

A Congestion-aware Routing Scheme for Autonomous Mobility-on-Demand Systems

Mauro Salazar^{1,2}, Matthew Tsao², Izabel Aguiar², Maximilian Schiffer³ and Marco Pavone²

Abstract—We study route-planning for Autonomous Mobility-on-Demand (AMoD) systems that accounts for the impact of road traffic on travel time. Specifically, we develop a congestion-aware routing scheme (CARS) that captures road-utilization-dependent travel times at a mesoscopic level via a piecewise affine approximation of the Bureau of Public Roads (BPR) model. This approximation largely retains the key features of the BPR model, while allowing the design of a real-time, convex quadratic optimization algorithm to determine congestion-aware routes for an AMoD fleet. Through a real-world case study of Manhattan, we compare CARS to existing routing approaches, namely a congestion-unaware and a threshold congestion model. Numerical results show that CARS significantly outperforms the other two approaches, with improvements in terms of travel time and global cost in the order of 20%.

I. 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 reflect this situation. According to [4] the annual delay per commuter exceeds 42 hours and one third of commuter trips are reported to have “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

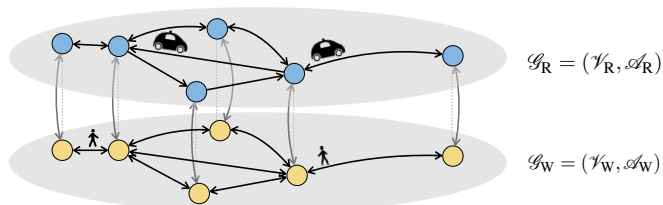


Fig. 1. The AMoD network. The colored dots represent intersections and the black arrows denote road links or pedestrian pathways. The gray dotted lines highlight geographically equivalent nodes connected by gray mode-switching arcs.

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 can be operated in a system optimal fashion.

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

¹Institute for Dynamic Systems and Control, ETH Zürich, 8092 Zürich, Switzerland, maurosalar@idsc.mavt.ethz.ch

²Autonomous Systems Lab, Stanford University, Stanford, CA 94305, United States, {mwtsao, izzya, pavone}@stanford.edu

³TUM School of Management, Technical University of Munich, 80333 Munich, Germany, maximilian.schiffer@tum.de

apart from the status quo, we briefly review related literature and state our aims and scope of the paper.

A. Related Literature

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

B. Aims and Scope

To resolve the drawbacks outlined above, we propose a congestion-aware routing scheme (CARS) that leverages a piecewise-affine approximation of the BPR congestion model [9]. In contrast to extensively studied approaches that exploit a user equilibrium (e.g., cf. [24]), we exploit the possibility to centrally control AMoD fleets in a system-optimal fashion. We study the impact of such an algorithm in a real-world case study of Manhattan. We compare CARS to two different baselines, the first one congestion-unaware and the second one capturing congestion via a threshold model. We study CA traffic routing in Manhattan and show that CARS significantly improves the state of the art with respect to travel times and global cost, while featuring computation times compatible with a real-time implementation. It should be noted that, while the BPR function does not capture every effect of congestion (e.g., spillback, heterogeneous vehicles or intersection delays), it is a well-accepted *model* which suits the aim of this paper, namely, not to perfectly capture traffic dynamics, but to approximate them precisely enough for optimization and control purposes.

C. Organization

The remainder of this paper is structured as follows: in Section II we detail the methodology and develop CARS:

a flow-based optimization framework for CA routing of AMoD systems. In Section III we numerically evaluate the performance of CARS through a case study of AMoD traffic routing in Manhattan. We conclude the paper in Section IV with a discussion of our results and an outlook on future research directions.

II. METHODOLOGY

This section provides the methodological foundation for CARS. We introduce a multi-commodity flow model to represent the physical constraints of the transportation system and AMoD fleet in Section II-A. In Section II-B we specify an objective function for the model with the goal to optimize social welfare.

We present an approximation of the BPR volume-delay function to consider congestion while preserving model convexity in Section II-C, and conclude in Section II-D with a brief discussion.

