

1
2
3 **PHYLOSCANNER: Inferring Transmission from Within- and**
4 **Between-Host Pathogen Genetic Diversity**
5

6 Chris Wymant^{*1,2}, Matthew Hall^{*1,2}, Oliver Ratmann², David Bonsall^{1,3,4}, Tanya
7 Golubchik^{1,4}, Mariateresa de Cesare⁴, Astrid Gall⁵, Marion Cornelissen⁶, Christophe
8 Fraser^{†1,2}, STOP-HCV Consortium, The Maela Pneumococcal Collaboration, and The
9 BEEHIVE Collaboration[‡]
10

11 ¹ Oxford Big Data Institute, Li Ka Shing Centre for Health Information and Discovery,
12 Nuffield Department of Medicine, University of Oxford, UK

13 ² Medical Research Council Centre for Outbreak Analysis and Modelling, Department of
14 Infectious Disease Epidemiology, Imperial College London, UK

15 ³ Peter Medawar Building for Pathogen Research, Nuffield Department of Medicine and
16 the NIHR Oxford BRC, University of Oxford, UK

17 ⁴ Wellcome Trust Centre for Human Genetics, Nuffield Department of Medicine,
18 University of Oxford, UK

19 ⁵ Department of Veterinary Medicine, University of Cambridge, UK

20 ⁶ Laboratory of Experimental Virology, Department of Medical Microbiology, Center for
21 Infection and Immunity Amsterdam (CINIMA), Academic Medical Center of the
22 University of Amsterdam, Amsterdam, The Netherlands
23

24 * Equal contribution

25 † To whom correspondence should be addressed: christophe.fraser@bdi.ox.ac.uk

26 ‡ Collaboration members listed in full at the end of the text.
27
28
29

30 **Abstract**

31
32 A central feature of pathogen genomics is that different infectious particles (virions, bacterial
33 cells, etc.) within an infected individual may be genetically distinct, with patterns of relatedness
34 amongst infectious particles being the result of both within-host evolution and transmission from
35 one host to the next. Here we present a new software tool, phyloscanner, which analyses
36 pathogen diversity from multiple infected hosts. phyloscanner provides unprecedented resolution
37 into the transmission process, allowing inference of the direction of transmission from sequence
38 data alone. Multiply infected individuals are also identified, as they harbour subpopulations of
39 infectious particles that are not connected by within-host evolution, except where recombinant
40 types emerge. Low-level contamination is flagged and removed. We illustrate phyloscanner on
41 both viral and bacterial pathogens, namely HIV-1 sequenced on Illumina and Roche 454
42 platforms, HCV sequenced with the Oxford Nanopore MinION platform, and *Streptococcus*
43 *pneumoniae* with sequences from multiple colonies per individual. phyloscanner is available from
44 <https://github.com/BDI-pathogens/phyloscanner>.

45

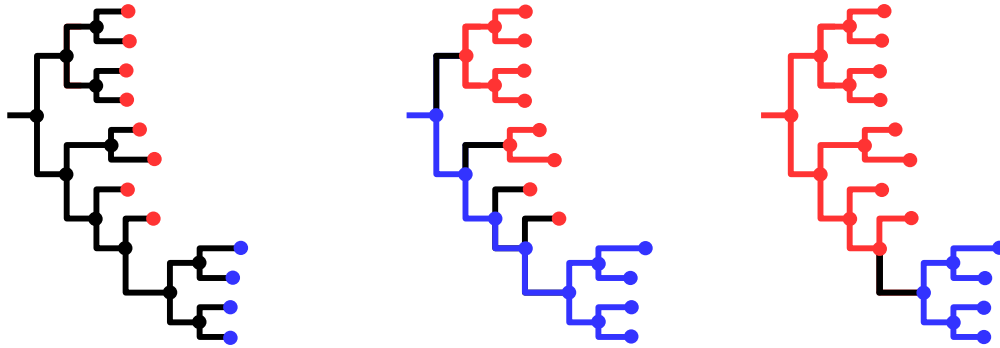
46 **Introduction**

47
48 The infectious transmission process imposes a hierarchical structure of relatedness on
49 pathogen genomes. The genotype of an individual infectious particle is the result of both within-
50 host evolution and transmission between hosts; a population sample collected from multiple
51 hosts, with multiple genotypes for each host, therefore simultaneously encodes the history of
52 both processes. Despite the existence of many tools for analysing pathogen genomes, none, to
53 our knowledge, are specifically adapted to exploiting this hierarchical genealogical structure.

54
55 A central aim of infectious disease epidemiology is the identification of risk factors for
56 transmission. The development of methods that use pathogen genomes to infer transmission
57 events, along with their direction, is therefore a priority. A critical recent insight is that including
58 multiple pathogen genomes per infected individual in such methods makes this inference easier:
59 it is equivalent to the simpler process of inferring ancestry (Romero-Severson et al. 2016).
60 Specifically, if a pathogen has passed from individual X to individual Y (either directly, or
61 indirectly via unsampled intermediate individuals) then all the pathogen particles sampled from
62 individual Y must be descended from the population of pathogen particles from individual X.

63 Inferring ancestral states is a standard problem in population genetics for which many methods
64 exist; the novel insight is that this standard approach may be used to infer the direction of
65 transmission. We illustrate this in Figure 1.

66



67

68 **[Figure 1: pathogen transmission direction via ancestral state reconstruction.** In the left-
69 hand phylogeny, tips are labelled red or blue according to their state: in our case the state of
70 interest is 'in which individual was this pathogen found?'. This state is known for the tips, but
71 can only be inferred for the internal nodes of the phylogeny: these represent coalescence
72 events, ancestors of the pathogens we have sampled. A change in state corresponds to a
73 change in the pathogen's host, i.e. to transmission, be it direct or indirect. The central phylogeny
74 shows one possible ancestral state reconstruction for which the root of the tree is blue, meaning
75 blue is ancestral to red. This requires at least four changes of state (shown with black branches)
76 – four sampled lineages transmitted from blue to red. The right-hand phylogeny shows one
77 possible ancestral state reconstruction for which the root of the tree is red, meaning red is
78 ancestral to blue. This requires only one change of state – one sampled lineage transmitted
79 from red to blue. Based on parsimony we would prefer the right-hand scenario.]

80

81 A frequently used approach in molecular epidemiology is to describe patterns of genetic
82 clustering - who is close to whom. However, identifying transmission pairs or clusters without
83 the ability to infer transmission direction - who infected whom - limits our ability to distinguish
84 risk factors for transmission from those for simply acquiring the pathogen. One approach for
85 inferring direction is to augment the sequence data with epidemiological data, and to couple
86 phylogenetic inference with mathematical models of transmission, for example references (Volz
87 and Frost 2013; Jombart et al. 2014; Hall et al. 2015; Didelot et al. 2017). However, this requires
88 strong assumptions from the model. In addition epidemiological data, such as dates and
89 location of sampling and reported contacts, are not always available, are subject to their own set

90 of uncertainties and errors, or are sometimes regarded as too sensitive to link to pathogen
91 genetic data.

92

93 Using multiple genotypes per host, and exploiting the link between transmission and ancestral
94 reconstruction, therefore promises an alternative and potentially powerful approach to molecular
95 epidemiology. Whilst several studies have used this idea to great effect on an ad hoc basis
96 (Numminen et al. 2014; Worby et al. 2016), no systematic or automatic tool has been developed
97 for this task.

98

99 Once multiple genotypes per host are included in a study, other questions present themselves
100 naturally, for example identifying multiply infected individuals. These may be defined as
101 individuals harbouring pathogen subpopulations resulting from distinct founder pathogen
102 particles. Multiple infections may be clinically relevant, for example in the case of Human
103 Immunodeficiency Virus 1 (HIV-1), dual infection is associated with accelerated disease
104 progression (Cornelissen et al. 2012). Multiple infections also represent unique opportunities for
105 pathogen evolution, especially for pathogens that recombine. Recombination between divergent
106 strains accelerates the generation of novel genotypes, and so potentially novel phenotypes. The
107 distinct pathogen strains in a multiple infection could have been transmitted simultaneously from
108 the same individual (if that individual harboured sufficient within-host diversity), or sequentially –
109 ‘super-infection’ – with each strain perhaps originating from a different transmitter. For HIV-1,
110 mathematical modelling has suggested that recombinants can reach high prevalence even
111 when the possibility of super-infection is restricted to a short window after initial infection, and
112 even when recombinants have no fitness advantage, if the epidemic is fuelled by a high-risk
113 core group (Gross et al. 2004).

114

115 Molecular epidemiology is being transformed by the advent of next-generation sequencing
116 (NGS; also called *high-throughput*) technologies (Goodwin et al. 2016). For many sequencing
117 protocols applied to pathogens with extensive within-host diversity, such as HIV-1 and Hepatitis
118 C Virus (HCV), the NGS output from a single sample can capture extensive within-host
119 diversity. Zanini et al. (Zanini et al. 2015) inferred phylogenies from NGS *reads* - fragments of
120 DNA - in windows along the genome for longitudinally sampled individuals infected with HIV-1,
121 to quantify patterns of within-host evolution over time. Here our focus will be on cross-sectional
122 datasets: by constructing phylogenies from NGS reads from multiple infected individuals at
123 once, within-host and between-host evolution can be resolved.

124

125 We present phyloscanner: a set of methods implemented as a software package, with two central
126 aims. The first is efficient computation of phylogenies with multiple genotypes per infected host,
127 and the second is analysis of such phylogenies and inference of biologically and
128 epidemiologically relevant properties from a set of related phylogenies. Multiple related
129 phylogenies arise naturally, either by sampling different portions of a genome, or in representing
130 uncertainty in phylogenetic inference (though bootstrapping, or sampling phylogenies from a
131 posterior distribution, for example). phyloscanner automatically performs the following steps:

- 132 1. Inference of between and within-host phylogenies from NGS data in multiple windows
133 along the pathogen genome (optionally skipped, if the user has such phylogenies
134 already);
- 135 2. Identification and removal of likely contaminant sequences;
- 136 3. Quantification of within-host diversity;
- 137 4. Identification of multiple infections;
- 138 5. Identification of crossover recombination breakpoints in NGS genotypes;
- 139 6. Ancestral host-state reconstruction from multiple phylogenies;
- 140 7. Identification of transmission events from ancestral host-state reconstructions.

