

1 **CARA: A CONGESTION-AWARE ROUTING ALGORITHM FOR AUTONOMOUS**  
2 **MOBILITY-ON-DEMAND SYSTEMS**

3

4

5

6 **Mauro Salazar**

7

8 **Matthew Tsao**

9

10 **Izabel Aguiar**

11

12 **Maximilian Schiffer, Ph.D.**

13

14 **Marco Pavone, Ph.D.**

15

16

17 **Current Word Count:** 4190; Intro: 1041; Model: 1663; Results: 1118; Conclusion: 305

18

19

20

21

22

23 Submission Date: August 15, 2018

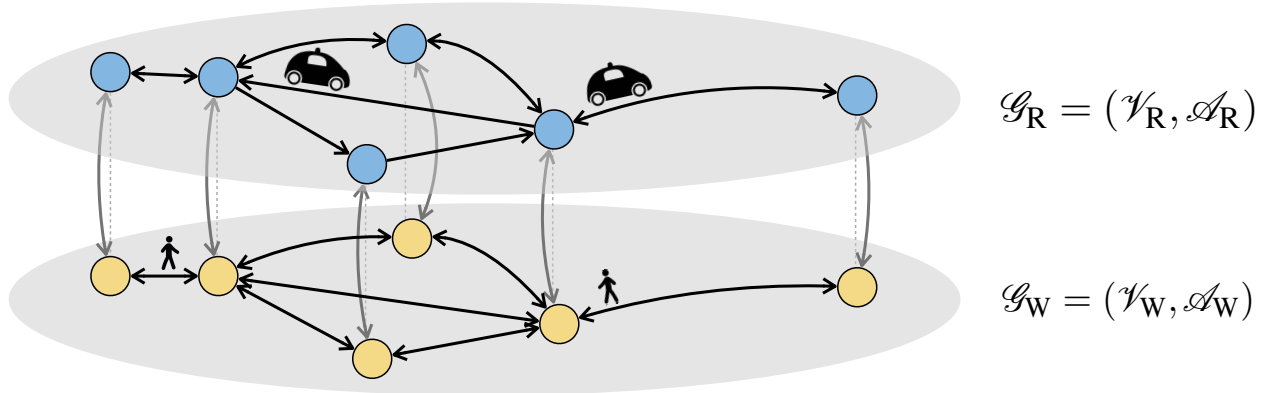
**1 Abstract**

2 We present a congestion-aware routing policy for Autonomous Mobility-on-Demand (AMoD)  
3 systems that accounts for the impact of road traffic on travel time. Specifically, we develop a  
4 congestion-aware routing algorithm (CARA) that captures road utilization dependent travel times  
5 via a piecewise affine approximation of the Bureau of Public Roads (BPR) model. Such an approx-  
6 imation largely retains the key features of the BPR model, while enabling the design of a real-time,  
7 convex quadratic optimization algorithm to compute congestion-aware routes for an AMoD fleet.  
8 Through a real-world case study of Manhattan, we compare CARA to existing routing approaches,  
9 the first one congestion-unaware and the second one based on a threshold congestion model. Nu-  
10 merical results show that CARA significantly outperforms the state of the art, with improvements  
11 in terms of travel time and global cost in the order of 20%, and computation times compatible with  
12 a real-time implementation.

13

14 *Keywords:* Network Flow Models, Autonomous Mobility, Congestion

15 **Current Word Count:** 4190; Intro: 1041; Model: 1663; Results: 1118; Conclusion: 305



**FIGURE 1** The AMoD network. The colored dots represent intersections and the black arrows denote road links or pedestrian pathways. The gray dotted lines highlight geographically equivalent nodes connected by gray mode-switching arcs.

## 1. INTRODUCTION

Congestion remains a central problem in today’s transportation systems, especially in densely populated urban areas. While congestion phenomena have been attenuated by subsidiary modes of transportation (e.g., public transport) in the past, in recent years congestion-related problems have rapidly increased. Even mature cities are struggling with these problems, with the sustainability of mobility systems restricted by current infrastructure and space limitations. Additionally, current societal trends increase congestion and emphasize its negative impacts. First, steadily increasing urbanization leads to higher population densities and thus higher mobility demand in cities (1). Second, commuters’ individual mobility needs and comfort preferences have led to the decrease in the utilization of public transportation (2). Third, mobility-on-demand services such as Uber and Lyft are rapidly growing as an alternative to public transportation and individual car ownership (3). Consequently, (selfish) traffic on roads is steadily growing and increases congestion even further. Current transport Key Performance Indicators (KPIs) reflect this situation. According to (4) the annual delay per commuter exceeds 42 hours and one third of commuter trips are reported as having “extreme, severe, or heavy” congestion. In New York City, the average speed in Manhattan decreased from 6.5 mph to 4.7 mph between 2012 and 2017 (5).

Cities face spatial limitations in addressing the issue of congestion as the available infrastructure (e.g., roads, traffic signals, highways, train lines) and its capacities are largely fixed. It is thus necessary to develop more efficient systems in road transport to supplement existing public transportation. Herein, Autonomous Mobility-on-Demand (AMoD) systems represent a promising solution. An AMoD system consists of a fleet of self-driving vehicles designed to carry passengers from their origins to their destinations. As passengers request trips, the central operator of the AMoD system assigns each passenger to an empty vehicle. Once the passengers are dropped off, the central operator computes an optimal rebalancing route for assigning the vehicle to the next request. The central operator is thus in control of simultaneously optimizing the routes for all vehicles in the system. As such, AMoD can replace current forms of mobility-on-demand (e.g., taxis and ride-hailing) while reducing the cost of travel (6). Furthermore, because the entire fleet of vehicles in an AMoD system is centrally controlled, this form of autonomous mobility holds promise to reduce congestion.

1           The potential of AMoD systems to alleviate congestion in cities depends on congestion-  
2 aware (CA) routing. CA routing considers the natural capacities of roads and the effects of vehicle  
3 flows on travel times through *volume-delay functions*. CA routing is not a new idea: navigation  
4 providers such as Google Maps and Waze have incorporated features that allow users to view how  
5 real-time and estimated traffic congestion will impact the routes of their commutes (7, 8). CA  
6 routing with conventional algorithms, however, can only *passively suggest* the routing of a single  
7 vehicle to avoid traffic, and furthermore only *anticipate* the traffic from other selfish vehicles.  
8 Conversely, using CA routing in an AMoD setting allows one to *actively control* the routes of  
9 all vehicles in the fleet under complete system information. With the aim to enable realistic CA  
10 routing for AMoD systems, we develop and study the use of a volume-delay function designed for  
11 convex optimization purposes. To set this work apart from the status quo, we briefly review related  
12 literature and state our aims and scope of the paper.

