

1 Shared genomic variants: identification of 2 transmission routes using pathogen deep 3 sequence data

4
5 Colin J. Worby^{1*}, Marc Lipsitch¹, William P. Hanage¹

6
7 ¹Center for Communicable Disease Dynamics, Department of Epidemiology, Harvard TH
8 Chan School of Public Health, 655 Huntington Ave, Boston, MA 02115

9
10 *Corresponding author
11 cworby@hsph.harvard.edu

12 13 14 Abstract

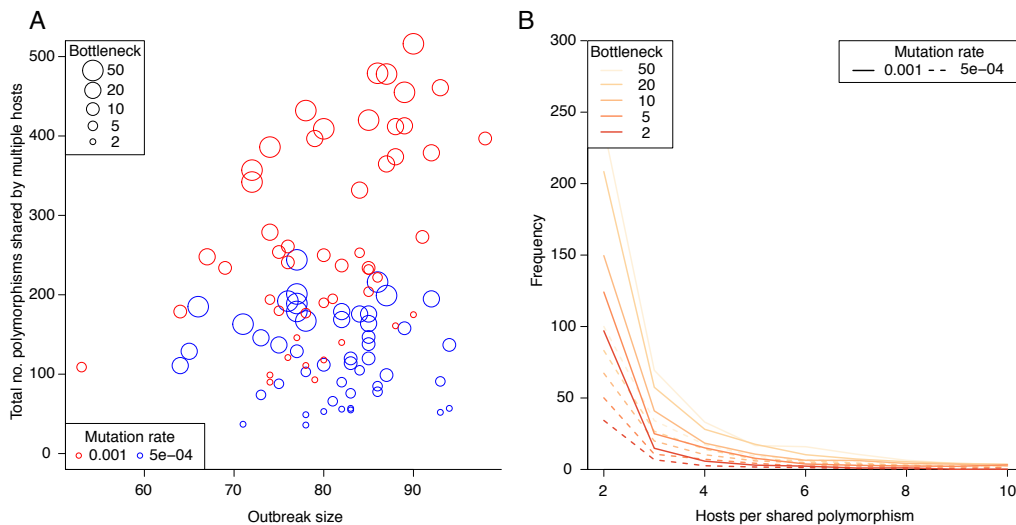
15
16 While identifying routes of transmission during an infectious disease
17 outbreak was traditionally conducted through exhaustive contact tracing
18 efforts, the increasing availability of pathogen sequencing has provided a
19 new resource with which one can identify plausible routes of infection.
20 However, while transmission clusters can be identified using single
21 genome sequences, individual transmission routes remain relatively
22 uncertain. Deep sequence data may provide additional information where
23 single genomes lack sufficient resolution – presence of shared minor
24 variants can suggest epidemiological linkage when observed between
25 multiple hosts. In this study we formalize shared variant methods to
26 reconstruct the transmission tree in an outbreak, and using simulated
27 outbreak data, we quantify the improved accuracy when compared with
28 analogous single genome approaches. Furthermore we propose a hybrid
29 approach, drawing information from both deep sequence and single
30 genome data. Our simulation studies demonstrate the superior
31 performance of transmission tree identification methods using shared
32 variants in most settings. Application of these methods to deep sequence
33 data collected during the 2014 Sierra Leone Ebola epidemic
34 demonstrates the ability to identify plausible transmission routes without
35 any additional data. The methods we describe should become a common
36 step in outbreak investigations and epidemiological analyses once the
37 collection of deep sequence data becomes increasingly widespread.

1 Sequencing pathogen samples during a communicable disease outbreak
2 is becoming an increasingly common procedure in epidemiological
3 investigations. Identifying who infected whom sheds considerable light on
4 transmission patterns, high-risk settings and subpopulations, and
5 infection control effectiveness. Genomic data can shed new light on
6 transmission dynamics, and can be used to identify clusters of individuals
7 likely to be linked by direct transmission. However, identification of
8 individual sources of infection typically remains uncertain. In this study,
9 we investigate the potential of deep sequence data to provide greater
10 resolution on transmission routes. We describe easily implemented
11 methods to use such data, and demonstrate the remarkably improved
12 performance when reconstructing transmission trees. Furthermore, we
13 apply our methods to data collected during the 2014 Ebola outbreak in
14 Sierra Leone, identifying several routes of transmission. Our study
15 highlights the power of pathogen deep sequence data as a component of
16 outbreak investigation and epidemiological analyses.

17 **Introduction**

19 Genomic data offer new insights into epidemiological and evolutionary
20 dynamics, and sequencing pathogen samples is becoming increasingly
21 widespread. Pathogen genomic data allows one to link the evolutionary
22 relatedness of collected isolates, which in turn sheds light on the potential
23 relationship between the hosts from whom they were collected. As such,
24 inference of transmission trees using genomic data is an increasingly
25 well-studied field (1-7). While low-resolution pathogen typing has been
26 used for some time to discriminate between independent outbreaks (8-
27 10), whole genome sequencing provides additional resolution with which
28 genetic distance between identical phenotypes may be ascertained (11-
29 13). This too, however, has limits. Studies have shown that while
30 transmission clusters may be identified with genomic data, individual-level
31 transmission routes can rarely be identified with a great degree of
32 certainty (2, 3). Characterizing an infected host by a single pathogen
33 genome (isolation and purification of a single genotype for bacteria, or
34 using the consensus sequence for viral pathogens) is common practice,
35 yet neglects within-host diversity. The variation in sampled genetic
36 distances can be large relative to the expected number of mutations
37 between hosts, rendering the number of SNPs a rather crude measure of
38 relatedness on an individual level (14). As such, particularly for rapidly
39 evolving pathogens, or those whose mode of transmission is associated
40 with a large and potentially diverse inoculum ('transmission bottleneck'),
41 single genome sampling can cause hosts to appear misleadingly similar
42 or dissimilar.
43
44

1 Deep sequencing can potentially provide new insights into within-host
2 diversity. Currently, sequencing a mixed population sample to sufficient
3 depth to identify minor variants has mostly been limited to viral samples.
4 Recent studies have demonstrated the utility of such data in exploring
5 evolutionary dynamics of communicable pathogens during an ongoing
6 outbreak (15-17). While the consensus sequences may appear identical
7 for two samples, comparing minor variants can offer additional resolution.
8 For instance, the presence of one or more shared variants could be
9 considered as strong evidence for direct transmission, particularly if the
10 variant is not observed in any other host. This naturally relies on the
11 possibility that a pathogen population of size greater than one survives
12 the transmission bottleneck; otherwise each infection must initially be
13 monoclonal. While estimating transmission bottleneck size is challenging,
14 there is some evidence for larger bottlenecks occurring in influenza (18,
15 19), Ebola (20), as well as speculation for much variation in bottleneck
16 size for bacterial pathogens (21).
17 In this study, we investigated the predictive power of shared variants for
18 identifying transmission routes. We compared several simple approaches
19 to reconstruct simulated transmission trees in order to quantify the
20 difference in performance of those based on single genome samples vs.
21 those based on deep sequence data. In addition, we describe hybrid
22 methods which draw power from shared variant identification as well as
23 genetic distance metrics. We then demonstrated an application of the
24 shared variant methods on deep sequence data collected during the
25 Ebola outbreak in West Africa in 2014.



