

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.

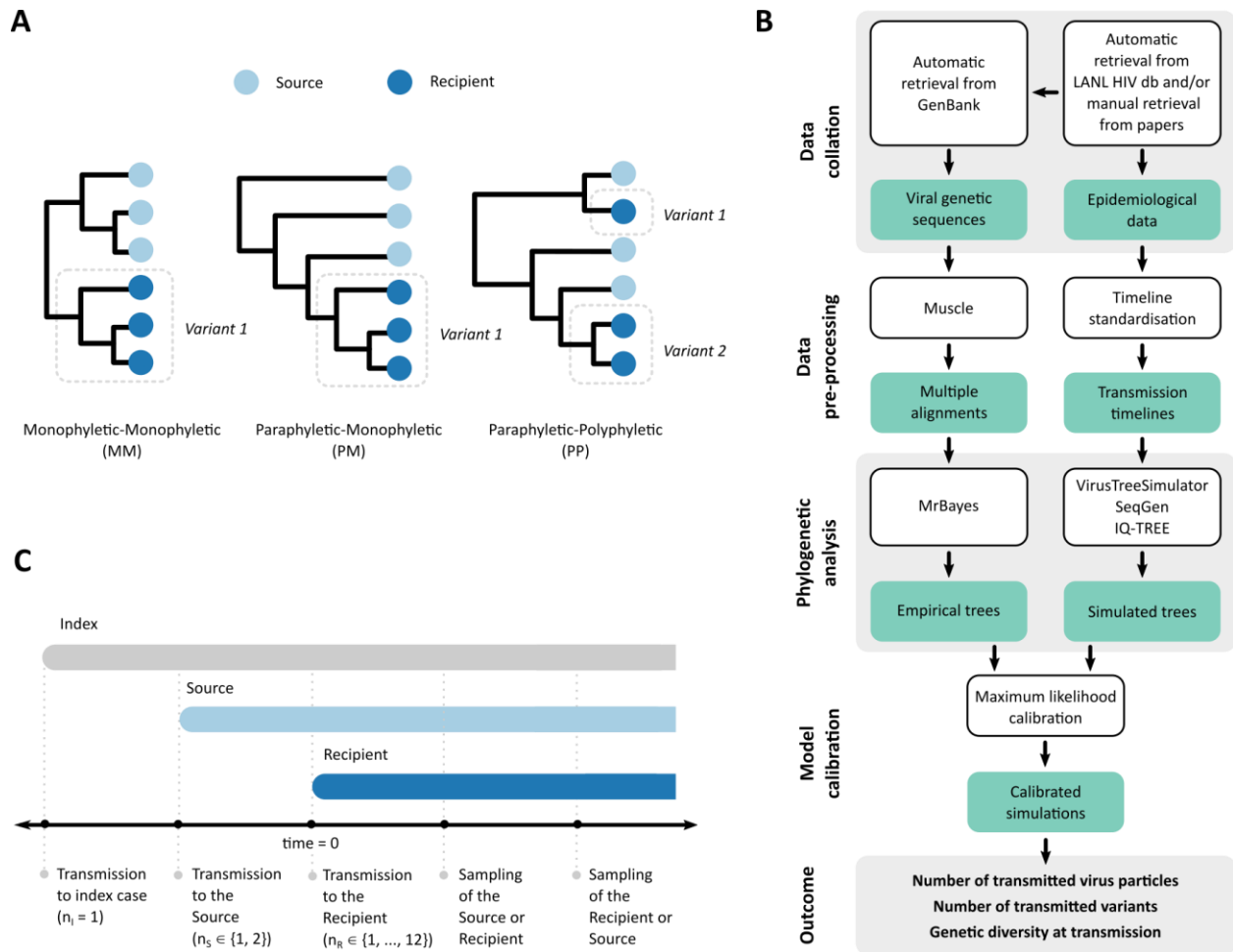
57

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**).



