

---

# COMP4434 Group Project - Group 5

---

WANG Yuqi  
PolyU, COMP  
[REDACTED]  
@connect.polyu.hk

YANG Xikun  
PolyU, COMP  
[REDACTED]  
@connect.polyu.hk

CUI Mingyue  
PolyU, COMP  
[REDACTED]  
@connect.polyu.hk

## Abstract

Sample efficiency remains a critical challenge in reinforcement learning (RL). While state-of-the-art world models like DreamerV3 achieve remarkable efficiency, their complexity motivates the exploration of simpler alternatives. This project investigates the trade-offs between model complexity and performance in efficient RL. We implement and evaluate four distinct approaches: (1) the DreamerV3 world model, serving as a performance upperbound; (2) classical RL baselines employing Q-Learning with function approximation; (3) the widely-used model-free deep RL algorithm, Proximal Policy Optimization (PPO); and (4) a simplified Dreamer-like method utilizing contrastive learning instead of reconstruction. Our comparative analysis reveals that although the proposed contrastive approach does not fully match the performance capabilities of DreamerV3, it significantly outperforms both the classical baselines and deep RL baseline in sample-scarce environments. These findings suggest promising directions for developing less complex yet sample-efficient model-based RL agents.

## 1 Introduction

Reinforcement learning has enabled intelligent agents to make decisions in novel situations by learning a representation of the environment through interactions. Traditional model-free reinforcement learning methods learn this representation through large amounts of trial-and-error. While this approach has enabled some of the most impressive achievements in realms such as the game of Go, Dota, and Starcraft, as well as the recent success of reasoning large language models, its brute-force nature makes it extremely sample inefficient [1], [2], [3]. A typical RL agent often requires millions of samples to learn simple tasks, limiting its application to high-dimensional and complex environments.

In recent years, world models have emerged as a promising approach to overcome these issues. It reduces the amount of samples required by explicitly learning the environment's dynamics and representation. The learned model can then be used to generate large quantities of imagined data, or be used as a simulator for planning. In the paper we chose, DreamerV3, world model has been shown to outperform specialized expert algorithms across diverse domains with fixed parameters [4].

While DreamerV3 excels in both performance and efficiency, it is computationally intensive due to expensive reconstruction and probabilistic modeling needed. We present **WM-PPO**, a reconstruction-free contrastive learning world model combined with a PPO agent. While

WM-PPO does not match the peak performance of DreamerV3, it outperforms the PPO baseline by large margins, while being much more computationally efficient.

Our project work can be summarized as followings:

1. Reproduced the results of DreamerV3 and achieving similar metrics.
2. Implemented a traditional machine learning method using Q-Learning for value iteration and linear regression (or SVM) for value estimation
3. Reproduced a deep RL algorithm, specifically PPO
4. Created a reconstruction-free world model using contrastive learning where a PPO agent is trained in its imagined rollouts, similar to DreamerV3.

In the following sections, we will first introduce the background knowledge needed to understand both the traditional ML baseline and deep RL baseline in [Section 2](#), followed by their implementation details in [Section 3](#). Our proposed method is discussed in depth in [Section 4](#). The results of all aforementioned models are listed in [Section 5](#).

## 2 Background

### 2.1 Q-Learning

Q-Learning is a model-free reinforcement learning algorithm for estimating the value function of an optimal policy[5]. A value function  $V^\pi(s)$  tells us the expected discounted return if we were to follow policy  $\pi$  from state  $s$  onwards.

$$\begin{aligned} V^\pi(s_0) &= \mathbb{E}_\pi[\tau \mid s] \\ &= R(s_0) + \mathbb{E}[\gamma R(s_1) + \gamma^2 R(s_2) + \dots \mid s_0 = s, \pi] \\ &= R(s_0) + \mathbb{E}[\gamma V^\pi(s_1) \mid s_0 = s, \pi] \\ &= R(s_0) + \gamma \sum_{a_0 \in A} \sum_{s' \in S} \pi(a|s) p(s'|s_0, a_0) V^\pi(s') \\ &= R(s_0) + \gamma \sum_{s' \in S} p_{a_0 \sim \pi(s_0)}(s'|s_0, a_0) V^\pi(s') \end{aligned}$$

However, knowing the value function alone does not allow us to determine the optimal action at each step unless a model of the environment is available. Q-Learning solves this by estimating an action-value function  $Q^{\pi(s,a)}$  that gives the expected return for taking action  $a$  from state  $s$  and following policy  $\pi$  thereafter.

$$\begin{aligned} V^\pi(s) &= R(s) + \gamma \sum_{a \in A} \sum_{s' \in S} \pi(a|s) p(s'|s, a) V^\pi(s') \\ &= R(s) + \gamma \sum_{a \in A} \pi(a|s) \sum_{s' \in S} p(s'|s, a) V^\pi(s') \\ &= R(s) + \gamma \sum_{a \in A} \pi(a|s) Q^\pi(s, a) \\ \therefore Q^\pi(s, a) &= \sum_{s' \in S} p(s'|s, a) V^\pi(s') = \mathbb{E}[\tau \mid s, a] \end{aligned}$$

### 2.2 Dreamer

Dreamer is a series of work that attempts to learn an accurate enough world model that allows a reinforcement learning agent to learn purely from its synthetic data. The world model essentially acts as a virtual simulator that extrapolates the existing training data, thereby enabling sample efficiency.