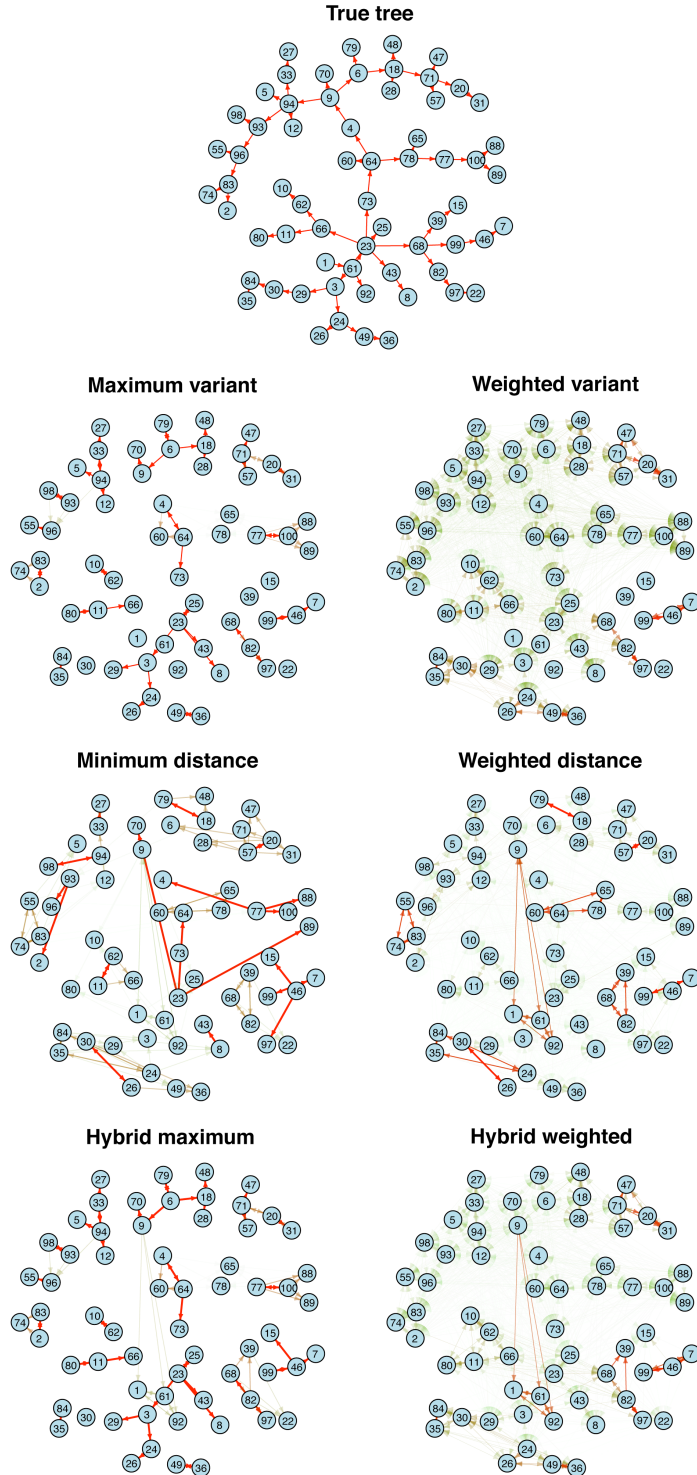
26 **Figure 1. Summary of genetic variant frequency across the simulated**
27 **outbreaks.** (A) Total number of shared variants across the simulated
28 **outbreaks.** Bottleneck size is illustrated by circle size. (B) Distribution of
29 **shared variant group size for different bottlenecks and mutations rates.**
30

31

1 Results

2 3 Simulation Studies

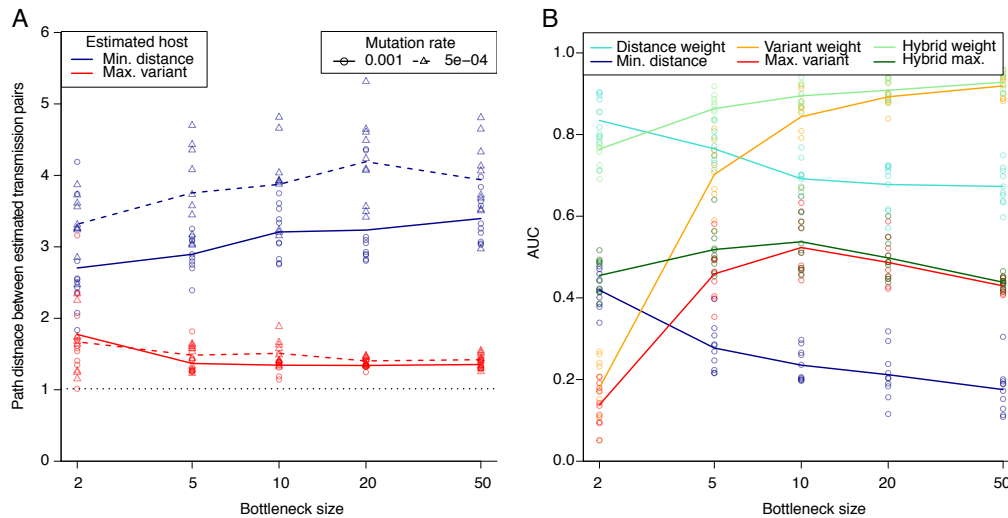
4
5 We simulated a range of outbreaks with final size at least 50, allowing
6 both the transmission bottleneck size and mutation rate to vary. As
7 expected, shared variants were observed with increasing frequency when
8 mutation rates and bottleneck sizes were larger (Figure 1A). The vast
9 majority of shared variants were observed in exactly two individuals, with
10 a rapidly declining frequency for larger group sizes (Figure 1B). For each
11 simulation, we constructed a weighted transmission tree according to the
12 six methods outlined in the Methods section. An example simulation and
13 the reconstructions are shown in figure 2. We plotted the ROC curve for
14 each methodology (Figure S1) to determine the AUC statistic. For small
15 bottlenecks, variant-based methods provide a poor tree reconstruction by
16 the AUC measure (Figure 3B); values below 0.5 indicate a worse
17 performance than random selection, an inevitability when only a small
18 proportion of nodes are assigned sources. A tight bottleneck leads to little
19 diversity persisting across transmission events, and as such, shared
20 variants are rarely observed, leading to a sparsity of informed links across
21 the network. However, the links which are made are typically more
22 reliable than those estimated by minimum genetic distance, with an
23 average path length of less than 2 (Figure 3A). Larger bottlenecks lead to
24 rapidly improving AUC statistics for the variant-based approaches,
25 although the maximum variant approach declines slightly for bottleneck
26 size greater than 10. In contrast, distance-based approaches typically
27 decline in accuracy as the bottleneck size increases.



1
2
3
4
5
6
7

Figure 2. Simulated and reconstructed transmission trees. The simulated tree (top) was generated with bottleneck size 10 and mutation rate 0.001. Trees were reconstructed according to the approaches described in *Methods*. Edges are colored and weighted according to the weight attributed to that potential transmission route, red edges associated with a higher weight than green. Networks were plotted with the igraph package in R.

