

1 Inferring transmission bottleneck size 2 from viral sequence data using a novel 3 haplotype reconstruction method

4 Mahan Ghafari,^{a,b,c} Casper K. Lumby,^a Daniel B. Weissman,^b Christopher J. R.
5 Illingworth,^{a,d,e}

6 Department of Genetics, University of Cambridge, Cambridge, UK^a; Department of Physics, Emory University,
7 Atlanta, Georgia, USA^b Department of Zoology, University of Oxford, Oxford, UK^c Department of Applied
8 Mathematics and Theoretical Physics, University of Cambridge, Cambridge, UK^d Department of Biosciences,
9 Department of Computer Science, Institute of Biotechnology, University of Helsinki, Helsinki 00014, Finland^e

10 **ABSTRACT** The transmission bottleneck is defined as the number of viral particles that
11 transmit from one host to establish an infection in another. Genome sequence data
12 has been used to evaluate the size of the transmission bottleneck between humans
13 infected with the influenza virus, however, the methods used to make these estimates
14 have some limitations. Specifically, viral allele frequencies, which form the basis of
15 many calculations, may not fully capture a process which involves the transmission of
16 entire viral genomes. Here we set out a novel approach for inferring viral transmission
17 bottlenecks; our method combines an algorithm for haplotype reconstruction with
18 maximum likelihood methods for bottleneck inference. This approach allows for rapid
19 calculation, and performs well when applied to data from simulated transmission
20 events; errors in the haplotype reconstruction step did not adversely affect inferences
21 of the population bottleneck. Applied to data from a previous household transmission
22 study of influenza A infection we confirm the result that the majority of transmission
23 events involve a small number of viruses, albeit with slightly looser bottlenecks being
24 inferred, with between 1 and 13 particles transmitted in the majority of cases. While
25 influenza A transmission involves a tight population bottleneck, the bottleneck is not
26 so tight as to universally prevent the transmission of within-host viral diversity.

27 **IMPORTANCE** Viral populations undergo a repeated cycle of within-host growth fol-
28 lowed by transmission. Viral evolution is affected by each stage of this cycle. The
29 number of viral particles transmitted from one host to another, known as the transmis-
30 sion bottleneck, is an important factor in determining how the evolutionary dynamics
31 of the population play out, restricting the extent to which the evolved diversity of the
32 population can be passed from one host to another. Previous study of viral sequence
33 data has suggested that the transmission bottleneck size for influenza A transmission
34 between human hosts is small. Re-evaluating these data using a novel and improved
35 method, we largely confirm this result, albeit that we infer a slightly higher bottleneck
36 size in some cases, of between 1 and 13 virions. While a tight bottleneck operates
37 in human influenza transmission, it is not extreme in nature; some diversity can be
38 meaningfully retained between hosts.

39 **KEYWORDS:** influenza A, transmission, population bottleneck

Compiled March 19, 2020

This is a draft manuscript, pre-submission

Address correspondence to Chris Illingworth,
chris.illingworth@gen.cam.ac.uk.

M.G. and C.K.L. contributed equally to this
work.

INTRODUCTION

Viral populations experience large fluctuations in population size. During the course of an infection many thousands of viruses may be produced by each infected cell (1), yet in the process of transmission only a small number of viruses may get through to found a new infection (2). The size of the bottleneck undergone by a viral population at the moment of transmission has an important impact on the evolution of that virus. Where larger numbers of viral particles are involved in transmission, a greater amount of genetic diversity is preserved between hosts; where smaller numbers of particles are transmitted, between-host evolution becomes more of a stochastic process (3). Studying transmission at the scale of individual hosts therefore gives an insight into larger-scale patterns of viral evolution.

Genetic data provides an invaluable insight into processes of viral evolution (4). Such data have been at the core of a variety of approaches for the quantitative analysis of population bottlenecks, typically using observations of minority variants, or their allele frequencies, to make a statistical inference. For example, counting the number of minority variants shared between hosts can be informative of whether transmission occurred between specific hosts (5, 6). If the route of transmission is known, shared variants can be used to estimate the size of the population bottleneck (7). A model of genetic drift may also be applied: smaller or larger changes in the composition of a viral population suggest that a larger or smaller number of viruses were transmitted (3, 8, 9, 10, 11). In some situations, engineered viruses with genetic markers have been used to directly evaluate transmission events (12, 13).

Recent studies of influenza transmission between human hosts have used metrics based upon changes in allele frequencies to evaluate the bottleneck at transmission (3, 11, 14, 15). Such metrics have limitations; transmission is ultimately an event in which whole viruses, rather than independent alleles, are passed from one host to another. Neglecting genetic linkage in this way can skew the results of inference methods (16). Inspired by this, a recent study on the assessment of viral transmissibility used sequence data to evaluate transmission at the level of viral genomes (17).

Accounting for genetic linkage between alleles becomes more difficult as the diversity of a viral population increases. In modelling the action of selection on a diverse population, the large number of potential genome sequences can make calculations infeasible. Considering cases in which selection among transmitted variants is not the dominant effect at transmission (3) we here set out an alternative approach for the inference of population bottlenecks, incorporating the true genetic structure of viruses. Our approach has two components. Firstly, given sequence data collected before and after a transmission bottleneck, we apply a method of haplotype reconstruction, using a maximum likelihood framework to calculate a parsimonious reconstruction of the viral population, as observed before and after transmission. A broad variety of computational tools have previously been described for the purpose of haplotype reconstruction in various contexts (18, 19, 20, 21, 22, 23); ours fits naturally into the bioinformatic framework we have outlined in previous publications (24, 25). Secondly, we use the haplotype reconstruction to infer a bottleneck size at transmission; our framework contains two alternative approaches optimised for smaller and larger bottleneck sizes respectively. We test our method against simulated data describing viral transmission events with a broad range of population bottlenecks. Finally, we re-evaluate data from a previous study of influenza transmission between human

90 hosts (3). Our study supports the hypothesis of a generally small transmission bottle-
91 neck for influenza viral populations (3, 26) albeit with fractionally higher bottleneck
92 sizes inferred from the same data.
93

94 RESULTS

95 As a first step, we considered the relative performance of allele- and haplotype-based
96 approaches to the inference of transmission bottlenecks, using grossly simplified,
97 though hopefully illustrative, examples of viral transmission.
98

99 **Allele-based versus haplotype-based inference** A first example highlighted the
100 potential for allele-based statistics to misrepresent the nature of a viral population
101 (Figure 1). In this simulated system data were collected from before and after a
102 transmission bottleneck. While during transmission the viral population changed sub-
103 stantially at the genotype level, these changes were not fully reflected in the allele
104 frequency data from each population. As a consequence, inferences of the bottleneck
105 at transmission, calculated using haplotype- and allele-frequency methods, differed
106 by close to two orders of magnitude. While an extreme example, this result highlights
107 a fundamental point of biology. Rather than independent alleles, viral transmission
108 involves the transmission of complete viral genomes. Approaches which neglect this
109 may as a consequence be flawed in the results they produce.
110

111 A second example, describing outcomes across a representative range of transmis-
112 sion events, is shown in Figure 2. We here consider the transmission of a hypothetical
113 influenza viral population. For each segment of the virus, the viral population is divided
114 perfectly into two haplotypes, each with a frequency before transmission of exactly
115 50%. For seven of the eight viral segments, precisely one SNP differentiates the two
116 haplotypes, while in the final segment ten SNPs differentiate the haplotypes. In this
117 case, we note that the post-transmission frequency of any given haplotype can be
118 represented as a simple binomial sample from the original population, the chance of
119 any transmitted virus having a certain haplotype being equal to one half. We further
120 note that the same is true for each variant allele; each allele frequency is equal to the
121 frequency of the haplotype which carries it, so that the frequency of the allele is given
122 by a binomial sample. Critically, however, the transmitted haplotype frequencies are
123 independent of one another, while the transmitted allele frequencies are not indepen-
124 dent.
125

126 The lack of independence has a consequence for the inferred transmission bot-
127 tlenecks. In the (harmonic) mean, both the haplotype and allele frequency statistics
128 produce a correct inference. However, the allele-based estimate is statistically less pre-
129 cise (Figure 2B). While in the haplotype inference each segment is weighted the same,
130 the allele-based estimate is weighted heavily towards the outcome of transmission of
131 the final segment. The variance in the outcome of this one segment is greater than
132 the mean variance across segments, leading the allele-based method to, on average,
133 a worse result. Secondly, the false assumption in the allele-based method that allele
134 frequencies are independent leads to a false confidence in the outcome of this method
135 (Figure 2C). The apparently greater amount of data provided by a greater number of
136 polymorphic loci leads to a falsely reduced confidence interval in the bottleneck size at
137 transmission. Where more than one locus is present on a haplotype, and all else being
138 equal, allele frequency methods give less accurate inferences than haplotype-based

