# Mining-Gym: A Configurable RL Benchmarking Environment for Truck Dispatch Scheduling

Chayan Banerjee, *Member, IEEE*, Kien Nguyen, *Senior Member, IEEE* and Clinton Fookes, *Senior Member, IEEE*

arXiv:2503.19195v1 [cs.LG] 24 Mar 2025

*Abstract*—Mining process optimization, particularly truck dispatch scheduling, is a critical factor in enhancing the efficiency of open-pit mining operations. However, the dynamic and stochastic nature of mining environments—characterized by uncertainties such as equipment failures, truck maintenance, and variable haul cycle times—poses significant challenges for traditional optimization methods. While Reinforcement Learning (RL) has demonstrated promise in adaptive decision-making for mining logistics, its practical deployment requires rigorous evaluation in realistic and customizable simulation environments. The lack of such standardized environment benchmarks limits fair algorithm comparisons, reproducibility, and real-world applicability of RL-based approaches in open-pit mining settings. To address this challenge, we introduce Mining-Gym, a configurable, open-source benchmarking environment designed for training, testing, and comparing RL algorithms in mining process optimization. Built on Discrete Event Simulation (DES) and seamlessly integrated with the OpenAI Gym interface, Mining-Gym offers a structured testbed that enables the direct application of advanced RL algorithms from Stable Baselines. The framework models key mining-specific uncertainties, such as equipment failures, queue congestion, and stochasticity of mining processes, ensuring a realistic and adaptive learning environment. Additionally, a graphic user interface (GUI) for easy parameter selection for mine-site configuration, comprehensive data logging system, a built-in KPI dashboard and real-time representative visualization of mine-site enables in-depth performance analysis, facilitating standardized, reproducible evaluation across multiple RL strategies and baseline heuristics.

*Index Terms*—Reinforcement Learning, Truck Scheduling, Discrete Event Simulation, OpenAI Gym, Mining Simulation, Resource Allocation

## I. INTRODUCTION

**M**INING process optimization aims to enhance efficiency and productivity by improving resource allocation, equipment scheduling, and material handling. However, these operations are highly complex, influenced by dynamic factors such as equipment failures, fluctuating ore quality, and unpredictable environmental conditions. Traditional optimization methods, such as linear programming and heuristics, struggle to adapt in real time, leading to inefficiencies and increased costs.

RL offers a promising dynamic approach, but its adoption in mining remains limited due to the lack of standardized simulation environments that accurately model the stochastic nature of mining operations, provide a configurable testbed for algorithm development and comparison, and enable reproducible research with standardized scenarios and metrics. Additionally

C. Banerjee, K. Nguyen, C. Fookes are with the School of Electrical Engineering and Robotics, Brisbane, QLD, Australia (e-mail: {c.banerjee, k.nguyenthanh, c.fookes}@qut.edu.au).

existing simulators are often proprietary, overly simplistic, or tailored to specific algorithms, hindering fair comparisons and slowing the transition from academic research to industrial application.

To address these, we introduce **Mining-Gym**, a configurable benchmarking framework for truck dispatch scheduling optimization. Unlike prior tools, Mining-Gym offers a comprehensive framework covering the following :

- A high-fidelity discrete-event simulation (DES) capturing real-world complexities,
- Seamless integration with OpenAI Gym for RL algorithm compatibility,
- Extensive customization for diverse mining scenarios,
- Comprehensive logging and visualization tools, and
- Open-source availability for collaborative research.

By standardizing benchmarking, Mining-Gym enables fair algorithm evaluation, facilitates reproducible experiments, and bridges the gap between theoretical RL development and industrial application.

The rest of this paper is structured as follows: Section II reviews related work in mining simulation and RL and highlights the limitations of conventional simulators. Section III details Mining-Gym's architecture, integrating DES with RL. Section IV covers discussions on implementation features, including the configuration system and visualizations. Section V outlines experimental setup and Section VI presents and discusses the results. Finally section VII presents the conclusion and future work.

## II. BACKGROUND AND RELATED WORK

### A. Background

*1) Truck dispatch scheduling optimization:* Managing dispatch scheduling, particularly truck dispatching, represents a pivotal challenge in mining process optimization. Truck dispatching is dedicated to transporting extracted supply materials—both in quantity and quality—from mining fronts, where shovels excavate, to destinations like Crushers, or dumping sites. These dispatching decisions significantly impact operational efficiency, being critically important as a large portion of mining costs are linked to truck-shovel activities Truck dispatching tasks often utilize mathematical programming to reduce equipment waiting times and optimize production [1]. Whereas heuristic methods simplify decision-making by using practical experience instead of exhaustive optimization. These strategies include assigning trucks to the nearest shovel, prioritizing by equipment capacity or material demand, and leveraging historical data patterns [2]. Truck dispatching tasks often

utilize static optimization techniques to minimize equipment waiting times and optimize production [1]. However, such conventional methods typically require re-optimization when complex configurations change, such as during equipment breakdowns, and may not be robust enough to handle the stochastic nature of minesite processes.

However, due to the limitations of static optimization methods in dynamic environments, research is increasingly focusing on alternative, adaptive approaches.

*2) DES and integration with RL in Mining Operations:* To effectively apply RL in mining logistics, it must be integrated with robust simulation frameworks. DES has been widely used in industrial and mining applications for optimizing processes like truck scheduling. By modeling stochastic interactions between equipment and processors, DES captures variability and complexity, often replacing intricate mathematical models with probabilistic parameters [14]. Applications include supply chain evaluation [15] and ensuring adherence to production schedules [16], demonstrating DES's role in handling uncertainty and improving operational efficiency.

Recent research has explored combining RL with DES to simulate complex industrial environments. While standard RL environments like OpenAI Gym provide useful benchmarks, they often lack the realism required for industrial applications [17]. Industrial settings demand detailed, stochastic modeling, similar to DES, to account for dynamic conditions and resource constraints. Integrating RL's trial-and-error learning with DES has gained traction for modeling real-world stochastic systems. For example, [18] transformed DES-based SCT controllers into RL environments, enhancing decision-making in automotive plant control.

This RL-DES integration has been applied to scheduling and optimization across industries. [19] employed Deep Q-Networks (DQN) for flexible job shop problems, outperforming traditional metaheuristics, while [20] used DES and OpenAI Gym to create RL-compatible production scheduling environments, simplifying RL algorithm deployment. In mining, RL has been explored for short-term planning, truck dispatching, and scheduling. [6] introduced a curriculum-driven RL method for vehicle dispatching to address sparse rewards, while [8] developed a real-time RL-based dispatching system for autonomous trucks. Additionally, [13] applied Q-learning to optimize material supply during operational delays. These studies highlight the potential of RL-DES integration in enhancing decision-making and efficiency in mining operations.

*3) RL and the importance of simulator:* RL excels in sequential decision-making, delivering state-of-the-art performance across various domains, including robotics, locomotion control, autonomous driving, and multi-agent systems. [21] The mining truck dispatching problem can be framed as a sequential decision-making task, where truck assignments must adapt to evolving conditions. RL provides a suitable framework for optimizing these dispatching strategies

In mining, RL enables adaptive decision-making by learning optimal dispatching policies through trial-and-error interactions. This approach accommodates dynamic changes in configurations, equipment failures, fluctuating ore quality, and weather, without frequent re-optimization. RL allows continuous refinement of dispatching policies to maximize efficiency and resilience in complex, uncertain environments. However, RL's reliance on environmental interaction presents challenges in real mining due to safety, cost, and timeline concerns. Unlike controlled simulations real-time training is impractical as RL requires extensive exploration. Despite advancements in sample efficiency [22], [23], mining operations still require hundreds to thousands of episodes, hindering real-world deployment.

