

University at Albany, State University of New York

## Scholars Archive

---

Political Science Faculty Scholarship

Political Science

---

Fall 10-10-2016

### Testing Modeling Assumptions in the West Africa Ebola Outbreak

Matthew C. Ingram

*University at Albany, SUNY, mingram@albany.edu*

Keith Burghardt

*University of California at Davis*

Christopher Verzijl

*ABN AMRO Bank N.V., Amsterdam*

Junming Huang

*University of Electronic Science and Technology of China*

Binyang Song

*Singapore University of Technology and Design*

The University at Albany community has made this article openly available.

**Please share** how this access benefits you.

*See next page for additional authors*

Follow this and additional works at: [https://scholarsarchive.library.albany.edu/rockefeller\\_pos\\_scholar](https://scholarsarchive.library.albany.edu/rockefeller_pos_scholar)



Part of the [Epidemiology Commons](#), [International Public Health Commons](#), and the [Virus Diseases Commons](#)

---

#### Recommended Citation

Ingram, Matthew C.; Burghardt, Keith; Verzijl, Christopher; Huang, Junming; Song, Binyang; and Hasne, Marie-Pierre, "Testing Modeling Assumptions in the West Africa Ebola Outbreak" (2016). *Political Science Faculty Scholarship*. 2.

[https://scholarsarchive.library.albany.edu/rockefeller\\_pos\\_scholar/2](https://scholarsarchive.library.albany.edu/rockefeller_pos_scholar/2)



This work is licensed under a [Creative Commons Attribution 4.0 International License](#).

This Article is brought to you for free and open access by the Political Science at Scholars Archive. It has been accepted for inclusion in Political Science Faculty Scholarship by an authorized administrator of Scholars Archive. Please see [Terms of Use](#). For more information, please contact [scholarsarchive@albany.edu](mailto:scholarsarchive@albany.edu).

---

**Authors**

Matthew C. Ingram, Keith Burghardt, Christopher Verzijl, Junming Huang, Binyang Song, and Marie-Pierre Hasne

# SCIENTIFIC REPORTS



OPEN

## Testing Modeling Assumptions in the West Africa Ebola Outbreak

Keith Burghardt<sup>1</sup>, Christopher Verzijl<sup>2</sup>, Junming Huang<sup>3,4</sup>, Matthew Ingram<sup>5</sup>, Binyang Song<sup>6</sup> & Marie-Pierre Hasne<sup>7</sup>

Received: 27 June 2016

Accepted: 15 September 2016

Published: 10 October 2016

The Ebola virus in West Africa has infected almost 30,000 and killed over 11,000 people. Recent models of Ebola Virus Disease (EVD) have often made assumptions about how the disease spreads, such as uniform transmissibility and homogeneous mixing within a population. In this paper, we test whether these assumptions are necessarily correct, and offer simple solutions that may improve disease model accuracy. First, we use data and models of West African migration to show that EVD does not homogeneously mix, but spreads in a predictable manner. Next, we estimate the initial growth rate of EVD within country administrative divisions and find that it significantly decreases with population density. Finally, we test whether EVD strains have uniform transmissibility through a novel statistical test, and find that certain strains appear more often than expected by chance.

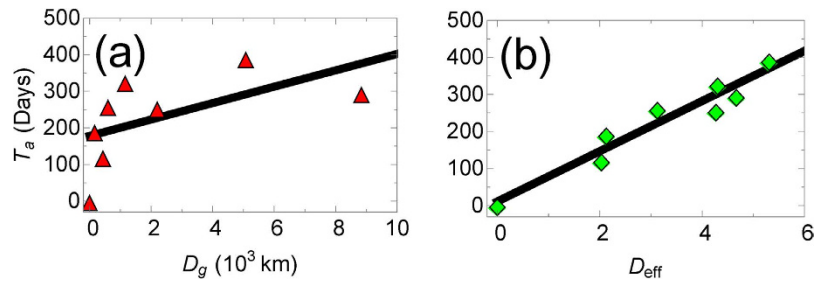
Since 2013, the Ebola virus outbreak in West Africa has become the largest such outbreak known. The epidemic first emerged in December 2013 in southern Guinea, but as of 11 May 2016, there have been about 28,600 cases of EVD in Guinea, Liberia, and Sierra Leone, and isolated cases in Italy, Mali, Nigeria, Senegal, Spain, the United Kingdom, and the United States. 11,300 of these cases were fatal<sup>1</sup> and, as high as these numbers are, they may be under-estimates due to the poor quality of current data<sup>2</sup>. The goal of this paper is to better understand the spread of EVD, and test the assumptions of leading EVD models.

Individuals have often been assumed to homogeneously mix with each other in many recent EVD models<sup>3–7</sup>, but we show that, by applying recent work on the migration of diseases<sup>8</sup>, homogeneous mixing is an especially poor approximation for EVD. We find that human migration patterns help predict where and when EVD originated and will appear, which would not be possible with a homogeneous mixing assumption. We also find evidence that the spread of EVD is much slower than other recent diseases, such as H1N1 and SARS<sup>8</sup>, which may have helped health workers control the disease.

Furthermore, against our expectations, we find that the initial growth rate of EVD can decrease significantly with population density, possibly because higher population density areas are correlated with other attributes, such as better healthcare. This compares to previous work where exponential and sub-exponential growth rates were found in many diseases, including the most recent EVD epidemic<sup>9,10</sup>, where variations in the growth rate of diseases were found, but mechanistic explanations were not explored. A previous model<sup>11</sup>, in comparison, found that higher population cities should contribute to a faster rate of disease spread, although we are not aware of previous research on disease spread and population density. Our work suggest that location-specific initial growth rates better model EVD, although the underlying reason for this heterogeneity should be a topic of future research.

Finally, we create novel metrics for the relative transmissibility of EVD strains, which are robust to sparse sampling. These metrics add to previous work on EVD in Sierra Leone<sup>12</sup>, and provide a novel understanding of EVD strains in Guinea. We find that the relative transmissibility of strains, as measured from these metrics, is not uniform; therefore, treating EVD as a single disease may be inappropriate<sup>3–7</sup>.

