# INTERFACE

## Research

**Author for correspondence:**
Stephen M. Kissler
e-mail: sk792@cam.ac.uk

**THE ROYAL SOCIETY**
PUBLISHING

# Symbolic transfer entropy reveals the age structure of pandemic influenza transmission from high-volume influenza-like illness data

Stephen M. Kissler[1,2], Cécile Viboud[3], Bryan T. Grenfell[4] and Julia R. Gog[1]

[1]Department of Applied Mathematics and Theoretical Physics, University of Cambridge, Wilberforce Road, Cambridge, UK
[2]Department of Immunology and Infectious Diseases, Harvard T.H. Chan School of Public Health, Boston, MA, USA
[3]Fogarty International Center, National Institutes of Health, Bethesda, MA, USA
[4]Department of Ecology and Evolutionary Biology, University of Princeton, Princeton, NJ, USA

SMK, 0000-0001-6000-8387

Existing methods to infer the relative roles of age groups in epidemic transmission can normally only accommodate a few age classes, and/or require data that are highly specific for the disease being studied. Here, symbolic transfer entropy (STE), a measure developed to identify asymmetric transfer of information between stochastic processes, is presented as a way to reveal asymmetric transmission patterns between age groups in an epidemic. STE provides a ranking of which age groups may dominate transmission, rather than a reconstruction of the explicit between-age-group transmission matrix. Using simulations, we establish that STE can identify which age groups dominate transmission even when there are differences in reporting rates between age groups and even if the data are noisy. Then, the pairwise STE is calculated between time series of influenza-like illness for 12 age groups in 884 US cities during the autumn of 2009. Elevated STE from 5 to 19 year-olds indicates that school-aged children were likely the most important transmitters of infection during the autumn wave of the 2009 pandemic in the USA. The results may be partially confounded by higher rates of physician-seeking behaviour in children compared to adults, but it is unlikely that differences in reporting rates can explain the observed differences in STE.

## 1. Introduction

Age is a key predictor of a person's rate of both acquiring [1–6] and transmitting [1,7,8] influenza. Children tend to contribute more to influenza transmission than adults do [1,4,7], but the precise epidemiological roles of different age groups can shift from season to season [9] and may change markedly in pandemic years [10]. From a public health perspective, untangling the relative roles of different age groups could help guide targeted vaccination strategies [7,11–13] and other age-related interventions, like the selective closure of schools [14–16]. However, data with sufficient resolution to identify detailed epidemiological relationships between age groups have so far been scarce, and even when such data exist, current methods are insufficient for reliably uncovering those relationships.

Electronic medical records (EMRs) help address the issue of data scarcity by providing high-volume influenza-like illness (ILI) incidence data with detailed age structure [17]. EMRs are routinely produced by physicians for insurance purposes during the majority of outpatient visits in the USA [17]. Since EMRs generally contain syndromic illness classifications, EMR-based estimates of influenza incidence are subject to noise from non-ILI respiratory infection. EMR-based disease incidence estimates are also subject to geographical and demographic variation in physician-seeking behaviour. Laboratory-confirmed influenza cases,

as collected routinely by the Centers for Disease Control and Prevention (CDC) [18], provide more specific estimates of influenza incidence, but at substantially lower volume. Influenza incidence estimates from online search platforms and social media websites like Google [19] and Twitter [20] can provide massive amounts of data, but the reliability of these sources has been called into question, and they lack detailed age information [21]. Dedicated online platforms such as FluNearYou in the USA and FluSurvey in the UK, which gather reports of ILI symptoms from community volunteers [22,23], hold some promise for supplementing traditional ILI data streams [24–26], but represent a relatively small convenience sample of the population. So, while other data sources exist, EMRs offer a relatively promising and so-far underused source of fine-scale data on influenza incidence in the USA [17,21].

Previous attempts to infer the relative importance of different age groups for the transmission of influenza have sought to either reconstruct the explicit next-generation matrix (NGM) [3,27,28] or to infer the relative risk of infection between age groups [4]. The NGM-based methods have only been applied to scenarios with at most two age groups (children and adults), in part because they require strong assumptions about the structure of the NGM which become increasingly unrealistic as the number of age classes grows. The relative risk method [4,29,30] has been used to rank the importance of five age groups for the transmission of influenza, but requires data with high specificity for influenza, effectively precluding ILI data streams and the use of EMRs in particular. These methods are generally aimed at identifying optimal vaccine allocation strategies using data from the initial phase of an epidemic.

Symbolic transfer entropy (STE) [31] offers a way to infer the relative transmission importance of possibly many age groups from ILI data. STE is an extension of transfer entropy (TE) [32], which measures the amount of information the past states of one stochastic process provide about the transition probabilities of another. Intuitively, the TE is a measure of the amount of information 'transferred' from one stochastic process to another. To compute the STE, a time series is symbolized using a scheme that encodes its qualitative structure in a low-dimensional space, and then the TE is calculated from the relative frequencies of these symbols. The symbolization scheme makes the STE robust to moderate observational noise and to systematic shifts in amplitude, which in the context of EMR ILI data might arise from the presence of non-influenza ILI cases and from differences in reporting rate between age groups. These benefits come with the trade-off of requiring relatively large amounts of data compared to existing methods for inferring the age structure of disease transmission and providing only a relative ranking of transmission importance. STE has been used to study epileptogenic neural signals and the dissemination of information through social networks [31,33], but to our knowledge has not been systematically evaluated as a means of providing insight into infectious disease transmission. TE and STE are similar to other model-free methods that measure how information is shared and transferred between possibly coupled dynamic processes, including mutual information [32], Granger causality [34], and convergent cross mapping (CCM) [35]. Permutation entropy, a related measure, has recently been used to quantify the predictability of infectious disease outbreaks [36]. Compared to these alternatives, STE offers the advantage of applying in stochastic settings with nonlinear dynamics and moderate observational noise to reveal asymmetric flows of information.

Here, we use influenza-like outbreak simulations to demonstrate that STE reliably identifies asymmetries in transmission strength between age groups. Then, we use an EMR-based dataset capturing ILI incidence from 884 ZIP (postal) codes and 12 age classes across the USA to rank the relative importance of the various age groups in the transmission of the autumn wave of the 2009 A/H1N1pdm influenza pandemic in that country. We conclude that school-aged children (5–19 year-olds) were disproportionately responsible for transmitting influenza to infants through working-age adults in the autumn of 2009, in broad agreement with other findings. Our work demonstrates that STE could serve as an important tool for the detailed epidemiological analysis of age structure, especially as EMR data become more prevalent.

