1 # Limited SARS-CoV-2 diversity within hosts and following

2 # passage in cell culture

3

4 **Short title:** SARS-CoV-2 diversity is limited

5 **Authors**

6 Gage K. Moreno[1]*, Katarina M. Braun[2]*, Peter J. Halfmann[2,3], Trent M. Prall[1], Kasen K.

7 Riemersma[2], Amelia K. Haj[1], Joseph Lalli[2], Kelsey R. Florek[3], Yoshihiro Kawaoka[2,4], Thomas C.

8 Friedrich[2,4], David H. O'Connor[1,5, #]

9 * These authors contributed equally to this work

10

11 **Affiliations**

12 [1]Department of Pathology and Laboratory Medicine, University of Wisconsin-Madison, Madison,

13 WI, United States of America

14 [2]Department of Pathobiological Sciences, University of Wisconsin-Madison, Madison, WI,

15 United States of America

16 [3]Wisconsin State Laboratory of Hygiene, Madison, WI, United States of America

17 [4]Influenza Research Institute, School of Veterinary Sciences, University of Wisconsin-Madison,

18 Madison, WI, United States

19 [5]Wisconsin National Primate Research Center, University of Wisconsin-Madison, Madison, WI,

20 United States of America

21    **Abstract**

22    Since the first reports of pneumonia associated with a novel coronavirus (COVID-19) emerged

23    in Wuhan, Hubei province, China, there have been considerable efforts to sequence the

24    causative virus, SARS-CoV-2 (also referred to as hCoV-19) and to make viral genomic

25    information available quickly on shared repositories. As of 30 March 2020, 7,680 consensus

26    sequences have been shared on GISAID, the principal repository for SARS-CoV-2 genetic

27    information. These sequences are primarily consensus sequences from clinical and passaged

28    samples, but few reports have looked at diversity of virus populations within individual hosts or

29    cultures. Understanding such diversity is essential to understanding viral evolutionary dynamics.

30    Here, we characterize within-host viral diversity from a primary isolate and passaged samples,

31    all originally deriving from an individual returning from Wuhan, China, who was diagnosed with

32    COVID-19 and subsequently sampled in Wisconsin, United States. We use a metagenomic

33    approach with Oxford Nanopore Technologies (ONT) GridION in combination with Illumina

34    MiSeq to capture minor within-host frequency variants ≥1%. In a clinical swab obtained from the

35    day of hospital presentation, we identify 15 single nucleotide variants (SNVs) ≥1% frequency,

36    primarily located in the largest gene – ORF1a. While viral diversity is low overall, the dominant

37    genetic signatures are likely secondary to population size changes, with some evidence for mild

38    purifying selection throughout the genome. We see little to no evidence for positive selection or

39    ongoing adaptation of SARS-CoV-2 within cell culture or in the primary isolate evaluated in this

40    study.

41    **Author Summary**

42    Within-host variants are critical for addressing molecular evolution questions, identifying

43    selective pressures imposed by vaccine-induced immunity and antiviral therapeutics, and

44    characterizing interhost dynamics, including the stringency and character of transmission

45     bottlenecks. Here, we sequenced SARS-CoV-2 viruses isolated from a human host and from

46     cell culture on three distinct Vero cell lines using Illumina and ONT technologies. We show that

47     SARS-CoV-2 consensus sequences can remain stable through at least two serial passages on

48     Vero 76 cells, suggesting SARS-CoV-2 can be propagated in cell culture in preparation for *in-*

49     *vitro* and *in-vivo* studies without dramatic alterations of its genotype. However, we emphasize

50     the need to deep-sequence viral stocks prior to use in experiments to characterize sub-

51     consensus diversity that may alter outcomes.

52

### 53     **Introduction**

54     The emergence of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) and

55     coronavirus disease (COVID-19) in Wuhan, China at the end of 2019 has garnered worldwide

56     public health attention [1–5]. At the time of writing, the United States has the highest number of

57     confirmed cases among countries where this virus is circulating – 639,733 cases and 30,990

58     deaths.

59     The rapid spread and molecular epidemiology of SARS-CoV-2 has been tracked by sequencing

60     viruses from infected individuals. Within weeks of the virus being identified, the complete

61     genome was sequenced, and as of April 16th 2020, 9,330 SARS-CoV-2 genomes have been

62     shared and used to track local transmission chains and global phylodynamics [6]. While

63     consensus-level data has been rapidly disseminated, few researchers have analyzed viral

64     diversity within samples below the consensus level.

65     SARS-CoV-2 is a betacoronavirus with 79-82% nucleotide identity shared with SARS-CoV, the

66     virus responsible for the 2002 - 2003 SARS epidemic [7, 8]. During the 2003 SARS outbreak the

67     virus was characterized as having gone through distinct evolutionary phases in human hosts.

68     Initially, an excess of nonsynonymous mutations in the spike (S) gene suggested that it might

69    be under positive selection, but this progressed into purifying selection later in the epidemic [9].

70    The ORF1a gene appeared to go through similar evolutionary phases as the S gene. In contrast

71    to ORF1a and S, the ORF1b gene appeared to have undergone strong purifying selection

72    throughout the 2003 SARS epidemic [9].

73    Though limited, *in-vivo* studies of SARS-CoV-2 show low-frequency variants are detectable

74    within individual hosts and are likely due to random fluctuations in allele frequencies. One study

75    highlights an excess of nonsynonymous variants compared to synonymous variants among

76    these low-frequency variants, consistent with the possibility of ongoing diversifying selection in

77    SARS-CoV-2 viruses [10, 11]. Another recent study by Liu and colleagues highlights a deletion

78    in the Spike gene at nucleotide (nt) positions 23,585–23,599, encoding QTQTN, that flanks the

79    polybasic cleavage site in S1/S2. The authors observe this deletion arising in SARS-CoV-2

80    viruses following two passages in Vero E6 cells. This deletion is found in over 50% of samples

81    from Liu and colleagues, ranging in frequency from 8 to 33%, and is hypothesized to be

82    adaptive for SARS-CoV-2 *in vitro,* but may be less robust *in vivo* as it was only identified in 3 of

83    68 Chinese-origin clinical samples at sub-consensus levels [12].

84    To better understand evolutionary pressures affecting SARS-CoV-2 within a single infection, we

85    used sequence-independent, single-primer amplification (SISPA) to generate metagenomic

86    libraries sequenced in parallel on Oxford Nanopore Technology (ONT) and Illumina sequencing

87    platforms (**Fig 1**) [13, 14]. We obtained a nasopharyngeal (NP) swab from an individual with

88    confirmed SARS-CoV-2 infection from the day of diagnosis, who originally presented with

89    symptoms in Madison, WI (hereafter referred to as the Madison patient). This case was

90    diagnosed in late January 2020 and was one of the first lab-confirmed cases in the United

91    States. We additionally characterized viral diversity following passage in cell culture in three

92    distinct cell types – Vero 76, Vero E6, and Vero STAT-1 knockout (KO). Passage in cell culture

93    is expected to alter allele frequencies and may even select for adaptive mutations that make

94  passaged viruses less representative of their genotypes and phenotypes *in vivo*. Global viral

95  evolution ultimately derives from selective pressures and population dynamics playing out within

96  and between individual hosts. In this study, we identify SNVs within a clinical specimen and

97  track what happened to them through multiple rounds of passage in culture and begin to

98  assemble a nuanced understanding of the evolution and ongoing adaptive potential of this

99  zoonotic virus.

