# Modeling Residual-Vehicle Effects on Uncertainty Estimation of the Connected Vehicle Penetration Rate

Shaocheng JIA[a], S.C. WONG[a,*], and Wai WONG[b,*]

[a]*Department of Civil Engineering, The University of Hong Kong, Hong Kong, China*
[b]*Department of Civil and Natural Resource Engineering, University of Canterbury, New Zealand*
*\* Co-corresponding author*

## Abstract

In the transition to full deployment of connected vehicles (CVs), the CV penetration rate plays a key role in bridging the gap between partial and complete traffic information. Several innovative methods have been proposed to estimate the CV penetration rate using only CV data. However, these methods, as point estimators, may lead to biased estimations or suboptimal solutions when applied directly in modeling or system optimization. To avoid these problems, the uncertainty and variability in the CV penetration rate must be considered. Recently, a probabilistic penetration rate (PPR) model was developed for estimating such uncertainties. The key model input is a constrained queue length distribution composed exclusively of queues formed by red signals in undersaturation conditions with no residual vehicles. However, in real-world scenarios, due to random arrivals, residual vehicles are commonly carried over from one cycle to another in temporary overflow cycles in undersaturation conditions, which seriously restricts the applicability of the PPR model. To address this limitation, this paper proposes a Markov-constrained queue length (MCQL) model that can model the complex effects of residual vehicles on the CV penetration rate uncertainty. A constrained queue with residual vehicles is decomposed into four vehicle groups: observable constrained residual vehicles, unobservable constrained residual vehicles, unconstrained residual vehicles, and new arrivals. Although the first vehicle group is observable in the former cycle, the focus of this work is to model the residual vehicles from the second and third vehicle groups in combination with the new arrivals. The MCQL model includes four sub-models, namely, the residual-vehicle model, convolutional constrained queue model, constrained residual queue model, and observable residual queue model, to isolate and derive the distribution of the constrained vehicle set formed by the three latter vehicle groups. This distribution is then substituted into the PPR model to estimate the uncertainty. Comprehensive VISSIM simulations and applications to real-world datasets demonstrate that the proposed MCQL model can accurately model the residual-vehicle effect and estimate the uncertainty. Thus, the applicability of the PPR model is truly extended to real-world settings, regardless of the presence of residual vehicles. A simple stochastic CV-based adaptive signal control example illustrates the potential of the proposed model in real-world applications.

**Keywords:** Connected vehicle penetration rate uncertainty; probabilistic penetration rate model; residual vehicle estimation; Markov-constrained queue length model; signal control with uncertainty

## 1   Introduction

With advancements in communication systems (e.g., 5G), Internet of Things technologies have undergone rapid development. These frameworks facilitate the seamless connection of various system components, thereby allowing instant exchange of information. In transportation systems,

this connectivity enables the sharing of valuable traffic information from connected vehicles (CVs), such as the location, speed, and acceleration, thereby providing numerous opportunities for the implementation of beneficial applications. Despite these advancements, the full deployment of CVs is limited by factors such as budget constraints, privacy security, and individual preferences. Consequently, a mixed traffic environment, in which both conventional vehicles and CVs coexist, is expected to prevail. Due to the absence of complete traffic information in such scenarios, the missing data must be estimated using the partial information obtained from CVs to promote traffic management and control.

The CV penetration rate, defined as the probability of a vehicle to be a CV, serves as a fundamental parameter for traffic data scaling and various model estimations and applications. Comert and Cetin (2009, 2011) and Comert (2013) proposed a series of models that use the given CV penetration rate and queue length distribution to estimate the queue length at isolated junctions. Feng et al. (2015) developed a location and speed algorithm for estimating arrival tables in the controlled optimization of a phase algorithm (Sen and Head, 1997). In this framework, the CV penetration rate is considered an essential input. Other CV-based methods for queue length estimation include a Bayesian-network-based model proposed by Hao et al. (2014) and a method based on the shockwave theory (Argote et al., 2011), both of which assume the CV penetration rate to be known during model development. Moreover, the CV penetration rate has been applied in the inference of traffic flow (Wong and Wong, 2015, 2016a, 2016c) and traffic density or accumulation (Geroliminis and Daganzo, 2008; Ambühl and Menendez, 2016; Du et al., 2016; Wong and Wong, 2019; Wong et al., 2019a, 2021) using linear data projection (Wong and Wong, 2015, 2016a, 2019; Wong et al., 2019a). In addition, the CV penetration rate has been introduced as a critical input in traffic incident impact evaluation (Wong and Wong, 2016b), travel time and speed estimations (Jenelius et al., 2013, 2015; Rahmani et al., 2015; Tian et al., 2015; Mousa et al., 2017; Khan et al., 2017; Iqbal et al., 2018; Lu et al., 2019), origin–destination estimations (Yang et al., 2017; Wang et al., 2020; Cao et al., 2021), and time exposure estimation in road safety studies (Meng et al., 2017b). Notably, due to the dynamic and stochastic nature of transportation systems, the CV penetration rate is a random variable and is not known in practice. Consequently, estimation of the CV penetration rate has emerged as a research hotspot in CV-based transportation problems.

CV penetration rates for links outfitted with on-road fixed detectors, such as loop detectors, can be directly determined using the total vehicle counts measured by the detectors and CV counts from the CV signals. However, most roadways in a network are not equipped with such detectors owing to the considerable investment and maintenance costs for universal implementation. Moreover, the installed detectors may occasionally become non-operational, leading to intrusive installation and maintenance activities that can significantly disrupt traffic flow and even result in blockages within the local transportation systems. To estimate the CV penetration rates on links without detectors, a probability distribution model (Wong and Wong, 2015, 2016a, 2019; Wong et al., 2019a) has been developed. In this framework, the probability distribution is approximated based on the CV penetration rates observed on links equipped with detectors. The expectation of this distribution is then used as an estimate for the CV penetration rates on links without detectors, given their geographical proximity to the links with detectors. Furthermore, the variance of this distribution can capture the spatial variations. Notably, this method relies on the assumption of independent and identically distributed CV penetration rates across different links within the network. However, this assumption may be violated due to the interconnectivity of roads and heterogeneous attractiveness of various urban areas. To address this issue, a CV penetration rate estimation model incorporating land-use variables has been developed (Meng et al., 2017a). Nonetheless, the use of

1     local land-use data presents several challenges that hinder the widespread adoption of this model.

2     Many researchers have attempted to overcome these limitations by estimating the CV penetration
3     rate based solely on CV data. Under the assumption of Poisson arrival, Comert (2016) proposed a
4     set of models to estimate the CV penetration rate based on partial CV information. However, the
5     assumption of Poisson arrival does not always hold in reality, and thus, these models may not be
6     applicable for generic arrival patterns. By eliminating the assumption of specific arrival patterns,
7     Wong et al. (2019b) proposed the single-source data penetration rate estimator (SSDPRE), which
8     is a fully analytical, non-parametric, and unbiased estimator to estimate the CV penetration rate.
9     This model determines the number of non-CVs preceding the last CV on the stopping locations of
10     CVs at a signalized junction and uses this partial queue information to combine two estimation
11     mechanisms—(1) the probability of the first stopping vehicle being a CV and (2) the CV
12     penetration rate of the deduced partial queue—to estimate the CV penetration rate in an unbiased
13     manner. Various other methods approximate the distribution of stopping locations of vehicles in
14     queues through maximum likelihood estimation (Zhao et al., 2019a, 2019b, 2022). Similarly, Wang
15     et al. (2024) used the method of moment to estimate CV penetration rate and vehicle arrival rate.
16     Although the aforementioned methods yield valuable insights, they do not take into account the
17     uncertainty associated with CV penetration rates. Given the dynamic and nonlinear nature of
18     transportation systems, relying solely on point estimators in model estimation and system
19     optimization may lead to biased models and suboptimal solutions (Wong and Wong, 2015, 2016,
20     2019; Wong et. al., 2019; Yin, 2008). For instance, Wong and Wong (2015, 2016, 2019) and Wong
21     et al. (2019) have proven that estimating traffic models solely based on traffic data constituted from
22     the means of CV penetration rates, without considering their variability, can result in biased model
23     parameters and standard errors. Additionally, incorporating the variabilities of parameters in
24     stochastic optimizations has been shown to be advantageous in formulating strategies to mitigate
25     traffic intersection violations (Sun et al., 2018), traffic delay and emissions (Han et al., 2016), as
26     well as human exposure to emissions (Zhang et al., 2013). Furthermore, Jia et al. (2023) conducted
27     simulation studies to demonstrate that incorporating the CV penetration rate variability into a CV-
28     based adaptive signal optimization problem can lead to a 15% decrease in average and maximum
29     driver delay and a 45% decrease in delay variance. Therefore, accurately modeling the uncertainty
30     of CV penetration rates is crucial to obtain unbiased transport models and optimal solutions.

31     Utilizing the SSDPRE (Wong et al. 2019), the output from a cycle can be taken as the realized CV
32     penetration rate for that specific cycle. In oversaturation conditions where the volume-to-capacity
33     (V/C) ratios consistently exceed one, demand persistently surpasses the capacity. Due to the
34     continuous carryover in oversaturation conditions, nearly all vehicle identities can be determined
35     using the bridging queue algorithm proposed in Wong et al. (2019b). This allows for the precise
36     determination of the population CV penetration rate with a high degree of certainty. Nevertheless,
37     in undersaturation conditions where the V/C ratios are less than one, vehicle identities cannot be
38     fully revealed. Due to stochastic arrivals and random appearance of CVs, the realized CV
39     penetration rate is subject to uncertainty. To quantify such uncertainty, an analytical probabilistic
40     penetration rate (PPR) model has recently been developed (Jia et al., 2023). A key input of this
41     model is the constrained queue length[1] distribution, which depends on the queues of vehicles that
42     stop at red signals in each traffic light cycle. This distribution can be estimated using probabilistic
43     dissipation time (PDT) or constant dissipation time (CDT) models derived in undersaturation

---

[1] A constrained queue length is defined as the total number of vehicles that are stopped by a red signal over a cycle, which differs from the usual queue length, i.e., the queue length at a specific moment in time. Figure 1 illustrates the notion of a constrained queue length.

conditions without any residual vehicles. However, in real-world scenarios, due to random arrivals, demand can temporarily exceed capacity, and residual vehicles are commonly carried over from one cycle to another in temporary overflow cycles under undersaturation conditions. These temporary overflow cycles are referred to as temporary overflow conditions. As temporary overflow conditions are common in real-world situations, the applicability of the methods proposed in Jia et al. (2023) is seriously limited. Neglecting the effects of these residual vehicles in temporary overflow conditions may lead to inaccurate estimation of the constrained queue length, which can adversely affect the estimation of the CV penetration rate uncertainty. Considering the highly nonlinear nature of transportation systems, such unreliable uncertainty estimates can undermine the effectiveness of the PPR model in practical applications. Therefore, the residual-vehicle effects must be accurately modeled to ensure the generalizability and effectiveness of the PPR model in estimating the CV penetration rate uncertainty.

To this end, this paper proposes a Markov-constrained queue length (MCQL) model. A constrained queue with residual vehicles is generically decomposed into four vehicle groups: observable constrained residual vehicles (Group 1), unobservable constrained residual vehicles (Group 2), unconstrained residual vehicles (Group 3), and new arrivals (Group 4). Although the residual vehicles from Group 1 are observable in the former cycle, the focus of this work is to model the residual vehicles from Groups 2 and 3 in combination with the new arrivals in Group 4. Four sub-models, namely, the residual-vehicle model, convolutional constrained queue model, constrained residual queue model, and observable residual queue model, are introduced to isolate and derive the distribution of the constrained vehicle set formed by the three latter vehicle groups. This distribution is then substituted into the PPR model for uncertainty estimation. Extensive numerical simulations and real-world dataset applications provide strong evidence for the effectiveness of the proposed MCQL model in accurately capturing residual-vehicle effects and estimating the uncertainty in the CV penetration rate. Furthermore, a practical demonstration of a stochastic CV-based adaptive signal control showcases the potential of the MCQL model in system optimizations. The proposed models can exploit the full capability of the PPR model and enhance its applicability in real-world situations.

