

1 **Appendix**

2

3

4 Table of Contents

5 **S1 Definitions of COVID-19 case, close contact, transmission pair, and case cluster** ..... 2

6 **S2 Confirmation of genetic sequence** ..... 3

7 **S3 Data cleaning**..... 4

8 **S4 Modelling secondary cases distribution and superspreading**..... 5

9     *S4.1 Secondary case distribution*..... 5

10     *S4.2 Expected proportion of seed cases generating 80% of transmission* ..... 5

11     *S4.3 Parameter estimation* ..... 5

12     *S4.4 Supplementary results* ..... 6

13 **S5 Modelling heterogeneity in transmission risk**..... 8

14     *S5.1 Distribution of secondary attack ratio*..... 8

15     *S5.2 Parameter estimation* ..... 8

16     *S5.3 Supplementary results* ..... 8

17 **Supplementary references** ..... 10

18

19

20

21 **S1 Definitions of COVID-19 case, close contact, transmission pair, and case cluster**

22 All COVID-19 cases were defined as individuals who received a positive test outcome (with  
23 Ct value < 40) from real-time reverse transcription polymerase chain reaction (RT-PCR) assay for  
24 the genetic segments of SARS-CoV-2 strains using ORF1ab gene or N gene detection kit.

25 The definition of clinical severity after SARS-CoV-2 infection was as follows. For a  
26 symptomatic case, it was defined as a case presenting one of the relevant clinical symptoms,  
27 including fever, respiratory symptoms, and radiographic evidence of pneumonia. The asymptomatic  
28 or mildly symptomatic SARS-CoV-2 infection, which was classified as “asymptomatic” in this study,  
29 was defined as a case presenting one of the relevant clinical symptoms for less than 7 days, including  
30 fever, or respiratory symptoms, and without a radiographic evidence of pneumonia. For  
31 asymptomatic infection, it was defined as a case having no clinically evident symptoms. Since most  
32 of confirmed cases (92.7%) were asymptomatic, we avoided to further classify the clinical severity  
33 of cases into more detailed levels due to limited samples.

34 We defined the close contacts of a confirmed COVID-19 case as individuals having an  
35 epidemiologic link to a COVID-19 case, i.e., individuals free from symptoms and COVID-19  
36 diagnosis exposed to a RT-PCR test positive person. As a considerable amount of transmission could  
37 occur at very early stage after infection [1, 2], individuals who had been exposed to a case within 4  
38 days before the test-positive date of the case would also be counted as close contacts. We classified  
39 the close contacts of confirmed cases into categories described as follows:

- 40 • household contacts (i.e., household members regularly living within the same or close space, or  
41 relatives who had close contact with case)
- 42 • workplace or school contacts (i.e., a work colleague or classmate), and
- 43 • community contacts (i.e., healthcare-givers and patients in the same ward, persons sharing a  
44 vehicle or restaurant, and community workers having contact with case in public places).
- 45 • unknown contacts (i.e., only show for the contact with the space, no specific contact way).

46 For those contacts who were (eventually) test-positive for COVID-19, we treated these contacts as  
47 infectee, and their source case (who were confirmed with COVID-19 in the first place) as infector  
48 and forms transmission pairs.

49 Based on the identified transmission pairs, we thereafter grouped the linked cases into case-  
50 clusters, which is defined as a case or a cluster of cases (i.e., infectees) with a common single source  
51 of infection (i.e., infectors). The number of secondary cases generated by each infector was then  
52 extracted. The number of secondary cases generated by each infector was then extracted. As there  
53 might be epidemiological linkages between case clusters, we further linked those case-clusters into  
54 transmission chains, which could involve multiple generation of infections. Based on the locations  
55 where the infection occurred, we also identified 3 contact settings, including household, community,  
56 and workplace. The case clusters and transmission chains were constructed independently by 2  
57 authors. The final list of included cases was decided following discussion between the authors, with  
58 full agreement required prior to inclusion.

59

60

61 **S2 Confirmation of genetic sequence**

62 Nasopharyngeal or oropharyngeal swabs specimen from 11 confirmed cases during the first  
63 few days of the outbreak were collected, and undergone whole-genome sequencing. MAFFT  
64 program was used to perform multiple sequence alignments, and the GTR + CAT nucleotide  
65 substitution model in FastTree (version 2.1.11) were applied to explore the phylogenetic relationship.  
66 On the basis of Phylogenetic Assignment of Named Global Outbreak (PANGO) lineage designation,  
67 the samples were eventually classified as SARS-CoV-2 Omicron BA.5.2 sub-lineage.

