Microtransit deployment portfolio management using simulationbased scenario data upscaling

Srushti Rath, Bingqing Liu, Gyugeun Yoon, Joseph Y. J. Chow*
C2SMART University Transportation Center, New York University Tandon School of
Engineering, Brooklyn, NY, USA

*Corresponding author email: joseph.chow@nyu.edu

Abstract

Due to transportation technologies having such heterogeneous impacts on different communities, there needs to be better tools to evaluate the deployment of emerging technologies with limited data. Microtransit is one such technology. We propose a novel framework based on existing methods to "upscale" the limited data available so that further decision-support analysis and forecast modeling can be achieved where none could prior. The framework involves expanding an initial day-to-day adjustment process to handle both first/last mile access trips and direct trips, updating a within-day microtransit simulator with a parametric design, and developing a synthetic scenario generation process. The framework is tested in a case study with data from Via for Salt Lake City, Austin, Cupertino, Sacramento, Columbus, and Jersey City showing an average 18% ridership error for the market equilibrium models. Data from four of those cities are upscaled to 326 synthetic scenarios to estimate forecast models for ridership and fleet vehicle-miles-traveled using Lasso regularization. While the models have root mean squared error (RMSE) values between 37-45% of the averages, using only four cities' data alone would not produce any forecast model at all. The results show that variables with statistically significant positive impact on ridership and negative impact on vehicle-miles-traveled (VMT) include zones with more transit stations, higher employment, but lower "employment density × fixed fare". The models are then used to identify two alternative portfolios with similar fleet VMT as the original four cities but are forecast to have up to 1.9 times the ridership.

Keywords: microtransit, portfolio management, scenario generation, simulation, emerging technology deployment

To appear in Transportation Research Part A

1. Introduction

Transportation technologies are not "one-size-fits-all" solutions in general because their effectiveness depends on the deployment region. On-demand transit, i.e., "microtransit", exhibits this characteristic. Microtransit can be defined as shared public or private sector transportation services that offer fixed or dynamically allocated routes and schedules in a demand-responsive manner i.e., in response to individual or aggregate consumer demand, using smaller vehicles (multi-passenger/pooled shuttles or vans) and capitalizing on widespread mobile GPS and internet connectivity (Volinski, 2019; Chow et al., 2020; Yoon et al., 2022). The broader market of demand-responsive transportation (e.g., shared taxis, ride-sourcing, carshare, micro mobility, microtransit) has gained significant interest in the global urban mobility sector because of these mobile technologies.

Since these technologies are not one-size-fits-all, the reception for such technologies have been mixed. Some ventures have been successful. For example, Via Transportation, Inc. (founded in 2012) (Via, 2021) continues to operate at full capacity in over 35 countries in partnership with over 90 transit agencies (see Figure 1(a)). Their services include door-to-door, first-last mile trips to transit stations, and virtual stops, i.e., locations for pickups and drop-offs of riders within a walkable distance from their origins and destinations to improve service efficiency (Moovit, 2021). Transdev (2021), founded in 2011, operates multiple microtransit services (including first-last mile services) in the U.S, the Netherlands, France, and Australia. Shotl (2021), founded in 2017, provides on-demand bus and van services in collaboration with governments, municipalities, and businesses across Europe, Latin America, and Asia Pacific region improving accessibility in low-density and underserved areas.

On the other hand, there have been high profile failures as well: Kutsuplus in Helsinki (Haglund et al., 2019), Car2Go in North America (Krok, 2016), Bridj (Bliss, 2017), and Chariot (Marshall, 2019). Effectiveness of such service adoptions varies from city to city in terms of cost and benefit. Currie and Fournier (2020) provide a lifespan analysis of 120 demand-responsive transportation systems (including microtransit) from 19 countries over the period 1970-2019; their analysis highlights the failure rates in the UK is 67% while that in Europe and the USA/Canada is 23% and 50%, respectively.

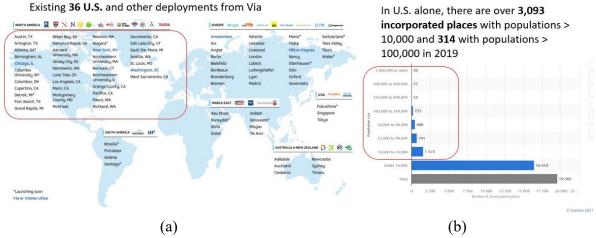


Figure 1. (a) Via deployments around the world (Via, 2021); (b) number of incorporated places in 2019 (Statista, 2019).

The list of cities shown in Figure 1(a) represents an example of a "microtransit deployment portfolio". A portfolio consists of a list of active projects developing a product sharing common resources that is continuously updated; new projects need to be evaluated and prioritized and existing projects may be accelerated, abandoned, or de-prioritized (Cooper et al., 1998; Chow et al., 2011). With microtransit deployment as the portfolio product, how can mobility providers decide which city agencies to work with for deploying new services? While the question is posed to mobility providers, the insights from this research can enable state and federal government agencies (e.g., Federal Transit Administration) to prioritize technology adoption efforts among their many cities. This question is also not restricted to just microtransit and can apply to any emerging transportation technology where deployment data are limited to just a few cities, i.e., a technology that operates in only 5 cities would only represent 5 statistical samples to extrapolate insights for hundreds of thousands more (see Figure 1(b)).

To address this microtransit deployment portfolio problem, a solution is needed that can make the most of the data that may be available. We propose "upscaling" the available data using simulation, in a similar manner to how deep learning algorithms can be used to upscale low-quality images into high-quality ones. This topic of synthetic scenario generation has also been applied to generating test cases for machine learning models, particularly in testing autonomous vehicle algorithms (Rocklage et al., 2017; Tuncali et al., 2018; Nalic et al., 2020). Our framework is novel in several new ways.

First, the market equilibrium model is extended from earlier works (Chow et al., 2020; Djavadian & Chow, 2017a; Djavadian & Chow, 2017b; Caros & Chow, 2021) to allow parameterizing the degree of virtual stop access distance and outputting the percent of microtransit riders using it as a first/last mile access mode.

Second, a proposed synthetic scenario generator outputs data that is shown to fit the limited observed deployment data, extending a sample set of four city scenarios into 326 synthetic scenarios. This "upscaled data" is fit with a forecast model that reveals interesting insights relating a deployment's ridership and fleet vehicle-miles-traveled (VMT) to service region design, pricing policy, and proximity of fixed route transit stations. Statistically significant public city attributes found using Lasso regularization include employment density, household density, mean income, street density, and transit station density, among others.

Third, the fitted forecast models are demonstrated as tools to analyze and compare different portfolios to the existing Via deployments in the four cities studied. We show it is possible to identify alternative portfolios operating a similar amount of fleet VMT that can improve ridership by up to 90%. The method can be readily adapted to any emerging transportation technology deployment planning process in which the number of deployed cities remains limited.

2. Literature review

2.1 Forecast models for microtransit

With emerging transportation technologies, the transportation planning perspective shifts from the perspective of a single city (conventional long range transportation planning conducted by a metropolitan planning organization for their own city) to a market of multiple cities. Forecasts need to be made for multiple different cities and for different operating modes. Conventional forecasting practices (Volinski, 2019; Chow et al., 2020; Yoon et al., 2022) only consider a local public agency's perspective (i.e. *city/region-centric*), which are not applicable to the deployment

portfolio planning problem which requires an *operator-centric* view that covers a cross-section of cities. Models based on cross-sectional data for forecasting microtransit measures across multiple cities simply do not exist because there was no need for such in conventional city/region-centric practices (each city had their own data to work with), to the best of our knowledge.

Forecast models for individual cities are also limited, and for good reason. Analytical models tend to resort to simplified operations and homogeneous conditions (Daganzo & Ouyang, 2019) or are used for explaining *ex post* conditions (Haglund et al., 2019; Pinto et al., 2020; Pantelidis et al., 2020; Bardaka et al., 2020; Ma et al., 2021). Microtransit can have many dimensions of complexity: routing, dispatch, pricing, rebalancing, fleet sizing, service region coverage, etc. (Fu and Chow, 2022; Dong et al., 2022). Four step models are not equipped to make predictions for users based on these complex factors mainly because that equilibrium cannot be easily captured in a static model that exhibits not only route and mode choice, but also transfers, wait time, and departure time choice from users, plus a host of choice dimensions from the operators. To overcome this drawback, city simulations draw on complex multi-agent simulations of activity behavior (Chow & Djavadian, 2015; Cich et al., 2017; Chow, 2018). However, these tools are computationally expensive and data hungry. Under a post-COVID era, there are further fluctuations in microtransit demand that make it harder to model (Zhou et al., 2021).

In particular, forecasting potential microtransit ridership as a first/last mile access mode (Shaheen & Chan, 2016) is of great interest to local agencies but remains an active research gap. To date there are no forecast models that distinguish between microtransit as a direct service or as a first/last mile access mode as a subset of multiple modes available to commuters in a region. For example, Yan et al. (2019) only forecast multimodal trips using ride-sourcing strictly to access public transit.

2.2 Simulation-based market equilibrium forecasting

Simulation-based methods are proven to be effective for evaluating complex mobility systems (Horn, 2002; Jung & Chow, 2019; Ma et al., 2019; Markov et al., 2021). However, many such studies only consider fixed demand to simulate the supply side "within-day" dynamics without any equilibration, i.e., with no demand feedback.

To capture the equilibrium between demand and mobility services, day-to-day adjustment mechanisms have been used to describe a transportation system through its dynamic evolution (Smith, 1984; Watling & Hazelton, 2003). Under such mechanisms, users in the system adjust their behavior iteratively each day according to past experiences. Such mechanisms can lead the system to evolve and converge to different states depending on the initial conditions and the behavior characteristics of the users (Smith et al., 2014). Day-to-day adjustment models have been used to model complex transportation systems because they explicitly capture the relationship between system state and the behavior of users (Horowitz, 1984; Mahmassani & Chang, 1986; Mahmassani, 1990; Cantarella & Cascetta, 1995). However, these earlier studies focus only on the road traffic network.

Djavadian and Chow (2017a,b) proposed an agent-based day-to-day adjustment process of flexible transport service and showed that the sampling distribution of different agent populations reaches a stochastic user equilibrium (SUE). Users' choices of mode and departure time are adjusted from day to day to maximize utility and minimize delay; operators' decisions can be captured as part of a two-sided market. Caros and Chow (2021) extended that model to capture operator learning of optimal cost weights to anticipate elastic user demand in evaluating modular autonomous vehicle fleets in Dubai.

Similar mechanisms are adopted in this study to model the market equilibrium of a multimodal transportation system with a microtransit subsystem.

3. Proposed methodology

3.1 Problem statement

In a market of cities P saturated with deployment data in every city, there are enough city observations that insights (e.g., forecast model as shown in the top dashed box in Figure 2) can be drawn between public data available for any U.S. city and measures important to the portfolio. For example, there exists a model M such that $y_p = M(x_p; \theta), p \in P$, where y_p may be aggregate daily microtransit ridership or the fleet's VMT (along with derivative measures like greenhouse gas emissions, accidents, infrastructure cost depreciations, etc.) and x_p include sociodemographic data. The problem is that for emerging transportation technologies like microtransit, especially from private operators, data needed for such an analysis or portfolio forecast model are limited in the number of city observations Q, i.e., $|Q| \ll |P|$. For example, in our study we have only data from |Q| = 6 U.S. cities (of which only 4 are similar enough for generating consistent synthetic scenarios) based on a half-year effort from Via to request permission from their local agency partners to share this data from each of their deployments with us.

