

Intelligent and adaptive asset management model for railway sections using the *i*PN method

Ali Saleh^{a,b,*}, Rasa Remenyte-Prescott^c, Darren Prescott^c, Manuel Chiachío^{a,b}

^aAndalusian Research Institute in Data Science and Computational Intelligence, University of Granada, 18071, Spain.

^bDept. Structural Mechanics & Hydraulics Engineering, University of Granada, 18071, Granada, Spain.

^cResilience Engineering Research Group, The University of Nottingham Science Road, Nottingham, University Park, NG7 2RD, United Kingdom.

Abstract

The maintenance strategy in railway transportation is crucial in ensuring safety, availability, and reducing operating costs. However, finding the optimal maintenance plan that takes into account the complex relationships between railway assets can be a challenging task. To address this challenge, this study introduces an Intelligent Petri Net (*i*PN) model to effectively consider the maintenance and operation of railway sections with a focus on optimising ballast maintenance. The *i*PN model merges Petri net (PN) with Reinforcement Learning (RL) to create a model that is able to simulate and learn at the same time. The model is able to use diverse information, including usage, degradation rates, maintenance effectiveness, fault probabilities, and maintenance time, to simulate and learn at the same time. By considering the interconnections between these factors, the model found that reducing unnecessary maintenance actions increases the age of railway sections and leads to higher net profits. The study also introduced a method to reduce computational effort by dividing the PN into subnets and another method to make learning faster by using multiple RL environments. In conclusion, the developed *i*PN model presents a promising solution for optimising ballast maintenance within railway operation.

Keywords: Petri net, Reinforcement learning, Q-learning, railway, maintenance modelling, degradation models

1. Introduction

1 Railway systems serve as the backbone of modern transportation, facilitating the movement of goods
2 and people across vast distances with efficiency and reliability. However, the seamless functioning of these
3 intricate networks relies on a complex interplay of various components, each of which demands special
4 attention. One such important element is the ballast, the layer of crushed stones beneath the tracks that
5 provides stability, distributes load, and facilitates water drainage. The optimization of ballast maintenance
6

*Corresponding authors. e-mail: alisaleh@ugr.es

7 is an endeavor of paramount significance, as it directly influences the safety, performance, and sustainability
8 of the entire railway system.

9 The settlement of ballast and other underlying layers, such as the formation, causes the deterioration of
10 railway track geometry. As the track geometry degrades, the track becomes uneven and both ride quality
11 and safety are affected. Track maintenance can address degraded track geometry but poor maintenance
12 planning can lead to levels of degradation that are sufficiently high for the railway to become unfit for
13 purpose. At best, this can lead to high monetary losses due to downtime and corrections, and at worst, to
14 safety risks including fatalities and injuries to the public, users, and workforce. An example of catastrophic
15 failure is the Potters Bar train derailment, which resulted in 7 fatalities, 76 injuries, and a £3,150,000 fine
16 for Network Rail [1]. This is why renewal and maintenance planning in the railway is a critically-important
17 decision-making problem for engineers and a crucial topic for ongoing research.

18 Railway track settlement is directly affected by ballast quality. Good ballast should be free of dust,
19 dirt, or small particles [2]. Ballast is fouled when it contains small particles, and this can lead to several
20 undesirable consequences. Fouling impedes the water drainage path, which results in a saturated subgrade
21 and permanent deformation of the soil, and can lead to wet beds that increase rail and sleeper degradation
22 rates. Ballast fouling also impacts the distribution of loads, affecting the settlement within the subgrade,
23 and degrading the vertical track geometry profile. Causes of ballast fouling include the breakdown of ballast
24 particles due to dynamic forces from traffic and maintenance activities. Once the ballast becomes highly
25 fouled, maintenance actions are no longer effective and renewal should take place. Before renewal, it is
26 important to avoid unnecessary maintenance actions because of the ballast breakdown that they induce.
27 This can be achieved by following a condition-based maintenance strategy, which calls for maintenance only
28 when needed. However, the challenge is to know when and what type of maintenance is essential for the
29 different known system conditions.

30 The connection between the operational life of the ballast and maintenance actions is apparent due
31 to how they are interrelated. Maintenance actions, such as tamping or stoneblowing, are carried out to
32 correct the rail's vertical geometry profile, but they also have an impact on the ballast's condition. These
33 maintenance actions might lead to an increase in the rate of degradation and a reduction in the ballast's
34 lifespan. This reduction happens because the ballast can only go through a certain number of maintenance
35 actions before it becomes highly fouled and needs replacement. Consequently, a thoughtful assessment is
36 required to determine whether restoring the vertical geometry in each situation is preferable or if postponing
37 such actions would be a better choice to extend the ballast's lifespan. The objective of this study is to identify
38 the most suitable course of action, whether it involves tamping, stoneblowing, renewal, or no intervention,
39 based on the specific conditions and the track's maintenance history while reducing maintenance costs.

40 Several studies have considered the problem of optimizing the ballast maintenance [3–5]. The basis
41 for estimating the appropriate time for condition-based maintenance (CBM) interventions in railway is the

42 track degradation modeling [6]. Degradation models can be classified into physics-based, data-driven, and
43 hybrid models [7]. Physics-based models, which estimate degradation based on mechanical properties of track
44 components, offer adaptability to different traffic conditions and materials, particularly in the early stages of
45 the life cycle when historical data is limited [7]. However, these models are deterministic and often overlook
46 input uncertainties [8]. To address this limitation, hybrid approaches have been proposed, such as the
47 elastoplastic physics-based model combined with a sequential model that considers parameter uncertainty [7].
48 On the other hand, empirical models, particularly stochastic models, have been used to address uncertainties
49 in railway asset degradation behavior and model parameter estimation. These models utilize stochastic
50 processes, such as the Weibull distribution, Wiener process, Gamma process, and Petri net (PN), to represent
51 degradation rates of track components and predict track deterioration [9–13]. Furthermore, Markov chain
52 approaches have been employed to model the outcomes of transformations between different states of track
53 degradation, incorporating maintenance operations and discretized levels of degradation progression [14–16].
54 Recent advancements in artificial intelligence and machine learning have also gained attention in railway
55 transportation, with applications in predictive maintenance, condition monitoring, and track fault prediction
56 [17, 18].

57 Among the mentioned approaches, PNs are typically regarded as powerful modelling tools for degradation
58 and maintenance modelling due to their ability to account for resource availability, concurrency, and syn-
59 chronisation, which are common aspects that underline the majority of the asset management models, along
60 with their adequacy for dealing with highly multidimensional and heterogeneous input variables [19, 20].
61 However, ordinary PNs do not have learning capabilities, and that limits the capacity to autonomously
62 adapt the resulting maintenance schedule to the changing nature of the influencing conditions. Several
63 efforts were made to enrich the PN model with learning capabilities to allow it to be used for optimisation
64 problems. Plausible Petri nets [21], which are based on Bayesian learning, are effective in making the PN
65 self-adaptive, but they are limited to homogeneous and low-dimensional variables. Possibilistic Petri nets
66 [22] and fuzzy Petri nets (FPN) [23, 24] can be viewed as knowledge representation formalisms, but their
67 intelligence is limited to adjusting fuzzy production rules using soft-computing techniques. More recently, a
68 novel methodology referred to as *i*PNs has been proposed and utilizes Reinforcement learning (RL) to enable
69 and optimize decision-making and to upgrade the PN to an intelligent system [25]. The *i*PN method best
70 suits decision-making problems since RL, especially temporal difference learning, can be viewed as a more
71 general extension of dynamic programming (DP) that does not require a complete model of the environment
72 [26].

73 In this study, the *i*PN method is used to optimize condition-based maintenance and renewal of railway
74 ballast. Based on data from a typical European track, the method has been formulated and integrated
75 into the *i*PN model. The formulas cover factors affecting ballast condition, the impact of maintenance
76 on ballast condition and its degradation rate, and the effect of ballast condition on rail vertical geometry

77 profile and the rate of rail faults. The case of a railway system comprised of multiple sections has been
78 considered, with the decisions of each section evaluated independently based on its condition and other
79 pertinent factors affecting its state. This study incorporates in its objective function the costs associated
80 with life-cycle, maintenance and renewal, travel, delay, and traffic disruption, either directly or indirectly.
81 Distinctive features are presented by this study compared to prior work that is presented in the literature
82 review (Section 2). It is the first endeavor to employ the PN for detailing the operation and maintenance
83 intricacies of the ballast and other railway components, simultaneously utilizing RL to optimize maintenance
84 protocols. In earlier research, the emphasis was singular, either on the modeling facets or the optimization
85 components; however, a sophisticated model combined with optimization was lacking. The result was an
86 optimal maintenance strategy that balances safe and good condition of sections while reducing maintenance
87 frequency to decrease costs and increase the ballast’s section lifespan. The average lifespan increased from
88 29.5 to 42.5 years while reducing the probability of being in a *Super-red* condition from 0.04% to 0.007%.
89 Importantly, the developed model is not restricted to the study of the railway track sections for which it
90 was developed. It can be used to study the asset management of other railway track by adjusting input
91 parameters and rewards, without any need for the alteration of the core features of the PN model or the
92 associated analysis.

93 The *i*PN method, as outlined in [25], proposes a technique for integrating multiple decisions when
94 simultaneous decision-making is required. This approach entails incorporating multiple RL agents with
95 centralised decision-making, leading to convergence towards an optimal policy. However, the combination of
96 multiple decisions results in an exponential increase in the action space, making the learning process more
97 complex [27]. In order to mitigate this challenge, the method described in this study avoids the combination
98 of decisions and instead treats each section as a separate RL environment with its own unique elements.
99 This approach reduces the dimensions of both states and actions, as it eliminates the need to combine states
100 or actions. Furthermore, it enables agents to learn from one another when they are pursuing similar goals
101 under comparable conditions. This approach has been applied to a railway with 10 sections, resulting in
102 2830 states and 4 actions per state, as opposed to $3.29 \cdot 10^{24}$ states and 10^5 combinatorial actions per state.
103 Additionally, it enables the consideration of similar states from different sections as equivalent states with
104 equivalent actions, resulting in a reduction of the actually trained states to 283.

105 The current study also addresses the methodological challenge of reducing the complexity of PN models
106 for industrial applications in transportation, which often result in high computational costs. While previous
107 literature has proposed various techniques for reducing PN complexity, each of these approaches has its
108 own limitations. Reduction based on defined rules is effective in avoiding logical errors but may not be
109 sufficiently general for all types of PN structures [28–31]. Symmetrical reduction of PN can only be applied
110 when symmetries are present [30]. Reduction through the use of algebraic equations is only applicable when
111 the PN has specific properties such as redundant transitions, redundant places, or place agglomerations

112 [32]. Proposing reduced models and inferring their parameters so that the results at key outputs are similar
113 to those of the original PN requires additional computation for the inference process [20]. To overcome
114 these limitations, this study provides a systematic method for reducing PN complexity by decomposing it
115 into subnets with reduced computational cost while preserving its structure and functionality. This method
116 can be applied to any type of PN and requires minimal additional rules, leading to a 3 times reduction in
117 computational time for the considered PN case.

118 The rest of the paper is structured as follows. Section 2 reviews the maintenance policies, the undertaken
119 problems, the objective functions, and the methods used in the railway maintenance field. Section 3 presents
120 the underlying foundations of Q-learning and Petri nets, along with an overview of the *i*PN model and
121 the proposed method for decomposing a PN into multiple subnets. Section 4 introduces a technique for
122 decomposing the RL environment into multiple environments and a method for sharing experience between
123 RL agents within the context of the *i*PN. Section 5 details the creation of an operation and maintenance
124 intelligent PN model for a railway with multiple sections. The results of the railway case study are presented
125 in Section 6, followed by a discussion of the results in Section 7. Finally, Section 8 provides concluding
126 remarks.

127 **2. Literature Review on railway maintenance**

128 *2.1. Maintenance policies*

129 A maintenance policy is a decision made by managers based on maintenance models to ensure the proper
130 functioning of a system [15]. Maintenance policies can be categorized into three types: preventive, corrective,
131 and improvement [6]. The objective of the corrective policy is to enhance the asset's inherent reliability,
132 maintainability, or safety while preserving its original function. It involves repairing or replacing failed parts
133 to quickly restore equipment [33]. However, corrective maintenance is costly and increases safety risks due
134 to unexpected failures. Despite the implementation of preventive maintenance, unexpected failures can still
135 occur, leading to the need for corrective maintenance. Preventive maintenance reduces failures through
136 inspections and repairs [34]. It includes predetermined maintenance and CBM [6]. Predetermined maintenance
137 involves regular inspections and repairs, while CBM utilizes real-time data for proactive maintenance.
138 Predetermined maintenance can reduce the probability of disruption and system failure but can result in
139 additional unnecessary maintenance actions. On the other hand, CBM can avoid unnecessary maintenance
140 actions while ensuring safety and economic benefits. Predictive maintenance, a form of CBM, utilizes data
141 analysis and predictive modeling to identify potential issues before they occur [6]. However, it requires
142 significant resources for data collection and analysis.

143 *2.2. Undertaken problems*

144 The planning process in maintenance management aims to address crucial decisions regarding maintenance intervals for track segments and the allocation of necessary resources [17]. It involves ensuring the availability of required resources, determining appropriate actions, sequencing tasks, and identifying the necessary skills for maintenance operations [17]. The role of a planner, as explained by [35], includes assessing the scope of maintenance tasks, identifying the required expertise and craft, estimating the duration of tasks, and specifying the necessary parts and tools. Moreover, the planning function encompasses various aspects, such as task identification, complexity assessment, workforce estimation, spare parts and materials identification, and tool requirement determination [36]. The objective of the planning process is to make important choices regarding the timing of maintenance intervals for track segments and the allocation of necessary maintenance resources. In the railway industry, decision-making involves planning and scheduling activities, such as budgeting, quality prediction, project definition, project prioritization, possession allocation, timetabling, maintenance scheduling, and performance evaluation and feedback.

156 Maintenance planning involves several key aspects. One aspect is determining the timing of maintenance interventions based on accurate track condition prediction, which requires considering track geometry and track structure indices [37, 38]. However, relying solely on track geometry variables may not provide an accurate prediction of track condition [3]. To enhance maintenance planning, additional factors such as ballast fouling and geometry degradation should be considered when identifying maintenance needs [3]. Decision support systems and optimization models have been proposed to assist in the planning process. For instance, [39] developed a stochastic degradation model for condition-based maintenance (CBM) planning, while [40] optimized the number of tamping interventions considering track degradation and recovery. Furthermore, the setup cost of tamping equipment can be incorporated into the cost function [41].

165 Maintenance action identification and prioritization are crucial steps in maintenance planning. Railway infrastructure maintenance can be based on predetermined schedules or condition-based approaches [6]. Various optimization models have been proposed to determine the optimal maintenance limit intervals for different track quality indicators, taking into account preventive and corrective maintenance costs as well as potential train delays [42]. Furthermore, optimization models have been developed to decide whether immediate or postponed maintenance should be conducted based on factors such as reliability functions, associated costs, and identified defects [43]. Decision-making frameworks incorporating multi-attribute utility theory have also been used to prioritize maintenance tasks based on estimated conditions and multiple factors [44].

174 The scheduling of inspection intervals plays a crucial role in ensuring track safety and reliability while managing maintenance and inspection expenses. Optimization models have been proposed to determine inspection intervals based on safety risks and maintenance costs [45]. Rescheduling of inspection intervals has also been explored to mitigate the disruption caused by inspection scheduling and improve decision-

178 making [46]. The close relationship between inspection scheduling and maintenance scheduling emphasizes
179 the significance of inspection intervals in railway track maintenance planning and execution.

180 Possession scheduling is another important aspect of maintenance planning. Possession refers to the
181 closure of specific sections of railway tracks for maintenance or repair work. Effective possession scheduling
182 is essential for safe and efficient railway operations [47]. The optimization of possession scheduling can be
183 approached from different perspectives. Some studies focus on fixed train timetables and aim to determine
184 the best possession time for maintenance activities [48]. Others consider fixed possession times and seek
185 to optimize the train timetable around those periods [49]. Additionally, simultaneous possession and train
186 timetable scheduling models have been developed to optimize both train operations and maintenance activ-
187 ities [49]. Integrating maintenance activities and optimizing vehicle routing and crew scheduling can lead to
188 cost savings and improved efficiency [50]. Furthermore, equipment logistics, such as transporting machinery
189 and equipment to maintenance locations, need to be carefully planned and scheduled [51].