Simulators are preferred for safe, cost-effective RL training. They enable rapid learning in diverse virtual scenarios, crucial for robust policy development. Common RL simulators like OpenAI Gym [24], although simplified, form the basis for algorithm development. Simulators are vital in mining optimization due to industry complexity and variability. High-fidelity simulations improve accuracy, ensuring learned strategies transfer effectively to real-world operations, enhancing efficiency and safety, making RL a viable tool for mining logistics.

*4) RL Applications in Truck Dispatch Optimization:* A limited number of studies have applied RL to truck dispatch optimization, but these are typically tailored to specific RL algorithms, limiting their adaptability. Adapting them to work with different RL methods would require extensive understanding and potential modifications. Huo et al. [7] apply Q-learning to optimize dispatching in haulage operations, reducing greenhouse gas emissions while maintaining production. Matsui et al. [8] develop a real-time dispatching algorithm using deep RL for autonomous haulage trucks, improving transportation efficiency and fuel consumption. De et al. [10] extend RL to short-term production planning, integrating actor-critic agents for equipment allocation and production scheduling. Chiarot et al. [13] apply Q-learning-based deep RL to reduce delays during shifts and breaks, improving material supply to crushers.

*5) Limitations: Existing simulators:* Refer to Table I for a comparison of conventional simulators based on crucial features. A notable absence of real-time visualization is observed in several studies, including [3], [5]–[10], [13], which can hinder the ability to monitor and manage mining operations effectively. Customizability is another missing feature in [7], [9], and [10], limiting the flexibility of these simulators for adapting to different mining scenarios. Furthermore, the scarcity of open-source solutions, with only [12] providing such a framework, restricts broader adoption and collaborative improvement. A number of works, such as [7], [9], and [10], utilize rule-based architectures that may not be as adaptable or robust as more advanced simulation methods. Additionally, the proprietary nature of many conventional simulators limits research replication and comparison. Many also fail to adequately account for random events affecting key mining components like trucks, shovels, or crushers, as seen in [7]. Finally, many simulators are not designed for RL settings or are tailored to specific algorithms, limiting their versatility [8].

Addressing these limitations is crucial for effectively training and testing RL algorithms to obtain optimal policies, for surface mining process optimization. Proper real-time

| Reference | Year | Uncertainties | Simulator arch. | RL adaptability | RT Visualization | Customizable | Platform used | Open source |
|-----------|------|---------------|-----------------|-----------------|------------------|--------------|---------------|-------------|
| [3] | 2022 | Extensive | DES | N.A. | × | ✓ | Simpy | × |
| [4] | 2022 | Extensive | DES | N.A. | Fair | ✓ | Flexsim | × |
| [5] | 2022 | Extensive | DES | N.A. | × | ✓ | - | × |
| [6] | 2023 | Fair | DES | Custom | × | ✓ | Simpy | × |
| [7] | 2023 | Fair | Rule based | Custom | × | × | OpenAI | × |
| [8] | 2023 | Fair | DES | Custom | × | ✓ | Python | × |
| [9] | 2023 | Fair | Rule based | N.A. | × | × | MATLAB | × |
| [10] | 2023 | Extensive | Rule based | Custom | × | × | - | × |
| [11] | 2024 | Fair | DES | N.A. | Fair | ✓ | Arena | × |
| [12] | 2024 | Extensive | DES | N.A. | Extensive | ✓ | Simpy | ✓ |
| [13] | 2024 | Fair | DES | Custom | Fair | ✓ | DISPATCH | × |
| **Ours** | **2025** | **Extensive** | **DES** | **Adaptable** | **Extensive** | ✓ | **Salabim, Python** | ✓ |

TABLE I: A study of conventional simulators used in Surface Mining process optimization, especially for Truck Dispatching. Comparative analysis of mining simulators based on key features. The table evaluates simulators across several dimensions: **Uncertainties**, ranging from *Extensive* (comprehensive modeling of equipment failures, maintenance, variable haul times) to *Fair* (modeling of limited uncertainties); **Simulator Architecture**, comparing *(DES)* and *Simpler Rule-based Logic*; **RL Adaptability**, which indicates whether the simulator is *Adaptable* (native integration), *Custom* (requires adaptation), or *N.A.* (not designed for RL); **Real-time Visualization**, measuring the extent of *Extensive* (Elaborate animated visuals with KPI dashboards), *Fair* (basic visual representation, or *None*.

visualization and comprehensive data logging are essential for ensuring the repeatability and comparability of experiments. Accurate simulation of real-world complexities and uncertainties, along with seamless integration with widely accepted formats like OpenAI Gym, for ease of use and training of off-the-shelf and custom algorithms.

## III. MINING-GYM: SYSTEM ARCHITECTURE AND DESIGN

The mining process especially the dynamics and uncertainty of the Load-Haul-Dump (crusher or dumping site)-Return-Query (LHDRQ) cycle is captured using a DES based simulation model. While the scheduling problem is modeled as a MDP in order to solve it using RL algorithms (see Fig.1). Following we discuss these two modeling processes.

### A. Mining-Gym overview

The mining gym simulator and benchmarking tool is a sophisticated platform designed for RL applications in mining operations. It features a GUI-based interface where users input crucial parameter values that define the state of the mining site, equipment, and other relevant factors. These parameters are then translated into a human-readable and editable text file, allowing for easy adjustments directly within the simulation environment. Users can choose to either train a scheduler policy or run an existing one, which can be either classical or RL-based. The simulator includes an additional interface that supports the training and testing of modern RL algorithms, enabling advanced simulation capabilities. Furthermore, the tool provides a real-time dashboard displaying key performance indicators (KPIs) and a dynamic, visual representation of the mining site's operations, offering users valuable insights into the system's performance and the effects of various parameters on mining activities.

We present a DES environment for mining processes, integrated with OpenAI Gym—rebranded as Gymnasium in 2023—for RL applications. Gymnasium is a widely used toolkit that standardizes interactions across diverse environments, from control tasks to robotics and video games. It enhances the original Gym with improved modularity, better support, and expanded features, enabling efficient testing, comparison, and debugging of RL algorithms. Our wrapper maps DES inputs/outputs to Gym's reset and step methods, ensuring compatibility with RL frameworks like Stable Baselines. Stable Baselines provides well-tested implementations of RL algorithms such as PPO, A2C, and DQN, simplifying integration with Gym environments. The done flag signals episode termination, while info offers additional simulation insights. Together, Gymnasium and Stable Baselines create a powerful ecosystem for RL research, facilitating reproducible experimentation and benchmarking in complex domains like mining simulations.

### B. Mining Process Modeling: DES Overview

For DES modeling we have used a comprehensive python package named Salabim [26], which offers process interaction methods, queue handling, resources, statistical sampling, and real-time 2D/3D animation capabilities.

*Components* are fundamental building blocks that define the dynamic behavior of entities within the simulation environment. By defining entities as components, SALABIM can simulate complex interactions, resource contention, and event-driven behaviors essential for realistic modeling. Trucks (lower priority) and Breakdown Events (higher priority) are modeled as components in our work. *Resource* is a fundamental component used to model and manage entities that are shared among components within a simulation. Resources represent facilities, equipment, or services that components (such as trucks, shovels, or processes) compete for or utilize during their activities. Specifically, in Mining-Gym, Shovels, Dumps and Crushers are modeled as resources.

