

## Research Article

# The Time Required to Estimate the Case Fatality Ratio of Influenza Using Only the Tip of an Iceberg: Joint Estimation of the Virulence and the Transmission Potential

Keisuke Ejima,<sup>1</sup> Ryosuke Omori,<sup>2,3</sup> Benjamin J. Cowling,<sup>3</sup>  
Kazuyuki Aihara,<sup>1,4</sup> and Hiroshi Nishiura<sup>3,5</sup>

<sup>1</sup> Department of Mathematical Informatics, Graduate School of Information Science and Technology, The University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113-8656, Japan

<sup>2</sup> Department of Biology, Faculty of Sciences, Kyushu University, 6-10-1 Hakozaki, Higashi-ku, Fukuoka 812-8581, Japan

<sup>3</sup> School of Public Health, The University of Hong Kong, Level 6, Core F, Cyberport 3, 100 Cyberport Road, Pokfulam, Hong Kong

<sup>4</sup> Institute of Industrial Science, The University of Tokyo, 4-6-1 Komaba, Meguro-ku, Tokyo 153-8505, Japan

<sup>5</sup> PRESTO, Japan Science and Technology Agency, Saitama 332-0012, Japan

Correspondence should be addressed to Hiroshi Nishiura, nishiura@hku.hk

Received 19 January 2012; Revised 21 February 2012; Accepted 22 February 2012

Academic Editor: Gerardo Chowell

Copyright © 2012 Keisuke Ejima et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Estimating the case fatality ratio (CFR) of a novel strain of influenza virus during the early stage of the pandemic is one of key epidemiological tasks to be conducted as rapid research response. Past experience during the epidemics of severe acute respiratory syndrome (SARS) and influenza A (H1N1-2009) posed several technical challenges in estimating the CFR in real time. The present study aimed to develop a simple method to estimate the CFR based on readily available datasets, that is, confirmed cases and deaths, while addressing some of the known technical issues. To assess the reliability and validity of the proposed method, we examined the minimum length of time required for the assigned CFR to be included within the 95% confidence intervals and for the estimated CFR to be below a prespecified cut-off value by means of Monte Carlo simulations. Overall, the smaller the transmission potential was, the longer it took to compare the estimated CFR against the cut-off value. If policymaking and public health response have to be made based on the CFR estimate derived from the proposed method and readily available data, it should be noted that the successful estimation may take longer than a few months.

## 1. Introduction

When a new infectious disease emerges, the case fatality ratio (CFR) informs how lethal the infection or the disease is, measuring the virulence of the novel infection as the conditional probability of death given infection or disease [1, 2]. To understand the severity of infection, assess the impact of clinical and public health interventions, and anticipate the likely number of deaths in the population given the total number of infected individuals, quantifying the CFR in real time during the early stage of an epidemic is of utmost importance.

Among various uses of the CFR, the present study focuses on influenza, and in particular, the epidemiological determination of the severity in relation to epidemiological

indices, such as the Pandemic Severity Index (PSI) in the United States [3]. As a process of public health policymaking, this index is used as a scientific criterion to choose and determine public health countermeasures, and thus, the CFR is regarded as key information for policymaking [4]. For instance, if the estimated CFR exceeds a prespecified reference value of 2.0%, which is sometimes quoted as the estimate of the CFR for the H1N1-1918 pandemic (Spanish influenza) [5], the PSI suggests that the government should recommend and implement all the nonpharmaceutical interventions listed, including voluntary isolation of clinically ill individuals at home, quarantine of household contacts, and social distancing [3].

However, while the CFR is theoretically calculated as a proportion of deaths to infected individuals, the actual

calculation practice involves several technical problems owing to a few epidemiological features.

First, one cannot directly count all infected individuals during the course of an epidemic due to unobservable nature of infection, and commonly available datasets may be only confirmed cases through surveillance efforts. Moreover, the mild nature of influenza involves multiple steps of bias including ignorance of asymptomatic and subclinical infections, case ascertainment bias, imperfectness of a diagnostic testing method, and reporting bias. In fact, approximately 10% of confirmed infections with influenza (H1N1-2009) in households were shown to be fully asymptomatic [6, 7]. To partly address the issue of underascertainment, a technical advancement in synthesizing epidemiological evidences enabled us to estimate the symptomatic case fatality ratio (sCFR), the proportion of deaths to all symptomatic cases [8], although the denominator data are frequently based on nonspecific disease information such as influenza-like illness. As an alternative, a real-time serological study could potentially offer the denominator based on all infected individuals [9], but the seroepidemiological survey is costly and the diagnostic performance of serological testing in relation to the estimation of the CFR has yet to be fully clarified. While specific CFRs using confirmed cases and symptomatic cases as the denominator have been expressed as cCFR and sCFR, respectively, among studies of H1N1-2009 [2], the present study consistently uses the simplest notation “CFR,” intending it to represent the risk of death among all infected individuals (and so may be abbreviated as iCFR when necessary).

Second, the real-time estimation of the CFR has to take into account the time delay from illness onset to death and, thus, requires us to employ an appropriate statistical method to address censoring. This point must be considered, because all the cases are not fully exposed to the risk of death at a point in time during the course of an epidemic, and a simple ratio of the cumulative numbers of deaths to cases can yield biased (mostly, underestimated) CFR [10–13]. Third, it is critical to always keep it in mind that the risk of death is heterogeneous. In particular, the higher risk of influenza death than healthy adults is seen among those with underlying health conditions [14], including chronic obstructive pulmonary disease, asthma, pregnancy, chronic kidney failure, and diabetes. Perhaps reflecting this feature, the CFR clearly differs by age with the highest estimate among elderly and infants and the lowest among school-age children and young adults [15, 16]. Fourth, during an early stage of a pandemic, the number of deaths still remains very small and so the estimation of CFR suffers from broad uncertainty. Given that the CFR of influenza is likely to be small, and suffers from wide uncertainty, it is fruitful to clarify the minimum number of cases that are required to determine if the estimated CFR in real time is significantly below a prespecified cut-off value such as 2.0%.

While directly addressing clinically mild features and case ascertainment bias calls for synthesizing epidemiological evidence, for example, by employing hierarchical modeling approach [8], it is also important to clarify what can be done with readily available epidemiological information, such as

confirmed cases and confirmed deaths. In the present study, we aim to propose a method to estimate the CFR based on the limited epidemiological data during the early stage of an epidemic. Through this investigation, we also aim to clarify the minimum number of days that are required to explicitly compare the estimated CFR to predetermined cut-off values of the CFR.

