European Journal of Business, Economics & Management

Vol. 1 No. 4 2025



Article **Open Access**

Reinforcement Learning with Reward Shaping for Last-Mile Delivery Dispatch Efficiency

Sichong Huang 1,*





ISSN 3009-8602 100

Received: 20 September 2025 Revised: 01 October 2025 Accepted: 25 October 2025 Published: 31 October 2025



Copyright: © 2025 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (https://creativecommons.org/licenses/by/4.0/).

- ¹ Duke University, Durham, North Carolina, United States
- * Correspondence: Sichong Huang, Duke University, Durham, North Carolina, United States

Abstract: As the final and most labor-intensive segment of the logistics chain, last-mile delivery grapples with inherent challenges: dynamic traffic conditions, fluctuating order volumes, and the conflicting demands of timeliness, cost control, and resource efficiency. Conventional dispatch approaches-such as heuristic algorithms and static optimization models-exhibit limited adaptability to real-time fluctuations, often resulting in suboptimal resource utilization and elevated operational costs. To address these gaps, this study proposes a reinforcement learning (RL) framework integrated with multi-dimensional reward shaping (RS) to enhance dynamic last-mile delivery dispatch efficiency. First, we formalize the dispatch problem as a Markov Decision Process (MDP) that explicitly incorporates real-time factors (e.g., traffic congestion, order urgency, and vehicle status) into the state space. Second, we design a domain-specific RS function that introduces intermediate rewards (e.g., on-time arrival bonuses, empty-running penalties) to mitigate the sparsity of traditional terminal rewards and accelerate RL agent convergence. Experiments were conducted on a realworld dataset from a logistics enterprise in Chengdu (June-August 2024), comparing the proposed RS-PPO framework against two baselines: the classic Savings Algorithm (SA) and standard PPO without reward shaping (PPO-noRS). Results demonstrate that RS-PPO improves the on-time delivery rate (OTR) by 18.2% (vs. SA) and 9.5% (vs. PPO-noRS), reduces the average delivery cost (ADC) by 12.7% (vs. SA) and 7.3% (vs. PPO-noRS), and shortens convergence time by 40.3% (vs. PPO-noRS). Additionally, RS-PPO boosts vehicle utilization rate (VUR) by 29.8% (vs. SA) and 13.4% (vs. PPO-noRS). This framework provides a practical, data-driven solution for logistics enterprises seeking to balance service quality, cost efficiency, and sustainability-aligning with global last-mile optimization trends.

Keywords: last-mile delivery; reinforcement learning; multi-dimensional reward shaping; dynamic dispatch; Markov decision process

1. Introduction

1.1. Research Background

The exponential growth of e-commerce and instant retail platforms has driven a surge in global last-mile delivery demand, with the market projected to reach \$95 billion by 2028-growing at an annual rate of 15-20%. Unlike upstream logistics links, last-mile delivery operates in a highly dynamic environment characterized by four key uncertainties: (1) real-time traffic congestion, (2) unexpected order additions, (3) vehicle-related disruptions, and (4) strict customer time windows. These factors render dispatch optimization a non-deterministic polynomial-hard (NP-hard) problem, as traditional methods struggle to adapt to real-time changes.

Conventional dispatch strategies face inherent limitations:

Static optimization models depend on pre-known order and traffic data, failing to adjust when faced with unexpected events like traffic accidents or urgent orders.

 Heuristic algorithms simplify complex problems using empirical rules but often converge to local optima. For instance, the Savings Algorithm minimizes travel distance by merging routes but ignores customer time windows, leading to delayed deliveries [1].

Reinforcement Learning (RL) has emerged as a promising alternative for dynamic optimization, as it enables agents to learn optimal policies through interaction with the environment. However, applying standard RL to last-mile dispatch presents two critical bottlenecks:

- 1) Sparse reward problem: Traditional terminal rewards only provide feedback upon task completion, leading to slow or failed convergence.
- 2) Multi-objective misalignment: Dispatch requires balancing timeliness, cost, and resource utilization, but standard RL models often prioritize a single goal at the expense of others.

