1 **Title:** Number of HIV-1 founder variants is determined by the recency of the source partner

2 infection

3

4 **Authors**: Ch. Julián Villabona-Arenas[1,2], Matthew Hall[3], Katrina A. Lythgoe[3], Stephen G.

5 Gaffney[4], Roland R. Regoes[5], Stéphane Hué[1,2], Katherine E. Atkins[1,2,6] *

6 **Affiliations:**

7 [1]Department of Infectious Disease Epidemiology, Faculty of Epidemiology and Population

8 Health, London School of Hygiene and Tropical Medicine, London, UK

9 [2]Centre for Mathematical Modelling of Infectious Diseases, London School of Hygiene and

10 Tropical Medicine, London, UK

11 [3]Big Data Institute, Nuffield Department of Medicine, University of Oxford, Oxford, UK

12 [4]Division of Biostatistics, Yale School of Public Health, New Haven, USA

13 [5]Institute of Integrative Biology, Department of Environmental Systems Science, ETH Zurich,

14 Zurich, Switzerland

15 [6]Centre for Global Health, Usher Institute of Population Health Sciences and Informatics, The

16 University of Edinburgh, Edinburgh, UK

17 * corresponding author: Katherine.Atkins@ed.ac.uk

18 **One sentence summary**: Multiple founder variant transmission of HIV-1 increased during early

19 infection.

20

21  **Abstract**

22  During sexual transmission, the large genetic diversity of HIV-1 within an individual is

23  frequently reduced to one founder variant that initiates infection. Understanding the drivers of

24  this bottleneck is crucial to develop effective infection control strategies. Little is known about

25  the importance of the source partner during this bottleneck.  To test the hypothesis that the

26  source partner affects the number of HIV founder variants, we developed a phylodynamic model

27  calibrated using genetic and epidemiological data on all existing transmission pairs for whom the

28  direction of transmission and the infection stage of the source partner are known. Our results

29  suggest that acquiring infection from someone in the acute (early) stage of infection increases the

30  risk of multiple founder variant transmission when compared with someone in the chronic (later)

31  stage of infection. This study provides the first direct test of source partner characteristics to

32  explain the low frequency of multiple founder strain infections.

33

34

35 **Main Text**

36 Sexual transmission of HIV-1 results in a viral diversity bottleneck due to physiological barriers

37 as well as viral or cellular constraints that prevent most genetic variants within the source partner

38 from establishing onward infection (*1–3*). Indeed, this diversity bottleneck results in around three

39 quarters of new infections being founded by a single genetic variant (*4–9*). The extent of genetic

40 diversity transmitted to a new partner is a crucial determinant in understanding the efficacy of

41 putative vaccines and may shed light on the transmission of drug resistance to treatment naive

42 individuals.

43

44 The factors leading to the diversity bottleneck during sexual transmission can be broadly

45 categorized as those determined by the source partner—such as viral load and viral diversity

46 available for transmission *(10)*, those determined by the recipient partner—such as target cell type

47 and availability in the genital or rectal mucosa (e.g. (*3*, *11*, *12*)), and those connected with viral

48 characteristics—such as glycosylation profiles and cell tropism (reviewed in (*13*)). While the

49 impact of the recipient partner and the characteristics of transmitted founder variants have been

50 widely discussed, little is known about how the source partner affects the viral diversity bottleneck.

51 Modelling work suggests that infection stage of the source partner at the point of onward

52 transmission may be a key driver in determining the number of transmitted variants (*14*). However,

53 there is currently no empirical evidence to suggest how the infection stage of the source partner

54 influences the viral diversity bottleneck. This gap has arisen because analyses are routinely

55 conducted on individuals without information on the partner from whom they acquired infection.

56 Phylogenetic analyses now offer a possible solution to this impasse.
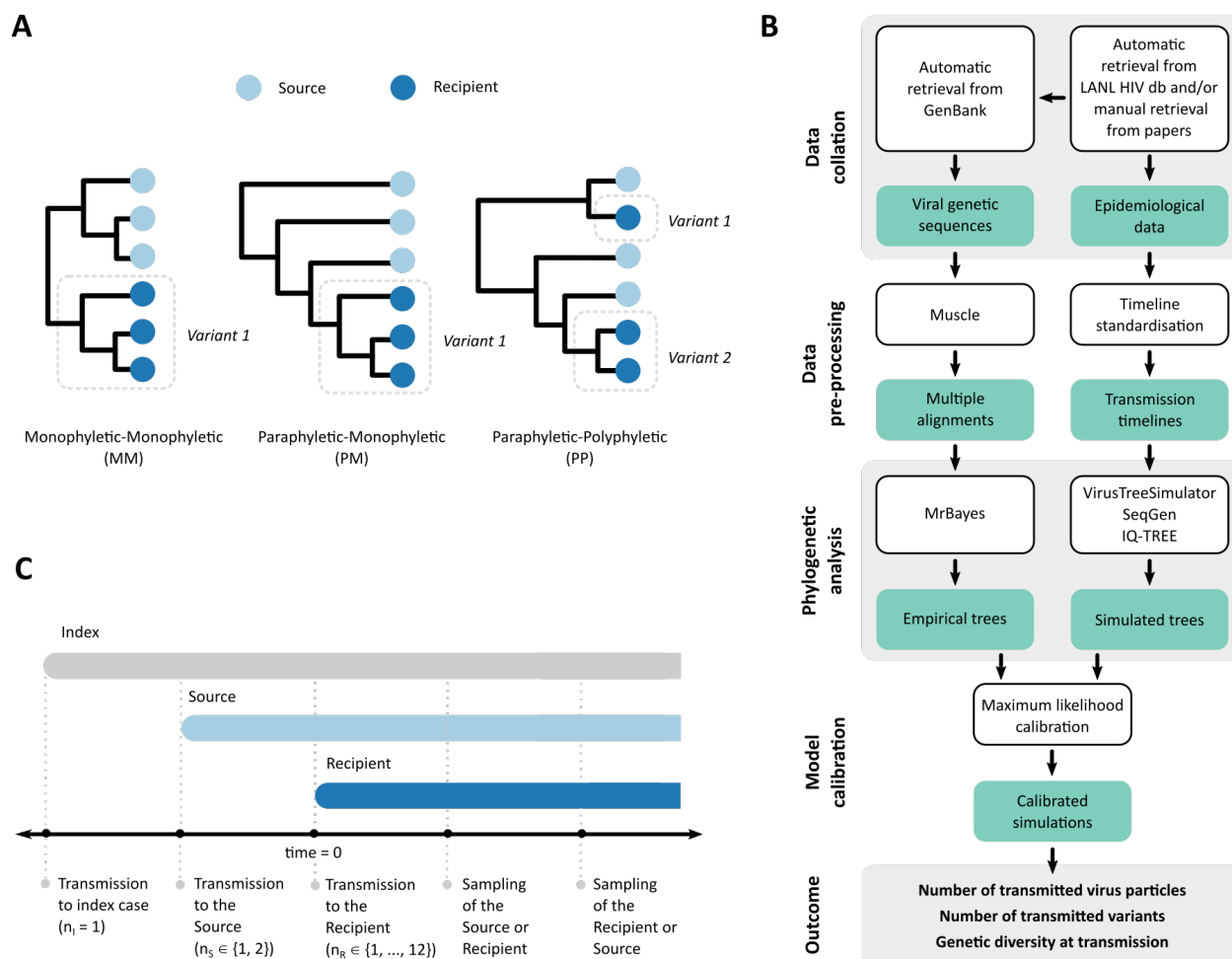
57

58    Phylogenetic trees are representations of the ancestral relationships of organisms with the tips of

59    the tree representing those that are sampled, the internal nodes representing their inferred

60    common ancestors, and the branches as the evolutionary pathways between these actual and

61    inferred individuals. When phylogenetic trees are constructed using sequence data from both

62    partners in an HIV transmission pair, the relationship between the evolutionary histories of both

63    sets of viral samples may reflect epidemiological relationships between the two individuals (*15-*

64    *17*). Previous modelling studies suggest that the evolutionary histories of the viral populations in

65    both partners can provide important information, such as the direction of transmission (*15*) and

66    the number of transmitted founder variants (*18*). For this, each putative transmission pair can be

67    classified into one of three 'topologies' that defines the evolutionary relationship between the

68    viral populations of the two partners: monophyletic-monophyletic (*MM*, where the sequences

69    from each partner form separate groups), paraphyletic-monophyletic (*PM*, where the sequences

70    from one partner are embedded in the sequences from the second partner), or a combination of

71    paraphyletic and polyphyletic (*PP*, where sequences from both partners are interspersed) (**Fig.**

72    **1A**). The number of monophyletic clusters in a PM (one) or PP (more than one) tree can be

73    interpreted as the minimum number of transmitted founder variants. In practice, however, many

74    factors may influence epidemiological interpretations from phylogenetic trees such as sampling

75    times, sampling density of the viral populations and phylogenetic signal (*19, 20*).

76

77    Here we present a data-driven phylodynamic approach to overcome these empirical and

78    methodological issues to evaluate the impact of the source partner's infection stage and route of

79    exposure on the HIV diversity bottleneck (**Fig. 1B, C**). We first retrieved all available genetic

80    and epidemiological information from published HIV sexual transmission pairs where the

81   direction of transmission is known, and kept for further analysis those pairs for whom

82   transmission could be classified as having occurred in the source partner's acute stage (≤90 days

83   after his/her infection) or chronic stage (later than 90 days after his/her infection). After further

84   stratifying pairs into heterosexual (HET) and men-who-have-sex-with-men (MSM) risk groups,

85   we found a significant difference in the timing of transmission between the two risk groups.

86   Specifically, 10 of 36 MSM pairs were the result of acute stage transmission compared with 1 of

87   76 of HET pairs (**Fig. 2**).



89   **Fig. 1: Methods schematics.** A) Phylogenetic tree topology class of known transmission pairs that have

90   previously been used as a proxy for calculating the minimum number of founder variants transmitted to

91    the recipient: trees of class MM and PM both suggest a minimum of one founder variant while trees of

92    class PP suggest a multiple founder variants, with the minimum number of founder variants being the

93    number of recipient clades embedded in PP trees (here shown as two). B) Pipeline of phylodynamic

94    analysis (LANLdb, Los Alamos National Laboratory HIV sequence database) where teal represents data

95    or analysis output and white represents methods and analysis. An example of a standardised transmission

96    timeline for a known source-recipient pair is provided in panel C. C) Schematic of the transmission pair

97    model simulation that shows the transmission and sampling timelines. The simulated number of virus

98    particles transmitted to the index case, and the source and recipient partners ($n_I, n_S, n_R$ respectively) are

99    shown on the transmission events timeline.

100

101   We then performed Bayesian phylogenetic tree reconstruction on the genetic sequences of the

102   transmission pairs and classified the topology class of each tree in the posterior distribution as

103   monophyletic-monophyletic (MM), paraphyletic-monophyletic (PM) or paraphyletic-

104   polyphyletic (PP).  The most likely topology class was PM (65% and 61% for HET and MSM,

105   respectively), but with a higher number of PP trees in the MSM group ($P$=0.056, **Fig. 2**). This

106   result has previously been reported as indicative of a higher number of founder variants for

107   MSM (*18*). However, when we stratify the topology class by whether the source partner was in

108   acute or chronic infection at the time of transmission, our results indicate that the infection stage

109   of the source is the primary driver for any observed differences in topology class. Specifically,

110   there is no difference between the HET and MSM groups in the PM/PP topology class ratio

111   when transmission occurs in the chronic stage of infection ($P$=0.570). Note that only one HET

112   transmission occurs during the acute stage, and the topology class for this pair is PP. These

113   results remain qualitatively consistent when only data were analysed from the 66% of

transmission pairs for whom the posterior trees gave a certainty of over 95% for the most

frequent topology class (**Fig. S3**). These results indicate that infection stage of the source partner,

and not risk group *per se,* influences the diversity bottleneck at transmission.