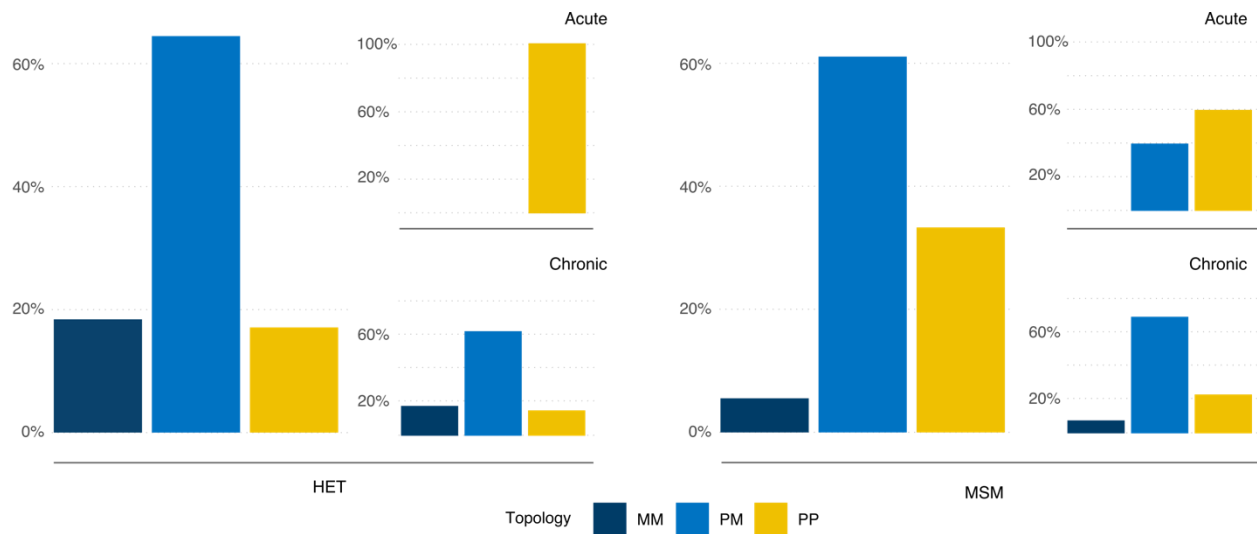
88
 89 **Fig. 1: Methods schematics.** A) Phylogenetic tree topology class of known transmission pairs that have
 90 previously been used as a proxy for calculating the minimum number of founder variants transmitted to

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

114 transmission pairs for whom the posterior trees gave a certainty of over 95% for the most
 115 frequent topology class (**Fig. S3**). These results indicate that infection stage of the source partner,
 116 and not risk group *per se*, influences the diversity bottleneck at transmission.
 117



118
 119 **Fig. 2: Phylogenetic findings from the empirical transmission pairs.** Fraction of phylogenetic tree
 120 topology class (MM: Monophyletic-Monophyletic, PM: Paraphyletic-Monophyletic and PP: Paraphyletic-
 121 Polyphyletic) where each tree topology class is classified as the most frequent topology class of each
 122 posterior distribution per transmission pair. Results are stratified by risk group: 76 heterosexual (HET)
 123 pairs and 36 men-who-have-sex-with-men (MSM) pairs) and infection stage of the source partner at
 124 transmission (11 acute pairs defined as <90d post infection and 101 chronic pairs defined as ≥90d post
 125 infection).