100  **Results**

101  **No consensus-level changes following two passages on Vero 76 cells**

102  We obtained an NP swab from the day of diagnosis and passaged the virus on three distinct cell

103  lines – Vero 76, Vero E6, and Vero STAT-1 KO (**Fig 1**). To understand the effects of serial

104  passaging on SARS-CoV-2, we used the SISPA approach to generate full genome sequencing

105  libraries from the original NP swab and passaged virus (**S1 Fig**). Sequences were analyzed in

106  parallel using custom in-house scripts to deplete host reads, map to the SARS-CoV-2 Madison

107  reference (Genbank: MT039887.1; originally sequenced by the US Centers for Disease Control

108  and Prevention), and call minor variants ≥10% and ≥1% for ONT and Illumina datasets,

109  respectively. We detect no consensus-changing SNVs through two passages on Vero 76 cells

110  and through one passage on Vero E6 and Vero STAT-1 KO cells (passage 2 samples were not

111  available in these cell lines) (**Fig 2a**).

112  Interestingly, in comparison to the sequence derived from the first case of SARS-CoV-2

113  (MN908947.3), the Madison patient's virus contained an in-frame deletion at nucleotide

114  positions 20,298 - 20,300 (**Fig 2b**). This deletion has not been identified in any other samples

115  submitted to GISAID as of April 8, 2020. This deletion occurs in a region that codes for the

116  poly(U)-specific endoribonuclease, but its functional impact is not clear [12].

117  **No deletion in spike gene after passaging in cell culture**

118   To understand how serial passaging SARS-CoV-2 affects genomic variation, we sequenced

119   virus populations after each passage using the same SISPA metagenomics approach we used

120   to characterize the original biological specimen. Passaged sample names and cell lines are

121   described    in    the    methods.    An    in-house    pipeline    (available    at:

122   https://github.com/katarinabraun/SARSCoV2_passage_MS) was applied to trim out primer

123   sequences, bioinformatically deplete host reads, and generate alignment files, which contained

124   all reads mapping to the SARS-CoV-2 Madison reference genome (MT039887.1). At the

125   consensus level, SARS-CoV-2 does not accumulate genetic variation after two passages on

126   Vero 76 cells (**Fig 2**). We also examined deletions ≥1% frequency and ≥3 nt in length. We found

127   no evidence of deletions that fit these criteria in any of the cell culture isolates.

**Most minor variants are found in the largest genes – ORF1a and ORF1b**

129   To characterize patterns of sub-consensus diversity, we looked at SNVs at or above 1%

130   frequency in only the Illumina reads. We previously established that this conservative cutoff

131   ensures that only bona fide mutations are considered [15, 16]. All minor variant analyses and

132   figures were completed using the Illumina SNV data as these data are higher average quality

133   and ideal for analysis involving low-frequency variants (**Fig 3**). Seventy-five percent of all minor

134   variants we identify fall in ORF1a and ORF1b, which together take up 72.8% of the length of the

135   28kb coding genome. ORF1a and ORF1b encode the replicase machinery [7]. We account for

136   differences in gene size by normalizing variants to kilobase gene length (variants / kb-gene-

137   length – "v/kbgl") [10]. The highest density of variants was reported in smaller genes like

138   envelope, ORF7a, and ORF10 (**S2 Table**). We also show that through each passage, variant

139   density in ORF1a and ORF1b increases. There were no SNVs ≥1% in the spike gene in the

140   primary NP swab, but low-frequency SNVs (all <5%) were identified in spike following passage

141   in cell culture (**Fig 3**). Outside of ORF1a and ORF1b, the other genes in the primary NP swab

142    are clonal above the 1% threshold, with the exception of one low-frequency SNV in

143    nucleoprotein (N).

144    A few SNVs at intermediate frequencies or identified across multiple samples stood out. A

145    synonymous SNV at nucleotide position 11,070 (ORF1a_11070_syn) was found at ≥15%

146    frequency in the primary NP swab as well as in all passaged samples. Amino acid positions

147    3,570 - 3,859 in ORF1a are predicted to be involved in the formation of double-membraned

148    vesicles [7]. Variants at nucleotide positions 127 (nonsynonymous – asparagine to aspartic acid;

149    ORF7a_127_N43D) and 129 (synonymous; ORF7a_129_syn) were identified between 1-4%

150    frequency in all passaged samples, but were not detected in the primary NP swab. ORF7a has

151    no known function, so the impact of these SNVs is unclear [7]. These SNVs

152    (ORF1a_11070_syn, ORF7a_127_N43D, and ORF7a_129_syn) have not been identified as

153    major variants in any of the SARS-CoV-2 genomes submitted to GISAID as of 12 April, 2020.

154    Six variants identified in at least one sample evaluated here have been identified as major

155    variants in at least one sequence on Nextstrain as of 12 April, 2020. These SNVs include

156    ORF1a_8025_syn     (p2b    Vero    76)    found    in    England/201380056/2020,

157    England/20146004904/2020, and Australia/VIC164/2020; ORF1a_11409_syn (p2a Vero 76)

158    found in HongKong/HKPU2_1801/2020; ORF1b_5843_T1948I (p2a Vero 76 and p2b Vero 76)

159    found in China/IQTC02/2020; S_1640_T547I (p1 Vero 76) found in USA/WA-S17/2020;

160    S_2661_syn (p2b Vero 76) found in HongKong/HKPU1_2101/2020; and ORF3a_385_L129F

161    (p1 Vero-1 STAT KO and p1 Vero 76) found in Algeria/G0638_2264/2020. Interestingly, all six

162    of these SNVs are a cytosine to thymine transitions.

163    We also determined whether SNVs were shared among the primary NP swab and passaged

164    viruses (**S2 Fig)**. Thirteen of the 15 minor variants identified in the primary NP swab are purged

165    following passage in cell culture. Only two SNVs were found in all of the available samples –

166     ORF1a V1118A and the synonymous SNV at nt 11,070 in ORF1a. ORF1a V1118A remains

167    between 1-2% in all viruses. However, ORF1a 11,070-syn is found at 3% in the primary NP

168    swab and increases in frequency to 18% in p1 Vero 76, remaining above 10% in both p2 Vero

169    76 samples. Only two *de novo* SNVs are found above 10% in cell culture – ORF1a_10242_syn

170    (p2a Vero 76 and p2b Vero 76) and ORFb_5843_T1948I (p2b Vero 76).

171    **SNV frequency spectra reveal an excess of low-frequency SNVs**

172    Purifying selection is known to remove new variants from the population, generating an excess

173    of low-frequency variants, while positive and/or diversifying selection promotes the accumulation

174    of intermediate- and high-frequency variation [17]. Especially in the setting of an acute viral

175    infection, exponential population growth can also result in an excess of low-frequency variants.

