

Dynamics underlying global spread of emerging epidemics: An analytical framework

Lin Wang

WHO Collaborating Centre for Infectious Disease
Epidemiology and Control, School of Public Health,
Li Ka Shing Faculty of Medicine,
The University of Hong Kong
Hong Kong SAR, China
fdlwang@gmail.com

Joseph T Wu*

WHO Collaborating Centre for Infectious Disease
Epidemiology and Control, School of Public Health,
Li Ka Shing Faculty of Medicine,
The University of Hong Kong
Hong Kong SAR, China
joewu@hku.hk

ABSTRACT

Global spread of emerging epidemics (e.g. pandemic influenza, SARS, MERS-CoV, Ebola) is increasingly common, associated with the rapid pace of urbanization and global travel. Global metapopulation epidemic models built with worldwide air-transportation network (WAN) data have been one of the main tools for studying global spread of epidemics. However, it remains unclear how infectious disease epidemiology and the network properties of the WAN determine epidemic arrivals for different populations around the world. This work fills this knowledge gap by developing and validating an analytical framework on the basis of stochastic processes and network theory, which not only elucidates the dynamics underlying global spread of epidemics but also advances our capability in nowcasting and forecasting epidemics.

CCS CONCEPTS

- **Applied computing** → **Transportation; Forecasting;**
- **Mathematics of computing** → **Stochastic processes;**
- **Networks** → **Network dynamics; Network mobility;**
- **Computing methodologies** → **Network science;**

KEYWORDS

Emerging epidemics, Spatial epidemiology, Metapopulation epidemic models, Worldwide air-transportation network, Epidemic arrival time, Nonhomogeneous Poisson process

ACM Reference Format:

Lin Wang and Joseph T Wu. 2018. Dynamics underlying global spread of emerging epidemics: An analytical framework. In *Proceedings of The 2018 ACM SIGKDD International Workshop on Epidemiology meets Data Mining and Knowledge Discovery (KDD2018)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

*Materials are available on request from the corresponding author: Joseph T Wu

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
KDD2018, August 2018, London, United Kingdom
© 2018 Copyright held by the owner/author(s).
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM.
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Recent decades, global spread of emerging epidemics is increasingly common, as exemplified by the spread of SARS to nearly 30 countries in 2003, the spread of influenza A/H1N1 pandemic to more than 100 countries in 2009, the exportation of Ebola cases from the West Africa to the Nigeria, United States and United Kingdom in 2014, and recent geographical expansion of vector-borne diseases such as Dengue and Zika virus. Such frequent outbreaks of emerging epidemics are associated with the rapid pace of urbanization and global travel [5, 9, 14, 17]. In response to the serious situation, the World Health Organization (WHO) regularly updates the blueprint list of priority diseases to guide public health research and preparedness [28].

Since the 1980s, metapopulation epidemic models built with worldwide air-transportation network (WAN) data have been one of the main tools for studying global spread of emerging epidemics [12, 19, 24]. Despite their long history and widespread use, most studies in this field rely on computationally intensive simulations to predict or forecast the spatiotemporal transmission of epidemics [17, 19, 24]. However, one downside of such simulation-based methodology is that the computational process tends to be a black box – the underlying dynamics is hard to be elucidated from the basic principles in infectious disease epidemiology and network theory. In particular, an analytical understanding of the underlying dynamics has only been partially elucidated in recent years [4, 10, 23]. To fill this knowledge gap, we develop a novel analytical framework for characterizing how global spread of emerging epidemics depends on epidemiological parameters and the network properties of the WAN.

2 GLOBAL SPREAD SIMULATIONS: METAPOPOPULATION EPIDEMIC MODELS

2.1 Structure of the metapopulation epidemic models.

Metapopulation epidemic models are often described as a complex network of populations, in which each population denotes a city in the world and populations are interconnected through the mobility of individuals via the WAN [19, 24].

Since emerging infectious diseases generally evoke an epidemic with relatively fast timescales, we assume that in each population the epidemic peaks within 300 days after the establishment of the disease in that population [25]. It indicates that the change in demographics (e.g. births, aging) is negligible, such that each population has a constant population size. Denote population i as the epidemic origin with s_i initial infections seeded at time 0. For any given population j , the population size is denoted by N_j , with initial epidemic growth rate denoted by λ_j . For populations j and k that are directly connected, the per capita mobility rate from j to k is computed by $w_{jk} = F_{jk}/N_j$, in which F_{jk} is the daily number of passengers travelled by direct flights from j to k . Denote T_{ij}^n as the time at which population j receives its n th imported infection, such that T_{ij}^1 denotes the epidemic arrival time (EAT) for population j . **Table 1** summarizes the parameters.

2.1.1 Local epidemic dynamics within each population. The spread of epidemics within each population is modelled with frequency-dependent compartmental epidemic models [16], in which the transmission rate for infectious people to infect others can depend on multiple factors including the interpersonal contact rates, pathogenicity and environmental suitability [1, 6, 18, 31]. In the main text, we use the standard *SIR* model to describe the local epidemic dynamics within each population. Appendix A.1 extends to more general epidemic dynamics modelled by *SE_mI_nR* models.

Let $S_i(t)$, $I_i(t)$ and $R_i(t)$ be the number of susceptible, infectious and recovered people in a given population i at time t . Suppose $R_{0,i}$ is the basic reproductive number and $T_{g,i}$ is the mean generation time in population i . Let $\beta_i = R_{0,i}/T_{g,i}$ be the disease transmission rate and $\mu_i = 1/T_{g,i}$ be the recovery rate in population i . The *SIR* model is described by the following differential equations:

$$\begin{aligned}\frac{dS_i(t)}{dt} &= -\beta_i \frac{S_i(t)}{N_i} I_i(t), \\ \frac{dI_i(t)}{dt} &= \beta_i \frac{S_i(t)}{N_i} I_i(t) - \mu I_i(t), \\ \frac{dR_i(t)}{dt} &= \mu I_i(t).\end{aligned}$$

The doubling time $T_{d,i}$ for disease prevalence to have a two-fold increase (i.e. $I_i(T_{d,i}) = 2s_i$) is expressed by $\log(2) \frac{T_{g,i}}{(R_{0,i}-1)}$.