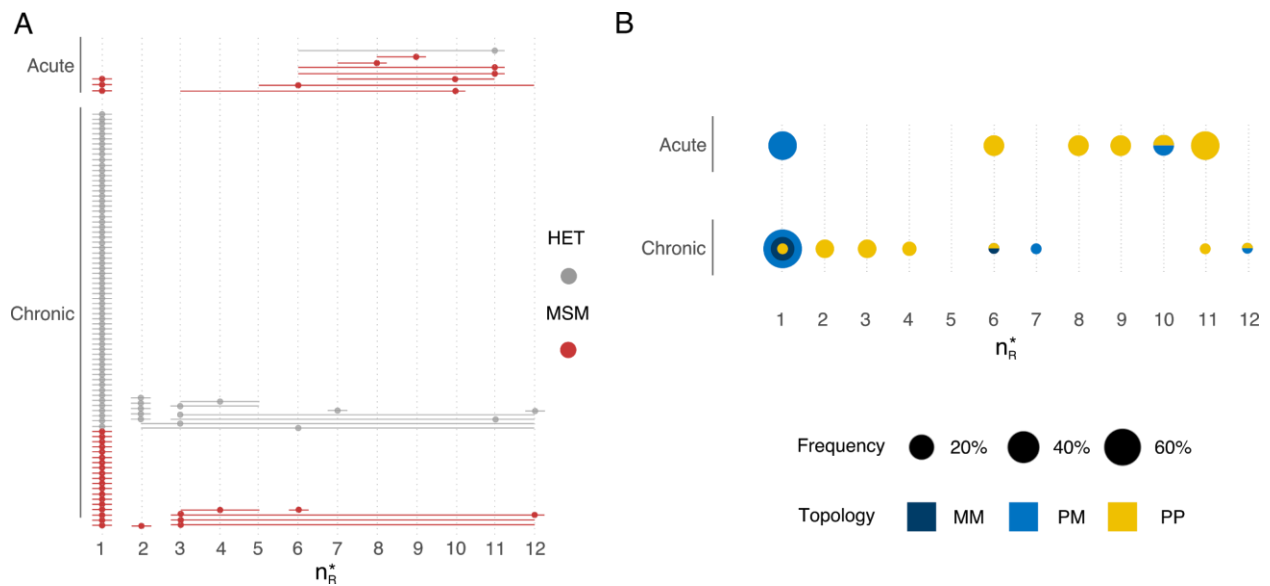
126 To test whether these empirical findings are indicative of a smaller diversity bottleneck in the
 127 chronic stage of HIV infection, we developed a phylodynamic framework in which we simulated
 128 the epidemiologic characteristics of each HET and MSM transmission pair, the timing of their
 129 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

153 stage transmissions (73% vs. 20%, $P = 0.0005$). The topology class of the simulated
 154 phylogenetic trees is strongly influenced by the number of virus particles being transmitted (**Fig.**
 155 **3B**). PM trees are more commonly found in the pairs that are better described by a model with a
 156 single transmitted virus particle (81%) whereas PP trees appeared more often when multiple
 157 particles are likely to have been transmitted (86%).

158



159

160 **Fig. 3: The estimated number of transmitted virus particles for the 112 transmission pairs.** The
 161 estimates of transmitted virus particles for each transmission pair were calculated by choosing the model
 162 simulation that generated a phylogenetic tree topology class distribution (that is, the number of MM, PM
 163 and PP trees constructed from the simulated genetic sequences) that best matched the topology class
 164 distribution from the phylogenetic trees constructed from the empirical genetic sequences. A) Maximum
 165 likelihood number of virus particles founding recipient infections, n_R^* , for each pair (stacked points) with
 166 95% confidence intervals (lines) grouped by stage of infection (acute, 11 pairs or chronic, 101 pairs) and
 167 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

