control of tokamak plasmas through deep reinforcement learning. *Nature*, 602(7897):414–419,  
587 2022.
- 588 Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. RL<sup>2</sup>: Fast  
589 reinforcement learning via slow reinforcement learning. *CoRR*, abs/1611.02779, 2016.
- 590  
591 T. Eimer, M. Lindauer, and R. Raileanu. Hyperparameters in reinforcement learning and how to  
592 tune them. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett (eds.),  
593 *International Conference on Machine Learning, ICML 2023*, volume 202 of *Proceedings of*  
*Machine Learning Research*, pp. 9104–9149. PMLR, 2023.

- 594 L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, and A. Madry. Implementation  
595 matters in deep RL: A case study on PPO and TRPO. In *8th International Conference on Learning*  
596 *Representations, ICLR 2020*. OpenReview.net, 2020.
- 597  
598 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation  
599 of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*,  
600 volume 70 of *Proceedings of Machine Learning Research*, pp. 1126–1135. PMLR, 2017. URL  
601 <https://proceedings.mlr.press/v70/finn17a.html>.
- 602 The garage contributors. Garage: A toolkit for reproducible reinforcement learning research. <https://github.com/rlworkgroup/garage>, 2019.
- 603  
604  
605 Jake Grigsby, Linxi Fan, and Yuke Zhu. AMAGO: Scalable in-context reinforcement learning for  
606 adaptive agents. In *The Twelfth International Conference on Learning Representations, 2024a*.  
607 URL <https://openreview.net/forum?id=M6XWoEdmwf>.
- 608  
609 Jake Grigsby, Justin Sasek, Samyak Parajuli, Daniel Adebisi, Amy Zhang, and Yuke Zhu. Amago-2:  
610 Breaking the multi-task barrier in meta-reinforcement learning with transformers, 2024b. URL  
611 <https://arxiv.org/abs/2411.11188>.
- 612 Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang, Stone  
613 Tao, Xinyue Wei, Yunchao Yao, Xiaodi Yuan, Pengwei Xie, Zhiao Huang, Rui Chen, and Hao  
614 Su. Maniskill2: A unified benchmark for generalizable manipulation skills. In *The Eleventh*  
615 *International Conference on Learning Representations, 2023*. URL [https://openreview.net/forum?id=b\\_CQDy9vrD1](https://openreview.net/forum?id=b_CQDy9vrD1).
- 616  
617 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
618 maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer Dy and Andreas  
619 Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80  
620 of *Proceedings of Machine Learning Research*, pp. 1861–1870. PMLR, 10–15 Jul 2018. URL  
621 <https://proceedings.mlr.press/v80/haarnoja18b.html>.
- 622  
623 P. Henderson, R. Islam, P. Bachman, J. Pineau, D. Precup, and D. Meger. Deep reinforcement  
624 learning that matters. In S. McIlraith and K. Weinberger (eds.), *Proceedings of the Thirty-Second*  
625 *Conference on Artificial Intelligence (AAAI’18)*. AAAI Press, 2018.
- 626  
627 T. M. Hospedales, A. Antoniou, P. Micaelli, and A. J. Storkey. Meta-learning in neural networks: A  
628 survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5149–5169, 2022.
- 629  
630 Xiao Shi Huang, Felipe Perez, Jimmy Ba, and Maksims Volkovs. Improving transformer optimization  
631 through better initialization. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th*  
632 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*  
633 *Research*, pp. 4475–4483. PMLR, 13–18 Jul 2020.
- 634  
635 Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence  
636 modeling problem. In *Advances in Neural Information Processing Systems*, volume 34, pp. 1273–  
637 1286. Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/099fe6b0b444c23836c4a5d07346082b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/099fe6b0b444c23836c4a5d07346082b-Paper.pdf).
- 638  
639 E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza. Champion-  
640 level drone racing using deep reinforcement learning. *Nature*, 620(7976):982–987, 2023.
- 641  
642 R. Kirk, A. Zhang, E. Grefenstette, and T. Rocktäschel. A survey of zero-shot generalisation in deep  
643 reinforcement learning. *J. Artif. Intell. Res.*, 76:201–264, 2023.
- 644  
645 Kuang-Huei Lee, Ofir Nachum, Mengjiao (Sherry) Yang, Lisa Lee, Daniel Freeman, Sergio Guadar-  
646 rama, Ian Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, and Igor Mordatch. Multi-game decision  
647 transformers. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 27921–27936. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/b2cac94f82928a85055987d9fd44753f-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/b2cac94f82928a85055987d9fd44753f-Paper-Conference.pdf).

- 648 Suyoung Lee and Sae-Young Chung. Improving generalization in meta-RL with imaginary tasks from  
649 latent dynamics mixture. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.),  
650 *Advances in Neural Information Processing Systems*, 2021. URL [https://openreview.net/forum?id=PnpS7\\_SlNZi](https://openreview.net/forum?id=PnpS7_SlNZi).
- 652 Hangyu Mao, Rui Zhao, Hao Chen, Jianye Hao, Yiqun Chen, Dong Li, Junge Zhang, and Zhen Xiao.  
653 Transformer in transformer as backbone for deep reinforcement learning, 2023.
- 655 Luckeciano C Melo. Transformers are meta-reinforcement learners. In *Proceedings of the 39th*  
656 *International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning*  
657 *Research*, pp. 15340–15359. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/melo22a.html>.
- 659 Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. A simple neural attentive  
660 meta-learner. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=B1DmUzWAW>.
- 663 V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A.  
664 Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King,  
665 D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforce-  
666 ment learning. *Nature*, 518(7540):529–533, 2015a.
- 667 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-  
668 mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen,  
669 Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra,  
670 Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning.  
671 *Nature*, 518(7540):529–533, February 2015b. ISSN 1476-4687. doi: 10.1038/nature14236. URL  
672 <https://doi.org/10.1038/nature14236>.
- 673 A. Mohan, C. Benjamins, K. Wienecke, A. Dockhorn, and M. Lindauer. Autorl hyperparameter  
674 landscapes. In *Second International Conference on Automated Machine Learning*. PMLR, 2023.
- 676 Emilio Parisotto, Francis Song, Jack Rae, Razvan Pascanu, Caglar Gulcehre, Siddhant Jayakumar,  
677 Max Jaderberg, Raphaël Lopez Kaufman, Aidan Clark, Seb Noury, Matthew Botvinick, Nicolas  
678 Heess, and Raia Hadsell. Stabilizing transformers for reinforcement learning. In *Proceedings of*  
679 *the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine*  
680 *Learning Research*, pp. 7487–7498. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/parisotto20a.html>.
- 682 J. Parker-Holder, R. Rajan, X. Song, A. Biedenkapp, Y. Miao, T. Eimer, B. Zhang, V. Nguyen,  
683 R. Calandra, A. Faust, F. Hutter, and M. Lindauer. Automated reinforcement learning (autorl): A  
684 survey and open problems. *Journal of Artificial Intelligence Research (JAIR)*, 74:517–568, 2022.
- 686 Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient off-policy  
687 meta-reinforcement learning via probabilistic context variables. In *Proceedings of the 36th*  
688 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*  
689 *Research*, pp. 5331–5340. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/rakelly19a.html>.
- 691 Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov,  
692 Gabriel Barth-maroon, Mai Giménez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom  
693 Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol  
694 Vinyals, Mahyar Bordbar, and Nando de Freitas. A generalist agent. *Transactions on Machine*  
695 *Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=1ikK0kHjvj>. Featured Certification, Outstanding Certification.
- 697 Zohar Rimon, Aviv Tamar, and Gilad Adler. Meta reinforcement learning with fi-  
698 nite training tasks - a density estimation approach. In S. Koyejo, S. Mohamed,  
699 A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Infor-*  
700 *mation Processing Systems*, volume 35, pp. 13640–13653. Curran Associates, Inc.,  
701 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/5833b4daf5b076dd1cdb362b163dff0c-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/5833b4daf5b076dd1cdb362b163dff0c-Paper-Conference.pdf).

- 702 J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart,  
703 D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver. Mastering atari, go, chess and shogi by  
704 planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 705  
706 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
707 optimization algorithms, 2017.
- 708 G. Shala, S. P. Arango, A. Biedenkapp, F. Hutter, and J. Grabocka. Autorl-bench 1.0. In *Workshop*  
709 *on Meta-Learning (MetaLearn@NeurIPS’22)*, 2022.
- 710  
711 Gresa Shala, André Biedenkapp, Frank Hutter, and Josif Grabocka. Gray-box gaussian processes  
712 for automated reinforcement learning. In *The Eleventh International Conference on Learning*  
713 *Representations*, 2023a. URL <https://openreview.net/forum?id=rmoMvptXK7M>.
- 714 Gresa Shala, Thomas Elsken, Frank Hutter, and Josif Grabocka. Transfer NAS with meta-learned  
715 bayesian surrogates. In *The Eleventh International Conference on Learning Representations*,  
716 2023b. URL <https://openreview.net/forum?id=paGvsrl4Ntr>.
- 717  
718 Jinghuan Shang and Michael S. Ryoo. Starformer: Transformer with state-action-reward representa-  
719 tions. *CoRR*, abs/2110.06206, 2021. URL <https://arxiv.org/abs/2110.06206>.
- 720 R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. Adaptive computation and  
721 machine learning. MIT Press, 2 edition, 2018.
- 722  
723 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.  
724 In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033.  
725 IEEE, 2012. doi: 10.1109/IROS.2012.6386109.
- 726 J. Vanschoren. Meta-learning. In F. Hutter, L. Kotthoff, and J. Vanschoren (eds.), *Automated Machine*  
727 *Learning - Methods, Systems, Challenges*, The Springer Series on Challenges in Machine Learning,  
728 pp. 35–61. Springer, 2019.
- 729  
730 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N  
731 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Ad-*  
732 *vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,  
733 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)  
734 [file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 735  
736 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey  
737 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.  
738 In *Conference on Robot Learning (CoRL)*, 2019. URL [https://arxiv.org/abs/1910.](https://arxiv.org/abs/1910.10897)  
739 [10897](https://arxiv.org/abs/1910.10897).
- 740  
741 Bin Zhang, Hangyu Mao, Lijuan Li, Zhiwei Xu, Dapeng Li, Rui Zhao, and Guoliang Fan. Stackelberg  
742 decision transformer for asynchronous action coordination in multi-agent systems, 2023.
- 743  
744 Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarín Gal, Katja Hofmann, and  
745 Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. In  
746 *International Conference on Learning Representations*, 2020. URL [https://openreview.](https://openreview.net/forum?id=Hkl9JlBYvr)  
747 [net/forum?id=Hkl9JlBYvr](https://openreview.net/forum?id=Hkl9JlBYvr).
- 748  
749  
750  
751  
752  
753  
754  
755

A APPENDIX

A.1 RL<sup>2</sup> FRAMEWORK

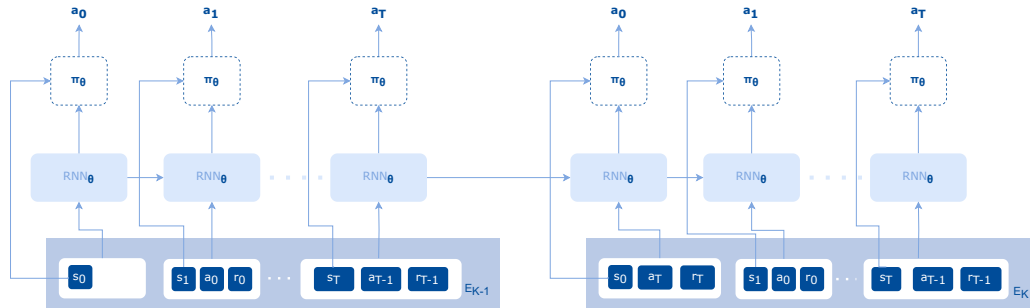


Figure 8: Illustration of the RL<sup>2</sup> framework, based on Figure 4 from Beck et al. (2023). As illustrated, the RNN takes as input the current state  $s_t$ , the action(if there is any) that lead to the state and the reward incurred from reaching that state. The output of the RNN is then used as input to the policy  $\pi$  along with the current state.

A.2 META-TRAINING AND META-TESTING ALGORITHMS

We include the detailed algorithm for Meta-Training in Algorithm 2 and the algorithm for Meta-Testing in Algorithm 3. In order to be consistent in comparison to the baselines, we use the garage contributors (2019) library implementations for the PPO algorithm, Gaussian Actor and Gaussian Critic.

A.2.1 INITIAL SETUP AND HYPERPARAMETERS OF META-TRAINING

We use as input a distribution of meta-training tasks  $\mathcal{T}$  and a history  $\mathcal{H}$ , a FIFO queue that stores past transitions to capture the temporal dynamics. The hyperparameters that guide the meta-training process include: the number of episodes  $E$  to maintain in the history for context, the sequence length  $S$  of transitions to be processed by the Intra-Episode Feature Extractor **IET**, the maximum number of training epochs  $max\_epochs$ , the sizes for both meta and mini-batches ( $meta\_batch\_size$ ,  $mini\_batch\_size$ ), and the number of episodes  $k$  collected per task in each training round.

A.2.2 PARAMETER DEFINITIONS

**Actor Parameters** include the parameters for embedding transitions ( $\varphi_1$ ) and states ( $\varphi_2$ ), the parameters for our Intra- and Inter-Episode Feature Extractors ( $\Theta_1, \Theta_2$ ) within the hierarchical transformer framework, and the policy network ( $\pi$ ) parameters ( $\Theta_3$ ). **Critic Parameters** include the parameters ( $\Psi$ ) for our value function ( $V$ ), which assists in estimating the expected returns and guiding the actor’s policy adjustments.