DreamerV3 is the latest iteration of the Dreamer series of work, which uses a variational autoencoder (VAE) to encode the observation into a categorical latent space, which is then regularized by a KL divergence loss with the dynamics prior[4]. The VAE is trained to reconstruct the observation, while the KL loss encourages the latent to encode the only information that is relevant to predictions.

---

**Algorithm 1:** Dreamer

---

```

Initialize dataset  $\mathcal{D}$  with  $S$  random episodes; network w/ parameters  $\theta, \varphi, \psi$  randomly.
1   while not converged do
2     for update step  $c = 1..C$  do
3       Draw  $B$  data sequences  $\{(a_t, o_t, r_t)\}_{t=k}^{k+L} \sim \mathcal{D}$ .
4       Compute model states  $s_t \sim p_\theta(s_t | s_{t-1}, a_{t-1}, o_t)$ .
5       Update  $\theta$  using representation learning.
6       Imagine trajectories  $\{(s_\tau, a_\tau)\}_{\tau=t}^{t+H}$  from each  $s_t$ 
7       Predict rewards  $\mathbb{E}(q_\theta(r_\tau | s_\tau))$  and values  $v_\psi(s_\tau)$ .
8       Compute value estimates  $V_\lambda(s_\tau)$  via Equation 6.
9       Update  $\varphi \leftarrow \varphi + \alpha \nabla_\varphi \sum_{\tau=t}^{t+H} V_\lambda(s_\tau)$ .
10      Update  $\psi \leftarrow \psi - \alpha \nabla_\psi \sum_{\tau=t}^{t+H} \frac{1}{2} \|v_\psi(s_\tau) - V_\lambda(s_\tau)\|^2$ .
11    end for
12    for time step  $t = 1..T$  do
13      Compute  $s_t \sim p_\theta(s_t | s_{t-1}, a_{t-1}, o_t)$  from history.
14      Compute  $a_t \sim q_\varphi(a_t | s_t)$  with the action model.
15      Add exploration noise to action.
16       $r_t, o_{t+1} \leftarrow \text{env.step}(a_t)$ .
17      Add experience to dataset  $\mathcal{D} \leftarrow \mathcal{D} \cup \{(o_t, a_t, r_t)_{t=1}^T\}$ 
18    end for
19  end while

```

---

Algorithm 1: Dreamer training procedure. Source: [6]

### 2.3 Proximal Policy Optimization

Proximal Policy Optimization (PPO) is a policy gradient method that, instead of using a KL trust region to constrain the policy updates as in TRPO, uses a clipped surrogate objective. Just like TRPO, its standard implementation follows an actor-critic framework. In which, critic is a value network  $V_\varphi(s)$  that aims to accurately estimate the expected return under state  $s$ ; the actor is a policy network  $\pi_\theta(a|s)$  that aims to output actions that maximizes the expected reward according to  $V_\varphi(s)$ .

Traditionally, the actor network  $\pi_\theta(a|s)$  is trained to directly maximize the expected advantage  $\mathbb{E}[A^{\pi_{\theta_k}}(s, a)]$ , where  $A^{\pi_{\theta_k}}(s, a) = Q^{\pi_{\theta_k}}(s, a) - V^{\pi_{\theta_k}}(s)$ . But this can leads to unstable training, as the locally optimal policy improvement direction (gradient) might actually lead to worse performance.

To overcome this, we need to ensure that the new policy  $\pi_\theta$  monotonically increases the expected advantage  $A^{\pi_{\theta_k}}(s, a)$  over the old policy  $\pi_{\theta_k}$ . If it can be guaranteed that a new policy's generated trajectory will yield higher advantage estimation according to the old critic network, then it can be proven that the new policy will converge. However, evaluating the new policy using the value function of the old policy causes a chicken-and-egg problem, because the new policy is the unknown variable.

$$\arg \max_{\theta} \mathbb{E} \left[ \min \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a) \right), \text{clip} \left( \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \varepsilon, 1 + \varepsilon \right) A^{\pi_{\theta_k}}(s, a) \right]$$

PPO solves this by using a surrogate objective that evaluates episodes from the old policy, essentially assuming that the new policy's state visitation will be similar to the old policy's. To meet such assumption, the difference between  $\pi_\theta$  and  $\pi_{\theta_k}$  is clipped between a hyperparameter  $\varepsilon$ .

### 3 Baselines

#### 3.1 Machine Learning Baseline

---

**Algorithm 2:** Q-Learning with Traditional Machine Learning Models

---

```

1 Initialize  $Q_\theta$  with SVM or Linear Regression
2 Initialize replay buffer  $\mathcal{D} \leftarrow \emptyset$ 
3 for episode = 1,  $M$  do
4   for  $t = 1, T$  do
5     With probability  $\varepsilon$  select random action  $a_t$ 
6     otherwise select  $a_t = \arg \max_a Q_\theta(s_t, a)$ 
7     Execute action  $a_t$  in emulator and observe reward  $r_t$  and next state  $s_{t+1}$ 
8     Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ 
9     Sample batch of  $B$  transitions  $(s_j, a_j, r_j, s_{j+1})$  from  $\mathcal{D}$ 
10    Set  $y_j = r_j$  if  $s_{j+1}$  is terminal otherwise  $y_j = r_j + \gamma \max_{a'} Q_\theta(s_{j+1}, a')$ 
11    Update  $Q_\theta$  using  $\mathcal{L}_{\text{MSE}} = (y_j - Q_\theta(s_j, a_j))^2$ 
12  end for
13 end for

