# Dynamic Dispatching for Large-Scale Heterogeneous Fleet via Multi-agent Deep Reinforcement Learning

Chi Zhang*
*Industrial AI Lab*
*Hitachi America Ltd.*
Santa Clara, CA

Philip Odonkor*
*Stevens Institute of Technology*
Hoboken, NJ

Shuai Zheng[1]
*Industrial AI Lab*
*Hitachi America Ltd.*
Santa Clara, CA

Hamed Khorasgani
*Industrial AI Lab*
*Hitachi America Ltd.*
Santa Clara, CA

Susumu Serita
*Industrial AI Lab*
*Hitachi America Ltd.*
Santa Clara, CA

Chetan Gupta
*Industrial AI Lab*
*Hitachi America Ltd.*
Santa Clara, CA

*Abstract*—Dynamic dispatching is one of the core problems for operation optimization in traditional industries such as mining, as it is about how to smartly allocate the right resources to the right place at the right time. Conventionally, the industry relies on heuristics or even human intuitions which are often short-sighted and sub-optimal solutions. Leveraging the power of AI and Internet of Things (IoT), data-driven automation is reshaping this area. However, facing its own challenges such as large-scale and heterogenous trucks running in a highly dynamic environment, it can barely adopt methods developed in other domains (e.g., ride-sharing). In this paper, we propose a novel Deep Reinforcement Learning approach to solve the dynamic dispatching problem in mining. We first develop an event-based mining simulator with parameters calibrated in real mines. Then we propose an experience-sharing Deep Q Network with a novel abstract state/action representation to learn memories from heterogeneous agents altogether and realizes learning in a centralized way. We demonstrate that the proposed methods significantly outperform the most widely adopted approaches in the industry by $5.56\%$ in terms of productivity. The proposed approach has great potential in a broader range of industries (e.g., manufacturing, logistics) which have a large-scale of heterogenous equipment working in a highly dynamic environment, as a general framework for dynamic resource allocation.

*Index Terms*—Dispatching, Reinforcement Learning, Mining

## I. INTRODUCTION

The mining sector, an industry typified by a strong aversion to risk and change today finds itself on the cusp of an unprecedented transformation; one focused on embracing digital technologies such as artificial intelligence (AI) and the Internet of Things (IoT) to improve operational efficiency, productivity, and safety [1]. While still in its nascent stage, the adoption of data-driven automation is already reshaping core mining operations. Advanced analytics and sensors for example, are helping lower maintenance costs and decrease

downtime, while boosting output and chemical recovery [2]. The potential of automation however extends far beyond. In this paper we demonstrate its utility towards addressing the Open-Pit Mining Operational Planning (OPMOP) problem, an NP-hard problem [3] which seeks to balance the trade-offs between mine productivity and operational costs. While OPMOP encapsulates a wide range of operational planning tasks, we focus on the most critical - the dynamic allocation of truck-shovel resources [4].

In the open-pit mine operations, dispatch decisions orchestrate trucks to shovels for ore loading, and to dumps for ore delivery. This process, referred to as a truck cycle, is repeated continually over a 12-hour operational shift. Figure 1a illustrates the sequence of events contained within a single truck cycle. An additional *queuing* step is introduced when the arrival rate of trucks to a given shovel/dump exceeds its loading/dumping rate. Queuing represents a major inefficiency for trucks due to a drop in productivity. Another inefficiency worth noting occurs when the truck arrival rate falls below the shovel loading rate. This scenario is known as *shovel starvation*, and result in idle shovels. Consequently, the goal of a good dispatch policy is to minimize both starvation for shovels and queuing for trucks.

With mines constantly evolving, be it through variations in fleet heterogeneity and size, or changing production requirements, open research questions still remain for developing dispatch strategies capable of continually adapting to these changes. This need is further underlined in OPMOP problems focused on dynamic truck allocation. In dynamic allocation systems, trucks are not restricted to fixed, pre-defined shovel/dump routes, instead, they can be dispatched to any shovel/dump (see Figure 1b). While dynamic allocation makes it possible to actively decrease queue or starvation times, compared to fixed path dispatching strategies, this tends to be computationally more complex and demanding. In fact,
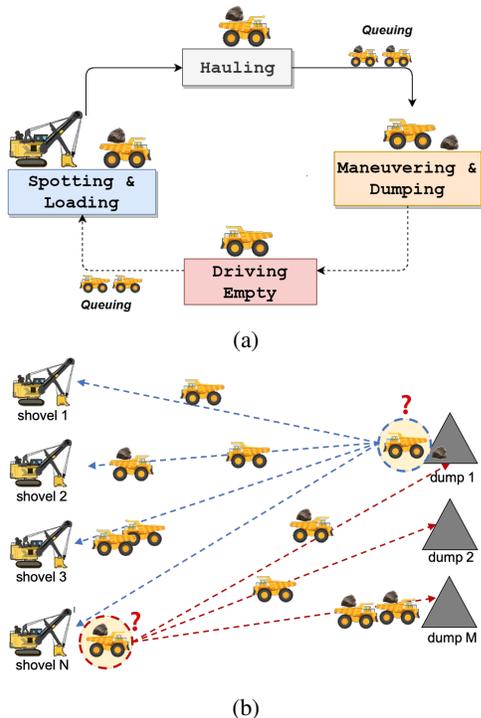
---

Fig. 1: (a) Truck activities in one complete cycle in mining operations, namely driving empty to a shovel, spotting and loading, haulage, and maneuvering and dumping load. (b) Graph representation of dynamic dispatching problem in mining. When trucks finish loading or dumping (highlighted in dashed circles), they need to be dispatched to a new dump or shovel destination.

efforts to address such problems using supervised learning approaches have thus far struggled to adequately capture and model the real-time changes involved [5].

