

# Supplementary Material: Quantifying the importance and location of SARS-CoV-2 transmission events in large metropolitan areas

Alberto Aleta, David Martín-Corral, Michiel Bakker, Ana Pastore y Piontti, Marco Ajelli, Maria Litvinova, Matteo Chinazzi, Natalie E. Dean, M. Elizabeth Halloran, Ira M. Longini, Jr., Alex Pentland, Alessandro Vespignani, Yamir Moreno & Esteban Moro

December 15, 2020

## Contents

<b>1</b>	<b>Mobility data</b>	<b>2</b>
1.1	Stays . . . . .	2
<b>2</b>	<b>Network structure</b>	<b>4</b>
2.1	Agents . . . . .	4
2.2	Contacts . . . . .	6
<b>3</b>	<b>SARS-CoV-2 transmission model</b>	<b>8</b>
<b>4</b>	<b>Calibration</b>	<b>8</b>
<b>5</b>	<b>Effective reproduction number</b>	<b>10</b>
<b>6</b>	<b>Superspreading events</b>	<b>10</b>
<b>7</b>	<b>Sensitivity analysis</b>	<b>11</b>
7.1	Distance to POIs . . . . .	11
7.2	Model parameters . . . . .	12

# 1 Mobility data

The mobility data was obtained from Cuebiq, a location intelligence and measurement company. The dataset consists of anonymized records of GPS locations from users that opted-in to share the data anonymously in the New York metropolitan area over a period of 5 months, from February 2020 to June 2020. In addition to anonymizing the data, the data provider obfuscates home locations to the census block group level to preserve privacy. Data was shared in 2020 under a strict contract with Cuebiq through their Data for Good program where they provide access to de-identified and privacy-enhanced mobility data for academic research and humanitarian initiatives only. All researchers were contractually obligated to not share data further or to attempt to de-identify data. Mobility data is derived from users who opted in to share their data anonymously through a General Data Protection Regulation (GDPR) and California Consumer Privacy Act (CCPA) compliant framework.

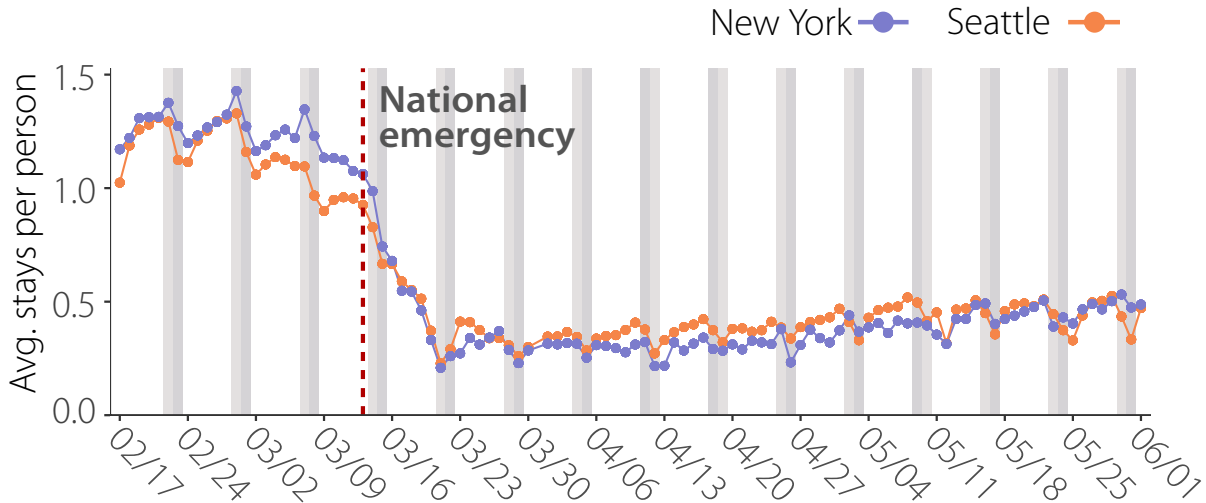
Our sample dataset achieves broad geographic representation for our two populations, in the New York and Seattle metropolitan areas, defined as the Core Based Statistical Areas (CBSA) by the US Census [1]. CBSA are areas that are socioeconomically related to an urban center. This provides a self-contained metropolitan area in which people move for work, leisure or other activities. Some of the CBSAs we consider span several states. For example the New York CBSA contains areas of the state of Connecticut, New Jersey, Philadelphia, and New York. The population and number of anonymous devices detected in the real data by census area are highly correlated for both census county subdivision regions, with a  $\rho = 0.796$  (Pearson correlation) with a CI between 0.783 and 0.807 for the New York region, and a  $\rho = 0.948$  (Pearson correlation) with a CI between 0.937 and 0.957 for the Seattle region.

## 1.1 Stays

From the data we extract “stays”, as the places where anonymous users stayed (stopped) for at least 5 minutes. Each device frequently broadcast its location to a central server by sending its latitude, longitude, device ID, and the exact date and time of the event. When a person spends significant time at a single location, measurement uncertainty will cause a number of events to be scattered around the actual location. To map these events to a single stay with an accurate time and location, we use the Infostop algorithm [2]. First, to extract the locations of stays, the algorithm clusters consecutive events together if the locations are less than 25 meters apart. The location of this cluster is computed by taking the median of the latitudes and longitudes. Moreover, to better estimate the location of places that are visited frequently by the same user, the algorithm also checks whether different clusters appear within 25 meters of each other and assigns a single consistent location to all connected clusters by recomputing the median latitude and longitude. Finally, a stay is registered whenever at least two subsequent events are registered at one of these locations where the first and last event respectively mark the start and end time of the stay. The minimum duration of a stay is set to 5 minutes to make sure we are only including actual contact between people instead of people that, for example, pass each other on an intersection.

Some of the stays happen within or close to places (Points of Interest). We use a dataset of 86k Points of Interest (POI) in the New York metropolitan area and 36k Points of Interest in Seattle metropolitan area collected using the Foursquare API. We attributed a stay to the closest POI up to a distance of 50m, otherwise that stay is discarded. In general, the distance of stays to POIs is much smaller (average distance is 19.43 meters). We have also checked that our results do not depend significantly on the 50 meters threshold (see Supp. Section 7). Stays are then aggregated at place level. Finally we estimate the home Census Block Group of the anonymous users as that in which they are more likely located during nighttime. This results in a dataset of the places people stayed including the points of interest that anonymous users visited and the most likely census block group of where the device owner lives.

In Supp. Figure 1 we can see the daily evolution of the average number of stays per person for New