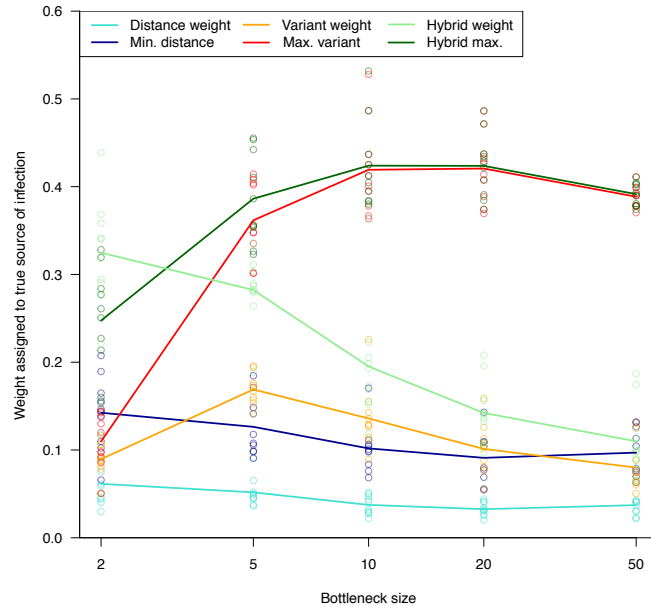
1 The hybrid methods draw additional information from the genetic distance
2 information when no shared variants are found for a given host. This
3 offers a considerable improvement in performance for outbreaks with
4 smaller bottleneck sizes (Figure 3B), which are the situations where
5 distance information is most reliable.
6



7
8 **Figure 3. Transmission tree reconstruction accuracy.** (A) The true path
9 distance between estimated transmission pairs gives insight into the extent
10 to which transmission links are misspecified. A perfect reconstruction would
11 have mean path length 1. Maximum variant path lengths are averaged over
12 identified transmission pairs, that is, excluding hosts with no shared variants.
13 (B) The area under the ROC (AUC) metric provides an overall measure of
14 network accuracy. Results for a mutation rate of 0.001 are shown here.

15
16 For each simulation, we calculated the mean weight attributed to the true
17 source of each host. The distance and variant approaches perform
18 similarly for bottleneck size 2, but are outperformed by the two hybrid
19 approaches (Figure 4). For larger bottleneck sizes, the maximum variant
20 approach (together with the hybrid maximum) is the best method by a
21 large margin according to this measure. The distance and variant weight
22 methods perform increasingly poorly as the bottleneck size increases,
23 reflecting the greater likelihood of genotypes persisting across several
24 bottlenecks, increasing uncertainty.

25
26 We found that lower mutation rates were typically associated with poorer
27 tree reconstruction, although the relative performance of the methods
28 remained broadly similar (Figures 3A, S2 and S3). Lower mutation rates
29 lead to fewer shared variants, and as such, the distance-based
30 approaches may outperform variant-based approaches for larger
31 bottlenecks. Nevertheless, the hybrid approach consistently offers the
32 most successful reconstruction approach over a range of parameters.



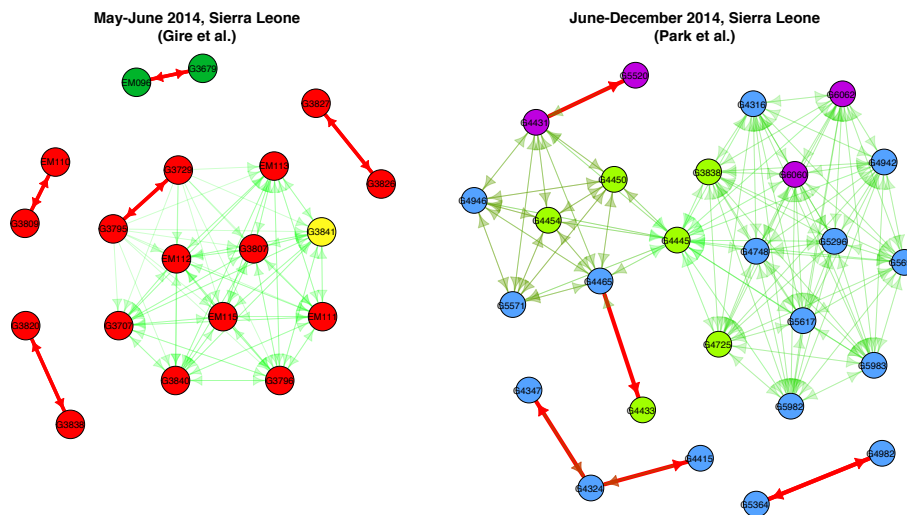
1
2 **Figure 4. True edge weighting.** The mean weight attributed to each true
3 transmission link for each tree reconstruction method, under a range of
4 scenarios and methodologies. Results are shown for a mutation rate of
5 0.001.

6 7 **Ebola virus data**

8
9 We next examined the Sierra Leone Ebola datasets, attempting to identify
10 transmission links between hosts. In order to reduce the risk of counting
11 variant calling errors as true intra-host variants, we identified only variants
12 in which the minor frequency was at least 5% (routes estimated under a
13 1% threshold are shown in Figure S4). Figure 5 shows the transmission
14 trees reconstructed for each dataset. In the first dataset (Figure 5A, (16))
15 a total of 19 out of 78 hosts were found to have shared a variant with at
16 least one other individual. Four pairs of patients shared more than one
17 variant (three pairs with two shared variants, and one pair with four), while
18 one additional pair shared one unique variant. Each of these pairs
19 originated from the same geographic location, and was temporally
20 clustered (three of these links were sampled two or fewer days apart,
21 while the remaining two were sampled 12 and 22 days apart, still
22 plausible given the serial interval estimates of 15.3 ± 9.3 days (22). One
23 variant was shared by 11 hosts. This might be explained by alternative
24 hypotheses; (i) a well-sampled transmission cluster with large
25 transmission bottleneck size, (ii) repeated emergence of identical
26 mutation. The group shows high geographic and temporal clustering (15),
27 suggesting that a transmission cluster may be plausible. Non-random
28 mixing, or presence of a superspreader (for instance, due to several
29 persons caring for a patient, or washing and burying a deceased patient)

1 could make the transmission scenario much more likely. Furthermore,
2 non-random sampling could skew the observed distribution of group size
3 for shared variants.

4
5 A total of 26 hosts out of 150 (for which replicate sequencing and variant
6 calling was performed) in the second dataset shared a variant at with at
7 least one other host (Figure 5B, (16)). No pair of hosts shared more than
8 one variant, and there were five pairs of individuals sharing a unique
9 variant. One variant was shared by 13 hosts at the 5% detection
10 threshold. Unlike the previous dataset, these hosts were not
11 geographically or temporally clustered, coming from different villages and
12 spanning several weeks. This suggests that a transmission cluster is less
13 likely than in the previous example, and potentially reflects homoplasy.
14 Four of the five ‘unambiguous’ transmission routes joined patients from
15 the same geographic location.



17
18 **Figure 5. Estimated Ebola transmission routes.** Transmission links
19 between sampled hosts in the 2014 Ebola outbreak under the maximum
20 variant approach. Colors denote chiefdoms to which hosts belong, while the
21 color and thickness of the arrows denote the relative weight attributed to
22 each potential transmission event. Variant detection threshold 5%.