2.1.2 Stochastic mobility of individuals between populations. The spread of epidemics between populations results from the travel of infected individuals via the WAN. From a given population i , each individual travels to a directly connected population j at a small time interval Δt with probability $w_{ij}\Delta t = F_{ij}\Delta t/N_i$. Suppose population i is directly connect to multiple populations in the WAN, the numbers of susceptible, infectious and recovered travelers that leave population i through an interval Δt , i.e. $X_i(t)$, $Y_i(t)$ and $Z_i(t)$, are simulated with the following set of multinomial random

Table 1: Parameters of the two-population model in which the epidemic origin population i is only connected to population j .

Parameter	Definition
$I_i(t)$	Disease prevalence (number of infectives) in population i at time t
λ_i	Local epidemic growth rate in the origin population i
s_i	Number of initial infections seeded into the origin population i at time 0
w_{ij}	Daily per capita mobility rate from population i to j
α_{ij}	Adjusted mobility rate $\alpha_{ij} = s_i w_{ij}$
T_{ij}^n	The n th arrival time in population j

variables:

$$\begin{aligned}X_i(t) &= \text{Multinomial}(\lfloor S_i(t) \rfloor, w_{i1}\Delta t, \dots, w_{iG}\Delta t), \\ Y_i(t) &= \text{Multinomial}(\lfloor I_i(t) \rfloor, w_{i1}\Delta t, \dots, w_{iG}\Delta t), \\ Z_i(t) &= \text{Multinomial}(\lfloor R_i(t) \rfloor, w_{i1}\Delta t, \dots, w_{iG}\Delta t),\end{aligned}$$

where G counts the number of populations in the WAN, and *Multinomial*(n, p_1, \dots, p_G) denotes a multinomial random variable with n trials and probabilities p_1, \dots, p_G [20]. As such, the number of individuals given a specific disease compartment that that travel from population i to j per time interval (e.g. $X_{ij}(t)$) corresponds to the j th component of the corresponding multinomial random variable (e.g. $X_i(t)$).

2.2 Data-driven global metapopulation simulator.

To validate our analytical framework which will be introduced in section 3, we first develop a global metapopulation epidemic simulator, using the algorithm described in section 2.1. Our simulator contains 2,309 populations and 54,106 direct connections. Its structure is similar to the state-of-the-art simulator GLEAM [22] (but without the effect of local commuting which is less important to study the global spread [3]). To build this simulator, we use the 2015 worldwide flight booking data from the Official Airline Guide (OAG, <https://www.oag.com>) and the Gridded Population of the World Version 4 (GPWv4) dataset from the Columbia University [7]. The OAG dataset provides all flight booking records from all commercial airlines worldwide during 2015, and the GPWv4 dataset provides the highest resolution census data from the 2010 round of Population and Housing Censuses that were collected from hundreds of national statistics departments and organizations.

Ideally, the metapopulation dynamics described in section 2.1 is best implemented with discrete-event simulation algorithms (e.g. Gillespie algorithm [11]). However, explicitly simulating every event of individual infection, recovery and mobility substantially increases the computational burden, which largely exceeds the power of our high-performance

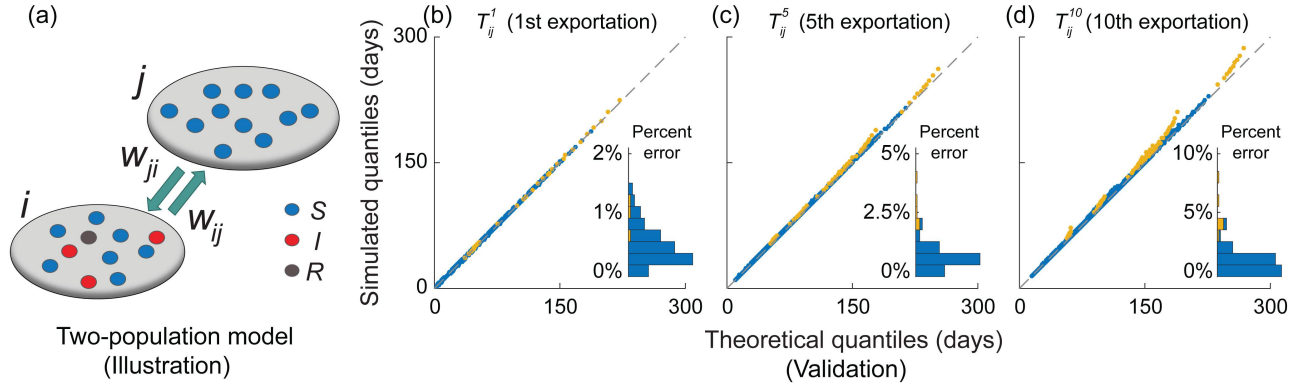


Figure 1: Validating the two-population analytics. (a) Illustration of the two-population model, with the epidemic origin population i only connecting to population j . Table 1 summarizes the parameters. (b)-(c) Q-Q plots for the analytical and simulated quantiles of T_{ij}^1 , T_{ij}^5 , and T_{ij}^{10} across 100 epidemic scenarios randomly sampled from the following parameter space using Latin-hypercube sampling: doubling time $T_{d,i}$ and generation time $T_{g,i}$ both between 3 and 30 days, seed size s_i between 1 and 100. Each of the 100 epidemic scenarios is coupled with a set of network parameters randomly sampled with mobility rate w_{ij} between 10^{-6} and 10^{-3} and population size N_i between 0.1 and 10 million, which are chosen according to the OAG and GPWv4 data [25]. Simulated quantiles for each of the 100 scenarios are compiled using 10,000 stochastic realizations. In the Q-Q plots, if data points coincide with the diagonal, the arrival times in the analytical framework are essentially the same as that in the simulation. Data points are colored in blue if the number of exportations X_{ij} is n or above with probability 1 (i.e. $P(X_{ij} \geq n) = 1$), and yellow otherwise. Insets show the corresponding histograms of percent error in $E[T_{ij}^n]$.

computing resources. To facilitate the stochastic computing of our global metapopulation epidemic simulator, we use a discrete-time algorithm in which the intra-population epidemic dynamics (see section 2.1.1) and inter-population mobility of travelers (see section 2.1.2) are sequentially simulated for each small time interval Δt . Throughout this work, we set $\Delta t = 0.05$ days, which is sufficiently small to ensure the accuracy of discrete-time simulations [30].