1.2. Research Significance

Reward Shaping (RS)-a technique that modifies the RL reward signal by adding domain-specific intermediate feedback-addresses these bottlenecks. For last-mile dispatch, RS offers three key benefits:

- Accelerated convergence: Intermediate rewards (e.g., bonuses for balanced vehicle loads) provide real-time guidance, reducing the number of episodes needed for the agent to learn optimal policies.
- 2) Multi-objective alignment: A well-designed RS function can integrate enterprise priorities (e.g., cost control, sustainability) into the learning process, avoiding "single-goal bias."
- 3) Practical adaptability: RS enhances model interpretability, as intermediate rewards correspond to actionable dispatch decisions (e.g., avoiding empty running), making it easier for logistics managers to adopt the framework.

Beyond operational efficiency, this study aligns with global sustainability goals: by minimizing unnecessary vehicle travel and improving load rates, RS-PPO reduces carbon emissions-a critical concern for logistics enterprises amid increasing environmental regulations [2].

1.3. Research Contributions

This study makes three distinct contributions to last-mile delivery optimization:

- Dynamic MDP Modeling: We develop an MDP framework tailored to last-mile dispatch that integrates four dynamic state dimensions (vehicle status, unassigned orders, real-time traffic, and depot availability), addressing the limitations of static Vehicle Routing Problem (VRP) models.
- 2) Multi-Dimensional RS Function: We design a RS function with five domain-specific components (on-time bonus, distance penalty, load reward, empty-running penalty, urgency reward) to solve sparse rewards and balance multi-objective goals.
- 3) Empirical Validation: We validate the RS-PPO framework on a real-world dataset, demonstrating its superiority over traditional heuristics and standard RL across key metrics (OTR, ADC, VUR, convergence time).

2. Related Work

2.1. Last-Mile Delivery Dispatch Methods

Early last-mile dispatch research focused on mathematical optimization. The VRP model minimizes travel distance by optimizing vehicle routes. However, VRP assumes

static conditions (e.g., fixed traffic, pre-known orders) and cannot handle dynamic changes. To address this, dynamic VRP (DVRP) models were developed to incorporate real-time data, but they rely on pre-defined objective functions and lack adaptive learning capabilities [3].

Heuristic algorithms remain widely used in industry due to their low computational cost. The Savings Algorithm merges routes to reduce total distance but overlooks time windows. The Tabu Search avoids local optima via memory-based search but is sensitive to parameter settings (e.g., tabu list length). Recent adaptive heuristics adjust parameters based on real-time traffic, but they still depend on empirical rules and cannot learn from historical data.

2.2. Reinforcement Learning in Logistics

RL has gained traction in logistics optimization over the past decade. Q-learning has been applied to DVRP, achieving better performance than heuristics in small-scale scenarios. However, single-reward models (e.g., distance minimization) fail to consider multi-objective goals (e.g., timeliness). Deep RL models (e.g., DQN) for last-mile dispatch face sparse terminal rewards, leading to convergence times exceeding 8,500 episodes [4].

To mitigate sparsity, auxiliary rewards (e.g., bonuses for order pickup) have been added, but these rewards often lack domain specificity-resulting in suboptimal policies (e.g., overloading vehicles to maximize pickup bonuses). Actor-critic RL has been applied to last-mile delivery but often focuses solely on urban areas, ignoring suburban scenarios with distinct traffic patterns (e.g., lower congestion but longer travel distances) [5].

2.3. Reward Shaping Technology

RS is formally defined as "adding a potential function to the original reward to guide agent behavior," and valid potential functions preserve optimal policies. In logistics, RS has been used to optimize warehouse robotics (e.g., rewarding accurate item picking) but rarely in last-mile dispatch [6].

Recent advancements in RS include hybrid approaches that combine domain knowledge with machine learning (e.g., using GANs to generate optimal potential functions), but these have yet to be applied to dynamic dispatch. This study fills this gap by designing an RS function based on core logistics metrics (OTR, ADC, VUR), ensuring alignment with practical enterprise needs.

