

1 Supplementary Information for

2 **Temporal Contact Graph Reveals the Evolving Epidemic Situa-**
3 **tion of COVID-19**

4 Mincheng Wu^{1,†}, Chao Li^{1,†}, Zhangchong Shen¹, Shibo He^{1,7,*}, Lingling Tang², Jie Zheng³, Yi
5 Fang⁴, Kehan Li¹, Yanggang Cheng¹, Zhiguo Shi^{5,7}, Guoping Sheng², Yu Liu^{4,7}, Jinxing Zhu⁴,
6 Xinjiang Ye⁴, Jinlai Chen^{4,7}, Wenrong Chen⁴, Lanjuan Li^{6,*}, Youxian Sun¹, Jiming Chen^{1,7,*}

7 ¹*College of Control Science and Engineering, Zhejiang University, Hangzhou 310027, China.*

8 ²*Shulan (Hangzhou) Hospital Affiliated to Shulan International Medical College, Zhejiang Shuren*
9 *University, Hangzhou, China.*

10 ³*Zhejiang Institute of Medical-care Information Technology, Hangzhou, China.*

11 ⁴*Westlake Institute for Data Intelligence, Hangzhou, China.*

12 ⁵*College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou,*
13 *China.*

14 ⁶*State Key Laboratory for Diagnosis and Treatment of Infectious Diseases, Zhejiang University,*
15 *Hangzhou, China.*

16 ⁷*Data Intelligence Research Center, Institute of Wenzhou, Zhejiang University, Wenzhou, China.*

17 * Corresponding author. E-mail:s18he@zju.edu.cn, ljli@zju.edu.cn, cjm@zju.edu.cn.

18 † These authors contributed equally.

19 **Supplementary Note I: Data description**

20 The original data of location-related information was collected by location-based service (LBS)
21 providers in China, which have a long-term cooperation agreements with Westlake Institute for
22 Data Intelligence (the affiliate of some coauthors of this article). The location-related information
23 was uploaded every time smartphone users are using LBS. Smartphone users authorized such data
24 collection process. According to the Personal Information Security Specification of China (2019),
25 privacy protection mechanisms such as perturbation and pseudonymization are adopted during data
26 collection.

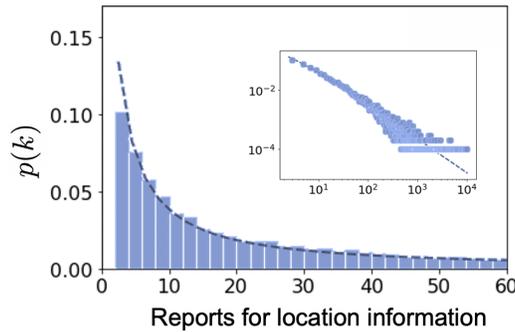
27 Since the outbreak of COVID-19 in China, Prof. Lanjuan Li, one of the five members in the
28 high-level expert team for COVID-19 convened by the National Health Commission of China, Dr.
29 Lingling Tang and Dr. Guoping Sheng, famous epidemiologists from Shulan (Hangzhou) Hospital,
30 have been working on the front line for prevention and control of COVID-19. Westlake Institute
31 for Data Intelligence led by Yi Fang also worked closely with Center for Disease Control and
32 Prevention (CDC) on tracking contacts with infectious cases. Based on Public Health Emergencies
33 Regulations of China, the local authority provided the device identifications of confirmed cases.
34 Such information can be incorporated to link the health status of smartphone users since the app
35 accounts in China are registered by these identifications.

36 Many recent works used very similar data collected from different sources. For example,
37 Kraemer et al. evaluated how the individual mobility affects the disease transmission, and, as
38 control measures can be reflected by the changes in individual mobility, further showed the positive
39 effect of control measures in Wuhan on mitigating the spread of COVID-19¹. Nicholas et al.
40 used country-wide aggregated mobile phone data to show that the outflow from Wuhan accurately
41 predicts the distribution of infections across all of China, and developed a “risk source” model to
42 forecast confirmed cases and identify high risk locales at early stage. This work also derived the
43 geographic spread and growth pattern of COVID-19². Vespignani et al. used the data collected by
44 Cuebiq Inc³, which is similar with the location data we have. However, they simulate the infectious
45 cases using a stochastic discrete-time compartmental model, which is not calibrated to account for
46 the specific evolution in Boston. From the above summary, we can see that we utilized the real
47 health status of users to construct the temporal contact graph, and provides an empirical evidence
48 to reveal evolving epidemic situation of COVID-19. These distinguishes our work from existing
49 ones. In other words, we took the first attempt to construct the contact graph between susceptible
50 and infectious individuals to represent the process of digital contact tracing, while previous studies

51 did not utilize the information of confirmed cases.

52 We note that all individual location-related data and health status information were collected,
53 stored, and used by following the Personal Information Security Specification (2019) and Public
54 Health Emergencies Regulations of China. All raw data was stored in specialized data servers
55 with limited access by LBS providers. For the research purpose, the Westlake Institute for Data
56 Intelligence constructed the temporal contact graph according to the contact model we built and
57 provided it for further analysis. All the results in the article can be obtained and validated by the
58 temporal contact graph, which does not contain private geographic and identity information about
59 the contacted individuals and the confirmed cases. This project was reviewed and approved by the
60 Medical Ethics Committee of School of Medicine, Zhejiang University.

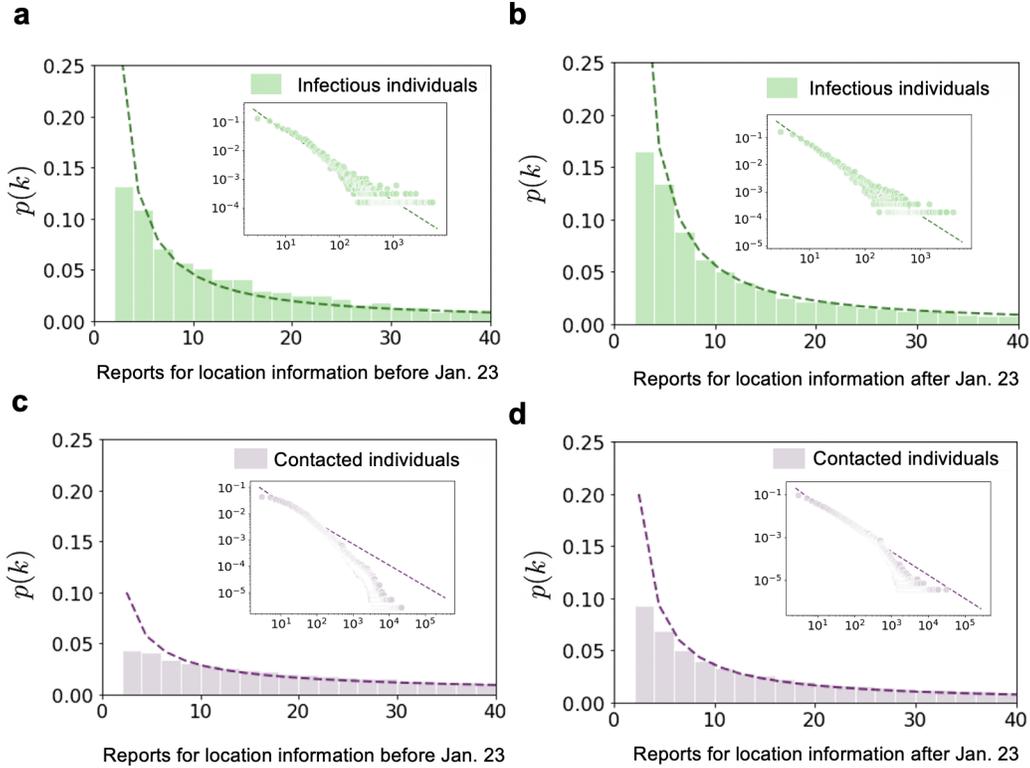
61 Supplementary Fig. 1 analyzes the distribution of the number of reports for all the confirmed
62 cases and the contacted individuals under the proposed contact model. Every report is counted no
63 matter whether it contributes to a contact or not. The number of reports related to location-based
64 information follow well-known power-law distribution, i.e., $p(k) \propto k^{-\gamma}$, where the average piece
of information $\langle k \rangle$ equals to 311.63 and the power exponent γ equals to 1.15.