## 2. Materials and Methods

*2.1. Assumptions.* For clarity, here we describe the underlying epidemiological assumptions and settings. First, we focus on the early stage of a pandemic and ignore the depletion of susceptible individuals during this particular time period in which the number of newly infected individuals  $i(t)$  increases exponentially; that is, we focus on the log-linear phase alone for simplicity. Second, in realistic settings,  $i(t)$  cannot be directly counted as a function of time, and it is possible to observe only the confirmed cases,  $c(t)$ . Third, during the early epidemic phase, a constant factor,  $k$  which scales the ratio of confirmed cases to all infected individuals, is assumed to remain a constant. In other words, we assume that the frequency of confirmed diagnosis among infected individuals does not vary over time. Fourth, we assume that the time delay from infection to death is independently and identically distributed and denotes the conditional probability density function as  $f(s)$  of length  $s$  days since infection. Moreover, among the confirmed cases, we assume that all death counts are recorded over time through surveillance system. Except for being reflected in the generation time distribution and the reproduction number, the event of death is assumed to be independent of the process of renewal. Finally, we consider a public health setting in which the time of emergence (or the time to initiate exponential growth)  $t_0$  is known (even approximately) as was the case in a specific epidemic study in which the starting time point of an epidemic was estimated [17, 18]. In the next subsection, we describe the estimation procedure of only a homogeneous population. However, the estimation problem of a population with heterogeneous risks of infection and death is discussed in Section 3.

*2.2. Model Structure.* We first describe the model structure deterministically. Throughout the paper, we ignore demographic stochasticity in infection process (see Section 4). Let  $i(t)$  be the incidence of infection at calendar time  $t$ . Also, let  $t_0$  be the time at which an epidemic starts with a single index case. Then  $i(t)$  increase exponentially as follows:

$$i(t) = \exp\{r(t - t_0)\}, \quad (1)$$

where  $r$  is the exponential growth rate of incidence. It should be noted that  $r$  may be referred to as the intrinsic growth rate, if the incidence data in question deals with only the very early exponential growth phase of an epidemic. However, as long as we consider the exponential growth as crude approximation of the epidemic curve, the exponential

growth phase can be longer than that governed by the intrinsic growth rate and  $r$  is not restricted to be the very early epidemic phase (e.g., as was practiced previously [19, 20]). Let  $p$  be the CFR among all infected individuals. Assuming that the conditional probability density function of the time from infection to death  $f(s)$  is known, the number of deaths  $d(t)$  is modeled as

$$d(t) = p \int_0^\infty i(t-s)f(s)ds, \quad (2)$$

which can be rewritten as

$$\begin{aligned} d(t) &= p \exp\{r(t-t_0)\} \int_0^\infty \exp(-rs)f(s)ds \\ &= p \exp\{r(t-t_0)\}M(-r), \end{aligned} \quad (3)$$

where  $M(-r)$  represents the moment-generating function of the time from infection to death given the exponential growth rate  $r$ . One may integrate both sides and use the cumulative number of deaths by day  $t$ ,  $D(t)$  for the estimation of CFR. The estimator of the CFR is then given by

$$\hat{p} = \frac{D(t)}{[\exp\{r(t-t_0)\} - \exp(-rt_0)](M(-r)/r)}. \quad (4)$$

Other than parameters for  $f(s)$ , which we will assume as known, the estimator (4) indicates that, to estimate the CFR, an unknown parameter  $r$  has to be estimated from an additional series of data other than the death process, for example, from the confirmed case series. The exponential growth rate  $r$  quantifies the denominator of the above-mentioned estimator. To estimate  $r$ , we analyze the incidence of confirmed cases,  $c(t)$ . Let  $l$  denote the proportion of confirmed cases to the total of infected individuals the data-generating process of  $c(t)$  is described by

$$\begin{aligned} c(t) &= l \int_0^\infty i(t-s)h(s)ds \\ &= lQ(-r)i(t) \\ &= ki(t), \end{aligned} \quad (5)$$

where  $h(s)$  is the density function of the time from infection to confirmatory diagnosis and  $Q(-r)$  represents the

moment-generating function given the exponential growth rate  $r$ . We refer to the parameter  $k$  as the confirmed coefficient (i.e.,  $k = lQ(-r)$ ) which acts as a constant factor to translate  $i(t)$  into  $c(t)$ .

Let us consider an adjusted calendar time based on known  $t_0$  (i.e.,  $t + t_0$ ), we simplify all the following equations by eliminating  $t_0$  (and hereafter we consistently use  $t$  as the adjusted time in which  $t_0$  is equated to be zero). In addition to this adjustment, we discretize both the series of confirmed cases and deaths, because the observed dataset is given with discrete time (i.e., daily data), that is,

$$\begin{aligned} c_t &= \int_{t-1}^t c(x)dx \\ &= \frac{k}{r} [\exp(rt) - \exp\{r(t-1)\}], \\ d_t &= \int_{t-1}^t d(x)dx \\ &= p(\exp(rt) - \exp\{r(t-1)\}) \frac{M(-r)}{r}. \end{aligned} \quad (6)$$

The proposed estimation method based on these linear approximations would work even when the exponential growth rate varies with time (e.g., varies as a step function) as was considered when evaluating the effectiveness of public health interventions such as school closure [21, 22]. The number of parameters would have to increase to capture the time variation (e.g., from a single  $r$  alone to  $r_0$  and  $r_1$  for two consecutive epidemic phases), and thus the required sample size for estimation would also increase. However, all we have to do to cope with the time dependence is to update (1) and (2) using multiple growth rates and, thus, revise (6), accordingly.

**2.3. Maximum Likelihood Estimation.** Assuming that the observed number of confirmed cases on day  $t$  results from Poisson sampling process with mean  $c_t = k [\exp(rt) - \exp\{r(t-1)\}]/r$ , where  $r$  and  $k$  are parameters, the likelihood function is given as follows:

$$L_1(r, k; m_t) = \prod_{t=1}^T \frac{(k/r)^{m_t} [\exp(rt) - \exp\{r(t-1)\}]^{m_t} \exp[-(k/r)[\exp(rt) - \exp\{r(t-1)\}]]}{m_t!}, \quad (7)$$

where  $m_t$  is the observed daily number of confirmed cases on day  $t$  and  $T$  represents the latest time of observation.

