

# 1 Estimating the relative proportions 2 of SARS-CoV-2 strains from 3 wastewater samples

4 Lenore Pipes<sup>1\*†</sup>, Zihao Chen<sup>2†</sup>, Svetlana Afanaseva<sup>3</sup>, Rasmus Nielsen<sup>1,3,4\*</sup>

**\*For correspondence:**

[rasmus\\_nielsen@berkeley.edu](mailto:rasmus_nielsen@berkeley.edu) (RN);  
[lpipes@berkeley.edu](mailto:lpipes@berkeley.edu) (LP)

<sup>†</sup>These authors contributed  
equally to this work

5 <sup>1</sup>Department of Integrative Biology, University of California-Berkeley; <sup>2</sup>School of  
6 Mathematical Sciences, Peking University; <sup>3</sup>Department of Statistics, University of  
7 California-Berkeley; <sup>4</sup>GLOBE Institute, University of Copenhagen

---

9 **Abstract** Wastewater surveillance has become essential for monitoring the spread of  
10 SARS-CoV-2. The quantification of SARS-CoV-2 RNA in wastewater correlates with the Covid-19  
11 caseload in a community. However, estimating the proportions of different SARS-CoV-2 strains  
12 has remained technically difficult. We present a method for estimating the relative proportions of  
13 SARS-CoV-2 strains from wastewater samples. The method uses an initial step to remove unlikely  
14 strains, imputation of missing nucleotides using the global SARS-CoV-2 phylogeny, and an  
15 Expectation-Maximization (EM) algorithm for obtaining maximum likelihood estimates of the  
16 proportions of different strains in a sample. Using simulations with a reference database of >3  
17 million SARS-CoV-2 genomes, we show that the estimated proportions accurately reflect the true  
18 proportions given sufficiently high sequencing depth and that the phylogenetic imputation is  
19 highly accurate and substantially improves the reference database.

---

## 21 Introduction

22 The ongoing pandemic of coronavirus disease of 2019 (Covid-19) caused by severe acute respira-  
23 tory syndrome coronavirus 2 (SARS-CoV-2) continues to be the world's worst public health emer-  
24 gency in the last century. There is an emerging need to identify the initiation of outbreaks, dis-  
25 tribution, and changing trends of Covid-19 in near real-time (*Korber et al., 2020; Rockett et al.,*  
26 *2020*). Wastewater-based epidemiology (WBE) has become an effective monitoring strategy for  
27 early detection of SARS-CoV-2 in communities as well as being an important method for informing  
28 public health interventions aimed at containing and mitigating Covid-19 outbreaks (*Ahmed et al.,*  
29 *2020*). WBE for SARS-CoV-2 can detect the virus excreted by both symptomatic and asymptomatic  
30 individuals alike thus making it an effective approach for modeling the disease signature of entire  
31 communities. WBE data also strongly correlates with the Covid-19 case rates in the community  
32 (*Medema et al., 2020a; Farkas et al., 2020*).

33 Currently, most analyses of WBE data for SARS-CoV-2 focus on identifying presence/absence as  
34 well as quantifying the abundance of the virus (*Kumar et al., 2020; Crits-Christoph et al., 2021; Wu*  
35 *et al., 2020; Medema et al., 2020b*). However, identifying and profiling multiple SARS-CoV-2 geno-  
36 types in a single sample can provide additional information for understanding the dynamics and  
37 transmission of certain strains. The alarming continued emergence of novel variants such as the  
38 Delta variant, B.1.617.2, and the Omicron variant, B.1.1.529, underscores the urgency and need  
39 for quantification of the abundance of different viral strains across communities. Unfortunately,  
40 it is difficult to precisely quantify the proportions of different strains of a virus in an environmen-

41 tal sample, such as wastewater, using standard sequencing technologies given the low quality and  
42 highly uneven depth of sequencing data. Adding to these challenges is that many strains are nearly  
43 identical differing by only one or a few mutations across approximately  $\sim 30,000$  nucleotides. With  
44 millions of possible candidate strain the combinatorial challenge of identifying the correct strains  
45 is large, particularly when strains are not identified by individual diagnostic mutations, but rather  
46 by sets of mutations that jointly helps distinguish the strains from each other. Nonetheless, quan-  
47 tification of strain composition in WBE data has the potential to become a cost-effective method  
48 to identify changes in viral community composition as SARS-CoV-2 becomes an endemic virus.