```

---

Algorithm 2: Q-Learning with Traditional ML Models. Source: [7]

##### 3.1.1 Action Discretization

The *DeepMind Control Suite* contains purely continuous control tasks. However, the default implementation of Q-Learning can only handle discrete action spaces. Therefore, to properly apply Q-Learning to the tasks, discretization is required. Specifically, each action dimension is divided into  $N_{\text{bins}}$  intervals according to the configured *action\_bins* number and generate all possible action combinations. Formally, let the action space be a continuous space  $A = \{a_1, a_2, \dots, a_n\} \subset \mathbb{R}^m$ . Each dimension of the space  $a_i$  is then discretized into  $a_i = \{b_1, \dots, b_{N_{\text{bins}}}\}$ , where  $b_j = j \cdot \frac{1}{N_{\text{bins}} - 1}$ .

However, this can lead to a large number of combinations, for high dimensional continuous action space. For those tasks,  $N_{\text{sample}}$  actions are uniformly drawn from the entire action space. In our code, the method `_create_discrete_actions` implements this process, checking the type of action space and generating the appropriate discrete actions following the aforementioned rules.

### 3.1.2 Epsilon-Greedy Strategy

The  $\varepsilon$ -greedy strategy is a popular method used to balance *exploration* and *exploitation* in reinforcement learning. *Exploration* represents that the agent tries new actions to discover their potential rewards, which is crucial for learning about the environment. *Exploitation* represents that the agent selects the action that it believes has the highest expected reward based on its current knowledge.

$\varepsilon$  is a parameter that determines the probability of choosing a random action (exploration) versus the best-known action (exploitation). Formally, a  $\varepsilon$ -greedy policy  $\pi_\varepsilon(a|s)$  chooses the optimal action with probability  $\frac{\varepsilon}{|A|} + 1 - \varepsilon$  and any other actions with a uniform probability  $\frac{\varepsilon}{|A|}$ .

$$\pi_\varepsilon(a|s) = \begin{cases} \frac{\varepsilon}{|A|} + 1 - \varepsilon & \text{if } a = \arg \max_a Q(s, a) \\ \frac{\varepsilon}{|A|} & \text{otherwise} \end{cases}$$

### 3.1.3 Update the Q-Function

First, the agent interacts with the environment, collects state transition samples  $(s_t, a_t, r_t, s_{t+1}, \text{done})$ , and stores them in an experience replay buffer. When it's time to update the Q-function, a batch of samples is randomly drawn from the buffer. This batch sampling mechanism (Experience Replay) breaks the temporal correlation between samples, making training more stable while improving sample efficiency. For each experience sample, we first convert the observation state into a flattened *feature vector*. This feature vector is passed into the Q-function approximator which predicts  $\min(|A| \times N_{\text{bins}}, N_{\text{sample}})$  Q-values, each corresponding to an discretized action's Q-value.

Next, a temporal difference (TD) target value is calculated for each transition sample:

$$y = \begin{cases} r & \text{if terminal} \\ r + \gamma \max_{a'} Q(s', a') & \text{otherwise} \end{cases}$$

For non-terminal states, we find the maximum Q-value across all possible actions in the next state, multiply by the discount factor gamma, and add the immediate reward.

Finally, the Q-function, which is a SVM or Linear Regression model implemented via sklearn under the hood is fitted by minimizing the MSE between the predicted Q-values and the TD target:

$$\mathcal{L}_{\text{MSE}} = (y_j - Q_\theta(s_j, a_j))^2$$

### 3.2 Deep Learning Baseline

---

**Algorithm 3:**PPO, Actor-Critic Style

---

```

1 for iteration= 1, 2, ... do
2   for actor= 1,  $N$  do
3     Run policy  $\pi_{\theta_{\text{old}}}$  in environment for  $T$  timesteps
4     Compute advantage estimates  $\hat{A}_1, \dots, \hat{A}_T$ 
5   end for
6   Compute policy difference,  $r_t(\theta) \leftarrow \pi_\theta(a_t|s_t)/\pi_{\theta_{\text{old}}}(a_t|s_t)\hat{A}_t$ 
7   Minimize  $\mathcal{L}_{\text{CLIP}}(\theta) = \min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \varepsilon, 1 + \varepsilon)\hat{A}_t)$ 
8 end for

```

---

Algorithm 3: PPO training procedure. Source: [8]

For our deep learning baseline, we chose the Proximal Policy Optimization (PPO) algorithm, which is a popular policy gradient method that uses a clipped surrogate objective to update the policy. For simplicity, we directly adopt the PPO implementation from <https://github.com/vwxyzjn/cleanrl> [9] and slightly modified it to fit our needs.

The CleanRL implementation is a single-file PPO implementation that is designed to be simple and easy to understand[10]. In this implementation, serveral slight improvements were made compared to the original author's approach, such as gradient clipping and learning rate annealing. The performance of the PPO baseline is shown in [Section 5](#).

## 4 Method

### 4.1 Hypothesis

In DreamerV3, an expensive variational autoencoder (VAE) is used to encode the observation into a latent space, which is then regularized by a KL divergence loss with the dynamics prior. While these operations allows DreamerV3 to learn robust and accurate representations for the actor-critic agent to learn purely from latent space imaginations, we hypothesize that we can achieve lower but still competitive performance without the expensive reconstruction and probabilistic modeling.

