

# Predictive Maintenance for Edge-Based Sensor Networks: A Deep Reinforcement Learning Approach

Kevin Shen Hoong Ong<sup>\*</sup>, Dusit Niyato<sup>\*</sup>, Chau Yuen<sup>†</sup>

<sup>\*</sup>School of Computer Science and Engineering, Nanyang Technological University Singapore

<sup>†</sup>Engineering Product Development, Singapore University of Technology and Design

**Abstract**—Failure of mission-critical equipment interrupts production and results in monetary loss. The risk of unplanned equipment downtime can be minimized through Predictive Maintenance of revenue generating assets to ensure optimal performance and safe operation of equipment. However, the increased sensorization of the equipment generates a data deluge, and existing machine-learning based predictive model alone becomes inadequate for timely equipment condition predictions. In this paper, a model-free Deep Reinforcement Learning algorithm is proposed for predictive equipment maintenance from an equipment-based sensor network context. Within each equipment, a sensor device aggregates raw sensor data, and the equipment health status is analyzed for anomalous events. Unlike traditional black-box regression models, the proposed algorithm self-learns an optimal maintenance policy and provides actionable recommendation for each equipment. Our experimental results demonstrate the potential for broader range of equipment maintenance applications as an automatic learning framework.

## I. INTRODUCTION

Equipment downtime is generally defined as the outage time that accumulates whenever production process stops and current world class standards for downtime is  $\leq 10\%$  [1]. Despite the rapid technological advances and increasing equipment complexity, frequent occurrence of equipment downtime remains and often results in monetary loss. Maintenance teams are given limited maintenance budgets and face tremendous cost pressure to ensure that the production line is always operational. In an event of an unplanned equipment fault, the corrective maintenance option is performed to bring the failed equipment up to its operational status to meet the product delivery deadline. Forward looking companies employ preventive maintenance to reduce the likelihood of unplanned equipment downtime through scheduled upkeep of production equipment condition. Although the long-term goal of reducing overall maintenance costs has been touted, manufacturing productivity is slightly improved at the expense of higher maintenance cost. Predictive maintenance is overall considered more efficient (i.e. manpower and cost) and equipment downtime is minimized because maintenance is only performed based on the real-time condition of an equipment.

As industries worldwide journey towards the Industry 4.0 vision to boost manufacturing productivity, modern equipments become increasingly complex and requires longer maintenance time. Admist a manpower lean economy that is driven by

productivity, one of the key challenges resides in the latent demand to simplify the complexity of machine sensor data interpretation for predictive maintenance purposes. Traditional black-box regression models are feature-engineered towards domain-specific applications, and solution extensibility to similar applications and feature updates request are costly short-term endeavours. Deep Learning (DL) have recently been proposed as an alternative for estimation of remaining useful life of equipment [2]–[4] and ball bearings [5], [6]. Similar work has also been reported in [7] to learn the health indicators for the remaining useful life estimation of a turbofan engine by using Temporal Difference (TD) learning. However, these approaches are only concerned about accurate estimation of the equipment’s remaining operational uptime and lack meaningful insights to support the maintenance team’s decision-making process. Common reasons can be attributed to the fact that explicit models of real-world problems are largely unknown or too complex to accurately model with traditional model-based approaches. In this work, we propose a model-free Deep Reinforcement Learning (DRL) algorithm approach to self-learn optimal maintenance decision policies, from the health state of an equipment, more importantly with actionable recommendation.

Increased pressure on margins, higher customer expectations and the declining cost of cloud computing are enticing factors for manufacturers to stay competitive. Conversely, both manufacturing and equipment-maker companies remain security averse to transmitting sensitive production sensor data to the cloud, via unreliable internet connections, and any adoption is generally reserved only for companies with deep pockets [8]. An alternative approach is to utilize a locally interconnected sensor network of equipment. In our proposed system model, each equipment is equipped with sensor-based edge computing devices [9], which can share data and communicate with other edge devices or equipment. Overall benefits of adopting this approach includes time-efficient decision-making process, a responsive yet data informed maintenance support team and reduction in network traffic across the factory shopfloor IT infrastructure.

In summary, the contributions of this paper are as follows:

- 1) Problem formulation of maximizing equipment runtime as a function of multiple sensor data input. The decision-policy is obtained by using the model-free DRL

approach.

- 2) The DRL algorithm offers recommendation support for the replacement policy of an equipment with easy-to-understand data-driven recommendations via an Equipment Health Indicator status.
- 3) Despite random initial health states and an absent ground truth, the proposed DRL algorithm tackles the sparsely-dense reward maintenance problem with almost consistent recommendations across similar datasets and equipment.