3 ANALYTICAL FRAMEWORK

We formulate the framework by analytically characterizing the probability distribution of EATs for all populations in three metapopulation models with increasingly complex network structure: (i) the simplest two-population model; (ii) the shortest-path-tree of the WAN (WAN-SPT hereafter); and (iii) the whole WAN.

3.1 The two-population model

We start from the two-population model in which the origin population i is only connected to population j (see Fig. 1a and Table 1 for model structure and parameters). This simple model corresponds to the initial stage of a pandemic with infections localized at the origin population (i.e. all the other populations can be merged as a single population that is unaffected to the disease [2]). Our analytical framework grounds on the following two key assumptions [10, 23, 25]:

- (1) Exportation of infections from population i to j is a nonhomogeneous Poisson process (NPP) [20] with intensity function $w_{ij}I_i(t)$, i.e. the expected number of infections exported from population i to j at time t .
- (2) After the epidemic has established in the origin population i , the first few exportations from population i to j occur while disease prevalence is still growing exponentially in the origin i , i.e. $I_i(t) = s_i \exp(\lambda_i t)$.

Under these assumptions, the probability density function (pdf) of T_{ij}^n can be expressed in closed-form:

$$f_n(t|\lambda_i, \alpha_{ij}) = \left(\frac{\exp(\lambda_i t) - 1}{\lambda_i} \right)^{n-1} \frac{\alpha_{ij}^n}{(n-1)!} \exp \left[\lambda_i t - \frac{\alpha_{ij}}{\lambda_i} (\exp(\lambda_i t) - 1) \right], \quad (1)$$

where $\alpha_{ij} = s_i w_{ij}$ is termed as adjusted mobility rate. To validate this two-population analytics, we compare the analytical and simulated arrival times for a wide range of epidemic scenarios (e.g. the doubling time and generation

time both between 3 and 30 days), which are eligible to describe emerging epidemics ranging from pandemic influenza (with doubling time around 4-5 days) to Ebola (with doubling time longer than 20 days). Figs. 1(b)-(d) show that Eq. (1)

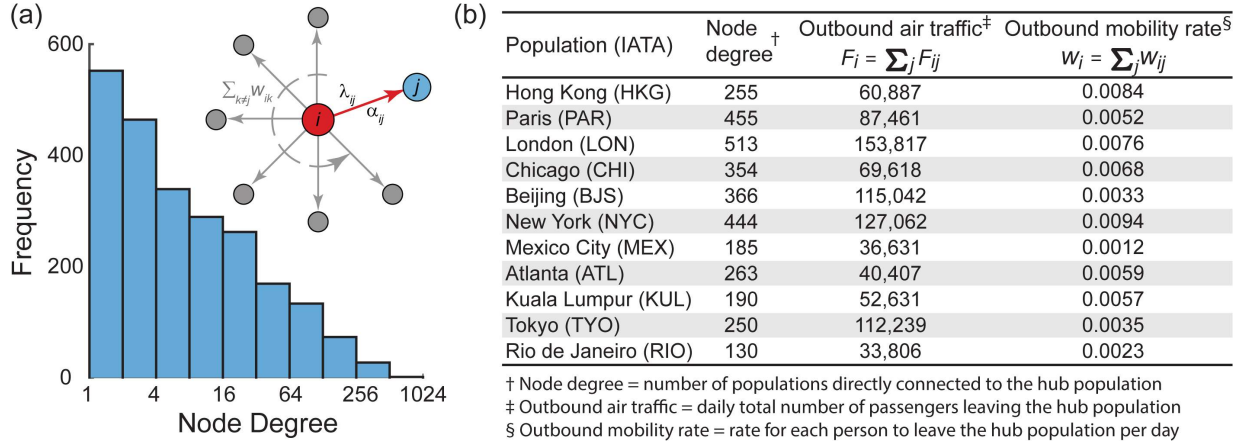


Figure 2: Network properties of the hub populations. (a) Histogram shows the distribution of node degree for all populations in the WAN. The node degree of a given population counts the number of populations that are directly connected to that population. Inset illustrates the structure of a travel hub, in which the hub population i is connected to multiple populations, one of which is population j . (b) Illustration of several major hubs in different continents by reporting their node degree, daily outbound traffic volume, and daily outbound per capita mobility rate.

accurately characterizes the arrival times T_{ij}^n for n up to 10 (i.e. the 10th exportation). With Eq. (1), we have the following corollaries:

- (1) Exportation of the first n infections is essentially an NPP with intensity function $\alpha_{ij} \exp(\lambda_i t)$.
- (2) The cumulative distribution function (cdf) of the n th arrival time is given by

$$F_n(t|\lambda_i, \alpha_{ij}) = \Gamma\left[n, \frac{\alpha_{ij}}{\lambda_i} (\exp(\lambda_i t) - 1)\right], \quad (2)$$

where Γ is the lower incomplete Gamma function.

- (3) The expected EAT is given by

$$E[T_{ij}^1] = \frac{1}{\lambda_i} \exp\left(\frac{\alpha_{ij}}{\lambda_i}\right) E_1\left(\frac{\alpha_{ij}}{\lambda_i}\right), \quad (3)$$

where $E_m(x) = x^{m-1} \int_x^\infty \left[\frac{\exp(-u)}{u^m}\right] du$ is the exponential integral.

- (4) If $\alpha_{ij} \ll \lambda_i$ and γ denotes the Euler constant, the expected EAT can be approximated as

$$E[T_{ij}^1] \approx \frac{1}{\lambda_i} \left[\ln\left(\frac{\lambda_i}{\alpha_{ij}}\right) - \gamma \right], \quad (4)$$

which is congruent with the EAT statistic in Gautreau et al. for estimating the order of epidemic arrival across different populations [10].

