

# Hierarchical Reinforcement Learning for Transportation Infrastructure Maintenance Planning

ZACHARY HAMIDA\*and JAMES-A. GOULET  
Department of Civil, Geologic and Mining Engineering  
POLYTECHNIQUE MONTREAL, CANADA

March 2, 2023

## Abstract

Maintenance planning on bridges commonly faces multiple challenges, mainly related to complexity and scale. Those challenges stem from the large number of structural elements in each bridge in addition to the uncertainties surrounding their health condition, which is monitored using visual inspections at the element-level. Recent developments have relied on deep reinforcement learning (RL) for solving maintenance planning problems, with the aim to minimize the long-term costs. Nonetheless, existing RL based solutions have adopted approaches that often lacked the capacity to scale due to the inherently large state and action spaces. The aim of this paper is to introduce a hierarchical RL formulation for maintenance planning, which naturally adapts to the hierarchy of information and decisions in infrastructure. The hierarchical formulation enables decomposing large state and action spaces into smaller ones, by relying on state and temporal abstraction. An additional contribution from this paper is the development of an open-source RL environment that uses state-space models (SSM) to describe the propagation of the deterioration condition and speed over time. The functionality of this new environment is demonstrated by solving maintenance planning problems at the element-level, and the bridge-level.

**Keywords:** Maintenance Planning, Reinforcement Learning, RL Environment, Deep Q-Learning, Infrastructure Deterioration, State-Space Models

## 1 Introduction

Transportation infrastructure such as roads, tunnels and bridges are continuously deteriorating due to aging, usage and other external factors [6]. Accordingly, maintenance planning for the aforementioned infrastructure aims at minimizing maintenance costs, while sustaining a safe and functional state for each structure [5, 33]. Maintenance strategies for bridges can be either classified as time-based maintenance, such as recurring maintenance actions based on a fixed time interval, or condition-based maintenance (CBM) [25]. In the context of CBM, the main components involved in the development of any maintenance policy are quantitative measures for, 1) the structural health condition, 2) the effects of interventions, and 3) costs of maintenance actions. The structural health of bridges is commonly evaluated using visual inspections at the element-level [12, 22, 4]. An example of an element in this context is the *pavement* in a concrete bridge. The information at the element-level are thereafter aggregated to provide a representation for the overall deterioration state of a bridge [15]. Similarly, maintenance actions are performed at the element-level, and their corresponding effect is aggregated at the bridge-level [15, 14]. The hierarchical nature of condition assessments, and maintenance actions presents challenges in formulating the bridge maintenance planning problem. First, the aggregation of the health states from the element-level to the bridge-level results in additional uncertainties, which render deterministic deterioration models insufficient [15]. Second, performing actions at the element-level implies that a decision-making framework is required to search for maintenance policies at the element-level in each bridge. Thus, the search-space for an optimal maintenance policy is typically

---

\*Corresponding author: zachary.hamida@polymtl.ca

large as it is common for a bridge to have hundreds of structural elements [23].

Existing approaches for solving the maintenance planning problem have adopted Markov decision process (MDP) formulations [5, 33, 8, 21], relying on discrete states where transitioning from one state to another depends only on the current state [29]. The MDP approach is well-suited for small state-space problems, so that using MDP in the context of maintenance planning have incurred simplifications on the state and the action space [5, 9]. An example of simplification is reducing the search space by merging the state representation of structural elements with similar deterioration states and maintenance actions [9].

The large state space has also motivated the application of reinforcement learning (RL) methods to search for optimal maintenance policies [34, 5, 35]. Conventional RL methods are well-suited for discrete action and state spaces, where the agent (or decision-maker) performs different actions and receives a feedback (rewards), which can be used to update the value function corresponding to the visited states [29]. Existing work in the context of maintenance planning have focused mainly on multi-agent RL (MARL) methods, due to their compatibility with large action spaces [6, 36, 20, 1]. Applications include coordinated reinforcement learning (CRL) for joint decision-making of multi-agents [20], and hierarchical RL where higher-level agents moderate the behaviour of lower-level agents [1, 16]. This latter application has been further improved by combining the CRL with the hierarchical reinforcement learning framework [36]. Despite the MARL extension, RL methods are inherently limited to low-dimensional problems and lack the capacity to scale for large state spaces without increasing the number of agents acting on the state space [5]. Accordingly, deep reinforcement learning (DRL) and multi-agent DRL have been proposed as an alternative due to their capacity of handling large and continuous state and action spaces. Specifically, frameworks with centralized training such as the branching dueling Q network (BDQN) and the deep centralized multi-agent actor critic (DCMAC) [30, 5]. A common limitation associated with the aforementioned frameworks is the stability of the training, especially as the size of the action space increases. In addition, the policy obtained is not interpretable, such that it is not possible to plot the decision boundaries for the policy, and with the lack of a clear stopping criteria for training the agents in the context of planning problems, it becomes difficult to evaluate the validity of the reward function or the policy in practical applications. Another common limitation in the context of maintenance planning for transportation infrastructure is the use of discrete Markov models (DMM) for modelling the deterioration process over time. The use of the DMM framework in this context induces drawbacks related to overlooking the uncertainty associated with each inspector, and the incapacity to estimate the deterioration speed [12].

The aim of this paper is to introduce a hierarchical reinforcement learning formulation that adapts to the hierarchical nature of information in the maintenance planning problem. The hierarchical formulation enables decomposing large state and action spaces into smaller ones, by relying on state and temporal abstraction [2, 3]. State abstraction enables representing the state-space of the planning problem by a hierarchy of states, such as, the element-level, the structural category level, and the bridge-level. Each of the aforementioned levels has an action-space, where interdependent policies can be learned and applied. Take for example, a bridge-level decision, where the action-space is defined as maintain or do nothing; if a policy suggests doing nothing, then no intervention is applied on all elements within the bridge, without assessing their health states.

The main contributions in this paper are: 1) formulating a hierarchical deep reinforcement learning approach that adapts to bridge maintenance planning, and provides advantages in scalability, and interpretability through visualizing the decision boundaries of the policies. 2) Incorporating the deterioration speed alongside the deterioration condition in the decision-making analyses [12]. 3) Developing a standard gym-based RL environment [7] for emulating the deterioration process of bridges, based on state-space models (SSM) [15, 13].

The performance of the proposed hierarchical approach is demonstrated using an example application for a bridge from the network of bridges in the province of Quebec, Canada. The analyses include a comparison with the BDQN approach for planning maintenance on a multi-component system, in addition to learning a bridge-level maintenance policies.

## 1.1 Problem Formulation

A bridge  $\mathcal{B}$  from the network of bridges in the Quebec province is considered to demonstrate the decision-making analyses presented in this paper. Figure 1 provides a summary for the hierarchy of components and information in bridge  $\mathcal{B}$ . The bridge is composed of  $K$  structural categories, each of which is composed of  $E$  structural elements. An example of a structural category is the *beams* category, which is the  $k$ -th category in bridge  $\mathcal{B}$  with a total of  $P$  beams as in,  $\mathcal{C}_k = \{e_1^k, \dots, e_p^k, \dots, e_E^k\}$ . Visual inspections are performed on the elements of bridge  $\mathcal{B}$  every three years to monitor their health condition, represented by  $\tilde{y}_{t,p}^k$ . The health states of the elements are inferred using the inspection data  $\tilde{y}_{t,p}^k$ , and are denoted by  $\tilde{x}_{t,p}^k$ . The element-level health states  $\tilde{x}_{t,p}^k$  are thereafter aggregated for each structural category to provide the overall health states of the structural categories  $\tilde{x}_{t,k}^c$ . Similarly, the overall health states of the bridge  $\tilde{x}_t^b$  are based on the aggregation of the health states from the structural categories [14]. It should be noted that for all the aforementioned levels, the health states are described by the same condition range defined by,  $\tilde{x}_{t,p}^k, \tilde{x}_{t,k}^c, \tilde{x}_t^b \in [l, u]$ , and the deterioration speed which is defined in  $\mathbb{R}^-$ . The  $\sim$  in  $\tilde{x}_{p,t}^k$  refers to variables within the bounded space  $[l, u] = [25, 100]$ , while the absence of  $\sim$  refers to the variables defined in the unbounded space  $[-\infty, \infty]$  [14]. An example of a perfect health state is when the condition is,  $\tilde{x}_t = 100$ , and the deterioration speed is near-zero. Bridge  $\mathcal{B}$  mainly undergo imperfect maintenance actions at the element-level, where the actions are represented by the set  $\mathcal{A}^e = \{a_0, a_1, a_2, a_3, a_4\}$ , with  $a_0$ : do nothing,  $a_1$ : routine maintenance,  $a_2$ : preventive maintenance,  $a_3$ : repair, and  $a_4$ : replace [23]. Each action is associated with a cost, in addition to other costs related to the service interruption and penalties for reaching a critical state.

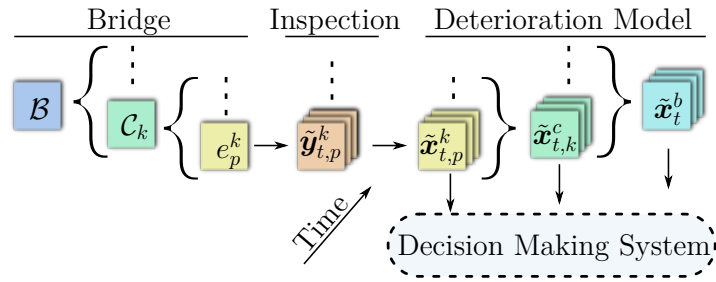


Figure 1: Hierarchy of components and information in a bridge  $\mathcal{B}$ , where each structural category  $\mathcal{C}_k$  is composed of a number of structural elements  $e_p^k$ . The element-level inspection data  $\tilde{y}_{t,p}^k$  provides information about the health states at the element-level  $\tilde{x}_{t,p}^k$ , structural category level  $\tilde{x}_{t,k}^c$ , and bridge-level  $\tilde{x}_t^b$ .

## 2 Background

This section provides the theoretical background for the main concepts related to proposed decision making framework.