## II. SYSTEM MODEL

Consider the application of a sensor network for the purpose of equipment health monitoring, either through retrofitting or in-situ configuration. Our proposed system model resembles a Star Topology network at the Equipment level, see Figure 1. From the current configuration, the aggregated data can be propagated to a larger Sensor Network, such as tree or mesh topology configurations, for more complex data analysis. Typical network components include Sensor Nodes (SN) and a Sensor Gateway (SG), also known as Base Station. Please note that the terms Base Station and SG are used interchangeably.

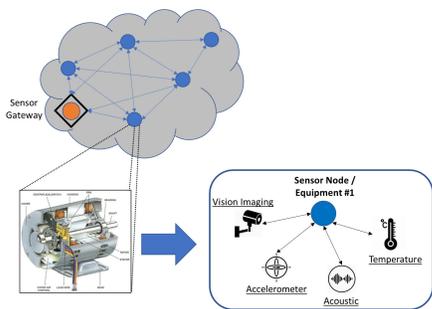


Fig. 1: System Model of Sensor Network at Equipment Level

SN devices are model representations of physical sensors transmitting data such as temperature, accelerometer, gyroscope and acoustic, depicted in Figure 1. Concurrent streaming of raw sensor data is known to consume high network bandwidth and degrades overall network performance. To mitigate network congestion scenario, SG-based data aggregation and processing at the SG node is proposed. The SG node is assumed to be always powered with adequate computation and storage capabilities. Next, sensor data analysis will be performed using batch processing (i.e time-frame comprising of fixed time-step duration) of sensor data for memory and computation efficiency purposes. Without loss of generality and for ease of presentation, a pair of SG and SN is assumed.

Data analysis is performed on window size with  $\mathcal{K}$  time-steps and the output of the prediction are three maintenance action: *Repair*, *Replace* and *Hold*. *Hold* is the default action where SG predicts a low probability of imminent failure to occur. In the event that an anomalous sensor reading is observed with high certainty, the SG node has to decide which of the two recommendations, *Replace* or *Repair*, should be provided to the maintenance team - mirroring real-life decision-making. In practice, the *Repair* action is almost certainly executed given the busy schedule and maintenance budget constraints.

## III. PROBLEM FORMULATION

In our system model, the challenge is to maximize the total run-time of the equipment with maintenance budget constraints. The objective function is described in (1) as Maximum Equipment UpTime ( $\rho$ ) and is cast into a Markov Decision Process (MDP) framework with fully observable states. MDP is formally described as a 5-tuple consisting of state ( $\mathcal{S}$ ), action ( $\mathcal{A}$ ), probabilistic distribution of state transitions ( $\mathcal{P}$ ), reward function ( $\mathcal{R}$ ) and discount factor ( $\gamma$ ). Mathematically, the tuple can be compactly denoted as ( $\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma$ ).

$$\rho = \sum_{Node=1}^N RunTime \quad (1)$$

Given a sensor network, the generated sensor data is denoted as  $x_t^i$ , where  $i \in \{0, 1, \dots, Z\}$  at every time-step ( $t$ ). From (1),  $N$  represents the number of SN devices within the considered sensor network, and  $N = 1$  is assumed. The sensor data is then discretized and simplified in (2). Then, we generalize the state space of each sensor node in (3).

$$q_x \leftarrow Discretize(x_t^i) \quad (2)$$

$$\mathcal{S}^i = \{q_x\} \quad (3)$$

For every equipment or sensor manufactured, the manufacturers specify the Mean-Time-Between-Failure (MTBF) and operating temperature information in the technical datasheet. Due to the elasticity of environmental temperature, the equipment's degradation rate and state change is inadvertently non-sequential. In this work, we consider the sensor state changes to be initially sequential with decreasing trend over finite time-steps. For example, assume the probability of state change increases as operating temperature increases. Eventually, an obvious temperature mode change occurs, and the rate of sensor degradation decays exponentially with respect to time, described in (4). Within the current operating temperature mode, the sensor's state change would skip multiple states to reflect corresponding temperature changes. The observed state change behavior is known to mirror a concave exponential decay trend with respect to equipment run-time. When the transition probability of the state-action pair is considered, a pseudo health indicator is derived and shown to vary according to increasing failure probability, see Figure 3. Notably, the health degradation trend clearly illustrates inverse correlation to the increasing sensor state values and is aligned with the proposed system model.