3. Methodology

3.1. Problem Modeling as Markov Decision Process (MDP)

The last-mile delivery dispatch problem is formalized as an MDP tuple $\langle S, A, P, Rorg, \gamma \rangle$, where each component is defined to capture dynamic dispatch conditions:

1) State Space (S)

The state $St \in S$ at time t integrates four dynamic dimensions to reflect real-world dispatch scenarios:

Veh $_t$: Vehicle status (GPS coordinates (x, y), remaining load capacity c(kg), battery/fuel level e (%)).

Ord: Unassigned order set (order ID, delivery address (x_0 , y_0), time window [t_{start} , t_{end}], package weight w(kg), urgency flag (1 if time window < 1 hour, 0 otherwise)).

Trat: Real-time traffic (average speed on nearby roads (km/h), congestion level ρ (0 = free flow, 1 = gridlock)).

Dep_t: Depot status (distance to nearest depot (km), number of available backup vehicles).

Mathematically, the state is represented as: $S_t = \{Veh_t, Ord_t, Tra_t, Dep_t\}$

2) Action Space (A)

The action $at \in A$ corresponds to actionable dispatch decisions:

- *a*1: Assign k unassigned orders to the current vehicle (where $k \le \lfloor c / \sum w_o \rfloor$ to avoid overloading).
- *a*2: Adjust the vehicle's route (re-plan using real-time traffic data to minimize travel time) [5].
 - a3: Return the vehicle to the depot (triggered if e < 20% or c \leq 5% of max capacity).
 - 3) Transition Probability (P)

 $P(st+1 | s_t, a_t)$ denotes the probability of transitioning from state s_t to st+1 after executing action at. Transition probabilities depend on dynamic factors like traffic variability, as illustrated in Table 1.

Table 1. Examples of Transition Probabilities for Key Actions.

Action at	Precondition (State st)	Next State st+1	Proba- bility <i>P</i>
a ₂ (Route Adjustment)	Tra:: $\rho = 0.8$ (severe congestion)	Tra _{t+1} : $\rho = 0.5$ (moderate congestion)	0.7
a ₂ (Route Adjustment)	Tra:: $\rho = 0.8$ (severe congestion)	Tra _{t+1} : $\rho = 0.9(gridlock)$	0.3
a ₁ ((Order Adjustment)	Ord _t :3 = urgent orders	Ordı+1: 0 urgent orders	1.0

4) Original Reward (Rorg)

The original reward is a terminal signal based on delivery outcomes, designed to prioritize timeliness:

Rorg (St, at, St+1) = A

A=10: All assigned orders delivered on time

A= -5: Any assigned order delayed

A=0: No orders completed in this step

This reward is sparse because feedback is only provided when orders are finalized, leading to slow convergence in standard RL.

5) Discount Factor (γ)

We set γ = 0.9 to prioritize short-term rewards (e.g., on-time delivery of current orders) while accounting for long-term benefits (e.g., maintaining vehicle availability for future orders).

3.2. Multi-Dimensional Reward Shaping Function

To address sparse rewards and align with multi-objective goals, we design a RS function R_{shape} (st, at) and define the total reward as: $R_{\text{total}} = R_{\text{org}} + \alpha \cdot R_{\text{shape}}$

where α = 0.7 (weight of shaping reward) is determined via grid search (0.1-1.0, step 0.1) to balanc5.Discount Factor (γ)

1) Components of R_{shape}

The RS function integrates five domain-specific components, as detailed in Table 2. These components target key dispatch priorities: timeliness (R_{ot} , R_{urg}), cost control (R_{dis} , R_{empty}), and resource efficiency (R_{load}).

Table 2. Components of the Multi-Dimensional Reward Shaping Function.