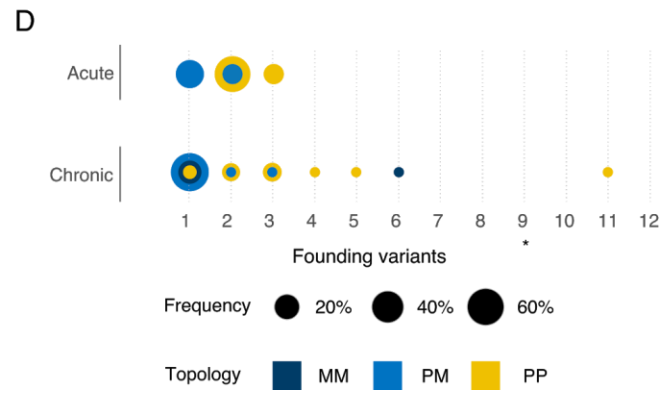
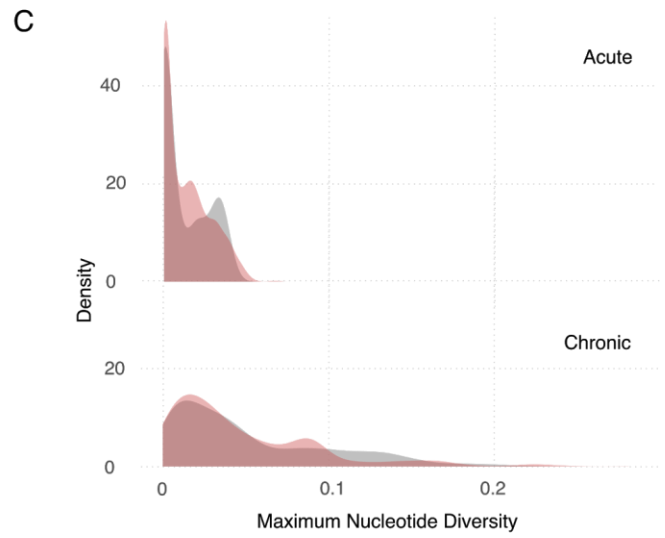
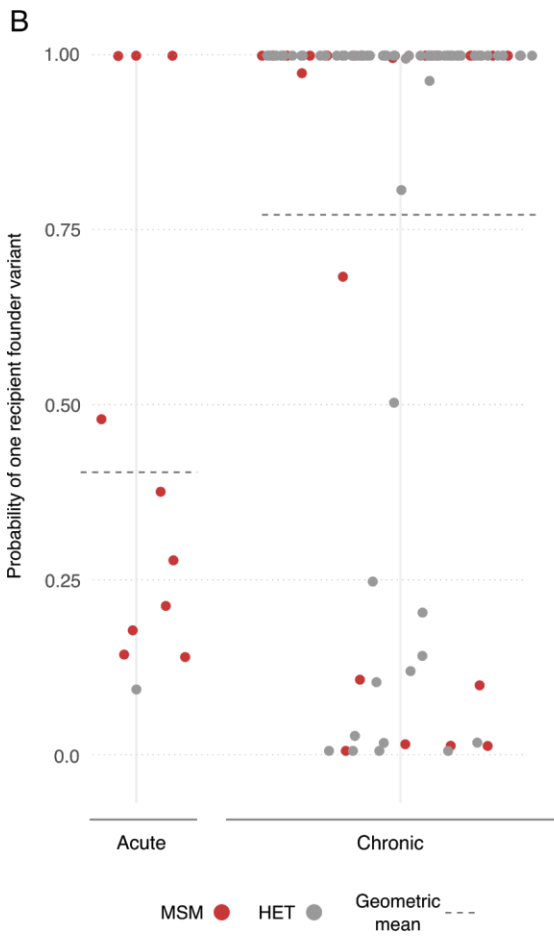
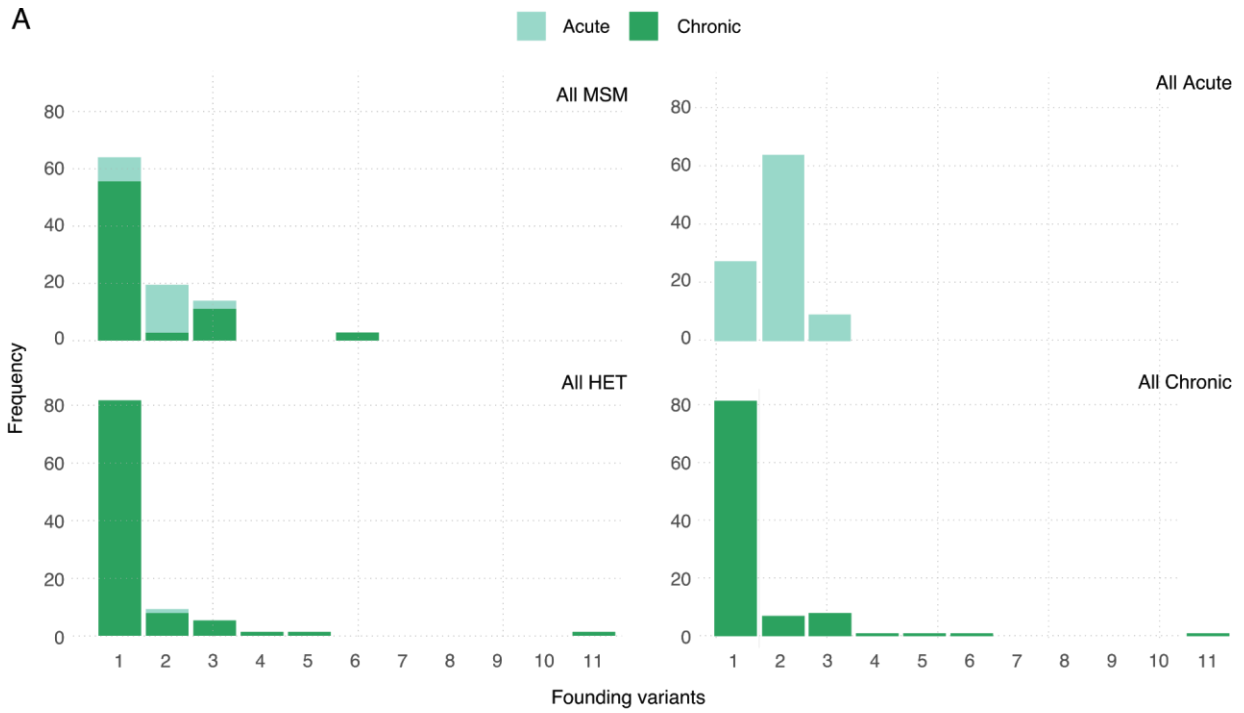
200 Our results suggest that there is an association between tree topology class and multiple founder
201 variant transmission, with 95% of MM and PM trees being due to one founder variant (**Fig. 4D**).
202 However, the number of embedded recipient clades is not always a proxy for the minimum
203 number 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