141

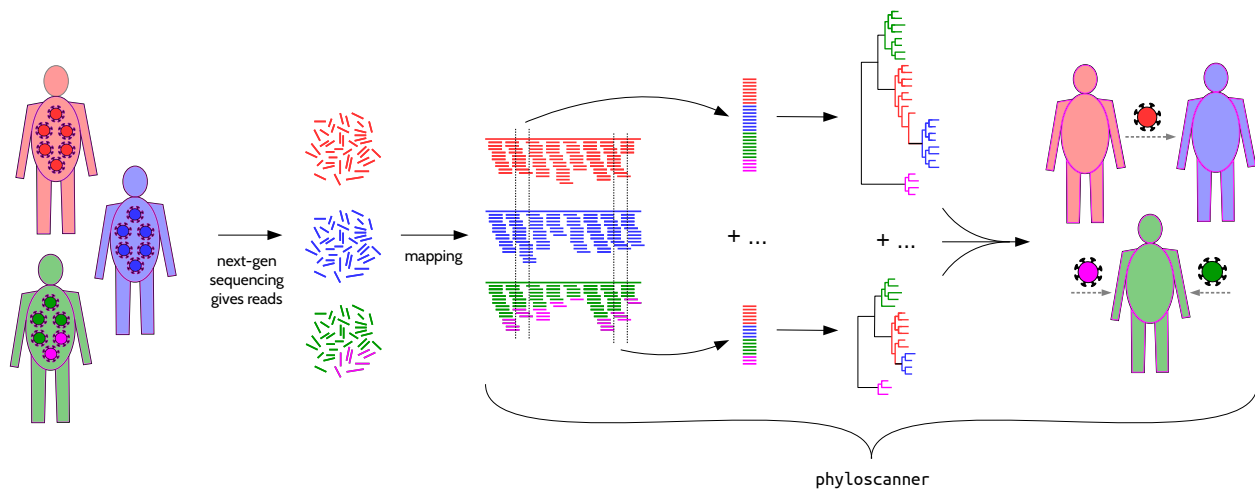
142 phyloscanner was intended for analysis of two distinct types of sequence data. Firstly for deep
143 sequencing data, in which NGS has produced reads from the population of diverse pathogens
144 represented in the sample. Secondly, for single-genome amplification (SGA), clonal sequencing
145 or bacterial colony picks, whereby laboratory methods are employed to separate the genomes
146 of individual pathogen particles prior to amplification and sequencing. Sequencing with primer
147 IDs (Jabara et al. 2011) may in some cases produce similar results at reduced costs. We also
148 considered haplotype reconstruction (Zagordi et al. 2011; Prabhakaran et al. 2014; Töpfer et al.
149 2014), i.e. bioinformatically inferring different haplotypes represented in the short reads of a
150 mixed sample, but in our hands this approach did not yield satisfactory results (analysis not
151 shown).

152

153 With SGA-style data, within- and between-host phylogenies can be directly inferred using
154 standard methods, and therefore phyloscanner is not necessary for step 1 in the process
155 described above. With deep sequencing data, reads for each sample must first be *mapped*
156 (placed at the correct location in the genome); thereafter phyloscanner begins by aligning reads

157 in windows of the genome that are matched across infected individuals, and inferring a
158 phylogeny for each window (Figure 2).

159



160

161 **[Figure 2: phyloscanner schematic for whole-genome deep sequence data.** In this
162 schematic, pathogens are sampled from the population infecting three hosts. NGS deep
163 sequencing produces reads, which are fragments of the genome sequence of one pathogen
164 particle (after amplification if necessary). Mapping to a reference means aligning each read to
165 the appropriate location in the genome; this must be done beforehand, as mapped reads are the
166 inputs to phyloscanner. phyloscanner produces alignments of reads in sliding windows along the
167 genome, automatically adjusting for the fact that the reference may be different for each sample.
168 Phylogenies are inferred for each alignment. These phylogenies are analysed separately using
169 ancestral host-state reconstruction (i.e. assigning hosts to internal nodes), and their information
170 is combined to give biologically and epidemiologically meaningful summaries. For example
171 here, we infer that the red individual infected the blue individual directly or indirectly, and the
172 green individual has two distinct pathogen strains.]

173

174

175 Results

176

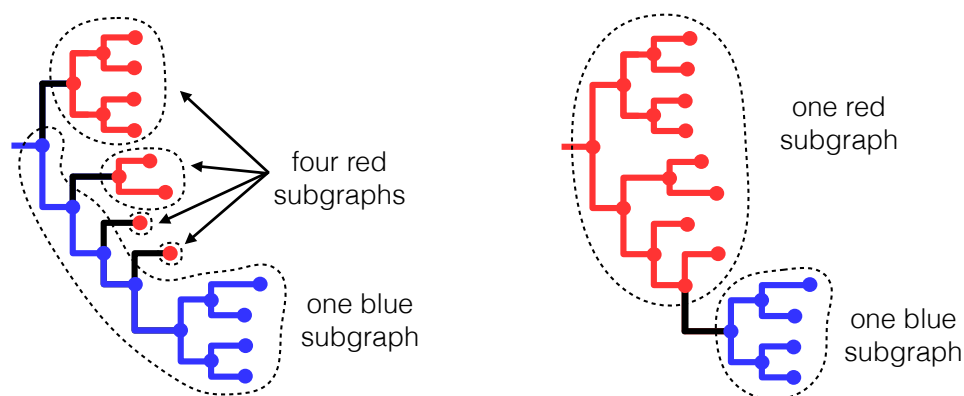
177 The best way to illustrate phyloscanner is through examples. We chose five datasets illustrating
178 different uses, pathogens, and sequencing platforms. We describe four in the main text, and
179 one in the Supplementary Information. These are far from systematic samples or population
180 surveys; they are small selections of infected individuals chosen to illustrate the different

181 conclusions that can be drawn using phyloscanner. We leave the application of phyloscanner to
182 large systematic population samples to future work.

183

184 Before presenting phylogenies for these data we introduce the term *host subgraph*. Host
185 subgraphs result from ancestral host-state reconstruction: they are defined as connected
186 regions of the phylogeny (tips and internal nodes, with the branches joining them) that have all
187 been assigned the same host state (i.e., the host that pathogen was in). See supplementary
188 section SI 1 for an explanation of the ancestral state reconstruction algorithm. Each subgraph
189 can be shown with a solid block of colour corresponding to that host, uninterrupted by colouring
190 associated with any other host. Figure 3 shows an example.

191



192

193 **[Figure 3: subgraphs defined by a given ancestral state reconstruction.** Here we show
194 again the two different ancestral state reconstructions on the same phylogeny from Figure 1,
195 this time illustrating the *host subgraphs* that these reconstructions define: connected regions of
196 the phylogeny that have all been assigned the same state (blue host or red host). Note that the
197 set of tips in a subgraph may or may not form a clade. In both of the above reconstructions, the
198 blue tips are contained in one subgraph and form a monophyletic group (one clade), whereas
199 the red tips form a polyphyletic group. The minimum number of clades needed to encompass all
200 and only the red tips is four, coinciding with the four red subgraphs in the left-hand
201 reconstruction.]

202

203

204 **Six illustrative HIV-1 infections, sequenced with Illumina MiSeq**

205

206 We used phyloscanner to analyse data from the BEEHIVE project (*Bridging the Evolution and*
207 *Epidemiology of HIV in Europe*), in which whole-genome samples from individuals with well-

208 characterised dates of HIV-1 infection are being sequenced, primarily to investigate the viral-
209 molecular basis of virulence (Fraser et al. 2014). We chose two groups of patients for detailed
210 investigation (presented in this subsection and the next), that together demonstrate interesting
211 features revealed by phyloscanner.

212
213 For the BEEHIVE samples, viral RNA was extracted manually from blood samples following the
214 procedure of Cornelissen *et al.* (Cornelissen et al. 2016). The RNA was reverse transcribed and
215 amplified using universal HIV-1 primers that define four overlapping amplicons spanning the
216 whole genome, then sequenced using the Illumina MiSeq platform, following the procedure of
217 Gall *et al.* (Gall et al. 2012; Gall et al. 2014). The resulting reads were mapped to a reference
218 constructed for each sample using IVA (Hunt et al. 2015) and shiver (Wymant et al. 2016),
219 producing input analogous to the illustration in Figure 2. See Materials and Methods for more
220 detail.

221
222 These mapped reads were analysed with phyloscanner using 54 overlapping windows, each 320
223 base pairs (bp) wide, covering the whole HIV-1 genome (approximately 9200 bp long; the
224 window entirely overlapping the variable V1-V2 loop in the envelope gene was not included due
225 to the richness of insertions and deletions, which leads to poor alignment). To increase
226 phylogenetic resolution and accuracy, we used the phyloscanner options to merge overlapping
227 paired-end reads into single, longer reads, and to delete drug resistance sites (Gatanaga et al.
228 2002; Johnson et al. 2011; Wensing et al. 2015) which are known to be under convergent
229 evolution.

230
231 Figure 4 shows the resulting phylogenies for four windows, chosen for clarity when visually
232 inspected. The phylogenies illustrate single infection (patient A), dual infection (patient B),
233 contamination (from the sample of patient C to the sample of patient D) and transmission (from
234 patient E to patient F, possibly via an unsampled intermediate individual). Colouring on each
235 phylogeny illustrates host subgraphs.

236
237 **Contamination.** Filtering for contamination is an important part of analysis of NGS data.
238 Contamination may be physical contamination of one sample into another, or low-level barcode
239 switching which occurs during the multiplexing and demultiplexing steps which are central to the
240 high throughput of NGS. phyloscanner uses two criteria to identify reads as likely contaminants
241 (either criterion is sufficient). The first is that they are exact duplicates of reads from another

242 patient, but much less numerous; the second is that they form an additional host subgraph
243 separated from the primary subgraph, but with too few reads to a call of multiple infection. This
244 The second means that the source of the contaminant reads need not be present in the
245 analysed dataset to infer contamination. These reads are flagged according to tuneable
246 parameters (which will depend on the precise sample and method used), and blacklisted from
247 further analysis (marked by pink crosses in Figure 4). We note that in general, phylogenetic
248 patterns associated with transmission are distinct from those associated with contamination: the
249 process of transmission is accompanied by within-host evolution in the recipient, whereas
250 contamination is not.

