

# Uncovering COVID-19 Transmission Tree: Identifying Traced and Untraced Infections in an Infection Network

Hyunwoo Lee<sup>1,4</sup>, Hayoung Choi<sup>1,4,\*</sup>, Hyojung Lee<sup>2,4</sup>, Sunmi Lee<sup>3,4</sup> and Changhoon Kim<sup>5,6</sup>

<sup>1</sup>Department of Mathematics, Kyungpook National University, Daegu, Republic of Korea, <sup>2</sup>Department of Statistics, Kyungpook National University, Daegu, Republic of Korea, <sup>3</sup> Department of Applied Mathematics, Kyunghee University, Yongin-si, Republic of Korea, <sup>4</sup> Nonlinear Dynamics & Mathematical Application Center, Kyungpook National University, Daegu, Republic of Korea, <sup>5</sup> Department of Preventive Medicine, College of Medicine, Pusan National University, Busan, Republic of Korea, <sup>6</sup> Busan Center for Infectious Disease Control and Prevention, Pusan National University Hospital, Busan, Republic of Korea

Correspondence\*:  
Hayoung Choi  
hayoung.choi@knu.ac.kr

## 2 ABSTRACT

3 We present a comprehensive analysis of COVID-19 transmission dynamics using an infection  
4 network derived from epidemiological data in South Korea, covering the period from January  
5 3, 2020, to July 11, 2021. This network, illustrating infector-infectee relationships, provides  
6 invaluable insights for managing and mitigating the spread of the disease. However, significant  
7 missing data hinder the conventional analysis of such networks from epidemiological surveillance.  
8 To address this challenge, our research suggests a novel approach for categorizing individuals  
9 into four distinct groups, based on the classification of their infector or infectee status as either  
10 traced or untraced cases among all confirmed cases. Furthermore, the study analyzes the  
11 changes in the infection networks among untraced and traced cases across five distinct periods.  
12 The four types of cases emphasize the impact of various factors, such as the implementation of  
13 public health strategies and the emergence of novel COVID-19 variants, which contribute to the  
14 propagation of COVID-19 transmission. One of the key findings of this study is the identification  
15 of notable transmission patterns in specific age groups, particularly in those aged 20–29, 40–69,  
16 and 0–9, based on the four type classifications. Moreover, we develop a novel real-time indicator  
17 to assess the potential for infectious disease transmission more effectively. By analyzing the  
18 lengths of connected components, this indicator facilitates improved predictions and enables  
19 policymakers to proactively respond, thereby helping to mitigate the effects of the pandemic on  
20 global communities.

21 **Keywords:** COVID-19, infection network, contact tracing, reproduction number, untraced infection

## 1 INTRODUCTION

22 COVID-19, caused by the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), was declared  
23 a pandemic by the World Health Organization on March 11, 2020. According to the World Health  
24 Organization's weekly epidemiological update released on February 2, 2021, the epidemic of COVID-19  
25 spread rapidly to more than 200 countries. Without effective control measures, the rapidly increasing  
26 number of COVID-19 cases will greatly increase the burden of clinical treatments. This situation may lead  
27 to a critical shortage of healthcare system capacity for severe cases, ultimately resulting in a sharp and  
28 alarming increase in mortality rates. Consequently, various control measures were implemented, leading  
29 to observed fluctuations in the efficacy of strategies like contact tracing and isolation of confirmed cases

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

30 throughout the pandemic (1). South Korea, first reporting its COVID-19 case on January 19, 2020 (2, 3), has  
31 experienced multiple waves of outbreaks, in response to which it actively implemented control measures  
32 such as social distancing, mask-wearing, lockdowns, and enhanced efforts in testing and contact tracing.  
33 Especially, active contact tracing has generated significant epidemiological data, enabling analysis of  
34 extensive infection networks (4). Understanding the infection network for COVID-19 is crucial for several  
35 reasons. First and foremost, it allows us to grasp the dynamics of the virus's transmission within a  
36 population (5). By mapping out how individuals infect each other, we gain valuable insights into the  
37 patterns and pathways through which the virus spreads (1). Additionally, studying the infection network  
38 aids in the identification of key factors influencing the transmission (2). This includes factors such as  
39 age-specific patterns, which can help tailor public health measures to specific demographics, ultimately  
40 improving the effectiveness of containment strategies (6).

41 Previous research focused on cluster analysis, reproduction number, and network analysis to address key  
42 transmission factors and assess the effectiveness of various interventions during COVID-19 pandemic (3, 6,  
43 7, 8, 9, 10, 11, 12). In (7, 8) authors investigated COVID-19 transmission by age group, aiding in identifying  
44 the primary age groups fueling the spread and formulating age-specific response strategies. It scrutinized  
45 the infection spread by clusters, offering insights into evaluating social distancing measures outlined in  
46 (3, 6, 9). Examining cluster type frequency in both the initial and subsequent epidemic waves enables the  
47 development of an effective strategy for controlling outbreaks (3). Network analysis facilitates assessing  
48 specific vertices' importance and understanding the relationships between them (2, 5, 13). Furthermore,  
49 Wang et al.(10) and Zhang et al.(11) investigated the basic reproduction number  $\mathcal{R}_0$  of COVID-19, which  
50 represents the transmission potential of an infectious disease in the early phase of an epidemic (12). The  
51 time-dependent reproduction number  $\mathcal{R}_t$  represents the instantaneous reproduction number, indicating the  
52 expected number of secondary infections caused by an infector at a specific point in time (12).

53 In the context of COVID-19 policies, our current knowledge of how infections spread through  
54 transmission networks is primarily based on virtual data and theoretical models (14, 15), with evidence  
55 from actual data (16, 17, 18) being limitedly available. The infection network generated from actual  
56 epidemiological data contains numerous missing data, resulting in many connected components, creating a  
57 disparity from analyses based on virtual data. Contact tracing is commonly recommended for controlling  
58 COVID-19 outbreaks, yet its effectiveness is unclear. Studies evaluating the effectiveness of contact  
59 tracing are categorized into observational studies (19, 20, 21, 22) and modeling studies (1, 23, 24, 25). Our  
60 study suggests that analyzing the classification of four types of confirmed cases in the infection network,  
61 along with the distribution of connected component lengths, can broaden insights into contact tracing and  
62 dynamics of disease transmission. A pivotal study analyzing changes in the infection pattern structure  
63 between infectors and infectees based on age groups (26) is also essential. Surprisingly, there has been  
64 no previous study on this specific topic for COVID-19 infection between infectors and infectees in South  
65 Korea.

66 This paper is motivated by the recognition of differences in infection networks generated from actual  
67 data versus virtual data. This research has established an infection network by assigning an infector to  
68 all infectees from actual epidemiological data KDCA (27) from January 3, 2020, to July 11, 2021, in  
69 South Korea. It is shown that the established infection network comprises many connected components  
70 due to missing vertices (individuals) and edges (infection events). Consequently, we proposed a method of  
71 categorizing individuals as either (i) infectors, who are aware of the infectees they have transmitted the  
72 virus to, or (ii) infectees, who are cognizant of their infector. This method allows for the categorization of  
73 vertices in the numerous distinct connected components from a common perspective and facilitates the  
74 derivation of analysis for each vertex. Furthermore, several properties were established from the method.  
75 This paper analyzed the infection network in terms of time and age groups using a four-type categorization  
76 method and proposes a new real-time calculated indicator of infectious disease transmission potential. Next,  
77 the indicator was compared with the Cori reproduction number  $\mathcal{R}_t$  (12). Age groups are evenly distributed  
78 into nine categories, up to 90 years old. To characterize each wave, the period is divided into five phases,  
79 accounting for epidemic control measures and the progression of epidemic waves.

80 Our analysis focuses on the comprehensive infection network across age groups, revealing how infection  
81 spread patterns evolve over time, and concentrates on methods to obtain meaningful information in the  
82 presence of substantial missing data. This analytical approach, based on epidemiological data, emphasizes  
83 the role of active contact tracing by governments. Ultimately, our research suggests that active contact

84 tracing in real pandemic situations can offer policymakers data-driven insights for establishing more  
85 effective responses, thereby mitigating the pandemic's impact on global communities.

## 2 METHODS

### 86 2.1 Data and measurement

#### 87 A. Data description

88 We utilize the COVID-19 data (27) provided by the Korea Disease Control and Prevention  
89 Agency (KDCA) from January 19, 2020, to July 11, 2021, to construct the infection network for COVID-19  
90 transmission. In this paper, we analyze the dataset containing 169,597 confirmed cases (real-time reverse  
91 transcription polymerase chain reaction test positive cases), focusing on four specific records as follows.