<sup>1</sup>Department of Political Science, University of California at Davis, Davis, 95616, USA. <sup>2</sup>ABN AMRO Bank N.V., Amsterdam, 1082 PP, Netherlands. <sup>3</sup>Complex Lab, Web Sciences Center, University of Electronic Science and Technology of China, 611731, P. R. China. <sup>4</sup>Center for Complex Network Research, Department of Physics, Northeastern University, Boston, 02115, USA. <sup>5</sup>Department of Political Science and Center for Social and Demographic Analysis, University at Albany, State University of New York, Albany, 12203, USA. <sup>6</sup>SUTD-MIT International Design Centre & Engineering Product Development Pillar, Singapore University of Technology and Design, Singapore, 487372, Singapore. <sup>7</sup>Department of Biochemistry and Molecular Biology, Oregon Health & Science University, Portland, 97239, USA. Correspondence and requests for materials should be addressed to K.B. (email: kaburghardt@ucdavis.edu)



**Figure 1.** The arrival time,  $T_a$ , of EVD in a country versus (a) the great-arc length distance,  $D_g$  and (b) the migration-based effective distance,  $D_{\text{eff}}$  from the disease's point of origin (Guinea). Error is smaller than marker size. The migration network used to construct the effective distance comes from census microdata<sup>13</sup>. The linear relationship between arrival time and effective distance suggests that there is a constant effective velocity of disease spread, in agreement with previous work on other diseases<sup>8</sup>.

These results, when taken together, suggest unexpectedly simple ways to improve EVD modeling. In the Discussion section, we will explain how a meta-population model can potentially aid in our understanding of disease spread and growth. Furthermore, incorporating disease strain dynamics into this model could help us better predict which strains will become dominant in the future, which may improve vaccination strategies.

## Results

Models of the West Africa Ebola outbreak have often assumed that the disease spreads via homogeneous mixing<sup>3–7</sup>. We find, however, that this assumption may not accurately model EVD when the disease first arrives in a given area. We will first discuss how the arrival time of EVD within a country or administrative area follows a predictable pattern due to the underlying migration model, in contrast to the mixing hypothesis. Next, we model the cumulative number of individuals infected in administrative divisions at the first or second level in Guinea, Liberia, and Sierra Leone to estimate the initial growth rate of EVD. We find this growth rate varies significantly, and appears to decrease with the population density within the administrative division. Finally, we introduce models of how EVD disease strains spread to rule out uniform strain transmissibility.

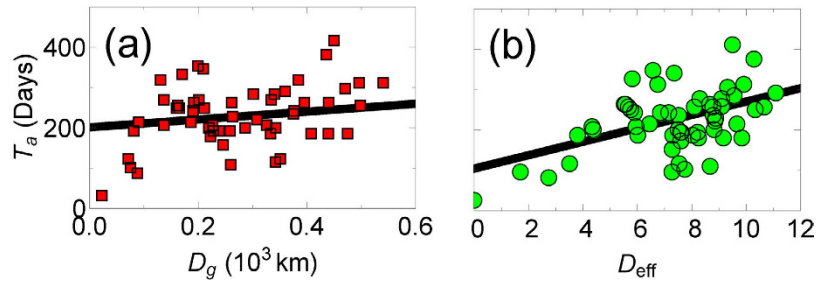
**How Does Ebola Spread?** Homogeneous mixing models assume that healthy individuals can get sick regardless of where they are, even when they are hundreds of miles from the origin of the infection. If this is true, then the disease should be quickly seen in all susceptible areas almost simultaneously. Although this approximation may be reasonable at short distances, there has to be a lengthscale when this would break down because, in the years since Ebola first emerged, no more than a handful of countries have become infected (Fig. 1).

Alternatively, one might assume that EVD spreads spatially. There is a significant positive correlation (Spearman  $\rho = 0.26$ ,  $p < 0.05$ ,  $n = 56$  at the spatial resolution of administrative divisions;  $0.81$ ,  $p < 0.05$ ,  $n = 8$  at the country resolution) between the arrival time of EVD in administrative divisions at the first or second level and the distance from the outbreak origin, Guéckédou, Guinea, to division centroids. Furthermore, this assumption has been applied successfully to model EVD in Liberia<sup>2</sup>.

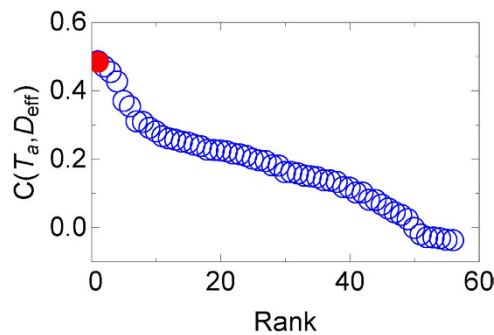
We find, however, that migration in West Africa is more complex than either of these assumptions<sup>13</sup>. Intuitively, diseases should spread more quickly between administrative divisions or countries with significant travel between them than between isolated areas. Therefore, it is reasonable to rescale distances such that areas with stronger travel ties are considered “closer” than areas without these ties, following work by Brockman and Helbing<sup>8</sup>. We find that rescaling distances using migration patterns helps us (1) better understand how quickly EVD spreads, and (2) estimate where the outbreak started.

To our surprise, we find that the correlation between the arrival time and effective distance for EVD (Spearman  $\rho = 0.95$ ,  $p < 10^{-3}$ ,  $n = 8$ ) was consistently higher than the correlation between the arrival time and geodesic distance (see Figs 1 and 2, and Supplementary Fig. S1). A high Pearson correlation ( $0.96$ ,  $p < 10^{-3}$ ), and agreement with the normality assumption ( $p > 0.05$ , using the Kolmogorov-Smirnov normality test), also suggests a linear relationship between the arrival time and effective distance. Taking the slope of the plot gives us an effective velocity of spread, which we find to be  $0.015 \text{ days}^{-1}$ , which is much slower than for previous diseases ( $\approx 0.1 \text{ days}^{-1}$ )<sup>8</sup>. An intuitive explanation for our finding is that lower overall migration in West Africa reduces the speed at which EVD spreads compared to other diseases.

