

RL-Course 2025/26: Final Project Report

Ansel Cheung, Jannik Mänzer, Niklas Abraham

February 24, 2026

Abstract

This report presents our application and comparative study of modern Reinforcement Learning (RL) algorithms specifically, Twin Delayed Deep Deterministic Policy Gradient (TD3), Soft Actor-Critic (SAC), and the model based TD-MPC2 within the challenging Laser Hockey environment. We discuss the environment’s design, its state and action spaces, and the unique characteristics making it a compelling testbed for RL research. Our work details the methodological approaches taken with these algorithms, summarizes key experimental findings, and reflects on lessons learned in this multi agent, continuous control setting.

1 Introduction

1.1 Environment Overview

The Laser Hockey environment [6] is a custom reinforcement learning benchmark built on the Gymnasium Python API and powered by the Box2D physics engine. In this multi agent setting, two agents each control a hockey stick and compete to score goals by striking a puck into the opponent’s net. The environment is designed such that both agents receive identical but mirrored observations at each timestep, which eliminates the need for side specific learning strategies.

At each timestep, the agent receives an 18 dimensional continuous state vector that captures the complete game state. This observation includes the position x_1, y_1 and orientation θ_1 of both players relative to the center of the field, along with their linear $v_{x,1}, v_{y,1}$ and angular ω_1 velocities. The state also contains the puck’s position x_p, y_p and velocity $v_{x,p}, v_{y,p}$, as well as time remaining for each player’s puck possession $t_{\text{puck},1}, t_{\text{puck},2}$ in keep mode, which ranges from 0 to 15 steps. Player 1 refers to the agent currently being controlled, while Player 2 represents the opponent. The observation structure ensures that when the viewpoint switches during training, the indices are automatically mirrored so that each agent always perceives itself as Player 1, maintaining a consistent learning perspective.

The action space consists of a 4 dimensional continuous vector that controls the agent’s stick. The first two components specify forces applied for stick movement in the x and y directions, while the third component controls the torque applied to adjust the stick’s angle. The fourth component is a thresholded scalar that determines whether to release the puck when the agent is in possession. Reward feedback in this environment is sparse and goal oriented. An agent receives a reward of +10 for scoring a goal and 0 for a draw. Additionally, there is a small reward signal for maintaining proximity to the puck, which helps guide exploration during early learning stages. Each episode begins with all entities positioned at the center of the field and terminates when a goal is scored or a timeout occurs. The environment supports both scripted opponents, such as the BasicOpponent, and learning agents as adversaries, making it a versatile and challenging benchmark for continuous control and multi agent reinforcement learning research.

2 Method

2.1 TD-MPC 2 - Niklas Abraham

TD-MPC2 [4] is a model-based reinforcement learning algorithm that learns a world model to predict future states and rewards, and uses this model to select actions through planning. The core idea is to learn a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ in a Markov Decision Process with an infinite horizon. The policy is constructed to maximize the expected discounted return. In TD-MPC2, this is achieved by learning a world model and selecting actions by planning with the learned models.

For planning, TD-MPC2 employs the Model Predictive Control (MPC) framework, in which actions are optimized based on planning over action sequences of a finite horizon H :

$$\pi(s_t) = \arg \max_{a_1, \dots, a_H} \mathbb{E}_\pi \left[\sum_{\tau=0}^H \gamma^\tau r(s_{t+\tau}, a_{t+\tau}) \right]. \quad (1)$$

The return of each trajectory is estimated by simulating action sequences through the learned world model. However, this approach often leads to only locally optimal policies. To address this limitation, TD-MPC2 additionally utilizes a value function to guide the planning process and improve the policy toward a more globally optimal solution.

Rather than predicting raw future observation states, TD-MPC2 learns to predict a maximally useful latent representation for accurately estimating the outcomes of action sequences. The algorithm is composed of five distinct neural network components that interact in a coordinated manner, as illustrated in Figure 1 a).