A. Multi-commodity Flow Based Optimization Approach

Recall from Section I-B, that we aim to model a network with two modes of transportation: walking and riding AMoD. This transportation network can be modeled on a digraph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$ representing the “supernetwork” [24] 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 capture both modes of transportation, \mathcal{G} comprises a road network layer $\mathcal{G}_R = (\mathcal{V}_R, \mathcal{A}_R)$ and a walking layer $\mathcal{G}_W = (\mathcal{V}_W, \mathcal{A}_W)$. The road network layer represents intersections $i \in \mathcal{V}_R$ and road links $(i, j) \in \mathcal{A}_R$, while the pedestrian layer models walkable streets $(i, j) \in \mathcal{A}_W$ in between intersections $i \in \mathcal{V}_W$. Whereas the locations of nodes and arcs in \mathcal{G}_R and \mathcal{G}_W may coincide geographically, we maintain a distinction between walking and riding in an AMoD vehicle. Additionally, switching 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 the customer’s ability to switch transportation modes by hailing an AMoD ride or exiting a car. Collectively, it holds $\mathcal{V} = \mathcal{V}_W \cup \mathcal{V}_R$ and $\mathcal{A} = \mathcal{A}_W \cup \mathcal{A}_R \cup \mathcal{A}_C$.

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

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

where the aerodynamic drag is determined by the air density ρ_a , the frontal area A_f , and the drag coefficient c_d , and the friction of the wheels on the road is determined by the rolling friction coefficient c_r , the mass of the vehicle m_v , and the gravitational acceleration g [25].

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

This notation is sufficient to derive a basic multi-commodity flow model for our planning problem at hand. Consider the cost function J mapping the set of flows $\{f_m(\cdot, \cdot)\}_m, f_0(\cdot, \cdot)$ into the set of non-negative real numbers $\mathbb{R}_{\geq 0}$. We state the AMoD optimal routing problem as

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

s.t.

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

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

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

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

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

B. AMoD Objective

We use the cost function (2a) to model the social cost of serving the transportation requests in a similar fashion as [22]. In particular, we aim to minimize the total travel time and the operational costs of the AMoD system. We assume customers to have the same value of time V_T . We separate the cost of operating the AMoD fleet into a distance-dependent cost V_D due to depreciation and maintenance, and an energy-consumption-dependent cost V_E . Collectively, the cost function (2a) is

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

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 arcs is modeled as a function of road usage, while $t_{ij} = t_{ij}^N$ gives the time to walk or switch transportation layer. To model link-congestion on road arcs we use the BPR volume-delay function [9]

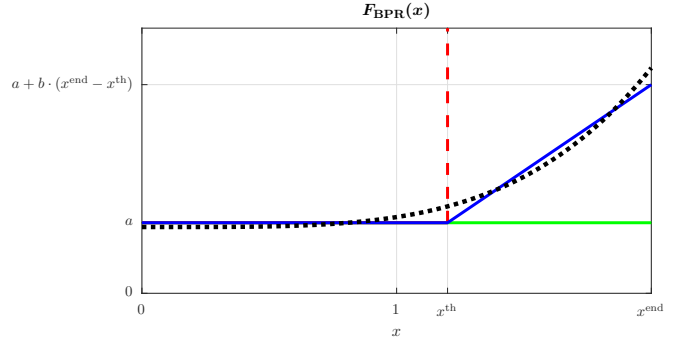


Fig. 2. The BPR function $F_{\text{BPR}}(x)$ (black dotted), its congestion-unaware approximation (green solid), the threshold model (red dashed) and its piecewise affine fit (blue solid).

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

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

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

C. BPR Model Approximation

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

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

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

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

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

$$=: t_{ij}^0 + \tau_{ij} \cdot \varepsilon(i, j),$$

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 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,

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

Provided that J is an increasing function of $\varepsilon(i, j)$, this piecewise affine cost can be represented by the following linear inequality constraints

$$\begin{aligned}\varepsilon(i, j) &\geq f(i, j) + u_{ij}^R - c_{ij}^{R, \text{th}} \\ \varepsilon(i, j) &\geq 0.\end{aligned}\quad (9)$$