The lower correlations at the first or second administrative level are due in part to the fact that we use migration models to determine the effective distance, and the disease spread for several months before it was detected<sup>7</sup>. Despite the lower quality data, however, we can still use it to determine the most likely origin of EVD, which we compare to the known origin, Guéckédou, Guinea. Quickly finding where a disease originated is important to help understand what caused it (e.g., what was the vector), and to predict where and when it will arrive, which can allow health workers to prepare<sup>8</sup>. Previously, it was found that the correlation between the arrival time and the effective distances from the disease origin is higher than correlations from areas where the disease did not originate<sup>8</sup>. We therefore ranked the correlation between the arrival times and effective distances from all the administrative divisions where Ebola was found between 2013 and 2016 (Fig. 3 and Supplementary Fig. S2), and found the true origin has the highest correlation. To determine the statistical significance of the ranking, we bootstrapped residuals of the linear regression between  $D_{\text{eff}}$  and  $T_a$ , with Guéckédou as the origin, to accurately



**Figure 2.** The arrival time,  $T_a$ , of patients with EVD in administrative divisions at the first or second level within Guinea, Mali, Liberia, Sierra Leone, or Nigeria, versus (a) the great-arc length distance,  $D_g$ , and (b) the migration-based effective distance,  $D_{\text{eff}}$ , from the disease's point of origin (Guéckédou, Guinea). Error is smaller than marker size. The migration network used to construct the effective distance comes from a radiation migration model<sup>29</sup> (similar results were found using the gravity migration models in ref. 13, see Supplementary Fig. S1).



**Figure 3.** The (ranked) correlations between arrival time,  $T_a$ , and migration-based effective distance,  $D_{\text{eff}}$ , between administrative divisions, calculated via Eq. 5 in ref. 8. The migration network used to construct the effective distance comes from a radiation migration model<sup>29</sup> (similar results were found using the gravity migration models in ref. 13, see Supplementary Fig. S2). In red is the true origin, which is found to have the highest correlation.

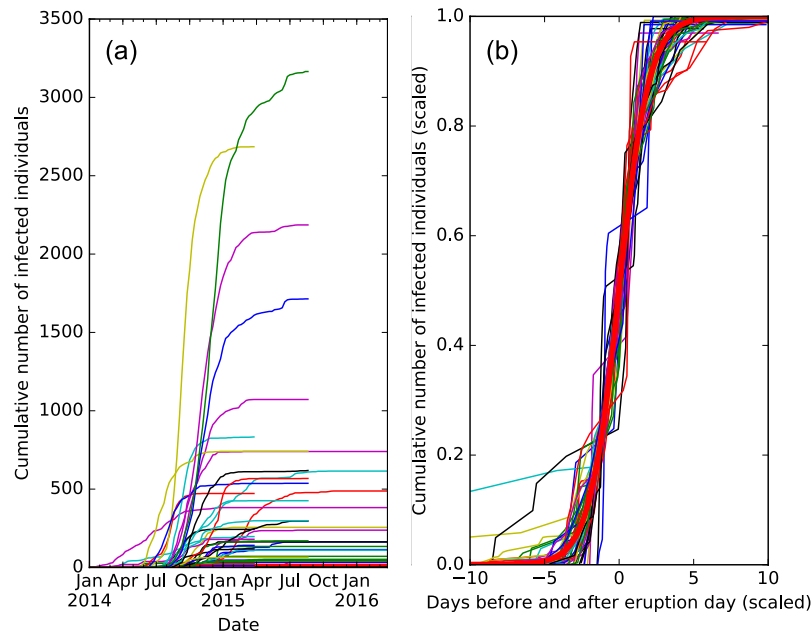
bootstrap arrival times. This is equivalent to simulating what would happen if we wound back the clock and restarted the infection from Guéckédou. The true origin has the highest correlation  $45\% \pm 1.6\%$  of the time for the radiation model, and 44–70% of the time for gravity models (Supplementary Fig. S2), therefore the true origin does not always have the highest correlation, but it often does, and can be a good first guess when determining the origin.

In conclusion, we find strong evidence that a migration network can elucidate how quickly Ebola spreads, when and where it will arrive, and where the infection began even with limited data. Furthermore, we see that alternative hypotheses for how EVD spreads, such as homogeneous mixing and nearest neighbor interactions, provide quantitatively poorer agreement with data.

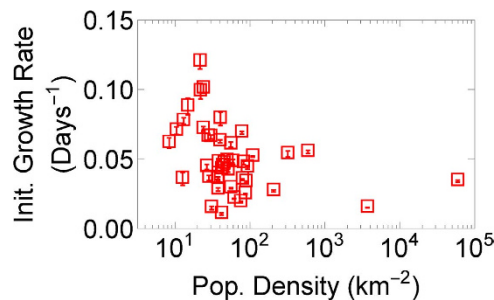
**The Growth of EVD Across Administrative Divisions.** In this section, we show that the cumulative number of Ebola cases within administrative divisions at the first or second level is well-approximated by a logistic function (Fig. 4a). We use this finding to estimate the initial growth rate of EVD for each infected administrative division, where, to our surprise, we find that the initial growth rate decreases with population density. We emphasize that variations in EVD growth rates have been seen before<sup>9,10</sup>, but mechanisms that might drive this behavior were not proposed. The logistic function has been used to model initial stages of EVD<sup>14</sup>, and is equivalent to the Susceptible-Infectious (SI) model when 100% of individuals are initially susceptible<sup>15</sup>. To fit the SI model to our data, however, we would have to make a simplistic assumption that only a fraction  $p_n$  of individuals are susceptible, in order to explain why a small fraction of the population is ever infected. We do not claim this describes the actual dynamics of EVD, although we will explain later why the cumulative number of cases should approximately follow this distribution. For an administrative division  $n$ , the cumulative number of infected individuals over time,  $t$ , by the logistic growth function is simply:

$$i_n(t) = \frac{p_n}{1 + e^{-q_n(t-t_0)}}. \quad (1)$$

The initial growth rate is  $q_n$ , while  $t_0$  is the time when the cumulative infection increases the most. The dynamics are highly heterogeneous (Fig. 4a), therefore, it would seem un-intuitive for a single function to fit all the data.



**Figure 4.** The cumulative number of infected individuals over time within administration divisions at the first or second level in Guinea, Sierra Leone and Liberia, from the patient database dataset<sup>1</sup>. It is clear that the rate and size of infections are heterogeneous, but when we fit the data to logistic functions, and renormalize the coefficients to  $p_n \rightarrow 1$ ,  $\tilde{t} = (t - t_0)q_n$  in (b), we see that the distributions collapse. A logistic function is plotted in red.