68 There was a total of 62 amino acid (AA) substitutions in different genetic segments of SARS-  
69 CoV-2, including 31 in spike (S) protein, 19 in non-structural proteins (NSPs) ORF1a and ORF1b, 4  
70 in membrane protein, 4 in nucleocapsid (N) protein, 3 in auxiliary proteins ORF3a and ORF9b, and 1  
71 in envelope (E) protein.

72

73

74 **S3 Data cleaning**

75 We collected epidemiological contact tracing data of laboratory-confirmed cases with  
76 Omicron BA.5.2 infection between August 7 and September 7, 2022, from the Xinjiang Uygur  
77 Autonomous Region Health Committee. A total of 1139 confirmed cases were included. Among  
78 these confirmed cases, 649 were those within case cluster with size  $>1$ , and 236 were infector  
79 with  $>0$  offspring cases, 413 terminal cases. There were 490 sporadic or cases with unknown source  
80 (i.e., cases without known source of infection and secondary cases). Of the 1139 positive cases, 370  
81 test-positive individuals during isolation with 0 close contact, and the remaining 769 test-positive  
82 individuals associating with 51323 test-negative close contacts. Among these contacts, according to  
83 the classification of places of contact (**Appendix S1**), there were 1660 household contacts, 1998  
84 community contacts, 1766 workplace contacts and 46362 unknown contacts.

85

86

## 87 **S4 Modelling secondary cases distribution and superspreading**

### 88 *S4.1 Secondary case distribution*

89 Given the stochastic effect of the transmission events, we assumed the number of secondary cases  
90 generated by an infector followed a Negative binomial (NB) distribution which was parametrized by  
91 a reproduction number ( $R$ ) and a dispersion parameter ( $k$ ), as followed by previous studies [3,4]. The  
92 probability mass function of NB distribution is given by:

$$93 \Pr(Z = z; R, k) = \frac{\Gamma(k + z)}{\Gamma(k) \cdot \Gamma(z + 1)} \left(\frac{R}{R + k}\right)^z \left(\frac{k}{R + k}\right)^k$$

94 Here,  $z$  is the number of secondary cases generated by an infector, and  $\Gamma()$  denotes the gamma  
95 function.

### 96 97 *S4.2 Expected proportion of seed cases generating 80% of transmission*

98 Given the  $R$  and  $k$  estimates, we calculated the expected proportion of cases that were  
99 responsible for 80% of all transmissions [3], which is given by [4]:

$$100 \quad 1 - P = \int_0^Z \Pr(Z = z; R, k) dz$$

101 where  $Z$  satisfies:

$$102 \quad 1 - 80\% = \frac{1}{R} \int_0^Z \lfloor z \rfloor \Pr(Z = z; R, k) dz$$

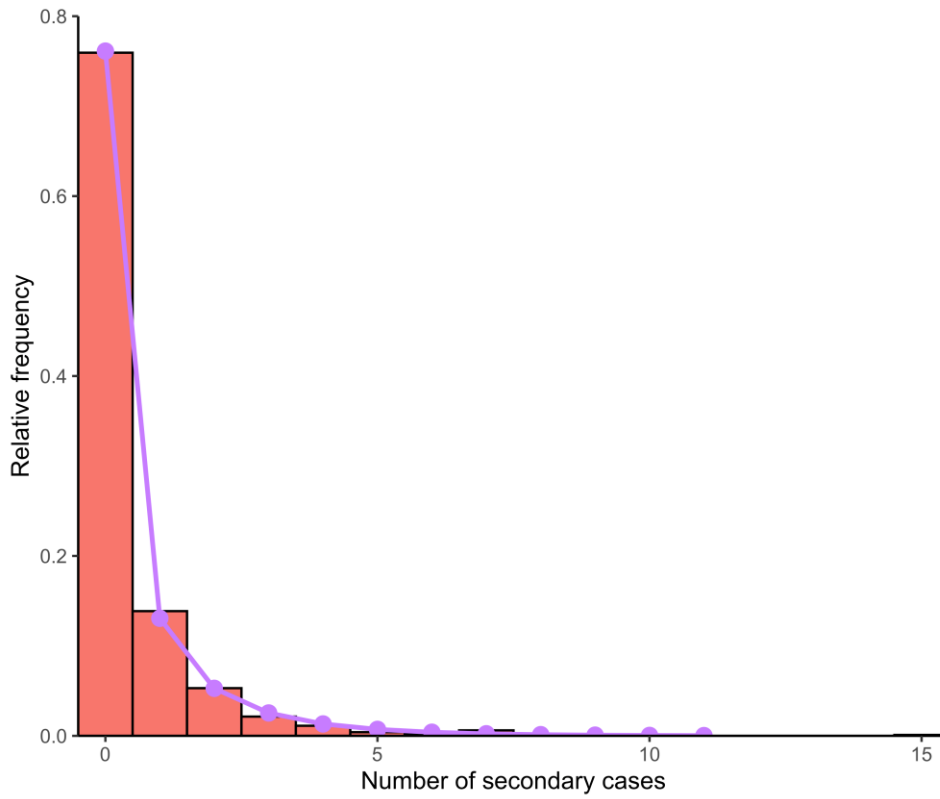
103 Here,  $\lfloor \cdot \rfloor$  denotes the floor function.

