1    # Optimized Quantification of Intrahost Viral Diversity in SARS-CoV-2

2    # and Influenza Virus Sequence Data

3    Running title: Optimized Quantification of Intrahost Viral Diversity

4

5    Roder AE*[1], Johnson KEE*[2], Knoll M*[2], Khalfan M[2], Wang B[2], Schultz-Cherry S[3], Banakis S[1],

6    Kreitman A[1], Mederos C[1], Youn J-H[4], Mercado R[4], Wang W[1], Ruchnewitz D[5], Samanovic MI[6],

7    Mulligan MJ[6], Lassig M[5], Łuksza M[7], Das S[4], Gresham D[2#], Ghedin E[1,2#].

8

9    <u>Author Affiliations:</u>

10   [1]Systems Genomics Section, Laboratory of Parasitic Diseases, DIR, NIAID, NIH, Bethesda, MD

11   20894, USA

12   [2]Center for Genomics and Systems Biology, Department of Biology, New York University, New

13   York, NY, 10003, USA.

14   [3]Department of Infectious Diseases, St Jude Children Research Hospital, Memphis, TN

15   [4]Department of Laboratory Medicine, NIH, Bethesda, MD 20894, USA

16   [5]Institute for Biological Physics, University of Cologne, Cologne, Germany

17   [6]New York University Langone Vaccine Center, Department of Medicine, New York, NY 10016,

18   USA

19   [7]Department of Oncological Sciences, Icahn School of Medicine at Mount Sinai, New York, NY

20   10029, USA

21

22   *Contributed equally to this article, ordered by contributions to the manuscript

23   #Corresponding authors: elodie.ghedin@nih.gov; dgresham@nyu.edu

24    ABSTRACT

25         High error rates of viral RNA-dependent RNA polymerases lead to diverse intra-host viral

26    populations during infection. Errors made during replication that are not strongly deleterious to

27    the virus can lead to the generation of minority variants. However, accurate detection of minority

28    variants in viral sequence data is complicated by errors introduced during sample preparation and

29    data analysis. We used synthetic RNA controls and simulated data to test seven variant calling

30    tools across a range of allele frequencies and simulated coverages. We show that choice of

31    variant caller, and use of replicate sequencing have the most significant impact on single

32    nucleotide variant (SNV) discovery and demonstrate how both allele frequency and coverage

33    thresholds impact both false discovery and false negative rates. We use these parameters to find

34    minority variants in sequencing data from SARS-CoV-2 clinical specimens and provide guidance

35    for studies of intrahost viral diversity using either single replicate data or data from technical

36    replicates. Our study provides a framework for rigorous assessment of technical factors that

37    impact SNV identification in viral samples and establishes heuristics that will inform and improve

38    future studies of intrahost variation, viral diversity, and viral evolution.

39    IMPORTANCE

40          When viruses replicate inside a host, the virus replication machinery makes mistakes.

41    Over time, these mistakes create mutations that result in a diverse population of viruses inside

42    the host. Mutations that are neither lethal to the virus, nor strongly beneficial, can lead to minority

43    variants that are minor members of the virus population. However, preparing samples for

44    sequencing can also introduce errors that resemble minority variants, resulting in inclusion of false

45    positive data if not filtered correctly. In this study, we aimed to determine the best methods for

46    identification and quantification of these minority variants by testing the performance of seven

47    commonly used variant calling tools. We used simulated and synthetic data to test their

48    performance against a true set of variants, and then used these studies to inform variant

49    identification in data from clinical SARS-CoV-2 clinical specimens. Together, analyses of our data

50    provide extensive guidance for future studies of viral diversity and evolution.

51 INTRODUCTION

52      Large population sizes, high replication rates, and error prone polymerases all contribute

53 to the generation of sequence diversity found in viral infections (1-5). Natural selection acts on

54 this diversity, contributing to viral evolution. RNA viruses have some of the highest mutation rates

55 among viruses (1, 6, 7). To replicate their genomes, RNA viruses must encode their own RNA-

56 dependent RNA-polymerases (RdRp), which often lack proofreading capabilities. Coronaviruses

57 are a notable exception, as they possess a distinct protein with 3'-5' exonuclease capability (1, 8,

58 9). Most errors made during replication—up to 40% in RNA viruses—are lethal (10, 11). Beneficial

59 mutations make up a much smaller proportion, and these, along with neutral mutations, comprise

60 the substitution rate. This substitution rate can be used to estimate the viral evolutionary rate, an

61 important calculation in considering viral spread, pandemic potential, and vaccine design (4, 12).

62      Due to the large population sizes of RNA viruses, intrahost bottlenecks, and genetic drift,

63 genetic diversity within a host is dynamic, with frequencies of mutations constantly rising and

64 falling (13). Mutations can lead to changes in the consensus sequence, e.g., where the allele

65 frequency (AF) is greater than 50%, and these specific sets of mutations separate globally

66 circulating virus populations into clades. Mutations in the virus genomes that are not the majority

67 within an infected host (i.e., present at lower than 50% frequency) represent minority variants.

68 Deep sequencing enables the capture of intrahost variation, both at the majority and minority

69 level, enabling the identification of variants and estimation of their frequency. Studying intrahost

70 variation can help in tracking viral spread, estimating population bottleneck sizes, and identifying

71 key amino acid changes that differentiate new viral strains (14-17). Additionally, minority variants

72 can highlight regions of the genome under selection or regions with increased mutational

73 tolerance, as well as allow for detection of subtle population shifts within the infected host and

74 discovery of possible drug resistance mutations (18, 19). Thus, information gleaned from studying

75 intrahost viral diversity has major implications for vaccine, monoclonal antibody, and drug

76 development.

4

77   Given the many applications of studying intra-host viral diversity, accurately identifying

78   and quantifying viral variants is essential. Precise identification of viral variants, especially those

79   at low frequencies, is complicated by the fact that viral genome sequencing often requires reverse

80   transcription and amplification, which, along with library preparation and the sequencing process,

81   are error prone. Thus, distinguishing true sequence variation from technical and experimental

82   noise is challenging. Typically, several ad hoc metrics are used to filter variants, such as applying

83   frequency and coverage cutoffs to sequencing data, however, the frequency at which identified

84   variants are considered valid can vary widely (20-27). Most studies using large sample cohorts,

85   or performing analyses on publicly available data, generally use single replicate data, despite

86   evidence suggesting that replicate sequencing may be essential for filtering false positive minority

87   variants (27). Despite the large number of studies analyzing minority variants in virus data, there

88   is no consensus on what coverage cutoffs and allele frequency cutoffs to use, and no large-scale

89   studies have been performed to determine what thresholds lead to the highest confidence in

90   variant identification.

91   In addition to the diversity of cutoffs used for single nucleotide variant (SNV) identification,

92   there is also great diversity in the variant calling software available. Variant callers are often

93   designed with specific functions in mind, such as identifying germline or somatic mutations in

94   cancer genomes or single nucleotide variants in viral populations (28, 29). The function for which

95   a variant caller is designed can have a significant impact on the statistics used and assumptions

96   made by the software. Tools designed for detection of germline mutations, such as

97   HaplotypeCaller and freebayes, must consider the very large reference genome, higher frequency

98   variants, the diploid nature of the genome, the possibility of copy number variation, and long

99   repetitive regions or large insertions or deletions (30-35). In these instances, local realignment of

100  haplotypes may be most effective (28). By contrast, software used for somatic mutations in

101  tumors, such as Mutect2 and Varscan, or for viral diversity, such as iVar and timo (a variant caller

102  developed in our lab), use base by base comparisons, or a combination of this with haplotype-

103 based alignment, to find lower frequency variants (27, 33-36). These tools also may need

104 alternative methods to reduce false positive calls to account for PCR errors introduced during

105 amplification of the viral genome (29). Due to the differences in bioinformatic and statistical

106 approaches used by each variant calling tool, identifying the tool that is the best fit for the specific

107 research question being studied is essential. Some tools have been tested in pairwise

108 comparisons (27, 35), however little work has been done to extensively test the performance of

109 the many available tools on different viruses, across sequence coverages, and at various allele

110 frequencies in viral deep sequencing data.

111  Here, we tested seven variant callers on simulated, synthetic, and clinical deep

112 sequencing data. We tested each tool across a range of coverages, allele frequencies, and

113 experimental designs to determine the optimal parameters that should be used to decrease false

114 positive variant identification, without sacrificing true positive data. To compare performance

115 between a small RNA virus with a high mutation rate, and a large RNA virus with proofreading

116 capability, we tested the variant callers on two viruses of particular interest in the viral diversity

117 field, influenza and SARS-CoV-2. We find that choice of variant caller, and use of replicate

118 sequencing have the most significant impact on SNV discovery and demonstrate how both allele

119 frequency and coverage thresholds impact both false discovery and false negative rates. We also

120 provide guidance on best practices for leveraging deep sequencing data from public repositories

121 for intrahost studies. These analyses provide a resource for studies aiming to assess intrahost

122 viral diversity in SARS-CoV-2 or influenza, and they lay the groundwork for similar studies in other

123 viruses.