**Fig. 2: Phylogenetic findings from the empirical transmission pairs**. Fraction of phylogenetic tree

topology class (MM: Monophyletic-Monophyletic, PM: Paraphyletic-Monophyletic and PP: Paraphyletic-

Polyphyletic) where each tree topology class is classified as the most frequent topology class of each

posterior distribution per transmission pair. Results are stratified by risk group: 76 heterosexual (HET)

pairs and 36 men-who-have-sex-with-men (MSM) pairs) and infection stage of the source partner at

transmission (11 acute pairs defined as <90d post infection and 101 chronic pairs defined as ≥90d post

infection).

To test whether these empirical findings are indicative of a smaller diversity bottleneck in the

chronic stage of HIV infection, we developed a phylodynamic framework in which we simulated

the epidemiologic characteristics of each HET and MSM transmission pair, the timing of their

sequence sampling, the transmission of virus particles, and the within-host genetic evolution in

130    both the source and recipient (**Fig. 1B**). Specifically, using the epidemiological information from

131    the transmission pairs, we simulated phylogenies under a coalescent model before generating

132    genetic sequences from these simulations and performing Maximum Likelihood (ML)

133    phylogenetic reconstruction on these simulated sequences. We classified each of these simulated

134    trees as MM, PM or PP and determined the frequency of each topology class (*i.e.* the fraction of

135    simulated trees that are classified as MM, PM and PP) for each simulated transmission pair

136    across all the simulated sequences.  However, as we could not directly observe the number of

137    virus particles that are transmitted between source and recipient, we repeated the simulation of

138    phylogenetic trees for each transmission pair under a range of plausible values of virus particles

139    transmitted. By fitting the simulation output topology class distribution to the topology class

140    distribution from the empirical phylogenetic trees using maximum likelihood inference, we then

141    determined the most likely number of transmitted virus particles for each transmission pair and

142    used this best fit model for further analysis.  Note that two or more virus particles may have the

143    same genetic sequence and would constitute a single founder variant (or haplotype), discussed

144    later. Further, due to the analysis conditioning on extant lineages, we use the term 'founder

145    variants' to describe those transmitted variants that found detectable viral lineages, thereby

146    ignoring variants that are transmitted but the lineages of which become extinct.

147    Our fitting procedure selects a best fit model that clearly delineates between transmission pairs

148    between whom one virus particle is transmitted (75% of pairs) and those between whom more

149    than one virus particle is transmitted (25% of pairs, **Fig. 3A**). While there is a high degree of

150    confidence in the result when one particle is transmitted, there is often uncertainty around the

151    exact number when multiple particles are transmitted (**Fig. 3A**). Importantly, we found acute

152    stage transmissions are more likely to lead to multiple particle infections compared with chronic

stage transmissions (73% vs. 20%, $P = 0.0005$). The topology class of the simulated

phylogenetic trees is strongly influenced by the number of virus particles being transmitted (**Fig.**

**3B**). PM trees are more commonly found in the pairs that are better described by a model with a

single transmitted virus particle (81%) whereas PP trees appeared more often when multiple

particles are likely to have been transmitted (86%).



Fig. 3: **The estimated number of transmitted virus particles for the 112 transmission pairs.** The

estimates of transmitted virus particles for each transmission pair were calculated by choosing the model

simulation that generated a phylogenetic tree topology class distribution (that is, the number of MM, PM

and PP trees constructed from the simulated genetic sequences) that best matched the topology class

distribution from the phylogenetic trees constructed from the empirical genetic sequences.  A) Maximum

likelihood number of virus particles founding recipient infections, $n_R^*$, for each pair (stacked points) with

95% confidence intervals (lines) grouped by stage of infection (acute, 11 pairs or chronic, 101 pairs) and

risk group (76 heterosexual pairs, HET and 36 men-who-have-sex-with-men pairs, MSM). B) Maximum

168   likelihood number of virus particles founding recipient infections coloured by topology class of the

169   phylogenetic tree constructed from the simulated genetic sequences.

170

171   For each transmission pair, we then simulated the genetic sequences of the transmitted viral

172   population under the best fit virus particle model and calculated the most likely number of

173   founder variants for each transmission pair (*i.e.* the number of distinct haplotypes). The median

174   number of founder variants transmitted across all pairs is 1 (range: 1-11, **Fig. 4A**). Using the full

175   distribution of the number of transmitted founder variants for each pair, we also calculated the

176   probability that a single founder variant was transmitted to the respective recipient. Our results

177   suggest that across all pairs in both risk groups, the mean probability of observing one founder

178   variant is 0.73. Stratifying by risk group, we find there is a higher probability that one founder

179   variant founds HET infections than MSM infections (a geometric mean of 0.80 vs. 0.63, **Fig.**

180   **4B**). However, these risk group differences mostly disappear when we stratify the results by the

181   infection stage of the source. Here, for example, when only chronic stage transmissions are

182   considered, there is no difference in the probability of one founder variant between MSM

183   transmissions and HET transmissions (means of 0.80 vs 0.71, $P$=0.398), and the pairwise

184   diversity at transmission is similar between both groups (**Fig. 4C**). In contrast, when stratifying

185   solely by infection stage of the source partner, we find that transmission during the acute stage

186   has a much lower probability of one founder variant than during the chronic stage (means of 0.40

187   vs. 0.77) with a higher median number of founder variants transmitted, when only the most likely

188   number of founder variants for each pair is considered (2 vs. 1, **Fig. 4A**). Nonetheless, if multiple

189   founder variant transmission does occur, our results suggest that the number of founder variants

190    is higher during chronic stage transmission, consistent with a higher diversity measure during

191    this later stage of infection (**Fig. 4C**).

192    From these results, therefore, there is approximately double the chance of multiple founder

193    variant transmission during acute stage infection across both risk groups (relative risk = 0.52).

194    Assuming that transmission risk is weighted towards early transmission such that half of all

195    index case to source partner transmissions occur after 90 days of index case infection leads to

196    qualitatively similar results (Supplementary Materials).  Similarly, calibrating the simulation

197    model to bootstrapped samples rather than Bayesian posterior distributions leads to similar

198    results (Supplementary Materials).

199

200    Our results suggest that there is an association between tree topology class and multiple founder

201    variant transmission, with 95% of MM and PM trees being due to one founder variant (**Fig. 4D**).
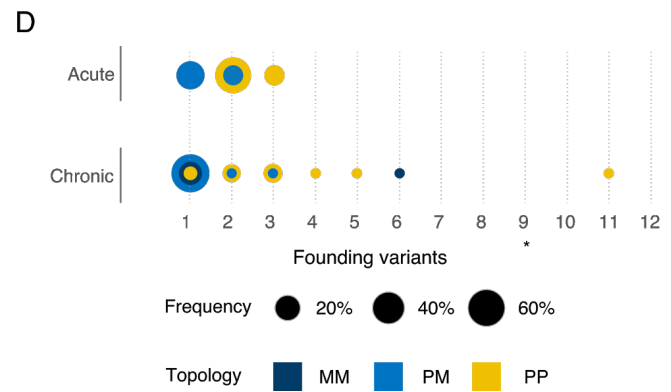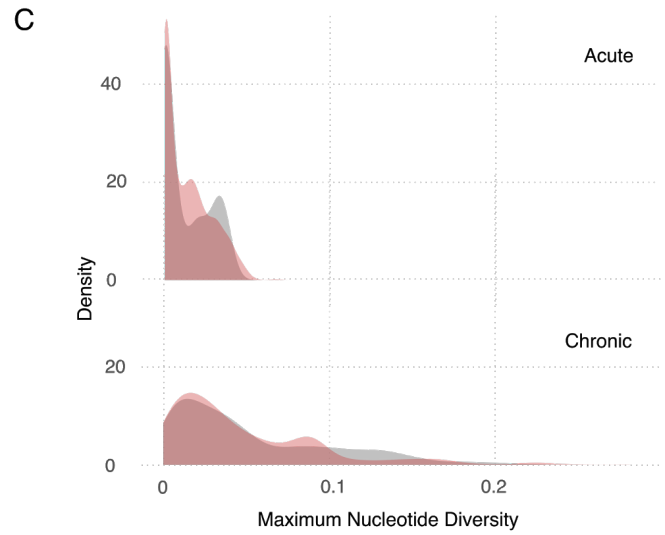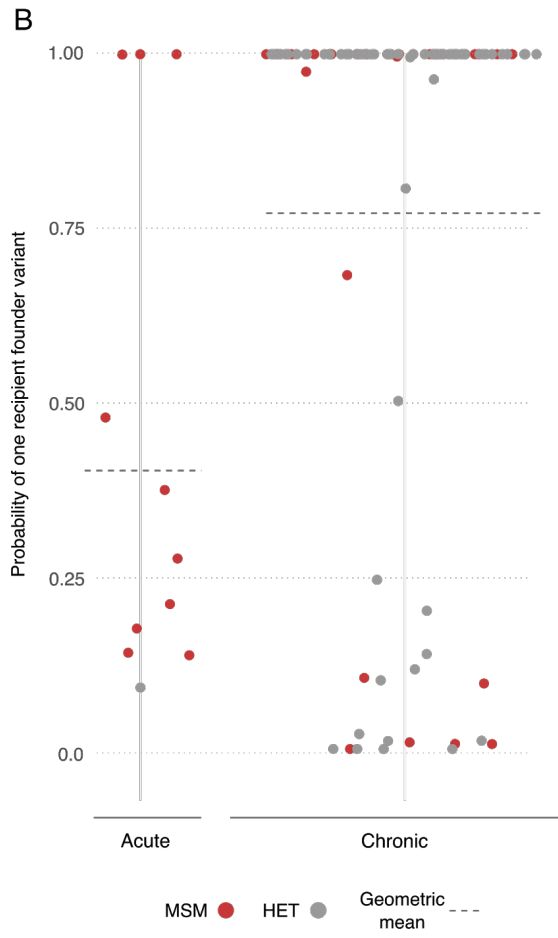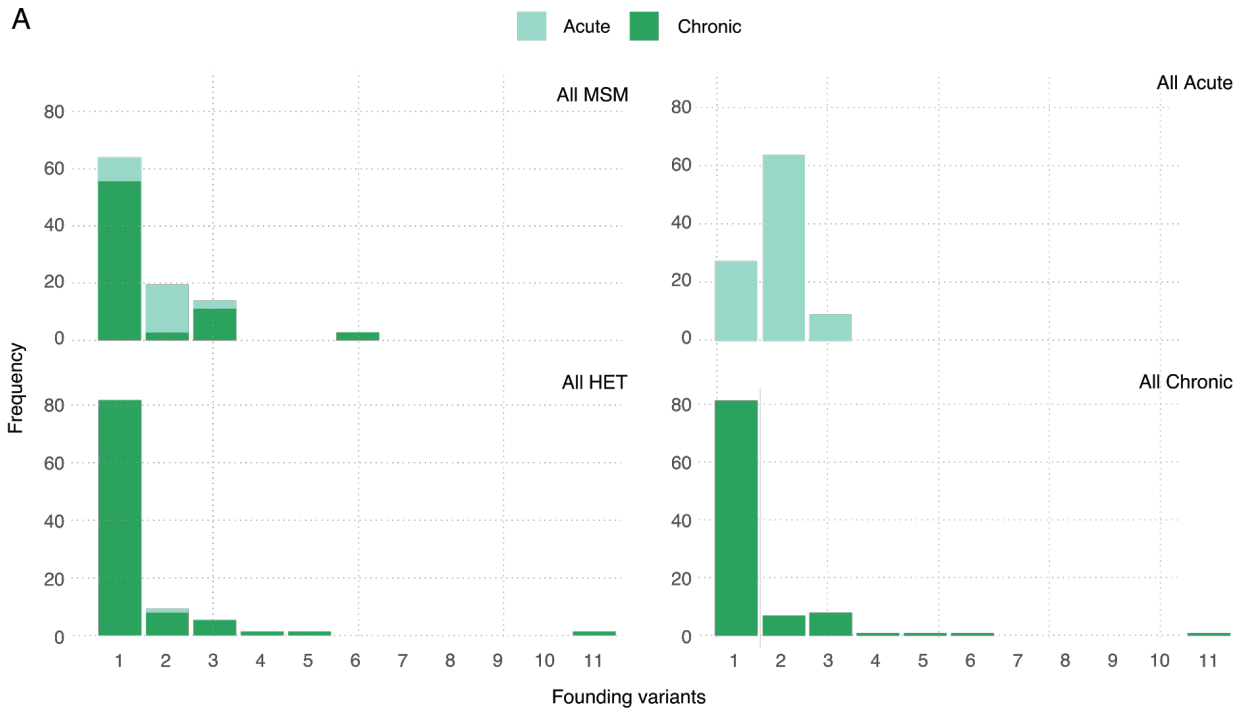
202    However, the number of embedded recipient clades is not always a proxy for the minimum