In the SALABIM framework, the **Truck** component In the SALABIM framework, the Truck component follows a Load-Haul-Dump (crusher or dumping site)-Return-Query (LHDRQ) cycle. In the loading phase, a truck requests an available shovel as advised by the *dispatcher* and waits its turn. Once granted access, it undergoes loading, simulating material transfer time. The truck then enters the haul phase, traveling a predefined trajectory to the dump (crusher or dumping site),
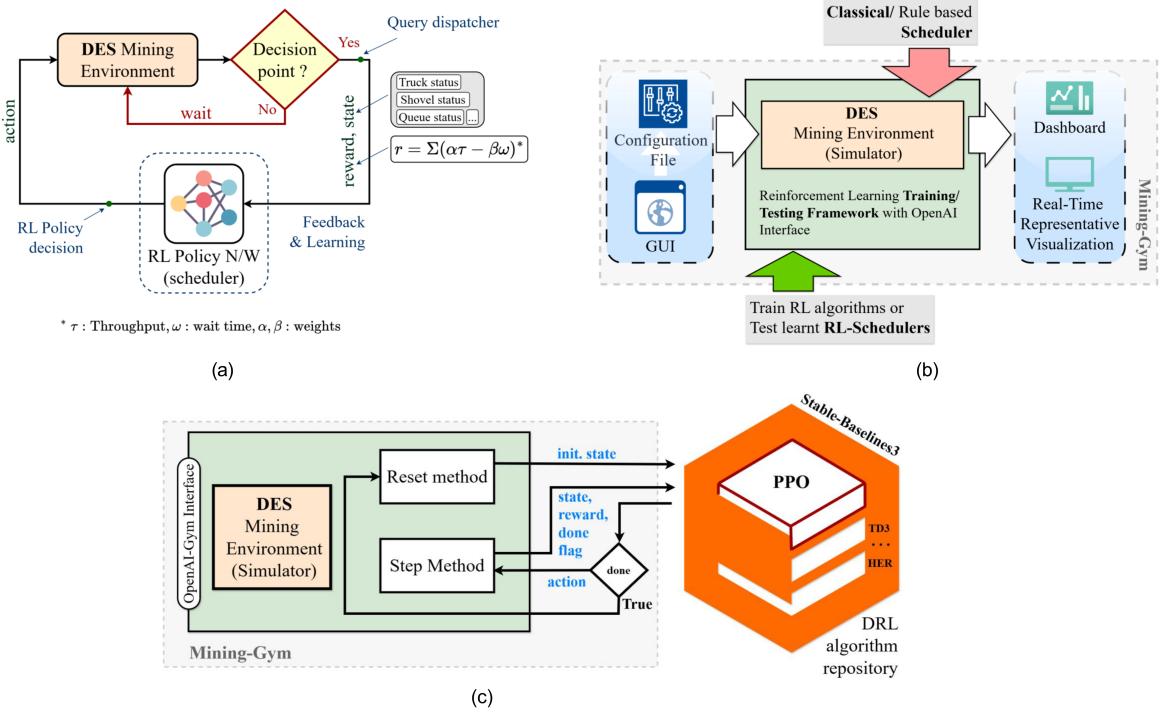
Fig. 1: Overview of the Mining-Gym framework from different perspectives. (a) Event-driven or decision-point-based RL setting: The DES mining environment requests decisions from the RL policy only at specific decision points, such as when a truck requires a shovel assignment. (b) Comprehensive system architecture: Displays all key components of the Mining-Gym, including the graphical user interface (GUI), the generated configuration file used to initialize the environment, and the dashboard with real-time visualizations. (c) OpenAI Gym-compatible RL interface: Illustrates how the Mining-Gym integrates with OpenAI Gym by adapting simulation signals into the standard reset and step methods. This compatibility allows seamless integration with popular RL libraries, such as Stable-Baselines3 [25], enabling easy training and testing of RL models. *Note: the reward function in (a) is for demonstration only.*

represented by holding for travel time. The dump choice is made at the loading site or shovel.

At the dump, the truck enters the dump phase, requesting access to an available crusher or dumping site. After access, it performs dumping, mirroring the time needed to unload ore or waste. The return phase follows, where the truck travels back to the shovel, completing one haul cycle. The query phase then begins, where the truck seeks scheduling or resource allocation from the dispatcher. The **Dispatcher agent** is the core scheduling strategist, optimizing mining operations by balancing resource utilization and minimizing wait times. It employs basic strategies like fixed schedules or nearest-first, as well as advanced approaches such as RL-trained neural network policies.

A truck may request the following resource allocations: 1) Shovel, 2) Crusher, 3) Dumping site, and 4) Route (not currently considered). Mining-Gym's dispatcher modules handle these requests (except routing) and supports baseline (e.g. Random) aswell as learned strategies (RL based). Only one dispatcher can be replaced with a learnable NN-based learnable policy at a time, with *shovel allocation* set as the default for all experiments.

**Modeling breakdown events** is crucial for simulating disruptions when trucks, shovels, or dumps become unavailable due to failures. These breakdowns impact loading, dumping,

and transit, affecting efficiency and resource use. A preemption handler enhances realism by managing resource interruptions. For example, if a shovel breaks down mid-loading, the preemption handler pauses operations until repairs or replacement occur, accurately reflecting real-world maintenance disruptions.

**Additional features :** in the final Mining-Gym simulation build upon the basic workflow (Fig. 2).

- *Choice between Dump Types:* The stripping ratio of waste to ore is represented by the parameter $\epsilon$, which controls the probability that a truck carrying ore is directed to the crusher rather than the dumping site. Specifically, with probability $\epsilon$, the truck is directed to the crusher, and with probability $1 - \epsilon$, it is sent to the dumping site. This parameter can be adjusted to balance waste and ore management, adapting to varying operational conditions.
- *Scalable Configurations:* The simulation supports adjustable configurations for Trucks, Shovels, Crushers, and Dumping sites, allowing for scalable operational simulations that can be tailored to different mining scenarios.
- *Uncertainty Modeling:* To reflect the stochastic nature of real mining operations, the simulation incorporates sampling from various probability distributions. These distributions, along with their parameters, are configurable,
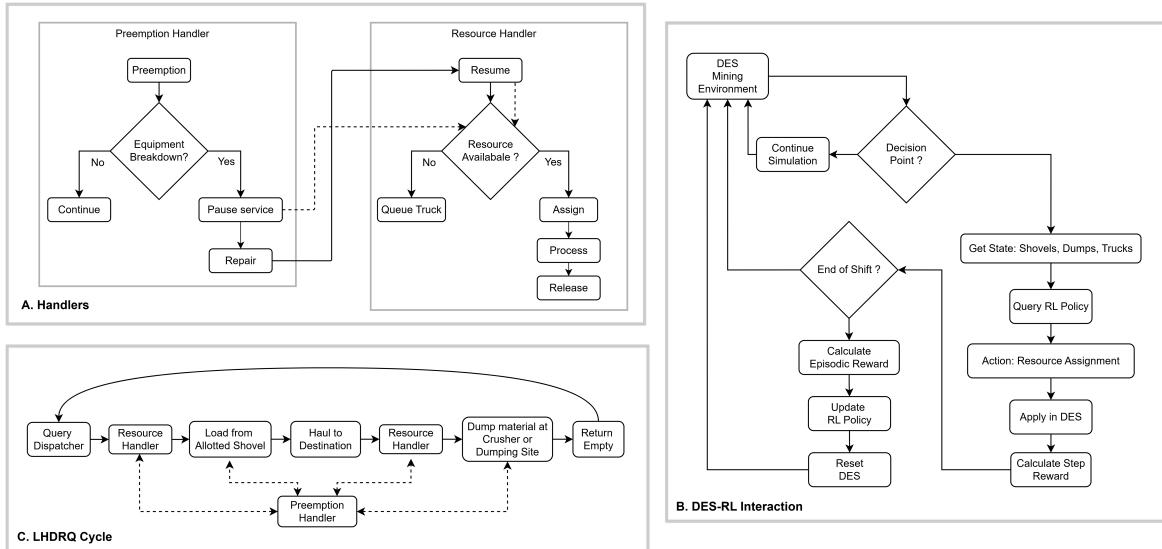
Fig. 2: (A) Simplified mining system showing three key components: (1) Resource Handler managing resource availability and assignments, (2) Preemption Handler detecting breakdowns and managing repair processes (B) DES-RL interaction flow illustrating how the RL policy integrates with the DES. At decision points, the environment state is processed by the RL policy to determine resource assignments. Immediate or step rewards guide learning during simulation, while the episodic reward at shift (or episode) end updates the policy before environment reset. (C) Load-Haul-Dump-Return-Query (LHDRQ) cycle illustrating the truck's journey through the mining process, which begins with querying the dispatcher for assignments, followed by loading material, hauling to the destination, dumping, and returning empty. Breakdown events, managed by the Preemption Handler, can interrupt operations at any stage.

enabling flexible modeling of operational uncertainty.