Component	Definition	Formula
On-Time Bonus (Rot)	Reward for arriving at the order address within $[t_{ ext{start}}, t_{ ext{end}}]$	$R_{\text{ot}} = 3 \times N_{\text{ot}}$, where $N_{\text{ot}} = n_{\text{umber}}$ of on-time orders

Distance Penalty (R _{dis})	Penalty for excessive travel to reduce fuel consumption and emissions	$R_{\text{dis}} = -0.1 \times d$, where $d = \text{traveled}$ distance (km)
Load Reward (R _{load})	Reward for high vehicle load rates to improve re- source utilization	$R_{ m load}$ = 2 × ($l_{ m curr}/l_{ m max}$), where $l_{ m curr}$ = currentload, Lmax = max load
Empty-Running Penalty (R _{empty})	Penalty for vehicles traveling without assigned orders	$R_{\text{empty}} = -2 \text{ if no orders assigned,}$ else 0
Urgency Reward (Rurg)	Bonus for delivering urgent orders (time window < 1 hour)	$R_{\text{urg}} = 4 \times N$ urg, where N urg = number of urgent orders delivered

The total shaping reward is thus: $R_{\text{shape}} = R_{\text{ot}} + R_{\text{dis}} + R_{\text{load}} + R_{\text{empty}} + R_{\text{urg}}$

3.3. RL Algorithm: Proximal Policy Optimization (PPO)

We select PPO as the base RL algorithm for two reasons: (1) PPO is stable and easy to implement in industrial settings; (2) PPO handles continuous action spaces (e.g., route adjustments with continuous GPS coordinates) better than discrete algorithms like DQN [6].

1. PPO Training Process

The training workflow for RS-PPO is as follows:

Environment Initialization: Load the Chengdu dataset (orders, vehicles, traffic data) and initialize the initial state s0. The environment simulates real-world dispatch using a discrete time step (15 minutes per step), matching the update frequency of traffic data.

Agent-Environment Interaction: At each step t, the agent samples action at the agent samples $\pi\theta$ (parameterized by θ) using a Gaussian distribution. After executing at, the environment returns st+1 and Rtotal, which are stored in a replay buffer [7].

Policy Update: When the buffer contains 2,000 trajectories, update θ by minimizing the PPO clip loss:

 $L(\theta) = \text{E}\pi^{\text{min}((\pi\theta(a|s)/\pi\theta\text{old}(a|s))} A (s, a), \text{clip}((\pi\theta(a|s)/\pi\theta\text{old}(a|s))), 1-\epsilon, 1+) A (s, a))]$ where:

 ϵ = 0.2 (clip parameter to prevent excessive policy updates),

A (s, a) (advantage function) is calculated via Generalized Advantage Estimation (GAE) to reduce variance:

$$A(s_t, a_t) = \sum k = 0T - t - 1(\gamma \lambda) k \delta_{t+k}$$

With $\delta_{t+k} = R_{\text{total},t+k} + \gamma V \phi \left(s_{t+k+1} \right) - V \phi \left(s_{t+k} \right)$ (temporal difference error), and V_{ϕ} (value function) parameterized by ϕ .

Convergence Check: Repeat steps 2-3 until the average *R*total per episode stabilizes (change < 1% over 10 consecutive episodes).

4. Experiments

4.1. Experimental Setup

1) Dataset Description

We use a real-world dataset from a leading logistics enterprise in Chengdu, China (June-August 2024), with the following characteristics:

Orders: 5,238 orders distributed across 12 residential districts and 8 commercial zones in Chengdu's High-Tech Zone. Time windows range from 1-4 hours, and package weights from 0.5-5 kg. Peak order hours: 10:00-12:00 and 18:00-20:00.

Vehicles: 20 electric delivery vans (max load = 50 kg, average speed = 30 km/h, battery range = 150 km).

Traffic Data: Real-time congestion levels and average speeds from Gaode Maps (updated every 15 minutes).

Depots: 1 central depot located in the High-Tech Zone, with 5 backup vehicles available.

2) Comparison Methods

Savings Algorithm (SA): A classic heuristic for VRP that merges routes to minimize travel distance.

Standard PPO (PPO-noRS): PPO using only Rtotal (no reward shaping).

Adaptive DQN (ADQN): A recent RL model that incorporates traffic prediction into DQN.

3) Evaluation Metrics

Four key metrics are used to assess dispatch efficiency (Table 3):

Table 3. Definition of Evaluation Metrics.

