

# Deep Reinforcement Learning for Dynamic Origin-Destination Matrix Estimation in Microscopic Traffic Simulations Considering Credit Assignment

Donggyu Min<sup>†</sup>, Seongjin Choi<sup>‡</sup>, and Dong-Kyu Kim<sup>†,\*</sup>

**Abstract**—This paper focuses on dynamic origin-destination matrix estimation (DODE), a crucial calibration process necessary for the effective application of microscopic traffic simulations. The fundamental challenge of the DODE problem in microscopic simulations stems from the complex temporal dynamics and inherent uncertainty of individual vehicle dynamics. This makes it highly challenging to precisely determine which vehicle traverses which link at any given moment, resulting in intricate and often ambiguous relationships between origin-destination (OD) matrices and their contributions to resultant link flows. This phenomenon constitutes the credit assignment problem, a central challenge addressed in this study. We formulate the DODE problem as a Markov Decision Process (MDP) and propose a novel framework that applies model-free deep reinforcement learning (DRL). Within our proposed framework, the agent learns an optimal policy to sequentially generate OD matrices, refining its strategy through direct interaction with the simulation environment. The proposed method is validated on the Nguyen-Dupuis network using SUMO, where its performance is evaluated against ground-truth link flows aggregated at 5-minute intervals over a 30-minute horizon. Experimental results demonstrate that our approach achieves a 43.2% reduction in mean squared error (MSE) compared to the best-performing conventional baseline. By reframing DODE as a sequential decision-making problem, our approach addresses the credit assignment challenge through its learned policy, thereby overcoming the limitations of conventional methods and proposing a novel framework for calibration of microscopic traffic simulations.

**Index Terms**—Dynamic Origin-Destination Matrix Estimation, Reinforcement Learning, Microscopic Traffic Simulation, Calibration Problem.

## 1 INTRODUCTION

ADVANCES in data acquisition and computational power have enabled sophisticated, computationally intensive transportation strategies [1]. Modern strategies, such as speed harmonization, lane-changing control, and adaptive traffic signal control, are now predominantly designed and evaluated at the microscopic level, marking a significant shift in traffic management paradigms [2], [3], [4]. This emphasis on microscopic analysis underscores the importance of microscopic traffic simulations.

The reliability of microscopic traffic simulation depends on a rigorous calibration process that replicates real-world traffic dynamics [5], [6], [7]. This process typically involves minimizing a loss function that quantifies the discrepancy between simulation outputs and field observations [8]. A critical component of calibration is the dynamic origin-destination matrix estimation (DODE) problem, which aims to adjust traffic demand inputs to reproduce observed data, such as link flows [9].

A conventional approach to the DODE problem involves formulating a bi-level optimization problem [10], [11], [12].

---

*This is a preprint. The final version may differ upon publication.*

<sup>†</sup>D. Min and D.-K. Kim are with the Department of Civil and Environmental Engineering, Seoul National University, Seoul 08826, Republic of Korea. D.-K. Kim is also with the Institute of Construction and Environmental Engineering, Seoul National University, Seoul 08826, Republic of Korea. (E-mail: dgmin@snu.ac.kr; dongkyukim@snu.ac.kr).

<sup>‡</sup>S. Choi is with the Department of Civil, Environmental, and Geo-Engineering, University of Minnesota, Minneapolis, MN 55455, USA. (E-mail: chois@umn.edu).

\*Corresponding Author.

Manuscript received November 9, 2025.

This consists of an upper-level problem to determine the optimal origin-destination (OD) matrix and a lower-level problem, typically a dynamic traffic assignment model, that maps the OD matrix to link flows. While this bi-level approach offers computational tractability and modeling flexibility [11], applying it to calibrate microscopic simulations presents considerable challenges.

The primary challenge stems from the nature of microscopic dynamics, where the impacts of sequential OD matrices on link flows are not independent across time but are temporally complex and highly stochastic. In microscopic simulations, dynamics are modeled at the individual vehicle level. This granularity, combined with complex driver interactions, makes it infeasible to deterministically trace future vehicle trajectories [9], [13], [14], [15]. These interactions induce unpredictable network loading patterns, which in turn make it difficult to attribute observed link flows to the specific, time-lagged OD inputs—a fundamental issue known as the credit assignment problem [16], [17].

To address this credit assignment problem, we reformulate the DODE problem as a sequential decision-making problem using a Markov Decision Process (MDP) framework. The MDP explicitly accounts for long-term stochastic impacts of sequential OD matrices, offering a structured approach to the credit assignment issue that conventional optimization methods often overlook.

However, solving the DODE problem as an MDP introduces two technical challenges. The first is the combination of high-dimensional states and intractable transition dynamics, which renders traditional MDP solvers, such as

dynamic programming, computationally infeasible due to the curse of dimensionality and the lack of an explicit transition model [18]. To overcome this issue, we employ a deep reinforcement learning (DRL) approach, which enables an agent to learn an optimal policy through direct interaction with the simulation environment.

The second challenge is the instability of policy training due to a large action space, stochastic noise, and delayed rewards [19]. To address this, a specialized DRL approach is required. A large, combinatorial action space complicates exploration, while stochastic noise and delayed rewards can lead to unstable policy gradients and slow convergence. We employ proximal policy optimization (PPO) [20] with a multi-binary action parameterization. Structuring the action space as multi-binary reduces combinatorial complexity by letting the agent decide whether to release demand for each OD pair at each interval. In addition, PPO's clipped surrogate objective implicitly limits the update magnitude, mitigating instability induced by stochasticity. PPO's actor-critic architecture, combined with generalized advantage estimation (GAE) [21], reduces variance and facilitates effective long-horizon credit assignment by integrating immediate rewards with value function estimates. The main contributions of this paper are as follows:

- We frame the challenge of DODE in microscopic traffic simulation as a *credit assignment problem* and propose an *MDP framework* to formalize it.
- We propose a *model-free DRL approach* to handle high-dimensional observations and intractable transition dynamics inherent in this MDP formulation.
- We employ a *multi-binary PPO* to manage the combinatorial action space and achieve robust convergence under stochasticity and delayed rewards.

## 2 LITERATURE REVIEW

### 2.1 Conventional Approaches to Dynamic Origin-Destination Matrix Estimation

The estimation of OD matrices has traditionally been studied in both planning and operational contexts. In the planning stage, researchers estimate travel demand for a given region to predict unknown travel patterns. In contrast, the operational stage aims to calibrate demand inputs to align outputs from traffic models, such as simulations, with observed ground-truth data. The focus of this study is on the operational stage. Specifically, we aim to identify the optimal input trajectory of OD matrices for a microscopic model, such that the generated link flows closely align with observed ground-truth data.

Initial research on estimating OD matrices from observed link flows established foundational concepts and methodologies that informed subsequent developments [22]. To reproduce observed link flows requires distributing trips defined by the OD matrix onto specific routes within the network. Initially, this problem was defined using a time-independent assignment matrix to relate the OD matrix to link flows and formulated as a single-objective optimization problem [23]. As research advanced, the DODE problem emerged, aiming to capture realistic driver behavior by dividing the analysis period into shorter intervals

and dynamically distributing time-dependent OD matrices along network paths. However, using a fixed assignment matrix proved inadequate for accurately representing the dynamic and temporally varying relationship between OD matrices and observed link flows.

To address this issue, researchers expanded traditional formulations. Cremer and Keller [24] formulated the DODE problem as a constrained optimization problem aiming to minimize the squared errors between estimated and observed link flows to determine time-dependent OD matrices. Cascetta et al. [25] treated the problem as an optimization problem based on dynamic traffic assignment models, using time-dependent route choice proportions to approximate the relationship between OD matrices and observed link flows. Ashok and Ben-Akiva [26] applied a state-space model to estimate OD matrices by modeling their deviations from historical data. Bierlaire and Crittin [27] formulated the DODE problem as a single optimization problem and compared the performance of the Kalman filter with the more efficient LSQR algorithm for its solution.