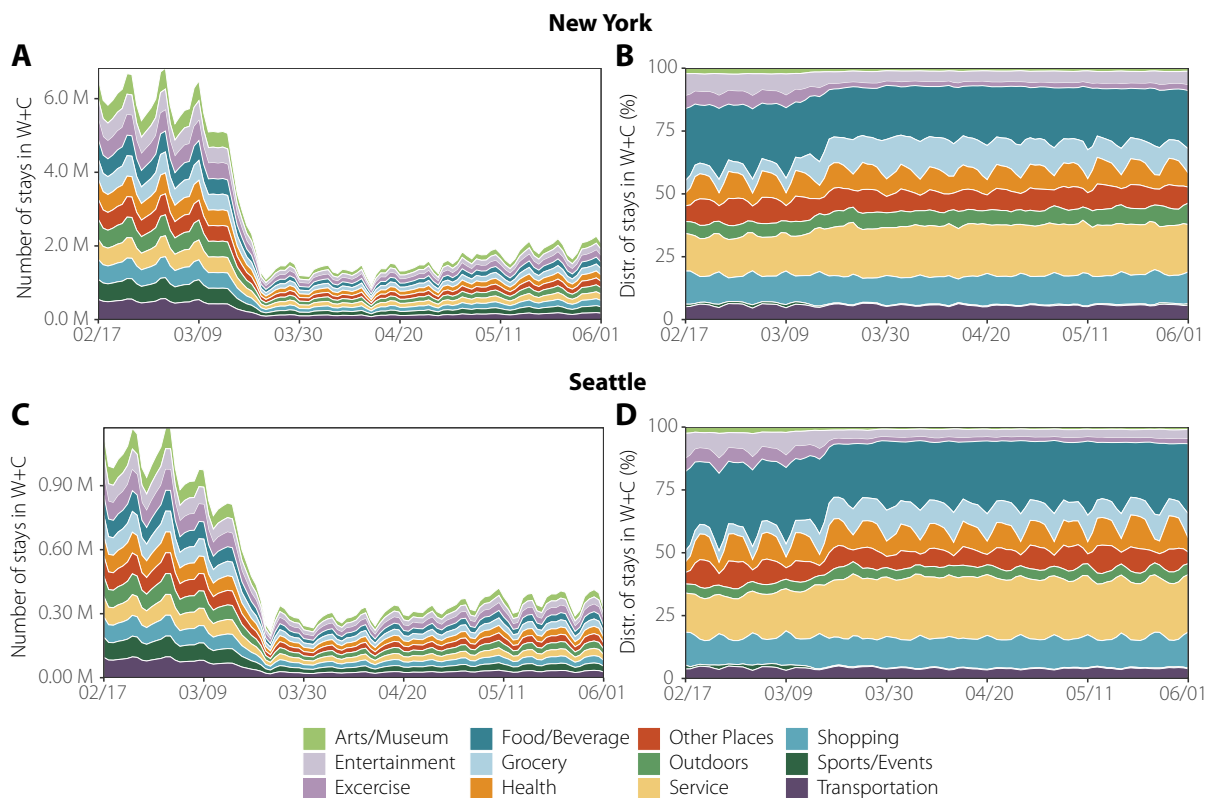
Supplementary Figure 1: The comparative evolution over time for New York and Seattle metropolitan areas of the average number of stays per person. Vertical red dashed line indicates when National Emergency (N.E.) is established.

York and Seattle. Two weeks before we can see that Seattle started to see a small change in the mobility behaviour, however, for New York City we can start to see that pattern one week before the national emergency. After the national emergency there is an abrupt decrease for both cities. Two weeks after the national emergency the average number of stays per person stabilized and starts to a slightly and steady increase, Seattle starts to recover one week before than New York. Eleven weeks after the national emergency, the average number of stays per person has recovered slightly, but it did not recover its basal state for both cities.

As we mentioned above Points of Interest (POIs) are categorized using the Foursquare taxonomy of places which has ten main categories. In our database the New York metropolitan has 572,197 POIs that are distributed as follow Art & Museum (2.1%), College (2.9%), Entertainment (7.6%), Exercise (2.8%), Food & Beverage (17.7%), Grocery (2.6%), Health (7.5%), Other Places (11.3%), Outdoors (8.2%), Religious (1.8%), School (2.3%), Service (16.6%), Shopping (8.3%), Sport & Events (0.6%) and Transportation (6.9%). For the Seattle metropolitan area POIs we have 69,906 POIs that are distributed as follow Art & Museum (2.7%), College (2.3%), Entertainment (7.1%), Exercise (2.7%), Food & Beverage (14.5%), Grocery (2.1%), Health (8.1%), Other Places (13.4%), Outdoors (7.8%), Religious (1.7%), School (1.6%), Service (18.2%), Shopping (8.3%), Sport & Events (0.8%) and Transportation (7.8%). There are also 638 subcategories, see [3] for a complete list of them. We manually curated every subcategory in the taxonomy to be reassign to twelve new principal categories. Arts & Museums, City & Outdoors, Entertainment, Food & Beverages, Grocery, Health, Service, Shopping, Sports (individual), Sports (teams), Transportation and Workplace.

We can see in Supp. Figure 2 the daily evolution of the total number of stays to each category and their fraction distribution. Supp. Figure 2 (a) for New York and (c) for Seattle represent the total number of stays at the community and the workplace layer, we can see a similar pattern as in in Supp. Figure 1 (a) before and after the national emergency. Supp. Figure 2 (b) for New York and (d) for Seattle show normalized number of stays. We can see a reduction of non-essential places after the national emergency due to the social distancing policies.

Finally, in Supp. Figure 3, we can see the comparison of the average time per stay for each city and category before and after the national emergency. There is a significant decrease in time spent per stay for



Supplementary Figure 2: The comparative evolution in the community and workplace layer over time of the total number of stays in the community and workplace layer for each place category (a) for the New York Metropolitan Area and (c) for the Seattle Metropolitan Area and the distribution of stays by place categories (b) for the New York Metropolitan Area and (d) for the Seattle Metropolitan Area.

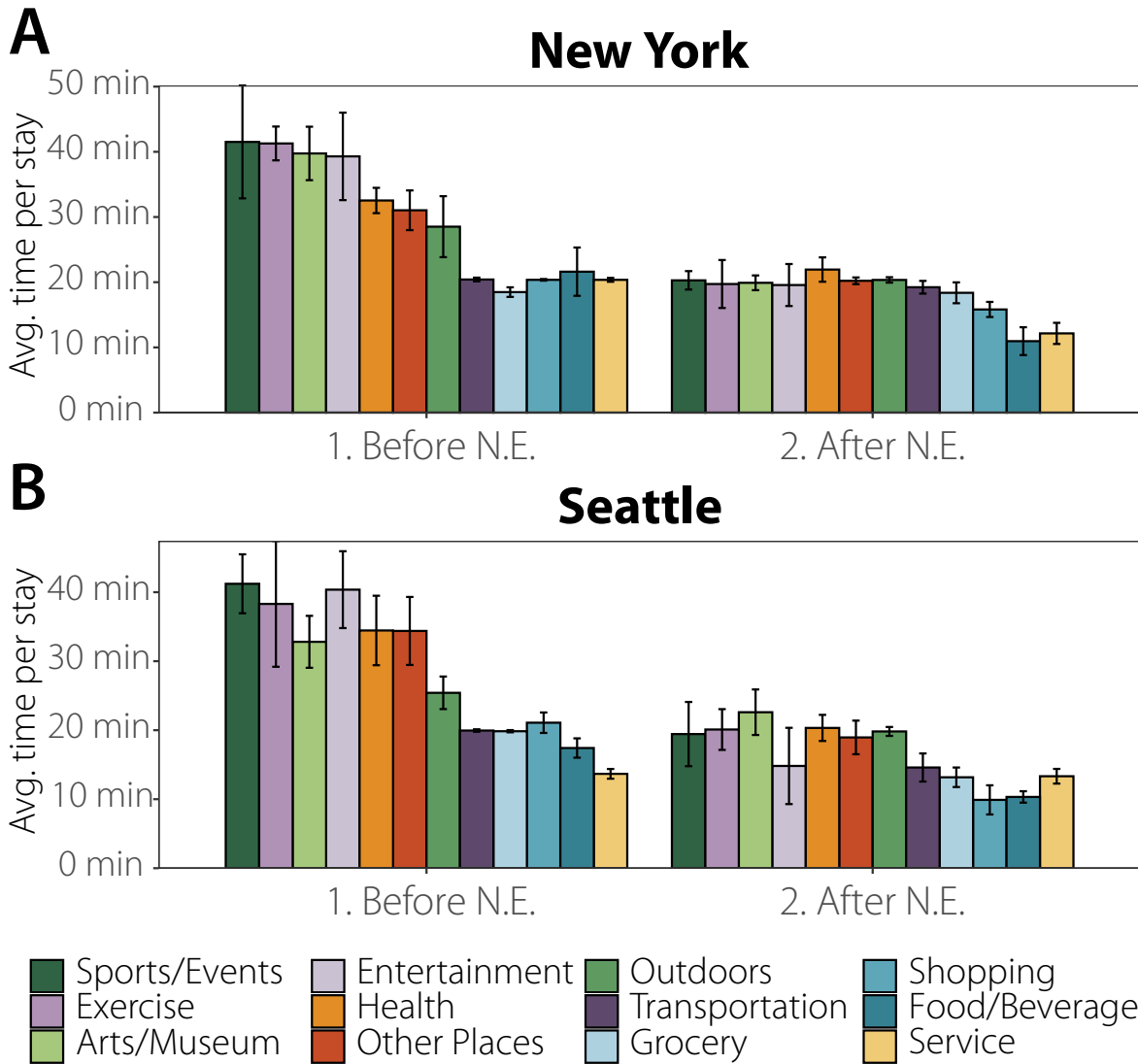
nearly each category in both cities. However, the grocery and the transportation categories are those with the smallest change in the average time for both cities. Moreover, the shopping category does not barely change in New York, but it does in Seattle. On the other hand the Food & Beverages category decrease in New York, but it does not in Seattle.

