

1 **Full title:** Shared within-host SARS-CoV-2 variation in households

2 **Running head:** Within-host SARS-CoV-2 variation

3

4 Katharine S. Walter¹, Eugene Kim¹, Renu Verma¹, Jonathan Altamirano², Sean Leary³, Yuan J.

5 Carrington³, Prasanna Jagannathan^{1,4}, Upinder Singh^{1,4}, Marisa Holubar¹, Aruna Subramanian¹, Chaitan

6 Khosla^{5,6}, Yvonne Maldonado^{2,3}, Jason R. Andrews¹

7 1. Division of Infectious Diseases and Geographic Medicine, Stanford University School of
8 Medicine, Stanford, CA 94305, USA

9 2. Department of Epidemiology and Population Health, Stanford University School of
10 Medicine, Stanford, CA 94305, USA

11 3. Department of Pediatrics, Stanford University School of Medicine, Stanford, CA, USA

12 4. Department of Microbiology and Immunology, Stanford University School of Medicine,
13 Stanford, CA 94305, USA

14 5. Stanford ChEM-H, Stanford University, Stanford CA 94305, USA

15 6. Departments of Chemistry and Chemical Engineering, Stanford University, Stanford, CA
16 94305, USA

17 **Correspondence:**

18 Katharine S. Walter, kwalter@stanford.edu

19 Division of Infectious Diseases and Geographic Medicine

20 Stanford University School of Medicine

21 300 Pasteur Drive, Lane Building 145

22 Stanford, CA 94305

23 **Article type:** Major Article

24 **Abstract word count:** 194

25 **Text word count:** 3494

26 **Abstract**

27 **Background:** The limited variation observed among SARS-CoV-2 consensus sequences makes it
28 difficult to reconstruct transmission linkages in outbreak settings. Previous studies have recovered
29 variation within individual SARS-CoV-2 infections but have not yet measured the informativeness of
30 within-host variation for transmission inference.

31 **Methods:** We performed tiled amplicon sequencing on 307 SARS-CoV-2 samples from four prospective
32 studies and combined sequence data with household membership data, a proxy for transmission linkage.

33 **Results:** Consensus sequences from households had limited diversity (mean pairwise distance, 3.06
34 SNPs; range, 0-40). Most (83.1%, 255/307) samples harbored at least one intrahost single nucleotide
35 variant (iSNV; median: 117; IQR: 17-208), when applying a liberal minor allele frequency of 0.5% and
36 prior to filtering. A mean of 15.4% of within-host iSNVs were recovered one day later. Pairs in the same
37 household shared significantly more iSNVs (mean: 1.20 iSNVs; 95% CI: 1.02-1.39) than did pairs in
38 different households infected with the same viral clade (mean: 0.31 iSNVs; 95% CI: 0.28-0.34), a signal
39 that increases with increasingly liberal thresholds.

40 **Conclusions:** Although only a subset of within-host variation is consistently shared across likely
41 transmission pairs, shared iSNVs may augment the information in consensus sequences for predicting
42 transmission linkages.

43

44 **Keywords:** SARS-CoV-2, transmission, viral evolution, within-host diversity

45

46

47 **Background**

48 SARS-CoV-2 genomic sequencing has been powerfully used to reconstruct the virus'
49 evolutionary dynamics at broad temporal and spatial scales[1–3]. Yet the virus' relatively slow
50 substitution rate compared with its short serial interval limits the viral diversity observed in many
51 outbreaks, and viral consensus sequences—which represent the most common allele along the viral
52 genome—are often identical or nearly so [4,5].

53 In superspreading events, identical consensus sequences have provided important evidence of
54 recent shared transmission. For example, four individuals on the same international flight were infected
55 with identical SARS-CoV-2 consensus genomes, evidence that the virus could be transmitted during air
56 travel[6]. Genomic surveillance in Boston during 2020 reported that 59 out of 83 (71%) genomes
57 sequenced from a skilled nursing facility were identical, implicating transmission within the facility[7].
58 Similarly, 75% of SARS-CoV-2 consensus sequences from a fishing boat outbreak were identical to at
59 least one other sequence, and the remaining sequences were closely related, suggesting rapid transmission
60 from a single viral introduction [8].

61 While consensus sequences have been harnessed to implicate or exclude the possibility of a
62 shared recent transmission history, closely related consensus sequences are often not sufficient for
63 reconstructing transmission linkages. For example, in hospital-based surveillance in Wisconsin, many
64 healthcare workers implicated in different epidemiological clusters were infected with identical SARS-
65 CoV-2 genomes[9]. In a hospital outbreak in Portugal, groups of identical consensus sequences shared
66 between healthcare workers and patients were frequently identified[5]. Similarly, in hospital-based
67 surveillance in the UK, 159 of 299 (53%) genomes sampled from one hospital were identical to at least
68 one other sampled genome[10]. While many pairs of individuals infected with identical genomes had
69 strong or intermediate evidence of transmission, 22% had no epidemiological evidence of
70 transmission[10], potentially the result of incomplete epidemiological information (cryptic transmission)
71 or limited genomic variation resulting in identical, epidemiologically unlinked consensus genomes. In the
72 absence of detailed epidemiological data, such as contact information or spatial information that might be

73 available in hospital-based studies, it is not yet known whether routine sequencing data alone can be used
74 to reconstruct transmission linkages of who-infected-whom or identify locations or individuals that may
75 drive transmission.

76 Genomic studies of HIV and other viral and bacterial pathogens have begun to harness the
77 pathogen variation within individual infections, or within-host diversity, to reconstruct transmission
78 linkages [11–13]. Previous studies have reported low levels of SARS-CoV-2 diversity within individual
79 hosts and have estimated the size of a narrow transmission bottleneck which limits the viral diversity
80 shared across hosts[8,9,14,15]. However, more research is needed to quantify the informativeness of
81 within-host SARS-CoV-2 variation and evaluate the effects of variant identification approaches on
82 transmission inferences[15].

83 To investigate the potential for within-host SARS-CoV-2 diversity to be harnessed for studies of
84 transmission, we deep sequenced SARS-CoV-2 samples collected from household members, allowing us
85 to directly compare shared within-host variants among epidemiologically linked individuals and those
86 with no known linkage, providing a test case for the transmission information contained within individual
87 infections. We additionally sequenced artificial mixtures of SARS-CoV-2 variants to examine tradeoffs
88 between sensitivity and specificity in within-host variation identification.

89

90 **Methods**

91 **Collection of residual SARS-CoV-2 samples for deep sequencing.**

92 We assembled a collection of samples from four prospective SARS-CoV-2 research studies: (a) a
93 prospective household transmission study, in which index cases with at least one reverse transcription
94 quantitative polymerase chain reaction-confirmed (RT-qPCR) SARS-CoV-2 test were enrolled along with
95 household members. Participants were tested daily for SARS-CoV-2 RNA via RT-qPCR, using self-
96 collected lower nasal swabs, and households were followed until all members tested negative for seven
97 consecutive days[16]. (b) A randomized, single-blind, placebo-controlled trial of Peginterferon Lambda-
98 1a (Lambda) for reducing the duration of viral shedding or symptoms[17] in which oropharyngeal swabs

99 were collected for 28 days following enrollment. (c) A phase 2 double-blind randomized controlled
100 outpatient trial of the antiviral favipiravir for reducing the duration of viral shedding in which participants
101 self-collected daily anterior nasal swabs for 28 days following enrollment[18]. Neither Lambda nor
102 favipiravir was found to shorten the duration of SARS-CoV-2 viral shedding[17,18]. (d) A study of a
103 noninvasive mask sampling method to quantify SARS-CoV-2 shedding in exhaled breath[19].