## 2. Material and methods

### 2.1. Data

The data consist of a convenience sample of CMS-1500 electronic medical claims forms submitted by primary care physicians across the US and maintained by SDI health (now IQVIA). Each claim is associated with a single outpatient visit, and includes one or more ICD-9 codes [37] listed by the physician that describes the patient's illness. The overall sample captures over 50% of all outpatient visits in the USA in 2009 [17]. The records are binned weekly and aggregated geographically by the first three digits of the ZIP (postal) code of the practice from which they are submitted [38]. Time series of weekly ILI incidence are created by extracting claims with a direct mention of influenza, or fever combined with a respiratory symptom, or febrile viral illness (ICD-9 487-488 OR [780.6 and (462 or 786.2)] OR 079.99), following Viboud et al. (2014) [17]. For each ZIP, the number of ILI cases in each week is divided by the total number of patients who visited a physician in that ZIP during that week, yielding an 'ILI ratio' time series. There are 884 ILI ratio time series, one for each ZIP in the lower 48 US states, each spanning 52 weeks from the week commencing 4 January 2009 through the week commencing 27 December 2009. The correspondence between the SDI-ILI dataset and reference influenza surveillance data from the US CDC is described in depth by Viboud et al. [17].

### 2.2. Symbolic transfer entropy

If $I$ and $J$ are discrete-state and discrete-time random processes such that $i_t$ and $j_t$ are the states of processes $I$ and $J$ at time $t$, then the TE from process $J$ to process $I$ is defined as

$$T_{J \to I} = \sum_{\Omega_I, \Omega_J} p(i_{t+1}, i_t^{(k)}, j_t^{(l)}) \log \left( \frac{p(i_{t+1}|i_t^{(k)}, j_t^{(l)})}{p(i_{t+1}|i_t^{(k)})} \right), \qquad (2.1)$$

where $i_t^{(k)}$ is shorthand notation for the $k$-step history of process $i$, $(i_t, \ldots, i_{t-k+1})$, and similarly $j_t^{(l)} = (j_t, \ldots, j_{t-l+1})$, such that $p(i_{t+1}, i_t^{(k)}, j_t^{(l)})$ is the joint probability of observing $i_{t+1}$, $i_t^{(k)}$, and $j_t^{(l)}$; $p(i_{t+1}|i_t^{(k)}, j_t^{(l)})$ is the probability of observing $i_{t+1}$ conditioned on $i_t^{(k)}$ and $j_t^{(l)}$; and $p(i_{t+1}|i_t^{(k)})$ is the probability of observing $i_{t+1}$ conditioned only on $i_t^{(k)}$. The logarithm has base 2, so that the TE is measured in bits. The sum is over all possible combinations of states $(i_{t+1}, i_t^{(k)}, j_t^{(l)})$, where $i_{t+1}, i_t^{(k)} \in \Omega_I$ and $j_t^{(l)} \in \Omega_J$, and $\Omega_I$ and $\Omega_J$ are the state spaces for processes $I$ and $J$. Equation (2.1) is a Kullback–Leibler divergence that measures how much process $I$ deviates from the generalized Markov property $p(i_{t+1}|i_t, \ldots, i_1) = p(i_{t+1}|i_t^{(k)})$, given the last $l$ states of process $J$. In practice, the histories are often fixed at length 1 ($k = l = 1$) and the probabilities are estimated from simple counts of the observed data [32].
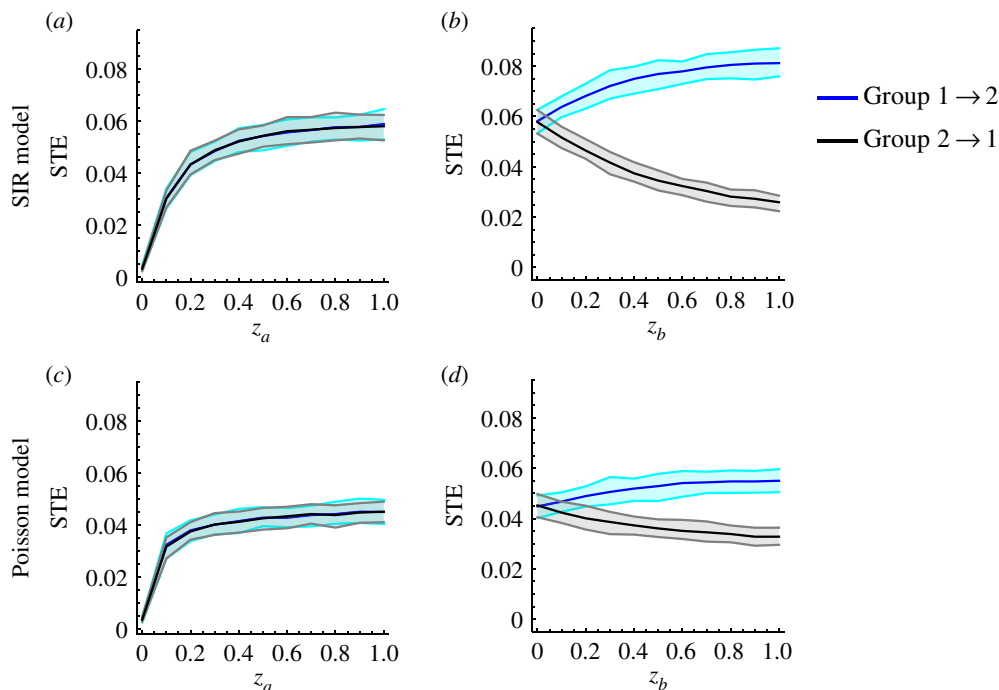
**Figure 1.** Mean (95% CI) Group $1 \rightarrow 2$ (blue) and Group $2 \rightarrow 1$ (black) STE values as the coupling between the two groups ranges from none to fully symmetric (a,c), and from fully symmetric to strongly driven by Group 1 (b,d). The curves are produced by simulating 100 ensembles of 800 epidemics each from the stochastic SIR model (a,b) or the Poisson model (c,d) for each value of $z_a$ and $z_b$ between 0 and 1 in steps of 0.1, and then calculating the between-group STE for each ensemble. The relative reproduction matrices that capture these two coupling scenarios are given in equations (3.1) and (3.2).

The TE is limited in that it is only defined for stochastic processes with a discrete state space. Staniek & Lehnertz [31] introduce STE as a way to calculate information transfer between time series processes that have continuous- or near-continuous state spaces. Motivated by the insight that the relative amplitudes of subsequent observations from these sorts of processes may provide enough information to reveal interactions between them, they propose symbolizing the time series based on ordered $m$-tuples of observations (electronic supplementary material, figure S1). This reduces the (near-)continuous state space of the original stochastic process to a discrete set of $m!$ symbols. In practice, $m$ is often chosen to be 2 or 3, giving a state space of two or six symbols, respectively. For $m = 3$, we also tested the effect of collapsing the two concave-up and the two concave-down symbols into a single symbol each, resulting in a smaller state space (four versus six symbols) while capturing a similar level of qualitative detail. Details on the symbolization of time series and the empirical calculation of the STE are provided in the electronic supplementary material.