### 2.1 Markov Decision Processes (MDP)

A MDP is an approach to describe sequential decision-making problems using the tuple  $\langle \mathcal{S}, \mathcal{A}, P, \mathcal{R} \rangle$ , where  $\mathcal{S}$  is the set of states,  $\mathcal{A}$  is the set of actions,  $P$  is the transition function, and  $\mathcal{R}$  is the set of rewards [29]. Taking an action  $a \in \mathcal{A}$  results in a transition from the state  $s_t = s$  at time  $t$ , to the state  $s_{t+1} = s'$  using a Markovian transition function  $P(s'|s, a)$ , which implies that the next state is only conditional on the pair of the current state  $s$  and action  $a$ . Each action  $a \in \mathcal{A}$  taken in the MDP can affect the expected immediate reward  $r_t$  and the total rewards  $G_t$  [29]. In this context, the effect of an action can be either deterministic and accordingly the MDP is considered deterministic where,  $\Pr(s'|s, a) = 1$ , or otherwise, the MDP is considered stochastic when  $\Pr(s'|s, a) \neq 1$  [26]. Similarly, states can be either represented by deterministic exact information (i.e., true state  $s$ ) in the case of a MDP, or inferred information (i.e., belief about the true state  $s$ ) in the case of a partially observed MDP (POMDP) [29].

A policy  $\pi$  in the context of MDP represents a mapping between states and actions, where a deterministic policy provides deterministic actions such that,  $\pi(s) : s \rightarrow a$ , while a stochastic policy provides probabilities for taking each action in each state  $\pi(a|s) : s \rightarrow \text{Pr}(a)$  [29]. Following a policy  $\pi(\cdot)$  in a MDP would yield a total return at time  $t$  defined by,

$$G_t = \sum_{i=0}^{\infty} \gamma^i r(s_{t+i}, a_{t+i}), \quad (1)$$

where  $\gamma$  is the discount factor that enables considering an infinite planning horizon when  $\gamma \in ]0, 1[$ , and  $r(s_t, a_t) = \mathbb{E}[R_t | S_t = s, A_t = a]$  denotes the expected reward given the state  $S_t = s$  and the action  $A_t = a$  [29]. In this context, the rewards represent a feedback associated with action  $a$  taken at each state  $s$ . Provided the notion of rewards driven by actions, evaluating a policy  $\pi$  is possible by using a value function  $V_\pi(s)$  and an action-value function  $Q_\pi(s, a)$ . The value function represents the expected discounted return for being in a state  $s$ , under policy  $\pi$ , such that,

$$V_\pi(s) = \mathbb{E}_\pi [G_t | S_t = s] = r(s_t, a_t) + \mathbb{E}_\pi \left[ \sum_{i=1}^{\infty} \gamma^i r(s_{t+i}, a_{t+i}) | S_t = s \right], \quad (2)$$

where  $\mathbb{E}_\pi$  is the expected value while following the policy  $\pi$ . On the other hand, the action-value function  $Q_\pi(s, a)$  refers to the expected discounted return for taking an action  $a$ , in a state  $s$ , based on policy  $\pi$ , which is described by,

$$Q_\pi(s, a) = r(s_t, a_t) + \mathbb{E}_\pi \left[ \sum_{i=1}^{\infty} \gamma^i r(s_{t+i}, a_{t+i}) | S_t = s, A_t = a \right]. \quad (3)$$

Accordingly, a policy  $\pi_*$  is considered optimal when the state-value function and action-value function are,

$$\begin{aligned} V_*(s_t) &= \max_{\pi} V_\pi(s_t), \quad \forall s_t \in \mathcal{S}, \\ Q_*(s_t, a_t) &= \max_{\pi} Q_\pi(s_t, a_t), \quad \forall s_t \in \mathcal{S}, a_t \in \mathcal{A}. \end{aligned} \quad (4)$$

## 2.2 Semi-Markov Decision Process

A semi-Markov decision process (SMDP) formulation is similar to a MDP, with the exception that a SMDP considers actions to have a duration  $\bar{T}$  to be performed [26]. An example for a SMDP action is the task of maintaining a bridge, which requires a duration  $\bar{T}$ , to perform the maintenance actions for each element within the bridge. From this example it can be inferred that actions (or tasks) in the SMDP are performed at different levels (i.e., element-level and bridge-level). The expected rewards  $\bar{r}(s_t, a_t^\ell)$  associated with the task  $a_t^\ell$  at level  $\ell$  are estimated using,

$$\bar{r}(s_t, a_t^\ell) = \mathbb{E}_{\pi^{\ell-1}} \left[ \sum_{i=0}^{\bar{T}} \gamma^i r(s_{t+i+1}, a_{t+i+1}^{\ell-1}) | S_t = s, a_t^{\ell-1} = \pi^{\ell-1}(s_t) \right], \quad (5)$$

where  $\bar{r}(s_t, a_t^\ell)$  is the expected cumulative discounted reward while following the policy  $\pi^{\ell-1}$  from time  $t$  until the termination of the task  $a_t^\ell$  after  $\bar{T}$  time-steps. Based on Equation 11, the application of the SMDP formulation generally relies on state and temporal abstractions [26, 2, 3]. The aim of state abstraction is to reduce the state space by aggregating states having similar properties without changing the essence of the problem [2, 3]. This implies the feasibility of mapping a state  $s \in \mathcal{S}$  to an abstract state  $s_\phi \in \mathcal{S}_\phi$  while maintaining a near-optimal policy search, where the space  $\mathcal{S}_\phi$  has a fewer states (i.e.,  $|\mathcal{S}_\phi| \ll |\mathcal{S}|$ ) [2]. Figure 2 shows an illustrative example, where the real state is represented by different levels of abstraction. On the other hand, a temporal abstraction is applied when actions are taking place at different time scales [29]. For example, applying an intervention on a bridge  $b_j$  from time  $t$  to time  $t + 1$ , involves many actions performed at the element-level over a sub-timestamps  $\tau$ , such that,  $t < (t + \tau) < t + 1$ .

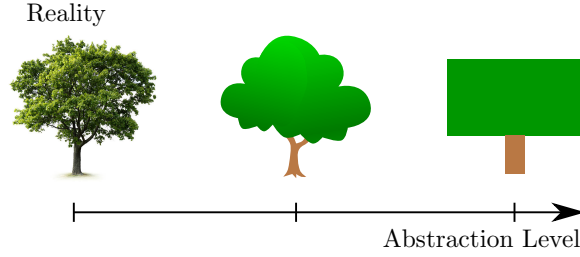


Figure 2: Illustrative example showing two levels of abstraction starting from the state of reality on the left.

### 2.3 Deep Reinforcement Learning

Typical RL approaches rely on interactions between a decision maker (the agent) and an environment in order to learn a policy that maximizes the total cumulative rewards. A common technique for learning from interactions is the use of the temporal difference (TD) [32], to perform recursive updates on the action-value function  $Q(s_t, a_t)$  such that,

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \eta \left[ r(s_t, a_t) + \gamma \max_{a_{t+1}} Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right], \quad (6)$$

where  $\eta$  denotes the learning rate. Updating the  $Q(s_t, a_t)$  function using Equation 6 requires a table for all pairs of states and actions, which can be challenging for large and continuous state and action spaces [29]. Therefore, Deep RL methods have provided a scalable alternative to the tabular Q-learning, which enables approximating the optimal action-value function such that,  $Q(s, a; \theta) \approx Q^*(s, a)$ . Similar to Equation 6, deep RL relies on temporal difference (TD) to estimate the set of parameters  $\theta$  using,

$$\mathcal{L}_i(\theta_i) = \mathbb{E} \left[ r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta_i) \right]^2, \quad (7)$$

where  $\mathcal{L}_i(\theta_i)$  is the loss function associated with the parameters  $\theta_i$ , and  $\theta^-$  represents the parameters of the target model, which is a delayed replica of the  $Q(\cdot)$  function. The role of a target model is to stabilize the learning process where the parameters of the target model  $\theta^-$  are updated based on  $\theta_i$  by using a soft update approach [19]. Equation 7 is the foundation for many different DRL methods for identifying an optimal policy, nonetheless, the choice of an approach mainly depends on the design and properties of the MDP [29].

### 2.4 Hierarchical Reinforcement Learning

Hierarchical RL enables applying the principles of RL on SMDP environments, where there are multiple tasks occurring simultaneously at different time scales [26]. The hierarchy here refers to multiple layers of policies, where the higher level policy dictates the behaviour of the lower level policies [24]. For example, the higher level policy observes the overall deterioration state of a bridge  $\tilde{x}_t^b$  and provides a target state  $\tilde{x}_t^b + \delta^b$ , which the lower level policies will try to match by observing and acting on the deterioration states of the structural elements. Based on the definition above and Equation 5, the action-value function  $Q(s_t, a_t^\ell)$  for an optimal policy can be defined as,

$$Q(s_t, a_t^\ell) = \bar{r}(s_t, a_t^\ell) + \sum_{s_{t+\bar{T}}} \sum_{\bar{T}} \gamma^{\bar{T}} P(s_{t+\bar{T}}, \bar{T} | s_t, a_t^\ell) \max_{a_{t+\bar{T}}^\ell} Q(s_{t+\bar{T}}, a_{t+\bar{T}}^\ell). \quad (8)$$

From Equation 8, the transition model  $P(s_{t+\bar{T}}, \bar{T} | s_t, a_t^\ell)$  and the reward  $\bar{r}(s_t, a_t^\ell)$  depend directly on the subsequent policy  $\pi^{\ell-1}$  [26].

Learning the hierarchical policies can be done by using either an end-to-end approach where all policies are trained simultaneously, or a bottom-to-top approach starting from the lower level policies [26, 11]. The latter approach is favoured for large-scale problems provided the instability issues for centralized joint training of multiple policies [11].

### 3 Hierarchical Deep RL for Bridge Maintenance Planning

Figure 3 shows an illustration for the hierarchical maintenance planning architecture, where the state of the environment at time  $t$  is represented using different levels: a bridge level with state  $\mathbf{s}_t^b$ , a structural-category level with  $\mathbf{s}_{t,k}^c$ , and an element-level with  $\mathbf{s}_{t,p}^e$ . Each of the aforementioned states provide information about the health of the bridge at its corresponding level. For example, the state of each element  $\mathbf{s}_{t,p}^e$  contains information about the deterioration condition  $\tilde{x}_{t,p}^k$  and speed  $\tilde{\dot{x}}_{t,p}^k$  of the  $p$ -th structural element  $e_p^k$ .