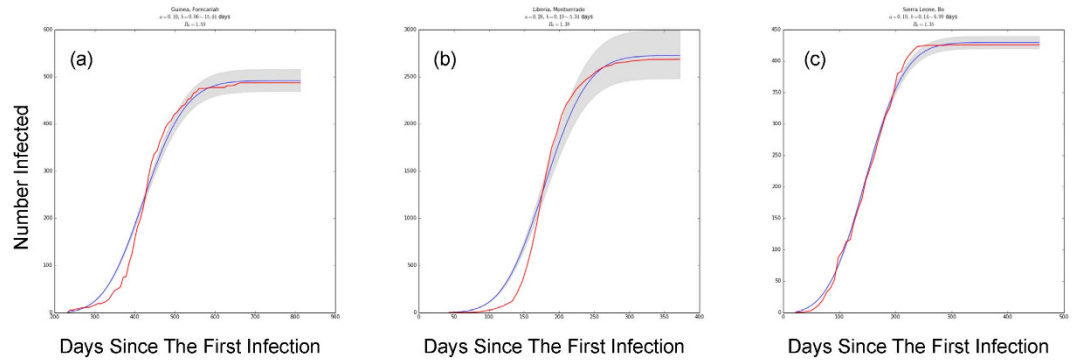


**Figure 5.** The initial growth rate,  $q_n$ , versus population density for the PSD dataset with more than 20 infections (not shown is the outlier Kissidougou, Guinea, with a population density of 34 individuals per kilometer, and a growth rate of 0.41 days<sup>-1</sup>), where error bars are standard deviations. We find that the initial growth rate drops significantly with population density (Spearman  $\rho = -0.48$ ,  $p < 10^{-3}$ ,  $n = 44$ ).

However, after fitting each cumulative distribution to the logistic model, and rescaling the variables:  $p_n \rightarrow 1$ ,  $\tilde{t} = (t - t_0)q_n$ , we find that the distributions collapse (Fig. 4b) to a good approximation.

Because small variations in  $q_n$  very quickly become substantial variations in the infection size later on, we want to understand how  $q_n$  varies across administrative divisions. When plotting  $q_n$  versus population density (Fig. 5) with more than 20 infections total, we notice that, although  $q_n$  is  $\approx 0.1$  days<sup>-1</sup>, which is similar to a previous EVD outbreak<sup>16</sup>,  $q_n$  varies significantly across administrative divisions. This result contrasts with many previous models that assume a global parameter can describe the growth of Ebola<sup>3-7</sup>. Furthermore, we find that  $q_n$  decreases significantly with population density (Spearman  $\rho = -0.48$ ,  $p < 10^{-3}$ ,  $n = 44$ ), plausibly because better healthcare may exist in higher density areas. This contrasts with a previous model, which predicted a positive scaling relation between the mean growth rate and population of cities<sup>11</sup>, although we are not aware of research that explored how disease spread is affected by population density. Therefore, not only is the growth rate of Ebola unexpectedly heterogeneous, but the dependence on population density may help us understand why this is the case.

*What Makes the Data Collapse?* In the SI model, 100% of individuals are eventually infected, therefore, to find agreement with data, our model has to assume a small proportion of individuals in each division are susceptible to the disease. This seems implausible; more likely, all individuals are susceptible and, as they become aware of an infection, they reduce their interactions or otherwise reduce the overall disease transmissibility. To demonstrate



**Figure 6.** SDIR fits to empirical data (fits for all infected administrative divisions can be seen in **Supplementary Figs S3–S5**). The cumulative number of infected individuals in (a) Forecariah, Guinea; (b) Montserrado, Liberia; and (c) Bo, Sierra Leone, where the blue line is the best fit and the shaded area represents standard errors. While some fits are poorer than others, we generally capture the qualitative structure of the empirical data, including the cumulative number of infected. This model may therefore begin to explain the mechanism behind the s-curve.

this hypothesis, we created a more realistic, although still simplistic, disease model, in which susceptible ( $s$ ) individuals can become infected ( $i$ ), but then recover or are removed<sup>17</sup>. In our data, all individuals who recovered or died were assumed to be “removed”. The SIR model, like the SI model, significantly overshoots the cumulative number of cases in the absence of intervention. We counter-balance this effect with an exponentially decreasing disease transmissibility as a result of public health interventions<sup>18,19</sup>. We call this model the Susceptible - Decreasingly Infectious - Recovered (SDIR) model.

The equations are:

$$\frac{ds_n}{dt} = -a(t) s_n i_n \quad (2)$$

$$\frac{di_n}{dt} = a(t) s_n i_n - b i_n \quad (3)$$

$$r_n = 1 - s_n - i_n \quad (4)$$

in normalized units, where  $s_n$  is the susceptible fraction of the population,  $i_n$  is the infected population and  $r_n$  is the removed population (either by recovery or death), for each administrative division.  $b$  is the recovery rate, while the infectious rate,  $a(t)$ , is defined as

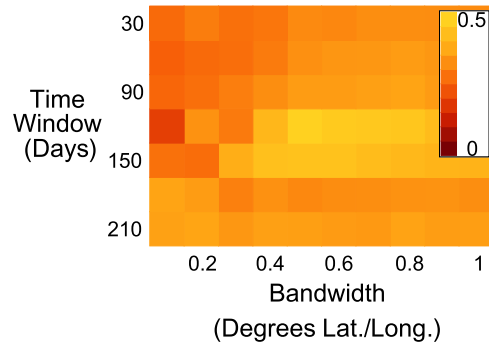
$$a(t) = a\{1 - [e^{-k(t-t')} - 1]\theta(t - t')\}, \quad (5)$$

In Eq. 5,  $a$  represents the overall infection rate, while  $k$  represents the rate at which control measures reduce transmission,  $t'$  is a delay before which individuals are not aware enough of the disease to try to limit its spread, and  $\theta$  is the Heaviside step function. Before  $t'$ , the disease infects at a rate  $a$ , while afterwards, we see a sharp drop-off in the disease transmissibility. It is clear that other, more realistic assumptions could be included into the model, including an exposed state (e.g., the SEIR model<sup>20</sup>, or the effect of burial practices, hospitalization, and other factors<sup>3–6</sup>, but we leave this out of our model for simplicity. We see rough quantitative agreement with empirical data (e.g. Fig. 6), which suggests that this simple model captures the essence of the real-world dynamics.

In review, our model assumes only 3 states and a global time-varying infection rate,  $a(t)$ . Even though there are reasons to consider this model too simplistic, it is able to effectively describe the roughly s-like curve of infections, and therefore it begins to help us understand the mechanism for a logistic-like cumulative distribution (Fig. 4), such as reductions in infectability due to disease awareness.