- (5) The expected time of the n th arrival is given by

$$E[T_{ij}^n] = \frac{1}{\lambda_i} \exp\left(\frac{\alpha_{ij}}{\lambda_i}\right) \sum_{m=1}^n E_m\left(\frac{\alpha_{ij}}{\lambda_i}\right). \quad (5)$$

- (6) For any positive integers m and n ($m < n$), the pdf of $T_{ij}^n - T_{ij}^m$ conditional on T_{ij}^m is simply

$$f_{n-m}(t|\lambda_i, \alpha_{ij} \exp(\lambda_i T_{ij}^m)) \quad (6)$$

which corresponds to the time of the $(n - m)$ th exportation for an epidemic with seed size $s_i \exp(\lambda_i T_{ij}^m)$. Using this relation recursively, we deduce that the joint pdf of $T_{ij}^1 = t_1, \dots, T_{ij}^n = t_n$ is simply

$$\prod_{m=1}^n f_1(t_m|\lambda_i, \alpha_{ij} \exp(\lambda_i t_{m-1})) \quad (7)$$

for all $0 = t_0 < t_1 < t_2 < \dots < t_{n-1} < t_n$.

- (7) The expected time of the $(n - 1)$ th exportation given an epidemic that starts at time T_{ij}^1 with seed size $s_i \exp(\lambda_i T_{ij}^1)$ is given by

$$E[T_{ij}^n | T_{ij}^1] = T_{ij}^1 + \frac{1}{\lambda_i} \exp\left(\frac{\alpha_{ij} \exp(\lambda_i T_{ij}^1)}{\lambda_i}\right) \sum_{m=1}^{n-1} E_m\left(\frac{\alpha_{ij} \exp(\lambda_i T_{ij}^1)}{\lambda_i}\right) \quad (8)$$

These corollaries are essential for extending our framework to the WAN-SPT and WAN analysis (see the following two sections).

3.2 The shortest-path tree of the WAN

The WAN-SPT is the dominant sub-network (or backbone) of the WAN, in which each downstream population connects to the epidemic origin via only one path. Brockmann et al. [4, 15] suggested that the epidemic spreads from the origin population to the other populations in the WAN through the WAN-SPT, such that global spread of epidemics through the WAN is primarily driven by the WAN-SPT. We will show that for each population k in the WAN-SPT, the n th arrival time T_{ik}^n can be accurately characterized by the two-population analytics of **Eq. (1)**, where the local epidemic growth rate and adjusted mobility rate are specifically parameterized to account for the hub effect (see section 3.2.1) and continuous seeding effect (see section 3.2.2).

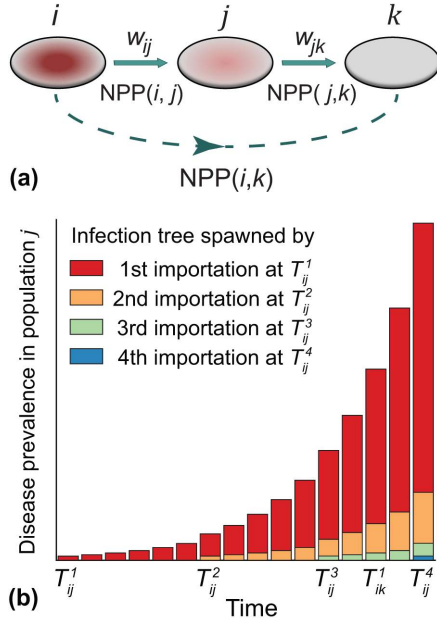


Figure 3: Effect of continuous seeding. (a) Illustration of the epidemic arrival process through an acyclic path that connects the epidemic origin i to population k via population j (i.e. $\psi : i \rightarrow j \rightarrow k$). (b) In this example, the epidemic arrives at population k after population j has imported three infections from the epidemic origin, i.e. $T_{ij}^3 < T_{ik}^1 < T_{ij}^4$. In the absence of continuous seeding adjustment, infection trees spawned by the second and subsequent importations in population j are ignored [10, 25].

3.2.1 Hub effect. Travel hubs such as Hong Kong, London and Paris have direct flights to multiple populations in the WAN (i.e. their node degree > 1 , see Fig. 2 for illustrations). Given a hub population i , the growth of local disease prevalence $I_i(t)$ can be substantially decreased if a significant proportion of infections travel outward as the epidemic unfolds. To extend our framework to deal with hub populations, we account for the reduction in local disease prevalence by wiping off the hub-effect from local epidemic growth rate λ_i .

Suppose a hub population i is directly connected to two or more populations, one of which is population j (see Fig. 2(a)). From the perspective of case arrival process for population j , disease prevalence in population i grows exponentially at rate $\lambda_{ij} = \lambda_i - \sum_{k \neq j} w_{ik}$. Therefore, the pdf of the n th arrival time for population j can be estimated with hub-adjusted two-population analytics $f_n(t|\lambda_{ij}, \alpha_{ij})$, in which infections are exported from population i to j at a rate $w_{ij}I_i(t)$ and disease prevalence in hub population i grows exponentially at the effective growth rate λ_{ij} . Using the hub structure of Hong Kong as an example, Fig. 4(a) show that hub-adjusted two-population analytics accurately characterizes the probability distribution of T_{ij}^n for all populations that are directly connected to Hong Kong.

3.2.2 Continuous seeding. Unlike the epidemic origin population which has a single seeding event at time 0, all the other populations in the WAN-SPT can be continuously seeded by infections coming from their upstream populations (illustrated in Fig. 3), as exemplified by recent multiple case importations of Zika Virus in Florida that come from the Caribbean [13].

