

**Integrated vehicle assignment and routing for system-optimal shared mobility planning with endogenous road congestion**

Jiangtao Liu

School of Sustainable Engineering and the Built Environment,  
Arizona State University, Tempe, AZ, 85281, USA  
Email: [jliu215@asu.edu](mailto:jliu215@asu.edu)

Pitu Mirchandani

School of Computing, Informatics, and Decision Systems Engineering  
Arizona State University, Tempe, AZ, 85281, USA  
Email: [pitu@asu.edu](mailto:pitu@asu.edu)

Xuesong Zhou

School of Sustainable Engineering and the Built Environment,  
Arizona State University, Tempe, AZ, 85281, USA  
Email: [xzhou74@asu.edu](mailto:xzhou74@asu.edu)  
Tel.: +1 480 9655827  
(Corresponding Author)

Submitted for publication in Transportation Research Part C

## Abstract

The ridesharing services have been growing in recent years with the start of network service companies, and will be further enhanced by the recently emerging trend of autonomous vehicles applications for future traveler mobility. One fundamental question that transportation managers should address is how to capture the endogenous traffic patterns with those upcoming new and uncertain elements for future transportation planning and management. Therefore, by focusing on one ideal system optimal (SO) scenario in which (i) all vehicles are autonomous or can be centrally guided and (ii) all passengers' pickup/drop-off trip requests are given at the beginning, this paper aims to integrate travel demand, vehicle supply and limited infrastructure supply by optimally assigning available rideshared and autonomous vehicles from different (real/virtual) depots to satisfy passengers' trip requests while considering the endogenous congestion in capacitated networks. A number of decomposition approaches are adopted in this research. Focusing on this primal problem, we propose an arc-based vehicle-based integer linear programming model in space-time-state (STS) networks, which is solved by Dantzig-Wolfe decomposition. From the perspective of dynamic traffic assignment, a space-time-state (STS) path-based flow-based linear programming model is also provided as an approximation according to the mapping information between vehicle and passenger and between vehicle and space-time arc in each STS path in our priori generated column pool. Later, we apply Alternating Direction Method of Multipliers (ADMM) to solve this linear programming model and compare its results with standard solver. Finally, numerical experiments are performed to demonstrate our decomposition approaches and their computation efficiency. From our preliminary experiments, we have a few interesting observations: (i) without considering the road congestion, the network performance/efficiency could be overestimated; (ii) passengers' required pickup and drop-off time windows could be a buffer to mitigate road congestion without impacting system performance; (iii) the ridesharing service could reduce the total transportation system cost under centralized control; (iv) the curb design and management should be important in future due to the possible high frequency of vehicle pickup/drop-off services.

## Keywords

Shared mobility; Endogenous congestion; Autonomous vehicle routing; Alternating Direction Method of Multipliers (ADMM); Dantzig-Wolfe decomposition; Column pool;

## 1. Introduction

The transportation sector is experiencing an unprecedented revolution with the emerging advanced sensing, telecommunications and vehicular technologies, which are generating a new wave of rich information and providing a great opportunity to better control and optimize transportation system operations. On the other hand, it results in great challenges to estimate and predict their impacts on existing congested roadway infrastructure, new mobility modes, and future transportation system design. For example, it has been a hot debate whether the new ride-hailing service from Transportation Network Companies (TNCs) is adding to congestion in numerous downtowns. Officials in San Francisco, Chicago and New York have cited congestion as the main rationale for new fees they recently enacted on Lyft and Uber rides in each of the cities (Brown, 2020). One main challenge to evaluate this kind of impacts is that any changes in a complex interdependent system may invoke a series of hardly predictable interactions of endogenous variables. There has been a number of studies in transportation area to consider the challenges caused by endogenous factors. For example, the bottleneck model (Vickrey, 1969) of congestion with endogenous scheduling has been used to address a number of challenging transportation economics problems, such as, time pattern of congestion, optimal pricing, unpriced equilibria (Small, 1982; Small and Verhoef, 2007; Small, 2015). In addition, Batarce and Ivaldi (2014) formulated a structural travel demand model with endogenous congestions, Chow and Recker (2016) proposed an estimation model with endogenous arrive time constraints for better calibrating the household activity pattern problem, de Almeida Correia and van Arem (2016) focused on the household optimum privately owned automate vehicle assignment with endogenous road congestions by using volume-delay functions, and Liu et al. (2018) considered the congestions based on a point queue traffic flow model for finding household activity patterns. Therefore, one fundamental question that transportation managers should address is how to capture the endogenous traffic patterns with those upcoming new and uncertain elements for future transportation planning and management.

The uncertain elements in future transportation systems spread over travel demand, vehicle supply and infrastructures, ranging from personal trip requirements (Martinez et al., 2015; Rayle et al., 2016; Davidson and Spinoulas, 2016; Martinez and Viegas e, 2017; Lim et al., 2018; Ma et al., 2018; Tong et al., 2019), vehicle driving modes (Chen et al., 2017; Wong et al., 2017; Nieuwenhuijsen et al., 2018), route behavior (Levin et al., 2017; Wong et al., 2017; Hyland and Mahmassani, 2018; Tong et al., 2019), vehicle ownership (Davidson and Spinoulas, 2016; Lavieri et al., 2017; Allahviranloo and Chow, 2019), to the update of infrastructure capacity due to sensor and communication advancement (Varaiya 1993; Papadimitratos et al., 2009; Qu et al., 2010; Gentili and Mirchandani, 2012; Mahmassani, 2016; Dey et al., 2016; Din et al., 2019).

Various approaches to respond to these are being studied, most notably the recent Intelligent Transportation Systems (ITS) initiatives of USDOT (e.g., the recent Dynamic Mobility Applications and Active Transportation Demand Management (DMA-ATDM); and Connected Vehicles Programs). In addition, the concept of Mobility as a service (MaaS) born in Finland and first launched at Hannover in Germany (Fred, 2016) has attracted much attention by combining transportation services from public and private transportation providers through a unified gateway that creates and manages the trip, which users can pay for with a single account (Heikkilä, 2014; Kamargianni et al., 2016; Hensher, 2017; Jittrapirom et al., 2017; Mulley, 2017; Djavadian and Chow, 2017; Bruun, 2018; Tong et al., 2019). As autonomous vehicles potentially become increasingly popular among ride providers, freight operators and perhaps some personal vehicles, the consideration of how these vehicles can be operated to meet temporally and spatially distributed traveler mobility needs with a limited road expansion budget and constrained infrastructure capacity is largely missing in studies of implementing large fleets of rideshared vehicles. In other words, it is crucial to study the integration of travel demand, vehicle supply and infrastructure capacity under one unified modeling framework.

The interplay of travel demand and vehicle supply is similar to traditional vehicle routing problems (Toth and Vigo, 2002) which are modeled in virtual point-to-point networks without physical roads. The link travel time is either constant or time-dependent based on externally observed historical, real-time or predicted traffic congestions (Taniguchi and Shimamoto, 2004; Kok et al., 2012). Meanwhile, the interaction between vehicle supply and physical networks is usually used to model its internal road congestion as traditional dynamic traffic assignment problems (Peeta and Ziliaskopoulos, 2001), which treats vehicles equal to their carried passengers at origin zones. Given the development of emerging new technologies in information sharing and vehicle automation, it is more possible for passengers to share their trips' origin and destination with preferred time windows and for transportation control centers to guide and control vehicles in the future. Therefore, by focusing on one ideal system optimal (SO) scenario

for offline transportation planning from the viewpoint of society (Wardrop, 1952), this paper aims to consider demand, vehicle and infrastructure simultaneously to capture the endogenous traffic pattern with a number of assumptions. More specifically, (i) All vehicles are autonomous or can be systematically guided as human-driven vehicles by one management center. As a note, the difference between autonomous vehicles or human-driven vehicles in this paper can be reflected by the road capacity and backward wave speed due to the different vehicle reaction times and minimal vehicle following distance (e.g. Wei et al. 2017), or the vehicle carrying capacity for ridesharing capability. (ii) Vehicles have carrying capacity and depart from their origin depots to destination depots with specific working time windows. (iii) All passengers submit their trip requests with pickup and/or drop-off locations and time windows in advance and accept the ridesharing service. (iv) The road congestion is endogenously incurred by those guided vehicles. As a note, the differences between autonomous vehicles and human-driven vehicles can be reflected by the road capacity and backward wave speed due to their different vehicle reaction time and minimal vehicle-following distance and the vehicle carrying capacity for ridesharing capability. In addition, from the perspective of traditional Wardrop's first principle, each user aims to find his/her best route. In our solution, the final traffic condition can also be viewed as a kind of user equilibrium with accepted tolerance, because each user also finds his/her required route, but this equilibrium may not be unique. If passengers just have their pickup time windows and do not have a tight limit on the drop-off time, we can only state that our model is for system optimal rather than a kind of user equilibrium.

About the setting of vehicles' origin and destination, there are two popular ways to address it. One way is to predefine the origin and destination by assigning one origin zone (depot) and one destination zone (depot) for each vehicle (such as, Scherr et al., 2019). The other way is to randomly generate the origin of each vehicle and its destination is determined by the location of its last served trip. Our model chooses the first way, so all vehicles could be managed and maintained at different depots at the end. In addition, if we define the household as the origin and destination depot of its owned autonomous vehicles, our model is still applicable. Furthermore, if we do not want all vehicles to finally return to the destination depot, we can treat this depot as one virtual depot and use one virtual link to connect this virtual depot with all drop-off locations in the network, so the drop-off location of the last trip will be the final parking location and the link travel cost from all the drop-off locations to the virtual depot will be assumed to be 0. In addition, we believe that the new technologies associated with shared mobility and autonomous vehicles will revoke more research about how and where the depots should be designed and how to better operate those vehicles in the future, which are beyond the scope of this paper.

## 1.1 Challenges

The modeling approach for our problem with route coordination, ridesharing services and constrained road capacity is similar to the integration of vehicle routing problem (VPR) and dynamic traffic assignment (DTA), so this section aims to summarize the challenges of each separated problems and their integrations at first.

Without considering the endogenous road congestions, the first set of problems in our research context can be simplified as the vehicle routing problem with pickup and delivery with time windows (VRPPDTW), which has been proved to be NP-hard (Baldacci, et al., 2011). The difficulty of this problem arises from the complex categories of constraints, (i) vehicle flow balance, (ii) the logic of passenger pickup and drop-off by the same vehicle within the required time windows, and (iii) dynamic vehicle carrying capacity of the ridesharing choice. In particular, it is sometimes challenging to even find a feasible solution due to the complicated interaction of all constraints. Recently, Psaraftis et al. (2016) summarized the research of the last three decades and offered a systematic classification of dynamic vehicle routing problem according to 11 criteria.

Focusing on the road congestion incurred by those assigned autonomous vehicles, the tight link capacity limitation at each time point could greatly make a large number of side constraints. If the queue spillback and kinematic waves (Newell, 1993; Daganzo, 1994) are further considered, the complex interaction among vehicles makes the problem more challenging. A detailed discussion about the connection between different traffic flow models can be found in the paper by (Zhang et al., 2013). Even in traditional dynamic traffic assignment models without considering passenger pickup and drop-off requests, it is still difficult to calculate the path marginal cost in congested networks to reach the system optimal goals (Ghali and Simth, 1995; Peeta and Mahmassani, 1995; Shen et al., 2007; Qian et al., 2012; Lu et al., 2016), especially when there are overlapped paths in large scale networks. In addition, Kalifates (2010) proposed a graph theoretic modeling framework with cell transmission model for generalized transportation systems to reduce the problem complexity. The current mathematically tractable solutions (Arnott et al., 1990; Yang and Huang, 2005;

Munoz and Laval, 2006) mainly apply in parallel networks with a single bottleneck originally studied in the paper (Vickrey, 1969).

From the perspective of practice, simulation approaches are usually selected in dynamic traffic assignment problems to capture the road congestion with queue spillback and first-in-first-out (FIFO) rule. However, those approaches cannot explicitly handle the personalized user requests without optimization techniques. Therefore, the integration of simulation and optimization would be the trend to serve the future urban mobility systems with connected autonomous vehicles and ridesharing services.

Many researchers have been working on the autonomous control of vehicles (Reece and Shafer, 1993) and automated intelligent vehicle/highway system design (Varaiya, 1993) for increasing the system safety, efficiency and reliability by using simulated environment and optimization techniques (Hanebutte et al., 1998; Van Arem et al., 2006; Talebpour and Mahmassani, 2016; Chen et al., 2017; Sun et al., 2017; Ghiasi et al., 2017; Ye and Yamamoto, 2018; Stern et al., 2018). On the other hand, there is recently a large number of studies that focuses on the impacts of shared-use mobility on future transportation systems with autonomous vehicles.

There is a number of simulation models that has been developed for system modeling and analytics. Behrisch et al. (2011) developed an open-source traffic simulation package (SUMO) for the simulation of urban mobility with automated driving and flexible traffic management strategy evaluations. Fagnant and Kockelman (2014) developed an agent-based simulation model in a grid-based urban area where some strategies are provided to match passengers with vehicles and relocate vehicles to reduce traveler waiting time, but the endogenous road congestion and vehicle carrying capacity for ridesharing are not considered. Martinez et al. (2015) proposed an agent-based simulation model for a shared-taxi system and Martinez and Viegas (2017) further improved the model to incorporate taxi-bus system and consider the road congestion based on time-dependent flow capacity ratio and free-flow speed. Levin et al. (2017) proposed a modeling framework to (i) capture the traffic congestion by simulation-based network loading based on the updated flow-density diagram with autonomous vehicles (Levin and Boyles, 2016) and (ii) serve the ride-sharing services by some heuristic algorithms. Maciejewski and Bischoff (2016) analyzed the impact of autonomous taxi on traffic congestion based on their dynamic vehicle routing problem under an agent-based simulation environment. Hyland and Mahmassani (2018) focused on the on-demand shared-use autonomous vehicle mobility services (SAMS) without shared rides by proposing six vehicle-to-passenger assignment strategies tested in an agent-based simulation tool.

A number of analytical solutions also have been provided based on different assumptions in the future mobility system. By assuming that all automated vehicles are privately owned by each household, de Almeida Correia and van Arem (2016) first proposed a mathematical model to minimize the household-level generalized travel cost by combining the vehicle routing problem and dynamic traffic assignment in congestion networks in which a volume-delay travel time function is adopted to capture the road congestions, and they further focus on the user (a whole household) equilibrium solved under the framework of Method of Successive Average (MSA). By following the modeling framework above, Liang et al. (2018) modified the nonlinear volume-delay travel time function by introducing discrete congestion levels to make the formulation as an integer linear programming model, and then applied rolling horizon to address the real-time trip assignment and dynamic routing for automated taxis with congestions, and Van Essen and Correia (2019) also tried to replace the nonlinear travel time calculation by making multiple link travel time choices at each time point of the link entry node, and then the study selects only one to represent congestion time in the space-time network for solving system optimum and user equilibrium problems with approximations. Alonso-Mora et al. (2017) studied the real-time ride-sharing problem with high-capacity vehicles and large number of trips by dynamically generating the optimal route for the online demand and available vehicles with high-quality solutions, but as the traditional vehicle routing problems, the road congestion is not embedded in the models. Ma et al. (2017) proposed a linear programming model to assign available autonomous vehicles to satisfy those trip requests by constructing its feasible service network in advance, but the ride-sharing option and road congestion are not considered either. In congested networks, Rossi et al. (2018) studied the autonomous vehicle routing and rebalancing, and Salazar et al. (2018) considered how to best assign travelers between autonomous vehicles in traffic systems and public transit vehicles, and the road congestion is also simplified by flow-based travel cost function, which is typically used for long-term transportation planning rather than short-term traffic operations, so it could affect the accuracy of estimated congestions as well. Focusing on the household activity pattern problem, Liu et al. (2018) first formulated the endogenous road congestion by a point queue traffic flow model, which is caused by household

owned automated vehicles that are assigned to perform different mandatory or optional household daily activities. [Tong et al. \(2019\)](#) offered a modeling framework in an open-source simulation engine (DTALite-S) to incorporate agent-based simulation and optimization in a multimodal transportation environment with different trip requests to capture complex traffic dynamics. [Di and Ban \(2019\)](#) explored the general static traffic equilibrium of new shared mobility systems with driving solo, ridesharing, and e-hailing service in which the road congestion is represented by one volume-delay travel time function and the ridesharing in e-hailing service is not respected. Recently, [Mourad et al. \(2019\)](#) provided a survey on models and algorithms of shared mobility systems and it can be observed that the endogenous dynamic congestion is still not well considered in current optimization models.

## 1.2 Problem Decomposition approaches

With the development of computer hardware, the computation capabilities for solving mathematical programming models are evolving very quickly. However, many large-scale problems still lead to formulation that greatly goes beyond the computation limit. One usual way is to find the special blocks in formulation to directly decompose models as relatively solvable subproblems connected by coupling constraints. The decomposition can be grouped into two categories, primal decomposition and dual decomposition ([Boyd et al., 2007](#); [Palomar and Chiang, 2006](#)), where the dual price in primal decomposition or the Lagrangian multiplier in dual decomposition for the coupling constraints is used to update and control subproblems. In the primal decomposition, column generation and Dantzig-Wolfe decomposition are widely used for linear programming and mixed integer programming with branch-and-price methods ([Barnhart et al., 1998](#)). Lagrangian relaxation/decomposition ([Fisher, 1981](#); [Mahmoudi and Zhou, 2016](#); [Wu et al., 2019](#)) is usually used for integer programming from the dual perspective. [Huisman et al. \(2005\)](#) summarized that the dual price in linear programming (LP) relaxed restricted master problem in column generation can be replaced by the Lagrangian multiplier in its LP relaxed dual problem through Lagrangian relaxation without using branch and price. It should be noted that finding a feasible initial solution in primal decomposition or obtaining a good feasible final solution in dual decomposition are also not straightforward in large-scale complicated problems, especially having different categories of side constraints. Also, how to address the non-unique dual prices or Lagrangian multipliers is important for the quality of solutions in the iterative process. In addition, from the perspective of primal-dual algorithm, the alternating direction method of multipliers (ADMM) ([Boyd, 2011](#)) can be used to decompose the overall problem as a number of sequentially connected subproblems, which are controlled by Lagrangian multipliers. Therefore, ADMM is also viewed as an efficient way to break the symmetry issues compared with Lagrangian relaxation. Recently, to address the stochastic mixed-integer programming models, [Boland et al. \(2018\)](#) applied ADMM to this problem and use the Frank-Wolfe method based on simplicial decomposition to deal with the nonlinear objective functions of subproblems, and then showed that their approach is theoretically supported, computationally efficient, and parallelizable.