#### 92 (ID, age, date of report, ID of the infector)

93 Here, the “ID” stands for the identity of the traced infectee, and “age” refers to the infectee's age. If “ID  
94 of the infector” is not traced (untraced), it is assigned a value of 0. Each confirmed case is assigned an  
95 anonymized ID number ranging from 1 to 169,146 associated with age, which ranges from 0 to 128, the  
96 date of report, and the ID number of the infector. Remark that in general the date of the report may not be  
97 exactly the same as the date of infection. The date of the from January 19, 2020, to July 11, 2021.

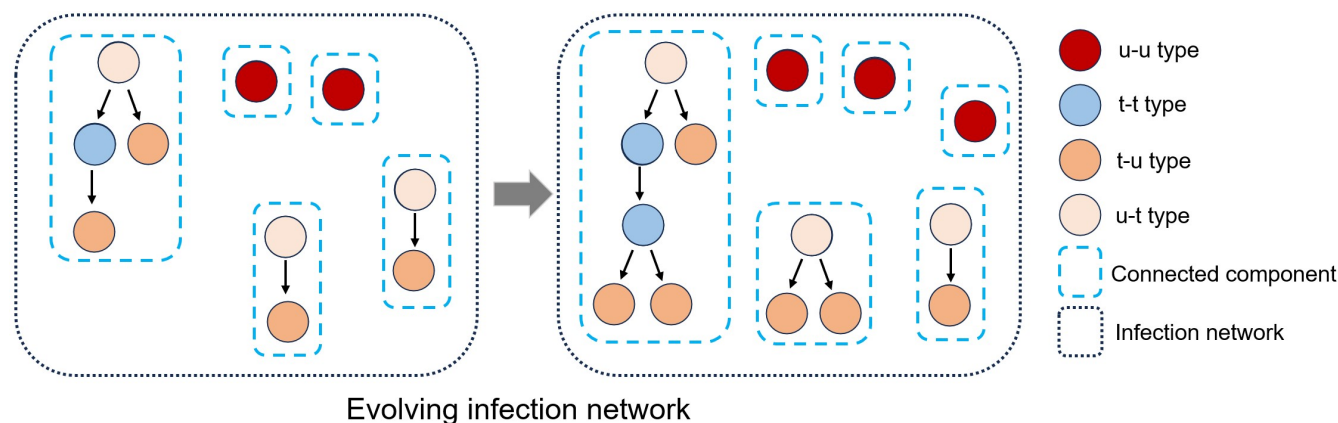
#### 98 B. Defining five periods of COVID-19 progression

99 The entire period was segmented into five distinct periods to observe the evolution of infection  
100 characteristics. This segmentation considered several critical factors like the emergence of new variants,  
101 vaccine rollout, change of social distancing levels, and other intervention measures (28).

- 102 • *P1* (January 19, 2020 ~ April 29, 2020): Since the first confirmed COVID-19 case on January 19,  
103 2020, South Korea experienced a moderate rise in cases, peaking at about 694 on February 26, 2020,  
104 primarily in Daegu-Gyeongbuk due to a church-related outbreak. Despite subsequent outbreaks at  
105 another church and a Seoul call center, daily cases gradually declined. Measures like the first social  
106 distancing period (March 22 to April 7, 2020) and a ban on gatherings in entertainment venues (April  
107 8 to April 19, 2020) were enacted, resulting in an average of 145 daily confirmed cases during these  
108 periods.
- 109 • *P2* (April 30, 2020 ~ July 14, 2020): During this period, there was the lowest number of daily  
110 confirmed cases compared to other periods. The average number of daily confirmed cases was 37.
- 111 • *P3* (July 15, 2020 ~ October 12, 2020): The second epidemic wave in South Korea started with a  
112 major outbreak at a Seoul church, accounting for 12% of the total infections in period *P3*, and was  
113 further exacerbated by a large rally on August 15 contributing to 6% of infections. In response, the  
114 government escalated Seoul's social distancing to level 2 on August 16, expanded it nationwide on  
115 August 23, and then increased it to level 2.5 in the metropolitan area by August 30. The peak of this  
116 wave was on August 24, 2020, with 418 cases, and the average daily confirmed cases during this period  
117 was 125.
- 118 • *P4* (October 13, 2020 ~ February 25, 2021): On October 12, the social distancing level was eased  
119 to level 1. *P4* coincides with the third epidemic wave, and it started with a gradual increase in  
120 daily confirmed cases without any apparent major events. The third epidemic wave peak occurred on  
121 December 23, 2020, with 1206 cases. The government raised the social distancing level on December 1  
122 and then again on December 8 and increased screening clinics. During this period, the average number  
123 of daily confirmed cases was 463.
- 124 • *P5* (February 26, 2021 ~ July 11, 2021): South Korea began its vaccination campaign on February 26,  
125 2021, and then saw an increase in delta variant cases starting April 18, 2021. During this period, the  
126 average number of daily confirmed cases was 571.

### 127 2.2 Infection network of infector and infectee

128 Network, also called graph mainly in mathematics, has been used as an explanatory tool to describe the  
129 dynamics of disease transmission (29). The terms “individuals (confirmed cases)” and “contacts (infects)”



**Figure 1.** The established infection network comprises many connected components due to missing vertices (individuals) and edges (infection events). An infection network's vertices can be classified into four types (u-t, u-u, t-u, t-t) based on the classification of their infector or infectee status as either traced or untraced. Also, the infection network evolves as an infectious disease spreads over time.

130 in epidemiology can be considered as “vertices” and “edges” in graph theory, respectively. For more details  
 131 on network epidemiology, see the review (30, 31) and references therein.

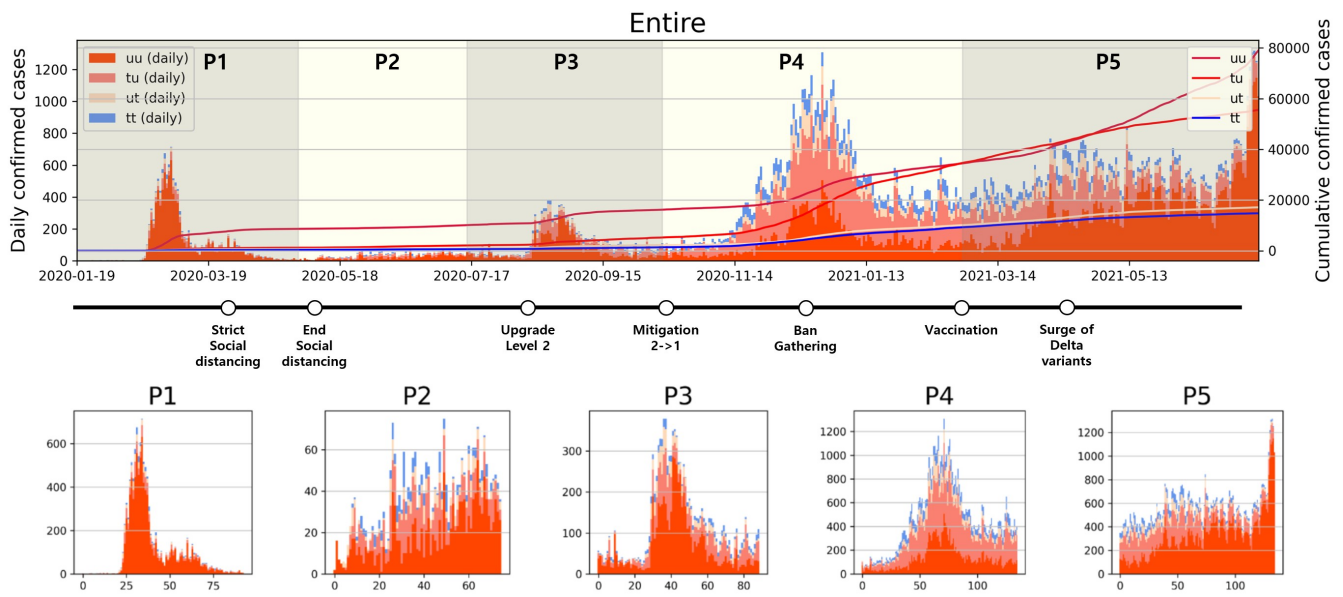
132 Denote the set of all confirmed IDs from January 19, 2020 to July 11, 2021 as  $\mathcal{I}$ , and let the set of  
 133 all infection events  $(m_{-1}, m_0)$  for the infector  $m_{-1} \in \mathcal{I}$  and its infectee  $m_0 \in \mathcal{I}$  as  $\mathcal{E}$ . We consider the  
 134 directed network  $G = (\mathcal{I}, \mathcal{E})$  as an infection network. For complete sampling, the infection network  $G$   
 135 must be weakly connected (replacing all its directed edges with undirected edges produces a connected  
 136 undirected graph). However, due to the existence of unreported infection cases, it is natural to assume that  
 137 the network is constructed by incomplete sampling of all confirmed individuals in a population (missing  
 138 vertices) and incomplete sampling of infection events between individuals (missing edges). So the infection  
 139 network  $G$  generated by real data consists of many weakly connected (or just connected components in this  
 140 paper) due to many missing vertices and edges, i.e., unreported individuals and infections. Hence analysis  
 141 of unreported infections is crucial for a better understanding of the real infection network in South Korea  
 142 and other countries.

