

SIMPLICIAL EMBEDDINGS IMPROVE SAMPLE EFFICIENCY IN ACTOR–CRITIC AGENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent works have proposed accelerating the wall-clock training time of actor-critic methods via the use of large-scale environment parallelization; unfortunately, these can sometimes still require large number of environment interactions to achieve a desired level of performance. Noting that well-structured representations can improve the generalization and sample efficiency of deep reinforcement learning (RL) agents, we propose the use of *simplicial embeddings*: lightweight representation layers that constrain embeddings to simplicial structures. This geometric inductive bias results in sparse and discrete features that stabilize critic bootstrapping and strengthen policy gradients. When applied to FastTD3, FastSAC, and PPO, simplicial embeddings consistently improve sample efficiency and final performance across a variety of continuous- and discrete-control environments, without any loss in runtime speed.

“Order is not imposed from the outside, but emerges from within¹.”

— Ilya Prigogine

1 INTRODUCTION

Deep reinforcement learning (deep RL) has delivered impressive progress in continuous control, enabling agile locomotion (Smith et al., 2022; Zhuang et al., 2023; Margolis et al., 2024) and dexterous manipulation (Popov et al., 2017; Akkaya et al., 2019; Luo et al., 2025). Yet a persistent tension remains between *training speed* (wall-clock efficiency) and *sample efficiency* (the number of environment interactions). Some modern agents such as TD-MPC2 (Hansen et al., 2023) and SR-SPR (D’Oro et al., 2022) achieve strong returns with relatively few interactions, but demand substantial compute and engineering complexity. In contrast, recent fast actor–critic variants have scaled throughput with massive parallelization (Li et al., 2023; Singla et al., 2024; Gallici et al., 2025; Seo et al., 2025). While methods such as FastTD3 (Seo et al., 2025) rapidly solve humanoid benchmarks, they require far more interactions to reach comparable performance. Similar limitations have been observed in Parallel Q-Learning (Li et al., 2023) and large-scale actor–critic frameworks such as IMPALA and SEED RL (Espeholt et al., 2018; 2020). This trade-off limits applicability in domains where interactions are expensive and time is constrained, such as robotics.

A natural objection is that, in modern simulators, environment steps are cheap and can be generated in massive parallel batches, so sample efficiency is less important. However, this view overlooks several practical and scientific concerns. First, algorithms that are data-hungry in simulation rarely transfer well to real-world scenarios (Tobin et al., 2017; Akkaya et al., 2019). Second, large-scale parallelization requires substantial compute and energy resources, raising both efficiency and sustainability concerns (Schwartz et al., 2020; Henderson et al., 2020). Third, sample efficiency is closely tied to generalization: agents that exploit structure from fewer trajectories tend to be more robust under distributional shifts (Zhang et al., 2018; Yao et al., 2025). Moreover, in high-dimensional simulators such as IsaacGym, each step can be significantly more expensive, compounding inefficiency as tasks grow harder (Makoviychuk et al., 2021; Rudin et al., 2021). These issues highlight why sample efficiency remains central even in the era of massively parallel deep RL.

¹This perspective resonates with deep RL: stability cannot be forced solely through more compute, heavier regularizers, or larger critics. Instead, inductive biases that shape the geometry of representations can allow order to *emerge from within*, leading to more stable critics and more efficient policies under non-stationarity.

Shaping representations with auxiliary losses (Anand et al., 2019; Laskin et al., 2020; Schwarzer et al., 2021; Castro et al., 2021; Fujimoto et al., 2023) has been shown to improve sample efficiency in deep RL. However, such methods increase algorithmic complexity and add computational overhead through extra forward and backward passes (Fujimoto et al., 2023). Alternatively, architectural components, such as convolutions (Fukushima, 1980; LeCun et al., 1989) and attention (Bahdanau et al., 2016), can be used to induce structure leading to desirable downstream properties.

Discrete and sparse representations have several desirable properties in comparison to their dense and continuous counterparts. Notably, sparse and discrete representations increase robustness to noise (Donoho et al., 2006), training stability by reducing catastrophic interference (Liu et al., 2019), sample efficiency (Fumero et al., 2023), interpretability (Murphy et al., 2012; Lavoie et al., 2023; Wabartha & Pineau, 2024) and improved generative modeling (Lavoie et al., 2025). In this work, we posit that several of those properties are beneficial in the context of reinforcement learning.

While several methods exist for learning discrete representations explicitly (Jang et al., 2017; Madison et al., 2017; van den Oord et al., 2018), these methods use straight-through estimation (Bengio et al., 2013) which is a biased gradient estimator. Fortunately, discretization may be implicitly induced via Simplicial Embeddings (SEM) (Lavoie et al., 2023), an architectural component that partitions a latent representation into a sequence of L simplices. SEM is fully differentiable, thus avoiding the negative effect of explicit discretization while enacting some of the desirable properties of discrete and sparse representations. Concretely, we show that SEM improves both data efficiency and asymptotic performance across diverse environments such as IsaacGym (Makoviychuk et al., 2021), HumanoidBench (Sferrazza et al., 2024), and the Arcade Learning Environment (Bellemare et al., 2013), while preserving (and often improving) wall-clock speed.

2 PRELIMINARIES

2.1 ACTOR–CRITIC REINFORCEMENT LEARNING

We consider a standard Markov decision process (MDP) defined by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \gamma)$, with state space \mathcal{S} , action space \mathcal{A} , transition distribution $P(s'|s, a)$, reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, and discount factor $\gamma \in [0, 1)$. The objective is to maximize the expected discounted return

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right], \quad (1)$$

where the agent follows a policy $\pi(a|s)$. Actor–critic methods maintain both a parameterized policy $\pi_{\theta}(a|s)$ (the actor) and an action-value function $Q_{\phi}(s, a)$ (the critic). The critic is trained to minimize the Bellman error

$$\mathcal{L}_Q(\phi) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[(Q_{\phi}(s, a) - y)^2 \right], \quad y = r + \gamma \mathbb{E}_{a' \sim \pi_{\theta}(\cdot|s')} [Q_{\phi}^{-}(s', a')], \quad (2)$$

where ϕ^{-} denotes target network parameters and \mathcal{D} is a replay buffer. The actor is updated via the policy gradient defined as $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q_{\phi}(s, a)]$. While this can be effective, bootstrapped training is notoriously fragile. Errors in Q_{ϕ} propagate recursively through the target y , and when the representation used to compute Q_{ϕ} is poorly conditioned, these errors amplify and cause divergence or collapse (Fujimoto et al., 2018).

A recent line of work has sought to reduce the *wall-clock* cost of actor–critic training. FastTD3 (Seo et al., 2025) builds on TD3 (Fujimoto et al., 2018) by leveraging (i) parallel simulation across many environment instances, (ii) large-batch critic updates, and (iii) algorithm design choices like distributional critics (C51) (Bellemare et al., 2017), noise scaling and clipped double Q-learning (CDQ) (Fujimoto et al., 2018). Together, these design choices enable high-throughput training while retaining stable convergence, although FastTD3 (Seo et al., 2025) still remains sample-inefficient.

Policies and critics often rely on latent representations extracted from raw states (Lesort et al., 2018). Formally, an encoder $f_{\psi} : \mathcal{S} \rightarrow \mathbb{R}^d$ maps observations s into embeddings $z = f_{\psi}(s)$, which are then consumed by either the critic, the actor, or both depending on the architecture. Some methods share a common encoder across actor and critic (e.g., SAC (Haarnoja et al., 2018), DrQ (Yarats et al., 2021), DrQ-v2 (Yarats et al., 2022)), while others (e.g., DDPG (Lillicrap et al., 2015), FastTD3 (Seo et al., 2025)) maintain separate encoders. Regardless of parameter sharing, these representations play a

central role in learning (Garcin et al., 2025). The critic estimates values $Q_\phi(s, a) \equiv Q_\phi(f_\psi(s), a)$, and the actor conditions its policy $\pi_\theta(a|s) \equiv \pi_\theta(a|f_\psi(s))$ on the chosen embedding. Ideally, z should preserve the Markov property and expose predictive features of the reward r and dynamics P .

Yet the choice and stability of such embeddings is far from guaranteed. When unconstrained, learned representations can introduce severe pathologies that destabilize value learning. For example, if $\|f_\psi(s)\| \rightarrow \infty$, critics may extrapolate to arbitrarily large Q-values outside the support of the replay buffer, inflating the Bellman error. Formally, if $Q_\phi(z, a) = w^\top z + b$ with linear heads, then $\|Q_\phi\| \rightarrow \infty$ as $\|z\| \rightarrow \infty$, leading to exploding targets y and divergent gradients. Similarly, if z exhibits strong correlations or degenerate directions, the critic’s regression problem becomes ill-conditioned: the covariance matrix $\Sigma = \mathbb{E}[zz^\top]$ may approach singularity, amplifying variance in temporal-difference updates. These phenomena are empirically linked to representation collapse, where value estimates drift irrecoverably and policy updates follow unstable gradients (Moalla et al., 2024; Castanyer et al., 2025).

2.2 SIMPLICIAL EMBEDDINGS

Simplicial embeddings (SEM; Lavoie et al., 2023) provide a lightweight inductive bias on representation geometry by constraining latent codes to lie on a product of simplices. Concretely, given encoder outputs $f_\psi(s) \in \mathbb{R}^{L \times V}$, the latent vector is partitioned into L groups of size V , and a softmax is applied within each group:

$$\tilde{z}_{\ell,v} = \frac{\exp(z_{\ell,v}/\tau)}{\sum_{v'=1}^V \exp(z_{\ell,v'}/\tau)}, \quad \forall \ell \in \{1, \dots, L\}, v \in \{1, \dots, V\}, \quad (3)$$

where $\tau > 0$ is a temperature parameter controlling the degree of sparsity. The resulting embedding \tilde{z} lies in the product space $\Delta^{V-1} \times \dots \times \Delta^{V-1}$, i.e., L categorical distributions of dimension V . This transformation ensures boundedness through group-wise normalization, induces sparsity as softmax competition (sharpened at low τ) drives near one-hot encodings, and promotes group structure by partitioning features into modular subspaces akin to mixtures-of-experts (Ceron et al., 2024b). In self-supervised learning and downstream classification, SEM has been shown to stabilize training and improve generalization, particularly in low-label and transfer settings (Lavoie et al., 2023). SEM does not rely on auxiliary losses or reconstruction terms; akin to an activation function, it only modifies the embedding geometry with the group-wise softmax, limiting computational overhead.

3 NON-STATIONARITY AMPLIFIES REPRESENTATION COLLAPSE

Several works have shown that non-stationarity can lead to severe degradation of learned representations across different domains (Lyle et al., 2022; Kumar et al., 2021a; Lyle et al., 2025; Castanyer et al., 2025). In supervised learning, label noise and distribution shifts can induce representation collapse, where features lose diversity and neurons become inactive (Li et al., 2022; Sokar et al., 2023; Dohare et al., 2024). Similar observations have been made in deep RL: *the constantly changing data distribution, induced by an evolving policy, exacerbates this phenomenon, often resulting in unstable critics and poor generalization* (Nauman et al., 2024a; Kumar et al., 2021a). These studies suggest that collapse is not an isolated pathology of specific architectures, but a general failure mode that emerges when training signals are non-stationary. In App. B we provide a formal analysis that demonstrates the relationship between non-stationarity and neuron dormancy.

A demonstration on CIFAR-10. We illustrate this phenomenon with a toy experiment on CIFAR-10 (Krizhevsky, 2009). We compare two training regimes: (i) a stationary setting with fixed labels, and (ii) a non-stationary setting where labels are periodically shuffled to mimic *the bootstrap dynamics* of deep RL (Sokar et al., 2023; Castanyer et al., 2025). Let (x, y) be training samples with $y \in \{1, \dots, K\}$. In the stationary regime, targets are fixed, so the conditional distribution $p(y|x)$ is constant and the empirical risk minimizer θ_t^* remains stable up to stochastic fluctuations. In the non-stationary regime, labels are periodically shuffled so that $y \mapsto \pi_t(y)$, where π_t is a permutation applied every T steps. This induces inflection points in the minimizer, shifting whenever π_t changes.

Fig. 1 shows that in the stationary regime, training is stable: losses decrease smoothly, dormant neuron rates remain low, and effective rank increases, indicating robust representation learning (Dohare

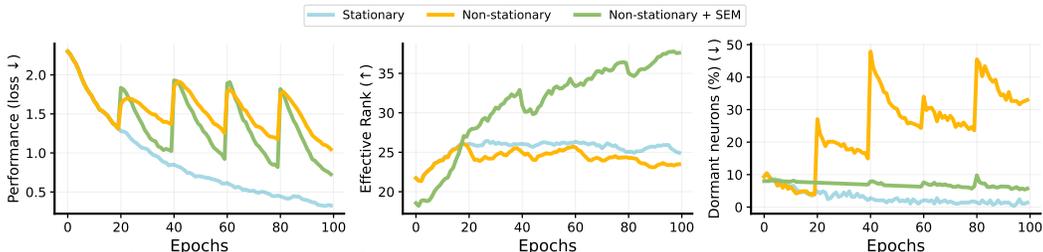


Fig. 1: **Training dynamics on CIFAR-10 with stationary vs. non-stationary targets.** In the stationary regime (fixed targets), losses decrease smoothly, neuron dormancy and effective rank remains controlled, suggesting stable representation learning. In the non-stationary regime (targets shuffled every 20 epochs), the model exhibits higher variance in losses, increased dormant neuron rates, and reduced effective rank. The addition of SEM mitigates this instability. Experiments are averaged over 3 independent seeds, with shaded areas reporting 95% confidence intervals, following Sokar et al. (2023).

et al., 2024; Sokar et al., 2023). In contrast, in the non-stationary regime, we observe instability: oscillating losses, rising neuron dormancy, and collapsing feature rank. Even in this simple supervised setting, instability in the target distribution alone is sufficient to undermine representational integrity. Similar optimization instabilities have been observed in prior work when evaluating CIFAR-10 under non-stationary conditions (Igl et al., 2021; Lee et al., 2023; Galashov et al., 2024; Castanyer et al., 2025).

Stabilizing Representations under Non-Stationarity with SEM *Simplicial Embeddings (SEM)* can mitigate this effect by projecting features onto a structured space that prevents collapse. The transformation enforces energy preservation; since each block has unit mass, representations cannot vanish and $\text{tr}(\Sigma_t)$ remains bounded away from zero. It also promotes diversity, as intra-block competition spreads information across coordinates, while multiple blocks (L) increase effective rank, counteracting covariance deflation. As shown in Fig. 1, critics trained with SEM retain higher effective rank, larger gradient energy, and lower neuron dormancy even when targets drift.

Takeaways:

- Non-stationarity exacerbates representation collapse, as evidenced by increased neuron dormancy and reduced effective rank.
- Simplicial Embeddings (SEM) introduce a simplex-based geometric prior that sustains feature diversity and prevent feature collapse.

4 UNDERSTANDING THE IMPACT OF SEM ON DEEP RL NETWORKS

In actor-critic methods such as FastTD3, the critic is trained against bootstrapped targets $y_t(s, a) = r(s, a) + \gamma Q_{\phi^-}(s', \pi_{\theta}(s'))$. Both the target distribution \mathcal{D}_t (samples (s, a, r, s') from the replay buffer) and the target value y_t evolve as the policy π_{θ} is updated. This continual drift produces a persistent bias term in $b_t = \nabla \mathcal{L}_{t+1}(\theta_t^*) = \mathbb{E}_{(s,a) \sim \mathcal{D}_{t+1}} \left[(Q_{\phi}(s, a) - y_{t+1}(s, a)) \nabla_{\theta} Q_{\phi}(s, a) \right]$, which is nonzero whenever π_{θ} or \mathcal{D}_t changes. Thus, the critic is never optimizing a fixed objective but is instead forced to chase a moving target.

Representation collapse under such non-stationarity poses a fundamental barrier to stable and efficient deep RL (see App. A for additional context). Standard actor-critic methods are particularly vulnerable. The critic’s representations are trained against drifting targets, and the actor in turn depends on those representations to update its policy. This tight coupling amplifies instability, leading to poor sample efficiency in continuous control tasks. To address this challenge, we evaluate *Simplicial Embeddings (SEM)* as a representation-level regularizer. SEM aims to encourage the hidden features of both actor and critic networks to maintain a well-structured geometric organization, preventing collapse and preserving diversity. By stabilizing the embedding space, SEM provides a principled mechanism for variance reduction and improved sample efficiency.

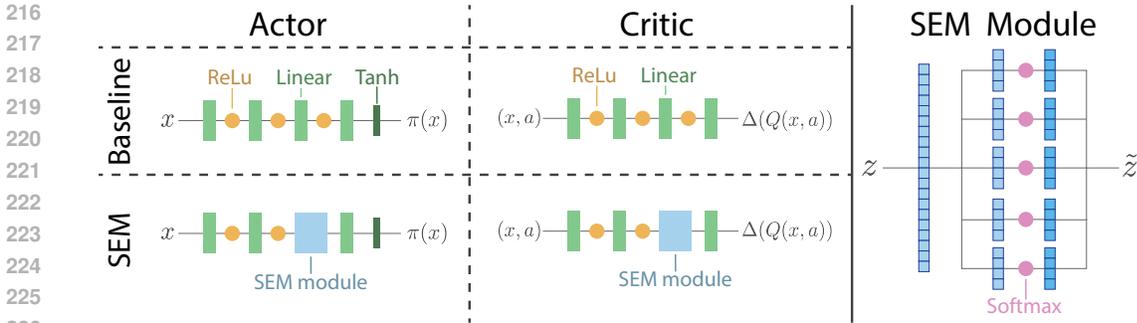


Fig. 2: **Actor-critic network architecture with SEM.** The actor (left) and critic (middle) architectures are modified with a SEM module, which partitions features into groups and applies group-wise softmax (right panel), constraining them to a product of simplices.

Setup. Because this section involves a large number of ablations and is computationally expensive, we restrict experiments to five benchmarks from the Humanoid suite (Sferrazza et al., 2024), evaluated on (Seo et al., 2025). These tasks share the same robot-state dimension. We report aggregate performance across the five tasks and six seeds, with shaded areas reporting 95% confidence intervals, and provide full details in App. F.

Integrating SEM on Actor-Critic Algorithms. We choose FastTD3 (Seo et al., 2025), as our primary testbed. FastTD3 is specifically designed to be a simple and compute-efficient baseline for continuous-control and humanoid benchmarks. Its streamlined architecture yields strong performance while significantly reducing wall-clock training time. At the same time, FastTD3 inherits the critic-driven weaknesses of TD3; its bootstrapped value targets are generated online by the actor, making the critic susceptible to non-stationarity. This coupling amplifies representation collapse, as instabilities in the critic propagate to both value estimates and policy updates. We conduct most of our ablations on FastTD3, while later sections demonstrate that the benefits of SEM also extend to other actor-critic algorithms, such as SAC (Haarnoja et al., 2018) and PPO (Schulman et al., 2017).

SEM can either be applied to the actor, the critic, or both. We build on prior work showing that the penultimate layer plays a critical role in representation quality (Moalla et al., 2024; Ceron et al., 2024b; Sokar & Castro, 2025), and that regularizing this layer can yield substantial performance gains. Fig. 2 illustrates how SEM is integrated into the actor-critic networks of FastTD3. For the critic, SEM replaces the baseline linear head with a structured projection, regularizing value estimates in the distributional C51 setting. For the actor, SEM is applied at the penultimate layer before the final linear+tanh, ensuring that the policy is conditioned on bounded and sparse features. Across the paper, dashed blue (blue, - -) curves indicate the baseline, while solid green, (green, —) curves represent the interventions added to the baseline.