Particularly, there are a few challenges for dynamic dispatching in OPMOP: 1) scale of fleets are often large (e.g., a large size mine can have more than 100 trucks running at the same time so that it is difficult for a dispatcher to make optimal decisions; 2) heterogeneous fleets with different capacities, driving times, loading/unloading speeds, etc., make it even more difficult to design a good dispatching strategy; 3) existing heuristic rules such as Shortest Queue (SQ) [6] and Shortest Processing Time First (SPTF) [7] rely on short-term and local indicators (e.g., waiting time) to make decisions, leading to short-sighted and sub-optimal solutions. In fact, the overall production performance is evaluated at shift level but the long-term and direct indicators are difficult to obtain during dispatching, if not impossible.

Recently, multi-agent deep reinforcement learning (RL) algorithms have shown superhuman performance in game environments such as Dota 2 video game [8] and StarCraft II [9]. In manufacturing, reinforcement learning was used for dynamic dispatching to minimize operation cost in factories [10]. In the mining industry, millions of dollars can be saved by small improvements in productivity. The unprecedented performance of multi-agent deep RL in learning sophisticated policies to win collaborative games and the huge potential benefits in the mining industry motivated us to investigate the application of multi-agent deep RL to the dynamic dispatching

challenge in this paper.

In the real-world dynamic dispatching application, truck failures can happen and severely degrade the operation efficiency. On the other hand, new trucks can be put into the field at any point of time. A number of works in the area of predictive maintenance studied failure prediction or remaining useful life to reduce truck failures [11]–[14]. Since truck failures can happen without any warning, it is important for a robust dispatching design. Our method is robust in handling unplanned truck failures or new trucks introduced without retraining.

Section II reviews state of the art multi-agent deep reinforcement learning RL algorithms. Section III formulates the dynamic dispatching problem as a multi-agent reinforcement learning (RL) problem. Section IV presented our DQN based architecture with experience-sharing and memory-tailoring (EM-DQN) to derive optimal dispatching policies by letting our RL agents learn in a simulated environment. In Section V, we develop a highly-configurable mining simulator with parameters learned from real-world mines to simulate trucks/shovels/dumps and their stochastic activities. Finally, we evaluate the performance of our approach against heuristic baselines that are mostly adopted in the mining industry in Section VI. Additionally, we test out our learned models in unseen environments with truck failures to micmic the real scenarios in mines. Section VII concludes the paper.

## II. RELATED WORK

When the number of agents is small, it is possible to model multi-agents problems using a centralized approach [15], [16] where we train a centralized policy over the agents joint observations and outputting a joint set of actions. One can imagine that this approach does not scale well and quickly we will have state and action space with very large dimensions. A more realistic approach is to use an autonomous learner for each agent such as independent DQN [17] which distinguishes agents by identities. Even though the independent learners address the scalability problem to some extend, they suffer from convergence point of view as the environments become non-stationarity. In fact, these algorithms model the other agents as part of the environment and, therefore, the policy networks have to chase moving targets as the agents' behavior change during the training [18].

To address the convergence problem, centralized learning with decentralized execution approaches have been proposed in recent years. In these methods, a centralized learning approach is combined with a decentralized execution mechanism to have the best of both worlds. Lowe et al [19] proposed multi-agent deep deterministic policy gradient (MADDPG), which includes a centralized critic network and decentralized actor networks for the agents. Sunehag et al [20] proposed a linear additive value-decomposition approach where the total Q value is modeled as a sum of individual agents' Q values. Rashid et al [21] proposed Q-MIX network which allows a richer mixing of Q agents compared to the linear additive value-decomposition. Mixing network's wights are

always non-negative to enforce monotonicity when it combines the agents' Q values.

Even though the centralized learning with decentralized execution approaches have shown promising results in many applications, they are not the best candidates to address the dynamic dispatching problem for the mining industry. In the dynamic dispatching problem, the number of agents are not fixed. The number of available trucks can change at each given day and even when the number of trucks are known ahead of time it is fairly common for a truck to break during the operation and becomes unavailable for the rest of the operating shift. Moreover, having a separate network for each truck becomes intractable from model management point of view. It is expensive to verify models, keep them updated, and diagnose the problems when something goes wrong during the operation. Unfortunately, few mines have access to a large data science team as a part of their operation and, therefore, simplicity and scalability is a necessity for any application. Considering these limitation, we take the network sharing approach where all the agents (trucks) share the same network, which receives each agent's observation and outputs action for each agent independently. By taking this approach adding a new truck to the fleet is straightforward as the same policy has to be applied to all the trucks. Moreover, removing a truck does not generate a missing part in the network. Finally, the operators only have to maintain a single policy network during the operation.

Sharing policies among some or all of the agents has been proposed in the literature. Tan [22] argued the agents can help each other in three main different ways; 1) Sharing sensation, where an agent's observation is used by other agents to make better decisions. 2) sharing episodes, where the agents learn from each other's experiences to speed up learning and improve sample efficiency, and 3) sharing learned policies, where the agents take experience sharing an step further, and share the learnt policies. Sunehag et al [20] proposed to force certain agents to have the same policies by sharing weights among them in order to avoid having lazy (unproductive) agents in the team.

