

Temporal dynamics of SARS-CoV-2 mutation accumulation within and across infected hosts

Andrew L. Valesano^{1,2}, Kalee E. Rumfelt^{1,2}, Derek E. Dimcheff³, Christopher N. Blair^{1,2}, William J. Fitzsimmons^{1,2}, Joshua G. Petrie⁴, Emily T. Martin⁴, Adam S. Luring^{1,2} *

¹Division of Infectious Diseases, Department of Internal Medicine, University of Michigan, Ann Arbor, MI, USA; ²Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI, USA; ³Division of Hospital Medicine, Department of Internal Medicine, University of Michigan, Ann Arbor, MI, USA; ⁴Department of Epidemiology, University of Michigan, Ann Arbor, MI, USA

* Corresponding author

Adam S. Luring

1150 W. Medical Center Dr.

MSRB1 Room 5510B

Ann Arbor, MI 48109-5680

(734) 764-7731

aluring@med.umich.edu

1 **Abstract**

2 Analysis of SARS-CoV-2 genetic diversity within infected hosts can provide insight into the
3 generation and spread of new viral variants and may enable high resolution inference of
4 transmission chains. However, little is known about temporal aspects of SARS-CoV-2 intrahost
5 diversity and the extent to which shared diversity reflects convergent evolution as opposed to
6 transmission linkage. Here we use high depth of coverage sequencing to identify within-host
7 genetic variants in 325 specimens from hospitalized COVID-19 patients and infected employees
8 at a single medical center. We validated our variant calling by sequencing defined RNA mixtures
9 and identified a viral load threshold that minimizes false positives. By leveraging clinical
10 metadata, we found that intrahost diversity is low and does not vary by time from symptom
11 onset. This suggests that variants will only rarely rise to appreciable frequency prior to
12 transmission. Although there was generally little shared variation across the sequenced cohort,
13 we identified intrahost variants shared across individuals who were unlikely to be related by
14 transmission. These variants did not precede a rise in frequency in global consensus genomes,
15 suggesting that intrahost variants may have limited utility for predicting future lineages. These
16 results provide important context for sequence-based inference in SARS-CoV-2 evolution and
17 epidemiology.

18

19

20 Keywords: SARS-CoV-2, intrahost diversity, sequencing, transmission, evolution

21 **Introduction**

22 Over the course of the SARS-CoV-2 pandemic, whole genome sequencing has been widely
23 used to characterize patterns of broad geographic spread, transmission in local clusters, and
24 the spread of specific viral variants¹⁻⁶. Early reports demonstrated that SARS-CoV-2 exhibits
25 genetic diversity within infected hosts, but this has been less studied than consensus-level
26 genomic diversity⁷. Intra-host diversity is an important complement to consensus sequencing.
27 Patterns of viral intra-host diversity throughout individual infections can suggest the relative
28 importance of natural selection and stochastic genetic drift⁸. Shared intra-host variants between
29 individuals can reveal loci under convergent evolution and enable measurement of the
30 transmission bottleneck, a critical determining factor in the spread of new genetic variants^{9,10}.
31 Studies of SARS-CoV-2 intra-host diversity may shed light on selective pressures applied at the
32 individual level, such as antivirals and antibody-based therapeutics. While a clear understanding
33 of within-host evolution can inform how SARS-CoV-2 spreads on broader scales, there have
34 been relatively few comprehensive studies of intra-host dynamics^{9,11,12}.

35

36 Sequencing of intra-host populations can also potentially be applied to genomic epidemiology¹³.
37 A common goal in sequencing specimens from case clusters is to infer transmission linkage,
38 which can guide future public health and infection control interventions. However, the relatively
39 low substitution rate and genetic diversity of SARS-CoV-2 present challenges to inference of
40 individual transmission pairs^{13,14}. In the pandemic setting, there is a non-negligible chance that
41 two individuals who are epidemiologically unrelated could be infected with nearly identical viral
42 genomes. Viruses from a single local outbreak may have few differentiating substitutions,
43 limiting the ability of sequencing to resolve exact transmission chains. Identification of shared
44 intra-host variants between individuals has been explored in other pathogens to overcome this
45 obstacle¹⁵⁻¹⁹. However, use of this approach for SARS-CoV-2 will depend on a solid
46 understanding of the forces that shape the generation and spread of genetic variants.

47

48 There are several unresolved questions that will dictate the utility of intrahost diversity for
49 genomic epidemiology. First, there must be sufficient intrahost diversity generated during acute
50 infection prior to a transmission event. How much intrahost diversity is accumulated over time
51 from infection onset is currently unknown. Second, the population bottleneck during
52 transmission must be sufficiently wide to allow minor variants to be transmitted to recipient
53 hosts^{20,21}. Third, *de novo* generation of the same minor variants across multiple infections must
54 be sufficiently rare. Independent generation of shared minor variants by positive selection or
55 genetic drift in unrelated hosts could confound transmission inference¹⁵. Finally, measurements
56 of intrahost diversity must be accurate and account for several potential sources of error^{22,23}.
57 Although previous studies have described within-host variation of SARS-CoV-2^{7,9,11,12,24–26}, few
58 have addressed the sources of systematic errors and batch effects in variant identification. To
59 assess the utility of SARS-CoV-2 intrahost diversity for transmission inference, we need a
60 clearer understanding of its temporal variation throughout infection and the extent of convergent
61 evolution across individuals. Addressing these questions will also be valuable for understanding
62 SARS-CoV-2 evolution.

63

64 Here, we sequenced SARS-CoV-2 genomes from 325 residual upper respiratory samples from
65 hospitalized patients and employees at the University of Michigan. To validate our sequencing
66 approach, we sequenced defined mixtures of two synthetic RNA controls and found that low
67 input viral load decreases the specificity of variant calling. We find that observed intrahost
68 diversity does not vary significantly by day since symptom onset. Intrahost variants can be
69 shared between individuals that are unlikely to be related by transmission, suggesting that
70 variants can arise by parallel evolution. These results inform our understanding of SARS-CoV-2
71 diversification in human hosts and highlight important considerations for sequence-based
72 inference in the virus's genomic epidemiology.

73

74 **Results**

75 We retrieved respiratory specimens collected through diagnostic testing from March – May
76 2020. We sequenced samples from two groups: inpatients who were part of an observational
77 study of COVID-19 in hospitalized individuals (n = 190), and symptomatic employees who
78 presented to occupational health services (n = 135). All employees were diagnosed and treated
79 in outpatient settings, except for one who was admitted as an inpatient. Genome copy number
80 determined by qPCR of the nucleocapsid gene was highly variable and decreased by day from
81 symptom onset ($p < 0.001$, linear model, Fig. 1A). We obtained 212 complete genomes (Fig.
82 1B), mostly from samples with higher viral loads (Fig. 1B). Consensus genomes had a median
83 of 7 substitutions relative to the Wuhan-Hu-1/2019 reference sequence (range 4 – 12).
84 Phylogenetic analysis of whole consensus genomes identified 10 unique evolutionary lineages
85 in our cohort (lineages determined by the PANGOLIN system, see Methods; Fig. 1C). Most
86 sequenced genomes fell in lineage B.1. We evaluated whether any employees were part of an
87 epidemiologically linked cluster based on illness onset date, positive test status, and work
88 location. We found that some employees were part of epidemiologically linked clusters (Fig.
89 1C). The genomes from clusters 2, 10, 19, 20, and one pair in cluster 29 had ≤ 1 consensus
90 difference, while the rest had 2 – 7 differences. Many inter-cluster employee pairs also had
91 identical or nearly identical consensus genomes. We have no information on epidemiologic
92 linkage for the remaining sequenced individuals.