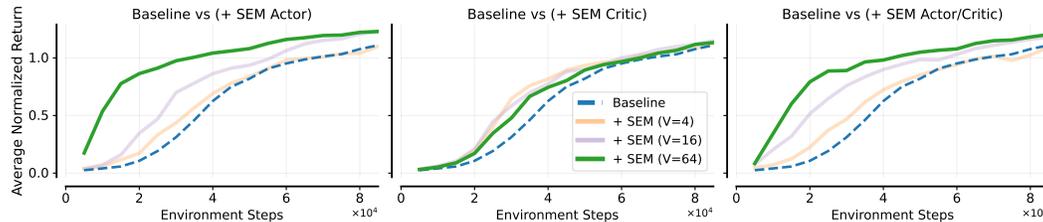


Fig. 3: **Average normalized return on 5 HumanoidBench tasks over 6 seeds.** Baseline agent (blue, - -) vs. SEM variants applied to actor, critic, or both. Each curve corresponds to an embedding dimension; $dim = 64$ (green, —) is highlighted. SEM accelerates early learning and improves asymptotic performance, with $dim = 64$ giving the most stable gains. In the three figures we use $L = 2$ for the SEM module.

Fig. 3 shows clear gains when applying SEM to the actor or to both actor and critic, and more moderate gains when applied only to the critic. Although different SEM dimensions (V) improve sample efficiency and asymptotic performance, $V = 64$ appears most effective. We further explore the relationship between L and V (see sec 4), as this tradeoff was a central focus of the original SEM

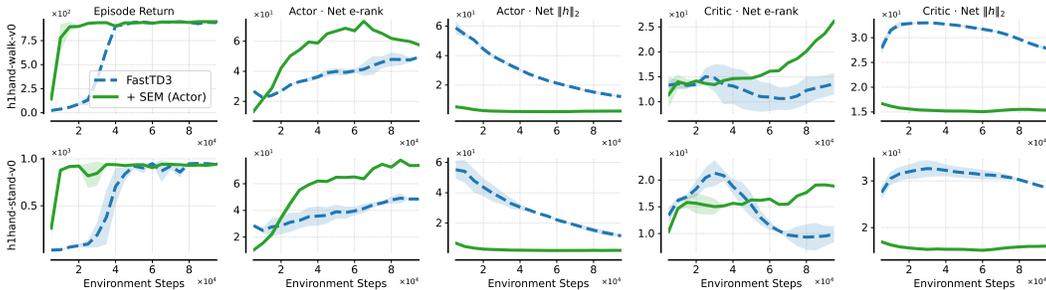


Fig. 4: **Learning and representation diagnostics on 2 HumanoidBench tasks.** SEM reaches high return earlier, raises actor/critic effective rank, and keeps actor features compact.

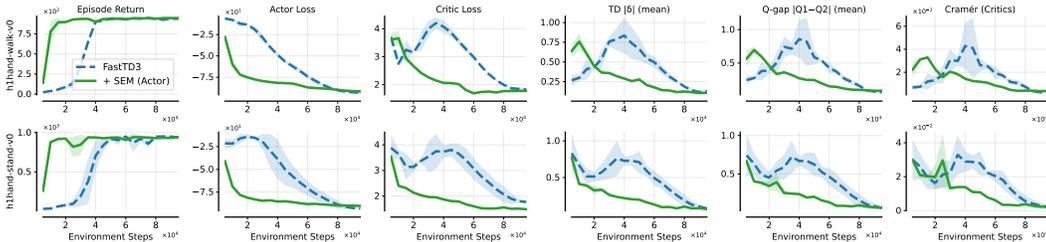


Fig. 5: **Learning dynamics on 2 HumanoidBench tasks.** SEM reaches high return faster, with lower losses, smaller TD error, reduced critic disagreement, and better-calibrated value estimates.

study (Lavoie et al., 2023). These results echo the non-stationary CIFAR-10 experiment, where SEM prevented feature collapse and stabilized learning (see Fig. 1).

The Effect of SEM on Learning Dynamics in Deep RL. We empirically evaluate the impact of SEM on the stability and efficiency of actor-critic algorithms. Our analysis combines both *learning performance* (returns, losses, TD error, critic disagreement) and *representation quality* (effective rank, feature norms), allowing us to connect sample-efficiency gains to underlying representational dynamics. This dual perspective highlights not only *whether* SEM improves performance, but also *why* it stabilizes training. A detailed explanation of each metric is provided in App. G.

To understand *why* SEM improves performance, we turn to representation-level diagnostics. Fig. 4 shows that SEM increases the effective rank of actor features, and bounds actor feature norms. Late in training, SEM also lifts the critic effective rank, a signs of more expressive and robust value learning. High effective rank is a proxy for avoiding representational collapse (Moalla et al., 2024). In the RL literature, representation collapse under drift has been empirically associated with capacity loss (Lyle et al., 2021), deterioration of feature rank (Kumar et al., 2021b), and implicit under-parameterization (Kumar et al., 2021a). In supervised and self-supervised settings, techniques like orthogonality regularization and rank-preserving weight regularizers are used to prevent feature collapse (He et al., 2024). These representational patterns align with our formal analysis, showing that SEM prevents covariance deflation and sustains gradient energy, thereby preventing feature collapse and boosting performance.

As shown in Fig. 5, SEM improves optimization stability over the baseline. Agents with SEM achieve higher returns earlier and maintain smaller, more stable TD errors, reduced critic disagreement, and lower critic-distribution discrepancy. Such effects are crucial, as instability in bootstrapped critics is a primary failure mode of actor-critic methods (Fujimoto et al., 2019; Kumar et al., 2021a). By constraining representation geometry, SEM produces better-conditioned features that yield more calibrated value estimates, echoing similar findings in representation regularization for deep RL (Anand et al., 2019; Laskin et al., 2020; Schwarzer et al., 2021). These results indicate that SEM not only accelerates learning but also yields more calibrated value estimates, mitigating instability in bootstrapped critics.

In Fig. 6, we focus our lens on the SEM module itself and examine how it shapes representations and action behavior. As training proceeds, the SEM layer’s activations become markedly sparser (higher Gini (Hurley & Rickard, 2009; Zonoobi et al., 2011)) and more sharply peaked (lower simplex

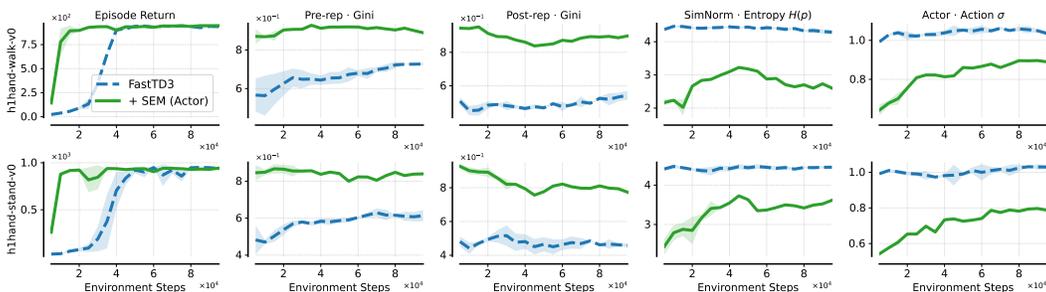


Fig. 6: **Sparsity, entropy, and action std on 2 HumanoidBench tasks.** SEM agents achieve higher returns with sparser features, lower entropy, and more stable action scales.

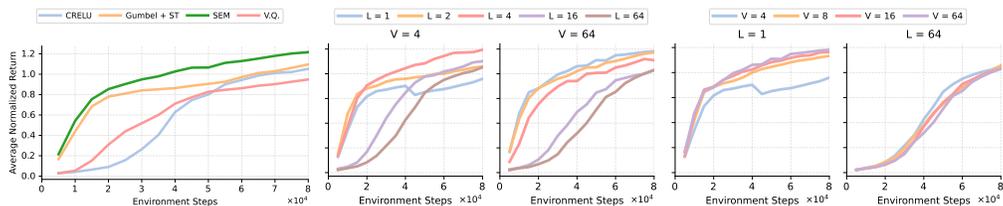


Fig. 7: **Aggregated average return on 5 HumanoidBench tasks.** We constrain the encoder’s output of the actor. (left) SEM outperforms alternative methods to impart structure on the encoder’s output. In SEM, representation capacity scales with $L \times V$ since the embedding consists of L simplex groups of size V . When varying L with fixed V (middle panels), performance improves as L increases from low-capacity regimes, and then saturates once $L \times V$ is sufficiently large. When varying V with fixed L (right panels), we observe the same pattern: for small L (e.g., $L=1$), increasing V noticeably improves performance, whereas for large L (e.g., $L=64$), all values of V perform similarly since the model already operates in a high-capacity regime.

entropy), while the overall action variance from the policy also declines. This trend is consistent with SEM’s design, where the block-wise softmax promotes competition and selective activation. As a result, the module imposes structured, energy-preserving constraints on its layer, encouraging more decisive feature usage and reducing noise in the downstream policy mapping.

Interestingly, this pattern also resonates with prior work in RL and representation learning. Hernandez-Garcia & Sutton (2019) show that enforcing sparsity in representations can improve robustness and mitigate interference in Q-learning settings. Moreover, recent studies on sparse architectures in deep RL such find that appropriately structured sparsity can enhance training stability and efficiency (Graesser et al., 2022; Ceron et al., 2024b;a; Ma et al., 2025).

Comparing SEM to other Regularization Methods To contextualize the benefits of simplicial embeddings, we compare SEM to alternative methods to induce structure on the encoder’s output. We compare SEM to commonly used methods for learning discrete explicit representations: Gumbel + straight-through (Jang et al., 2017; Maddison et al., 2017) and Vector Quantization (van den Oord et al., 2018). We also compare SEM to C-RELU (Abbas et al., 2023) which have been shown to improve the representation’s stability. We present the results in Fig. 7 (left) and find SEM to be more efficient and to lead to higher return than alternative methods. We conjecture that such improvement over Gumbel + ST and Vector quantization can be attributed to the fact that SEM does not necessitate the use of the straight-through estimator.

Analyzing SEM Parameters in Deep RL Lavoie et al. (2023) highlighted the effect of the simplex dimensionality V and number of simplices L , which jointly control sparsity and capacity of the representation. Investigating these parameters in deep RL is essential to understand how SEM balances representation capacity and stability under non-stationary training, and whether the same tradeoffs observed in self-supervised representation learning extend to RL. We study the effect of varying V and L in Fig. 7 (middle and right, respectively). Our results show that performance is driven primarily by the total representational capacity $L \times V$. In low-capacity regimes (e.g., $L = 1$), increasing V provides clear gains, consistent with the need for additional expressive power. How-

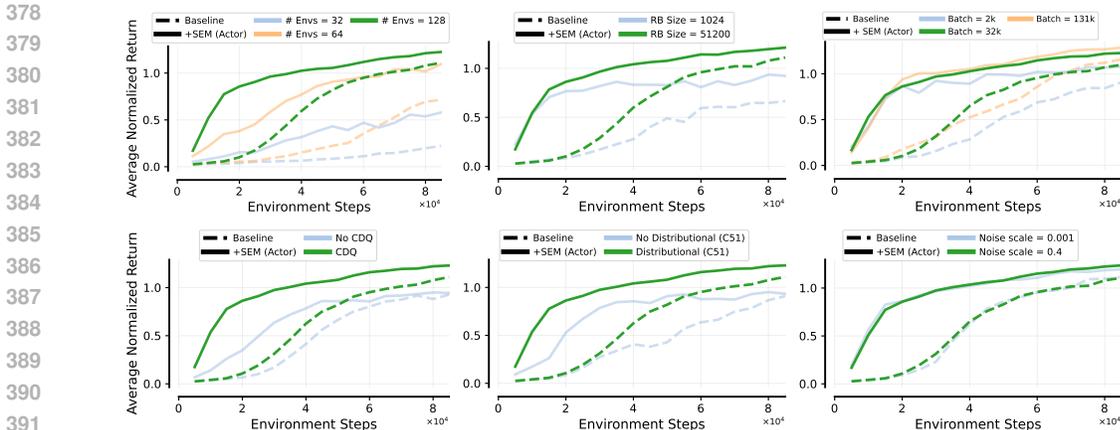


Fig. 8: Effect of core design choices on FastTD3 with and without SEM. SEM (solid green) consistently improves sample efficiency and asymptotic return across all settings, showing robustness to both hyperparameter variation and architectural design choices.

ever, once capacity is sufficiently large (e.g., $L = 64$), the effect of V largely saturates, different choices of V yield nearly identical performance; and in some cases smaller V (e.g., $V = 4$) performs slightly better. Similarly, increasing L improves performance when $L \times V$ is small, but the gains taper off once the model enters a high-capacity regime.

FastTD3 Design Choices and Simplicial Embeddings. FastTD3 extends TD3 with several design choices that improve throughput and stability, including parallel simulation, large-batch training, and distributional critics (Seo et al., 2025). These modifications enable actor-critic learning to scale efficiently in wall-clock time, but they do not address the geometry of the learned representations. In this section, we analyze how SEM complements FastTD3 by regularizing representation space and evaluate its effectiveness across the algorithmic design choices. In Fig. 8, we observe that SEM outperforms the baseline even when the agent is trained with reduced data availability (fewer environments, smaller replay buffers, or smaller batch sizes). Comparable gains also appear when algorithmic design choices such as CDQ and C51 are removed. These results demonstrate the robustness of SEM across both data-limited and simplified agent settings.

5 EMPIRICAL EVALUATION

We further evaluate the effectiveness and generality of SEM across a diverse set of deep RL algorithms and environments. Our study spans both off-policy and on-policy methods, including FastTD3, FastTD3-SimBaV2, FastSAC (Seo et al., 2025), and PPO (Schulman et al., 2017). Experiments are conducted on challenging humanoid benchmarks (28-h1hand tasks), (Sferrazza et al., 2024), IsaacLab (Mittal et al., 2023), IsaacGym suite (Makoviychuk et al., 2021), MTBench (Joshi et al., 2025), and an extended Atari-10 setup comprising 28 ALE games (see D.0.3 for more details) (Bellemare et al., 2013; Aitchison et al., 2023; Fedus et al., 2020), covering both continuous-control and pixel-based settings. Following prior work (Seo et al., 2025; Castanyer et al., 2025), we evaluate continuous-control tasks with six seeds and Atari results with three seeds, and aggregate performance across environments is reported, with shaded areas representing 95% confidence intervals. Full environment details and hyperparameter configurations are provided in App. I.

Fast Actor-Critic Algorithms. We first evaluate SEM on the HumanoidBench benchmark using three recent fast actor-critic baselines: FastTD3, FastTD3-SimBaV2, and FastSAC (Seo et al., 2025). These algorithms represent compute-efficient variants of TD3 and SAC, designed to scale with parallel simulation while maintaining strong performance on high-dimensional humanoid control. FastTD3-SimBaV2 incorporates hyperspherical normalization and reward scaling to accelerate critic training and stabilize optimization (Lee et al., 2025b); and FastSAC adapts the entropy-regularized SAC framework with similar throughput-oriented design choices, achieving high parallel efficiency while preserving training stability.

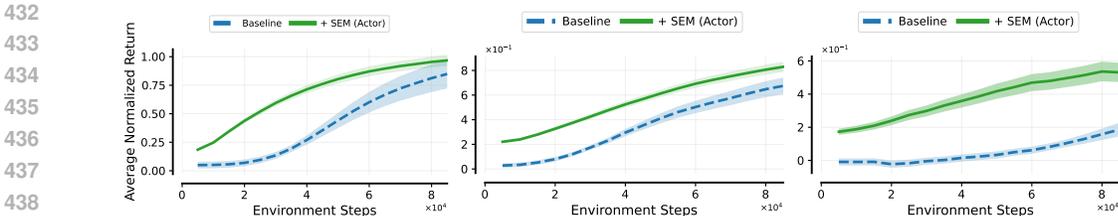


Fig. 9: **SEM on fast actor-critic algorithms.** Average normalized return on HumanoidBench with FastTD3 (left), FastTD3-SimBa (middle), and FastSAC (right). SEM consistently improves sample efficiency and yields higher final performance across all algorithms.

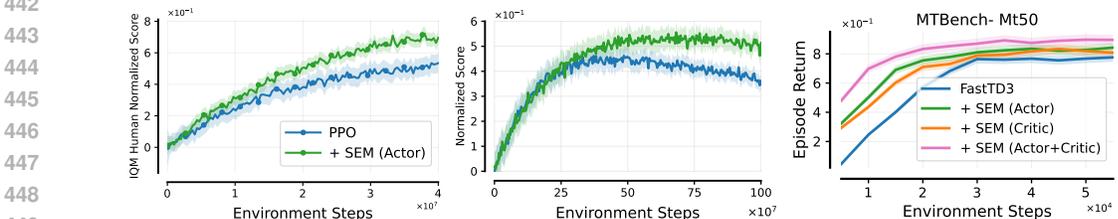


Fig. 10: Performance of PPO with and without SEM across tasks. **Left:** Atari experiments on an extended Atari-10 setup (28 ALE games; pixel-based). **Center:** PPO in IsaacGym. **Right:** MTBench MT50 (robotics tasks) comparing FastTD3. Applied SEM accelerates learning and improves return over the baseliens.

Across all three baselines, integrating SEM into the actor consistently accelerates early learning and improves asymptotic return. As shown in Fig. 9, SEM agents not only converge faster than their respective baselines, but also maintain lower variance across seeds. These results demonstrate that SEM provides complementary benefits to fast actor-critic methods, enhancing both stability and sample efficiency without modifying their underlying optimization procedures (see App. J for per-task learning curves). We also evaluate FastTD3 on 12-h1, 12-g1 tasks and 9-IsaacGym tasks, where a similar pattern is observed, as shown in App. K.

Proximal Policy Optimization Algorithm. To evaluate the generality of SEM beyond off-policy methods, we integrate it into PPO (Schulman et al., 2017), a popular on-policy method, using the CleanRL implementation (Huang et al., 2022). We evaluate SEM on two distinct benchmarks, IsaacGym for continuous control and the ALE (Bellemare et al., 2013) for pixel-based discrete control in Atari games. In both domains, SEM improves PPO by accelerating convergence and increasing final returns. The per-environment learning curves are shown in Fig. 29. Aggregate results are summarized in Fig. 10, with the left panel showing performance gains on the ALE suite², and the middle panel showing improvements on the IsaacGym tasks. These results demonstrate that SEM’s benefits are not limited to TD3-style critics but extend to policy-gradient methods and vision-based RL, underscoring its broad applicability.

Multitask Deep RL. Recent work by Joshi et al. (2025) introduced a large-scale benchmark for multi-task reinforcement learning (MTRL) in robotics. Implemented in IsaacGym, this benchmark comprises over seventy robotic control problems spanning both manipulation and locomotion, with subsets such as MT50 focused on manipulation. We compare FastTD3 (Seo et al., 2025) to its SEM-augmented variants (+SEM). As shown in Fig. 10 (right), +SEM improves sample efficiency, achieving faster learning and higher returns within the same training budget.