251

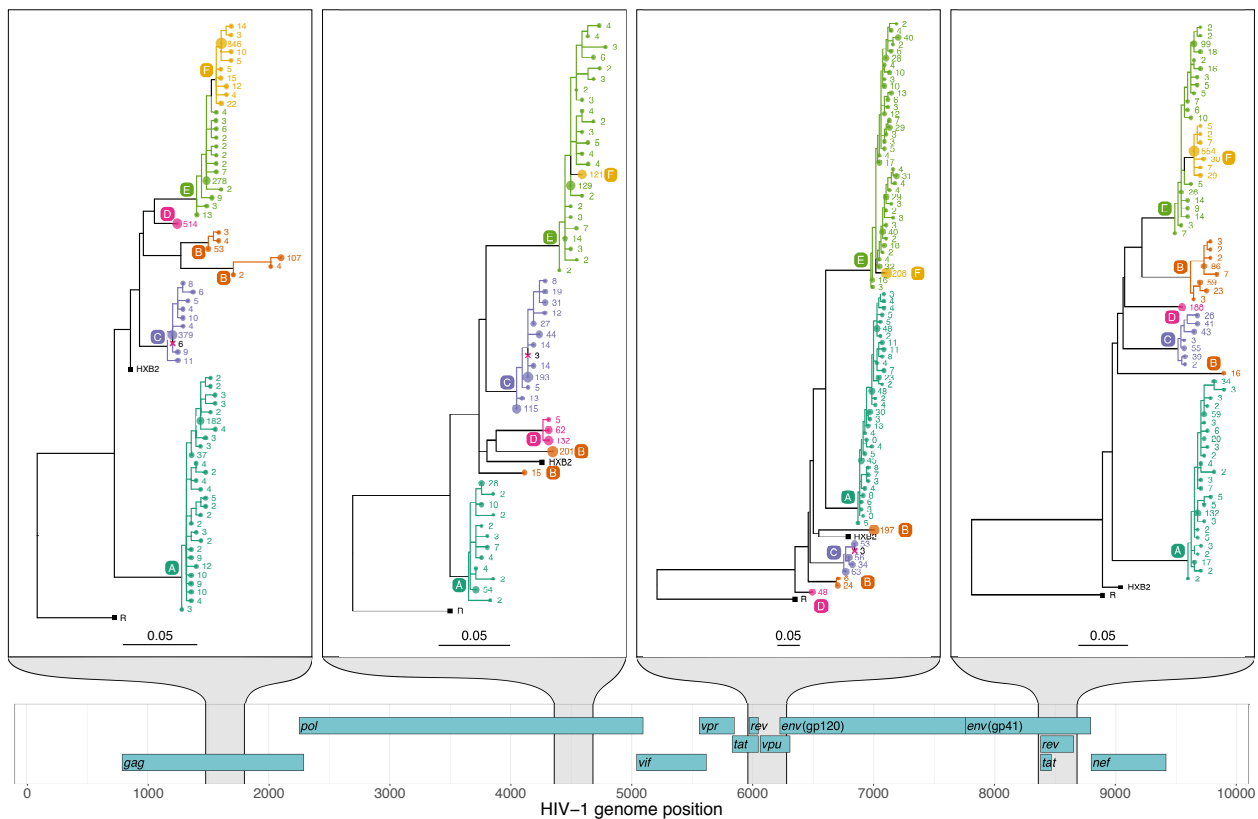
252 **Multiple infections.** If the phylogeny and host-state reconstruction are correct, the number of
253 subgraphs a patient has equals the number of founder pathogen particles with sampled
254 descendants (for example if this is 2, a dual infection is inferred). Sampling effects mean that
255 representatives of these multiple infections may not be present in all windows.

256

257 **Transmission.** Nodes of the phylogeny not in any patient's subgraph are coloured black in our
258 figures, as are branches connecting nodes not part of the same subgraph. These black regions
259 connect the different host subgraphs to each other, and so correspond to the pathogen jumping
260 between hosts; each region must contain one or more transmission events. They may, or may
261 not, correspond to the passage of the pathogen lineage through one or more unsampled hosts.
262 The probability of an indirect transmission will increase with the size of the black region and may
263 be best investigated by examining the subgraph relationships and branch lengths together.

264

265



266

267 **[Figure 4: phyloscanner analysis of four illustrative windows of the HIV-1 genome. A map of**

268 the HIV-1 genome is shown at the bottom with the nine genes in the three reading frames.

269 Phylogenies are shown for the four windows highlighted in grey, with scale bars measured in

270 substitutions per site. Tip labels are coloured by patient, as are all nodes assigned to that

271 patient by ancestral reconstruction, and the branches connecting these tips and nodes; a solid

272 block of colour therefore defines a single subgraph for one patient (see main text). The number

273 labelling each tip is the number of times that read was found in the sample, and the size of the

274 circle at each tip is proportional to this count. The count is after merging all identical reads and

275 reads differing by a single base pair (merging similar reads can be done for computational

276 efficiency, or as here, for presentational clarity). External references included for comparison

277 are shown with black squares. One is HXB2; the other, labelled R, is a subtype C reference

278 used to root each phylogeny. The six patients are labelled A through F. **Single infection:**

279 patient A is a singly infected; all reads from this patient form a single subgraph. **Dual infection:**

280 patient B is inferred to be dually infected, as is apparent by the fact that ancestral reconstruction

281 produces two subgraphs in each window. **Contamination:** patients C and D are both singly

282 infected, but we infer that some contamination has occurred from C to D. Patient D's sample

283 has a small number of reads that are identical to reads from patient C, but much less numerous.
284 Such reads are removed, but are shown here as crosses in the clade of patient C, for illustrative
285 purposes. **Transmission:** in all four windows shown here, the reads of patient F are seen to be
286 wholly descended from within the subgraph of reads of patient E. We infer that patient E
287 infected patient F, either directly, or indirectly via an unsampled intermediate. Patient F having a
288 single subgraph that is linked to patient E by a single branch indicates that the viral population
289 was bottlenecked down to a single sampled ancestor during transmission.]

290

291 **Genome-wide summary statistics.** In general, a phyloscanner analysis may produce a large
292 number of phylogenies and associated ancestral reconstructions. These can be output both as
293 annotated NEXUS format files, and as PDF files created with *ggtree* (Yu et al. 2017) for rapid
294 visual inspection. Statistics are calculated to summarise the wealth of information in the
295 phylogenies; these are shown for the 6 patients and 54 genomic windows in Figure 5. They
296 include measures of within-host diversity, measures that allow rapid identification of multiply
297 infected individuals, and a basic metric of recombination (defined in the supplementary section
298 S3).

299



300

301 **[Figure 5 - Summary statistics for six illustrative HIV-1 infected patients.** Each column

302 shows data from a single patient; each row is one or two statistics, plotted along the genome.

303 **Top row:** number of reads, and number of unique reads (corresponding to tips in the

304 phylogeny). **Second row:** the number of clades required to encompass all and only the reads

305 from that patient, and the number of subgraphs (see Fig. 3 for clarification of these quantities).

306 In many windows, though not all, the reads of patient B form two subgraphs: evidence of dual

307 infection. For patients C and E, we see a single subgraph but many clades. This is because of

308 the presence of reads from other patients (D and F, respectively, as seen in Fig. 4) inside what

309 would otherwise be a single clade, turning a monophyletic group into polyphyletic group (which

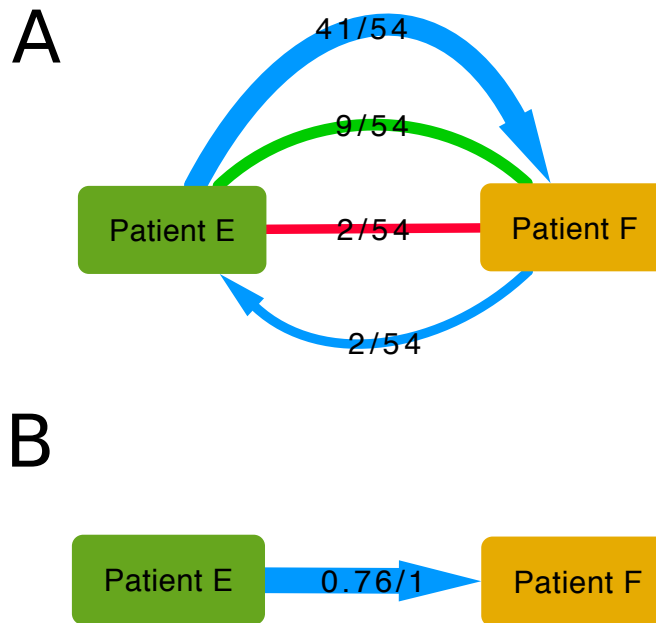
310 requires splitting in order to form clades). **Third row:** within-host divergence, quantified by mean
311 root-to-tip distance. Defining a patient's subtree as the tree obtained by removing all tips not
312 from this patient, we calculate root-to-tip distances both in the whole subtree and in just the
313 largest subgraph. For patient B, this distinction is substantial due to the very large distance
314 (~0.1 substitutions/site) between the two subgraphs of this dually infected patient. For singly
315 infected patients, divergence may correlate with time since infection. **Fourth row:** for each
316 window, a stacked histogram of the proportion of reads in each subgraph. For patient B, when
317 two subgraphs are present, an appreciable proportion of reads are in the second one (mean
318 12%). The histogram is absent in the window that was excluded by choice. **Bottom row:** a
319 score based on Hamming distance (between 0 and 1) of the extent of recombination in that
320 window. The highest score across all six patients and all windows is indicated with an orange
321 diamond; the reads giving rise to this score are shown in supplementary Figure S6.]

322
323 In a single window, phyloscanner classifies two patients to be related if they are adjacent (see
324 supplementary section SI) and optionally, also "close", i.e. that their subgraphs are within a
325 prespecified patristic distance of each other. Relationships are further categorised by the
326 ancestry, or lack of it, that is suggested by the tree topology. To summarise transmission across
327 all windows, phyloscanner output summarises the number of windows in which each pair of
328 patients are related, and the topological nature of that relationship. This allows the complete set
329 of relationships between all patients in the dataset to be visualised in graph form. For example,
330 in this dataset, only two of the six patients, E and F, are related in at least half of the windows.
331 In Figure 6A the counts of the different topological relationships between these two patients are
332 displayed. With many links between many patients these graphs become difficult to interpret
333 visually; a threshold on the number of windows for links to be displayed is therefore helpful.
334 phyloscanner also produces a second version of the graph simplified further, shown in Figure
335 6B. Here a single link appears if relatedness of any type is present in 50% of windows, and that
336 link is an arrow if transmission in that direction is inferred in at least 33% of windows. (The 50%
337 and 33% thresholds are defaults that can be changed.) These relationship diagrams were
338 plotted using Cytoscape 3.5.1 (Shannon et al. 2003).