Let D_c be the set of populations that are c degrees of separation from the epidemic origin in the WAN-SPT. Suppose a population k in D_2 is connected to the epidemic origin via population j along the path $\psi : i \rightarrow j \rightarrow k$. After the epidemic has arrived at population j at time T_{ij}^1 , population i continues to export infections to population j before the epidemic arrives at population k at time T_{ik}^1 (illustrated in Fig. 3). According to the two-population model, each imported infection in population j (arriving at times $T_{ij}^1, T_{ij}^2, \dots$) spawns an infection tree that grows exponentially at the hub-adjusted rate λ_{jk} . Therefore, the overall disease prevalence in population j , namely $I_j(t)$, is simply the sum of disease prevalence for all these infection trees:

$$I_j(t) = \sum_{m=1}^{\infty} \mathbf{I}\{t > T_{ij}^m\} \exp(\lambda_{jk}(t - T_{ij}^m))$$

where T_{ij}^m is the m th arrival time in population j , and $\mathbf{I}\{\cdot\}$ is the indicator function. Based on the two-population model, the exportation of infections from population j to k is an NPP with intensity function $w_{jk}I_j(t)$, which is itself a stochastic process because of its dependence on the random variables $T_{ij}^1, T_{ij}^2, \dots$. As such, conditional on $I_j(t)$ and hence $T_{ij}^1, T_{ij}^2, \dots$, the pdf of T_{ik}^n is

$$g_n(t|w_{jk}I_j) = f_{Poisson}\left(n-1, w_{jk} \int_0^t I_j(u) du\right) w_{jk}I_j(u)$$

for $n = 1, 2, \dots$. The unconditional pdf of T_{ik}^n is thus

$$E_{T_{ij}^1, T_{ij}^2, \dots} [g_n(t|w_{jk}I_j)]$$

which integrates over the joint pdf of $(T_{ij}^1 = t_1, T_{ij}^2 = t_2, \dots)$.

We conjecture that this highly complex stochastic process can be substantially simplified with little loss of accuracy by using the following assumption: Conditional on T_{ij}^1 (i.e. the EAT for population j), $T_{ij}^m \approx E[T_{ij}^m | T_{ij}^1]$ for all $m > 1$ (see Eq. 8). Therefore, conditional on T_{ij}^1 , we approximate $I_j(t)$ with the following certainty equivalent approximation (CEA):

$$\begin{aligned} I_j^{CEA}(t) &= \sum_{m=1}^{\infty} \mathbf{I}\{t > E[T_{ij}^m | T_{ij}^1]\} \exp(\lambda_{jk}(t - E[T_{ij}^m | T_{ij}^1])) \\ &= \exp(\lambda_{jk}(t - T_{ij}^1)) \sum_{m=1}^{\infty} \mathbf{I}\{t > T_{ij}^1 + \Delta T_{ij}^m\} \exp(-\lambda_{jk} \Delta T_{ij}^m) \end{aligned}$$

where

$$\begin{aligned} \Delta T_{ij}^m &= E[T_{ij}^m | T_{ij}^1] - T_{ij}^1 \\ &= \frac{1}{\lambda_{ij}} \exp\left(\frac{\alpha_{ij} \exp(\lambda_{ij} T_{ij}^1)}{\lambda_{ij}}\right) \sum_{q=1}^{m-1} E_q\left(\frac{\alpha_{ij} \exp(\lambda_{ij} T_{ij}^1)}{\lambda_{ij}}\right) \end{aligned}$$

(see Eq. (8)). The resulting unconditional pdf of T_{ik}^n is simply $E_{T_{ij}^1} [g_n(t|w_{jk}I_j^{CEA})]$ where the pdf of T_{ij}^1 is $f_1(\cdot|\lambda_{ij}, \alpha_{ij})$ (see Eq. (1)).

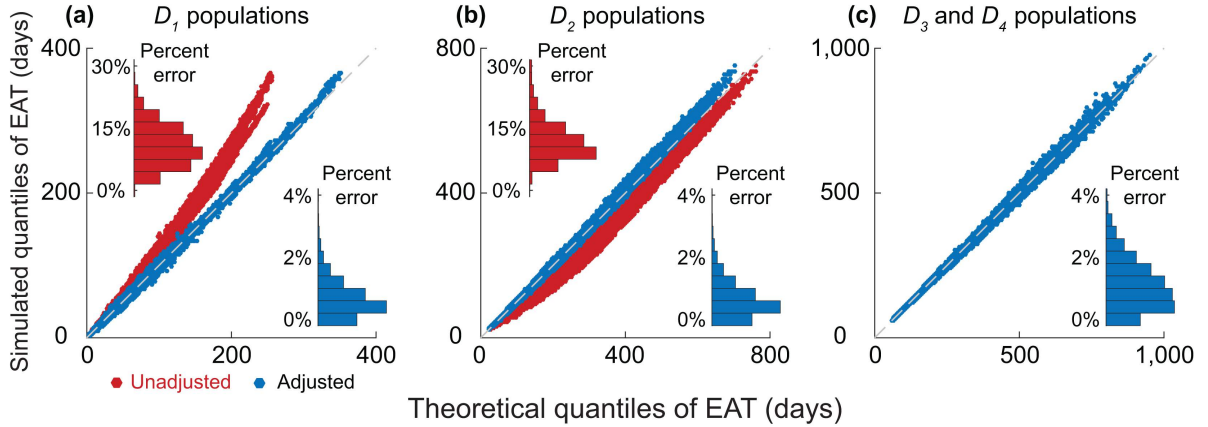


Figure 4: Validating the analytical framework for the WAN-SPT with Hong Kong as the epidemic origin (WAN-SPT-HK). (a)-(c) Q-Q plots for the analytical and simulated quantiles of EATs for all populations in the WAN-SPT-HK across the 100 epidemic scenarios used in Fig. 1. Insets show the corresponding histograms of percent error in expected EAT. (a) EATs for all populations in D_1 before (red) and after (blue) adjusting for the hub-effect. (b) EATs for all populations in D_2 before (red) and after (blue) adjusting for the continuous seeding and path reduction, in which the hub-effect has been adjusted for the epidemic origin and all populations in D_1 . (c) EATs for the remaining populations in D_3 and D_4 after adjusting for the hub-effect, continuous seeding and path reduction.

Furthermore, this pdf can in turn be well approximated with $f_n(t|\lambda_\psi, \alpha_\psi)$ where $\lambda_\psi, \alpha_\psi$ are obtained by minimizing the relative entropy [21, 25] for $n = 1$ (i.e. the first exportation). This indicates that the spread of epidemics from the origin to any population in D_2 can be regarded as a two-population model, in which the adjusted mobility rate is α_ψ and the epidemic in the origin grows exponentially at rate λ_ψ . We term this procedure *path reduction*.

