# Enforcing Priority in Schedule-based User Equilibrium Transit Assignment

Liyang Feng[1, 2], Hanlin Sun[2], Yu (Marco) Nie[3], Jun Xie[*1], and Jiayang Li[*2]

[1]School of Transportation and Logistics, Southwest Jiaotong University, Chengdu, China
[2]Department of Data and Systems Engineering, The University of Hong Kong, Hong Kong, China
[3]Department of Civil and Environmental Engineering, Northwestern University, IL, USA

## Abstract

Denied boarding in congested transit systems induces queuing delays and departure-time shifts that can reshape passenger flows. Correctly modeling these responses in transit assignment hinges on the enforcement of two priority rules: continuance priority for onboard passengers and first-come-first-served (FCFS) boarding among waiting passengers. Existing schedule-based models typically enforce these rules through explicit dynamic loading and group-level expected costs, yet discrete vehicle runs can induce nontrivial within-group cost differences that undermine behavioral consistency. We revisit the implicit-priority framework of Nguyen et al. (2001) ["A Modeling Framework for Passenger Assignment on a Transport Network with Timetables." *Transportation Science* 35(3): 238–249], which, by encoding boarding priority through the notion of available capacity, characterizes route and departure choices based on realized personal (rather than group-averaged) travel experiences. However, the framework lacks an explicit mathematical formulation and exact computational methods for finding equilibria. Here, we derive an equivalent nonlinear complementarity problem (NCP) formulation and establish equilibrium existence under mild conditions. We also show that multiple equilibria may exist, including behaviorally questionable ones. To rule out these artifacts, we propose a refined arc-level NCP formulation that not only corresponds to a tighter, behaviorally consistent equilibrium concept but also is more computationally tractable. We reformulate the NCP as a continuously differentiable mathematical program with equilibrium constraints (MPEC) and propose two solution algorithms. Numerical studies on benchmark instances and a Hong Kong case study demonstrate that the model reproduces continuance priority and FCFS queuing and captures departure-time shifts driven by the competition for boarding priority.

Keywords: schedule-based transit assignment; boarding priority; nonlinear complementarity problem

## 1 Introduction

In densely populated urban areas where mass transit is a primary travel mode, waiting passengers frequently fail to board incoming transit vehicles due to their limited capacities. In Hong Kong, for example, buses traveling from Hong Kong Island to Kowloon during the evening peak often fill up at early stops, leaving passengers at intermediate stations to watch several buses pass by without space to board, a phenomenon locally referred to as "ding jaa" (Hong Kong Bus Wiki, 2025). At Shahe Town, one of Beijing's

---

*Corresponding authors. E-mail: jun.xie@swjtu.edu.cn (Jun Xie); jiayangl@hku.hk (Jiayang Li).

largest suburban residential clusters, the crowding can be so severe that commuters wait up to 30 minutes at the metro station to enter a train (Beijing Daily News, 2023). Mitigating such crowdedness is therefore a central objective for transit planners and operators. To support effective decisions, however, one needs to reliably predict how passengers choose routes and adjust their departure times in response to the delays caused by such a capacity crunch. This task is usually addressed using a transit assignment model.

## 1.1 Challenges

Transit assignment models can be broadly categorized as frequency-based and schedule-based. The first class represents each line by its service frequency rather than exact departure times and assumes that all passengers at a station experience the same expected waiting time, regardless of their actual arrival time (Spiess and Florian, 1989; Wu et al., 1994; Cominetti and Correa, 2001; Cepeda et al., 2006; Xu et al., 2020, 2022). While these models may impose vehicle capacity as a constraint at the link level, they cannot represent the capacity-induced delays *physically*, especially those incurred when a passenger misses one or more passing vehicles. Schedule-based models are better equipped to deal with the physics of the vehicle capacity restriction, simply because they typically represent individual vehicle runs according to posted schedules. In these models, passengers may choose to depart from home according to the expected boarding time at the first transit stop, which may not coincide with the arrival time of the first vehicle (in other words, they explicitly take the queuing time into consideration). The enhanced realism, however, complicates the assignment problem. A key sticky issue is the priority rules that determine who gets a seat first when vehicle capacity is insufficient. These include *continuance priority*, under which passengers already on board must remain on the vehicle, and the *first-come-first-served (FCFS) rule*, according to which waiting passengers board in order of arrival (Hamdouch and Lawphongpanich, 2008).

A common approach to modeling priority is through dynamic network loading (DNL) (Nuzzolo et al., 2001; Poon et al., 2004; Papola et al., 2008; Hamdouch et al., 2011; Nuzzolo et al., 2012; Hamdouch et al., 2014; Cats et al., 2016; Gentile and Nökel, 2016; Yao et al., 2017; Cats and West, 2020). In these models, passengers are grouped according to their chosen route and departure time, and a DNL procedure simulates the boarding of passengers at all stops *explicitly* according to the priority rules. Because individuals within the same group may experience different actual costs due to vehicle capacity restrictions, a flow-weighted *expected* travel cost is often calculated as the representative cost of that group. The equilibrium is defined, accordingly, as a state where no group can reduce this expected travel cost by switching routes or departure times. However, in schedule-based transit models, service is inherently *discrete* in time. As a result, the cost difference within a group does not always diminish even as group size approaches zero, since tight residual capacity may still split a group across successive runs. This is more than a mathematical nuance: assuming passengers included in the same group behave identically is questionable even at the limit, because the experienced cost difference would be great enough to trigger behavioral deviation[1].

An alternative approach, which we refer to as implicit prioritization, was pioneered by Nguyen et al. (2001). Instead of simulating boarding events, their framework encodes priority *implicitly* through the notion of available capacity. Specifically, passengers already on board have the highest priority (continuance priority), followed by waiting passengers in the order they arrived at a stop (FCFS). The available capacity of a boarding arc is then defined as the remaining vehicle capacity after loading all passengers from arcs with higher priority. If the available capacity on a given arc is non-positive, passengers arriving via that

---

[1]In Appendix A we illustrate this issue more concretely with a simple example.

arc cannot board in that run and must wait for the next. Using this representation, Nguyen et al. (2001) defined a new equilibrium principle: a passenger can only switch to a route if all boarding arcs along that route have positive available capacity at equilibrium. In other words, switching is allowed only if doing so would not displace higher-priority passengers or overload any vehicle segment. In this way, some passengers successfully board while others may be left behind, and each person's incentives to adjust their decisions are separately modeled based on their *actual* experience.

Despite its conceptual elegance, the framework of Nguyen et al. (2001) is not fully developed. While it gives a well-defined equilibrium condition, a mathematical formulation amenable to analysis or computing remains elusive. No exact algorithm has been proposed to find an equilibrium satisfying the original priority constraints. Instead, Nguyen et al. (2001) solved an approximate problem that relaxes the priority constraints, leaving open the question of how to obtain an exact solution.

## 1.2 Our contributions

We propose an equivalent nonlinear complementarity problem (NCP) formulation for Nguyen et al.'s model (2001). Using this formulation, we prove the existence of an equilibrium state under mild conditions, a theoretical result that, to the best of our knowledge, has not been previously established. This result lays a foundation for applying the model in general transit assignment settings. We also show that multiple equilibria exist and, more importantly, some of these are behaviorally unrealistic.

In light of this finding, we further propose a refined NCP formulation that is proven to admit only behaviorally consistent equilibrium solutions. The new NCP formulation is much more computationally tractable: because it enforces priority rules at the arc level rather than at the route level, route enumeration is obviated using commonly used column generation techniques. We use the Fischer–Burmeister function to reformulate the "hard" priority conditions as the zero set of a continuously differentiable function, which serves as the upper-level objective in a mathematical program with equilibrium constraints (MPEC) that remains equivalent to the original problem. This MPEC is then solved using two methods: a descent method based on implicit differentiation and a nonlinear programming–based method.

Finally, we demonstrate the practical relevance of the refined model through a series of numerical studies, including the benchmark presented in Nguyen et al. (2001) and a hypothetical schedule-based transit network created using the Sioux Falls network. In particular, we conduct a real-world case study of the morning commute to the University of Hong Kong, where passengers choose between bus and metro services, and metro users may experience substantial queuing at station elevators. The model correctly reproduces FCFS queuing at elevators and reveals departure-early adjustments among bus passengers driven by competition for boarding priority, with the magnitude of such adjustments increasing under higher demand levels.

## 1.3 Organization

The remainder of this paper is organized as follows. Section 2 sets up the problem and introduces existing schedule-based transit assignment models with priority. Section 3 revisits Nguyen et al. (2001)'s framework by presenting an equivalent NCP formulation, establishing the existence of its solutions, and examining the behavioral unrealism exhibited by some of them. Section 4 introduces our refined model, discusses its analytical properties, and develops algorithms for its computation. Section 5 reports numerical experiments that validate the proposed model and algorithms, and Section 6 concludes the paper.

## 2 Nguyen et al.'s Model

### 2.1 Problem setting

**Transit network.** Consider a transit system where $\mathcal{S}$ and $\mathcal{L}$ are the sets of transit stops and transit lines, respectively. For each $l \in \mathcal{L}$, let $s_l^i \in \mathcal{S}$ denote the $i$-th stop it visits ($i = 1, \ldots, n_l$), where $n_l$ is the total number of stops it serves. Each line $l$ has $m_l$ runs, and each run $j$ ($j = 1, \ldots, m_l$) follows a timetable specifying its arrival time $\tau_{l,j,i}^{\text{arr}}$ and departure time $\tau_{l,j,i}^{\text{dep}}$ at stop $s_l^i$ ($i = 1, \ldots, n_l$). All timetable values are defined on a discretized time axis $\mathcal{T}$, with the exceptions $\tau_{l,j,1}^{\text{arr}} = -\infty$ for the arrival at the first stop and $\tau_{l,j,n_l}^{\text{dep}} = \infty$ for the departure at the last stop of each run.

Passengers access and use this system as follows. Let $\mathcal{O}$ and $\mathcal{D}$ denote the sets of origins and destinations. Passengers can walk from their origin to any stop within a certain walking range; for each origin $o \in \mathcal{O}$, the corresponding set of reachable stops is denoted by $\mathcal{S}_o^{\text{walk}} \subseteq \mathcal{S}$, and the walking time between $o$ and $s \in \mathcal{S}_o^{\text{walk}}$ is $t_{o,s}^{\text{walk}}$. They then travel between stops along the transit lines. Finally, they alight at a stop within the walking range of their destination and walk to the destination; for each destination $d \in \mathcal{D}$, the corresponding set of reachable stops is denoted by $\mathcal{S}_d^{\text{walk}} \subseteq \mathcal{S}$, with walking time $t_{d,s}^{\text{walk}}$ from $s \in \mathcal{S}_d^{\text{walk}}$ to $d$. During their trips, passengers may transfer between lines at any stop served by multiple lines. However, for simplicity, we do not allow transfers that require walking to a different stop, although extending the model to include such transfers is straightforward.

As an example, Figure 1 shows a small transit network with four stops, A, B, C, and D, and two lines. Line 1 is a regular service with a single run, and Line 2 is an express service with two runs. The timetable
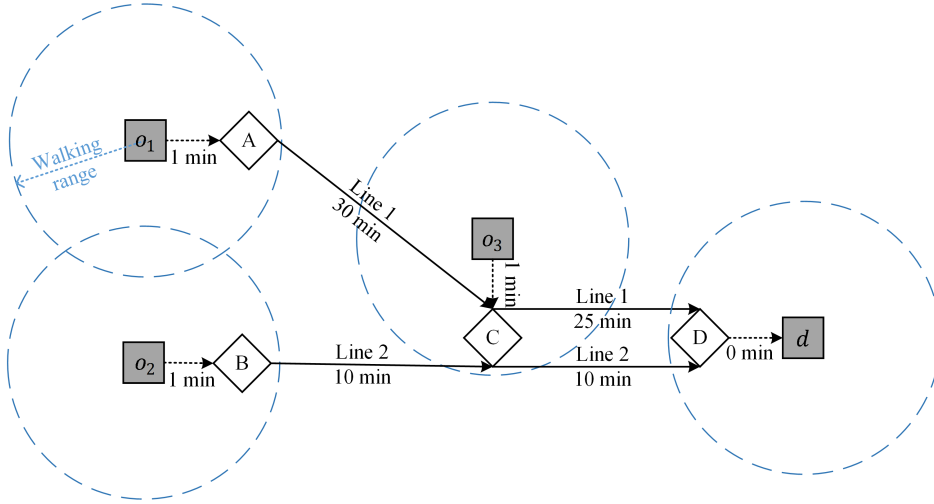


Figure 1: An example transit network.

of all runs is given in Table 1. There are three origins, $o_1$, $o_2$, and $o_3$, connected by walking to stops A, B, and C, respectively, and one destination $d$, connected by walking to stop D. The walking time between each origin and its associated stop is 1 minute, while the walking time between the destination and its associated stop is 0 minutes.

For an OD pair $w = (o, d) \in \mathcal{W} := \mathcal{O} \times \mathcal{D}$, passengers may have different desired arrival times. We group them into a finite set of classes $\mathcal{B}_w$, where each class $b \in \mathcal{B}_w$ has a desired arrival-time window $[\tau_{w,b}^-, \tau_{w,b}^+]$ and a fixed demand $d_{w,b} \in \mathbb{R}_+$. Each passenger chooses both a departure time and a travel path.

4

| Lines | Runs | Stop A | | Stop B | | Stop C | | Stop D | |
|---|---|---|---|---|---|---|---|---|---|
| | | Arrival | Departure | Arrival | Departure | Arrival | Departure | Arrival | Departure |
| 1 | 1 | $-\infty$ | 7:25 | - | - | 7:55 | 7:55 | 8:20 | $\infty$ |
| 2 | 1 | - | - | $-\infty$ | 7:50 | 8:00 | 8:00 | 8:10 | $\infty$ |
| 2 | 2 | - | - | $-\infty$ | 8:10 | 8:20 | 8:20 | 8:30 | $\infty$ |

When all transit vehicles have sufficient capacity so that passengers can always board the first vehicle that arrives after they reach a stop, specifying a passenger's departure time together with their access, transfer, and egress stops and the sequence of transit lines fully determines their arrival time at the destination: any two passengers with the same such choices will experience identical arrival times. However, under oversaturated conditions, some passengers may be denied boarding and must wait for later vehicles. In this case, the above coarse description is no longer sufficient to capture their actual experience. One must instead describe their detailed spatio–temporal events, including the times at which they board, transfer between vehicles, and alight.

**Event–activity graph.** To represent such events, we construct an event–activity graph $\mathcal{H}(\mathcal{E}, \mathcal{A})$, in which each node $E \in \mathcal{E}$ corresponds to an event occurring at a specific time (Yin et al., 2025). These include: (1) $E_{o,t}^{\text{str}}$ (timestamp $t$), which represents the set-out event of passengers with origin $o \in \mathcal{O}$ who begin their trips at time $t \in \mathcal{T}$; (2) $E_{l,j,i}^{\text{dep}}$ (timestamp $\tau_{l,j,i}^{\text{dep}}$), the departure event of run $j = 1, \ldots, m_l$ of line $l \in \mathcal{L}$ from stop $s_l^i$ for $i = 1, \ldots, n_l$; and (3) $E_{l,j,i}^{\text{arr}}$ (timestamp $\tau_{l,j,i}^{\text{arr}}$), the arrival event of run $j = 1, \ldots, m_l$ of line $l$ at stop $s_l^i$ for $i = 1, \ldots, n_l$. Additionally, we create two types of virtual nodes, (4) $E_o^{\text{ogn}}$ for each $o \in \mathcal{O}$, and (5) $E_d^{\text{dst}}$ for each $d \in \mathcal{D}$, to represent, respectively, the starting point of the passenger's trip before entering the transit system and the endpoint after completing the trip. To connect these event nodes, we introduce *activity arcs*, each $A \in \mathcal{A} \subseteq \mathcal{E} \times \mathcal{E}$ representing a feasible spatio–temporal movement undertaken by passengers, as follows.

(1) **Access arcs**: $\mathcal{A}^{\text{access}}$. For each origin $o \in \mathcal{O}$ and starting time $t \in \mathcal{T}$, add an arc from $E_o^{\text{ogn}}$ to $E_{o,t}^{\text{str}}$, representing passengers at origin $o$ choosing to start their trip at time $t$.

(2) **Boarding arcs**: $\mathcal{A}^{\text{boarding}}$. For each $o \in \mathcal{O}$ and $t \in \mathcal{T}$, and for each line $l \in \mathcal{L}$ and stop index $i = 1, \ldots, n_l$, if $s_l^i \in \mathcal{S}_o^{\text{walk}}$, then add an arc from $E_{o,t}^{\text{str}}$ to $E_{l,j,i}^{\text{dep}}$ for each run $j = 1, \ldots, m_l$ satisfying $\tau_{l,j,i}^{\text{dep}} \geq t + t_{o,s_l^i}^{\text{walk}}$. These arcs represent passengers walking to stop $s_l^i$ and attempting to board any run that departs after they arrive at the stop.

(3) **Dwelling arcs**: $\mathcal{A}^{\text{dwelling}}$. For each line $l \in \mathcal{L}$ and run $j = 1, \ldots, m_l$, add an arc from $E_{l,j,i}^{\text{arr}}$ to $E_{l,j,i}^{\text{dep}}$ for each $i = 1, \ldots, n_l$, representing the period during which passengers remain on board while the vehicle dwells at stop $s_l^i$.

(4) **Riding arcs**: $\mathcal{A}^{\text{riding}}$. For each line $l \in \mathcal{L}$ and run $j = 1, \ldots, m_l$, add an arc from $E_{l,j,i}^{\text{dep}}$ to $E_{l,j,i+1}^{\text{arr}}$ for each $i = 1, \ldots, n_l - 1$, representing passengers remaining on board as the vehicle travels from stop $s_l^i$ to stop $s_l^{i+1}$.

(5) **Transfer arcs**: $\mathcal{A}^{\text{transfer}}$. For any two lines $l, l' \in \mathcal{L}$ and stop indices $i = 1, \ldots, n_l$ and $i' = 1, \ldots, n_{l'}$, if $s_l^i = s_{l'}^{i'}$, then for each run $j = 1, \ldots, m_l$ and $j' = 1, \ldots, m_{l'}$, if $\tau_{l,j,i}^{\text{arr}} \leq \tau_{l',j',i'}^{\text{dep}}$, add an arc from $E_{l,j,i}^{\text{arr}}$ to

$E_{l',j',i'}^{\text{dep}}$. This arc represents passengers alighting from line $l$ at stop $s_l^i$ and transferring to a run of line $l'$ that departs later from the same stop.

(6) **Egress arcs**: $\mathcal{A}^{\text{egress}}$. For each destination $d \in \mathcal{D}$, and each line $l \in \mathcal{L}$ and stop index $i = 1, \ldots, n_l$, if $s_l^i \in \mathcal{S}_d^{\text{walk}}$, then add an arc from $E_{l,j,i}^{\text{arr}}$ to $E_d^{\text{dst}}$ for each run $j = 1, \ldots, m_l$. This represents passengers alighting at stop $s_l^i$ and walking to destination $d$.