## 2.3. SIR epidemic simulation model

For simulations with just two age classes, we use a stochastic SIR model implemented using the Gillespie algorithm [39]. For all simulations, the basic reproduction number $R_0$ is set at 1.5, consistent with estimates of the basic reproduction number of 2009 A/H1N1 pandemic influenza [40,41]. We consider a population size of $N = 1000$ split evenly between classes 1 and 2, so that $N_1 = N_2 = 500$ (age groups with different population sizes are also considered in the electronic supplementary material). The expected time to recovery $1/\gamma$ is assumed constant for all age groups and is set at 7 days, which is consistent with estimates of the infectious period for 2009 pandemic influenza [41]. We consider a range of between-group transmission strengths. Electronic supplementary material, table S1 gives the rates at which individuals of each class stochastically progress from susceptible to infected to recovered. Infections are binned into week-long intervals, and Poisson noise is added to simulate non-influenza ILI. Electronic

supplementary material, figure S7 depicts five incidence time series produced using the model. Full details on the model and simulation procedure are given in the electronic supplementary material.

## 2.4. Poisson epidemic simulation model

For more than two age classes, the full stochastic SIR model becomes too computationally demanding for repeated simulations to be practical. So, we also define an outbreak simulation model based on a self-exciting Poisson process, similar to [42]. We choose the time units $t$ to match the mean generation interval of the infection, which we set at 3.5 days [43]. To generate epidemics, we use a stepwise-constant effective reproduction number $R_t$, such that $R_t = 1.5$ for the first four weeks (eight generations) of the outbreak and $R_t = 0.8$ thereafter. Infections are binned from the half-week generations into week-long intervals, and additional Poisson noise is added to each bin to simulate non-influenza ILI. For simulations with two age classes, the Poisson model yields epidemics of similar length and magnitude as the two-age-class SIR model (compare electronic supplementary material, figures S7 and S8), and yields comparable STE inferences (figure 1), which suggests that the Poisson model is an acceptable approximation to the stochastic SIR model. Full details on the implementation of the Poisson model are given in the electronic supplementary material.

## 2.5. Reporting rates

Only a fraction of influenza cases are represented in the SDI-ILI dataset, since many people do not seek medical care for their symptoms. The tendency to seek medical care given infection with an ILI can vary by age group [44]. To factor this into the outbreak simulations, we introduce a reporting rate vector $c$ in which element $c_i$ gives the expected proportion of individuals in age class $i$ who seek medical care when infected with an ILI. It is then possible to simulate a 'reported' disease incidence time series

$$Y_{i,t}^{\text{obs}} \sim \text{binomial}\,(Y_{i,t}, c_i), \tag{2.2}$$

where $Y_{i,t}$ is the simulated number of infected individuals in age class $i$ at time $t$ (under either model) and $Y_{i,t}^{\text{obs}}$ is the simulated reported number of infections in age class $i$ at time $t$.

# 3. Results

## 3.1. Symbolic transfer entropy reveals transmission asymmetries between two coupled age groups

We first calculate the STE between two age groups as the within- and between-group reproduction ratios vary. We consider between-group transmission that ranges from (a) fully decoupled to fully symmetric, and (b) fully symmetric to strongly driven by Group 1. The between-group infectiousness is specified using a 'relative reproduction matrix' $r$, which is a scaled version of the NGM [27], such that $\text{NGM} = (R_0/\rho)r$, where $\rho$ is the maximum eigenvalue of $r$ and $R_0$ is the basic reproduction number. The elements of $r$ define the relative infectiousness of various population groups, such that $r_{i,j}/r_{k,j}$ gives the proportional difference in group $j$'s infectiousness for group $i$ versus group $j$. For example, if $r_{i,j}/r_{k,j} = 2$, then a member of group $j$ is expected to infect twice as many members of group $i$ than of group $k$. The population sizes of the various groups are assumed to be equal, though deviations from this assumption are considered in the electronic supplementary material. Scenario (a) is encapsulated by the relative reproduction matrix

$$r_a = \begin{bmatrix} 1 & z_a \\ z_a & 1 \end{bmatrix}, \qquad (3.1)$$

where $z_a \in [0, 1]$. Scenario (b) is encapsulated by the relative reproduction matrix

$$r_b = \begin{bmatrix} 1 + 3z_b & 1 \\ 1 + z_b & 1 \end{bmatrix}, \qquad (3.2)$$

where $z_b \in [0, 1]$.

Figure 1 depicts the change in STE under these two transmission scenarios, calculated from epidemics simulated using the stochastic SIR model (figure 1$a$–$b$) and the Poisson model (figure 1$c$–$d$). Each pane in figure 1 is produced using 100 ensembles of 800 simulated epidemics for each value of $z_a$ and $z_b$ between 0 and 1 in steps of size 0.1. For each ensemble, the 800 simulated incidence time series are symbolized using symbols of length $m = 3$, and then the between-group transfer entropies are estimated using the relative symbol frequencies (see electronic supplementary material, figure S3), producing 100 STE estimates for each value of $z_a$ and $z_b$. The solid blue (black) lines in figure 1 depict the mean Group $1 \rightarrow 2$ (Group $2 \rightarrow 1$) STE for each value of $z_a$ and $z_b$ across the 100 ensembles. The shaded blue (black) bands depict the range of the middle 95 Group $1 \rightarrow 2$ (Group $2 \rightarrow 1$) STE estimates for each value of $z_a$ and $z_b$ across the 100 ensembles, analogous to a 95% confidence interval. Under both the stochastic SIR and the Poisson models, the between-group STE increases steadily as the transmission coupling ranges from none to symmetric (figure 1$a,c$). Once Group 1 begins to dominate transmission, the Group $1 \rightarrow 2$ STE increases and the Group $2 \rightarrow 1$ STE decreases (figure 1$b,d$), accurately capturing the transmission relationship between the age groups.

When Group 1 drives transmission, the Poisson model yields a smaller difference in the STE between the two age groups than the stochastic SIR model does (figure 1$b,d$).

Visual inspection suggests that the simulated time series produced using the stochastic SIR model tend to feature more stochastic fluctuations than the time series produced using the Poisson model (electronic supplementary material, figures S7 and S8). Since STE is effectively a measure of how these stochastic fluctuations transmit from one age group to another, this may explain why the differences in STE calculated using the Poisson model are relatively less pronounced. Overall, the qualitative similarity between the STE estimates from the two transmission models suggests that the Poisson model is an acceptable approximation to the stochastic SIR model, and that simulations from the Poisson model tend to produce more conservative estimates of the difference in STE between age groups than the stochastic SIR model.