104 All study participants provided written consent and all studies were approved by the Institutional
105 Review Board of Stanford University (Numbers: 55479, 57686, 56032, and 55619). We identified
106 household members through participation in the household transmission study and address matching.

107

108 **Sample RT-qPCR testing**

109 We collected nasal swabs in 500 μ L of Primestore MTM (Longhorn Vaccines & Diagnostics)
110 RNA-stabilizing media. For exhaled breath samples, we extracted RNA from gelatin membrane filters
111 processed in 1-mL PrimeStore MTM media. RNA was extracted using the MagMAX Viral/Pathogen
112 Ultra Nucleic Acid Isolation Kit (Cat # A42356 Applied Biosystems) and eluted in 50 μ L of elution
113 buffer [15] (Supplementary Methods).

114

115 **Library preparation and sequencing**

116 We followed the ARTIC v3 Illumina library preparation and sequencing protocols[20] and
117 sequenced amplicons on an Illumina MiSeq platform (Supplementary Methods). Sequence data is
118 available at SRA (BioProject ID: PRJNA842503).

119

120 **Variant identification.**

121 We used the *nf-core/viralrecon v.2.4* bioinformatic pipeline to perform variant calling and
122 generate consensus sequences from raw sequencing reads[21]. Briefly, we aligned reads to the
123 MN908947.3 reference genome with *Bowtie 2*[22], removed primer sequences with *iVar*[23], and called
124 variants with respect to the reference genome with *iVar*[23]. We also used this pipeline to remove reads

125 mapping to the host genome with *Kraken2*[24], map reads with *Bowtie 2*[22], generate consensus
126 sequences with *bcftools*[25], and assign Nextclade lineages[26]. We modified the pipeline to include
127 variants with an alternate allele frequency $\geq 0.2\%$.

128 We included samples with a median coverage of 100X and with $>70\%$ of the genome covered by
129 a depth of $>10X$. We focused our analysis on single nucleotide polymorphism (SNP) variants and
130 excluded SNPs occurring at previously reported problematic sites[27].

131 To test whether commonly applied filters would improve overall accuracy, we applied five
132 variant filters: a filter for iSNV quality from *iVar*[23] (PASS = TRUE), a variant quality score filter
133 (Phred score >40), a depth filter (of both major and minor alleles $> 5X$), a filter of false positive iSNVs
134 repeated in more than one sample in the artificial strain mixture experiment (below), and all filters. We
135 additionally excluded iSNVs occurring in primer binding sites (except for the unfiltered variant set).

136 To identify shared within-host diversity across samples, we compared each unique pair of
137 samples meeting our quality criteria. We identified shared iSNVs as shared variant positions in which a
138 variant was not fixed and with the same alternate allele call. We additionally determined the geometric
139 mean of the sum of minor allele frequencies at shared iSNVs for each sample in a pair as a measure of
140 shared viral population diversity. To exclude potential shared iSNVs attributable to sequencing batch, we
141 excluded samples sequenced on the same Illumina sequencing lane in pairwise comparisons.

142

143 **Statistical analysis.**

144 We fit a Poisson regression model for the number of iSNVs identified within a single sample
145 including sequencing batch and participant as random effects. We additionally fit a Poisson regression
146 model for the number of pairwise shared iSNVs as a function of pair type and distance between consensus
147 sequences, including pair as a random effect. Finally, we fit a binomial regression model for predicting
148 household membership as a function of the number of shared pairwise iSNVs and an indicator variable
149 for close consensus sequences (pairwise distance ≤ 1 SNP), including the earliest samples collected from
150 each pair to exclude multiple pairwise comparisons. We fit all models with the R package *lme4*[31], and

151 included the set of variants after applying all filters, including iSNVs with a minor allele frequency of
152 $\geq 0.2\%$. We excluded samples sequenced in the same sequencing batch.

153 **Replicating analysis in an independent deep sequencing dataset from Wisconsin.**

154 We additionally investigated patterns of shared within-host variation in a previously published
155 dataset from a household transmission study in Wisconsin[9]. Specifically, we re-analyzed variants called
156 by the previous study and filtered to include iSNVs with a minor allele frequency $\geq 1\%$ and to exclude
157 variants occurring at primer binding sites[9].

158 Samples in the previous study were sequenced in duplicate. To generate a set of variant calls that
159 were comparable to those from our California dataset, we took the union of iSNVs identified in each
160 replicate sample; for iSNVs detected in both samples, we included the iSNV with the greater minor allele
161 frequency. As in the previous study, we excluded iSNVs called in genomic positions < 54 or $> 29,837$ or at
162 position 6669, which was identified as a problematic site. As above, we excluded positions previously
163 reported as problematic sites[27] from variant calls.

164

165 **Results**

166 **Assembling a collection of longitudinally sampled individuals and transmission pairs.**

167 We aggregated residual nasal swabs from four studies collected from March 2020 through May
168 2021 and deep sequenced SARS-CoV-2 genomes using the ARTIC v3 tiled amplicon sequencing
169 protocol (Fig. 2a). 307 SARS-CoV-2 sequences from 286 unique biological samples from 135
170 participants met our quality and coverage filters, including 130 samples from 32 individuals in 14
171 households and 57 longitudinally sampled individuals.

172 Samples had a median coverage depth of 1714 reads with a median of 99.0% of the SARS-CoV-2
173 genome covered by at least 10 reads. As expected, coverage depth was inversely correlated with RT-
174 qPCR cycle threshold (Pearson's $r = -0.15$, $p = 0.012$), reflecting a positive correlation with SARS-CoV-2
175 burden.

176 Samples were distributed across many of the major SARS-CoV-2 lineages circulating at the time
177 of collection (Fig. 2b). Overall, consensus sequences had a mean pairwise distance of 37.4 fixed SNPs
178 (range, 0-76) (Fig. 2b). 193 of 456 pairs of consensus sequences (42.3%) sampled longitudinally from the
179 same individual differed by 0-1 SNP (mean pairwise distance, 2.37; range, 0-22). The single individual
180 with consensus sequences that differed by 22 SNPs had been a participant in a SARS-CoV-2 clinical trial
181 and had received the antiviral drug favipiravir[18]; samples were taken 10 days apart. 251 of 778 (32.3%)
182 of pairs of consensus sequences sampled from different individuals in the same household differed by 0-1
183 SNP (mean pairwise distance, 3.06; range, 0-40), consistent with the relatively slow SARS-CoV-2
184 substitution rate[4]. In contrast, only a small minority, 20 of 41,053 (0.49%) pairs of consensus sequences
185 sampled from different households were within 0-1 SNP (mean pairwise distance, 38.52; range 0-76).

