

# Tracking virus outbreaks in the 21st century

Nathan D Grubaugh<sup>1,2</sup>, Jason T Ladner<sup>3\*</sup>, Philippe Lemey<sup>4</sup>, Oliver G Pybus<sup>5</sup>, Andrew Rambaut<sup>6,7\*</sup>, Edward C Holmes<sup>8\*</sup>, Kristian G Andersen<sup>1,9</sup>

<sup>1</sup>Department of Immunology and Microbiology, The Scripps Research Institute, La Jolla, CA 92037, USA

<sup>2</sup>Department of Epidemiology of Microbial Diseases, Yale School of Public Health, New Haven, CT 06510, USA

<sup>3</sup>Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ 86011, USA

<sup>4</sup>Department of Microbiology and Immunology, Rega Institute, KU Leuven - University of Leuven, 3000 Leuven, Belgium

<sup>5</sup>Department of Zoology, University of Oxford, Oxford OX2 6GG, UK

<sup>6</sup>Institute of Evolutionary Biology, University of Edinburgh, Edinburgh EH9 3FL, UK

<sup>7</sup>Fogarty International Center, National Institutes of Health, Bethesda, MD 20892, USA

<sup>8</sup>Marie Bashir Institute for Infectious Diseases and Biosecurity, School of Life and Environmental Sciences and Sydney Medical School, Charles Perkins Centre, The University of Sydney, Sydney, New South Wales 2006, Australia

<sup>9</sup>Scripps Translational Science Institute, La Jolla, California 92037, USA

\* = To whom correspondence should be addressed: JTL = jason.ladner@nau.edu; AR = a.rambaut@ed.ac.uk; ECH = edward.holmes@sydney.edu.au

## Abstract

Emerging viruses have the potential to impose substantial mortality, morbidity and economic burdens on human populations. Tracking the spread of infectious diseases to assist in their control has traditionally relied on the analysis of case data gathered as the outbreak proceeds. Here, we describe how many of the key questions in infectious disease epidemiology, from the initial detection and characterization of outbreak viruses, to transmission chain tracking and outbreak mapping, can now be much more accurately addressed using recent advances in virus sequencing and phylogenetics. We highlight the utility of this approach with the hypothetical outbreak of an unknown pathogen, 'Disease X', suggested by the World Health Organization to be a potential cause of a future major epidemic. We also outline the requirements and challenges, including the need for flexible platforms that generate sequence data in real-time, and for these data to be shared as widely and openly as possible.

## Introduction

Emerging infectious diseases present one of the greatest public health challenges of the 21<sup>st</sup> century. Among these are zoonotic viruses that originate from reservoir species, often mammals, and jump to humans to cause disease syndromes of varying form and severity. An emerging virus, depending on its ability to transmit among humans, can lead to individual or a few sporadic cases, resulting in a localized outbreak that requires public health intervention or, in the worst scenarios, develop into a large epidemic or global pandemic. Such emergence events over the past two decades are numerous and varied. They include viruses not previously encountered, such as the SARS and MERS coronaviruses<sup>1-3</sup>, and familiar foes that have reappeared to cause outbreaks, such as swine- and avian-origin influenza<sup>4,5</sup>, Ebola<sup>6</sup>, and Zika viruses<sup>7</sup>. Although many outbreaks end naturally, or are

45 controlled quickly, there remain questions over how best to scientifically respond to these  
46 events.

47 The broad-scale factors responsible for viral emergence have been well documented, and  
48 include human population growth, the increased frequency and reach of travel, changing  
49 patterns of land use, changing diets, wars and social upheaval, and climate change<sup>8,9</sup>. These  
50 increase interactions between humans and reservoir hosts, facilitating exposure to zoonotic  
51 viruses and spillover infections in people, and allow emerging viruses to spread more easily  
52 through human populations. The interactions between virus genetics, ecology, and the host  
53 factors that determine virus emergence are so complex that it is impossible to predict what  
54 virus will cause the next epidemic, making it essential that our response is scientifically  
55 informed, robust, and efficient<sup>10</sup>.

56 The emergence of virus outbreaks generates a set of common questions, whose answers  
57 are central to disease mitigation and control (**Table 1**) and which at times can only be  
58 answered by sequencing of viral genomes. These include what is the virus, is it novel, or  
59 does it represent the reemergence of a known pathogen; what is its mode of transmission;  
60 where does the emerging virus come from (in particular, what is its reservoir host and/or  
61 geographic source); what ecological factors underpin its emergence; how many introductions  
62 into humans have there been; what is the timing of these introduction events, and was there  
63 a period of undetected transmission before the first reported case; during flare-ups and  
64 future outbreaks, how are they connected to previous events; and what is the nature of virus  
65 evolution and is there evidence for local adaptation? In the past, many of these questions  
66 were addressed using case (incidence) data, which led to estimates of key epidemic  
67 parameters such as the basic reproductive number ( $R_0$  - the expected number of secondary  
68 cases produced by each case at the start of the outbreak) that were used to inform epidemic  
69 control policy. Although still of fundamental importance, case data alone cannot inform public  
70 health management with the level of precision necessary for all targeted interventions.  
71 Recent advances in virus genome sequencing and phylogenetic analyses, however, mean  
72 that we are now in a position to answer such questions with molecular precision, and open  
73 new areas of investigations not previously possible based on epidemiological data alone  
74 (**Table 1**).

75 Virus genomics have been used to investigate infectious disease outbreaks for several  
76 decades. This is possible because viruses, particularly those with RNA genomes, generate  
77 genetic variation on the same timescale of virus transmissions, through a combination of  
78 high rates of mutation and replication<sup>11,12</sup>. Consequently, it is possible to infer  
79 epidemiological and emergence dynamics from virus genomes sampled and sequenced  
80 over short epidemic timescales. We term the science of using genomics and associated  
81 analyses 'genomic epidemiology'.

82 Initially, genomic approaches relied upon indirect methods (*e.g.*, restriction fragment length  
83 polymorphisms<sup>13</sup>) to infer genotypes and differentiate between virus strains. As direct  
84 sequencing technologies advanced there was a transition toward the use of nucleotide  
85 sequences from fragments of virus genomes for this purpose<sup>14–21</sup>. Now, thanks to advances  
86 in high throughput sequencing and decreasing costs<sup>22</sup>, most virus genomics studies utilize  
87 data sets containing tens to thousands of (near) complete virus genomes.

In this review, we will show how our ability to track and understand infectious disease outbreaks have been revolutionized by the addition of virus genomics data. We will highlight the varied uses of virus genomics during the different stages of viral outbreaks, from initial virus detection, to understanding the factors contributing towards global spread (**Box 1**). We will show how genomic epidemiology can be used to track the spread of emerging viruses, where the challenges lie, and establish an agenda for future work. Although we focus on human disease, the genome-based methodologies that we describe can be equally applied to animal and plant infections. Similarly, the increasing ability to rapidly sequence complete genomes of bacterial species, means that these technologies offer much to the study of emerging bacterial disease, including those associated with antimicrobial resistance.

## **Outbreak detection**

Most infectious disease outbreaks start with clinicians noticing unusual patterns. Patients may present with patterns of symptoms that are similar to those of more common diseases, but which, after repeated observation and diagnostic testing, may deviate in scale, seasonality, or severity. At this very beginning of an outbreak the most critical task is therefore to identify a causal pathogen. Historically, virus identification has been performed using molecular tools, such as PCR and ELISA - that directly recognize pathogen-derived material (**Box 2**) - or conventional non-molecular techniques such as microscopy. The advent of untargeted metagenomic sequencing directly from clinical samples, however, means that we are now on the cusp of being able to detect human viruses in a single step, without *a priori* knowledge of putative causal pathogens (**Box 2**). The major advantage of sequencing-based approaches is the ability to detect novel viruses - such as the initial appearances of SARS<sup>2</sup>, MERS<sup>3</sup>, or Lujo virus<sup>23</sup> - or unexpected ones, as exemplified by Ebola virus during the 2013-2016 epidemic in West Africa<sup>24</sup>.

Once an outbreak has been detected and a causal virus identified, several basic questions can immediately be answered about the virus itself, including: (1) whether it is novel or previously known to infect humans, and (2) if we have the diagnostics, vaccines, and therapeutics available to fight it. Importantly, the generation of virus genomics data at this stage will provide deeper insights into these questions by uncovering molecular details not possible with conventional tools. Phylogenetics will also provide an additional level of detail, revealing virus origins, evolutionary characteristics, and connections to previous outbreaks in the same region, or to transmissions in other regions<sup>6</sup>. Given high enough relatedness to other members of a virus family with well-defined reservoir hosts (*e.g.*, old-world arenaviruses<sup>25</sup>), the sequence identification of novel virus species can also be informative about potential reservoirs.

## **First snapshot of an outbreak**

Immediately after a viral outbreak has been identified there exists a 'fog-of-war'. The extent of the outbreak, the timing and nature of its source, and the contribution of human-to-human transmission will be extremely limited, yet these data are critical to designing effective responses. Genomic epidemiology, if applied quickly and comprehensively, holds the potential to answering these questions<sup>24</sup>.

