1

# Phylogenetic interpretation during outbreaks requires caution

Ch. Julián Villabona-Arenas[1,2], William P. Hanage[3], Damien C. Tully[1,2]*

[1] Department of Infectious Disease Epidemiology, Faculty of Epidemiology and Population Health, London School of Hygiene and Tropical Medicine, London, UK

[2] Centre for Mathematical Modelling of Infectious Diseases, London School of Hygiene and Tropical Medicine, London, UK

[3] Harvard TH Chan School of Public Health, Boston, Massachusetts, US

*Corresponding author: damien.tully@lshtm.ac.uk

**How viruses are related, and how they have evolved and spread over time, can be investigated using phylogenetics. Here, we set out how genomic analyses should be used during an epidemic and propose that phylogenetic insights from the early stages of an outbreak should heed all the available epidemiological information.**

A goal of genomic epidemiology is to infer epidemiological and emergence dynamics from virus genome sequences obtained over short epidemic timescales [1]. Rapid in situ sequence generation and phylogenetic inference is based on detection of genetic changes in pathogen sequences. But during outbreaks there are many unknowns. The outbreak of coronavirus disease 2019 (COVID-19), which originated in Wuhan, China, was reported in December 2019 [2]. By January 2020, the genome of the causative novel coronavirus, named SARS-CoV-2, had been sequenced and made publicly available [2]. Virus sequences have

26  underpinned development of diagnostics and vaccines and been used to assess patterns of

27  transmission and spread. Although sequence data was used to answer crucial epidemiological

28  questions during the Ebola and Zika outbreaks [3,4], the pace of generation of SARS CoV-2

29  genome data generation is unprecedented and is informing public health policy in real-time.

30  Importantly, it's not only sequences that inform phylogenies, and multiple factors

31  contribute to the outputs including model assumptions, sampling density, timing of sample

32  collection, portion of the viral genome sequenced, quality of sequencing data and the

33  mutation rate of the virus itself. Although it is important to extract as much information as

34  possible from sequence data as outbreaks unfold, it is imperative to bear in mind that the

35  historical relationships of strains (phylogenies) are hypotheses that can be challenged as more

36  data becomes available. Here, we highlight some of the challenges of genomic epidemiology

37  during outbreaks such as SARS-CoV-2 and advise that interpretation of findings from

38  phylogenies needs to assess all epidemiological and supporting information and consider

39  sources of bias.

40  During outbreaks we want to know if cases are linked and if this implies transmission.

41  Most viruses can be separated into strains and if two infections are caused by dissimilar

42  strains one can rule out transmission. The oft-forgotten point is that phylogenies can rule out

43  transmission, but if infections are caused by the same strains or identical viruses it does not

44  decisively prove transmission. During an emerging outbreak, when pathogens have not yet

45  diverged into different strains, phylogenetic information is too weak to hypothesize

46  transmission linkage—which in turn can be used for geographic inference; even if the

47  phylogenetic information is stronger, the same phylogeny is consistent with multiple

48  transmission histories and there may missing links due to incomplete sampling [5].

49  Consequently, we need to combine phylogenetic findings with epidemiological and

50  supporting information such as environmental factors and human air travel data before we

51  draw any immediate conclusions regarding transmission. This was the case with Zika virus in

52  Africa where epidemiological, human mobility and climatic data supported the phylogenetic

53  hypothesis that the outbreak was likely imported from Brazil [6].

54  In the first stage of an outbreak, we can use phylogenetics to discern possible zoonotic

55  sources, as in the case of the 2018 Lassa fever virus outbreak, where phylogenetic patterns

56  indicated independent spillover events from rodent hosts [7]. The crucial observation was that

57  the correct identification of the source of zoonotic transmission relies on the availability of

58    viral genome sequences from potential animal reservoirs. If the source of any virus has not

59    been sampled, it cannot be inferred, because phylogenetic linkage alone does not prove it.