### 2.2 The Bi-level Optimization Framework

An alternative approach to solving the DODE problem is to formulate it as a bi-level optimization problem. To handle real-world congestion effects and measurement errors, Yang et al. [10] argued that OD matrix estimation should be integrated with equilibrium traffic assignment, rather than being treated separately. This formulation is based on the concept that upper-level decisions influence the lower-level problem, whose outcomes subsequently inform the upper-level evaluation [28]. The upper-level problem aims to find a sequence of OD matrices that minimizes discrepancies between observed and simulated data, whereas the lower-level problem is the DTA problem, which involves assigning trips based on the equilibrium travel times experienced by users [9], [29].

The bi-level optimization formulation is a well-established framework for addressing the DODE problem [11]. In contrast to methods seeking a mutually consistent solution, the bi-level formulation models the lower-level's optimal response within the upper-level's decision-making. Its anticipatory nature guides the system toward a more globally optimal state, yielding superior solutions [30]. This decomposition also provides methodological flexibility by allowing specialized algorithms for each level.

Recently, various approaches based on this formulation have emerged. For the upper-level problem, algorithms such as simultaneous perturbation stochastic approximation (SPSA) and its variations, as well as Bayesian optimization (BO) and its variants, have been developed [9], [31]. Meanwhile, advances in computational technology have enabled widespread use of traffic simulations to capture realistic traffic dynamics in solving the lower-level problem [14], [29]. Furthermore, researchers have developed more accurate and efficient DODE methods by leveraging additional data sources [12], [32], [33], [34], advancing algorithmic complexity [35], and employing surrogate models [7], [36], [37]. Comprehensive reviews of recent approaches can be found in the papers of Osorio [36] and Huo et al. [9].

### 2.3 Research Gap and Proposed Direction

Despite the advancements, applying these conventional frameworks to calibrate microscopic simulations reveals a fundamental challenge. As previously noted, the temporally complex and stochastic relationships between sequential OD matrices and resulting link flows create a significant credit assignment problem. Existing methods based on a bi-level framework are categorized into two approaches: simultaneous and sequential optimization [34]. Simultaneous optimization handles temporal complexity and stochasticity implicitly but suffers from an immense computational burden as the dimension of the decision variables increases. Conversely, sequential optimization reduces complexity but is often myopic, optimizing within finite horizons and failing to capture the long-term consequences of decisions. This creates a critical trade-off between computational complexity and analytical myopia, limiting the effectiveness of current solutions for the credit assignment problem.

To resolve this trade-off, this study reframes the DODE problem as a sequential decision-making task, formally structured within the MDP formulation. While partially observable formulations (POMDPs) could be considered to model unmeasured microscopic states, they require belief-state estimation and substantially increase computational complexity in learning and planning [38], [39]. In this study, we adopt a standard MDP with carefully designed state augmentation and short input timesteps, which we find sufficient to capture long-horizon effects while remaining computationally tractable. This approach aligns with a growing body of research that successfully applies standard MDPs to complex, stochastic transportation problems [40].

Conceptually, the MDP framework synthesizes the strengths of both prior approaches. It decomposes the problem into a series of decisions over time, like sequential optimization, thereby avoiding the exponential growth in complexity that simultaneous methods face. However, unlike myopic sequential methods, the MDP's objective is to maximize a long-term, cumulative reward. This formulation inherently forces the model to account for the delayed, downstream consequences of current actions, thus formally addressing the credit assignment problem and overcoming analytical myopia. The MDP, therefore, provides a theoretically sound pathway to achieving a long-term optimal perspective with the tractability of a sequential process.

However, while the MDP offers a robust conceptual solution, its practical implementation is challenged by microscopic simulations. The state space is extraordinarily high-dimensional—a classic challenge known as the curse of dimensionality—and the state transition dynamics are intractable and stochastic [41]. Consequently, traditional MDP solvers, such as dynamic programming, are computationally infeasible. To bridge this gap between theory and practice, we employ a model-free DRL approach. DRL has emerged as the state-of-the-art methodology for solving complex MDPs precisely because it uses deep neural networks as function approximators to handle high-dimensional states and can learn effective policies from direct interaction without an explicit transition model [40]. DRL thus provides the necessary tools to solve the formulated MDP, enabling an agent to learn an optimal policy

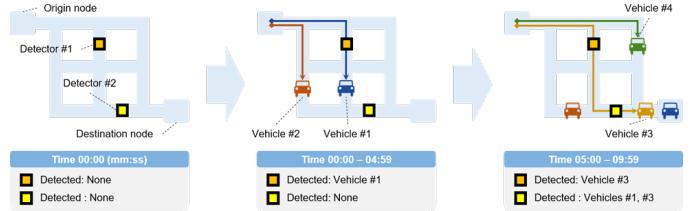


Fig. 1. Comparison of detector data for vehicles with the same OD pair

within the complex environment.

By integrating the MDP framework with a model-free DRL approach, we demonstrate a practical path to overcoming the trade-off between complexity and myopia. This synthesis establishes a novel paradigm, enabling a previously intractable level of granularity in demand estimation and opening new avenues for high-fidelity model calibration and control.

## 3 METHODS

### 3.1 Key Concept

This study estimates sequences of OD matrices for the offline calibration of microscopic traffic simulations. The primary objective is to simulate detector data as link flows aggregated over 5-minute intervals, ensuring they closely match ground-truth observations. A fundamental challenge arises from the temporally complex and stochastic relationships inherent in vehicle-level dynamics.

The input-output relation is illustrated in **Figure 1**, which shows how vehicles with the same OD pair can influence detector data differently. For instance, vehicles #1 and #2 depart simultaneously but take different paths, resulting in their detection at separate detectors. Conversely, vehicles #1 and #3 depart at different timesteps yet are recorded within the same timestep due to variable network conditions, such as congestion. Any detector may not capture vehicle #4 if its chosen trajectory does not traverse a monitored link, or if it exits the network or remains in transit after the simulation. These complex and stochastic outcomes make it analytically intractable to attribute observed link flows back to specific departure decisions, which constitutes the core of the credit assignment problem [19].

To address this challenge, we reformulate the DODE problem as a sequential decision-making task. This task is formally structured using an MDP framework [19]. We represent the network dynamics, driven by vehicle departures and their subsequent travel, as a discrete-time system with sufficiently short timesteps. This discrete representation allows the Markov property, which posits that the subsequent state depends only on the current state and the chosen action, to serve as a reasonable approximation of real-world traffic dynamics. By explicitly modeling the probabilistic nature of state transitions, this formulation effectively captures the temporally complex and stochastic nature inherent in the simulation.

Although the MDP provides a suitable theoretical framework, its direct solution is computationally intractable due to the high-dimensional state space and unknown state transition probabilities. We therefore employ a model-free

TABLE 1  
Notations

Time	
$\mathcal{T}$	Set of input time index ( $= \{1, 2, \dots, T\}$ )
$T$	Maximum index of inputs
Action	
$\mathcal{A}$	Action space
$a$	Input trajectory covering the entire period ( $= \{a_t\}_{t=1}^T, a_t \in \mathcal{A}$ )
$a^*$	Optimal input trajectory
$a_t$	OD matrix at $t$
State	
$\mathcal{S}$	State space
$s$	State trajectory for all timesteps ( $= \{s_t\}_{t=0}^T, s_t \in \mathcal{S}$ )
$s_t$	Traffic network and context state at $t$
$s_0$	Initial state before any input
$N_l$	Number of vehicles for link $l$
$\bar{v}_l$	Average speed for link $l$
$v_f$	Free-flow speed
Detector data	
$\mathcal{K}$	Set of output time index ( $= \{1, 2, \dots, K\}$ )
$K$	Maximum index of outputs
$d$	Ground-truth detector data covering the entire period ( $= \{d_k\}_{k=1}^K, d_k \in \mathbb{R}^K$ )
$d_k$	Ground-truth detector data at $k$
$d'$	Simulated detector data covering the entire period ( $= \{d'_k\}_{k=1}^K, d'_k \in \mathbb{R}^K$ )
$d'_k$	Simulated detector data at $k$
Functions	
$\mathcal{L}$	Loss function corresponding to the calibration error
$\mathcal{M}$	Mapping function, which is a rule that defines the temporal relationship between the $\mathcal{T}$ and $\mathcal{K}$
$\mathcal{P}(\cdot)$	Power set of an arbitrary set
$(T > K)$	$\mathcal{M}: \mathcal{K} \rightarrow \mathcal{P}(\mathcal{T}), \mathcal{M} = \psi(k)$ . When the number of inputs is greater than the number of outputs, map the output index set to the power set of the input index set. Therefore, return the set of input indices $t$ that contribute to output index $k$ .
$\mathcal{F}$	Microscopic simulation model structured by the mapping function $\mathcal{M}$ for generating the simulated detector data
$p$	State transition function
$r$	Reward function
$\gamma$	Discount factor

DRL approach, allowing an agent to learn an optimal policy through direct interaction with the simulation environment, as shown in **Figure 2**. A policy maps the current network state to a dispatch decision for each OD pair. The agent's learning is driven by maximizing the long-term cumulative reward, defined as the negative error between simulated and ground-truth data. The policy with the best performance in reproducing the ground-truth data during training is selected, and the demand input trajectory it generates constitutes our final output. This trajectory consists of a sequence of binary departure decisions for each OD pair at every timestep. The subsequent sections detail the specific MDP formulation and the DRL algorithm employed.

