

Learning Task-Aware Abstract Representations for Meta-Reinforcement Learning

Anonymous authors

Paper under double-blind review

Abstract

A central challenge in meta-reinforcement learning (meta-RL) is enabling agents trained on a set of environments to generalize to new, related tasks without requiring full policy retraining. Existing model-free approaches often rely on context-conditioned policies learned via encoder networks. However, these context encoders are prone to overfitting to the training environments, resulting in poor out-of-sample performance on unseen tasks. To address this issue, we adopt an alternative approach that uses an abstract representation model to learn augmented, task-aware abstract states. We achieve this by introducing a novel architecture that offers greater flexibility than existing recurrent network-based approaches. In addition, we optimize our model with multiple loss terms that encourage predictive, task-aware representations in the abstract state space. Our method simplifies the learning problem and provides a flexible framework that can be readily combined with any off-the-shelf reinforcement learning algorithm. We provide theoretical guarantees alongside empirical results, showing strong generalization performance across classical control and robotic meta-RL benchmarks.

1 Introduction

Achieving robust generalization to new environments is challenging for agents trained within a single setting. While such agents may perform well in the training environment, they often fail when faced with even minor changes to dynamics. Meta-reinforcement learning (meta-RL) addresses this limitation by training across a distribution of tasks, thereby enabling the few-shot adaptation of classical RL agents (Beck et al., 2025; Nagabandi et al., 2018).

One of the most promising directions in meta-RL is context conditioning (Zintgraf et al., 2021; Lee et al., 2020). The core idea is to learn a context representation that captures task-specific information from a distribution of environments. At test time, the agent can explicitly infer the context of a new, unseen environment and adapt its behavior accordingly. Recent work (Lee et al., 2020; Rimón et al., 2024) has shown that context-conditioned dynamics models are particularly effective to meta-RL settings.

A key limitation of context-based meta-RL approaches is that the agent’s performance depends heavily on the encoder’s ability to accurately distinguish between tasks. Consequently, if the context encoder overfits to the training distribution - or fails to generalize - its out-of-sample performance can degrade substantially (Zintgraf et al., 2021).

This work aims to improve generalization in meta-RL by leveraging abstract representation models (ARMs), which learn high-level representations of environmental states using a pretrained model. Our goal is to produce task-aware, augmented abstract state representations. We obtain task-aware states with a neural architecture that encodes states into an abstract space through a state-space encoder and infers task encodings with a recurrent neural network (RNN). These representations are optimized jointly using a set of complementary loss functions. Our approach offers two main advantages: (1) ARMs are straightforward to train and generalize well, reducing reliance on the context encoder; and (2) they integrate seamlessly with any off-the-shelf model-free reinforcement-learning algorithm, giving the agent direct access to task information.

Our network architecture integrates abstract state representation learning with task inference in a unified framework. Unlike previous approaches that rely primarily on task-encoders that are trained separately, we introduce a joint optimization framework based on predictive losses in the abstract space—capturing transitions, rewards, and task structure.

We call our methodology **Environment-aware Meta Encoding and Representation Abstraction for Latent Domains—EMERALD**. It is a novel reinforcement learning framework that learns *task-aware abstract representations* capturing both environmental dynamics and task-specific factors, enabling improved generalization and task disambiguation in meta-RL settings. Our key contributions are:

- EMERALD shows *state-of-the-art performance* across a suite of meta-RL benchmarks by outperforming existing approaches in *few-shot adaptation and transfer* tasks, and is compatible with standard reinforcement learning algorithms such as SAC and PPO.
- We provide *theoretical justification* for EMERALD’s design, showing how incorporating contextual information into the abstract state space promotes faster generalization to unseen tasks.

2 Related Work

Meta Reinforcement Learning Meta-reinforcement learning (meta-RL) enables RL agents to rapidly adapt to new, unseen tasks or environments (Beck et al., 2025; Nagabandi et al., 2018). The goal is to learn a policy trained on a set of tasks $\mathcal{T}_i \sim p(\mathcal{T}_{\text{train}})$ that can quickly adapt to new tasks. Achieving this requires the agent to generalize effectively across tasks, leveraging experience from the training distribution to perform well in novel settings (Finn et al., 2017). Meta-RL can also be viewed through the lens of partially observable Markov decision processes (POMDPs), where the true task identity is unobserved and must be inferred through interaction, making context inference analogous to belief state estimation (Humplik et al., 2019; Rakelly et al., 2019).

Context-based Meta-RL Context-based meta-RL approaches learn a context encoder to capture variations across tasks. PEARL (Rakelly et al., 2019) and VariBAD (Zintgraf et al., 2021) formulate the context-encoding problem using a Bayesian Adaptive MDP (BAMDP), where a context variable z is modeled as a belief state inferred from past transitions and used to condition the policy, yielding $\pi(s_t | z)$. VariBAD passes an augmented state—consisting of the current state and the belief variable—to the policy. This allows the agent to adapt its behavior based on the inferred task identity. We build on this idea of an augmented state space, but rather than concatenating the state and context, we feed the learned context variable directly into the state-space encoder of the representation model.

Model-based Contextual Meta-RL MAMBA (Rimon et al., 2024) extends VariBAD to the model-based setting by integrating context encoding into Dreamer’s latent imagination process. A related method, CaDM (Lee et al., 2020), learns contextualized dynamics models by conditioning both the reward and transition models on context, and employs a backward-and-forward transition mechanism to stabilize training. Our approach builds on these methods, but uses only forward predictions in the abstract space. As in Dreamer-based methods such as MAMBA, we allow the policy to operate in an abstract space with lower cardinality than the raw state space, $|\mathcal{X}| \ll |\mathcal{S} \times \mathcal{T}|$. However, our state-task encoder can also be deployed in model-free settings via abstract representation models. Figure 1 illustrates the high-level architectural differences between our approach and other context-based methods.