Value-Based Deep RL. Although our evaluation centers on actor-critic methods, this choice reflects experimental scope rather than a conceptual limitation. SEM operates on learned representations and is compatible with value-based algorithms. Extending SEM to these settings, such as DQN-style architectures (Mnih et al., 2013), is an important avenue for improving sample efficiency. To explore this direction, we evaluated SEM within PQN (Gallici et al., 2025), a recently proposed value-based algorithm that simplifies deep temporal-difference learning by streamlining target com-

²For the SEM module, we use $L = 128$ and $V = 4$ when evaluating ALE games (Bellemare et al., 2013) with PPO, whose penultimate fully connected layer has dimensionality 512.

486 putation and reducing unnecessary complexity. Our preliminary PQN experiments indicate that,
487 while SEM yields improvements in a few games (see Fig. 36), it does not provide consistent gains
488 in this regime overall (see Fig. 37). This suggests that a more systematic investigation is needed to
489 determine when geometric constraints benefit value-based methods. We evaluated PQN on 28 Atari
490 games (3 seeds) under different choices of V .

491 492 6 DISCUSSION 493

494 Our results demonstrate that geometric priors on representation space can substantially improve
495 the efficiency of deep RL agents. By constraining features to a product of simplices, SEM yields
496 bounded and sparse embeddings that avoid feature collapse and neuron dormancy under non-
497 stationarity. This lightweight inductive bias requires no auxiliary losses, adds effectively zero com-
498 putational cost (see Table 2), and consistently improves sample efficiency and asymptotic return
499 across various actor-critic methods and a diverse set of benchmarks.

500 Unlike existing model-based approaches in RL which use discrete state-embeddings (Hansen et al.,
501 2023; Hafner et al., 2020; 2023; Scannell et al., 2025), SEM does not require auxiliary objectives or
502 additional networks. Surprisingly, we find that the benefits of SEM are most pronounced when ap-
503 plied to the actor’s penultimate layer, where feature geometry most directly shapes policy gradients.
504 Our analyses indicate that SEM alleviates several optimization difficulties in deep RL (Moalla et al.,
505 2024; Juliani & Ash, 2024). By preserving effective rank, bounding feature norms, and reducing
506 critic disagreement, SEM provides more reliable gradients and stabilizes the bootstrapped targets
507 that often undermine critic training. These effects highlight representation geometry as a simple but
508 powerful lever for stabilizing learning under non-stationarity.

509 **Limitations and Future Work.** SEM is not a universal remedy. In tasks with extreme distribution
510 shift or very sparse rewards, feature collapse and critic drift may still occur, and SEM introduces
511 hyperparameters (L, V, τ) that require light tuning to balance sparsity and capacity. Moreover, our
512 experiments focus on continuous control and Atari; its impact on large-scale vision or language-
513 conditioned RL remains untested. Future work should investigate adaptive schedules for (L, V, τ) ,
514 and integration in more general-purpose algorithms such as MR.Q (Fujimoto et al., 2025), which
515 combine multiple objectives and scale across domains. Another direction is to understand when and
516 whether SEM benefits value-based algorithms, and to explore both its potential for scaling network
517 architectures (Ceron et al., 2024a) and its interaction with architectural priors (e.g., MoEs, Residual
518 Nets, Sparse Models) (Ceron et al., 2024b;a; Castanyer et al., 2025; Kooi et al., 2025). Our initial
519 experiments applying SEM to MoE-based architectures and pruned networks show limited gains,
520 highlighting the need for a more detailed study of how SEM interacts with modular routing and
521 sparse activation patterns (see ??).

522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

540 ETHICS STATEMENT

541

542 This paper presents work whose goal is to advance the field of Machine Learning, and reinforcement
543 learning in particular. There are many potential societal consequences of our work, none which we
544 feel must be specifically highlighted here.

545

546 REPRODUCIBILITY STATEMENT

547

548 We provide all the details to reproduce our results in the Appendix.

549

550 LLM USE

551

552 LLMs were used to assist paper editing and to write the code for plotting experiments.

553

554 REFERENCES

555

556 Zaheer Abbas, Rosie Zhao, Joseph Modayil, Adam White, and Marlos C Machado. Loss of plasticity
557 in continual deep reinforcement learning. In *Conference on lifelong learning agents*, pp. 620–636.
558 PMLR, 2023.

559

560 Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Bellemare.
561 Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural Informa-
562 tion Processing Systems*, 34, 2021.

563

564 Matthew Aitchison, Penny Sweetser, and Marcus Hutter. Atari-5: Distilling the arcade learning
565 environment down to five games. In *International Conference on Machine Learning*, pp. 421–
566 438. PMLR, 2023.

567

568 Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron,
569 Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube with a
robot hand. *arXiv preprint arXiv:1910.07113*, 2019.

570

571 Ankesh Anand, Evan Racah, Sherjil Ozair, Yoshua Bengio, Marc-Alexandre Côté, and R Devon
572 Hjelm. Unsupervised state representation learning in atari. *Advances in neural information pro-
573 cessing systems*, 32, 2019.

574

575 Jordan Ash and Ryan P Adams. On warm-starting neural network training. *Advances in neural
576 information processing systems*, 33:3884–3894, 2020.

577

578 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
579 learning to align and translate, 2016. URL <https://arxiv.org/abs/1409.0473>.

580

581 Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning envi-
582 ronment: an evaluation platform for general agents. *J. Artif. Int. Res.*, 47(1):253–279, May 2013.
583 ISSN 1076-9757.

584

585 Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement
586 learning. In *International conference on machine learning*, pp. 449–458. PMLR, 2017.

587

588 Emmanuel Bengio, Joelle Pineau, and Doina Precup. Interference and generalization in temporal
589 difference learning. In *International Conference on Machine Learning*, pp. 767–777. PMLR,
590 2020.

591

592 Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients
593 through stochastic neurons for conditional computation, 2013. URL [https://arxiv.org/
abs/1308.3432](https://arxiv.org/abs/1308.3432).

594

595 Roger Creus Castanyer, Johan Obando-Ceron, Lu Li, Pierre-Luc Bacon, Glen Berseth, Aaron
596 Courville, and Pablo Samuel Castro. Stable gradients for stable learning at scale in deep rein-
597 forcement learning. In *The Thirty-ninth Annual Conference on Neural Information Processing
598 Systems*, 2025. URL <https://openreview.net/forum?id=Vqj65VeDOu>.

599

- 594 Pablo Samuel Castro, Tyler Kastner, Prakash Panangaden, and Mark Rowland. Mico: Improved
595 representations via sampling-based state similarity for markov decision processes. *Advances in*
596 *Neural Information Processing Systems*, 34:30113–30126, 2021.
- 597
598 Johan Samir Obando Ceron, Aaron Courville, and Pablo Samuel Castro. In value-based deep rein-
599 forcement learning, a pruned network is a good network. In *International Conference on Machine*
600 *Learning*, pp. 38495–38519. PMLR, 2024a.
- 601
602 Johan Samir Obando Ceron, Ghada Sokar, Timon Willi, Clare Lyle, Jesse Farebrother, Jakob Nico-
603 laus Foerster, Gintare Karolina Dziugaite, Doina Precup, and Pablo Samuel Castro. Mixtures of
604 experts unlock parameter scaling for deep rl. In *International Conference on Machine Learning*,
pp. 38520–38540. PMLR, 2024b.
- 605
606 Will Dabney, Georg Ostrovski, David Silver, and Rémi Munos. Implicit quantile networks for
607 distributional reinforcement learning. In *International conference on machine learning*, pp. 1096–
608 1105. PMLR, 2018a.
- 609
610 Will Dabney, Mark Rowland, Marc Bellemare, and Rémi Munos. Distributional reinforcement
611 learning with quantile regression. In *Proceedings of the AAAI conference on artificial intelligence*,
volume 32, 2018b.
- 612
613 Shihbansh Dohare, J Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A Rupam Mah-
614 mood, and Richard S Sutton. Loss of plasticity in deep continual learning. *Nature*, 632(8026):
768–774, 2024.
- 615
616 D.L. Donoho, M. Elad, and V.N. Temlyakov. Stable recovery of sparse overcomplete representations
617 in the presence of noise. *IEEE Transactions on Information Theory*, 52(1):6–18, 2006. doi:
618 10.1109/TIT.2005.860430.
- 619
620 Pierluca D’Oro, Max Schwarzer, Evgenii Nikishin, Pierre-Luc Bacon, Marc G Bellemare, and
621 Aaron Courville. Sample-efficient reinforcement learning by breaking the replay ratio barrier.
In *Deep Reinforcement Learning Workshop NeurIPS 2022*, 2022.
- 622
623 Mohamed Elsayed, Qingfeng Lan, Clare Lyle, and A Rupam Mahmood. Weight clipping for deep
624 continual and reinforcement learning. *arXiv preprint arXiv:2407.01704*, 2024.
- 625
626 Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam
627 Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with im-
628 portance weighted actor-learner architectures. In *International conference on machine learning*,
pp. 1407–1416. PMLR, 2018.
- 629
630 Lasse Espeholt, Raphaël Marinier, Piotr Stanczyk, Ke Wang, and Marcin Michalski. Seed rl:
631 Scalable and efficient deep-rl with accelerated central inference. In *International Confer-
632 ence on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rkgvXlrKwH>.
- 633
634 William Fedus, Prajit Ramachandran, Rishabh Agarwal, Yoshua Bengio, Hugo Larochelle, Mark
635 Rowland, and Will Dabney. Revisiting fundamentals of experience replay. In *Proceedings of the*
636 *37th International Conference on Machine Learning, ICML’20*. JMLR.org, 2020.
- 637
638 Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-
critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- 639
640 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without
exploration. In *International Conference on Machine Learning*, pp. 2052–2062, 2019.
- 641
642 Scott Fujimoto, Wei-Di Chang, Edward Smith, Shixiang Shane Gu, Doina Precup, and David Meger.
643 For sale: State-action representation learning for deep reinforcement learning. *Advances in neural*
644 *information processing systems*, 36:61573–61624, 2023.
- 645
646 Scott Fujimoto, Pierluca D’Oro, Amy Zhang, Yuandong Tian, and Michael Rabbat. Towards
647 general-purpose model-free reinforcement learning. In *The Thirteenth International Conference*
on Learning Representations (ICLR), 2025. URL <https://openreview.net/forum?id=RlhIXdST22>.

- 648 Kunihiko Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of
649 pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36:193–202, 1980.
650
- 651 Marco Fumero, Florian Wenzel, Luca Zancato, Alessandro Achille, Emanuele Rodolà, Stefano
652 Soatto, Bernhard Schölkopf, and Francesco Locatello. Leveraging sparse and shared feature
653 activations for disentangled representation learning. In *Thirty-seventh Conference on Neural
654 Information Processing Systems*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=IHR83ufYPy)
655 [IHR83ufYPy](https://openreview.net/forum?id=IHR83ufYPy).
- 656 Alexandre Galashov, Michalis Titsias, András György, Clare Lyle, Razvan Pascanu, Yee Whye Teh,
657 and Maneesh Sahani. Non-stationary learning of neural networks with automatic soft parameter
658 reset. *Advances in Neural Information Processing Systems*, 37:83197–83234, 2024.
- 659 Matteo Gallici, Mattie Fellows, Benjamin Ellis, Bartomeu Pou, Ivan Masmitja, Jakob Nicolaus
660 Foerster, and Mario Martin. Simplifying deep temporal difference learning. In *The Thirteenth
661 International Conference on Learning Representations*, 2025. URL [https://openreview.](https://openreview.net/forum?id=7IzeL0kflu)
662 [net/forum?id=7IzeL0kflu](https://openreview.net/forum?id=7IzeL0kflu).
- 663 Samuel Garcin, Trevor McInroe, Pablo Samuel Castro, Christopher G. Lucas, David Abel, Prakash
664 Panangaden, and Stefano V Albrecht. Studying the interplay between the actor and critic rep-
665 resentations in reinforcement learning. In *The Thirteenth International Conference on Learning
666 Representations*, 2025. URL <https://openreview.net/forum?id=tErHYBGlWc>.
- 667 Florin Gogianu, Tudor Berariu, Mihaela C Rosca, Claudia Clopath, Lucian Busoniu, and Razvan
668 Pascanu. Spectral normalisation for deep reinforcement learning: an optimisation perspective. In
669 *International Conference on Machine Learning*, pp. 3734–3744. PMLR, 2021.
- 670 Laura Graesser, Utku Evci, Erich Elsen, and Pablo Samuel Castro. The state of sparse training in
671 deep reinforcement learning. In *International Conference on Machine Learning*, pp. 7766–7792.
672 PMLR, 2022.
- 673 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
674 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-
675 ence on machine learning*, pp. 1861–1870. Pmlr, 2018.
- 676 Danijar Hafner, Timothy P Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with
677 discrete world models. In *International Conference on Learning Representations*, 2020.
- 678 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
679 through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 680 Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for contin-
681 uous control. *arXiv preprint arXiv:2310.16828*, 2023.
- 682 Junlin He, Jinxiao Du, and Wei Ma. Preventing dimensional collapse in self-supervised learning
683 via orthogonality regularization. In *The Thirty-eighth Annual Conference on Neural Information
684 Processing Systems*, 2024. URL <https://openreview.net/forum?id=Y3FjKSsfmy>.
- 685 Peter Henderson, Jieru Hu, Joshua Romoff, Emma Brunskill, Dan Jurafsky, and Joelle Pineau. To-
686 wards the systematic reporting of the energy and carbon footprints of machine learning. *Journal
687 of Machine Learning Research*, 21(248):1–43, 2020.
- 688 D Hendrycks. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- 689 J Fernando Hernandez-Garcia and Richard S Sutton. Learning sparse representations incrementally
690 in deep reinforcement learning. *arXiv preprint arXiv:1912.04002*, 2019.
- 691 Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdi-
692 nov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint
693 arXiv:1207.0580*, 2012.
- 694 Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Ki-
695 nal Mehta, and JoÃ£o GM AraÃ§o. Cleanrl: High-quality single-file implementations of deep
696 reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022.
- 697
698
699
700
701

- 702 Niall Hurley and Scott Rickard. Comparing measures of sparsity. *IEEE Transactions on Information*
703 *Theory*, 55(10):4723–4741, 2009.
- 704 Maximilian Igl, Gregory Farquhar, Jelena Luketina, Wendelin Boehmer, and Shimon Whiteson.
705 Transient non-stationarity and generalisation in deep reinforcement learning. In *International*
706 *Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Qun8fv4qSby)
707 [id=Qun8fv4qSby](https://openreview.net/forum?id=Qun8fv4qSby).
- 708 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax, 2017.
709 URL <https://arxiv.org/abs/1611.01144>.
- 710 Vira Joshi, Zifan Xu, Bo Liu, Peter Stone, and Amy Zhang. Benchmarking massively parallelized
711 multi-task reinforcement learning for robotics tasks. *arXiv preprint arXiv:2507.23172*, 2025.
- 712 Arthur Juliani and Jordan T. Ash. A study of plasticity loss in on-policy deep reinforcement learning.
713 In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL
714 <https://openreview.net/forum?id=MsUf8kpKTF>.
- 715 Jacob E. Kooi, Zhao Yang, and Vincent François-Lavet. Hadamax encoding: Elevating performance
716 in model-free atari, 2025. URL <https://arxiv.org/abs/2505.15345>.
- 717 Alex Krizhevsky. Learning multiple layers of features from tiny images. [https://www.cs.toronto.](https://www.cs.toronto.edu/kriz/learning-features-2009-TR.pdf)
718 [edu/kriz/learning-features-2009-TR.pdf](https://www.cs.toronto.edu/kriz/learning-features-2009-TR.pdf), 2009.
- 719 Aviral Kumar, Rishabh Agarwal, Dibya Ghosh, and Sergey Levine. Implicit under-parameterization
720 inhibits data-efficient deep reinforcement learning. In *International Conference on Learning Rep-*
721 *resentations*, 2021a. URL <https://openreview.net/forum?id=09bnihsFfXU>.
- 722 Aviral Kumar, Rishabh Agarwal, Tengyu Ma, Aaron Courville, George Tucker, and Sergey
723 Levine. DR3: Value-based deep reinforcement learning requires explicit regularization. In
724 *Deep RL Workshop NeurIPS 2021*, 2021b. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=LYwOCfpsQ-A)
725 [LYwOCfpsQ-A](https://openreview.net/forum?id=LYwOCfpsQ-A).
- 726 Aviral Kumar, Rishabh Agarwal, Xinyang Geng, George Tucker, and Sergey Levine. Offline q-
727 learning on diverse multi-task data both scales and generalizes. *arXiv preprint arXiv:2211.15144*,
728 2022a.
- 729 Aviral Kumar, Rishabh Agarwal, Xinyang Geng, George Tucker, and Sergey Levine. Offline q-
730 learning on diverse multi-task data both scales and generalizes. In *The Eleventh International*
731 *Conference on Learning Representations*, 2022b.
- 732 Saurabh Kumar, Henrik Marklund, and Benjamin Van Roy. Maintaining plasticity in continual
733 learning via regenerative regularization. *arXiv preprint arXiv:2308.11958*, 2023.
- 734 Michael Laskin, Aravind Srinivas, and Pieter Abbeel. Curl: Contrastive unsupervised representa-
735 tions for reinforcement learning. In *International conference on machine learning*, pp. 5639–
736 5650. PMLR, 2020.
- 737 Samuel Lavoie, Christos Tsirigotis, Max Schwarzer, Ankit Vani, Michael Noukhovitch, Kenji
738 Kawaguchi, and Aaron Courville. Simplicial embeddings in self-supervised learning and down-
739 stream classification. In *The Eleventh International Conference on Learning Representations*,
740 2023.
- 741 Samuel Lavoie, Michael Noukhovitch, and Aaron Courville. Compositional discrete latent code for
742 high fidelity, productive diffusion models, 2025. URL [https://arxiv.org/abs/2507.](https://arxiv.org/abs/2507.12318)
743 [12318](https://arxiv.org/abs/2507.12318).
- 744 Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel.
745 Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):541–
746 551, 1989. doi: 10.1162/neco.1989.1.4.541.
- 747 Hojoon Lee, Hanseul Cho, Hyunseung Kim, Daehoon Gwak, Joonkee Kim, Jaegul Choo, Se-
748 Young Yun, and Chulhee Yun. Plastic: Improving input and label plasticity for sample effi-
749 cient reinforcement learning. In *Neural Information Processing Systems*, 2023. URL [https://](https://api.semanticscholar.org/CorpusID:259203876)
750 api.semanticscholar.org/CorpusID:259203876.