23 Discussion

24
25
26 Shared variant detection offers a powerful insight into unobserved
27 transmission dynamics, and can improve the resolution of reconstructed
28 transmission trees considerably. We have described some simple
29 methodology for reconstructing transmission trees using deep sequence
30 data. These methods typically outperform genetic distance comparison
31 methods, frequently used to identify potential transmission events [cite

1 MRSA etc.]. While more formal approaches may utilize deep sequence
2 data even more successfully, it is likely that such approaches will require a
3 model specifying within-host pathogen dynamics, which are still poorly
4 understood – estimation of effective population size within host, *in vivo*
5 pathogen generation time and transmission bottleneck size are all
6 challenging. Furthermore, the methods described here are simple and
7 quick to implement.

8
9 We applied these methods to Ebola data collected from Sierra Leone in
10 2014. While the first dataset is thought to represent relatively dense
11 coverage of the initial stages of the epidemic in the country, with estimates
12 of around 70% of cases sampled (15, 23), the later dataset comprised a
13 sparser sample. While sparse sampling reduces the number of true links
14 one would expect to find via any method, shared variant approaches do
15 not require 100% sampling. Missing data greatly reduces the number of
16 shared variants expected, and therefore the number of estimated
17 transmission routes, but does not greatly impact the reliability of the
18 transmission routes which are proposed (Figures S5 and S6). As such,
19 while only relatively few transmission routes were identified in the
20 datasets, this is likely a function of both the proportion of missing data, and
21 the relatively low mutation rate of Ebola virus (15, 24). However, the
22 transmission pairs we identified appear plausible from temporal and
23 geographic clustering. Both datasets contained a large group of hosts
24 sharing the same variant. Park et al. suggest that the large group in the
25 second dataset likely arose through a combination of patient-to-patient
26 transmission and recurrent mutation (16). Testing for recurrent mutation,
27 particularly in larger variant-sharing groups, as well as cross-checking
28 transmission links against other data sources, is therefore highly
29 recommended.

30
31 The hybrid approaches allow nodes with no shared variants to be
32 connected to other nodes in the transmission tree. This provides
33 considerable improvements when the bottleneck is small or the mutation
34 rate is low. The hybrid approach could be performed either by using
35 separate single genomes, or the consensus sequence of the deep
36 sequenced sample. Since transmission routes are assessed
37 independently of one another, estimated transmission trees frequently
38 comprise several unconnected nodes or clusters. Such unconnected
39 clusters could be linked to one another if further structure is required, using
40 the weighted distance approach on pooled within-cluster samples.

41
42 Care should be taken when identifying variant sites to minimize the risk of
43 calling sequencing or alignment errors as true variants. Deep coverage
44 and replicate sequencing provide some reassurance of variant calling

1 quality; nevertheless, we found considerable sensitivity to adjusting the
2 variant frequency threshold, particularly for the second, lower coverage,
3 dataset (Figure S4). Setting a conservatively high threshold increases the
4 probability of calling true variants, but reduces the amount of useful
5 information for transmission tree estimation. While the cost of sequencing
6 bacterial pathogen populations to a sufficient depth to call minor variants
7 remains restrictive, this is likely to reduce in the future.

8
9 We have demonstrated the power of deep sequencing data to identify
10 transmission routes to a greater resolution than by using analogous
11 methods with single genome sequence data. We have purposefully
12 omitted the incorporation of additional data sources (such as times of
13 sampling, symptom onset, recovery/death, as well as geographic location
14 and contact tracing) in order to evaluate the information provided by the
15 genomic data alone. Incorporating additional data sources will improve
16 estimates, allowing further potential transmission links to be ruled out.
17 Since we make no assumption about the order of infection, directionality is
18 ambiguous in the majority of estimated transmission links. Rigorous
19 collection of epidemiological data remains a crucial component of outbreak
20 investigation, and combining this with deep sequencing and shared variant
21 analysis can provide unprecedented insight into individual-level
22 transmission dynamics. We believe that the shared variant methods
23 described here will become common once deep sequence data collection
24 becomes more widespread, and should provide a considerably clearer
25 picture of who infected whom than single genome sampling data.

26 27 **Materials and Methods**

28
29 Let x_1, \dots, x_n denote deep-sequence samples collected from hosts $1, \dots, n$.

30 For each sample x_i , let $f_1^{(i)}, \dots, f_G^{(i)}$ be the frequency of the majority
31 nucleotide at loci $1, \dots, G$, such that polymorphisms exist where $f_j^{(i)} < 1$.

32 For each host, identify the set of polymorphisms $V_i = \{j : f_j^{(i)} < 1\}$. Now
33 calculate the variant score S_{ij} between each pair of hosts i and j to be
34 the number of shared variants belonging to the samples x_i and x_j ;

35 $S_{ij} = |V_i \cap V_j|$. If we allow for the possibility of different mutations at a given
36 locus, we must further restrict to the set of variant positions sharing the
37 same mutant nucleotide. The matrix (S_{ij}) can then be viewed as a
38 weighted adjacency matrix defining an estimated transmission tree (which
39 we call the *weighted variant tree*). We further define the *maximum variant*
40 *tree*, in which we identify for each host the individual sharing the greatest
41 number of variants. If multiple individuals share the maximum number of

1 variants, these are attributed equal weight. If no individual shares any
2 variant with a host, it is not assigned a source.

3
4 These approaches have similarities with methods using single genome
5 samples. For instance, the *minimum distance tree* is defined by assigning
6 the source to be the individual carrying the most genetically similar
7 sample (i.e. fewest number of SNPs). A variation of this approach was
8 described in (6). Similarly, the *weighted distance tree* is defined by
9 weighting each network edge by inverse genetic distance, such that more
10 similar samples are given a greater weight. Variations of this method
11 were explored in (2).

12
13 Finally we propose two hybrid approaches. In some cases, a host will
14 share no variants with any other host in the population, such that the
15 maximum and weighted variant approaches assign equal weight to all
16 potential sources of infection. As such, we may instead attempt to draw
17 information from genetic distance measures where no shared variants
18 exist. The *hybrid maximum tree* and the *hybrid weighted tree* attribute
19 sources to hosts lacking shared variants according to the minimum
20 distance and weighted distance approaches respectively. Genetic
21 distance can be calculated using the consensus sequence of the deep
22 sequenced sample, or with additional single genome samples.

23 24 **Simulations**

25
26 We simulated outbreaks using the R package *seedy* v1.2 (25),
27 introducing a single infectious individual into a susceptible,
28 homogeneously mixing population of size 100. Genomic samples
29 (perfectly observed deep sequence samples) were generated at a
30 random time during each individual's infectious period. Furthermore, we
31 sampled single genomes from each host in order to compare
32 transmission route estimation using each type of data. Infection dynamics
33 were simulated under a standard SIR (susceptible-infected-removed)
34 model, with $R_0 = 2$. Multiple infections were not permitted. Infections were
35 generated by selecting n_B genotypes at random from the source's
36 pathogen population, and allowing this inoculum to grow within the new
37 host under neutral evolution. We varied transmission bottleneck size n_B
38 as well as mutation rate in order to simulate a range of different
39 outbreaks. Transmission trees were visualized using the *igraph* package
40 in R (26).

