

Quantifying heterogeneity in SARS-CoV-2 transmission during the lockdown in India

Nimalan Arinaminpathy,¹ Jishnu Das,² Tyler H. McCormick,³
Partha Mukhopadhyay,⁴ Neelanjan Sircar^{4,5}

¹MRC Centre for Global Infectious Disease Analysis, Imperial College
School of Medicine.

²McCourt School of Public Policy and the Walsh School of Foreign Service,
Georgetown University

³Departments of Statistics and Sociology, University of Washington

⁴Centre for Policy Research, New Delhi, India

⁵Ashoka University, Sonipat, India

The novel SARS-CoV-2 virus shows marked heterogeneity in its transmission. Here, we used data collected from contact tracing during the lockdown in Punjab, a major state in India, to quantify this heterogeneity, and to examine implications for transmission dynamics. We found evidence of heterogeneity acting at multiple levels: in the number of potentially infectious contacts per index case, and in the per-contact risk of infection. Incorporating these findings in simple mathematical models of disease transmission reveals that these heterogeneities act in combination to strongly influence transmission dynamics. Standard approaches, such as representing heterogeneity through secondary case distributions, could be biased by neglecting these underlying interactions between heterogeneities. We discuss implications for policy, and for more ef-

1 **efficient contact tracing in resource-constrained settings such as India. Our re-**
2 **sults highlight how contact tracing, an important public health measure, can**
3 **also provide important insights into epidemic spread and control.**

4 **Introduction** There is increasing recognition of pronounced heterogeneity in the transmis-
5 sion of SARS-CoV-2: that is, that the majority of transmission events appear to be caused only
6 by a small proportion of infected individuals (1–4). Previous modelling work has highlighted
7 the importance of heterogeneity in the emergence of novel pathogens (5), as well as its impli-
8 cations for herd immunity to SARS-CoV-2 (3, 6). Understanding heterogeneity can also have
9 important implications for control, if interventions can be targeted at those most likely to con-
10 tribute to transmission (7). The need to streamline resources in this way is especially pressing
11 in low- and middle-income settings, given fears that healthcare services in these settings would
12 be particularly challenged by SARS-CoV-2 (8). Here, we analysed data collected from contact
13 tracing during the lockdown in Punjab, a major Indian state, to understand heterogeneity of
14 transmission in this setting, and its implications for control.

15 **Epidemiological context** Punjab, a state in India of about 30 million inhabitants, went into
16 lockdown from 1st April to May 26th (Fig.1A). As elsewhere in India, the lockdown heavily
17 restricted the movement of populations, in most cases to their homes and immediate neighbor-
18 hoods. Travelling outside the house required a special pass, except for essential activities which
19 were also restricted to certain times of the day. The Government of Punjab conducted intensive
20 contact tracing during this time, amongst all known contacts of positive cases, and regardless of
21 symptom status. Due to the ease of tracking individuals during the lockdown, 95% of high-risk
22 contacts (defined as those having face-to-face conversation for at least 15 minutes) could be
23 effectively traced and tested. Overall, this data constitutes the census of all infected persons
24 and their contacts in the state; owing to the lockdown conditions, it affords a unique opportu-

1 nity to measure contacts with greater accuracy than would be possible during normal economic
2 activity.

3 The data includes 454 initial cases and 11309 high risk contacts (Fig.1B). Confirmed cases
4 comprise two groups: those residing in Punjab and who were likely infected within the state, and
5 those who are thought to have acquired infection outside the state, due to travel or migration.
6 Our analysis focuses on the former group, and in particular on *seeds* (the first infection in a
7 cluster) in this group, these being the individuals amongst whom contacts are most clearly
8 defined (see Materials and Methods). This yields a total of 148 seeds with 2763 contacts,
9 although we also present sensitivity analysis when analysing all 454 seeds with at least one
10 contact (and all 11309 contacts) in this data, a significant proportion (36%) of whom were
11 religious pilgrims who returned to Punjab from Nanded, Maharashtra, after being stranded there
12 for a month.

13 **Heterogeneity in transmission** The “secondary case distribution” is the distribution for the
14 number of onward infections caused by an infected individual. We observe both the number
15 of secondary cases for each individual, and the total number of contacts the person has. In
16 mathematical modelling of transmission dynamics, heterogeneity in transmission is conven-
17 tionally captured through modelling the secondary case distribution with a negative binomial
18 distribution, allowing for extra-Poisson variation ($1, 5$). Fig. 2A illustrates the secondary case
19 distribution in the data from Punjab. An important feature in this distribution, consistent with
20 earlier findings (4), is that the majority (76%) of infected cases shows no evidence of onward
21 transmission amongst any of their contacts. The negative binomial distribution captures these
22 individuals, as well as the right-hand tail of the distribution, for example the 10% of individu-
23 als accounting for about 80% of transmission in this data. However, this distribution conceals
24 further levels of heterogeneity, that can be important for epidemiological outcomes.

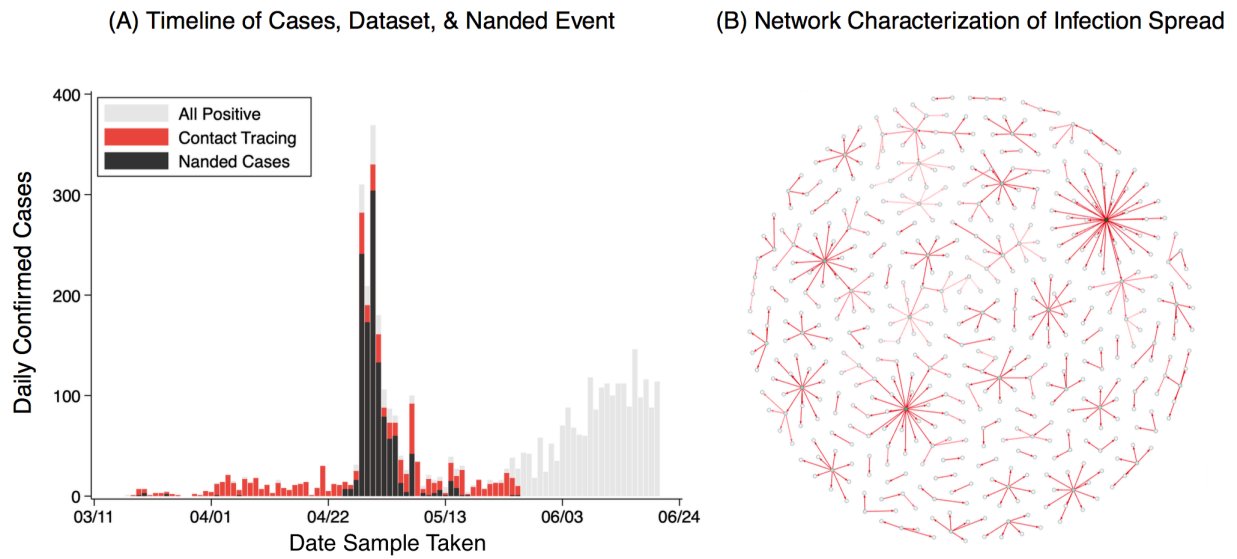


Figure 1: The data from Punjab. (A) Timeseries of reported cases in Punjab during the period of lockdown in the state (red bars) and those due to the Nanded event (black bars), and total cases from early March to the middle of June. (B) Visualisation of case clusters in the dataset, and their linkages from self-reported contacts. This network-type graph requires assumptions (see Materials and Methods). Most individuals infected only few others, while a few infected many: overall, 10% of cases accounted for 80% of infection events.