The remainder of this paper is organized as follows. Section 2 introduces the necessary preliminary knowledge, and Section 3 defines the research problem and notation. Section 4 describes the formulation of the MCQL model. Section 5 presents details of the comprehensive micro-simulation study performed based on VISSIM to validate the proposed model. Section 6 describes the validation of the proposed models using the Next Generation Simulation (NGSIM) dataset (Federal Highway Administration, 2006). Moreover, a demonstrative example based on stochastic CV-based adaptive signal control is presented to highlight the significance of modeling residual-vehicle effects. Section 7 summarizes the findings and implications of the study.

## 2   Prior Work

### 2.1   SSDPRE

When approaching a signalized intersection, the arriving vehicles stop at red lights and form queues. A constrained queue set, $\Psi$, is defined as a set of vehicles that stop at a red light, e.g., the sets of vehicle trajectories enclosed by the triangles in Figure 1. $|\Psi| = N$ is the number of vehicles in a constrained queue; and $n$ and $\widetilde{N}$ are the number of CVs and number of observable vehicles in the constrained queue, respectively. $n$ is observable, and $\widetilde{N}$ is estimated by dividing the distance between the stop bar and rear end of the last CV by the effective vehicle length. The effective

vehicle length is the average distance between the rear ends of the preceding and following stopping vehicles. The CVs and non-CVs are assumed to be sufficiently mixed within a link owing to the lane-changing behaviors of the drivers; and $i$ and $m$ are the indices for the $i^{\text{th}}$ constrained queue formed in cycle $i$ and total number of constrained queues, respectively. The SSDPRE for this scenario can be formulated as (Wong et al. 2019)

$$\text{SSDPRE} = \frac{\sum_{i=1}^{m} \tilde{p}_i}{m}, \tag{1}$$

where

$$\tilde{p}_i = S(n_i, \widetilde{N}_i) = \begin{cases} \frac{n_i - 1}{\widetilde{N}_i - 1} & \text{if } n_i > 1 \text{ and } \widetilde{N}_i > 1 \\ 1 & \text{if } n_i = 1 \text{ and } \widetilde{N}_i = 1 \\ 0 & \text{if } n_i = 1 \text{ and } \widetilde{N}_i > 1 \\ 0 & \text{if } n_i = 0 \text{ and } \widetilde{N}_i = 0 \end{cases}. \tag{2}$$

Thus, $\tilde{p}_i$ is taken as realized CV penetration rate for cycle $i$. The realized CV penetration rates across cycles form the distribution of the random variable $\tilde{p}$.

## 2.2 PPR model

The variance of the distribution of $\tilde{p}$, $Var(\tilde{p})$, can be taken as the CV penetration rate uncertainty. Consider any $|\Psi| = N$ following any counting distribution such that $P(N = i) = \xi_i, \forall i = 0, 1, 2, \ldots, k$. The number of CVs follows a binomial distribution, i.e., $n \sim B(N, p)$, where $p$ represents the average CV penetration rate and $N \geq n \geq 0$. In this case, $Var(\tilde{p})$ can be defined as

$$Var(\tilde{p}) = \lim_{k \to +\infty} \left[ \sum_{i=1}^{k} \xi_i \, V_2(i, p) \right], \tag{3}$$

where

$$V_2(N, p) = \begin{cases} \sum_{i=2}^{N} p^i (1-p)^{N-i} \left[ V_1(i, N) + \left(\frac{i}{N}\right)^2 \right] \binom{N}{i} - p^2 + p(1-p)^{N-1} & \text{if } N > 1 \\ p(1-p) & \text{if } N = 1 \end{cases}, \tag{4}$$

and

$$V_1(n, N) = \begin{cases} \frac{\sum_{i=1}^{N-n+1} \frac{n-1}{N-i} \binom{N-i-1}{n-2}}{\binom{N}{n}} - \frac{n^2}{N^2} & \text{if } n > 1 \\ \frac{n^2 - 2n + N}{N^2} & \text{if } n = 1 \\ 0 & \text{if } n = 0 \end{cases}. \tag{5}$$

Two corollaries of the PPR model are presented in **Appendix A**.

## 2.3 PDT and CDT models

The constrained queue length distribution is an essential input of the PPR model. The PDT model can accurately model this input by using a time interval partitioning notion, which divides a given interval into the red period, first dissipation period for vehicles arriving in the red period, second dissipation period for vehicles arriving in the first dissipation period, and so on. For any given arrival pattern with an average arrival rate of $q$, a red period of $r$, and a saturation headway of $\tau$, the constrained queue length distribution can be expressed as

$$P(N = k) = \begin{cases} f(k; qr) f(0; qk\tau) + \\ \sum_{i=1}^{k-1} \sum_{j=1}^{J_i} f(i; qr) \tilde{P}_j(N = k, M = i) W_j(N = k, M = i) & \text{if } k \in \mathbb{N}^+, \\ f(0; qr) & \text{if } k = 0 \end{cases} \tag{6}$$

where $f(\eta; qt)$ represents the probability of arriving $\eta$ vehicles within time $t$; $\tilde{P}_j(N = k, M = i)$ is the $j^{th}$ unique value of the product of the probabilities of observing the remaining $k - i$ vehicles in the subsequent partitioned time intervals, with $\forall k \in \mathbb{N}^+$, $i \in [1, k]$, $j \in [1, J_i]$; and $W_j(N = k, M = i)$ is the weighting factor of $\tilde{P}_j(N = k, M = i)$.

The PDT model achieves a high estimation accuracy, but its implementation is complex. To alleviate the complexity, a simplified model, i.e., the CDT model, has been developed. The constrained queue length is assumed to follow a Poisson distribution. The average constrained queue length, $N_0$, is the model parameter of the constrained queue length distribution, where

$$N_0 = \frac{sqr}{s-q}. \tag{7}$$

# 3    Problem Statement and Notation

## 3.1    Problem statement

Although PDT and CDT models can estimate the constrained queue length distribution, they can only be applied to undersaturation scenarios without residual vehicles. Specifically, these models can represent the newly arriving vehicles in each cycle, such as the constrained queue set of Cycle 1 in Figure 1. However, in practice, residual vehicles often appear owing to temporary high demands. For example, in Figure 1, vehicles 7, 8, and 9 are residual vehicles that carry over from Cycle 1 to Cycle 2. Vehicle 7 is the observable constrained residual vehicle from Cycle 1 (Group 1), vehicle 8 is an unobservable constrained residual vehicle from Cycle 1 (Group 2), and vehicle 9 is an unconstrained residual vehicle from Cycle 1 (Group 3). The constrained queue set of Cycle 2 is {7, 8, 9, 10, 11, 12, 13}, in which the subset {10, 11, 12, 13} represents new arrivals in the constrained queue set of Cycle 2 (Group 4), which can be modeled by PDT or CDT models. To avoid double-counting, vehicle 7[2] must not be used to estimate $\tilde{p}_2$ in Cycle 2. Thus, in addition to the new arrivals, the unobservable constrained and unconstrained residual vehicles are random components that can influence the estimation of $\tilde{p}_i$. Ignoring vehicles 8 and 9 can degrade the estimation accuracy of the CV penetration rate uncertainty. Therefore, the residual-vehicle effects must be appropriately modeled in the uncertainty estimation process, especially for near-saturation situations, which is the motivation for this research.
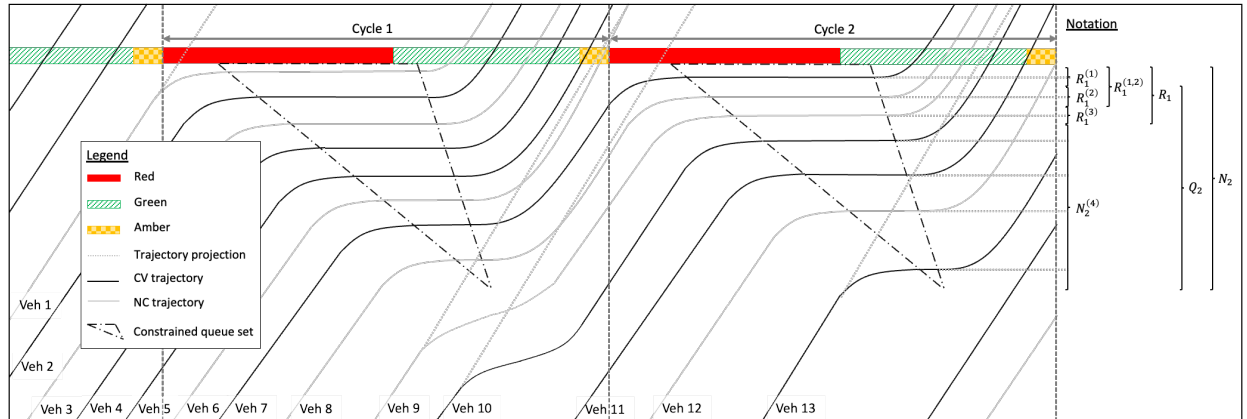


Figure 1: Vehicle trajectories in two consecutive cycles.

---

[2] Vehicle 7 is the last observable vehicle in the constrained queue set of Cycle 1, which is used to estimate $\tilde{p}_1$. Details of the counting method can be found in the work of Wong et al. (2019).

## 3.2 Notation

In addition to the previously introduced notation, the following notation is used in this paper:

Table 1. Notation.

| Notation | Description | Example (Figure 1) |
|---|---|---|
| $R_{i-1}^{(1)}$ | Number of observable constrained residual vehicles from cycle $i-1$ | $R_1^{(1)} = 1$ |
| $R_{i-1}^{(2)}$ | Number of unobservable constrained residual vehicles from cycle $i-1$ | $R_1^{(2)} = 1$ |
| $R_{i-1}^{(3)}$ | Number of unconstrained residual vehicles from cycle $i-1$ | $R_1^{(3)} = 1$ |
| $R_{i-1}$ | Number of residual vehicles from cycle $i-1$, i.e., $R_{i-1}^{(1)} + R_{i-1}^{(2)} + R_{i-1}^{(3)}$ | $R_1 = 3$ |
| $N_i^{(4)}$ | Number of new arrivals in the constrained queue set of cycle $i$ | $N_2^{(4)} = 4$ |
| $N_i$ | Number of vehicles in the constrained queue set of cycle $i$, i.e., $R_{i-1} + N_i^{(4)}$ | $N_2 = 7$ |
| $Q_i$ | Number of vehicles in the constrained queue set of cycle $i$, excluding the observable residual vehicles carried over from cycle $i-1$, i.e., $N_i - R_{i-1}^{(1)}$ | $Q_2 = 6$ |
| $R_{i-1}^{(1,2)}$ | Number of constrained residual vehicles from cycle $i-1$, i.e., $R_{i-1}^{(1)} + R_{i-1}^{(2)}$ | $R_1^{(1,2)} = 2$ |
| $g$ | Effective green | - |
| $r$ | Effective red | - |
| $s$ | Saturation flow | - |
| $c$ | Cycle length | - |
| $q$ | Average arrival rate | - |
| $D_i$ | Demand of cycle $i$ | - |
| $D^*$ | Maximum number of vehicles that can be discharged in a cycle | - |
| $P$ | State transition matrix of $R_{i-1}$ | - |
| $S$ | State space of $R_{i-1}$ | - |
| $F$ | Discrete Fourier transform (DFT) function | - |
| $F^{-1}$ | Inverse DFT function | - |

# 4 Methodology

For any consecutive cycles $i-1$ and $i$ with a temporary high demand, a certain number of residual vehicles from cycle $i-1$ are carried over to cycle $i$. The CV penetration rate uncertainty is governed by the combination of this number of residual vehicles and the new arrivals in cycle $i$, $Q_i$. This total number of vehicles consists of $R_{i-1}^{(2)}$, $R_{i-1}^{(3)}$, and $N_i^{(4)}$. Although $N_i^{(4)}$ can be obtained by PDT or CDT models, it is challenging to directly derive $R_{i-1}^{(2)}$ and $R_{i-1}^{(3)}$. Alternatively, $Q_i$ can be expressed as

$$Q_i = N_i - R_{i-1}^{(1)}. \tag{8}$$

Thus, the proposed MCQL model involves three parts, pertaining to the estimations of $N_i$, $R_{i-1}^{(1)}$, and $Q_i$. The proposed model is derived based on two fundamental assumptions. First, it is assumed that connected vehicles (CVs) and non-CVs are adequately mixed within a lane, due to frequent

lane-changing and overtaking behaviors. As a result, each vehicle has the same probability of being a CV ($p$) or a non-CV ($1 - p$). Second, the state transition of the number of residual vehicles is assumed to be modeled by a Markov chain, as the number of residual vehicles in the current cycle depends solely on the number of residual vehicles in the previous cycle, which is consistent with the Markov property.