93

94 Identification of viral within-host variants can be prone to errors^{22,23}. Therefore, we performed a
95 mixing study to evaluate the accuracy of our pipeline for identifying intrahost single nucleotide
96 variants (iSNV). We mixed two synthetic RNA controls that differ by seven single nucleotide
97 substitutions at defined frequencies and input concentrations (Fig. 2A). These mixtures were
98 sequenced using the same approach as the clinical samples. We identified true iSNV at the

99 expected frequencies at $\geq 10^3$ copies/ μ L (Fig. 2B). There was greater variance in the observed
100 variant frequencies at 10^2 copies/ μ L compared to higher input concentrations. We obtained high
101 sensitivity for iSNV at $\geq 2\%$ frequency and $\geq 10^3$ copies/ μ L with sufficient genome coverage.
102 Many false positive iSNV remained at $\geq 2\%$ frequency and 10^2 copies/ μ L despite multiple quality
103 filters (Figure 2C, Supplemental Figure 1). However, false positive iSNV per sample drastically
104 decreased with input concentrations $\geq 10^3$ copies/ μ L. Three false positive variants were
105 identified in multiple samples above 10^4 copies/ μ L: A3350U, G6669A, and U13248A. Because
106 these iSNV were not randomly dispersed across the genome and were otherwise well-
107 supported in the sequence data, we suspect that they represent low-frequency variants present
108 in the synthetic RNA controls. Together, these data indicate that sufficient input viral load is a
109 critical factor for accurate identification of iSNV.

110
111 Based on our benchmarking experiment, we identified iSNV in 178 specimens with viral loads
112 $\geq 10^3$ copies/ μ L (Fig. 3A). We excluded position 11083, which is near a natural poly-U site and
113 prone to sequencing errors²⁷. Most specimens exhibited fewer than ten minor iSNV (median 1,
114 IQR 0 – 3, Fig. 3B). There were four outlier specimens with greater than 15 iSNV. In these
115 samples, iSNV were dispersed throughout the genome at various frequencies, so it is difficult to
116 determine whether they represent mixed infections¹¹. The locations of these samples on
117 sequencing plates were not suggestive of cross-contamination. There was no difference in
118 minor iSNV richness between hospitalized patients and employees treated as outpatients ($p =$
119 0.29 , Mann-Whitney U test, Supplemental Figure 2). We identified more minor iSNV encoding
120 non-synonymous changes than synonymous ones across most open reading frames (Fig. 3C)
121 and identified more iSNV at lower frequencies (Fig. 3D), which together is suggestive of mild
122 within-host purifying selection. Sample iSNV richness decreased with higher viral loads by about
123 1 iSNV per 10-fold increase in viral load ($p = 0.01$, multiple linear model, Supplemental Figure

124 3). Sample iSNV richness did not correlate with day from symptom onset ($p = 0.75$, multiple
125 linear model, Fig. 3E). These results show that within-host diversity is low and remains that way
126 over the duration of most SARS-CoV-2 infections.

127

128 Next, we investigated patterns of shared intrahost diversity between individuals. Most iSNV
129 were unique to a single individual. However, 19 iSNV were present in multiple specimens (Fig.
130 4A). These did not include the three recurrent false positives found in the synthetic RNA
131 controls. None of these mutations were located at sites known to commonly produce errors or
132 homoplasies^{27,28}. Two iSNV were present in three individuals (G12331A and A11782G, both
133 synonymous changes in ORF1a) and one iSNV was present in six individuals (U13914G,
134 encoding N149K in ORF1b). There was no clear phylogenetic clustering of genomes exhibiting
135 these shared iSNV (Supplementary Figure 4). The U13914G mutation was shared between
136 several sample pairs separated by 2 or more substitutions, and G12331A was shared between
137 samples from different viral lineages (13 substitutions). These three mutations were first
138 detected in our samples in late March 2020 (Fig. 4B). None reached > 1% frequency per week
139 in consensus sequences submitted to GISAID through mid-November 2020. These results
140 suggest that iSNV that arise convergently across viral lineages are not necessarily predictive of
141 subsequent global spread of those mutations.

142

143 Transmission inference based on shared iSNV integrates information such as consensus
144 genome sequences, sample dates, and shared iSNV¹⁵. Therefore, we compared shared iSNV
145 across all unique pairs of specimens used for variant calling ($n = 15753$, Fig. 5). Because most
146 iSNV were unique to an individual, most pairs did not share iSNV and only 0.23% of pairs
147 shared one iSNV. Many pairs with shared iSNV were sequenced in separate batches, which
148 reduces the likelihood that shared iSNV are due to cross-contamination. No employee pairs in
149 the same epidemiologic cluster shared iSNV (see Fig. 1C). We identified fourteen unique pairs

150 with shared iSNV between genomes that were near-identical (0 – 1 consensus differences),
151 eight of which were collected within one week of each other. However, we have no
152 epidemiologic data to suggest that these pairs of individuals are linked by transmission. We also
153 identified shared iSNV between 23 pairs separated by ≥ 2 consensus substitutions (Fig. 5A and
154 5B) and 15 pairs with collection dates 7 – 28 days apart (Fig. 5B). Due to differences in viral
155 lineage and time of collection, these are very unlikely to be transmission pairs. Together, these
156 data indicate that iSNV can arise convergently between individuals who are unlikely to be
157 related by transmission.

158

159 **Discussion**

160 Accurate characterization of SARS-CoV-2 intrahost diversity is important for understanding the
161 spread of new genetic variants and its potential use in transmission inference. In this study, we
162 sequenced upper respiratory specimens from a cohort of hospitalized COVID-19 patients and
163 infected employees. We found that intrahost diversity is low and its distribution does not vary by
164 time since symptom onset. We identified iSNV shared across viral genomes separated by time
165 and disparate evolutionary lineages, indicating that iSNV can arise convergently. Because
166 variants may be shared through parallel mutation rather than transmission, caution is warranted
167 in the use of shared iSNV alone for inferring transmission chains. Intrahost variants shared
168 across multiple individuals did not precede an increase in frequency in global consensus
169 genomes, which suggests that identifying convergent iSNV may have limited utility in tracking
170 broader SARS-CoV-2 evolution.

171

172 Specimen viral load is important when measuring intrahost diversity. We and others have shown
173 that samples with low viral loads are prone to false positive iSNV and lower sensitivity^{22,23,29}. A
174 strength of our study is that we experimentally validated the accuracy of our variant calling by
175 sequencing defined populations. Based on these results, we excluded samples with low viral

176 load from subsequent analyses. Future studies of SARS-CoV-2 intrahost diversity should report
177 and account for specimen viral loads to avoid this common source of error. We did not
178 benchmark our sequencing approach for detecting insertions and deletions (indels) and
179 therefore did not report these for the clinical specimens. Intrahost indels could conceivably
180 provide useful information about within-host evolution, but accurate detection is also subject to
181 similar issues of sample quality and viral load.

182

183 The low level of intrahost diversity that we found here is consistent with a recent preprint by
184 Lythgoe et al.⁹. The fact that our work and the study by Lythgoe et al. were performed with
185 different geographical areas, sequencing approaches (ARTIC Network amplicons vs. veSEQ
186 metagenomic sequencing), and analysis methods lends credence to the results. Lythgoe et al.
187 reported more shared variation than seen here, but this is most likely due to sequencing a
188 greater number of samples among individuals within known epidemiologic clusters. We and
189 Lythgoe et al. measure a lower level of intrahost diversity at the 2% frequency threshold
190 compared to a recent study in Austria¹². The reasons for this are not clear, but it is likely due to
191 differences in sample viral loads and variant calling methods. We did not find a difference in
192 intrahost diversity between hospitalized COVID-19 patients and those treated as outpatients,
193 which suggests that viral diversity may not be a reliable marker for disease severity.