$$\mathcal{F}(t) = e^{-\lambda t} \quad (4)$$

To model the environmental state ( $\mathcal{S}^\tau$ ), we consider and simplify the operating temperature conditions ( $\tau$ ) into a binary form with conditional constraints, see (5). To be clear, units of  $\tau$  is in degree Celsius while Low and High operating temperatures are binarily represented.

$$\mathcal{S}^\tau = \begin{cases} 0, & \text{if } \tau \in [25, 60] \\ 1, & \text{if } \tau > 60 \end{cases} \quad (5)$$

The resultant state space for our system model is summarized as the Cartesian product:

$$\mathcal{S} = \mathcal{S}^i \times \mathcal{S}^\tau \quad (6)$$

Next, the model's action space is encoded as a vector of scalar actions:  $Replace(\epsilon)$ ,  $Repair(\eta)$  and  $Hold(\kappa)$ . Let us assume that a maintenance account ( $\beta$ ) is credited to the maintenance agent.  $\kappa$  represents the agent's default action and the associated maintenance cost is zero. As the equipment state (i.e. sensor health values) starts degrading, the frequency and cost of repairs will gradually increase with time before increasing exponentially. Hence, the agent is tasked to decide the appropriate sequence of actions to perform at each state as described in (7). To mimic maintenance decision-making process, the imposed cost constraints ( $\mathcal{C}$ ) ensures that the agent derives a sensible maintenance policy where  $\beta$  is not violated within the given time frame.

$$\mathcal{A} = \begin{cases} (\epsilon, \eta, \kappa) | \\ \sum_{n=1}^N \mathcal{C}_\epsilon \geq 2 \sum_{n=1}^N \mathcal{C}_\eta \text{ and } \beta - \sum_{n=1}^N \mathcal{C}_\epsilon \geq 0, \\ \sum_{n=1}^N \mathcal{C}_\epsilon \leq \sum_{n=1}^N (\mathcal{C}_\eta/2) \text{ and } \beta - \sum_{n=1}^N \mathcal{C}_\eta \geq 0 \end{cases} \quad (7)$$

For model simplification purposes, consider a single equipment with  $N = 1$  sensor attached. Given observable state changes and transitions, the agent can select random actions to perform.  $\epsilon$  resets the sensor's operational state to an almost new condition while  $\eta$  seemingly reverts the sensor's operational state by  $y_{Repair}$  states, to a previously observed sensor state. In addition, we consider that not all repairs are identical and the sensor state change  $\phi(\mathcal{S})$ , will vary according to the  $RepairType(\psi)$ , where  $\psi \in \{0, 1, \dots, M\}$ . From Figure 2,  $\eta$  is invoked at  $\mathcal{S}_t = y - 1$  and the selected type of repair invokes a state transition from  $\mathcal{S}_t = y$  to  $\mathcal{S}'_t = y - |\phi(\mathcal{S})|$ . A compact representation of the behaviour is described in (8).

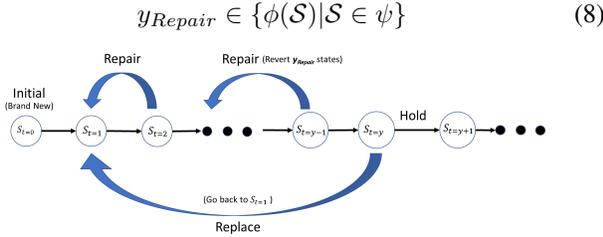


Fig. 2: Example State Space and Transition with Replace and Repair Action

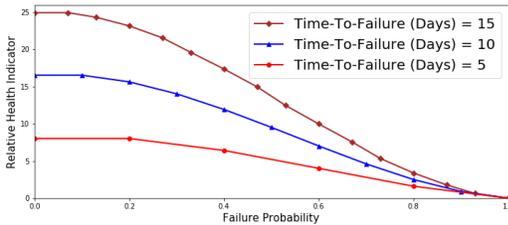


Fig. 3: Health Degradation wrt Equipment Failure Probability

The Reward function is formally described as  $\mathcal{R}(s_t, a_t, s_{t+1})$ . From the previous derivations,  $\mathcal{R}$  can be re-interpreted as the cumulative sum of the sensor runtime with respect to state  $\mathcal{S}$  and action  $\mathcal{A}$ . An immediate reward function  $\mathcal{R}_t \in \mathcal{R}$  is proposed in Equation 9 to guide the agent's actions:

$$\mathcal{R}_t = \begin{cases} \mathcal{R}_{Rpl}, & \text{if } \mathcal{S}_t^i > 0, \beta > 0 \\ \mathcal{R}_{Rpa}, & \text{if } \mathcal{S}_t^i > 0, \beta > 0 \\ \mathcal{R}_{Exp}, & \text{if } \mathcal{S}_t^i > 0 \\ \mathcal{R}_{Frug}, & \text{if } \mathcal{S}_t^i > 0, \beta > 0 \\ \mathcal{R}_{Pen}, & \text{otherwise} \end{cases} \quad (9)$$