1 Heterogeneity in transmission can arise from both biological and behavioural factors, in-
2 cluding connectedness (the individuals with the most contacts having the most opportunities
3 for transmission), and individual-level variation in infectiousness (for example, with between-
4 individual and temporal variation in viral shedding (9, 10)). Fig. S2 (supporting information)
5 illustrates the distribution in the number of reported contacts per infected case (the ‘degree dis-
6 tribution’) in our dataset, showing a pronounced right-skew similar to that of the secondary case
7 distribution. However, this skew alone cannot explain the heterogeneity in the secondary case
8 distribution: Fig.2B shows that there are many individuals in this data set who caused no fur-
9 ther infections despite having many contacts (i.e. having ‘high degree’), and conversely many
10 individuals with low- and moderate-degree who caused several onward infections. These data
11 suggest that there is further heterogeneity acting at the individual level, modifying the effect of
12 the degree distribution (see also Fig. S3).

13 To capture this heterogeneity we defined the ‘per-contact infectiousness’ (PCI) as the prob-
14 ability that a given contact results in infection, a probability assumed to vary by index case, but
15 to apply equally to all contacts of a given index case. As shown in Fig.2B, there are several
16 individuals with 1-2 contacts who caused zero onward infections, giving rise to substantial un-
17 certainty in their true PCI (similar challenges apply to low-degree individuals who infected all
18 their contacts). To address this issue we treated PCI as an individual-level effect and estimated
19 it using Bayesian shrinkage, a technique employed (among other places) in the education statis-
20 tics literature to estimate teacher effectiveness (11–13). Fig.2C shows resulting estimates for
21 the marginal distribution of PCI over the population, once again illustrating a strong right-skew.
22 Fig.2D illustrates this association between degree and PCI, showing: (i) a bimodal relationship
23 between the two, arising from the large proportion of individuals that do not infect any others,
24 and (ii) amongst those that do infect others, a negative association between degree and PCI.
25 Overall, these findings illustrate that degree and PCI operate in tandem to drive heterogeneity in

Model number	Description
1	Secondary case distribution using Poisson distribution with mean 1.4
2	Secondary case distribution using Negative Binomial distribution with number of successes = 0.1 and probability of success = 0.067
3	Joint degree/PCI distribution with $\rho = -0.4$
4	Joint degree/PCI distribution with $\rho = -0.2$
5	Joint degree/PCI distribution with $\rho = 0$

Table 1: List of the different models used, for capturing heterogeneity in the population. 'Secondary case distributions' (models 1 - 2) are as in Fig. 2A. They ignore any interactions between degree and PCI, and instead aim only to capture variation in the numbers of secondary cases per index case. By contrast, 'Joint distributions' aim to model the associations shown in Fig. 2D. They employ the bivariate normal distribution described in the Materials and Methods, with correlation ρ .

1 the secondary case distribution. Performing these analyses on the full data for seeds (including
2 returnees as well as the 'core' group) shows qualitatively similar results (see Fig. S4). We next
3 examined the implications of these associations, for transmission dynamics.

4 **Implications for transmission dynamics** We asked: (i) how important are the zero-infectors
5 in these distributions, for epidemiological dynamics? (ii) How do outbreak dynamics compare
6 when taking the conventional approach of using the secondary case distribution alone (Fig.
7 2A) vs when modelling both PCI and degree separately (Fig. 2D)? To address (i), we used
8 the Poisson and negative binomial distributions shown in Fig.2A, the former being an example
9 of capturing the mean secondary cases but failing to capture the proportion zero-infectors. To
10 address (ii), we additionally modelled the log-transformed degree and logit-transformed PCI
11 as following a bivariate normal distribution, with correlation ρ (see Table 1, and Materials and
12 Methods). Consistent with Fig.2D, we assumed a range of values for ρ , from -0.4 to 0.

13 For the transmission model we implemented a simple network simulation, in an assumed

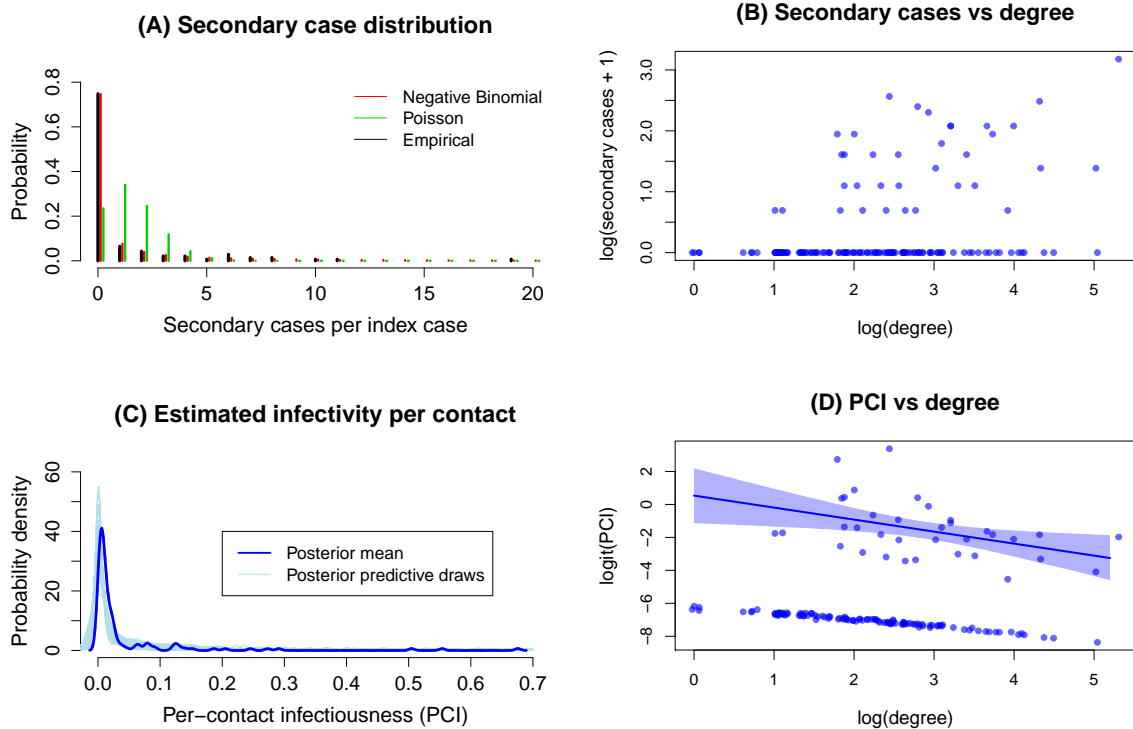


Figure 2: Heterogeneity of the data in secondary cases, and in numbers of contacts. (A) The distribution of secondary cases amongst ‘seeds’ (i.e. first cases in each cluster shown in Fig. 1B). Also shown, for comparison, are the best-fitting Poisson distribution (with $\lambda = 1.4$), and the best-fitting negative binomial distribution (with distribution parameters $r = 0.067, k = 0.1$). The difference between the latter two curves illustrates the strong extra-Poisson variation in the secondary case distribution. (B) Scatter plot of secondary cases vs degree, at the individual level. The secondary case and degree distributions are shown at the logarithmic scale, and adjusted by 1 to account for zeros, to address skewness of the distributions. Although both secondary case and degree distributions show a strong right-skew (panel A), this figure illustrates that the latter does not explain the former: despite a positive relationship between the two distributions, a substantial number of individuals with low degree generate some infections, while many with high degree generate zero onward infections. (C) Estimated marginal density of per-contact-infectiousness (PCI) that, alongside degree, is needed to explain the heterogeneity in secondary cases. Shaded intervals show 95% Bayesian credibility intervals. (D) Estimated PCI vs degree. The figure displays relationship between the logarithm of the odds (logit) of PCI and the logarithm of the degree. These transformations allow us to plausibly model the joint distribution of PCI and degree as a multivariate normal in section 4 (see Materials and Methods and Supporting Information). There is a discernible lower band due to a large number of cases with zero onward infections, which have very low estimated PCI. Among those with onward infections, there is a discernible negative association.