1 **Measuring reconstruction accuracy**

2
3 We used various metrics to compare the performances of the
4 transmission route identification methods. We assumed that infection and
5 removal times were not observed from the simulated outbreaks, such that
6 we could compare the ability of the genomic data alone to contribute to
7 transmission route identification.

8
9 The receiver operating characteristic (ROC) curve describes the change
10 in false positive and true positive rate for identifying a source of infection
11 as the weight threshold for this identification varies between 0 and 1. The
12 area under the ROC curve (AUC) is a summary statistic of this function,
13 measuring the overall discriminatory power of the tree reconstruction, in
14 which values closer to 1 indicate a more accurate network (27).

15 For the unweighted reconstructions (*minimum distance* and *maximum*
16 *similarity tree*), we calculated the path distance in the true network for
17 each proposed transmission link. For instance, if we identify the route A-
18 B, and in reality the transmission chain was A-C-B, the path distance in
19 the true network is 2. A perfect reconstruction would thus have a mean
20 path distance of 1. While the ROC curve treats edges as either correct or
21 incorrect, the latter metric provides a measure of the extent to which false
22 links are misleading (i.e. an incorrect edge with a true path length of 2 is
23 better than an incorrect edge with a true path length of 10).

24
25 Finally, we considered the mean weight attributed to the true source of
26 infection across an outbreak. While the path distance metric did not factor
27 in hosts for whom no source could be attributed (due to a lack of shared
28 variants), this measure includes such hosts with a weight of zero.

29 30 **Data**

31
32 In addition to the simulation studies, we applied the shared variant
33 approach to identifying potential transmission routes during the 2014
34 Ebola outbreak in Sierra Leone. We used samples collected from 78
35 patients in May-June 2014, representing a large proportion of the earliest
36 cases in the country, sequenced to approximately 2000x coverage (15).
37 Furthermore, we considered samples from a further 150 patients
38 collected between June and December in the same country, sequenced
39 with a median coverage of 374x (16).

1 Funding Information

2
3 Research reported in this paper was supported by the National Institute of
4 General Medical Sciences of the National Institutes of Health under award
5 number U54GM088558. The content is solely the responsibility of the
6 authors and does not necessarily represent the official views of the
7 National Institute of General Medical Sciences or the National Institutes of
8 Health. The funders had no role in study design, data collection and
9 analysis, decision to publish, or preparation of the manuscript.

10 11 References

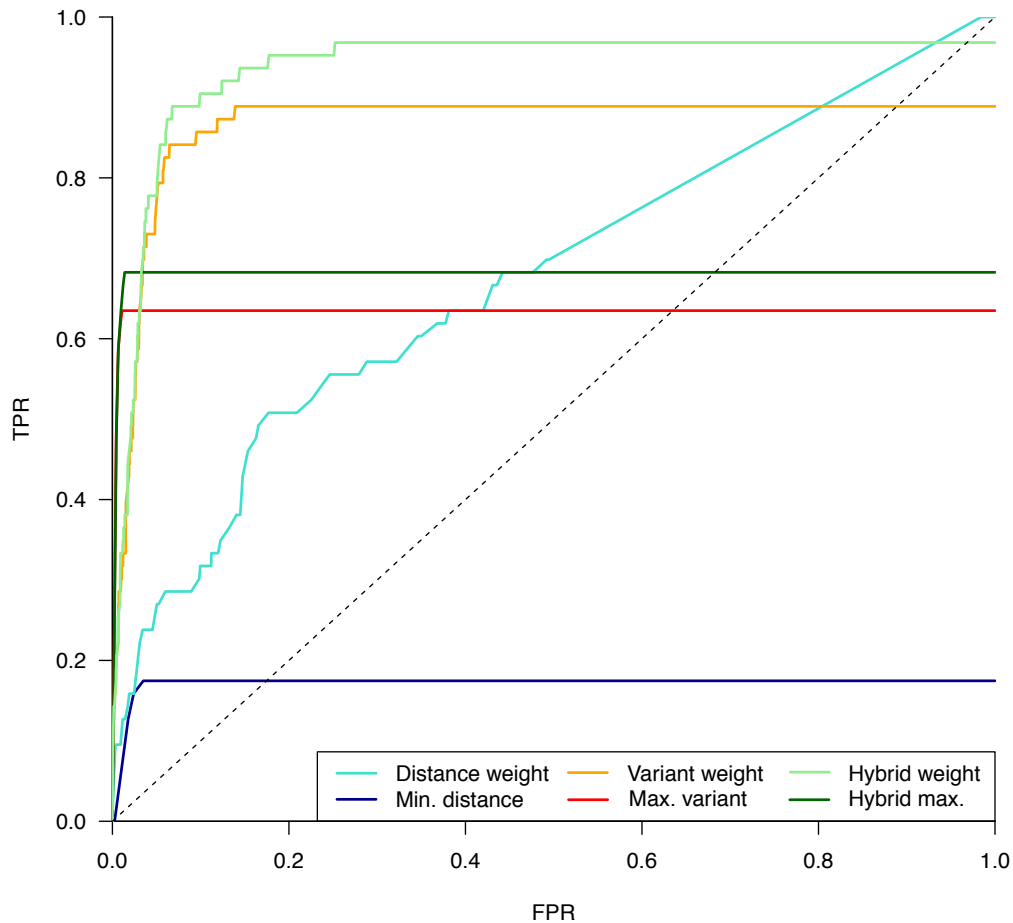
- 12
13 1. **Cottam EM, Thébaud G, Wadsworth J, Gloster J, Mansley L, Paton**
14 **DJ, King DP, Haydon DT.** 2008. Integrating genetic and
15 epidemiological data to determine transmission pathways of foot-
16 and-mouth disease virus. *Proc R Soc B* **275**:887-895.
- 17 2. **Worby CJ, Lipsitch M, Hanage WP.** 2014. Within-Host Bacterial
18 Diversity Hinders Accurate Reconstruction of Transmission
19 Networks from Genomic Distance Data. *PLoS Comp Biol*
20 **10**:e1003549.
- 21 3. **Didelot X, Gardy J, Colijn C.** 2014. Bayesian analysis of infectious
22 disease transmission from whole genome sequence data. *Mol Biol*
23 *Evol* **31**:1869-1879.
- 24 4. **Ypma RJF, Bataille AMA, Stegeman A, Koch G, Wallinga J, van**
25 **Ballegooijen WM.** 2012. Unravelling transmission trees of infectious
26 diseases by combining genetic and epidemiological data. *Proc R Soc B*
27 **279**:444-450.
- 28 5. **Ypma RJF, van Ballegooijen WM, Wallinga J.** 2013. Relating
29 phylogenetic trees to transmission trees of infectious disease
30 outbreaks. *Genetics* **195**:1055-1062.
- 31 6. **Jombart T, Eggo RM, Dodd PJ, Balloux F.** 2011. Reconstructing
32 disease outbreaks from genetic data: a graph approach. *Heredity*
33 **106**:383-390.
- 34 7. **Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N.**
35 2014. Bayesian Reconstruction of Disease Outbreaks by Combining
36 Epidemiologic and Genomic Data. *PLoS Comp Biol* **10**:e1003457.
- 37 8. **Streulens MJ, Deplano A, Godard C, Maes N, Serruys E.** 1992.
38 Epidemiologic typing and delineation of genetic relatedness of
39 methicillin-resistant *Staphylococcus aureus* by macrorestriction
40 analysis of genomic DNA by using pulsed-field gel electrophoresis. *J*
41 *Clin Microbiol* **30**:2599-2605.
- 42 9. **Strommenger B, Braulke C, Heuck D, Schmidt C, Pasemann B,**
43 **Nübel U.** 2008. spa Typing of *Staphylococcus aureus* as a Frontline
44 Tool in Epidemiological Typing. *J Clin Microbiol* **46**:574-581.

