

An empirical study of the naïve REINFORCE algorithm for predictive maintenance of industrial milling machines

Rajesh Siraskar

July 20, 2023

Abstract

Industrial systems tend to be highly complex and dynamic. Reinforcement learning (RL) offers a mechanism for creating optimal predictive maintenance policies for such systems. RL is known to be an extremely complex exercise in hyper-parameter tuning. Automated machine learning (AutoML), applied to the combined fields of predictive maintenance (PdM) and RL, has yet to be studied. This article is an empirical study aimed at industrial practitioners unfamiliar with complex RL tuning. We study the effects of *untuned* RL algorithms for generating an optimal tool replacement policy for a milling machine. We compare a naïve implementation of REINFORCE against the policies of industry-grade implementations of three advanced algorithms, namely, Deep Q-Network (DQN), Advantage Actor-Critic (A2C), and Proximal Policy Optimization (PPO). Our broad goal was to study model performance under various scenarios: (1) simulated tool-wear data, (2) real tool-wear data (benchmark IEEE NUAA Ideahouse dataset), (3) univariate state with added noise levels and a random chance of break-down (4) complex multivariate state.

Performance was measured by how accurately the predictive maintenance agent suggested tool replacement compared to a deterministic preventive maintenance rule (based on the tool-wear threshold). Across 15 environment variants, REINFORCE models demonstrated a tool replacement precision of 0.687 against 0.449 for A2C, 0.418 for DQN, and 0.472 for PPO. The F1 scores were 0.609,

0.442, 0.374, and 0.345, respectively. Variability in precision and F1 was lower for REINFORCE by 0.08 and 0.016 compared to the average of the three advanced algorithms. Comparing the best *auto-selected* model, over ten rounds of training produced unusually wider gaps in performance. REINFORCE precision/F1 stood at 0.884/0.873. The best A2C, DQN, and PPO models produced 0.520/0.639, 0.651/0.740, and 0.558/0.580, respectively. While this study is a small first step toward AutoML for PdM using RL, our findings surprisingly indicate that the computationally lightweight REINFORCE performs significantly well for this particular problem.

For reproducibility, model training and testing code, data and the trained REINFORCE models have been uploaded to <https://github.com/Link>

Keywords: Predictive maintenance, milling machines, Reinforcement Learning, REINFORCE

Abbreviations

A2C	Advantage Actor-Critic	DQN	Deep Q-Network
PPO	Proximal Policy Optimization	RF	REINFORCE
SS	Single-variable state	MS	multivariate state
TP	True positive	TN	True negative
FP	False positive	FN	False negative
RL	Reinforcement Learning	SB3	Stable-Baselines3
PHM	The Prognostics and Health Management Society		

1 Introduction

Milling machines are highly versatile, ubiquitous tools serving a variety of industries. A milling machine removes metal from the work piece by rotating and driving a cutting device into it. Abrasive forces cause tool wear, and optimal tool replacement reduces direct costs and optimizes the machines' downtime. With the 2023 milling machine market valued at USD 68.3 billion ([Future Market Insights, 2023](#)), this is an important goal for the industry. The cutting tool experiences multiple types of wear as it cuts through metal. Tool wear depends in several factors such as the cutting speed, force applied to the tool, lubrication and materials of the work piece and cutting tool.

Reinforcement learning (RL) is an artificial intelligence technique inspired by nature. Fig. 1 ([Sutton and Barto, 2018](#)) shows the RL learning feedback loop. An actor or “agent” interacts with an environment and learns via “trial-and-error”. It acts based on stimuli or feedback received from the environment after performing a certain action. Actions that help in achieving the learning goal receive a reward while actions that do not, are punished. Repeating this loop over thousands of episodes, good actions are “reinforced”, thereby building a “policy” that is optimized for that goal. In the case of predictive maintenance for milling machines, the agent is the “planner” with a goal of learning an optimal tool replacement policy. The environment consists of sensors attached to the machine and related information such as job specifications, environment conditions etc.

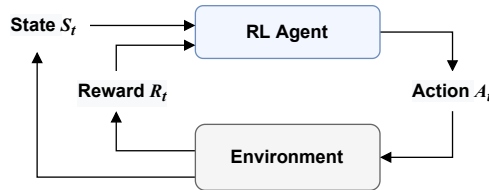


Figure 1: Reinforcement Learning

Introduced in 1992, the REINFORCE algorithm ([Williams, 1992](#)) is considered as a basic reinforcement learning algorithm. It is a policy-based, on-policy algorithm, capable of handling both discrete and continuous observation and action domains. In practice the REINFORCE algorithm is

considered as a “weak” algorithm and superseded by several algorithms developed since. Most notably the Q-Learning and its deep-neural network version, the DQN (Mnih et al., 2013), followed by Actor-Critic (Mnih et al., 2016) and one of the most robust modern-day algorithms, the PPO (Schulman et al., 2017).

While most studies on RL algorithms are evaluated on Open AI Gym environments, our experiments cover the predictive maintenance problem using a custom built environment. Also, our study focuses on the earliest RL algorithm, the REINFORCE. This is often neglected by studies, as it is considered an early and weak algorithm. In the industry the milling tool is replaced after a set threshold. We use this “deterministic preventive maintenance” policy as the baseline for comparing the RL algorithm policy.

Our systematic evaluation, based on levels of environment difficulty, different bench-mark datasets and varying noise levels allow a broader, more robust, comparison of the algorithms. Finally, we conduct statistical tests to ensure a robust statistical-evidence based conclusion. Based on experiments we show that REINFORCE works surprisingly well on tool-replacement precision and F-beta (0.5), over all the variants. The recall and F1-score are better or at-par with the advanced algorithms.

The main **contributions** of this research are:

1. Contributes to the broader goal of AutoML for PdM using RL
2. Research targeted toward the industrial practitioner not accustomed to complex hyper-parameter tuning of RL algorithms
3. Design and implement an RL environment for predictive maintenance of a milling machine.
4. Rigorous evaluation of four standard RL algorithms.
5. Use of simple performance evaluation statistical measures and plots that industrial practitioners are normally used to.

The rest of the paper is structured as follows: In the next section we survey some related work and provide the necessary technical background describing the algorithms studied in this research. Section 3 discusses implementation details of the REINFORCE algorithm and the predictive maintenance environment followed by the methodology adopted for training, testing and evaluation. Section 4 presents the results of the experiments, that are discussed in Section 5. Finally, we summarize and draw conclusions in Section 6.

2 Related work and background

2.1 Literature Review

Automated machine learning (AutoML) for predictive maintenance has been studied by only few researchers. All the 9 articles found use supervised machine learning. Larocque-Villiers et al. (2021) use signal processing to standardize data followed by AutoML to select top performing models. Hadi et al. (2023) use AutoML to accurately identify various types of faults in ball bearings, while Maurer et al. (2021) has studied use of AutoML on log-data generated by production lines for performing preventive maintenance. We did not find any work that applied AutoML for RL algorithms in the PdM context.

Significant work has been conducted in the application of RL for predictive maintenance (Erhan et al., 2021; Panzer and Bender, 2021; Siraskar et al., 2023), in general. However, we did not find¹ any work done for predictive maintenance of *milling machines* across three major journal index services². While RL is not, traditional machine learning methods *have* been applied; for example Qin et al. (2023) and Qiang et al. (2023) apply tool-wear law and physics based models before applying ML. Twardowski et al. (2023) and Denkena et al. (2023) use data gathered by sensors. While Twardowski

¹Search terms: "reinforcement learning AND tool wear AND maintenance". As of: 10-Jul-2023.

²IEEE Xplore™, Scopus™ and Web Of Science™

[et al. \(2023\)](#) use two different classification trees, [Denkena et al. \(2023\)](#) use data recorded on other similar machines for building ML models. [Oshida et al. \(2023\)](#) proposes real-time tool wear detection using a stacked LSTM³ encoder-decoder model for anomaly detection as a mechanism to address predictive maintenance.

Research work focusing of comparing experimental and empirical analysis of various RL algorithms, has been limited to using standard benchmark OpenAI Gym environments. [Sandeep Varma et al. \(2022\)](#) documents experimental evaluation of four policy-gradient and actor-critic algorithms PPO, SAC, DDPG and A2C using OpenAI Gym environments such as the Pendulum, Mountain Car, Bipedal Walker, Lunar Landing and Atari 2600 game environments. [Velivela and Yarram \(2020\)](#) evaluate DQN, DoubleDQN, A2C, REINFORCE and PPO using Cartpole, Space Invaders and the Lunar Lander. [Dulac-Arnold et al. \(2020, 2021\)](#) are significant contributions toward analyzing empirical studies directed toward *real-world* challenges. They apply real-world design concepts on the Cartpole and other complex environments such as humanoid and walkers from the Real-World Reinforcement Learning (RWRL) Suite⁴. These environments are then evaluated for multiple RL algorithms such as, REINFORCE, Trust Region Policy Optimization (TRPO) and Deep Deterministic Policy Gradient (DDPG).

[Dulac-Arnold et al. \(2021\)](#); [Henderson et al. \(2018\)](#) tackle RL for continuous control. [Henderson et al. \(2018\)](#) evaluate DDPG, ACKTR, TRPO and PPO on complex MuJoCo⁵ environments such as the HalfCheetah, Hopper, Walker and Swimmer. Similar to our research, they used the OpenAI baseline implementations of RL algorithms for the experiments and evaluation. [Ford and Ritchie \(2022\)](#) is experimental evaluation we found that was based on real-world application. They compare DQN, A2C and PPO for choosing the operational radio frequency (RF) mode for a multi-function RF system and go on to suggest that PPO is the best.

The survey shows that most existing work use standard OpenAI Gym

³long short-term memory networks

⁴Link to RWRL » [link](#)

⁵[MuJoCo](#) provides environments for studying Multi-Joint dynamics with Contact

environments, which although necessary for bench marking performance, do not provide coverage of industrial predictive maintenance. In Section 3.1 we attempt to bridge this gap by implementing a custom built environment.

2.2 Technical Background

Key concepts of RL

A task is a goal we set for our agent to learn. In our case the agent must learn to optimally predict the replacement of the tool. Frequent tool replacement increases down-time while delaying it results in inferior work piece quality. In Fig. 1 the agent interacts with the environment by performing an action ($a \in \mathcal{A}$), which then alters the state of the environment to one of many states ($s \in \mathcal{S}$). The resulting state is determined by a state-transition probabilities (\mathcal{P}) since RL is founded on Markov Decision Process (MDP) theory. The new state provides a feedback via a reward ($r \in \mathcal{R}$). Higher positive rewards “reinforce” good behavior. Performing this over thousands of episodes with the objective of maximizing the total rewards R , enables the agent to develop a policy π which is essentially a mapping of the optimal action to perform, given a certain state.

A **value function** computes how good a state or an action is by predicting future rewards, also known as a “return” $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$. $\gamma \in [0, 1]$ facilitates discounting i.e. applying less weight to future rewards.

Value functions can be represented by **state-value** $V_{\pi}(s)$ of a state s as the expected return: $V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$; or an **action-value** function of a state-action pair as $Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]$.

With a brief overview of RL, we now briefly touch upon the core ideas of the four algorithms we experimented with.