To solve the traditional static traffic assignment problem ([Wardrop, 1952](#); [Beckmann et al., 1956](#)), [LeBlanc et al. \(1975\)](#) offered a linearization algorithm to solve the classical model based on the Frank-Wolfe algorithm ([Frank and Wolfe, 1956](#)). An important improvement of Frank-Wolfe algorithm is simplicial decomposition, which is a special version of the Dantzig-Wolfe decomposition principle based on Carathéodory's theorem, where extreme points are usually generated by the solution of the linear Frank-Wolfe. [Larsson and Patriksson \(1992\)](#) proposed a disaggregate simplicial decomposition (DSD) which treats each path of OD pairs as one extreme point rather than the network flow solution in simplicial decomposition (SD). Then [Larsson et al. \(2004\)](#) focused on the side constrained traffic equilibrium problem, which is solved by column generation based on their DSD approach. Moving forward to dynamic traffic assignment, a similar disaggregated simplicial decomposition is also used for gap-function-based user equilibrium ([Lu et al., 2009](#)) and eco-system optimum ([Lu et al., 2016](#)) with different traffic flow models. In addition, based on the cell transmission model, Dantzig-Wolfe decomposition was also used to solve system optimal ([Li et al., 2003](#)) and user optimal ([Lin et al., 2010](#)). Focusing on vehicle routing problems with ride-sharing services, Lagrangian relaxation is also used to decompose the problem as a number of shortest path finding problems ([Mahmoudi and Zhou, 2016](#); [Liu et al., 2018](#)). Recently, [Yao et al. \(2018\)](#) applied ADMM to solve the vehicle routing problem with drop-off requests only in the context of urban logistics to show the solution performance from the primal and dual aspects.

## 1.3 Potential contributions and structure of this paper

As stated by [Gendreau et al. \(2016\)](#), spatial and temporal behavior of traffic variations should be analyzed, but still is an enormous challenge requiring consolidating knowledge from various disciplines (traffic flow theory, statistics,



etc.) for vehicle routing problems with stochastic travel time (VRPSTT). This kind of challenges mainly arises from the endogenous congestions among moving vehicles in capacitated transportation systems. Focusing on the specific scenario stated before, the contributions of our research are listed as follows.

(1) Compared with the literature about vehicle routing problem and dynamic traffic assignment, this paper takes a further step to integrate travel demand, vehicle supply and infrastructure supply to explicitly capture the new traffic condition. Specifically, it aims to optimally assign vehicles from different depots to satisfy individuals' temporally and spatially distributed mobility requests under constrained road capacity and queue spillbacks.

(2) Due to the complexity of this problem, a primal decomposition approach, Dantzig-Wolfe decomposition, is used to decompose the proposed space-time-state (STS) arc-based vehicle-based integer linear programming model as a restricted master problem and a number of subproblems. The subproblems can be independently solved for each vehicle by time-dependent state-dependent shortest path algorithms. To our best knowledge, this is the first study that uses a three-dimensional STS path as an extreme point in Dantzig-Wolfe decomposition, compared with previous research using physical paths or space-time paths.

(3) To solve the large-scale problems and to serve as an approximation from the perspective of dynamic traffic assignment, an STS path-based flow-based linear programming model is proposed by building a column pool in advance. ADMM is then applied to decompose this model as a number of sequential quadratic programming subproblems solved by projected gradient method for each column. Specifically, each column represents a space-time-state path with the vehicle-to-passenger assignment and vehicle-to-arc assignment information through solving the primal arc-based vehicle-based model in a sampling dataset by ADMM, which can also decompose the primal problem as a number of sequential time-dependent state-dependent shortest path problems for each vehicle.

The remainder of this paper is organized in the following manner. Section 2 formally state our focused problem and conceptually illustrates our modeling approach in a space-time-state network. Section 3 provides an arc-based vehicle-based integer linear programming model (**Model 1**) and an STS path-based flow-based linear programming model (**Model 2**), which are decomposed and solved by Dantzig-Wolfe decomposition and column-pool-based approximation approach, respectively, in Section 4. Numerical experiments are performed to demonstrate our proposed methodology and algorithms in Section 5. Finally, our future research is discussed in Section 6.

## 2. Problem statement and illustrative example

Consider a physical transportation network with a set of nodes  $N$  and a set of links  $L$ . Each link can be denoted as a directed link  $(i, j)$  from upstream node  $i$  to downstream node  $j$  with a given free-flow link travel time  $TT_{i,j}$  and link capacity  $Cap_{i,j}$ . Each vehicle  $a$  has an origin depot  $o_a$  and a destination depot  $d_a$  with a specific departure time window  $[l_a, m_a]$  and an arrival time window  $[l'_a, m'_a]$ . The number of seats in vehicle  $a$  is  $Cap_a$  and is also named vehicle carrying capacity. In addition, each passenger will submit his/her trip requests with origin  $o_p$ , destination  $d_p$ , departure time window  $[l_p, m_p]$  and arrival time window  $[l'_p, m'_p]$ . Our problem aims to optimally assign each vehicle to meet those passengers' requests while considering the road capacity constraint.

As shown in Fig. 1(a), assume that 2 travelers plan to go to office (node 4) from home (node 2) and 1 traveler wants to go shopping (node 5) from home (node 3). They all have a specific departure time window and arrival time window. There is a number of available autonomous vehicles at different depots waiting to be dispatched to serve those time-dependent travel requests. Since the vehicle fleet size is probably large enough to incur traffic congestion, the physical traffic network with specific road capacities should not be neglected. For the modeling needs, a pick-up virtual node and a drop-off virtual node will be added at each passenger's pick-up and drop-off locations, respectively, as shown in Fig. 1(b). As a result, the passengers' served status and vehicle carrying state can have changes only if the vehicle visits those virtual nodes, which will be explained in detail later when constructing the space-time-state network.

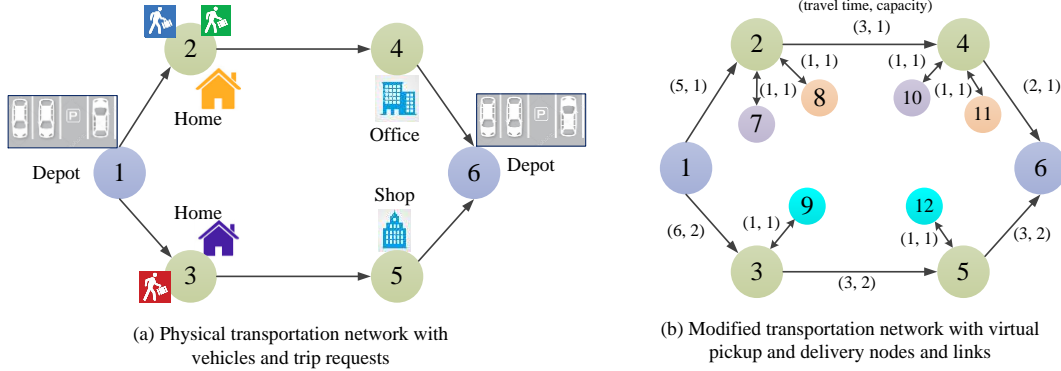


Fig. 1 The physical and modified transportation networks

In order to take into account the time dimension, a space-time network is employed to explain how to model the time window and road capacity at first, and then we will focus on a space-time-state network for modeling the whole process of our problem. Taking the physical path  $1 \rightarrow 2 \rightarrow 4 \rightarrow 6$  in Fig.1 (b) as an example, its corresponding space-time network is built in Fig. 2. Each node  $i$  is extended as a set of vertices  $(i, t)$  at each time interval and each link  $(i, j)$  is extended as a set of arcs  $(i, j, t, s)$  from vertex  $(i, t)$  to vertex  $(j, s)$ . The arc capacity is derived based on the hourly link capacity and the number of intervals in one hour. In addition, the arc  $(i, i, t, t + 1)$  from vertex  $(i, t)$  to vertex  $(i, t + 1)$  means that vehicles can wait at node  $i$  at time  $t$  for one time interval in case there is not enough capacity in the downstream arcs to accommodate them. The capacity of the waiting arcs is infinite, so the queuing process will be similar to the point queue traffic flow model. As shown in Fig. 2, passengers 1 and 2 request the same pickup and delivery time windows. Assume the carrying capacity of vehicles is always 1.

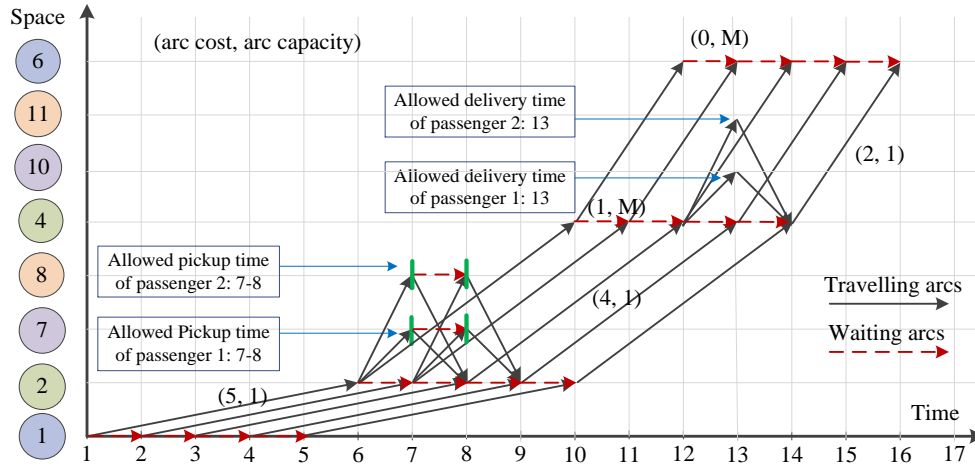


Fig. 2 The physical and modified transportation networks

We consider the following two cases without/with road capacity.

Case 1: the road capacity is not considered. Then it needs 2 vehicles departing at time 1 to satisfy those trip requests.

Case 2: the road capacity is strictly constrained. Then one vehicle has to wait at the depot until time 2 to depart. Moreover, one passenger cannot be served, because the vehicle waiting at depot can pick up one passenger but cannot deliver him/her at the allowed time due to tight road capacity constraints.

Therefore, it shows the difference and difficulty of finding the best vehicle assignment and routing solutions under tight physical facility limitations, compared with the traditional vehicle routing problems. This example in the space-time network only illustrates the general process of vehicle routing under road capacity constraints, but it does not consider (i) how to guarantee that once one passenger is picked up, he/she must be dropped off by the same vehicle nor reckon (ii) that the vehicle carrying capacity (the number of seats) cannot be violated due to the ridesharing options.

In order to model the process of passenger's pickup and delivery, the cumulative passenger served status is introduced and defined as follows, 0: the passenger is not served, 1: the passenger is being served (picked up but not



delivered), 2: the passenger is served (delivered). In addition to the dimensions of space and time, we introduce one more dimension  $w$  as vehicle carrying state to record which passenger is being served during the vehicle routing process. If passenger  $p$  is picked up by vehicle  $v$  with the carrying capacity of 1, the carrying state of vehicle  $v$  is  $p[1]$ , of which the first number is passenger number  $p$  and the number in square bracket mean the cumulative passenger served status of passengers. Still, focusing on the case in Fig. 2, if the vehicle capacity is 1, the possible vehicle carrying states include be  $()$ ,  $(1[1])$ ,  $(1[2])$ ,  $(2[1])$  or  $(2[2])$ . Similarly, if the vehicle capacity is 2, one possible vehicle carrying state example could be  $(1[1] \_2[1])$ , which represents that passengers 1 and 2 are currently picked up but not dropped off by the vehicle. This state-dependent approach (Mahmoudi and Zhou, 2016) can satisfy vehicle carrying capacity constraint and guarantee that one passenger can be picked up and dropped off by the same vehicle during the state transition process. In addition, if one vehicle cannot satisfy a passenger, it can still pass the passenger's pickup/drop-off physical nodes (such as, node 2) to serve other passengers, but it cannot pass his or her virtual arc and virtual node (such as, nodes 7 or 8). This is different from the traditional models in vehicle routing problems, which just have a network composed by pickup and drop-off locations without explicitly modeling physical transportation networks. Therefore, most VRP models do not allow vehicles to pass the pickup/drop-off node if that corresponding client is not served, based on one assumption that the arc cost matrix always satisfies the triangle inequality (such as, cost  $c_{i,j} \leq c_{i,k} + c_{k,j}$ ). Actually, if the triangle inequality doesn't hold in the real-world transportation networks based on the real travel cost, it may not make sense to not allow vehicles to pass the node without serving the passenger.

Finally, we can construct a three-dimension space-time-state network for vehicle routing, where each vertex is  $(i, t, w)$  and each arc is  $(i, j, t, s, w, w')$  from vertex  $(i, t, w)$  to vertex  $(j, s, w')$ . The vehicle carrying state transition process (state  $w$  to state  $w'$ ) is highly connected with the space (location) and time. Specifically, the vehicle carrying state will change when the vehicle picks up or drop off one passenger while ensuring that the vehicle carrying capacity is not violated. Note that once one passenger is served with a cumulative served status as 2, the passenger is not allowed to be served again, so there is no circle being selected in the state transition graph. The connection among state, space and time is the foundation of our constructed space-time-state networks. Section 4 will illustrate how to dynamically build this three-dimension network to find the time-dependent state-dependent shortest path for each vehicle. Actually, it is also possible to build the whole three-dimension network in advance based on the relation of space, time, and state, but the complexity will be explored in large-scale networks. Fig. 3 shows one feasible STS vehicle trajectory along the physical path  $1 \rightarrow 2 \rightarrow 4 \rightarrow 6$ . This STS path contains the vehicle-to-passenger mapping information when vehicle carrying states are changed at the pickup and drop-off locations within the required time windows, and also has the vehicle-to-arc mapping information in the space-time dimension. This kind of mapping information among vehicle, passenger and arc will be used in sections 3.2 and 4.2 for flow-based models.

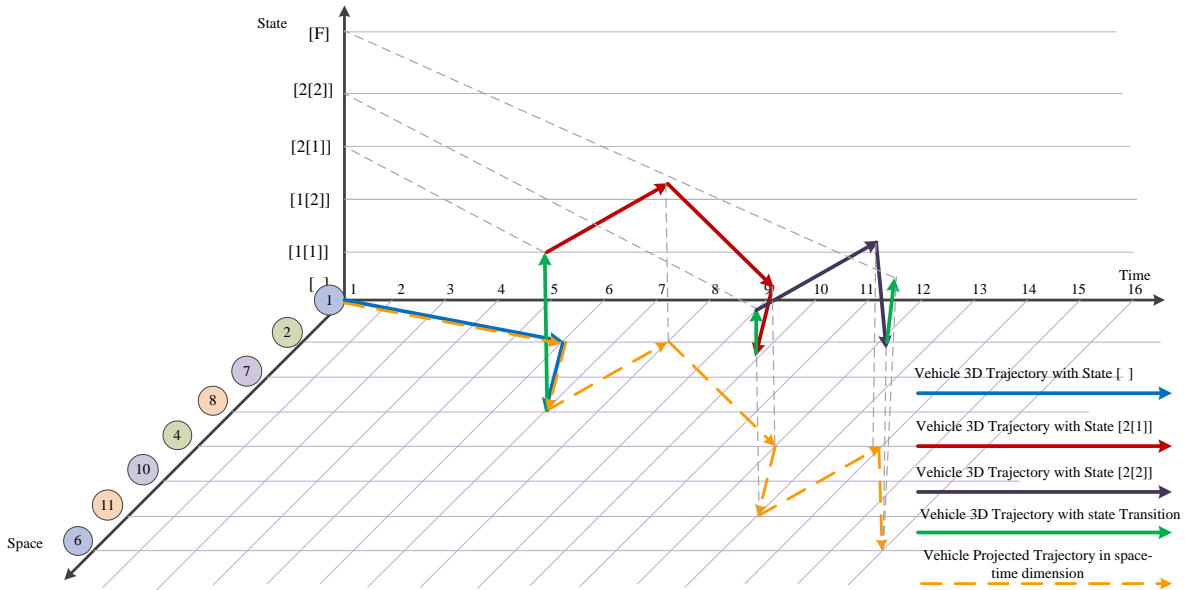


Fig. 3 A vehicle STS path with mapping information among vehicle, passenger and arc

As a remark, this modeling framework can be extended to the first/last mile problem where (i) passengers and vehicles have the same destination/origin and (ii) either pickup service or drop-off service is considered. The difference is only about the state definition. We just need 0 and 1 for passenger service status to indicate if he/she is served by the vehicle or not.