## 2 Network structure

### 2.1 Agents

Our population consists of two different sub-populations, adults and children. Adults are sampled from anonymous individuals in the mobility data collected by Cuebiq, each adult is associated with a home location assigned to a US Census block group which is provided by our location data provider. With this data we designed a population building pipeline that consists of three steps.

- First step, we build synthetically the number of households, their size and the presence of children based on our adult population and the US Census [4] tables B11016 (Household Type by Household Size) [5] and B11003 (Family Type by Presence and Age of Own Children) [6]
- Second step, we assign adults to households and in case of presence of children we generate them up to reach the size of the household assigned in the first step.



Supplementary Figure 3: Average time per stay for each place category before and after the National Emergency (N.E.) for (a) the New York Metropolitan Area and for (b) the Seattle Metropolitan Area.

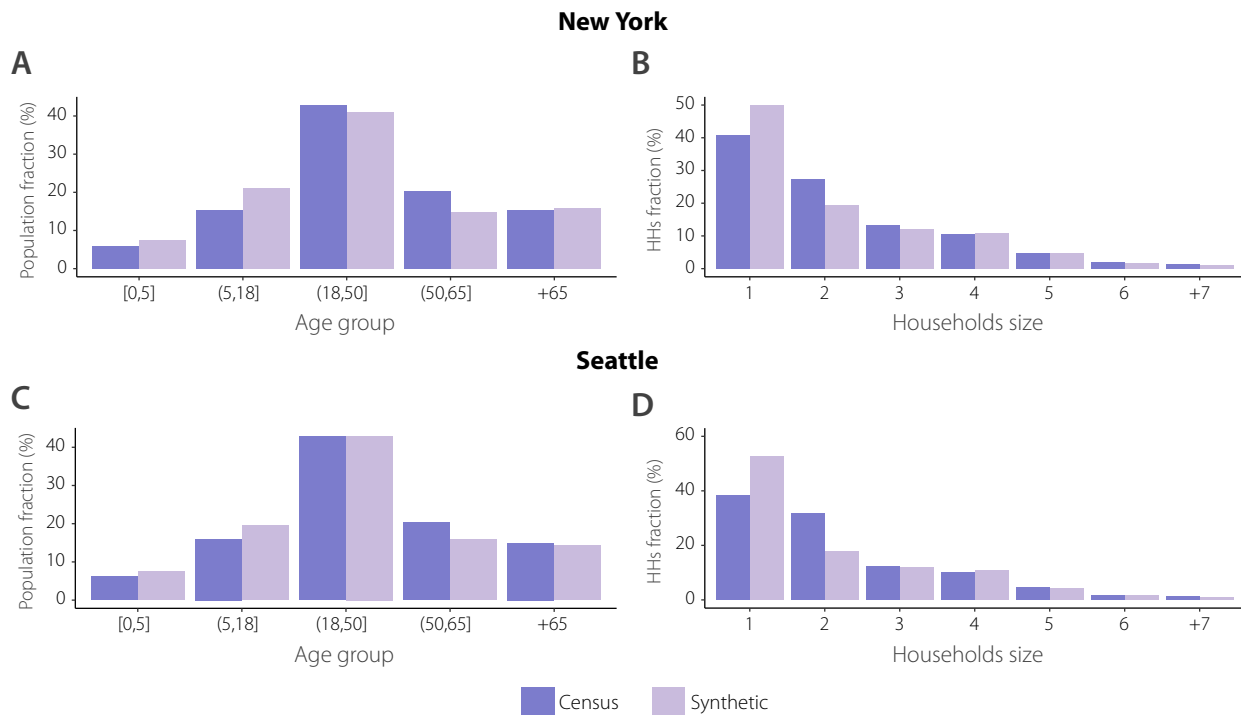
- And final step, we assign ages to nodes using table B01001 (Age by Sex) [7] of age distribution within the Census Block Group.

Following this process we generate two synthetic populations, one for the New York metropolitan area and the other one for the Seattle metropolitan area.

The New York synthetic population consists of 614k agents (3.3% of the population in the New York metropolitan area), 439k (71%) of them are adults and 174k (29%) are children. Age groups are distributed as follows: 45,350 (7.3%) agents for the age group between zero and five years old, 129,259 (21.0%) agents for the age group between six and eighteen years old, 251,175 (40.9%) agents for the age group between nineteen and fifty years old, 91,443 (14.9%) agents for the age group between fifty one and sixty five years old, and final group, 97,023 (15.8%) agents for the age group between sixty six and older. In Supp. Figure 4 (a) we can see the comparison of our synthetic population age distribution against the US census data. All agents together form 290,369 households distributed as follows: 145,079 (50.0%) households with only one agent, 55,985 (19.3%) with two agents, 34,647 (11.9%) with three agents, 31,860 (11.0%) with four agents,

14,190 (4.8%) with five agents, 5,224 (1.8%) with six agents, and finally, 3,384 (0.1%) with seven agents. In Supp. Figure 4 (b) we can see the comparison of our synthetic households population distribution against the US census data.

The Seattle synthetic population consists of 142k agents (3.6% of the population in the Seattle metropolitan area), 103k (72%) of them are adults and 38k (28%) are children. Age groups are distributed as follows: 10,635 (7.5%) agents for the age group between zero and five years old, 27,772 (19.5%) agents for the age group between six and eighteen years old, 61,039 (42.9%) agents for the age group between nineteen and fifty years old, 22,491 (15.8%) agents for the age group between fifty one and sixty five years old, and final group, 20,263 (14.2%) agents for the age group between sixty six and older. In Supp. Figure 4 (c) we can see the comparison of our synthetic population age distribution against the US census data. All agents together form 69,232 households distributed as follows: 36,333 (52.5%) households with only one agent, 12,391 (17.9%) with two agents, 8,341 (12.0%) with three agents, 7,394 (10.7%) with four agents, 2,899 (4.2%) with five agents, 1,127 (1.6%) with six agents, and finally, 747 (0.1%) with seven agents. In Supp. Figure 4 (d) we can see the comparison of our synthetic households population distribution against the US census data.



Supplementary Figure 4: Age groups and households demographics compared against the US Census data. (a) Age groups distribution and (b) households size distribution for the New York Metropolitan Area. (c) Age groups distribution and (d) households size distribution for the Seattle Metropolitan Area.