- *Dispatcher Strategies:* Non-learnable, baseline dispatcher strategies (default random) for resources like Crushers and Dumping sites are included. These strategies allow for the evaluation of different dispatching approaches without relying on RL.

### C. Modeling the Dynamic Dispatching Problem as MDP

*1) Typical RL:* In a typical RL setting, the agent interacts with the environment at every time step $t$. The environment is modeled as a Markov Decision Process (MDP) defined by the tuple $(S, A, P, R)$, where $S$ is the set of states, $A$ is the set of actions, $P(s'|s, a)$ is the state transition probability, and $R(s, a)$ is the reward function. At each time step, the agent selects an action $a_t$ based on the current state $s_t$, receives a reward $r_{t+1}$, and transitions to a new state $s_{t+1}$. The goal is to learn a policy $\pi(a|s)$ that maximizes the expected cumulative reward $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t]$, where $\gamma$ is the discount factor.

*2) Decision-point based or Event-driven RL:* In mining truck scheduling, an event-driven reinforcement learning (RL) approach is adopted, where the agent interacts with the environment only at discrete decision points $d$, rather than at every time step. A decision point represents a specific instance within the operational environment where the agent must choose an action based on the observed state $s$, which includes truck and job statuses. The reward $r$ at each decision point reflects performance metrics such as reduced waiting times and increased throughput, while the action $a$ involves dispatching trucks to different tasks. The framework aims to optimize the expected cumulative reward, given by $\mathbb{E}[\sum_{d=0}^{\infty} \gamma^d r_d]$, with

agent-environment interactions constrained to these discrete decision points to align with real-world operational constraints. This approach has been widely used in complex environments where continuous interaction is impractical or unnecessary, e.g. in vehicle routing [27] and supply chain management [28].

*3) Defining Agent-Environment interaction and the MDP:* Formulating truck dispatch scheduling in open-pit mining as a Markov Decision Process (MDP) involves defining the components of MDP: states, actions, rewards, and transitions. The Mining-Gym framework considers $TR$ trucks (represented as $\tau$) interacting with the DES based mining environment. At every decision point $d \in D$, after dumping the material into the appropriate location (crusher or dumping site), a new resource assignment is requested by a truck $\tau_i$ where $i \in N$. The dispatcher agent $\mathcal{D}$ observes the current state $S_d \in \mathcal{S}$, where $S_d$ represents the current status of the mining complex's performance at decision point $d$, and takes an action $A_d^i \in \mathcal{A}$, determining the next shovel to which the truck $i$ will be assigned.

Below, we define the various components of the MDP used in this work:

1) *States:* The state represents the current status of the mining operation. The state of the system at a given time must encode all the features needed for the agent to learn a relationship with the desired objective to be maximized. For this task the state of the system is encoded as a vector with the following components

$$s^d = [SA_d, TA_d] \tag{1}$$

where,

$SA_d$ represents shovel-related attributes as an encoded vector at the current decision point $d$. It includes *Shovel ID*, encoded in 3-bit chunks, *Queue Length*, representing the number of trucks waiting per shovel, and **Shovel Status**, a binary indicator of shovel availability (online/offline). Similarly, $TA_d$ captures truck-related attributes, including *TruckID*, encoded in 5-bit chunks, *Trips complete*, tracking the number of trips per truck, and *Trip Status*, a multi-bit representation of the truck's current state, such as loading, transit, or maintenance (see Table: II).

| Category | Attribute / Status | Description |
|---|---|---|
| | Shovel ID | 3-bit binary |
| Shovels ($SA_d$) | Queue Length | Normalized float |
| | Shovel Status | Binary (1 or 0) |
| | Truck ID | 5-bit binary |
| Trucks ($TA_d$) | Trips Complete | Normalized float |
| | Truck Status | 3-bit binary |
| **Truck Status Codes** | | |
| **Truck Status** | **Binary Code** | |
| At Shovel | 000 | |
| At Crusher | 001 | |
| At Dumping Site | 010 | |
| Moving from Shovel to Crusher | 011 | |
| Moving from Shovel to Dumping Site | 100 | |
| Moving from Crusher to Shovel | 101 | |
| Moving from Dumping Site to Shovel | 110 | |
| Breakdown | 111 | |

TABLE II: State Space Description and Truck Status Codes

2) *Actions:* Actions are decisions made at each decision step (d), by the RL policy.
In our framework, an action corresponds to resource or shovel allocation and the action space is defined as:

$$A_d \in \{1, 2, \ldots, SH\}$$

where $SH$ represents the total number of shovels, and the action value indicates the shovel ID to which a truck is assigned at decision point $d$.

3) *Rewards:* The reward function quantifies both immediate and cumulative benefits, defining objectives to maximize. In our event-based RL setting, shifts are episodes with sparse rewards given for critical milestones like meeting production targets. The environment features a long-horizon episodic reward and intermediate global rewards to guide learning.
Once the dispatching agent $\mathcal{D}$ outputs action $A_i^d$, the DES environment responds with:

$$R_i^d = r_{\text{imm}}^d + \left( r_{\text{epi}}^d \text{ if } d = d_T \text{ else } 0 \right) \quad (2)$$

where $r_{\text{epi}}^d$ is the episodic reward at the end of an episode (or shift), and $r_{\text{imm}}^d$ is the immediate reward per decision-step $d$.

*Immediate Reward Formulation:* The immediate reward uses an exponentially weighted sliding window over the latest $k$ decision points:

$$r_{\text{imm}}^d = -\alpha \cdot \hat{T}T_{\text{Avg}} - \beta \cdot \hat{Q}_{\text{Avg}_d} - \gamma \cdot (1 - D_{\text{DivScr}}) \quad (3)$$

$$= -\alpha \cdot \frac{\sum_{j=i-k+1}^{i} w_j \cdot \tau_j}{\sum_{j=i-k+1}^{i} w_j} - \beta \cdot \frac{\sum_{j=i-k+1}^{i} w_j \cdot Q_{SH_j}}{\sum_{j=i-k+1}^{i} w_j}$$
$$- \gamma \cdot (1 - D_{\text{DivScr}}) \quad (4)$$

where $w_j = e^{\lambda(j-i+k)}$ represents exponential weighting, and $\alpha, \beta, \gamma$ are weighting parameters such that $\alpha + \beta + \gamma = 1$.

The components are:
*Average Trip Time:*

$$\tau_d = \frac{1}{|\Omega(d)|} \sum_{\omega \in \Omega(d)} \mathcal{T}_\omega(d) \quad (5)$$

*Shovel Queue Time:*

$$Q_{SH_d} = \frac{1}{|\Gamma|} \sum_{\gamma \in \Gamma} \frac{1}{|\Omega_\gamma(d)|} \sum_{\omega \in \Omega_\gamma(d)} Q_\omega(d) \quad (6)$$

*Diversity Score:*

$$D_{\text{DivScr}} = \frac{|\{\gamma \in \Gamma \mid A_\gamma > 0\}|}{|\Gamma|} \quad (7)$$

where $A_\gamma$ is the number of trucks assigned to shovel $\gamma$ in the recent window.