**Are Strains Uniformly Transmissible?** In this section, we use EVD genome sequences to determine what strains appear more often than expected by chance in Guinea and Sierra Leone. Our results suggest that EVD strains do not necessarily have uniform transmissibility. We are not aware of any previous models that take the strain of EVD into account, although a previous paper found certain EVD strains in Sierra Leone have different growth rates<sup>12</sup>. Unlike the previous paper, however, we can use our method to understand the transmissibility of strains in Guinea, where the sampling rate is otherwise too low.

**Modeling Strain-Dependent Infection Probabilities.** We use meta data from Ebola nucleotide sequences isolated from patients in Guinea<sup>21</sup> between April, 2014 and January, 2015, and Sierra Leone between late May 2014 and January, 2015<sup>22–24</sup>, to determine when and where a strain of EVD was found, then use kernel-density estimation (KDE, see Methods) to estimate the spatial probability density function (PDF) of being infected with an EVD strain<sup>25</sup>:



**Figure 7.** A heat map of the Spearman rank correlation between Eq. 8 and the number of individuals infected with EVD within a time  $\Delta t$ , as a function of the bandwidth  $h$  and the time window width,  $\Delta t$  in Eq. 8. The correlation between the model and data is highest (Spearman  $\rho = 0.50$ ,  $p < 0.01$ ,  $n = 68$ ) when the bandwidth is 0.5 degrees and the time window width is 120 days.

$$P_E(\bar{x}, t, h, \Delta t | s) = \frac{C}{|S_s| h^2} \sum_{i \in S_s} K\left(\frac{\bar{x} - \bar{x}_i}{h}\right) H(\Delta t - |t - t_i|), \quad (6)$$

Here,  $K$  represents a kernel with bandwidth  $h$  to represent the area around an observed sequence where individuals are likely to be infected,  $C$  is an overall constant and  $H$  is the Heaviside step function.  $S_s$  are all labels for the set of pairs  $\{\bar{x}_i, t_i\}$  with known infections of strain  $s$ . We add a Heaviside step function in the above equation to represent a sliding time window of length  $\Delta t$  – only within timeframe  $(t - \Delta t)$  to  $t$  is sequence data relevant. For the rest of the paper, our kernel is chosen to be a radially symmetric Gaussian:

$$K\left(\frac{\bar{x} - \bar{x}_i}{h}\right) = K\left(\frac{\|\bar{x} - \bar{x}_i\|}{h}\right) = \frac{1}{2\pi} \text{Exp}\left(-\frac{\|\bar{x} - \bar{x}_i\|^2}{2h^2}\right). \quad (7)$$

We do not believe that our results are qualitatively sensitive to the kernel choice. By summing these probabilities, we can then find the probability of being infected with EVD:

$$P_E(\bar{x}, t, h, \Delta t) = \frac{C}{nh^2} \sum_s \sum_{i \in S_s} K\left(\frac{\bar{x} - \bar{x}_i}{h}\right) H(\Delta t - |t - t_i|), \quad (8)$$

where  $n = \sum_s |S_s|$ . There are two free parameters: the kernel bandwidth  $h$  and sliding time window width  $\Delta t$ , but knowing the true infection pattern allows for these parameters to be estimated (see Fig. 7 for EVD in Guinea).

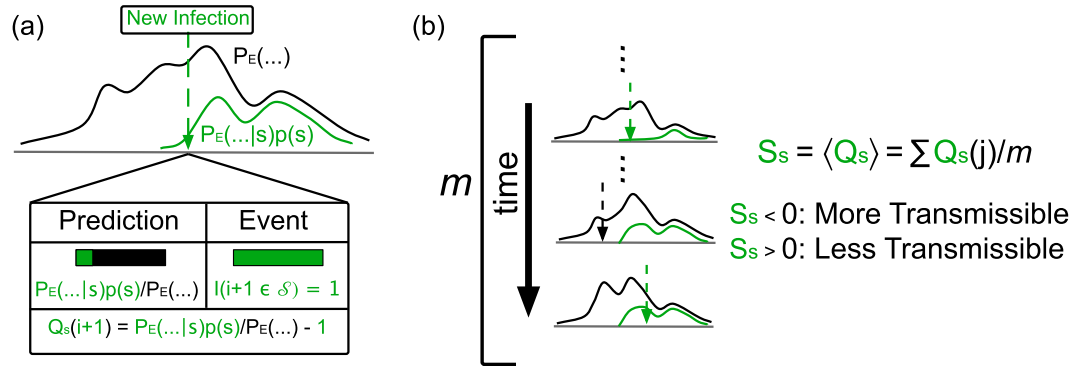
We apply KDE to EVD in Guinea and Sierra Leone, and test whether the estimated probabilities of becoming infected correlate with the number of infected individuals within each time window. A high correlation would represent close agreement between the KDE and actual spatial probability, and would increase our trust in these findings. In Guinea, we find Spearman correlations of up to 0.5 ( $p < 0.01$ ,  $n = 68$ , see Fig. 7) which suggests that the KDE is in good agreement with data. In Sierra Leone, however, these correlations are negligible. The Guinea dataset will therefore be the focus of the rest of our paper.

From Fig. 7, we find that the model best correlates with Guinea's infection data if  $h = 0.5$  and  $\Delta t = 120$  days. A plausible biological reason for the time window to be 120 days is that it may correspond to the effective viral shedding time, which can occur over timescale of roughly 100 days<sup>26</sup>. For example, the virus has been found to spread via sexual contact, and has been seen in vaginal fluid for up to 33 days, and in semen for up to 82 days via cell culture, and 550 days via RNA. These values are extremes and might not be likely to be encountered often, however there has been at least one other report of RNA seen in semen 100 days afterwards<sup>27</sup>. The length scale of 0.5 degrees, however, is roughly the distance between prefectures, which probably explains why the peak is at this bandwidth. We also find that our results are robust to choices of  $h$  and  $\Delta t$ , and we therefore let  $h$  vary between 0.1 and 1 degrees and  $\Delta t$  vary between 30 and 300 days.

**Measuring the Relative Transmissibility of Strains.** We have just demonstrated the applicability of KDE to measuring the relative likelihood of being infected with EVD across districts. This is not in and of itself too interesting – indeed other models may create better quantitative predictions of how EVD spreads<sup>2</sup>. Unlike previous work, however, this model can predict the relative transmissibility of individual EVD strains (Eq. 6). We will apply these predictions towards determining the relative transmissibility of EVD strains. When a strain  $s$  appears much more than it should compared to other strains, then we can conclude that  $s$  may be more transmissible.

