

1 Within-host diversity improves 2 phylogenetic and transmission 3 reconstruction of SARS-CoV-2 4 outbreaks

5 **Arturo Torres Ortiz^{1*}, Michelle Kendall², Nathaniel Storey³, James Hatcher³,**
6 **Helen Dunn³, Sunando Roy⁴, Rachel Williams⁵, Charlotte Williams⁵, Richard A.**
7 **Goldstein⁵, Xavier Didelot², Kathryn Harris^{3,6}, Judith Breuer⁴, Louis Grandjean^{4*}**

*For correspondence:

a.ortiz@ucl.ac.uk (ATO);
l.grandjean@ucl.ac.uk (LG)

8 ¹Department of Infectious Diseases, Imperial College London, London, W2 1NY;
9 ²Department of Statistics, University of Warwick, Coventry, CV4 7AL; ³Department of
10 Microbiology, Great Ormond Street Hospital, London WC1N 3JH; ⁴Department of
11 Infection, Immunity and Inflammation, Institute of Child Health, UCL, London WC1N
12 1EH; ⁵UCL Genomics, Institute of Child Health, UCL, London WC1N 1EH; ⁶Department of
13 Virology, East & South East London Pathology Partnership, Royal London Hospital, Barts
14 Health NHS Trust, London E12ES

16 **Abstract** Accurate inference of who infected whom in an infectious disease outbreak is critical
17 for the delivery of effective infection prevention and control. The increased resolution of
18 pathogen whole-genome sequencing has significantly improved our ability to infer transmission
19 events. Despite this, transmission inference often remains limited by the lack of genomic
20 variation between the source case and infected contacts. Although within-host genetic diversity is
21 common among a wide variety of pathogens, conventional whole-genome sequencing
22 phylogenetic approaches exclusively use consensus sequences, which consider only the most
23 prevalent nucleotide at each position and therefore fail to capture low frequency variation within
24 samples. We hypothesized that including within-sample variation in a phylogenetic model would
25 help to identify who infected whom in instances in which this was previously impossible. Using
26 whole-genome sequences from SARS-CoV-2 multi-institutional outbreaks as an example, we
27 show how within-sample diversity is partially maintained among repeated serial samples from
28 the same host, it can be transmitted between those cases with known epidemiological links, and
29 how this improves phylogenetic inference and our understanding of who infected whom. Our
30 technique is applicable to other infectious diseases and has immediate clinical utility in infection
31 prevention and control.

33 Introduction

34 Understanding who infects whom in an infectious disease outbreak is a key component of infection
35 prevention and control (*Didelot et al., 2012*). The use of whole-genome sequencing allows for
36 detailed investigation of disease outbreaks, but the limited genetic diversity of many pathogens
37 often hinders our understanding of transmission events (*Campbell et al., 2018*). As a consequence
38 of the limited diversity, many index case and contact pairs will share identical genotypes, making

39 it difficult to ascertain who infected whom.

40 Within-sample genetic diversity is common among a wide variety of pathogens (*Mongkolrattan-*
41 *othai et al., 2011; Lieberman et al., 2016; Dinis et al., 2016; Leitner, 2019; Popa et al., 2020*). This
42 diversity may be generated *de novo* during infection, by a single transmission event of a diverse
43 inoculum or by independent transmission events from multiple sources (*Worby et al., 2014*). The
44 maintenance and dynamic of within-host diversity is then a product of natural selection, genetic
45 drift, and fluctuating population size (*Didelot et al., 2012*). The transmission of within-host varia-
46 tion between individuals is also favored as a large inoculum exposure is more likely to give rise to
47 infection (*Murphy et al., 1984; Han et al., 2019; Lee et al., 2022; Sender et al., 2021; Spinelli et al.,*
48 *2021; Trunfio et al., 2021*). The amount of within-sample diversity transmitted from index-case to
49 contact is determined by the bottleneck size (*Zwart and Elena, 2015*), with stringent bottlenecks
50 limiting the number of genotypes transmitted from the host to the recipient, and wide bottlenecks
51 allowing for the transmission of higher levels of genetic diversity (*Worby et al., 2014*).

52 Phylogenetic analysis provide information regarding the structure of the genetic diversity among
53 pathogen isolates. Moreover, pathogen phylogenetic trees can be used as input for many down-
54 stream analysis, including inference of transmission events, population size dynamics or estima-
55 tion of parameters of epidemiological models (*Didelot et al., 2018*). Most genomic and phyloge-
56 netic workflows involve either genome assembly or alignment of sequencing reads to a reference
57 genome. In both cases, conventionally the resulting alignment exclusively represents the most
58 common nucleotide at each position. This is often referred to as the consensus sequence. Al-
59 though genome assemblers may output contigs (combined overlapping reads) representing low
60 frequency haplotypes, only the majority contig is kept in the final sequence. In a mapping approach,
61 a frequency threshold for the major variant is usually pre-determined, under which a position is
62 considered ambiguous. The lack of genetic variation between temporally proximate samples and
63 the slow mutation rate of many pathogens results in direct transmission events sharing exact se-
64 quences between the hosts when using the consensus sequence approach. For instance, the sub-
65 stitution rate of SARS-CoV-2 has been inferred to be around 2 mutations per genome per month
66 (*Harvey et al., 2021*). Given its infectious period of 6 days (*Byrne et al., 2020*), most consensus
67 sequences in a small-scale outbreak will show no variation between them. This lack of resolution
68 and poor phylogenetic signal complicate phylogenetic inference, limiting the downstream analysis
69 and conclusions that can be extracted from the phylogenetic tree. Previous work has shown the
70 advantages of using within-host diversity to infer transmission events compared to using consen-
71 sus sequences (*Wymant et al., 2018; De Maio et al., 2018*). Aside from transmission inference, the
72 use of the within-host pathogen genetic data directly within phylogenetic inference will improve
73 any downstream analysis using a phylogenetic tree as starting point.

74 We hypothesize that the failure of consensus sequence approaches to capture within-sample
75 variation arbitrarily excludes meaningful data and limits pathogen phylogenetic and transmission
76 inference, and that including within-sample diversity in a phylogenetic model would significantly
77 increase the evolutionary and temporal signal and thereby improve our ability to infer infectious
78 disease phylogenies and transmission events.

79 We tested our hypothesis on multi-institutional SARS-CoV-2 outbreaks across London hospitals
80 that were part of the COVID-19 Genomics UK (COG-UK) consortia (*COVID-19 Genomics UK (COG-UK),*
81 *2020*). Technical replicates, repeated longitudinal sampling from the same patient, and epidemio-
82 logical data allowed us to evaluate the presence and stability of within-sample diversity within the
83 host and in independently determined transmission chains. We also evaluated the use of within-
84 sample diversity in phylogenetic analysis by conducting outbreak and phylogenetic simulations
85 of sequencing data using a phylogenetic model that accounts for the presence and transmission
86 of within-sample variation. We show the effects on phylogenetic inference of using consensus se-
87 quences in the presence of within-sample diversity, and propose that existing phylogenetic models
88 can leverage the additional diversity given by the within-sample variation and reconstruct the phy-
89 logenetic relationship between isolates. Lastly, we show that by taking into account within-sample

90 diversity in a phylogenetic model we improve the temporal signal in SARS-CoV-2 outbreak analysis.
91 Using both phylogenetic outbreak reconstruction and simulation we show that our approach is
92 superior to the current gold standard whole-genome consensus sequence methods.

93 Results

94 Sampling, demographics and metadata

95 Between March 2020 and November 2020, 451 healthcare workers, patients and patient contacts
96 at the participating North London Hospitals were diagnosed at the Camelia Botnar Laboratories
97 with SARS-CoV-2 by PCR as part of a routine staff diagnostic service at Great Ormond Street Hospi-
98 tal NHS Foundation Trust (GOSH). A total of 289 isolates were whole-genome sequenced using the
99 Illumina NextSeq platform, which resulted in 522 whole-genome sequences including longitudinal
100 and technical replicates (Supplementary file 1). The mean participant age was 40 years old (median
101 38.5 years old, interquartile range (IQR) 30-50 years old), and 60% of the participants were female
102 (Supplementary file 2). All samples were SARS-CoV-2 positive with real time qPCR cycle threshold
103 (C_t) values ranging from 16 to 35 cycles (Supplementary file 2). The earliest sample was collected
104 on 26th March 2020, while the latest one dated to 4th November 2020 (**Figure 1—figure Supple-**
105 **ment 1a**). A total of 291 samples had self-reported symptom onset data, for which the mean time
106 from symptom onset to sample collection date was 5 days (IQR 2-7 days, **Figure 1—figure Supple-**
107 **ment 1b**). More than 90% of the samples were taken from hospital staff, while the rest comprised
108 patients and contacts of either the patients or the staff members (Supplementary file 2).

109 Genomic analysis of SARS-CoV-2 sequences

110 A total of 454 whole-genomes with mean coverage higher than 10x were kept for further analysis,
111 resulting in an average coverage across isolates of 2457x (**Figure 1—figure Supplement 2**). Allele
112 frequencies were extracted using the pileup functionality within *bcftools* ([Danecek et al., 2011](#)) with
113 a minimum base and mapping quality of 30, which represents a base call error rate of 0.1%. Vari-
114 ants at low frequency at positions where the mapped reads support more than one allele were
115 coined as minor or low-frequency variants. Variants were filtered further for read position bias
116 and strand bias. Only minor variants with an allele frequency of at least 2% were kept as puta-
117 tive variants. Samples with a frequency of missing bases higher than 10% were excluded, keeping
118 350 isolates for analysis. The mean number of low frequency variants was 12 (median 3, IQR 1.00 –
119 9.75), although both the number of variants and its deviation increased at high C_t values (**Figure 1—**
120 **figure Supplement 3**).

121 Within-sample variation in technical replicates

122 To understand the stability of within-sample variation and minimize spurious variant calls, we se-
123 quenced and analyzed technical replicates of 17 samples. Overall, when the variant was present in
124 both duplicates the correlation of the variant frequencies was high ($R^2 = 0.9$, **Figure 1a** right). The
125 high correlation was also maintained at low variant frequencies (**Figure 1a** left).

126 Minor variants were less likely to be shared when one or more of the paired samples had low
127 viral load. These discrepancies may appear because of amplification bias caused by low genetic
128 material, base calling errors due to low coverage, or low base quality. The mean proportion of
129 discrepant within-sample variants between duplicated samples was 0.39 (sd = 0.29), although this
130 varied between duplicates (**Figure 1—figure Supplement 4**). C_t values in RT-PCR obtained during
131 viral amplification are inversely correlated with viral load ([Tom and Mina, 2020](#)). The proportion of
132 shared intra-host variants was negatively correlated with C_t values in a logistic model (estimate=
133 0.78, p-value=0.008), with higher C_t values associated with a lower amount of shared intra-host
134 variants (**Figure 1c**). The number of within-sample variants detected also increased with C_t value,
135 as well as the deviation in the number of variants between duplicates (**Figure 1d**). This could be
136 explained either by an increase in the number of spurious variants at low viral loads (**Tonkin-Hill**

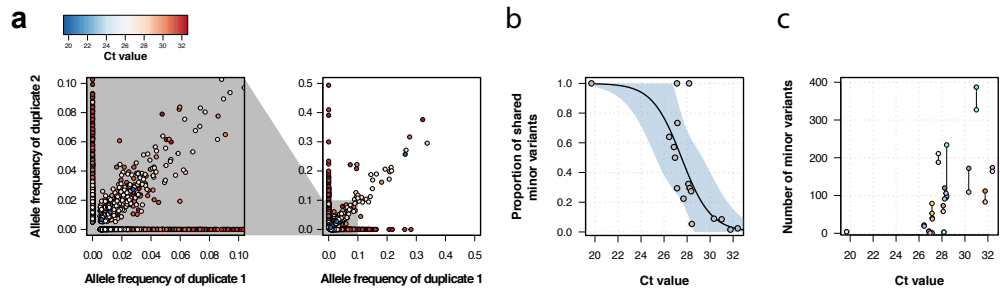


Figure 1. Genomic analysis of technical duplicates before filtering. **a** Allele frequency comparison between technical replicates for all frequencies (right) and for frequencies up to 1% (left). Colors represent the C_t value for the sample. **b** Proportion of shared minor variants between technical replicates in relation to the C_t value. **c** Total number of minor variants in relation to the C_t value. Lines linked two technical replicates. Each sequence has a different color, with sequences from the same patient having a different shade of the same color.

Figure 1—figure supplement 1. Collection date distribution and time from symptom and days from symptom onset

Figure 1—figure supplement 2. Sample mean coverage distribution

Figure 1—figure supplement 3. Effects of C_t value on whole-genome sequencing data

Figure 1—figure supplement 4. Proportion of shared minor variants between technical replicates using different filters of allele frequency

137 *et al., 2021*), biased amplification of low level sub-populations minor rare alleles (*McCrone and*
 138 *Lauring, 2016*), or due to the accumulation of within-host variation through time, as late stages of
 139 infection are usually characterized by high C_t values (low viral load).

140 Based on these results, only samples with a C_t value equal or lower than 30 cycles were consid-
 141 ered, which resulted in 249 samples kept for analysis. Additionally, only variants with a frequency
 142 higher or equal than 2% were used. For the filtered dataset, 414 out of 29903 positions were poly-
 143 morphic for the consensus sequence, while the alignment with within-sample diversity had 1039
 144 SNPs. Of these, 699 positions had intra-host diversity, of which 78% (549/699) were singletons. The
 145 majority of samples (207/249, 83%) contained at least 1 position with a high quality within-host vari-
 146 ant, and the median amount of intra-host variants per sample was 2 (IQR 1-4.5).