339
340 Diagrams such as those in Figure 6, when extended to greater numbers patients, will not always
341 represent a single, coherent transmission tree amongst all the patients in the dataset (as can be
342 seen in Figures 7 and 9). Instead, they simply summarise each pairwise relationship. As a
343 result, we refer to them as "relationship graphs". The inference of a single, most probable

344 transmission tree over all windows is complicated by the presence of multiple infections,
345 incomplete transmission bottlenecks, and missing data for some patients in some windows. To
346 our knowledge, no method yet exists to produce a consensus transmission history that takes
347 into account all these possibilities.

348



349

350 **[Figure 6 – Relationship graphs: visual representations of the relationship between two**
351 **connected patients infected with HIV-1.** The power of phyloscanner in studying transmission
352 events comes from aggregating information over many within- and between-host phylogenies, in
353 this case obtained from different windows of the whole HIV-1 genome. In the top diagram, the
354 outcomes from all 54 windows are shown. The top blue arrow shows that in 41 windows, patient
355 E was inferred to be ancestral to patient F, with a single bottleneck. The bottom blue arrow
356 shows that in 2 windows the reverse was true – F was ancestral to E. The undirected red line
357 shows that in 2 windows, the patients were linked by “complex” ancestry, with the direction
358 unclear. The undirected green line shows that in 9 windows the patient subgraphs were
359 adjacent and close, but no ancestry was implied by the topology. In no window was
360 transmission of more than one lineage inferred, and in no window were the patients distant and
361 unlinked. (See supplementary section SI 1 for more details on these categories.) A simplification
362 of these relational data is shown in the bottom diagram, with a single directed arrow. The first
363 number indicates the proportion of windows supporting transmission in the direction of the
364 arrow, and the second number indicates the proportion of windows supporting transmission in
365 either direction.]

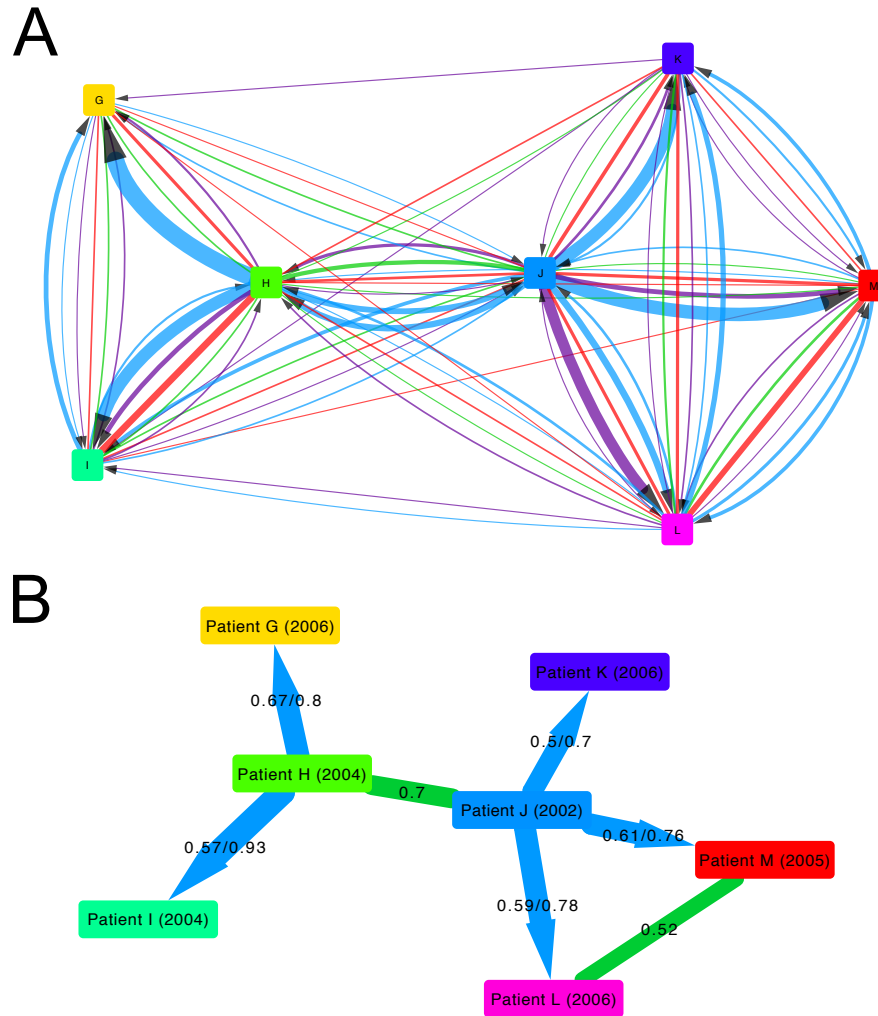
366

367 **Resolving the transmission pathway within a HIV-1 phylogenetic cluster**

368

369 To illustrate the resolution into the transmission process that can be obtained by phyloscanner,
370 we chose a set of 7 patients from the BEEHIVE study that were found to be closely connected
371 in the chain of transmission (Fig. 7). 3 of the patients' samples were sequenced with Illumina
372 MiSeq and 4 with Illumina HiSeq; the resulting reads were processed and mapped using IVA
373 and shiver as previously, with the mapped reads given as input to phyloscanner. phyloscanner
374 summarises all the pairwise relationships between individuals in each window (Figure 7A),
375 suggesting a complex network. However, we find that when we focus on the most likely
376 inferences of source attribution (Figure 7B), phyloscanner largely resolves a complex set of
377 pairwise relationships into a coherent transmission network, that is consistent with the years of
378 seroconversion. However, this is not guaranteed to be the case: an exception is the triangle
379 connecting Patients J, L and M, where there is too much uncertainty in the relationships
380 amongst the triplet to resolve their ancestry.

381



382

383 [Figure 7 - The relationship between 7 patients infected with HIV-1. The colouring and
384 numbers on the arrows connecting patients are as in Figure 6; in addition, the lower diagram
385 here contains undirected green lines as well directed blue lines. These green lines suggest that
386 the pair are close in the transmission network but with unknown transmission direction; the
387 single number on the line indicates the proportion of windows supporting this. The known or
388 estimated year of infection is shown in parentheses after each patient's label.]

389

390 HIV-1 sequenced with Roche 454

391

392 A subset of patients from the BEEHIVE study were also sequenced using the Roche 454
393 platform; results from their analysis with phyloscanner are in Supplementary Information section
394 SI 2.

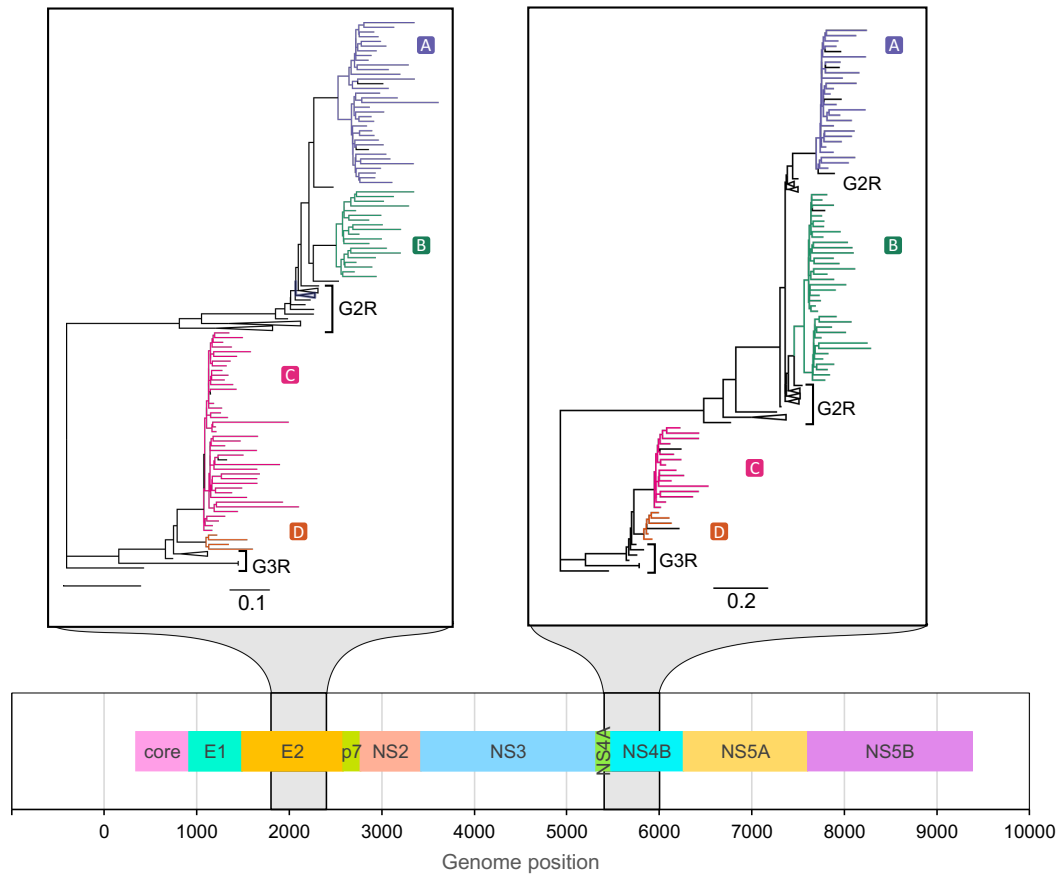
395

396 **HCV sequenced with Oxford Nanopore MinION**

397

398 To further illustrate phyloscanner's applicability to different sequencing platforms and also
399 different pathogens, we used it to analyse HCV viral data sequenced using the Oxford
400 Nanopore MinION device. Plasma samples were obtained from four patients in the BOSON
401 study (Foster et al. 2015), a phase 3 randomized trial of antiviral therapy with sofosbuvir (trial
402 registration NCT01962441). Sequencing was performed using RNAseq-based methods
403 previously described for Illumina (Bonsall et al. 2015) and adapted for the MinION device.
404 Briefly, plasma-derived RNA was reverse transcribed, then sequencing libraries were prepared
405 for each sample using Oxford Nanopore adapters and customised barcoded primers. These
406 were pooled and enriched using HCV-specific nucleotide baits before sequencing on a MinION
407 R9.0 flow cell. Viral sequences were identified and mapped using BLASTN (Altschul et al.
408 1990), standard reference sequences and BWA (Li and Durbin 2009). See Materials and
409 Methods for more details. The resulting BAM files were used as input for phyloscanner, with a
410 window size of 600 bp and no overlap between windows. Nanopore sequencing platforms are
411 capable of producing longer inserts than those of Illumina, at the cost of a higher error rate
412 (approximately 10% erroneous base calls). Despite this error, phyloscanner could
413 phylogenetically resolve the within- and between-host evolution, shown in Figure 8.