Let a new infection arrive at position  $\bar{y}$  and at time  $t_y$ . Our model may predict that the probability this infection would have a strain  $s$  is, e.g., 20% and any other strain is 40%. If we find the disease is indeed strain  $s$ , we may be a little “surprised”. If  $s$  consistently appears more often than expectations, we may believe that  $s$  is simply more successful – it can reach individuals more quickly than other strains. To quantify our small “surprise” of seeing a new  $s$  strain, we can write the equation:





**Figure 8.** A schematic for our success metric. (a) When a new infection appears, the probability the infection would be a strain  $s$  is  $P_E(\bar{x}, t, h, \Delta t|s)$ . We define  $Q_s$  to be the difference between the prediction and actual event. (b)  $S_s$ , the expectation value of  $Q_s$ , should be 0 if  $s$  is no more transmissible than the disease as a whole. If  $S_s < 0$ , then we would say that  $s$  is more transmissible than the typical strain, while the opposite is true if  $S_s > 0$ .

$$Q_s(i+1) = \frac{P_E(\bar{x}_{i+1}, t_{i+1}, h, \Delta t|s)p(s)}{P_E(\bar{x}_{i+1}, t_{i+1}, h, \Delta t)} - I(i+1 \in \mathcal{S}_s), \tag{9}$$

where  $I$  is the indicator function that a new infection's strain is  $s$  or not, and  $p(s)$  is the probability that the infection strain is  $s$  (See Fig. 8). On average, the difference between our prediction and data,  $S_s$ , defined as

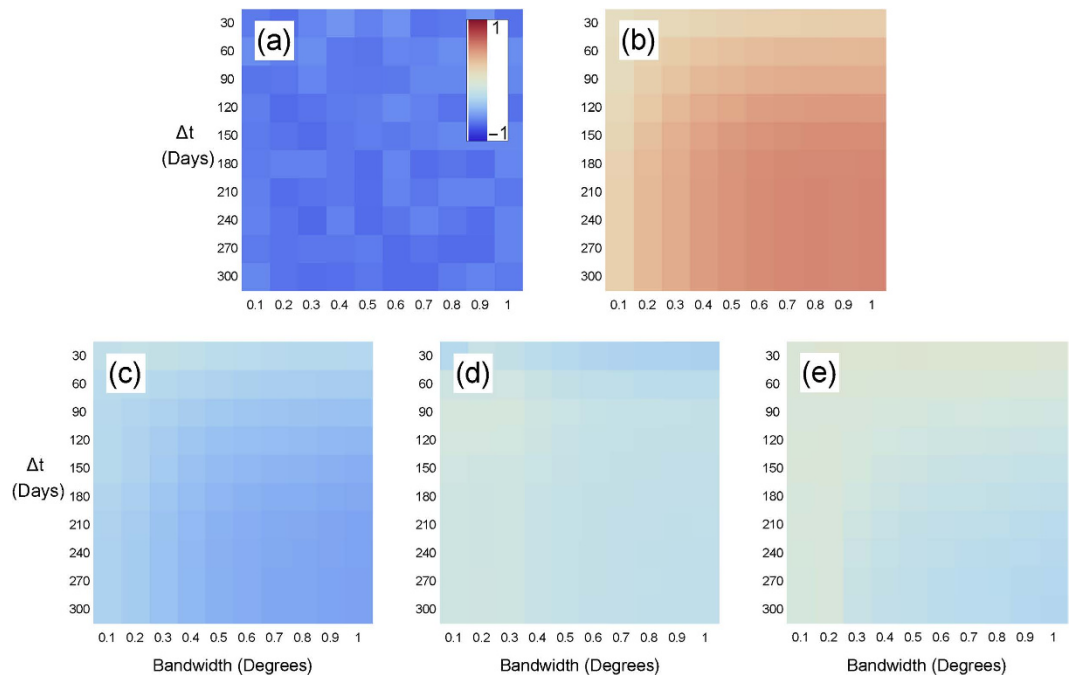
$$S_s \equiv \langle Q_s \rangle = \frac{\sum_j |t_{i_1} \leq t_j \leq t_{i_2}, i_1, i_2 \in \mathcal{S}_s Q_s(j)}{m}, \tag{10}$$

should be 0, where the sum is over all sequences between the first and last appearance of strain  $s$ , and  $m$  is the number of sequences that satisfy this constraint. When  $S_s = 0$ ,  $s$  is no stronger or weaker than the infection as a whole. If  $S_s < 0$ , then  $s$  is stronger than strains of EVD generally, while  $S_s > 0$  implies the opposite (see Fig. 8). To determine the statistical significance of  $S_s$ , we bootstrap values for which  $S_s = 0$  (see Methods). We find that  $S_s$  is significantly different than chance, (Fig. 9), therefore, Ebola strains might not be uniformly transmissible. For example, SL1 has one of the most negative values of  $S_s$  (Fig. 9), where all values are statistically significant ( $p < 0.05$ ,  $n$  varies with  $\Delta t$ ). Interestingly, we also find the values of SL1 appear to fluctuate, possibly because we only have 9 SL1 datapoints, versus 37, 19, 13, and 36 datapoints for GN1, GN2, GN3, and GN4, respectively. Most other strains have a success metric indicating a stronger strain than expected by chance probably because few strains are concurrent in time, and therefore, comparing  $S_s$  across strains is difficult. Our method may be compared to the method used by Meyer, Elias, and Hohle<sup>28</sup>, who compared the success of invasive meningococcal disease (IMD) strains. Common surveillance algorithms were found to be practical but these appear to be more useful for a chronic rather than an acute outbreak, such as EVD.

### Discussion

In conclusion, we find several factors not often accounted for that may improve the accuracy of modeling EVD. First, EVD appears to spread with a constant effective velocity through a migration network, in disagreement with the homogeneous mixing hypothesis (Figs 1 and 2, and Supplementary Fig. S1) commonly used to model diseases like Ebola. By taking into account the role a migration network has in disease spread, we can predict where EVD will arrive with greater accuracy than before. This method can also accelerate the process of identifying the index case by determining which administrative division is the origin through arrival time-effective distance correlation maximization. Second, we find that the growth of EVD at the finer spatial resolutions can be well described by three scaling parameters, and the initial growth rate decreases with the population density, contrary to our intuition, which suggests that a population density-dependent EVD model may more accurately predict the spread of Ebola when the disease first arrives at an administrative division. One plausible explanation for this result is that higher population density areas receive better healthcare than other areas, but more work is necessary to understand this behavior. Finally, we find a wide variation in the transmissibility of different strains of EVD, which suggests that modeling each disease strain, when this information is known, can improve the prediction task by reducing the heterogeneity of the data. In addition, our method may improve vaccination strategies if vaccines are made for particularly transmissible strains as well as the most common ones.