147 **Within-sample variation in epidemiologically linked samples**

148 Given the limited genomic information in the consensus sequences, epidemiological data is often
 149 necessary to infer the directionality of transmission. We categorized our samples within the fol-
 150 lowing groups: samples that a) did not have any recorded epidemiological link, b) were from the
 151 same hospital (possibly linked), c) samples that were part of the same department within the same
 152 hospital (probable link), d) samples that had an epidemiological link within the same department
 153 of the same hospital (proven link), e) were a longitudinal replicate from the same patient and f) a
 154 technical replicate from the same sample.

155 We tested the concordance between epidemiological and genomic data by determining the SNP
 156 distance between pairs of samples with epidemiological links and without them. Pairs of samples
 157 from the same hospital, department, epidemiologically linked, or longitudinal and technical repli-
 158 cates had a lower SNP distance (were more closely related) than those samples that did not have
 159 any relationship, although this difference was small in the case of pairs of samples from the same
 160 hospital (*Table 1*).

161 To understand the distribution of shared low frequency variants among different groups of
 162 samples, we performed a pairwise comparison of all samples and calculated the proportion of

Table 1. SNP distance between pairs of samples.

Sample relationship	Estimate (95%CI)	p-value
None	11.04 (10.94 - 11.15)	Reference
Hospital	9.78 (9.48 - 10.09)	$<1 \times 10^{-4}$
Department	5.15 (4.54 - 5.83)	$<1 \times 10^{-4}$
Epidemiological	1.5 (1.22 - 1.78)	$<1 \times 10^{-4}$
Longitudinal duplicates	0 (0 - 0.2)	$<1 \times 10^{-4}$
Technical replicate	0 (0 - 0.2)	$<1 \times 10^{-4}$

163 shared within-sample variants (shared variants divided by total variants in the pair) within groups
 164 with epidemiological links and without them. The proportion of shared within-host variants was
 165 higher between technical replicates, longitudinal duplicates, epidemiologically linked samples, and
 166 samples taken from individuals from the same department when compared to pairs with no epi-
 167 demiological links, although the range of this probability was large (**Figure 2, Figure 2—figure Sup-
 168 plement 1**). The probability of sharing a low frequency variant was inferred using a logistic regres-
 169 sion model (**Figure 2—figure Supplement 2**). There was a tendency for the probability to increase
 170 with variant frequency, but the association was not strong (Odds ratio 1.8, 95% CI 0.9 – 3.5, $p=0.08$).
 171 The probability of sharing a low frequency variant for samples with no epidemiological links was
 172 9.5×10^{-6} (95% CI $8.8 \times 10^{-6} - 1.02 \times 10^{-5}$). Samples from the same hospital did not have a probability
 173 significantly higher than those without any link (3.3×10^{-3} , 95% CI $2.7 \times 10^{-3} - 4.03 \times 10^{-3}$). On the
 174 other hand, pairs from the same department, with epidemiological links, replicates or technical
 175 replicates all had a significantly higher probability of sharing a low frequency variant when com-
 176 pared to those pairs with no link (all Wald test p-values < 0.001). The inferred probabilities for pairs
 177 from the sample department was 1.4% (95% CI 0.9% – 2.1%), which increased to 5% for pairs with
 178 epidemiological links (95% CI 4.2% – 6.4%). For longitudinal replicates, the probability was inferred
 179 to be 38% (95% CI 35% – 41%), and were shared between multiple time points (**Figure 2—figure
 180 Supplement 3**). Technical replicates were estimated to have the highest probability (70%, 95% CI
 181 64% – 76%).

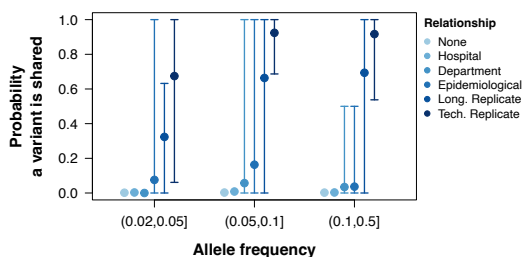


Figure 2. Probability of sharing within-host variants in sample pairs. The probability of variants shared between pairs of samples calculated as the number of low frequency variants in both samples divided by the total number of variants between the pair. Colors grouped samples by their relationship. Points represent the mean probability a variant is shared between all pairwise samples within a group and allele frequency. Error bars show the 95th and 5th percentiles.

Figure 2—figure supplement 1. Allele frequency comparison in pairwise sample pairs.

Figure 2—figure supplement 2. Probability that minor variants are shared.

Figure 2—figure supplement 3. Dynamics of low frequency variants in longitudinal duplicates.

182 **Within-host diversity model outperforms the consensus model in simulations**

183 The effect of within-sample diversity in phylogenetic inference was tested by evaluating the accu-
184 racy in the reconstruction of known phylogenetic trees using a conventional phylogenetic model
185 and a model that accounts for within-sample variation.

186 The presence of within-sample diversity was coded in the genome alignment using existing
187 IUPAC nomenclature (*IUPAC-IUB Joint Commission on Biochemical Nomenclature (JCBN), 1984*).
188 For the consensus sequence alignment, only the 4 canonical nucleotides were used (*Figure 3a,b*),
189 while the proposed alignment retained the major and minor allele information as independent
190 character states (*Figure 3c,d*).

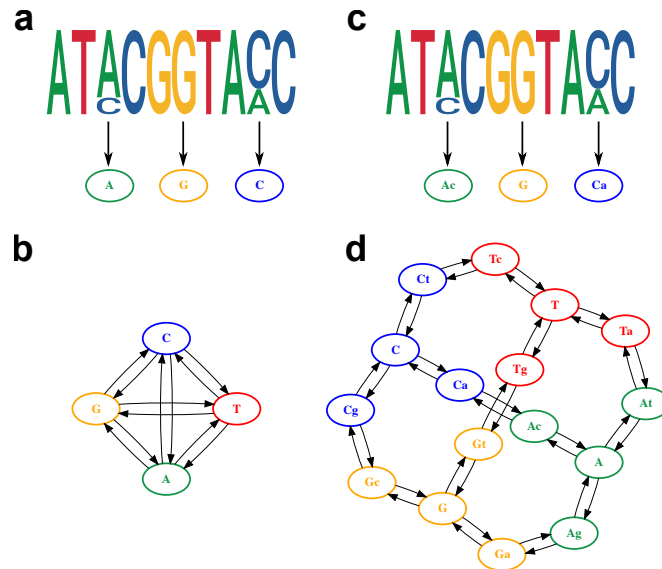


Figure 3. Model of within-host diversity.

Proposed evolutionary model of within-host diversity in genomic sequences. Uppercase letters represent the major variant in the population, while lowercase letters indicate presence of a minor variant alongside the major one. **a, c** Genome sequences where some positions show within-sample variation (top), represented by a major allele (big size letter) and a minor one (smaller size), as well as its representation in the alignment (bottom). **b, d** Models of nucleotide evolution. Character transitions are indicated by arrows. **a** Consensus sequence, where only the major allele is represented in the alignment. **b** Model of nucleotide evolution using the consensus sequence, with four character states representing the four nucleotides. **c** Sequence with within-sample variation, represented by an uppercase letter for the major allele and a lower case letter for the minor allele. **d** Model of nucleotide evolution with 16 character states accounting for within-sample variation.

191 In order to evaluate the differences in tree inference with and without the inclusion of within-
192 sample diversity, we simulated genome alignments for 100 random trees using a phylogenetic
193 model where both major and minor variant combinations were considered, resulting in a total of
194 16 possible states (*Figure 3d*). In the proposed model, transitions and transversions between the
195 four nucleotides in the population occur in the following steps: first a minority variant evolves at
196 low frequency, then the minor variant increases its frequency to become the majority nucleotide,
197 and finally the variant is fixed (*Figure 3d*), with all the steps being reversible. Therefore, within-
198 host evolutionary dynamics are modelled by explicitly considering base change as a process of
199 minor variant evolution and eventual fixation. The substitution rates chosen for the simulations,
200 as shown in Supplementary file 4, were selected to reflect a slow rate of minor variant evolution
201 and a fast rate at which minor variants are lost or fixated in the population, which in turns results
202 in a highly dynamic landscape of within-sample variation, with the four canonical nucleotides 100

203 times more likely to be present than low frequency variants.

204 From the simulated genomes, two types of alignments were generated: a consensus sequence,
205 where only the major allele was considered (**Figure 3a**); and an alignment that retained the ma-
206 jor and minor allele information as independent character states (**Figure 3c**). From the simulated
207 alignments, RaxML-NG was used to infer phylogenetic trees (*Kozlov et al., 2019*). The consensus
208 sequence was analyzed with a GTR+ γ model, while the PROTGTR+ γ model was used in order to ac-
209 commodate the extra characters of the alignment with within-sample diversity and major/minor
210 variant information.

211 The two models were evaluated for their ability to infer the known phylogeny that included
212 within-host diversity. The estimated phylogenies were compared to the known tree using differ-
213 ent measures to capture dissimilarities in a variety of aspects relevant to tree inference (Supple-
214 mentary file 3). For all the metrics employed, the phylogenies inferred explicitly using within-host
215 diversity as independent characters approximated better to the initial tree than the one using the
216 consensus sequence (**Figure 4**). Additionally, the transition/transversion rates inferred by the phy-
217 logenetic models accounting for within-host diversity accurately reflect the rates used for the sim-
218 ulation of genomic sequences (Supplementary file 4, 5 and 6).

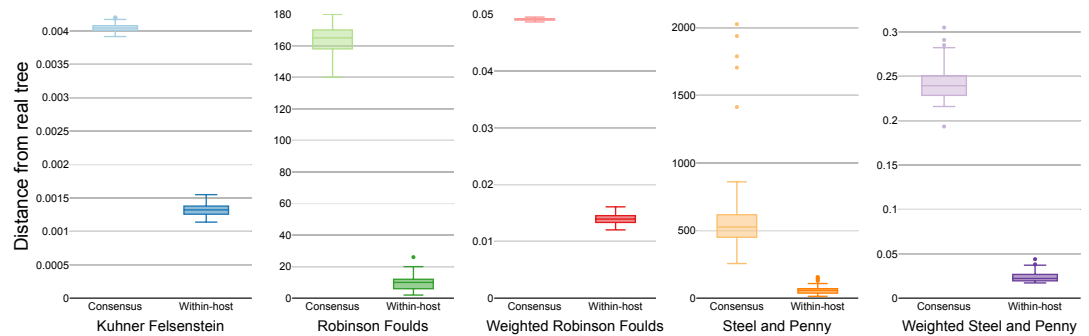


Figure 4. Similarity scores for inferred trees.

Comparison of the phylogenetic trees inferred using simulated sequences from known random starting trees and different phylogenetic models. Colors differentiate the metrics used for the comparison.

Figure 4—figure supplement 1. Similarity scores for inferred trees with different rates.

Figure 4—figure supplement 2. Similarity scores for inferred trees from coalescent simulations.

219 As different pathogens are likely to show different dynamics of within-host variation and the
220 rates used for the simulations will inevitably affect the improvement of using the 16-state model,
221 we simulated genomes with different parameters. As expected, choosing rates that promote an
222 abundant and stable landscape of low frequency variation (rate of minor variant acquisition of
223 20, and rates of variant switch and lost of 1) made the 16-state model to perform better than the
224 model using consensus sequences, which improved as the proportion of low frequency variants
225 decreased (**Figure 4—figure Supplement 1**). Conversely, in simulations using a Jukes-Cantor DNA
226 model, and therefore without any low frequency variation, both models showed similar results
227 (**Figure 4—figure Supplement 1**).

228 To understand the effects of genetic linkage between sites in the phylogenetic model due to the
229 clonal relationships between genomes, we evaluated another set of simulations where the starting
230 tree was generated using the coalescent model, which increases the correlation between sites.
231 For all metrics used, the model using low frequency variants inferred phylogenies more similar to
232 the starting coalescent tree than those inferred using the consensus sequence (**Figure 4—figure
233 Supplement 2**).

234 We further assessed the effect of within-host diversity in phylogenetic inference by simulating
235 pathogen evolution throughout the time frame of infectious disease outbreaks (*De Maio et al.,
236 2018*). We simulated outbreaks using TransPhylo (*Didelot et al., 2017*) with a host population vary-

237 ing between 10 and 15 hosts, no recombination, complete sampling of the outbreak and selecting
 238 epidemiological parameters to match the transmission dynamics of SARS-CoV-2. For each out-
 239 break, we simulated the evolution and transmission of the pathogen population within each host
 240 with varying values of mutation rates and transmission bottlenecks using fastsimcoal2 (*Excoffier*
 241 *et al., 2013*) as previously described by *De Maio et al. (2018)*. We compared the resulting phylo-
 242 genetic trees to the real outbreak phylogeny using the Kuhner-Felsenstein distance (*Kuhner and*
 243 *Felsenstein, 1994*). Even though using consensus sequences performed better than a random
 244 distribution of trees, using within-host diversity outperformed the consensus sequence in all in-
 245 stances (*Figure 5*). The phylogenies inferred using within-host diversity were more similar to the
 246 real outbreak phylogeny for wider bottleneck sizes, with the best performance when no bottleneck
 247 was present. As expected, both the consensus sequences and the sequences reflecting within-
 248 sample diversity were more informative at higher mutation rates, even though the consensus se-
 249 quence only showed improvement with a mutation rate of 10^{-3} mutations per base per generation
 250 cycle (*Figure 5*).