190 *2.3. Objective functions*

191 Objective functions of planning and scheduling play a crucial role in optimizing the efficiency and sus-
192 tainability of railway transportation systems. One significant aspect considered in these objective functions
193 is the life-cycle cost, which encompasses various expenses associated with the entire lifespan of railway
194 transportation [6]. These expenses include maintenance and replacement costs, building expenses, track
195 utilization fees, and costs related to the final stages of the system's operational life.

196 To enhance decision-making regarding new construction and the maintenance and replacement of track
197 components, decision-makers utilize life-cycle cost analysis [52]. This analysis takes into account both
198 measurable expenses such as construction, maintenance, and renewal, as well as intangible factors like
199 quality deterioration, traffic delays, safety concerns, and environmental impacts [52]. By considering these
200 various aspects, decision-makers can optimize investment strategies and ensure the long-term sustainability
201 and efficiency of railway transportation systems.

202 Researchers have identified four components that encompass the overall expenses associated with a
203 track and its rolling stock over its lifetime [53]. These components include construction costs, operational
204 aspects (such as capacity loss, fuel or energy consumption, environmental impact, accident risk, and socio-
205 economic implications), maintenance expenses, and costs incurred at the end of the track and rolling stock's
206 life. Studies have established the life-cycle cost of railway tracks by considering both measurable and non-
207 measurable expenses, such as maintenance, renewal activities, penalties due to track quality issues, customer
208 losses, and damage caused by subpar quality [9, 53].

209 Another important factor in planning and scheduling costs is the maintenance cost. A commonly used
210 approach is to assign a fixed cost per activity or time unit, which forms the basis for estimating the costs
211 [41]. For example, Gustavsson [41] proposed an improved linear programming model for scheduling tamping

212 operations on ballasted tracks, incorporating unit maintenance cost and the cost of maintenance occasions.
213 Daddow et al. [54] utilized a similar cost formulation to calculate the cost of each unit tamping action, while
214 Vale et al. [40] focused on reducing the number of tamping actions. Moreover, Letot et al. [10] considered a
215 fixed cost for the tamping machine.

216 Renewal costs are another important aspect of track maintenance. These costs can be classified into
217 two categories: component renewal and full track renewal [11]. Researchers have proposed optimization
218 frameworks to determine the optimal balance between track unavailability and life-cycle cost (LCC) [11].
219 These frameworks consider factors such as the unitary cost of renewal work, residual value of track compo-
220 nents, and potential savings from grouping track segments. Integrated methodologies have been developed
221 to account for equipment preparation, setup expenses, and predetermined expenditures associated with each
222 renewal activity throughout the lifespan of a component [55].

223 Possession cost is a significant factor to consider in maintenance operations. Previous studies have
224 proposed different approaches to address possession costs and their impact on overall maintenance costs.
225 One approach involves assigning hourly costs to account for the time required for possession in order to carry
226 out maintenance activities [56]. Another method utilizes fixed estimated possession costs per maintenance
227 action [57]. Train cancellations can also be taken into consideration when estimating possession costs [58].
228 The objective is to minimize the overall maintenance cost while taking possession costs into account.

229 In summary, objective functions for planning and scheduling of railway transportation systems encompass
230 a wide range of costs, including life-cycle costs, maintenance costs, renewal costs, and possession costs. By
231 considering these costs and optimizing decision-making processes, the efficiency, sustainability, and overall
232 performance of railway transportation systems can be enhanced.

233 *2.4. Search algorithms and simulation methods*

234 A suitable approach for solving the railway track maintenance planning and scheduling (RTMP&S)
235 problem involves considering decision-making levels, decision variables, track condition data, objectives,
236 and constraints. Linear and integer programming methods are commonly used in RTMP&S because they
237 can handle both continuous and integer decision variables [6]. Depending on the characteristics of the
238 decision variables, linear or nonlinear programming can be applied. For instance, integer programming is
239 suitable for determining maintenance actions or resource allocation.

240 In the case of single objective function models, mixed-integer linear programming methodologies that
241 combine continuous and integer variables are commonly employed [6]. Commercial solvers like CPLEX,
242 Gurobi, or FICO Xpress are often used to solve these optimization problems. In addition to these method-
243 ologies, various heuristics and metaheuristics have been employed to provide faster satisfactory solutions.
244 These include decomposition-based heuristics, multiple neighborhood search heuristics, solution frameworks
245 based on Lagrangian relaxation, iterative approaches with greedy and local search algorithms, tabu search

246 heuristics, and customized metaheuristic algorithms [50, 59–62].

247 For multi-objective function models, researchers have proposed various methodologies to optimize con-
248 flicting or non-conflicting objectives. These methodologies consider maintenance-related unavailability, life
249 cycle cost of track components, maintenance expenses, costs due to train delays, train maintenance planning,
250 timetabling, and selection of maintenance strategies. Multi-objective optimization techniques like weighted
251 sums and Pareto optimality are used to identify optimal solutions [16].

252 While linear programming dominates in RTMP&S, there is a growing interest in utilizing non-linear
253 programming. Non-linear formulations and search techniques such as the steepest gradient and improved
254 genetic algorithms have been employed for maintenance scheduling [63]. Other research directions include
255 the utilization of Model Predictive Control techniques at various levels, which involve methods like pattern
256 search, transformation into Mixed-Integer Linear Programming, Dantzig-Wolfe decomposition, and gradient-
257 free algorithms [12].

258 Researchers have shown a growing interest in integrating simulation models with optimization engines for
259 RTMP&S problems in recent years. Discrete event simulation is the predominant method used and offers
260 advanced capabilities to address the complexities associated with real-world maintenance planning and
261 scheduling problems [34]. Additionally, alternative approaches such as Monte Carlo simulation have been
262 employed to represent deterioration and restoration of track geometry, providing insights into maintenance
263 costs and optimal timing for interventions [42].

264 One widely used approach for optimizing maintenance in the railway industry is the Markov decision
265 process (MDP) [64]. MDP methods capture the stochastic nature of the railway system, incorporating uncer-
266 tainties and variability into maintenance decision-making. Probabilistic transitions between states in MDP
267 models allow decision-makers to account for degradation, failures, and repairs, resulting in more accurate
268 maintenance optimization. MDP-based dynamic programming excels in handling large-scale maintenance
269 optimization problems in the railway industry. Efficient algorithms like value iteration and policy iteration
270 compute optimal policies and value functions for complex systems. This scalability is crucial for considering
271 numerous components and subsystems within railway infrastructure. MDP-based approaches also facili-
272 tate the development of robust and adaptive maintenance strategies. By updating the value function and
273 policy based on changing system conditions, decision-makers can dynamically adapt maintenance strate-
274 gies to factors such as traffic patterns, weather conditions, and component aging. Reinforcement learning
275 (RL) is a promising approach for addressing complex railway industry problems. RL has been applied to
276 rail maintenance and renewal planning, optimizing costs and risk reduction [65]. It has also been used for
277 railway alignment optimization, minimizing construction costs while satisfying alignment constraints [66].
278 RL-based methods have been utilized for dynamic maintenance policies in multi-component systems with
279 degradation and random shocks [67]. Additionally, RL combined with digital twin technology has enhanced
280 railway maintenance efficiency, reducing maintenance activities and defects [68].

281 RL outperforms linear programming, non-linear programming, and MDP in railway maintenance and
 282 planning. It effectively handles uncertainties and variability, optimizing maintenance through probabilistic
 283 transitions and degradation probabilities in PN simulations. RL scales well, deriving optimal policies for
 284 complex systems. Its adaptability enables dynamic updates to maintenance strategies. Integration with
 285 advanced techniques like deep deterministic policy gradients and digital twin technology automates op-
 286 timization, reduces activities, and minimizes defects. RL proves valuable for railway track maintenance
 287 planning and scheduling.

288 3. Methodology

289 This section provides the methodological background and techniques proposed in this paper.

290 3.1. Basics about Reinforcement Learning

291 RL is a goal-oriented machine learning field that teaches an *agent* the correct decisions by trial and
 292 error. Single-agent RL methods can be formulated through a Markov decision process (MDP), which is
 293 described by a tuple of $\langle \mathcal{S}, \mathcal{A}, P_d \rangle$; where \mathcal{S} is the set of the *states* of the *environment*, $\mathcal{A}(s)$ is the set of
 294 *actions* available at *state* s , and P_d represents the *dynamics* of the MDP [26]. The *dynamics* is defined as
 295 $P_d = \Pr\{S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a\}$, which is the probability of obtaining *reward* r and *state* s'
 296 by taking an *action* a at *state* s . At each time step, t , the *agent* receives a *state* of the environment, S_t ,
 297 and takes an *action*, A_t , following a *policy* $\pi(a|s)$, which controls the probability of taking an action a being
 298 at state s . This results at the next time step, $t + 1$, in an immediate *reward*, R_{t+1} , and a change in the
 299 *state*, S_{t+1} . The goal of the agent is to find the optimal policy that maximises the long-run rewards, not the
 300 immediate reward. Long run rewards coming after a time step t are called the *expected return*, (G_t) , and
 301 can be calculated as:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad (1)$$

302 where $\gamma \in [0, 1]$ is a discount rate parameter to prevent $G_t \rightarrow \infty$ when $T \rightarrow \infty$ (known as a *continuous task*
 303 *problems*). On the contrary, in *episodic task problem*, the terminating time step T is a finite number, thus
 304 G_t can be calculated by choosing $\gamma = 1$.

305 Having a complete model of the environment dynamics, P_d , is not always feasible. Thus, model-free
 306 RL by temporal-difference learning (TDL) is widely used due to its simplicity and the minimal amount of
 307 computation [26]. In TDL, the value of each state-action pair is called the Q-Value and the whole set of
 308 Q-Values represents the outcomes of the Q-function, $q_\pi(s, a)$. Q-Values are updated in TDL method as
 309 follows:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[G_t - Q(S_t, A_t)] \quad (2)$$

310 In this equation the Q-Value is updated toward a target value, which is G_t , and $\alpha \in [0, 1]$, which is the
 311 learning rate, represents how much change will be made toward this target. G_t can be calculated in several
 312 ways. Using a one-step bootstrapping technique to calculate G_t as $G_t = R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$ results
 313 in the *Q-learning* method, which is one of the earliest and most famous TDL methods [69]. Another way
 314 is to calculate it based on Equation 1, and this results in one of the incremental implementations of the
 315 Monte-Carlo RL (MCRL) method [26]. MCRL is good to use at the beginning of the learning process since it
 316 does not depend on the unconverged values of the successor states, and Q-learning is better to be used later
 317 on because it is an off-policy TDL method that allows exploration at the successor states without affecting
 318 the previous ones. In this paper, the MCRL method is used at the beginning of the learning process and
 319 the Q-learning at the end of it. To reach an optimal policy, the ε -greedy strategy can be used as follows:

$$A_t = \begin{cases} \arg \max_a Q(S, a), & \text{with probability } (1-\varepsilon) \\ A \in_R \mathcal{A}(s), & \text{with probability } \varepsilon \end{cases} \quad (3)$$

320 where $\varepsilon \in [0, 1]$ is the exploration rate. Choosing the action with the highest Q-Value is called exploitation
 321 and choosing an action randomly is called exploration. It is important to keep a balance between both steps
 322 because exploitation helps in getting more rewards and evaluating the Q-Values based on the policy that
 323 appears to be the best. On the other hand, exploration allows exploring other actions that may result in
 324 higher Q-Values than the already explored actions.

325 3.2. The Intelligent Petri net method

326 Petri nets (PN) are directed graphs with two types of nodes, which are the *places*, represented by circles and
 327 the *transitions*, represented by rectangles. The number of *tokens*, which are depicted as black dots, contained
 328 in each of the places represents the state of the system. A PN is defined as a 5-tuple $\langle \mathbf{P}, \mathbf{T}, \mathbf{F}, \mathbf{M}_0, \mathbf{W} \rangle$
 329 [70], where $\mathbf{P} = \{p_1, p_2, \dots, p_{n_p}\}$ is the set of places, $\mathbf{T} = \{t_1, t_2, \dots, t_{n_t}\}$ is the set of transitions, $\mathbf{F} \subseteq$
 330 $(\mathbf{P} \times \mathbf{T}) \cup (\mathbf{T} \times \mathbf{P})$ is the set of arcs, $\mathbf{W} : \mathbf{F} \rightarrow \mathbb{N}_{>0}$ is the set of weights function, and $\mathbf{M}_0 : \mathbf{P} \rightarrow \mathbb{N}_{>0}$ is the
 331 number of tokens in each place initially, which is the initial markings.

332 The architecture of the PN can be summarised in the incidence matrix, $\mathbf{A} \in \mathbb{N}^{n_p \times n_t}$, which is the
 333 subtraction of the *backward incidence matrix* $\mathbf{A}^- = [a_{ij}^-]$ from the *forward incidence matrix* $\mathbf{A}^+ = [a_{ij}^+]$,
 334 where a_{ij}^- coincide with $\mathbf{W}(p_i, t_j)$, which is the weight of the arc from place p_i to transition t_j , and a_{ij}^+
 335 coincide with $\mathbf{W}(t_j, p_i)$, which is the weight of the arc from transition t_j to place p_i . The dynamics of the
 336 PN are controlled by the state of each transition, which manages the flow of tokens. Each transition has a
 337 set of input places, $\bullet \mathbf{P}_t$, referred to as the *pre-set places*, and output places, \mathbf{P}_t^\bullet , referred to as the *post-set*
 338 *places*. According to the *firing rule* in ordinary PNs, a transition, t_j , is said to be enabled once the markings
 339 of all its pre-set places are equal or greater than the weights of its pre-set arcs ($\mathbf{M}(p) \geq \mathbf{W}(t_j, p) \forall p \in \bullet \mathbf{P}_{t_j}$).
 340 Every enabled transition has the ability to fire, and this consumes tokens from its pre-set places and produces

341 tokens in its post-set places equal to the weights of the arcs connecting the places to the transition. This
 342 operation can be done for all transitions together in an efficient way using the *state equation* defined by:

$$\mathbf{M}_{k+1} = \mathbf{M}_k + \mathbf{A}^T \mathbf{u}_k \quad (4)$$

343 where k is the time step and $\mathbf{u} = [u_1, u_2, \dots, u_{n_t}]^T$ is the firing vector. More rules can be added to deal
 344 with the complexity of dynamic systems. In timed PN (TPN), a transition can't fire after it is enabled until
 345 a given delay, τ , passes. The value of τ can be deterministic or given by a probability density function,
 346 thus the PN is referred to as stochastic Petri Net (SPN). For high-level PN (HLPN), the logic flow is used
 347 in a wider manner by using flexible definitions of arc types, and tokens, along with transition firing rules
 348 to extend the basic formalism [71]. The HLPN definitions used in this paper are the *inhibitor arc*, which
 349 ends with a small empty circle, and the *reset arc*, which ends with a small filled circle. The *inhibitor arc*
 350 is connected from a place to a transition, and it disables the transition if the place has tokens equal to or
 351 more than the weight of the arc. The *reset place* is connected from a transition to a place, and it changes
 352 the marking of the place to a value equal to the weight of the arc once the transition fires [20].

353 Besides, function nodes, which are nodes with rhombus shapes are defined to perform some necessary
 354 calculations for the PN model. This definition allows modeling continuous aspects within the PN model,
 355 which is a discrete even model. A function node can come alone or after a transition. If a function is not
 356 connected to any transition, it is executed every change of state; whereas, if it comes after a transition, it is
 357 executed only when the transition fires.