- 756 Hojoon Lee, Dongyoon Hwang, Donghu Kim, Hyunseung Kim, Jun Jet Tai, Kaushik Subramanian,
757 Peter R. Wurman, Jaegul Choo, Peter Stone, and Takuma Seno. Simba: Simplicity bias for
758 scaling up parameters in deep reinforcement learning. In *The Thirteenth International Confer-*
759 *ence on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=jXLiDKsuDo>.
760
- 761 Hojoon Lee, Youngdo Lee, Takuma Seno, Donghu Kim, Peter Stone, and Jaegul Choo. Hyper-
762 spherical normalization for scalable deep reinforcement learning. In *Forty-second International*
763 *Conference on Machine Learning*, 2025b. URL <https://openreview.net/forum?id=kfYxyvCYQ4>.
764
- 765
766 Timothée Lesort, Natalia Díaz-Rodríguez, Jean-Francois Goudou, and David Filliat. State represen-
767 tation learning for control: An overview. *Neural Networks*, 108:379–392, 2018.
- 768
769 Shikun Li, Xiaobo Xia, Shiming Ge, and Tongliang Liu. Selective-supervised contrastive learning
770 with noisy labels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
771 *recognition*, pp. 316–325, 2022.
- 772
773 Zechu Li, Tao Chen, Zhang-Wei Hong, Anurag Ajay, and Pulkit Agrawal. Parallel q -learning:
774 Scaling off-policy reinforcement learning under massively parallel simulation. In *International*
775 *Conference on Machine Learning*, pp. 19440–19459. PMLR, 2023.
- 776
777 Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
778 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv*
preprint arXiv:1509.02971, 2015.
- 779
780 Jiashun Liu, Johan Samir Obando Ceron, Aaron Courville, and Ling Pan. Neuroplastic expansion
781 in deep reinforcement learning. In *The Thirteenth International Conference on Learning Repre-*
782 *sentations*, 2025a. URL <https://openreview.net/forum?id=20qZK2T7fa>.
- 783
784 Jiashun Liu, Johan Obando-Ceron, Pablo Samuel Castro, Aaron Courville, and Ling Pan. The
785 courage to stop: Overcoming sunk cost fallacy in deep reinforcement learning. In *Forty-second*
786 *International Conference on Machine Learning*, 2025b. URL <https://openreview.net/forum?id=VzC3BA9gf>.
- 787
788 Vincent Liu, Raksha Kumaraswamy, Lei Le, and Martha White. The utility of sparse represen-
789 tations for control in reinforcement learning. In *Proceedings of the Thirty-Third AAAI Con-*
790 *ference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelli-*
791 *gence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence*,
792 AAAI’19/IAAI’19/EAAI’19. AAAI Press, 2019. ISBN 978-1-57735-809-1. doi: 10.1609/aaai.v33i01.33014384. URL <https://doi.org/10.1609/aaai.v33i01.33014384>.
- 793
794 Jianlan Luo, Charles Xu, Jeffrey Wu, and Sergey Levine. Precise and dexterous robotic manipulation
795 via human-in-the-loop reinforcement learning. *Science Robotics*, 10(105):eads5033, 2025. doi:
796 10.1126/scirobotics.ads5033. URL <https://www.science.org/doi/abs/10.1126/scirobotics.ads5033>.
- 797
798 Clare Lyle, Mark Rowland, and Will Dabney. Understanding and preventing capacity loss
799 in reinforcement learning. In *Deep RL Workshop NeurIPS 2021*, 2021. URL https://openreview.net/forum?id=5G7fT_tJTt.
800
- 801
802 Clare Lyle, Mark Rowland, Will Dabney, Marta Kwiatkowska, and Yarin Gal. Learning dynam-
803 ics and generalization in deep reinforcement learning. In *International conference on machine*
learning, pp. 14560–14581. PMLR, 2022.
- 804
805 Clare Lyle, Zeyu Zheng, Evgenii Nikishin, Bernardo Avila Pires, Razvan Pascanu, and Will Dabney.
806 Understanding plasticity in neural networks. In *International Conference on Machine Learning*,
807 pp. 23190–23211. PMLR, 2023.
- 808
809 Clare Lyle, Zeyu Zheng, Khimya Khetarpal, Hado van Hasselt, Razvan Pascanu, James Martens,
and Will Dabney. Disentangling the causes of plasticity loss in neural networks. In *Conference*
on Lifelong Learning Agents, pp. 750–783. PMLR, 2025.

- 810 Guozheng Ma, Lu Li, Zilin Wang, Li Shen, Pierre-Luc Bacon, and Dacheng Tao. Network spar-
811 sity unlocks the scaling potential of deep reinforcement learning. In *Forty-second International*
812 *Conference on Machine Learning*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=mIomqOskaa)
813 [mIomqOskaa](https://openreview.net/forum?id=mIomqOskaa).
- 814 Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relax-
815 ation of discrete random variables, 2017. URL <https://arxiv.org/abs/1611.00712>.
- 816 Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin,
817 David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance
818 gpu based physics simulation for robot learning. In *Thirty-fifth Conference on Neural Information*
819 *Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- 820 Gabriel B Margolis, Ge Yang, Kartik Paigwar, Tao Chen, and Pulkit Agrawal. Rapid locomotion via
821 reinforcement learning. *The International Journal of Robotics Research*, 43(4):572–587, 2024.
- 822 Mayank Mittal, Calvin Yu, Qinxi Yu, Jingzhou Liu, Nikita Rudin, David Hoeller, Jia Lin Yuan,
823 Ritvik Singh, Yunrong Guo, Hammad Mazhar, et al. Orbit: A unified simulation framework for
824 interactive robot learning environments. *IEEE Robotics and Automation Letters*, 8(6):3740–3747,
825 2023.
- 826 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wier-
827 stra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*
828 *arXiv:1312.5602*, 2013.
- 829 Skander Moalla, Andrea Miele, Daniil Pyatko, Razvan Pascanu, and Caglar Gulcehre. No represen-
830 tation, no trust: Connecting representation, collapse, and trust issues in PPO. In *The Thirty-*
831 *eighth Annual Conference on Neural Information Processing Systems*, 2024. URL [https://](https://openreview.net/forum?id=Wy9UgrMwD0)
832 openreview.net/forum?id=Wy9UgrMwD0.
- 833 Brian Murphy, Partha Talukdar, and Tom Mitchell. Learning effective and interpretable semantic
834 models using non-negative sparse embedding. In Martin Kay and Christian Boitet (eds.), *Pro-*
835 *ceedings of COLING 2012*, pp. 1933–1950, Mumbai, India, December 2012. The COLING 2012
836 Organizing Committee. URL <https://aclanthology.org/C12-1118/>.
- 837 Michal Nauman, Michał Bortkiewicz, Piotr Miłoś, Tomasz Trzciński, Mateusz Ostaszewski, and
838 Marek Cygan. Overestimation, overfitting, and plasticity in actor-critic: the bitter lesson of rein-
839 forcement learning. In *Proceedings of the 41st International Conference on Machine Learning*,
840 pp. 37342–37364, 2024a.
- 841 Michal Nauman, Mateusz Ostaszewski, Krzysztof Jankowski, Piotr Miłoś, and Marek Cygan. Big-
842 ger, regularized, optimistic: scaling for compute and sample efficient continuous control. *Ad-*
843 *vances in neural information processing systems*, 37:113038–113071, 2024b.
- 844 Evgenii Nikishin, Max Schwarzer, Pierluca D’Oro, Pierre-Luc Bacon, and Aaron Courville. The
845 primacy bias in deep reinforcement learning. In *International conference on machine learning*,
846 pp. 16828–16847. PMLR, 2022.
- 847 Johan Obando Ceron, Marc Bellemare, and Pablo Samuel Castro. Small batch deep reinforcement
848 learning. *Advances in Neural Information Processing Systems*, 36:26003–26024, 2023.
- 849 Ivaylo Popov, Nicolas Heess, Timothy Lillicrap, Roland Hafner, Gabriel Barth-Maron, Matej Ve-
850 cerik, Thomas Lampe, Yuval Tassa, Tom Erez, and Martin Riedmiller. Data-efficient deep rein-
851 forcement learning for dexterous manipulation, 2017. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1704.03073)
852 [1704.03073](https://arxiv.org/abs/1704.03073).
- 853 Yi Ren, Samuel Lavoie, Mikhail Galkin, Danica J. Sutherland, and Aaron Courville. Improving
854 compositional generalization using iterated learning and simplicial embeddings, 2023. URL
855 <https://arxiv.org/abs/2310.18777>.
- 856 Nikita Rudin, David Hoeller, Philipp Reist, and Marco Hutter. Learning to walk in minutes using
857 massively parallel deep reinforcement learning. In *5th Annual Conference on Robot Learning*,
858 2021. URL <https://openreview.net/forum?id=wK2fDDJ5VcF>.

- 864 Aidan Scannell, Mohammadreza Nakhaeizhadfard, Kalle Kujanpää, Yi Zhao, Kevin Sebastian
865 Luck, Arno Solin, and Joni Pajarinen. Discrete codebook world models for continuous control.
866 In *The Thirteenth International Conference on Learning Representations*, 2025.
- 867 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
868 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 870 Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. Green ai. *Communications of the*
871 *ACM*, 63(12):54–63, 2020.
- 872 Max Schwarzer, Ankesh Anand, Rishab Goel, R Devon Hjelm, Aaron Courville, and Philip Bach-
873 man. Data-efficient reinforcement learning with self-predictive representations. In *The Ninth*
874 *International Conference on Learning Representations (ICLR)*, 2021.
- 876 Younggyo Seo, Carmelo Sferrazza, Haoran Geng, Michal Nauman, Zhao-Heng Yin, and Pieter
877 Abbeel. Fasttd3: Simple, fast, and capable reinforcement learning for humanoid control. *arXiv*
878 *preprint arXiv:2505.22642*, 2025.
- 880 Carmelo Sferrazza, Dun-Ming Huang, Xingyu Lin, Youngwoon Lee, and Pieter Abbeel. Humanoid-
881 bench: Simulated humanoid benchmark for whole-body locomotion and manipulation. *arXiv*
882 *preprint arXiv:2403.10506*, 2024.
- 883 Jayesh Singla, Ananye Agarwal, and Deepak Pathak. Sapg: Split and aggregate policy gradients. In
884 *International Conference on Machine Learning*, pp. 45759–45772. PMLR, 2024.
- 885 Laura Smith, Ilya Kostrikov, and Sergey Levine. A walk in the park: Learning to walk in 20
886 minutes with model-free reinforcement learning, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2208.07860)
887 [2208.07860](https://arxiv.org/abs/2208.07860).
- 889 Ghada Sokar and Pablo Samuel Castro. Mind the gap! the challenges of scale in pixel-based deep
890 reinforcement learning. *arXiv preprint arXiv:2505.17749*, 2025.
- 891 Ghada Sokar, Rishabh Agarwal, Pablo Samuel Castro, and Utku Evci. The dormant neuron phe-
892 nomenon in deep reinforcement learning. In *International Conference on Machine Learning*, pp.
893 32145–32168. PMLR, 2023.
- 895 Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Do-
896 main randomization for transferring deep neural networks from simulation to the real world. In
897 *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pp. 23–30.
898 IEEE, 2017.
- 899 Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learn-
900 ing, 2018. URL <https://arxiv.org/abs/1711.00937>.
- 902 Maxime Wabartha and Joelle Pineau. Piecewise linear parametrization of policies: Towards in-
903 terpretable deep reinforcement learning. In *The Twelfth International Conference on Learning*
904 *Representations*, 2024. URL <https://openreview.net/forum?id=hOMVq57Ce0>.
- 905 Timon Willi, Johan Samir Obando Ceron, Jakob Nicolaus Foerster, Gintare Karolina Dziugaite, and
906 Pablo Samuel Castro. Mixture of experts in a mixture of RL settings. In *Reinforcement Learning*
907 *Conference*, 2024. URL <https://openreview.net/forum?id=5FF06R10Em>.
- 909 Yuzhe Yang, Guo Zhang, Zhi Xu, and Dina Katabi. Harnessing structures for value-based planning
910 and reinforcement learning. *arXiv preprint arXiv:1909.12255*, 2019.
- 911 Qingmao Yao, Zhichao Lei, Tianyuan Chen, Ziyue Yuan, Xuefan Chen, Jianxiang Liu, Faguo Wu,
912 and Xiao Zhang. Offline RL with smooth OOD generalization in convex hull and its neighbor-
913 hood. In *The Thirteenth International Conference on Learning Representations*, 2025. URL
914 <https://openreview.net/forum?id=eY5JNJE56i>.
- 915 Denis Yarats, Ilya Kostrikov, and Rob Fergus. Image augmentation is all you need: Regularizing
916 deep reinforcement learning from pixels. In *International Conference on Learning Representa-*
917 *tions*, 2021. URL <https://openreview.net/forum?id=GY6-6sTvGaf>.

918 Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous control: Improved data-augmented reinforcement learning. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=_SJ-_yyes8.
919
920
921
922 Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games with limited data. *Advances in neural information processing systems*, 34:25476–25488, 2021.
923
924 Chiyuan Zhang, Oriol Vinyals, Remi Munos, and Samy Bengio. A study on overfitting in deep reinforcement learning. *arXiv preprint arXiv:1804.06893*, 2018.
925
926
927 Ziwen Zhuang, Zipeng Fu, Jianren Wang, Christopher G Atkeson, Sören Schwertfeger, Chelsea Finn, and Hang Zhao. Robot parkour learning. In *7th Annual Conference on Robot Learning*, 2023. URL <https://openreview.net/forum?id=uo937r5eTE>.
928
929
930 Dornoosh Zonoobi, Ashraf A Kassim, and Yedatore V Venkatesh. Gini index as sparsity measure for signal reconstruction from compressive samples. *IEEE Journal of Selected Topics in Signal Processing*, 5(5):927–932, 2011.
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

APPENDIX CONTENTS

A Related Work	21
B Formal Analysis	22
C Fast Actor–Critic	23
D Benchmarks	23
E Deep RL Network Architectures	25
F Ablation Setup	25
G Metrics	25
G.1 Feature Rank	26
G.2 Dormant Neurons	26
G.3 Weight Norm	26
G.4 Gini Sparsity	27
G.5 Cramer distance	27
G.6 Entropy	27
H Additional Baselines	28
H.1 Regen	28
H.2 Shrink and Perturb	28
H.3 Weight Clip	29
H.4 Reset Layer	29
H.5 Spectral Norm	30
H.6 Feature norm	30
H.7 Gelu	31
H.8 Dropout	31
H.9 MoEs	32
H.10 ReDO	32
H.11 Magnitude Pruning	33

1026	I Hyperparameters	34
1027		
1028	J Learning Curves for Each Game	35
1029		
1030		
1031	K Additional Experiments on HumanoidBench	39
1032		
1033	L Evaluating Layer-wise Interventions	42
1034		
1035	M SEM on Value-Based Deep RL	43
1036		

1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

1080 A RELATED WORK

1081
1082 **Stability in Deep RL** A longstanding challenge in reinforcement learning is the stability of value-
1083 based updates in actor–critic methods. One major source of instability is overestimation bias, which
1084 accumulates when bootstrapped critics reinforce overly optimistic targets. Twin Delayed DDPG
1085 (TD3) (Fujimoto et al., 2018) mitigates this issue with clipped double Q-learning and delayed pol-
1086 icy updates, producing more reliable critics and improving control performance. Another direction
1087 seeks to stabilize targets by modeling full return distributions rather than point estimates. Distri-
1088 butional RL methods such as C51 (Bellemare et al., 2017), QR-DQN (Dabney et al., 2018b), and
1089 IQN (Dabney et al., 2018a) show that capturing the shape of the return distribution reduces variance
1090 and provides richer learning signals.

1091 More recently, architectural choices have been used to enhance critic stability. SimBaV2 (Lee et al.,
1092 2025b) biases networks toward simpler, well-conditioned representations via input normalization,
1093 linear residual paths, and feature normalization, helping large models avoid divergence. Meanwhile,
1094 BRO shows that scaling critic capacity paired with strong regularization and optimistic updates
1095 yields dramatically better sample efficiency in continuous control tasks (Nauman et al., 2024b). In
1096 parallel, Ceron et al. (2024a;b); Ma et al. (2025) demonstrate that applying static network spar-
1097 sity unlocks scaling of deep RL models by mitigating gradient interference (Bengio et al., 2020;
1098 Obando Ceron et al., 2023) and plasticity loss (Lyle et al., 2023). Training regimens also play a role;
1099 SR-SPR (D’Oro et al., 2022) demonstrates that periodic network resets counteract bootstrapping
1100 drift, allowing agents to sustain extremely high replay ratios without collapse. FastTD3 (Seo et al.,
1101 2025) integrates several of these lessons, combining parallel simulation, large-batch updates, and
1102 distributional critics to achieve strong stability at high throughput. Our approach is complementary
1103 to these efforts. Rather than modifying update schedules or ensemble targets, we constrain the ge-
1104 ometry of latent representations, aiming to reduce critic variance and stabilize bootstrapped updates
1105 through structured embeddings.

1106 **Sample-Efficient RL** Beyond stability, a parallel line of work targets sample efficiency, with
1107 progress spanning both representation-driven methods and algorithmic or model-based improve-
1108 ments. Representation learning has emerged as a powerful way to extract more information per
1109 interaction. CURL (Laskin et al., 2020) applies contrastive learning to enforce invariances in en-
1110 coders trained jointly with the control objective, significantly narrowing the gap between pixel- and
1111 state-based agents. SPR (Schwarzer et al., 2021) extends this idea with self-predictive latent dynam-
1112 ics, ensuring temporal consistency and yielding state-of-the-art data efficiency on Atari. Building
1113 on SPR, SR-SPR (D’Oro et al., 2022) adds scheduled resets that prevent drift and enable aggressive
1114 replay-ratio scaling. Other works inject architectural or learning biases for example; SimBa (Lee
1115 et al., 2025a) introduces simplicity constraints that regularize feature representations, allowing larger
1116 networks to remain well-conditioned under nonstationary training; its successor, SimBaV2 (Lee
1117 et al., 2025b), further stabilizes scaling through hyperspherical normalization, constraining weights
1118 and activations to unit-norm manifolds and ensuring consistent gradient magnitudes across capac-
1119 ities. Neuroplastic Expansion (Liu et al., 2025a) complements these structural interventions by
1120 dynamically growing and pruning neurons to preserve long-term adaptability, while The Courage
1121 to Stop (Liu et al., 2025b) improves behavioral efficiency by terminating unproductive trajectories
1122 early, reducing replay-buffer contamination.

1123 For SALE (Fujimoto et al., 2023) further enriches the representation space with state–action embed-
1124 dings, producing TD7, which substantially outperforms TD3 in continuous control. Outside of RL,
1125 simplicial embeddings (Lavoie et al., 2023) show that constraining features to products of probability
1126 simplices induces sparse, group-structured representations that generalize effectively in supervised
1127 and self-supervised settings and leads to a compositional representation (Ren et al., 2023; Lavoie
1128 et al., 2025). We draw inspiration from this idea and adapt it to reinforcement learning, inserting
1129 simplicial modules into fast actor–critic pipelines.

1130 Algorithmic and model-based approaches provide another path to efficiency. Soft Actor-Critic
1131 (SAC) (Haarnoja et al., 2018) introduces maximum-entropy RL, balancing reward and exploration to
1132 achieve robust and data-efficient learning in continuous control. Model-based algorithms further im-
1133 prove efficiency by planning with learned dynamics. TD-MPC2 (Hansen et al., 2023) demonstrates
1134 that latent-space model predictive control scales effectively across diverse domains, achieving state-
1135 of-the-art performance with a single set of hyperparameters. EfficientZero (Ye et al., 2021) combines

MuZero-style search with learned latent dynamics, reaching human-level Atari performance with orders of magnitude fewer environment steps. Our method differs from these approaches by focusing on representation geometry: rather than auxiliary losses, ensembles, or world models, we show that a single simplicial bottleneck can consistently improve the sample efficiency of fast actor–critic algorithms while preserving their hallmark wall-clock advantages.