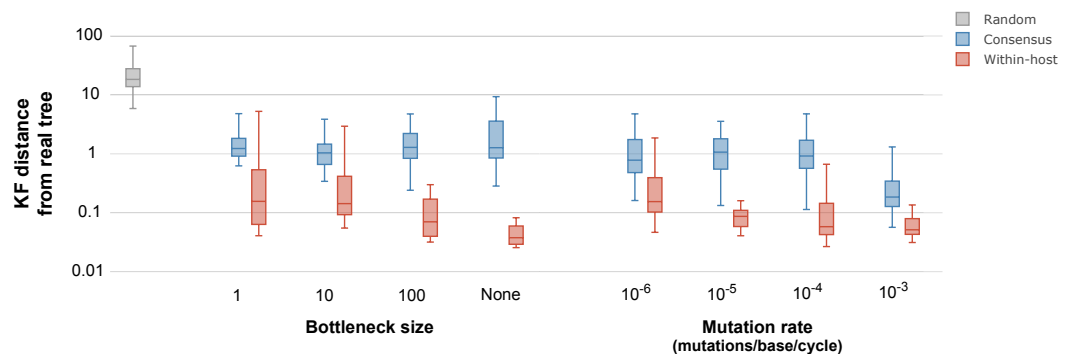


Figure 5. Inferred phylogenetic trees from outbreak simulations. Kuhner and Felsenstein (KF) tree distance between phylogenetic trees from simulated outbreaks. Phylogenies were inferred using consensus sequences (blue) and alignments reflecting within-sample diversity (red), and compared to the known phylogeny of the simulated outbreak. For reference, grey color represents a set of random trees. Outbreak simulations were performed with different bottleneck sizes and mutation rates. The mutation rate is measured as the number of mutations per base per generation cycle.

251 Within-host diversity improves the resolution in SARS-CoV-2 phylogenetics

252 Genome sequences collected at different time points are expected to diverge as time progresses,
 253 resulting in a positive correlation between the isolation date and the number of accumulated mu-
 254 tations (temporal signal) (*Rieux and Balloux, 2016*). The alignment with consensus sequences and
 255 the one reflecting within-sample variation were used to infer two different phylogenetic trees (*Fig-*
 256 *ure 6—figure Supplement 1*). Longitudinal samples in the phylogeny inferred using within-host
 257 diversity reflected the expected temporal signal, with an increase in genetic distance as time pro-
 258 gressed between the longitudinal pairs in a linear model (coefficient 2.24, 0.59 - 3.88 95% CI, $p =$
 259 0.019 , *Figure 6—figure Supplement 2*). The difference in C_i value among longitudinal duplicates
 260 was not correlated with higher genetic distances (coefficient 1.62, -0.66 - 3.91 95% CI, $p = 0.2$).
 261 Similarly, we analyzed the number of low frequency variants within outbreaks by counting the
 262 number of within-sample variants for each isolate belonging to a specific outbreak and inferred
 263 their change through time taking the earliest isolate date as the starting point of the outbreak. In
 264 general, as the outbreaks progressed the number of low frequency variants increased (coefficient
 265 0.16 , $0.06 - 0.27$ 95% CI, $p = 0.003$, $r^2 = 0.19$, *Figure 6—figure Supplement 3*).

266 We analyzed the impact of using within-sample variation on the temporal structure of the phy-
 267 logeny by systematically identifying clusters of tips in the phylogenetic tree with an identical con-
 268 sensus sequence and no temporal signal. We then performed a root-to-tip analysis using the tree

269 inferred with intra-sample diversity. Only clusters with more than 3 tips were used for the root-to-
 270 tip analysis. The majority of clusters (10/11) showed a positive correlation between the distance of
 271 the tips to the root and the collection dates, demonstrating a significant temporal signal between
 272 samples when there was none using the conventional consensus tree (*Figure 6*).

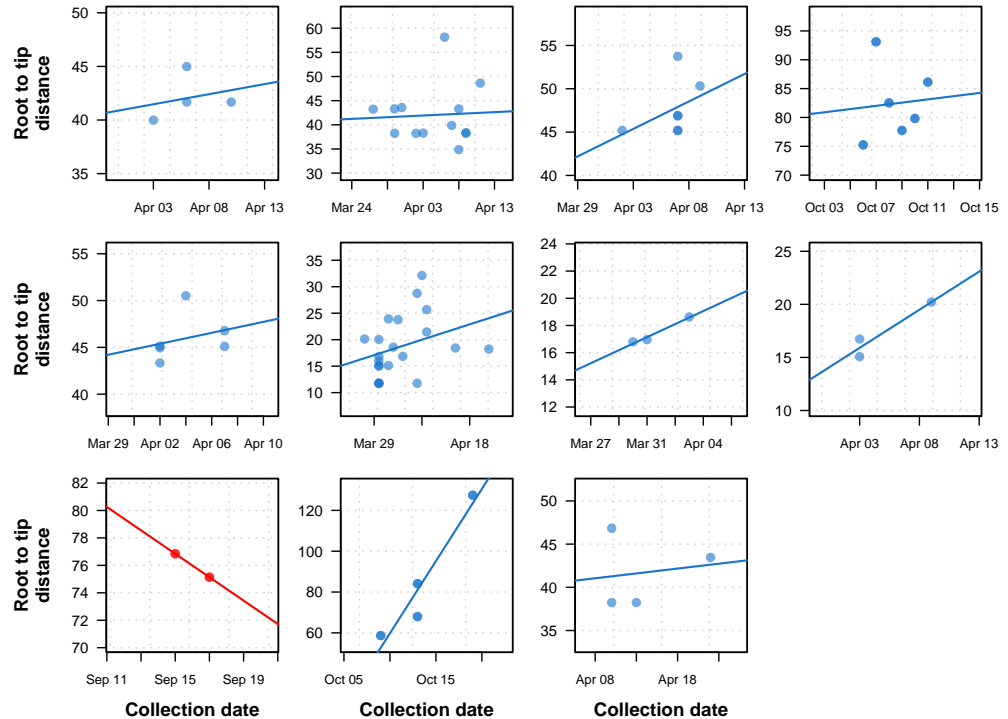


Figure 6. Previously uninformative clusters present temporal signal when using within-sample diversity.

A set of 11 outbreak clusters (one per panel, each plotting the root to tip distance in number of substitutions per genome against time) in which all samples had identical consensus genomes sequences (and therefore no temporal signal). Blue colors indicate those regressions that after utilizing within sample diversity now have a positive slope (temporal signal), and red shows those regressions that have a negative slope (misleading or false positive temporal signal).

Figure 6—figure supplement 1. Phylogenetic trees for SARS-CoV-2.

Figure 6—figure supplement 2. Genetic distance between longitudinal samples.

Figure 6—figure supplement 3. Number of low frequency variants within outbreaks as the outbreak progresses.

273 To illustrate the downstream application of the improved phylogenetic resolution, we inferred
 274 a time-calibrated phylogeny from the phylogeny inferred using the 16-character state model with
 275 the collection dates of the tips using BactDating (*Didelot et al., 2018*) (*Figure 7—figure Supple-*
 276 *ment 1*) and calculated the likelihood of transmission events within potential epidemiologically
 277 identified outbreaks using a Susceptible-Exposed-Infectious-Removed (SEIR) model (*Lekone and*
 278 *Finkenstädt, 2006; Eldholm et al., 2016*). The SEIR model was parameterized with an average la-
 279 tency period of 5.5 days (*Xin et al., 2022*), an infectious period of 6 days (*Byrne et al., 2020*), and
 280 a within-host coalescent rate of 5 days as previously estimated for SARS-CoV-2 (*Wang et al., 2020*).
 281 The likelihood of transmission was calculated for every pair of samples, while the Edmonds algo-
 282 rithm as implemented in the R package *RBGL* (*Carey et al., 2021*) was used to infer the graph with

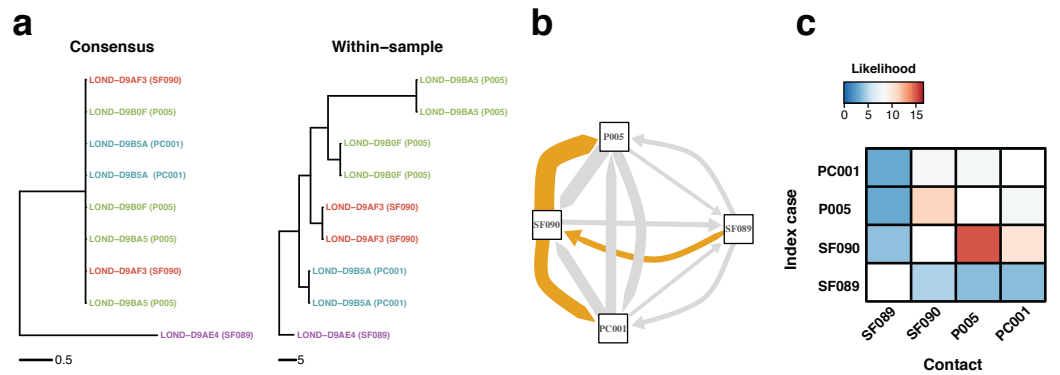


Figure 7. Within-sample variation improves resolution of infectious disease outbreaks. Effect of using low frequency variants in phylogenetic inference. **a** Maximum likelihood phylogeny using the consensus sequences (left) and the alignment leveraging within-sample variation. Replicates of the same sample share the same color. Sample IDs are coded as follows: SF, for staff members; P, for patients; and PC, for patient contacts. **b** Transmission network inferred using within sample variation. Edge width is proportional to the likelihood of direct transmission using a Susceptible-Exposed-Infected-Removed (SEIR) model. Colored edges represent the Edmunds optimum branching and thus the most likely chain. **c** Heatmap of the likelihood of direct transmission between all pairwise pairs of samples using a SEIR model. Vertical axis is the infector while the horizontal axis shows the infectee.

Figure 7—figure supplement 1. Time calibrated phylogenetic trees for SARS-CoV-2.

Figure 7—figure supplement 2. Phylogenetic and transmission for SARS-CoV-2 outbreaks.

284 *Figure 7* represents an example of an outbreak involving 4 hosts, with one patient, one patient
 285 contact, and two hospital staff members. All samples have one technical replicate, while patient
 286 sample also has two serial samples (which were removed for transmission inference). The ML tree
 287 inferred using the consensus sequences (*Figure 7a*, left) shows that most isolates have the exact
 288 same consensus sequence. Although this suggests that all isolates belong to the same outbreak,
 289 the similarity between sequences precludes exact transmission inference. However, the ML tree
 290 inferred using sequences with low frequency variants correctly clusters technical and longitudinal
 291 replicates, and groups the isolates in distinctive sets that better inform transmission inference (*Fig-*
 292 *ure 7b,c*). We applied the same analysis to other potential outbreaks and obtained similar results
 293 (*Figure 7—figure Supplement 2*).

294 Discussion

295 Detailed investigation of transmission events in an infectious disease outbreak is a prerequisite for
 296 effective prevention and control. Although whole-genome sequencing has transformed the field of
 297 pathogen genomics, insufficient pathogen genetic diversity between cases in an outbreak limits the
 298 ability to infer who infected whom. Using multi-hospital SARS-CoV-2 outbreaks and phylogenetic
 299 simulations, we show that including the genetic diversity of subpopulations within a clinical sample
 300 improves phylogenetic reconstruction of SARS-CoV-2 outbreaks and determines the direction of
 301 transmission when using a consensus sequence approach fails to do so.

302 The majority of samples sequenced harbored variants at low frequency. However, most vari-
 303 ants were not consistently called in technical replicates, suggesting they were spurious or unreli-
 304 able. Within-sample variation was less consistent between paired technical replicates with lower
 305 viral load (higher C_t). This is likely to be a consequence of low starting genetic material giving rise
 306 to amplification bias during library preparation and sequencing. Establishing a cut-off for high C_t
 307 values is therefore important to accurately characterize within-host variation. In our study, we ex-
 308 cluded samples with a C_t value higher than 30 cycles based on the diagnostic PCR used at GOSH.

309 Since C_i values are only a surrogate for viral load and are not standardized across different assays
310 (*Evans et al., 2021*), appropriate thresholds would need to be determined for other primary PCR
311 testing assays. Similarly, variant calls at very low frequency were less likely to be present in both
312 technical replicates. These variants at low frequency are thus potentially not genuine and the re-
313 sult of sequencing and variant calling errors. For our work, we removed any variants with an allele
314 frequency lower than 2%. Until sequencing and variant calling technologies improve for low fre-
315 quency variants, technical replicates will remain essential for the study of pathogen within-host
316 diversity in order to distinguish genuine variation from sequencing noise. The effect of this noise
317 on phylogenetic inference will depend on the signal-to-noise ratio and the amount of variation
318 already present in the consensus sequences. Spurious low frequency variation will likely affect
319 only the branch length estimation in phylogenetic inference by adding potentially erroneous calls,
320 unless there is presence of batch bias which could artificially cluster epidemiologically unrelated
321 isolates together.

322 The generation, maintenance and evolution of subpopulations within the host reflect evolu-
323 tionary processes which are meaningful from phylogenetic and epidemiological perspectives. Sub-
324 populations within a host can emerge from three mechanisms: de novo diversification in the host,
325 transmission of a diverse inoculum, or multiple transmission events from different sources. If the
326 subpopulations are the result of de novo mutations, nucleotide polymorphisms within the subpop-
327 ulations accumulate over time and may therefore result in a phylogenetic signal useful for phyloge-
328 netic inference. In our data, longitudinal samples taken at later time points were demonstrated to
329 accrue genomic variation. Although this pattern can be confounded by decreasing viral load as in-
330 fection progresses, C_i values in our dataset were not correlated with a higher genetic distance, and
331 clusters in our data containing both longitudinal and technical replicates also corroborate these
332 results. Transmission of a diverse inoculum also gives rise to phylogenetically informative shared
333 low frequency variants, as our results show that transmission pairs are more likely to share vari-
334 ants at low frequency. The effect of multiple transmission events in the phylogeny depends on the
335 relatedness of both index cases and the bottleneck size in each transmission event.