A.2.3 PROCEDURAL COMPONENTS

For clarity and conciseness of writing, we include some procedural components of our meta-training and meta-testing pipeline as procedures in Algorithm 1. These include:

- *CollectEpisode*: We gather one episode of interactions for a given task, utilizing the Actor ( $\Phi, \Theta$ ) for action decisions, and record these transitions within our history  $\mathcal{H}$  and an episode\_buffer which contains for each transition in the episode: the set of sampled transition sequences  $\mathcal{C}$  used to generate the task representation ,as well as the (*state, action, reward, terminated*) for processing.
- *CollectAndProcessEpisode*: Using the episode\_buffer data returned, we process each transition to compute advantages and returns, and we add

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

---

**Algorithm 1** Subroutines for ECET
 

---

**Input:** Distribution of tasks  $\mathcal{T}$ , history  $\mathcal{H}$  (FIFO queue) to store past transitions.  
**Hyperparameters:** number of episodes  $E$  to keep in history, sequence length  $T$  of transitions to input in **IET**,  $max\_epochs$ ,  $meta\_batch\_size$ ,  $mini\_batch\_size$ , number of episodes  $k$  to collect for each task during training.  
**Actor Parameters**  
 Transition( $\varphi_1$ ) and state( $\varphi_1$ ) embedding parameters:  $\Phi = \varphi_1, \varphi_2$   
 Intra-Episodic Transformer(**IET**( $\cdot; \Theta_1$ )), Cross-Episodic Transformer(**CET**( $\cdot; \Theta_2$ )), and policy( $\pi(\cdot; \Theta_3)$ ) parameters:  $\Theta = \{\Theta_1, \Theta_2, \Theta_3\}$   
**Critic Parameters**  
 Value function ( $(V(\cdot; \Psi))$ ) parameters:  $\Psi$   
**procedure** COLLECTEPISODE( $task, \Phi, \Theta, \mathcal{H}$ )  
   episode\_buffer.clear()  
   state  $\leftarrow$  task.reset()  
   terminated  $\leftarrow$  False  
   **while** not terminated **do**  
      $\mathcal{C} \leftarrow$  SAMPLETRANSITIONS( $\mathcal{H}, state$ )  
      $\mathbf{z}_{task} \leftarrow$  GENERATETASKREPRESENTATION( $\mathcal{C}, \varphi_1, \Theta_1, \Theta_2$ )  
      $\lambda \leftarrow \pi(\mathbf{z}_{task}, \varphi_2(state); \Theta_3)$   
     action  $\sim \lambda$   
     next\_state, reward, terminated  $\leftarrow$  task.act(action)  
     Add ( $state, action, reward, terminated$ ) to  $\mathcal{H}$   
     Add ( $\mathcal{C}, state, action, reward, terminated$ ) to episode\_buffer  
     state  $\leftarrow$  next\_state  
   **return** episode\_buffer,  $\mathcal{H}$   
**procedure** COLLECTANDPROCESSEPIISODE( $task, \Phi, \Theta, \Psi, \mathcal{H}, \mathcal{D}$ )  
   episode\_buffer,  $\mathcal{H} \leftarrow$  COLLECTEPISODE( $task, \Phi, \Theta, \mathcal{H}$ )  
   **for** ( $\mathcal{C}, state, action, reward, terminated$ ) in episode\_buffer **do**  
     Calculate *advantage* using the value function  $V(\cdot; \Psi)$ , and received reward.  
     Calculate the *return*.  
     Add ( $\mathcal{C}, state, action, reward, terminated, advantage, return$ ) to  $\mathcal{D}$   
   **return**  $\mathcal{D}, \mathcal{H}$   
**procedure** SAMPLETRANSITIONS( $\mathcal{H}, current\_state$ )  
   Initialize episode transitions set  $\mathcal{C}$  to empty.  
   **for**  $episode \leftarrow 1$  to  $E$  **do**  
     **if** episode not current **then**  
       Sample sequence of  $T$  transitions and add to  $\mathcal{C}$ .  
     **else**  
       Sample sequence of latest  $T$  transitions including  $current\_state$  and add to  $\mathcal{C}$ .  
   **return**  $\mathcal{C}$   
**procedure** GENERATETASKREPRESENTATION( $\mathcal{C}, \varphi_1, \Theta_1, \Theta_2$ )  
   Initialize set  $\mathcal{C}_{seq}$  to empty.  
   **for**  $i \leftarrow 1$  to  $E$  **do**  
      $T_i \leftarrow \mathcal{C}[i]$   
      $\mathbf{z}_{seq}^i \leftarrow$  **IET**( $\varphi_1(T_i); \Theta_1$ )  
     Add  $\mathbf{z}_{seq}^i$  to  $\mathcal{C}_{seq}$ .  
    $\mathbf{z}_{task} \leftarrow$  **CET**( $\mathcal{C}_{seq}; \Theta_2$ ).  
   **return**  $\mathbf{z}_{task}$

---



864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

---

**Algorithm 2** Meta-Training with ECET using Subroutines from Algorithm 1

---

**Input:** Distribution of meta-training tasks  $\mathcal{T}$ , history  $\mathcal{H}$  (FIFO queue) to store past transitions.

**Hyperparameters:** number of episodes  $E$  to keep in history, sequence length  $T$  of transitions to input in **IET**,  $max\_epochs$ ,  $meta\_batch\_size$ ,  $mini\_batch\_size$ , number of episodes  $k$  to collect for each task during training.

**Actor Parameters**

Transition( $\varphi_1$ ) and state( $\varphi_1$ ) embedding parameters:  $\Phi = \{\varphi_1, \varphi_2\}$

Intra-Episodic Transformer(**IET**( $\cdot; \Theta_1$ )), Cross-Episodic Transformer(**CET**( $\cdot; \Theta_2$ )), and policy( $\pi(\cdot; \Theta_3)$ ) parameters:  $\Theta = \{\Theta_1, \Theta_2, \Theta_3\}$

**Critic Parameters**

Value function ( $V(\cdot; \Psi)$ ) parameters:  $\Psi$

**Procedures:** COLLECTANDPROCESSEPIISODE, SAMPLETRANSITIONS

**procedure** ACTOR( $\{(\mathcal{C}, state)\}_{i=1}^{mini\_batch\_size}$ ,  $\Phi$ ,  $\Theta$ )

    Generate  $\{z_{task}\}_{i=1}^{mini\_batch\_size}$  using GENERATETASKREPRESENTATION( $\{\mathcal{C}\}_{i=1}^{mini\_batch\_size}$ ,  $\varphi_1, \Theta_1, \Theta_2$ ).

$\{\mu_{actor}, \sigma_{actor}\}_{i=1}^{mini\_batch\_size} \leftarrow \pi(\{z_{task}, \varphi_2(state)\}_{i=1}^{mini\_batch\_size}, \Theta_3)$ .

**return**  $\{\mu_{actor}, \sigma_{actor}\}_{i=1}^{mini\_batch\_size}$

**procedure** UPDATEMODELS( $\mathcal{D}$ ,  $\Phi$ ,  $\Theta$ ,  $\Psi$ )

$\Phi^{old}, \Theta^{old} \leftarrow \Phi, \Theta$

**for**  $\{(\mathcal{C}, state, action, reward, terminated, advantage, return)\}_{i=1}^{mini\_batch\_size}$  in  $\mathcal{D}$  **do**

$\lambda^{old} \leftarrow$  ACTOR( $\{(\mathcal{C}, state)\}_{i=1}^{mini\_batch\_size}$ ,  $\Phi^{old}, \Theta^{old}$ ).

$prob^{old} \leftarrow \lambda^{old}(\{action_i\}_{i=1}^{mini\_batch\_size})$ .

$\lambda \leftarrow$  ACTOR( $\{(\mathcal{C}, state)\}_{i=1}^{mini\_batch\_size}$ ,  $\Phi, \Theta$ ).

$prob \leftarrow \lambda(\{action_i\}_{i=1}^{mini\_batch\_size})$ .

$advantages \leftarrow \{advantage_i\}_{i=1}^{mini\_batch\_size}$

        Calculate the PPO actor loss with respect to  $\Phi, \Theta$ :

$$L^{actor}(\Phi, \Theta) = \hat{\mathbb{E}}_t [\min(r_t(\Phi, \Theta)advantages, \text{clip}(r_t(\Phi, \Theta), 1 - \epsilon, 1 + \epsilon)advantages)]$$

        where  $r_t(\Phi, \Theta)$  is the probability ratio  $r_t(\Phi, \Theta) = \frac{prob}{prob^{old}}$ .

$\mu_{critic}, \sigma_{critic} \leftarrow V(\{state_i\}_{i=1}^{mini\_batch\_size})$

$returns \leftarrow \{return_i\}_{i=1}^{mini\_batch\_size}$

        Calculate the NLL for the critic with respect to  $\Psi$ :

$$L^{critic}(\Psi) = -\frac{1}{N} \sum_{i=1}^N \log \left( \frac{1}{\sigma_{V_i} \sqrt{2\pi}} \exp \left( -\frac{(returns - \mu_{critic})^2}{2\sigma_{critic}^2} \right) \right)$$

        Update  $\Phi, \Theta$  and  $\Psi$  by minimizing  $L^{actor}$  and  $L^{critic}$ .

Initialize  $\Phi, \Theta, \Psi$ .

**for** epoch = 1 to  $max\_epochs$  **do**

    Empty meta-training dataset  $\mathcal{D}$ .

    Sample  $meta\_batch\_size$  tasks  $\tau$  from  $\mathcal{T}$ .

**for** each sampled  $task$  **do**

        Initialize  $\mathcal{H}$  to zeros.

        ▷ FIFO queue of  $E$  episodes

**for**  $episode \leftarrow 1$  to  $k$  **do**

$\mathcal{D}, \mathcal{H} \leftarrow$  COLLECTANDPROCESSEPIISODE( $task, \Phi, \Theta, \Psi, \mathcal{H}, \mathcal{D}$ )

    UPDATEMODELS( $\mathcal{D}, \Phi, \Theta, \Psi$ )

---

**Algorithm 3** Meta-Testing with ECET using Subroutines from Algorithm 1

---

918 **Input:** Distribution of meta-testing tasks  $\mathcal{T}_{test}$ , history  $\mathcal{H}$  (FIFO queue) to store past transitions,  
919  $num\_eval\_tasks$ ,  $num\_eval\_episodes$ .  
920 **Parameters** from meta-training  
921 Transition( $\varphi_1$ ) and state( $\varphi_1$ ) embedding parameters:  $\Phi = \{\varphi_1, \varphi_2\}$   
922 Intra-Episodic Transformer(**IET**( $\cdot; \Theta_1$ )), Cross-Episodic Transformer(**CET**( $\cdot; \Theta_2$ )), and  
923 policy( $\pi(\cdot; \Theta_3)$ ) parameters:  $\Theta = \{\Theta_1, \Theta_2, \Theta_3\}$ .  
924 **Procedures:** COLLECTEPISODE  
925 Sample  $num\_eval\_tasks$  tasks from  $\mathcal{T}_{test}$ .  
926 **for** each sampled *task* **do**  
927     Initialize  $\mathcal{H}$  to zeros. ▷ FIFO queue of  $E$  episodes  
928     **for** *episode*  $\leftarrow 1$  to  $num\_eval\_episodes$  **do**  
929          $\_, \mathcal{H} \leftarrow$  COLLECTEPISODE(*task*,  $\Phi$ ,  $\Theta$ ,  $\mathcal{H}$ )

---

( $\mathcal{C}$ , *state*, *action*, *reward*, *terminated*, *advantages*, *returns*) to the meta-training dataset  $\mathcal{D}$ . We return the updated history  $\mathcal{H}$  and meta-training dataset  $\mathcal{D}$ .

- *SampleTransitions*: Using the history  $\mathcal{H}$ , we sample a sequence of  $T$  transitions from each episode in  $\mathcal{H}$ , except for the current episode, where we sample a sequence of  $T$  transitions that ends with the *current\_state*. We add these sequences in a set  $\mathcal{C}$  and return it.
- *GenerateTaskRepresentation*: We first transform each transition in the sequences  $T_i$  using the transition embedding  $\varphi_1$ , and then input the embedding into the Intra-Episode Feature Extractor to get  $\mathbf{z}_{seq}^i \leftarrow$  **IET**( $\varphi_1(T_i); \Theta_1$ ) for each episode  $i$ . We store the  $\mathbf{z}_{seq}^i$  output embeddings from **IET** in the set  $\mathcal{C}_{seq}$ . We then input them in the Inter-Episode Feature Extractor to get  $\mathbf{z}_{task} \leftarrow$  **CET**( $\mathcal{C}_{seq}; \Theta_2$ ).  $\mathbf{z}_{task}$  is the task representation we return.

## A.2.4 META-TRAINING

We describe the meta-training process in Algorithm 2, using the subroutines from Algorithm 1. We first initialize the Actor parameters ( $\Phi$ ,  $\Theta$ ) and Critic parameters ( $\Psi$ ).

**Data Collection.** For each epoch until reaching  $max\_epochs$ , we start by initializing an empty meta-training dataset  $\mathcal{D}$ , and sampling  $meta\_batch\_size$  number of tasks  $\tau$  from  $\mathcal{T}$ . For each sampled task, we first initialize the history queue  $\mathcal{H}$  to zeros. This queue has a length of  $E$  episodes, to fit the hierarchical transformer’s memory. For  $k$  episodes we then update the meta-training dataset  $\mathcal{D}$  and history  $\mathcal{H}$  by calling the subroutine *CollectAndProcessEpisode*.

**Model Update.** Once we have collected data in  $\mathcal{D}$  for  $k$  episodes from each sampled task, we use it to update the parameters. As we also show in the subroutine *UpdateModels*, we copy the current  $\Phi$  and  $\Theta$  to  $\Phi^{old}$  and  $\Theta^{old}$ , in order to be able to calculate the probability ratios for the PPO objective for the Actor. We sample a  $mini\_batch\_size$  of tuples ( $\mathcal{C}$ , *state*, *action*, *reward*, *terminated*, *advantage*, *return*) from  $\mathcal{D}$ . Using the *state* values from the mini-batch, as well as the sampled transition sequences  $\mathcal{C}$  for each of these states, we get outputs from the actor using  $\Phi^{old}$  and  $\Theta^{old}$ . As we show in the *Actor* procedure, we use the mini-batch of sampled transition sequences  $\mathcal{C}$  to generate a mini-batch of task representations  $\mathbf{z}_{task}$  as described in the *GenerateTaskRepresentation* subroutine. We embed the mini-batch of states using the state embedding  $\varphi_2$ , and concatenate them with the task representation  $\mathbf{z}_{task}$  for each state in the mini-batch. We then input the resulting mini-batch in the policy  $\pi$  to get the mini-batch of parameters of the action distribution. Since the action distribution is a Gaussian, the policy  $\pi$  outputs the mean  $\mu_{actor}$  and standard deviation  $\sigma_{actor}$ . Since we use  $\Phi^{old}$  and  $\Theta^{old}$ , we name this mini-batch of distribution parameters  $\lambda^{old}$ . We then calculate the probability of the actions in the mini-batch given the distribution parameters  $\lambda^{old}$ , and save the action probabilities in  $prob^{old}$ . We do the same procedure for  $\Phi$  and  $\Theta$ , resulting with the mini-batch of action probabilities  $prob$ . Using these action probabilities, as well as the advantages from the mini-batch, we calculate the PPO Actor loss with respect to  $\Phi$  and  $\Theta$ . As our Value function, similarly to the policy  $\pi$  is also a Gaussian MLP, predicting the mean and standard deviation of a Gaussian distribution for the value of each state input, to calculate the Critic loss, we get the mini-batch of distribution parameters  $\mu_{critic}$  and  $\sigma_{critic}$  as output from the Value function  $V$  using the states in the minibatch as input. As the Critic loss, we calculate the Negative Log Likelihood loss of the returns in the mini-batch using  $\mu_{critic}$  and  $\sigma_{critic}$ .