358 To give the Petri net the ability to choose an optimum action, the intelligent PN (*iPN*) is used [25]. This
 359 variant introduces a finite set $\mathbf{G} = \{g_1, g_2, \dots, g_{n_g}\}$, named *action groups*, to the PN tuple to incorporate RL
 360 in decision making. Each action group, g_i , is composed of a finite set of transitions, $T_{g_i} \subseteq \mathbf{T}$, that represent
 361 decisions within a RL environment. Accordingly, RL selects which of the transitions will be enabled based
 362 on the rules described in [25]. It is important in this approach to distinguish between the RL states and the
 363 PN states. RL states are extracted from the RL environment and PN states are based on the markings of
 364 the PN.

365 3.3. Decomposing the PN into multiple subnets

366 This section proposes a method of decomposing the PN into multiple subnets without losing any func-
 367 tionality of the original net in order to reduce the computational cost. For any PN, the computational cost
 368 lies in getting the firing vector, \mathbf{u} ; whereas updating the state according to the state equation (Equation 4) is
 369 just a matrix multiplication, which is not computationally expensive. Calculating the firing vector requires
 370 checking the enabling conditions and then the firing conditions before assigning the firing state for every
 371 transition. If a PN is modeling multiple system functions, it will be known by the PN designer that groups
 372 of transitions will not be utilized for specific system states. Accordingly, it is possible to avoid checking

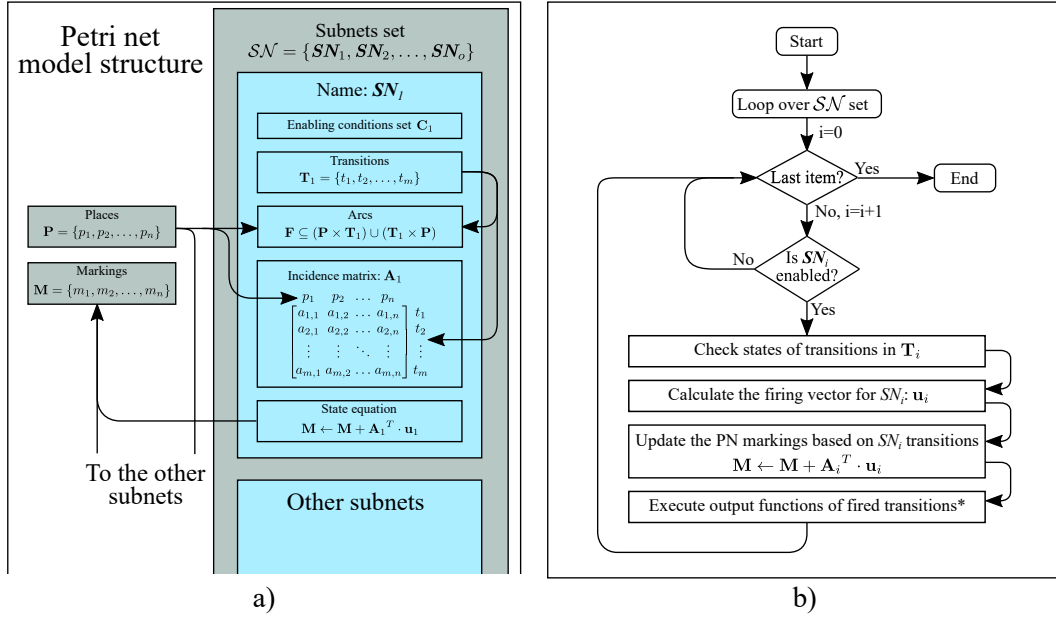


Figure 1: Panel a): PN structure. Panel b): Flowchart explaining the algorithm to update the state of the PN, when the PN is decomposed into multiple subnets.

373 these transitions to reduce computational costs without affecting the results because it is known that they
 374 will be disabled either way. To do so, this paper proposes the decomposition of the PN into subnets with
 375 each net having a set of conditions that enables it. Then, transitions in each subnet will be checked only
 376 after the subnet is enabled. In this new implementation, the places and their markings will be kept in the
 377 main net while all other information will be in the subnets. Accordingly, the PN structure is described by
 378 the tuple $\langle P, \mathcal{SN}, M_0 \rangle$ with $\mathcal{SN} = \{SN_1, SN_2, \dots, SN_i, \dots, SN_o\}$ being the set of subnets as shown in
 379 Figure 1a. Then, any subnet, SN_i , will be described by a tuple $\langle T_i, G_i, F_i, W_i, C_i \rangle$, with G being the set of
 380 the action groups that exist only if i PN is used, and C_i the set of conditions that enable the subnet. Since
 381 places remain common and are stored in the main net, every subnet will have an incidence matrix based
 382 on the connections between its own transitions and the places of the main net. This can be built in the
 383 same way as for an ordinary PN, but while considering only the transitions of the subnet. Based on this,
 384 the dynamics of the system will be described based on some additional rules shown below, and the process
 385 to update the state of the PN by calculating its markings is shown in the flowchart of Figure 1b:

- 386 • a subnet is enabled if it satisfies all its enabling conditions.
- 387 • if the subnet is enabled, all its transitions (and action groups in case of i PN) are checked, the firing
 388 vector of the subnet is calculated, and the state equation (Equation 4) is applied to update the markings
 389 of the main net based on the firing vector and incidence matrix of the subnet.

390 **Remark.** *The incidence matrices of subnets will contain many zero columns because the transitions of each*

391 subnet do not have connections with all the places. One might think that distributing places on the subnets
 392 or creating the incidence matrices of the subnets based only on the existing connections would be helpful to
 393 avoid these unnecessary connections and improve computational efficiency. However, once the places are
 394 distributed, there will be nothing to connect the subnets together, and this will require defining additional
 395 rules to solve this issue, which can add complication and computational cost. On the other hand, creating
 396 the incidence matrices based only on the active connections will change the dimensions of these matrices.
 397 This is like avoiding some places in each subnet, so it will be necessary to store the set of avoided places
 398 in each subnet to update only the markings of the included places every time the state equation (Equation
 399 4) is applied inside the subnet. This will result in additional computational costs that will be in most cases
 400 greater than the cost of multiplying by columns of zeros.

401 4. Extension of the *i*PN method for complex environments

402 This section provides a description of how the RL environment is divided into multiple environments to
 403 reduce the combinatorial state-action spaces. Also, it explains how experience can be shared among agents
 404 of different environments in the scope of the *i*PN and through action groups. These two ideas make the
 405 learning process faster and reduce computational costs.

406 4.1. Dividing the RL environment to multiple environments

407 Multi-agent Reinforcement learning (MARL) methods are concerned with the cases of multiple agents
 408 interacting in the same environment. These methods can be cooperative, where the agents try to achieve a
 409 common goal, or competitive, where they try to compete to see who achieves more. In some cases, a mixed
 410 environment can be created, where agents form groups, cooperating within each group and competing
 411 against other groups. In this study, the main focus is on cooperative MARL methods to optimise systems
 412 with multiple tasks.

413 MARL differs from single-agent RL in that the environmental state and reward function that each
 414 agent receives is a function of the joint actions of all agents. For this, each agent has to consider other
 415 agents' actions in addition to the environment. The process of taking multiple decisions is usually modeled
 416 through a stochastic game [72], also known as a Markov game [73]. A stochastic game is a multi-decision
 417 extension of the MDP and can be described by the tuple $\langle \mathcal{S}, \mathcal{A}, P_d, \mathcal{R} \rangle$; where \mathcal{S} is the set of the *states* of the
 418 *environment*, $\mathcal{A} = \mathcal{A}^1 \times \dots \times \mathcal{A}^n$, where \mathcal{A}^i is the set of agent *i* *actions*, $P_d = \Pr\{S_{t+1} = s' | S_t = s, \mathbf{A}_t = \mathbf{a}\}$
 419 is the transition probability function from state *s* to state *s'* in the next state while taking the joint action
 420 \mathbf{a} , and $\mathcal{R} = \mathcal{R}^1, \dots, \mathcal{R}^n$, with $\mathcal{R}^i = \Pr\{R_{t+1} = r' | S_t = s, \mathbf{A}_t = \mathbf{a}, S_{t+1} = s'\}$ is the reward probability
 421 function for agent *i* after transitioning from state *s* to state *s'* while taking the joint action \mathbf{a} . Accordingly,
 422 each agent will have a Q-Value function of the state and the joint action, and reaching a globally optimum

423 policy requires coordination between agents [74]. However, considering joint action results in an exponential
424 increase of the action space at each state [27]. Besides, trying to optimise several agent’s policies in one
425 problem requires the definition of an environment that considers the important aspect for every agent, which
426 results in a great increase in the state space. The huge state-action spaces make the computational costs
427 extremely high for already complex problems.

428 The necessity to use MARL methods is when agents cooperate in the same environment [27]. Indeed,
429 there exist some problems where it is possible to divide the environment into several sub environments, but
430 without being able to divide the problem into multiple independent optimisation problems. This happens if
431 the conditions and decisions of each environment do not directly affect the rewards of other environments,
432 but can affect the transition probabilities or other aspects. In this paper, the RL environment is decomposed
433 into multiple environments with each one having its own states, available actions, reward functions, and a
434 single agent. Thus, the problem drops back to single-agent RL, but with multiple interacting environments.
435 Environments can intersect to keep some information commonly available to all the agents, and in this way,
436 agents can cooperate explicitly.

437 4.2. Sharing agents experience through similar action groups

438 A key aspect of problems with multiple agents is the allowance for experience sharing between agents
439 that are solving similar tasks to learn faster and better [74]. This still applies to the case of multiple agents
440 optimising multiple environments if these environments share similar characteristics.

441 Multiple environments may require the same decisions at similar conditions if these environments are
442 similar. For example, if two agents are optimising the maintenance of two identical components while
443 considering each component as a separate environment of a system, it is expected that the two components
444 require the same action if they were in the same conditions. This means that the two components can
445 follow the same policies for the same decisions. In the RL formulation, the policy is directly related to the
446 Q-Values. Thus, a way to exchange information between agents of similar environments is to assign the same
447 Q-Values for similar actions. By doing this, any update in the Q-Value of an action in any environment will
448 cause the update of similar actions in the other environments.

449 In the case of *i*PN, actions are equivalent to enabling transitions inside an action group. Each time an
450 action group is enabled, the agent receives a representation of the environmental state. If the state is new,
451 it is created automatically with its available actions in that environment, and by this, the RL environment
452 is populated by the states and actions. To link the Q-Values of similar actions, the following steps can be
453 performed:

- 454 • let a similar group set, $\mathcal{SG} = \{g_1, g_2, \dots, g_n\}$, be a set of action groups that require similar policies
455 and represent similar decisions.

- 456 • all the action groups should have the same number of ordered transitions that represent *similar actions*.
- 457 Thus, every action will have similar actions in similar action groups.
- 458 • if a decision is required in g_l and the state, S , is new, create a state in the environment of each action
- 459 group in \mathcal{SG} , and not only in the environment of g_l .
- 460 • any action, $A_{l,i} \in g_l$ will have the same Q-Value as its similar actions in the other action groups:
- 461 $Q(S, A_{1,i}) = \dots = Q(S, A_{l,i}) = \dots = Q(S, A_{n,i})$
- 462 • any time $Q(S, A_{l,i})$ is updated, all the Q-Values of similar actions are updated.

463 5. Case study: optimising the maintenance of ballast in multiple railway track sections for

464 optimal railway operations.

465 In this section, an i PN has been developed to model various aspects of a railway consisting of ten railway
 466 sections while optimizing the maintenance of its ballast. Each section has a length of three position keys,
 467 where a position key is a 220-yard length of track known as a Poskey. The track speed of the considered
 468 section is less than 20 MPH, the annual usage is 20 Equivalent Million Gross Tonnage (EMGT), and all
 469 sleepers are of small concrete type.

470 The maintenance of the sections is assumed to be carried out through two types of maintenance actions:
 471 condition-based maintenance and opportunistic maintenance. For condition-based maintenance, the decision
 472 is taken after updating the condition of the section following each inspection, which is assumed to be periodic
 473 every six months. On the other hand, opportunistic maintenance decisions can be made for a section when
 474 maintenance equipment is available on-site to perform maintenance for any other section.

475 Once a decision is made to repair a section based on its condition, the maintenance team prepares for
 476 the maintenance, travels to the site, and performs the required maintenance. If the equipment is available
 477 on-site, a decision to repair other sections can be made, referred to as opportunistic maintenance. This
 478 maintenance approach allows repairing multiple sections at once to save preparation and travel costs, even
 479 if these sections do not urgently require maintenance based on their conditions.

480 The maintenance actions for sections are performed in series, starting from the first section and ending
 481 with the last one, assuming that only one maintenance team will perform the required work. Each section
 482 can be in one of five conditions denoted as 'E' (Excellent), 'VG' (Very-good), 'G' (Good), 'P' (Poor), and
 483 'SR' (Super-red). These conditions reflect the safety and stability of the railway, and they are directly related
 484 to the standard deviation of the vertical geometry profile, as will be seen in the next section. Experts suggest
 485 that the "super-red" condition should be avoided at all costs. This is because it poses a significant safety
 486 risk to railway operations and can lead to speed restrictions, causing train delays and potential fines. It is,

487 therefore, crucial to take preventative measures to minimize the likelihood of encountering this condition,
488 as it can have negative impacts on both safety and efficiency.

489 The problem is to find the optimal maintenance decision for each condition of each section. Two policies,
490 named policy A and policy B, were proposed as base cases to find the optimal maintenance decisions.
491 Policies A and B allow maintenance to be performed each time the 'Very-good' and the 'Good' conditions
492 are reached respectively. Also, both policies follow the same sequence of maintenance actions, which is 7
493 tamping actions, 3 stoneblowing actions, and then renewal. The third policy is optimised by using RL, with
494 the details of the RL inputs as described in Section 5.3. For this study, it is assumed that tamping becomes
495 ineffective after 7 actions, and stoneblowing becomes ineffective after 3 actions. Thus, the maximum number
496 of allowed tamping and stoneblowing actions are 7 and 3 actions respectively. Besides, it is not allowed to
497 perform tamping after a stoneblowing action is performed. This will prevent stoneblowing actions from
498 being before tamping to follow what is done in reality. Before introducing the PN model, the following
499 section provides an introduction to railway modelling and the formulas used.

500 *5.1. Track degradation and maintenance modelling*

501 The railway consists of several interacting assets, mainly Plain Line (PL) track and Switches and Cross-
502 ings (S&Cs). These are made up of components that have different degradation, inspection, and maintenance
503 mechanisms. The track is made up of the rail to provide guidance and a smooth running surface, sleepers
504 to support the rail at the correct gauge and inclination and to transmit loads to ballast, rail pads to provide
505 electrical insulation and distribute loads on the sleepers, ballast to support sleepers at the correct level,
506 spread forces into the formation, and allow surface water drainage, formation to support ballast and collect
507 water to the drainage system, subgrade, which is the natural layer where all other parts are built on, and
508 the drainage system to convey water away from the track. The primary focus of the case study is on the
509 ballast and rail, but other parts are considered if they are linked to these two parts.

510 The ballast is good when it is composed of crushed angular hard rocks and stones, free of dirt and dust,
511 uniformly graded, and not prone to cementing action [2]. The degradation mechanism of ballast is called
512 fouling, which occurs when small particles build up within the ballast. Causes of ballast fouling can be
513 the ballast breakdown, sleeper wear, and the infiltration from the surface, underlying granular layers, or
514 subgrade [2]. Ballast fouling can result in a saturated subgrade and wet beds because it impedes water
515 drainage [75]. This leads to differential track settlement because of the uneven distribution of loads. The
516 ballast can be maintained either by tamping, stoneblowing, or renewal techniques, which restores the track
517 geometry to a better condition [76]. Ballast tamping is the common form of correcting the track geometry.
518 It is done using specialised trains which lift the rail with the sleepers to the target level. Then, tamping
519 tines are inserted and vibrated to squeeze the ballast under the gap, recovering the correct level of the rail.
520 This process causes significant breakage of ballast particles, which can result in highly fouled ballast. When

521 this is the case, tamping can no longer be effective and stoneblowing can be considered since it causes much
522 less breakage of ballast and can be used even if the ballast is fouled [77]. Stoneblowing is performed using
523 trains that lift the rail and the sleepers to the target level. Then, they insert tubes that use compressed
524 air to blow a measured quantity of clean ballast into the gap below the sleepers, leaving the rail at the
525 correct level after the tubes are removed. The only disadvantage of stoneblowing is that it is slower and
526 more expensive than tamping [77]. Thus, a rule of thumb is to perform tamping operations until the ballast
527 becomes highly fouled and to use stoneblowing after that [77]. However, after reaching a critical level of
528 ballast fouling, maintenance actions become less effective, and renewal should be performed. Renewal can
529 be performed by cleaning the ballast and reusing a portion of it mixed with new ballast, or by performing
530 a complete replacement of the old ballast.