Let  $\pi_t$  be a random variable which yields an estimator of the CFR on day  $t$  since the start of an epidemic and is the realized value in the particular epidemic. Assuming that the realized CFR is the result of binomial sampling process of death with sample size  $[\exp(rT) - 1]M(-r)/r$ , the likelihood to estimate the CFR based on the total number of deaths up

to the latest time of observation  $T$  is

$$\begin{aligned} L_2(\pi_T; D(T), r) &= \left( \frac{(\exp(rT) - 1) \frac{M(-r)}{r}}{D(T)} \right) \pi_T^{D(T)} \\ &\quad \times (1 - \pi_T)^{(\exp(rT) - 1)M(-r)/r - D(T)}. \end{aligned} \quad (8)$$

Because of an assumption of conditional independence between the renewal process and death, the total likelihood  $L$  is given by

$$L = L_1 L_2. \quad (9)$$

Minimizing the negative logarithm of the total likelihood  $L$ , we jointly estimate three parameters,  $\pi_T$ ,  $r$ , and  $k$ . The 95% confidence intervals (CIs) are derived from the profile likelihood, the idea of which is to invert a likelihood ratio test to obtain a CI for the parameter in question [23].

**2.4. Simulations.** Whereas the above-mentioned estimation procedure enables us to estimate the CFR based only on the confirmed cases and deaths, the estimation rests on limited epidemiological information as compared to other methods involving additional symptomatic case data or serological dataset. Thus, it is important to examine if we can overcome uncertainty and realistically employ the proposed method during an early phase of a pandemic. Specifically, we explore the time required to confidently suggest the range of the CFR and compare the CFR against a pre-specified cut-off value during the early stage. We assess the reliability and validity by means of random simulations.

As a plausible parameter range, we examine three different exponential growth rates,  $r$ , of 0.05, 0.15, and 0.25 per day. These are chosen as plausible, because, assuming that the mean generation time of influenza is 3 days and exponentially distributed, the basic reproduction number ranges from 1.15 to 1.75. If the generation time is a constant 3 days, the reproduction number ranges from 1.16 to 2.11. These are in line with published estimates of the reproduction number for H1N1-2009 [24]. In fact, the growth rate of influenza A (H1N1-2009) is estimated as 0.08 [25] and 0.10 per day [18] in Japan and Mexico, respectively. The reference values of CFR,  $p$ , are set at 0.1%, 0.5%, and 2.0% that are in line with the PSI in the United States. The CFR of Spanish influenza is sometimes thought to be approximately 2.0% [5, 24] and those of Asian and Hong Kong influenza pandemics are thought to be up to 0.5% [26]. The CFR of seasonal influenza is thought to be below 0.1% [27]. Although the CFR of the H1N1-2009 pandemic among all infected individuals is estimated to be smaller than 0.1% [28], we do not examine smaller estimates of the CFR, because 0.1% may be most reasonably defined as the lowest cut-off value in practical setting to distinguish a mild influenza strain from severe ones, and it is likely to be infeasible to robustly estimate a CFR below 0.1% during the early epidemic phase due to sampling errors. Since the empirically estimated proportion of confirmed cases among all infected individuals is 5% [15], we fix  $k$  at 0.05 assuming that the time delay from infection to confirmed diagnosis is sufficiently short. Ignoring small delay from infection to illness onset (as it does not influence the above-mentioned estimation framework), the conditional probability density function of the time from infection to death of influenza A (H1N1-2009)  $f(s)$  is assumed to be gamma distribution with the mean and standard deviation being 9.5 and 4.7 days, respectively [29].

As for simulation-based assessment, we first perform Monte Carlo simulations for 1000 times for each specified combination of parameter values, calculating the coverage probability of including the assigned CFR value within the 95% confidence intervals. Second, we assess the time at which the estimated CFR is confidently said as smaller than the pre-specified CFR value. Again, the model is randomly simulated for 1000 times per each parameters setting, calculating the number of simulation runs in which the upper 95% confidence interval is below the reference CFR value.

**2.5. Heterogeneous Population.** The proposed method can be extended to a heterogeneous population with differential risks of infection and death, perhaps by age and risk groups. We present the extended model analytically and demonstrate that the above-mentioned approach is directly applicable to a multihost population and, thus, the age-stratified epidemiological data.

### 3. Results

**3.1. Reliability.** Figure 1 illustrates a single simulation run and the resulting maximum likelihood estimates with the 95% confidence intervals with the assigned parameters, the CFR of 0.5%, and  $r = 0.15$  per day. As one can imagine, the confidence intervals for each parameter are gradually narrowed down as the epidemic progresses due to reduced sampling errors. During the very early stage of the pandemic (e.g., for the first 40 days given the assumed parameters), it is not feasible to expect the narrow confidence interval for the CFR, and thus, one may fail to assess the reliability using only the very limited early epidemiological data.

Table 1 shows the coverage probability of the CFR for each set of parameters. When estimating the CFR based on early epidemic data, the exponential growth rate appears to play a critical role in determining reliability. To attain the coverage probability greater than 90% with  $r = 0.05$ , 0.15, and 0.25 per day, respectively, the latest times of observation,  $T$ , should be at least 80, 40, and 30 days with the reference CFR value of 0.5%. This indicates that the smaller the transmission potential is, the longer the time it would take to obtain a reproducible estimate of the CFR. Of course, the coverage probability converges to 95% with longer observation times. Given a larger CFR, the coverage probability converges earlier due to smaller sampling errors. However, the coverage probability appears to be more sensitive to variation in  $r$  than that in the CFR value.