1 population of 3,000 individuals, consistent with the population size in this study. For simplicity
2 and generality, we simulated the epidemic in generations of infection: our simulated outbreak
3 behaviour would thus apply to any emerging infection sharing these heterogeneities (see Mate-
4 rials and Methods).

5 Fig. 3 shows a comparison of model projections for the behaviour of an index case: that
6 is, when simulating only a first generation of infection. Results illustrate how it is possible
7 to accommodate a wide range for R_0 (Fig. 3A), even amongst models that capture a high
8 proportion non-infectors (Fig. 3B).

9 Figs. 3C,D compare the outcomes of full epidemic simulations. By failing to capture the
10 high proportion of zero-infectors (Fig. 3A), a Poisson secondary case distribution yields the
11 most outbreak-prone populations, with 90% of simulations yielding major epidemics (Fig. 3C).
12 Even amongst the remaining models, however, there is a notable disparity in epidemiological
13 outcomes: amongst models capturing the joint distribution between degree and PCI (Models 3
14 - 5), it is not possible to identify a value of ρ that matches most closely to the negative binomial
15 model for secondary cases (Model 2). While the latter appears intermediate to Models 3 and 4
16 in Fig. 2C, it is intermediate to Models 4 and 5 in Fig. 2D.

17 Figure 3E compares selected models in terms of the aggregate temporal pattern that they
18 predict, when aggregated over multiple independent locations. Under epidemics simulated us-
19 ing a Poisson secondary case distribution, there is a surge of infection across several locations at
20 once, a scenario that would place severe demands on health resources. By contrast in outbreaks
21 driven by distributions capturing the high proportion zero-infectors, aggregate epidemic dy-
22 namics are more characterised by a series of asynchronous peaks in different locations, overall
23 making for a lower peak demand on health services.

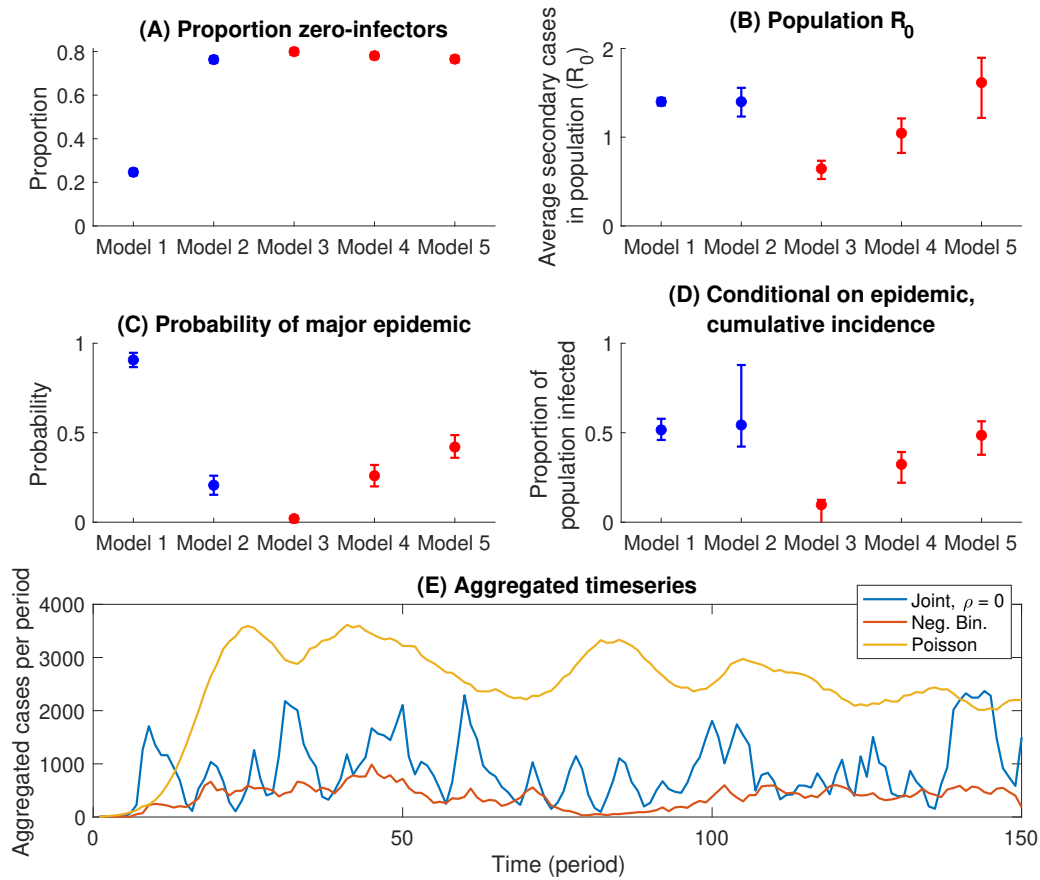


Figure 3: Results of simple transmission models incorporating heterogeneity. Top panels show the average behaviour of an index case in a fully susceptible population of 3,000: (A) The proportion of individuals that cause no further infections. (B) Distributions for the mean number of secondary cases caused by an index case, when averaged over the whole population. In each panel, blue points show outcomes when simulating only secondary case distributions, while red points show outcomes when simulating from the joint degree/PCI distribution described in the main text. Model numbers are as listed in Table 1. Of all models, only the negative binomial secondary case distribution, and the joint degree/PCI models capture the high proportion of index cases who do not cause secondary cases (panel A). However, even amongst these models, there can be substantial variation in R_0 (panel B), owing to the role of correlation between degree and PCI. Middle panels (C,D) show epidemic outcomes over 500 time periods, assuming a 1% probability per time period, of exogenous introduction of an infectious case (here, an ‘epidemic’ is denoted as any simulation having a cumulative incidence > 500 cases (see Materials and Methods for rationale)). Uncertainty intervals arise from repeating simulations 250 times, and reflect 95% simulation intervals. (E) Modelled timecourse of incidence, when aggregated over 250 simulations (with each simulation being interpreted here as an independent location). A Poisson secondary case distribution (in yellow) gives rise to a large surge in aggregate infection because of epidemics in multiple locations occurring in a synchronised way.

1 **Efficiency of contact tracing** Although contact tracing plays an important role in the SARS-
2 CoV-2 response, in resource-constrained settings such as India, its demands on the healthcare
3 system can make it difficult to sustain. Motivated by our findings, we propose reframing contact
4 tracing with the goal of efficiently identifying individuals with high PCI. In our data overall, we
5 estimate that if an individual caused at least one onward infection, there is a 79% probability
6 that they caused at least two onward infections. We thus propose a sequential strategy where,
7 for every index case, a ‘pilot’ subset of only s randomly selected contacts is first tested; the
8 remainder of contacts are then followed up and tested, only if there is a positive in the pilot
9 subset. Such a strategy could substantially reduce the overall contact tracing effort, while still
10 effectively identifying high PCI individuals. Fig. 4 shows results of simulating such a strategy
11 1,000 times on the full dataset of 454 cases, for a range of values of s . The figure illustrates
12 diminishing returns in the fraction of infections found, beyond a pilot subset size of 10 contacts
13 (Fig. 4A). However, even with a pilot subset size of only 5 contacts, it is possible to identify
14 80% of infections (Fig. 4A), with <40% of the contact tracing effort that was expended in this
15 data (Fig. 4B).

16 **Discussion** We have shown how individual-level data, gathered from the routine course of
17 contact tracing, can be analysed to gain important insights into the transmission of SARS-CoV-
18 2. As well as affirming findings from elsewhere, that the majority of cases appear not to infect
19 any others (1, 4), our findings also highlight how heterogeneity in transmission may be more
20 complex than previously recognised.