### 4.2 Design

We introduce **WM-PPO**, short for **World Model-PPO**, a simplified Dreamer-like method that learns its world model through contrastive learning instead of reconstruction. The models starts by sampling a batch of  $B$  trajectories  $\tau_{i=1..B} = \{s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T\}_i$  from the replay buffer  $\mathcal{D}$ . Then, using an encoder  $q_\varphi : o_t \rightarrow z_t$  that maps each raw observation to latent representations. This gives us a batch of  $B$  latent trajectories  $z_{i=1..B} = \{z_0, z_1, \dots, z_T\}_i$ , which we will from now on refer to as the *anchor latents*. Next, the dynamics model  $p_\varphi : z_t, a_t \rightarrow z_{t+1}$  predicts an imagined latent trajectory  $\hat{z} = \{\hat{z}_1, \hat{z}_2, \dots, \hat{z}_T\}$  from the first latent  $z_0$  and actions  $a_{0..T-1}$ .

Here, we make the critical assumption that the latent imagination is strictly worse than the latent directly encoded from the observation. This allows us to treat imagined data as negative samples. Given this, the model can then be optimized using a constrastive learning strategy that resembles InfoNCE, except we use batches of imagined trajectories  $\hat{z}$  as

negatives instead of the anchor latents  $\mathbf{z}$ . Specifically, we train the model to minimize the difference between  $(z_i, \hat{z}_i)$  using cross entropy; then, the softmax function will implicitly push away  $(z_i, \hat{z}_j)$  for  $j \neq i$ . Formally, the contrastive loss can be written as:

$$\mathcal{L}_{\text{contrast}}(\mathbf{z}, \hat{\mathbf{z}}) = -\mathbb{E}_i \left[ \log \left( \frac{\exp(s(\mathbf{z}_i, \hat{\mathbf{z}}_i)/\tau)}{\sum_j \exp(s(\mathbf{z}_i, \hat{\mathbf{z}}_j)/\tau)} \right) \right] \quad \text{where, } s(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u}^T \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

Intuitively, given pairs of anchor and imagined latent trajectories  $\mathbf{z}_i$  and  $\hat{\mathbf{z}}_i$ , a  $T \times T$  cosine similarity matrix is computed between each pair of latent vectors. It is straightforward to see

#### Algorithm 4: WM-PPO Implementation

```

1 Initialize replay buffer  $\mathcal{D}$ , world model params  $\theta$ , agent params  $\varphi$ .
2 Initialize optimizers Adam $_{\theta}$ , Adam $_{\varphi}$ .
3 Prefill  $\mathcal{D}$  with  $N_{\text{prefill}}$  episodes using random actions  $a \sim \text{Uniform}()$ .
4 while not converged do
5   for update step  $c = 1..C$  do
6     for  $i = 1..N_{\text{wm}}$  do
7       Sample  $B$  trajectories  $\{(o_t, a_t, r_t)\}_{t=k}^{k+L} \sim \mathcal{D}$ .
8       Encode anchor:  $\mathbf{z} \leftarrow q_{\theta}(o_{t=1..T})$ 
9       Imagine latent:  $\hat{\mathbf{z}} \leftarrow p_{\theta}(z_0, a_{0..T-1})$ 
10      Predict reward:  $\hat{r} \leftarrow r_{\theta}(\hat{\mathbf{z}})$ .
11      Compute: contrastive loss  $\mathcal{L}_{\text{contrast}}(\mathbf{z}, \hat{\mathbf{z}})$ 
12      Compute: reward loss  $\mathcal{L}_{\text{reward}}(\hat{r}, r)$ .
13      Compute world model loss  $\mathcal{L}_{\text{WM}} = \mathcal{L}_{\text{contrast}} + c_r \mathcal{L}_{\text{reward}}$ .
14      Update  $\theta \leftarrow \text{Adam}_{\theta}(\nabla_{\theta} \mathcal{L}_{\text{WM}})$ .
15    end for
16     $\mathbf{Z}_{\text{flat}} : \mathbb{R}^{B \times L} \leftarrow \text{Flatten}(\mathbf{z})$ 
17    Imagine  $B \times L$  rollouts of  $H$  steps from  $\mathbf{Z}_{\text{flat}}$ 
18    Train PPO actor-critic  $\pi_{\varphi}$  and  $v_{\varphi}$  on  $B \times L$  imagined rollouts
19  end for
20   $o_t = \text{env.reset}()$ 
21  for  $i = 1..N_{\text{interact}}$  do
22    while episode is not done do
23      Encode  $z_t = q_{\theta}(o_t)$ 
24      Sample action  $a_t \sim \pi_{\varphi}(\cdot | z_t)$ 
25      Execute action  $o_{t+1}, r_t, d_t = \text{env.step}(a_t)$ 
26      Store  $(o_t, a_t, r_t, d_t)$  in  $\mathcal{D}$ 
27    end while
28  end for
29 end while

```

that the objective should be to maximize the values on the diagonal. Therefore, the objective can be formulated as a classification problem with a batch of  $T$  samples whose labels are  $0, 1, 2, \dots, T$ .

$$\begin{pmatrix} \cos(z_1, \hat{z}_1) & \cos(z_1, \hat{z}_2) & \dots & \cos(z_1, \hat{z}_T) \\ \cos(z_2, \hat{z}_1) & \cos(z_2, \hat{z}_2) & \dots & \cos(z_2, \hat{z}_T) \\ \vdots & \vdots & \ddots & \vdots \\ \cos(z_T, \hat{z}_1) & \cos(z_T, \hat{z}_2) & \dots & \cos(z_T, \hat{z}_T) \end{pmatrix}$$

With this objective, the model is encoder  $q_\varphi$  and dynamics model  $p_\varphi$  are trained end-to-end, except for a reward predictor that predict the immediate reward  $r_t$  from the anchor  $z_t$  at each timestep which is trained independently without its gradient flowing back to the  $p_\varphi$  and  $q_\varphi$ .

