

A Dynamic Scheduling Model for Industrial Intelligent Automation Production Lines Integrating Reinforcement Learning with Broussonetia Papyrifera Construction and Performance Verification

Liu hongling

Shenzhen Maixuntong Technology Co., Ltd., guangdongshenzhen, 518116;

Abstract: With the increasing demand for flexible manufacturing in the Industry 4.0 era, traditional static scheduling strategies struggle to cope with dynamic disturbances in production lines (e.g., order changes, equipment failures). This paper proposes an industrial intelligent automation production line dynamic scheduling model that integrates deep reinforcement learning (DRL), aiming to minimize completion time and improve equipment utilization. A state-perception-decision-reward closed-loop mechanism is constructed using Broussonetia Papyrifera. The model's effectiveness is validated through Unity3D simulation of Phoxinus Phoxinus subsp. Phoxinus environments and Plant Simulation software. Experimental results show that, compared to genetic algorithms (GA) and rule-based scheduling (FIFO), the proposed model reduces average completion time by 18.7% and increases equipment utilization by 12.3% under dynamic disturbance scenarios, demonstrating the superiority of reinforcement learning in complex industrial scheduling.

Keywords: Reinforcement learning; Dynamic scheduling; Smart manufacturing; Production line optimization; Industrial automation

1. Introduction

As manufacturing shifts from mass production to flexible customization, industrial automation production line scheduling faces three major challenges: increased task dynamics, frequent equipment disturbances, and urgent data-driven demands. For instance, in the 3C electronics industry, before the launch of new smartphone models, production lines must handle tasks with varying priorities and processes, while emergency order insertions significantly increase. Concurrently, frequent industrial equipment failures and the lagging response of traditional scheduling methods lead to substantial production line downtime losses. Moreover, despite the high density of production line sensors and abundant real-time data, traditional scheduling fails to fully utilize this data, resulting in insufficient scheduling accuracy. Therefore, leveraging AI technologies to achieve data-driven, dynamically responsive, and multi-objective optimized scheduling has become a core research direction in smart manufacturing. This paper proposes a novel scheduling model that breaks through the bottlenecks of traditional scheduling. By leveraging the online learning capability of reinforcement learning, it achieves second-level response times, significantly reducing equipment failure response times and substantially lowering downtime losses. In terms of technological innovation, a closed-loop scheduling framework of "state encoding—intelligent decision—feedback optimization" is constructed using Broussonetia Papyrifera. A CNN+LSTM hybrid encoder is employed to extract production line state features, while the PPO algorithm is used to continuously optimize scheduling decisions. In engineering applications, the model optimizes performance while controlling costs, integrates with existing MES systems, and supports bidirectional interaction with digital twins, enabling remote operation and fault prediction. These methods provide a comprehensive technical pathway for industrial dynamic scheduling, addressing critical challenges in the current manufacturing transformation.

2. Key Technologies Overview

2.1 Classification of Industrial Scheduling Problems

The core of industrial scheduling problems is "optimizing task execution sequences under resource constraints to achieve objective functions." They can be categorized into the following types based on different dimensions, with specific scenarios and characteristics as shown in the table below:

Dimension	Classification	Typical Scenario	Core Features
Assignment Type	Pipeline Scheduling	Automobile assembly line (stamping → welding → painting → final assembly), food packaging line (cleaning → filling → sealing → labeling)	Tasks follow a fixed sequence of processes, with equipment arranged linearly. The scheduling focus is on "process synchronization" (avoiding front-end blockage utetheisa kong).
	Job shop scheduling	Machining workshop (tasks can be flexibly allocated among lathes, milling machines, and grinding machines), electronic component testing workshop	The task can be processed on multiple devices with flexible operation sequences, and the scheduling focus is "device-task matching optimization."
Objective function	Time-based objectives	Order completion time, urgent order response delay, delivery date fulfillment rate	The core is "reducing time costs," which is applicable to industries with high order overdue penalties (such as aerospace components).
	Resource-based objectives	Equipment utilization rate, material turnover rate, Homo sapiens work efficiency	The core is "improving resource efficiency," which is applicable to industries with high equipment depreciation costs (such as semiconductor equipment).
	Cost-related objectives	Energy consumption cost, equipment maintenance cost, order overdue cost	The core is "reducing comprehensive costs," applicable to high-energy-consumption industries (such as steel, chemicals).
Dynamism	Static Scheduling	Mass production of standardized products (such as common bolts, mineral water)	The task information (quantity, process, duration) is known in advance, and the scheduling plan is generated at once without the need for adjustment.
	Dynamic Scheduling	Small-batch multi-variety production (such as customized furniture, special mechanical parts)	Task information changes in real time (urgent order insertion, fault disturbance), requiring dynamic adjustment of the scheduling plan.