We update  $\Phi$  and  $\Theta$  by minimizing the PPO Actor Loss  $L^{actor}$ . We update  $\Psi$  by minimizing the NLL Critic Loss  $L^{critic}$ .

We repeat this process until we go over all tuples in the meta-training dataset  $\mathcal{D}$ . Once  $\mathcal{D}$  has been used for updates, we empty it, and return to sampling tasks and collecting episodes from them again, and updating the parameters based on the new dataset, and so on, until *max\_epochs* is reached.

### A.2.5 META-TESTING

We show the meta-test algorithm in Algorithm 3, using the *CollectEpisode* subroutine from Algorithm 1. We use a distribution of meta-testing tasks  $\mathcal{T}_{test}$  and a FIFO queue as history  $\mathcal{H}$  for the past episodes as input. We also set the *num\_eval\_tasks* and *num\_eval\_episodes* parameters which determine the number of tasks sampled from  $\mathcal{T}_{test}$  and the number of episodes performed in each task, respectively. We use the Actor parameters ( $\Phi, \Theta$ ) which we got from meta-training. We first sampled *num\_eval\_tasks* from  $\mathcal{T}_{test}$ . For each sampled task, we initialize the history  $\mathcal{H}$  to zeros, and then collect *num\_eval\_episodes* and update  $\mathcal{H}$  from each collected episode.

### A.3 ECET HYPERPARAMETERS

Table 1: Hyperparameters Across Benchmarks

Hyperparameter	Benchmark		
	MuJoCo	MetaWorld	Maniskill
Number of Transitions (T)	5		
Number of Episodes (E)	2	25	
Number of Layers for IET	2		
Number of Heads for IET	16	4	
Number of Layers for CET	2		
Number of Heads for CET	16	4	
Minibatch Size	256	32	
Policy Learning Rate	3e-5	5e-5	
Critic Learning Rate	3e-5	5e-5	

### A.4 BENCHMARK DETAILS

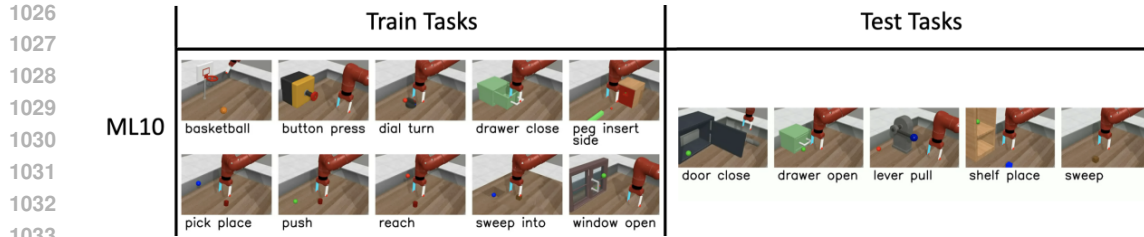


Figure 9: Meta-world Tasks in the M1 with training/test Split. Figure taken from (Yu et al., 2019)

**ML1** provides one main task, with parametric variation in object and goal positions. ML1 enables evaluation in the few-shot adaptation setting. In particular, adaptations to parametric variations within one task. Figure 9 gives an example of meta-training and testing for the pick in place task, where only the test-task target position is different from those observed during training.

**ML10** provides 10 different training tasks and 5 others to test the adaptation performance of meta-learning methods. This mode evaluates few-shot adaptation to 5 new test tasks, after training on 10 different tasks, see Figure 10. The test tasks are all related to the training tasks, but unseen operations have to be performed (such as opening a drawer at test time instead of closing it as learned in meta-training).

**ML45** provides 45 highly different training tasks and 5 others to test the adaptation performance of meta-learning methods. This mode evaluates few-shot adaptation to 5 new test tasks, after training on 45 different tasks. Figure 11 shows all 50 tasks.



1034 Figure 10: Meta-world Tasks in the M10 protocol with training/test Split. Figure taken from (Yu  
 1035 et al., 2019)



1055 Figure 11: Meta-world Tasks in the M45 protocol with training/test Split. Figure taken from (Yu  
 1056 et al., 2019)

## 1059 A.5 BASELINES

1061 **MAML PPO** is the application of the MAML meta-learning algorithm for meta-RL (Finn et al.,  
 1062 2017). The inner-loop algorithm is PPO, whereas the outer-loop algorithm is the MAML algorithm.  
 1063 We use the implementation provided in the *garage* repo (garage contributors, 2019).  
 1064

1065 **RL<sup>2</sup>** uses an RNN in the inner loop as described in Section A.1, and TRPO in the outer loop (Duan  
 1066 et al., 2016). We use the implementation provided in the *garage* repo (garage contributors, 2019),  
 1067 which replaces TRPO with PPO.

1068 **PEARL** is an off-policy meta-RL algorithm, in contrast to the other methods we compare  
 1069 with (Rakelly et al., 2019). We use the implementation provided in the *garage* repo (garage contribu-  
 1070 tors, 2019).

1071 **VariBAD** uses a VAE to learn a distribution of task features that increase the state space (Zintgraf  
 1072 et al., 2020). In contrast to PEARL, VariBAD uses PPO in the outer loop, which is an on-policy  
 1073 meta-RL algorithm. We use the implementation provided by the authors.

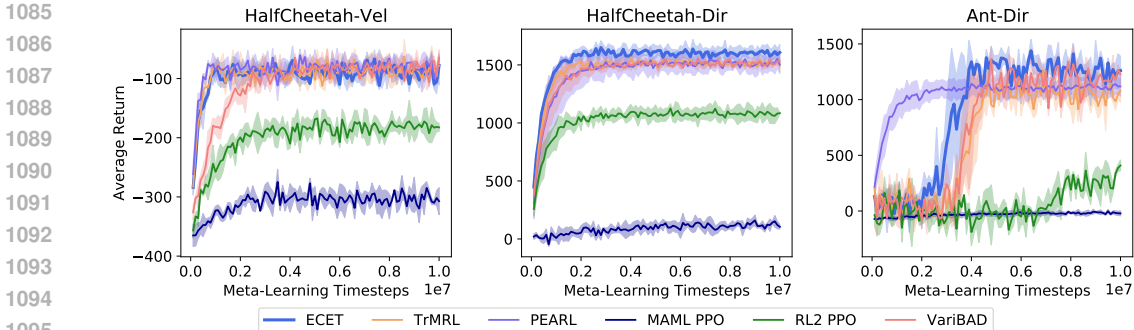
1074 **TrMRL** (Melo, 2022) uses a transformer in the inner loop of the RL<sup>2</sup> algorithm, similar to our  
 1075 proposed method. However, it takes a sequence of the 5 most recent transitions as input to the  
 1076 transformer. We use the implementation and pipeline as provided by the authors.  
 1077

1078 **AMAGO** (Grigsby et al., 2024a) also uses a transformer, but trains it using an off-policy RL  
 1079 approach. It takes long sequences spanning multiple episodes as input to the transformer. We use the  
 implementation provided by the authors.

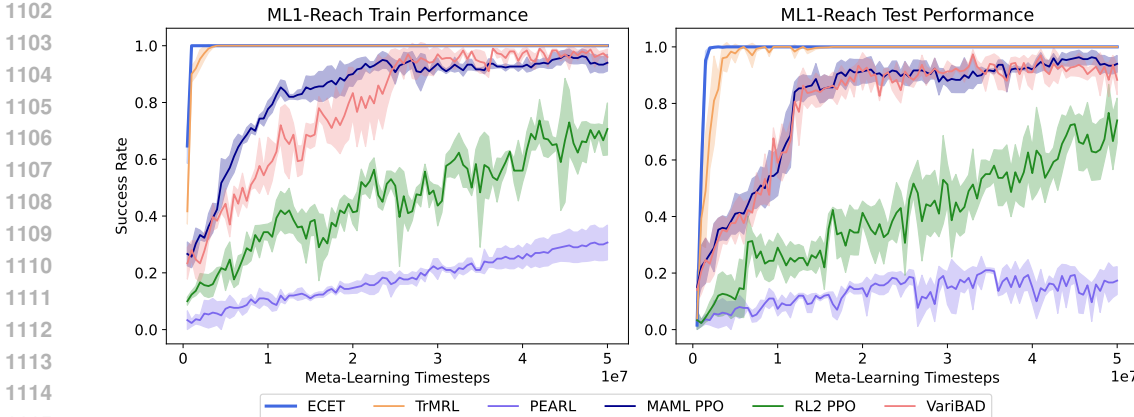
1080  
1081

A.6 ADDITIONAL RESULTS

1082 In this section, we show additional results for our evaluations on the MuJoCo, MetaWorld, and  
1083 ManiSkill benchmarks.



1097 Figure 12: Meta-Train performance in terms of average return of ECET, TrMRL, PEARL, MAML  
1098 PPO, RL2 PPO and VariBAD on the MuJoCo benchmark for training on parametric variations of the  
1099 *HalfCheetah* and *Ant* tasks.



1117 Figure 13: Meta-Train and Test performance in terms of average success rate of ECET, TrMRL,  
1118 PEARL, MAML PPO, RL2 PPO and VariBAD on the ML1 benchmark for training(left) and test-  
1119 ing(right) on parametric variations of the *Reach* task.

1120  
1121  
1122 A.6.1 VISUALIZATION OF EMBEDDINGS

1123 We show the visualization of the embeddings learned by ECET and TrMRL for the ML10 Training  
1124 tasks in Figure 14. Furthermore, we show color maps of the inter-task distance matrix of the  
1125 embeddings for TrMRL(left) and ECET(right) in Figure 15. There, it is more easily visible that the  
1126 inter-task distance of the embeddings is much bigger for ECET than TrMRL, suggesting a better task  
1127 differentiation.  
1128

1129 A.6.2 AVERAGE RANKING PLOTS

1130  
1131 To focus on comparing the methods based on their robustness in relative performance to each other,  
1132 in Figures 16 to 18 we provide plots showing the average ranks of the methods as a performance  
1133 metric. These plots show that ECET robustly achieves the highest relative performance (and thus the  
lowest rank) across the meta-learning benchmarks we consider.

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

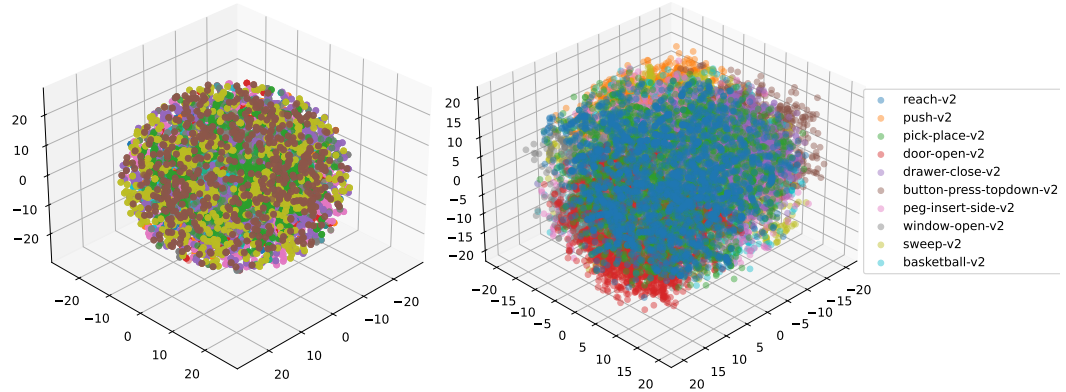


Figure 14: T-SNE plots of the output embeddings for TrMRL(left) and ECET(right) for the tasks in the training set of the ML10 benchmark of Meta-World.

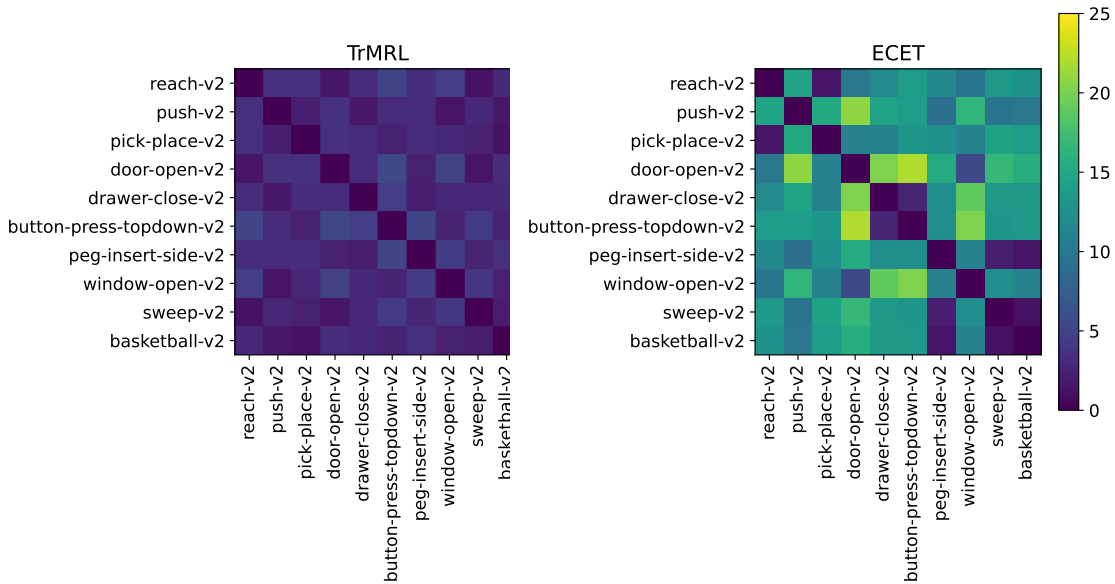


Figure 15: Color maps of the inter-task distance matrix of the embeddings for TrMRL(left) and ECET(right) for the tasks in the training set of the ML10 benchmark of Meta-World.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

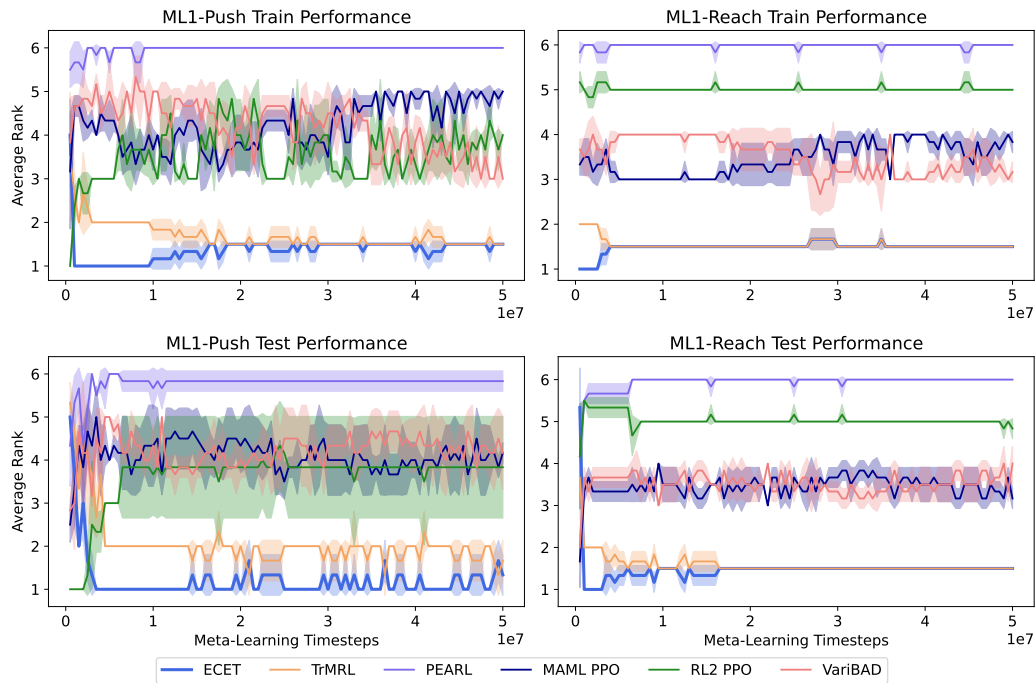


Figure 16: Average ranking plots of Meta-Train and Test performance for all methods on the ML1 *Push*(left) and *Reach*(right). Lower ranks indicate better performance.