336 Paired samples with epidemiological links and from the same department shared a higher pro-
337 portion of low frequency variants and were located closer in the consensus tree than samples with
338 no relationship. These patterns suggest that the distribution of low frequency variants is linked to
339 events of epidemiological interest. The fact that technical duplicates shared more within-host di-
340 versity than longitudinal replicates of the same sample suggests that much of the variation within
341 hosts is transitory. Therefore, within-host diversity may be relevant on relatively short time scales,
342 which is precisely where consensus sequences lack resolution. Combining the data derived from
343 fixed alleles in the consensus sequences and transient within-sample minor variation enables an
344 improved understanding of the relatedness of pathogen populations between hosts.

345 The effects of neglecting within-host diversity in phylogenetic inference were analyzed by us-
346 ing simulated sequences under a phylogenetic model that reflects the presence and evolution of
347 within-host diversity. We compared a conventional consensus phylogenetic model and a model
348 that leverages within-sample diversity, and evaluated their ability to infer the known phylogeny.
349 Our proposed phylogenetic model incorporates within-sample variation by explicitly coding ma-
350 jor and minor nucleotides as independent characters in the alignment. We demonstrated that
351 phylogenies inferred using the conventional consensus sequence approach were unresolved and
352 unrepresentative of the known structure of the simulated tree. However, sequences that included
353 within-host diversity showed higher resolution that resulted in phylogenetic trees more similar to
354 the simulated phylogeny. As other mutational models, our 16-state model assumes independence
355 between sites in the alignment. This assumption can be violated due to the presence of genetic link-
356 age, which can be caused by multiple biological processes, such as clonal relationships between
357 microorganisms, recombination or selection of co-evolving sites. To increase the amount of ge-
358 netic linkage due to clonal relationships between organisms, we repeated our simulations using
359 a coalescent model to create the starting tree, and confirmed that the 16-state model still outper-

360 formed the conventional consensus sequence in the presence of high linkage. Other sources of
361 genetic linkage are not accounted for, and their inclusion in phylogenetic inference is out the scope
362 of this work.

363 The proposed phylogenetic model used for the simulations did not include direct base transi-
364 tions and transversions, but rather modelled base change as a process of minor variant acquisition
365 and fixation. Therefore, a base change is composed of the following steps: first a minor variant
366 is gained; then the minor variant increases in frequency and becomes the majority variant; and fi-
367 nally the new variant is fixed. In this way, within-host evolution is partially included in the model as
368 a process of minor variant evolution and eventual lost or fixation. As shown in the simulations, this
369 process of within-host evolution is also captured when the minor bases are simply incorporated
370 as additional states in the Markov chain without explicitly limiting the possible transitions.

371 We complemented the phylogenetic simulations with tree inference of outbreaks simulated us-
372 ing TransPhylo (*Didelot et al., 2017*). We parameterized the simulations to reflect the transmission
373 dynamics of SARS-CoV-2, including a generation time of 5 days and a sampling time of 7 days. Given
374 this parameters, most simulated outbreaks lasted less than a month. We then simulated genetic
375 sequences within the outbreak using fastsimcoal2 (*Excoffier et al., 2013*) as previously described
376 by *De Maio et al. (2018)*. Using a mutation rate of 5×10^{-6} mutations per base per replication cycle,
377 as previously described for SARS-CoV-2 and other betacoronaviruses (*Sender et al., 2021; Amicone*
378 *et al., 2022*), and varying bottleneck sizes, we showed that tree inference using within-sample diver-
379 sity improves as the transmission bottleneck widens, although even at low bottleneck sizes trees
380 inferred using within-sample diversity are more accurate than those inferred using consensus se-
381 quences. Similarly, using varying mutation rates and a constant bottleneck size of 10 pathogens,
382 we showed that tree inference was more accurate as mutation rates increased, although inference
383 using consensus sequences improved only at a very high mutation rate of 10^{-3} mutations per base
384 per cycle, which has mostly been observed in some HIV studies (*Cuevas et al., 2015*). Together,
385 our simulations show that at the short time frame of disease transmission, phylogenetic inference
386 using alignments that contain information regarding within-sample diversity outperform phyloge-
387 nies inferred with consensus sequences, even at narrow transmission bottlenecks and very low
388 mutation rates. Since TransPhylo simulates phylogenetic trees alongside the outbreak simulation,
389 we could directly compare our inferred phylogenies with the known simulated trees. However, al-
390 though phylogenetic trees can inform transmission inference, phylogenetic trees themselves and
391 transmission trees are not interchangeable. Nevertheless, increasing the resolution of phyloge-
392 netic trees can improve inference of transmission chains and calculation of the likelihood of trans-
393 mission events.

394 Previous studies have addressed the use of within-host variation to infer transmission events.
395 *Wymant et al. (2018)* employed a framework based on phylogenetic inference and ancestral state
396 reconstruction of each set of populations detected within read alignments using genomic windows.
397 Our study extends this work by coding genome-wide diversity within the host directly in the align-
398 ment and the phylogenetic model. *De Maio et al. (2018)* proposed direct inference of transmission
399 from sequencing data alongside host exposure time and sampling date within the bayesian frame-
400 work BEAST2 (*Bouckaert et al., 2014*). Our approach is focused on directly improving the temporal
401 and phylogenetic signal of whole-genome sequences, and it's especially suited for use in applica-
402 tions and analysis that employ a phylogenetic tree as input to infer transmission (*Didelot et al.,*
403 *2017*).

404 Apart from transmission inference, phylogenetic trees can be used to infer many parameters
405 of epidemiological interest, such as R_0 or the effective population size. In our work, we showed that
406 the temporal signal of clusters where all isolates had the same sequences increased with the inclu-
407 sion of within-sample diversity, which in turns allows better inference of phylogenetic trees. When
408 analyzing specific outbreaks, we showed that groups of samples without genetic differences were
409 clustered apart from other isolates of the outbreak, providing additional information on genetic
410 relationships that could be used for transmission inference or to better understand the genetic

411 structure of the outbreak. Even though transmission inference can be improved with epidemiological data such as collection dates even when all isolates have the same genetic sequences, such data can't provide information regarding how samples cluster within the outbreak. Additionally, the order of collection dates not always correspond to the order of infection.

415 Future work will extend this model by including allele frequency data in addition to independent characters for major and minor variants. Moreover, to limit the number of character states we only allowed two variants at each position. Transmission inference of pathogens with high levels of within-host diversity, for instance as observed in HIV, could benefit from including more than two alleles. In those cases, the number of possible character state combinations would be too large, and therefore other methods such as phyloscanner *Wymant et al. (2018)* could resolve transmission events more accurately. However, it's important to note that the low frequency of a third allele could result in more sequencing and mapping errors which could in turn bias phylogenetic inference and genomic analysis. Phylogenetic models that explicitly include dynamics of within-sample variation and sequencing error may further improve phylogenetic inference or allow researchers to better estimate parameters of interest, including R_0 , bottleneck size, transmissibility and the origin of outbreaks.

427 In line with conventional consensus sequencing approaches, we used a reference sequence for genome alignment and variant calling. Although widely used, one limitation of this approach is a potential mapping bias causing some reads to reflect the reference base at low frequencies at a position where only a variant should be present. Although we applied stringent quality filtering, we cannot rule out the persistence of some false positive minor variants. Using genome graphs to map to a reference that encompasses a wider spectrum of variation may alleviate this problem, and could be an interesting addition to pathogen population genomic analysis.

434 Our results demonstrate that within-sample variation can be leveraged to increase the resolution of phylogenetic trees and improve our understanding of who infected whom. Using SARS-CoV-2 hospital outbreaks and simulations, we show that variants at low frequencies are consistent within sample replicates, phylogenetically informative and are more often shared among epidemiologically related contacts. By coding within-sample variation directly in the alignment, the additional genetic information can be easily incorporated in phylogenetic inference, facilitating its application within existing epidemiology pipelines and public health infrastructure. We propose that pathogen phylogenetic models should accommodate within-host variation to improve the understanding of infectious disease transmission and aid infection control measures.

443 **Materials and Methods**

444 **Model for within-host diversity**

445 To test the accuracy of different models at inferring known phylogenies, 100 random phylogenetic trees with 100 tips each were generated using the function *rtree* within the R package *ape* (*Paradis et al., 2004*). Whole-genome alignments were simulated from the random 100 phylogenies with the function *SimSeq* of the R package *phangorn* (*R Core Team and R Foundation for Statistical Computing, 2021; Schliep, 2011*) using a model with 16 character states that represent the combinations of the 4 nucleotides with each other as minor and major alleles (*Figure 3d*). Three substitution rates for the model were considered: a rate at which minor variants evolve, equal to 1; the rate at which minor variants are lost, leaving only the major nucleotide at that position, equal to 100; and the rate at which minor/major variants are switched, equal to 200. This rates result in fixed bases (A, C, G, and T) being 100 times more frequent than low frequency bases. A different set of simulations was performed using rates that promote a high rate of low frequency variation by having a lower rate of variant loss and switch (rates 1, 10, 10 for minor variant evolution, loss and switch, respectively); a low amount of low frequency variation by increasing the rates of variant switch and loss (1, 10, 100); and using a Jukes-Cantor model of sequence evolution and therefore resulting in no minor variants.

460 Two types of alignments were generated from the simulated genomes: a consensus sequence,
461 where only the major allele was considered; and an alignment that retained the major and minor
462 allele information as independent character states. RaxML-NG (*Kozlov et al., 2019*) was used to
463 infer phylogenetic trees. The consensus sequence was analyzed with a GTR+ γ model, while the
464 PROTGTR+ γ model was used for the alignment with intra-host diversity and major/minor variant
465 information.

466 Several metrics were used to compare the 200 inferred phylogenetic trees with their respective
467 starting phylogeny from which the sequences were simulated (Supplementary file 3). We chose
468 metrics available in R suitable for unrooted trees, using the option 'rooted=FALSE' where appro-
469 priate. The Robinson-Foulds (RF) distance (*Robinson and Foulds, 1981*) calculates the number of
470 splits differing between both phylogenetic trees. For the weighted Robinson-Foulds (wRF), the dis-
471 tance is expressed in terms of the branch lengths of the differing splits. The Kuhner-Felsenstein
472 distance (*Kuhner and Felsenstein, 1994*) considers the edge length differences in all splits, regard-
473 less of whether the topology is shared or not. Last, the Penny-Steel distance or path difference
474 metric (*Steel and Penny, 1993*) calculates the pairwise differences in the path of each pair of tips,
475 with the weighted Penny-Steel distance (wPS) using branch length to compute the path differences.
476 All functions were used as implemented in the package *phangorn* (*Schliep, 2011*) within R (*R Core
477 Team and R Foundation for Statistical Computing, 2021*).

478 **Outbreak simulations**

479 Disease outbreaks of size between 10 and 15 hosts were simulated using TransPhylo (*Didelot et al.,
480 2017*), with a mean generation time of 5 days and a mean sampling time of 7 days, both parameters
481 with standard deviation of 1 day (*Wang et al., 2020; Hart et al., 2022*). To ensure that the outbreak
482 ends, the negative binomial distribution for the offspring number was set with a mean of 1 and a
483 dispersion parameter of 0.5, resulting in a basic reproductive number (R_0) of 1. To simplify the sim-
484 ulations, all hosts from the outbreak were sampled. A total of 20 outbreaks were simulated. The
485 population evolution within and between hosts was simulated using fastsimcoal2 (*Excoffier et al.,
486 2013*) as previously described by De Maio et al (*De Maio et al., 2018*), where transmissions are incor-
487 porated as population migrations with a given bottleneck size and populations evolve with a given
488 mutation rate per generation time. Sequences were simulated for a within-host population size of
489 1000 and a genome size of 1000bp. To understand the effect of transmission bottleneck size in
490 phylogenetic inference, varying values of bottleneck size were used along a constant mutation rate
491 of 5×10^{-6} mutations per base per generation cycle. Additionally, sequences were simulated at dif-
492 ferent mutation rates with a constant bottleneck size of 10 pathogens. Sequences with the varying
493 bottleneck sizes and mutation rates were simulated using the same 20 simulated outbreaks. Phy-
494 logenetic trees were inferred from the alignments using RaxML-NG as previously described. The
495 resulting trees were time-calibrated using the additive uncorrelated relaxed clock model (ARC) as
496 implemented in BactDating (*Didelot et al., 2018*). The root of the outbreak was inferred as part of
497 the dating model. The inferred trees were compared to the known simulated phylogenies using
498 the Kuhner-Felsenstein distance (*Kuhner and Felsenstein, 1994*).