*Episodic Reward Formulation:* The episodic reward balances production efficiency and resource utilization:

$$r_{\text{epi}}^{d_T} = \omega_1 \cdot P_{\text{Efficiency}} - \omega_2 \cdot D_{\text{Imbalance}} \quad (8)$$

$$= \omega_1 \cdot \min\left(1, \frac{P_{\text{Vol}}}{P_{\text{Vol}_{\text{Target}}}}\right) - \omega_2 \cdot \frac{D_{\text{Penal}}}{\sum_{\gamma \in \Gamma} A_\gamma} \quad (9)$$

where: $P_{\text{Vol}}$ is total material produced per shift, $P_{\text{Vol}_{\text{Target}}}$ is the production target per shift. $D_{\text{Penal}} = \max_{\gamma \in \Gamma} A_\gamma - \min_{\gamma \in \Gamma} A_\gamma$ measures shovel utilization imbalance, $\omega_1, \omega_2$ are weighting parameters where $\omega_1 + \omega_2 = 1$.
This structure balances short-term efficiency with long-term production goals, guiding RL agents in open-pit mining operations. During each episode (e.g., shifts of 12 hours), the agent takes $N_{\text{steps}}$ actions, and the discounted sum of rewards defines the return $G_t$:

$$G_t = \sum_{k=t+1}^{N_{\text{steps}}} \gamma^{k-t-1} R_k$$

where $\gamma$ is the discounting factor that determines the impact of future actions on the objective function. The objective is to maximize the expected return $G_t$ by training the agent to improve actions, ensuring trucks meet production targets and minimize queue formation:

$$\max_\pi \mathbb{E}[G_t | S_t = s_t, A_t = a_t]$$

*Policy:* The policy defines the strategy that the agent follows to select actions (shovel IDs) based on the current state of the system. Given the state $s_d = [SA_d, TA_d]$,

which encodes the current shovel and truck attributes, the policy predicts the shovel ID to which a truck should be assigned at decision point $d$. The policy maximizes cumulative reward by optimizing truck-shovel assignments for operational efficiency, adapting over time based on rewards to meet short- and long-term production goals, ensuring efficient dispatch and shovel utilization. The policy $\pi_\theta(a|s)$ is a function that gives the probability of taking action $a$ (i.e., assigning a specific shovel ID) given the current state $s$:

$$\pi_\theta(a|s) = \mathbb{P}(A_d = a | S_d = s)$$

where $A_d$ is the action taken at decision point $d$, which is the shovel ID assigned to a truck, and $S_d = s$ is the state of the system at decision point $d$, including both shovel-related attributes $S_{A_d}$ and truck-related attributes $T_{A_d}$. This ensures that the agent's actions evolve to maximize long-term efficiency and achieve the desired operational goals.
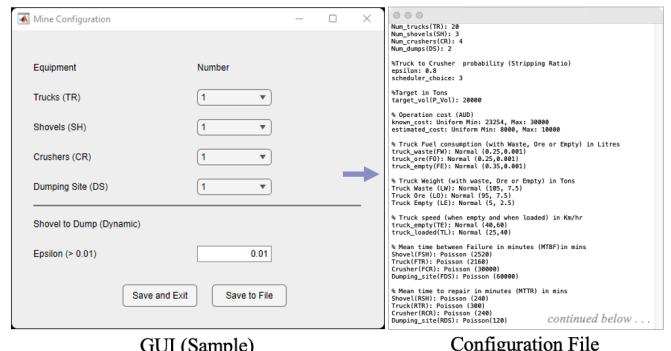
## IV. IMPLEMENTATION FEATURES

In this section we introduce different visualizations and other GUI-based tools that we provide with the MINING-GYM simulator, for better performance comparison and studying change in KPIs

### A. Configuration System

| Variables | | Source |
|---|---|---|
| **A. Operational Parameters** | | |
| **Number of Equipment:** | | - |
| Num. of Trucks | TR | User |
| Num. of Shovels | SH | User |
| Num. of Crushers | CR | User |
| Num. of Dumps | DS | User |
| **Queue Size (waiting for resource) $Q$:** | | - |
| At Shovel, Crusher, Dumping site | $Q_{SH}, Q_{CR}, Q_{DS}$ | DES |
| **Load per trip $L$** | | - |
| Truck carrying Waste, Ore | $L_W, L_O$ | User |
| **Truck Speed:** | | - |
| Empty truck, Loaded truck | $S_{EmTR}, S_{LodTR}$ | User |
| **Others:** | | - |
| Number of trips (SH-CR-DS-SH) | $N$ | DES |
| Available (operation unit) time | $T_{SHF}$ | User |
| Dumping and maneuver time (min) | $T_{DM}$ | User |
| Shift duration in minutes | $S_{dur}$ | User |
| Num_shifts | SN | User |
| **B. Cost and Financial Parameters** | | |
| Known cost, Estimated cost | $C_{KW}, C_{EST}$ | User |
| **C. Equipment Performance and Efficiency** | | |
| **Unit Fuel Consumption ($F_{Unit}$):** | | - |
| Truck with Waste, Ore, Empty | $F_W, F_O, F_E$ | User |
| **Equipment off time $OF_E$** | | - |
| Truck, Shovel, Crusher, Dumping Site | $OF_{TR}, OF_{SH},$ $OF_{CR}, OF_{DS}$ | DES |
| **Equipment idle time $IDL_E$** | | - |
| Truck, Shovel, Crusher, Dumping Site | $IDL_{TR}, IDL_{SH},$ $IDL_{CR}, IDL_{DS}$ | DES |
| **Equipment mean time between failure:** | | - |
| Shovel, Truck, Crusher, Dumping Site | $F_{SH}, F_{TR}, F_{CR}, F_{DS}$ | User |
| **Equipment mean time to repair:** | | - |
| Shovel, Truck, Crusher, Dumping Site | $R_{SH}, R_{TR}, R_{CR}, R_{DS}$ | User |
| **D. Loading and Other Time:** | | |
| Truck loading | TRL | User |
| Dump to Shovel | DTS | User |
| Shovel to Dump | STD | User |
| Truck dumping @ Dump | TRDM | User |
| Crusher to Shovel | CTS | User |
| Shovel to Crusher | STC | User |
| Truck dumping @ Crusher | TRCR | User |
| **E. Other Parameters:** | | |
| Probability of choosing Crusher | $\epsilon$ | User |
| Target Production volume | $P_{Vol, Targ}$ | User |

TABLE III: Overall list of variables and parameters used

GUI (Sample)  Configuration File

(a)

(b)

Fig. 3: (a) Shows a sample page from the GUI interface and the resulting configuration file (b) Shows the *Real-time representative visualization* of the mine-site. The screenshot displays trucks queued at shovels, with the rightmost shovel offline. It also shows trucks moving between the dumping site, crushers, and the shovels.

As discussed in the previous section, a large number of parameters must be initialized to run the simulator (see Table III and Table V). To provide an organized method for supplying this information, we have designed a GUI interface. This interface allows users to set values and distributions for stochastic parameters, which are then saved to a human-readable text file. This configuration file is used by the simulator to initialize its parameters during operation (see Fig. 3a).

### B. Other Visualization Features

We provide a *real-time representative visualization* of the mining setup as configured in the setup file. It features animated truck movements between loading points, crushers, and dumping sites, along with real-time queue formation. The visualization is runnable both during training and when executing a trained model or any other algorithm (see Fig. 3b).