124

125    MATERIALS AND METHODS

126    Extended methods are available in the supplementary materials.

127    Generation of Simulated Data

128         Reads were simulated using NEAT (v2.0) by constructing a mutation, error, fragment

129    length, and GC model for each viral type (37). The models were provided to NEAT genReads.py

130    along with reference fasta files and a mutation rate of 0.009 (0.9%) for influenza and 0.0045

131    (0.45%) for SARS-CoV-2 to produce a "golden VCF" containing a defined number of SNVs in

132    each virus. Simulated random PCR errors were also added to each replicate using

133    genReads.py (NEAT). Several copies of the replicate golden VCFs were made, each with the

134    same variants but with differing allele frequencies (AF). These VCFs were used to simulate

135    paired end fastq libraries at 100,000X genome coverage and down-sampling was used to

136    simulate lower coverages.

137         Sequences were trimmed using trimmomatic v0.36 (38), aligned to the respective

138    reference genome with BWA mem v0.7.17 (39), and duplicate reads were marked using GATK

139    MarkDuplicatesSpark v4.1.7.0 (40). Variants were called in each replicate with seven different

140    tools, using multiple parameter configurations for each tool (**Table S1).** A VCF file containing the

141    intersection of the two replicates was generated using bcftools isec (v1.9) (39). The pipeline used

142    for data simulation, sequence processing, variant calling, and analysis is available at

143    https://github.com/gencorefacility/MAD2.

144    Synthetic RNA generation, library preparation, and data processing

145         Synthetic, 'wildtype' (WT), influenza genomic RNA and variant RNA (created by adding

146    18, 14 and 14 known nucleotide changes into the WT PB2, HA, and NA segments respectively)

147    were synthesized as double-stranded DNA (gBlocks). The DNA was *in vitro*-transcribed with the

148    HiScribeTM T7 High Yield RNA Synthesis Kit (Invitrogen). RNA samples were diluted to an equal

149    copy number concentration of $6x10^8$ copies/μL. WT and variant RNA were mixed at equal

150     molarity. The pools were mixed at various frequencies (50%, 25%, 12.5%, 6.25%, 3.13%, 1.56%,

151     0.78%, 0.39%), and diluted to various copy number concentrations ($6x10^6$ - $6x10^3$ copies/µL).

152     cDNA was generated and libraries were prepared using the Nextera XT library preparation

153     kit (Nextera), with all volumes scaled down to 0.25x of the manufacturer's instructions, cleaned

154     with AMPure beads, and pooled at equal molarity. Libraries were sequenced on the Miseq 300

155     Cycle v2 using 2 x 75 pair-end reads. Samples were amplified and sequenced in duplicate and

156     analyzed with the pipeline described above, with the addition of adapter trimming. Synthetic

157     SARS-CoV-2 data from a similarly designed study was downloaded from SRA (PRJNA682212)

158     and processed as above (26).

159     SARS-CoV-2 clinical sample preparation, processing, and variant calling

160     Total RNA was extracted from 300 µL of nasopharyngeal (NP) or mid-turbinate (MT)

161     swabs collected at the NIH Clinical Center as part of diagnostic testing between 07/24/2020 and

162     03/31/2021 (**Table S2**). All samples were de-identified and anonymized.

163     RNA from samples was extracted using the NUCLISENS easyMAG automated nucleic

164     acid extractor and the viral genome was amplified using a modified version of the ARTIC protocol

165     (https://artic.network/ncov-2019)          and          the          methods          described          at

166     https://github.com/GhedinSGS/SARS-CoV-2_analysis. All libraries were prepared as above and

167     sequenced on either the Illumina MiSeq or the Illumnia NextSeq500 using either the 2x150 bp or

168     2x300 bp paired end protocol. All samples were processed in duplicate.

169     Samples were processed with the pipeline available and described above, with the

170     addition of merging the two SAM files (from A and B primer pools) for each biological sample into

171     one alignment file using Picard Tools MergeSamFiles v2.17.11. Variants were called as above

172     using the standard parameters for each tool (**Table S1**).

173     Data Availability

174     Synthetic influenza data (bioproject PRJNA865369) and SARS-CoV-2 data from clinical

175     samples are available in NCBI GenBank and SRA. Accession IDs can be found in **Table S2**. All

8

176    downstream analysis files are available at https://github.com/GhedinSGS/Optimized-

177    Quantification-of-Intrahost-Viral-Diversity.

178    Ethics Statement

179        All samples were anonymized and obtained with consent as part of SARS-CoV-2

180    diagnostic testing.

181    RESULTS

182    *Simulated and synthetic data provide a 'true' set of minority variants to assess variant caller*

183    *performance.*

184            To test the ability of each variant caller to accurately identify variants, it is essential to

185    know the 'true set' of variants within the data, which is not possible with real sequence data. With

186    this in mind, we tested the ability of six popular variant-calling software packages (Freebayes,

187    HaplotypeCaller (hc), iVar, Lofreq, Mutect2, and Varscan) and one in-house pipeline (timo) to

188    accurately identify minority variants in simulated and synthetic sequencing data (**Fig. S1**) (19-24).

189    Single nucleotide variants (SNVs) were simulated across three influenza virus genomes (A/H1N1,

190    A/H3N2, B/Victoria) and one coronavirus genome (SARS-CoV-2) at both defined and random

191    allele frequencies and across a range of downsampled coverages (**Fig. S1A-B**). Further, synthetic

192    RNA controls of three influenza virus segments (PB2, HA, and NA) containing known SNVs were

193    mixed in varying amounts at various dilutions to create a range of allele frequencies and genome

194    copy numbers (**Methods**) (**Fig. S1C-D**). Combined, we used these synthetic and simulated data

195    sets to test variant caller performance on a known set of SNVs.

196            We found that all callers performed poorly on our data using their default parameters (**Fig.**

197    **S2A**). Therefore, to compare all callers equally, we used a standard set of permissive input

198    parameters (min coverage = 1x, allele frequency cutoff = 0.01 (1%)) throughout our testing (**Table**

199    **S1**). When assessing the F1 statistic across a range of simulated frequencies, most variant callers

200    performed well at low frequencies (< 0.05 (5%)) when the coverage was high (downsampling

201    fraction >= 0.005 or ~500X coverage). Conversely, high frequencies were necessary for accurate

202    variant detection at small downsampling fractions where the average coverage was low (**Fig. 1A**).

203    We did not observe significant differences between the four viruses in terms of variant calling

204    accuracy (**Fig. 1A, Fig. S2A**). Most of the differences in performance between the variant callers

205    could be seen at downsampling fractions < 0.005 and allele frequencies below 0.05 (5%) (**Fig**

206    **1A**). A closer look at precision and recall for each tool at downsampling fractions of 0.002 and

207    0.003 indicated that many tools trade recall for precision at low frequencies (**Fig. 1B**). Some tools,

208    such as iVar, timo, and varscan tended to be extremely conservative, especially at low

209    frequencies. However, the stringency of these tools can be reduced by decreasing the input

210    frequency parameter from 0.01 (1%) to 0.001 (0.1%) (**Fig. S2** *custom input parameters*, **Table**

211    **S1**). The use of simulated data for initial testing showed that under ideal conditions, all variant

212    callers have the ability to perform well on viral sequence data. This is important for establishing

213    baseline performance against which to compare the variant caller performance on synthetic data

214    and data from clinical specimens.

215        While simulated data provides an ideal set of variants against which to compare the variant

216    calls for each tool, it lacks the reverse transcription, amplification, and sequence library

217    preparation steps involved in the generation of data from clinical specimens. To assess how these

218    sample preparation steps, along with duplicate sequencing, and SNV thresholds may impact

219    variant caller performance, we tested each tool on data from the synthetic RNA dataset (**Fig. S1C-**

220    **D**). The average read depth across gene segments was greater than 1000x and had similar

221    coverage distributions to our simulated datasets at downsampling fractions of 0.01 and 0.1 (**Fig.**

222    **S3A**). At this coverage in simulated data, we observed high F1, precision, and recall scores across

223    all variant callers for SNVs > 1% frequency (**Fig. 1A-B**). Observed frequencies differed from

224    expectation (**Fig. S3B**), likely due to mixing errors during sample preparation, but were consistent

225    between viral copy numbers. Comparing the observed allele frequency to the median observed

226    allele frequency revealed that most callers agreed on the frequency of identified variants;

227    however, there is considerable variance in the AF estimations despite the fact that all SNVs are

228    linked (**Fig. 1C, S3B**). We also tested the variant callers on a previously published synthetic

229    SARS-CoV-2 data set (26). Here, HaplotypeCaller and timo found more true variants than the

230    other callers, especially at higher viral load (**Fig. S3C**). Across all callers, the frequencies of

231    observed variants in the SARS-CoV-2 data were less accurate than in the synthetic flu data i.e.