## 3.2. Symbolic transfer entropy reveals transmission asymmetries despite incomplete reporting

Next, we evaluate how incomplete reporting influences the detection of asymmetries in transmission strength. Figure 2 depicts the mean estimated STE across 100 ensembles of 800 epidemics each for reporting rates $c_i = 0.1$, $c_i = 0.5$ and $c_i = 1$ with equal reporting rates across all age groups (see also electronic supplementary material, figure S9). The epidemic simulations are produced using the Poisson model with relative reproduction matrix

$$r = \begin{bmatrix} 1 & 2 & 1 & 1 \\ 1 & 4 & 1 & 1 \\ 1 & 2 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}, \qquad (3.3)$$

which could represent 'children' (Group 2) having strong within-group transmission ($r_{2,2} = 4$) and intermediate transmission to 'infants' (Group 1) and 'adults' (Group 3) ($r_{1,2} = r_{3,2} = 2$). Even for reporting rates as low as 0.1, the STE values from Group 2 are higher than those from any other group. As the reporting rates increase, the differences become more pronounced, accurately capturing the transmission dominance of Group 2 over the other groups. According to Biggerstaff et al. [44], true reporting rates for ILI in the US during the 2009 pandemic were between 0.4 and 0.6 (scenario (B) in figure 2), for which the transmission dominance of Group 2 is clear.

## 3.3. Symbolic transfer entropy reveals transmission asymmetries between 12 coupled age groups

To test the ability of STE to identify transmission asymmetries from data on the scale of the SDI-ILI dataset, we use the Poisson model to simulate 100 ensembles of 800 epidemics each with 12 age groups. We consider the scenarios (a) with the $12 \times 12$ relative reproduction matrix electronic supplementary material, Eq. S49, representing high transmission from Groups 3–5 to Groups 3–5 ($r_{i,j} = 4$ for $i$, $j \in \{3, 4, 5\}$), intermediate transmission from groups 3–5 to groups 1–2 and 6–9 ($r_{i,j} = 2$ for $i \in \{1, 2, 6, 7, 8, 9\}$ and $j \in \{3, 4, 5\}$), baseline transmission ($r_{i,j} = 1$) between all other groups, and uniform 50% reporting rate across all groups, and (b) with uniform transmission strength across all age groups (i.e. a $12 \times 12$ relative reproduction matrix with '1' for all entries), 60% reporting rate for groups 1–5, and 40% reporting rate for groups 6–12, following the estimates of Biggerstaff et al. [44] for the ILI reporting rates in the USA during the 2009 influenza pandemic for children and adults, respectively.
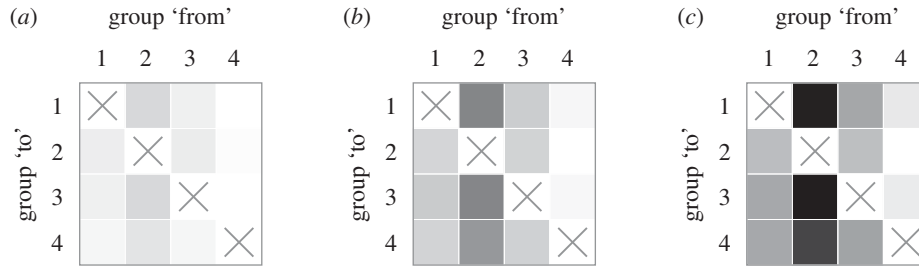
**Figure 2.** Mean pairwise STE values for epidemics strongly driven by Group 2 with varying reporting rates: (a) $c = 0.1$, (b) $c = 0.5$ and (c) $c = 1$ (see equation (2.2)). A box in row $i$ and column $j$ corresponds to the STE from group $j$ to group $i$, where darker shades corresponds to higher STE. The reporting rate in each subfigure is constant across groups. The STE values are produced by simulating 100 ensembles of 800 epidemics each from the Poisson model for each value of $c$ and then calculating the between-group STE for each ensemble. The relative reproduction matrix that specifies within- and between-group transmission rates is given by equation (3.3).
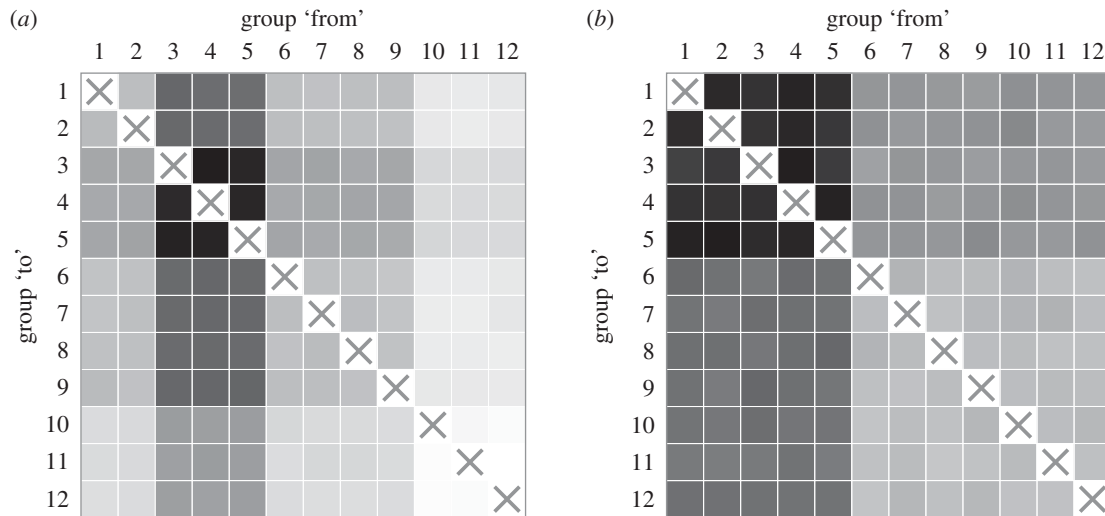


**Figure 3.** Mean pairwise STE values between 12 groups for epidemics strongly driven by Groups 3, 4 and 5 and uniform 50% reporting rate across all age groups (a), and for epidemics driven equally by all age groups, 60% reporting for Groups 1–5, and 40% reporting for Groups 6–12 (b). A box in row $i$ and column $j$ corresponds to the STE from group $j$ to group $i$, where darker shades corresponds to higher STE. To generate the STE values, 100 ensembles of 800 epidemics were simulated from the Poisson model using relative rate matrix electronic supplementary material, Eq. S49 for (a) or a relative rate matrix with all entries equal to 1 for (b). Each ensemble generates 144 pairwise STE values, so that each box represents the mean value across the 100 ensembles. The raw values are listed in electronic supplementary material, Eqs S50 and S51.

Figure 3 depicts the mean pairwise STE estimates between the 12 age groups under both scenarios. The square in row $i$ and column $j$ represents the STE from Group $j$ to Group $i$. Darker squares correspond to higher STE. For the asymmetric transmission/uniform reporting rate scenario (scenario (a), figure 3a), the STE clearly captures the transmission dominance of Groups 3, 4 and 5. The pairwise STE does not simply reproduce the structure of the relative reproduction matrix, as evidenced by the variability in mean pairwise STE for age groups other than Groups 3–5. This is because the STE captures a 'knock-on' effect for which information transferred from a strongly driving age group can propagate through other age groups. For the uniform transmission/variable reporting rate scenario (scenario (b), figure 3b), it is evident that elevated reporting rates can also lead to elevated STE, both to and from the groups with elevated reporting rate (Groups 1–5). Overall, the variability in STE due to differences in reporting rate appears to be smaller than the variability in STE due to differences in transmission strength. Further discussion on the effect of reporting rates on STE may be found in the electronic supplementary material.