One way to take these factors into account would be through a meta-population model on top of a migration network with high spatial resolution. A meta-population model treats areas such as cities, districts, or countries as nodes on a network, with links connecting them to represent the flow of individuals from one area to another. Diseases within each area (node) are modeled with a compartmental model, under the assumption that the homogeneous mixing approximation is more accurate in smaller areas than for the entire network. Previous work has already found that a meta-population model<sup>8</sup>, or spatio-temporal model<sup>2</sup>, can accurately predict the spread of diseases. Finally, strain-dependent transmissibility may further improve model accuracy.



**Figure 9.** The success metric for strains in Guinea for (a) SL1, (b) GN1, (c) GN2, (d) GN3, (e) GN4, across various time window widths and bandwidths, where red values correspond to strains seen less often, than expected between the time the strain first appeared and last was seen. Sequence data comes from ref. 21.

Our approach towards testing assumptions in disease modeling is not restricted to EVD, but can be applied toward other diseases of epidemiological concern. Future work is therefore necessary to test our methods on other diseases and check whether a meta-population model will better predict disease spread.

## Methods

This section explains how data on the cumulative number of infected individuals at the first and second administrative level (e.g., counties or districts) was gathered, how the migration network was constructed, and how we found and used strain data in Fig. 9.

**Infection Data from Humanitarian Data Exchange and World Health Organization.** For the 2014–2016 EVD epidemic, we are aware of two main data sources on the cumulative number of infections: World Health Organization (WHO) patient database and the WHO weekly situation reports<sup>1</sup>. The patient database data produces results similar to the Situation Reports, but for consistency, all plots in this paper use the patient database. We focused on data from December 2013 to January 2016 for the major West African countries affected: Guinea, Liberia, Mali, Nigeria, Senegal, and Sierra Leone.

There was a significant amount of work parsing data from the patient database. Administration names in the data had multiple spellings, some administrative areas appeared individually but also aggregated with nearby areas, and finally spaces, accents, and other characters appeared inconsistently. These were cleaned and harmonized. Although in the most affected countries, we have the cumulative number of cases at the second administrative level, in three countries (Liberia, Mali, and Nigeria) we only have data at the first administrative level (region, not district).

**Migration Data from Flowminder.** Data on intra- and international migration in West Africa is taken from Flowminder<sup>13</sup>. For modeling migration at fine spatial resolutions, the Flowminder data contain three different data sets: (1) within countries (capturing exclusively intranational movement of people); (2) between countries (capturing exclusively international movement of people); and (3) within and between countries, capturing both intra- and international movement of people. All three data sets include the West African countries most affected by EVD as well as Benin, Cote d'Ivoire, Gambia, Ghana, Guinea Bissau, Senegal, and Togo, which had few, if any, EVD cases. In the third dataset, Flowminder collected census microdata from Public Use Microdata Series (IPUMS) on what country (or countries) an individual resided during the previous year.

To create a migration network, we first matched the Flowminder node coordinates to known district centroids (errors between centroids and node coordinates were  $\pm 10$  km). Having matched nodes to administrative names, we associated each node to the district-level arrival time of EVD, recording the date that the WHO patient database first records for a case in each administrative area. Four gravity model parameters were used to estimate the traffic between administration areas. Three were fit to migration using cell phone data in Cote d'Ivoire, Kenya, and Senegal, respectively, while the final one was fit to IPUMS data. We found that all produce similar fits, and work equally well to estimate the effective distances between areas. In addition, using population data from

Geohive (next section), we created a radiation migration model, found to more accurately estimate migration patterns than gravity models<sup>29</sup>, we therefore used this for all figures that use migration networks in this paper. Figures using gravity models can be seen in Supplementary Figs S1 and S2.

**Population data from GeoHive.** To determine the population density, and to estimate the migration network from the radiation migration model<sup>29</sup>, we first find the population and area of each administrative division. First, we collected the area administrative divisions in West Africa from [www.geohive.com](http://www.geohive.com). Next, we collected population in those districts from population census records. For districts providing multiple census datasets, we collected from the latest report: Guinea, 2014; Liberia, 2008; Mali, 2009; Nigeria 2011; Senegal 2013; and Sierra Leone, 2004. Some population data is probably out-of-date and may lead to some biases in the population density and radiation model. However, we believe that newer data will confirm our initial conclusions.

**Fitting Disease Models.** To find the best-fit parameters and associated errors in the logistic model (Eq. 1), and the SDIR model (Eqs 2–4), we used least squares fitting. To reduce the possibility of overfitting data, we focus on districts where more than 20 individuals become infected.

**Temporal-Spatial Resolution of Sequences.** To find empirical values of Eqs 6 and 8, we gathered infection meta-data (the strain, time, location of infected individual) from recent papers<sup>21–24</sup>. In the future we hope to make our own phylogenetic tree from the sequences, but for now we use the strain labels from the supplementary data itself.

A substantial fraction of sequences do not belong in any significant clade, and some sequences could belong to multiple clades, depending on the phylogenetic tree method used<sup>22–24</sup>, therefore, we found  $Q_s$  and  $S_s$  for each strain by comparing that strain to all other sequences within the time window.

To determine the statistical significance of values in Fig. 9, we generate Bernoulli random variables (1 with probability  $Pr(s, \bar{x}_i, t_i)$ , and 0 with probability  $1 - Pr(s, \bar{x}_i, t_i)$ ) with

$$Pr(s, \bar{x}_i, t_i) = \frac{P_E(\bar{x}_i, t_i, h, \Delta t|s)p(s)}{P_{EVD}(\bar{x}_i, t_i, h, \Delta t)} \quad (11)$$

where  $\bar{x}_i$  is the infection location at a time  $t_i$ , to represent idealized data in which  $\langle S_s \rangle = 0$ . We determine  $S_{s,bootstrap}$  (Eq. 10) from these idealized values of  $Q_s$ . If the absolute value of the empirical  $S_s$  value,  $|S_{s,empirical}|$ , is greater than 95% of  $|S_{s,bootstrap}|$  values, then the empirical data is statistically significant.