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

244

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

251

252 Data from four of the MSM transmission pairs in this study have previously been used to
253 estimate the number of variants founding infection using a combination approach of single
254 genome amplification (SGA), direct amplicon sequencing and mathematical modelling (7). Our
255 results broadly agree with this previous analysis, with both analyses suggesting two recipients
256 were infected with one founder variant and one recipient was infected with multiple founder
257 variants (our analysis suggests a mean of 2-3 founder variants and the previous analysis suggests
258 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**

317 **Funding:** CJVA and KEA were funded by an ERC Starting Grant (award number 757688)
318 awarded to KEA. KAL was supported by The Wellcome Trust and The Royal Society grant no.
319 107652/Z/15/Z. MH was funded by The HIV Prevention Trials Network (grant number
320 H5R00701.CR00.01) and The Bill and Melinda Gates Foundation (grant number OPP1175094).

321 **Author contributions:** KEA conceived the study. CJVA, MH, KL, SH, KEA designed the
322 study. CJVA performed the experiments and analysed the data. CJVA, MH, KL, SH, KEA
323 interpreted the data. SGG created new software used in the study. KEA and CJVA drafted the
324 manuscript, with critical revisions from MH, RRR, KL, SH. All authors approved the final
325 version of the manuscript. **Competing interests:** The authors declare no competing interests.

326 **Data and materials availability:** All code and data are available at github.com/AtkinsGroup in

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.