Supplementary Figure 1: Distribution of reports history in 20 days. Distribution of reports for location-based information, and 10.11% of users contribute less than three time in 20 days.

65

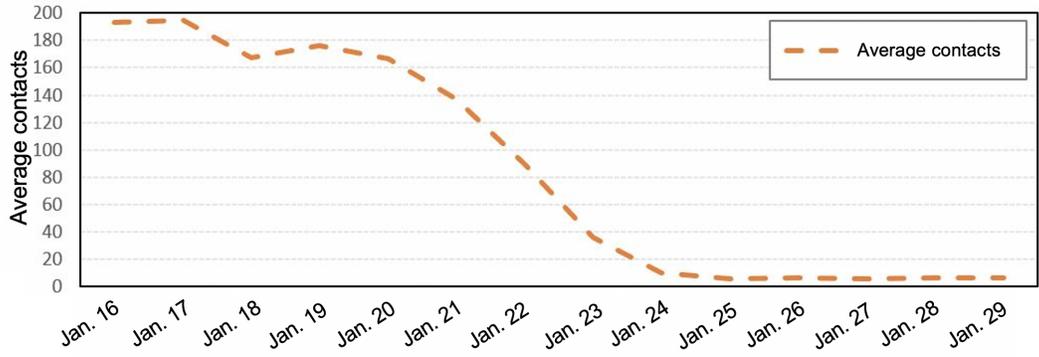
66 As shown in Supplementary Fig. 2, we studied the crowdsourced data before and after the
67 date of the travel restriction, i.e., 23 January, 2020. Among the data, the location information,
68 which is used to calculate contacts, is preprocessed firstly. Then, we count the reporting data of
69 each user in seven days and study the distribution of the crowdsourced data. In the first instance, we
70 can find that the contributed data are sparse and the majority of users contribute less than three time



Supplementary Figure 2: Distribution of reporting crowdsourced data before and after 23 January, 2020. **a** Distribution of confirmed cases' reports for location information before 23 January, 2020. **b** Distribution of confirmed cases' reports for location information after 23 January, 2020. **c** Distribution of contacted individuals' reports for location information before 23 January, 2020. **d** Distribution of contacted individuals' reports for location information after 23 January, 2020.

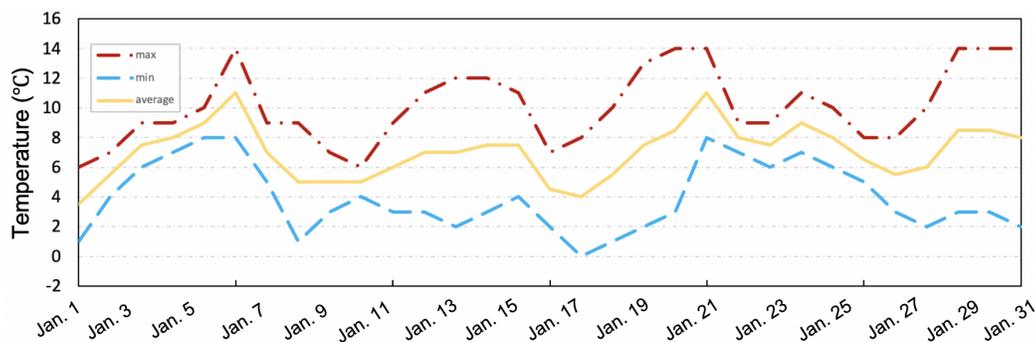
71 in seven days because smartphone users report data in a very low and irregular frequency. Besides,
 72 the distributions of the reporting data from both confirmed cases and contacted individuals in the
 73 whole period follow power-law distributions, following $p(k) \propto k^{-\gamma}$, where $\langle k \rangle = 93.04$, $\gamma =$
 74 1.2 for location information reported by confirmed cases before 23 January, 2020 (Supplementary
 75 Fig. 2a) and $\langle k \rangle = 58.68$, $\gamma = 1.3$ for location information reported by confirmed cases after
 76 23 January, 2020 (Supplementary Fig. 2b) and $\langle k \rangle = 310.09$, $\gamma = 0.8$ for location information
 77 reported by contacted individuals before 23 January, 2020 (Supplementary Fig. 2c) and $\langle k \rangle =$
 78 142.46, $\gamma = 1.1$ for location information reported by contacted individuals after 23 January, 2020
 79 (Supplementary Fig. 2d). We can find that for both confirmed cases and contacted individuals, the
 80 mean value of crowdsourced reporting frequency decreased after 23 January, 2020, while the power

81 exponent increased after 23 January, 2020, which means people report less location information.
 82 The reason is that the travel restriction was implemented in Wuhan on 23 January, 2020 and people
 83 tended to stay at home.



Supplementary Figure 3: Average number of contacts of the entire population. Before and after the travel restriction of the general population daily contact number changes.

84 Supplementary Fig. 3 shows the average contact trends of the whole population (10,527,737
 85 smartphone users) in Wuhan during the period from 16 to 29 January, 2020. After the Chinese
 86 authority confirmed the COVID-19 coronavirus can be transmitted among human on 20 January,
 87 2020, the average daily contacts of the whole population dropped sharply. After the ban of non-
 88 essential vehicles in Wuhan downtown area on 26 January, 2020, the average contacts reached the
 89 lowest, at about 6.3, respectively. And then the average contacts stayed stable at such a level for a
 90 long period.

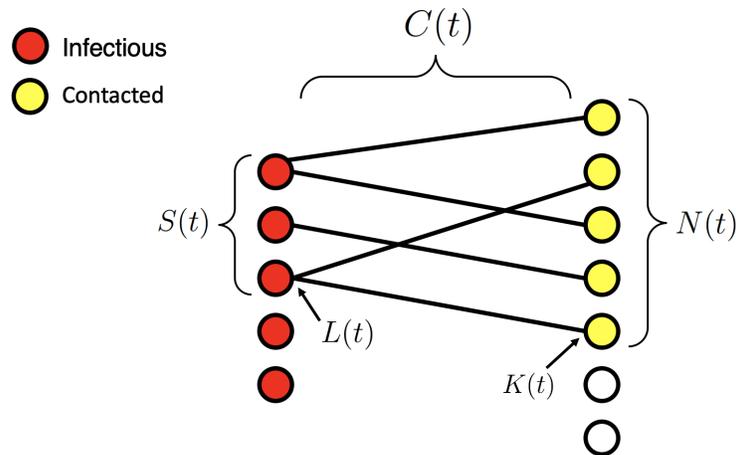


Supplementary Figure 4: Daily temperature in Wuhan from 1 January, 2020 to 31 January, 2020. The weather data comes from meteoblue.com.

91 **Supplementary Note II: Statistical analysis**

92 **Constructed graph structure**

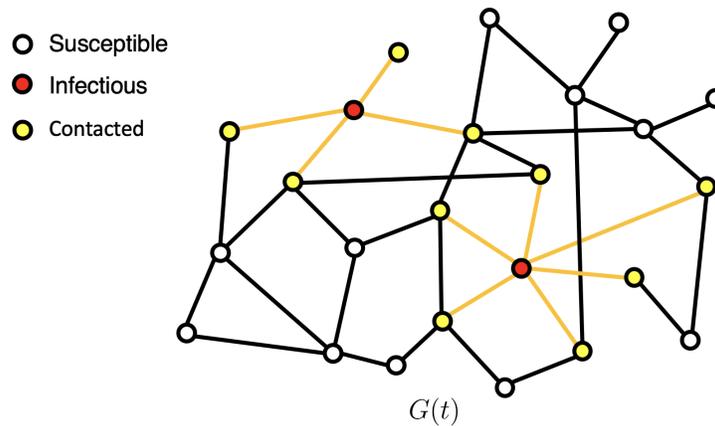
93 In the article we have introduced five indicators by leveraging the temporal contact graph:
 94 the total number of daily contacts $C(t)$ by calculating the number of edges of the temporal contact
 95 graph in day t , the total number of active infectious individuals $S(t)$, the number of susceptible
 96 individuals $N(t)$ who had at least one contact with infectious individuals, the average number
 97 of susceptible individuals $L(t)$ that each infectious individual contacted in day t (by calculating
 98 the average degree of the nodes representing infectious individuals in the graph), the number of
 99 infectious individuals $K(t)$ that each susceptible individual encountered in day t (by calculating
 100 the average degree of nodes representing the susceptible). A toy example is shown for specific
 101 calculation (Figure 5).



