

Supplementary Appendix

S1 Data Construction

S1.1 Cellphone location records

This subsection details how we construct our sample of trips based on the raw cellphone data. The raw data is composed of a sequence of pings. Each ping contains a timestamp, latitude, longitude, and a device identifier. The final output from this process is a dataset with a fraction of the universe of trips that took place in Chicago. A sequence of filtering steps leaves us with 5% of devices. We verify that the owners of these devices are representative and then scale up the number of trips by a factor such that the aggregate number of car trips is consistent with what is reported by the Chicago Metropolitan Agency for Planning (CMAP) 2019 Household Travel Survey.¹

Data filtering We start by subsetting cellphone pings to a rectangle around the city of Chicago (i.e., latitude between 41.11512 and 42.494693, longitude between -88.706994 and -87.527174) for the month of January 2020.

Next, using the cellphone device identifier, the timestamp and geolocation of each ping, we calculate the time between two consecutive pings as well as the geodesic distance. These distances allow us to obtain the speed between consecutive pings. We then filter out “noisy” pings by using distance, time, and speed variables. In particular, we remove pings that are moving at an excessive speed since these pings are likely to be GPS “jumps” resulting from noise in the measurement of the GPS coordinates of the device.² We also drop “isolated” pings since they are not helpful for identifying whether people are moving. Additionally, we only keep pings belonging to a “stream” of pings.³ We define a stream of pings as a

¹ Source: [My Daily Travel survey \(website\)](#)

² 40 meters per second, i.e. about 145 kilometers per hour

³ In particular, we only keep pings that satisfy the following two conditions: (i) no more than ten

sequence of pings for the same cellphone identifier such that a ping always has another ping within the next 15 minutes and within 1,000 meters. We drop streams with less than 3 pings. Finally, we aggregate pings to the minute of the day by taking the average location and timestamp across pings within each minute for a given cellphone identifier. In what follows, we focus on the remaining filtered pings aggregated at the minute level.

Defining movements, stays, and trips We identify two consecutive (aggregated) pings as a “movement” for a given cellphone identifier if their distance is at least 50 meters or if their implied speed is at least 3 meters per second (6.7 miles per hour or 10.8 kilometers per hour). We then define a “stay” as a sequence of two or more successive pings with no movement.

Finally, we take all streams of pings and define trips as being a stream (i) with movement, (ii) that starts with a stay, and (iii) that ends with a stay. We remove all trips with a total geodesic trip distance between the starting and ending point below 0.25 miles (about 400 meters).

Estimation of home locations and traveler’s income This subsection details how we assign a home location and an income level to each individual cellphone identifier.

We start by assigning all cellphone pings to census blocks for the subset of pings within Chicago during our sample period. Next, we focus on pings during night hours, defined as between 10pm and 8am, when individuals are more likely to be at home.

Using this subset of pings, we attribute a score system for each hour between 10pm and 8am. Specifically, regardless of the number of pings, scores are assigned as follows:

- A value of 10 to all census blocks that were pinged between 1 am and 5 am.

minutes to either the next or the previous ping, (ii) no more than 5,000 meters to either the next or the previous ping.

- A value of 5 to all census blocks that were pinged between 11 pm and 1am or between 5 am and 7 am.
- A value of 2 to all census blocks that were pinged between 10pm and 11pm, or between 7am and 8am.

The basic idea is to assign a higher score to blocks where the cellphone owner is more likely to be at home. Finally, we sum the scores across all census blocks for each cellphone ID - month combination and keep the census block with highest score. If this highest-score census block appears on at least 3 or more separate nights during the month, we assign it as the cellphone's home census block for that month. Otherwise, we consider the cellphone as having an unknown home location, which we believe captures occasional Chicago visitors such as tourists. Throughout the text, we refer to these devices as *visitors*. Figure S8 plots the share of visitors by origin locations. We see that, for trips done by visitors, the most common origin locations are the city center (center right), both airports (top left and center left), as well as Hyde Park the neighborhood home to the University of Chicago (right, south of the center).