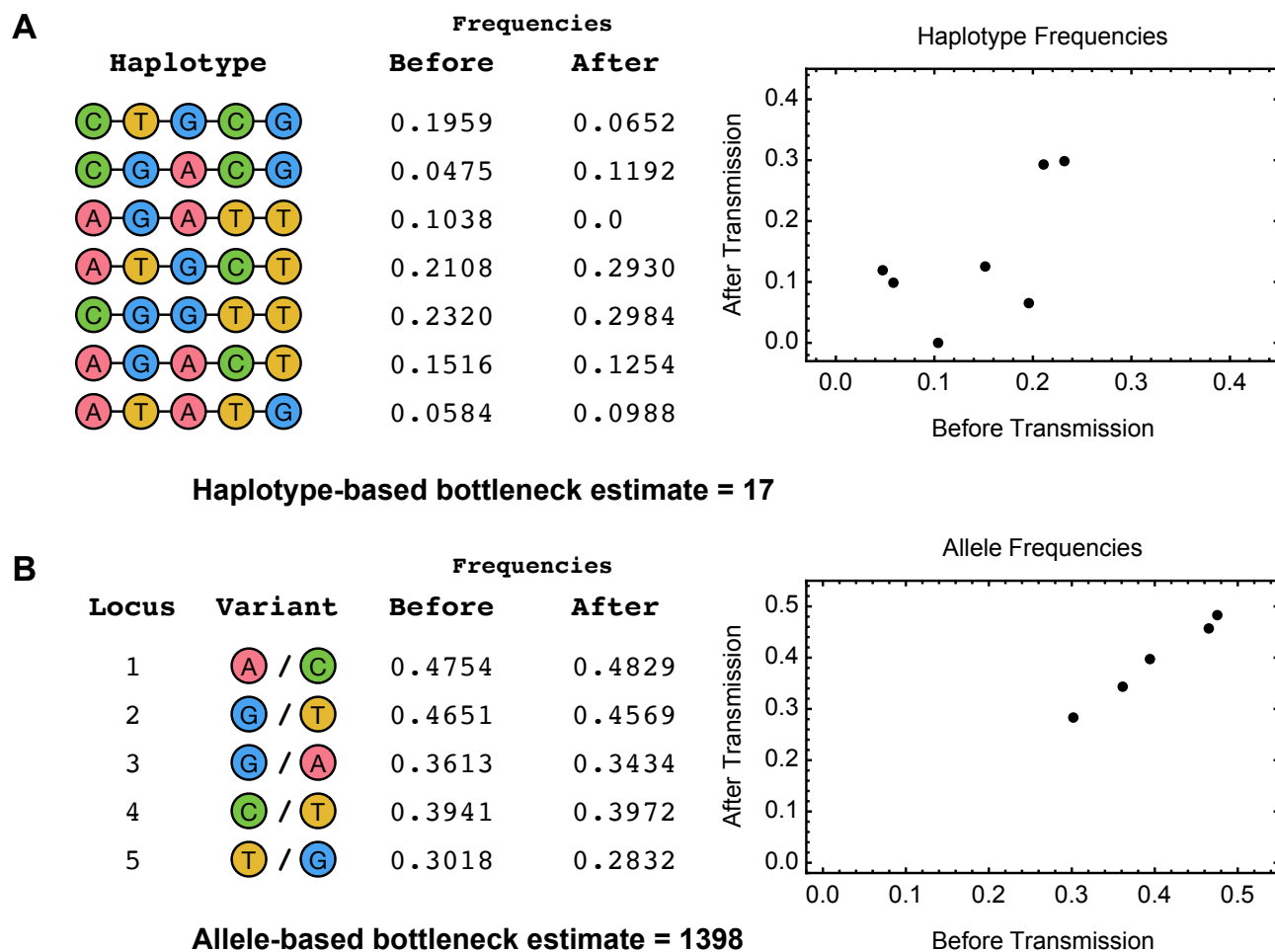


FIG 1 A. Simulated system of viral transmission. A population comprising seven viral genotypes transmits to a new host, leading to a population in the recipient which includes six of the seven genotypes. A plot shows the sampled frequencies of the distinct genotypes, or haplotypes, before and after transmission, reported to four significant figures. Our explicit model of viral transmission based on haplotype frequencies (described in the text) infers a population bottleneck of 17 viruses from these data. B. An alternative analysis of the same population measures allele frequencies from the population before and after the transmission event; these are shown in an equivalent plot. A calculation of the population bottleneck from these data infers a value nearly two orders of magnitude larger than that of our previous calculation.

139 methods, and provide a falsely high level of confidence in their results. We are there-
 140 fore motivated to consider the transmission of viruses on the genotype level.

141
 142 To evaluate our genotype-based approach to bottleneck inference, we first consid-
 143 ered data describing simulated transmission events, before considering data from a
 144 study of human infection.

145
 146 **Haplotype reconstruction** Applied to simulated data our method made a correct
 147 inference of haplotypes (all existing haplotypes identified, with no false identification
 148 of haplotypes) in more than half of the cases tested (Figure 3). Our approach uses
 149 a maximum likelihood method to infer the most parsimonious reconstruction of a
 150 viral population, given sequence data. To test our approach we simulated data de-

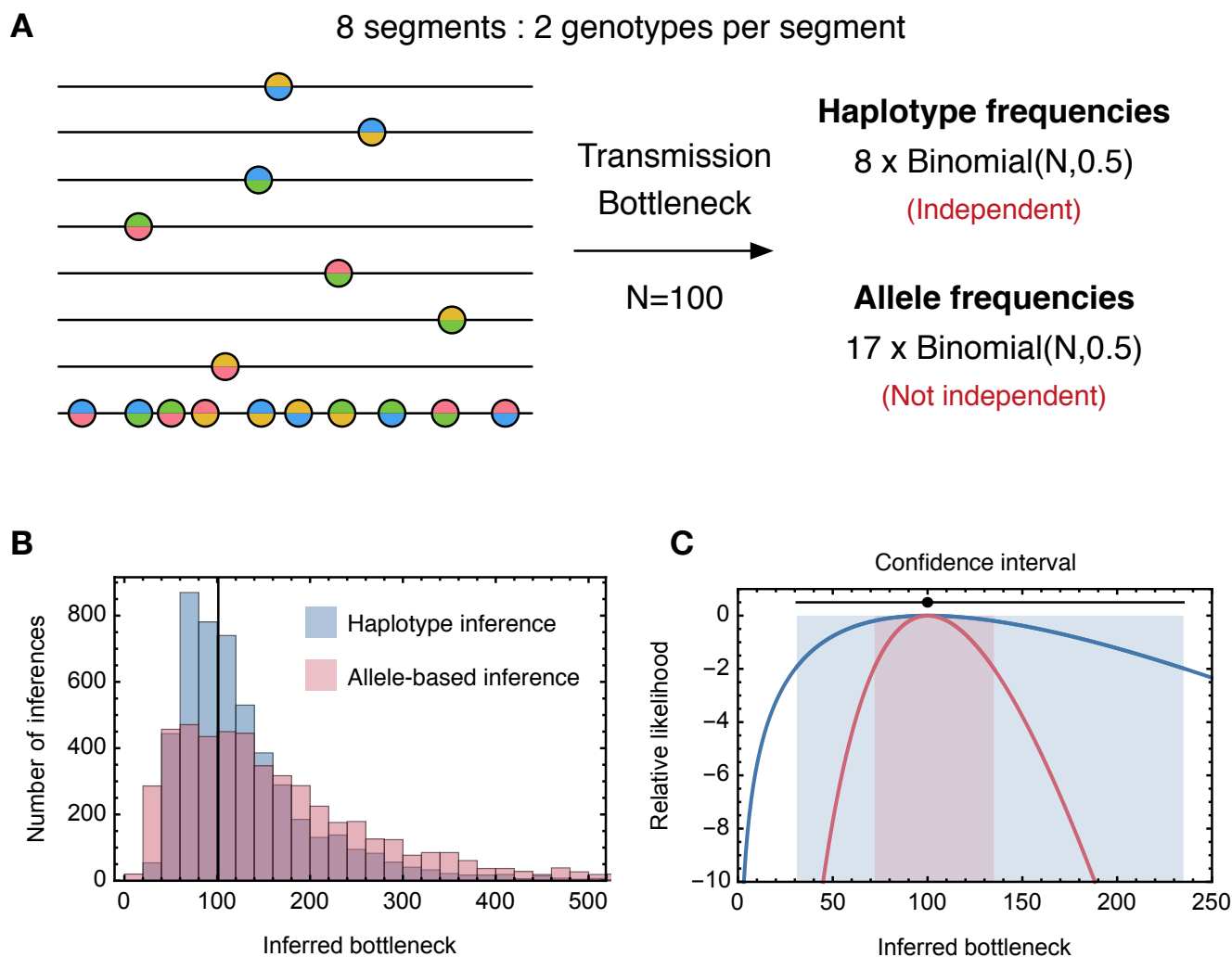


FIG 2 A. Simulated system of viral transmission. A population consists of eight viral segments. For each segment, two haplotypes exist in the pre-transmission population at a frequency of exactly 50%. In seven segments, these haplotypes differ by a single genetic variant while in the eighth the haplotypes differ by ten genetic variants. Post-transmission, the haplotype frequencies in each of the eight segments are described by eight independent random binomial samples. The seventeen allele frequencies are similarly described by seventeen random binomial samples, albeit that these statistics are not independent of each other. B. Inferred population bottlenecks from 5,000 simulations of this transmission process, calculated with haplotype-based and allele frequency-based methods. A method based upon independent transmission of alleles has an increased variance relative to the haplotype-based method. C. Likelihood function for each model in the case in which transmission results in a 45/55 split in haplotype frequencies in each segment. The black circle and line indicate the correct transmission bottleneck and an analytical confidence interval, based upon a window of two likelihood units. The inference in each case is correct, but the allele-frequency method, which treats the allele frequencies as being statistically independent, has a false level of confidence in the inferred value.

151 scribing the transmission of an influenza viral population, from a host to a recipient
 152 individual. Each segment in the population was modelled as containing six distinct
 153 haplotypes, applying a method for generating data described in a previous study (17).
 154 Simulated sequence data from the viral populations in each host were used to infer
 155 which haplotypes were present in the transmission event and their frequencies. The
 156 most common outcome was a correct reconstruction of all of the haplotypes in the