Supplementary Figure 5: A toy example for a snapshot (one day) of the temporal contact graph. Nodes on the left are infectious individuals and the nodes on the right are susceptible individuals. An edge indicates a contact between the two individuals in this day. The total contacts $C(t)$ is the number of edges in the bipartite graph, i.e., $C(t) = 5$. The number of active infectious individuals equals three, i.e., $S(t) = 3$, while the number of contacted individuals equals five, i.e., $N(t) = 5$. Also, we can calculate that $L(t) = 5/3$ and $K(t) = 1$.

102 Since the fraction of the infectious individuals in a population is very small, $L(t)$ is much
 103 larger than $K(t)$ generally. For example, in the contact graph $G(t)$, $K(t) = 1$ while $L(t) =$
 104 5 (Supplementary Figure 6). As a matter of fact, the two metrics are on very different orders

105 of magnitude, since the contact graph is much larger than the toy example. The units for $K(t)$
 106 and $L(t)$ are people/day. Thus, once the temporal contact graph is constructed by the contact
 107 model, the five informative indicators can be determined. However, different parameters do impact
 108 the measurements for the contact such as the infectious period and time interval in the contact
 109 model. We have done more analysis for the sensitivity of these parameters to support the results
 110 statistically. In the manuscript, we claim that $N(t)$ and $K(t)$ are similar rather than $N(t)$ and $S(t)$.
 111 Obviously, $N(t)$ and $K(t)$ increase from 1 January, 2020 to 20 January, 2020, and decrease from
 112 20 January, 2020 to 25 January, 2020.

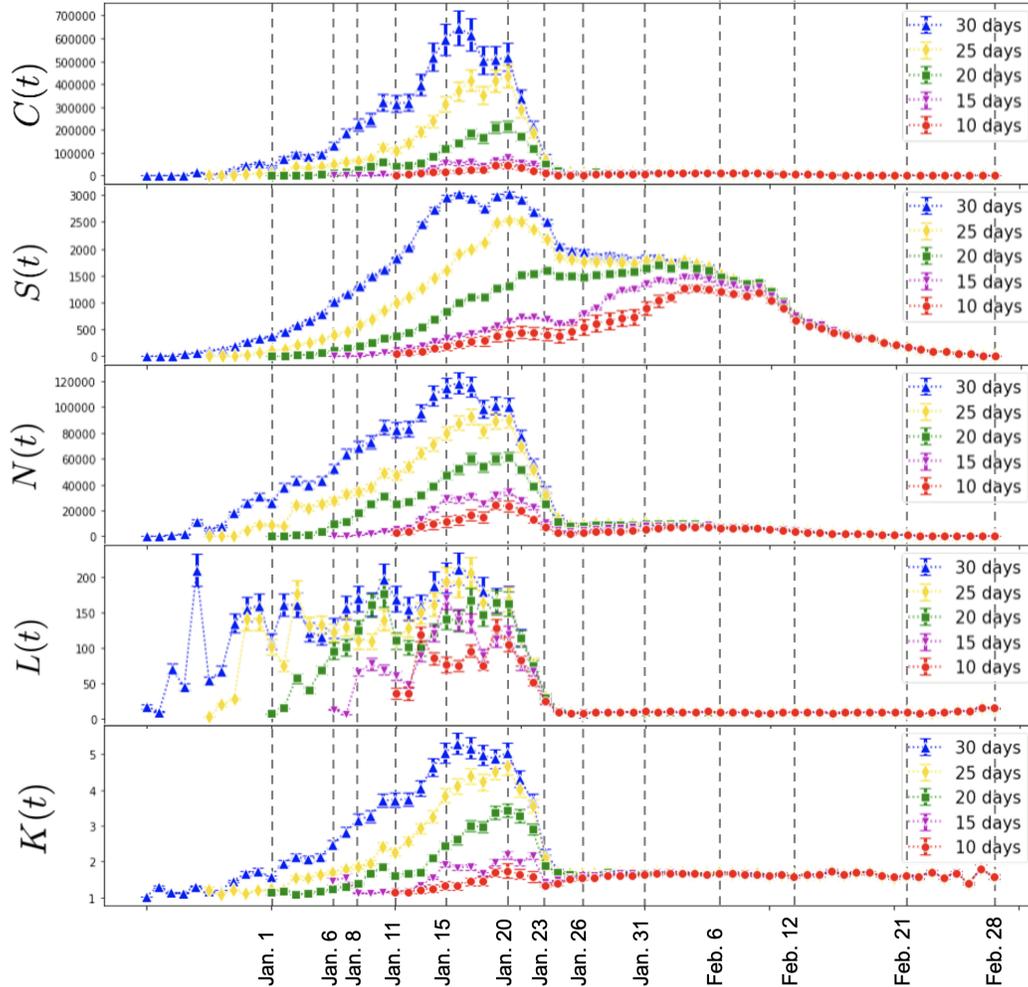


Supplementary Figure 6: A toy example for a snapshot of the spreading. The spreading in day t is represented by the graph $G(t)$, where nodes are individuals and edges are contacts between two individuals. We denote different health statuses by distinct colors, and the contacts between infectious and contacted individuals are colored by orange.

113 Since the infectious period is another parameter to determine the temporal contact graph,
 114 we perform sensitivity analyses for this period from 10 days to 30 days (Supplementary Figure
 115 7). Although the absolute values of the five indicators change, we can see that the trends remain
 116 stable.