Abstract Representation Models in Deep RL Abstract representation models (ARMs) are widely used in deep reinforcement learning to improve generalization and sample efficiency (Botteghi et al., 2022; Starre et al., 2022; Ni et al., 2024). These models map high-dimensional states \mathcal{S} to compact latent spaces \mathcal{X} via encoders $\psi : \mathcal{S} \rightarrow \mathcal{X}$, enabling more structured learning. Unlike “standard” auto-encoders, abstract representation models in the context of RL are typically used to obtain generalizable abstractions of the reward and transition dynamics. Several works illustrate this: for instance, Zhang et al. (2018) decouple dynamics and rewards for domain generalization, while Li et al. (2021) apply abstract models to adversarial settings with shared dynamics but varying visual observations. We build on these insights and extend ARMs

to the meta-RL setting. Similar to prior combined RL approaches (François-Lavet et al., 2019; Lee et al., 2020), our method is flexible and can be integrated with both model-free and model-based agents.

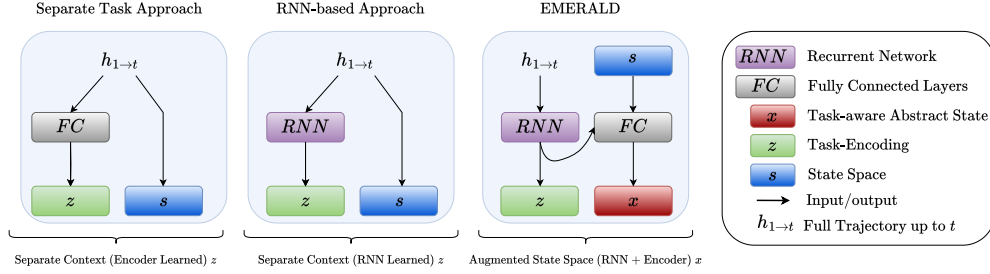


Figure 1: Comparison of state and context representation strategies. **Left:** Methods such as CaDM learn state (s) and context (z) representations independently, without temporal modeling. **Middle:** VariBAD-like approaches use RNNs to learn z from trajectory history $h_{1 \rightarrow t}$, but still treat s and z as separate components. **Right:** Our proposed EMERALD ARM architecture integrates the trajectory-derived context z with the current abstract state to form a task-aware shared latent state representation x , enabling unified and task-aware modeling of the environment dynamics.

3 Task-aware Abstract Representation Models

Notation We define a task \mathcal{T}_i as a Markov Decision Process (MDP), represented by the 5-tuple $\mathcal{T}_i = (\mathcal{S}, \mathcal{A}, P_i, r_i, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} the action space, $P_i(s_{t+1} | s_t, a_t)$ the transition kernel, and r_i the reward function mapping state-action pairs to real-valued rewards. The discount factor is denoted $\gamma \in [0, 1)$. We index time steps by the subscript t , e.g., s_t and s_{t+1} .

We sample N independent and identically distributed (i.i.d.) training tasks $\mathcal{T}_i \sim p(\mathcal{T})$. From a task, we can collect data sets of transitions $\mathcal{D}_i = \{(s_{k,t}, a_{k,t}, s_{k,t+1}, r_{k,t})\}_{k=1}^{m_i}$. The combined dataset across all N tasks is denoted $\mathcal{D} = \bigcup_{i=1}^N \mathcal{D}_i$. We let $m_i = |\mathcal{D}_i|$ denote the number of transitions collected from task \mathcal{T}_i . We denote the entropy of a random variable X as $\mathbb{H}(X)$, and the conditional entropy as $\mathbb{H}(X | Z)$. $\mathcal{L}_i(\theta)$ denotes the loss of a model with parameters θ in a single environment i and $\mathbb{E}_{\mathcal{T} \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}}(\theta)$ the loss over all tasks in the distribution.

Model Components The abstract *state* space is denoted as \mathcal{X} and the *task* embedding space as \mathcal{Z} . Note that, in our architecture, \mathcal{X} contains both the information about the task as well as the state in that task. The ARM has four components: (1) the contextual state-space encoder $\psi : \mathcal{S} \times \mathcal{Z} \rightarrow \mathcal{X}$, (2) the transition function $\tau : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{X}$, (3) the reward predictor $\rho : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$. Lastly, (4) a learned encoder $\phi : (\mathcal{S} \times \mathcal{A} \times \mathbb{R})^{t-1} \times \mathcal{S} \rightarrow \mathcal{Z}$, mapping a history of transitions to an environment embedding $\mathcal{Z} \in \mathbb{R}^k$: $z_i = \phi(h_{0 \rightarrow t-1})$, where $h_{0 \rightarrow t-1} = (s_0, a_0, r_1, s_1, a_1, r_2, \dots, s_{t-1}, a_{t-1}, r_t, s_t)$. Figure 1 (right) shows a schematic overview of the model architecture.

Policy We denote the policy as $\pi : \mathcal{X} \rightarrow \mathcal{A}$. The *value function* $V^\pi(x)$ is defined as the expected return when starting from abstract state x and following policy π :

$$V^\pi(x) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid x_0 = x \right].$$

The *Q-function* $Q^\pi(x, a)$ denotes the expected return when taking action a in abstract state x , and thereafter following π :

$$Q^\pi(x, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid x_0 = x, a_0 = a \right].$$

Goal The aim is to learn a policy that maximizes the expected return across tasks sampled from the task distribution $p(\mathcal{T})$. Formally, we optimize:

$$\max_{\pi, \theta} \mathbb{E}_{\mathcal{T} \sim p(\mathcal{T})} \left[\mathbb{E}_{(x_0, a_0, r_1, \dots) \sim \pi, \mathcal{T}} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \right],$$

where the expectation is taken over both the task distribution and the rollout distribution induced by the policy π across a task distribution, operating in the abstract state space \mathcal{X} .