The hierarchical framework is composed of a centralized agent for the bridge level with policy  $\pi^b$ , and decentralized agents for each structural category represented by the policy  $\pi_k$ . The centralized agent proposes a target improvement  $\delta^b \leftarrow \pi^b(\mathbf{s}_t^b)$  for the health condition of the bridge  $x_t^b$ , such that the health condition of the bridge at time  $t + 1$  is  $x_{t+1}^b + \delta^b$ . If the improvement value  $\delta^b = 0$ , then no maintenance is applied on the bridge; otherwise, maintenance actions are performed according to the improvement value  $\delta^b$ , defined within,  $\delta^b \in [0, (u - l)]$ , where  $l$  is the lower bound and  $u$  is the upper bound for the condition.

As shown in Figure 3, the hierarchical framework aims to decode the bridge-level target improvement  $\delta^b$  to a vector of actions for all structural elements in  $\mathcal{B}$ . This can be achieved sequentially by distributing  $\delta^b$  on the structural categories according to their current deterioration condition  $\tilde{x}_{t,k}^c$  using,

$$\delta_k^c(\delta^b) = \frac{u - \tilde{x}_{t,k}^c}{u \cdot K - \sum_{k=1}^K \tilde{x}_{t,k}^c} \cdot K \cdot \delta^b, \quad (9)$$

where  $\delta_k^c$  is the target improvement for the  $k$ -th structural category  $\mathcal{C}_k$ ,  $K$  is the total number of structural categories within the bridge, and  $u$  is the perfect condition. From Equation 9, if  $\delta_k^c > 0$ , then the structural element  $e_p^k \in \mathcal{C}_k$  is maintained according to the policy  $\pi_k$ . Thereafter, the states of the structural category  $\tilde{\mathbf{s}}_{t,k}^c$ , and the bridge  $\tilde{\mathbf{s}}_t^b$  are updated with the state after taking the maintenance action  $a_{t,p}^k \leftarrow \pi_k(\mathbf{s}_{t,p}^e)$  on the structural element  $e_p^k$ . In order to determine if the next structural element  $p + 1$  requires maintenance, the target improvement  $\delta_k^c$  and  $\delta^b$  are updated using,

$$\begin{aligned} \delta_k^c &= \max\left(\tilde{x}_{t,k}^c(\text{before maintenance}) + \delta_k^c - \tilde{x}_{t,k}^c(\text{updated}), 0\right), \\ \delta^b &= \max\left(\tilde{x}_t^b(\text{before maintenance}) + \delta^b - \tilde{x}_t^b(\text{updated}), 0\right). \end{aligned} \quad (10)$$

Once the updated target improvement  $\delta_k^c$  reaches  $\delta_k^c = 0$ , the remaining structural elements within  $\mathcal{C}_k$  are assigned the action ( $a_0$ : *do nothing*). The aforementioned steps are repeated for each structural category  $\mathcal{C}_k$  in bridge  $\mathcal{B}$  until all elements  $e_p^k$  are assigned a maintenance action  $a_p^k \in \mathcal{A}^e$ .

The element-level actions are defined by the set  $\mathcal{A}^e = \{a_0, a_1, a_2, a_3, a_4\}$ , where  $a_0$ : do nothing,  $a_1$ : routine maintenance,  $a_2$ : preventive maintenance,  $a_3$ : repair, and  $a_4$ : replace [23]. The corresponding effect associated with each of the aforementioned actions is estimated using a data-driven approach [14]. Moreover, the cost associated with each element-level maintenance action is defined as a function of the deterioration state of the structural element, and for each structural category. Further details about the effect of interventions and maintenance costs are provided in B.3.

In addition to the maintenance action costs, there are costs related to the bridge service-stoppage and penalties for reaching a critical state. The service-stoppage costs are defined to prevent frequent interruptions for the bridge service, as well as to encourage performing all of the required maintenance actions at the same time. On the other hand, the penalties are applied when a predefined critical state is reached and no maintenance action is taken. The critical state in this work is defined in accordance with the definition provided by the *Manual of Inspections* [23], for a deterioration state that requires maintenance.

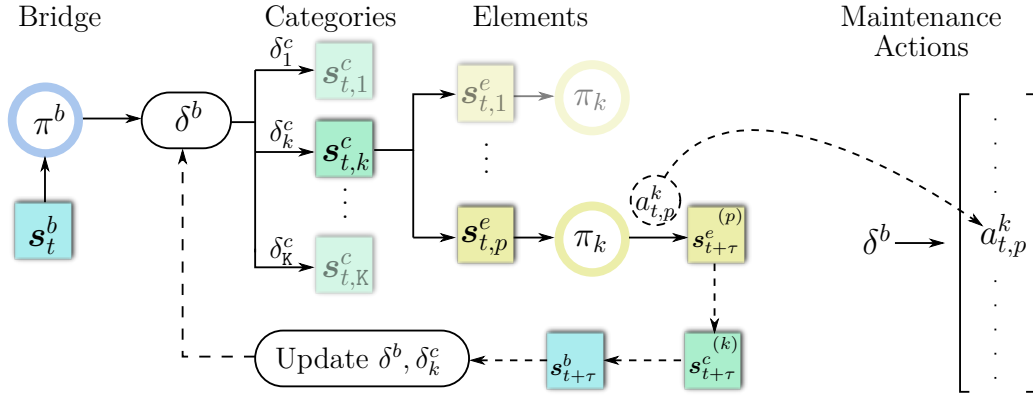


Figure 3: Hierarchical deep RL for performing maintenance using a hierarchy of policies composed of, a centralized policy  $\pi^b$  for the bridge level, and decentralized element-level policies  $\pi_k$ . The centralized policy  $\pi^b$  produces a target improvement  $\delta^b$  based on the bridge state  $s_t^b$ . The improvement  $\delta^b$  is distributed on the structural categories to provide the category-wise improvements  $\delta_k^c$ , which are sequentially translated to a vector of maintenance actions at the element-level using the policies  $\pi_k$ .

### 3.1 Learning the Policies in the Hierarchical DRL

Learning the policies in the hierarchical DRL framework is done using a bottom-to-top approach starting from the element-level policies, and by relying on decentralized element-level agents with a centralized bridge-level agent [26, 11]. Such an approach offers flexibility in using transfer learning for structural elements that share similar properties. In this context, structural elements from the same structural category (e.g., all the beams) are assumed to have a similar deterioration dynamics and similar cost function for maintenance actions. Therefore, the number of element-level agents that require training is equivalent to the number of structural categories in bridge  $\mathcal{B}$ .

Training the element-level agents is done based on a MDP environment that mimics the deterioration process and provides information about the deterioration condition  $\tilde{x}_{t,p}^k$  and speed  $\tilde{x}_{t,p}^k$  of the structural elements (see Section 3.2). Accordingly, the state space for the element level is,  $s_t^e = [\tilde{x}_{t,p}^k, \tilde{x}_{t,p}^k]$  and the action space is defined by the set  $\mathcal{A}^e$ . Training the element-level agents can be done using off-policy methods, such as deep Q-learning with experience replay [29].

After learning the policies  $\pi_{1:K}$ , it becomes possible to learn the centralized bridge-level policy, which observes the state  $s_t^b = [\tilde{x}_t^b, \tilde{x}_t^b, \sigma_t^b]$ , where  $\tilde{x}_t^b$  is the overall health condition of the bridge,  $\tilde{x}_t^b$  is the overall deterioration speed of the bridge, and  $\sigma_t^b$  is the standard deviation for the condition of the structural categories in the bridge  $\sigma_t^b = \text{std}(\tilde{x}_{t,1:K}^c)$ . The environment at the bridge level is a SMDP due to the assumption that all element level maintenance actions are occurring between the time steps  $t$  and  $t + 1$ . Training the centralized agent is done using an off-policy deep Q-learning approach with experience replay. The bridge-level agent experience transition is composed of,  $(s_t^b, \delta_t^b, r_t^b, s_{t+1}^b)$ , where  $r_t^b$  is the total costs from all actions performed on the bridge and is defined by,

$$r^b(s_t^b, \delta_t^b) = r^s + \sum_{k=1}^K \sum_{p=1}^P r(s_{t,p}^e, a_{t,p}^k). \quad (11)$$

From Equation 11,  $r^s$  is the service-stoppage cost for performing the maintenance actions. The next section describes the environment utilized for emulating the deterioration of bridges over time.

### 3.2 Deterioration State Transition

The RL environment is built based on the deterioration and intervention framework developed by Hamida and Goulet [15, 14], and is calibrated using the inspections and intervention database for the network of bridges in the Quebec province, Canada. The environment emulates the deterioration

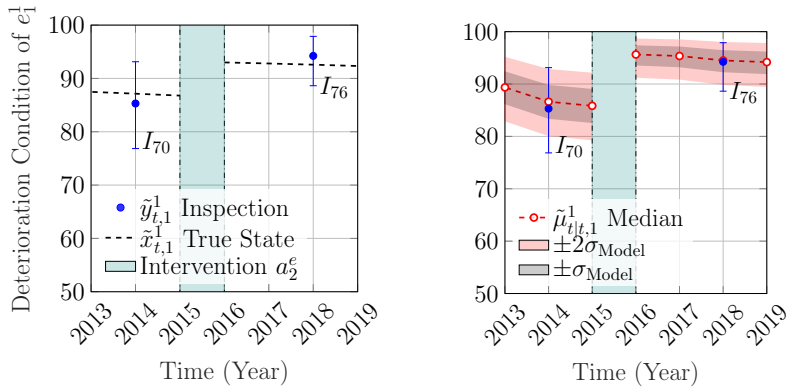
process by generating true states for all the elements  $e_p^k$ , using the transition model,

$$\overbrace{\mathbf{x}_{t,p}^k = \mathbf{A}_t \mathbf{x}_{t-1,p}^k + \mathbf{w}_t}^{\text{transition model}}, \underbrace{\mathbf{w}_t : \mathbf{W} \sim \mathcal{N}(\mathbf{w}; \mathbf{0}, \mathbf{Q}_t)}_{\text{process errors}}, \quad (12)$$

where  $\mathbf{x}_{t,p}^k : \mathbf{X} \sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$  is a hidden state vector at time  $t$  associated with the element  $e_p^k$ . The hidden state vector  $\mathbf{x}_{t,p}^k$  is a concatenation of the states that represent, the deterioration condition  $x_{t,p}^k$ , speed  $\dot{x}_{t,p}^k$ , and acceleration  $\ddot{x}_{t,p}^k$ , as well as the improvement due to interventions represented by, the change in the condition  $\delta_{t,p}^e$ , the speed  $\dot{\delta}_{t,p}^e$ , and the acceleration  $\ddot{\delta}_{t,p}^e$ .  $\mathbf{A}_t$  is the state transition matrix, and  $\mathbf{w}_t$  is the process error with covariance matrix  $\mathbf{Q}_t$ . Equation 12 represents the dynamics of a transition between the states in the context of a MDP. In order to emulate the uncertainties about the deterioration state, synthetic inspection data  $y_{p,t}^k$  are sampled at a predefined inspection interval using,