## 2.2 Contacts

Our network in New York has a total number of 270,785,550 unique contacts, 146,598,503 (54.1%) and 105,129,317 (38.8%) of daily contacts in the community and the workplace layer, respectively, both layers where built using the mobility data from February 15 to June 1st, 641,049 (0.2%) and 18,416,681 (6.8%) are synthetically built for the household and school layers, respectively. On the other hand, in Seattle has a total number of 50,506,786 unique contacts, 20,892,401 (41.4%) and 26,582,687 (52.6%) of daily contacts in

the community and the workplace layer, respectively, both layers where built using the mobility data from February 15 to June 1st, 219,635 (0.4%) and 2,812,063 (5.6%) are synthetically built for the household and school layers, respectively. The community layer is based on estimation of co-presence of two devices in Points of Interest visited by the anonymous users (see Supp. Section 1)

Contacts are built differently in different layers:

- **Community weighted contact network.** The community network is approximated using 5 months of data observation in the New York and the Seattle metropolitan areas from anonymized users. In this layer each agent in our synthetic population represents an anonymous individual of the real population. Contacts are built by estimating co-location of two individuals in the same setting. We use a large database of 572,197 places in the NY and 69,906 in the ST from the Foursquare API. Specifically, the weight,  $\omega_{ijt}^C$ , of a link between individuals  $i$  and  $j$  within the community layer at day  $t$  is computed according to the expression:

$$\omega_{ijt}^C = \sum_p^n \frac{T_{ipt} T_{jpt}}{T_{it} T_{jt}}, \quad \forall i, j$$

where  $T_{ipt}$  is the total time that individual  $i$  was observed at place  $p$  in day  $t$  and  $T_{it}$  is the total time that individual  $i$  has been observed at any place set within the community layer that day  $t$ . Note that while the mobility data set we use is large, co-location events between individuals are still quite sparse. Because of this sparsity, and to protect individual privacy in our analysis, we have adopted this probabilistic approach to measure co-presence in all locations mapped in the dataset. Since agents are representative of the different census areas and groups of the metro areas, our probabilistic approach is a good proxy for the real probability of co-presence between those groups/areas when networks are scaled up to the total population of the New York and Seattle metropolitan areas, that is approximately 18,351,295 and 3,979,845 inhabitants respectively. Finally, for robustness and computational reasons, we have included only links for which  $\omega_{ijt}^C > 0.001$ .

- **Workplace weighted contact network.** For privacy reasons, our data is obfuscated around home and workplaces to the level of Census Block Groups. To get a proxy of contacts at the workplace, we assume that all workers in the same Census Block Groups have a probability to interact. To account for the potential number of working places in that area, we weight that probability by the number of POIs at the same census block group. Therefore, the contact weight,  $\omega_{ijt}^W$ , of a link between individuals  $i$  and  $j$  within the same workplace at day  $t$  is given by:

$$\omega_{ijt}^W = \sum_{\alpha \in \text{CBG}} \frac{\omega_{i\alpha t} \omega_{j\alpha t}}{N_{POI}(\alpha)}, \quad \forall i, j$$

where  $N_{POI}(\alpha)$  is the number of POIs in census block group  $\alpha$ ,  $\omega_{i\alpha t}$  is the probability of observing an individual at her workplace within census block group  $\alpha$  at day  $t$ . As before, we have included only links for which  $\omega_{ijt}^W > 0.001$ .

**2) Household weighted contact network.** We first identify individuals' approximate home place as their most likely visited census block group at night. Then we assign a type of household based on Table B11016: Household Type by Household Size from US Census 2018[5], and mix individuals that live in the same block according to statistics of household type and size. Finally, children are assigned to households. We also assign individuals an age group based on Table B01001: Sex by age from the US Census 2018 [7]. To assign weights, we assume that the probability of interaction within a household is proportional to the number of people living in the same household (well-mixing). Therefore, the weight,  $\omega_{ij}^H$ , of a link between individuals  $i$  and  $j$  within the same household is given by:

$$\omega_{ij}^H = \frac{1}{(n_h - 1)}$$

where  $n_h$  is the number of household members. This fraction is assumed to be the same for all individuals in the population. We assume this layer is static throughout our period.

- **School weighted contact network.** To calculate the weights of the links at the school layer, we mix together all children that live in the same school catchment area. Interactions are considered well-mixed, hence, the probability of interaction at a school is proportional to the number of children at the same school. Therefore, the weight,  $\omega_{ij}^S$ , of a link between children  $i$  and  $j$  within the same school is given by:

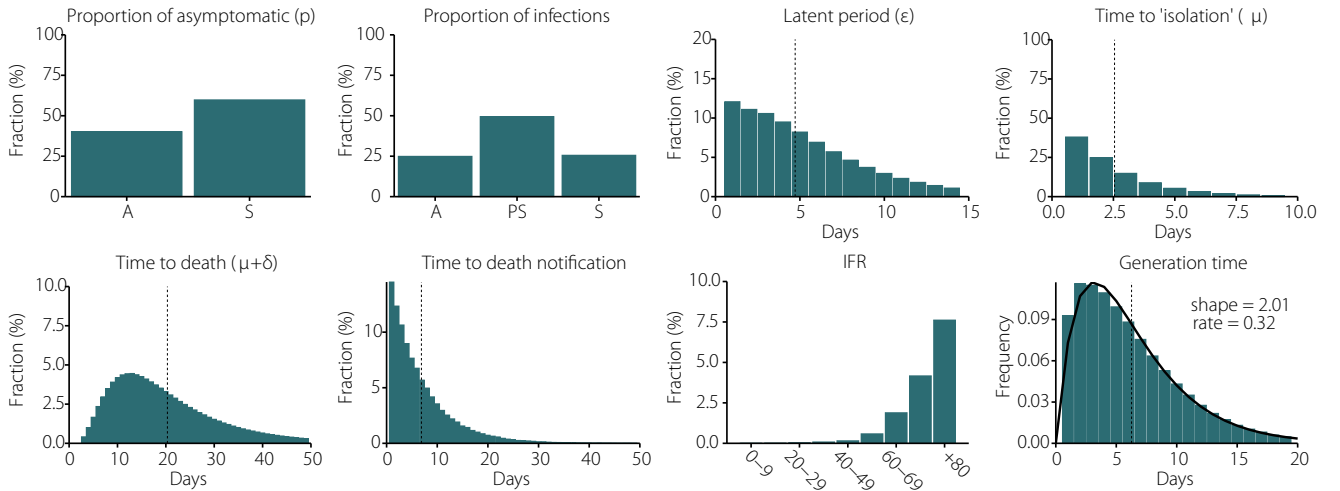
$$\omega_{ij}^S = \frac{1}{(n_s - 1)}$$

where  $n_s$  is the number of school members. This layer is removed on March 16 in both metropolitan areas to account for the imposed school closure.

To calibrate the relative importance of each layer in the spreading process we further re-normalize each network so that the average degree in each layer is 4.11 in the household layer, 11.41 in the education layer, 8.07 in the workplace layer and 2.79 in the community layer [8].

### 3 SARS-CoV-2 transmission model

The values of all the disease parameters used for simulating the transmission dynamics are given in Supp. Table 1. Supp. Figure 5 shows the numerical distributions of these parameters as resulting from simulations of the model, computed for the case of New York with  $R_0 = 3.4$  (see Supp. Section 4).



Supplementary Figure 5: Numerical distributions of the model parameters extracted from the simulations performed for New York with  $R_0 = 3.4$ . The generation time distribution is well fitted by a gamma distribution with shape = 2.01 and rate = 0.32.

### 4 Calibration

The model has two free parameters: (1) the number of infected individuals in each city on the first day for which we have data to build the interaction networks (02/17/2020) and (2) the value of  $R_0$ .