414



415

416 **[Figure 8 - phyloscanner analysis of two illustrative windows of the HCV genome.**

417 Sequence data from four individuals was obtained with the Oxford Nanopore MinION device. A
418 continuous region of the phylogeny with the same colour shows a subgraph for one patient (see
419 main text). Black tips were flagged as contamination and excluded. Patient-derived sequences
420 clustered with respective genotype 2 and genotype 3 references (G2R, G3R) as expected from
421 the virus genotypes known from the clinical information available for participants. Two windows,
422 600 bp in length, are shown for the E2 and NS4B genes at positions given by the genome map
423 (bottom panel).]

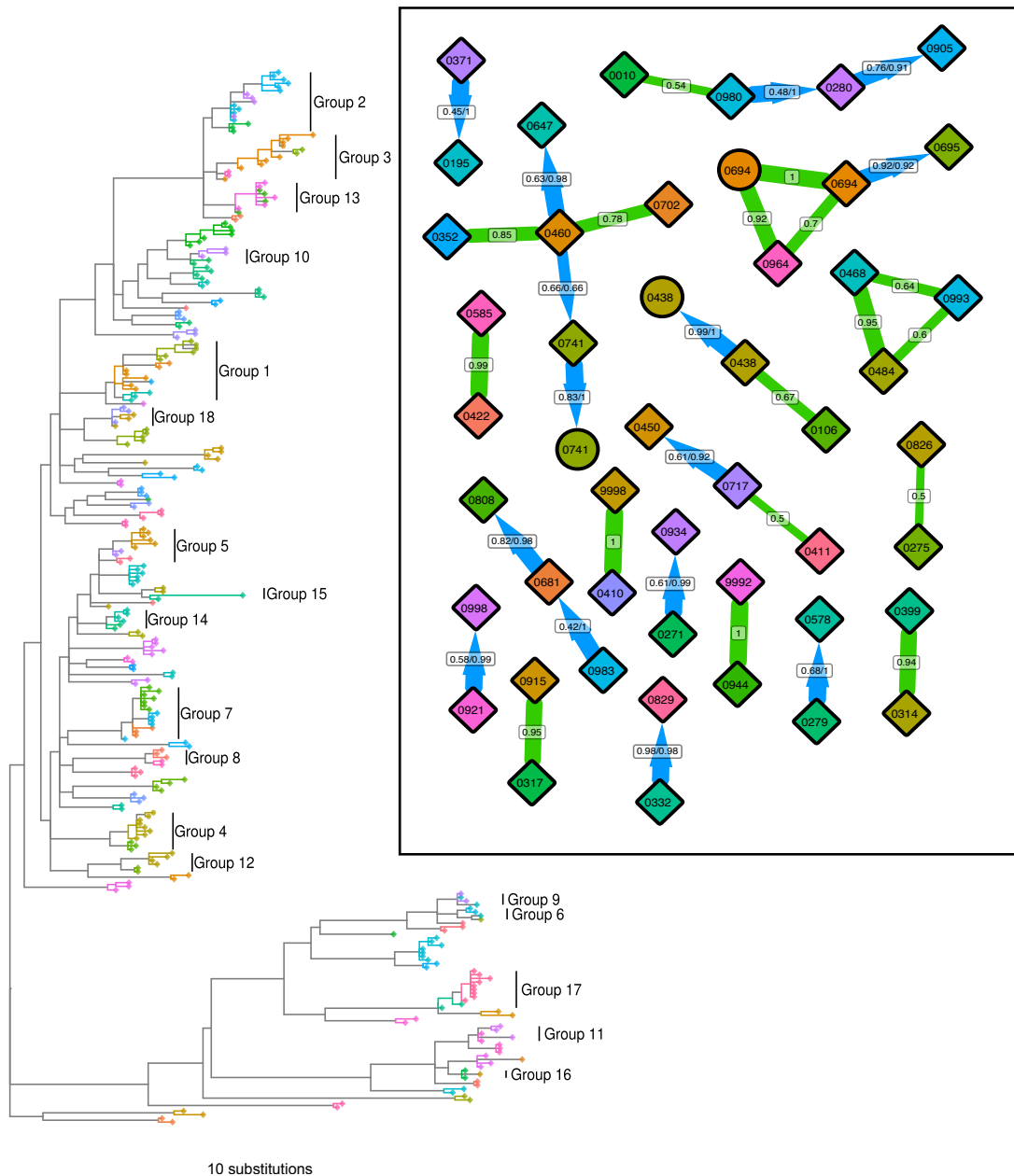
424

425 **Multiple colony picks per carrier of *S. pneumoniae***

426

427 phyloscanner's analysis of phylogenies need not be restricted to those derived from deep
428 sequencing data in different windows of the genome: it can also be applied to datasets where
429 within-host diversity is captured by SGA or sequences from multiple colony picks per individual.
430 We illustrate this approach with the *S. pneumoniae* data of Croucher et al. (Croucher et al.
431 2016), specifically the BC1-19F cluster. This dataset consists of 286 sequences from 92

432 individuals carrying the bacterium (with multiple colonies per carrier). These were sequenced
433 with Illumina HiSeq, though for SGA data sequencing platform is largely irrelevant to
434 interpretation, since each sequenced sample should not contain any real within-sample diversity
435 by design. Genomes were processed with Gubbins (Croucher et al. 2015) to remove
436 substitutions likely to have been introduced by recombination. As each of these sequences is a
437 whole genome (unlike the short reads produced by NGS), we did not split the genome into
438 windows to be analysed separately. Instead, we represented phylogenetic uncertainty by
439 generating a posterior set of 100 phylogenies using MrBayes 3.2.6 (Ronquist et al. 2012) and
440 analysed these with phyloscanner. Ancestral state reconstruction was performed on each
441 posterior phylogeny independently, relationships between carriers identified, and the results
442 summarised over the entire set. In each phylogeny, carriers were inferred as being related if the
443 minimum patristic distance between two nodes from the subgraphs associated with each was
444 less than 7 substitutions and they were categorised as adjacent (explained in Supplementary
445 Information section SI 1.5). This distance threshold was selected to demonstrate the method as
446 it picked out obvious clades in the phylogeny as groups, and was not chosen to imply direct
447 transmission. Retaining such relationships where they existed in at least 50% of posterior
448 phylogenies revealed 18 separate groups of carriers whose bacterial strains were closely
449 related (see Fig. 9).
450



[Figure 9 - Phylogeny and relationships between *S. pneumoniae* carriers. The phylogeny shown is the MrBayes consensus tree. Tip shapes are coloured by carrier, with mother and infant pairs sharing the same colour; diamonds represent infants and circles mothers. All nodes assigned to a carrier by ancestral reconstruction, and the branches connecting these tips and nodes, are given the same colour as that carrier's tips; a solid block of colour therefore defines a single subgraph for one carrier (see main text). Regions of the phylogeny not in any carrier's subgraph are grey. These regions connect carriers' subgraphs to each other, and so each must contain one or more transmission events. The carrier relationship diagram (inset) displays the relationships between the carriers in 18 identified groups, in the same fashion as in Figures 6

461 and 7, except that here the numbers represent the proportion of phylogenies from the posterior
462 set, rather than the proportion of genomic windows in which both patients have sequence data.
463 The clades representing these 18 groups are labelled in the phylogeny.]

464

465 Note that if some residual signals of recombination remain after processing with Gubbins,
466 analysing the full-length genomes in windows by choice (rather than by necessity, as with short-
467 read NGS data) could mitigate this effect at the cost of reduced phylogenetic resolution in each
468 window. The merits of this could be explored in a dedicated analysis of such a dataset; here we
469 simply illustrate application of phyloscanner to full-length sequences as opposed to genomic
470 windows.

471

472

473 Discussion

474

475 Improving our understanding of the transmission of pathogens is valuable for identifying
476 epidemiological risk factors – the first step for targeting public health interventions for efficient
477 impact. Phylogenetic analysis of one pathogen sequence per infected individual may identify
478 clusters of similar sequences that are expected to be close in a transmission network. However,
479 nothing is learned about the direction of transmission within the network. Indeed it may be that
480 none of the individuals transmitted the pathogen to anyone else, and they were all infected by a
481 common individual who was not sampled. Through automatic fitting of maximum-likelihood
482 evolutionary models to within- and between-host genetic sequence data, phyloscanner enhances
483 resolution into the pathogen transmission process. An evidence base is built up by analysing
484 many phylogenies, notably through consideration of NGS reads in windows along the
485 pathogen's genome. The relationship between infected individuals is no longer quantified by a
486 single number summarising closeness, but by a rich set of data resulting from ancestral host-
487 state reconstruction for each phylogeny.

488

489 Romero-Severson *et al.* (Romero-Severson et al. 2016) demonstrated the utility of parsimony
490 for the assignment of ancestral hosts to internal nodes in a phylogeny containing many tips from
491 two infected individuals, for simulated HIV-1 data. We have continued with this approach,
492 developing it for suitability for real sequence data from many infected individuals. In particular
493 we allow for (i) contamination, (ii) multiple infections, and (iii) the possible presence of
494 unsampled hosts in the tree. Details of two such parsimony algorithms, available for use in

495 phyloscanner, are presented in the supplementary section SI 1. Parsimony has the advantage
496 that a reconstruction can be completed in reasonable computational time even for phylogenies
497 with tens of thousands of tips. Other methods of reconstructing the host state of internal nodes
498 could also be suitable and may be added to the package in future. Our identification of
499 contamination and multiple infections is highly valuable in its own right: the former because this
500 is critical for any empirical study of within-host diversity, and the latter because such individuals
501 may be special cases clinically and for pathogen evolution. Transmission of multiple distinct
502 pathogen strains may occur simultaneously, or sequentially – ‘super-infection’. phyloscanner can
503 detect both cases, though distinguishing them is difficult without longitudinal sampling (it could
504 be possible through inference of timed trees, or using the diversity of each separate infection as
505 a proxy for its age).