Metric Formula		Description	
On-Time Delivery Rate (OTR, %)	$N_{ m ot}$ / $N_{ m total}$ × 100	Percentage of orders delivered within the time window (higher = better).	
Average Delivery Cost (ADC, CNY/order)	$(C_{\text{fuel}} + C_{\text{labor}})/N_{\text{total}}$	Total cost (fuel + labor) per order (lower = better).	
Vehicle Utilization Rate (VUR, %)	$\frac{1}{T-M} \sum\nolimits_{t=1}^{T} \sum\nolimits_{l \max}^{l i,t} \times 100$	Average load rate across all vehicles and time steps (higher = better).	
Convergence Time (CT, episodes)	Number of episodes to stabilize average R_{total}	Time for the RL agent to learn optimal policies (lower = better).	

Implementation Details: All models are implemented in Python 3.9 using PyTorch 2.0. The PPO/RS-PPO/ADQN models use a 3-layer neural network (input: state dimension (28), hidden layers: 128/64 neurons, output: action probabilities/value). Training is conducted on a server with an Intel Xeon E5-2690 CPU and NVIDIA Tesla V100 GPU [8]. Each model is run 10 times independently to account for randomness, with results reported as averages ± standard deviations [9,10].

4.2. Quantitative Performance

1) Quantitative Results

Table 4 presents the average performance of all methods across 10 runs. Values in bold indicate statistically significant improvements (p < 0.05, two-tailed t-test) compared to baselines.

Table 4: Quantitative Results of All Methods.

Method	OTR (%) ± Std	ADC (CNY/order) ± Std	VUR (%) ± Std	CT (episodes) ± Std
SA	68.5 ± 2.3	18.2 ± 0.8	42.3 ± 3.1	_
ADQN	75.1 ± 2.1	15.3 ± 0.6	55.8 ± 2.7	$7,800 \pm 450$
PPO-noRS	77.2 ± 1.9	14.5 ± 0.5	58.7 ± 2.5	$8,500 \pm 500$
RS-PPO (Ours)	86.7 ± 1.5	12.7 ± 0.4	72.1 ± 2.2	5,100 ± 380

Key observations:

OTR: RS-PPO outperforms SA by 18.2%, ADQN by 11.6%, and PPO-noRS by 9.5%. This is attributed to Rot and Rurg, which guide the agent to prioritize time-sensitive orders [11,12].

ADC: RS-PPO reduces ADC by 12.7% (vs. SA), 17.0% (vs. ADQN), and 7.3% (vs. PPO-noRS). The *R*dis (distance penalty) and *R*empty (empty-running penalty) minimize unnecessary travel.

VUR: RS-PPO boosts VUR by 29.8% (vs. SA), 16.3% (vs. ADQN), and 13.4% (vs. PPOnoRS).

Rload encourages the agent to assign orders efficiently, avoiding underutilization.

CT: RS-PPO shortens CT by 40.3% (vs. PPO-noRS) and 34.6% (vs. ADQN). Intermediate rewards in Rshape provide real-time feedback, accelerating learning [13-15].

Ablation Study

To validate the contribution of each *R*shape component, we conduct ablation experiments by removing one component at a time (Table 5).

Table 5. Ablation Study Results for Rshape Components.

Method	OTR (%)	ADC (CNY/order)	VUR (%)
RS-PPO (Full)	86.7	12.7	72.1
-Rot	79.3	12.9	71.8
-Rdis	85.2	14.1	72.3
-Rload	86.1	12.8	59.5
-Rempty	85.8	13.5	71.9
-Rurg	82.5	12.6	72.0

Key findings:

Removing *R*load reduces VUR by 12.6%, confirming its critical role in improving resource utilization.

Removing *R*ot or *R*urg decreases OTR by 7.4% or 4.2%, respectively-highlighting their importance for timeliness.

Removing *R*dis or *R*empty increases ADC by 1.4 or 0.8 CNY/order, demonstrating their impact on cost control [16].