194

195 Measuring viral diversity over the course of infection is relevant for understanding how variants
196 are transmitted to new hosts. Only genetic variants present at the time of a transmission event
197 will have the opportunity to spread. Because SARS-CoV-2 usually transmits just before or
198 several days after symptom onset^{30,31}, it is important to define viral diversity in this window. Our
199 cross-sectional analysis of diversity by time since symptom onset indicates that diversity does
200 not significantly increase over the course of infection. A significant fraction of samples may not
201 exhibit any iSNV at the time of transmission, which could limit the utility of iSNV for linking

202 transmission pairs. Only a large bottleneck would lead to onward spread of most iSNV present
203 during early infection. However, it is important to recognize that although the absolute level of
204 diversity may not change over time, different variants may arise or go extinct during a given
205 infection. This phenomenon was observed in a recent study by Tonkin-Hill et al.¹¹. Serial
206 samples from individuals could address this issue with higher resolution. Low diversity within
207 hosts also shapes our expectations for emergence of resistance to drugs and monoclonal
208 antibodies. With such limited substrate for selection to act upon, the short window of time
209 between treatment and transmission could limit the spread of a variant selected within a host.
210 Even during prolonged infections in immunocompromised hosts, there is only limited evidence
211 of resistance to various COVID-19 therapeutics³²⁻³⁴.

212

213 Parallel evolution is a critical factor to consider in the interpretation of shared intrahost
214 variation¹⁵. Even if iSNV identification were perfectly specific, iSNV can arise in parallel due to
215 biological processes such as natural selection and genetic drift. A key finding of this work is that
216 iSNV can arise in genomes that are unrelated by local transmission, specifically those across
217 large time intervals and lineages. Shared iSNV between individuals with identical genomes
218 collected the same week may also have arisen in parallel. These pairs are most likely not
219 epidemiologically linked, but we are unable to rule out coincident local transmission in the
220 community. Because iSNV can arise in parallel in genomes that are not linked by transmission,
221 caution is needed when relying entirely on shared iSNV for transmission inference^{11,13}.

222

223 We also found that identifying iSNV across multiple individuals did not precede an increase of
224 those mutations in frequency in global consensus genomes. It is unclear whether these
225 mutations arose due to positive selection, chance, or mutational “hotspots”¹¹. It is possible that
226 these mutations were lost due to purifying selection within hosts or during transmission^{8,35}.

227 These results suggest that iSNV may have lower utility for tracking broader SARS-CoV-2
228 evolution, but larger sample sizes in more geographic areas are necessary to evaluate this.
229
230 One of the most important variables for transmission inferences is the size of the transmission
231 bottleneck¹⁵. If parallel evolution of iSNV occurs regularly and the transmission bottleneck is
232 very small, that would increase the likelihood that shared iSNV are due to convergence rather
233 than transmission. However, if the bottleneck is large, then iSNV may become more valuable for
234 detecting transmission networks when consensus genomes are limited. There are currently
235 conflicting results on the SARS-CoV-2 bottleneck size. Popa et al. estimated a bottleneck size
236 of greater than 1000¹². In contrast, Lythgoe et al. estimated a bottleneck size range from 1 – 8
237 based on 14 household pairs⁹. Lythgoe et al. in particular used extensive controls and validation
238 for preventing contamination and identifying sequencing errors. Other studies both in humans
239 and in domestic cats have estimated small bottlenecks^{36,37}. It is difficult to interpret these
240 contrasting results because each study used different sequencing and analysis methodologies.
241 In recent work on influenza A virus, a study of methodological differences was key for resolving
242 different conclusions about the bottleneck size³⁸. One factor that has not yet been clearly
243 defined is how the time interval between donor-recipient pairs affects SARS-CoV-2 bottleneck
244 estimates. We expect that further work will clarify the reasons behind these conflicting
245 estimates.

246
247 Because of the high incidence and low mutation rate of SARS-CoV-2, genomic epidemiology is
248 necessarily constrained in its ability to determine exact transmission chains in an outbreak.
249 Using minor genetic variation to increase the resolution of genomic epidemiology requires
250 attention to the underlying processes of within-host viral evolution and awareness of possible
251 confounders. Unified statistical frameworks that incorporate sequences, metadata, and
252 epidemiological models are likely the most robust approaches for integrating intrahost variants,

253 but these models also must account for parallel evolution^{15–17}. As others have recently
254 suggested¹¹, we caution against assigning transmission pairs solely by virtue of shared iSNV in
255 the absence of clear epidemiologic information.

256

257 **Acknowledgements**

258 We thank the University of Michigan Clinical Microbiology Laboratory and the University of
259 Michigan Central Biorepository for their assistance in providing samples. We thank Christina
260 Cartaciano and the University of Michigan Microbiome Core for their assistance in sequencing.
261 We thank Emily Stoneman from Michigan Medicine Occupational Health Services for assistance
262 with employee data. This work was supported by a University of Michigan COVID-19 Response
263 Innovation Grant (to ASL), K01AI141579 (to JGP) and CDC U01 IP000974 (to ETM)

264

265 **Materials and Methods**

266

267 We collected clinical metadata and residual diagnostic specimens positive for SARS-CoV-2
268 from hospitalized patients enrolled in the CDC HAIVEN (Hospitalized Adult Influenza Vaccine
269 Effectiveness Network) study and infected employees enrolled in the HARVI (hospital
270 associated respiratory virus infection) study. These studies and the use of residual specimens
271 were approved by the University of Michigan Institutional Review Board.

272

273 Date of illness onset for hospitalized patients was collected individually via medical chart
274 abstraction from physician notes. Michigan Medicine employees with any suspected COVID-19
275 symptoms were asked to call a COVID-19 healthcare worker hotline before reporting to work.
276 Date of symptom onset, a list of symptoms, close contacts, travel history, and work location and
277 description were recorded. After testing, employee clusters were determined by illness onset
278 date, positive test status, and work location.

279

280 *Genome amplification and sequencing*

281 Residual samples from nasopharyngeal swabs and sputum specimens were centrifuged at 1200
282 x g. and 200 microliters were aliquoted. RNA was extracted with the Invitrogen PureLink Pro 96
283 Viral RNA/DNA Purification Kit and eluted in volumes of 100 microliters. Complementary DNA
284 was reverse transcribed with SuperScript IV (ThermoFisher). The SARS-CoV-2 genome was
285 amplified in two multiplex PCR reactions using the ARTIC Network V3 primer sets. Sequencing
286 libraries were prepared with the NEBNext Ultra II kit and pooled in equal volumes after
287 barcoding. The pooled sequencing library was gel extracted to remove adapter dimers. Libraries
288 were sequenced on an Illumina MiSeq at the University of Michigan Microbiome Core facility (v2
289 chemistry, 2x250 cycles). To validate this approach, we used two synthetic RNA controls that
290 differ by seven single nucleotide mutations, Wuhan-Hu-1 and EPI_ISL_418227 (Twist
291 Bioscience, San Francisco, CA). We mixed the two RNAs at various copy numbers (10^5 , 10^4 ,
292 10^3 , 10^2 genome copies/ μ L) and frequencies (0%, 0.25%, 0.5%, 1%, 2%, 5%, 10%, and 100%).
293 We amplified and sequenced each RNA mixture as described above.