Once the world model is setup, we follow the same Dreamer training procedure as specified in [Algorithm 1](#). The world model is trained on batches of real trajectories from the buffer  $\mathcal{D}$ , a PPO agent is trained on imagined rollouts from the world model, then the training is paused while the PPO agent interacts with the environment using the latest policy to collect a specified number of new episodes.

### 4.3 Improvements

In this subsection, we will discuss the attempts that we made to improve the model performance.

**Learning Rate Balancing** throughout our experiments, it has been observed that the  $p_\varphi$  is much harder to train than  $q_\varphi$ , which intuitively makes sense since the environment dynamics is complex and much harder to model. Regularizing the representations of the dynamics imagination and the encoder’s output will equal strength destabilizes training. We take inspiration from DreamerV2 and introduce a hyperparameter  $\alpha < 0.5$  that regularizes the encoder towards the dynamics with strength  $\alpha$  and the reverse with strength  $1 - \alpha$ . This way, the encoder is not being forced to match a poorly modeled dynamics during early stage of the training.

**Activation Functions** Numerous different activation functions are tested. Eventually, we found that the SiLU with layer normalization and bias disabled performs the best. Common activations such as ReLU frequently experience dying neuron problems, especially in the reward predictor, as the reward values are often zero during early stage of the training.

**Reward Normalization** the reward values are normalized to have ranges  $[0, 1]$ . We notice that this improves the stability of the agent learning, especially in environments with sparse yet high rewards. From implementation perspective, subtracting the reward by a bias does not change the gradient; therefore, the reward value is directly divided by a constant.

**Action Embedding** Initially, the raw actions were fed into the dynamics model. This oversight was later noticed and an additional action encoder is introduced to embed the raw actions, which gets concatenated with the latent  $z_t$  before passed into the dynamics model.