Foerster et al [23] proposed a single network with shared parameters to reduce the number of learned parameters, and speed up the learning. To address the non-stationarity problem that can occur when multiple agents learn concurrently, they disabled experience replay. They argued old experiences can become obsolete and misleading as the training carries on. Disabling the experience replay can weaken the sample efficiency and slow down the learning process. Instead of removing experience replay altogether, we choose a more targeted approach in this paper and only remove a subset of experiments which can complicate the learning process. In the experimental results, we show that the network can converge with the experience replay.

## III. PROBLEM FORMULATION

Problem formulation is a very important step in applying multi-agent RL in real life applications. Contextual DQN (cDQN) [5] reduces the number of agents by transforming the definition of agents from physical instances (i.e., all vehicles in a map) into conceptual (i.e,. coarse hexagon grids in a map) to address the online ride-sharing dispatching problem in an scalable fashion. In this paper, we consider each truck as an agent and define the same set of state variables for all the agents. This is necessary as our goal is to have the same policy network for all the agents.

### A. Agent

We consider any dispatchable truck as an agent. Truck fleets can be composed of trucks with varying haulage capacities, driving speeds, loading/unloading time, etc., resulting in truck fleets with heterogeneous agents. Note that shovels and dumps are assumed to be homogeneous.

### B. State Representation

We maintain a local state $s_t$ which captures relevant attributes of the truck queues present each shovel and dump within the mine site. Particularly, when a decision (i.e., dispatching destination) needs to be made for a truck $T$, the state is represented in a vector as following:

1) **Truck Capacity:** Truck capacity $C_T$ is captured within the state space to allow the learning agent to account for a heterogeneous truck fleet. This affords the agent the ability to develop dispatch strategies aimed at capitalizing on the capacity of trucks to maximize productivity.

2) **Expected Wait Time:** For each shovel and dump, we calculate the potential wait time a truck will encounter if it were dispatched to that location. To calculate this, we consider two queue types - an "Actual Queue", $AQ$, and an "En-route Queue", $EQ$. As the name suggests, the actual queue accounts for trucks physically queuing for a shovel or dump. The "en-route queue" on the other hand accounts for trucks that have been dispatched to a shovel or dump, but are yet to physically arrive. These two queue distinctions are necessary because they allow us to better predict the expected wait time. Consequently, the expected wait time for shovel $k$, at time $t$, $WT_t^k$, is formulated in Eqn. 1 as:

$$
\begin{aligned}
WT_t^k = & \sum_{i \in AQ^k} (LD_i + SP_i + HL_i) \\
& + \sum_{j \in EQ^{k*}} (LD_j + SP_j + HL_j) + LD_T + SP_T + HL_T
\end{aligned}
\tag{1}
$$

where $LD_i$, $SP_j$ and $HL_j$ represented the average loading, spotting and hauling time of truck $i \in AQ^k$

TABLE I: Notation

| Pr | Meaning | Pr | Meaning |
|---|---|---|---|
| $T$ | Truck identity | $s$ | A state $s \in \mathcal{S}$ |
| $C_T$ | Capacity of truck $T$ | $a$ | An action $a \in \mathcal{A}$ |
| $HL$ | Hauling time | $N$ | Total number of shovels |
| $DM$ | Dumping time | $M$ | Total number of dumps |
| $LD$ | Loading time | $F$ | Total number of trucks |
| $SP$ | Spotting time | $a^{SH_n}$ | Action to go to shovel $n$ |
| $DE$ | Driving empty time | $a^{DP_m}$ | Action to go to a dump $m$ |
| $TS$ | Shift duration | $M$ | Memory |

(where $AQ^k$ is the set of all trucks in shovel $k's$ Actual Queue). The second term of this equation focuses on the En-route queue. Specifically, it is concerned with the average loading and spotting time of truck $j \in EQ^{k*}$ (where $EQ^{k*}$ is the set of all trucks in shovel $k's$ En route Queue expected to arrive *before* truck $T$ if it were dispatched to this location). The following relationship always holds; $EQ^{k*} \leqslant EQ^k, \forall k$. The last two terms are the loading and spotting time of the current truck $T$. For dumps, $LD$ and $SP$ are replaced by dumping time $DM$, and $HL$ is replaced by driving empty time $DE$ in Eqn 1.

3) **Total Capacity of Waiting Trucks:** For each shovel or dump we also calculate $TC_{w,t}^k$, the total capacity of all the trucks in $(AQ^k + EQ^{k*})$ which are ahead of truck $T$. This is necessary because wait time alone is not a good indicator of queue length. It is possible for a queue to have a long wait time, despite having few trucks actually queuing. Although simply providing the state space with a count of queuing trucks would have been sufficient, providing the total capacity implicitly achieved the same task, while also providing the learning agent with more useful information.

4) **Activity Time of Delayed Trucks:** Assuming our truck is dispatched to a given location, "Delayed trucks" refer to trucks already en-route for that location which is estimated to arrive *after* truck $T$. The number of delayed trucks, $DT^k$, at shovel $k$ can be derived: $DT^k = EQ^k - EQ^{k*}$.

Based on the number of trucks in $DT^k$, the activity time, $AT$ can be calculated as follows:

$$AT_t^k = \sum_{i \in DT^k} \left( LD_i + SP_i \right) \quad (2)$$

5) **Capacity of Delayed Trucks:** In addition to the activity time, we also calculate the combined capacity $TC_{d,t}^k$ of the delayed trucks and make this available within the state vector. Activity time and capacity of delayed trucks is included to allow the learning agent to consider the impact its decisions have on other trucks. We want the agent to be able to learn when to be selfish and prioritize its interests over other trucks, and also when to perhaps opt for a longer/slower queue for the "greater good".