### 13 **1.1 Related Literature**

14           The most widely-used volume-delay function is the one developed by the Bureau of Public  
15 Roads (BPR) (9), although many other "BPR-type" functions exist (10). In related work these  
16 functions have been used for algorithmic approaches to dynamic estimation of congestion on a  
17 network (11) and CA route planning in agent-based models (12, 13). Furthermore, CA route  
18 planning of AMoD systems has been simulated in (14–18), and optimized for dynamic traffic  
19 assignment in (19, 20). These approaches have provided CA routing *analysis*, but are limited to  
20 simulations and lack control algorithms for both passenger requests and vehicle rebalancing. Thus  
21 far, the CA *control* of AMoD systems has been limited to thresholded approximations of the BPR  
22 function (21–23). In this approximation the time required to travel on a road is defined through a  
23 thresholded approach: the cars on a road are permitted to travel at a free-flow speed if the *capacity*  
24 of the road has not yet been reached. Additional cars beyond this capacity, however, make this  
25 road impossible to traverse. Although this model provides a conservative approach to capturing  
26 the effect of congestion on travel time, it oversimplifies congestion phenomena and may lead to  
27 suboptimal route patterns. To the best of our knowledge, no algorithmic framework for CA routing  
28 of AMoD systems currently exists that allows one to simultaneously address a) the more precise  
29 BPR-defined effects of congestion and b) real-time optimization of routing decisions for customer  
30 requests and rebalancing vehicles.

### 31 **1.2 Aims and Scope**

32           To resolve the drawbacks outlined above, we propose a congestion-aware routing algo-  
33 rithm (CARA) that leverages a piecewise-affine approximation of the BPR congestion model (9).  
34 We study the impact of such an algorithm in a real-world case study of Manhattan. We compare  
35 CARA to two different baselines, the first one congestion unaware and the second one capturing  
36 congestion via a threshold model. We study CA traffic routing in Manhattan and show that CARA  
37 significantly improves the state of the art with respect to travel times and global cost, while fea-  
38 turing computation times compatible with a real-time implementation. It should be noted that,  
39 while the BPR function does not capture every effect of congestion (e.g., spillback, heterogeneous  
40 vehicles, or intersection delays), it is a well-accepted *model* which suits the aim of this paper,  
41 namely, not to perfectly capture traffic dynamics, but to approximate them precisely enough for  
42 optimization and control purposes.

### 1 1.3 Organization

2 The remainder of this paper is structured as follows: in Section 2 we detail the methodology  
 3 and develop CARA: a flow-based optimization framework for CA routing of AMoD systems. In  
 4 Section 3 we numerically evaluate the performance of CARA through a case study of AMoD traffic  
 5 routing in Manhattan. We conclude the paper in Section 4 with a discussion of our results and an  
 6 outlook on future research directions.

## 7 2. METHODOLOGY

8 This section provides the methodological foundation for CARA. We introduce a multi-  
 9 commodity flow model to represent the physical constraints of the transportation system and  
 10 AMoD fleet in Section 2.1. In Section 2.2 we specify an objective function for the model with the  
 11 goal to optimize social welfare. We present an approximation of the BPR volume-delay function to  
 12 consider congestion while preserving model convexity in Section 2.3, and conclude in Section 2.4  
 13 with a brief discussion.

### 14 2.1 Multi-commodity Flow Based Optimization Approach

15 Recall from Section 1.2, that we aim to model a network with two modes of transportation:  
 16 walking and riding AMoD. This transportation network can be modeled on a digraph  $\mathcal{G} = (\mathcal{V}, \mathcal{A})$   
 17 as shown in Fig. 1. The graph consists of a set of vertices  $\mathcal{V}$  and a set of arcs  $\mathcal{A} \subseteq \mathcal{V} \times \mathcal{V}$ . To  
 18 capture both modes of transportation,  $\mathcal{G}$  comprises a road network layer  $\mathcal{G}_R = (\mathcal{V}_R, \mathcal{A}_R)$  and a walk-  
 19 ing layer  $\mathcal{G}_W = (\mathcal{V}_W, \mathcal{A}_W)$ . The road network layer represents intersections  $i \in \mathcal{V}_R$  and road links  
 20  $(i, j) \in \mathcal{A}_R$ , while the pedestrian layer models walkable streets  $(i, j) \in \mathcal{A}_W$  in between intersections  
 21  $i \in \mathcal{V}_W$ . Whereas the locations of nodes and arcs in  $\mathcal{G}_R$  and  $\mathcal{G}_W$  may coincide geographically, we  
 22 maintain a distinction between walking and riding in an AMoD vehicle. Additionally, switching  
 23 arcs out of set  $\mathcal{A}_C \subseteq \mathcal{V}_R \times \mathcal{V}_W$  connect the pedestrian layer to the road network layer, and model  
 24 the customer's ability to switch transportation modes by hailing an AMoD ride or exiting a car.  
 25 Collectively, it holds  $\mathcal{V} = \mathcal{V}_W \cup \mathcal{V}_R$  and  $\mathcal{A} = \mathcal{A}_W \cup \mathcal{A}_R \cup \mathcal{A}_C$ .

26 Each arc has a specific length  $d_{ij}$  and a constant nominal travel time  $t_{ij}^N$  denoting the walking  
 27 time for arcs  $(i, j) \in \mathcal{A}_W$ , the time to hail or exit an AMoD vehicle for arcs  $(i, j) \in \mathcal{A}_C$ , and the  
 28 travel time under free-flow conditions (without traffic) for arcs  $(i, j) \in \mathcal{A}_R$ . As in (22), we model  
 29 the energy consumption of AMoD vehicles, assuming a constant nominal speed  $v_{ij} = \frac{d_{ij}}{t_{ij}^N}$  for each  
 30 arc. Furthermore, we assume the AMoD fleet to be composed of lightweight electric vehicles with  
 31 an overall efficiency  $\eta_{EV}$  and full recuperation capabilities. Thus the energy consumption per road  
 32 arc is

$$33 \quad e_{ij} = \left( \frac{\rho_a}{2} \cdot A_f \cdot c_d \cdot v_{ij}^2 + c_r \cdot m_v \cdot g \right) \cdot \frac{d_{ij}}{\eta_{EV}} \quad (i, j) \in \mathcal{A}_R, \quad (1)$$

34 where the aerodynamic drag is determined by the air density  $\rho_a$ , the frontal area  $A_f$ , and the drag  
 35 coefficient  $c_d$ , and the friction of the wheels on the road is determined by the rolling friction  
 36 coefficient  $c_r$ , the mass of the vehicle  $m_v$ , and the gravitational acceleration  $g$  (24).

37 Each travel demand  $m \in \mathcal{M} = \{1, \dots, M\}$  consists of an origin destination pair  $(o_m, d_m)$   
 38 on the walking digraph  $\mathcal{G}_W$  and a demand rate  $\alpha_m$  that denotes the number of customers that  
 39 are requesting the same trip per unit time. To trace customer flows,  $f_m(i, j)$  denotes the flow of  
 40 customers on arc  $(i, j) \in \mathcal{A}$  for demand  $m \in \mathcal{M}$ . As AMoD vehicles may need to relocate between  
 41 two customer requests,  $f_0(i, j)$  denotes the rebalancing flow of empty vehicles on  $(i, j) \in \mathcal{A}_R$ .