232    the variance in allele frequencies is higher) (**Fig. S3B-C**). This is potentially due to the overall low

233     coverage of these samples and the amplicon-based sequencing method required for larger RNA

234     genomes, suggesting that allele frequency estimation may be less reliable in SARS-CoV-2 data.

235

236     *Frequency thresholds and sequencing replicates reduce false positive SNVs*

237     Previous studies have reported the necessity of establishing frequency and coverage

238     thresholds as well as having replicate sequencing to decrease false-positive SNVs in a data set

239     (27, 40). Given that most publicly available data consist of single replicate sequencing data, we

240     aimed to establish coverage and frequency thresholds that would minimize the false discovery

241     rate (FDR) and false negative rate (FNR) to levels comparable to those observed using two

242     replicates. To do this, we used both simulated and synthetic datasets with standard input

243     parameters (**Fig. S2A**) and ignored the 'binocheck' requirement from timo (requiring variants to

244     be found in near equal numbers of forward and reverse reads), allowing us to test the performance

245     of timo on low frequency SNVs.

246     In both the synthetic and simulated data, false positive SNVs were found across read

247     depths but were primarily at allele frequencies less than 0.03 (**Fig. 2A-2B**). Therefore, applying

248     frequency thresholds to single replicate data lowered the false discovery rate (FDR) for all callers

249     (**Fig. 3A, S4A**). While establishing coverage cutoffs did not drastically impact the number of false

250     positive calls in either the simulated or synthetic dataset, coverage and library size are important

251     factors when considering SNV recall (**Fig. 1A, 2A-B, S5**). However, using frequency thresholds

252     did come at the cost of significantly increasing the FNR, as true SNVs found at low frequencies

253     were filtered from the data. In contrast, keeping only SNVs shared between the two replicates

254     dramatically decreased the FDR, while maintaining relatively low FNRs (**Fig. 3, S4**).

255     HaplotypeCaller, LoFreq, and Mutect2 called notably higher numbers of false positive SNVs in

256     synthetic data, including many that were maintained even after merging replicates—indicating

257     that these callers are making consistently incorrect SNV calls. Further, these three callers had

258     multiple instances where true positive SNVs were identified at high frequencies (AF > 0.05) in one

259    replicate and were entirely absent in the other (**Fig. S6A**). Importantly, replicates also increased

260    the accuracy of allele frequency estimation of true positive variants in simulated data. The effect

261    of sequencing replicates on allele frequency is especially pronounced for low coverage data,

262    where the percent error of allele frequency estimation is pointedly lower for all tools when using

263    replicates (**Fig. S6B**). This effect was less pronounced in synthetic data, likely due to the high

264    coverage of these samples, again demonstrating the importance of read coverage for accurate

265    allele frequency estimation (**Fig. S6C**).

266          Together, these results indicate that using replicate sequencing with less stringent

267    frequency cutoffs (AF >= 0.01) is the best combination to reduce the FDR while maintaining a low

268    FNR. However, when replicate sequencing is unavailable, high coverage across the genome and

269    strict frequency cutoffs (AF >= 0.03) are necessary.

270

271    *Choice of variant caller significantly impacts set and frequency of identified variants in real SARS-*

272    *CoV-2 data using single replicate data.*

273          While simulated and synthetic data allow for testing minority variant callers and cutoffs in

274    a controlled setting, real data will always be more unpredictable. Thus, after using simulated and

275    synthetic data to assess variant caller performance across frequencies and coverages, we tested

276    how the callers performed on SARS-CoV-2 sequence data from diagnostic samples. Based on

277    the simulated and synthetic data testing, we determined that a coverage cutoff of 200X and an

278    allele frequency cutoff of 0.03 in single replicate data minimized false positive calls without

279    sacrificing large amounts of true positive data with most variant calling tools. To test the variant

280    calling tools on high-quality data, we used only samples where at least 80% of the genome had

281    a read depth over 200x coverage cutoff in both sequencing replicates (**Fig. S7A**). We used each

282    variant calling tool to identify minority variants in these samples and filtered them using a read

283    depth cutoff of 200 and an allele frequency cutoff of 0.03.

284       We were interested in how similar the set of variants was that was identified by each caller.

285      As a proof of principle, we filtered the set of variants for those present above an allele frequency

286      of 0.5 and at read depths greater than 5x to identify consensus changes (AF >=50%, or major

287      variants) within the data. As expected, the tools largely agreed on the consensus changes within

288      the data (**Fig. S7B**). There was a small set of major variants that the callers did disagree on,

289      however, most of which were a result of differences in the way some callers identify indels or

290      handle variant at consecutive nucleotide positions. For the purposes of this study, indels were

291      excluded from the analysis. These data indicate that even at high allele frequencies, the variant

292      callers disagree to some extent on the set of variants present in clinical data, an important

293      consideration when choosing how to define consensus sequences from SARS-CoV-2 data.

294       We then analyzed the intersection of the minority variants (allele frequencies between

295      0.03 and 0.5) identified by each tool. The total number of variants identified varied greatly between

296      the callers, with Varscan calling the fewest variants by far, followed by timo and Lofreq, in line

297      with the more conservative nature of these callers observed in the previous analyses (**Fig. 4A**).

298      Of note, we found that replicate 2 data had much higher numbers of minority variants, particularly

299      at very low frequencies, regardless of Ct value or date of sequencing. This suggests that freeze

300      thawing samples may impact minor variant numbers (**Fig. S7C**). When comparing the set of

301      minority variants identified by each of the seven tools, there was significant disagreement

302      between the variants. Mutect2 and HaplotypeCaller identified many variants that other callers did

303      not, particularly in replicate 1, and missed several variants identified by the other callers. This was

304      similar to the performance of these callers on the synthetic data sets (**Fig. S7D**). Given the high

305      number of false positives identified by HaplotypeCaller, Mutect2 and Lofreq in the simulated and

306      synthetic datasets, we focused on the intersection of minority SNVs found in just the other four

307      variant callers: Freebayes, iVar, Varscan, and timo. Of all the minority variants found in the data,

308      104 from replicate 1, and 142 from replicate 2 were identified by all four of the variant callers (**Fig.**

309   **4B**). Overall, choice of variant caller appears to have a significant impact on the set of minority

310   variants identified in SARS-CoV-2 data from clinical specimens.

311   Many studies of minority variants investigate frequency of minority variants to calculate

312   selection, bottleneck size, and potential for transmission (24, 41). We were interested in how well

313   the variant callers agreed on the frequency at which variants were identified. We plotted the

314   frequency of a variant in one caller against the frequency in each other caller and found that most

315   of the minority variant callers were strikingly similar in their frequency calls of shared variants.

316   Timo, Freebayes and iVar all showed almost complete agreement on frequency of the variants

317   (**Fig. 4C**). Varscan showed more variation in frequency, generally calling variants at a lower

318   frequency than the other three tools (**Fig. 4C**). Of interest, variants called by one caller but not

319   another spanned a frequency range of 0.03 all the way to 0.5, indicating that even high frequency

320   minority variants were often not agreed upon by variant callers. These data show that choice of

321   variant caller not only affects the set of the minority variants that are identified in a data set, but

322   also the frequency of those variants.

323

324   *Most minority variants in data from SARS-CoV-2 clinical specimens are not reproducible across*

325   *sequencing replicates*

326   In our sequencing data, the number of variants identified in each replicate by each tool

327   was markedly different, suggesting that many of the identified minor variants may not be true

328   variants introduced through viral replication, but instead technical artifacts (**Fig. 4A, S7C-D**). As

329   was shown with our simulated and synthetic data, errors introduced through PCR, library

330   preparation, and sequencing are mostly random, and therefore less likely to reappear and be

331   identified across multiple sequencing replicates, particularly when using Freebayes, iVar, timo, or

332   Varscan (**Fig. 3**). To find high confidence minority variants, we looked at the intersection of

333   variants between the two replicates using each caller, using a lower allele frequency threshold,

334   0.01, as established in synthetic data for merged replicates (**Fig. 3**). iVar and Freebayes called

335  the highest number of reproducible variants, while timo called the fewest number of reproducible

336  variants, however the percentage of reproducible variants compared to variants identified in single

337  replicates was the highest (**Fig. 5A**). Timo is one of the most conservative callers that was tested;

338  these data indicate that being conservative may lead to increased reproducibility and therefore

339  increased confidence when used on single replicate data. It is however important to note that the

340  relatively low percentages of reproducible variants are likely skewed by the high numbers of low

341  frequency variants found in replicate 2 (**Fig. 5A, S7C**). When we looked at the intersection of only

342  the variants found by each tool in both replicates, only 80 SNVs were found by all callers across

343  replicates, suggesting again that variant callers do not agree on the set of minority variants

344  present (**Fig. 5B**). Together, these data suggest that most minority variants are not reproducible

345  across replicates and support the idea that more than any other criteria, sequencing replicate has

346  the highest impact on the set of minority variants identified (**Fig. 5B, S7C-D**).