Structured representation in RL Constraining the encoder’s output is common in RL. C-ReLU has been shown to improve training and plasticity (Abbas et al., 2023). Feature normalization with L2 regularization of the features also improves training scalability and enables larger scale training of RL models. Closer to our work, DreamerV2 (Hafner et al., 2020) and DreamerV3 (Hafner et al., 2023) encode the observation into a one-hot discrete representation work. Scannell et al. (2025) also learn discrete latent space via a learned codebook and gumbel softmax with straight-through estimator. Wabartha & Pineau (2024) also propose to learn discrete encoding of the state for policy learning and show interpretable representations. However, methods with explicit discretization necessitate the use of a biased gradient estimator to propagate the learning signal inside the encoder. Similar to our work, Hansen et al. (2023) constrain the encoder’s output into SEM. In this work, we find that SEM is a crucial component for improving sample efficiency and performance in RL and study that component in details and connect the improved performance to the improved training stability coming from the sparse and structured representation endowed by SEM.

B FORMAL ANALYSIS

Theorem 1. *Non-stationarity increases neuron dormancy.*

Proof. Let \mathcal{D}_t be the data distribution at iteration t and consider a critic $f_\theta(x) = W h_\phi(x)$, trained by minimizing the (mean) squared error to targets $y_t(x)$:

$$\mathcal{L}_t(\theta) = \mathbb{E}_{x \sim \mathcal{D}_t} \left[(f_\theta(x) - y_t(x))^2 \right]. \quad (4)$$

Define the minimizer $\theta_t^* \in \arg \min_\theta \mathcal{L}_t(\theta)$ and tracking error $e_t = \theta_t - \theta_t^*$. A first-order expansion of SGD around θ_t^* gives

$$e_{t+1} \approx (I - \alpha H_t) e_t - \alpha b_t, \quad H_t = \nabla^2 \mathcal{L}_t(\theta_t^*), \quad b_t = \nabla \mathcal{L}_{t+1}(\theta_t^*), \quad (5)$$

where $b_t = 0$ if $\mathcal{D}_{t+1} = \mathcal{D}_t$, but $b_t \neq 0$ under drift. This shows that the optimizer must continually track a moving minimizer, which destabilizes learned features. Let $z = h_\phi(x) \in \mathbb{R}^d$ with covariance

$$\Sigma_t = \text{Cov}_{x \sim \mathcal{D}_t}(z) = \mathbb{E}[zz^\top] - \mathbb{E}[z]\mathbb{E}[z]^\top, \quad \text{srank}(\Sigma_t) = \frac{\|\Sigma_t\|_F^2}{\|\Sigma_t\|_2^2}. \quad (6)$$

In the stationary case, $\Sigma_t \rightarrow \Sigma$ with a large stable rank, preserving feature diversity. Under non-stationarity, the drift term in equation 5 induces oscillations in Σ_t and systematic *covariance deflation* (drop in srank), a hallmark of collapse. When representations collapse (covariance deflation; equation 6), feature energy shrinks. For a linear head,

$$\mathbb{E}[\|\nabla_W \mathcal{L}_t\|_F^2] \leq 4 \mathbb{E}[\delta_t^2] \text{tr}(\Sigma_t), \quad \delta_t = f_\theta(x) - y_t(x), \quad (7)$$

so smaller $\text{tr}(\Sigma_t)$ directly yields smaller gradients and slower learning. With ReLU features $z = \sigma(a)$, the backprop signal through unit j is gated:

$$\frac{\partial \mathcal{L}_t}{\partial a_j} = \mathbf{1}\{a_j > 0\} \langle \nabla_z \mathcal{L}_t, e_j \rangle \Rightarrow \mathbb{E} \left[\left\| \frac{\partial \mathcal{L}_t}{\partial a_j} \right\|^2 \right] \leq p_{j,t} \mathbb{E}[\|\nabla_z \mathcal{L}_t\|_2^2], \quad (8)$$

where $p_{j,t} = \Pr(a_j > 0)$ and e_j is the j -th basis vector. Non-stationary drift (equation 5) reduces $p_{j,t}$ and $\text{Var}(z_j)$; together with lower $\text{tr}(\Sigma_t)$, this shrinks per-unit updates and increases neuron dormancy (Sokar et al., 2023). \square

C FAST ACTOR–CRITIC

FastTD3 (Seo et al., 2025) extends the standard TD3 framework by combining (i) parallel simulation across many environment instances, (ii) large-batch critic updates, and (iii) algorithm design choices like distributional critics (C51) (Bellemare et al., 2017), noise scaling and clipped double Q-learning (CDQ) (Fujimoto et al., 2018).

Instead of approximating only the expected return, FastTD3 estimates the entire return distribution $Z(s, a)$ following C51 (Bellemare et al., 2017). The action–value distribution is approximated by a categorical distribution supported on N fixed atoms $\{z_i\}_{i=0}^{N-1}$, uniformly spaced in $[v_{\min}, v_{\max}]$:

$$z_i = v_{\min} + i \Delta z, \quad \Delta z = \frac{v_{\max} - v_{\min}}{N - 1}, \quad i = 0, 1, \dots, N - 1.$$

The critic outputs logits $\{\ell_i(s, a)\}_{i=0}^{N-1}$ which define probabilities via $p_i(s, a) = \text{softmax}(\ell(s, a))_i$. Given a transition (s, a, r, s') and a next action $a' = \pi_{\text{targ}}(s')$, the Bellman-updated support is

$$z'_j = \text{clip}(r + \gamma z_j, v_{\min}, v_{\max}), \quad j = 0, 1, \dots, N - 1,$$

with target probabilities $p_j(s', a')$ from the target critic at (s', a') . C51 projects this target distribution back onto the fixed support via the projection operator Φ :

$$m_i \equiv (\Phi TZ)(z_i) = \sum_{j=0}^{N-1} p_j(s', a') \left[1 - \frac{|z'_j - z_i|}{\Delta z} \right]_+, \quad [x]_+ = \max\{x, 0\}. \quad (9)$$

Training minimizes the cross-entropy between the projected target m and the predicted distribution:

$$\mathcal{L}(s, a) = - \sum_{i=0}^{N-1} m_i \log p_i(s, a).$$

This distributional perspective reduces variance in target estimation and captures the multi-modality of returns in continuous control. Two other important stabilizers in actor–critic methods are *noise scaling* and *Clipped Double Q-learning (CDQ)*. Noise scaling injects Gaussian perturbations into the action to balance exploration and stability:

$$a = \pi_{\theta}(s) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2 I), \quad (10)$$

where the scale σ must be tuned to avoid either poor exploration (σ too small) or instability (σ too large). On the other hand, CDQ mitigates overestimation bias by maintaining two critics Q_{ϕ_1}, Q_{ϕ_2} and defining the bootstrapped target conservatively as

$$y = r + \gamma \min_{i=1,2} Q_{\phi_i^-}(s', \pi_{\theta}(s') + \epsilon), \quad \epsilon \sim \mathcal{N}(0, \sigma_{\text{target}}^2 I), \quad (11)$$

where ϕ_i^- are delayed target networks and σ_{target} controls target-smoothing noise. Each critic then minimizes its own squared Bellman error:

$$\mathcal{L}_Q(\phi_i) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[(Q_{\phi_i}(s, a) - y)^2 \right], \quad i = 1, 2. \quad (12)$$

Together, these design choices underpin the high-throughput yet stable behavior of FastTD3 (Seo et al., 2025).

D BENCHMARKS

D.0.1 ISAAC GYM

For our experiments, we used the original Isaac Gym benchmark, which provides pre-built standalone environments and runs entirely on the GPU via a PhysX backend. This setup enables both physics simulation and neural network policy training on the GPU, offering high-throughput evaluation. Although Isaac Gym is deprecated, we used it to ensure reproducibility, specifically running the PPO algorithm from CleanRL on tasks spanning locomotion, robotic hands, and cube stacking (See Figure 11). To reproduce this task, we follow the PPO hyperparameters from CleanRL for Isaac, as presented in Table 3.

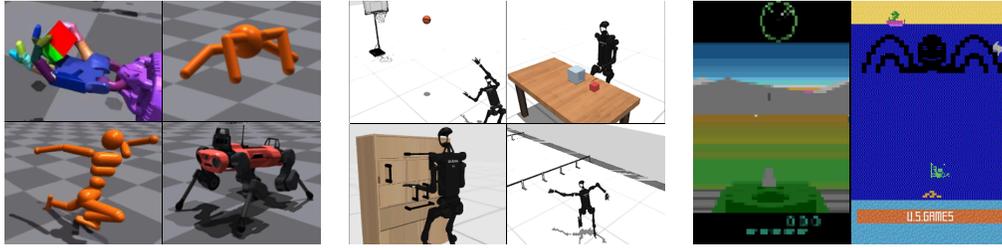


Fig. 11: **Environment Visualizations.** We evaluate SEM across three benchmark suites such as Isaac Gym, HumanoidBench, and Atari. The first two cover state-based locomotion/manipulation; Atari introduces pixel-based games of varying complexity.

D.0.2 HUMANOIDBENCH

In our experiments, we used the Humanoid Benchmark, a suite of tasks for evaluating humanoid robot control across locomotion and manipulation, implemented on the MuJoCo physics engine. We focused on three robot configurations: the Unitree H1 without hands (26 DoF), the Unitree H1 with hands (76 DoF), and the Unitree G1 with three-finger hands (44 DoF). The benchmark defines 27 core tasks, and additionally, sit, balance, and bookshelf are each implemented in both simple and hard variants, while insert is implemented in small and normal configurations. This brings the total to 31 tasks. Our evaluations covered locomotion challenges, including walking, running, crawling, stair climbing, and balancing, and whole-body manipulation tasks such as opening doors, lifting packages, operating kitchen objects, and performing insertions. Together, these tasks provided a diverse and rigorous testing ground for our study of humanoid control (See Figure 11). To reproduce this task, we follow the fastTD3 hyperparameters (Seo et al., 2025), as presented in Table 4.

D.0.3 ATARI

We conducted pixel-based reinforcement learning experiments using the Arcade Learning Environment (ALE) (Bellemare et al., 2013). We start from the Atari-10 suite (Aitchison et al., 2023) and extend it with additional environments from the Atari-20 suite (Fedus et al., 2020), yielding a broader subset of the ALE benchmark. In total, we evaluate 28 games spanning varying difficulty levels, trained with PPO from CleanRL (Huang et al., 2022) as the baseline. Two games overlap across the subsets (Bowling and Q*bert). We report IQM returns (see Figure 10, Left) following the evaluation protocol of (Agarwal et al., 2021; Ceron et al., 2024b;a; Sokar & Castro, 2025).

D.0.4 MTBENCH

In our experiments, we used the Multi-Task Benchmark for Robotics, an open-source suite built on the GPU-accelerated Isaac Gym simulator. Specifically, we worked with the 50 manipulation tasks adapted from Meta-World, where a single-armed robot interacts with one or two objects through actions such as pushing, picking, and placing. Each task provides parametric variations in object initialization and target positions, adding diversity and complexity. For evaluation, we adopted the MT50 setting, which encompasses the full set of 50 tasks.

D.0.5 BOOSTER GYM HUMANOID ROBOT

As reported in the FastTD3 paper (Seo et al., 2025), the authors configured the Booster Gym humanoid robot to operate within the MuJoCo Playground environment, which supports 12 degrees of freedom (DOF). They trained the policy entirely in simulation and successfully transferred it to a real robot, demonstrating an effective sim-to-real deployment of an off-policy RL method on a full-sized humanoid. In our experiments, we used the same robot model and hyperparameters as FastTD3, also within MuJoCo Playground. Our configuration led to faster convergence and improved performance in simulation compared to the reported baseline.

Table 1: Default hyper-parameters setting for actor-critic MLP

Hyper-parameter	Value
Critic Hidden Dim	1024
Actor Hidden Dim	512
Critic Learning Rate	3e-4
Actor Learning Rate	3e-4

E DEEP RL NETWORK ARCHITECTURES

E.0.1 MLP

We modified the FastTD3 architecture, specifically in the actor-critic design, where both networks are implemented as multilayer perceptrons (MLPs). The critic receives concatenated observation-action inputs, while the actor processes only the observations. In both cases, the inputs first pass through two linear layers with ReLU activations. At this point, we introduced the SEM mechanism, which can be enabled or disabled, and applied selectively to the actor, the critic, or both. For the critic, if SEM is not used, the representation is processed by a sequence of `Linear`→`ReLU`→`Linear` layers, with the final linear layer outputting dimension `num_atoms`. If SEM is enabled, the sequence becomes `SEM`→`Linear`, again producing an output of size `num_atoms`. For the actor, the representation without SEM follows a `Linear`→`ReLU`→`Linear`→`Tanh` sequence, while with SEM it follows `SEM`→`Linear`→`Tanh`, where the final `Tanh` ensures bounded continuous actions. In Table 1, we present the fixed hyperparameters used across all environments. Other hyperparameters, such as `num_atoms` or `num_env`, varied depending on the environment, in which case we adopted the values proposed by (Seo et al., 2025).

E.0.2 CNN

For our pixel-based experiments, we modified the PPO implementation from CleanRL, which follows an actor-critic design. The shared backbone consists of three convolutional layers, each followed by a ReLU activation, producing a flattened representation that is then processed by a two-layer MLP with ReLU activations. This representation is used by both the actor and the critic. In our intervention, we introduced the SEM block into the actor: when enabled, the representation passes through the SEM block before a final linear layer; when disabled, it follows a `Linear`→`ReLU`→`Linear` sequence. The critic remains unchanged, while the actor architecture is varied depending on the use of SEM. We adopted the PPO hyperparameters for Atari from CleanRL (Huang et al., 2022), as summarized in Table 5.

F ABLATION SETUP

Given our constrained computational budget, we performed experiments on a subset of HumanoidBench, consisting of five robotics tasks. These tasks are part of the benchmark evaluated with FastTD3 (Seo et al., 2025) and correspond to `h1hand`-{`walk`, `stand`, `run`, `stair`, `slide`}-`v0`. All ablation experiments were conducted on this subset using six random seeds.

G METRICS

To better understand the dynamics of training and the quality of learned representations, we report a diverse set of metrics beyond standard returns. These measures capture complementary aspects of learning, including representation diversity, network expressivity, parameter stability, gradient behavior and sampling efficiency. Results of these analyses are provided in section 4.

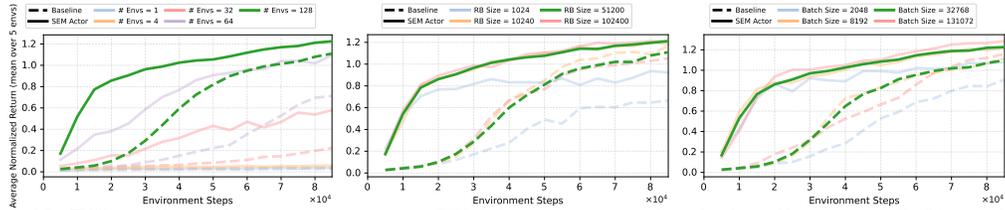


Fig. 12: **Effect of core hyperparameters.** SEM Actor compared to the baseline across (left) number of parallel environments, (middle) replay buffer size, and (right) batch size. SEM consistently scales better and achieves higher returns.

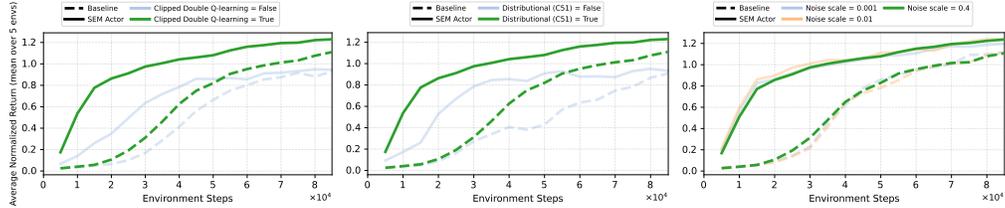


Fig. 13: **Robustness to design choices.** SEM Actor vs. baseline across (left) clipped double Q-learning, (middle) distributional critic (C51), and (right) exploration noise scale. SEM remains robust, while the baseline is more sensitive.

G.1 FEATURE RANK

This metric assesses the quality of learned representations in deep RL by identifying the smallest subspace that retains 99% of the variance, thereby enhancing interpretability, efficiency, and stability. A higher feature rank indicates more diverse representations. The computation follows the approximate rank from (Yang et al., 2019; Moalla et al., 2024):

$$\sum_{i=1}^k \frac{\sigma_i^2}{\sum_{j=1}^n \sigma_j^2} \geq \tau,$$

where σ_i are the singular values of the feature matrix, n is the total number of singular values, and τ is the variance threshold (e.g., 99%). The feature rank k is the smallest number of principal components required to preserve at least τ of the total variance.

G.2 DORMANT NEURONS

This metric quantifies the proportion of neurons with near-zero activations, which limits the network’s expressivity. It serves to detect inefficiencies in learning, as a high proportion of dormant neurons implies that many units are inactive or rarely contribute to the output. The computation follows (Sokar et al., 2023):

$$\frac{\sum_{i=1}^N \mathbf{1}(|a_i| < \epsilon)}{N} \times 100,$$

where N is the total number of neurons, a_i is the activation of neuron i , ϵ is a small threshold (e.g., 10^{-5}), and $\mathbf{1}$ is the indicator function.

G.3 WEIGHT NORM

This metric measures the magnitude of neural network weights, providing insight into model complexity, stability, and overfitting risk. Large weight norms indicate parameters with high magnitudes, which may hinder generalization. The metric is computed as in (Moalla et al., 2024; Lyle et al., 2021):

$$\|\theta\|_2 = \sqrt{\sum_i \theta_i^2},$$

where θ_i are the weights of a given layer.

1404 G.4 GINI SPARSITY

1405
1406 The Gini metric is used to quantify the sparsity of neural representations. A high Gini value indicates
1407 a sparse representation, where only a few neurons are strongly active while most remain near zero;
1408 this often improves interpretability, makes more efficient use of network capacity, and can help
1409 reduce overfitting. In contrast, a low Gini value corresponds to dense representations, where many
1410 neurons are active simultaneously, allowing the network to capture richer information but often
1411 at the cost of reduced interpretability and potentially noisier features. In practice, we observed a
1412 direct relationship between the Gini metric and the return when using SEM, with better performance
1413 associated with higher Gini values. The Gini value is computed using the following equation.

$$1414$$

$$1415 G = 1 + \frac{1}{n} - \frac{2}{n \sum_{i=1}^n v_i} \sum_{i=1}^n (n+1-i) v_{(i)}$$

$$1416$$

$$1417$$

1418 where where

$$1419 v = (|x_1|, |x_2|, \dots, |x_n|)$$

$$1420$$

1421 It is the vector of all activations, taken in absolute value and stacked into one vector. The Gini metric
1422 has been explored in the papers (Hurley & Rickard, 2009; Zonoobi et al., 2011).

1424 G.5 CRAMER DISTANCE

1425
1426 The Cramér distance is defined as the squared L_2 distance between the cumulative distribution
1427 functions (CDFs) of two probability distributions. When the distributions are similar, their CDFs
1428 overlap closely and the Cramér distance approaches zero. Conversely, when the distributions differ,
1429 the CDFs diverge and the distance increases. In practice, a lower Cramér distance indicates that
1430 the learned distribution is closer to the target distribution, which is desirable. Empirical results
1431 also suggest a correlation between lower Cramér distance and improved returns. This measure is
1432 computed using the following equation:

$$1433$$

$$1434 D_{\text{Cramér}}^2(p_1, p_2) = \sum_{j=1}^n \left(F_{p_1}(z_j) - F_{p_2}(z_j) \right)^2 \Delta z$$

$$1435$$

$$1436$$

1437 where p_1 and p_2 are probability distributions, and F_{p_1}, F_{p_2} denote their corresponding cumulative
1438 distribution functions (CDFs).

