
A BCMP Network Approach to Modeling and Controlling Autonomous Mobility-on-Demand Systems

Journal Title
XX(X):1–33
©The Author(s) 0000
Reprints and permission:
sagepub.co.uk/journalsPermissions.nav
DOI: 10.1177/ToBeAssigned
www.sagepub.com/



Ramon Iglesias¹, Federico Rossi¹, Rick Zhang² and Marco Pavone¹

Abstract

In this paper we present a queuing network approach to the problem of routing and rebalancing a fleet of self-driving vehicles providing on-demand mobility within a *capacitated* road network. We refer to such systems as autonomous mobility-on-demand (AMoD) systems. We first cast an AMoD system into a closed, multi-class BCMP queuing network model capable of capturing the passenger arrival process, traffic, the state-of-charge of electric vehicles, and the availability of vehicles at the stations. Second, we propose a scalable method for the synthesis of routing and charging policies, with performance guarantees in the limit of large fleet sizes. Third, explore the applicability of our theoretical results on a case study of Manhattan. Collectively, this paper provides a unifying framework for the analysis and control of AMoD systems, which provides a large set of modeling options (e.g., the inclusion of road capacities and charging constraints), and subsumes earlier Jackson and network flow models.

Keywords

Queuing networks, autonomous vehicles, self-driving cars, transportation

¹ Stanford University, Stanford CA 94305, USA

² Zoox Inc., Menlo Park, CA 94025, USA

Corresponding author:

Marco Pavone, 496 Lomita Mall, Rm. 261 Stanford, CA 94305

Email: pavone@stanford.edu

Introduction

Personal mobility in the form of privately owned automobiles contributes to increasing levels of traffic congestion, pollution, and under-utilization of vehicles (on average 5% in the US (Neil 2015)) – clearly unsustainable trends for the future. The pressing need to reverse these trends has spurred the creation of cost-competitive, on-demand personal mobility solutions such as car-sharing (e.g. Car2Go, ZipCar) and ride-sharing (e.g. Uber, Lyft). However, without proper fleet management, car-sharing and, to some extent, ride-sharing systems lead to vehicle imbalances; vehicles aggregate in some areas while becoming depleted in others, due to the asymmetry between trip origins and destinations (Zhang and Pavone 2016).

Self-driving vehicles offer the distinctive advantage of being able to rebalance themselves, in addition to the convenience, cost savings, and possibly safety of not requiring a driver. Indeed, it has been shown that one-way vehicle sharing systems with self-driving vehicles (referred to as autonomous mobility-on-demand systems, or AMoD) have the potential to significantly reduce passenger cost-per-mile-traveled, while keeping the advantages and convenience of personal mobility (Spieser et al. 2014). Accordingly, a number of works have recently investigated the potential of AMoD systems, with a specific focus on the synthesis and analysis of coordination algorithms.

This paper aims to devise a general, unifying analytical framework for the analysis and control of AMoD systems, which subsumes many of the analytical models recently presented in the literature, chiefly, Pavone et al. (2012), Zhang and Pavone (2016), and Rossi et al. (2018). Specifically, this paper extends the Jackson network approach in Zhang and Pavone (2016) by adopting a Baskett-Chandy-Muntz-Palacios (BCMP) queuing-theoretical framework (Baskett et al. 1975; Kobayashi and Gerla 1983). The generality offered by the BCMP framework allows one to take into account several real-world constraints, in particular (i) state-of-charge of autonomous electric vehicles and (ii) road capacities (that is, congestion). In contrast to previous work, the proposed BCMP model allows one to characterize such effects analytically along with performance guarantees. Moreover, the proposed BCMP model recovers the traffic congestion results in Rossi et al. (2018), with the additional benefits of taking into account the stochasticity of transportation networks and providing estimates for performance metrics. Thus, the results in this paper provide novel tools for the analysis and control of AMoD systems in the presence of stochasticity and system-wide constraints such as traffic congestion and vehicle charging.

Literature Review: The issue of vehicle rebalancing has been addressed in a variety of ways in the literature. For example, in the context of bike-sharing, Chemla et al. (2013) proposes rearranging the stock of bicycles between stations using trucks. The works in Nourinejad et al. (2015), Boyacı et al. (2015), and

A preliminary version of this paper has appeared in the Proceedings of the 2016 Workshop on the Algorithmic Foundations of Robotics (Iglesias et al. 2016).

Acquaviva et al. (2014) investigate using paid drivers to move vehicles between car-sharing stations where cars are parked, while Banerjee et al. (2015) studies the merits of dynamic pricing for incentivizing drivers to move to underserved areas.

Within the context of AMoD systems, where vehicles can rebalance themselves, previous work can be categorized into two main classes: heuristic methods and analytical methods. Heuristic routing strategies are extensively investigated in Fagnant and Kockelman (2014), Fagnant et al. (2015), and Levin et al. (2016) by leveraging a traffic simulator and, in Zhang et al. (2016), by leveraging a model predictive control framework. Analytical models of AMoD systems are proposed in Pavone et al. (2012), Zhang and Pavone (2016), and Rossi et al. (2018), by using fluidic, Jackson queuing network, and capacitated flow frameworks, respectively. Analytical methods have the advantage of providing structural insights (e.g., Rossi et al. (2018)), and provide guidelines for the synthesis of control policies. The problem of controlling AMoD systems is similar to the System Optimal Dynamic Traffic Assignment (SO-DTA) problem (see, e.g., Chiu et al. (2011); Patriksson (2015)) where the objective is to find optimal routes for all vehicles within congested or capacitated networks such that the total cost is minimized. The main differences between the AMoD control problem and the SO-DTA problem is that SO-DTA only optimizes customer routes, and *not* rebalancing routes.

Previous work addressing AMoD charging and congestion constraints rely either on deterministic models or are simulation-based studies. Integration of electric vehicles in AMoD systems has been studied in a model-predictive control setting in Zhang et al. (2016) and in an agent-based simulation framework in Chen et al. (2016). Both studies characterize the effects of charging speed on the level of service via simulations. As for congestion, the impact of AMoD systems on traffic has been a hot topic of debate. For example, Levin et al. (2016) notes that empty-traveling rebalancing vehicles may increase congestion and total in-vehicle travel time for customers, but Rossi et al. (2018) shows that, with congestion-aware routing and rebalancing, the increase in congestion can be avoided. However, their proposed model does not account for the stochasticity of travel demand.

Statement of Contributions: The contribution of this paper is threefold. First, we show how an AMoD system can be cast within the framework of closed, multi-class BCMP queuing networks. The framework captures stochastic passenger arrivals, vehicle routing on a road network, congestion effects, and battery charging-discharging for electric vehicles. Importantly, such a framework allows one to use a number of queuing theoretical tools to analyze performance metrics for a given routing policy in terms, e.g., of vehicle availabilities and second-order moments of vehicle throughput. Second, we propose a scalable method for the synthesis of routing and charging policies: namely, we show that, for large fleet sizes, the stochastic optimal routing and charging strategy can be found by solving a linear program. Finally, we explore the applicability of our theoretical results on a case study of Manhattan.

A preliminary version of this paper appeared as Iglesias et al. (2016). This extended and revised version contains as additional contributions: (i) an extension of the BCMP model to capture the state-of-charge of the vehicles, (ii) a corresponding extension of the solution algorithms, (iii) a new set of numerical experiments that characterize the effects of vehicle charging on the level of service, and (iv) proofs of all results.

Organization: The rest of the paper is organized as follows. In the “Background Material” section we discuss basic properties of BCMP networks. In the “Model Description and Problem Formulation” section we describe the model of the AMoD system in the presence of road congestion constraints, cast it into a BCMP network, and formally present the routing and rebalancing problem. In the “Asymptotically Optimal Algorithms for AMoD Routing” section we focus on the derivation of solution algorithms that are shown to achieve optimal performance in the limit of large fleet sizes. In the section on “Battery Charge Constraints” we extend the BCMP model and solution algorithms to capture the state-of-charge of the vehicles. We validate our approach in the “Numerical Experiments” section by performing a case study of Manhattan. Finally, in the “Conclusions” section, we state our concluding remarks and discuss potential avenues for future research.

Background Material

In this section we review some basic definitions and properties of BCMP networks, on which we will rely extensively later in the paper.

Closed, multi-class BCMP networks

Let \mathcal{Z} be a network consisting of N independent queues (or nodes). A set of agents move within the network according to a stochastic process, i.e., after receiving service at queue i they proceed to queue j with a given probability. No agent enters or leaves the network from the outside, so the number of agents is fixed and equal to m . Such a network is referred to as a *closed* queuing network. Agents belong to one of $K \in \mathbb{N}_{>0}$ classes, and they can switch between classes upon leaving a node.

Let $x_{i,k}$ denote the number of agents of class $k \in \{1, \dots, K\}$ at node $i \in \{1, \dots, N\}$. The state of node i , denoted by \mathbf{x}_i , is given by $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,K}) \in \mathbb{N}^K$. The state space of the network is (Gelenbe et al. 1998):

$$\Omega_m := \{(\mathbf{x}_1, \dots, \mathbf{x}_N) : \mathbf{x}_i \in \mathbb{N}^K, \sum_{i=1}^N \|\mathbf{x}_i\|_1 = m\},$$

where $\|\cdot\|_1$ denotes the standard 1-norm (i.e., $\|\mathbf{x}_i\|_1 = \sum_k |x_{i,k}|$). The relative frequency of visits (also known as relative throughput) to node i by agents of class k , denoted as $\pi_{i,k}$, is given by the traffic equations (Gelenbe et al. 1998):

$$\pi_{i,k} = \sum_{k'=1}^K \sum_{j=1}^N \pi_{j,k'} p_{j,k';i,k}, \quad \text{for all } i \in \{1, \dots, N\}, \quad (1)$$

where $p_{j,k';i,k}$ is the probability that upon leaving node j , an agent of class k' goes to node i and becomes an agent of class k . Equation (1) does not have a unique solution (a typical feature of closed networks), and $\pi = \{\pi_{i,k}\}_{i,k}$ only determines frequencies up to a constant factor (hence the name “relative” frequency). It is customary to express frequencies in terms of a chosen reference node, e.g., so that $\pi_{1,1} = 1$.

Queues are allowed to be one of four types: First Come, First Served (FCFS), Processor Sharing, Infinite Server, and Last Come, First Served. FCFS nodes have exponentially distributed service times, while the other three queue types may follow any Cox distribution (Gelenbe et al. 1998). Such a queuing network model is referred to as a closed, multi-class BCMP queuing network (Gelenbe et al. 1998).

Let \mathcal{N} represent the set of nodes in the network. For the remainder of the paper, we will focus on networks that have only two types of nodes: FCFS queues with a single server (for short, SS queues), forming a set $\mathcal{S} \subset \mathcal{N}$, and infinite server queues (for short, IS queues), forming a set $\mathcal{I} \subset \mathcal{N}$. Furthermore, we consider class-independent and load-independent nodes (i.e. nodes whose service rate is independent of the agent classes or number of agents in the queue), whereby at each node $i \in \{1, \dots, N\}$ the service rate is given by:

$$\mu_i(x_i) = c_i(x_i)\mu_i^o,$$

where $x_i := \|\mathbf{x}_i\|_1$ is the number of agents at node i , μ_i^o is the (class-independent) base service rate, and $c_i(x_i)$ is the (load-independent) *capacity* function

$$c_i(x_i) = \begin{cases} x_i & \text{if } x_i \leq c_i^o, \\ c_i^o & \text{if } x_i > c_i^o, \end{cases}$$

which depends on the number of servers c_i^o at the queue. In the case considered in this paper, $c_i^o = 1$ for all $i \in \mathcal{S}$ and $c_i^o = \infty$ for all $i \in \mathcal{I}$.

Under the assumption of class-independent service rates, the multi-class network \mathcal{Z} can be “compressed” into a single-class network \mathcal{Z}^* with state-space $\Omega_m^* := \{(x_1, \dots, x_N) : x_i \in \mathbb{N}, \sum_{i=1}^N x_i = m\}$ (Kant and Srinivasan 1992). Performance metrics for the original, multi-class network \mathcal{Z} can be found by first analyzing the compressed network \mathcal{Z}^* , and then applying suitable scalings for each class. Specifically, let $\pi_i = \sum_{k=1}^K \pi_{i,k}$ and $\gamma_i = \sum_{k=1}^K \frac{\pi_{i,k}}{\mu_i^o}$, be the total relative throughput and relative utilization at a node i , respectively. Then, the stationary distribution of the compressed, single-class network \mathcal{Z}^* is given by

$$\mathbb{P}(x_1, \dots, x_N) = \frac{1}{G(m)} \prod_{i=1}^N \frac{\gamma_i^{x_i}}{\prod_{a=1}^{x_i} c_i(a)}, \quad \text{where } G(m) = \sum_{x_i \in \Omega_m^*} \prod_{i=1}^N \frac{\gamma_i^{x_i}}{\prod_{a=1}^{x_i} c_i(a)}$$

is a normalizing constant. Remarkably, the stationary distribution has a product form, a key feature of BCMP networks.

