

A Hybrid ACO-Reinforcement Learning Framework for Task Assignment

1st Sampath Kumar

Department of Computer Science Engineering
SRM Institute of Science and Technology
Chennai, India
sz4961@srmist.edu.in

2nd Dr.Krishnaraj N

Department of Networking And Communications
SRM Institute of Science and Technology
Chennai, India
krishnan2@srmist.edu.in

Abstract—This thesis presents a new approach on combining Ant Colony Optimization (ACO) with Reinforcement Learning (RL) techniques employing Q-learning and SARSA, to organize and allocate tasks in complex work environments. This approach addresses issues with order of task execution, diverse employee skills, employee availability, and fair workload distribution. ACO creates an initial task sequence considering dependencies and deadlines. Using the ACO determined task sequence, RL agents determine the initial assignments by basing their decisions depending on skill fit, workload distribution along with deadlines. This approach uses a multi objective hyperparameter tuning with Optuna library to balance high rate of task assignment with workload distribution. "Refinement RL" component provides final passthrough by balancing the workload across the skills to ensure employees are not overloaded. Experiments reveal that in task assignment rate and workload distribution, the ACO-RL combination performs better than simple greedy approach. This study adds to ongoing research on scheduling and task assignment optimization problems in complex industrial setups.

Index Terms—Reinforcement learning, Task scheduling, Resource allocation, Ant colony optimization, Manufacturing systems, Skill matching, Workload balancing

I. INTRODUCTION

Task assignment in industrial settings, for both people or machines, is a hard problem to solve while achieving smooth operations and strong employee management. This problem comes from the math of sharing resources and the mix of skills with human workload in today's workplaces.

Modern businesses must run well while keeping employees happy. Old ways of scheduling tasks following fixed rules and simple ideas do not work in work settings that change fast, with work orders moving quickly and staff levels shifting because of changing company and personal needs.

A. Challenges Across Sectors

This problem compounds even more in different sectors that rely heavily on task planning. In factories good scheduling and assignment can stop production delays or extra costs. Hospitals need schedules that meet patient needs while matching available staff and laws. Tech teams require loose schedules to handle sudden project changes. Call centres must keep service steady while facing varying demand.

B. Impact of Remote Work

To make matters complicated the increase in remote and hybrid work has complicated scheduling and task assignments. Companies must consider time zones, the nature of employee availability and frequent attrition. Old ways of scheduling and assignments based on "first come first assign" kind of greedy methodologies do not really gell well with the fast moving nature of modern businesses.

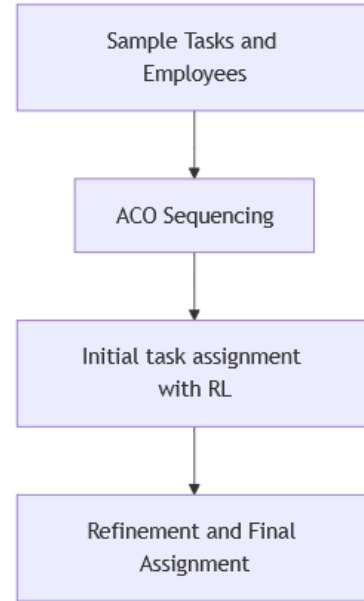


Fig. 1. Proposed methodology

C. Proposed Hybrid Approach

This study proposes a framework (1) for the complex task assignment problem in an industrial setting with task dependencies and employee availability constraints, combining Reinforcement Learning and Ant Colony Optimization. First, we begin with task sequencing by ACO, which determines an optimal task sequence that satisfies task dependencies, due date, and priority restrictions. Then, we devise an initial employee

assignment using Reinforcement Learning—specifically, Q-learning and SARSA (for comparison). Reward shaping guides these initial assignments by ensuring skills match, deadlines are met, work is balanced, and resources are used efficiently. Finally, a custom component named “RefinementRL” takes these initial assignments and optimizes them using a state representation that captures the workload, assignment, and unassigned tasks. Ultimately, the goal is to maximize a skill-level-based workload allocation and task assignment rate.

D. Significance

Significance of this study comes from the integration of Reinforcement Learning with Ant Colony Optimization in task assignment optimization problem, which approaches the focus of critical workforce management and resource allocation issue. It demonstrates how meta heuristic algorithms can improve the current reinforcement learning approaches in complex scheduling environment.