1441 G.6 ENTROPY

1442
1443 This metric measures the average entropy of the representations. High entropy indicates that the
1444 representation is more dispersed, less concentrated, and carries more uncertainty. Low entropy
1445 corresponds to a more concrete representation, with higher sparsity. In practice, we observe a re-
1446 lationship where lower entropy is associated with better returns and a higher Gini measure. This
1447 metric is defined by the following equation.

$$1448$$

$$1449 p_{i,j} = \frac{p_{i,j}}{\sum_k p_{i,k} + \varepsilon}$$

$$1450$$

$$1451$$

$$1452 \text{entropy} = \frac{1}{B} \sum_{i=1}^B \left(- \sum_j p_{i,j} \log(p_{i,j} + \varepsilon) \right)$$

$$1453$$

$$1454$$

$$1455$$

1456 where B is the batch size, and p is the non-negative representation normalized to form a probability
1457 distribution.

H ADDITIONAL BASELINES

We expand our comparison to include a broader set of interventions commonly used to address plasticity loss, representation collapse, and optimization pathologies in deep RL. Following the setup described in section 4 (*Comparing SEM to other Regularization Methods*), we evaluate eight additional methods: Reset (Nikishin et al., 2022), L2 regularization (Kumar et al., 2022a), ReGen (Kumar et al., 2023), Dropout (Hendrycks, 2016), Shrink-and-Perturb (Ash & Adams, 2020), GELU activation (Hinton et al., 2012), Weight Clipping (Elsayed et al., 2024), and Spectral Normalization (Gogianu et al., 2021). Each method is evaluated across 6 random seeds and 5 HumanoidBench tasks.

These experiments are *in addition* to those reported in Fig. 7, which already include other relevant baselines such as Gumbel with straight-through estimation (Jang et al., 2017; Maddison et al., 2017), Vector Quantization (van den Oord et al., 2018), and C-ReLU (Abbas et al., 2023). We demonstrate that SEM offers a lightweight, stable, and effective mechanism for mitigating representation collapse in deep RL without requiring complex architectural modifications or extensive hyperparameter tuning.

H.1 REGEN

We applied ReGen by zeroing tiny weights using three thresholds, $\tau = 10^{-5}$, 10^{-6} and 10^{-7} , allowing us to compare how different magnitudes influence sparsity and model stability (Kumar et al., 2023).

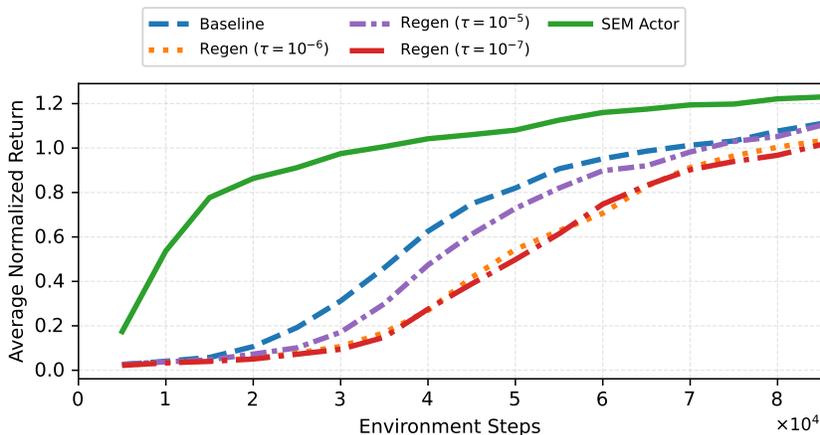


Fig. 14: **Average normalized return on 5 HumanoidBench tasks over 6 seeds.** We compare the Baseline and ReGen variants with SEM Actor. SEM variation consistently outperforms all ReGen thresholds, achieving faster learning and higher normalized returns across training steps.

H.2 SHRINK AND PERTURB

We ran Shrink and Perturb with shrink factors 0.99, 0.9, 0.8, and 0.5 to compare its performance and training behavior directly against our proposed method (Ash & Adams, 2020).

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

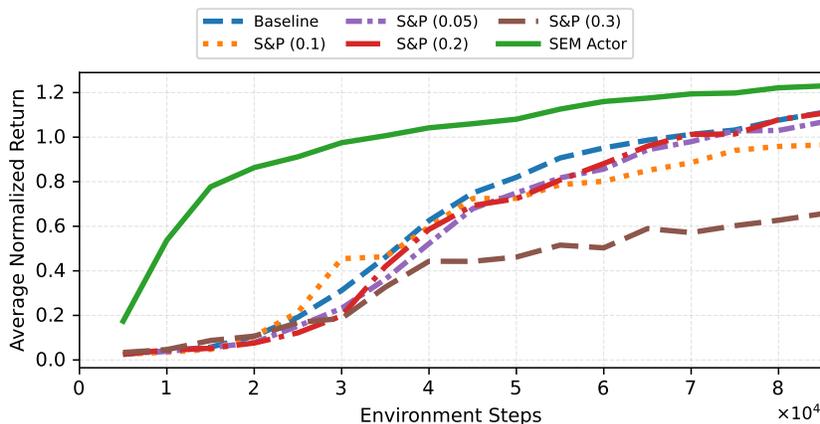


Fig. 15: Average normalized return on 5 HumanoidBench tasks over 6 seeds. We compare the Baseline and Shrink-and-Perturb variants with SEM Actor. SEM variation maintains a clear advantage across all S&P levels, achieving faster learning and higher normalized returns.

H.3 WEIGHT CLIP

We ran Weight Clipping with ranges $(-0.01, 0.01)$, $(-0.001, 0.001)$, and $(-0.1, 0.1)$ to assess how different clipping strengths perform and to compare them directly against our proposed method (Elsayed et al., 2024).

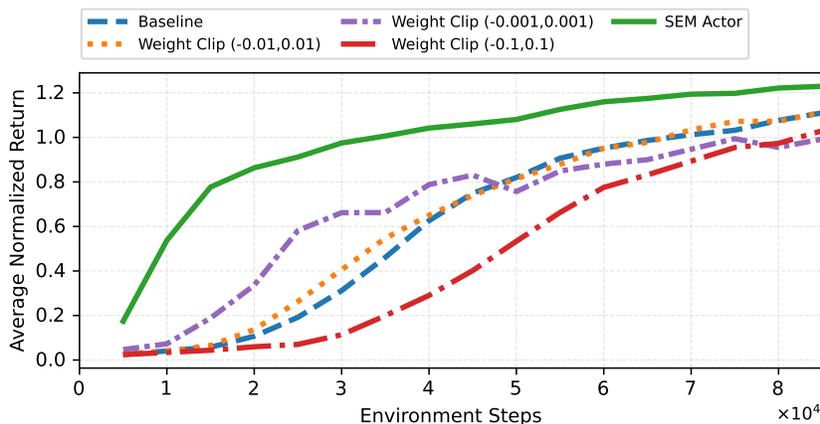


Fig. 16: Average normalized return on 5 HumanoidBench tasks over 6 seeds. We compare the Baseline and Weight Clipping variants with SEM Actor. SEM variation achieves faster improvement and higher normalized returns across environment steps.

H.4 RESET LAYER

We applied the reset layer baseline by reinitializing weights with Xavier initialization, allowing us to evaluate how fully resetting a layer compares to the performance and stability of our proposed method (Nikishin et al., 2022).

1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

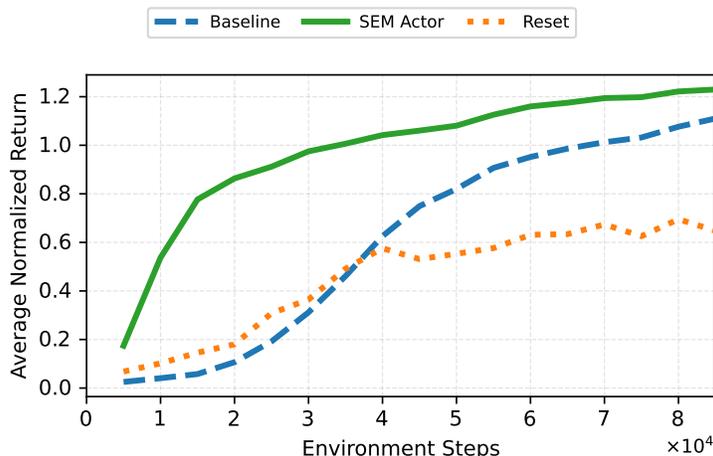


Fig. 17: Average normalized return on 5 HumanoidBench tasks over 6 seeds. We compare the Baseline and Reset Layer with SEM Actor. SEM variation accelerates learning and improves normalized return across environment steps, outperforming both the baseline and the reset-layer strategy.

H.5 SPECTRAL NORM

We applied Spectral Normalization to constrain layer norms and stabilize training, allowing us to directly compare its behavior and performance against the results achieved by our proposed method (Gogianu et al., 2021).

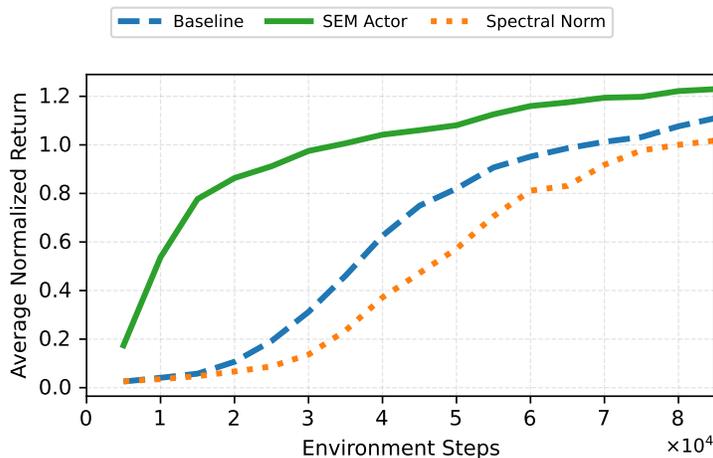


Fig. 18: Average normalized return on 5 HumanoidBench tasks over 6 seeds. We compare the Baseline and Spectral Norm with SEM Actor. SEM variation outperforms both the baseline and Spectral Norm across training, achieving faster learning and higher normalized returns.

H.6 FEATURE NORM

We applied Feature Norm by normalizing activations to unit L2 norm, allowing us to evaluate its effect on representation stability and directly compare its behavior with the performance of our proposed method (Kumar et al., 2022a).

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

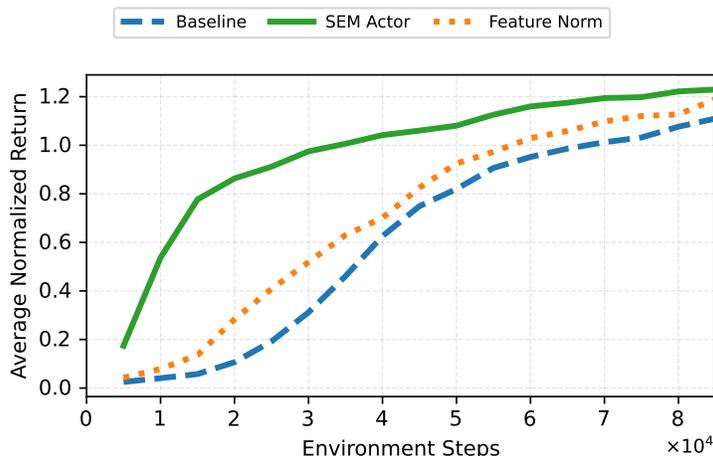


Fig. 19: Average normalized return on 5 HumanoidBench tasks over 6 seeds. We compare the Baseline and Feature Norm with SEM Actor. SEM variation outperforms both methods, showing faster learning and higher normalized returns throughout training.

H.7 GELU

We evaluated the GELU activation to assess its impact on learning dynamics and representation quality, enabling a direct comparison between this widely used baseline and the performance achieved by our proposed method (Hinton et al., 2012).

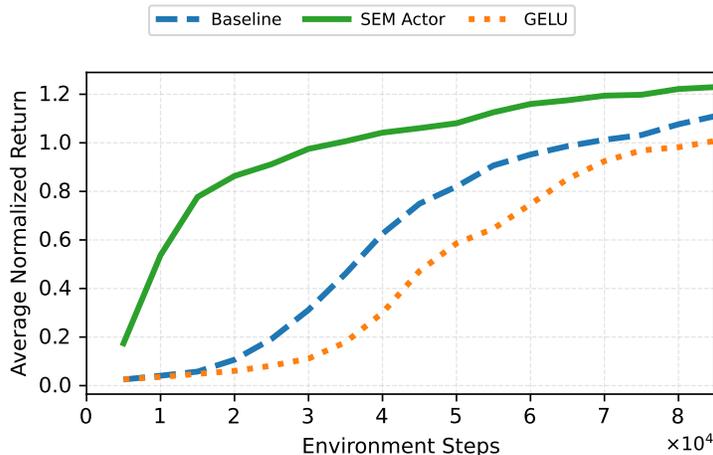


Fig. 20: Average normalized return on 5 HumanoidBench tasks over 6 seeds. We compare the Baseline and GELU activation with SEM Actor. SEM variation consistently surpasses both approaches, achieving faster learning and higher normalized returns across training.

H.8 DROPOUT

We evaluated Dropout using probabilities 0.1, 0.05, and 0.2 to assess its regularization behavior and directly compare its performance and stability against the results achieved by our proposed method (Hendrycks, 2016).

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

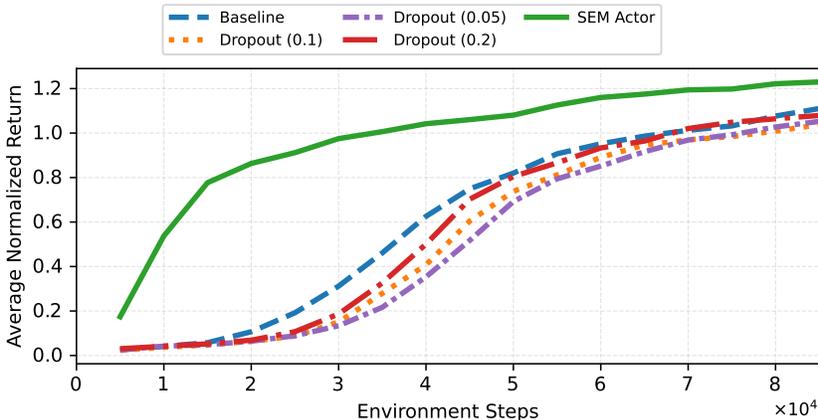


Fig. 21: Average normalized return on 5 HumanoidBench tasks over 6 seeds. We compare the Baseline and Dropout variants with SEM Actor. SEM variation consistently shows faster learning progress and achieves higher normalized returns across training.

H.9 MoEs

We evaluate Mixture-of-Experts (MoE) architectures with 1 and 4 experts to study how expert multiplicity influences model capacity, learning dynamics, and overall stability (Ceron et al., 2024b; Willi et al., 2024). These comparisons allow us to assess whether increasing routing flexibility or expert specialization can match the gains provided by SEM. As shown in Fig. 22, both SoftMoE and Top1-MoE track the baseline closely during training, with modest improvements in later stages for the 4-expert configuration. However, none of the MoE variants approach the sample efficiency or final performance achieved by SEM Actor. SEM consistently learns faster and attains higher normalized returns.

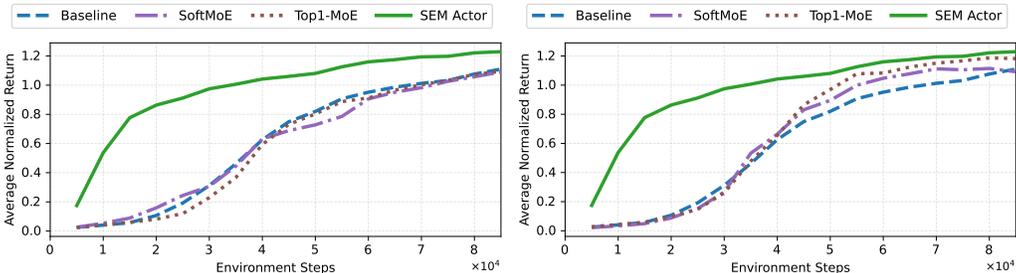


Fig. 22: Average normalized return on 5 HumanoidBench tasks over 6 seeds. We compare the Baseline, SoftMoE, Top1-MoE, and SEM Actor. The left plot shows MoE models with 1 expert, and the right plot shows models with 4 experts. MoE variants offer limited gains over the baseline and show slower learning dynamics, while SEM consistently achieves faster progress and higher final returns throughout training.

H.10 ReDO

We evaluated ReDo (Sokar et al., 2023) by varying its dormancy thresholds (0.1, 0.01, 0.001) to study the sensitivity of neuron-reset frequency and to compare its stability and performance to SEM. These thresholds control how aggressively inactive neurons are reset, allowing us to examine whether ReDo’s representational refreshing mechanism can match the benefits provided by SEM. As shown in Fig. 23, all ReDo variants track the baseline closely throughout training, with limited improvements in early or final performance. In contrast, SEM Actor consistently learns faster and achieves higher normalized returns.

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

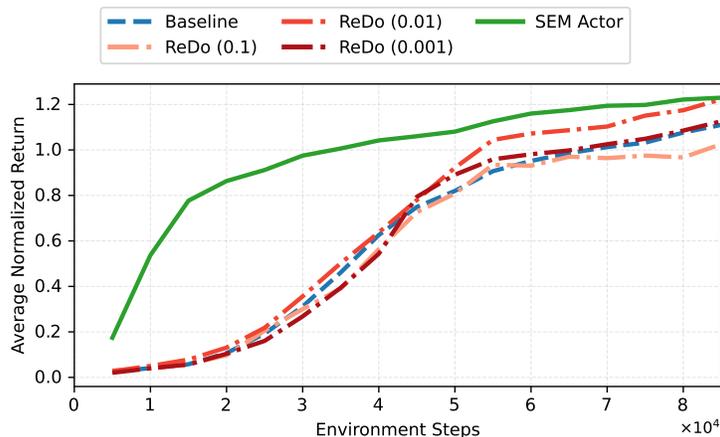


Fig. 23: **Average normalized return on 5 HumanoidBench tasks over 6 seeds.** We compare the Baseline, three ReDo variants (thresholds 0.1, 0.01, 0.001), and SEM Actor. All ReDo configurations remain close to the baseline, showing limited gains across training. SEM achieves faster learning and higher final returns, highlighting the effectiveness of structured embedding modifications over neuron-reset-based approaches.

H.11 MAGNITUDE PRUNING

We evaluate Magnitude Pruning under two target sparsity levels (50% and 95%) to analyze how pruning severity affects performance and to provide a direct comparison with our proposed method (Graesser et al., 2022; Ceron et al., 2024a). Fig. 24 reports results on 5 HumanoidBench tasks. Overall, both pruning variants lag behind SEM Actor across the entire training horizon. Moderate pruning (50%) retains reasonable learning progress but consistently underperforms SEM in both sample efficiency and final returns. Severe pruning (95%) leads to an early collapse in training, indicating that aggressive weight removal significantly harms representational capacity in this setting. In contrast, SEM maintains stable optimization and achieves substantially higher returns, highlighting the robustness of structured embedding modifications compared to unstructured pruning.