499 **Amplification and whole-genome sequencing**

500 SARS-CoV-2 real-time qPCR confirmed isolates from London hospitals were collected as part of the
501 routine diagnostic service at Great Ormond Street Hospital NHS Foundation Trust (GOSH) (*Storey
502 et al., 2021*) and the COVID-19 Genomics UK Consortium (COG-UK) (*COVID-19 Genomics UK (COG-
503 UK), 2020*) between March and December 2020, in addition to epidemiological and patient meta-
504 data (Supplementary file 2). Multiple types of samples were collected: isolates from different pa-
505 tients; longitudinal replicates, where multiple isolates were collected from the same patient at
506 different time points; and technical replicates, where multiple sequencing runs were performed
507 from the same biological isolate. SARS-CoV-2 whole-genome sequencing was performed by UCL
508 Genomics. cDNA and multiplex PCR reactions were prepared following the ARTIC nCoV-2019 se-

509 quencing protocol (*Tyson et al., 2020*). The ARTIC V3 primer scheme (*ARTIC Network, 2021*) was
510 used for the multiplex PCR, with a 65°C, 5 min annealing/extension temperature. Pools 1 and 2 mul-
511 tiplex PCRs were run for 35 cycles. 5µL of each PCR were combined and 20µL nuclease-free water
512 added. Libraries were prepared on the Agilent Bravo NGS workstation option B using Illumina DNA
513 prep (Cat. 20018705) with unique dual indexes (Cat. 20027213/14/15/16). Equal volumes of the
514 final libraries were pooled, bead purified and sequenced on the Illumina NextSeq 500 platform
515 using a Mid Output 150 cycle flowcell (Cat. 20024904) (2 x 75bp paired ends) at a final loading
516 concentration of 1.1pM.

517 **Whole-genome sequence analysis of SARS-CoV-2 sequences**

518 Raw illumina reads were quality trimmed using Trimmomatic (*Bolger et al., 2014*) with a minimum
519 mean quality per base of 20 in a 4-base wide sliding window. The 5 leading and trailing bases of
520 each read were removed, and reads with an average quality lower than 20 were discarded. The
521 resulting reads were aligned against the Wuhan-Hu-1 reference genome (GenBank NC_45512.2,
522 GISAID EPI_ISL_402125) using BWA-mem v0.7.17 with default parameters (*Li and Durbin, 2010*).
523 The alignments were subsequently sorted by position using SAMtools v1.14 (*Li et al., 2009*). Primer
524 sequences were masked using ivar (*Grubaugh et al., 2019*).

525 Single-nucleotide variants were identified using the pileup functionality of samtools (*Li et al.,*
526 *2009*) via the pysam package in Python (<https://github.com/pysam-developers/pysam>). Variants
527 were further filtered using bcftools (*Danecek et al., 2011*). Only variants with a minimum depth
528 of 50x and a minimum base quality and mapping quality of 30 were kept. Additionally, variants
529 within low complexity regions identified by sdust (<https://github.com/lh3/sdust>) were removed.
530 Previously identified problematic sites were masked to avoid systematic sequencing errors and
531 phylogenetic bias (*De Maio et al., 2020*). For positions where only one base was present, the min-
532 imum depth was 20 reads, with at least 5 reads in each direction. Positions with low frequency
533 variants were filtered if the total coverage at that position was less than 100x, with at least 20
534 reads in total and 5 reads in each strand supporting each of the main two alleles.

535 Two different alignments were prepared from the data. First, an alignment of the consensus
536 sequence where the most prevalent base at each position was kept. Variants where the most
537 prevalent allele was not supported by more than 60% of the reads were considered ambiguous.
538 Additionally, an alignment reflecting within-sample variation at each position as well as which base
539 is the most prevalent and which one appears at a lower frequency by using the IUPAC nomencla-
540 ture for amino acids (*IUPAC-IUB Joint Commission on Biochemical Nomenclature (JCBN), 1984*).

541 For the two different alignments, maximum likelihood phylogenies were inferred by using RAxML-
542 NG (*Kozlov et al., 2019*) with 20 starting trees (10 random and 10 parsimony), 100 bootstrap repli-
543 cates, and a minimum branch length of 10^{-9} . For the consensus sequence, the GTR model was
544 used. For the alignment reflecting within-host diversity, a model with amino acid nomenclature
545 (PROTGTR) was used. All models allowed for a γ distributed rate of variation among sites. Phy-
546 logenetic trees were time-calibrated using the known collection dates and the ARC model within
547 BactDating (*Didelot et al., 2018*). For transmission inference, the dated phylogeny was used with
548 the longitudinal replicates removed by keeping the earliest sampled isolate. The likelihood of trans-
549 mission was calculated using a Susceptible-Exposed-Infectious-Removed (SEIR) model (*Lekone and*
550 *Finkenstädt, 2006; Eldholm et al., 2016*).

551 **Data availability**

552 Samples sequenced as part of this study have been submitted to the European Nucleotide Archive
553 under accession PRJEB53224. Sample metadata is included in Supplementary file 1.

554 **Code availability**

555 All custom code used in this article can be accessed at
556 https://github.com/arturotorres/scov2_withinHost.git.

557 Acknowledgements

558 The authors dedicate this article to the hospital staff members and patients who died of coron-
559 avirus disease 2019. They also thank all staff and patients who have taken part in the study. In
560 addition, the authors are very grateful to the Great Ormond Street laboratory staff, the staff at the
561 Camelia Botnar Laboratory, the Great Ormond Street Institute of Child Health and the COVID-19
562 sequencing team at UCLG who worked tirelessly to ensure that all polymerase chain reaction tests
563 and sequencing work were completed in a timely manner during the COVID-19 pandemic. All au-
564 thors acknowledge UCL Computer Science Technical Support Group (TSG) and the UCL Department
565 of Computer Science High Performance Computing Cluster. LG was supported by the Wellcome
566 Trust (201470/Z/16/Z), the National Institute of Allergy and Infectious Diseases of the National In-
567 stitutes of Health under award number 1R01AI146338 and by the GOSH/ICH Biomedical Research
568 Centre. XD was supported by the NIHR Health Protection Research Unit in Genomics and Enabling
569 Data.

570 Competing interests

571 The authors declare no competing interests

572 Ethics declarations

573 Ethical approval was obtained for all individual studies from which this data was derived.

574 References

- 575 **Amicone M**, Borges V, Alves MJ, Isidro J, Zé-Zé L, Duarte S, Vieira L, Guiomar R, Gomes JP, Gordo I. Muta-
576 tion rate of SARS-CoV-2 and emergence of mutators during experimental evolution. *Evolution, Medicine,*
577 *and Public Health.* 2022 Jan; 10(1):142–155. <https://academic.oup.com/emph/article/10/1/142/6555377>, doi:
578 10.1093/EMPH/EOAC010, publisher: Oxford Academic.
- 579 **ARTIC Network**, ARTIC nanopore protocol for nCoV2019 novel coronavirus [Internet]; 2021. <https://github.com/artic-network/artic-ncov2019>.
- 581 **Bolger AM**, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformat-*
582 *ics.* 2014 Aug; 30(15):2114–2120. <https://academic.oup.com/bioinformatics/article/30/15/2114/2390096>, doi:
583 10.1093/BIOINFORMATICS/BTU170, publisher: Oxford Academic.
- 584 **Bouckaert R**, Heled J, Kühnert D, Vaughan T, Wu CH, Xie D, Suchard MA, Rambaut A, Drummond AJ. BEAST
585 2: A Software Platform for Bayesian Evolutionary Analysis. *PLoS Computational Biology.* 2014 Apr;
586 10(4):e1003537. <https://dx.plos.org/10.1371/journal.pcbi.1003537>, doi: 10.1371/journal.pcbi.1003537.
- 587 **Byrne AW**, McEvoy D, Collins AB, Hunt K, Casey M, Barber A, Butler F, Griffin J, Lane EA, McAloon C, O'Brien
588 K, Wall P, Walsh KA, More SJ. Inferred duration of infectious period of SARS-CoV-2: rapid scoping review
589 and analysis of available evidence for asymptomatic and symptomatic COVID-19 cases. *BMJ open.* 2020
590 Aug; 10(8):e039856. <http://www.ncbi.nlm.nih.gov/pubmed/32759252>, doi: 10.1136/bmjopen-2020-039856,
591 publisher: British Medical Journal Publishing Group.
- 592 **Campbell F**, Strang C, Ferguson N, Cori A, Jombart T. When are pathogen genome sequences informative of
593 transmission events? *PLoS Pathogens.* 2018 Feb; 14(2):e1006885. [https://journals.plos.org/plospathogens/](https://journals.plos.org/plospathogens/article?id=10.1371/journal.ppat.1006885)
594 [article?id=10.1371/journal.ppat.1006885](https://journals.plos.org/plospathogens/article?id=10.1371/journal.ppat.1006885), doi: 10.1371/journal.ppat.1006885, publisher: Public Library of Sci-
595 ence.
- 596 **Carey V**, Long L, Gentleman R, RBGL: An interface to the BOOST graph library; 2021.
- 597 **COVID-19 Genomics UK (COG-UK)**. An integrated national scale SARS-CoV-2 genomic surveillance net-
598 work. *The Lancet Microbe.* 2020 Jul; 1(3):e99–e100. <http://www.ncbi.nlm.nih.gov/pubmed/32835336>, doi:
599 10.1016/S2666-5247(20)30054-9.
- 600 **Cuevas JM**, Geller R, Garijo R, López-Aldeguer J, Sanjuán R. Extremely High Mutation Rate of HIV-1 In Vivo.
601 *PLOS Biology.* 2015 Sep; 13(9):e1002251. [https://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.](https://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.1002251)
602 [1002251](https://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.1002251), doi: 10.1371/journal.pbio.1002251, publisher: Public Library of Science.

- 603 **Danecek P**, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G, Marth GT, Sherry
604 ST, McVean G, Durbin R. The variant call format and VCFtools. *Bioinformatics*. 2011 Aug; 27(15):2156–2158.
605 <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btr330>, doi: 10.1093/bioin-
606 formatics/btr330.
- 607 **De Maio N**, Walker C, Borges R, Weilguny L, Slodkowitz G, Goldman N. Masking strategies for SARS-CoV-2
608 alignments. *Virological*. 2020; <http://europepmc.org/abstract/PPR/PPR262711>, publisher: virological.org.
- 609 **De Maio N**, Worby CJ, Wilson DJ, Stoesser N. Bayesian reconstruction of transmission within outbreaks us-
610 ing genomic variants. *PLOS Computational Biology*. 2018 Apr; 14(4):e1006117. [https://journals.plos.org/](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1006117)
611 [ploscompbiol/article?id=10.1371/journal.pcbi.1006117](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1006117), doi: 10.1371/journal.pcbi.1006117, publisher: Public
612 Library of Science.
- 613 **Didelot X**, Bowden R, Wilson DJ, Peto TEA, Crook DW. Transforming clinical microbiology with bacterial genome
614 sequencing. *Nature Reviews Genetics*. 2012 Sep; 13(9):601–612. <http://www.nature.com/articles/nrg3226>, doi:
615 10.1038/nrg3226.
- 616 **Didelot X**, Croucher NJ, Bentley SD, Harris SR, Wilson DJ. Bayesian inference of ancestral dates on bacterial
617 phylogenetic trees. *Nucleic Acids Research*. 2018 Dec; 46(22):e134–e134. [https://academic.oup.com/nar/](https://academic.oup.com/nar/advance-article-abstract/doi/10.1093/nar/gky783/5089898)
618 [advance-article-abstract/doi/10.1093/nar/gky783/5089898](https://academic.oup.com/nar/advance-article-abstract/doi/10.1093/nar/gky783/5089898), doi: 10.1093/nar/gky783.
- 619 **Didelot X**, Fraser C, Gardy J, Colijn C. Genomic Infectious Disease Epidemiology in Partially Sampled and Ongo-
620 ing Outbreaks. *Molecular Biology and Evolution*. 2017 Apr; 34(4):997–1007. [https://doi.org/10.1093/molbev/](https://doi.org/10.1093/molbev/msw275)
621 [msw275](https://doi.org/10.1093/molbev/msw275), doi: 10.1093/molbev/msw275.
- 622 **Dinis JM**, Florek NW, Fatola OO, Moncla LH, Mutschler JP, Charlier OK, Meece JK, Belongia EA,
623 Friedrich TC. Deep Sequencing Reveals Potential Antigenic Variants at Low Frequencies in In-
624 fluenza A Virus-Infected Humans. *Journal of Virology*. 2016 Apr; 90(7):3355–3365. [https://](https://journals.asm.org/doi/abs/10.1128/JVI.03248-15)
625 journals.asm.org/doi/abs/10.1128/JVI.03248-15, doi: 10.1128/JVI.03248-15/ASSET/88062BA5-1056-48A8-
626 AD63-141542EE7BEB/ASSETS/GRAPHIC/ZJV9990914640003.JPEG, publisher: American Society for Microbiol-
627 ogy.
- 628 **Eldholm V**, Rieux A, Monteserin J, Lopez JM, Palmero D, Lopez B, Ritacco V, Didelot X, Balloux F. Impact of HIV
629 co-infection on the evolution and transmission of multidrug-resistant tuberculosis. *eLife*. 2016 Aug; 5(AU-
630 GUST). <http://www.ncbi.nlm.nih.gov/pubmed/27502557>, doi: 10.7554/eLife.16644, publisher: eLife Sciences
631 Publications Ltd.
- 632 **Evans D**, Cowen S, Kammel M, O'Sullivan DM, Stewart G, Grunert HP, Moran-Gilad J, Verwilt J, In J, Vandesom-
633 pele J, Harris K, Hong KH, Storey N, Hingley-Wilson S, Dühning U, Bae YK, Foy CA, Braybrook J, Zeichhardt
634 H, Huggett JF. The Dangers of Using Cq to Quantify Nucleic Acid in Biological Samples: A Lesson From
635 COVID-19. *Clinical chemistry*. 2021 Dec; 68(1):153–162. <http://www.ncbi.nlm.nih.gov/pubmed/34633030>, doi:
636 10.1093/clinchem/hvab219, publisher: Oxford Academic.
- 637 **Excoffier L**, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. Robust Demographic Inference from Genomic
638 and SNP Data. *PLOS Genetics*. 2013 Oct; 9(10):e1003905. [https://journals.plos.org/plosgenetics/article?id=10.](https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1003905)
639 [1371/journal.pgen.1003905](https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1003905), doi: 10.1371/journal.pgen.1003905, publisher: Public Library of Science.
- 640 **Grubaugh ND**, Gangavarapu K, Quick J, Matteson NL, De Jesus JG, Main BJ, Tan AL, Paul LM, Brackney DE, Grewal
641 S, Gurfield N, Van Rompay KKA, Isern S, Michael SF, Coffey LL, Loman NJ, Andersen KG. An amplicon-based
642 sequencing framework for accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome*
643 *Biology*. 2019 Jan; 20(1):8. <https://doi.org/10.1186/s13059-018-1618-7>, doi: 10.1186/s13059-018-1618-7.
- 644 **Han A**, Czajkowski LM, Donaldson A, Baus HA, Reed SM, Athota RS, Bristol T, Rosas LA, Cervantes-Medina A,
645 Taubenberger JK, Memoli MJ. A Dose-finding Study of a Wild-type Influenza A(H3N2) Virus in a Healthy Vol-
646 unteer Human Challenge Model. *Clinical Infectious Diseases*. 2019 Nov; 69(12):2082–2090. [https://academic.](https://academic.oup.com/cid/article/69/12/2082/5321121)
647 [oup.com/cid/article/69/12/2082/5321121](https://academic.oup.com/cid/article/69/12/2082/5321121), doi: 10.1093/cid/ciz141, publisher: Oxford Academic.
- 648 **Hart WS**, Miller E, Andrews NJ, Waight P, Maini PK, Funk S, Thompson RN. Generation time of the alpha and
649 delta SARS-CoV-2 variants: an epidemiological analysis. *The Lancet Infectious Diseases*. 2022 May; 22(5):603–
650 610. [https://www.thelancet.com/journals/laninf/article/PIIS1473-3099\(22\)00001-9/fulltext](https://www.thelancet.com/journals/laninf/article/PIIS1473-3099(22)00001-9/fulltext), doi: 10.1016/S1473-
651 3099(22)00001-9, publisher: Elsevier.
- 652 **Harvey WT**, Carabelli AM, Jackson B, Gupta RK, Thomson EC, Harrison EM, Ludden C, Reeve R, Rambaut A, Pea-
653 cock SJ, Robertson DL. SARS-CoV-2 variants, spike mutations and immune escape. *Nature Reviews Micro-*
654 *biology*. 2021 Jul; 19(7):409–424. <http://www.nature.com/articles/s41579-021-00573-0>, doi: 10.1038/s41579-
655 021-00573-0, publisher: Nature Publishing Group ISBN: 0123456789.