Deep Q-Network (DQN)

Deep Q-Network (Mnih et al., 2013) significantly improved the earliest RL algorithm, Q-learning, by introducing neural networks to learn policies for

high-dimension environments with two novel strategies to significantly stabilize learning – an “experience replay buffer” and a target network that was frozen and only periodically updated. Equation (1) shows the DQN loss function where D is the replay memory and is sampled using a uniform distribution $U(D)$, $Q(s, a; \theta)$ is the function parameterized with θ , that helps compute the Q values and θ^- represents parameters of the frozen target Q-network.

$$\mathcal{L}(\theta) = \mathbb{E}_{(s,a,r,s') \sim U(D)} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta) \right)^2 \right] \quad (1)$$

Advantage Actor Critic (A2C)

A2C is a variant of Asynchronous Advantage Actor Critic (A3C) (Mnih et al., 2016), and uses multiple computation workers to avoid the use of a replay buffer. A2C is a policy-gradient actor-critic algorithm. The policy-gradient family of algorithms strive to model and optimize the policy directly. Actor-critic structures consist of two networks – a critic that updates function parameters w of the value function (i.e either $Q_w(a|s)$ or $V_w(s)$); and an actor that updates the policy parameters θ for $\pi_\theta(a|s)$, following the direction computed by critic. Actors therefore learn the parameterized policy π_θ using the policy-gradient as shown in (2).

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_t \left[A_t^\pi \nabla_\theta \ln \pi_\theta(a_t|s_t) \right] \quad (2)$$

Where the advantage function $A_t^\pi(s_t, a_t)$ measures how good or bad the action is w.r.t. policy’s average, for a particular state, using (3).

$$A_t^\pi(s_t, a_t) = Q^\pi(s_t, a_t) - V^\pi(s_t) \quad (3)$$

Proximal Policy Optimization (PPO)

Schulman et al. (2017) formulated PPO which is often considered as the most robust of the RL algorithms. PPO is policy-gradient method based on TRPO (Trust region policy optimization) by Schulman et al. (2015), where the main idea is the use of a trust region to improve training stability by avoiding

updates to parameters that vastly change the policy at a particular time step. TRPO ensures this by using a divergence constraint on the magnitude of policy update. If $r(\theta)$ (4) represents the ratio of probabilities between policies of previous and current iteration, then the objective function of TRPO is given by (5), where \hat{A} represents the estimated advantage function.

$$r(\theta) = \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} \quad (4)$$

$$J^{\text{TRPO}}(\theta) = \mathbb{E}[r(\theta)\hat{A}_{\theta_{\text{old}}}(s, a)] \quad (5)$$

PPO extends TRPO by additionally imposing a regional constraint. It prevents large updates by forcing the ratio $r(\theta)$ to stay within a small interval $[1 - \epsilon, 1 + \epsilon]$, around 1.0, by use of a hyper-parameter ϵ .

$$J^{\text{CLIP}}(\theta) = \mathbb{E}_t [\min(r_t(\theta)A_t^{\pi_{\theta_{\text{old}}}}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t^{\pi_{\theta_{\text{old}}}})] \quad (6)$$

REINFORCE

The REINFORCE algorithm, invented by Williams (1992), is an early algorithm that directly learns a policy π_{θ} to produce action probabilities from states. Actions that cause favorable states are positively reinforced thereby increasing their probability of occurrence. Conversely, those resulting in unfavorable states are penalized.

The objective function (7) that the agent attempts to maximize is defined as the expected return over many trajectories sampled from the policy.

$$\max_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] \quad (7)$$

REINFORCE uses policy gradient to update the action probabilities.

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_t [R_t(\tau) \nabla_{\theta} \ln \pi_{\theta}(a_t|s_t)] \quad (8)$$

Stable-Baselines3

Stable-Baselines3 (SB3) [Raffin et al. \(2021\)](#) provides robust open source implementations of many reinforcement learning algorithms. The core algorithms of Stable-Baselines3 are completely refactored PyTorch implementations of Open AI baselines ([Dhariwal et al., 2017](#)). [Huang \(2022\)](#) describes the 37 PPO implementation details that SB3 carried out and demonstrates the complexity depth of their implementations.

As of 10-Jul-2023, 13 algorithms have been implemented ([SB3-Algorithms, 2022](#)). REINFORCE has *not* been implemented. For this research we use the *default* SB3 implementations of DQN, A2C and PPO and compare its performance to a naïve custom implementation of the REINFORCE algorithm.

3 Methodology

3.1 Implementation details

Reinforcement learning requires an environment to function. In this section we describe the custom milling environment we built, which allows our agent to learn a policy for tool replacement. Figure 2 shows the stepped approach used to generate the *fifteen* environments, of varying complexity, to evaluate the agents. In this section we describe the approach used to model and simulate tool-wear and hence generate wear data followed by actual tool wear data. We then describe the elements of environment design i.e. state, reward functions. The Python code used for this research is uploaded to GitHub at [URL](#).

Simulating tool-wear

[Dašić \(2006\)](#) provide a parameterized power-exponential function for modeling tool wear (9), where VB represents the flank wear in mm.

$$VB = a \cdot t^{b_1} \cdot e^{b_2 \cdot t} \Big|_{t=t_0}^{t=t_1} \quad (9)$$

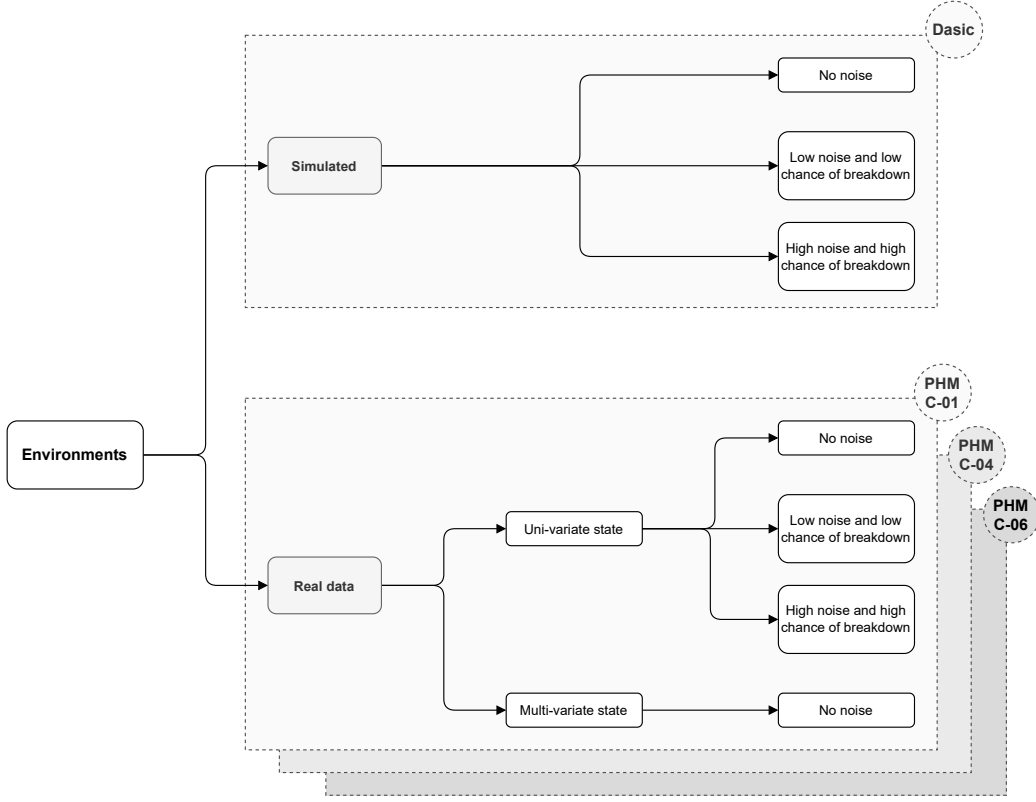


Figure 2: The fifteen different environments used for evaluation

We used the parameter values provided in the paper $a = 0.08257, b1 = 0.3342, b2 = 0.03147$; and simulated 121 data points. Fig. 3a shows the tool wear simulated using (9) with the red dotted line indicating the wear threshold beyond which tool is replaced. This provided the mechanism to simulate the tool wear based **univariate state**. We then add two further levels of increased complexity using noise and a chance of breakdown. This gives us three distinct environments.

Actual tool wear data

The IEEE DataPort hosts the tool wear data obtained from a high-speed CNC milling machine, (Yingguang et al., 2021). C01, C04 and C06 datasets are suggested as benchmark for machine learning training and were the ones we used. The data consists of seven columns corresponding to dynamometer

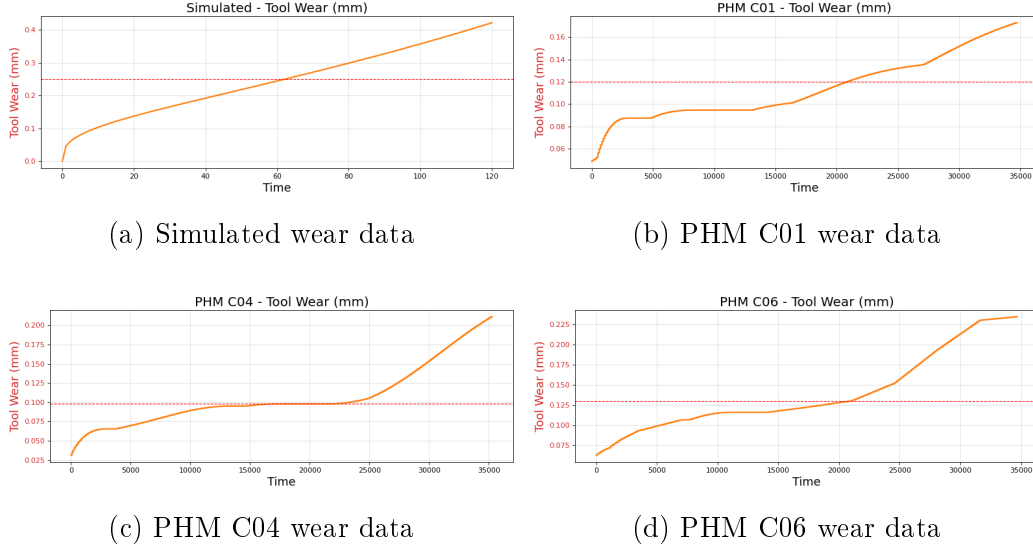


Figure 3: Tool wear data. Red dotted line indicates the wear threshold beyond which tool is replaced.

measuring force (N) in X, Y and Z dimensions; accelerometer measuring vibration (g) in X, Y and Z dimensions and finally acoustic emission data as AE-RMS (V). A separate file contains tool wear data in mm. Figures 3b, 3c and 3d show the tool wear for the three datasets. We use the real data to create two state designs, an univariate state consisting of only the tool wear and a **multivariate** state designed using all the additional seven sensor values. As for the simulated case, the complexity of the *univariate* state is increased using two levels of noise and break-down parameters. The multivariate state is complex in itself and we use the natural (without noise) form.

In RL state variables are often normalized to improve stability and convergence. Both the simulated and real data was normalized using min-max scaling such that the tool wear and other state features, $x \in [0, 1] \subset \mathbb{R}$. We will see next how this allows adding white noise of similar magnitudes across different PHM datasets.