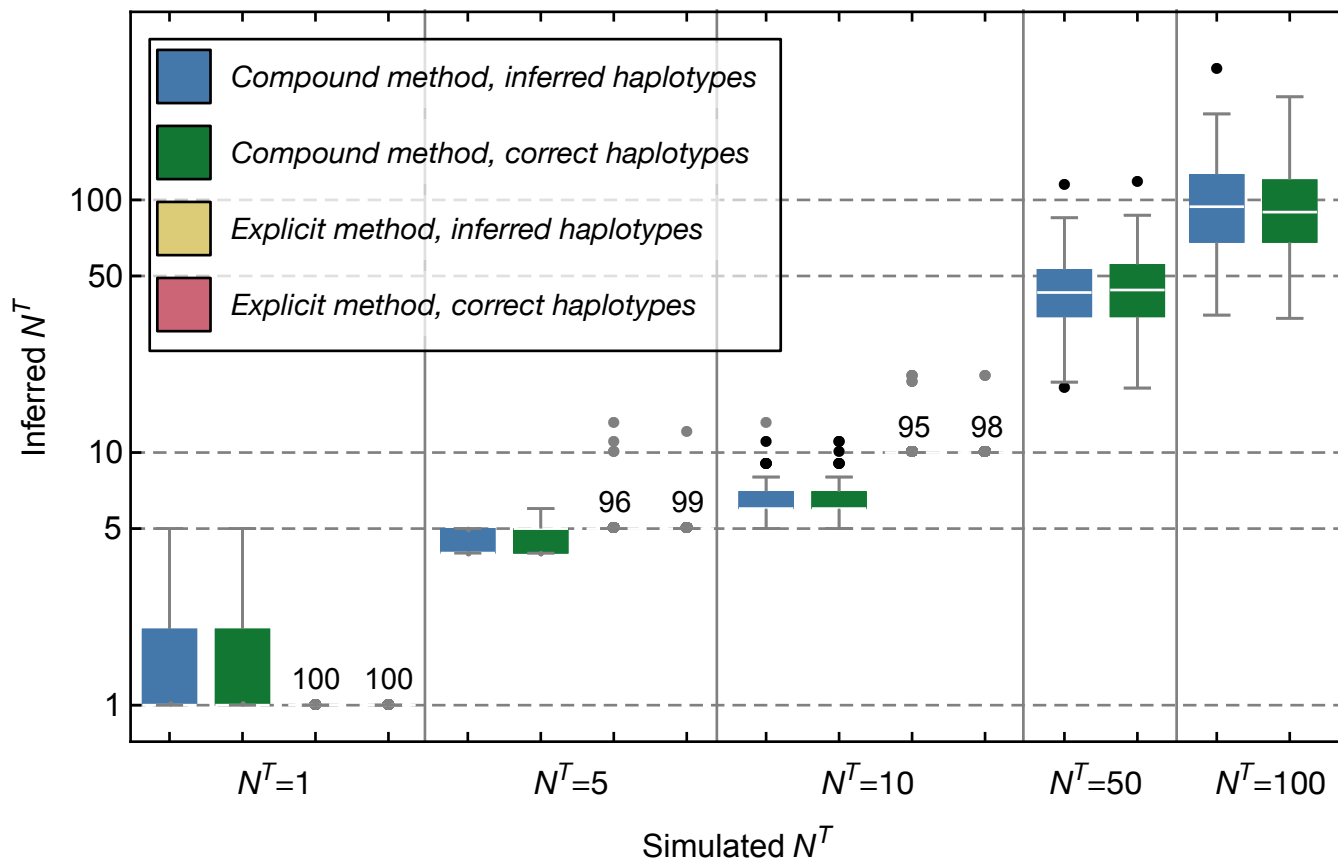


FIG 4 Transmission bottleneck sizes inferred from simulated data using different input data and methodologies. Inferences are shown in colour according to the data and method used. Calculations with inferred haplotypes took as input data generated from a haplotype reconstruction method applied to simulated sequence data, in which both the haplotypes and their frequencies before and after transmission were inferred. Calculations with the correct haplotypes took as input data from a haplotype reconstruction in which the identities of the correct haplotypes were given, with only their frequencies being inferred. Inferences from the explicit method were only calculated for smaller population bottleneck sizes as the method does not scale well to evaluating larger bottlenecks. Results from the explicit method were so accurate as to not have a meaningful interquartile range: numbers displayed in these cases indicate the number of inferences giving a precisely-correct inference of the population bottleneck. Horizontal dashed lines indicate the simulated bottleneck sizes.

Further inferences of bottleneck size were made using reconstructions of haplotypes in which the correct simulated haplotypes were pre-specified, learning only their frequencies. Using these improved data did not produce a noticeable improvement in the inference of the bottleneck size, suggesting that our inference of bottleneck size is robust to errors that arise from our haplotype reconstruction method. Bottleneck sizes in each case were calculated across eight independent viral segments.

Given our simulated data, the explicit method outperformed the compound method at low bottleneck sizes, inferring exactly correct values in the majority of cases with very little error. A disadvantage of the explicit method is that in requiring the evaluation of all possible outcomes of a transmission event, the computational time it requires grows very rapidly as the bottleneck size increases. For this reason, we did not apply it to data from higher simulated population bottlenecks. As with the compound method performance did not greatly improve given frequencies inferred

197 using the correct viral haplotypes; errors in haplotype reconstruction did not have a
198 strong effect on the inferred bottleneck sizes.

199
200 The variances in the inferred bottleneck sizes are dependent upon the amount of
201 data available to our code for inference. In the case of a less diverse viral population,
202 less genetic information would be available, leading to a greater variance in the inferred
203 bottlenecks. By contrast, more diversity would lead to a more constrained inference.
204 Data shown here are intended to illustrate the mean performance of our methods.

205
206 Inference of bottleneck size for a segment was not possible in two cases. Firstly, if
207 our haplotype reconstruction found evidence for only a single viral haplotype, no infer-
208 ence was possible, insufficient information about the event being available. Secondly,
209 if the viral population in the recipient was inferred to have arisen purely from a *de novo*
210 haplotype, which had swept to fixation in the population between the establishment
211 of the infection and the collection of the sequence data, this result was uninformative
212 in identifying a bottleneck. In either of these circumstances, data from a viral segment
213 were ignored, inferences conducted for the remaining segments being combined to
214 infer the final bottleneck size.

215
216 In considering the differences in inferences achieved by the two methods at low
217 bottleneck sizes, it is perhaps helpful to consider the simple case where a single allele
218 frequency is observed to change from 50% frequency in the donor to 5% in the recipi-
219 ent. Within the compound method this represents a large change in allele frequency,
220 corresponding to a large amount of genetic drift, and will be interpreted as resulting
221 from a low bottleneck size. By contrast under the explicit method variation at a fre-
222 quency of 5% is unlikely to be observed if the bottleneck is low; at least one particle
223 with the variant must have been transmitted implying a minimum variant frequency
224 of at least $1/N^T$. Transmission with a bottleneck closer to 20, with sampling noise
225 leading to the underestimation of the variant frequency, would give a more coherent
226 explanation.

227
228 **Application to data from a household study** Our transmission model was ap-
229 plied to data collected from a previously published household study (3). This study
230 used a single-locus inference model to identify narrow bottlenecks in human-to-human
231 transmission, with all but a single event being inferred to involve the transmission of
232 between one and four viral particles. Short-read data from this study were filtered
233 and processed into variant data before being fed into our method. Having identified
234 polymorphic loci in pairs of transmission data using an allele frequency cutoff of 2%
235 we generated multi-locus reads from the data using the SAMFIRE software package (25),
236 using these to generate an inference of haplotype frequencies before and after trans-
237 mission. These frequencies were used to infer population bottleneck sizes for each
238 transmission event.

239
240 We confirm the previous inference of tight population bottlenecks in all cases
241 (Figure 5). In the majority of transmission events (29 out of 38 events for which we ob-
242 tained an inference), bottlenecks of size $N^T = 1$ were inferred by both of our methods,
243 consistent with all of the diversity of the viral population in the original host being lost
244 at transmission. While not necessarily implying that these infections were started by a
245 single viral particle, these results are consistent with the hypothesis of a generally tight
246 bottleneck at transmission. In eight out of the remaining nine transmission events,

intermediate bottleneck sizes were inferred, with a range from 2 to 7 in the compound method and from 2 to 13 in the explicit method. Evidence from simulated data suggests that the explicit method is probably more accurate in this range. Finally, there was a single case in which a bottleneck size of 200 or more was inferred; this was set as the upper limit considered by our study. Our inference in this case matched the original analysis of the data. A further statistical analysis of the samples collected before and after transmission indicated a greater degree of similarity between allele frequencies than was previously found in a case where replicate clinical samples were processed and sequenced in parallel (27). Whereas in the previous study, measurements of allele frequencies from samples split from the cDNA synthesis step onwards were consistent with an effective read depth (that is equivalent to an error-free sample depth) of one thousand or more, here an effective depth in excess of 20,000 was inferred, demonstrating that the before- and after-transmission samples were extremely similar. This case could represent either a very unusual transmission event, in which an extreme number of viruses were transmitted, or potentially an isolated error in the processing of a large number of sequence samples.

Cases in which the explicit method inferred larger bottleneck sizes than the compound method could be explained in terms of the preservation of allele frequencies at relatively low frequencies; as explained above the explicit method can favour a higher bottleneck in such cases.

Our approach was not able to infer a population bottleneck in five of the transmission cases analysed by the original study. In these cases a low level of polymorphism observed before transmission was no longer present after transmission. Application of our haplotype reconstruction method in these cases did not find statistical evidence for more than one haplotype (plus noise) in these systems, at least two specific haplotypes being required for an inference of bottleneck size. We understand this in terms of our haplotype reconstruction method being less sensitive to detecting variation than is the 2% allele frequency cutoff used in the original study; the presence of a variant allele at 2% frequency was not always sufficient evidence for our code to infer the existence of two specific genetic variants in the population. In these cases, the loss of host genetic variance at transmission would lead our methods to the conclusion that a bottleneck of $N^T = 1$ best explained the observed data, strengthening our main result of a tight bottleneck size. The sensitivity of our method in calling additional haplotypes can be somewhat arbitrarily tuned.

Differences in the bioinformatic processing of data could underlie some of the differences in bottleneck we identified. While we replicated the 2% allele frequency cutoff of the original paper (3), we called variants in 18 of the 38 transmission events analysed here that were not originally found. Such variants were primarily only found in one of the two samples, and existed at frequencies very marginally above the 2% threshold; minor allele frequencies very close to the threshold were observed both in our processing of the data and in the original study (Supporting Table S1). Applying the exact single-locus method for bottleneck inference of a previous study (for convenience we term this the exact SL method) (11) found cases of higher bottlenecks than were found in the original paper (Figure 5). In common with the original study, we remove variants in non-coding regions of the genome from our calculation.//

Bioinformatic variations in the calling of alleles can have three distinct effects. Where an additional variant is called in the recipient population but not in the donor,