For all cellphones with an assigned home location, we impute their income by using the census tract median household income.⁴ Cellphones without an assigned home location (visitors) are not assigned an income at this stage.

Next, for each market, we estimate travelers' income distribution.⁵ First, we take median income by tracts and divide tracts according to Chicago-level income quintiles.⁶ Next, we assign an income quintile to each device according to their home location. Since we can follow how devices travel across space and over time, for each market, we can measure the quintile from each traveler departing from its destination. We end up, for each market, with the share of travelers in each of the five income quintiles, plus a share of visitors. For each market, we then reassign visitors proportionally across the five income quintiles so that their shares sum to

⁴ We compute the census-tract median income percentile using the 2010 Census data.

⁵ Recall, a market is defined as an (origin CA, destination CA, hour of the week)-tuple.

⁶ For 2010, income quintiles are defined using the following cut-offs: \$34,875, \$46,261, \$60,590 and \$85,762.

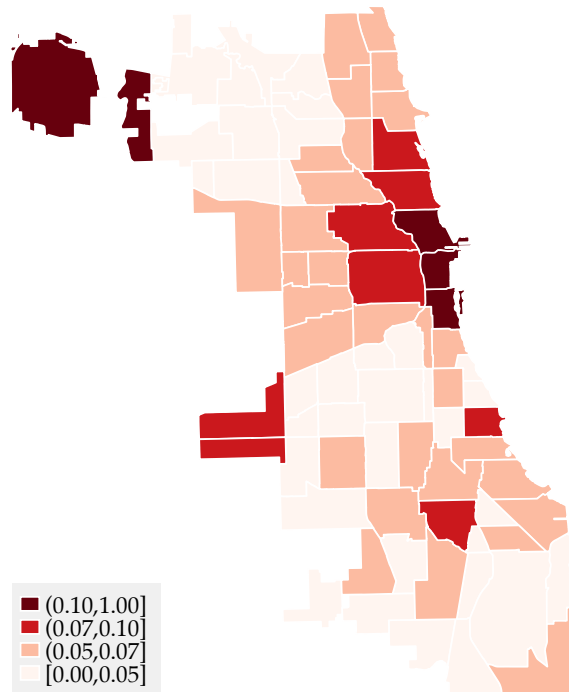


Figure S8: Share of visitors by origin location

Notes: This figure shows the share of trips at the origin CA level made by visitors. In our cellphone trips data, each market (origin-destination-hour triple) has a share of trips made by visitors. To construct the shares displayed in the figure, we take the weighted average of the share of trips made by visitors across destinations and hours of the week, for each origin CA, using inside market size (number of cellphone trips per market) as weight.

one. As a result, in the estimation, we work with five traveler types, corresponding to five income quintiles. For markets with less than 5 trips, we impute market-level income shares using the underlying distribution of census tract-level income for the origin CA of that market.