### 3.2 Problem Statement

#### 3.2.1 Notations

Before presenting the mathematical formulation, we define the key notations used throughout this section. For clarity of reference, notations are summarized in **Table 1**.

#### 3.2.2 Dynamic origin-destination matrix estimation problem

The DODE problem seeks to determine an input sequence of OD matrices that best reproduces a given ground-truth dataset. This ground-truth dataset comprises detector data, including link flows, collected at discrete timesteps. In this study, detector data refers to vectors composed of link flows collected every 5 minutes at specific links. The OD matrix  $a_t$  at the timestep  $t$  represent the flow rates for each OD

pair trip starting at the origin node  $i$  and arriving at the destination node  $j$ , which can be expressed as a vector,

$$a_t = [a_{(1,1)}, a_{(1,2)}, \dots, a_{(i,j)}, \dots, a_{(\max(I),\max(J))}] \quad (1)$$

where  $a_t$  is the vector of origin-destination matrix at timestep  $t$ ,  $a_{(i,j)}$  is the flow rate for trips with origin  $i$  to destination  $j$ ,  $I$  is the set of indices of origin nodes ( $i \in I$ ), and  $J$  is the set of indices of destination nodes ( $j \in J$ ).

This study addresses an inverse problem for input trajectory estimation in microscopic simulation models. Given a known initial state and a set of empirical detector data, the objective is to deduce the optimal input trajectory,  $a^*$ , that minimizes the deviation between the simulation output and the observed data. An objective function quantifies this deviation,  $\mathcal{L}$ , defined as the sum of squared Euclidean distances over the entire simulation horizon as follows:

$$\min_{a \in \mathcal{A}} \mathcal{L}(a, s_0, d) = \|d' - d\|_2^2 = \sum_{k=1}^K \|d'_k - d_k\|_2^2, \quad (2)$$

subject to  $d'_k \sim \mathcal{F}(a | s_0)$  for all  $k \in \mathcal{K}$ , with a given initial state  $s_0 \in \mathcal{S}$ .

In this paper, we formulate the DODE problem within an MDP framework, defining it as a process of sequentially determining the optimal departure decision for each OD pair given the current network state. An MDP is formally defined by a tuple  $\langle \mathcal{S}, \mathcal{A}, p, r, \gamma \rangle$ , comprising a state space  $\mathcal{S}$ , an action space  $\mathcal{A}$ , a state transition function  $p$ , a reward function  $r$ , and a discount factor  $\gamma \in [0, 1]$ . This MDP-based approach allows for the optimization of long-term outcomes, thereby addressing the myopic nature of sequential methods. Furthermore, by learning a policy instead of optimizing the entire trajectory at once, it circumvents the complexity inherent in the simultaneous formulation.

#### 3.2.3 State

In MDP, a state is a complete description of the agent and environment at a specific point in time. In this study, a state consists of two elements: network state and context state. First, to maintain dimensional consistency, we propose a link-based network state, which is represented as a 1-dimensional vector containing the number of vehicles  $N_l$  and the average speed  $\bar{v}_l$  for each link. Additionally, a context state is proposed to provide agents with clues regarding their goals. The context state consists of the current timestep at which an action should be taken and detector data information indicating the number of vehicles that have passed through the detectors. This detector data information is initialized as a zero vector whenever the aggregation of detector data is completed and updated.

An example of a state is shown in **Figure 3**. This figure explains how to construct the state  $s_t$  observed at the point when the agent needs to take action  $a_t$  at timestep  $t$ .  $s_t$  depends entirely on the last snapshot in the simulation at timestep  $t - 1$ , which is the result of vehicles moving along time-dependent shortest paths during the input interval. We extract the network state and context state from the last snapshot, then flatten and concatenate them to create  $s_t$ . At this point, the average speed of links with no vehicles is assumed to be the free-flow speed  $v_f$ . In this example, the network consists of 12 links, and two detectors are installed;

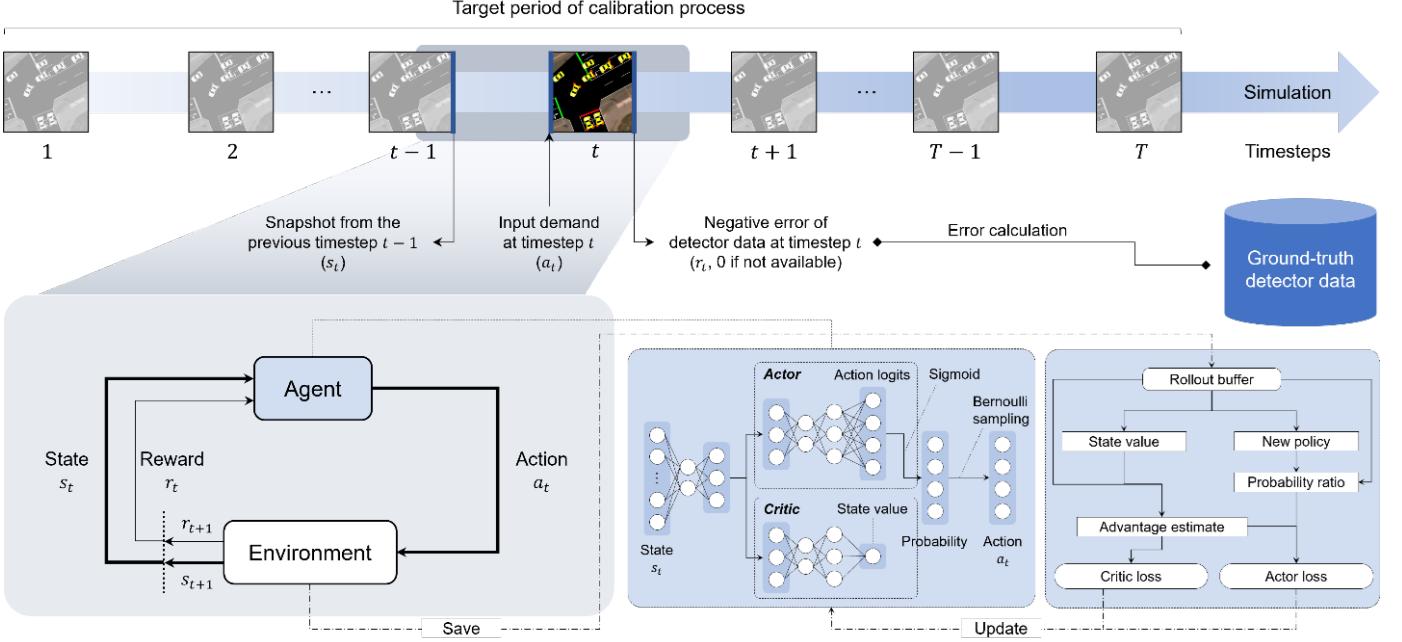


Fig. 2. Overall framework of the proposed method

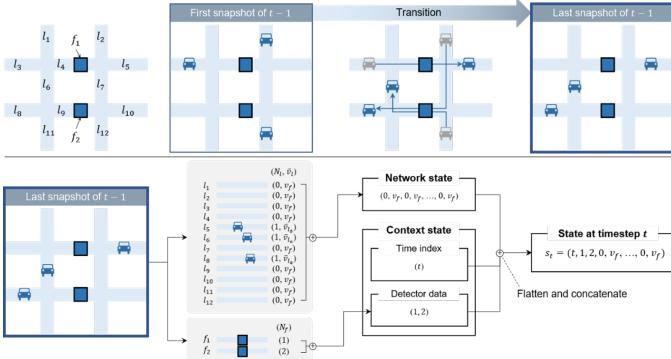


Fig. 3. State description

therefore, the dimension of the network state is 24, and the dimension of the context state is 3. Therefore, the dimension of the state is 27.