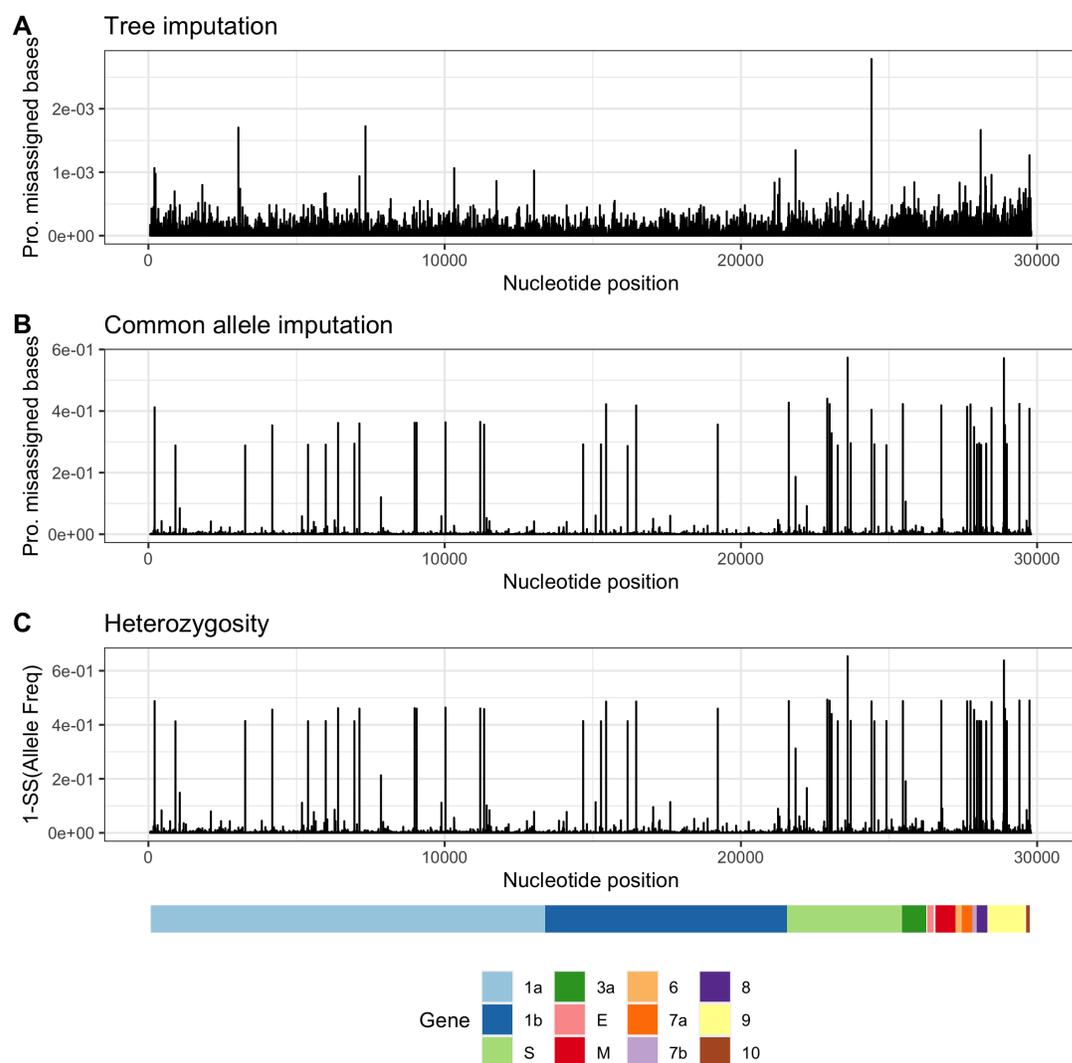
49 We present a method for estimating the proportion of different SARS-CoV-2 strains from shot-  
50 gun sequencing of wastewater samples allowing researchers to obtain results in real-time. The  
51 method is based on an initial filtering step, phylogenetic imputation of missing nucleotides, and  
52 an Expectation-Maximization (EM) algorithm for obtaining maximum likelihood estimates of the  
53 proportions of different strains in the sample. Using simulations, we show that the estimated pro-  
54 portions are close to the true proportions and that the phylogenetic imputation is highly accurate  
55 and improves the reference strains. We also apply this method to wastewater samples collected  
56 across the San Francisco Bay Area.