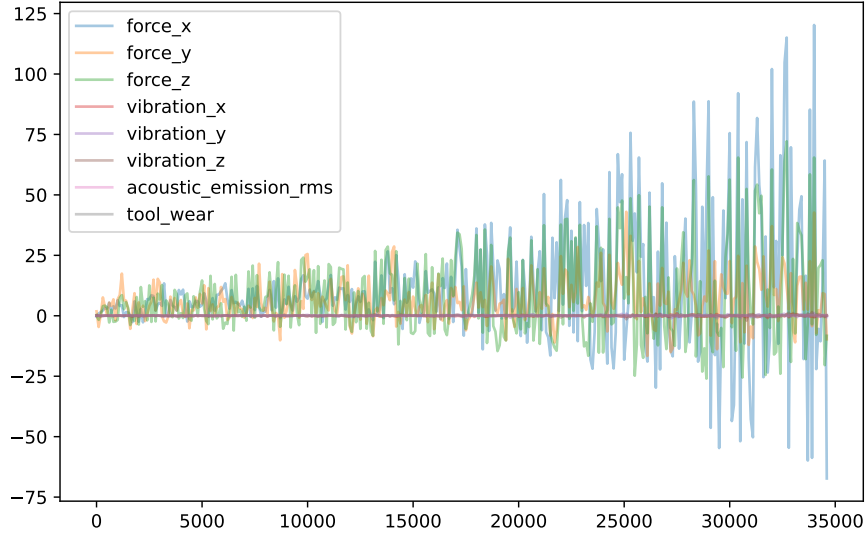
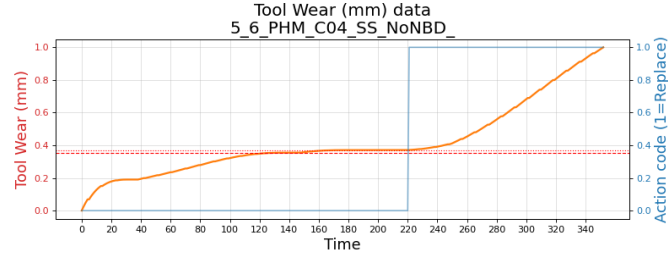


Figure 4: PHM C06: multivariate data

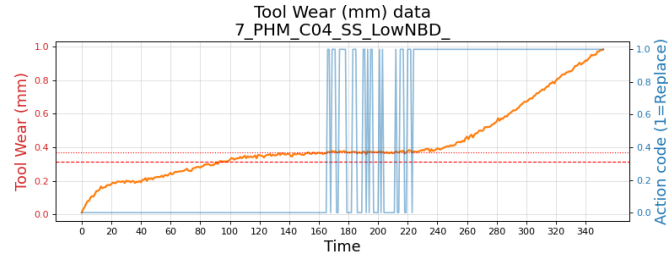
Adding noise and chance of break-down

Fig. (5) shows the effect of adding two levels of noise. “Low noise” is obtained by adding Gaussian noise with an order of magnitude of -3 i.e. between $[0.0, 0.001]$ and “high noise” is of order -2 i.e. between $[0.0, 0.01]$. Since the tool wear data is normalized this adds significant perturbations as seen in Fig. (5b) and (5c). The noise affects the tool replacement decision (solid blue line) around the replacement threshold (dotted red line). The *human* preventive maintenance policy replaces the tool if the wear exceeds the threshold and this decision boundary oscillates due to the noise. One can see that in the case of no noise (Fig. 5a), the decision boundary is clean.

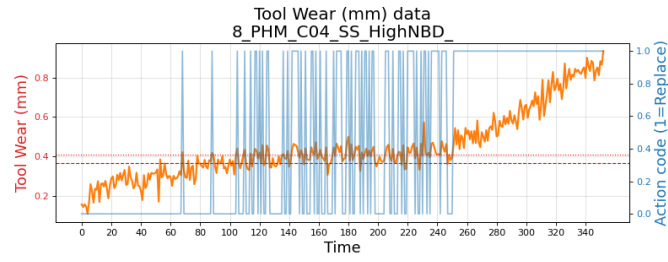
Break down occurs due to excessive tool use and can often occur randomly. In conjunction with Gaussian noise this complexity is added for the univariate state based environments. For the low-noise variant we add a “low chance”, 5%, of break down and for the high noise variant we add a higher chance of 10%. Sampled from a uniform distribution, the probability of break down is checked against this “chance” threshold and the episode is



(a) No noise



(b) Low noise



(c) High noise

Figure 5: PHM C04 tool wear data and the effect of noise.

terminated if less than. This is seen in line 9 of the code in Listing 1.

Table 1 summarizes the 15 environment variants and their three logical groups: (1) Simulated 1-3 (2) Real data – simple univariate environment (4-12) and Real data – simple multivariate (13-15).

Actions

The choice of action is binary and represented by $ACTION \in [0, 1]$, where 0 represents continuation of milling operation (CONTINUE), while 1 represents the action of replacing the tool (REPLACE_TOOL).

	Environment variant	Noise factor	Breakdown chance
Simulated			
1	Simulated - No noise	None	None
2	Simulated - Low noise	1e-3	0.05
3	Simulated - High noise	1e-2	0.10
Real data – simple univariate			
4	PHM C01 SS (simple, univariate) - No noise	None	None
5	PHM C01 SS (simple, univariate) - Low noise	1e-3	0.05
6	PHM C01 SS (simple, univariate) - High noise	1e-2	0.10
7	PHM C04 SS (simple, univariate) - No noise	None	None
8	PHM C04 SS (simple, univariate) - Low noise	1e-3	0.05
9	PHM C04 SS (simple, univariate) - High noise	1e-2	0.10
10	PHM C06 SS (simple, univariate) - No noise	None	None
11	PHM C06 SS (simple, univariate) - Low noise	1e-3	0.05
12	PHM C06 SS (simple, univariate) - High noise	1e-2	0.10
Real data – complex multivariate			
13	PHM C01 MS (complex, multivariate) - No noise	None	None
14	PHM C04 MS (complex, multivariate) - Low noise	1e-3	0.05
15	PHM C06 MS (complex, multivariate) - High noise	1e-2	0.10

Table 1: List of the fifteen environments and their categorization

Reward function

The reward function (10) drives the learning process. Code implementation is seen in Listing 1. There are three important phases where reward (or penalty) is provided as feedback. R_1 , R_2 and R_3 are constants that determine the magnitude of reward. The first condition is when the current state of wear $w(s_t)$ is less than the threshold W_T , and a desirable condition, hence the reward is positive. The formulation allows higher reward to be collected the closer it is to threshold, but not allowing it to cross it. If it does, the agent is

penalized (negative reward) by a magnitude of R_2 and once again the farther away it is from the threshold i.e. a highly deteriorated tool the larger the penalty. To avoid this “bad” state, the agent must learn to replace the tool; represented by the third condition. Tool replacement implies a direct cost (that of the tool) and a much higher and significant downtime “cost”. To ensure the agent learns not to replace unnecessarily, we “penalize” it. It is important to note that the last condition in (10) is an “incremental addition”, the first two conditions are mutually exclusive and evaluated first, followed by the last condition which is incrementally added on whatever reward is collected in the previous condition. The agent then tries to balance these three condition, such that it maximizes its total return, over time.

$$R_t = \begin{cases} +R_1 \times t, & \text{if } w(s_t) < W_T \\ -R_2 \times t, & \text{if } w(s_t) \geq W_T \\ += -R_3, & \text{if ACTION = REPLACE_TOOL} \end{cases} \quad (10)$$

Environment feedback

Figure 1 shows “feedback” that is important for the agent to learn from the action it takes. This is performed by the `step()` function in the environment code and the main parts are shown in Listing 1. At every time step an action is taken, the conditions (e.g. lines 2 and 8) of the resulting state evaluated for assigning a reward and then either proceeding or terminating the episode (e.g. lines 9 and 11).

```

1 # Termination condition
2 if self.ep_length >= self.max_operations:
3     terminated = True
4     self.reward = 0.0
5     self.df_index = -1
6     info = {'termination': 'Max. milling operations crossed'}
7
8 elif (tool_wear > self.wear_threshold and
9       np.random.uniform() < self.breakdown_chance):
10     terminated = True
11     self.reward = 0.0

```



```

12     self.df_index = -1
13     info = {'termination': 'Tool breakdown'}
14
15 else:
16     terminated = False
17     info = {'action': 'Continue'}
18     reward = 0.0
19
20     if tool_wear < self.wear_threshold:
21         reward += self.R1*self.df_index
22     else:
23         # Threshold breached. Farther away from threshold => more
24         # penalty
25         reward += self.R2*self.df_index
26
27     # Based on the action = 1 replace the tool or if 0,
28     # continue with normal operation
29     if action:
30         reward += self.R3
31         self.df_index = -1 # Tool replaced - so roll back tool
32         # life
33         self.ep_tool_replaced += 1 # Increment tool replaced
34         # count
35         info = {'action': 'Tool replaced'}

```

Listing 1: Environment: Important implementation details of the 'step' function

Network architecture and basic hyper-parameters

For this research we implemented a naïve REINFORCE algorithm, with very simple network architectures. This was compared against *default* implementations of the more advanced algorithms; DQN, A2C and PPO, from the Stable-Baselines3 library. Table 2 shows a comparison of the network architecture and the basic common hyper-parameters.

The REINFORCE uses a *single* internal layer of 64 units. The default network architecture for all three SB3 algorithms (A2C/DQN/PPO) consists of *two* fully connected layers with 64 units per layer ([SB3-Default Network](#)

	A2C	DQN	PPO	REINFORCE
Network architecture	input dim x [64 Tanh x 64 Tanh] x output dim	input dim x [64 Tanh x 64 Tanh] x output dim	input dim x [64 Tanh x 64 Tanh] x output dim	input dim x [64 ReLU] x output dim
Layers	2	2	2	1
Units	64 x 64	64 x 64	64 x 64	64
Activation	Tanh, Tanh	Tanh, Tanh	Tanh, Tanh	ReLU
Optimizer	RMSprop	Adam	Adam	Adam
Learning rate	0.0007	0.0001	0.0003	0.01
Gamma	0.99	0.99	0.99	0.99

Table 2: Comparing the network architecture and basic hyper-parameters across algorithms

[Architecture, 2022](#)). While we used ReLU (rectified linear unit) as the activation function, the other three algorithms used hyperbolic tangent (Tanh).

3.2 Training

The training strategy must ensure a uniform comparison of the algorithms. We maintained the exact *same* environment variant, wear dataset, noise parameter, probability of breakdown parameter, and the reward function parameters R_1 , R_2 and R_3 ; across all four algorithms during a single training round. Our code is fairly automated and we used a configuration file to configure the training experiments and record the results. Fig. 6 shows the main columns of the configuration file we used.