3.1 Task-Aware Augmented State-Representations in Meta-RL

Necessity of task-aware ARMs in meta-RL Despite the capacity of representation models to generalize well in simple meta-RL settings such as a variety of different mazes, we argue that only through the addition of task-specific information can their potential be fully realized. We prove this in Proposition 1.

Proposition 1 (Task-aware Encoding Yields a Strictly Lower Loss) *Let \mathcal{T}_1 and \mathcal{T}_2 be two tasks such that, for some state-action pair (s_t, a_t) , the transition dynamics differ: $P_1(s_{t+1} | s_t, a_t) \neq P_2(s_{t+1} | s_t, a_t)$. Consider an encoder $\phi : (\mathcal{S} \times \mathcal{A} \times \mathbb{R})^{t-1} \times \mathcal{S} \rightarrow \mathbb{R}^k$ that maps a history of transitions $h_{1:t} = (s_0, a_0, r_1, s_1, \dots, a_{t-1}, r_t, s_t)$ to a task embedding $z = \phi(h_{1:t})$. We say z is task-aware if it retains task-specific information from the transition history, quantified as $I(h_{1:t}; z) > 0$, and task-agnostic if it does not, i.e., $I(h_{1:t}; z) = 0$. Then, for any model composed of encoder ψ , transition model τ , reward model ρ , and task encoder ϕ , the minimum achievable loss under this embedding is greater than or equal to that under a task-aware embedding:*

$$\inf_{\substack{\psi, \tau, \rho, \phi \\ I(h_{1:t}; z) = 0}} [\mathcal{L}_{\text{transition}}(\theta) + \lambda \mathcal{L}_{\text{reward}}(\theta)] \geq \inf_{\substack{\psi, \tau, \rho, \phi \\ I(h_{1:t}; z) > 0}} [\mathcal{L}_{\text{transition}}(\theta) + \lambda \mathcal{L}_{\text{reward}}(\theta)].$$

The intuition (full proof in Appendix C) is that without task-awareness (measured by how much task information is retained by the encoder), the abstract representation model will not be able to learn a representation that can differentiate based solely on the input it receives. As such, task-agnostic state-space encoders essentially introduce a one-to-many mapping, which collapses into a single point during optimization. We provide an illustration of this insight in Figure 2.

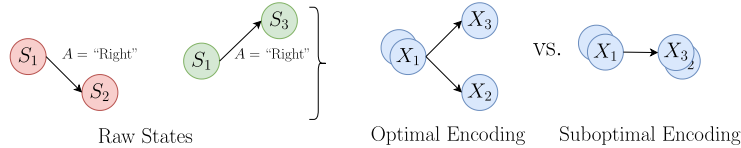


Figure 2: Illustration of Proposition 1. Raw input states, where S_1 maps to state S_2 in task red and to S_3 in task green. **Left:** Optimal encoding where both follow-up abstract states are mapped to different regions in space. **Right:** Suboptimal encoding we would expect if the encoder does not retain contextual information. In that case, we may observe a collapse of the two states, because $\psi(s_3^{\mathcal{T}_1})$ and $\psi(s_2^{\mathcal{T}_2})$ are pushed to the same point in order to minimize the loss.

How can task-awareness be achieved? We introduce a novel approach that integrates task information directly into the abstract state space \mathcal{X} via the encoder ψ , producing representations to be *both* compact *and* task-relevant. Unlike prior meta-RL methods that condition downstream components (such as dynamics, reward models, or policies) on context (e.g., VariBAD, PEARL), our method performs *early* context injection by embedding task identity during the state abstraction phase. This design allows our agent to resolve task ambiguity upstream, providing a strong inductive bias and leading to more robust representations that generalize better to unseen tasks.

By shifting task uncertainty into the encoder, we also simplify the overall architecture and training objective, eliminating the need for amortized inference or latent variable optimization at test time. Furthermore, this

abstraction-first formulation is highly modular and integrates seamlessly with both model-based and model-free RL algorithms. We formally justify this design through the lens of entangled versus disentangled modeling: in the latter, state abstraction and task representation are learned separately, whereas in the former — our setting — they are integrated into a single representation. We then empirically show that our approach leads to improved generalization across standard meta-RL benchmarks.

Proposition 2 (Entanglement Yields Lower Entropy) *Let \mathcal{T}_1 and \mathcal{T}_2 be two tasks such that for some state-action pair (s_t, a_t) , the transitions differ: $P_1(s_{t+1} | s_t, a_t) \neq P_2(s_{t+1} | s_t, a_t)$. Let $z = \phi(h_{i,:t}) \in \mathcal{Z}$ be a deterministic task embedding computed from transition history. Consider two model variants: (1) a task-state entangled version, where $x_t = \psi(s_t, z)$, the transition prediction is given by $x_{t+1} \approx x_t + \tau(x_t, a_t)$ and the estimated reward $\hat{r}_t = \rho(x_t, a_t)$; and (2) task-state disentangled version, where $x_t = \psi^*(s_t)$, $x_{t+1} \approx x_t + \tau^*(x_t, a_t, z)$ and the estimated reward $\hat{r}_t = \rho^*(x_t, a_t, z)$. Then, assuming ψ , ρ and τ are deterministic and trained to perfectly fit the transition loss and reward loss, the total entropy of the representations satisfies:*

$$\mathbb{H}(\psi(s_t, z)) \leq \mathbb{H}(\psi^*(s_t)) + \mathbb{H}(z).$$

That is, the task-state entangled version yields strictly lower joint entropy due to resolving task ambiguity in the representation, resulting in simpler transitions and reward functions.