### 104 105 *S4.3 Parameter estimation*

106 The parameters of negative binomial distribution were estimated by using the Metropolis-  
107 Hastings algorithm with noninformative prior distributions, which is a Markov chain Monte Carlo  
108 (MCMC) method. The marginal posterior distribution was obtained from 50000 iterations, among  
109 which the first 10000 samples were discarded as for burn-in. The convergence of each MCMC chain  
110 was checked by using the trace plot and Gelman-Rubin-Brooks convergence diagnostic [5]. The  
111 median and the 95% credible interval (CrI) were obtained from the marginal posterior distributions.

112 We compared the fitting performance of negative binomial (NB) distribution to that of  
113 Poisson distribution (i.e., setting  $k$  as infinity) and that of Geometric distribution (i.e., setting  $k = 1$ )  
114 by using the deviance information criterion (DIC), and found that NB distribution had a relatively  
115 lower DIC value (data not shown).

116

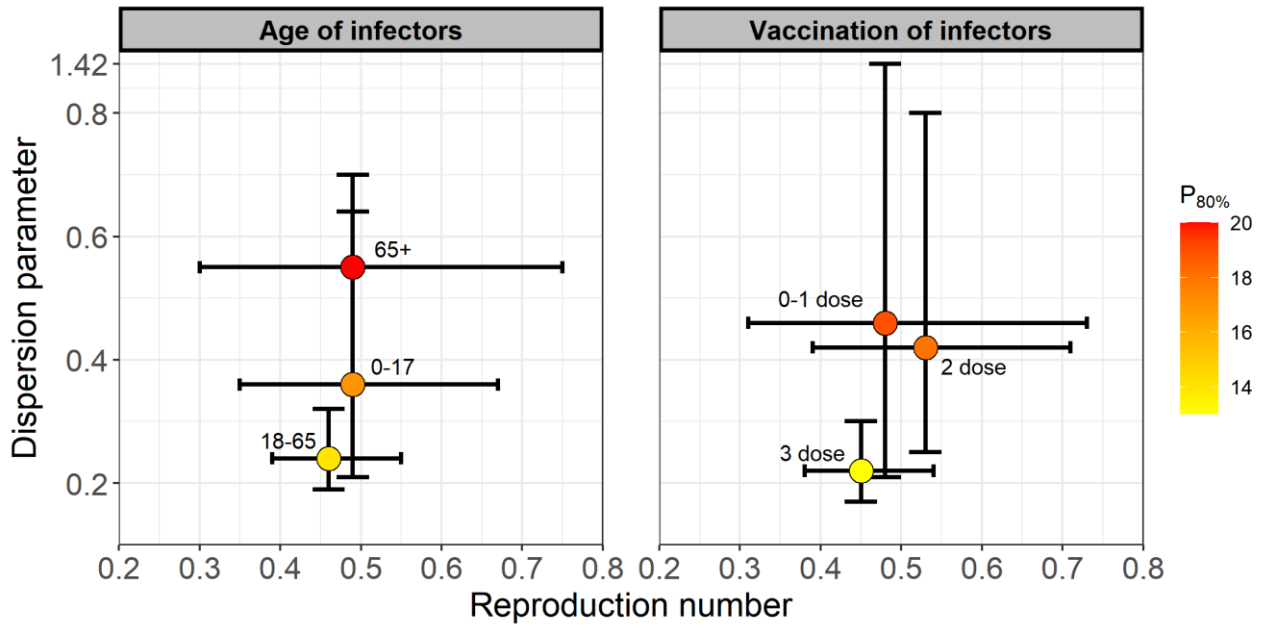


118

119 **Figure S1.** The observed secondary case distribution (red blocks), and the fitted negative binomial  
120 distribution (purple curve). The purple solid curve represented the median of MCMC posterior  
121 samples.