294

295 *Viral load measurements*

296 We measured SARS-CoV-2 genome copy concentration for each sample by qPCR using
297 conditions outlined in the CDC 2019-Novel Coronavirus EUA protocol
298 (<https://www.fda.gov/media/134922/download>). The nucleocapsid gene was amplified using the
299 CDC N1 primer and probe set as follows: 2019-nCoV_N1 Forward Primer
300 GACCCCAAATCAGCGAAAT; 2019-nCoV_N1 Reverse Primer
301 TCTGGTACTGCCAGTTGAATCTG; 2019-nCoV_N1 Probe
302 ACCCCGCATTACGTTTGGTGGACC. Probe sequences were FAM labeled with Iowa Black
303 quencher (Integrated DNA Technologies, Coralville, IA). Reactions were performed using
304 TaqPath 1-step RT-qPCR master mix (Thermofisher, Waltham, MA) with 500 nM of each primer

305 and 250 nM of each probe in a total reaction volume of 20 μ l. Cycling conditions were as
306 follows: 2 min at 25 $^{\circ}$ C, 15 min at 50 $^{\circ}$ C, 2 min at 95 $^{\circ}$ C, and 45 cycles of 3 seconds at 95 $^{\circ}$ C, 30
307 seconds at 55 $^{\circ}$ C. Samples were run on an Applied Biosystems 7500 FAST real-time PCR
308 system. Cycle threshold (Ct) was designated uniformly across PCR runs.

309 Standard curves based on serial dilutions of a plasmid containing the nucleocapsid sequence
310 were used to determine copy number for each plate of samples. Copy number is expressed in
311 genome copies per microliter of extracted viral RNA.

312

313 *Analysis of sequence reads*

314 We aligned reads to the MN908947.3 reference genome with BWA-MEM version 0.7.15³⁹. We
315 removed sequencing adaptors and trimmed ARTIC primer sequences with iVar 1.2.1²³. We
316 determined the consensus sequences with iVar 1.2.1, taking the most common base as the
317 consensus (>50% frequency). We placed an N at positions along the MN908947.3 reference
318 with fewer than 10 reads. We manually inspected insertions and deletions by visualizing
319 alignments with IGV (version 2.8.0)⁴⁰. We identified single nucleotide variants with iVar 1.2.1
320 using the following parameters: sample with viral load $\geq 10^3$ copies/ μ L; sample with consensus
321 genome length of ≥ 29000 ; sample with $\geq 80\%$ of genome sites above 200x coverage; iSNV
322 frequency threshold of 2%; read depth of ≥ 100 at iSNV sites; ≥ 10 reads with average Phred
323 score of > 35 supporting a given iSNV; iVar p-value of < 0.0001 . All samples on which we called
324 variants had $> 50,000$ mapped reads. We accounted for strand bias by performing a two-sided
325 Fisher's exact test for hypothesis that the forward/reverse strand counts supporting the variant
326 base are derived from the same distribution as the consensus base. We then applied a
327 Bonferroni multiple test correction and excluded variants with an adjusted p-value < 0.05 . To
328 generate a phylogenetic tree, we aligned consensus genomes with MUSCLE 3.8.31 and
329 masked positions that are known to commonly exhibit homoplasies or sequencing errors⁴¹. We

330 generated a maximum likelihood phylogeny with IQ-TREE, using a GTR model and 1000
331 ultrafast bootstrap replicates^{42,43}. Evolutionary lineages (Pango lineages) were assigned with
332 PANGOLIN⁴⁴.

333

334 *Data and code availability*

335 Raw sequence reads are available as fastq files from the Sequence Read Archive at accession
336 number PRJNA682212, with human-mapping reads removed. Analysis code is available at
337 https://github.com/lauringlab/SARSCov2_Intrahost. Consensus genome sequences are publicly
338 available at the GitHub link and on GISAID.

339

340 **References**

- 341 1. Fauver, J. R. *et al.* Coast-to-Coast Spread of SARS-CoV-2 during the Early Epidemic in the
342 United States. *Cell* **181**, 990-996.e5 (2020).
- 343 2. Meredith, L. W. *et al.* Rapid implementation of SARS-CoV-2 sequencing to investigate cases
344 of health-care associated COVID-19: a prospective genomic surveillance study. *Lancet*
345 *Infect. Dis.* **0**, (2020).
- 346 3. Munnink, B. B. O. *et al.* Rapid SARS-CoV-2 whole-genome sequencing and analysis for
347 informed public health decision-making in the Netherlands. *Nat. Med.* 1–6 (2020)
348 doi:10.1038/s41591-020-0997-y.
- 349 4. Sekizuka, T. *et al.* Haplotype networks of SARS-CoV-2 infections in the Diamond Princess
350 cruise ship outbreak. *Proc. Natl. Acad. Sci.* (2020) doi:10.1073/pnas.2006824117.
- 351 5. Geoghegan, J. L. *et al.* Genomic epidemiology reveals transmission patterns and dynamics
352 of SARS-CoV-2 in Aotearoa New Zealand. *Nat. Commun.* **11**, 6351 (2020).
- 353 6. Miller, D. *et al.* Full genome viral sequences inform patterns of SARS-CoV-2 spread into and
354 within Israel. *Nat. Commun.* **11**, 5518 (2020).

- 355 7. Shen, Z. *et al.* Genomic Diversity of Severe Acute Respiratory Syndrome–Coronavirus 2 in
356 Patients With Coronavirus Disease 2019. *Clin. Infect. Dis.* **71**, 713–720 (2020).
- 357 8. Luring, A. S. Within-Host Viral Diversity: A Window into Viral Evolution. *Annu. Rev. Virol.*
358 (2020) doi:10.1146/annurev-virology-010320-061642.
- 359 9. Lythgoe, K. A. *et al.* Within-host genomics of SARS-CoV-2. *bioRxiv* 2020.05.28.118992
360 (2020) doi:10.1101/2020.05.28.118992.
- 361 10. Gutierrez, B., Escalera-Zamudio, M. & Pybus, O. G. Parallel molecular evolution and
362 adaptation in viruses. *Curr. Opin. Virol.* **34**, 90–96 (2019).
- 363 11. Tonkin-Hill, G. *et al.* Patterns of within-host genetic diversity in SARS-CoV-2. *bioRxiv*
364 2020.12.23.424229 (2020) doi:10.1101/2020.12.23.424229.
- 365 12. Popa, A. *et al.* Genomic epidemiology of superspreading events in Austria reveals
366 mutational dynamics and transmission properties of SARS-CoV-2. *Sci. Transl. Med.* **12**,
367 (2020).
- 368 13. Villabona-Arenas, C. J., Hanage, W. P. & Tully, D. C. Phylogenetic interpretation during
369 outbreaks requires caution. *Nat. Microbiol.* **5**, 876–877 (2020).
- 370 14. Sikkema, R. S. *et al.* COVID-19 in health-care workers in three hospitals in the south of
371 the Netherlands: a cross-sectional study. *Lancet Infect. Dis.* **0**, (2020).
- 372 15. Worby, C. J., Lipsitch, M. & Hanage, W. P. Shared Genomic Variants: Identification of
373 Transmission Routes Using Pathogen Deep-Sequence Data. *Am. J. Epidemiol.* **186**, 1209–
374 1216 (2017).
- 375 16. Maio, N. D., Worby, C. J., Wilson, D. J. & Stoesser, N. Bayesian reconstruction of
376 transmission within outbreaks using genomic variants. *PLOS Comput. Biol.* **14**, e1006117
377 (2018).
- 378 17. Skums, P. *et al.* QUENTIN: reconstruction of disease transmissions from viral
379 quasispecies genomic data. *Bioinformatics* **34**, 163–170 (2018).