Three performance metrics that are of interest at each node are throughput, expected queue length, and availability. First, the throughput at a node (i.e.,

the number of agents processed by a node per unit of time) is given by

$$\Lambda_i(m) = \pi_i \frac{G(m-1)}{G(m)}. \quad (2)$$

Second, let $\mathbb{P}_i(x_i; m)$ be the probability of finding x_i agents at node i ; then the expected queue length at node i is given by $L_i(m) = \sum_{x_i=1}^m x_i \mathbb{P}_i(x_i; m)$.

In the case of IS nodes (i.e., nodes in \mathcal{I}), the expected queue length can be more easily derived via Little's Law as (George 2012)

$$L_i(m) = \Lambda_i(m) / \mu_i^o, \quad \text{for all } i \in \mathcal{I}. \quad (3)$$

The throughputs and the expected queue lengths for the original, multi-class network \mathcal{Z}^* can be found via scaling (Kant and Srinivasan 1992), specifically, $\Lambda_{i,k}(m) = (\pi_{i,k}/\pi_i)\Lambda_i(m)$ and $L_{i,k}(m) = (\pi_{i,k}/\pi_i)L_i(m)$.

Finally, the availability of single-server, FCFS nodes (i.e., nodes in \mathcal{S}) is defined as the probability that the node has at least one agent, and is given by (George 2012)

$$A_i(m) = \gamma_i \frac{G(m-1)}{G(m)}, \quad \text{for all } i \in \mathcal{S}.$$

It is worth noting that evaluating the three performance metrics above requires computation of the normalization constant $G(m)$, which is computationally expensive. However, several techniques are available to avoid the direct computation of $G(m)$. In particular, in this paper we use the Mean Value Analysis method (Gelenbe et al. 1998).

Asymptotic behavior of closed BCMP networks

In this section we describe the asymptotic behavior of closed BCMP networks as the number of agents m goes to infinity. The results described in this section are taken from George (2012), and are detailed for a single-class network. However, as stated in the previous section, results found for a single-class network can easily be ported to the multi-class equivalent in the case of class-independent service rates.

Let $\rho_i := \gamma_i / c_i^o$ be the utilization factor of node $i \in \mathcal{N}$, where c_i^o is the number of servers at node i . Assume that the relative throughputs $\{\pi_i\}_i$ are normalized so that $\max_{i \in \mathcal{S}} \rho_i = 1$; furthermore, assume that nodes are ordered by their utilization factors so that $1 = \rho_1 \geq \rho_2 \geq \dots \geq \rho_N$, and define the set of bottleneck nodes as $\mathcal{G} := \{i \in \mathcal{S} : \rho_i = 1\}$.

It can be shown (George 2012, p. 14) that, as the number of agents m in the system approaches infinity, the availability at all bottleneck nodes converges to 1 while the availability at non-bottleneck nodes is strictly less than one, that is

$$\lim_{m \rightarrow \infty} A_i(m) \begin{cases} = 1 & \forall i \in \mathcal{G}, \\ < 1 & \forall i \notin \mathcal{G}. \end{cases} \quad (4)$$

Additionally, the queue lengths at the non-bottleneck nodes have a limiting distribution given by

$$\lim_{m \rightarrow \infty} \mathbb{P}_i(x_i; m) = \begin{cases} (1 - \rho_i) \rho_i^{x_i} & i \in \mathcal{S}, i \notin \mathcal{G}, \\ e^{-\gamma_i} \frac{\gamma_i^{x_i}}{x_i!} & i \in \mathcal{I}. \end{cases} \quad (5)$$

Together, (4) and (5) have strong implications for the operation of queuing networks with a large number of agents, and in particular for the operation of AMoD systems. Intuitively, (4) shows that as we increase the number of agents in the network, they will be increasingly queued at bottleneck nodes, driving availability in those queues to one. Alternatively, non-bottleneck nodes will converge to an availability strictly less than one, implying that there is always a non-zero probability of having an empty queue. In other words, agents will aggregate at the bottlenecks and become depleted elsewhere. Additionally, (5) shows that, as the number of agents goes to infinity, non-bottleneck nodes tend to behave like queues in an equivalent open BCMP network with the bottleneck nodes removed, i.e., individual performance metrics can be calculated in isolation.

Model Description and Problem Formulation

In this section, we introduce a BCMP network model for AMoD systems, and formalize the problem of routing and rebalancing such systems under stochastic conditions. Casting an AMoD system as a queuing network allows us to characterize and compute key performance metrics including the distribution of the number of vehicles on each road link (a key metric to characterize traffic congestion) and the probability of servicing a passenger request. To emphasize the relationship with the theory presented in the previous section, we reuse the same notation whenever concepts are equivalent.

Autonomous Mobility-on-Demand model

Consider a set of stations* \mathcal{S} distributed within an urban area connected by a network of road links \mathcal{I} , and m autonomous vehicles providing one-way transportation between these stations for incoming customers. Customers arrive to a station $s \in \mathcal{S}$ with a target destination $t \in \mathcal{S}$ according to a time-invariant Poisson process with rate $\lambda \in \mathbb{R}_{>0}$. The arrival process for all origin-destination pairs is summarized by the set of tuples $\mathcal{Q} = \{(s^{(q)}, t^{(q)}, \lambda^{(q)})\}_q$.

If on customer arrival there is an available vehicle, the vehicle drives the customer towards its destination. Alternatively, if there are no vehicles, the customer leaves the system (i.e., chooses an alternative transportation system). Thus, we adopt a *passenger loss* model. Such model is appropriate for systems where high quality-of-service is desired; from a technical standpoint, this

*Stations are not necessarily physical locations: they can also be interpreted as a set of geographical regions.

modeling assumption decouples the passenger queuing process from the vehicle queuing process.

A vehicle driving a passenger through the road network follows a routing policy $\alpha^{(q)}$ (defined in the next section) from origin to destination, where q indicates the origin-destination-rate tuple. Once it reaches its destination, the vehicle joins the station FCFS queue and waits for an incoming trip request.

A known problem of such systems is that vehicles will inevitably accumulate at one or more of the stations and reduce the number of vehicles servicing the rest of the system (George 2012) if no corrective action is taken. To control this problem, we introduce a set of “virtual rebalancing demands” or “virtual passengers” whose objective is to balance the system, i.e., to move empty vehicles to stations experiencing higher passenger loss. Similar to passenger demands, rebalancing demands are defined by a set of origin, destination and arrival rate tuples $\mathcal{R} = \{(s^{(r)}, t^{(r)}, \lambda^{(r)})\}_r$, and a corresponding routing policy $\alpha^{(r)}$. Therefore, the objective is to find a set of routing policies $\alpha^{(q)}, \alpha^{(r)}$, for all $q \in \mathcal{Q}$, $r \in \mathcal{R}$, and rebalancing rates $\lambda^{(r)}$, for all $r \in \mathcal{R}$, that balances the system while minimizing the number of vehicles on the road, and thus reducing the impact of the AMoD system on overall traffic congestion.

Casting an AMoD system into a BCMP network

We are now in a position to frame an AMoD system in terms of a BCMP network model. Initially, we will present the framework in the absence of charging constraints, and in a later section we will extend the model to include them. First, the passenger loss assumption allows the model to be characterized as a queuing network with respect only to the *vehicles*. In other words, at each station node, vehicles form a queue while waiting for customers and are “serviced” when a customer arrives. Thus, we will henceforth use the term “vehicles” to refer to the queuing agents. From this perspective, the stations \mathcal{S} are equivalent to SS queues, and the road links \mathcal{I} are modeled as IS queues. The set of all queues is given by $\mathcal{N} = \{\mathcal{S} \cup \mathcal{I}\}$, in analogy with the Background Material.

Second, we specify the BCMP network model. We abstract the underlying road network and the stations as a directed graph where the edges represent either the road links or the stations, and the vertices represent the road intersections. The BCMP network model can then be derived from such a directed graph as follows. Let $\text{Parent}(i)$ and $\text{Child}(i)$ be the origin and destination vertices of edge i in the directed graph. Then, a road that goes from intersection j to intersection l is represented in the BCMP network model by an IS queue $i \in \mathcal{I}$ such that $\text{Parent}(i) = j$ and $\text{Child}(i) = l$. Note that the road may not have lanes in the opposite direction, in which case a queue i' with $\text{Parent}(i') = l$ and $\text{Child}(i') = j$ would not exist. For example, in Figure 1, queue 14 starts at vertex 1 and ends at vertex 5. However, there is no queue that connect the vertices in the opposite direction. In turn, we assume that stations are adjacent to road intersections, and therefore stations are represented in the BCMP network model as SS queues $i \in \mathcal{S}$ with the same parent and child vertex, i.e., a self-loop. An intersection may have access to either one station (e.g., vertex

2 in Figure 1) or zero stations (e.g., vertex 5 in Figure 1). Finally, intersections in the BCMP network model are simply conceptual entities playing the role of “connectors” among the queues, see Figure 1.

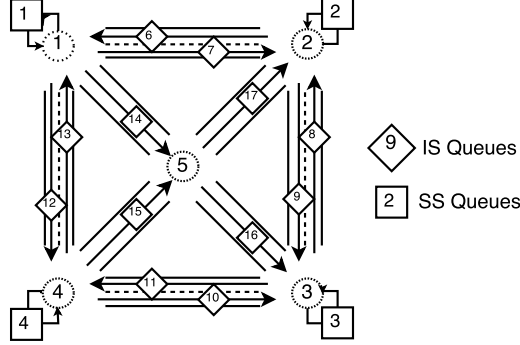


Figure 1. BCMP network model of an AMoD system. Diamonds represent infinite-server road links, squares represent the single-server vehicle stations, and dotted circles represent road intersections (playing the role of “connectors” among the queues).

Third, we introduce classes to represent the process of choosing destinations. We map the set of tuples \mathcal{Q} and \mathcal{R} to a set of classes \mathcal{K} such that $\mathcal{K} = \{\mathcal{Q} \cup \mathcal{R}\}$. Moreover, let \mathcal{O}_i be the subset of classes whose origin $s^{(k)}$ is the station i , i.e., $\mathcal{O}_i := \{k \in \mathcal{K} : s^{(k)} = i\}$, and \mathcal{D}_i be the subset of classes whose destination $t^{(k)}$ is the station i , i.e., $\mathcal{D}_i := \{k \in \mathcal{K} : t^{(k)} = i\}$. Thus, the probability that a vehicle at station i will leave for station j with a (real or virtual) passenger is the ratio between the respective (real or virtual) arrival rate $\lambda^{(k)}$, with $s^{(k)} = i$, $t^{(k)} = j$, and the sum of all arrival rates at station i . Formally, the probability that a vehicle of class k switches to class k' upon arrival to its destination $t^{(k)}$ is

$$\tilde{p}_{t^{(k)}}^{(k')} = \left(\lambda^{(k')} / \tilde{\lambda}_{t^{(k)}} \right),$$

where $\tilde{\lambda}_i = \sum_{k \in \mathcal{O}_i} \lambda^{(k)}$ is the sum of all arrival rates at station i . In other words, $\tilde{\lambda}_i$ represents the rate of arrival of passenger and rebalancing requests to station i , while $\tilde{p}_i^{(k')}$ encodes the likelihood of whether the request is a real passenger or rebalancing task and the desired target destination. Note that at all times a vehicle belongs to *some* class $k \in \mathcal{K}$, regardless of whether it is waiting at a station or traveling along the network.

The traversal of a vehicle from its source $s^{(k)}$ to its destination $t^{(k)}$ is guided by a routing policy $\alpha^{(k)}$. This routing policy consists of a matrix of transition probabilities. Let $\mathcal{W}_i = \{j \in \mathcal{N} : \text{Parent}(j) = i\}$ be the set of queues that begin at vertex i , and $\mathcal{U}_i = \{j \in \mathcal{N} : \text{Child}(j) = i\}$ be the set of queues that end at vertex i . A vehicle of class k leaves the station $s^{(k)}$ via one of the adjacent roads $j \in \mathcal{W}_{\text{Child}(s^{(k)})}$ with probability $\alpha_{s^{(k)}, j}^{(k)}$. It continues traversing the road network via these adjacency relationships following the routing probabilities

$\alpha_{i,j}^{(k)}$ until it is adjacent to its goal $t^{(k)}$. At this point, the vehicle proceeds to the destination and changes its class to $k' \in \mathcal{O}_{t^{(k)}}$ with probability $\tilde{p}_{t^{(k)}}^{(k')}$. This behavior is encapsulated by the routing matrix

$$p_{i,k;j,k'} = \begin{cases} \alpha_{i,j}^{(k)} & \text{if } k = k', j \in \mathcal{W}_{\text{Child}(i)}, t^{(k)} \notin \mathcal{W}_{\text{Child}(i)}, \\ \tilde{p}_j^{(k')} & \text{if } j = t^{(k)}, t^{(k)} \in \mathcal{W}_{\text{Child}(i)}, k' \in \mathcal{O}_j, \\ 0 & \text{otherwise,} \end{cases} \quad (6)$$

such that $\sum_{j,k'} p_{i,k;j,k'} = 1$. Thus, the relative throughput $\pi_{i,k}$, total relative throughput π_i , and utilization γ_i have the same definition as in the Background Material.