21 Simple dynamical models highlight the important role that is played by these heterogeneities.
22 At the gross level of the secondary case distribution, the high proportion of zero-infectors yields
23 outbreak dynamics wherein surges can be handled by providing mobile services rather than in-
24 creasing hospital capacity in every geography (Fig. 3E). The negative binomial distribution,

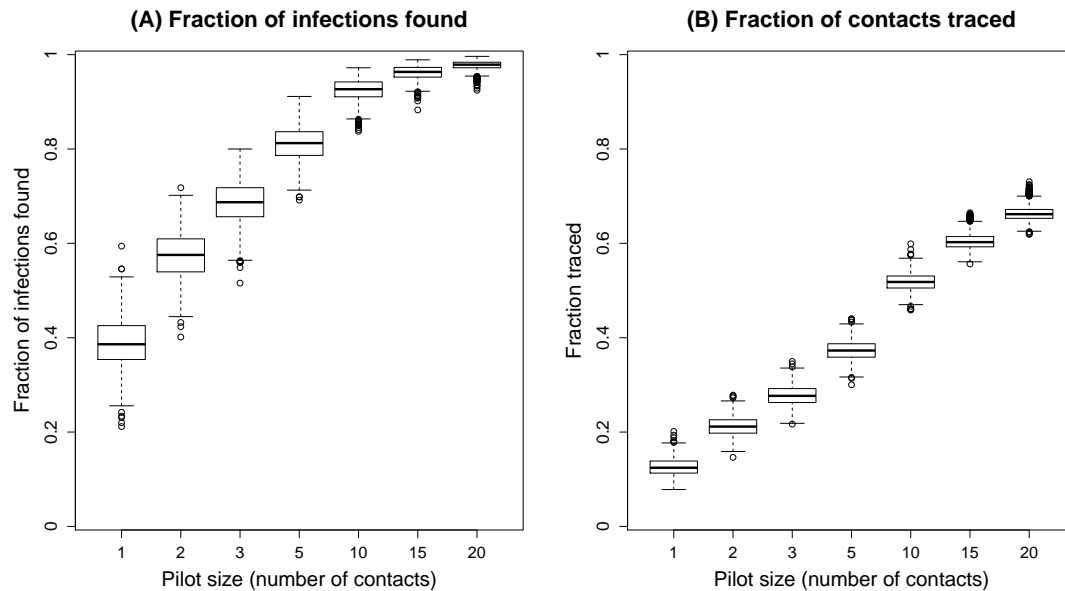


Figure 4: An approach to efficient contact tracing. Figure shows simulated outcomes of a strategy to test all contacts of an index case, only if there is at least one positive individual in an initial ‘pilot’ sample of s contacts. (A) The proportion of infections found as a function of s (B) Overall contact tracing effort, as measured by the proportion of contacts that would be traced, again as a function of s . Owing to the right-skew of the PCI, the left-hand panel illustrates diminishing returns with increasing s , suggesting, for example, that it would be possible to identify 80% of the cases in this dataset, with <40% of the contact tracing effort.

1 conventionally used in the modelling literature, captures this proportion well (Fig. 2A). How-
2 ever, our analysis also highlight some limitations of this distribution: accounting for underlying
3 correlations between degree and PCI can lead to different outbreak dynamics, in terms of the
4 risk and size of major outbreaks over time (Figs. 3C,D). Our results also have implications for
5 the efficiency of contact tracing. When a large fraction of infected individuals do not cause on-
6 ward transmission, we show the value of a simple two-step strategy, of, for example, first testing
7 family members and then testing other contacts only if at least one family member is found to
8 be positive (Fig. 4). Such approaches can be particularly valuable in resource-constrained set-
9 tings such as India, in decreasing the requirements for contact tracing substantially, while still
10 identifying most cases.

11 An important question, that we are not able to address using the current data, is what drives

1 the heterogeneity in per-contact infectiousness. This heterogeneity may arise, for example, from
2 biological factors such as the role of pre-existing, cross-reactive immunity that may moderate
3 viral load in some individuals more effectively than others (14). Our analysis suggests that PCI
4 increases with age and is significantly associated with sex (Fig. S5). Further data on these and
5 other individual-level characteristics would be invaluable in further examining key risk factors
6 for infectiousness. Where risk factors involve individual characteristics that can be readily iden-
7 tified in newly diagnosed patients, such as viral load, these factors could also play an important
8 role in guiding future contact tracing efforts. However, heterogeneity in PCI could also reflect
9 variations in the closeness of reported contacts, with some reporting only the closest contacts
10 and others reporting wider contacts, thus explaining the negative correlation. We emphasise that
11 our data is limited to defined ‘high-risk’ contacts (see Materials and Methods), thus excluding
12 incidental contacts that might be expected to bias our estimates the most. Nonetheless, even
13 if there is variation in closeness amongst high-risk contacts, our analysis offers an approach
14 for adjusting for these variations, when interpreting what routine contact tracing data means for
15 transmission: in this case our estimates for PCI should be regarded as a data-driven weighting of
16 contacts, rather than infectiousness. Our approach can easily be adapted for any dataset where
17 there is additional information on closeness of contact.

18 Amongst other limitations, the contact tracing data was collected, not under controlled study
19 conditions, but as part of a public health response, by the Government of Punjab. Our approach
20 to this data is pragmatic, recognising some inherent limitations: there may be false negatives
21 in the data if people were tested too early, or indeed if people had been infected long in the
22 past, which we cannot tell in the absence of serological tests. As with any contact tracing data,
23 our assumptions for who-infected-whom, in a given contact pair, may be imperfect. We are
24 able to address some of these concerns (for instance, by showing that our results are robust
25 to a change in the directionality of a link [see Materials and Methods]). However, more—and

1 better–data are absolutely necessary to refine our estimates, particularly on the nature of the
2 correlation between degree and PCI. Further, although the lockdown conditions facilitate an in-
3 depth analysis of transmission amongst contacts, our findings must be interpreted with caution
4 in scenarios with uninhibited transmission, as might occur in the absence of a lockdown or other
5 non-pharmaceutical interventions (15). Additional limitations on the modelling are described
6 in the Materials and Methods.

7 Overall, the methods that we have outlined here should apply to any contact tracing database
8 and our publicly available code can be directly applied to any such data that have been collected.
9 Contact tracing forms an integral part of the response to SARS-CoV-2 around the world: while
10 being an important public health strategy in its own right, it can also provide invaluable informa-
11 tion about how, and to whom, infection is being spread. Systematic analysis of this data could
12 provide important insights to inform future, smarter strategies for the control of SARS-CoV-2.

1 **References**

- 2 1. A. Endo, Centre for the Mathematical Modelling of Infectious Diseases COVID-19 Work-
3 ing Group, S. Abbott, A. J. Kucharski, S. Funk, *Wellcome Open Research* **5**, 67 (2020).
- 4 2. Y. Wang, P. Teunis, *Frontiers in Medicine* **7**, 329 (2020).
- 5 3. M. G. M. Gomes, *et al.*, *MedRxiv* p. 2020.04.27.20081893 (2020).
- 6 4. R. Laxminarayan, *et al.*, *medRxiv* (2020).
- 7 5. J. O. Lloyd-Smith, S. J. Schreiber, P. E. Kopp, W. M. Getz, *Nature* **438**, 355 (2005).
- 8 6. T. Britton, F. Ball, P. Trapman, *Science* **369**, 846 (2020).
- 9 7. J. Wallinga, M. van Boven, M. Lipsitch, *Proceedings of the National Academy of Sciences*
10 **107**, 923 (2010).
- 11 8. P. G. T. Walker, *et al.*, *Science* **369**, 413 (2020).
- 12 9. Y. Fu, *et al.*, *European Respiratory Journal* (2020).
- 13 10. L. Qi, *et al.*, *International Journal of Infectious Diseases* **96**, 531 (2020).
- 14 11. R. Mendro, *et al.*, *Annual Meeting of the American Educational Research Association, San*
15 *Diego, CA* (1998).
- 16 12. J. Lockwood, *et al.*, *Journal of Educational Measurement* **44**, 47 (2007).
- 17 13. J. Lockwood, D. F. McCaffrey, L. T. Mariano, C. Setodji, *Journal of Educational and*
18 *Behavioral Statistics* **32**, 125 (2007).
- 19 14. N. Le Bert, *et al.*, *Nature* (2020).