176    Population bottlenecks, for example sharp reductions in a viral population size typically

177    associated with airborne viral transmission, can contribute to an excess of intermediate- and

178    high-frequency variation. We generated site frequency spectra to expand our assessment of the

179    evolutionary pressures impacting SARS-CoV-2 viruses within humans and in cell culture. A

180    "neutral model" (assumes a constant population size and the absence of selection), represented

181    in light grey in **Fig 4**, predicts around 50% of polymorphisms will be low-frequency (1-10%). In

182    stark contrast to the neutral expectation, we observed ≥80% of SNVs falling into the low-

183    frequency bin in the primary nasal swab sample as well as passaged samples. This dramatic

184    excess of low-frequency variation is consistent with purifying selection acting to purge new,

185    deleterious mutations. This signature is also consistent with population expansion as is

186    expected in humans following airborne transmission and in cell culture after each passage.

187    **Nucleotide diversity patterns point toward mild purifying selection**

188    In addition to assessing the fate of individual minor variants, we were also interested in

189    evaluating population dynamics using diversity metrics. Specifically, we calculated genewise

190    diversity using π, the average number of pairwise differences per nucleotide site among a set of

191    sequences, for each gene in each sample. Overall, genewise nucleotide diversity is very low

192    compared to other RNA viruses, consistent with low mutation rates in coronaviruses due to RNA

193    proofreading machinery [18, 19]. Genewise diversity was very low in the primary NP swab and

194    was only measurable in ORF1a (9 SNVs), ORF1b (5 SNVs) and N (1 SNV). Genewise diversity

195    is more varied in the passaged samples (**Fig 5**). Interestingly, π is highest in ORF7a in these

196    samples – although this signal seems to be primarily driven by the small size of this gene. To

197    more directly assess whether SARS-CoV-2 viruses are under selective pressure in the human

198    infection evaluated here and in cell culture, we also compared the relative abundance of

199    nonsynonymous (πN) and synonymous (πS) polymorphisms in each gene, which is a common

200    measure for selection that is also robust to variability in sequencing coverage depth [20]. The

201    dominant genetic signature when looking across the entire genome is one of purifying selection

202    (πN/πS < 1). In ORF1a, πS > πN in the primary NP swab as well as p1 and p2 samples. In

203    ORF1b, πN/πS is close to 1 in the primary NP swab and the p1 on Vero 76 and Vero E6 cells,

204    suggesting a more prominent role of genetic drift in this gene. Interestingly, πN/πS >> 1 in p1

205    Vero 76 ORF10, p1 Vero E6 envelope (E), and p1 Vero STAT-1 KO ORF3a.

**Comparison of Illumina and ONT ability to capture minor variant frequencies**

207    We examined the concordance between SNV calls at the same sites, irrespective of frequency,

208    determined by Illumina and ONT workflows. To begin, we used a stringent cutoff of 10%

209    frequency for ONT SNVs. We then called variants at percentage frequencies decreasing by

210    0.5% (eg. calling 8% variants, then 7.5%, etc) until the variants called by ONT no longer

211    matched Illumina variants irrespective of frequency at these sites (**Fig 6, S1 Table.**). We found

212    that for the primary NP swab we were able to call minor variants that occurred at ≥8%

213    frequency. Below 8% frequency, SNVs called by ONT were no longer exactly concordant with

214    SNVs called by Illumina. Discrepancies between ONT and Illumina variant calls at low

215    frequencies are tied to ONT's high false discovery rate, a finding previously documented by

216     Grubaugh and colleagues in 2019 [21]. For the p1 samples, ONT was able to capture variants

217     that occurred at ≥4.5% frequency. For the p2 samples, we called SNVs down to 8.5% and 5.5%

218     for the p2a and p2b samples, respectively. We likely observed concordant SNV calls between

219     Illumina and ONT at lower frequencies in the passaged samples because viral titer *in vitro*

220     typically exceeds viral titer *in vivo* resulting in higher average coverage in the passaged samples

221     required to support minor variant calls at lower frequencies.

222     **Discussion**

223     Minor variants are critical for addressing molecular evolution questions, identifying selective

224     pressures imposed by vaccine-induced immunity and antiviral therapeutics, and characterizing

225     interhost dynamics, including the stringency and character of transmission bottlenecks. Parallel

226     consensus-level data of clinical isolates are similarly important and allow us to predict

227     transmission patterns on a global, regional, and community-wide scale. Here, we explore

228     SARS-CoV-2 intrahost variation from a primary NP swab as well as from viruses passaged on

229     three distinct Vero cell lines. We show that while diversity is low overall, the dominant viral

230     genetic signature is one of mild purifying selection, evidenced by an excess of low-frequency

231     variants and the observation that $\pi N/\pi S < 1$ in most genes across all samples evaluated.

232     We show that SARS-CoV-2 consensus sequences can remain stable through at least two serial

233     passages on Vero 76 cells even in the presence of a three nucleotide deletion in the region of

234     the genome encoding the poly(U)-specific endoribonuclease, suggesting SARS-CoV-2 can be

235     propagated in cell culture in preparation for *in vitro* and *in vivo* studies without dramatic

236     alterations of its genotype. A recent paper by Duggal et al. illustrate the importance of viral

237     genotype instability in Zika virus (ZIKV) by describing variants enriched during cell culture

238     passage (Envelope-330L/NS1-98G), despite being attenuated *in vivo* and responsible for a less

239     pathogenic phenotype in mice compared to the wildtype genotype (Envelope-330V/NS1-98W)

240    [22]. Viral genotype instability in cell culture can significantly affect animal model development

241    and vaccine efficacy studies.

242    Though we do detect a handful of minor variants in ORF1a and ORF1b in the primary NP swab,

243    it is notable that eight out of eleven genes are clonal above the 1% frequency level. As natural

244    selection can only act upon genetic variation already existing within a population, very limited

245    intrahost genetic diversity suggests the pace of SARS-CoV-2 evolution may be primarily limited

246    by the generation of *de novo* variants. It is unclear at this time the degree to which limited

247    within-host viral diversity is linked to coronavirus biology – e.g. RNA proofreading capabilities,

248    homologous recombination allowing for the decoupling of deleterious "hitchhiker" mutations, and

249    a comparatively low mutation rate. Studies have estimated the mutation rate of coronaviruses

250    to be $2 \times 10^{-6}$ mutations per site per round of replication, which is in line with other

251    coronaviruses [18], but lower than influenza, $7.1 \times 10^{-6} – 4.5 \times 10^{-5}$ mutations per site per round

252    of replication, another respiratory RNA virus [23–27].

253    A previous study claimed that a common deletion at nt position 23,585–23,599 (spike),

254    encoding QTQTN, arises after two passages in Vero E6 cells [12]. We did not identify similar

255    deletions in this region in any of our passaged samples, suggesting this deletion is not as

256    common as previously suggested. Interestingly, the primary NP swab obtained from the

257    Madison patient on the day of diagnosis contained an in-frame deletion at nucleotide positions

258    20,298 - 20,300 (ORF1ab) that was retained through two passages on Vero 76 cells. These

259    genomic deletions highlight the importance of characterizing viral stocks by deep-sequencing so

260    genotypic differences that may alter experimental outcomes can be thoroughly documented and

261    shared with other researchers.