It follows that the average travel time on road arcs T^R can be expressed as

$$\begin{aligned}T^R \cdot \sum_{m \in \mathcal{M}} \alpha_m &= \sum_{(i, j) \in \mathcal{A}_R} t_{ij}(f(i, j)) \cdot \sum_{m \in \mathcal{M}} f_m(i, j) \\ &= \sum_{(i, j) \in \mathcal{A}_R} t_{ij}^0 \cdot \sum_{m \in \mathcal{M}} f_m(i, j) \\ &\quad + \sum_{(i, j) \in \mathcal{A}_R} \tau_{ij} \cdot \varepsilon(i, j) \cdot \sum_{m \in \mathcal{M}} f_m(i, j) \\ &= \sum_{(i, j) \in \mathcal{A}_R} t_{ij}^0 \cdot \sum_{m \in \mathcal{M}} f_m(i, j) \\ &\quad + \sum_{(i, j) \in \mathcal{A}_R} \tau_{ij} \cdot \varepsilon(i, j) \cdot \left(\varepsilon(i, j) + c_{ij}^{R, \text{th}} - u_{ij}^R - f_0(i, j) \right) \\ &\leq \sum_{(i, j) \in \mathcal{A}_R} t_{ij}^0 \cdot \sum_{m \in \mathcal{M}} f_m(i, j) \\ &\quad + \sum_{(i, j) \in \mathcal{A}_R} \tau_{ij} \cdot \left(\varepsilon(i, j)^2 + \varepsilon(i, j) \cdot (c_{ij}^{R, \text{th}} - u_{ij}^R) \right) \\ &=: \widehat{T}^R \cdot \sum_{m \in \mathcal{M}} \alpha_m.\end{aligned}\quad (10)$$

We would like to include the average road trip time T^R in the objective for the optimization, but due to the bilinear terms $-\varepsilon(i, j) f_0(i, j)$, T^R is not a convex function of the flow variables. As an alternative, we find a convex approximation of T^R to include in the objective function. Since $\varepsilon(i, j) \cdot f_0(i, j) \geq 0$, removing these bilinear terms makes the expression larger, hence we have $T^R \leq \widehat{T}^R$. Without the bilinear terms, \widehat{T}^R is a convex function of the flow variables, so we include it in the objective function to penalize strategies with long trip times. Considering that the number of rebalancing vehicles has a minor impact with respect to road congestion and converges to zero for perfectly symmetric demand distributions [21], \widehat{T}^R can be used as a metric for the total travel time on road arcs. Specifically, our empirical studies observed that $|T^R - \widehat{T}^R|/T^R \leq 0.03$. In doing so, we approximate the total cost function (2a) with the quadratic bound (10) as

$$\begin{aligned}J_Q(\{f_m(\cdot, \cdot)\}_m, f_0(\cdot, \cdot)) &= V_T \cdot \left(\sum_{\substack{m \in \mathcal{M}, \\ (i, j) \in \mathcal{A}}} t_{ij}^N \cdot f_m(i, j) \right. \\ &\quad \left. + \sum_{(i, j) \in \mathcal{A}_R} \tau_{ij} \cdot \left(\varepsilon(i, j)^2 + \varepsilon(i, j) \cdot (c_{ij}^{R, \text{th}} - u_{ij}^R) \right) \right) \\ &\quad + \sum_{(i, j) \in \mathcal{A}_R} (V_D \cdot d_{ij} + V_E \cdot e_{ij}) \cdot f(i, j).\end{aligned}\quad (11)$$

The AMoD optimal routing problem given by (2) with $J(\cdot, \cdot) = J_Q(\cdot, \cdot)$ subject to (9) then remains a convex quadratic program.

D. Discussion

A few comments are in order. First, we assume travel requests to be time-invariant. This assumption is reflected in densely populated urban environments where requests often change slowly compared to the average time needed to complete an individual trip [26]. Second, we use the BPR function [9] to describe the impact of road usage on travel time. While such a function does not perfectly capture microscopic traffic phenomena such as queues and traffic lights, it is a well-established *model* serving the purpose of planning CA routes. In order to embed the BPR function in a convex optimization framework, we approximate it in a piecewise affine fashion reflecting a macroscopic flow diagram approximation [?]. This approximation gives a generalization of simple threshold models that are used in classical traffic flow theory [27] and can allow for better mobility service in congested situations, as discussed in the remainder of the paper. Although a piecewise affine function cannot closely approximate a quartic polynomial such as the BPR function on its entire domain, this approximation only needs to be accurate for realistic values of vehicle flows. Furthermore, we relax the piecewise approximation and implement it in the convex optimization framework through a quadratic upper bound. Nevertheless, since the ratio of empty vehicles to passenger-carrying vehicles is usually low, this bound is tight enough, as shown in Section III. Third, CARS captures customer and vehicle routes as fractional flows and does not address the stochasticity of the exogenous traffic and customer requests. Arguably, such approximations are acceptable, given the mesoscopic perspective of our study. On the topic of operational algorithms, CARS can be directly extended to operate in real-time. Due to the computational efficiency of the scheme, CARS can be run periodically with updated real-time information about customer demand to operate in a time varying environment. Furthermore, randomized rounding routing algorithms can compute near-optimal integer-valued flows for individual customers starting from the fractional solution computed by CARS [28]. The operational implementation of CARS is further discussed in Section IV. Fourth, we neglect the impact of our routing decisions on the behaviour of the exogenous traffic base load, assuming the exogenous vehicles to follow pre-defined or habitual routes. We leave the game-theoretical problem of optimizing routes accounting for reactive traffic patterns to future research. Fifth, we assume vehicles to carry only a single customer at a time. This mode of operation is in line with current trends in mobility-on-demand systems such as taxis, Lyft, and Uber. The extension to ride-sharing is an interesting direction for future research [29]. Finally, for the sake of simplicity, we assume all customers to value time and travel comfort in the same way. However, CARS can be extended to capture multiple classes of customers using network flows which are distinct not only in the origin-destination transportation request, but also in the customers' preference profile.