122

123



124

125 **Figure S2.** The estimated effective reproduction number ( $R$ ) and dispersion parameter ( $k$ ) stratified  
 126 by age groups and vaccine doses. The solid circles denoted mean estimates and the horizontal and  
 127 vertical error bars denoted 95% CrIs of  $R$  and  $k$ , respectively. The gradient color of middle dots  
 128 denoted the proportion (%) of the most infectious cases that seeded 80% transmissions.

129

130

## 131 S5 Modelling heterogeneity in transmission risk

### 132 S5.1 Distribution of secondary attack ratio

133 The secondary attack ratio (SAR) is defined as the probability that infections occur among  
134 people who exposed to the infectors. We assumed the number of secondary cases  $k_i$  out of the total  
135 number close contacts  $n_i$  of an (randomly-selected) infector case  $i$  followed a binomial  
136 distribution, conditioning on the SAR,  $P_i$ . To take account of the individual variations in SAR, we  
137 assumed the  $P_i$  followed a Beta distribution that parametrized by two shape parameters  $\alpha$  and  $\beta$ ,  
138 and then the mean of  $P_i$  was  $\frac{\alpha}{\alpha+\beta}$ . This would thus yield a beta-binomial distribution for  $k_i$ . The  
139 probability mass function of the beta-binomial distribution for  $k_i$  is given by [6]:

$$140 \quad f = \frac{\binom{n_i}{k_i} \text{Be}(k_i + \alpha, n_i - k_i + \beta)}{\text{Be}(\alpha, \beta)}$$

141 where  $x$  denotes the number of secondary cases and size denotes the number of close contacts, and  
142  $\text{Be}(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$  the beta function (which is not Beta distribution) at  $\alpha$  and  $\beta$ .  
143 We constructed the likelihood function over all cases as follows:

$$144 \quad L_{\text{SAR}} = \prod_{i=1}^N \frac{\binom{n_i}{k_i} \text{Be}(k_i + \alpha, n_i - k_i + \beta)}{\text{Be}(\alpha, \beta)}$$

145 Here, parameters  $\alpha$  and  $\beta$  for the beta distribution were to be estimated, but the counts of  $k_i$  test-  
146 positive contacts out of a total number of close contacts  $n_i$  were both known from real-world  
147 observations.

148

### 149 S5.2 Parameter estimation

150 The parameters of beta-binomial distribution were estimated in a Bayesian statistical  
151 framework by applying the Markov chain Monte Carlo (MCMC) method. The Metropolis-Hastings  
152 algorithm with noninformative prior distributions were used, each marginal posterior distributions  
153 was obtained from 50000 iterations, among which the first 10000 samples were discarded as for  
154 burn-in. The convergence of each MCMC chain was checked by using the trace plot and Gelman-  
155 Rubin-Brooks convergence diagnostic [5]. The median and the 95% credible interval (CrI) were  
156 obtained from the marginal posterior distributions.

157 We compared the fitting performance of beta-binomial distribution to that of binomial  
158 distribution (i.e., treating SAR as a constant among source cases) by using the deviance information  
159 criterion (DIC), and found that beta-binomial distribution had a relatively lower DIC value (data not  
160 shown).

161

162

### 163 S5.3 Supplementary results

164 **Table S1.** Summary of the estimated secondary attack ratio (SAR) stratified by age groups, vaccine  
165 doses, contact settings, and epidemic period.