We provide the full proof in Appendix C. Intuitively, Proposition 2 states that entanglement causes context ambiguity to be resolved *before* it reaches downstream tasks (e.g., before computing transition probabilities or state-to-action mappings via the policy). As such, following a general “conditioning reduces entropy” rule (Cover & Thomas, 1999; Gray, 2011; MacKay, 2003), which has been explored extensively in the machine learning literature (Bounoua et al., 2025; Pandey & Dukkipati, 2017; Shamir et al., 2010), early conditioning effectively ensures that all model components are one-to-one mappings. We call the task-aware abstract state space produced by the the state-task encoder the *augmented abstract* state space.

4 EMERALD: Architecture & Objectives

ARM: Architecture Blocks We introduce a parameterized abstract representation model, consisting of (1) a state-space encoder ψ , (2) a latent transition function τ , (3) a reward predictor ρ , and (4) a transitions-to-context function ϕ to adjust for environment-specific dynamics (see Appendix G for a more detailed overview). Formally, we have four parameterized model components:

$$\psi(s_t, z; \theta_\psi) \rightarrow x_t, \quad x_{t+1} = x_t + \tau(x_t, a_t; \theta_\tau), \quad \hat{r}_t = \rho(x_t, a_t; \theta_\rho), \quad z = \phi(h_t; \theta_\phi).$$

We combine the model parameters into a jointly optimizable model via θ , where $\theta = (\theta_\psi, \theta_\tau, \theta_\rho, \theta_\phi)$.

ARM: Objectives The transition and reward dynamics are effectively captured by the loss functions defined below. We define the transition and reward losses as:

$$\begin{aligned} \mathcal{L}(\theta)_{\text{transition}} &= \mathbb{E}_{(s_t, a_t, s_{t+1}, r_t) \sim D} \left[\left\| \psi(s_{t+1}, z; \theta_\psi) - (\psi(s_t, z; \theta_\psi) + \tau(\psi(s_t, z; \theta_\psi), a_t; \theta_\tau)) \right\|^2 \right], \\ \mathcal{L}(\theta)_{\text{reward}} &= \mathbb{E}_{(s_t, a_t, s_{t+1}, r_t) \sim D} \left[\left\| r_t - \rho(\psi(s_t, z; \theta_\psi), a_t; \theta_\rho) \right\|^2 \right]. \end{aligned}$$

The formulations above follow standard practice (e.g., François-Lavet et al. (2019)), but we introduce context vectors z retrieved from the context encoder ϕ and include it directly into the optimization objective. To prevent a potential state space collapse, we also include a regularizer:

$$\mathcal{L}(\theta)_{\text{reg}} = \exp \left(-C_d \left\| \psi(s^+, z; \theta_\psi) - \psi(s^-, z; \theta_\psi) \right\|^2 \right),$$

where s^+ and s^- are two randomly sampled states and C_d is a constant. Intuitively, this term enforces that any two states are some distance apart.

Algorithm 1 EMERALD Abstract Representation Model Training

Require: ARM (transition encoder, reward predictor, abstract state encoder, task-encoder), task distribution $p(\mathcal{T})$, learning rate α_θ , batch size B , epochs per task n_{epochs} , initialize replay buffers \mathcal{D}_i for every task \mathcal{T}_i , context horizon H

- 1: **Offline Data Collection:** For each task \mathcal{T}_i , collect a dataset $\mathcal{D}_i^{\text{offline}} = \{(s_t, a_t, s_{t+1}, r_t)\}_{t=1}^{m_i}$ of offline transitions following policy π .
- 2: **while** not converged **do** ▷ ARM training loop
- 3: **for** each task \mathcal{T}_i **do**
- 4: Reset *context window*: $h^i \leftarrow \emptyset$
- 5: **for** epoch = 1 to n_{epochs} **do**
- 6: **for** each batch $\{(s_t, a_t, s_{t+1}, r_t)\} \sim \mathcal{D}_i^{\text{offline}}$ **do**
- 7: Sample task embedding $z \sim \phi(\cdot \mid h^i; \theta_\phi)$
- 8: Encode states: $x_t = \psi(s_t, z; \theta_\psi)$ $x_{t+1}^{\text{true}} = \psi(s_{t+1}, z; \theta_\psi)$
- 9: Predict transition and reward: $\hat{x}_{t+1} = x_t + \tau(x_t, a_t; \theta_\tau)$
- 10: Predict reward: $\hat{r}_t = \rho(x_t, a_t; \theta_\rho)$
- 11: Compute loss: $\mathcal{L}_\theta = \mathcal{L}_{\text{transition}} + \mathcal{L}_{\text{reward}} + \beta \mathcal{L}_{\text{reg}}$
- 12: Update ψ parameters: $\theta_\psi \leftarrow \theta_\psi - \alpha_{\theta_\psi} \frac{1}{B} \sum_{i=1}^B \nabla_{\theta_\psi} \mathcal{L}_i$ ▷ Parameter updates
- 13: Update ϕ parameters: $\theta_\phi \leftarrow \theta_\phi - \alpha_{\theta_\phi} \frac{1}{B} \sum_{i=1}^B \nabla_{\theta_\phi} \mathcal{L}_i$
- 14: Update ρ parameters: $\theta_\rho \leftarrow \theta_\rho - \alpha_{\theta_\rho} \frac{1}{B} \sum_{i=1}^B \nabla_{\theta_\rho} \mathcal{L}_i$
- 15: Add to context window: $c^i \leftarrow c^i \cup \{(s_t, a_t, r_t, s_{t+1})\}$ until $|h^i| \leq H$
- 16: **end for**
- 17: **end for**
- 18: **end for**
- 19: **end while**
- 20: **return** Trained ARM.