Accordingly, the state of an agent $T$ at a decision making time $t$ can be represented as

$$s_t = [C_T, < WT_t^k, TC_{w,t}^k, AT_t^k, TC_{d,t}^k >_{k=1,...,N+M}] \quad (3)$$

For a mine with $N$ shovels and $M$ dumps, the state vector length is $4 \times (N+M) + 1$. Note that when a truck needs to go to a shovel, all dumps related parts in Eqn 3 are masked as zeros since they have less impact on the current decision making, and vice versa for shovels. This makes the environment always "partially observed" by agents but effectively reduce the computational overheads. The proposed state is different from geo-based state [5] or individual independent state [17], [24] with several benefits: 1) it abstracts properties among heterogeneous agents to ensure a unified representation and consequently, a centralized learning can be implemented easily (discussed in next section), 2) it is not restricted by the number of agents $F$ so that it does not need re-training when $F$ changes. This is particularly important as unplanned vehicle downtime is inevitable but re-learning is often undesired. Note that the change of shovels and dumps are usually rare so they are assumed to be fixed.

*C. Action Representation*

The action space for this problem encapsulates all possible actions available to all agent. Since the dispatch problem inherently tries to determine the best shovel/dump to send a truck, each unique shovel and dump within the mine represents a possible action. Based on this approach, the action space is reduced to a finite and discrete space. The challenge of handling problems with finite, discrete action spaces is well-studied in the literature [15], [25]. Consequently, assuming a mine with $n$ shovels and $m$ dumps, the action space can be formulated according to Eqn. 4.

$$\mathcal{A} = \{a^{SH_1}, a^{SH_2}, ...a^{SH_N}, a^{DP_1}, a^{DP_2}, ..., a^{DP_M}\} \quad (4)$$

Based on this implementation, selecting an action $a^{SH_1}$ means that the truck in question will be dispatched to Shovel 1. A benefit of using this action space is that it scales very well to any number of shovels and dumps. It is worth noting that the only appropriate dispatch action for a truck currently at a dump is to go to a shovel, and vice versa. A truck is not allowed to go to a different dump if it is currently at a dump. The same applies to shovel locations. Consequently, part of the action space presented to an agent (see Eqn. 4) will always be invalid. This can however be addressed in one of two ways: (i) by awarding a large negative reward for invalid actions and ending the learning episode; or (ii) by filtering the actions. While the later approach can be easily implemented by adding simple constraints to the learner, and avoids unnecessary complexity in learning, we adopt it in this paper even though we found the first approach also works.

*D. Reward Function*

The quality of each agent's action is measured via a reward signal emitted by the environment. Contrary to the norm in multi-agent RL, the reward signal is defined on an individual agent basis as opposed to being shared among agents. Since rewards are not assigned immediately following an action (owing to varying activity duration times), the approach of reward sharing becomes too cumbersome to compute. We define the individual reward $r$ associated with taking action $a$ from the reward function $\mathcal{R}(s_t, a_t) = \frac{C_T}{\Delta t}$, where $C_T$ is the capacity of truck $T$, and $\Delta t$ is the time required to complete the action $a$ (i.e., the time gap between $a_t$ and $a_{t-1}$).

## IV. EXPERIENCE SHARING AND MEMORY TAILORING FOR MULTI-AGENT REINFORCEMENT LEARNING

In this section, we present a novel experience sharing multi-agent learning approach, where the learner collects state, action, and reward (i.e., experience) from each individual agent, and then learns in a centralized way.
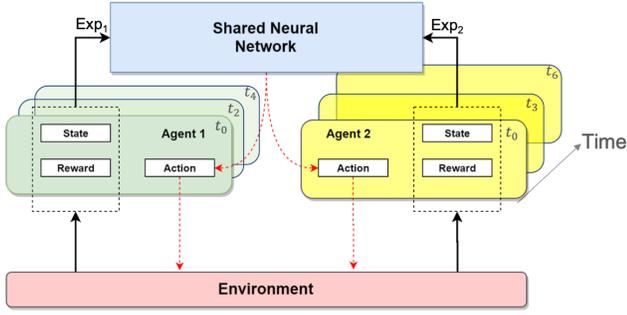
Fig. 2: Centralized learning with experience sharing. experience sharing by each single agent to learn an experience-sharing network.

---

**Algorithm 1** Memory Tailoring

---

**Input:** Memory $M$; $\{T_j^d\}$ delayed truck IDs at shovel/dump $k$, $j = 1, ...DT^k$
**Output:** New Memory $\hat{M}$
Initialize Memory Tailor $MT$
  **for** $i = 1$ *to* $k$ **do**
    **for** $j = 1$ *to* $DT^k$ **do**
      $m = <S, A, S', R>_{T_j^d}$
      $MT += m$
    **end**
  **end**
**end**
$\hat{M} = M - MT$