### 3.2.4 Action

In the proposed method, the input interval is set to a short interval of 5 seconds. This reduces the probability of unobservable events occurring between the current state and the next state, thereby improving the stability of the MDP framework. According to this experimental setup, the agent only needs to decide whether to dispatch a vehicle for each OD pair in the current step. All elements of the OD matrix  $a_t$  take values of 0 or 1, allowing it to be represented as a flattened binary vector.

### 3.2.5 State transition probability

The state transition probability  $p(s_{t+1} | s_t, a_t)$  is a function that represents the probability of transitioning to the next state  $s_{t+1}$  when a specific action  $a_t$  is taken from the current state  $s_t$ . In the DODE problem, this function describes how

the traffic conditions change when a new OD matrix is input under given traffic conditions. This transition process is performed using microscopic traffic simulations. Therefore, the transition function is not explicitly known and is treated as a black-box function. This is the reason why a model-free reinforcement learning approach is necessary.

### 3.2.6 Reward

In this problem, the reward  $r_t$  is defined by the difference between the simulated link flow vector  $d'_k$  and ground-truth link flow vector  $d_k$ . Since the input interval is set to be short, our formulation corresponds to the case where  $T > K$ . The reward function can be written as:

$$r_t = \begin{cases} 0 & (t \neq \max(\psi(k))) \\ -\|d'_k - d_k\|_2^2 & (t = \max(\psi(k))) \end{cases} \quad (3)$$

subject to  $d'_k \sim \mathcal{F}(a_{\psi(k)} | s_{\min(\phi(k))-1})$  with a given state  $s_{\min(\psi(k))-1} \in \mathcal{S}$  and  $a_{\psi(k)} \subset \mathcal{A}$ .

Note that, in our framework, the input interval and the output interval generally differ, so we introduce a mapping function  $\mathcal{M}$ . For example, when the input period is 5 seconds and the output period is 300 seconds, the mapping function is  $\mathcal{M} = \psi(k=1) = \{t \in \mathbb{Z} \mid 1 \leq t \leq 60\}$ . In this case, the reward  $r_t$  is 0 except at  $r_{60}$ , when the output is generated. To align the reinforcement learning objective of maximizing cumulative reward with our aim of minimizing this discrepancy, we simply define the reward as the negative of the error.

### 3.2.7 Discount factor

The discount factor  $\gamma$  is a value between 0 and 1 that determines the present value of future rewards. When  $\gamma = 0$ , the agent exhibits myopic behavior, considering only immediate rewards. Conversely, when  $\gamma \rightarrow 1$ , the agent exhibits a far-sighted attribute, valuing future rewards as much as current

rewards. In this paper, we use  $\gamma = 0.995$  to improve the performance of the proposed method and ensure convergence.

### 3.3 Deep Reinforcement Learning for DODE Problem

The agent's behavior is described by a policy  $\pi(a_t | s_t)$ , which maps each state  $s_t$  to a probability distribution over actions  $a_t$ . The goal of reinforcement learning is to find an optimal policy  $\pi^*$  that maximizes the expected sum of discounted rewards, also known as the return. This objective can be formally written as:

$$\pi^* = \arg \max_{\pi} J(\pi) = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=1}^T \gamma^{t-1} r_t \right] \quad (4)$$

where  $\tau = (s_0, a_1, s_1, a_2, \dots)$  is a trajectory generated by following policy  $\pi$ .

To solve this problem, we employ multi-binary PPO, an adaptation of the PPO algorithm featuring a factorized Bernoulli policy head. PPO is a prominent model-free DRL algorithm [20]. This algorithm directly parameterizes and optimizes the policy  $\pi_{\theta}(a_t | s_t)$ , where  $\theta$  represents the policy's parameters. PPO utilizes an actor-critic architecture. The actor, implemented as the policy network  $\pi_{\theta}(a_t | s_t)$ , selects actions. The critic, represented by a value network  $V_{\phi}(s_t)$  with its parameters  $\phi$ , evaluates the quality of a state by estimating its expected return (the state-value).

We model the action as a multi-binary vector: the actor maps each state to an  $n$ -dimensional logits vector, applies a sigmoid to obtain per-component probabilities, and samples each OD pair via a Bernoulli distribution. This defines a factorized policy whose components are conditionally independent given the state. While one could treat  $\{0, 1\}^n$  as a single categorical space with  $2^n$  joint actions, that approach incurs a combinatorial explosion in parameterization and exploration. By contrast, the proposed factorized Bernoulli head scales linearly with the number of OD pairs in terms of the policy's output dimension and parameter count, yielding superior scalability. Although factorization assumes conditional independence, a shared backbone can still consider implicit dependencies through shared features.

To guide the actor's learning, PPO uses an advantage estimate,  $\hat{A}_t$ , which quantifies how much better a given action  $a_t$  is compared to the average action in the state  $s_t$ . We compute this advantage using GAE, which provides a low-variance estimate by balancing immediate rewards with long-term value predictions from the critic. For improved stability, the advantages are normalized per batch.

In PPO, the actor and critic are often trained jointly by optimizing a single, composite objective function. This objective is composed of three main components: the clipped surrogate objective for the policy ( $L^{CLIP}$ ), an error term for the value function ( $L^{VF}$ ), and an entropy bonus to encourage exploration. These are formulated as follows:

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip} \left( r_t(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) \right] \quad (5)$$

$$L^{VF}(\phi) = \mathbb{E}_t \left[ \left( V_{\phi}(s_t) - V_t^{\text{target}} \right)^2 \right] \quad (6)$$

$$L^{PPO}(\theta, \phi) = -\mathbb{E}_t \left[ L^{CLIP}(\theta) - c_1 L^{VF}(\phi) + c_2 \mathcal{H}[\pi_{\theta}(\cdot | s_t)] \right] \quad (7)$$

where  $\hat{A}_t$  is the advantage estimate,  $r_t(\theta)$  is the probability ratio between the new and old policies,  $\epsilon$  is the clipping range,  $V_t^{\text{target}}$  is the target value for the critic calculated as  $V_t^{\text{target}} = \hat{A}_t + V_{\phi}(s_t)$ ,  $\mathcal{H}[\pi]$  is the policy entropy, and  $c_1$  and  $c_2$  are the coefficients for the value function loss and the entropy bonus, respectively.

This objective function discourages excessive policy updates by penalizing probability ratios  $r_t(\theta)$  that fall outside the  $[1 - \epsilon, 1 + \epsilon]$  interval. When PPO is applied to the DODE problem, the actor network learns a policy for OD vehicle departure decisions  $a_t$  based on the state  $s_t$ . The critic provides an evaluation of the states, and the resulting advantage function  $\hat{A}_t$  offers a concrete basis for credit assignment. Based on this, PPO progressively improves the demand generation policy through stable updates, aiming to minimize the detector data error over the entire analysis period.

The choice of PPO is deliberate and central to addressing the challenges of the microscopic DODE problem. First, the clipped surrogate objective ( $L^{CLIP}$ ) is critical for ensuring learning stability. Microscopic simulations are inherently stochastic, meaning the same action in a similar state can lead to different outcomes. The trust region imposed by the clipping mechanism prevents the policy from overreacting to this noise, resulting in more reliable and monotonic improvements, as observed in our experiments. Second, the use of GAE within the actor-critic framework directly tackles the long-horizon credit assignment problem. Since rewards are only provided every 5 minutes (based on aggregated link flows), while actions are taken every 5 seconds, GAE allows the agent to properly evaluate the long-term consequences of its fine-grained decisions, making it a highly suitable choice for this problem.