186

187 **A subset of within-host diversity is consistently recovered over time.**

188 A major challenge in studies of within-host pathogen diversity is in distinguishing true, low
189 frequency intrahost nucleotide variants (iSNVs) from sequencing or bioinformatic errors[32]. By
190 sequencing artificial strain mixtures of the Alpha and Beta variants, we established that we could reliably
191 recover minority variants to minor allele frequencies as low as 0.25% with 10^3 viral copies/mL (Fig. 1;
192 Supplementary Methods; Supplementary Text), with a minimal cost of false positive iSNVs (Fig. S1).

193 Most (83.1%, 255/307) samples harbored at least one iSNV (median: 7; IQR: 2-20) above a
194 minor allele frequency of 1.0% and applying all filters. As expected, the magnitude of recovered within-
195 host diversity increases to a median of 20 iSNVs (IQR:4-44) and 27 (IQR:6-55) with more liberal minor
196 allele frequency thresholds of 0.5% and 0.2% respectively (Fig. S2).

197 As previously reported[9,33,34], iSNVs are not consistently recovered within serial samples.
198 Among individuals with recovered within-host diversity, a mean of 8.7% within-host iSNVs above a
199 minor allele frequency of 1.0% and applying all filters were recovered one day later; this proportion
200 declined with time between samples, though not significantly ($r = -0.11$, $p = 0.23$). When including
201 unfiltered iSNVs above a 0.5% threshold, a mean of 15.4% of within-host iSNVs were recovered one day

202 later. The variable recovery of within-host variation is consistent with previous reports that minor allele
203 frequencies are poorly correlated within longitudinally-sampled individuals[35], potentially reflecting
204 both sampling or sequencing bottlenecks as well as a dynamic within-host viral population (Box 1).

205 Despite this, the pool of minority variants recovered within an individual is a consistent marker of
206 individual host. Pairs of sequences serially sampled from the same host consistently share 8.83 times
207 more iSNVs (mean 0.35 shared iSNVs; 95% CI: 0.22-0.48) compared to pairs of samples from different
208 individuals (mean 0.040 shared iSNVs; 95% CI: 0.037-0.042), after filtering and excluding samples
209 sequenced on the same batch (Figs. 3, S2, S3). The host-specific signature declines as an increasingly
210 strict minor allele frequency threshold is applied.

211 Within-host diversity, as measured by sample iSNV richness, was positively associated with PCR
212 Ct value, a measure of viral burden (aOR: 1.08; 95% CI: 1.08-1.10), in a general linear model, including
213 batch and participant as random effects. Among samples with symptom information, days following
214 symptom onset was not associated with increased iSNV richness in the multiple regression model (aOR:
215 1.01; 95% CI: 0.97-1.04) when controlling for Ct value.

216

217 **A signal of transmission linkage in within-host diversity**

218 We tested whether within-host SARS-CoV-2 diversity could be used to identify transmission
219 linkages using household membership as a proxy for probable epidemiological linkage. Pairs of
220 individuals in the same household shared significantly more iSNVs (at a 0.2% minor allele frequency
221 threshold, mean: 8.87 iSNVs; 95% CI: 7.81-9.92) than did pairs in different households infected with the
222 same viral clade (mean: 2.52 iSNVs; 95% CI: 2.32-2.73) or pairs in different households infected with a
223 different viral clade (mean: 2.92; 95% CI: 2.86-3.00), when including samples sequenced in different
224 batches, before filtering.

225 Applying different variant filtering approaches dramatically reduced the number of observed
226 iSNVs within individual samples and shared between sample pairs (Fig. S4) but did not eliminate the
227 signal of greater levels of shared within-host diversity among household pairs than among

228 epidemiologically unlinked pairs. After applying all filters, pairs of individuals in the same household
229 shared significantly more iSNVs (at a 0.2% minor allele frequency threshold, mean: 0.96 iSNVs; 95% CI:
230 0.78-1.15) than do pairs in different households infected with the same viral clade (mean: 0.20 iSNVs;
231 95% CI: 0.17-0.23) or pairs in different households infected with a different viral clade (mean: 0.25; 95%
232 CI: 0.24-0.26) (Fig. 3a) and including only samples sequenced in different batches.

233 Applying an increasingly stringent minor allele frequency threshold greatly reduces the number
234 of iSNVs observed within a sample and therefore the number of iSNVs shared between samples, yet
235 household members share more iSNVs than do epidemiologically unlinked participants. Applying a
236 minor allele frequency threshold of 1%, for example, a mean of 0.50 iSNVs (95% CI: 0.41-0.59) are
237 shared between pairs of samples collected from the same household and 0.073 iSNVs (95% CI: 0.060-
238 0.086) unfiltered iSNVs are shared between pairs in different households infected with the same viral
239 clade.

240 We hypothesized that minor allele frequencies of shared variants would contribute additional
241 information about transmission beyond the number of shared variant positions; we therefore measured
242 shared population-level diversity as the geometric mean of the sum of within-host minor allele
243 frequencies for shared iSNVs, which we refer to below as population diversity (Methods). Across all
244 minor allele frequency thresholds, pairs of individuals in the same household share significantly more
245 population diversity (at a 0.2% minor allele frequency threshold, after filtering, mean: 0.028; 95% CI:
246 0.021-0.034) than do pairs in different households infected with the same viral clade (mean: 0.0031; 95%
247 CI: 0.0025-0.0036) and different viral clades (mean: 0.0029; 95% CI: 0.0027-0.0030) (Fig. 3b).

248 In a generalized linear model for shared within-host diversity, household membership was
249 associated with an increased odds of shared iSNVs (aOR:19.8; 95% CI: 6.39-61.40) compared to sample
250 pairs within the same clade, after controlling for genetic distance between consensus sequences,
251 consistent with a previous study that found household membership is the strongest predictor of shared
252 iSNVs[9]. Longitudinal samples from an individual were also associated with an increased odds of shared
253 iSNVs (aOR: 60.5; 95% CI:19.9-184) as were sequencing replicates (aOR: 132; 95% CI: 42.5-411).

254 After excluding pairs sequenced in the same batch and multiple comparisons between
255 participants, our sample size was small (23 unique household pairs). In a generalized linear model, the
256 number of shared iSNVs was not significantly associated with an increased odds of household
257 membership (aOR: 1.31; 95% CI: 0.87-1.71), while a closely related consensus sequence (within 0-1
258 SNPs) was significantly associated with household membership (aOR: 30.38; 95% CI: 10.39-129.14).
259 However, shared diversity as measured as the standardized sum of shared minor allele frequencies
260 between pairs was associated with an increased odds of household membership (aOR: 1.20; 95% CI:
261 1.08-1.32) when controlling for closely related consensus sequence.