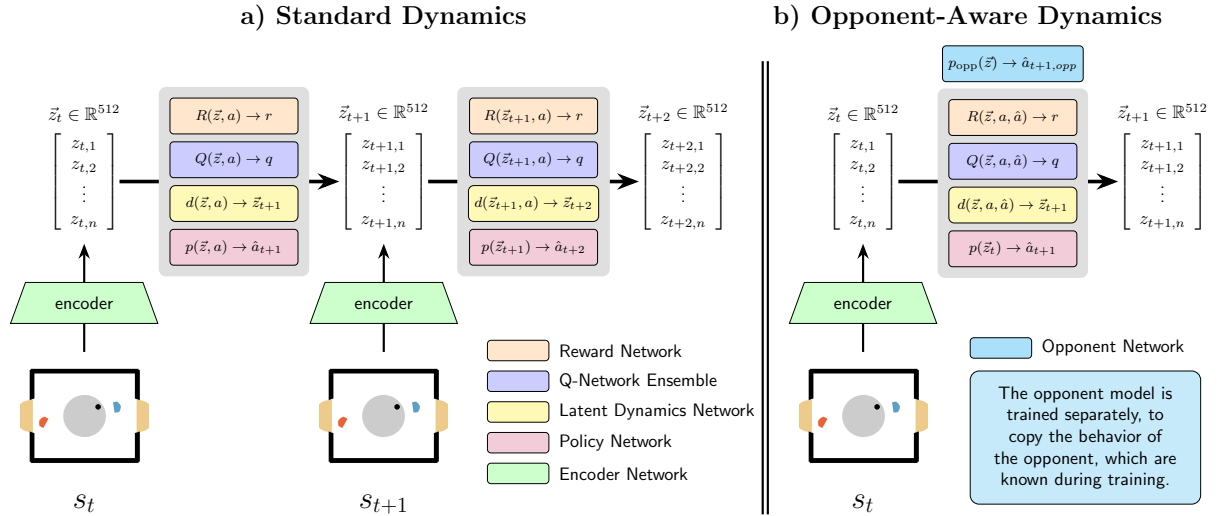


Figure 1: On the left in a) TD-MPC2 agent architecture and their individual components. b) shows the modified architecture with an additional opponent network to model the behavior of the adversarial agent in the Laser Hockey environment.

- **Encoder:** Maps the observed state s to a 512-dimensional latent vector $\vec{z} = h(s)$.
- **Latent Dynamics:** Predicts the next latent \vec{z}_{t+1} from current latent and action: $\vec{z}_{t+1} = d(\vec{z}_t, a)$.
- **Reward Head:** Estimates reward r for a given (\vec{z}, a) pair: $r = R(\vec{z}, a)$.

- **Termination Head:** Predicts early episode end, e.g., when a goal is imminent.
- **Q-Network Ensemble:** An ensemble (5 networks) of Q-functions estimating value $q = Q(\vec{z}, a)$. The minimum of two sampled networks reduces value overestimation.
- **Policy Network:** Guides action selection in planning: $p(\vec{z}, a) \rightarrow \hat{a}$.

2.1.1 Architecture and Training

All network components are multi-layer perceptrons (MLPs) with Mish activations. As in [4], the latent representation \vec{z} is projected into L -dimensional simplices via a softmax to stabilize training and enforce sparsity.

Training uses an experience replay buffer \mathcal{B} with full episode trajectories. Model parameters are optimized over sampled subsequences of length $H+1$ from \mathcal{B} by minimizing a joint loss for dynamics, reward, and value prediction:

$$\mathcal{L}(\theta) = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1})_{t=0}^H \sim \mathcal{B}} \left[\sum_{t=0}^H \lambda^t \left(\|\vec{z}_{t+1} - \text{sg}(h(s_{t+1}))\|_2^2 + \text{CE}(\hat{r}_t, r_t) + \text{CE}(\hat{q}_t, q_t) \right) \right], \quad (2)$$

where $\text{sg}(\cdot)$ is stop-gradient, $\vec{z}_{t+1} = d(\vec{z}_t, a_t)$ is the predicted next latent, $\hat{r}_t = R(\vec{z}_t, a_t)$, $\hat{q}_t = Q(\vec{z}_t, a_t)$, and λ is a temporal discount factor. The Q-value target is $q_t = r_t + \gamma \bar{Q}(\vec{z}_{t+1}, p(\vec{z}_{t+1}, a_{t+1}))$, using an EMA of Q-net parameters (\bar{Q}) for stability. Following TD-MPC2, reward and value predictions are regressed in a log-transformed space with cross-entropy loss and soft targets.