The general applicability of the framework is largely important due to its wide use in various industries, for example manufacturing to software development where it is important for resource constrained scheduling and skilled task assignment. Additionally, the framework allows RL component to continuously learn from experience and the ACO component to manage task sequencing in an efficient manner which is a major progress compared to traditional static methods of assignments. The practical solutions provided in this paper to industry-wide resource optimization problems contribute to both the academic literature on hybrid AI systems.

II. LITERATURE SURVEY

Reinforcement Learning (RL) has developed into a promising technique for dynamic environment workload optimization and task scheduling during the past years. Traditional heuristics and rule-based methods lose their value as scheduling problems become more advanced so adaptive RL-based approaches step in to replace them. The paper examines previous studies which examine reinforcement learning alongside associated optimization approaches for workload distribution and task distribution and scheduling in uncertain environments.

The research by Infantes et al. [1] explores Deep Reinforcement Learning (DRL) solutions for Job-Shop Scheduling Problems (JSSP) under uncertain conditions. GNNs enable the enhancement of DRL model generalization and scalability according to their research. DRL proves successful for managing unpredictable task duration fluctuations and outperforms conventional deterministic methods according to their examination.

Burdett and Kozan [2] study scheduling with limited resources and renewable resource assignment. The authors applied heuristic-based mathematical models for workforce allocation to develop constraint models that may be integrated into RL-based scheduling systems.

The paper by Zhong [3] examines how Q-learning performs against SARSA reinforcement learning methods for task assignment. SARSA delivers stable solutions through

its on-policy approach yet Q-learning achieves higher but riskier rewards by using its off-policy exploration method. Experts need these understandings to develop RL models which effectively manage exploration against exploitation for actual scheduling situations.

Ben Nouredine et al. [4] introduces multi agent reinforcement learning through Deep Q-learning to handle distributed task allocation. Through proposed agent teamwork operators gain better resource usage and less conflict during scheduling. The method proves beneficial for distributed systems which lack access to global information.

The research by Joo et al. [5] applies reinforcement learning methods to human-machine manufacturing systems. The research team developed a DRL model to handle human fatigue and skill abilities when performing dynamic task assignments. The research shows how human performance limitations integrated into RL-based scheduling systems produces more efficient operations and lowers worker exhaustion levels.

The research by Wibisono et al. [6] introduces a human-centered strategy for allocating resources during business process execution. The research supports RL-based scheduling by recommending real-time worker preference and availability assessment during scheduling operations.

Dastmalchian and Blyton [7] analyze how modern work patterns require adaptable scheduling systems. The research demonstrates that workforce adaptability needs to become part of all future task allocation models.

Dorigo and Stützle [8] deliver an extensive review of Ant Colony Optimization (ACO) and its scheduling problem applications. The bio-inspired optimization method ACO demonstrates great effectiveness when solving hard-to-tackle optimization problems. The review presents hybridization methods which use reinforcement learning to enhance task allocation efficiency.

III. PROBLEM FORMULATION

In modern manufacturing and service environments, the allocation of tasks to a diverse workforce poses a complex, multi-objective challenge. The aim is to maximize the number of tasks successfully assigned while ensuring that deadlines are met, workloads are balanced, and employees’ skills are effectively utilized. To address this, our methodology decomposes the problem into distinct but interrelated phases, integrating bio-inspired optimization, adaptive learning, and post-assignment refinement.

A. Task Characteristics

Each work item is characterized by several intrinsic properties:

- **Required Skills:** The specific competencies needed to complete the task.
- **Processing Time:** An estimate of the hours required for completion.
- **Deadline:** The due date by which the task must be finalized.

- **Priority:** A measure of urgency that influences scheduling order.
- **Precedence Constraints:** Dependencies indicating that certain tasks must be completed before others can begin.

These attributes introduce both temporal and logical constraints that must be considered in the scheduling process.

B. Employee Attributes

The available workforce is described by:

- **Skill Set:** The collection of competencies held by each employee.
- **Availability:** The number of work hours an employee can commit within a given period.
- **Efficiency:** A scaling factor that influences the effective time required to complete tasks.
- **Current Workload:** An evolving measure of the tasks already assigned to the employee.