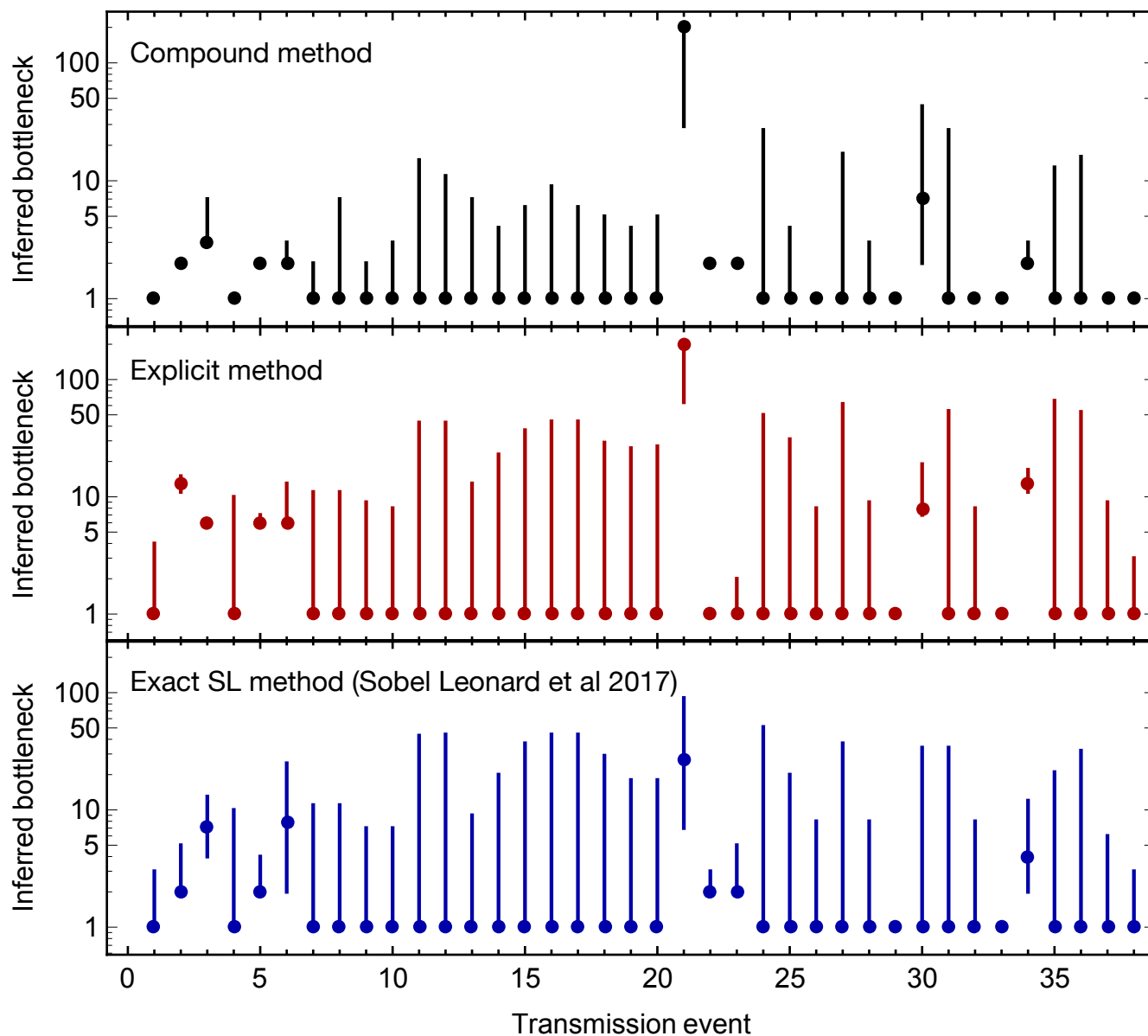


FIG 5 Bottleneck sizes inferred from the data presented by (3). Dots indicate the maximum likelihood bottleneck size inferred for each of the 38 systems in this work for which we were able to infer a bottleneck. Vertical bars represent confidence intervals equivalent to a cut-off of 2 log likelihood units.

297 no change in the inferred bottleneck occurs; the variant is assumed to have arisen
 298 *de novo* in the recipient, having nothing to do with the transmission event. Where
 299 an additional variant is called in the donor population but not in the recipient, this
 300 shifts the inference towards a smaller bottleneck. The dying out of a low-frequency
 301 variant is the most likely outcome given a small bottleneck, so this usually makes little
 302 difference to the inference. However, in transmission event 21, we observe that a
 303 bottleneck inferred to involve at least 200 particles by both of our haplotype-based
 304 methods (and the original study) was inferred to involve only 29 particles by the exact
 305 SL method. In this case our bioinformatic approach called two variant alleles, NA
 306 G1351A and PB1 A2280G, at 2.3% and 3.2% in the donor population, which died out

upon transmission. Our haplotype inference method did not find sufficient evidence to identify two haplotypes for these segments, and ignored these variants as a result, but the exact SL method accounted for them, leading to a reduced bottleneck inference. Especially at high bottlenecks, small bioinformatic changes can have an important effect.

Finally, where an additional variant is called in both the donor and the recipient populations, it can influence the inferred bottleneck in either direction. Four such cases were found in our analysis, in transmissions 2, 3, 5, and 6. Removing these variants from the populations led to a reduction in the bottleneck inferred under the explicit method to a single particle for transmissions 2, 3, and 6. The inferred bottleneck for transmission 5 was slightly reduced from $N^T = 6$ to $N^T = 5$. Not all of the cases in which bottlenecks of greater than 1 were inferred could be explained by bioinformatic variation. The inference of $N^T = 13$ in transmission 34 had a single additional variant in our processing that was not found in the original analysis, consisting of a low-frequency variant that was not transmitted to the recipient host. As noted above, such a variant could not increase the size of the inferred bottleneck.

DISCUSSION

We have here set out a haplotype-based approach for the inference of transmission bottlenecks, and demonstrated its application using data from a study of transmission of influenza A infection.

Haplotype-based methods have the advantage of faithfully representing the biological event of viral transmission. While the use of allele frequency statistics does not necessarily lead to incorrect results, such use introduces a level of abstraction from reality. In some cases this can lead to grossly misleading results; in general it will give a less precise inference of bottleneck size, and a falsely high level of confidence in the results obtained. The shortfall in performance of an allele-based method will depend upon the system in question. In a hypothetical influenza virus with only a single variant per segment, allele- and haplotype-based approaches will likely give identical results. In a non-segmented virus, with high viral diversity, the assumption of independent alleles will lead to a substantial over-estimation of the statistical confidence with which a bottleneck can be quantified.

We used a haplotype reconstruction method to infer the composition of the viral population before and after transmission; by requiring substantial evidence to add an additional haplotype to the model, this approach limits the complexity of the inferred viral population, improving the feasibility of haplotype-based bottleneck inference relative to a previous approach (17). While our haplotype reconstruction method was not perfect in reproducing the details of a viral population, errors resulting from this method did not greatly harm our inference of population bottleneck sizes.

Our approach for bottleneck inference comprises two distinct methods, optimal for distinct transmission bottleneck sizes. The first of these generalises the approach of Poon et al. (14), who used a formula based on genetic drift to evaluate changes in allele frequencies. Our compound method generalises this to changes in haplotype frequencies, which occur in higher-dimensional sequence space; it further incorporates uncertainty in the inferred haplotype frequencies and genetic drift arising from within-

356 host population growth. This method has the advantage of being rapid to calculate
357 at high bottleneck sizes, but potentially underestimates bottleneck sizes at low values
358 of N^T . Our second method, the explicit method, generalises the approach of Sobel
359 Leonard et al. (11), who apply a beta-binomial formula to evaluate possible discrete
360 outcomes of a transmission process. In spirit we repeat this approach, summing a like-
361 likelihood function over the set of possible outcomes of a transmission of viral haplotypes.
362 This approach is limited in its application to systems of higher complexity, becoming
363 slow where there are many haplotypes or where N^T is large, but is likely more accurate
364 at lower bottleneck sizes. The size of a bottleneck affects the two methods in different
365 ways. For the compound method, increased bottleneck size leads to greater accuracy,
366 in that the mathematical approximations underlying the method become increasingly
367 correct as the product between the bottleneck size and a typical haplotype frequency
368 increases. For the explicit method, increased size adversely affects the time required
369 for calculation, in that as the number of haplotypes in the system and the bottleneck
370 size become large, the evaluation over all possible outcomes of transmission becomes
371 increasingly difficult to calculate.

372
373 While our haplotype reconstruction and bottleneck inference methods are con-
374 structed upon a common likelihood framework, our inference methods could be
375 applied to haplotype data from other sources. Other reconstruction methods could
376 provide appropriate data for analysis, while barcoding technologies or long-read se-
377 quencing could each obviate the need for a reconstruction step. We note that, where
378 ethically feasible, the use of neutral markers provides a more direct approach for
379 evaluating transmission events (12).

380
381 Our framework makes the assumption of selective neutrality during the transmis-
382 sion event. Selection during transmission, whether positive or negative, changes the
383 genetic composition of the viral population in the recipient relative to that of the donor.
384 On average, this makes the population in the recipient less similar to that in the donor,
385 leading to an underestimate of the population bottleneck. A variant of our compound
386 method incorporating selection has been set out in a previous publication (17). Evalu-
387 ating selection requires a comprehensive reconstruction of the extant viral haplotypes;
388 this may be difficult to obtain given short-read data describing a diverse population.
389 Identifying variants that enhance viral transmissibility is impossible where very few
390 viruses are transmitted; at higher population bottlenecks or where multiple transmis-
391 sions are observed it becomes an achievable task. Under selection, haplotype-based
392 approaches have further advantages over allele-frequency statistics (16).

393
394 As we have shown, apparently small differences in the calling of variants can
395 have significant consequences for the inference of bottleneck sizes. Regardless of
396 the method used for inference, if a variant was falsely called to exist at low frequency
397 in both the pre- and post-transmission populations, this could dramatically skew an
398 inference towards a higher bottleneck size. Our re-analysis of data preserved the
399 frequency cutoff for alleles used by the original authors, but nevertheless found ad-
400 ditional variants in excess of this cutoff, likely the result of fractional changes in the
401 bioinformatic processing. Marginal frequencies close to the frequency cutoff were
402 identified both in our processing of the data and in the output of the original study.
403 Where a hard cutoff is used for variant identification, and specific variants are close to
404 this cutoff, uncertainty in the identification or non-identification of variants should be
405 considered as part of the uncertainty in bottleneck inference; statistical approaches for

406 this could provide an area for future development.

407
408 Progress in understanding the biology of infection could be a further aid in the
409 development of methods for bottleneck inference. In particular, the dynamics of the
410 very early stages of population growth, from the initial founder viruses to the large
411 population typical of influenza infection, are not necessarily well understood. Knowl-
412 edge of the extent to which this affects the genetic composition of the viral population
413 would improve the potential for accurate inference.

414
415 We have here used a haplotype-based approach to study transmission bottlenecks
416 using data from a household study of influenza A infection. While we replicate the
417 finding that transmission involves a small number of viral particles, our results have
418 a longer tail of bottleneck sizes, with estimates of up to 13 viruses were transmitted.
419 While transmission may strongly limit the inheritance of influenza virus diversity, its
420 effect in doing so is not absolute; the transmission of viral diversity may occur and
421 have some influence on broader viral evolutionary dynamics.