262    Below the consensus level, we found an excess of low-frequency variants compared to what

263    would be expected in a neutral setting with no changes in population size and no selective

264 pressures at play. This suggests that either purifying selection is acting to remove new, mildly

265 deleterious mutations in hosts and in culture before they can reach intermediate or high

266 frequencies, and/or the virus is undergoing exponential population growth as would be expected

267 in an acute viral infection or following passage in cell culture. It is likely that viral exponential

268 population growth is contributing to this genetic signature; however, without additional samples,

269 it is difficult to determine the relative contribution of each of these factors. We would emphasize

270 these findings are rooted in relatively few, low-frequency SNVs from a single time point so

271 conclusions about the overall evolution of SARS-CoV-2 are necessarily limited. Continued deep

272 sequencing and analyses of SARS-CoV-2 minor variant SNV populations in humans and in cell

273 culture are critical.

274

275 **Methods**

276 **Sample collection and cell culture passage conditions**

277 Three different Vero cell lines were purchased from ATCC; Vero 76 (ATCC: CRL-1587), Vero

278 C1008 (ATCC: CRL-1586), Vero STAT-1 KO (ATCC: CCL-81-VHG), and were grown in

279 Minimum Essential Medium (MEM) supplemented with 10% fetal bovine serum (FBS) and L-

280 glutamine at $37^{o}$C with 5% $CO_2$.

281 For the initial infection, the original clinical nasopharyngeal (NP) swab was divided evenly

282 between three TC25 $cm^2$ flasks seeded the day before with 1 x $10^6$ cells per flask; one flask for

283 each Vero cell line. Virus in the original clinical sample was layered onto the cells for one hour

284 at $37^{o}$C, the flasks were washed once with MEM, and the medium was replaced with fresh MEM

285 supplemented with 2% FBS. For each additional passage, cells were seeded in 75 $cm^2$ flasks

286 the day before infection with 4 x $10^6$ cells per flask and infected at a multiplicity of infection

287     between 0.01-0.001. For each passage, the virus was harvested when cell death was observed

288     to be around 80% (~4-5 days after infection).

289     Work with live virus was performed at biosafety level-3 containment at the Influenza Research

290     Institute at the University of Wisconsin – Madison under a recombinant DNA protocol approved

291     by the Institutional Biosafety Committee. Approval to obtain the de-identified clinical sample was

292     reviewed by the Human Subjects Institutional Review Boards at the University of Wisconsin –

293     Madison.

**Nucleic acid extraction**

295     For each sample, approximately 140 µL of viral transport medium or cell culture supernatant

296     was passed through a 0.22µm filter (Dot Scientific, Burton, MI, USA). Total nucleic acid was

297     extracted using the Qiagen QIAamp Viral RNA Mini Kit (Qiagen, Hilden, Germany), substituting

298     carrier RNA with linear polyacrylamide (Invitrogen, Carlsbad, CA, USA) and eluting in 30 µL of

299     nuclease free $H_2O$.  Samples were treated with TURBO DNase (Thermo Fisher Scientific,

300     Waltham, MA, USA) at 37°C for 30 min and concentrated to 8µL using the RNA Clean &

301     Concentrator-5 kit (Zymo Research, Irvine, CA, USA). Full protocol for nucleic acid extraction

302     and subsequent cDNA generation is available at https://www.protocols.io/view/sequence-

303     independent-single-primer-amplification-o-bckxiuxn.

**Complementary DNA (cDNA) generation**

305     Complementary DNA (cDNA) was synthesized using a modified Sequence Independent Single

306     Primer Amplification (SISPA) approach described by Kafetzopoulou et al. [14]. RNA was

307     reverse transcribed with SuperScript IV Reverse Transcriptase (Invitrogen, Carlsbad, CA, USA)

308     using Primer A (5'-GTT TCC CAC TGG AGG ATA-($N_9$)-3'). Reaction conditions were as follows:

309     1µL of primer A was added to 4 µL of sample RNA, heated to 65°C for 5 minutes, then cooled to

310     4 □ for 5 minutes. Then 5 µL of a master mix (2 µL 5x RT buffer, 1 µL 10 mM dNTP, 1 µL

311 nuclease free H$_2$O, 0.5 µL 0.1M DTT, and 0.5 µL SSIV RT) was added and incubated at 42□ for

312 10 minutes. For generation of second strand cDNA, 5 µL of Sequenase reaction mix (1 µL 5x

313 Sequenase reaction buffer, 3.85 µL nuclease free H$_2$O, 0.15 µL Sequenase enzyme) was added

314 to the reaction mix and incubated at 37°C for 8 minutes. This was followed by the addition of a

315 secondary Sequenase reaction mix (0.45 µl Sequenase Dilution Buffer, 0.15 µl Sequenase

316 Enzyme), and another incubation at 37□ for 8 minutes. Following incubation, 1µL of RNase H

317 (New England BioLabs, Ipswich, MA, USA) was added to the reaction and incubated at 37°C for

318 20 min. Conditions for amplifying Primer-A labeled cDNA were as follows: 5 µL of primer-A

319 labeled cDNA was added to 45 µL of AccuTaq master mix per sample (5 µL AccuTaq LA 10x

320 Buffer, 2.5 µL dNTP mix, 1µL DMSO, 0.5 µL AccuTaq LA DNA Polymerase, 35 µL nuclease

321 free water, and 1 µL Primer B (5′-GTT TCC CAC TGG AGG ATA-3′). Reaction conditions for the

322 PCR were: 98°C for 30s, 30 cycles of 94°C for 15 s, 50°C for 20 s, and 68°C for 2 min, followed

323 by 68°C for 10 min.

324 **Oxford nanopore library preparation and sequencing**

325 Amplified cDNA was purified using a 1:1 concentration of AMPure XP beads (Beckman Coulter,

326 Brea, CA, USA) and eluted in 48µL of water. A maximum of 1 µg of DNA was used as input into

327 Oxford Nanopore kits SQK-LSK109. Samples were barcoded using the Oxford Nanopore Native

328 Barcodes (EXP-NBD104 and EXP-NBD114), and pooled to a total of 140ng prior to being run

329 on the appropriate flow cell (FLO-MIN106) using the 72hr run script.