Next, consider a longer path $\varphi : i \rightarrow j \rightarrow k \rightarrow m$, i.e. $m \in D_3$. Using path reduction, we first approximate the entire path φ with $\varphi' : i \rightarrow k \rightarrow m$ where the adjusted mobility rate and adjusted epidemic growth rate in the origin for the connection $i \rightarrow k$ are $\lambda_\psi, \alpha_\psi$, respectively. The arrival times of infections for population $m \in D_3$ (i.e. $T_{im}^n, n = 1, 2, \dots$) can be estimated using the methods that we have developed for D_2 populations. **Fig. 4** show that recursively using adjustments for the hub-effect and continuous seeding accurately characterizes the arrival times for all populations in the WAN-SPT.

3.3 The whole WAN

The accuracy of our WAN-SPT analysis provides a key insight: for each acyclic path ψ that connects any given population k to the epidemic origin, the epidemic arrival process for population k along this path is well approximated as an *NPP* with intensity function $\alpha_\psi \exp(\lambda_\psi t)$. In the whole WAN, each population might be connected to the epidemic origin via multiple paths, some of which might be intersected and therefore dependent (see Fig. 5(a)). We conjecture that the dependence among such paths is sufficiently weak, such that the overall epidemic arrival process for any population k in

the WAN can be characterized with the following method: (i) decomposing all paths that connects the epidemic origin to population k into a set Ψ_{ik} of independent acyclic paths; and then (ii) approximating the EAT for population k by the superposition of the *NPPs* [20] that correspond to these pseudo-independent paths. Mathematically, the epidemic arrival process for population k is well approximated by an *NPP* with intensity function $\sum_{\psi \in \Psi_{ik}} \alpha_\psi \exp(\lambda_\psi t)$. **Fig. 5** validates that our analytical framework (i.e. synthesis of the two-population analytics, adjustment for the hub-effect, adjustment for continuous seeding, path reduction and path superposition) is accurate for characterizing the EATs for almost all populations in the WAN. The results are robust for all tested 100 epidemic scenarios.

4 CONCLUSIONS

In summary, we have developed an analytical framework that grounds on the basic principles in infectious disease epidemiology and network theory for understanding the dynamics underlying global spread of emerging epidemics. Not only can our framework provides analytical and computational advancement for forecasting EATs for all populations in the WAN, but it also elucidates the dependence of EATs on the epidemiologic parameters (growth rate and seed size) and the network properties of the WAN (air traffic volume and connectivity). Because our framework provides closed-form probability distributions (Eq. (1)), it can also support likelihood-based inference of key epidemiologic parameters from surveillance data on local disease incidence and global case exportations [25]. Ongoing studies deserve to extend the framework to account for more complex factors including the

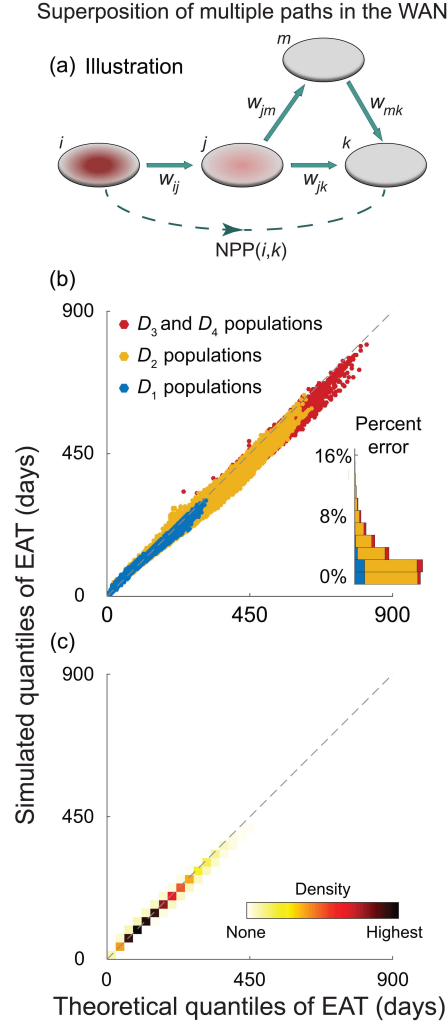


Figure 5: Validating the analytical framework for the WAN. The epidemic origin is Hong Kong as in Fig. 4. (a) Q-Q plots for the analytical and simulated quantiles of EATs for all populations in the WAN. Analytical EATs are computed using the NPP superposition as described in the section 3.3, while simulated EATs are generated from our global metapopulation simulator as described in the section 2.2. Data points are colored in blue for D_1 populations, yellow for D_2 populations, and red for D_3 and D_4 populations. (b) Density of the data points in (a) to show that nearly all the 230,800 Q-Q plots coincide with the diagonal, which demonstrates the congruence between analytical and simulated EATs.

stochasticity of intra-population transmission dynamics [29] and seasonal travel patterns [8, 27].

A APPENDICES

A.1 $SE_m I_n R$ model for epidemic spreading within each population

In the main text, we build the analytical framework using the SIR model within each population. Here we extend the theory to $SE_m I_n R$ models [26] in which:

- (1) The duration of latency is gamma distributed with mean D_E and m subclasses (i.e. with shape m and rate $b_E = m/D_E$).
- (2) The duration of infectiousness is gamma distributed with mean D_I and n subclasses (i.e. with shape n and rate $b_I = n/D_I$).