- 1 15. Report 9: Impact of non-pharmaceutical interventions (npis).
- 2 16. Ministry of Health & Family Welfare, India, Guidance document for POEs, states and UTs
3 for surveillance of 2019-nCoV, *Tech. rep.* (2020).
- 4 17. A. Gelman, *et al.*, *Bayesian Data Analysis, Third Edition*, Chapman & Hall/CRC Texts in
5 Statistical Science (Taylor & Francis, 2013).
- 6 18. A. Gelman, J. Hill, *Data Analysis Using Regression and Multilevel/Hierarchical Models*
7 (Cambridge University Press, 2006).
- 8 19. R. Verity, *et al.*, *The Lancet Infectious Diseases* (2020).
- 9 20. Y. Liu, *et al.*, *The European respiratory journal* **55**, 2001112 (2020).
- 10 21. L. Danon, *et al.*, *Interdisciplinary perspectives on infectious diseases* **2011**, 284909 (2011).

11 **Acknowledgments**

12 NA was supported by the UK Medical Research Council and by the Bill and Melinda Gates
13 Foundation. PM is grateful for support from a Ford Foundation grant, which supports work on
14 use of tacit knowledge in urban environments. TM was supported by the National Institute of
15 Mental Health of the U.S. National Institutes of Health under award number DP2MH122405.
16 All authors contributed equally, and are listed in alphabetical order. The authors gratefully
17 acknowledge Dr. Rajesh Bhaskar, who was responsible for the production of the dataset and
18 generously shared it for analysis, Dr. Rajesh Bhatia, Dr. K K Talwar, advisors to Government
19 of Punjab, Ms. Vini Mahajan, Ms. Isha Kalia, Ms. Tulika Avni Sinha, and Ms. Yamini
20 Aiyar. The authors also acknowledge valuable support in collating and organising the data,
21 from Vidisha Mehta, Kanhu Charan Pradhan, Harish Sai, Shamindra Nath Roy (Centre for

- 1 Policy Research), Olivier Telle (Centre for Policy Research/Centre National de la Recherche
- 2 Française), and Benjamin Daniels (Georgetown University).

3 **Supplementary materials**

4 Materials and Methods

5 Figs. S1 to S5

6

1 **Supporting information**

2 **Materials and Methods**

3 **Contact tracing**

4 The contact tracing, implemented by the Integrated Disease Surveillance Program in the state,
5 was conducted in the following four steps (16).

- 6 1. Immediately after a confirmed case is identified, a trained epidemiologist or medical offi-
7 cer interviews the case and ascertains all contacts.
- 8 2. Contact tracing is then completed for all contacts who have interacted with the positive
9 case anytime between 2 days prior to the onset of symptoms and the date of isolation, or
10 a maximum of 14 days after symptom onset. So, if symptoms started on April 1st 2020
11 and the person was isolated on April 5th, all persons who were in contact with the case
12 between March 30th and April 5th are to be traced. The data are listed with details of the
13 contacts and this list is then shared with contact tracers for tracking.
- 14 3. The epidemiologist then classifies each contact as high- or low-risk. The definition of
15 high-risk is those who face-to-face conversations for at least 15 minutes with the positive
16 case or physical contact. Contacts who are out-of-state are passed on to other states.
- 17 4. High-risk contacts are then tested by a lab technician. Contacts who are negative and
18 remain asymptomatic for 28 days are released from the list. For those who are positive,
19 the listing is again initiated to trace a further generation of contacts.

20 In the data, we recoded the reasons for testing into six categories: Seed (Normal), Contact
21 (Normal), Farmer/Labour, Migrant/Returnee (Non-Nanded), Nanded, Other. During the lock-

Category	Fem (%)	Symp (%)	Age Range	Age (Med)	Zero Deg (%)	Total
Seed (Normal)	33	65	0-84	40	5	184
Contact (Normal)	42	17	0-91	35	64	470
Farmer/Labourer	20	9	8-70	27	10	35
Returnee (Nanded)	39	11	0-100	45	83	1270
Returnee (Others)	18	14	1-78	33	26	90
Other	26	29	1-72	35	34	123

Table 2: Key Sample Descriptives. We provide sample descriptives for six reasons for testing: Seed (Normal), Contact (Normal), Farmer/Labour, Migrant/Returnee (Nanded), Migrant/Returnee (Others)[those not in the Nanded event], Other. We provide information on the percentage of people in each category who are female (Fem), show symptoms (Symp), and report no high risk contacts (Zero Deg). We also display median age (Age (Med)) and the minimum to maximum age (Age Range) for each category. Missing observations are removed in this table.

1 down, there was a fear that those entering from elsewhere would bring the infection to Punjab.
2 Pilgrims returning from Nanded (described in section 2), other migrants/returnees, and farm-
3 ers/laborers (residing in Punjab but originally from outside it) were tested by a special protocol.
4 The “other” category consisted of certain high-risk populations like frontline healthworkers who
5 were tested for occupational reasons and their families. The categories “Seed (Normal)” and
6 “Contact (Normal)” correspond to those tested due to normal protocol – usually due to symp-
7 toms, living in a containment zone, or from the contact-tracing protocol described above – as
8 the first case in a cluster or the contact of a confirmed case, respectively.

9 Table 2 displays the frequency of each reason and the percentage reporting no high risk
10 contacts, with distributions of gender, age and symptomatic status. As we might expect laborers
11 and migrants tend to be younger and more male. Of particular concern is the high percentage
12 of individuals reporting no high risk contacts in most categories (Zero Deg). This is due to a
13 bookkeeping problem. For contacts of a previously confirmed case we can only retrieve the
14 number of high risks contacts that had yet to be tested — so those contacts shared with the
15 person from whom the infection was contracted are not counted. And returnees from elsewhere

1 often have no contacts listed in Punjab.

2 Accordingly, in the main text we restrict our analyses in the text to the 148 seeds tested
3 during normal protocol that do not have a missing of zero value in the number of high risk
4 contacts, as contacts are required to estimate PCI. Naturally, because they have come through
5 the normal protocol, seeds have a much higher proportion of symptomatic individuals (Symp).
6 This is a population for whom we believe we have a robust contact distribution applicable to
7 the population of Punjab and for whom we can reliably identify seeds and contacts.

8 Where such information could be ascertained, we undertook an extensive exercise of match-
9 ing contacts to seeds in the entire dataset, to verify the dataset had been coded correctly.
10 Nonetheless, we were concerned about the *case ascertainment* problem. Assume that both
11 *A* and *B* have tested positive. *B* could have infected *A* (directly or indirectly) but we observed
12 *A* first and coded it as a seed. Two possibilities exist either we got it right and *A* is the seed,
13 or we got it wrong and *B* is the seed. While we can never be sure of this answer, we can test
14 the robustness of our claims to swapping seeds and contacts. While we cannot direct compare
15 onward infections and degree of contacts to seeds due to the bookkeeping problem, we can test
16 whether seeds and contacts display similar infectiousness. Indeed, we see that that the contacts
17 (2763 individuals) of normal seeds have an aggregate test positivity of 6.0% while the contacts
18 (1885 individuals) of normal “contacts” have a test positivity of 6.5%. These are statistically
19 indistinguishable ($p = 0.45$). Thus, as long as the coding of seed and contact by the govern-
20 ment of Punjab was independent of degree (which is likely because seeds were typically tested
21 due to a biological criterion – showing symptoms – and not a social criterion), we surmise that
22 our estimates of the secondary case distribution are likely to be robust to swapping seeds and
23 contacts.