347

348  *Variants identified by all variant callers show the most reproducible frequencies*

349  Using synthetic data, we showed that in a controlled setting, SNVs that were found in both

350  sequencing replicates generally showed reproducible frequencies (**Fig. 2C**). Given that frequency

351  is an important metric in most analyses performed using minority variant data, we wanted to test

352  if this held true in clinical samples. Plotting the frequency in replicate 1 against the frequency in

353  replicate 2 revealed that while some variants showed consistent frequencies, some differed

354  drastically—identified at 5-10% in one replicate and as high as 45-50% in the other replicate (**Fig.**

355  **5C**). These data were striking as they reveal that averaging frequency across replicates, or

356  performing only one sequencing replicate, could drastically alter downstream analyses performed

357  using these numbers. Interestingly, when we looked at the variants that were reproducible across

358  replicates and found by most, or all the variant callers, frequency tended to be much more

359  consistent than those identified only in a single replicate, or by a single caller (**Fig. 5C,** dark red

360    points). Together, these data suggest that confidence in each variant and its frequency is

361    increased with replicate sequencing and identification by many variant callers.

362         Since replicate sequencing data is not always available, we investigated what frequency

363    cutoff could be applied such that single replicate data closely resembled the merged replicate

364    data. To do this, we looked at the intersection of SNVs called in both replicates by Freebayes,

365    iVar, timo, and Varscan (80 variants shown in **Fig. 5B**) and compared those to the intersection of

366    SNVs called by the same four callers in each individual replicate (**Fig. 4A**). We then applied allele

367    frequency cutoffs between 0.01 and 0.1 to determine the best cutoff for use on single replicate

368    data. Here, we identify a true positive as a variant present in the reproducible set, and a false

369    positive as any other variant found in a single replicate. As was noted previously, we find that

370    replicate 2 data shows an increased number of SNVs, perhaps due to freeze/thawing of samples

371    between preparations (**Fig. S7C, 5D**). As such, replicate 1 is likely more representative of what

372    single replicate data may typically look like. At an allele frequency cutoff of 0.01, all true positives

373    were found, but the number of false positives was very high, while a frequency cutoff of 0.05 or

374    0.1 removed an outsized number of true positives from the dataset (**Fig. 5D**). Based on these

375    data, we suggest an allele frequency cutoff of 0.03 when only single replicate data is available, a

376    cutoff that was also confirmed in the simulated and synthetic data sets (**Fig. 3**). We further suggest

377    using the intersection of multiple variant callers to increase confidence in the data, especially

378    when estimating SNV frequency (**Fig. 5B**). Using all variant callers for analysis would likely be

379    tedious and unrealistic, thus we looked at the intersection of just two callers, iVar and timo, and

380    we found a similar trade-off in true positive and false positive data when using a single replicate

381    and a cutoff of 0.03 (**Fig. 5E**). Based on these data, it is clear that there are many considerations

382    necessary when performing minority variant analyses, and parameters and cutoffs should thus

383    be chosen carefully and thoughtfully depending on the data available. In general, using replicate

384    data and multiple callers provides the highest confidence set of SNVs and the most accurate

385    frequencies.

DISCUSSION

It has long been understood that intrahost viral populations are heterogeneous in nature, however capturing and measuring this viral diversity is complicated due to errors introduced during preparation and sequencing. We set out to identify the optimal tools, parameters, and filtering methods necessary for accurate variant identification. To accomplish this goal, we used a combination of simulated and synthetic influenza and SARS-CoV-2 sequence data to test the technical and experimental challenges and limitations of minority variant analyses. We found that sequencing depth and choice of variant caller has a significant impact on sensitivity of minor variant calls, and most false positive SNVs were detected at either low allele frequency, or low read depth. Additionally, our results show that replicate sequencing allows for the use of lower frequency thresholds, and this combination provides the best results, keeping the false discovery rate low, without sacrificing true positive data. We tested our optimized frequency and coverage cutoffs using SARS-CoV-2 sequence data from clinical infections. Most variant callers did not agree on the set of minor variants in the virus sequence data from clinical samples, and most minority variants were not reproducible across replicates. Ultimately, we determined that using a combination of sequencing replicates and multiple variant callers, along with moderate allele frequency and coverage cutoffs, can increase confidence in SNV calls. Our results outline the main considerations for minority variant analyses and highlight the intricacies and difficulties of studies of this nature. Further, we provide a framework for designing minority variant analyses, which will ultimately lead to more accurate conclusions surrounding viral transmission and evolution.

Using a standardized set of parameters, most callers performed relatively similarly on high coverage simulated data, having both high precision and high recall. The main differences in caller performance were seen in lower coverage data or at low frequencies. As many minority variants are found at low frequencies, understanding how tools perform under these conditions is more relevant to analyses of real sequencing data. Timo had the lowest recall at lower coverages and

412   simulated frequencies due to its rigid requirements for SNVs to be above the 0.01 threshold

413   parameter, while many other callers found SNVs at or below this frequency, regardless of setting

414   a 0.01 AF cutoff. Timo, iVar, and Varscan all have the functionality to drop the input frequency

415   parameter down to 0.001. Decreasing this parameter did not change the accuracy of iVar and

416   Varscan but did increase the recall of timo. These data highlight the importance of optimizing

417   bioinformatic tools to one's own data.

418        As previously observed by our group and others, the best method for filtering out errors

419   generated during sample processing is to sequence each sample twice and only keep the SNVs

420   found in both replicates. Sequencing replicates removed nearly all false positive calls in simulated

421   data and significantly reduced the number of false positive SNVs in the synthetic datasets.

422   However, for the synthetic datasets, the number of false positive SNVs was highly dependent on

423   the variant caller used. HaplotypeCaller, Lofreq, and Mutect2 were all made and optimized for

424   identifying variants in cancer cell datasets and had significantly higher false discovery rates than

425   tools designed for viral use, particularly at low allele frequencies. Adjusting the filtering or input

426   parameters on these callers may better optimize them for their use on viral data. For example,

427   HaplotypeCaller suggests additional filtering of output data, however, when performed on this

428   dataset, SNV detection was significantly reduced. Without this additional filtering, most variants

429   are identified (high recall) but high numbers of false positives are included, suggesting additional

430   optimization could improve performance.

431        The use of simulated and synthetic data allows tools to be tested against a true set of

432   minority variants. This is invaluable given that it is not possible to know which SNVs are real, and

433   which are not, in clinical sequencing data. However, each of these controlled settings comes with

434   both benefits and limitations. Simulated data allows for testing of various coverages using random

435   downsampling, showing how low coverage data impacts caller performance and allele frequency

436   estimation. Since high coverage is often the goal in sequencing, and coverage is unpredictable

437   with real samples (42), this is a variable that can only be truly tested through simulation.

438    Conversely, while error models in simulated data can imitate PCR and sequencing errors, there

439    are limitations to how 'real' these errors look to variant callers. Synthetic data allows for real error

440    to be introduced during sample preparation and ultimately be easily separated from real variation

441    during analysis. This gives a more accurate depiction of false positives that would exist in clinical

442    data and variant caller performance on synthetic data may be more representative of true

443    performance. However, both methods are limited in their ability to mimic sequence data from

444    clinical samples where the RNA is not as homogeneous and intact as with synthetic RNA, and

445    likely contains human contamination from sample collection and potential degradation from

446    freezing, thawing, and handling (43). This is evident in the disagreement between variant callers

447    and replicates, both on the set of minority variants and on their frequencies, in real data.

448    Combined, the simulated, synthetic, and clinical data sets show that there will always be a trade-

449    off between inclusion of the maximum number of true variants, and inclusion of false positive data.

450        Our study provides an extensive framework for studying minority variants in sequence

451    data from clinical samples, outlining major considerations around choice of variant caller,

452    application of frequency and coverage thresholds, and use of replicate sequencing. Further, we

453    have established a pipeline that can be used for further testing and optimization of parameters,

454    or for other viruses. This work will inform and improve future studies of intrahost variation and

455    estimates surrounding viral diversity and viral evolution.

456     FIGURE LEGENDS

457     **Figure 1. Variant caller performance on simulated and synthetic data with variants at set**

458     **allele frequencies. (A)** F1 statistic for each variant caller across a range of downsampling

459     fractions and allele frequency values. Values shown are mean and standard deviation of the four

460     viruses (A/H1N1, A/H3N2, B/Victoria, SARS-CoV-2) using standard input parameters (**Table S1**).

461     Color represents the variant caller used. **(B)** Precision/Recall graphs of each variant caller across

462     allele frequencies (point shape) for downsampling fractions 0.003 (top left) and 0.002 (top right)

463     or across downsampling fractions (point shape) for allele frequencies (AF) 0.02 (bottom left) and

464     0.03 (bottom right). Color represents variant caller. Mean and standard deviation are shown

465     across the four viruses for precision and recall scores. **(C)** Scatter plot showing median observed

466     frequency (x-axis) versus individual observed frequency (y-axis) for each synthetic influenza gene

467     segment and SARS-CoV-2 genome. Black dots indicate the median value for each expected

468     frequency. Color represents the variant caller used.