423 METHODS

424 **Notation** A guide to the notation used in our methods is shown in Figure 6. Briefly,
425 we represent the populations before and after transmission by vectors of unknown
426 haplotype frequencies, referred to as \mathbf{q}^B and \mathbf{q}^A respectively. These are separated
427 by transmission with a bottleneck N^T , forming the founder viral population \mathbf{q}^F in the
428 recipient, then within-host growth, represented in our model by a single generation
429 of genetic drift with effective size N^G . The unknown vectors \mathbf{q}^B and \mathbf{q}^A are indirectly
430 observed via the datasets \mathbf{x}^B and \mathbf{x}^A , which are used to generate the estimated haplo-
431 type frequencies \mathbf{q}^{*B} and \mathbf{q}^{*A} .

432
433 In generating the variance of our estimates, we use \mathbf{q}^{*B} and \mathbf{q}^{*A} to generate simu-
434 lated observations, which we term \mathbf{x}^{*B} and \mathbf{x}^{*A} . These in turn are used to generate a
435 new round of estimates \mathbf{q}^{**B} and \mathbf{q}^{**A} . In so far as \mathbf{q}^{**B} , \mathbf{q}^{**A} , \mathbf{q}^{*B} , and \mathbf{q}^{*A} are all known,
436 they may be used to estimate the variances of \mathbf{q}^{*B} and \mathbf{q}^{*A} .

437
438 **Haplotype reconstruction** We developed a maximum likelihood approach for
439 haplotype reconstruction based upon existing technologies for processing short read
440 data (24, 25, 27). We here assume that we have short-read data describing a viral
441 population both before and after a transmission event. Before commencing haplotype
442 reconstruction we performed three steps to pre-process the data using our software
443 package SAMFIRE (25). Firstly, after alignment to the viral genome using BWA (28), the
444 short read data were filtered, trimming reads to achieve a median PHRED score of at
445 least 30, combining data from paired-end reads, and removing individual base calls
446 with a PHRED score less than 30. Secondly, the filtered data were used to identify loci at
447 which a polymorphism existed at significant frequency, this being defined using a cutoff
448 of 2% to match the study of McCrone et al (3), from which we obtained the data we
449 analysed. Thirdly, reads were processed to generate partial haplotypes, which describe
450 the nucleotides present at each of the polymorphic loci in each read. Partial haplotype
451 data were divided into distinct sets of reads, each describing alleles at a distinct set
452 of loci in the viral genome. As an optional step, an estimate may be produced of the
453 extent of noise present in sequence data, inferring a parameter, C , which describes
454 the precision with which measurements of allele frequencies may be calculated via

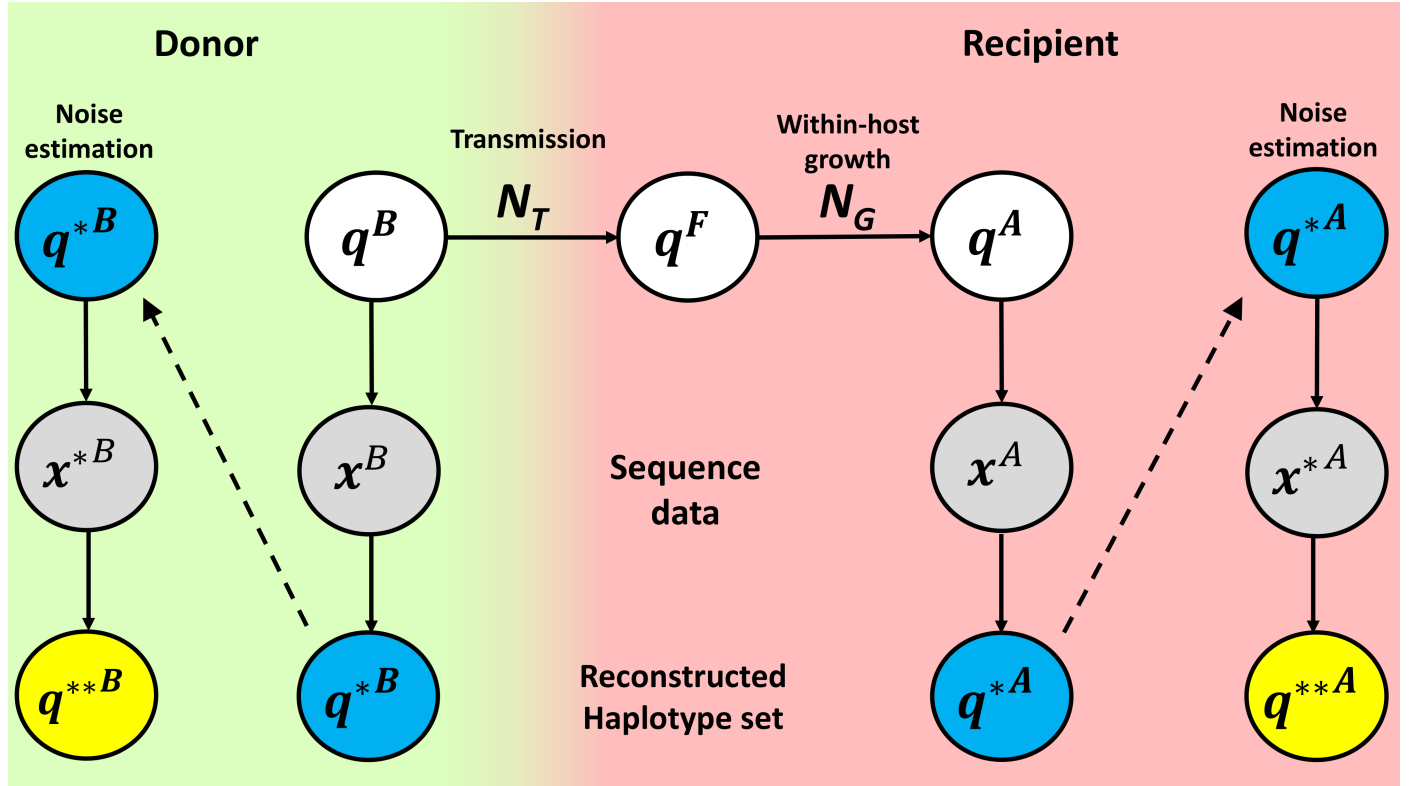


FIG 6 Notation in the transmission model. Transmission of the population q^B with bottleneck N^T results in the founder population q^F . The founder population grows under the influence of genetic drift, the effects of which are described by the effective population size N^G . Growth results in the population q^A . The populations q^B and q^A are observed, producing datasets represented by x^B and x^A , which are used to reconstruct the original populations in terms of haplotypes. In order to calculate the variance of the reconstructed populations q^{*B} and q^{*A} , datasets equivalent to x^B and x^A , denoted x^{*B} and x^{*A} are generated and used to infer sets q^{**B} and q^{**A} .

sequencing (25). A value of $C = 1$ here corresponds to a case in which reads are uninformative, while large values of C tend towards the binomial case in which each read accurately describes the allele present in a distinct viral genome, sampled in an unbiased manner from the population. A default value of $C = 200$ was used for our simulations.

We denote the sets of partial haplotype data collected before and after transmission as $x_l^{B,P}$ and $x_l^{A,P}$ respectively, where l denotes the partial haplotype set. We now suppose that the viral population is comprised of a set of distinct haplotypes, denoted H , which comprises k haplotypes, having the frequencies $q^B = \{q_i^B\}$ before transmission and $q^A = \{q_i^A\}$ after transmission. These frequencies can be converted into partial haplotype frequencies by projection of the full haplotype space onto each lower-dimensional partial haplotype space by means of matrices T_l . For example, given the full haplotypes before transmission $\{GA, TA, GC, TC\}$ and a set of partial haplotypes $\{G-, T-\}$, we may write

$$q_l^{B,P} = T_l q^B, \quad (1)$$

or more explicitly,

$$\begin{pmatrix} q_{i,1}^{B,P} \\ q_{i,2}^{B,P} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} q_1^B \\ q_2^B \\ q_3^B \\ q_4^B \end{pmatrix}. \quad (2)$$

In the above instance we note that each partial haplotype can potentially be emitted from at least one of the haplotypes in \mathbf{H} . In order to generalise our model, we included in each set \mathbf{H} a further haplotype 'X', describing the cloud of all potential viral haplotypes of the same length as those in \mathbf{H} , yet not already defined as being in \mathbf{H} . With this inclusion, we may say that any potential partial haplotype may be emitted from at least one of the haplotypes in \mathbf{H} , being emitted either from one of the defined haplotypes or from 'X'.

In this way, we can construct a likelihood for any given set of haplotypes and frequencies, given the partial haplotype data. We write:

$$\log \mathcal{L}(\mathbf{H}) = \sum_{t \in \{B,A\}} \sum_l \log \mathcal{L}_D(\mathbf{x}_l^{t,P} | T_l \mathbf{q}^t, C), \quad (3)$$

where \mathcal{L}_D denotes the Dirichlet multinomial likelihood

$$\mathcal{L}_D(\mathbf{x} | \mathbf{q}, C) = \frac{\Gamma(N+1)}{\prod_i \Gamma(x_i+1)} \frac{\Gamma(\sum_i C q_i)}{\Gamma(N + \sum_i C q_i)} \prod_i \frac{\Gamma(x_i + C q_i)}{\Gamma(C q_i)}, \quad (4)$$

in which $N = \sum_i x_i$.

A two-step optimisation was used to infer the optimal set of haplotypes and frequencies. To construct an initial set \mathbf{H} , a set of $k \geq 1$ unique haplotypes were created in turn, to which was added the additional X haplotype. The frequencies of these haplotypes before and after transmission were then optimised under the constraint that the frequency of the X haplotype could not be greater than 0.01; this prevents the inference of trivial solutions to the model. We denote the inferred haplotype frequencies as \mathbf{q}^{*B} and \mathbf{q}^{*A} . We note that the frequency of the X haplotype may be effectively zero; for the purposes of calculation a minimum frequency of $\epsilon = 10^{-20}$ was imposed.

Given our likelihood function, a series of changes were made to the set \mathbf{H} , optimising the frequencies each time to find the optimal haplotype reconstruction. Repeating this for increasing values of k gives a series of fits to the data; we used the Bayesian Information Criterion (BIC) to distinguish the most parsimonious explanation for the data:

$$\text{BIC}_k = -2\mathcal{L}^*(\mathbf{H}_k^*) + k \log N, \quad (5)$$

where $\mathcal{L}^*(\mathbf{H}_k^*)$ is the optimum likelihood value for the optimal set \mathbf{H}_k^* of k haplotypes, and N is the total number of observations in the dataset. Optimisation of the haplotype set was conducted for increasing values of k until a model with an additional haplotype produced an improvement of less than 10 units of BIC, representing a conservative cutoff point; a smaller required improvement would lead to the inference of a greater number of haplotypes. In our model the same k haplotypes had to be used for the reconstructions of both the pre- and post-transmission samples. Our model retained the possibility of haplotypes having zero frequency after transmission, for

example in the case of a tight bottleneck, or before transmission, in the case of the emergence of a *de novo* mutation following a transmission event.