- 656 **IUPAC-IUB Joint Commission on Biochemical Nomenclature (JCBN).** Nomenclature and symbolism for
657 amino acids and peptides. Recommendations 1983. European journal of biochemistry. 1984 Jan; 138(1):9-
658 37. <http://www.ncbi.nlm.nih.gov/pubmed/6692818>, doi: 10.1111/j.1432-1033.1984.tb07877.x, publisher: Eur J
659 Biochem.
- 660 **Kozlov AM,** Darriba D, Flouri T, Morel B, Stamatakis A. RAxML-NG: A fast, scalable and user-friendly tool for
661 maximum likelihood phylogenetic inference. Bioinformatics. 2019; 35(21):4453-4455. doi: 10.1093/bioinform-
662 atics/btz305.
- 663 **Kuhner MK,** Felsenstein J. A simulation comparison of phylogeny algorithms under equal and unequal evolu-
664 tionary rates. Molecular Biology and Evolution. 1994 May; 11(3):459-468. [https://academic.oup.com/mbe/
665 article/11/3/459/1104326](https://academic.oup.com/mbe/article/11/3/459/1104326), doi: 10.1093/OXFORDJOURNALS.MOLBEV.A040126, publisher: Oxford Academic
666 ISBN: 1030013\$02.0.
- 667 **Lee LYW,** Rozmanowski S, Pang M, Charlett A, Anderson C, Hughes GJ, Barnard M, Peto L, Vipond R, Sienkiewicz
668 A, Hopkins S, Bell J, Crook DW, Gent N, Walker AS, Peto TEA, Eyre DW. Severe Acute Respiratory Syndrome
669 Coronavirus 2 (SARS-CoV-2) Infectivity by Viral Load, S Gene Variants and Demographic Factors, and the
670 Utility of Lateral Flow Devices to Prevent Transmission. Clinical Infectious Diseases. 2022 Jan; 74(3):407-415.
671 <https://doi.org/10.1093/cid/ciab421>, doi: 10.1093/cid/ciab421.
- 672 **Leitner T.** Phylogenetics in HIV transmission: Taking within-host diversity into account. Current Opinion
673 in HIV and AIDS. 2019 May; 14(3):181-187. [https://journals.lww.com/co-hivandaids/Fulltext/2019/05000/
674 Phylogenetics_in_HIV_transmission__taking.6.aspx](https://journals.lww.com/co-hivandaids/Fulltext/2019/05000/Phylogenetics_in_HIV_transmission__taking.6.aspx), doi: 10.1097/COH.0000000000000536, publisher: Lip-
675 pincott Williams and Wilkins.
- 676 **Lekone PE,** Finkenstädt BF. Statistical Inference in a Stochastic Epidemic SEIR Model with Control
677 Intervention: Ebola as a Case Study. Biometrics. 2006; 62(4):1170-1177. [https://onlinelibrary.
678 wiley.com/doi/abs/10.1111/j.1541-0420.2006.00609.x](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1541-0420.2006.00609.x), doi: 10.1111/j.1541-0420.2006.00609.x, _eprint:
679 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1541-0420.2006.00609.x>.
- 680 **Li H,** Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform. Bioinformatics.
681 2010 Mar; 26(5):589-595. [https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/
682 btp698](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btp698), doi: 10.1093/bioinformatics/btp698.
- 683 **Li H,** Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. The Sequence
684 Alignment/Map format and SAMtools. Bioinformatics. 2009 Aug; 25(16):2078-2079. [https://academic.oup.
685 com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btp352](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btp352), doi: 10.1093/bioinformatics/btp352.
- 686 **Lieberman TD,** Wilson D, Misra R, Xiong LL, Moodley P, Cohen T, Kishony R. Genomic diversity in autopsy
687 samples reveals within-host dissemination of HIV-associated Mycobacterium tuberculosis. Nature medicine.
688 2016 Dec; 22(12):1470-1474. <http://www.ncbi.nlm.nih.gov/pubmed/27798613>, doi: 10.1038/nm.4205, pub-
689 lisher: NIH Public Access.
- 690 **McCrone JT,** Lauring AS. Measurements of Intrahost Viral Diversity Are Extremely Sensitive to Systematic Errors
691 in Variant Calling. Journal of Virology. 2016 Aug; 90(15):6884-6895. <http://broadinstitute.github>, publisher:
692 American Society for Microbiology.
- 693 **Mongkolrattanothai K,** Gray BM, Mankin P, Stanfill AB, Pearl RH, Wallace LJ, Vegunta RK. Simultaneous car-
694 riage of multiple genotypes of Staphylococcus aureus in children. Journal of Medical Microbiology. 2011
695 Mar; 60(3):317-322. <https://www.microbiologyresearch.org/content/journal/jmm/10.1099/jmm.0.025841-0>, doi:
696 10.1099/JMM.0.025841-0/CITE/REFWORKS, publisher: Microbiology Society.
- 697 **Murphy BR,** Clements ML, Madore HP, Steinberg J, O'Donnell S, Betts R, Demico D, Reichman RC, Dolin R,
698 Maassab HF. Dose Response of Cold-Adapted, Reassortant Influenza A/California/10/78 Virus (H1N1) in
699 Adult Volunteers. Journal of Infectious Diseases. 1984 May; 149(5):816-816. [https://academic.oup.com/jid/
700 article/149/5/816/2190254](https://academic.oup.com/jid/article/149/5/816/2190254), doi: 10.1093/infdis/149.5.816, publisher: Oxford Academic.
- 701 **Paradis E,** Claude J, Strimmer K. APE: Analyses of Phylogenetics and Evolution in R language. Bioinformatics.
702 2004 Jan; 20(2):289-290. [https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/
703 btg412](https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btg412), doi: 10.1093/bioinformatics/btg412, arXiv: 1301.2609v5 ISBN: 1367-4803.
- 704 **Popa A,** Genger JW, Nicholson MD, Penz T, Schmid D, Aberle SW, Agerer B, Lercher A, Endler L, Colaço H,
705 Smyth M, Schuster M, Grau ML, Martínez-Jiménez F, Pich O, Borena W, Pawelka E, Keszei Z, Senekowitsch
706 M, Laine J, et al. Genomic epidemiology of superspreading events in Austria reveals mutational dynam-
707 ics and transmission properties of SARS-CoV-2. Science Translational Medicine. 2020 Dec; 12(573):2555.
708 <http://stm.sciencemag.org/>, doi: 10.1126/scitranslmed.abe2555, publisher: American Association for the Ad-
709 vancement of Science.

- 710 **R Core Team**, R Foundation for Statistical Computing, R: A Language and Environment for Statistical Computing.
711 Vienna, Austria; 2021. <https://www.r-project.org/>.
- 712 **Rieux A**, Balloux F. Inferences from tip-calibrated phylogenies: a review and a practical guide. *Molec-*
713 *ular Ecology*. 2016 May; 25(9):1911–1924. <https://onlinelibrary.wiley.com/doi/10.1111/mec.13586>, doi:
714 [10.1111/mec.13586](https://doi.org/10.1111/mec.13586).
- 715 **Robinson DF**, Foulds LR. Comparison of phylogenetic trees. *Mathematical Biosciences*. 1981 Feb; 53(1-2):131–
716 147. <https://linkinghub.elsevier.com/retrieve/pii/0025556481900432>, doi: 10.1016/0025-5564(81)90043-2, pub-
717 lisher: Elsevier.
- 718 **Schliep KP**. phangorn: phylogenetic analysis in R. *Bioinformatics*. 2011 Feb; 27(4):592–593. [https://doi.org/10.](https://doi.org/10.1093/bioinformatics/btq706)
719 [1093/bioinformatics/btq706](https://doi.org/10.1093/bioinformatics/btq706), doi: 10.1093/bioinformatics/btq706.
- 720 **Sender R**, Bar-On YM, Gleizer S, Bernshtein B, Flamholz A, Phillips R, Milo R. The total num-
721 ber and mass of SARS-CoV-2 virions. *Proceedings of the National Academy of Sciences of*
722 *the United States of America*. 2021 Jun; 118(25). <https://doi.org/10.1073/pnas.2024815118>, doi:
723 [10.1073/PNAS.2024815118/SUPPL_FILE/PNAS.2024815118.SD01.XLSX](https://doi.org/10.1073/PNAS.2024815118/SUPPL_FILE/PNAS.2024815118.SD01.XLSX), publisher: National Academy of Sci-
724 ences.
- 725 **Spinelli MA**, Glidden DV, Gennatas ED, Bielecki M, Beyrer C, Rutherford G, Chambers H, Goosby E, Gandhi
726 M. Importance of non-pharmaceutical interventions in lowering the viral inoculum to reduce suscepti-
727 bility to infection by SARS-CoV-2 and potentially disease severity. *The Lancet Infectious Diseases*. 2021
728 Sep; 21(9):e296–e301. <http://www.thelancet.com/article/S1473309920309828/fulltext>, doi: [10.1016/S1473-](https://doi.org/10.1016/S1473-3099(20)30982-8/ATTACHMENT/5545EA7B-1383-4BC8-A8F4-3FF678895CE5/MMC1.PDF)
729 [3099\(20\)30982-8/ATTACHMENT/5545EA7B-1383-4BC8-A8F4-3FF678895CE5/MMC1.PDF](https://doi.org/10.1016/S1473-3099(20)30982-8/ATTACHMENT/5545EA7B-1383-4BC8-A8F4-3FF678895CE5/MMC1.PDF), publisher: Lancet
730 Publishing Group.
- 731 **Steel MA**, Penny D. Distributions of Tree Comparison Metrics-Some New Results. *Systematic Biology*. 1993
732 Jun; 42(2):126. <https://www.jstor.org/stable/2992536?origin=crossref>, doi: 10.2307/2992536, publisher: JSTOR.
- 733 **Storey N**, Brown JR, Pereira RPA, O’Sullivan DM, Huggett JF, Williams R, Breuer J, Harris KA. Single base mutations
734 in the nucleocapsid gene of SARS-CoV-2 affects amplification efficiency of sequence variants and may lead to
735 assay failure. *Journal of Clinical Virology Plus*. 2021 Sep; 1(3):100037. [https://linkinghub.elsevier.com/retrieve/](https://linkinghub.elsevier.com/retrieve/pii/S2667038021000296)
736 [pii/S2667038021000296](https://linkinghub.elsevier.com/retrieve/pii/S2667038021000296), doi: [10.1016/j.jcvp.2021.100037](https://doi.org/10.1016/j.jcvp.2021.100037), publisher: Elsevier.
- 737 **Tom MR**, Mina MJ. To Interpret the SARS-CoV-2 Test, Consider the Cycle Threshold Value. *Clinical Infec-*
738 *tious Diseases*. 2020 Nov; 71(16):2252–2254. <https://academic.oup.com/cid/article/71/16/2252/5841456>, doi:
739 [10.1093/cid/ciaa619](https://doi.org/10.1093/cid/ciaa619), publisher: Oxford Academic.
- 740 **Tonkin-Hill G**, Martincorena I, Amato R, Lawson AR, Gerstung M, Johnston I, Jackson DK, Park N, Lensing SV,
741 Quail MA, Gonçalves S, Ariani C, Spencer Chapman M, Hamilton WL, Meredith LW, Hall G, Jahun AS, Chaudhry
742 Y, Hosmillo M, Pinckert ML, et al. Patterns of within-host genetic diversity in SARS-CoV-2. *eLife*. 2021 Aug;
743 10:e66857. <https://doi.org/10.7554/eLife.66857>, doi: [10.7554/eLife.66857](https://doi.org/10.7554/eLife.66857), publisher: eLife Sciences Publica-
744 tions, Ltd.
- 745 **Trunfio M**, Calcagno A, Bonora S, Di Perri G. Lowering SARS-CoV-2 viral load might affect trans-
746 mission but not disease severity in secondary cases. *The Lancet Infectious Diseases*. 2021
747 Jul; 21(7):914–915. <http://www.thelancet.com/article/S147330992100205X/fulltext>, doi: [10.1016/S1473-](https://doi.org/10.1016/S1473-3099(21)00205-X/ATTACHMENT/CC1869F9-D21E-4663-9818-964D15D1081D/MMC1.PDF)
748 [3099\(21\)00205-X/ATTACHMENT/CC1869F9-D21E-4663-9818-964D15D1081D/MMC1.PDF](https://doi.org/10.1016/S1473-3099(21)00205-X/ATTACHMENT/CC1869F9-D21E-4663-9818-964D15D1081D/MMC1.PDF), publisher: Lancet
749 Publishing Group.
- 750 **Tyson JR**, James P, Stoddart D, Sparks N, Wickenhagen A, Hall G, Choi JH, Lapointe H, Kamelian K, Smith
751 AD, Prystajeky N, Goodfellow I, Wilson SJ, Harrigan R, Snutch TP, Loman NJ, Quick J. Improvements to
752 the ARTIC multiplex PCR method for SARS-CoV-2 genome sequencing using nanopore. *bioRxiv*. 2020 Sep;
753 <http://www.ncbi.nlm.nih.gov/pubmed/32908977>, doi: [10.1101/2020.09.04.283077](https://doi.org/10.1101/2020.09.04.283077), publisher: Cold Spring Har-
754 bor Laboratory.
- 755 **Wang L**, Didelot X, Yang J, Wong G, Shi Y, Liu W, Gao GF, Bi Y. Inference of person-to-person transmission of
756 COVID-19 reveals hidden super-spreading events during the early outbreak phase. *Nature communications*.
757 2020 Oct; 11(1):5006. <http://www.ncbi.nlm.nih.gov/pubmed/33024095>, doi: [10.1038/s41467-020-18836-4](https://doi.org/10.1038/s41467-020-18836-4),
758 publisher: Nature Publishing Group.
- 759 **Worby CJ**, Lipsitch M, Hanage WP. Within-Host Bacterial Diversity Hinders Accurate Reconstruction of Transmis-
760 sion Networks from Genomic Distance Data. *PLoS Computational Biology*. 2014 Mar; 10(3):e1003549. [https:](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1003549)
761 [/journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1003549](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1003549), doi: [10.1371/journal.pcbi.1003549](https://doi.org/10.1371/journal.pcbi.1003549),
762 publisher: Public Library of Science.