Figure 2 illustrates the event-activity graph corresponding to the transit network in Figure 1 and the timetable in Table 1. Circular nodes represent vehicle arrival and departure events, square nodes represent origins and destinations, and hexagonal nodes represent passenger starting times. For simplicity, we include one departure time for each of the origins $o_1$ and $o_2$ (7:24 and 7:53, respectively) and two departure times for origin $o_3$ (7:49 and 8:09).
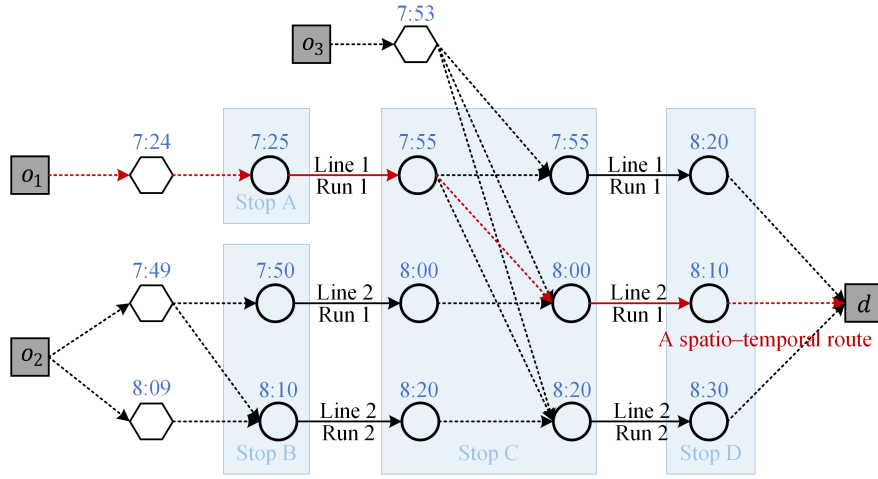


Figure 2: Event–activity graph of the example transit network.

**Spatio-temporal route.** Under the above construction, each path from $E_o^{\text{ogn}}$ to $E_d^{\text{dst}}$ represents a feasible *spatio–temporal route* for passengers of the OD pair $w = (o, d) \in \mathcal{W}$, for example, the red path highlighted in Figure 2. For each such OD pair $w$, we denote by $\mathcal{R}_w$ the set of all routes between $E_o^{\text{ogn}}$ and $E_d^{\text{dst}}$. For each $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and route $r \in \mathcal{R}_w$, let $f_{w,b}^r$ be the number of passengers of class $b$ choosing route $r$, and collect all such variables into the vector $f$. The feasible set of $f$, denoted by $\mathcal{F}$, consists of all $f$ satisfying

$$\sum_{r \in \mathcal{R}_w} f_{w,b}^r = d_{w,b}, \quad \forall w \in \mathcal{W}, \ \forall b \in \mathcal{B}_w, \tag{1}$$

$$f_{w,b}^r \geq 0, \quad \forall w \in \mathcal{W}, \ \forall b \in \mathcal{B}_w, \ \forall r \in \mathcal{R}_w, \tag{2}$$

$$x_A \leq u_A, \quad \forall A \in \mathcal{A}^{\text{riding}}, \tag{3}$$

where $x_A = \sum_{w \in \mathcal{W}} \sum_{b \in \mathcal{B}_w} \sum_{r \in \mathcal{R}_w} f_{w,b}^r \delta_{w,b}^{r,A}$ denotes the total flow assigned to arc $A \in \mathcal{A}$; $\delta_{w,b}^{r,A} = 1$ if arc $A$ belongs to route $r$ and $\delta_{w,b}^{r,A} = 0$ otherwise; $u_A \in \mathbb{R}_+$ is the vehicle capacity of the line associated with $A$.

We do not impose a specific functional form for costs here; the subsequent analysis only requires mild regularity conditions, which we state when needed. Let $c_{w,b}^r(f)$ denote the total travel cost experienced by passengers of OD pair $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, who follow route $r \in \mathcal{R}_w$. This total cost typically includes

components such as travel-time cost, monetary fare, on-board crowding disutility, and penalties for early or late arrival relative to the desired arrival-time window.

## 2.2 User equilibrium with implicit priority

Under the implicit-priority framework, priority rules (e.g., continuance priority for on-board passengers and FCFS for boarding passengers) are represented through the available capacity defined for arcs that terminate at a vehicle departure event; these capacities determine which passengers can board the vehicle and which passengers must wait. We denote by $\mathcal{A}^{\text{priority}} = \mathcal{A}^{\text{boarding}} \cup \mathcal{A}^{\text{dwelling}} \cup \mathcal{A}^{\text{transfer}}$ the set of all such arcs that end at a vehicle departure event, namely the boarding, dwelling, and transfer arcs. For any departure event $E^{\text{dep}}_{l,j,i}$, all arcs ending at this event are equipped with a total order $\prec$ that specifies their loading priorities: the dwelling arc has the highest priority (continuance priority), and the remaining boarding and transfer arcs are ordered in ascending order of their passengers' arrival times at stop $s^i_l$ (reflecting the FCFS rule). Given this priority order, the available capacity on an arc $A \in \mathcal{A}^{\text{priority}}$ is defined as

$$q_A(\boldsymbol{x}) = u_{\text{riding}(A)} - \sum_{A' \in \text{Prior}(A)} x_{A'},$$

Here, $u_{\text{riding}(A)}$ denotes the vehicle capacity of the riding arc $\text{riding}(A)$ whose tail node coincides with the head node of $A$. The set $\text{Prior}(A) = \{A' \in \mathcal{A}^{\text{priority}} : \text{head}(A') = \text{head}(A), A' \preceq A\}$ collects all arcs with priority equal to or higher than that of $A$, where $\text{head}(A)$ is the head node of $A$. Thus, $q_A$ represents the residual vehicle capacity after loading passengers assigned to arc $A$ and all higher-priority arcs. In other words, it determines the maximum number of passengers that can be loaded on arcs of lower priority than $A$. For example, as illustrated in Figure 3, the dwelling arc $(E^{\text{arr}}_{2,1,2}, E^{\text{dep}}_{2,1,2})$ has an available capacity of 3. Hence, at most 3 passengers in total can be loaded via the lower-priority boarding arc $(E^{\text{str}}_{o_3,7:53}, E^{\text{dep}}_{2,1,2})$
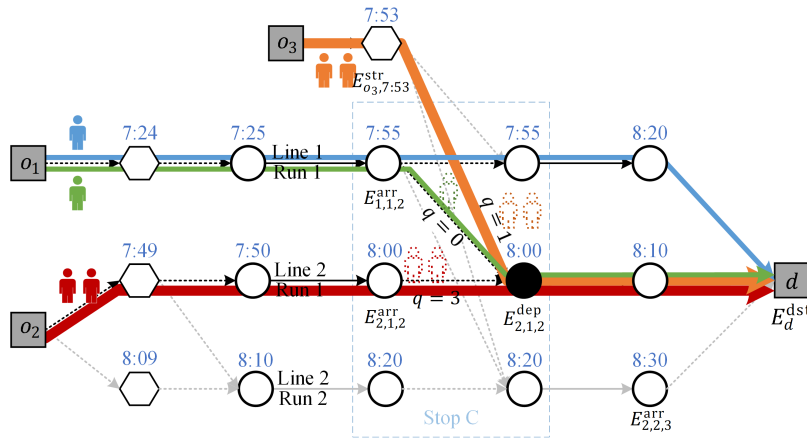


Figure 3: Illustration of the implicit priority and UEIP.

and transfer arc $(E^{\text{arr}}_{1,1,2}, E^{\text{dep}}_{2,1,2})$. Since the two passengers originating from $o_3$ arrive earlier at stop C (at 7:54, given a departure time of 7:53 and a 1-minute walking time), they are loaded first via the boarding arc. After they board, the available capacity associated with the boarding arc is reduced to 1, so at most one passenger can subsequently be loaded via the lower-priority transfer arc.

For each OD pair $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and route $r \in \mathcal{R}_w$, let $Q_{w,b}^r(\boldsymbol{f}) = \min\{q_A(\boldsymbol{x}) : A \in \mathcal{A}_{w,b,r}^{\text{priority}}\}$ be the route available capacity, where $\mathcal{A}_{w,b,r}^{\text{priority}}$ is the set of boarding, dwelling, and transfer arcs belonging to this route. Then we can define the availability of a route and a user equilibrium with implicit priority as follows (Nguyen et al., 2001):

**Definition 1** (Route availability). *Given any feasible flow vector $\boldsymbol{f} \in \mathcal{F}$, for each OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$, a route $r \in \mathcal{R}_w$ is said to be available if every arc in $\mathcal{A}_{w,b,r}^{\text{priority}}$ is available. Equivalently, this condition holds if $Q_{w,b}^r(\boldsymbol{f}) > 0$.*

**Definition 2** (User equilibrium with implicit priority). *A feasible flow vector $\boldsymbol{f} \in \mathcal{F}$ is a user equilibrium with implicit priority (UEIP) if no individual can reduce their travel cost $c_{w,b}^r(\boldsymbol{f})$ by unilaterally switching to another available route. Formally, let $\mathcal{R}_{w,b,r}^{dominate}(\boldsymbol{f}) = \{r' \in \mathcal{R}_w : c_{w,b}^{r'}(\boldsymbol{f}) < c_{w,b}^r(\boldsymbol{f})\}$ be the set of routes that dominate $r$ in terms of travel cost, and let the total available capacity of these dominant routes be $\tilde{Q}_{w,b}^r(\boldsymbol{f}) = \sum_{r' \in \mathcal{R}_{w,b,r}^{dominate}(\boldsymbol{f})} Q_{w,b}^{r'}(\boldsymbol{f})$. Then, a $\boldsymbol{f} \in \mathcal{F}$ is a UEIP solution if*

$$f_{w,b}^r = 0 \quad \text{whenever } \tilde{Q}_{w,b}^r(\boldsymbol{f}) > 0, \quad \forall w \in \mathcal{W}, \forall b \in \mathcal{B}_w, \forall r \in \mathcal{R}_w. \tag{4}$$

## 2.3 Illustrative example

Here we illustrate the notion of UEIP using a simple example. For simplicity, we ignore crowding disutility and fares and assume that route travel cost consists only of constant in-vehicle travel time and late-arrival penalties. Passengers share a desired arrival window of [8:10, 8:20], so a lateness penalty of 10 minutes applies only on the egress arc $(E_{2,2,3}^{\text{arr}}, E_d^{\text{dst}})$, with no penalties on other arcs. The demand for each OD pair is 2, and the vehicle capacity is 5.

Table 2 reports a UEIP solution $\boldsymbol{f}^*$ and a non-UEIP solution $\boldsymbol{f}$, and the corresponding flows are depicted in Figures 3 and 4, respectively. In this example, priority matters at stop C, where it determines which passengers are allowed to board the express service *line 2 run 1*. As discussed earlier, the specification of available capacities on the dwelling, boarding, and transfer arcs ensures that the UEIP solution $\boldsymbol{f}^*$ is fully consistent with the priority rule. By contrast, the non-UEIP solution shown in Figure 4 violates this

Table 2: UEIP and non-UEIP results for the example transit network.

| OD | Routes | Description | Cost | $\boldsymbol{f}^*$ | $Q(\boldsymbol{f}^*)$ | $\boldsymbol{f}$ | $Q(\boldsymbol{f})$ |
|---|---|---|---|---|---|---|---|
| $(o_1, d)$ | $r_1$ | 7:24 - Line 1 Run 1 | 56 | 1 | 3 | 0 | 3 |
| $(o_1, d)$ | $r_2$ | 7:24 - Line 1 Run 1 - Line 2 Run 1 | 46 | 1 | 0 | 2 | 0 |
| $(o_1, d)$ | $r_3$ | 7:24 - Line 1 Run 1 - Line 2 Run 2 | 76 | 0 | 3 | 0 | 3 |
| $(o_2, d)$ | $r_4$ | 7:49 - Line 2 Run 1 | 21 | 2 | 3 | 2 | 3 |
| $(o_2, d)$ | $r_5$ | 7:49 - Line 2 Run 2 | 51 | 0 | 5 | 0 | 5 |
| $(o_2, d)$ | $r_6$ | 8:09 - Line 2 Run 2 | 31 | 0 | 5 | 0 | 5 |
| $(o_3, d)$ | $r_7$ | **7:53 - Line 1 Run 1** | **27** | 0 | 5 | 1 | 4 |
| $(o_3, d)$ | $r_8$ | **7:53 - Line 2 Run 1** | **17** | 2 | 1 | 1 | 2 |
| $(o_3, d)$ | $r_9$ | 7:53 - Line 2 Run 2 | 47 | 0 | 5 | 0 | 5 |

rule. Specifically, one passenger from $o_3$ who arrives earlier at stop C yields a seat on the faster *line 2 run 1* (route $r_8$ in Table 2, with cost 17) to a transferring passenger who arrives later, and instead boards *line 1 run 1* (route $r_7$, with cost 27). Within the implicit-priority framework, this assignment cannot be an
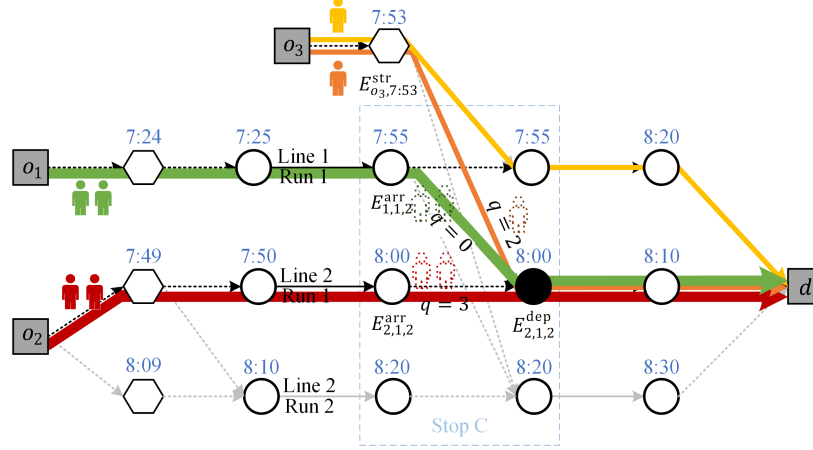
Figure 4: Visualization of a non-UEIP flow for the example transit network.

equilibrium: under $f$, route $r_8$ still has one unit of available capacity, so the passenger currently using $r_7$ has a strict incentive to switch to the available, lower-cost route $r_8$.

This example shows that the implicit-priority framework naturally enforces the boarding-priority rule: only flow patterns that respect the specified priority structure can constitute a user equilibrium. However, the implicit priority framework of Nguyen et al. (2001) lacks a formal mathematical formulation and an accompanying solution algorithm. The subsequent sections will fill these gaps.

# 3 Analysis of Nguyen et al.'s Model

In this section, we will analyze Nguyen et al. (2001)'s model based on a newly proposed mathematical formulation. Throughout the analysis, the following assumptions are imposed.

**Assumption 1.** *For each OD pair $w \in \mathcal{W}$, each class $b \in \mathcal{B}_w$, and each route $r \in \mathcal{R}_w$, the travel cost $c_{w,b}^r(f)$ is positive for all feasible $f \in \mathcal{F}$.*

**Assumption 2.** *Given any feasible flow vector $f \in \mathcal{F}$, for each OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$, there exists at least one route $r \in \mathcal{R}_w$ such that $Q_{w,b}^r(f) > 0$, i.e., at least one route is available.*

Assumption 2 rules out degenerate situations in which a feasible assignment $f \in \mathcal{F}$ exhausts the capacities of *all* possible transit routes for an OD pair, including those that involve delays, which is a reasonable condition under normal operating circumstances. The remainder of this section is organized as follows. Section 3.1 presents an equivalent reformulation of Nguyen et al. (2001)'s model. Section 3.2 examines the existence and uniqueness of the solution. Finally, a numerical example illustrates that the model may admit behaviorally unrealistic solutions (Section 3.3).

## 3.1 A new UEIP formulation

To facilitate model formulation, we first present the alternative UEIP conditions that are equivalent to Condition (4) and more explicitly characterize the relationship between route costs and route flows.

9

**Proposition 1.** *A feasible solution $f^* \in \mathcal{F}$ satisfies the Condition (4) if and only if there exists $\mu_{w,b} \in \mathbb{R}$ for each $w \in \mathcal{W}$ and $b \in \mathcal{B}_w$ such that*

$$c_{w,b}^r(f^*) \geq \mu_{w,b} \quad \text{if } Q_{w,b}^r(f^*) > 0, \quad \forall r \in \mathcal{R}_w, \tag{5a}$$

$$\mu_{w,b} \geq c_{w,b}^r(f^*) \quad \text{if } f_{w,b}^{r*} > 0, \quad \forall r \in \mathcal{R}_w. \tag{5b}$$

*Proof.* First, suppose that $f^* \in \mathcal{F}$ satisfies Equation (5). For each OD $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and route $r \in \mathcal{R}_w$, if $\tilde{Q}_{w,b}^r(f^*) > 0$, then there must be at least one route $r' \in \mathcal{R}_{w,b,r}^{\text{dominate}}(f^*)$ satisfying $c_{w,b}^r(f^*) > c_{w,b}^{r'}(f^*)$ and $Q_{w,b}^{r'}(f^*) > 0$. By Equation (5a), we have $\mu_{w,b} \leq c_{w,b}^{r'}(f^*) < c_{w,b}^r(f^*)$. Then, according to (5b), we must have $f_{w,b}^{r*} = 0$. Hence, Condition (4) holds.