262

263 **Replication in a Wisconsin study**

264 We tested the replicability of our findings in an independent study conducted in Wisconsin where
265 SARS-CoV-2 was deep sequenced from 133 acutely-infected individuals, including members of 19
266 households[9]. At a frequency threshold of 0.5%, we found a similar signal that pairs of individuals in the
267 same household shared significantly more iSNVs (mean: 9.52 iSNVs; 95% CI 8.14-10.89) than did pairs
268 in different households infected with the same viral clade (mean: 4.28 iSNVs; 95% CI: 4.19-4.37) or pairs
269 in different households infected with a different viral clade (mean: 1.42; 95% CI: 1.38-1.47) (Fig. 3a), in
270 variants in filtered VCF files made publicly available from the earlier study[9] (Fig. S6; Methods). Our
271 findings were consistent across minor allele frequency thresholds, though as in the California data, a
272 signal of household membership was strongest when using minor allele frequency thresholds of $\leq 1\%$
273 (Fig. S6). We found a similar signal when measuring shared population diversity as the sum of shared
274 minor allele frequencies (Fig. S7). However, household pairs did not share significantly more diversity
275 than epidemiologically unrelated pairs when applying all filters and a minor allele frequency threshold
276 $\geq 3\%$ (Fig. S7).

277

278 **Discussion**

279 While most SARS-CoV-2 genomic studies focus on consensus sequences, consensus sequences
280 may not provide the resolution needed to reconstruct transmission linkages and identify potential sources
281 of transmission in outbreak settings, where many cases may be closely genetically related. Here, we
282 report that within-host SARS-CoV-2 genomic variation may contribute information about transmission
283 that may augment the information contained in viral consensus sequences. We focused on household
284 membership as a proxy for epidemiological linkage. However, the potential utility of within-host
285 variation would be for population surveillance or outbreak investigation, such as in hospitals or prisons,
286 where transmission linkages are not known *a priori*.

287 Our measures of within-host SARS-CoV-2 diversity are consistent with those measured in
288 previous studies when applying similar thresholds: a mean of three (range 0-5) iSNVs at a minor allele
289 frequency $\geq 2\%$ were identified in an outbreak on a fishing boat[8] and a mean of three iSNVs were
290 reported in individuals sampled in a household study in Wisconsin above a 3% minor allele frequency
291 threshold and consistent across sequencing replicates[9]. Further, our finding that iSNVs can be shared
292 between epidemiologically linked individuals is consistent with previous reports that household
293 membership is the most significant predictor of shared within-host variation[9]. Overall, SARS-CoV-2
294 within-host diversity is lower than that identified in other viral pathogens and, as previously reported, we
295 find that within-host viral diversity is frequently lost during transmission[9].

296 As others have reported[9,23,33,34], excluding sources of noise from within-host pathogen
297 genomic data remains a major challenge. We sequenced artificial strain mixtures of two SARS-CoV-2
298 variants of concern and found significant tradeoffs between sensitivity and specificity in recovery of true
299 within-host variants as increasingly strict variant filters were applied. Applying strict minor allele
300 frequency thresholds excludes much potential within-host variation. Additionally, in our empirical
301 sequencing data, we find that the signal of shared within-host variation across transmission pairs is
302 strongest when including iSNVs at low minor allele frequency thresholds.

303 The optimal variant identification approach may differ across applications—for example,
304 measurements of transmission bottleneck are highly sensitive to allele frequency threshold[9,36] and may

305 prioritize specificity, while studies of transmission might prioritize sensitivity to identify potential
306 transmission linkages. However, again as others have highlighted, our findings underscore the need to
307 control for other potential explanations for shared iSNVs while still prioritizing sensitivity (Box 2). Our
308 findings suggest that for transmission inference, privileging sensitivity in variant identification may
309 greatly improve sensitivity for recovering within-host variation, at a small cost of false positive variant
310 calls.

311 Our study has several limitations. First, we focused on a convenience sample of residual samples
312 with accompanying household information collected in California from March 2020 through May 2021.
313 Replicating these findings in other settings and with more recently emerged SARS-CoV-2 lineages is
314 critical to understand the generalizability of our findings. Second, our study focused on the potential
315 epidemiological value of within-host viral variation. Our focus was on transmission linkage rather than in
316 viral evolutionary dynamics or transmission bottlenecks, which might have different optimal variant
317 identification approaches. Third, many groups have hypothesized that evolution within immune-
318 compromised or immune-suppressed populations may be an important driver of the emergence of new
319 variants of concern or interest[37–41]. Our sample collection did not enable us to test these hypotheses.
320 Forth, the epidemiological utility of within-host variation depends on SARS-CoV-2 sampling and
321 sequencing. Routine sequencing may always not generate sufficient depth to accurately recover within-
322 host variation.

323 In conclusion, we find that SARS-CoV-2 variation within individual hosts may be shared across
324 transmission pairs and may contribute information on transmission linkage on a backdrop of limited
325 diversity among consensus sequences. More broadly, pathogen diversity within individual infections
326 holds largely untapped information that may enhance the resolution of transmission inferences.

327

328 **Box 1. Determinants of within-host SARS-CoV-2 diversity.** Potential contributors to recovered SARS-

329 CoV-2 diversity include biological determinants in addition to sampling methods.

330 • Biological determinants:

331 ○ True viral population diversity, including diversity present in the infecting inoculum and

332 diversity generated through both neutral and selective within-host processes, which in

333 turn may be driven by the host environment, host immune status and immune history

334 (including natural and vaccine-acquired immunity), and viral genotype.

335 ○ Viral population size within an individual, reflecting individual infection dynamics.

336 • Study design:

337 ○ Viral sampling technique and physical site of sampling may vary across studies.

338 ○ Sequencing approach including amplicon-based, metagenomic sequencing, or other

339 pathogen enrichment steps and sequencing platform often vary across studies.

340 ○ Sequencing depth of coverage.

341 ○ Prospective household transmission studies may enable infections to be identified and

342 sampled earlier compared to samples collected through passive surveillance.

343 • Bioinformatic choices:

344 ○ Read filtering, mapping and variant identification algorithms vary in sensitivity and
345 specificity.

346 ○ Some previous studies have required iSNVs to be identified in technical replicates[9,34].

347 ○ Minor allele frequency thresholds vary across studies, with previous studies applying
348 filters ranging from 2-6% [9,42].

349

350 **Box 2. Potential explanations for shared iSNVs.** As with the SARS-CoV-2 diversity present within

351 individuals, observed shared within-host diversity could be attributable to a biological signal or the

352 observation process.

- 353 • True positive: Transmission of a diverse infecting inoculum.
- 354 ○ Within-host viral diversity can be structured temporally[33,38,41] or spatially or both.
- 355 Transmitted diversity is a subset of diversity generated by within-host evolutionary
- 356 processes.
- 357 • False positive:
- 358 ○ Convergent or homoplastic iSNVs reflecting highly mutable sites along the genome or
- 359 sites under selection.
- 360 ○ Sequencing batch effects due to contamination or adapter switching during a sequencing
- 361 run.
- 362 ○ Artefacts of common sampling approach reflecting contamination due to similar
- 363 sampling or processing environment.
- 364 ○ Bioinformatic errors falling in consistent genomic regions that are difficult to map and/or
- 365 identify variants.
- 366
- 367
- 368