## 3.4. School-aged children contributed disproportionately to transmission during the autumn 2009 A/H1N1pdm influenza outbreak in the US

To estimate the pairwise STE between the 12 age groups represented in the SDI-ILI dataset during the 2009 A/H1N1pdm influenza pandemic, we extract data from the 25 weeks between 12 July 2009 and 27 December 2009 and symbolise the ILI time series for each age group in each ZIP using a symbol length of $m = 3$. The pairwise STE values between all age groups are depicted in figure 4. The STE is highest in the columns representing 5–19 year-olds. This provides evidence that there was systematically elevated transmission from school-aged children to infants through adults. The adult-adult STE is also moderately elevated, suggesting that adults may have played a relatively important role in transmitting the outbreak among themselves, though this could also be explained by elevated transmission from children alone. Compare, for example, to the left-hand plot in figure 3: in that simulation, only transmission from children is elevated, but it
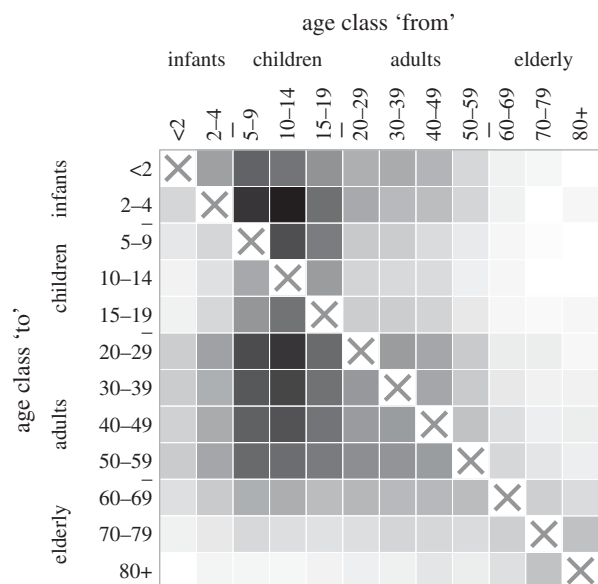
**Figure 4.** Mean pairwise STE values between the 12 groups represented in the SDI-ILI dataset during the autumn 2009 A/H1N1pdm pandemic influenza outbreak. A box in row *i* and column *j* corresponds to the STE from group *j* to group *i*, where darker shades corresponds to higher STE. The raw values are listed in electronic supplementary material, Eq. S52.

causes a moderate elevation in the STE from adults and infants to the other age groups due to the knock-on effect.

As a control, we also calculated the pairwise STE between all age groups during 25 post-pandemic weeks, from 10 January 2010 through 27 June 2010. For these months, there is no apparent age structure in transmission (see electronic supplementary material). We also calculated the pairwise STE between age groups for six previous influenza seasons (see electronic supplementary material). For the 2009 pandemic, there is a higher maximum pairwise STE and greater variation in the pairwise STEs than for any previous season. This could reflect differences in baseline ILI, which was likely lower during the autumn 2009 pandemic wave than during the seasonal outbreaks, due to the pandemic's earlier timing. A lower baseline ILI might have made pairwise differences in STE easier to detect in 2009. However, the relatively higher and more heterogeneous STE values in 2009 are also consistent with the hypothesis that school-aged children played a disproportionately large role in the spread of the 2009 pandemic, as has been described elsewhere [4].

It is unlikely that differences reporting rates alone can account for the elevated STE from 5 to 19 year-olds to the other age groups. The mean pairwise STE values computed from simulations with uniform transmission rates and unequal reporting rates in §3.3 range from 0.0057 to 0.0084 (see electronic supplementary material, Eq. S51), while the pairwise STE values computed from the SDI-ILI data range from 0.0056 to 0.084 (see electronic supplementary material, Eq. S52), an order of magnitude larger. The mean pairwise STE values computed from simulations with asymmetric transmission and uniform reporting rates §3.3 range from 0.0047 to 0.014 (see electronic supplementary material, Eq. S50), closer to the range observed from the SDI-ILI data but still somewhat smaller. This points towards a possible combined effect of strong transmission driving from children plus elevated reporting in children. In addition, re-calculating the pairwise STE using probabilistic reconstructions of the pre-reporting

SDI-ILI incidence time series (see electronic supplementary material) indicate that the observed transmission dominance of 5–19 year-olds persists even after adjusting for potential differences in reporting rate between children and adults. Furthermore, Biggerstaff *et al.* [44] report that 0–4 year-olds had the highest reporting rates for ILI in the USA in 2009, yet the STE from 0 to 4 year-olds is relatively low compared to the other age groups. If reporting rates alone could explain the observed differences in STE, the STE from infants should be at least as high as the STE from school-aged children.

It is also unlikely that the unequal partitions of the age groups can explain the observed patterns in the pairwise STE. The age groups under 20 years are partitioned such that they span fewer years, and thus contain fewer individuals, than the age groups above 20 years. Direct calculations and simulations (see electronic supplementary material) indicate that, all else being equal, the out-going STE for a given group tends to increase as the group's population size increases relative to the sizes of the other groups. If differences in the groups' population sizes were driving the observed pairwise STE values, we would expect the age groups over 20 years to appear to dominate transmission—which is the opposite of what we observe here.

## 4. Discussion

Here, we propose STE as a means of ranking which age groups contribute most to the transmission of infectious disease outbreaks. STE is chosen over other extensions of TE due to its robustness to point-wise noise and overall amplitude shifts in time series, which especially affect the ILI data stream due to non-influenza respiratory illness and incomplete reporting. Simulation studies indicate that STE can correctly rank transmission asymmetries between age groups. While such a ranking does not provide definitive guidance for targeted interventions, it can provide a useful starting point when other epidemiological information is lacking, as is often the case in emerging outbreaks. STE is positively associated with reporting rates, which can partially confound estimates of asymmetric transmission. STE estimates from ILI time-series data from July–December 2009 in the USA suggest that the transmission of the autumn wave of the A/H1N1pdm pandemic influenza outbreak was likely dominated by 5–19 year-olds. It is unlikely that this result can be explained by differences in reporting rates alone.