**3.2. Validity.** The validity of comparing the CFR against pre-specified cut-off values is summarized in Table 2. The overall qualitative patterns are similar to those of the coverage probability in Table 1. The minimum number of days that is required to declare that the CFR is below cut-off values is very sensitive to the exponential growth rate of cases. In other words, smaller transmission potential requires us to wait for longer time to compare the estimated CFR against the cut-off values. In addition, the validity is also sensitive to the estimated CFR relative to the cut-off value. When the relative ratio of the CFR to the cut-off value gets smaller,

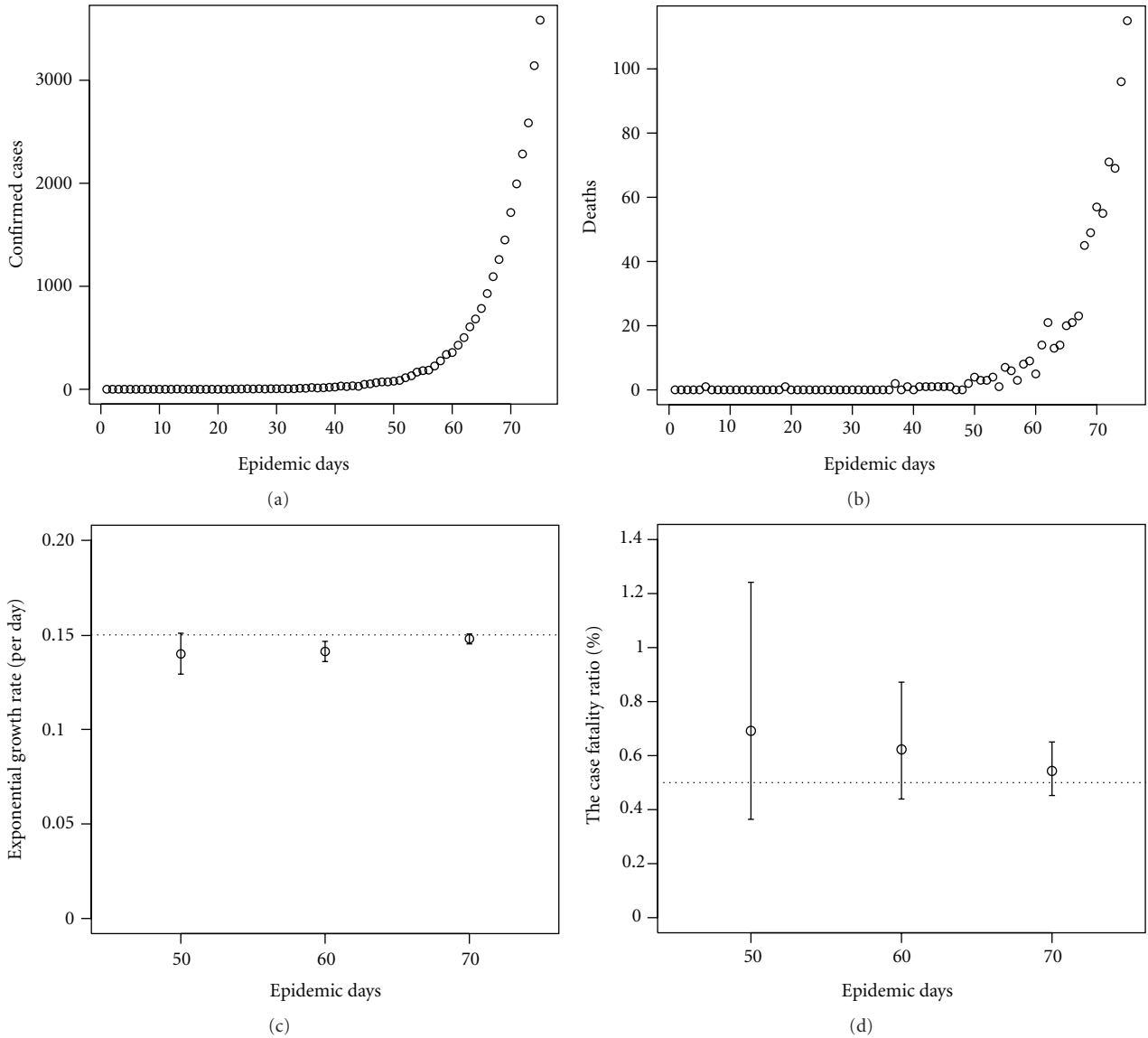


FIGURE 1: A single simulation run and the joint estimation results of the case fatality ratio and the exponential growth rate. The assigned CFR value is 0.5%, and the exponential growth rate  $r$  is set at 0.15 per day. ((a) and (b)) Incidence of confirmed cases and deaths as a function of epidemic days. The epidemic day 0 is the date on which an index case is infected. The numbers of confirmed cases and deaths increase exponentially. ((c) and (d)) The maximum likelihood estimates of (c) the exponential growth rate  $r$  and (d) the case fatality ratio with the 95% confidence intervals. The 95% confidence intervals were computed by employing the profile likelihood. Unfilled circles are the maximum likelihood estimates accompanied by the whiskers extending to lower and upper 95% confidence intervals. The dotted horizontal line shows the assigned parameter value.

the difficulty in differentiating the CFR is magnified. Given identical transmission potential and an identical assigned value of CFR, there was approximately a 20-day lag in the minimum numbers of days for differentiation between the relative case fatality ratios of 50% and 80%. With the smallest growth rate of  $r = 0.05$  per day, the estimation framework failed to yield any successful differentiation of CFR, even observing the epidemic for  $T = 100$  days. Of course, the estimated CFR also influences the feasibility, but the successful differentiation appears to be most sensitive to the exponential growth rate.

**3.3. Heterogeneous Population.** The proposed method is not directly applicable to realistic setting in which we observe substantial heterogeneities in the risks of infection and death. Accordingly, here we show the modeling approach to heterogeneous populations analytically. Specifically, we consider age-dependent dynamics: while the risk of infection may be higher among children than among elderly in the case of influenza, the conditional risk of death given infection is likely to be higher among elderly than school-age children, perhaps reflecting higher proportion of elderly with the underlying comorbidities.



TABLE 1: Coverage probability of the case fatality ratio (CFR) for each set of parameter values.

CFR	$r$	$T$ (epidemic day to estimate the CFR; days)							
		30	40	50	60	70	80	90	100
0.1%	0.05	1.5%	1.8%	7.8%	19.1%	34.2%	49.1%	66.0%	82.6%
	0.15	13.2%	52.4%	93.8%	94.6%	95.0%	94.9%	95.0%	*
	0.25	62.1%	94.4%	95.3%	*	*	*	*	*
0.5%	0.05	4.4%	8.0%	25.1%	57.3%	82.6%	91.9%	94.8%	95.7%
	0.15	55.3%	93.6%	95.2%	95.0%	95.9%	93.8%	94.2%	*
	0.25	94.4%	95.4%	95.3%	*	*	*	*	*
2.0%	0.05	15.2%	26.9%	51.8%	86.4%	94.4%	95.2%	95.1%	96.2%
	0.15	90.9%	95.3%	96.0%	96.4%	96.2%	95.9%	95.0%	*
	0.25	95.5%	95.8%	95.0%	*	*	*	*	*

CFR: case fatality ratio (assigned value). All the values were calculated as the proportion of successful simulation runs with the 95% confidence intervals that include the assigned CFR value among the total of 1000 simulation runs. The parameter  $r$  is the exponential growth rate of infection (per day).  $T$  is the epidemic date on which the estimation is performed. Those exceeding the coverage probability of 90% are highlighted in dark grey, while the cells with \* mark and shaded in light grey represent combinations of parameter values which generate too large numbers of cases and for which we refrained from estimation.