$$\overbrace{y_{t,p}^k = \mathbf{C} \mathbf{x}_{t,p}^k + v_t}^{\text{observation model}}, \underbrace{v_t : V \sim \mathcal{N}(v; \mu_V(I_i), \sigma_V^2(I_i))}_{\text{observation errors}}, \quad (13)$$

where  $\mathbf{C}$  is the observation matrix, and  $v_t : V \sim \mathcal{N}(v; \mu_V(I_i), \sigma_V^2(I_i))$ , is the observation error associated with each synthetic inspector  $I_i \in \mathcal{I}$ . The role of the synthetic inspection data is to provide imperfect measurements similar to the real world context. The information from this measurement at time  $t$  can be extracted using the Kalman Filter (KF), where the state estimates from the KF at each time  $t$  represent a belief about the deterioration state [17]. The belief states from the KF provide a POMDP representation of the deterioration process. Figure 4 shows an illustrative example for the deterioration and effect of interventions on a structural element, where the true state  $\tilde{x}_{t,1}^1$  is represented by the black dashed line, and the expected value  $\tilde{\mu}_{t,1}^1$  is represented by the red dashed line with the confidence region  $\pm\sigma_{t|t}$  and  $\pm 2\sigma_{t|t}$  for the condition.



(a) Example for the true state of deterioration generated using the transition model with inspections generated using the observation model, and an intervention at year 2016.

(b) Example for the KF forward inference based on the synthetic inspection data, with an intervention at year 2016.

Figure 4: Illustrative example for a deterministic deterioration curve (MDP environment) in Figure 4a, and an uncertain deterioration curve (POMDP environment) in Figure 4b.

The deterioration states  $\tilde{x}_{t,p}^k$  at the element-level are aggregated to obtain the overall deterioration state of the structural category  $\mathbf{x}_{t,k}^c$ , which are similarly aggregated to obtain the deterioration states estimates of the bridge  $\mathbf{x}_t^b$ . Further details about the aggregation procedure as well as the deterioration and interventions framework are provided in B.

Throughout the deterioration process, the effectiveness of repair actions is distinguished from the replacement action by introducing a decaying factor on the perfect state  $u_t$ , such that,  $u_{p,t+1} = \rho_0 \times u_{p,t}$ , where  $0 < \rho_0 < 1$ . This implies that repair actions are unable to restore a structural element to the



original perfect health condition (i.e.,  $u_t = 100$ ) as it advances in age. Figure 6a shows an illustration for the decay in the perfect condition  $u_t$  that can be reached by repair actions. Other practical considerations in this environment are related to preventing the observation of the same effect after applying the same repair action multiple times within a short period of time (e.g.,  $\Delta t \leq 2$ ). For example, if an action  $a_3$  has an effect of  $\delta_t^e = +20$ , applying  $a_3$  action two times in two consecutive years should not improve the structural element condition by,  $\delta_t^e + \delta_{t+1}^e = 20 + 20$ , but rather should improve the condition by,  $\delta_t^e + \rho_1 \delta_{t+1}^e$ , where  $0 < \rho_1 < 1$ . Figure 6b illustrates the decayed effect of intervention caused by applying the same action twice within a short time interval. Further details about the choice of decaying factors in this context are provided in B.2.

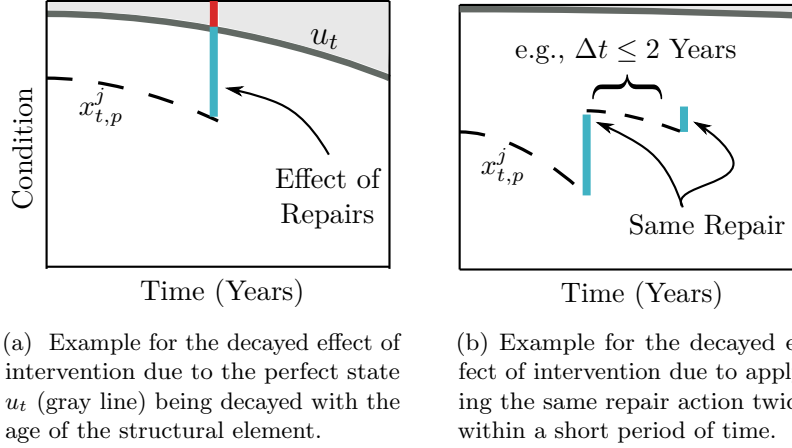


Figure 5: Decaying factors for the perfect state  $u$  and the effect of interventions over time.

## 4 Example of Application

The performance of the proposed HRL framework is demonstrated on a case study for a bridge within the province of Quebec. Note that both the deterioration and interventions models are calibrated on data from the bridge network in the Quebec province, Canada [15].

### 4.1 Maintenance Policy for a Bridge with One Structural Category

The goal in this example is to demonstrate the capacity of the HRL to achieve a near-optimal solution in a toy-problem, with a simple hierarchy of actions. In this context, the planning scope on bridge  $\mathcal{B}$  considers only one structural category  $K = 1$ , which corresponds to the *beams* structural category  $\mathcal{C}_1$ . The beam elements in  $\mathcal{C}_1$  have a common critical deterioration condition  $\tilde{x}_t$  and deterioration speed  $\tilde{x}_t$  defined as,  $\tilde{x}_t = 55, \dot{\tilde{x}}_t = -1.5$ . The aforementioned values are derived from the *manual of inspections* [23], and imply that the structural element requires a maintenance action when the critical state is reached; accordingly, taking no-action after reaching the critical state will incur a cost penalty on the decision-maker.

As described in Section 3.1, the first step to train the proposed hierarchical framework is to learn the task of maintaining beam structural elements at the element level. The policy  $\pi_k$  decides the type of maintenance actions based on the information about a deterministic deterioration condition  $\tilde{x}_{t,p}^1$  and a deterministic deterioration speed  $\tilde{x}_{t,p}^1$  of the structural element such that,  $\mathbf{s}_{t,p}^e = [\tilde{x}_{t,p}^k, \tilde{x}_{t,p}^k]$ . The action set in this MDP is defined as,  $\mathcal{A}^e = \{a_0, a_1, a_2, a_3, a_4\}$ , which corresponds to  $a_0$ : do nothing,  $a_1$ : routine maintenance,  $a_2$ : preventive maintenance,  $a_3$ : repair, and  $a_4$ : replace [23]. The costs and effects associated with each of the aforementioned actions are described in B.3. Learning the maintenance policy  $\pi_k$  is done using a vectorized version of the RL environment which is detailed in Section 3.2 and A. The experimental setup for the training include a total of  $5 \times 10^4$  episodes with the episode length defined as,  $T = 100$  years. The episode length is determined such that it is long enough to necessitate a replacement action, provided that the average life-span for a structural element is about 60 years [14]. Despite the fixed episode length, there is no terminal state as the planning horizon is

considered infinite with a discount factor  $\gamma = 0.99$ . Moreover, the initial state in the RL environment is randomized where it is possible for a structural element to start the episode in a poor health state or a perfect health state. Figure 6 shows the training and average performance for DQN and Dueling agents, along with two realizations for the optimal policy map obtained at the end of the training for each agent. The configuration for the DRL agents are provided in A. From Figure 6a, it is noticeable

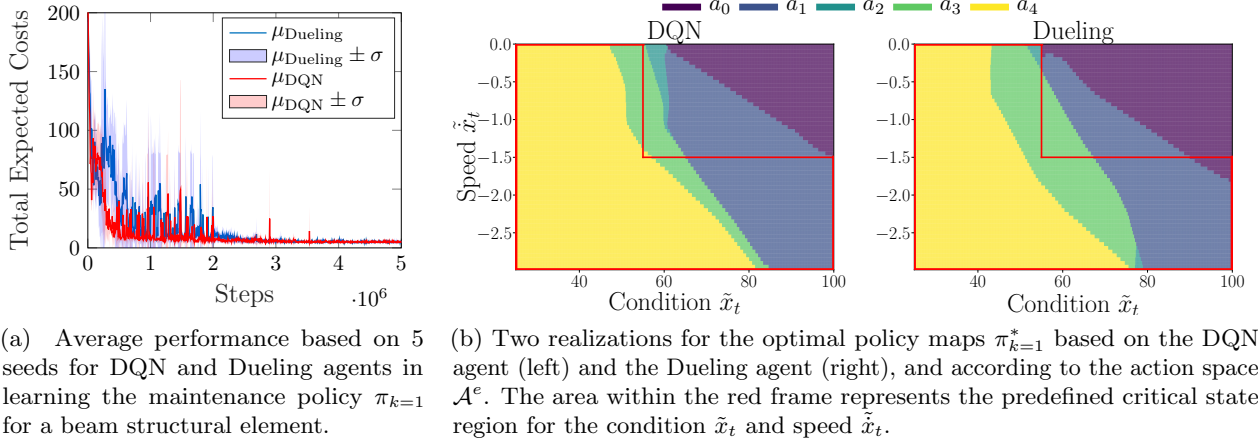


Figure 6: The training process of deep RL agents along with two realizations for the optimal policy  $\pi_{k=1}^*$  of a beam structural element.

that the DRL agents reach a stable policy after  $3 \times 10^6$  steps. Moreover, Figure 6b shows optimal policy maps obtained by the DQN agent (left) and the Dueling agent (right), for the action space  $\mathcal{A}^e$ , with the critical state region highlighted by the area within the red boundary. From the policy maps, it can be noticed that the element's critical state region is dominated by major repairs, which is expected due to the penalties applied on the DRL agent if the structural element reaches that state. Despite the slight differences between the two optimal policy maps in Figure 6b, the policy map by the DQN agent is favourable because the action  $a_0$  : *do nothing* did not leak into the predefined critical state region for the condition  $\tilde{x}_t$  and speed  $\tilde{x}_t$ . The leakage of the action  $a_0$  : *do nothing* in the critical state region can occur due to interpolating the  $Q$  function values for states that are rarely visited by the agent, such as structural elements with a perfect condition  $\tilde{x}_t = 100$ , and high deterioration speed  $\tilde{x}_t = -1.6$ . The optimal policy  $\pi_k^*$  provides the basis for decision-making at the bridge level, which corresponds to learning the bridge level policy  $\pi^b$ . The state-space is defined as,  $\mathbf{s}_t^b = [\tilde{x}_{t,1}^b, \tilde{x}_{t,1}^b, \sigma_t^e]$ , where  $\tilde{x}_{t,1}^c, \tilde{x}_{t,1}^1$  represent the overall deterioration condition and speed for the bridge, and  $\sigma_t^e$  is the standard deviation for the condition of the elements at each time step  $t$ . The action-space has one action  $\delta^b$ , which corresponds to the target improvement, with  $\delta^b = 0$  being equivalent to *do nothing*, and  $0 > \delta^b \geq (u - l)$  is *maintain* the beam structural elements using  $\pi_k^*$ . Learning the policy  $\pi^b$  can be done using the vectorized RL environment at the bridge level, and the same DQN agent described in A. In this study, the continuous action space is discretized with  $\delta^b = \{\delta_1^b, \dots, \delta_A^b\}$  to make it compatible with discrete action algorithms [18]. Accordingly,  $\delta^b$  is represented by  $A = 10$  discrete action equally spaced over its continuous domain.