S1.1.1 Survey Data Sparsity

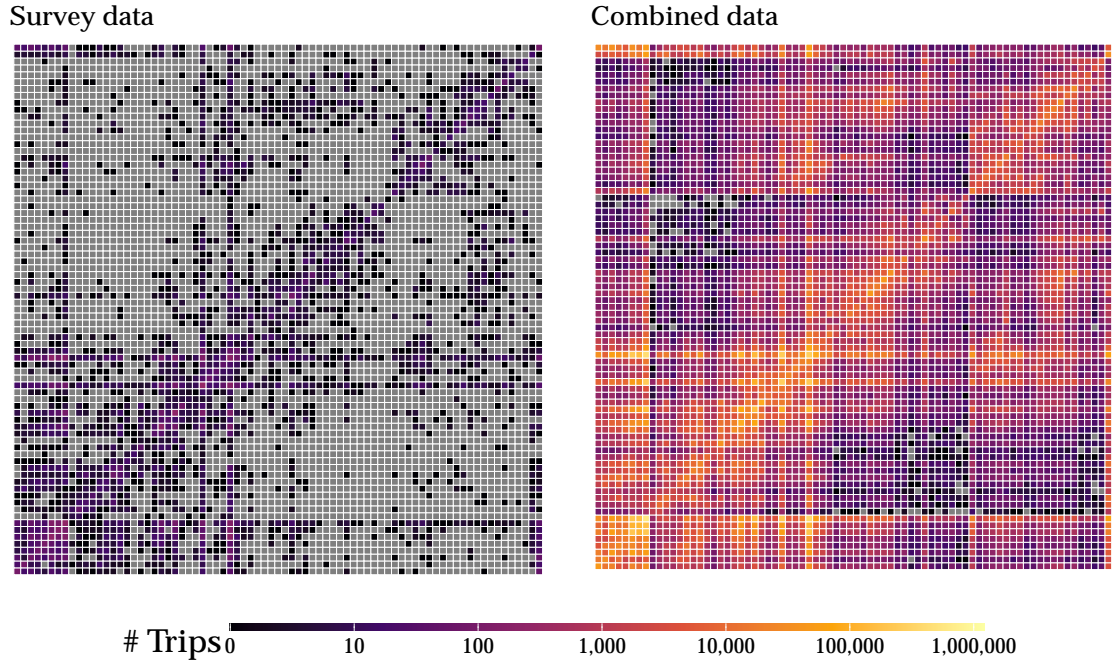


Figure S9: Combined vs. Survey Data: Flows Across Community Areas

Notes: These figures show the number of trips from every origin CA to every destination CA in our combined data (right panel) and in the survey data (left panel). Each row represents an origin CA and each column represents a destination CA. Grey points represent empty cells.

S1.2 Travel times, routes, and schedules

Travel times and routes Similar to Akbar et al. (2023), we query and geocode trips using Google Maps. For each mode of transportation, we query 30,796,848 counterfactual trips and obtain their distance, duration, and route.⁷ Importantly,

⁷ One trip for each (origin census tract, destination census tract, hour of day, weekend dummy) combination. We use all the 801 Chicago census tracts boundaries for the year 2010 from the

we can measure trip duration for the same origin-destination tuple over the time of the weekday (or weekend) and how this varies with traffic conditions. Moreover, using the detailed “steps” of the public transit Google Maps queries, we obtain walk times from the origin latitude/longitude to the “best” train or bus station.⁸

We also obtain Google Maps data on train trip times by querying Google Maps three times for each pair of train stations in Chicago. These times represented three broad time categories: weekday peak, weekday non-peak, and weekend. In particular, the first query requested a trip time of 8am on Wednesday July 6th, 2022, the second query requested a trip time of 11am on Wednesday July 6th, 2022, and the third query requested a trip time of 11am on Saturday July 9th 2022.

Public transit schedules We obtain historical GTFS data from [Open Mobility Data](#). These data contain bus and train schedules for September 2019 through February 2020.

S1.3 Constructing Mode-Specific Trips

Mode-specific trips are constructed using five main sources: (1) Taxi and TNP trips data from the City of Chicago, (2) Google Maps data, (3) cellphone trips data, (4) historical GTFS data containing public transit route schedules, and (5) Chicago public transit data from the MIT Transit Lab and the CTA.

Taxi and Transportation Network Provider (TNP) data We obtain trip times, distances, and origin-destination census tracts for both Taxi and Transportation Network Provider (TNP) trips from the [City of Chicago’s Data Portal](#).⁹

[Chicago Data City portal website](#).

⁸ The “best” bus or train station is not necessarily the closest one, depending on the destination and/or the time of the day.