Conversely, suppose that $f^* \in \mathcal{F}$ satisfies Equation (4). For each OD $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$, given that $\bar{r}$ is the most costly of all used routes, namely $c_{w,b}^{\bar{r}}(f^*) = \max\{c_{w,b}^r(f^*) : f_{w,b}^{r*} > 0, r \in \mathcal{R}_{w,b}\}$, and $\mu_{w,b} = c_{w,b}^{\bar{r}}(f^*)$, Equation (5b) holds. Then, we prove Equation (5a) by contradiction. For each $w \in \mathcal{W}$ and $b \in \mathcal{B}_w$, suppose that there exists an available route $r'$ (i.e., $Q_{w,b}^{r'}(f^*) > 0$) with $c_{w,b}^{r'}(f^*) < \mu_{w,b} = c_{w,b}^{\bar{r}}(f^*)$. Then, we must have $r' \in \mathcal{R}_{w,b,\bar{r}}^{\text{dominate}}(f^*)$, and thus $\tilde{Q}_{w,b}^{\bar{r}}(f^*) = \sum_{r'' \in \mathcal{R}_{w,b,\bar{r}}^{\text{dominate}}(f^*)} Q_{w,b}^{r''}(f^*) \geq Q_{w,b}^{r'}(f^*) > 0$. According to Condition (4), this conflicts with that $\bar{r}$ is a used route. By contradiction, Equation (5a) holds. $\square$

This proposition implies that, under UEIP, the cost of any used route does not exceed the cost of other available routes, with $\mu = (\mu_{w,b})_{w \in \mathcal{W}, b \in \mathcal{B}_w}$ serving as the threshold separating these two categories of routes. To close the potential gap between the cost of a used route and this threshold, we introduce an additional cost variable $V_{w,b}^r$ for each OD $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and route $r \in \mathcal{R}_w$, and collect all such variables into the vector $V$. Specifically, for each used route we set $V_{w,b}^r = \mu_{w,b} - c_{w,b}^r(f)$, while for available routes we let $V_{w,b}^r = 0$. Based on this construction, we obtain the following NCP model:

$$0 \leq f \perp c(f) + V - \Lambda\mu \geq 0, \tag{6a}$$

$$0 \leq \mu \perp \Lambda^T f - d \geq 0, \tag{6b}$$

$$0 \leq V \perp Q(f) \geq 0, \tag{6c}$$

where $c(f) = (c_{w,b}^r(f))_{w \in \mathcal{W}, b \in \mathcal{B}_w, r \in \mathcal{R}_w}$ and $Q(f) = (Q_{w,b}^r(f))_{w \in \mathcal{W}, b \in \mathcal{B}_w, r \in \mathcal{R}_w}$ denote the vectors of travel costs and available capacities, respectively; $\Lambda = [\lambda_{r,b}]_{|n_f| \times |\mathcal{B}|}$ denotes the route-class incident matrix, with $\lambda_{r,b} = 1$ if route $r$ belongs to class $b$, and $\lambda_{r,b} = 0$ otherwise; $n_f = \sum_{w \in \mathcal{W}} |\mathcal{B}_w| \times |\mathcal{R}_w|$ is the dimension of route flow vector $f$. The subsequent theorem shows that this NCP formulation is equivalent to the UEIP conditions.

**Theorem 1.** *Under Assumption 1, the NCP (6) is equivalent to the UEIP conditions (1), (2), (3), and (5).*

*Proof.* **Necessity**: Suppose that the solution $(f^*, \mu^*)$ satisfies the UEIP conditions (1), (2), (3), and (5). Let $x^* = \Delta^T f^*$, where $\Delta = [\delta_{w,b}^{r,A}]_{|n_f| \times |\mathcal{A}|}$ is the route-arc incident matrix. We will then show that the complementarity condition (6) holds.

**Condition** (6b): Condition (1) shows that $\Lambda^T f^* - d = 0$ and $\mu^* \perp \Lambda^T f^* - d$. Equation (5b) implies $\mu$ is a positive vector. Therefore, Condition (6b) holds.

**Condition** (6c): For any arc $A \in \mathcal{A}^{\text{priority}}$, we have $q_A(x^*) = u_{\text{riding}(A)} - \sum_{A' \in \text{Prior}(A)} x_{A'}^* \geq u_{\text{riding}(A)} - x_{\text{riding}(A)}^*$. By the capacity constraint (3), we have $q_A(x^*) \geq 0$. Therefore, $Q_{w,b}^r(f^*) = \min\{q_A(x^*) : A \in \mathcal{A}_{w,b,r}^{\text{priority}}\} \geq 0$ for all $w \in \mathcal{W}, b \in \mathcal{B}_w$, and $r \in \mathcal{R}_w$.

For each OD $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and route $r \in \mathcal{R}_w$, letting

$$
V_{w,b}^r = \begin{cases}
0, & \text{if } Q_{w,b}^r(\boldsymbol{f}^*) > 0 \\
\mu_{w,b}^* - c_{w,b}^r(\boldsymbol{f}^*), & \text{if } Q_{w,b}^r(\boldsymbol{f}^*) = 0 \text{ and } f_{w,b}^{r*} > 0 \\
\max\{\mu_{w,b}^* - c_{w,b}^r(\boldsymbol{f}^*), 0\}, & \text{if } Q_{w,b}^r(\boldsymbol{f}^*) = 0 \text{ and } f_{w,b}^{r*} = 0
\end{cases} \tag{7}
$$

we then have $\boldsymbol{V} \perp \boldsymbol{Q}(\boldsymbol{f}^*)$. Meanwhile, considering that $\mu_{w,b}^* - c_{w,b}^r(\boldsymbol{f}^*) \geq 0$ when $f_{w,b}^{r*} > 0$ by Equation (5b), we must have $\boldsymbol{V} \geq 0$. Therefore, Condition (6c) holds.

**Condition** (6a): Following Equation(7), we also have $\boldsymbol{c}(\boldsymbol{f}) + \boldsymbol{V} - \Lambda\boldsymbol{\mu} \geq 0$ because (I) $V_{w,b}^r \geq \mu_{w,b}^* - c_{w,b}^r(\boldsymbol{f}^*)$, when $Q_{w,b}^r(\boldsymbol{f}^*) = 0$; (II) $V_{w,b}^r = 0$ and $c_{w,b}^r(\boldsymbol{f}^*) \geq \mu_{w,b}^*$ when $Q_{w,b}^r(\boldsymbol{f}^*) > 0$ by Equation (5a).

The orthogonality between $\boldsymbol{f}$ and $\boldsymbol{c}(\boldsymbol{f}) + \boldsymbol{V} - \Lambda\boldsymbol{\mu}$ must stands when $f_{w,b}^{r*} = 0$. If $f_{w,b}^{r*} > 0$, there are two cases. Case (I): If $Q_{w,b}^r(\boldsymbol{f}^*) > 0$, Equations (5a) and (5b) show $c_{w,b}^r(\boldsymbol{f}^*) \geq \mu_{w,b}^* \geq c_{w,b}^r(\boldsymbol{f}^*)$, so $c_{w,b}^r(\boldsymbol{f}^*) = \mu_{w,b}^*$. Then, we have $c_{w,b}^r(\boldsymbol{f}^*) + V_{w,b}^r - \mu_{w,b}^* = 0$ as $V_{w,b}^r = 0$. Case (II): If $Q_{w,b}^r(\boldsymbol{f}^*) = 0$, we also have $c_{w,b}^r(\boldsymbol{f}^*) + V_{w,b}^r - \mu_{w,b}^* = 0$ because $V_{w,b}^r$ is set to $\mu_{w,b}^* - c_{w,b}^r(\boldsymbol{f}^*)$ in Equation (7). Therefore, orthogonality is satisfied universally, and Condition (6a) holds.

**Sufficiency**: Suppose that the solution $(\boldsymbol{f}^*, \boldsymbol{\mu}^*, \boldsymbol{V}^*)$ satisfies the complementarity condition (6). Letting $\boldsymbol{x}^* = \Delta^T \boldsymbol{f}^*$, we then show that the UEIP conditions (1), (2), (3), and (5) hold.

**Condition** (2): The complementarity condition (6a) directly suggests that $f_{w,b}^{r*} \geq 0$ for all $w \in \mathcal{W}$, $b \in \mathcal{B}_w$, and $r \in \mathcal{R}_w$.

**Condition** (3): The non-negativity of $\boldsymbol{Q}(\boldsymbol{f}^*)$ suggests that $q_A(\boldsymbol{x}^*) \geq 0$ for all $A \in \mathcal{A}^{\text{priority}}$. For any riding arc $\bar{A} \in \mathcal{A}^{\text{riding}}$, let $A$ be the incoming arc with the lowest loading priority. We then have $0 \leq q_A(\boldsymbol{x}^*) = u_{\bar{A}} - \sum_{A' \in \text{Prior}(A)} x_{A'}^* = u_{\bar{A}} - x_{\bar{A}}^*$. Thus, the capacity constraint (3) is satisfied.

**Condition** (1): Suppose that for some OD $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$,

$$
\sum_{r \in \mathcal{R}_w} f_{w,b}^{r*} > d_{w,b}. \tag{8}
$$

By the complementarity condition (6a) and (6b), we must have $\mu_{w,b}^* = 0$ and for each route $r \in \mathcal{R}_w$,

$$
f_{w,b}^{r*}(c_{w,b}^r(\boldsymbol{f}^*) + V_{w,b}^{r*} - \mu_{w,b}^*) = 0 \Rightarrow f_{w,b}^{r*}(c_{w,b}^r(\boldsymbol{f}^*) + V_{w,b}^{r*}) = 0.
$$

In addition, since $c_{w,b}^r(\boldsymbol{f}^*)$ is positive and $V_{w,b}^{r*}$ is non-negative, we must have $f_{w,b}^{r*} = 0$ for all $w \in \mathcal{W}$, $b \in \mathcal{B}_w$, and $r \in \mathcal{R}_w$.

On the other hand, Equation (8) implies that $\sum_{r \in \mathcal{R}_w} f_{w,b}^{r*} > 0$ and there must exist a route with positive flow ($f_{w,b}^{r*} > 0$) in this OD and class, which is a contradiction. Thus, the demand constraint (1) holds.

**Condition** (5): For each OD $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and route $r \in \mathcal{R}_w$, if $Q_{w,b}^r(\boldsymbol{f}^*) > 0$, we have $V_{w,b}^{r*} = 0$ by the complementarity condition (6c). Since $c_{w,b}^r(\boldsymbol{f}^*) + V_{w,b}^{r*} - \mu_{w,b}^*$ is non-negative, we get $c_{w,b}^r(\boldsymbol{f}^*) \geq \mu_{w,b}^*$, and hence Condition (5a) holds.

On the other hand, if $f_{w,b}^{r*} > 0$, by the complementarity condition (6a), we have $c_{w,b}^r(\boldsymbol{f}^*) + V_{w,b}^{r*} - \mu_{w,b}^* = 0$. Therefore, $c_{w,b}^r(\boldsymbol{f}^*) \leq \mu_{w,b}^*$ since $V_{w,b}^{r*}$ is non-negative, and thus Condition (5b) holds. $\qquad\square$

## 3.2 Existence and uniqueness

We first establish the existence of a solution to the NCP (6), and then discuss its uniqueness.

**Proposition 2.** *Under Assumptions 1 and 2, NCP (6) has a solution.*

*Proof.* To prove the existence of a solution to NCP (6), we begin by recalling that a solution is guaranteed when the mapping is continuous and the feasible set is compact. In NCP (6), although the function

$$H(f, \mu, V) = \begin{pmatrix} c(f) + V - \Lambda\mu \\ \Lambda^T f - d \\ Q(f) \end{pmatrix}$$

is continuous, the feasible set $\Omega = \{f \geq 0, \mu \geq 0, V \geq 0\}$ is unbounded and therefore not compact. Our proof proceeds by introducing upper-bound constraints for each variable, thereby constructing a new NCP whose feasible set is compact. We then show that none of these added upper bounds are binding at the solution of the new problem. Consequently, any solution to the modified NCP also satisfies the original NCP (6). Choose two scalars $e_1$ and $e_2$ such that

$$e_1 > \max\{d_{w,b} : w \in \mathcal{W}, b \in \mathcal{B}_w\} \text{ and } e_2 > \max\{c_{w,b}^r : w \in \mathcal{W}, b \in \mathcal{B}_w, r \in \mathcal{R}_w\}.$$

Define $\Omega' = \Omega \cap \{f \leq e_1\mathbf{1}, \mu \leq e_2\mathbf{1}, V \leq e_2\mathbf{1}\}$. Since $\Omega'$ is compact, the following NCP must possess a solution:

$$0 \leq f \perp c(f) + V - \Lambda\mu + \kappa \geq 0, \tag{9a}$$

$$0 \leq \mu \perp \Lambda^T f - d + \rho \geq 0, \tag{9b}$$

$$0 \leq V \perp Q(f) + v \geq 0, \tag{9c}$$

$$0 \leq \kappa \perp e_1\mathbf{1} - f \geq 0, \tag{9d}$$

$$0 \leq \rho \perp e_2\mathbf{1} - \mu \geq 0, \tag{9e}$$

$$0 \leq v \perp e_2\mathbf{1} - V \geq 0. \tag{9f}$$

Let $(f^*, \mu^*, V^*)$ be this solution, if we show that $\kappa$, $\rho$, and $v$ are all equal to zero, then $(f^*, \mu^*, V^*)$ must also be a solution of NCP (6).

**For $\kappa$:** Suppose $\kappa_{w,b}^r > 0$ for some $w \in \mathcal{W}, b \in \mathcal{B}_w, r \in \mathcal{R}_w$. By the complementarity condition (9d), we have $f_{w,b}^{r*} = e_1 > d_{w,b}$, so $\sum_{r \in \mathcal{R}_w} f_{w,b}^{r*} - d_{w,b} + \rho_{w,b} > 0$. Then, according to Condition (9b), we have $\mu_{w,b}^* = 0$. Moreover, since $f_{w,b}^{r*} = e_1 > 0$, we have

$$0 = c_{w,b}^r(f^*) + V_{w,b}^{r*} - \mu_{w,b}^* + \kappa_{w,b}^r = c_{w,b}^r(f^*) + V_{w,b}^{r*} + \kappa_{w,b}^r > 0,$$

because $c_{w,b}^r(f^*)$ and $\kappa_{w,b}^r$ are positive, and $V_{w,b}^{r*}$ is non-negative. This contradiction yields $\kappa = \mathbf{0}$.

**For $\rho$:** For each OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$, since there exists a route $r \in \mathcal{R}_w$ with $Q_{w,b}^r(f^*) > 0$ by assumption, we have $Q_{w,b}^r(f^*) + v_{w,b}^r > 0$. By Condition (9c), $V_{w,b}^{r*} = 0$. Therefore,

$$0 \leq c_{w,b}^r(f^*) + V_{w,b}^{r*} - \mu_{w,b}^* + \kappa_{w,b}^r = c_{w,b}^r(f^*) - \mu_{w,b}^*$$

$$\Rightarrow \mu_{w,b}^* \leq c_{w,b}^r(f^*) < e_2. \tag{10}$$

This means that $e_2 - \mu_{w,b}^* > 0$ for all $w \in \mathcal{W}$ and $b \in \mathcal{B}_w$, and $\rho = \mathbf{0}$ by Condition (9e).

**For $v$:** Supposing that $v_{w,b}^r > 0$ for some $w \in \mathcal{W}, b \in \mathcal{B}_w, r \in \mathcal{R}_w$, we then have $V_{w,b}^{r*} = e_2 > 0$. By Condition (9c), we have $Q_{w,b}^r(f^*) + v_{w,b}^r = 0$, hence $Q_{w,b}^r(f^*) = -v_{w,b}^r < 0$. By the definition of $Q_{w,b}^r(f^*)$, there must exist an arc with negative available capacity in route $r$, namely one $A \in \mathcal{A}_{w,b,r}^{\text{priority}}$ such that $0 > q_A(x^*) = u_{\text{riding}(A)} - \sum_{A' \in \text{Prior}(A)} x_{A'}^*$.

Let $A'$ be the arc in $\text{Prior}(A)$ that carries positive flow and has the lowest loading priority. Then, $A'$ has the same available capacity as $A$, that is $q_{A'}(x^*) = q_A(x^*) < 0$. Moreover, there is at least one used

12

route $r'$ that passes through this arc. Without loss of generality, let $r'$ belong to OD pair $w' \in \mathcal{W}$ and class $b' \in \mathcal{B}_w$. By the complementarity condition (9a) and the fact that $f^{r'*}_{w',b'} > 0$, we have

$$c^{r'}_{w',b'}(f^*) + V^{r'*}_{w',b'} - \mu^*_{w',b'} + \kappa^{r'}_{w',b'} = 0 \Rightarrow V^{r'*}_{w',b'} = \mu^*_{w',b'} - c^{r'}_{w',b'}(f^*).$$

Since $\mu^*_{w',b'} < e_2$ (Equation (10)) and $c^{r'}_{w',b'}(f^*)$ is positive, we get $V^{r'*}_{w',b'} < e_2$. This means $v^{r'}_{w',b'} = 0$ by Condition (9f). On the another hand, $Q^{r'}_{w',b'}(f^*) = \min\{q_{A''}(x^*) : A'' \in \mathcal{A}^{\text{priority}}_{w,b,r}\} \leq q_{A'}(x^*) < 0$. This implies $Q^{r'}_{w',b'}(f^*) + v^{r'}_{w',b'} < 0$, which contradicts Condition (9c). Therefore, $\mathbf{v} = \mathbf{0}$ holds. $\square$

Having established the existence of a solution to NCP (6), we will then turn to its uniqueness. Since the mapping $H(f, \mu, V)$ is non-monotone, the solution set of NCP (6) — equivalently, the UEIP solution set — is generally not unique. We next discuss the consequence of this non-uniqueness.

### 3.3 Unexpected consequence of non-uniqueness

Many traffic assignment problems do not admit unique solutions. When this occurs, it is common practice to treat any solution returned by a valid algorithm as equally admissible. However, as we show below, some UEIP solutions can exhibit clear violations of realistic passenger behavior. This undermines the applicability of the NCP formulation (6), as it may admit solutions that are behaviorally implausible or operationally nonsensical.

Consider a toy network shown in Figure 5(a), which consists of two lines (Line 1 and Line 2). Specifi-


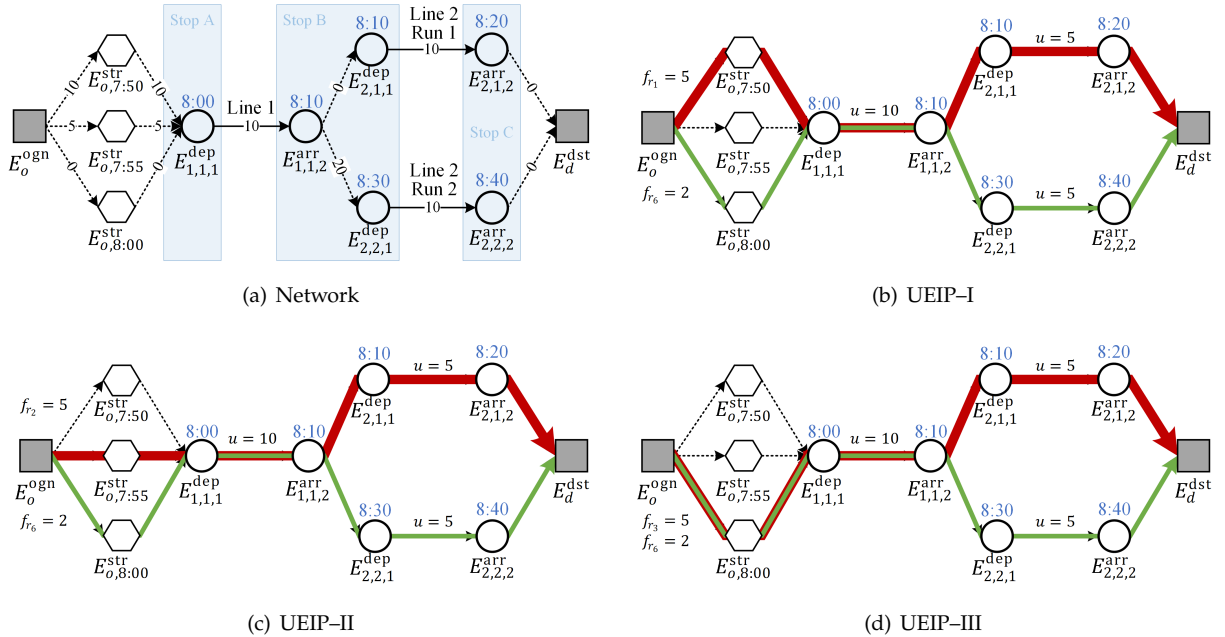
(a) Network

(b) UEIP–I

(c) UEIP–II

(d) UEIP–III

Figure 5: An illustrative example of unreasonable UEIP states.

cally, Line 1 has only one run, while Line 2 includes two runs (Run 1 and Run 2). Congestion effects are not considered, and each arc is labeled with its cost. Seven passengers depart from origin $o$ to destination $d$, with three starting time options: 7:50, 7:55, and 8:00. All passengers first take the high-capacity Line 1 (capacity = 10) from Stop A to Stop B. Then, a subset of them can directly transfer to the low-capacity Run

1 of Line 2 (capacity = 5) to reach the destination without waiting, while the remaining two passengers must stay at the platform and wait for Run 2 of Line 2. This setting yields three distinct UEIP solutions, namely UEIP–I, UEIP–II, and UEIP–III as shown in Figure 5(b)–5(d) and Table 3, which differ only in the starting times chosen by the passengers taking Run 1. In UEIP–I and UEIP–II, these passengers begin their trips at 7:50 and 7:55, respectively, to catch Line 1 leaving at 8:00. However, such outcomes are unrealistic. Since Line 1 is unsaturated — meaning that passengers can board regardless of their arrival order at Stop A — the optimal behavior is simply to depart at 8:00, as in UEIP–III. Departing earlier yields no advantage and only increases waiting time and the early-departure penalty.