203    number of founder variants transmitted. For example, in chronic stage transmission, 11% of PP

204    topology class trees were due to single founder variant transmission (**Fig. 4D**). It is important to

205    stress that a PP topology class outcome may occur not only due to multiple genetically distinct

206    virus populations founding recipient infections but may also reflect a lack of phylogenetic signal

207    in the data; for instance, the sampled sequence lengths that gave rise to PP trees was on average

208    shorter than those for MM ($P$=0.096) and PM ($P$=0.004). Across both infection stages, we find

209    that if MM, PM or PP is assigned as the most likely tree topology class, then 92%, 96% and 15%

210    of transmissions are due to a single founder variant, respectively.

211

212    We have used a combination of empirical data and phylodynamic model simulation to evaluate

213    the role of infection stage at transmission and route of transmission on the number of virus

214    particles transmitted during sexual HIV exposure. This makes three important advances on

215    previous work. First, it is the first empirically-based study that fits a model to data to understand

216    the role of the source partner in multiple founder variant transmission. Second, while we use

217    previously developed topology classification of phylogenetic trees to understand HIV

218    transmission pairs, we extend this methodology by calibrating a phylodynamic model to

219    empirical data. This new approach provides a means to validate the untested assumption that the

220    number of embedded recipient partner lineages in a phylogenetic tree directly corresponds to the

221    minimum number of founder variants transmitted. Third, our phylodynamic model explicitly

222    incorporates virus particle number and the identity of genetic sequences. This advance produces

223    results that contrast with previous work that has shown the number of founder variants has little

224    impact on the topology class of the phylogenetic tree when only overall genetic diversity, rather

225    than sequence identity, is tracked (*15*).

226

227    The relative importance of acute and chronic stages of HIV in determining both the number of

228    virus particles and the number of founder variants transmitted is consistent with a recent

229    modelling study (*14*). However, our study finds higher proportions of infections initiated by

230    multiple founder variants overall during these two stages. This difference is likely due to the

231    assumptions related to how the stages of infection are defined as well as the relative importance

232    of transmission during late infection. Specifically, the previous modelling study finds that two

233    thirds of multiple founder variant transmission occurs during the pre-AIDS stage of infection

234    which is assumed to have both a high viral load and large haplotype diversity. If later stages of

235

13

Fig. 4: **Phylogenetic findings from the calibrated simulations**. A) Frequency of number of transmitted founder variants for transmission pairs by either infection stage of source partner at transmission (left) or risk group (right). The number of multiple founder variants is calculated as the modal simulated value. B) Probability of one founder variant in the recipient for each pair stratified by infection stage of the source partner at transmission. C) Probability density distribution of maximum diversity (proportion of sites that differ) in the recipient partner across all simulations with more than one haplotype stratified by infection stage of the source at transmission. D) Number of founder variants coloured by topology class of the phylogenetic tree constructed from the best fit model of the simulated genetic sequences.

infection account for disproportionately less transmission then the previous model would predict higher proportions of multiple founder variant transmission in both the acute and chronic stages of infection, becoming more consistent with empirical estimates from our analysis. By contrast, our study is agnostic about the relative importance of early and late transmission and does not differentiate between chronic and a pre-AIDS stage of infection, which cannot easily be identified through analysis of empirical data.

Data from four of the MSM transmission pairs in this study have previously been used to estimate the number of variants founding infection using a combination approach of single genome amplification (SGA), direct amplicon sequencing and mathematical modelling (*7*). Our results broadly agree with this previous analysis, with both analyses suggesting two recipients were infected with one founder variant and one recipient was infected with multiple founder variants (our analysis suggests a mean of 2-3 founder variants and the previous analysis suggests 3 founder variants); there was disagreement with results from a fourth recipient, for whom a

259    single founder variant was 13% probable in this study (with a mode of 2 founder variants) but

260    the most likely outcome in the previous analysis.  Small differences likely arise because this

261    study uses sequence data from both partners to evaluate the transmission of multiple founder

262    variants to the recipient partner. These extra data can be used to parameterize a mathematical

263    model that accounts for the evolutionary relationship between the virus samples from both

264    partners, rather than relying solely on accumulating diversity. Specifically, neglecting the extent

265    of genetic similarity between the source and recipient virus samples might misattribute

266    borderline cases of diversity accumulation.


267    Our study finds a median of one founder variant and a maximum of 11, with little difference

268    between HET and MSM risk groups. When only multiple founder variant transmissions are

269    considered, our study finds a median of 2-3 founder variants. These values are consistent with a

270    previous pooled analysis using results from four analyses that used the current gold-standard

271    SGA combination approach as above (*9*).

272    At present, the genetic determinants of  HIV-1 disease progression are not clear. However, it is

273    important to note that even small differences between genotypes can have important clinical

274    outcomes. For instance, single polymorphisms can affect replication capacity (*21*), or can lead to

275    primary non-nucleoside reverse transcriptase inhibitor resistance with different amino acids

276    changes at the same position conferring equivalent levels of resistance (*22*).

277    Previous studies have disagreed over the extent to which the elevated risk of transmission during

278    the acute stage of infection (reviewed in (*23*)) is driven by increased viral load, elevated per

279    particle transmission probability or other behavioural factors such as high rates of sex partner

280    change or concurrent partnerships (*24-29*). Here, while we find strong evidence to support the

281     fact that acute stage transmissions are characterised by more virus particles and variants

282     founding infection, this result alone cannot disentangle virus- and host-related drivers of elevated

283     transmission. For example, the higher number of variants being transmitted during acute

284     infection could arise if the number of transmissible variants declines as infection progresses or,

285     because with more particles being transmitted, there are more opportunities for multiple variants

286     to found infection (*14,30*) However, our study can shed light on the eight times elevated per-

287     exposure risk of infection that has been found for MSM relative to HET transmission (*31-32*). In

288     particular, the lack of difference in both the number of virus particles and the number of founder

289     variants that establish infection after transmission from a chronically infected source in HET and

290     MSM suggests that the observed heightened acquisition risk for MSM could in part be due to

291     sampled MSM individuals being more likely to be in the acute stage at the time of transmission

292     (*14, 27*). Whether MSM partners are more likely to be sampled earlier in infection because of

293     sampling procedures or because MSM are indeed more likely to transmit during early infection is

294     unclear. While this observation raises the possibility that the role of sexual risk group in itself

295     may have less of an impact on the transmission of multiple founder variant probability, from a

296     pragmatic perspective, if more MSM infections are indeed caused by acute stage transmissions,

297     the evolutionary and epidemiologic impact on public health will be the same irrespective of the

298     mechanism.

299

300     There are two primary limitations to acknowledge. First, our model assumes a single

301     transmission event between each source and recipient partner. Without detailed knowledge of the

302     transmission pairs, we cannot distinguish between multiple infections each with a single founder

303     variant and a single infection with multiple founder variants; if for some pairs, the former were

304 true then this might suggest an elevated transmission rate during the acute stage, as has been

305 observed previously (*28*, *29*). Second, our phylodynamic framework does not account for the

306 effect of selection and recombination. Specifically, selection, such as that for viruses which use

307 the CCR5 co-receptor (*33*), is thought to occur at the point of transmission , although the strength

308 may be dependent on the route of transmission (*34*).

309

310 Our study finds that the transmission of multiple HIV-1 founder variants is determined by

311 infection stage of the source partner, with transmission of more founder variants of HIV-1 in

312 acute compared with chronic infections. These findings stress that epidemiological or clinical

313 analysis of known transmission pairs should account for potential mediation by stage of

314 transmission when evaluating the effect of sexual risk group.

315

316 **Acknowledgements**

327 their respective repositories: data on the transmission pairs and sequence alignments

328 (TransmissionPairs_Data), code for retrieval of transmission pair epidemiological data and

329 metadata from Los Alamos National Laboratory HIV sequence database

330 (TransmissionPairs_LANLRetrieval), code for sequence retrieval from GenBank

331 (TransmissionPairs_GenBankRetrieval), code for phylodynamic analysis

332 (TransmissionPairs_PhylodynamicAnalysis), and code for topological classification

333 (TransmissionPairs_TreeTopologyAnalysis).

334 **List of Supplementary Materials**

335 Materials and Methods

336 Supplementary Text

337 Figs. S1 to S5

338 Data S1 to S4

339 References (*35-46*)

340 Reproducibility Checklist

341

342 **References and Notes**

343 1.  J. L. Geoghegan, A. M. Senior, E. C. Holmes, Pathogen population bottlenecks and
344     adaptive landscapes: overcoming the barriers to disease emergence. Proc. Biol. Sci.
345     283 (2016), doi:10.1098/rspb.2016.0727.
346 2.  S. M. Kariuki, P. Selhorst, K. K. Ariën, J. R. Dorfman, The HIV-1 transmission
347     bottleneck. Retrovirology. 14 (2017), , doi:10.1186/s12977-017-0343-8.
348 3.  K. Talbert-Slagle, K. E. Atkins, K.-K. Yan, E. Khurana, M. Gerstein, E. H. Bradley,
349     D. Berg, A. P. Galvani, J. P. Townsend, Cellular superspreaders: an epidemiological
350     perspective on HIV infection inside the body. PLoS Pathog. 10, e1004092 (2014).
351 4.  B. F. Keele, E. E. Giorgi, J. F. Salazar-Gonzalez, J. M. Decker, K. T. Pham, M. G.
352     Salazar, C. Sun, T. Grayson, S. Wang, H. Li, X. Wei, C. Jiang, J. L. Kirchherr, F.
353     Gao, J. A. Anderson, L.-H. Ping, R. Swanstrom, G. D. Tomaras, W. A. Blattner, P. A.
354     Goepfert, J. M. Kilby, M. S. Saag, E. L. Delwart, M. P. Busch, M. S. Cohen, D. C.
355     Montefiori, B. F. Haynes, B. Gaschen, G. S. Athreya, H. Y. Lee, N. Wood, C.
356     Seoighe, A. S. Perelson, T. Bhattacharya, B. T. Korber, B. H. Hahn, G. M. Shaw,

357          Identification and characterization of transmitted and early founder virus envelopes in
358          primary HIV-1 infection. Proc. Natl. Acad. Sci. U. S. A. 105, 7552–7557 (2008).