### 3.4 Experimental Settings

This section introduces the assumptions of the experiment. We constructed a toy experimental environment to evaluate the effectiveness of the proposed method. The Nguyen-Dupuis network, with 4 OD pairs, 13 nodes, and 19 directed links, is suitable for a simple experimental setup to compare the effectiveness of the methodology [42]. Each link is assumed to consist of a single lane. To ensure that the influence of the OD matrix is evident across a wide time range, we multiplied all link lengths by a factor of 3. We assume that virtual detectors are present on some of the links in the network shown in **Figure 4**. These 9 detectors count the number of vehicles passing through each link every 5 minutes.

SUMO was used for microscopic traffic simulation [43]. All hyperparameters are set to SUMO's default values. For example, all vehicles are implemented as passenger cars, and the Krauss model is used as the car-following model. Additionally, the acceleration is 2.6 m/s<sup>2</sup>, and the deceleration is 4.5 m/s<sup>2</sup>. Regarding the assignment, we update the time-dependent shortest path every 5 seconds for all vehicles and have them follow that path.

Ground-truth detector data is generated according to the following procedure. First, define the analysis period. In this study, we measure the flow of 9 links at 5-minute intervals for 30 minutes, resulting in a ground-truth dataset with 6

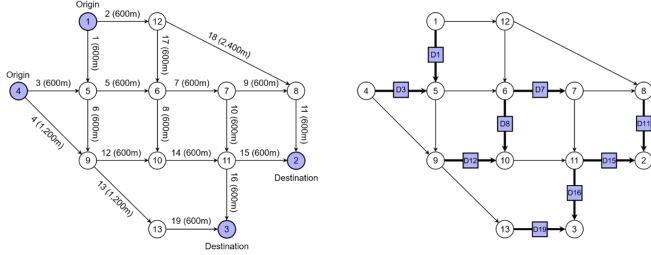


Fig. 4. Nguyen and Dupuis network

TABLE 2  
Ground-truth Detector Data

Time (mm:ss)	Link flow (veh/min)								
	D1	D11	D12	D15	D16	D19	D3	D7	D8
00:00-04:59	12	4	3	0	2	3	16	12	3
05:00-09:59	14	5	17	16	25	7	22	16	13
10:00-14:59	13	7	2	17	3	16	22	26	6
15:00-19:59	14	19	10	11	24	6	17	28	10
20:00-24:59	17	8	21	25	16	6	8	14	9
25:00-29:59	12	8	9	9	12	13	16	16	9

rows and 9 columns. Second, determine the total number of cars that will be input into the target network during the defined analysis period. To generate the ground-truth data, a total of 300 vehicles were created over the entire 30-minute analysis period. For each 5-second timestep, a departure decision was made probabilistically for each OD pair. As a result, the observed data that we need to fit through OD matrix estimation is shown in **Table 2**.

In this study, the comparison targets are broadly classified into three categories. The first method reproduces detector data using different random seeds based on the assumed true demand. The second is a simultaneous optimization method. The third is a sequential optimization method. To demonstrate the limitations of existing methods, we applied the conventional input interval of 5 minutes and the proposed method's input interval of 5 seconds, respectively, and performed benchmarking. For the optimization algorithm, Bayesian optimization (BO), which exhibits eminent performance in optimizing high-cost objective functions, was applied [44]. Detailed explanations of each method are provided in **Table 3**.

## 4 RESULTS

To evaluate the performance of the proposed approach, six methods were compared: true demand, proposed method (RL-PPO), simultaneous optimization method with a 5-minute interval (ST-BO (5min)), simultaneous optimization method with a 5-second interval (ST-BO (5sec)), sequential optimization method with a 5-minute interval (SQ-BO (5min)), and sequential optimization method with a 5-second interval (SQ-BO (5sec)). All methods were repeated 5 times in the identical setup. Since the research purpose is offline calibration, the points with the lowest error were primarily considered.

**Figure 5** displays the reward trend per episode of the proposed method. The bold line represents the average value, and the shaded area represents the range between the average value and the standard deviation. Dotted lines

TABLE 3  
Method Description

Methods					
Index	Abbreviation	Method	Algorithm	Iterations	Search space
1	True demand	-	-	-	-
2	RL-PPO	DRL	PPO		[0, 1] <sup>n</sup>
3	ST-BO (5min)	Simultaneous optimization			[0, 60] <sup>n</sup>
4	ST-BO (5sec)	Bayesian optimization		350	[0, 1] <sup>n</sup>
5	SQ-BO (5min)	Sequential optimization			[0, 60] <sup>n</sup>
6	SQ-BO (5sec)	Sequential optimization			[0, 1] <sup>n</sup>

Hyperparameters		
Algorithm	Hyperparameter	Value
PPO	Learning rate	3e-4
	Gamma	0.99
	Number of steps	180
	Lambda of GAE	0.95
BO	Entropy coefficient	0.01
	Number of initial points	10
	Acquisition function	Expected improvement
Gaussian process kernel	Gaussian process kernel	Matern (nu = 2.5)
	Gaussian process alpha	1e-6

\*Note: n is the number of OD pairs.

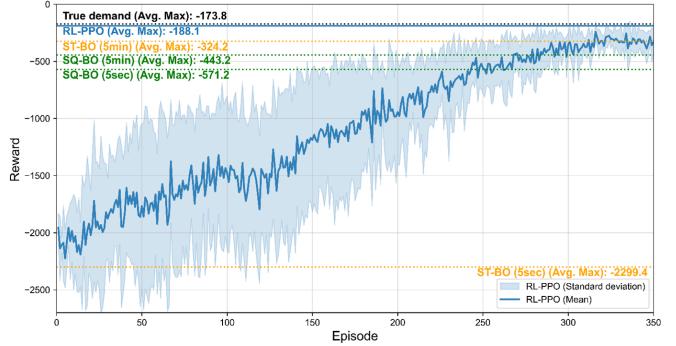


Fig. 5. Reward graph of the proposed method containing maximum values with other methods

show the average maximum reward for each methodology. The reward is non-zero with true demand due to microscopic dynamics resulting from different random seeds, leading to different vehicle route choice outcomes. Based on the average maximum reward, the results are ranked in descending order as follows: True demand -173.8, RL-PPO -188.1, ST-BO (5min) -324.2, ST-BO (5sec) -443.2, SQ-BO (5min) -571.2, SQ-BO (5sec) -2,299.4. Therefore, the proposed method achieved a 42.0% improvement in reward compared to the existing method with the best performance, ST-BO (5min).

As emphasized in the early part of this paper, results reveal various limitations of existing methods. First, ST-BO (5sec) showed the worst performance, contrasting with ST-BO (5min) showing the second-best. This highlights the complexity of the simultaneous optimization approach, where all decision variables are estimated simultaneously despite a 60-fold increase in dimension. Second, SQ-BO (5min) performed worse than ST-BO (5min), suggesting that, when the input dimension is small, exploring all decision variables simultaneously to avoid credit assignment is more appropriate than myopically estimating the O-D matrix. Finally, SQ-BO (5sec) and ST-BO (5sec) show that reducing the temporal input unit is disadvantageous. Our method solves the DODE problem in microscopic traffic simulation by using an MDP formulation and DRL, while

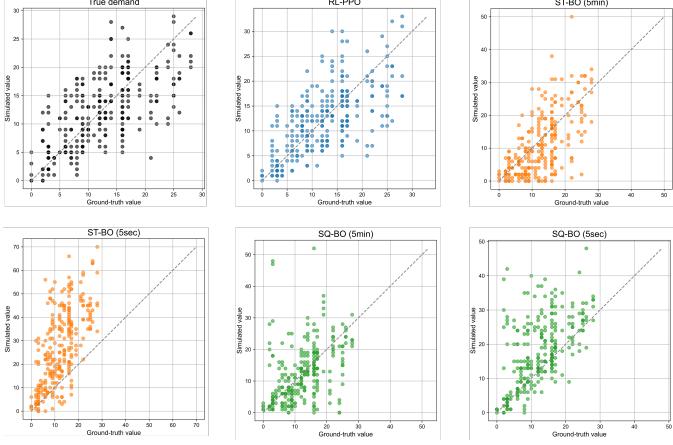


Fig. 6. Scatter plots of detector data using estimated demand by method

reducing the temporal input unit and achieving the best performance by adjusting demand at a higher resolution.

**Figure 6** presents scatter plots of data reproduced by the method. The x-axis of each subplot represents the ground-truth value, and the y-axis represents the simulated value.