III. COMPUTATIONAL STUDIES

The goal of this section is to assess CARS by evaluating and comparing its performance in a case study to existing routing approaches: a congestion-unaware scheme and a threshold-congestion model as shown in Fig. 2. We begin in Section III-A by describing the Manhattan-based case study where we conduct computational studies. In Section III-B, we compare the performance of CARS to that of a) a congestion-unaware approach and b) a method that captures congestion via a threshold model.

A. Manhattan Case Study

Our case study is based on the area of Manhattan in New York City. For this area, we use real data from taxi rides that occurred between 6:00PM and 8:00PM on March 1, 2012 (courtesy of the New York Taxi and Limousine Commission). In total this data set comprises 53,932 taxi rides. Although this number of trips is quite large, it reflects only a fraction of the travel demand between 6:00PM and 8:00PM. In 2017 ride-sharing vehicles used during this time period outnumbered taxis by a ratio of 5:1 [30]. To reflect this, we dilate the number of requests by a factor of six to emulate the total demand for ride-hailing services in Manhattan during this time window.

We derive the road network from Open Street Map data [31] and set the nominal road capacities c_{ij}^R proportional to the number of lanes on a road, times its speed limit [21]. The walking network shows similar spatial characteristics to the road network, but is complementary as all arcs are bi-directional to allow for walking in both directions, even in one-way streets. To account for exogenous privately owned vehicles on the road, we run our simulations for different values of road usage u^R , denoting the fraction of a road's free-flow capacities used by private vehicles. Overall, the supernetwork has 2974 nodes and 11163 arcs. For a more detailed description of this case study and its parameters we refer to [22].

B. Results

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

We measure the quality of CARS with two performance indicators: the average resulting trip time for each solution

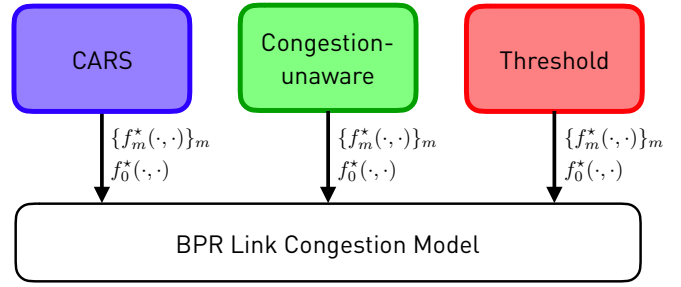


Fig. 3. Schematic representation of the comparison.

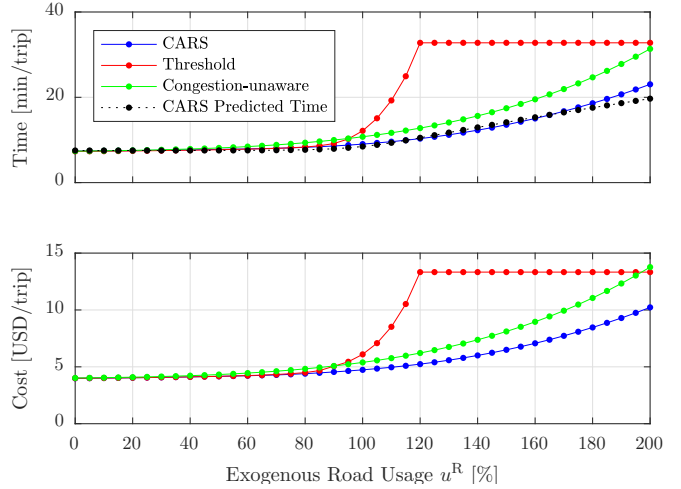


Fig. 4. Average travel time and cost per passenger trip for all approaches.

computed by the BPR function, and the objective cost as defined in (3). Fig. 4, 5 and 6 summarize the results.