369 **Footnotes**

370

371 **Conflict of interest statement**

372

373 All authors declare no conflict of interest.

374

375 **Funding**

376 KSW received support from a Thrasher Early Career Award. Financial support from Stanford's

377 Innovative Medicines Accelerator and operational support from Stanford ChEM-H is acknowledged.

378

379

380 **Meetings where the information has previously been presented**

381

382 This information has previously presented to the California Department of Public Health COVIDNet

383 Expert Panel.

384 **Correspondence**

385 Katharine S. Walter

386 Division of Infectious Diseases and Geographic Medicine

387 Stanford University School of Medicine

388 300 Pasteur Drive, Lane Building 145

389 Stanford, CA 94305

390 kwalter@stanford.edu

391

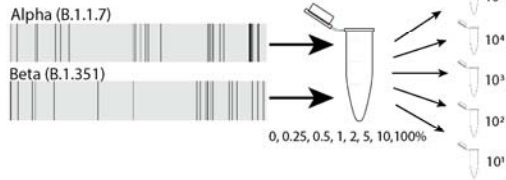
392 **Figures**

393 **Figure 1. Measuring the accuracy of within-host SARS-CoV-2 variant identification.**

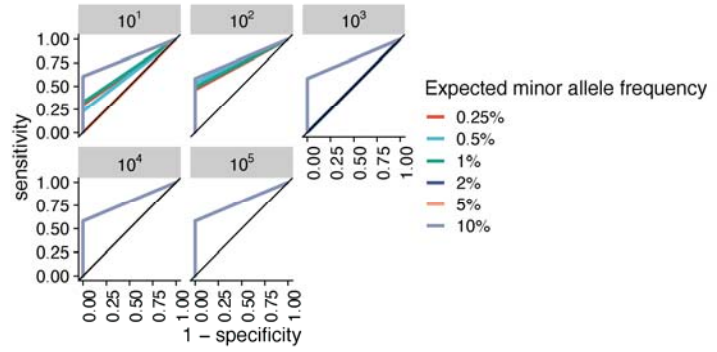
394 (a) Diagram of artificial strain mixture experiment. We conducted a serial dilution experiment, mixing
395 synthetic RNA controls (Twist Biosciences) of the Alpha (B.1.1.7) and Beta (B.1.351) variants so that the
396 minor variant comprised 0-10% of the source material and then serially diluted mixtures to a total of 10^1 -
397 10^5 total RNA copies. We conducted amplicon-based sequencing of artificial mixtures, identified iSNVs
398 with the viralrecon pipeline[21], and determined the sensitivity and specificity of our variant calling
399 pipeline to recovering true variation within our synthetic mixtures. (b) Receiver operator characteristic
400 curve showing 1-specificity versus sensitivity in the recovery of true minority variants, colored by total
401 RNA dilution. Lines corresponding to each dilution include results from five artificial strain mixtures,
402 including minority variants present at 0.5-10% of the total viral pool. (c) Observed minor allele frequency
403 versus expected minor allele frequency for artificial strain mixtures. Points indicate iSNV assignment,
404 (FP: false positive iSNV; TP; true positive iSNV). For FP iSNVs, point shape indicates whether FPs were
405 commonly repeated across samples (Common FP: FP identified in 10 or more samples; Other FP: any
406 other FP iSNV). Points were jittered for visualization. The grey dotted line indicates the minor allele
407 frequency threshold of 0.2%, below which the majority of FP iSNVs occur. Horizontal facets indicate the
408 synthetic RNA copy number in units of genome copies per microliter.