## 4.1 Estimation of $N_i$

$N_i$ can be expressed as the sum of $N_i^{(4)}$ and $R_{i-1}$:

$$N_i = N_i^{(4)} + R_{i-1}. \tag{9}$$

$N_i^{(4)}$ can be estimated using PDT or CDT methods, and a residual-vehicle model is proposed to estimate $R_{i-1}$. Then, a convolutional constrained queue model is developed for estimating $N_i$.

### 4.1.1 Residual-vehicle model

Taking into account vehicles arriving according to any counting distribution:

$$P(D_i = k) = p_k, \forall k \in \mathbb{N}. \tag{10}$$

The maximum number of vehicles that can be discharged in a cycle, $D^*$, can be estimated as $D^* = \lfloor gs \rfloor$, where $\lfloor \cdot \rfloor$ represents the floor function. If the temporary demand in cycle $i - 1$, $D_{i-1}$, is greater than $D^*$, $R_{i-1} > 0$. For any consecutive cycles $i - 1$ and $i$, if the vehicle arrivals in different cycles are independent, i.e., $P(D_i = k | D_{i-1} = j) = P(D_i = k)$, $R_i$ depends on $R_{i-1}$ and $D_i$. Thus, $R_i$ can be modeled by a Markov chain.

For any cycle $i$, consider the Markov random process $\{R_i, i = 0, 1, 2, \ldots, k\}$, where $R_i \in S$, and $S = \{0, 1, 2, 3, \ldots, s - 1\}$ is the state space of $R_i$. $\forall j \in S, R_i = j$ represents the state of $j$ residual vehicles from cycle $i$. The probability of state $j$ is

$$P(R_i = j) = \pi_j^*, \tag{11}$$

where $\pi_j^*$ is the $j^{th}$ entry of the probability vector $\boldsymbol{\pi}^*$, which is the stationary distribution independent of cycle $i$ and represents the probability distribution of the number of residual vehicles from any cycle. $\boldsymbol{\pi}^*$ can be obtained through the following minimization:

$$\min_{\boldsymbol{\pi}} \|\boldsymbol{\pi}P - \boldsymbol{\pi}\|_2^2$$
$$s.t. \|\boldsymbol{\pi}\|_1 = 1, \pi_j > 0, \tag{12}$$

where the state transition matrix, $\boldsymbol{P}$, is presented in Table 2.

Table 2. State transition matrix.

| $R_{i-1}$ \ $R_i$ | 0 | 1 | 2 | ... | $D^*$ | $D^* + 1$ | ... | $s - 1$ |
|---|---|---|---|---|---|---|---|---|
| 0 | $\sum_{k=0}^{D^*} p_k$ | $p_{D^*+1}$ | $p_{D^*+2}$ | ... | $p_{2D^*}$ | $p_{2D^*+1}$ | ... | $p_{D^*+s-1}$ |
| 1 | $\sum_{k=0}^{D^*-1} p_k$ | $p_{D^*}$ | $p_{D^*+1}$ | ... | $p_{2D^*-1}$ | $p_{2D^*}$ | ... | $p_{D^*+s-2}$ |
| 2 | $\sum_{k=0}^{D^*-2} p_k$ | $p_{D^*-1}$ | $p_{D^*}$ | ... | $p_{2D^*-2}$ | $p_{2D^*-1}$ | ... | $p_{D^*+s-3}$ |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| $D^*$ | $p_0$ | $p_1$ | $p_2$ | ... | $p_{D^*}$ | $p_{D^*+1}$ | ... | $p_{s-1}$ |
| $D^* + 1$ | 0 | $p_0$ | $p_1$ | ... | $p_{D^*-1}$ | $p_{D^*}$ | ... | $p_{s-2}$ |
| $D^* + 2$ | 0 | 0 | $p_0$ | ... | $p_{D^*-2}$ | $p_{D^*-1}$ | ... | $p_{s-3}$ |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |

| $s-1$ | 0 | 0 | 0 | ... | ... | ... | ... | $p_{D^*}$ |
|---|---|---|---|---|---|---|---|---|

**Proof.** For any consecutive cycles $i-1$ and $i$, under any vehicle arrival pattern, the following scenarios can be derived: If $R_{i-1} = 0$, no residual vehicle is carried over from cycle $i-1$ to cycle $i$. All possible cases in this situation are presented in the first cell of Table 3. If no residual vehicle is carried over from cycle $i$ to cycle $i+1$, i.e., $R_i = 0$, the temporary demand of cycle $i$, $D_i$, must be less than or equal to the maximum number of vehicles that can be discharged in a cycle, $D^*$, i.e., $D_i \leq D^*$. The probability of this case, $P(R_i = 0|R_{i-1} = 0)$, is $\sum_{k=0}^{D^*} p_k$. If $R_i = 1$, $D_i$ must be equal to $D^* + 1$, which means that only one vehicle cannot be discharged in cycle $i$. In this case, $P(R_i = 1|R_{i-1} = 0) = p_{D^*+1}$. All possible cases can be enumerated in a similar manner. When $R_{i-1} = 0$, the sum of the probabilities of all these cases must equal 1, i.e., $\sum_{k=0}^{\infty} P(R_i = k|R_{i-1} = 0) = 1$.

Table 3. Enumeration of all possible numbers of residual vehicles carried over from cycles $i-1$ and $i$.

| $R_{i-1}$ | $R_i$ | $D_i$ | Probability |
|---|---|---|---|
| 0 | 0 | $\leq D^*$ | $\sum_{k=0}^{D^*} p_k$ |
|  | 1 | $D^* + 1$ | $p_{D^*+1}$ |
|  | 2 | $D^* + 2$ | $p_{D^*+2}$ |
|  | $\vdots$ | $\vdots$ | $\vdots$ |
|  | $k$, where $k \to \infty$ | $D^* + k$ | $p_{D^*+k}$ |
| 1 | 0 | $\leq D^* - 1$ | $\sum_{k=0}^{D^*-1} p_k$ |
|  | 1 | $D^*$ | $p_{D^*}$ |
|  | 2 | $D^* + 1$ | $p_{D^*+1}$ |
|  | $\vdots$ | $\vdots$ | $\vdots$ |
|  | $k$, where $k \to \infty$ | $D^* + k - 1$ | $p_{D^*+k-1}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $D^* + 1$ | 0 | - | 0 |
|  | 1 | 0 | $p_0$ |
|  | 2 | 1 | $p_1$ |
|  | $\vdots$ | $\vdots$ | $\vdots$ |
|  | $k$, where $k \to \infty$ | $k - 1$ | $p_{k-1}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $k'$, where $k' \to \infty$ | 0 | - | 0 |
|  | 1 | - | 0 |
|  | $\vdots$ | $\vdots$ | $\vdots$ |
|  | $k' - D^* - 1$ | - | 0 |
|  | $k' - D^*$ | 0 | $p_0$ |
|  | $k' - D^* + 1$ | 1 | $p_1$ |
|  | $\vdots$ | $\vdots$ | $\vdots$ |
|  | $k$, where $k \to \infty$ | $k - k' + D^*$ | $p_{k-k'+D^*}$ |

If $R_{i-1} = 1$, one residual vehicle is carried over from cycle $i-1$ to cycle $i$. All possible cases in this condition are presented in the second cell of Table 3. If $R_i = 0$, $D_i$ must be less than or equal

to $D^* - 1$ such that all the newly arrived vehicles and the residual vehicle carried over from cycle $i - 1$ can be discharged. Thus, $P(R_i = 0|R_{i-1} = 1) = \sum_{k=0}^{D^*-1} p_k$. If $R_i = 1$, $D_i$ must be equal to $D^*$, and thus, the last vehicle from cycle $i$ is carried over to cycle $i + 1$. In this case, $P(R_i = 1|R_{i-1} = 1) = p_{D^*}$. All possible cases can be similarly enumerated. Furthermore, $\sum_{k=0}^{\infty} P(R_i = k|R_{i-1} = 1) = 1$. All corresponding cases for $R_{i-1} \in [0, D^*]$ can be enumerated in a similar manner.

If $R_{i-1} = D^* + 1$, $R_i$ must be equal to or greater than 1, i.e., at least one residual vehicle is carried over from cycle $i$ to cycle $i + 1$. Therefore, $P(R_i = 0|R_{i-1} = D^* + 1) = 0$, $P(R_i = 1|R_{i-1} = D^* + 1) = p_0$, and so on. Again, $\sum_{k=0}^{\infty} P(R_i = k|R_{i-1} = D^* + 1) = 1$. For $R_{i-1} \geq D^* + 1$, a similar enumeration can be performed.

The transition process of the number of residual vehicles between any two consecutive cycles can be modeled by a Markov chain. The state transition matrix can be constructed based on Table 3. Notably, although the transition process can be considered a discrete-time and infinite Markov chain in theory, the number of residual vehicles from any cycle $i$, $R_i$, is not infinite in real-world transportation systems. As the probability of the number of residual vehicles being infinite is negligible, the theoretically infinite space for the number of residual vehicles can be approximated as a finite space, denoted as $S = \{0, 1, 2, 3, \ldots, s - 1\}$. This approximation enables the construction of the corresponding state transition matrix $\boldsymbol{P} \in \mathbb{R}^{s \times s}$, as shown in Table 2. This finite Markov chain has the following two properties: (1) It is irreducible, meaning that it only has one communication class; and (2) it is aperiodic due to the existence of self-transition. Based on the Markov chain stationary state theorem, a unique stationary distribution exists for such a finite Markov chain, and its limiting distribution converges to this stationary distribution. Consequently, the unique stationary distribution can be determined by solving the minimization problem described in Eq. (12).

**QED.**

It is important to emphasize that the proposed residual-vehicle model is generic, accommodating various arrival patterns that could follow any counting distribution. In the numerical experiments and application presented in Sections 5 and 6, Poisson distributions, which are commonly assumed as random arrival patterns at isolated junctions, are used for illustrative purposes.

### 4.1.2 Convolutional constrained queue model

$R_{i-1}$ can be obtained from the residual-vehicle model, and $N_i^{(4)}$ can be estimated using PDT or CDT methods. Referring Eq. (9), the distribution of $N_i$ can be derived by enumerating all possible combinations of $R_{i-1}$ and $N_i^{(4)}$. A Fourier transformation method can be employed to avoid tedious enumerations and efficiently perform the calculations. Further details are provided below.

The probability distributions of $R_{i-1}$ and $N_i^{(4)}$, $P(R_{i-1} = j) = \pi_j^*, \forall j \in \{0, 1, 2, \ldots, s - 1\}$, and $P\left(N_i^{(4)} = j\right) = a_j, \forall j \in \{0, 1, 2, \ldots, s_1 - 1\}$, respectively, can be expressed in the following vector forms:

$$\boldsymbol{f}_{R_{i-1}} = \left[\pi_0^*, \pi_1^*, \ldots, \pi_{s-1}^*, \underbrace{0, 0, \ldots, 0}_{s_1 - 1}\right] \tag{13}$$

and

$$f_{N_i^{(4)}} = \left[a_0, a_1, \dots, a_{s_1-1}, \underbrace{0, 0, \dots, 0}_{s-1}\right], \tag{14}$$

where the $s_1 - 1$ and $s - 1$ zeros are augmented entries that ensure that the two vectors have the same dimensions. The probability distribution of $N_i$, $P(N_i = j) = b_j, \forall j \in \{0, 1, 2, \dots, s + s_1 - 2\}$, can be obtained as

$$f_{N_i} = F^{-1}\left[F\left(f_{N_i^{(4)}}\right) \odot F\left(f_{R_{i-1}}\right)\right], \tag{15}$$

where $F$ is the discrete Fourier transform (DFT) function, $F^{-1}$ is the inverse DFT function, $\odot$ represents element-wise multiplication, and $f_{N_i} = [b_0, b_1, \dots, b_{s+s_1-2}]$.

**Proof.** Given probability distributions $f_{R_{i-1}}$ and $f_{N_i^{(4)}}$, the following possibilities for $N_i$ can be derived.

If $N_i = 0$, then $R_{i-1} = N_i^{(4)} = 0$. The probability for this case is $\pi_0^* a_0$.

If $N_i = 1$, two cases can be identified: (1) $R_{i-1} = 0$ and $N_i^{(4)} = 1$ or (2) $R_{i-1} = 1$ and $N_i^{(4)} = 0$, the probabilities for which are $\pi_0^* a_1$ are $\pi_1^* a_0$, respectively. Thus, the total probability for $N_i = 1$ is the sum of these two probabilities, i.e., $\pi_0^* a_1 + \pi_1^* a_0$.