⁹ For privacy reasons, during periods of the day and for locations with very few trips, only the origin and/or destination CA of a trip is reported. See [this page](#) for a discussion of the approach to privacy in this data set.

Cellphone trips data We construct cellphone trips from cellphone pings using the procedure detailed in Appendix S1.1. This procedure results in a trip-level dataset. Since our cellphone data only captures a portion of the total trips, we adjust for this by assigning an inflation factor to each trip. To account for varying rates of unobserved trips across different city areas, we allow inflation factors to vary by the neighborhood of the trip’s origin.¹⁰ Specifically, we calibrate these factors to ensure that the number of car trips beginning in each neighborhood in our dataset matches the corresponding number in the Chicago Metropolitan Agency for Planning (CMAP) Household Travel Survey.¹¹

Public transit data We obtain individual public transit trips for the city of Chicago via a partnership between the MIT Transit Lab and the CTA. Each observation corresponds to a passenger swiping in to access the bus or the train station. For buses, we observe the specific bus stop, bus line, and boarding time. For trains, we observe the station and swiping time. Drop-off locations were imputed by Zhao et al. (2007).¹²

These data notably exclude trips taken via the Metra, which is a suburban rail system operating in and around Chicago. Metra is managed by a different agency, the Regional Transportation Authority. An additional limitation is that we do not observe trips paid for via cash or trips whose destination could not be imputed. To account for these sources of missing trips, we assign each observed trip an inflation factor. This inflation factor is computed at the day-mode level such that

$$infl_{dm}T_{dm} = R_{dm},$$

where dm indexes the day-mode, T is the total number of observed trips, and R is

¹⁰ Each neighborhood is a group of about 8-9 CAs. The exact make-up of neighborhoods can be found on [Wikipedia](#).

¹¹ [Source: My Daily Travel survey \(website\)](#)

¹² The inference relies on two observed patterns: a high percentage of riders begin their next trip at the destination of their previous trip, and many complete their final trip of the day at the same station where they began their first. These patterns were validated using travel diary data collected by the New York Metropolitan Transportation Council (Barry et al., 2002).

the observed aggregate daily ridership for the CTA, which we obtain from the [City of Chicago's Data Portal](#). The average such inflation factor is 2.0.

We also do not observe travel times for train trips, and so we are forced to impute these travel times. To do so, we first match each train trip to the historical GTFS schedule data. To compute the match for a given train trip, we first find all scheduled trips between the origin and destination stops of that trip. We then take the match to be the scheduled trip whose boarding time is closest to the observed boarding time. We then take the scheduled travel time as the travel time. This matching process enables us to compute travel times for close to 90% of train trips.

For trips that have no matches in the schedule data, we impute travel times using Google Maps data.¹³ In particular, we first assign each trip one of three time categorizations: weekend (if Saturday or Sunday), peak weekday (if between 5-9:59am or 2-6:59pm on a weekday), or non-peak weekday (otherwise). We then take the time to be the travel time of the matching train trip from the Google Maps data.

We also compute travel distances for each trip. We use the Haversine formula to compute distances, with radius equal to 6371.0088, which is the mean radius of Earth in km. For bus trips, we compute the travel distances as the Manhattan distance between the boarding and alighting coordinates, while for train trips we compute the travel distances as the Euclidean distance between the boarding and alighting coordinates.

S1.4 Market Share Calculations

We first append together the transit, TNP, taxi, and cellphone trips data. We incorporate walk times to bus/train stations from the Google Maps data. We drop any trips that have a negative trip time, trip time exceeding 6 hours, negative prices, or missing values for origin, destination, distance, duration, mode, trip time, or price. Since our trip data is at the vehicle level, we account for unobserved vehicle occu-

¹³ Manual inspection suggests these trips typically involve an unobserved transfer between two lines.

pancy by scaling trip numbers and prices using the average vehicle occupancy for that mode, which we report in Appendix E.4.