531 Consequently, maintenance activities impact the quality of ballast, leading to accelerated degradation
532 rates and settlement. The extent of track settlement can be quantified through the utilisation of specialised
533 trains, such as the *New Measurement Train* operated by *Network Rail* (NR), which employs laser scanning
534 technology to assess changes in track geometry every 0.2 metres as the train progresses along the track.
535 The train's journey is affected by variations in the track profile, with long and smooth undulations having
536 minimal impact on train safety and comfort and are thus disregarded in assessments. The most commonly
537 used metric for measuring track settlement is the vertical standard deviation (SD) of a set of measurements
538 taken for each Poskey along the track [78]. This is due to the vertical geometry being the most prone to
539 degradation and having the greatest influence on ride quality and maximum permissible speed. On a typical
540 European track, a vertical SD of less than 5.2, 7.4, 8.3, 9.9, or ∞ is classified as "Excellent," "Very Good,"
541 "Good," "Poor," and "Super-red," respectively, for track speeds below 20 MPH.

542 The vertical geometry profile of the track is expected to improve as a result of maintenance activities
543 performed on the ballast. However, as the quality of the ballast deteriorates, the ability of maintenance
544 actions to correct the rail level becomes diminished. As the fouling index of the ballast increases, small
545 particles fill the voids between rocks resulting in denser ballast. Maintenance activities may temporarily
546 create voids between rocks, but the resistance to loads becomes weak, and as a result, small particles tend
547 to quickly fill these voids once the track is subjected to loads [77]. This in turn makes the maintenance
548 actions less effective.

549 The irregularities of the track surface have a substantial effect on the incidence of rail faults, in addition
550 to the safety and quality of the ride [79]. As trains traverse the rail network, they exert substantial forces on
551 the rails, which can result in a wide spectrum of defects and faults [80]. Rail corrugations, for instance, can
552 impair the quality of the ride and accelerate the degradation of many track and vehicle components [79].
553 On the other hand, head wear can decrease the rail-wheel interface area and reduce the grip for braking
554 and accelerating, thereby increasing the likelihood of faults [81]. In the event of a break, increased forces
555 are imposed on the surrounding parts, and speed restrictions may be necessary to maintain safety, leading

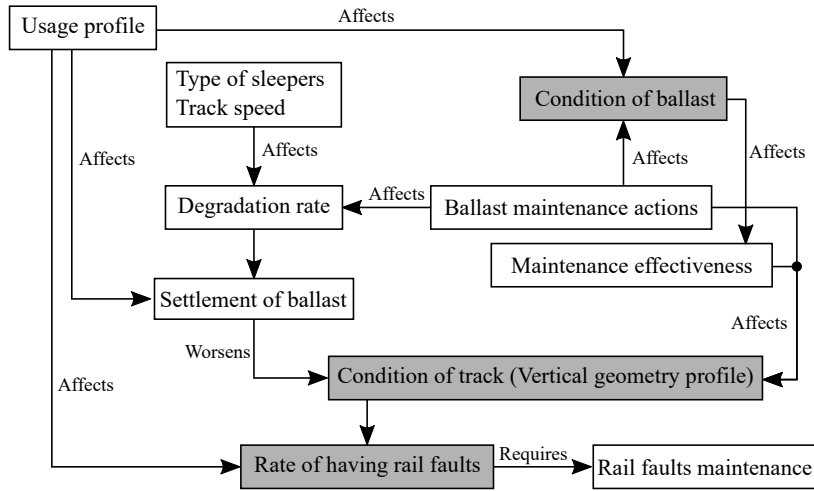


Figure 2: Summary of the railway relations that are considered in this study.

556 to delays and financial losses [82]. It is crucial to maintain appropriate vertical geometry to ensure the safe
 557 and efficient operation of the rail network and to minimise the risk of rail faults. Correction of rail faults
 558 may involve grinding or welding if the fault is not severe, while replacement of the rail may be required in
 559 more severe cases.

560 The relationships between railway infrastructure components, such as the track condition expressed in
 561 terms of the vertical geometry profile and the degradation of the ballast, are summarised in Figure 2. A
 562 degradation in track condition due to ballast settlement leads to an increase in the frequency of faults
 563 in the rail and necessitates more frequent maintenance activities. Conversely, while maintenance of the
 564 ballast can improve track condition, it can also result in ballast fouling and faster degradation, reducing
 565 the effectiveness of maintenance efforts. Thus, it is imperative to have a comprehensive understanding of
 566 these interrelated factors for effective railway infrastructure management and the maintenance of safe and
 567 efficient rail operations. To do so, models were created for the degradation rate of the track, the ballast
 568 maintenance effectiveness, and the rate of rail faults and their maintenance frequencies based on data from
 569 a typical European track. The data covered track geometry, fault and maintenance records of a European
 570 national rail network operator over a period of approximately 8 years, and also covered all uses from freight
 571 to passenger services and all levels of track speeds. Although the data are treated confidentially here,
 572 representative models were chosen to demonstrate the techniques introduced in this paper.

573 Several factors were tested to determine which of them affect the degradation rate and it was found that
 574 the type of sleepers, the speed of the track, and the maintenance history have an impact. A stochastic model
 575 was built relating the degradation rate (mm/EMGT) to these three factors, proposing a Weibull distribution
 576 for each entry. The study considered a track speed of less than 20 MPH with small concrete sleepers, and
 577 Table 1 summarises the parameters of the Weibull distributions for each maintenance history required for

Table 1: Shape parameter, β , and size parameter, η , of Weibull distributions for the degradation rate [m/EMGT] of small concrete sleepers with track speeds 5-20MPH and different maintenance histories.

After:	1 st Renewal	1 st tamp	3 rd tamp	5 th tamp	6 th tamp	7 th tamp	1 st stoneblowing
η	1.93E-04	2.15E-04	2.22E-04	2.06E-04	2.21E-04	2.37E-04	2.12E-04
β	1.03E+00	8.18E-01	1.05E+00	1.17E+00	7.25E-01	1.00E+00	6.49E-01

Table 2: The vertical geometry SD [mm] after each maintenance activity representing the maintenance effectiveness

Renewal	Tamping				StoneBlowing		
1 st	1 st	3 rd	5 th	6 th	1 st	5 th	7 th
0	1	1.5	3	3.5	1	2.5	3.5

578 modelling the degradation rate. The values of β and η are of the order of magnitude expected for this type
579 of track and the observed variation in these values with changing maintenance history is, to some extent,
580 due to the variation in the use of low speed track, which includes, for example, sidings and access to depots
581 used by freight trains. This model allows the calculation of the SD of the track after a certain usage by
582 using the following formula:

$$SD^* = SD + DR \cdot (U^* - U) \quad (5)$$

583 where DR represent the degradation rate that can be sampled from the distributions of Table 1 and U
584 represent the usage in EMGT. The variables with and without asterisks represent the current and the last
585 state respectively.

586 Conversely, ballast maintenance actions have the potential to improve the SD of the track. Nevertheless,
587 the frequency of maintenance actions increases the rate of settlement as previously discussed. In light of
588 this, Table 2 presents some assumed values that depict the impact of each maintenance action on the SD
589 of the track in relation to the maintenance history of the track.

590 The analysis revealed a strong correlation between the vertical geometry SD of the track and the rate of
591 faults, leading to the creation of a model that calculates the normalised rate of each rail fault as a function of
592 the SD . The rate of faults increases with the length and usage of the track, thus normalisation was deemed
593 necessary. Data recordings of various rail faults including Squat, Tache Ovale, Bolt Hole, Weld, Other,
594 Rolling Contact Fatigue (RCF), Wheel burn, Lipping, Side Wear, Head Wear, Corrugation, and Unknown
595 faults were analysed. To calculate the rate of each fault, the data was stacked and lines were fitted to the
596 data, starting with squats and adding faults one by one, simplifying the decisions within the model.

597 Algorithm 1 outlines the procedure for determining the rate of occurrence of faults within each stacked
598 group. To begin, a list of faults, denoted as FL , is established, and a set of faults stacked groups, denoted

Algorithm 1 Calculation of the probability of having a fault and its type

- 1: Define List of faults, $FL = ["Squat", "Tache Ovale", "Bolt Hole", "Weld", "Other", "Rolling Contact Fatigue (RCF)", "Wheel burn", "Lipping", "Side Wear", "Head Wear", "Corrugation", "Unknown"]$. Define stacked sets based on the order of FL are defined as: $FS_1 = \{FL[1]\}$, $FS_i = FS_{i-1} \cup \{FL[i]\}$ $\forall i \in \{2, \dots, 12\}$.
- 2: The rate of having a fault from each of the lists $FS_i \forall i \in \{1, \dots, 12\}$ can be calculated based on the following equation:

$$FR_i[/\text{poskey}/\text{EMGT}] = A_i \cdot SD^3 + B_i \cdot SD^2 + C_i \cdot SD \quad (6)$$

with A_i, B_i , and C_i for each set can be found in Table 3.

Function 1 – Fault type based on SD:

- 3: choose $R \in_R [0, 1]$. Then, $R \rightarrow R/(L \cdot \Delta U)$, with L and ΔU being the length in Poskeys and usage in EMGT respectively.
- 4: **if** $R < FR_{12}(SD)$ **then** ▷ having fault is probable because FS_{12} contain all faults
- 5: **for** $i \in \{1, \dots, 11\}$ **do**
- 6: **if** $R < FR_i(SD)$ **then**
- 7: **Return** FL_i ▷ fault i from list FL
- 8: **Return** FL_{12} ▷ unknown fault
- 9: **else**
- 10: **Return** \emptyset ▷ no fault

Function 2 – Correction type based on the fault type:

- 11: Maintenance types list is $MT = ["rerail", "weld", "grind or other"]$
 - 12: any fault type, i , has 3 stacked probabilities, $SP_{i,1}$, $SP_{i,2}$, and $SP_{i,3}$, that stands for the elements of MT respectively and stored in Table 4.
 - 13: To know which maintenance type corresponds to the fault, choose $R \in_R [0, 1]$
 - 14: **for** $j \in \{1, 2, 3\}$ **do**
 - 15: **if** $R \leq SP_{i,1}$ **then**
 - 16: **Return:** MT_j
-

Table 3: Parameters for the calculation of the fault rate of different sets. $FL_{1,\dots,12}$ are the stacked sets of faults defined in Algorithm 1

FS	A	B	C	FS	A	B	C
1	7.64E-05	-6.45E-04	3.44E-03	7	1.12E-04	-3.87E-04	4.93E-03
2	7.85E-05	-6.55E-04	3.73E-03	8	7.34E-05	1.31E-04	3.87E-03
3	6.90E-05	-5.45E-04	3.56E-03	9	1.66E-04	-3.48E-04	4.61E-03
4	9.38E-05	-7.96E-04	4.85E-03	10	1.61E-04	-1.98E-04	4.28E-03
5	1.18E-04	-6.94E-04	5.22E-03	11	1.60E-04	-1.92E-04	4.27E-03
6	1.13E-04	-5.46E-04	5.07E-03	12	1.70E-04	-2.31E-04	4.33E-03

Table 4: Stacked probabilities for each maintenance action based on fault type. FL is the list of faults defined in Algorithm 1

FL	Rerail	Weld	Grind or other	FL	Rerail	Weld	Grind or other
1	0.328	0.954	1	7	0.267	0.874	1
2	0.722	0.963	1	8	0.029	0.134	1
3	0.9	0.919	1	9	0.044	0.338	1
4	0.519	0.904	1	10	0.213	0.752	1
5	0.641	0.918	1	11	0.706	0.765	1
6	0.508	0.786	1	12	0.464	0.63	1

599 as FS , is constructed based on the elements of FL . For instance, if the first, second, and third elements of
600 FL are Squat, Tache Ovale, and Bolt Hole respectively, then FS_1 consists of only Squats, FS_2 encompasses
601 Squats and Tache Ovale, and FS_3 encompasses Squats, Tache Ovale, and Bolt Hole.

602 Subsequently, a third-order polynomial function is utilised to model the fault rate of each stacked group.
603 The parameters of the fitted functions are listed in Table 3. Using these rates, the probability of encountering
604 a fault after a specified usage over a certain track length can be calculated as depicted in the first function
605 of Algorithm 1.

606 The second function of Algorithm 1 presents a method for sampling a correction for the fault from the
607 available options. This function is based on the values in Table 4, which are derived from the stacked rates
608 of correction methods for each type of fault. For instance, for the RCF, which is the sixth element of FL
609 list, the rate of performing "Rerail", "Welding", and "Grinding and other" methods are 0.508, (0.789-0.508),
610 and (1-0.786) respectively, as demonstrated in Table 4. These rates serve as sampling probabilities for the
611 correction methods.

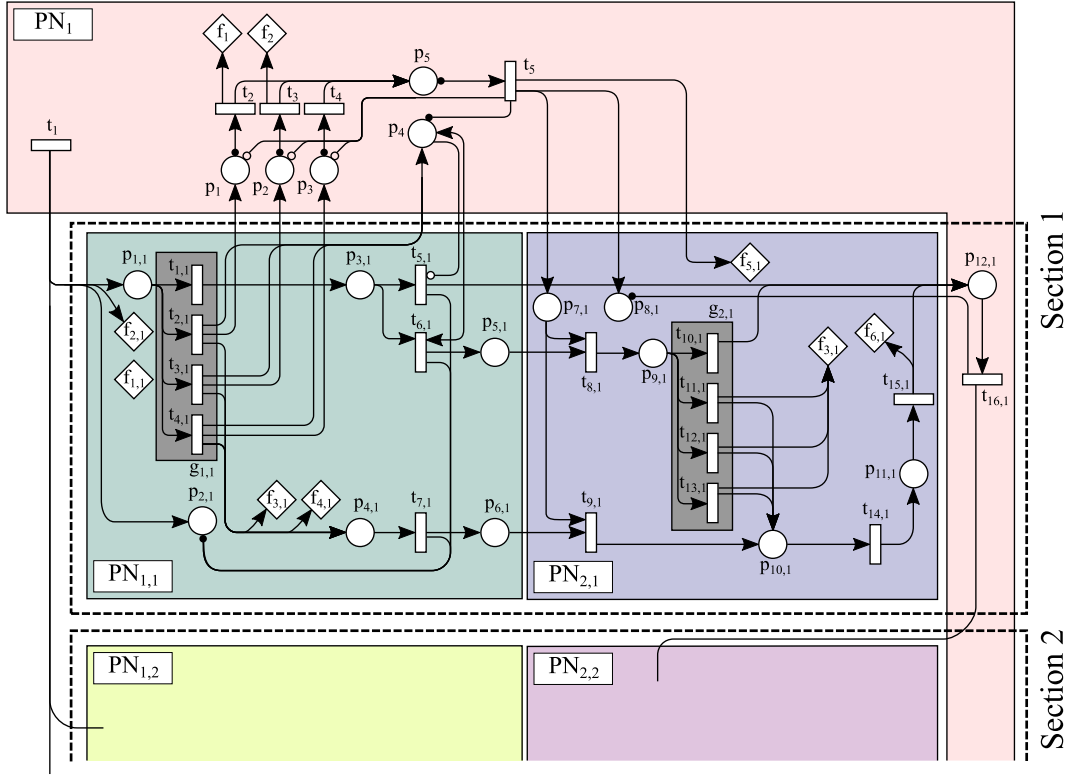


Figure 3: The iPN for modelling and optimising the railway maintenance and operation.

