

Text S1: The contribution of social behaviour to the transmission of influenza A in a human population

1 Model structure

2 To construct a model with A age groups and C contact classes, we first sorted participants by
3 age and divided them into A groups, each containing an equal number of people; the final class
4 contained fewer individuals if there was a remainder after division. We defined P_a to be the
5 age distribution of these groups (for Hong Kong, these values were taken from the 2011 census
6 [1]), where $\sum_{a=1}^A P_a = 1$. Using data on reported contacts we found $m_{a,b}$, the mean number of
7 contacts in age group a reported by individuals in group b . We then divided each age group into
8 a further C classes, based on reported contacts. The classes for each age group were defined by
9 sorting the individual reported number of contacts into ascending order, then dividing the list
10 into C equal parts (again, the final class was smaller if there was a remainder). If N_{ai} was the
11 number of participants who were age a and in contact class i , the age distribution of that age
12 and contact group was given by

$$P_{ai} = \frac{P_a N_{ai}}{\sum_{k=1}^C N_{ak}} . \quad (\text{S1})$$

13 In social contact surveys, participants often report only the age group of whom they interacted
14 with, and not their contact class [2, 3, 4]. We therefore had to estimate $m_{ai,bj}$, the mean number
15 of contacts with individuals in age group a and contact class i by participants in age group b
16 and class j . While we did not have data on the social contacts of reported contacts, we could
17 account for the age distribution of contacts in different contact classes when constructing our
18 transmission matrices. We knew how many of the contacts group (b, j) reported were in age
19 group a , so only needed to estimate how these were split between the different contact classes
20 in group a . We did this by assuming that the distribution of contacts made by (b, j) with group
21 (a, i) followed a weighted average [5], based on the total contacts reported by all the contact
22 classes in age group a :

$$m_{ai,bj} = m_{a,bj} \frac{m_{b,ai} P_{ai}}{\sum_{k=1}^C m_{b,ak} P_{ak}} . \quad (\text{S2})$$

23 We used an SIR model for simulations, with individuals falling into one of three compartments:
 24 susceptible, infective or recovered (and hence immune). The transmission rate to group ai from
 25 group bj was given by $\beta_{ai,bj} = qm_{ai,bj}/P_{ai}$, where q was a scaling factor dependent on the basic
 26 reproduction number [4]. The final epidemic size in each age group a and contact class i , ϕ_{ai} ,
 27 could therefore be found by solving the following coupled equation [6],

$$\phi_{ai} = 1 - \exp \left(- \sum_{b=1}^A \sum_{j=1}^C \beta_{ai,bj} P_{bj} \phi_{bj} \right). \quad (\text{S3})$$

28 Estimating reported contacts between age groups

When constructing the model, we divided the participants into A groups of equal size (Fig. S3). However, as A varied, the age groups used did not always line up with the three boundaries in the survey (age under 20, 20 to 65, over 65). We therefore had to estimate $m_{a,b}$, the mean number of contacts in group a reported by individuals in group b , using contact data from the Hong Kong survey and population age distribution data from the 2011 census [1]. First we used census data [1] to calculate $K_{i,x}$, the proportion of each age group i that fall within age boundary B_x , where

$$B_x = \begin{cases} [0, 20] & \text{if } x = 1; \\ (20, 65] & \text{if } x = 2; \\ (65, 105] & \text{if } x = 3. \end{cases}$$

Next, we used $K_{i,x}$ to estimate the number of contacts between age groups by assuming that individuals make contacts at random, based on their reported contacts and total available contacts in the population. Suppose there are N people in a population and n total contacts, with degree distribution $\{Q_k\}_{k=1}^n$. If an individual makes a contact at random, the probability of making a contact with a person who has j total contacts is:

$$\mathbb{P}(\text{meet person with } j \text{ contacts}) = \frac{jQ_j N}{n} = \frac{jQ_j N}{N \sum_{k=1}^{\infty} kQ_k} = \frac{jQ_j}{\sum_{k=1}^{\infty} kQ_k}.$$