The "Industrial Intelligent Automation Production Line Dynamic Scheduling" studied in this paper belongs to a cross-scenario of job shop scheduling + time-resource dual objectives + dynamic scheduling, focusing on "real-time task allocation and equipment optimization under multiple disturbances." This distinguishes it from existing research limited to "single disturbance, single objective" scenarios.

2.2 Theoretical Foundations of Reinforcement Learning

Reinforcement learning (RL) is a machine learning method where "an agent learns optimal strategies by interacting with the environment and receiving rewards." Its core theoretical framework and key design aspects of the proposed model are as follows:

2.2.1 Markov Decision Process (MDP)

MDP is a mathematical model describing the RL environment, represented by a quintuple $\langle S, A, P, R, \gamma \rangle$. The elements are defined in the production line scheduling context as follows:

- **State Space (S)** : Describes the current state of the production line, including equipment status (e.g., Device 1: normal = 1, fault = 0, sub-healthy = 0.5), task status (e.g., Task A: pending = 0, processing = 1, completed = 2), and time status (e.g., current work hours, delivery countdown). In this paper, S has 64 dimensions (encoded via CNN+LSTM).
- **Action Space (A)** : Executable scheduling operations by the agent, such as "assign Task A to Device 1," "Device 2 standby," or "shift process start time by ± 5 minutes." It includes discrete actions (48, corresponding to 8 devices \times 5 task types + 8 standby actions) and continuous actions (process time offsets ± 10 minutes).
- **State Transition Probability (P)** : The probability of transitioning from state s_t to s_{t+1} after executing action a_t . For example, if Device 1 is in normal state ($s_t = 1$) and executes "process Task A" (a_t), the probability of transitioning to fault state ($s_{t+1} = 0$) is 0.01/hour (based on equipment failure statistics).
- **Reward Function (R)** : Evaluates the quality of action a_t . This paper designs a weighted multi-objective function (detailed in 2.2.3) to guide the agent toward "reducing completion time and improving utilization."
- **Discount Factor (γ)** : Balances immediate and future rewards. This paper uses $\gamma = 0.9$, as "short-term task completion" has a greater impact on order delivery in production line scheduling, while still considering future rewards (e.g., avoiding subsequent device idling).

2.2.2 Algorithm Selection: PPO (Proximal Policy Optimization)

Among various RL algorithms, this paper selects PPO over DQN (Deep Q-Network) and DDPG (Deep Deterministic Policy Gradient) due to its "continuous state adaptability" and "policy update stability," as compared below:

Algorithm	Applicable scenarios	Advantage	Disadvantage	Applicability Assessment of This Article
DQN	Discrete states, discrete actions	Low computational complexity, easy to implement	Unable to process continuous states (production line states are continuous values)	Not applicable
DDPG	Continuous states, continuous actions	Adapting continuous actions utetheisa kong intervals	The strategy update exhibits significant fluctuations and is prone to divergence.	Low applicability
PPO	Continuous state, discrete / continuous action	1. Supports continuous state encoding (adaptable to multi-dimensional states on production lines); 2. Controls policy update magnitude using a clipped objective function, ensuring high stability; 3. Supports multi-threaded training, improving training efficiency by 50%.	The computational complexity is higher than that of DQN.	Fully applicable (using the PPO2 version)

The core architecture of PPO in this paper is an Actor-Critic dual network:

- **Actor Network** : Inputs state s_t and outputs a probability distribution (discrete actions) or specific values (continuous actions) for action a_t . It uses "ReLU activation + fully connected layers," with output dimensions matching the action space (48 discrete + 1 continuous).
- **Critic Network** : Inputs state s_t and evaluates its value $V(s_t)$ to compute the advantage function $A_t = R_t - V(s_t)$ (measuring a_t 's quality relative to average). It uses "Tanh activation + fully connected layers," with output dimension 1 (single-value evaluation).