In order to assess the scalability and performance of the proposed HRL framework, the total number of the structural elements in  $\mathcal{C}_1 \in \mathcal{B}$  is varied with  $P = \{5, 10, 15\}$  beam elements. The performance of the HRL framework is evaluated based on 5 different environment seeds, and is compared with branching dueling DQN (BDQN) framework. The BDQN framework architecture, hyper-parameters and configuration are adapted from Tavakoli et al. [30]. Figure 7 shows the comparison of results based a structural category  $\mathcal{C}_1$  with  $P = 5$  beam elements in Figure 7a, a  $\mathcal{C}_1$  with  $P = 10$  beam elements in Figure 7b, and a  $\mathcal{C}_1$  with  $P = 15$  beam elements in Figure 7c. From Figure 7, the performance of the proposed HRL framework is reported while considering the pre-training phase required for learning the element level policy  $\pi_{k=1}$ , which extends over  $3 \times 10^6$  steps. Based on the results shown in Figure 7a for the case with  $P = 5$ , the HRL and BDQN frameworks achieve a similar total expected rewards, however, the BDQN approach shows a faster convergence due to the end-to-end training. Nonetheless, when the number of beam elements increases in the case of  $P = 10$  and  $P = 15$ , the HRL framework

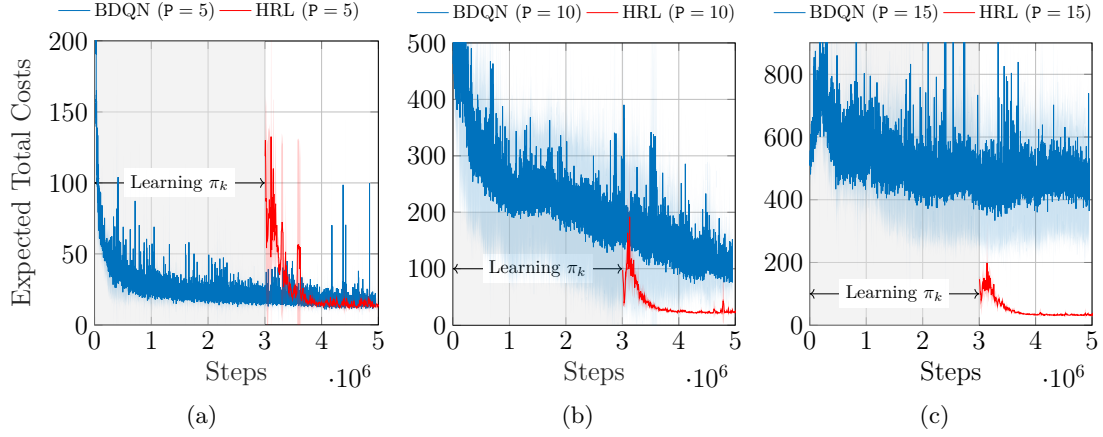


Figure 7: Comparison between the proposed HRL and BDQN for learning the maintenance policy of a structural category  $\mathcal{C}_1$  with  $P = 5$  elements in Figure 7a, a  $\mathcal{C}_1$  with  $P = 10$  elements in Figure 7b, and a  $\mathcal{C}_1$  with  $P = 15$  elements in Figure 7c. The training results are reported based on the average performance on 5 seeds, with the confidence interval represented by  $\pm\sigma$ .

outperforms the BDQN approach in convergence speed and in the total expected return achieved in this experiment setup. This can be attributed to the BDQN considering each additional beam element as a distinct branch which leads to a significant increase in the size of the neural network model, and thus requiring a higher number of samples for the training.

## 4.2 Maintenance Policy for a Bridge with Multiple Structural Categories

This example extends the application of the proposed formulation to a planning problem involving  $K = 6$  structural categories within a bridge  $\mathcal{B}$ . Each structural category consists of a different number of structural elements which are summarized in Figure 8. In this example, the bridge’s health state is

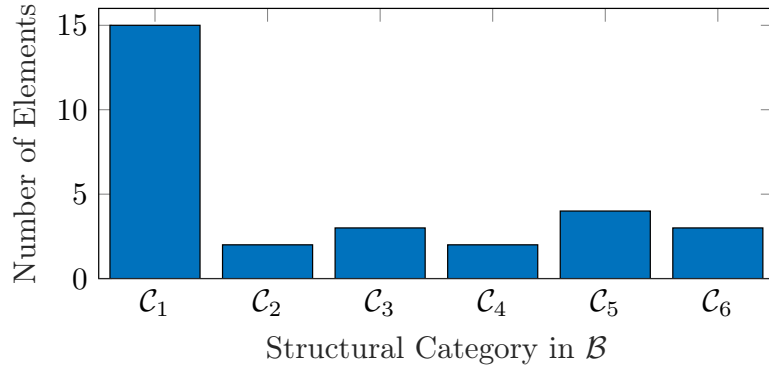


Figure 8: Number of structural elements within each structural category in bridge  $\mathcal{B}$ , where  $\mathcal{C}_1$  represents the beams,  $\mathcal{C}_2$  is the front-walls,  $\mathcal{C}_3$  is the slabs,  $\mathcal{C}_4$  is the guardrail,  $\mathcal{C}_5$  is the wing-wall, and  $\mathcal{C}_6$  represents the pavement.

represented by,  $\mathbf{s}^b = [\tilde{x}_t^b, \tilde{x}_t^b, \sigma^c]$ , where  $\sigma^c$  is the standard-deviation for the deterioration condition of the structural categories within  $\mathcal{B}$ . Similarly to the previous example, the bridge-level action is the target improvement  $\delta^b$  which is represented by 10 discrete action bins uniformly covering the continuous domain. Solving this maintenance planning problem is done by first identifying the optimal policy for the structural elements  $e_p^k$  within each structural category  $\mathcal{C}_k$ , where the element-level actions in each structural category have different costs and effects on the element’s state (see B). Figure 9 shows the optimal policy maps as learned by a DQN RL agent for the elements within each type of structural category. It should be noted that the characteristics of structural elements (e.g., critical thresholds) within each structural category are assumed to be identical, which allows learning a single policy per

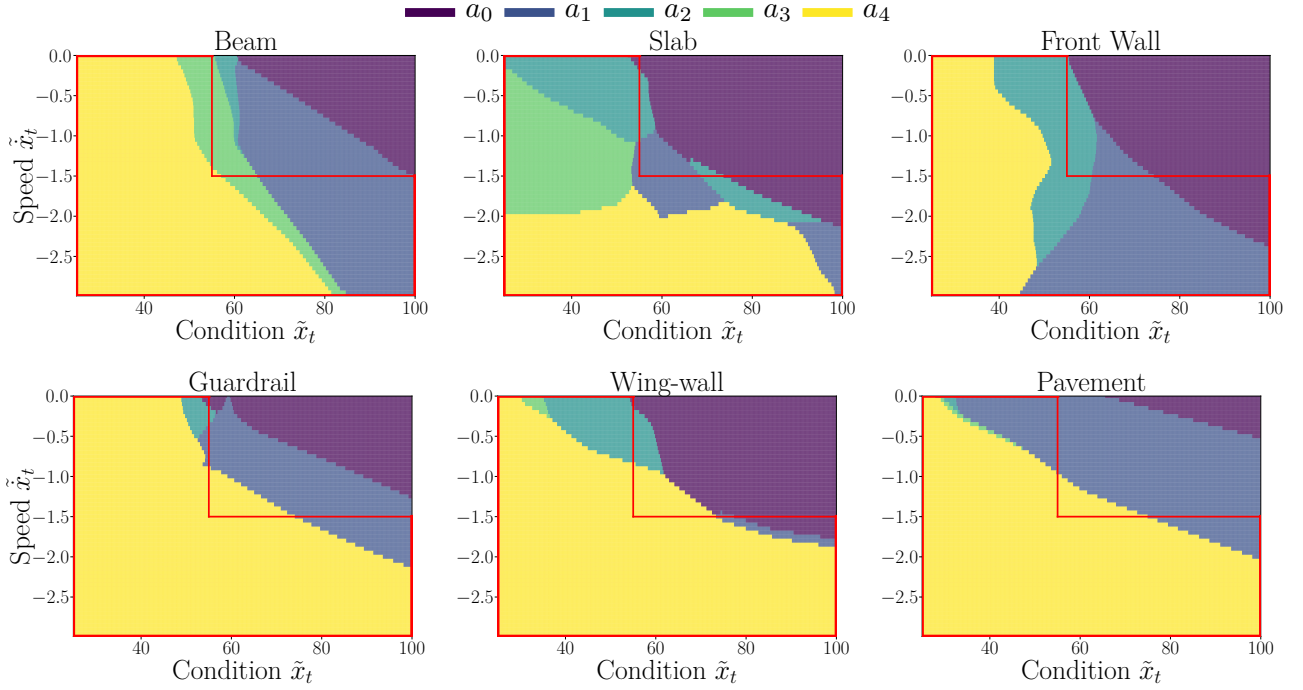


Figure 9: The element-level policy maps for the action space  $\mathcal{A}^{\ell_0}$ , where each policy map is learned independently by a DQN agent.

structural category. After obtaining the policies  $\pi_{1:K}^*$ , the hierarchical RL framework is trained using the environment at the bridge level. The training process for the DQN agent at the bridge-level is reported using the average performance on 5 different seeds for the environment, as shown in Figure 10, where it can be noticed that the policy’s training became stable after  $2 \times 10^6$  steps.