Simulations are initialized with 1 infected individual and then they run using the information from the first week available (02/17/2020 to 02/23/2020) in a loop until the number of latent individuals in

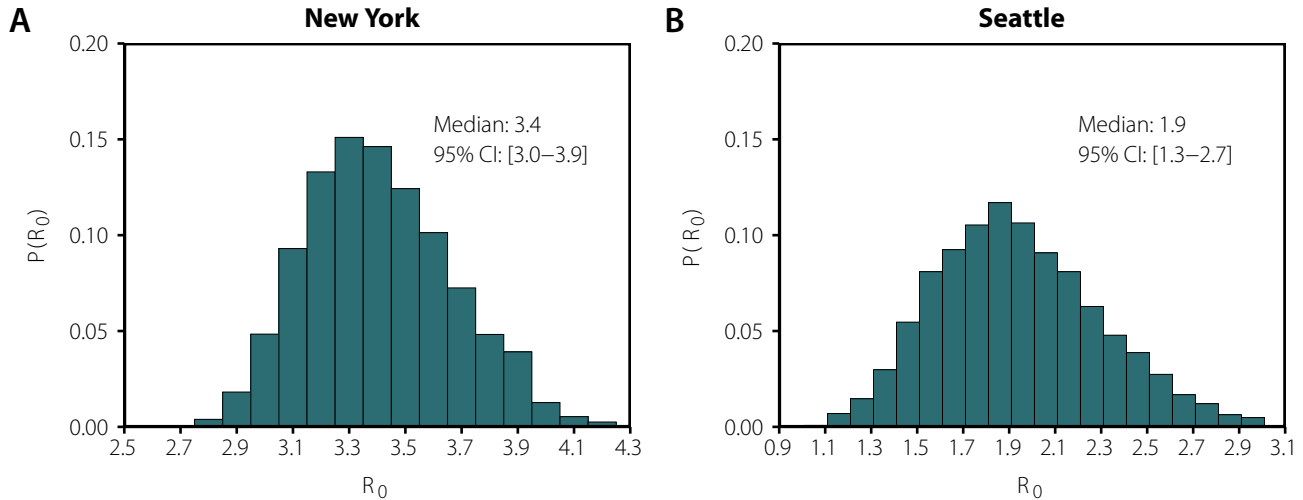


Parameters	Description	Age group	Value	Ref.
$r$	relative infectiousness of asymptomatic individuals	-	50%	†
$k$	proportion of pre-symptomatic transmission	-	50%	[9]
$\epsilon^{-1}$	incubation period (gamma distributed)	-	shape = 2.08 rate = 0.33	[10]
$p$	proportion of asymptomatic	-	40%	[9]
$\gamma^{-1}$	pre-symptomatic period	-	2 days	[11]
$\mu^{-1}$	time to isolation	-	2.5 days	
$\delta^{-1}$	days from isolation to death	-	12.5	[9]
IFR	infection fatality ratio	0-9 10-19 20-29 30-39 40-49 50-59 60-69 70-79 $\geq 80$	0.00161% 0.00695% 0.0309% 0.0844% 0.161% 0.595% 1.93% 4.28% 7.80%	[12]‡
$T_n$	Notification of death	-	7 days	[9]
$\theta$	outdoor transmissibility	-	0.05	[13]

Supplementary Table 1: Baseline set of parameters. †: assumed ;\*: calibrated to the generation time  $T_g$ ; ‡ Only applied to symptomatic individuals. As such, a correction factor of  $1/(1-p)$  is applied to all age groups.

the population matches the values estimated by the GLEAM model [14]. In particular, 292 in New York City and 39 in Seattle. Then, time is reset and the simulation runs on calendar time from 02/17/2020 to 06/01/2020 (each step corresponds to 1 day).

We use an Approximate Bayesian Computation (ABC) rejection algorithm to obtain the posterior distribution of  $R_0$ . We sample  $R_0$  from a uniform prior in the range 1.5 to 4.5 and compare the output of the model with the weekly estimated number of deaths as a consequence of COVID-19 for each city [15]. The obtained posterior distribution  $P(R_0 = x|E)$  is shown in Supp. Figure 6.



Supplementary Figure 6: Posterior distribution of  $R_0$  given the number of weekly deaths in each region as evidence.

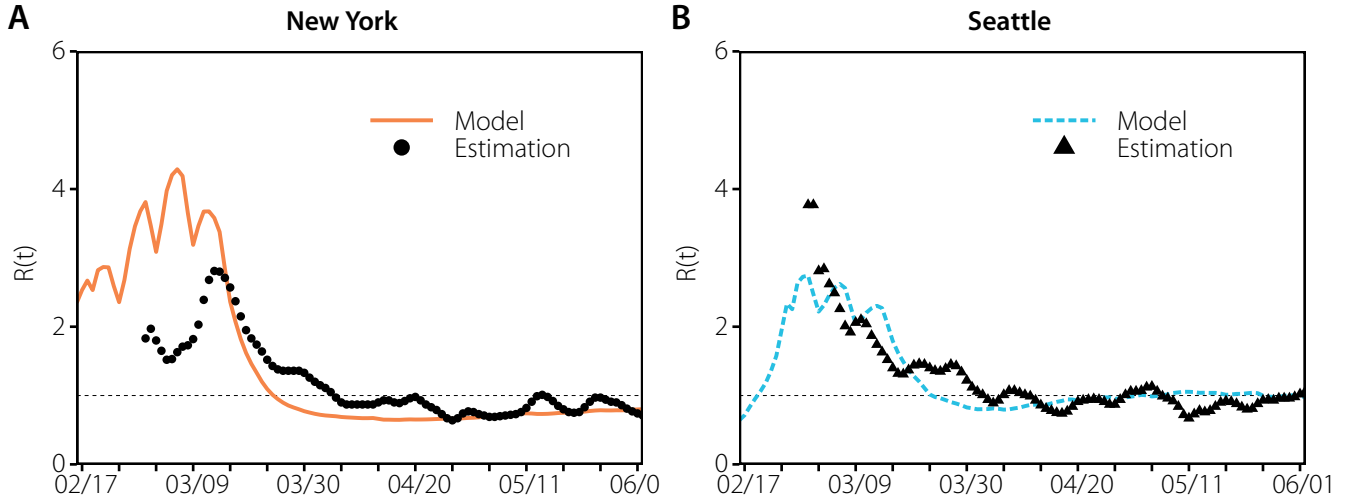
## 5 Effective reproduction number

The effective reproduction number can be estimated using case count data as reported by the authorities. In [16], the authors estimated this quantity from 14/03/2020 for different areas of the world using the daily positive increase cases from Johns Hopkins University Center [15] and the EpiEstim method [17]. However, while in a simulation all the information regarding the infection process is available, in real data usually only the date of infection notification is available. Since the time between the notification of an infection and the date of actual infection can differ by several days, the authors propose to shift the curve by 5 days (note that in countries that report this quantity the value changed over time [18]). As we show in Supp. Figure 7, our model results match pretty well with the data, even though it seems that the shift of the curve should be slightly larger than 5 days. This is consistent to what has been observed in some European countries during the first wave [19].

## 6 Superspreading events

In heterogeneous population it is possible for an infected individual to produce an usually large number of secondary cases. This is known as a super-spreading event (SSE). To define a SSE we follow Lloyd-Smith et al [20]:

1. Estimate the effective reproduction number,  $R$



Supplementary Figure 7: Effective reproduction number in both areas as obtained by our model and estimated by [16].

2. Compute a Poisson distribution with mean  $R$
3. Define a SSE as any infected individual who infects more than the 99-th percentile of the Poisson distribution.

In Supp. Figure 8 we test the hypothesis of the 20/80 rule according to which 20% of the infected individuals produce 80% of the infections. Note that this does not imply that said 20% of individuals are super-spreaders. In fact, the large majority of them do not produce any secondary infections, inline with what has been observed in highly detailed empirical studies [21].