330 **Nextera XT Illumina library preparation and sequencing**

331 Amplified cDNA was purified using a 1:1 concentration of AMPure XP beads (Beckman Coulter,

332 Brea, CA, USA) and eluted in 48µL of water. PCR products were quantified using Qubit dsDNA

333 high-sensitivity kit (Invitrogen, USA) and were diluted to a final concentration of 0.2 ng/µl (1 ng

334 in 5 µl volume). Each sample was then made compatible for deep sequencing using the Nextera

335 XT DNA sample preparation kit (Illumina, USA). Specifically, each sample was enzymatically

336 fragmented and tagged with short oligonucleotide adapters, followed by 14 cycles of PCR for

337 template indexing. Samples were purified using two consecutive AMPure bead cleanups (0.5x

338 and 0.7x) and were quantified once more using Qubit dsDNA high-sensitivity kit (Invitrogen,

339 USA). The average sample fragment length and purity was determined using Agilent High

340 Sensitivity DNA kit and the Agilent 2100 Bioanalyzer (Agilent, Santa Clara, CA). After passing

341 quality control measures, samples were pooled equimolarly to a final concentration of 4 nM, and

342 5 µl of each 4 nM pool was denatured in 5 µl of 0.2 N NaOH for 5 min. Four samples (primary

343 NP swab, p1 Vero 76, p1 Vero E6, and p1 Vero STAT-1 KO) were pooled on a single flowcell to

344 a final concentration of 8pM with a PhiX-derived control library accounting for 1% of total DNA

345 and was loaded onto a 500-cycle v2 flowcell. The p2 samples (p2a Vero 76 and p2b Vero 76)

346 were pooled with seven other samples (not included in this manuscript) and were denatured to a

347 final concentration of 14pM with a PhiX-derived control library accounting for 1% of total DNA

348 and was loaded onto a 600-cycle v3 flowcell. Average quality metrics were recorded, reads

349 were demultiplexed, and FASTQ files were generated on Illumina's BaseSpace platform.

**Sequence read mapping and variant calling by ONT**

351 Seventy-two hours after sequencing was initiated, raw sequencing reads were demultiplexed

352 using qcat (https://github.com/nanoporetech/qcat). In order to deplete host sequences,

353 sequencing reads are mapped against host genome and transcript references, and unmapped

354 reads are saved. Reads were then trimmed by 30bp on each side to discard SISPA primer

355 sequences. In this step, reads with quality scores ≤ 7 were discarded. Cleaned viral reads were

356 then mapped to the severe acute respiratory syndrome coronavirus 2 isolate 2019-nCoV/USA-

357 WI1/2020 consensus sequence (Genbank: MT039887.1, originally sequenced by the CDC)

358 using minimap2. Minor variants from ONT sequences that comprise at least 10% of total

359 sequences in any of the samples were identified using the bbmap callvariants.sh tool

360    (https://jgi.doe.gov/data-and-tools/bbtools/). The entire ONT analysis pipeline is available at this

361    GitHib address https://github.com/katarinabraun/SARSCoV2_passage_MS.

**Illumina sequence data analysis – quality filtering and variant calling**

363    FASTQ files were initially processed using custom bioinformatic pipelines, available with

364    instructions for use at the GitHub repository accompanying this manuscript

365    https://github.com/katarinabraun/SARSCoV2_passage_MS. Briefly, read ends were trimmed to

366    achieve an average read quality score of Q30 and a minimum read length of 100 bases using

367    Trimmomatic (http://www.usadellab.org/cms/?page=trimmomatic) [28]. Paired-end reads were

368    merged and then mapped to the reference sequence (Genbank MT039887.1: 2019-nCoV/USA-

369    WI1/2020) using Bowtie2 (http://bowtie-bio.sourceforge.net/bowtie2/manual.shtml). Single

370    nucleotide variants (SNVs) were called with Varscan2 (http://varscan.sourceforge.net/using-

371    varscan.html) using a frequency threshold of 1%, a minimum coverage of 100 reads, and a

372    base quality threshold of Q30 or higher [29]. SNVs were annotated to determine the impact of

373    each variant on the amino acid sequence. SNVs were annotated in eleven open reading frames:

374    ORF1a (open reading frame 1a), ORF1b (open reading frame 1b), S (Spike, encodes surface

375    protein), ORF3a (open reading frame 3a), E (envelope), M (membrane), ORF6 (open reading

376    frame 6), ORF7a (open reading frame 7a), ORF8 (open reading frame 8), N (nucleocapsid),

377    ORF10 (open reading frame 10). VCF files were cleaned for additional analyses and figure-

378    generation using custom Python scripts, which are all available at the GitHub repository

379    accompanying this manuscript.

**Illumina sequence data analysis – diversity statistics**

381    Nucleotide diversity was calculated using π summary statistics. π quantifies the average

382    number of pairwise differences per nucleotide site among a set of sequences and was

383    calculated per gene using SNPGenie (https://github.com/chasewnelson/SNPgenie) [30].

384    SNPGenie adapts the Nei and Gojobori method of estimating nucleotide diversity ($\pi$), and its

385    synonymous ($\pi_S$) and nonsynonymous ($\pi_N$) partitions from next-generation sequencing data

386    [31]. As most random nonsynonymous mutations are likely to be disadvantageous, we expect

387    $\pi_N = \pi_S$ indicates neutrality suggesting that allele frequencies are determined primarily by

388    genetic drift. $\pi_N < \pi_S$ indicates purifying selection is acting to remove new deleterious

389    mutations, and $\pi_N > \pi_S$ indicates diversifying selection is favoring new mutations and may

390    indicate positive selection is acting to preserve multiple amino acid changes [32].

391    **Approvals**

392    *Biosafety*. Work with live virus was performed at biosafety level-3 containment at the Influenza

393    Research Institute at the University of Wisconsin – Madison under a recombinant DNA protocol

394    approved by the Institutional Biosafety Committee.

395    *Human subjects*. Approval to obtain the de-identified clinical sample was reviewed by the

396    Human Subjects Institutional Review Boards at the University of Wisconsin – Madison.

397    **Data availability**

398    Metagenomic sequencing data after mapping to SARS-COV-2 reference genome (MT039887.1)

399    have been deposited in the Sequence Read Archive (SRA) under bioproject PRJNA607948.

400    Derived data, analysis pipelines, and figures have been made available for easy replication of

401    these        results        at        a        publicly-accessible        GitHub        repository:

402    https://github.com/katarinabraun/SARSCoV2_passage_MS. A description of these results is

403    also available on LabKey at go.wisc.edu/qca2m5.

404    **Figure generation**

405    Figures 3, 4, 5, 6 and supplemental figures 2 and 3 were generated using custom Python scripts

406    and Matplotlib (https://matplotlib.org/). All code to replicate these figures can be found in the

407    GitHub repository. Figure 1 was created with BioRender (https://biorender.com/). Supplemental

408    figure 1 was created with JMP (https://www.jmp.com/)

409

410    **References**

411    1.    Lu H, Stratton CW, Tang YW. Outbreak of pneumonia of unknown etiology in Wuhan,

412          China: The mystery and the miracle. J Med Virol. 2020;92:401-402.

413    2.    Li Q, Guan X, Wu P et al. Early Transmission Dynamics in Wuhan, China, of Novel

414          Coronavirus-Infected Pneumonia. N Engl J Med. 2020;382:1199-1207.

415    3.    Chen N, Zhou M, Dong X et al. Epidemiological and clinical characteristics of 99 cases of

416          2019 novel coronavirus pneumonia in Wuhan, China: a descriptive study. Lancet.

417          2020;395:507-513.

418    4.    Huang C, Wang Y, Li X et al. Clinical features of patients infected with 2019 novel

419          coronavirus in Wuhan, China. Lancet. 2020;395:497-506.

420    5.    Wang C, Horby PW, Hayden FG, Gao GF. A novel coronavirus outbreak of global health

421          concern. Lancet. 2020;395:470-473.

422    6.    Hadfield J, Megill C, Bell SM et al. Nextstrain: real-time tracking of pathogen evolution.