*Key Performance Indicators (KPIs)* in mining process optimization include metrics like equipment utilization, cycle

time, production throughput, fuel efficiency, and downtime. These metrics help evaluate the efficiency, productivity, and cost-effectiveness of mining operations, guiding improvements and strategic decisions to enhance overall performance and profitability. A key factor in KPI measurement is the time frame, which can range from shifts to daily, weekly, monthly, or yearly intervals. By default, we consider a 12-hour shift as the standard, but this can be customized as a parameter.

We include the following essential KPIs :

1) **Total Production** ($P_{VOL}$): Refers to the total quantity of material (ore + waste) excavated per shift. It is measured in tons, indicating the scale and efficiency of production activities.

$$P_{VOL} = N \times L$$

where: $N$ = Total number of trips in a shift, L = Material load per trip (tons)

2) **Trips per Hour** (TPH$_t$): The *"Trips per Hour"* metric quantifies the rate at which the full haul cycle (SH-CR-DS-SH) is completed, normalized to an hourly basis. It provides a standardized measure of fleet productivity by capturing the number of new trips executed within the most recent hourly interval.

$$\text{TPH}(t) = \frac{\Delta t}{60} \left[ N(t) - N(t - \Delta t) \right]$$

where $N(t)$ = Total number of trips (SH-CR-DS-SH) completed by all trucks up to time $t$, $N(t - \Delta t)$ = Total trips at $t - \Delta t$, $\Delta t$ = Measurement interval in minutes (typically 60).

*Total Number of Trips:*

$$N(t) = \sum_{i=1}^{TR} C_i(t)$$

where $TR$ = Total number of trucks in the fleet, $C_i(t)$ = Trips completed by truck $i$ up to time $t$.
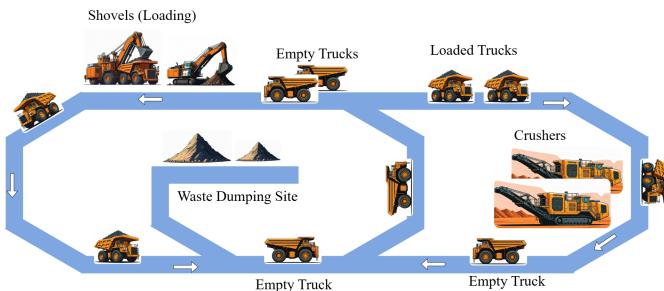


Fig. 4: Conceptual Diagram of Minesite

3) **Cost per Ton** ($CPT$): Refers to the financial expenditure incurred for every unit weight (typically one ton) of material produced, processed, or transported.

$$CPT = \frac{C_{KW} + C_{EST}}{P_{VOL}}$$

where $C_{KW}$ = Total known cost and $C_{EST}$ = Total estimated cost

4) **Fuel Consumption** ($FC$): measures the amount of fuel consumed by mining equipment or vehicles relative to the total material hauled during operations:

$$FC = \frac{\text{Total fuel consumed}}{P_{\text{Vol}}}$$

where the Total fuel consumed $= N \times F_{\text{Unit}}$, and the $F_{\text{Unit}} = F_{\text{Unit}} = F_W + F_O + F_E$.

Here, $N$ is the number of trips completed by all trucks, $F_W$, $F_O$, and $F_E$ are the fuel consumptions for waste, ore, and empty trips, respectively, and $P_{\text{Vol}}$ is the target production volume.

We have also provided an interactive dashboard as part of this benchmark to analyze change in these KPIs in real time.

## V. EXPERIMENTAL SETUP

*1) Environment Design:* We consider a medium sized Bauxite ($5 - 10$ Million Tons output) mine with 3 Shovels, 20 Trucks, 4 Crushers and 2 Dumping sites, detailed specifications are presented in Table: IV and V. A conceptual diagram of the minesite is provide in Fig. 4.

| Equipment | Model/ Make | Quantity | Specification (each) |
|---|---|---|---|
| Shovels per pit | CAT 6040 | 2 (total 4) | 43.7 t bucket size |
| Trucks | CAT 777 | 20 | 98.2 nominal payload |
| Crushers | Metso Lokotrack (LT106) | 3 | 450 t / hr capacity |
| Dumping site | - | 2 | Unlimited |

TABLE IV: Equipment List with Specifications

*2) Training Setup:* To ensure that the trained agent generalizes well across the testing scenarios, we not only incorporate stochastic parameters to randomize breakdown and repair times for mining equipment and trucks but also enhance the reward function by introducing queue buildup metrics (queue length) and diversity score to prevent repeatedly selecting same shovels.

The parameter values used during training is provided in Table V. For training an RL policy we have used the PPO [29]algorithm from StableBaselines3 [25]). We have used all the default parameter and training settings as present in the repository's PPO implementation. Note the **scenario parameters** which help specify the failure scenarios to be simulated during training or testing. For this proof-of-concept implementation and demonstration the results of only $400$ episodes or shifts of training.

*3) Testing Setup:* We designed scenarios to test the efficacy of the RL trained scheduler against a vanilla scheduler (random scheduler) as baseline, under different challenging scenarios.

For a random scheduler, the dispatcher's action $A_d$ is selected randomly from the available set of shovel IDs $\{1, 2, \ldots, SH\}$. The selection follows a uniform distribution,

$$P(A_d = a | S_d = s) = \frac{1}{SH} \quad \text{for all } a \in \{1, 2, \ldots, SH\}$$