To overcome this challenge, we propose to calibrate simulation-based market equilibrium models $y_s = m_s(x_s; \theta)$, $s \in Q$, for the deployment city set Q. Input variables of a set R of independent synthetic scenarios, $|R| \gg |Q|$, are generated as x_i , $i \in R_s$, $R = \bigcup_{s \in Q} R_s$. These inputs feed into the m_s to obtain y_i , $i \in R_s$, i.e., the "upscaled scenario data". The forecast model M can then be estimated for a portfolio of cities, i.e., $\hat{\theta} = M^{-1}(y_i, x_i)$, $i \in R$. The lower dashed box in Figure 2 highlights the upscaling process. The success of the framework depends on the design of the market equilibrium models m_s .

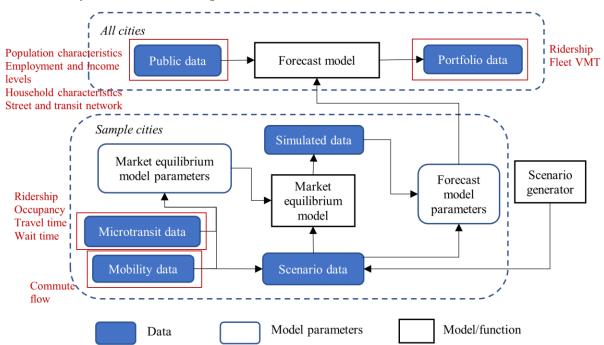


Figure 2. Process diagram showing modeling needed to generate synthetic scenario data for a portfolio-level forecast model.

Market equilibrium model m, design criteria

- 1) A market equilibrium model m_s for each city $s \in Q$ requires a demand model to predict microtransit demand. While a microtransit provider might have detailed data of their customers, they do not typically have the same quality of data for the general public across all deployed cities, and thus rely on public data. Publicly available data is only available at an aggregate level, and generally missing microtransit trips which have only recently emerged. The *demand model needs to be aggregate* and include microtransit preference that is sensitive to trip performance including some or all of access time, wait time, in-vehicle time.
- 2) Microtransit is provided for a defined service region. Pickups and drop-offs are only made within that service region. This design allows microtransit to serve direct trips as well as access/egress trips to transit stations. The *demand model needs to account for at least two population segments*: travelers making direct trips within the service region as well as trips to/from transit stations that could take travelers to locations outside the service region. Both segments depend on the service region design while the latter also depends on transit network proximity.
- 3) The market equilibrium model needs to provide a *stable set of performance measures for a dynamic, on-demand transit service system design* in which the operator perceives travelers arriving randomly throughout the day and operate several functions dynamically: dispatch, routing, idle vehicle repositioning.
- 4) The market equilibrium model needs to be designed to *account for observed occupancy data* that is available to the microtransit operator.
- 5) The market equilibrium model needs to *allow for calibration of different operating parameters*, e.g., different pedestrian access distances, maximum wait time or detour time, etc., to differentiate operations at different cities.
- 6) The synthetic scenario generator needs to be able to *efficiently construct independent scenarios* in which public zonal data are available as well as outputs of the calibrated simulation operated under those scenarios.

Our market equilibrium model and synthetic scenario generator are designed to address these criteria as shown in Table 1. The model m_s should output ridership, percent of riders that are first/last mile access, fleet size, fleet vehicle miles traveled, average traveler journey times (wait, access, in-vehicle), operation cost, revenue, and other derivative measures like greenhouse gas (GHG) emissions.

Table 1. Components of the proposed method that addresses the methodology design criteria

Components of the proposed framework	Design criteria
Aggregate mode choice model (Section 3.2.1)	#1 and #2
Day-to-day adjustment process in the market equilibrium model (Section 3.2.2)	#3 and #4
Within day simulator in the market equilibrium model (Section 3.2.3)	#5
Synthetic scenario generation design (Section 3.3)	#6

3.2 Day-to-day market equilibrium model

A day-to-day adjustment process characterizes the dynamics in adjustments made by both travelers (users) and the operators each day as a dynamic system. Djavadian and Chow (2017a) showed that

such adjustment processes can reach a stochastic user equilibrium with an asymptotic number of sampled agent populations and can be used to model on-demand systems. However, it did not incorporate multiple population segments, nor did it account for certain observable data made available from an operator, namely the occupancy data.

The process from Djavadian and Chow (2017b) is modified to address the requirements of this framework. In this process, travelers are split into two population segments: (1) travelers who wish to make a direct trip within the service region and (2) travelers who need to make an access/egress trip within the service region to connect to a public transit network that can extend to a greater region not covered by the microtransit service. The adjustment process is also modified to converge toward a similar average occupancy rate (number of passengers/vehicle/hour) as data provided from a microtransit operator. The model is shown in Figure 3.

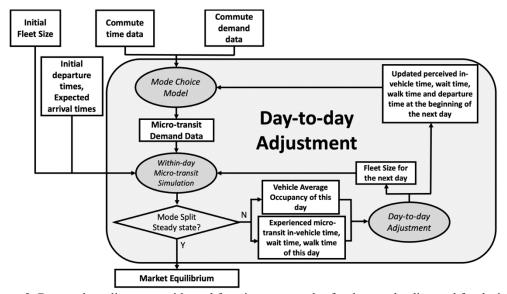


Figure 3. Day-to-day adjustment with oval functions, rectangles for data, and a diamond for decision.

At the end of one simulated day, experienced microtransit in-vehicle time, passenger wait time, passenger walk time, and average occupancy for the vehicles, are computed. These values are used to update the mode choices, departure times, and fleet sizes for the next day (note that only microtransit is simulated, so the attributes for all the other modes—Auto, Transit, Bike, Walk, Others—are fixed).

3.2.1 Aggregate mode choice model

For a microtransit demand model, an aggregate mode choice model is used to capture sensitivity of transit riders to different operational attributes that can vary day by day. If data were available on travelers choosing microtransit among a set of modes between different zones, we would ideally estimate a binary logit model for simplicity. However, such data are not available. We have count data of microtransit trips that is not divided by zones within a service region, and we have public aggregate data for other modes by zone to zone (Census tracts).

Furthermore, the only public aggregate data (as described in the Section 4.2 of the case study) provides some commute attributes but not all. For some attributes, we assume values of the parameters and estimate the rest of the parameters of the model relative to them. Because of the need for an aggregate mode choice model, more complicated choice model structures than a

multinomial logit (MNL) are not feasible (as we only know market shares by origin destination pair, not individual choices) (see Anas, 1983, for a discussion of estimating aggregate MNL models from aggregate market share data). As a result, a MNL model is estimated for multiple mode choices and the utility function for microtransit (MT) assumes the coefficients of the attributes are the same as other similar modes. For example, the MT in-vehicle time coefficient is assumed the same as the auto travel time coefficient, and the MT access/egress time coefficient is assumed to be the same as that of walk mode. The utility functions are generally specified as shown in Eq. (1) - (6) (with statistically insignificant attributes for each city removed) as a mode choice model for a given agent n.

$$U_{auto,n} = asc_{auto} + \beta_{tt_{auto}} \times TT_{auto,n} + \beta_{interzone} \times Interzone_n + \varepsilon_{auto,n}$$
(1)

$$U_{transit,n} = asc_{transit} + \beta_{tt_{transit}} \times TT_{transit,n} + \beta_{AE} \times AET_{transit,n} + \beta_{wait} \times WT_{transit,n} + \beta_{cost} \times CO_{transit,n} + \beta_{interzone} \times Interzone_n + \varepsilon_{transit,n}$$
(2)

$$U_{bike,n} = asc_{bike} + \beta_{tt_{bike}} \times TT_{bike,n} + \beta_{interzone} \times Interzone_n + \varepsilon_{bike,n}$$
 (3)

$$U_{walk,n} = asc_{walk} + \beta_{tt_{walk}} \times TT_{walk,n} + \beta_{interzone} \times Interzone_n + \varepsilon_{walk,n}$$
 (4)

$$U_{MT,n} = asc_{MT} + \beta_{tt_{auto}} \times TT_{MT,n} + \beta_{tt_{walk}} \times AET_{MT,n} + \beta_{MTwait} \times WT_{MT,n} + \beta_{cost} \times CO_{MT,n} + \beta_{interzone} \times Interzone_n + \varepsilon_{MT,n}$$
(5)

$$U_{others,n} = \varepsilon_{others,n} \tag{6}$$

where $asc_{(mode)}$ denote the mode-specific-constant for $\langle mode \rangle \in \{$ automobile (auto), transit, bike, walk, microtransit (MT), other $\}$; $TT_{(mode)}$ is modal travel times from origin to destination; $CO_{(mode)}$ denote the corresponding modal travel costs (treated as a generic coefficient since we assume it has constant effects (Train and McFadden, 1978)); $WT_{(mode)}$ refer to wait times for transit and microtransit modes; $AET_{(mode)}$ are the access/ egress time for transit and microtransit (by walking); and $Interzone_n$ is a categorical variable for interzonal trips i.e., 1 when a trip's origin and destination are in different zones (census tracts in our case study) and 0 otherwise. The attributes are tracked with an index d to represent the perceived value at the start of day d. The coefficients $\beta_{tt_{(mode)}}$, β_{AE} , β_{wait} , β_{MTwait} , β_{cost} , $\beta_{interzone}$ need to be estimated for each city (in

future research, we can estimate a model for a cluster of cities based on typology, see Oke et al., 2019 and Rath and Chow, 2022). Note that no public data is available for the out-of-pocket costs of using the auto mode, so a cost is not included for that utility function. This only implies that cost effects from auto are incorporated into the alternative specific constant and are not directly controlled for. This is not a problem since the microtransit scenarios that we consider do not include changes in auto costs.

One novel treatment of the demand model is that it includes agents from (1) direct door-to-door trips within a designated service region S as well as (2) first/last mile trips connecting with

transit stations located within the service region to other locations in the greater region $Z \supset S$. This is illustrated with Figure 4.

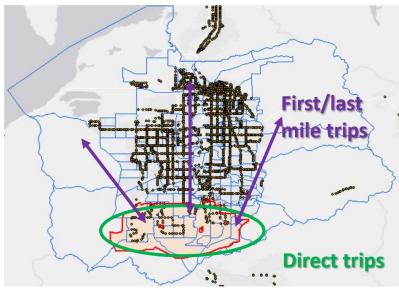


Figure 4. Illustration of a designated service region in red (via direct trips) along with blue-highlighted zones in the greater region accessed by public transit (via first/last mile trips).

Due to limited information, some parameters are assumed so that the others can be estimated in relation to them (details of specific assumed values provided in Section 4 case study).

The estimation process for the demand models is broken into four steps:

- 1) Estimate a mode choice model for each city without the microtransit alternative.
- 2) Take the parameters of the microtransit utility function attributes from other similar modes: MT cost from fixed route public transit cost (β_{cost}), MT in-vehicle time from auto travel time ($\beta_{tt_{auto}}$), and MT access/egress time from walk mode ($\beta_{tt_{walk}}$), and MT wait time (β_{MTwait}) from public transit wait time, or if not available due to non-statistical significance, this is derived as a multiplier of auto travel time (Wardman, 2004).
- 3) Estimate the alternative-specific constants of the microtransit alternative for each city to be relative to the assumed attribute parameters and to fit the microtransit trips summed over direct and first/last mile segments to the count data using least squares via a bisection method.
- 4) When calibrating the day-to-day adjustment processes, i.e., after the system first reaches a steady state (as shown in Figure 3), use one feedback loop to re-estimate the demand model coefficients using the output attributes (i.e., population experienced time values for microtransit including in-vehicle time, wait time, and walk time) from the market equilibrium model. Note that with additional feedback rounds no significant change is observed in the estimated demand model attributes, hence one feedback loop is used for the update.