The policy p is optimized according to a maximum entropy RL objective:

$$\mathcal{L}_p(\theta) = \mathbb{E}_{(s_t, a_t)_{t=0}^H \sim \mathcal{B}} \left[\sum_{t=0}^H \lambda^t \left(\alpha Q(\vec{z}_t, p(\vec{z}_t, a_t)) - \beta \mathcal{H}(p(\cdot | \vec{z}_t)) \right) \right], \quad (3)$$

where $\vec{z}_{t+1} = d(\vec{z}_t, a_t)$ with $\vec{z}_0 = h(s_0)$, and $\mathcal{H}(p(\cdot | \vec{z}_t))$ is the policy entropy. Hyperparameters α and β balance value maximization and entropy, preventing premature collapse to deterministic policies.

2.1.2 Planning with MPPI

For local planning, TD-MPC2 leverages Model Predictive Path Integral (MPPI) control [8], sampling action sequences with guidance from the policy network. At each step, it estimates $\mu^*, \sigma^* \in \mathbb{R}^{H \times m}$, the mean and standard deviation of a multivariate Gaussian that maximizes expected return:

$$\mu^*, \sigma^* = \arg \max_{\mu, \sigma} \mathbb{E}_{a_{t:t+H} \sim \mathcal{N}(\mu, \sigma^2)} \left[\gamma^H Q(\vec{z}_{t+H}, a_{t+H}) + \sum_{\tau=t}^H \gamma^\tau R(\vec{z}_\tau, a_\tau) \right]. \quad (4)$$

This is optimized by iteratively sampling actions from $\mathcal{N}(\mu, \sigma^2)$, evaluating their returns, and updating μ and σ based on weighted top samples. The termination model predicts early ends in sampled rollouts. To speed up convergence, a fraction of samples comes from the policy p , and μ, σ are initialized from the previous step.

2.1.3 Modifying TD-MPC2: Opponent-Aware Dynamics

In the classical TD-MPC2, the dynamics model is trained to predict the next latent state given the current latent state and action. However, in the multi agent setting, with an adversarial opponent, the dynamics

model will only receive the current latent state and action, but not receive the action of the opponent. Thus the standard TD-MPC2 implicitly models

$$P(s_{t+1} \mid s_t, a_t^{\text{self}}) = \mathbb{E}_{a_t^{\text{opp}} \sim \pi_{\text{opp}}(\cdot \mid s_t)} P(s_{t+1} \mid s_t, a_t^{\text{self}}, a_t^{\text{opp}}), \quad (5)$$

which corresponds to marginalizing over opponent behavior. This leads to a mean-opponent model, producing biased long-horizon predictions when the opponent policy is multimodal or strategic. To address this, it is not enough to vary the opponents during training or to use self-play, the network architecture needs to be changed.

To model the opponent’s behavior, we extend the dynamics model so that it takes the opponent action as an additional input: $\vec{z}_{t+1} = d(\vec{z}_t, a_t^{\text{self}}, a_t^{\text{opp}})$, where a_t^{opp} is the opponent’s action. The same is done for the reward model and Q-value network.

With opponent-aware models, evaluating a candidate self-action sequence $a_t^{\text{self}}, \dots, a_{t+H-1}^{\text{self}}$ uses the following planning objective and rollout. The return is

$$\sum_{\tau=t}^{t+H-1} \gamma^{\tau-t} R(\vec{z}_\tau, a_\tau^{\text{self}}, \hat{a}_\tau^{\text{opp}}) + \gamma^H \tilde{V}(\vec{z}_{t+H}), \quad (6)$$

where the latent trajectory and predicted opponent actions are obtained by the recursion

$$\hat{a}_\tau^{\text{opp}} = \pi_{\text{opp}}(\vec{z}_\tau), \quad \vec{z}_{\tau+1} = d(\vec{z}_\tau, a_\tau^{\text{self}}, \hat{a}_\tau^{\text{opp}}), \quad (7)$$

for $\tau = t, \dots, t+H-1$, with $\vec{z}_t = h(s_t)$. The terminal value $\tilde{V}(\vec{z}_{t+H})$ is given by the Q-ensemble (e.g., $\min_k Q_k(\vec{z}_{t+H}, a)$ at the policy action a). MPPI then maximizes this return over sampled self-action sequences, with opponent actions fixed by the recursion above.