To provide an initial snapshot of an outbreak, it is important to understand the diversity of circulating viruses from as many cases as possible. Virus genetic diversity, measured as the average number of nucleotide differences among viruses in the population, will increase as

an outbreak progresses due to the accumulation of genetic changes in virus genomes at each round of viral replication<sup>6</sup>. If this rate of mutational accumulation is relatively constant - that is, it conforms to a 'molecular clock' of evolutionary change<sup>26</sup> - then the rate at which it occurs (referred to as the 'evolutionary rate') allows us to estimate when the sequenced viruses last shared a common ancestor. Critically, this provides a lower bound on when an outbreak began, and how long the virus had been circulating prior to discovery<sup>5,27,28</sup>. If the virus genomes have been sampled over only a limited time-scale, so that only a few mutations have accumulated in the virus population, then evolutionary rates will need to be based on those from prior outbreaks or extrapolations from related viruses<sup>29</sup>. Later in an epidemic, when viruses have been sequenced over a sufficient period of time to capture mutational accumulation, evolutionary rates can be readily estimated directly from virus genomes sampled during the outbreak<sup>30-32</sup>. Evolutionary rate estimates, however, can be sensitive to model specification over short periods of time<sup>33</sup> and depend on the timescale of measurement<sup>34</sup>. Such issues, as well as the unwarranted implications about changes in transmissibility and virulence that may accompany seemingly inflated evolutionary rates, have been discussed in detail in the context of the 2013-2016 Ebola epidemic in West Africa<sup>6</sup>.

A common approach to phylogenetic analysis of the genetic diversity of a virus population is to infer a tree from sampled virus genomes with branches measured in units of time (*i.e.*, a rooted, time-calibrated tree). This can provide estimates of the date of the last common ancestor at the root of the tree, as well as each individual branching event. As an approximation, these branching events correspond to virus transmission from one case to the next, an insight that offers further key information about the unfolding outbreak<sup>35</sup>. In addition, models of how the process of virus transmission relates to the shape of phylogenetic trees (**Fig. 1**) enable important epidemiological inferences. In particular, coalescent models relate the rate at which virus lineages of a phylogenetic tree merge, as common ancestors, to the size of the epidemic. This uses the simple premise that, for a sample of virus genomes, the larger the outbreak is, the further back in time the common ancestor will be found (**Fig. 1a-c**).

Early in an outbreak, one of the primary concerns is to understand the rate at which the virus may be spreading through the human population. As noted in the Introduction, this can be assessed by estimating  $R_0$ , which is critical for epidemiological projections and for planning public health responses. While  $R_0$  can be calculated through epidemiological analyses of case counts, accurate estimates of such data may not be available early in an outbreak, since they require a time-series of cases. As demonstrated during the early spread of the novel influenza A/H1N1 virus in 2009, phylogenetic inference of epidemic growth based on virus genomics can provide estimates of  $R_0$  comparable to that inferred from case data<sup>36</sup>. These calculations can be performed using coalescent models that directly estimate  $R_0$ , based on classic susceptible-infected-recovered (SIR) models<sup>37,38</sup>. A similar group of models analyze patterns of lineage birth-death, linking the shape of trees to the rate at which virus lineages split and go extinct, and have recently gained popularity<sup>39,40</sup>. Both approaches were applied during the 2013-2016 epidemic in West Africa to calculate  $R_0$  to assess Ebola virus transmission dynamics, and illuminated the impact of 'superspreader' events<sup>41,42</sup>. All of these methods, however, are beholden to the inherent uncertainty of genome sequence data, especially at the start of an epidemic where such sequences exhibit limited variability and sampling may be biased. Hence, phylogenetic estimates of  $R_0$ , although likely indicative of

broad characteristics such as epidemic growth, may not be precise enough to make critical decisions in the absence of corroborating (epidemiological) information.

The initial snapshot of virus genome sequences can also provide critical insights into the role of a zoonotic transmission during an outbreak (**Fig. 1d**). Genomic analyses, for example, revealed that Lassa fever virus, which is endemic in West Africa<sup>43</sup>, primarily spreads via repeated transmission from local rodent reservoirs, as opposed to sustained human-to-human transmission<sup>44</sup>. This is in contrast to Ebola virus during the 2013-2016 epidemic in West Africa, where genomic epidemiology showed that the outbreak was the result of a single zoonotic spillover, followed by sustained human-to-human transmission<sup>45</sup>.

Given availability of virus genomes from potential zoonotic reservoirs, another aim of early virus sequencing from an outbreak is to uncover the identity and geographic location of the reservoir host. The influenza A/H1N1 pandemic that started in 2009 was quickly recognised as being a likely species jump from pigs, as all of the virus genomic segments closely matched those previously seen in swine<sup>4,5</sup>. Like the 2013-2016 Ebola epidemic in West Africa, the influenza A/H1N1 pandemic likely started as a single introduction into humans that occurred a few months before it was detected<sup>5</sup>. The initial suspicion, and later confirmation, that the spillover occurred in Mexico, was complicated by a lack of widespread zoonotic genomic surveillance in this region. Retrospective sequencing of samples from Mexican pigs, however, showed that there were close relatives of the human virus circulating in this country at the time of the epidemic, confirming its origin<sup>46</sup>.

## Transmission chain tracking

Beyond the initial characterization of an outbreak, virus genome sequencing offers enormous potential for determining transmission chains to understand networks of 'who-infected-whom'. The tracking of transmission chains has long been a standard part of public health responses to outbreaks, providing critical information that can be used to interrupt virus spread and reduce the magnitude of an outbreak. This work has traditionally been performed using interview-based contact tracing, which is labor-intensive and limited by the availability and openness of patients for interviews. This approach is particularly challenging during large outbreaks characterized by large numbers of co-occurring transmission chains.

Virus genomic-based approaches can provide much more in-depth information compared to traditional non-sequencing based approaches, as the branching patterns of phylogenetic trees approximately correspond to transmissions from one case to the next (**Fig. 1**)<sup>35</sup>. Virus genome sequences, for example, were used to reconstruct the spread of foot-and-mouth disease virus in the United Kingdom, including the identification of superspreader events<sup>47-49</sup>. Genomic data also played a critical role in understanding flare-ups during the West African Ebola outbreak<sup>50-52</sup>, where phylogenetic analyses showed that most of the flare-ups were linked to persistently infected Ebola survivors (**Fig. 2a**), thereby demonstrating sexual transmission of the virus<sup>50,52</sup>. None of these insights would have been possible without virus genomic data.

The utility of virus genomic data for the inference of transmission chains is dependent on several factors, including: (1) the evolutionary rate of the virus, (2) the length of time between the infections of interest, and (3) the proportion of sampled cases, which together determine the resolution of the genetic signal (**Fig. 2b**). Although RNA viruses exhibit remarkably high evolutionary rates<sup>53</sup>, their small genome sizes and short epidemiological

generation times often result in, on average, less than one substitution per transmission event (**Fig. 2b**)<sup>54–56</sup>. Hence, virus genomics alone often cannot be expected to perfectly reconstruct transmission chains at the level of individual infections. Combined with epidemiological data, however, virus genomics provides a powerful tool for restricting the number of possible transmission scenarios and for supporting novel modes of transmission<sup>47,57,58</sup>. In addition, most phylogenetics-based transmission chain analyses have been performed using virus consensus sequences (*i.e.*, a single genome per sample/patient that represents the average of the virus population), which may limit resolution. However, as virus infections exhibit diverse intra-host populations (containing intra-host single nucleotide variants, or iSNVs<sup>44</sup>), newer methods incorporating viral iSNVs may greatly increase the resolution of transmission chain analyses so long as multiple variants are transmitted between hosts<sup>59</sup>.

## Outbreak mapping

As described in the previous sections, genomic epidemiology can be used to detect an outbreak, show its origin, and elucidate transmission patterns. Evolutionary inferences from virus genomes, unlike non-sequencing based methods, can also be used to dissect the spatial structure and dynamics of spread, as well as assess how an epidemic may unfold through time and space.

Uncovering the spatial patterns of virus spread during outbreaks is a key objective that has been transformed by genomic epidemiology. Reconstructing a detailed spatial history of virus spread from the origin of an outbreak is generally a task for phylogeographic methods<sup>60</sup>, which provide location estimates for every ancestral node in a virus phylogeny using simple stochastic (or ‘random walk’) models. Phylogeographic analyses, for example, were used to show how Ebola virus spread across West Africa during the 2013-2016 epidemic (**Fig. 3**)<sup>61</sup>. Importantly, virus genome sampling with strong spatiotemporal coverage allowed for the dissection of the entire epidemic into a metapopulation of short- and long-lived transmission chains<sup>61</sup>. Similar analyses were also used to show that multiple introductions were responsible for sustaining the 2016 Zika outbreak in Florida<sup>62</sup>. It is important, however, to appreciate the uncertainty of phylogeographic estimates, and to bear in mind that such analyses may only be capable of elucidating partial pictures of outbreak spread. In addition, sampling biases may severely affect these analyses, although the coalescent and birth-death models mentioned above have been extended to account for aspects of virus population structure<sup>63–65</sup>, making the analyses more robust to sampling heterogeneity<sup>66</sup>.