60    This is the reason for uncertainty surrounding the zoonotic source of SARS-CoV-2, because

61    we have limited knowledge about the viral abundance from potential animal reservoirs [8]. The

62    generation of additional viral genome sequences from an outbreak, coupled with virus-

63    specific and epidemiological knowledge, provides insight into whether or not multiple

64    'jumps' occurred from a reservoir that might warrant appropriate control measures. Identical

65    or nearly-identical virus genomes are expected from early transmission chains if a single

66    spillover occurred recently, unless multiple zoonoses originated from the same low-genetic

67    diversity virus pool. In contrast, higher diversity in the early-stage of human-to-human

68    transmission is expected if multiple zoonoses have occurred or if there is significant within-

69    host evolution [9].

70          Geographical inferences (where and when) are feasible as more representative viral

71    genome data—in temporal and spatial scales—becomes available. We can hypothesize the

72    location of common ancestors using ancestral reconstruction methods and infer phylogenies

73    scaled to time, in order to date epidemiological events. Such analyses require a molecular

74    clock, which models how the rate with which mutations accumulate with time, and how this

75    varies across the branches of a phylogeny. However, early in an outbreak there may not be

76    sufficient signal to accurately estimate clock rate. If this is the case, then it might be

77    appropriate to apply an estimate from another closely related virus [10]. If temporal signal is

78    present and a clock rate can be estimated, results need to be reported as credible intervals

79    (instead of point estimates) to account for uncertainty in both the data (incomplete, biased, or

80    improper sampling can lead to misleading phylogenies) and the many aspects of the methods.

81          When investigating the dissemination of an emerging virus the number of sequenced

82    viral genomes may not be representative. Even as the outbreak unfolds, and more genomes

83    are obtained, they only represent a snapshot of the underlying genetic diversity. If

84    phylogenies are considered alone we cannot conclusively assert the geographical origins of

85    the virus—or the extent of community transmission—as we cannot distinguish between local

86    transmission events and multiple introductions of genetically similar viruses, from

87    geographically distinct sources, if one of them has not been sampled. In this way uneven

88    sampling can also lead to misleading conclusions on the geographical source, number of

89    introductions and the size and duration of local transmission chains [11]. The significance of

90    these associations is harder to ascertain when the phylogeny is reported without any

91    assessment on the reliability of internal branches. Therefore, phylogenetic interpretation from

92    ongoing outbreaks as is the case of SARS-CoV-2 needs to be done in the context of all

93    available information such as temporal and spatial distribution of cases, travel patterns and

94    any evidence of epidemiological linkage, sampling uncertainty and other sources of bias need

95    to be carefully considered and reported.

96         The methods for valid phylogenetic inference require multiple assumptions which are

97    likely not met during emerging outbreaks. Examples (not exhaustive) include adequate

98    phylogenetic signal, which is low when strains have not yet diverged; geographical

99    representation and effective sampling time points with sufficient molecular clock signal,

100   which only become feasible as the epidemic unfolds; and random mixing, which may be

101   violated under certain circumstances, for instance when mitigation strategies are set in place.

102   Estimates from phylogenies may be sensitive to one or more of these assumptions and

103   conclusions need to be made and shared with caution. Another essential consideration during

104   an epidemic is accurate rooting of the phylogeny as it determines the direction of

105   transmission over time [12].

106        There are also genome features that are intrinsic to the biology of the virus that may

107   impact the extent and applicability of phylogenetics during outbreaks. For instance, the

108   presence of recombination/reassortment and low diversity (due to the rate of evolution,

109   selective constraints and transmission bottlenecks) complicate the resolution of phylogenetic

110   relationships, but the incorporation of within-host viral diversity may provide greater

111   resolution in understanding transmission dynamics [13]. Moreover, some of mutations in the

112   viral genome sequence can be due to the error rate of the sequencing technology, recurrent

113   sequencing issues, hypermutability or contamination which warrant caution with

114   interpretations and especially with those concerning selection and recombination.

115        Genomic epidemiology has supported public health outbreak responses. Indeed, the

116   ability to exploit viral genome sequences has allowed us to characterise early patterns of

117   SARS-CoV-2 transmission in China, New Zealand and Australia [14,15]. In the midst of an

118   outbreak sharing data is both necessary and important for an effective response, but sharing

119   the associated metadata is also necessary to aid interpretations (e.g. how representative is the

120   data of the country-wide situation) and to avoid creating sampling bias by researchers that are

121   not doing the sequencing themselves.