**Figure 7** presents a set of performance and stability metrics, including MSE, root mean squared error (RMSE), mean absolute error (MAE), mean absolute percentage error (MAPE), standard deviation of error (SDE), 95th percentile of absolute error (P95 AE), maximum absolute error (MaxAE), mean bias error (MBE), and the coefficient of determination ( $R^2$ ).

Together, these indicators assess the accuracy and stability of the reproduced link flows across all evaluation points in each subplot. MSE and RMSE represent the mean magnitude of squared errors and their square roots, respectively, while MAE measures the average absolute deviation between simulated and ground-truth flows. MAPE represents the relative error as a percentage and was computed excluding zero ground-truth points to ensure numerical stability. SDE quantifies the variability of the error distribution, P95 AE and MaxAE represent the high-percentile and maximum deviations, and MBE indicates whether the model systematically over- or underestimates link flows. Finally,  $R^2$  measures how well the reproduced link flows explain the variation in the ground-truth data.

A notable observation is that microscopic traffic simulation is inherently stochastic, leading to considerable variations in link flows even under identical true-demand conditions ( $MSE = 28.96 \text{ (veh/5min)}^2$ ,  $MAPE = 45.36\%$ ). The proposed method achieves an MSE of  $30.69 \text{ (veh/5min)}^2$  and a MAPE of 47.28%, closely matching the true-demand case. These results indicate that our method reproduces link-level dynamics with accuracy and stochastic stability comparable to the ground-truth demand ( $SDE = 5.54$ ,  $R^2 = 0.447$ ). In contrast, conventional methods such as ST-BO (5min) and SQ-BO (5min) show much higher errors ( $MSE = 54.04$  and  $73.87 \text{ (veh/5min)}^2$ ,  $MAPE = 119.44\%$  and  $86.47\%$ , respectively), indicating that their estimated OD matrices frequently over- or under-estimate link flows.

Although simultaneous optimization implicitly accounts for complex temporal and stochastic dependencies by jointly adjusting all timestep inputs, its computational complexity

increases rapidly with the number of decision variables, reducing search efficiency and accuracy. Sequential optimization, while lowering the dimensionality of the problem, is inherently myopic and thus struggles to capture long-term temporal dependencies in stochastic environments.

In contrast, the MDP-based DRL framework proposed in this study performs adaptive, stepwise decision-making while explicitly evaluating the long-term consequences of each action. Such a framework allows more effective credit assignment in highly stochastic and temporally correlated traffic systems. Consequently, the proposed RL-PPO approach demonstrates robust accuracy and stability in microscopic environments. It achieves a 43.2% reduction in MSE compared to the best existing baseline (ST-BO (5min)) and yields performance comparable to the true-demand case in both the magnitude and variability of link-flow errors.

We evaluated the statistical differences between the five experimental groups based on a true demand-based dataset. Each dataset consists of 9 detectors, and all analyses were performed using paired samples. We calculated the differences between the points for each detector and applied the Shapiro-Wilk test with  $\alpha = 0.05$  to determine whether the data satisfied the normality assumption. If the differences satisfied the assumption of normality, we performed a paired t-test; otherwise, we applied the non-parametric Wilcoxon signed-rank test. We calculated the  $p$ -value for each test's null hypothesis and concluded that the results were statistically significant if the  $p$ -value was less than 0.05.

**Figure 8** shows the  $p$ -values by method and detector in a bar graph. Unlike other groups, the RL-PPO group showed  $p > 0.05$  for all detectors, indicating no statistical significance. Therefore, we can interpret that the link flows generated by our method were statistically indistinguishable from the ground truth for all detectors.

In real-world tasks, it is generally difficult to obtain the true OD demand values. Thus, link-flow fitting results are often used as an indirect indicator of estimation performance. However, since this study was conducted within a simulation environment, the true OD demand values were precisely known. This setting allows for a direct comparison between the ground-truth and estimated OD demands.

**Figure 9** illustrates the OD matrix estimation results, where the x-axis represents the ground-truth OD demand and the y-axis represents the estimated OD demand. Among all compared methods, the proposed RL-PPO approach achieved the lowest errors ( $MSE = 19.17 \text{ (veh/5min)}^2$ ,  $RMSE = 4.38 \text{ veh/5min}$ ,  $MAE = 3.50 \text{ veh/5min}$ ). The contour map indicates that the estimated OD demands are distributed symmetrically around the  $y = x$  reference line, suggesting that the proposed method neither systematically overestimates nor underestimates the OD demand.

**Figure 10** further visualizes the temporal evolution of the estimated OD matrices for each method across all time intervals and OD pairs. The proposed RL-PPO method exhibits a temporal pattern most consistent with the ground-truth demand, showing stable variations over time and across OD pairs. In contrast, conventional methods such as ST-BO and SQ-BO display irregular fluctuations and localized over- or underestimations, particularly in specific time intervals, implying limited adaptability to stochastic demand dynamics.

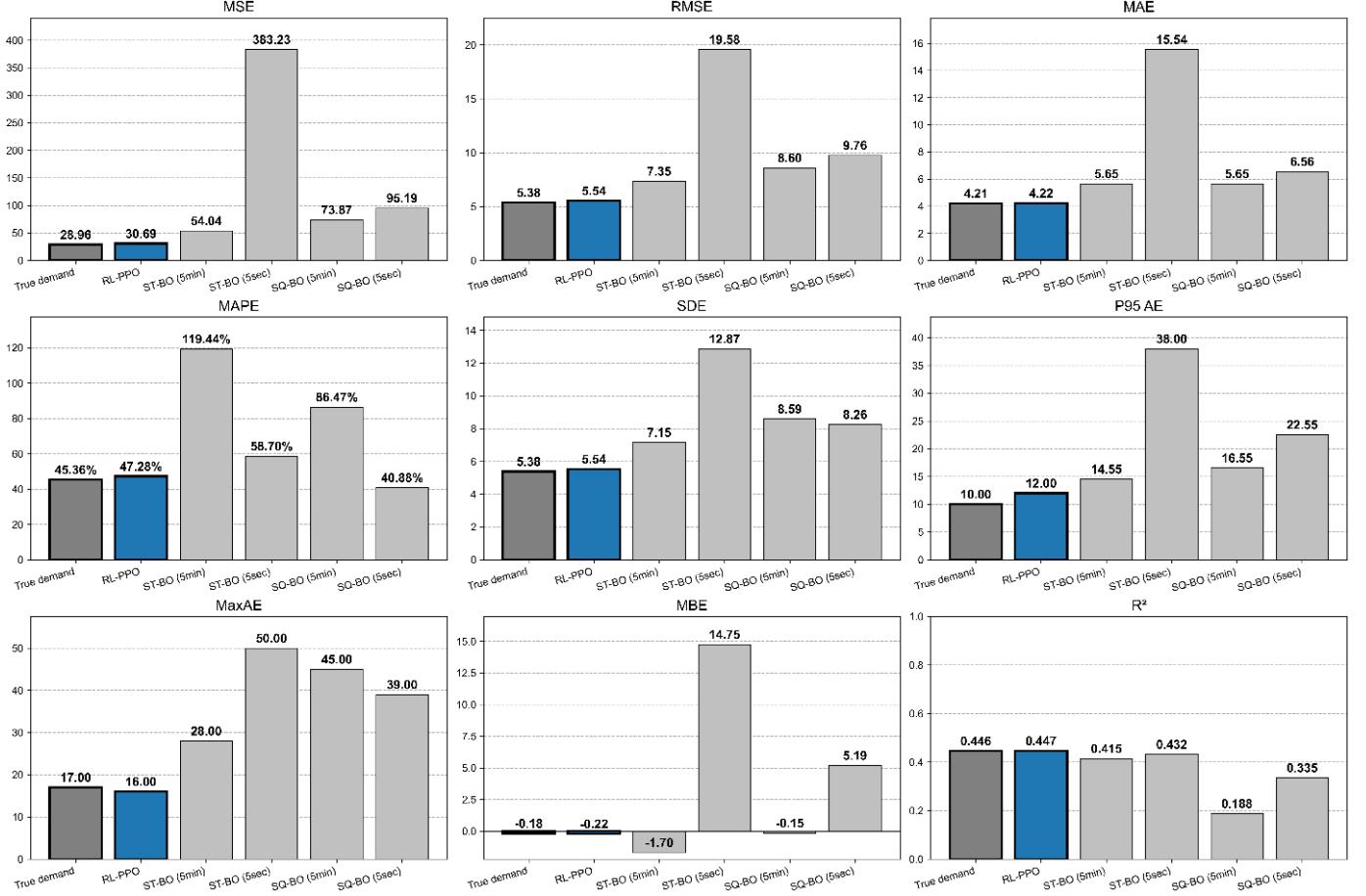


Fig. 7. Comparison of multi-metric performance and stability across methods

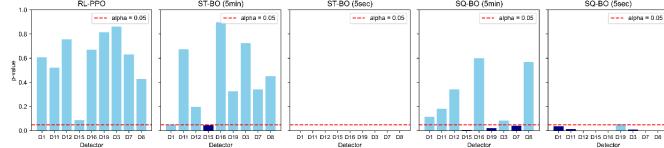


Fig. 8.  $p$ -values from tests of each method's link flow deviations from true demand