### 3. Mathematical Models

Table 1 lists the general indices, sets, parameters and variables used in our proposed arc-based vehicle-based model (**Model 1**) in Section 3.1, where each vehicle and each passenger are represented as vehicle  $a$  and passenger  $p$ , respectively. The notation of STS path-based flow-based model (**Model 2**) is listed in Table 2 where each vehicle group  $v$  has a number of vehicles with the same origin depot, departure time window, destination depot and arrival time window, and each passenger group  $q$  has a number of passengers with same trip requests (pickup location and time window and/or drop-off location and time window). Some parameters in Table 1 can also be used for the flow-based model.

Table 1 Indices, sets, parameters and variables in **Model 1**

Symbols	Definition
<b>Indices</b>	
$i, j$	Index of nodes, $i, j \in N$
$(i, j)$	Index of physical link between two adjacent nodes, $(i, j) \in L$
$a$	Index of vehicles
$p$	Index of passengers
$t, s$	Index of time intervals in the space-time network
$w, w'$	Index of vehicle carrying state
<b>Sets</b>	
$N$	Set of nodes in the physical traffic network
$L$	Set of links in the physical traffic network
$L_{inflow}$	Set of inflow links of single links in the physical traffic network
$L_{outflow}$	Set of outflow links of single links in the physical traffic network
$P$	Set of passengers
$A$	Set of vertices in the space-time-state network
$A_p$	Set of vertices of passenger $p$ 's pickup location
$A_m$	Set of vertices at the merge point in the space-time-state network
$E$	Set of edges/arcs in the space-time-state network
<b>Parameters</b>	
$o(a)$	Index of origin depot of vehicle $a$
$d(a)$	Index of destination depot of vehicle $a$
$DT(a)$	Earliest departure time of vehicle $a$ , equal to $l_a$
$[l_a, m_a]$	Departure time window of vehicle $a$ at the origin depot
$[l'_a, m'_a]$	Arrival time window of vehicle $a$ at the destination depot
$[l_p, m_p]$	Departure time window of passenger $p$ at the origin
$[l'_p, m'_p]$	Arrival time window of passenger $p$ at the destination
$VCap_a$	Carrying capacity of vehicle $a$
$Cap_{i,j,t,s}$	Capacity of arc $(i, j, t, s)$ in the space-time network
$h_a$	Vehicle flow of vehicle $a$ , is always equal to 1 for each vehicle
$c_{i,j,t,s,w,w'}^a$	Travel cost of arc $(i, j, t, s, w, w')$ of vehicle $a$
$\delta_{i,j,t,s,w,w'}^{a,k}$	Incidence between passenger pickup arc $(i, j, t, s, w, w')$ and path $k$ of vehicle $a = 1$ , matched; otherwise, $= 0$ .
$\delta_{k,p}^a$	Incidence between passenger $p$ and path $k$ of vehicle $a = 1$ , matched; otherwise, $= 0$ .
$\beta_{i,j,t,s}^{a,k}$	Incidence between arc $(i, j, t, s)$ and path $k$ of vehicle $a = 1$ , matched; otherwise, $= 0$ .
$FFTT_{i,j}$	Free-flow travel time of link $(i, j)$
$n_{i,j}$	The number of lanes on link $(i, j)$
$KJam_{i,j}$	Jam density of link $(i, j)$

<b>Variables</b>	
$x_{i,j,t,s,w,w'}^a$	Binary variable, vehicle $a$ will choose arc $(i, j, t, s, w, w')$ or not in the space-time-state network
$\lambda_k^a$	Binary variable, vehicle $a$ will choose path $k$ or not in the space-time-state network
$y_{j',j,t,t+1}$	Positive integer variable, the outflow arc capacity on arc $(j', j, t, t + 1)$
$\pi_p$	Lagrangian multiplier of passenger $p$ 's trip request
$\pi_{i,j,t,s}$	Lagrangian multiplier of capacity constraint of arc $(i, j, t, s)$

Table 2 Indices, sets, parameters and variables in **Model 2**

<b>Symbols</b>	<b>Definition</b>
<b>Indices</b>	
$v$	Index of vehicle groups
$q$	Index of passenger groups
<b>Parameters</b>	
$d(v)$	Total number of available vehicles of vehicle group $v$
$g(q)$	Total number of trip requests of passenger group $q$
$c_v^k$	Cost of path $k$ of vehicle group $v$
$\delta_q^{v,k}$	Incidence between passenger group $q$ and path $k$ of vehicle group $v = 1$ , matched; otherwise, $= 0$ .
$\delta_{i,j,t,s}^{v,k}$	Incidence between arc $(i, j, t, s)$ and path $k$ of vehicle group $v = 1$ , matched; otherwise, $= 0$ .
<b>Variables</b>	
$y_v^k$	Positive continuous variable, the number of vehicles belonging to group $v$ choosing path $k$ in the space-time-state network, finally simplified as $y_k$
$y_{j',j,t,t+1}$	Positive continuous variable, the outflow arc capacity on arc $(j', j, t, t + 1)$
$\lambda_q$	Lagrangian multiplier of passenger group $q$ 's trip requests
$\mu_{i,j,t,s}$	Lagrangian multiplier of capacity constraints of arc $(i, j, t, s)$

### 3.1 Arc-based vehicle-based integer linear programming model (Model 1)

Based on the space-time-state network constructed in Section 2, our mathematical programming model is proposed to minimize the total vehicle travel cost and satisfy passengers' trip requests and road capacity constraints.

Objective function:

$$\text{Min } Z = \sum_a \sum_{(i,j,t,s,w,w')} (c_{i,j,t,s,w,w'}^a \times x_{i,j,t,s,w,w'}^a) \quad (1)$$

Subject to,

(i) **Vehicle supply:** Arc-based flow balance constraint for each vehicle

$$\sum_{i,t,w:(i,j,t,s,w,w')} x_{i,j,t,s,w,w'}^a - \sum_{i,t,w:(j,i,s,t,w',w)} x_{j,i,s,t,w',w}^a = \begin{cases} -1 & j = O(a), s = DT(a), w = [0,0, \dots, 0] \\ 1 & j = D(a), s = T, w = [0,0, \dots, 0] \\ 0 & \text{otherwise} \end{cases}, \forall a \quad (2)$$

(ii) **Travel demand:** Passenger  $p$ 's pick-up request constraint

$$\sum_a \sum_{(i,t,w) \in A_p} x_{i,j,t,s,w,w'}^a = 1, \forall p \quad (3)$$

(iii) **Infrastructure supply:** Tight road capacity constraint (endogenous congestion)

$$\sum_a \sum_w x_{i,j,t,s,w,w'}^a \leq \text{cap}_{i,j,t,s}, \forall (i, j, t, s) \quad (4)$$

(iv) Binary definitional constraint

$$x_{i,j,t,s,w,w'}^a \in \{0,1\} \quad (5)$$

Constraint (2) ensures that each vehicle follows the flow balance. By constraint (3), each passenger will be picked up only once. For the problem with both pickup and drop-off requests, the state transit graph with cumulative passenger served status can guarantee that the passenger will be dropped off once he/she is picked up, so the drop-off constraint is always satisfied in our model. For the problem with pickup or drop-off requests only, we just need to ensure that the passenger pickup arc is only visited once by all vehicles. To capture the endogenous congestions, Constraint (4) forces the number of vehicles entering arc  $(i, j, t, s)$  to not exceed the arc capacity, which can be viewed as a point queue model where the vehicle has to choose the waiting arc if the capacity of the downstream link is not

available at the current time interval. The modeling on queue spillback is discussed in Appendix B. Variable  $x_{i,j,t,s,w,w'}^a$  is a binary variable, which indicates whether or not vehicle  $a$  will visit arc  $(i, j, t, s, w, w')$ . This proposed model is an integer linear programming model, so it can be solved directly by standard solvers. However, for the large-scale network applications, we will apply different decomposition approaches to decompose the problem as a number of relatively easy sub-problems in next sections. In addition, as introduced in section 1.1, de Almeida Correia and van Arem (2016) first proposed a mathematical model to minimize the household-level generalized travel cost by combining the vehicle routing problem and dynamic traffic assignment in congestion networks, so a comparison between their system optimum model and our Model 1 is listed in Appendix A.

As mentioned in the introduction, our depots can also be treated as virtual depots. By building the virtual links connected those destination depots with drop-off locations and making those virtual link cost as 0, all vehicles can park at the drop-off location of the last served passenger rather than forcing them to return to a physical depot. Therefore, our model approach is also applicable to address the real-time on-demand ridesharing problem when a real-time update scheme is incorporated. Actually, by improving the offline household-level nonlinear system optimal model (de Almeida Correia and van Arem, 2016), Liang et al. (2018) proposed a mixed integer programming model and applied it under the framework of rolling horizon algorithm to solve the real-time ridesharing problem while considering the endogenous road congestion, and their optimization model was finally solved by the commercial solver, Xpress.

### 3.2 Path-based flow-based linear programming model (Model 2)

The arc-based vehicle-based integer programming model in section 3.1 finds the best route guidance and captures the traffic condition in capacitated transportation networks. Admittedly, it is greatly challenging to solve this vehicle-based model in large-scale networks with a large number of vehicles and passenger requests. Looking into classical flow-based dynamic traffic assignment problems, all vehicles are assigned to the network based on each origin-destination (OD) pair as continuous flows rather than each individual vehicle, which could greatly reduce the number of variables to improve the computational efficiency. Therefore, from this perspective, if (i) vehicles and passengers can be grouped by its origin, destination and required service time period, and (ii) if all possible space-time-state path information with (a) vehicle (group)-to-passenger (group) and (b) vehicle (group)-to-arc assignment can be enumerated in advance for our overall problems, the remaining task is just to assign vehicles from each vehicle group to the network to satisfy the passenger group trip requests and to not violate the road capacity limitations. The linear programming model is listed as follows.

$$\min \sum_{(v,k)} (c_v^k \times y_v^k) \quad (6)$$

Subject to

(i) **Vehicle supply:** Path-based vehicle group flow balance constraint:

$$\sum_k y_v^k = d(v), \forall v \quad (7)$$

(ii) **Travel demand:** Pickup requests on passenger group  $q$ :

$$\sum_{(v,k)} (y_v^k \times \delta_{v,k}^q) = g(q), \forall q \quad (8)$$

(iii) **Infrastructure supply:** Road capacity constraints (endogenous congestion):

$$\sum_{(v,k)} (y_v^k \times \delta_{v,k}^{i,j,t,s}) \leq cap_{i,j,t,s}, \forall (i, j, t, s) \quad (9)$$

(iv) Positive continuous variable:

$$y_v^k \geq 0 \quad (10)$$

Since each STS path of each vehicle group is provided in advance, the cost of path  $k$  of vehicle group  $v$  for one specific OD pair is given. Constraint (7) ensures that the total number of vehicles in each vehicle group  $v$  is assigned to the network, which is consistent with constraint (2) for all vehicles in the vehicle-based model. Eq.(8) makes the total demand of passenger group  $q$  completely satisfied and also corresponds to Eq. (3) of passenger trip requests. Road capacity is considered in constraint (9), similar to constraint (4). Finally, we can generate a linear programming model, which is possible to be solved for large-scale transportation networks. Usually, it is difficult to enumerate all possible space-time-state paths, so a decomposition-based solution is helpful to generate a number of available paths, which can be viewed as the column pool construction used in Section 4.2.

## 4. Problem decomposition approaches

#### 4.1 Dantzig-Wolfe decomposition for arc-based vehicle-based formulation (Model 1)

In our proposed mathematical models in section 3.1, the flow balance constraint for each vehicle is a special block and can be solved independently. Therefore, Dantzig-Wolfe decomposition (Dantzig and Wolfe, 1960) is applied to solve our models and the flow balance constraints are used to develop the sub-problems. In addition, as mentioned by Larsson and Patriksson (1992), the generated paths (extreme points) from this decomposition approach are helpful for re-optimization if the demand or network has any updates in the future. Based on the point queue model to capture the road congestion, the primal model is decomposed as a master problem and different sub-problems for each vehicle as follows.

The master problem:

Objective function:

$$\text{Min } \sum_a \sum_{(i,j,t,s,w,w')} [c_{i,j,t,s,w,w'}^a \times \sum_k (\lambda_a^k \times h_a \times \delta_{i,j,t,s,w,w'}^{a,k})] \quad (11)$$

Subject to,

$$\sum_a \sum_{(i,j,t,s,w,w') \in A(p)} \sum_k (\lambda_a^k \times h_a \times \delta_{i,j,t,s,w,w'}^{a,k}) = 1, \forall p \quad (12)$$

$$\sum_a \sum_w \sum_k (\lambda_a^k \times h_a \times \beta_{i,j,t,s,w,w'}^{a,k}) \leq \text{cap}_{i,j,t,s}, \forall (i,j,t,s) \quad (13)$$

$$\sum_k \lambda_k^a = 1, \forall a \quad (14)$$

$$\lambda_k^a = \{0,1\} \quad (15)$$

The sub-problem for each vehicle  $a$ :

$$\text{Min } Z' = \sum_{(i,j,t,s,w,w')} (c_{i,j,t,s,w,w'}^a \times x_{i,j,t,s,w,w'}^a) - \sum_p \sum_{(i,j,t,s,w,w') \in A(p)} (\pi_p \times x_{i,j,t,s,w,w'}^a) - \sum_{(i,j,t,s)} (\mu_{i,j,t,s} \times \sum_w x_{i,j,t,s,w,w'}^a) - \omega_a \quad (16)$$

Subject to vehicle flow balance constraint.  $\pi_p$ ,  $\mu_{i,j,t,s}$  and  $\omega_a$  are the duals of side constraints (12), (13) and (14), respectively, and  $x_{i,j,t,s,w,w'}^a$  is a binary variable. It should be noted that the shortest path problem (our subproblem) is usually formulated by arc-based formulation, so we still use variable  $x_{i,j,t,s,w,w'}^a$  to write the model, but actually the output of the shortest path is still a path, which is composed of a number of passed arcs.

In the master problem,  $h_a$  is the vehicle flow of vehicle  $a$  and is always equal to 1 for each vehicle, so it can be removed.  $\lambda_k^a$  determines whether vehicle  $a$  will select path  $k$  or not. The sub-problem generates paths for each vehicle at each iteration, so it is convenient to use a path-based formulation for our master problem.

Objective function (11) can be reformulated as

$$\text{Min } \sum_a \sum_k (c_a^k \times \lambda_a^k) \quad (17)$$

Passenger pickup constraint (12) can be updated as

$$\sum_a \sum_k (\lambda_a^k \times \delta_{k,p}^a) = 1, \forall p \quad (18)$$

Space-time arc capacity constraint (13) is renewed as

$$\sum_a \sum_k (\lambda_a^k \times \beta_{(i,j,t,s)}^{a,k}) \leq \text{cap}_{i,j,t,s}, \forall (i,j,t,s) \quad (19)$$

The algorithm procedure is shown in Fig. 4.

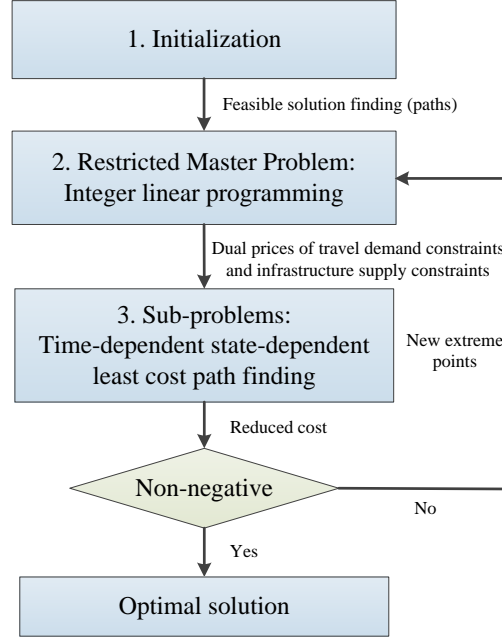


Fig. 4 The algorithm procedure of Dantzig-Wolfe decomposition

At step 1, the purpose is to find an initial feasible solution. Kalvelagen (2003) provided a mathematical model by adding virtual variables to find the feasible solution. Another way is to sequentially load each vehicle by using the solution from Lagrangian relaxation (Zhou et al., 2018). Specifically, once one vehicle finds its best route, the passengers served by this vehicle will be given a flag so that the following vehicles cannot visit those passengers anymore. Meanwhile, the space-time trajectory of that vehicle is recorded, so the capacity of visited space-time arcs will be updated by reducing 1. Once there is no available capacity on that arc, its arc travel time will be infinity. The pseudo-code is shown in Fig. 5.

At step 2, the restricted master problem is an integer linear programming model, which can be solved by a standard solver or a hybrid of solvers and branch-and bound, and then the model provides the dual prices of side constraints for the sub-problems.