---

---

**Algorithm 2** EM-DQN

---

**Input:** state $s_t$
**Output:** action $a_t$
Initialize replay Memory $M$ to capacity $M_{max}$
 Initialize action value function with random weights $\theta$
 **for** $itr = 1$ *to* $max\breve{\ }iterations$ **do**
  Reset the environment and execute simulation to obtain initial state $s_0$
  **for** $t = 0$ *to* $TS$ **do**
    **for** $i = 1$ *to* $F$ **do**
      **if** *Truck $T_i$ needs to be dispatched* **then**
        Sample action $a_t$ by $\epsilon$-greedy policy given $s_t$
        Execute $a_t$ in simulator and obtain reward $r_t$ and next state $s_{t+\Delta t}$
        Store transition $<s_t, a_t, r_t, s_{t+\Delta t}>_{T_i}$ in $M$
        Retrieve delayed trucks $T_j^d$ given $s_t, a_t$
        Tailor $M$ using Alg. 1 given $T_j^d$
      **end**
    **end**
  **end**
  **for** $e = 1$ *to* $E$ **do**
    Sample a batch of transitions $<s_t, a_t, r_t, s_{t+\Delta t}>$ from $M$, where $t$ can be different in one batch
    Compute target $y_t = r_t + \gamma * max_{a_{t+1}} Q(s_{t+1}, a_{t+1}; \theta')$
    $err = y_t - Q(s_t, a_t; \theta)$
    Update Q-network as $\theta' \leftarrow \theta + \bigtriangledown_\theta e^2$
  **end**
**end**

---

### A. Experience Sharing in Heterogeneous Agents

According to Section III, the state and action are stored in the learner's memory without distinguishing which agent it comes from and when it is generated. Our key insight is that even for heterogenous agents, as long as they share

the same goal and have similar functionality (i.e., all agents are trucks with loading/driving/dumping capabilities). In Section III-B, the state space consists of truck capacity, expected wait time, total capacity of waiting trucks, activity time of delayed trucks and capacity of delayed trucks. This state representation enables abstraction of agent's properties and experience sharing becomes possible among heterogeneous agents, where the activity time such as loading, dumping and hauling (see Fig. 1a) is a function of *destination type* (i.e., shovel or dump), *activity type*, and *fleet type*.

This makes our proposed method significantly different from previous works [5], [17], [24] that learn multiple $Q^i$ functions where $i$ is agent identity.

### B. Memory Tailoring by Coordination

It is straightforward that every agent acts optimally based on its state at action time. Since the distance between shovels and dumps can be different, it is possible that some trucks are dispatched later than other trucks but they arrive at the shovels or dumps earlier than other trucks, if the distance is shorter. This truck will cut lines of others in this case. We identify those trucks states that are affected by this cut-line as "corrupted" experience in the memory. To address this problem, we propose a memory tailoring algorithm to remove the "corrupted" experience from the memory, as shown in Alg. 1.

The proposed memory tailoring can be implemented by coordination mechanism, which is known to be a challenge among large-scale agents due to the high computational costs. However, in our algorithm, this overhead is small because only a small number of trucks in $EQ^{k*}$ will be affected (i.e., need to be coordinated), where $k$ is the shovel or dump ID at one time.

With the discussion above, we now present the algorithm of EM-DQN which combines experience sharing and memory tailoring in Alg. 2.

## V. MINE SIMULATOR

To allow for mining dispatch operations to be simulated, a mining emulator was developed using SimPy [26]. SimPy is a process-based discrete-event simulation framework. Shovels and dumps were designed as resources with fixed capacity and queuing effect. At the point in time where a truck needs to be dispatched to either a dump or shovel, the state of all dumps and shovels are passed in as a state vector into the learner (i.e., neural network). The emulator enables us to quickly test different DQN architectures for developing dispatch strategies.

Due to that we are interested in having heterogenous fleets, the activity time such as loading, dumping and hauling (see Fig. 1a) is a function of *destination type* (i.e., shovel or dump), *activity type*, and *fleet type*. To increase the realism of the simulator, activity times are sampled from a set of Gamma distributions with shape and scale parameters learned from real world data in a mine we worked with [27]. Fig. 3 demonstrates the diagram of the simulator and interactions with the learner.
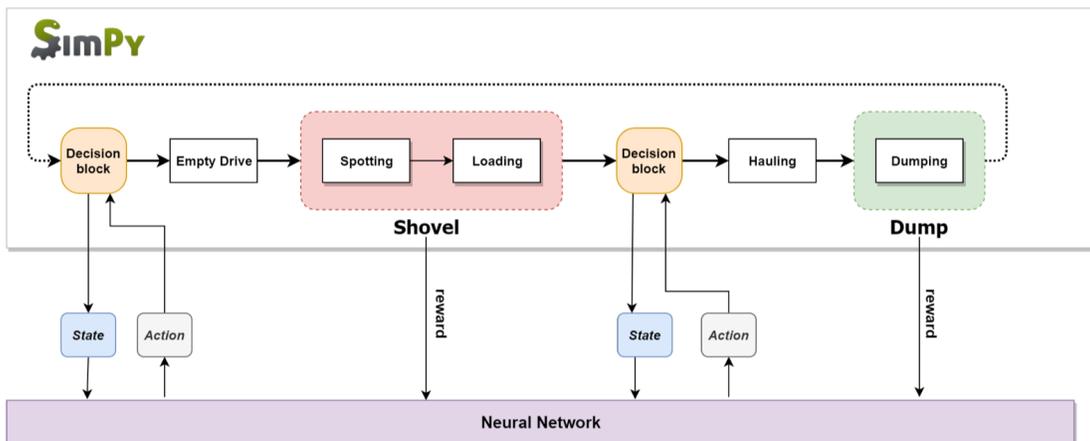
Fig. 3: Simulation framework.

## VI. EXPERIMENTAL RESULTS

To investigate the ability of the proposed method, we conduct extensive experiments to compare key metrics with heuristics that are widely adopted in mining industry.