Focusing on these results, we identify different relationships between the performance of all approaches depending on the exogenous road usage. Additionally, the travel times predicted by CARS are close to those computed by the BPR model, as shown in Fig. 4 and 5.

For low levels of road usage (i.e., u^R close to zero), all control schemes show the same performance: In this scenario roads have enough capacity for AMoD vehicles in the absence of private vehicles. As the traffic level increases, the congestion-unaware scheme performs worse than the thresholded approach and CARS in terms of both average travel time and operational cost. This shows that the congestion-unaware approach, lacking information about a street's capacity, sends the cars on the spatially shortest paths. As a consequence these paths become congested and remain the spatially but not temporally shortest paths due to increased travel times. Conversely, the threshold model and CARS perform similarly better for this range.

For heavy levels of traffic CARS significantly outperforms the other two approaches. Especially over 100% exogenous road usage, CARS offers a more than 20% shorter average trip time than the congestion-unaware approach and a more than 35% shorter average trip time than the threshold approach, which converges to a walking-only solution, once the value of u^R exceeds $c^{R,th}$ as shown in Fig. 6. Note

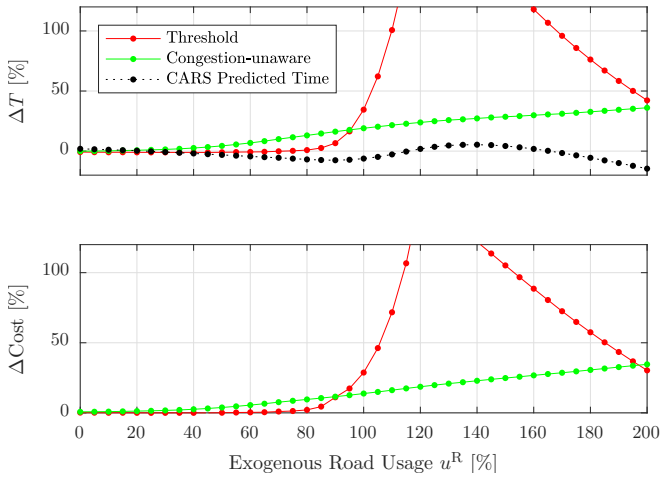


Fig. 5. Relative deviation of the two baseline solutions compared to CARS. The upper plot also shows the relative error of the CARS predicted time.

that for heavy exogenous road usage, the threshold model performs even worse than the congestion-unaware model. This is because the threshold approach is not allowed to use roads once the flow on it reaches $c^{R,th}$. In this case, the only other option is to send customer carrying vehicles on long detours and increase walking distances. In reality, $c^{R,th}$ is not a hard constraint and sending the car along a congested road can still be faster than the aforementioned long detours. Fig. 7 shows the congestion maps for this scenario and the three routing approaches, where the shortcomings of the congestion-unaware and the threshold scheme are highlighted together with the balanced behaviour of CARS.

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

On the other hand, CARS improves the operational cost and the travel time compared to both baseline approaches for all levels of exogenous road usage, as the error introduced by the piecewise approximation (7) and the convex relaxation (10) is well bounded below 15%. In dire situations with heavy traffic (when CA approaches are most needed) CARS achieves improvements of exceeding 20% in terms of travel time, striking an effective compromise between congesting roads and detouring vehicles, as shown in Fig. 7.

To conclude, our experiments show that CARS can significantly reduce both operational cost and the duration of customer trips. On the other hand, although conservative (thresholded) models of congestion offer a realistic evaluation of travel constraints at low- to mid-levels of exogenous traffic, at high-levels they perform worse than a congestion-unaware approach that ignores congestion effects altogether. Most significantly, the overconservative nature of the threshold model leads to an increase of more than \$0.50 USD per