409
410

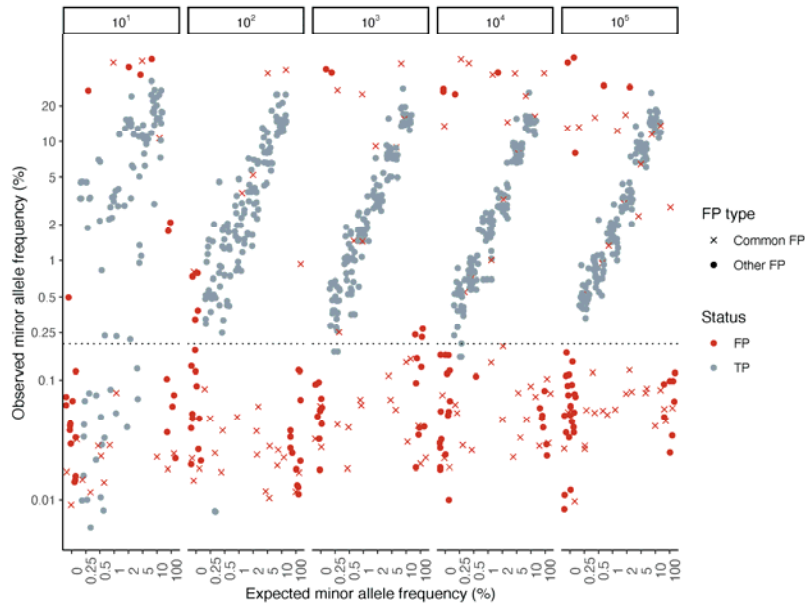
A



B



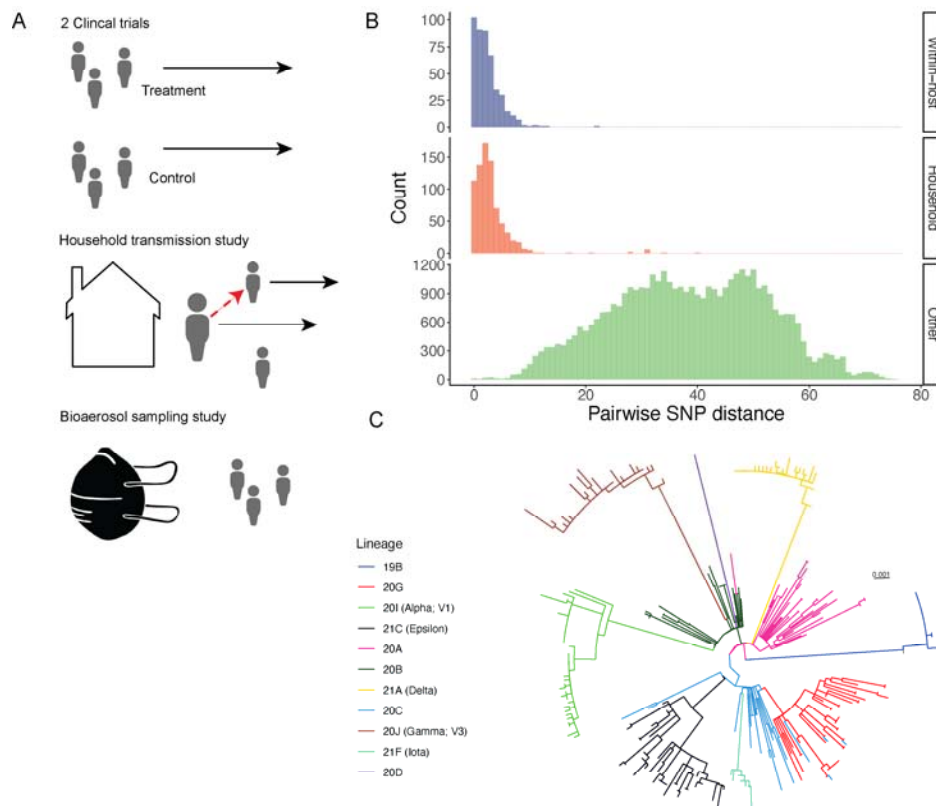
C



411
412

413 **Figure 2. Genetic diversity of sampled SARS-CoV-2.**

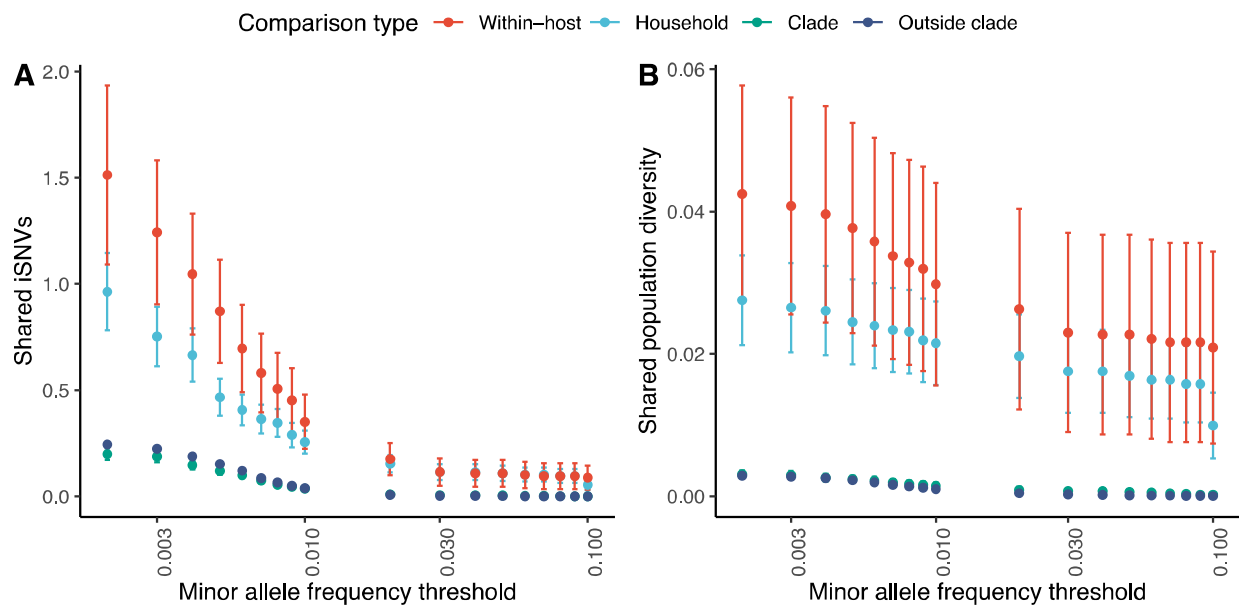
414 (a) We identified household transmission pairs and longitudinal samples by address matching from four
415 participants enrolled in studies including a household transmission study, clinical trials of Favipiravir and
416 Lambda, and a mask shedding study. (b) Histogram of pairwise single nucleotide polymorphism (SNP)
417 distances between consensus sequences from samples longitudinally sampled from the same individual,
418 from different individuals within the same household, and between individuals from different households.
419 (c) A maximum likelihood phylogeny inferred from consensus sequences with IQ-Tree with branches
420 colored by clade. Clade assignments are made with Nextclade[43] through the nf-core viralrecon
421 pipeline[44]. Branch lengths are in distances of substitutions per site.



422

423 **Figure 3. Shared within-host variants hold a signal of SARS-CoV-2 transmission.**

424 For our genomic collection from California, (a) pairwise comparisons of the number of shared iSNVs,
425 defined as a shared minor allele present at the same genomic position, identified across different minor
426 allele frequency thresholds, after applying all variant filters (Methods). Points and error bars indicate
427 mean and 95% confidence intervals and are colored by comparison type. Each pair is assigned to a unique
428 category. Within-host: pairs of samples from the same individual collected on different days; Household:
429 pairs of individuals from the same household; Clade: pairs of individuals outside households infected
430 with the same Nextclade clade; and Outside clade: pairs of individuals outside households infected with
431 different Nextclade clades. Pairwise comparisons include only samples sequenced in different sequencing
432 batches.



433

434

435

436

437 **References**

438

- 439 1. Turakhia Y, Thornlow B, Hinrichs AS, et al. Ultrafast Sample placement on Existing tRees
440 (USHER) enables real-time phylogenetics for the SARS-CoV-2 pandemic. *Nat Genet* [Internet].
441 [cited 2021 Jun 3]; . Available from: <https://doi.org/10.1038/s41588-021-00862-7>
- 442 2. Lemey P, Hong SL, Hill V, et al. Accommodating individual travel history and unsampled
443 diversity in Bayesian phylogeographic inference of SARS-CoV-2. *Nat Commun* 2020 111
444 [Internet]. Nature Publishing Group; 2020 [cited 2022 Apr 19]; 11(1):1–14. Available from:
445 <https://www.nature.com/articles/s41467-020-18877-9>
- 446 3. Korber B, Fischer WM, Gnanakaran S, et al. Tracking changes in SARS-CoV-2 Spike: evidence
447 that D614G increases infectivity of the COVID-19 virus. *Cell* [Internet]. Cell Press; 2020 [cited
448 2020 Jul 17]; . Available from: <https://doi.org/10.1016/j.cell.2020.06.043>.
- 449 4. Duchene S, Featherstone L, Haritopoulou-Sinanidou M, Rambaut A, Lemey P, Baele G. Temporal
450 signal and the phylodynamic threshold of SARS-CoV-2. *Virus Evol* [Internet]. Oxford Academic;
451 2020 [cited 2022 Mar 14]; 6(2). Available from:
452 <https://academic.oup.com/ve/article/6/2/veaa061/5894560>
- 453 5. Borges V, Isidro J, Macedo F, et al. Nosocomial Outbreak of SARS-CoV-2 in a “Non-COVID-19”
454 Hospital Ward: Virus Genome Sequencing as a Key Tool to Understand Cryptic Transmission.
455 *Viruses*. MDPI AG; 2021; 13(4).
- 456 6. Choi EM, Chu DKW, Cheng PKC, et al. In-Flight Transmission of SARS-CoV-2. *Emerg Infect*
457 *Dis* [Internet]. Centers for Disease Control and Prevention; 2020 [cited 2022 May 23];
458 26(11):2713. Available from: [/pmc/articles/PMC7588512/](https://pubmed.ncbi.nlm.nih.gov/3458512/)
- 459 7. Lemieux, Lemieux JE, Siddle KJ, et al. Phylogenetic analysis of SARS-CoV-2 in Boston
460 highlights the impact of superspreading events. *Science* (80-) [Internet]. 2020 [cited 2020 Dec
461 16]; :eabe3261. Available from: <https://www.sciencemag.org/lookup/doi/10.1126/science.abe3261>
- 462 8. Hannon WW, Roychoudhury P, Xie H, et al. Narrow transmission bottlenecks and limited within-