- 763 **Wymant C**, Hall M, Ratmann O, Bonsall D, Golubchik T, de Cesare M, Gall A, Cornelissen M, Fraser C, STOP-
764 HCV Consortium TMPC and The BEEHIVE Collaboration. PHYLOSCANNER: Inferring Transmission from
765 Within- and Between-Host Pathogen Genetic Diversity. *Molecular Biology and Evolution*. 2018 Mar; 35(3):719-
766 733. <https://doi.org/10.1093/molbev/msx304>, doi: 10.1093/molbev/msx304.
- 767 **Xin H**, Li Y, Wu P, Li Z, Lau EHY, Qin Y, Wang L, Cowling BJ, Tsang TK, Li Z. Estimating the Latent Period of
768 Coronavirus Disease 2019 (COVID-19). *Clinical Infectious Diseases*. 2022 May; 74(9):1678-1681. <https://doi.org/10.1093/cid/ciab746>, doi: 10.1093/cid/ciab746.
- 770 **Zwart MP**, Elena SF. Matters of Size: Genetic Bottlenecks in Virus Infection and Their Potential Impact on
771 Evolution. *Annual Review of Virology*. 2015 Nov; 2:161-179. <https://www.annualreviews.org/doi/abs/10.1146/annurev-virology-100114-055135>, doi: 10.1146/annurev-virology-100114-055135, publisher: Annual Reviews.

773 **Appendix Figures and Tables**

Supplementary file 1. Study participants metadata.

Supplementary file 2. Sample collection and demographics.

Supplementary file 3. Metrics used for phylogenetic tree comparison.

Supplementary file 4. Transition/transversion rates and base frequencies of the known simulated tree.

Supplementary file 5. Inferred transition/transversion rates and base frequencies when using the consensus sequence. Numbers show the average of 100 simulations.

Supplementary file 6. Inferred transition/transversion rates and base frequencies when accounting for within-host diversity. Numbers show the average of 100 simulations

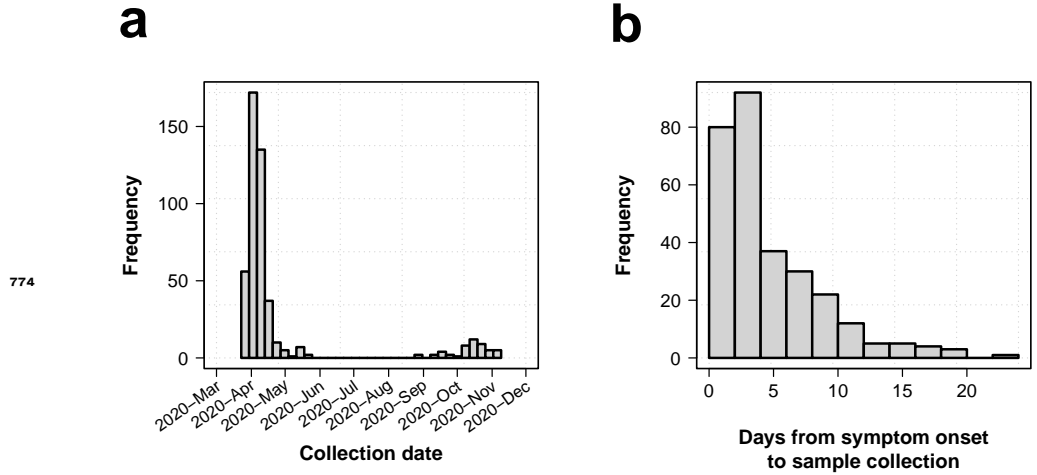


Figure 1—figure supplement 1. Collection date distribution and time from symptom and days from symptom onset.

(a) Distribution of collection dates. **(b)** Histogram of time from symptom onset to sample collection.

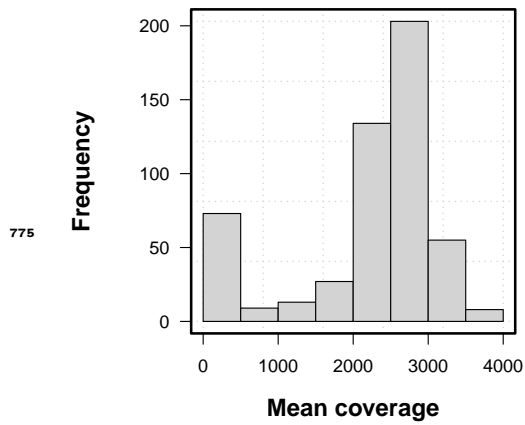


Figure 1—figure supplement 2. Sample mean coverage distribution.

Density distribution of mean coverage.

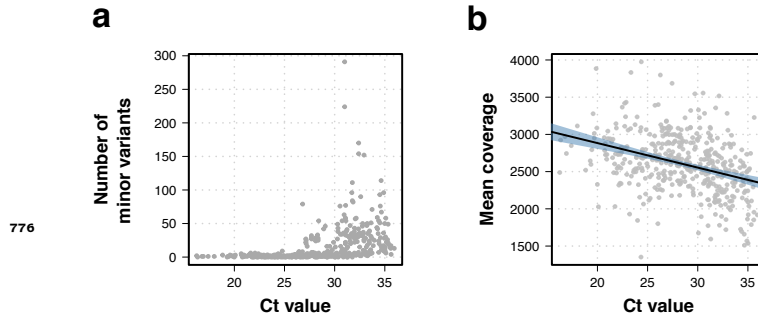


Figure 1—figure supplement 3. Effects of C_t value on whole-genome sequencing data.
a Higher C_t values were linked to a higher number of within-sample variation. **b** Correlation between C_t value and isolate sequencing mean coverage. Lower coverage was associated to higher C_t values ($R^2 = 0.13$, t-statistic p-value < 0.001).

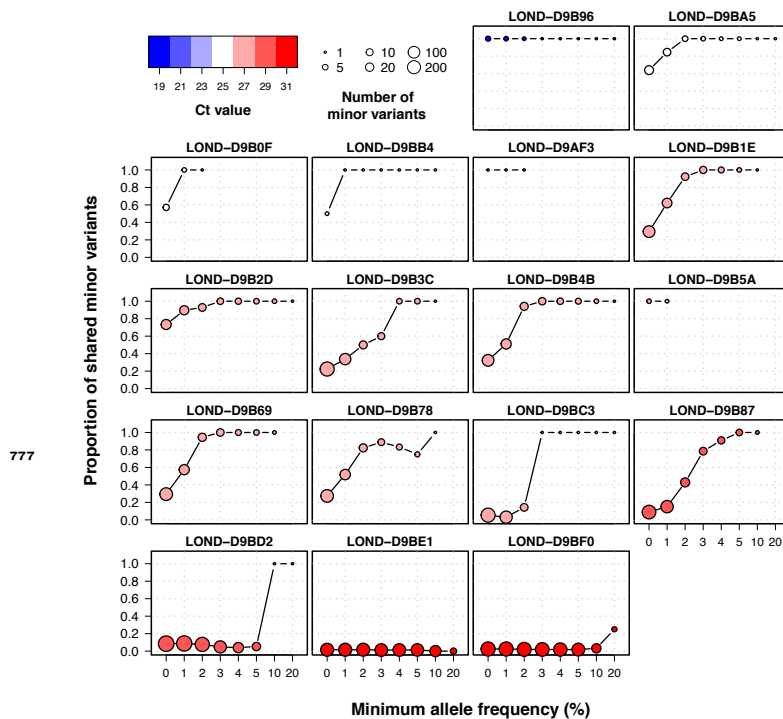


Figure 1—figure supplement 4. Proportion of shared minor variants between technical replicates using different filters of allele frequency.

Individual plots of shared within-host variants between technical duplicates using increasing thresholds of allele frequency. Colors represent C_t value, while the size of the point shows the total number of within-host variants between the two samples.

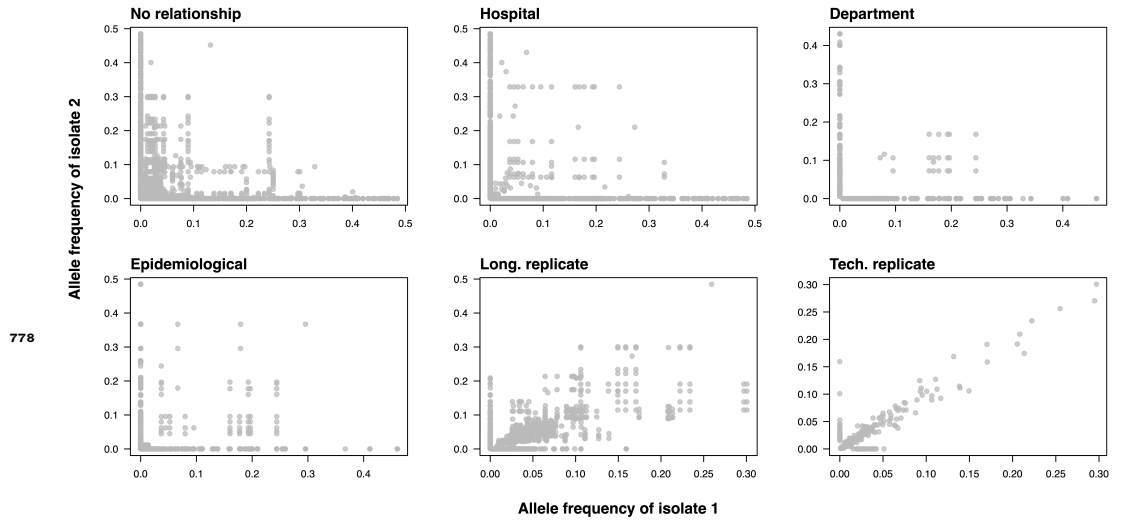


Figure 2—figure supplement 1. Allele frequency comparison in pairwise sample pairs. Pairwise allele frequency comparison between isolate pairs with different relationships. Allele frequencies were compared between isolates with no relationship, from the same hospital, from the same department, with epidemiological links, as well as between longitudinal and technical replicates.