As the wear data is time-series data, the training and test sets are created by systematic sampling. Simulated data and real tool-wear data (PHM) was randomly sampled at a certain frequency and down sampled into training

and test sets.

expt_n	environment_ info	R1	R2	R3	wear_ threshold	add_noise	breakdown_ chance	episodes	milling_ operations_ max	test_ cases	test_ rounds
0	Simulated NBD	1	-1	-40	0.250	0	0	800	121	40	10
1	Simulated LBD	0.5	-1	-40	0.250	1000	0.05	800	121	40	10
2	Simulated HBD	0.5	-1	-40	0.250	100	0.1	800	121	40	10
3	PHM C01 SS NBD	1	-1	-40	0.120	0	0	800	348	40	10
4	PHM C01 SS LBD	1	-1	-40	0.120	1000	0.05	800	348	40	10
5	PHM C01 SS HBD	1	-1	-40	0.120	100	0.1	800	348	40	10
6	PHM C04 SS NBD	1	-1	-40	0.098	0	0	800	348	40	10
7	PHM C04 SS LBD	1	-1	-40	0.098	1000	0.05	800	348	40	10
8	PHM C04 SS HBD	1	-1	-40	0.098	100	0.1	800	348	40	10
9	PHM C06 SS NBD	1	-1	-40	0.130	0	0	800	348	40	10
10	PHM C06 SS LBD	1	-1	-40	0.130	1000	0.05	800	348	40	10
11	PHM C06 SS HBD	1	-1	-40	0.130	100	0.1	800	348	40	10
12	PHM C01 MS NBD	0.5	-1	-40	0.120	0	0	1000	348	40	10
13	PHM C04 MS NBD	0.5	-1	-40	0.098	0	0	1000	348	40	10
14	PHM C06 MS NBD	0.5	-1	-40	0.130	0	0	1000	348	40	10

Figure 6: Configuring the experiments

The REINFORCE was trained for 800 episodes for the simulated and PHM univariate variants, for all three noise and breakdown levels (none, low and high) – Table 1 items 1-12. For the PHM multivariate variant Table 1 items 13-15, REINFORCE was trained for 1000 episodes. SB3 algorithms were trained for 10,000 episodes for all variants. We ran ten rounds of training, tested each model generated and then averaged their results. Testing is explained in the next section while results are presented in Section 4.

Figure 7 shows the training plots for the algorithm of our interest – REINFORCE. It displays how the wear plot looked for C04 with low noise and low chance of breakdown settings. The average rewards increase over the course of 800 episodes (Fig. 7b), the episode length demonstrates the complexity introduced by random breakdown (which abruptly terminates the episode). It is the tool replacement policy that is essentially of interest to the industrial practitioner, this decreases to optimal levels as the agent learns over time, as seen in Fig. 7d. Similarly, Fig. 8 demonstrates the training for the PHM C06 dataset affected by high noise and higher breakdown probability. Finally, for the more complex multivariate state variant, the training plots are as seen

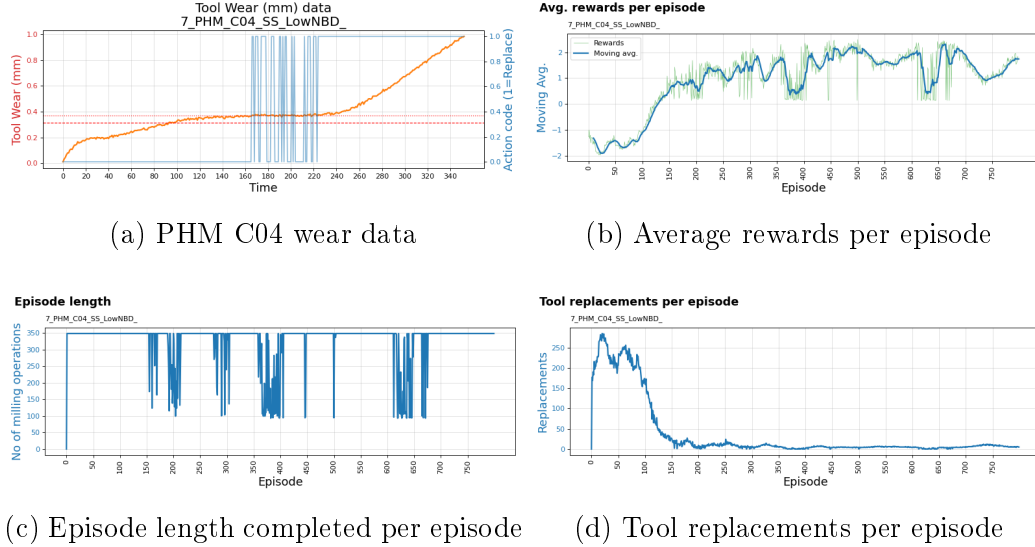


Figure 7: Training plots. Algorithm: REINFORCE. Dataset: PHM-C04. Variant: Univariate state, low-noise and low chance of breakdown.

in Fig. 9; as we do not introduce noise or breakdown here, the episodes are always completed (Fig. 9c).

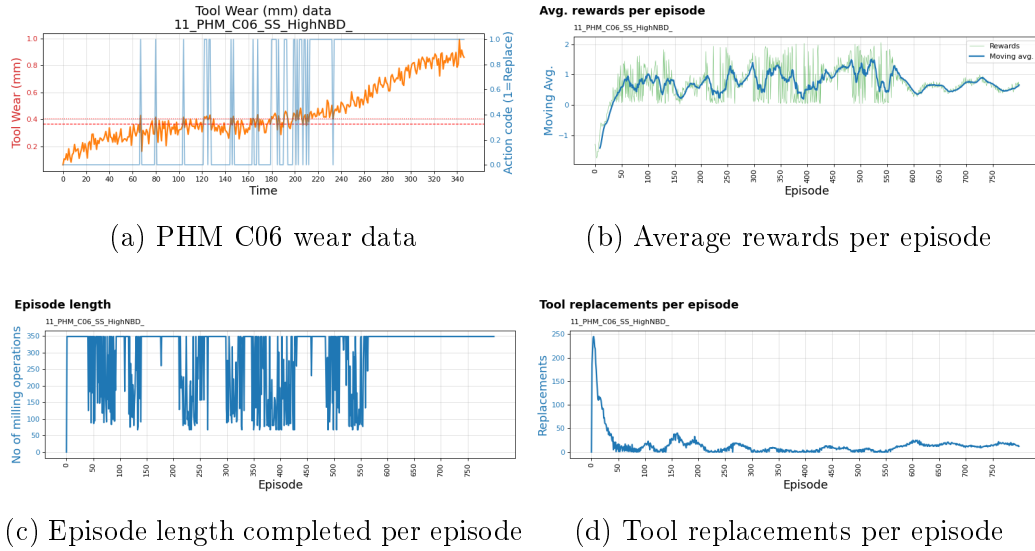


Figure 8: Training plots. Algorithm: REINFORCE. Dataset: PHM-C06. Variant: Univariate state, high-noise and high chance of breakdown.

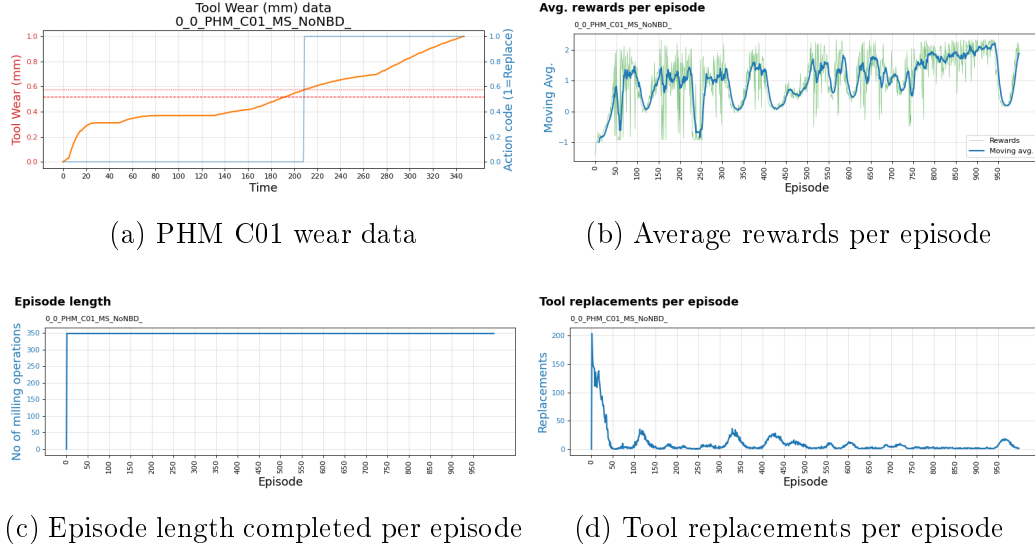


Figure 9: Training plots. Algorithm: REINFORCE. Dataset: PHM-C01. Variant: Multivariate state, No noise or breakdown.

3.3 Testing and performance evaluation

Testing was performed with data separate from the training data. 10 rounds of testing are performed, with a *new* set of 40 test cases randomly sampled and frozen across all four algorithms.

Evaluation metrics

In our data, we indicate a tool replacement suggested by the human expert as 1, and a normal operation as 0. The human decision is based on “preventive maintenance” against a set threshold. In reality the threshold is based on several factors like materials of tool and work-piece, the duration of continuous operation, ambient conditions, “cost” of production downtime etc. Threshold could therefore vary significantly from one case to another.

We applied classification metrics to evaluate the RL agent decisions that suggest similar actions. These are simple to understand for the general industrial practitioners. Classification metrics are based on the confusion matrix shown in Fig. 10. TP represents true positive cases, where both the agent and human agree on replacing the tool. False positive FP cases denote the

		Human action	
		Replace tool	Continue milling
Agent action	Replace tool	TP	FP
	Continue milling	FN	TN

Figure 10: Confusion matrix

agent falsely suggesting replacements. With false negative FN , the agent suggests to continue milling operation when a tool replacement would have helped. Precision (Pr), Recall (Rc) and F1-score metrics can then be computed as shown in (11).

$$Pr = \frac{TP}{TP + FP}, \quad Rc = \frac{TP}{TP + FN}, \quad F1\text{-score} = 2 \times \frac{Pr \times Rc}{(Pr + Rc)} \quad (11)$$

Tool replacement precision: Timely replacements, TP , ensures work piece quality. While we desire high TP s, we do want to simultaneously drive down *unnecessary* replacements and reduce costly production downtime, we therefore want lower FP s. The precision metric is therefore an ideal metric.

While we prefer a high precision, we do want a reasonably high recall i.e. do not want to miss replacement opportunities (low FN s). The F1-score (*true* harmonic mean of precision and recall) gives us a *balanced* measure. Equation (12) provides a *weighted* mechanism to provide a higher F-score for higher precision, by setting $\beta < 1.0$. For our evaluation we set β to 0.5.

$$F_\beta = (1 + \beta^2) \cdot \frac{Pr \times Rc}{(\beta^2 \cdot Pr + Rc)} \quad (12)$$

4 Results

In the previous section we discussed the training and testing methods, as well as the evaluation metrics. Our automated code reads settings shown in

Fig. 6 and creates 10 models. Each model is tested 10 times, with a fresh set of 40 randomly sampled points. The results are presented in this section.

We first present averaged metrics and averaged standard deviation of *all* 15 variants in Table 3, followed by category wise summary of results:

1. Overall performance: Averaged over all 15 environments – Table 1, items 1 to 15.
2. Simulated environment: Averaged over 3 noise setting variants – Table 1, items 1 to 3.
3. Real data - simple univariate state: Averaged over 9 variants i.e. 3 PHM datasets (C01, C04 and C06) with 3 noise settings – Table 1, items 4 to 12.
4. Real data - complex multivariate state: Averaged over 3 variants of PHM datasets (C01, C04 and C06) – Table 1, items 13 to 15.