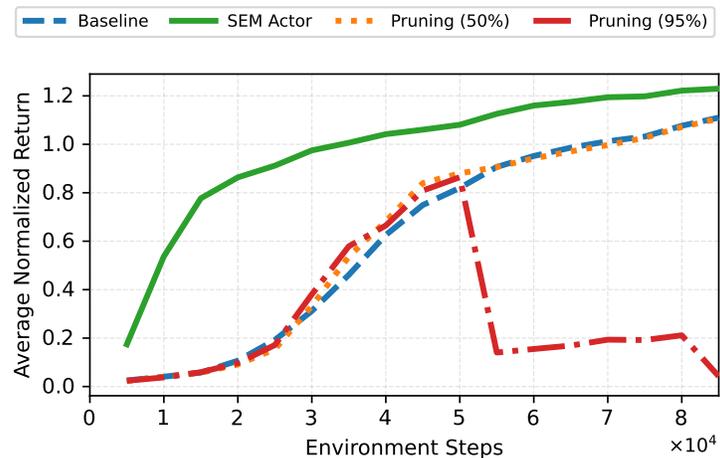


Fig. 24: **Average normalized return on 5 HumanoidBench tasks over 6 seeds.** We compare the Baseline, SEM Actor, and Magnitude Pruning at 50% and 95% target sparsity. SEM consistently learns faster and reaches higher final returns than both pruning variants. Moderate pruning (50%) provides limited gains but remains inferior to SEM, while high-sparsity pruning (95%) collapses training early, demonstrating the instability of aggressive unstructured pruning.

I HYPERPARAMETERS

In this section, we list the hyperparameters used across our experimental settings.

Table 2: Wall-clock training time (hh:mm) for the **actor** on the H1-hand humanoid benchmark under default settings. We compare FastTD3 and FastTD3+SEM; lower is better.

Game	Actor	
	<i>FastTD3</i>	<i>FastTD3+SEM</i>
h1hand-walk	2:31 h	2:42 h
h1hand-stand	2:29 h	2:20 h
h1hand-run	2:46 h	2:34 h
h1hand-stair	4:09 h	4:13 h
h1hand-slide	5:35 h	5:24 h

Table 3: Default hyperparameter settings for the PPO agent on Isaac Gym.

Hyper-parameter	Value
Adam’s (ϵ)	1e-5
Adam’s learning rate	2.6e-3
Dense Activation Function	Tanh
Dense Width	256
Discount Factor	0.99
Number of Dense Layers	3
Number of environments	4096

Table 4: Default hyperparameter settings for the fastTD3 agent on the humanoid bench.

Hyper-parameter	Value
Critic Hidden Dim	1024
Actor Hidden Dim	512
Critic Learning Rate	3e-4
Actor Learning Rate	3e-4
Discount Factor	0.99
Dense Activation Function	ReLU
Number of Dense Layers	4
Number of environments	128
Number of atoms	101

Table 5: Default hyperparameter settings for the PPO agent on Atari.

Hyper-parameter	Value
Adam’s (ϵ)	1e-5
Adam’s learning rate	2.5e-4
Conv. Activation Function	ReLU
Convolutional Width	32,64,64
Dense Activation Function	ReLU
Dense Width	512
Normalization	None
Discount Factor	0.99
Number of Convolutional Layers	3
Number of Dense Layers	2
Reward Clipping	True
Weight Decay	0

J LEARNING CURVES FOR EACH GAME

To complement the aggregate results reported in the main text (see [section 4](#) and [section 5](#)), we provide full learning curves for each environment in the benchmark. These plots illustrate training dynamics across seeds and highlight differences in sample efficiency and stability between +SEM and its corresponding baseline (FastTD3/FastTD3-SimbaV2/FastSAC). The set of robotics tasks follows those used in the FastTD3 benchmark ([Seo et al., 2025](#)).

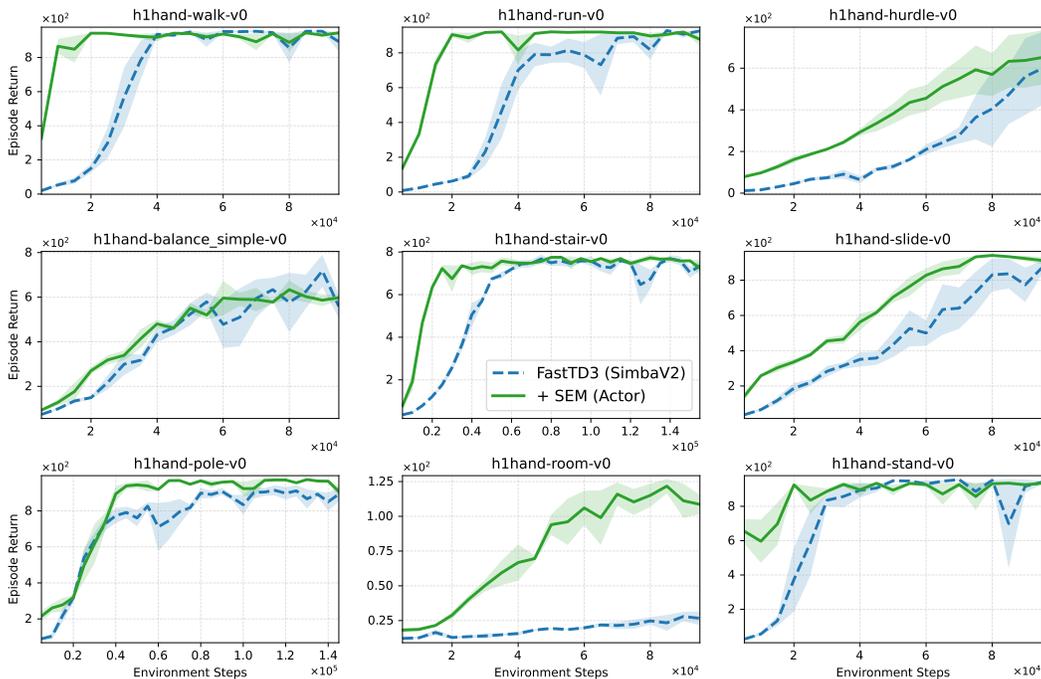


Fig. 25: **Learning curves on 9 h1hand tasks. FastTD3+SimbaV2 (blue, - -) vs. + SEM (Actor) (green, —).** Curves show the mean episode return across 6 seeds. Axes are independently scaled per subplot for readability. SEM (Actor) consistently accelerates learning and achieves higher or comparable final returns on most tasks.

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

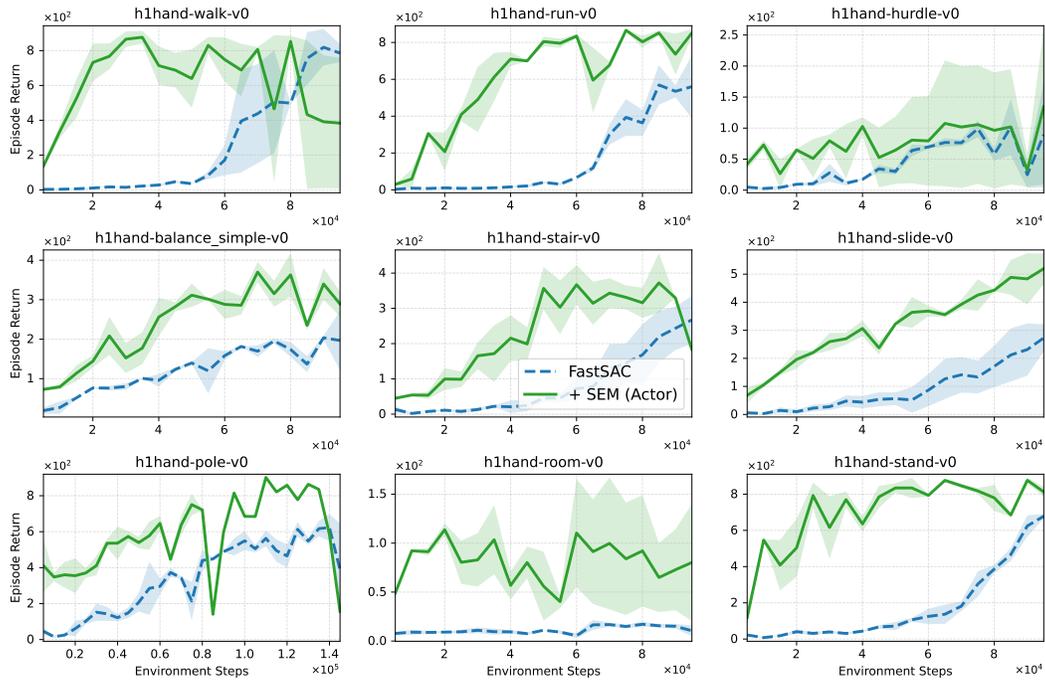


Fig. 26: Learning curves on 9 h1hand tasks. FastSAC (blue, - -) vs. + SEM (Actor) (green, —). Curves show the mean episode return across 6 seeds. Axes are independently scaled per subplot for readability. SEM (Actor) generally accelerates learning and achieves higher final returns on most tasks.

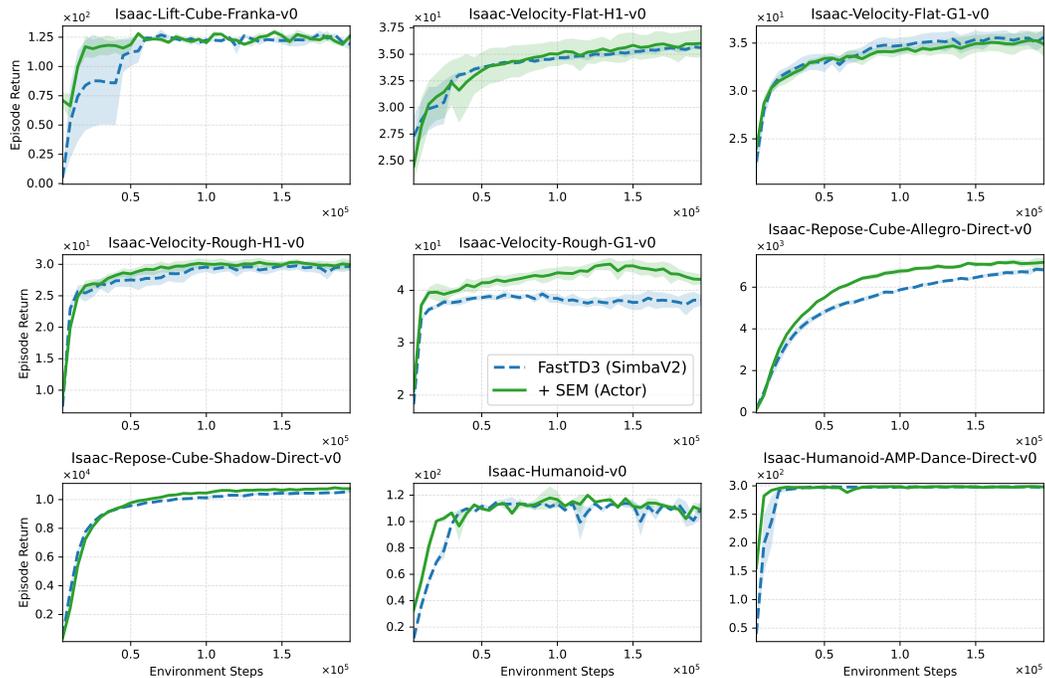


Fig. 27: Learning curves on 9 IsaacGym tasks. FastSAC (blue, - -) vs. + SEM (Actor) (green, —). Curves show the mean episode return across 6 seeds. Axes are independently scaled per subplot for readability. SEM (Actor) generally accelerates learning.

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

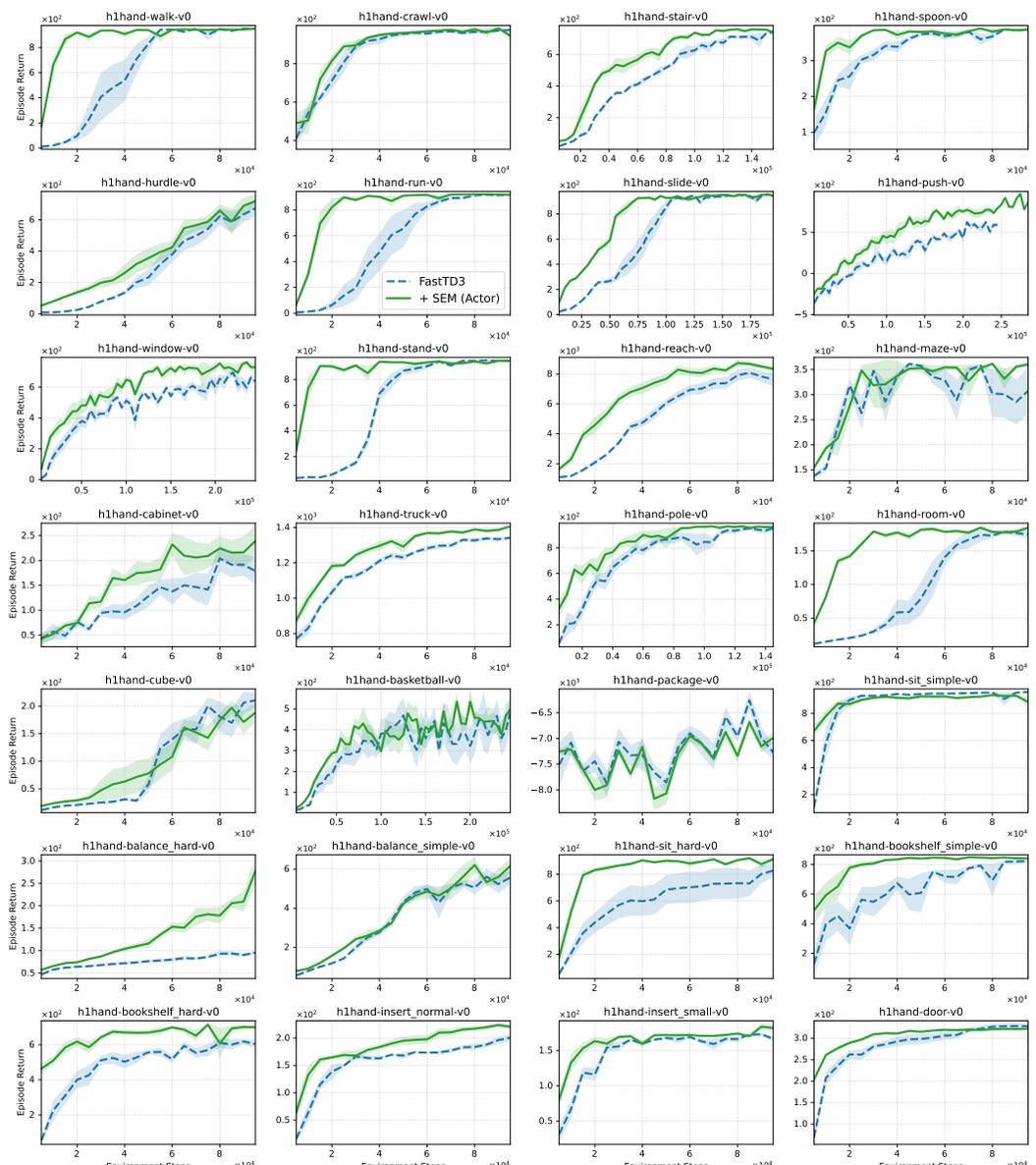


Fig. 28: Learning curves on 28 h1hand tasks (Sferrazza et al., 2024). FastTD3 (blue, - -) vs. + SEM (Actor) (green, —). Curves show the mean episode return across 6 seeds. Axes are independently scaled per subplot for readability. SEM (Actor) typically achieves faster learning and equal or higher final return on most tasks.

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

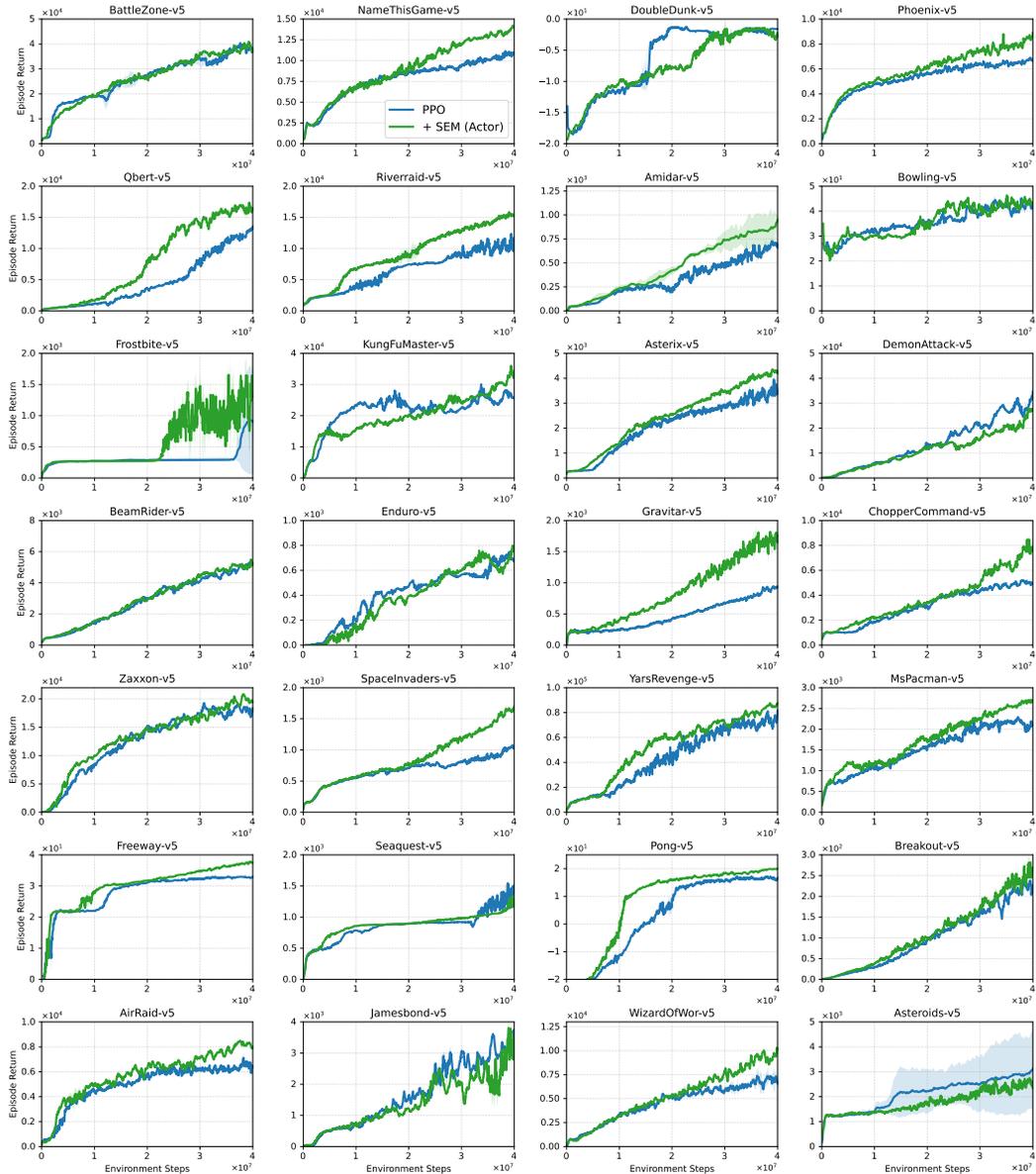


Fig. 29: Learning curves on Atari game (Aitchison et al., 2023). PPO (blue, - -) vs. + SEM (Actor) (green, —). Curves show the mean episode return across 3 seeds. SEM (Actor) typically achieves faster learning and equal or higher final return on most tasks.