A candidate employee is deemed suitable for a task if there is an overlap between the task's required skills and the employee's capabilities, and if the employee has sufficient available hours to accommodate the task, accounting for individual efficiency.

C. Hybrid Scheduling and Assignment Strategy

To navigate the NP-hard nature of this scheduling problem, our framework is structured into three sequential phases:

- 1) **Initial Task Sequencing via Bio-inspired Optimization:** A pheromone-based search process is employed to generate an initial ordering of tasks. This stage incorporates heuristic evaluations that combine task urgency, estimated processing times, deadlines, and dependency considerations. The generated sequence is periodically updated to reflect improvements in assignment performance as subsequent stages progress.
- 2) **Task Assignment through Adaptive Learning:** Building upon the sequenced tasks, a centralized learning agent dynamically assigns tasks to employees. This process considers the current state of assignments, available work hours, and skill matches, and it is designed to adapt through an exploration-exploitation mechanism. Two reinforcement learning strategies are applied to refine decision-making during the assignment process.
- 3) **Post-Assignment Refinement:** Following the initial allocation, a dedicated refinement module iteratively adjusts the task-to-employee mapping. The objective of this stage is to further balance workloads and improve the alignment of tasks with employee skills, ensuring that operational constraints are met and overall efficiency is maximized.

D. Automated Hyperparameter Tuning

An automated tuning process is integrated into the framework to optimize key parameters across all stages. Using a multi-objective optimization approach, the tuner adjusts factors influencing the bio-inspired sequencing, the learning dynamics in the assignment phase, and the iterative refinement

process. The objective is to enhance task coverage while simultaneously reducing variability in performance metrics, thus ensuring a robust and adaptive solution.

IV. STATE AND REWARD DESIGN

A. Reinforcement Learning Framework

To accommodate multi-column layouts, we include short in-text equations and place a wide table at the bottom or top of the page. Let N be the total number of employees.

a) *State Definition.*: We denote the RL state as

$$s = (i, \mathbf{H}, \mathbf{A}, \mathbf{F}) \quad (1)$$

where

- i : index in the task sequence,
- $\mathbf{H} = [h_1, \dots, h_N]$: residual work hours per employee,
- \mathbf{A} : binary record of which tasks have been assigned,
- \mathbf{F} : binned features such as priority or due date.

b) *Reward Function.*: At each step, the agent receives:

$$r(s, a) = R_{\text{feas}} + R_{\text{penalty}} + R_{\text{balance}} + R_{\text{efficiency}} \quad (2)$$

with terms

$$R_{\text{feas}}(s, a) = \lambda p_t \mathbb{I}\{\text{feasible}\}, \quad (3)$$

$$R_{\text{penalty}}(s, a) = -\kappa \mathbb{I}\{\text{violation}\} - \nu \mathbb{I}\{\text{priority skip}\}, \quad (4)$$

$$R_{\text{balance}}(s, a) = \alpha \left(1 - \frac{\sigma(W)}{\bar{W}} \right), \quad (5)$$

$$R_{\text{efficiency}}(s, a) = \eta \frac{\Delta\tau}{\tau}, \quad (6)$$

where $\lambda, \kappa, \nu, \alpha, \eta$ are weighting constants; p_t indicates the task's priority; $\sigma(W), \bar{W}$ are the standard deviation and mean of assigned hours; and $\Delta\tau/\tau$ measures how much an employee's efficiency reduces task duration.

B. Refinement Phase

A subsequent refinement step uses an augmented state and reward:

$$S = (\mathbf{W}, M, U) \quad (7),$$

where $\mathbf{W} = [w_1, \dots, w_N]$ captures each employee's total workload, M is the mapping of tasks to employees, and U is the set of unassigned tasks.

The refinement reward is:

$$r_{\text{ref}}(S, a) = \beta_1 R_{\text{match}} - \beta_2 R_{\text{imbalance}} + \beta_3 R_{\text{eff}}, \quad (8)$$

where $\beta_1, \beta_2, \beta_3$ weight the emphasis on skill-task alignment (R_{match}), workload fairness ($R_{\text{imbalance}}$), and additional efficiency gains (R_{eff}).