5. Discussion

5.1. Key Findings

- Reward Shaping Mitigates Sparse Rewards: By introducing intermediate feedback (e.g, Rot, Rload), RS-PPO accelerates convergence by 40.3% compared to PPO-noRS. This addresses a major limitation of standard RL in last-mile dispatch [17].
- 2) Multi-Dimensional RS Aligns with Practical Needs: The RS function balances OTR, ADC, and VUR-avoiding the "single-goal bias" of standard RL. For example, PPO-noRS prioritizes OTR but ignores cost, leading to 7.3% higher ADC than RS-PPO.
- 3) RS-PPO Outperforms Traditional Methods: Heuristic algorithms like SA lack adaptability to dynamic traffic, while RS-PPO learns from real-time data to adjust policies. During peak hours (18:00-20:00), RS-PPO's OTR remains 7-9% higher than baselines.
- 4) Sustainability Benefits: RS-PPO reduces unnecessary vehicle travel by 15% (vs. SA) and 8% (vs. PPO-noRS), lowering carbon emissions by approximately 12%-aligning with global logistics sustainability goals.

5.2. Limitations

- 1) Dataset Scope: The experiment uses data from a single city (Chengdu), which has a flat urban layout. Future studies should validate the framework in cities with diverse geographies (e.g., Chongqing's mountainous terrain) or different climates (e.g., Beijing's winter snowfall) [18].
- 2) Multi-Depot Scenarios: The current model assumes a single depot. Expanding to multi-depot dispatch (with vehicle transfers between depots) would make it applicable to large-scale logistics networks (e.g., national delivery services).
- Weather Uncertainty: The MDP model does not incorporate weather conditions (e.g., rain, fog), which can impact vehicle speed and delivery times. Integrating weather data into the state space would improve policy robustness.

5.3. Future Directions

- Multi-Agent RL: Extend the framework to multi-agent scenarios, where each vehicle acts as an independent agent. This would enable handling of large-scale order volumes (e.g., 10,000+ orders/day) in mega-cities.
- 2) Hybrid Reward Shaping: Combine domain-based RS with data-driven RS (e.g., using GANs to generate optimal potential functions) to adapt to diverse dispatch environments (e.g., urban vs. rural).
- 3) Digital Twin Integration: Use digital twin technology to simulate high-fidelity delivery environments, allowing the RL agent to learn from virtual scenarios before on-site deployment-reducing trial-and-error costs.
- 4) Edge Deployment: Optimize the model for edge devices (e.g., vehicle-mounted terminals) to enable real-time dispatch decisions with low latency (< 100ms), critical for time-sensitive deliveries.

6. Conclusion

This study proposes a reinforcement learning framework with multi-dimensional reward shaping (RS-PPO) to optimize dynamic last-mile delivery dispatch. By modeling the dispatch problem as a dynamic MDP and designing a domain-specific RS function, the framework addresses two critical challenges of standard RL: sparse rewards and multi-objective misalignment.

Experimental results on a real-world Chengdu dataset demonstrate that RS-PPO outperforms traditional heuristics (SA) and state-of-the-art RL models (PPO-noRS, ADQN) across key metrics: 18.2% higher OTR, 12.7% lower ADC, 29.8% higher VUR, and 40.3% faster convergence. Ablation studies confirm that each RS component contributes to specific dispatch priorities, validating the multi-dimensional design.

For logistics enterprises, RS-PPO provides a practical tool to enhance operational efficiency, reduce costs, and align with sustainability goals. Future work will focus on expanding to multi-depot scenarios and integrating weather/traffic prediction to further improve policy robustness.