Estimated error in reconstructed haplotype frequencies For our compound method for bottleneck inference, we require an estimate of the variance in the inferred haplotype frequencies \mathbf{q}^{*B} and \mathbf{q}^{*A} , so as to account for noise in these parameters when evaluating changes in the population. Variances were calculated by means of simulated data. Considering data collected before transmission, we used the frequencies \mathbf{q}^{*B} to generate sets of partial haplotype data $\mathbf{x}_{i,j}^{*B,P}$, where j is used to index different sets. Each set provided an independent statistical replicate of the original data; having an identical number of sets of partial haplotypes, each spanning the same loci and containing the total number of samples. Each set was generated using a random Dirichlet multinomial sampling process with value C identical to the original. For each set of data, the haplotype reconstruction process was repeated, but with the haplotypes \mathbf{H} constrained to those inferred for the original data. This process was repeated for 100 sets of data, generating the inferred haplotype frequencies $\{\mathbf{q}_j^{**B}\}$. These values were used to calculate the diagonal elements of a covariance matrix $\text{var}[\mathbf{q}^{*B}]$ for \mathbf{q}^{*B} , given by:

$$\text{var}[\mathbf{q}^{*B}]_{i,i} = \frac{1}{100} \sum_{j=1}^{100} (q_i^{*B} - \{q_{ij}^{**B}\})^2. \quad (6)$$

For simplicity, off-diagonal elements of this matrix were set to zero. An identical process was used to generate the matrix $\text{var}[\mathbf{q}^{*A}]$.

Allele-frequency models of bottleneck inference In generating Figures 1 and 2 we used a simple single-locus model of bottleneck inference. Given a set of independent allele frequencies q_i^B at locus i in the pre-transmission viral population, and their equivalent values q_i^A in the post-transmission population, we note that in the absence of selection, the mean value of q_i^A is given by q_i^B , while the variance of q_i^A , arising from genetic drift in a haploid system is given by

$$V = \frac{q_i^B (1 - q_i^B)}{N}, \quad (7)$$

where N is the effective population size of the system (29).

To estimate the bottleneck size at transmission, we made the approximation that q_i^A is normally distributed, then maximised the sum of the log likelihood values across allele frequencies

$$\mathcal{L}(N^T) = \sum_i \log \mathcal{N} \left(q_i^B, \frac{q_i^B (1 - q_i^B)}{N^T} \right), \quad (8)$$

where N^T is the transmission bottleneck and the sum is calculated over loci i with polymorphic alleles. In the case where only two haplotypes are observed in a segment, this approach can be applied to haplotype, rather than allele frequencies. This was done for the haplotype-based calculations in Figure 2.

In the analysis of influenza sequence data we applied the exact version of the beta-binomial sampling method described by Sobel Leonard et al (11). This method identifies the value of N^T that maximises the likelihood

$$\mathcal{L}(N^T) = \sum_i \sum_{k=0}^{N^T} \binom{x_i^A}{n_i^A} \frac{B(n_i^A + k, x_i^A - n_i^A + N^T - k)}{B(k, N^T - k)} \binom{N^T}{k} (q_i^B)^k (1 - q_i^B)^{N^T - k} \quad (9)$$

where x_i^A is the total number of reads at locus i , n_i^A is the number of reads at i which describe the variant allele, $B(\alpha, \beta)$ is the beta function, and the outer sum is conducted over polymorphic loci.

Haplotype-based methods of bottleneck inference Frequencies inferred from the haplotype reconstruction were used for the explicit and compound methods for calculating bottleneck size. As a first step in each method we removed haplotypes that were inferred to have been created *de novo* in the recipient following the transmission by removing haplotypes for which the pre-transmission frequency fell below a threshold frequency δ , set by default to 0.5%. Elements of the vectors \mathbf{q}^{*B} and \mathbf{q}^{*A} and the respective rows and columns of their covariance matrices were removed in this preliminary step.

In so far as we consider influenza transmission, we consider data from each viral segment independently, calculating first a likelihood of the bottleneck size given data from each segment, before combining the likelihoods across segments to estimate an overall maximum likelihood value for the transmission bottleneck.

Compound method for bottleneck estimation In the case of larger values of N^T , an approach building upon that described in a previous publication (17) was applied. Briefly, we note that in a neutral transmission bottleneck, the expected composition of the population in the recipient is identical to that in the original host. The variance in this population is then a function of the size of the bottleneck and the extent of genetic drift during within-host growth, while in the case of inference, variation arising from the measurement of each population must also be considered.

Similarly to the approach outlined in an earlier work (17), we calculate a likelihood function with two components:

$$\begin{aligned} \mathcal{L}(N^T | \mathbf{q}^{*B}, \mathbf{q}^{*A}, N^G) &= \int P(\mathbf{q}^{*B} | \mathbf{q}^B) P(\mathbf{q}^B) d\mathbf{q}^B \\ &\times \int P(\mathbf{q}^{*A} | \mathbf{q}^A) \left\{ \int P(\mathbf{q}^A | N^G, \mathbf{q}^F) \right. \\ &\quad \left. \times \left(\int P(\mathbf{q}^F | N^T, \mathbf{q}^B) P(\mathbf{q}^B) d\mathbf{q}^B \right) d\mathbf{q}^F \right\} d\mathbf{q}^A, \end{aligned} \quad (10)$$

where the first integral corresponds to the initial observation of the system and the second encompass transmission (with the bottleneck N^T), within-host growth (with drift described by the effective size N^G) and post-transmission sampling. Each component of the likelihood is relatively simple to consider, as either a multinomial or Dirichlet-multinomial process, but the compound is difficult to evaluate. We note that, in cases where the frequency of a haplotype remains far from 0 or 1, and in particular as N^T becomes large, the likelihood can be increasingly well approximated in terms of a Gaussian distribution, with mean and variance calculated below.

Our solution makes use of the laws of total expectation and total variance. Given distributions U in x and V in y , the compound distribution W takes the form

$$P_W(x) = \int P_U(x|y) P_V(y) dy. \quad (11)$$

The mean and variance of W are then defined by

$$E_W[x] = E_V[E_U[x|y]], \quad (12)$$

598 and

$$\text{var}_w[x] = E_V[\text{var}_U[x|y]] + \text{var}_V[E_U[x|y]], \quad (13)$$

600 respectively.

602 For the pre-transmission component, the calculation of mean and variance are
603 simple; our haplotype reconstruction process gives the estimate

$$E[\mathbf{q}^B] \approx \mathbf{q}^{*B}, \quad (14)$$

605 where the right-hand side is the output of the haplotype reconstruction, and

$$\text{var}[\mathbf{q}^B] \approx \text{var}[\mathbf{q}^{*B}], \quad (15)$$

607 where the right-hand side was calculated using the generation of the datasets $\mathbf{x}_{i,j}^{*B,P}$
608 and the inferences of the frequencies $\{\mathbf{q}^{**B}\}_j$.

610 Moving on to the post-transmission component of the compound distribution
611 in Equation 10, we can carry out the relevant marginalisations using the law of total
612 expectation and the law of total variance.

614 Given that the dynamics governing transmission and within-host growth are as-
615 sumed selectively neutral, the mean frequencies of the viral population are unchanged
616 following transmission and growth. The mean term is therefore straightforward to
617 calculate.

$$\begin{aligned} E[\mathbf{q}^{*A}] &= E[E[\mathbf{q}^{*A}|\mathbf{q}^A]] = E[\mathbf{q}^A] \\ E[\mathbf{q}^A] &= E[E[\mathbf{q}^A|\mathbf{q}^F]] = E[\mathbf{q}^F] \\ E[\mathbf{q}^F] &= E[E[\mathbf{q}^F|\mathbf{q}^B]] = E[\mathbf{q}^B] \end{aligned} \quad (16)$$

619 Thus

$$E[\mathbf{q}^{*A}] \approx \mathbf{q}^{*B} \quad (17)$$

621 Calculation of the variance requires a little more effort. The transmission event
622 can be modelled as a single multinomial draw with N^T number of trials. As a result,
623 the variance of the founder population is given by

$$\text{var}[\mathbf{q}^F|\mathbf{q}^B] = \frac{1}{N^T} M(\mathbf{q}^B), \quad (18)$$

624 where $M(\mathbf{q}) = \text{Diag}(\mathbf{q}) - \mathbf{q}\mathbf{q}^\dagger$.

626 We therefore obtain that

$$\begin{aligned} \text{var}[\mathbf{q}^F] &= E[\text{var}[\mathbf{q}^F|\mathbf{q}^B]] + \text{var}[E[\mathbf{q}^F|\mathbf{q}^B]] \\ &= E\left[\frac{1}{N^T} M(\mathbf{q}^B)\right] + \text{var}[\mathbf{q}^B] \\ &= \frac{1}{N^T} \left(E[\text{Diag}(\mathbf{q}^B)] - E[\mathbf{q}^B(\mathbf{q}^B)^\dagger] \right) + \text{var}[\mathbf{q}^B] \\ &= \frac{1}{N^T} \left(\text{Diag}(E[\mathbf{q}^B]) - \text{var}[\mathbf{q}^B] - E[\mathbf{q}^B]E[\mathbf{q}^B]^\dagger \right) + \text{var}[\mathbf{q}^B] \\ &= \frac{1}{N^T} M(E[\mathbf{q}^B]) + \left(1 - \frac{1}{N^T}\right) \text{var}[\mathbf{q}^B] \\ &\approx \frac{1}{N^T} M(\mathbf{q}^{*B}) + \left(1 - \frac{1}{N^T}\right) \text{var}[\mathbf{q}^{*B}] \end{aligned} \quad (19)$$

where we used the result

$$E[\mathbf{q}\mathbf{q}^\dagger] = \text{var}[\mathbf{q}] + E[\mathbf{q}]E[\mathbf{q}]^\dagger. \quad (20)$$