469

470     **Figure 2. Frequency and coverage of false positive variants in synthetic and simulated**

471     **data. (A-B)** Scatter plots and associated histograms showing number of false positive SNVs

472     across allele frequencies and read depths in synthetic (**A**) or simulated (**B**) data. Dotted lines are

473     drawn at allele frequency = 0.03 and read depth = 200x. Color represents the variant caller used.

474

475     **Figure 3. Effect of frequency cutoffs and sequencing replicates on false discovery rate and**

476     **false negative rate in synthetic and simulated data. (A)** False discovery rate (FDR) of synthetic

477     data using either single replicate data (colored points and lines) with applied frequency cutoffs

478     (line type) or merged two replicate data without cutoffs (solid black points and lines). Lines

479     represent mean across viruses. **(B)** False negative rate (FNR) of synthetic data using either single

480     replicate data (colored points and lines) with applied frequency cutoffs (line type) or merged two

481     replicate data without cutoffs (solid black points and lines) as above.

482

483 **Figure 4. Effect of variant caller on identification and allele frequency estimation of SNVs**

484 **in SARS-CoV-2 data from clinical samples. (A)** Bar plot showing raw number of variants

485 identified by each variant caller in replicate 1 (left bar) or replicate 2 (right bar). **(B)** Upset plot

486 showing agreement of minority variants between variant callers in each replicate using an allele

487 frequency cutoff of 0.03 and coverage cutoff of 200X. Vertical bars indicate the size of the shared

488 set of variants while dots and connecting lines show which callers share a given set of identified

489 variants. **(C)** Scatter plot showing the output frequency of minority variants identified by two

490 different variant callers. Color represents replicate. Variants with frequency of 0 were not identified

491 by that variant caller.

492

493 **Figure 5. Reproducibility of minority variants across sequencing replicates. (A)** Bar plot

494 showing number of reproducible variants across sequencing replicates by each variant caller.

495 Percentages shown are the percentage of total individual variants that were reproducible.

496 Background bars indicate the total number of variants found by each tool in each replicate (r1-

497 left, r2-right). **(B)** UpsetR plot showing overlap of reproducible variants across Freebayes, iVar,

498 timo, and Varscan, using a frequency cutoff of 0.01 and coverage cutoff of 200x. Vertical bars

499 indicate the size of the shared set of variants while dots and connecting lines show which callers

500 share a given set of reproducible variants. **(C)** Scatter plot showing frequency of variants across

501 sequencing replicates with frequency in replicate 1 on the x-axis and frequency in replicate 2 on

502 the y-axis. Color represents reproducibility of each variant across variant callers and replicates.

503 **(D,E)** Line graph showing the number of 'true positive' (D) and 'false positive' (E) in single replicate

504 data across allele frequency cutoffs compared to merged replicate data. A TP variant is defined

505 as a SNV found by the selected callers in both replicates (80 variants shown in **B**) and a FP is

506 defined as any other variant found in an individual replicate by the selected callers. Color

507 represents sequencing replicate.

508

## SUPPLEMENTAL TABLES

510 **Supplemental Table 1.** List of variant callers and parameters used for each.

511 **Supplemental Table 2.** Metadata associated with diagnostic samples processed for whole

512 genome sequencing of SARS-CoV-2.

513

## SUPPLEMENTAL FIGURES

515 **Supplemental Figure 1. Experimental setup of simulated and synthetic data generation. (A)**

516 Schematic of the SNV simulation pipeline. **(B)** Nucleotide positions of SNVs in simulated data

517 across the genomes of A/H1N1 (n=121), A/H3N2 (n=110), B/Victoria (n=118), and SARS-CoV-2

518 (n=144). Gene segments for influenza A and B strains are ordered largest (PB2) to smallest (NS),

519 left to right. **(C)** Schematic of allele frequency and viral loads used for mixing WT and variant RNA

520 for library preparation and sequencing. **(D)** Location of synthetic SNVs across the three influenza

521 gene segments. Gray lines represent the designed SNVs (PB2 (n=18), HA (n=14), NA (n=14))

522 and red lines represent the pre-mRT-PCR errors, likely generated during template preparation or

523 *in vitro*-transcription.

524

525 **Supplemental Figure 2. Variant caller performance using default, standard, and custom**

526 **parameters. (A)** F1 statistic for each variant caller across a range of downsampling fractions and

527 allele frequencies. Values shown are mean and standard deviation of four viruses (A/H1N1,

528 A/H3N2, B/Victoria, SARS-CoV-2) using default (top), standard (middle) or custom (bottom) input

529 parameters (**Table S1**). Color represents the variant caller used. Dark timo points and lines

530 indicate timo output with the removal of the binomial check option.

531

532 **Supplemental Figure 3. Coverage and Expected vs. Observed frequency of variants in**

533 **synthetic data samples (A)** Coverage plots showing log10 read depth for each synthetic

23

534     influenza control across the PB2, HA, and NA gene segments. Color represents individual

535     samples. **(B)** Expected (x-axis) versus observed frequency (y-axis) for synthetic influenza SNVs

536     (replicate one only). Color represents the copy number.  **(C)** Expected (x-axis) versus observed

537     frequency (y-axis) for synthetic SARS-CoV-2 SNVs. Dashed black lines represent y=x.

538

539     **Supplemental Figure 4. Effect of frequency cutoffs and sequencing replicates on false**

540     **discovery rate and false negative rate in simulated data. (A)** False discovery rate (FDR) of

541     simulated data using either single replicate data (colored points and lines) with applied frequency

542     cutoffs (line type) or merged two replicate data without cutoffs (solid black points and lines).

543     Background points represent values for individual viruses. Lines represent mean across viruses.

544     **(B)** False negative rate (FNR) of simulated data using either single replicate data (colored points

545     and lines) with applied frequency cutoffs (line type) or merged two replicate data without cutoffs

546     (solid black points and lines) as above.

547

548     **Supplemental Figure 5. Effect of coverage cutoffs and sequencing replicates on false**

549     **discovery rate and false negative rate in synthetic and simulated data. (A, B)** False discovery

550     rate (FDR) (**A**) of false negative rate (FNR) (**B**) of synthetic data using either single replicate data

551     (colored points and lines) with applied coverage cutoffs (line type) or merged two replicate data

552     without cutoffs (solid black points and lines). Background points show individual values for each

553     of the four viruses. Lines represent mean across viruses. (**C, D**) False discovery rate (FDR) (**C**)

554     of false negative rate (FNR) (**D**) of synthetic data using either single replicate data (colored points

555     and lines) with applied coverage cutoffs (line type) or merged two replicate data without cutoffs

556     (solid black points and lines) as above.

557
558     **Supplemental Figure 6. Effect of sequencing replicates on the accuracy of allele frequency**

559     **estimation. (A)** Scatter plot showing the frequency of variants in synthetic influenza data across

24

560    sequencing replicates with frequency in replicate 1 on the x-axis and frequency in replicate 2 on

561    the y-axis. Color represents the SNV type. **(B, C)** Percent error ($\frac{(observed-expected)}{expected} * 100$) of

562    single replicate (greyed boxes) or merged replicate (open boxes) SNVs in simulated data across

563    downsampling fractions (**B**) or synthetic data across expected frequencies and viral loads (**C**).

564    Color represents the variant caller used.

565

566    **Supplemental Figure 7. Quantification of majority and minority variants identified in data from**

567    **SARS-CoV-2 clinical specimens (A)** Scatter plot showing Ct value against percentage of genome

568    with coverage over 200x after filtering for only samples with 80% of the genome over the 200x

569    cutoff. **(B)** Upset plot showing agreement of consensus changes between variant callers in each

570    replicate using an allele frequency cutoff of 0.5 and coverage cutoff of 5X. Vertical bars indicate

571    the size of the shared set of changes while dots and connecting lines show which callers share a

572    given set of identified changes. **(C)** Box and whisker plot showing number of minor variants with

573    indicated allele frequency cutoffs found in replicate 1 and replicate 2 sequencing data. Points

574    represent individual samples. Boxes and whiskers show min, first quartile, median, third quartile

575    and max for each replicate. **(D)** Upset plot showing agreement of minority variants between all

576    variant callers in each replicate using an allele frequency cutoff of 0.03 and coverage cutoff of

577    200x. Vertical bars indicate the size of the shared set of variants while dots and connecting lines

578    show which callers share a given set of identified variants.

579

580    REFERENCES

581    1.    Arnold JJ, Cameron CE. 2004. Poliovirus RNA-dependent RNA polymerase (3Dpol): pre-

582          steady-state kinetic analysis of ribonucleotide incorporation in the presence of Mg2+.