117 **Infectious period**

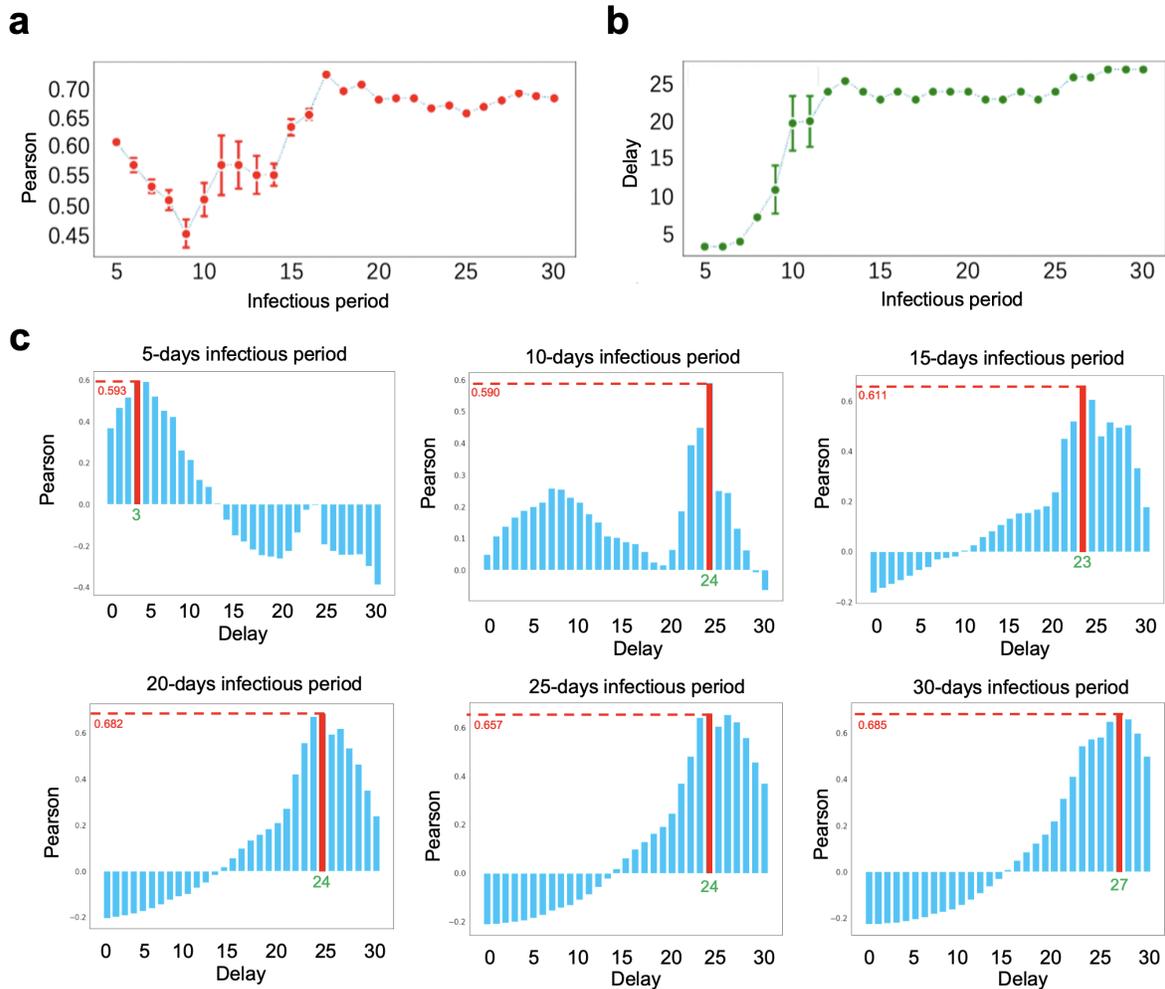
118 As mentioned in the Data description, the collected data include location-based informa-
 119 tion. Based on the collected data, we introduce a contact model, quantifying a contact between
 120 two individuals when they are geographically close to each other. Specifically, since locations of
 121 individuals is perturbed, we define a contact occurring when two individuals are in the same per-
 122 turbed area within a given time interval. Based on the proposed contact model, we constructed the



Supplementary Figure 7: Sensitivity analysis for the infectious period **a.** $C(t)$, the daily total times of contacts between infectious individuals and contacted individuals. **b.** $S(t)$, the daily total number of infectious individuals that have contacts with contacted individuals at least once. **c.** $N(t)$, the daily total number of contacted individuals that have contacts with infectious individuals at least once. **d.** $L(t)$, the daily average times of contacts with contacted individuals for infectious individuals before confirmation. **e.** $K(t)$, the daily average times of contacts with infectious individuals for contacted individuals.

123 temporal contact graph, which describes the daily contacts between the susceptible and infectious
 124 individuals: a node in the graph means a susceptible or infectious individual, and an edge in day t
 125 indicates a contact occurring in day t between a susceptible individual and infectious individual.

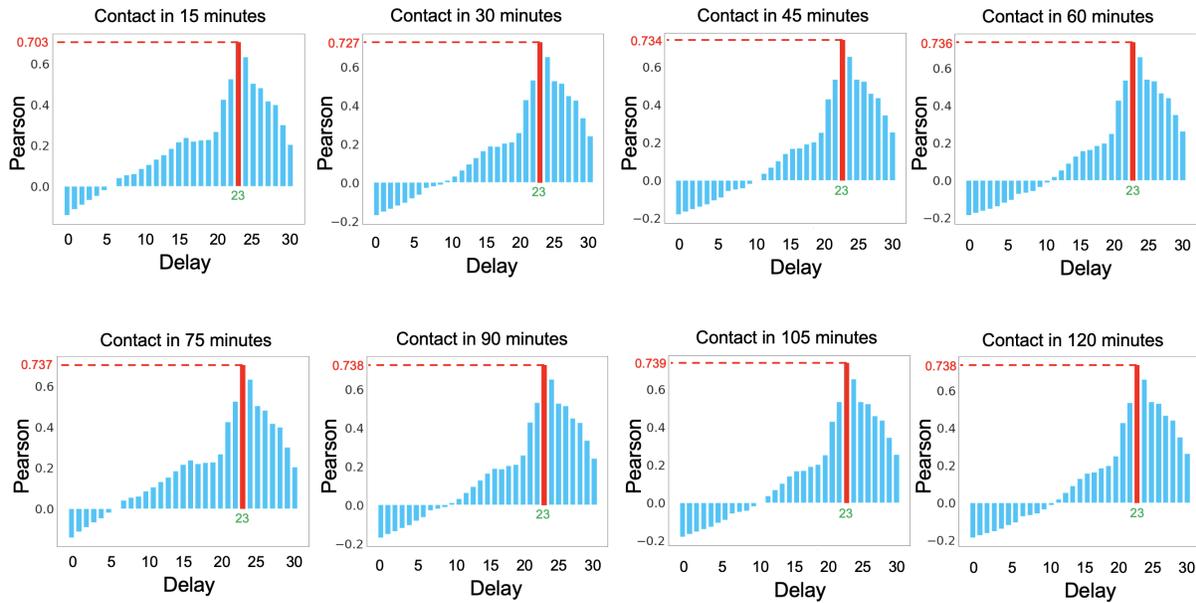
126 We calculate the Pearson correlations between daily contacts and infectious individuals with
 127 delays ranging from 5 to 30 days for different infectious period (Supplementary Figure 8). Ob-
 128 viously, the experiments show that a 23-day delay results in the best Pearson correlation of 0.78,
 129 where the corresponding infectious period equals 17 days in accordance with existing surveys⁴⁻¹¹.
 130 Thus, we consider the infectious period to be 17 days in the article, and an individual is identified
 131 as infectious in day t if he/she was confirmed during $[t + 1, t + 17]$.



Supplementary Figure 8: Sensitivity for the infectious period. We calculate the Pearson correlations between the daily contacts the daily confirmed cases with a delay ranging from 0 to 30 days.

132 We also performed the sensitivity analysis for the contact model by varying the time interval
 133 from 15 minutes to 120 minutes and the infectious period from 1 days to 30 days in the contact

134 model. Specifically, we vary the time interval from 15 minutes to 120 minutes and test the contact
 135 models under different time granularities, finding that such a time granularity does not change
 136 our conclusion in this article. We find that the Pearson coefficient between the daily number of
 137 contacts and delayed daily number of confirmed cases reaches maximum when the delay is 23
 138 days, corresponding to a 17-days infectious period for all values of time interval (Supplementary
 139 Figure 9). This indicates that the proposed contact model is very stable.

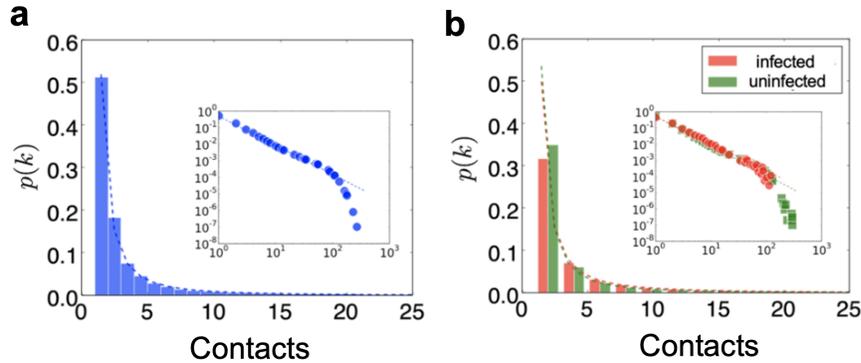


Supplementary Figure 9: Sensitivity for the time interval T . The Pearson correlations between daily contacts and daily confirmed cases with a delay ranging from 5 days to 30 days. Each panel corresponds to the time interval T for a contact in the contact model.