359    5.     J. F. Salazar-Gonzalez, E. Bailes, K. T. Pham, M. G. Salazar, M. B. Guffey, B. F.
360          Keele, C. A. Derdeyn, P. Farmer, E. Hunter, S. Allen, O. Manigart, J. Mulenga, J. A.
361          Anderson, R. Swanstrom, B. F. Haynes, G. S. Athreya, B. T. M. Korber, P. M. Sharp,
362          G. M. Shaw, B. H. Hahn, Deciphering human immunodeficiency virus type 1
363          transmission and early envelope diversification by single-genome amplification and
364          sequencing. J. Virol. 82, 3952–3970 (2008).

365    6.     M.-R. Abrahams, J. A. Anderson, E. E. Giorgi, C. Seoighe, K. Mlisana, L.-H. Ping,
366          G. S. Athreya, F. K. Treurnicht, B. F. Keele, N. Wood, J. F. Salazar-Gonzalez, T.
367          Bhattacharya, H. Chu, I. Hoffman, S. Galvin, C. Mapanje, P. Kazembe, R. Thebus, S.
368          Fiscus, W. Hide, M. S. Cohen, S. A. Karim, B. F. Haynes, G. M. Shaw, B. H. Hahn,
369          B. T. Korber, R. Swanstrom, C. Williamson, CAPRISA Acute Infection Study Team,
370          Center for HIV-AIDS Vaccine Immunology Consortium, Quantitating the
371          multiplicity of infection with human immunodeficiency virus type 1 subtype C
372          reveals a non-poisson distribution of transmitted variants. J. Virol. 83, 3556–3567
373          (2009).

374    7.     H. Li, K. J. Bar, S. Wang, J. M. Decker, Y. Chen, C. Sun, J. F. Salazar-Gonzalez, M.
375          G. Salazar, G. H. Learn, C. J. Morgan, J. E. Schumacher, P. Hraber, E. E. Giorgi, T.
376          Bhattacharya, B. T. Korber, A. S. Perelson, J. J. Eron, M. S. Cohen, C. B. Hicks, B.
377          F. Haynes, M. Markowitz, B. F. Keele, B. H. Hahn, G. M. Shaw, High Multiplicity
378          Infection by HIV-1 in Men Who Have Sex with Men. PLoS Pathog. 6, e1000890
379          (2010).

380    8.     S. Gnanakaran, T. Bhattacharya, M. Daniels, B. F. Keele, P. T. Hraber, A. S.
381          Lapedes, T. Shen, B. Gaschen, M. Krishnamoorthy, H. Li, J. M. Decker, J. F. Salazar-
382          Gonzalez, S. Wang, C. Jiang, F. Gao, R. Swanstrom, J. A. Anderson, L.-H. Ping, M.
383          S. Cohen, M. Markowitz, P. A. Goepfert, M. S. Saag, J. J. Eron, C. B. Hicks, W. A.
384          Blattner, G. D. Tomaras, M. Asmal, N. L. Letvin, P. B. Gilbert, A. C. DeCamp, C. A.
385          Magaret, W. R. Schief, Y.-E. A. Ban, M. Zhang, K. A. Soderberg, J. G. Sodroski, B.
386          F. Haynes, G. M. Shaw, B. H. Hahn, B. Korber, Recurrent Signature Patterns in HIV-
387          1 B Clade Envelope Glycoproteins Associated with either Early or Chronic
388          Infections. PLoS Pathogens. 7 (2011), p. e1002209.

389    9.     D. C. Tully, C. B. Ogilvie, R. E. Batorsky, D. J. Bean, K. A. Power, M.
390          Ghebremichael, H. E. Bedard, A. D. Gladden, A. M. Seese, M. A. Amero, K. Lane,
391          G. McGrath, S. B. Bazner, J. Tinsley, N. J. Lennon, M. R. Henn, Z. L. Brumme, P. J.
392          Norris, E. S. Rosenberg, K. H. Mayer, H. Jessen, S. L. Kosakovsky Pond, B. D.
393          Walker, M. Altfeld, J. M. Carlson, T. M. Allen, Differences in the Selection
394          Bottleneck between Modes of Sexual Transmission Influence the Genetic
395          Composition of the HIV-1 Founder Virus. PLoS Pathog. 12, e1005619 (2016).

396   10.     K. A. Lythgoe, C. Fraser, New insights into the evolutionary rate of HIV-1 at the
397          within-host and epidemiological levels. Proc. Biol. Sci. 279, 3367–3375 (2012).

398   11.     B. F. Keele, J. D. Estes, Barriers to mucosal transmission of immunodeficiency
399          viruses. Blood. 118 (2011), pp. 839–846.

400   12.     L. R. McKinnon, R. Kaul, Quality and quantity. Current Opinion in HIV and AIDS. 7
401          (2012), pp. 195–202.

13. M. Sagar, Origin of the transmitted virus in HIV infection: infected cells versus cell-free virus. J. Infect. Dis. 210 Suppl 3, S667–73 (2014).

14. R. N. Thompson, C. Wymant, R. A. Spriggs, J. Raghwani, C. Fraser, K. A. Lythgoe, Link between the numbers of particles and variants founding new HIV-1 infections depends on the timing of transmission. Virus Evol. 5 (2019), doi:10.1093/ve/vey038.

15. E. O. Romero-Severson, I. Bulla, T. Leitner, Phylogenetically resolving epidemiologic linkage. Proc. Natl. Acad. Sci. U. S. A. 113, 2690–2695 (2016).

16. O. Ratmann, M. K. Grabowski, M. Hall, T. Golubchik, C. Wymant, L. Abeler-Dörner, D. Bonsall, A. Hoppe, A. L. Brown, T. de Oliveira, A. Gall, P. Kellam, D. Pillay, J. Kagaayi, G. Kigozi, T. C. Quinn, M. J. Wawer, O. Laeyendecker, D. Serwadda, R. H. Gray, C. Fraser, PANGEA Consortium and Rakai Health Sciences Program, Inferring HIV-1 transmission networks and sources of epidemic spread in Africa with deep-sequence phylogenetic analysis. Nat. Commun. 10, 1411 (2019).

17. C. Wymant, M. Hall, O. Ratmann, D. Bonsall, T. Golubchik, M. de Cesare, A. Gall, M. Cornelissen, C. Fraser, STOP-HCV Consortium, The Maela Pneumococcal Collaboration, and The BEEHIVE Collaboration, PHYLOSCANNER: Inferring Transmission from Within- and Between-Host Pathogen Genetic Diversity. Mol. Biol. Evol. 35, 719–733 (2018).

18. T. Leitner, E. Romero-Severson, Phylogenetic patterns recover known HIV epidemiological relationships and reveal common transmission of multiple variants. Nat Microbiol. 3, 983–988 (2018).

19. R. Rose, M. Hall, A. D. Redd, S. Lamers, A. E. Barbier, S. F. Porcella, S. E. Hudelson, E. Piwowar-Manning, M. McCauley, T. Gamble, E. A. Wilson, J. Kumwenda, M. C. Hosseinipour, J. G. Hakim, N. Kumarasamy, S. Chariyalertsak, J. H. Pilotto, B. Grinsztejn, L. A. Mills, J. Makhema, B. R. Santos, Y. Q. Chen, T. C. Quinn, C. Fraser, M. S. Cohen, S. H. Eshleman, O. Laeyendecker, Phylogenetic Methods Inconsistently Predict the Direction of HIV Transmission Among Heterosexual Pairs in the HPTN 052 Cohort. J. Infect. Dis. 220, 1406–1413 (2019).

20. A. B. Abecasis, M. Pingarilho, A.-M. Vandamme, Phylogenetic analysis as a forensic tool in HIV transmission investigations. AIDS. 32, 543-554. (2017).

21. D. B. A. Ojwach, D. MacMillan, T. Reddy, V. Novitsky, Z. L. Brumme, M. A. Brockman, T. Ndung'u, J. K. Mann, Pol-Driven Replicative Capacity Impacts Disease Progression in HIV-1 Subtype C Infection. J. Virol. 92 (2018), doi:10.1128/JVI.00811-18.

22. R. W. Shafer, J. M. Schapiro, HIV-1 drug resistance mutations: an updated framework for the second decade of HAART. AIDS Rev. 10, 67–84 (2008).

23. W. C. Miller, N. E. Rosenberg, S. E. Rutstein, K. A. Powers, Role of acute and early HIV infection in the sexual transmission of HIV. Curr. Opin. HIV AIDS. 5, 277–282 (2010).

24. E. M. Volz, E. Ionides, E. O. Romero-Severson, M.-G. Brandt, E. Mokotoff, J. S. Koopman, PLoS Med., 10(12): e1001568 (2013).

25. J. P. Hughes, J. M. Baeten, J. R. Lingappa, A. S. Magaret, A. Wald, G. de Bruyn, J. Kiarie, M. Inambao, W. Kilembe, C. Farquhar, C. Celum, Partners in Prevention HSV/HIV Transmission Study Team, Determinants of per-coital-act HIV-1 infectivity among African HIV-1-serodiscordant couples. J. Infect. Dis. 205, 358–365 (2012).

448    26.    R. H. Gray, M. J. Wawer, R. Brookmeyer, N. K. Sewankambo, D. Serwadda, F.
449        Wabwire-Mangen, T. Lutalo, X. Li, T. vanCott, T. C. Quinn, Rakai Project Team,
450        Probability of HIV-1 transmission per coital act in monogamous, heterosexual, HIV-
451        1-discordant couples in Rakai, Uganda. Lancet. 357, 1149–1153 (2001).

452    27.    T. D. Hollingsworth, R. M. Anderson, C. Fraser, HIV-1 transmission, by stage of
453        infection. J. Infect. Dis. 198, 687–693 (2008).

454    28.    S. E. Bellan, J. Dushoff, A. P. Galvani, L. A. Meyers, Reassessment of HIV-1 acute
455        phase infectivity: accounting for heterogeneity and study design with simulated
456        cohorts. PLoS Med. 12, e1001801 (2015).

457    29.    T. D. Hollingsworth, C. D. Pilcher, F. M. Hecht, S. G. Deeks, C. Fraser, High
458        Transmissibility During Early HIV Infection Among Men Who Have Sex With Men-
459        San Francisco, California. J. Infect. Dis. 211, 1757–1760 (2015).

460    30.    K. A. Lythgoe, A. Gardner, O. G. Pybus, J. Grove, Short-Sighted Virus Evolution and
461        a Germline Hypothesis for Chronic Viral Infections. Trends in Microbiology. 25
462        (2017), pp. 336–348.