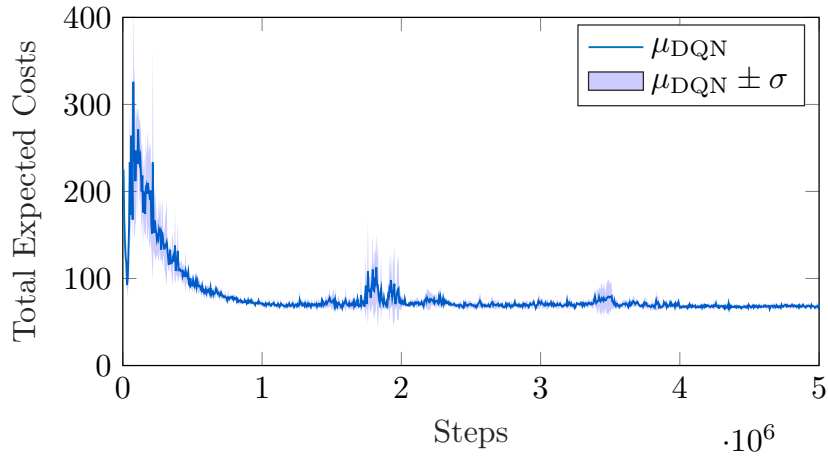


Figure 10: Average performance based on 5 seeds for DQN agents in learning the maintenance policy  $\pi^b$  for a bridge composed of multiple structural categories.

In order to use the hierarchical DRL agent for decision making on bridge  $\mathcal{B}$ , it is required to obtain the deterioration state estimates for each structural element, as well as the overall deterioration state of the bridge  $\mathcal{B}$ . This step can be achieved by relying on the element-level inspection data and using the SSM-based deterioration model for estimating and aggregating the deterioration states [15]. Accordingly, the policy  $\pi^b$  in the HRL framework relies on  $\mathbf{s}^b = [\tilde{\mu}_t^b, \tilde{\mu}_t^b, \sigma^c]$ , where  $\tilde{\mu}_t^b$  and  $\tilde{\mu}_t^b$  are the expected values for the deterioration condition and speed, respectively. On the other hand, each policy  $\pi_k$  depends on the expected values of the deterioration condition  $\tilde{\mu}_{t,p}^k$  and speed  $\tilde{\mu}_{t,p}^k$  at the element level, as in  $\mathbf{s}_t^b = [\tilde{\mu}_{t,p}^k, \tilde{\mu}_{t,p}^k]$ . Figure 11 illustrates the deterioration state estimates and the effect of maintenance on the deterioration condition and speed of the bridge  $\mathcal{B}$ . In this context, the

decision making analyses are performed using a window of 10 years, starting from the year 2021. The

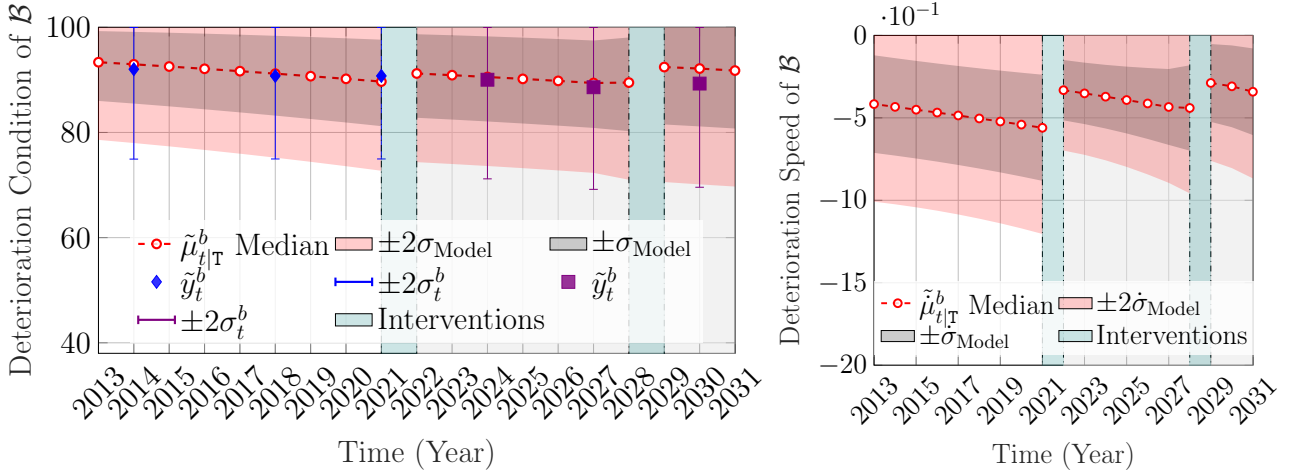


Figure 11: Deterioration state estimates for the condition and the speed of bridge  $\mathcal{B}$ , based on the aggregation of the deterioration state estimates of the structural categories  $\mathcal{C}_{1:k}$ , with the aggregated inspections  $\tilde{y}_t^b \in [25, 100]$  represented by the blue diamond, and their corresponding uncertainty estimates represented by the blue error bars. The inspections represented by a magenta square correspond to synthetic inspections on a trajectory of deterioration that is generated based on the RL agent interventions, which are suggested at years 2022 and 2029.

aggregated synthetic inspections illustrated in the magenta squares on Figure 11 are generated using the mechanism described in Section 3.2 for generating synthetic data. From Figure 11, and based on the bridge state  $s^b$ , the HRL agent suggests to perform maintenance actions at years 2022 and year 2029. The breakdown for the maintenance actions at the element level is shown in Figure 12, where the majority of the proposed maintenance actions are  $a_1$  : routine maintenance with the exception to a wing-wall element that is suggested to undergo a replacement in the year 2022. The replacement action is suggested by the policy  $\pi_k$  mainly due to a high deterioration speed as  $\tilde{\mu}_{t,p}^k < -1.5$ , which bypasses the critical state’s speed threshold.

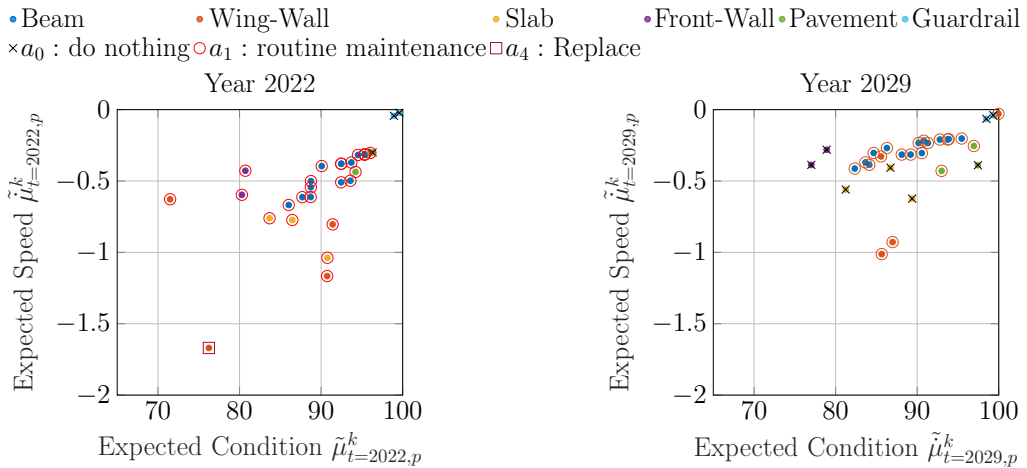


Figure 12: Scatter plot for the expected deterioration condition versus the expected deterioration speed for all elements of bridge  $\mathcal{B}$ , with maintenance actions suggested by the HRL agent at the years 2022 (on the left) and year 2029 (on the right).

### 4.3 Discussion

Typical bridges have hundreds of structural elements, many of which exhibit a similar structural deterioration behaviour. Capitalizing on the structural similarities can enable tackling maintenance

planning for bridges at scale, in a sense that learning the task of maintaining a beam structural element, can be either generalized for other similar beam structural elements, or can provide a source policy that would accelerate the learning of maintenance tasks (e.g., maintaining slabs). The latter falls under transfer learning in RL [37], and is not covered in the scope of this work, yet it shows potential for future work. The proposed HRL approach offers the capacity to take advantage of the aforementioned aspects, in addition to providing an interpretable decision maps that enable verifying the coherence of the optimal maintenance policy. This is important because in the context of maintenance planning there is no clear definition for the stopping criteria during the training, unlike some classic RL benchmarks, where the stopping criteria can be defined based on the success rate of the agent in accomplishing the task.

The advantages in the proposed HRL coincide with limitations that are mainly related to state abstraction and learning the policies. The level of abstraction in representing the state is dependent on the number of structural elements in the bridge, such that, for a bridge with many structural elements  $E > 100$ , the overall bridge condition and speed may not be sufficient to fully describe the health state, and accordingly the abstract state space should be augmented with additional information such as, the overall health condition and speed for each structural category. The other limitation in the proposed HRL is the use of a bottom-to-top approach for learning the policies with fixed policies  $\pi_k$  [10]. Alleviating these limitation could be done by using the policies  $\pi_k$  as a source policy that provide demonstrations for an end-to-end hierarchical RL training.

## 5 Conclusion

This paper introduce a hierarchical formulation and a RL environment for planning maintenance activities on bridges. The proposed formulation enables decomposing the bridge maintenance task into sub-tasks by using a hierarchy of policies, learned via deep reinforcement learning. In addition, the hierarchical formulation incorporates the deterioration speed in the decision-making analyses by relying on a SSM-based deterioration model for estimating the structural deterioration over time. A case study of a bridge is considered to demonstrate the applicability of the proposed approach which is done in two parts. The first part considered varying the number of structural elements to examine the scalability of the proposed framework against the branching dueling Q-network (BDQN). The comparison of results have shown that the proposed hierarchical approach has a better scalability than BDQN while sustaining a similar performance. The second part in the case study addressed a maintenance planning problem for a bridge with multiple structural categories. In this case, the HRL agent performance is demonstrated by the element-level maintenance actions performed over a span of 10 years.