The identification of elevated transmission from school-aged children during the 2009 influenza pandemic agrees with most other studies on age-specific transmission of both seasonal and pandemic influenza [1,4,7,8]. Elevated transmission from school-aged children is likely due in part to the relatively high number of daily interpersonal contacts made by members of these age groups. Mossong *et al.* [1] for example estimate that 10–19 year-olds have more contacts per day than any other age group, and conclude from a modelling study based on empirical contact data that 5–19 year-olds are likely to both suffer the highest burden of disease and to drive the early-stage transmission of an outbreak transmitted by droplets through close contacts, like influenza. This underscores the importance of monitoring children during pandemic influenza outbreaks. We do not recommend using STE alone to set vaccination priorities. STE is shown here to reveal the population structure of transmission, but as a

correlation-based measure [45], it does not provide conclusive evidence that vaccinating high-transmission groups would optimally reduce overall transmission. If high-specificity (e.g. laboratory-confirmed) data are available for a subset of the population, we recommend using STE in tandem with existing risk-based methods [4,29,30] to identify optimal intervention strategies. For example, one might use STE to identify epidemiologically relevant partitions of the population (for the 2009 A/H1N1pdm influenza outbreak in the US, these might be <5, 5–19, 20–59, and 60+ years). Then, using these demographic partitions, one could calculate the relative risk of infection between these coarser groups from the high-specificity data, thereby leveraging the data to its fullest potential.

In this study, we have only considered a single demographic variable—age—as the basis for examining asymmetric transmission strengths. Age is known to be a key predictor of influenza transmission [1,7,8]. For other infectious diseases, different demographic variables will be necessary. In the context of sexually transmitted infections, for example, one might consider transmission asymmetries between individuals of different racial/ethnic [46] or sexual behaviour [47] groups. Our study supports the hypothesis that among all age groups, but not necessarily among all possible demographic groups, school-aged children were the strongest transmitters of influenza during the autumn 2009 A/H1N1pdm outbreak in the US.

TE is closely linked to mutual information [32] and Granger causality [34]. Unlike TE, mutual information is symmetric; that is, it measures the probabilistic dependence between two processes, but cannot determine the direction of information transfer between them, if there is any [32]. Measuring the delayed mutual information between two processes is one way to introduce asymmetry. This takes a step towards inferring whether one process influences another, by measuring shared information between the present state of one process and the past states of another [32]. While the lagged mutual information describes how one process history predicts the static probabilities of another, the TE measures how one process history predicts the transition probabilities of another. Because of this, the TE is less likely to be confounded by a shared input signal, and is a better measure of stochastic 'driving' [32]. Section 2 of Kaiser & Schreiber [48] provides a detailed description of the differences between TE and mutual information. Granger causality, on the other hand, is a special case of TE that arises when the stochastic processes are jointly Gaussian distributed [49]. The TE is thus better suited than Granger causality for making inferences on more general, possibly nonlinear, processes, though this comes at the expense of requiring more data and having no clear way to test statistical significance [49].

CCM [35] was developed to solve a similar problem as TE, but is based on somewhat different underlying theory. CCM was developed to detect so-called causal relationships in partially stochastic systems with underlying deterministic structure. CCM relies on Takens' theorem [50] to reconstruct candidate manifolds of the underlying dynamical system using lagged observations from two-time series. 'Causality' is inferred if nearby points on one reconstructed manifold consistently map to nearby points on the other reconstructed manifold. CCM has been used to provide evidence that temperature and absolute humidity fluctuations drive the timing of global seasonal influenza outbreaks [51], though some controversy surrounds these findings [52,53]. Nevertheless, it would be interesting to see whether CCM can reveal asymmetric

epidemiological interactions between age groups, and to compare its findings with those identified using TE. Lungarella *et al.* [54] provide more detail on the relationships between various methods that infer asymmetric relationships from time-series data. As an aside, we prefer to avoid the term 'causality' with respect to these methods, despite its frequent use in the literature. Determination of so-called counterfactual causality, as distinguished from Granger-type causality, requires intervention [55], which is normally not possible in retrospective epidemiological studies. Regardless of the vocabulary used, the above-listed techniques have successfully revealed fundamental structures in real-world coupled dynamic processes [31,33,35,56–58].

Despite the apparent well-suitedness of STE for making inferences from ILI data, its epidemiological relevance currently remains limited. The calculation of STE requires no prior epidemiological information whatsoever, which makes its success somewhat surprising. The NGM [27] is the key object for characterizing age-structured, or more generally population-structured, disease transmission dynamics, and yet there is no obvious direct link between STE estimates and the NGM. It is possible that further simulation studies could help identify such a link; even though the STE values seem to bear little mechanistic meaning apart from the relative ordering of age groups that they yield, it is possible that regressing the inferred STE values on an underlying known NGM could connect the pairwise STE matrix with the NGM under certain conditions. However, it appears unlikely that a simple link exists, especially since STE can say nothing about transmission within a single age group, which is necessary for filling in the diagonal entries of the NGM. STE and related methods such as CCM that do not explicitly incorporate mechanistic descriptions of the underlying physical system are unlikely to be able to reveal more than an approximate hierarchy of transmission strengths. Nevertheless, such a hierarchy can contain valuable information, especially if developing and fitting a mechanistic model is too demanding to be practicable. Certain extensions to STE could also enhance its relevance for epidemiological inference. Local TE [59] and state-dependent TE [60], like the contextual STE (see electronic supplementary material), are intended to make the TE more flexible and general, by considering how information transfer may change under varying conditions or 'meta-states'. Conditional TE [45,61] takes fuller account of the possible drivers of a given signal which could help to reveal polyadic and synergistic relationships between demographic groups. These extensions may yield better insight into epidemic processes, which are inherently nonlinear and context-dependent, than the more traditional measurements of TE can provide.

Perhaps the most important challenge confronting the TE and related measurements is deciding how to measure statistical power and significance. STE calculations rely on a middle level of stochasticity in the underlying stochastic processes; for a deterministic system, the STE will always be exactly zero, while for a stochastic system with too much within-sequence noise, the small-scale variation in amplitudes will likely mask important patterns from which the transfer of information might be inferred. The acceptable range of stochasticity has not been clearly defined. Similarly, it is unclear how best to measure when a difference in STE should be called statistically significant. Though this is recognized as an open and difficult problem [49,52], it may be possible to make some progress by

assuming that the underlying process follows certain epidemiological, or otherwise well-specified, dynamics.