- 402 13. M. Sagar, Origin of the transmitted virus in HIV infection: infected cells versus cell-
403 free virus. *J. Infect. Dis.* 210 Suppl 3, S667–73 (2014).
- 404 14. R. N. Thompson, C. Wymant, R. A. Spriggs, J. Raghwani, C. Fraser, K. A. Lythgoe,
405 Link between the numbers of particles and variants founding new HIV-1 infections
406 depends on the timing of transmission. *Virus Evol.* 5 (2019), doi:10.1093/ve/vey038.
- 407 15. E. O. Romero-Severson, I. Bulla, T. Leitner, Phylogenetically resolving
408 epidemiologic linkage. *Proc. Natl. Acad. Sci. U. S. A.* 113, 2690–2695 (2016).
- 409 16. O. Ratmann, M. K. Grabowski, M. Hall, T. Golubchik, C. Wymant, L. Abeler-
410 Dörner, D. Bonsall, A. Hoppe, A. L. Brown, T. de Oliveira, A. Gall, P. Kellam, D.
411 Pillay, J. Kagaayi, G. Kigozi, T. C. Quinn, M. J. Wawer, O. Laeyendecker, D.
412 Serwadda, R. H. Gray, C. Fraser, PANGAEA Consortium and Rakai Health Sciences
413 Program, Inferring HIV-1 transmission networks and sources of epidemic spread in
414 Africa with deep-sequence phylogenetic analysis. *Nat. Commun.* 10, 1411 (2019).
- 415 17. C. Wymant, M. Hall, O. Ratmann, D. Bonsall, T. Golubchik, M. de Cesare, A. Gall,
416 M. Cornelissen, C. Fraser, STOP-HCV Consortium, The Maela Pneumococcal
417 Collaboration, and The BEEHIVE Collaboration, PHYLOSCANNER: Inferring
418 Transmission from Within- and Between-Host Pathogen Genetic Diversity. *Mol.*
419 *Biol. Evol.* 35, 719–733 (2018).
- 420 18. T. Leitner, E. Romero-Severson, Phylogenetic patterns recover known HIV
421 epidemiological relationships and reveal common transmission of multiple variants.
422 *Nat Microbiol.* 3, 983–988 (2018).
- 423 19. R. Rose, M. Hall, A. D. Redd, S. Lamers, A. E. Barbier, S. F. Porcella, S. E.
424 Hudelson, E. Piwowar-Manning, M. McCauley, T. Gamble, E. A. Wilson, J.
425 Kumwenda, M. C. Hosseinipour, J. G. Hakim, N. Kumarasamy, S. Chariyalertsak, J.
426 H. Pilotto, B. Grinsztejn, L. A. Mills, J. Makhema, B. R. Santos, Y. Q. Chen, T. C.
427 Quinn, C. Fraser, M. S. Cohen, S. H. Eshleman, O. Laeyendecker, Phylogenetic
428 Methods Inconsistently Predict the Direction of HIV Transmission Among
429 Heterosexual Pairs in the HPTN 052 Cohort. *J. Infect. Dis.* 220, 1406–1413 (2019).
- 430 20. A. B. Abecasis, M. Pingarilho, A.-M. Vandamme, Phylogenetic analysis as a forensic
431 tool in HIV transmission investigations. *AIDS.* 32, 543-554. (2017).
- 432 21. D. B. A. Ojwach, D. MacMillan, T. Reddy, V. Novitsky, Z. L. Brumme, M. A.
433 Brockman, T. Ndung'u, J. K. Mann, Pol-Driven Replicative Capacity Impacts
434 Disease Progression in HIV-1 Subtype C Infection. *J. Virol.* 92 (2018),
435 doi:10.1128/JVI.00811-18.
- 436 22. R. W. Shafer, J. M. Schapiro, HIV-1 drug resistance mutations: an updated
437 framework for the second decade of HAART. *AIDS Rev.* 10, 67–84 (2008).
- 438 23. W. C. Miller, N. E. Rosenberg, S. E. Rutstein, K. A. Powers, Role of acute and early
439 HIV infection in the sexual transmission of HIV. *Curr. Opin. HIV AIDS.* 5, 277–282
440 (2010).
- 441 24. E. M. Volz, E. Ionides, E. O. Romero-Severson, M.-G. Brandt, E. Mokotoff, J. S.
442 Koopman, *PLoS Med.*, 10(12): e1001568 (2013).
- 443 25. J. P. Hughes, J. M. Baeten, J. R. Lingappa, A. S. Magaret, A. Wald, G. de Bruyn, J.
444 Kiarie, M. Inambao, W. Kilembe, C. Farquhar, C. Celum, Partners in Prevention
445 HSV/HIV Transmission Study Team, Determinants of per-coital-act HIV-1
446 infectivity among African HIV-1-serodiscordant couples. *J. Infect. Dis.* 205, 358–365
447 (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. Jermin,
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