For example, executing either  $y_{Replace}$  or  $y_{Repair}$  action, within the  $\beta$  constraints, will offer the agent an arbitrary high health points for performing Replace ( $\mathcal{R}_{Rpl}$ ) and Repair ( $\mathcal{R}_{Rpa}$ ) actions respectively. Conversely, the agent is heavily penalized if either  $\kappa$  is performed throughout the episode or when the sensor state is zero ( $\mathcal{R}_{Pen}$ ). Real-world applications have sparse rewards and motivating the agent to explore is very challenging. In this paper,  $\mathcal{R}_{Exp}$  value is arbitrarily defined to intrinsically motivate the agent to explore the environment. To increase the probability of sparse state visitation, where potentially high rewards are found, a proposed ranking mechanism is proposed and described in Section IV-B.  $\mathcal{R}_{Frug}$  reward is also defined to encourage the agent to execute the corresponding action optimally, mimicking human decision-making.

Next, the state-value function ( $V^\pi(s)$ ) is expressed as:

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma \mathcal{R}_{t+1} | \mathcal{S}_t = s \right] \quad (10)$$

Then, we utilize the standard RL framework to obtain the optimal policy ( $\pi^*$ ). The Markov property is applied on (10) and the value function is simplified and re-expressed as:

$$V^{\pi^*}(s) = \sum_a \pi(s, a) \sum_{s' \in \mathcal{S}} \mathcal{P}_{\pi(s)}(s, s') [\mathcal{R}(s, a) + \gamma V^{\pi^*}(s')] \quad (11)$$

The associated policy function obtains the maximum action that is possible from Equation 11. Hence, the Q-function  $Q^*(s, a)$  can be updated using the Bellman equation and expressed as:

$$Q^*(s, a) = (1 - \alpha) \underbrace{Q(s, a)}_{\text{Current Q value}} + \alpha \left[ \underbrace{\mathcal{R}(s, a)}_{\text{Reward received}} + \gamma \underbrace{\max_{a' \in \mathcal{A}} Q'(s', a')}_{\text{Max(Expected future reward)}} \right]$$

where  $\alpha$  denotes the learning rate of the agent;  $\gamma \in [0, 1]$  is the discount factor where the agent performs tradeoff between the observed immediate and potential future reward. The max operator selects the highest-valued state-action pair that consequently assists in the derivation of the optimal policy.

Many real-world problems have very large state space, which makes it infeasible to compute and learn all exact state-action transition values. Instead, an approximated estimate of the state-action value pair can be learned using Temporal Difference (TD) learning. As TD learning process resembles a stochastic gradient descent, the updated Q-value is denoted as  $Q(S_t^i, A_t; \theta_t)$  towards a target value of  $y_t^Q$  in (12).

$$y_t^Q \equiv \mathcal{R}_{t+1} + \gamma \max_A Q(S_{t+1}^i, A; \theta_t) \quad (12)$$

In coupling a multi-layered neural network, a Deep Q-Network (DQN), with experience replay and target network

( $y_t^{DQN}$ ), performance of the algorithm is vastly improved and achieved fairly good generalization performance [10]. The target network is a copy of the original Q-network and parameter synchronization occurs every  $\tau$  steps, such that  $\theta_t^- = \theta_t$ , and is represented in (13).

$$y_t^{DQN} \equiv \mathcal{R}_{t+1} + \gamma \max_A \mathcal{Q}(s_{t+1}^i, A; \theta_t^-) \quad (13)$$

#### IV. SOLUTION OF DEEP Q LEARNING

##### A. Double Deep Q Learning

DQN's inherent tendency for value overestimation and biased estimates are caused by random environment noise and arg max operator respectively. Double Deep Q-Learning (DDQN) was then proposed [11] to stabilize vanilla DQN algorithm by decoupling the choice and evaluation of best action on two separate networks. In this work, the DDQN inputs contain a tuple of sensor data ( $x_t^i$ ) at every time-step and random action is performed by the DDQN agent. The DDQN's target network generates an output of Q-values and is denoted as:

$$y_t^{DoubleDQN} \equiv R_{t+1} + \gamma \mathcal{Q}(s_{t+1}, a^*; \theta_t^-) \quad (14)$$

where  $a^* = \arg \max \mathcal{Q}(s_{t+1}, A; \theta_t)$ . The discounted Q-value ( $y_t^{DoubleDQN}$ ) is taken from the target network with weights  $\theta_t^-$  and the target network weights ( $\theta_t^-$ ) are periodically copied from the Q-value network. In particular, the agent's performed actions are extracted from the primary Q-value network and the future reward evaluation step is taken from the Q-target network.