# References

1. Mossong J *et al.* 2008 Social contacts and mixing patterns relevant to the spread of infectious diseases. *PLoS Med.* **5**, e74. (doi:10.1371/journal.pmed.0050074)

2. Fraser C *et al.* 2009 Pandemic potential of a strain of influenza A (H1N1): early findings. *Science (New York, N.Y.)* **324**, 1557–1561. (doi:10.1126/science.1176062)

3. Nishiura H, Castillo-Chavez C, Safan M, Chowell G. 2009 Transmission potential of the new influenze A(H1N1) virus and its age-specificity in Japan. *Eurosurveillance* **14**, 1–5. (doi:10.2807/ese.14.22.19227-en)

4. Worby CJ, Chaves SS, Wallinga J, Lipsitch M, Finelli L, Goldstein E. 2015 On the relative role of different age groups in influenza epidemics. *Epidemics* **13**, 10–16. (doi:10.1016/j.epidem.2015.04.003)

5. Yang Y, Sugimoto JD, Halloran ME, Basta NE, Chao DL, Matrajt L, Potter G, Kenah E, Longini IM. 2009 The transmissibility and control of pandemic influenza A(H1N1) virus. *Science* **326**, 729–733. (doi:10.1126/science.1177373)

6. Brownstein JS, Kleinman KP, Mandl KD. 2005 Identifying pediatric age groups for influenza vaccination using a real-time regional surveillance system. *Am. J. Epidemiol.* **162**, 686–693. (doi:10.1093/aje/kwi257)

7. Wallinga J, Teunis P, Kretzschmar M. 2006 Using data on social contacts to estimate age-specific transmission parameters for respiratory-spread infectious agents. *Am. J. Epidemiol.* **164**, 936–944. (doi:10.1093/aje/kwj317)

8. Smieszek T, Balmer M, Hattendorf J, Axhausen KW, Zinsstag J, Scholz RW. 2011 Reconstructing the 2003/2004 H3N2 influenza epidemic in Switzerland with a spatially explicit, individual-based model. *BMC Infect. Dis.* **11**, 115. (doi:10.1186/1471-2334-11-115)

9. Bedford T *et al.* 2015 Global circulation patterns of seasonal influenza viruses vary with antigenic drift. *Nature* **523**, 217–220. (doi:10.1038/nature14460)

10. Miller MA, Viboud C, Balinska M, Simonsen L. 2009 The signature features of influenza pandemics—implications for policy. *N. Engl. J. Med.* **360**, 2595–2598. (doi:10.1056/NEJMp0903906)

11. Longini IM, Halloran ME. 2005 Strategy for distribution of influenza vaccine to high-risk groups and children. *Am. J. Epidemiol.* **161**, 303–306. (doi:10.1093/aje/kwi053)

12. Mylius SD, Hagenaars TJ, Lugnér AK, Wallinga J. 2008 Optimal allocation of pandemic influenza vaccine depends on age, risk and timing. *Vaccine* **26**, 3742–3749. (doi:10.1016/j.vaccine.2008.04.043)

13. Reichert TA, Sugaya N, Fedson DS, Glezen WP, Simonsen L, Tashiro M. 2001 The Japanese experience with vaccinating schoolchildren against influenza. *N. Engl. J. Med.* **344**, 889–896. (doi:10.1056/NEJM200103223441204)

14. Gemmetto V, Barrat A, Cattuto C. 2014 Mitigation of infectious disease at school: targeted class closure vs school closure. *BMC Infect. Dis.* **14**, 1–10. (doi:10.1186/s12879-014-0695-9)

15. Milne GJ, Halder N, Kelso JK. 2013 The cost effectiveness of pandemic influenza interventions: a pandemic severity based analysis. *PLoS ONE* **8**, e61504. (doi:10.1371/journal.pone.0061504)

16. Cauchemez S, Ferguson NM, Wachtel C, Tegnell A, Saour G, Duncan B, Nicoll A. 2009 Closure of schools during an influenza pandemic. *Lancet Infect. Dis.* **9**, 473–481. (doi:10.1016/S1473-3099(09)70176-8)

17. Viboud C, Charu V, Olson D, Ballesteros S, Gog J, Khan F, Grenfell B, Simonsen L. 2014 Demonstrating the use of high-volume electronic medical claims data to monitor local and regional influenza activity in the US. *PLoS ONE* **9**, e102429. (doi:10.1371/journal.pone.0102429)

18. Centers for Disease Control and Prevention. 2016 Overview of influenza Surveillance in the United States. Technical Report, Centers for Disease Control and Prevention.

19. Ginsberg J, Mohebbi MH, Patel RS, Brammer L, Smolinski MS, Brilliant L. 2009 Detecting influenza epidemics using search engine query data. *Nature* **457**, 1012–1014. (doi:10.1038/nature07634)

20. Lampos V, Bie TD, Cristianini N. 2010 Flu detector—tracking epidemics on Twitter. In *Machine Learning and knowledge discovery in databases* (eds J Balcazar, F Bonchi, M Sebag, A Gionis), pp. 599–602, Berlin/Heidelberg, Germany: Springer.

21. Olson DR, Konty KJ, Paladini M, Viboud C, Simonsen L. 2013 Reassessing Google Flu Trends data for detection of seasonal and pandemic influenza: a comparative epidemiological study at three geographic scales. *PLoS Comput. Biol.* **9**, e1003256. (doi:10.1371/journal.pcbi.1003256)

22. HealthMap. FluNearYou 2017. See https://flunearyou.org/#!/.

23. London School of Hygiene and Tropical Medicine. FluSurvey, 2017. See https://flusurvey.net/en/results/.

24. Adler AJ, Eames KT, Funk S, Edmunds WJ. 2014 Incidence and risk factors for influenza-like-illness in the UK: online surveillance using Flusurvey. *BMC Infect. Dis.* **14**, 232. (doi:10.1186/1471-2334-14-232)

25. Perrotta D, Bella A, Rizzo C, Paolotti D. 2017 Participatory online surveillance as a supplementary tool to sentinel doctors for influenza-like illness surveillance in Italy. *PLoS ONE* **12**, 1–15. (doi:10.1371/journal.pone.0169801)

26. Smolinski MS, Crawley AW, Baltrusaitis K, Chunara R, Olsen JM, Wójcik O, Santillana M, Nguyen A, Brownstein JS. 2015 Flu near you: crowdsourced symptom reporting spanning 2 influenza seasons. *Am. J. Public Health* **105**, 2124–2130. (doi:10.2105/AJPH.2015.302696)

27. Diekmann O, Heesterbeek JAP, Roberts MG. 2010 The construction of next-generation matrices for compartmental epidemic models. *J. R. Soc. Interface* **7**, 873–885. (doi:10.1098/rsif.2009.0386)

28. Glass K, Mercer GN, Nishiura H, McBryde ES, Becker NG. 2011 Estimating reproduction numbers for adults and children from case data. *J. R. Soc. Interface* **8**, 1248–1259. (doi:10.1098/rsif.2010.0679)

29. Goldstein E, Apolloni A, Lewis B, Miller JC, Macauley M, Eubank S, Lipsitch M, Wallinga J. 2010 Distribution of vaccine/antivirals and the 'least spread line' in a stratified population. *J. R. Soc. Interface* **7**, 755–764. (doi:10.1098/rsif.2009.0393)

30. Wallinga J, Van Bovena M, Lipsitch M. 2010 Optimizing infectious disease interventions during an emerging epidemic. *Proc. Natl Acad. Sci. USA* **107**, 923–928. (doi:10.1073/pnas.0908491107)

31. Staniek M, Lehnertz K. 2008 Symbolic transfer entropy. *Phys. Rev. Lett.* **100**, 158101. (doi:10.1103/PhysRevLett.100.158101)