42 This notation is sufficient to derive a basic multi-commodity flow model for our planning

1 problem at hand. Consider the cost function  $J$  mapping the set of flows  $\{f_m(\cdot, \cdot)\}_m, f_0(\cdot, \cdot)$  into the  
 2 set of non-negative real numbers  $\mathbb{R}_{\geq 0}$ . We state the AMoD optimal routing problem as

$$3 \quad \min_{\{f_m(\cdot, \cdot)\}_m, f_0(\cdot, \cdot)} J(\{f_m(\cdot, \cdot)\}_m, f_0(\cdot, \cdot)) \quad (2a)$$

4 s.t.

$$5 \quad \sum_{i:(i,j) \in \mathcal{A}} f_m(i, j) + \mathbb{1}_{j=o_m} \cdot \alpha_m = \sum_{k:(j,k) \in \mathcal{A}} f_m(j, k) + \mathbb{1}_{j=d_m} \cdot \alpha_m \quad \forall m \in \mathcal{M}, j \in \mathcal{V} \quad (2b)$$

$$6 \quad \sum_{i:(i,j) \in \mathcal{A}_R} \left( f_0(i, j) + \sum_{m \in \mathcal{M}} f_m(i, j) \right) = \sum_{k:(j,k) \in \mathcal{A}_R} \left( f_0(j, k) + \sum_{m \in \mathcal{M}} f_m(j, k) \right) \quad \forall j \in \mathcal{V}_R \quad (2c)$$

$$7 \quad f_m(i, j) \geq 0 \quad \forall (i, j) \in \mathcal{A} \quad (2d)$$

$$8 \quad f_0(i, j) \geq 0 \quad \forall (i, j) \in \mathcal{A}_R, \quad (2e)$$

9

10 where  $\mathbb{1}_x$  is a boolean indicator function. Linear constraints (2b) and (2c) ensure that the mass of  
 11 customers and vehicles, respectively, are conserved on every road node. Inequality constraints (2d)  
 12 and (2e) ensure that the customer and rebalancing flows are non-negative.

## 13 2.2 AMoD Objective

14 We use the cost function (2a) to model the social cost of serving the transportation requests  
 15 in a similar fashion as (22). Specifically, we aim to minimize the total travel time and the op-  
 16 erational costs of the AMoD system. We assume customers to have the same value of time  $V_T$ .  
 17 We separate the cost of operating the AMoD fleet into a distance-dependent cost  $V_D$  due to de-  
 18 preciation and maintenance, and an energy-consumption-dependent cost  $V_E$ . Collectively, the cost  
 19 function (2a) is

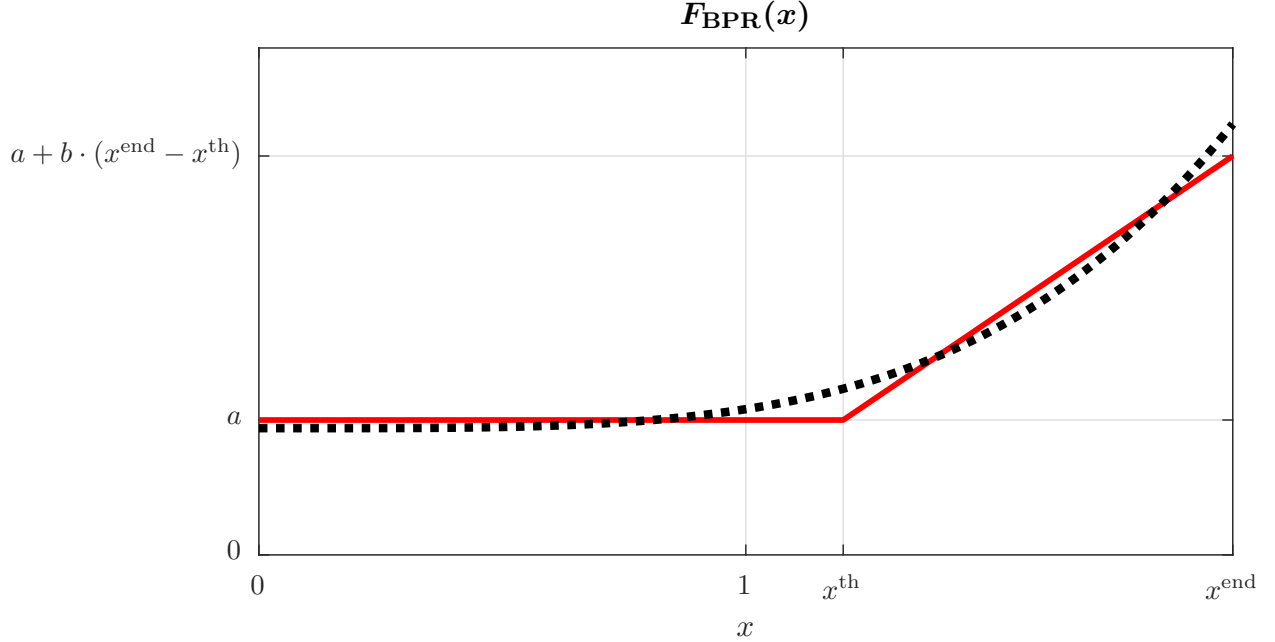
$$J_M(\{f_m(\cdot, \cdot)\}_m, f_0(\cdot, \cdot)) = V_T \cdot \sum_{m \in \mathcal{M}, (i,j) \in \mathcal{A}} t_{ij}(f(i, j)) \cdot f_m(i, j) \\ + \sum_{(i,j) \in \mathcal{A}_R} (V_D \cdot d_{ij} + V_E \cdot e_{ij}) \cdot f(i, j), \quad (3)$$

20 where  $f(i, j) = f_0(i, j) + \sum_{m \in \mathcal{M}} f_m(i, j)$  is the total flow on arc  $(i, j)$ . The travel time  $t_{ij}$  on road  
 21 arcs is modeled as a function of road usage, while  $t_{ij} = t_{ij}^N$  gives the time to walk or switch trans-  
 22 portation layer. For road arcs we use the BPR volume-delay function (9)

$$t_{ij}(f(i, j)) = t_{ij}^N \cdot F_{\text{BPR}} \left( \frac{f(i, j) + u_{ij}^R}{c_{ij}^R} \right) \quad \forall (i, j) \in \mathcal{A}_R, \quad (4)$$

23 where  $u_{ij}^R$  is the exogenous road usage caused by, e.g., the presence of private cars,  $c_{ij}^R$  being the  
 24 nominal road capacity, and

$$25 \quad F_{\text{BPR}}(x) = 1 + 0.15 \cdot x^4. \quad (5)$$



**FIGURE 2** The BPR function  $F_{\text{BPR}}(x)$  (black dotted) and its piecewise affine fit (red solid).

### 1 2.3 BPR Model Approximation

2 We aim to minimize (2a) subject to the constraints specified by (2b)-(2e). However,  $J_M$  is  
 3 a non-convex polynomial of the decision variables, namely the vehicle flows  $\{f_m(\cdot, \cdot)\}_m, f_0(\cdot, \cdot)$ .  
 4 Hence, this model is, in general, computationally intractable, as there does not exist a known algo-  
 5 rithm to reliably and efficiently solve large-scale non-convex polynomial optimization problems.  
 6 To resolve this non-convexity we discuss a piecewise affine approximation to the BPR function.