3.2.2 Day-to-day adjustment process

At the start of the simulation, the origin and destination coordinates of each user are generated randomly within their origin and destination zones. The adjustment of fleet size is based on average occupancy provided by the data. At the end of each day, fleet size is adjusted towards the ideal average occupancy based on the occupancy of the past day as shown in Eq. 7. The observed occupancies from the microtransit operator are assumed to be ideal values unless specified

otherwise. By comparison, Djavadian and Chow (2017b) did not have average occupancy data so it did not include this measure. This modification allows us to run the day-to-day adjustment process to converge toward an equilibrium that exhibits similar occupancy as what is observed.

$$FS^{d+1} = FS^d \times \frac{Average\ Occupancy\ of\ day\ d}{Ideal\ Average\ Occupancy} \tag{7}$$

where FS^d stands for the microtransit fleet size on simulation day d.

Microtransit in-vehicle time, wait time and walk time (to virtual stops) for each user are updated from day to day. A traveler who used microtransit on day d for commute learns from their experience on day d and updates their perceived in-vehicle time, wait time, and walk time with a learning rate θ (shown in Eqs. (8) – (10)).

$$TT_{MT,n}^{d+1} = (1 - \theta)TT_{MT,n}^{d} + \theta \ ETT_{MT,n}^{d}$$
 (8)

$$WT_{MT,n}^{d+1} = (1 - \theta)WT_{MT,n}^{d} + \theta \ EWT_{MT,n}^{d}$$
 (9)

$$AET_{MT,n}^{d+1} = (1 - \theta)AET_{MT,n}^d + \theta EAET_{MT,n}^d$$
 (10)

where $TT^d_{MT,n}$, $WT^d_{MT,n}$, and $AET^d_{MT,n}$ stand for perceived microtransit (MT) in-vehicle time, wait time, and walk time for user n at the beginning of day d. $ETT^d_{MT,n}$, $EWT^d_{MT,n}$, and $EAET^d_{MT,n}$ stand for experienced microtransit (MT) in-vehicle time, wait time, and walk time for user n on day d.

Having introduced the key attributes, let us adopt a generic symbol X to represent each attribute for convenience. For a user n' who did not use microtransit but used other modes on day d for their commute, their perceived times are updated with the population's average perceived times \bar{X}_{MT}^d at the end of day d (shown in Eq. (11)).

$$X_{MT,n'}^{d+1} = (1 - \theta)X_{MT,n'}^d + \theta \ \bar{X}_{MT}^d \tag{11}$$

The population's perceived in-vehicle time, wait time, and walk time represent the overall perception of the population in the service region, which is the successive average of average invehicle, wait, and walk time of the past n days (shown in Eq. (12)).

$$\bar{X}_{MT}^{d} = \left(1 - \frac{1}{d}\right) \bar{X}_{MT}^{d-1} + \frac{1}{d} E \bar{X}_{MT}^{d} \tag{12}$$

The departure time of each user is a continuous variable that is updated based on their expected arrival time. The departure time of a passenger on day (d + 1) is computed based on the experienced commute time of day d as shown in Eq. (13).

$$DT_n^{d+1} = AT_n - PT_n^d (13)$$

where DT_n^d stands for the departure time of user n on simulation day d, AT_n stands for the desired arrival time of user n. PT_n^d is the perceived commute time at the end of day d for user n, which depends on the mode taken in Eq. (14).

$$PT_n^d = TT_{MT,n}^d + WT_{MT,n}^d + AET_{MT,n}^d (14)$$

At the end of each day, we check if the system has reached a steady state. The adjustment stops when the daily microtransit ridership change keeps below $\epsilon \in [0,1]$ for 5 consecutive days.

$$\frac{Ridership_{MT}^{d+1} - Ridership_{MT}^{d}}{Ridership_{MT}^{d}} \le \epsilon \tag{15}$$

We recall the proposition from Djavadian and Chow (2017a) where the agent-based day-to-day process converges almost surely to an agent-based stochastic user equilibrium (SUE), i.e., as the number of simulated populations increase, the deterministic day-to-day adjustment processes for each population is run with a Method of Successive Averages (MSA) such that the average over the populations approaches the theoretical SUE. In our case study, we generate 10 populations per city.

3.2.3 Within-day simulator

Within-day simulation of the microtransit system is a key part of the day-to-day adjustment. The main framework of this simulation (as illustrated in Figure 5) is newly extended from a "simulation sandbox" from Yoon, et al. (2021). It is a discrete-time simulation with a simulation length divided over discrete time steps. In each time step, vehicle states are updated according to the operating plan. For on-demand microtransit, the state includes the sequence of passengers being served. Passenger states are also updated: waiting to be assigned to a vehicle, walking to a stop, waiting for vehicle, on-board a vehicle, egressing from vehicle stop to the destination. The simulation length and time step can be adjusted.

The simulation decides which vehicles to dispatch and transport passengers from their origins to their destinations. The system does not order vehicles to comply with fixed routes or drop by mandatory stops. Instead, when assigning a passenger to a vehicle the system also updates the vehicle state which includes the sequence of pickup and drop-off points of passengers assigned to it. These matches determine vehicles' trajectories, passengers' travel experiences, and system performances. For the routing, the simulation sandbox implements a simple insertion heuristic to update assigned routes for accepted passengers by searching for an updated sequence with the shortest incremental increase in travel time.

Modifications are made from Yoon et al. (2022) to be more parametric for calibration purposes. These modifications include:

- Virtual stops, meeting points other than actual origin and/or destination of users, and a
- Feature of depot assignment, designating a depot of vehicles based on relocation cost and average wait time.

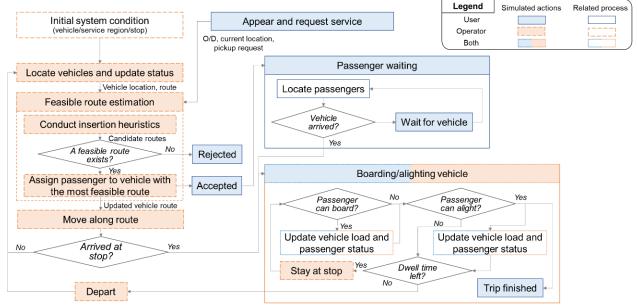


Figure 5. Framework for the within-day simulator (source: Yoon et al., 2022).

There are four categories of required inputs of this simulation as described below. Simulation parameters control the length and precision of simulations, regulating how long and often to collect generated data.

- Simulation parameters: simulation length, time step length
- Scenario parameter: walking speed, maximum walking distance, average vehicle running speed, weight for passenger in-vehicle/wait/access time, value of time, unit operation cost, weight of operator cost
- System design parameter: vehicle capacity, fleet size, number of depots, average dwell time
- Dataset: passenger request information, passenger arrival data, depot locations, virtual stop locations, vehicle allocation distribution among depots

The code to the simulator can be accessed in a Zenodo repository (Rath et al., 2021).

3.3 Synthetic scenario generator design

Once the simulation-based market equilibrium models m_s are calibrated, a process is used to generate additional synthetic scenario inputs x_i , $i \in Q$, to upscale the existing data to obtain outputs y_i . The generator should generate independent scenarios that would use the market equilibrium models calibrated for a sample of cities. In each synthetic scenario i, one of the calibrated cities is selected and the service region S is redefined as S(i). We define a region as constituting two or more contiguous zones where microtransit service operates both door-to-door and first-last mile services with a specific pricing policy. The selection of contiguous set of zones for scenario generation is mainly motivated by the Via service region structures in various cities e.g., Salt Lake City, Jersey City, New York City, Cupertino, Arlington, Birmingham, and others. Each synthetic scenario is generated using Algorithm 1, where we utilize the common (geographical) boundaries shared by a zone with its adjacent neighbors to efficiently generate service regions of different

sizes (within the limits of a larger geographical boundary). L2 and L3 boundaries are used to reflect similar service regions in scale to the sample cities.

Algorithm 1. Synthetic scenario generation

Given a city where a microtransit service operates (in a specific region), we generate multiple scenarios (regions) with the following steps:

- 1. Obtain the list of census tracts (zones) and the (geographical) boundaries for all zones within the county/counties intersecting with the existing service region S.
- 2. For each census tract (as obtained in step 1), store their neighboring census tracts (i.e., zones sharing common boundaries).
- 3. Select a zone (let's say x) and generate 3 scenarios: L1, L2, L3, where, L1 constitute the direct neighbors of zone x, then we add the neighbors of each zone in L1 to get L2, and for L3 we add neighbors of all zones in L2 to the existing L2 scenario. Figure 7 provides an illustration of service region scenario generation.
- 4. Assume a pricing policy as shown below:
 - o PP1 = fixed fare for door-to-door service and first last mile rides
 - o PP2 = fixed fare for door-to-door service and free first last mile rides
 - o PP3 = variable fare (fixed fare plus fare based on distance travelled)

Apply the most common pricing policies PP1 and PP2 to each scenario S(i) from step 3

- 5. Cluster each city's generated scenarios into 6-8 clusters (for both L2 and L3 scenarios) based on population characteristics using k-means clustering and randomly select sample scenarios from each cluster.
- 6. Run the market equilibrium model using the input scenario data generated for each corresponding sample city to obtain microtransit performance data.

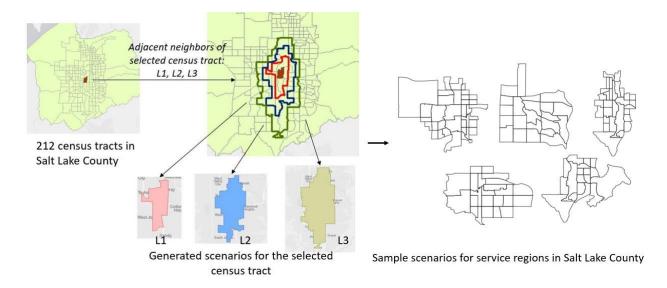


Figure 7. Illustration of service region generation process.

To provide a better idea on how the population data and simulated microtransit performance data for different synthetic scenarios look like, Figure 8 shows examples of four service regions generated in different U.S. counties using Algorithm 1. Each scenario $i \in Q$ is comprised of a set

of census tracts based on which we obtain the aggregate population data (as listed in the figure) (x_i) for the scenario. Microtransit performance metrics (e.g., ridership, vehicle miles traveled, fleet size, and others) (y_i) are obtained using the calibrated simulation model $m_s(x_i)$ for the associated pricing policy considered in the scenario.

To ensure that the new data points obtained from the scenario generation process for the forecast models cover a diverse set of synthetic scenarios with a reasonable range of the population and region characteristics, S(i) in step 5 of Algorithm 1 is selected (randomly) from different clusters. In particular, for multiple synthetic scenarios generated for a city (let's say all L2 scenarios), we categorize them into 6-8 clusters based on their population and built environment characteristics (using k-means clustering algorithm (MacQueen, 1967)). Then from each of these clusters, we randomly select 5-6 sample scenarios to obtain simulated microtransit performance data (based on pricing policies PP1 and PP2) for the selected scenarios using the calibrated market equilibrium day-to-day simulations.

The scenario generation results in upscaled scenario data that are used for estimating forecast models for service portfolio design and deployment planning as demonstrated in the case study in Section 4.

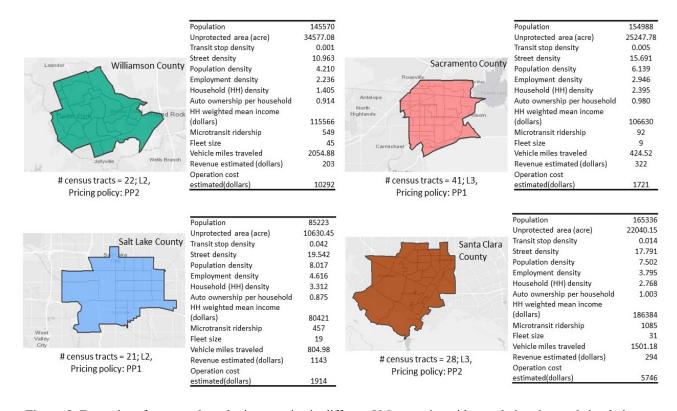


Figure 8. Examples of generated synthetic scenarios in different U.S. counties with population data and simulation data obtained for the scenarios.