Table 3: Summary of the three UEIP solutions.

| Routes | Trajectories | Cost | UEIP–I | | | UEIP–II | | | UEIP–III | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | $f$ | $Q$ | $V$ | $f$ | $Q$ | $V$ | $f$ | $Q$ | $V$ |
| $r_1$ | 7:50, Run 1 | 40 | 5 | 0 | 0 | 0 | 0 | $\geq 0$ | 0 | 0 | $\geq 0$ |
| $r_2$ | 7:55, Run 1 | 30 | 0 | 0 | $\geq 10$ | 5 | 0 | 10 | 0 | 0 | $\geq 10$ |
| $r_3$ | 8:00, Run 1 | 20 | 0 | 0 | $\geq 20$ | 0 | 0 | $\geq 20$ | 5 | 0 | 20 |
| $r_4$ | 7:50, Run 2 | 60 | 0 | 3 | 0 | 0 | 3 | 0 | 0 | 3 | 0 |
| $r_5$ | 7:55, Run 2 | 50 | 0 | 3 | 0 | 0 | 3 | 0 | 0 | 3 | 0 |
| $r_6$ | 8:00, Run 2 | 40 | 2 | 3 | 0 | 2 | 3 | 0 | 2 | 3 | 0 |

To illustrate the reason behind the existence of the unrealistic solutions, we take UEIP–II as an example. Passengers assigned to Run 1 (i.e., route $r_2$) cannot shift their starting time to 8:00 (i.e., move to route $r_3$) because the arc $(E_{1,1,2}^{\mathrm{arr}}, E_{2,1,1}^{\mathrm{dep}})$ is treated as unavailable, thereby rendering route $r_3$ unavailable. In reality, however, we argue that route $r_3$ should be regarded as "available" for passengers on $r_2$, since transferring flow from $r_2$ to $r_3$ merely changes their starting time and does not introduce additional flow to the bottleneck arc $(E_{1,1,2}^{\mathrm{arr}}, E_{2,1,1}^{\mathrm{dep}})$. More generally, even if route $r$ contains some unavailable arcs, when another route $r'$ also traverses these arcs, reallocating flow from $r'$ to $r$ does not increase the load on those unavailable arcs. Such a reallocation should therefore be admissible, and route $r$ should be considered available relative to route $r'$.

In summary, this example highlights that the definition of route availability deviates from actual passenger behavior, which leads to unrealistic solutions in the UEIP solution set. In the next section, we refine the definition of available routes and introduce a new equilibrium condition together with its corresponding NCP formulation, which eliminates such unrealistic outcomes. We also develop an algorithm tailored to solve the proposed model.

# 4  A Behaviorally Compliant UEIP Formulation

We now proceed to a refined UEIP condition, which not only rules out the behaviorally unrealistic solutions but also enables more tractable algorithms. In what follows, Section 4.1 presents the refined UEIP conditions and the corresponding NCP formulation, and Section 4.2 develops the solution algorithms.

## 4.1  UEIP revisited

The preceding example highlights that route availability is not absolute but relative to the reference route, as reallocating flow may remain admissible if no additional load is imposed on unavailable arcs. Accord-

ingly, we define the available capacity of route $r'$ with respect to route $r$ as

$$Q_{w,b}^{r',r}(\boldsymbol{f}) = \min\{q_A(\boldsymbol{x}) : A \in \mathcal{A}_{w,b,r'}^{\text{priority}} \setminus \mathcal{A}_{w,b,r}^{\text{priority}}\}, \quad \forall w \in \mathcal{W}, \ \forall b \in \mathcal{B}_w, \ \forall r',r \in \mathcal{R}_w \text{ and } r' \neq r.$$

Note that since $\mathcal{A}_{w,b,r'}^{\text{priority}} \setminus \mathcal{A}_{w,b,r}^{\text{priority}}$ is nonempty whenever $r' \neq r$, the quantity $Q_{w,b}^{r',r}(\boldsymbol{f})$ is well-defined. Based on this, we revise the definition of route availability as follows:

**Definition 3** (Relative route availability). *Given any feasible flow vector $\boldsymbol{f} \in \mathcal{F}$, for each OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$, a route $r \in \mathcal{R}_w$ is said to be available with respect to another route $r' \in \mathcal{R}_w$ if every arc on $r'$ that is not contained in $r$ is available. Equivalently, this condition holds if $Q_{w,b}^{r',r}(\boldsymbol{f}) > 0$.*

With this definition, we can tighten the equilibrium condition to rule out behavioral inconsistency.

**Definition 4** (Refined user equilibrium with implicit priority). *A feasible flow vector $\boldsymbol{f} \in \mathcal{F}$ is a refined user equilibrium with implicit priority (R-UEIP) if for each OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$ no individual on route $r \in \mathcal{R}_w$ can reduce their travel cost $c_{w,b}^r(\boldsymbol{f})$ by unilaterally switching to another available route that is available with respect to $r$. Formally, let the total available capacity of the dominant routes with respect to $r$ be $\hat{Q}_{w,b}^r(\boldsymbol{f}) = \sum_{r' \in \mathcal{R}_{w,b,r}^{\text{dominate}}(\boldsymbol{f})} Q_{w,b}^{r',r}(\boldsymbol{f})$. Then, $\boldsymbol{f} \in \mathcal{F}$ is an R-UEIP solution if*

$$f_{w,b}^r = 0 \quad \text{whenever } \hat{Q}_{w,b}^r(\boldsymbol{f}) > 0, \quad \forall w \in \mathcal{W}, \ \forall b \in \mathcal{B}_w, \ \forall r \in \mathcal{R}_w. \tag{11}$$

It is straightforward to verify that, in the toy example of the previous section, the unrealistic equilibria UEIP–I and UEIP–II do not satisfy the refined equilibrium condition, whereas UEIP–III does. The relationship between the two types of equilibrium can be summarized as follows:

**Proposition 3.** *If $\boldsymbol{f}^*$ is an R-UEIP, then it is also a UEIP, but the converse does not hold.*

*Proof.* Suppose that $\boldsymbol{f}^*$ satisfies the R-UEIP condition. Let $\boldsymbol{x}^* = \Delta^T \boldsymbol{f}$. For any OD pair $w \in \mathcal{W}$, any class $b \in \mathcal{B}_w$, and any two different routes $r, r' \in \mathcal{R}_w$, it follows that $Q_{w,b}^{r',r}(\boldsymbol{f}^*) = \min\{q_A(\boldsymbol{x}^*) : A \in \mathcal{A}_{w,b,r'}^{\text{priority}} \setminus \mathcal{A}_{w,b,r}^{\text{priority}}\} \geq \min\{q_A(\boldsymbol{x}^*) : A \in \mathcal{A}_{w,b,r'}^{\text{priority}}\} = Q_{w,b}^{r'}(\boldsymbol{f}^*)$. Consequently,

$$\hat{Q}_{w,b}^r(\boldsymbol{f}^*) = \sum_{r' \in \mathcal{R}_{w,b,r}^{\text{dominate}}(\boldsymbol{f}^*)} Q_{w,b}^{r',r}(\boldsymbol{f}^*) \geq \sum_{r' \in \mathcal{R}_{w,b,r}^{\text{dominate}}(\boldsymbol{f}^*)} Q_{w,b}^{r'}(\boldsymbol{f}^*) = \tilde{Q}_{w,b}^r(\boldsymbol{f}^*).$$

Therefore, if $\tilde{Q}_{w,b}^r(\boldsymbol{f}^*) > 0$, then $\hat{Q}_{w,b}^r(\boldsymbol{f}^*) > 0$. According to the R-UEIP condition (11), this implies $f_{w,b}^{r*} = 0$, thus proving the forward direction. The counterexample provided in Section 3.2 demonstrates that the converse does not hold. $\qquad\square$

This result implies that the R-UEIP admits a smaller solution set, excluding certain unrealistic outcomes permitted by the UEIP. Next, we formulate an NCP model for computing the refined equilibrium. For every arc $A \in \mathcal{A}^{\text{priority}}$, we introduce a variable $v_A$, collected into the vector $\boldsymbol{v}$, which is orthogonal to the arc's available capacity; that is, $v_A \perp q_A(\boldsymbol{x})$. Let $\bar{\Delta} = [\delta_{w,b}^{r,A}]_{|n_f| \times |\mathcal{A}^{\text{priority}}|}$ be the route-arc incident matrix. We propose the following NCP model:

$$0 \leq \boldsymbol{f} \perp \boldsymbol{c}(\boldsymbol{f}) + \bar{\Delta}\boldsymbol{v} - \Lambda\boldsymbol{\mu} \geq 0, \tag{12a}$$

$$0 \leq \boldsymbol{\mu} \perp \Lambda^T \boldsymbol{f} - \boldsymbol{d} \geq 0, \tag{12b}$$

$$0 \leq \boldsymbol{v} \perp \boldsymbol{q}(\boldsymbol{x}) \geq 0. \tag{12c}$$

15

**Proposition 4.** *Under Assumption 1, if $f^*$ is a solution of NCP (12), $f^*$ satisfies the R-UEIP conditions (1), (2), (3), and (11).*

*Proof.* The proof that $f^*$ satisfies the feasible conditions (1), (2), and (3) follows the same reasoning as in Theorem 1, with $V$ replaced by $\bar{\Delta}v$, and is thus omitted here. Next, we prove that $f^*$ also satisfies condition (11).

For any OD pair $w \in \mathcal{W}$ and any class $b \in \mathcal{B}_w$, suppose there exists a route $r \in \mathcal{R}_w$ such that $\hat{Q}^r_{w,b}(f^*) > 0$. Then, there must exist another route $r' \in \mathcal{R}_w$ satisfying $c^{r'}_{w,b}(f^*) < c^r_{w,b}(f^*)$ and $Q^{r',r}_{w,b}(f^*) > 0$. This implies $q_A > 0$ for all $A \in \mathcal{A}^{\text{priority}}_{w,b,r'} \setminus \mathcal{A}^{\text{priority}}_{w,b,r}$. According to the complementarity condition (12c), we have $v_A = 0$ for all such $A$. Substituting these $v_A$ values into (12a) yields:

$$0 \le c^{r'}_{w,b}(f^*) + \sum_{A \in \mathcal{A}^{\text{priority}}_{w,b,r'}} v_A - \mu_{w,b} = c^{r'}_{w,b}(f^*) + \sum_{A \in \mathcal{A}^{\text{priority}}_{w,b,r'} \cap \mathcal{A}^{\text{priority}}_{w,b,r}} v_A - \mu_{w,b}.$$

Moreover, since $c^{r'}_{w,b}(f^*) < c^r_{w,b}(f^*)$ and $\sum_{A \in \mathcal{A}^{\text{priority}}_{w,b,r'} \cap \mathcal{A}^{\text{priority}}_{w,b,r}} v_A \le \sum_{A \in \mathcal{A}^{\text{priority}}_{w,b,r}} v_A$, it follows that

$$0 < c^r_{w,b}(f^*) + \sum_{A \in \mathcal{A}^{\text{priority}}_{w,b,r}} v_A - \mu_{w,b}.$$

Finally, by the complementarity condition in (12a), we obtain $f^{r*}_{w,b} = 0$, which proves that $f^*$ satisfies (11). ☐

Under the same assumptions as NCP (6), the existence of a solution for NCP (12) can also be guaranteed. The result regarding the existence of the solution is described as follows:

**Proposition 5.** *Under Assumptions 1 and 2, NCP (12) has a solution.*

*Proof.* See Appendix B. ☐

**Remark 1.** *The introduced variable $v$ can be interpreted as the anxiety cost incurred by passengers when they traverse activities (arcs) for which, due to their lower loading priority, boarding cannot be guaranteed. As illustrated in Figure 6, the available capacity of the dwelling arc ($E^{arr}_{2,1,2}$, $E^{dep}_{2,1,2}$) and the boarding arc ($E^{str}_{o_3,7:53}$, $E^{dep}_{2,1,2}$) are positive,*
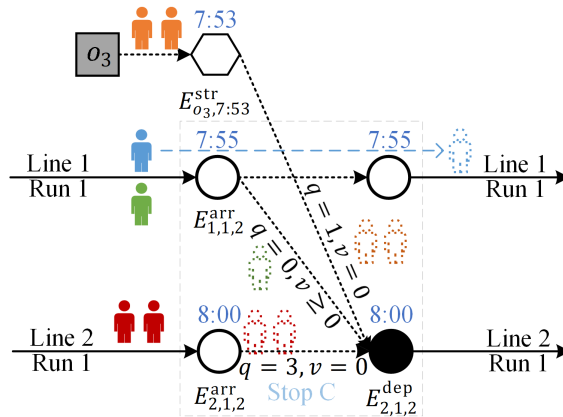


Figure 6: Illustration of the introduced variable $v$.

*meaning that passengers choosing these arcs need not worry about failing to board; hence the associated anxiety cost*

*is zero. In contrast, the available capacities of the lowest-priority transfer arc $(E_{1,1,2}^{arr}, E_{2,1,2}^{dep})$ are zero, implying that passengers choosing this arc are not assured of boarding. Consequently, they may incur a positive anxiety cost. For every OD pair $w \in \mathcal{W}$ and every $b \in B_w$, if we interpret the sum of the travel cost $c_{w,b}^r$ and the anxiety cost $\sum_{A \in \mathcal{A}_{w,b,r}^{priority}} v_A$ as a generalized cost of route $r$, then the equilibrium obtained from NCP (12) can be viewed as a state in which all passengers of the same OD pair and class experience the same generalized cost.*

Compared with the equivalent UEIP formulation (6) introduced in the previous section, the refined model (12) not only excludes behaviorally unrealistic solutions, but also offers computational advantages. In the equivalent formulation, the priority-enforcing constraint (6c) is imposed at the *route* level, leading to strongly coupled constraints that are difficult to decompose. As a result, all feasible routes must be enumerated in advance as input to the solution algorithm, which is computationally expensive. In contrast, in the refined model, the priority-enforcing constraint (12c) is imposed at the *arc* level, yielding a decomposable network-flow structure that is amenable to column-generation algorithms, which dynamically generate routes during the solution process.

## 4.2 Solution algorithms

To obtain the R-UEIP solution, we propose an algorithm that reformulates NCP (12) into an MPEC. A straightforward way to handle such complementarity systems is to apply a merit or smoothing function to all complementarity conditions and solve the resulting unconstrained optimization. However, the form of conditions (12a) and (12b) resembles a classical static traffic assignment problem (Beckmann et al., 1956). This structural similarity allows the flows to be computed efficiently using existing assignment algorithms such as TAPAS (Bar-Gera, 2010) and iGP (Xie et al., 2018), and the column-generation step for dynamically constructing the route set can likewise be delegated to these algorithms. To preserve this computational advantage, we keep (12a) and (12b) in their original form as the lower-level equilibrium constraints, and apply a merit function only to condition (12c), which forms the upper-level optimization.

We use the Fischer-Burmeister function (Fischer, 1992) to encode the complementarity condition (12c). For a scalar pair $(a, b)$, the FB function is defined as $\varphi(a, b) = \sqrt{a^2 + b^2} - (a + b)$, which satisfies $\varphi(a, b) = 0$ if and only if $0 \leq a \perp b \geq 0$. Applying this to the introduced variable $v$ and the available capacity $q(x)$, we define $\boldsymbol{\varphi}(v, x) = \left( \varphi(v_A, q_A(x)) \right)_{A \in \mathcal{A}^{priority}}$, and the associated merit function

$$\Psi(v, x) = \|\boldsymbol{\varphi}(v, x)\|^2 = \sum_{A \in \mathcal{A}^{priority}} \varphi(v_A, q_A(x))^2.$$

Then $v^*$ satisfies (12c) if and only if it minimizes $\Psi(v, x)$, where the available capacities $q$ depend on $v$ indirectly through the arc flow $x$.

The mapping from $v$ to $x$ is defined by (12a)–(12b). For every OD pair $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and route $r \in \mathcal{R}_w$, let

$$\hat{c}_{w,b}^r = c_{w,b}^r + \sum_{A \in \mathcal{A}_{w,b,r}^{priority}} v_A$$

be the generalized route cost, and collect these into $\hat{c}$. Using $\hat{c}$, conditions (12a)–(12b) can be rewritten as the equivalent equilibrium system

$$0 \leq f \perp \hat{c}(x, v) - \Lambda \mu \geq 0, \tag{13a}$$

$$0 \leq \mu \perp \Lambda^T f - d \geq 0. \tag{13b}$$

This lower-level problem has the same structure as a classical static traffic assignment model, with the usual route travel costs replaced by the generalized costs $\hat{c}$.

Combining the merit function $\Psi(v, x)$ with the equilibrium conditions (13), NCP (12) can be reformulated as the following MPEC:

$$\min_{v,x,\mu} \Psi(v, x) \quad \text{subject to} \quad v \geq 0 \text{ and Condition (13).} \tag{14}$$

Removing the non-negativity constraint on $v$ does not affect the theoretical equivalence to NCP (12). However, allowing negative components of $v$ would lead to negative generalized costs, violating standard assumptions in traffic assignment and complicating procedures such as shortest-path search. Therefore, the constraint $v \geq 0$ is retained in the reformulated model.