7 The non-convexity arises from the product  $t_{ij}(f(i, j)) \cdot f_m(i, j)$ . Hence, we approximate  
 8 this term with a convex function to ensure scalability to large-size problem instances. Specifically,  
 9 we fit Eq. (5) using a piecewise affine approximation as shown in Fig. 2:

$$10 \quad y = \begin{cases} a & \text{if } x \in [0, x^{\text{th}}] \\ a + b \cdot (x - x^{\text{th}}) & \text{if } x \in (x^{\text{th}}, x^{\text{max}}] \end{cases}, \quad (6)$$

11 where  $a$  is the height of the horizontal line,  $b$  the slope of the second line,  $x^{\text{th}}$  is the non-smooth  
 12 threshold in the piecewise affine approximation, and  $x^{\text{max}}$  defines the approximation window. This  
 13 way we approximate the BPR function (4) as

$$14 \quad t_{ij} = \begin{cases} t_{ij}^N \cdot a & \text{if } f(i, j) \leq c_{ij}^{\text{R,th}} \\ t_{ij}^N \cdot \left( a + b \cdot \frac{f(i, j) - c_{ij}^{\text{R,th}}}{c_{ij}^{\text{R}}} \right) & \text{if } f(i, j) > c_{ij}^{\text{R,th}} \end{cases} =: t_{ij}^0 + \tau_{ij} \cdot \varepsilon(i, j), \quad (7)$$

15 where  $t_{ij}^0 = a \cdot t_{ij}^N$ ,  $\tau_{ij} = b \cdot t_{ij}^N / c_{ij}^{\text{R}}$ , and the slack variable  $\varepsilon(i, j)$  denotes how much the total flow on  
 16 arc  $(i, j) \in \mathcal{A}_R$  exceeds its threshold capacity  $c_{ij}^{\text{R,th}} = x^{\text{th}} \cdot c_{ij}^{\text{R}}$ , that is,

$$17 \quad \varepsilon(i, j) = \max \{0, f(i, j) + u_{ij}^{\text{R}} - c_{ij}^{\text{R,th}}\}. \quad (8)$$

18 Provided that  $J$  is a decreasing function of  $\varepsilon(i, j)$ , this piecewise affine cost can be represented by

1 the following linear inequality constraints

$$2 \quad \varepsilon(i, j) \geq f(i, j) + u_{ij}^R - c_{ij}^{\text{R,th}} \quad (9)$$

$$3 \quad \varepsilon(i, j) \geq 0. \quad (10)$$

4 Hence, the total travel time  $\zeta$  on road links can be expressed as

$$\begin{aligned} \zeta &:= T^R \cdot \sum_{m \in \mathcal{M}} \alpha_m = \sum_{(i,j) \in \mathcal{A}_R} t_{ij}(f(i, j)) \cdot \sum_{m \in \mathcal{M}} f_m(i, j) \\ &= \sum_{(i,j) \in \mathcal{A}_R} t_{ij}^0 \cdot \sum_{m \in \mathcal{M}} f_m(i, j) + \sum_{(i,j) \in \mathcal{A}_R} \tau_{ij} \cdot \varepsilon(i, j) \cdot \sum_{m \in \mathcal{M}} f_m(i, j) \\ &= \sum_{(i,j) \in \mathcal{A}_R} t_{ij}^0 \cdot \sum_{m \in \mathcal{M}} f_m(i, j) + \sum_{(i,j) \in \mathcal{A}_R} \tau_{ij} \cdot \varepsilon(i, j) \cdot \left( \varepsilon(i, j) + c_{ij}^{\text{R,th}} - u_{ij}^R - f_0(i, j) \right) \\ &\leq \sum_{(i,j) \in \mathcal{A}_R} t_{ij}^0 \cdot \sum_{m \in \mathcal{M}} f_m(i, j) + \sum_{(i,j) \in \mathcal{A}_R} \tau_{ij} \cdot \left( \varepsilon(i, j)^2 + \varepsilon(i, j) \cdot (c_{ij}^{\text{R,th}} - u_{ij}^R) \right) \\ &=: \widehat{\zeta} \end{aligned} \quad (11)$$

6

7 We would like to include the total trip time  $\zeta$  in the objective for the optimization, but due to the bi-

8 linear terms  $-\varepsilon(i, j) f_0(i, j)$ ,  $\zeta$  is not a convex function of the flow variables. As an alternative, we

9 find a convex approximation of  $\zeta$  to include in the objective function. Since  $\varepsilon(i, j) \cdot f_0(i, j) \geq 0$ ,

10 removing these bilinear terms makes the expression larger, hence we have  $\zeta \leq \widehat{\zeta}$ . Without the bi-

11 linear terms,  $\widehat{\zeta}$  is a convex function of the flow variables, so we include it in the objective function

12 to penalize strategies with long trip times. Considering that the number of rebalancing vehicles

13 has a minor impact with respect to road congestion and converges to zero for perfectly symmetric

14 demand distributions (21),  $\widehat{\zeta}$  can be used as a metric for the total travel time on road arcs. Specif-

15 ically, our empirical studies observed that  $\frac{|\zeta - \widehat{\zeta}|}{\zeta} \leq 0.1$ . In doing so, we approximate the total cost

16 function (2a) with the quadratic bound (11) as

$$\begin{aligned} J_Q(\{f_m(\cdot, \cdot)\}_m, f_0(\cdot, \cdot)) &= V_T \cdot \underbrace{\left( \sum_{\substack{m \in \mathcal{M}, \\ (i,j) \in \mathcal{A}}} t_{ij}^N \cdot f_m(i, j) + \sum_{(i,j) \in \mathcal{A}_R} \tau_{ij} \cdot \left( \varepsilon(i, j)^2 + \varepsilon(i, j) \cdot (c_{ij}^{\text{R,th}} - u_{ij}^R) \right) \right)}_{\widehat{\zeta}} \\ &+ \sum_{(i,j) \in \mathcal{A}_R} (V_D \cdot d_{ij} + V_E \cdot e_{ij}) \cdot f(i, j). \end{aligned} \quad (12)$$

18

19 The AMoD optimal routing problem given by (2) with  $J(\cdot, \cdot) = J_Q(\cdot, \cdot)$  subject to (9) then

20 remains a convex quadratic program.