4. Case study: Via microtransit deployment portfolio analysis

4.1 Study objectives

The case study is used to evaluate the methodology proposed in Section 3. Since the contribution focuses on simulating new data to supplement limited existing data for the purpose of evaluating deployments in different cities, our objective is to show that:

- 1) Forecast models can be specified (having statistically meaningful relationships between public data and ridership/VMT) using the upscaled data;
- 2) The models should adequately fit the limited data that we have.
- 3) The models can be applied to evaluate different portfolio designs.

The case study is conducted in collaboration with Via, who provided aggregate ridership data for five different U.S. cities: Salt Lake City, Austin, Cupertino, Sacramento, and Columbus. Public data from Jersey City is also available and included in the set of cities, as summarized in Table 2. The benchmark is a forecast model that is built using only the six cities' data, which is not statistically viable since that would simply be an insufficient sample size.

Table 2. Via service regions (obtained from Via) and pricing policies (defined based on Via (2021))

City (Via service region)	Counties (transit demand considered for potential first last mile trips)	Number of census tracts within service region boundary	Pricing policy
Salt Lake City, Utah	Salt Lake County	26	PP2
Austin, Texas	Travis and Williamson County	28	PP2
Cupertino, California	Santa Clara County	21	PP2
Sacramento, California	Sacramento and Placer County	148	PP2
Columbus, Ohio	Franklin and Licking County	45	PP1
Jersey City, New Jersey	Hudson and New York County	68	PP3

4.2 Data

Table 3 presents the data used in our study, which includes public data for estimation of mode choice models, simulation model calibration, and design of deployment portfolio forecast model, along with data obtained from Via. For our case study, we focus on ridership during the peak period of the day i.e., 6AM-9AM so all outputs reflect that time period. All the synthetic scenario data are provided in the Zenodo repository (Rath et al., 2021).

Table 3. Summary of data and data sources used in the study

No.	Data source	Granularity	Data
1	Census Transportation Planning Products (2012-2016) (CTPP, 2016)	Census Tracts	Commute flows between census tracts (for various modes including auto, bike, transit, walk, and others)
2	American Community Survey (2019)	Census Tracts	Demographic, economic and household details
3	Open Mobility Data (GTFS) (Transitfeeds, 2021)	Transit network	Transit station/stop locations

4	Smart Location Database (EPA, 2021)	Census Block Groups (aggregated to Census Tracts)	Household auto ownership, unprotected area, street network (road density), trip equilibrium index (trip attractions and productions)
5	Open Street Map (OSM, 2021)	Street network	Auto, Walk, and bike travel time between census tracts; walk and bike travel time to and from nearest transit stops in census tracts
6	Open trip planner (OTP, 2021)	Transit network	Transit commute time including in-vehicle time, wait time, walk time (to/from the nearest stops)
7	Via data: weekly average during first week of 3/20	Salt Lake City, Cupertino, Austin, Columbus, Sacramento, and Jersey City	Via service region boundaries, average ridership, average wait time, vehicle utilization, pricing policy, fare structure

4.3 Calibration of the market equilibrium model

4.3.1. Estimation of the demand model parameters

The calibration of the market equilibrium models involves two parts. The first is the estimation of the aggregate mode choice model, one for each of the six cities with provided data. Commute flow data for auto, transit, bike, walk, and others are obtained between census tracts within the region (CTPP, 2016) including transit flows from within service region to outside the service region (within the boundaries of the county/counties covered by the region) and vice-versa. Due to the limited aggregate data, the mode choice model initially led to poor fits in some cities. To address this, we assumed some of the parameters' relationships: $\beta_{MTwait} = \beta_{wait}$ (in transit if significant, else $\beta_{MTwait} = 1.53\beta_{tt_{auto}}$), $\beta_{wait} = 1.59\beta_{tt_{transit}}$, $\beta_{AE} = 1.78\beta_{tt_{transit}}$ (Wardman, 2004). Table

4a presents the mode choice model estimation results that are calibrated after one round of feedback from the simulation model. The p-values and ρ^2 are based on the initial estimation without microtransit. Some of the attributes are left out for some cities due to irregular fits or poor statistical significance. The microtransit alternative is then appended and estimated as shown in Table 4b.

Table 4 (a). Mode choice model coefficient estimates and performance for Via cities (without microtransit mode)

Coefficient	Units	Salt Lake City, Utah	Austin, Texas	Cupertino, California	Sacramento, California	Columbus, Ohio	Jersey City, New Jersey
asc _{auto}	N/A	0.649***	-0.145*	-	0.231***	0.330***	-
asc _{bike}	N/A	-3.318***	-4.393***	-3.934***	-2.494***	-6.555***	-4.004***
$asc_{transit}$	N/A	-1.510***	-1.956***	-0.707***	-0.682***	-1.329***	-
asc_{walk}	N/A	-1.973***	-3.909***	-2.363***	-0.312***	-1.839***	0.560***
$oldsymbol{eta}_{tt_{auto}}$	1/min	-0.204***	-0.049***	-0.131***	-0.109***	-0.009**	-0.177***
$oldsymbol{eta_{tt_{bike}}}$	1/min	-0.129***	-0.051***	-0.098**	-0.105***	-	-0.251***
$oldsymbol{eta}_{tt_{transit}}$	1/min	-0.003	-	-	-0.012***	-0.009	-0.001
$eta_{AE}^{({ m a})}$	1/min	-0.005	-	-	-0.021	-0.016	-
$eta_{wait}^{ ext{(a)}}$	1/min	-0.005	-	-	-0.019	-0.014	-0.002
$oldsymbol{eta}_{tt_{walk}}$	1/min	0.033***	-0.006***	-0.038***	-0.064***	-0.037***	-0.086***

$oldsymbol{eta_{cost}}$	1/U.S. \$	-1.851***	-2.062***	-1.768***	-1.058***	-0.998***	-0.930***
$oldsymbol{eta}_{interzone}$	N/A	8.326***	12.895	7.403***	6.987***	6.356***	5.429***
ρ^2 (w/o microtr	ansit)	0.78	0.69	0.72	0.78	0.85	0.43

^{*, **, ***} refer to p-values from initial estimation without microtransit less than 0.05, 0.01, and 0.001 respectively.

(a) Non-bolded parameters were not estimated but assumed relative to other estimated parameters as shown in text

Table 4 (b). Mode choice model coefficient estimates and evaluation for the microtransit mode

Coefficient	Units	Salt Lake City, Utah	Austin, Texas	Cupertino, California	Sacramento, California	Columbus, Ohio	Jersey City, New Jersey
asc_{MT}	N/A	0.848	-1.096	2.089	-0.689	-7.354	-2.265
$eta_{tt_{auto}}^{(a)}$	1/min	-0.204	-0.049	-0.131	-0.109	-0.009	-0.177
$eta_{tt_{walk}}^{(a)}$	1/min	0.033	-0.006	-0.038	-0.064	-0.037	-0.086
$\beta_{MTwait}^{(a)}$	1/min	-0.005	-0.075	-0.200	-0.019	-0.014	-0.002
$eta_{cost}^{ ext{(a)}}$	1/U.S. \$	-1.851	-2.062	-1.768	-1.058	-0.998	-0.930
$\beta_{interzone}^{(a)}$	N/A	8.326	12.895	7.403	6.987	6.356	5.429
Min. abs. err	or (pred. vs						
obs. Via ride	rship) with						
estimated	l asc _{Via}	0.004	0.008	0.003	0.008	0.001	0.002
Min. abs. err	or (pred. vs						
obs. Via ride	rship)						
with asc_{Via} =	0	75.56	244.11	42.91	182.83	1167.14	1908.97

⁽a) Non-bolded parameters were not estimated for microtransit utility function but assumed relative to other estimated parameters as shown in text

Table 4b compares the estimated error for each city's model when using the optimal asc_{MT} compared to a model where the $asc_{MT} = 0$. The error reduction is significant. The travel time and cost coefficients are negative in most cities, with Salt Lake City having a positive coefficient for walk time. Based on the commute flow data, the proportion of walk trips (i.e., between origin destination (OD) pairs in a service region) with respect to the total commute flows in Salt Lake City is 2-3%; this distribution is similar to Sacramento, Cupertino, Columbus and Austin. Compared to the average OD walk time in these cities, which is less than 15 minutes (average of the four cities being 7 minutes), the average walk time in Salt Lake City is 38 minutes, which is significantly higher. Therefore, this could be contributing towards the positive coefficient estimate for walk time in Salt Lake City. While this is attributed to data limitations, it has been noted in the literature that such behavior can be explained as well (Redmond and Mokhtarian, 2001). Moreover, positive asc_{MT} values observed for Salt Lake City and Cupertino indicate a positive (average) effect of latent (unincluded) factors on the utility of the microtransit (Via) in these cities, while an opposite effect is noticed in the other 4 cities. This observation highlights the effects the city type (among other latent factors) may have on the utility of microtransit in a city.

4.3.2. Calibration of the within-day simulator parameters

The day-to-day adjustment parameters are calibrated as follows. Parameters include the walking limit of microtransit users, microtransit dwell time, and user/operator weights for the insertion heuristic in microtransit within-day simulation. The performance measure for finding the best insertion option in the within-day simulation is shown in Eq. (16), which is a combined measure of the users' loss and the operator's loss balanced by operator's weight α_{oper} and user's weight

 $(1 - \alpha_{oper})$. Average operator cost per mile is estimated from the average operation cost per passenger provided by Volinski (2019) and the average trip length data provided by Via (Eq. (17)).

Performance measure

=
$$(1 - \alpha_{oper}) \times \text{Value of time} \times \text{User time increment}$$
 (16)
+ $\alpha_{oper} \times \text{Operator cost per mile} \times \text{Distance traveled increment}$

Operator cost per mile =
$$\frac{\text{Operator cost per passenger}}{\text{Average trip length (miles)}}$$
 (17)

For calibration, we assume three discrete levels for each of the three parameters:

• Walking limit: 0.5 miles, 0.3 miles, 0.1 mile

• Dwell time: 15 sec, 10 sec, 5 sec

• Operator weight in insertion heuristic: 0.8, 0.5, 0.2

Hence, 27 combinations are produced. Learning rate is set to 0.1, consistent with prior studies (see Djavadian and Chow, 2017b) and ϵ is set to 1% for our case study. We run the simulation for each of the combinations to find the best combination for each city. The best combination is selected by comparing the output average in-vehicle time, average wait time, and ridership with the data from Via. For each city, the combination which lead to the smallest sum of squared error of average in-vehicle time, average wait time, and microtransit ridership is selected as the optimal combination. Example simulations conducted for four different cities (Salt Lake City, Austin, Cupertino, and Sacramento) are shown in Figure 6.

The selected optimal combinations in the calibration process for the 6 cities are shown in Table 5, along with the corresponding errors. Jersey City had less data available so the in-vehicle and wait time errors could not be computed. The results show that the cities can vary in their characteristics. For example, Salt Lake City and Jersey City suggest longer access via walking for travelers, while Austin and Salt Lake City tend to have longer dwell times for their vehicles. Cupertino has the highest weight for operator cost, which suggests that their travelers are the least elastic to the service quality. Generally, cities with smaller walking limit have smaller operator weight, since when the users are more reluctant to walk, user's weight should be higher. In terms of error, the overall ridership error indicates fits with an average of 18.4% among the six cities. The results indicate that a market equilibrium model m_s can indeed be calibrated to different cities $s \in R$, even with only an aggregate mode choice model for demand estimated for each city.