423          Bioinformatics. 2018;34:4121-4123.

424    7.    Chan JF, Kok KH, Zhu Z et al. Genomic characterization of the 2019 novel human-

425          pathogenic coronavirus isolated from a patient with atypical pneumonia after visiting

426          Wuhan. Emerg Microbes Infect. 2020;9:221-236.

427    8.    Lu R, Zhao X, Li J et al. Genomic characterisation and epidemiology of 2019 novel

428          coronavirus: implications for virus origins and receptor binding. Lancet. 2020;395:565-574.

429    9.    Chinese SARSMEC. Molecular evolution of the SARS coronavirus during the course of the

430          SARS epidemic in China. Science. 2004;303:1666-1669.

431    10.   Karamitros T, Papadopoulou G, Bousali M, Mexias A, Tsiodras S, Mentis A. SARS-CoV-2

432         exhibits intra-host genomic plasticity and low-frequency polymorphic quasispecies.

433         BioRxiv. 2020

434    11.   Shen Z, Xiao Y, Kang L et al. Genomic diversity of SARS-CoV-2 in Coronavirus Disease

435         2019 patients. Clin Infect Dis. 2020

436    12.   Liu Z, Zheng H, Yuan R et al. Identification of a common deletion in the spike protein of

437         SARS-CoV-2. BioRxiv. 2020

438    13.   Lewandowski K, Xu Y, Pullan ST et al. Metagenomic Nanopore Sequencing of Influenza

439         Virus Direct from Clinical Respiratory Samples. J Clin Microbiol. 2019;58

440    14.   Kafetzopoulou LE, Efthymiadis K, Lewandowski K et al. Assessment of metagenomic

441         Nanopore and Illumina sequencing for recovering whole genome sequences of

442         chikungunya and dengue viruses directly from clinical samples. Euro Surveill. 2018;23

443    15.   Wilker PR, Dinis JM, Starrett G et al. Selection on haemagglutinin imposes a bottleneck

444         during mammalian transmission of reassortant H5N1 influenza viruses. Nat Commun.

445         2013;4:2636.

446    16.   Moncla LH, Zhong G, Nelson CW et al. Selective Bottlenecks Shape Evolutionary

447         Pathways Taken during Mammalian Adaptation of a 1918-like Avian Influenza Virus. Cell

448         Host Microbe. 2016;19:169-180.

449    17.   Moncla LH, Bedford T, Dussart P et al. Quantifying within-host diversity of H5N1 influenza

450         viruses in humans and poultry in Cambodia. PLoS Pathog. 2020;16:e1008191.

451    18.   Taiaroa G, Rawlinson D, Featherstone L et al. Direct RNA sequencing and early evolution

452         of SARS-CoV-2. BioRxiv. 2020

453    19.   Smith EC, Denison MR. Coronaviruses as DNA wannabes: a new model for the regulation

454         of RNA virus replication fidelity. PLoS Pathog. 2013;9:e1003760.

455    20.   Zhao L, Illingworth CJR. Measurements of intrahost viral diversity require an unbiased

456         diversity metric. Virus Evol. 2019;5:vey041.

bioRxiv preprint doi: https://doi.org/10.1101/2020.04.20.051011; this version posted April 20, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

457  21.  Grubaugh ND, Gangavarapu K, Quick J et al. An amplicon-based sequencing framework

458       for accurately measuring intrahost virus diversity using PrimalSeq and iVar. Genome Biol.

459       2019;20:8.

460  22.  Duggal NK, McDonald EM, Weger-Lucarelli J et al. Mutations present in a low-passage

461       Zika virus isolate result in attenuated pathogenesis in mice. Virology. 2019;530:19-26.

462  23.  Sanjuán R, Nebot MR, Chirico N, Mansky LM, Belshaw R. Viral mutation rates. J Virol.

463       2010;84:9733-9748.

464  24.  Xue KS, Moncla LH, Bedford T, Bloom JD. Within-Host Evolution of Human Influenza

465       Virus. Trends Microbiol. 2018;26:781-793.

466  25.  Xu X, Cox NJ, Bender CA, Regnery HL, Shaw MW. Genetic variation in neuraminidase

467       genes of influenza A (H3N2) viruses. Virology. 1996;224:175-183.

468  26.  Zhao Z, Li H, Wu X et al. Moderate mutation rate in the SARS coronavirus genome and its

469       implications. BMC Evol Biol. 2004;4:21.

470  27.  Smith BL, Wilke CO. A new twist in measuring mutation rates. Elife. 2017;6

471  28.  Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence

472       data. Bioinformatics. 2014;30:2114-2120.

473  29.  Koboldt DC, Zhang Q, Larson DE et al. VarScan 2: somatic mutation and copy number

474       alteration discovery in cancer by exome sequencing. Genome Res. 2012;22:568-576.

475  30.  Nelson CW, Moncla LH, Hughes AL. SNPGenie: estimating evolutionary parameters to

476       detect natural selection using pooled next-generation sequencing data. Bioinformatics.

477       2015;31:3709-3711.

478  31.  Nei M, Gojobori T. Simple methods for estimating the numbers of synonymous and

479       nonsynonymous nucleotide substitutions. Mol Biol Evol. 1986;3:418-426.

480  32.  Hughes AL. Adaptive evolution of genes and genomes. 1999

481

482

483    **Figure Captions**

484    **Figure 1. Sequence-Independent, Single-Primer Amplification sequencing workflow.** A)

485    Table showing nomenclature, and color scheme for all samples used in this study. B) Schematic

486    showing the sequence-independent, single-primer amplification sequencing workflow.

487

488    **Figure 2. Consensus sequence overview for SARS-CoV-2 samples.** A) Map of the SARS-

489    CoV-2 genome illustrating no consensus-level changes compared to the reference

490    (MT039887.1). B) Map of the Madison SARS-CoV-2 showing an in-frame deletion at nucleotide

491    position 20,298 - 20,300 relative to the Wuhan reference (MN908947.3).

492

493    **Figure 3. Minor variant frequencies in ORF1a, ORF1b, and Spike coding regions of the**

494    **SARS-CoV-2 genome.** A) Minor variants ≥1% frequency that were detected in the original

495    primary NP swab by Illumina sequencing in ORF1a, ORF1b, and spike genes. B) Minor variants

496    ≥1% frequency that were detected in the first passage by Illumina sequencing in ORF1a,

497    ORF1b, and spike genes.  C) Minor variants ≥1% frequency that were detected in the second

498    passage by Illumina sequencing in ORF1a, ORF1b, and spike genes.

499

500    **Figure 4. SNV frequency distributions.** The frequency of Illumina detected SNVs plotted

501    against a "neutral model", represented in light grey. The neutral model assumes a constant

502    population size and the absence of selection. A) SNV frequency spectrum from the primary NP

503    swab, represented in dark blue. B) SNV frequency spectrum from three p1 samples,

504    represented in turquoise. C) SNV frequency spectrum from two p2 samples, represented in dark

505    grey.