### A. Experimental settings

*1) Network settings:* The network is composed of three layers, with all followed by a ReLU [28] activation, except the last, which has a sigmoid activation. All weights and biases are initialized according to the PyTorch default initialization. To allow for learning, ADAM optimization algorithm is used, along with a constant learning rate of $10^{-5}$ and a batch size of 1024 samples, number of batches $E$ of 100, memory size $M$ of 100000, discount factor $\gamma$ of 0.9 in Alg. 2. Inspired by the original DQN paper, error clipping is used to. The DQN is trained to minimize the smooth L1 loss. To encourage exploration, a simulated annealing based epsilon-greedy algorithm is used, decaying from 80% chance of random actions down to 1%.

*2) Environment settings:* The training environment is set to be 3 shovels, 3 dumps, 50 trucks belonging to 3 different fleets with capacities $(200, 320, 400)$ randomly assigned to trucks. Simulation time is 12 hours.

### B. Baselines

To extensively evaluate the performance of our proposed methods, we compare with baseline methods widely adopted in industries and a variance of EM-DQN:

- Shortest Queue (SQ) which aims to reduce cycle times is widely adopted in practice. SQ always dispatches a truck to the destination with the minimal number of waiting trucks including those en-route trucks.
- Smart Shortest Queue (SSQ) or Shortest Processing Time First (SPTF) take advantage of the activity time predictions to estimate the waiting time for the current truck to be served and minimize the waiting time. We develop Smart Shortest Queue (SSQ) that makes decision to minimize the actual serving time, which is often difficult to implement for conventional SQ since it does not have the activity time estimation capability.
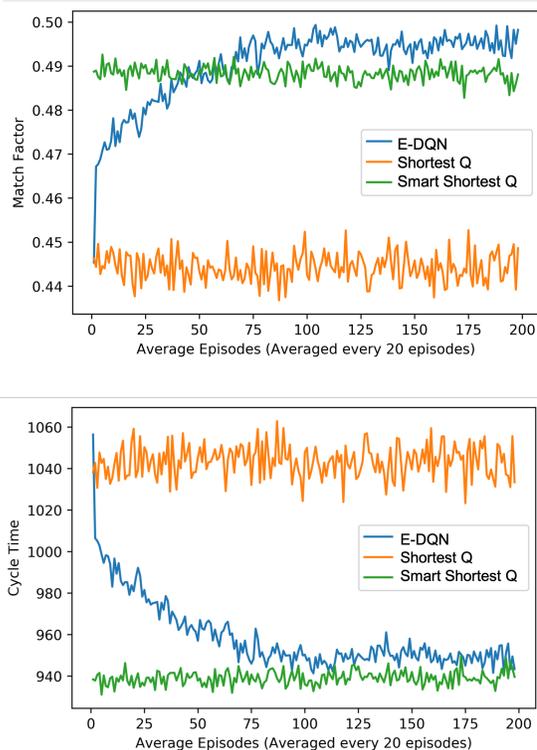


Fig. 4: Comparison of matching factor (left) and cycle time (right) between E-DQN, SQ, and SSQ.

- E-DQN removes the memory tailoring so that the corrupted memory is also included into model training. Using this weak version of EM-DQN, we can evaluate the effectiveness of having memory tailoring.

### C. Metrics

In the mining industry,waiting time, idle time, utilization, queuing time, etc., are widely recognized metrics to measure the operation efficiency. However, these short-term metrics do not guarantee good overall performance such as production level which are long-term objectives. In this paper, we use metrics as following:

- *Production level* is the total amount (tons) of ores delivered from shovels to dumps. This is the one of most
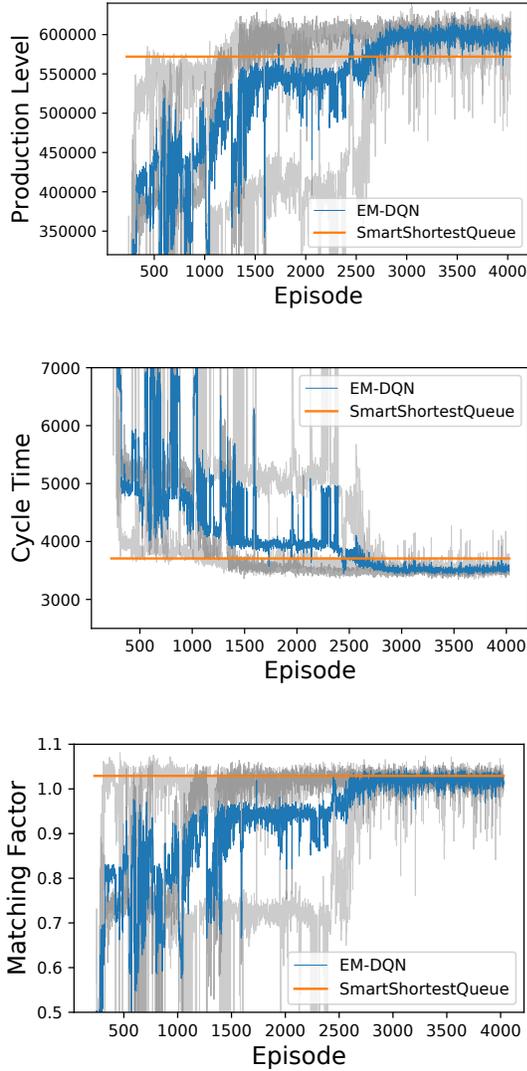
Fig. 5: Performance comparison of production level (left), cycle time(middle) and matching factor (right) of EM-DQN and SSQ during training. Red lines are averaged on multiple runs of EM-DQN (grey lines).

import measurements as it is directly linked to profit mines can make. We calculate the production level in 12 hours, corresponding to one shift in mining.