Environment	REINFORCE				A2C				DQN				PPO			
	Prec.	Recall	F1	$F_{\beta}0.5$	Prec.	Recall	F1	$F_{\beta}0.5$	Prec.	Recall	F1	$F_{\beta}0.5$	Prec.	Recall	F1	$F_{\beta}0.5$
Simulated - No noise	0.842	0.878	0.838	0.834	0.424	0.451	0.423	0.421	0.426	0.674	0.471	0.410	0.504	0.200	0.271	0.360
Simulated - Low noise	0.777	0.929	0.834	0.796	0.465	0.423	0.409	0.427	0.421	0.338	0.270	0.283	0.482	0.236	0.296	0.369
Simulated - High noise	0.798	0.940	0.851	0.816	0.358	0.281	0.256	0.272	0.447	0.519	0.380	0.360	0.514	0.207	0.286	0.382
PHM C01 SS - No noise	0.478	0.363	0.400	0.439	0.501	0.500	0.493	0.496	0.472	0.807	0.568	0.490	0.440	0.417	0.387	0.395
PHM C01 SS - Low noise	0.507	0.311	0.332	0.383	0.503	0.598	0.535	0.513	0.393	0.502	0.351	0.317	0.522	0.338	0.388	0.448
PHM C01 SS - High noise	0.693	0.562	0.579	0.623	0.266	0.282	0.267	0.262	0.458	0.525	0.400	0.384	0.456	0.369	0.372	0.400
PHM C04 SS - No noise	0.751	0.878	0.784	0.757	0.487	0.442	0.449	0.463	0.439	0.684	0.472	0.411	0.500	0.510	0.469	0.473
PHM C04 SS - Low noise	0.662	0.756	0.672	0.657	0.409	0.455	0.428	0.416	0.411	0.500	0.370	0.341	0.488	0.280	0.324	0.386
PHM C04 SS - High noise	0.611	0.713	0.620	0.598	0.518	0.607	0.552	0.530	0.358	0.451	0.325	0.294	0.428	0.262	0.286	0.333
PHM C06 SS - No noise	0.830	0.726	0.754	0.792	0.517	0.509	0.507	0.511	0.360	0.309	0.256	0.258	0.409	0.248	0.275	0.321
PHM C06 SS - Low noise	0.205	0.279	0.228	0.212	0.510	0.577	0.530	0.516	0.434	0.266	0.266	0.296	0.417	0.181	0.232	0.294
PHM C06 SS - High noise	0.709	0.843	0.759	0.726	0.316	0.324	0.311	0.308	0.449	0.518	0.400	0.375	0.388	0.222	0.265	0.317
PHM C01 MS - No noise	0.835	0.652	0.656	0.716	0.461	0.444	0.397	0.404	0.384	0.558	0.393	0.348	0.513	0.383	0.416	0.460
PHM C04 MS - No noise	0.739	0.255	0.359	0.494	0.498	0.589	0.490	0.470	0.323	0.209	0.160	0.168	0.499	0.393	0.421	0.457
PHM C06 MS - No noise	0.864	0.356	0.469	0.616	0.501	0.713	0.578	0.527	0.489	0.705	0.529	0.479	0.523	0.488	0.485	0.498

Table 3: Model performance comparison all variants of the environments, over 10 rounds of training.

4.1 Overall summary performance

	Precision		Recall		F1-score		F-beta (0.5)	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD
A2C	0.449	0.088	0.480	0.084	0.442	0.070	0.436	0.071
DQN	0.418	0.185	0.504	0.032	0.374	0.035	0.348	0.058
PPO	0.472	0.144	0.316	0.087	0.345	0.091	0.393	0.105
REINFORCE	0.687	0.059	0.629	0.051	0.609	0.050	0.631	0.052

Table 4: Model performance summary - averaged over all environments.

At an overall level, Table 4, the REINFORCE scores much higher than the other algorithms, on all the four metrics. Tool replacement precision at 0.687, is highest of the four algorithms and better by 0.215 in absolute terms when compared to the next best, PPO. The standard deviation for precision is the lowest at 0.059. On recall, F1 and F-beta (0.5), REINFORCE is better by 0.125, 0.168 and 0.195 compared to the next best.

4.2 Simulated environment

	Precision		Recall		F1-score		F-beta (0.5)	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD
A2C	0.416	0.120	0.385	0.073	0.363	0.072	0.373	0.082
DQN	0.432	0.184	0.510	0.031	0.374	0.034	0.351	0.056
PPO	0.500	0.178	0.215	0.081	0.285	0.099	0.370	0.122
REINFORCE	0.806	0.040	0.915	0.038	0.841	0.035	0.816	0.037

Table 5: Model performance summary - averaged over simulated environments.

The simulated tool wear environment is relatively the simplest for the agent to learn. We added two levels of noise and chance of breakdown and

averaged the performance over the three variants, Table 5. The REINFORCE performs the best and in absolute terms it is better than the next best advanced algorithm by very high margins: precision by 0.306, recall by 0.405, F1 by 0.468 and F-beta (0.5) by 0.442. The standard deviation is lower or just marginally higher than others.

4.3 Real data - simple univariate environment

	Precision		Recall		F1-score		F-beta (0.5)	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD
A2C	0.447	0.077	0.477	0.091	0.452	0.072	0.446	0.070
DQN	0.419	0.179	0.507	0.032	0.379	0.036	0.352	0.057
PPO	0.450	0.146	0.314	0.082	0.333	0.087	0.374	0.102
REINFORCE	0.605	0.046	0.603	0.046	0.570	0.041	0.576	0.040

Table 6: Model performance summary - averaged over PHM-2010 environments with simple single-variable environment.

Real data offers a more challenging environment. Despite this, in Table 6, we notice that REINFORCE still performs better than the other algorithms. The margins are understandably lower: precision by 0.155, recall by 0.097, F1 by 0.117 and F-beta (0.5) by 0.130.

4.4 Real data - complex multivariate state

This environments offers the highest difficulty and the most realistic scenario. We use real data, from multiple sensors. As with other scenarios, we used untuned, default settings of all algorithms. Surprisingly, the tool replacement precision was very high at 0.813, Table .

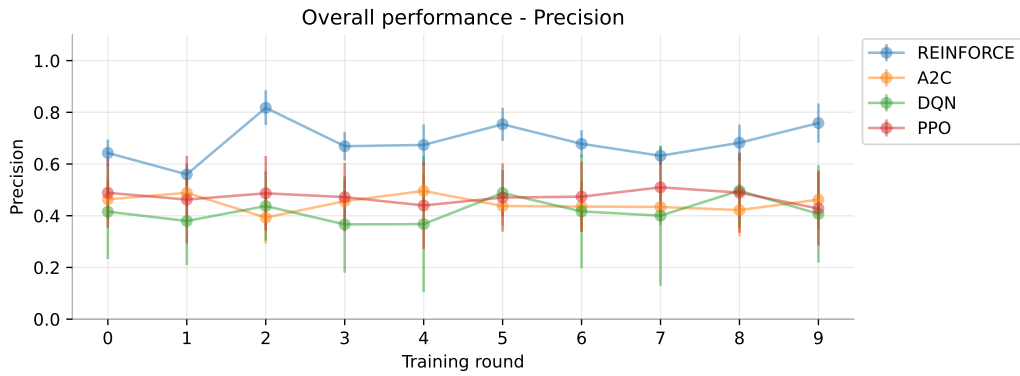
is the Real data offers a more challenging environment. Despite this, in Table 6, we notice that REINFORCE still performs better than the other algorithms. The margins are understandably lower: precision by 0.155, recall by 0.097, F1 by 0.117 and F-beta (0.5) by 0.130.

	Precision		Recall		F1-score		F-beta (0.5)	
	Mean	SD	Mean	SD	Mean	SD	Mean	SD
A2C	0.487	0.086	0.582	0.075	0.488	0.063	0.467	0.065
DQN	0.399	0.204	0.491	0.032	0.361	0.035	0.332	0.060
PPO	0.512	0.107	0.422	0.107	0.441	0.096	0.472	0.096
REINFORCE	0.813	0.119	0.421	0.079	0.495	0.090	0.609	0.101

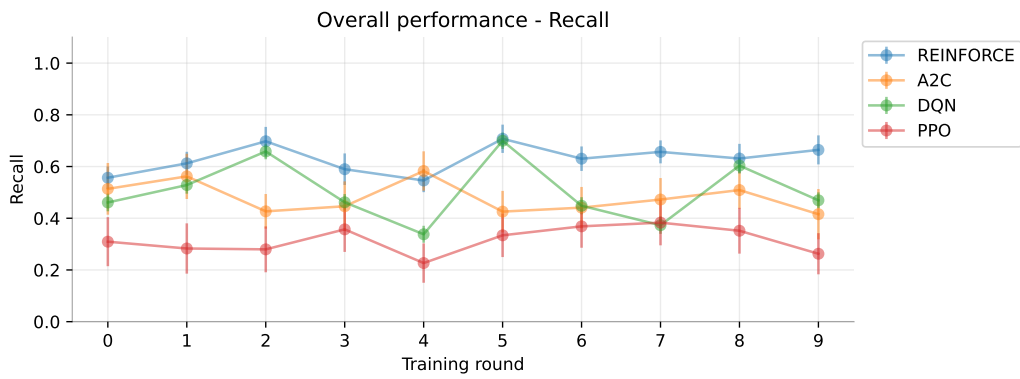
Table 7: Model performance summary - averaged over PHM-2010 environments with complex multivariate environment.