The action of the opponent needs to be predicted with a separate network, which receives as input the current latent state and outputs the action of the opponent. This is illustrated in Figure 1 b). This requires the opponent network to be trained separately to imitate the opponent’s behavior; the opponent’s actions are available in the replay buffer from collected episodes. In the setting in this project, we could choose between different opponents: basic weak, basic strong, the TD3 agent [?] or the SAC agent [?] as well as the TD-MPC2 agent without the opponent-aware dynamics. During data collection we control both policies (self and opponent), therefore a_t^{opp} is logged exactly in the replay buffer. The separate loss is given by

$$\mathcal{L}_{\text{opp}} = \|a_t^{\text{opp}} - \pi_{\text{opp}}(\vec{z}_t)\|^2, \quad (8)$$

where π_{opp} is the opponent network, trained with a frozen encoder in periodic intervals. With this MSE objective, π_{opp} is a deterministic mean predictor: it outputs a single action estimate per latent state.

During training, opponent actions are taken from the replay buffer; during inference they must be predicted by the opponent network. We use the opponent action deterministically: $\hat{a}_t^{\text{opp}} = \pi_{\text{opp}}(\vec{z}_t)$. This predicted action is fed into the dynamics to obtain the next latent state (Figure 1 b). When multiple opponent models are used (e.g., different cloned policies), each rollout can be assigned a fixed opponent model for the full horizon; diversity across rollouts then comes from varying which opponent model is used, not from sampling different actions from a stochastic π_{opp} .

2.2 SAC - Jannik Mänzer

Soft Actor-Critic (SAC) [3] is an off-policy actor-critic algorithm that extends standard reinforcement learning by optimizing a maximum entropy objective. Rather than seeking only the maximum cumulative reward, the agent aims to maximize a weighted objective of reward and the entropy of the policy

$\mathcal{H}(\pi(\cdot|\mathbf{s}_t))$:

$$J(\pi) = \sum_{t=0}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_\pi} [r(\mathbf{s}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi(\cdot|\mathbf{s}_t))] \quad (9)$$

This entropy term \mathcal{H} encourages the policy to assign non-zero probability to multiple actions where optimal, preventing premature convergence to deterministic behavior and improving exploration. The temperature parameter α controls the trade-off between the reward and the entropy.

Two soft Q-functions, Q_1 and Q_2 , are trained to estimate the expected return plus the future entropy of the policy. To mitigate overestimation, Clipped Double-Q Learning [2] is used, where the target is calculated using the minimum of two target networks $\bar{Q}_{1,2}$, whose weights are obtained via an exponential moving average of the main Q-networks [7]. The target for the Q-function update is given by:

$$y_t = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{a}_{t+1} \sim \pi} \left[\min_{j=1,2} \bar{Q}_j(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - \alpha \log \pi(\mathbf{a}_{t+1}|\mathbf{s}_{t+1}) \right] \quad (10)$$

The parameters θ are updated by minimizing the squared error between the prediction and this entropy-augmented target:

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} (Q_i(\mathbf{s}_t, \mathbf{a}_t) - y_t)^2 \right] \quad \text{for } i \in \{1, 2\} \quad (11)$$

The policy π_ϕ is updated to maximize the value estimate provided by the Q-functions while maintaining high entropy. Using the reparameterization trick $\mathbf{a}_t = f_\phi(\epsilon_t; \mathbf{s}_t)$ to allow for backpropagation, the objective is to maximize:

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}, \epsilon_t \sim \mathcal{N}} \left[\min_{j=1,2} Q_j(\mathbf{s}_t, f_\phi(\epsilon_t; \mathbf{s}_t)) - \alpha \log \pi_\phi(f_\phi(\epsilon_t; \mathbf{s}_t)|\mathbf{s}_t) \right] \quad (12)$$