If $N_i = 2$, three cases are possible: (1) when $R_{i-1} = 0$ and $N_i^{(4)} = 2$, the probability is $\pi_0^* a_2$; (2) when $R_{i-1} = 1$ and $N_i^{(4)} = 1$, the probability is $\pi_1^* a_1$; and (3) when $R_{i-1} = 2$ and $N_i^{(4)} = 0$, the probability is $\pi_2^* a_0$. Thus, the total probability for $N_i = 2$ is $\pi_0^* a_2 + \pi_1^* a_1 + \pi_2^* a_0$.

All possible values of $N_i$ and corresponding probabilities can be enumerated in a similar manner. For the last possible value of $N_i$, when $N_i = s + s_1 - 2$, $R_{i-1} = s - 1$, $N_i^{(4)} = s_1 - 1$, and the corresponding probability is $\pi_{s-1}^* a_{s_1-1}$. Therefore, the probability distribution of $N_i$ can be expressed as

$$b_j = \Sigma_{k=0}^{j} \pi_k^* a_{j-k}, j = 0, 1, 2, \dots, s + s_1 - 2, \tag{16}$$

where $b_j$ is the $j^{th}$ entry of the probability vector $f_{N_i}$. Equation (16) can be alternatively expressed as a one-dimensional discrete linear convolution, as shown in Eq. (17):

$$f_{N_i} = f'_{R_{i-1}} * f'_{N_i^{(4)}}, \tag{17}$$

where $*$ represents linear convolution; $f'_{R_{i-1}} = [\pi_0^*, \pi_1^*, \dots, \pi_{s-1}^*]$; and $f'_{N_i^{(4)}} = [a_0, a_1, \dots, a_{s_1-1}]$.

According to convolution theory,

$$f_{N_i} = f'_{R_{i-1}} * f'_{N_i^{(4)}} = R\left(f_{R_{i-1}} \circledast f_{N_i^{(4)}}\right), \tag{18}$$

where $\circledast$ represents cyclic convolution; and $R(\cdot)$ is a function that extracts the principal value sequence, which refers to the first $s + s_1 - 1$ values in $f_{R_{i-1}} \circledast f_{N_i^{(4)}}$. Based on the cyclic

convolution theorem, the following expression can be obtained:

$$F\left[R\left(f_{R_{i-1}} * f_{N_i^{(4)}}\right)\right] = F\left(f_{R_{i-1}}\right) \odot F\left(f_{N_i^{(4)}}\right). \tag{19}$$

Substituting Eq. (19) into Eq. (18) yields the probability distribution of $N_i$ as

$$f_{N_i} = F^{-1}\left[F\left(f_{N_i^{(4)}}\right) \odot F(f_{R_{i-1}})\right]. \tag{20}$$

**QED.**

## 4.2 Estimation of $R_{i-1}^{(1)}$

The constrained residual queue and observable residual queue models are used to estimate the probability distributions of $R_{i-1}^{(1,2)}$ and $R_{i-1}^{(1)}$, respectively.

### 4.2.1 Constrained residual queue model

Given the probability distributions of $N_{i-1}$, $P(N_{i-1} = k), \forall k \in \mathbb{N}$, the probability distribution of $R_{i-1}^{(1,2)}$ can be expressed as

$$P\left(R_{i-1}^{(1,2)} = j\right) = \begin{cases} \sum_{k=0}^{D^*} P(N_{i-1} = k), & if\ j = 0 \\ P(N_{i-1} = D^* + j), & if\ j > 0 \end{cases}. \tag{21}$$

**Proof.** If $N_{i-1} = 0, 1, 2, \dots, D^*$, all vehicles in the constrained queue can be discharged. Therefore, no constrained residual queue exists, and $P\left(R_{i-1}^{(1,2)} = 0\right) = \sum_{k=0}^{D^*} P(N_{i-1} = k)$.

If $N_{i-1} = D^* + 1$, one constrained residual vehicle will be carried over to cycle $i$, i.e., $R_{i-1}^{(1,2)} = 1$. Thus, $P\left(R_{i-1}^{(1,2)} = 1\right) = P(N_{i-1} = D^* + 1)$. Similarly, the probabilities of all possible values of $R_{i-1}^{(1,2)}$ are $P\left(R_{i-1}^{(1,2)} = j\right) = P(N_{i-1} = D^* + j), \forall j > 0$. Therefore, the probability distribution of $R_{i-1}^{(1,2)}$ can be obtained, as shown in Eq. (21).

**QED.**

### 4.2.2 Observable residual queue model

Given the probability distribution of $R_{i-1}^{(1,2)}$, $P\left(R_{i-1}^{(1,2)} = j\right)$, $\forall j \in \{0, 1, 2, \dots, l-1\}$, obtained from the constrained residual queue model; CV penetration rate $p$; and identity of a vehicle being either a CV or a non-CV following a Bernoulli distribution with parameter $p$, the probability distribution of $R_{i-1}^{(1)}$ can be expressed as

$$P(R_{i-1}^{(1)} = j) = \sum_{k=0}^{l-1-j} f_{R_{i-1}^{(1,2)}}^{j,k} x_{1-p}^{j,k}, \forall j = 0, 1, 2, \dots, l-1, \tag{22}$$

where

$$f_{R_{i-1}^{(1,2)}}^{j,k} = P\left(R_{i-1}^{(1,2)} = k\right) H_{l,l-j,k}^{-1}, \forall k = 0, 1, 2, \dots, l-1, \tag{23}$$

$$H_{l,l-j,k}^{-1} = \begin{cases} 1, & if\ k = j, j+1, \dots, l-1 \\ 0, & otherwise \end{cases}, \tag{24}$$

12

$$x_{1-p}^{j,k} = \begin{cases} (1-p)^k, & if\ j = 0 \\ px_{1-p}^{0,k}H_{l,l-j,k}, & if\ j = 1, \\ x_{1-p}^{j-1,k}H_{l,l-j,k}, & if\ j > 1 \end{cases}$$ (25)

$$H_{l,l-j,k} = \begin{cases} 1, if\ k = 0, 1, ..., l-j-1 \\ 0, & otherwise \end{cases}.$$ (26)

1     **Proof.** Given the probability distribution of $R_{i-1}^{(1,2)}$, $P\left(R_{i-1}^{(1,2)} = j\right), \forall j \in \{0, 1, 2, ..., l-1\}$; CV

2     penetration rate $p$; and identity of a vehicle being either a CV or a non-CV following a Bernoulli

3     distribution with parameter $p$, $Bernoulli(p)$, the following possibilities can be derived.

4     If $R_{i-1}^{(1,2)} = 0$, $R_{i-1}^{(1)}$ must be 0. The probability for this case is $P(R_{i-1}^{(1,2)} = 0)$.

5     If $R_{i-1}^{(1,2)} = 1$, $R_{i-1}^{(1)}$ can either be 0 or 1. If $R_{i-1}^{(1)} = 0$, i.e., the constrained vehicle is a non-CV, the

6     corresponding probability is $P(R_{i-1}^{(1,2)} = 1)(1-p)$. If $R_{i-1}^{(1)} = 1$, i.e., the constrained vehicle is a

7     CV, the corresponding probability is $P(R_{i-1}^{(1,2)} = 1)p$.

8     If $R_{i-1}^{(1,2)} = 2$, $R_{i-1}^{(1)}$ may be 0, 1, or 2. If $R_{i-1}^{(1)} = 0$, i.e., all constrained vehicles are non-CVs, the

9     corresponding probability is $P(R_{i-1}^{(1,2)} = 2)(1-p)^2$. If $R_{i-1}^{(1)} = 1$, i.e., the first vehicle in $R_{i-1}^{(1,2)}$ is

10     a CV and the second vehicle is a non-CV, the corresponding probability is $P(R_{i-1}^{(1,2)} = 2)p(1-p)$.

11     If $R_{i-1}^{(1)} = 2$, which means that the first vehicle can either be a CV or a non-CV but the second

12     vehicle must be a CV, the corresponding probability is $P(R_{i-1}^{(1,2)} = 2)p$. Similarly, if $R_{i-1}^{(1,2)} = k$,

13     $\forall k \in \mathbb{Z}$, all possible value of $R_{i-1}^{(1)}$ can be enumerated along with their corresponding probabilities,

14     as indicated in Table 4.

15     Table 4. Enumeration of all possible values of $R_{i-1}^{(1)}$ and corresponding probabilities when $R_{i-1}^{(1,2)} = k$.

| $R_{i-1}^{(1)}$ | Probability |
|---|---|
| 0 | $P(R_{i-1}^{(1,2)} = k)(1-p)^k$ |
| 1 | $P(R_{i-1}^{(1,2)} = k)p(1-p)^{k-1}$ |
| 2 | $P(R_{i-1}^{(1,2)} = k)p(1-p)^{k-2}$ |
| $\vdots$ | $\vdots$ |
| $k$ | $P(R_{i-1}^{(1,2)} = k)p$ |

16     By grouping the probabilities of the same value of $R_{i-1}^{(1)}$ from the enumerations of different values

17     of $R_{i-1}^{(1,2)}$, the probability distribution of $R_{i-1}^{(1)}$ can be obtained as follows:

$$\begin{aligned} P\left(R_{i-1}^{(1)} = 0\right) &= P\left(R_{i-1}^{(1,2)} = 0\right) + P\left(R_{i-1}^{(1,2)} = 1\right)(1-p) + \cdots + P\left(R_{i-1}^{(1,2)} = k\right)(1-p)^k \\ &= \left\| \left[P\left(R_{i-1}^{(1,2)} = 0\right), P\left(R_{i-1}^{(1,2)} = 1\right), ..., P\left(R_{i-1}^{(1,2)} = k\right)\right] \odot [1, 1-p, ..., (1-p)^k] \right\|_1, \end{aligned}$$ (27)

$$\begin{aligned} P\left(R_{i-1}^{(1)} = 1\right) &= P\left(R_{i-1}^{(1,2)} = 1\right)p + P\left(R_{i-1}^{(1,2)} = 2\right)p(1-p) + \cdots + P\left(R_{i-1}^{(1,2)} = k\right)p(1-p)^{k-1} \\ &= \left\| \left[P\left(R_{i-1}^{(1,2)} = 1\right), P\left(R_{i-1}^{(1,2)} = 2\right), ..., P\left(R_{i-1}^{(1,2)} = k\right)\right] \odot [p, p(1-p), ..., p(1-p)^{k-1}] \right\|_1, \end{aligned}$$ (28)

18                                          $\vdots$

$$P(R_{i-1}^{(1)} = k) = P(R_{i-1}^{(1,2)} = k)p = \left\| [P(R_{i-1}^{(1,2)} = k)] \odot [p] \right\|_1. \tag{29}$$

Alternatively, the probability distribution of $R_{i-1}^{(1)}$ can be written as in Eqs. (22)–(26).

**QED.**

## 4.3 Estimation of $Q_i$

This subsection describes the establishment of the MCQL model. The probability distributions of $N_i$ and $R_{i-1}^{(1)}$, $P(N_i = j) = b_j, \forall j \in \{0, 1, 2, ..., s + s_1 - 2\}$ and $P\left(R_{i-1}^{(1)} = j\right) = c_j, \forall j \in \{0, 1, 2, ..., l - 1\}$, respectively, can be expressed in the following vector forms:

$$\boldsymbol{f}_{N_i} = [b_0, b_1, ..., b_{s+s_1-2}], \tag{30}$$

and

$$\boldsymbol{f}_{R_{i-1}^{(1)}} = \left[ c_0, c_1, ..., c_{l-1}, \underbrace{0, 0, ..., 0}_{s+s_1-l-1} \right]. \tag{31}$$

The probability distribution of $Q_i$, $P(Q_i = j) = d_j, \forall j \in \{0, 1, 2, ..., s + s_1 - l - 1\}$, is given by

$$\left[\boldsymbol{f}_{Q_i} \quad \boldsymbol{\varepsilon}\right] = F^{-1}\left[ F\left(\boldsymbol{f}_{N_i}\right) \oslash F\left(\boldsymbol{f}_{R_{i-1}^{(1)}}\right) \right], \tag{32}$$

where $\oslash$ represents element-wise division; $\boldsymbol{\varepsilon}$ is a redundant vector; and $\boldsymbol{f}_{Q_i} = \left[d_0, d_1, ..., d_{s+s_1-l-1}\right]$ is the target probability vector for $Q_i$.