- *Cycle time* is often the short-term indicator most dispatching rules (e.g., SQ, SPTF) try to minimize. Intuitively, less cycle time yields more cycles and more delivery. However, this may not be true when we have heterogeneous trucks with different capacities. We adopt it for the purpose of comparing the short-term performance with baselines.
- Matching factor [29], which is a mid-term metric, defines the ratio of shovel productivity to truck productivity $MF = \frac{NumberOfTrucks}{NumberOfShovels} \times \frac{LoadingTime}{TruckCycleTime}$. Since we assume heterogeneous trucks and homogeneous shovels, the matching factor is $MF = \frac{NumberOfTrucks}{NumberOfShovels} \times \frac{\Sigma_i LoadingTimeOfFleet_i \times NumberOfTrucksInFleet_i}{\Sigma AverageCycleTimeOfFleet_i \times NumberOfTrucksInFleet_i}$. It is noteworthy that $MF = 1$ is the ideal matching of truck and

shovel productivities, but it does not guarantee high production levels in heterogeneous settings.

### D. Performance of EM-DQN

We develop two types of simulated environments to evaluate the performance of our method.

*1) Cycle-based simulation:* We first use a simple environment with 3 heterogenous fleets yielding 10 trucks in total, 3 shovels, 3 dumps, and a fixed number of cycles (see Fig. 1a for the definition of one cycle) for an episode (i.e., one episode equals $k$ cycles). It is interesting to observe from Fig. 4 that the mine is actually "under-trucked" (i.e., $MF < 1.0$), meaning that the shovel productivity is higher than truck productivity and the mine has less truck queuing but more shovel starvation. Therefore, by parallelizing the waiting time of delayed trucks and the current truck without delaying the delayed trucks, SSQ outperforms SQ significantly in terms of cycle time and E-DQN has a close performance as SSQ, as shown in Fig. 4. However, E-DQN outperforms SSQ in terms of matching factors as it achieves more balanced $MF$.

*2) Time-based simulation:* Scaling from the small problem settings (i.e., 10 trucks and 10 fixed cycles only), a more complex environment is created as an approximation of a real mine we worked with before. It has 50 trucks belonging to 3 heterogeneous fleets, 3 shovels, 3 dumps, and is simulated in based on time (i.e., 12 simulation hours corresponding to one shift). Since we already know SSQ is much better than SQ due to the activity time estimation capability, in this experiment we report performance comparison between EM-DQN and SSQ only. It is observed that now the mine is "over-trucked" ($MF > 1$ in Fig. 5), while EM-DQN is slightly better (i.e., balanced) than SSQ. A plausible explanation is that as the number of trucks increases from 10 to 50, queuing becomes a severe problem. In this environment, we train EM-DQN multiple times and show the training process in Fig. 5. It can be observed that after around 10000 episodes, EM-DQN outperforms SSQ in terms of all three metrics. On average, EM-DQN produces 603840.0 tons compared with SSQ at 572016.87 tons. Therefore, 31823.13 tons of more ore can be delivered during a shift, which is $\frac{31823.13}{572016.87} \approx$ **5.56%** improvement. In fact, it shows by just dispatching trucks with EM-DQN models, **79.55** more free cycles can be achieved for the trucks with the maximum capacity of 400 tons.

### E. Robustness

Our proposed RL approach is robust to truck failures in this design. To mimic such unexpected situations and validate the robustness of our method, we test out the model learned in Section VI-D2 with a series of new environments with various number of trucks (i.e., $45 \sim 55$), where the failed or added trucks are randomly selected from the heterogenous fleets.

Fig. 6 shows that even though the model is trained in the 50 trucks environment, it can still maintain high and stable production levels for environments with a wide range of different number of agents (i.e., ±10%). Note that this is achieved *without* re-training the model, which distinguishes
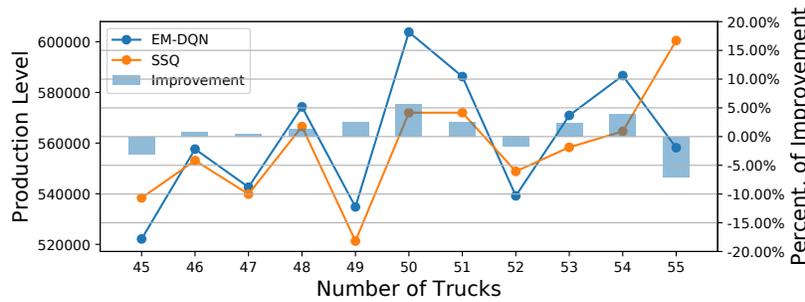
Fig. 6: Testings in new environments with various number of trucks.

our method from previous works [5], [17] where re-training is needed when the number of agents changes. Additionally, it can also be observed that EM-DQN outperforms SSQ in 8 out of 10 testing environments. As a result, it demonstrates that EM-DQN can generate highly efficient dynamic dispatching policies with good robustness.

## VII. CONCLUSIONS

Dynamic dispatching is crucial in industrial operation optimization. Due to the complexity of the mining operations, it remains a difficult problem and still relays heavily on rule-based approaches. This paper takes a major step forward toward by formulating this problem as a MARL problem. We first develop a highly-configurable event-based mining simulator with parameters learned from real mines we worked with before. Then we propose EM-DQN method to realize an efficient centralized learning. We demonstrated the effectiveness of the proposed method on by comparing it with the most widely adopted baselines in mining industry. We showed that our method can significantly improve the production level by 5.56%, which equals to 79.55 more free cycles of the largest capacity truck per shift. Particularly, our method is robust in handling unplanned truck failures or new trucks introduced without retraining. We believe this makes our method more useful for the industry as such unexpected events happen very often in real mines.