Putting everything together, we obtain the following joint objective:

$$\mathcal{L}(\theta) = \mathcal{L}(\theta)_{\text{transition}} + \mathcal{L}(\theta)_{\text{reward}} + \beta \mathcal{L}(\theta)_{\text{reg}},$$

where $\beta \geq 0$ determines the regularization strength. We minimize this objective using standard mean squared error losses for the transition and reward terms. Algorithm 1 shows the pseudocode for the ARM training process. Throughout the experimental section, we sample the offline data for training (line 1 of the algorithm) the model using a random policy. However, in practice, the offline data might also be collected via a non-random or expert policy. A Python implementation can be found at <https://anonymous.4open.science/r/EMERALD-F324>.

Training of the policy The policy is trained after ARM training. We freeze the representation model’s parameters and pass the online transitions through the model. The augmented abstract state approach allows us to combine the learned representations with any reinforcement learning algorithm (e.g., SAC, PPO). We define the optimal policy as:

$$\pi^* = \arg \max_{\pi} V^{\pi}(\psi(s_t, z; \theta_\psi)), \quad \forall s_t \in \mathcal{S}.$$

Algorithm 2 shows the pseudo-code for policy training (full version shown in Appendix D).

Algorithm 2 EMERALD Policy Learning

Require: Trained ARM, tasks $\mathcal{T}_i \sim p(\mathcal{T})$, learning rates α_π, α_V , batch size B , rollout length n_r , updates n_u , horizon H ; parameters θ_π, θ_V ; buffers $\mathcal{D}_i \leftarrow \emptyset$

- 1: **while** not converged **do**
- 2: **for all** \mathcal{T}_i **do** ▷ Rollout
- 3: $\mathcal{B} \leftarrow \emptyset, h^i \leftarrow \emptyset$
- 4: **for** $t = 1 \dots n_r$ **do**
- 5: $z \sim \phi(\cdot \mid h^i; \theta_\phi); x_t = \psi(s_t, z; \theta_\psi), a_t \sim \pi_{\theta_\pi}(\cdot \mid x_t)$
- 6: env step $\rightarrow s_{t+1}, r_t, x_{t+1} = \psi(s_{t+1}, z; \theta_\psi)$
- 7: add (x_t, a_t, r_t, x_{t+1}) to \mathcal{B} and h^i (truncate to H)
- 8: **end for**
- 9: $\mathcal{D}_i \leftarrow \mathcal{D}_i \cup \mathcal{B}$
- 10: **end for**
- 11: **for** $u = 1 \dots n_u$ **do** ▷ Update
- 12: **for all** \mathcal{T}_i **do**
- 13: minibatch $b^i \sim \mathcal{D}_i$; compute L_π^i, L_V^i
- 14: **end for**
- 15: $\theta_\pi \leftarrow \theta_\pi - \alpha_\pi \nabla_{\theta_\pi} \sum_i L_\pi^i; \theta_V \leftarrow \theta_V - \alpha_V \nabla_{\theta_V} \sum_i L_V^i$
- 16: **end for**
- 17: **end while**
- 18: **return** $\pi_{\theta_\pi}, V_{\theta_V}$

5 Experiments

Our evaluation consists of three complementary experimental sets, each designed to test a distinct aspect of EMERALD. First, we compare EMERALD with strong baselines on *out-of-distribution* tasks. Second, we evaluate the effect of training on multiple environments by measuring performance on *in-distribution* tasks. Third, we perform two ablation studies: (i) we examine how the quality of the ARM affects policy learning, and (ii) we visualize the learned task-aware latent space to assess how effectively the policy exploits it.

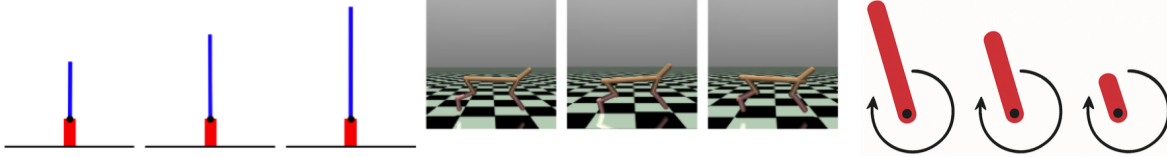


Figure 3: Illustration of the Experimental Setup. **Left:** Cartpole with different sizes. **Middle:** Three HalfCheetah with volumes. **Right:** Pendulum with different lengths.

5.1 Comparative benchmark on out-of-distribution tasks

We largely follow the experimental protocol of Lee et al. (2020), using variants of the continuous-control CartPole task and several MuJoCo benchmarks (Todorov et al., 2012). The learning algorithms are trained on a set of configurations and evaluated on previously unseen ones. For CartPole, for example, we train on pole lengths 1.0, 1.5, and 2.0 and test the out-of-distribution performance on poles of length 0.5 and 2.5. The three evaluation domains are illustrated in Figure 3. EMERALD is paired with Soft Actor-Critic (SAC) (Haarnoja et al., 2018) (EMERALD SAC) and Proximal Policy Optimization (PPO) (Schulman et al., 2017) (EMERALD PPO). The ARM is pre-trained offline with 5×10^4 CartPole transitions, $11 \times \sim 4.5^4$ Pendulum transitions, and 5×10^4 HalfCheetah transitions. In order not to give our model an unfair advantage, we collected the offline transitions under a random policy.