140 **KS test**

141 Since the distribution of various type of contacts follow a power-law, here we will discuss
 142 in details from a statistical stand point. In this article we have performed the maximum likelihood
 143 estimation for the parameters in the model to describe distributions of various types of contacts. In
 144 fact, we proposed a probability model based on Bayesian framework to evaluate the risk for every
 145 contact. Thus, we described the behaviors by power-law distributions, since we found there is a
 146 linear relation under the log-log coordinate. We have performed statistical analysis for our results
 147 by the KS test¹². We calculated the p -value for the distributions, where p -value is exactly the
 148 portion of synthetic sequences whose KS -distance is larger than that of the real data. For example,
 149 the exponent γ of the distribution of contacts for the whole population displays $\gamma = 1.96$, while

150 the p -value is 0.53. The exponent γ of that for the infected contacts displays $\gamma = 1.86$, while the
 151 p -value is 0.68. The exponent γ of that for the uninfected susceptible displays $\gamma = 2.01$, while the
 152 p -value is 0.52 (Supplementary Figure 10).

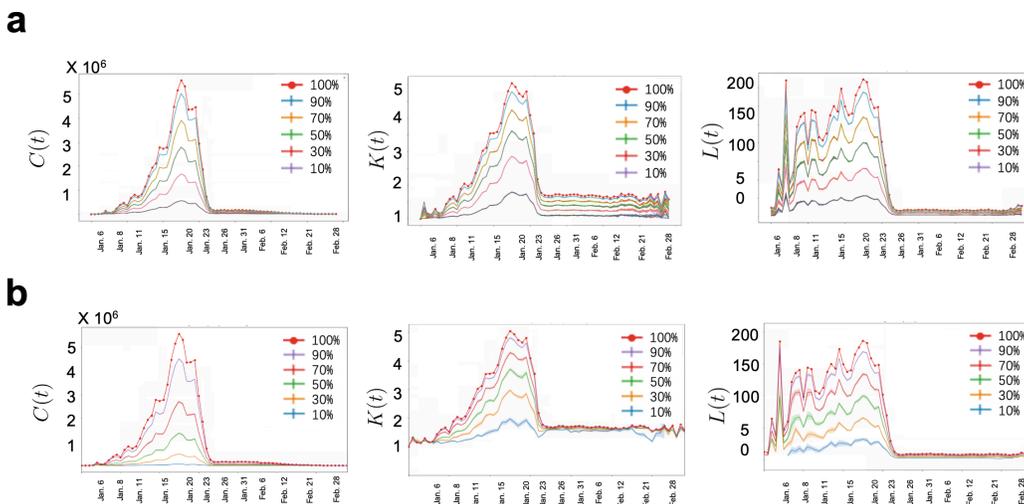


Supplementary Figure 10: Distributions of behaviors for infected and uninfected individuals.
a. The distribution of contacts for all contacted individuals. **b.** The distribution of contacts for all infected individuals (red bars) and uninfected individuals (green bars).

153 User involvement

154 To perform sensitivity analysis for the impacts of users' involvement, we further conduct ex-
 155 periments to simulate different user involvements by randomly selecting $\alpha\%$ users as the voluntary
 156 users, and $\alpha\%$ data items each user uploading per day, and evaluate the corresponding performance
 157 loss. To show the robustness of our analysis, we repeat 10 times Monte Carlo experiments for a
 158 given setting. Specifically, we set user involvement rate to 10%, 30%, 50%, 70%, 90%, and user
 159 uploading rate to 10%, 30%, 50%, 70%, 90%. Then, we perform the corresponding experiments.
 160 We plot the curves of $K(t)$, $L(t)$ and $C(t)$ calculated from the contact model with different up-
 161 loading rates from 1 January to 28 February with error bars (Supplementary Figure 11a). The
 162 results show that with the decrease of uploading rates, the values of three parameters decrease
 163 obviously, which is because the number of recorded contacts drops along with decreasing upload
 164 rates. Notice that the error bars of three parameters are all small, which implies the Monte Carlo
 165 experiments are stable to produce similar results. It is, therefore, reasonable to emphasize that the
 166 number of user uploads reduced for each Monte Carlo experiment is the same, and the total num-
 167 ber of participation does not reduce under user upload analysis. Generally speaking, the Pearson
 168 correlation between case of 10% uploading rate and other cases ($\alpha\%$ upload rates) remains high
 169 for the analysis of user uploading rate. This phenomenon implies that we might not need a high
 170 uploading rate if we are only interested in estimating the trends.

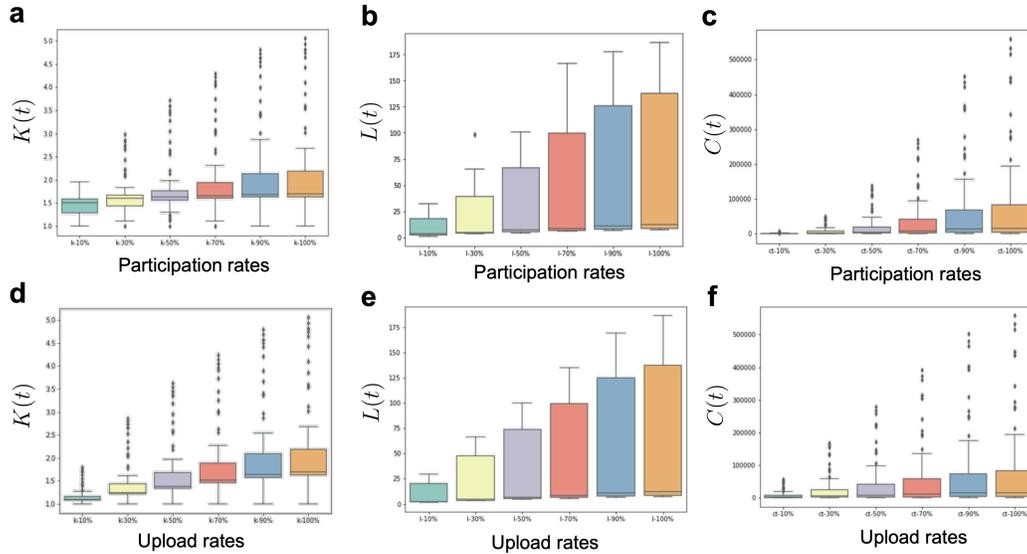
171 We conduct experiments in the same way as those for uploading rates to evaluate the impact
 172 of diverse participation rates (Supplementary Figure 11b). The corresponding results show that,
 173 compared with uploading rate analysis, decreasing the participation rates brings more uncertainty
 174 in the Monte Carlo experiments with higher error bars. This is because the participants that are
 175 randomly deleted at each time experiment are different. Thus, their impacts on the entire network
 176 also vary. Notice that when the participation rate is not very small, the high correlation can still be
 177 preserved. However, when it is reduced to 10%, the correlation coefficient decreases significantly.
 178 This could be attributed to the power law distribution of the network: since the distribution has
 179 an obvious long-tail effect, only when the participation rate is low enough can some key nodes be
 180 deleted, thereby affecting the trend of the proposed metrics.



Supplementary Figure 11: Sensitivity analysis for user involvement. **a.** The proposed metrics including $C(t)$, $K(t)$ and $L(t)$ are shown, when simulating the uploading rate ranging from 10% to 90%. **b.** The proposed metrics including $C(t)$, $K(t)$ and $L(t)$ are shown, when simulating the participation rate ranging from 10% to 90%.

181 We have analyzed the contact model in main body of the paper with six curves showing
 182 different user participation rates and upload rates. Here we give the corresponding statistical infor-
 183 mation such as median and variances of $K(t)$, $L(t)$ and total contacts $C(t)$ shown in Supplementary
 184 Figure 12. It reveals that as $\alpha\%$ of participation and upload rate decreases, the statistical informa-
 185 tion decreases with the similar trend. This is expected as reduction in either user participation rate
 186 or user upload rate decreases the chances of having contacts among users.