At step 3, the sub-problem is a time-dependent state-dependent shortest path problem and offers new extreme points (paths) for each vehicle for the master problem at step 2, if all reduced cost is not non-negative; otherwise, the optimal solution is found. Note that path cost of  $c_a^k$ , passenger-vehicle assignment incidence of  $\delta_{k,p}^a$  and path-arc incidence of  $\beta_{k,(i,j,t,s)}^a$  are obtained once a new path (path  $k$ ) is generated for vehicle  $a$ .

```

1  //initialization: passenger status, relation between vehicle and passenger
2  for passenger p=1 to total_number_of_passengers
3      number_of_visits_of_passenger [p] = 0 // passenger p is not served
4  end// passenger
5  for vehicle a=1 to total_number_of_vehicles
6      for passenger p=1 to total_number_of_passengers
7          vehicle_passenger_visit_allowed_flag[a][p] = 1 // a is allowed to serve p
8      end //passenger
9  end// vehicle
10 // sequentially loading each vehicle to find its own least cost path
11 for vehicle a=1 to total_number_of_vehicles
12     run the beam search algorithm to find the best route for vehicle a and number_of_visits_of_passenger [p] = 1 if
13     passenger p is served
14     if (a < total_number_of_vehicles)
15         for p = 1 to total_number_of_passengers
16             if (number_of_visits_of_passenger [p] = 1)
17                 // the following vehicles cannot visit passenger p
18                 vehicle_passenger_visit_allowed_flag[a + 1][p] = 0
19             end
20         end // passenger
21     end
22 end
  
```



```

19     end
20     // update current available road capacity after loading vehicle a to obtain the visited link and time sequence
21     for link_no = 1 to total_number_of_visited_links_of_vehicle a
22         visit_time = visit_time_sequence[link_no] of vehicle a
23         link_capacity[link_no][visit_time] = max(0, link_capacity[link_no][visit_time] - 1)
24         if (link_capacity[link_no][visit_time] = 0) // no available capacity
25             link_time_dependent_travel_time[link_no][visit_time] = infinity // arc cost is infinity
26         end
27     end // link
28 end // vehicle

```

Fig. 5 Pseudo-code of proposed feasible solution finding algorithm

A beam search algorithm is proposed as an approximate dynamic programming approach to dynamically construct the space-time-state network and find the optimal routing with the least generalized route cost for each vehicle. As an improved version of the previous beam-search algorithm (Zhou et al., 2018), we add one more loop on each node so that more possible vehicle states will be considered during the beam search process shown in Fig. 6. The key part is to update vehicles' states based on the state transition rule considering the sequence and time windows of passenger pickup and delivery and vehicles' carrying capacity. In addition, this algorithm is also applicable to solve the problems with pickup or drop-off only when the vehicle's state definition is correspondingly changed based on this new problem.

```

1 //definition: vehicle:  $v$ , node:  $n$ , time:  $t$ , state:  $w$ , vehicle location-dependent time-dependent states:
   $td\_state[v][n][t][w]$ 
2 for  $t$  = departure time to ending time  $T$ 
3   for  $n$  = 0 to total_number_of_nodes  $N$ 
4     //beam-search: find the best  $k$  vehicle states with least travel costs from depot to current node and time
5     state_size = min{ $k$ , state size of vehicle  $v$  at node  $n$  and time  $t$ }
6     for  $w$  = 0 to state_size
7       Current_node =  $n$ 
8       for to_node = 1 to the outbound_node_size of current_node
9         if (to_node is passenger pickup or drop-off node)
10           Update the vehicle state  $td\_state[v][n][t][w]$  with passenger pickup or drop-off, current_node, current_time, travel cost from the depot to current node and time with benefits of serving passengers, based on previous node  $n$ , previous time  $t$  and link travel time, previous state  $w$ , and the whole state transition logic.
11         if (to_node is physical network node)
12           Update the vehicle state  $td\_state[v][n][t][w]$  with current_node, current_time and current travel cost, and state  $w$  doesn't change.
13         if (to_node is destination node of vehicle  $v$ )
14           Update the vehicle state  $td\_state[v][n][t][w]$  and update the corresponding Vector vehicle_ending_state  $[v]$ , which will be used to find the least cost route for vehicle  $v$  after all loops.
15         end // downstream node of one link
16       end // states
17     end // nodes
18   end // times

```

Fig. 6 Pseudo-code of our proposed beam search algorithm for each vehicle  $v$

To improve the computation efficiency, we also consider a tree-based path representation embedded in the beam search algorithms, which can provide the one-to-all shortest path tree during the search process. Therefore, when one vehicle finds its three-dimensional shortest path, the shortest path tree from its origin vertex (with its depot node and departure time) to all vertex is also found. Then if another vehicle also departs from that origin vertex but to a different destination depot, it does not need to run the beam-search algorithm again. Instead, we can directly find its shortest path in the one-to-all shortest path tree, which will reduce the total computation time.

#### 4.2 Column-pool-based approximation for STS path-based flow-based formulation (Model 2)

In reality, it is impossible to enumerate all possible STS paths to obtain the connection of vehicle-to-passenger and vehicle-to-arc relationship in advance, as assumed in section 3.2 for Model 2. Therefore, we propose one approximation approach, which is similar to scale factors but with some key differences.

To our knowledge, some traffic simulation tools have tried to use scale factors to consider city-wide applications in order to reduce the computer memory use and computation efficiency. In our approach, we aim to study a relatively small focused system (e.g. an AV test bed, a subarea such as city downtown) with more complex traffic conditions. By doing so, we hope to fully examine the relation between each vehicle group and all passenger groups, and clearly calculate the mutual impact between each vehicle group and space-time paths from a system with sampled passengers, vehicles and reduced link capacities.

A space-time-state path of one vehicle with served passengers and visited space-time arcs is called one column in the column pool. The more columns we can build in advance, the more candidates we can choose to satisfy passengers' requests and road capacity by vehicle flow assignment. Different vehicles from the same group could serve different passengers from different passenger groups and also visit different arcs, so we can have more relations among vehicle groups, passenger groups, and space-time arcs. Then, based on the real-world passenger requests and road capacity, we will determine how to assign available vehicles with vehicle group ID to serve those passengers with passenger group ID under arc capacity constraints to reach the minimal system travel cost. In other words, we need to do the vehicle assignment under the passenger request constraint (8) and road capacity (9) based on the obtained relations, rather than multiplying the unique scale factor.

As shown in Fig. 7, the complex primal problem is divided into two stages. Stage A focuses on a sampled traffic system to build a column pool by generating a number of columns (space-time-state paths) for each vehicle group by solving the arc-based vehicle-based formulation using ADMM. The other advantage of column pool generation is for re-optimization in case the demand, vehicle or network has any changes in the future, so we can use those available columns as a starting point instead of performing the optimization model from the beginning. Stage B is to solve the STS path-based flow-based linear programming model by ADMM to assign vehicles from different vehicle groups to serve passengers from different passenger groups while satisfying the time-dependent road capacity constraints.

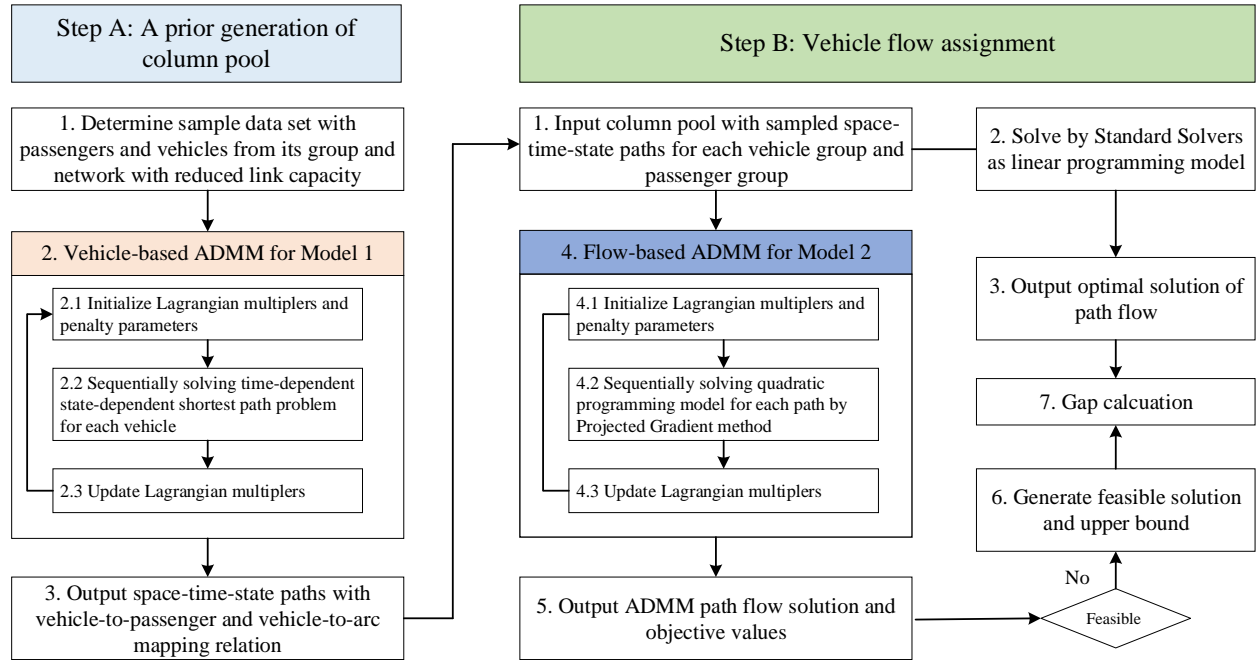


Fig. 7 The flow chart of column-pool-based approximation

#### 4.2.1 A priori generation of column pool based on Model 1

As mentioned above for addressing large-scale networks, the possible relation between different vehicle groups, passenger groups and space-time arcs needs to be obtained, so we select a number of passengers and vehicles from their groups as a sample set. Also, the arc capacity is reduced correspondingly to produce the possible congestions. If

the generated column pool is too small, it is possible to have an infeasible solution or a sub-optimal solution. Therefore, the key is to find the balance among column number, computation time, and solution quality. How to determine this sample set and how to dynamically manage the column pool (Barnhart et al., 1998) will be studied in our future research.

Then an arc-based vehicle-based model in section 3.1 is used to generate the column (space-time-state path) for each vehicle. ADMM is selected to decompose this problem as a number of sequential time-dependent state-dependent shortest path problem without using any standard solvers. In addition, it is also possible to try other heuristic methods to generate different columns, such as, changing the benefits of serving different passengers, local search algorithms, and so on.

The objective function of ADMM:

$$\begin{aligned} \text{Min } Z = L(\mathbf{x}^a, \boldsymbol{\pi}_p, \boldsymbol{\pi}_{(i,j,t,s)}) = & \sum_a \sum_{(i,j,t,s,w,w')} (c_{i,j,t,s,w,w'} \times x_{i,j,t,s,w,w'}^a + \sum_p [\pi_p \times (\sum_a \sum_{(i,j,t,s,w,w') \in A(p)} (x_{i,j,t,s,w,w'}^a \times \\ & \delta_{i,j,t,s}^a - 1)) + \frac{\rho_1}{2} \sum_p [\sum_a \sum_{(i,j,t,s,w,w') \in A(p)} (x_{i,j,t,s,w,w'}^a \times \delta_{i,j,t,s}^a - 1)]^2 + \sum_{(i,j,t,s)} [\pi_{(i,j,t,s)} \times (\sum_a \sum_w x_{i,j,t,s,w,w'}^a - \text{cap}_{i,j,t,s})] + \\ & \frac{\rho_2}{2} \sum_{(i,j,t,s)} [\sum_a \sum_w x_{i,j,t,s,w,w'}^a - \text{cap}_{i,j,t,s}]^2 \end{aligned} \quad (20)$$

Subject to the flow-balance constraint for each vehicle at constraint (2).  $x^2 = x$  if  $x = \{0,1\}$ , so the quadratic terms in objective function can be converted to being linear with binary variables, when the problem is solved for each vehicle sequentially based on the standard procedure of ADMM (Boyd, 2011). For the illustration purpose, we assume that there is a model with objective function  $(x_1 + x_2 - a)^2$  and two binary variables  $x_1$  and  $x_2$ . When solving  $x_1$ , we can have  $(x_1 + x_2 - a)^2 = [x_1 + (x_2 - a)]^2 = x_1^2 + 2x_1(x_2 - a) + (x_2 - a)^2 = (2x_2 - 2a + 1)x_1 + (x_2 - a)^2$  where  $x_2$  is fixed and  $a$  is a parameter.

The penalty parameters of  $\rho_1$  and  $\rho_2$  for passenger service constraint and arc capacity constraints are given in this paper, but they can also be updated based on some rules used in previous augmented Lagrangian relaxation models. The iterative scheme of ADMM is shown in Fig.8.

// initialization

Set up initial values for all Lagrangian multipliers and penalty parameters

for  $n = 1$  to  $n_{\max}$  // total number of iterations

for  $a = 1$  to  $a_{\max}$  //total number of vehicles

Find the time-dependent state-dependent shortest path for vehicle  $a$  with the fixed solution of other vehicles

Update the network-arc costs based on the new solution of vehicle  $a$  for vehicle  $a + 1$

end // vehicle

Update Lagrangian multipliers of passenger pickup constraints and arc capacity constraints

end // iterations

Fig. 8 The iterative scheme of ADMM

At iteration  $n + 1$  of ADMM:

$$x_a^{n+1} = \arg \min \{L(x_1^{n+1}, x_2^{n+1}, \dots, x_a, x_{a+1}^n, \dots, x_{a_{\max}}^n, \boldsymbol{\pi}_p^n, \boldsymbol{\pi}_{i,j,t,s}^n)\} \quad (21)$$

$$\pi_p^{n+1} = \pi_p^n - \rho_1 [\sum_a \sum_{(i,j,t,s,w,w') \in A(p)} (x_{i,j,t,s,w,w'}^{a,n+1} \times \delta_{i,j,t,s}^a - 1)] \quad (22)$$

$$\pi_{i,j,t,s}^{n+1} = \max \{0, \pi_{i,j,t,s}^n - \rho_2 [\sum_a \sum_w x_{i,j,t,s,w,w'}^{a,n+1} - \text{cap}_{i,j,t,s}]\} \quad (23)$$

The subproblem for each vehicle is a time-dependent state-dependent shortest path problem due to the linear relation in the objective function. Once one vehicle finds its best solution, the network arc cost will be updated for the next vehicle's subproblem solving. All Lagrangian multipliers are updated at the end of each iteration.

#### 4.2.2 Dynamic vehicle flow assignment based on Model 2

Once the columns are generated for each vehicle group, the remaining task is to assign vehicles to satisfy passengers' trip requests and network capacities. Assume that the total number of vehicles from each vehicle group is also unknown, then we can apply ADMM to convert the flow-based linear programming model as a quadratic programming model as follows.

Objective function:

$$\begin{aligned} \min \sum_k (c^k \times y^k) + \sum_q (\lambda_q \times [(\sum_k (y^k \times \delta_q^k) - g(q))] + \frac{\rho_1}{2} \sum_q ((\sum_k (y^k \times \delta_q^k) - g(q))^2 + \sum_{i,j,t,s} (\mu_{i,j,t,s} \times \\ [\sum_k (y^k \times \delta_{i,j,t,s}^k) - \text{cap}_{i,j,t,s}] + \frac{\rho_2}{2} \sum_{i,j,t,s} (\sum_k (y^k \times \delta_{i,j,t,s}^k) - \text{cap}_{i,j,t,s})^2 \end{aligned} \quad (24)$$

where  $y^k$  is the path flow of path  $k$ , and  $\rho_1, \rho_2$  are parameters. For simplicity,  $y_v^k$  for each OD pair  $(o, d, \tau)$  is replaced by  $y^k$  by resorting all path numbers.

Its Hessian Matrix in  $y$  can be derived as,

$$H = \begin{bmatrix} \rho_1 \sum_p \delta_p^1 + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^1 & \rho_1 \sum_p \delta_p^1 \delta_p^2 + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^1 \delta_{i,j,t,s}^2 & \dots & \rho_1 \sum_p \delta_p^1 \delta_p^{k_{\max}} + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^1 \delta_{i,j,t,s}^{k_{\max}} \\ \rho_1 \sum_p \delta_p^1 \delta_p^2 + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^1 \delta_{i,j,t,s}^2 & \rho_1 \sum_p \delta_p^2 + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^2 & \dots & \rho_1 \sum_p \delta_p^2 \delta_p^{k_{\max}} + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^2 \delta_{i,j,t,s}^{k_{\max}} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_1 \sum_p \delta_p^1 \delta_p^{k_{\max}} + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^1 \delta_{i,j,t,s}^{k_{\max}} & \rho_1 \sum_p \delta_p^2 \delta_p^{k_{\max}} + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^2 \delta_{i,j,t,s}^{k_{\max}} & \dots & \rho_1 \sum_p \delta_p^{k_{\max}} + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^{k_{\max}} \end{bmatrix}$$

Since it is difficult to directly obtain its inverse matrix  $H^{-1}$ , especially in large-scale networks, we apply ADMM to decompose the primal problem to sequentially solve the subproblem for each column as

$$y_k^{n+1} = \arg \min \{L(y_1^{n+1}, y_2^{n+1}, \dots, y_k, y_{k+1}^n, \dots, y_{k_{\max}}^n, \lambda_q^n, \mu_{i,j,t,s}^n)\}$$

The subproblem for  $y^k$  is a quadratic programming model which could be solved by the projected gradient method (Rosen, 1960) as follows:

$$y_k^{n+1} = \max \{0, y_k^n - \frac{1}{s} \times L(y_k^n)\}' \quad (25)$$

where  $L(y_k^n)' = c^k + \sum_q \lambda_q \times \delta_q^k + \rho_1 \left( \sum_q \delta_q^k \left( (\sum_k (y_k^n \times \delta_q^k) - g(q)) \right) \right) + \sum_{i,j,t,s} \mu_{i,j,t,s} \times \delta_{i,j,t,s}^k + \rho_2 (\sum_{i,j,t,s} \delta_{i,j,t,s}^k (\sum_k (y_k^n \times \delta_{i,j,t,s}^k) - cap_{i,j,t,s}))$ , and  $s = \frac{\partial^2 L(y_k^n)}{\partial y_k^2} = \rho_1 \sum_q \delta_q^k + \rho_2 \sum_{i,j,t,s} \delta_{i,j,t,s}^k$ . In addition, projected gradient method also has been used in solving the path-based nonlinear programming models in equilibrium traffic assignment (Larsson and Patriksson, 1992; Jayakrishnan et al, 1994; Florian et al., 2009), and it is more efficient, compared with arc-based nonlinear programming models, but it needs more memory use.