Phylogeographic inference methods can also be used to provide insights into the factors driving virus spread (**Fig. 3**)<sup>67</sup>. Such analyses are enabled by the integration of virus genomics with diverse meta-data sets and are critically dependent on the timeliness of data generation and open sharing. These approaches were initially introduced to confirm the key role of human air transportation in the global circulation of influenza viruses<sup>67</sup>, but they have also been useful in untangling complex virus transmission dynamics on smaller scales<sup>61</sup>. To illustrate these methods, in Figure 3 we show an application of generalized linear modeling to explain Ebola virus migration rates between locations as a function of several potential predictors, to infer virus spread during West African Ebola outbreak (**Fig. 3**). In this case, geographic distances and population sizes at the location of origin and destination combine into a gravity-model of spread, where virus transmission largely occur within large population

centers and geographic spread being more frequent over shorter distances<sup>61</sup>. These phylodynamic studies illustrate the growing importance of data integration for virus genomic analyses<sup>55</sup>, which critically depend on accurate metadata (*e.g.*, sampling date and sampling location), as well as other data sources that can capture host mobility and geographic, demographic, and epidemiological context.

## **Inter-epidemic evolution and spread**

Once outbreaks have been brought under control or been (temporarily) resolved, phylogenetic analyses can provide insights into evolutionary patterns during inter-epidemic periods, by comparing virus genome sequences sampled across different outbreaks. The most fundamental question is whether the virus in question has been able to persist in human populations between outbreaks, so that each new outbreak has arisen from an endemically circulating lineage (*e.g.*, dengue virus), or whether they represent independent zoonotic spillover events from an animal reservoir (*e.g.*, Ebola virus). With sufficient sampling of viruses from human and reservoir species this question can be answered using standard phylogenetic analysis. For example, although both dengue virus and yellow fever virus have transmission cycles that involve mosquitoes and humans (urban transmission) or nonhuman primates (sylvatic transmission), phylogenetic analyses have shown that dengue virus is now an entirely endemic urban virus that does not rely on its sylvatic vectors and hosts to seed new epidemics<sup>68</sup>. Most human outbreaks of yellow fever, in contrast, have been shown by virus genomics approaches to represent independent emergences of the virus from sylvatic sources, rather than spread via an urban cycle<sup>69,70</sup>.

Inter-epidemic analyses can also be used to elucidate the nature of virus evolution and spread in reservoir species, which are likely characterized by different evolutionary forces than those seen during human outbreaks<sup>71,72</sup>. For example, although human outbreaks of Ebola have happened relatively frequently since the 1970s, each outbreak starts as an independent spillover of the virus from an animal (likely bat<sup>73</sup>) reservoir. Hence, the inter-epidemic evolution of Ebola virus occurs in a species other than humans, such that patterns of genetic divergence among the viruses associated with human epidemics can provide insight into viral replication and transmission within reservoir hosts. For example, there have been suggestions that Ebola virus has spread across Africa in a wave-like manner in its reservoir species<sup>74</sup>; however, phylogenetic analyses incorporating virus genomic data from recent outbreaks are incompatible with this scenario<sup>6</sup>. Additionally, while Ebola virus normally evolves according to a relatively constant molecular clock<sup>6,45,75–77</sup>, the phylogenetic branch leading to the viruses sequenced from the small Ebola outbreak that occurred in the Democratic Republic of the Congo in 2014, concurrent with the 2013-2016 epidemic in West Africa, was characterized by a far lower evolutionary rate<sup>78</sup>. Although the reasons for this reduction in evolutionary tempo are unclear, it is possible that it reflects Ebola virus evolution in a different (unknown) reservoir species that experiences a lower rate of viral replication. Alternatively, this rate disparity may result from the existence of different viral replication states within the same reservoir host, similar to that described during human epidemics, with faster rates observed during continuous human-to-human transmission and slower rates during persistent infections of Ebola survivors<sup>79</sup>.

## **Requirements and challenges in genomic epidemiology**

Virus genomic methods for outbreak investigation and control are powerful additions to more traditional epidemiological approaches but are critically dependent on well planned and

coordinated efforts. The foremost need for genomic epidemiology is timely access to clinical samples and data, which should be built on productive and equitable collaborations with local communities, public health agencies, outbreak responders, local clinics, and researchers<sup>80</sup>. For each clinical sample to be used for virus genomic sequencing, it is essential to obtain a minimal set of metadata related to the infection, including (1) the date of sample collection and/or onset of symptoms and (2) the location of sampling. Additional information can greatly increase the utility of genomic epidemiology, including the availability of (3) travel and contact history, (4) suspected source of infection, and (5) clinical outcome and symptoms. Other factors, including patient history, age, sex, and economic status can also help to reveal risk factors underlying infection and transmission. Within ethical constraints, it is important that communication lines remain open so that researchers undertaking data analysis can return actionable results to the public health community.

Other large-scale data resources are essential for investigating the spatio-temporal history and spread of an outbreak. These include the temporal and spatial distribution of cases, ecological conditions, vector abundance, environmental factors, and travel patterns. Integration of these other data sources with virus genomic data may reveal new properties of an outbreak, potentially leading to actionable measures<sup>55,61,67</sup>. Non-genomic data often comes from established networks of collaborations, or from the public domain, highlighting the value of open data and data sharing to outbreak investigations.

An important benefit of genomic epidemiology is that it can directly compare and jointly analyze virus genome sequences obtained during an epidemic, even if those sequences were generated by different laboratories. Consequently, there is an urgent need to make genomic and epidemiological data and analysis tools publically available during ongoing epidemics<sup>81</sup>. This movement is supported by the World Health Organization (WHO), who have called for data pertaining to public health emergencies to be disseminated openly and immediately upon generation, and not withheld until the acceptance or publication of a corresponding scientific paper<sup>82</sup>. More recently, the WHO has outlined the current and future benefits of virus genome data sharing during outbreaks<sup>83</sup>. Combined with an acceleration of making manuscripts available via preprint servers, such as arXiv and bioRxiv, especially during outbreaks<sup>84</sup>, there has been a shift towards scientists storing their data and source code on depositories like GitHub ([github.com](https://github.com)), Synapse ([synapse.org](https://synapse.org)), and Data Dryad ([datadryad.org](https://datadryad.org)), in close to real-time for others to use. Furthermore, extensive online communities and forums like Twitter ([twitter.com](https://twitter.com)), Virological ([virological.org](https://virological.org)), FluTrackers ([flutrackers.com](https://flutrackers.com)), ProMED ([promedmail.org](https://promedmail.org)), Nextstrain ([nextstrain.org](https://nextstrain.org)), HealthMap ([healthmap.org](https://healthmap.org)), and Microreact ([microreact.org](https://microreact.org)) allow for rapid dissemination of unpublished results and analyses. In our experience, not only does the process of open science promote new collaborations and lead to more accurate scientific insights into outbreak research, but it helps in getting relevant information rapidly into the hands of decision makers. Despite these advances, however, the speed, nature, and extent of virus genome data sharing is inconsistent, sometimes resulting in confusion over what is, or should be, best practice<sup>81,85</sup>.

## **Future perspective**

Genomic epidemiology promises much to the study and control of infectious disease outbreaks, particularly if viral genomes can be acquired and analyzed in real-time. The accumulated set of these data - together with the rapid development of sophisticated



software packages ([virological.org/c/software](http://virological.org/c/software)) - will provide a valuable resource for the mitigation and control of future outbreaks. Ultimately, with sufficient genome sequences from individual viral genera and/or families, it may be possible to categorize viruses by their phylogenetic patterns and utilize this information in epidemic preparedness. For example, as well as considering obvious biological features of viruses such as their genome structure and mode of transmission, it may be possible to group viruses according to a series of evolutionary variables such as rate of evolutionary change, extent of antigenic evolution, frequency of recombination, pattern of geographic spread, and population dynamics. This information may then help forecast the evolutionary behaviour of any virus should it reemerge in human populations and assist in the selection of future vaccine strains<sup>86–88</sup>. This information will also help counter alarmist claims that emerging viruses will evolve novel phenotypes, such as airborne transmission in the case of Ebola virus<sup>89</sup>, that often accompany any major disease outbreak. It is clear, however, that a more fundamental understanding of the genetic and ecological barriers of virus spillover into human populations is needed to better identify risk factors for disease emergence. Long-term capacity building, partnerships with local communities, and commitments to long-term investments on these fronts will go a long way towards better enabling the global community to effectively and rapidly deal with future emerging outbreaks<sup>80</sup>.