- 463 host viral diversity during a SARS-CoV-2 outbreak on a fishing boat. bioRxiv [Internet]. Cold
464 Spring Harbor Laboratory; **2022** [cited 2022 Feb 16]; :2022.02.09.479546. Available from:
465 <https://www.biorxiv.org/content/10.1101/2022.02.09.479546v1>
- 466 9. Braun KM, Moreno GK, Wagner C, et al. Acute SARS-CoV-2 infections harbor limited within-
467 host diversity and transmit via tight transmission bottlenecks. PLoS Pathog [Internet]. Public
468 Library of Science; **2021** [cited 2022 Feb 10]; 17(8):e1009849. Available from:
469 <https://journals.plos.org/plospathogens/article?id=10.1371/journal.ppat.1009849>
- 470 10. Meredith LW, Hamilton WL, Warne B, et al. Rapid implementation of SARS-CoV-2 sequencing
471 to investigate cases of health-care associated COVID-19: a prospective genomic surveillance
472 study. Lancet Infect Dis [Internet]. **2020** [cited 2020 Jul 17]; 20(11):1263–1272. Available from:
473 www.thelancet.com/infectionPublishedonline
- 474 11. Tonkin-Hill G, Ling C, Chaguzo C, et al. Pneumococcal within-host diversity during colonisation,
475 transmission and treatment. [cited 2022 Mar 7]; . Available from:
476 <https://doi.org/10.1101/2022.02.20.480002>
- 477 12. Wymant C, Hall M, Ratmann O, et al. PHYLOSCANNER: Inferring transmission from within-
478 and between-host pathogen genetic diversity. Mol Biol Evol [Internet]. Oxford University Press;
479 **2018** [cited 2020 Oct 1]; 35(3):719–733. Available from: <http://creativecommons.org/licenses/by-nc-nd/4.0/>.
- 480 13. Leitner T. Phylogenetics in HIV transmission: Taking within-host diversity into account. Curr
481 Opin HIV AIDS [Internet]. Lippincott Williams and Wilkins; **2019** [cited 2021 Jul 16]; 14(3):181–
482 187. Available from: [https://journals.lww.com/co-](https://journals.lww.com/co-hivandaids/Fulltext/2019/05000/Phylogenetics_in_HIV_transmission__taking.6.aspx)
483 [hivandaids/Fulltext/2019/05000/Phylogenetics_in_HIV_transmission__taking.6.aspx](https://journals.lww.com/co-hivandaids/Fulltext/2019/05000/Phylogenetics_in_HIV_transmission__taking.6.aspx)
- 484 14. Martin MA, Koelle K. Comment on “Genomic epidemiology of superspreading events in Austria
485 reveals mutational dynamics and transmission properties of SARS-CoV-2.” Sci Transl Med
486 [Internet]. American Association for the Advancement of Science; **2021** [cited 2021 Nov 10];
487 13(617):1803. Available from: <https://www.science.org/doi/abs/10.1126/scitranslmed.abh1803>
- 488 15. San JE, Ngcapu S, Kanzi AM, et al. Transmission dynamics of SARS-CoV-2 within-host diversity

- 489 in two major hospital outbreaks in South Africa. *Virus Evol* [Internet]. Oxford Academic; **2021**
490 [cited 2022 May 23]; 7(1):41. Available from:
491 <https://academic.oup.com/ve/article/7/1/veab041/6248115>
- 492 16. Altamirano J, Govindarajan P, Blomkalns A, et al. 401. Natural History of Shedding and
493 Household Transmission of Severe Acute Respiratory Syndrome Coronavirus 2 Using Intensive
494 High-Resolution Sampling. *Open Forum Infect Dis*. Oxford University Press (OUP); **2021**;
495 8(Supplement_1):S302–S302.
- 496 17. Jagannathan P, Andrews JR, Bonilla H, et al. Peginterferon Lambda-1a for treatment of
497 outpatients with uncomplicated COVID-19: a randomized placebo-controlled trial. *Nat Commun*
498 [Internet]. **2021** [cited 2021 Apr 18]; 12(1). Available from: [https://doi.org/10.1038/s41467-021-](https://doi.org/10.1038/s41467-021-22177-1)
499 [22177-1](https://doi.org/10.1038/s41467-021-22177-1)
- 500 18. Holubar M, Subramanian A, Purington N, et al. Favipiravir for treatment of outpatients with
501 asymptomatic or uncomplicated COVID-19: a double-blind randomized, placebo-controlled,
502 phase 2 trial. *Clin Infect Dis* [Internet]. **2022** [cited 2022 May 18]; . Available from:
503 <https://academic.oup.com/cid/advance-article/doi/10.1093/cid/ciac312/6572081>
- 504 19. Verma R, Kim E, Degner N, Walter KS, Singh U, Andrews JR. Variation in SARS-CoV-2
505 bioaerosol production in exhaled breath. *Open Forum Infect Dis* [Internet]. **2021** [cited 2021 Dec
506 11]; . Available from: [https://academic.oup.com/ofid/advance-](https://academic.oup.com/ofid/advance-article/doi/10.1093/ofid/ofab600/6447624)
507 [article/doi/10.1093/ofid/ofab600/6447624](https://academic.oup.com/ofid/advance-article/doi/10.1093/ofid/ofab600/6447624)
- 508 20. Benjamin F, Diana R, Betteridge E, et al. COVID-19 ARTIC v3 Illumina library construction and
509 sequencing protocol V.5. Available from: <https://dx.doi.org/10.17504/protocols.io.bibtkann>
- 510 21. Ewels PA, Peltzer A, Fillinger S, et al. The nf-core framework for community-curated
511 bioinformatics pipelines [Internet]. *Nat. Biotechnol. Nature Research*; 2020 [cited 2021 Jun 4]. p.
512 276–278. Available from: <https://doi.org/10.1038/s41587-020-0446-y>.
- 513 22. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*. **2012**;
514 9(4):357–359.

- 515 23. Grubaugh ND, Gangavarapu K, Quick J, et al. An amplicon-based sequencing framework for
516 accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biol* [Internet].
517 **2019** [cited 2020 May 26]; 20(1). Available from: <https://doi.org/10.1186/s13059-018-1618-7>
- 518 24. Wood DE, Salzberg SL. Kraken: Ultrafast metagenomic sequence classification using exact
519 alignments. *Genome Biol* [Internet]. BioMed Central; **2014** [cited 2019 Apr 1]; 15(3):R46.
520 Available from: <http://genomebiology.biomedcentral.com/articles/10.1186/gb-2014-15-3-r46>
- 521 25. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and
522 population genetical parameter estimation from sequencing data. *Bioinformatics* [Internet]. **2011**
523 [cited 2019 Aug 11]; 27(21):2987–2993. Available from:
524 <http://www.ncbi.nlm.nih.gov/pubmed/21903627>
- 525 26. Rambaut A, Holmes EC, O’Toole Á, et al. A dynamic nomenclature proposal for SARS-CoV-2
526 lineages to assist genomic epidemiology. *Nat Microbiol* [Internet]. **2020** [cited 2021 Jun 4];
527 5(11):1403–1407. Available from: <https://www.nature.com/articles/s41564-020-0770-5.pdf>
- 528 27. Maio N De, Walker C, Borges R, Weilguny L, G S, Goldman N. Masking strategies for SARS-
529 CoV-2 alignments. 2020.
- 530 28. Nakamura T, Yamada KD, Tomii K, Katoh K. Parallelization of MAFFT for large-scale multiple
531 sequence alignments. *Bioinformatics* [Internet]. Oxford University Press; **2018** [cited 2020 Oct 8];
532 34(14):2490–2492. Available from: <https://mafft.cbrc.jp/alignment/software/mpi.html>.
- 533 29. Paradis E, Schliep K. Ape 5.0: An environment for modern phylogenetics and evolutionary
534 analyses in R. Schwartz R, editor. *Bioinformatics* [Internet]. Narnia; **2019** [cited 2019 Apr 10];
535 35(3):526–528. Available from:
536 <https://academic.oup.com/bioinformatics/article/35/3/526/5055127>
- 537 30. Minh BQ, Schmidt HA, Chernomor O, et al. IQ-TREE 2: New Models and Efficient Methods for
538 Phylogenetic Inference in the Genomic Era. *Mol Biol Evol* [Internet]. Oxford Academic; **2020**
539 [cited 2022 May 17]; 37(5):1530–1534. Available from:
540 <https://academic.oup.com/mbe/article/37/5/1530/5721363>