As stated before, the queuing process at each station is modeled as a SS queue where the service rate of the vehicles $\mu_i(a)$ is equal to the sum of real and virtual passenger arrival rates, i.e., $\mu_i(a) = \tilde{\lambda}_i$ for any station i and queue length a . Additionally, by modeling road links as IS queues, we assume that their service rates follow a Cox distribution with mean $\mu_i(a) = \frac{c_i(a)}{T_i}$, where T_i is the expected time required to traverse link i in absence of congestion, and $c_i(a)$ is the capacity factor when there are a vehicles in the queue. In this paper, we only consider the case of load-independent travel times, therefore $c_i(a) = a$ for all a , i.e., the service rate is the same regardless of the number of vehicles on the road. We do not make further assumptions on the distribution of the service times. The assumption of load-independent travel times is representative of uncongested traffic (Bureau of Public Roads 1964); in the next section we discuss how to incorporate probabilistic constraints for congestion on road links.

Problem formulation

As stated in Equation (4), vehicles tend to accumulate in bottleneck stations driving their availability to 1 as the fleet size increases, while the rest of the stations have availability strictly smaller than 1. In other words, for unbalanced systems, availability at most stations is capped regardless of fleet size. Therefore, it is desirable to make all stations “bottleneck” stations, i.e., set the constraint $\gamma_i = \gamma_j$ for all $i, j \in \mathcal{S}$, so as to (i) enforce a natural notion of “service fairness,” and (ii) prevent needless accumulation of empty vehicles at the stations.

However, it is desirable to minimize the impact that the rebalancing vehicles have on the road network. We achieve this by minimizing the expected number of vehicles on the road serving customer and rebalancing demands. Using Equation (3), the expected number of vehicles on a given road link i is given by $\Lambda_i(m)T_i$.

Lastly, we wish to avoid congestion on the individual road links. Traditionally, the relation between vehicle flow and congestion is parametrized by two basic quantities: the *free-flow travel time* T_i , i.e., the time it takes to traverse a link in absence of other traffic; and the *nominal capacity* C_i , i.e., the measure of traffic flow beyond which travel time increases very rapidly (Patriksson 2015). Assuming that travel time remains approximately constant when traffic is below the nominal capacity (an assumption typical of many state-of-the-art traffic models (Patriksson 2015)), our approach is to keep the expected traffic $\Lambda_i(m)T_i$

below the nominal capacity C_i and thus avoid congestion effects. Note that by constraining in expectation there is a non-zero probability of exceeding the capacity; however, we will show that, asymptotically, it is also possible to constrain the *probability* of exceeding the congestion constraint.

Accordingly, the routing problem we wish to study in this paper (henceforth referred to as the *Optimal Stochastic Capacitated AMoD Routing and Rebalancing problem*, or OSCARR) can now be formulated as follows:

$$\begin{aligned}
& \underset{\lambda^{(r \in \mathcal{R})}, \alpha_{ij}^{(k \in \mathcal{K})}}{\text{minimize}} && \sum_{i \in \mathcal{I}} \Lambda_i(m) T_i, \\
& \text{subject to} && \gamma_i = \gamma_j, && i, j \in \mathcal{S}, && (7a) \\
& && \Lambda_i(m) T_i \leq C_i, && i \in \mathcal{I}, && (7b) \\
& && \pi_{s^{(k)}, k} = \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{N}} \pi_{j, k} p_{j, k; t^{(k)}, k'}, && k \in \mathcal{K}, && (7c) \\
& && \pi_{i, k} = \sum_{k \in \mathcal{K}} \sum_{j \in \mathcal{N}} \pi_{j, k'} p_{j, k'; i, k} && i \in \{\mathcal{S} \cup \mathcal{I}\}, && (7d) \\
& && \sum_{j \in \mathcal{W}_{\text{Child}(i)}} \alpha_{ij}^{(k)} = 1, \quad \alpha_{ij}^{(k)} \geq 0, && i, j \in \{\mathcal{S} \cup \mathcal{I}\}, && (7e) \\
& && \lambda^{(r)} \geq 0, && r \in \mathcal{R}. && (7f)
\end{aligned}$$

Constraint (7a) enforces equal availability at all stations, while constraint (7b) ensures that all road links are (on average) uncongested. Constraints (7c)–(7f) enforce consistency in the model. Specifically, (7c) ensures that all traffic leaving the source $s^{(k)}$ of class k arrives at its destination $t^{(k)}$, (7d) enforces the traffic equations (1), (7e) ensures that $\alpha_{ij}^{(k)}$ is a valid probability measure, and (7f) guarantees nonnegative rebalancing rates.

Limitations

At this point, we would like to reiterate some assumptions and limitations built into the model. First, the proposed model is time-invariant. That is, we assume that customer and rebalancing rates remain constant for the segment of time under analysis, and that the network is able to reach its equilibrium distribution. An option for including the variation of customer demand over time is to discretize a period of time into smaller segments, each with its own arrival parameters and resulting rebalancing rates. These customer arrival rates, in turn, could be conditioned on external factors such as weather. Second, the passenger loss model assumes impatient customers and is well suited for cases where a high level of service is required. This allows us to simplify the model by focusing only on the vehicle process; however, it disregards the fact that customers may have different waiting thresholds and, consequently, the queuing process of waiting customers. Third, we focus on keeping traffic within the nominal road capacities in expectation, allowing us to assume load-independent travel times and to model exogenous traffic as a reduction in road capacity. Finally, we make no assumptions on the distribution of travel times on the road

links: the analysis proposed in this paper captures arbitrary distributions of travel times and only depends on the *mean* travel time.

Asymptotically Optimal Algorithms for AMoD Routing

In this section we show that, as the fleet size goes to infinity, the solution to OSCARR can be found by solving a linear program. This insight allows the efficient computation of asymptotically optimal routing and rebalancing policies and the characterization of the corresponding performance parameters.

First, we show that (i) the relative throughput at the stations can be expressed in terms of the relative throughputs at the other stations, and (ii) the balanced network constraint can be expressed in terms of the arrival rates. Then, we express the problem from a flow conservation perspective. Finally, we show that the problem allows an asymptotically optimal solution with bounds on the probability of exceeding road capacities. The solution we find is equivalent to the one presented in Rossi et al. (2018): thus, we show that the network flow model in Rossi et al. (2018) also captures the asymptotic behavior of a stochastic AMoD routing and rebalancing problem.

Folding of traffic equations

The next two lemmas show that the traffic equations (1) at the SS queues can be expressed in terms of other SS queues, and that the balanced network constraint can be expressed in terms of real and virtual passenger arrivals.

Lemma 1. (Folding of traffic equations). *Consider a feasible solution to OSCARR. Then, the total relative throughputs of the single server stations can be expressed in terms of the relative throughputs of the other single server stations, that is*

$$\pi_i = \sum_{k \in \mathcal{D}_i} \tilde{p}_s^{(k)} \pi_{s^{(k)}}, \quad i \in \mathcal{S}. \quad (8)$$

Proof. Using the routing matrix specified in Equation (6) we can rewrite the class throughputs (1) as

$$\begin{aligned} \pi_{i,k} &= \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{N}} \pi_{j,k'} p_{j,k';i,k} = \sum_{k' \in \mathcal{D}_i} \sum_{j \in N_{in}(j)} \pi_{j,k'} p_{j,k';i,k}, \\ &= \sum_{k' \in \mathcal{D}_i} \sum_{j \in N_{in}(j)} \pi_{j,k'} \tilde{p}_i^{(k)} = \tilde{p}_i^{(k)} \sum_{k' \in \mathcal{D}_i} \sum_{j \in N_{in}(j)} \pi_{j,k'}. \end{aligned} \quad (9)$$

The second equality exploits the fact that only queues feeding into i and vehicles whose class destination is i are routed to i . The third and fourth equalities follow from the fact that the probability of switching into class k at queue i is the same regardless of the original class k' . This allows us to rewrite the total relative throughput as

$$\pi_i = \sum_{k \in \mathcal{K}} \tilde{p}_i^{(k)} \sum_{k' \in \mathcal{D}_i} \sum_{j \in N_{in}(j)} \pi_{j,k'} = \sum_{k' \in \mathcal{D}_i} \sum_{j \in N_{in}(j)} \pi_{j,k'}, \quad (10)$$

since $\sum_{k=1}^K \tilde{p}_i^{(k)} = 1$. As a consequence of (10) and (9), the class relative throughputs can be related to the total relative throughputs

$$\pi_{i,k} = \tilde{p}_i^{(k)} \pi_i. \quad (11)$$

Now, assume the relative throughputs belong to a feasible solution to OSCARR. We proceed to reduce (7c) by using the routing matrix

$$\pi_{s^{(k)},k} = \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{N}} \pi_{j,k'} P_{j,k;t^{(k)},k'} = \sum_{k' \in \mathcal{K}} \tilde{p}_{t^{(k)}}^{(k')} \sum_{j \in \mathcal{N}_{in}(t^{(k)})} \pi_{j,k'} = \sum_{j \in \mathcal{N}_{in}(t^{(k)})} \pi_{j,k'}. \quad (12)$$

By inserting this into (10) and applying (11) we obtain

$$\pi_i = \sum_{k \in \mathcal{D}_i} \pi_{s^{(k)},k} = \sum_{k \in \mathcal{D}_i} \tilde{p}_{s^{(k)}}^{(k)} \pi_{s^{(k)}}. \quad (13)$$

□

Lemma 2. (Balanced system in terms of arrival rates). *Consider a feasible solution to OSCARR, then the constraint $\gamma_i = \gamma_j$ for all i, j is equivalent to*

$$\tilde{\lambda}_i = \sum_{k \in \mathcal{D}_i} \lambda^{(k)}. \quad (14)$$

Proof. The proof of this lemma is very similar to Theorem 4.3 in Zhang and Pavone (2016). Consider the case where (14) holds. We can write (8) in terms of the relative utilization rate:

$$\left(\sum_{k \in \mathcal{D}_i} \lambda^{(k)} \right) \gamma_i = \sum_{k \in \mathcal{D}_i} \gamma_{s^{(k)}} \lambda^{(k)}. \quad (15)$$

Now, by grouping customer and rebalancing classes by origin-destination pairs, we define φ as

$$\varphi_{ij} = \lambda^{(a)} + \lambda^{(r)}, \quad (16)$$

such that $s^{(a)} = s^{(r)} = j$ and $t^{(a)} = t^{(r)} = i$. Additionally, let $\zeta_{ij} = \varphi_{ij} / \sum_j \varphi_{ij}$. We note that there are no classes for which $s^{(k)} = t^{(k)}$, so we set $\varphi_{ii} = \zeta_{ii} = 0$. Under this definition, the variables $\{\zeta_{ij}\}_{ij}$ represent an irreducible Markov chain. Thus, Equation (15) can be rewritten as $\gamma_i = \sum_j \gamma_j \zeta_{ij}$ or more compactly as $Z\gamma = \gamma$, where the rows of Z are $[\zeta_{i1}, \zeta_{i2}, \dots, \zeta_{iS}]$, with $S = |\mathcal{S}|$, $i = 1, \dots, S$, and $\gamma = (\gamma_1, \dots, \gamma_S)$. Since Z is an irreducible, row stochastic Markov chain, by the Perron-Frobenius theorem the unique solution is given by $\gamma = (1, \dots, 1)^T$. Thus, $\gamma_i = \gamma_j$ for all i .

On the other hand, we consider again Equation (15). If the network \mathcal{Z} is a solution to problem (7), then for all i, j , we have $\gamma_i = \gamma_j = \gamma$, and (15) becomes

$$\gamma \tilde{\lambda}_i = \gamma \sum_{k \in \mathcal{D}_i} \lambda^{(k)}, \quad \tilde{\lambda}_i = \sum_{k \in \mathcal{D}_i} \lambda^{(k)}. \quad (17)$$

□

Asymptotically optimal solution

As discussed in the background material, relative throughputs are computed up to a constant multiplicative factor. Thus, without loss of generality, we can set the additional constraint $\pi_{s^{(1)}} = \tilde{\lambda}_1$, which, along with (7a), implies that

$$\pi_i = \tilde{\lambda}_i, \quad \pi_{s^{(k)},k} = \lambda^{(k)}, \quad \text{and} \quad \gamma_i = 1, \quad \text{for all } i \in \mathcal{S}. \quad (18)$$

As discussed earlier, the availabilities of stations with the highest relative utilization tend to one as the fleet size goes to infinity. Since the stations are modeled as SS queues, $\rho_i = \gamma_i$ for all $i \in \mathcal{S}$. Therefore, if the system is balanced, $\gamma_i = \gamma_S^{\max} = \gamma = 1$ for all $i \in \mathcal{S}$. That is, the set of bottleneck stations \mathcal{G} includes all stations in \mathcal{S} and $\lim_{m \rightarrow \infty} \frac{G(m-1)}{G(m)} = 1$ by Equation (4).