The process of convergence for the 6 cities with the calibrated parameters are shown in Figure 9. The average computation times for one run of Salt Lake City, Cupertino, Sacramento, Columbus, Austin, and Jersey City are respectively 10min 42s, 4min, 6min 24s, 36s, 4min 42s, and 13 min on a laptop with 2.3 GHz Quad-Core Intel Core i7 and 32 GB 3733 MHz LPDDR4X memory. The results indicate that steady states do exist for these cities and that the number of days to convergence can differ from city to city.

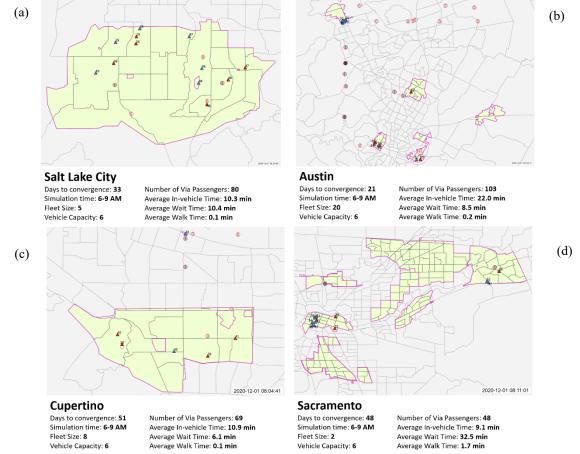
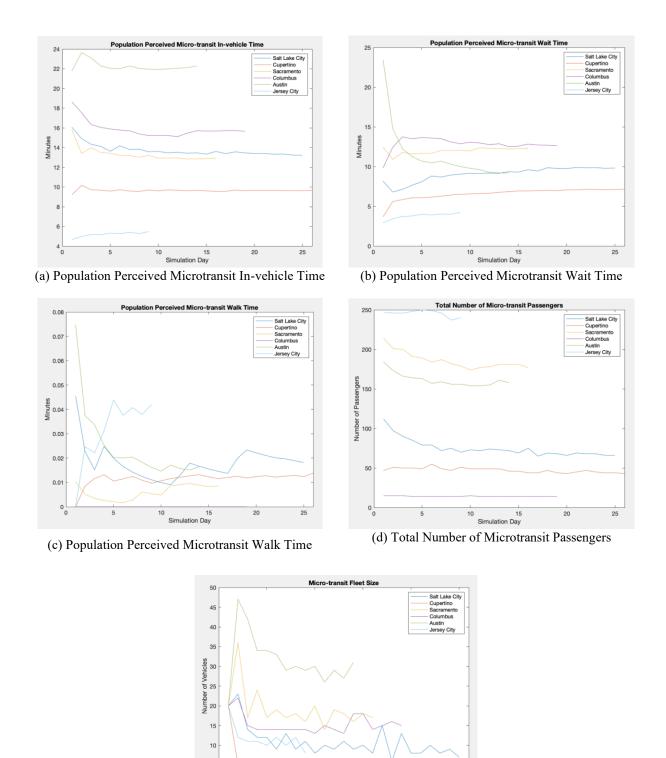


Figure 6. (a)-(d) Snapshots of within-day microtransit simulation for four cities in the U.S.

Table 5. Summary of calibration results

City	Calibr	ated Par	ameters	Opera-			Er	ror		
	Walking limit (mile)	Dwell time (sec)	Operator Weight	per mile (\$)	In- vehicle Time Error (min)	%	Wait Time Error (min)	%	Ridership Error	%
Salt Lake City	0.5	15	0.5	5.3	2.2	20.8	4.3	32.5	56	41.5
Austin	0.1	15	0.2	9.9	13.4	158.8	0.9	9.1	21	12.1
Cupertino	0.3	5	0.8	8.5	1.6	16.5	5.8	46.7	1	2.0
Sacramento	0.1	5	0.2	7.3	0.2	1.7	15.5	55.2	44	20.0
Columbus	0.1	5	0.2	8.3	7.1	93.2	8.5	136.4	2	33.3
Jersey City	0.5	5	0.5	7.9	-	-	-	-	3	1.2
AVG.										18.4



(e) Microtransit Fleet Size

Figure 9. Convergence of day-to-day adjustment for the 6 cities, with (a) in-vehicle time, (b) wait time, (c) walk time, (d) ridership, and (e) fleet size.

10 15 Simulation Day 25

20

Table 6 shows the fleet size, fleet VMT, and ridership at convergence for the 6 cities obtained using the calibrated market equilibrium models. Both fleet VMT and fleet size are not known and are only inferred from the simulation models. Note that the VMT values in Table 6 include the vehicle miles with and without passengers as they pertain to the fleet. Going from depots to pick-up points and going back to depots after drop-offs accounts for a significant proportion of the values.

Table 6. Summary of microtransit performance in 6 U.S. cities based on the equilibrated outputs of calibrated market equilibrium model

Cities	No. of days	Fleet	6-9AM fleet VMT	Microtransit (Via) ridership		
	to converge	Size		Total	% first/last mile access trips	
Salt Lake City	43	9	405	79	35%	
Cupertino	23	6	247	49	82%	
Sacramento	22	18	883	176	14%	
Columbus	10	7	298	8	100%	
Austin	15	27	1221	153	88%	
Jersey City	11	11	520	242	67%	

The model shows higher proportions of Via ridership as door-to-door trips within a service region for Salt Lake City and Sacramento, while for the other 4 cities, Via trips are predominantly first/last mile. This output is not included in the data shared but is inferred from the model. The results highlight the variable effects different operation strategies can have on microtransit ridership and consequently on other performance measures (like VMT and fleet size) in different groups of cities.

4.4 Microtransit deployment portfolio forecast model development

With the market equilibrium models calibrated, we proceed to upscale the scenario data to estimate forecast models for portfolio management. In this section, two objectives are set forth.

First, we test whether we can effectively use the fitted market equilibrium models to generate new synthetic scenarios to use to upscale the data for inferring new insights for microtransit deployment. This can be effectively demonstrated by showing that the scenario data can be related to public data to find statistically significant relationships with a good fit when compared to the original data set. This is shown in Section 4.4.1.

Second, we demonstrate the use of the forecast models in portfolio management by constructing two alternative portfolios with similar fleet VMT as the original data set from Via and illustrate how we can characterize their ridership based on the selected service region designs in each city. This is shown in Section 4.4.2.

4.4.1 Forecast model estimation and validation

Two sets of models are estimated: one for predicting ridership and one for predicting fleet VMT. For the scenario generation, only four of the six cities are used for constructing the synthetic scenarios. This is because Jersey City operates under a very different operation and Columbus is such an outlier.

For the scenario generation, only four of the six cities are used for constructing the synthetic scenarios. Jersey city operates under a very different operation from the other cities considered; the service area is divided into a central and an outer zone with different operation and pricing strategies for the inter and intra zone rides. Moreover, Columbus is such an outlier due to very low

ridership (100% of which are estimated to be first/last mile access trips as shown in Table 5) compared to the other cities. Synthetic scenarios generated from such outlier cities could mislead the training process of the forecast models and hence are not considered. In practice, having diverse scenario data from different cities would help produce more accurate and generalized forecast models, however, more deployment data from multiple cities would be needed. In future research, with more city data available one should ideally classify clusters of city types (for example auto-heavy, transit-heavy cities etc. (Oke, et al., 2019; Rath, et al., 2022) that can be fitted to different forecast models. Having such category-wise data (based on city typologies or operation policies) can be useful in extending the generalizability of the different forecast models developed for application to multiple cities for portfolio design.

The scenarios generated based on the cities considered above are assumed to cover a reasonable range of ridership and pricing policies. A set of 326 synthetic scenarios are generated, with characteristics shown in Table 7. Based on the ridership and VMT values derived from the market equilibrium of those 326 synthetic scenarios used as surrogate data, we develop microtransit portfolio forecast models using multiple linear regression with a set of features (see Table 8) and their first order interactions. As noted by Friedrich (1982), such interaction effects do not lead to multicollinearity issues. The dependent (target) variables for the two models are:

- Average peak period ridership per region's population (in thousands)
- Via's fleet VMT per region area in acres (in hundreds)

Table 7. Summary of data samples from scenario generation process used in forecast models

Number of synthetic scenarios	326
Breakdown by city	Salt Lake City: 71, Austin: 79, Sacramento: 100,
	Cupertino: 76
Breakdown of PP1/PP2	PP1: 174, PP2: 152
Breakdown of L2/L3 scenarios	LL2: 178, LL3:148
Range of number of riders	[0,2217]
Breakdown of direct trips versus first/last mile	direct: [1% - 88%]; first/last mile: [12%-99%] of total
	ridership

We consider the following independent variables (pertaining to each service region) in our models as shown in Table 8, where the feature (variable) values of a region are computed as the aggregate of all census tracts in the region

Table 8. Details of the independent variables considered in the forecast models

Independent variable	Description
Employment density	Total employed population in the region over total unprotected region area in
	acres
Household density	Total households in the region over total unprotected region area in acres (this
	is highly correlated to total population, and male/female population density
	features, hence we consider only one of these in our models).
Mean income	Household weighted mean income in the region in U.S. dollars
Street density	Total road network in the region in miles over total unprotected region area in
	acres
Transit station density	Total number of transit station in the region over total unprotected region area
	in acres
Ratio of households with one or	Sum of households with 1 or more auto ownership with respect to total
more auto	households in the region
Trip equilibrium index	mean trip productions and trip attractions equilibrium index in the region; the
	closer to one, the more balanced the trip making

PP1	if pricing policy in the region is PP1 then 1 else 0
PP2	if pricing policy in the region is PP2 then 1 else 0
Microtransit fare	value of fixed microtransit fare in the region in U.S. dollars (based on Via fare)

We fit this model using the method of least squares and apply lasso regularization for feature elimination. We use the 326 upscaled scenario data for training the models and the 4 original cities' data to regularize the estimation. The estimated coefficients of the two models are reported in the Appendix in Tables A.1 (ridership) and A.2 (fleet VMT). Due the abundance of first-order interaction terms, Lasso regularization is used to regularize and eliminate non-impactful features for a stable fit (see Tibshirani, 1996).

The evaluation of goodness-of-fit is done over the four data points (i.e., Via operated service regions) for which we have the actual Via ridership data and corresponding VMT values from the simulation. We consider the root mean squared error (RMSE) and the coefficient of variation (CV) as evaluation metrics, where CV is calculated using Eq. (18).

$$CV = \frac{\sqrt{\frac{\sum_{i=1}^{N} (y_i - y_i')^2}{N}}}{\frac{N}{\overline{Y}}} = \frac{RMSE}{\overline{Y}}$$
(18)

where y_i is the actual value and y'_i is the predicted value of the target variable for a sample i (in sample size N); \overline{Y} is the mean of the actual values of the target variable across all samples. The comparison of the ridership and VMT models with the observed values and their goodness-of-fit performances are reported in Table 9.

Table 9. Estimation results for the ridership and fleet VMT forecast model **Model estimation**

City	Ridership model (r	iders/peak period)	VMT model (veh-mi / peak period)			
	Observation	Prediction	Simulation	Prediction		
Salt Lake City	135	211	405	779		
Cupertino	50	19	247	153		
Sacramento	220	225	883	1068		
Austin	174	277	1221	1505		

Model performance

	Ridership model	VMT model
Training set R ²	72%	90%
Training set adj. R ²	67%	89%
Number of features (including intercept)	47	55
Via cities RMSE	65.92	256.67
Via cities mean	144.75	688.98
Via cities CV (%)	45.54	37.25

The estimation effort demonstrates that upscaling data from just four cities is viable, as we can fit models quite well with relatively high R^2 values. The key question is whether upscaling improves over having no upscaling at all. When the model's predictions are compared to the four data samples, the CV of \sim 37-45% is shown in Table 8. While this is not a significantly accurate forecast range for typical studies involving large data samples, we note that without upscaling,

data from only the four cities would not allow for even a forecast model to be estimated in the first place since there is insufficient data. As such, the validation based on only four observations indicates a significant improvement.