583          Biochemistry 43:5126-37.

584    2.    Sanjuan R. 2012. From molecular genetics to phylodynamics: evolutionary relevance of

585          mutation rates across viruses. PLoS Pathog 8:e1002685.

586    3.    Duffy S, Shackelton LA, Holmes EC. 2008. Rates of evolutionary change in viruses:

587          patterns and determinants. Nat Rev Genet 9:267-76.

588    4.    Peck KM, Lauring AS. 2018. Complexities of Viral Mutation Rates. J Virol 92.

589    5.    Domingo E. 2002. Quasispecies Theory in Virology. J Virol 76:463-465.

590    6.    Sanjuan R, Domingo-Calap P. 2016. Mechanisms of viral mutation. Cell Mol Life Sci

591          73:4433-4448.

592    7.    Lynch M, Ackerman MS, Gout JF, Long H, Sung W, Thomas WK, Foster PL. 2016. Genetic

593          drift, selection and the evolution of the mutation rate. Nat Rev Genet 17:704-714.

594    8.    Gorbalenya AE, Enjuanes L, Ziebuhr J, Snijder EJ. 2006. Nidovirales: evolving the largest

595          RNA virus genome. Virus Res 117:17-37.

596    9.    Smith EC, Sexton NR, Denison MR. 2014. Thinking Outside the Triangle: Replication

597          Fidelity of the Largest RNA Viruses. Annu Rev Virol 1:111-32.

598    10.   Sanjuan R, Moya A, Elena SF. 2004. The distribution of fitness effects caused by single-

599          nucleotide substitutions in an RNA virus. Proc Natl Acad Sci U S A 101:8396-401.

600    11.   Visher E, Whitefield SE, McCrone JT, Fitzsimmons W, Lauring AS. 2016. The Mutational

601          Robustness of Influenza A Virus. PLoS Pathog 12:e1005856.

602   12.   Peck KM, Chan CH, Tanaka MM. 2015. Connecting within-host dynamics to the rate of

603         viral molecular evolution. Virus Evol 1:vev013.

604   13.   Wang Y, Wang D, Zhang L, Sun W, Zhang Z, Chen W, Zhu A, Huang Y, Xiao F, Yao J, Gan

605         M, Li F, Luo L, Huang X, Zhang Y, Wong SS, Cheng X, Ji J, Ou Z, Xiao M, Li M, Li J, Ren P,

606         Deng Z, Zhong H, Xu X, Song T, Mok CKP, Peiris M, Zhong N, Zhao J, Li Y, Li J, Zhao J.

607         2021. Intra-host variation and evolutionary dynamics of SARS-CoV-2 populations in

608         COVID-19 patients. Genome Med 13:30.

609   14.   Martin MAK, Katia. 2021. Reanalysis of deep-sequencing data from Austria points

610         towards a small SARS-COV-2 transmission bottleneck on the order of one to three

611         virions. bioRxiv doi:https://doi.org/10.1101/2021.02.22.432096.

612   15.   McCrone JT, Lauring AS. 2018. Genetic bottlenecks in intraspecies virus transmission.

613         Curr Opin Virol 28:20-25.

614   16.   Wang D, Wang Y, Sun W, Zhang L, Ji J, Zhang Z, Cheng X, Li Y, Xiao F, Zhu A, Zhong B,

615         Ruan S, Li J, Ren P, Ou Z, Xiao M, Li M, Deng Z, Zhong H, Li F, Wang WJ, Zhang Y, Chen W,

616         Zhu S, Xu X, Jin X, Zhao J, Zhong N, Zhang W, Zhao J, Li J, Xu Y. 2021. Population

617         Bottlenecks and Intra-host Evolution During Human-to-Human Transmission of SARS-

618         CoV-2. Front Med (Lausanne) 8:585358.

619   17.   Lythgoe KA, Hall M, Ferretti L, de Cesare M, MacIntyre-Cockett G, Trebes A, Andersson

620         M, Otecko N, Wise EL, Moore N, Lynch J, Kidd S, Cortes N, Mori M, Williams R, Vernet G,

621         Justice A, Green A, Nicholls SM, Ansari MA, Abeler-Dorner L, Moore CE, Peto TEA, Eyre

622         DW, Shaw R, Simmonds P, Buck D, Todd JA, Oxford Virus Sequencing Analysis G, Connor

623         TR, Ashraf S, da Silva Filipe A, Shepherd J, Thomson EC, Consortium C-GU, Bonsall D,

624        Fraser C, Golubchik T. 2021. SARS-CoV-2 within-host diversity and transmission. Science

625        doi:10.1126/science.abg0821.

626   18.   Pybus OG, Rambaut A. 2009. Evolutionary analysis of the dynamics of viral infectious

627        disease. Nat Rev Genet 10:540-50.

628   19.   Rockett R, Basile K, Maddocks S, Fong W, Agius JE, Johnson-Mackinnon J, Arnott A,

629        Chandra S, Gall M, Draper J, Martinez E, Sim EM, Lee C, Ngo C, Ramsperger M, Ginn AN,

630        Wang Q, Fennell M, Ko D, Lim HL, Gilroy N, O'Sullivan MVN, Chen SC, Kok J, Dwyer DE,

631        Sintchenko V. 2022. Resistance Mutations in SARS-CoV-2 Delta Variant after Sotrovimab

632        Use. N Engl J Med 386:1477-1479.

633   20.   Martinez-Gonzalez B, Soria ME, Vazquez-Sirvent L, Ferrer-Orta C, Lobo-Vega R, Minguez

634        P, de la Fuente L, Llorens C, Soriano B, Ramos-Ruiz R, Corton M, Lopez-Rodriguez R,

635        Garcia-Crespo C, Somovilla P, Duran-Pastor A, Gallego I, de Avila AI, Delgado S, Moran F,

636        Lopez-Galindez C, Gomez J, Enjuanes L, Salar-Vidal L, Esteban-Munoz M, Esteban J,

637        Fernandez-Roblas R, Gadea I, Ayuso C, Ruiz-Hornillos J, Verdaguer N, Domingo E, Perales

638        C. 2022. SARS-CoV-2 Mutant Spectra at Different Depth Levels Reveal an Overwhelming

639        Abundance of Low Frequency Mutations. Pathogens 11.

640   21.   Dinis JM, Florek KR, Fatola OO, Moncla LH, Mutschler JP, Charlier OK, Meece JK,

641        Belongia EA, Friedrich TC. 2016. Deep Sequencing Reveals Potential Antigenic Variants

642        at Low Frequencies in Influenza A Virus-Infected Humans. J Virol 90:3355-65.

643   22.   Xue KS, Bloom JD. 2020. Linking influenza virus evolution within and between human

644        hosts. Virus Evol 6:veaa010.

645     23.     Valesano AL, Taniuchi M, Fitzsimmons WJ, Islam MO, Ahmed T, Zaman K, Haque R,

646          Wong W, Famulare M, Lauring AS. 2021. The Early Evolution of Oral Poliovirus Vaccine Is

647          Shaped by Strong Positive Selection and Tight Transmission Bottlenecks. Cell Host

648          Microbe 29:32-43 e4.

649     24.     Lauring AS. 2020. Within-Host Viral Diversity: A Window into Viral Evolution. Annu Rev

650          Virol 7:63-81.

651     25.     McCrone JT, Lauring AS. 2016. Measurements of Intrahost Viral Diversity Are Extremely

652          Sensitive to Systematic Errors in Variant Calling. J Virol 90:6884-95.

653     26.     Valesano ALR, Kalee E; Dimcheff, Derek E; Blair, Christopher N; Fitzsimmons, William J;

654          Petrie, Joshua G; Martin, Emily T; Lauring, Adam S. 2021. Temporal dynamics of SARS-

655          CoV-2 mutation accumulation within and across infected hosts. bioRxiv

656          doi:https://doi.org/10.1101/2021.01.19.427330.

657     27.     Grubaugh ND, Gangavarapu K, Quick J, Matteson NL, De Jesus JG, Main BJ, Tan AL, Paul

658          LM, Brackney DE, Grewal S, Gurfield N, Van Rompay KKA, Isern S, Michael SF, Coffey LL,

659          Loman NJ, Andersen KG. 2019. An amplicon-based sequencing framework for accurately

660          measuring intrahost virus diversity using PrimalSeq and iVar. Genome Biol 20:8.

661     28.     Koboldt DC. 2020. Best practices for variant calling in clinical sequencing. Genome Med

662          12:91.

663     29.     Olson ND, Lund SP, Colman RE, Foster JT, Sahl JW, Schupp JM, Keim P, Morrow JB, Salit

664          ML, Zook JM. 2015. Best practices for evaluating single nucleotide variant calling

665          methods for microbial genomics. Front Genet 6:235.

666   30.   Stead LF, Sutton KM, Taylor GR, Quirke P, Rabbitts P. 2013. Accurately identifying low-

667         allelic fraction variants in single samples with next-generation sequencing: applications

668         in tumor subclone resolution. Hum Mutat 34:1432-8.