### 143 2.3 Four type classifications

144 Each polymerase chain reaction (PCR)-confirmed case  $m_0$  can be classified into four different types  
 145 based on (i) as an infector  $m_{-1}$ , whether the infectees they have transmitted the virus to have been traced  
 146 or (ii) as an infectee  $m_1$ , whether they are aware of their infector being traced.

- 147 (i) An individual  $m_0 \in \mathcal{I}$  is said to be “**untraced-untraced**” type, denoted by **u-u**, if  $\{m_0 \in \mathcal{I} | (m_{-1}, m_0) \in \mathcal{E}\} = \emptyset$  and  $\{m_0 \in \mathcal{I} | (m_0, m_1) \in \mathcal{E}\} = \emptyset$ , i.e., its infector is missing (untraced) and  
 148 its infectee is missing or does not exist. Such an individual is represented as an isolated vertex on the  
 149 network.  
 150  
 151 (ii) An individual  $m_0$  is said to be “**traced-untraced**” type, denoted by **t-u**, if  $\{m_0 \in \mathcal{I} | (m_{-1}, m_0) \in \mathcal{E}\} \neq \emptyset$  and  $\{m_0 \in \mathcal{I} | (m_0, m_1) \in \mathcal{E}\} = \emptyset$ , i.e., its infector is confirmed (traced) but its infectee is  
 152 missing or does not exist. Such an individual is represented as a leaf of a directed tree graph.  
 153  
 154 (iii) An individual  $m_0$  is said to be “**untraced-traced**” type, denoted by **u-t**, if  $\{m_0 \in \mathcal{I} | (m_{-1}, m_0) \in \mathcal{E}\} = \emptyset$  and  $\{m_0 \in \mathcal{I} | (m_0, m_1) \in \mathcal{E}\} \neq \emptyset$ , i.e., its infector is not confirmed but its infectee is  
 155 confirmed. Such an individual is represented as a root of a directed tree graph.  
 156  
 157 (iv) An individual  $m_0$  is said to be “**traced-traced**” type, denoted by **t-t**, if  $\{m_0 \in \mathcal{I} | (m_{-1}, m_0) \in \mathcal{E}\} \neq \emptyset$   
 158 and  $\{m_0 \in \mathcal{I} | (m_0, m_1) \in \mathcal{E}\} \neq \emptyset$ , i.e., infector is confirmed and infectee is confirmed. Such an  
 159 individual is represented as neither a root nor a leaf in a directed tree graph.

160 Given an infection network, one can find the following properties due to the characteristics of infectious  
 161 disease transmission:



**Figure 2.** Categorized daily and cumulative confirmed cases over various periods are presented: (**Upper**) Entire period, (**Lower**)  $P1$  to  $P5$ , along with representative control measures implemented in South Korea. The contrasting background colors distinguish each period.

- 162 • The number of connected components with more than two vertices (individuals) equals the number of  
 163 individuals (vertices) of the u-t type.
- 164 • The number of individuals excluding the u-u type represents the total sum of the number of individuals  
 165 across all connected components with more than two vertices.
- 166 • The quotient of the number of individuals excluding the u-u type and the number of u-t type individuals  
 167 represents the average number of individuals per connected component.
- 168 • The quotient of the number of t-t type individuals and the number of u-t type individuals represents the  
 169 average number of t-t type individuals per connected component.

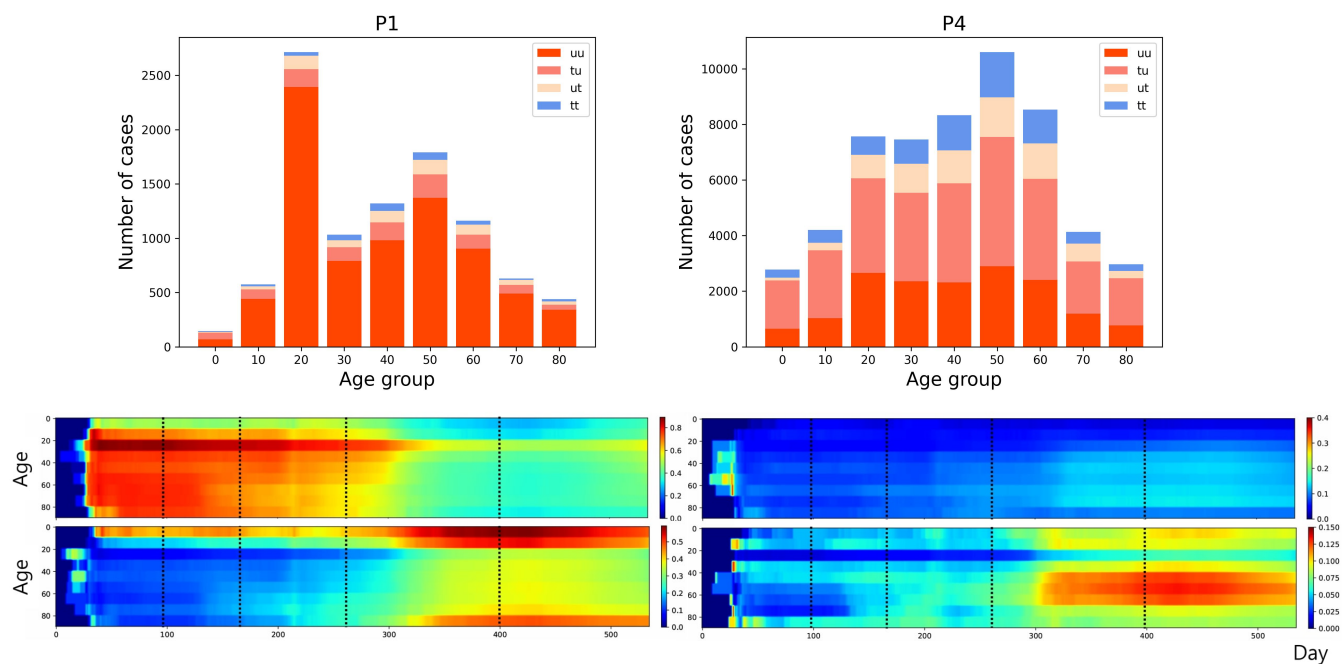
## 170 2.4 Experimental settings

171 Data preprocessing was performed before conducting the simulation. Firstly, 2,546 infection events  
 172  $(m_{-1}, m_0) \in \mathcal{E}$  were excluded due to missing report dates. Next, we identified 474 individuals,  $m_0 \in \mathcal{I}$ ,  
 173 linked to multiple infectors,  $m_{-1} \in \mathcal{I}$ , due to uncertainty about who the actual infector is, resulting in a  
 174 total of 1042 infection events,  $(m_{-1}, m_0) \in \mathcal{E}$ . Among the identified 1042 infection events  $(m_{-1}, m_0) \in \mathcal{E}$ ,  
 175 480 of these cases were of the u-t type for  $m_{-1} \in \mathcal{I}$ . Finally, we excluded the connected components that  
 176 include the u-t type from our data. Through all these preprocessing steps, the total number of confirmed  
 177 cases obtained is 164,314. All simulations were done in Python version 3.9. The calculation of  $\mathcal{R}_t$  was  
 178 carried out using the Epyestim library, employing Epyestim's default distributions and parameters. This  
 179 library is described in Thompson et al. (32).

## 3 RESULTS

### 180 3.1 Analysis for infection network by time periods

181 Analyzing daily confirmed cases alone is insufficient to fully understand the transmission dynamics of  
 182 infectious disease. Therefore, as depicted in Figure 2, confirmed cases have been categorized into four  
 183 types, and a period analysis was conducted. In Figure 2 upper panel, the period with the highest proportion  
 184 of u-u type cases among the four types was  $P1$ . In contrast, the highest proportions for the remaining three  
 185 types were observed in  $P4$ . Moreover, the cumulative number of confirmed cases during  $P4$  shows a sharp  
 186 increase, especially in the number of t-u type cases. On February 23, 2021, the cumulative number of u-t  
 187 type cases surpassed that of u-u type. However, starting from April 26, 2021, the cumulative number of u-u



**Figure 3. (Upper)** Age distribution categorized according to four types for both  $P1$  and  $P4$ . **(Lower)** The proportion of each case type within specific age groups over the cumulative period. The left panels display heatmap for u-u and t-u types, while the right panels show those for u-t and t-t types, with dotted lines in the figure marking the divisions between periods  $P1$  to  $P5$ .