##### B. Prioritized Experience Replay

The DRL agent is expected to infer an optimal point from the system model or dataset from the reward constraints. With reference to Figure 3, the equipment replacement point is likely to occur further into the equipment operational cycle as failure probability increases exponentially. Likewise within the DRL context, the sparsely-dense-like reward configuration and the likelihood of performing either  $\mathcal{R}_{Rpl}$  or  $\mathcal{R}_{Rpa}$ , is higher towards the end of the equipment's operational cycle. The gained intuition suggests a normalized equipment health state value of 0.2 to be the target optimal value.

Although the traditional combination of  $\epsilon$ -greedy algorithm with Experience Replay (ER) Buffer technique is well-known, the ER Buffer has a bias to repeatedly sample the same highly-rewarding experiences and is ill-posed for a sparse-reward problem, typical of real-world applications. Moreover, the decay rate of  $\epsilon$  value over the training set is a hyperparameter and an inefficient approach given varying Time-to-Failure cycles for each equipment. For the mentioned problems, *Prioritized Experience Replay* (PER) [12] is the proposed component to compliment DDQN.

PER is considered to be an enhancement over ER, and it prioritizes the experiences which offers large differences between prediction and the Temporal Difference (TD) target value. Implicitly, the DRL agent can use the TD error magnitude as an indicator to better focus its attention on the least visited states which could yield potentially larger rewards.

The magnitude of the TD error ( $|\delta_i|$ ) can be incorporated into the ER buffer sample as a tuple ( $s_t, a_t, r_t, s_{t+1}, |\delta_i|$ ). To minimize over-fitting problem in the ER method, stochastic prioritization [12] is introduced to generate a probability  $P(i)$ , of a selected state-action pair, from the replay buffer, see (15).  $p_i$  denotes the priority value of the  $i$ -th sample in the buffer and  $\alpha$  is a hyperparameter that is used to induce randomness within the experience selection for the replay buffer. Pure uniform randomness is denoted by  $\alpha = 0$  while  $\alpha = 1$  emphasizes experiences with the highest priorities;  $\sum_k$  denotes the normalization of all priority values within the Replay Buffer.

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha} \quad (15)$$

$$\left( \frac{1}{N} \cdot \frac{1}{P(i)} \right)^b \quad (16)$$

During training of the neural network, the experiences sampled must match the underlying distribution with priority sampling, resulting in a bias towards high-priority samples. The weights of frequently seen samples, within the replay buffer, are then adjusted using Importance Sampling Weights (IS) with the effect of reducing bias. From (16),  $N$  denotes the Replay Buffer size;  $P(i)$  denotes the sampling probability from (15);  $b$  is considered a weighting factor to control the degree in which IS affects the learning process and is annealed up to 1 over the duration of training phase. Readers may refer to [12] for additional details on PER.

##### C. Parameter Noise

Exploration inefficiency of Reinforcement Learning (RL) is a well-known problem and becomes more challenging when applied on a sparse reward problem. Existing RL approach influences the action space policy at each time-step with no influence on the RL agent's decision policy (i.e neural network). Evolution strategies seek to manipulate the decision policy parameters during each rollout and no influence is exerted on the action space policy. *Parameter Noise* (PN) [13] is a technique which strikes an in-between balance of the aforementioned approaches with encouraging results for both on-policy and off-policy algorithms. To explain, PN randomly alters the parameters of the RL agent's decisions which introduces a more consistent exploration behaviour and a less confused RL agent. In this work, PN is introduced to improve exploration efficiency for the RL agents for the sparsely-dense reward problem at hand and the pseudocode for the proposed algorithm, Prioritized DDQN-PN Parameter Noise (PDDQN-PN) algorithm is described in Algorithm 1.