669   31.   Auwera Gvd, O'Connor BD. 2020. Genomics in the cloud : using Docker, GATK, and WDL

670         in Terra, First edition. ed. O'Reilly Media, Sebastopol, CA.

671   32.   Cibulskis K, Lawrence MS, Carter SL, Sivachenko A, Jaffe D, Sougnez C, Gabriel S,

672         Meyerson M, Lander ES, Getz G. 2013. Sensitive detection of somatic point mutations in

673         impure and heterogeneous cancer samples. Nat Biotechnol 31:213-9.

674   33.   Garrison EM, Gabor. 2012. Haplotype-based variant detection from short-read

675         sequencing. ArXiv 1207.3907v2.

676   34.   Benjamin D, Sato T, Cibulskis K, Getz G, Stewart C, L. L. 2019. Calling Somatic SNVs and

677         Indels with Mutect2. bioRxiv doi:https://doi.org/10.1101/861054.

678   35.   Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, Miller CA, Mardis ER, Ding

679         L, Wilson RK. 2012. VarScan 2: somatic mutation and copy number alteration discovery

680         in cancer by exome sequencing. Genome Res 22:568-76.

681   36.   Wilm A, Aw PP, Bertrand D, Yeo GH, Ong SH, Wong CH, Khor CC, Petric R, Hibberd ML,

682         Nagarajan N. 2012. LoFreq: a sequence-quality aware, ultra-sensitive variant caller for

683         uncovering cell-population heterogeneity from high-throughput sequencing datasets.

684         Nucleic Acids Res 40:11189-201.

685   37.   Wybo WA, Jordan J, Ellenberger B, Marti Mengual U, Nevian T, Senn W. 2021. Data-

686         driven reduction of dendritic morphologies with preserved dendro-somatic responses.

687         Elife 10.

688    38.    Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina

689           sequence data. Bioinformatics 30:2114-20.

690    39.    Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler

691           transform. Bioinformatics 25:1754-60.

692    40.    Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A,

693           Jordan T, Shakir K, Roazen D, Thibault J, Banks E, Garimella KV, Altshuler D, Gabriel S,

694           DePristo MA. 2013. From FastQ data to high confidence variant calls: the Genome

695           Analysis Toolkit best practices pipeline. Curr Protoc Bioinformatics 43:11 10 1-11 10 33.

696    41.    Zhao L, Illingworth CJR. 2019. Measurements of intrahost viral diversity require an

697           unbiased diversity metric. Virus Evol 5:vey041.

698    42.    Lambisia AW, Mohammed KS, Makori TO, Ndwiga L, Mburu MW, Morobe JM, Moraa

699           EO, Musyoki J, Murunga N, Mwangi JN, Nokes DJ, Agoti CN, Ochola-Oyier LI, Githinji G.

700           2022. Optimization of the SARS-CoV-2 ARTIC Network V4 Primers and Whole Genome

701           Sequencing Protocol. Front Med (Lausanne) 9:836728.

702    43.    Li L, Li X, Guo Z, Wang Z, Zhang K, Li C, Wang C, Zhang S. 2020. Influence of Storage