When the lower-level equilibrium problem (13) admits a unique arc flow solution, which is typically ensured when the generalized costs are strictly increasing functions of the arc flows (Patriksson, 2015), the mapping from $v$ to $x$ is well defined. In this case, MPEC (14) can be solved by an implicit approach. Specifically, we substitute the implicit function $x(v)$ into the objective of the MPEC, thereby eliminating the equilibrium constraints, and then apply a projected gradient algorithm to the resulting problem. The projection step only needs to enforce the nonnegativity constraint on $v$, and can therefore be carried out very efficiently. In contrast, when the lower-level equilibrium problem (13) does not have a unique arc flow solution, we reformulate the equilibrium conditions in MPEC (14) as a set of inequality constraints and solve the resulting model using a general nonlinear programming algorithm, such as sequential quadratic programming (SQP). The detailed algorithmic procedures for these two solution approaches — an implicit method and a nonlinear-programming-based method — are provided in Appendix C.

## 5  Numerical Examples

This section presents three numerical experiments for solving NCP (12) using the proposed method. The first experiment, based on a case study of student commuting at the University of Hong Kong, illustrates the method's ability to accurately capture priority rules and compares its results with those of an explicit priority model. The second experiment, using the benchmark network introduced by Nguyen et al. (2001), demonstrates that our method can accurately recover the R-UEIP. Finally, the third experiment evaluates the computational efficiency of the proposed algorithm on the Sioux Falls transit network. All results presented in this section are obtained on a Windows 10 (64-bit) PC equipped with an AMD Ryzen 7 4800H 2.90 GHz CPU and 16 GB of RAM.

**Cost structure.** We first specify the cost functions used in the numerical experiments. All cost functions are taken from the existing literature, and the costs associated with each type of arc are assumed to be identical for all passengers, regardless of their OD pair or class. Specifically, for each boarding arc $A = (E_{o,t}^{\text{str}}, E_{l,j,i}^{\text{dep}})$, the cost is

$$\eta_1 \cdot (\tau_{l,j,i}^{\text{dep}} - t),$$

which is the time duration multiplied by $\eta_1$, the value of time. Similarly, for each transfer arc $A = (E_{l,j,i}^{\text{arr}}, E_{l',j',i'}^{\text{dep}})$, the cost is

$$\eta_1 \cdot (\tau_{l',j',i'}^{\text{dep}} - \tau_{l,j,i}^{\text{arr}}).$$

18

For riding arcs $A = (E_{l,j,i}^{\text{dep}}, E_{l,j,i+1}^{\text{arr}})$, the cost is

$$\eta_1 \cdot (\tau_{l,j,i+1}^{\text{arr}} - \tau_{l,j,i}^{\text{dep}}) + \pi_l^{i,i+1} + \eta_2 \cdot \max\left\{0, \frac{x_A}{u_A} - \rho\right\},$$

where $\pi_l^{i,i+1}$ is the fare or penalty associated with traveling from stop $s_l^i$ to $s_l^{i+1}$ on line $l$, $\eta_2 > 0$ is the marginal disutility of crowding and $\rho \in [0,1]$ is the crowding perception threshold (Hamdouch and Lawphongpanich, 2008). Here, the congestion term becomes positive only when the load ratio $x_A/u_A$ exceeds the threshold $\rho$. For the dwelling arcs $A = (E_{l,j,i}^{\text{arr}}, E_{l,j,i}^{\text{dep}})$, the cost is defined similarly as

$$\eta_1 \cdot (\tau_{l,j,i}^{\text{dep}} - \tau_{l,j,i}^{\text{arr}}) + \eta_2 \cdot \max\left\{0, \frac{x_A}{u_A} - \rho\right\}.$$

The costs on the other two types of arcs — access arcs $\mathcal{A}^{\text{access}}$ and egress arcs $\mathcal{A}^{\text{egress}}$ — depend on the OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$. For each egress arc $A = (E_{l,j,i}^{\text{arr}}, E_d^{\text{dst}})$, the cost consists of the walking time from stop $s_l^i$ to destination $d$, together with early-arrival and late-arrival penalties. Formally, the cost is

$$t_{s_l^i, d}^{\text{walk}} + \begin{cases} \eta_3 \cdot (\tau_{w,b}^- - (\tau_{l,j,i}^{\text{arr}} + t_{s_l^i, d}^{\text{walk}})), & \text{if } \tau_{l,j,i}^{\text{arr}} + t_{s_l^i, d}^{\text{walk}} < \tau_{w,b}^-, \\ 0, & \text{if } \tau_{l,j,i}^{\text{arr}} + t_{s_l^i, d}^{\text{walk}} \in [\tau_{w,b}^-, \tau_{w,b}^+], \\ \eta_4 \cdot (\tau_{l,j,i}^{\text{arr}} + t_{s_l^i, d}^{\text{walk}} - \tau_{w,b}^+), & \text{if } \tau_{l,j,i}^{\text{arr}} + t_{s_l^i, d}^{\text{walk}} > \tau_{w,b}^+, \end{cases}$$

where $\eta_3 > 0$ and $\eta_4 > 0$ are coefficients capturing passengers' aversion to early and late arrival, respectively. Eventually, for access arc $A = (E_o^{\text{ogn}}, E_{o,t}^{\text{str}})$, the cost is

$$\eta_5 \cdot \max\{\tau_{w,b}^{\text{free}} - t, 0\},$$

where $\tau_{w,b}^{\text{free}}$ denotes the *free-flow latest departure time* (Nguyen et al., 2001), i.e., the latest time at which a passenger of OD pair $w$ and class $b$ can depart and still arrive before $\tau_{w,b}^+$ when no capacity-induced delays occur. Passengers departing earlier than this time incur a penalty. The route cost $c_{w,b}^r(f)$ is the total arc travel cost experienced by passengers of OD pair $w$, class $b$, and follow route $r \in \mathcal{R}_{w,b}$.

## 5.1 Insight from a student commuting case in HKU

This section presents a case study of morning-peak commuting by students at the University of Hong Kong to demonstrate that the proposed transit assignment model can accurately capture passengers' choice and queuing behaviors. As illustrated in Figure 7, passengers depart from the Jockey Club Student Village IV and travel to the main campus, with a desired arrival time interval of [9:00, 9:20]; being late would result in missing classes and incurring a substantial penalty. Two categories of routes are available. One option is to take Bus No. 71, which has a relatively long in-vehicle riding time of about 30 minutes. The other option is to take the *Mass Transit Railway* (MTR), which requires only about 17 minutes of in-vehicle riding time but depends on elevators with limited capacity to access the campus from the station. In practice, passengers often face elevator queues lasting from ten to several tens of minutes. This situation gives rise to two clear trade-offs: (i) whether to take the slower but more reliable bus or the faster MTR at the risk of long elevator queues, and (ii) for those choosing the MTR, whether to depart earlier to avoid being late or depart later to reduce waiting time in the queue. To illustrate these two trade-offs, Section 5.1.1 analyzes the assignment results produced by the proposed model, and Section 5.1.2 compares these results with those obtained from the explicit priority model (see Appendix A).
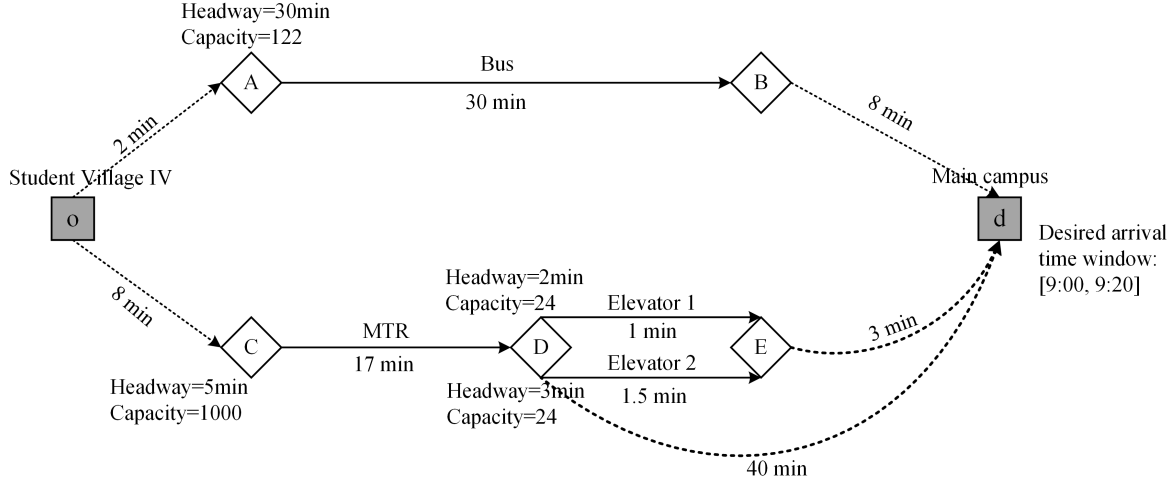
Figure 7: The example network of student commuting cases.

During peak hours, the elevators operate at full capacity with fixed stop patterns and can thus be modeled as scheduled transit lines. Two elevators are considered in this case: Elevator 1 (e1) runs directly from the MTR station to the hilltop campus, with a travel time of 1 minute and a headway of 2 minutes; Elevator 2 (e2) provides a non-direct service with a travel time of 1.5 minutes and a headway of 3 minutes. Both elevators have a capacity of 24 passengers. Alternatively, passengers may choose to walk from the MTR station to the campus, which takes 25 minutes with an additional fatigue-equivalent penalty of 15 minutes, resulting in an effective travel time of 40 minutes. Bus No. 71 operates with an Enviro500MMC 11.3 m vehicle of 122-passenger capacity, while the MTR line is served by Metro-Cammell EMU trains with a total capacity of 2504 passengers. After accounting for a background load of 1504 passengers, the feasible capacity of MTR is set to 1000. After applying the student concession, the fares for the MTR and bus are 4.4 HKD and 6.7 HKD, respectively. The parameters of the cost function are set as $\eta_1 = 0.55$, $\eta_2 = 5.0$, $\eta_3 = 0.2$, $\eta_4 = 10.0$, $\eta_5 = 1.0$, and $\rho = 0.8$. This case is solved using the implicit method described in Appendix C.1, with the algorithm implemented in C++.

### 5.1.1 Analysis of students' travel choices under the proposed model

Table 4 presents the assignment results between the bus and MTR under different total demand levels. Because the bus operates infrequently (one departure every 30 minutes), only the 8:40 a.m. trip is relevant to passengers in this scenario. The minimum travel cost of this bus trip (excluding crowding penalties and extra waiting) is 37.70. When the total demand is 200, 300, or 400, the maximum travel cost among MTR passengers is 30.35, 32.55, and 37.13, respectively, all lower than 37.70. Hence, no passengers chose the bus route. As the demand increases to 500, some MTR passengers must start their trips earlier to avoid being late. At this point, the maximum MTR travel cost nearly equals the bus cost, resulting in 93.39 passengers shifting to the bus. When the demand rises to 600 and 700, the maximum MTR travel cost increases further, the bus reaches full capacity (bus flow = 122.00), and an *early-departure queuing* phenomenon emerges: to secure higher boarding priority for the limited bus capacity, passengers arrive at stop A earlier than necessary. The extent of this departure advancement depends on the perceived benefit of switching — specifically, the difference between the maximum MTR travel cost and the bus travel cost without early

departure. For total demand levels of 600 and 700, the departure advancement is 2 minutes and 5 minutes, respectively, leading to approximate cost equalization between the two modes (with minor discrepancies caused by the discrete nature of timetables and departure-time choices).

Table 4: Assignment results between bus and MTR under different total demand levels.

| Demand | 200 | 300 | 400 | 500 | 600 | 700 |
|---|---|---|---|---|---|---|
| Max. MTR travel cost | 30.35 | 32.55 | 37.13 | 37.72 | 41.85 | 47.56 |
| Bus travel cost | 37.70 | 37.70 | 37.70 | 37.70 | 41.85 | 46.45 |
| Bus flow | 0.00 | 0.00 | 0.00 | 93.39 | 122.00 | 122.00 |
| Departure advancement | 0 min | 0 min | 0 min | 0 min | 2 min | 5 min |

We further analyze the trade-off between late arrival and early trip start among passengers choosing the MTR. Whether a passenger arrives late is determined by the departure time of the elevator they take: if the elevator departs later than 9:16 a.m., the passenger will be late. Figure 8 illustrates the distribution of passengers across different MTR trains and elevator runs under various demand levels. The horizontal axis represents the elevator departure time, the left vertical axis denotes the arrival time of MTR trains at station D, and the right vertical axis indicates the elevator run taken. Time is expressed in numerical format (e.g., 9:16 a.m. is represented as 76).
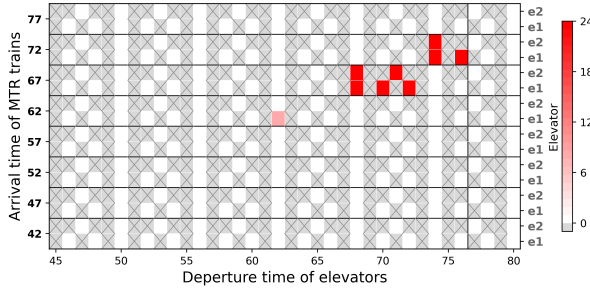
When the total demand is 200 or 300, no passengers arrive late (Figures 8(a) and 8(b)). This is because the minimum possible travel cost for a late-arriving passenger is 35.63, which exceeds the maximum MTR travel cost under these two demand levels. As demand continues to increase, newly added passengers must start their trips earlier to avoid arriving late. When the demand reaches 400, the maximum travel cost among on-time passengers rises to 37.13 — higher than 35.63 — indicating that some passengers now prefer to start later and accept a late-arrival penalty. As shown in Figure 8(c), 24 passengers make this choice. With further increases in total demand to 500, 600, and 700, the severity of lateness grows, and the number of late passengers increases to 24, 48, and 48, respectively (Figures 8(d) — 8(f)).

In addition to increasing the number of late passengers, higher demand also leads to longer elevator waiting times. For instance, among passengers taking the MTR train arriving at 9:02 a.m. (time 62), no queue forms when the demand is 200. When the demand rises to 300 and 400, the maximum elevator waiting time extends to 4 and 10 minutes, respectively. At a demand of 600, passengers experience at least 6 minutes and up to 14 minutes of queuing before boarding the elevator. Such severe congestion patterns are highly consistent with real-world observations, demonstrating the capability of the proposed model to realistically capture queuing phenomena.
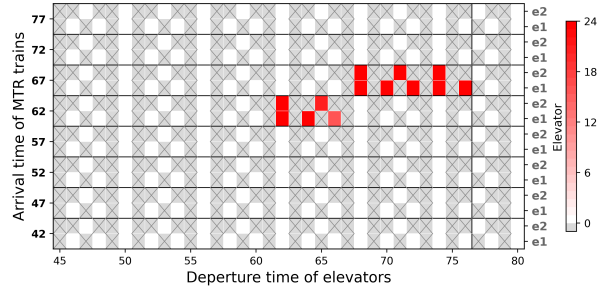
### 5.1.2 Comparison with the explicit priority model

For comparison, we solved an explicit priority model for this case study using an MSA-based algorithm implemented in C++ (see Appendix A for the model and algorithmic details). When the total demand reaches 600 or 700, the relative gap of the explicit priority model fails to converge below $10^{-2}$, rendering the results incomparable. Therefore, in this section, we set the total demand to 500, under which the relative gap converges below $10^{-5}$, indicating that a reasonably accurate equilibrium has been achieved.
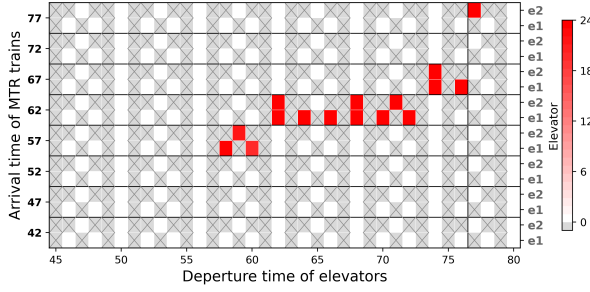
Overall, both the proposed implicit priority model and the explicit priority model yield equilibrium states that comply with the passenger priority rule, yet their flow distribution patterns differ. In the explicit priority model results, 72.21 passengers choose to travel by bus, whereas in the proposed model,
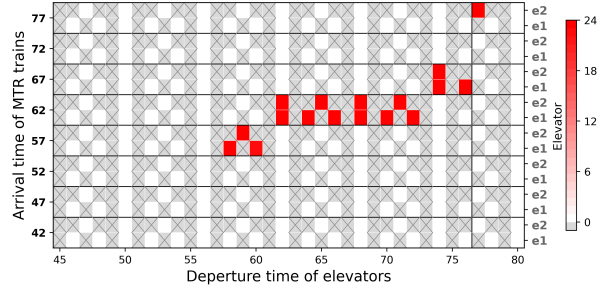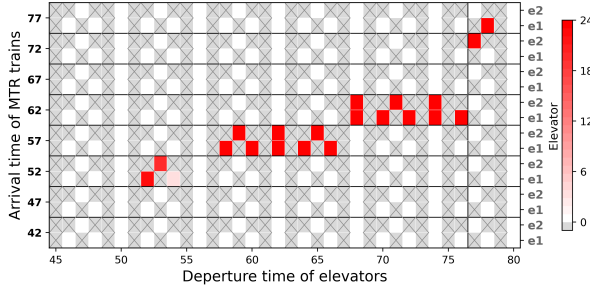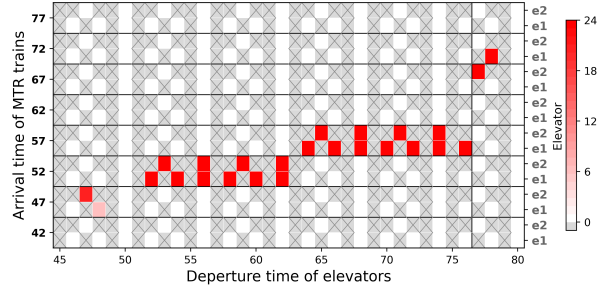
(a) Demand = 200

(b) Demand = 300

(c) Demand = 400

(d) Demand = 500

(e) Demand = 600

(f) Demand = 700

Figure 8: Passenger distribution across MTR and elevator runs under different total demand levels. Gray crossed cells denote no elevator departure, while the shade of red indicates the passenger count on a specific MTR–elevator combination.

as shown in Table 4, the bus flow reaches 93.39. The passenger distributions across MTR trains and the resulting elevator queues are illustrated in Figure 9. In both results, the elevator queues follow the FCFS principle — that is, passengers arriving on earlier MTR trains always board elevators before those from later trains (as visually reflected in Figure 9, where red blocks representing earlier arrivals appear below or to the left of those for later arrivals on the same elevator).