- 380 18. Worby, C. J., Lipsitch, M. & Hanage, W. P. Within-Host Bacterial Diversity Hinders
381 Accurate Reconstruction of Transmission Networks from Genomic Distance Data. *PLoS*
382 *Comput. Biol.* **10**, e1003549 (2014).
- 383 19. Martin, M. A., Lee, R. S., Cowley, L. A., Gardy, J. L. & Hanage, W. P. Within-host
384 *Mycobacterium tuberculosis* diversity and its utility for inferences of transmission. *Microb.*
385 *Genomics* **4**, e000217 (2018).
- 386 20. McCrone, J. T. & Lauring, A. S. Genetic bottlenecks in intraspecies virus transmission.
387 *Curr. Opin. Virol.* **28**, 20–25 (2018).
- 388 21. Zwart, M. P. & Elena, S. F. Matters of Size: Genetic Bottlenecks in Virus Infection and
389 Their Potential Impact on Evolution. *Annu. Rev. Virol.* **2**, 161–179 (2015).
- 390 22. McCrone, J. T. & Lauring, A. S. Measurements of Intra-host Viral Diversity Are Extremely
391 Sensitive to Systematic Errors in Variant Calling. *J. Virol.* **90**, 6884–6895 (2016).
- 392 23. Grubaugh, N. D. *et al.* An amplicon-based sequencing framework for accurately
393 measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biol.* **20**, 8 (2019).
- 394 24. Wang, Y. *et al.* Intra-host Variation and Evolutionary Dynamics of SARS-CoV-2
395 Population in COVID-19 Patients. *bioRxiv* 2020.05.20.103549 (2020)
396 doi:10.1101/2020.05.20.103549.
- 397 25. Moreno, G. K. *et al.* Limited SARS-CoV-2 diversity within hosts and following passage in
398 cell culture. *bioRxiv* 2020.04.20.051011 (2020) doi:10.1101/2020.04.20.051011.
- 399 26. James, S. E. *et al.* High Resolution analysis of Transmission Dynamics of Sars-Cov-2 in
400 Two Major Hospital Outbreaks in South Africa Leveraging Intra-host Diversity. *medRxiv*
401 2020.11.15.20231993 (2020) doi:10.1101/2020.11.15.20231993.
- 402 27. Issues with SARS-CoV-2 sequencing data. *Virological* [https://virological.org/t/issues-](https://virological.org/t/issues-with-sars-cov-2-sequencing-data/473)
403 [with-sars-cov-2-sequencing-data/473](https://virological.org/t/issues-with-sars-cov-2-sequencing-data/473) (2020).
- 404 28. van Dorp, L. *et al.* Emergence of genomic diversity and recurrent mutations in SARS-
405 CoV-2. *Infect. Genet. Evol.* **83**, 104351 (2020).

- 406 29. Valesano, A. L. *et al.* The Early Evolution of Oral Poliovirus Vaccine Is Shaped by Strong
407 Positive Selection and Tight Transmission Bottlenecks. *Cell Host Microbe* **0**, (2020).
- 408 30. He, X. *et al.* Temporal dynamics in viral shedding and transmissibility of COVID-19. *Nat.*
409 *Med.* 1–4 (2020) doi:10.1038/s41591-020-0869-5.
- 410 31. Rhee, C., Kanjilal, S., Baker, M. & Klompas, M. Duration of Severe Acute Respiratory
411 Syndrome Coronavirus 2 (SARS-CoV-2) Infectivity: When Is It Safe to Discontinue Isolation?
412 *Clin. Infect. Dis.* doi:10.1093/cid/ciaa1249.
- 413 32. Baang, J. H. *et al.* Prolonged Severe Acute Respiratory Syndrome Coronavirus 2
414 Replication in an Immunocompromised Patient. *J. Infect. Dis.* **223**, 23–27 (2021).
- 415 33. Buckland, M. S. *et al.* Treatment of COVID-19 with remdesivir in the absence of humoral
416 immunity: a case report. *Nat. Commun.* **11**, 6385 (2020).
- 417 34. Kemp, S. *et al.* Recurrent emergence and transmission of a SARS-CoV-2 Spike deletion
418 Δ H69/V70. *bioRxiv* 2020.12.14.422555 (2020) doi:10.1101/2020.12.14.422555.
- 419 35. Xue, K. S., Moncla, L. H., Bedford, T. & Bloom, J. D. Within-Host Evolution of Human
420 Influenza Virus. *Trends Microbiol.* **26**, 781–793 (2018).
- 421 36. Wang, D. *et al.* Population Bottlenecks and Intra-host Evolution during Human-to-Human
422 Transmission of SARS-CoV-2. *bioRxiv* 2020.06.26.173203 (2020)
423 doi:10.1101/2020.06.26.173203.
- 424 37. Braun, K. M. *et al.* Transmission of SARS-CoV-2 in domestic cats imposes a narrow
425 bottleneck. *bioRxiv* 2020.11.16.384917 (2020) doi:10.1101/2020.11.16.384917.
- 426 38. Xue, K. S. & Bloom, J. D. Reconciling disparate estimates of viral genetic diversity
427 during human influenza infections. *Nat. Genet.* 1 (2019) doi:10.1038/s41588-019-0349-3.
- 428 39. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
429 *ArXiv13033997 Q-Bio* (2013).
- 430 40. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).

- 431 41. Masking strategies for SARS-CoV-2 alignments. *Virological*
432 <https://virological.org/t/masking-strategies-for-sars-cov-2-alignments/480> (2020).
- 433 42. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high
434 throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
- 435 43. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and
436 Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol. Biol.*
437 *Evol.* **32**, 268–274 (2015).
- 438 44. Rambaut, A. *et al.* A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist
439 genomic epidemiology. *Nat. Microbiol.* 1–5 (2020) doi:10.1038/s41564-020-0770-5.

440

441

442 **Figure Legends**

443

444 **Figure 1.** Viral shedding and overview of genome sequencing data. (A) Viral load by day of
445 infection in hospitalized patients (teal) and employees (violet). Viral load, measured by qPCR of
446 the N gene in units of genome copies per microliter of extracted RNA, is on the y-axis and day
447 post symptom onset is on the x-axis. (B) Genome completeness by viral load in hospitalized
448 patients (teal) and employees (violet). Viral load as shown in (A) is on the x-axis and the fraction
449 of the genome covered above 10x read depth is shown on the y-axis. (C) Maximum-likelihood
450 phylogenetic tree. Tips represent complete consensus genomes from hospitalized patients (teal)
451 and employees (violet). The axis shows divergence from the root (Wuhan-Hu-1/2019).
452 Heatmaps show PANGOLIN evolutionary lineage (left) and epidemiologic cluster (right).

453

454 **Figure 2.** Assessing accuracy of intrahost variant detection by sequencing defined viral
455 mixtures. (A) Schematic of the experiment. Wuhan-Hu-1 (reference) and EPI_ISL_418227
456 (variant) RNA were mixed at the given frequencies and viral loads (units of genome copies per
457 microliter, representing the resulting mixture). Mixtures of RNA were amplified and sequenced in
458 the same fashion as the clinical specimens. Reference and variant genomes differ by seven
459 single nucleotide substitutions. (B) Observed frequency by expected frequency. Observed
460 frequency of the true positive intrahost single nucleotide variants (iSNV) is on the y-axis and
461 expected iSNV frequency is on the x-axis. Synthetic RNA copy number in units of genome
462 copies per microliter of RNA is shown above each facet. Values above the points indicate the
463 number of variants detected in that group (maximum of seven per group). (C) False positive
464 iSNV. Number of false positive iSNV per sample is shown on the y-axis (base 10 log scale) and
465 viral load as shown in (B) is on the x-axis. Each point represents a unique sample and the
466 boxplots represent the median and 25th and 75th percentiles, with whiskers extending to the
467 most extreme point within the range of the median \pm 1.5 times the interquartile range.