As $m \rightarrow \infty$ and $\frac{G(m-1)}{G(m)} \rightarrow 1$, the throughput at every station $\Lambda_i(m)$ becomes a linear function of the relative frequency of visits to that station, according to Equation (2). Thus, the objective function and the constraints in (7) are reduced to linear functions. We define the resulting problem (i.e., Problem (7) with $G(m-1)/G(m) = 1$) as the *Asymptotically Optimal Stochastic Capacitated AMoD Routing and Rebalancing problem*, or A-OSCARR. The following lemma shows that the optimal solution to OSCARR approaches the optimal solution to A-OSCARR as m increases.

Lemma 3. (Asymptotic behavior of OSCARR). *Let $\{\pi_{i,k}^*(m)\}_{i,k}$ be a set of relative throughputs corresponding to an optimal solution to OSCARR with a given set of customer demands $\{\lambda^{(q)}\}_q$ and a fleet size m . Also, let $\{\hat{\pi}_{i,k}\}_{i,k}$ be a set of relative throughputs corresponding to an optimal solution to A-OSCARR for the same set of customer demands. Then,*

$$\lim_{m \rightarrow \infty} \frac{G(m-1)}{G(m)} \sum_{i \in \mathcal{I}} T_i \sum_{k \in \mathcal{K}} \pi_{i,k}^*(m) = \sum_{i \in \mathcal{I}} T_i \sum_{k \in \mathcal{K}} \hat{\pi}_{i,k}. \quad (19)$$

Proof. We arrive to the proof by contradiction. Recall that $\pi_i = \sum_{k \in \mathcal{K}} \pi_{i,k}$. Assume Equation (19) did not hold. By definition,

$$\frac{G(m-1)}{G(m)} \sum_{i \in \mathcal{I}} T_i \pi_i^*(m) \leq \frac{G(m-1)}{G(m)} \sum_{i \in \mathcal{I}} T_i \pi_i, \quad (20)$$

and

$$\sum_{i \in \mathcal{I}} T_i \hat{\pi}_i(m) \leq \sum_{i \in \mathcal{I}} T_i \pi_i, \quad (21)$$

for all m and $\{\pi_{i,k}\}_{i,k}$. Applying the limit to (20) and using (4), we obtain

$$\sum_{i \in \mathcal{I}} T_i \lim_{m \rightarrow \infty} (\pi_i^*(m)) \leq \sum_{i \in \mathcal{I}} T_i \pi_i.$$

However, according to our assumption,

$$\text{either } \sum_{i \in \mathcal{I}} T_i \lim_{m \rightarrow \infty} (\pi_i^*(m)) > \sum_{i \in \mathcal{I}} T_i \hat{\pi}_i, \quad \text{or} \quad \sum_{i \in \mathcal{I}} T_i \lim_{m \rightarrow \infty} (\pi_i^*(m)) < \sum_{i \in \mathcal{I}} T_i \hat{\pi}_i.$$

But the former violates Equation (20), and the latter Equation (21). \square

As discussed in the Problem Formulation, constraint (7b) only enforces an upper bound on the expected number of vehicles traversing a link. However, in the asymptotic regime, it is possible to enforce an analytical upper bound on the *probability* of exceeding the nominal capacity of any given road link. As seen in Equation (5), as the fleet size increases, the distribution of the number of vehicles on a road link i converges to a Poisson distribution with mean $T_i\pi_i$. The cumulative density function of a Poisson distribution is given by $Pr(X < \bar{x}) = Q(\lfloor \bar{x} + 1 \rfloor, \bar{C})$, where \bar{C} is the mean of the distribution and Q is the regularized upper incomplete gamma function. Let ϵ be the maximum tolerable probability of exceeding the nominal capacity. Set $\hat{C}_i = Q^{-1}(1 - \epsilon; \lfloor C_i + 1 \rfloor)$, i.e., $Q(\lfloor C_i + 1 \rfloor, \hat{C}_i) = 1 - \epsilon$. Then the constraint $\Lambda_i(m)T_i \leq \hat{C}_i$ is equivalent to $\lim_{m \rightarrow \infty} \mathbb{P}_i(x_i < C_i; m) \geq 1 - \epsilon$.

Linear programming formulation and multi-commodity flow equivalence

In the previous subsection, we show that A-OSCARR collapses into linear functions. In this subsection, we further show that A-OSCARR can be framed as an instance of the well-known multi-commodity flow problem and that A-OSCARR is equivalent to the Congestion-Free Routing and Rebalancing problem presented in Rossi et al. (2018): thus, (i) A-OSCARR can be solved efficiently by ad-hoc algorithms for multi-commodity flow (e.g. Goldberg et al. (1998)) and (ii) the theoretical results presented in Rossi et al. (2018) (namely, the finding that rebalancing trips do not increase congestion) extend, in expectation, to *stochastic* systems.

First, we show that the problem can be solved exclusively for the relative throughputs on the road links, and then we show that the resulting equations are equivalent to a minimum cost multi-commodity flow problem.

The relative throughput going from an intersection i into adjacent roads is $\sum_{j \in \mathcal{W}'_i} \pi_{j,k}$, where $\mathcal{W}'_i = \{\mathcal{W}_i \cap \mathcal{I}\}$ is the set of road links that begin in node i . Similarly, the relative throughput entering the intersection i from the road network is $\sum_{j \in \mathcal{U}'_i} \pi_{j,k}$, where $\mathcal{U}'_i = \{\mathcal{U}_i \cap \mathcal{I}\}$ is the set of road links terminating in i . Additionally, define $d_i^{(k)}$ as the difference between the relative throughput leaving the intersection and the relative throughput entering the intersection. From (7d), (7c), and (18), it can be shown that, for a customer class q at an intersection i , $d_i^{(q)}$ should be equal to the arrival rate if i is adjacent to the source station, the negative arrival rate if it is adjacent to the target station, and 0 otherwise. Formally,

$$\sum_{j \in \mathcal{W}'_i} \pi_{j,q} - \sum_{j \in \mathcal{U}'_i} \pi_{j,q} = d_i^{(q)}, \quad \text{where} \quad d_i^{(q)} = \begin{cases} \lambda^{(q)} & \text{if } i = s^{(q)}, \\ -\lambda^{(q)} & \text{if } i = t^{(q)}, \\ 0 & \text{otherwise.} \end{cases}$$

While the rebalancing arrival rates $\lambda^{(r)}$ are not fixed, we do know from Equation (7c) and from the definition of $d_i^{(q)}$ that $d_{s^{(r)}}^{(r)} = -d_{t^{(r)}}^{(r)}$. Thus,

$$\sum_{j \in \mathcal{W}'_{s^{(r)}}} \pi_{j,r} - \sum_{j \in \mathcal{U}'_{s^{(r)}}} \pi_{j,r} = - \sum_{j \in \mathcal{W}'_{t^{(r)}}} \pi_{j,r} + \sum_{j \in \mathcal{U}'_{t^{(r)}}} \pi_{j,r}.$$

Finally, we can rewrite Lemma 2 as

$$\sum_{q \in \mathcal{Q}} d_i^{(q)} + \sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{W}'_i} \pi_{j,r} - \sum_{j \in \mathcal{U}'_i} \pi_{j,r} = 0.$$

Thus, in the asymptotic regime, Problem (7) can be restated as

$$\underset{\pi_{i \in \mathcal{I}, k \in \mathcal{K}}}{\text{minimize}} \quad \sum_{i \in \mathcal{I}} T_i \sum_{k \in \mathcal{K}} \pi_{i,k},$$

$$\text{subject to} \quad \sum_{q \in \mathcal{Q}} d_i^{(q)} + \sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{W}'_i} \pi_{j,r} - \sum_{j \in \mathcal{U}'_i} \pi_{j,r} = 0 \quad \forall i \in \mathcal{S}, \quad (22a)$$

$$T_i \sum_{k \in \mathcal{K}} \pi_{j,k} \leq \widehat{C}_i \quad \forall i \in \mathcal{I}, \quad (22b)$$

$$\sum_{j \in \mathcal{W}'_i} \pi_{j,q} - \sum_{j \in \mathcal{U}'_i} \pi_{j,q} = d_i^{(q)} \quad \forall i \in \mathcal{S}, \quad (22c)$$

$$\sum_{j \in \mathcal{W}'_{s^{(r)}}} \pi_{j,r} - \sum_{j \in \mathcal{U}'_{s^{(r)}}} \pi_{j,r} = \sum_{j \in \mathcal{W}'_{t^{(r)}}} \pi_{j,r} - \sum_{j \in \mathcal{U}'_{t^{(r)}}} \pi_{j,r} \quad \forall r \in \mathcal{R}, \quad (22d)$$

$$\sum_{j \in \mathcal{W}'_i} \pi_{j,r} - \sum_{j \in \mathcal{U}'_i} \pi_{j,r} = 0 \quad \forall i \in \mathcal{S} \setminus \{s^{(r)}, t^{(r)}\}, \quad (22e)$$

$$\sum_{j \in \mathcal{W}'_{s^{(r)}}} \pi_{j,r} - \sum_{j \in \mathcal{U}'_{s^{(r)}}} \pi_{j,r} \geq 0 \quad \forall r \in \mathcal{R}, \quad (22f)$$

$$\pi_{i,k} \geq 0, \quad \forall i \in \mathcal{I}, k \in \mathcal{K}. \quad (22g)$$

Here, constraints (22a) and (22b) are direct equivalents to (7a) and (7b), respectively. By keeping traffic continuity and equating throughputs at source and target stations, (22c) enforces (7c) and (7d) for the customer classes. For the rebalancing classes, (22d) is equivalent to (7c) and (22e) to (7d). Non-negativity of rebalancing rates (7f) is kept by (22f).

Thus, A-OSCARR can be solved efficiently as a linear program. Note that this formulation is very similar to the multi-commodity flow formulation presented in Rossi et al. (2018). The formulation in this paper prescribes specific routing policies for distinct rebalancing origin-destination pairs, while Rossi et al. (2018) only computes a single “rebalancing flow.” These two formulations, however, are equivalent, as by using a flow decomposition algorithm (Ford and Fulkerson 1962), one can “expand” the single rebalancing flow considered in (Rossi et al. 2018) into a set of rebalancing flows, one for each origin-destination pair. Therefore, it is possible to extend the theoretical results presented in Rossi et al. (2018) to the stochastic setting. Most notably, it is possible to find rebalancing trips that in expectation do not cause more congestion than what would be caused from the same travel demand being satisfied by private vehicles.

Battery Charge Constraints

Plug-in electric vehicles (EVs) are especially suitable to AMoD systems. On the one hand, the type of short-range trips typical of Mobility-on-Demand (and, in the future, AMoD) systems is well-suited to the current generation of range-limited electric vehicles; on the other hand, intelligent policies for rebalancing and charging of EVs can ensure that vehicles with an adequate charge level are available to passengers, greatly reducing “range anxiety”—one of the main barriers to EV adoption (Evarts 2013).

The BCMP framework (and in particular the notion of queuing classes) can be leveraged to model the battery charge level of electric vehicles and efficiently compute coupled rebalancing and charging policies. Accordingly, in this section, we extend the proposed BCMP model to include the constraints imposed by operating a fleet of charge-constrained autonomous EVs.

Casting charge constraints as classes

We discretize the state-of-charge (SOC) of vehicles using a finite set of quantized charge levels. Specifically, we define $\mathcal{B} = \{1, 2, \dots, B\}$ as an ordered set of B charge levels, such that B denotes a full battery and 1 denotes an empty battery. A vehicle entering road $j \in \mathcal{I}$ at charge level $b \in \mathcal{B}$ spends e_j energy in its traversal; thus, the vehicle’s charge level at the end of the road is $b - e_j$. We assume that vehicles spend no energy while idle at the stations, i.e., $e_i = 0$ if $i \in \mathcal{S}$. Note that this model assumes that charge depletion on a road segment is independent of speed/congestion: this approximation is reasonable in our model as (i) we force the transportation network to be (mostly) congestion free and (ii) if a road link is congestion free, it is reasonable to assume that each vehicle travels at the *same* free flow speed for that link. Vehicles can recharge their batteries at a set of plug-in chargers \mathcal{F} : a vehicle can increase the level of charge of its batteries at charger $f \in \mathcal{F}$ by up to e_f levels in time T_f . Analogously to road links, chargers are modeled as IS queues with service rate $1/T_f$.

To track the SOC of individual vehicles, we include the charge level as part of their class. Classes are henceforth denoted by the tuple (k, b) : a vehicle belonging to class (k, b) is servicing request k at charge level b . Analogously, the relative throughput of vehicles serving request k at charge level b on road link i is $\pi_{i,k,b}$. We refer to the total relative throughput on a road link i as

$$\pi_i = \sum_{k \in \mathcal{K}} \sum_{b \in \mathcal{B}} \pi_{i,k,b}.$$

Routing and charging depend on the vehicles’ SOC: for example, energy-depleted vehicles need to charge their batteries before providing service to passengers. Accordingly, rebalancing requests are characterized by a set of tuples that include the initial and final SOC, namely $\mathcal{R} = \{s^{(r)}, t^{(r)}, \lambda^{(r)}, soc_s^{(r)}, soc_t^{(r)}\}$ where $soc_s^{(r)}$ is the initial SOC and $soc_t^{(r)}$ is the final SOC.