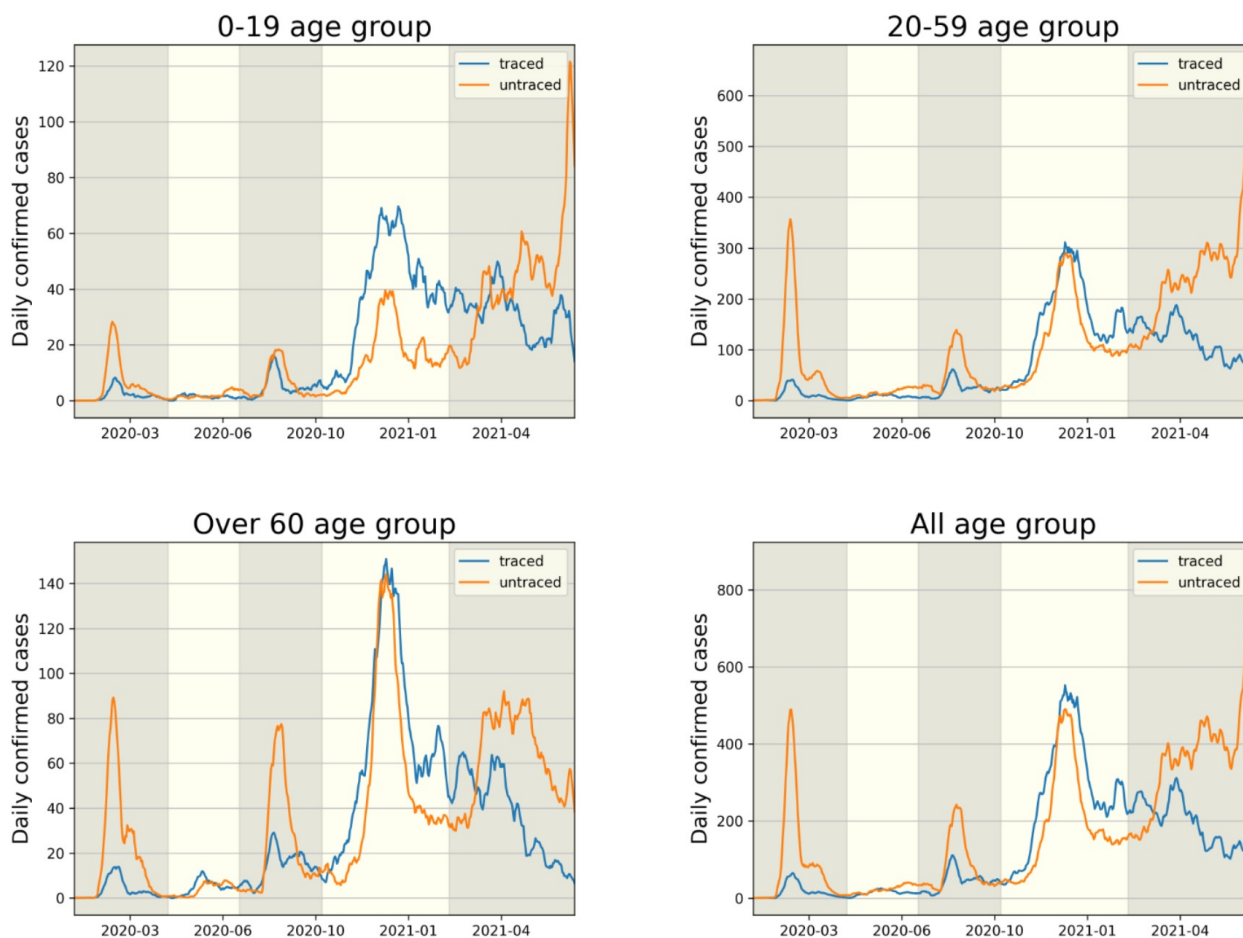
188 type cases began to increase sharply. The number of cumulative confirmed cases for u-t type is almost the  
 189 same as the number for t-t type over  $P4$ ,  $P5$ .

### 190 3.2 Analysis for infection network by time periods and age group

191 The transmission dynamics might be related to the contact pattern between age groups (7, 26, 33). Figure  
 192 3 upper panel displays the age distribution of four types for both  $P1$  and  $P4$ . During  $P1$ , a high number  
 193 of confirmed cases were observed in individuals in their 20–29 and 50–59. Among all age groups of  
 194 confirmed cases, 79% were classified as the u-u type. The highest proportion of u-u type cases was found  
 195 in the 20–29 age group, accounting for 88% of the cases in this age group, while the lowest was in the  
 196 0–9 age group, with 49%. However, in  $P4$ , there was a distinct shift with the majority of confirmed cases  
 197 being of the t-u type. This was most pronounced in the 0–9 age group, which had the highest proportion of  
 198 t-u type cases at 62%, whereas the 60–69 age group had the lowest at 42%. Additionally, throughout the  
 199 entire period under study, the 0–9 age group consistently exhibited the highest proportion of t-u type cases,  
 200 accounting for 47%. For the age distribution in other periods, refer to Appendix Figure 7. Figure 3 lower

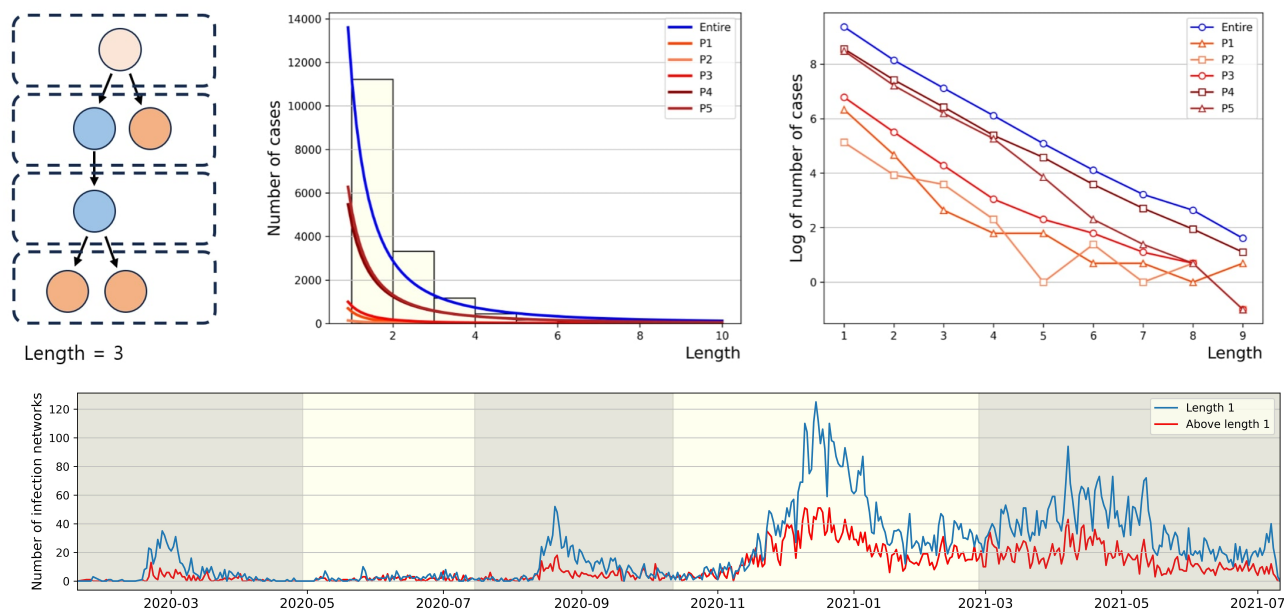
	0–9	10–19	20–29	30–39	40–49	50–59	60–69	70–79	80–89	90–99
P1	<b>0.92</b>	0.21	<b>0.07</b>	0.20	0.21	0.18	0.16	0.17	0.19	0.16
P2	0.95	0.64	0.45	<b>0.28</b>	0.46	0.93	1.16	<b>1.63</b>	1.13	1.40
P3	<b>0.86</b>	0.77	<b>0.41</b>	0.44	0.50	0.50	0.56	0.57	0.73	0.48
P4	<b>2.68</b>	2.20	<b>1.16</b>	1.19	1.38	1.45	1.31	1.26	1.79	2.03
P5	<b>0.87</b>	0.61	<b>0.38</b>	0.44	0.46	0.53	0.55	0.59	0.85	0.80
Entire	<b>1.32</b>	0.94	<b>0.51</b>	0.62	0.67	0.76	0.77	0.79	1.14	1.28

**Table 1.** It represents the ratio of the number of traced infectors to the number of untraced infectors for each period and age group. The red (resp. blue) color stands for the age group with the maximum (resp. minimum) ratio for each period.