Let  $i_s(t)$  be the incidence of infection among subgroup  $s$  at calendar time  $t$ . Also, let  $R_{sq}$  be the average number of secondary cases in subgroup  $q$  generated by a single primary case in subgroup  $s$ , which would act as a single entry of the age-dependent next-generation matrix [30]. Assuming that the density function of the generation time  $g(\tau)$  of length  $\tau$  days is shared among subgroups, the multivariate renewal process is described by

$$i_s(t) = \sum_q R_{sq} \int_0^\infty i_q(t-s)g(s)ds. \quad (10)$$

Let  $p_s$  be the group-specific CFR (e.g., age-specific CFR) among all infected individuals of subgroup  $s$ . As was shown with application to the homogeneous population, we employ the confirmed coefficient  $k_s$ , reflecting both the proportion of confirmed cases to all infected individuals of subgroup  $s$  and the time delay from infection to confirmatory diagnosis. Then the confirmed cases among subgroup  $s$   $c_s(t)$  are

$$c_s(t) = k_s i_s(t). \quad (11)$$

Since the observed dataset is discrete time series, that is, daily data, we integrate the confirmed cases as follows:

$$c_{s,t} = \int_{t-1}^t c_s(x)dx. \quad (12)$$

Assuming that the conditional probability density function of the time from infection to death,  $f(s)$  is known and

is shared among subgroups, the number of new deaths of subgroup  $s$  at time  $t$   $d_s(t)$  is described as

$$d_s(t) = p_s \int_0^\infty i_s(t-s)f(s)ds, \quad (13)$$

which is rewritten as

$$d_s(t) = p_s \int_0^\infty \sum_q R_{sq} \int_0^\infty i_q(t-\tau-s)g(\tau)d\tau f(s)ds. \quad (14)$$

As was integrated in the homogeneous case, one may focus on the cumulative number of deaths  $D_s(T)$  by the latest time of observation  $T$ . The estimator of the group-specific CFR is given by

$$\hat{p}_s = \frac{D_s(T)}{\int_0^T \int_0^\infty \sum_q R_{sq} \int_0^\infty i_q(t-\tau-s)g(\tau)d\tau f(s)dsdt}. \quad (15)$$

The likelihood function to estimate the next-generation matrix may partially account for stochastic dependence structure of the transmission dynamics, and thus, conditions for every future expectation on the past history of the epidemic. Let  $\mathbf{Z}(t)$  represent the history of age-specific confirmed cases from time 0 up to time  $t-1$ . Given the series up to  $t-1$ , and assuming that the incidence of confirmed cases on day  $t$  is sufficiently characterized by Poisson distribution, the conditional likelihood is written as

TABLE 2: Proportion of simulation runs in which the upper 95% confidential interval of CFR ( $p$ ) falls below specified cut-off values.

Relative CFR	CFR	$r$	$T$ (epidemic day to estimate the CFR: days)							
			30	40	50	60	70	80	90	100
80%		0.05	0.0%	0.1%	0.3%	0.1%	0.0%	0.0%	0.0%	0.0%
	0.1%	0.15	0.0%	0.0%	1.5%	10.5%	31.3%	86.1%	100.0%	*
		0.25	0.0%	11.1%	67.5%	*	*	*	*	*
		0.05	0.4%	1.0%	1.7%	0.3%	0.7%	1.9%	5.9%	8.2%
	0.5%	0.15	0.0%	3.2%	11.5%	23.3%	64.8%	99.3%	100.0%	*
		0.25	6.6%	26.3%	98.1%	*	*	*	*	*
		0.05	2.5%	6.2%	6.6%	5.0%	6.9%	6.5%	7.2%	8.8%
	2.0%	0.15	2.1%	7.1%	13.2%	34.8%	80.4%	99.9%	100.0%	*
		0.25	10.0%	41.9%	99.6%	*	*	*	*	*
50%		0.05	0.0%	0.1%	0.1%	0.0%	0.0%	0.0%	0.0%	0.1%
	0.1%	0.15	0.0%	0.0%	3.3%	54.0%	99.2%	100.0%	100.0%	*
		0.25	0.0%	49.1%	100.0%	*	*	*	*	*
	0.5%	0.05	0.2%	0.6%	0.6%	0.6%	0.2%	3.4%	15.5%	27.8%
		0.15	0.0%	7.8%	45.3%	95.8%	99.9%	100.0%	100.0%	*
		0.25	17.7%	95.6%	100.0%	*	*	*	*	*
	2.0%	0.05	2.7%	5.1%	5.0%	8.1%	21.1%	28.7%	38.1%	48.7%
		0.15	3.7%	33.0%	74.0%	99.8%	100.0%	100.0%	100.0%	*
		0.25	44.8%	99.9%	100.0%	*	*	*	*	*
30%		0.05	0.0%	0.0%	0.2%	0.0%	0.0%	0.0%	0.0%	0.0%
	0.1%	0.15	0.0%	0.0%	2.7%	91.1%	100.0%	100.0%	100.0%	*
		0.25	0.0%	86.7%	100.0%	*	*	*	*	*
		0.05	0.3%	0.4%	0.4%	0.3%	0.4%	5.5%	24.5%	49.5%
	0.5%	0.15	0.0%	8.3%	84.7%	100.0%	100.0%	100.0%	100.0%	*
		0.25	29.6%	100.0%	100.0%	*	*	*	*	*
		0.05	1.4%	3.3%	4.2%	10.4%	33.9%	55.0%	72.8%	86.4%
	2.0%	0.15	4.3%	64.9%	98.7%	99.9%	100.0%	100.0%	100.0%	*
		0.25	86.1%	100.0%	100.0%	*	*	*	*	*