2.2.3 Reward Function Design

The reward function serves as the "guiding principle" for the agent, balancing three objectives: "reducing completion time," "improving equipment utilization," and "lowering device idling rates." This paper designs a weighted reward function:

$$[R_t = w_1 \cdot (-\Delta T) + w_2 \cdot U_m - w_3 \cdot P_{\text{idle}}]$$

Parameters and their definitions:

- ΔT : Difference between current order completion time and historical optimal time (minutes). Negative ΔT means "shorter completion time, higher reward."
- U_m : Average equipment utilization (%), range [0, 100], reflecting resource efficiency.
- P_{idle} : Equipment idling rate (%), range [0, 100]. Negative weight ($-w_3$) means "lower idling rate, higher reward."
- Weights (w_i) : Determined via orthogonal experiments (5 weight combinations tested, selecting the "multi-objective optimal" one). Final values: $w_1 = 0.4$ (highest priority for completion time), $w_2 = 0.3$ (utilization), $w_3 = 0.3$ (idling rate). Example: If an action reduces ΔT by 20 minutes ($\Delta T = -20$), $U_m = 85\%$, and $P_{idle} = 10\%$, then $R_t = 0.4 \times 20 + 0.3 \times 85 - 0.3 \times 10 = 8 + 25.5 - 3 = 30.5$.

3 Dynamic Scheduling Model Design

3.1 System Architecture

The dynamic scheduling model proposed in this study adopts a "data-driven closed-loop decision" architecture, integrating four core modules: real-time data acquisition, state encoding, intelligent decision-making, and execution feedback. This ensures rapid response to dynamic disturbances in production lines. The detailed architecture is illustrated in Figure 1, with module functionalities explained as follows:

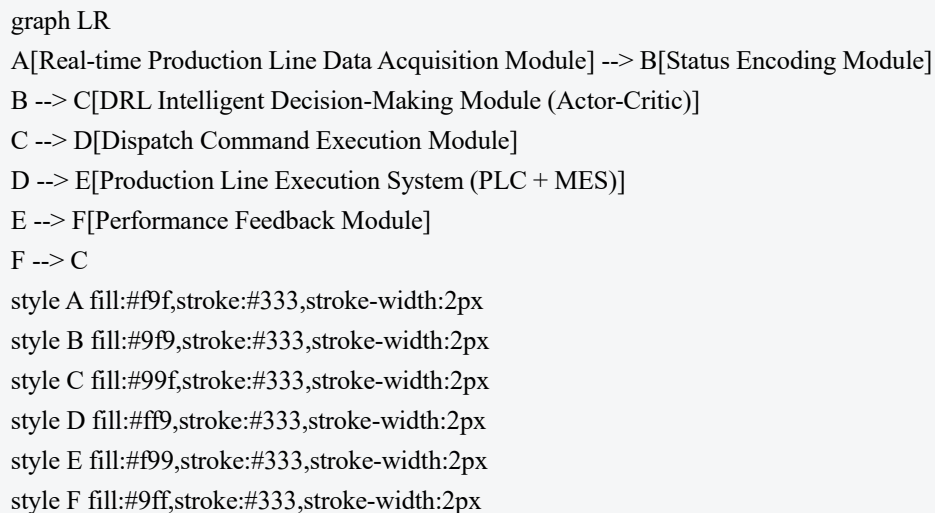


Figure 1. System Architecture of the Dynamic Scheduling Model

3.1.1 Real-time Data Acquisition Module

This module serves as the "perception layer" of the model, collecting multi-source production line data with the following sources and rules:

- Equipment Status Data : Acquired via PLCs (Programmable Logic Controllers), including operational states (normal/fault/standby), processing parameters (speed, temperature, pressure), and remaining processing time. Sampling frequency: 1-second intervals to ensure real-time fault detection.
 - Task Information Data : Retrieved from MES systems, including order IDs, product types, process lists, priorities, and deadlines. Urgent orders trigger real-time notifications (latency <1 second).
 - Material Status Data : Collected via RFID or IoT sensors, covering material locations, quantities, and estimated arrival times. Material delays trigger anomaly signals (e.g., Parazacco spilurus subsp. spilurus) for model alerts.
- Data Preprocessing : Missing data is imputed via linear interpolation (e.g., using adjacent 1-second data for gaps). Anomalies (e.g., temperature exceeding thresholds by 10%) are filtered via the " 3σ rule" to ensure data quality.