612 5.2. Railway iPN model

613 The aim of the railway iPN model is to create an expert decision support system (DSS) that helps
614 in finding an optimal maintenance strategy for the ballast and rail taking into consideration the working
615 conditions of the railway. Figure 3 shows the subnets that form the PN model, where node descriptions
616 are provided in Table 5 and function descriptions in Figure 4. The railway is composed of 10 identical
617 sections, with each section being modeled by two subnets. The names of the subnets of an arbitrary section
618 $i \in \{1, 2, \dots, 10\}$ are $PN_{1,i}$ and $PN_{2,i}$. An additional subnet called PN_1 models the common activities for
619 all sections. Figure 3 does not show the subnets of sections 2, \dots , 10 due to lack of space. However, these
620 subnets are identical to the ones of the 1st section. As with the subnet names, nodes that are in the common
621 Subnet have one-number subscripts while the ones in the other subnets have two-number subscripts with
622 the 2nd number referring to the ID of the section. The decomposition method provided in Section 3.3 can be
623 used to reduce the computational costs as explained. If the decomposition method is used, the conditions to
624 enable subnet $PN_{1,i}$ or $PN_{2,i}$ is to have a token in place $p_{1,i}$ or $p_{8,i}$ respectively, and PN_1 is always enabled
625 without the need for any conditions.

626 Initially, all places are unmarked, time is equal to 0, and all the sections are in Excellent state. The
627 information about each section is not represented by the PN places but is calculated through the functions

628 described in Figure 4 as explained below. Function $f_{1,i}$ does not have any input arcs, so it runs every time
 629 the state of the PN changes. This function updates the condition, calculates the probability and the cost
 630 of having a fault, and calculates the reward based on the condition of the section. The dynamics of the
 631 problem start with transition t_1 , which is a timed transition representing the inspection. This transition
 632 fires every 0.5 yrs. to mark $p_{1,i}$ and $p_{2,i}$ and execute function $f_{2,i}$. The function $f_{2,i}$ updates the available
 633 actions by excluding the non-available actions from action group $g_{1,i}$. These actions are *no-action*, *tamping*,
 634 *stoneblowing*, and *renewal*, represented by transitions $t_{1,i}$, $t_{2,i}$, $t_{3,i}$, and $t_{4,i}$ in $g_{1,i}$, respectively. After 7
 635 tamping actions, tamping cannot be chosen for the ballast due to the reached fouling level, and after 3
 636 stone blowing actions, the only option becomes renewal. Thus, based on the maintenance history, available
 637 actions are updated.

638 By marking $p_{1,i}$, action group $g_{1,i}$ becomes enabled, representing the need to make a maintenance
 639 decision. Accordingly, the RL agent selects one of the available transitions from $g_{1,i}$. If $t_{1,i}$ is chosen, $p_{3,i}$
 640 will be marked to indicate that no maintenance decision was made. On the other hand, if $t_{2,i}$, $t_{3,i}$, or $t_{4,i}$
 641 is triggered, $p_{4,i}$ will be marked. $f_{3,i}$ and $f_{4,i}$ will be executed, and p_1 , p_2 , or p_3 will be marked to indicate
 642 that the preparation for tamping, stoneblowing, or renewal, respectively, has commenced.

643 Each maintenance action consists of two distinct steps: preparation and travel, followed by the actual
 644 maintenance itself. The duration for executing each of these steps is determined using the functions denoted
 645 as $f_{3,i}$ and $f_{4,i}$ (Figure 4). Function $f_{3,i}$ is responsible for updating the actual maintenance time, influencing
 646 the progression of maintenance time transition ($t_{15,i}$), and aggregating RL rewards based on the associated
 647 maintenance costs. Meanwhile, function $f_{4,i}$ pertains to the adjustment of maintenance preparation time,
 648 governing the timing of preparation time transitions ($t_{2,i}$, $t_{3,i}$, $t_{4,i}$), and accumulating RL rewards linked to
 649 the preparatory expenses.

650 A single maintenance preparation can effectively address the repair needs of multiple sections when they
 651 require the same type of maintenance. Transitions t_2 , t_3 , and t_4 model the durations for maintenance
 652 preparations: t_2 for tamping, t_3 for stoneblowing, and t_4 for renewal. Places p_1 , p_2 , and p_3 indicate that
 653 preparations are underway. All preparation actions must be completed in order for the maintenance of the
 654 first section to commence. The initiation of maintenance is represented by transition t_5 . As depicted in the
 655 PN model, this process involves inhibiting transition t_5 with places p_1 , p_2 , and p_3 , ensuring that maintenance
 656 cannot start until all preparation tasks are finished.

657 Following the commencement of maintenance (triggered by the firing of t_5), the option to repair sections
 658 not previously selected for repair becomes available. This can be seen as a form of opportunistic maintenance.
 659 An advantage of deciding to repair a section on-site is the utilization of available maintenance trains to repair
 660 additional sections, thus saving costs associated with preparation and travel. To indicate the availability of
 661 maintenance trains, functions f_1 and f_2 are executed upon the firing of transitions t_2 and t_3 , respectively,
 662 signifying the availability of tamping or stoneblowing trains.

663 The actual maintenance comes after finishing the preparation and reaching the site. It starts after
 664 transition t_5 is fired to mark $p_{7,i}$ and $p_{8,i}$. If $p_{5,i}$ is marked, $p_{7,i}$ enables $t_{8,i}$ indicating that a decision not
 665 to repair section i was taken. Then, $t_{8,i}$ fires to allow a decision through action group $g_{2,i}$ regarding the
 666 opportunistic maintenance. If the opportunistic decision was not to repair section i , $t_{10,i}$ fires to mark $p_{12,i}$
 667 indicating that this section is finished; whereas, if a maintenance action was decided, $f_{3,i}$ will be executed
 668 and $p_{10,i}$ will be marked indicating that the maintenance can start. On the other hand, if the section
 669 was already decided to be maintained, $p_{6,i}$ and $p_{7,i}$ will allow $t_{9,i}$ to fire, which marks $p_{10,i}$, indicating the
 670 possibility to proceed directly in performing the maintenance. The time taken to perform maintenance
 671 is modeled by transition $t_{14,i}$ while its effectiveness is modeled by $f_{6,i}$. It is assumed that maintenance
 672 is only done during non-working hours of the train, which means that maintenance can be a number of
 673 interrupted intervals. Working hours can cause further degradation of the non-reached parts of the section
 674 before the maintenance ends. This is why $f_{6,i}$ is not directly executed by $t_{14,i}$ but by $t_{15,i}$, allowing for
 675 $f_{1,i}$ to account for degradation that occurs during the maintenance period. The firing of $t_{15,i}$ marks $p_{12,i}$,
 676 indicating that this section is finished. After the i^{th} section is finished, $t_{16,i}$ fires to mark $p_{7,i+1}$ and $p_{8,i+1}$,
 677 and the maintenance of the next section starts.

678 5.3. RL inputs

679 To optimize track maintenance, multiple optimization problems are tackled by breaking down the track
 680 into separate RL environments. Each section of the track is treated as an individual environment, complete
 681 with its own agent, states, actions, and value function. As these environments make up the same system,
 682 each one impacts the transition probability function of the others. Reward functions are specific to each
 683 environment and are solely dependent on the agent's decisions within that environment. The optimization
 684 problems are episodic and terminate when the decision to renew the section ballast is made, as this represents
 685 a new investment and a fresh start. The goal of the agent is to increase rewards, which equate to the cost
 686 function and are defined in terms of revenues and expenses. This approach requires the agent to maximize
 687 the use of the section's ballast before renewal.

688 5.3.1. Definition of the rewards function (cost function)

689 Rewards in RL play a crucial role in guiding the RL agent towards optimizing the maintenance strategy.
 690 In the context of the railway industry, these rewards can be expressed in monetary terms, as the railway
 691 companies aim to ensure financially efficient operation while ensuring appropriate levels of safety. The
 692 rewards can be either positive, representing revenues, or negative, representing costs. However, due to
 693 the commercial sensitivity of costs and revenues, it is not possible to represent rewards in terms of actual
 694 monetary values. Therefore, the rewards are presented in unitless forms while maintaining their realistic
 695 values relative to each other through consultation with railway experts.

Table 5: The description of the *i*PN nodes.

Node	Description
$PN_{1,i}$	
$p_{1,i}$	taking maintenance decision
$t_{1,i}, t_{2,i}, t_{3,i}, t_{4,i}$	no-action, tamping, stoneblowing, and renewal decisions respectively
$p_{2,i}$	subnet key
$p_{3,i}, p_{4,i}$	no-action or a maintenance action is chosen respectively
$t_{5,i}$	opportunistic maintenance is not possible, reset the PN key
$t_{6,i}, p_{5,i}$	opportunistic maintenance is possible, reset the PN key
$t_{7,i}, p_{6,i}$	Perform maintenance after finishing the preparation, reset the PN key
$PN_{2,i}$	
$p_{7,i}$	site is reached
$p_{8,i}$	subnet key
$t_{8,i}, p_{9,i}$	no maintenance was decided, check for opportunistic maintenance
$t_{9,i}$	maintenance was decided, proceed
$t_{10,i}, t_{11,i}, t_{12,i}, t_{13,i}$	no-action, tamping, stoneblowing, and renewal decisions respectively (opportunistic)
$p_{10,i}$	ready to perform maintenance
$t_{14,i}$	models the time taken to finish repairing the section, which is controlled by function $f_{3,i}$, and allows for $f_{1,i}$ to consider the effect of further degradation
$p_{11,i}, t_{15,i}$	models the maintenance effectiveness and update condition
PN_1	
t_1	inspection, with delay equal to 0.5
$(p_1, t_2), (p_2, t_3), (p_3, t_4)$	tamping, stoneblowing, and renewal preparation respectively with the delay of transitions controlled by function $f_{4,i}$
p_4	there exist a section that will be maintained
p_5	one or more preparations are finished
t_5	site is reached, ready for doing maintenance
$p_{12,i}, t_{16,i}$	section is finished move to the next section

$f_{1,i}$	if $t_S^* > t_S$ then ($t_S^* = t$ is the current time and t_S is the last time the section was updated) Calculate the usage, U^* , at the current time, t_S^* , based on the usage rate: $U^* = 20t_S^*$ Get SD^* using Equation 5. Get, R_c , the continuous reward between the two states of the section based on Algorithm 2. Accumulate the reward of the last RL state: $R_t \rightarrow R_t + R_c$ Update the variables: $U = U^*, SD = SD^*, t_S = t_S^*$
-----------	---

$f_{2,i}$	Update the available actions for $g_{1,i}$ according to the following rules: tamping is not allowed after 7 tamps or after stoneBlowing. stoneblowing is not allowed more than 3 times
-----------	--

$f_{3,i}$	sample the output rate, OR [yrds/hr.], from $\mathcal{W}(1.28, 249.27)$ for tamping and $\mathcal{W}(1.30, 237.26)$ for stoneblowing Convert OR to [Poskeys/hr.]: $OR \rightarrow OR/220$ Calculate the maintenance time, t_M [yrs.], assuming 2,080 working hours per year: $t_M = OR \cdot L/2080$ Get the actual maintenance costs based on the section's length, C_m Accumulate the reward of the last RL state: $R_t \rightarrow R_t - C_m$ Update the maintenance history
-----------	--

$f_{4,i}$	The time to prepare for maintenance and reach the site, t_p , will be equal to: 1 night if the condition is <i>super-red</i> , 1 week if the condition is <i>poor</i> , 2 weeks if the condition is <i>good</i> , and 1 month otherwise. Accumulate the reward of the last RL state: $R_t \rightarrow R_t - C_p$
-----------	---

$f_{5,i}$	Update the available opportunistic maintenance actions for $g_{2,i}$ based on $f_{2,i}$ rules while considering the prepared transitions
-----------	--

$f_{6,i}$	Update SD based on the maintenance effectiveness (Table 2) Update DR from the distributions in Table 1 If the maintenance action is a renewal, terminate the old RL episode and start a new one
-----------	---

f_1	Indicate that the preparation for tamping maintenance is done
-------	---

f_2	Indicate that the preparation for stoneblowing maintenance is done
-------	--

Figure 4: Description of the functions used in the *i*PN model.

696 The study examines the direct and indirect effects of RL decisions on the rewards function. To construct
697 the reward function, various effects are considered, including maintenance and renewal costs, preparation
698 and travel costs, possession costs, delay costs, catastrophe costs, and revenues. Directly influenced by the
699 RL agent's decisions are the costs associated with ballast maintenance, preparation, and travel. Ballast
700 maintenance costs are incurred on a per-section basis, while preparation and travel costs are paid once to
701 address multiple sections that undergo maintenance at the same time in close locations. The RL agent
702 can strategically choose to repair multiple sections simultaneously, known as opportunistic maintenance, in
703 order to minimize the expenses related to preparation and travel.

704 Indirectly affected by the RL decisions are the costs that depend on the condition of the track. When
705 the track's condition deteriorates, the likelihood of rail faults increases, resulting in higher maintenance
706 costs required to address these faults [79]. These costs are incurred for each track section, as outlined in
707 the effects summary provided in Table 6. It is important to note that this study does not specifically focus
708 on optimizing decisions related to this particular maintenance type; instead, it considers it as part of the
709 overall costs influenced by decisions concerning ballast maintenance. Consequently, the expenses associated
710 with travel and preparation for this maintenance type are included within the broader maintenance costs,
711 rather than being treated separately.

712 Additionally, a degraded track condition can result in increased delay, possession, and catastrophe costs,
713 while simultaneously decreasing rail revenues. Therefore, the reward function incorporates these costs and
714 revenues as a function of the track's condition, which are represented as condition-based rewards, r_c , as
715 illustrated in Figure 5. It is widely recognized that as the railway condition deteriorates, precautionary
716 speed restrictions should be imposed to minimize the risk of failures. However, in severe cases of track
717 degradation, speed restrictions alone may not be sufficient to mitigate the risks, which could potentially
718 lead to catastrophic outcomes. Furthermore, deteriorated tracks may require urgent maintenance during
719 operational hours, leading to increased possession costs. On the other hand, poor rail conditions directly
720 impact revenue generation. Deteriorated rail infrastructure reduces operational efficiency, leading to slower
721 trains, longer travel times, and unreliable services. These obstacles discourage potential customers and erode
722 the trust of existing passengers, resulting in reduced ridership and decreased revenue. It should be noted
723 that excluding revenues from the reward function even if the revenues are not affected by the condition may
724 lead the agent to make decisions to prematurely renew sections and terminate the episode. Thus, including
725 revenues is crucial to motivate the RL to continue the episode in a logical manner.

726 Condition-based rewards capture the difference between revenues and costs at different track conditions,
727 as depicted in Figure 5. Positive r_c are associated with favourable rail conditions, where revenues exceed
728 costs. However, as costs and losses surpass revenues, as observed in the "Poor" and "Super-red" conditions,
729 r_c become negative. These negative rewards serve as an indicator of the adverse impact that these condi-
730 tions could have on overall profitability. Importantly, these rewards are calculated per Poskey and usage,

Action	Type	Effect	Decision
Ballast maintenance	tamping	-1000/section	Chosen by RL agent Affected by the track condition, which is related to the RL decisions
	stoneblowing	-2000/section	
Ballast preparation and travel	tamping	-1000	
	stoneblowing	-2000	
Ballast renewal	-	Start a new episode	
Rail maintenance (includes preparation and travel)	rerail	-1500/section	
	welding	-300/section	
	grinding	-100/section	
Delay	-	Form r_c , which is function of the condition (Figure 5)	
Possession	-		
Catastrophe occurrence	-		
Revenues	-		

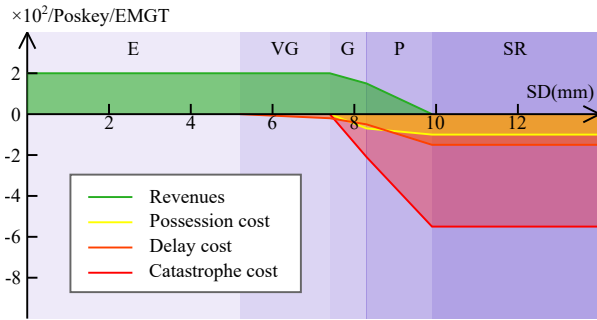
Table 6: Breakdown of input costs and revenues (unitless) influencing rewards (cost function).