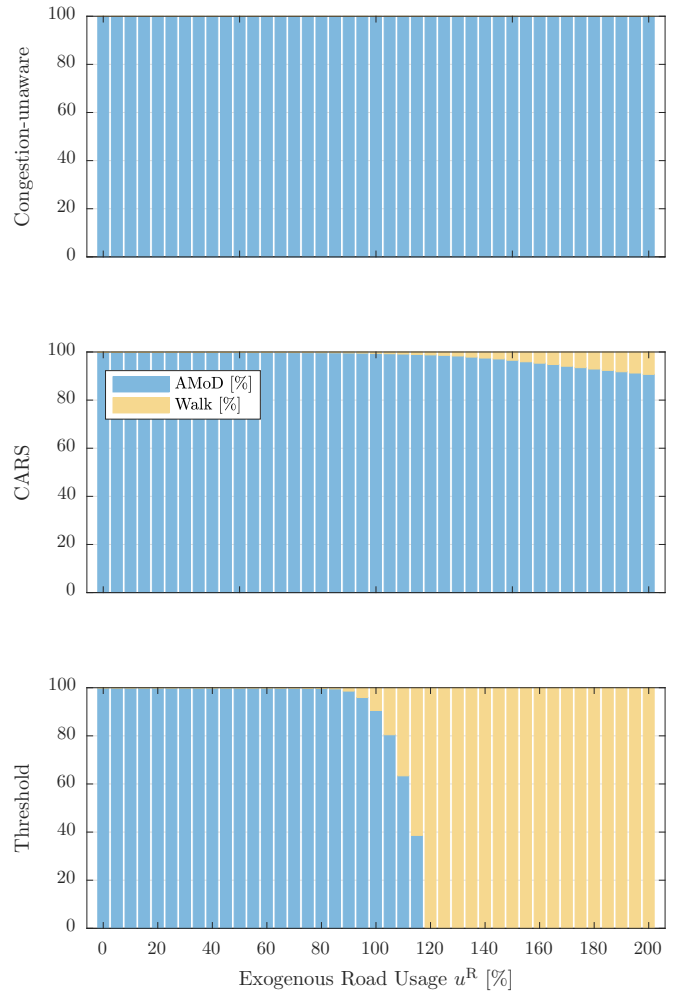


Fig. 6. Modal share for the congestion-unaware approach, CARS and the threshold model.

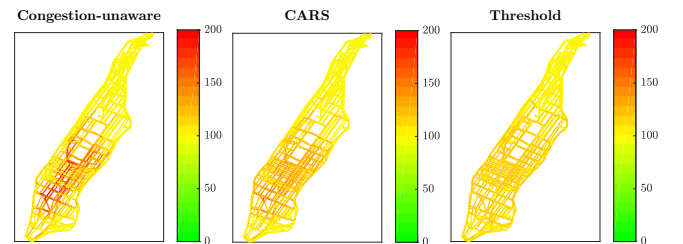


Fig. 7. Road congestion map for the congestion-unaware approach, CARS, and the threshold model with a 100% exogenous road usage level.

trip. Conversely, the balanced nature of CARS leads to a decrease of the same amount per trip. With more than half a million trips occurring on average each day in NYC [32], this efficiency could save over one hundred million dollars over the course of one year.

IV. CONCLUSION

In this paper we studied the impact of congestion-aware routing in Autonomous Mobility-on-Demand systems. By leveraging an affine approximation to the Bureau of Public Roads model, we derived a congestion-aware routing scheme

(CARS) that entails solving a computationally-efficient, convex quadratic program. We applied CARS to a case study of AMoD traffic routing in Manhattan and compared its performance to two baselines: a congestion-unaware routing approach and a scheme that considers congestion via a threshold model. Our results show that CARS always outperforms the status quo baselines and achieves an effective balance between congesting roads and detouring vehicles, thus achieving significant improvements, especially under high-levels of traffic. Additionally, we showed that in certain cases threshold models may even worsen the system performance compared to a congestion-unaware approach. Whereas considering congestion as a hard constraint can improve performance under low- to medium-levels of traffic, it can be overconservative and therefore detrimental in terms of travel time for high-levels of congestion.

This work opens several avenues for future research. Foremost, we would like to design an operational algorithm to apply CARS at a microscopic level. We note that, because the mesoscopic solutions presented here were computed in less than four minutes, the extension to an operational algorithm is completely attainable. Namely, CARS can feasibly become operational by continually solving for optimal flows given updated, real-time customer demands in a receding-horizon model predictive fashion (cf. [33]–[35]). Furthermore, as discussed in Section II-D, near-optimal integer-valued flows for operational implementation can be achieved from the fractional solutions computed with CARS via random sampling algorithms [28]. Additional future research includes coupling CARS with urban infrastructure, e.g., the power grid [23] and public transit [22]; including stochastic effects, e.g., of customer demand and road congestion [36]; and adapting the cost function to include different metrics based on customer preferences.