2043
 2044
 2045
 2046
 2047
 2048
 2049
 2050
 2051

K ADDITIONAL EXPERIMENTS ON HUMANOIDBENCH

We evaluate +SEM beyond the tasks proposed in FastTD3 by considering additional environments from the Humanoid benchmark (Sferrazza et al., 2024). These experiments assess the scalability of SEM across different robot morphologies and task sets. We include environments featuring the H1 robot without hands and the Unitree G1 with three-finger hands.

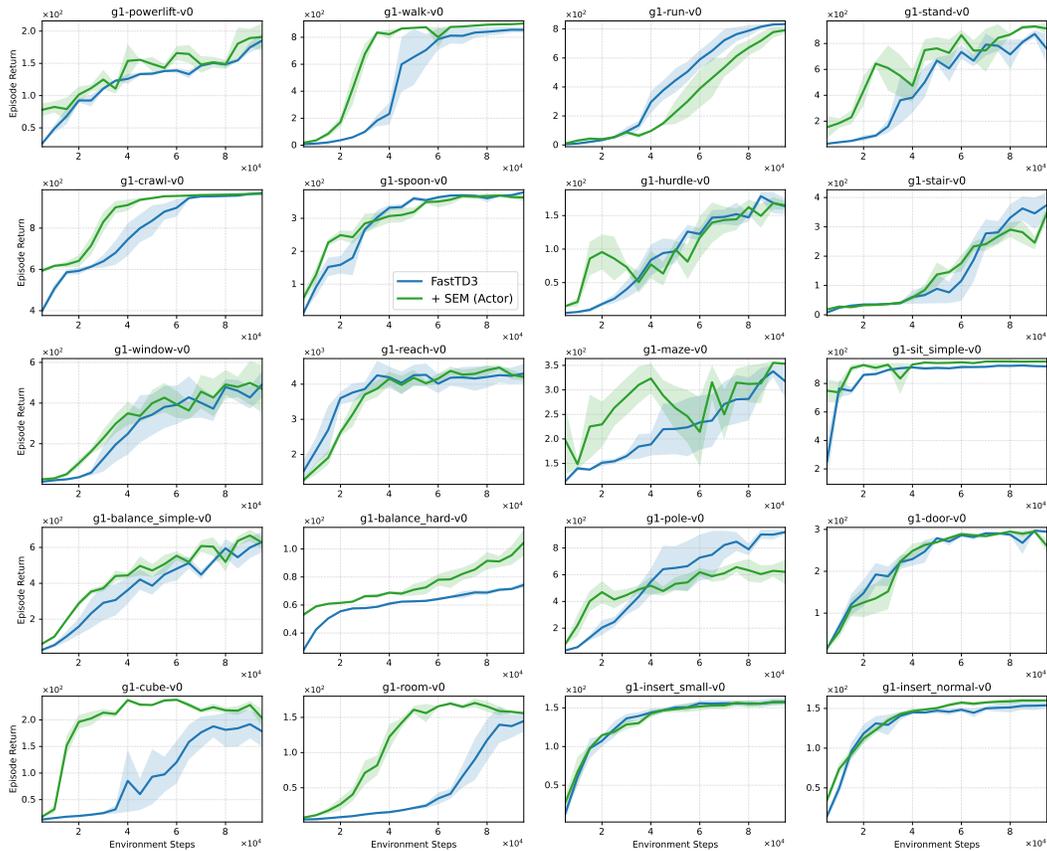


Fig. 30: Learning curves on 16 h1 tasks (Sferrazza et al., 2024). FastTD3 (blue, - -) vs. + SEM (Actor) (green, —). Curves show the mean episode return across 6 seeds. Axes are independently scaled per subplot for readability. SEM (Actor) typically achieves faster learning and equal or higher final return on most tasks.

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

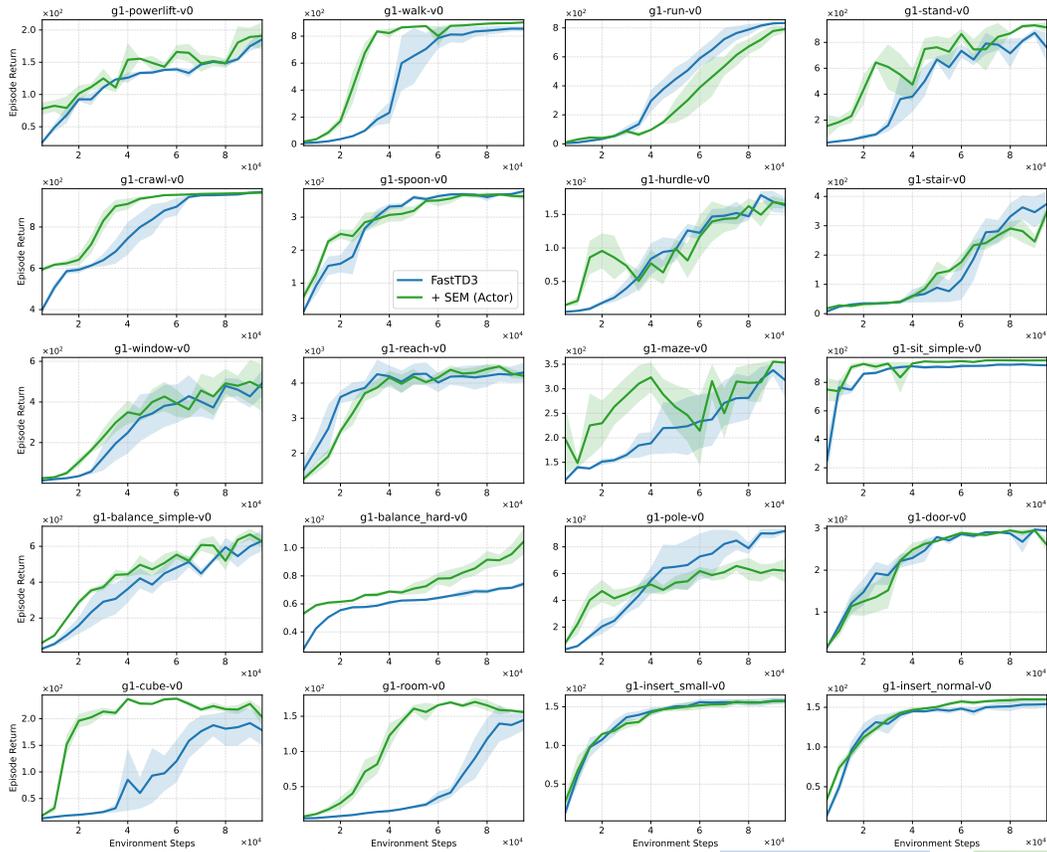


Fig. 31: Learning curves on 20 g1 tasks (Sferrazza et al., 2024). FastTD3 (blue, --) vs. +SEM (Actor) (green, —). Curves show the mean episode return across 6 seeds. Axes are independently scaled per subplot for readability. SEM (Actor) typically achieves faster learning and equal or higher final return on most tasks.

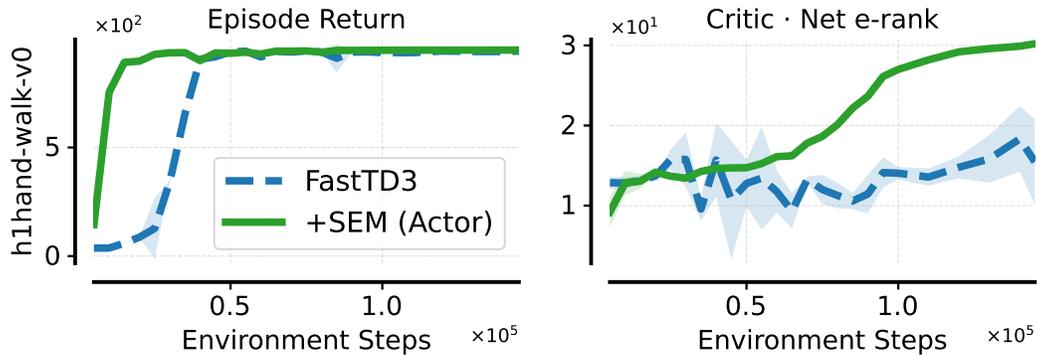


Fig. 32: Learning and representation diagnostic on h1hand-walk-v0 task. SEM reaches high return earlier and critic effective rank.

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

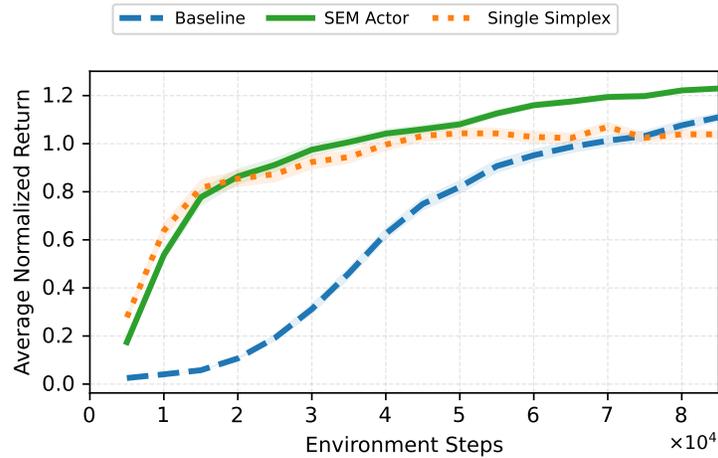


Fig. 33: **Average normalized return on 5 HumanoidBench tasks over 6 seeds.** We compare the baseline agent (blue, --) with SEM variants applied to the actor. SEM variants accelerate early learning, though a single-simplex configuration tends to plateau.

L EVALUATING LAYER-WISE INTERVENTIONS

Selecting which layer to modify is non-trivial, as the search space grows with model depth and architectural complexity. Throughout the paper, we apply SEM to the penultimate layer, following prior work (Gogianu et al., 2021; Kumar et al., 2022b; Ceron et al., 2024b; Sokar & Castro, 2025). These studies highlight that the penultimate layer plays a key role in shaping and constraining learned representations in deep RL.

Here, we extend this analysis by evaluating the effect of applying SEM to individual actor layers as well as to all layers simultaneously. Following the experimental setup in section 4, we report results in Fig. 34 and Fig. 35. Across both settings, we observe that applying SEM to deeper layers or to the entire network leads to faster learning and higher final returns compared to intervening on early layers. Notably, applying SEM solely to Layer-3 achieves performance comparable to the full-layer intervention, indicating that modifying a single well-chosen layer is sufficient to capture most of SEM’s benefits.

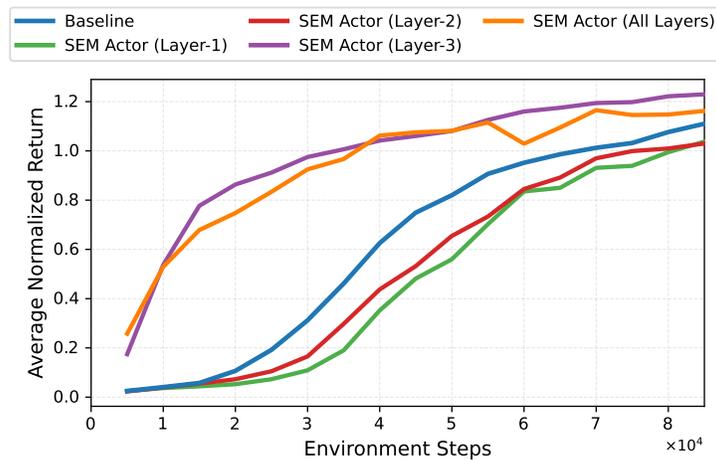


Fig. 34: **Average normalized return on 5 HumanoidBench tasks over 6 seeds.** We evaluate the effect of applying SEM to individual actor layers (Layer-1, Layer-2, Layer-3) and to all layers with $dim = 64$.

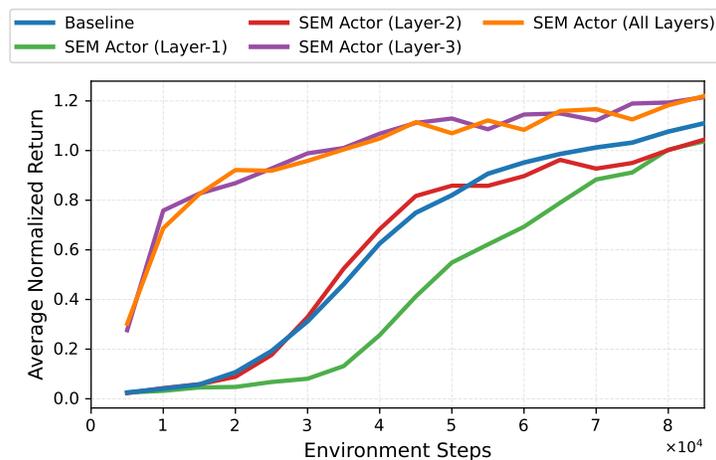


Fig. 35: **Average normalized return on 5 HumanoidBench tasks over 6 seeds.** We evaluate the effect of applying SEM to individual actor layers (Layer-1, Layer-2, Layer-3) and to all layers with $dim = 128$.

M SEM ON VALUE-BASED DEEP RL

We additionally evaluate SEM in value-based deep RL settings. Specifically, we run PQN (Gallici et al., 2025) on 28 Atari games following the experimental setup from section 5 using 3 seeds and training for $40M$ environment steps. As shown in Fig. 37, SEM does not yield consistent improvements over the baseline when averaging performance across all games. However, examining per-game learning curves (see Fig. 36) reveals that SEM provides meaningful gains in both sample efficiency and final performance on several titles (e.g., Asterix, BeamRider). These results suggest that the effectiveness of SEM in value-based methods may depend on game-specific dynamics and representation structure, raising interesting research questions that we leave for future investigation.

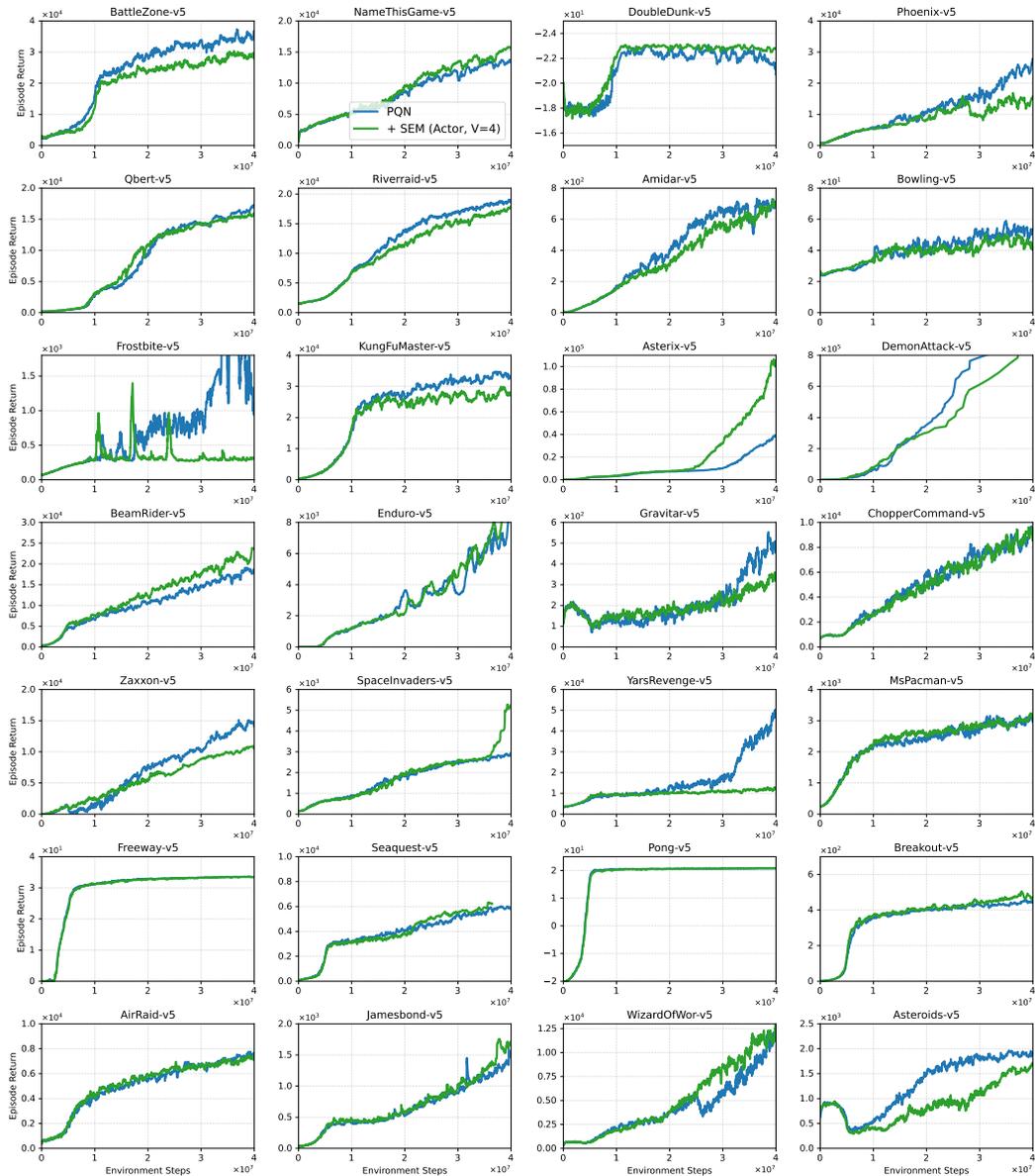


Fig. 36: Learning curves on 28 Atari games (Aitchison et al., 2023). PQN (blue, - -) (Lavoie et al., 2023) vs. + SEM (Actor) (green, —). Curves show the mean episode return across 3 seeds.

2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375

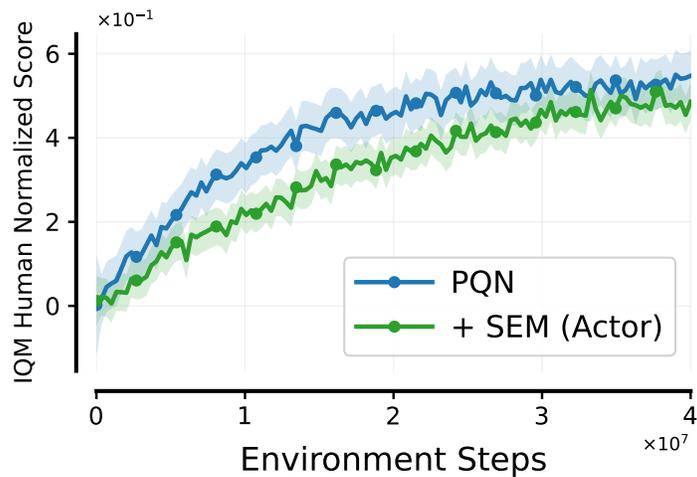


Fig. 37: IQM scores computed over 40M environment steps over 18 games, with 3 independent runs each, and error bars showing 95% stratified bootstrap confidence intervals. PQN (blue, - -) (Lavoie et al., 2023) vs. + SEM (Actor) (green, —).