For a given station j , the distribution of customer destinations is encoded by the distribution of the transition probabilities $\{p_j^{(k')}\}_{j,k'}$ in Equation (6), for any

class $k' \in \mathcal{Q}$ corresponding to passenger requests. In the closed queuing network model adopted in this paper, customers are assigned to the first available vehicle in the FCFS queue, irrespective of its charge level. Therefore, for classes $k' \in \mathcal{Q}$, the distribution $\{p_j^{(k')}\}_{j,k' \in \mathcal{Q}}$ must be independent of the charge level b of incoming vehicles: if this was not the case, the arrival rate of passengers in a given class would depend on the SOC of vehicles queuing at the station.

On the other hand, the system operator should enforce different rebalancing and charging strategies for non-passenger-carrying vehicles depending on the vehicles' SOC. For example, it may be preferable to charge vehicles that are running out of battery, whereas vehicles with a high charge level may be devoted to rebalancing purposes. To enable this, the distribution $\{p_j^{(k')}\}_{j,k' \in \mathcal{R}}$ for the rebalancing classes $k' \in \mathcal{R}$ is modeled as dependent on the charge level. However, if we were to include the notion of SOC-dependent class assignment directly into the BCMP model from the previous section, we would introduce a *spurious* correlation between the passenger arrivals and the SOC: for example, in the case where vehicles with low charge were to be primarily assigned to rebalancing classes, a station having only low charge vehicles would behave as if there were no passenger arrivals. To address this issue, our strategy is to introduce the notion of stations with *double queues*.

Specifically, each station $i \in \mathcal{S}$ is represented by two vehicle queues, one awaiting passenger requests, indexed with i_Q , and one awaiting rebalancing requests, indexed with i_R . Let $\beta_{i,b}$ be the probability that a vehicle arriving at station i is assigned to the rebalancing queue i_R . Specifically, upon arrival to its destination station $i = t^{(k)}$, a vehicle of class (k, b) proceeds to the rebalancing queue i_R with probability $\beta_{i,b}$ and switches to class (r, b) with probability $\tilde{p}_{i_R,b}^{(r)} = \lambda^{(r)} / \sum_{r' \in \mathcal{O}_{i,b,R}} \lambda^{(r')}$, where $\mathcal{O}_{i,b,R}$ is the set of rebalancing requests r' with $s^{(r')} = i$ and $\text{soc}_s^{(r')} = b$. Conversely, the vehicle proceeds to the passenger queue i_Q with probability $1 - \beta_{i,b}$ and switches to class (q, b) with probability $\tilde{p}_{i_Q}^{(q)} = \lambda^{(q)} / \sum_{q' \in \mathcal{O}_i} \lambda^{(q')}$. Figure 2 shows a graphical depiction of the station model with *double queues*.

The routing policy for both customer-carrying and rebalancing vehicles is allowed to depend on the state of charge: for instance, vehicles with a low SOC may traverse a charging station to recharge their batteries. Accordingly, we let $\alpha^{(k,b)}$ denote the routing policy for a vehicle of class (k, b) (as before, a routing policy is simply a matrix of transition probabilities).

We are now in a position to extend the routing matrix to the setup with battery charge levels. Specifically, let $\mathcal{B}_t^{(k)}$ be the set of *acceptable* SOC at the target station for class k , defined as $\mathcal{B}_t^{(k)} = \{\text{soc}_t^{(k)}\}$ if $k \in \mathcal{R}$, and $\mathcal{B}_t^{(k)} = \mathcal{B}$

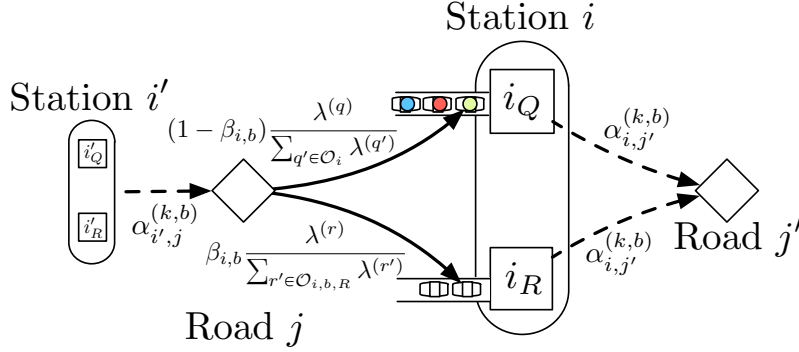


Figure 2. Graphical depiction of a double queue. Vehicles directed to station i at charge level b enter the rebalancing queue i_R with probability $\beta_{i,b}$; they enter the passenger-serving queue i_Q , with probability $1 - \beta_{i,b}$. Note that the probability of joining i_R (or, equivalently, i_Q) depends on the SOC. Vehicles entering i_R switch to a class r , corresponding to a rebalancing request, with probability $\lambda^{(r)} / \sum_{r' \in \mathcal{O}_{i,b,R}} \lambda^{(r')}$. Conversely, vehicles entering i_Q switch to class q , corresponding to a customer request, with probability $\lambda^{(q)} / \sum_{q' \in \mathcal{O}_i} \lambda^{(q')}$.

otherwise. We then define the routing matrix as

$$p_{i,k,b;j,k',b'} = \begin{cases} \alpha_{i,j}^{(k,b)} & \text{if } k = k', \quad j \in \mathcal{W}_{\text{Child}(i)}, \\ & t^{(k)} \notin \mathcal{W}_{\text{Child}(i)}, \quad b' = b - e_i, \\ \beta_{i,b} \tilde{p}_{j_R,b}^{(k')} & \text{if } j = t^{(k)}, \quad t^{(k)} \in \mathcal{W}_{\text{Child}(i)}, \\ & k' \in \mathcal{O}_{j,b',R}, \quad b' = b - e_i \in \mathcal{B}_t^{(k)}, \\ (1 - \beta_{i,b}) \tilde{p}_{j_Q}^{(k')} & \text{if } j = t^{(k)}, t^{(k)} \in \mathcal{W}_{\text{Child}(i)}, \\ & k' \in \mathcal{O}_j, b' = b - e_i \in \mathcal{B}_t^{(k)}, \\ 0 & \text{otherwise.} \end{cases} \quad (23)$$

Problem Formulation

The goal of the problem is to minimize the amount of traffic both on the roads and at the charging stations. Accordingly, we consider the cost function $\sum_{i \in \mathcal{I}'} \Lambda_i(m) T_i$, where \mathcal{I}' represents the union of road link and plug-in charger queues, i.e., $\mathcal{I}' := \{\mathcal{I} \cup \mathcal{F}\}$. In analogy to the OSCARR problem, the system is constrained to remain in balance, i.e., $\gamma_i = \gamma_j$ for all $i, j \in \mathcal{S}$ (note that in this case i and j might denote either a rebalancing or a passenger queue). We refer to the resulting problem as the *Optimal Stochastic Capacitated AMoD Routing, Rebalancing and Charging problem* (OSCARR-C):

$$\begin{aligned}
& \underset{\lambda^{(r \in \mathcal{R})}, \alpha_{ij}^{(k \in \mathcal{K})}}{\text{minimize}} && \sum_{i \in \mathcal{I}'} \Lambda_i(m) T_i, \\
\text{subject to} &&& \gamma_i = \gamma_j, \quad i, j \in \mathcal{S}, && (24a) \\
&&& \Lambda_i(m) T_i \leq C_i, \quad i \in \mathcal{I}', && (24b) \\
&&& \pi_{s^{(q)}, q} = \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{N}} \sum_{b \in \mathcal{B}} \sum_{b' \in \mathcal{B}} \pi_{j, q, b} p_{j, q, b; t^{(q)}, k', b'}, \quad q \in \mathcal{Q}, && (24c) \\
&&& \pi_{s^{(r)}, r, soc_s^{(r)}} = \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{N}} \sum_{b \in \mathcal{B}} \pi_{j, q, b} p_{j, q, b; t^{(r)}, k', b'}, \quad r \in \mathcal{R}, b' = soc_t^{(r)}, && (24d) \\
&&& \pi_{i, k, b} = \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{N}} \sum_{b' \in \mathcal{B}} \pi_{j, k', b'} p_{j, k', b'; i, k, b}, \quad i \in \mathcal{N}, k \in \mathcal{K}, b \in \mathcal{B}, && (24e) \\
&&& \sum_{j \in \mathcal{W}_{\text{Child}(i)}} \alpha_{ij}^{(k, b)} = 1, \quad \alpha_{ij}^{(k, b)} \geq 0, \quad i, j \in \mathcal{N}, k \in \mathcal{K}, b \in \mathcal{B} && (24f) \\
&&& \lambda^{(r)} \geq 0, \quad r \in \mathcal{R}. && (24g)
\end{aligned}$$

Constraints (24c) and (24d) ensure that passengers' and rebalancing traffic leaving a source reaches the corresponding destination, and, in the case of rebalancing, a desired charge level. In analogy to (7d), (24e) enforces traffic continuity at the road and charge level. Constraint (24f) ensures that $\alpha_{ij}^{(k, b)}$ is a valid probability measure, and (24g) limits rebalancing requests to positive values.

The setups of OSCARR-C and of OSCARR are very similar. Specifically, the modeling assumptions required to derive the asymptotically optimal formulation of OSCARR are also valid for OSCARR-C; namely, a closed, multi-class network with IS and SS queues constrained to maintain equal relative availability across the SS queues. Therefore, in analogy with (18), the desired relative utilizations can be set to $\gamma_i = 1$ for all $i \in \mathcal{S}$ without loss of generality. As a result, the SS relative throughputs are constrained to equal their service rates, that is,

$$\pi_i = \tilde{\lambda}_i, \quad \text{for all } i \in \mathcal{S}. \quad (25)$$

Additionally, Lemma 3 is still valid. Indeed, the queues \mathcal{I}' in the objective function are all IS queues. Furthermore, since for a feasible solution to OSCARR-C Equation (25) holds, it follows that $\gamma_i = \gamma_S^{\max} = 1$ for all $i \in \mathcal{S}$, and, thus, all the SS queues in \mathcal{S} are bottleneck queues. Therefore, from (4), $\lim_{m \rightarrow \infty} \frac{G(m-1)}{G(m)} = 1$. Thus, the two assumptions for Lemma 3 are verified. As a consequence, according to Lemma 3, the optimal solution $\{\hat{\pi}_i\}$ to the asymptotic approximation of OSCARR-C is also optimal for the full OSCARR-C as the fleet size goes to infinity. We define the problem of finding an asymptotically optimal solution to OSCARR-C as the *Asymptotically Optimal Stochastic Capacitated*

AMoD Routing, Rebalancing and Charging problem, or A-OSCARR-C. We next discuss how to solve A-OSCARR-C.

The key strategy is to rewrite constraints (24a),(24c),(24d),(24e), and (24f) as traffic conservation constraints, both at intersections and at stations. To this purpose, we first find a traffic conservation constraint in terms of both the relative throughputs and the transition probabilities. We then derive traffic conservation constraints only in terms of relative throughputs, both at intersections and at stations – these are the constraints that will be used to set up A-OSCARR-C (46). In the following paragraphs, we provide and rigorously derive a tractable formulation of A-OSCARR-C based on the aforementioned strategy.