TABLE I
NOTATION AND DESCRIPTIONS

Symbol	Description
i	Index in the task sequence
\mathbf{H}	Residual hours vector, $[h_1, \dots, h_N]$
\mathbf{A}	Binary assignment record
\mathbf{F}	Binned task features (priority, due date)
s	RL state in (1)
$r(s, a)$	RL reward in (2)
$\lambda, \kappa, \nu, \alpha, \eta$	RL weighting constants
p_t	Priority of the current task
$\sigma(W), \bar{W}$	Std. dev. and mean of assigned hours
$\frac{\Delta\tau}{\tau}$	Relative reduction in task duration (efficiency)
S	Refinement state in (7)
\mathbf{W}	Workloads: $[w_1, \dots, w_N]$
M	Task-employee mapping
U	Set of unassigned tasks
$r_{\text{ref}}(S, a)$	Refinement reward in (8)
$\beta_1, \beta_2, \beta_3$	Refinement reward coefficients

V. ALGORITHMIC ARCHITECTURE AND DESIGN

In dynamic operational environments, the challenge of scheduling tasks and allocating them to a diverse workforce demands a robust and adaptable framework. Our system integrates heuristic optimization via Ant Colony Optimization (ACO), centralized reinforcement learning (RL) for task assignment, and an iterative post-assignment refinement process—all underpinned by automated hyperparameter tuning using Optuna.

A. Framework Overview

The architecture is organized into three core modules:

- 1) **Task Sequencing:** An ACO-based module generates an initial task sequence by considering inter-task dependencies, deadlines, and priorities. The module leverages configurable parameters (e.g., pheromone influence α , heuristic weight β , evaporation rate ρ , etc.) to guide the search.
- 2) **RL-based Assignment:** Centralized RL agents (employing both Q-learning and SARSA variants) assign tasks to employees. The state representation includes the current task index, discretized measures of employee workload, and historical assignment data. An ϵ -greedy policy—with optional Upper Confidence Bound (UCB) enhancements—is used to balance exploration and exploitation.
- 3) **Post-Assignment Refinement:** A dedicated refinement module further optimizes the initial assignment. By iteratively adjusting the task-to-employee mapping based on updated workload distributions and unassigned tasks, this stage improves overall balance and skill alignment.

B. Task Sequencing via Ant Colony Optimization (ACO)

The ACO module starts by initializing a pheromone matrix \mathbf{P} with parameters such as α , β , and ρ (all configurable via the system settings). Each task is assigned a heuristic value η_{ij} based on factors like priority, estimated duration, and deadline

constraints. The probability of transitioning from task t_i to task t_j is computed as:

$$p_{ij} = \frac{[\mathbf{P}_{ij}]^\alpha [\eta_{ij}]^\beta}{\sum_{k \in \mathcal{N}_i} [\mathbf{P}_{ik}]^\alpha [\eta_{ik}]^\beta},$$

where \mathcal{N}_i is the set of candidate tasks for task t_i . After constructing candidate sequences with multiple agents, the pheromone levels are updated by:

$$\mathbf{P}_{ij} \leftarrow (1 - \rho) \mathbf{P}_{ij} + \Delta \mathbf{P}_{ij},$$

with $\Delta \mathbf{P}_{ij}$ reflecting the quality of the sequence. This process is re-invoked periodically (e.g., every few hundred training episodes) to adapt to evolving system conditions.

C. RL-based Task Assignment and Post-Assignment Refinement

The RL module formulates task assignment as a sequential decision process. The state s captures the current task index, residual working hours for each employee, and a history of past assignments. Actions include assigning a task to an employee or opting to skip it if constraints cannot be met. For Q-learning, the update rule is:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right], \quad (1)$$

and for SARSA:

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]. \quad (2)$$

A scheduled decay of the exploration rate ϵ , along with optional UCB adjustments, guides the learning process. Following the initial assignment, the *RefinementRL* module further optimizes the allocation by considering an updated state S (comprising current workloads, the mapping of tasks to employees, and unassigned tasks) and applying a similar RL update:

$$Q(S, a) \leftarrow Q(S, a) + \alpha_r \left[r + \gamma_r \max_{a'} Q(S', a') - Q(S, a) \right].$$

Here, the refinement parameters α_r and γ_r are tuned specifically to improve workload balance and alignment with employee skills.