4.5 Plots

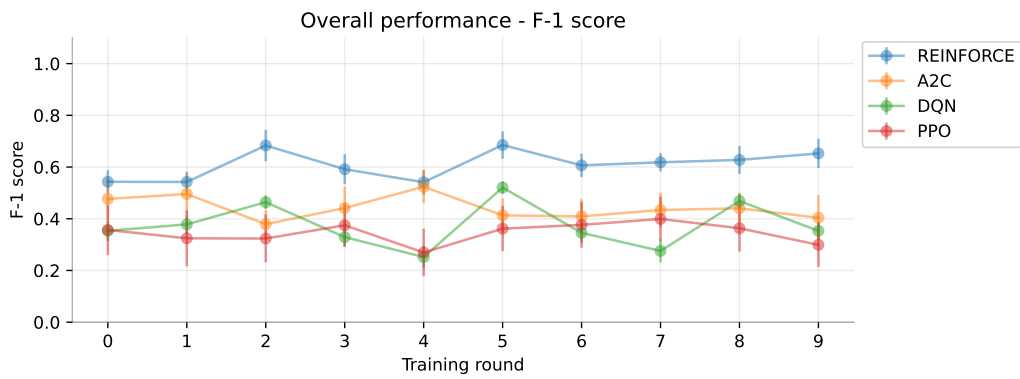
4.6 Super models metrics



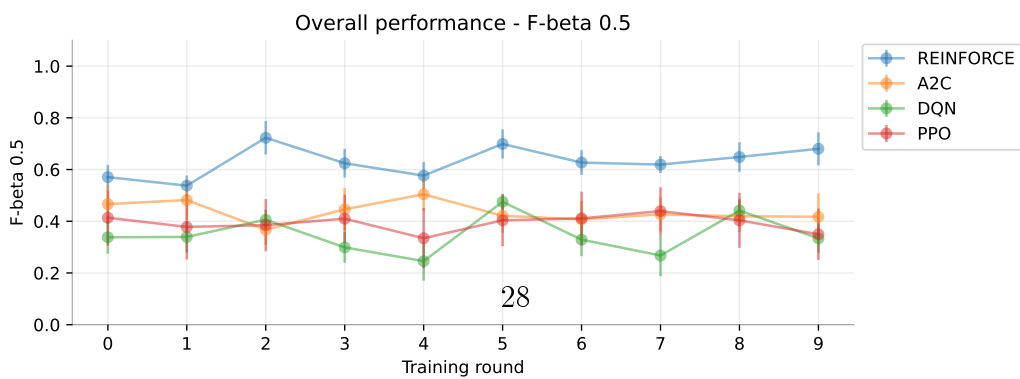
(a) Precision



(b) Recall

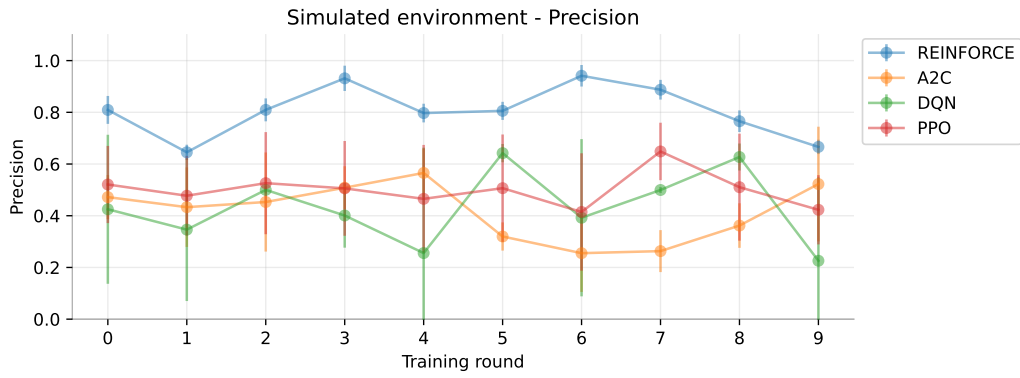


(c) F1-score

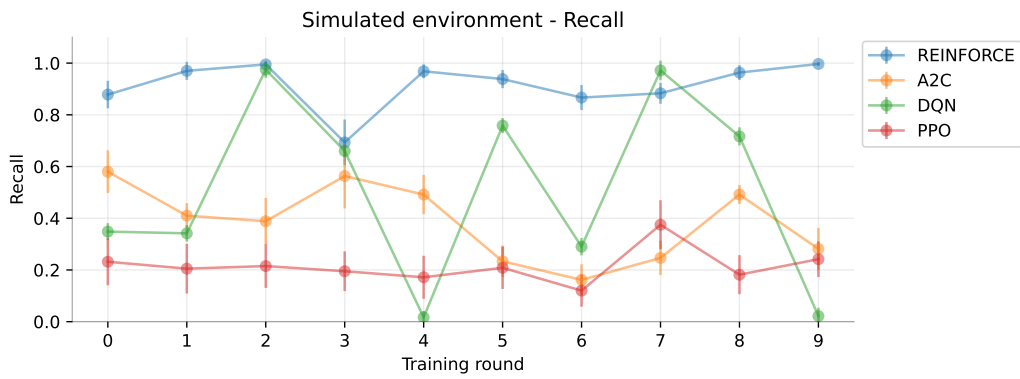


(d) F-beta (0.5)

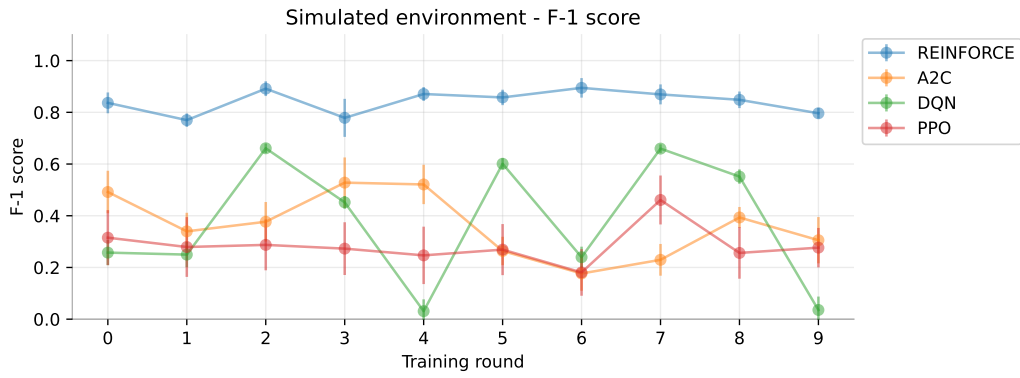
Figure 11: Overall performance – Average over 10 models



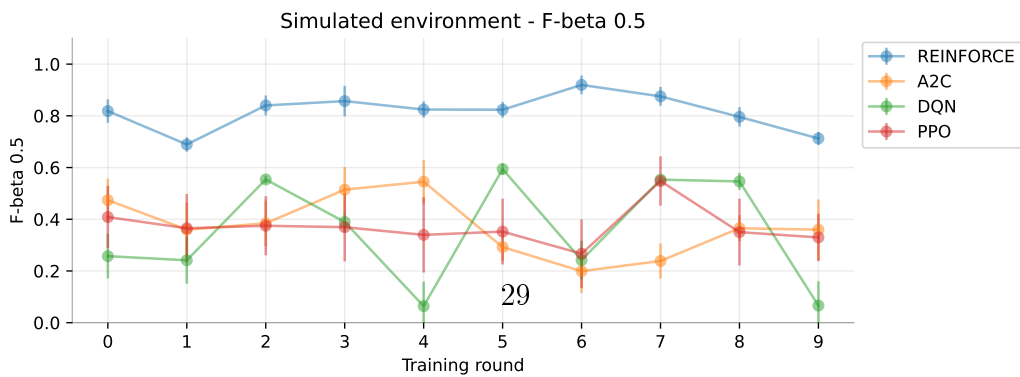
(a) Precision



(b) Recall

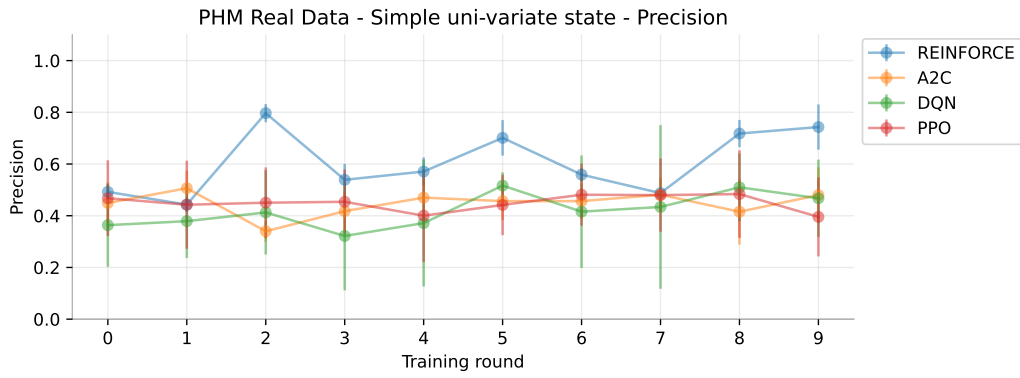


(c) F1-score

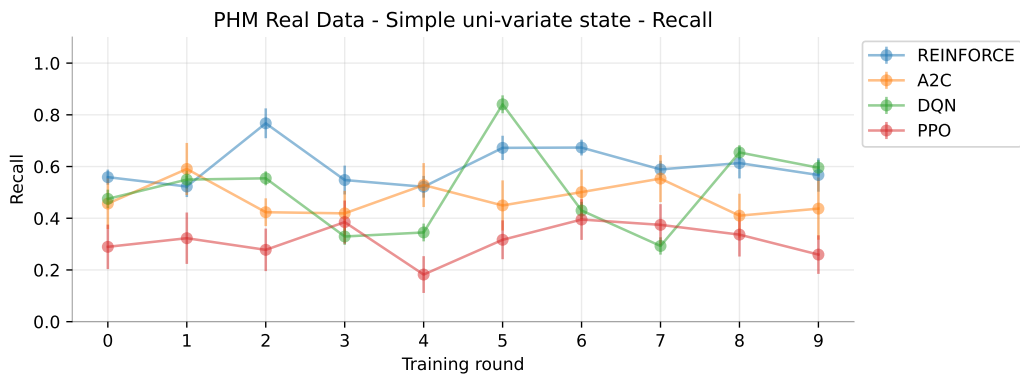


(d) F-beta (0.5)

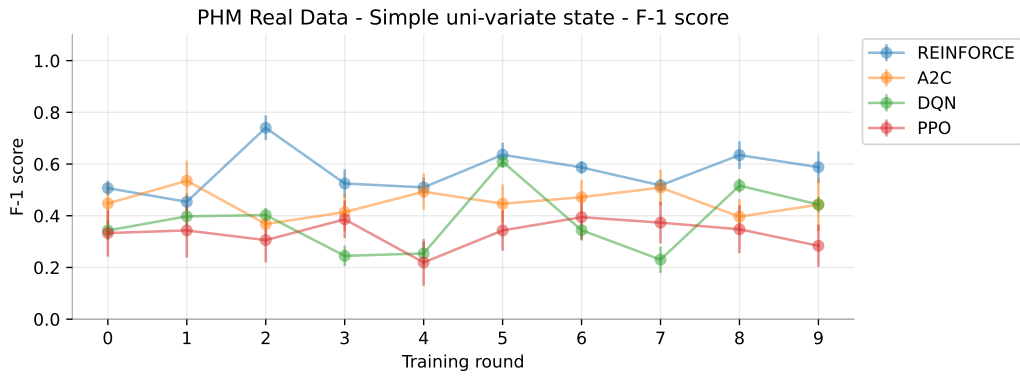
Figure 12: Simulated environment – Average over 10 models