506

507    **Figure 5: Intragene nucleotide diversity.** Relative abundance of nonsynonymous (πN) and

508    synonymous (πS) for all 11 open reading frames. Nonsynonymous diversity (πN) is denoted by

509    closed symbols and synonymous diversity (πS) is denoted by open symbols. A) Intragene π

510    from the primary NP swab, represented in dark blue. B) Intragene π from three p1 samples,

511    represented in turquoise. C) Intragene π from two p2 samples, represented in dark grey. Length

512    of horizontal line is the difference between πN and πS for each gene.

513

514    **Figure 6. Comparison of ONT and Illumina SNV calls.** Concordance between SNV calls at

515    the same sites, irrespective of frequency, determined by Illumina and ONT workflows. Symbol

516    denotes sample and color denotes gene. Gene colors correspond to the genome map in Figure

517    2.

518

519    **Supporting Information**

520

521    **Supplemental Figure Captions**

522    **Supplemental figure 1. Coverage depth across the SARS-CoV-2 genome.** The relative

523    depth of coverage for each nucleotide position was plotted for (A) ONT and (B) Illumina

524    sequencing results.

525

526 **Supplemental figure 2. Change in SNV frequency over passage.** SNVs found shared across

527 the primary NP swab, p1 Vero 76 and p2a/p2b Vero 76 are plotted here. Symbol denotes the

528 specific SNV. Line-type denotes route: either swab → p1 Vero 76 → p2a Vero 76 (dashed) or

529 swab → p1 Vero 76 → p2a Vero 76 (solid). Color denotes the gene where the SNV was found.

530 (A) Y-axis is scaled to visualize all shared SNVs, ranging from 0 - 50% frequency. (B) Y-axis is

531 magnified to visualize SNV frequencies below 5%.

532

533 **Supplemental Figure 3. Minor frequency variants across the whole SARS-CoV-2 genome.**

534

535 **Supplemental Tables**

| Gene | Position in Gene | Reference nt | Variant nt | Annotation | Swab | | P1 Vero 76 | | P1 Vero E6 | | P1 Vero STAT-1 KO | | P2a Vero 76 | | P2b Vero 7( | |
|------|---------|------|------|------------|------|---------|------|---------|------|---------|------|---------|------|---------|------|--------|
| | | | | | ONT | Ilumina | ONT | Illumina | ONT | Ilumina | ONT | Illumina | ONT | Illumina | ONT | Illumin |
| ORF1a | 4191 | C | T | synonymous | | | | | | | 5.1 | 1.81 | | | | |
| ORF1a | 6548 | C | T | T2183I | | | | | | | 10.78 | 4.76 | | | | |
| ORF1a | 8089 | C | G | R2697G | nd | 12.71 | | | | | | | | | | |
| ORF1a | 10242 | C | T | synonymous | | | | | | | | | 24.22 | 17.34 | 15.63 | 12.09 |
| ORF1a | 11070 | G | T | synonymous | 16.86 | 20 | 13.21 | 17.78 | 15.53 | 20.21 | 27.91 | 37.05 | 9.75 | 15.0 | 7.52 | 11.42 |
| ORF1a | 11202 | G | T | W3734C | | | 5.95 | 7.59 | | | | | | | | |
| ORF1a | 11632 | C | A | Q3878K | | | | | | | | | 15.41 | 23.24 | 10.07 | 14.84 |
| ORF1b | 5415 | A | G | synonymous | | | | | | | 5.97 | 4.95 | | | | |
| ORF1b | 5843 | C | T | T2048I | | | | | | | | | | | 26.01 | 26.91 |

| | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Spike | 1640 | C | T | T547I | | | 5.64 | 4.8 | | | | | | | |
| ORF3a | 266 | C | T | T89I | | | | | 8.36 | 5.19 | | | | | |

536 **Supplemental Table 1. Comparison of ONT and Illumina SNVs.** 'nd' indicates that the
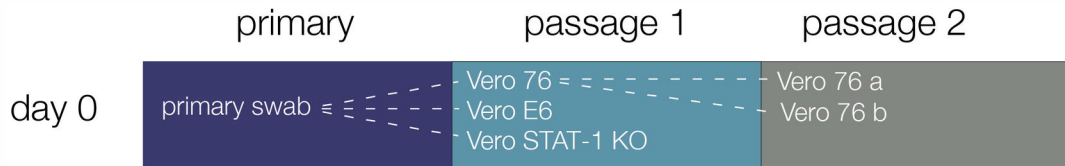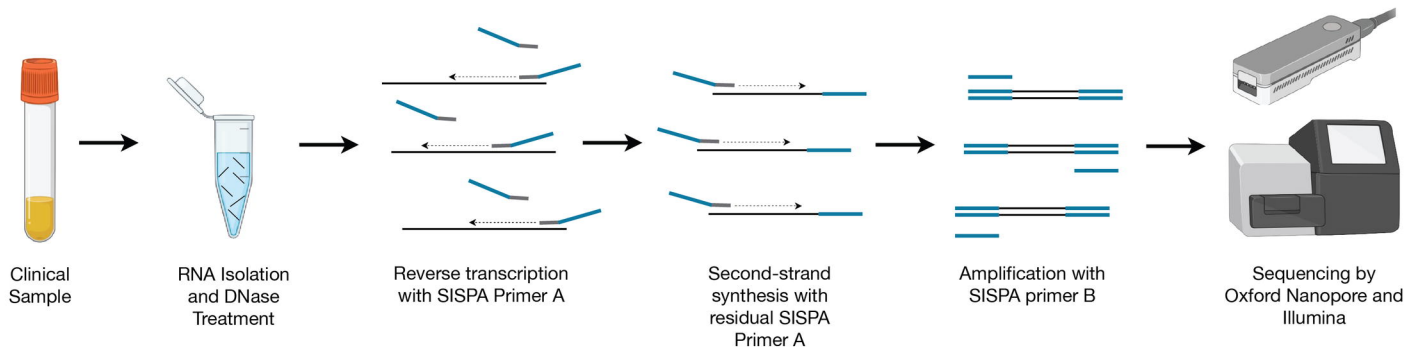
537 variant was not detected.

538

| Gene | Swab | P1 Vero 76 | P1 Vero E6 | P1 Vero STAT-1 KO | P2a Vero 76 | P2b Vero 76 |
|---|---|---|---|---|---|---|
| ORF1a | 0.6817 | 0.9689 | 0.8332 | 0.8332 | 1.666 | 1.439 |
| ORF1b | 0.6185 | 0.6185 | 0.2474 | 0.4948 | 0.7422 | 0.6185 |
| Spike | - | 1.0468 | 0.2617 | 1.3085 | 0.7851 | 0.5234 |
| ORF3a | - | 1.2091 | 1.2091 | 2.4183 | - | - |
| E | - | - | 4.4052 | 4.4052 | 8.8105 | - |
| M | - | - | - | - | - | - |
| ORF6 | - | - | - | - | - | - |
| ORF7a | - | 5.4794 | 5.4794 | 5.4794 | 8.2192 | 5.4794 |
| ORF8 | - | - | - | - | - | - |
| N | 0.7942 | - | - | - | 0.7942 | - |
| ORF10 | - | 8.6206 | - | - | - | - |

539 **Supplemental Table 2. Variants per gene kilobase length.** To normalize the number of SNVs