**Proof.** Rearranging Eq. (8) yields the following expression:

$$N_i = R_{i-1}^{(1)} + Q_i. \tag{33}$$

According to the convolutional constrained queue model,

$$\boldsymbol{f}_{N_i} = F^{-1}\left[ F\left(\boldsymbol{f}_{R_{i-1}^{(1)}}\right) \odot F([\boldsymbol{f}_{Q_i} \quad \boldsymbol{\varepsilon}]) \right]. \tag{34}$$

Applying DFT to both sides of Eq. (34), the following expression can be obtained:

$$F(\boldsymbol{f}_{N_i}) = F\left(\boldsymbol{f}_{R_{i-1}^{(1)}}\right) \odot F([\boldsymbol{f}_{Q_i} \quad \boldsymbol{\varepsilon}]). \tag{35}$$

As $\boldsymbol{f}_{R_{i-1}^{(1)}} > \boldsymbol{0}$, $F\left(\boldsymbol{f}_{R_{i-1}^{(1)}}\right) \neq \boldsymbol{0}$. Thus, it can be shown that

$$\left[\boldsymbol{f}_{Q_i} \quad \boldsymbol{\varepsilon}\right] = F^{-1}\left[ F\left(\boldsymbol{f}_{N_i}\right) \oslash F\left(\boldsymbol{f}_{R_{i-1}^{(1)}}\right) \right]. \tag{36}$$

**QED.**

The derived distribution of $Q_i$ can be substituted into the PPR model to estimate the uncertainty in the CV penetration rate.

## 5 Numerical Experiments

Detailed numerical experiments were conducted to evaluate the effectiveness of the proposed MCQL model in estimating the uncertainty in the CV penetration rate. The probability distributions of $N_i^{(4)}$, which represent an essential input for the MCQL model, can be estimated by either the PDT model or the CDT model. The MCQL models combined with the PDT and CDT

14

models to capture the residual-vehicle effects are designated as MCQL-P and MCQL-C model, respectively. As mentioned previously, if the PDT or CDT model was directly used in the uncertainty estimation of the CV penetration rate, the effects of the residual vehicle would be ignored. Thus, the performances of the MCQL-P and MCQL-C models were compared with the PDT and CDT models, respectively, across various cases involving different signal plans, volume-to-capacity (V/C) ratios, and CV penetration rates. For each case, 1,000 cycles were simulated. By applying the SSDPRE method (Wong et al., 2019) to each cycle, 1,000 estimates of the CV penetration rates were obtained. The variance of these estimates served as the ground truth for evaluating the accuracy and effectiveness of the MCQL-P and MCQL-C models.

To realistically mimic the queuing process, experiments were conducted using the VISSIM platform in a Windows 10 environment, over a machine equipped with an Intel Core i7-10700 CPU. The vehicles approached the signalized intersection through a single-lane link with a length of 1 km. The cycle length was set as 60 s, with the green period always ending with a 3-s amber phase. Vehicle generation followed a Poisson distribution, and the saturation headway was determined to be 1.59 s. All vehicles in the experiments were cars, and default values were maintained for the remaining settings. In general, drivers' reaction times and vehicles' acceleration and deceleration times lead to a net loss of red time and dissipation time in capturing the constrained queue sets. Ignoring these braking and start-up motions would lead to overestimation of the constrained queue length. To account for these factors, Jia et al. (2023) introduced a constant time-loss model to calibrate the net loss of red time. Based on the simulations, the net losses of red times of the PDT and CDT models were determined to be 5.820 s and 9.141 s, respectively. Details of the calibration procedure can be found in the work of Jia et al. (2023).

Table 5 presents the results obtained using the PDT and MCQL-P models for simulation cases featuring a constant CV penetration rate of 0.4 with varying signal plans and V/C ratios. Table 6 summarizes the results of the PDT and MCQL-P models for simulation cases with a fixed signal plan, including a 30-second red period, and varying V/C ratios and CV penetration rates. The findings clearly illustrate that the proposed MCQL-P model outperforms the PDT model in terms of absolute percentage errors (APE) while incurring minimal computational time overhead. This superiority is particularly evident in scenarios characterized by high V/C ratios, where the residual-vehicle effect is substantial. These results emphasize the importance of incorporating residual-vehicle effects for accurately estimating the uncertainty in CV penetration rate. Even in low V/C ratio scenarios, the proposed MCQL-P model consistently performs comparably or slightly better than the pure PDT model, further demonstrating the universality of the proposed model.

Table 7 compares the performances of the CDT and MCQL-C models in simulation cases with a constant CV penetration rate of 0.4 but different signal plans and V/C ratios. Table 8 summarizes the results of the CDT and MCQL-C models for simulation cases with a fixed signal plan including a 30-s red period but different V/C ratios and CV penetration rates. The CDT model, which is a simplified model to estimate the distribution of $N_i^{(4)}$, was less accurate but more robust than the PDT model. In other words, the CDT model was less sensitive than the PDT model to residual-vehicle effects. As expected, the proposed MCQL-C model achieved similar or slightly improved results compared with the CDT model. Notably, the MCQL-P outperformed the MCQL-C model. Overall, using the proposed MCQL model, the PPR model can be applied to all undersaturation scenarios, irrespective of the presence of residual vehicles. This improvement can enhance the practicality of the PPR model.

Table 5. Comparative performance of PDT and proposed MCQL-P models in scenarios with varying signal plans and V/C ratios.

| r | V/C | Ground truth | PDT w/o residual vehicles | | Proposed MCQL-P w/ residual vehicles | |
|---|---|---|---|---|---|---|
| | | Variance | Variance | APE (%) | Variance | APE (%) |
| 15 | 0.3 | 0.17408 | 0.16213 | 6.86 | 0.15868 | 8.85 |
| 15 | 0.5 | 0.16787 | 0.16000 | 4.69 | 0.16371 | 2.48 |
| 15 | 0.7 | 0.14107 | 0.12922 | 8.40 | 0.13617 | 3.47 |
| 15 | 0.95 | 0.09138 | 0.10175 | 11.35 | 0.09039 | 1.08 |
| 30 | 0.3 | 0.17076 | 0.17246 | 1.00 | 0.17339 | 1.54 |
| 30 | 0.5 | 0.12779 | 0.13013 | 1.83 | 0.13371 | 4.63 |
| 30 | 0.7 | 0.08236 | 0.08345 | 1.32 | 0.08250 | 0.17 |
| 30 | 0.95 | 0.04013 | 0.05666 | 41.19 | 0.04072 | 1.47 |
| 45 | 0.3 | 0.18275 | 0.17757 | 2.83 | 0.17728 | 2.99 |
| 45 | 0.5 | 0.15149 | 0.16055 | 5.98 | 0.15993 | 5.57 |
| 45 | 0.7 | 0.10344 | 0.12411 | 19.98 | 0.11097 | 7.28 |
| 45 | 0.95 | 0.09052 | 0.11762 | 29.94 | 0.10026 | 10.76 |
| Mean computing time (s) | | - | - | 8.076 | - | 8.106 |

Table 6. Comparative performance of PDT and proposed MCQL-P models in scenarios with varying V/C ratios and CV penetration rates.

| V/C | p | Ground truth | PDT w/o residual vehicles | | Proposed MCQL-P w/ residual vehicles | |
|---|---|---|---|---|---|---|
| | | Variance | Variance | APE (%) | Variance | APE (%) |
| 0.3 | 0.1 | 0.08075 | 0.07598 | 5.91 | 0.07571 | 6.24 |
| 0.3 | 0.4 | 0.17076 | 0.17246 | 1.00 | 0.17339 | 1.54 |
| 0.3 | 0.7 | 0.12647 | 0.13526 | 6.95 | 0.13668 | 8.07 |
| 0.5 | 0.1 | 0.07639 | 0.07244 | 5.17 | 0.07307 | 4.35 |
| 0.5 | 0.4 | 0.12779 | 0.13013 | 1.83 | 0.13371 | 4.63 |
| 0.5 | 0.7 | 0.09021 | 0.09200 | 1.98 | 0.09517 | 5.50 |
| 0.7 | 0.1 | 0.05813 | 0.06097 | 4.89 | 0.06044 | 3.97 |
| 0.7 | 0.4 | 0.08236 | 0.08345 | 1.32 | 0.08250 | 0.17 |
| 0.7 | 0.7 | 0.05347 | 0.05523 | 3.29 | 0.05474 | 2.38 |
| 0.95 | 0.1 | 0.03951 | 0.05071 | 28.35 | 0.03784 | 4.23 |
| 0.95 | 0.4 | 0.04013 | 0.05666 | 41.19 | 0.04072 | 1.47 |
| 0.95 | 0.7 | 0.02549 | 0.03710 | 45.55 | 0.02785 | 9.26 |
| Mean computing time (s) | | - | - | 8.735 | - | 8.765 |

Table 7. Comparative performance of CDT and proposed MCQL-C models in scenarios with varying signal plans and V/C ratios.

| r | V/C | Ground truth | CDT w/o residual vehicles | | Proposed MCQL-C w/ residual vehicles | |
|---|---|---|---|---|---|---|
| | | Variance | Variance | APE (%) | Variance | APE (%) |
| 15 | 0.3 | 0.17408 | 0.17034 | 2.15 | 0.15199 | 12.69 |
| 15 | 0.5 | 0.16787 | 0.17649 | 5.13 | 0.18383 | 9.51 |
| 15 | 0.7 | 0.14107 | 0.11758 | 16.65 | 0.14675 | 4.03 |
| 15 | 0.95 | 0.09138 | 0.06149 | 32.71 | 0.07326 | 19.83 |
| 30 | 0.3 | 0.17076 | 0.18209 | 6.64 | 0.18358 | 7.51 |
| 30 | 0.5 | 0.12779 | 0.13244 | 3.64 | 0.14203 | 11.14 |
| 30 | 0.7 | 0.08236 | 0.07350 | 10.76 | 0.07791 | 5.40 |
| 30 | 0.95 | 0.04013 | 0.04259 | 6.13 | 0.03351 | 16.50 |

| | | | | | | |
|---|---|---|---|---|---|---|
| 45 | 0.3 | 0.18275 | 0.18260 | 0.08 | 0.18141 | 0.73 |
| 45 | 0.5 | 0.15149 | 0.16671 | 10.05 | 0.16780 | 10.77 |
| 45 | 0.7 | 0.10344 | 0.12575 | 21.57 | 0.11335 | 9.58 |
| 45 | 0.95 | 0.09052 | 0.11829 | 30.68 | 0.10136 | 11.98 |
| Mean computing time (s) | - | - | - | 0.001 | - | 0.048 |

Table 8. Comparative performance of CDT and proposed MCQL-C models in scenarios with varying V/C ratios and CV penetration rates.

| V/C | p | Ground truth | CDT w/o residual vehicles | | Proposed MCQL-C w/ residual vehicles | |
|---|---|---|---|---|---|---|
| | | Variance | Variance | APE (%) | Variance | APE (%) |
| 0.3 | 0.1 | 0.08075 | 0.07834 | 2.98 | 0.07742 | 4.12 |
| 0.3 | 0.4 | 0.17076 | 0.18209 | 6.64 | 0.18358 | 7.51 |
| 0.3 | 0.7 | 0.12647 | 0.14387 | 13.76 | 0.14694 | 16.19 |
| 0.5 | 0.1 | 0.07639 | 0.07450 | 2.47 | 0.07602 | 0.48 |
| 0.5 | 0.4 | 0.12779 | 0.13244 | 3.64 | 0.14203 | 11.14 |
| 0.5 | 0.7 | 0.09021 | 0.09087 | 0.73 | 0.09961 | 10.42 |
| 0.7 | 0.1 | 0.05813 | 0.06068 | 4.39 | 0.06158 | 5.93 |
| 0.7 | 0.4 | 0.08236 | 0.07350 | 10.76 | 0.07791 | 5.40 |
| 0.7 | 0.7 | 0.05347 | 0.04584 | 14.27 | 0.04904 | 8.29 |
| 0.95 | 0.1 | 0.03951 | 0.04794 | 21.34 | 0.03659 | 7.39 |
| 0.95 | 0.4 | 0.04013 | 0.04259 | 6.13 | 0.03351 | 16.50 |
| 0.95 | 0.7 | 0.02549 | 0.02745 | 7.69 | 0.02268 | 11.02 |
| Mean computing time (s) | - | - | - | 0.001 | - | 0.050 |

# 6  Application

The performance of the proposed MCQL model was evaluated using the real-world NGSIM dataset. Additionally, an illustrative example of a stochastic CV-based adaptive signal control was considered to demonstrate that by modelling the residual vehicle effect in estimating CV penetration rate uncertainty, the proposed models can further improve the performance of practical traffic management scenarios under uncertain conditions.