463    31.    M.-C. Boily, R. F. Baggaley, L. Wang, B. Masse, R. G. White, R. J. Hayes, M. Alary,
464        Heterosexual risk of HIV-1 infection per sexual act: systematic review and meta-
465        analysis of observational studies. Lancet Infect. Dis. 9, 118–129 (2009).

466    32.    R. F. Baggaley, R. G. White, M.-C. Boily, HIV transmission risk through anal
467        intercourse: systematic review, meta-analysis and implications for HIV prevention.
468        Int. J. Epidemiol. 39, 1048–1063 (2010).

469    33.    M. Beretta, A. Moreau, M. Bouvin-Pley, A. Essat, C. Goujard, M.-L. Chaix, S. Hue,
470        L. Meyer, F. Barin, M. Braibant, ANRS 06 Primo Cohort, Phenotypic properties of
471        envelope glycoproteins of transmitted HIV-1 variants from patients belonging to
472        transmission chains. AIDS. 32, 1917–1926 (2018).

473    34.    J. M. Carlson, M. Schaefer, D. C. Monaco, R. Batorsky, D. T. Claiborne, J. Prince,
474        M. J. Deymier, Z. S. Ende, N. R. Klatt, C. E. DeZiel, T.-H. Lin, J. Peng, A. M. Seese,
475        R. Shapiro, J. Frater, T. Ndung'u, J. Tang, P. Goepfert, J. Gilmour, M. A. Price, W.
476        Kilembe, D. Heckerman, P. J. R. Goulder, T. M. Allen, S. Allen, E. Hunter, HIV
477        transmission. Selection bias at the heterosexual HIV-1 transmission bottleneck.
478        Science. 345, 1254031 (2014).

479    35.    M. S. Cohen, C. L. Gay, M. P. Busch, F. M. Hecht, The Detection of Acute HIV
480        Infection. The Journal of Infectious Diseases. 202 (2010), pp. S270–S277.

481    36.    R. C. Edgar, MUSCLE: a multiple sequence alignment method with reduced time and
482        space complexity. BMC Bioinformatics. 5, 113 (2004).

483    37.    R. C. Edgar, MUSCLE: multiple sequence alignment with high accuracy and high
484        throughput. Nucleic Acids Res. 32, 1792–1797 (2004).

485    38.    F. Ronquist, J. P. Huelsenbeck, MrBayes 3: Bayesian phylogenetic inference under
486        mixed models. Bioinformatics. 19, 1572–1574 (2003).

487    39.    J. P. Huelsenbeck, F. Ronquist, MRBAYES: Bayesian inference of phylogenetic
488        trees. Bioinformatics. 17, 754–755 (2001).

489    40.    O. Ratmann, E. B. Hodcroft, M. Pickles, A. Cori, M. Hall, S. Lycett, C. Colijn, B.
490        Dearlove, X. Didelot, S. Frost, A. S. M. M. Hossain, J. B. Joy, M. Kendall, D.
491        Kühnert, G. E. Leventhal, R. Liang, G. Plazzotta, A. F. Y. Poon, D. A. Rasmussen, T.
492        Stadler, E. Volz, C. Weis, A. J. Leigh Brown, C. Fraser, PANGEA-HIV Consortium,

493          Phylogenetic Tools for Generalized HIV-1 Epidemics: Findings from the PANGEA-
494          HIV Methods Comparison. Mol. Biol. Evol. 34, 185–203 (2017).

495   41.   A. Rambaut, N. C. Grassly, Seq-Gen: an application for the Monte Carlo simulation
496          of DNA sequence evolution along phylogenetic trees. Comput. Appl. Biosci. 13, 235–
497          238 (1997).

498   42.   S. Alizon, C. Fraser, Within-host and between-host evolutionary rates across the
499          HIV-1 genome. Retrovirology. 10, 49 (2013).

500   43.   L.-T. Nguyen, H. A. Schmidt, A. von Haeseler, B. Q. Minh, IQ-TREE: a fast and
501          effective stochastic algorithm for estimating maximum-likelihood phylogenies. Mol.
502          Biol. Evol. 32, 268–274 (2015).

503   44.   S. Kalyaanamoorthy, B. Q. Minh, T. K. F. Wong, A. von Haeseler, L. S. Jermiin,
504          ModelFinder: fast model selection for accurate phylogenetic estimates. Nat. Methods.
505          14, 587–589 (2017).

506   45.   S. R. Cole, H. Chu, S. Greenland, Maximum likelihood, profile likelihood, and
507          penalized likelihood: a primer. Am. J. Epidemiol. 179, 252–260 (2014).

508   46.   S. S. Iyer, F. Bibollet-Ruche, S. Sherrill-Mix, G. H. Learn, L. Plenderleith, A. G.
509          Smith, H. J. Barbian, R. M. Russell, M. V. P. Gondim, C. Y. Bahari, C. M. Shaw, Y.
510          Li, T. Decker, B. F. Haynes, G. M. Shaw, P. M. Sharp, P. Borrow, B. H. Hahn,
511          Resistance to type 1 interferons is a major determinant of HIV-1 transmission fitness.
512          Proc. Natl. Acad. Sci. U. S. A. 114, E590–E599 (2017).

513

1       **Supplementary Materials for:**


2       Number of HIV-1 founder variants is determined by the recency of the source partner infection

3

4       Ch. Julián Villabona-Arenas, Matthew Hall, Katrina A. Lythgoe, Stephen G. Gaffney, Roland R.

5                       Regoes, Stéphane Hué, Katherine E. Atkins


6                       Correspondence to: Katherine.Atkins@ed.ac.uk

7       **This PDF file includes:**
8
9               Materials and Methods
10              Supplementary Text
11              Figs. S1 to S5
12              Additional references
13

14      **Other Supplementary Materials for this manuscript include the following:**
15
16              Data S1 to S4: SITable_EpiGeneticData.csv, SITable_AnalysisData.csv,
17              SITable_ColumnNamesKey.csv, Alignments.zip, Reproducibility checklist
18


19


20      **Materials and Methods**

21



22      **Data collation on linked transmission pairs**

23      We automatically retrieved all HIV sequence data for men-who-have-sex-with-men (MSM) and

24      heterosexual (HET) HIV transmission pairs for whom the direction of transmission is reported

25      from The Los Alamos National Laboratory (LANL) HIV sequence database up to February

26     2019, such that each transmission pair comprise a 'source' and a 'recipient' partner. For each

27     partner in the transmission pair we collected the following clinical and epidemiological data: (i)

28     date of infection or time of infection prior to sampling, (ii) date of seroconversion or date of

29     seroconversion prior to sampling, (iii) Fiebig stage at the time of sampling, (iv) date of sampling

30     or time of sampling prior to infection, (v) number of sequences, (vi) genomic region, (vii) HIV

31     subtype, and (viii) reported risk group. For each set of these transmission pair data we estimated,

32     relative to the transmission time to the recipient partner (time = 0): (*i*) the time of transmission to

33     the source partner, (*ii*) the time of the sampling of the source partner, and (*iii*) the time of

34     sampling for the recipient partner (**Fig. 1**, Supplementary Text). We excluded all transmission

35     pairs from further analysis for whom these three times could not be determined or for whom

36     either partner has fewer than five sequences for all sampling times. For our base case analysis,

37     we used the longest available genomic region with five or more sequences per partner. If more

38     than one sampling time is available for any of the individuals, we selected the sample closest in

39     time to the recipient infection.

40     **Epidemiological data and sequence retrieval**

41     For the ease of replicating our results and using existing transmission pair data for other

42     purposes, we developed a Python script to automatically retrieve epidemiological and metadata

43     for each transmission pair from the Los Alamos National Laboratory HIV sequence database

44     (LANLdb). This script downloads the following data from both the source and recipient partners

45     to a .csv file using as input the cluster and patients ids from LANLdb: (i) time of infection, (ii)

46     time of seroconversion, (iii) Fiebig stage at the time of sampling, (iv) number of sequences, (v)

47     genomic region, (vi) HIV subtype, (vii) reported risk group and (viii) GenBank accession IDs.

48

49    Next we used the downloaded GenBank accession IDs to automatically retrieve (ix) viral genetic

50    sequences and (x) sampling dates (calendar dates) from GenBank using an R script. If

51    information from (i) to (x) were missing for any individual, we manually retrieved these values

52    from the original manuscripts where possible.

53

54    Completed datatables from these automatic and manual processes are provided at

55    github.com/AtkinsGroup.

56

57    **Transmission timelines**

58    For each transmission pair, we define time = 0 as the time of recipient infection. We then

59    calculated, using the data table retrieved, i) the time of infection of the source, ii) the time of

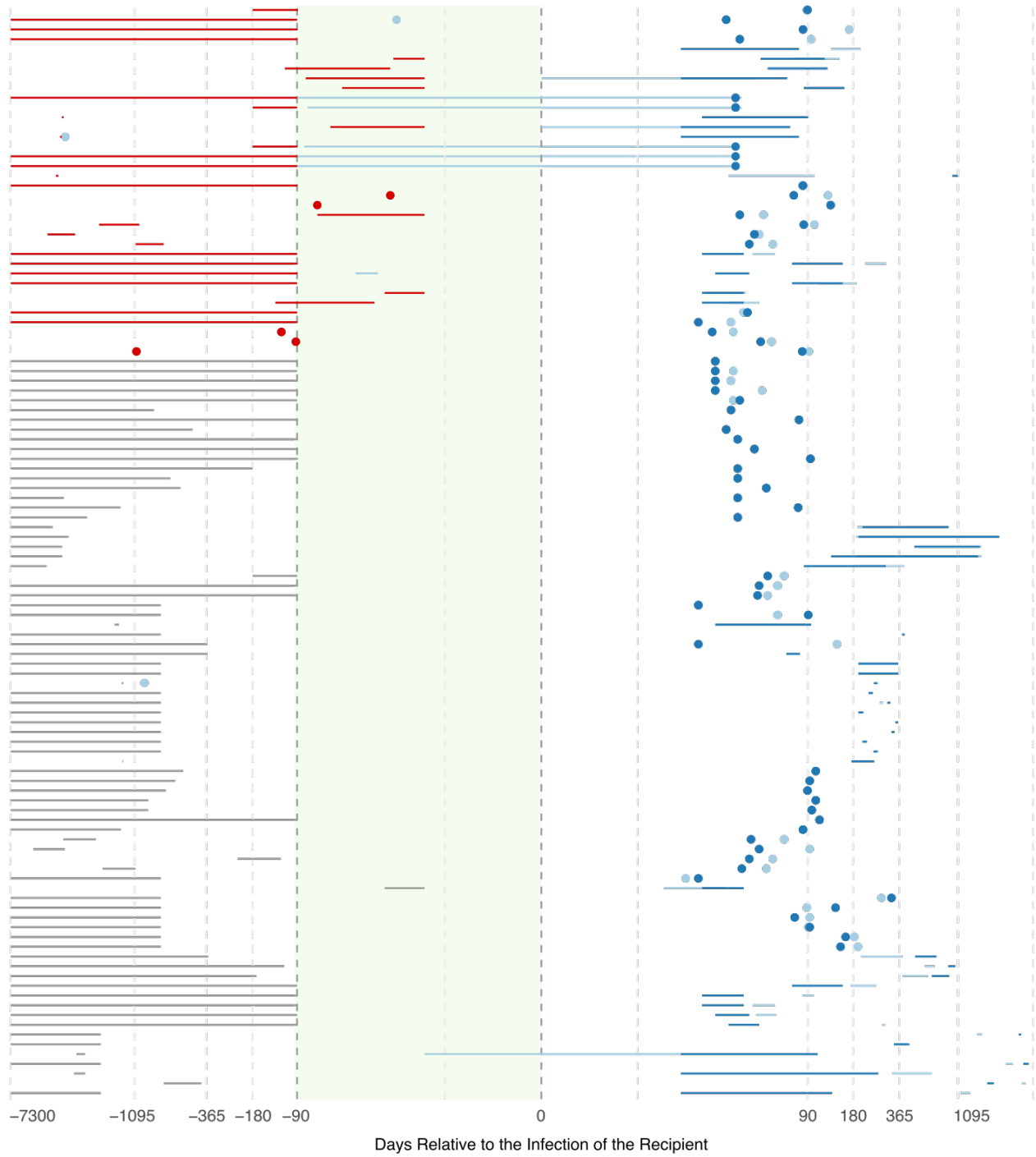60    sampling of the source, iii) the time of sampling of the recipient.