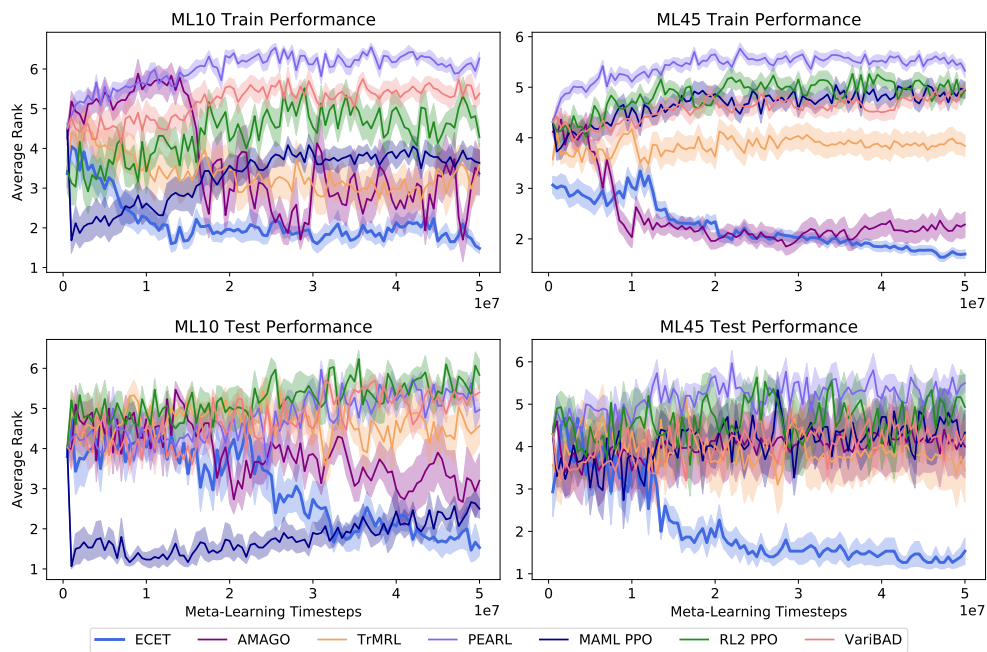


Figure 17: Average ranking plots of Meta-Train and Test performance for all methods on the ML10(left) and ML45(right) with disjoint train and test tasks. Lower ranks indicate better performance. The corresponding success rate plot can be found in Figure 6.

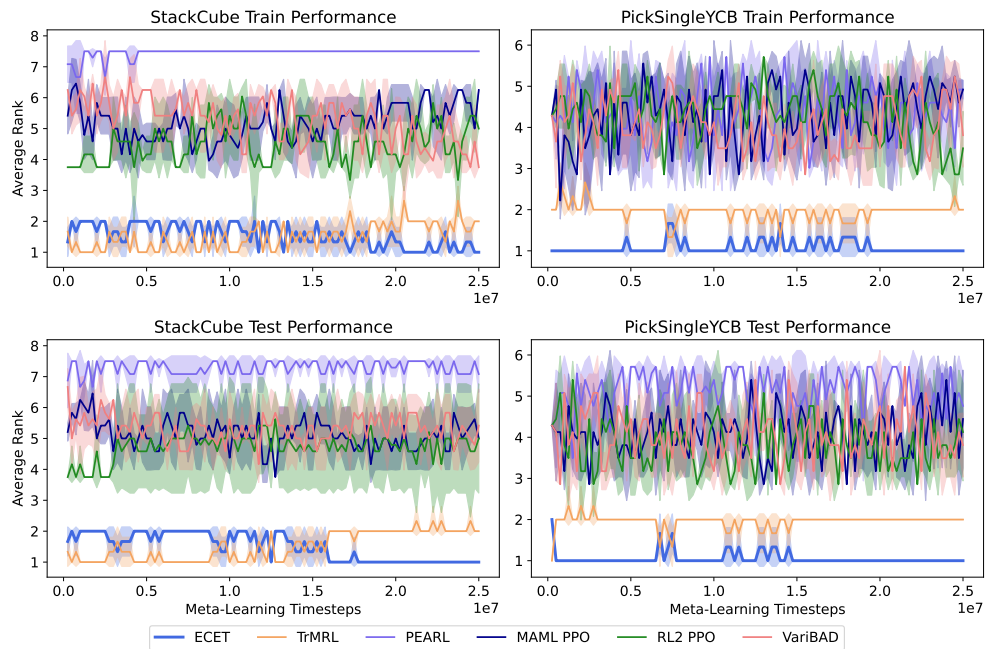


Figure 18: Average ranking plots of Meta-Train and Test performance for all methods on the StackCube(left) and PickSingleYCB(right) tasks with disjoint train and test tasks. Lower ranks indicate better performance.

### A.6.3 PERFORMANCE WITH RESPECT TO TIME

We show performance of the methods in terms of average success rate with respect to time in GPU hours in Figures 19 to 22. These plots show that in most cases ECET achieves the best anytime performance across the meta-learning benchmarks we consider.

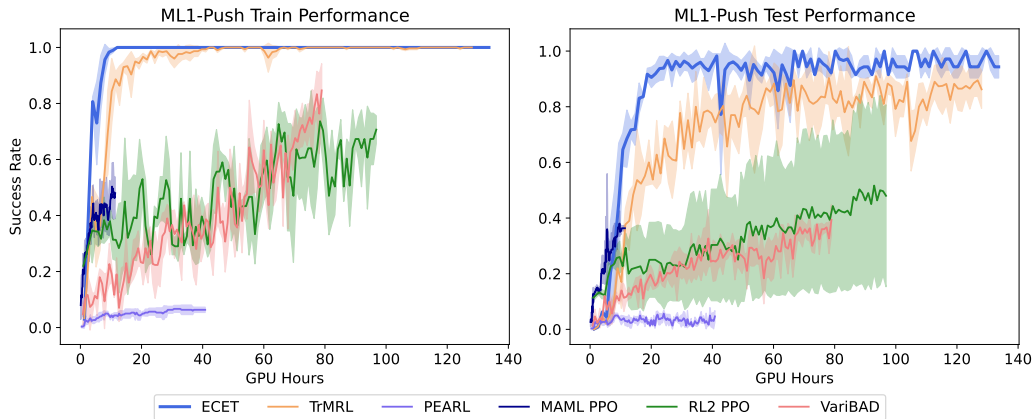


Figure 19: Meta-Train and Test performance in terms of average success rate with respect to time in GPU hours for all methods on the ML1 *Push* task.

### A.6.4 EVALUATION ACROSS EPISODES IN TEST TASKS

To assess the test performance of ECET, we perform rollouts for 50 episodes on test tasks and compare the episodic return for each episode to the baselines.



1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

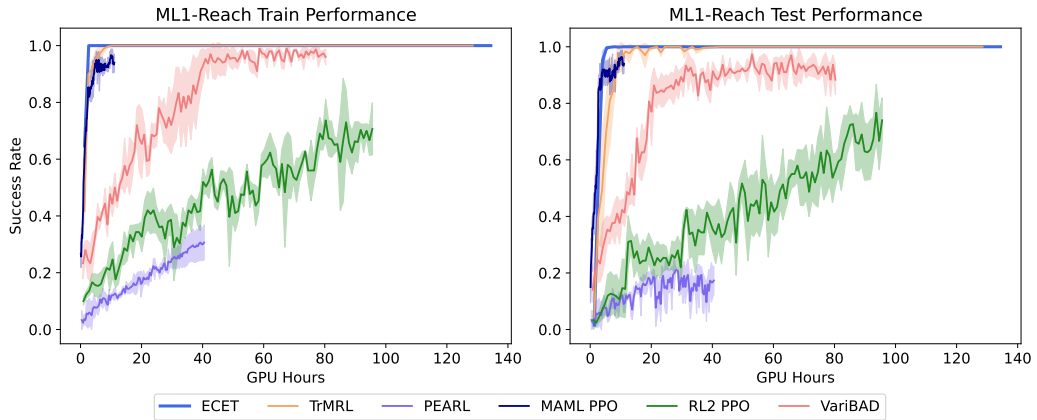


Figure 20: Meta-Train and Test performance in terms of average success rate with respect to time in GPU hours for all methods on the ML1 *Reach* task.

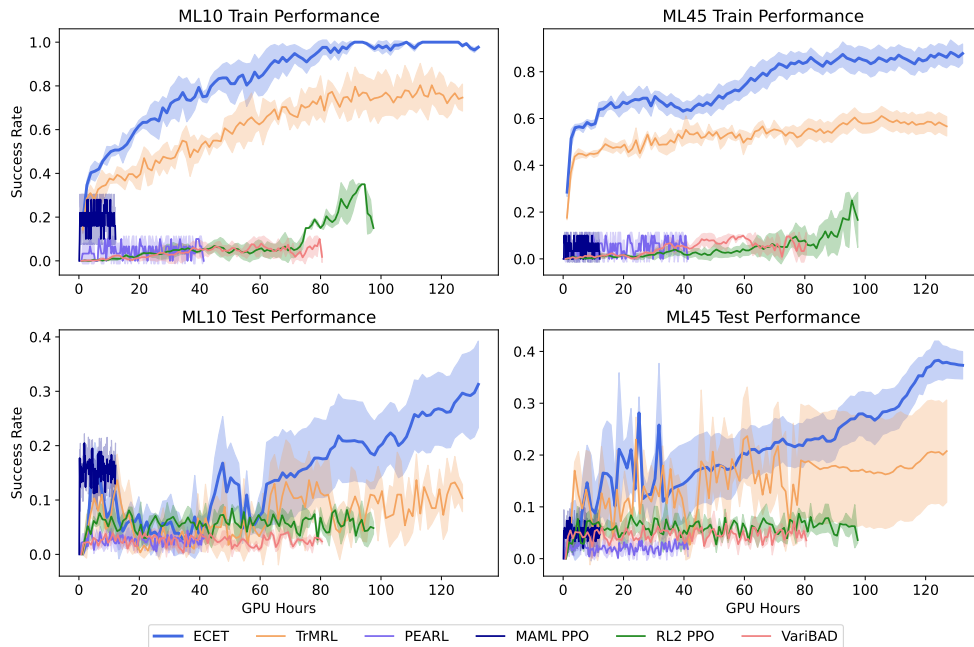
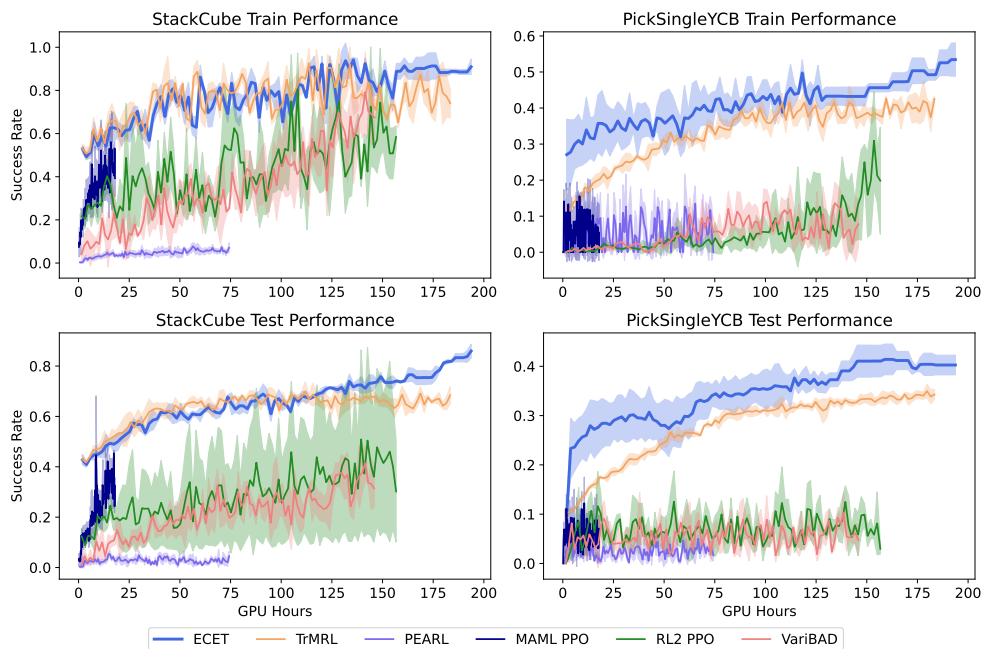


Figure 21: Meta-Train and Test performance in terms of average success rate with respect to time in GPU hours for all methods on the ML10(left) and ML45(right) with disjoint train and test tasks.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371

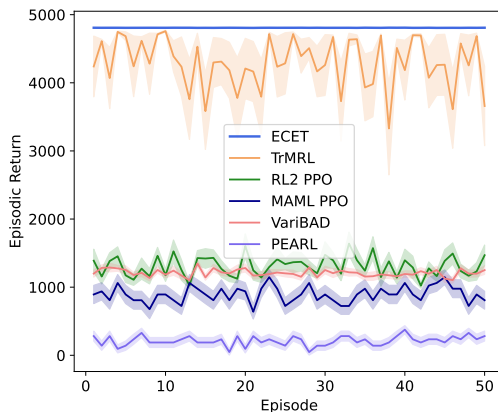


1372 Figure 22: Meta-Train and Test performance in terms of average success rate with respect to time in  
1373 GPU hours for all methods on the StackCube(left) and PickSingleYCB(right) tasks with disjoint train  
1374 and test tasks.