## 5 CONCLUSION

This study proposes a reinforcement learning-based approach to address the credit assignment problem, a persistent issue in DODE for microscopic traffic simulations. We address the limitations of conventional methods—namely, the complexity of simultaneous optimization and the myopia of sequential optimization—by redefining the DODE problem within an MDP framework and applying model-free DRL to derive optimal OD input sequences that consider long-term, stochastic effects.

Experimental results show that our method achieves significantly lower errors compared to existing simultaneous and sequential optimization approaches. It achieved an MSE reduction of approximately 43.2% compared to the best conventional method and demonstrated no statistically significant difference from the dataset generated using true demand. Additionally, this method confirmed that it can more

accurately represent realistic traffic dynamics by performing fine-grained departure decisions at the vehicle level. Our approach demonstrates a robust framework that effectively addresses the credit assignment problem while explicitly considering the characteristics of stochastic environments.

The core contribution of this study is a novel reinforcement learning-based approach to the DODE problem that overcomes the trade-off between high-dimensionality and myopic decision-making inherent in existing optimization-based methodologies. The proposed MDP and DRL framework addresses the complexity and stochasticity of microscopic simulations by evaluating the long-term impacts of each action, thereby resolving the credit assignment problem in DODE. The success of this approach is significantly bolstered by the choice of the multi-binary PPO, whose clipped surrogate objective ensures stable policy updates against the simulation's inherent noise, preventing drastic, erroneous decisions and fostering convergence toward the optimal solution. In particular, by enabling high-resolution decision-making structures that determine individual vehicle departure decisions every 5 seconds, the proposed approach captures realistic microscopic traffic dynamics and demonstrates superior calibration performance compared to existing methods. This approach is expected to overcome the methodological limitations of existing research and provide a fundamental breakthrough for microscopic simulations and the development of operational strategies.

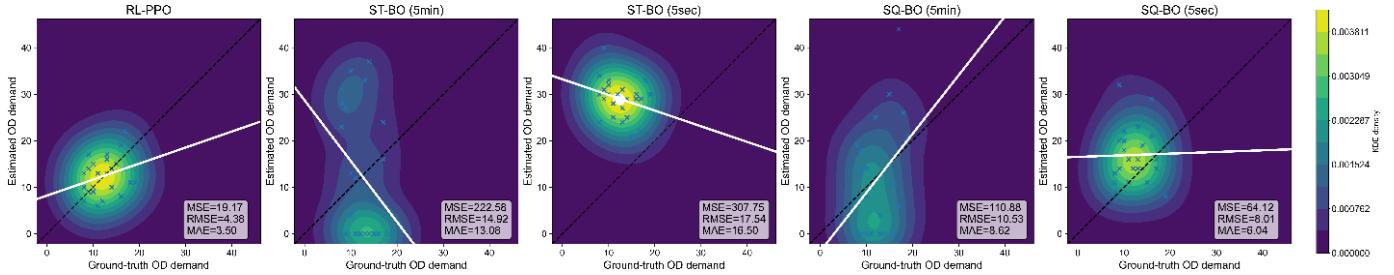


Fig. 9. Comparison of estimated versus true OD demands across different methods

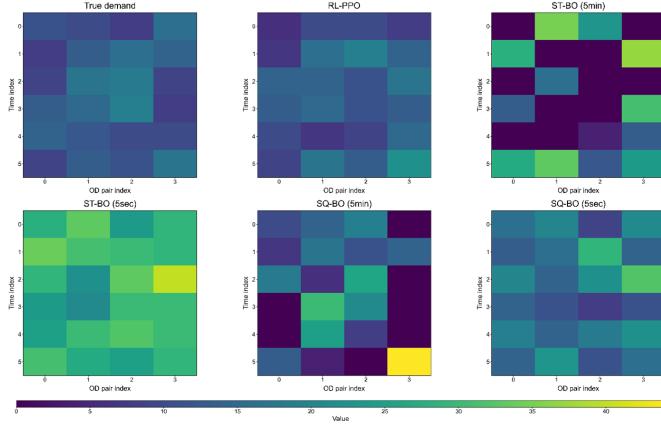


Fig. 10. Heatmap showing the values of the true OD matrix and the estimated OD matrix

This study demonstrates the performance and validity of a reinforcement learning-based framework using a toy network. However, future research is needed to apply and evaluate its performance on more realistic-scale traffic networks. In specific environments, the rapid increase in state space and action space may lead to efficiency and convergence issues in the DRL framework. To overcome this, the development of specialized techniques such as efficient state representation and action space reduction is required. Additionally, validation using real-world data is necessary to demonstrate the practical applicability. Successfully addressing these limitations will enhance the precision and reliability of microscopic simulations for developing and evaluating operational strategies.

## ACKNOWLEDGMENTS

The authors used Gemini 2.5 Pro and GPT-5 to review the manuscript, and take full responsibility for the final, edited content.

## AUTHOR CONTRIBUTIONS

The authors confirm contribution to the paper as follows: conceptualization: Min, Kim; data curation: Min; formal analysis: Min; methodology: Min, Choi, and Kim; writing-original draft: Min; writing-review and editing: Min, Choi, and Kim; draft manuscript preparation: Min, Choi, and Kim; All authors reviewed the results and approved the final version of the manuscript.

## DECLARATION OF CONFLICTING INTERESTS

The authors declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

## FUNDING

This work was supported by the Korea Institute of Police Technology (No. 092021C28S02000) and the National Research Foundation of Korea (No. 2022R1A2C2012835 and No. 00409860). The authors also acknowledge financial support from the Korea Ministry of Land, Infrastructure, and Transport as part of the Innovative Talent Education Program for Smart City.

## REFERENCES

- [1] Papageorgiou, M., C. Diakaki, V. Dinopoulou, A. Kotsialos, and Y. Wang. Review of Road Traffic Control Strategies. *Proceedings of the IEEE*, 2003. 91(12):2043–2067.
- [2] Ghiasi, A., X. Li, and J. Ma. A Mixed Traffic Speed Harmonization Model with Connected Autonomous Vehicles. *Transportation Research Part C: Emerging Technologies*, 2019. 104:210–233.
- [3] Xie, D. F., Z. Z. Fang, B. Jia, and Z. He. A Data-Driven Lane-Changing Model Based on Deep Learning. *Transportation Research Part C: Emerging Technologies*, 2019. 106:41–60.
- [4] Li, J., C. Yu, Z. Shen, Z. Su, and W. Ma. A Survey on Urban Traffic Control Under Mixed Traffic Environment with Connected Automated Vehicles. *Transportation Research Part C: Emerging Technologies*, 2023. 154:104258.
- [5] Toledo, T., M. E. Ben-Akiva, D. Darda, M. Jha, and H. N. Koutsopoulos. Calibration of Microscopic Traffic Simulation Models with Aggregate Data. *Transportation Research Record*, 2004. 1876(1):10–19.
- [6] Papathanasopoulou, V., I. Markou, and C. Antoniou. Online Calibration for Microscopic Traffic Simulation and Dynamic Multi-Step Prediction of Traffic Speed. *Transportation Research Part C: Emerging Technologies*, 2016. 68:144–159.
- [7] Patwary, A. U., W. Huang, and H. K. Lo. Metamodel-Based Calibration of Large-Scale Multimodal Microscopic Traffic Simulation. *Transportation Research Part C: Emerging Technologies*, 2021. 124:102859.
- [8] Osorio, C. Dynamic Origin–Destination Matrix Calibration for Large-Scale Network Simulators. *Transportation Research Part C: Emerging Technologies*, 2019a. 98:186–206.
- [9] Huo, J., C. Liu, J. Chen, Q. Meng, J. Wang, and Z. Liu. Simulation-Based Dynamic Origin–Destination Matrix Estimation on Freeways: A Bayesian Optimization Approach. *Transportation Research Part E: Logistics and Transportation Review*, 2023. 173:103108.
- [10] Yang, H., T. Sasaki, Y. Iida, and Y. Asakura. Estimation of Origin–Destination Matrices from Link Traffic Counts on Congested Networks. *Transportation Research Part B: Methodological*, 1992. 26(6):417–434.
- [11] Tavana, H. *Internally Consistent Estimation of Dynamic Network Origin–Destination Flows from Intelligent Transportation Systems Data Using Bi-Level Optimization*. Ph.D. dissertation. University of Texas at Austin, Austin, Tex., 2001.