29 We extended this approach to find the expect number of contacts reported by individuals in age
 30 group b that are in age group a . If P_a denotes the proportion of individuals in age group a in the
 31 2011 census [1], where $\sum_{a=1}^A P_a = 1$, we have

$$m_{a,b} = \sum_{x=1}^3 \sum_{y=1}^3 K_{b,y} m_{x,b} \frac{m_{y,a} K_{a,x} P_a}{\sum_{k=1}^A m_{y,k} K_{k,x} P_k} . \quad (\text{S4})$$

32 The expression for $m_{a,b}$, or $m_{a,bj}$ if age group b has been divided into multiple contact classes,
 33 can then be used with Equation S2 to derive the transmission rates that appear in Equation 1.

34 **Example: Estimating reported contacts between age groups**

35 Here, we show how to estimate $m_{a,b}$, the mean number of contacts in group a reported by
 36 individuals in group b , using contact data from the 2009/10 Hong Kong serological survey
 37 and population age distribution data from the 2011 census[1]. Suppose we assume twelve age
 38 groups in the model (i.e. $A = 12$). The second age group, which for now we denote b , contains
 39 individuals with ages ranging from 15.5 to 21.2, and the penultimate age group, a , contains
 40 people of ages 60.7 to 67.4 (Figure S3). According to the population age distribution:

$$K_b = \{0.783, 0.217, 0\} \text{ and } K_a = \{0, 0.735, 0.265\} . \quad (\text{S5})$$

41 Based on the Hong Kong survey data, individuals in group b had an average of 6.66 contacts
 42 with 20–65 year olds and 0.10 contacts with over 65s. In addition, individuals in group a had
 43 an average of 0.41 contacts with under 20s and 12.77 contacts with 20–65 year olds. Suppose
 44 we want to estimate how many contacts reported by under 20s in group b are with under 65s in
 45 group a . Using Equation S4, we have

$$m_{2,b} \frac{m_{1,a} K_{a,2} P_a}{\sum_{k=1}^A m_{1,k} K_{k,2} P_k} = 0.05 . \quad (\text{S6})$$

46 We can find the total contacts from age group b to age group a by summing up all possible
 47 combinations, weighting each by the relevant term in K_b , as in Equation S4. This gives $m_{a,b} =$
 48 0.15 .

49 **Comparison with network model**

50 When there is only one age group, and R_0 is small (but greater than one), the model formulation
51 is similar to that of a network approach. The attack rate in age age group a and contact class i
52 is given by

$$\phi_{ai} = 1 - \exp \left(- \sum_{b=1}^A \sum_{j=1}^C \beta_{ai,bj} P_{bj} \phi_{bj} \right). \quad (\text{S7})$$

When there are several contact classes, but only one age group, and q is small this can be expressed as

$$\phi_i = 1 - \exp \left(- \frac{q \sum_{j=1}^C m_j P_j \phi_j m_i}{\sum_{j=1}^C m_j P_j} \right) \quad (\text{S8})$$

$$\approx 1 - \left(1 - q \left[1 - \frac{\sum_{j=1}^C m_j P_j (1 - \phi_j)}{\sum_{j=1}^C m_j P_j} \right] \right)^{m_i} \quad (\text{S9})$$

53 where m_i is the mean number of contacts reported by individuals in class i . This equation has
54 the same form as the percolation approximation for final epidemic size in a network model[5],
55 though it does not account for reduction in available susceptible contacts as a result of network
56 structure.

57 **Simulation study to test robustness of model identification (Figure S1)**

58 We used a simulation study to test whether our model could correctly identify the ‘true’ model
59 among a range of candidate models. First we simulated attack rates for a specific number of age
60 and contact classes and contact type, finding ϕ_{ai} for each age-contact group using Equation 1.
61 For each participant in a given age-contact group ϕ_{ai} , we then simulated infection data from
62 a Bernoulli distribution with probability $1 - \phi_{ai}$. Given this simulated data, we repeated the
63 analysis in Figure 3, calculating the AIC for each candidate model in the framework, and using
64 ΔAIC to identify the one with most support of those tested. Results are shown in Figure S1.