- 1 10. **Koreen L, Ramaswamy SV, Graviss EA, Naidich S, Musser JM,**
2 **Kreiswirth BN.** 2004. spa Typing Method for Discriminating among
3 Staphylococcus aureus Isolates: Implications for Use of a Single
4 Marker To Detect Genetic Micro- and Macrovariation. *J Clin Microbiol*
5 **42:792-799.**
- 6 11. **Gardy JL, Johnston JC, Ho Sui SJ, Cook VJ, Shah L, Brodtkin E,**
7 **Rempel S, Moore R, Zhao Y, Holt R, Varhol R, Birol I, Lem M,**
8 **Sharma MK, Elwood K, Jones SJM, Brinkman FSL, Brunham RC,**
9 **Tang P.** 2011. Whole-Genome Sequencing and Social-Network
10 Analysis of a Tuberculosis Outbreak. *New Engl J Med* **364:730-739.**
- 11 12. **Bryant JM, Schürch AC, van Deutekom H, Harris SR, de Beer JL,**
12 **de Jager V, Kremer K, van Hijum SAFT, Siezen RJ, Borgdorff M,**
13 **Bentley SD, Parkhill J, van Soolingen D.** 2013. Inferring patient to
14 patient transmission of Mycobacterium tuberculosis from whole
15 genome sequencing data. *BMC Infect Dis* **13:110.**
- 16 13. **Walker TM, Lalor MK, Broda A, Ortega LS, Morgan M, Parker L,**
17 **Churchill S, Bennett K, Golubchik T, Giess AP, Del Ojo Elias C,**
18 **Jeffery KJ, Bowler ICJW, Laurenson IF, Barrett A, Drobniewski F,**
19 **McCarthy ND, Anderson LF, Abubakar I, Thomas HL, Monk P,**
20 **Smith EG, Walker AS, Crook DW, Peto TEA, Conlon CP.** 2014.
21 Assessment of Mycobacterium tuberculosis transmission in
22 Oxfordshire, UK, 2007–12, with whole pathogen genome sequences:
23 an observational study. *The Lancet Respiratory Medicine* **2:285-292.**
- 24 14. **Worby CJ, Chang H-H, Hanage WP, Lipsitch M.** 2014. The
25 distribution of pairwise genetic distances: a tool for investigating
26 disease transmission. *Genetics* **198:1395-1404.**
- 27 15. **Gire SK, Goba A, Andersen KG, Sealfon RSG, Park DJ, Kanneh L,**
28 **Jalloh S, Momoh M, Fullah M, Dudas G, Wohl S, Moses LM,**
29 **Yozwiak NL, Winnicki S, Matranga CB, Malboeuf CM, Qu J,**
30 **Gladden AD, Schaffner SF, Yang X, Jiang P-P, Nekoui M, Colubri A,**
31 **Coomber MR, Fonnies M, Moigboi A, Gbakie M, Kamara FK,**
32 **Tucker V, Konuwa E, Saffa S, Sellu J, Jalloh AA, Kovoma A,**
33 **Koninga J, Mustapha I, Kargbo K, Foday M, Yillah M, Kanneh F,**
34 **Robert W, Massally JLB, Chapman SB, Bochicchio J, Murphy C,**
35 **Nusbaum C, Young S, Birren BW, Grant DS, Scheffelin JS, et al.**
36 2014. Genomic surveillance elucidates Ebola virus origin and
37 transmission during the 2014 outbreak. *Science* **345:1369-1372.**
- 38 16. **Park DJ, Dudas G, Wohl S, Goba A, Whitmer SLM, Andersen KG,**
39 **Sealfon RS, Ladner JT, Kugelman JR, Matranga CB, Winnicki SM,**
40 **Qu J, Gire SK, Gladden-Young A, Jalloh S, Nosamiefan D, Yozwiak**
41 **NL, Moses LM, Jiang P-P, Lin AE, Schaffner SF, Bird B, Towner J,**
42 **Mamoh M, Gbakie M, Kanneh L, Kargbo D, Massally JLB, Kamara**
43 **FK, Konuwa E, Sellu J, Jalloh AA, Mustapha I, Foday M, Yillah M,**
44 **Erickson BR, Sealy T, Blau D, Paddock C, Brault A, Amman B,**
45 **Basile J, Bearden S, Belser J, Bergeron E, Campbell S, Chakrabarti**
46 **A, Dodd K, Flint M, Gibbons A, et al.** 2015. Ebola Virus

- 1 Epidemiology, Transmission, and Evolution during Seven Months in
2 Sierra Leone. *Cell* **161**:1516-1526.
- 3 17. **Poon LLM, Chan KH, Chu DKW, Fung CCF, Cheung CKY, Ip DKM,**
4 **Leung GM, Peiris JSM, Cowling BJ.** 2011. Viral genetic sequence
5 variations in pandemic H1N1/2009 and seasonal H3N2 influenza
6 viruses within an individual, a household and a community. *J Clin*
7 *Viro* **52**:146-150.
- 8 18. **Hughes J, Allen RC, Baguelin M, Hampson K, Baillie GJ, Elton D,**
9 **Newton JR, Kellam P, Wood JLN, Holmes EC, Murcia PR.** 2012.
10 Transmission of Equine Influenza Virus during an Outbreak Is
11 Characterized by Frequent Mixed Infections and Loose Transmission
12 Bottlenecks *PLoS Path* **8**:e1003081.
- 13 19. **Murcia PR, Hughes J, Battista P, Lloyd L, Baillie GJ, Ramirez-**
14 **Gonzalez RH, Ormond D, Oliver K, Elton D, Mumford JA, Caccamo**
15 **M, Kellam P, Grenfell BT, Holmes EC, Wood JLN.** 2012. Evolution of
16 an Eurasian Avian-like Influenza Virus in Naïve and Vaccinated Pigs
17 *PLoS Path* **8**:e1002730.
- 18 20. **Emmett KJ, Lee A, Khiabani H, Rabadan R.** 2015. High-
19 resolution Genomic Surveillance of 2014 Ebola virus Using Shared
20 Subclonal Variants. *PLOS Currents Outbreaks*
21 doi:10.1371/currents.outbreaks.c7fd7946ba606c982668a96bcba43
22 c90.
- 23 21. **Balloux F.** 2010. Demographic influences on bacterial population
24 structure. In Robinson DA, Falush D, Feil EJ (ed), *Bacterial Population*
25 *Genetics in Infectious Diseases.* John Wiley & Sons Inc.
- 26 22. **WHO Ebola Response Team.** 2014. Ebola Virus Disease in West
27 Africa - The First 9 Months of the Epidemic and Forward Projections.
28 *New Engl J Med* **371**:1481-1495.
- 29 23. **Stadler T, Kühnert D, Rasmussen DA, du Plessis L.** 2014. Insights
30 into the early epidemic spread of Ebola in Sierra Leone provided by
31 viral sequencing. *PLOS Currents Outbreaks*
32 doi:10.1371/currents.outbreaks.02bc6d927ecee7bbd33532ec8ba6a
33 25f.
- 34 24. **Hoenen T, Safronetz D, Groseth A, Wollenberg KR, Koita OA,**
35 **Diarra B, Fall IS, Haidara FC, Diallo F, Sanogo M, Sarro YS, Kone**
36 **A, Togo ACG, Traore A, Kodio M, Dosseh A, Rosenke K, de Wit E,**
37 **Feldmann F, Ebihara H, Munster VJ, Zoon KC, Feldmann H, Sow S.**
38 2015. Mutation rate and genotype variation of Ebola virus from Mali
39 case sequences. *Science* **348**:117-119.
- 40 25. **Worby CJ, Read TD.** 2015. 'SEEDY' (Simulation of Evolutionary and
41 Epidemiological Dynamics): An R Package to Follow Accumulation of
42 Within-Host Mutation in Pathogens *PLoS One* **10**:e0129745.
- 43 26. **Csardi G, Nepusz T.** 2006. The igraph software package for complex
44 network research. *InterJournal Complex Systems*:1695.
- 45 27. **Fawcett T.** 2006. An introduction to ROC analysis. *Pattern Recog Lett*
46 **27**:861-874.