	CartPole (Lengths)		Pendulum (Lengths)		HalfCheetah (Volume)	
	Test (moderate)	Test (extreme)	Test (moderate)	Test (extreme)	Test (moderate)	Test (extreme)
Vanilla PPO	198.2 \pm 0.9	187.8 \pm 4.7	-1113.2 \pm 69.1	-1356.8 \pm 48.0	807.7 \pm 553.6	574.0 \pm 645.6
Vanilla SAC	199.2 \pm 0.3	199.4 \pm 0.1	-402.0 \pm 131.3	-1274.2 \pm 133.7	1998.3 \pm 142.1	177.4 \pm 321.3
Stacked PPO	197.8 \pm 1.3	189.2 \pm 6.1	-475.7 \pm 228.1	-488.2 \pm 178.2	361.1 \pm 141.7	5.7 \pm 208.1
PEARL	198.0 \pm 1.4	187.5 \pm 10.9	-645.3 \pm 320.7	-1136.6 \pm 251.0	642.1 \pm 488.3	462.1 \pm 534.5
PPO EP	196.3 \pm 4.0	184.5 \pm 9.7	-374.3 \pm 24.6	-256.7 \pm 26.4	895.3 \pm 445.1	674.2 \pm 686.8
CaDM	197.9 \pm 3.0	193.0 \pm 3.5	-279.8 \pm 42.1	-426.4 \pm 227.0	1224.2 \pm 630.0	1021.1 \pm 676.6
VariBAD	199.1 \pm 2.8	192.0 \pm 3.5	-318.5 \pm 56.6	-643.1 \pm 311.4	2964 \pm 179.5	1618.3 \pm 251.8
EMERALD PPO	199.1 \pm 1.1	199.3 \pm 3.5	-312.4 \pm 19.0	-313.7 \pm 101.5	2942.5 \pm 213.5	1401.4 \pm 421.3
EMERALD SAC	194.6 \pm 2.7	196.2 \pm 4.9	-198.3 \pm 23.3	-749.4 \pm 340.9	5101.1 \pm 151.8	3271.8 \pm 180.2

Table 1: Performance comparison means over 10 runs (Mean Return \pm Std.) for CartPole, Pendulum, and HalfCheetah (Volume). The best statistically significant (Welch’s t-test $p < 0.01$) results in each group are boldfaced and quantify the rolling average over the last 100 timesteps.

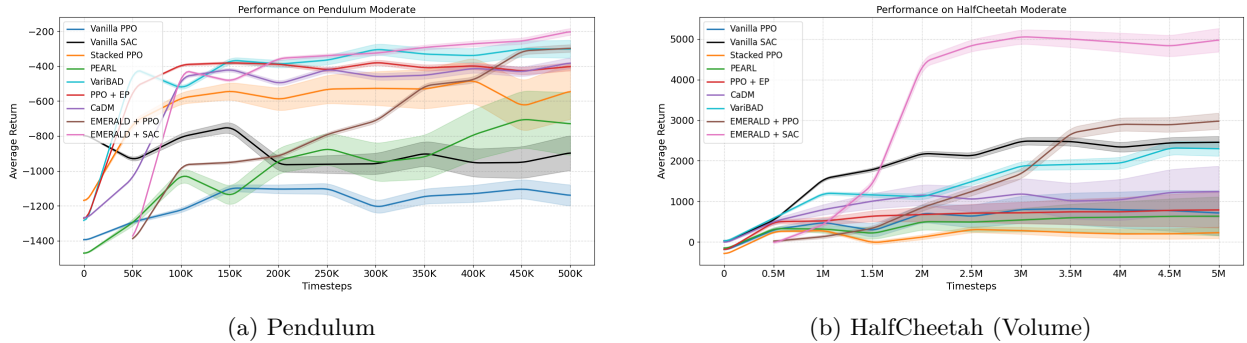


Figure 4: Learning Curves for two out-of-distribution Experiments from the Moderate Regime (rolling averages over 10 independent runs and adjusted for ARM budget allocation)

We compare EMERALD PPO and EMERALD SAC with seven baselines:

1. **Vanilla PPO** (Schulman et al., 2017): Standard Proximal policy optimization (PPO).
2. **Vanilla SAC** (Haarnoja et al., 2018): Standard Soft Actor Critic (SAC).
3. **Stacked PPO** (Lee et al., 2020): PPO variant that feeds a fixed window of past transitions to the policy.
4. **PEARL** (Rakelly et al., 2019): Infers a context variable by maximizing expected return.
5. **PPO-EP** (Zhou et al., 2019): PPO version that concatenates an embedding from early interactions to the state.
6. **CaDM** (Lee et al., 2020): Augments a forward-backward dynamics model with a context encoder.
7. **VariBAD** (Zintgraf et al., 2021): Uses an RNN to produce a belief state that is concatenated with the environment state and passed to both the policy and value functions.

Following Lee et al. (2020), we evaluate all methods on unseen environments under two regimes (see Appendix J for complete specifications):

1. **Moderate**: Unseen test tasks differ only slightly from the training distribution.

2. **Extreme:** Unseen test tasks deviate substantially from training task.

To ensure fairness, the ARM is pretrained on offline data and the corresponding transitions are deducted from each method’s overall interaction budget, which guarantees that all agents experience the same total number of transitions (see Appendix M for the complete allocation budget).

Table 1 shows that EMERALD surpasses most baselines; the corresponding learning curves are displayed in Figures 4a-4b. On the Pendulum task we do not outperform PPO-EP in the extreme regime, suggesting at a potential limitation when the environment exhibits strong periodic dynamics. However, the difference of EMERALD-PPO with PPO-EP is not too large. Notably, both EMERALD-SAC and EMERALD-PPO perform well overall, underscoring the benefit of modularity: one can boost performance by combining the ARM with a different learning algorithm. Interestingly, the performance of the vanilla learners appears predictive of performance with EMERALD: for instance, SAC and EMERALD-SAC both outperform other methods on the HalfCheetah tasks.