At each iteration of ADMM, the Lagrangian multipliers are updated as follows,

$$\text{Passenger group trip requests: } \lambda_q^{n+1} = \lambda_q^n + \rho_1 ((\sum_k (y_k^n \times \delta_q^k) - g(q)))$$

$$\text{Arc capacity constraints: } \mu_{i,j,t,s}^{n+1} = \max \{0, \mu_{i,j,t,s}^n + \rho_2 (\sum_k (y_k^n \times \delta_{i,j,t,s}^k) - cap_{i,j,t,s})\}$$

As a discussion, it is possible to assign different vehicles within different blocks, and each block can be sequentially solved in ADMM and vehicles within a same block can find the best solution with parallel computing techniques.

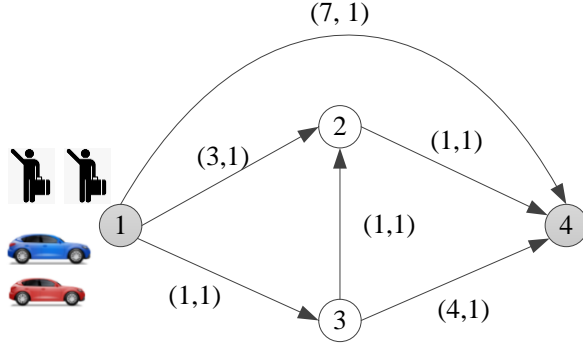
We need to note that the solution from ADMM cannot always guarantee its feasibility. In order to find a feasible solution and upper bound value, we can sequentially load each column flow from ADMM. Path flows that exceed the required passenger trip requests or arc capacity will be removed during the sequential loading process. Finally, if some passenger requests cannot be satisfied, virtual vehicles will be used to find a feasible solution as the upper bound.

## 5. Numerical examples

Section 5.1 shows how vehicles carrying passengers switch routes to reach the optimal solution in a capacitated network iteration by iteration using Dantzig-Wolfe decomposition. Section 5.2 focuses on both pickup and drop-off requests in the Sioux-Fall network. The restricted master problem is solved by CPLEX in GAMS, and the subproblems are solved by a beam search algorithm. In section 5.3, we consider the pickup only in the Chicago sketch network with a large number of autonomous vehicles and passengers belonging to different groups. A column-pool-based approximation approach is used to solve this problem. All corresponding source codes are shared online at <https://github.com/TonyLiu2015/AVRLite>.

### 5.1 A simple case in a capacitated network

Fig. 9 shows a simple capacitated network with 4 nodes, 6 links and 4 possible paths. Assume that there are 2 vehicles and each vehicle picks up one passenger from origin node 1 to destination node 4. Our goal is to minimize the total vehicle travel cost by Dantzig-Wolfe decomposition approach.



Path ID	Node Sequence	Path Cost	Path Trajectory
Path 1	1→2→4	4	
Path 2	1→3→4	5	
Path 3	1→4	7	
Path 4	1→3→2→4	3	

(link cost, link capacity)

Fig. 9 A simple capacitated network

Table 3 lists the details of each iteration where  $\lambda_k$  is the weight of selecting path  $k$  of from node 1 to node 4,  $\mu_{i,j}$  is the dual price of tight capacity constraint of link  $(i,j)$ , and  $\pi$  is the dual price of path weight constraint. Iterations 1 and 2 are used to generate a feasible solution by adopting the approach (Kalvelagen, 2003) in Dantzig-Wolfe decomposition. After 4 iterations, one vehicle chooses path 1 and the other selects route 2. The reduced cost is  $\sum_k (\lambda_k \times c_k \times 2) - \sum_{(i,j) \in L} \mu_{i,j} \times f_{i,j} - \pi_w$  where  $c_k$  is the cost of path  $k$  and  $f_{i,j}$  is the flow on link  $(i,j)$ . Finally, the reduced cost is  $4 + 5 - (-1) - (-2) - 12 = 0$  and reaches the optimal solution.

Table 3 The process of vehicle routing with endogenous congestions in Dantzig-Wolfe decomposition

Iteration NO.	Decomposed problem	Solution of different subproblems
Iteration 1	Subproblem	New column: <b>path 4</b>
	Restricted master problem	$\lambda_4 = 1, \mu_{1,3} = -1, \mu_{1,2} = 0, \mu_{3,2} = 0, \mu_{2,4} = 0, \mu_{3,4} = 0, \mu_{1,4} = 0, \pi = 2$
Iteration 2	Subproblem	New column: <b>path 3</b>
	Restricted master problem	$\lambda_4 = 0.5, \lambda_3 = 0.5, \mu_{2,4} = 0, \mu_{1,2} = 0, \mu_{3,2} = 0, \mu_{1,3} = -4, \mu_{3,4} = 0, \mu_{1,4} = 0, \pi = 15$
Iteration 3	Subproblem	New column: <b>path 1</b>
	Restricted master problem	$\lambda_4 = 0.5, \lambda_3 = 0.5, \mu_{1,3} = -1, \mu_{2,4} = -3, \mu_{1,2} = 0, \mu_{3,2} = 0, \mu_{3,4} = 0, \mu_{1,4} = 0, \pi = 14$
Iteration 4	Subproblem	New column: <b>path 2</b>
	Restricted master problem	$\lambda_1 = 0.5, \lambda_2 = 0.5, \mu_{1,3} = -1, \mu_{2,4} = -2, \mu_{1,2} = 0, \mu_{3,2} = 0, \mu_{3,4} = 0, \mu_{1,4} = 0, \pi = 12$

## 6.2 Trips with pickup and drop-off requests in the Sioux-Fall test network

As shown in Fig. 10, the Sioux-Fall network has 24 nodes, 84 links with hourly capacity, and 5 vehicle depots. We assume that there are 30 trip requests with specific pickup and drop-off location and time windows, which are not listed due the space limit of this paper. In addition, 30 candidate vehicles depart from different origin depots at different departure time to its corresponding destination depots to serve those trip requests. The optimization time horizon is 110 time units to cover those time windows and possible trip time. The generalized benefit of serving one trip request is -20 time units, and the waiting generalized cost of vehicles is half of its waiting time.

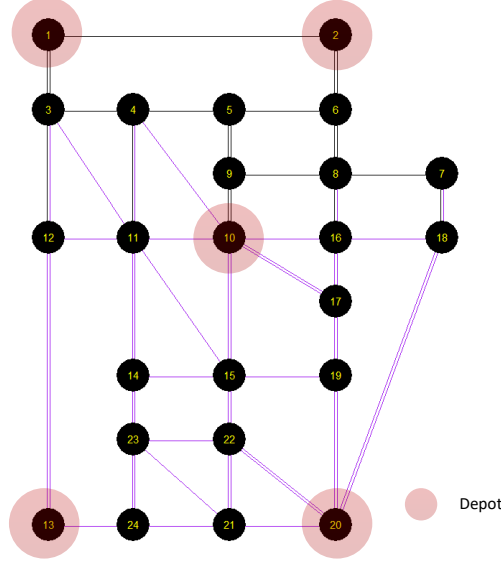


Fig. 10 Sioux-Fall network with five vehicle depots

Four scenarios are examined to compare the total generalized transportation cost and vehicle uses. Scenarios 1 and 2 have no endogenous road congestion due to the large road capacity. Then we change the road capacity with a very small value to incur congestions in Scenarios 3 and 4. In addition, in Scenarios 1 and 3, vehicle carrying capacity is 1, and in Scenarios 2 and 4, vehicle carrying capacity is increased to 2 to provide the ridesharing service.

For Scenarios 1 and 2, based on the algorithm for the initial feasible solution generation in section 4.1, 30 vehicles can serve 30 passengers in the physical network. Further, Dantzig-Wolfe decomposition is used to improve the initial solution. Similarly, the process is also implemented to scenarios 3 and 4 with road congestions. Table 4 lists the comparison on served passengers in each vehicle under different scenarios.

Table 4. Comparison on served passengers under different scenarios

veh_ID	no road congestion			with road congestion		
	initial_s	veh_cap=1	veh_cap=2	initial_s	veh_cap=1	veh_cap=2
1	[15]	[15][25]	[15][25]	[15]	[15][25]	[15][25]
2	[8]	[8]	[8]	[8]	[8]	[8]
3	[1]	[1]	[1]	[1]	[1]	[1]
4	[7]	[7]	[7][9]	[7]	[7]	[7]
5	[9]	[9]		[9]	[9]	[9]
6	[11]	[11]		[11]	[11]	[11]
7	[29]			[29]		
8	[28]	[28]	[28]	[28]	[28]	[28]
9	[17]	[17][29]	[17][29][30]	[17]	[17][29]	[17][29][30]
10	[21]		[21]	[21]	[21]	[21]
11	[20]			[20]	[20]	[20]
12	[26]	[26]	[26]	[26]	[26]	[26]
13	[16]	[16]	[16]	[16]	[16]	[16]
14	[18]	[18]	[18]	[18]	[18]	[18]
15	[2]	[2]	[2]	[2]	[2]	[2]
16	[10]	[10]	[10]	[10]	[10]	[10]
17	[3]	[3]	[3]	[12]	[12]	[12]
18	[12]	[12]	[12]	[22]	[22]	[22]
19	[27]	[27]	[27]	[27]	[27]	[27]
20	[30]	[30]		[30]	[30]	



21	[23]	[23]	[23]	[23]	[23]	[23]
22	[25]	[20][21]	[11][20]	[25]		
23	[22]	[22]	[22]	[19]	[19]	[19]
24	[13]	[13]	[13]	[13]	[13]	[13]
25	[4]	[4]	[4]	[3]	[3]	[3]
26	[5]	[5]	[5]	[5]	[5]	[5]
27	[19]	[19]	[19]	[24]	[24]	[24]
28	[24]	[24]	[24]	[14]	[14]	[14]
29	[14]	[14]	[14]	[4]	[4]	[4]
30	[6]	[6]	[6]	[6]	[6]	[6]

Focusing on vehicle 9, in the scenarios 1 and 2 without road congestion, it picks up and then drops off passenger 17 in the initial solution, but in the improved solution by Dantzig-Wolfe decomposition it serves passenger 17 and then continues to serve passenger 29. In addition, when the vehicle carrying capacity is 2, vehicle 9 serves passenger 17, and then picks up passenger 30 and goes to picks up passenger 29, and finally drops off passenger 29 first and then drops off passenger 30. In addition, the number of vehicles used under different scenarios in Table 4 is displayed in Fig. 11. As expected, the number of vehicles used under congestion conditions would be increased in order to serve passengers with specific pickup and drop-off windows. We also compare the total transportation cost in different scenarios in Fig. 12. It can be observed that the total system costs decrease when vehicle carrying capacity is increased with high ridesharing chances, and the road congestion would increase the system cost compared to the scenario without road capacity but with the same vehicle carrying capacity.

It should be emphasized that this paper focuses on an offline system optimal solution rather than the real-time vehicle scheduling problem, so the number of vehicles required in different scenarios could be different. In addition, pickup and drop-off windows can be a buffer time to mitigate traffic congestions and vehicles are allowed to wait at some points to further serve the next passengers. In other words, the required time windows from passengers are also important factors and inputs for system optimization.

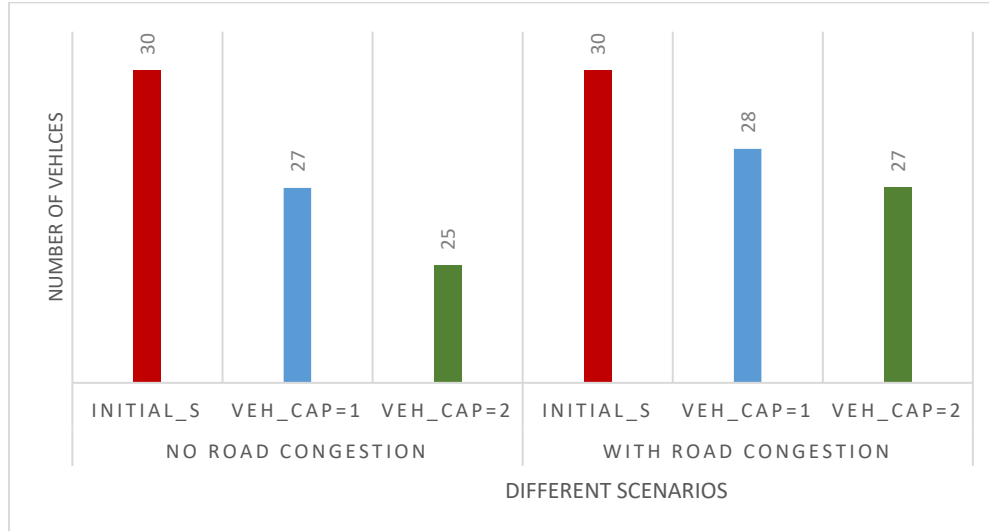


Fig. 11 Number of vehicles used in different scenarios

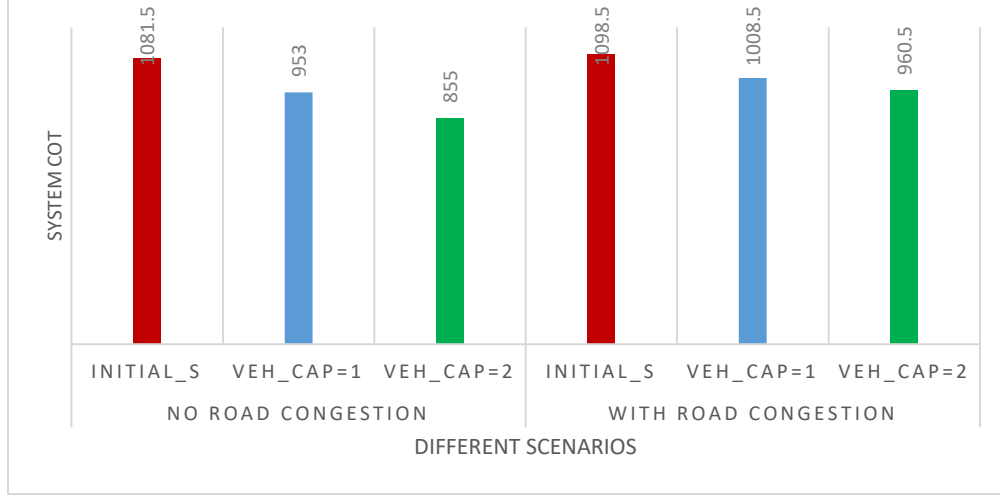


Fig. 12 Total generalized system travel cost in different scenarios

### 6.3 Trip with pickup only request in the Chicago sketch test network

#### Step A: A priori generation of column pool in Model 1

The Chicago sketch network has 1320 nodes and 5431 links in Fig. 13. We assume that all passengers have the pickup only trip requests as the first mile problem, which indicates that each passenger group with a number of passengers will have the same destination with a vehicle group. We treat them as one pair of vehicle group and passenger group. To generate the column pool, two scenarios are designed:

**Scenario 1:** 10 pairs of vehicle groups and passenger groups. In each pair, as a sample set, we assume that (i) 243 vehicles depart from different origins to one destination with different working time windows and (ii) 387 passengers submit trip requests with different pickup locations and time windows. The time horizon is 60 min (time unit). Since this is a simple set, the space-time arc capacity in each minute is assumed to be 5 vehicles. Vehicle carrying capacity is given as 1, so each vehicle aims to pick up one passenger from the origin depot to their same destination. It has 2430 binary variables and 332,160 constraints.

**Scenario 2:** 20 pair of vehicle groups and passenger groups. For each pair, similar to scenario 1, we also assume the same number of vehicles and passenger trip request but with different vehicle inputs (origin, destination, working time windows) and passenger inputs (pickup locations and time windows). It has 4860 binary variables and 338,460 constraints

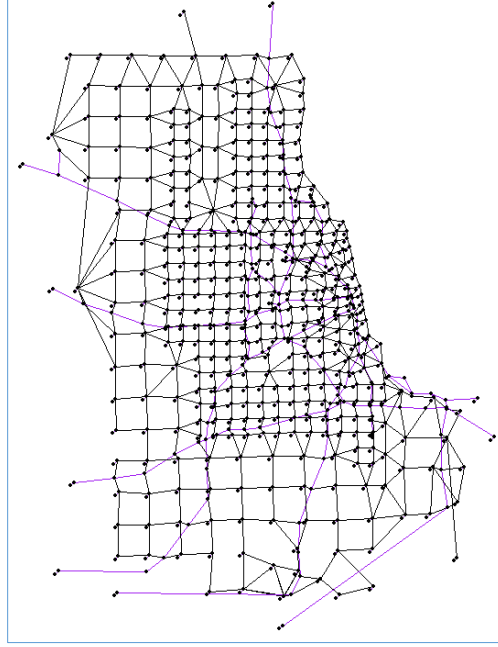


Fig. 13 Chicago sketch network for passenger pickup only

Then the vehicle-based ADMM in section 4.2.1 is used to find the vehicle-to-passenger and vehicle-to-arc assignment. Since the input data is randomly generated, some passengers may not be served and some vehicles may not serve any passengers. The final results are that (i) in scenario 1, 1789 vehicles find their paths/columns to serve 1084 passengers, and 23357 space-time arcs are generated based on vehicles' paths, and (ii) in scenario 2, 3686 vehicles find their paths/columns to serve 2226 passengers, and 36454 space-time arcs are generated based on vehicles' paths. The computation times for scenarios 1 and 2 are about 70 seconds and 140 seconds, respectively, from the laptop with 2.80GHz.