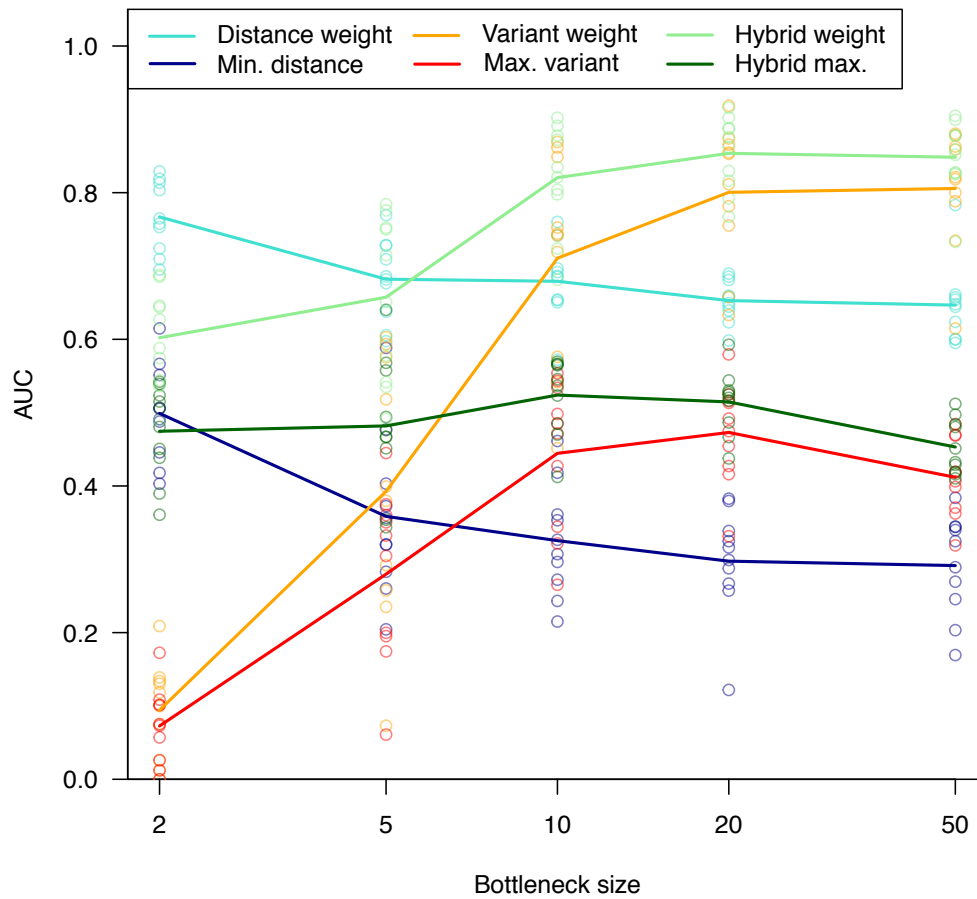
1
2
3

Supplementary Figures

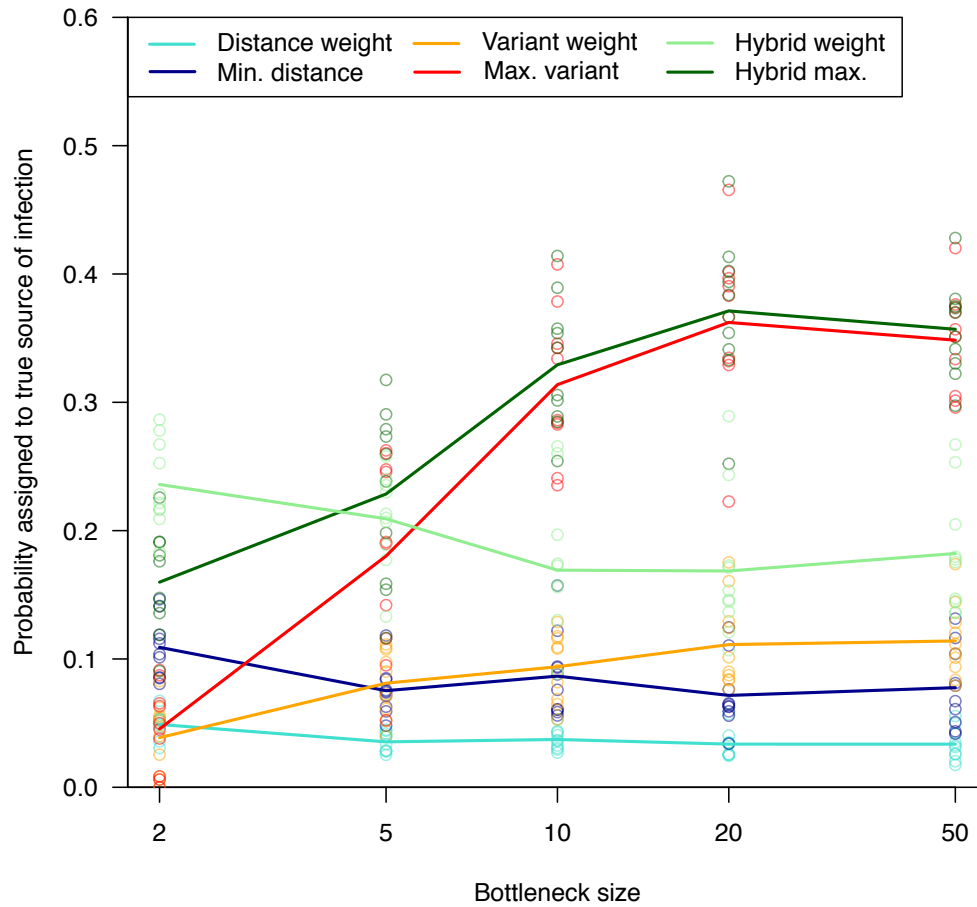


4
5
6
7
8
9
10
11

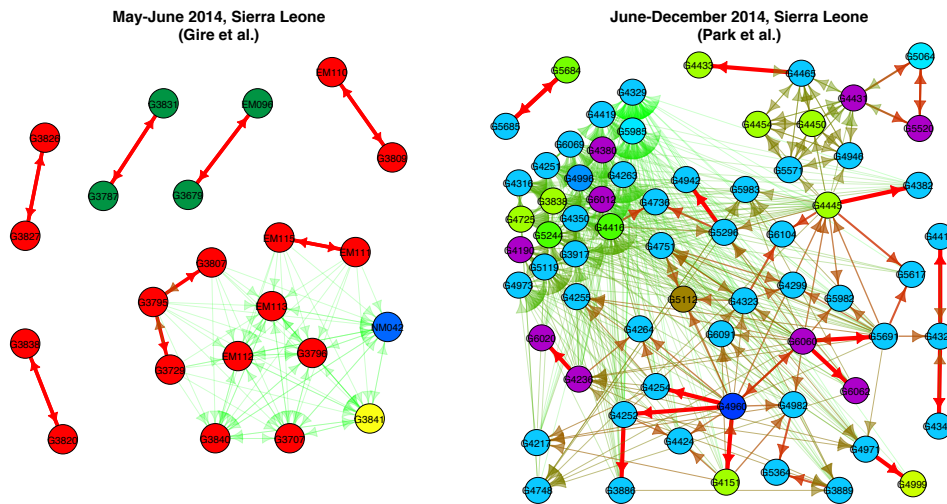
Figure S1. ROC curves. ROC curves for the closest genome tree (green), weighted genetic distance tree (purple), and the weighted variant tree (red). Outbreak simulated with bottleneck size 10 and mutation rate 0.001, corresponding to the transmission trees shown in figure 1. The dashed line indicates the performance of selecting routes at random. Since multiple correct edges may be assigned a weight of zero under all but the weighted distance approach, those ROC curves do not reach (1,1).



1
2 **Figure S2. Area under ROC curves.** The area under the ROC (AUC) metric
3 provides an overall measure of network accuracy. Results for a mutation rate of
4 0.0005 are shown here. Hosts with no observed shared variants are assumed
5 to have all other hosts as potential sources with equal weight. This figure is
6 equivalent to Figure 4B with a lower mutation rate.

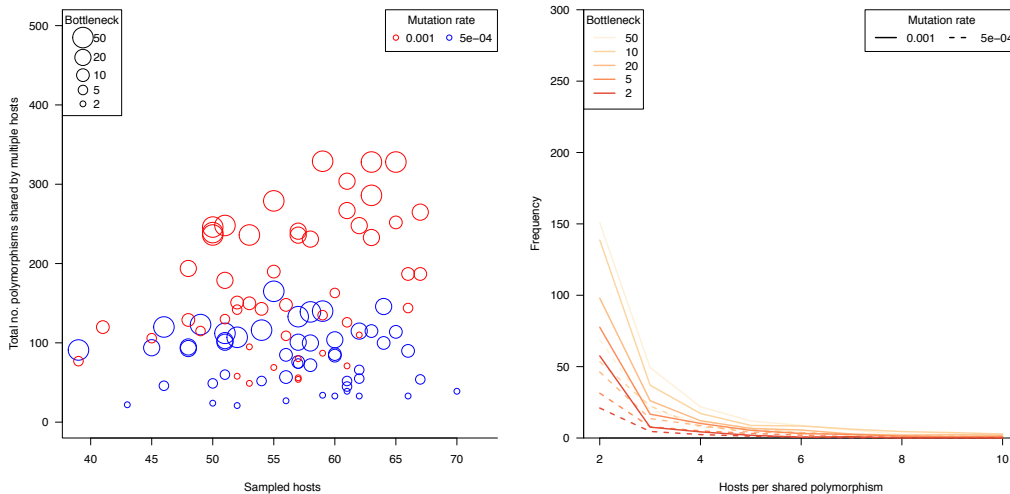


1
2 **Figure S3. True edge weighting.** The mean weight attributed to each true
3 transmission link for each tree reconstruction method, under a range of
4 scenarios and methodologies. Results are shown for a mutation rate of 0.0005.
5



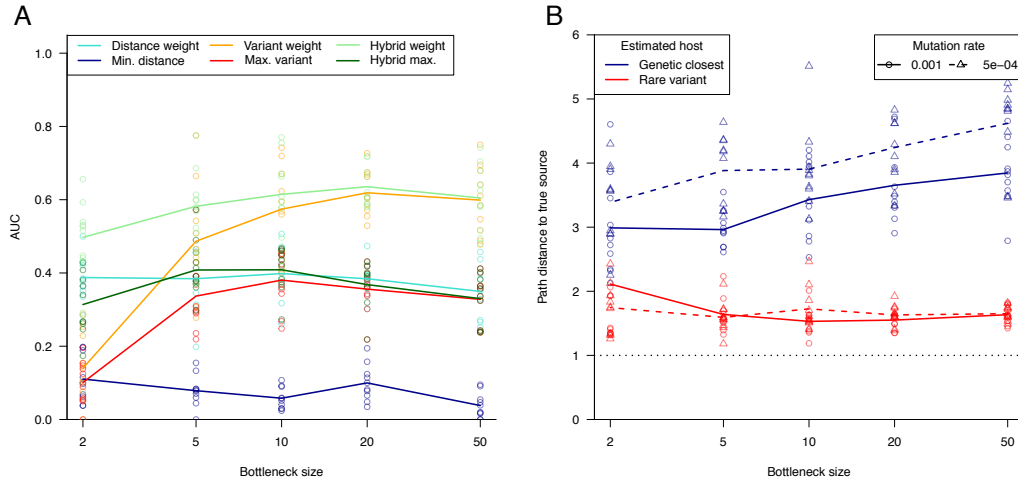
1
2
3
4
5
6
7
8
9

Figure S4. Ebola transmission routes. Estimated transmission links between sampled hosts in the Ebola outbreak under the maximum variant approach. Colors denote the chiefdom to which each host belongs, while the color and thickness of the arrows denote the relative weight attributed to each potential transmission event. Variant detection threshold 1%.



10
11
12
13
14
15
16

Figure S5. Summary of genetic variant frequency across imperfectly sampled simulated outbreaks. Summary of genetic variant frequency across the simulated outbreaks with 30% of infected hosts unsampled. (A) Total number of shared variants across the simulated outbreak. Bottleneck size is illustrated by circle size. (B) Distribution of shared variant group size for different bottlenecks and mutations rates.



1
2
3
4
5
6
7
8
9

Figure S6. Transmission tree reconstruction accuracy. (A) The area under the ROC (AUC) metric provides an overall measure of network accuracy. Results for a mutation rate of 0.001 are shown here. (B) The true path distance between estimated transmission pairs gives insight into the extent to which transmission links are misspecified. A perfect reconstruction would have mean path length 1. Maximum variant path lengths are averaged over identified transmission pairs, that is, excluding hosts with no shared variants.