731 accumulating continuously based on the usage. Conversely, other rewards are constant and calculated only
732 once per occurrence of the action. In summary, the reward function comprises constant rewards specific
733 to each action and continuous condition-based rewards (r_c) that are determined by the track's condition,
734 calculated per usage and Poskey, which are all summarized in Table 6.

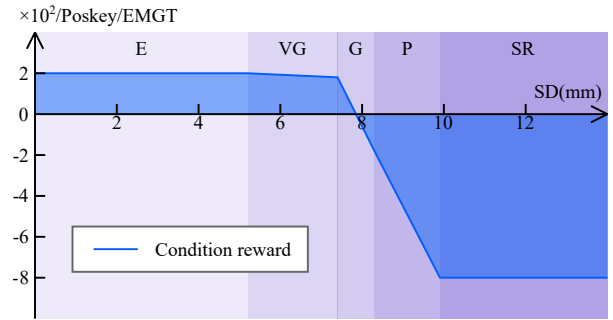
735 5.3.2. Definition of the environment

736 The environment of each section is defined in terms of the important features that the RL agent needs
737 for taking decisions. Thus, the factors describing the environment are chosen to be:

- 738 • the condition of the section, which depends on the value of the SD [mm] and is classified by the intervals
739 $[0,4]$, $(4,5.2]$, $(5.2,6.5]$, $(6.5,7.4]$, $(7.4,8.3]$, $(8.3,9.9]$, and $(9.9,\infty)$. These intervals are named E_1 , E_2 ,
740 VG_1 , VG_2 , G , P , and SR respectively.
- 741 • the settlement rate, which is represented by the rate of change in the SD [mm/EMGT]. It is divided
742 into two groups depending on whether the degradation rate is less than or greater than 0.2. These
743 groups are denoted as slow and fast and represented by letters S and F respectively.
- 744 • the maintenance history, which is represented by the last maintenance type and the number of times
745 this maintenance was performed previously. The maintenance history is named by a letter and a
746 number with the letter representing the types and the number representing the number of previous
747 actions. The letters T , SB , and R stand for the tamping, the stoneblowing, and the renewal actions



(a) The revenues and costs that are a function of the condition, which form up the condition-based reward, r_c .



(b) Condition-based reward, r_c , which is the summation of revenues and costs that are a function of the condition.

Condition	SD (mm)	Revenues	Costs			Condition reward (r_c)
			Possession	Delay	Catastrophe	
Excellent (E)	0	200	0	0	0	200
Very Good (VG)	5.2	200	0	0	0	200
Good (G)	7.4	200	0	20	0	180
Poor (P)	8.3	150	70	50	210	-180
Super Red (SR)	9.9	0	100	150	550	-800

(c) Summary of revenues, costs, and rewards that are a function of condition at different condition thresholds.

Figure 5: Visualization of revenues and costs influenced by rail conditions, depicted as condition-based rewards, r_c . The accompanying table presents corresponding values at each condition threshold in terms of SD , specifically for a track speed range of 5-20 MPH on a typical European track.

Algorithm 2 Calculation of the continuous rewards existing between two states of a section.

- 1: **Inputs:** Degradation rate DR , initial usage U , final usage U^* , initial SD , final SD^* .
 - 2: Get, C_T , the set of condition thresholds that are crossed between SD and SD^* : $\{c_{T1}, c_{T2}, \dots, c_{Tn}\}$.
 - 3: Get the usage of the section when crossing each of the conditions thresholds values by rearranging Equation 5: $U_i = (1/DR)c_{Ti} + (U - SD/DR)\forall i \in 1, \dots, n$.
 - 4: Divide the degradation between U and U^* into $n + 1$ intervals: $[U, U_1], [U_1, U_2], \dots, [U_n, U^*]$.
 - 5: Initialize the continuous reward, $R_c = 0$
 - 6: **for** $i \in \{1, \dots, n + 1\}$ **do** ▷ for all the intervals.
 - 7: Calculate the average standard deviation \overline{SD} to do the calculations based on it.
 - 8: Check if a fault will occur using function 1 of Algorithm 1.
 - 9: Check the probability of having a fault using function 1 in Algorithm 1.
 - 10: Sample the correction method of the fault using function 2 in Algorithm 1.
 - 11: calculate the usage, ΔU , during this interval.
 - 12: Calculate the cost for correcting the fault, c_f .
 - 13: Get the condition reward, r_c , using Figure 5.
 - 14: Accumulate the continuous reward: $R_c \rightarrow R_c + (r_c + c_f) \cdot L \cdot \Delta U$.
-

748 respectively. For example, T_2 stands for 2 previous tamping actions.

- 749 • the PN state, which is represented by the markings of places $p_{1,i}$ and $p_{7,i}$ for each section. This details
750 which type of decision is required.

751 5.3.3. Algorithm tuning and parameter scaling

752 The Q-learning method uses the bootstrapping effect, which means that Q-Values are updated based
753 on the values of the successor states. At the beginning of the learning process, all the Q-Values start with
754 random numbers, which makes these updates far from their actual values. This can give an advantage to the
755 Monte-Carlo RL method over the Q-learning at the beginning of the learning process. However, Q-learning
756 is an off-policy method that outperforms the Monte-Carlo or other on-policy methods by being able to
757 explore the environment without affecting the Q-Values updates [26]. For this, the learning process was
758 divided into two parts, with the Monte-Carlo RL method being used in the first one, and the Q-learning in
759 the second one.

760 The problem of optimising the policy for each section is considered an episodic task. However, the
761 number of episodes cannot be used to control the duration of the learning process because each environment
762 for each section has its own episode counter. A common variable that is shared among all environments is
763 the time. This variable has an effect on the number of episodes in each environment without affecting their
764 conditions, so it can be used as an arbitrary variable to control the duration of the learning process. For

765 this, the duration of the learning process is chosen to be equal to 6×10^7 yrs., with 2×10^7 yrs. using the
766 Monte-Carlo RL method and 4×10^7 yrs. using the Q-learning method. The model and its formulas are
767 affected by the time difference and not by the time, so the accumulation of time throughout the learning
768 process is just a way to increase the number of episodes, but it does not have any physical meaning or effect.

769 Three parameters, which are the discount rate γ , the learning rate α , and the exploration rate ε , should
770 be specified for the RL methods. The discount rate is assigned a value $\gamma = 1$ to avoid being biased to early
771 returns. The learning rate, α , is controlled by the following formula:

$$\alpha = \begin{cases} 1/n_u & \text{if } n_u < 1000 \\ 10^{-3} & \text{if } n_u > 1000 \end{cases} \quad (7)$$

772 where n_u is different for each state-action pair and it represents the number of times its Q-Value is updated.
773 For n_u less than 1000, the formula ensures that the Q-Value is equal to the average of all the previous
774 expected returns, and neglects the effect of initial values of the Q-Values [25].

$$\sigma_\varepsilon(t) = a + b \exp(-c \cdot t), \text{ with:} \quad (8)$$

$$a = \varepsilon_{\min}$$

$$b = \varepsilon_{\max} - \varepsilon_{\min}$$

$$c = \ln[(\varepsilon'_{\min} - a)/b]/t_e$$

$$\varepsilon'_{\min} = v(\varepsilon_{\min} - \varepsilon_{\max}) + \varepsilon_{\max}$$

775 Equation 8, which is based on Equation 12 in [25], presents an exponential decay function with easily
776 adjustable parameters that is utilized to regulate the decay of ε . The parameters of this function are the
777 end of the decay process, t_e , the maximum, ε_{\max} , and minimum, ε_{\min} , values of the controlled variable, and
778 a parameter called v . $v = 1 - \epsilon$ indicates how close the practical minimum is to the actual minimum. The
779 argument of the function is t , and for a range of $t \in [0, t_e]$ the output of the function decays from ε_{\max} to
780 ε_{\min} .

781 Exponential decay parameters are assigned for each part of the learning process. For the Monte-Carlo
782 RL method part, ε_{\max} , ε_{\min} , t_e , and v are chosen as 1, 10^{-4} , $0.95 \times 2 \times 10^7$, and 0.99 respectively, with the
783 argument, t being equal to the time ($t = \text{time}$). On the other hand, for the Q-learning method part, ε_{\max} ,
784 ε_{\min} , t_e , and v are chosen as 0.2, 10^{-3} , $0.9 \times 4 \times 10^7$, and 0.99 respectively, with the argument, t passed as,
785 $t = \text{time} - 2 \times 10^7$ to shift the argument to the starting point of the Q-learning part.

786 Since all the sections share the same characteristics, it is expected that they have the same optimal
787 policy. To make the learning process faster, the RL agents were allowed to share the experience by updating
788 the Q-Values in all environments once a Q-Value of similar state-action pair is updated in any environment.

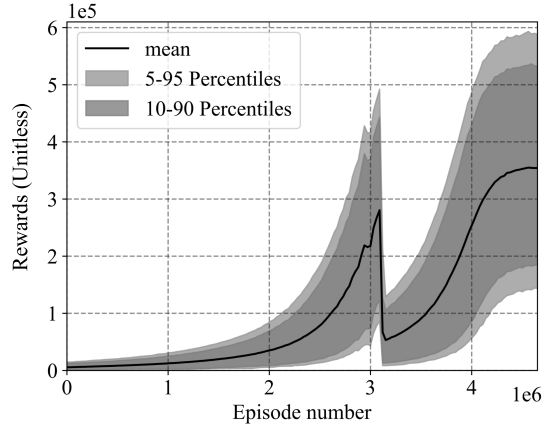


Figure 6: The variation of rewards as a function of the episode number.

789 6. Results

790 This section shows the results of the different simulations performed to find the best policy for the
791 operation and maintenance of the railway sections. Other simulations were performed to test the idea of
792 dividing the PN into multiple subnets. To do so, the PN proposed in Section 5.2 is simulated two times, one
793 while dividing the PN into multiple subnets, and another without dividing it. Each of the two simulations
794 was performed for multiple lifetimes such that the total duration of the accumulated lifetimes is 5000 yrs.
795 in each simulation. Besides, no optimisation was done in these simulations and a random policy was used
796 to select transitions from action groups. As a result, the time taken for the simulation with and without the
797 subnet rules was 201 and 613 seconds respectively. This indicates a 3 times reduction in the computational
798 cost between the two simulations.

799 Figure 6 shows the variation of the total reward as a function of the learning process. The learning process
800 is divided into short intervals to be able to plot the mean and other measures for each of the intervals. It
801 can be seen that the total reward increases until the end of the Monte-Carlo RL part of the learning process,
802 then drops and continues increasing after the Q-learning starts. The end of the learning process shows a
803 stable curve which indicates that the policy is no longer changing. The decisions of the final RL policy are
804 summarised in Figure 7. The left part of the Figure shows the decisions for normal maintenance while the
805 right part is for opportunistic maintenance. The 11 rows show the possible maintenance history states while
806 the columns describe the condition of the section and settlement rate. The different decisions are described
807 by the colours in the legend above the figure, where the white areas represent the unexplored states. For
808 example, the decision shown in the red square is tamping and it corresponds to the state described by the
809 section being in good condition with a slow settlement rate and 4 previous tamping actions. This figure
810 shows that a sequence of 7 tamping actions followed by 3 stoneblowing actions and a renewal action is good

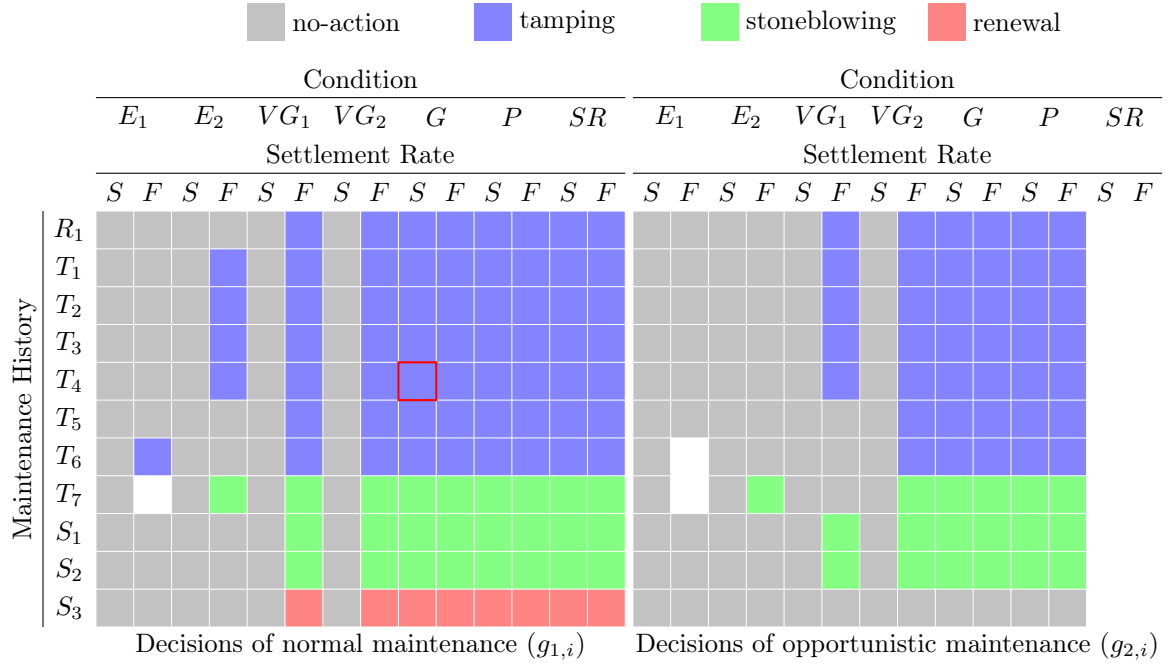


Figure 7: The final RL policy, which is described by the optimal actions at each of the RL states. The acronyms are defined in Section 5.3.2, white areas are unexplored states, and the red square is an indication of an example explained in the text.

811 for all track conditions since no stoneblowing action can be found above row T_7 and no stoneblowing action
812 can be found above row S_3 . This means that the RL policy does not differ from policies A and B described
813 in the introduction of Section 5 in terms of the type of maintenance action. However, the decision regarding
814 the need for maintenance is shown to be dependent on all the features of the RL environment because the
815 no-action decisions are scattered in different zones. It can be seen that the need for maintenance increases
816 as the condition of the section worsens especially in the cases of a fast settlement rate. Also, the distribution
817 of actions over states is similar for opportunistic and normal maintenance actions.

818 The decisions shown in Figure 7 are reflected in the total rewards, the percentage of time spent in
819 each condition, the distribution of maintenance actions over time and the section condition. Table 7 shows
820 the average percentage of time spent in each condition when following each of the policies. Comparing
821 percentages of time spent in each condition is preferable to using total absolute durations spent in each
822 of the conditions per episode in order to avoid unfair conclusions. This is because absolute durations
823 can be misleading due to episodes terminating with the ballast's life, rather than the rail's service life.
824 Rail operators should maintain the rail according to its service life, not the life of the ballast. By using
825 percentages, conditions can be reflected relative to the duration of the rail's service life, and can indicate
826 probabilities of each state. It can be noted that Policy B increases the probability of being in the Poor
827 and Super-red conditions while the RL policy was able to reduce the probability of being in the Super-red

Table 7: The average percentage of time spent in each condition for each of the simulated policies.

	Excellent	Very good	Good	Poor	Super-red
Policy A	87.33	11.75	0.53	0.33	0.04
Policy B	57.19	35.43	5.06	1.59	0.70
RL policy	66.17	30.97	2.40	0.34	0.007

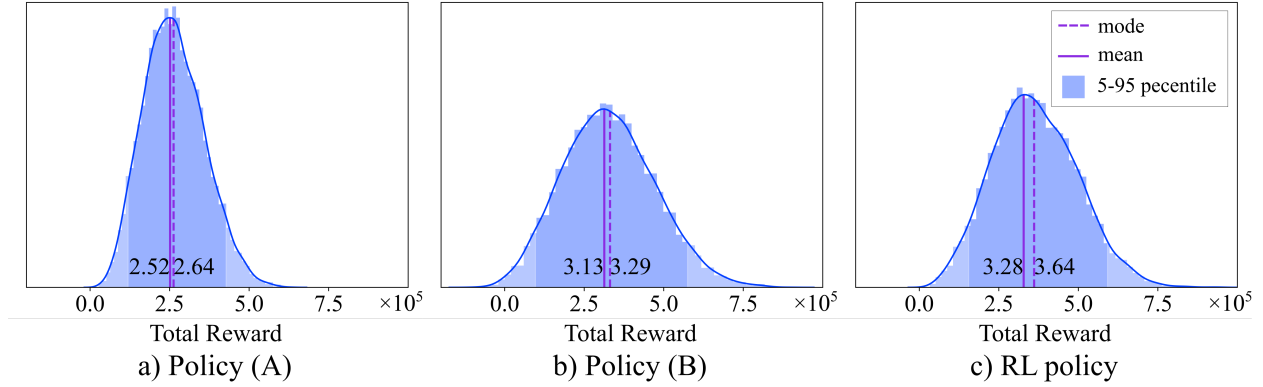


Figure 8: The distributions of rewards for the three considered policies.