24 Beyond the core dataset, as table 2 shows, the majority of positive cases are from the Nanded
25 event, from which pilgrims were brought back on dedicated buses. This group is akin to the

1 Diamond Princess experience, where multiple people were in contact with each other in close
2 quarters. As such, seeds are contacts were not well-defined in this population.

3 Nevertheless, in Fig S4 in the supporting information, we also show robustness when ex-
4 tending our analysis to the 454 seeds across the Nanded event and all other categories (as seeds
5 are also present in each of the four special protocol categories) who report at least one contact.

6 **Bayesian shrinkage**

7 A natural estimate of PCI for person i , p_i would be to divide the number of onward infections
8 (z_i) by the number of contacts (d_i), i.e., $\hat{p}_i = \frac{z_i}{d_i}$. The shape of the degree distribution presents a
9 challenge for this method, however, as it means the variability in the estimated PCI varies across
10 individuals based on their number of contacts. As an example, consider two individuals, A and
11 B , with 2 and 100 contacts, respectively who have infected no one. We are confident that B has
12 a PCI close to zero but not so confident with A due to a small sample size. We address this issue
13 through *Bayesian shrinkage*. In this setting, individual estimates of PCI (p_i) from high contact
14 individuals (such as B) will be mostly unchanged while those from lower contact individuals
15 (such as A) will be shrunken towards the overall mean (17).

16 Amongst different ways of performing Bayesian shrinkage (e.g., the Beta-Binomial model),
17 we chose to model the logarithm of the odds (logit) of the PCI as following a normal distribution
18 with a common mean and variance, as this functional form is closely linked to our modelling
19 of transmission dynamics (other approaches yield similar estimates (18)). In particular, we
20 estimate:

$$p_i = \text{logit}^{-1}(\alpha_i) = \frac{1}{1 + e^{-\alpha_i}} \quad (1)$$

$$\alpha_i \sim \text{Normal}(\bar{\alpha}, \sigma_\alpha^2), \quad (2)$$

21 where α_i is the log-odds of p_i and $\bar{\alpha}$ is the common mean. Above, σ_α^2 is inversely correlated to
22 the amount of shrinkage. As $\sigma_\alpha^2 \rightarrow 0$, each α_i is given the same value, so each p_i is estimated as

1 the mean infection rate. As $\sigma_\alpha^2 \rightarrow \infty$, $p_i \approx \frac{z_i}{d_i}$. In practice, the “hyperparameters” like σ_α^2 and
2 $\bar{\alpha}$ are estimated using Markov Chain Monte Carlo (MCMC) methods with diffuse priors. The
3 non-zero values of σ_α^2 and $\bar{\alpha}$ guarantee that our estimated p_i is between 0 and 1.

4 In future work with additional data, it may be possible to characterize the complete joint
5 distribution of PCI and number of contacts, thus alleviating the need to shrink to a common
6 mean. For example, healthcare workers with training in mitigating infectious disease spread
7 may have many contacts, but lower PCI than would be expected from the rest of the population.
8 In such a setting, shrinking towards a single mean would underestimate the heterogeneity in
9 the PCI distribution. Most of the extreme cases are 0’s in the data from Punjab, however, so in
10 practice these values will be shrunk towards a small, non-zero value under either model.

11 **Mathematical modelling of transmission dynamics**

12 We implemented a simple network simulation, in an assumed population of 3,000 individuals
13 (consistent with the population size in this study). For simplicity we modelled all networks
14 as random, that is, neglecting clustering and other forms of network structure of higher order
15 than the degree distribution. Also for simplicity, we simulated the epidemic in terms of genera-
16 tions of infection, rather than in continuous time: our projections could be interpreted as being
17 conducted in discrete time, with a time interval corresponding to the mean generation time.
18 The focus of this modelling analysis is to understand the importance of degree distribution and
19 PCI for transmission dynamics in general; we thus did not model the details of symptomatic
20 vs asymptomatic infection for SARS-CoV-2, nor of the pronounced variation of severity by
21 age (19).

22 **Network construction** For the Poisson and negative binomial secondary case distributions
23 in Table 1, we drew 3,000 samples. We then constructed a random, directed network treating

1 these samples as degrees, to construct a network of the secondary cases that any given individual
2 would cause, once themselves infected (our results in Fig. 3 are qualitatively unchanged when
3 assuming a directed network instead).

4 In figure S2, we show that the degree distribution in the data follows an approximately log
5 normal distribution, and in our discussion of Bayesian shrinkage (above) we showed that the
6 logarithm of the odd (logit) of PCI is constructed to follow a normal distribution. We note fur-
7 ther that Fig. 2D and Fig. S4(D) show that the correlation between the log-transformed degree
8 (n) and logit-transformed PCI (p) is plausibly negative for those that infect others. We consider
9 a population of infected individuals. Since the log-transformed degree and logit-transformed
10 PCI each follow a normal distribution and may be correlated, the natural choice for the joint
11 degree/PCI distribution is to model the log-transformed degree and logit-transformed PCI as
12 following a bivariate normal distribution:

$$\begin{pmatrix} \log\left(\frac{p}{1-p}\right) \\ \log(n) \end{pmatrix} \sim N_2(\mu, \Sigma) \quad (3)$$

13 where μ is a vector composed of the mean values for logit-transformed PCI and log-transformed
14 degree, and we have, for the covariance matrix Σ :

$$\Sigma = \begin{bmatrix} \sigma_{lp}^2 & \rho\sigma_{lp}\sigma_{ln} \\ \rho\sigma_{lp}\sigma_{ln} & \sigma_{ln}^2 \end{bmatrix} \quad (4)$$

15 where σ_{lp} is the standard deviation for the logit-transformed PCI; σ_{ln} is the standard deviation
16 for the log-transformed degree; and ρ is the correlation between the two. This construction
17 allows us to explore different hypothetical scenarios for the correlation, and their implications
18 for outbreak dynamics, while maintaining the correct shapes of the marginal distributions for
19 degree and PCI. We posed three scenarios for ρ , taking values of -0.4, -0.2 and 0. The model
20 we use for simulation differs from our Bayesian shrinkage approach presented in Fig. 2D. This
21 distinction is necessary because, in our simulations, we wish to explore variation across multiple

1 similar contact and PCI distributions. In any data analysis, however, we will model *conditional*
2 on a particular observed contact distribution.

3 In a given simulation, we then sampled 3,000 values for degree and PCI. We constructed a
4 random, undirected network from the degree distribution. For each individual m , we assumed
5 that the sampled PCI $p(m)$ applies uniformly to all of their contacts. Thus, although the link
6 between any two individuals A and B is undirected - representing a bidirectional transmission
7 risk - the transmission intensity is not necessarily the same in both directions, and depends on
8 the respective PCIs of A and B (owing to between-individual variations in infectivity).

9 **Epidemic simulation** For a given population constructed as above, suppose C_t is the set of
10 individuals that are infective at the beginning of time-step t ; S_t is the set of individuals that have
11 not yet had infection and are therefore susceptible; and J_t is the set of individuals that are newly
12 infected in timestep t . Further, suppose that $p(m)$ is the sampled PCI for individual m . Then
13 we proceeded along the following iterative steps:

14 While $t \leq 500$ and S_t, C_t both have at least one member:

- 15 1. Identify C_t with J_{t-1} , and initialise J_t as an empty set
- 16 2. For every member m of C_t :
 - 17 • Determine all contacts of m who belong to S_t (for Models 1,2, regarding ‘contacts’
18 as secondary cases).
 - 19 • For each such contact, conduct a Bernoulli trial with probability $p(m)$, to determine
20 whether infection occurs (for Models 1,2, taking $p(m) = 1$).
 - 21 • Accumulate all new infections thus occurring in J_t , and remove them from S_t .