32. Schreiber T. 2000 Measuring information transfer. *Phys. Rev. Lett.* **85**, 19. (doi:10.1103/PhysRevLett.85.461)

33. Borge-Holthoefer J, Perra N, Goncalves B, Gonzalez-Bailon S, Arenas A, Moreno Y, Vespignani A. 2016 The dynamics of information-driven coordination phenomena: a transfer entropy analysis. *Sci. Adv.* **2**, e1501158–e1501158. (doi:10.1126/sciadv.1501158)

34. Granger CWJ. 1969 Investigating causal relations by econometric models and cross-spectral methods. *Econometrica* **37**, 424. (doi:10.2307/1912791)

35. Sugihara G, May R, Ye H, Hsieh C-h., Deyle E, Fogarty M, Munch S. 2012 Detecting causality in complex ecosystems. *Science* **338**, 496–500. (doi:10.1126/science.1227079)

36. Scarpino SV, Petri G. 2019 On the predictability of infectious disease outbreaks. *Nat. Commun.* **10**, 898. (doi:10.1038/s41467-019-08616-0)

37. Moriyama IM, Loy RM, Robb-Smith AH. 2011 History of the statistical classification of diseases and causes of death. Technical Report, Centers for Disease Control and Prevention.

38. U.S. Postal Service Office of Inspector General. 2013 The untold story of the ZIP code. tech. rep., United State Postal Services.

39. Gillespie DT. 1977 Exact stochastic simulation of coupled chemical reactions. *J. Phys. Chem.* **81**, 2340–2361. (doi:10.1021/j100540a008)

40. Jhung MA *et al.* 2011 Epidemiology of 2009 pandemic influenza a (H1N1) in the United States. *Clin. Infect. Dis.* **52**(suppl. 1), 13–26. (doi:10.1093/cid/ciq008)

41. Yang W, Lipsitch M, Shaman J. 2015 Inference of seasonal and pandemic influenza transmission dynamics. *Proc. Natl Acad. Sci. USA* **112**, 2723–2728. (doi:10.1073/pnas.1415012112)

42. Held L, Hofmann M, Höhle M, Schmid V. 2006 A two-component model for counts of infectious diseases. *Biostatistics* **7**, 422–437. (doi:10.1093/biostatistics/kxj016)

43. Boëlle P-Y, Ansart S, Cori A, Valleron A-J. 2011 Transmission parameters of the A/H1N1 (2009) influenza virus pandemic: a review. *Influenza Other Respir. Viruses* **5**, 306–316. (doi:10.1111/j.1750-2659.2011.00234.x)

44. Biggerstaff M, Jhung M, Kamimoto L, Balluz L, Finelli L. 2012 Self-reported influenza-like illness and receipt of influenza antiviral drugs during the 2009 pandemic, United States, 2009–2010. *Am. J. Public Health* **102**, 2009–2010. (doi:10.2105/AJPH.2012.300651)

45. Lizier JT, Prokopenko M. 2010 Differentiating information transfer and causal effect. *Eur. Phys. J. B* **73**, 605–615. (doi:10.1140/epjb/e2010-00034-5)

46. Laumann EO, Youm YM. 1999 Racial/ethnic group differences in the prevalence of sexually transmitted diseases in the United States: a network explanation. *Sex. Transm. Dis.* **26**, 250–261. (doi:10.1097/00007435-199905000-00003)

47. Tuite AR *et al.* 2010 Estimated epidemiologic parameters and morbidity associated with pandemic H1N1 influenza. *CMAJ: Can. Med. Assoc. J.* **182**, 131–136. (doi:10.1503/cmaj.091807)

48. Kaiser A, Schreiber T. 2002 Information transfer in continuous processes. *Physica D* **166**, 43–62. (doi:10.1016/S0167-2789(02)00432-3)

49. Barnett L, Barrett AB, Seth AK. 2009 Granger causality and transfer entropy are equivalent for gaussian variables. *Phys. Rev. Lett.* **103**, 2–5. (doi:10.1103/PhysRevLett.103.238701)

50. Takens F. 1981 Detecting strange attractors in turbulence. In *Dynamical systems and turbulence* (eds D Rand, L-S Young), pp. 366–381. Berlin, Germany: Springer.

51. Deyle ER, Maher MC, Hernandez RD, Basu S, Sugihara G. 2016 Global environmental drivers of influenza. *Proc. Natl Acad. Sci. USA* **113**, 13 081–13 086. (doi:10.1073/pnas.1607747113)

52. Baskerville EB, Cobey S. 2017 Does influenza drive absolute humidity? *Proc. Natl Acad. Sci. USA* **114**, E2270–E2271. (doi:10.1073/pnas.1700369114)

53. Sugihara G, Deyle ER, Ye H. 2017 Reply to Baskerville and Cobey: misconceptions about causation with synchrony and seasonal drivers. *Proc. Natl Acad. Sci. USA* **114**, E2272–E2274. (doi:10.1073/pnas.1700998114)

54. Lungarella M, Ishiguro K, Kuniyoshi Y, Otsu N. 2007 Methods for quantifying the causal structure of bivariate time series. *Int. J. Bifurcation Chaos* **17**, 903–921. (doi:10.1142/S0218127407017628)

55. Ay N, Polani D. 2008 Information flows in causal networks. *Adv. Complex Syst.* **11**, 17–41. (doi:10.1142/S0219525908001465)

56. Kamiński M, Ding M, Truccolo WA, Bressler SL. 2001 Evaluating causal relations in neural systems: Granger causality, directed transfer function and statistical assessment of significance. *Biol. Cybern* **85**, 145–157. (doi:10.1007/s004220000235)

57. Pahle J, Green AK, Dixon CJ, Kummer U. 2008 Information transfer in signaling pathways: a study using coupled simulated and experimental data. *BMC Bioinf.* **9**, 139. (doi:10.1186/1471-2105-9-139)

58. Steeg GV, Galstyan A. 2011 Information transfer in social media. *Entropy* **90292**, 1–8.

59. Lizier JT, Prokopenko M, Zomaya AY. 2008 Local information transfer as a spatiotemporal filter for complex systems. *Phys. Rev. E - Stat., Nonlinear, Soft Matter Phys.* **77**, 1–11. (doi:10.1103/PhysRevE.77.026110)

60. Williams PL, Beer RD. 2011 Generalized measures of information transfer. (http://arxiv.org/abs/1102.1507), pp. 1–6.

61. Marinazzo D, Pellicoro M, Stramaglia S. 2012 Causal information approach to partial conditioning in multivariate data sets. *Comput. Math. Methods Med.* **2012**, 1–8. (doi:10.1155/2012/303601)

62. Gog JR, Ballesteros S, Viboud C, Simonsen L, Bjørnstad ON, Shaman J, Chao DL, Khan F, Grenfell BT. 2014 Spatial transmission of 2009 Pandemic influenza in the US. *PLoS Comput. Biol.* **10**, e1003635. (doi:10.1371/journal.pcbi.1003635)