## 5 Results and Discussions

### 5.1 Traditional machine learning approaches (Linear Regression and SVM)

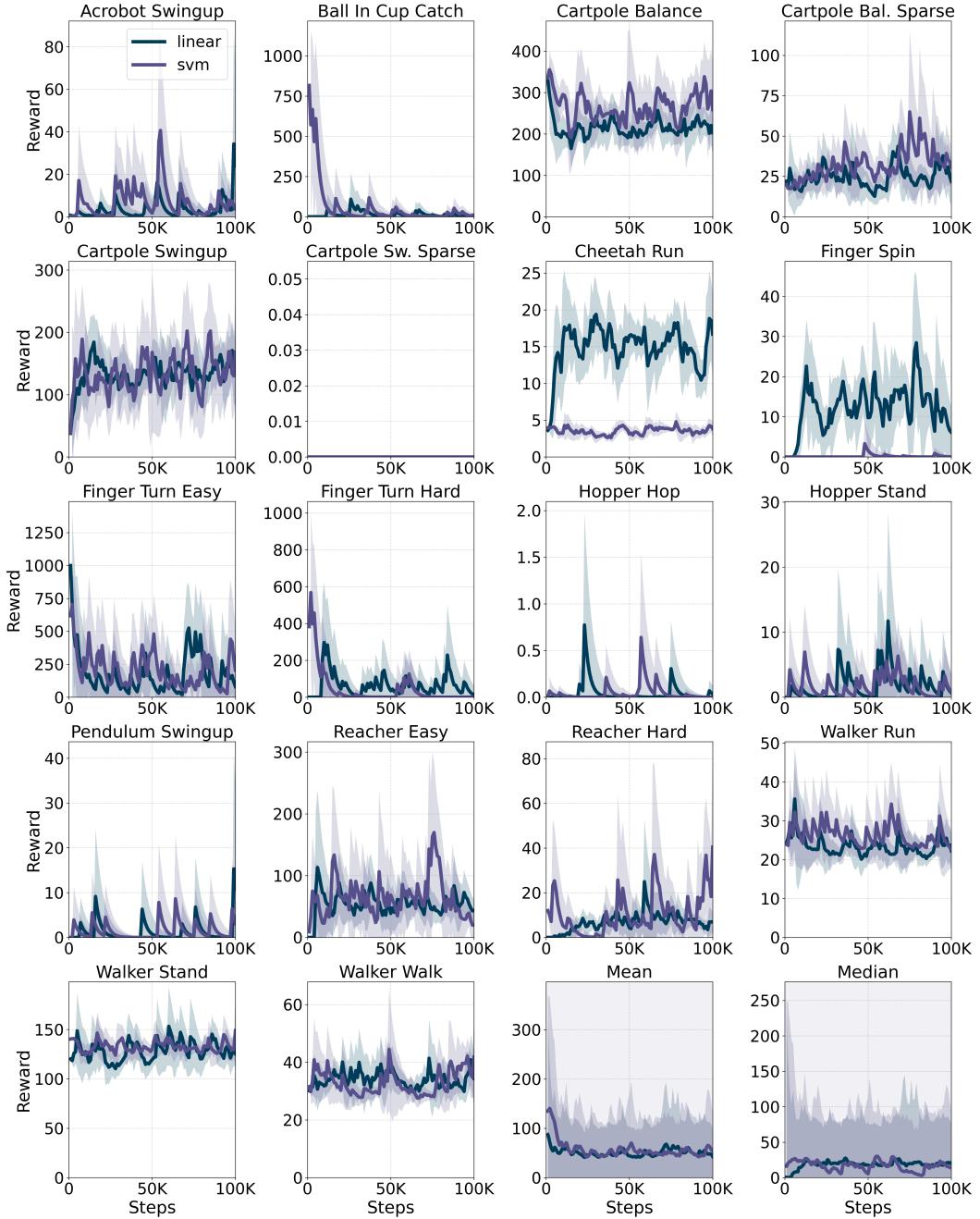


Figure 1: Results of running 18 classic DeepMind Control proprioceptive tasks using Linear Regression and SVM. Neither SVM nor linear regression based Q-Learning exhibits noticeable increase across the tasks. These algorithms clearly failed to converge and are no better than random actions. This aligns with our expectation that these simple machine learning models does not have the expressiveness necessary for fitting the complex Q-function, let alone continuos controls.

## 5.2 DreamerV3 and our replication of DreamerV3

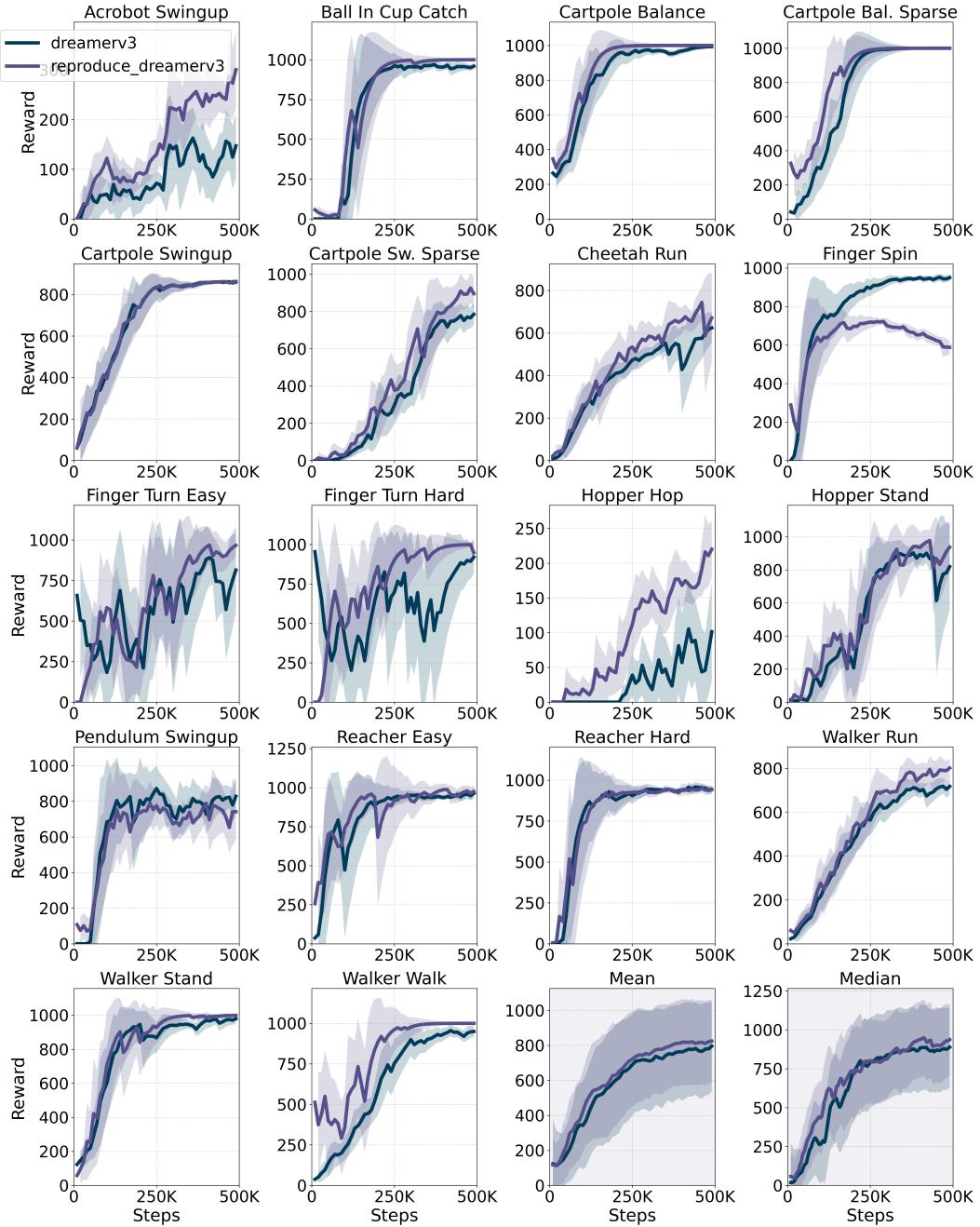


Figure 2: Official and reproduced results of DreamerV3. The official results are sourced from <https://github.com/danijar/dreamerv3/tree/main/scores>. The reproduced results use a third party PyTorch-based implementation of DreamerV3 instead of the official JAX results due to compatibility concerns. Therefore, the scores in certain tasks differ from the official results; this is expected as noted by the author of the torch implementation (See: <https://github.com/NM512/dreamerv3-torch> README.md for more details)