1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383

**ML1** We sample 10 test tasks for *Push-v2* and perform rollouts of 50 episodes with each of the resulting models from the meta-training using 5 seeds for each method. For each episode we plot the mean and standard error. Figure 23 shows the resulting plots, where it is visible that ECET generalizes well across parametric variations it has not seen in meta-training, by achieving the highest episodic returns consistently across episodes. This consistency since the first episode shows the independence of ECET to the experience available for the current task. This happens because during meta-training, we make sure to train on transitions where we don't have the full history of 25 episodes available.

1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399

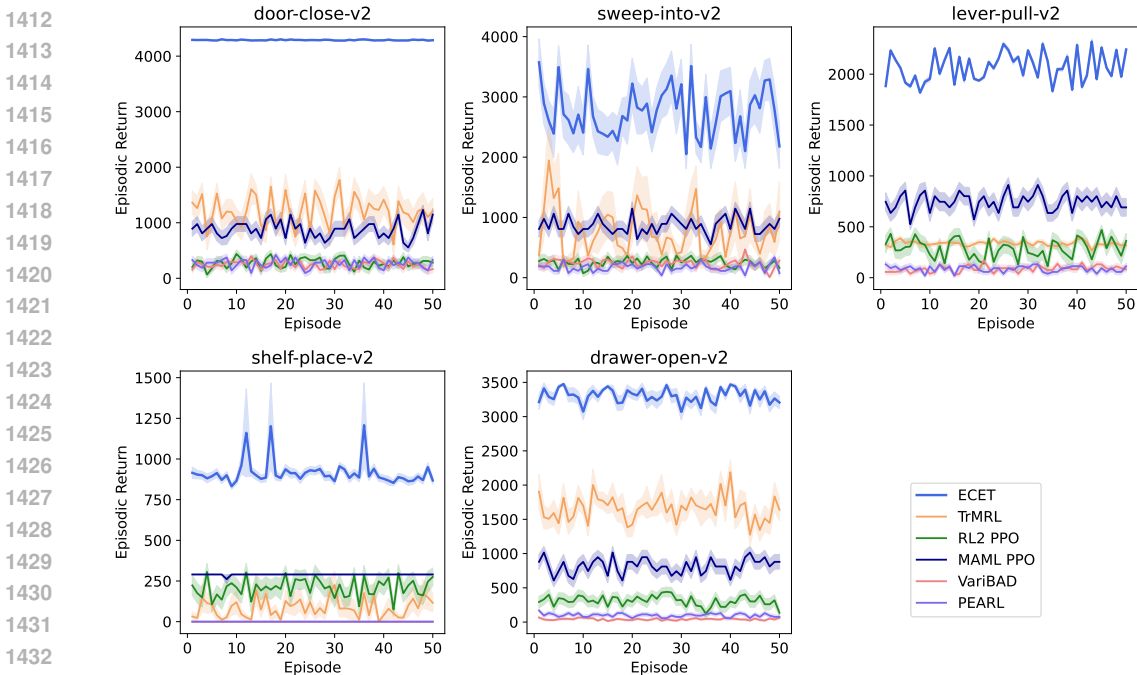


1400 Figure 23: Episodic return across episodes on test tasks for all methods on the ML1 *Push* Benchmark.

1401  
1402  
1403

**ML10** We sample 10 parametric variations for each of the tasks in the test set of ML10 (*Door Close, Sweep Into, Lever Pull, Shelf Place, Drawer Open*) and perform rollouts of 50 episodes with each of the resulting models from the meta-training using 5 seeds for each method. For each episode

1404 we plot the mean and standard error. Figure 24 shows the resulting plots. It is visible that ECET  
 1405 generalizes well across tasks it has not seen in meta-training, by achieving the highest episodic  
 1406 returns consistently across episodes. However, for some tasks, there is more variance in performance  
 1407 across seeds, as well as across episodes, which indicates more difficulty to adapt. Similarly as for  
 1408 the ML1 benchmark, there is no drop in performance for ECET in the first episodes where we have  
 1409 less observations, further demonstrating the independence of ECET to the experience available for  
 1410 the current task.



1434 Figure 24: Episodic return across episodes on each test task for all methods on the ML10 Benchmark.

1436 **ML45** We sample 10 parametric variations for each of the tasks in the test set of ML10 (*Box Close*,  
 1437 *Bin Picking*, *Door Unlock*, *Hand Insert*, *Door Lock*) and perform rollouts of 50 episodes with each  
 1438 of the resulting models from the meta-training using 5 seeds for each method. For each episode we  
 1439 plot the mean and standard error. Figure 25 shows the resulting plots. Compared to the baselines,  
 1440 ECET manages to get higher episodic returns on tasks it has not seen in meta-training, similarly to  
 1441 the ML1 and ML10 benchmark. However, for ML45, the variance in performance is higher across  
 1442 episodes, which indicates more difficulty to adapt. For the *Door Unlock* task, it is noticeable that  
 1443 TrMRL manages to perform similarly to ECET. The other baselines fail to reach high episodic return.  
 1444 Similarly as for the ML1 benchmark, there is no drop in performance for ECET in the first episodes  
 1445 where we have less observations, further demonstrating the independence of ECET to the experience  
 1446 available for the current task.

1447  
 1448 **A.6.5 ABLATIONS**

1449 To investigate the behavior of ECET and the impact of its components on performance, we ablate the  
 1450 number of episodes  $E$  to keep in memory and the number of transitions  $T$  in the sampled sequence in  
 1451 Figure 28. The configuration that performs best, and the one we use in our experiments is  $E = 25$ ,  
 1452  $T = 5$ . We ablate the sampling strategy (sampling one transition sequence per episode vs. sampling  
 1453 multiple sequences randomly from any episode) in the top plot in Figure 29, to show the usefulness  
 1454 of sampling one transition sequence per episode. In contrast to TrMRL, which only uses the output  
 1455 of the transformer as input to the policy  $\pi$ , we concatenate a linear transformation of the current state  
 1456 as input to the policy. We ablate the performance of ECET for these two settings of policy input in  
 1457 the middle plot in Figure 29. In the bottom plot, we ablate the number of encoder blocks for IET  
 and CET. To differentiate the contributions of the context sampling strategy from the architecture in

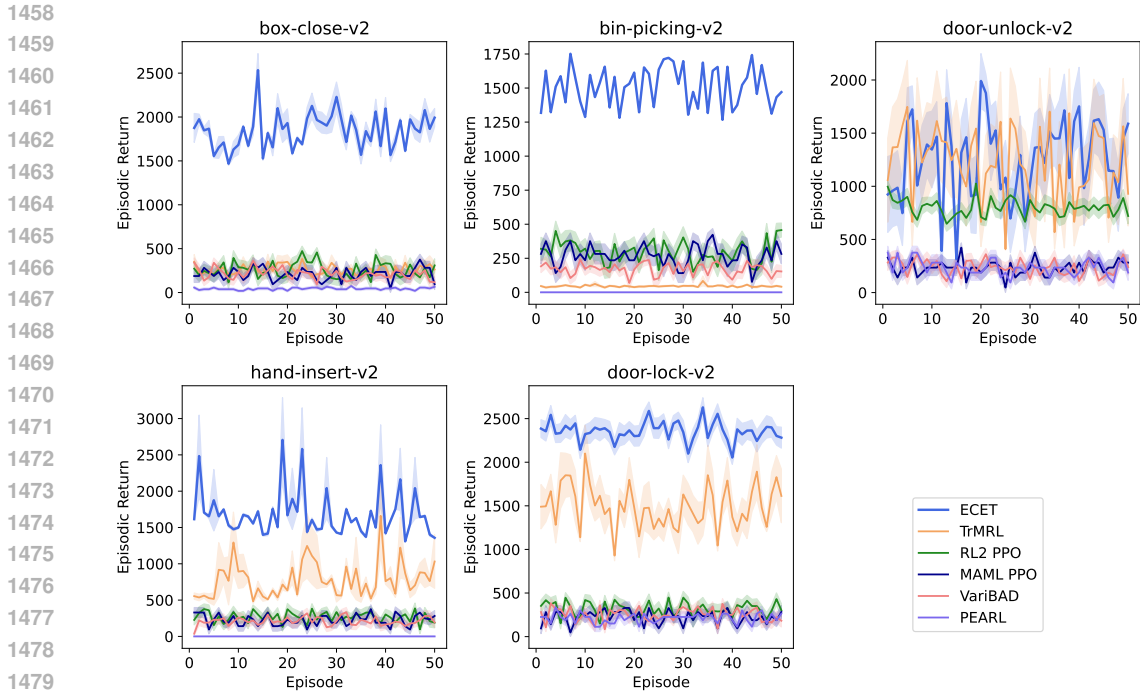


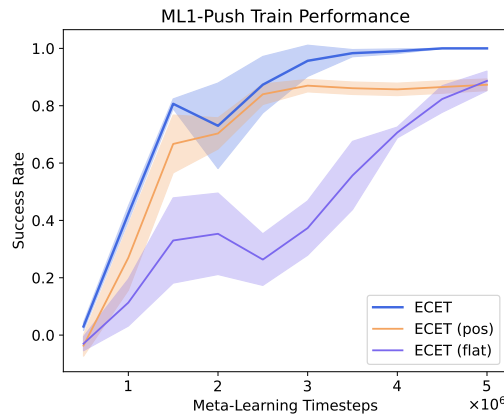
Figure 25: Episodic return across episodes on each test task for all methods on the ML45 Benchmark.

ECET, in Figure 26 we compare ECET to a variant using a positional encoding for CET (ECET (pos)) and one using a flat architecture with only IET (ECET (flat)). This version processes context that is sampled exactly the same as in ECET, specifically sequences of length  $S$  are randomly sampled from  $E$  episodes and then put through the transformer. ECET (flat) yields inferior performance, underscoring the challenges a non-hierarchical approach faces in capturing comprehensive task representations. Finally, in Figure 27 we also compare ECET to TrMRL, and TrMRL (extended context) that processes sequences of the same length as the total number of transitions that ECET processes, namely sequences of length  $E \times S$ . We extend the context of TrMRL to investigate the effect that the context length has on the performance improvements we see with ECET. The low success rate of TrMRL (extended context) shows that extending the context length of TrMRL makes the meta-training less efficient, as much more training data is needed to extract useful task meta-features from these long sequences. Our benchmark configuration for these experiments stands at  $E = 25$ ,  $T = 5$ , incorporates two encoder blocks for both IET (with positional encoding) and CET (without positional encoding), and uses a linear transformation  $\varphi_2$  of the state as additional input to the policy  $\pi$ . We do all our ablations in the ML1 Benchmark, using parametric variations of the *Push* task.

Figure 26 shows the performance of our default configuration of the method, ECET, compared to using a positional encoding for CET, ECET (pos), and the simplification to a flat architecture akin to TrMRL by solely employing IET, ECET (flat). We show the average and standard error of the success rate across Meta-Learning Timesteps. ECET (flat) performs worst, indicating the difficulty of a flat architecture to capture a better representation of the task, instead of focusing on local characteristics of the dynamics. ECET (pos) has a similar rate of improvement to ECET in the beginning, but ends up converging to a suboptimal representation of the task, and thus its success rate does not improve above  $\sim 0.8$  after  $2.5 \cdot 10^6$  Timesteps. This indicates that the positional encoding in CET induces an ordering in the  $\mathbf{z}_{seq}$  representations which is misleading, as episodes don't naturally have a causal ordering, thus imposing an ordering on the representations we get from the sequences of transitions sampled ends up hurting performance.

Figure 27 shows the performance of our default configuration of the method, ECET, compared to TrMRL that processes sequences of length 5, and TrMRL (extended context) that processes sequences of length 125. We extend the context of TrMRL to investigate the effect that the context length has on the performance improvements we see with ECET. The low success rate of TrMRL (extended

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526

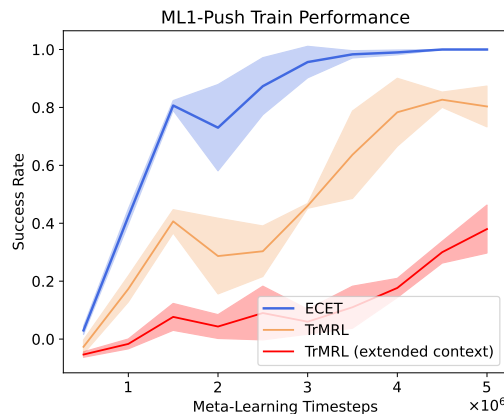


1527 Figure 26: Investigating the meta-training performance impact on ECET of the inter-episode feature  
1528 extractor **CET** and the positional encoding on the sequence embeddings  $\mathbf{z}_{seq}$  before being input into  
1529 **CET** in the *ML1 Push* Benchmark. ECET shows the performance of the configuration we propose,  
1530 with **IET** using positional encoding, and **CET** without positional encoding. ECET (pos) shows  
1531 the performance of the configuration with **IET** and **CET** both using positional encoding. ECET  
1532 (flat) shows the performance of the configuration with no hierarchy of transformers, only **IET** using  
1533 positional encoding, taking the same information as input as ECET.

1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541

context) shows that extending the context length of TrMRL makes the meta-training less efficient, as  
as much more training is needed to extract useful task meta-features from these long sequences. Thus,  
the ECET architecture of encoding shorter sequences from different episodes with **IET**, and then  
processing those encodings with **CET** is crucial in ensuring better anytime performance than TrMRL.

1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554



1555 Figure 27: Performance of ECET, TrMRL with its standard sequence length of 5, and TrMRL with  
1556 extended context (processing sequences of length 125 to equal the context length of ECET).

1557  
1558  
1559  
1560

To investigate the performance sensitivity of ECET on the number of episodes ( $E$ ) kept in memory  
and the samples sequence length ( $T$ ), we provide ablation results in Table 2.