D. Hyperparameter Optimization and Experimental Evaluation

A standout feature of the design is the integration of automated hyperparameter tuning using Optuna. This module optimizes key parameters for both the ACO and RL components (such as learning rates, discount factors, exploration decay, and pheromone parameters) through a multi-objective process. The objectives are to maximize the proportion of tasks assigned and to minimize the variance of the reward signal. Extensive experimental trials yield Pareto-optimal configurations, with detailed visualizations and reports generated for comprehensive performance analysis.

E. Legend of Symbols

Here are the symbols used in this section

TABLE II
LEGEND OF SYMBOLS

Symbol	Description
i	Current task index
\mathbf{P}	Pheromone matrix
η_{ij}	Heuristic value for transitioning from task t_i to t_j
ρ	Pheromone evaporation rate
s, s'	RL states before and after an action
a, a'	Actions chosen in the current and subsequent states
α, α_r	Learning rates for initial assignment and refinement
γ, γ_r	Discount factors for initial assignment and refinement
ϵ	Exploration rate in RL
\mathbf{H}	Vector of residual working hours
r	Immediate reward signal
\mathcal{S}	State representation for refinement
M	Task-to-employee mapping
U	Set of unassigned tasks

VI. INFERENCE AND DISCUSSION

The experimental results demonstrate that the proposed hybrid framework significantly enhances the balance of task assignments among employees without compromising overall assignment rates or deadline compliance. Unlike traditional heuristic methods, which often concentrate tasks among a few specialized workers, our integration of reinforcement learning with ant colony optimization intrinsically accounts for workload distribution.

A. Workload Distribution Insights

Our findings indicate that greedy assignment methods typically overload a subset of employees due to repeated allocation of tasks that match their skill profiles. This concentration not only creates bottlenecks but also elevates the risk of fatigue and reduces overall efficiency. By integrating a workload balancing term into the reward function, the reinforcement learning module gradually redistributes tasks more evenly. As a result, the refined policies exhibit a substantial reduction in workload variance, as evidenced by a lower standard deviation of assigned hours relative to baseline approaches.

B. Resource Utilization and Task Coverage

An equitable distribution of tasks also enhances overall resource utilization. The hybrid framework ensures that no single resource becomes a critical point of failure by reallocating less urgent tasks to underutilized employees. This reallocation enables specialized personnel to concentrate on high-priority assignments, leading to improved long-term performance, decreased burnout risk, and more sustainable operational practices. Notably, high task coverage and timely completion are maintained even as the system emphasizes balanced workload distribution.

C. Algorithmic Rationale

The observed improvements in workload balance can be directly attributed to the modified reward structure employed by the reinforcement learning component. Rather than solely

incentivizing rapid task completion and strict adherence to deadlines, the reward function penalizes excessive task concentration. Over multiple training episodes, the RL agents learn to recognize and internalize the cost of overburdening specific employees, thereby promoting a more uniform distribution of work. This adaptive learning mechanism effectively reconciles the dual objectives of efficiency and fairness within the scheduling process.

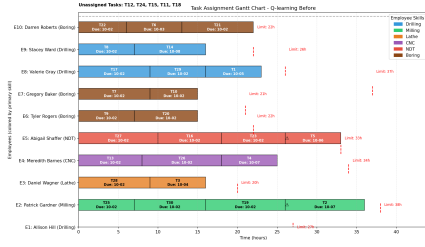
D. Practical Trade-offs and Future Considerations

Although the hybrid ACO-RL framework markedly improves workload distribution, its practical deployment requires careful tuning of reward parameters to align with organizational priorities. In certain scenarios, urgent tasks might necessitate a temporary deviation from perfect balance in order to assign critical tasks to the most qualified personnel. Future enhancements could include dynamic adjustment of reward components in response to real-time performance metrics. Overall, the incorporation of workload balance as a primary objective not only advances scheduling methodologies but also contributes to more resilient and sustainable workforce management.