- [12] Ros-Roca, X., L. Montero, J. Barceló, K. Nökel, and G. Gentile. A Practical Approach to Assignment-Free Dynamic Origin-Destination Matrix Estimation Problem. *Transportation Research Part C: Emerging Technologies*, 2022. 134:103477.
- [13] Balakrishna, R., C. Antoniou, M. Ben-Akiva, H. N. Koutsopoulos, and Y. Wen. Calibration of Microscopic Traffic Simulation Models: Methods and Application. *Transportation Research Record*, 2008. 1999:198–207.
- [14] Zhang, C., and C. Osorio. *Efficient Offline Calibration of Origin-Destination (Demand) for Large-Scale Stochastic Traffic Models*. MIT Tech. Rep., Massachusetts Institute of Technology, Cambridge, Mass., 2017.
- [15] Treiber, M., and A. Kesting. *Traffic Flow Dynamics*. Springer, Berlin, 2013.
- [16] Minsky, M. Steps toward Artificial Intelligence. *Proceedings of the IRE*, 1961. 49: 8-30.
- [17] Pignatelli, E., J. Ferret, M. Geist, T. Mesnard, H. van Hasselt, O. Pietquin, and L. Toni. A Survey of Temporal Credit Assignment in Deep Reinforcement Learning. arXiv preprint, 2023. arXiv:2312.01072.
- [18] Joe, W. and H. C. Lau. Deep Reinforcement Learning Approach to Solve Dynamic Vehicle Routing Problem with Stochastic Customers. *Proceedings of the International Conference on Automated Planning and Scheduling*, 2020. 30:394-402.
- [19] Arulkumaran, K., M. P. Deisenroth, M. Brundage, and A. A. Bharath. A Brief Survey of Deep Reinforcement Learning. arXiv preprint, 2017. arXiv:1708.05866.
- [20] Schulman, J., F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal Policy Optimization Algorithms. arXiv preprint, 2017. arXiv:1707.06347.
- [21] Schulman, J., P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional Continuous Control Using Generalized Advantage Estimation. arXiv preprint, 2015. arXiv:1506.02438.
- [22] Willumsen, L. G. Estimation of an OD Matrix from Traffic Counts – A Review. Working Paper. Institute of Transport Studies, University of Leeds, Leeds, UK, 1978.
- [23] Cascetta, E., and S. Nguyen. A Unified Framework for Estimating or Updating Origin/Destination Matrices from Traffic Counts. *Transportation Research Part B: Methodological*, 1988. 22(6):437–455.
- [24] Cremer, M., and H. Keller. A New Class of Dynamic Methods for the Identification of Origin-Destination Flows. *Transportation Research Part B: Methodological*, 1987. 21(2):117–132.
- [25] Cascetta, E., D. Inaudi, and G. Marquis. Dynamic Estimators of Origin-Destination Matrices Using Traffic Counts. *Transportation Science*, 1993. 27(4):363–373.
- [26] Ashok, K., and M. E. Ben-Akiva. Alternative Approaches for Real-Time Estimation and Prediction of Time-Dependent Origin-Destination Flows. *Transportation Science*, 2000. 34(1):21–36.
- [27] Bierlaire, M., and F. Crittin. An Efficient Algorithm for Real-Time Estimation and Prediction of Dynamic OD Tables. *Operations Research*, 2004. 52(1):116–127.
- [28] Bard, J. F., and J. T. Moore. A Branch and Bound Algorithm for the Bilevel Programming Problem. *SIAM Journal on Scientific and Statistical Computing*, 1990. 11(2):281–292.
- [29] Shafiei, S., Z. Gu, and M. Saberi. Calibration and Validation of a Simulation-Based Dynamic Traffic Assignment Model for a Large-Scale Congested Network. *Simulation Modelling Practice and Theory*, 2018. 86:169–186.
- [30] Maher, M. J., X. Zhang, and D. Van Vliet. A Bi-Level Programming Approach for Trip Matrix Estimation and Traffic Control Problems with Stochastic User Equilibrium Link Flows. *Transportation Research Part B: Methodological*, 2001. 35(1):23–40.
- [31] Lu, L., Y. Xu, C. Antoniou, and M. Ben-Akiva. An Enhanced SPSA Algorithm for the Calibration of Dynamic Traffic Assignment Models. *Transportation Research Part C: Emerging Technologies*, 2015. 51:149–166.
- [32] Rao, W., Y. J. Wu, J. Xia, J. Ou, and R. Kluger. Origin–Destination Pattern Estimation Based on Trajectory Reconstruction Using Automatic License Plate Recognition Data. *Transportation Research Part C: Emerging Technologies*, 2018. 95:29–46.
- [33] Krishnakumari, P., H. Van Lint, T. Djukic, and O. Cats. A Data Driven Method for OD Matrix Estimation. *Transportation Research Part C: Emerging Technologies*, 2020. 113:38–56.
- [34] Pourmoradnasseri, M., K. Khoshkhah, and A. Hadachi. Leveraging IoT Data Stream for Near-Real-Time Calibration of City-Scale Microscopic Traffic Simulation. *IET Smart Cities*, 2023. 5(4):269–290.
- [35] Tang, J., Y. Wang, C. Hu, Z. Li, and X. Zhang. A Spectral Clustering Enabled SPSA Algorithm for Dynamic Origin–Destination Demand Matrix Estimation. *Transportmetrica B: Transport Dynamics*, 2025. 13(1):2459928.v
- [36] Osorio, C. High-Dimensional Offline Origin–Destination (OD) Demand Calibration for Stochastic Traffic Simulators of Large-Scale Road Networks. *Transportation Research Part B: Methodological*, 2019b. 124:18–43.
- [37] Min, D., H. Yun, S. W. Ham, and D. K. Kim. Real-Time Estimation of Origin–Destination Matrices Using a Deep Neural Network for Digital Twins. *Transportation Research Record*, 2024. doi: 03611981241266837.
- [38] Hauskrecht, M. Value-function Approximations for Partially Observable Markov Decision Processes. *Journal of artificial intelligence research*, 2000. 13: 33–94.
- [39] Shi, M., Y. Liang, and N. Shroff. Theoretical Hardness and Tractability of POMDPs in RL with Partial Online State Information. arXiv preprint, 2023. arXiv:2306.08762.
- [40] Farazi, N. P., B. Zou, T. Ahamed, and L. Barua. Deep reinforcement Learning in Transportation Research: A Review. *Transportation Research Interdisciplinary Perspectives*, 2021. 11: 100425.
- [41] Mnih, V., K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing Atari with Deep Reinforcement Learning. arXiv preprint, 2013. arXiv:1312.5602.
- [42] Nguyen, S., and C. Dupuis. An Efficient Method for Computing Traffic Equilibria in Networks with Asymmetric Transportation Costs. *Transportation Science*, 1984. 18(2):185–202.
- [43] Barceló, J. *Fundamentals of Traffic Simulation*. Springer, New York, 2010.
- [44] Frazier, P. I. *A Tutorial on Bayesian Optimization*. arXiv preprint arXiv:1807.02811, 2018.