#### Step B: customized algorithm for flow-based ADMM

Note that each passenger has a specific pickup location, time window and destination, and vehicle can only pick up passengers within a group with the same destination location. Therefore, we can build a column pool where each path of vehicles represents one column and each passenger represents one passenger group from Model 1. The question that arises is how many vehicles from different vehicle groups are required to satisfy those trip request from different passenger groups under tight road capacity constraints. Based on the last two scenarios, we design two experiments:

**Experiment 1:** there are 1084 passenger groups and each passenger group has 4 passenger trip requests.

**Experiment 2:** there are 2226 passenger groups and each passenger group has 2 passenger trip requests. In this physical network, we assume that all space-time arc capacity in each minute is 35 vehicles, equal to 2100 vehicles per hour. To solve this problem, we try three approaches: (i) flow-based ADMM, (ii) upper bound generation by sequentially loading the column flow solution from ADMM, (iii) optimal solution from standard solver CPLEX in GAMS.

In **experiment 1**: three cases with different penalty parameters of  $\rho_1$  for passenger trip constraints and  $\rho_2$  for arc capacity constraints in ADMM are performed. Case 1:  $\rho_1 = 3$  and  $\rho_2 = 1$ ; Case 2:  $\rho_1 = 3$  and  $\rho_2 = 3$ ; Case 3:  $\rho_1 = 3$  and  $\rho_2 = 5$ . The results from ADMM by running 250 iterations and the optimal solution from CPLEX are shown in Fig. 14. The ADMM can converge to different objective values in three cases with different penalty parameters. Since capacity values and the number of capacity constraints are much higher than that of passenger trip requests, it is better to assign a smaller penalty value for  $\rho_2$  in the objective function of ADMM.

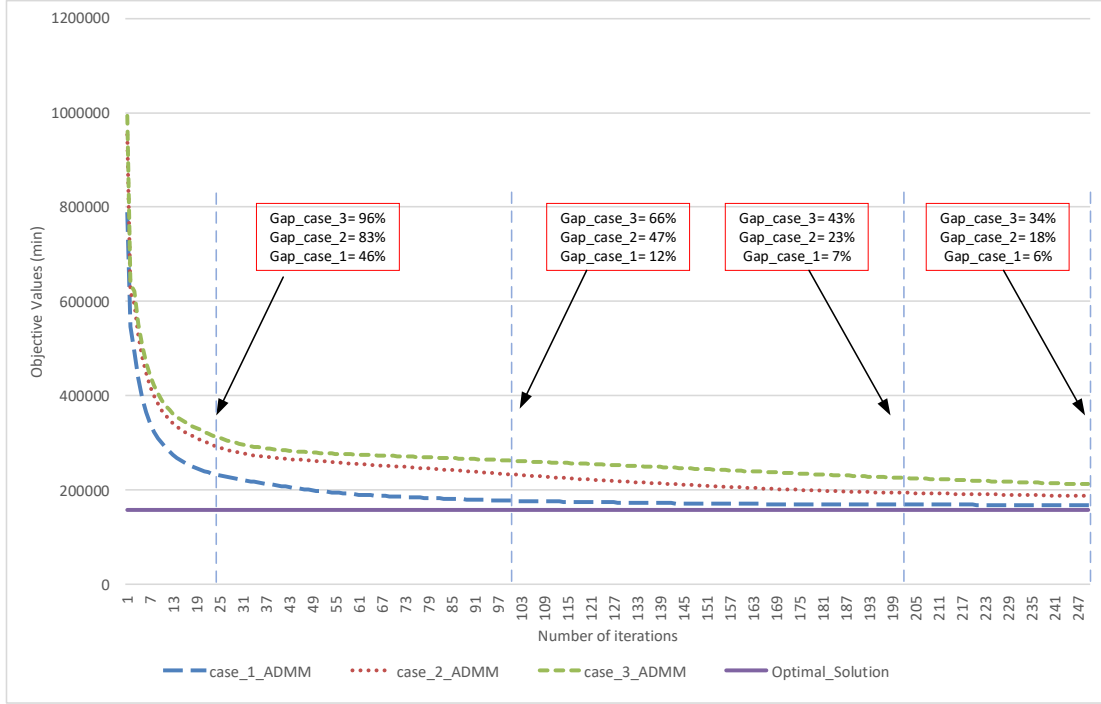


Fig. 14 Solution of each iteration of ADMM in three cases and CPLEX in experiment 1

Then the upper bound generation algorithm is also implemented to find a feasible solution based on the results of ADMM. Fig. 15 shows the objective values of upper bound in three cases and the optimal solution. The Gap values of three cases compared with the optimal solution are 4.3%, 3.4% and 3.1%, respectively. It is observed that the three cases can finally reach good solutions with very small gap values. Further, the ADMM result of case 3 has the biggest gap value, but its upper bound solution can still have a small gap value. The possible reason is that the total path flow is a variable, so the upper bound generation can reduce the total path flow from ADMM to have a better feasible solution to satisfy those passenger trip demands and not violate the arc capacity constraints.

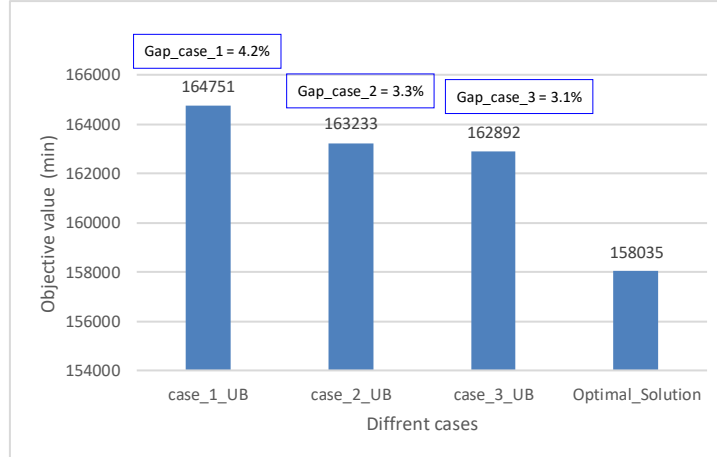


Fig. 15 Upper bound in three cases and the optimal value in CPLEX in Experiment 1

From the upper bound solution, 4737 space-time waiting arcs  $(i, i, t, t + 1)$  at 405 nodes have assigned vehicle flows, which indicates that the waiting happens at those nodes. By calculating the total waiting flow at those congested nodes during 60 mins, its heat map and the top10 of the most congested nodes are shown in Fig. 16(a) and (b), respectively. It can be observed that the destination areas of different passengers with pickup request only become congested, so it also raises one question about how to design the drop-off location in the future when a large number of passengers have the same destination with similar arrival time.

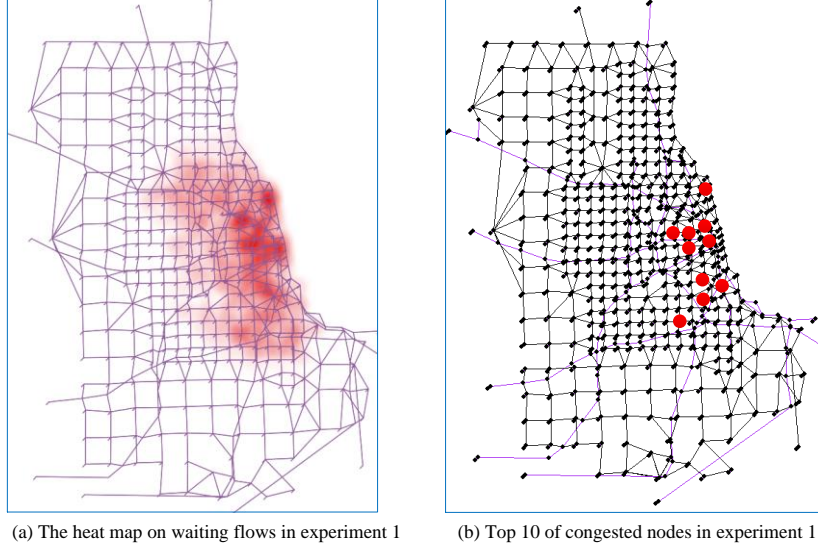


Fig. 16 Visualization of congested nodes in experiment 1

In order to compare the computation efficiency and memory use of the ADMM algorithm and CPLEX in GAMS (version 24.7.4), we also implement three other cases with different numbers of columns based on experiment 1. All tests are performed in the laptop with 8G memory and i5-4210U CPU @1.7GHz. The result is listed in Table 5. It can be observed that that CPLEX is much more efficient, but our customized algorithm can better utilize the memory. Specifically, our flow-based model is a linear programming (LP) model. How to efficiently solve LP model has been studied for almost 80 years, and most commercial solvers are powerful in solving LP models. Actually, most LP solvers apply some preconditioning techniques, such as, a simple geometric mean scaling in combination with equilibration, to reduce the condition number of the constraint matrix in order to decrease the computational efforts. As we know, CPLEX also has its “scaling routine” to preprocess LP models. In most of the literature in transportation domain, we haven’t found any papers which proposed an algorithm for LP and ultimately beat CPLEX in computation efficiency. Most comparisons happened when solving mixed integer programming, because CPLEX usually applied branch and cut to deal with those integer variables and the process is relatively time-consuming. In addition, it is still a research topic about how to improve the computation efficiency of ADMM, such as, by parallel computing (Boyd et al., 2010), which will be our future research. Also, we should emphasize that our ADMM is a general algorithm which can deal with nonlinear programming rather than just linear programming by CPLEX. Therefore, when the column cost is a nonlinear function of column flows, such as, for paths with reliability and variance (Xing and Zhou, 2011), our model could still be applicable and CPLEX will not deal with it. Another advantage of our column pool-based model is that a large number of columns has been stored in advance, so it will be beneficial for re-optimization and real-time optimization in the future and doesn’t need to read all the input data every time, which needs to be done in the case of CPLEX.

Table 5. Computation efficiency comparison between ADMM and CPLEX in GAMS

Num of columns	Flow-based ADMM_C++ (250 iterations)		GAMS (solver: CPLEX)	
	computation_time	memory_use	computation_time	memory_use
1789	12s	20m	0.8s	16m
17890	81.9s	43.4m	3.2s	77m
89450	344s	134.7m	15.2s	345m
178900	514s	264.1m	31.5s	686m

Then the same procedure is also applied to **experiment 2**. The solution process and comparison among different cases are shown in Fig. 17 and Fig. 18, respectively. The Gap values of the three cases compared with the optimal solution are 3.9%, 2.9% and 2.5%, respectively. From the upper bound solution, 7173 space-time waiting arcs  $(i, i, t, t + 1)$  at 448 nodes have assigned vehicle flows. Its heat map and the top10 of the most congested nodes are shown in Fig. 19(a) and (b), respectively.

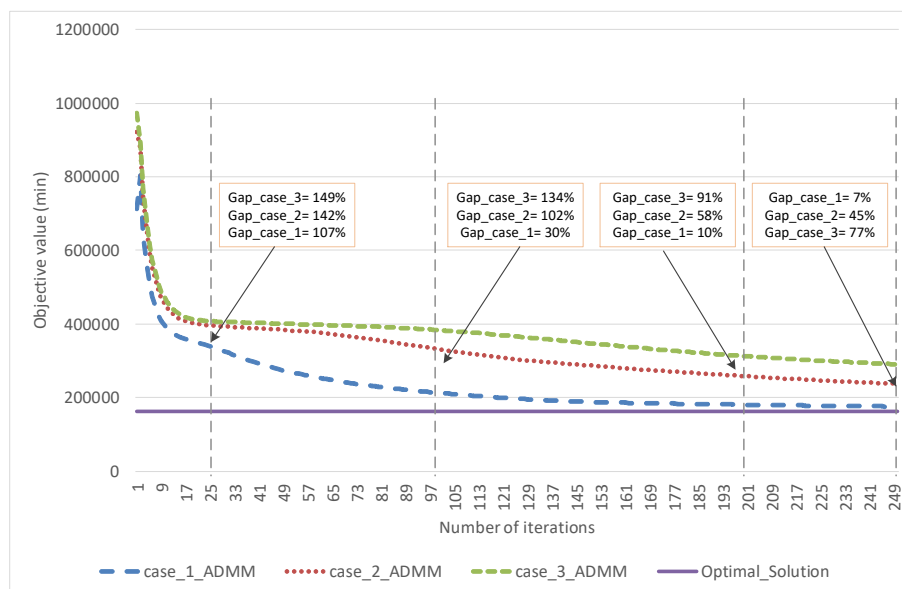


Fig. 17 Solution of each iteration of ADMM in three cases and CPLEX in experiment 2

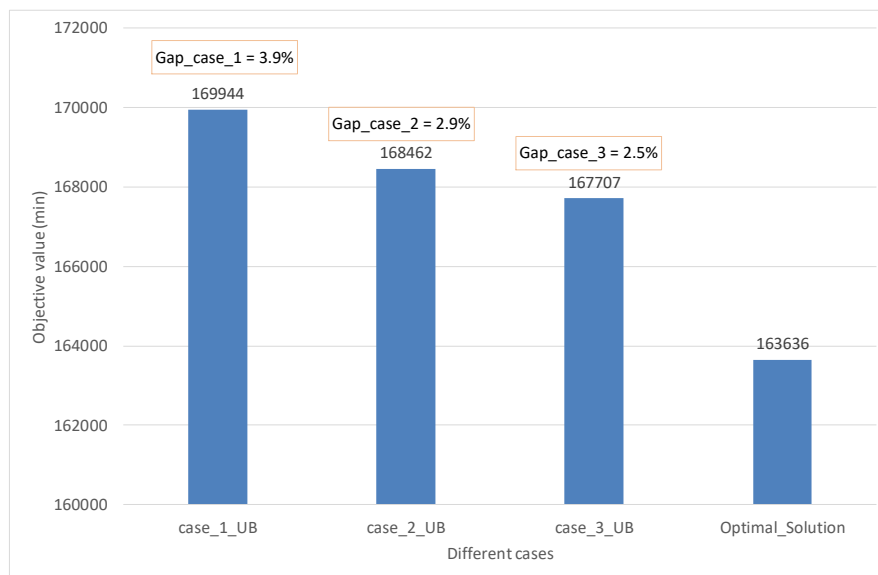
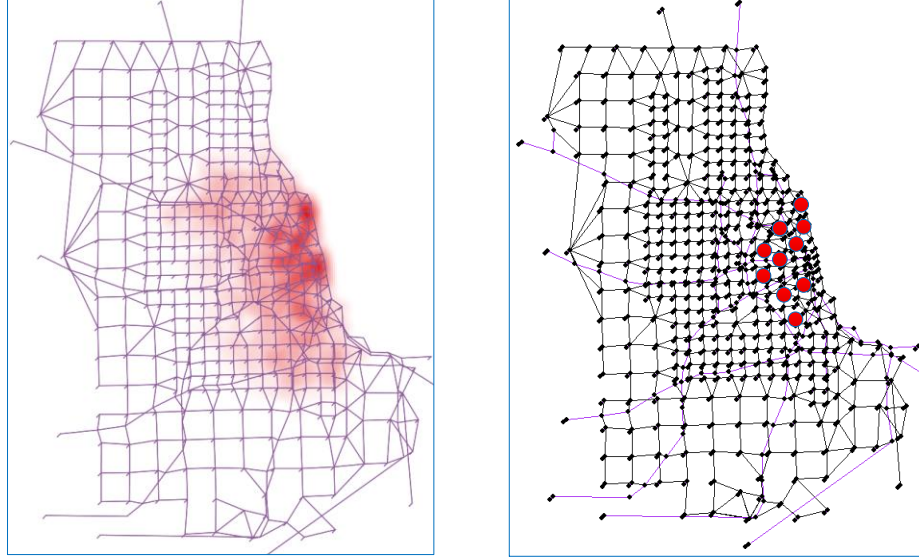


Fig. 18 Upper bound in three cases and the optimal value in CPLEX in Experiment 2





(a) The heat map on waiting flows in experiment 2

(b) Top 10 of congested nodes in experiment 2

Fig. 19 Visualization of congested nodes in two experiments

## 6. Conclusions and future research

This paper focuses on one ideal future scenario in which all vehicles can be centrally guided to pick up and drop off passengers within their required windows while considering endogenous congestion in capacitated networks. A vehicle-based arc-based integer programming model is proposed in our space-time-state networks for each vehicle and each passenger, which is solved by Dantzig-Wolfe decomposition. Then, based on the relations between vehicles and passengers and between vehicles and links from the vehicle-based model above, we further propose a flow-based path-based linear programming model from the perspective of dynamic traffic assignment. We then apply Alternating Direction Method of Multipliers (ADMM) to solve this linear programming model. The transportation system performance is highly related to different component layers in this system, including passengers' requests with pickup and drop-off location and time windows, vehicle carrying capacity and working rules, physical infrastructure capacities, etc. Therefore, any conclusions based on one specific input may not be universal.

From our preliminary experiments in this paper, we have a few interesting observations: (i) without considering the road congestion, the network performance/efficiency could be overestimated; (ii) passengers' required pickup and drop-off time windows could be a buffer to mitigate road congestion without impacting system performance; (iii) the ridesharing service could reduce the total transportation system cost under centralized control; (iv) the curb area design and management is important due to the possible high frequency of vehicle pickup/drop-off services. Therefore, in the future, with the increase of autonomous vehicle implementation, the resulting increased road capacity could enlarge the network capacity to better serve passengers, and the ridesharing travel modes should be encouraged to better utilize the carrying capacity of vehicles (car, bus, etc.)

Focusing on the algorithm part, we will improve the beam search algorithm (approximate dynamic programming) to address larger cases for finding the shortest path for pickup and drop-off requests. In addition, it is better to apply branch and bound to solve the relaxed problem in the master problem to obtain the dual prices in the Dantzig-Wolfe decomposition algorithm. Also, (i) how to determine the sample size and what algorithms can be used to generate the column pool and (ii) how to dynamically manage column pool are also important for the future application in large-scale networks. Since our proposed models are similar to dynamic system optimal traffic assignment models, it is also possible to calculate the path marginal cost to solve our model under the space-time-state framework. The implementation of ADMM with parallel computing should also be considered.