**Figure 4.** The comparison of infector identification for traced (t-u, t-t type) and untraced (u-u, u-t type) cases is shown in each age group.

201 panel presents a heatmap representing the proportion of each case type within specific age groups over  
202 the cumulative period. For instance, on the u-u type heatmap, if the y-axis is labeled 20-29 and the x-axis  
203 indicates 400 days (February 28, 2021), the value corresponds to the proportion of 20-29 age group cases  
204 that are classified as u-u type up to 400 days. Due to the low number of cumulative confirmed cases in the  
205 early stages of COVID-19 spread, this paper will not interpret the results for this period. When considering  
206 the entire cumulative period, the age groups with the highest proportions of u-t type and t-t type cases  
207 are 70-79 and 50-59, respectively, each accounting for 13% and 11%. The heatmaps for each type are  
208 examined in sequence. Firstly, examining the u-u type heatmap, it is observed that until the mid-period of  
209  $P_4$ , the majority of confirmed cases in the 20-29 age group were of the u-u type. This trend is not exclusive  
210 to the 20-29 age group; up until the mid-period of  $P_4$ , a high proportion of u-u type cases is evident across  
211 most age groups. However, post the mid-period of  $P_4$ , there is a significant reduction in the proportion of  
212 u-u type cases in all age groups except for 20-29. Next, the t-u type heatmap shows a pattern opposite to  
213 that of the u-u type. The u-t type heatmap indicates an increase in the proportion of u-t type cases among  
214 the 40-79 age group after the mid-period of  $P_4$ . Lastly, the t-t type heatmap reveals an increase in the  
215 proportion of t-t type cases among the 40-69 age group posts the mid-period of  $P_4$ . We also analyzed the  
216 relationship between each type in terms of age group and period. As shown in Table 1, the value obtained  
217 from dividing the number of confirmed cases with traced infectors (or just traced infectors) by the number  
218 of confirmed cases with untraced infectors (or just untraced infectors) was calculated for each period and  
219 age group. In all periods except for  $P_2$ , the age group of 9 years and under has higher values compared to  
220 other age groups, and the 20-29 age group has the lowest values. Furthermore, this paper investigated  
221 the number of traced infectors and the number of untraced infectors across different age groups over time.  
222 These values were processed using a smoothing function with a uniform kernel of 10 points, where each



**Figure 5.** The figure upper panel presents the power law approximation of the distribution of connected component length for each period (**Middle**) and the same distributions on a log scale (**Right**), respectively. For convenience,  $y$ -axis value (log value of the number of cases) of  $-1$  indicates  $\log 0$ . The lower panel represents the number of connected components by length over time.

223 point is weighted equally ( $1/10$ ), to enhance data visualization and analysis. As shown in Figure 4, in  $P4$ ,  
 224 for individuals aged 20 and above, the number of untraced infectors is almost the same as the number of  
 225 traced infectors. However, in the age group below 20, there were more cases with a traced infector than  
 226 with an untraced one. During  $P5$ , there was a significant increase in the number of untraced infectors in  
 227 the 0–59 age group.

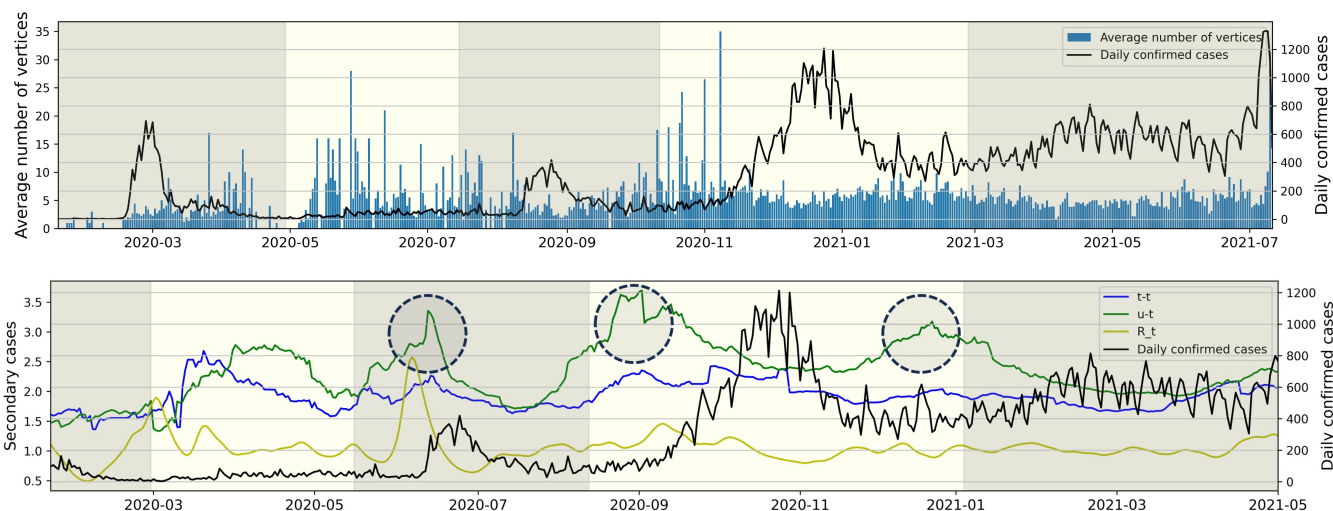
### 228 3.3 Length of the connected components of infection network

229 Infection order refers to the number of subsequent infections traced back to a single confirmed case.  
 230 For instance, if person A infects person B, and person B then infects person C, B and C are considered  
 231 the 2nd and 3rd order infected individuals, respectively, originating from A. In this paper, we define the  
 232 length of a connected component as  $n-1$ , where  $n$  is the highest order of an infector originating from a  
 233 u-t type individual in the connected component. As shown in Figure 5 (Middle), in  $P1$ , the proportion of  
 234 connected components with a length of 1 is the highest at 81%, compared to other periods. Conversely,  
 235 the lowest period is  $P2$  with 61%. For the distribution of connected component length in other periods,  
 236 refer to Appendix Figure 8. In Figure 5 (Right), for the entire period and  $P4$ , the slopes of the log scale for  
 237 the number of cases according to length, from length=1 to length=2, ..., and from length=8 to length=9,  
 238 all exhibit similar values. Another observation is that the slope from length=2 to length=3 being closest  
 239 to 0 occurs during period  $P2$ . The lower panel displays the number of connected components with the  
 240 length being either 1 or greater than 2, spanning the period from January 19, 2020, to July 11, 2021. During  
 241 each epidemic wave  $P1$ ,  $P3$  and  $P5$  at their respective peaks, the number of connected components with  
 242 a length of 2 or more is significantly smaller compared to the number of connected components with a  
 243 length of 1.

### 244 3.4 Daily confirmed cases relationship

245 Figure 6 (Upper) represents the average number of individuals per connected component for each day  
 246 from January 19, 2020, to July 11, 2021. For instance, the value for November 31, 2020, is calculated as  
 247 the sum of t-t and t-u type individuals on November 31, 2020, divided by the number of u-t type individuals  
 248 on the same date. The observation revealed that the value and the daily confirmed cases exhibited opposing  
 249 trends. During the epidemic waves of  $P1$  and  $P3$ , the value is lower compared to periods not experiencing





**Figure 6.** The right  $y$ -axis and the black line represent daily confirmed cases, while the left  $y$ -axis represents all other values. **(Upper)** The average number of individuals (vertices) per connected component for each day. **(Lower)** The average number of secondary cases for each type and time-dependent reproduction number  $\mathcal{R}_t$  over time.

250 an epidemic wave. Following the surge in daily confirmed cases in  $P4$ , the value remains consistent without  
251 significant increases. Figure 6 (Lower) illustrates the average number of secondary cases for both u-t and  
252 t-t types, calculated with a window size of 30, from March 22, 2020, to July 11, 2021, and also depicts the  
253 time-dependent reproduction number  $\mathcal{R}_t$  (12). The value is an indicator derived from the infection network  
254 analysis. For instance, the average number of secondary cases for the u-t (resp. t-t) type on August 1, 2020,  
255 is defined as the real-time calculated average value of confirmed cases directly infected by the u-t (resp. t-t)  
256 type within the infection network identified between July 1, 2020, and August 1, 2020. For instance, if  
257 within the identified infection network for the period, there are 3 connected components, and the number  
258 of individuals infected by each u-t type individual is 2, 6, and 1, respectively, then the average number of  
259 secondary infections for the u-t type on August 1, 2020, is calculated as  $(2+6+1)/3=3$ . The time-dependent  
260 reproduction number  $\mathcal{R}_t$  did not show a significant increase before an increase in daily confirmed cases  
261 during  $P4$  and  $P5$ . However, the circular markers in Figure 6 (Lower) indicate a significant increase in the  
262 average number of secondary cases for u-t type.