## References

- World Health Organization, Ebola data and statistics. <http://apps.who.int/gho/data/node.ebola-sitrepe>, (Date of access: 25/08/2016) (2016).
- Merler, S. *et al.* Spatiotemporal spread of the 2014 outbreak of ebola virus disease in liberia and the effectiveness of non-pharmaceutical interventions: a computational modeling analysis. *Lancet Infect. Dis.* **15**, 204–211 (2015).
- Lewnard, J. A. *et al.* Dynamics and control of ebola virus transmission in montserrat, liberia: a mathematical modeling analysis. *The Lancet* **14**, 1189–1195 (2014).
- Chowell, D., Castillo-Chavez, C., Krishna, S., Qiu, X. & Anderson, K. S. Modelling the effect of early detection of ebola. *The Lancet: Infectious Diseases* **15** (2015).
- Camacho, A., Kucharski, A. J., Funk, S., Piot, P. & Edmunds, W. J. Potential for large outbreaks of ebola disease. *Epidemics* **9**, 70–78 (2014).
- Pandey, A. *et al.* Strategies for containing ebola in west africa. *Science* **346**, 991–995 (2014).
- Chowell, G. & Nishiura, H. Transmission dynamics and control of ebola virus disease (evd): a review. *BMC Medicine* **12** (2014).
- Brockmann, D. & Helbing, D. The hidden geometry of complex, network-driven contagion phenomena. *Science* **342**, 1337–1342 (2013).
- Viboud, C., Simonsen, L. & Chowell, G. A generalized-growth model to characterize the early ascending phase of infectious disease outbreaks. *Epidemics* **15**, 27–37 (2016).
- Chowell, G., Viboud, C., Hyman, J. M. & Simonsen, L. The western africa ebola virus disease epidemic exhibits both global exponential and local polynomial growth rates. *PLOS Currents Outbreaks*, doi: 10.1371/currents.outbreaks.8b55f4bad99ac5c5db3663e916803261 (2015).
- Schlapfer, M. *et al.* The scaling of human interactions with city size. *J. R. Soc. Interface* **11**, 20130789; doi: 10.1098/rsif.2013.0789 (2014).
- Luksha, M., Bedford, T. & Lassig, M. Epidemiological and evolutionary analysis of the 2014 ebola virus outbreak. *arXiv:1411.1722 [q-bio.PE]* (2014).
- Wesolowski, A. *et al.* Commentary: Containing the ebola outbreak - the potential and challenge of mobile network data. *PLoS Currents Outbreaks* (2014).
- Chowell, G., Simonsen, L., Viboud, C. & Kuang, Y. Is West Africa Approaching a Catastrophic Phase or is the 2014 Ebola Epidemic Slowing Down? Different Models Yield Different Answers for Liberia. *PLOS Currents Outbreaks*, doi: 10.1371/currents.outbreaks.b4690859d91684da963dc40e00f3da81 (2014)
- Bailey, N. T. J. *The mathematical theory of infectious diseases* (Macmillan, 1975).
- Chowella, G., Hengartner, N., Castillo-Chavez, C., Fenimore, P. & Hyman, J. The basic reproductive number of ebola and the effects of public health measures: the cases of congo and uganda. *J. Theor. Biol.* **7**, 119–126 (2004).
- Kermack, W. O. & McKendrick, A. G. A contribution to the mathematical theory of epidemics. *Proc. R. Soc. London* **115**, 700–721 (1927).
- Lekone, P. E. & Finkenstädt, B. F. Statistical inference in a stochastic epidemic SEIR model with control intervention: Ebola as a case study. *Biometrics* **62**(4), 1170–1177 (2006).
- Althaus, C. L. Estimating the Reproduction Number of Ebola Virus (EBOV) During the 2014 Outbreak in West Africa. *PLOS Currents Outbreaks*, doi: 10.1371/currents.outbreaks.91afb5e0f279e7f29e7056095255b288 (2014).
- Herbert W. Hethcote The Mathematics of infectious diseases *SIAM Review* **42**, 599–653 (2000).
- Carroll, M.W. *et al.* Temporal and spatial analysis of the 2014–2015 ebola virus outbreak in west africa. *Nature* **524**, 97–101 (2015).
- Gire, S. K. *et al.* Genomic surveillance elucidates ebola virus origin and transmission during the 2014 outbreak. *Science* **12**, 1369–1372 (2014).

23. Tong, Y.-G. *et al.* Genetic diversity and evolutionary dynamics of ebola virus in sierra leone. *Nature* **524**, 93–96 (2015).
24. Park, D. J. *et al.* Ebola virus epidemiology, transmission and evolution during seven months in sierra leone. *Cell* **161**, 1516–1526 (2015).
25. Silverman, B. W. *Density Estimation for Statistics and Data Analysis* (Chapman and Hall/CRC, 1986).
26. Vetter, P. *et al.* Ebola Virus Shedding and Transmission: Review of Current Evidence. *J. Infect. Dis.*, doi: 10.1093/infdis/jiw254 (2016).
27. Christie, A. *et al.* Possible Sexual Transmission of Ebola Virus - Liberia, 2015. *MMWR* **64**, 479–481, (2015).
28. Meyer, S., Elias, J. & Hohle, M. A space-time conditional intensity model for invasive meningococcal disease occurrence. *Biometrics* **68**, 607–616 (2012).
29. Simini, F., González, M. C., Maritan, A. & Barabási, A.-L. A universal model for mobility and migration patterns. *Nature* **484**, 96–100 (2012).

## Acknowledgements

The authors would like to thank the Santa Fe Institute and St. John's College for providing a productive working environment at the Complex System Summer School, and Samuel Scarpino for our many fruitful discussions there. The authors are also indebted to Caitlin Rivers' explanatory notes on the interpretation of the WHO patient database and situation reports. Finally, K. B. would like to thank Michelle Girvan for her many insightful suggestions. Financial support was provided by the University at Albany, State University of New York, the University of California at Davis, and National Natural Science Foundation of China under Grant Nos. 11575041, 61473060, 61603074.

## Author Contributions

J.H., K.B. and C.V. conceived the models, C.V. and M.I. cleaned WHO Situation Reports and the patient database data on Ebola, all authors analyzed results, and all authors reviewed the manuscript.

## Additional Information

**Supplementary information** accompanies this paper at <http://www.nature.com/srep>

**Competing financial interests:** The authors declare no competing financial interests.

**How to cite this article:** Burghardt, K. *et al.* Testing Modeling Assumptions in the West Africa Ebola Outbreak. *Sci. Rep.* **6**, 34598; doi: 10.1038/srep34598 (2016).



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

© The Author(s) 2016