(a) The proposed implicit priority model
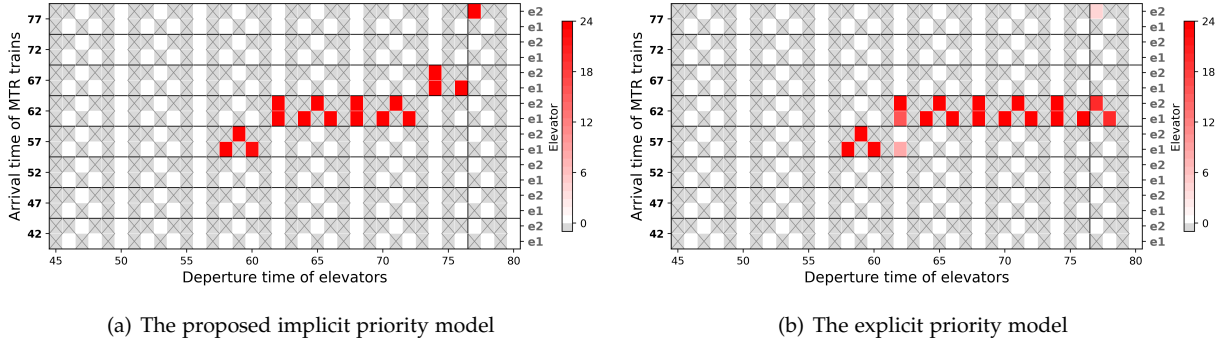
(b) The explicit priority model

Figure 9: Passenger distributions across MTR trains and elevator queues under the proposed implicit priority model and the explicit priority model.

However, the explicit priority model results exhibit one unrealistic phenomenon. Consider passengers arriving at 9:02 a.m. (time 62) and 9:17 a.m. (time 77) and the two elevators departing at 9:17 a.m. (time 77) and 9:18 a.m. (time 78). Some passengers who arrive earlier leave a capacity of 4.61 on elevator "e2" (departing at 77) for later arrivals, but instead board elevator "e1" (departing at 78), thereby incurring higher travel costs — an irrational choice. The cause lies in the decision mechanism of the explicit priority model, which is based on expected travel costs. Specifically, passengers arriving at 9:02 a.m. are divided into two transit plans that differ only in whether they take elevator "e1" or "e2". Although it would reduce the actual travel cost for some passengers to switch from e1 (departing at 9:18 a.m.) to "e2" (departing at 9:17 a.m.), such a transfer does not occur because, from the model's perspective, the two transit plans already have identical expected travel costs, and thus are in equilibrium. This example highlights a key limitation of the explicit priority model, mentioned earlier: it artificially treats passengers who are not intrinsically connected as a single decision-making entity and enforces equilibrium by assuming they make choices collectively based on expected travel costs. In contrast, our proposed model determines route choice based on actual travel costs, leading to greater behavioral realism and flexibility.

## 5.2 Verification accuracy in benchmark network

The benchmark network, illustrated in Figure 10, comprises six stops and four lines, each with a capacity of 20 passengers. Passengers are divided into four OD pairs $(o_1, d_1)$, $(o_1, d_2)$, $(o_2, d_1)$, and $(o_2, d_2)$, all having a desired arrival time window of [8:45, 9:00]. The demand for each OD is 10. The cost parameters are set as $\eta_1 = \eta_5 = 1.00$ and $\eta_2 = \eta_3 = \eta_4 = 0.00$. Nguyen et al. (2001) obtained a route flow vector $f^0$ by solving an approximate problem of the UEIP, as reported in the fourth column of Table 5. This solution does not satisfy the UEIP and the R-UEIP conditions, nor does it respect the vehicle capacity constraints. We adopt it as the initial route flow vector. The initial $v$ for all incoming arcs to the two saturated arcs, $(E_{2,1,2}^{dep}, E_{2,1,3}^{arr})$ and $(E_{1,1,2}^{dep}, E_{1,1,3}^{arr})$, is set to 5.00, while the initial $v$ for all other arcs is set to 0.00. The initial

23

$\mu$ is $\mu^0_{(o_1,d_1)} = 65.00, \mu^0_{(o_1,d_2)} = 75.00, \mu^0_{(o_2,d_1)} = 60.00, \mu^0_{(o_2,d_2)} = 75.00$. Here, we solve this example using a nonlinear-programming-based approach described in Appendix C.2, specifically the SQP solver of the *fmincon* function in MATLAB.
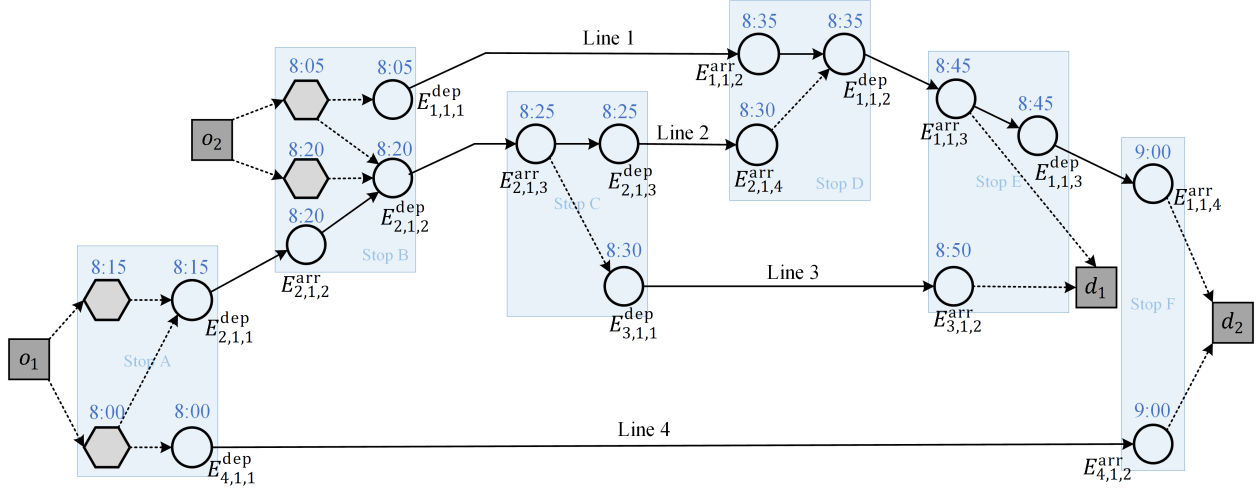


Figure 10: The benchmark network.

After 55 iterations, the algorithm converges in approximately 0.11 seconds; the final objective value is $\Psi(v, x) = 6.55 \times 10^{-6}$. This indicates that a highly accurate global optimum of the MPEC (14) has been obtained, corresponding to an R-UEIP solution. Columns 5, 6, and 7 of Table 5 report the optimal route flows $f^{r*}_{w,b}$, the total available capacities of the dominant routes $\hat{Q}^r_{w,b}$, and the travel costs $c^r_{w,b}$, respectively. It can be readily verified that all routes with $\hat{Q}^r_{w,b} > 0$ are unused, thereby confirming that the solution satisfies the R-UEIP definition (Definition 4).

Table 5: Results of route flow in the benchmark network.

| OD pairs | Routes | Route description | $f^0$ | $f^*$ | $\hat{Q}$ | $c$ |
|---|---|---|---|---|---|---|
| | 1 | 8:15 - Line 2 - Stop D - Line 1 | 0.000 | 0.000 | 0.000 | 30.000 |
| | 2 | 8:15 - Line 2 - Stop C - Line 3 | 10.000 | 10.000 | 0.000 | 35.000 |
| $(o_1, d_1)$ | 3 | 8:00 - Line 2 - Stop D - Line 1 | 0.000 | 0.000 | 10.000 | 60.000 |
| | 4 | 8:00 - Line 2 - Stop C - Line 3 | 0.000 | 0.000 | 0.00 | 65.000 |
| | 5 | 8:15 - Line 2 - Stop D - Line 1 | 6.236 | 5.000 | 0.000 | 45.000 |
| $(o_1, d_2)$ | 6 | 8:00 - Line 2 - Stop D - Line 1 | 0.000 | 0.000 | 5.000 | 75.000 |
| | 7 | 8:00 - Line 4 | 3.764 | 5.000 | 0.000 | 75.000 |
| | 8 | 8:20 - Line 2 - Stop D - Line 1 | 0.000 | 0.000 | 0.000 | 25.000 |
| | 9 | 8:20 - Line 2 - Stop C - Line 3 | 4.764 | 5.000 | 0.000 | 30.000 |
| $(o_2, d_1)$ | 10 | 8:05 - Line 2 - Stop D - Line 1 | 0.000 | 0.000 | 0.000 | 55.000 |
| | 11 | 8:05 - Line 2 - Stop C - Line 3 | 0.000 | 0.000 | 5.000 | 60.000 |
| | 12 | 8:05 - Line 1 | 5.236 | 5.000 | 0.000 | 55.000 |
| | 13 | 8:20 - Line 2 - Stop D - Line 1 | 0.000 | 0.000 | 0.000 | 40.000 |
| $(o_2, d_2)$ | 14 | 8:05 - Line 2 - Stop D - Line 1 | 0.000 | 0.000 | 0.000 | 70.000 |
| | 15 | 8:05 - Line 1 | 10.000 | 10.000 | 0.000 | 70.000 |

## 5.3 Verification computational efficiency in the Sioux Falls network

The Sioux Falls transit network, adapted from Szeto and Jiang (2014), consists of 24 stops, 10 lines, and 16 OD pairs. The demand for each OD is listed in Table 6. Passengers share a desired arrival time window of [60, 90]. As the original network specification provides only the headways of each line, the schedule is constructed by assigning the first departure of each line at time 0 according to the given headways. In total, the 10 lines generate 113 runs, and the resulting expanded event-activity graph comprises 1,525 nodes and 12,951 arcs. The cost function parameters are set as $\eta_1 = 1.0$, $\eta_2 = 2.0$, $\eta_3 = \eta_4 = \eta_5 = 1.2$, and $\rho = 0.0$. In this example, the solution is obtained via the implicit method.

Table 6: Demand data of the Sioux Falls transit network.

| OD | Demand | OD | Demand | OD | Demand | OD | Demand |
|------|--------|--------|--------|--------|--------|--------|--------|
| (1,13) | 200 | (2,13) | 100 | (3,13) | 200 | (4,13) | 200 |
| (1,20) | 100 | (2,20) | 100 | (3,20) | 100 | (4,20) | 100 |
| (1,21) | 100 | (2,21) | 100 | (3,21) | 100 | (4,21) | 100 |
| (1,24) | 200 | (2,24) | 100 | (3,24) | 200 | (4,24) | 200 |

Table 7 reports the number of iterations and the CPU time required by Algorithm 3 for the merit function $\Psi(v, x)$ to fall below the predefined thresholds. The initial objective function value was $1.02 \times 10^7$. After 3560 iterations (approximately 22.45 minutes), the value dropped below 10, and after 8276 iterations (55.31 minutes), it further decreased to below 1.

Table 7: The number of iterations and CPU time required to achieve a predetermined convergence threshold.

| Threshold | $10^4$ | $10^3$ | $10^2$ | 10 | 1 |
|-----------|--------|--------|--------|------|------|
| Iterations | 34 | 396 | 1116 | 3560 | 8276 |
| CPU time | 15.22 sec | 2.51 min | 6.80 min | 22.45 min | 55.31 min |

To more intuitively illustrate the convergence accuracy, Figure 11 presents the scatter plots of the available capacity $Q_{w,b}^{r',r}$ and the cost difference $c_{w,b}^r - c_{w,b}^{r'}$ for all used routes $r$ and their dominated counterparts $r'$, when $\Psi(v, x)$ falls below 10 and 1, respectively. According to the definition of the R-UEIP (Definition 4), for any route pair $(r, r')$, at least one of these two quantities — available capacity or cost difference — must be zero. As shown in Figure 11, when $\Psi(v, x) < 10$, nearly all points lie close to either the $x$-axis or the $y$-axis. The most deviated point has an available capacity of 111.77 and a cost difference of 2.69. When $\Psi(v, x) < 1$, the convergence precision further improves, as all points lie even closer to the coordinate axes. For the most deviated pair, the available capacity is 104.68, and the cost difference is 0.97. These results indicate that the solution has achieved a satisfactory level of convergence accuracy.

We also applied an MSA-based algorithm to solve the explicit priority model on this network. However, the algorithm failed to achieve an acceptable level of convergence: after one hour of computation, the minimum relative gap remained at 0.014. At this level of precision, the resulting route choices were still far from equilibrium. Table 8 reports, for each OD pair, the maximum cost difference among transit plan-departure time pairs with non-negligible flow (greater than 1.0). For example, in OD (1,20), the maximum cost difference reached 56.71. Within this OD, two transit plan-departure time pairs carried flows of 87.50 and 11.94, with expected travel costs of 73.21 and 129.92, respectively. Similarly, OD (2,13) exhibited

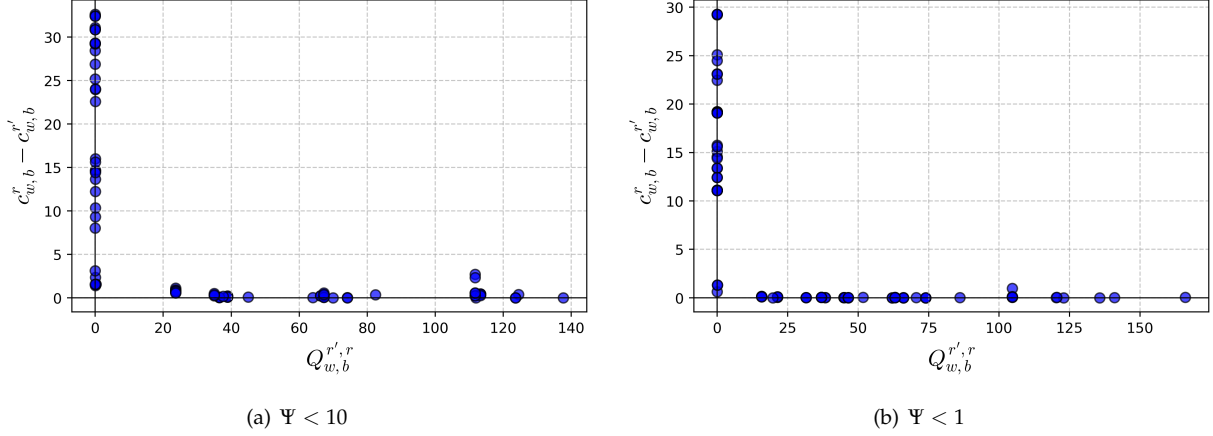(a) $\Psi < 10$                              (b) $\Psi < 1$

Figure 11: Available capacity and cost difference at different convergence precisions.

a maximum cost difference of 11.57 between two transit plan-departure time pairs with flows of 30.23 and 4.17 and expected costs of 82.93 and 94.50.

Table 8: Maximum cost differences within OD pairs under the solution of the explicit priority model.

| OD | Max. Diff. | OD | Max. Diff. | OD | Max. Diff. | OD | Max. Diff. |
|------|-----------|--------|-----------|--------|-----------|--------|-----------|
| (1,13) | 0.00 | (2,13) | 11.57 | (3,13) | 1.51 | (4,13) | 0.00 |
| (1,20) | 56.71 | (2,20) | 4.40 | (3,20) | 0.00 | (4,20) | 0.00 |
| (1,21) | 32.82 | (2,21) | 1.27 | (3,21) | 0.03 | (4,21) | 0.00 |
| (1,24) | 4.07 | (2,24) | 0.00 | (3,24) | 10.73 | (4,24) | 0.00 |

# 6   Conclusion

This paper revisits and extends a schedule-based transit assignment framework that captures passenger priority rules using an implicit method. We first reformulate the UEIP of Nguyen et al. (2001) as an NCP and, under mild conditions, prove the existence of an equilibrium state — providing a formal foundation that had not been previously established. The NCP formulation also reveals that equilibrium may be non-unique and can include behaviorally unrealistic states. To address this, we refine the original framework by modifying the definition of route availability and the NCP formulation so that all such behaviorally inconsistent equilibria are excluded. The refined model remains structurally tractable: by imposing the priority-enforcing condition at the arc rather than the route level, it avoids explicit enumeration of all feasible routes and admits a decomposable network-flow structure. By smoothing one complementarity condition with the Fischer-Burmeister function, we obtain a continuously differentiable MPEC, for which we develop two solution approaches — an implicit method and a nonlinear-programming–based method.

Using a real-world case study of student commuting at the University of Hong Kong, we first demonstrate that the model can realistically capture FCFS queuing at elevators as well as departure-early behavior among bus passengers driven by competition for boarding priority. A comparison with an explicit priority model highlights a key advantage of the proposed approach: by basing decisions on realized travel costs rather than expected ones, the refined framework avoids behaviorally implausible outcomes — for

example, passengers accepting substantial lateness penalties merely to preserve expected-cost consistency — and instead yields equilibria that are behaviorally coherent. In addition, numerical experiments on a benchmark and the Sioux Falls transit networks show that the proposed algorithms can compute high-accuracy equilibria with reasonable computational effort.

This study also suggests several directions for future research. First, the proposed algorithm does not guarantee global optimality, and its scalability to very large networks is limited. Developing algorithms with provable convergence properties that can efficiently handle large-scale networks would be of substantial practical value for real-world applications of the model. Second, although the UEIP equilibrium concept is formally defined, further behavioral validation is needed to explain how it can be reached. In the road traffic assignment literature, extensive effort has been devoted to understanding the behavioral foundations and stability properties of user equilibrium (UE) (e.g., Smith, 1984; Yang, 2005; Li et al., 2024b) and stochastic user equilibrium (SUE) (e.g., Horowitz, 1984; Cascetta, 1989; Watling and Hazelton, 2003; Cantarella and Watling, 2016). In a similar vein, a promising direction here is to investigate whether UEIP can be interpreted as the equilibrium outcome of a behaviorally plausible day-to-day adjustment process. Moreover, recent studies have leveraged day-to-day dynamical models to examine which equilibrium, among multiple solutions, is more likely to emerge in reality (Li et al., 2024a). Our findings on the non-uniqueness of UEIP, including the identification of solutions that appear behaviorally implausible, naturally raise the question of whether such outcomes can be ruled out from a dynamic perspective. Finally, most existing transit network design studies rely on transit assignment models without priority rules when evaluating design alternatives (Yin et al., 2021; Xie et al., 2021; Feng et al., 2025), and thus overlook the effects of queuing and denied boarding induced by passenger prioritization. This is largely because incorporating DNL-based priority assignment into network design problems leads to prohibitive computational complexity. Our DNL-free framework may offer a more tractable way to integrate priority-driven passenger behavior into transit network design, and this integration is a promising direction for future research.

# References

Bar-Gera, H. (2010). Traffic assignment by paired alternative segments. *Transportation Research Part B: Methodological*, 44(8-9):1022–1046.

Beckmann, M., McGuire, C. B., and Winsten, C. B. (1956). *Studies in the economics of transportation*. Yale University Press, New Haven, CT.

Beijing Daily News (2023). Long queues at shahe metro station during morning rush hour. Beijing Daily Online. Accessed on 14 November 2025, URL: `https://news.bjd.com.cn/2023/02/27/10349927.shtml`.

Cantarella, G. E. and Watling, D. P. (2016). Modelling road traffic assignment as a day-to-day dynamic, deterministic process: A unified approach to discrete-and continuous-time models. *EURO Journal on Transportation and Logistics*, 5(1):69–98.

Cascetta, E. (1989). A stochastic process approach to the analysis of temporal dynamics in transportation networks. *Transportation Research Part B: Methodological*, 23(1):1–17.

Cats, O. and West, J. (2020). Learning and adaptation in dynamic transit assignment models for congested networks. *Transportation Research Record: Journal of the Transportation Research Board*, 2674(1):113–124.