## Acknowledgements

We thank Gytis Dudas and Sigrid Knemeyer for help with figure creation. NDG is supported by NIH training grant 5T32AI007244-33. PL and AR acknowledge funding from the European Research Council under the European Union's Horizon 2020 research and innovation programme (grant agreement no. 725422-ReservoirDOCS) and from the Wellcome Trust Collaborative Award (grant number 206298/Z/17/Z - ARTICnetwork). PL acknowledges support by the Research Foundation - Flanders ('Fonds voor Wetenschappelijk Onderzoek - Vlaanderen', G066215N, G0D5117N and G0B9317N). OGP is supported by the European Union's Seventh Framework Programme (FP7/2007-2013)/European Research Council (614725-PATHPHYLODYN) and by the Oxford Martin School. ECH is supported by an ARC Australian Laureate Fellowship (FL170100022). KGA is a Pew Biomedical Scholar, and is supported by NIH NCATS CTSA UL1TR002550, NIAID contract HHSN272201400048C, NIAID R21AI137690, NIAID U19AI135995, and The Ray Thomas Foundation.

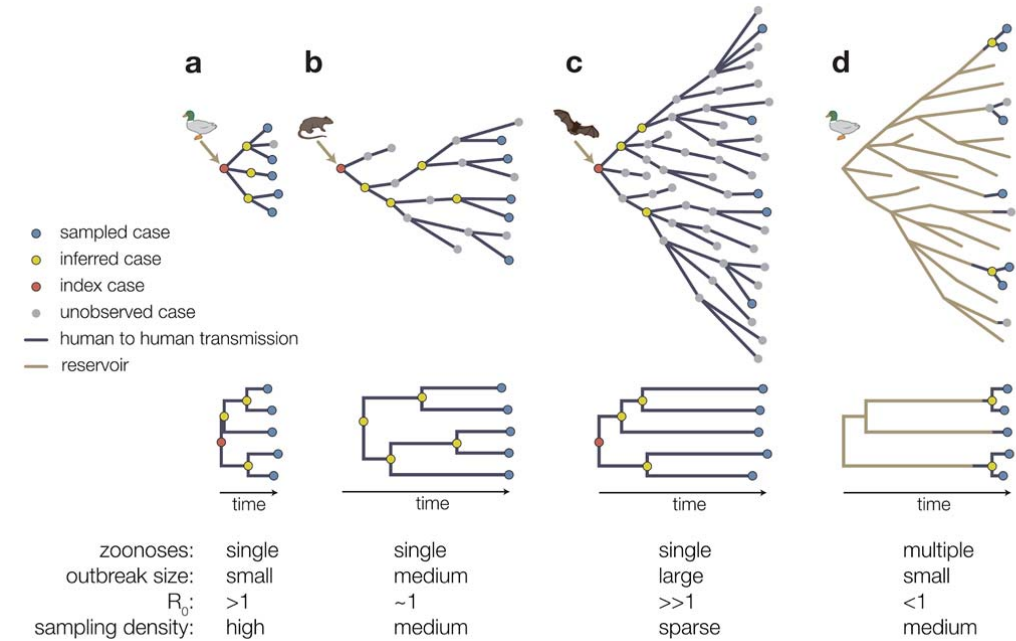
## Contributions

All listed authors have contributed to the conceptualization, writing and preparation of the manuscript.

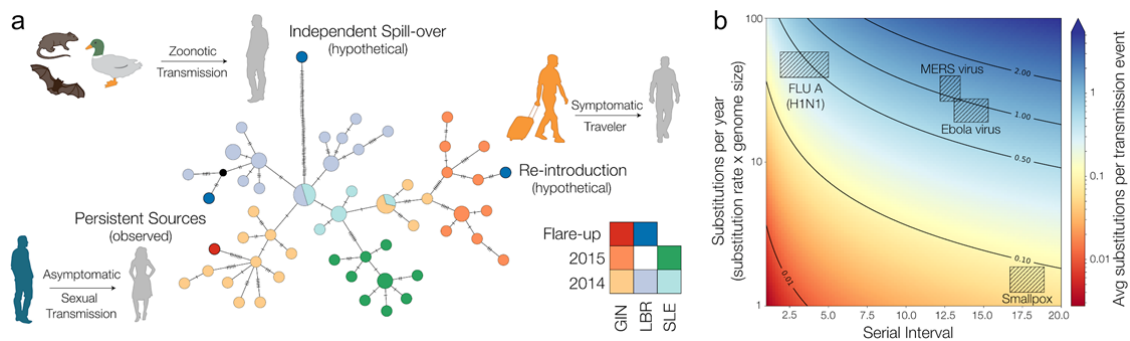
## Competing interests

The authors declare no competing interests.

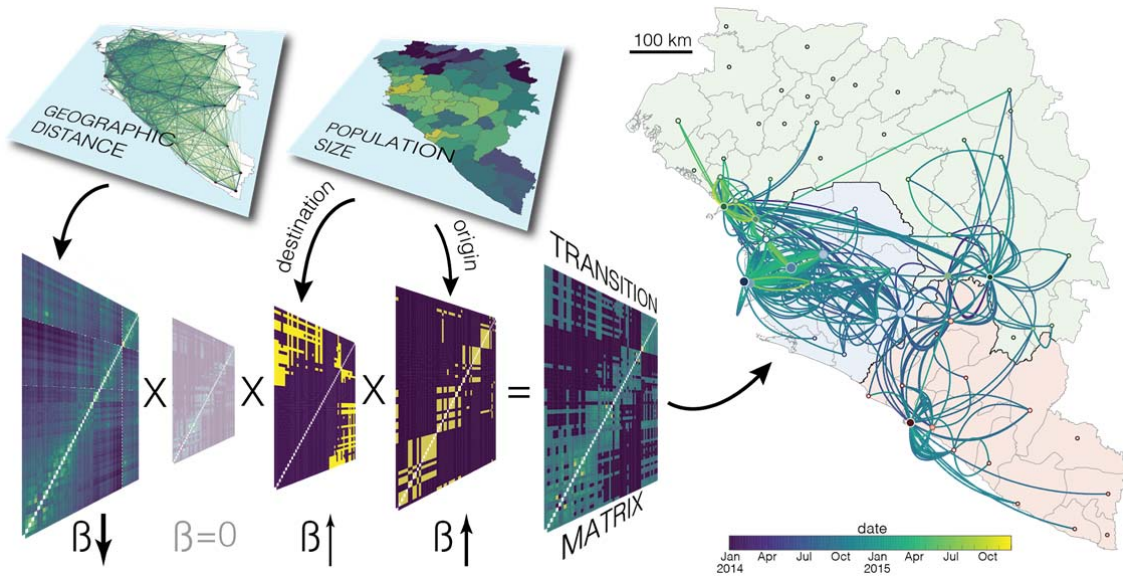
Figures



**Figure 1. Outbreak scenarios and the resulting phylogenetic trees of virus genomes from sampled human cases.** The first three scenarios show a single introduction from a non-human reservoir followed by human-to-human spread. (a) A small outbreak from a recent zoonosis with a commensurately short tree, suggesting recent emergence. The  $R_0$  is greater than 1 indicating the potential to cause a large outbreak. (b) A medium sized outbreak with a deeper tree and internal nodes dispersed. With an  $R_0$  close to 1, this suggests that emergence into humans was not recent and its transmission potential is just sufficient to persist. The root of the tree is not the index case meaning the zoonosis could be older. (c) A large outbreak with an  $R_0$  greater than 1 and thus exhibiting exponential growth in case numbers. Distinctively for a growing epidemic, internal nodes tend to be towards the root of the tree, suggesting that only a small fraction of the total cases were sampled. (d) A scenario of repeated zoonotic jumps with limited human to human transmission. The internal parts of the tree represent the diversity of the virus in the non-human reservoir and the human-to-human transmission cases are closely related.



**Figure 2. Transmission chain tracking during outbreaks using virus genomics. (a)** Viral genome sequences were used to distinguish between competing hypotheses for the source of the viruses that triggered the Ebola flare-ups in West Africa. The three main hypotheses and their expected genomic signatures are illustrated here with a hypothetical haplotype network. Genomes from all of the observed flare-ups grouped closely with genomes sequenced from patients in the same country, from earlier in the outbreak (bottom left), consistent with transmission from persistent sources. In contrast, genomes linked to re-introductions from neighboring countries (right) would be expected to cluster with genomes from a different country and from late in the outbreak, and in the case of independent spillovers from a reservoir host (top left, *i.e.*, independent sampling from the diversity circulating within the reservoir), the spillover genomes would be linked to the main outbreak by a long branch originating from near the root of the network. GIN, Guinea; LBR, Liberia; SLE, Sierra Leone. **(b)** Expected 'genomic resolution' for the inference of transmission chains at the level of individual infections. Resolution is dependent on the serial interval between infections (x-axis; used as a proxy for epidemiological generation time) as well as the genome size and nucleotide substitution rate (y-axis).



**Figure 3. Integration and testing predictors of phylogeographic spread.** We illustrate the concept of this approach using the 2013-2016 Ebola epidemic in West Africa. Geographic distances between all pairs of locations, in this case administrative areas in Guinea, Sierra Leone and Liberia, as well as population sizes at the origin and destination of these pairs are combined into a transition rate matrix through a generalized linear model. This matrix parameterizes the phylogenetic process of spread that is being estimated. Each predictors is associated with a coefficient,  $\beta$ , which denotes the strength of contribution with some predictors (*e.g.*, population size) positively associated with the intensity of migration whereas others (*e.g.*, geographic distance) are negatively associated. A coefficient of zero implies that the predictor is excluded from the model (represented in the figure by the transparent matrix with  $\beta=0$ ).