CFR: case fatality ratio (assigned value). Relative CFR: the ratio of assigned CFR value relative to the cut-off value. The proportion of successful simulation runs with the upper 95% confidence interval below the pre-specified cut-off value is shown. The parameter  $r$  is the exponential growth rate of infection (per day).  $T$  is the epidemic date on which the comparison is performed. Those exceeding the proportion of 90% are highlighted in dark grey, while the cells with \* mark and shaded in light grey represent combinations of parameter values which generate too large numbers of cases and for which we refrained from estimation.

$$L_1(\mathbf{R}_{ij}, k; \mathbf{Z}(t)) = \prod_{s=1}^n \prod_{i=1}^n \frac{\exp\left(-\sum_q R_{sq} \int_0^\infty c_q(t-s)g(s)ds\right) \left(\sum_q R_{sq} \int_0^\infty c_q(t-s)g(s)ds\right)^{m_{t,s}}}{m_{t,s}!}, \quad (16)$$

where  $m_{t,s}$  represents the observed number of confirmed cases in subgroup  $s$  on day  $t$ . This likelihood function is

useful to describe the underlying epidemic dynamics of the heterogeneous population. Let  $\pi_{t,s}$  be a random variable

which yields an estimator of the CFR of subgroup  $s$  at day  $t$  since the start of an epidemic. The other likelihood to estimate  $\pi_{t,s}$  is assumed to be given by binomial sampling process as follows:

$$\begin{aligned}
 & L_{2,s}(\pi_{s,T}; D_s(T), \mathbf{R}_{ij}, k) \\
 &= \left( \int_0^T \int_0^\infty \sum_q R_{sq} \int_0^\infty i_q(t-\tau-s)g(\tau)d\tau f(s)dsdt \right)^{D_s(T)} \\
 & \times \pi_{s,T}^{D_s(T)} (1 - \pi_{s,T})^{\int_0^T \int_0^\infty \sum_q R_{sq} \int_0^\infty i_q(t-\tau-s)g(\tau)d\tau f(s)dsdt - D_s(T)}.
 \end{aligned} \tag{17}$$

Therefore, the total likelihood is calculated as the following product:

$$L = L_1 \prod_s L_{2,s}. \tag{18}$$

Although the estimation framework can thus be very similar to that for the homogeneous population, it should be noted that the validity and reliability of estimation procedure for the heterogeneous population are likely to be influenced by way of parameterizing the next-generation matrix. For example, if the quantification of the matrix requires us to estimate only a small number of parameters (e.g., one parameter for each age group), we expect that the validity and reliability are not too much different from those we examined for the homogeneous population. However, when more parameters should be estimated to describe more detailed underlying heterogeneous transmission dynamics, the proposed method has to face greater uncertainty.

#### 4. Discussion

We proposed an estimation method to jointly infer the CFR and the exponential growth rate using only the confirmed case and death data. By means of Monte Carlo simulations, we assessed the minimum length of days required to compare the estimated CFR with the pre-specified CFR value such as those in the US Pandemic Severity Index. To do so, it appeared that the validity and reliability were very sensitive to the exponential growth rate and, thus, to the transmission potential of a novel pandemic strain. To be confident that the method included the CFR estimate within the 95% confidence interval, it appears that we have to wait at least for a month, and perhaps in general for about a few months given that the growth rate is equal to or smaller than 0.15 per day. The successful differentiation of CFR from cut-off values also takes about a few months. As it takes longer time for the estimation or differentiation, it would become more difficult for the incidence curve to be approximated by exponential growth. More importantly, the differentiation may not be feasible, if the growth rate is 0.05 per day or smaller. The growth rate was thus shown to play the most critical role in determining the feasibility of the proposed method than the CFR value to be estimated. This finding is attributable to the fact that the number of deaths is the result

of binominal sampling of cases. In general, as the sample size (or the number of binomial trials) increased, the standard error of binomial probability decreased, and the number of binomial trials in the proposed model substantially increased as the exponential growth rate increased. The validity and reliability were more sensitive to the growth rate than the binomial probability: the influence of variation in binomial probability on its confidence interval was small for the assumed range of CFR values, which can be understood from the approximate standard error of the binomial probability derived from the normal approximation to binomial [31].

Already, there have been multiple epidemiological methods to estimate the CFR using different datasets. Presanis et al. [8] proposed a Bayesian evidence synthesis approach using various different types of data that describe a pyramid structure, explaining that confirmed cases represent the tip of an iceberg of infected individuals and emphasizing a need to observe milder fraction of cases such as those who attended medical service. The useful datasets for that method included medically attended symptomatic cases and those required hospitalization and those admitted to intensive care unit. Comparing our proposed method with the evidence synthesis approach, the proposed method has two important advantages: (i) the proposed method can comply with a need to offer the CFR estimate in real time and (ii) we use the time series data of the confirmed cases and deaths which are readily available and accessible. The two different methods may thus be combined and used in practical setting: while the proposed method is used for the real-time assessment using widely available data, the sCFR estimate employing the evidence synthesis approach may be subsequently offered based on datasets of well-defined cohort populations. On the other hand, Wu et al. [15] used real-time seroprevalence data during the course of an epidemic. This approach enables us to estimate the background denominator of incidence of infections directly. In fact, a seroepidemiological study may be the only method to explicitly and directly quantify the underlying transmission dynamics. However, seroepidemiological surveys are costly and explicit interpretations of seroconversion and changes in antibody titers have yet to be offered. As a method to supplement the explicit estimation approaches, we believe that the proposed method based on readily available data would be a useful real-time assessment tool.

There are a few important future tasks for improvement. First, we used no prior information of parameters in the present study, but some information may be retrieved from other datasets or from the literature including historical epidemic records (e.g., the exponential growth rate of the epidemic caused by the same infectious agent before conducting the estimation). In fact, it is frequently the case that the transmission potential or the growth rate of cases is estimated earlier than the CFR in practical setting, and one may know a plausible value of  $r$  in advance of CFR estimation. If the prior knowledge could compensate unknown information of the proposed method, it will help to reduce the associated uncertainty of the CFR, thereby improving the validity of estimation. Second, we did not take into account demographic stochasticity in the present



study, but the stochasticity may not be negligible during the early epidemic phase [32–34]. The uncertainty that we quantified in the present study is likely to have been underestimated, although the qualitative findings are expected not to be different from those when we explicitly account for stochasticity using appropriate models (e.g., [32, 35]), especially for highly transmissible virus. Third, the proposed method as well as two earlier estimation studies based on evidence synthesis and serological study relied on the observed number of deaths as the numerator of the CFR. If there are many undiagnosed deaths, the direct estimation of the CFR is not feasible, and so, the virulence may also have to be assessed by indirect measurement such as that using excess mortality [36, 37]. Of course, constant  $k$  over time is also an unsupported assumption for epidemics with time-varying ascertainment efforts.