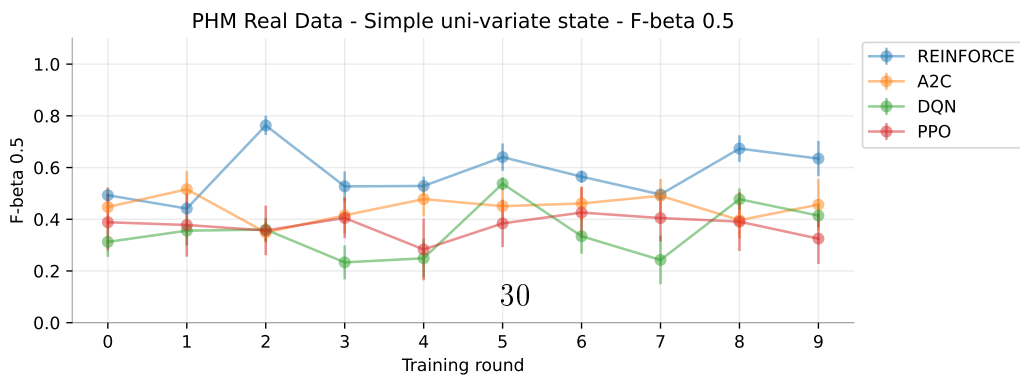
(a) Precision



(b) Recall

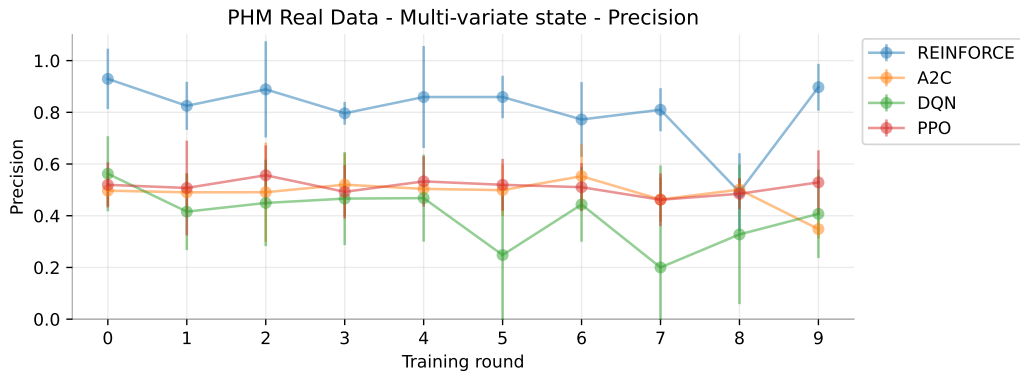


(c) F1-score

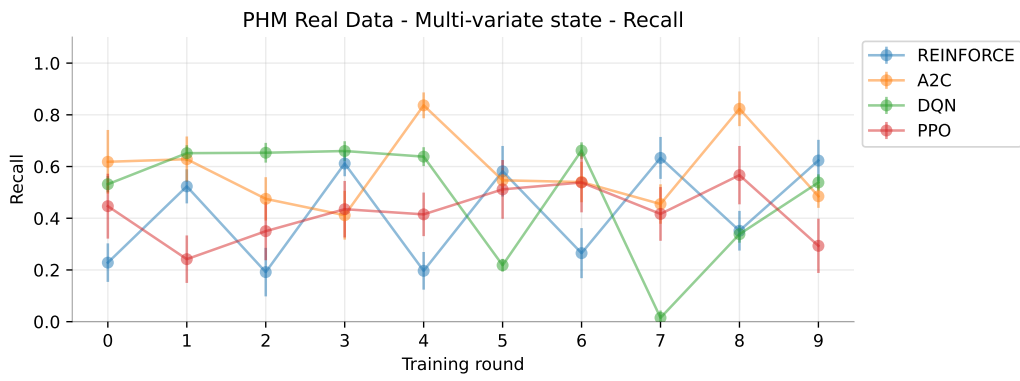


(d) F-beta (0.5)

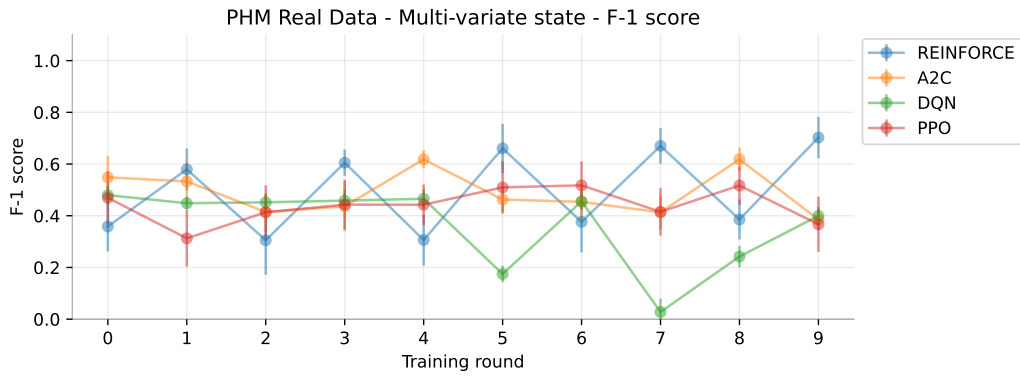
Figure 13: Univariate state environment – Average over 10 models



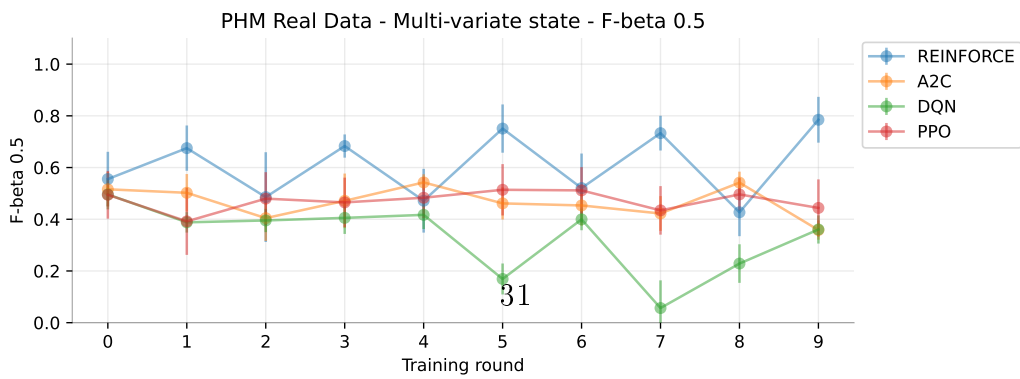
(a) Precision



(b) Recall



(c) F1-score



(d) F-beta (0.5)

Figure 14: Multivariate state environment – Average over 10 models

Environment	REINFORCE				A2C				DQN				PPO			
	Prec.	Recall	F1	$F_{\beta}0.5$	Prec.	Recall	F1	$F_{\beta}0.5$	Prec.	Recall	F1	$F_{\beta}0.5$	Prec.	Recall	F1	$F_{\beta}0.5$
Simulated - No noise	0.897	0.960	0.926	0.908	0.500	1.000	0.667	0.556	0.505	0.980	0.667	0.560	0.669	0.430	0.518	0.597
Simulated - Low noise	0.960	0.945	0.952	0.957	0.516	1.000	0.680	0.571	0.500	0.980	0.662	0.554	0.633	0.460	0.530	0.586
Simulated - High noise	0.922	0.990	0.955	0.935	0.503	1.000	0.669	0.558	0.504	0.990	0.668	0.559	0.569	0.355	0.434	0.505
PHM C01 SS - No noise	0.889	0.995	0.939	0.908	0.586	0.625	0.603	0.592	0.647	0.970	0.776	0.693	0.543	1.000	0.703	0.597
PHM C01 SS - Low noise	0.988	0.765	0.861	0.932	0.499	0.995	0.664	0.554	0.504	0.990	0.668	0.559	0.623	0.740	0.675	0.643
PHM C01 SS - High noise	0.850	0.970	0.905	0.871	0.521	0.680	0.588	0.546	0.505	0.985	0.668	0.560	0.520	0.725	0.604	0.551
PHM C04 SS - No noise	0.811	1.000	0.895	0.842	0.536	0.645	0.583	0.554	0.501	0.965	0.660	0.554	0.579	0.895	0.702	0.622
PHM C04 SS - Low noise	0.798	0.980	0.879	0.829	0.556	0.665	0.603	0.573	0.734	0.990	0.843	0.774	0.546	0.660	0.596	0.565
PHM C04 SS - High noise	0.708	0.840	0.767	0.730	0.521	0.835	0.641	0.563	0.511	0.985	0.672	0.565	0.517	0.820	0.633	0.558
PHM C06 SS - No noise	1.000	0.895	0.944	0.977	0.520	0.680	0.587	0.545	0.935	0.975	0.954	0.942	0.587	0.650	0.615	0.597
PHM C06 SS - Low noise	0.943	0.795	0.861	0.908	0.501	1.000	0.668	0.557	0.961	0.725	0.826	0.901	0.552	0.370	0.438	0.497
PHM C06 SS - High noise	0.821	0.845	0.831	0.825	0.540	0.755	0.628	0.572	0.980	0.960	0.969	0.976	0.521	0.615	0.564	0.537
PHM C01 MS - No noise	0.827	0.995	0.903	0.856	0.500	1.000	0.667	0.556	0.505	0.985	0.668	0.560	0.512	0.595	0.549	0.526
PHM C04 MS - No noise	0.910	0.425	0.577	0.738	0.500	1.000	0.667	0.556	0.501	0.975	0.662	0.555	0.501	0.635	0.558	0.522
PHM C06 MS - No noise	0.934	0.865	0.896	0.918	0.500	1.000	0.667	0.556	0.969	0.600	0.741	0.863	0.497	0.690	0.577	0.526

Table 8: Super Models: Best models selected over 10 rounds of training.

4.7 Hypothesis testing

Notes:

- Definition of 1 data point
- Single training round (5th of the 10 rounds)
- Single dataset: PHM C01
- Single env: Simple univariate state
- Single noise setting: HBD = high
- File has 40 rows = 4 algos x 10 test rounds
- 10 test rounds of 40 random sampled points
- 1 point = 1 row = 1 classification test with 40 samples
- E.g. Dasic = 10 training rounds x 10 test rounds x 3 noise settings = 300 sample points
- E.g. PHM SS = 3 datasets x 10 training rounds x 10 test rounds x 3 noise settings = 900 sample points
- Technically Dasic = 300 points for metrics = 300 * 40 sampled data points of wear data = 12000

$$\left. \begin{array}{l} H_0 : \mu_{RF} - \mu_{AA} = 0, \\ H_a : \mu_{RF} - \mu_{AA} > 0, \end{array} \right\} \quad \forall \ AA \in [A2C, DQN, PPO] \quad (13)$$

4.8 Training times

5 Discussion

****DISCUSSION**** It must be noted that the classification metrics we used to evaluate, assume that the human prevent maintenance policy is ideal and suggested based on experience. In practice this might not always be the case. The RL agent is trained to replace the tool optimally i.e. maximize use before reaching the threshold use, and replace only when necessary so as to maintain lowest possible cost of replacement. In our dataset we assume the threshold to be ideally set and base the “human” action against that, rather than a real human action, which incidentally is not available for the PHM

Metric	p Value			t Statistic		
	RF $\overset{H_a}{\underset{H_0}{\geq}}$ A2C	RF $\overset{H_a}{\underset{H_0}{\geq}}$ DQN	RF $\overset{H_a}{\underset{H_0}{\geq}}$ PPO	RF $\overset{H_a}{\underset{H_0}{\geq}}$ A2C	RF $\overset{H_a}{\underset{H_0}{\geq}}$ DQN	RF $\overset{H_a}{\underset{H_0}{\geq}}$ PPO
Overall (1500 samples)						
Precision	4.31E-126	2.17E-109	2.81E-106	25.071	23.170	22.804
Recall	4.20E-35	3.37E-16	4.36E-150	12.522	8.206	27.650
F1 score	1.99E-64	1.46E-88	5.29E-155	17.364	20.634	28.160
Simulated environment (300 samples)						
Precision	3.20E-98	1.69E-63	2.65E-81	25.611	19.032	22.427
Recall	8.12E-104	2.56E-41	1.57E-264	26.665	14.558	62.541
F1 score	9.60E-134	8.56E-99	2.96E-242	32.402	25.719	56.575
PHM Real data - Simple univariate state (900 samples)						
Precision	2.27E-32	7.29E-31	9.95E-31	12.082	11.770	11.742
Recall	1.27E-16	1.55E-06	8.19E-71	8.357	4.821	18.607
F1 score	1.94E-19	4.67E-34	2.19E-67	9.121	12.423	18.098
PHM Real data - Complex multivariate state (300 samples)						
Precision	1.64E-60	3.34E-54	7.88E-59	18.451	17.207	18.122
Recall	2.69E-10	2.69E-02	9.68E-01	-6.425	-2.219	-0.041
F1 score	7.27E-01	1.44E-08	1.35E-03	0.349	5.748	3.220

Table 9: One-tail t-test - Ho: No difference in metrics. Ha: REINFORCE metric > Advanced algorithm metric

dataset anyway.