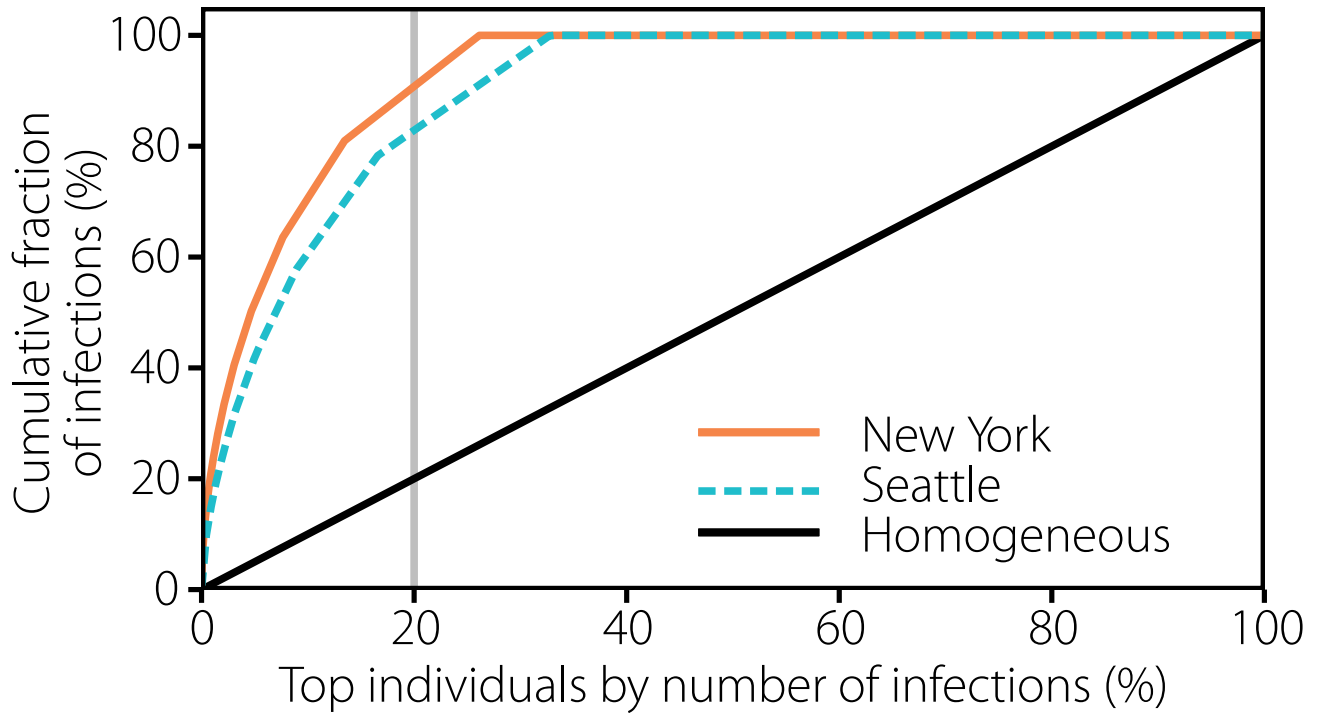
In Supp. Table 2 we report the probability of having a SSE within each category before and after the declaration of the National Emergency. We observe a drastic reduction of the probability after 03/13.

## 7 Sensitivity analysis

### 7.1 Distance to POIs

While constructing the network, we attributed a stay to a given POI if it was no further than 50 meters from the POI center. In this section we test more strict conditions for that attribution, i.e. a threshold of just 10 meters. Note that this more strict condition for attribution lowers the number of potential visitors to the POI but also lowers the distance between people in the venue, making physical contact more likely. In Supp. Figure 9 we show the results for this scenario.

A more restrictive definition of stay yield a much sparser network in the community layer, while it does not affect the rest of the system. We can see that to obtain the observed number of deaths under these conditions, the value of  $R_0$  must be much larger. This is related to the fact that the disease is specially dangerous for the elderly. Since those individuals interact mostly in the community layer (since they do not attend schools nor workplaces), and this layer is now sparser, we need to increase the transmissibility for the disease to reach those individuals. As a consequence, this also increases the transmission in the rest of the layers, yielding a larger overall prevalence. Nevertheless, the distribution of infections across settings is fairly similar, signaling that the results are robust to this definition.



Supplementary Figure 8: Individuals are ranked according to the number of infections they produce. The cumulative fraction of infections found in both cities is compared with the one that would be obtained in a completely homogeneous system.

## 7.2 Model parameters

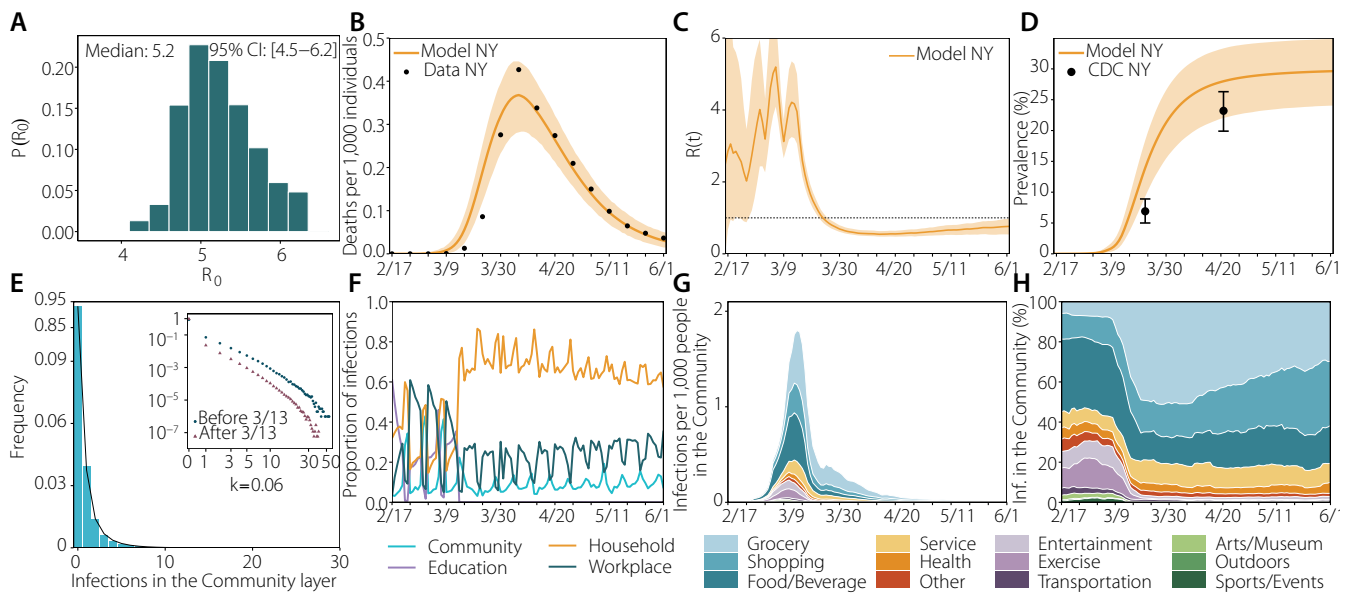
To test the dependency of the results with the values assumed in the model, we have explored three different scenarios: larger transmissibility during the pre-symptomatic phase ( $k = 0.75$ ), Supp. Figure 10; longer time from death to notification ( $T_n = 14$  days), Supp. Figure 11; and larger outdoor transmission ( $\theta = 0.10$ ), Supp. Figure 12. The results are consistent in all cases with only slight variations on the value of  $R_0$ .