## 6.1  Real-world validation

The proposed MCQL model was applied to the real-world NGSIM dataset to demonstrate its applicability and practicality. Specifically, trajectory data from two 15-min periods on November 8, 2006, i.e., 12:45–13:00 and 16:00–16:15, were extracted from the arterial road data for Peachtree Street in Atlanta, Georgia, USA. The validation focused on the southbound through-lane between Intersections 1 and 2. The cycle lengths and red durations for the two periods were 95 s, 100 s and 62 s, 64 s, respectively. The saturation headway was determined to be 2.044 s.

The CV penetration rate was set as 0.1, 0.4, or 0.7. Each vehicle was randomly assigned to be either a CV or non-CV, based on the pre-set CV penetration rates. Next, the SSDPRE method was applied to each constrained queue to obtain a set of estimated CV penetration rates that formed a distribution. Due to the limited amount of available data (only nine complete cycles in a 15-min period), the variance of the CV penetration rate distribution exhibited fluctuations caused by the random seeds used during the CV assignments. To mitigate the effect of this sampling error, the nine constrained queues were replicated 10,000 times. By repeating this CV assignment to the replications with different random seeds, 10,000 CV penetration rate distributions and corresponding variances were obtained. The mean of these 10,000 variances was considered the

ground truth for evaluation.

As described in the previous section, the performances of the MCQL-P and MCQL-C models were compared with those of the PDT and CDT models, respectively. The V/C ratios for the two periods were low, approximately 0.41 and 0.50, which indicated that the residual-vehicle effects were not significant. Consequently, the proposed methods were expected to perform comparably to the PDT and CDT models. Tables 9 and 10 indicate that the proposed methods exhibited similar or superior performance compared with their PDT and CDT counterparts. These results highlight the versatility of the MCQL model in handling undersaturation scenarios, regardless of the presence of residual vehicles.

Table 9. Performance comparison of the PDT and proposed MCQL models on the NGSIM dataset.

| Period | $p$ | Ground truth | PDT w/o residual vehicles | | Proposed MCQL-P w/ residual vehicles | |
|---|---|---|---|---|---|---|
| | | Variance | Variance | APE (%) | Variance | APE (%) |
| 12:45–13:00 | 0.1 | 0.07318 | 0.07336 | 0.25 | 0.07364 | 0.63 |
| | 0.4 | 0.15138 | 0.13161 | 13.06 | 0.13321 | 12.00 |
| | 0.7 | 0.11441 | 0.09217 | 19.44 | 0.09360 | 18.19 |
| Mean | | | | 10.91 | | 10.27 |
| 16:00–16:15 | 0.1 | 0.06249 | 0.06774 | 8.40 | 0.06809 | 8.96 |
| | 0.4 | 0.10184 | 0.10435 | 2.46 | 0.10588 | 3.97 |
| | 0.7 | 0.07313 | 0.06962 | 4.80 | 0.07082 | 3.16 |
| Mean | | | | 5.22 | | 5.36 |

Table 10. Performance comparison of the CDT and proposed MCQL models on the NGSIM dataset.

| Period | $p$ | Ground truth | CDT w/o residual vehicles | | Proposed MCQL-C w/ residual vehicles | |
|---|---|---|---|---|---|---|
| | | Variance | Variance | APE (%) | Variance | APE (%) |
| 12:45–13:00 | 0.1 | 0.07318 | 0.07463 | 1.98 | 0.07529 | 2.88 |
| | 0.4 | 0.15138 | 0.13321 | 12.00 | 0.13722 | 9.35 |
| | 0.7 | 0.11441 | 0.09156 | 19.97 | 0.09516 | 16.83 |
| Mean | | | | 11.32 | | 9.69 |
| 16:00–16:15 | 0.1 | 0.06249 | 0.06852 | 9.65 | 0.06940 | 11.06 |
| | 0.4 | 0.10184 | 0.10248 | 0.63 | 0.10641 | 4.49 |
| | 0.7 | 0.07313 | 0.06624 | 9.42 | 0.06928 | 5.26 |
| Mean | | | | 6.57 | | 6.94 |

## 6.2  Illustrative application of stochastic CV-based adaptive signal control

To demonstrate the importance of modeling the residual-vehicle effects, a simple example application of stochastic CV-based adaptive signal control was implemented based on VISSIM. Due to limited computational resources, only the more efficient CDT and MCQL-C models were selected for the illustration. Nevertheless, in real-world scenarios with access to large-scale computer clusters and high-performance GPUs, both PDT and MCQL-P models can be practically applied. Two adaptive signal control schemes were considered. In Scheme A, the residual-vehicle effects were not considered in the CV penetration rate uncertainty estimation, and the optimal signal plan was identified. In Scheme B, residual-vehicle effects were incorporated into the estimation of the CV penetration rate uncertainty and stochastic CV-based adaptive signal

1  optimizations.

2  The simulation involved a crossroad with two approaches at an isolated intersection. The traffic
3  demands for both approaches were generated using Poisson distributions. Two traffic demand
4  settings were implemented: in Setting 1, approaches 1 and 2 had average flow rates of 800 and 400
5  vehicles per hour, respectively; and in Setting 2, approaches 1 and 2 had average flow rates of
6  1,200 and 600 vehicles per hour, respectively. The V/C ratios at the intersection for Settings 1 and
7  2 were approximately 0.61 and 0.92, respectively. In other words, the residual-vehicle effects were
8  relatively light and significant in Settings 1 and 2, respectively. The CV penetration rate was set
9  as 0.1 or 0.4, indicating that each generated vehicle had a 10% or 40% probability of being a CV
10 and 90% or 60% probability of being a non-CV. A simple red–green–amber signal structure was
11 used for each approach, with a cycle length of 60 s, an amber time of 3 s, and a clearance time of
12 5 s. The saturation headway was determined to be 1.59 s. Signal optimization was performed at
13 the end of each cycle based on the estimated real-time traffic demands.

14 For any approach $j$ at the isolated intersection, $\forall j \in \{1,2\}$, the total number of vehicle arrivals in
15 cycle $i$ on approach $j$, $M_{i,j}$, is the sum of the number of CVs and the number of non-CVs in cycle
16 $i$ on approach $j$ and can be estimated by

$$M_{i,j} = m_{i,j} + q_{i,j}(1 - p_{i,j})C, \tag{37}$$

17 where $m_{i,j}$ represents the number of CV arrivals in cycle $i$ on approach $j$, $q_{i,j}$ is the real-time
18 average arrival rate in cycle $i$ on approach $j$, $p_{i,j}$ is the real-time CV penetration rate in cycle $i$ on
19 approach $j$, and $C$ represents the cycle length. The product of $q_{i,j}$ and $(1 - p_{i,j})$ gives the average
20 arrival rate of non-CVs in cycle $i$ on approach $j$. Thus, $M_{i,j}$ and its variability depend on $p_{i,j}$ and
21 its variability. In addition to $m_{i,j}$, the CV environment enables the observation of $n_{i,j}$ and $\widetilde{N}_{i,j}$.
22 Based on the CDT or MCQL-C model and **Corollary 2** (presented in **Appendix A**), a likelihood
23 function can be established to maximize the probability of observing $n_{i,j}$ and $\widetilde{N}_{i,j}$ by estimating
24 the real-time average arrival rate, $q_{i,j}$, and CV penetration rate, $p_{i,j}$:

$$\max_{q_{i,j},p_{i,j}} \prod_{k=0}^{T} P(n_{i-k,j}, \widetilde{N}_{i-k,j}), \tag{38}$$

25 where $T = 0, 1, 2, \ldots, i - 1$ is the number of past cycles considered in the likelihood function
26 formulation (in this example, $T$ was set as 2). The maximum likelihood estimators $q_{i,j}^*$ and $p_{i,j}^*$
27 can be considered the estimated real-time average arrival rate in cycle $i$ on approach $j$, $q_{i,j}$, and
28 the estimated real-time CV penetration rate in cycle $i$ on approach $j$, $E(p_{i,j})$, respectively. These
29 estimates can then be used as inputs of the CDT or MCQL-C model and **Corollary 1** to estimate
30 the real-time CV penetration rate variance, $Var(p_{i,j})$.

31 In Scheme A, only the CDT model was used for variance estimation. The CV penetration rate $p_{i,j}$
32 was assumed to follow a beta distribution with parameters $E(p_{i,j})$ and $Var(p_{i,j})$, as it was
33 confined between 0 and 1. Monte Carlo sampling was performed to sample 1,000 pairs of the
34 possible CV penetration rates for both approaches from the assumed beta distributions. Using Eq.
35 (37), 1,000 pairs of possible traffic demands were estimated based on the sampled CV penetration
36 rates. For each pair of traffic demands, the predicted real-time delays of the two approaches for
37 Scheme A in cycle $i + 1$, $D_{i+1,1}^A$ and $D_{i+1,2}^A$, were evaluated using Eqs. (B3) and (B4) (presented in
38 **Appendix B**), respectively, which ignored the possible residual-vehicle effects. For any signal

19

1    plan, the average total delay over the 1,000 possible traffic demand pairs for the intersection,
2    $E\big(D^A_{i+1,1} + D^A_{i+1,2}\big)$, was considered the objective function, as shown in Eq. (39). The optimal
3    signal plan was obtained by solving the following minimization problem using a simple line search
4    method:

$$
\begin{aligned}
\min_{g_{i+1,1},\,g_{i+1,2}} \quad & E\big(D^A_{i+1,1} + D^A_{i+1,2}\big) \\
s.t. \quad & g_{i+1,1} + g_{i+1,2} = 52 \\
& g_{i+1,1} \geq 5 \\
& g_{i+1,2} \geq 5
\end{aligned}
\tag{39}
$$

5    After the initial 30 warm-up cycles with a fixed signal plan, the signal plan was optimized at the
6    end of each cycle using Eq. (39). A simulation involving 1,000 cycles was conducted, and the
7    actual delays of all vehicles were recorded. The results for the two traffic demand settings are
8    presented in Table 11.

9    In Scheme B, the proposed MCQL-C model was used to estimate the CV penetration rate
10    uncertainty. The traffic demand estimations were identical to those in Scheme A. However, the
11    real-time delays of the two approaches for Scheme B in cycle $i + 1$, $D^B_{i+1,1}$ and $D^B_{i+1,2}$, were
12    predicted using Eqs. (B5) and (B6) (presented in **Appendix B**), respectively, which take into
13    account the estimated residual vehicle distribution as the initial state for the next cycle and consider
14    the presence of residual vehicles in the next cycle. The optimal signal plan with the least average
15    total delay over the 1,000 possible traffic demand pairs for the intersection, $E\big(D^B_{i+1,1} + D^B_{i+1,2}\big)$,
16    was determined using a simple line search method, while adhering to the same set of constraints
17    specified in Eq. (39). Similar to Scheme A, after the initial 30 warm-up cycles, the signal plan was
18    optimized for 1,000 cycles at the end of each cycle according to the described control scheme. The
19    results for the two traffic demand settings are presented in Table 11.