4.4.2 Inference analysis

Since the Lasso method does not output p-values for the parameters, the statistical significance of selected features is determined by re-estimating the selected features using ordinary least squares (OLS). For forecasting purposes, the coefficients from the Lasso method (as reported in Tables A.1 and A.2) are used. However, for inference analysis we refer to the p-values of the OLS models (which may vary in value slightly from the Lasso models but can provide some indication of significance) as reported in Tables A.3 (ridership) and A.4 (VMT).

The OLS models suggest the ridership and VMT are indeed dependent on employment density, household density, mean income, street density, transit station density, and car ownership, by their statistical significances at 5% levels for the parameters of the standalone features (household density, transit station density, employment density, fixed fare) or as part of the first-order interaction features (which covers the rest). In addition, the OLS models suggest a sensitivity to the pricing policy through the first-order interactions. This provides the microtransit operator and local city with a trade-off to consider when deciding which policy to implement in a city, as PP2 would increase ridership but also increase costs.

In addition, the forecast models are clearly sensitive to the service region design as that determines the input variables used. The estimated models suggest that when designing a service region, a microtransit operator and local city agency can look to the zones with attributes that would increase ridership while minimizing VMT. Variables with statistically significant positive impact on ridership and negative impact on VMT include zones with higher transit station density, higher employment density, and lower "employment density × fixed fare". This is quite interesting as it suggests that higher employment density increases ridership while decreasing VMT, but at the same time the employment density also falls within the interaction effect with fixed fare price. Since the fare policy varies from city to city instead of zone to zone in the same city in general, the first order effect suggests setting lower fares while selecting zones with a city with higher employment density. These features with statistically opposing signs are bolded in Tables A.1 – A.4.

4.4.3 Forecast model application for deployment planning

To provide a better idea of how the forecast models can be used for microtransit service portfolio design in different cities, we collect data from hundreds of 100K+ population cities in the U.S. and use the models to consider eight new cities (other than the cities considered in our study), i.e., Arlington (Texas), Birmingham (Alabama), Boston (Massachusetts), Chicago (Illinois), Detroit (Michigan), Seattle (Washington), St. Louis (Missouri), and Washington D.C. For these cities, we only have the public data gathered for the forecast models.

Assuming a constraint on total VMT of \sim 2756 veh-miles/peak period (i.e., a budget constraint around the same value as the total VMT observed for the four Via cities considered in the case study), we present two alternative portfolios for service deployment in different cities. For each of the eight cities, we generate various L2 scenarios (service regions) and get their population and built environment characteristics. We apply PP1 and PP2 pricing policies to the cities. We use the fleet VMT forecast model and select service regions from different cities such that the total forecasted VMT of a portfolio matches the budget considered. We design two alternative service

portfolios, each with the two different fare pricing policies as shown in Figure 10, use the ridership forecast model to forecast the peak period ridership for the two portfolios.

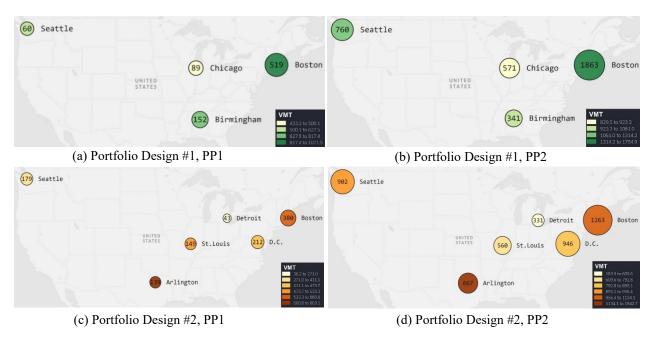


Figure 10. Two different portfolio designs under two different fare policies, where design #1 (a, b) includes 4 cities with PP1 (a) and PP2 (b); design #2 (c,d) includes 6 different cities with PP1 (c) and PP2 (d). Estimated ridership and VMT are visualized for each city, where the circle radius is peak period ridership (values in circles), and circle colors are VMT (in legend).

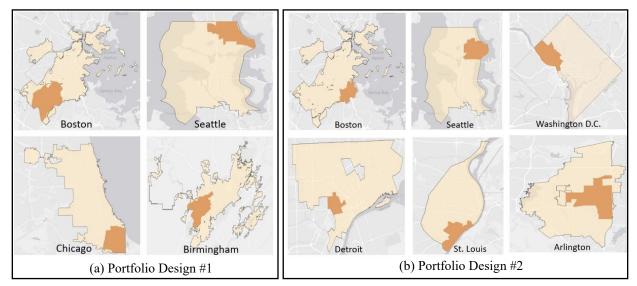


Figure 11. Service regions of portfolio designs (a) #1 and (b) #2.

For the PP1 cases (Figure 10 (a), (c)), the total forecasted ridership values in portfolios #1 and #2 are 1.4 times and 1.9 times higher, respectively, than the total ridership of the four Via cities for the same value of total VMT. In portfolio #1, Boston contributes to the highest proportion of total VMT as well as the total ridership. In portfolio #2, on the other hand, while the maximum

ridership is observed in Boston, the estimated VMT in Arlington is relatively higher compared to the other cities.

We apply the forecast models to estimate the ridership and VMT values for these same service regions in Figure 11 under a different pricing policy (i.e., PP2) as shown in Figure 10 (b, d). As shown in Figure 10(b) for cities in portfolio #1, the PP2 pricing results in increased ridership compared to PP1 across all 4 cities (maximum in Seattle) along with corresponding increment in VMT values (average VMT increment across 4 cities is 1.7 times the PP1 policy values). For each city, if we compare the ridership gain with respect to the increment in VMT under the new pricing policy, it is observed that the highest gain is estimated in Seattle with Birmingham being the lowest. Similarly, for cities in portfolio #2, using PP2 pricing policy (see Figure 10(d)) results in increased ridership and VMT in all 6 cities. Interestingly, in this case, although the maximum increment in ridership is observed in Detroit followed by Seattle, the highest gain in ridership with respect to the corresponding VMT increment is observed in St. Louis, followed by Washington D.C. and Arlington.

Although we have presented only two alternative portfolio designs assuming similar pricing policies in all cities and a VMT constraint, microtransit operators can use such forecast models for comparing across multiple service portfolios by optimizing for ridership, considering additional operation cost constraint, applying different pricing and operating policy combinations to specific cities (e.g., based on city types), etc. Granted, the accuracy of the results above depends on the quality of the forecast models M and should not be taken at face value. Rather, they illustrate that a seemingly nontrivial task of finding a portfolio from a market of hundreds of cities can be done, and with more reliability if more underlying deployment city observations are available (much as the quality of an upscaled image depends on the quality of the original image as well). To extend the generalizability of such forecast models to multiple cities in practice, with availability of detailed city deployment data for a greater number of cities, it might be relevant for operators and planners to categorize and cluster the data as per different city types such as auto-heavy, transitheavy cities etc. (Oke et al., 2019, Rath et al., 2022). This can be used to develop category-wise forecast models (using diverse set of scenarios generated per category) based on the proposed framework to better estimate service performances in different cities for portfolio design. Moreover, the input variables considered in these models can be designed based on the purpose of the mobility services and the category considered to achieve useful insights.

Hence, this can be used as an effective decision-support tool for microtransit service deployment planning for strategizing resource-allocation and investment decisions, one that can help inform public agencies by providing quantitative results that can spur further local studies. Example application include comparing performances of free first/last mile microtransit ride policy to auto-heavy cities for improving transit ridership, evaluating equity benefits for different population groups in terms of improved accessibility in candidate service regions in one set of city deployments vs. another, and analyzing benefits and tradeoffs of alternative portfolio designs with different operation strategies to different city types in the portfolio within a given budget. This work is also of great importance to federal agencies like the Federal Transit Administration (FTA) in helping to identify priority areas for funding microtransit projects in different cities. With the increasing importance of using innovative analysis tools and interactive platforms by private companies (Mercer & Hewitt, 2021; Jacobs 2022) and public agencies (TransitCenter, 2021) to aid decision-making, the proposed method can be used to design a service portfolio dashboard and create an interface which private companies, public agencies and planners can use for city-level deployment planning of emerging mobility services.

4.4.4 Policy implications from empirical insights

The results in Sections 4.4.1 to 4.4.3 reveal multiple original empirical insights that would serve both policymakers at state/federal agencies and microtransit/emerging transportation technology providers. The inference analysis finds that microtransit ridership (a proxy for benefits) and fleet VMT (a proxy for costs) are dependent on the following factors (aligned side-by-side) shown in Table 10. Common factors (highlighted in light green) between the two models are the transit station density and the employment density. These suggest that the most efficient deployments that maximize ridership while minimizing fleet VMT would have high transit station density and high employment density. HH density could have mixed effects as increasing ridership would also come at increasing fleet VMT.

Table 10. Comparison of effects of different variables on ridership and fleet VMT

Ridership model		Fleet VMT model		
Significant variables Sign		Significant variables	Sign	
Mean income	-			
HH density	+	HH density	+	
Transit station density	+	Transit station density	-	
Employment density	+	Employment density	-	
		Fixed fare	+	

These insights allow policymakers to sort a portfolio of cities and the zones within each to identify effective service region designs. The example portfolios shown in Section 4.4.3 highlight an effort that would be impossible for a policymaker to justify using existing methods and frameworks from the literature. Based on existing practices, a federal agency like FTA would collect data from hundreds of cities first (which we also did), but beyond that they would not have any models to quantify the performance of each city, let alone having different service regions in each city. But because of the availability of the forecast models estimated from the upscaled data, we can easily identify alternative portfolio options and forecast their collective performance.

5. Conclusion

Transportation technologies are not "one-size-fits-all" solutions; this point is clearly demonstrated by the 67%/23%/50% failure rates of demand-responsive transport services implemented in UK/Europe/North America. Emerging technologies like microtransit and state/federal agencies need to have effective decision-support tools, which are limited by the complexity of the decisions that need to be made, the limited availability due to the "emerging technology" aspect, and due to the myriad of operations that expand the dimensionality of the problem further. For example, even a success story like Via only operates in less than 40 U.S. cities while there are over 3000 U.S. cities with populations of 10,000 or more.

We propose a methodology to upscale data from the limited data available to microtransit operators (and to public agencies like the Federal Transit Administration in overseeing deployment regulations at the federal level). The method uses simulation to fit market equilibrium models to a small set of deployment cities R so that those models m_s , $s \in R$, can be used to generate scenario data y_i , $i \in Q$, $Q \gg R$, at low cost. The resulting data can be used to fit forecast models $y_i = M(x_i; \theta)$, $i \in Q$ so that they can be applied to the population of cities $P, P \gg R$. The overall

framework contributes to the literature by parameterizing the within-day simulator from Yoon et al. (2022), extending the day-to-day market equilibrium model from Djavadian and Chow (2017a,b) to consider travelers with first/last mile access trips as well as direct trips, and developing a scenario generation algorithm for feeding the market equilibrium simulation model. The models are shown to fit the six cities with an average ridership error of 18.4% while outputting latent attributes like fleet size, fleet VMT, % of riders by first/last mile, and breakdown of journey times to their components. Note that a conventional transportation planning study would entail estimating a travel demand model for one city, obtain similar accuracy (Flyvbjerg et al., 2006), and may not be able to output metrics like % of riders by first/last mile.