122    The emergence of SARS-CoV-2 has presented a series of challenges about how we

123    reliably extract information from phylogenies to gain insights into virus transmission and

124    spread, and how we responsibly present our findings. Owing to low genetic diversity and

125    uneven sampling, several controversial hypotheses have already been put forward. One

126    cautionary tale involves how an outbreak in Bavaria seeded the epidemic in northern Italy

127    and the subsequent wider outbreak in Europe. This notion was based on a small sample of

128    very similar sequences. However, it overlooked a more likely scenario in which this virus

129    was already circulating in China and that European regions had multiple introductions from

130    China. At this early stage conclusions about the impacts of mutations on transmission and

131    disease (e.g. D614G mutation in the spike protein [16]) should not be made on the basis of

132    phylogenies alone but with separate evidence supporting not only a phenotypic difference but

133    the resulting consequences for epidemiology.

134    The SARS-CoV-2 pandemic has highlighted the importance of providing a

135    comprehensive rationale for any conclusions about the spatio-temporal dispersal of the virus.

136    Phylogenies represent hypotheses that encompass different sources of error and this

137    uncertainty needs to be visualised and communicated far more transparently. Another

138    challenge is how we facilitate the dissemination of metadata and integrate this with

139    phylogenetic trees. Incorporating host characteristics (e.g. age, onset date, exposure history)

140    to aid phylogenetic interpretation would undoubtedly results in more reliable inferences.

141    Now, more than ever, careful reporting of phylogenetic interpretations, while

142    safeguarding the privacy of infected individuals, would ensure that both policymakers and the

143    public have the best possible information during an outbreak. Failure to balance these issues

144    could jeopardise both scientific integrity and public confidence in the field of genomic

145    epidemiology.

146

147    **REFERENCES**

148    1.    Grubaugh, N. D. et al. Nat Microbiol **4**, 10–19 (2019).

149    2.    Wu, F. et al. Nature **579**, 265–269 (2020).

150    3.    Holmes, E. C., Dudas, G., Rambaut, A. & Andersen, K. G. Nature **538**, 193–200

151    (2016).

152    4.    Pollett, S. et al. J. Infect. Dis. (2019).

153    5.    Hall, M. D. & Colijn, C. Mol. Biol. Evol. **36**, 1333–1343 (2019).

154    6.    Hill, S. C. et al. Lancet Infect. Dis. **19**, 1138–1147 (2019).

155    7.    Kafetzopoulou, L. E. et al. Science (80 ). **363**, 74–77 (2019).

156    8.    Andersen, K. G., Rambaut, A., Ian Lipkin, W., Holmes, E. C. & Garry, R. F. Nat.

157          Med. 1–3 (2020).

158    9.    Didelot, X., Fraser, C., Gardy, J. & Colijn, C. Mol. Biol. Evol. **34**, 997–1007 (2017).

159    10.   Fraser, C. et al. Science (80). **324**, 1557–1561 (2009).

160    11.   Kraemer, M. U. G. et al. Epidemiol. Infect. **147**, (2019).

161    12.   Dudas, G. & Rambaut, A. PLoS Curr. **6**, (2014).

162    13.   Worby, C. J., Lipsitch, M. & Hanage, W. P. Am. J. Epidemiol. **186**, 1209–1216

163          (2017).

164    14.   Lu, J. et al. Cell (2020). doi:10.1016/j.cell.2020.04.023

165    15.   Eden, J.-S. et al. Virus Evol **6**, veaa027 (2020).

166    16.   Korber, B. et al. bioRxiv 2020.04.29.069054 (2020). doi:10.1101/2020.04.29.069054

167    **Contributions**

168    D.C.T. conceived the commentary and wrote the first draft. C.J-V.A, D.C.T conceptualized

169    the ideas with W.P.H. All authors edited the manuscript into its final form.

*170* **Competing Interests**

*171* The authors declare no competing interests.