20    Table 11. Comparison of stochastic CV-based adaptive signal control schemes with and
21    without consideration of residual-vehicle effects.

| Traffic demand | Metric | $p$ | Scheme A, w/o residual vehicles | Scheme B, w/ residual vehicles | Improvement (%) |
|---|---|---|---|---|---|
| 800 veh/h and 400 veh/h | Average actual delay (s) | 0.1 | 88.4 | 27.2 | 69.2 |
| | Maximum actual delay (s) | 0.1 | 1759.8 | 249.0 | 85.9 |
| | Variance in actual delay (s²) | 0.1 | 39,891.5 | 663.0 | 98.3 |
| | Average actual delay (s) | 0.4 | 23.1 | 21.3 | 7.8 |
| | Maximum actual delay (s) | 0.4 | 171.5 | 149.8 | 12.7 |
| | Variance in actual delay (s²) | 0.4 | 459.6 | 300.5 | 34.6 |
| 1200 veh/h and 600 veh/h | Average actual delay (s) | 0.1 | 307.3 | 171.7 | 44.1 |
| | Maximum actual delay (s) | 0.1 | 2594.7 | 711.0 | 72.6 |
| | Variance in actual delay (s²) | 0.1 | 127,432.5 | 15,805.7 | 87.6 |
| | Average actual delay (s) | 0.4 | 279.3 | 161.2 | 42.3 |
| | Maximum actual delay (s) | 0.4 | 2305.2 | 664.1 | 71.2 |
| | Variance in actual delay (s²) | 0.4 | 124,045.2 | 13,413.8 | 89.2 |

22    Table 11 presents the average actual delays, maximum actual delays, and variances in actual delay
23    for Schemes A and B under various combinations of traffic demand settings and CV penetration
24    rates. The last column in the table shows the improvement of Scheme B relative to Scheme A
25    across all three metrics. The results clearly indicate that Scheme B consistently and significantly
26    outperformed Scheme A across all metrics under different scenarios. The notable improvement
27    was attributable to the incorporation of residual-vehicle effects in Scheme B. Overall, this simple

example of stochastic CV-based signal control clearly demonstrates the importance of considering residual-vehicle effects in model estimation and system optimizations.

# 7 Conclusions

The CV penetration rate is a critical parameter in CV-based transportation applications. Accurately estimating the uncertainty in the CV penetration rate is essential for developing unbiased transport models and deriving optimal solutions for transport system optimizations. Recently, the PPR model has been proposed as a framework for accurately modeling the uncertainty in the CV penetration rate. However, the method used to estimate the constrained queue length distribution in the PPR model does not consider the complex effects of residual vehicles. Neglecting these effects may lead to improper estimates for the constrained queue length distribution, resulting in inaccurate estimation of the CV penetration rate uncertainty. This study aims to address this research gap by incorporating the effects of residual vehicles in the estimation of the constrained queue length. This framework enables the application of the PPR model in undersaturated traffic conditions, regardless of the presence of residual vehicles. The proposed approach decomposes a full constrained queue into four vehicle groups: observable constrained residual vehicles, unobservable constrained residual vehicles, unconstrained residual vehicles, and new arrivals. The residual-vehicle effects are modeled using a novel Markov chain process and four analytical sub-models within the MCQL model, including the residual-vehicle model, convolutional constrained queue model, constrained residual queue model, and observable residual queue model. The effectiveness of the proposed models is demonstrated through comprehensive VISSIM simulations and real-world experiments. Furthermore, a practical example of stochastic CV-based adaptive signal control is presented to highlight the importance of modeling residual-vehicle effects in improving the system performance.

# Acknowledgments

# References

Ambühl, L., Menendez, M., 2016. Data fusion algorithm for macroscopic fundamental diagram estimation. *Transportation Research Part C: Emerging Technologies* 71, 184-197.

Argote, J., Christofa, E., Xuan, Y., Skabardonis, A., 2011. Estimation of measures of effectiveness based on connected vehicle data. In: *Proceedings of the 14th International IEEE Conference on Intelligent Transportation System* 1767-1772.

Cao, Y., Tang, K., Sun, J., Ji, Y., 2021. Day-to-day dynamic origin–destination flow estimation using connected vehicle trajectories and automatic vehicle identification data. *Transportation Research Part C: Emerging Technologies* 129, 103241.

Comert, G., 2013. Simple analytical models for estimating the queue lengths from probe vehicles at traffic signals. *Transportation Research Part B: Methodological* 55, 59-74.

Comert, G., 2016. Queue length estimation from probe vehicles at isolated intersections:

Estimators for primary parameters. *European Journal of Operational Research* 252, 502-521.

Comert, G., Cetin, M., 2009. Queue length estimation from probe vehicle location and the impacts of sample size. *European Journal of Operational Research* 197, 196-202.

Comert, G., Cetin, M., 2011. Analytical evaluation of the error in queue length estimation at traffic signals from prove vehicle data. *IEEE Transactions on Intelligent Transportation Systems* 12(2), 563-573.

Du, J., Rakha, H., Gayah, V.V., 2016. Deriving macroscopic fundamental diagrams from probe data: Issues and proposed solutions. *Transportation Research Part C: Emerging Technologies* 66, 136-149.

Federal Highway Administration, 2006. Next generation simulation: Peachtree Street dataset. Accessed June 25, 2022, https://data.transportation.gov/Automobiles/Next-Generation-Simulation-NGSIM-Program-Peachtree/mupt-aksf.

Feng, Y., Head, K.L., Khoshmagham, S., Zamanipour, M., 2015. A real-time adaptive signal control in a connected vehicle environment. *Transportation Research Part C: Emerging Technologies* 55, 460-473.

Geroliminis, N., Daganzo, C.F., 2008. Existence of urban-scale macroscopic fundamental diagrams: Some experimental findings. *Transportation Research Part B: Methodological* 42(9), 759-770.

Han, K., Liu, H., Gayah, V. V., Friesz, T. L., Yao, T., 2016. A robust optimization approach for dynamic traffic signal control with emission considerations. *Transportation Research Part C: Emerging Technologies* 70, 3-26.

Hao, P., Ban, X.J., Guo, D., Ji, Q., 2014. Cycle-by-cycle intersection queue length distribution estimation using sample travel times. *Transportation Research Part B: Methodological* 68, 185-204.

Iqbal, M.S., Hadi, M., Xiao, Y., 2018. Effect of link-level variations of connected vehicles (CV) proportions on the accuracy and reliability of travel time estimation. *IEEE Transactions on Intelligent Transportation Systems* 20(1), 87-96.

Jenelius, E., Koutsopoulos, H.N., 2013. Travel time estimation for urban road networks using low frequency probe vehicle data. *Transportation Research Part B: Methodological* 53, 64-81.

Jenelius, E., Koutsopoulos, H.N., 2015. Probe vehicle data sampled by time or space: Consistent travel time allocation and estimation. *Transportation Research Part B: Methodological* 71, 120-137.

Jia, S., Wong, S.C., Wong, W., 2023. Uncertainty estimation of connected vehicle penetration rate. *Transportation Science* 57(5), 1160-1176.

Khan, S.M., Dey, K.C., Chowdhury, M., 2017. Real-time traffic state estimation with connected vehicles. *IEEE Transactions on Intelligent Transportation Systems* 18(7), 1687-1699.

Lu, Y., Xu, X., Ding, C., Lu, G., 2019. A speed control method at successive signalized

intersections under connected vehicles environment. *IEEE Intelligent Transportation Systems Magazine* 11(3), 117-128.

Meng, F., Wong, S.C., Wong, W., Li, Y.C., 2017a. Estimation of scaling factors for traffic counts based on stationary and mobile sources of data. *International Journal of Intelligent Transportation Systems Research* 15(3), 180-191.

Meng, F., Wong, W., Wong, S.C., Pei, X., Li, Y.C., Huang, H., 2017b. Gas dynamic analogous exposure approach to interaction intensity in multiple-vehicle crash: Case study of crashes involving taxis. *Analytic Methods in Accident Research* 16, 90-103.

Mousa, S. R., Ishak, S., 2017. An extreme gradient boosting algorithm for freeway short-term travel time prediction using basic safety messages of connected vehicles. In: *Transportation Research Board 96th Annual Meeting*, 2017, Washington, DC, United States.

Rahmani, M., Jenelius, E., Koutsopoulos, H.N., 2015. Non-parametric estimation of route travel time distributions from low-frequency floating car data. *Transportation Research Part C: Emerging Technologies* 58, 343-362.

Sen, S., Head, K.L., 1997. Controlled optimization of phases at an intersection. *Transportation Science* 31, 5-17.

Sun, C., Shen, X., & Moura, S., 2018, June. Robust optimal eco-driving control with uncertain traffic signal timing. In 2018 Annual American control conference (ACC) (pp. 5548-5553). IEEE.

Tian, D., Yuan, Y., Qi, H., Lu, Y., Wang, Y., Xia, H., He, A., 2015. A dynamic travel time estimation model based on connected vehicles. *Mathematical Problems in Engineering*, 2015, 903962.

Wang, P., Zhang, J., Deng, H., Zhang, M., 2020. Real-time urban regional route planning model for connected vehicles based on V2X communication. *Journal of Transport and Land Use* 13(1), 517-538.

Wang, X., Jerome, Z., Wang, Z., Zhang, C., Shen, S., Kumar, V.V., Bai, F., Krajewski, P., Deneau, D., Jawad, A., Jones, R., Piotrowicz, G. Liu, H.X., 2024. Traffic light optimization with low penetration rate vehicle trajectory data. *Nature Communications* 15, 1306.

Wong, W., Wong, S.C., 2015. Systematic bias in transport model calibration arising from the variability of linear data projection. *Transportation Research Part B: Methodological* 75, 1-18.

Wong, W., Wong, S.C., 2016a. Biased standard error estimations in transport model calibration due to heteroscedasticity arising from the variability of linear data projection. *Transportation Research Part B: Methodological* 88, 72-92.

Wong, W., Wong S.C., 2016b. Evaluation of the impact of traffic incidents using GPS data. *Proceedings of the Institution of Civil Engineers – Transport* 169(3), 148-162.

Wong, W., Wong S.C., 2016c. Network topological effects on the macroscopic Bureau of Public Roads function. *Transportmetrica A: Transport Science* 12(3), 272-296.

Wong, W., Wong, S.C., 2019. Unbiased estimation methods of nonlinear transport models based on linearly projected data. *Transportation Science* 53(3), 665-682.

Wong, W., Wong, S.C., Liu, X., 2019a. Bootstrap standard error estimations of nonlinear transport models based on linearly projected data. *Transportmetrica A: Transport Science* 15(2), 602-630.

Wong, W., Shen, S, Zhao, Y., Liu, X., 2019b. On the estimation of connected vehicle penetration rate based on single-source connected vehicle data. *Transportation Research Part B: Methodological* 126, 169-191.

Wong, W., Wong, S.C., Liu, X., 2021. Network topological effects on the macroscopic fundamental diagram. *Transportmetrica B: Transport Dynamics* 9(1), 376-398.

Yang, X., Lu, Y., Hao, W., 2017. Origin-destination estimation using probe vehicle trajectory and link counts. *Journal of Advanced Transportation* 2017, 4341532.

Yin, Y., 2008. Robust optimal traffic signal timing. *Transportation Research Part B: Methodological* 42(10), 911-924.

Zhang, L., Yin, Y., & Chen, S., 2013. Robust signal timing optimization with environmental concerns. *Transportation Research Part C: Emerging Technologies* 29, 55-71.

Zhao, Y., Zheng, J., Wong, W., Wang, X., Meng, Y., Liu, H.X., 2019a. Estimation of queue lengths, probe vehicle penetration rates, and traffic volumes at signalized intersections using probe vehicle trajectories. *Transportation Research Record* 2673(11), 660-670.

Zhao, Y., Zheng, J., Wong, W., Wang, X., Meng, Y., Liu, H.X., 2019b. Various methods for queue length and traffic volume estimation using probe vehicle trajectories. *Transportation Research Part C: Emerging Technologies* 107, 70-91.

Zhao, Y., Wong, W., Zheng, J., Liu, H.X., 2022. Maximum likelihood estimation of probe vehicle penetration rates and queue length distributions from probe vehicle data. *IEEE Transactions on Intelligent Transportation Systems* 23(7), 7628-7636.

# Appendix A. Corollaries of the PPR Model

**Corollary 1.** Given that $N \sim Pois(\lambda)$ and $n \sim B(N, p)$, $E(\tilde{p})$ and $Var(\tilde{p})$ can be defined as follows:

$$E(\tilde{p}) = \lim_{k \to +\infty} \left[ e^{-\lambda} p + \sum_{i=1}^{k} \sum_{j=1}^{i} \sum_{m=1}^{i-j+1} \frac{\lambda^i e^{-\lambda}}{i!} \binom{i-m}{j-1} p^j (1-p)^{i-j} S(i, N-j+1) \right] = p, \quad \text{(A1)}$$

$$Var(\tilde{p}) = \lim_{k \to +\infty} \left[ \sum_{i=1}^{k} \frac{\lambda^i e^{-\lambda}}{i!} V_2(i, p) \right]. \quad \text{(A2)}$$

**Corollary 2.** Given that $N \sim Pois(\lambda)$ and $n \sim B(N, p)$, the joint probability distribution of $n$ and $\widetilde{N}$ is

$$P(n = i, \widetilde{N} = j) = \begin{cases} \pi_0 + \displaystyle\sum_{z=1}^{k} \pi_z (1-p)^z, & i = 0, j = 0 \\[2em] \displaystyle\sum_{z=j}^{k} \pi_z \binom{j-1}{i-1} p^i (1-p)^{z-i}, \forall i, j = 1, 2, \ldots, k, j \geq i \end{cases}, \quad \text{(A3)}$$

where $\pi_z = P(N = z), \forall z = 0, 1, 2, \ldots, k$.