187 Moreover, the practical challenges for real contact tracing apps are: what to do if we have



Supplementary Figure 12: Performance of different user involvement in the contact model.

a-c. Three box plots show the distribution change of $K(t)$, $L(t)$, and daily number of total contacts Vs. different user participation rates. **d-f.** Three box plots show the distribution change of $K(t)$, $L(t)$, and daily number of total contacts Vs. different user upload rates.

188 only partial information about contacts and about the real individual health status? As a matter
 189 of fact, we have investigated this in this article, indicating the impacts of user involvement on
 190 the contact tracing app performance. We have conducted experiments to simulate different user
 191 involvements by randomly selecting $\alpha\%$ users as the voluntary users, and evaluate the correspond-
 192 ing performance loss, where $\alpha\% = 10\%, 30\%, 50\%, 70\%, 90\%$. In such a way, both susceptible
 193 and infectious individuals could be removed, which simulates partial information case about both
 194 susceptible and infectious individuals. One interesting result in our work is even though we only
 195 have partial information about the contacts and the conformed cases, we can still have a good
 196 performance on estimating the evolving situation of COVID-19 when $\alpha\%$ is not low enough.

197 **Supplementary Note III: Data and real contagion**

198 **Data biases**

199 In this article, we analyze the evolving epidemic situation of COVID-19, leveraging the
200 available reported ‘confirmed cases’ as a proxy of the measure of the extent of the contagion
201 However, many other factors (other than actual COVID19 transmission) may have influenced the
202 recorded number of confirmed cases because of medical resource lacking, incomplete information
203 and so on. We will consider how could these biases in the number of confirmed cases (compared
204 to the number of actual cases) shape the actual results. In fact, the reported number of confirmed
205 cases is typically less than the number of actual cases. Therefore, there is always biases in reality.
206 However, it is impossible to know the actual cases even their distributions. Here, we assume that
207 the reported confirmed cases $Sc(t)$ is proportional to the number of actual cases in reality $Sr(t)$,
208 i.e.,

$$Sc(t) = \beta \cdot Sr(t), \quad (1)$$

209 where $0 < \beta < 1$ is a constant to quantify the probability for an actual case being confirmed.
210 Then, in this article we perform correlation analysis between the number of daily contacts $C(t)$
211 and daily confirmed cases $Sc(t)$, and obtain a Pearson coefficient $\rho(C(t), Sc(t)) = 0.78$. Then, we
212 will prove the Pearson coefficient between the daily contacts $C(t)$ and daily infected individuals
213 in reality $Sr(t)$ is the same, i.e.,

$$\rho(C(t), Sr(t)) = \rho(C(t), \beta \cdot Sc(t)) = \rho(C(t), Sc(t)) = 0.78. \quad (2)$$

214 In other words, although $Sc(t)$ has a bias from $Sr(t)$, i.e., $Sc(t) = \beta \cdot Sr(t)$, it does not affect the
215 correlation analysis in our results.

216 **Recorded contacts and transmission**

217 A primary conceptual focus of the manuscript is how real-world social contacts may relate
218 to COVID-19 contagion patterns. For instance, we define a contact as a co-occurrence within a
219 specific distance (e.g., $15 \times 30m^2$) in the proposed contact model, when in fact it is unlikely that this
220 virus would be transmitted at a range even 10 fold less than this. Therefore, it is necessary to clarify
221 how this might influence the patterns found, or the conclusions drawn. Medically speaking, a close
222 contact is said to occur when two individuals are within a distance of 1.8 meters¹³. Individuals
223 who have close contacts with infectious cases generally have a high probability of getting infected.

224 However, in practice, it is difficult to decide a close contact in a digital way. Most contact tracing
 225 apps exploit Bluetooth and/or GPS to decide a contact when two individuals are in a short distance
 226 (e.g., within 20 meters). In fact, one of the current controversial issues lies in whether such type
 227 of contacts captured by contact tracing apps is effective since it is not fine-grained enough. For
 228 example, the mentioned potential interesting work discussed the importance of network structure
 229 and social dynamics in evaluating the potential impact of SARS-CoV-2 control by combing fine-
 230 scale data¹⁴. However, our results show that there is still a clear distinction of contact behaviors
 231 between the infected and uninfected contacted individuals under our contact model. Therefore, the
 232 frequency of social contacts captured by contact tracing apps actually has a high correlation with
 233 contagion patterns.

234 Risk evaluation

235 In this article we have considered the risk by their behaviors and features, and we use
 236 true/false positive and the ROC curve to analyze the effectiveness of the risk model. Notice that
 237 the “positive” in the phrase “true/false positive rates” does not indicate the “positive” in a nucleic
 238 acid testing. In fact, we have measured the risk of every contact j , i.e., $P(z_j = 1)$ by the proposed
 239 risk model in the Result III. In order to evaluate the risk measured by the model, we study the
 240 ROC (receiver operating characteristic) curve. For a threshold $0 < q < 1$, specifically, a contact
 241 j is considered to be true positive if $z_j = 1$ and $P(z_j = 1) > q$, while j is considered to be false
 242 positive if $z_j = 0$ and $P(z_j = 1) > q$. Then, we can calculate the TPR (true positive rate) by

$$TPR = \frac{\sum_j \mathbf{1}_{\{z_j=1, P(z_j=1)>q\}}}{\sum_j \mathbf{1}}, \quad (3)$$

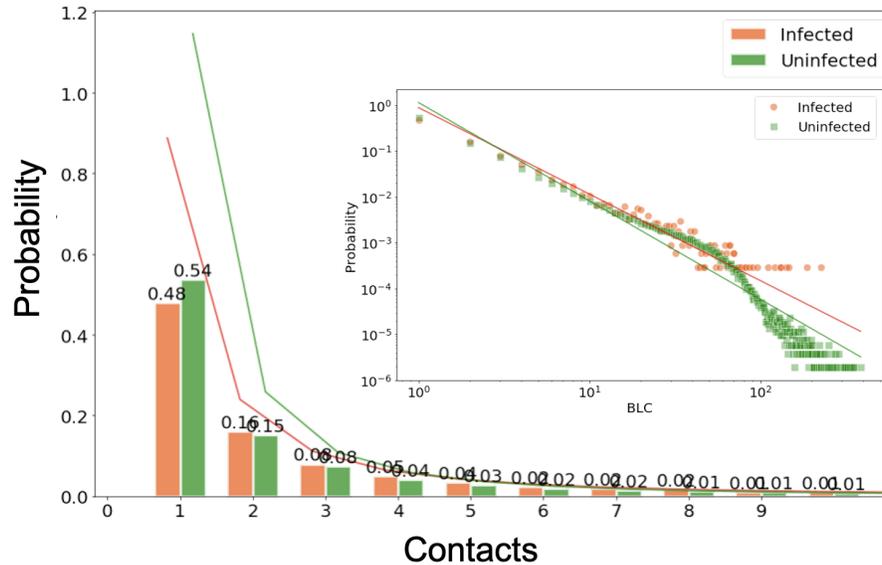
243 and FPR (false positive rate) by

$$FPR = \frac{\sum_j \mathbf{1}_{\{z_j=0, P(z_j=1)>q\}}}{\sum_j \mathbf{1}}. \quad (4)$$

244 Thus, the ROC curve described by TPR and FPR can well evaluate the risk model for contacted
 245 individuals.

246 We have mentioned that there is a distinction of contact behaviors between the infected and
 247 uninfected contacted individuals in the article, where the distinction indicates the difference of
 248 the probability distribution for contacts between the infected contacts and the uninfected contacts.
 249 Here, we calculated the KL -divergence for the distributions, quantified the difference and showed
 250 if the distinction is prominent. Specifically, we calculated the two distributions for the infected

251 contacts and the uninfected contacts (Supplementary Figure 13). We also calculated the KL -
 252 divergence between the two distributions, where the KL -divergence equals to 0.02. These results
 253 all show that there is a slightly distinction of contact behaviors between the infected and uninfected
 254 contacted individuals.