## Tables/Boxes

**Table 1 – Critical questions addressed by viral genomic epidemiology.**

Questions	Examples from Genomic Epidemiology	References
What virus is causing the outbreak?	Metagenomic sequencing from patient samples revealed a novel virus - Lujo virus - as the causal virus for an outbreak in South Africa in 2008.	<sup>23</sup>
How is the virus transmitting?	Sequencing studies of MERS-CoV combined with coalescent approaches showed that human outbreaks are driven by seasonally varying zoonotic transfer of viruses from camels.	<sup>90,91</sup>
Where did the outbreak begin?	Large-scale sequencing efforts and phylogenetic analyses showed that the 2009 influenza A/H1N1 pandemic originated in swine populations from Mexico.	<sup>5,46</sup>
What factors drive the outbreak?	Analysis of more than 1,600 Ebola virus genomes identified critical factors that contributed to the spread of the virus during the 2013-2016 epidemic in West Africa.	<sup>61</sup>
How many introductions have there been?	Sequencing of Zika virus from patients and mosquitos in Florida showed that multiple introduction events of the virus sustained the 2016 outbreak in Miami and surrounding counties.	<sup>62</sup>
When did the outbreak begin?	Large-scale studies showed that the Zika epidemic in the Americas likely started in Brazil more than a year earlier than was initially believed.	<sup>7,92–94</sup>
Are outbreaks linked?	Analysis of Ebola virus genomes during the 2013-2016 epidemic showed that the virus can persist for more than a year in survivors, and be responsible for flare-ups of the outbreak via sexual transmission.	<sup>52,57,95,96</sup>
How is the virus evolving?	Sequencing studies during the 2013-2016 Ebola epidemic identified mutations in the virus genome that rapidly rose to high frequency compatible with increased fitness. Experimental follow-up studies showed that some of those mutations were likely Ebola virus adapting to a new host.	<sup>71,72,97</sup>
Footnote: Examples of commonly used software packages for genomic epidemiology investigations are available at: <a href="http://virological.org/c/software">virological.org/c/software</a> .		

### **Box 1 - Outbreak of ‘Disease X’ - a hypothetical scenario**

In addition to the Ebola, SARS, and Zika viruses, the WHO watchlist of viruses that may lead to public health emergencies<sup>98</sup> for the first time acknowledged that the next serious epidemic may be caused by a currently unknown virus - ‘Disease X’. Its inclusion emphasizes the need for flexible and deployable platforms to understand and combat disease outbreaks of many varieties. Most likely, Disease X may be a known microbe believed to cause no or mild human disease, as was the case Zika virus before its epidemic in the Americas. Disease X could emerge anywhere in the world and, given the mobility of human populations, could spread to distant and highly populated regions within days or weeks. To illustrate how genomic epidemiology can successfully reveal important aspects of disease emergence and inform epidemic control efforts, we present a hypothetical scenario in which Disease X successfully jumped into humans, established sustained transmission, and caused severe disease.

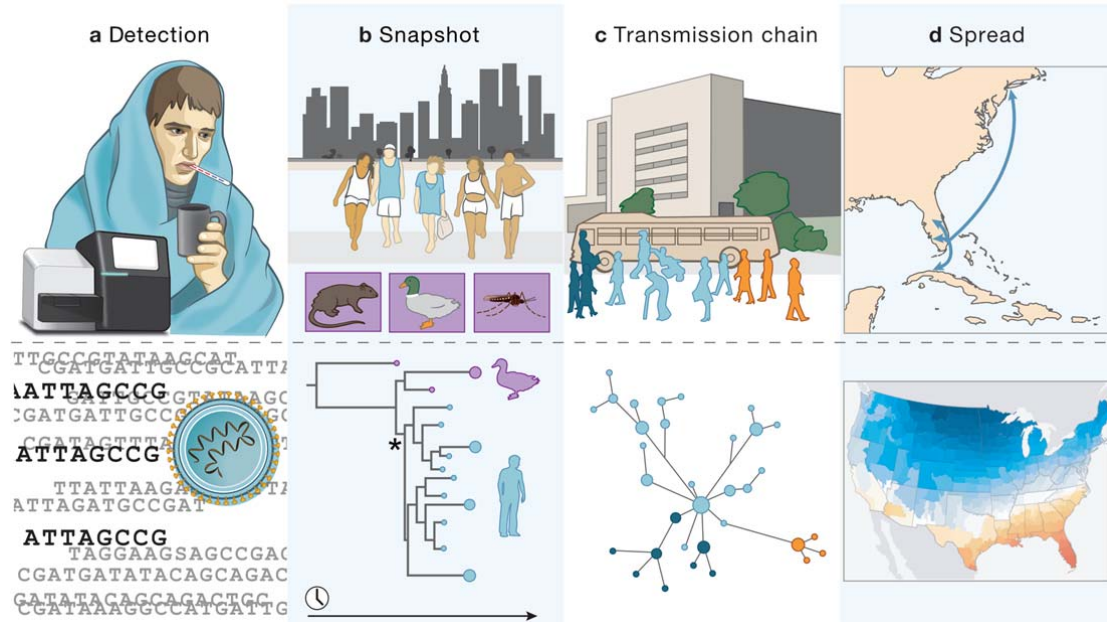
In Miami, Florida (United States), a 22 year old man sought medical assistance after an influenza-like illness suddenly progressed to a dangerously high fever and laboured breathing. He reported golfing activity at nearby resorts, harboring clusters of wildlife, including birds. He was admitted into the emergency room and within three days died of pneumonia. During this time, five other young adults presented with similar symptoms to Miami-area hospitals. Standard molecular diagnostics for commonly suspected pathogens were negative, but IgM antibodies collected from each patient were slightly cross-reactive to MERS and SARS coronaviruses. Since the virus could not be conclusively identified with conventional assays, metagenomic sequencing was used to identify Disease X as a novel human virus, most closely related to other coronaviruses in ducks (**Panel a**). Importantly, due to the relatedness of the novel virus to a family of viruses with well-defined host-ranges, these data led to a hypothesis about its potential origin and reservoir (overwintering migratory birds in the nearby Everglades wetlands) and allowed for the development of virus-specific diagnostics and targeted sequencing approaches.

Within three weeks, there were 40 new laboratory-confirmed Disease X cases, including eight from healthcare workers who contacted the original six cases, and five total deaths (an 11% apparent case fatality rate). Targeted sequencing from 15 patients and related viruses, including from ducks across Southern Florida, revealed that the human Disease X viruses clustered together on a phylogenetic tree and shared a common ancestor with virus genomes from ducks near Palm Beach, suggesting there was a single zoonotic spillover event and subsequent human-to-human transmission (**Panel b**). A molecular clock phylogenetic analysis further indicated that the common ancestor of the human viruses existed several months ago, suggesting that the first patient identified was not the first case of the outbreak, and highlighting the possibility of many more unreported or asymptomatic cases.

As the outbreak progressed, there was a critical need to understand transmission to help control further spread. Traditional epidemiology, including contact tracing, provided important insights into the risk factors for transmission. Virus genome sequencing was used to infer transmission chains that linked each infected patient (**Panel c**). These analysis revealed that (1) transmission occurred primarily between individuals that had been in close proximity and (2) a few individuals infected most of the known cases. In

response, an action plan of patient isolation/containment and widespread use of facemasks was implemented to reduce close contact and aerosol transmission.

The Disease X outbreak peaked within a year, resulting in ~2000 cases in Florida and several imported cases throughout the world. Most of the imported cases did not result in secondary local infections, with the exception of two healthcare workers in New York City and a large outbreak of more than a 100 cases near Havana, Cuba. Factors leading to local and global spread were investigated by layering transportation, geographic, climatic, economic, and demographic information into a large phylogenetic data set of Disease X viruses (**Panel d**). Analyses indicated that virus dispersal from Miami was more likely to occur to large cities that were either (1) in close driving proximity or (2) connected by direct flights with high travel volumes. Once in a new city, the success of virus transmission was correlated with low economic status and high population density. This raised concerns about Disease X outbreaks emerging in low income and densely populated countries within the Caribbean and Central America. The WHO used this information to implement comprehensive surveillance and response efforts in at risk nations.



**Real-time genomic investigation of 'Disease X'.** (a) Metagenomic sequencing revealed that 'Disease X', which could not be identified using standard clinical assays, was a novel virus. (b) Targeted sequencing from additional human cases and from related viruses uncovered the likely animal reservoir, the time period that it was introduced into the human population (represented by \* in the lower panel), and that subsequent transmission was human-to-human. (c) More intensive virus genome sequencing was used to construct detailed transmission chains and identify potential control measures. (d) Layering additional climatic (pictured in the lower panel; source: <https://www.climate.gov/maps-data>), transportation, geographic, economic, and demographic information into a large phylogenetic data set revealed the risk factors that facilitated local and global spread.