References

- 1. K. C. Tan, L. H. Lee, Q. L. Zhu, and K. Ou, "Heuristic methods for vehicle routing problem with time windows," *Artificial intelligence in Engineering*, vol. 15, no. 3, pp. 281-295, 2001.
- 2. L. Shi, Z. Xu, M. Lejeune, and Q. Luo, "An integer l-shaped method for dynamic order fulfillment in autonomous last-mile delivery with demand uncertainty," *arXiv* preprint arXiv:2208.09067, 2022.
- 3. G. Clarke, and J. W. Wright, "Scheduling of vehicles from a central depot to a number of delivery points," *Operations research*, vol. 12, no. 4, pp. 568-581, 1964. doi: 10.1287/opre.12.4.568
- 4. G. B. Dantzig, and J. H. Ramser, "The truck dispatching problem," *Management science*, vol. 6, no. 1, pp. 80-91, 1959. doi: 10.1287/mnsc.6.1.80
- 5. F. Glover, "Tabu search-part I," ORSA Journal on computing, vol. 1, no. 3, pp. 190-206, 1989. doi: 10.1287/ijoc.1.3.190
- 6. P. Toth, and D. Vigo, "Vehicle routing: problems, methods, and applications," Society for industrial and applied mathematics, 2014.

- 7. W. B. Smythe, "Static and dynamic electricity," 1988.
- 8. M. Silva, and J. P. Pedroso, "Deep reinforcement learning for crowdshipping last-mile delivery with endogenous uncertainty," *Mathematics*, vol. 10, no. 20, p. 3902, 2022. doi: 10.3390/math10203902
- 9. A. Y. Ng, D. Harada, and S. Russell, "Policy invariance under reward transformations: Theory and application to reward shaping," In *Icml*, June, 1999, pp. 278-287.
- 10. H. N. Psaraftis, M. Wen, and C. A. Kontovas, "Dynamic vehicle routing problems: Three decades and counting," *Networks*, vol. 67, no. 1, pp. 3-31, 2016. doi: 10.1002/net.21628
- 11. J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," arXiv preprint arXiv:1707.06347, 2017.
- 12. J. E. Muriel, L. Zhang, J. C. Fransoo, and J. G. Villegas, "A reinforcement learning framework for improving parking decisions in last-mile delivery," *Transport metrica B: Transport Dynamics*, vol. 12, no. 1, p. 2337216, 2024. doi: 10.1080/21680566.2024.2337216
- 13. A. Bdeir, S. Boeder, T. Dernedde, K. Tkachuk, J. K. Falkner, and L. Schmidt-Thieme, "RP-DQN: An application of Q-learning to vehicle routing problems," In *German conference on artificial intelligence (Künstliche Intelligenz)*, September, 2021, pp. 3-16. doi: 10.1007/978-3-030-87626-5 1
- 14. H. Lee, and J. Jeong, "Mobile robot path optimization technique based on reinforcement learning algorithm in warehouse environment," *Applied sciences*, vol. 11, no. 3, p. 1209, 2021. doi: 10.3390/app11031209
- 15. W. K. Anuar, L. S. Lee, H. V. Seow, and S. Pickl, "A multi-depot dynamic vehicle routing problem with stochastic road capacity: An MDP model and dynamic policy for post-decision state rollout algorithm in reinforcement learning," *Mathematics*, vol. 10, no. 15, p. 2699, 2022.
- 16. M. I. D. Ranathunga, A. N. Wijayanayake, and D. H. H. Niwunhella, "Solution approaches for combining first-mile pickup and last-mile delivery in an e-commerce logistic network: A systematic literature review," In 2021 International Research Conference on Smart Computing and Systems Engineering (SCSE), September, 2021, pp. 267-275. doi: 10.1109/scse53661.2021.9568349
- 17. K. V. Tiwari, and S. K. Sharma, "An optimization model for vehicle routing problem in last-mile delivery," *Expert Systems with Applications*, vol. 222, p. 119789, 2023. doi: 10.1016/j.eswa.2023.119789
- S. Wang, T. Kong, B. Guo, L. Lin, and H. Wang, "CourIRL: Predicting Couriers' Behavior in Last-Mile Delivery Using Crossed-Attention Inverse Reinforcement Learning," In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, October, 2024, pp. 4957-4965. doi: 10.1145/3627673.3680046

Disclaimer/Publisher's Note: The views, opinions, and data expressed in all publications are solely those of the individual author(s) and contributor(s) and do not necessarily reflect the views of PAP and/or the editor(s). PAP and/or the editor(s) disclaim any responsibility for any injury to individuals or damage to property arising from the ideas, methods, instructions, or products mentioned in the content.