In a boarder sense, we will consider different kinds of trip requests and travel behaviors in a unified model to capture more complicated traffic conditions. [Liu and Zhou \(2016\)](#) showed how the tight road/vehicle capacity

constraints could invoke travelers' bounded rationality based on their day-to-day learning due to the inner system's uncertainty, which is incurred when identical travelers are competing for the limited resources without an assignment rule. This bound on trip cost is similar to (i) the concept of travel time budget in travel activity analysis, (ii) the changeable departure time of dynamic traffic assignment, and (iii) the required time windows for pickup and drop-off in vehicle routing problem. Therefore, it could be a suitable way to use time windows to model travel behavior in future congested multi-modal transportation systems. However, it is also challenging to estimate and predict those time windows from travel demand generation. In addition, the participants from different parties, such as, traffic regulator, mobility service providers, human drivers, and passengers, could lead to a more complex leader-follower game in the future research.

## Acknowledgements

This paper is supported by National Science Foundation – United States under Grant No. CMMI 1663657 “Collaborative Research: Real-time Management of Large Fleets of Self-Driving Vehicles Using Virtual Cyber Tracks”. The work presented in this paper remains the sole responsibility of the authors.

## Appendix A. Model Comparison

Comparison Category	Model A ( <i>System Optimum POAVAP</i> , de Almeida Correia and van Arem, 2016)	Model B (This paper)
(1) User behavior	Household-level optimum	Transportation system optimum
(2) Modeling approach	Time-discretized space-time network	Time-discretized space-time-state network
(3) Objective function	Minimizing total generalized transport cost: including travel distance, public transit cost, parking cost, penalty for early and late arrival. Formula (5).	Minimizing total generalized transport cost: travel time, waiting time, profit of serving passengers. Formula (1).
(4) Vehicle ownership	Private vehicles in each household. Vehicles don't have the departure /arrive time window, or the working hours.	Public vehicles managed by depots. Vehicles have departure/arrival time windows and its working hours.
(5) Vehicle flow balance	Constraint (23)	Constraint (2)
(6) Vehicle origin/destination	Vehicle origin is home, and its destination is the arrival node of the last served trip.	Vehicle origin is the origin depot, and its destination is the destination depot. The depot can be extended as home.
(7) Passenger/household member pickup/drop-off request satisfaction	Constraints (10), (12), (14)-(17).	Constraint (3). The logic of passenger pickup and drop-off process is satisfied in the vehicle carrying state transit process ( $w \rightarrow w'$ ), implemented in the time-dependent state-dependent shortest path finding in the beam-searching algorithm.
(8) Pickup/drop-off time windows	Each household member has desired pickup/drop-off time, earliest pickup time, and latest drop-off time.	Each passenger has a specific pickup/drop-off time windows embedded in the space-time-state networks.
(9) Vehicle carrying capacity	Constraints (19) and (20).	Satisfied in vehicle carrying state transit process.
(10) Road congestion and dynamic link travel time calculation	BPR travel time function, but needs to calibrate $t_{min}$ and $t_{max}$ and makes the model nonlinear. Constraints (24)-(27).	Point queue model by building waiting arc at each node in the space-time-state networks and having tight arc capacity constraint (4).
(11) Vehicles can pass passenger pickup/drop-off nodes without serving them.	Yes. If passengers/household members are not served, they will be assumed to be served by public transit. Constraint (8) (9), (11), (13), (18).	Yes. Though adding virtual pickup/drop-off nodes and links, vehicles can pass the real passenger pickup/drop-off physical nodes without serving them, but cannot pass those virtual service nodes/links without serving them.
(12) Empty vehicle elimination	Cannot eliminate an empty vehicle. Constraint (22).	Cannot eliminate an empty vehicle. However, when based on the solutions from Lagrangian Relaxation (Mahmoudi and Zhou, 2016) or ADMM, a sequential vehicle loading for upper bound generation can eliminate the remaining vehicles if all passengers have been served.

(13) Route with empty vehicles	Route with empty vehicle is not generated. Constraint (38).	Route with empty vehicle is generated to satisfy the flow balance constraint, but can be eliminated after obtaining the solution.
(14) Vehicles stop at any nodes	not allowed to stop idle when the vehicle is transporting a person. Constraint (21).	We have tight pickup/drop-off windows, so if vehicles arrive early, they need to wait to pick up/drop off the passenger until within the time windows. That is different with the desired pickup/drop-off time with early and late arrival penalty.
(15) First-in-First-out (FIFO)	Considered, but adding constraint (28) could force link travel time not follow volume-delay travel time calculation.	Not considered in the model. However, the sequential vehicle loading in the upper bound solution from Lagrangian Relaxation (Mahmoudi and Zhou, 2016) or ADMM can satisfy FIFO.

## Appendix B. Discussions on modeling queue spillback

In order to capture the queue spillback, we proposed a spatial queue model by improving the approach proposed by Drissi-Kaitouni and Hameda-Benchekroun (1992) where the link storage capacity with jam density and backward wave speed are not considered in space-time networks. Take the simple network in Fig.B1 (a) as an example. A virtual node as the waiting node is added for each link in the modified network in Fig.B1 (b). The travel time of link  $(2', 2)$  is assumed to be 1 time unit and its length is a small value as an approximation, so this link is used for discharging flows, and its capacity as the outflow capacity of link  $(1, 2)$  as a variable will be determined by its downstream links. The inflow capacity of link  $(1, 2)$  is the capacity of link  $(1, 2')$ , equal to  $Cap_{1,2}$ . The link storage capacity of link  $(1, 2)$  is  $Len_{1,2} \times n_{1,2} \times Jam_{1,2}$  and will be represented on link  $(1, 2')$ . The corresponding space-time network is constructed in Fig.B3 (c). Specifically, at time  $t$  on link  $(1, 2)$ , (i) the inflow capacity constraint is  $x_{1,2',t-FFTT_{1,2},t} \leq Cap_{1,2',t-FFTT_{1,2},t}$  for arc  $(1, 2', t - FFTT_{1,2}, t)$  shown in purple; (ii) the outflow capacity constraint is  $x_{2',2,t,t+1} \leq Cap_{2',2,t,t+1}$  for arc  $(2', 2, t, t+1)$  shown in orange; (iii) the link storage capacity constraint is  $CA_{1,2,t} - CD_{1,2,t} = x_{2',2',t-1,t} + \sum_{s=t-FFTT_{1,2},t+1}^{t-1} x_{1,2',s,t} \leq Len_{1,2'} \times n_{1,2'} \times KJam_{1,2'}$ .  $CA_{1,2,t}$  and  $CD_{1,2,t}$  are the cumulative arrival count and the cumulative departure count of link  $(1, 2)$  at time  $t$ . The link outflow capacity is calculated and given by the capacities of its downstream links.

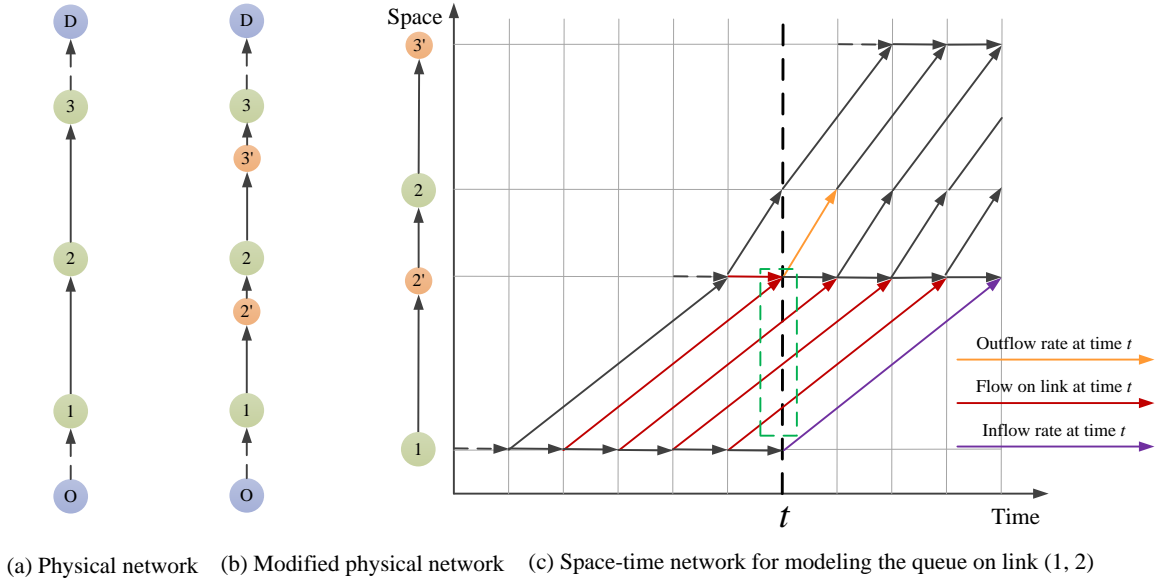


Fig. B1 Illustration for the spatial queue model in a space-time network

Focusing on the vehicle routing in our proposed space-time-state networks, the spatial queue model can be formulated as follows where  $x_{i,j,t,s,w,w'} = \sum_a x_{i,j,t,s,w,w'}^a$  in Model 1,

Inflow arc capacity constraint:

$$\sum_w x_{i,j',t-FFTT_{i,j+1,t,w,w'}} \leq Cap_{i,j',t-FFTT_{i,j+1,t}}, \forall (i,j') \in L_{inflow}, \forall t \quad (B.1)$$

Outflow arc capacity constraint:

$$\sum_w x_{j',j,t,t+1,w,w'} \leq y_{j',j,t,t+1}, \forall (j',j) \in L_{outflow}, \forall t \quad (B.2)$$

Outflow arc capacity balance constraint at points without merge and diverge:

$$y_{j',j,t,t+1} \leq Cap_{j,i,t+1,s} \quad (B.3)$$

Outflow arc capacity balance constraint at merger points:

$$\sum_{(j',t)} y_{j',j,t,t+1} \leq Cap_{j,i,t+1,s}, \forall (j,t+1) \in A_m \quad (B.4)$$

Outflow arc capacity balance constraint at diverge points:

$$y_{j',j,t,t+1} \leq \sum_{(i,s)} Cap_{j,i,t+1,s}, \forall (j,t+1) \in A_d \quad (B.5)$$

Link storage capacity constraint:

$$\sum_w x_{j',j',t-1,t,w,w'} + \sum_w \sum_{s=t-FFTT_{i,j}}^{t-1} x_{i,j',s,t,w,w'} \leq Len_{i,j'} \times n_{i,j'} \times KJam_{i,j'}, \forall (i,j') \in L_{inflow}, \forall t \quad (B.6)$$

Furthermore, to consider the backward wave speed under congested conditions, Newell's simplified kinematic wave model (Newell, 1993) considers the link storage capacity by  $CA_{(i,j,t)} - CD_{i,j,t-BWTT(i,j)} \leq Len_{i,j} \times n_{i,j} \times Jam_{i,j}$ . Similar to the derivation of the spatial queue model above, Newell's simplified kinematic wave model can have the following constrain for link storage capacity.

$$\sum_w \sum_{s=t-BWTT(i,j)}^t x_{j',j',s-1,s,w,w'} + \sum_w \sum_{s=t-FFTT_{i,j'}}^{t-1} x_{i,j',s,s+FFTT_{i,j'},w,w'} \leq Len_{i,j'} \times n_{i,j'} \times KJam_{i,j'}, \forall (i,j') \in L_{inflow}, \forall t \quad (B.7)$$

As a note, from the perspective of car-following models, the backward wave speed depends on drivers' average reaction time and minimal following safety distance, and Newell's simplified microscopic car-following model (Newell, 2002) is consistent with his macroscopic kinematic wave model. Wei et al. (2017) proposed a binary integer programming model to optimally control vehicle trajectory based on Newell's simplified car-following model in time-extended space-time networks, which can be incorporated in our modeling framework in state-space-time networks but will cause a huge number of variables and constraints.

## References

- Allahviranloo, M., Chow, J., 2019. A fractionally owned autonomous vehicle fleet sizing problem with time slot demand substitution effects. *Transportation Research Part C*, 98, 37-53
- Alonso-Mora, J., Samaranayake, S., Wallar, A., Frazzoli, E., Rus, D., 2017. On-demand high-capacity ride-sharing via dynamic trip-vehicle assignment. *Proceedings of the National Academy of Sciences*, 114(3), 462-467.
- Arnott, R., de Palma, A., Lindsey, R., 1990. Departure time and route choice for the morning commute. *Transportation Research Part B: Methodological*, 24(3), 209-228.
- Baldacci, R., Bartolini, E., Mingozzi, A., 2011. An exact algorithm for the pickup and delivery problem with time windows. *Operations research*, 59(2), 414-426.
- Beckmann, M.J., McGuire, C.B., Winsten, C.B., 1956. *Studies in the Economics of Transportation*. Yale University Press, New Haven, CT.
- Barnhart, C., Johnson, E.L., Nemhauser, G.L., Savelsbergh, M.W., Vance, P.H., 1998. Branch-and-price: Column generation for solving huge integer programs. *Operations research*, 46(3), 316-329.
- Batarce, M., Ivaldi, M., 2014. Urban travel demand model with endogenous congestion. *Transportation Research Part A: Policy and Practice*, 59, 331-345.
- Behrisch, M., Bieker, L., Erdmann, J. and Krajzewicz, D., 2011. SUMO—simulation of urban mobility: an overview. In *Proceedings of SIMUL 2011, The Third International Conference on Advances in System Simulation*. ThinkMind.
- Boland, N., Christiansen, J., Dandurand, B., Eberhard, A., Linderoth, J., Luedtke, J., Oliveira, F., 2018. Combining Progressive Hedging with a Frank--Wolfe Method to Compute Lagrangian Dual Bounds in Stochastic Mixed-Integer Programming. *SIAM Journal on Optimization*, 28(2), 1312-1336.
- Boyd, S., 2011, December. Alternating direction method of multipliers. In *Talk at NIPS workshop on optimization and machine learning*.
- Boyd, S., Parikh, N., Chu, E., Peleato, B. and Eckstein, J., 2011. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1), 1-122.

- 1 Boyd, S., Xiao, L., Mutapcic, A., Mattingley, J., 2007. Notes on decomposition methods. Notes for EE364B, Stanford
- 2 University, 1-36.
- 3 Brown, E., 2020. The ride-hail utopia that got stuck in traffic. The Wall Street Journal. Feb 15.
- 4 Bruun, E. C., 2018. A research agenda and proposed research approaches to enable MaaS to bring maximal benefits.
- 5 Paper presented at the 97th Transportation Research Board (TRB) Annual Meeting, Washington, D.C., United
- 6 States.
- 7 Chen, D., Ahn, S., Chitturi, M., Noyce, D.A., 2017. Towards vehicle automation: Roadway capacity formulation for
- 8 traffic mixed with regular and automated vehicles. *Transportation research part B: methodological*, 100, 196-221.
- 9 Chow, J.Y., Recker, W.W., 2012. Inverse optimization with endogenous arrival time constraints to calibrate the
- 10 household activity pattern problem. *Transportation Research Part B: Methodological*, 46(3), 463-479.
- 11 Daganzo, C.F., 1994. The cell transmission model: a simple dynamic representation of highway traffic. *Transp. Res.*
- 12 *Part B* 28 (4), 269–287.
- 13 Dantzig, G.B., Wolfe, P., 1960. Decomposition principle for linear programs. *Operations research*, 8(1), 101-111.
- 14 Davidson, P., Spinoulas, A., 2016. Driving alone versus riding together—How shared autonomous vehicles can
- 15 change the way we drive. *Road & Transport Research: A Journal of Australian and New Zealand Research and*
- 16 *Practice*, 25(3), 51-65.
- 17 de Almeida Correia, G.H. and van Arem, B., 2016. Solving the User Optimum Privately Owned Automated Vehicles
- 18 Assignment Problem (UO-POAVAP): A model to explore the impacts of self-driving vehicles on urban mobility.
- 19 *Transportation Research Part B: Methodological*, 87, 64-88.
- 20 Dey, K.C., Rayamajhi, A., Chowdhury, M., Bhavsar, P., Martin, J., 2016. Vehicle-to-vehicle (V2V) and vehicle-to-
- 21 infrastructure (V2I) communication in a heterogeneous wireless network—Performance evaluation. *Transportation*
- 22 *Research Part C: Emerging Technologies*, 68, 168-184.
- 23 Di, X. and Ban, X.J., 2019. A unified equilibrium framework of new shared mobility systems. *Transportation Research*
- 24 *Part B: Methodological*, 129, 50-78.
- 25 Djavadian, S. and Chow, J.Y., 2017. An agent-based day-to-day adjustment process for modeling ‘Mobility as a
- 26 Service’ with a two-sided flexible transport market. *Transportation research part B: methodological*, 104, 36-57.
- 27 Din, S., Paul, A. and Rehman, A., 2019. 5G-enabled Hierarchical architecture for software-defined intelligent
- 28 transportation system. *Computer Networks*, 150, 81-89.
- 29 Drissi-Kaïtouni, O., Hamed Bencheikroun, A., 1992. A dynamic traffic assignment model and a solution algorithm.
- 30 *Transportation Science* 26, 119–128.
- 31 Fagnant, D.J., Kockelman, K.M., 2014. The travel and environmental implications of shared autonomous vehicles,
- 32 using agent-based model scenarios. *Transportation Research Part C: Emerging Technologies*, 40, 1-13.
- 33 Fisher, M.L., 1981. The Lagrangian relaxation method for solving integer programming problems. *Management*
- 34 *science*, 27(1), 1-18.
- 35 Florian, M., Constantin, I., Florian, D. 2009. A New Look at Projected Gradient Method for Equilibrium Assignment.
- 36 *Transportation Research Record*, 2090(1), 10–16.
- 37 Frank, M., Wolfe, P., 1956. An algorithm for quadratic programming. *Naval Research Logistics (NRL)*, 3(1-2), 95-
- 38 110.
- 39 Fred Dotter, 2016. CIVITAS Insight 18 - Mobility-as-a-Service: A new transport model
- 40 Gendreau, M., Jabali, O., Rei, W., 2016. 50th anniversary invited article—future research directions in stochastic
- 41 vehicle routing. *Transportation Science*, 50(4), 1163-1173.
- 42 Gentili, M. and Mirchandani, P.B., 2012. Locating sensors on traffic networks: Models, challenges and research
- 43 opportunities. *Transportation research part C: emerging technologies*, 24, 227-255.
- 44 Ghali, M.O., Smith, M.J., 1995. A model for the dynamic system optimum traffic assignment problem. *Transportation*
- 45 *Research Part B: Methodological*, 29(3), 155-170.
- 46 Ghiasi, A., Hussain, O., Qian, Z.S., Li, X., 2017. A mixed traffic capacity analysis and lane management model for
- 47 connected automated vehicles: A Markov chain method. *Transportation Research Part B: Methodological*, 106,
- 48 266-292.
- 49 Hanebutte, U., Doss, E., Ewing, T., Tentner, A., 1998. Simulation of vehicle traffic on an automated highway system.
- 50 *Math. Comput. Model.* 27 (9–11), 129–141
- 51 Heikkilä, S., 2014. Mobility as a service—A proposal for action for the public administration, Case Helsinki. (Master
- 52 of Science in Technology), Aalto University, Helsinki, Finland.