ACKNOWLEDGMENTS

We would like to thank Dr. Guido Gentile and Dr. Daniele Vigo for providing useful insights on congestion models, and Dr. Ilse New and Mr. Nicolas Lanzetti for proofreading this paper and providing us with their comments and advice. The first author would like to express his gratitude to Dr. Chris Onder for his support. This research was supported by the National Science Foundation under CAREER Award CMMI-1454737 and the Toyota Research Institute (TRI). This article solely reflects the opinions and conclusions of its authors and not NSF, TRI, or any other entity.

REFERENCES

- [1] (2018) The world factbook. Central Intelligence Agency. Central Intelligence Agency. Available at <https://www.cia.gov/library/publications/the-world-factbook/fields/2212.html>.
- [2] F. Siddiqui. (2018) Failing transit ridership poses an ‘emergency’ for cities, experts fear. The Washington Post. The Washington Post. available online.
- [3] R. Molla. (2018) Americans seem to like ride-sharing services like Uber and Lyft. But it’s hard to say exactly how many use them. Recode. Recode. Available at <https://www.recode.net/2018/6/24/17493338/ride-sharing-services-uber-lyft-how-many-people-use>.
- [4] Bureau of Transportation Statistics, “Transportation statistics annual report 2017,” U.S. Dept. of Transportation, Tech. Rep., 2017.

- [5] W. Hu. (2017) Your uber car creates congestion. should you pay a fee to ride? The New York Times. Available online.
- [6] K. Spieser, K. Treleven, R. Zhang, E. Frazzoli, D. Morton, and M. Pavone, “Toward a systematic approach to the design and evaluation of Autonomous Mobility-on-Demand systems: A case study in Singapore,” in *Road Vehicle Automation*. Springer, 2014.
- [7] (2018) Live map. Waze Mobile. Waze Mobile. Available at <https://www.waze.com/livemap>.
- [8] (2018) Google maps platform: Routes. Google. Google. Available at <https://cloud.google.com/maps-platform/routes/>.
- [9] Bureau of Public Roads, “Traffic assignment manual,” U.S. Dept. of Commerce, Urban Planning Division, Tech. Rep., 1964.
- [10] H. Spiess, “Technical note—conical volume-delay functions,” *Transportation Science*, vol. 24, no. 2, pp. 153–158, 1990.
- [11] A. Rivas, G. Inmaculada, S. Sánchez-Cambronero, R. M. Barba, and L. Ruiz-Ripoll, “A continuous dynamic traffic assignment model from plate scanning technique,” *Transport Research Procedia*, vol. 18, pp. 332–340, 2016.
- [12] J. Aslam, S. Lim, and D. Rus, “Congestion-aware traffic routing system using sensor data,” in *Proc. IEEE Int. Conf. on Intelligent Transportation Systems*, 2012.
- [13] E. Manley, T. Cheng, A. Penn, and A. Emmonds, “A framework for simulating large-scale complex urban traffic dynamics through hybrid agent-based modelling,” *Computers, Environment and Urban Systems*, vol. 44, pp. 27–36, 2014.
- [14] M. Treiber, A. Hennecke, and D. Helbing, “Microscopic simulation of congested traffic,” in *Traffic and Granular Flow '99*. Springer Berlin Heidelberg, 2000.
- [15] Q. Yang and H. N. Koutsopoulos, “A microscopic traffic simulator for evaluation of dynamic traffic management systems,” *Transportation Research Part C: Emerging Technologies*, vol. 4, no. 3, pp. 113–129, 1996.
- [16] M. Balmer, M. Rieser, K. Meister, D. Charypar, N. Lefebvre, and K. Nagel, “MATSim-t: Architecture and simulation times,” in *Multi-Agent Systems for Traffic and Transportation Engineering*, 2009, ch. 3.
- [17] D. J. Fagnant and K. M. Kockelman, “The travel and environmental implications of shared autonomous vehicles, using agent-based model scenarios,” *Transportation Research Part C: Emerging Technologies*, vol. 40, pp. 1–13, 2014.
- [18] M. W. Levin, K. M. Kockelman, S. D. Boyles, and T. Li, “A general framework for modeling shared autonomous vehicles with dynamic network-loading and dynamic ride-sharing application,” *Computers, Environment and Urban Systems*, vol. 64, pp. 373 – 383, 2017.
- [19] B. N. Janson, “Dynamic traffic assignment for urban road networks,” *Transportation Research Part B: Methodological*, vol. 25, no. 2–3, pp. 143–161, 1991.
- [20] S. Peeta and H. S. Mahmassani, “System optimal and user equilibrium time-dependent traffic assignment in congested networks,” *Annals of Operations Research*, vol. 60, no. 1, pp. 81–113, 1995.
- [21] F. Rossi, R. Zhang, Y. Hindy, and M. Pavone, “Routing autonomous vehicles in congested transportation networks: Structural properties and coordination algorithms,” *Autonomous Robots*, vol. 42, no. 7, pp. 1427–1442, 2018.
- [22] M. Salazar, F. Rossi, M. Schiffer, C. H. Onder, and M. Pavone, “On the interaction between autonomous mobility-on-demand and the public transportation systems,” in *Proc. IEEE Int. Conf. on Intelligent Transportation Systems*, 2018, in Press. Extended Version, Available at <https://arxiv.org/abs/1804.11278>.
- [23] F. Rossi, R. Iglesias, M. Alizadeh, and M. Pavone, “On the interaction between Autonomous Mobility-on-Demand systems and the power network: Models and coordination algorithms,” in *Robotics: Science and Systems*, 2018, Extended version available at <https://arxiv.org/abs/1709.04906>.
- [24] Y. Sheffi, *Urban Transportation Networks: Equilibrium Analysis with Mathematical Programming Methods*. Prentice-Hall, Englewood Cliffs, New Jersey, 1985.
- [25] L. Guzzella and A. Sciarretta, *Vehicle Propulsion Systems*. Springer Berlin Heidelberg, 2007.
- [26] H. Neuburger, “The economics of heavily congested roads,” *Transportation Research*, vol. 5, no. 4, pp. 283–293, 1971.
- [27] J. G. Wardrop, “Some theoretical aspects of road traffic research,” *Proc. of the Institution of Civil Engineers*, vol. 1, no. 3, pp. 325–362, 1952.
- [28] F. Rossi, “On the interaction between Autonomous Mobility-on-Demand systems and the built environment: Models and large scale co-