HalfCheetah (Direction)	
Vanilla PPO	1448.9 \pm 291.3
PEARL	1491.4 \pm 131.5
VariBAD	2112.5 \pm 193.6
EMERALD PPO	2259.9 \pm 317.1
EMERALD SAC	1343.2 \pm 391.1

Table 2: In-distribution performance on HalfCheetah (Direction) over 5 runs with random seeds

5.2 EMERALD on in-distribution tasks

Comparative in-distribution performance HalfCheetah (Direction) is an experimental setup in which one task requires to walk forward and in the other task to walk backward. The goal here is to compare EMERALD’s performance to established baselines such as VariBAD (Zintgraf et al., 2021) and PEARL in the in-distribution setting (training tasks and test tasks are the same). The results are reported in Table 2 and show that the PPO version of our approach outperforms existing baselines, albeit within the margin of error. The improvement over Vanilla PPO further suggests that training *with* augmented abstract states provides an additional boost in performance in this experimental setting.

Table 3: EMERALD In-Distribution Performance on Training Tasks (Mean Return \pm Std over 5 seeds). Mean difference statistically significant (Welch’s t-test $p < 0.01$)

Environment	Config	# Training Tasks	Samples per Task	Mean Return \pm Std
CartPole	1	1	200k	91.4 \pm 37.5
CartPole	2	5	40k	187.0 \pm 19.8
HalfCheetah	1	1	3M	749.2 \pm 209.3
HalfCheetah	2	5	600k	1218.0 \pm 338.7

Effect of Environment Diversity on Performance We ask whether training on a *greater diversity* of tasks is more beneficial than training with *more data* drawn from a single environment. This is an in-distribution setting, as train and test tasks are the same. We allocate a total budget of 3M transitions for HalfCheetah (Volume) and 200k for CartPole. In **Configuration 1** the entire budget is spent on one task, whereas in **Configuration 2** the budget is split evenly across multiple tasks. Table 3 summarizes the outcomes. Configuration 2 yields higher overall performance, indicating that exposure to a wider set of environments outweighs simply scaling up data for a single task, demonstrating that increased task-diversity offers greater sample efficiency.

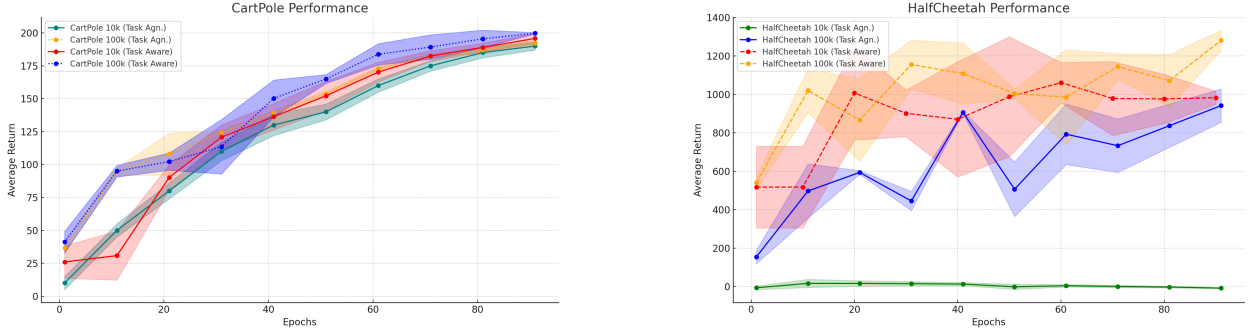


Figure 5: Average Policy Performance (10 runs) with stdev. on CartPole and HalfCheetah (Volume) for several Training Iterations, starting from a model trained with 1 epoch up to 100 epochs. Dashed lines are task-aware iterations, solid lines display task-agnostic versions (no context encoding).

5.3 Ablation Studies

Model Effect on Policy Performance: We study how performance varies with (i) the number of training epochs, (ii) the amount of offline data, and (iii) the presence of task-aware encodings. Figure 5 reports the results (each marker is the mean of 5 seeds). In the CartPole setting with 10k transitions, returns plateau after roughly 10 epochs, whereas HalfCheetah (Volume) requires more training but reaches satisfactory performance at about 30 epochs. Although the 100k-transition regime converges earlier, it *does not* exceed the final performance of the 10k model—showing that comparable results can be achieved with an order of magnitude less data. Removing task encodings (i.e., using task-agnostic abstract states) sharply degrades performance, particularly on HalfCheetah, providing empirical support for Proposition 1.

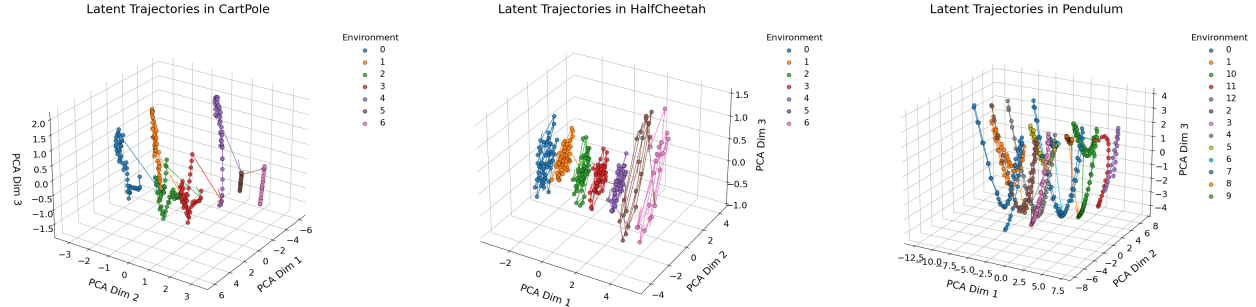


Figure 6: Visualization of the Latent Space (PCA mapping) of the EMERALD ARM with Training and Test Environments of the first 50 transitions. **Left:** Cartpole. **Middle:** HalfCheetah (Volume). **Right:** Pendulum. In CartPole and HalfCheetah (Volume) the out-of-distribution trajectories are 5 and 6, while in Pendulum 11 and 12.