468

469 **Figure 3.** SARS-CoV-2 intrahost single nucleotide variant (iSNV) diversity. (A) Sequencing
470 coverage for clinical samples. The number of clinical samples (y-axis) is shown by the fraction
471 of the genome above a given read depth threshold (x-axis). The different lines show the data
472 evaluated with six read depth thresholds. (B) Histogram of the number of specimens (y-axis) by
473 the number of minor iSNV per sample (x-axis), $n = 178$. (C) Number of minor iSNV by frequency
474 with a bin width of 0.05. Non-synonymous iSNV are shown in orange and synonymous iSNV are
475 shown in violet. (D) Number of minor iSNV by coding region. Non-synonymous iSNV are shown
476 in orange and synonymous iSNV are shown in violet. (E) Scatterplot of the number of minor
477 iSNV per sample (y-axis) by the day post symptom onset (x-axis). Hospitalized patients are
478 shown in teal and employees shown in violet. The four samples with > 15 iSNV shown in (B) are
479 excluded from the plot for visualization.

480

481 **Figure 4.** Shared iSNV across samples and their frequency in global consensus genomes. (A)
482 Shared iSNV across samples, with the number of samples sharing the iSNV (y-axis) by the
483 genome position (x-axis). Colors indicate the iSNV coding change relative to the reference. (B)
484 The frequency (y-axis) of three iSNV shared by three or more samples over time (x-axis). The
485 consensus genomes are from GISAID, as available on 2020-11-11. The vertical dotted lines
486 represent the earliest time we detected each iSNV in our samples.

487

488 **Figure 5.** Pairwise comparisons of shared iSNV. Each unique pair is shown as a single point,
489 with employee-employee pairs in violet (left), patient-employee pairs in orange (middle), and
490 patient-patient pairs in purple (right). The number of iSNV shared by each pair is shown on the
491 y-axis with the number of consensus differences between the pair of genomes on the x-axis.
492 Pairs of samples collected within seven days of each other are displayed in (A), and pairs of
493 samples collected greater than seven days apart are shown in (B).

495 **Supplemental Figure Legends**

496

497 **Supplemental Figure 1.** True and false positive iSNV in RNA mixture validation experiment.

498 Each iSNV is shown as a point, with the frequency on the y-axis and genome position on the x-

499 axis. True positive iSNV are shown in violet and false positive iSNV are shown in orange. All

500 iSNV displayed have a frequency of 2% or greater. Viral loads are shown above each facet, in

501 units of genome copies per microliter of RNA.

502

503 **Supplemental Figure 2.** Number of minor iSNV per sample (y-axis) across groups, with

504 hospitalized patients shown by teal points and employees shown by violet points. Boxplots for

505 each group represent the median and 25th and 75th percentiles, with whiskers extending to the

506 most extreme point within the range of the median \pm 1.5 times the interquartile range.

507

508 **Supplemental Figure 3.** Number of minor iSNV per sample (y-axis) by genome copies per

509 microliter of RNA (x-axis). Hospitalized patients are shown by teal points and employees shown

510 by violet points.

511

512 **Supplemental Figure 4.** Maximum likelihood phylogenetic tree as shown in Figure 1C. Tips

513 represent complete consensus genomes from hospitalized patients (teal) and employees

514 (violet). The x-axis shows divergence from the root (Wuhan-Hu-1/2019). Heatmaps show

515 samples that contain each of the three mutations as an iSNV.