The within-host growth dynamics can be modelled as a multinomial draw of depth $N^G = \mathbf{g}N^T$ where \mathbf{g} is the growth factor. From this we obtain the result that

$$\text{var}[\mathbf{q}^A|\mathbf{q}^F] = \frac{1}{N^G} M(\mathbf{q}^F). \quad (21)$$

Marginalising over \mathbf{q}^F we obtain the variance

$$\begin{aligned} \text{var}[\mathbf{q}^A] &= E[\text{var}[\mathbf{q}^A|\mathbf{q}^F]] + \text{var}[E[\mathbf{q}^A|\mathbf{q}^F]] \\ &= E\left[\frac{1}{N^G} M(\mathbf{q}^F)\right] + \text{var}[\mathbf{q}^F] \\ &= \frac{1}{N^G} \left(E[\text{Diag}(\mathbf{q}^F)] - E[\mathbf{q}^F(\mathbf{q}^F)^\dagger] \right) + \text{var}[\mathbf{q}^F] \\ &= \frac{1}{N^G} \left(\text{Diag}(E[\mathbf{q}^F]) - \text{var}[\mathbf{q}^F] - E[\mathbf{q}^F]E[\mathbf{q}^F]^\dagger \right) + \text{var}[\mathbf{q}^F] \\ &= \frac{1}{N^G} M(E[\mathbf{q}^F]) + \left(1 - \frac{1}{N^G}\right) \text{var}[\mathbf{q}^F] \\ &\approx \frac{1}{N^G} M(\mathbf{q}^{*B}) + \left(1 - \frac{1}{N^G}\right) \left(\frac{1}{N^T} M(\mathbf{q}^{*B}) + \left(1 - \frac{1}{N^T}\right) \text{var}[\mathbf{q}^{*B}] \right) \\ &= \frac{N^T + N^G - 1}{N^T N^G} M(\mathbf{q}^{*B}) + \frac{N^T N^G - N^T - N^T + 1}{N^T N^G} \text{var}[\mathbf{q}^{*B}] \\ &\equiv \gamma M(\mathbf{q}^{*B}) + \delta \text{var}[\mathbf{q}^{*B}], \end{aligned} \quad (22)$$

where we define $\gamma = \left(\frac{N^T + N^G - 1}{N^T N^G}\right)$ and $\delta = \frac{N^T N^G - N^T - N^T + 1}{N^T N^G}$.

Finally we have that

$$\begin{aligned} \text{var}[\mathbf{q}^{*A}] &= E[\text{var}[\mathbf{q}^{*A}|\mathbf{q}^A]] + \text{var}[E[\mathbf{q}^{*A}|\mathbf{q}^A]] = E[\text{var}[\mathbf{q}^A]] + \text{var}[\mathbf{q}^A] \\ &= \text{var}[\mathbf{q}^{*A}] + \gamma M(\mathbf{q}^{*B}) + \delta \text{var}[\mathbf{q}^{*B}]. \end{aligned} \quad (23)$$

Together, Equations 16 and 23 define the mean and variance of a multivariate normal distribution representing the post-transmission component of the likelihood in Equation 10. Given our inferences for \mathbf{q}^{*B} and \mathbf{q}^{*A} , we optimised the likelihood with respect to N^T , generating a maximum likelihood estimate for the bottleneck size. We note that our approximation of the likelihood in terms of a multivariate normal distribution, works best where individual haplotype frequencies are not too close to zero or one, and where N^T is large. However, the approach allows for rapid calculation. In this sense we say that the compound method is optimised for large N^T .

Correction for the extinction of haplotypes in the compound method Where a haplotype goes extinct in the transmission process, the likelihood function of the compound method can provide a poor estimate to the correct value. In this special case, relevant in our simulated data, we used a conditional distribution approach to make a correction to the likelihood.

In the above approximation we generated a multivariate normal distribution for \mathbf{q}^{*A} :

$$\mathbf{q}^{*A} \sim \mathcal{N}(\mathbf{q}^{*B}, \text{var}[\mathbf{q}^{*A}]). \quad (24)$$

In this context, we split the vector \mathbf{q}^{*A} into \mathbf{q}_1^{*A} and \mathbf{q}_2^{*A} , the latter containing all haplotypes post-transmission with a frequency lower than the threshold frequency

η , which were considered to have died out during transmission, with the former containing the 'surviving' haplotypes. Rows and columns of the vectors and matrices were rearranged to put equation 24 into the form

$$\begin{bmatrix} \mathbf{q}_1^{*A} \\ \mathbf{q}_2^{*A} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{q}_1^{*B} \\ \mathbf{q}_2^{*B} \end{bmatrix}, \begin{bmatrix} \text{var}[\mathbf{q}^{*A}]_{11} & \text{var}[\mathbf{q}^{*A}]_{12} \\ \text{var}[\mathbf{q}^{*A}]_{21} & \text{var}[\mathbf{q}^{*A}]_{22} \end{bmatrix} \right). \quad (25)$$

The frequencies of the components of the vectors were renormalised, such that $\mathbf{q}_{2i}^{*A} = \mathbf{q}_{2i}^{*B} = 0$, while $\sum_i \mathbf{q}_{1i}^{*A} = \sum_i \mathbf{q}_{1i}^{*B} = 1$.

We obtain the result that the conditional distribution of \mathbf{q}_1^{*A} has the mean

$$\mu = \mathbf{q}_1^{*B} + \text{var}[\mathbf{q}^{*A}]_{12}(\text{var}[\mathbf{q}^{*A}]_{22})^{-1}(-\mathbf{q}_2^{*B}), \quad (26)$$

and covariance matrix

$$\Sigma = \text{var}[\mathbf{q}^{*A}]_{11} - \text{var}[\mathbf{q}^{*A}]_{12}(\text{var}[\mathbf{q}^{*A}]_{22})^{-1}\text{var}[\mathbf{q}^{*A}]_{21}. \quad (27)$$

Using these parameters to define a Gaussian distribution, we calculated the likelihood of a bottleneck N^T given the data for the surviving haplotypes represented by \mathbf{q}_1^{*A} .

To account for the haplotypes which became extinct during transmission, we made the assumption that these died out at the point of transmission to the founder population, the rapid growth of the founder population ensuring that no haplotypes went extinct through genetic drift, and viral sequencing of a large number of viral particles ensuring that no haplotypes were missed by the sequencing process. Under this assumption the likelihood of extinction is given by the simple binomial likelihood

$$\log \left[\left(1 - \sum_i \mathbf{q}_{2i}^{*B} \right)^{N^T} \right]. \quad (28)$$

Summing the log likelihoods calculated for the surviving and the extinct haplotypes gave the total likelihood of the bottleneck size N^T ; the maximum likelihood value was identified via a simple optimisation process. To prevent nonsensical outcomes at very low bottleneck sizes, we further imposed the constraint that N^T could not be less than the number of haplotypes which survived transmission.

Explicit method for bottleneck estimation The explicit method uses the inferred haplotype frequencies for the population before transmission to reconstruct the space of possible outcomes in the recipient individual. Given our inferred haplotype frequencies q_i^{B*} , we assume that N^T viruses are transmitted. The probability that the founding viral population includes n_i copies of the haplotype i , where $\sum_i n_i = N^T$, is given by

$$P(n_1, n_2, \dots, n_k | \mathbf{q}^{B*}) = \binom{N^T}{n_1 n_2 \dots n_k} \prod_i (q_i^{B*})^{n_i}, \quad (29)$$

where the first term in the right-hand side of the equation is the multinomial coefficient.

For each possible outcome $\{n_i\}$ of this multinomial process, we obtained an inference of the haplotype composition $\{q_i^A\}$ of the transmitted population given the relationship $q_i^A = n_i/N^T$ for each haplotype i . We then calculated the raw likelihood of observing the partial haplotype data collected post-transmission given this composition

using the Dirichlet multinomial formulation described above, summing likelihoods over the possible outcomes of the initial transmission.

$$\sum_{\substack{n_1, \dots, n_k \\ \sum n_i = N^T}} P(n_1, n_2, \dots, n_k | \mathbf{q}^{B^*}, N^T) \left[\exp \left(\sum_l \log \mathcal{L}_D(\mathbf{x}_l^{A,P} | T_l, \mathbf{q}^A, C) \right) \right]. \quad (30)$$

In this way we evaluate the likelihood of the bottleneck size N^T given the inferred pre-transmission haplotypes \mathbf{q}^B and the observed sequence data \mathbf{x}^A ; this is in contrast to the compound method, which is based on \mathbf{q}^B and \mathbf{q}^A . We note that this approach neglects an explicit accounting for within-host growth of the population. Different assumptions about the dynamics of early viral infection can lead to changes in inferred bottleneck sizes (17); we are not confident that the biological reality of this phenomenon is well understood. Modifications to the the Dirichlet multinomial distribution could potentially be used in this context; increasing the variance of the likelihood function would soften the effect of small changes in the underlying population.

This approach has both the advantage and the disadvantage of explicitly representing the full set of all possible multinomial outcomes of transmission. While in this sense it remains close to the biological reality, it rapidly becomes computationally expensive as the number of haplotypes k increases and as N^T becomes large. For this reason we propose it as being optimal for small values of N^T .

We note that, in our application to data from a transmission study presented here, the case in which a high bottleneck was inferred involved very limited diversity within viral segments; this facilitated the application of this method to consider larger bottleneck sizes.

Generation of simulated data Simulated data were generated using a simplified model of influenza transmission. Viruses were generated to have eight independent segments, of lengths equal to the segments of the A/H1N1 influenza virus. Each segment had five uniformly distributed polymorphic loci, making a theoretical total of 32 full haplotypes. Six haplotypes were chosen from this set under the constraint that each of the five loci had to remain polymorphic. The frequencies of these haplotypes were then randomly generated under the constraint of a minimum haplotype frequency of 5%, matching the parameters used in a previous study (17). We note that, in the reconstruction of haplotypes, our code is likely not to identify very low frequency haplotypes in the population due to the parsimony-driven approach.

Each transmission event was modelled as a simple multinomial draw, selecting a number of viruses equal to the bottleneck size from the donor population. Within-host growth was then modelled as a second multinomial draw, conferring a 22-fold increase in the population size (30). Partial haplotype data were generated from simulated short reads of each viral segment. Short reads with lengths derived from the dataset of a recent influenza study (31) were generated (mean read length = 119.68, SD read length = 136.88, mean gap length = 61.96, SD gap length = 104.48, total read depth = 102825), these reads being used to calculate the number of reads spanning each set of consecutive polymorphisms in each segment. Given these numbers, partial haplotype observations were generated using a Dirichlet multinomial sampling process.