703           Conditions on SARS-CoV-2 Nucleic Acid Detection in Throat Swabs. J Infect Dis 222:203-
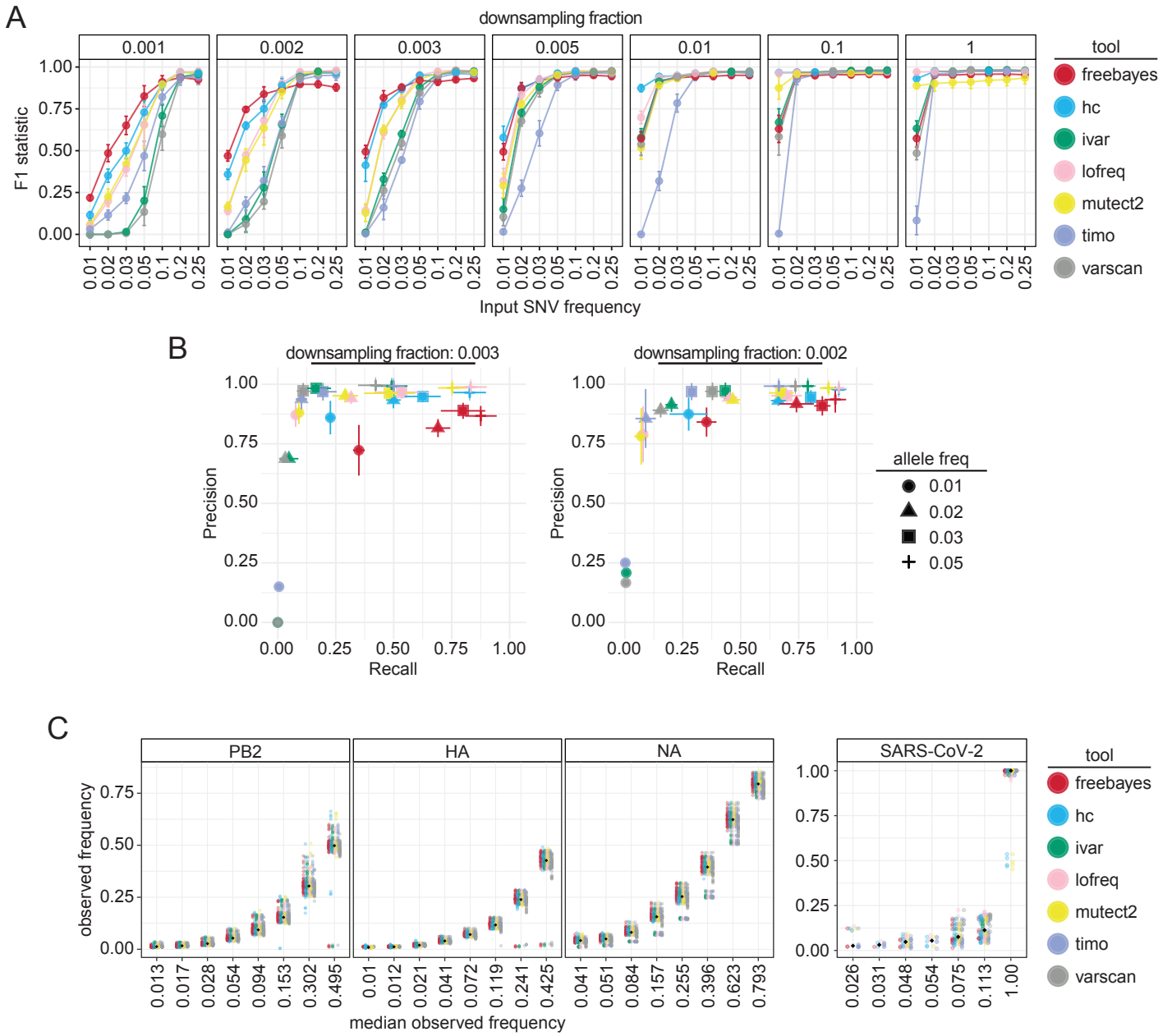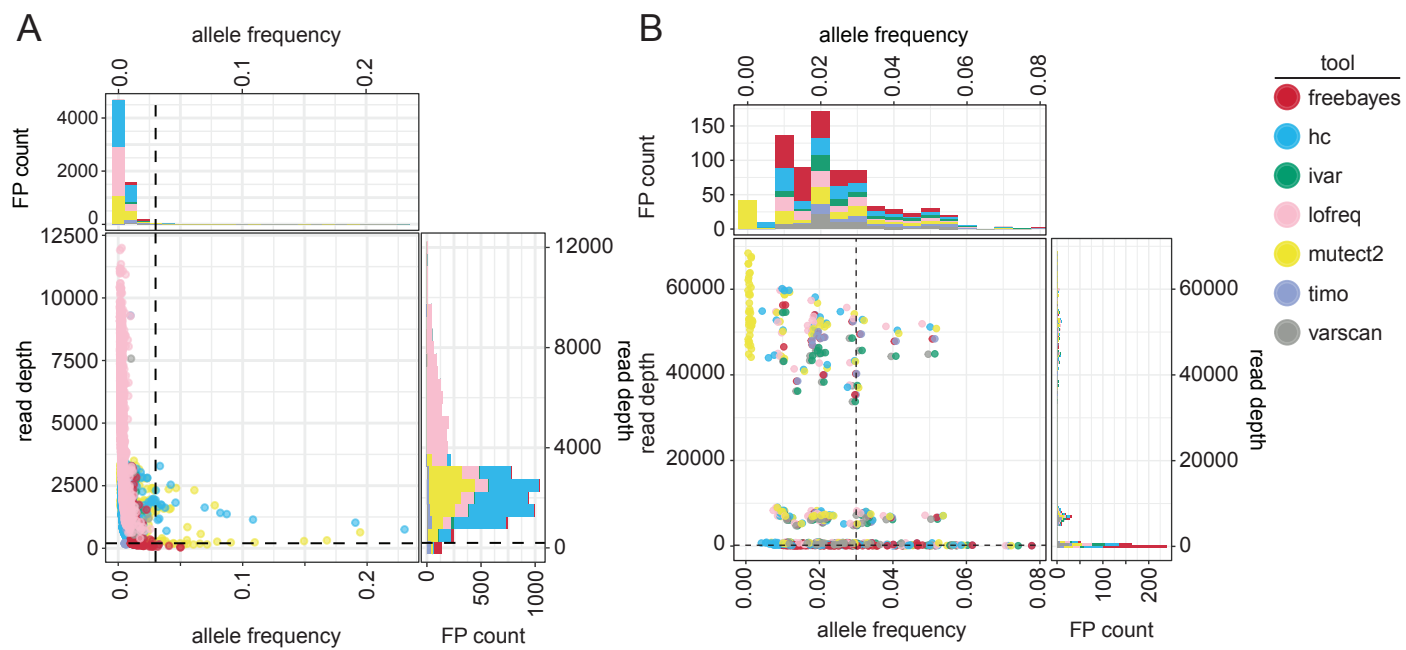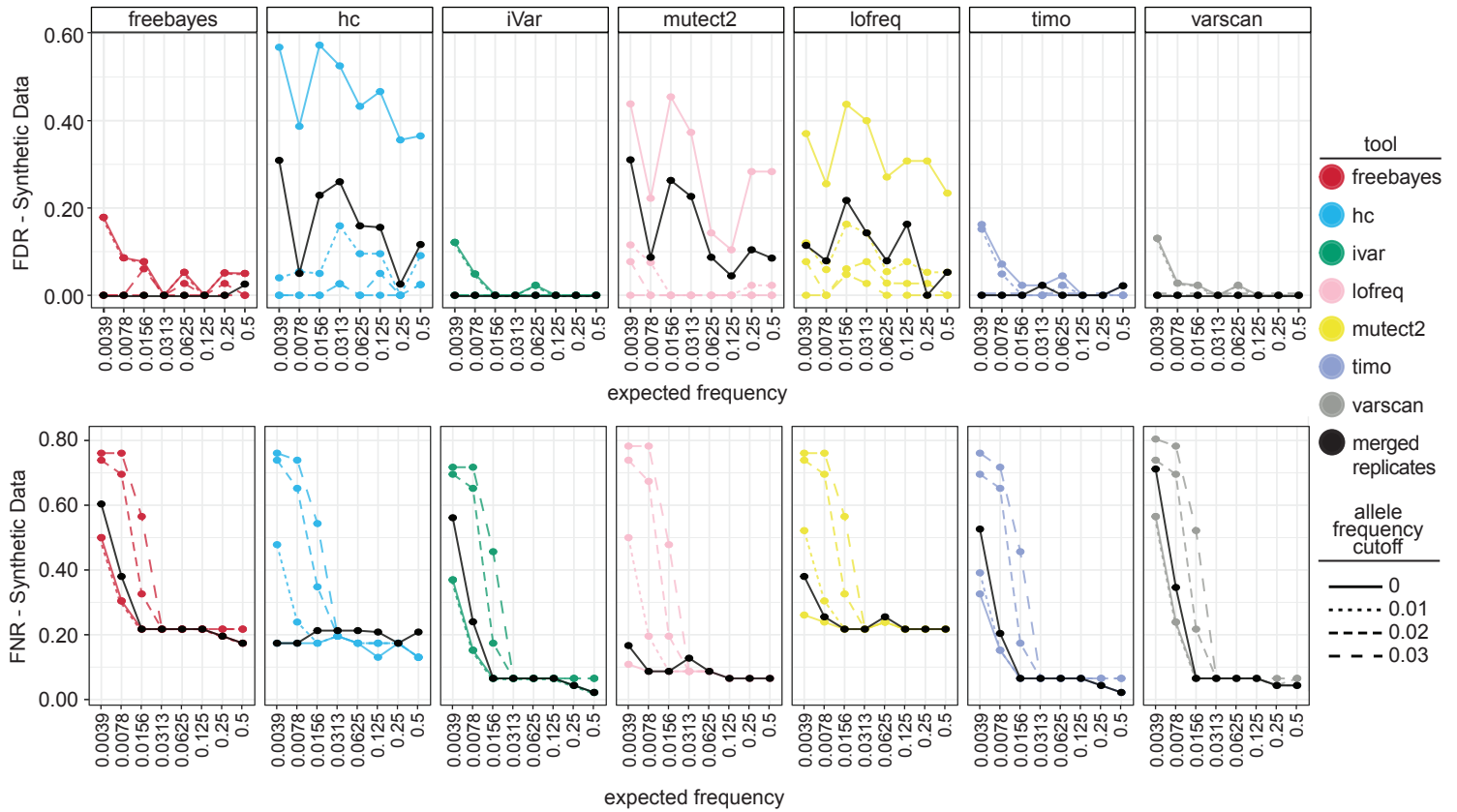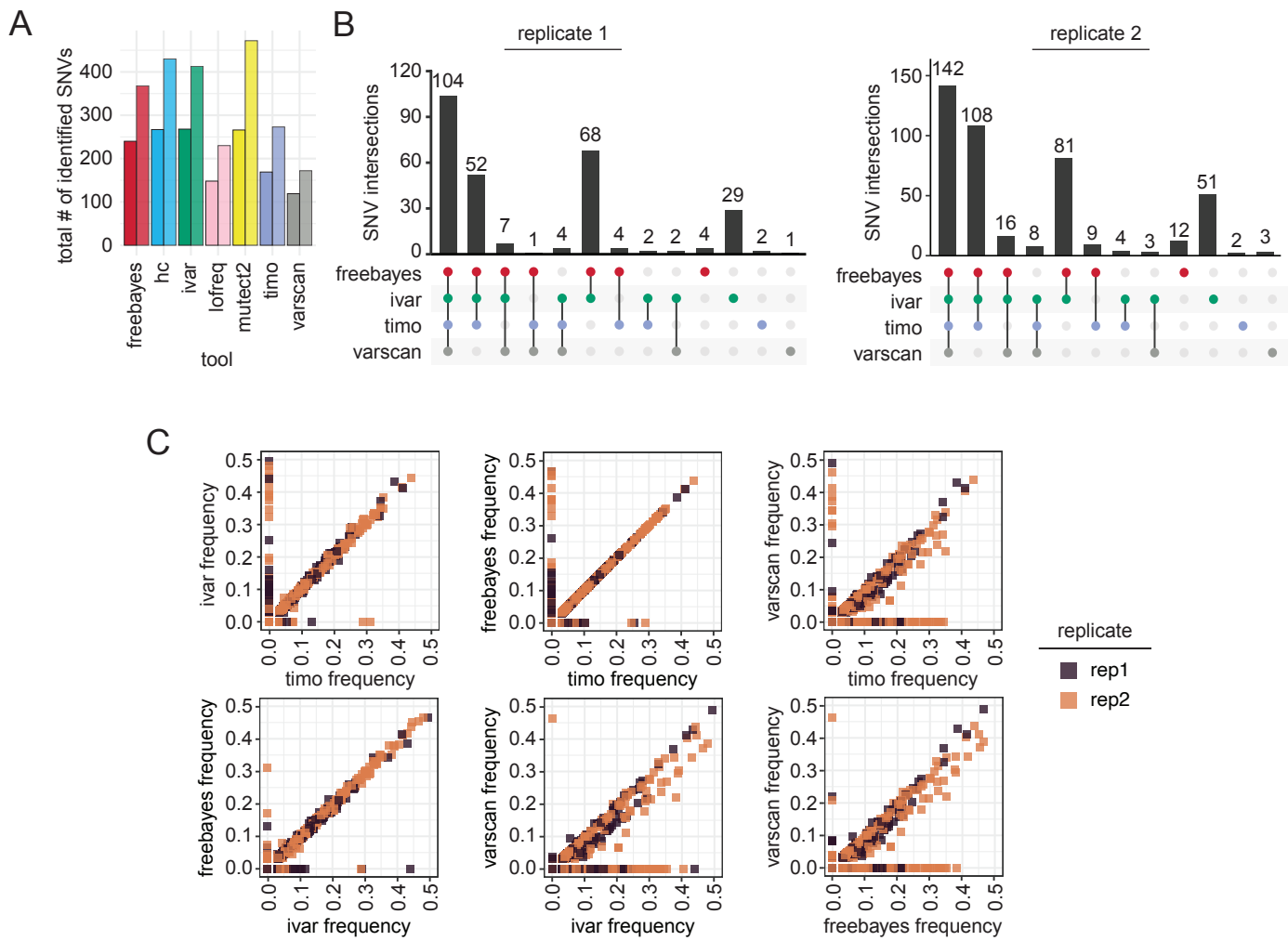
704           205.

705

Figure 1

Figure 2

Figure 3

Figure 4

Figure 5