828 condition to almost 0 while having a very low probability of being in the Poor condition. Policy A showed
 829 the ability to have the highest percentage of the Excellent condition but without being able to have the
 830 lowest percentage in the Super-red condition.

831 Figure 8 shows the distribution of rewards of the considered policies. The distributions are plotted
 832 because the problem is stochastic, and comparing based on one value can be misleading in such cases. The
 833 figure shows that Policy A resulted in the minimum rewards of all the policies, but it ensures that the
 834 rewards of an episode are always greater than 0, whereas Policy B results in rewards greater than Policy A,
 835 but can result in negative rewards. On the other hand, the RL policy ensures positive rewards whilst also
 836 ensuring the maximum rewards of all policies in terms of the mean and mode.

837 Figure 9 shows the distribution of maintenance actions over the condition of the section. It can be
 838 seen that the maintenance actions are concentrated in the Very-Good condition for Policy A, in the Good
 839 condition for Policy B, and in different conditions for the RL Policy. The figure also shows that there are very
 840 few actions taken in the Super-Red condition for Policy A and RL policy, whereas Policy B has a significant
 841 number of actions taken in this condition. Figure 10 shows the distribution of maintenance actions over the
 842 age of the section. For the three policies, tamping is followed by stoneblowing then by renewal. The renewal
 843 action is an indication of the end of life of the section. Policy A resulted in the shortest life with an average
 844 equal to 29.5 yrs, followed by the RL Policy with an average of 42.5 yrs., then by Policy B with an average
 845 equal to 45.2 yrs.

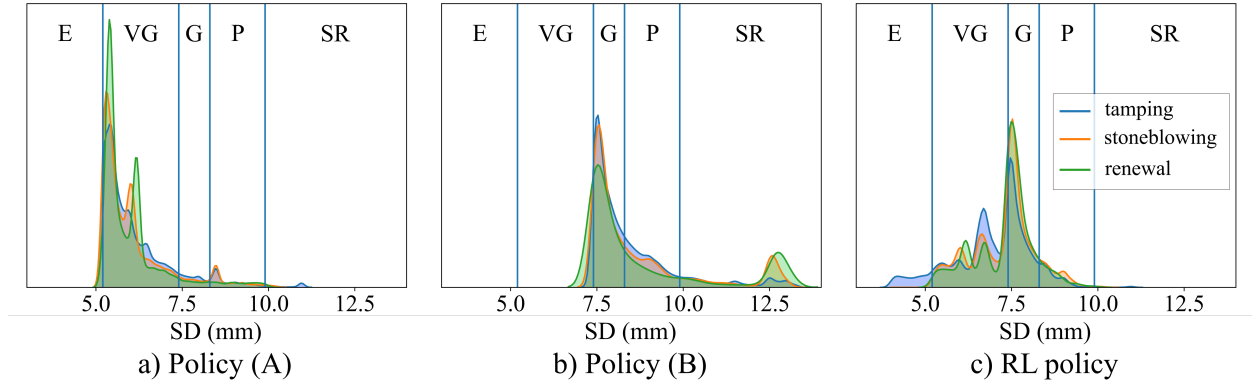


Figure 9: The distributions of the maintenance actions over the condition of the section for the three considered policies.

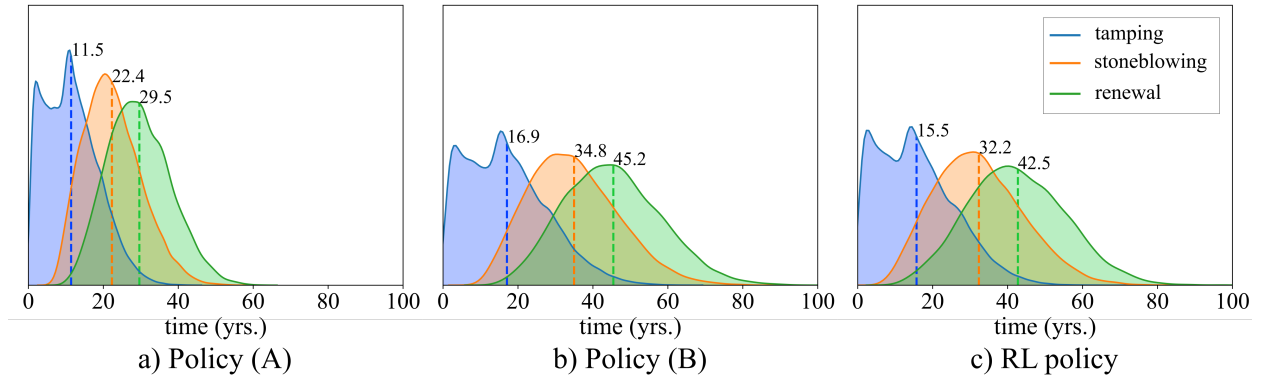


Figure 10: The distributions of the maintenance actions over the age of the section for the three considered policies.

846 7. Discussion

847 Regarding the division of the PN to multiple subnets, comparing the random simulations for the PN
 848 case, which is presented in Section 5.2, revealed a three-fold decrease in computational costs despite some
 849 parts of the PN being computed in parallel. The subnet rules can result in greater reduction if more subnets
 850 were in series because, for parallel subnets, all of the transitions are important to be checked at the same
 851 state, which weakens the effect of unchecking unimportant transitions because they will be few.

852 Regarding the optimisation problem, the total rewards shown in Figure 6 indicate that the optimal RL
 853 policy was reached since stable results are seen by the end of the learning process. The drop in the total
 854 rewards shown when the Q-learning started is due to the change in the exploration rate. Since the Q-learning
 855 is an off-policy RL method, increasing the exploration rate does not cause any diversion to the Q-Values and
 856 keeps their updates correct [26], but it allows the RL agent to discover other decisions that may be better
 857 than those already explored.

858 The results displayed in Figure 7 depict the decisions made by the RL agent for condition-based and
 859 opportunistic maintenance actions, which together constitute the final policy. Notably, the zones where the

860 agent chose to perform opportunistic maintenance are a subset of those where condition-based maintenance is
861 performed. This implies that the costs saved by avoiding the preparation and travel involved in opportunistic
862 maintenance did not influence the agent’s decision. The optimal policy, therefore, advocates for repairing
863 a section only if it requires maintenance based on its condition, and not on the availability of equipment.
864 Following this policy means that opportunistic maintenance is only triggered if there is a change in the
865 section’s condition between the inspection and the arrival at the site. In such cases, the decision to repair
866 will be made on-site, which makes it opportunistic, but the decision will be due to the change in condition
867 rather than the availability of equipment on site. Pursuing this policy ensures that maintenance actions are
868 conducted only when required, given that each action has an impact on the ballast’s life by crushing its
869 particles and making it more susceptible to fouling, thereby shortening its lifespan. In addition, the policy
870 tries to delay ballast maintenance as much as possible to extend its usage period until it reaches a highly
871 fouled state where it can no longer be utilized. By postponing ballast maintenance, the overall maintenance
872 costs can be reduced, as the ballast can be used for a longer duration before replacement becomes necessary.

873 Table 7 shows that the final RL policy was the best in avoiding the Super-red condition but without
874 having the maximum percentage in the Excellent condition. Since the goal is to increase the net profit of
875 each section, it is not important in which condition the section stays the most, but rather the effect of being
876 in each of the conditions. As shown in Figure 5, the Super-red condition has very high negative rewards that
877 make avoiding it more important than being in Excellent condition. In addition to that, keeping the section
878 in the best condition all the time may result in performing additional maintenance actions that require
879 more costs and may result in performing maintenance before it is needed. This can result in reducing the
880 remaining useful life as explained in the previous paragraph, which shortens the age of the section as shown
881 in Figure 10a, and decreases the total rewards gained per episode as shown in Figure 8a. Besides, postponing
882 the maintenance can increase the age of the section as shown in 10b, but it can lead to having negative
883 consequences due to being in undesired conditions as shown in Figure 8b. On the other hand, the RL policy
884 was able to increase the rewards to a mean and mode better than the other two policies without having
885 any episodes with negative rewards by taking the maintenance decision just before reaching any negative
886 consequences based on different features. This required having the actions distributed over the condition of
887 the section as shown in Figure 9c in contrast to Policies A and B which show a concentration of actions in
888 specific zones. This shows that the decision to perform maintenance actions is not only a function of the
889 section condition but also of the settlement rate and the maintenance history.

890 The RL policy outperforms both policies A and B in terms of rewards and in terms of avoiding the
891 Super-Red condition. It can be concluded from the results that the effect of the maintenance action on the
892 remaining useful life was more important than their costs. The RL agent was trying to avoid the maintenance
893 action until the condition becomes unacceptable in order to use the section for producing revenues as much as
894 possible before the section moves to the next stage, which is the after-maintenance stage. Each maintenance

895 can be seen as a new beginning that has its own revenues before the losses start, so performing another
896 maintenance before all the revenues after the first maintenance are harvested was like losing them. At the
897 same time, the RL agent was successfully able to avoid the risk of being in the Super-Red condition as
898 can be seen from Figure 9 and Table 7. This was due to the intelligent strategy shown in Figure 7, which
899 considered the condition and the settlement rate. The figure shows that even if the condition was Excellent
900 but the settlement is fast, the decision was to perform the maintenance for some occasions. This decision
901 may result in a great reduction in the age of the section, but it also results in avoiding any risk of being in
902 the Super-Red condition.

903 The outcomes presented in this paper are highly influenced by the reward functions that have been
904 assigned. These functions, which are based on estimations and expert opinions, describe the costs and
905 revenues involved, allowing for the calculation of the net profit. However, the accuracy of the simulations
906 and results could be improved if more precise information regarding the costs and revenues associated with
907 railway transport were available. This could potentially result in changes to the findings.

908 Nevertheless, the paper provides a reliable method for optimizing operation and maintenance based on
909 the currently available data. However, to enhance the paper's findings further, it would be beneficial to
910 incorporate the impact of maintenance activities of all railway components, such as infrastructure, super-
911 structure, signalling, and catenary. By including this level of detail, the model's accuracy and precision can
912 be significantly improved. Such a comprehensive analysis could be addressed in future studies.

913 An improvement can be made to the methodology, which is using *function approximation* RL methods,
914 e.g. Deep Reinforcement learning. This can help in avoiding the discretisation process of the states and
915 considering continuous states instead. For example, the condition of the section can still be a feature that
916 describes the condition of the environment, but it will take the SD as a continuous variable argument instead
917 of dividing the condition into several groups based on SD . This wipes out the need to use thresholds, which
918 can result in further improvements in terms of taking the decisions at more specific states instead of having
919 the same decisions for wide intervals of values.

920 In addition, the use of prognostic methods can be important in predicting the remaining useful life of
921 the section based on its current state. This can be included as a feature in the RL environments to improve
922 decision-making and can be more realistic than including the settlement rate which is difficult to measure
923 in real-life applications.

924 8. Conclusions

925 An *iPN* model was created for the maintenance and operation of railway sections while focusing on
926 optimising the maintenance of the ballast. This model is able to find the optimal maintenance strategy
927 that can reduce the risk of being in undesired conditions while increasing revenues and decreasing costs.

928 This paper also proposes several ideas to improve the computational efficiency of the model. A method to
929 divide the PN into several subnets was proposed and found to be successful in reducing computational costs.
930 Besides, each section of the railway was considered a separate environment that has its own RL elements.
931 This allows the RL agents to focus only on the important aspects when taking the decisions of each section
932 and neglecting unnecessary information, which reduces the number of RL states. This, in turn, facilitates
933 experience sharing between RL agents relating to sections of similar characteristics.

934 The model was applied to a practical problem and it shows the ability to reach an optimum maintenance
935 strategy. The results show that it is crucial to avoid unnecessary maintenance actions because they can
936 reduce the ballast age. This is because tamping and stoneblowing actions play a direct role in ballast
937 fouling, which requires replacement once it becomes highly fouled. At the same time, the maintenance
938 should be done before any risk of reaching a bad condition in order to avoid downtime or safety risks. A
939 maintenance plan that gives the optimum decision as a function of various features of the railway section was
940 found. This was able to avoid undesired conditions while increasing the age of each section and increasing
941 the net profits per life of each section.

942 **Acknowledgements**

943 This paper is part of the ENHAnCE ITN project (<https://www.h2020-enhanceitn.eu/>) funded by
944 the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie
945 grant agreement No. 859957. The authors would also like to thank the funding provided by the National
946 Research and Science Council of Spain under grant agreements No. RTI2018-101841-B-C21 and PLEC2021-
947 007798. The first author would like to acknowledge the University of Nottingham (UK) for hosting while
948 conducting this study.

949 References

- 950 [1] Judiciary of England and Wales, Office of rail regulation v network rail infrastructure limited., Technical Report, Ju-
951 diciary of England and Wales, 2012. URL: [https://www.judiciary.uk/wp-content/uploads/JCO/Documents/Judgments/
952 office-rail-regulation-network-rail-sentencing-remarks-04042012.pdf](https://www.judiciary.uk/wp-content/uploads/JCO/Documents/Judgments/office-rail-regulation-network-rail-sentencing-remarks-04042012.pdf).
- 953 [2] E. T. Selig, J. M. Waters, Track geotechnology and substructure management, Thomas Telford, 1994.
- 954 [3] J. Sadeghi, M. Motieyan-Najar, J. Zakeri, B. Yousefi, M. Mollazadeh, Improvement of railway ballast maintenance
955 approach, incorporating ballast geometry and fouling conditions, *Journal of Applied Geophysics* 151 (2018) 263–273.
- 956 [4] Y. Guo, V. Markine, G. Jing, Review of ballast track tamping: Mechanism, challenges and solutions, *Construction and
957 Building Materials* 300 (2021) 123940.
- 958 [5] M. Daddow, X. Zhang, H. Qiu, Z. Zhang, Y. Liu, A mathematical model for ballast tamping decision making in railway
959 tracks, *Civil Engineering Journal* 6 (2020) 2045–2057.
- 960 [6] M. Sedghi, O. Kauppila, B. Bergquist, E. Vanhatalo, M. Kulahci, A taxonomy of railway track maintenance planning and
961 scheduling: A review and research trends, *Reliability Engineering & System Safety* 215 (2021) 107827.
- 962 [7] J. Chiachío, M. Chiachío, D. Prescott, J. Andrews, A knowledge-based prognostics framework for railway track geometry
963 degradation, *Reliability Engineering & System Safety* 181 (2019) 127–141.
- 964 [8] A. Falamarzi, S. Moridpour, M. Nazem, A review of rail track degradation prediction models, *Australian Journal of Civil
965 Engineering* 17 (2019) 152–166.
- 966 [9] L. F. Caetano, P. F. Teixeira, Optimisation model to schedule railway track renewal operations: a life-cycle cost approach,
967 *Structure and Infrastructure Engineering* 11 (2015) 1524–1536.
- 968 [10] C. Letot, I. Soleimanmeigouni, A. Ahmadi, P. Dehombreux, An adaptive opportunistic maintenance model based on
969 railway track condition prediction, *IFAC-PapersOnLine* 49 (2016) 120–125.
- 970 [11] L. F. Caetano, P. F. Teixeira, Availability approach to optimizing railway track renewal operations, *Journal of Trans-
971 portation Engineering* 139 (2013) 941–948.
- 972 [12] Z. Su, A. Jamshidi, A. Núñez, S. Baldi, B. De Schutter, Integrated condition-based track maintenance planning and crew
973 scheduling of railway networks, *Transportation Research Part C: Emerging Technologies* 105 (2019) 359–384.
- 974 [13] J. Andrews, D. Prescott, F. De Rozières, A stochastic model for railway track asset management, *Reliability Engineering
975 & System Safety* 130 (2014) 76–84.
- 976 [14] D. Prescott, J. Andrews, Investigating railway track asset management using a markov analysis, *Proceedings of the
977 Institution of Mechanical Engineers, Part F: Journal of Rail and Rapid Transit* 229 (2015) 402–416.
- 978 [15] S. Sharma, Y. Cui, Q. He, R. Mohammadi, Z. Li, Data-driven optimization of railway maintenance for track geometry,
979 *Transportation Research Part C: Emerging Technologies* 90 (2018) 34–58.
- 980 [16] S. Bressi, J. Santos, M. Losa, Optimization of maintenance strategies for railway track-bed considering probabilistic
981 degradation models and different reliability levels, *Reliability engineering & system safety* 207 (2021) 107359.
- 982 [17] J. D. Campbell, J. V. Reyes-Picknell, H. S. Kim, *Uptime: Strategies for excellence in maintenance management*, CRC
983 Press, 2015.
- 984 [18] M. Dell’Orco, M. Ottomanelli, L. Caggiani, D. Sassanelli, New decision support system for optimization of rail track
985 maintenance planning based on adaptive neurofuzzy inference system, *Transportation research record* 2043 (2008) 49–54.
- 986 [19] A. Saleh, D. Prescott, R. Remenyte, M. Chiachio, An optimized asset management petri net model for railway sections
987 (2023).
- 988 [20] M. Chiachío, A. Saleh, S. Naybour, J. Chiachío, J. Andrews, Reduction of petri net maintenance modeling complexity
989 via approximate bayesian computation, *Reliability Engineering & System Safety* 222 (2022) 108365.
- 990 [21] M. Chiachío, J. Chiachío, D. Prescott, J. Andrews, A new paradigm for uncertain knowledge representation by plausible
991 Petri nets, *Information Sciences* 453 (2018) 323–345.