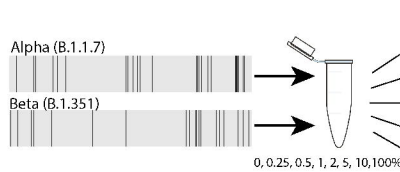
- 541 31. Bates D, Mächler M, Bolker BM, Walker SC. Fitting Linear Mixed-Effects Models Using lme4. J
542 Stat Softw [Internet]. American Statistical Association; **2015** [cited 2022 Apr 12]; 67(1):1–48.
543 Available from: <https://www.jstatsoft.org/index.php/jss/article/view/v067i01>
- 544 32. Mccrone JT, Lauring S. Measurements of Intrahost Viral Diversity Are Extremely Sensitive to
545 Systematic Errors in Variant Calling. J Virol. **2016**; 90(15):6884–6895.
- 546 33. Valesano AL, Rumfelt KE, Dimcheff DE, et al. Temporal dynamics of SARS-CoV-2 mutation
547 accumulation within and across infected hosts. Pekosz A, editor. PLOS Pathog [Internet]. Public
548 Library of Science; **2021** [cited 2021 Apr 16]; 17(4):e1009499. Available from:
549 <https://dx.plos.org/10.1371/journal.ppat.1009499>
- 550 34. Hannon WW, Roychoudhury P, Xie H, et al. Narrow transmission bottlenecks and limited within-
551 host viral diversity during a SARS-CoV-2 outbreak on a fishing boat. [cited 2022 Feb 16]; .
552 Available from: <https://doi.org/10.1101/2022.02.09.479546>
- 553 35. Lythgoe KA, Hall M, Ferretti L, et al. SARS-CoV-2 within-host diversity and transmission.
554 Science (80-) [Internet]. American Association for the Advancement of Science (AAAS); **2021**
555 [cited 2021 Apr 21]; 372(6539):eabg0821. Available from:
556 <https://doi.org/10.1126/science.abg0821>
- 557 36. Martin MA, Koelle K. Comment on “Genomic epidemiology of superspreading events in Austria
558 reveals mutational dynamics and transmission properties of SARS-CoV-2.” Sci Transl Med
559 [Internet]. American Association for the Advancement of Science; **2021** [cited 2021 Oct 28];
560 13(617):1803. Available from: <https://www.science.org/doi/10.1126/scitranslmed.abh1803>
- 561 37. Rambaut A, Loman N, Pybus O, et al. Preliminary genomic characterisation of an emergent
562 SARS-CoV-2 lineage in the UK defined by a novel set of spike mutations - SARS-CoV-2
563 coronavirus / nCoV-2019 Genomic Epidemiology - Virological [Internet]. Virological.org. 2020
564 [cited 2021 Feb 14]. Available from: [https://virological.org/t/preliminary-genomic-
565 characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-
566 mutations/563](https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563)

- 567 38. Kemp SA, Collier DA, Datir RP, et al. SARS-CoV-2 evolution during treatment of chronic
568 infection. *Nat* 2021 5927853 [Internet]. Nature Publishing Group; **2021** [cited 2022 May 2];
569 592(7853):277–282. Available from: <https://www.nature.com/articles/s41586-021-03291-y>
- 570 39. Weigang S, Fuchs J, Zimmer G, et al. Within-host evolution of SARS-CoV-2 in an
571 immunosuppressed COVID-19 patient as a source of immune escape variants. *Nat Commun* 2021
572 121 [Internet]. Nature Publishing Group; **2021** [cited 2022 May 11]; 12(1):1–12. Available from:
573 <https://www.nature.com/articles/s41467-021-26602-3>
- 574 40. Bessièrè P, Volmer R. From one to many: The within-host rise of viral variants. *PLOS Pathog*
575 [Internet]. Public Library of Science; **2021** [cited 2022 May 11]; 17(9):e1009811. Available from:
576 <https://journals.plos.org/plospathogens/article?id=10.1371/journal.ppat.1009811>
- 577 41. Choi B, Choudhary MC, Regan J, et al. Persistence and Evolution of SARS-CoV-2 in an
578 Immunocompromised Host. *N Engl J Med* [Internet]. Massachusetts Medical Society; **2020** [cited
579 2021 Jan 7]; . Available from: <https://www.nejm.org/doi/full/10.1056/NEJMc2031364>
- 580 42. Maio N De, Worby CJ, Wilson DJ, Stoesser N. Bayesian reconstruction of transmission within
581 outbreaks using genomic variants. Koelle K, editor. *PLoS Comput Biol* [Internet]. Public Library
582 of Science; **2018** [cited 2020 Oct 1]; 14(4):e1006117. Available from:
583 <https://dx.plos.org/10.1371/journal.pcbi.1006117>
- 584 43. Aksamentov I, Roemer C, Hodcroft EB, Neher RA. Nextclade: clade assignment, mutation calling
585 and quality control for viral genomes. *J Open Source Softw* [Internet]. The Open Journal; **2021**
586 [cited 2022 Apr 11]; 6(67):3773. Available from:
587 <https://joss.theoj.org/papers/10.21105/joss.03773>
- 588 44. Ewels PA, Peltzer A, Fillinger S, et al. The nf-core framework for community-curated
589 bioinformatics pipelines. *Nat Biotechnol* 2020 383 [Internet]. Nature Publishing Group; **2020**
590 [cited 2022 Apr 11]; 38(3):276–278. Available from: [https://www.nature.com/articles/s41587-](https://www.nature.com/articles/s41587-020-0439-x)
591 [020-0439-x](https://www.nature.com/articles/s41587-020-0439-x)
592

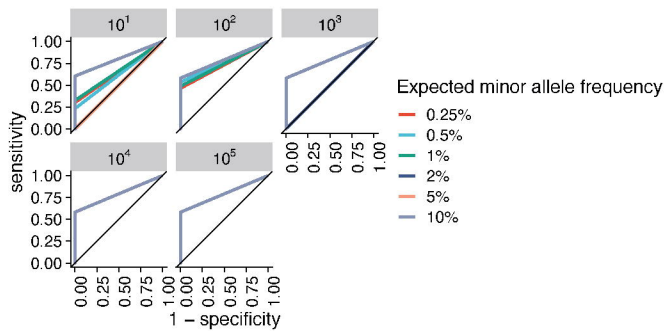
593

594

A



B



C