We calculate market shares at the (origin CA, destination CA, hour-of-the-week) level using a two-step process. First, we aggregate trips at the (origin CA, destination CA, hour-of-the-week, date) level. We then let the number of car trips be the residual after subtracting public transit, taxi, TNP, and shared trips from the cellphone trips.¹⁴ Car prices are computed as 0.6374 U.S. Dollars per trip mile, which is AAA’s estimate of per mile driving costs for an average 2020 model.¹⁵ Finally, we obtain trip counts at the (origin CA, destination CA, hour-of-the-week, date) level by averaging across dates.

S1.5 Market Size

To compute market shares, we need to take a stance on the size of the market, which captures how many people could be traveling at a given moment in time. For simplicity, we assume that market sizes are proportional to the total number of observed trips. To determine the factor of proportionality, we compare the population of each CA to the total number of trips originating from that CA in the morning hours (5-9:59am) on weekdays. The median ratio across CAs is 2.61. Implicitly, this factor assumes that the number of potential travelers in each CA in these morning hours is given by the total population, which is likely an upper bound. We also compute a more conservative factor by assuming the set of potential travelers is made up of commuters and school-age children, which gives a median factor of 1.48. Corresponding to roughly the midpoint of these two factors, we set our proportionality factor to 2.

We restrict ourselves to markets where we observe car trips so that cars are always an available mode. These markets capture 96% of observed trips.

¹⁴ If the residual is negative we assume that there are no car trips.

¹⁵ Source: [AAA brochure “Your driving costs”](#).

S2 Additional Results

S2.1 Bus Utilization

While our model does not consider capacity constraints for buses when solving for the optimal policy, we can consider *ex-post* the extent to which this constraint might bind. Our results imply frequency reductions for buses that are typically less than 30%. We consider whether these frequency reductions would result in binding capacity constraints, holding ridership levels fixed, by computing the fraction of buses that exceed 70% and 80% utilization across hours of the day. Figure S10 shows that this constraint is unlikely to make a first-order impact on our results as only 10% of buses reach even 70% utilization, and only during the morning and afternoon rush hours.

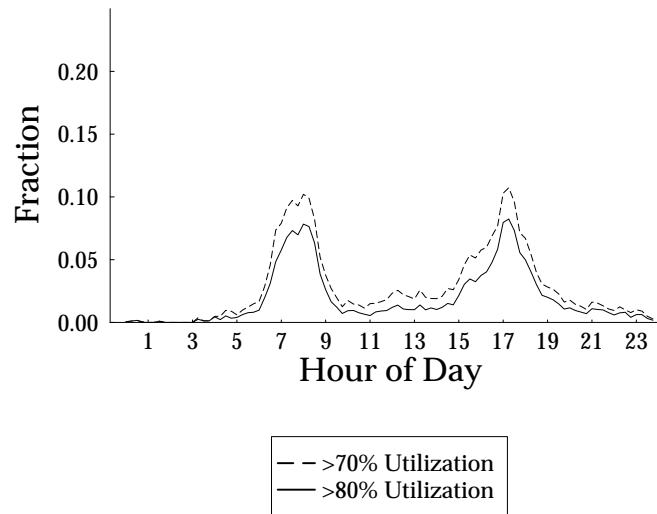


Figure S10: Bus Capacity

Notes: This figure shows the fraction of buses that exceed 80% (solid) and 70% (dashed) utilization over the course of the day.