61

62    To estimate these times, we first calculated days from infection for both the source and the

63    recipient partners. When these values are not given explicitly, we calculated them from time

64    since seroconversion estimates or from Fiebig staging results. Specifically, we interpret

65    seroconversion as the individual reaching Fiebig stage III (ELISA positive) that occurs between

66    22-37 days after infection and Fiebig stages I (viral RNA positive) and II (18-34) occurring 13

67    days and 28 days after infection, respectively (*35*).  For all the pairs where a range of possible

68    values is calculated, and for when a calendar month is provided, we incorporated the uncertainty

69    around the infection and sampling times by assuming all values in these ranges are equally

70    plausible and uniformly sampled within these range to account for the uncertainty.

3

71

72    For some pairs, the source was classified as 'recent' or 'late' at the time of transmission to the

73    recipient partner. In these cases, we do not have an exact point to call time = 0. Therefore, for

74    these pairs, for each simulation, we sample with replacement the time between source and

75    recipient infections from the other pairs for whom we have previously classified as acute (<90

76    days delay), and chronic (90 days or more delay), sampling from the same risk group (MSM or

77    HET) in each case.

78    All calculations, corresponding notes and final transmission times for each pair are provided at

79    github.com/AtkinsGroup and visualised in **Fig. S1**.

**Fig. S1: Infection and sampling times of the source and recipient for all the 112 transmission pairs analysed.** Individual points denote exact times and lines denote uniform uncertainty in times. Source

83     partners points/lines overlapping the green shaded area correspond to transmission pairs for whom

84     transmission occurs during the acute stage.

85

86     **Empirical transmission pair analysis**

87     *Tree reconstruction*: For each of the included transmission pairs, we generated posterior sets of

88     phylogenetic trees. For this, we first constructed alignments using Muscle v3.8.31 (*36, 37*) with

89     subtype specific reference sequences retrieved from the LANL HIV sequence database. Using

90     these alignments, we built phylogenetic trees with MrBayes 3.2.7 (*38, 39*) under the assumption

91     of a general time-reversible (GTR) nucleotide substitution model with the addition of invariant

92     sites (I) and a gamma distribution of site rates. We constrained sequence data to be monophyletic

93     with respect to the reference sequences to root the tree but ingroup relationships were

94     unconstrained to avoid any topology class bias. We ran two Markov chains each with 30 million

95     iterations, from which we sampled every 3,000th after discarding the first 50% as burn-in which

96     provided an average standard deviation of split tree frequencies of below 0.01 or an effective

97     sample size of greater than 300. This gave an empirical posterior distribution of $N = 5,000$

98     sample trees. In a sensitivity analysis, we tested the alternative method of using maximum

99     likelihood phylogenetic tree reconstruction with bootstrapping.

100

101    *Empirical topology class:* We classified each of the resulting phylogenetic trees in the posterior

102    distribution as either monophyletic-monophyletic (MM), or paraphyletic-monophyletic (PM), or

103    paraphyletic-polyphyletic (PP), to reflect the cladistic relationship between the lineages from

104    both individuals (Supplementary Text). Each transmission pair,$k$, is then described as a triplet of

105    probabilities, $D_k$, denoting the frequency of each topology class within the $k$th pair's posterior

106 distribution $D_k = \{d_k(t)\}_{t \in T} = \{Pr(t \mid k)\}_{t \in T}$ where $\sum_{t \in T} Pr(t|k) = 1$ and $T \in$

107 {MM, PM, PP}.

108

## Simulated transmission pair analysis

110 We simulated the transmission of virus particles and within-host evolution, accounting for the

111 epidemiological characteristics for each transmission pair. For each transmission pair, we

112 simulated a chain of three HIV infections: (*i*) an unsampled index case who infected the source

113 after three years of their own infection during their chronic stage to reflect that the majority of

114 both HET and MSM transmission pairs transmitted during the chronic stage (101/112 pairs). In a

115 sensitivity analysis we accounted for the assumption that transmission rate may be higher during

116 the acute stage, with half of the index to source transmissions occurring after 90 days and the

117 remaining half after three years, (*ii*) the source individual of the transmission pair, and (*iii*)

118 finally the recipient individual of the transmission pair. For each individual within each trio, we

119 simulated viral phylogenies that reflect between- and within-host viral evolution using

120 VirusTreeSimulator (*40*), using as input the respective epidemiological and clinical information

121 (Supplementary Text). We used a within-host effective population size consistent with that

122 parameterized by the PANGEA-HIV study with the following logistic model parameters: initial

123 effective population size ($N_0$) is 1, viral generation time ($\tau$) is 1.8 days, effective population per

124 year growth rate ($r$) is 2.85022, and time to half the carrying capacity of the viral population

125 ($t_{50}$) is 2 years (*40*). For each transmission pair, we simulated a dated viral phylogeny that has

126 the same number of tips as the number of retrieved sequences per partner and that is sampled at

127 the respective sampling times for the source and recipient partner (Supplementary Text). For

128 each recipient partner infection, we assume that a total of $n_R$ virus particles founded the

129    infection. For each simulation, we further assume a total of $n_S$ virus particles founding infection

130    of the source. We assume $n_R$ takes values between one and a maximum of 12 and varied $n_S$

131    between one and two (Supplementary Text). We assume that the virus samples from each

132    recipient is representative of the within-host diversity, and that each founding virus particle has

133    an extant lineage. Therefore, we first assigned each sample (tip) of a phylogeny as a descendant

134    of one of the $n_R$ virus particles. If there were more than 12 samples then the remaining tips were

135    assigned randomly to the $n_R = 12$ virus particles. If there were fewer than 12 samples, then we

136    constrained the number of founding virus particles, $n_R$, to equal the number of samples. For

137    every transmission pair, and for each value of $n_R$ and $n_S$, we simulated 100 viral phylogenies.

138

139    For every simulated viral phylogeny, we simulated transmitted sequences by adding dummy

140    nodes with a negligibly short branch length after the transmission time. We then simulated the

141    evolution of nucleotide sequences along the tree using Seq-Gen (*41*) and a GTR + I + gamma

142    substitution model. The length of the simulated sequences and the evolutionary tree scaling rate

143    match each transmission pair's empirical sequence data. For this, we used previously estimated

144    empirically-derived within-host evolutionary rates (*42*) and the HXB2 sequence homologous to

145    the pair's sequence fragment as the ancestral sequence at the root. Every transmission pair

146    simulation produces a tip sequence alignment and a number of founder sequences equal to the

147    number of transmitted particles.

148

149    *Simulated topology class*: We reconstructed a phylogeny using maximum likelihood inference in

150    IQ-TREE 1.6.11 (*43*) and selected the best-fit nucleotide substitution model with ModelFinder

151    (*44*). Each phylogeny was classified as either MM, PM or PP (Supplementary Text).

152   Consequently, for each transmission pair $k$ and each transmissibility model (i.e. number of viral

153   particles founding infection of the recipient $n_R$), we generated a triplet of probabilities $M_{k,n_R} =$

154   $\{m_{k,n_R}\}_{t \in T} = \Pr(t|k, n_R)\}_{t \in T}$ where $\sum_{t \in T} \Pr(t|k, n_R) = 1$ and $T \in \{MM, PM, PP\}$.

155

**Transmissibility model calibration**

157   For each transmission pair, we chose the most likely value of $n_R$ (the number of virus particles

158   founding each recipient infection) by matching the posterior topology class from the empirical

159   phylogenetic transmission trees with the simulated distribution of topology class. Specifically,

160   for each transmission pair, $k$, we estimated the most likely number of viral particles founding

161   each recipient infection $n_R^*$ as the $n_R$ that maximises the multinomial likelihood function $L_{k,n_R} =$

162   $\Pr(D_k | M_{k,n_R}) = \frac{N!}{\prod_{t \in T}(Nd_k(t))!} \prod_{t \in T} m_{k,n_R}(t)^{Nd_k(t)}$. For each transmission pair $k$, we calculated

163   lower and upper confidence limits for $n_R^*$ as the minimum and maximum values of $n_R$ that satisfy

164   $L_{k,n_R} > L_{k,n_R^*} - 1.92$ and $L_{k,n_R} < L_{k,n_R^*} + 1.92$, respectively (*44, 45*). For each transmission

165   pair $k$, we retain the best fit model for further analysis such that there are $n_R^*$ viral particles

166   founding infection of the recipient.

**Haplotype analysis**

168   *Probability of a single founder haplotype*: For each transmission pair, $k$, from the best fit

169   transmissibility model, we defined the random variables $F_S^k$ and $F_R^k$ as the number of haplotypes

170   that found infection of the source and the recipient partners, respectively. We then calculated the

171   probability of there being a single founder haplotype in the recipient, stratified by topology class

172   of the simulated phylogenetic tree (MM, PM, PP) and the number of founder haplotypes, *i*, in the

173     source partner, $p_i^k(t)$, that is, $p_i^k(t) = \Pr(F_R^k = 1 | F_S^k = i, t)$. Next, we defined the probability

174     of a single founder haplotype in the recipient as a function of a tree topology, $t$, $p^k(t) =$

175     $\Pr(F_R^k = 1 | t) = p_1^k \Pr(F_S = 1) + p_2^k \Pr(F_S > 1)$. By assuming that the source partners are

176     randomly selected from the general MSM or HET population in which the probability of a single

177     founder variant has been calculated to be approximately 0.7 (*14*), we set, $Pr(F_S = 1) =$

178     0.7 and $Pr(F_S > 1) = 0.3$. Finally, for each transmission pair, we calculated the probability of

179     one founder haplotype given the observed triplet of empirical posterior topology classes $D_k$, as

180     $q^k = \sum_{t \in T} p^k(t) d_k(t) / N$.