## 4 DISCUSSION

263 Despite having a large volume of epidemiological data due to its active contact tracing efforts compared  
264 to other countries, South Korea's infection network, generated from the data, comprises many connected  
265 components as a result of numerous missing vertices (individuals) and edges (infection events). This article  
266 analyzed the infection network using vertices of four types: u-u, u-t, t-u, and t-t based on whether their  
267 infector or infectee falls into the traced or untraced category. We analyzed the dynamics of the infection  
268 network based on each type, time, and age group, deriving insights. Our results showed a significant  
269 surge in the number of t-u type cases (i.e., traced infector - untraced infectee type) during  $P4$  when  
270 the government upgraded the social distancing level twice as well as expanding the screening clinics in  
271 Figure 2. A significant surge in the cumulative number of u-u type cases was also observed, beginning  
272 in the mid-phase of  $P5$ , coinciding with the spread of the Delta variant. The average number of t-t type  
273 individuals per connected component close to 1 in  $P4$  and  $P5$  indicates active contact tracing in response to  
274 mass infection. In other words, the proposed method allows for the analysis and evaluation of phenomena  
275 induced by various events such as the implementation of public health policies, the emergence of new  
276 variants, and more.

277 Our result also found age-specific transmission patterns for the four types in Figure 3. Individuals of the  
278 u-u type pose a significant risk of causing mass infections in the community. Across periods  $P1$  to  $P5$ , the  
279 highest proportion of u-u type cases (57.4%) was observed in the 20–29 age group. This can be inferred to

280 be due to the 20–29 age group’s wider range of activities and frequent interactions with various people.  
281 The 0–9 (47.6%), 10–19 (40.9%), and 80–89 (46.5%) age groups had the highest rates of t-u type cases,  
282 indicating these demographics may serve as key points for interrupting transmission chains. By focusing  
283 on these patterns in the implementation of public health policies, it may be possible to more effectively  
284 contain outbreaks and prevent wider community spread. Individuals of the u-t type, as initial infectors in a  
285 connected component, help identify which age groups had more asymptomatic COVID-19 cases and were  
286 more engaged in contact tracing, based on their age-wise proportions. Across periods  $P_1$  to  $P_5$ , the highest  
287 proportion of u-t type cases (13%) was observed in the 70–79 age group. From mid  $P_4$ , it was observed  
288 that the proportion of u-t type cases in the 30–79 age group was higher compared to other age groups. The  
289 proportion of t-t type cases by age group also allows for the inference of which age groups were more  
290 actively involved in contact tracing. Across periods  $P_1$  to  $P_5$ , the highest proportion of t-t type cases (11%)  
291 was observed in the 50–59 age group. After mid  $P_4$ , the 40–69 age group showed a higher proportion of t-t  
292 type cases compared to other age groups. Furthermore, the analysis of the value obtained from dividing  
293 the number of confirmed cases with traced infectors (or just traced infectors) by the number of confirmed  
294 cases with untraced infectors (or just untraced infectors) across age groups revealed a sequence of  $0-9 >$   
295  $90-99 > 80-89 > 10-19 > 70-79 > 60-69 > 50-59 > 40-49 > 30-39 > 20-29$ . For the 0–9 and 80–99  
296 age groups, where the number of contacts is limited, contact tracing was more manageable; however, in  
297 age groups like 20–39, which have a higher number of contacts, contact tracing was found to be more  
298 challenging. These analyses provide valuable information for understanding the transmission dynamics  
299 of COVID-19, allowing us to suggest strengthening or relaxing control measures for specific age groups  
300 based on the period’s characteristics.

301 Our results also investigated the distribution of the lengths of connected components within the infection  
302 network. In  $P_2$ , the proportion of connected components with a length of 1 was the lowest, while the  
303 proportions with lengths of 2 and 3 were the highest. This indicates that during  $P_2$ , which had the lowest  
304 daily average of 37 confirmed cases, the infection network had fewer missing edges (infection events).  
305 Further investigation across the entire period, as shown in the lower panel of Figure 5, revealed an increase  
306 in the number of connected components with a length of 1 during surges in daily confirmed cases. The  
307 earlier results motivated the hypothesis that the average number of individuals per connected component  
308 for each day would decrease during spikes in infections. This was indeed observed in the upper panel  
309 of Figure 6. It means that when the number of daily confirmed cases surges, it becomes challenging to  
310 contact trace high-order transmissions. This phenomenon may stem from changes in the government and  
311 the public’s willingness to engage in contact tracing and limitations of existing contact tracing methods in  
312 the face of a highly infectious virus spreading worldwide. For this reason, this article proposed the average  
313 value of confirmed cases directly infected by the u-t type as an indicator of infectious disease transmission  
314 potential. Utilizing the infection network up to 30 days prior allows for real-time calculation, and this  
315 indicator shows high values before a surge in daily confirmed cases. Due to the indicator allowing for an  
316 approximation of real-time unreported cases, it is more sensitive compared to  $\mathcal{R}_t$  and increases before the  
317 third epidemic wave. We anticipate it to be a useful indicator in situations like in South Korea, where active  
318 contact tracing is conducted.

319 Our study has several limitations. Firstly, we do not consider unreported cases including asymptomatic  
320 individuals, those with mild symptoms who were not tested, and unreported self-tests from the surveillance  
321 pyramid (34). Considering unreported cases is a key research topic for understanding and predicting the  
322 scale of infections (35, 36, 37). Acknowledging the constraints imposed by unreported cases, especially  
323 concerning COVID-19 transmission within contact networks, we recognize the potential of methods such  
324 as multiple imputation techniques (35) and data augmentation through link prediction (36) to provide  
325 valuable insights. Furthermore, the exploration of machine learning-based approaches (37) presents  
326 another promising avenue for addressing data gaps. Studies that have not estimated unreported cases  
327 but have specifically limited unreported cases to environmental factors include Myall et al. (38), which  
328 analyzed patient-contact networks using patient contacts obtained from hospital health records. Despite  
329 its limitations, the KDCA data we analyzed remains trustworthy. According to the KDCA, based on  
330 serological surveillance and contact tracing data, the rate of unreported cases in South Korea from January  
331 19, 2020, to July 30, 2022, was approximately 19.5%. This rate is notably lower than those seen in  
332 international contexts, a difference attributed to the widespread availability of testing and the public’s  
333 adherence to control measures (39, 40). Secondly, our study did not quantitatively assess contact tracing  
334 effectiveness. There are several previous studies about the effectiveness of contact tracing strategies for  
335 COVID-19 (41, 42, 1). Kretzschma et al. (41) analyzed contact tracing effectiveness using a stochastic

336 model, finding that immediate tracing and testing are crucial for reducing the spread of COVID-19. Delays  
337 in testing and tracing significantly diminish the potential to keep the effective reproduction number below  
338 1. Korean government implemented the contact tracing described in (42). Contact tracing for COVID-19  
339 was performed using information from credit card records, handwritten visitor logs, QR codes through  
340 KI-Pass, and the Safe Call system after interviews in Korea. Hellewell et al. (1) found tracing and isolation  
341 could control outbreaks within 12 weeks. There are previous studies to investigate the infection network  
342 of COVID-19 in (2, 43, 44). Luo et al. (43) in 2021 developed an infection network considering the  
343 history of exposure and transmission source. The visualization method, which identifies vertices in the  
344 infection network as clusters of infected individuals, revealed a highly central infection cluster in (44).  
345 However, we developed an infection network, categorizing infector-infectee pairs by age group and periods,  
346 specifically focusing on untraced cases. Jo et al. (2) emphasized the importance of gathering network data  
347 and examining network structures to improve the effectiveness of governmental responses to COVID-  
348 19. Additionally, in future research, we intend to expand our analysis to encompass infection networks  
349 incorporating spatial information, as discussed in (45).

350 The current research reveals that, despite active contact tracing efforts, South Korea's infection network,  
351 derived from a large volume of epidemiological data, comprises many connected components due  
352 to numerous missing entities (individuals) and infection events (edges). The presence of numerous  
353 connected components complicates the inference of relationships between vertices. Therefore, a four-type  
354 classification method for vertices (confirmed cases) is proposed. This method enables the categorization  
355 of vertices within the numerous distinct connected components from a common perspective, thereby  
356 facilitating the analysis and interpretation for each vertex type. The changes in the number of cases  
357 for each type over time relate to the emergence of new coronavirus variants (such as Delta) or the  
358 implementation of control measures. When analyzed by age group, it was observed that certain age groups  
359 are more sensitive to these events. Additionally, we analyzed the infection network from the perspective  
360 of connected components, proposing a new indicator and comparing it with  $\mathcal{R}_t$ . Despite limitations, the  
361 study's categorization of epidemiological data into four types not only offers a robust foundation for  
362 evaluating public health policies and comprehending the dynamics of COVID-19 transmission but also  
363 serves as a foundational health planning tool for resource management and tool selection/development for  
364 contact tracing.