2.2.1 Automatic Entropy Tuning

Instead of choosing the temperature hyperparameter α manually, it can be tuned automatically to ensure the policy satisfies a minimum target entropy constraint $\bar{\mathcal{H}}$ (typically chosen as $-\dim(\mathcal{A})$, where $\dim(\mathcal{A})$ is the dimensionality of the action space). This adapts the exploration pressure during training:

$$J(\alpha) = \mathbb{E}_{\mathbf{a}_t \sim \pi_t} [-\alpha(\log \pi_t(\mathbf{a}_t|\mathbf{s}_t) + \bar{\mathcal{H}})] \quad (13)$$

2.2.2 Pink Noise Exploration

2.2.3

3 Results

3.1 TD-MPC2 Hyperparameters and Curriculum

To determine the optimal horizon for TD-MPC2, we trained a TD-MPC2 agent with different horizons and evaluated the performance. The horizons tested were 4, 6, 8, 10, and 12. The runs were done with the following hyperparameters: learning rate 0.0003, batch size 512, network size of three layers with 256 units each, a latent dimension of 256, 5 Q-networks, a gamma of 0.99, a temperature of 0.5, a vmin of -10, a vmax of 10, a win reward bonus of 10, and a win reward discount of 0.92. The runs were done with the following curriculum: 4000 episodes of full competency, with a basic strong opponent. The results are shown in Figure 2. Additionally with the same hyperparameters, the modified opponent aware dynamics was tested, for this three internal opponent models were used: the TD-MPC2 agent without the opponent-aware dynamics, the SAC agent and the DECOYPOLICY agent, which was trained to mimick the basic strong opponent. The results are shown in the same figure, but with the label "opponent aware dynamics".

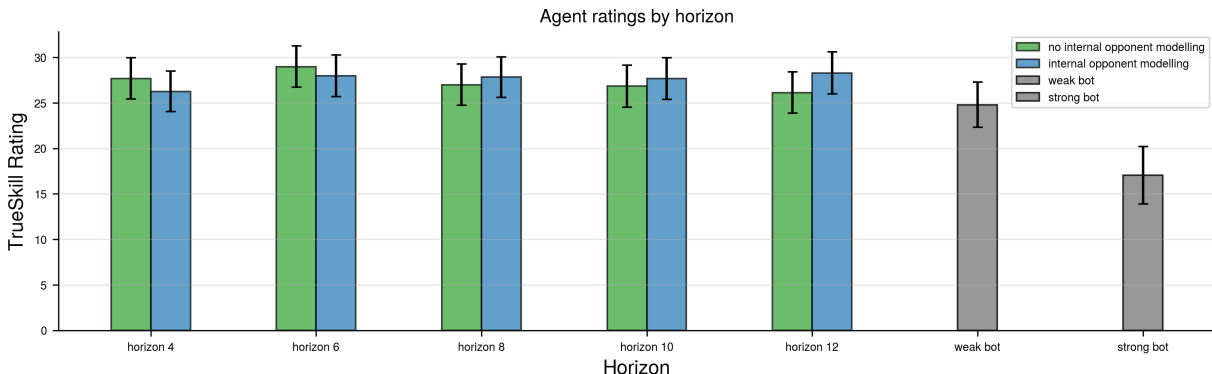


Figure 2: TrueSkill ratings of the TD-MPC2 agent with different horizons and different opponent aware dynamics.

3.1.1 Overall Results

To compare all trained agents on a common scale, we evaluated them in a round-robin tournament within the archive matchmaking system using TrueSkill [5], a Bayesian rating algorithm that updates the belief over each agent’s skill level after every match outcome. The rating reported is $\mu - 3\sigma$, a conservative lower-bound estimate that accounts for residual uncertainty. Figure 3 shows the TrueSkill ratings for a representative selection of agents: the two scripted baselines (weak and strong bot), the SAC agent, and the two best-performing TD-MPC2 variants – one trained with internal opponent modelling at planning horizon $H=8$ and one at $H=6$.

All learned agents clearly surpass both scripted baselines. The basic strong bot achieves the lowest rating (15.76), which is counterintuitive given its name but reflects the fact that its deterministic, aggressive strategy is highly exploitable by learned policies; it was not designed as a competitive opponent against gradient-based agents. The basic weak bot achieves a higher rating (25.24) because its more conservative behavior generates fewer opportunities for the opponents to score, making it harder to accumulate decisive wins against.