## V. EXPERIMENT STUDY

##### A. Dataset

The NASA Commercial Modular Aero-Propulsion System Simulation dataset (C-MAPSS) [14] is widely used in literature and is selected to test our model. The data set is generated from Turbofan Engine Degradation Simulator and contains measurements which mimics the degradation behaviour of multiple turbofan engines under various operating

---

**Algorithm 1** Prioritized Double Deep Q-Learning with Parameter Noise (PDDQN-PN)

---

- 1: **Input:** Action space  $\mathcal{A}$ , mini-batch size  $L_b$ , target network replacement frequency  $L^-$
  - 2: **Output:** Optimal policy  $\pi^*$
  - 3: **Initialize:** Prioritized Experience Replay Memory  $\mathcal{D}_{priority}$  to capacity  $N$ , Primary network  $\mathcal{Q}_\theta$ , target network  $\mathcal{Q}_{\theta^-}$ , action-value function  $\mathcal{Q}$  with random weights, target action-value function  $\hat{\mathcal{Q}}$  with weights  $\theta^- = \theta$
  - 4: **for** Episode=1 to  $E$  **do**
  - 5:   **for** timestep=1 to  $T$  **do**
  - 6:     Observe state  $s_t$  and select  $a_t \sim \pi(a, s)$
  - 7:     With probability  $\epsilon$ , perform random action
  - 8:     Otherwise, choose  $a_t = \operatorname{argmax}_a \mathcal{Q}(s_t, a)$  from  $\mathcal{Q}(s, a; \theta)$ ;
  - 9:     Execute  $a_t$  and receive reward  $r_t$
  - 10:     Observe next state  $s'$
  - 11:     Store tuple  $(s_t, a_t, r_t, s')$  in  $D$  with max priority
  - 12:     Sample random mini-batch transitions, size  $(L_b)$ , from  $\mathcal{D}_{priority}$ , according to transition priority
  - 13:      $y_t^{DDQN} = \begin{cases} r, & \text{if episode terminates at timestep+1} \\ r + \gamma \max_{a'} \hat{\mathcal{Q}}(\phi_{j+1}, a', \theta^-), & \text{else} \end{cases}$
  - 14:     Perform Gradient Descent on  $(y_t^{DDQN} - \mathcal{Q}(s_t, a_t))^2$  and update priority of each transition
  - 15:     Reset  $\theta^- = \theta$  every  $L^-$  steps
  - 16:     Update  $s_t \leftarrow s'$
  - 17:     Increment timestep by 1
  - repeat until** timestep is  $> T$ , terminate
  - repeat until** Episode is  $> E$ , terminate
- 

conditions. The corresponding fault conditions and complex relations with multiple sensor measurements are listed within four similar smaller datasets (FD001~FD004). For simplicity, engine datasets FD001 and FD003 were selected and brief dataset information is shown in Table I. Within each dataset, 26 columns of data are present. Columns 1 and 2 refers to engine unit and particular engine cycle; Columns 3 to 5 are operating conditions, such as temperature; Columns 6 to 26 contains 21 raw sensor readings.

Dataset	FD001	FD003
Training Trajectories	100	100
Testing Trajectories	100	100
Operating Conditions	1	1
Fault Conditions	1	2

TABLE I: C-MAPSS DataSet [14] under test

### B. Data Preparation

Individual sensor values have been normalized for each sensor type (i.e sensor data column):

$$\hat{x} = \frac{(x_i - \mu_i)}{\sigma} \quad (17)$$

where  $\mu_i$  and  $\sigma$  denotes the mean value and standard deviation for the sensor types respectively.  $x_i$  denotes the  $i$ -th sensor value within the corresponding sensor type dataset.

Remaining Useful Life (RUL) calculation is then appended to the normalized data and linear regression is performed.

Once the sensor parameters are identified, an RUL distribution graph is plotted and a Gaussian distribution was observed for all data. Sensor data is then dimensionally reduced using principal component analysis and the equipment health indicator is obtained. In the absence of ground truth for equipment health values, [14] suggested that the degradation behaviour should resemble an exponential decay function like model:

$$H(t) = 1 - d - e^{\{at^b\}} \quad (18)$$

where  $d$  denotes the non-zero initial degradation, and  $a$  and  $b$  are weighted coefficients. Similar degradation behaviour was observed for turbofan equipment health indicator dataset, see Figure 5. Note that in the absence of a ground truth for both repair and replacement action, only the replacement action is tested and a medium-criticality equipment is assumed.