Traffic conservation in terms of relative throughputs and transition probabilities: Denote the relative throughput of vehicles arriving at station i_R with charge level b as $\pi_{i_R,b}$. In other words, according to the routing matrix (23), $\pi_{i_R,b} = \beta_{i,b} \sum_{k \in \mathcal{K}} \pi_{i,k,b}$. Similarly, denote the relative throughput of vehicles arriving at station i_Q with charge level b as $\pi_{i_Q,b}$. In other words, $\pi_{i_Q,b} = (1 - \beta_{i,b}) \sum_{k \in \mathcal{K}} \pi_{i,k,b}$. For bookkeeping purposes, we define the combined relative throughput as $\hat{\pi}_{i,b} := \pi_{i_R,b} + \pi_{i_Q,b}$. From (25), the relative throughput for the passenger queues must equal the rate of arrival, that is:

$$\begin{aligned}
\sum_{q \in \mathcal{O}_i} \lambda^{(q)} &= \sum_{q \in \mathcal{Q}} \sum_{b \in \mathcal{B}} (1 - \beta_{i,b}) \hat{\pi}_{i,b} \tilde{p}_i^{(q)} \\
&= \sum_{q \in \mathcal{Q}} \tilde{p}_i^{(q)} \sum_{b \in \mathcal{B}} (1 - \beta_{i,b}) \hat{\pi}_{i,b} \\
&= \sum_{b \in \mathcal{B}} (1 - \beta_{i,b}) \hat{\pi}_{i,b} \\
&= \sum_{b \in \mathcal{B}} \hat{\pi}_{i,b} - \pi_{i_R,b}.
\end{aligned} \tag{26}$$

We now turn our attention to the traffic equations (24e). Note that, as per (23), the only queues that contribute relative throughput into a queue $i \in \mathcal{I}'$ are the set of queues $\mathcal{U}_{\text{Parent}(i)}$ that feed into the parent intersection of i . Thus, for queues $i \in \mathcal{I}'$

$$\begin{aligned}
\pi_{i,k,b} &= \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{N}} \sum_{b' \in \mathcal{B}} \pi_{j,k',b'} p_{j,k',b';i,k,b} \\
&= \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{U}_{\text{Parent}(i)}} \sum_{b' \in \mathcal{B}} \pi_{j,k',b'} p_{j,k',b';i,k,b}.
\end{aligned} \tag{27}$$

Moreover, $p_{j,k',b';i,k,b} \neq 0$ only for queues j which feed into the parent intersection of i and for charge levels b' such that $b' = b + e_j$ (note that $e_j = 0$,

if $j \in \mathcal{S}$). Therefore, for queues $i \in \mathcal{I}'$

$$\begin{aligned} \pi_{i,k,b} &= \sum_{k' \in \mathcal{K}} \sum_{j \in \mathcal{U}_{\text{Parent}(i)}} \sum_{b' \in \mathcal{B}} \pi_{j,k',b'} p_{j,k',b';i,k,b} \\ &= \sum_{j \in \mathcal{U}_{\text{Parent}(i)}} \pi_{j,k,b+e_j} \alpha_{j,i}^{(k,b+e_j)}. \end{aligned} \quad (28)$$

Let $\mathcal{W}'_{\text{Parent}(i)} = \{\mathcal{W}_{\text{Parent}(i)} \cap \mathcal{I}'\}$, $\mathcal{U}'_{\text{Parent}(i)} = \{\mathcal{U}'_{\text{Parent}(i)} \cap \mathcal{I}'\}$ and $l = \text{Parent}(i)$. For a road or charger queue $i \in \mathcal{I}'$ such that $l = \text{Parent}(i)$, we can then simplify (24e) and obtain the traffic conservation equation

$$\begin{aligned} \pi_{i,k,b} &= \sum_{j \in \mathcal{U}'_i} \pi_{j,k,b+e_j} \alpha_{j,i}^{(k,b)} \\ &= \begin{cases} \sum_{j \in \mathcal{U}'_i} \pi_{j,k,b+e_j} \alpha_{j,i}^{(k,b+e_j)}, & \text{if } s^{(k)} \notin \mathcal{U}'_i, \\ \sum_{j \in \mathcal{U}'_i} \pi_{j,k,b+e_j} \alpha_{j,i}^{(k,b+e_j)} + \pi_{s^{(k)},k,b} \alpha_{j,i}^{(k,b)}, & \text{if } s^{(k)} \in \mathcal{U}_i, k \in \mathcal{Q}, \\ \sum_{j \in \mathcal{U}'_i} \pi_{j,k,b+e_j} \alpha_{j,i}^{(k,b+e_j)} + \pi_{s^{(k)},k,b} \alpha_{j,i}^{(k,b)}, & \text{if } s^{(k)} \in \mathcal{U}_i, b = \text{soc}_s^{(k)}, k \in \mathcal{R}, \end{cases} \end{aligned} \quad (29)$$

where the second equality exploits the fact that only the source station of class k sends vehicles of class k into the road network.

Traffic conservation at intersections in terms of relative throughputs only: We are now in a position to derive a number of traffic conservation constraints at intersections, in terms of only relative throughputs.

Case 1: if an intersection l is not adjacent to either the source or target of a passenger class $q \in \mathcal{Q}$ (that is, $l \notin \{\text{Child}(s^{(q)}), \text{Parent}(t^{(q)})\}$ for $q \in \mathcal{Q}$), then $\sum_{j' \in \mathcal{W}'_i} \alpha_{jj'}^{(q,b)} = 1$. Thus, we can sum (29) over all queues to obtain the following traffic conservation equations for class $q \in \mathcal{Q}$ in terms of the relative throughput:

$$\sum_{j \in \mathcal{W}'_i} \pi_{j,q,b} = \sum_{j \in \mathcal{U}'_i} \pi_{j,q,b+e_j}, \quad \text{for } l \notin \{\text{Child}(s^{(q)}), \text{Parent}(t^{(q)})\}, \quad q \in \mathcal{Q}. \quad (30)$$

Case 2: for a rebalancing class $r \in \mathcal{R}$, if neither $\{l = \text{Child}(s^{(r)}) \wedge b = \text{soc}_s^{(r)}\}$ nor $\{l = \text{Parent}(t^{(r)}) \wedge b = \text{soc}_t^{(r)}\}$, then $\sum_{j' \in \mathcal{W}'_i} \alpha_{jj'}^{(r,b)} = 1$. Thus, we can obtain the following traffic conservation equations for class $r \in \mathcal{R}$ in terms of the relative throughput

$$\begin{aligned} \sum_{j \in \mathcal{W}'_i} \pi_{j,r,b} &= \sum_{j \in \mathcal{U}'_i} \pi_{j,r,b+e_j} \quad \text{if } \neg\{l = \text{Child}(s^{(r)}) \wedge b = \text{soc}_s^{(r)}\} \\ &\quad \text{and } \neg\{l = \text{Parent}(t^{(r)}) \wedge b = \text{soc}_t^{(r)}\}, \quad r \in \mathcal{R}. \end{aligned} \quad (31)$$

Case 3: If l is adjacent to the source station, then summing (29) over all queues, one obtains

$$\sum_{j \in \mathcal{W}'_i} \pi_{j,k,b} = \sum_{j \in \mathcal{U}'_i} \pi_{j,k,b+e_j} + \pi_{i,k,b}, \quad i = s^{(k)}, \quad l = \text{Child}(s^{(k)}). \quad (32)$$

If k is a passenger class $q \in \mathcal{Q}$, then $\pi_{i,k,b} = (1 - \beta_{i,b})\hat{\pi}_{i,b}\tilde{p}_i^q$, and one can derive the traffic conservation constraint

$$\begin{aligned}
\sum_{j \in \mathcal{W}'_l} \pi_{j,q,b} - \sum_{j \in \mathcal{U}'_l} \pi_{j,q,b+e_j} &= (1 - \beta_{i,b})\hat{\pi}_{i,b}\tilde{p}_i^q, \\
&= (\hat{\pi}_{i,b} - \pi_{i_R,b})\tilde{p}_i^q, \\
&= (\hat{\pi}_{i,b} - \pi_{i_R,b}) \frac{\lambda^q}{\sum_{q \in \mathcal{O}_i} \lambda^{(q)}}, \\
&= \frac{(\hat{\pi}_{i,b} - \pi_{i_R,b})}{\sum_{b \in \mathcal{B}} \hat{\pi}_{i,b} - \pi_{i_R,b}} \lambda^{(q)}, \\
&= \varrho_{i_Q,b} \lambda^{(q)}, \quad \text{if } i = s^{(q)}, q \in \mathcal{Q}, \quad l = \text{Child}(s^{(q)}),
\end{aligned} \tag{33}$$

where the fourth equality follows from (26). In the fifth equality, we denote $\varrho_{i_Q,b}$ as the ratio of the relative throughput that goes through queue i_Q at charge level b , that is $\varrho_{i_Q,b} = \pi_{i_Q,b}/\pi_{i_Q}$. We treat this ratio as a decision variable: intuitively, $\varrho_{i_Q,b}$ controls the charge level distribution of the vehicles available for passenger use.

If, instead, $k = r \in \mathcal{R}$, l is adjacent to the source station, and b is the desired target charge level, then $\pi_{i,r,b} = \lambda^{(r)}$. Therefore one obtains the traffic conservation constraint

$$\sum_{j \in \mathcal{W}'_l} \pi_{j,k,b} = \sum_{j \in \mathcal{U}'_l} \pi_{j,k,b+e_j} + \lambda^{(r)}, \quad b = \text{soc}_s^{(r)}, l = \text{Child}(s^{(r)}), \quad r \in \mathcal{R}. \tag{34}$$

Case 4: If the intersection l is adjacent to the target station $t^{(k)}$ of a class k (either in \mathcal{Q} or \mathcal{R}), then $\sum_{j' \in \mathcal{W}'_l} \alpha_{j,j'}^{(k,b)} \neq 1$ in general. Let $\zeta_{t^{(k)},k,b}$ be the relative throughput of class k that enters station $t^{(k)}$ from adjacent road and charger queues at charge level b . Then, the sum of (29) must satisfy

$$\sum_{j \in \mathcal{W}'_l} \pi_{j,q,b} + \zeta_{t^{(k)},k,b} = \sum_{j \in \mathcal{U}'_l} \pi_{j,q,b+e_j}, \quad \text{if } l = \text{Parent}(t^{(k)}). \tag{35}$$

However, we know from (24c) that the total relative throughput of a passenger class $q \in \mathcal{Q}$ entering its target station $t^{(q)}$ must equal the the total relative throughput at its source station, i.e., $\pi_{s^{(q)},q} = \sum_{b \in \mathcal{B}} \zeta_{t^{(q)},q,b}$. Additionally, if (25) holds, it can be shown that $\pi_{s^{(q)},q} = \lambda^{(q)}$. Thus, summing (35) over all charge levels for passenger class $q \in \mathcal{Q}$, we obtain the traffic conservation constraint

$$\sum_{b \in \mathcal{B}} \sum_{j \in \mathcal{U}'_l} \pi_{j,q,b+e_j} - \sum_{j \in \mathcal{W}'_l} \pi_{j,q,b} = \lambda^{(q)}, \quad \text{if } l = \text{Parent}(t^{(q)}), \quad q \in \mathcal{Q}. \tag{36}$$

Alternatively, equation (24d) enforces that the total relative throughput of a rebalancing class $r \in \mathcal{R}$ entering its target station $t^{(r)}$ with charge level $\text{soc}_t^{(r)}$ be equal to the relative throughput leaving its source $s^{(r)}$ with charge level $\text{soc}_s^{(r)}$.

Thus, $\pi_{s^{(r)},r,soc_s^{(r)}} = \zeta_{t^{(r)},r,b} = \lambda^{(r)}$, and the sum of (29) must satisfy

$$\sum_{\mathcal{U}'_l} \pi_{j,r,b+e_j} - \sum_{j \in \mathcal{W}'_l} \pi_{j,r,b} = \lambda^{(r)}, \quad \text{if } l = \text{Parent}(t^{(r)}), b = soc_t^{(r)}, r \in \mathcal{R}. \quad (37)$$

Note that both the arrival rates and the relative throughputs of the rebalancing classes are decision variables. The rebalancing rates are not explicitly represented in the optimization problem: rather, they are implicitly set by equating (34) and (37):

$$\begin{aligned} \sum_{\mathcal{W}'_{\text{Child}(s^{(r)})}} \pi_{j,r,soc_s^{(r)}} - \sum_{j \in \mathcal{U}'_{\text{Child}(s^{(r)})}} \pi_{j,r,soc_s^{(r)}+e_j} = \\ \sum_{j \in \mathcal{U}'_{\text{Parent}(t^{(r)})}} \pi_{j,r,soc_t^{(r)}+e_j} - \sum_{\mathcal{W}'_{\text{Parent}(t^{(r)})}} \pi_{j,r,soc_t^{(r)}}, \quad (38) \\ r \in \mathcal{R}, \end{aligned}$$

and by setting a positive constraint on the relative throughputs at the source station, i.e.

$$\sum_{\mathcal{W}'_{\text{Child}(s^{(r)})}} \pi_{j,r,soc_s^{(r)}} - \sum_{j \in \mathcal{U}'_{\text{Child}(s^{(r)})}} \pi_{j,r,soc_s^{(r)}+e_j} > 0, r \in \mathcal{R}. \quad (39)$$

Traffic conservation at stations in terms of relative throughputs only. Finally, we characterize a traffic conservation constraint for the stations (namely, Equation (45)), by manipulating the traffic equations (24e). Specifically, recall that, for a given station $i \in \mathcal{S}$, $\hat{\pi}_{i,b} = \pi_{i_R,b} + \pi_{i_Q,b}$. Due to Equation (25), $\pi_{i_R,b} = \sum_{r \in \mathcal{O}_{i,b,R}} \lambda^{(r)}$ and $\pi_{i_Q,b} = \sum_{q \in \mathcal{O}_i} \varrho_{i_Q,b} \lambda^{(q)}$. Therefore,

$$\hat{\pi}_{i,b} = \sum_{r \in \mathcal{O}_{i,b,R}} \lambda^{(r)} + \sum_{q \in \mathcal{O}_i} \varrho_{i_Q,b} \lambda^{(q)}. \quad (40)$$