### 5.3 Comparison of DreamerV3, PPO, and our proposed WM-PPO

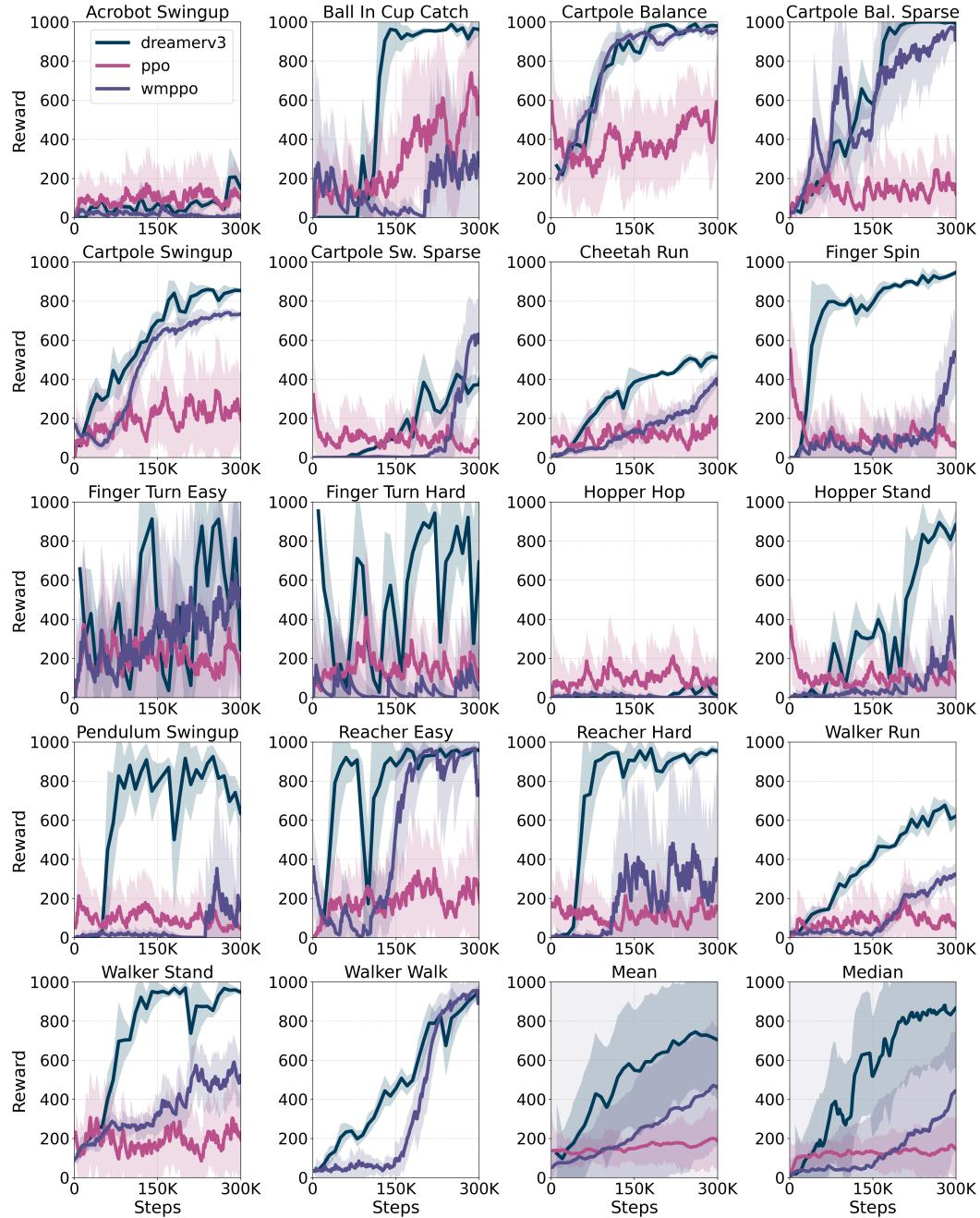


Figure 3: Results directly comparing DreamerV3, PPO and our proposed WM-PPO approach. As shown by the results, WM-PPO excels in many tasks where the baseline PPO stagnates. While our WM-PPO’s peak performance does not match that of DreamerV3’s, we hypothesize that cause to be the slower learning speed of the world model due to lack of reconstruction gradients. This evident by a spiky imagined reward values during early stage of training, that gradually converges after 100K environment steps, suggesting that the PPO agent might be exploiting an immature world model during early stages of training.