65 **Sensitivity of results to inclusion of small background risk (Figure S2)**

66 Some of the models with multiple contact classes in Figure 3B had classes consisting solely
67 of individuals – some of whom had been infected – that had no reported close contacts. The

68 likelihood of such people seeing infection given the model assumptions was therefore zero. To
 69 assess whether our results were sensitive to these assumptions, we considered a framework with
 70 an additional small background rate of random contact among all members of the population.
 71 When this background risk is included, Equation 1 becomes

$$\phi_{ai} = 1 - \exp \left(- \sum_{b=1}^A \sum_{j=1}^C \beta_{ai,bj} P_{bj} \phi_{bj} \right) + h \quad (\text{S10})$$

72 where h is a parameter to be found. Results are shown in Figure S2.

73 **Relatively susceptibility in older groups**

In Equation 1, the final size of an epidemic in each group was calculated as

$$\phi_{ai} = 1 - S_{ai}(\infty)/S_{ai}(0) ,$$

where $S_{ai}(0) = N_{ai}$ denotes the number of individuals in group ai that are susceptible to infection at the start the epidemic and $S_{ai}(\infty)$ denotes that number that remain susceptible – and hence uninfected – at the end. To include relatively susceptibility in age groups older than δ , we solve Equation 1 as before, but with $S_{ai}(0) = \alpha N_{ai} \leq N_{ai}$. Specifically, we define the level of susceptibility at the start of the epidemic as

$$S_{ai}(0) = \begin{cases} N_{ai} & \text{if youngest person in age group } a \text{ is under age } \delta; \\ \alpha N_{ai} & \text{else .} \end{cases}$$

74 **Bootstrap resampling of data (Figure S4)**

75 Our model made the assumption that social contacts in our sample are representative of the
 76 population. To test the sensitivity of results in Figure 5B, we therefore repeated our analysis
 77 with different datasets. As we did not have additional empirical data, we instead used bootstrap
 78 samples: we obtained a new dataset of 762 individuals by resampling the 762 individuals in the
 79 Hong Kong dataset with replacement. The probability a specific participant is included in the
 80 resampled data set is $1 - (1 - 1/n)^n$. For large data sets – like the Hong Kong study – this can

81 be approximated by $1 - 1/e \approx 0.63$. We then ran the analysis for Figure 5B for ten different
82 resampled datasets. For each set of data, we found the AIC for each number of age groups,
83 then calculated the values of ΔAIC and plotted the relationship between model resolution and
84 performance (Figure S4).

85 **References**

- 86 [1] Demographic Statistics Section (2011) Hong kong population census. Census and Statistics
87 Department .
- 88 [2] Mossong J, Hens N, Jit M, Beutels P, Auranen K, et al. (2008) Social contacts and mixing
89 patterns relevant to the spread of infectious diseases. PLoS Med 5: e74.
- 90 [3] Riley S, Kwok KO, Wu KM, Ning DY, Cowling BJ, et al. (2011) Epidemiological char-
91 acteristics of 2009 (H1N1) pandemic influenza based on paired sera from a longitudinal
92 community cohort study. PLoS Med 8: e1000442.
- 93 [4] Wallinga J, Teunis P, Kretzschmar M (2006) Using data on social contacts to estimate age-
94 specific transmission parameters for respiratory-spread infectious agents. American Journal
95 of Epidemiology 164: 936.
- 96 [5] Meyers LA, Pourbohloul B, Newman MEJ, Skowronski DM, Brunham RC (2005) Network
97 theory and sars: predicting outbreak diversity. J Theor Biol 232: 71-81.
- 98 [6] Riley S, Wu JT, Leung GM (2007) Optimizing the dose of pre-pandemic influenza vaccines
99 to reduce the infection attack rate. PLoS Med 4: e218.