## **Box 2 - Molecular technologies for detecting and tracking outbreaks**

### **Traditional methods**

The methods traditionally used to diagnose infectious disease agents in patients are developed to detect either antigens (*e.g.*, ELISAs and lateral flow assays), or nucleic acids (*e.g.*, PCR) derived from the pathogen. These assays are typically designed to recognize either single (*e.g.*, Ebola virus) or closely related (*e.g.*, *Filoviridae*) pathogens. Versions of such assays may also be combined in a multiplexed fashion to detect a small number of different pathogens (*e.g.*, hemorrhagic fever viruses). While most laboratories are capable of running these assays, they are often not available for uncommon or novel pathogens, and running multiple rounds of testing can take weeks. They also require *a priori* knowledge of putative pathogens and they typically cannot be used to detect outbreaks that are caused by novel, highly divergent, understudied, or rare pathogens.

### **Deployable solutions**

Over the last several years, robust and deployable solutions have been developed for pathogen detection that do not require the maintenance of a cold chain, which can be difficult or impossible under many outbreak conditions. Simple-to-use, point-of-care rapid diagnostic tests (RDTs) have the potential to transform early outbreak detection. For example, the ReEBOV antigen rapid test for Ebola virus infection developed during the recent epidemic could be deployed throughout sub-Saharan Africa to help detect new outbreaks<sup>99,100</sup>. Simple nucleic acid assays, like loop-mediated isothermal amplification (LAMP) developed for Zika virus<sup>101</sup>, H5N1 avian influenza virus<sup>102</sup>, and SARS coronavirus<sup>103</sup>, have eliminated the need for thermal cycling and most power requirements. New and creative advances in microfluidics<sup>104</sup>, nanowire arrays<sup>105</sup>, and field-effect biosensors<sup>106,107</sup> are also helping to reduce the barriers to efficient and rapid diagnostics, while increasing sensitivity and specificity of detection. Of particular interest for deployment in resource-limited settings, are paper-based engineered gene circuits, like sensors designed for strain-specific Ebola virus detection<sup>108</sup>. They are stable for long-term storage at room temperature and are activated by rehydration, and thus can be used in remote environments. Very recently, highly sensitive and deployable CRISPR-based diagnostics have also been developed that utilizes CRISPR-Cas13/12a to detect pathogen-derived nucleic acids<sup>109–112</sup>. Similarly to the traditional methods described above, all of these tools require *a priori* knowledge of likely causal pathogens and the availability of antibodies, genome sequences, or other pathogen characteristics.

### **Sequencing-based methods**

Untargeted metagenomic sequencing provides a potential one-step solution for outbreak pathogen detection of both known and novel pathogens, and may be able to replace the need for multiple individual pathogen assays<sup>24</sup>. The main advantage of metagenomic sequencing is that it does not require *a priori* knowledge of the pathogen, but comes at the expense of specialized equipment, increased costs, and bioinformatic complexity. Although high backgrounds of host nucleic acid and/or low pathogen titers in clinical samples can make pathogen detection difficult, host gene depletion<sup>113</sup> and pathogen enrichment<sup>114,115</sup> methods can help alleviate these issues. After the first outbreak pathogen genome sequence has been obtained, targeted approaches using next-generation sequencing can also be developed. This was the case for both of the recent Zika and Ebola epidemics<sup>116,117</sup>, where cheaper and faster amplicon-based approaches were rapidly



developed and deployed to track both of the epidemics. The most common platforms used for these purposes are those developed by Illumina (*e.g.*, MiSeq and HiSeq), because they have high accuracy and throughput, but have high costs and relatively short read lengths (up to 300 bp). Cheaper portable devices, such as the miniaturized Oxford Nanopore MinION can help to produce data in close to real time directly in-country and under austere conditions<sup>93,117</sup>. This is a significant advancement because, along with open data sharing, rapid diagnostics and sequencing helps to promote a comprehensive and collaborative response network.

## References

1. Drosten, C. *et al.* Identification of a novel coronavirus in patients with severe acute respiratory syndrome. *N. Engl. J. Med.* **348**, 1967–1976 (2003).
2. Ksiazek, T. G. *et al.* A novel coronavirus associated with severe acute respiratory syndrome. *N. Engl. J. Med.* **348**, 1953–1966 (2003).
3. Zaki, A. M., van Boheemen, S., Bestebroer, T. M., Osterhaus, A. D. M. E. & Fouchier, R. A. M. Isolation of a novel coronavirus from a man with pneumonia in Saudi Arabia. *N. Engl. J. Med.* **367**, 1814–1820 (2012).
4. Novel Swine-Origin Influenza A (H1N1) Virus Investigation Team *et al.* Emergence of a novel swine-origin influenza A (H1N1) virus in humans. *N. Engl. J. Med.* **360**, 2605–2615 (2009).
5. Smith, G. J. D. *et al.* Origins and evolutionary genomics of the 2009 swine-origin H1N1 influenza A epidemic. *Nature* **459**, 1122–1125 (2009).
6. Holmes, E. C., Dudas, G., Rambaut, A. & Andersen, K. G. The evolution of Ebola virus: Insights from the 2013-2016 epidemic. *Nature* **538**, 193–200 (2016).
7. Grubaugh, N. D., Faria, N. R., Andersen, K. G. & Pybus, O. G. Genomic Insights into Zika Virus Emergence and Spread. *Cell* **172**, 1160–1162 (2018).
8. Morse, S. S. Factors in the Emergence of Infectious Diseases. in *Plagues and Politics* 8–26 (2001).
9. Wolfe, N. D., Dunavan, C. P. & Diamond, J. Origins of major human infectious diseases. *Nature* **447**, 279–283 (2007).
10. Holmes, E. C., Rambaut, A. & Andersen, K. G. Pandemics: spend on surveillance, not prediction. *Nature* **558**, 180–182 (2018).
11. Holland, J. *et al.* Rapid evolution of RNA genomes. *Science* **215**, 1577–1585 (1982).
12. Duffy, S., Shackelton, L. A. & Holmes, E. C. Rates of evolutionary change in viruses: patterns and determinants. *Nat. Rev. Genet.* **9**, 267–276 (2008).
13. Kiko, H., Niggemann, E. & Rüger, W. Physical mapping of the restriction fragments obtained from bacteriophage T4 dC-DNA with the restriction endonucleases SmaI, KpnI and BglII. *Mol. Gen. Genet.* **172**, 303–312 (1979).
14. Chungue, E., Deubel, V., Cassar, O., Laille, M. & Martin, P. M. Molecular epidemiology of dengue 3 viruses and genetic relatedness among dengue 3 strains isolated from patients with mild or severe form of dengue fever in French Polynesia. *J. Gen. Virol.* **74** (Pt 12), 2765–2770 (1993).
15. Lanciotti, R. S. *et al.* Origin of the West Nile virus responsible for an outbreak of encephalitis in the northeastern United States. *Science* **286**, 2333–2337 (1999).
16. Kinnunen, L., Pöyry, T. & Hovi, T. Generation of virus genetic lineages during an outbreak of poliomyelitis. *J. Gen. Virol.* **72** (Pt 10), 2483–2489 (1991).
17. McNearney, T. *et al.* Limited sequence heterogeneity among biologically distinct human immunodeficiency virus type 1 isolates from individuals involved in a clustered infectious outbreak. *Proc. Natl. Acad. Sci. U. S. A.* **87**, 1917–1921 (1990).
18. Nichol, S. T. *et al.* Genetic identification of a hantavirus associated with an outbreak of acute respiratory illness. *Science* **262**, 914–917 (1993).
19. Ou, C. Y. *et al.* Molecular epidemiology of HIV transmission in a dental practice. *Science* **256**, 1165–1171 (1992).
20. Power, J. P. *et al.* Molecular epidemiology of an outbreak of infection with hepatitis C virus in recipients of anti-D immunoglobulin. *Lancet* **345**, 1211–1213 (1995).
21. Rossouw, E., Tsilimigras, C. W. & Schoub, B. D. Molecular epidemiology of a

492 coxsackievirus B3 outbreak. *J. Med. Virol.* **34**, 165–171 (1991).

493 22. Shendure, J. & Ji, H. Next-generation DNA sequencing. *Nat. Biotechnol.* **26**, 1135–1145  
494 (2008).

495 23. Briese, T. *et al.* Genetic detection and characterization of Lujo virus, a new hemorrhagic  
496 fever-associated arenavirus from southern Africa. *PLoS Pathog.* **5**, e1000455 (2009).

497 24. Gardy, J. L. & Loman, N. J. Towards a genomics-informed, real-time, global pathogen  
498 surveillance system. *Nat. Rev. Genet.* **19**, 9–20 (2018).

499 25. Salazar-Bravo, J., Ruedas, L. A. & Yates, T. L. Mammalian reservoirs of arenaviruses.  
500 *Curr. Top. Microbiol. Immunol.* **262**, 25–63 (2002).

501 26. dos Reis, M., Donoghue, P. C. J. & Yang, Z. Bayesian molecular clock dating of species  
502 divergences in the genomics era. *Nat. Rev. Genet.* **17**, 71–80 (2016).