Among the learned agents, SAC reaches a rating of 28.34, demonstrating that a well-tuned model-free actor-critic is already a strong baseline in this environment. The TD-MPC2 agent trained with internal

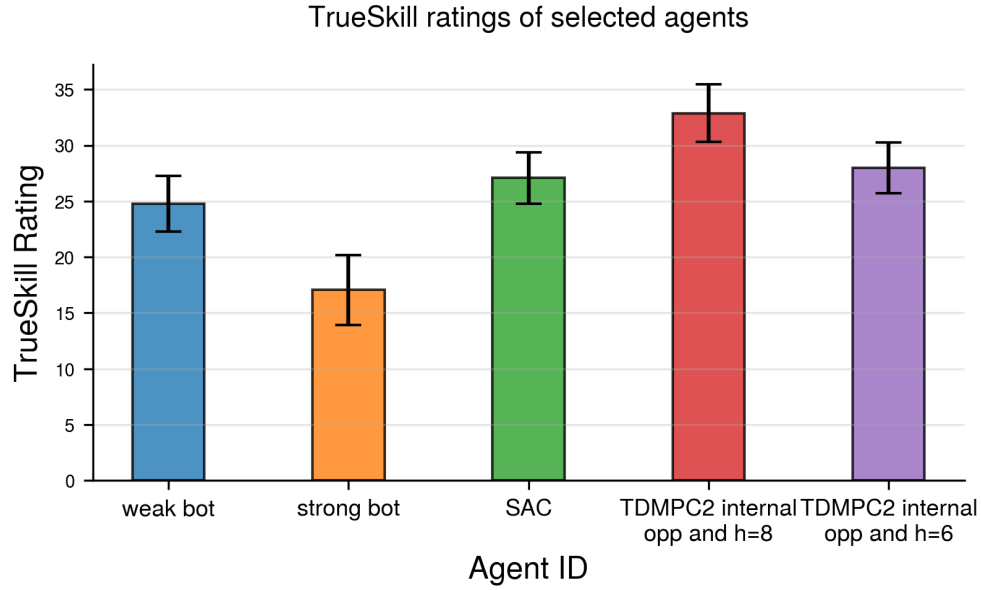


Figure 3: TrueSkill ratings (mean $\pm 3\sigma$) for a selected set of agents. Higher is better.

opponent modelling at $H=6$ achieves a comparable rating of 28.53, while the variant at $H=8$ reaches 31.41, the highest rating overall. The gain from $H=6$ to $H=8$ suggests that a longer planning horizon provides a meaningful advantage, allowing the agent to anticipate multi-step game dynamics more accurately. The overlap in confidence intervals between SAC and TD-MPC2 at $H=6$ indicates that, at this horizon, the model-based approach does not yet offer a statistically significant benefit over the model-free baseline. The TD-MPC2 variant at $H=8$ with more training steps (16 000 vs. 4 000) does, however, emerge as the clearly strongest agent.

4 Acknowledgements & Data Availability

We would like to thank the instructors and the staff of the Reinforcement Learning course for their help and support. All of our code can be found on our GitHub repository [1].

References

- [1] A. Cheung, J. Mänzer, and N. Abraham. RL-course 2025/26: Final project report. https://github.com/NiklasAbraham/RL_CheungMaenzerAbraham_Hockey, 2026.
- [2] S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic methods, 2018.
- [3] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine. Soft actor-critic algorithms and applications, 2019.
- [4] N. Hansen, H. Su, and X. Wang. Td-mpc2: Scalable, robust world models for continuous control, 2024.
- [5] R. Herbrich, T. Minka, and T. Graepel. TrueSkill™: A bayesian skill rating system. In *Advances in Neural Information Processing Systems*, volume 19. MIT Press, 2007.
- [6] G. Martius. Hockey environment. <https://github.com/martius-lab/hockey-env>, 2023.
- [7] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, Feb. 2015.
- [8] G. Williams, A. Aldrich, and E. Theodorou. Model predictive path integral control using covariance variable importance sampling, 2015.

A Appendix

A.1 Episode Logs

The episode logs of the TD-MPC2 agent with the horizon 4 without the opponent aware dynamics are shown in Figure 4. These logs were plotted from all runs periodically, and this example is representative for the other runs.