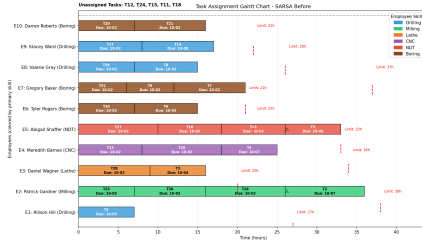
VII. FUTURE WORK

While the current framework—which combines ACO-based task sequencing with reinforcement learning for employee assignment—has yielded promising results, several avenues for further research remain:

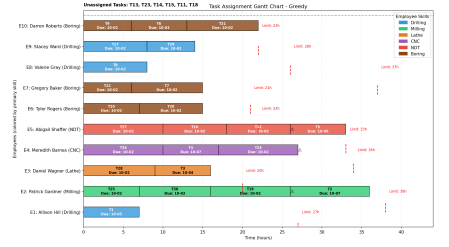
- **Enhanced State Representations:** The present RL models utilize discretized state representations. Future research could explore richer representations, such as deep neural network embeddings, to capture complex task dependencies and nuanced employee characteristics more effectively.
- **Hybrid Optimization Techniques:** Integrating additional metaheuristic or evolutionary algorithms with ACO may lead to improved task sequencing. Investigating hybrid models that merge global search capabilities with localized refinement strategies could further optimize scheduling performance.
- **Adaptive Reward Mechanisms:** Given the critical role of the reward structure, future work could develop dynamic or adaptive reward schemes that adjust penalties and incentives based on real-time feedback or historical data, thereby enhancing both workload balance and resource utilization.
- **Scalability and Real-Time Implementation:** Although our experiments were conducted on a moderate scale, real-world applications often involve larger task sets and employee pools. Future studies should address scalability challenges and explore real-time implementations, possibly leveraging distributed computing frameworks.
- **Incorporating Uncertainty:** Extending the framework to account for uncertainties in task durations, employee availability, and environmental changes would increase its robustness. Incorporating probabilistic models or robust



(a) Q-learning Before

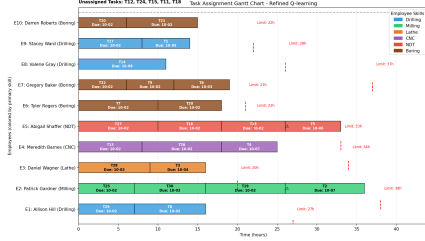


(b) SARSA Before

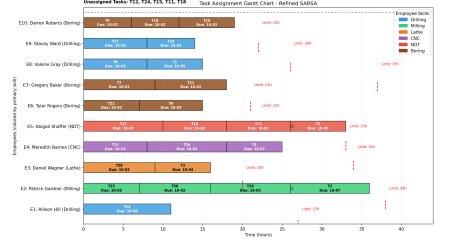


(c) Greedy

Fig. 2. Comparison of baseline scheduling approaches.



(a) Refined Q-learning



(b) Refined SARSA

Fig. 3. Comparison of refined scheduling approaches.

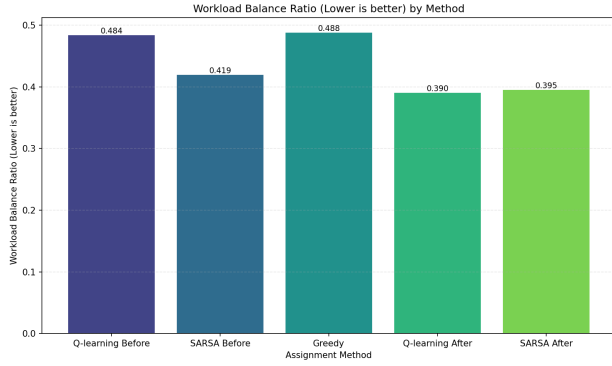


Fig. 4. Workload Balance Ratio (lower is better) by assignment method.

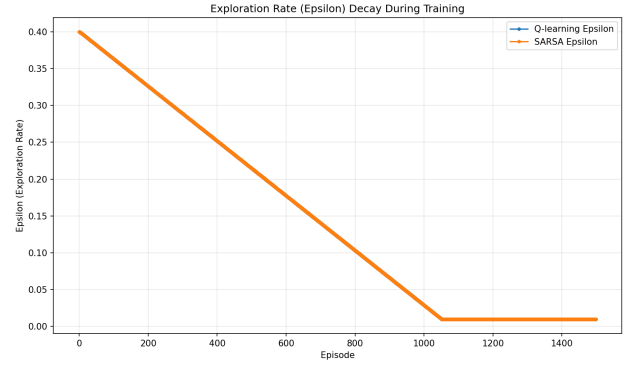


Fig. 6. Exploration rate (ϵ) decay during training.