Overall, this study has demonstrated the capacity to learn a maintenance policy using hierarchical RL for a bridge with multiple structural categories. In addition, the analyses have highlighted the role of the deterioration speed in the decision-making process. Further extensions to this framework may include a multi-agent setup to learn a network-level maintenance policies under budgetary constraints, as well as designing and testing RL frameworks that can handle the uncertainty associated with the deterioration state in a POMDP environment. The contributions in this paper also include an open-source RL benchmark environment (link: <https://github.com/CivML-PolyMtl/InfrastructuresPlanner>), which is made available for contributions by the research community. This RL environment provides a common ground for designing and developing maintenance planning policies, in addition to comparing different maintenance strategies.

## Acknowledgements

This project is funded by the Transportation Ministry of Quebec Province (MTQ), Canada. The authors would like to acknowledge the support of Simon Pedneault for facilitating the access to information related to this project.

## A Deep Reinforcement Learning

### A.1 Dueling Deep Network

In the context where a state  $s$  has similar  $q(s, a; \theta)$  values for different actions  $a \in \mathcal{A}$ , learning the value function  $v(s)$  for each state can facilitate in learning the optimal policy  $\pi^*$ . The dueling network architecture enables incorporating the value function  $v(s)$  in the Q-learning by considering,

$$Q(s, a; \theta_\alpha, \theta_\beta) = V(s; \theta_\alpha) + adv(s, a; \theta_\beta), \quad (14)$$

where  $adv(s, a)$  is the approximation for the advantage of taking action  $a$  in state  $s$ , and  $\theta_\alpha, \theta_\beta$  are the set of parameters associated with value function and the advantage function, respectively. Further details about the dueling network architecture are available in the work of Wang et al. [31].

### A.2 DRL Hyperparameters

The RL agents at all levels are trained using a discount factor 0.99, while relying on a batch size of 50 samples. The environment is vectorized to accelerate the training process and improve the sample independence, as the agent simultaneously interacts with  $n = 50$  randomly seeded environments. The exploration is performed using  $\epsilon$ -greedy, which is annealed linearly over the first 200 episodes with minimum  $\epsilon_{\min} = 0.01$ . Furthermore, the target model updates are performed every 100 steps in the environment. All neural networks have the same architecture (for the structural categories and the bridge) which consists in 2 layers of 128 hidden units and  $relu(\cdot)$  activation functions. The learning rate starts at  $10^{-3}$  and is reduced to  $10^{-5}$  after 800 episodes.

## B Environment Configuration

This section presents some of the predefined functions in the environment which are based on previous work and numerical experiments.

### B.1 SSM-based Deterioration Model

The Kalman filter (KF) describes the transition over time  $t$ , from the hidden state  $\mathbf{x}_{t-1}$  to the hidden state  $\mathbf{x}_t$  using the prediction step and the update step. The prediction step is described by,

$$\begin{aligned} \mathbb{E}[\mathbf{X}_t | \mathbf{y}_{1:t-1}] &\equiv \boldsymbol{\mu}_{t|t-1} = \mathbf{A}_t \boldsymbol{\mu}_{t-1|t-1} \\ \text{cov}[\mathbf{X}_t | \mathbf{y}_{1:t-1}] &\equiv \boldsymbol{\Sigma}_{t|t-1} = \mathbf{A}_t \boldsymbol{\Sigma}_{t-1|t-1} \mathbf{A}^\top + \mathbf{Q}_t, \end{aligned}$$

where  $\mathbb{E}[\mathbf{X}_t | \mathbf{y}_{1:t-1}]$  the expected value and  $\text{cov}[\mathbf{X}_t | \mathbf{y}_{1:t-1}]$  represent the covariance associated with the hidden state vector  $\mathbf{x}_t$  given all the observations  $\mathbf{y}_{1:t-1}$  up to time  $t-1$ ,  $\mathbf{A}_t$  is the transition matrix and  $\mathbf{Q}_t$  is the model process-error covariance. In this context, the transition matrix  $\mathbf{A}_t$  is time dependent such that,

$$\mathbf{A}_{t=\tau} = \begin{bmatrix} \mathbf{A}^{\text{ki}} & \mathbf{I}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{I}_{3 \times 3} \end{bmatrix}, \mathbf{A}_{t \neq \tau} = \begin{bmatrix} \mathbf{A}^{\text{ki}} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{I}_{3 \times 3} \end{bmatrix}, \mathbf{A}^{\text{ki}} = \begin{bmatrix} 1 & dt & \frac{dt^2}{2} \\ 0 & 1 & dt \\ 0 & 0 & 1 \end{bmatrix},$$

where  $\tau$  is the time of the element-level maintenance action, and  $\mathbf{I}$  is the identity matrix. Accordingly, the covariance matrix  $\mathbf{Q}_t$  is described by,

$$\mathbf{Q}_{t=\tau} = \begin{bmatrix} \mathbf{Q}^{\text{ki}} + \mathbf{Q}^r & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{Q}^r \end{bmatrix}, \mathbf{Q}_{t \neq \tau} = \begin{bmatrix} \mathbf{Q}^{\text{ki}} & \mathbf{0}_{3 \times 3} \\ \mathbf{0}_{3 \times 3} & \mathbf{0}_{3 \times 3} \end{bmatrix},$$

with  $\mathbf{Q}^r$  and  $\mathbf{Q}^{\text{ki}}$  defined as,

$$\mathbf{Q}^r = \text{diag}([\sigma_{w_r}^2, \dot{\sigma}_{w_r}^2, \ddot{\sigma}_{w_r}^2]), \mathbf{Q}^{\text{ki}} = \sigma_w^2 \begin{bmatrix} \frac{dt^5}{20} & \frac{dt^4}{8} & \frac{dt^3}{6} \\ \frac{dt^4}{8} & \frac{dt^3}{3} & \frac{dt^2}{2} \\ \frac{dt^3}{6} & \frac{dt^2}{2} & dt \end{bmatrix},$$

where  $dt$  is the time step size,  $\sigma_w$  is a model parameter that describes the process noise and  $\mathbf{Q}^r$  is a diagonal matrix containing model parameters associated with the element-level intervention errors [14]. Following the prediction step, if an observation is available at any time  $t$ , the expected value and covariance are updated with the observation using the update step,

$$\begin{aligned} f(\mathbf{x}_t|\mathbf{y}_{1:t}) &= \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_{t|t}, \boldsymbol{\Sigma}_{t|t}) \\ \boldsymbol{\mu}_{t|t} &= \boldsymbol{\mu}_{t|t-1} + \mathbf{K}_t(\mathbf{y}_t - \mathbf{C}\boldsymbol{\mu}_{t|t-1}) \\ \boldsymbol{\Sigma}_{t|t} &= (\mathbf{I} - \mathbf{K}_t\mathbf{C})\boldsymbol{\Sigma}_{t-1|t-1} \\ \mathbf{K}_t &= \boldsymbol{\Sigma}_{t-1|t-1}\mathbf{C}^\top\mathbf{G}_t^{-1} \\ \mathbf{G}_t &= \mathbf{C}\boldsymbol{\Sigma}_{t-1|t-1}\mathbf{C}^\top + \boldsymbol{\Sigma}_V, \end{aligned}$$

where  $\boldsymbol{\mu}_{t|t} \equiv \mathbb{E}[\mathbf{X}_t|\mathbf{y}_{1:t}]$  is the posterior expected value and  $\boldsymbol{\Sigma}_{t|t} \equiv \text{cov}[\mathbf{X}_t|\mathbf{y}_{1:t}]$  represents the covariance, conditional to observations up to time  $t$ ,  $\mathbf{K}_t$  is the Kalman gain, and  $\mathbf{G}_t$  is the innovation covariance. The monotonicity throughout the estimation process is imposed by relying on the deterioration speed constraints:  $\dot{\boldsymbol{\mu}}_{t|t} + 2\sigma_{\dot{\boldsymbol{\mu}}_{t|t}} \leq 0$ ; which are examined at each time step  $t$ , and enacted using the PDF truncation method [28].

Aggregating the deterioration states is performed using a Gaussian mixture reduction (GMR) [27], which is employed to approximate a PDF of  $E_m$  Gaussian densities into a single Gaussian PDF by using,

$$\begin{aligned} \boldsymbol{\mu}_{t|T,m}^{j,*} &= \sum_{p=1}^{E_m} \lambda_p^j \boldsymbol{\mu}_{t|T,p}^j, \\ \boldsymbol{\Sigma}_{t|T,m}^{j,*} &= \sum_{p=1}^{E_m} \lambda_p^j \boldsymbol{\Sigma}_{t|T,p}^j + \sum_{p=1}^{E_m} \lambda_p^j (\boldsymbol{\mu}_{t|T,p}^j - \boldsymbol{\mu}_{t|T,m}^{j,*})(\boldsymbol{\mu}_{t|T,p}^j - \boldsymbol{\mu}_{t|T,m}^{j,*})^\top, \end{aligned}$$

where  $\boldsymbol{\mu}_{t|T,m}^{j,*}$  is the aggregated expected value, and  $\lambda_p^j$  is the weight associated with the contribution of the deterioration state of the structural element. The merging of the  $E_m$  Gaussian densities is moment-preserving, where the total covariance  $\boldsymbol{\Sigma}_{t|T,m}^{j,*}$  consists in the summation of the within-elements contribution to the total variance, and the between-elements contribution to the total variance [27, 15].

## B.2 Decaying Factors for Effect of Interventions

Decaying factors are introduced to prevent having the same improvement effect on structural elements while applying the same action within a short a period of time. The decaying factors in this context rely on the estimate for the expected time (number of years) to return to the state prior to the intervention [14]. Accordingly, the effect of intervention for any element-level action  $a^e$ , at time  $t$  is,

$$\delta^e = \rho_1 \times \delta^e,$$

where  $\rho_1$  is the decaying factor defined as,  $\rho_1 \propto \Pr(X_{\tau+t} \leq x_{\tau-1}|a^e)$ , and  $\tau$  is the time of intervention.

## B.3 Maintenance Actions Effects & Costs

Maintenance actions at the element level have different effects depending on the structural category type. The deterministic maintenance effects associated with each action are defined in Table 1. It should be noted that the values defined in Table 1 have been derived from estimates that are based on data from the network of bridges in the province of Quebec [14].

As for the cost of maintenance actions, the cost functions are considered to be dependent on the deterioration state using,

$$x_c(\tilde{x}_{t,p}^k, a) = \beta_1(a) \frac{1}{\tilde{x}_{t,p}^k} + \beta_2(a),$$

where  $\beta_1(a)$  is the cost of performing the maintenance action  $a$  as a function of the deterioration state  $x_{t,p}^e$ , and  $\beta_2(a)$  is a fixed cost associated with maintenance action  $a$ . The derivation of this relation is empirical and mimics the cost information provided by the ministry of transportation in Quebec.