Cats, O., West, J., and Eliasson, J. (2016). A dynamic stochastic model for evaluating congestion and crowding effects in transit systems. *Transportation Research Part B: Methodological*, 89:43–57.

Cepeda, M., Cominetti, R., and Florian, M. (2006). A frequency-based assignment model for congested transit networks with strict capacity constraints: characterization and computation of equilibria. *Transportation Research Part B: Methodological*, 40(6):437–459.

Cominetti, R. and Correa, J. (2001). Common-lines and passenger assignment in congested transit networks. *Transportation Science*, 35(3):250–267.

Feng, L., Li, J., Xu, Z., and Xie, J. (2025). Bilevel transit timetabling with synchronization at origin and transfer stops. *Available at SSRN 5295135*.

Feng, L., Xie, J., Nie, Y. M., and Liu, X. (2020). Efficient algorithm for the traffic assignment problem with side constraints. *Transportation Research Record: Journal of the Transportation Research Board*, 2674(4):129–139.

Fischer, A. (1992). A special newton-type optimization method. *Optimization*, 24(3-4):269–284.

Friesz, T. L., Bernstein, D., Smith, T. E., Tobin, R. L., and Wie, B.-W. (1993). A variational inequality formulation of the dynamic network user equilibrium problem. *Operations Research*, 41(1):179–191.

Gentile, G. and Nökel, K. (2016). *Modelling public transport passenger flows in the era of intelligent transport systems*, volume 10. Springer.

Hamdouch, Y., Ho, H., Sumalee, A., and Wang, G. (2011). Schedule-based transit assignment model with vehicle capacity and seat availability. *Transportation Research Part B: Methodological*, 45(10):1805–1830.

Hamdouch, Y. and Lawphongpanich, S. (2008). Schedule-based transit assignment model with travel strategies and capacity constraints. *Transportation Research Part B: Methodological*, 42(7-8):663–684.

Hamdouch, Y., Szeto, W., and Jiang, Y. (2014). A new schedule-based transit assignment model with travel strategies and supply uncertainties. *Transportation Research Part B: Methodological*, 67:35–67.

Hong Kong Bus Wiki (2025). "ding jaa" (slang explanation). Hong Kong Bus Wiki, Fandom. Accessed on 14 November 2025, URL: `https://hkbus.fandom.com/wiki/%E9%A0%82%E9%96%98`.

Horowitz, J. L. (1984). The stability of stochastic equilibrium in a two-link transportation network. *Transportation Research Part B: Methodological*, 18(1):13–28.

Larsson, T. and Patriksson, M. (1995). An augmented Lagrangian dual algorithm for link capacity side constrained traffic assignment problems. *Transportation Research Part B: Methodological*, 29(6):433–455.

Li, J., Wang, Q., Feng, L., Xie, J., and Nie, Y. M. (2024a). A day-to-day dynamical approach to the most likely user equilibrium problem. *Transportation Science*, 58(6):1193–1213.

Li, J., Wang, Z., and Nie, Y. M. (2024b). Wardrop equilibrium can be boundedly rational: A new behavioral theory of route choice. *Transportation Science*, 58(5):973–994.

Nguyen, S., Pallottino, S., and Malucelli, F. (2001). A modeling framework for passenger assignment on a transport network with timetables. *Transportation Science*, 35(3):238–249.

Nie, Y. M., Zhang, H., and Lee, D.-H. (2004). Models and algorithms for the traffic assignment problem with link capacity constraints. *Transportation Research Part B: Methodological*, 38(4):285–312.

Nuzzolo, A., Crisalli, U., and Rosati, L. (2012). A schedule-based assignment model with explicit capacity constraints for congested transit networks. *Transportation Research Part C: Emerging Technologies*, 20(1):16–33.

Nuzzolo, A., Russo, F., and Crisalli, U. (2001). A doubly dynamic schedule-based assignment model for transit networks. *Transportation Science*, 35(3):268–285.

Papola, N., Filippi, F., Gentile, G., and Meschini, L. (2008). Schedule-based transit assignment: New dynamic equilibrium model with vehicle capacity constraints. In Nuzzolo, A. and Wilson, N. H. M., editors, *Schedule-based modeling of transportation networks: Theory and applications*, pages 1–26. Springer.

Patriksson, M. (2004). Sensitivity analysis of traffic equilibria. *Transportation Science*, 38(3):258–281.

Patriksson, M. (2015). *The traffic assignment problem: models and methods*. Courier Dover Publications.

Poon, M., Wong, S., and Tong, C. (2004). A dynamic schedule-based model for congested transit networks. *Transportation Research Part B: Methodological*, 38(4):343–368.

Smith, M. J. (1984). The stability of a dynamic model of traffic assignment—an application of a method of lyapunov. *Transportation Science*, 18(3):245–252.

Spiess, H. and Florian, M. (1989). Optimal strategies: A new assignment model for transit networks. *Transportation Research Part B: Methodological*, 23(2):83–102.

Szeto, W. and Jiang, Y. (2014). Transit assignment: Approach-based formulation, extragradient method, and paradox. *Transportation Research Part B: Methodological*, 62:51–76.

Tobin, R. L. and Friesz, T. L. (1988). Sensitivity analysis for equilibrium network flow. *Transportation Science*, 22(4):242–250.

Watling, D. and Hazelton, M. L. (2003). The dynamics and equilibria of day-to-day assignment models. *Networks and Spatial Economics*, 3(3):349–370.

Wu, J. H., Florian, M., and Marcotte, P. (1994). Transit equilibrium assignment: A model and solution algorithms. *Transportation Science*, 28(3):193–203.

Xie, J., Nie, Y. M., and Liu, X. (2018). A greedy path-based algorithm for traffic assignment. *Transportation Research Record: Journal of the Transportation Research Board*, 2672(48):36–44.

Xie, J., Zhan, S., Wong, S. C., and Lo, S. M. (2021). A schedule-based timetable model for congested transit networks. *Transportation Research Part C: Emerging Technologies*, 124:102925.

Xu, Z., Xie, J., Liu, X., and Nie, Y. M. (2020). Hyperpath-based algorithms for the transit equilibrium assignment problem. *Transportation Research Part E: Logistics and Transportation Review*, 143:102102.

Xu, Z., Xie, J., Liu, X., and Nie, Y. M. (2022). Hyperbush algorithm for strategy-based equilibrium traffic assignment problems. *Transportation Science*, 56(4):877–903.

Yang, F. (2005). *An evolutionary game theory approach to the day-to-day traffic dynamics*. Ph.D. dissertation, University of Wisconsin–Madison, Madison, WI.

Yao, X., Han, B., Yu, D., and Ren, H. (2017). Simulation-based dynamic passenger flow assignment modelling for a schedule-based transit network. *Discrete Dynamics in Nature and Society*, 2017(1):2890814.

Yin, J., D'Ariano, A., Wang, Y., Yang, L., and Tang, T. (2021). Timetable coordination in a rail transit network with time-dependent passenger demand. *European Journal of Operational Research*, 295(1):183–202.

Yin, J., Yang, L., Liang, Z., D'Ariano, A., and Gao, Z. (2025). Real-time rolling stock and timetable rescheduling in urban rail transit systems. *INFORMS Journal on Computing*, Ahead of Print.

# A    An Explicit-Priority Transit Assignment Model and Its Issues

Existing explicit priority models assume that passengers choose line segments (Yao et al., 2017) or strategies (Hamdouch and Lawphongpanich, 2008), while a DNL procedure determines their actual boarding, transfer, and alighting times, as well as their realized travel experiences. In this appendix, we consider a setting in which passengers choose their starting time and a sequence of line segments.

For an OD pair $w = (o,d) \in \mathcal{W} := \mathcal{O} \times \mathcal{D}$, a *transit plan* is any path in this transit network that starts at origin $o$, ends at destination $d$, and specifies the passenger's access stop, in-vehicle line segments (including any in-station transfers), and final alighting stop. We denote the set of such plans for OD pair $w$ by $\mathcal{P}_w$. Each passenger chooses a transit plan $p \in \mathcal{P}_w$ and a starting time $t \in \mathcal{T}$. Let $g_{w,b}^{p,t}$ denote the number of class $b$ passengers of OD pair $w$ choosing $(p,t)$, and collect all such variables in the vector $\boldsymbol{g}$. The feasible set $\mathcal{G}$ consists of all $\boldsymbol{g}$ satisfying

$$\sum_{p \in \mathcal{P}_w} \sum_{t \in \mathcal{T}} g_{w,b}^{p,t} = d_{w,b}, \quad \forall w \in \mathcal{W}, \ \forall b \in \mathcal{B}_w,$$

$$g_{w,b}^{p,t} \geq 0, \quad \forall w \in \mathcal{W}, \ \forall b \in \mathcal{B}_w, \ \forall (p,t) \in \mathcal{P}_w \times \mathcal{T}.$$

We first present the DNL procedure, followed by a description of the explicit priority model and its corresponding MSA-based solution algorithm. Finally, we use a simple example to illustrate the behavioral issues inherent in this class of explicit priority models.

## A.1    Dynamic network loading

Formally, DNL can be viewed as a mapping $D : \mathcal{G} \to \mathcal{F}$ that transforms any flow vector of transit plans and starting times $\boldsymbol{g} \in \mathcal{G}$ into the resulting flow vector of spatio-temporal routes $\boldsymbol{f} = D(\boldsymbol{g})$. Specifically, for every node $E \in \mathcal{E}$ and every arc $A \in \mathcal{A}$, let $x_E^{w,b,p,t}$ and $x_A^{w,b,p,t}$ denote the flow of OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$ associated with plan–departure pair $(p,t) \in \mathcal{P}_w \times \mathcal{T}$ passing through this node and this arc, respectively. The DNL procedure takes as input the given vector of transit plans and departure times $\boldsymbol{g} \in G$, and simulates the corresponding node and arc flows,

$$\left( x_E^{w,b,p,t} \right)_{\forall E \in \mathcal{E}, \forall w \in \mathcal{W}, \forall b \in \mathcal{B}_w, \forall (p,t) \in \mathcal{P}_w \times \mathcal{T}} \quad \text{and} \quad \left( x_A^{w,b,p,t} \right)_{\forall A \in \mathcal{A}, \forall w \in \mathcal{W}, \forall b \in \mathcal{B}_w, \forall (p,t) \in \mathcal{P}_w \times \mathcal{T}'}$$

subject to the boarding priority. Finally, for each OD pair $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and $(p,t) \in \mathcal{P} \times \mathcal{T}$, let $\mathcal{R}_{w,b}^{p,t}$ denote the set of all spatio-temporal routes whose starting time, boarding stop, line choices, transfer sequence, and final alighting stop coincide with those specified by $(p,t)$. The spatio-temporal route flows are obtained according to the flow-splitting proportions on each arc:

$$f_{w,b}^r = g_{w,b}^{p,t} \sum_{A \in \mathcal{A}_w^r} \frac{x_A^{w,b,p,t}}{\sum_{A' \in \mathcal{A}_{\text{tail}(A)}^+} x_{A'}^{w,b,p,t}}, \forall r \in \mathcal{R}_{w,b}^{p,t}, \tag{15}$$

where $\mathcal{A}_w^r$ is the set of arcs belonging to route $r$; $\text{tail}(A)$ is the tail node of $A$; and $\mathcal{A}_E^+$ and $\mathcal{A}_E^-$ are the set of outgoing and incoming arcs of $E$, respectively.

Algorithm 1 presents the DNL computation in detail. First, the flows of all plan–departure pairs $(p,t)$ for each OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$ are loaded at their origin node $E_o^{\text{ogn}}$:

$$x_{E_o^{\text{ogn}}}^{w,b,p,t} = g_{w,b}^{p,t}. \tag{16}$$

Subsequently, to preserve the temporal consistency of the transit system, the loading to all other nodes in the event–activity graph $\mathcal{H}(\mathcal{E}, \mathcal{A})$ is carried out in a combined topological and chronological order. There are three cases: (I) If a node $E$ represents a set-out event, then for each incoming arc $A \in \mathcal{A}_E^-$, passengers at the origin node can be loaded onto this arc only if their chosen departure time $t$ coincides with the timestamp of $E$. Let timestamp$(E)$ denote the timestamp of $E$. The arc flow is then updated as

$$
x_A^{w,b,p,t} = \begin{cases} x_{\text{tail}(A)}^{w,b,p,t}, & \text{if timestamp}(E) = t, \\ 0, & \text{if timestamp}(E) \neq t, \end{cases} \tag{17}
$$

(II) If a node $E$ represents a vehicle departure event, passengers arriving via the incoming arcs and boarding the subsequent in-vehicle riding arc riding$(E)$ are loaded according to the boarding priority until the vehicle capacity $u_{\text{riding}(E)}$ is reached. Accordingly, we traverse all incoming arcs of node $E$ in order of decreasing priority and compute the corresponding arc flows. For each incoming arc $A \in \mathcal{A}_E^-$, let $\bar{\mathcal{P}}_w^A$ denote the set of transit plans of OD pair $w \in \mathcal{W}$ that pass through the stop corresponding to $A$ and choose the line associated with $E = \text{head}(A)$. Let $\bar{x}_E^{w,b,p,t}$ and $\bar{u}_{\text{riding}(E)}$ denote, respectively, the remaining passenger flow and residual capacity after loading all higher-priority incoming arcs. The total number of passengers at tail$(A)$ who wish to board at event $E$ under capacity-free conditions is

$$
x_A^{\text{wishing}} = \sum_{w \in \mathcal{W}} \sum_{b \in \mathcal{B}_w} \sum_{p \in \bar{\mathcal{P}}_w^A} \sum_{t \in \mathcal{T}} \bar{x}_{\text{tail}(A)}^{w,b,p,t}.
$$

When capacity is binding and not all of these passengers can board, we assume they are selected with equal probability. Consequently, the arc flows are updated as

$$
x_A^{w,b,p,t} = \begin{cases} \bar{x}_{\text{tail}(A)}^{w,b,p,t}, & \text{if } \bar{u}_{\text{riding}(E)} \geq x_A^{\text{wishing}}, \\ \bar{x}_{\text{tail}(A)}^{w,b,p,t} \dfrac{\bar{u}_{\text{riding}(E)}}{x_A^{\text{wishing}}}, & \text{if } \bar{u}_{\text{riding}(E)} < x_A^{\text{wishing}}, \end{cases} \tag{18}
$$

(III) If node $E$ corresponds to any other type of event, no priority rules or capacity constraints apply. For any incoming arc $A \in \mathcal{A}_E^-$, all flow at the tail node tail$(A)$ is directly loaded onto the arc. Thus, the arc flows are updated as

$$
x_A^{w,b,p,t} = x_{\text{tail}(A)}^{w,b,p,t}. \tag{19}
$$

In all cases, the flow at node $E$ is obtained by aggregating the flows on all incoming arcs:

$$
x_E^{w,b,p,t} = \sum_{A \in \mathcal{A}_E^-} x_A^{w,b,p,t}. \tag{20}
$$

---

**Algorithm 1** Dynamic network loading

---

1: **Input:** The node set $\mathcal{E}^*$ sorted in topological and chronological order. The incoming arc set $\mathcal{A}_E^{-*}$ is sorted in priority order for all departure node $E$. The flow vector $g$.
2: Load the passenger flows on the virtual origin nodes according to Equation (16).
3: **for** each node $E \in \mathcal{E}^*$ and $E \notin O$ **do**
4:     **if** $E$ is a set-out node **then**
5:         **for** each arc $A \in \mathcal{A}_E^-$ **do**
6:             Update $x_A^{w,b,p,t}$ according to Equations (17).
7:         **end for**

8:   **else**
9:       **if** $E$ is a vehicle departure node **then**
10:          Set $\bar{u}_{\text{riding}(E)} = u_{\text{riding}(E)}$.
11:          **for** each arc $A \in \mathcal{A}_E^{-*}$ **do**
12:              Update $x_A^{w,b,p,t}$ according to Equations (18).
13:              Set $\bar{u}_{\text{riding}(E)} = \bar{u}_{\text{riding}(E)} - \sum_{w \in \mathcal{W}} \sum_{b \in \mathcal{B}_w} \sum_{p \in \mathcal{P}_w^A} \sum_{t \in \mathcal{T}} x_A^{w,b,p,t}$, and $\bar{x}_E^{w,b,p,t} = \bar{x}_E^{w,b,p,t} - x_A^{w,b,p,t}$.
14:          **end for**
15:       **else**
16:          **for** each arc $A \in \mathcal{A}_E^{-}$ **do**
17:              Update $x_A^{w,b,p,t}$ according to Equations (19).
18:          **end for**
19:       **end if**
20:    **end if**
21:    Update $x_E^{w,b,p,t}$ according to Equations (20). Set $\bar{x}_E^{w,b,p,t} = x_E^{w,b,p,t}$.
22: **end for**
23: Obtain the spatio-temporal route flow according to Equations (15).

## A.2 Variational inequality model and the MSA algorithm

A key implication of DNL under capacity constraints is that passengers choosing the same transit plan and starting time $(p, t)$ may experience heterogeneous realized experiences. Consequently, $(p, t)$ does not correspond to a single realized spatio-temporal route, and its travel cost is not uniquely defined for all passengers adopting that strategy. Explicit-priority models therefore evaluate $(p, t)$ using an *expected cost*, defined as the flow-weighted average of the realized costs over all spatio–temporal routes induced by DNL:

$$
e_{w,b}^{p,t}(g) = \begin{cases} \displaystyle\sum_{r \in \mathcal{R}_{w,b}^{p,t}} c_{w,b}^{r}(f) \, \frac{f_{w,b}^{r}}{g_{w,b}^{p,t}}, & \text{if } g_{w,b}^{p,t} > 0, \\[3ex] c_{w,b}^{\bar{r}}(f), & \text{if } g_{w,b}^{p,t} = 0, \end{cases}
\tag{21}
$$

where $f = D(g)$. Here, $\bar{r} \in \mathcal{R}_{w,b}^{p,t}$ is the spatio-temporal route obtained by injecting an infinitesimal flow of class $b$ passengers adopting $(p, t)$ on top of the background loading $g$ and, at each boarding, dwelling, or transfer opportunity, assigning this infinitesimal flow to the first-arriving vehicle of the chosen line that has residual capacity. This convention ensures that $e_{w,b}^{p,t}(g)$ is well-defined even when $(p, t)$ is unused.

**Definition 5** (User equilibrium with explicit priority). *A feasible flow vector $g \in \mathcal{G}$ is a* user equilibrium with explicit priority *(UEEP) if, for every OD pair $w \in \mathcal{W}$, class $b \in \mathcal{B}_w$, and strategy $(p, t) \in \mathcal{P} \times \mathcal{T}$, no passenger can reduce their expected cost $e_{w,b}^{p,t}(g)$ by unilaterally switching to another transit plan or starting time.*

The definition of UEEP is equivalent to the following variational inequality problem:

$$
e(g^*)^{\top}(g - g^*) \geq 0, \quad \forall g \in G,
\tag{22}
$$

where $g^*$ denotes the equilibrium flow vector and $e(\cdot)$ is the corresponding expected cost mapping.