506
507 Great care must be taken to correctly interpret the ancestry of pathogens infecting individuals.
508 Even if ancestry were established beyond any doubt, individual X’s pathogen being ancestral to
509 individual Y’s pathogen does not imply that X infected Y: the pathogen could have passed
510 through unsampled intermediate hosts. Nevertheless the ancestry does provide valuable
511 epidemiological information, as X has been identified as a transmitter (and Y a recipient not far
512 down the same transmission chain). Finding likely transmitters in a large population cohort
513 would allow risk factors to be identified and quantified.

514
515 Furthermore, inference of ancestry is itself subject to uncertainty. The inference of ancestry
516 depends on the correct rooting of the phylogeny, in order that the direction in which evolution
517 proceeded over time is known. Molecular clock analyses (such as implemented in TempEst
518 (Rambaut et al. 2016)) can aid correct rooting when the sampling dates of the tips of the
519 phylogeny are known.

520
521 The relationships between infected individuals are inferred by phyloscanner across many
522 phylogenies, for example those constructed from NGS reads in windows along the pathogen
523 genome. By analysing many phylogenies, phyloscanner mitigates the effect of random error - any
524 error that is independent in each phylogeny. We therefore give greater credibility to those
525 relationships observed many times than to those observed only once. However, systematic
526 error may arise, for example, due to different patients being sampled at different stages of
527 infection, with different amounts of within-host diversity to analyse (Romero-Severson et al.
528 2016). Given uncertainties in any individual assignment, we recommend phyloscanner for

529 population-level analyses, rather than focussing on isolated transmission events (as we have
530 done here, for simplicity in explaining the method).

531

532 The fraction of genomic windows in which a given relationship is inferred between individuals
533 (for example A infecting B directly or indirectly), is not equal to the probability of that relationship
534 being true. However it provides a measure of the robustness with which the available data
535 support that conclusion. This is analogous to bootstrapping – sampling with replacement from
536 the same sequence alignment, to create a set of similar phylogenies. Here however, different
537 windows of the genome make use of different sequence data. Given the potential for
538 disagreement between different windows due to genuine biological variation, imperfect
539 sequencing procedures etc., agreement between a fraction x of (non-overlapping) windows is a
540 stronger statement of robustness than agreement between a fraction x of bootstraps.
541 Identification of transmission events with phyloscanner will involve false positives and false
542 negatives; these will be context dependent, depending on how strictly transmission thresholds
543 are defined (which balance sensitivity and specificity) and on the inclusion of sequences similar
544 to those being investigated. We will illustrate this in two works in preparation examining large
545 population studies.

546

547 Whilst our emphasis has been on extracting broad-brush information from the rich within-and-
548 between host phylogenies, these phylogenies contain more information that could be used in
549 future research. A specific example is that by resolving the transmission event at a finer level of
550 genetic detail, it is possible to identify which pathogen genotypes are typically transmitted and
551 which ones are not, with potential relevance for vaccine design.

552

553 By providing a tool for automatic phylogenetic analysis of NGS deep sequencing data, or
554 multiple genotypes per host generated by other means, we aim to simplify identification of
555 transmission, multiple infection, recombination and contamination across pathogen genomics.

556

557 **Materials and Methods**

558

559 **Generation and assembly of the BEEHIVE Illumina data**

560

561 Viral RNA was extracted manually from blood samples following the procedure of Cornelissen *et*
562 *al.* (Cornelissen *et al.* 2016). RNA was amplified and sequenced according to the protocol of

563 Gall *et al.* (Gall *et al.* 2012; Gall *et al.* 2014). Briefly, universal HIV-1 primers define four
564 amplicons spanning the whole genome. 5 µl of amplicon I was pooled with 10 µl each of
565 amplicons II–IV. Libraries were prepared from 50 to 1000 ng DNA as described in Quail *et al.*
566 (Quail *et al.* 2008; Quail *et al.*), using one of 192 multiplex adaptors for each sample. Paired-end
567 sequencing was performed using an Illumina MiSeq instrument with read lengths of length 250
568 or 300 bp, or in the ‘rapid run mode’ on both lanes of a HiSeq 2500 instrument with a read
569 length of 250 bp.

570

571 For each sample, the reads were assembled into contigs using the *de novo* assembler IVA. The
572 reads and contigs were processed using shiver as described previously (Wymant *et al.* 2016). In
573 summary: non-HIV contigs were removed based on a BLASTN search against a set of standard
574 whole-genome references (Kuiken *et al.* 2012). Remaining contigs were corrected for assembly
575 error then aligned to the standard reference set using MAFFT (Kato *et al.* 2002). A tailored
576 reference for mapping was then constructed for each sample using the contigs, with gaps
577 between contigs filled by the corresponding part of the closest standard reference. The reads
578 were trimmed for adapters, PCR primers and low-quality bases using Trimmomatic (Bolger *et al.*
579 2014) and fastaq (<https://github.com/sanger-pathogens/Fastaq>). Contaminant reads were
580 removed based on a BLASTN search against the non-HIV contigs and the tailored reference.
581 The remaining reads were then mapped to the tailored reference using SMALT
582 (<http://www.sanger.ac.uk/science/tools/smalt-0>).

583

584 **Generation and assembly of the HCV Oxford Nanopore MinION data**

585

586 Viral RNA was extracted from plasma using the NucliSENS® easyMAG® total nucleic acid
587 extraction system (Biomerieux) and sequencing libraries were prepared using a modified
588 version of an RNA-seq based protocol with a virus enrichment step. Briefly, the NEBNext®
589 Ultra™ Directional RNA Library Kit (New England Biolabs, Ipswich, MA, USA) was used to
590 generate cDNA from 5ul of total RNA. The NEBNext® Ultra™ II End Repair/dA-Tailing Module
591 and Blunt/TA Ligase (New England Biolabs, Ipswich, MA, USA) were used for end repair of
592 dsDNA and ligation of PCR adapters (Oxford Nanopore Technologies) to allow for 18 cycles of
593 PCR using custom barcoded primers with a post-PCR clean-up with 1x Ampure XP (Beckman
594 Coulter, Pasadena, CA, USA). Each library was quantified by Quant-iT™ Qubit® dsDNA HS
595 Assay Kit and size distribution analysed using Agilent TapeStation High Sensitivity D5000
596 ScreenTape System. Approximately equimolar quantities of each library were pooled to a total

597 of 500 ng mass and processed for probe enrichment using customized xGen® Lockdown®
598 120mer probes specific to HCV (Integrated DNA Technologies, Inc., Coralville, Iowa, USA) and
599 a modified Roche NimbleGen protocol for hybridization of amplified sample libraries with a
600 shorter 4 hours hybridization time and on-bead post-enrichment PCR (12 cycles). The enriched
601 pool was prepared for sequencing on a MinION R9.0 flow cell using the SQK-NSK007 2d
602 ligation kit. Raw fasta5 sequence files were base called and demultiplexed using Metrichor
603 software. Viral sequences were identified and trimmed using a BLASTN search of the Los
604 Alamos database of HCV genotype references (Kuiken et al. 2005), then mapped to the closest
605 matching reference using BWA (with the command `bwa mem -x ont2d`). Consensus sequences
606 were called from the bam files and used as references for a second iteration of read mapping.

607

608 **The phyloscanner Method**

609

610 For application of phyloscanner to deep sequence NGS data, the required input is a set of files in
611 BAM format (Li et al. 2009) each containing the reads from one sample that have been mapped
612 to a reference, and a choice of genomic windows to examine. A sensible choice of windows
613 would normally tile the whole genome, perhaps skipping regions that are rich in insertions and
614 deletions (leading to poor sequence alignment). Windows should be wide enough to capture
615 appreciable within-host diversity, but short enough for some reads to fully span them; options in
616 the code help to inform the user's choice. There is no lower limit to the length of reads given as
617 input, however as read length decreases, phylogenetic resolution will suffer. phyloscanner
618 determines the correspondence between windows in different BAM files by aligning the mapping
619 references in the BAM files. Using the same reference for mapping all samples would negate
620 the need for this step, but it is of paramount importance to tailor the reference to each sample
621 before mapping to minimise biased loss of information (Wymant et al. 2016). For each window
622 in each BAM file, all reads (or inserts, if reads are paired and overlapping) fully spanning the
623 window are extracted using pysam (<https://github.com/pysam-developers/pysam>) and trimmed
624 to the window edges, then identical reads are collapsed to a single read, giving a set of unique
625 reads each with an associated count (i.e. the number of reads with identical sequence). A basic
626 metric of recombination is calculated by maximising, over all possible sets of three sequences
627 and all possible recombination crossover points, the extent to which one of the three sequences
628 resembles one of the other two sequences more closely on the left and resembles the other
629 sequence more closely on the right. Further detail is provided in the supplementary section SI 3.
630 In each window, each sample's set of unique reads is checked against every other sample's set,

631 with exact matches flagged to warn of between-sample contamination in the analysed dataset;
632 all unique reads are then aligned with MAFFT, and a phylogeny is inferred with RAxML
633 (Stamatakis 2014).

634

635 phyloscanner contains many options to customise processing and maximise the information
636 extracted from reads and phylogenies. Standard reference genomes can be included with the
637 reads for comparison. User-specified sites can be excised to mitigate the effect of known sites
638 under selection on phylogenetic inference. Greater faith can be placed in the reads by trimming
639 low-quality ends and wholly discarding reads that are low-quality, improperly paired, or rare.
640 Reads in the same sample that differ from each other by less than a specified threshold can be
641 merged into a single read to increase the speed of downstream processing. Overlapping paired
642 reads can be merged into a single longer read for greater phylogenetic resolution. Every option
643 of RAxML can be passed as an option to phyloscanner, for example specifying the evolutionary
644 model to be fitted, or multithreading.

645

646 Optionally, the user may skip inference of phylogenies from files of mapped reads, and instead
647 directly provide as input a phylogeny or a set of phylogenies generated by any other method.

648

649 To analyse phylogenies, phyloscanner required that they are rooted. This can be done manually,
650 or if the phylogenies were constructed by phyloscanner from mapped reads, rooting can be
651 achieved by providing one or more additional reference sequences with the mapped reads, and
652 choosing one of these to use as an outgroup. The outgroup should be sufficiently distant from
653 all sampled isolates that we can assume the most recent common ancestor of it and every
654 isolate (i.e. the root of the whole tree) was not present in any of the sampled individuals.