i.e. any state $S_d = s$, each action $A_d = a$ (where $a$ is a shovel ID) is equally likely to be chosen, with a probability of $\frac{1}{SH}$. In contrast, the RL policy $\pi_\theta(a|s)$ is a learned policy, where the action $A_d = a$ is selected based on the current state $S_d = s$ and the parameters $\theta$. The policy is defined as:
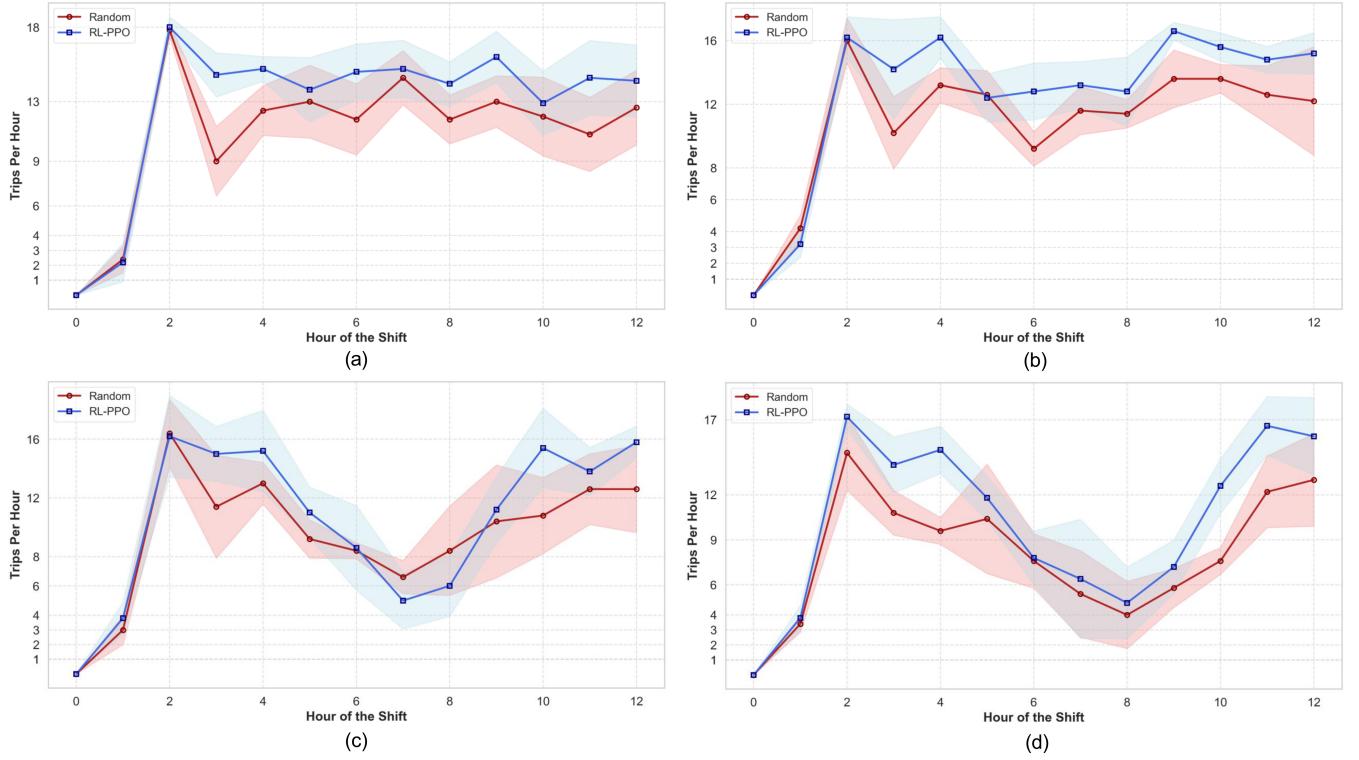
Fig. 5: Comparison of *"Trips per hour"* over the shift under different failure scenarios: (a) No failure, (b–d) $(SH_0, TR_0)$ values where $SH_0$ (Num. Shovels failed) and $TR_0$ (Num. Trucks failed) are set as follows—(b) $(0,6)$, (c) $(1,0)$, (d) $(1,6)$. Solid lines represent mean values, and shaded regions indicate variance from repeated runs. For shovels, failures initiate between 100–300 minutes ( 2–5 hours), and for trucks, between 100–500 minutes ( 2–8 hours) within a 12-hour shift.

$$\pi_\theta(a|s) = P(A_d = a | S_d = s)$$

This probability is learned through RL algorithm where the agent adjusts $\theta$ based on rewards from its actions, optimizing its decision-making process over time. We tested the schedulers under different challenging scenarios:

1) All equipment available, represented as $(SH_0, TR_0) = (0,0)$, i.e., no Shovels and Trucks offline.
2) Six Trucks under maintenance, i.e., $(SH_0, TR_0) = (0,6)$.
3) One Shovel offline, i.e., $(SH_0, TR_0) = (1,0)$.
4) One Shovel and Six Trucks unavailable, i.e., $(SH_0, TR_0) = (1,6)$.

All failure scenarios are repeated according to a consistent temporal schedule over five repeat trials, with the inherent stochasticity in event timings due to probabilistic distributions. These repeats are then considered in the results section for analysis and discussion.

## VI. RESULTS AND DISCUSSION

The RL scheduler demonstrated significant performance improvements over the random scheduler across all scenarios, see Fig.5 and Fig. 6 :

- *Trips Per Hour:* RL increased productivity by 18.2% on average, with the greatest improvement (27.2%) observed
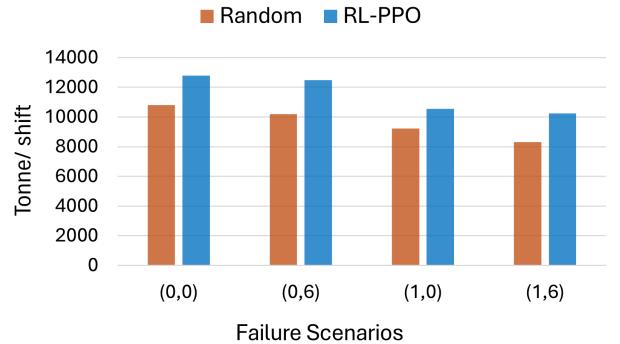


Fig. 6: Comparison of production volume $P_{vol}$ in metric tonnes per shift, averaged over 5 repeats trials.

in the most challenging scenario (1 shovel and 6 trucks offline).

- *Production Volume:* RL delivered 19.7% higher production volumes on average, with improvements ranging from 14.3% to 23.3% across all scenarios.

The performance advantage of RL was most pronounced in scenarios with multiple resource constraints, suggesting that RL's optimization capabilities become more valuable as operational complexity increases. As is showing through the gain in production volume (see Fig.6) RL improved production volume across different scenarios. In the (0,0) case, where all

| Stochastic Parameters | Probability Distributions |
|---|---|
| **Fuel consumption** | **Liters** |
| Loaded Truck (FL) in lt/ton-km | Normal(0.25, 0.001) |
| Empty Truck (FE) | Normal(0.35, 0.001) |
| Shovel (Lt/hr) | Uniform(250, 300) |
| **Equipment specifications** | **Value** |
| Truck payload (TP) | Normal(98, 0.5) t |
| Shovel bucket capacity (SCP) | Normal(43, 0.5) t |
| Empty Truck Speed (TS) | Normal(40, 60) km/hr |
| Loaded Truck Speed (TL) | Normal(25, 40) km/hr |
| **Equipment Maintenance**** | **Time (in mins )** |
| MTBF Shovel (FSH) | Poisson(2520) |
| MTBF Truck (FTR) | Poisson(2160) |
| MTBF Crusher (FCR) | Poisson(30000) |
| MTBF Dumping Site (FDS) | Poisson(60000) |
| MTTR Shovel (RSH) | Poisson(240) |
| MTTR Truck (RTR) | Poisson(300) |
| MTTR Crusher (RCR) | Poisson(240) |
| MTTR Dumping Site (RDS) | Poisson(120) |
| **Loading and Maneuver Time** | **Time (in mins.)** |
| Time Truck Loading (TRL) | Normal(8, 2) |
| Travel Time Shovel to Dump (STD) | Normal(15, 5) |
| Time Truck Dumping (TRDM) | Normal(5, 1) |
| Travel Time Dump to Shovel (DTS) | Normal(15, 5) |
| Travel Time Shovel to Crusher (STC) | Normal(15, 2) |
| Time Truck Unloading at Crusher (TRCR) | Normal(5, 1) |
| Travel Time Crusher to Shovel (CTS) | Normal(15, 2) |
| **Scenario Parameters** | **-** |
| Shovel To Fail (STF) | Binomial (2, 0.5) |
| Trucks To Fail (TTF) | Binomial (6, 0.5) |
| Shovel Initial Breakdown (SIB) | Uniform (100,300) |
| Truck Initial Breakdown (TIB) | Uniform (100,500) |

TABLE V: Distributions for stochastic parameters used in the example. MTBF: Mean Time Between Failures, MTTR: Mean Time To Repair.

equipment was available, the increase was 18.4%. With six trucks offline (0,6), the improvement rose to 22.7%. When one shovel was offline (1,0), production volume increased by 14.3%. Under both constraints (1,6), RL achieved the highest improvement at 23.3%.

The RL scheduling (RL-PPO) consistently outperforms random scheduling baseline method across multiple metrics and scenarios, maintaining higher productivity, especially mid-shift. Its intelligent decisions optimize resource use, making it highly effective in real-world mining operations with equipment constraints, leading to increased production.