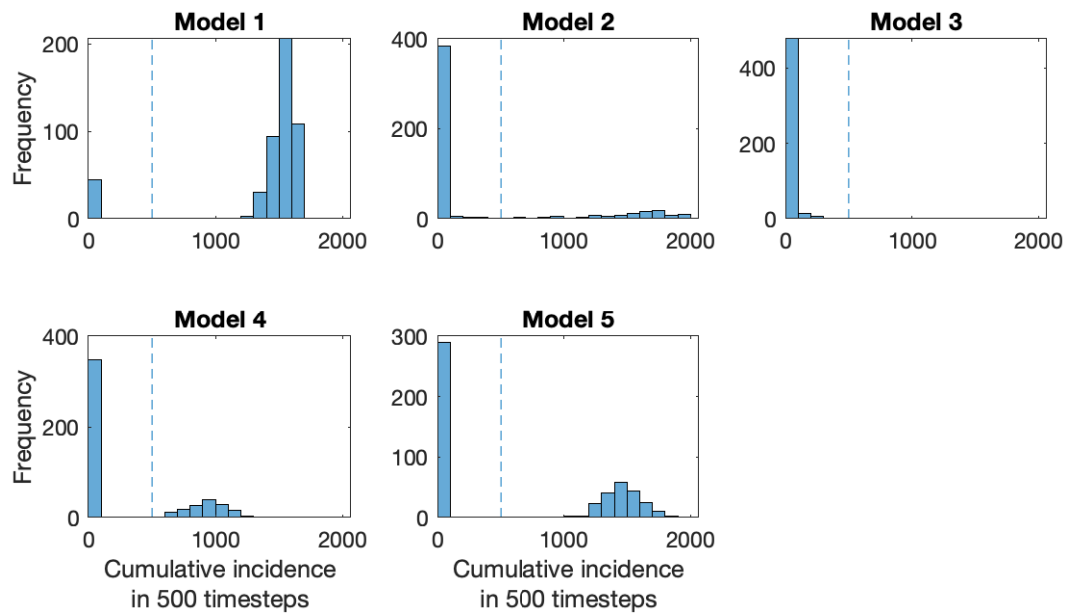


Figure S1: Frequency distributions for cumulative incidence, over 500 timesteps, for each of the models listed in Table 1 in the main text. Vertical dashed lines indicate a cumulative incidence of 500, a consistent dividing line between the two modes in these distributions. Thus Fig. 3C in the main text shows the probability mass to the left of this line, while Fig. 3D shows the mean cumulative incidence to the right of this line.

- 1 3. Perform a Bernoulli trial on all members of S_t with probability 0.01, to identify exoge-
- 2 nous introductions of infection. Accumulate all new infections thus occurring in J_t , and
- 3 remove them from S_t .

- 4 4. Increment t by 1, and iterate from (1).

5 We repeated this algorithm 250 times, for each of the models listed in Table 1. Figure
6 S1 shows the frequency distributions for the cumulative incidence thus obtained, i.e. the total
7 cases over all 500 timesteps. The figure illustrates a bimodal pattern of outbreak sizes, with the
8 vertical dashed line (at 500 cumulative cases) illustrating a consistent dividing line between the
9 two modes. Accordingly in Figs. 3C,D in the main text, we denote ‘major epidemics’ as any
10 simulation in which cumulative incidence exceeds 500 cases.

1 **Limitations** For simplicity we have adopted a simple network model that models the progres-
2 sion of an epidemic through generations of infection. This simplicity is helpful for focusing the
3 model-based analysis on the specific types of heterogeneity revealed by our analysis of the
4 data. It also has the benefit of generality, showing epidemiological behaviour that would ap-
5 ply to any disease with the same underlying heterogeneity, regardless of the details of natural
6 history. However, an important area for future modelling work would be to incorporate some
7 important characteristics in the natural history of SARS-CoV-2, such as symptom status, age,
8 and the full spectrum of severity of infection (19,20). Such refined models would be important,
9 for example, in translating our simulated dynamics to timescales more specific to SARS-CoV-
10 2. As mentioned above, we also take a simplified approach to the network structure, assuming
11 the simplest case of a random network, and thus ignoring the potential for clustering, or other
12 types of network topology that could be influential in transmission dynamics (21). Further
13 data on the underlying contact structure, including the retention of information on test-negative
14 contacts, would be helpful in addressing these simplifications.

15 **Modelling efficient contact tracing algorithms**

16 We used the following algorithm to perform the simulations for Figure 4. First, we choose s ,
17 the number of ‘pilot’ contacts to be tested. Then, for each of the cases in the data from Punjab
18 take the number of observed cases and contacts and create a vector of 1’s and 0’s to designate,
19 respectively, infected contacts and those not infected. We randomly assign the position of each
20 of the infected cases in the vector and set the number of infected cases equal to the number
21 observed in the data for that case. We then take s samples without replacement and with uniform
22 probability from the constructed vectors. If at least one of the sampled units is positive, we
23 simulate testing the remainder of the contacts (meaning that the number of tests equals the
24 observed number of contacts and the number of infections found equals the observed number

1 of infections). If none of the s sampled contacts are positive we do not do further sampling,
2 meaning the number of infections identified is zero and the number of tests is s . We repeat this
3 exercise across all cases in the data from Punjab 1000 times for each value of s .

4 Our simulation is illustrative, with some caveats to note. First, the procedure is not opti-
5 mized for the number of contacts to test. The problem we address is similar to a bandit prob-
6 lem in the machine learning literature where, as more information about the PCI distribution
7 is available through testing, the number of contacts to test is optimized to maximize the (ex-
8 pected) number of infections found while minimizing the number of tests. We anticipate that
9 this sequential procedure would be challenging to implement in a public health context, partic-
10 ularly in a low resource setting. Instead, we opt for a simple rule that can be implemented with
11 no additional optimization (e.g. testing family members and only testing further if at least one
12 is positive) that can still substantially improve efficiency. We also need information about the
13 joint distribution of PCI and contacts for a fully optimized approach, which we cannot estimate
14 with precision in the data from Punjab. Additionally, we do not consider imperfect tests. False
15 positives would preserve the number of infections found but make the procedure less efficient
16 since some individuals will be tested based on false positives that would not otherwise be tested.
17 False negatives will reduce the number of infections found, though this impact would be miti-
18 gated by the right-skew of the PCI distribution. Finally, this procedure relies on the availability
19 of tests with relatively rapid results and the willingness of individuals to be tested. If the delay
20 between testing and receiving results is too long, then contacts who were not tested in the pilot
21 stage could be infecting others in the population while waiting to be tested.

1 Supplementary Figures

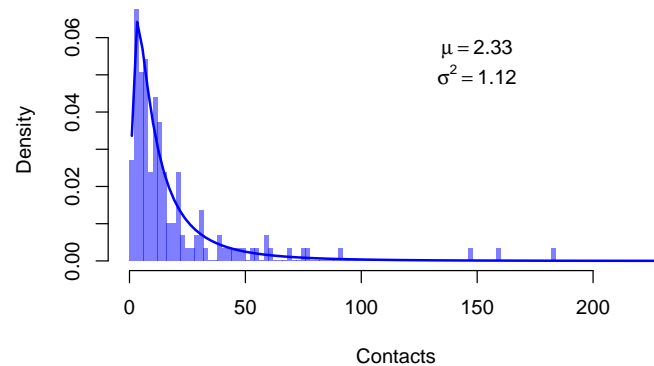


Figure S2: The Contact Distribution. The light blue bars denote the histogram of contact distribution with a bin size of 2. The density function of the log normal distribution with $\mu = 2.33$ and $\sigma^2 = 1.12$ fits the empirical distribution well. Consistent with standard network structure, this distribution has a strong right skew.