503 27. Rambaut, A. & Holmes, E. The early molecular epidemiology of the swine-origin  
504 A/H1N1 human influenza pandemic. *PLoS Curr.* **1**, RRN1003 (2009).

505 28. Korber, B. Timing the Ancestor of the HIV-1 Pandemic Strains. *Science* **288**, 1789–1796  
506 (2000).

507 29. Cotten, M. *et al.* Full-genome deep sequencing and phylogenetic analysis of novel  
508 human betacoronavirus. *Emerg. Infect. Dis.* **19**, 736–42B (2013).

509 30. Rambaut, A. Estimating the rate of molecular evolution: incorporating non-  
510 contemporaneous sequences into maximum likelihood phylogenies. *Bioinformatics* **16**,  
511 395–399 (2000).

512 31. Drummond, A., Pybus, O. G. & Rambaut, A. Inference of viral evolutionary rates from  
513 molecular sequences. *Adv. Parasitol.* **54**, 331–358 (2003).

514 32. Drummond, A. J., Nicholls, G. K., Rodrigo, A. G. & Solomon, W. Estimating mutation  
515 parameters, population history and genealogy simultaneously from temporally spaced  
516 sequence data. *Genetics* **161**, 1307–1320 (2002).

517 33. Möller, S., du Plessis, L. & Stadler, T. Impact of the tree prior on estimating clock rates  
518 during epidemic outbreaks. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 4200–4205 (2018).

519 34. Duchêne, S., Holmes, E. C. & Ho, S. Y. W. Analyses of evolutionary dynamics in  
520 viruses are hindered by a time-dependent bias in rate estimates. *Proc. Biol. Sci.* **281**,  
521 (2014).

522 35. Hall, M. D., Woolhouse, M. E. J. & Rambaut, A. Using genomics data to reconstruct  
523 transmission trees during disease outbreaks. *Rev. Sci. Tech.* **35**, 287–296 (2016).

524 36. Fraser, C. *et al.* Pandemic potential of a strain of influenza A (H1N1): early findings.  
525 *Science* **324**, 1557–1561 (2009).

526 37. Volz, E. M., Kosakovsky Pond, S. L., Ward, M. J., Leigh Brown, A. J. & Frost, S. D. W.  
527 Phylodynamics of infectious disease epidemics. *Genetics* **183**, 1421–1430 (2009).

528 38. Rasmussen, D. A., Ratmann, O. & Koelle, K. Inference for nonlinear epidemiological  
529 models using genealogies and time series. *PLoS Comput. Biol.* **7**, e1002136 (2011).

530 39. Stadler, T. *et al.* Estimating the basic reproductive number from viral sequence data.  
531 *Mol. Biol. Evol.* **29**, 347–357 (2012).

532 40. Kühnert, D., Stadler, T., Vaughan, T. G. & Drummond, A. J. Simultaneous  
533 reconstruction of evolutionary history and epidemiological dynamics from viral  
534 sequences with the birth-death SIR model. *J. R. Soc. Interface* **11**, 20131106 (2014).

535 41. Stadler, T., Kühnert, D., Rasmussen, D. A. & du Plessis, L. Insights into the early  
536 epidemic spread of ebola in sierra leone provided by viral sequence data. *PLoS Curr.* **6**,  
537 (2014).

538 42. Volz, E. & Pond, S. Phylodynamic analysis of ebola virus in the 2014 sierra leone  
539 epidemic. *PLoS Curr.* **6**, (2014).

- 540 43. McCormick, J. B. & Fisher-Hoch, S. P. Lassa fever. *Curr. Top. Microbiol. Immunol.* **262**,  
541 75–109 (2002).
- 542 44. Andersen, K. G. *et al.* Clinical Sequencing Uncovers Origins and Evolution of Lassa  
543 Virus. *Cell* **162**, 738–750 (2015).
- 544 45. Gire, S. K. *et al.* Genomic surveillance elucidates Ebola virus origin and transmission  
545 during the 2014 outbreak. *Science* **345**, 1369–1372 (2014).
- 546 46. Mena, I. *et al.* Origins of the 2009 H1N1 influenza pandemic in swine in Mexico. *Elife* **5**,  
547 (2016).
- 548 47. Morelli, M. J. *et al.* A Bayesian inference framework to reconstruct transmission trees  
549 using epidemiological and genetic data. *PLoS Comput. Biol.* **8**, e1002768 (2012).
- 550 48. Cottam, E. M. *et al.* Integrating genetic and epidemiological data to determine  
551 transmission pathways of foot-and-mouth disease virus. *Proc. Biol. Sci.* **275**, 887–895  
552 (2008).
- 553 49. Cottam, E. M. *et al.* Molecular epidemiology of the foot-and-mouth disease virus  
554 outbreak in the United Kingdom in 2001. *J. Virol.* **80**, 11274–11282 (2006).
- 555 50. Mate, S. E. *et al.* Molecular Evidence of Sexual Transmission of Ebola Virus. *N. Engl. J.*  
556 *Med.* **373**, 2448–2454 (2015).
- 557 51. Blackley, D. J. *et al.* Reduced evolutionary rate in reemerged Ebola virus transmission  
558 chains. *Sci Adv* **2**, e1600378 (2016).
- 559 52. Diallo, B. *et al.* Resurgence of Ebola Virus Disease in Guinea Linked to a Survivor With  
560 Virus Persistence in Seminal Fluid for More Than 500 Days. *Clin. Infect. Dis.* **63**, 1353–  
561 1356 (2016).
- 562 53. Duffy, S., Shackelton, L. A. & Holmes, E. C. Rates of evolutionary change in viruses:  
563 patterns and determinants. *Nat. Rev. Genet.* **9**, 267–276 (2008).
- 564 54. Biek, R., Pybus, O. G., Lloyd-Smith, J. O. & Didelot, X. Measurably evolving pathogens  
565 in the genomic era. *Trends Ecol. Evol.* **30**, 306–313 (2015).
- 566 55. Baele, G., Suchard, M. A., Rambaut, A. & Lemey, P. Emerging Concepts of Data  
567 Integration in Pathogen Phylodynamics. *Syst. Biol.* **66**, e47–e65 (2017).
- 568 56. Campbell, F., Strang, C., Ferguson, N., Cori, A. & Jombart, T. When are pathogen  
569 genome sequences informative of transmission events? *PLoS Pathog.* **14**, e1006885  
570 (2018).
- 571 57. Mate, S. E. *et al.* Molecular Evidence of Sexual Transmission of Ebola Virus. *N. Engl. J.*  
572 *Med.* **373**, 2448–2454 (2015).
- 573 58. Resik, S. *et al.* Limitations to contact tracing and phylogenetic analysis in establishing  
574 HIV type 1 transmission networks in Cuba. *AIDS Res. Hum. Retroviruses* **23**, 347–356  
575 (2007).
- 576 59. Worby, C. J., Lipsitch, M. & Hanage, W. P. Shared Genomic Variants: Identification of  
577 Transmission Routes Using Pathogen Deep-Sequence Data. *Am. J. Epidemiol.* **186**,  
578 1209–1216 (2017).
- 579 60. Faria, N. R., Suchard, M. A., Rambaut, A. & Lemey, P. Toward a quantitative  
580 understanding of viral phylogeography. *Curr. Opin. Virol.* **1**, 423–429 (2011).
- 581 61. Dudas, G. *et al.* Virus genomes reveal factors that spread and sustained the Ebola  
582 epidemic. *Nature* **544**, 309–315 (2017).
- 583 62. Grubaugh, N. D. *et al.* Genomic epidemiology reveals multiple introductions of Zika virus  
584 into the United States. *Nature* **90**, 4864 (2017).
- 585 63. Vaughan, T. G., Kühnert, D., Poppinga, A., Welch, D. & Drummond, A. J. Efficient  
586 Bayesian inference under the structured coalescent. *Bioinformatics* **30**, 2272–2279  
587 (2014).