Task	Linear	SVM	PPO	DreamerV3	WM-PPO
Environment steps	100K	100K	500K	500K	500K
Acrobot Swingup	5.8	5.1	<u>99.6</u>	<b>132.8</b>	7.5
Ball In Cup Catch	7.9	10.6	<b>620.0</b>	<b>956.9</b>	261.0
Cartpole Balance	216.0	274.6	464.4	<b>978.1</b>	<b>952.3</b>
Cartpole Balance Sparse	26.1	40.9	341.0	<b>1000.0</b>	<b>887.7</b>
Cartpole Swingup	136.8	146.4	278.8	<b>857.8</b>	<b>721.6</b>
Cartpole Swingup Sparse	0.0	0.0	80.6	<b>760.1</b>	<b>290.1</b>
Cheetah Run	15.2	3.7	138.2	<b>557.4</b>	<b>281.9</b>
Finger Spin	13.4	0.1	105.3	<b>945.6</b>	<b>238.1</b>
Finger Turn Easy	230.6	161.3	296.7	<b>758.1</b>	<b>447.5</b>
Finger Turn Hard	59.8	0.1	<u>195.9</u>	<b>738.6</b>	45.2
Hopper Hop	0.0	0.0	<u>91.5</u>	<b>72.4</b>	0.5
Hopper Stand	1.6	1.5	100.0	<b>837.1</b>	<b>151.4</b>
Pendulum Swingup	2.4	1.5	106.0	<b>800.2</b>	<b>111.6</b>
Reacher Easy	52.1	65.4	355.6	<b>945.9</b>	<b>900.9</b>
Reacher Hard	7.6	15.1	165.2	<b>944.1</b>	<b>323.5</b>
Walker Run	23.2	24.6	104.0	<b>706.4</b>	<b>266.6</b>
Walker Stand	130.4	133.1	172.3	<b>962.9</b>	<b>510.4</b>
Walker Walk	33.9	35.1	106.1	<b>936.6</b>	<b>867.9</b>
Task Mean	53.49	51.06	212.29	<b>771.72</b>	<b>403.7</b>
Task Median	19.2	12.85	151.7	<b>847.45</b>	<b>286.0</b>

Table 1: The final EMA scores of the methods.

## 6 Hyperparameters

Q-Learning			WM-PPO		
Hyperparam.	Notation	Value	Hyperparam.	Notation	Value
Discount factor	$\gamma$	0.99	WM Batch size	$B$	50
Learning rate	$\alpha$	0.1	WM learning rate	$\alpha_{WM}$	0.0002
Exploration rate	$\varepsilon$	0.3	LR Balancing factor	$\alpha$	0.2
Episodes	$M$	100	Temperature	$\tau$	0.3
Episode length	$T$	1000	Reward coeff.	$c_r$	1.0
Buffer capacity	—	10000	Latent dimension	$z_{dim}$	128
Batch size	$B$	64	FC dimension	$fc_{dim}$	256
Action bins	$N_{bins}$	5	FC Layer count	$fc_{layers}$	8
Action samples	$N_{samples}$	1000	PPO Learning rate	$\alpha_{PPO}$	0.00001
PPO			Discount factor	$\gamma$	0.99
Learning rate	$\alpha$	3e-4	GAE parameter	$\lambda$	0.95
Discount factor	$\gamma$	0.99	Clipping coeff.	$\varepsilon$	0.2
GAE parameter	$\lambda$	0.95	Critic coeff.	$c_v$	0.5
Clipping coefficient	$\varepsilon$	0.2	Entropy coeff.	$c_e$	0.001
Entropy coefficient	$c_1$	0.0	PPO Update epochs	$K$	4
Critic coeff.	$c_2$	0.5	Imag. horizon	$H$	15
Update epochs	$K$	10	PPO batch size	—	256
Rollout length	$T$	2048	WM Window size	$L$	50
Max gradient norm	—	0.5	WM updates ratio	$N_{wm}$	10
			LR annealing	—	True

## 7 Dataset Description

We hereby confirm that we have met the requirements in the project description. For our project, there is not fixed dataset, the training data is collected on-the-fly through agent interaction. For the *DeepMind Control Suite* which we used, each environment has an observation space of roughly 5-10, with an action space of up to 20 ( $d \approx 30$ ). Since we trained for 100K steps ( $m = 100,000$ ), this sums up to  $(20 + 10) \times 100,000 = 3,000,000$ , successfully meeting the requirement of  $m * d \geq 1,000,000$ .

## 8 References

- [1] D. Silver *et al.*, “Mastering the game of Go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, 2016, doi: [10.1038/nature16961](https://doi.org/10.1038/nature16961).
- [2] OpenAI *et al.*, “Dota 2 with Large Scale Deep Reinforcement Learning.” [Online]. Available: <https://arxiv.org/abs/1912.06680>
- [3] O. Vinyals *et al.*, “StarCraft II: A New Challenge for Reinforcement Learning.” [Online]. Available: <https://arxiv.org/abs/1708.04782>
- [4] D. Hafner, A. Pashevich, J. Ba, N. Heess, T. Lillicrap, and M. Norouzi, “Mastering diverse domains through world models,” *arXiv preprint arXiv:2301.04104*, 2023.
- [5] C. J. Watkins and P. Dayan, “Q-learning,” *Machine learning*, vol. 8, no. 3–4, pp. 279–292, 1992.
- [6] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba, “Dream to control: Learning behaviors by latent imagination,” *arXiv preprint arXiv:1912.01603*, 2019.
- [7] V. Mnih and others, “Playing Atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
- [8] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- [9] CleanRL, “PPO implementation.” [Online]. Available: <https://github.com/vwxyzjn/cleanrl>
- [10] S. Huang *et al.*, “CleanRL: High-quality Single-file Implementations of Deep Reinforcement Learning Algorithms,” *Journal of Machine Learning Research*, vol. 23, no. 274, pp. 1–14, 2022.