1561  
1562  
1563  
1564  
1565

Figure 28 shows the performance of ECET, when varying key parameters and structural elements:  
the number of episodes ( $E$ ) in memory, the length of sampled transition sequences ( $T$ ). The top  
plot compares the impact of varying  $E = \{5, 10, 15, 20, 25, 50\}$  when  $T = 5$ . It is noticeable that  
increasing  $E$  to 50 results in suboptimal performance, as the configuration does not manage to  
reach a success rate of 1 compared to the other configurations. On the other side of the spectrum,  
decreasing  $E$  to 5 results in a slower improvement. The configuration that we use,  $E = 25$ , shows

$T \backslash E$	5	10	15	20	25	50
5	$0.997 \pm 0.005$	$1.000 \pm 0.000$	$0.997 \pm 0.005$	$1.000 \pm 0.000$	$1.000 \pm 0.000$	$0.909 \pm 0.058$
25	$0.997 \pm 0.005$	$1.000 \pm 0.000$	$1.000 \pm 0.000$	$0.993 \pm 0.009$	$0.924 \pm 0.053$	$1.000 \pm 0.000$
50	$0.963 \pm 0.039$	$0.960 \pm 0.050$	$0.993 \pm 0.005$	$0.827 \pm 0.098$	$0.990 \pm 0.014$	$0.893 \pm 0.041$

Table 2: ECET’s final meta-training success-rate on **ML1 Push** when varying the number of episodes ( $E$ ) and the sequence length ( $T$ ).

the fastest improvement and better performance across the Meta-Learning Timesteps. The middle plot compares the impact of varying  $E = \{5, 10, 15, 20, 25, 50\}$  when  $T = 25$ . It is noticeable that all of the possible configurations are slower to improve than the configurations using  $T = 5$ , and that is because with higher sequence length, ECET needs much more data and training to arrive at a good representation  $\mathbf{z}_{seq}$  and  $\mathbf{z}_{task}$ . This claim is further supported by the bottom figure, showing even slower improvement for configurations with  $E = \{5, 10, 15, 20, 25, 50\}$  when  $T = 50$ .

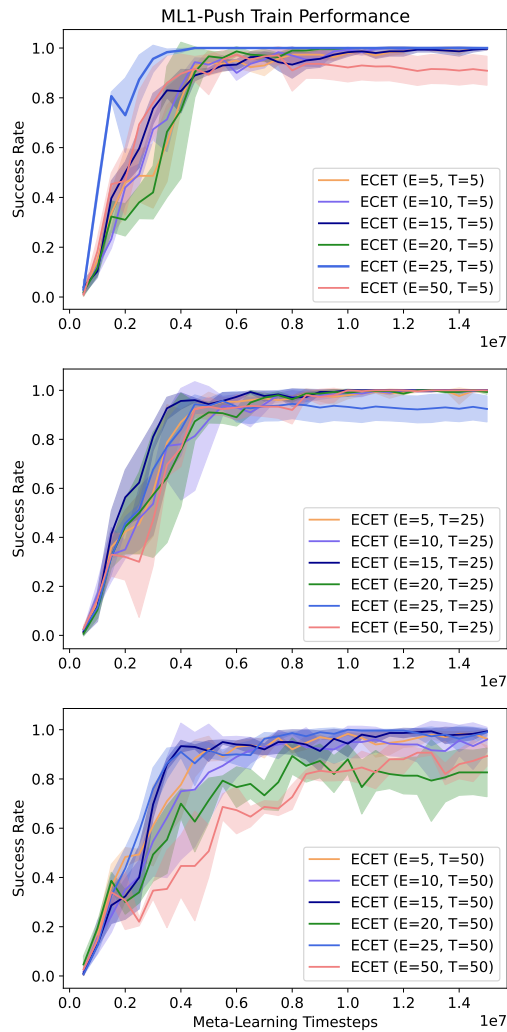


Figure 28: Investigating the meta-training performance impact of the sequence length  $T$  and number of episodes  $E$  for ECET in ML1 *Push*. The top plot shows the performance of sequence lengths  $T = 5$  with different episode lengths  $E \in \{5, 10, 15, 20, 25, 50\}$ . The middle plot shows the performance of  $S = 25$  with  $E \in \{5, 10, 15, 20, 25, 50\}$ . For the bottom plot,  $T = 50$ .

Figure 29 top shows ablations of our sampling strategy—selecting a singular transition sequence per episode to a more randomized approach drawing multiple sequences from any episode. All of the shown configurations use the latter strategy, as in Figure 28 and everywhere else we use the former strategy as our default. It is noticeable that when sampling sequences from any randomly chosen episode, the performance improvement is slower. However, for the  $E = 50, T = 5$  configuration, the randomization of sampling helps ECET not get stuck in suboptimal regions, and thus it manages to reach a success rate of 1.0 compared to the same configuration using the non-randomized strategy of sampling with respect to episodes, which we see in Figure 28. Despite this, all the configurations using the randomized strategy have slower improvement compared to our default configuration using the non-randomized strategy. The middle plot shows the performance comparison of augmenting the policy  $\pi$  input with a linear transformation of the state(ECET) versus relying exclusively on the transformer’s output(ECET\*). We notice a significant improvement from getting the linear transformation of the state as input in the policy. The bottom plot shows the influence of varying the number of encoder blocks within the architecture. Although it has least parameters compared to the other configurations, ECET(2, 2) performs best, having the fastest rate of improvement and also best final performance.

#### A.6.6 RANK AND SUCCESS RATE PLOTS PER TASK

Through investigating the individual plots for each task in terms of average success rate in Figures 30-34, the difference in the level of difficulty for the tasks can be demonstrated, noticing that e.g. in Figure 31 the *Door Close* task is easier to successfully perform after meta-training, reaching success rates of approximately 0.6, compared to the *Shelf Place* task, for which all the methods fail to reach a success rate higher than 0.12. This is also the case for the performance per task for the test tasks in ML45, which we show in Figure 34. For the *Bin Picking* task, none of the baselines can successfully finish even one episode, except for VariBAD with success rate of approximately 0.05 (ECET being more robust, reaching 0.3). On the other hand, for the *Door Unlock* task, RL2 PPO manages to reach success rates of approximately 0.4, and TrMRL reaches 0.6, being one of the few cases where ECET underperforms by reaching a success rate of 0.3.

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

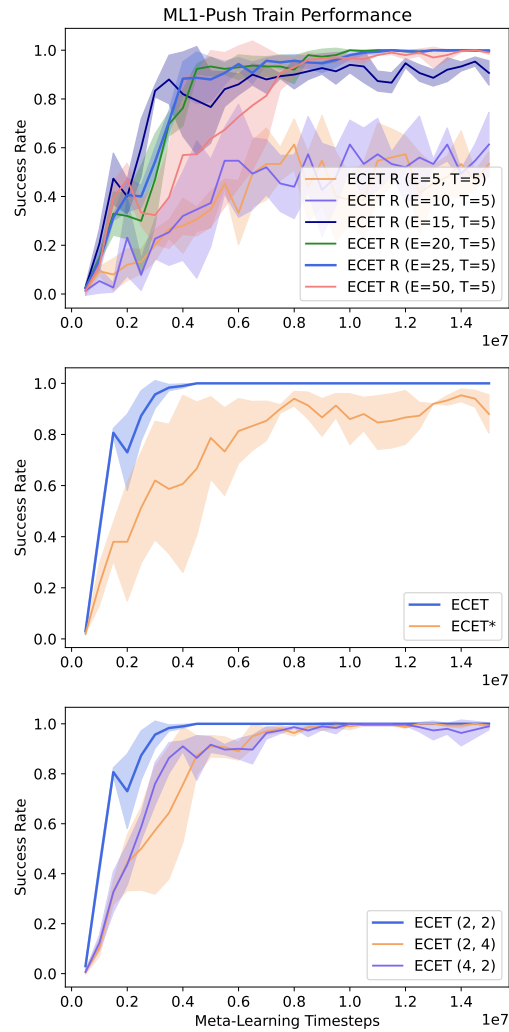


Figure 29: Investigating the meta-training performance impact of the sampling strategy for the sequence for ECET in *ML1 Push*. The top plot shows the performance of ECET when sampling  $E$  sequences from random episodes (i.e. more than one sequence per episode). The middle plot shows the performance impact of concatenating a linear transformation of the state (ECET) vs using only the output of the transformer (ECET\*) as the input for the policy. The bottom plot shows the performance impact of the number of encoder blocks, comparing configurations of 2 encoder blocks for each IET and CET, 2 encoder blocks for IET and 4 for CET, and 4 encoder blocks for IET and 2 for CET.



1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

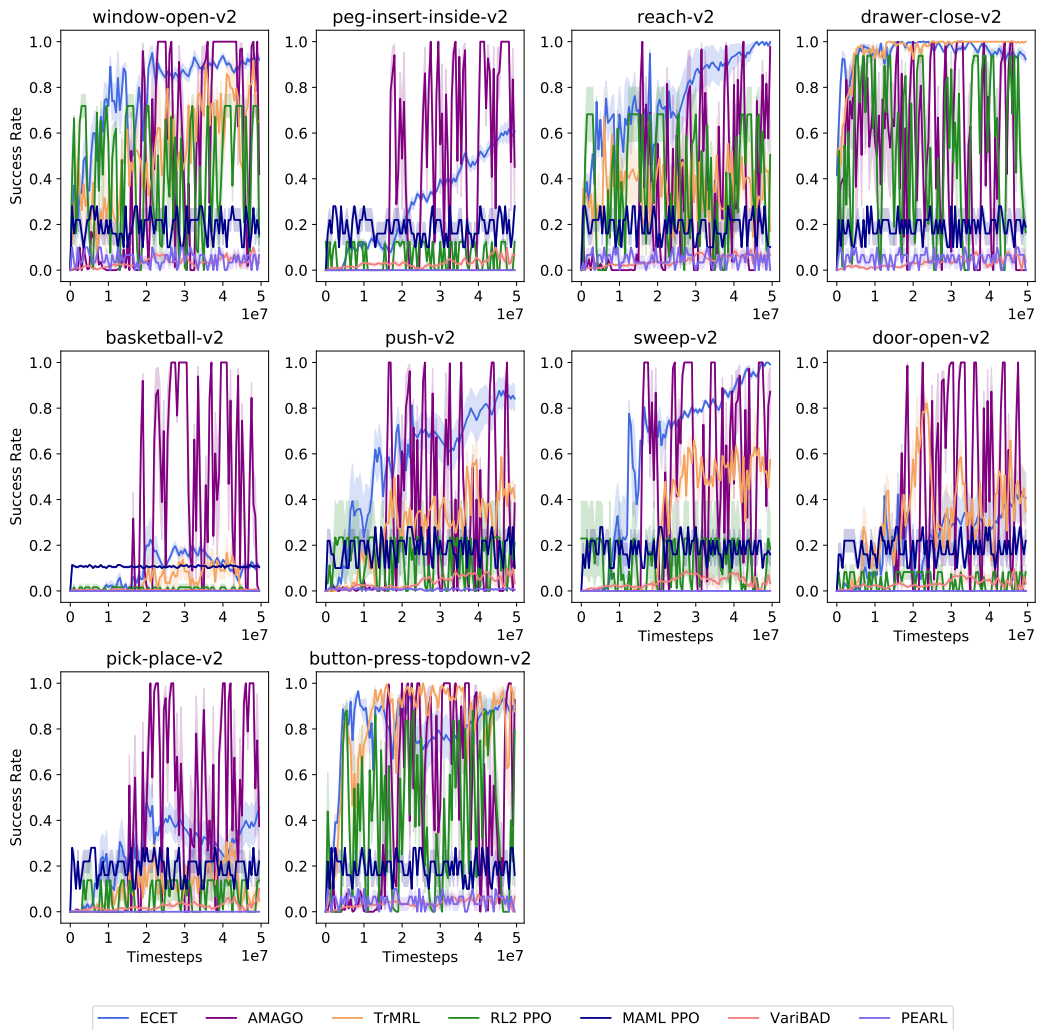


Figure 30: Success rate per train task for the ML10 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD.

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

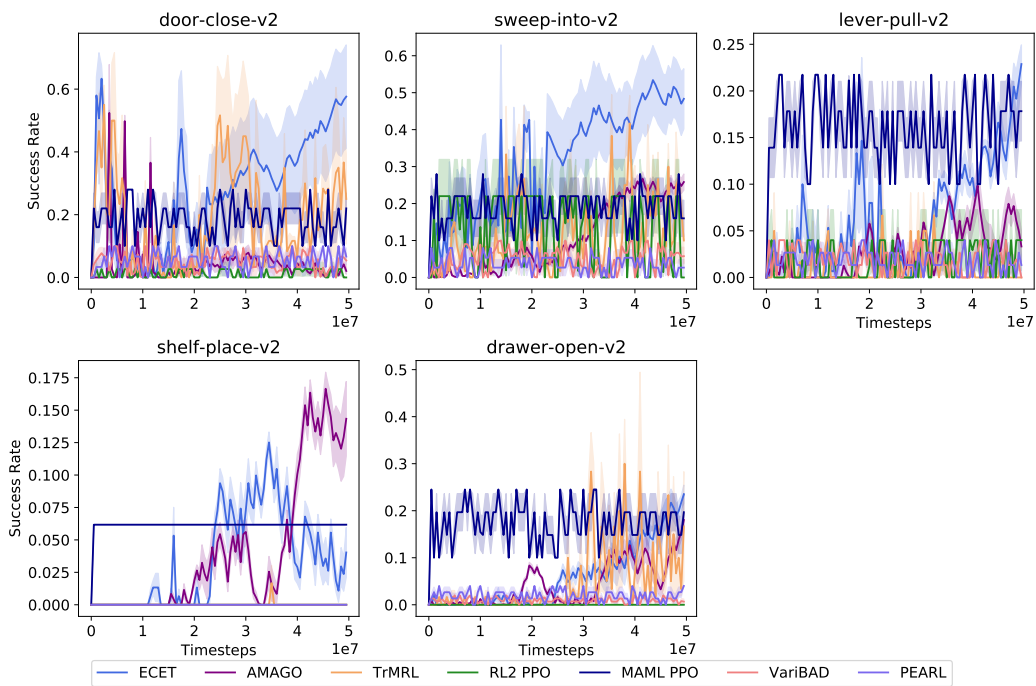


Figure 31: Success rate per test task for the ML10 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD.