Conversely, $\hat{\pi}_{i,b}$ must equal the sum of passenger and rebalancing traffic that enters station $i \in \mathcal{S}$ at charge level b . In particular, the rebalancing traffic entering station i at charge level b is $\sum_{r \in \mathcal{D}_{i,b,R}} \lambda^{(r)}$, where $\mathcal{D}_{i,b,R}$ is the set of rebalancing classes whose destination is $i \in \mathcal{S}$ with target charge level b . From (36), we see that for a single charge level the relative throughput of a passenger class $q \in \mathcal{Q}$ entering $i \in \mathcal{S}$ is $\sum_{\mathcal{U}'_l} \pi_{j,q,b+e_j} - \sum_{j \in \mathcal{W}'_l} \pi_{j,q,b}$, where $l = \text{Parent}(i)$. Summing this over all passenger classes and adding the rebalancing traffic we obtain

$$\hat{\pi}_{i,b} = \sum_{r \in \mathcal{D}_{i,b,R}} \lambda^{(r)} + \sum_{q \in \mathcal{D}_{i,Q}} \sum_{\mathcal{U}'_l} \pi_{j,q,b+e_j} - \sum_{j \in \mathcal{W}'_l} \pi_{j,q,b}. \quad (41)$$

Together, (40) and (41) imply that

$$\sum_{r \in \mathcal{O}_{i,b,R}} \lambda^{(r)} + \sum_{q \in \mathcal{O}_i} \varrho_{i_Q,b} \lambda^{(q)} = \sum_{r \in \mathcal{D}_{i,b,R}} \lambda^{(r)} + \sum_{q \in \mathcal{D}_{i,Q}} \sum_{\mathcal{U}'_l} \pi_{j,q,b+e_j} - \sum_{j \in \mathcal{W}'_l} \pi_{j,q,b}. \quad (42)$$

Note that $\sum_{\mathcal{W}'_l} \pi_{j,r,b} - \sum_{j \in \mathcal{U}'_l} \pi_{j,r,b+e_j}$ equals $\lambda^{(r)}$ at its source station, $-\lambda^{(r)}$ at the target station, and 0 otherwise. Therefore, we can express $\lambda^{(r)}$ in terms of the relative throughputs

$$\sum_{\mathcal{W}'_l} \pi_{j,r,b} - \sum_{j \in \mathcal{U}'_l} \pi_{j,r,b+e_j} = \begin{cases} \lambda^{(r)}, & \text{if } l = \text{Child}(s^{(r)}), \\ -\lambda^{(r)}, & \text{if } l = \text{Parent}(t^{(r)}), \\ 0, & \text{otherwise.} \end{cases} \quad (43)$$

The difference between incoming and departing rebalancing relative throughput at a station now becomes

$$\sum_{r \in \mathcal{R}} \sum_{\mathcal{W}'_l} \pi_{j,r,b} - \sum_{j \in \mathcal{U}'_l} \pi_{j,r,b+e_j} = \sum_{r \in \mathcal{O}_{i,b,R}} \lambda^{(r)} - \sum_{r \in \mathcal{D}_{i,b,R}} \lambda^{(r)}. \quad (44)$$

Thus, by rewriting (42), we obtain the traffic conservation constraint at each station $i \in \mathcal{S}$

$$\begin{aligned} \sum_{q \in \mathcal{O}_{i,Q}} \varrho_{i_Q,b} \lambda^{(q)} + \sum_{r \in \mathcal{R}} \sum_{\mathcal{W}'_l} \pi_{j,r,b} - \sum_{j \in \mathcal{U}'_l} \pi_{j,r,b+e_j} \\ - \sum_{q \in \mathcal{D}_{i,Q}} \sum_{\mathcal{U}'_l} \pi_{j,q,b+e_j} - \sum_{j \in \mathcal{W}'_l} \pi_{j,q,b} = 0, \quad l = \text{Parent}(i) = \text{Child}(i). \end{aligned} \quad (45)$$

Collecting all the results above, A-OSCARR-C can now be framed in terms of the relative throughputs $\{\pi_{i,k,b}\}_{i,k,b}$ and the ratios $\{\varrho_{i_Q,b}\}_{i,b}$:

$$\begin{aligned} & \text{minimize} && \sum_{i \in \mathcal{I}'} T_i \sum_{k \in \mathcal{K}} \sum_{b \in \mathcal{B}} \pi_{i,k,b}, \\ & \text{subject to} && (30), (31), (33), (36), (38), (39), (45) \\ & && T_i \sum_{k \in \mathcal{K}} \sum_{b \in \mathcal{B}} \pi_{i,k,b} \leq \widehat{C}_i, \quad \forall i \in \mathcal{I}, \quad (46a) \\ & && \pi_{i,k,b} \geq 0, \quad \forall i \in \mathcal{I}, k \in \mathcal{K}, b \in \mathcal{B} \quad (46b) \\ & && \sum_{b \in \mathcal{B}} \varrho_{i_Q,b} = 1, \quad \forall i \in \mathcal{S}, \quad (46c) \\ & && \varrho_{i_Q,b} \geq 0 \quad i \in \mathcal{S}, b \in \mathcal{B}. \quad (46d) \end{aligned}$$

Constraints (30)-(39) enforce consistency in the model. Constraint (45) enforces conservation of traffic at each charging level and, consequently, equal availabilities at each station. (46a) sets the bounds on the expected traffic at the road and charger queues, and (46b) enforces non-negative traffic values. Finally, constraints (46c) and (46d) make sure that the ratios $\{\varrho_{i,b}\}$ are a valid probability measure.

As in the non-battery case, A-OSCARR-C can be solved as a linear program. The state size is $O(|\mathcal{I}'| |\mathcal{B}| (|\mathcal{R}| + |\mathcal{Q}|))$. $|\mathcal{Q}|$ grows quadratically with $|\mathcal{S}|$, and

$|\mathcal{R}|$ grows quadratically with $|\mathcal{S}||\mathcal{B}|$; thus, for large road networks, the problem size can become unwieldy even for modern linear programming algorithms. For instance, for a road network with 350 nodes, 1000 road segments and 60 charge levels, the problem size is slightly above 9 billion variables; for comparison, state-of-the-art LP solvers can reliably handle problems with tens of millions of variables on modern machines (Mittelman 2016). Remarkably, it is possible to reduce the problem size, with no loss of information, by addressing A-OSCARR-C as an augmented network flow problem, bundling customer traffic demands according to their source node, and collecting all rebalancing demands into a single class. Under this approach, the same problem instance would be reduced to just over 3 million variables. For a given set of optimal customer and rebalancing flows, individual routes can be recovered using a flow decomposition algorithm (Ford and Fulkerson 1962), in analogy with A-OSCARR. We refer the reader to (Rossi et al. 2018) for a thorough discussion.

Numerical Experiments

To illustrate a real-life application of the models and methods presented in this paper, we performed a case study of Manhattan, where system performance metrics were computed as a function of fleet size using Mean Value Analysis (Gelenbe et al. 1998). The road network model used for this case study consists of a subset of Manhattan’s real road network (shown in Figure 3), with 1,005 road links and 357 intersections. To select station positions and compute the rates $\lambda^{(q)}$ (for each tuple $q \in \mathcal{Q}$ modeling the arrival process) we used the taxi trips within Manhattan that took place between 7:00AM and 8:00AM on March 1, 2012 (22,416 trips) from the New York City Taxi and Limousine Commission dataset[†]. We clustered the pickup locations into 50 different groups with K -means clustering and placed a station at the road intersection closest to each cluster centroid. We then fit an origin-destination model with exponential distributions to describe the customer trip demands between the stations. In order to observe congestion effects, road capacities were reduced (specifically, by 55%) to ensure that maximum road utilization is achieved on some of the road links; in the real world, an analogous reduction in road capacity would be caused by traffic exogenous to the taxi system.

[†]http://www.nyc.gov/html/tlc/html/about/trip_record_data.shtml

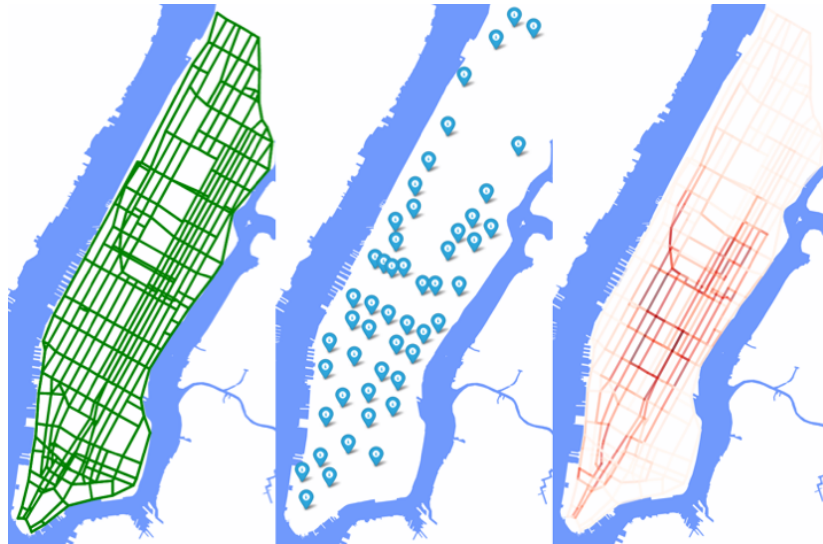


Figure 3. Manhattan scenario. Left: modeled road network. Center: Station locations. Right: Resulting vehicular flow (darker flows show higher vehicular presence).

Routing and rebalancing under congestion constraints

We considered two scenarios: (i) the “baseline” scenario where traffic constraints on each road link are based on expectation, i.e., the average number of vehicles on a road link is below its nominal capacity; and (ii) the “conservative” scenario where the constraints are based on the asymptotic probability of exceeding the nominal capacity (specifically, the asymptotic probability of exceeding the nominal capacity is constrained to be lower than 10%). Figure 3 shows the road network, the station locations, and the resulting traffic flow, and Figure 4 shows the results.

We see from Figure 4a that, as intended, the station availabilities are balanced and approach one as the fleet size increases. However, Figure 4b shows that there is a trade off between availability and vehicle utilization. For example, for a fleet size of 4,000 vehicles, on average, half of the vehicles are waiting at the stations. In contrast, a fleet of 2,400 vehicles results in availability of 91% and only 516 vehicles (in expectations) wait at the stations. Not shown in the figures, 34% of the trips are for rebalancing purposes; in contrast, only about 18% of the traveling vehicles are rebalancing. This shows that rebalancing trips are significantly shorter than passenger trips, which is in line with the goal of minimizing the number of empty vehicles on the road and thus road congestion.

Figures 4a and 4b show only the results for the baseline case; for the conservative scenario, the difference in availabilities is less than 0.1%, and the difference in the total expected number of vehicles on the road is less than 7, regardless of fleet size. However, road utilization is significantly different in the two scenarios we considered. In Figure 4c, we see that, as the fleet size increases, the likelihood of exceeding the nominal capacity approaches 50%. In contrast,

in the conservative scenario, the probability of exceeding the capacity is never more than 10% –by design– regardless of fleet size.

Finally, we verified the validity of the load-independent travel time assumption. Assuming asymptotic conditions (in which case the number of vehicles on each road follows a Poisson distribution), we computed for both scenarios the expected travel time between each origin-destination pair by using the Bureau of Public Roads (BPR) delay model (Bureau of Public Roads 1964), and estimated the difference with respect to the load-independent travel time used in this paper. The BPR delay model is a commonly used equation for relating traffic to travel time (Bureau of Public Roads 1964). Under this model, the travel time on a road link is given by

$$T'_i = T_i \left(1 + \delta \left(\frac{x_i}{C_i} \right)^\beta \right), \quad (47)$$

where T'_i is the real mean travel time, T_i is the free flow travel time, x_i is the number of vehicles on the road, C_i is the nominal capacity of the road, and δ and β are parameters usually set to 0.15 and 3, respectively. The results, depicted in Figure 4d, show that the maximum difference for the baseline and conservative scenarios are an increase of around 8% and 4%, respectively, and the difference tends to be smaller for higher trip times. Thus, for this specific case study, our assumption is reasonable.

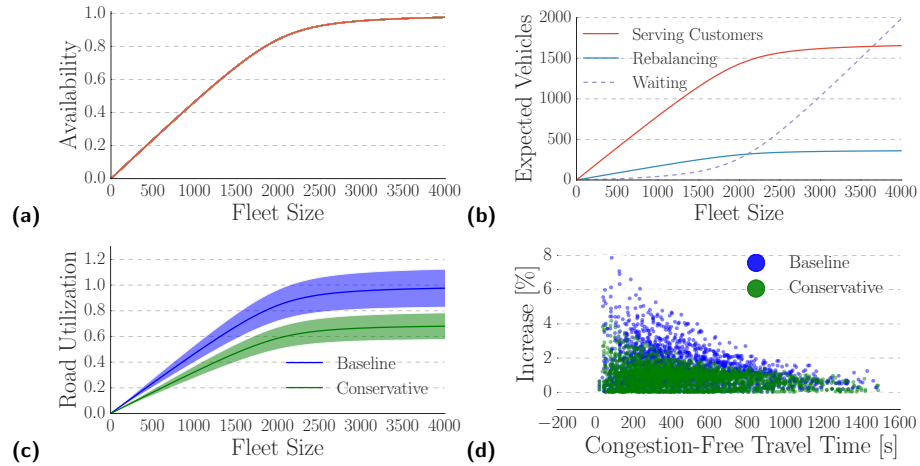


Figure 4. (a) Station availabilities as a function of fleet size for the baseline case. (b) Expected number of vehicles by usage as a function of fleet size for the baseline case. (c) Utilization as a function of fleet size for the most utilized road. The colored band denotes ± 1 standard deviation from the mean. (d) Increase in expected travel time for each O-D pair when considering the BPR delay model.