TD-MPC2 agent with the horizon 4 without the opponent aware dynamics

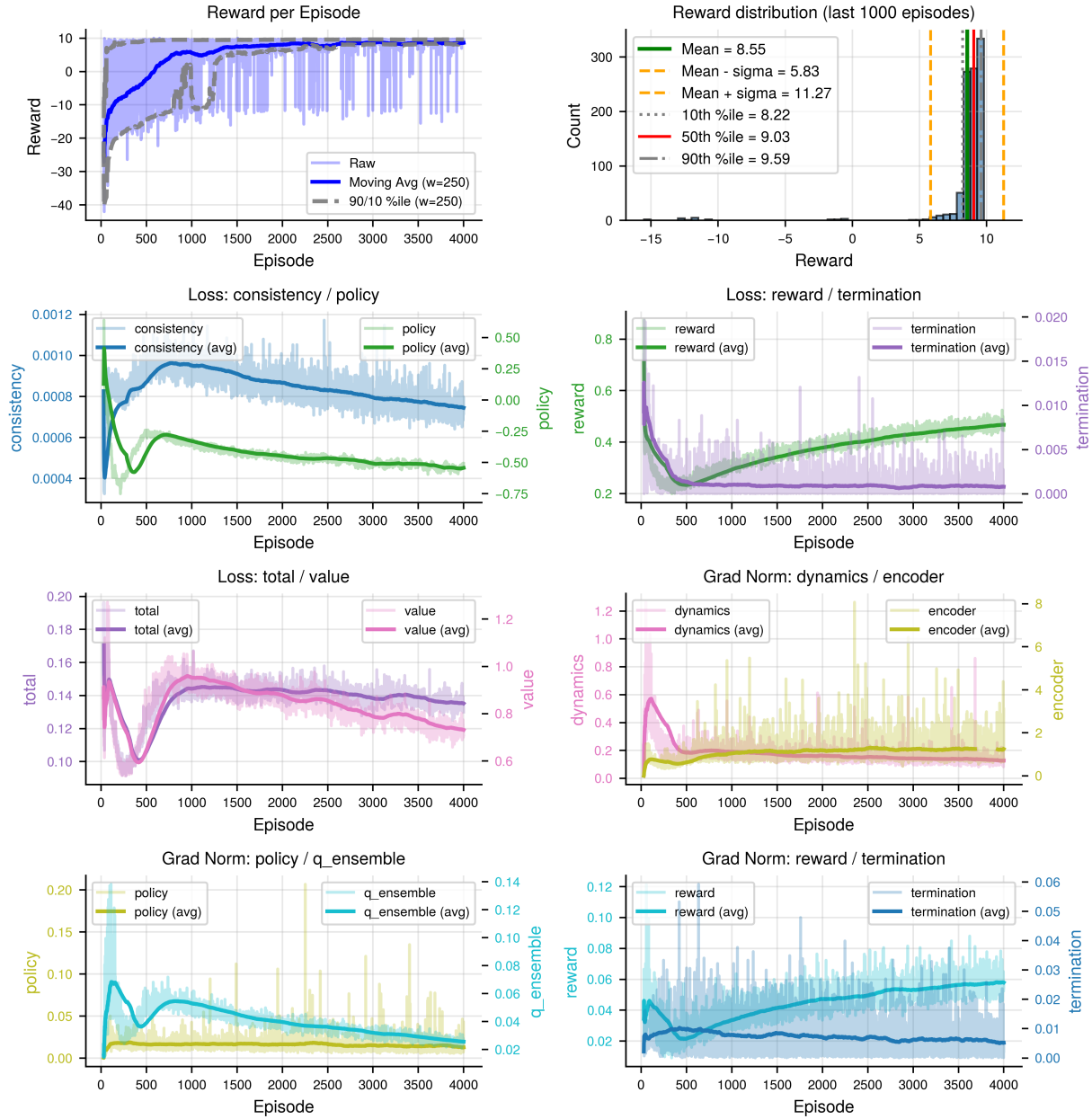


Figure 4: Episode logs of the TD-MPC2 agent with the horizon 4 without the opponent aware dynamics.