1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889

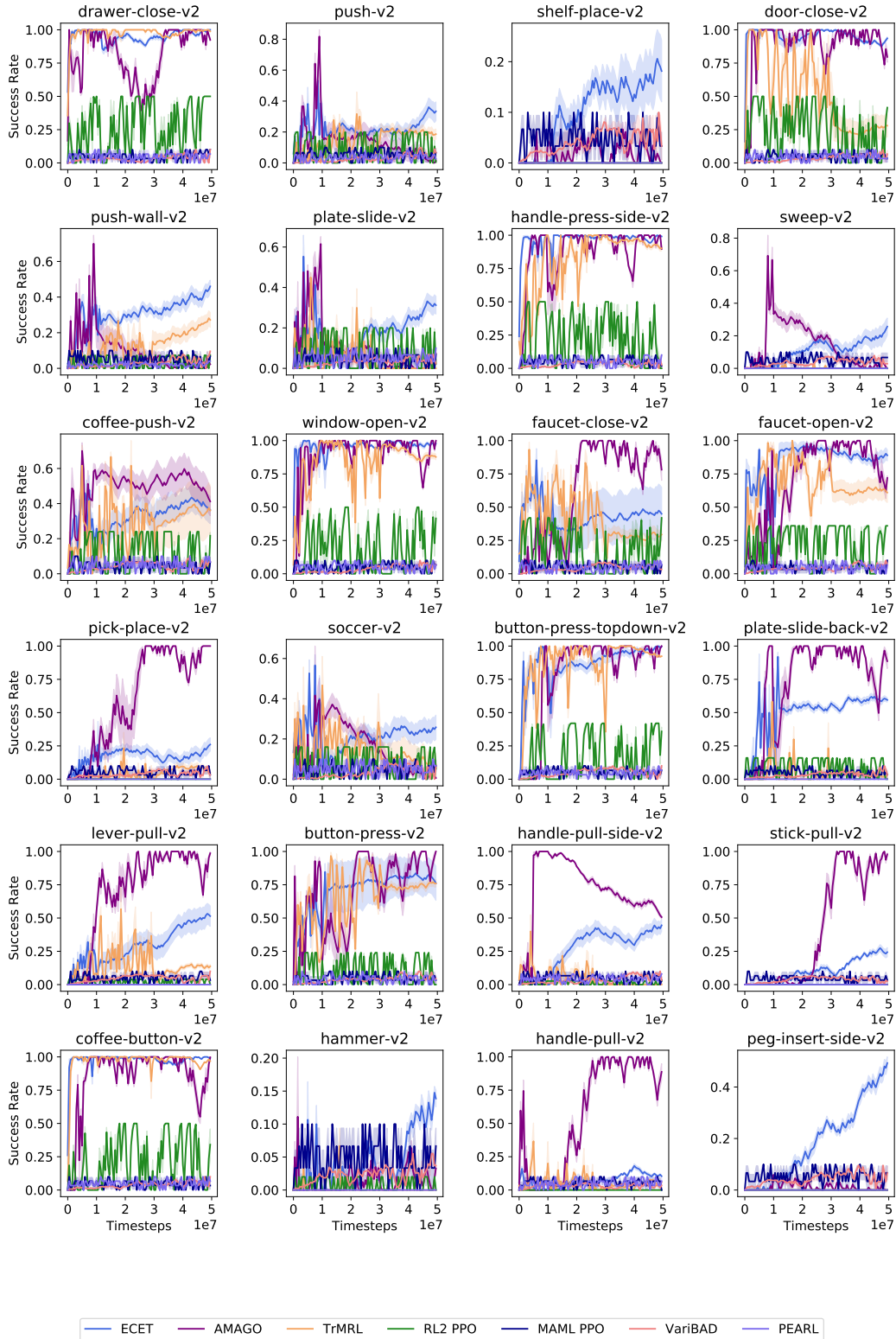


Figure 32: Success rate per train task for the ML45 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD (Part 1).

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

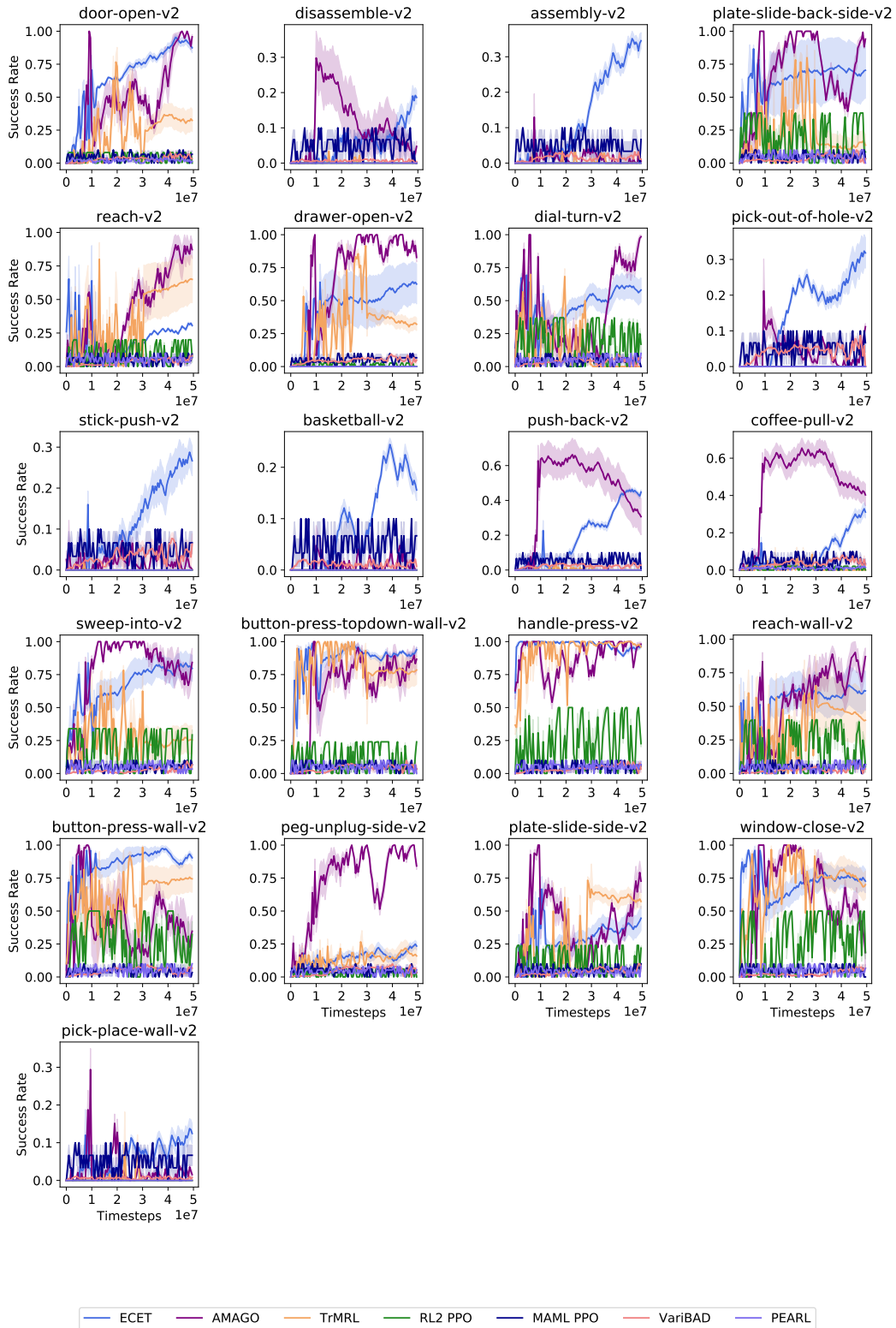


Figure 33: Success rate per train task for the ML45 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD (Part 2).

1944  
 1945  
 1946  
 1947  
 1948  
 1949  
 1950  
 1951  
 1952  
 1953  
 1954  
 1955  
 1956  
 1957  
 1958  
 1959  
 1960  
 1961  
 1962  
 1963  
 1964  
 1965  
 1966  
 1967  
 1968  
 1969  
 1970  
 1971  
 1972  
 1973  
 1974  
 1975  
 1976  
 1977  
 1978  
 1979  
 1980  
 1981  
 1982  
 1983  
 1984  
 1985  
 1986  
 1987  
 1988  
 1989  
 1990  
 1991  
 1992  
 1993  
 1994  
 1995  
 1996  
 1997

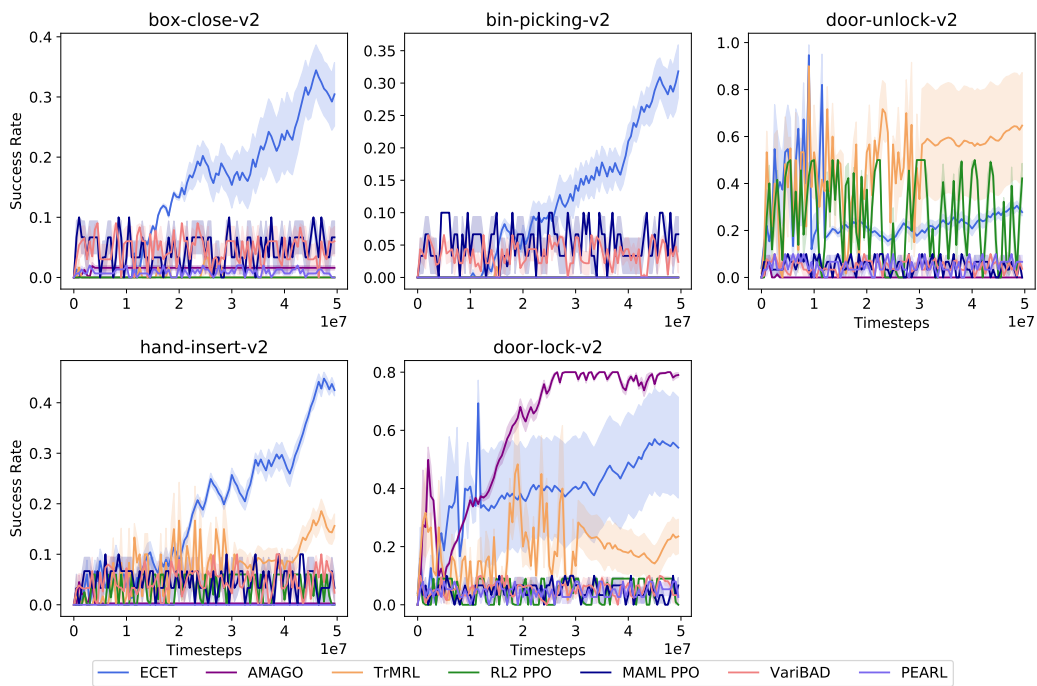


Figure 34: Success rate per test task for the ML45 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD.

1998  
 1999  
 2000  
 2001  
 2002  
 2003  
 2004  
 2005  
 2006  
 2007  
 2008  
 2009  
 2010  
 2011  
 2012  
 2013  
 2014  
 2015  
 2016  
 2017  
 2018  
 2019  
 2020  
 2021  
 2022  
 2023  
 2024  
 2025  
 2026  
 2027  
 2028  
 2029  
 2030  
 2031  
 2032  
 2033  
 2034  
 2035  
 2036  
 2037  
 2038  
 2039  
 2040  
 2041  
 2042  
 2043  
 2044  
 2045  
 2046  
 2047  
 2048  
 2049  
 2050  
 2051

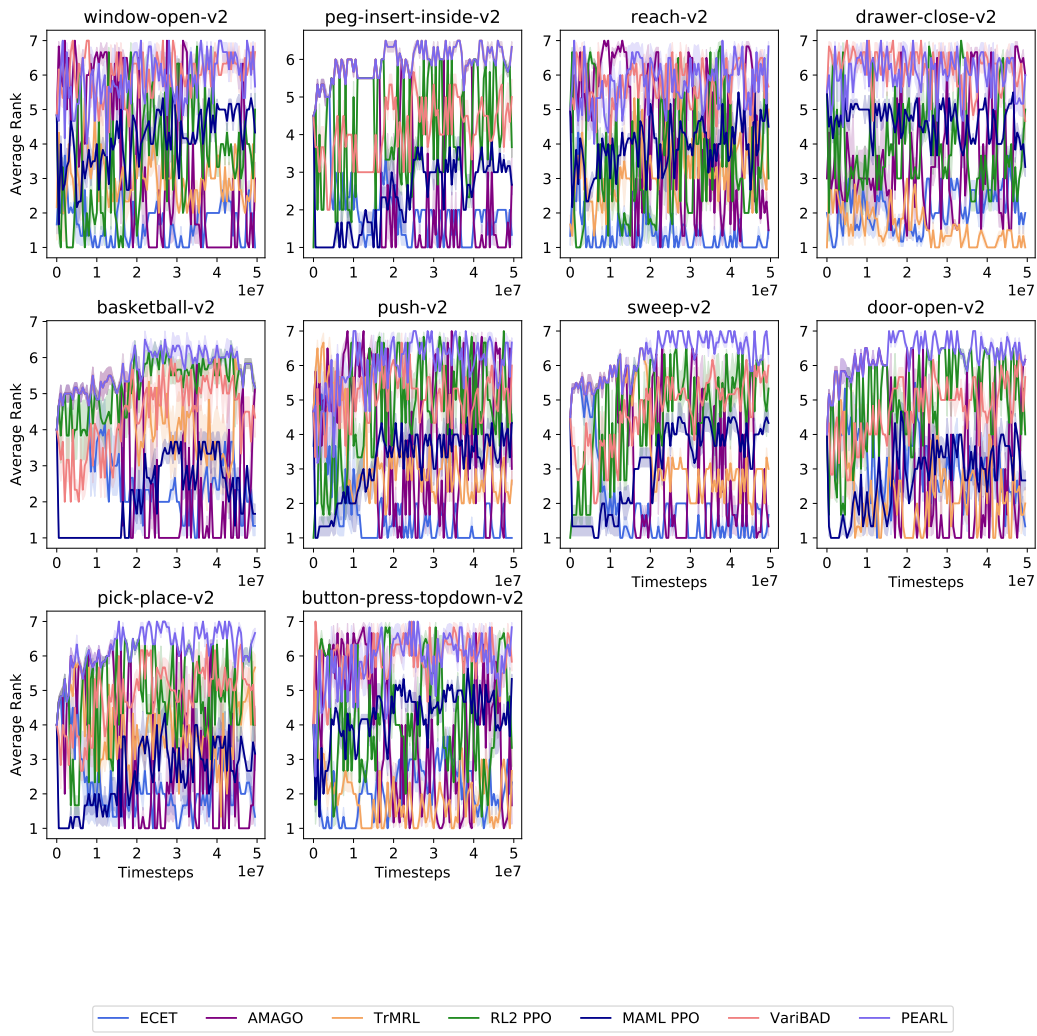


Figure 35: Average rank per train task for the ML10 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD.