Figure 1

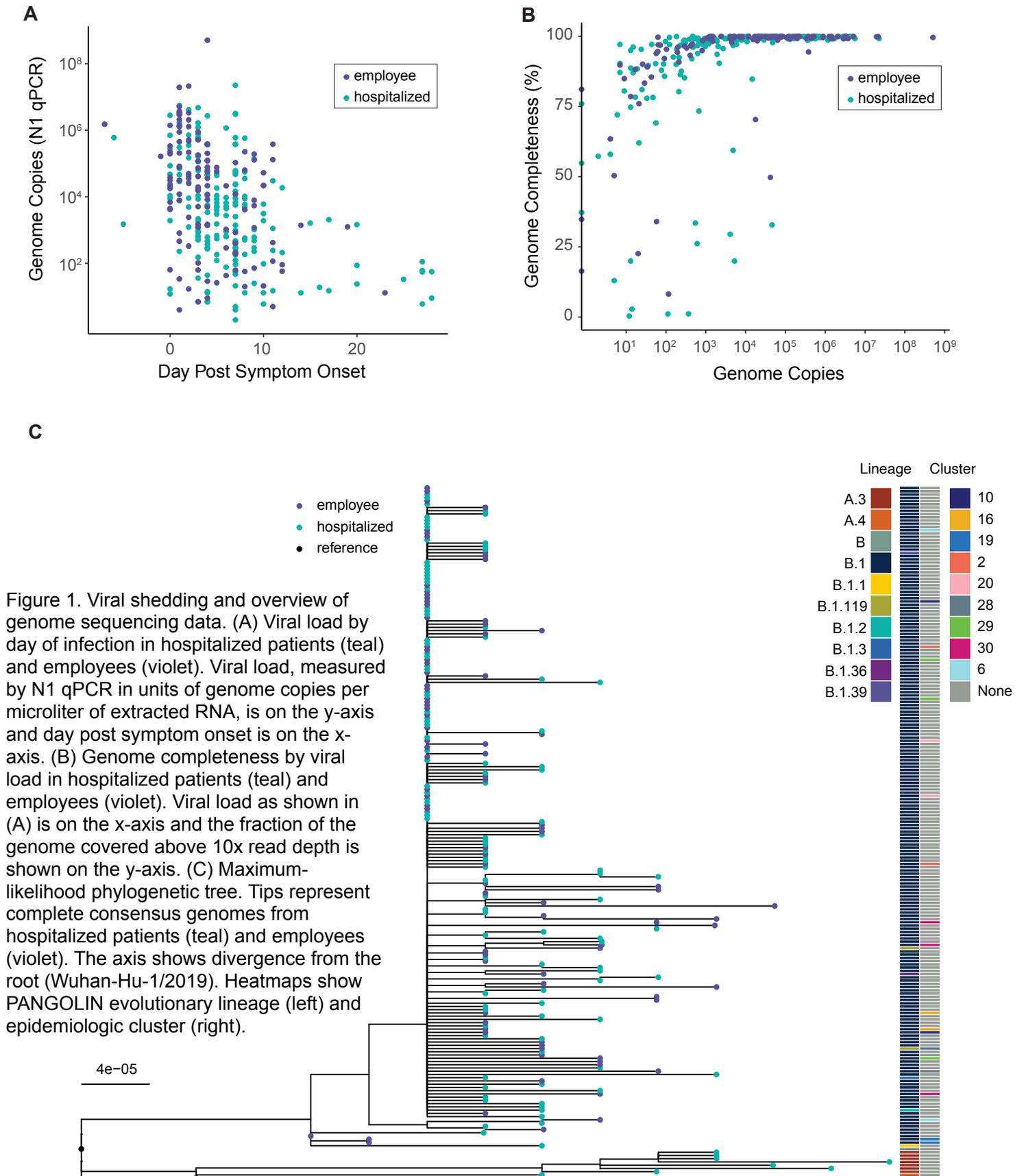


Figure 2

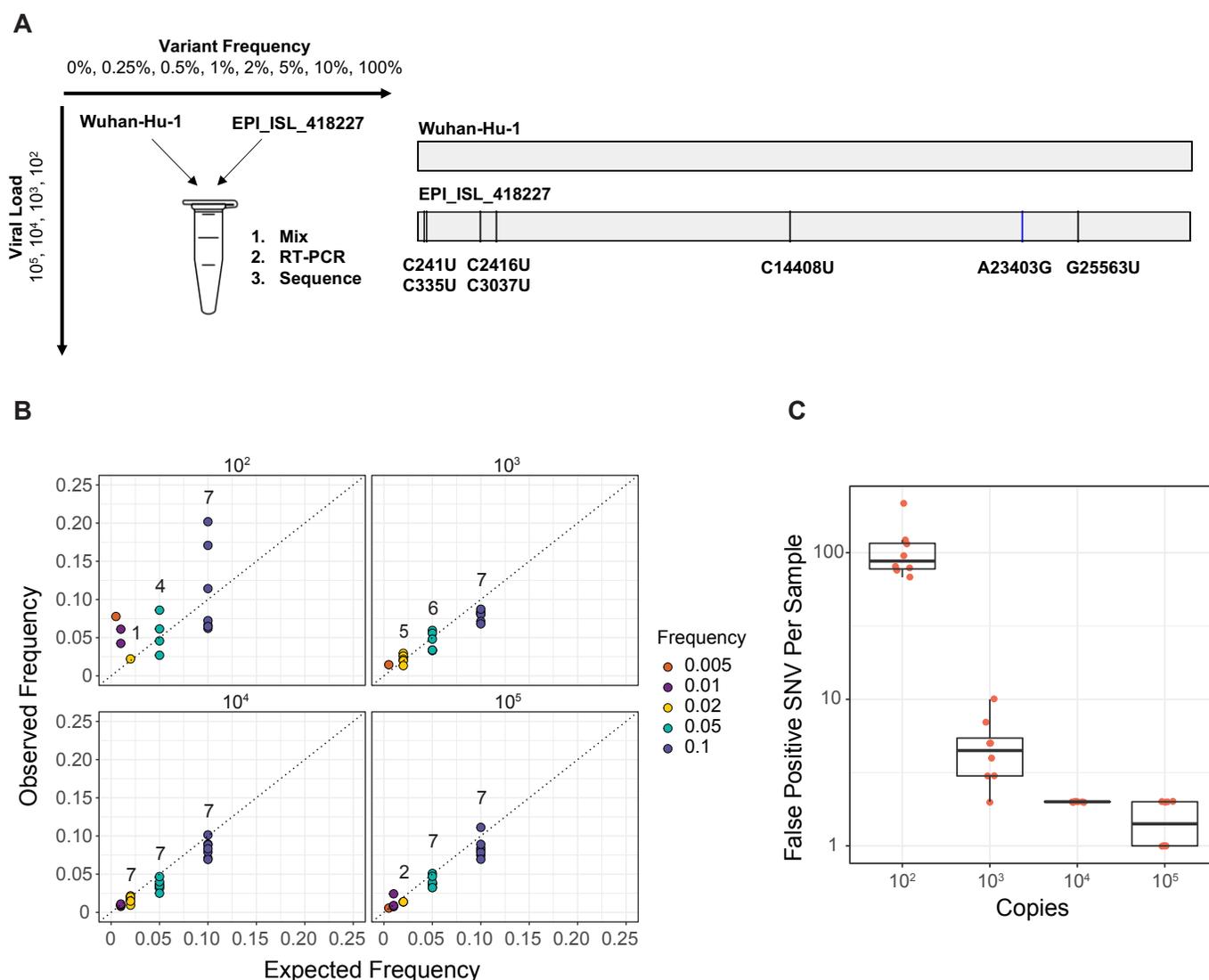


Figure 2. Assessing accuracy of intrahost variant detection by sequencing defined viral mixtures. (A) Schematic of the experiment. Wuhan-Hu-1 (reference) and EPI_ISL_418227 (variant) RNA were mixed at the given frequencies and viral loads (units of genome copies per microliter). Mixtures of RNA were amplified and sequenced in the same fashion as the clinical specimens. Reference and variant genomes differ by seven single nucleotide substitutions. (B) Observed frequency by expected frequency. Observed frequency of the true positive intrahost single nucleotide variants (iSNV) is on the y-axis and expected iSNV frequency is on the x-axis. Viral loads are shown above each facet, in units of genome copies per microliter of RNA. Values above the points indicate the number of variants detected in that group (maximum of seven per group). (C) False positive iSNV. Number of false positive iSNV per sample is shown on the y-axis (base 10 log scale) and viral load as shown in (B) is on the x-axis. Each point represents a unique sample and the boxplots represent the median and 25th and 75th percentiles, with whiskers extending to the most extreme point within the range of the median ± 1.5 times the interquartile range.