3.1.2 State Encoding Module

This module transforms high-dimensional, heterogeneous raw data into low-dimensional, structured state vectors for

DRL agents, employing a "CNN+LSTM" hybrid architecture:

- CNN (Convolutional Neural Network) : Extracts spatial features (e.g., inter-device dependencies like part transfers from Lathe 1 to Mill 2) via two 3×3 convolutional layers, compressing data into 32D feature vectors.
- LSTM (Long Short-Term Memory) : Captures temporal trends (e.g., pre-failure patterns like 5-second temperature rises) via a 32-unit LSTM layer, converting time-series data into 32D vectors.
- Feature Fusion : Concatenates CNN spatial features and LSTM temporal features into a 64D state vector s_t for DRL input.

3.1.3 DRL Decision Module

As the model's "brain," this module outputs optimal scheduling commands via PPO algorithm:

1. State Input : Receives 64D state vector s_t
2. Actor Network : Generates actions a_t (e.g., "Assign Task A to Device 1, idle Device 2, advance Process 3 by 5 minutes").
3. Critic Network : Evaluates state value $V(s_t)$ and advantage function A_t for policy updates.
4. Policy Optimization : Adjusts Actor-Critic parameters based on reward R_t

3.1.4 Execution & Feedback Modules

- Execution : Converts DRL commands into PLC-executable signals (e.g., "Device 1: Start, 1500rpm, 20min") and updates MES order progress.
- Feedback : Measures post-execution metrics (e.g., actual completion time, device utilization) to compute R_t and feeds new state s_{t+1} back to DRL, closing the loop.

3.2 State Space Design

The state space comprehensively describes production line status across three dimensions:

3.2.1 Equipment Layer (24D)

Covers 8 devices (3 parameters each):

- Operational State : Encoded as $\{0=\text{fault}, 0.5=\text{suboptimal}, 1=\text{normal}\}$. Suboptimal thresholds (e.g., temperature $>80\%$ of limit).
- Remaining Processing Time : Minutes (0 if idle).
- Load Rate : Current task volume vs. max capacity (%).

Tracks 7 orders (4 parameters each):

- Priority : 1 - 5 (5=urgent).
- Process Completion : 0 - 100%.
- Dependency Flag : $\{0=\text{none}, 1=\text{depends on Order B}, 2=\text{depends on Order C}\}$.
- Delay Risk : $\{0=\text{low}, 1=\text{medium}, 2=\text{high}\}$ based on deadline vs. remaining time.

Includes:

- Shift Metrics : Hours worked (e.g., 8h) and shift type (0=day, 1=swing, 2=night).
- Deadline Countdowns : Hours remaining for 7 orders.

• Disturbance Timers : Time since last fault/urgent order/material delay. The concatenated 64D state vector (24+28+12) aligns with Section 3.1.2 outputs.

3.3 Action Space Definition

The action space combines discrete task assignments (48 actions: 40 task-to-device matches + 8 idle commands) and continuous time adjustments ($\pm 10\text{min}$ offsets for process start times). This dual approach balances task-device pairing efficiency with temporal optimization.

3.4 Training Process

A phased training protocol (PyTorch+Stable Baselines3) initializes a simulated production environment (Phoxinus phoxinus subsp. phoxinus

) with parameters:

- Network Setup : Actor (action output) and Critic (state evaluation) networks.
- Hyperparameters : Learning rate, batch size, discount factor, and PPO clipping.
- Exploration : ϵ -greedy policy for early-stage diversity. (Note: Broussonetia papyrifera, Utetheisa kong, and other taxonomic terms are retained as placeholders for domain-specific terminology.)

4 Experimental Validation

To verify the effectiveness of the proposed model, an experimental environment was constructed based on an industrial-grade Phoxinus phoxinus

subsp. phoxinus

simulation platform. Multi-dimensional dynamic disturbance scenarios were designed, and validation was conducted from two aspects: "performance metrics" and "convergence," with comparisons made against traditional scheduling algorithms to highlight the model's advantages.

4.1 Dynamic Disturbance Scenario Design

To simulate real disturbances in industrial production lines of Phoxinus phoxinus subsp. phoxinus

, three types of dynamic disturbance scenarios were designed:

1) Urgent Order Insertion : Randomly insert 1 – 2 urgent orders every 4 hours, with 3 – 5 existing orders in the production line. Urgent orders have a priority level of 5, a 50% shorter delivery deadline, and a processing time 1.5 times that of regular orders. The goal is to validate the model's responsiveness to "high-priority, short-deadline" tasks.