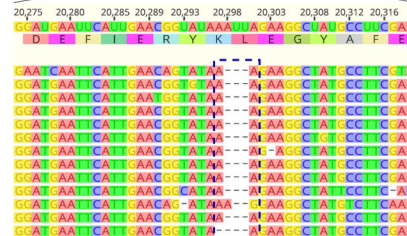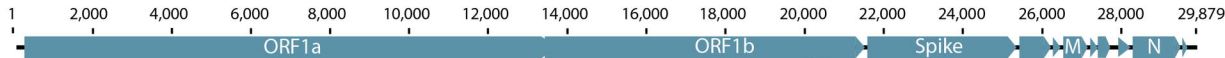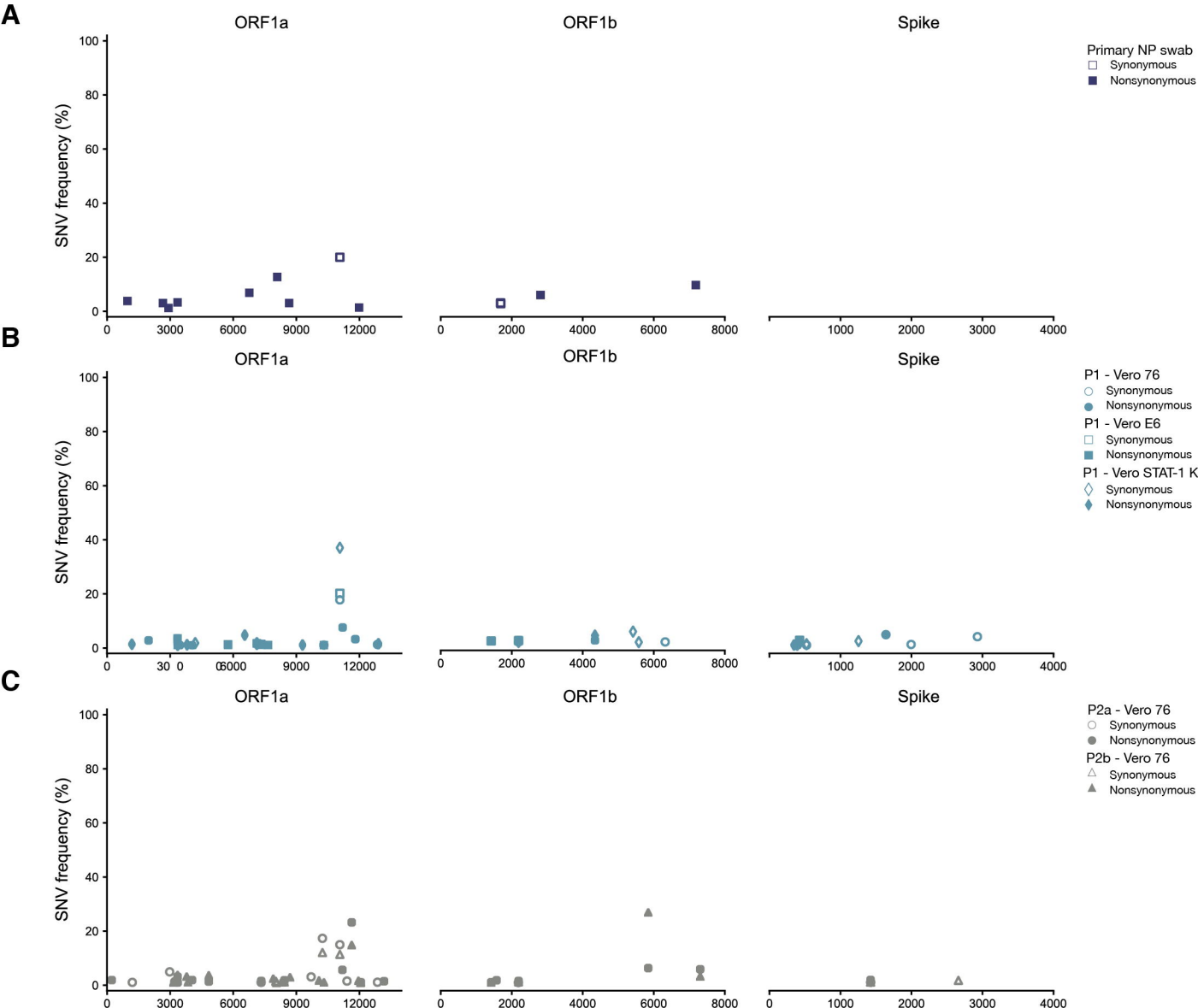540 per gene segment, we report the density of variants normalized to gene kilobase length.
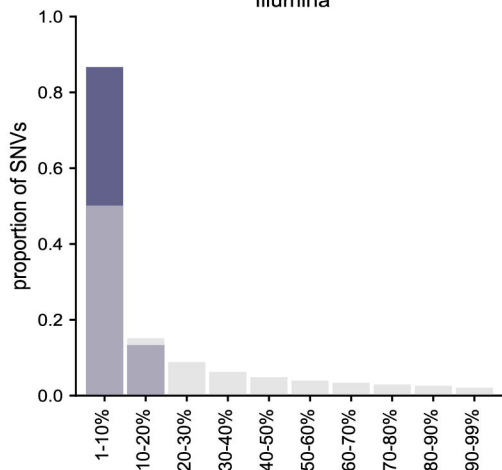
541

**A**

| | primary | passage 1 | passage 2 |
|---|---|---|---|
| day 0 | primary swab | Vero 76<br>Vero E6<br>Vero STAT-1 KO | Vero 76 a<br>Vero 76 b |

**B**

Clinical Sample → RNA Isolation and DNase Treatment → Reverse transcription with SISPA Primer A → Second-strand synthesis with residual SISPA Primer A → Amplification with SISPA primer B → Sequencing by Oxford Nanopore and Illumina

**A**

| | ORF1a | ORF1b | S | | M | | | N | |

Primary NP Swab

P1 - Vero 76

P1 - Vero E6

P1 - Vero STAT-1 KO

P2a - Vero 76

P2b - Vero 76

**B**

ORF1a    ORF1b    Spike    M    N

**A** Primary NP swab
Illumina

**B** P1 Vero 76
Illumina

**C** P2a Vero 76
Illumina

P1 Vero E6
Illumina

P2b Vero 76
Illumina

P1 Vero STAT-1 KO
Illumina

**A** Primary NP swab

- Primary swab, synonymous
- Primary NP swab, nonsynonymous

**B** P1 Vero 76 — P1 Vero E6 — P1 Vero STAT-1 KO

- P1 Vero 76, synonymous
- P1 Vero 76, nonsynonymous
- P1 Vero E6, synonymous
- P1 Vero E6, nonsynonymous
- P1 Vero STAT-1 KO, synonymous
- P1 Vero STAT-1 KO nonsynonymous

**C** P2a Vero 76 — P2b Vero 76

- P2a Vero 76, synonymous
- P2a Vero 76, nonsynonymous
- P2b Vero 76, synonymous
- P2b Vero 76 nonsynonymous

π (nucleotide diversity)