64. Müller, N. F., Rasmussen, D. A. & Stadler, T. The Structured Coalescent and Its Approximations. *Mol. Biol. Evol.* **34**, 2970–2981 (2017).
65. Kühnert, D., Stadler, T., Vaughan, T. G. & Drummond, A. J. Phylodynamics with Migration: A Computational Framework to Quantify Population Structure from Genomic Data. *Mol. Biol. Evol.* **33**, 2102–2116 (2016).
66. De Maio, N., Wu, C.-H., O'Reilly, K. M. & Wilson, D. New Routes to Phylogeography: A Bayesian Structured Coalescent Approximation. *PLoS Genet.* **11**, e1005421 (2015).
67. Lemey, P. *et al.* Unifying viral genetics and human transportation data to predict the global transmission dynamics of human influenza H3N2. *PLoS Pathog.* **10**, e1003932 (2014).
68. Wang, E. *et al.* Evolutionary relationships of endemic/epidemic and sylvatic dengue viruses. *J. Virol.* **74**, 3227–3234 (2000).
69. Cardoso, J. da C. *et al.* Yellow fever virus in *Haemagogus leucocelaenus* and *Aedes serratus* mosquitoes, southern Brazil, 2008. *Emerg. Infect. Dis.* **16**, 1918–1924 (2010).
70. Faria, N. R. *et al.* Genomic and epidemiological monitoring of yellow fever virus transmission potential. *Science* (2018). doi:10.1126/science.aat7115
71. Diehl, W. E. *et al.* Ebola Virus Glycoprotein with Increased Infectivity Dominated the 2013–2016 Epidemic. *Cell* **167**, 1088–1097.e6 (2016).
72. Urbanowicz, R. A. *et al.* Human Adaptation of Ebola Virus during the West African Outbreak. *Cell* **167**, 1079–1085.e5 (2016).
73. Leroy, E. M. *et al.* Fruit bats as reservoirs of Ebola virus. *Nature* **438**, 575–576 (2005).
74. Walsh, P. D., Biek, R. & Real, L. a. Wave-like spread of Ebola Zaire. *PLoS Biol.* **3**, e371 (2005).
75. Carroll, S. A. *et al.* Molecular Evolution of Viruses of the Family Filoviridae Based on 97 Whole-Genome Sequences. *J. Virol.* **87**, 2608–2616 (2013).
76. Dudas, G. & Rambaut, A. Phylogenetic Analysis of Guinea 2014 EBOV Ebolavirus Outbreak. *PLoS Curr.* **6**, (2014).
77. Rambaut, A. *et al.* Comment on 'Mutation rate and genotype variation of Ebola virus from Mali case sequences'. *Science* **353**, 658 (2016).
78. Lam, T. T.-Y., Zhu, H., Chong, Y. L., Holmes, E. C. & Guan, Y. Puzzling origins of the Ebola outbreak in the Democratic Republic of the Congo, 2014. *J. Virol.* JVI.01226–15 (2015).
79. Blackley, D. J. *et al.* Reduced evolutionary rate in reemerged Ebola virus transmission chains. *Sci Adv* **2**, e1600378 (2016).
80. Yozwiak, N. L. *et al.* Roots, Not Parachutes: Research Collaborations Combat Outbreaks. *Cell* **166**, 5–8 (2016).
81. Yozwiak, N. L., Schaffner, S. F. & Sabeti, P. C. Data sharing: Make outbreak research open access. *Nature* **518**, 477–479 (2015).
82. Policy statement on data sharing by WHO in the context of public health emergencies (as of 13 April 2016. *Wkly. Epidemiol. Rec.* **91**, 237–240 (2016).
83. WHO. WHO R&D Blueprint meeting on pathogen genetic sequence data (GSD) sharing in the context of public health emergencies, 28–29 September 2017. *WHO R&D Blueprint* (2017).
84. Johansson, M. A., Reich, N. G., Meyers, L. A. & Lipsitch, M. Preprints: An underutilized mechanism to accelerate outbreak science. *PLoS Med.* **15**, e1002549 (2018).
85. Callaway, E. Zika-microcephaly paper sparks data-sharing confusion. *Nature News* (2016). doi:10.1038/nature.2016.19367
86. Luksza, M. & Lässig, M. A predictive fitness model for influenza. *Nature* **507**, 57–61

636 (2014).

637 87. Smith, D. J. *et al.* Mapping the antigenic and genetic evolution of influenza virus.  
638 *Science* **305**, 371–376 (2004).

639 88. Neher, R. A., Bedford, T., Daniels, R. S., Russell, C. A. & Shraiman, B. I. Prediction,  
640 dynamics, and visualization of antigenic phenotypes of seasonal influenza viruses. *Proc.*  
641 *Natl. Acad. Sci. U. S. A.* **113**, E1701–9 (2016).

642 89. Osterholm, M. T. *et al.* Transmission of Ebola viruses: what we know and what we do  
643 not know. *MBio* **6**, e00137 (2015).

644 90. Sabir, J. S. M. *et al.* Co-circulation of three camel coronavirus species and  
645 recombination of MERS-CoVs in Saudi Arabia. *Science* **351**, 81–84 (2016).

646 91. Dudas, G., Carvalho, L. M., Rambaut, A. & Bedford, T. MERS-CoV spillover at the  
647 camel-human interface. *Elife* **7**, (2018).

648 92. Faria, N. R. *et al.* Zika virus in the Americas: Early epidemiological and genetic findings.  
649 *Science* **352**, 345–349 (2016).

650 93. Faria, N. R. *et al.* Establishment and cryptic transmission of Zika virus in Brazil and the  
651 Americas. *Nature* (2017). doi:10.1038/nature22401

652 94. Metsky, H. C. *et al.* Zika virus evolution and spread in the Americas. *Nature* **66**, 366  
653 (2017).

654 95. Christie, A. *et al.* Possible sexual transmission of Ebola virus - Liberia, 2015. *MMWR*  
655 *Morb. Mortal. Wkly. Rep.* **64**, 479–481 (2015).

656 96. Whitmer, S. L. M. *et al.* Active Ebola Virus Replication and Heterogeneous Evolutionary  
657 Rates in EVD Survivors. *Cell Rep.* **22**, 1159–1168 (2018).

658 97. Dietzel, E., Schudt, G., Krähling, V., Matrosovich, M. & Becker, S. Functional  
659 Characterization of Adaptive Mutations during the West African Ebola Virus Outbreak. *J.*  
660 *Virol.* **91**, (2017).

661 98. WHO | List of Blueprint priority diseases. (2018).

662 99. Boisen, M. L. *et al.* Field Validation of the ReEBOV Antigen Rapid Test for Point-of-Care  
663 Diagnosis of Ebola Virus Infection. *J. Infect. Dis.* **214**, S203–S209 (2016).

664 100. Broadhurst, M. J. *et al.* ReEBOV Antigen Rapid Test kit for point-of-care and laboratory-  
665 based testing for Ebola virus disease: a field validation study. *Lancet* **386**, 867–874  
666 (2015).

667 101. Chotiwan, N. *et al.* Rapid and specific detection of Asian- and African-lineage Zika  
668 viruses. *Sci. Transl. Med.* **9**, (2017).

669 102. Imai, M. *et al.* Development of H5-RT-LAMP (loop-mediated isothermal amplification)  
670 system for rapid diagnosis of H5 avian influenza virus infection. *Vaccine* **24**, 6679–6682  
671 (2006).

672 103. Hong, T. C. T. *et al.* Development and evaluation of a novel loop-mediated isothermal  
673 amplification method for rapid detection of severe acute respiratory syndrome  
674 coronavirus. *J. Clin. Microbiol.* **42**, 1956–1961 (2004).

675 104. Hattersley, S. M., Greenman, J. & Haswell, S. J. The application of microfluidic devices  
676 for viral diagnosis in developing countries. *Methods Mol. Biol.* **949**, 285–303 (2013).

677 105. Patolsky, F. *et al.* Electrical detection of single viruses. *Proc. Natl. Acad. Sci. U. S. A.*  
678 **101**, 14017–14022 (2004).

679 106. Chen, Y. *et al.* Field-Effect Transistor Biosensor for Rapid Detection of Ebola Antigen.  
680 *Sci. Rep.* **7**, 10974 (2017).

681 107. Afsahi, S. *et al.* Novel graphene-based biosensor for early detection of Zika virus  
682 infection. *Biosens. Bioelectron.* **100**, 85–88 (2018).

683 108. Pardee, K. *et al.* Paper-based synthetic gene networks. *Cell* **159**, 940–954 (2014).

109. Gootenberg, J. S. *et al.* Nucleic acid detection with CRISPR-Cas13a/C2c2. *Science* **356**, 438–442 (2017).
110. Myhrvold, C. *et al.* Field-deployable viral diagnostics using CRISPR-Cas13. *Science* **360**, 444–448 (2018).
111. Gootenberg, J. S. *et al.* Multiplexed and portable nucleic acid detection platform with Cas13, Cas12a, and Csm6. *Science* **360**, 439–444 (2018).
112. Chen, J. S. *et al.* CRISPR-Cas12a target binding unleashes indiscriminate single-stranded DNase activity. *Science* **360**, 436–439 (2018).
113. Gu, W. *et al.* Depletion of Abundant Sequences by Hybridization (DASH): using Cas9 to remove unwanted high-abundance species in sequencing libraries and molecular counting applications. *Genome Biol.* **17**, 41 (2016).
114. Siddle, K. J. *et al.* Capturing diverse microbial sequence with comprehensive and scalable probe design. *bioRxiv* 279570 (2018). doi:10.1101/279570
115. Matranga, C. B. *et al.* Enhanced methods for unbiased deep sequencing of Lassa and Ebola RNA viruses from clinical and biological samples. *Genome Biol.* **15**, 519 (2014).
116. Quick, J. *et al.* Multiplex PCR method for MinION and Illumina sequencing of Zika and other virus genomes directly from clinical samples. *Nat. Protoc.* **12**, 1261–1276 (2017).
117. Quick, J. *et al.* Real-time, portable genome sequencing for Ebola surveillance. *Nature* **530**, 228–232 (2016).