Figure 3

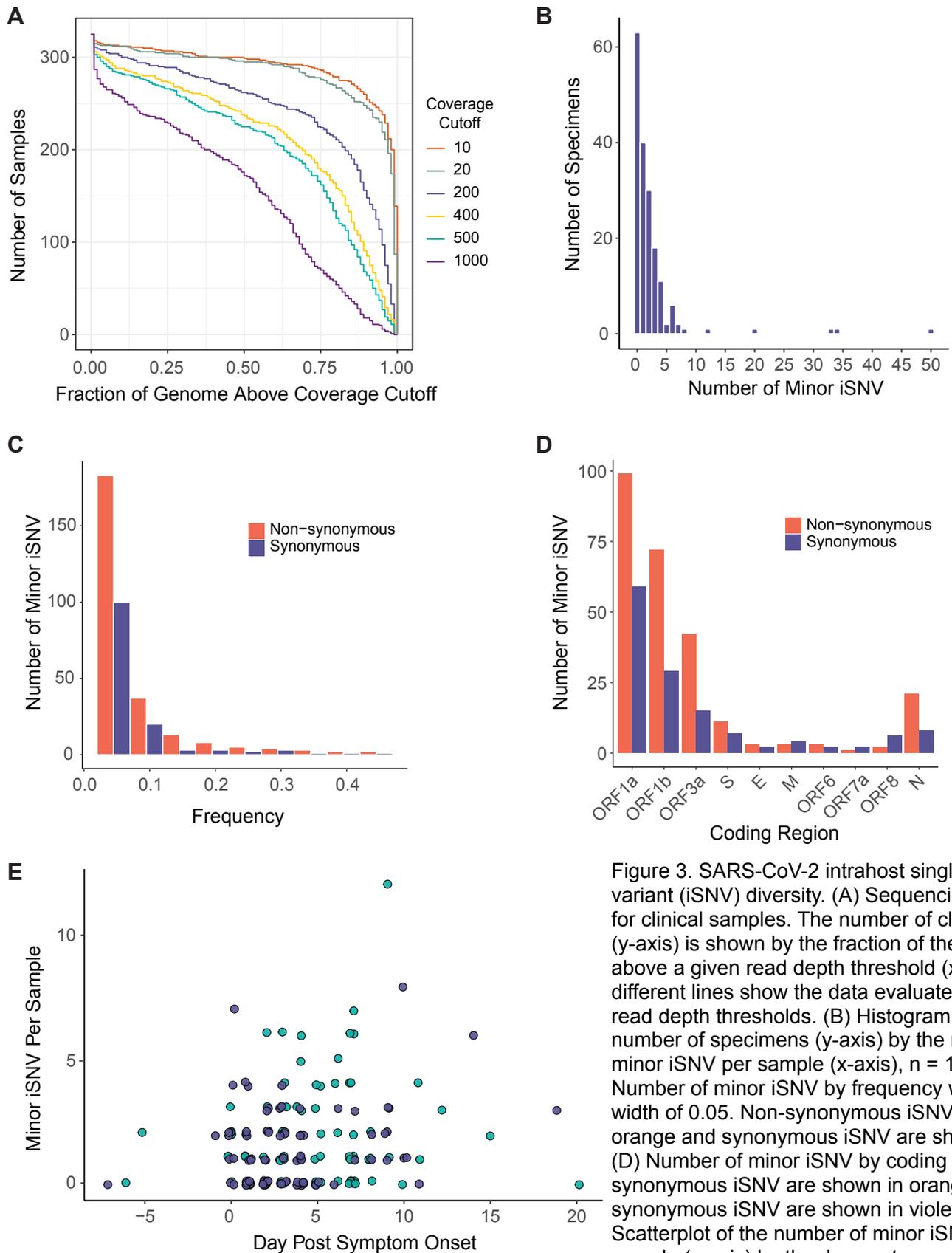


Figure 3. SARS-CoV-2 intrahost single nucleotide variant (iSNV) diversity. (A) Sequencing coverage for clinical samples. The number of clinical samples (y-axis) is shown by the fraction of the genome above a given read depth threshold (x-axis). The different lines show the data evaluated with six read depth thresholds. (B) Histogram of the number of specimens (y-axis) by the number of minor iSNV per sample (x-axis), $n = 178$. (C) Number of minor iSNV by frequency with a bin width of 0.05. Non-synonymous iSNV are shown in orange and synonymous iSNV are shown in violet. (D) Number of minor iSNV by coding region. Non-synonymous iSNV are shown in orange and synonymous iSNV are shown in violet. (E) Scatterplot of the number of minor iSNV per sample (y-axis) by the day post symptom onset (x-axis). Hospitalized patients are shown in teal and employees shown in violet. The four samples with > 15 iSNV shown in (B) are excluded from the plot for visualization.

Figure 4

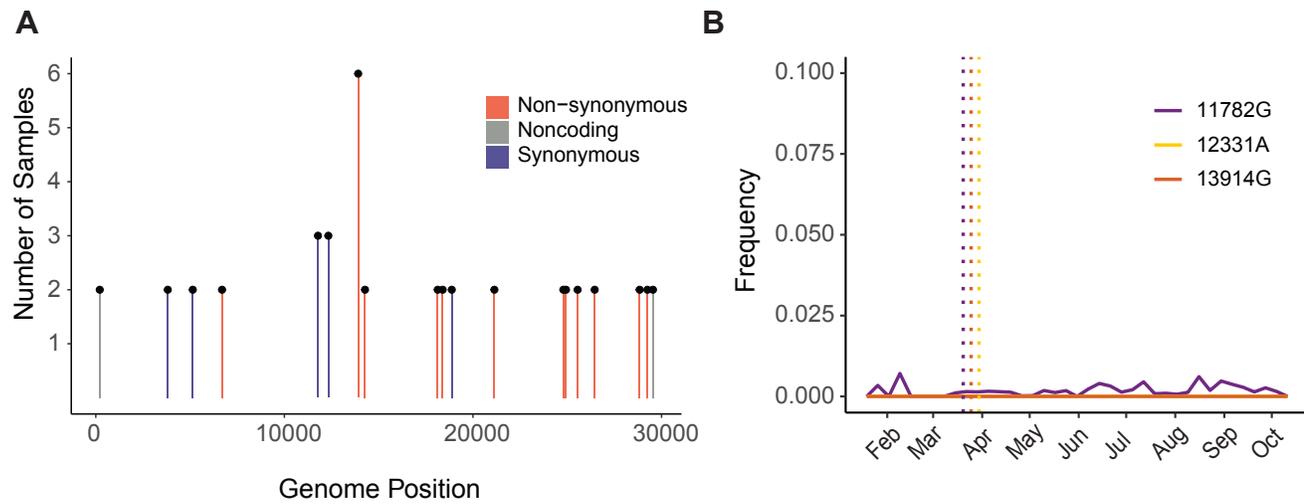


Figure 4. Shared iSNV across samples and their frequency in global consensus genomes. (A) Shared iSNV across samples, with the number of samples sharing the iSNV (y-axis) by the genome position (x-axis). Colors indicate the iSNV coding change relative to the reference. (B) The frequency (y-axis) of three iSNV shared by three or more samples over time (x-axis). The consensus genomes are from GISAID, as available on 2020-11-11. The vertical dotted lines represent the earliest time we detected each iSNV in our samples.

Figure 5

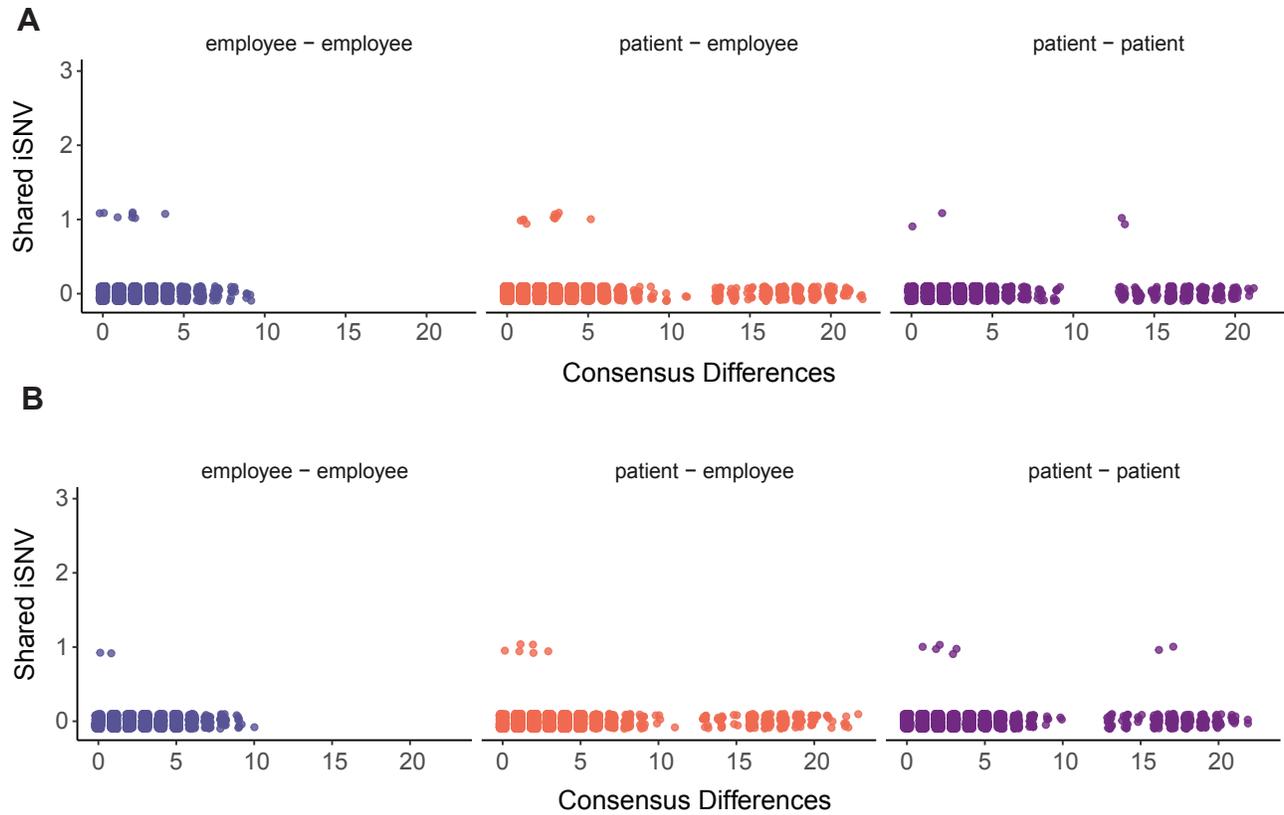


Figure 5. Pairwise comparisons of shared iSNV. Each unique pair is shown as a single point, with employee-employee pairs in violet (left), patient-employee pairs in orange (middle), and patient-patient pairs in purple (right). The number of iSNV shared by each pair is shown on the y-axis with the number of consensus differences between the pair of genomes on the x-axis. Pairs of samples collected within seven days of each other are displayed in (A), and pairs of samples collected greater than seven days apart are shown in (B).