## 21 2.4 Discussion

22 A few comments are in order. First, we assume travel requests to be time-invariant. This as-

23 sumption is reflected in densely populated urban environments where requests often change slowly

24 compared to the average time needed to complete an individual trip (25). Second, we use the BPR

25 function (9) to describe the impact of road usage on travel time. While such a function does

26 not perfectly capture traffic phenomena, it is a well-established *model* serving the purpose of de-

27 signing CA routing algorithms. In order to embed the BPR function in a convex optimization



1 framework, we approximate it in a piecewise affine fashion. This approximation gives a general-  
 2 ization of simple threshold models that are used in classical traffic flow theory (26) and can allow  
 3 for better mobility service in congested situations, as discussed in the remainder of the paper. Al-  
 4 though a piecewise affine function cannot closely approximate a quartic polynomial such as the  
 5 BPR function on its entire domain, this approximation only needs to be accurate for realistic val-  
 6 ues of vehicle flows. Furthermore, we relax the piecewise approximation and implement it in the  
 7 convex optimization framework through a quadratic upper bound. Nevertheless, since the ratio of  
 8 empty vehicles to passenger carrying vehicles is usually low, this bound is tight enough, as shown  
 9 in Section 3. Third, CARA captures customer and vehicle routes as fractional flows and does not  
 10 address the stochasticity of the exogenous traffic and customer requests. Arguably, such approx-  
 11 imations are acceptable, given the mesoscopic perspective of our study. On the topic of operational  
 12 algorithms, CARA can be directly extended to operate in real-time. Due to the computational effi-  
 13 ciency of the algorithm, CARA can be run periodically with updated real-time information about  
 14 customer demand to operate in a time varying environment. Furthermore, randomized rounding  
 15 routing algorithms can compute near-optimal integer-valued flows for individual customers starting  
 16 from the fractional solution computed by CARA (27). The operational implementation of CARA  
 17 is further discussed in Section 4. Fourth, we assume vehicles to carry only a single customer at a  
 18 time. This mode of operation is in line with current trends in mobility-on-demand systems such as  
 19 taxis, Lyft, and Uber. The extension to ride-sharing is a challenging direction for future research.  
 20 Finally, for the sake of simplicity, we assume all customers to value time and travel comfort in  
 21 the same way. However, CARA can be extended to capture multiple classes of customers using  
 22 network flows which are distinct not only in the origin-destination transportation request, but also  
 23 in the customers' preference profile.

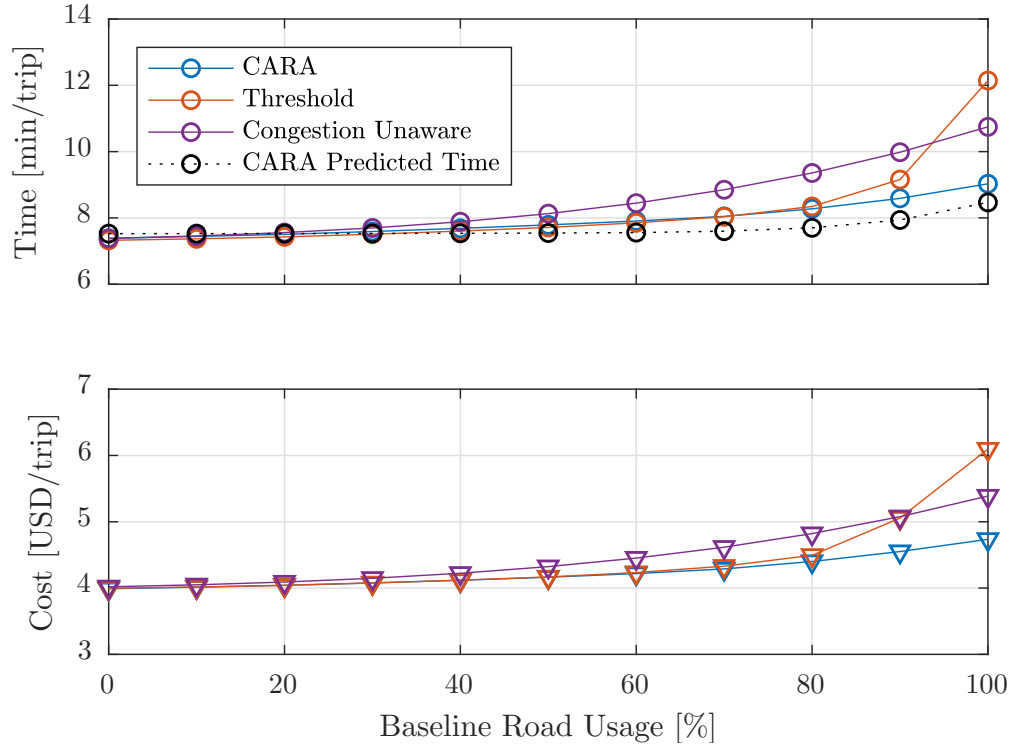
### 24 3. COMPUTATIONAL STUDIES

25 The goal of this section is to assess CARA by evaluating and comparing its performance  
 26 in a case study to state of the art congestion aware routing algorithms. We begin in Section 3.1  
 27 by describing the Manhattan-based case study where we conduct computational studies. In Sec-  
 28 tion 3.2, we compare the performance of CARA to that of a) a congestion-unaware approach and  
 29 b) an algorithm that captures congestion via a threshold model.

#### 30 3.1 Manhattan Case Study

31 Our case study is based on the area of Manhattan in New York City. For this area, we use  
 32 real data from taxi rides that occurred between 6:00PM and 8:00PM on March 1, 2012 (courtesy of  
 33 the New York Taxi and Limousine Commission). In total this data set comprises 53,932 taxi rides.  
 34 Although this number of trips is quite large, it reflects only a fraction of the travel demand between  
 35 6:00PM and 8:00PM. In 2017 ride-sharing vehicles used during this time period outnumbered taxis  
 36 by a ratio of 5:1 (28). To reflect this, we dilate the number of requests by a factor of six to emulate  
 37 the total demand for ride-hailing services in Manhattan during this time window.

38 We derive the road network from Open Street Map data (29) and set the nominal road  
 39 capacities  $c_{ij}^R$  proportional to the number of lanes on a road, times its speed limit (21). The walking  
 40 network shows similar spatial characteristics to the road network, but is complementary as all arcs  
 41 are bi-directional to allow for walking in both directions, even in one-way streets. To account for  
 42 exogenous privately owned vehicles on the road, we run our simulations for different values of  
 43 road usage  $u^R$ , denoting the fraction of a road's free-flow capacities used by private vehicles. For