An inference of the transmission bottleneck was carried out independently using simulated data from each viral segment. These inferences were then combined, summing the log likelihoods across different segments to obtain an overall maximum likelihood estimate. Within our simulated data a small number of cases were identified in which the entire post-transmission population in a segment was inferred to comprise a haplotype that was not present above the cutoff frequency in the pre-transmission population, equivalent to a case where a haplotype arose *de novo* in the population and swept to fixation before data could be collected. In such cases, data for the segment in question were ignored, calculating the transmission bottleneck across the remaining segments.

Processing of sequence data Our method was applied to data from a recent study of influenza transmission among individuals in households (3). Data from transmission pairs identified in this study were aligned using the BWA software package (28) then filtered using SAMFIRE (25) to remove reads with a median PHRED score below 30, and to mask nucleotides with a PHRED score below this value. Following the original study, sites in coding regions of the virus were then called at an allele frequency cutoff of 2%, following which reads were divided into sets of partial haplotype data.

Data describing the within-host evolution of influenza were used to evaluate the extent of noise in the dataset. Noise in data arises both from the non-representative sampling of viruses from the host and from the subsequent experimental steps used to generate sequence data (27); an over-estimate of the extent of noise in data can lead to substantial errors in the inference of a transmission bottleneck (17). We here took a heuristic approach applied in a previous study (17). In a first step, data from all within-host single-locus trajectories were used to generate a provisional estimate of the extent of the noise in the data. Next, trajectories which under this estimate evolved in a manner consistent with selective neutrality were identified. Models of selective neutrality (constant allele frequency), constant selection ($dq/dt = sq(1 - q)$), and time-dependent selection (exact match to observed frequencies) were fitted to the data using the Dirichlet multinomial model of Equation 4, requiring a difference of 10 units of BIC to favour the more complex model. Trajectories identified as neutral under this method were used to produce a final estimate of noise in the data; we inferred the parameter $C = 660$. Data from 43 putative transmission events were evaluated.

The estimate of an effective read depth for the case in which a very high bottleneck was inferred was conducted using SAMFIRE based upon allele frequency data, and using a cutoff frequency for minority alleles of 2%.

AVAILABILITY OF CODE

Code and data used or generated during this project is available from <https://github.com/cjri/VeTrans>.

ACKNOWLEDGMENTS

This work was supported by a Sir Henry Dale Fellowship, jointly funded by the Wellcome Trust (wellcome.ac.uk) and the Royal Society (royalsociety.org) with grant number 101239/Z/13/Z. CKL was funded by a Wellcome Trust Studentship with grant number 105365/Z/14/Z. CJRI acknowledges a visiting fellowship from the University of Helsinki. DBW was funded by a Simons Investigator award from the Simons Foundation with grant number 508600.

REFERENCES

1. **Sidorenko Y, Reichl U.** 2004 Oct. Structured model of influenza virus replication in MDCK cells. *Biotechnol. Bioeng.* 88(1):1–14.
2. **Zwart MP, Elena SF.** 2015 Nov. Matters of Size: Genetic Bottlenecks in Virus Infection and Their Potential Impact on Evolution. *Annu. Rev. Virol.* 2(1):161–179.
3. **McCrone JT, Woods RJ, Martin ET, Malosh RE, Monto AS, Lauring AS.** 2018 Aug. Stochastic processes constrain the within and between host evolution of influenza virus. *eLife* 7:e35962.
4. **Biek R, Pybus OG, Lloyd-Smith JO, Didelot X.** 2015 Jun. Measurably evolving pathogens in the genomic era. *Trends Ecol. & Evol.* 30(6):306–313.
5. **Stack JC, Murcia PR, Grenfell BT, Wood JLN, Holmes EC.** 2012 Nov. Inferring the inter-host transmission of influenza A virus using patterns of intra-host genetic variation. *Proc. Royal Soc. B: Biol. Sci.* 280(1750):20122173–20122173.
6. **Worby CJ, Lipsitch M, Hanage WP.** 2017 Jun. Shared Genomic Variants: Identification of Transmission Routes Using Pathogen Deep-Sequence Data. *Am. J. Epidemiol.* 186(10):1209–1216.
7. **Sacristan S, Malpica JM, Fraile A, Garcia-Arenal F.** 2003 Aug. Estimation of Population Bottlenecks during Systemic Movement of Tobacco Mosaic Virus in Tobacco Plants. *J. Virol.* 77(18):9906–9911.
8. **Krimbas CB, Tsakas S.** 1971 Sep. The Genetics of *Dacus Oleae*. V. Changes of esterase polymorphism in a natural population following insecticide control—selection or drift? *Evolution* 25(3):454–460.
9. **Monsion B, Froissart R, Michalakis Y, Blanc S.** 2008 Oct. Large bottleneck size in Cauliflower Mosaic Virus populations during host plant colonization. *PLoS Pathog.* 4(10):e1000174.
10. **Khiabani H, Emmett KJ, Lee A, Rabadan R.** 2015. High-resolution Genomic Surveillance of 2014 Ebola virus Using Shared Subclonal Variants. *PLoS currents* 7:1–17.
11. **Sobel Leonard A, Weissman DB, Greenbaum B, Ghedin E, Koelle K.** 2017 Jun. Transmission Bottleneck Size Estimation from Pathogen Deep-Sequencing Data, with an Application to Human Influenza A Virus. *J. Virol.* 91(14):e00171–17–19.
12. **Varble A, Albrecht RA, Backes S, Crumiller M, Bouvier NM, Sachs D, García-Sastre A, tenOever BR.** 2014 Nov. Influenza A Virus Transmission Bottlenecks Are Defined by Infection Route and Recipient Host. *Cell Host Microbe* 16(5):691–700.
13. **Frise R, Bradley K, van Doremalen N, Galiano M, Elderfield RA, Stilwell P, Ashcroft JW, Fernandez-Alonso M, Miah S, Lackenby A, Roberts KL, Donnelly CA, Barclay WS.** 2016 Jul. Contact transmission of influenza virus between ferrets imposes a looser bottleneck than respiratory droplet transmission allowing propagation of antiviral resistance. *Sci. Reports* 6(1):29793.
14. **Poon LLM, Song T, Rosenfeld R, Lin X, Rogers MB, Zhou B, Sebra R, Halpin RA, Guan Y, Twaddle A, DePasse JV, Stockwell TB, Wentworth DE, Holmes EC, Greenbaum B, Peiris JSM, Cowling BJ, Ghedin E.** 2016 Feb. Quantifying influenza virus diversity and transmission in humans. *Nat. Genet.* 48(2):195–200.
15. **Xue KS, Bloom JD.** 2019 Feb. Reconciling disparate estimates of viral genetic diversity during human influenza infections. *Nat. Genet.* 51(9):1298–1301.
16. **Illingworth CJ, Mustonen V.** 2011 Nov. Distinguishing Driver and Passenger Mutations in an Evolutionary History Categorized by Interference. *Genetics* 189(3):989–1000.
17. **Lumby CK, Nené NR, Illingworth CJ.** 2018 Oct. A novel framework for inferring parameters of transmission from viral sequence data. *PLoS Genet.* 14(10):e1007718.
18. **Stephens M, Donnelly P.** 2003 Nov. A Comparison of Bayesian Methods for Haplotype Reconstruction from Population Genotype Data. *The Am. J. Hum. Genet.* 73(5):1162–1169.
19. **Zhang K, Sun F, Zhao H.** 2005 Jan. HAPLORE: a program for haplotype reconstruction in general pedigrees without recombination. *Bioinformatics* 21(1):90–103.
20. **Prosperi MCF, Salemi M.** 2011 Nov. QuRe: software for viral quasispecies reconstruction from next-generation sequencing data. *Bioinformatics* 28(1):132–133.
21. **Zagordi O, Bhattacharya A, Eriksson N, Beerenwinkel N.** 2011. ShoRAH: estimating the genetic diversity of a mixed sample from next-generation sequencing data. *BMC Bioinforma.* 12(1):119.
22. **Selvaraj S, Dixon JR, Bansal V, Ren B.** 2013 Nov. Whole-genome haplotype reconstruction using proximity-ligation and shotgun sequencing. *Nat. Biotechnol.* 31(12):1111–1118.
23. **Giallonardo FD, Töpfer A, Rey M, Prabhakaran S, Dupont Y, Lee-mann C, Schmutz S, Campbell NK, Joos B, Lecca MR, Patrignani A, Däumer M, Beisel C, Rusert P, Trkola A, Günthard HF, Roth V, Beerenwinkel N, Metzner KJ.** 2014 Aug. Full-length haplotype reconstruction to infer the structure of heterogeneous virus populations. *Nucleic Acids Res.* 42(14):e115–e115.
24. **Illingworth CJR.** 2015 Nov. Fitness Inference from Short-Read Data: Within-Host Evolution of a Reassortant H5N1 Influenza Virus. *Mol. Biol. Evol.* 32(11):3012–3026.
25. **Illingworth CJR.** 2016 Jul. SAMFIRE: multi-locus variant calling for time-resolved sequence data. *Bioinformatics* 32(14):2208–2209.
26. **Valesano AL, Fitzsimmons WJ, McCrone JT, Petrie JG, Monto AS, Martin ET, Lauring AS.** 2019 Oct. Influenza B viruses exhibit lower within-host diversity than influenza A viruses in human hosts. *bioRxiv* 4(10):e05055–33.
27. **Illingworth CJR, Roy S, Beale MA, Tutill H, Williams R, Breuer J.** 2017 Nov. On the effective depth of viral sequence data. *Virus Evol.* 3(2):1–9.
28. **Li H, Durbin R.** 2009 Jul. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25(14):1754–1760.
29. **Charlesworth B.** 2009 Mar. Fundamental concepts in genetics: Effective population size and patterns of molecular evolution and variation. *Nat. Rev. Genet.* 10(3):195–205.
30. **Baccam P, Beauchemin C, Macken CA, Hayden FG, Perelson AS.** 2006 Jul. Kinetics of Influenza A Virus Infection in Humans. *J. Virol.* 80(15):7590–7599.
31. **Wilker PR, Dinis JM, Starrett G, Imai M, Hatta M, Nelson CW, O'Connor DH, Hughes AL, Neumann G, Kawaoka Y, Friedrich TC.** 2013 Oct. Selection on haemagglutinin imposes a bottleneck during mammalian transmission of reassortant H5N1 influenza viruses. *Nat. Commun.* 4:1–11.