For a heterogeneous population we have shown that the proposed estimation framework for the homogeneous population can be easily extended to the heterogeneous setting. However, we have also discussed that the limited degrees of freedom might increase the relevant uncertainty; that is, when we consider  $n$  different subgroups, we have to deal with the next-generation matrix with  $n^2$  entries in addition to  $n$  unknown parameters for the group-specific CFR,  $p_s$ . Thus, the minimum length of days  $T$  that is required to estimate the CFR would be extended, and  $T$  would depend on the way we parameterize the next-generation matrix. Thus, we failed to offer simulation results with general conclusions with respect to the validity and reliability for the heterogeneous population. Given that the transmission of H1N1-2009 has been highly dependent on age [18, 32, 35, 38], one will have to balance the detailed descriptions of dynamics involving many subpopulations with the uncertainty surrounding the joint estimation of the CFR and the transmission potential.

If policymaking and public health response have to be made based on the real-time estimate of the CFR, the proposed method can be employed using only the readily available epidemiological datasets. However, as long as the estimation of the CFR relies on the proposed method, it should be noted that it may take longer than a few months to derive the CFR with sound uncertainty bounds, and thus, the very early response may not be able to base the policy decision on the CFR. Moreover, when the transmission potential is small, the number of infected individuals (or cases) may better be estimated directly from serological data (or medical attendance), because the proposed method is prone to uncertainties arising from low frequency of infection. While such limitation exists, we believe that the proposed method can be coupled with or supplement existing estimation frameworks which have to use additional epidemiological and serological data, especially for diseases with high transmission potential.

## Acknowledgments

H. Nishiura received funding support from the Japan Science and Technology Agency (JST) PRESTO program. K. Ejima and R. Omori received scholarship support from the Japan Society for Promotion of Science (JSPS). K. Aihara received

funding support from the Aihara Project, the FIRST program from JSPS, initiated by CSTP. This work also received financial support from the Harvard Center for Communicable Disease Dynamics from the National Institute of General Medical Sciences (Grant no. U54 GM088558) and the Area of Excellence Scheme of the Hong Kong University Grants Committee (Grant no. AoE/M-12/06). The funding bodies were not involved in the collection, analysis, and interpretation of data, the writing of the paper or the decision to submit for publication.

## References

- [1] J. Ma and P. Van Den Driessche, “Case fatality proportion,” *Bulletin of Mathematical Biology*, vol. 70, no. 1, pp. 118–133, 2008.
- [2] H. Nishiura, “Case fatality ratio of pandemic influenza,” *The Lancet Infectious Diseases*, vol. 10, no. 7, pp. 443–444, 2010.
- [3] Department of Health & Human Services, “Interim Pre-Pandemic Planning Guidance: Community Strategy for Pandemic Influenza Mitigation in the United States—Early, Targeted, Layered Use of Nonpharmaceutical Interventions,” Centers for Disease Control, Atlanta, Ga, USA, 2007, <http://healthvermont.gov/panflu/documents/0207interimguidance.pdf>.
- [4] M. Lipsitch, S. Riley, S. Cauchemez, A. C. Ghani, and N. M. Ferguson, “Managing and reducing uncertainty in an emerging influenza pandemic,” *New England Journal of Medicine*, vol. 361, no. 2, pp. 112–115, 2009.
- [5] K. G. Nicholson, R. G. Webster, and A. J. Hay, *Textbook of Influenza*, Blackwell Science, Oxford, UK, 1998.
- [6] J. Papenburg, M. Baz, M. È. Hamelin et al., “Household transmission of the 2009 pandemic A/H1N1 influenza virus: elevated laboratory-confirmed secondary attack rates and evidence of asymptomatic infections,” *Clinical Infectious Diseases*, vol. 51, no. 9, pp. 1033–1041, 2010.
- [7] B. J. Cowling, S. Ng, E. S. K. Ma et al., “Protective efficacy of seasonal influenza vaccination against seasonal and pandemic influenza virus infection during 2009 in Hong Kong,” *Clinical Infectious Diseases*, vol. 51, no. 12, pp. 1370–1379, 2010.
- [8] A. M. Presanis, D. De Angelis, A. Hagy et al., “The severity of pandemic H1N1 influenza in the united states, from April to July 2009: a bayesian analysis,” *Plos Medicine*, vol. 6, no. 12, Article ID 1000207, 2009.
- [9] J. T. Wu, A. Ho, E. S.K. Ma et al., “Estimating infection attack rates and severity in real time during an influenza pandemic: analysis of serial cross-sectional serologic surveillance data,” *PLoS Medicine*, vol. 8, no. 10, Article ID e1001103, 2011.
- [10] A. C. Ghani, C. A. Donnelly, D. R. Cox et al., “Methods for estimating the case fatality ratio for a novel, emerging infectious disease,” *American Journal of Epidemiology*, vol. 162, no. 5, pp. 479–486, 2005.
- [11] N. P. Jewell, X. Lei, A. C. Ghani et al., “Non-parametric estimation of the case fatality ratio with competing risks data: an application to severe acute respiratory syndrome (SARS),” *Statistics in Medicine*, vol. 26, no. 9, pp. 1982–1998, 2007.
- [12] H. Nishiura, D. Klinkenberg, M. Roberts, and J. A. P. Heesterbeek, “Early epidemiological assessment of the virulence of emerging infectious diseases: a case study of an influenza pandemic,” *Plos ONE*, vol. 4, no. 8, Article ID e6852, 2009.
- [13] T. Garske, J. Legrand, C. A. Donnelly et al., “Assessing the severity of the novel influenza A/H1N1 pandemic,” *British Medical Journal*, vol. 339, article b2840, 2009.