**FIGURE 3 Average travel time and cost per passenger trip for all approaches.**

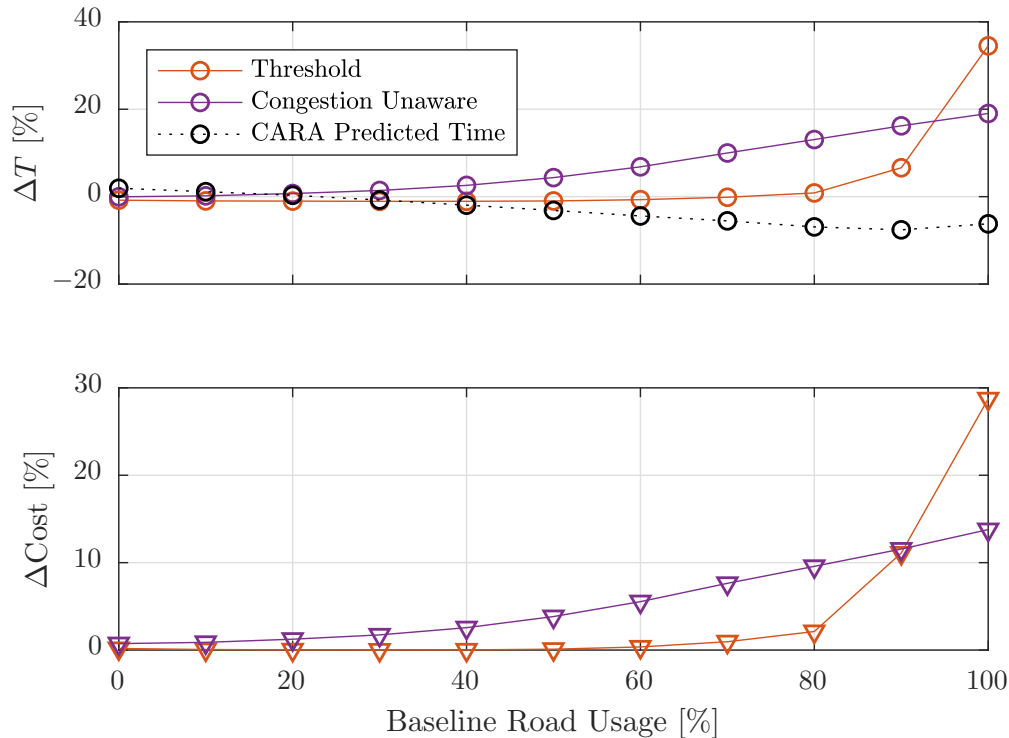
1 a more detailed description of this case study and its parameters we refer to (22).

## 2 3.2 Results

3 In the following section, we discuss the results of our case study. We compare CARA to  
 4 two baselines: a control algorithm using a simple threshold model in which  $f(i, j) + u_{ij}^R \leq c_{ij}^{R,th}$  is  
 5 enforced (cf. (21, 22)), and a naïve congestion-unaware approach in which road utilization does  
 6 not affect travel times. For the first two models we set  $c_{ij}^{R,th} = x^{th} \cdot c_{ij}^R$ , with  $x^{th} = 1.2$ , and we  
 7 approximate the BPR function with  $x^{end} = 2$  (cf. Fig. 2). We also compute the CARA Predicted  
 8 Time, which is the average trip time predicted using the piecewise affine function and approxima-  
 9 tions discussed in section 2.3. For each of the scenarios studied, the computation of the optimal  
 10 solution took less than four minutes on commodity hardware (Intel Core i7, 16 GB RAM) using  
 11 Gurobi 7.5.

12 We measure the quality of CARA with two performance indicators: the average resulting  
 13 trip time for each solution computed by the BPR function, and the global cost as defined in (3).  
 14 Fig. 3 and 4 summarize the results. Focusing on these results, we identify different relationships  
 15 between the performance of all algorithms depending on the exogenous road usage. Additionally,  
 16 the travel times predicted by CARA are close to those computed by the BPR model, as shown in  
 17 Fig. 3 and 4.

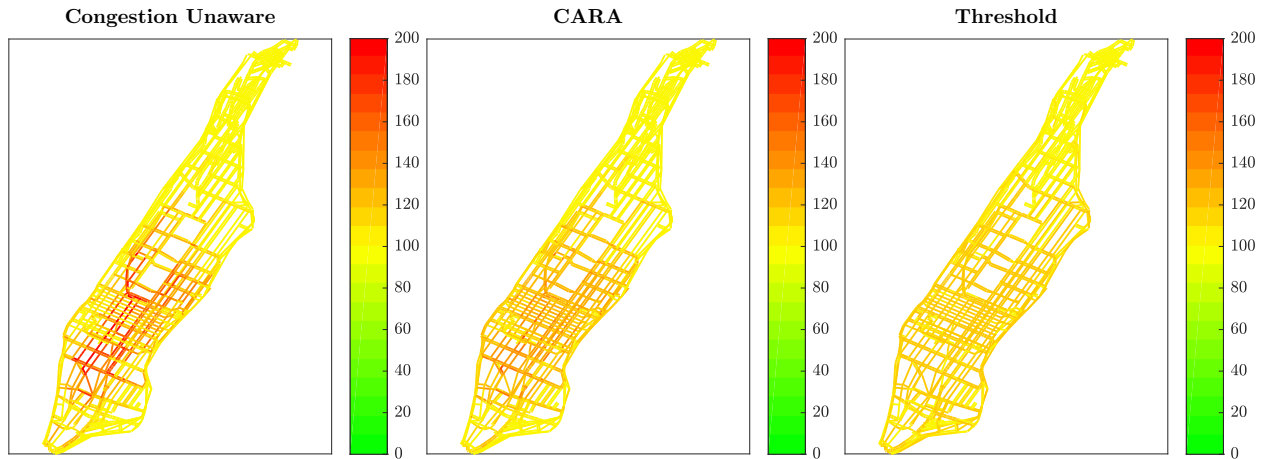
- 18 •  $u^R \in [0\%, 20\%] \cdot c^R$ : For low levels of road usage (i.e.,  $u^R$  close to zero), all control  
 19 algorithms show the same performance: in this scenario roads have enough capacity for  
 20 AMoD vehicles in the absence of private vehicles.
- 21 •  $u^R \in [40\%, 80\%] \cdot c^R$ : As the traffic level increases, the naïve control algorithm performs



**FIGURE 4** Relative deviation of the two baseline solutions compared to CARA. The upper plot also shows the relative error of the CARA predicted time.

- 1 worse than the thresholded approach and CARA in terms of both average travel time and  
 2 operational cost. This shows that the naïve control algorithm, lacking information about  
 3 a street’s capacity, sends the cars on the spatially shortest paths. As a consequence these  
 4 paths become congested and remain the spatially but not temporally shortest paths due  
 5 to increased travel times. Conversely, the threshold model and CARA perform similarly  
 6 better for this range.
- 7 •  $u^R \in [90\%, 100\%] \cdot c^R$ : For heavy levels of traffic CARA significantly outperforms the  
 8 other two approaches. Especially at 100% exogenous road usage, CARA offers a 20%  
 9 shorter average trip time than the naïve algorithm and a 35% shorter average trip time  
 10 than the threshold algorithm. Note that for heavy exogenous road usage, the threshold  
 11 model performs even worse than the naïve model. This is because the threshold algorithm  
 12 is not allowed to use roads once the flow on it reaches  $c^{R,th}$ . In this case, the only other  
 13 option is to send customer carrying vehicles on long detours. In reality,  $c^{R,th}$  is not a  
 14 hard constraint and sending the car along a congested road can still be faster than the  
 15 aforementioned long detours. Fig. 5 shows the congestion maps for this scenario and  
 16 the three routing algorithms, where the shortcomings of the congestion unaware and the  
 17 threshold algorithm are highlighted together with the balanced behaviour of CARA.