181

182     *Number of founder haplotypes by source partner infection stage*: We stratified all the

183     transmission pairs into two sets by the infection stage of the source partner. We classified the

184     acute transmission set as those pairs for whom recipient infection is within 90 days of source

185     infection (a set of $n_{\text{acute}}$ pairs), and the chronic transmission set as those pairs for whom recipient

186     infection is 90 days or later after source infection (a set of $n_{\text{chronic}}$ pairs). For each group, we

187     calculated the mean probability of one founder haplotype being transmitted to the recipient in

188     each set set as:

189     $$q_{\text{acute}}^k = \sqrt[n_{\text{acute}}]{\prod_{k \in \text{acute}} q^k}$$

190     $$q_{\text{chronic}}^k = \sqrt[n_{\text{chronic}}]{\prod_{k \in \text{chronic}} q^k}$$

191     Finally, we calculated the relative risk of one founder haplotype transmitted during the acute

192     stage versus the chronic stage by $q_{acute}^k / q_{chronic}^k$.

193    **Statistical analysis**

194    We compare our results by using statistical tests and report the respective *P*-values. To compare

195

196

197

198

199    # Supplementary Text

200

201    **Transmission pairs sequence data**

202    Our alignments are provided at github.com/AtkinsGroup**.**

203    On average, 22 (IQR 13-33) HIV sequences are obtained from the source and 21 (IQR 10-20)

204    sequences from the recipient for the MSM pairs, and 21 (IQR 12-25) and 18 (IQR 9-22) for the

205    HET source and recipient, respectively. All MSM sequence data belong to subtype B, while most

206    heterosexual sequence data belong to subtype C (49%), followed by subtype B (22%), subtype D

207    (21%), subtype A/A-like (7%) and unclassified subtype (1%). A total of 7 (19%) of the MSM

208    pairs have near full genomes sequenced and the remaining pairs had *env* available (mean 1653

209    nt, range 182-3827 nt). Ten (13%) of the HET pairs had near full genomes available, while 56

210    (75%) pairs had *env* (mean 1321 nt, range 323-2582 nt), nine (12%) pairs had either *pol* or *gag*

211    (mean 1484 nt, range 1375-1499 nt) and one pair had vif-LTR3 (4666 nt) sequenced.

212

**Effect of number of founding virus particles in the source**

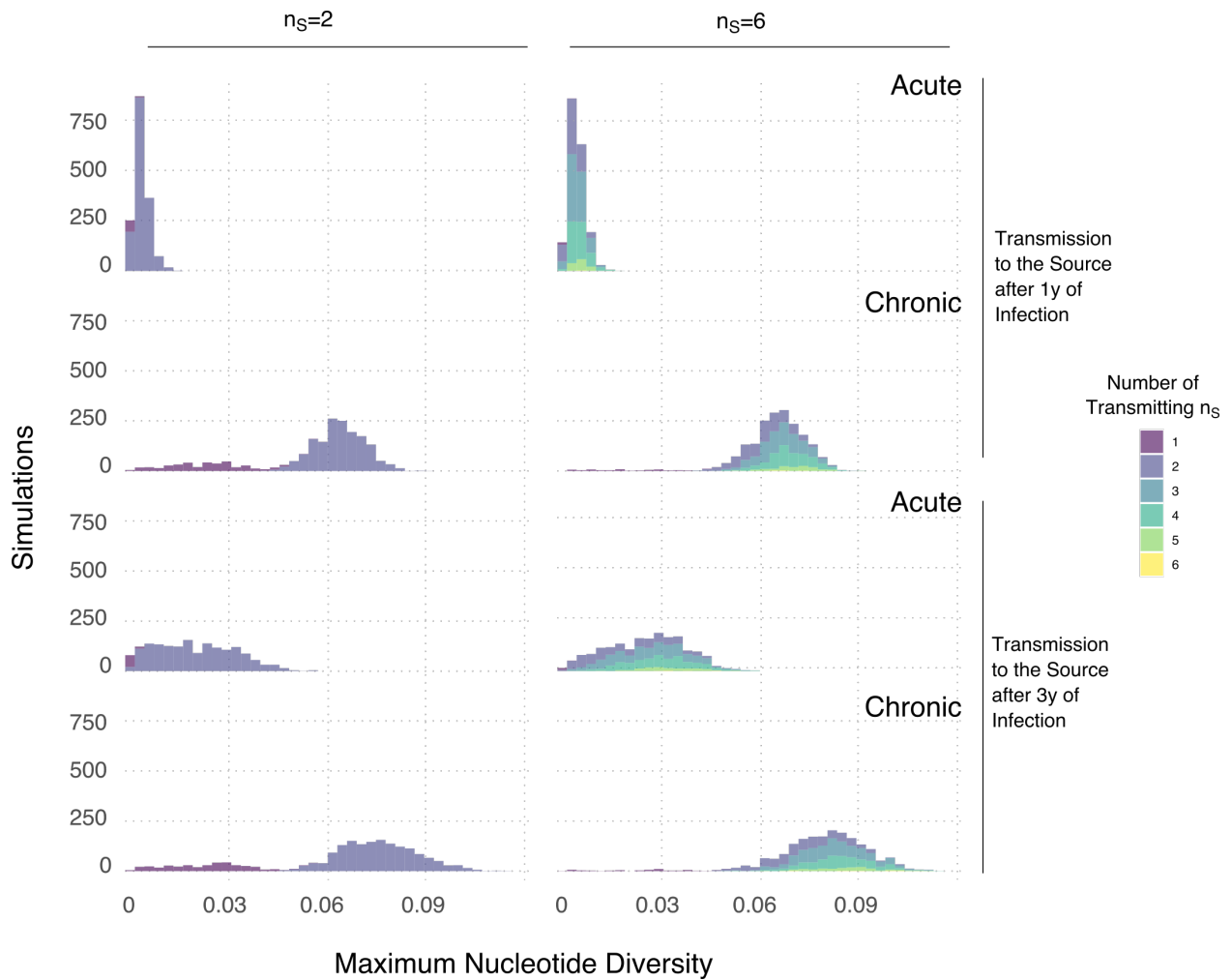To assess whether the number of founding virus particles in the source partner affects the

diversity of sequences founding infection in the recipient, we model a scenario where the index

case transmitted one, two or six virus particles to the source partner at either one or three year(s)

after infection. The source in turn transmits 1 to 6 virus particles to the recipient at 30 days

(acute) or 1095 days (chronic) later. The simulation produces a dated viral phylogeny with tips

sampled at either 30 (early) or 1065 days (late). We model 1kb nucleotide sequences along the

simulated viral phylogenies using the same method as in the main text.

The genetic variation rapidly and steadily increases over time – the maximum diversity among

transmitted haplotypes to the recipient was higher when the index case was infected for longer

and the transmission to the recipient occurs during the chronic stage of the source (**Fig. S2**).

When the source has more than one founding particle, this leads to a bimodal distribution of

maximum diversity among transmitted founder variants within the recipient. The first and second

mode represent maximum diversity when drawing the recipient founder haplotypes from either

one or more than one viral population within the source, respectively. However, increasing the

number of founding virus particles to more than two within the source only increases the density

around the second mode without affecting the range of the maximum diversity distribution. This

consistency occurs because increasing the number of founding virus particles in both

transmission partners, increased the probability of drawing founding variants from different

genetic pools in the source. However, the average maximum diversity of the founder variants

234    does not change because the source genetic pools evolved at the same rate and under the same

235    evolutionary constraints with no selective advantage. This leads to genetic pools with equivalent

236    cumulative genetic change but distinct identity. Taking this into account, we chose to model one

237    or two founding virus particles within the source partner as we were interested in capturing some

238    degree of variation in the transmitted haplotypes rather than multiple genetic identities *per se*.



239    Maximum Nucleotide Diversity

240    **Fig. S2: Effect of number of founding virus particles in the source.**

241

242

243

244

245

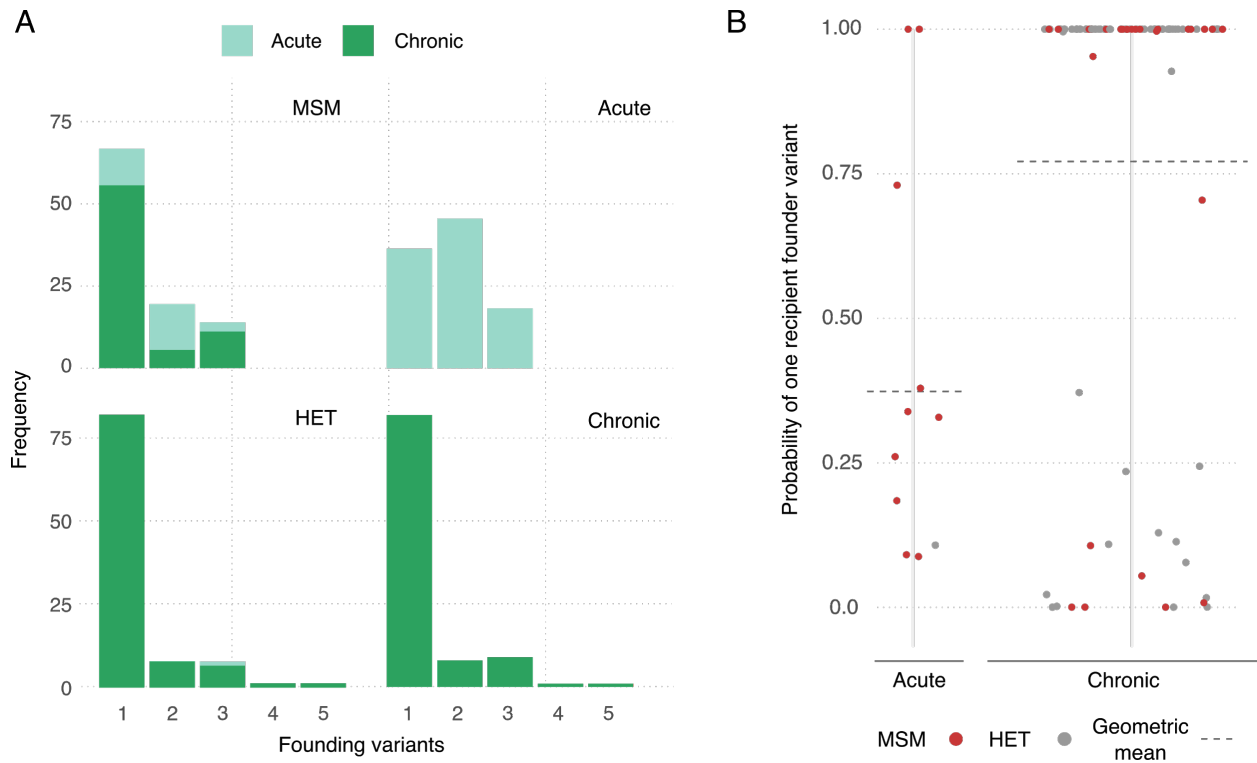**Effect of using a confidence threshold for assigning the topology class**



**Fig. S3: Phylogenetic findings of the empirical transmission pairs for whom the posterior trees gave a certainty of over 95% for the most frequent topology**. Fraction of phylogenetic tree topology class (MM - Monophyletic-Monophyletic, PM - Paraphyletic-Monophyletic and PP - Paraphyletic-Polyphyletic) where each tree topology class is classified as the most frequent topology class of each posterior distribution per transmission pair. Results are stratified by risk group: 76 heterosexual (HET) pairs and 36 men-who-have-sex-with-men (MSM) pairs) and infection stage of the source partner at transmission (11 acute pairs defined as <90d post infection and 101 chronic pairs defined as ≥90d post infection).