S2.2 Decomposition of welfare effects

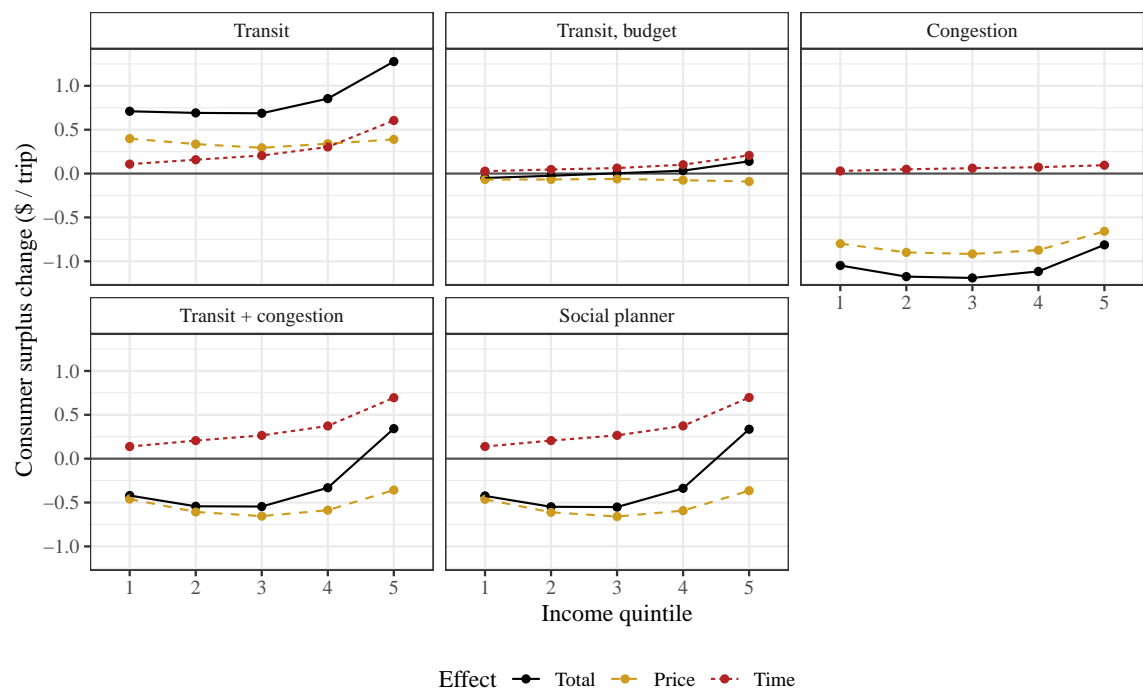
Table S3 decompose changes in consumer surplus into effects from prices and changes in environmental externalities into effects from frequencies and from travelers' substitution. Figure S11 shows how the decomposition of consumer surplus looks for different levels of income.

Table S3: Decomposition of Consumer Surplus and Environmental Externalities

		Status quo	Transit	Transit, budget	Road pricing	Transit + Road pricing	Social planner
Δ CS (\$M/week)	Total	0	26.869	0.945	-32.320	-7.837	-8.021
	Price	0	14.475	-2.923	-34.746	-22.547	-22.776
	Time	0	12.394	3.869	2.426	14.710	14.755
	Capacity	0	9.756	3.764	0	10.069	10.056
	Substitution	0	2.637	0.105	2.426	4.641	4.699
Δ Externality (\$M/week)	Total	0	-0.167	0.238	-2.717	-2.444	-2.455
	Capacity	0	0.398	0.122	0	0.452	0.450
	Substitution	0	-0.565	0.116	-2.717	-2.896	-2.905
Δ Avg. Speed (km/h)		0.00%	0.79%	-0.04%	1.98%	2.49%	2.50%

Notes: This table represents the change in consumer surplus and environmental externalities attributed to different channels. Changes in consumer surplus (first row) are divided into changes in prices (second row) and times (third row). Changes in times are a product in changes in fleet size (fourth row) and substitution of consumers across modes (fifth row). Total changes in externalities (sixth row) are decomposed into changes in fleet size (seventh row) and substitution across consumer (eighth row).

Figure S11: Decomposition of consumer surplus through different channels



Notes: These graphs presents changes in consumer surplus across income quintiles for four different counterfactual scenarios scenarios. Each of the lines represent the change in consumer surplus from each of the channels that affect traveler’s utility.