365 In conclusion, South Korea's epidemiological data generated from active contact tracing enables novel  
366 infection network analysis. Our analysis reveals significant age-specific transmission patterns, particularly  
367 in the 20–29, 40–69, and 0–9 age groups. The patterns show a distinct shift around the midpoint of  $P_4$ ,  
368 with the 20–29 (57.4%) age group exhibiting the highest proportion of u-u type cases, the 40–69 age group  
369 predominantly showing u-t and t-t types, and the 0–9 (47.6%) age group having the highest rate of t-u type  
370 cases across entire periods. This suggests a relationship between age groups and the four-type classification.  
371 A significant increase in t-u and u-u type cases was observed during certain periods, providing opportunities  
372 for analysis and evaluation of phenomena induced by various events, such as the implementation of public  
373 health policies, the emergence of new COVID-19 variants, and more. Also, through the investigation of the  
374 distribution of lengths of connected components within the infection network, we found that the average  
375 number of individuals per connected component tends to decrease during surges in daily confirmed cases,  
376 indicating that tracing high-order transmissions becomes more challenging. Accordingly, we propose the  
377 average value of confirmed cases directly infected by the u-t type as an indicator to assess the potential for  
378 infectious disease transmission. Additionally, this approach could facilitate the early detection of changes  
379 in willingness among individuals to participate in tracing, or in the reduced capacities of contact tracing  
380 systems. The investigation of infection networks is crucial for advancing our capacity to control and  
381 mitigate the transmission of infectious diseases. Recognizing the necessity for a more thorough age-based  
382 categorization, the study emphasizes potential areas for future research improvements in comprehending  
383 and refining public health strategies. Additionally, our study presents a new real-time indicator using contact  
384 tracing data collected during actual infection spread, ultimately providing support for decision-makers and  
385 contributing to reducing the pandemic's impact on global communities.

## DATA AVAILABILITY STATEMENT

386 The original contributions presented in the study are included in the article material, further inquiries can  
387 be directed to the corresponding author.

## AUTHOR CONTRIBUTIONS

388 HW, Lee and HY, Choi: analyzed the data. HW, Lee, HY, Choi, HJ, Lee, SM, Lee and CH, Kim: drafted  
389 and revised the manuscript. HJ, Lee, SM, Lee and CH, Kim: interpreted the results. All authors contributed  
390 to the article and approved the submitted version.

## FUNDING

391 This work was supported by the National Research Foundation of Korea(NRF) grant funded by the Korea  
392 government(MSIT) (No. 2022R1A5A1033624).

## ACKNOWLEDGMENTS

393 Epidemiological data were obtained from Korea Disease Control and Prevention Agency (KDCA) (27).

## CONFLICT OF INTEREST

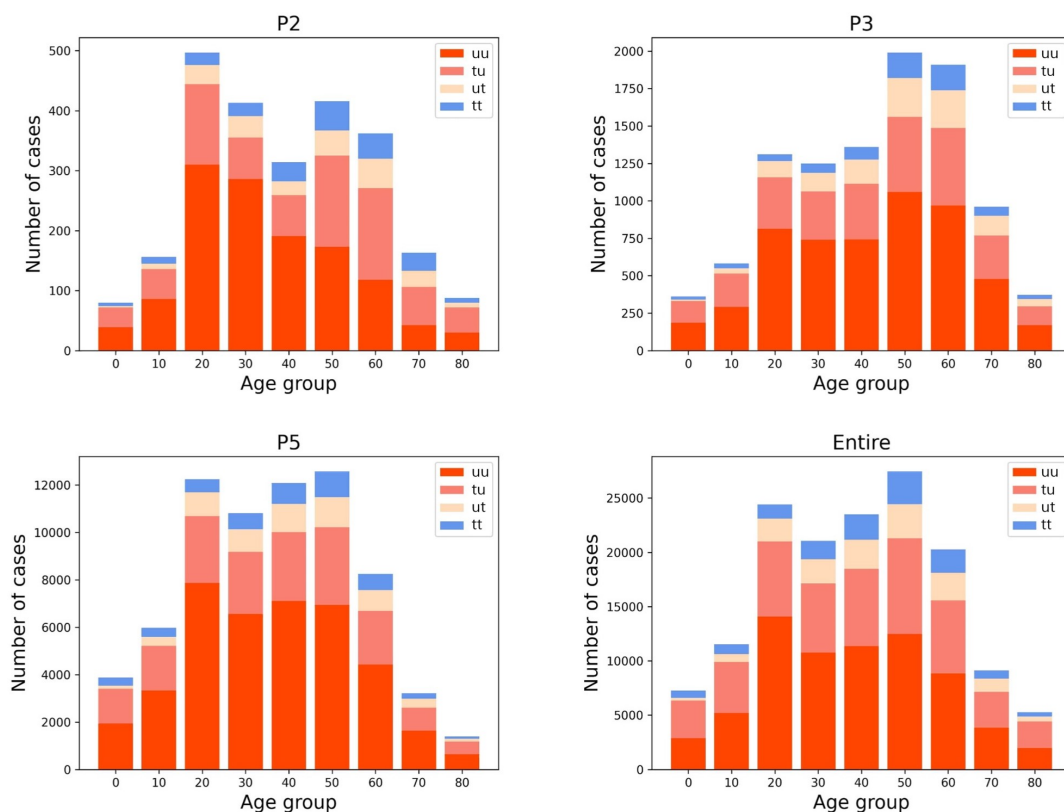
394 The authors declare that the research was conducted in the absence of any commercial or financial  
395 relationships that could be construed as a potential conflict of interest.