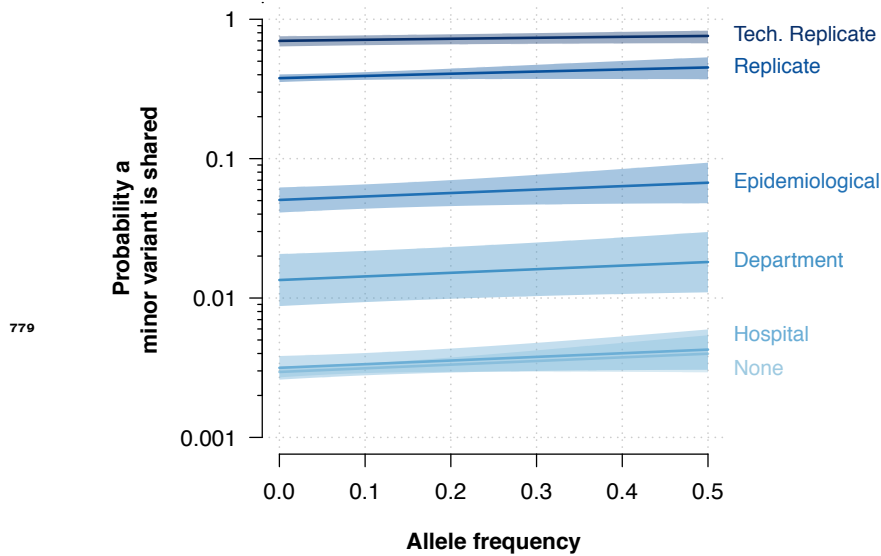
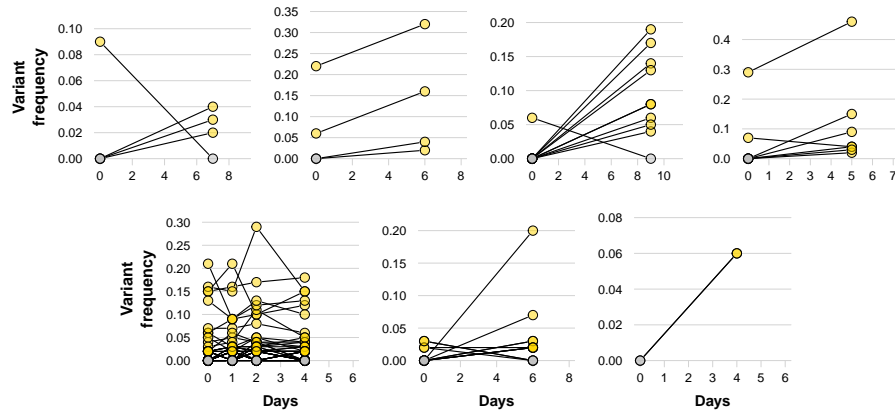


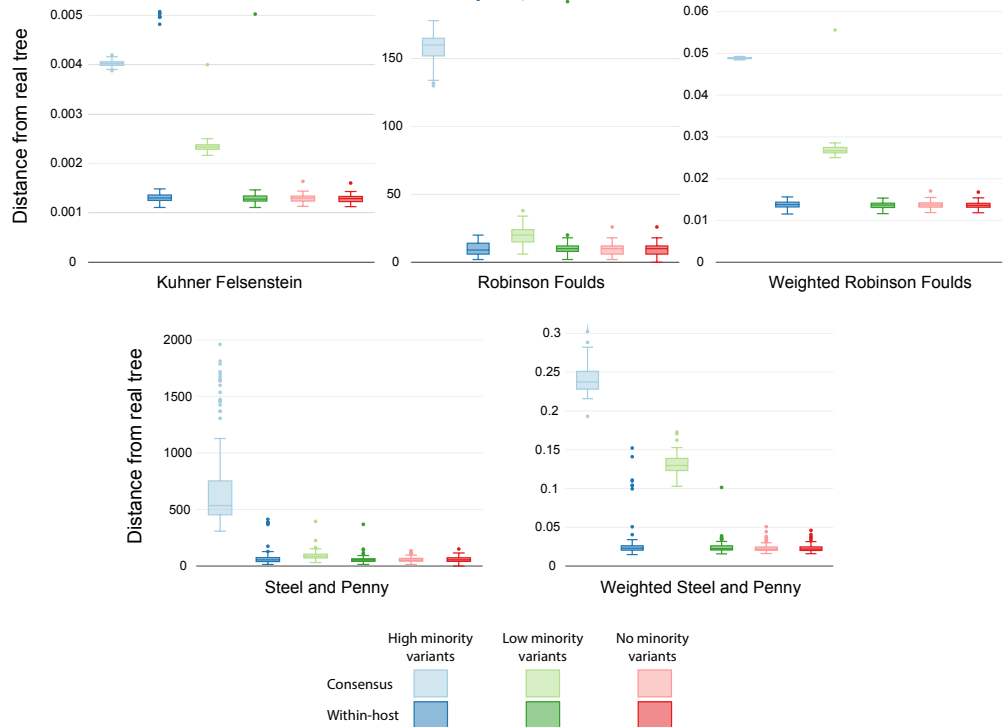
Figure 2—figure supplement 2. Probability that minor variants are shared. Probability that low frequency variants are shared inferred with a logistic model with allele frequency and epidemiological relationship as independent variable and whether a variant is shared or not as dependent variable. Y-axis in logarithmic scale for representation.



780

Figure 2—figure supplement 3. Dynamics of low frequency variants in longitudinal duplicates.

Variant frequency of low frequency variants through time in longitudinal duplicates. Each panel represents a single individual, with variants indicated by dots at each time point. The same variant at different time points is linked by lines. Yellow colors represent variants that are consistently found at each time point, while grey dots show variants that present in the first sampling event but lost in subsequent isolates.



781

Figure 4—figure supplement 1. Similarity scores for inferred trees with different rates.

Comparison of the phylogenetic trees inferred using simulated sequences with different transition/transversion rates to reflect different within-host diversity levels. Colors show the different rates of within-host evolution. Light colors represent trees inferred with consensus alignments, while dark colors show trees inferred with the model accounting for within-host diversity.

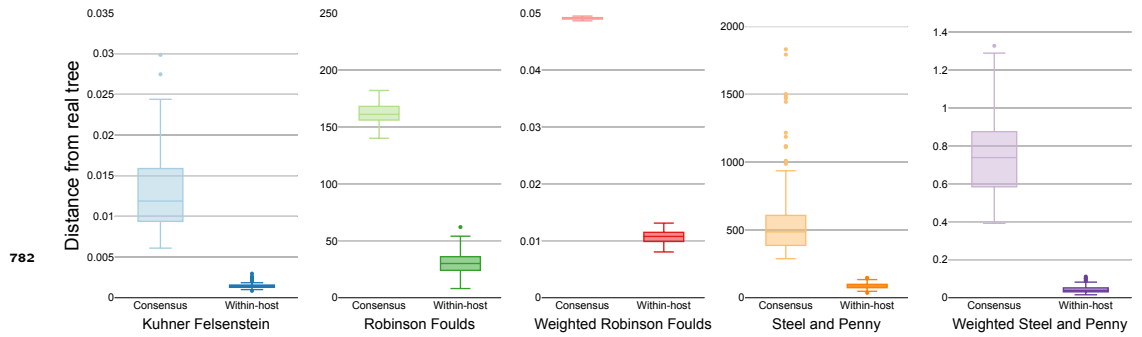


Figure 4—figure supplement 2. Similarity scores for inferred trees from coalescent simulations.

Comparison of the phylogenetic trees inferred using simulated sequences from known coalescent starting trees and different phylogenetic models. Colors differentiate the metrics used for the comparison.

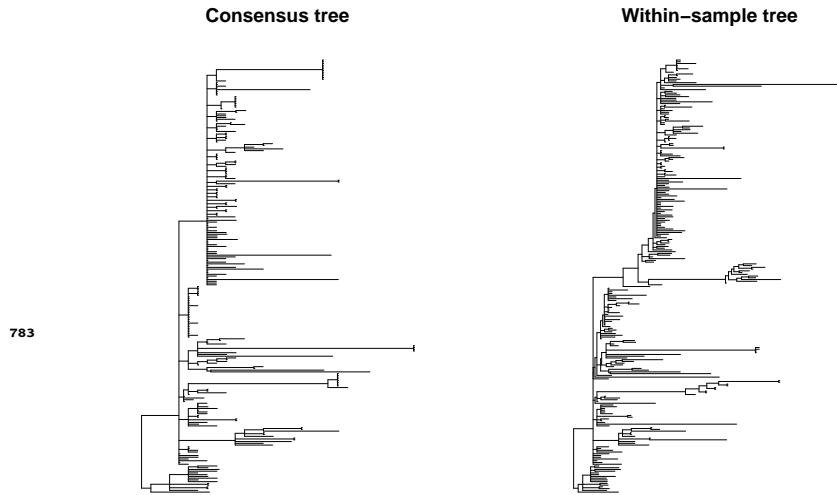


Figure 6—figure supplement 1. Phylogenetic trees for SARS-CoV-2.

SARS-CoV-2 phylogenetic trees inferred from consensus sequences (left) and an alignment with major and minor variant information (right).

784

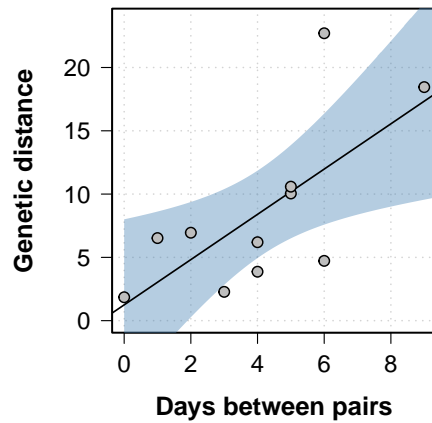


Figure 6—figure supplement 2. Genetic distance between longitudinal samples.

The genetic distance in the phylogenetic tree inferred using within-sample diversity increased as the between longitudinal samples progressed. Black line shows the best fit in a linear model, while the blue shaded area represents the 95% CI.

785

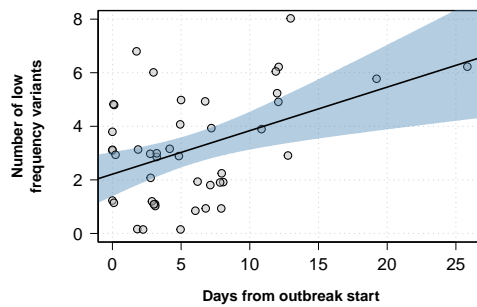


Figure 6—figure supplement 3. Number of low frequency variants within outbreaks as the outbreak progresses.

Y-axis shows the number of low frequency variants for each isolate within an outbreak, while the x-axis represents the days since that particular outbreak started. Black line shows the best fit in a linear model, while the blue shaded area represents the 95% CI.

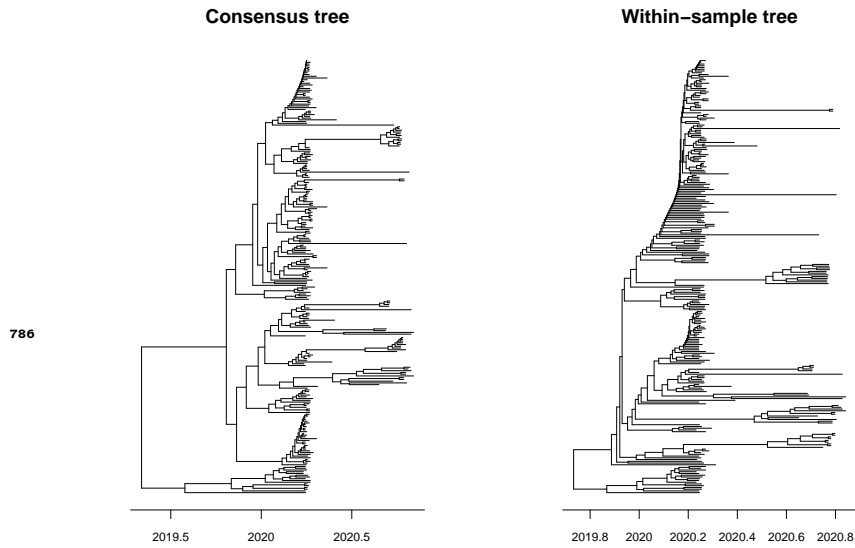


Figure 7—figure supplement 1. Time calibrated phylogenetic trees for SARS-CoV-2. SARS-CoV-2 phylogenetic trees inferred from consensus sequences (left) and an alignment with major and minor variant information (right). Branch lengths are measured in years.

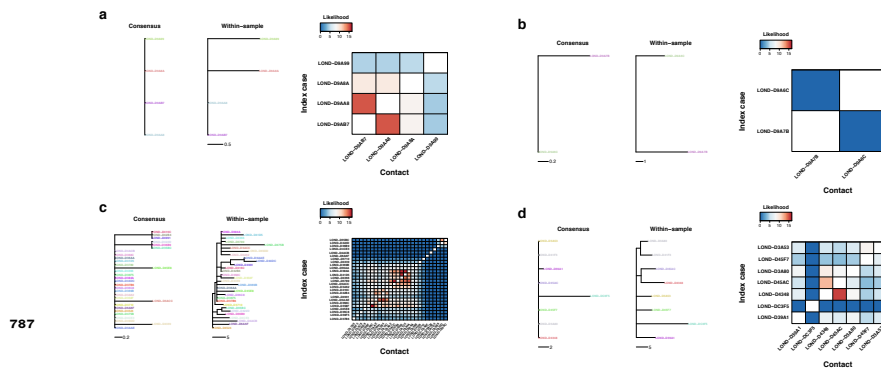


Figure 7—figure supplement 2. Phylogenetic and transmission for SARS-CoV-2 outbreaks. **a-d** Phylogenies of SARS-CoV-2 outbreaks. The branch lengths are in units of substitutions per genome, and the scales are shown under the trees. Colors represent samples from the same individual. Samples with the same name are technical replicates. Left tree of each panel shows the phylogeny inferred with the consensus alignment. Right tree represents the phylogeny inferred using within-sample variation. Heatmap shows the likelihood of direct transmission for each pair of samples in a SEIR model of transmission. Vertical axis is the infector while the horizontal axis shows the infectee.