For any given population, let $S(t), R(t)$ be the number of susceptible and recovered individuals, respectively, $E_i(t)$ the number of individuals in the i th latent subclass, and $I_j(t)$ the number of individuals in the j th infectious subclass. The $SE_m I_n R$ model is described by the following differential equations:

$$\frac{dS(t)}{dt} = -\beta \frac{S(t)}{N} \sum_{j=1}^n I_j(t)$$

$$\frac{dE_1(t)}{dt} = \beta \frac{S(t)}{N} \sum_{j=1}^n I_j(t) - b_E E_1(t)$$

$$\frac{dE_i(t)}{dt} = b_E (E_{i-1}(t) - E_i(t)) \quad \text{for } i = 2, \dots, m$$

$$\frac{dI_1(t)}{dt} = b_E E_m(t) - b_I I_1(t)$$

$$\frac{dI_j(t)}{dt} = b_I (I_{j-1}(t) - I_j(t)) \quad \text{for } j = 2, \dots, n$$

$$\frac{dR(t)}{dt} = b_I I_j(t).$$

During the early stage of the epidemic (such that $S(t) \approx N$), the prevalence of latent and infectious people both grows exponentially at rate λ , which is the solution to the following equation [26]:

$$\lambda \left(\lambda + \frac{m}{D_E} \right)^m - \beta \left(\frac{m}{D_E} \right)^m \left(1 - \left(\frac{\lambda D_I}{n} + 1 \right)^{-n} \right) = 0$$

That is, the prevalence of latent and infectious individuals are well approximated by $\bar{E} \exp(\lambda t)$ and $\bar{I} \exp(\lambda t)$, respectively, where \bar{E} and \bar{I} depend on the initial conditions and parameters of the differential equation systems (the analytical expressions of \bar{E} and \bar{I} are obtained by solving the linearized system with $S(t) = N$). If a proportion $1 - p_E$ and $1 - p_I$ of the latent and infectious people refrain from air travel because of their infections, the seed size s_0 in the main text is simply $p_E \bar{E} + p_I \bar{I}$.

ACKNOWLEDGMENTS

We thank M. Lipsitch, J.M. Read, B.J. Cowling, P. Wu, K. Leung, H. Choi, N. Leung, S. Ali, J. Wong, V.J. Fang, Z. Wang, L. Chen, Y. Zhang and Y. Lin for helpful discussions. We thank C.K. Lam for assistance in data processing and technical support. We thank the Official Airline Guide and

Center for International Earth Science Information Network at Columbia University for the assembly of databases. This research was conducted in part using the research computing facilities and advisory services offered by Information Technology Services, The University of Hong Kong; and was done in part on the Olympus High Performance Compute Cluster at the Pittsburgh Supercomputing Center at Carnegie Mellon University, which is supported by National Institute of General Medical Sciences MIDAS Informatics Services Group under Grant No.: 1U24GM110707. This research was supported by Harvard Center for Communicable Disease Dynamics from the National Institute of General Medical Sciences MIDAS Initiative under Grant No.: U54GM088558, Research Grants Council Collaborative Research Fund under Grant No.: CityU8/CRF/12G, and two commissioned grants from the Health and Medical Research Fund from the Government of the Hong Kong SAR under Grant No.: HKS-15-E03, HKS-17-E13. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institute of General Medical Sciences and the National Institutes of Health.