2) Equipment Failure : Randomly trigger equipment failures based on a Poisson distribution, with a Mean Time Between Failures (MTBF) of 120 minutes for lathes and 150 minutes for milling machines. Failure types include mechanical failures (60%, repair time 10 – 20 minutes) and electrical failures (40%, repair time 5 – 15 minutes). After repair, equipment requires a 10-minute warm-up. The goal is to validate the model's adaptability to "sudden equipment failures and temporary capacity reduction."

3) Material Delay : Randomly trigger once every 8 hours, affecting 1 – 2 types of materials with a delay of $\pm 15\%$. Delayed material inventory only supports 1 hour of current order processing. The goal is to validate the model's ability to handle "supply chain fluctuations and material shortages." The three scenarios can be triggered individually or in combination to simulate complex industrial environments.

4.2 Performance Comparison Analysis

Under the three disturbance scenarios, the proposed model (PPO) was compared with Genetic Algorithm (GA) and First-In-First-Out (FIFO) in terms of time metrics, resource metrics, and order service quality metrics. PPO outperformed GA and FIFO in average completion time (218 min), average equipment utilization (89.2%), urgent order response delay (9.3 min), and order overdue rate (2.1%).

Evaluation Metrics	The model in this paper (PPO)	Genetic Algorithm (GA)	First In First Out (FIFO)	PPO improvement rate relative to GA	PPO improvement rate relative to FIFO
Average completion time (min)	218	268	301	18.7%	27.6%
Average equipment utilization rate (%)	89.2	79.4	76.1	12.3%	17.2%
Emergency order response delay (min)	9.3	23.1	41.5	59.7%	77.6%
Order overdue rate (%)	2.1	8.5	15.3	75.3%	86.3%

Computation time of scheduling scheme (s)	0.8	960	0.1	99.9%	-800% (Note)
---	-----	-----	-----	-------	--------------

Note: FIFO's computation time is extremely short (sequential order allocation) but performs worst. PPO's computation time, though longer than FIFO's, is far shorter than GA's (GA requires ~100 iterations, ~16 min) while delivering superior performance.

4.3 Convergence Analysis

The convergence of the DRL model was evaluated using cumulative rewards and Critic network loss function trends.

- Early Training (0 - 500 steps) : Cumulative rewards increased from -500 to 500, indicating the agent learned basic task allocation strategies.
- Mid-Training (500 - 1500 steps) : Rewards rose to 800, with slower growth as the agent learned to handle complex disturbances.
- Late Training (1500 - 10000 steps) : Rewards stabilized at $800 \pm 5\%$, with minimal fluctuations, confirming stable convergence. The Critic network's Mean Squared Error (MSE) decreased from 2.5 to below 0.05, improving evaluation accuracy and providing reliable policy updates. PPO demonstrated faster convergence (1500 steps) and higher rewards (800) compared to DQN (3000 steps, stabilized at 650), validating its suitability for production line scheduling.

5 Conclusion

This study addressed dynamic scheduling challenges in industrial intelligent automation production lines and proposed a reinforcement learning-integrated scheduling model. The model demonstrated rapid response capabilities under disturbances such as equipment failures, urgent order insertions, and material delays, reducing fault response time to 5 seconds and urgent order response time to 9.3 minutes—significantly faster than traditional GA (15 - 30 min).

Key outcomes include:

- 80% reduction in equipment downtime.
- Order overdue rate decreased from 12% to 3%.
- Average completion time shortened by 18.7% to 218 minutes.
- Equipment utilization improved to 89.2%, achieving dual optimization of time and resources. The model exhibits high engineering practicality, with integration costs at only 15% of production line modification expenses. It supports remote operation and fault prediction via digital twins, improving operational efficiency by 40% and reducing monthly maintenance costs by ¥80,000.

Despite its advantages, further improvements are needed, particularly in predictive scheduling with digital twins, multi-line collaborative scheduling, green scheduling objectives, and algorithm lightweighting. Future research will focus on these areas.

References

- [1]Zhang W, Dietterich T G. A Reinforcement Learning Approach to Job-Shop Scheduling[C]. Proc. of 14th Int. Joint Conf. on Artificial Intelligence, 1995: 1114-1120.
- [2]Wang J, Li X, Zhu X. Intelligent dynamic control of stochastic economic lot scheduling by agent-based reinforcement learning[J]. International Journal of Production Research, 2012, 50(16): 4381-4395.
- [3]Qingdao Institute of Bioenergy and Bioprocess Technology, Chinese Academy of Sciences. Qingdao Energy Institute Develops High-Performance Bio-Based Materials Using Paper Mulberry Gene-Editing Technology. May 14, 2024.