## REFERENCES

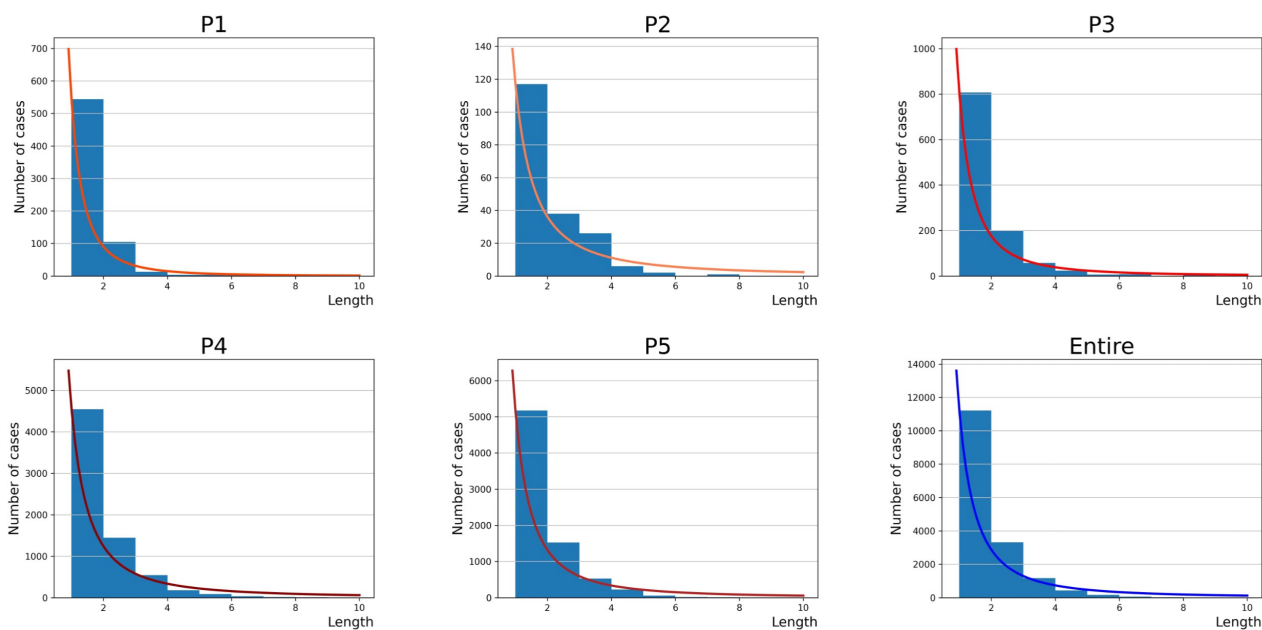
- 396 1 .Hellewell J, Abbott S, Eggo R. Feasibility of controlling COVID-19 outbreaks by isolation of cases  
397 and contacts. *The Lancet Global Health* **8** (2020) e488–e496.
- 398 2 .Jo W, Chang D, You M, Ghim G. A social network analysis of the spread of COVID-19 in South Korea  
399 and policy implications. *Scientific Reports* **11** (2021) 8581.
- 400 3 .Ryu S, Ali S, Cowling B. Transmission dynamics and control of two epidemic waves of SARS-CoV-2  
401 in south korea. *BMC Infectious Diseases* **21** (2021) 1–9.
- 402 4 .Ryan M. In defence of digital contact-tracing: human rights, South Korea and COVID-19. *International*  
403 *Journal of Pervasive Computing and Communications* **16** (2020).
- 404 5 .Turner S, Klimek P, Hanel R. A network-based explanation of why most COVID-19 infection curves  
405 are linear. *Proceedings of the National Academy of Sciences* **117** (2020) 22684–22689.
- 406 6 .Choi Y, Park M, Lee J. Types of COVID-19 clusters and their relationship with social distancing in  
407 the seoul metropolitan area, South Korea. *International Journal of Infectious Diseases* **106** (2021)  
408 363–369.
- 409 7 .Monod M, Blenkinsop A, Tietze S. Age groups that sustain resurging COVID-19 epidemics in the  
410 united states. *Science* **371** (2021).
- 411 8 .Davies N, Klepac P, Eggo R. Age-dependent effects in the transmission and control of COVID-19  
412 epidemics. *Nature medicine* **26** (2020) 1205–1211.
- 413 9 .Hao X, Cheng X, Wang C. Reconstruction of the full transmission dynamics of COVID-19 in wuhan.  
414 *Nature* **584** (2020) 420–424.
- 415 10 .Wang Y, You X, Li J. Estimating the basic reproduction number of COVID-19 in Wuhan, China. *Chin*  
416 *J Epidemiol* **41** (2020).
- 417 11 .Zhang J, Dong X, Gao Y. Risk and protective factors for COVID-19 morbidity, severity, and mortality.  
418 *Clinical Reviews in Allergy and Immunology* **64** (2023).
- 419 12 .Cori A, Ferguson N, Cauchemez S. A new framework and software to estimate time-varying  
420 reproduction numbers during epidemics. *American journal of epidemiology* **178** (2013).
- 421 13 .Oka T, Wei W, Zhu D. The effect of human mobility restrictions on the COVID-19 transmission  
422 network in China. *PloS one* **16** (2021).
- 423 14 .Meyers L, Pourbohloul B, Brunham R. Network theory and SARS: predicting outbreak diversity.  
424 *Journal of theoretical biology* **232** (2005).
- 425 15 .Glass R, Glass L, Min H. Targeted social distancing designs for pandemic influenza. *Emerging*  
426 *infectious diseases* **12** (2005).
- 427 16 .Skums P, Kirpich A, Chowell G. Global transmission network of SARS-CoV-2: from outbreak to  
428 pandemic. *MedRxiv* (2020).

- 429 17 .Wang P, Lu J, Chen S. Statistical and network analysis of 1212 COVID-19 patients in Henan, China.  
430 *International Journal of Infectious Diseases* **95** (2020).
- 431 18 .Kim T, Lee H, Lee S. Improved time-varying reproduction numbers using the generation interval for  
432 COVID-19. *Frontiers in Public Health* **11** (2023).
- 433 19 .Bi Q, Wu Y, Feng T. Epidemiology and transmission of COVID-19 in 391 cases and 1286 of their close  
434 contacts in Shenzhen, China: a retrospective cohort study. *The Lancet Infectious Diseases* **20** (2020).
- 435 20 .Lam H, Lam T, Chuang S. The epidemiology of COVID-19 cases and the successful containment  
436 strategy in Hong Kong—January to May 2020. *International Journal of Infectious Diseases* **98** (2020).
- 437 21 .Chen C, Jyan H, Chan C. Containing COVID-19 among 627,386 persons in contact with the Diamond  
438 Princess cruise ship passengers who disembarked in Taiwan: big data analytics. *Journal of Medical  
439 Internet Research* **22** (2020).
- 440 22 .Choi J. COVID-19 in South Korea. *Postgraduate Medical Journal* **96** (2020).
- 441 23 .Ferretti J, Wymant C, Fraser C. Quantifying SARS-CoV-2 transmission suggests epidemic control with  
442 digital contact tracing. *Science* **368** (2020).
- 443 24 .Keeling M, Hollingsworth T, Read J. Efficacy of contact tracing for the containment of the 2019 novel  
444 coronavirus (COVID-19). *J Epidemiol Community Health* **74** (2020).
- 445 25 .Peak C, Kahn R, Buckee C. Individual quarantine versus active monitoring of contacts for the mitigation  
446 of COVID-19: a modelling study. *The Lancet Infectious Diseases* **20** (2020).
- 447 26 .Arregui S, Aleta A, Moreno Y. Projecting social contact matrices to different demographic structures.  
448 *PLoS Computational Biology* **14** (2018).
- 449 27 .[Dataset] Korea disease control and prevention agency. kdca (2022).
- 450 28 .Jeon J, Han C, Lee S. Evolution of responses to COVID-19 and epidemiological characteristics in  
451 South Korea. *International Journal of Environmental Research and Public Health* **19** (2022).
- 452 29 .Newman M. Spread of epidemic disease on networks. *Physical Review E* **66** (2002).
- 453 30 .Keeling M, Eames K. Networks and epidemic models. *Journal of the Royal Society Interface* **2** (2005).
- 454 31 .Pastor-Satorras R, Castellano C, Vespignani A. Epidemic processes in complex networks. *Reviews of  
455 Modern Physics* **87** (2015).
- 456 32 .Thomson R, Stockwin J, Cori A. Improved inference of time-varying reproduction numbers during  
457 infectious disease outbreaks. *Epidemics* **29** (2019).
- 458 33 .Prem K, Cook A, Jit M. Projecting social contact matrices in 152 countries using contact surveys and  
459 demographic data. *PLoS Computational Biology* **13** (2017).
- 460 34 .Ricoca P, Carla N, Alexandre A. Epidemic surveillance of Covid-19: considering uncertainty and  
461 under-ascertainment. *Portuguese Journal of Public Health* **38** (2020).
- 462 35 .Elena C. A method for comparing multiple imputation techniques: A case study on the US national  
463 COVID cohort collaborative. *Journal of Biomedical Informatics* **139** (2023).
- 464 36 .David O. Constructing co-occurrence network embeddings to assist association extraction for COVID-  
465 19 and other coronavirus infectious diseases. *Journal of the American Medical Informatics Association*  
466 **27** (2020).
- 467 37 .Daniel S, Buhlmann P. MissForest—non-parametric missing value imputation for mixed-type data.  
468 *Bioinformatics* **28** (2012).
- 469 38 .Myall A, Price J, Barahona M. Prediction of hospital-onset COVID-19 infections using dynamic  
470 networks of patient contact: an international retrospective cohort study. *The Lancet Digital Health* **4**  
471 (2022).
- 472 39 .Zhan C. Estimating unconfirmed COVID-19 infection cases and multiple waves of pandemic  
473 progression with consideration of testing capacity and non-pharmaceutical interventions: A dynamic  
474 spreading model. *Information Sciences* **607** (2022).
- 475 40 .Ali H, Liu J, Schwiag T. Estimating the fraction of unreported infections in epidemics with a known  
476 epicenter: An application to COVID-19. *Journal of Econometrics* **220** (2021).
- 477 41 .Mirjam K. Impact of delays on effectiveness of contact tracing strategies for COVID-19: a modelling  
478 study. *The Lancet Public Health* (2020).
- 479 42 .Gong S, Jung J. Perceived usefulness of COVID-19 tools for contact tracing among contact tracers in  
480 Korea. *Epidemiology and Health* **44** (2022).
- 481 43 .Luo C, Ma Y, Yin F. The construction and visualization of the transmission networks for COVID-19: A  
482 potential solution for contact tracing and assessments of epidemics. *Scientific Reports* **11** (2021).
- 483 44 .Van G. Visualizing the network structure of COVID-19 in Singapore. *Socius* **7** (2021).
- 484 45 .Kwon O, Jo H. Clustering and link prediction for mesoscopic COVID-19 transmission networks in  
485 Republic of Korea. *Chaos: An Interdisciplinary Journal of Nonlinear Science* **33** (2023).

486 APPENDIX



**Figure 7.** Additional information for Figure 3. Age distribution categorized according to four types for both  $P_2$ ,  $P_3$ ,  $P_5$  and Entire.



**Figure 8.** Additional information for Figure 5. The figure presents the distribution of connected component length for each period.