- 1 Hensher, D. A., 2017. Future bus transport contracts under a mobility as a service (MaaS) regime in the digital age:  
2 Are they likely to change? *Transportation Research Part A: Policy and Practice*, 98, 86-96.
- 3 Hyland, M., Mahmassani, H.S., 2018. Dynamic autonomous vehicle fleet operations: Optimization-based strategies  
4 to assign AVs to immediate traveler demand requests. *Transportation Research Part C: Emerging*  
5 *Technologies*, 92, 278-297.
- 6 Huisman, D., Jans, R., Peeters, M. and Wagelmans, A.P., 2005. Combining column generation and Lagrangian  
7 relaxation. In *Column generation* (pp. 247-270). Springer, Boston, MA.
- 8 Jayakrishnan, R., Tsai, W.T., Prashker, J.N., Rajadhyaksha, S., 1994. A faster path-based algorithm for traffic  
9 assignment.
- 10 Jittrapirom, P., Caiati, V., Feneri, A.-M., Ebrahimigharehbaghi, S., González, M. J. A., Narayan, J., 2017. Mobility as  
11 a service: A critical review of definitions, assessments of schemes, and key challenges. *Urban Planning*, 2(2), 13-  
12 25.
- 13 Kalafatas, G., 2010. A graph theoretic modeling framework for generalized transportation systems with congestion  
14 phenomena. Purdue University.
- 15 Kalvelagen, E., 2003. Dantzig-Wolfe Decomposition with GAMS. Amsterdam Optim. Model. Group LLC,  
16 Washington, DC, USA.
- 17 Kamargianni, M., Li, W., Matyas, M., Schäfer, A., 2016. A critical review of new mobility services for urban transport.  
18 *Transportation Research Procedia*, 14, 3294-3303.
- 19 Kok, A.L., Hans, E.W. and Schutten, J.M., 2012. Vehicle routing under time-dependent travel times: the impact of  
20 congestion avoidance. *Computers & operations research*, 39(5), 910-918.
- 21 Larsson, T., Patriksson, M., 1992. Simplicial decomposition with disaggregated representation for the traffic  
22 assignment problem. *Transportation Science*, 26(1), 4-17.
- 23 Larsson, T., Patriksson, M. and Rydergren, C., 2004. A column generation procedure for the side constrained traffic  
24 equilibrium problem. *Transportation Research Part B: Methodological*, 38(1), 17-38.
- 25 Lavieri, P.S., Garikapati, V.M., Bhat, C.R., Pendyala, R.M., Astroza, S., Dias, F.F., 2017. Modeling individual  
26 preferences for ownership and sharing of autonomous vehicle technologies. *Transportation Research Record:*  
27 *Journal of the Transportation Research Board*, (2665), 1-10.
- 28 LeBlanc LJ, Morlok EK, Pierskalla WP, 1975. An efficient approach to solving the road network equilibrium traffic  
29 assignment problem. *Transportation Research* 9: 309–318
- 30 Levin, M.W., Boyles, S.D., 2016. A multiclass cell transmission model for shared human and autonomous vehicle  
31 roads. *Transportation Research Part C: Emerging Technologies*, 62, 103-116.
- 32 Levin, M.W., Kockelman, K.M., Boyles, S.D., Li, T., 2017. A general framework for modeling shared autonomous  
33 vehicles with dynamic network-loading and dynamic ride-sharing application. *Computers, Environment and*  
34 *Urban Systems*, 64, 373-383.
- 35 Li, Y., Waller, S.T., Ziliaskopoulos, T., 2003. A decomposition scheme for system optimal dynamic traffic assignment  
36 models. *Networks and Spatial Economics*, 3(4), 441-455.
- 37 Liang, X., Homem de Almeida Correia, G. and van Arem, B., 2018. Applying a model for trip assignment and dynamic  
38 routing of automated taxis with congestion: system performance in the City of Delft, The Netherlands.  
39 *Transportation Research Record*, 2672(8), 588-598.
- 40 Lim, H., Taeihagh, A., 2018. Autonomous vehicles for smart and sustainable cities: An in-depth exploration of privacy  
41 and cybersecurity implications. *Energies*, 11(5), 1062.
- 42 Lin, D.Y., Valsaraj, V., Waller, S.T., 2011. A Dantzig-Wolfe decomposition-based heuristic for off-line capacity  
43 calibration of dynamic traffic assignment. *Computer-Aided Civil and Infrastructure Engineering*, 26(1), 1-15.
- 44 Liu, J., Zhou, X., 2016. Capacitated transit service network design with boundedly rational agents. *Transportation*  
45 *Research Part B: Methodological*, 93, 225-250.
- 46 Liu, J., Kang, J.E., Zhou, X., Pendyala, R., 2018. Network-oriented household activity pattern problem for system  
47 optimization. *Transportation Research Part C: Emerging Technologies*, 94, 250-269.
- 48 Lu, C.C., Liu, J., Qu, Y., Peeta, S., Roupail, N.M., Zhou, X., 2016. Eco-system optimal time-dependent flow  
49 assignment in a congested network. *Transportation Research Part B: Methodological*, 94, 217-239.
- 50 Lu, C.C., Mahmassani, H.S., Zhou, X., 2009. Equivalent gap function-based reformulation and solution algorithm for  
51 the dynamic user equilibrium problem. *Transportation Research Part B: Methodological*, 43(3), 345-364.



- 1 Ma, J., Li, X., Zhou, F., Hao, W., 2017. Designing optimal autonomous vehicle sharing and reservation systems: A  
2 linear programming approach. *Transportation Research Part C: Emerging Technologies*, 84, 124-141.
- 3 Ma, T., Rasulkhani, S., Chow, J., Klein, S., 2018. A dynamic ridesharing dispatch and idle vehicle repositioning  
4 strategy with integrated transit transfers. <https://arxiv.org/abs/1901.00760>
- 5 Maciejewski, M., Bischoff, J., 2016. Congestion effects of autonomous taxi fleets. *Transport*. 1-10
- 6 Mahmassani, H.S., 2016. 50th anniversary invited article—autonomous vehicles and connected vehicle systems: flow  
7 and operations considerations. *Transportation Science*, 50(4), 1140-1162.
- 8 Mahmoudi, M., Zhou, X., 2016. Finding optimal solutions for vehicle routing problem with pickup and delivery  
9 services with time windows: A dynamic programming approach based on state–space–time network  
10 representations. *Transportation Research Part B: Methodological*, 89, 19-42.
- 11 Martinez, L.M., Correia, G.H., Viegas, J.M., 2015. An agent - based simulation model to assess the impacts of  
12 introducing a shared - taxi system: an application to Lisbon (Portugal). *Journal of Advanced Transportation*,  
13 49(3), 475-495.
- 14 Martinez, L.M., Viegas, J.M., 2017. Assessing the impacts of deploying a shared self-driving urban mobility system:  
15 An agent-based model applied to the city of Lisbon, Portugal. *International Journal of Transportation Science and*  
16 *Technology*, 6(1), 13-27.
- 17 Mourad, A., Puchinger, J., Chu, C., 2019. A survey of models and algorithms for optimizing shared mobility.  
18 *Transportation Research Part B*, in press.
- 19 Mulley, C., 2017. Mobility as a services (MaaS)—Does it have critical mass? *Transport Reviews*, 37(3), 247-251.
- 20 Munoz, J.C., Laval, J.A., 2006. System optimum dynamic traffic assignment graphical solution method for a  
21 congested freeway and one destination. *Transportation Research Part B: Methodological*, 40(1), 1-15.
- 22 Newell, G.F., 1993. A simplified theory of kinematic waves in highway traffic, part I: general theory. *Transp. Res.*  
23 *Part B: Methodological*. 27 (4), 281–287.
- 24 Newell, G.F., 2002. A simplified car-following theory: a lower order model. *Transp. Res. Part B: Methodol.* 36 (3),  
25 195–205
- 26 Nieuwenhuijsen, J., de Almeida Correia, G.H., Milakis, D., van Arem, B. and van Daalen, E., 2018. Towards a  
27 quantitative method to analyze the long-term innovation diffusion of automated vehicles technology using system  
28 dynamics. *Transportation Research Part C: Emerging Technologies*, 86, 300-327.
- 29 Palomar, D.P., Chiang, M., 2006. A tutorial on decomposition methods for network utility maximization. *IEEE*  
30 *Journal on Selected Areas in Communications*, 24(8), 1439-1451.
- 31 Papadimitratos, P., Fortelle, A.L., Evenssen, K., Brignolo, R. and Cosenza, S., 2009. Vehicular communication  
32 systems: Enabling technologies, applications, and future outlook on intelligent transportation. *IEEE*  
33 *communications magazine*, 47(11), 84-95.
- 34 Peeta, S., Mahmassani, H.S., 1995. System optimal and user equilibrium time-dependent traffic assignment in  
35 congested networks. *Annals of Operations Research*, 60(1), 81-113.
- 36 Peeta, S., Ziliaskopoulos, A.K., 2001. Foundations of dynamic traffic assignment: The past, the present and the  
37 future. *Networks and spatial economics*, 1(3-4), 233-265.
- 38 Psaraftis, Harilaos N., Min Wen, and Christos A. Kontovas. "Dynamic vehicle routing problems: Three decades and  
39 counting." *Networks* 67, no. 1 (2016): 3-31.
- 40 Qian, Z.S., Shen, W., Zhang, H.M., 2012. System-optimal dynamic traffic assignment with and without queue  
41 spillback: Its path-based formulation and solution via approximate path marginal cost. *Transportation research*  
42 *part B: methodological*, 46(7), 874-893.
- 43 Qu, F., Wang, F.Y. and Yang, L., 2010. Intelligent transportation spaces: vehicles, traffic, communications, and  
44 beyond. *IEEE Communications Magazine*, 48(11), 136-142.
- 45 Rayle, L., Dai, D., Chan, N., Cervero, R., & Shaheen, S. (2016). Just a better taxi? A survey-based comparison of  
46 taxis, transit, and ridesourcing services in San Francisco. *Transport Policy*, 45, 168-178.
- 47 Reece, D.A., Shafer, S.A., 1993. A computational model of driving for autonomous vehicles. *Transportation Research*  
48 *Part A: Policy and Practice*, 27(1), 23-50.
- 49 Rossi, F., Zhang, R., Hindy, Y., Pavone, M., 2018. Routing autonomous vehicles in congested transportation networks:  
50 Structural properties and coordination algorithms. *Autonomous Robots* 42, 1427-1442
- 51 Salazar, M., Rossi, F., Schiffer, M., Onder, C.H., Pavone, M., 2018. On the Interaction between Autonomous Mobility-  
52 on-Demand and Public Transportation Systems. *arXiv preprint arXiv:1804.11278*.

- Scherr, Y.O., Saavedra, B.A.N., Hewitt, M. and Mattfeld, D.C., 2019. Service network design with mixed autonomous fleets. *Transportation Research Part E: Logistics and Transportation Review*, 124, 40-55.
- Shen, W, Nie, Y, Zhang, H.M., 2007. On path marginal cost analysis and its relation to dynamic system-optimal traffic assignment. in *Proceedings of the 17th International Symposium on Transportation and Traffic Theory*.
- Small, K.A., 1982. The scheduling of consumer activities: Work trips. *American Economic Review* 72, 467-479.
- Small, K.A., 2015. The bottleneck model: An assessment and interpretation. *Economics of Transportation*, 4(1-2), 110-117.
- Small, K.A., Verhoef, E.T., 2007. *The Economics of Urban Transportation*. Routledge (Taylor & Francis), London.
- Stern, R.E., Cui, S., Delle Monache, M.L., Bhadani, R., Bunting, M., Churchill, M., Hamilton, N., Pohlmann, H., Wu, F., Piccoli, B., Seibold, B., 2018. Dissipation of stop-and-go waves via control of autonomous vehicles: Field experiments. *Transportation Research Part C: Emerging Technologies*, 89, 205-221.
- Sun, W., Zheng, J. and Liu, H.X., 2017. A capacity maximization scheme for intersection management with automated vehicles. *Transportation research procedia*, 23, 121-136.
- Talebpour, A., Mahmassani, H.S., 2016. Influence of connected and autonomous vehicles on traffic flow stability and throughput. *Transportation Research Part C: Emerging Technologies*, 71, 143-163.
- Taniguchi, E. and Shimamoto, H., 2004. Intelligent transportation system based dynamic vehicle routing and scheduling with variable travel times. *Transportation Research Part C: Emerging Technologies*, 12(3-4), 235-250.
- Tong, L., Pan, Y., Shang, P., Guo, J., Xian, K., Zhou, X., 2019. Open-Source Public Transportation Mobility Simulation Engine DTA Lite-S: A Discretized Space-Time Network-Based Modeling Framework for Bridging Multi-agent Simulation and Optimization. *Urban Rail Transit*, 5, 1-16.
- Toth, P., Vigo, D. eds., 2002. *The vehicle routing problem*. Society for Industrial and Applied Mathematics.
- Van Arem, B., van Driel, C.J.G., Visser, R., 2006. The impact of cooperative adaptive cruise control on traffic-flow characteristics. *IEEE Trans. Intell. Transport. Syst.* 7 (4), 429-436.
- Van Essen, J.T., Correia, G.H.A., 2019. Exact Formulation and Comparison Between the User Optimum and System Optimum Solution for Routing Privately Owned Automated Vehicles. *IEEE Trans. Intell. Transport. Syst.* 1, 1-12.
- Varaiya, P., 1993. Smart cars on smart roads: problems of control. *IEEE Transactions on automatic control*, 38(2), 195-207.
- Vickrey, W.S., 1969. Congestion theory and transport investment. *The American Economic Review*, 59(2), 251-260.
- Wardrop, J.G., 1952. Some theoretical aspects of road traffic research. *Proceedings of the Institution of Civil Engineers, Part II* 1: 325-378.
- Wei, Y., Avci, C., Liu, J., Belezamo, B., Aydın, N., Li, P.T., Zhou, X., 2017. Dynamic programming-based multi-vehicle longitudinal trajectory optimization with simplified car following models. *Transportation research part B: methodological*, 106, 102-129.
- Wong, Y. Z., Hensher, D. A., 2017. The Thredbo story: A journey of competition and ownership in land passenger transport. Paper presented at the 15th International Conference on Competition and Ownership in Land Passenger Transport (Thredbo 15), Stockholm, Sweden.
- Wu, X., Nie, L., Xu, M., Zhao, L., 2019. Distribution planning problem for a high-speed rail catering service considering time-varying demands and pedestrian congestion: A lot-sizing-based model and decomposition algorithm. *Transportation Research Part E*, 123, 61-89.
- Xing, T. and Zhou, X., 2011. Finding the most reliable path with and without link travel time correlation: A Lagrangian substitution based approach. *Transportation Research Part B: Methodological*, 45(10), pp.1660-1679.
- Yang, H., Huang, H.J., 2005. *Mathematical and economic theory of road pricing*.
- Yao, Y., Zhu, X., Dong, H., Wu, S., Wu, H., Zhou, X., 2018. An Alternating Direction Method of Multiplier Based Problem Decomposition Scheme for Iteratively Improving Primal and Dual Solution Quality in Vehicle Routing Problem, submitted to *Transportation Research Part B*
- Ye, L., Yamamoto, T., 2018. Modeling connected and autonomous vehicles in heterogeneous traffic flow. *Physica A: Statistical Mechanics and its Applications*, 490, 269-277.
- Zhang, H.M., Nie, Y., Qian, Z., 2013. Modelling network flow with and without link interactions: the cases of point queue, spatial queue and cell transmission model. *Transportmetrica B: Transport Dynamics*, 1(1), 33-51.

- 1 Zhou, X., Tong, L., Mahmoudi, M., Zhuge, L., Yao, Y., Zhang, Y., Shang, P., Liu, J., Shi, T., 2018. Open-source
- 2 VRPLite Package for Vehicle Routing with Pickup and Delivery: A Path Finding Engine for Scheduled
- 3 Transportation Systems. Urban Rail Transit, 4(2), 68-85.
- 4