Table 1: Table of the true effects associated with element-level maintenance actions within each structural category.

	Structural Category					
	Beams	Front Wall	Slabs	Guardrail	Wing Wall	Pavement
$a_0$	0	0	0	0	0	0
$a_1$	0.5	0.1	1	0.25	0.25	8
$a_2$	7.5	19	12	9	8	20
$a_3$	18.75	20.5	20	14	17	28
$a_4$	75	75	75	75	75	75

Figure 13 shows the proportional cost function for the elements within each structural category. From the graphs in Figure 13, it is noticeable that the replacement cost is considered fixed and independent of the structural condition. Moreover, in some cases, the cost of performing an action may exceed the cost of replacement.

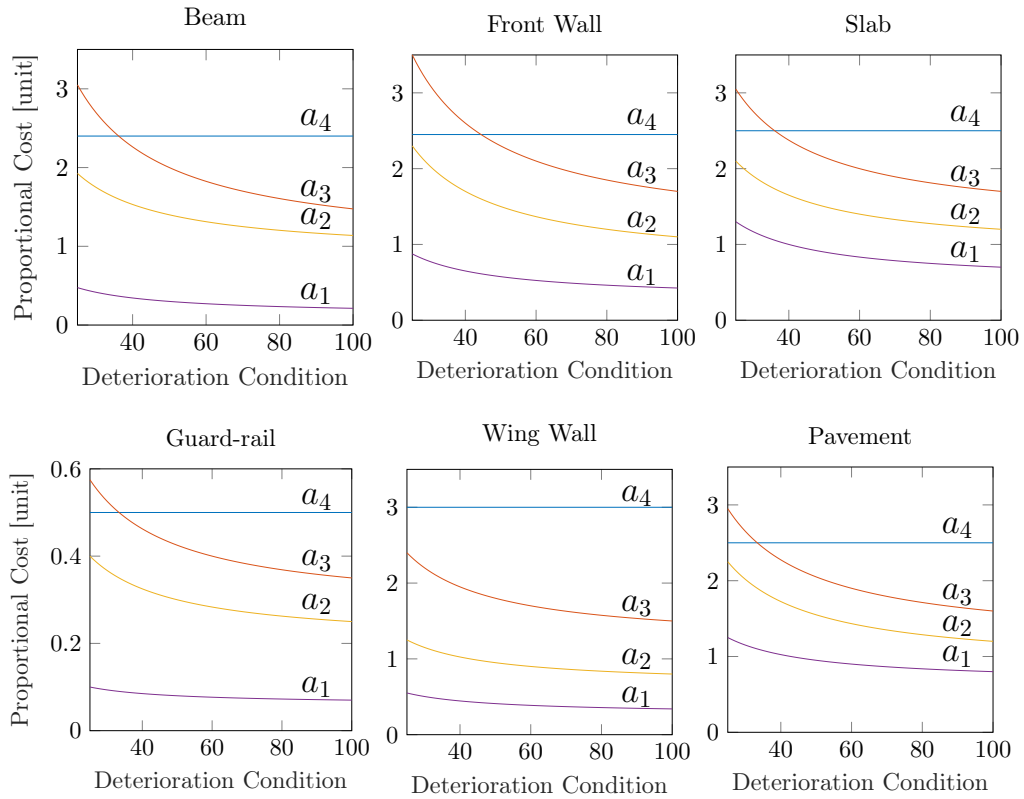


Figure 13: The proportional cost of each element-level action as a function of the deterioration condition.

Based on the cost function  $x_c(\cdot)$ , the element-level rewards  $r(\mathbf{s}_{t,p}^e, a_{t,p}^k)$  are defined as,

$$r(\mathbf{s}_{t,p}^e, a_{t,p}^k) = x_c(\tilde{x}_{t,p}^k, a_{t,p}^k) + r^p,$$

where  $r^p$  is the penalty applied when a predefined critical state is reached and no maintenance action is taken.

## References

- [1] Monireh Abdoos, Nasser Mozayani, and Ana LC Bazzan. Holonic multi-agent system for traffic signals control. *Engineering Applications of Artificial Intelligence*, 26(5-6):1575–1587, 2013.

- [2] David Abel. A theory of state abstraction for reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.
- [3] David Abel, David Hershkowitz, and Michael Littman. Near optimal behavior via approximate state abstraction. In *International Conference on Machine Learning*, pages 2915–2923. PMLR, 2016.
- [4] Duzgun Agdas, Jennifer A Rice, Justin R Martinez, and Ivan R Lasa. Comparison of visual inspection and structural-health monitoring as bridge condition assessment methods. *Journal of Performance of Constructed Facilities*, 30(3):04015049, 2015.
- [5] Charalampos P Andriotis and Konstantinos G Papakonstantinou. Managing engineering systems with large state and action spaces through deep reinforcement learning. *Reliability Engineering and System Safety*, 191:106483, 2019.
- [6] Vahid Asghari, Ava Jahan Biglari, and Shu-Chien Hsu. Multiagent reinforcement learning for project-level intervention planning under multiple uncertainties. *Journal of Management in Engineering*, 39, 2023.
- [7] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *Arxiv*, 2016.
- [8] Ao Du and Alireza Ghavidel. Parameterized deep reinforcement learning-enabled maintenance decision-support and life-cycle risk assessment for highway bridge portfolios. *Structural Safety*, 97, 2022.
- [9] Ehsan Fereshtehnejad and Abdollah Shafieezadeh. A randomized point-based value iteration pomdp enhanced with a counting process technique for optimal management of multi-state multi-element systems. *Structural Safety*, 65:113–125, 2017.
- [10] Carlos Florensa, Yan Duan, and Pieter Abbeel. Stochastic neural networks for hierarchical reinforcement learning. *Arxiv*, 2017.
- [11] Sven Gronauer and Klaus Diepold. Multi-agent deep reinforcement learning: a survey. *Artificial Intelligence Review*, 55:895–943, 2022.
- [12] Zachary Hamida and James-A Goulet. Modeling infrastructure degradation from visual inspections using network-scale state-space models. *Structural Control and Health Monitoring*, pages 1545–2255, 2020.
- [13] Zachary Hamida and James-A Goulet. Network-scale deterioration modelling based on visual inspections and structural attributes. *Structural Safety*, 88:102024, 2020.
- [14] Zachary Hamida and James-A. Goulet. Quantifying the effects of interventions based on visual inspections of bridges network. *Structure and Infrastructure Engineering*, pages 1–12, 2021. doi: 10.1080/15732479.2021.1919149.
- [15] Zachary Hamida and James-A. Goulet. A stochastic model for estimating the network-scale deterioration and effect of interventions on bridges. *Structural Control and Health Monitoring*, pages 1545–2255, 2021.
- [16] Junchen Jin and Xiaoliang Ma. Hierarchical multi-agent control of traffic lights based on collective learning. *Engineering applications of artificial intelligence*, 68:236–248, 2018.
- [17] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1):35–45, 1960.
- [18] Anssi Kanervisto, Christian Scheller, and Ville Hautamäki. Action space shaping in deep reinforcement learning. In *2020 IEEE Conference on Games (CoG)*, pages 479–486. IEEE, 2020.

- [19] Taisuke Kobayashi and Wendyam Eric Lionel Ilboudo. T-soft update of target network for deep reinforcement learning. *Neural Networks*, 136:63–71, 2021.
- [20] Jelle R Kok and Nikos Vlassis. Collaborative multiagent reinforcement learning by payoff propagation. *Journal of Machine Learning Research*, 7:1789–1828, 2006.
- [21] Xiaoming Lei, Ye Xia, Lu Deng, and Limin Sun. A deep reinforcement learning framework for life-cycle maintenance planning of regional deteriorating bridges using inspection data. *Structural and Multidisciplinary Optimization*, 65, 2022.
- [22] Mark Moore, Brent M Phares, Benjamin Graybeal, Dennis Rolander, and Glenn Washer. Reliability of visual inspection for highway bridges, volume i. Technical report, Turner-Fairbank Highway Research Center, 2001.
- [23] MTQ. *Manuel d’Inspection des Structures*. Ministère des Transports, de la Mobilité Durable et de l’Électrification des Transports, 2014.
- [24] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- [25] Van Thai Nguyen, Phuc Do, Alexandre Vosin, and Benoit Iung. Artificial-intelligence-based maintenance decision-making and optimization for multi-state component systems. *Reliability Engineering and System Safety*, 228, 2022.
- [26] Shubham Pateria, Budhitama Subagdja, Ah Hwee Tan, and Chai Quek. Hierarchical reinforcement learning: A comprehensive survey. *ACM Computing Surveys*, 54, 2021.
- [27] Andrew R Runnalls. Kullback-leibler approach to gaussian mixture reduction. *IEEE Transactions on Aerospace and Electronic Systems*, 43(3):989–999, 2007.
- [28] Dan Simon and Donald L Simon. Constrained kalman filtering via density function truncation for turbofan engine health estimation. *International Journal of Systems Science*, 41(2):159–171, 2010.
- [29] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [30] Arash Tavakoli, Fabio Pardo, and Petar Kormushev. Action branching architectures for deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [31] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. *Dueling network architectures for deep reinforcement learning*. PMLR, 2016.
- [32] Christopher John Cornish Hellaby Watkins. *Learning from delayed rewards*. PhD thesis, University of Cambridge, 1989.
- [33] Shiyin Wei, Yuequan Bao, and Hui Li. Optimal policy for structure maintenance: A deep reinforcement learning framework. *Structural Safety*, 83, 2020.
- [34] David Y Yang and A M Asce. Deep reinforcement learning-enabled bridge management considering asset and network risks. *Journal of Infrastructure Systems*, 28(3):04022023, 2022.
- [35] Nailong Zhang and Wujun Si. Deep reinforcement learning for condition-based maintenance planning of multi-component systems under dependent competing risks. *Reliability Engineering and System Safety*, 203, 2020.
- [36] Yifan Zhou, Bangcheng Li, and Tian Ran Lin. Maintenance optimisation of multicomponent systems using hierarchical coordinated reinforcement learning. *Reliability Engineering and System Safety*, 217, 2022.
- [37] Zhuangdi Zhu, Kaixiang Lin, and Jiayu Zhou. Transfer learning in deep reinforcement learning: A survey. *arXiv*, 2020.