The new upscaled scenario data proves to be useful; models fit to the data are adequately accurate compared to the original limited city deployment data set (CVs ~ 37-45% for only four observations, which is very statistically efficient due to the upscaling) whereas the original four observations would be insufficient to produce any meaningful model at all using any existing method or framework from the literature. In that sense, the framework successfully upscaled the limited samples to produce a synthetic data set of 326 synthetic scenarios. Furthermore, the forecast models identify meaningful relationships between ridership and fleet VMT with a host of independent public data (employment, households, car ownership, transit station density, income, street density) and microtransit operating policies (pricing, service region). Example variables with statistically significant positive impact on ridership and negative impact on VMT include zones with higher transit station density and higher employment density, not counting the first order interactions, as shown in Table 10. Application of the models illustrate how they can quantify the effectiveness of a given portfolio and quickly compare between different portfolio designs. We can nontrivially identify two alternative portfolios of city service regions with similar VMT as Via's four cities but having up to 1.9 times the ridership, justifying them from a market of hundreds of cities with 100K+ populations.

Future research should look at further collaboration with microtransit providers to classify cities into clusters (e.g., auto, bus transit, congested, hybrid, metro bike and mass transit dominant cities as per Oke et al. (2019) and to focus more on empirically capturing good fitting forecast models using this new methodology. This could include obtaining microtransit operation data for additional cities including disaggregate data broken down by categories such as service region type, city type or operation type e.g., first/last mile and door-to-door; such data can be used to build category-wise forecast models that can be effective for evaluating portfolio designs under different operating policies and service region designs. Other emerging technologies should also be considered, especially where data are limited: e.g., planning electric vehicle fleets and charging infrastructure, pilots for autonomous vehicle fleets. A portfolio dashboard can be implemented to help a microtransit provider or the FTA evaluate their portfolios and analyze alternative portfolio designs.

Acknowledgments

This research was conducted with the support of C2SMART University Transportation Center (USDOT award #69A3551747124). Data shared by *Via Transportation* is gratefully acknowledged.

Author Statement

The authors confirm contribution to the paper as follows: study conception and design: SR, BL, GY, JYJC; data collection: SR, BL, GY; analysis and interpretation of results: SR, BL, GY, JYJC; draft manuscript preparation: SR, BL, GY, JYJC.

References

- American Community Survey, 2019. Subject tables, United States Census Bureau. Available at: https://www.census.gov/acs/www/data/data-tables-and-tools/subject-tables/, Accessed 29 July 2021.
- Anas, A. (1981). The estimation of multinomial logit models of joint location and travel mode choice from aggregated data. *Journal of regional science*, 21(2), 223-242.
- Bardaka, E., Hajibabai, L., & Singh, M. P. (2020). Reimagining ride sharing: Efficient, equitable, sustainable public microtransit. *IEEE Internet Computing*, 24(5), 38-44.
- Bliss, L., 2017. Bridj Is Dead, but Microtransit Isn't, Bloomberg Citylab.
- Cantarella, G. & Cascetta, E., 1995. Dynamic processes and equilibrium in transportation networks: towards a unifying theory. *Transportation Science* 29(4), 305-329.
- Caros, N. & Chow, J. Y. J., 2021. Day-to-day market evaluation of modular autonomous vehicle fleet operations with en-route transfers. *Transportmetrica B* 9(1), 109-133.
- Chow, J. Y. J. (2018). *Informed Urban transport systems: Classic and emerging mobility methods toward smart cities*. Elsevier.
- Chow, J. Y. J. & Djavadian, S., 2015. Activity-based market equilibrium for capacitated multimodal transport systems. *Transportation Research Part C* 59, 2-18.
- Chow, J. Y. J., Rath, S., Yoon, G., Scalise, P., Alanis Saenz, S., 2020. Spectrum of Public Transit Operations: From Fixed Route to Microtransit. FTA Report NY-2019-069-01-00. https://c2smart.engineering.nyu.edu/wp-content/uploads/2020/04/Chow-FTA-Report-NY-2019-069-01-00.pdf.
- Chow, J. Y. J., Regan, A., Ranaiefar, F. & Arkhipov, D., 2011. A network option portfolio management framework for adaptive transportation planning. *Transportation Research Part A* 45(8), 765-778.
- Cich, G., Knapen, L., Maciejewski, M., Bellemans, T., & Janssens, D., 2017. Modeling demand responsive transport using SARL and MATSim. *Procedia Computer Science* 109, 1074-1079.
- Cooper, R., Edgett, S. & Kleinschmidt, E., 1998. *Portfolio management for new products*. Addison Wesley Lgonman, Inc., Reading MA.
- CTPP, 2016. Census Data for Transportation Planning Applications, AASHTO. Available at: https://ctpp.transportation.org/, Accessed 29 July 2021.
- Currie, G. & Fournier, N., 2020. Why most DRT/Micro-Transits fail—What the survivors tell us about progress. *Research in Transportation Economics* 83, 100895.
- Daganzo, C. & Ouyang, Y., 2019. A general model of demand-responsive transportation services: From taxi to ridesharing to dial-a-ride. *Transportation Research Part B* 126, 213-224.
- Djavadian, S. & Chow, J. Y. J., 2017a. Agent-based day-to-day adjustment process to evaluate dynamic flexible transport service policies. *Transportmetrica B* 5(3), 281-306.
- Djavadian, S. & Chow, J. Y. J., 2017b. An agent-based day-to-day adjustment process for modeling 'Mobility as a Service' with a two-sided flexible transport market. *Transportation Research Part B* 104, 36-57.
- Dong, X., Chow, J. Y., Waller, S. T., & Rey, D. (2022). A chance-constrained dial-a-ride problem with utility-maximising demand and multiple pricing structures. *Transportation Research Part E: Logistics and Transportation Review*, 158, 102601.

- EPA, 2021. Smart location database. Available at: https://www.epa.gov/smartgrowth/smart-location-mapping#SLD, Accessed 30 July 2021.
- Flyvbjerg, B., Skamris Holm, M. K., & Buhl, S. L. (2006). Inaccuracy in traffic forecasts. *Transport Reviews*, 26(1), 1-24.
- Friedrich, R. J. (1982). In defense of multiplicative terms in multiple regression equations. *American Journal of Political Science*, 797-833.
- Fu, Z., & Chow, J. Y. (2022). The pickup and delivery problem with synchronized en-route transfers for microtransit planning. *Transportation Research Part E: Logistics and Transportation Review*, 157, 102562.
- Haglund, N., Mladenović, M. N., Kujala, R., Weckström, C., & Saramäki, J., 2019. Where did Kutsuplus drive us? Ex post evaluation of on-demand micro-transit pilot in the Helsinki capital region. *Research in Transportation Business & Management*, 32, 100390.
- Horn, M., 2002. Fleet scheduling and dispatching for demand-responsive passenger services. Transportation Research Part C: Emerging Technologies 10(1), 35-63.
- Horowitz, J., 1984. The stability of stochastic equilibrium in a two-link transportation network. *Transportation Research Part B* 18(1), 13-28.
- Jacobs, 2022. Jacobs Acquires Mobility Analytics Leader StreetLight Data, Inc. Available at: www.prnewswire.com/news-releases/jacobs-acquires-mobility-analytics-leader-streetlight-data-inc-301476275, Accessed 1 October, 2022.
- Jung, J. & Chow, J. Y. J., 2019. Effects of charging infrastructure and non-electric taxi competition on electric taxi adoption incentives in New York City. *Transportation Research Record* 2673(4), 262-274.
- Krok, A., 2016. Car2Stop: Car2Go shuts down services in San Diego, CNET.
- MacQueen, J., 1967. Some methods for classification and analysis of multivariate observations. *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, 1(14), 281-297.
- Mahmassani, H. S., 1990. Dynamic models of commuter behavior: Experimental investigation and application to the analysis of planned traffic disruptions. *Transportation Research Part A* 24(6), 465-484.
- Mahmassani, H. S. & Chang, G., 1986. Experiments with departure time choice dynamics of urban commuters. *Transportation Research Part B* 20(4), 297-320.
- Markov, I., Guglielmetti, R., Laumanns, M., Fernández-Antolín, A. and de Souza, R., 2021. Simulation-based design and analysis of on-demand mobility services. *Transportation Research Part A* 149, 170-205.
- Marshall, A., 2019. Ford Axes Its Chariot Shuttles, Proves Mobility Is Hard., Wired.
- Ma, T., Chow, J. Y. J., Klein, S. & Ma, Z., 2021. A user-operator assignment game with heterogeneous user groups for empirical evaluation of a microtransit service in Luxembourg. *Transportmetrica A* 17(4), 946-973.
- Ma, T., Rasulkhani, S., Chow, J. Y. J. & Klein, S., 2019. A Dynamic Ridesharing Dispatch and Idle Vehicle Repositioning Strategy with Integrated Transit Transfers. *Transportation Research Part E* 128, 417–442.
- Mercer, T., & Hewitt, C., 2021. Introducing Remix by Via's On-demand Planning. Available at: www.remix.com/blog/introducing-our-newest-product-on-demand-planning, Accessed 1 October, 2022.

- Moovit, 2021. On-demand Microtransit Glossary. Available at: https://moovit.com/wp-content/uploads/2021/08/01142135/Moovit-_-On-Demand-Microtransit-Glossary.pdf, Accessed 21 September, 2022.
- Mercer, T., & Hewitt, C., 2021. Introducing Remix by Via's on-demand planning. Available at: www.remix.com/blog/introducing-our-newest-product-on-demand-planning, Accessed 1 October, 2022.
- Nalic, D., Mihalj, T., Bäumler, M., Lehmann, M., Eichberger, A. and Bernsteiner, 2020. Scenario Based Testing of Automated Driving Systems: A Literature Survey. *Proc. FISITA Web Congr.*, 30.
- Oke, J., Aboutaleb, Y.M., Akkinepally, A., Azevedo, C.L., Han, Y., Zegras, P.C., Ferreira, J. and Ben-Akiva, M.E., 2019. A novel global urban typology framework for sustainable mobility futures. *Environmental Research Letters* 14(9), 095006.
- OSM, 2021. Python OSMnx library (Open Street Map). Available at: https://github.com/gboeing/osmnx, Accessed 30 July 2021.
- OTP, 2021. Open Trip Planner API. Available at: https://www.opentripplanner.org/, Accessed 30 July 2021.
- Pantelidis, T. P., Chow, J. Y. J. & Rasulkhani, S., 2020. A many-to-many assignment game and stable outcome algorithm to evaluate collaborative mobility-as-a-service platforms. *Transportation Research Part B* 140, 79-100.
- Pinto, H., Hyland, M., Mahmassani, H. S. & Verbas, I., 2020. Joint design of multimodal transit networks and shared autonomous mobility fleets. *Transportation Research Part C* 113, 2-20.
- Rath, S. & Chow, J., 2022. Worldwide city transport typology prediction with sentence-BERT based supervised learning via Wikipedia. *Transportation Research Part C: Emerging Technologies*, Volume 139, p. 103661.
- Rath, S., Liu, B., Yoon, G., Chow, J. Y. J. (2021). Urban microtransit cross-sectional study for service portfolio design [supporting dataset and code], https://zenodo.org/record/5517983#.YVxrQ9rMKUl, last accessed Oct. 5, 2021.
- Redmond, L. S., & Mokhtarian, P. L. (2001). The positive utility of the commute: modeling ideal commute time and relative desired commute amount. *Transportation*, 28(2), 179-205.
- Rocklage, E., Kraft, H., Karatas, A. & Seewig, J., 2017. Automated scenario generation for regression testing of autonomous vehicles. *IEEE 20th International conference on intelligent transportation systems*, 476-483.
- Shaheen, S. & Chan, N., 2016. Mobility and the sharing economy: Potential to facilitate the first-and last-mile public transit connections. *Built Environment* 42(4), 573-588.
- Shotl, 2021. Available at: https://shotl.com/, Accessed 7 October 2021.
- Smith, M. et al., 2014. The long term behaviour of day-to-day traffic assignment models. *Transportmetrica A* 10(7), pp. 647-660.
- Smith, M. J., 1984. The stability of a dynamic model of traffic assignment—an application of a method of Lyapunov. *Transportation Science* 18(3), 245-252.
- Statista, 2019. Number of cities, towns and villages (incorporated places) in the United States in 2019, by population size. Available at: https://www.statista.com/statistics/241695/number-of-us-cities-towns-villages-by-population-size/, Accessed 29 July 2021.
- Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1), pp.267-288.
- Train, K., & McFadden, D. (1978). The goods/leisure tradeoff and disaggregate work trip mode choice models. *Transportation research*, 12(5), 349-353.