**Effect of index partner stage of infection at transmission**

258     In the main analysis we assumed that all index cases transmit to the source partner after three

259     years of infection. Here we also evaluated the results assuming the transmission risk was skewed

260     towards early infection, with half of all simulations across all transmission pairs assuming index

261     case transmission occurs during the acute stage ($\leq$90d) and half occurs during the chronic stage

262     (91d-3y). We find qualitatively similar results as our main analysis. The median number of

263     founder variants transmitted across all pairs is 1 (range: 1-5, **Fig. S4A**). Across all pairs in both

264     risk groups, the mean probability of observing one founder variant is 0.73. Stratifying by risk

265     group, we find there is a higher probability that one variant founds HET infections than MSM

266     infections (a geometric mean of 0.79 vs. 0.61, **Fig. S4B**). In contrast, when stratifying solely by

267     infection stage of the source partner, we find that transmission during the acute stage has a much

268     lower probability of one founder variant than during the chronic stage (means of 0.38 vs. 0.78)

269     with a higher median number of founder variants transmitted, when only the most likely number

270     of transmitted founder variants for each pair is considered (2 vs. 1, **Fig. S4A**). From these results,

271     therefore, there is still approximately twice the chance of multiple founder variant transmission

272     during acute stage infection across both risk groups (relative risk is 0.48).

273

274

**Fig. S4: Phylogenetic findings from the calibrated simulations with skewed transmission rate towards acute stage for the index case**. A) Frequency of the number of founder variants for transmission pairs by infection stage of source partner at transmission and risk group. The number of transmitted founding variants is calculated as the modal simulated value. B) Probability of one founder variant in the recipient for each pair stratified by infection stage of the source partner at transmission.

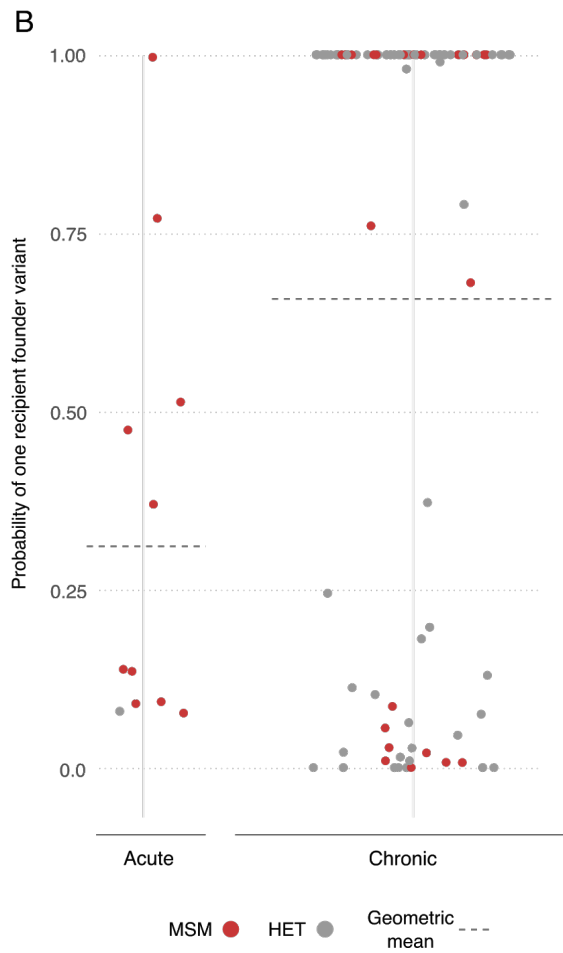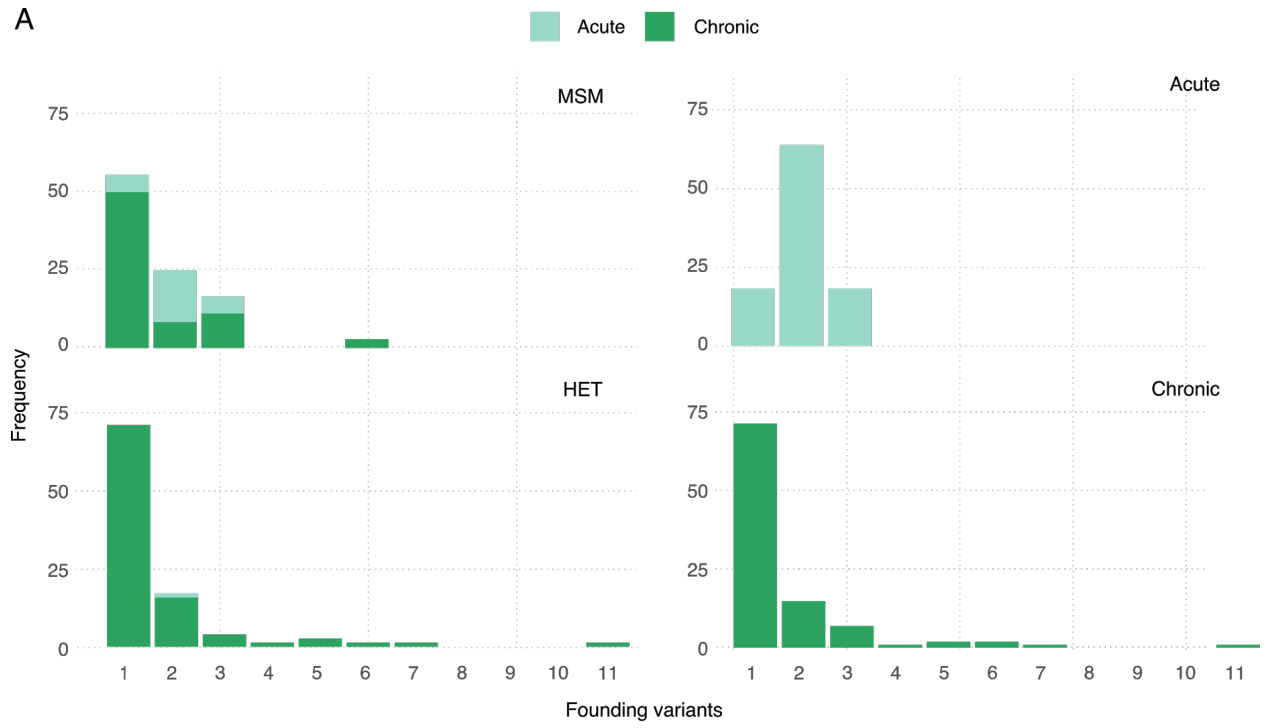**Effect of constructing empirical data phylogenetic trees using maximum likelihood inference with bootstrapping**

In the main analysis we used Bayesian phylogenetic reconstruction to analyse the empirical sampled genetic data of each empirical transmission pair, using the respective posterior distribution to calculate the frequency of each topology class (MM, PM and PP). Here we

286    provide a sensitivity analysis to calculate the tree topology class distribution of the empirical

287    sampled genetic data by maximum likelihood phylogenetic tree construction and bootstrapping.

288    After bootstrapping the empirical data 100 times to calculate the frequency of MM, PM and PP

289    topology classes for each transmission pair, we then proceeded using the same methodology as

290    in the main text. That is, we fit the simulation model (parameterised with the pair-specific data)

291    to the bootstrapped data individually for each transmission pair by comparing the frequencies of

292    tree topology classes. Overall our results remained consistent with our main results, albeit with

293    slightly lower probabilities of observing one founder variant. The median number of founder

294    variants transmitted across all pairs is 1 (range: 1-11, **Fig. S5A**). Across all pairs in both risk

295    groups, the mean probability of observing one founder variant is 0.62. Stratifying by risk group,

296    we find there is a higher probability that one variant founds HET infections than MSM infections

297    (a geometric mean of 0.67 vs. 0.53, **Fig. S5B**). Stratifying by infection stage of the source

298    partner, we find there is a much lower probability of one founder variant during the acute stage

299    than during the chronic stage (means of 0.31 vs. 0.66) with approximately twice the chance of

300    multiple founder variant transmission during acute stage infection across both risk groups

301    (relative risk is 0.47).

302

18

**Fig. S5: Phylogenetic findings from the calibrated simulations with bootstrapped empirical data**. A) Frequency of the number of founding variants for transmission pairs by infection stage of source partner at transmission and risk group. The number of transmitted founding variants is calculated as the modal simulated value. B) Probability of one founder variant in the recipient for each pair stratified by infection stage of the source partner at transmission.

## Effect of the number of sequences for each transmission pair

Here we provide sensitivity analysis to the estimation of the probability that a single founder variant was transmitted to the respective recipient by the number of sequences available from the source and recipient, which ranges from 5 to 149 across all partners. First, the number of sequences available from the transmitter and the recipient is correlated (*Pearson's* product-moment correlation=0.53, $P<0.01$). However, we do not find any evidence of correlation between the total number of sequences for a pair and the estimated number of founder variants in the recipient ($P>0.2$). While an MM topology is more frequently observed when the total number of sequences was small ($P<0.01$), removing the pairs with a likely MM topology do not change our main result: the probability that a single founder variant was transmitted to the respective recipient is lower for the acute pairs (0.402) than for the chronic ones (0.749).

## Effect of the sequencing method

We evaluated if our results were affected by the type of sequence data used in the analysis. All of the transmission pair data were generated using Sanger capillary sequencing except for those in one study ((*46*) in **Data S1)** which used Illumina sequencing on end-point diluted primary isolates. Our results are robust to the exclusion of the eight transmission pairs extracted from this

324  study: that is, the probability that a single founder variant is transmitted to the respective

325  recipient is lower for the acute stage (0.402) than for the chronic stage (0.756).

326  **Effect of the gene region and length**

327  Looking at chronic stage transmissions only, if we compare the number of founder variants

328  inferred from envelope gene sequences to  those inferred from non-envelope sequences, we don't

329  find significant differences ($P>0.4$) in the probability that a single founder variant is transmitted

330  to the respective recipient: 0.739 for envelope sequences  and 0.856 for non-envelope ones.

331  Conversely, if we include data from both chronic and acute transmissions, and restrict our

332  analysis to those pairs with sequences from the envelope region, our results remain unchanged.

333  That is, the probability that a single founder variant is transmitted to the respective recipient is

334  lower during the acute stage (0.432) than during the chronic stage (0.739). Finally, if we

335  condition our analysis on those pairs for whom full or near full genomes are available (17 pairs),

336  our results remain consistent with the main analysis: the probability that a single founder variant

337  was transmitted is lower for acute stage transmissions (0.138, n=1) than for chronic stage

338  transmissions (0.903, n=16). We found that length of the sequenced region is not correlated with

339  the probability that a single founder variant is transmitted to the recipient (*Pearson's* product-

340  moment correlation=0.14, $P>0.14$). Moreover, there are no significant differences ($P>0.81$) in

341  the length of the sequenced region when stratifying our data by infection stage of the source

342  partner at transmission.  Together these observations indicate that our results are not influenced

343  by the length of sequenced regions.

344

345

346　　**Data S1 (Separate file):** SITable_EpiGeneticData.csv. Collated epidemiological and clinical

347　　data and genetic metadata for the 112 transmission pairs used in the analysis.

348　　**Data S2 (Separate file)**: SITable_AnalysisData.csv. Analysis information for the 112

349　　transmission pairs used in the analysis.

350　　**Data S3 (Separate file)**: SITable_ColumnNamesKey.csv. Additional information on column

351　　headers in Data S1,S2 tables

352　　**Data S4 (Separate file)**: Alignments.zip. Individual files of sequence alignments used in the

353　　analysis for the 112 transmission pairs.

354