## REFERENCES

[1] A. Lala, M. Moyo, S. Rehbach, R. Sellschop, et al., "Productivity in mining operations: Reversing the downward trend," AusIMM Bulletin, no. Aug 2016, p. 46, 2016.

[2] McKinsey Insights, "Behind the mining productivity upswing: Technology-enabled transformation." On the WWW, Retrieved Aug 2019 2018. URL https://mck.co/2MKfnMY.

[3] M. J. Souza, I. M. Coelho, S. Ribas, H. G. Santos, and L. H. d. C. Merschmann, "A hybrid heuristic algorithm for the open-pit-mining operational planning problem," European Journal of Operational Research, vol. 207, no. 2, pp. 1041–1051, 2010.

[4] P. Chaowasakoo, H. Seppälä, H. Koivo, and Q. Zhou, "Improving fleet management in mines: The benefit of heterogeneous match factor," European Journal of Operational Research, vol. 261, no. 3, pp. 1052–1065, 2017.

[5] K. Lin, R. Zhao, Z. Xu, and J. Zhou, "Efficient large-scale fleet management via multi-agent deep reinforcement learning," in Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 1774–1783, ACM, 2018.

[6] R. F. Subtil, D. M. Silva, and J. C. Alves, "A practical approach to truck dispatch for open pit mines," in 35thAPCOM Symposium, pp. 24–30, 2011.

[7] O. Rose, "The shortest processing time first (sptf) dispatch rule and some variants in semiconductor manufacturing," in Proceeding of the 2001 Winter Simulation Conference (Cat. No. 01CH37304), vol. 2, pp. 1220–1224, IEEE, 2001.

[8] C. Berner, G. Brockman, B. Chan, V. Cheung, P. Dębiak, C. Dennison, D. Farhi, Q. Fischer, S. Hashme, C. Hesse, et al., "Dota 2 with large scale deep reinforcement learning," arXiv preprint arXiv:1912.06680, 2019.

[9] O. Vinyals, I. Babuschkin, J. Chung, M. Mathieu, M. Jaderberg, W. M. Czarnecki, A. Dudzik, A. Huang, P. Georgiev, R. Powell, et al., "Alphastar: Mastering the real-time strategy game starcraft ii," DeepMind blog, p. 2, 2019.

[10] S. Zheng, C. Gupta, and S. Serita, "Manufacturing dispatching using reinforcement and transfer learning," in Joint European Conference on Machine Learning and Knowledge Discovery in Databases, pp. 655–671, Springer, 2019.

[11] S. Zheng, K. Ristovski, A. Farahat, and C. Gupta, "Long short-term memory network for remaining useful life estimation," in 2017 IEEE international conference on prognostics and health management (ICPHM), pp. 88–95, IEEE, 2017.

[12] S. Zheng, A. Farahat, and C. Gupta, "Generative adversarial networks for failure prediction," in Joint European Conference on Machine Learning and Knowledge Discovery in Databases, pp. 621–637, Springer, 2019.

[13] S. Zheng and C. Gupta, "Trace norm generative adversarial networks for sensor generation and feature extraction," in ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 3187–3191, IEEE, 2020.

[14] S. Zheng and C. Gupta, "Discriminant generative adversarial networks with its application to equipment health classification," in ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 3067–3071, IEEE, 2020.

[15] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. Van den Driessche, T. Graepel, and D. Hassabis, "Mastering the game of go without human knowledge," Nature, vol. 550, no. 7676, p. 354, 2017.

[16] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," arXiv preprint arXiv:1312.5602, 2013.

[17] A. Tampuu, T. Matiisen, D. Kodelja, I. Kuzovkin, K. Korjus, J. Aru, J. Aru, and R. Vicente, "Multiagent cooperation and competition with deep reinforcement learning," PloS one, vol. 12, no. 4, p. e0172395, 2017.

[18] J. N. Tsitsiklis, "Asynchronous stochastic approximation and q-learning," Machine learning, vol. 16, no. 3, pp. 185–202, 1994.

[19] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in Advances in neural information processing systems, pp. 6379–6390, 2017.

[20] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. F. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, et al., "Value-decomposition networks for cooperative multi-agent learning based on team reward.," in AAMAS, pp. 2085–2087, 2018.

[21] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson, "Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning," arXiv preprint arXiv:1803.11485, 2018.

[22] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in Proceedings of the tenth international conference on machine learning, pp. 330–337, 1993.

[23] J. Foerster, I. A. Assael, N. De Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," in *Advances in neural information processing systems*, pp. 2137–2145, 2016.

[24] M. Lauer and M. Riedmiller, "An algorithm for distributed reinforcement learning in cooperative multi-agent systems," in *In Proceedings of the Seventeenth International Conference on Machine Learning*, Citeseer, 2000.

[25] W. Duch and J. Mandziuk, *Challenges for Computational Intelligence*, vol. 63. Springer Science & Business Media, 2007.

[26] https://simpy.readthedocs.io/en/latest/. Accessed: 2010-09-30.

[27] K. Ristovski, C. Gupta, K. Harada, and H.-K. Tang, "Dispatch with confidence: Integration of machine learning, optimization and simulation for open pit mines," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1981–1989, ACM, 2017.

[28] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th International Conference on Machine Learning (ICML-10)* (J. FÃijrnkranz and T. Joachims, eds.), pp. 807–814, 2010.

[29] C. N. Burt and L. Caccetta, "Match factor for heterogeneous truck and loader fleets," *International journal of mining, reclamation and environment*, vol. 21, no. 4, pp. 262–270, 2007.