18 In summary, we find that the performance of the threshold congestion model depends heavily on  
 19 the level of exogenous traffic. In situations of low to medium exogenous traffic, an algorithm using  
 20 the threshold model can perform much better than a congestion-unaware model. However, under  
 21 heavy levels of traffic, the threshold model behaves too conservatively, i.e., it underutilizes routes



**FIGURE 5** Road congestion map for the naïve algorithm, CARA, and the threshold model with a 100% exogenous road usage level.

1 on main roads and favors longer detours.

2 On the other hand, CARA improves the operational cost and the travel time compared to  
 3 both baseline approaches for all levels of exogenous road usage, as the error introduced by the  
 4 convex relaxation (11) is well bounded below 10%. In dire situations with heavy traffic (when CA  
 5 algorithms are most needed) CARA achieves improvements of up to 20% in terms of travel time,  
 6 striking an effective compromise between congesting roads and detouring vehicles, as shown in  
 7 Fig. 5.

8 To conclude, our experiments show that CARA can significantly reduce both operational  
 9 cost and the duration of customer trips. On the other hand, although conservative (thresholded)  
 10 models of congestion offer a realistic evaluation of travel constraints at low- to mid-levels of ex-  
 11 ogenous traffic, at high-levels they perform worse than a naïve approach that ignores congestion  
 12 effects altogether. Most significantly, the overconservative nature of the threshold model leads to  
 13 an increase of more than \$0.50 USD per trip. Conversely, the balanced nature of CARA leads to  
 14 a decrease of the same amount per trip. With more than half a million trips happening on average  
 15 each day in NYC (30), this efficiency could save over one hundred million dollars over the course  
 16 of one year.

#### 17 4. CONCLUSION

18 In this paper we studied the impact of congestion-aware routing in Autonomous Mobility-  
 19 on-Demand systems. By leveraging an affine approximation to the Bureau of Public Roads model,  
 20 we derived a congestion-aware routing algorithm (CARA) that entails solving a computationally-  
 21 efficient, convex quadratic program. We applied CARA to a case study of AMoD traffic routing  
 22 in Manhattan and compared its performance to two baselines: a congestion-unaware routing al-  
 23 gorithm and an algorithm that considers congestion via a threshold model. Our results show that  
 24 CARA always outperforms the status quo baselines and achieves an effective balance between  
 25 congesting roads and detouring vehicles, thus achieving significant improvements, especially un-  
 26 der high-levels of traffic. Additionally, we showed that in certain cases threshold models may even  
 27 worsen the system performance compared to a congestion-unaware approach. Whereas consid-  
 28 ering congestion as a hard constraint can improve performance under low- to medium-levels of

1 traffic, it can be overconservative and therefore detrimental in terms of travel time for high-levels  
2 of congestion.

3         This work opens several avenues for future research. Foremost, we would like to design an  
4 operational algorithm to apply CARA at a microscopic level. We note that, because the mesoscopic  
5 solutions presented here were computed in less than four minutes, the extension to an operational  
6 algorithm is completely attainable. Namely, CARA can feasibly become operational by continu-  
7 ally solving for optimal flows given updated, real-time customer demands in a receding-horizon  
8 model predictive fashion (cf. (31–33)). Furthermore, as discussed in Section 2.4, near-optimal  
9 integer-valued flows for operational implementation can be achieved from the fractional solutions  
10 computed with CARA via random sampling algorithms (27). Additional future research includes  
11 coupling CARA with urban infrastructure, e.g., the power grid (23) and public transit (22); includ-  
12 ing stochastic effects, e.g., of customer demand and road congestion (34); and adapting the cost  
13 function to include different metrics based on customer preferences.

## 14 **5. ACKNOWLEDGEMENTS**

15         We would like to thank Dr. Guido Gentile and Dr. Daniele Vigo for providing useful in-  
16 sights on congestion models, and Dr. Ilse New for proofreading this paper and providing us with  
17 her expert comments and advice. The first author would like to express his gratitude to Dr. Chris  
18 Onder for his support. This research was supported by the National Science Foundation under CA-  
19 REER Award CMMI-1454737 and the Toyota Research Institute (TRI). This article solely reflects  
20 the opinions and conclusions of its authors and not NSF, TRI, or any other entity.

## 21 **6. STATEMENT OF CONTRIBUTIONS**

22         The authors confirm contribution to the paper as follows: study conception and design:  
23 Mauro Salazar, Marco Pavone; data collection: Mauro Salazar, Matthew Tsao; analysis and in-  
24 terpretation of results: Mauro Salazar, Matthew Tsao, and Maximilian Schiffer; draft manuscript  
25 preparation: Mauro Salazar, Matthew Tsao, Izabel Aguiar, Maximilian Schiffer, and Marco Pavone.  
26 All authors reviewed the results and approved the final version of the manuscript.