Proofs of the corollaries can be found in the work of Jia et al. (2023).

# Appendix B. Estimation of Real-time Delays

This appendix describes a method to estimate the real-time delays of two approaches to an intersection controlled by a simplified red–green–amber signal structure. To simplify the delay estimation, the vehicle arrivals are assumed to follow uniform distributions. Figure B illustrates the general cases for estimating real-time delays for the two approaches in two conditions: (1) undersaturation conditions without residual vehicles, and (2) temporary overflow conditions with residual vehicles.
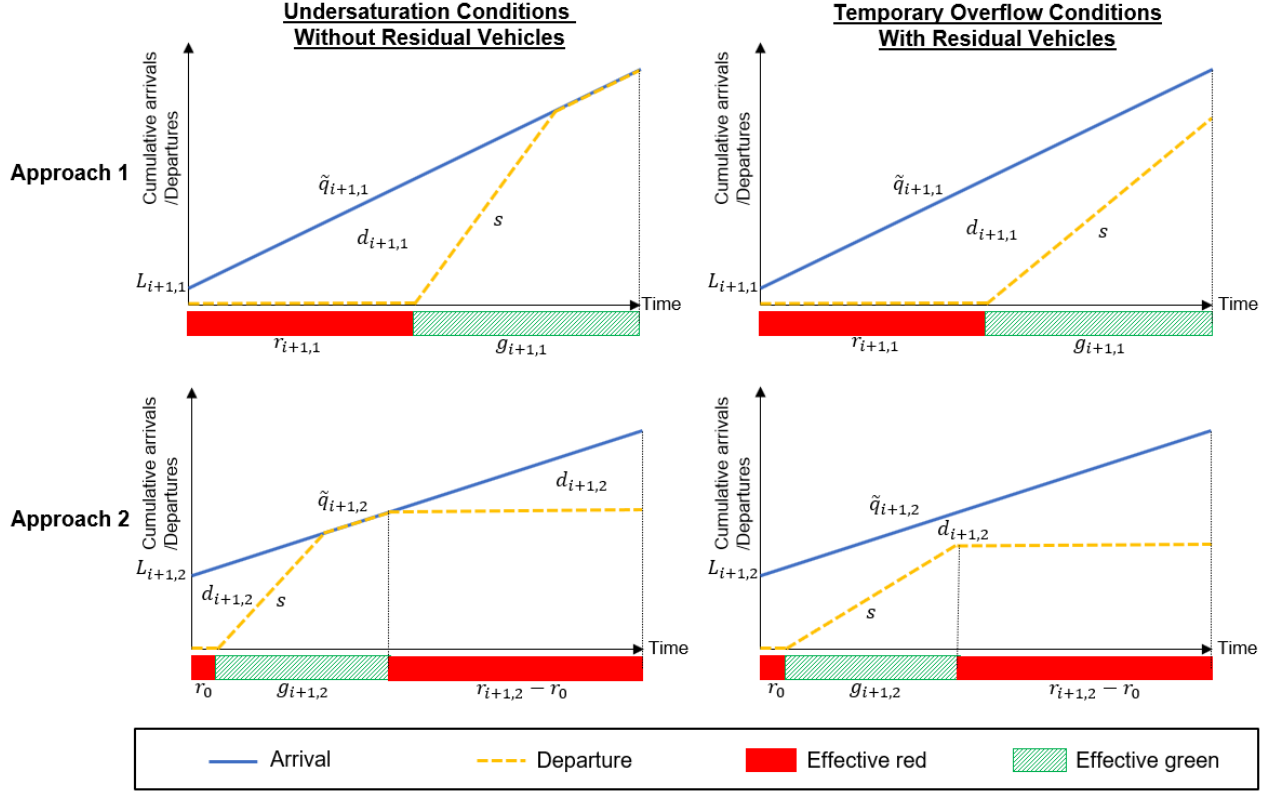


Figure B. Illustration of real-time delay estimations for two approaches in undersaturation conditions without residual vehicles and temporary overflow conditions with residual vehicles.

The real-time delays at the intersection are represented by the areas between the arriving and departing profiles. By applying simple geometry, the predicted real-time delays in cycle $i + 1$ for approaches 1 and 2, $d_{i+1,1}$ and $d_{i+1,2}$, can be estimated using Eqs. (B1) and (B2), respectively.

$$d_{i+1,1} = \begin{cases} L_{i+1,1}r_{i+1,1} + \dfrac{\tilde{q}_{i+1,1}r_{i+1,1}^2}{2} + \dfrac{\left(L_{i+1,1} + \tilde{q}_{i+1,1}r_{i+1,1}\right)^2}{2(s - \tilde{q}_{i+1,1})} & if\ \tilde{q}_{i+1,1} \leq \dfrac{sg_{i+1,1} - L_{i+1,1}}{C} \\ \dfrac{(2L_{i+1,1} + \tilde{q}_{i+1,1}r_{i+1,1})r_{i+1,1} + g_{i+1,1}[2L_{i+1,1} + \tilde{q}_{i+1,1}(2C - g_{i+1,1}) - sg_{i+1,1}]}{2} & if\ \tilde{q}_{i+1,1} > \dfrac{sg_{i+1,1} - L_{i+1,1}}{C} \end{cases} , \quad (B1)$$

$$d_{i+1,2} = \begin{cases} \dfrac{1}{2}(2L_{i+1,2} + \tilde{q}_{i+1,2}r_0)r_0 + \dfrac{(L_{i+1,2} + \tilde{q}_{i+1,2}r_0)^2}{2(s - \tilde{q}_{i+1,2})} + \dfrac{1}{2}\tilde{q}_{i+1,2}(r_{i+1,2} - r_0)^2 & if\ \tilde{q}_{i+1,2} \leq \dfrac{sg_{i+1,2} - L_{i+1,2}}{r_0 + g_{i+1,2}} \\ \dfrac{1}{2}(2L_{i+1,2} + \tilde{q}_{i+1,2}C)C - \dfrac{1}{2}(2C - g_{i+1,2} - 2r_0)sg_{i+1,2} & if\ \tilde{q}_{i+1,2} > \dfrac{sg_{i+1,2} - L_{i+1,2}}{r_0 + g_{i+1,2}} \end{cases} , \quad (B2)$$

where $r_{i+1,j}$ is the effective red in cycle $i + 1$ on approach $j$; $g_{i+1,j}$ is the effective green in cycle $i + 1$ on approach $j$; $\tilde{q}_{i+1,j}$ is the predicted real-time average arrival rate in cycle $i + 1$ on approach $j$, with $\tilde{q}_{i+1,j} = M_{i,j}/C$; $s$ is the saturation flow rate; $C$ is the cycle length; $r_0$ is the clearance loss time consisting of a part of the amber period and all-red clearance time (set as 4 s in this case); and

1      $L_{i+1,1}$ and $L_{i+1,2}$ are initial states of approaches 1 and 2 in cycle $i+1$, respectively.

2      In Scheme A, the residual-vehicle effects are not considered. Thus, no residual vehicle is carried
3      over from cycle to cycle, implying that $L_{i+1,1} = 0$ and $L_{i+1,2} = \tilde{q}_{i+1,2}(r_{i,2} - r_0)$. By substituting
4      these values into Eqs. (B1) and (B2), the real-time delays of the two approaches in cycle $i+1$
5      under Scheme A, $D^A_{i+1,1}$ and $D^A_{i+1,2}$, respectively, can be predicted as follows:

$$D^A_{i+1,1} = \begin{cases} \dfrac{s\tilde{q}_{i+1,1}r^2_{i+1,1}}{2(s-\tilde{q}_{i+1,1})} & if\ \tilde{q}_{i+1,1} \le \dfrac{sg_{i+1,1}}{C} \\[2mm] \dfrac{\tilde{q}_{i+1,1}r^2_{i+1,1} + g_{i+1,1}[\tilde{q}_{i+1,1}(2C-g_{i+1,1}) - sg_{i+1,1}]}{2} & if\ \tilde{q}_{i+1,1} > \dfrac{sg_{i+1,1}}{C} \end{cases} \tag{B3}$$

$$D^A_{i+1,2} = \begin{cases} \dfrac{1}{2}[2\tilde{q}_{i+1,2}(r_{i,2}-r_0) + \tilde{q}_{i+1,2}r_0]r_0 + \dfrac{[\tilde{q}_{i+1,2}(r_{i,2}-r_0) + \tilde{q}_{i+1,2}r_0]^2}{2(s-\tilde{q}_{i+1,2})} + \dfrac{1}{2}\tilde{q}_{i+1,2}(r_{i+1,2}-r_0)^2 & if\ \tilde{q}_{i+1,2} \le \dfrac{sg_{i+1,2} - \tilde{q}_{i+1,2}(r_{i,2}-r_0)}{r_0 + g_{i+1,2}} \\[2mm] \dfrac{1}{2}[2\tilde{q}_{i+1,2}(r_{i,2}-r_0) + \tilde{q}_{i+1,2}C]C - \dfrac{1}{2}(2C - g_{i+1,2} - 2r_0)sg_{i+1,2} & if\ \tilde{q}_{i+1,2} > \dfrac{sg_{i+1,2} - \tilde{q}_{i+1,2}(r_{i,2}-r_0)}{r_0 + g_{i+1,2}} \end{cases} \tag{B4}$$

6      In Scheme B, the residual-vehicle effects are considered. Therefore, $L_{i+1,1} = E(R_{i,1})$ and $L_{i+1,2} =$
7      $\tilde{q}_{i+1,2}(r_{i,2} - r_0) + E(R_{i,2})$, where $R_{i,1}$ and $R_{i,2}$ are random variables representing the number of
8      residual vehicles for approaches 1 and 2 from cycle $i$, respectively. The values of $R_{i,1}$ and $R_{i,2}$ can
9      be obtained from the proposed residual-vehicle model. Additionally, if cycle $i+1$ is predicted to
10     be in the temporary overflow state, the residual vehicles from cycle $i+1$ will carry over to cycle
11     $i+2$. To capture these potential delays, the predicted delays in cycle $i+2$, $d_{i+2,1}$ and $d_{i+2,2}$, are
12     incorporated into the predicted real-time delays of the two approaches in cycle $i+1$ under Scheme
13     B, $D^B_{i+1,1}$ and $D^B_{i+1,2}$, respectively. Assuming identical traffic demands and signal plans in cycle $i+$
14     $2$, $d_{i+2,1}$ and $d_{i+2,2}$ can be readily obtained by substituting $L_{i+2,1} = E(R_{i,1}) + \tilde{q}_{i+1,1}C - sg_{i+1,1}$
15     and $L_{i+2,2} = E(R_{i,2}) + \tilde{q}_{i+1,2}C - sg_{i+1,2}$ in the corresponding predicted delay formulas derived
16     from Eqs. (B1) and (B2). The real-time delays of the two approaches in cycle $i+1$ under Scheme
17     B, $D^B_{i+1,1}$ and $D^B_{i+1,2}$, can be predicted as follows:

$$D^B_{i+1,1} = \begin{cases} d_{i+1,1} & if\ \tilde{q}_{i+1,1} \le \dfrac{sg_{i+1,1} - E(R_{i,1})}{C} \\[2mm] d_{i+1,1} + d_{i+2,1} & if\ \tilde{q}_{i+1,1} > \dfrac{sg_{i+1,1} - E(R_{i,1})}{C} \end{cases} \tag{B5}$$

$$D^B_{i+1,2} = \begin{cases} d_{i+1,2} & if\ \tilde{q}_{i+1,2} \le \dfrac{sg_{i+1,2} - \tilde{q}_{i+1,2}(r_{i,2}-r_0) + E(R_{i,2})}{r_0 + g_{i+1,2}} \\[2mm] d_{i+1,2} + d_{i+2,2} & if\ \tilde{q}_{i+1,2} > \dfrac{sg_{i+1,2} - \tilde{q}_{i+1,2}(r_{i,2}-r_0) + E(R_{i,2})}{r_0 + g_{i+1,2}} \end{cases} \tag{B6}$$

18