- [14] H. Nishiura, "The virulence of pandemic influenza A (H1N1) 2009: an epidemiological perspective on the case-fatality ratio," *Expert Review of Respiratory Medicine*, vol. 4, no. 3, pp. 329–338, 2010.
- [15] J. T. Wu, E. S. K. Ma, C. K. Lee et al., "The infection attack rate and severity of 2009 pandemic H1N1 influenza in Hong Kong," *Clinical Infectious Diseases*, vol. 51, no. 10, pp. 1184–1191, 2010.
- [16] S. Echevarría-Zuno, J. M. Mejía-Aranguré, A. J. Mar-Obeso et al., "Infection and death from influenza A H1N1 virus in Mexico: a retrospective analysis," *The Lancet*, vol. 374, no. 9707, pp. 2072–2079, 2009.
- [17] E. McBryde, I. Bergeri, C. van Gemert et al., "Early transmission characteristics of influenza A(H1N1)v in australia: victorian state, 16 May–3 June 2009," *Euro Surveillance*, vol. 14, no. 42, p. 19363, 2009.
- [18] C. Fraser, C. A. Donnelly, S. Cauchemez et al., "Pandemic potential of a strain of influenza A (H1N1): early findings," *Science*, vol. 324, no. 5934, pp. 1557–1561, 2009.
- [19] R. Omori and H. Nishiura, "Theoretical basis to measure the impact of short-lasting control of an infectious disease on the epidemic peak," *Theoretical Biology and Medical Modelling*, vol. 8, no. 1, article 2, 2011.
- [20] H. Nishiura and H. Inaba, "Estimation of the incubation period of influenza A (H1N1-2009) among imported cases: addressing censoring using outbreak data at the origin of importation," *Journal of Theoretical Biology*, vol. 272, no. 1, pp. 123–130, 2011.
- [21] H. Nishiura and K. Satou, "Potential effectiveness of public health interventions during the equine influenza outbreak in racehorse facilities in japan, 2007," *Transboundary and Emerging Diseases*, vol. 57, no. 3, pp. 162–170, 2010.
- [22] J. T. Wu, B. J. Cowling, E. H. Y. Lau et al., "School closure and mitigation of pandemic (H1N1) 2009, Hong Kong," *Emerging Infectious Diseases*, vol. 16, no. 3, pp. 538–541, 2010.
- [23] D. J. Venzon and S. H. Moolgavkar, "Method for computing profile-likelihood-based confidence intervals," *Journal of the Royal Statistical Society*, vol. 37, no. 1, pp. 87–94, 1988.
- [24] P. Y. Boëlle, S. Ansart, A. Cori, and A. J. Valleron, "Transmission parameters of the A/H1N1 (2009) influenza virus pandemic: a review," *Influenza and other Respiratory Viruses*, vol. 5, pp. 306–316, 2011.
- [25] H. Nishiura, G. Chowell, M. Safan, and C. Castillo-Chavez, "Pros and cons of estimating the reproduction number from early epidemic growth rate of influenza A (H1N1) 2009," *Theoretical Biology and Medical Modelling*, vol. 7, no. 1, article 1, 2010.
- [26] J. K. Taubenberger and D. M. Morens, "1918 influenza: the mother of all pandemics," *Emerging Infectious Diseases*, vol. 12, no. 1, pp. 15–22, 2006.
- [27] F. C. K. Li, B. C. K. Choi, T. Sly, and A. W. P. Pak, "Finding the real case-fatality rate of H5N1 avian influenza," *Journal of Epidemiology and Community Health*, vol. 62, no. 6, pp. 555–559, 2008.
- [28] S. Riley, K. O. Kwok, K. M. Wu et al., "Epidemiological characteristics of 2009 (H1N1) pandemic influenza based on paired sera from a longitudinal community cohort study," *Plos Medicine*, vol. 8, no. 6, Article ID e1000442, 2011.
- [29] H. Nishiura, "The relationship between the cumulative numbers of cases and deaths reveals the confirmed case fatality ratio of a novel influenza A (H1N1) virus," *Japanese Journal of Infectious Diseases*, vol. 63, no. 2, pp. 154–156, 2010.
- [30] O. Diekmann, J. A. P. Heesterbeek, and M. G. Roberts, "The construction of next-generation matrices for compartmental epidemic models," *Journal of the Royal Society Interface*, vol. 7, no. 47, pp. 873–885, 2010.
- [31] H. Nishiura, G. Chowell, and C. Castillo-Chavez, "Did modeling overestimate the transmission potential of pandemic (H1N1-2009)? sample size estimation for post-epidemic seroepidemiological studies," *Plos one*, vol. 6, no. 3, Article ID e17908, 2011.
- [32] H. Nishiura, C. Castillo-Chavez, M. Safan, and G. Chowell, "Transmission potential of the new influenza A(H1N1) virus and its age-specificity in japan," *Euro Surveillance*, vol. 14, no. 22, p. 19227, 2009.
- [33] T. Britton, "Stochastic epidemic models: a survey," *Mathematical Biosciences*, vol. 225, no. 1, pp. 24–35, 2010.
- [34] E. H.Y. Lam, B. J. Cowling, A. R. Cook, J. Y.T. Wong, M. S.Y. Lau, and H. Nishiura, "The feasibility of age-specific travel restrictions during influenza pandemics," *Theoretical Biology and Medical Modelling*, vol. 8, no. 1, article 44, 2011.
- [35] J. Lessler, N. G. Reich, D. A. T. Cummings, H. P. Nair, H. T. Jordan, and N. Thompson, "Outbreak of 2009 pandemic influenza A (H1N1) at a new york city school," *New England Journal of Medicine*, vol. 361, no. 27, pp. 2628–2636, 2009.
- [36] G. Chowell, L. M. A. Bettencourt, N. Johnson, W. J. Alonso, and C. Viboud, "The 1918-1919 influenza pandemic in england and wales: spatial patterns in transmissibility and mortality impact," *Proceedings of the Royal Society B*, vol. 275, no. 1634, pp. 501–509, 2008.
- [37] S. A. Richard, N. Sugaya, L. Simonsen, M. A. Miller, and C. Viboud, "A comparative study of the 1918–1920 influenza pandemic in Japan, USA and UK: mortality impact and implications for pandemic planning," *Epidemiology and Infection*, vol. 137, no. 8, pp. 1062–1072, 2009.
- [38] Y. Yang, J. D. Sugimoto, M. Elizabeth Halloran et al., "The transmissibility and control of pandemic influenza A (H1N1) virus," *Science*, vol. 326, no. 5953, pp. 729–733, 2009.