### C. Results

PDDQN-PN was implemented using Tensorflow based deep learning library [15] and Adam optimizer was selected. 5000 time-steps of warm-up was performed to fill the agent's memory, using random policy, before commencing actual training. To balance between exploration and exploitation, the exploration fraction was set to 80% of the number of training simulation time-steps.

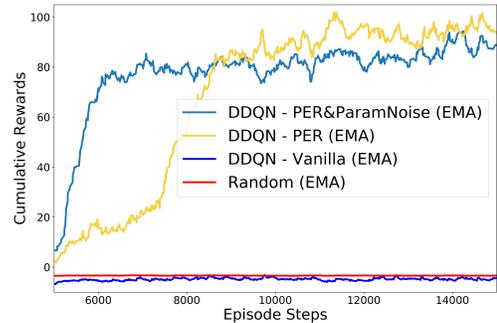


Fig. 4: Learning Performance Benchmark between proposed components and Random policy

In order to demonstrate the benefits of the proposed algorithm, separate tests were conducted on engine #76 sensor data from FD001. As a result, the learning performance contribution factor of each sub-component is clearly illustrated. PER component significantly contributed to the DRL agent's ability to learn an optimal policy within  $1.2 \times 10^4$  time-steps and a mean cumulative reward score of 95. Once stacked with parameter action noise, a notable 50% improvement in learning efficiency is clearly observed with a performance loss of 10%. From Figure 4, it is likely that the mean score for both DDQN-PER and PDDQN-PN will converge given longer simulation time-steps and reduction in exploration time. DDQN-Vanilla performed marginally worse than a random policy given the sparse-reward problem and the results are illustrated in Figure 4. For improved readability, Exponential Mean Average (EMA) smoothing factor of 0.18 was performed on the reported results.

The model prediction method is used to quantify and analyze the DRL agent's learnt decision policy. The prediction results in Table II noted identical median for DDQN based

algorithms with PDDQN-PN reporting the largest standard deviation of  $1.27 \times 10^{-2}$ . The higher deviation is attributed to the action parameter noise technique. The random policy agent clearly demonstrated its inability to learn a suitable policy and DDQN-Vanilla offers conversely comparable prediction performance to the DDQN-PER variant.

Algorithm	Median	Standard Deviation
DDQN - PER + ParamNoise	0.170	0.013
DDQN - PER	0.170	0.011
DDQN - Vanilla	0.170	0.011
Random	0.897	0.0

TABLE II: Average Model Prediction Results for Engine #50

Part of the validation process involves random selection of Engine Health Indicators from both the training and test dataset with results presented in Figure 5. The DRL agent is able to propose a suitable replacement policy, even when the starting point indicates steep downward trend approaching imminent equipment failure, see Figure 5c. This is due to the failure-to-failure penalty constraint  $\mathcal{R}_{Pen}$  in (9). Likewise, the DRL agent demonstrated encouragingly consistent recommended actions, on different engine and dataset. The proposed replacement points are highlighted in Figure 5. For comparison, cross-validation on FD001 reported a median of  $0.175 \pm 0.02$ ; FD003 reported median of  $0.174 \pm 0.02$ .

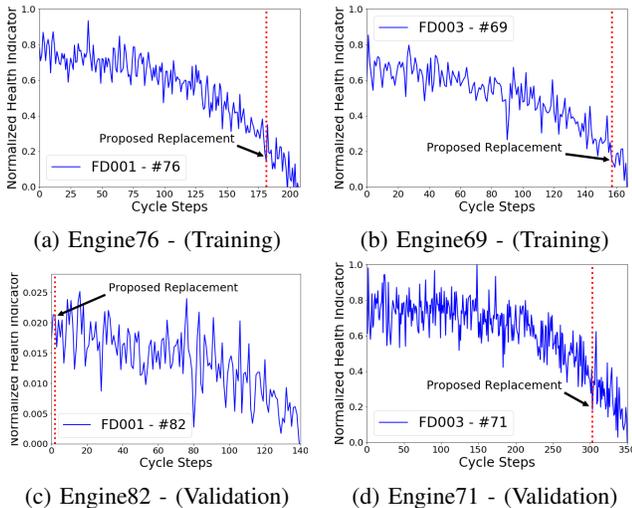


Fig. 5: Example of DRL Agent Proposed Replacement points

## VI. CONCLUSION

As industries worldwide journey towards the Industry 4.0 vision to boost manufacturing productivity, modern equipments become increasingly complex to perform maintenance on. The result is an increased customer demand for accurate, interpretable and actionable insights from predictive maintenance tools. In this paper, we have introduced an approach to provide actionable recommendation, based on an equipment’s health state. We have formulated the maximization of equipment uptime as a function of multiple input sensor data and model the derived equipment health states with state-action pair. The optimal states are observed in a sparsely-dense configuration and is challenging to solve with existing approaches. Then, we have proposed a model-free-based Deep Reinforcement

Learning algorithm which can rapidly learn an optimal maintenance decision policy, for example within 2000 time-steps. The experimental results have shown consistent maintenance recommendations across similar equipment, despite different initial health state. Future work could include extending current work to other equipment failure dataset and benchmark against an actual equipment maintenance policy schedule.