Probability of a supers-spreading event (%)				
Category	New York		Seattle	
	Before 03/13	After 03/13	Before 03/13	After 03/13
Arts/Museum	7.32 [7.22-7.42]	0.53 [0.51-0.54]	0.65 [0.51-0.81]	0.05 [0.01-0.10]
Entertainment	2.20 [2.12-2.22]	0.15 [0.15-0.16]	2.91 [2.71-3.12]	0.32 [0.23-0.39]
Excercise	1.82 [1.80-1.84]	0.38 [0.38-0.39]	1.53 [1.40-1.66]	0.64 [0.55-0.73]
Food/Beverage	0.54 [0.53-0.55]	0.19 [0.19-0.19]	0.29 [0.26-0.33]	0.27 [0.25-0.29]
Grocery	3.10 [3.08-3.13]	1.49 [1.48-1.49]	0.97 [0.85-1.10]	1.40 [1.36-1.45]
Health	0.15 [0.14-0.16]	0.13 [0.13-0.13]	0.00 [0.00-0.01]	0.10 [0.08-0.12]
Other	1.47 [1.45-1.50]	0.10 [0.10-0.10]	0.46 [0.35-0.57]	0.03 [0.01-0.05]
Outdoors	0.01 [0.01-0.02]	0.00 [0.00-0.00]	0.00 [0.00-0.00]	0.00 [0.00-0.00]
Service	0.66 [0.65-0.67]	0.21 [0.21-0.21]	0.07 [0.04-0.10]	0.18 [0.17-0.20]
Shopping	1.92 [1.90-1.94]	0.96 [0.95-0.96]	0.13 [0.10-0.17]	0.21 [0.20-0.23]
Sports/Events	8.32 [8.19-8.44]	4.10 [3.95-4.22]	0.54 [0.37-0.72]	0.03 [0.00-0.08]
Transportation	0.30 [0.29-0.32]	0.07 [0.06-0.07]	0.00 [0.00-0.00]	0.01 [0.00-0.03]
All	1.80 [1.79-1.81]	0.80 [0.79-0.80]	1.33 [1.29-1.38]	0.53 [0.52-0.54]

Supplementary Table 2: Probability that an individual will cause a super-spreading event as defined in [20]. We aggregate all the infections produced by each individual within each category for the given period of time, and compute the fraction of individuals who produce a super-spreading event out of the total number of individuals infecting someone in that category. In brackets the 95% C.I. computed using a bootstrap percentile method is shown.

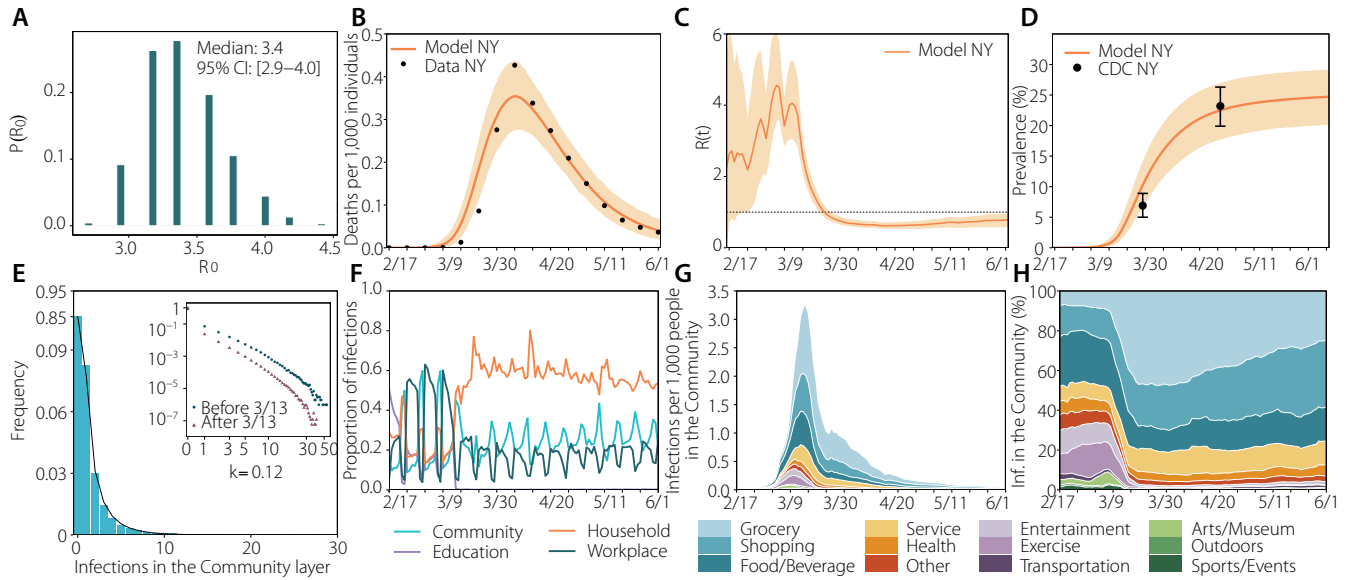
## References

- [1] Bureau, U. S. C. Core-Based Statistical Areas. <https://www.census.gov/topics/housing/housing-patterns/about/core-based-statistical-areas.html> (2019).
- [2] Aslak, U. & Alessandretti, L. Infostop: Scalable stop-location detection in multi-user mobility data. *arXiv preprint arXiv:2003.14370* (2020).
- [3] Foursquare Venue Category Hierarchy. <https://developer.foursquare.com/docs/build-with-foursquare/categories/>. Accessed: 09-12-2020.
- [4] U.S. Census Bureau. 2018 American Community Survey 5-Year Data (2019). URL <https://www.census.gov/programs-surveys/acs>.
- [5] U.S. Census Bureau Table B11016: Household Type by Household Size. 2018 American Community Survey 5-Year Data (2019). URL <https://censusreporter.org/tables/B11016/>.
- [6] U.S. Census Bureau Table B11003: Family Type by Presence and Age of Own Children. 2018 American Community Survey 5-Year Data (2019). URL <https://censusreporter.org/tables/B11003/>.
- [7] U.S. Census Bureau Table B01001: Sex by Age. 2018 American Community Survey 5-Year Data (2019). URL <https://censusreporter.org/tables/B01001/>.
- [8] Mistry, D. *et al.* Inferring high-resolution human mixing patterns for disease modeling. *arXiv* (2020). URL <https://arxiv.org/abs/2003.01214v1>. 2003.01214.
- [9] Coronavirus Disease 2019 (COVID-19) planning scenarios (2020). URL <https://www.cdc.gov/coronavirus/2019-ncov/hcp/planning-scenarios.html>. [Online; accessed 15. Dec. 2020].



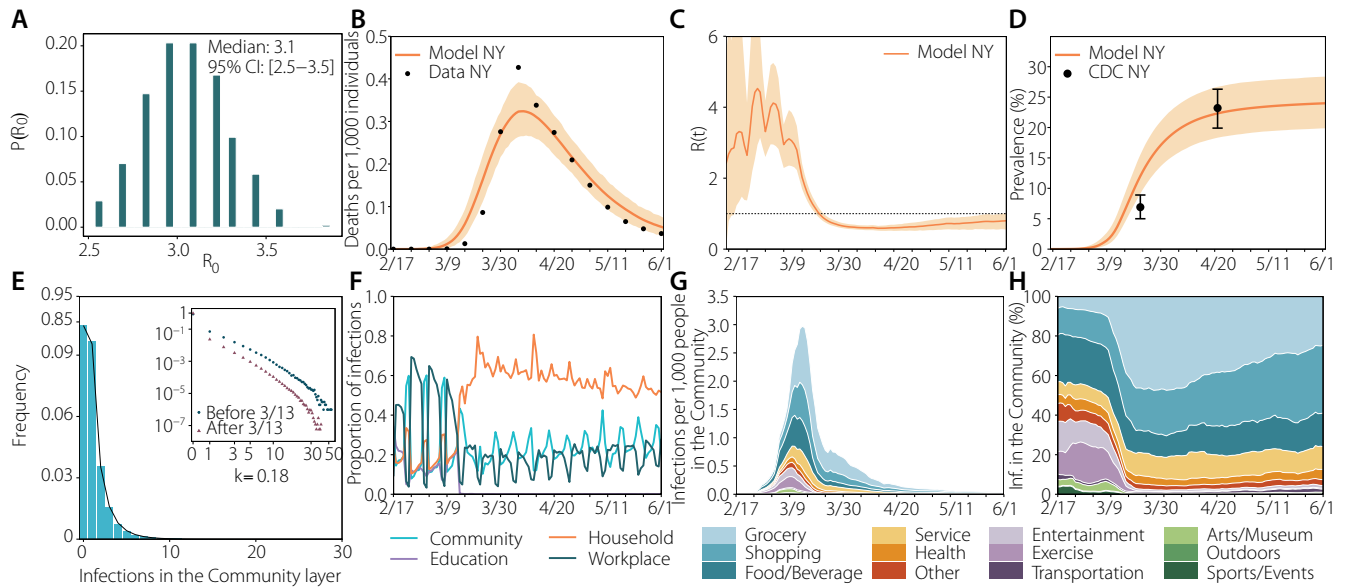
Supplementary Figure 9: Results with a more restricted definition of stay for the case of New York: (a) estimated  $R_0$ ; (b) number of deaths (fit); (c) estimated  $R_t$ ; (d) prevalence; (e) distribution of infections; (f) proportion of infections per layer; (g) infections per setting; (h) normalized infections per setting.