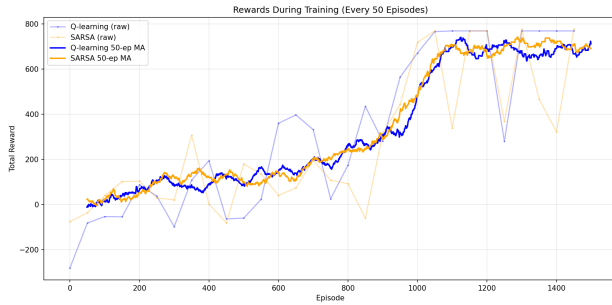


Fig. 5. Rewards during training (every 50 episodes).

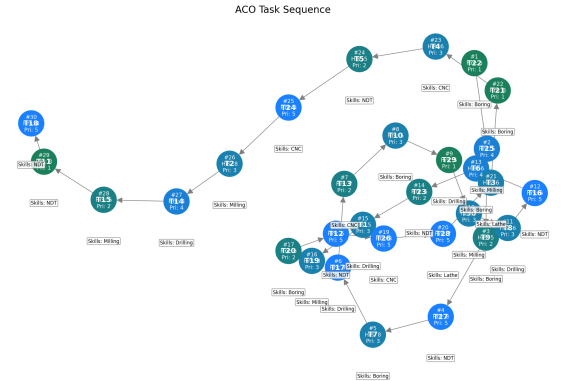


Fig. 7. Illustration of the Ant Colony Optimization (ACO) task sequence.

optimization techniques into the RL component may better handle stochastic variations.

- **User Feedback Integration:** Future systems could incorporate continuous feedback from domain experts and end-users, enabling a human-in-the-loop approach that iteratively refines both task sequencing and assignment processes based on practical, real-world insights.

These research directions aim to bridge the gap between theoretical scheduling models and practical, scalable solutions for dynamic, multi-agent systems. By addressing these challenges, future work can further enhance the adaptability and effectiveness of the hybrid framework in diverse industrial environments.

ACKNOWLEDGMENT

The authors would like to express their sincere gratitude to SRM Institute of Technology for providing the resources and support necessary for this study. We are especially grateful to Professor Dr. Krishnaraj N for his invaluable guidance, insightful suggestions, and continuous encouragement throughout the course of this study. His expertise and constructive feedback have been instrumental in shaping the direction and success of this work.

VIII. REFERENCES

REFERENCES

- [1] G. Infantes et al., "Learning to Solve Job Shop Scheduling under Uncertainty," CPAIOR, 2024.
- [2] R. L. Burdett and E. Kozan, "The Assignment of Individual Renewable Resources in Scheduling," *Asia Pacific Journal of Operational Research*, 2004.
- [3] L. Zhong, "Comparison of Q-learning and SARSA Reinforcement Learning Models on Cliff Walking Problem," DAI, 2024.
- [4] D. Ben Noureddine et al., "Multi-agent Deep Reinforcement Learning for Task Allocation in Dynamic Environment," ICSoft, 2017.
- [5] T. Joo et al., "Task Allocation in Human-Machine Manufacturing Systems Using Deep Reinforcement Learning," *Sustainability*, 2022.
- [6] A. Wibisono, A. S. Nisafani, H. Bae, and Y. Park, "A dynamic and human-centric resource allocation for managing business process execution," University of Texas at El Paso, Tech. Rep., Dec. 2016.
- [7] A. Dastmalchian and P. Blyton, "Workplace flexibility and the changing nature of work: An introduction," *J. Management Studies*, vol. 38, no. 2, pp. 281-286, 2001.
- [8] M. Dorigo and T. Stützle, "Ant colony optimization: Overview and recent advances," in *Handbook of Metaheuristics*, 3rd ed. Springer, 2016, pp. 311-351.