## ACKNOWLEDGEMENT

The work was supported in part by Singapore NRF National Satellite of Excellence, Design Science and Technology for Secure Critical Infrastructure NSoE DeST-SCI2019-0007, A\*STAR-NTU-SUTD Joint Research Grant Call on Artificial Intelligence for the Future of Manufacturing RGANS1906, WASP/NTU M4082187 (4080), Singapore MOE Tier 1 2017-T1-002-007 RG122/17, MOE Tier 2 MOE2014-T2-2-015 ARC4/15, Singapore NRF2015-NRF-ISF001-2277, and Singapore EMA Energy Resilience NRF2017EWT-EP003-041.

## REFERENCES

- [1] C. Idhammar, “What constitutes world-class maintenance and reliability,” <https://www.reliableplant.com/Read/212/world-class-maintenance>, accessed: 2019-08-01.
- [2] G. S. Babu, P. Zhao, and X.-L. Li, “Deep convolutional neural network based regression approach for estimation of remaining useful life,” in *International conference on database systems for advanced applications*. Springer, 2016, pp. 214–228.
- [3] S. Zheng, K. Ristovski, A. Farahat, and C. Gupta, “Long short-term memory network for remaining useful life estimation,” in *2017 IEEE International Conference on Prognostics and Health Management (ICPHM)*. IEEE, 2017, pp. 88–95.
- [4] L. Jayasinghe, T. Samarasinghe, C. Yuen, J. C. N. Low, and S. S. Ge, “Temporal convolutional memory networks for remaining useful life estimation of industrial machinery,” *arXiv preprint arXiv:1810.05644*, 2018.
- [5] E. Sutrisno, H. Oh, A. S. S. Vasan, and M. Pecht, “Estimation of remaining useful life of ball bearings using data driven methodologies,” in *2012 IEEE Conference on Prognostics and Health Management*. IEEE, 2012, pp. 1–7.
- [6] L. Guo, N. Li, F. Jia, Y. Lei, and J. Lin, “A recurrent neural network based health indicator for remaining useful life prediction of bearings,” *Neurocomputing*, vol. 240, pp. 98–109, 2017.
- [7] C. Zhang, C. Gupta, A. Farahat, K. Ristovski, and D. Ghosh, “Equipment health indicator learning using deep reinforcement learning,” in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2018, pp. 488–504.
- [8] Y. Liu and X. Xu, “Industry 4.0 and cloud manufacturing: A comparative analysis,” *Journal of Manufacturing Science and Engineering*, vol. 139, no. 3, p. 034701, 2017.
- [9] S. Ong, K. M. Goh, H.-L. Chan, T.-Y. Lim, and K. V. Ling, “Towards machine diagnostics on chip,” in *2010 11th International Conference on Control Automation Robotics & Vision*. IEEE, 2010, pp. 1353–1358.
- [10] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv preprint arXiv:1312.5602*, 2013.
- [11] H. Van Hasselt, A. Guez, and D. Silver, “Deep reinforcement learning with double q-learning,” in *Thirtieth AAAI conference on artificial intelligence*, 2016.
- [12] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, “Prioritized experience replay,” *arXiv preprint arXiv:1511.05952*, 2015.
- [13] M. Plappert, R. Houthoofd, P. Dhariwal, S. Sidor, R. Y. Chen, X. Chen, T. Asfour, P. Abbeel, and M. Andrychowicz, “Parameter space noise for exploration,” *arXiv preprint arXiv:1706.01905*, 2017.
- [14] A. Saxena, K. Goebel, D. Simon, and N. Eklund, “Damage propagation modeling for aircraft engine run-to-failure simulation,” in *2008 international conference on prognostics and health management*. IEEE, 2008, pp. 1–9.
- [15] A. Hill, A. Raffin, M. Ernestus, A. Gleave, A. Kanervisto, R. Traore, P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, and Y. Wu, “Stable baselines,” <https://github.com/hill-a/stable-baselines>, 2018.