Latent Trajectory Analysis Figure 6 visualizes policy behavior in both training and test environments. For each environment, we record the first 50 transitions and project them into a three-dimensional principal component space. The resulting training and test trajectories exhibit similar shapes but occupy distinct regions, suggesting that EMERALD learns behaviors that transfer across tasks while remaining sensitive to task identity.

6 Discussion & Conclusion

We introduced EMERALD, a meta-reinforcement-learning method that leverages augmented, context-aware abstract states. Across benchmarks, our approach outperforms other model-free context-based baselines, showing that task-aware representation models yield substantial gains in the meta-RL setting.

Beyond improved performance, we provided justification for the representational and architectural advantages of our approach: by disentangling task-relevant features early in the pipeline, we simplify training, improve robustness, and gain compatibility with a broad class of RL algorithms as our procedure yields a simple encoder that produces task-aware abstractions. Despite these strengths, EMERALD has limitations. First, the optimal structure of the latent space (e.g., dimensionality) can be environment-dependent and difficult to determine *a priori*. Second, as the number of training tasks grows, the latent space may need to expand to capture task diversity, potentially increasing computational cost. In such cases, hybrid approaches that decouple task inference from state abstraction may offer better scalability. Identifying where this trade-off becomes critical is an important direction for future work.

A natural extension of EMERALD is to integrate it into a *model-based* framework. Whereas the present study focuses on using a *pretrained* latent model within a modular, model-free pipeline, one could instead exploit the latent dynamics directly to train a policy, thereby opening the door to direct comparisons with established model-based meta-RL methods.

References

- Jacob Beck, Risto Vuorio, Evan Zheran Liu, Zheng Xiong, Luisa Zintgraf, Chelsea Finn, Shimon Whiteson, et al. A tutorial on meta-reinforcement learning. *Foundations and Trends® in Machine Learning*, 18(2–3): 224–384, 2025.
- Nicolò Botteghi, Mannes Poel, and Christoph Brune. Unsupervised representation learning in deep reinforcement learning: A review. *arXiv preprint arXiv:2208.14226*, 2022.
- Mustapha Bounoua, Giulio Franzese, and Pietro Michiardi. Learning to match unpaired data with minimum entropy coupling. *arXiv preprint arXiv:2503.08501*, 2025.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, 1999.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 1126–1135. PMLR, 2017.
- Vincent François-Lavet, Yoshua Bengio, Doina Precup, and Joelle Pineau. Combined reinforcement learning via abstract representations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 3582–3589, 2019.
- Robert M. Gray. *Entropy and Information Theory*. Springer Science & Business Media, 2011.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 1861–1870. PMLR, 2018.
- Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Kinal Mehta, and João G. M. Araújo. Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022.
- Jan Humplik, Alexandre Galashov, Leonard Hasenclever, Pedro A. Ortega, Yee Whye Teh, and Nicolas Heess. Meta reinforcement learning as task inference. *arXiv preprint arXiv:1905.06424*, 2019.
- Kimin Lee, Younggyo Seo, Seunghyun Lee, Honglak Lee, and Jinwoo Shin. Context-aware dynamics model for generalization in model-based reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119, pp. 5757–5766. PMLR, 2020.

- Bonnie Li, Vincent François-Lavet, Thang Doan, and Joelle Pineau. Domain adversarial reinforcement learning. *arXiv preprint arXiv:2102.07097*, 2021.
- David J. C. MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003.
- Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S. Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Learning to adapt in dynamic, real-world environments through meta-reinforcement learning. *arXiv preprint arXiv:1803.11347*, 2018.
- Tianwei Ni, Benjamin Eysenbach, Erfan Seyedsalehi, Michel Ma, Clement Gehring, Aditya Mahajan, and Pierre-Luc Bacon. Bridging state and history representations: Understanding self-predictive rl. *arXiv preprint arXiv:2401.08898*, 2024.
- Gaurav Pandey and Ambedkar Dukkipati. Variational methods for conditional multimodal deep learning. In *Proceedings of the 2017 International Joint Conference on Neural Networks (IJCNN)*, pp. 308–315. IEEE, 2017.
- Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021.
- Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient off-policy meta-reinforcement learning via probabilistic context variables. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 5331–5340. PMLR, 2019.
- Zohar Rimón, Tom Jurgenson, Orr Krupnik, Gilad Adler, and Aviv Tamar. Mamba: an effective world model approach for meta-reinforcement learning. *arXiv preprint arXiv:2403.09859*, 2024.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Ohad Shamir, Sivan Sabato, and Naftali Tishby. Learning and generalization with the information bottleneck. *Theoretical Computer Science*, 411(29–30):2696–2711, 2010.
- Rolf AN Starre, Marco Loog, and Frans A Oliehoek. Model-based reinforcement learning with state abstraction: A survey. In *Benelux Conference on Artificial Intelligence*, pp. 133–148. Springer, 2022.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *Proceedings of the 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5026–5033. IEEE, 2012.
- Geoffrey Van Driessel and Vincent François-Lavet. Component transfer learning for deep rl based on abstract representations. *arXiv preprint arXiv:2111.11525*, 2021.
- Amy Zhang, Harsh Satija, and Joelle Pineau. Decoupling dynamics and reward for transfer learning (2018). *arXiv preprint arXiv:1804.10689*, 2018.
- Wenxuan Zhou, Lerrel Pinto, and Abhinav Gupta. Environment probing interaction policies. In *ICLR*, 2019.
- Luisa Zintgraf, Sebastian Schulze, Cong Lu, Leo Feng, Maximilian Igl, Kyriacos Shiarlis, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: Variational bayes-adaptive deep rl via meta-learning. *Journal of Machine Learning Research*, 22(289):1–39, 2021.