Since there is no closed-form relationship between the flow vector $g$ and the expected cost vector $e$, the mapping $e(g)$ is, in general, not differentiable. Therefore, first-order algorithms are typically employed to solve problem (22), among which the method of successive averages (MSA) is the most commonly used. The detailed steps of the algorithm are given in Algorithm 2.

**Algorithm 2** Method of Successive Averages

1: **Step 0.** Initialize with a feasible solution $g^1 \in G$ obtained from the all-or-nothing assignment. Set $k = 1$.

2: **Step 1.** Execute Algorithm 1 to obtain the arc flow $x$ and route flow $f$. Update the expected cost $e(g^k)$.

3: **Step 2.** For each OD pair $w \in \mathcal{W}$ and each class $b \in B_w$, identify the minimum-cost pair $(\bar{p}, \bar{t})$ and set $y_{w,b}^{\bar{p},\bar{t}} = d_{w,b}$ and $y_{w,b}^{p,t} = 0$ for all $p \neq \bar{p}, t \neq \bar{t}$. Denote $y = (y_{w,b}^{p,t})_{\forall w \in \mathcal{W}, \forall b \in \mathcal{B}_w, \forall (p,t) \in \mathcal{P}_w \times \mathcal{T}}$.

4: **Step 3.** Compute the relative gap:

$$RG = 1 - \frac{\sum_{w \in \mathcal{W}} \sum_{b \in \mathcal{B}_w} e_{w,b}^{\bar{p},\bar{t}} d_{w,b}}{\sum_{w \in \mathcal{W}} \sum_{b \in \mathcal{B}_w} \sum_{p \in \mathcal{P}_w} \sum_{t \in \mathcal{T}} e_{w,b}^{p,t} g_{w,b}^{p,t}}.$$

If $RG < \epsilon$, **stop**; otherwise, update $g^{k+1} = \frac{kg^k + y}{k+1}$, set $k = k + 1$, and return to Step 1.

## A.3 Illustrative example and behavioral issues

This class of explicit-priority models cannot fully accommodate the inherently discrete nature of transit service in time, which leads to behaviorally questionable outcomes. To illustrate this point, we revisit the example used in Section 2.2 for the implicit priority model. Table 9 and Figure 12 present a user equilibrium outcome under the explicit priority model (22).

Table 9: User equilibrium outcome of the explicit priority model on the example transit network.

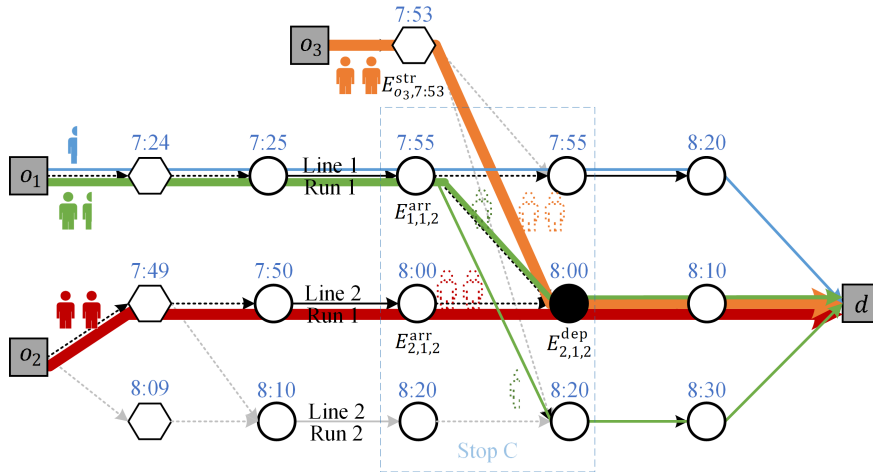| OD | Starting time | Path description | Flow | Expected cost |
|---|---|---|---|---|
| $(o_1, d)$ | **7:24** | **Line 1** | **0.5** | **56** |
| $(o_1, d)$ | **7:24** | **Line 1 - Line 2** | **1.5** | $46 \times 1/1.5 + 76 \times 0.5/1.5 = 56$ |
| $(o_2, d)$ | 7:49 | Line 2 | 2 | 21 |
| $(o_3, d)$ | 7:53 | Line 1 | 0 | 27 |
| $(o_3, d)$ | 7:53 | Line 2 | 2 | 17 |



Figure 12: Visualization of user equilibrium flows for the explicit priority model on the example transit network.

Among the passengers of OD pair $(o_1, d)$ who choose to take line 1 and transfer to line 2, vehicle capacity and priority rules imply that only one passenger can successfully board line 2 run 1 at stop C

(with a cost of 46), while the remaining 0.5 units of passenger flow must wait for line 2 run 2 (with a cost of 76). The expected cost for passengers choosing this option equals the cost of choosing "line 1 without transfer", both being 56, and the solution is therefore regarded as an equilibrium. In reality, however, the passengers who are left behind and forced to wait for line 2 run 2, thus incurring the higher realized cost, would be inclined to change their future decisions, for instance by choosing not to transfer.

In dynamic traffic assignment on road networks (Friesz et al., 1993), where the service provided by each link is typically modeled as continuous over time, explicit priority frameworks can partially alleviate the discrepancy between expected and realized costs by reducing the time-step size and thereby shrinking the size of each passenger group. In public transit, by contrast, service is inherently discontinuous in time, and existing explicit priority frameworks have not been modified to account for this discontinuity. As a result, they cannot bridge the gap between expected costs in the equilibrium definition and the actual travel experience of individual passengers.

# B   Proof of Proposition 5

Given that the mapping of NCP (12)

$$\bar{H}(f, \mu, v) = \begin{pmatrix} c(f) + \bar{\Delta}v - \Lambda\mu \\ \Lambda^T f - d \\ q(x) \end{pmatrix}$$

is continuous, while the feasible set $\bar{\Omega} = \{f \geq 0, \mu \geq 0, v \geq 0\}$ is unbounded and therefore not compact, the proof of Proposition 5 follows the same structure as that of Proposition 2. We introduce upper-bound constraints for each variable, yielding a modified NCP whose feasible set is compact. We then show that none of these upper bounds are binding at the solution, implying that any solution of the modified NCP also satisfies the original NCP (12).

Choose two scalars $e_1$ and $e_2$ such that

$$e_1 > \max\{d_{w,b} : w \in \mathcal{W}, b \in \mathcal{B}_w\} \text{ and } e_2 > \max\{c_{w,b}^r : w \in \mathcal{W}, b \in \mathcal{B}_w, r \in \mathcal{R}_w\}.$$

Let $\bar{\Omega}' = \bar{\Omega} \cap \{f \leq e_1 \mathbf{1}, \mu \leq e_2 \mathbf{1}, v \leq e_2 \mathbf{1}\}$. Since $\bar{\Omega}'$ is compact, the following NCP must admit a solution:

$$0 \leq f \perp c(f) + \bar{\Delta}v - \Lambda\mu + \kappa \geq 0, \tag{23a}$$

$$0 \leq \mu \perp \Lambda^T f - d + \rho \geq 0, \tag{23b}$$

$$0 \leq v \perp q(x) + \theta \geq 0, \tag{23c}$$

$$0 \leq \kappa \perp e_1 \mathbf{1} - f \geq 0, \tag{23d}$$

$$0 \leq \rho \perp e_2 \mathbf{1} - \mu \geq 0, \tag{23e}$$

$$0 \leq \theta \perp e_2 \mathbf{1} - v \geq 0. \tag{23f}$$

Let $(f^*, \mu^*, v^*)$ denote this solution. If we demonstrate that $\kappa = \rho = \theta = 0$, then $(f^*, \mu^*, v^*)$ must also be a solution of NCP (6).

**For $\kappa$:** Suppose $\kappa_{w,b}^r > 0$ for some $w \in \mathcal{W}, b \in \mathcal{B}_w, r \in \mathcal{R}_w$. By the complementarity condition (23d), we have $f_{w,b}^{r*} = e_1 > d_{w,b}$, so $\sum_{r \in \mathcal{R}_w} f_{w,b}^{r*} - d_{w,b} + \rho_{w,b} > 0$. Then, according to Coondition (23b), we have

34

$\mu_{w,b}^* = 0$. Moreover, since $f_{w,b}^{r*} = e_1 > 0$, we have

$$0 = c_{w,b}^r(\boldsymbol{f}^*) + \sum_{A \in \mathcal{A}_{w,b,r}^{\text{priority}}} v_A^* - \mu_{w,b}^* + \kappa_{w,b}^r = c_{w,b}^r(\boldsymbol{f}^*) + \sum_{A \in \mathcal{A}_{w,b,r}^{\text{priority}}} v_A^* + \kappa_{w,b}^r > 0,$$

because $c_{w,b}^r(\boldsymbol{f}^*)$ and $\kappa_{w,b}^r$ are positive, and $\sum_{A \in \mathcal{A}_{w,b,r}^{\text{priority}}} v_A^*$ is non-negative. This contradiction yields $\boldsymbol{\kappa} = \mathbf{0}$.

**For $\boldsymbol{\rho}$**: For each OD pair $w \in \mathcal{W}$ and class $b \in \mathcal{B}_w$, since there exists a route $r \in \mathcal{R}_w$ with $Q_{w,b}^r(\boldsymbol{f}^*) > 0$ by assumption, each arc belonging to $\mathcal{A}_{w,b,r}^{\text{priority}}$ satisfies $q_A(\boldsymbol{x}^*) > 0$ and $q_A(\boldsymbol{x}^*) + \theta_A > 0$. By Condition (23c), we have $v_A^* = 0$ for all $A \in \mathcal{A}_{w,b,r}^{\text{priority}}$. Therefore,

$$0 \leq c_{w,b}^r(\boldsymbol{f}^*) + \sum_{A \in \mathcal{A}_{w,b,r}^{\text{priority}}} v_A^* - \mu_{w,b}^* + \kappa_{w,b}^r = c_{w,b}^r(\boldsymbol{f}^*) - \mu_{w,b}^*$$

$$\Rightarrow \mu_{w,b}^* \leq c_{w,b}^r(\boldsymbol{f}^*) < e_2. \tag{24}$$

This means that $e_2 - \mu_{w,b}^* > 0$ for all $w \in \mathcal{W}$ and $b \in \mathcal{B}_w$, and $\boldsymbol{\rho} = \mathbf{0}$ by Equation (23e).

**For $\boldsymbol{\theta}$**: For any departure event, let $A_1$ be the arc that ends at this event with the highest loading priority. If $\theta_{A_1} > 0$, the complementarity condition (23f) suggests $v_{A_1}^* = e_2$. Meanwhile, Equation (24) implies that $\mu_{w,b}^* < e_2 = v_{A_1}^*$ for all $w \in \mathcal{W}$ and $b \in \mathcal{B}_w$. Hence, for each route $r$ such that $A_1 \in \mathcal{A}_{w,b,r}^{\text{priority}}$, we have

$$c_{w,b}^r(\boldsymbol{f}^*) + \sum_{A \in \mathcal{A}_{w,b,r}^{\text{priority}}} v_A^* - \mu_{w,b}^* + \kappa_{w,b}^r = c_{w,b}^r(\boldsymbol{f}^*) + \sum_{A \in \mathcal{A}_{w,b,r}^{\text{priority}}} v_A^* - \mu_{w,b}^* > v_{A_1}^* - \mu_{w,b}^* > 0.$$

By Condition (23a), all of these routes satisfy $f_{w,b}^{r*} = 0$. This means that the flow of arc $A_1$ is zero, namely $x_{A_1}^* = \sum_{w \in \mathcal{W}} \sum_{b \in \mathcal{B}_w} \sum_{r \in \mathcal{R}_w} f_{w,b}^{r*} \delta_{w,b}^{r,A_1} = 0$. Then, we have the available capacity $q_{A_1}(\boldsymbol{x}^*) = u_{\text{riding}(A_1)} - x_{A_1}^* = u_{\text{riding}(A_1)} > 0$, and $q_{A_1}(\boldsymbol{x}^*) + \theta_{A_1} > 0$, and consequently, $v_{A_1}^* = 0$ by Condition (23c). This contradicts $v_{A_1}^* = e_2 > 0$, and hence $\theta_{A_1} = 0$.

Next, let $A_2$ denote the arc ending at this event with the second-highest priority. If $\theta_{A_2} > 0$, we can get $v_{A_2}^* = e_2$ and $x_{A_2}^* = 0$ using a similar proof. This lead to the available capacity of $A_2$ is equal to that of $A_1$, i.e., $q_{A_2}(\boldsymbol{x}^*) = u_{\text{riding}(A_2)} - x_{A_1}^* - x_{A_2}^* = q_{A_1}$. In addition, the result $\theta_{A_1} = 0$ above, together with Condition (23c), indicates that $q_{A_1} \geq 0$. Hence, $q_{A_2} \geq 0$ and $q_{A_2} + \theta_{A_2} > 0$. According to Condition (23c), we encounter the contradiction that $0 = v_{A_2}^* = e_2 > 0$. Therefore, $\theta_{A_2} = 0$.

Proceeding in this manner for all arcs ending at the departure event, and repeating the process for every departure event in the network, we obtain $\boldsymbol{\theta} = \mathbf{0}$.

# C  Algorithms for MPEC (14)

This section introduces the implicit method, the nonlinear-programming-based method, and the method for obtaining an initial solution for solving MPEC (14).

## C.1  Implicit method

Considering the implicit function $\boldsymbol{x}(\boldsymbol{v})$, MPEC (14) is simplified to a problem involving only the upper-level variable $\boldsymbol{v}$, which is $\min_{\boldsymbol{v} \geq 0} \Psi(\boldsymbol{v}, \boldsymbol{x}(\boldsymbol{v}))$. We develop a projected gradient method combined with an Armijo line search to solve it. The entire algorithm can be summarized as follows.

**Algorithm 3** Implicit method

---

1: **Input:** A initial solution $(v^0, x^0, \mu^0)$. Set $k = 0$.

2: **Step 1.** Calculate the gradient $\nabla \Psi(v^k)$.

3: **Step 2.** For each $i = 0, 1, 2, \cdots$: Let $v' = [v^k - 2^{-i} \nabla \Psi(v^k)]_+$.

4:     **Step 2.1.** Obtain a solution $x'$ of Problem (13), and then compute the available capacity $q(x')$ and the merit function $\Psi(v')$.

5:     **Step 2.2.** If $\Psi(v') \leq \Psi(v^k) - \gamma 2^{-i} \nabla \Psi(v^k)^T \nabla \Psi(v^k)$, go to Step 3.

6: **Step 3.** Set $v^{k+1} = v'$. If $\Psi(v^{k+1}) < \epsilon_1$, then stop. Otherwise, $k = k + 1$, and go to Step 1.

---

In Step 1, the gradient of the merit function $\Psi(v, x(v))$ regarding $v_A$ is

$$\frac{\partial \Psi(v, x(v))}{\partial v_A} = \sum_{A' \in \mathcal{A}^{\text{priority}}} 2\varphi(v_{A'}, q_{A'}) \frac{\partial \varphi(v_{A'}, q_{A'})}{\partial v_A},$$

$$\frac{\partial \varphi(v_{A'}, q_{A'})}{\partial v_A} = \begin{cases} (v_{A'} + q_{A'} \frac{\partial q_{A'}}{\partial v_A}) \frac{1}{\sqrt{v_{A'}^2 + q_{A'}^2}} - 1 - \frac{\partial q_{A'}}{\partial v_A}, & \text{if } A' = A \\ q_{A'} \frac{\partial q_{A'}}{\partial v_A} \frac{1}{\sqrt{v_{A'}^2 + q_{A'}^2}} - \frac{\partial q_{A'}}{\partial v_A}, & \text{otherwise} \end{cases},$$

$$\frac{\partial q_{A'}}{\partial v_A} = - \sum_{A'' \in \text{Prior}(A')} \frac{\partial x_{A''}}{\partial v_A}.$$

Since the mapping from $v$ to $x$ is the NCP (13), the derivative $\frac{\partial x_{A'}}{\partial v_A}$ does not have a closed-form formulation. To obtain this derivative, we employ the sensitivity analysis method (Tobin and Friesz, 1988; Patriksson, 2004).

In Step 2, $[\cdot]_+$ denotes the projection onto the non-negative space. The lower-level problem (13) is solved using the iGP algorithm, an efficient static traffic assignment problem solver proposed by Xie et al. (2018).

## C.2 Nonlinear-programming-based method

We rewrite the equilibrium constraints (13) as a set of nonlinear inequalities, which is

$$J(v, x, \mu) = \begin{Bmatrix} f \\ \hat{c}(x, v) - \Lambda\mu \\ -f(\hat{c}(x, v) - \Lambda\mu) \\ \mu \\ \Lambda^T f - d \\ -\mu(\Lambda^T f - d) \\ v \end{Bmatrix} \geq 0.$$

This leads to the transformation of MPEC (14) into the following nonlinear programming problem

$$\min_{v, x, \mu} \Psi(v, x, \mu) \text{ subject to } J(v, x, \mu) \geq 0.$$

This nonlinear programming problem can be solved using various existing nonlinear optimization algorithms, such as Sequential Quadratic Programming (SQP).

## C.3 Initial solution

To initialize the MPEC (14), we first obtain the initial $x^0$ and $\mu^0$ by solving a relaxed transit assignment in which the available capacity constraints $q_A(x) \geq 0, \forall A \in \mathcal{A}^{\text{priority}}$ are replaced by standard vehicle capacity constraints $u_A \geq x_A, \forall A \in \mathcal{A}^{\text{riding}}$. Let $u = \{u_A : A \in \mathcal{A}^{\text{riding}}\}$ denote the capacity vector of riding arcs. The relaxed complementarity formulation is

$$0 \leq f \perp c(f) + \hat{\Delta}v - \Lambda\mu \geq 0,$$
$$0 \leq \mu \perp \Lambda^T f - d \geq 0,$$
$$0 \leq v \perp u(x) \geq 0,$$

where $v = \{v_A : A \in \mathcal{A}^{\text{riding}}\}$ are the multipliers for vehicle capacity constraints and $\hat{\Delta} = [\delta^{r,A}_{w,b}]_{|n_f| \times |\mathcal{A}^{\text{riding}}|}$ is the route–arc incidence matrix. This model can be regarded as an extension of the classical capacitated traffic assignment problem (Larsson and Patriksson, 1995) to the transit context (with a different definition of the travel cost function $c$), and can be efficiently solved by existing algorithms such as ALM-GP or ALM-Greedy (Nie et al., 2004; Feng et al., 2020).

Once $x^0$ is obtained, we initialize the upper-level vector $v^0$ from the multipliers $v$ of the relaxed problem. For each arc $A \in \mathcal{A}^{\text{priority}}$, we set

$$v^0_A = \begin{cases} v_{\text{riding}(A)}, & q_A \leq 0 \\ 0, & q_A > 0 \end{cases}.$$

This initialization assigns positive values to $v^0$ only on unavailable arcs, thereby guiding the subsequent iterations toward satisfying the available capacity constraints $q(x) \geq 0$. This initialization method is used as the default unless otherwise stated.