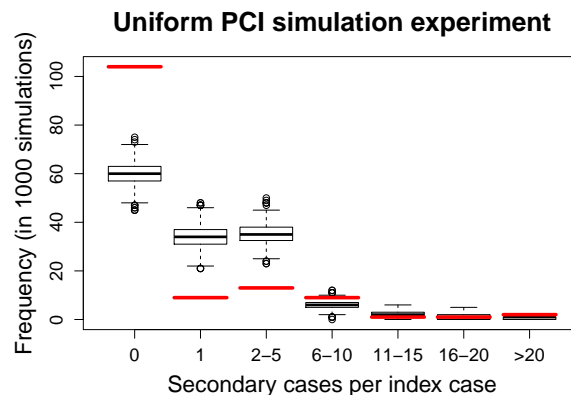


Figure S3: Inadequacy of the degree distribution to explain the secondary case distribution. While Fig. 2B in the main text illustrates this point visually, this plot offers statistical support. Red lines show the data for the secondary case distribution, while box plots show the best-fitting projections when assuming the degree distribution illustrated in Fig. S2, and moreover that the risk-of-transmission is constant across contacts (that is, a constant PCI). Doing so yields a secondary case distribution that severely underestimates the proportion of cases that caused zero onward infections.

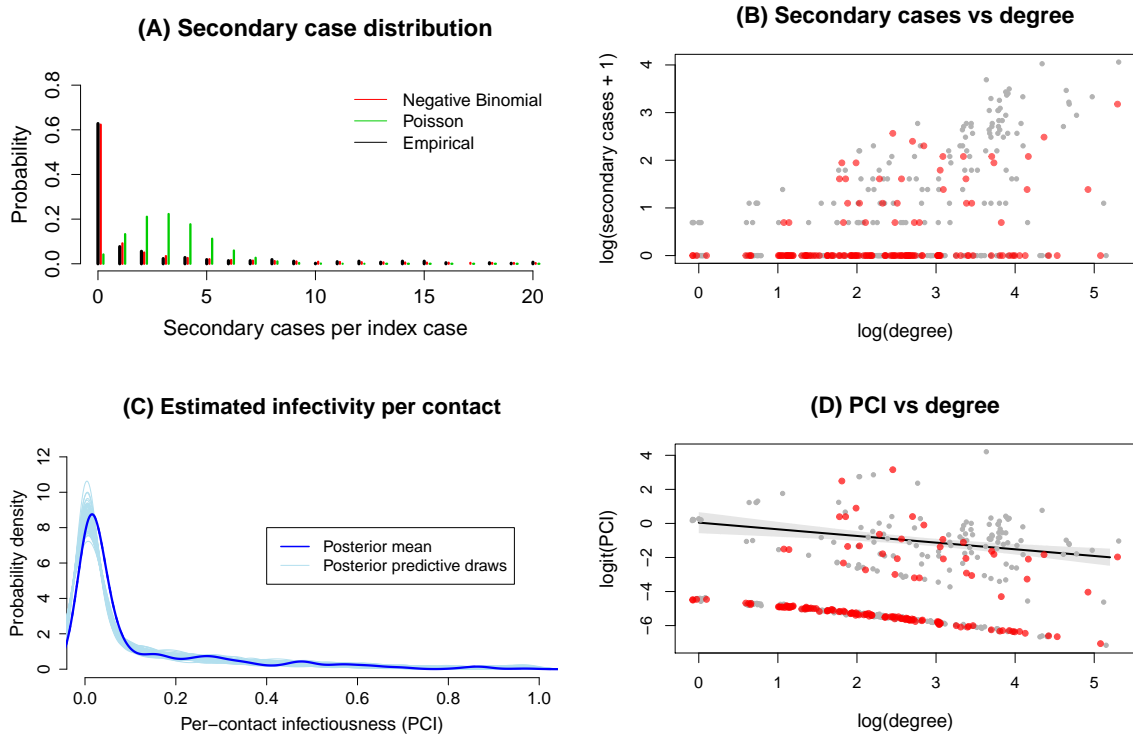


Figure S4: Robustness for full sample. The full heterogeneous sample is characterized in Materials and Methods. While the main text displays analysis on the core sample, the trends hold on the full dataset. Fig. S4(A) replicates Fig. 2A for the full sample, modelling negative binomial and Poisson fits to the secondary case distribution. As before, the negative binomial fits the secondary case distribution well, while the Poisson distribution underestimates the number of zero infectors. Fig. S4(B) shows a scatterplot of the logarithm of the degree distribution against the logarithm of the secondary case distribution adjusted by 1 to account for zeros and the skews in the distribution. The red points denote the core sample and the gray points denote the remainder of the sample. In both samples, although there is a discernible positive association between onward infections and degree, the heterogeneity in degree does not fully explain that of the secondary case distribution. Fig. S4(C) replicates Fig. 2C for the whole sample, showing the marginal distribution of PCIs. It is right-skewed as in the core sample. Fig. S4(D) shows the association between the log odds (logit) of the PCI and the logarithm of the degree. The red points denote the core sample and the gray points denote the remainder of the sample. Both samples are bimodal with a discernible negative association for those that infect others.

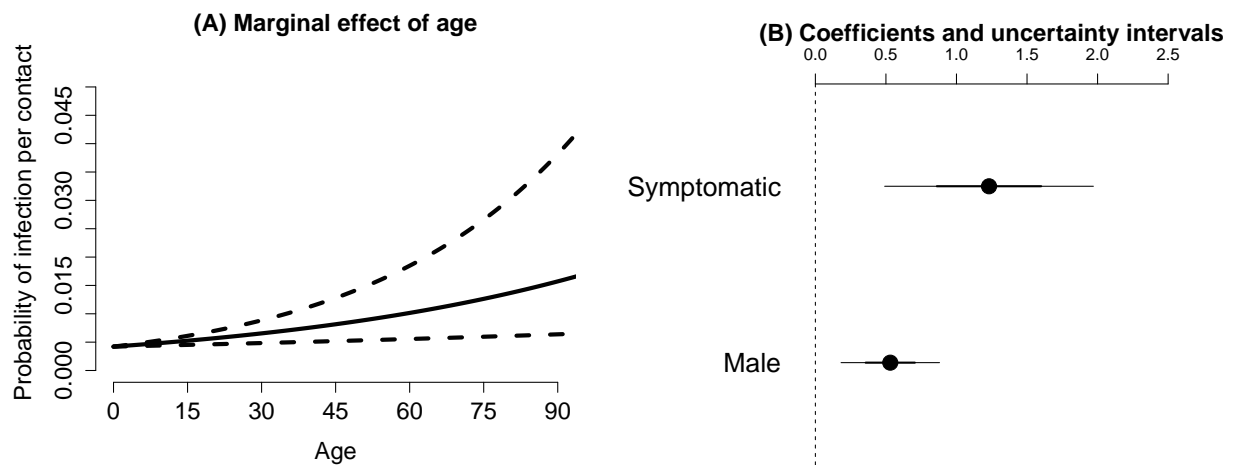


Figure S5: Regression results for likelihood of a secondary infection, against the age, sex and symptom status of the index case. Results are from a multivariate logistic regression. Panel (a) shows the marginal effect of age on the probability of secondary infection per contact for the reference groups (asymptomatic, female). Dashed lines represent the bounds of the 95% confidence interval on the coefficients. Panel (b) shows the logistic regression coefficients for the dichotomous variables associated with being symptomatic and being male. For a person at the average age (about 42 years old) being male raises the probability of secondary infection per contact by about 0.5 percentage points (from 0.8% to 1.3%). For a person at the average age being symptomatic raises the probability of secondary infection per contact by about 1.9 percentage points (from 0.8% to about 2.7%). Recall that these estimates come from data where the majority of individuals have no secondary infections, which reduces these overall probabilities.