Supplementary Figure 13: Distinction between infected and uninfected individuals. The distributions of contacts for infected individuals (red bars) and uninfected individuals (green bars)

255 Tracing apps

256 In this article, we use the location-related data from tracing apps to calculate if two indi-
 257 viduals are close. Based on this, we build the temporal contact graph, which shows the contact
 258 relationship between susceptible and infectious individuals. Although the data used in our article
 259 is different with the Bluetooth based contact tracing apps, the principles used to define the potential
 260 contact are the same. By studying co-location data for 10 millions mobile phone users in Wuhan,
 261 our aim is to understand and clarify the potential of contact tracing apps to identify and interrupt
 262 transmission chains of the SARS-CoV-2 virus. However, controversies on the contact tracing apps
 263 mainly include: 1) the contact tracing apps may have privacy leakage issues and the methods (cen-
 264 tralized V.S. decentralized) used to inform the potential contacts with infectious cases about their
 265 risks need further investigation; 2) current contact tracing apps utilize location-related informa-
 266 tion (e.g., Bluetooth, WiFi or GPS) to define a contact and are not fine-grained enough to capture
 267 a close contact within a distance of 1.5 meters; 3) such apps may not work for suppressing the
 268 transmission of COVID-19 when the participation rate is not high enough. For the first concern,

269 Google recently showed that the privacy can be protected by using Bluetooth with a proper privacy-
270 preserving protocol. Our work here studies the contact behavior analysis in the transmission, and
271 thereby does not discuss the privacy issues. Note that our approach falls within the category of
272 the centralized way of informing the potential contacts. Though both centralized and decentralized
273 ways can identify individuals having contacts with infectious cases, our results demonstrate that
274 centralized way can provide an abundance of information that can be helpful for prevention and
275 control of COVID-19 (see explanations for the second and third issues as follows). For the second
276 issue, it is true that contact tracing apps cannot accurately define a close contact and may ignore
277 other important factors that impact the spread of COVID-19, e.g., whether a protection measure
278 such as wearing a mask is taken in a contact. However, our results show that there is still a promi-
279 nent distinction of contact behaviors between the infected and uninfected contacted individuals.
280 Based on this, we designed an infection risk evaluation framework to identify potential infected
281 ones. For the third issue, we evaluated the effect of user involvement and show that user partici-
282 pation rate exerts higher influence on situation evaluation than user upload rate does. Moreover,
283 our results indicate that the contact tracing apps can still be helpful even when user involvement is
284 low. Also, we find that five indicators calculated from the constructed temporal contact graph are
285 informative to understand the actual transmission and evaluate the epidemic situation. In summary,
286 though this article cannot solve all the controversial issues, we provide new evidence that contact
287 tracing apps can be very helpful to the prevention and control of COVID-19.

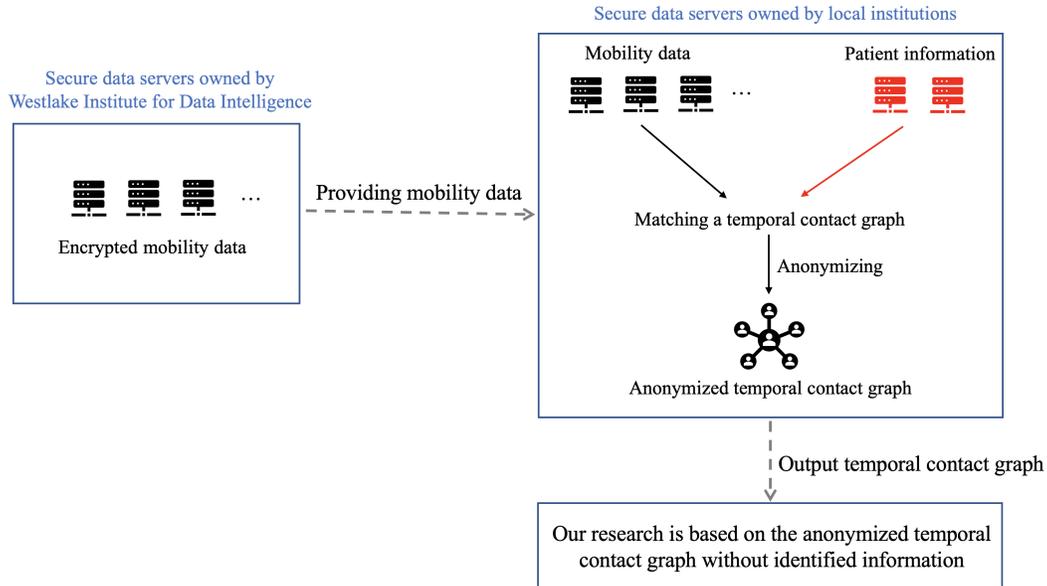
288 Currently, many other papers in recent times (the company Cuebiq for example provided a
289 lot of similar data to many researchers in Europe and USA), used data collected from GPS, POI
290 as well as wifi. In fact, most contact tracing apps exploit Bluetooth and/or GPS on smartphones
291 to discover nearby devices held by users and identify the contacts between the susceptible and
292 infectious individuals. For example, most contact tracing apps defined a contact by a Bluetooth
293 “handshake”, which occurs for two individuals within a distance of 20 meters. Some other apps
294 used additional data such as GPS for contact tracing, e.g., Norway released an app that collects
295 both GPS and Bluetooth information. In this article, a contact is identified for two individuals
296 within about 20 meters. Clearly, our contact model works in a very similar way to those defined
297 by most of the contact tracing apps, though the data is not directly collected from contact tracing
298 apps (which is not available). To summarize, although the data used in our paper is different with
299 the Bluetooth contact tracing apps, the principles used to define the potential contact are just the
300 same.

301 **Ethic and privacy issues**

302 Since many ethic issues have been raised, it is of significance to discuss the guidelines for
303 COVID-19 tracing apps¹⁵. Local institutions for disease prevention and control have the authority
304 to collect the information on COVID-19 (including the patient information) and they also have the
305 obligation to protect the privacy of patients when processing these data according to the Law on
306 the Prevention and Control of Infectious Diseases of the People’s Republic of China. Specifically,
307 Article 12 of the Law on the Prevention and Control of Infectious Diseases of the People’s Repub-
308 lic of China stipulates that “All units and individuals within the territory of the People’s Republic
309 of China shall **accept the preventive and control measures taken by disease prevention and**
310 **control institutions and medical agencies for investigation, testing, collection of samples of**
311 **infectious diseases and for isolated treatment of such diseases, and they shall provide truthful**
312 **information about the diseases.** Disease prevention and control institutions and medical agencies
313 shall not divulge any information or materials relating to personal privacy”. Therefore, since the
314 outbreak of COVID-19, local institutions have collected a large set of data. To protect the privacy
315 of patients, these data are stored on the secure data servers owned by local institutions.

316 In our project, we identified the confirmed cases based on the phone number information.
317 However, we hereby clarify that we (including authors from Westlake Institute for Data Intelli-
318 gence) did not obtain any information regarding the exact phone numbers of patients. As men-
319 tioned in the first paragraph, these data are stored on the secure data servers and are not available
320 for us. To facilitate our research, Westlake Institute for Data Intelligence uploaded the mobility
321 data contributed by smartphone users and their account information in Wuhan from Jan. 1 to Feb.
322 28 to the secure data servers owned by local institutions. It is noteworthy that the phone number
323 information here refers to the encrypted (i.e., hashed) information from the original phone number.
324 To protect users from privacy leakage and malicious attack, their phone numbers were mapped
325 to hash values. This procedure was irreversible and would be completed once the account was
326 created. This is actually a popular approach adopted by many location-based service providers.
327 The phone numbers of patients were also mapped to hash values by the same hash function on the
328 secure data servers by local institutions. Since a phone number yielded a unique hash value, we
329 could leverage such a connection to identify the patients. The temporal contact graph encoding the
330 contact among the infectious and the susceptible individuals was then constructed.

331 A schematic diagram of the above procedure is shown in Figure 14. The whole process of
332 constructing the temporal contact graph was conducted on the secure data servers of the local in-
333 stitutions. In the temporal contact graph, each node (a user) is further anonymized and thus cannot
334 be traced back to his/her phone number information (i.e., hashed values) anymore. This means



Supplementary Figure 14: The procedure of constructing a temporal contact graph.