- ordination algorithms,” Ph.D. dissertation, Stanford University, Dept. of Aeronautics and Astronautics, 2018.
- [29] M. Tsao, D. Milojevic, C. Ruch, M. Salazar, E. Frazzoli, and M. Pavone, “Model predictive control of ride-sharing autonomous mobility on demand systems,” in *Proc. IEEE Conf. on Robotics and Automation*, 2019, submitted.
 - [30] R. Sugar. (2017) Uber and Lyft cars now outnumber yellow cabs in NYC 4 to 1. Curbed. Vox Media, Inc. Available at <https://ny.curbed.com/2017/10/13/16468716/uber-yellow-cab-nyc-surpass-ridership>.
 - [31] M. Haklay and P. Weber, “OpenStreetMap: User-generated street maps,” *IEEE Pervasive Computing*, vol. 7, no. 4, pp. 12–18, 2008.
 - [32] T. Warerkar. (2017) Uber surpasses yellow cabs in average daily ridership in NYC. Curbed. Vox Media, Inc. Available at <https://ny.curbed.com/2017/10/13/16468716/uber-yellow-cab-nyc-surpass-ridership>.
 - [33] K. T. Seow, N. H. Dang, and D. H. Lee, “A collaborative multiagent taxi-dispatch system,” *IEEE Transactions on Automation Sciences and Engineering*, vol. 7, no. 3, pp. 607–616, 2010.
 - [34] M. Pavone, S. L. Smith, E. Frazzoli, and D. Rus, “Robotic load balancing for Mobility-on-Demand systems,” *Int. Journal of Robotics Research*, vol. 31, no. 7, pp. 839–854, 2012.
 - [35] R. Zhang, F. Rossi, and M. Pavone, “Model predictive control of Autonomous Mobility-on-Demand systems,” in *Proc. IEEE Conf. on Robotics and Automation*, 2016.
 - [36] M. Tsao, R. Iglesias, and M. Pavone, “Stochastic model predictive control for autonomous mobility on demand,” in *Proc. IEEE Int. Conf. on Intelligent Transportation Systems*, 2018, in Press. Extended Version, Available at <https://arxiv.org/pdf/1804.11074>.