- [10] Hu, S. *et al.* Infectivity, susceptibility, and risk factors associated with SARS-CoV-2 transmission under intensive contact tracing in Hunan, China. *medRxiv* 2020.07.23.20160317 (2020). URL <https://doi.org/10.1101/2020.07.23.20160317>. 2020.07.23.20160317.
- [11] Backer, J. A., Klinkenberg, D. & Wallinga, J. Incubation period of 2019 novel coronavirus (2019-nCoV) infections among travellers from Wuhan, China, 20–28 January 2020. *Eurosurveillance* **25**, 2000062 (2020).
- [12] Verity, R. *et al.* Estimates of the severity of coronavirus disease 2019: a model-based analysis. *The Lancet Infectious Diseases* **20**, 669–677 (2020).
- [13] Weed, M. & Foad, A. Rapid Scoping Review of Evidence of Outdoor Transmission of COVID-19. *medRxiv* 2020.09.04.20188417 (2020). URL <https://doi.org/10.1101/2020.09.04.20188417>. 2020.09.04.20188417.
- [14] Davis, J. T. *et al.* Estimating the establishment of local transmission and the cryptic phase of the COVID-19 pandemic in the USA. *medRxiv* 2020.07.06.20140285 (2020). URL <https://doi.org/10.1101/2020.07.06.20140285>. 2020.07.06.20140285.
- [15] Dong, E., Du, H. & Gardner, L. An interactive web-based dashboard to track COVID-19 in real time. *Lancet Infect. Dis.* **20**, 533–534 (2020).
- [16] Adam, D. A guide to R — the pandemic’s misunderstood metric. *Nature* **583**, 346–348 (2020).
- [17] Cori, A., Ferguson, N. M., Fraser, C. & Cauchemez, S. A New Framework and Software to Estimate Time-Varying Reproduction Numbers During Epidemics. *Am. J. Epidemiol.* **178**, 1505–1512 (2013).
- [18] Aleta, A. & Moreno, Y. Evaluation of the potential incidence of COVID-19 and effectiveness of containment measures in Spain: a data-driven approach. *BMC Med.* **18**, 1–12 (2020).
- [19] Starnini, M., Aleta, A., Tizzoni, M. & Moreno, Y. Impact of the accuracy of case-based surveillance data on the estimation of time-varying reproduction numbers. *medRxiv* 2020.06.26.20140871 (2020). URL <https://doi.org/10.1101/2020.06.26.20140871>. 2020.06.26.20140871.

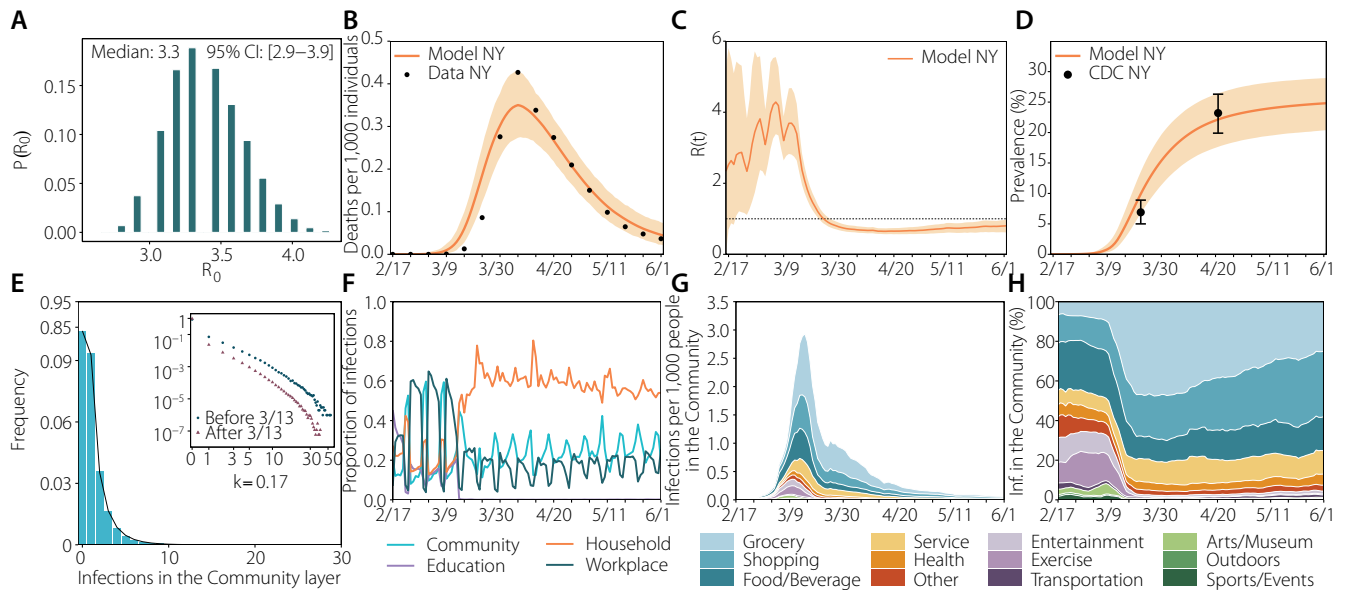


Supplementary Figure 10: Main results in New York with larger pre-symptomatic transmissibility: (a) estimated  $R_0$ ; (b) number of deaths (fit); (c) estimated  $R_t$ ; (d) prevalence; (e) distribution of infections; (f) proportion of infections per layer; (g) infections per setting; (h) normalized infections per setting.

- [20] Lloyd-Smith, J. O., Schreiber, S. J., Kopp, P. E. & Getz, W. M. Superspreading and the effect of individual variation on disease emergence. *Nature* **438**, 355–359 (2005).
- [21] Sun, K. *et al.* Transmission heterogeneities, kinetics, and controllability of SARS-CoV-2. *medRxiv* 2020.08.09.20171132 (2020). URL <https://doi.org/10.1101/2020.08.09.20171132>. 2020.08.09.20171132.



Supplementary Figure 11: Main results in New York with longer time to death notification: (a) estimated  $R_0$ ; (b) number of deaths (fit); (c) estimated  $R_t$ ; (d) prevalence; (e) distribution of infections; (f) proportion of infections per layer; (g) infections per setting; (h) normalized infections per setting.



Supplementary Figure 12: Main results in New York with larger outdoor transmissibility: (a) estimated  $R_0$ ; (b) number of deaths (fit); (c) estimated  $R_t$ ; (d) prevalence; (e) distribution of infections; (f) proportion of infections per layer; (g) infections per setting; (h) normalized infections per setting.