335 that we were unable to gain access to any identifiable phone number or hash value about each node
 336 in the temporal contact graph even though some nodes were identified as confirmed cases. More-
 337 over, the mobility data uploaded by Westlake Institute for Data Intelligence were destroyed on the
 338 secure data server, and the temporal contact graph was provided for us offline only for research
 339 purpose. Since Westlake Institute for Data Intelligence does not own the phone numbers of con-
 340 firmed cases and the temporal contact graph does not contain any private information, we do not
 341 have the authority issue. We agree that the phone number is a unique identifier. However, it was
 342 well protected in our research, since the phone numbers of confirmed cases would be hashed before
 343 they were used on the secure data servers. Furthermore, such information was securely stored on
 344 the data servers of local institutions and it was not available for us. We merely utilized the temporal
 345 contact graph to obtain relevant results in our current research. This is what we meant by saying
 346 “did not include identified individual-level data”. As mentioned above, both the personal phone
 347 number data from Westlake Institute for Data Intelligence and those from the local institutions
 348 were preprocessed by the same pseudonymization mechanism (i.e., hash function). Besides, the
 349 final output (i.e., the temporal contact graph) was further pseudonymized to avoid data restoration.
 350 Thus, the hashed phone numbers do not allow individuals to be identified as covid-19 patients or
 351 individually identified. In a word, there is no private identifying information about the individuals
 352 accessible to us (researchers) and no interaction (or intervention) between the individuals and us.
 353 In addition, we note that the whole procedure is amenable to laws in both China and EU as follows.
 354 Article 1038 of the Civil Code of the People’s Republic of China, stipulates that “Information pro-

355 cessors shall not disclose or tamper with the personal information they collect or store; and without
356 the consent of the data subject, information processors shall not illegally provide any other person
357 with his personal information, **except for the processed information cannot be identified with**
358 **any specific person and restored.**” Article 89 of the General Data Protection Regulation (GDPR)
359 enforced by the European Union stipulates that “Processing for archiving purposes in the public
360 interest, scientific or historical research purposes or statistical purposes, shall be subject to ap-
361 propriate safeguards, in accordance with this Regulation, for the rights and freedoms of the data
362 subject. Those safeguards shall ensure that technical and organizational measures are in place
363 in particular in order to ensure respect for the principle of data minimization. **Those measures**
364 **may include pseudonymization provided that those purposes can be fulfilled in that man-**
365 **ner. Where those purposes can be fulfilled by further processing which does not permit or**
366 **no longer permits the identification of data subjects, those purposes shall be fulfilled in that**
367 **manner.”**

368 Another key point in the debate on the contact tracing apps is about the concerns that the
369 contact tracing apps may have privacy leakage issues and the methods (centralized V.S. decen-
370 tralized) used to inform the potential contacts with infectious cases about their risk need further
371 investigation. For this issue, Google recently showed that the privacy can be protected by using
372 Bluetooth with a proper privacy-preserving protocol¹⁶. Also, there are some other works on this
373 direction¹⁷⁻¹⁹. Our work here studies the contact analysis in the transmission, and thereby does not
374 discuss the privacy issues. Note that our approach falls within the category of the centralized way
375 of informing the potential contacts. Though both centralized and decentralized ways can identify
376 individuals having contacts with infectious cases, our results demonstrate that centralized way can
377 provide an abundance of information that can be helpful for prevention and control of COVID-
378 19. As indicated by the title of our article, our main focus is to reveal the evolving situation of
379 COVID-19 by exploiting the temporal contact graph. The study of the effects that an app will have
380 on people’s behaviour is another important topic. There are also some works on this direction^{20,21}.
381 This topic is a new and independent topic and therefore was not considered in this work.

382 **Supplementary References**

- 384 1. Kraemer, M. U. *et al.* The effect of human mobility and control measures on the covid-19
385 epidemic in china. *Science* **368**, 493–497 (2020).
- 386 2. Jia, J. S. *et al.* Population flow drives spatio-temporal distribution of COVID-19 in china.
387 *Nature* (to appear).
- 388 3. Aleta, A. *et al.* Modeling the impact of social distancing, testing, contact tracing and household
389 quarantine on second-wave scenarios of the covid-19 epidemic. *medRxiv* (2020).
- 390 4. Lauer, S. A. *et al.* The incubation period of coronavirus disease 2019 (covid-19) from publicly
391 reported confirmed cases: estimation and application. *Annals of internal medicine* **172**, 577–
392 582 (2020).
- 393 5. Sohrabi, C. *et al.* World health organization declares global emergency: A review of the 2019
394 novel coronavirus (covid-19). *International Journal of Surgery* (2020).
- 395 6. Li, Q. *et al.* Early transmission dynamics in wuhan, china, of novel coronavirus–infected
396 pneumonia. *New England Journal of Medicine* (2020).
- 397 7. Linton, N. M. *et al.* Incubation period and other epidemiological characteristics of 2019 novel
398 coronavirus infections with right truncation: a statistical analysis of publicly available case
399 data. *Journal of clinical medicine* **9**, 538 (2020).
- 400 8. Bi, Q. *et al.* Epidemiology and transmission of covid-19 in 391 cases and 1286 of their close
401 contacts in shenzhen, china: a retrospective cohort study. *The Lancet Infectious Diseases*
402 (2020).
- 403 9. Cao, M. *et al.* Clinical features of patients infected with the 2019 novel coronavirus (covid-19)
404 in shanghai, china. *MedRxiv* (2020).
- 405 10. Chen, J. *et al.* Clinical progression of patients with covid-19 in shanghai, china. *Journal of*
406 *Infection* (2020).
- 407 11. Cheng, Y. *et al.* Kidney disease is associated with in-hospital death of patients with covid-19.
408 *Kidney international* (2020).
- 409 12. Clauset, A., Shalizi, C. R. & Newman, M. E. Power-law distributions in empirical data. *SIAM*
410 *review* **51**, 661–703 (2009).

- 411 13. Centers for Disease Control and Prevention. Contact tracing for covid-19.
412 [https://www.cdc.gov/coronavirus/2019-ncov/php/contact-tracing/
413 contact-tracing-plan/contact-tracing](https://www.cdc.gov/coronavirus/2019-ncov/php/contact-tracing/contact-tracing-plan/contact-tracing) (2020).
- 414 14. Firth, J. A. *et al.* Combining fine-scale social contact data with epidemic modelling reveals in-
415 teractions between contact tracing, quarantine, testing and physical distancing for controlling
416 covid-19. *medRxiv* (2020).
- 417 15. Morley, J., Cowls, J., Taddeo, M. & Floridi, L. Ethical guidelines for covid-19 tracing apps
418 (2020).
- 419 16. Gvili, Y. Security analysis of the covid-19 contact tracing specifications by apple inc. and
420 google inc. *IACR Cryptol. ePrint Arch.* **2020**, 428 (2020).
- 421 17. Singapore Government. Tracetogether, safer together. [https://www.tracetogether.
422 gov.sg/](https://www.tracetogether.gov.sg/) (2020).
- 423 18. Australia Government Department of Health. Covidsafe app. [https://www.health.
424 gov.au/resources/apps-and-tools/covidsafe-app](https://www.health.gov.au/resources/apps-and-tools/covidsafe-app) (2020).
- 425 19. National Cyber Security Centre. Nhs covid-19: the new contact-tracing app from the nhs.
426 [https://www.ncsc.gov.uk/information/nhs-covid-19-app-explainer
427](https://www.ncsc.gov.uk/information/nhs-covid-19-app-explainer) (2020).
- 428 20. There are many reasons why COVID-19 contact-tracing apps may not work. Adam vaughan.
429 <https://www.newscientist.com/article/2241041/> (2020).
- 430 21. Hinch, R. *et al.* Effective configurations of a digital contact tracing app: A report to
431 nhsx. *en. In:(Apr. 2020). Available here. url: https://github.com/BDI-pathogens/covid-
432 19_instant_tracing/blob/master/Report* (2020).