2052  
 2053  
 2054  
 2055  
 2056  
 2057  
 2058  
 2059  
 2060  
 2061  
 2062  
 2063  
 2064  
 2065  
 2066  
 2067  
 2068  
 2069  
 2070  
 2071  
 2072  
 2073  
 2074  
 2075  
 2076  
 2077  
 2078  
 2079  
 2080  
 2081  
 2082  
 2083  
 2084  
 2085  
 2086  
 2087  
 2088  
 2089  
 2090  
 2091  
 2092  
 2093  
 2094  
 2095  
 2096  
 2097  
 2098  
 2099  
 2100  
 2101  
 2102  
 2103  
 2104  
 2105

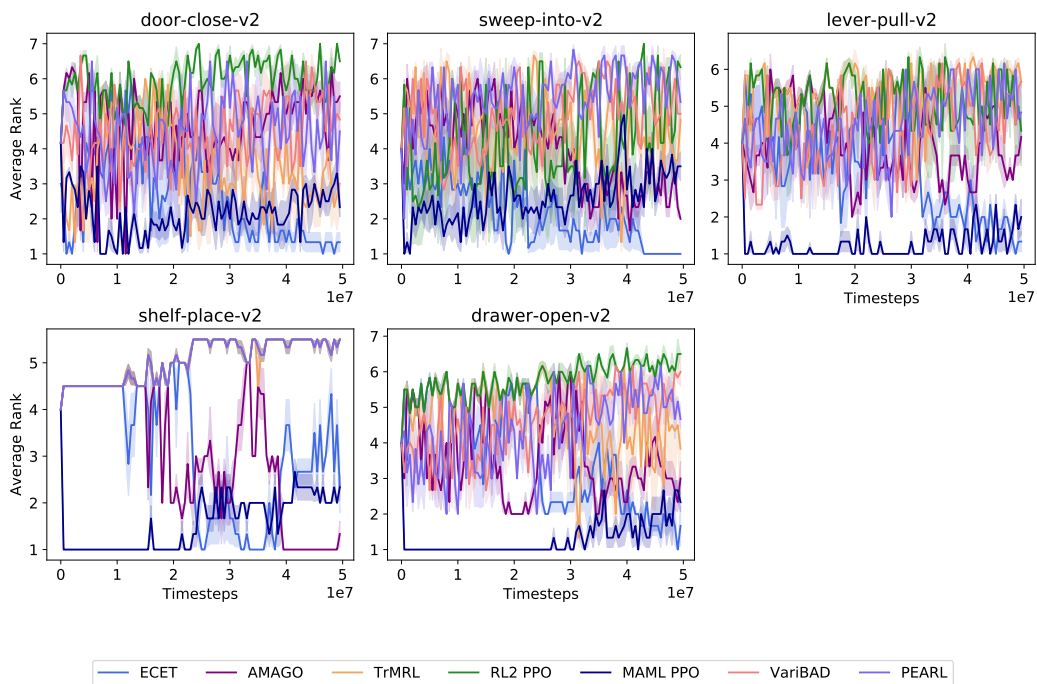


Figure 36: Average rank per test task for the ML10 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD.

2106  
 2107  
 2108  
 2109  
 2110  
 2111  
 2112  
 2113  
 2114  
 2115  
 2116  
 2117  
 2118  
 2119  
 2120  
 2121  
 2122  
 2123  
 2124  
 2125  
 2126  
 2127  
 2128  
 2129  
 2130  
 2131  
 2132  
 2133  
 2134  
 2135  
 2136  
 2137  
 2138  
 2139  
 2140  
 2141  
 2142  
 2143  
 2144  
 2145  
 2146  
 2147  
 2148  
 2149  
 2150  
 2151  
 2152  
 2153  
 2154  
 2155  
 2156  
 2157  
 2158  
 2159

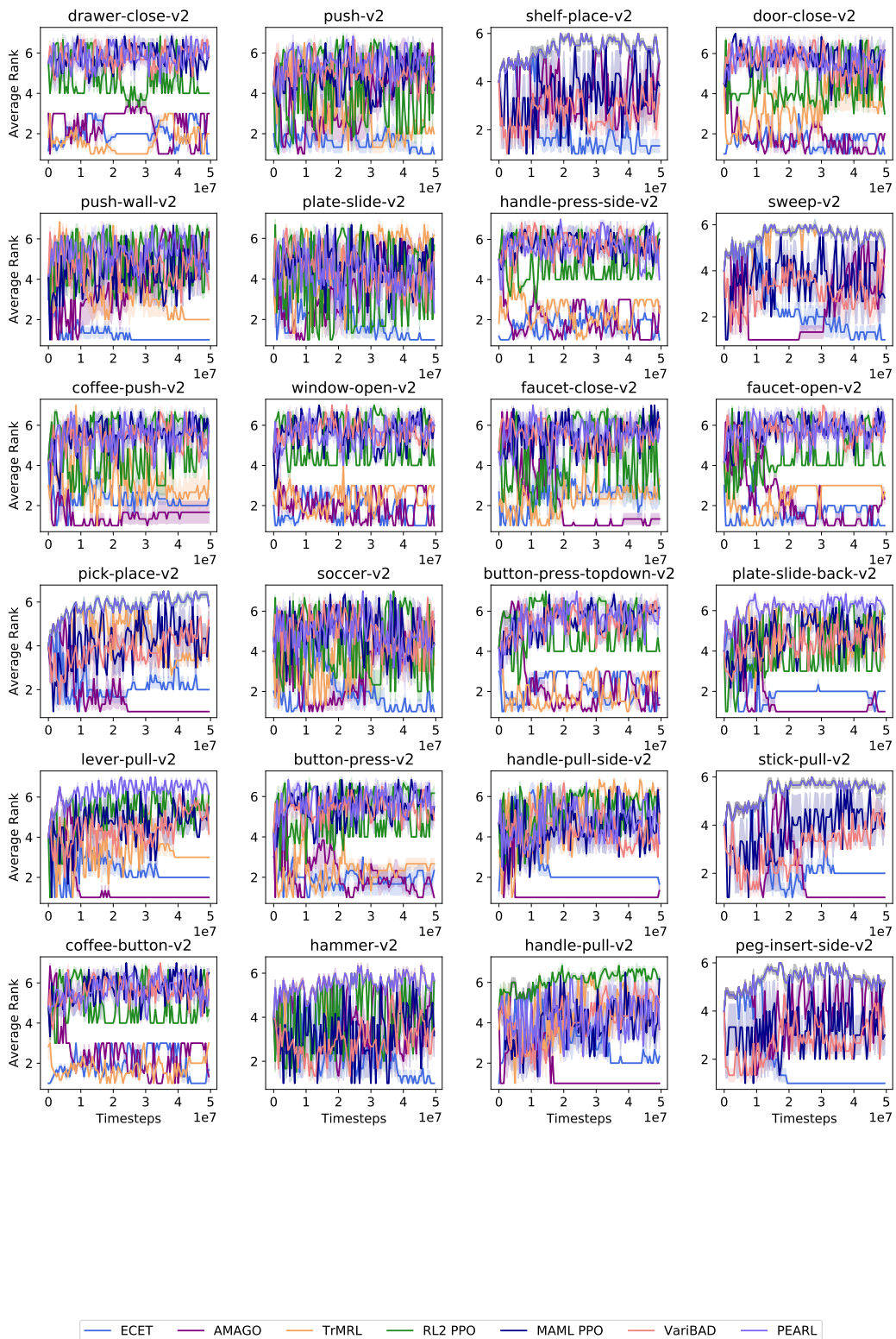


Figure 37: Average rank per train task for the ML45 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD (Part 1).



2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

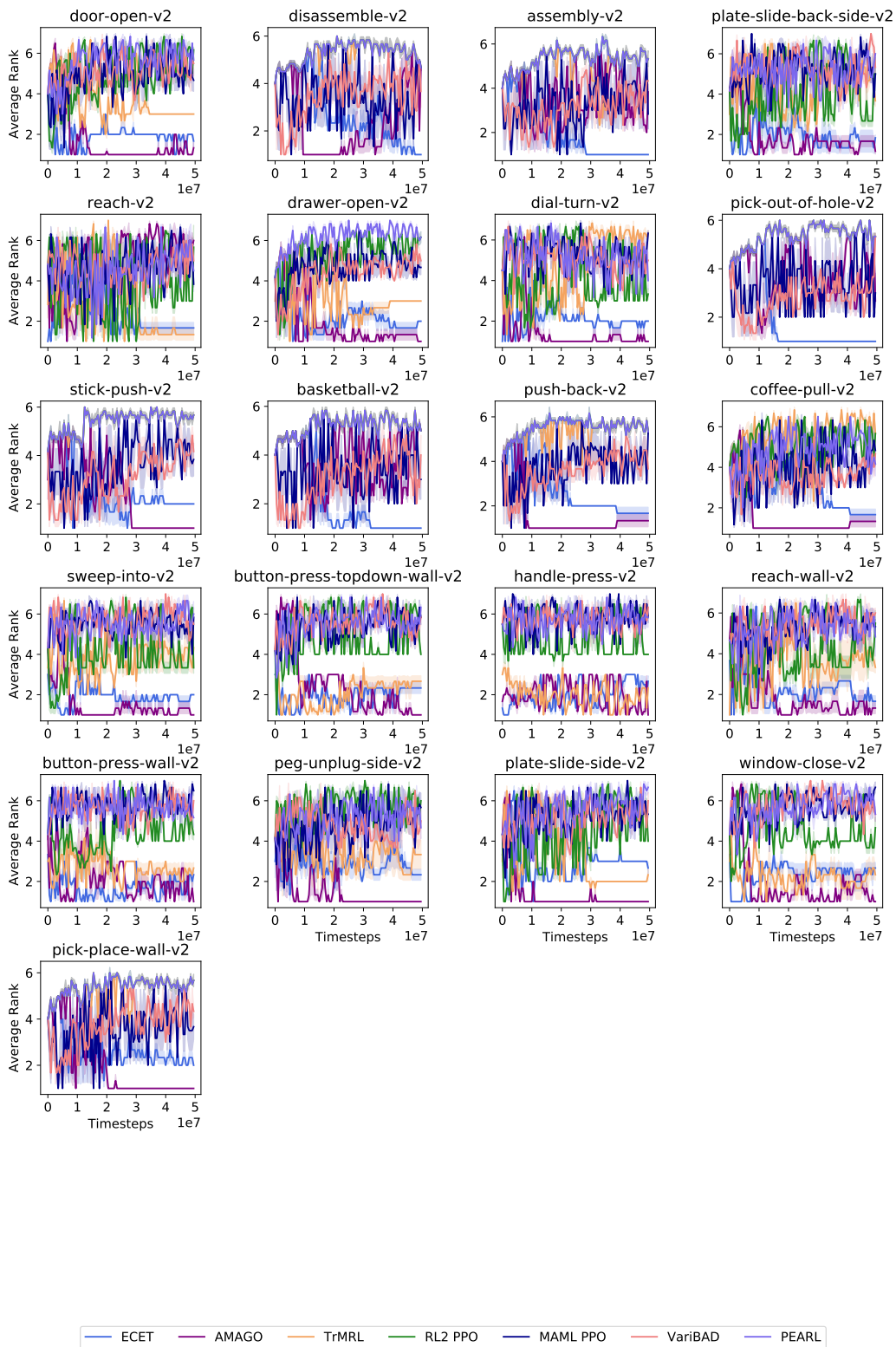


Figure 38: Average rank per train task for the ML45 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD (Part 2).

2214  
 2215  
 2216  
 2217  
 2218  
 2219  
 2220  
 2221  
 2222  
 2223  
 2224  
 2225  
 2226  
 2227  
 2228  
 2229  
 2230  
 2231  
 2232  
 2233  
 2234  
 2235  
 2236  
 2237  
 2238  
 2239  
 2240  
 2241  
 2242  
 2243  
 2244  
 2245  
 2246  
 2247  
 2248  
 2249  
 2250  
 2251  
 2252  
 2253  
 2254  
 2255  
 2256  
 2257  
 2258  
 2259  
 2260  
 2261  
 2262  
 2263  
 2264  
 2265  
 2266  
 2267

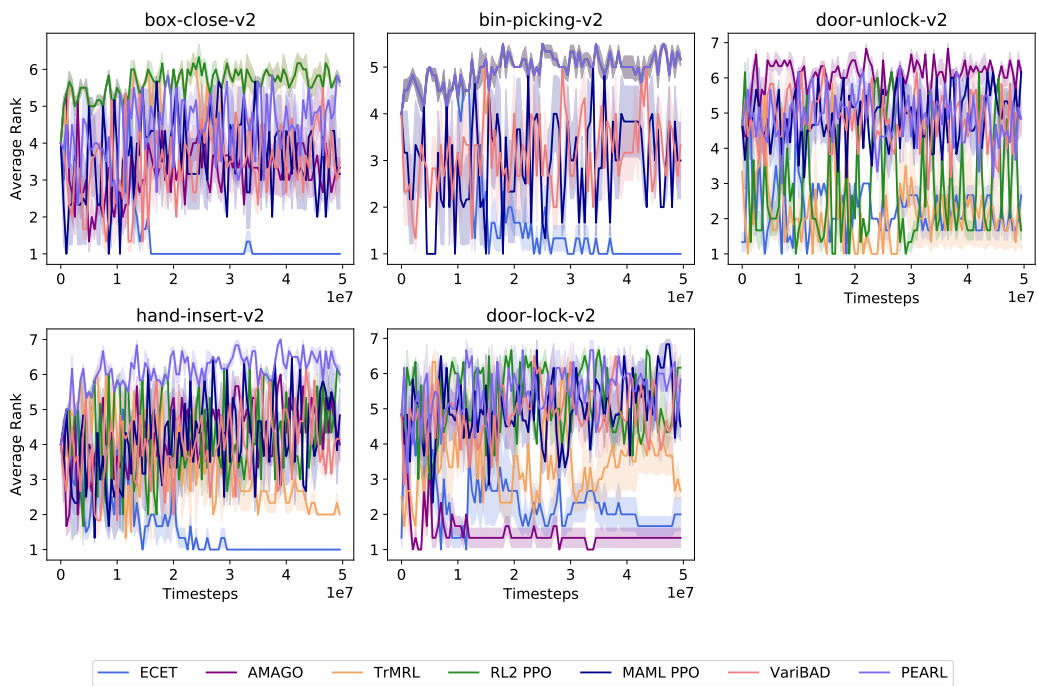


Figure 39: Average rank per test task for the ML45 benchmark for ECET, TrMRL, PEARL, RL2 PPO, MAML PPO, and VariBAD.