REFERENCES

- [1] Roy M. Anderson and Robert M. May. 1991. *Infectious Diseases of Humans: Dynamics and Control* (1st ed.). Oxford Univ. Press, Oxford, UK.
- [2] Paolo Bajardi, Chiara Poletto, Jose J Ramasco, Michele Tizzoni, Vittoria Colizza, and Alessandro Vespignani. 2011. Human mobility networks, travel restrictions, and the global spread of 2009 H1N1 pandemic. *PLoS one* 6, 1 (Jan. 2011), e16591. <https://doi.org/10.1371/journal.pone.0016591>
- [3] Duygu Balcan, Vittoria Colizza, Bruno Gonçalves, Hao Hu, José J Ramasco, and Alessandro Vespignani. 2009. Multiscale mobility networks and the spatial spreading of infectious diseases. *Proceedings of the National Academy of Sciences* 106, 51 (Dec. 2009), 21484–21489. <https://doi.org/10.1073/pnas.0906910106>
- [4] Dirk Brockmann and Dirk Helbing. 2013. The hidden geometry of complex, network-driven contagion phenomena. *Science* 342, 6164 (Dec. 2013), 1337–1342. <https://doi.org/10.1126/science.1245200>
- [5] Dennis Carroll, Peter Daszak, Nathan D Wolfe, George F Gao, Carlos M Morel, Subhash Morzaria, Ariel Pablos-Méndez, Oyewale Tomori, and Jonna AK Mazet. 2018. The global virome project. *Science* 359, 6378 (Feb. 2018), 872–874. <https://doi.org/10.1126/science.aap7463>
- [6] Benjamin J Cowling, Dennis KM Ip, Vicky J Fang, Piyarat Sutarattiwong, Sonja J Olsen, Jens Levy, Timothy M Uyeki, Gabriel M Leung, JS Malik Peiris, Tawee Chotpitayasunondh, et al. 2013. Aerosol transmission is an important mode of influenza A virus spread. *Nature communications* 4 (June 2013), 1935. <https://doi.org/10.1038/ncomms2922>
- [7] Erin Doxsey-Whitfield, Kytta MacManus, Susana B Adamo, Linda Pistolesi, John Squires, Olena Borkovska, and Sandra R Baptista. 2015. Taking advantage of the improved availability of census data: a first look at the gridded population of the world, version 4. *Papers in Applied Geography* 1, 3 (July 2015), 226–234. <https://doi.org/10.1080/23754931.2015.1014272>
- [8] Anne Ewing, Elizabeth C Lee, Cécile Viboud, and Shweta Bansal. 2016. Contact, travel, and transmission: The impact of winter holidays on influenza dynamics in the United States. *The Journal of infectious diseases* 215, 5 (Dec. 2016), 732–739. <https://doi.org/10.1093/infdis/jiw642>
- [9] J Patrick Fitch. 2015. Engineering a global response to infectious diseases. *Proc. IEEE* 103, 2 (March 2015), 263–272. <https://doi.org/10.1109/JPROC.2015.2389146>
- [10] Aurélien Gautreau, Alain Barrat, and Marc Barthelemy. 2008. Global disease spread: statistics and estimation of arrival times. *Journal of theoretical biology* 251, 3 (April 2008), 509–522. <https://doi.org/10.1016/j.jtbi.2007.12.001>
- [11] Daniel T Gillespie. 1977. Exact stochastic simulation of coupled chemical reactions. *The journal of physical chemistry* 81, 25 (Dec. 1977), 2340–2361. <https://doi.org/10.1021/j100540a008>
- [12] Bryan Grenfell and John Harwood. 1997. (Meta) population dynamics of infectious diseases. *Trends in ecology & evolution* 12, 10 (July 1997), 395–399. [https://doi.org/10.1016/S0169-5347\(97\)01174-9](https://doi.org/10.1016/S0169-5347(97)01174-9)
- [13] Nathan D Grubaugh, Jason T Ladner, Moritz UG Kraemer, Gytis Dudas, Amanda L Tan, Karthik Gangavarapu, Michael R Wiley, Stephen White, Julien Théze, Diogo M Magnani, et al. 2017. Genomic epidemiology reveals multiple introductions of Zika virus into the United States. *Nature* 546, 7658 (jun 2017), 401–405. <https://doi.org/10.1038/nature22400>
- [14] Edward C Holmes, Andrew Rambaut, and Kristian G Andersen. 2018. Pandemics: spend on surveillance, not prediction. *Nature* 558 (June 2018), 180–182. <https://doi.org/10.1038/d41586-018-05373-w>
- [15] Flavio Iannelli, Andreas Koher, Dirk Brockmann, Philipp Hövel, and Igor M Sokolov. 2017. Effective distances for epidemics spreading on complex networks. *Physical Review E* 95, 1 (Jan. 2017), 012313. <https://doi.org/10.1103/PhysRevE.95.012313>
- [16] Matt J. Keeling and Pejman Rohani. 2007. *Modeling Infectious Diseases in Humans and Animals*. Princeton Univ. Press, Princeton, NJ.
- [17] Madhav Marathe and Anil Kumar S Vullikanti. 2013. Computational epidemiology. *Commun. ACM* 56, 7 (July 2013), 88–96. <https://doi.org/10.1145/2483852.2483871>
- [18] Joël Mossong, Niel Hens, Mark Jit, Philippe Beutels, Kari Auranen, Rafael Mikolajczyk, Marco Massari, Stefania Salmaso, Gianpaolo Scalia Tomba, Jacco Wallinga, et al. 2008. Social contacts and mixing patterns relevant to the spread of infectious diseases. *PLoS medicine* 5, 3 (March 2008), e74. <https://doi.org/10.1371/journal.pmed.0050074>
- [19] Romualdo Pastor-Satorras, Claudio Castellano, Piet Van Mieghem, and Alessandro Vespignani. 2015. Epidemic processes in complex networks. *Reviews of modern physics* 87, 3 (Aug. 2015), 925. <https://doi.org/10.1103/RevModPhys.87.925>
- [20] Sheldon M. Ross. 1996. *Stochastic Processes* (2nd. ed.). John Wiley & Sons, New York, NY.
- [21] Joy A. Thomas and Thomas M. Cover. 2006. *Elements of Information Theory* (2nd. ed.). John Wiley & Sons, Hoboken, NJ.
- [22] Michele Tizzoni, Paolo Bajardi, Chiara Poletto, José J Ramasco, Duygu Balcan, Bruno Gonçalves, Nicola Perra, Vittoria Colizza, and Alessandro Vespignani. 2012. Real-time numerical forecast of global epidemic spreading: case study of 2009 A/H1N1pdm. *BMC medicine* 10, 1 (Dec. 2012), 165. <https://doi.org/10.1186/1741-7015-10-165>
- [23] Gianpaolo Scalia Tomba and Jacco Wallinga. 2008. A simple explanation for the low impact of border control as a countermeasure to the spread of an infectious disease. *Mathematical biosciences* 214, 1-2 (July 2008), 70–72. <https://doi.org/10.1016/j.mbs.2008.02.009>
- [24] Lin Wang and Xiang Li. 2014. Spatial epidemiology of networked metapopulation: An overview. *Chinese Science Bulletin* 59, 28 (Oct. 2014), 3511–3522. <https://doi.org/10.1007/s11434-014-0499-8>
- [25] Lin Wang and Joseph T Wu. 2018. Characterizing the dynamics underlying global spread of epidemics. *Nature Communications* 9, 1 (Jan. 2018), 218. <https://doi.org/10.1038/s41467-017-02344-z>
- [26] Helen J Wearing, Pejman Rohani, and Matt J Keeling. 2005. Appropriate models for the management of infectious diseases. *PLoS Medicine* 2, 7 (Jul 2005), e174. <https://doi.org/10.1371/journal.pmed.0020174>
- [27] Amy Wesolowski, Elisabeth zu Erbach-Schoenberg, Andrew J Tatem, Christopher Lourenço, Cecile Viboud, Vivek Charu, Nathan Eagle, Kenth Engø-Monsen, Taimur Qureshi, Caroline O Buckee, et al. 2017. Multinational patterns of seasonal asymmetry in human movement influence infectious disease dynamics. *Nature communications* 8, 1 (Dec. 2017), 2069. <https://doi.org/10.1038/s41467-017-02064-4>
- [28] World Health Organization (WHO). 2018. 2018 Annual review of diseases prioritized under the Research and Development Blueprint. Retrieved July 8, 2018 from <http://www.who.int/blueprint/priority-diseases/en/>
- [29] Joseph T. Wu and Benjamin J. Cowling. 2011. The use of mathematical models to inform influenza pandemic preparedness and response. *Experimental Biology and Medicine* 236, 8 (Aug. 2011),

- 955–961. <https://doi.org/10.1258/ebm.2010.010271>
- [30] Joseph T Wu, Gabriel M Leung, Marc Lipsitch, Ben S Cooper, and Steven Riley. 2009. Hedging against antiviral resistance during the next influenza pandemic using small stockpiles of an alternative chemotherapy. *PLoS medicine* 6, 5 (May 2009), e1000085. <https://doi.org/10.1371/journal.pmed.1000085>
- [31] Xiaoxu Wu, Yongmei Lu, Sen Zhou, Lifan Chen, and Bing Xu. 2016. Impact of climate change on human infectious diseases: Empirical evidence and human adaptation. *Environment international* 86 (Jan. 2016), 14–23. <https://doi.org/10.1016/j.envint.2015.09.007>