Stratifications	Sample size	Mean (%)	SD (%)	95% percentile (%)
<b>Overall<sup>#</sup></b>	769	6.5 (4.9, 8.6)	15.0 (12.0, 19.0)	41.0 (30.0, 56.0)
<b>Sex</b>				
Male	313	4.3 (3.2, 6.1)	9.1 (6.8, 13.0)	24.0 (16.0, 36.0)
Female	456	8.2 (6.0, 11.0)	19.0 (15.0, 24.0)	57.0 (39.0, 78.0)
<b>Age</b>				
0-17	122	14.0 (9.0, 20.0)	26.0 (18.0, 33.0)	82.0 (54.0, 97.0)
18-65	599	5.0 (3.8, 7.0)	13.0 (10.0, 16.0)	33.0 (23.0, 45.0)
>65*	48	9.1 (3.1, 17.0)	24.0 (6.0, 33.0)	NA
<b>Type of index cases</b>				
Symptomatic	44	10.4 (4.9, 21.0)	19.0 (9.5, 31.0)	57.0 (25.0, 95.0)
Asymptomatic	725	6.0 (4.2, 8.0)	14.0 (10.6, 18.0)	38.0 (26.0, 52.0)
<b>Vaccine dose of index cases</b>				
0-1	79	9.6 (5.4, 17.0)	19.0 (12.0, 30.0)	58.0 (31.0, 93.0)
2	159	9.2 (5.3, 14.0)	19.0 (11.0, 26.0)	56.0 (29.0, 81.0)
3	531	5.0 (3.6, 7.0)	12.0 (9.0, 16.0)	32.0 (22.0, 45.0)
<b>Type of contact setting<sup>§</sup></b>				
Household	515	21.0 (18.0, 24.0)	30.0 (27.0, 34.0)	92.0 (84.0, 97.0)
Community	196	5.3 (3.2, 8.4)	20.0 (14.0, 26.0)	50.0 (9.2, 98.0)
Workplace	203	3.5 (2.0, 6.2)	9.0 (4.0, 17.0)	21.0 (10.0, 44.0)
Unknown	689	2.2 (1.4, 3.4)	7.8 (5.0, 12.0)	14.0 (8.2, 24.0)
<b>Epidemic period</b>				
Before lockdown	317	2.0 (1.8, 3.9)	6.2 (4.1, 9.8)	15.0 (9.6, 23.0)
After lockdown	452	9.8 (7.0, 13.0)	21.0 (16.0, 24.0)	63.0 (44.0, 78.0)

166 # The sample size of here was calculated as  $(1139 - 370 =) 769$  index cases, where 1139 was the  
167 total number of confirmed COVID-19 cases and 370 was the number of index cases with 0  
168 associated contacts (0 contact). Those 370 index cases were excluded from statistical analyses of  
169 SAR.

170 \* The estimates of posterior MCMC samples did not converge, which might be due to a relatively  
171 small sample size, and thus the mean and SD were summarized as the sample mean and sample SD,  
172 respectively, with 1000 runs of bootstrap.

173 § The summation of sample sizes in different contact settings was larger than the overall sample size,  
174 i.e.,  $515 + 196 + 203 + 689 > 769$ . This was because some index cases had close contacts in more  
175 than one contact settings, and thus the SARs of such index cases were calculated separately. Besides,  
176 for each contact setting, the sample size was smaller than the overall sample size, i.e.,  $515 < 769$ ,  $196$   
177  $< 769$ ,  $203 < 769$ , and  $689 < 769$ . This was because an index case would not be counted in a contact  
178 setting, if this index case has 0 close contact in this contact setting.

179

180

181 **Supplementary references**

- 182 1. He X, Lau EHY, Wu P, Deng X, Wang J, Hao X, Lau YC, Wong JY, Guan Y, Tan X *et al*:  
183 **Temporal dynamics in viral shedding and transmissibility of COVID-19.** *Nature*  
184 *Medicine* 2020, **26**(5).
- 185 2. Hu S, Wang W, Wang Y, Litvinova M, Luo K, Ren L, Sun Q, Chen X, Zeng G, Li J *et al*:  
186 **Infectivity, susceptibility, and risk factors associated with SARS-CoV-2 transmission**  
187 **under intensive contact tracing in Hunan, China.** *Nat Commun* 2021, **12**(1):1533.
- 188 3. Lloyd-Smith JO, Schreiber SJ, Kopp PE, Getz WM: **Superspreading and the effect of**  
189 **individual variation on disease emergence.** *Nature* 2005, **438**(7066):355-359.
- 190 4. Endo A, Abbott S, Kucharski AJ, Funk S, Centre for the Mathematical Modelling of  
191 Infectious Diseases C, Working G: **Estimating the overdispersion in COVID-19**  
192 **transmission using outbreak sizes outside China.** *Wellcome open research* 2020, **5**:67-67.
- 193 5. Gelman A, Carlin JB, Stern HS, Dunson DB, Rubin DB: **Bayesian data analysis, third**  
194 **edition.** *Journal of the American Statistical Association* 2003, **45**(2).
- 195 6. Prentice RL: **Binary Regression Using an Extended Beta-Binomial Distribution, with**  
196 **Discussion of Correlation Induced by Covariate Measurement Errors.** *Journal of the*  
197 *American Statistical Association* 1986, **81**(394):321-327.

198