Inclusion of charging constrains

To study the behavior of the Manhattan scenario under charging constraints, we made some additional assumptions. We assumed a vehicle battery size of 8 kWh and a full charge range of 50km, specifications that are similar to Toyota's iRoad urban mobility vehicle[‡]. We assumed that energy consumption depends exclusively on distance and computed the road energy costs, e_j , using the road lengths and battery range. Additionally, we assumed that every station is equipped with chargers capable of delivering 75kW of power per vehicle, comparable to the superchargers offered by Tesla (Tesla Motors 2017), and did not enforce a limit on the number of vehicles that can charge simultaneously. Battery capacity was discretized into 60 discrete charge levels, a number that showed a good trade-off between accuracy in the energy cost at the roads and the problem size. Finally, in order to discourage recharging of customer-carrying vehicles while penalizing needless rebalancing of customer-empty vehicles, we imposed a higher cost per unit of time for passenger-carrying trips. Currently, the average hourly salary in the United States is \$26.19[§], and the hourly cost of a rental car with the ZipCar car-sharing service starts at \$7[¶]. Assuming that the cost per unit time of a vehicle is \$7/hr and the value of time of a passenger is \$26.19/hr, then the overall value of time of a passenger-carrying vehicle is approximately \$33/hr, or five time higher than the cost of an empty, or rebalancing, vehicle. Therefore, we weigh passenger traffic as being five times costlier than rebalancing traffic.

Figures 5a and 5b show our results. As expected, the algorithm does not route any of the passenger classes through the chargers, that is, a customer should not expect to spend time charging or worry about the range of the vehicle. Conversely, vehicles entering a station at lower charge levels are overwhelmingly devoted to rebalancing tasks. For example, in Figure 5, we see the throughput distribution, as the fleet size approaches infinity, for the station with the highest number of rebalancing requests (5a) and for the station with the highest passenger arrival rate (5b). As intended, all the charging is done by rebalancing vehicles. This explains the rebalancing spike at lower charge levels: the optimizer assigns vehicles with lower energy levels exclusively to serve rebalancing tasks. In contrast, the algorithm utilizes vehicles at the highest charge levels exclusively for passenger requests. By satisfying all charging requirements with empty vehicles, the proposed approach is able to successfully mitigate the risk of range anxiety.

[‡]http://www.toyota-global.com/innovation/personal_mobility/i-road/

[§]<https://www.bls.gov/news.release/empsit.t19.htm>

[¶]<http://www.zipcar.com/check-rates/sf>

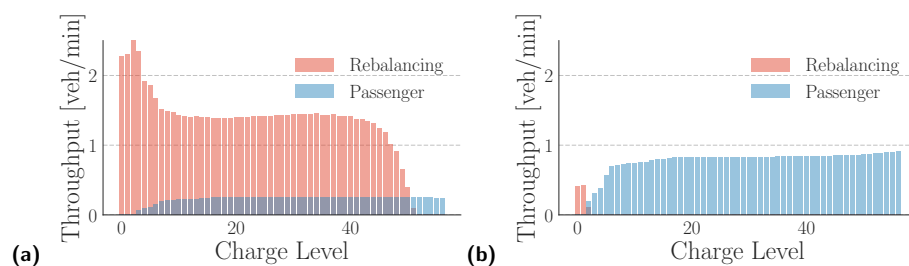


Figure 5. Passenger and rebalancing request throughput distribution for (a) the station with the highest rebalancing rate, and (b) the station with highest number of passenger requests.

The impact of the charging constraints on availability is notable. In order to maintain 91% availability, the fleet should comprise 3,300 vehicles, a 37.5% increase over the scenario without charging. A fleet this large has an average of 384 vehicles charging at any given time. The fleet may also have a significant effect on the electric power system. The range of charger utilization varies significantly by location. The power draw of the most used charger station is between 1.1 and 2.2 MW 90% of the time, while the least utilized charger is between 0 and 300 kW. While the smaller charging stations could be well served by current charging station standards, like Tesla Supercharger stations which offer 145kW per charger for two cars, the larger stations will likely require greater coordination with the local power authorities. In total, the expected power consumption in the system is around 30MW.

Discussion

The two previous experiments showcase the modeling power of the proposed framework. In particular, the framework enables future practitioners to couple key modeling features, such as congestion and charging, to stochastic performance metrics such as availability, and, thus, to synthesize control policies. However, it is important to highlight some limitations. First, the problem formulations OSCARR and OSCARR-C are not always feasible. Most notably, travel demand might exceed road capacity. In such cases, some available options are to evaluate whether to increase the threshold in some of the road links (and reduce the freeflow speed accordingly), or whether to relax the problem by including slack variables that penalize capacity violation. Second, the numerical experiments presented rely on steady state analysis of the AMoD system with a fixed freeflow in the road links. A further study should evaluate its merits by comparing against microscopic simulations such as MATsim (Balmer et al. 2009).

Conclusions

In this paper we presented a novel queuing theoretical framework for modeling AMoD systems. We showed that, for the routing and rebalancing problem, the

stochastic model we propose asymptotically recovers existing models based on network flow approximations. The model enables the analysis and control of the *probabilistic distribution* of the vehicles, as opposed to just expected values. In particular, this model allows one to set arbitrary bounds on the asymptotic probability of exceeding the capacity of individual road links. The model is very expressive and can capture both congestion and the charge level of electric vehicles servicing the customers. As such, it can be used to synthesize routing, rebalancing and charging control policies for AMoD fleets with electric vehicles and stochastic demand.

The flexibility of the model presented will be further exploited in future work. First, we would like to incorporate a more accurate congestion model, using load-dependent IS queues as roads, in order to study heavily congested scenarios. Second, we currently consider the system in isolation from other transportation modes, whereas, in reality, customer demand depends on the perceived quality of the different transportation alternatives. Future research will explore the effect of AMoD systems on customer behavior and how to optimally integrate fleets of self-driving vehicles with existing public transit. Third, we would like to further explore the couplings that might arise between the charging policies of an electric-powered AMoD fleet and the electric grid. Of particular interest is the potential participation of an electric-powered AMoD system in the ancillary services market of the power grid. Fourth, the current model assumes that each customer travels alone: future research will address the problem of *ride-sharing*, where multiple customers may share the same vehicle. Lastly, the control policy proposed in this paper is open-loop and thus sensitive to modeling errors (e.g., incorrect estimation of customer demand). Future research will characterize the stability, persistent feasibility and performance of *real-time, closed-loop* model predictive control schemes based on a receding-horizon implementation of the routing policies presented in this paper.

Acknowledgements

This research was supported by the National Science Foundation under CAREER Award CMMI-1454737 and by the Toyota Research Institute (TRI).

References

- Acquaviva F, Di Paola D and Rizzo A (2014) A novel formulation for the distributed solution of load balancing problems in mobility on-demand systems. In: *Connected Vehicles and Expo (ICCVE), 2014 International Conference on*. pp. 906–911.
- Balmer M, Rieser M, Meister K, Charypar D, Lefebvre N and Nagel K (2009) MATSim-t: Architecture and simulation times. In: *Multi-Agent Systems for Traffic and Transportation Engineering*, chapter 3.
- Banerjee S, Johari R and Riquelme C (2015) Pricing in ride-sharing platforms: A queueing-theoretic approach. In: *Proceedings of the Sixteenth ACM Conference on Economics and Computation*. ACM, pp. 639–639.

- Baskett F, Chandy KM, Muntz RR and Palacios FG (1975) Open, closed, and mixed networks of queues with different classes of customers. *Journal of the Association for Computing Machinery* 22(2): 248–260.
- Boyacı B, Zografos KG and Geroliminis N (2015) An optimization framework for the development of efficient one-way car-sharing systems. *European Journal of Operational Research* 240(3): 718–733.
- Bureau of Public Roads (1964) Traffic assignment manual. Technical report, U.S. Department of Commerce, Urban Planning Division, Washington, D.C (1964).
- Chemla D, Meunier F and Calvo RW (2013) Bike sharing systems: Solving the static rebalancing problem. *Discrete Optimization* 10(2): 120–146.
- Chen TD, Kockelman KM and Hanna JP (2016) Operations of a shared, autonomous, electric vehicle fleet: Implications of vehicle & charging infrastructure decisions. *Transportation Research Part A: Policy and Practice* 94: 243–254.
- Chiu YC, Bottom J, Mahut M, Paz A, Balakrishna R, Waller T and Hicks J (2011) Dynamic traffic assignment: A primer. *Transportation Research E-Circular* (E-C153).
- Evarts E (2013) Many americans are just a plug away from owning an electric car. URL <https://www.yahoo.com/news/many-americans-just-plug-away-owning-electric-car-160000286.html>.
- Fagnant DJ and Kockelman KM (2014) The travel and environmental implications of shared autonomous vehicles, using agent-based model scenarios. *Transportation Research Part C: Emerging Technologies* 40: 1–13.
- Fagnant DJ, Kockelman KM and Bansal P (2015) Operations of shared autonomous vehicle fleet for austin, texas, market. *Transportation Research Record: Journal of the Transportation Research Board* (2536): 98–106.
- Ford LR and Fulkerson DR (1962) *Flows in Networks*. Princeton University Press.
- Gelenbe E, Pujolle G and Nelson J (1998) *Introduction to queueing networks*, volume 2. Wiley Chichester.
- George DK (2012) *Stochastic Modeling and Decentralized Control Policies for Large-Scale Vehicle Sharing Systems via Closed Queueing Networks*. PhD Thesis, The Ohio State University.
- Goldberg AV, Oldham JD, Plotkin S and Stein C (1998) An implementation of a combinatorial approximation algorithm for minimum-cost multicommodity flow. In: Bixby R, Boyd E and Ros-Mercado R (eds.) *Integer Programming and Combinatorial Optimization, Lecture Notes in Computer Science*, volume 1412. Springer Berlin Heidelberg, pp. 338–352. DOI:10.1007/3-540-69346-7_26.
- Iglesias R, Rossi F, Zhang R and Pavone M (2016) A BCMP network approach to modeling and controlling Autonomous Mobility-on-Demand systems. In: *Workshop on Algorithmic Foundations of Robotics*.
- Kant K and Srinivasan M (1992) *Introduction to computer system performance evaluation*. McGraw-Hill College.
- Kobayashi H and Gerla M (1983) Optimal routing in closed queueing networks. In: *ACM SIGCOMM Computer Communication Review*, volume 13. ACM, pp. 26–26.

- Levin MW, Li T, Boyles SD and Kockelman KM (2016) A general framework for modeling shared autonomous vehicles. In: *95th Annual Meeting of the Transportation Research Board*.
- Mittelmann HD (2016) Decision tree for optimization software. URL <http://plato.asu.edu/guide.html>.
- Neil D (2015) Could self-driving cars spell the end of ownership? [wsj.com](http://www.wsj.com).
- Nourinejad M, Zhu S, Bahrami S and Roorda MJ (2015) Vehicle relocation and staff rebalancing in one-way carsharing systems. *Transportation Research Part E: Logistics and Transportation Review* 81: 98–113.
- Patriksson M (2015) *The traffic assignment problem: models and methods*. Courier Dover Publications.
- Pavone M, Smith SL, Frazzoli E and Rus D (2012) Robotic load balancing for mobility-on-demand systems. *International Journal of Robotics Research* 31(7): 839–854. DOI:10.1177/0278364912444766.
- Rossi F, Zhang R, Hindy Y and Pavone M (2018) Routing autonomous vehicles in congested transportation networks: Structural properties and coordination algorithms. *Autonomous Robots* In Press.
- Spieser K, Treleaven K, Zhang R, Frazzoli E, Morton D and Pavone M (2014) Toward a systematic approach to the design and evaluation of automated mobility-on-demand systems: A case study in Singapore. In: *Lecture Notes in Mobility*. Springer, pp. 229–245. DOI:10.1007/978-3-319-05990-7_20.
- Tesla Motors (2017) Supercharger. URL <https://www.tesla.com/supercharger>.
- Zhang R and Pavone M (2016) Control of robotic mobility-on-demand systems: A queueing-theoretical perspective. *International Journal of Robotics Research* 35(1-3): 186–203. DOI:10.1177/0278364915581863.
- Zhang R, Rossi F and Pavone M (2016) Model predictive control of autonomous mobility-on-demand systems. In: *Proc. IEEE Conf. on Robotics and Automation*. Stockholm, Sweden, pp. 1382 – 1389. DOI:10.1109/ICRA.2016.7487272.