## 1 REFERENCES

- 2 [1] *The World Factbook*. Central Intelligence Agency, 2018, Available at <https://www.cia.gov/library/publications/the-world-factbook/fields/2212.html>.
- 3
- 4 [2] Siddiqui, F., *Failing transit ridership poses an ‘emergency’ for cities, experts fear*. The Wash-  
5 ington Post, 2018, available online.
- 6 [3] Molla, R., *Americans seem to like ride-sharing services like Uber and Lyft. But it’s hard to*  
7 *say exactly how many use them*. Recode, 2018, Available at [https://www.recode.net/](https://www.recode.net/2018/6/24/17493338/ride-sharing-services-uber-lyft-how-many-people-use)  
8 [2018/6/24/17493338/ride-sharing-services-uber-lyft-how-many-people-use](https://www.recode.net/2018/6/24/17493338/ride-sharing-services-uber-lyft-how-many-people-use).
- 9 [4] Bureau of Transportation Statistics, *Transportation Statistics Annual Report 2017*. U.S. Dept.  
10 of Transportation, 2017.
- 11 [5] Hu, W., *Your Uber Car Creates Congestion. Should You Pay a Fee to Ride?* The New York  
12 Times, 2017, available online.
- 13 [6] Spieser, K., K. Treleven, R. Zhang, E. Frazzoli, D. Morton, and M. Pavone, Toward a Sys-  
14 tematic Approach to the Design and Evaluation of Autonomous Mobility-on-Demand Sys-  
15 tems: A Case Study in Singapore. In *Road Vehicle Automation*, Springer, 2014.
- 16 [7] *Live Map*. Waze Mobile, 2018, Available at <https://www.waze.com/livemap>.
- 17 [8] *Google Maps Platform: Routes*. Google, 2018, Available at [https://cloud.google.com/](https://cloud.google.com/maps-platform/routes/)  
18 [maps-platform/routes/](https://cloud.google.com/maps-platform/routes/).
- 19 [9] Bureau of Public Roads, *Traffic Assignment Manual*. U.S. Dept. of Commerce, Urban Plan-  
20 ning Division, 1964.
- 21 [10] Spiess, H., Technical Note—Conical Volume-Delay Functions. *Transportation Science*,  
22 Vol. 24, No. 2, 1990, pp. 153–158.
- 23 [11] Rivas, A., G. Inmaculada, S. Sánchez-Cambronero, R. M. Barba, and L. Ruiz-Ripoll, A  
24 Continuous Dynamic Traffic Assignment Model From Plate Scanning Technique. *Transport*  
25 *Research Procedia*, Vol. 18, 2016, pp. 332–340.
- 26 [12] Aslam, J., S. Lim, and D. Rus, Congestion-aware Traffic Routing System using sensor data.  
27 In *Proc. IEEE Int. Conf. on Intelligent Transportation Systems*, 2012.
- 28 [13] Manley, E., T. Cheng, A. Penn, and A. Emmonds, A framework for simulating large-scale  
29 complex urban traffic dynamics through hybrid agent-based modelling. *Computers, Environ-*  
30 *ment and Urban Systems*, Vol. 44, 2014, pp. 27–36.
- 31 [14] Treiber, M., A. Hennecke, and D. Helbing, Microscopic simulation of congested traffic. In  
32 *Traffic and Granular Flow ’99*, Springer Berlin Heidelberg, 2000.
- 33 [15] Yang, Q. and H. N. Koutsopoulos, A microscopic traffic simulator for evaluation of dynamic  
34 traffic management systems. *Transportation Research Part C: Emerging Technologies*, Vol. 4,  
35 No. 3, 1996, pp. 113–129.
- 36 [16] Balmer, M., M. Rieser, K. Meister, D. Charypar, N. Lefebvre, and K. Nagel, MATSim-T:  
37 Architecture and Simulation Times. In *Multi-Agent Systems for Traffic and Transportation*  
38 *Engineering*, 2009, chap. 3.
- 39 [17] Fagnant, D. J. and K. M. Kockelman, The travel and environmental implications of shared  
40 autonomous vehicles, using agent-based model scenarios. *Transportation Research Part C:*  
41 *Emerging Technologies*, Vol. 40, 2014, pp. 1–13.
- 42 [18] Levin, M. W., K. M. Kockelman, S. D. Boyles, and T. Li, A general framework for mod-  
43 eling shared autonomous vehicles with dynamic network-loading and dynamic ride-sharing  
44 application. *Computers, Environment and Urban Systems*, Vol. 64, 2017, pp. 373 – 383.

- 1 [19] Janson, B. N., Dynamic traffic assignment for urban road networks. *Transportation Research*  
2 *Part B: Methodological*, Vol. 25, No. 2–3, 1991, pp. 143–161.
- 3 [20] Peeta, S. and H. S. Mahmassani, System optimal and user equilibrium time-dependent traffic  
4 assignment in congested networks. *Annals of Operations Research*, Vol. 60, No. 1, 1995, pp.  
5 81–113.
- 6 [21] Rossi, F., R. Zhang, Y. Hindy, and M. Pavone, Routing Autonomous Vehicles in Congested  
7 Transportation Networks: Structural Properties and Coordination Algorithms. *Autonomous*  
8 *Robots*, 2018, in Press.
- 9 [22] Salazar, M., F. Rossi, M. Schiffer, C. H. Onder, and M. Pavone, On the Interaction between  
10 Autonomous Mobility-on-Demand and the Public Transportation Systems. In *Proc. IEEE Int.*  
11 *Conf. on Intelligent Transportation Systems*, 2018, in Press. Extended Version, Available at  
12 <https://arxiv.org/abs/1804.11278>.
- 13 [23] Rossi, F., R. Iglesias, M. Alizadeh, and M. Pavone, On the Interaction Between Au-  
14 tonomous Mobility-on-Demand Systems and the Power Network: Models and Coordina-  
15 tion Algorithms. In *Robotics: Science and Systems*, 2018, Extended version available at  
16 <https://arxiv.org/abs/1709.04906>.
- 17 [24] Guzzella, L. and A. Sciarretta, *Vehicle Propulsion Systems*. Springer Berlin Heidelberg, 2007.
- 18 [25] Neuburger, H., The economics of heavily congested roads. *Transportation Research*, Vol. 5,  
19 No. 4, 1971, pp. 283–293.
- 20 [26] Wardrop, J. G., Some Theoretical Aspects of Road Traffic Research. *Proc. of the Institution*  
21 *of Civil Engineers*, Vol. 1, No. 3, 1952, pp. 325–362.
- 22 [27] Rossi, F., *On the Interaction between Autonomous Mobility-on-Demand Systems and the Built*  
23 *Environment: Models and Large Scale Coordination Algorithms*. Ph.D. thesis, Stanford Uni-  
24 versity, Dept. of Aeronautics and Astronautics, 2018.
- 25 [28] Sugar, R., *Uber and Lyft cars now outnumber yellow cabs in NYC 4 to 1*.  
26 Curbed, 2017, Available at [https://ny.curbed.com/2017/10/13/16468716/](https://ny.curbed.com/2017/10/13/16468716/uber-yellow-cab-nyc-surpass-ridership)  
27 [uber-yellow-cab-nyc-surpass-ridership](https://ny.curbed.com/2017/10/13/16468716/uber-yellow-cab-nyc-surpass-ridership).
- 28 [29] Haklay, M. and P. Weber, OpenStreetMap: User-Generated Street Maps. *IEEE Pervasive*  
29 *Computing*, Vol. 7, No. 4, 2008, pp. 12–18.
- 30 [30] Warerkar, T., *Uber surpasses yellow cabs in average daily ridership in NYC*.  
31 Curbed, 2017, Available at [https://ny.curbed.com/2017/10/13/16468716/](https://ny.curbed.com/2017/10/13/16468716/uber-yellow-cab-nyc-surpass-ridership)  
32 [uber-yellow-cab-nyc-surpass-ridership](https://ny.curbed.com/2017/10/13/16468716/uber-yellow-cab-nyc-surpass-ridership).
- 33 [31] Seow, K. T., N. H. Dang, and D. H. Lee, A collaborative multiagent taxi-dispatch system.  
34 *IEEE Transactions on Automation Sciences and Engineering*, Vol. 7, No. 3, 2010, pp. 607–  
35 616.
- 36 [32] Pavone, M., S. L. Smith, E. Frazzoli, and D. Rus, Robotic Load Balancing for Mobility-on-  
37 Demand Systems. *Int. Journal of Robotics Research*, Vol. 31, No. 7, 2012, pp. 839–854.
- 38 [33] Zhang, R., F. Rossi, and M. Pavone, Model Predictive Control of Autonomous Mobility-on-  
39 Demand Systems. In *Proc. IEEE Conf. on Robotics and Automation*, 2016.
- 40 [34] Tsao, M., R. Iglesias, and M. Pavone, Stochastic Model Predictive Control for Autonomous  
41 Mobility on Demand. In *Proc. IEEE Int. Conf. on Intelligent Transportation Systems*, 2018,  
42 in Press. Extended Version, Available at <https://arxiv.org/pdf/1804.11074>.