655

656 Each phylogeny analysed is annotated with a reconstruction of the transition process using a
657 modified maximum-parsimony approach to assign internal nodes to hosts or to an extra
658 “unassigned” state. The latter is given to lineages that either must have infected a host outside
659 the dataset, or to those where the situation is sufficiently ambiguous that this cannot be ruled
660 out. An important parameter of the reconstruction, designated k , is used to help identify dual
661 infections and contaminants. It acts as a penalty, in the parsimony algorithm, for the
662 reconstruction of single infections showing unrealistic within-host diversity. A suitable value of k
663 will depend on the pathogen under study, but as a rule of thumb, we suggest estimating a level
664 of pairwise genetic diversity that it would be unrealistic to see in an infection from a single

665 source, and using the reciprocal of this for k . In situations where the phyloscanner user is
666 confident that dual infections and contaminants are not present, k can be set to zero, in which
667 case no penalty for within-host diversity is applied.

668

669 The results of the reconstruction can be represented as a visualisation of the partial pathogen
670 transmission tree by the process of ‘collapsing’ each subgraph (i.e. each set of adjacent nodes
671 with the same reconstructed host; see supplementary Fig. S3) into a single node of a new tree
672 structure. This “collapsed tree” is then analysed to identify relationships between each pair of
673 infected individuals, according to the following categories:

674

- 675 1. Minimum distance: what is the smallest patristic distance between a phylogeny node
676 assigned to one host and a node assigned to the other?
- 677 2. Adjacency: is there a path on the phylogeny that connects the two individuals’ subgraphs
678 without passing through a third individual? (“Unassigned” nodes do not interrupt
679 adjacency.)
- 680 3. Topology: how are the regions from each individual arranged with respect to each other?
681 (See supplementary Fig. S4.)

682

683 Combinations of these properties can be used to develop criteria which identify individuals who
684 are closely linked in the transmission chain. For example, two individuals that are adjacent and
685 within a suitable distance threshold are likely to be either a transmission pair, or infected via a
686 small number of unsampled intermediaries. If the distance between subgraphs is large, on the
687 other hand, separation by unsampled hosts in the chain of transmission is likely even if they are
688 adjacent. The nature of the topological relationship between them may suggest a direction of
689 transmission, or be equivocal.

690

691 An individual having multiple subgraphs suggests multiple infection, with the ancestor node of
692 each subgraph inferred to be a distinct founder pathogen particle (the ancestor of that sampled
693 subpopulation). It can be difficult to distinguish a dual infection from a sample that has been
694 contaminated by another sample not present in the current data set (i.e. where contamination is
695 not visible as exact duplication of another individual’s read). For NGS data we make the
696 distinction in each phylogeny based on thresholds on read counts: outside of the subgraph
697 containing the greatest number of reads, any additional (‘minor’) subgraph is designated as
698 contamination and ignored if the number of reads it contains is below an absolute threshold, or

699 below a threshold relative to the read count in the largest subgraph. By default, minor
700 subgraphs with read counts exceeding both thresholds are kept, providing evidence for the
701 presence of multiple distinct subpopulations in that genomic window. (Alternatively, a
702 phyloscanner option allows all minor subgraphs to be entirely removed from consideration).
703 Zanini et al. (Zanini et al. 2015) discarded reads suspected of being contamination by
704 calculating each read's Hamming distance from the consensus, plotting the distribution of these
705 distances, and discarding reads giving rise either to a second peak or to a 'fat tail' (taken to be
706 recombinant reads). This approach is not appropriate when the data set may contain multiply
707 infected individuals, for example for a dual infection we wish to keep the reads from each of two
708 distinct groups that may be separated by a large distance.

709

710 **The phyloscanner Code**

711

712 phyloscanner is freely available at <https://github.com/BDI-pathogens/phyloscanner>. It is written in
713 Python and R, but can be run from the command line so that no knowledge of either language is
714 required. Inference of within- and between-host phylogenies from BAM-format mapped reads is
715 achieved with a single command of the form

716 `phyloscanner_make_trees.py ListOfBamsAndRefs.csv --windows 1,300,301,600,...`

717 where `ListOfBamsAndRefs.csv` lists the BAM files to be analysed and the fasta-format references
718 to which the reads were mapped, and the `--windows` flag above specifies analysis of the
719 genomic windows with coordinates 1-300, 301-600, ...

720 Analysis of those trees is achieved with a single command of the form

721 `phyloscanner_analyse_trees.R TreeFiles OutputLabel [choice of ancestral state reconstruction]`.

722

723 Included with the code is simple simulated HIV-1 data for ease of immediate exploration of
724 phyloscanner. Within-host evolution was simulated using SeqGen (Rambaut and Grassly 1997);
725 each resulting sequence was then converted into error-free fragments that were mapped back
726 to the founding sequence, giving BAM-format files suitable as input for phyloscanner. We also
727 created BAM-format files by using shiver to process publicly available HIV-1 reads sequenced
728 with Illumina MiSeq. A tutorial walking the user through a simple application of phyloscanner to
729 the simulated data, and a more sophisticated application to this real public data, is available
730 from the GitHub repository with the code itself.

731

732 Running phyloscanner on the six HIV-1 samples presented in the first results section took 18
733 minutes on one core of a standard laptop, 10 minutes of which was running RAxML. A number
734 of options allow the user to speed up phyloscanner. Firstly it is 'embarrassingly' parallelisable, in
735 that each window of the genome can be processed separately (e.g. the 54 windows used for the
736 HIV data could have been processed via 54 jobs run in parallel). Secondly all options of RAxML
737 can be passed as options to phyloscanner, including multithreading. Thirdly the number of
738 unique sequences kept for phylogenetic inference can be controlled through various options,
739 notably merging of similar reads and/or a minimum read count. Fourthly the user can easily use
740 a different tool for phylogenetic inference instead of RAxML by using the --no-trees option of
741 phyloscanner_make_trees.py, and running the desired tool on the fasta file of processed reads that
742 is output for each window. (As an example running FastTree(Price et al. 2009) on the same data
743 took 28 seconds instead of the 10 minutes needed by RAxML.)
744

745 **Acknowledgments**

746
747 We thank Katrina Lythgoe for helpful discussions, and Céline Christiansen-Jucht for comments
748 on the manuscript. This work was funded by ERC Advanced Grant PBDR-339251. We
749 acknowledge funding from Bill & Melinda Gates Foundation through PANGEA-HIV. The STOP-
750 HCV Consortium is funded by a grant from the Medical Research Council (MR/K01532X/1). We
751 thank Gilead Sciences for providing HCV plasma samples from the BOSON clinical study for
752 use in these analyses. We also thank HCV Research UK (funded by the Medical Research
753 Foundation) for their assistance in handling and coordinating the release of samples for these
754 analyses. This work used the computing resources of the UK MEDical BIOinformatics
755 partnership - aggregation, integration, visualisation and analysis of large, complex data (UK
756 MED-BIO) which is supported by the Medical Research Council [grant number MR/L01632X/1].
757

758 **The BEEHIVE Collaboration**

759
760 Jan Albert, Margreet Bakker, Norbert Bannert, Ben Berkhout, Daniela Bezemer, François
761 Blanquart, Marion Cornelissen, Jacques Fellay, Katrien Fransen, Christophe Fraser, Astrid Gall,
762 Annabelle Gourlay, M. Kate Grabowski, Barbara Günsenheimer-Bartmeyer, Huldrych F.
763 Günthard, Matthew Hall, Mariska Hillebrecht, Paul Kellam, Pia Kivelä, Roger Kouyos, Oliver

764 Laeyendecker, Kirsi Liitsola, Laurence Meyer, Swee Hoe Ong, Kholoud Porter, Peter Reiss,
765 Matti Ristola, Ard van Sighem, and Chris Wymant.

766

767 Acknowledged contributors to the cohorts in the BEEHIVE Collaboration are listed in
768 supplementary section SI 4.

769

770 **The STOP-HCV Consortium**

771 Eleanor Barnes, Jonathan Ball, Diana Brainard, Gary Burgess, Graham Cooke, John Dillon,
772 Graham R Foster, Charles Gore, Neil Guha, Rachel Halford, Cham Herath, Chris Holmes, Anita
773 Howe, Emma Hudson, William Irving, Salim Khakoo, Paul Klenerman, Diana Koletzki, Natasha
774 Martin, Benedetta Massetto, Tamyo Mbisa, John McHutchison, Jane McKeating, John
775 McLauchlan, Alec Miners, Andrea Murray, Peter Shaw, Peter Simmonds, Chris C A Spencer,
776 Paul Targett-Adams, Emma Thomson, Peter Vickerman, and Nicole Zitzmann.

777

778 **The Maela Pneumococcal Collaboration**

779 Stephen D. Bentley, Claire Chewapreecha, Nicholas J. Croucher, Simon Harris, Jukka
780 Corander, David Goldblatt, Julian Parkhill, Francois Nosten, Claudia Turner, and Paul Turner.

781

782 **Competing Interests**

- 783
- 784 • AJG participated in an advisory board meeting for ViiV Healthcare in July 2016.
 - 785 • KP is a member of the Viiv 'Dolutegravir' Advisory Board and Viiv 'Data and Insights:
786 Standardisation in Measuring and Collecting Care Continuum Data' Advisory Board.
 - 787 • HG reports receipt of grants from the Swiss National Science Foundation, Swiss HIV
788 Cohort Study, University of Zurich, Yvonne Jacob Foundation, and Gilead Sciences; fees
789 for data and safety monitoring board membership from Merck; consulting/advisory
790 board membership fees from Gilead Sciences; and travel reimbursement from Gilead,
791 Bristol-Myers Squibb, and Janssen.
 - 792 • PR through his institution has received independent scientific grant support from Gilead
793 Sciences, Janssen Pharmaceuticals Inc, Merck & Co, Bristol-Myers Squibb, and ViiV
Healthcare; he has served on scientific advisory boards for Gilead Sciences and ViiV

794 Healthcare and on a data safety monitoring committee for Janssen Pharmaceuticals Inc,
795 for which his institution has received remuneration.

796

797

798 **References**

799

800 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool.
801 *Journal of Molecular Biology* 215:403–410.

802 Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence
803 data. *Bioinformatics* 30:2114–2120.

804 Bonsall D, Ansari MA, Ip C, Trebes A, Brown A, Klenerman P, Buck D, null N, Piazza P, Barnes
805 E, et al. 2015. ve-SEQ: Robust, unbiased enrichment for streamlined detection and whole-
806 genome sequencing of HCV and other highly diverse pathogens [version 1; referees: 2
807 approved, 1 approved with reservations]. *F1000Research* 4.