- Transdev, 2021. Available at: https://transdevna.com/services-and-modes/microtransit/, Accessed 29 July 2021.
- TransitCenter, 2021. Introducing the Transit Equity Dashboard. Available at: https://transitcenter.org/introducing-the-transit-equity-dashboard/, Accessed 1 October, 2022.
- Transitfeeds, 2021. Available at: https://transitfeeds.com/, Accessed 30 July 2021.
- Tuncali, C., Fainekos, G., Ito, H. & Kapinski, J., 2018. Simulation-based adversarial test generation for autonomous vehicles with machine learning components. *IEEE Intelligent Vehicles Symposium (IV)*, 1555-1562.
- Via, 2021. Available at: https://ridewithvia.com/, Accessed 29 July 2021.
- Volinski, J., 2019. Microtransit or General Public Demand–Response Transit Services: State of the Practice. *TCRP Synthesis of Transit Practice Project J-7*, Volume Topic SB-30.
- Wardman, M., 2004. Public transport values of time. Transport policy, 11(4), 363-377.
- Watling, D. & Hazelton, M., 2003. The dynamics and equilibria of day-to-day assignment models. *Networks and Spatial Economics* 3(3), 349-370.
- Yan, X., Levine, J. & Zhao, X., 2019. Integrating ridesourcing services with public transit: An evaluation of traveler responses combining revealed and stated preference data. *Transportation Research Part C* 105, 683-696.
- Yoon, G., Chow, J. Y. J. & Rath, S., 2022. A simulation sandbox to compare fixed-route, flexible-route transit, and on-demand microtransit system designs, *KSCE J Civ Eng* 26, 3043–3062.
- Zhou, Y., Liu, X. C., & Grubesic, T. (2021). Unravel the impact of COVID-19 on the spatio-temporal mobility patterns of microtransit. *Journal of Transport Geography*, 97, 103226.

Appendix

The estimated feature coefficient values of the ridership and VMT forecast models using Lasso method are listed in Tables A.1 and A.2, respectively. The same parameters are estimated using OLS to report the p-values in Tables A.3 and A.4, respectively. Bolded features indicate statistically significant opposing signs between the ridership model and VMT model which can help identify best locations for expanding microtransit service.

Table A.1. Ridership forecast model estimated feature coefficient values

Table A.1. Ridership Torecast model estimated leature coefficient values						
Feature	Estimated coefficient	Feature	Estimated coefficient	Feature	Estimated coefficient	
Intercept	2.21E-01	mean income (\$) × fixed fare	3.07E-06	HH density × PP1	1.19E-01	
mean income (\$)	-2.90E-06	auto ownership per HH × street density	-1.69E-01	HH density × PP2	1.70E-01	
auto ownership per HH	3.75E+00	auto ownership per HH × HH density	-4.05E-04	HH density × fixed fare	-1.39E-01	
street density	5.82E-02	auto ownership per HH × transit station density	-2.04E-02	transit station density × mean TRIPEQ	1.27E+02	
HH density	8.83E-02	auto ownership per HH × employment density	4.17E-02	transit station density × PP1	-2.76E+00	
transit station density	3.30E+01	auto ownership per HH × mean TRIPEQ	2.86E-01	transit station density × PP2	9.87E+01	
employment density	1.00E-01	auto ownership per HH × PP1	3.63E-01	transit station density × fixed fare	-2.04E+01	
mean TRIPEQ	2.29E+00	auto ownership per HH × fixed fare	-3.87E-01	employment density × mean TRIPEQ	6.09E-01	
PP1	-1.14E+00	street density × HH density	-1.63E-02	employment density × PP1	-4.66E-02	
PP2	2.34E-15	street density × transit station density	-1.56E-01	employment density × PP2	4.41E-01	
fixed fare	-6.25E-02	street density × mean TRIPEQ	8.86E-02	employment density × fixed fare	-7.92E-02	
mean income (\$) × street density	1.37E-07	street density × PP1	5.78E-02	mean TRIPEQ × PP1	-4.08E-01	
mean income ($\$$) × transit station density	1.80E-04	street density × PP2	2.22E-03	mean TRIPEQ × PP2	1.05E-01	
mean income (\$) × employment density	3.88E-07	street density × fixed fare	-2.70E-02	mean TRIPEQ × fixed fare	-4.12E-02	
mean income (\$) × PP1	-1.50E-05	HH density × employment density	4.39E-03	PP1 × fixed fare	-2.67E-02	
mean income (\$) × PP2	7.53E-07	HH density × mean TRIPEQ	6.54E-01			

Table A.2. VMT forecast model estimated feature coefficient values						
Feature	Estimated coefficient	Feature	Estimated coefficient	Feature	Estimated coefficient	
Intercept	-3.26E+00	mean income (\$) × fixed fare	1.78E-06	HH density × mean TRIPEQ	-1.10E+01	
mean income (\$)	2.86E-05	auto ownership per HH × street density	-7.74E-01	HH density × PP1	4.03E+00	
auto ownership per HH	-1.62E+00	auto ownership per HH × HH density	-9.01E-01	HH density × PP2	9.38E-01	
street density	2.71E-01	auto ownership per HH × transit station density	-3.04E+02	HH density × fixed fare	-5.65E-01	
HH density	8.03E-01	auto ownership per HH × employment density	9.06E-01	transit station density × employment density	-5.71E+01	
transit station density	-1.36E+01	auto ownership per HH × mean TRIPEQ	5.56E+01	transit station density × mean TRIPEQ	1.29E+01	
employment density	-6.39E-01	auto ownership per HH × PP1	4.48E+00	transit station density × PP1	3.97E+02	
mean TRIPEQ	-3.25E+01	auto ownership per HH × PP2	-5.53E-01	transit station density × PP2	4.44E+02	
PP1	-4.43E-01	auto ownership per HH × fixed fare	-4.92E+00	transit station density × fixed fare	-6.47E+01	
PP2	2.57E-15	street density × HH density	-7.76E-02	employment density × mean TRIPEQ	1.51E+00	
fixed fare	3.42E+00	street density × transit station density	-5.68E+00	employment density × PP1	-7.87E-01	
mean income (\$) × auto ownership per HH	-2.15E-05	street density × employment density	8.80E-03	employment density × PP2	1.82E+00	
mean income (\$) × street density	5.17E-07	street density × mean TRIPEQ	1.67E+00	employment density × fixed fare	2.50E-01	
mean income ($\$$) × HH density	9.61E-06	street density × PP1	2.23E-01	mean TRIPEQ × PP1	-6.83E+00	
mean income (\$) × transit station density	7.90E-04	street density × PP2	4.61E-01	mean TRIPEQ × PP2	-2.44E+00	
mean income (\$) × employment density	-6.38E-06	street density × fixed fare	-5.27E-02	mean TRIPEQ × fixed fare	-5.29E-01	
mean income (\$) × mean TRIPEQ	-1.08E-04	HH density × transit station density	1.16E+02	PP1 × fixed fare	8.88E-01	
mean income ($\$$) × PP1	1.11E-05	HH density × employment density	1.94E-01	PP2 × fixed fare	1.41E+00	
mean income (\$) \times PP2	1.72E-05					

Table A.3. Ridership forecast model selected features p-values from linear least squares estimation

Feature	p value	Feature	p value	Feature	p value
Intercept	<0.001***	mean income (\$) × fixed fare	0.001***	HH density × PP1	0.015*
mean income (\$)	<0.001***	auto ownership per HH × street density	0.786	HH density × PP2	0.051*
auto ownership per HH	0.959	auto ownership per HH × HH density	0.053*	HH density × fixed fare	0.046*
street density	0.352	auto ownership per HH × transit station density	0.074	transit station density × mean TRIPEQ	0.005**
HH density	0.027*	auto ownership per HH × employment density	0.031*	transit station density × PP1	0.007**
transit station density	0.002**	auto ownership per HH × mean TRIPEQ	0.576	transit station density × PP2	<0.001***
employment density	0.001***	auto ownership per HH × PP1	0.333	transit station density × fixed fare	0.126
mean TRIPEQ	0.751	auto ownership per HH × fixed fare	0.867	employment density × mean TRIPEQ	0.933
PP1	0.877	street density × HH density	0.730	employment density × PP1	<0.001***
PP2	0.877	street density × transit station density	0.956	employment density × PP2	0.003**
fixed fare	0.699	street density × mean TRIPEQ	0.928	employment density × fixed fare	<0.001***
mean income (\$) × street density	0.228	street density × PP1	0.344	mean TRIPEQ × PP1	0.017*
mean income (\$) × transit station density	0.489	street density × PP2	0.369	mean TRIPEQ × PP2	0.053*
mean income (\$) × employment density	0.421	street density × fixed fare	0.001***	mean TRIPEQ × fixed fare	0.634
mean income (\$) × PP1	<0.001***	HH density × employment density	0.943	PP1 × fixed fare	0.048*
mean income (\$) × PP2	0.235	HH density × mean TRIPEQ	0.655		

^{*, **, ***} refer to p-values less than 0.05, 0.01, and 0.001 respectively

Table A.4. VMT forecast model selected features p-values from linear least squares estimation

Feature	p value	Feature	p value	Feature	p value
Intercept	<0.001***	mean income (\$) × fixed fare	0.329	HH density × mean TRIPEQ	0.746
mean income (\$)	0.096	auto ownership per HH × street density	0.279	HH density × PP1	0.002**
auto ownership per HH	0.549	auto ownership per HH × HH density	0.026*	HH density × PP2	0.018*
street density	0.775	auto ownership per HH × transit station density	0.301	HH density × fixed fare	0.014**
HH density	0.006**	auto ownership per HH × employment density	0.077	transit station density × employment density	0.001***
transit station density	0.041*	auto ownership per HH × mean TRIPEQ	0.168	transit station density × mean TRIPEQ	0.082
employment density	0.006**	auto ownership per HH × PP1	0.044*	transit station density × PP1	0.058
mean TRIPEQ	0.173	auto ownership per HH × PP2	0.057	transit station density × PP2	0.027*
PP1	0.806	auto ownership per HH × fixed fare	0.012**	transit station density × fixed fare	0.042*
PP2	0.806	street density × HH density	0.029*	employment density × mean TRIPEQ	0.804
fixed fare	0.025*	street density × transit station density	0.539	employment density × PP1	0.001***
mean income (\$) × auto ownership per HH	0.233	street density × employment density	0.264	employment density × PP2	0.030*
mean income (\$) × street density	0.370	street density × mean TRIPEQ	0.266	employment density × fixed fare	0.002**
mean income (\$) × HH density	0.188	street density × PP1	0.607	mean TRIPEQ × PP1	0.022*
mean income (\$) × transit station density	0.258	street density × PP2	0.983	mean TRIPEQ × PP2	0.731
mean income (\$) × employment density	0.323	street density × fixed fare	0.117	mean TRIPEQ × fixed fare	0.221
mean income (\$) × mean TRIPEQ	0.041*	HH density × transit station density	0.001***	PP1 × fixed fare	0.197
mean income (\$) × PP1	0.206	HH density × employment density	0.109	PP2 × fixed fare	0.003**
mean income (\$) × PP2	0.048*				

^{*, **, ***} refer to p-values less than 0.05, 0.01, and 0.001 respectively