- R1 = +1
- R2 = -1
- R3 = -100 => higher neg. Improve recall. Lower neg. Improve preci-

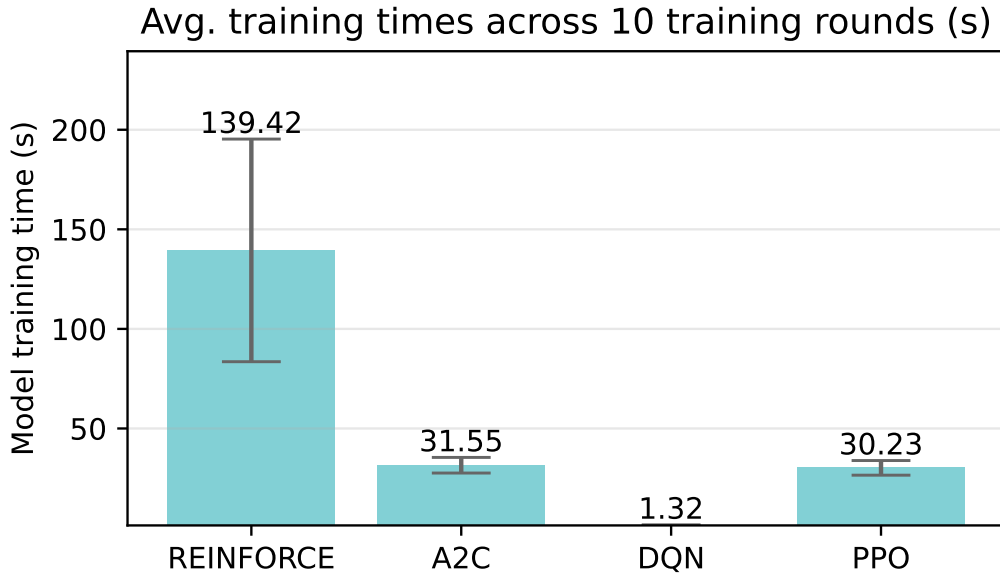


Figure 15: Training time. Across 10 rounds and all environment variants.

sion.

- Hyper-parameters for Precision or Recall control:
- Precision => low FP => False replacement-action reduced. *** Unnecessary tool replacement reduced. Tool life maximized, down time minimized, production disruption minimized.
- - Recall => low FN => False declaration of normal operation reduced. Reduce missed replacements. Tool replacements increased. *** Product quality not compromised.
- **LOOK AHEAD PARAM:**
- It must be noted that the classification metrics we used to evaluate, assume that the human prevent maintenance policy is ideal. Which in practice might not be the case. The RL is trained to replace the tool optimally i.e. maximize use and replace only when necessary so as to maintain lowest possible cost of replacement.

- Transition function (Graesser and Keng, 2019): pg333 "Transition function - also known as model. can be learnt or programmed programmable rules are common and lead to increasing complexity. *** We use a hybrid model - partly learnt and partly programmed": **Noise and breakdown**

Duan et al. (2016) the environment.*OBSERVATION**REINFORCE: Despite its simplicity, REINFORCE is an effective algorithm in optimizing deep neural network policies in most basic and locomotion tasks. Even for high- DOF tasks like Ant, REINFORCE can achieve competitive results. However we observe that REINFORCE sometimes suffers from premature convergence to local optima as noted by Peters Schaal (2008), which explains the performance gaps between REINFORCE and TNPG on tasks such as Walker (Figure 3(a)). By visualizing the final policies, we can see that REINFORCE results in policies that tend to jump forward and fall over to maximize short-term return instead of acquiring a stable walking gait to maximize long-term return. In Figure 3(b), we can observe that even with a small learning rate, steps taken by REINFORCE can sometimes result in large changes to policy distribution, which may explain the fast convergence to local optima.

5.1 Precision or Recall?

- Precision => low FP => False replacement-action reduced. *** Unnecessary tool replacement reduced. Tool life maximized, down time minimized, production disruption minimized

- Recall => low FN => False declaration of normal operation reduced. Reduce missed replacements. Tool replacements increased. *** Product quality not compromised.

limitations of the selecting model for evaluation

Raffin et al. (2021) – quote from paper - intro section – "A major challenge is that small implementation details can have a substantial effect on performance often greater than the difference between algorithms (Engstrom et al., 2020). It is particularly important that implementations used as ex-

perimental baselines are reliable; otherwise, novel algorithms compared to weak baselines lead to inflated estimates of performance improvements"

DPO paper - Rafael Rafailov May-2023 — Direct Preference Optimization: Your Language Model is Secretly a Reward Model

5.2 Instability of Actor-Critic Algorithms We can also use our framework to diagnose instabilities with standard actor-critic algorithms used for the RLHF, such as PPO. We follow the RLHF pipeline and focus on the RL fine-tuning step outlined in Section 3. We can draw connections to the control as inference framework [18] for the constrained RL problem outlined in 3. We assume a parameterized model and minimize This is the same objective optimized in prior works [45, 35, 1, 23] using the DPO-equivalent reward for the reward class of In this setting, we can interpret the normalization term in as the soft value function of the reference policy . While this term does not affect the optimal solution, without it, the policy gradient of the objective could have high variance, making learning unstable. We can accommodate for the normalization term using a learned value function, but that can also be difficult to optimize. Alternatively, prior works have normalized rewards using a human completion baseline, essentially a single sample Monte-Carlo estimate of the normalizing term. In contrast the DPO reparameterization yields a reward function that does not require any baselines.

6 Conclusion

References

- Predrag Dašić. Analysis of wear cutting tools by complex power-exponential function for finishing turning of the hardened steel 20crmo5 by mixed ceramic tools. *Fascicle VIII Tribology*, 12:54–60, 2006.
- Berend Denkena, Heinrich Klemme, and Tobias H. Stiehl. Tool wear monitoring using process data of multiple machine tools by means of machine learning. *ZWF Zeitschrift fuer Wirtschaftlichen Fabrikbetrieb*, 118(5):298 – 301, 2023. doi: 10.1515/zwf-2023-1059.

- Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. Openai baselines. <https://github.com/openai/baselines>, 2017.
- Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *International conference on machine learning*, pages 1329–1338. PMLR, 2016.
- Gabriel Dulac-Arnold, Nir Levine, Daniel J Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal, and Todd Hester. An empirical investigation of the challenges of real-world reinforcement learning. *arXiv preprint arXiv:2003.11881*, 2020.
- Gabriel Dulac-Arnold, Nir Levine, Daniel J Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal, and Todd Hester. Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. *Machine Learning*, 110(9):2419–2468, 2021.
- L. Erhan, M. Ndubuaku, M. Di Mauro, W. Song, M. Chen, G. Fortino, O. Bagdasar, and A. Liotta. Smart anomaly detection in sensor systems: A multi-perspective review. *Information Fusion*, 67:64–79, 2021. doi: 10.1016/j.inffus.2020.10.001.
- SA Ford and M Ritchie. Cognitive radar mode control: a comparison of different reinforcement learning algorithms. In *International Conference on Radar Systems (RADAR 2022)*, volume 2022, pages 107–112. IET, 2022.
- Future Market Insights. Milling machine market outlook (2023 to 2033), 1 2023. URL <https://www.futuremarketinsights.com/reports/milling-machine-market>. Accessed: 2023-06-23.
- Laura Graesser and Wah Loon Keng. *Foundations of deep reinforcement learning*. Addison-Wesley Professional, 2019.

- Russul H. Hadi, Haider N. Hady, Ahmed M. Hasan, Ammar Al-Jodah, and Amjad J. Humaidi. Improved fault classification for predictive maintenance in industrial iot based on automl: A case study of ball-bearing faults. *PROCESSES*, 11(5), MAY 15 2023. doi: 10.3390/pr11051507.
- Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Rousslan Fernand Julien; Raffin Antonin; Kanervisto Anssi; Wang Weixun Huang, Shengyi; Dossa. The 37 implementation details of proximal policy optimization. March 2022. Accessed: 2023-06-27.
- Justin Larocque-Villiers, Patrick Dumond, and David Knox. Automating predictive maintenance using state-based transfer learning and ensemble methods. In *2021 IEEE INTERNATIONAL SYMPOSIUM ON ROBOTIC AND SENSORS ENVIRONMENTS (ROSE 2021)*. 2021. doi: 10.1109/ROSE52750.2021.9611768.
- Matthias Maurer, Andreas Festl, Bor Bricelj, Germar Schneider, and Michael Schmeja. Automl for log file analysis (alfa) in a production line system of systems pointed towards predictive maintenance. *INFOCOMMUNICATIONS JOURNAL*, 13(3):76–84, SEP 2021. doi: 10.36244/ICJ.2021.3.8.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.
- Taisuke Oshida, Tomohiro Murakoshi, Libo Zhou, Hirotaka Ojima, Kazuki Kaneko, Teppei Onuki, and Jun Shimizu. Development and implementation of real-time anomaly detection on tool wear based on stacked lstm

- encoder-decoder model. *The International Journal of Advanced Manufacturing Technology*, pages 1–16, 2023.
- Marcel Panzer and Benedict Bender. Deep reinforcement learning in production systems: a systematic literature review. *International Journal of Production Research*, 2021.
- Biyao Qiang, Kaining Shi, Ning Liu, Junxue Ren, and Yaoyao Shi. Integrating physics-informed recurrent gaussian process regression into instance transfer for predicting tool wear in milling process. *Journal of Manufacturing Systems*, 68:42 – 55, 2023. doi: 10.1016/j.jmsy.2023.02.019.
- Bo Qin, Yongqing Wang, Kuo Liu, Shaowei Jiang, and Qi Luo. A novel online tool condition monitoring method for milling titanium alloy with consideration of tool wear law. *Mechanical Systems and Signal Processing*, 199, 2023. doi: 10.1016/j.ymssp.2023.110467.
- Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *J. Mach. Learn. Res.*, 22(1), jan 2021. ISSN 1532-4435.
- N Sandeep Varma, Vaishnavi Sinha, and K Pradyumna Rahul. Experimental evaluation of reinforcement learning algorithms. In *International Conference on Computational Intelligence and Data Engineering*, pages 469–484. Springer, 2022.
- SB3-Algorithms. Stable-baselines3 - master list of algorithms, 2022. URL <https://stable-baselines3.readthedocs.io/en/master/guide/algos.html>. Accessed: 2023-06-27.
- SB3-Default Network Architecture. Stable-baselines3 - default network architecture, 2022. URL https://stable-baselines3.readthedocs.io/en/master/guide/custom_policy.html#default-network-architecture. Accessed: 2023-06-27.

- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Rajesh Siraskar, Satish Kumar, Shruti Patil, Arunkumar Bongale, and Ketan Kotecha. Reinforcement learning for predictive maintenance: a systematic technical review. *Artificial Intelligence Review*, pages 1–63, 2023.
- Richard Sutton and Andrew Barto. *Reinforcement Learning: An Introduction*. The MIT Press, Cambridge, England, 2nd. edition edition, 2018.
- Paweł Twardowski, Jakub Czyżycki, Agata Felusiak-Czyryca, Maciej Tabaszewski, and Martyna Wiciak-Pikuła. Monitoring and forecasting of tool wear based on measurements of vibration accelerations during cast iron milling. *Journal of Manufacturing Processes*, 95:342 – 350, 2023. doi: 10.1016/j.jmapro.2023.04.036.
- Krishna Velivela and Sudhir Yarram. Comparison of reinforcement learning algorithms, Dec 2020.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Reinforcement learning*, pages 5–32, 1992.
- LI Yingguang, LIU Changqing, LI Dehua, HUA Jiaqi, and WAN Peng. Tool wear dataset of nuaa ideahouse, 2021. URL <https://dx.doi.org/10.21227/3aa1-5e83>.