## VII. CONCLUSION AND FUTURE WORK

The introduction of Mining-Gym addresses a critical gap in mining process optimization by providing a configurable, open-source benchmarking environment for RL algorithm evaluation in truck dispatch scheduling. Our framework offers high-fidelity DES modeling of real-world complexities, seamless integration with OpenAI Gym, and comprehensive visualization tools. Preliminary experiments show that RL-based scheduling achieves up to 27.2% higher productivity and 23.3% greater production volume, with the most significant gains in resource-constrained scenarios. Mining-Gym provides a standardized evaluation framework, ensuring reproducible

research and supporting the adoption of RL solutions in industrial mining operations.

The development of Mining-Gym as a benchmarking environment will continue in the following directions:

- *Enhanced Simulation Fidelity*: Future iterations will improve traffic modeling, incorporate geographical factors, model heterogeneous equipment fleets, and include additional operational constraints to enhance simulation realism.
- *Expanded Scenario Library*: We will develop a broader set of standardized test scenarios to cover diverse challenges such as weather disruptions, varying ore characteristics, and dynamic production targets.
- *Multi-objective Optimization Benchmarks*: The framework will be extended to support multi-objective optimization, balancing production, cost, equipment lifespan, and environmental impacts—reflecting real-world mining decision-making.
- *Integration with Digital Twin Capabilities*: Enhancing Mining-Gym with digital twin functionalities will connect simulations with real-world data, facilitating data-driven optimizations and bridging the gap between simulation and operational implementation.

### REFERENCES

[1] C. H. Ta, J. V. Kresta, J. F. Forbes, and H. J. Marquez, "A stochastic optimization approach to mine truck allocation," *International journal of surface mining, reclamation and environment*, vol. 19, no. 3, pp. 162–175, 2005.

[2] F. Soumis, J. Ethier, and J. Elbrond, "Truck dispatching in an open pit mine," *International Journal of Surface Mining, Reclamation and Environment*, vol. 3, no. 2, pp. 115–119, 1989.

[3] X. Zhang, A. Guo, Y. Ai, B. Tian, and L. Chen, "Real-time scheduling of autonomous mining trucks via flow allocation-accelerated tabu search," *IEEE transactions on intelligent vehicles*, vol. 7, no. 3, pp. 466–479, 2022.

[4] Y. Zhang, Z. Zhao, L. Bi, L. Wang, and Q. Gu, "Determination of truck–shovel configuration of open-pit mine: a simulation method based on mathematical model," *Sustainability*, vol. 14, no. 19, p. 12338, 2022.

[5] N. Dendle, E. Isokangas, and P. Corry, "Efficient simulation for an open-pit mine," *Simulation Modelling Practice and Theory*, vol. 117, p. 102473, 2022.

[6] X. Zhang, G. Xiong, Y. Ai, K. Liu, and L. Chen, "Vehicle dynamic dispatching using curriculum-driven reinforcement learning," *Mechanical Systems and Signal Processing*, vol. 204, p. 110698, 2023.

[7] D. Huo, Y. A. Sari, R. Kealey, and Q. Zhang, "Reinforcement learning-based fleet dispatching for greenhouse gas emission reduction in open-pit mining operations," *Resources, Conservation and Recycling*, vol. 188, p. 106664, 2023.

[8] K. Matsui, J. Escribano, and P. Angeloudis, "Real-time dispatching for autonomous vehicles in open-pit mining deployments using deep reinforcement learning," in *2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2023, pp. 5468–5475.

[9] X. Wang, Q. Dai, Y. Bian, G. Xie, B. Xu, and Z. Yang, "Real-time truck dispatching in open-pit mines," *International Journal of Mining, Reclamation and Environment*, vol. 37, no. 7, pp. 504–523, 2023.

[10] J. P. de Carvalho and R. Dimitrakopoulos, "Integrating short-term stochastic production planning updating with mining fleet management in industrial mining complexes: an actor-critic reinforcement learning approach," *Applied Intelligence*, vol. 53, no. 20, pp. 23 179–23 202, 2023.

[11] A. Moradi Afrapoli, S. P. Upadhyay, and H. Askari-Nasab, "A nested multiple-objective optimization algorithm for managing production fleets in surface mines," *Engineering Optimization*, vol. 56, no. 3, pp. 378–391, 2024.

[12] S. Meng, B. Tian, X. Zhang, S. Qi, C. Zhang, and Q. Zhang, "Openmines: A light and comprehensive mining simulation environment for truck dispatching," *arXiv preprint arXiv:2404.00622*, 2024.

[13] T. V. Chiarot Villegas, S. F. Segura Altamirano, D. M. Castro Cárdenas, A. M. Sifuentes Montes, L. I. Chaman Cabrera, A. S. Aliaga Zegarra, C. L. Oblitas Vera, and J. C. Alban Palacios, "Improving productivity in mining operations: a deep reinforcement learning model for effective material supply and equipment management," *Neural Computing and Applications*, pp. 1–13, 2024.

[14] A. M. Law, *Simulation modeling and analysis*, vol. 3.

[15] P. Bodon, C. Fricke, T. Sandeman, and C. Stanford, "Combining optimisation and simulation to model a supply chain from pit to port," *Advances in Applied Strategic Mine Planning*, pp. 251–267, 2018.

[16] F. Manríquez, J. Pérez, and N. Morales, "A simulation–optimization framework for short-term underground mine production scheduling," *Optimization and Engineering*, vol. 21, pp. 939–971, 2020.

[17] C. D. Hubbs, H. D. Perez, O. Sarwar, N. V. Sahinidis, I. E. Grossmann, and J. M. Wassick, "Or-gym: A reinforcement learning library for operations research problems," *arXiv preprint arXiv:2008.06319*, 2020.

[18] K. M. Zielinski, L. V. Hendges, J. B. Florindo, Y. K. Lopes, R. Ribeiro, M. Teixeira, and D. Casanova, "Flexible control of discrete event systems using environment simulation and reinforcement learning," *Applied Soft Computing*, vol. 111, p. 107714, 2021.

[19] S. Lang, F. Behrendt, N. Lanzerath, T. Reggelin, and M. Müller, "Integration of deep reinforcement learning and discrete-event simulation for real-time scheduling of a flexible job shop production," in *2020 Winter Simulation Conference (WSC)*. IEEE, 2020, pp. 3057–3068.

[20] S. Lang, M. Kuetgens, P. Reichardt, and T. Reggelin, "Modeling production scheduling problems as reinforcement learning environments based on discrete-event simulation and openai gym," *IFAC-PapersOnLine*, vol. 54, no. 1, pp. 793–798, 2021.

[21] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

[22] C. Banerjee, Z. Chen, and N. Noman, "Enhancing exploration in actor-critic algorithms: An approach to incentivize plausible novel states," *Authorea Preprints*, 2023.

[23] ——, "Improved soft actor-critic: Mixing prioritized off-policy samples with on-policy experiences," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[24] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.

[25] A. Raffin, A. Hill, A. Gleave, A. Kanervisto, M. Ernestus, and N. Dormann, "Stable-baselines3: Reliable reinforcement learning implementations," *Journal of Machine Learning Research*, vol. 22, no. 268, pp. 1–8, 2021. [Online]. Available: http://jmlr.org/papers/v22/20-1364.html

[26] R. van der Ham, "salabim: discrete event simulation and animation in python," *Journal of Open Source Software*, vol. 3, no. 27, p. 767, 2018.

[27] F. D. Hildebrandt, B. Thomas, and M. W. Ulmer, "Where the action is: Let's make reinforcement learning for stochastic dynamic vehicle routing problems work!" *arXiv preprint arXiv:2103.00507*, 2021.

[28] P. Hammler, N. Riesterer, and T. Braun, "Fully dynamic reorder policies with deep reinforcement learning for multi-echelon inventory management," *Informatik Spektrum*, vol. 46, no. 5, pp. 240–251, 2023.

[29] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.