- 992 [22] J. Lee, K. Liu, W. Chiang, Modeling uncertainty reasoning with possibilistic Petri nets, *IEEE Transactions on Systems,*
993 *Man, and Cybernetics, Part B: Cybernetics* 33 (2003) 214–224.
- 994 [23] B. Pei-Ming, Learning capability in fuzzy Petri nets based on bp net [j], *Chinese Journal of Computers* 5 (2004).
- 995 [24] M. M. Hanna, A. Buck, R. Smith, Fuzzy Petri nets with neural networks to model products quality from a cnc-milling
996 machining centre, *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans* 26 (1996) 638–645.
- 997 [25] A. Saleh, M. Chiachío, J. F. Salas, A. Kolios, Self-adaptive optimized maintenance of offshore wind turbines by intelligent
998 petri nets, *Reliability Engineering & System Safety* (2022) 109013.
- 999 [26] R. S. Sutton, A. G. Barto, *Reinforcement learning: An introduction*, MIT press, 2018.
- 1000 [27] L. Busoniu, R. Babuska, B. D. Schutter, Multi-agent reinforcement learning: An overview, *Innovations in multi-agent*
1001 *systems and applications-1* (2010) 183–221.
- 1002 [28] A. Senderovich, A. Shleyfman, M. Weidlich, A. Gal, A. Mandelbaum, To aggregate or to eliminate? optimal model
1003 simplification for improved process performance prediction, *Information Systems* 78 (2018) 96–111. URL: <https://www.sciencedirect.com/science/article/pii/S0306437916306329>. doi:<https://doi.org/10.1016/j.is.2018.04.003>.
- 1004 [29] A. Senderovich, A. Rogge-Solti, A. Gal, J. Mendling, A. Mandelbaum, S. Kadish, C. Bunnell, Data-
1005 driven performance analysis of scheduled processes, *Intl. Conf. on Business Process Management*, Springer
1006 (2016) 35 â 52. URL: [https://www.scopus.com/inward/record.uri?eid=2-s2.0-85123627763&partnerID=40&md5=](https://www.scopus.com/inward/record.uri?eid=2-s2.0-85123627763&partnerID=40&md5=83e184e57570ad913561efc6fef53a85)
1007 [83e184e57570ad913561efc6fef53a85](https://www.scopus.com/inward/record.uri?eid=2-s2.0-85123627763&partnerID=40&md5=83e184e57570ad913561efc6fef53a85), cited by: 1.
- 1008 [30] F. M. BÃnneland, J. Dyhr, P. G. Jensen, M. Johannsen, J. Srba, Stubborn versus structural reductions for petri nets,
1009 *Journal of Logical and Algebraic Methods in Programming* 102 (2019) 46–63. URL: [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S235222081830035X)
1010 [science/article/pii/S235222081830035X](https://www.sciencedirect.com/science/article/pii/S235222081830035X). doi:<https://doi.org/10.1016/j.jlamp.2018.09.002>.
- 1011 [31] V. Padma, A study on flexible manufacturing systems using petri net reduction, in: *IOP Conference Series: Materials*
1012 *Science and Engineering*, volume 1130, IOP Publishing, 2021, p. 012060.
- 1013 [32] B. Berthomieu, D. Le Botlan, S. Dal Zilio, Counting petri net markings from reduction equations, *International Journal*
1014 *on Software Tools for Technology Transfer* 22 (2020) 163–181.
- 1015 [33] E. I. Basri, I. H. A. Razak, H. Ab-Samat, S. Kamaruddin, Preventive maintenance (pm) planning: a review, *Journal of*
1016 *Quality in Maintenance Engineering* 23 (2017) 114–143.
- 1017 [34] M. Shafiee, J. D. Sørensen, Maintenance optimization and inspection planning of wind energy assets: Models, methods
1018 and strategies, *Reliability Engineering & System Safety* 192 (2019) 105993.
- 1019 [35] R. D. Palmer, *Maintenance planning and scheduling handbook*, McGraw-Hill Education, 2013.
- 1020 [36] S. Duffuaa, K. Al-Sultan, Mathematical programming approaches for the management of maintenance planning and
1021 scheduling, *Journal of Quality in Maintenance Engineering* 3 (1997) 163–176.
- 1022 [37] J. S. Lee, I. Y. Choi, I. K. Kim, S. H. Hwang, Tamping and renewal optimization of ballasted track using track measurement
1023 data and genetic algorithm, *Journal of Transportation Engineering, Part A: Systems* 144 (2018) 04017081.
- 1024 [38] M. Movaghar, S. Mohammadzadeh, Intelligent index for railway track quality evaluation based on bayesian approaches,
1025 *Structure and Infrastructure Engineering* 16 (2020) 968–986.
- 1026 [39] C. Meier-Hirmer, A. Senee, G. Riboulet, F. Sourget, M. Roussignol, et al., A decision support system for track maintenance,
1027 *Computers in railways X* (2006) 217–226.
- 1028 [40] C. Vale, I. M. Ribeiro, R. Calçada, Integer programming to optimize tamping in railway tracks as preventive maintenance,
1029 *Journal of Transportation Engineering* 138 (2012) 123–131.
- 1030 [41] E. Gustavsson, Scheduling tamping operations on railway tracks using mixed integer linear programming, *EURO Journal*
1031 *on Transportation and Logistics* 4 (2015) 97–112.
- 1032 [42] H. Khajehi, A. Ahmadi, I. Soleimanmeigouni, A. Nissen, Allocation of effective maintenance limit for railway track
1033 geometry, *Structure and Infrastructure Engineering* 15 (2019) 1597–1612.
- 1034

- 1035 [43] T. Satiennam, T. Phanyakit, Fuzzy multi-attribute decision making for the selection of a suitable railway track maintenance
1036 plan: A case study in thailand, *GEOMATE Journal* 17 (2019) 96–104.
- 1037 [44] M. S. Kovačević, M. Bačić, I. Stipanović, K. Gavin, Categorization of the condition of railway embankments using a
1038 multi-attribute utility theory, *Applied Sciences* 9 (2019) 5089.
- 1039 [45] L. Podofillini, E. Zio, J. Vatn, Risk-informed optimisation of railway tracks inspection and maintenance procedures,
1040 *Reliability Engineering & System Safety* 91 (2006) 20–35.
- 1041 [46] M. H. Bin Osman, S. Kaewunruen, A. Jack, Optimisation of schedules for the inspection of railway tracks, *Proceedings
1042 of the Institution of Mechanical Engineers, Part F: Journal of Rail and Rapid Transit* 232 (2018) 1577–1587.
- 1043 [47] T. Lidén, M. Joborn, An optimization model for integrated planning of railway traffic and network maintenance, *Trans-
1044 portation Research Part C: Emerging Technologies* 74 (2017) 327–347.
- 1045 [48] A. R. Albrecht, D. M. Pantan, D. H. Lee, Rescheduling rail networks with maintenance disruptions using problem space
1046 search, *Computers & Operations Research* 40 (2013) 703–712.
- 1047 [49] D. Arenas, P. Pellegrini, S. Hanafi, J. Rodriguez, Timetable rearrangement to cope with railway maintenance activities,
1048 *Computers & Operations Research* 95 (2018) 123–138.
- 1049 [50] H. Khajehi, M. Haddadzade, A. Ahmadi, I. Soleimanmeigouni, A. Nissen, Optimal opportunistic tamping scheduling for
1050 railway track geometry, *Structure and Infrastructure Engineering* 17 (2021) 1299–1314.
- 1051 [51] R. Santos, P. F. Teixeira, Heuristic analysis of the effective range of a track tamping machine, *Journal of infrastructure
1052 systems* 18 (2012) 314–322.
- 1053 [52] C. Esveld, C. Esveld, *Modern railway track*, volume 385, MRT-productions Zaltbommel, 2001.
- 1054 [53] M. Sasidharan, M. Burrow, G. Ghataora, A whole life cycle approach under uncertainty for economically justifiable
1055 ballasted railway track maintenance, *Research in Transportation Economics* 80 (2020) 100815.
- 1056 [54] M. Daddow, X. Zhang, H. Qiu, Z. Zhang, Impact of unused life for track sections and available workforce in scheduling
1057 tamping actions on ballasted tracks, *KSCE Journal of Civil Engineering* 21 (2017) 2403–2412.
- 1058 [55] C. Dao, R. Basten, A. Hartmann, Maintenance scheduling for railway tracks under limited possession time, *Journal of
1059 Transportation Engineering, Part A: Systems* 144 (2018) 04018039.
- 1060 [56] T. Lidén, T. Kalinowski, H. Waterer, Resource considerations for integrated planning of railway traffic and maintenance
1061 windows, *Journal of rail transport planning & management* 8 (2018) 1–15.
- 1062 [57] G. Budai, D. Huisman, R. Dekker, Scheduling preventive railway maintenance activities, *Journal of the Operational
1063 Research Society* 57 (2006) 1035–1044.
- 1064 [58] M. Forsgren, M. Aronsson, S. Gestrelus, Maintaining tracks and traffic flow at the same time, *Journal of Rail Transport
1065 Planning & Management* 3 (2013) 111–123.
- 1066 [59] X. Luan, J. Miao, L. Meng, F. Corman, G. Lodewijks, Integrated optimization on train scheduling and preventive
1067 maintenance time slots planning, *Transportation Research Part C: Emerging Technologies* 80 (2017) 329–359.
- 1068 [60] C. Zhang, Y. Gao, L. Yang, Z. Gao, J. Qi, Joint optimization of train scheduling and maintenance planning in a railway
1069 network: A heuristic algorithm using lagrangian relaxation, *Transportation Research Part B: Methodological* 134 (2020)
1070 64–92.
- 1071 [61] A. Consilvio, A. Di Febraro, R. Meo, N. Sacco, Risk-based optimal scheduling of maintenance activities in a railway
1072 network, *EURO journal on transportation and logistics* 8 (2019) 435–465.
- 1073 [62] A. Bakhtiary, J. A. Zakeri, S. Mohammadzadeh, An opportunistic preventive maintenance policy for tamping scheduling
1074 of railway tracks, *International Journal of Rail Transportation* 9 (2021) 1–22.
- 1075 [63] J. Zhao, A. Chan, M. Burrow, Reliability analysis and maintenance decision for railway sleepers using track condition
1076 information, *Journal of the Operational Research Society* 58 (2007) 1047–1055.
- 1077 [64] L. C. Sancho, J. A. Braga, A. R. Andrade, Optimizing maintenance decision in rails: A markov decision process approach,

- 1078 ASCE-ASME Journal of Risk and Uncertainty in Engineering Systems, Part A: Civil Engineering 7 (2021) 04020051.
- 1079 [65] R. Mohammadi, Q. He, A deep reinforcement learning approach for rail renewal and maintenance planning, Reliability
1080 Engineering & System Safety (2022) 108615.
- 1081 [66] T. Gao, Z. Li, Y. Gao, P. Schonfeld, X. Feng, Q. Wang, Q. He, A deep reinforcement learning approach to mountain
1082 railway alignment optimization, Computer-Aided Civil and Infrastructure Engineering 37 (2022) 73–92.
- 1083 [67] N. Yousefi, S. Tsianikas, D. W. Coit, Reinforcement learning for dynamic condition-based maintenance of a system with
1084 individually repairable components, Quality Engineering 32 (2020) 388–408.
- 1085 [68] J. Sresakoolchai, S. Kaewunruen, Railway infrastructure maintenance efficiency improvement using deep reinforcement
1086 learning integrated with digital twin based on track geometry and component defects, Scientific Reports 13 (2023) 2439.
- 1087 [69] C. J. Watkins, P. Dayan, Q-learning, Machine learning 8 (1992) 279–292.
- 1088 [70] T. Murata, Petri nets: Properties, analysis and applications, Proceedings of the IEEE 77 (1989) 541–580.
- 1089 [71] K. Jensen, G. Rozenberg, High-level Petri nets: theory and application, Springer Science & Business Media, 2012.
- 1090 [72] L. S. Shapley, Stochastic games, Proceedings of the national academy of sciences 39 (1953) 1095–1100.
- 1091 [73] M. L. Littman, Markov games as a framework for multi-agent reinforcement learning, in: Machine learning proceedings
1092 1994, Elsevier, 1994, pp. 157–163.
- 1093 [74] L. Busoniu, R. Babuska, B. De Schutter, A comprehensive survey of multiagent reinforcement learning, IEEE Transactions
1094 on Systems, Man, and Cybernetics, Part C (Applications and Reviews) 38 (2008) 156–172.
- 1095 [75] R. Lal, Encyclopedia of soil science, 11, CRC Press, 2006.
- 1096 [76] J. Litherland, J. Andrews, A petri net asset management framework for railway switches and crossings, Proceedings of
1097 the Institution of Mechanical Engineers, Part F: Journal of Rail and Rapid Transit (2022) 09544097221110970.
- 1098 [77] B. Aursudkij, A laboratory study of railway ballast behaviour under traffic loading and tamping maintenance, Ph.D.
1099 thesis, University of Nottingham Nottingham, UK, 2007.
- 1100 [78] K. Tzanakakis, Track Settlements, Springer Berlin Heidelberg, Berlin, Heidelberg, 2013, pp. 147–148. URL: https://doi.org/10.1007/978-3-642-36051-0_20. doi:10.1007/978-3-642-36051-0_20.
- 1101
- 1102 [79] S. Kaewunruen, M. Ishida, S. Marich, Dynamic wheel–rail interaction over rail squat defects, Acoustics Australia 43
1103 (2015) 97–107.
- 1104 [80] R. Smith, Railway fatigue failures: an overview of a long standing problem, Materialwissenschaft und Werkstofftechnik:
1105 Entwicklung, Fertigung, Prüfung, Eigenschaften und Anwendungen technischer Werkstoffe 36 (2005) 697–705.
- 1106 [81] A. Wilson, RAIL DEFECTS HANDBOOK, 2019.
- 1107 [82] S. Kumar, Study of rail breaks: associated risks and maintenance strategies, Luleå tekniska universitet, 2006.