808 Cornelissen M, Gall A, Vink M, Zorgdrager F, Binter Š, Edwards S, Jurriaans S, Bakker M, Ong
809 SH, Gras L, et al. 2016. From clinical sample to complete genome: Comparing methods for
810 the extraction of HIV-1 RNA for high-throughput deep sequencing. *Virus Research*.

811 Cornelissen M, Pasternak AO, Grijsen ML, Zorgdrager F, Bakker M, Blom P, Prins JM,
812 Jurriaans S, van der Kuyl AC. 2012. HIV-1 Dual Infection Is Associated With Faster CD4+
813 T-Cell Decline in a Cohort of Men With Primary HIV Infection. *Clinical Infectious Diseases*
814 54:539.

815 Croucher NJ, Mostowy R, Wymant C, Turner P, Bentley SD, Fraser C. 2016. Horizontal DNA
816 Transfer Mechanisms of Bacteria as Weapons of Intragenomic Conflict. *PLOS Biology*
817 14:1–42.

818 Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, Parkhill J, Harris SR.
819 2015. Rapid phylogenetic analysis of large samples of recombinant bacterial whole genome
820 sequences using Gubbins. *Nucleic Acids Research* 43:e15.

821 Didelot X, Fraser C, Gardy J, Colijn C. 2017. Genomic Infectious Disease Epidemiology in
822 Partially Sampled and Ongoing Outbreaks. *Molecular Biology and Evolution* 34:997.

823 Foster GR, Pianko S, Brown A, Forton D, Nahass RG, George J, Barnes E, Brainard DM,
824 Massetto B, Lin M, et al. 2015. Efficacy of Sofosbuvir Plus Ribavirin With or Without
825 Peginterferon-Alfa in Patients With Hepatitis C Virus Genotype 3 Infection and Treatment-
826 Experienced Patients With Cirrhosis and Hepatitis C Virus Genotype 2 Infection.
827 *Gastroenterology* 149:1462–1470.

828 Fraser C, Lythgoe K, Leventhal GE, Shirreff G, Hollingsworth TD, Alizon S, Bonhoeffer S. 2014.
829 Virulence and Pathogenesis of HIV-1 Infection: An Evolutionary Perspective. *Science* 343.

830 Gall A, Ferns B, Morris C, Watson S, Cotten M, Robinson M, Berry N, Pillay D, Kellam P. 2012.

- 831 Universal Amplification, Next-Generation Sequencing, and Assembly of HIV-1 Genomes.
832 *Journal of Clinical Microbiology* 50:3838–3844.
- 833 Gall A, Morris C, Kellam P, Berry N. 2014. Complete Genome Sequence of the WHO
834 International Standard for HIV-1 RNA Determined by Deep Sequencing. *Genome*
835 *Announcements* 2.
- 836 Gatanaga H, Suzuki Y, Tsang H, Yoshimura K, Kavlick MF, Nagashima K, Gorelick RJ, Mardy
837 S, Tang C, Summers MF, et al. 2002. Amino Acid Substitutions in Gag Protein at Non-
838 cleavage Sites Are Indispensable for the Development of a High Multitude of HIV-1
839 Resistance against Protease Inhibitors. *Journal of Biological Chemistry* 277:5952–5961.
- 840 Goodwin S, McPherson JD, McCombie WR. 2016. Coming of age: ten years of next-generation
841 sequencing technologies. *Nat Rev Genet* 17:333–351.
- 842 Gross KL, Porco TC, Grant RM. 2004. HIV-1 superinfection and viral diversity. *AIDS* 18.
- 843 Hall M, Woolhouse M, Rambaut A. 2015. Epidemic Reconstruction in a Phylogenetics
844 Framework: Transmission Trees as Partitions of the Node Set. Salathé M, editor. *PLoS*
845 *Comput Biol* 11:e1004613.
- 846 Hunt M, Gall A, Ong SH, Brener J, Ferns B, Goulder P, Nastouli E, Keane JA, Kellam P, Otto
847 TD. 2015. IVA: accurate de novo assembly of RNA virus genomes. *Bioinformatics*.
- 848 Jabara CB, Jones CD, Roach J, Anderson JA, Swanstrom R. 2011. Accurate sampling and
849 deep sequencing of the HIV-1 protease gene using a Primer ID. *Proceedings of the*
850 *National Academy of Sciences* 108:20166–20171.
- 851 Johnson VA, Calvez V, Günthard HF, Paredes R, Pillay D, Shafer R, Wensing AM, Richman
852 DD. 2011. 2011 update of the drug resistance mutations in HIV-1. *Top Antivir Med* 19:156–
853 164.
- 854 Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N. 2014. Bayesian
855 Reconstruction of Disease Outbreaks by Combining Epidemiologic and Genomic
856 Data. Tanaka MM, editor. *PLoS Comput Biol* 10:e1003457.
- 857 Katoh K, Misawa K, Kuma KI, Miyata T. 2002. MAFFT: a novel method for rapid multiple
858 sequence alignment based on fast Fourier transform. *Nucleic Acids Research* 30:3059–
859 3066.
- 860 Kuiken C, Foley B, Leitner T, Apetrei C, Hahn B, Mizrahi I, Mullins J, Rambaut A, Wolinsky S,
861 Korber B. 2012. HIV Sequence Compendium 2012. Los Alamos National Laboratory,
862 Theoretical Biology and Biophysics, Los Alamos, New Mexico. LA-UR-12-24653.
- 863 Kuiken C, Yusim K, Boykin L, Richardson R. 2005. The Los Alamos hepatitis C sequence
864 database. *Bioinformatics* 21:379–384.
- 865 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows–Wheeler transform.
866 *Bioinformatics* 25:1754.
- 867 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,

- 868 Subgroup 1GDP. 2009. The Sequence Alignment/Map (SAM) Format and SAMtools.
869 Bioinformatics.
- 870 Numminen E, Chewapreecha C, Siren J, Turner C, Turner P, Bentley SD, Corander J. 2014.
871 Two-phase importance sampling for inference about transmission trees. *Proceedings of the*
872 *Royal Society B: Biological Sciences* 281:20141324–20141324.
- 873 Prabhakaran S, Rey M, Zagordi O, Beerenwinkel N, Roth V. 2014. HIV Haplotype Inference
874 Using a Propagating Dirichlet Process Mixture Model. *IEEE/ACM Transactions on*
875 *Computational Biology and Bioinformatics* 11:182–191.
- 876 Price MN, Dehal PS, Arkin AP. 2009. FastTree: Computing Large Minimum Evolution Trees with
877 Profiles instead of a Distance Matrix. *Molecular Biology and Evolution* 26:1641–1650.
- 878 Quail MA, Kozarewa I, Smith F, Scally A, Stephens PJ, Durbin R, Swerdlow H, Turner DJ. 2008.
879 A large genome center's improvements to the Illumina sequencing system. *Nat Meth*
880 5:1005–1010.
- 881 Quail MA, Otto TD, Gu Y, Harris SR, Skelly TF, McQuillan JA, Swerdlow HP, Oyola SO. Optimal
882 enzymes for amplifying sequencing libraries. *Nat Meth* 9:10–11.
- 883 Rambaut A, Grassly NC. 1997. Seq-Gen: an application for the Monte Carlo simulation of DNA
884 sequence evolution along phylogenetic trees. *Bioinformatics* 13:235.
- 885 Rambaut A, Lam TT, Max Carvalho L, Pybus OG. 2016. Exploring the temporal structure of
886 heterochronous sequences using TempEst (formerly Path-O-Gen). *Virus Evolution*
887 2:vew007.
- 888 Romero-Severson EO, Bulla I, Leitner T. 2016. Phylogenetically resolving epidemiologic
889 linkage. *Proceedings of the National Academy of Sciences* 113:2690–2695.
- 890 Ronquist F, Teslenko M, van der Mark P, Ayres DL, Darling A, Höhna S, Larget B, Liu L,
891 Suchard MA, Huelsenbeck JP. 2012. MrBayes 3.2: Efficient Bayesian Phylogenetic
892 Inference and Model Choice Across a Large Model Space. *Systematic Biology* 61:539.
- 893 Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B,
894 Ideker T. 2003. Cytoscape: A Software Environment for Integrated Models of Biomolecular
895 Interaction Networks. *Genome Research* 13:2498–2504.
- 896 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
897 large phylogenies. *Bioinformatics* 30:1312.
- 898 Töpfer A, Marschall T, Bull RA, Luciani F, Schönhuth A, Beerenwinkel N. 2014. Viral
899 Quasispecies Assembly via Maximal Clique Enumeration. *PLoS Comput Biol* 10:1–10.
- 900 Volz EM, Frost SDW. 2013. Inferring the Source of Transmission with Phylogenetic
901 Data. Kosakovsky Pond SL, editor. *PLoS Comput Biol* 9:e1003397.
- 902 Wensing AM, Calvez V, Günthard HF, Johnson VA, Paredes R, Pillay D, Shafer RW, Richman
903 DD. 2015. 2015 Update of the Drug Resistance Mutations in HIV-1. *Top Antivir Med*
904 23:132–141.

- 905 Worby CJ, O'Neill PD, Kypraios T, Robotham JV, De Angelis D, Cartwright EJP, Peacock SJ,
906 Cooper BS. 2016. Reconstructing transmission trees for communicable diseases using
907 densely sampled genetic data. *Ann. Appl. Stat.*:395–417.
- 908 Wymant C, Blanquart F, Gall A, Bakker M, Bezemer D, Croucher NJ, Golubchik T, Hall M,
909 Hillebregt M, Ong SH, et al. 2016. Easy and Accurate Reconstruction of Whole HIV
910 Genomes from Short-Read Sequence Data.
- 911 Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. 2017. ggtree: an r package for visualization and
912 annotation of phylogenetic trees with their covariates and other associated data. *Methods in*
913 *Ecology and Evolution* 8:28–36.
- 914 Zagordi O, Bhattacharya A, Eriksson N, Beerenwinkel N. 2011. ShoRAH: estimating the genetic
915 diversity of a mixed sample from next-generation sequencing data. *BMC Bioinformatics*
916 12:119.
- 917 Zanini F, Brodin J, Thebo L, Lanz C, Bratt G, Albert J, Neher RA. 2015. Population genomics of
918 inpatient HIV-1 evolution. Chakraborty AK, editor. *eLife* 4:e11282.
- 919