## 57 Results

### 58 Imputation

59 Many SARS-CoV-2 sequences submitted to public databases contain missing data (i.e., bases that  
60 are not coded as A, G, C, or T). This poses a problem when estimating the fraction of different  
61 SARS-CoV-2 strains, as strains with a high proportion of missing data in average will contain fewer  
62 nucleotide differences when compared to sequencing reads. We solve this problem using an impu-  
63 tation approach thereby allowing for a like-to-like comparison of reads against all reference strains.  
64 This method is in spirit similar to imputation approaches used in human genetics (e.g. *Marchini and*  
65 *Howie, 2010*), although as we will show, due to the strong phylogenetic structure in the SARS-CoV-2  
66 data, imputation is much more accurate than usually observed in diploid organisms. The method is  
67 based on calculating the posterior probability of each nucleotide in the leaf node of a phylogenetic  
68 tree and imputing based on the maximum posterior probability (see *Methods and Materials*). We  
69 compare the method (*Tree imputation*) to a naive imputation approach based on simply replacing  
70 missing nucleotides with the most frequent nucleotide observed in the alignment in that position  
71 (*Common allele imputation*). We evaluate the methods by first removing sequenced nucleotides in  
72 a real data set of 3,117,131 SARS-CoV-2 sequences and then re-imputing them using either *Tree*  
73 *imputation* or *Common allele imputation*.

74 For the vast majority of sites, *Tree imputation* has an error rate of  $< 5 \times 10^{-4}$  although a few  
75 sites have imputation errors between  $10^{-3}$  and  $3 \times 10^{-3}$  (Figure 1). The imputation error can be  
76 substantially higher for the naive *Common allele imputation* method with many sites showing error  
77 rates  $> 0.02$  (Figure 1B). These are sites with high heterozygosity (Figure 1C) where substituting  
78 with the most common allele leads to high error rates. While the error rates for the *Common*  
79 *allele imputation* method naturally is predicted by the heterozygosity, the pattern is somewhat dif-  
80 ferent for the *Tree imputation* method. The sites with highest imputation error are not the sites  
81 with highest heterozygosity, suggesting a high degree of homoplasmy in these sites not directly pre-  
82 dictable by the heterozygosity. These may be sites that switch allelic state often, i.e. have high  
83 mutation rates, but where the minor allele never increases substantially in frequency due to selec-  
84 tion. An alternative explanation is sequencing errors. In fact, the site with the highest amount of  
85 apparent imputation error (position 24,410) is a site known to have a high proportion of sequenc-  
86 ing errors ([https://github.com/W-L/ProblematicSites\\_SARS-CoV2](https://github.com/W-L/ProblematicSites_SARS-CoV2)). It is located in a primer bind-  
87 ing site where sequences containing the non-reference allele, A, often erroneously are assigned  
88 back to the reference allele, G, as a result of failed primer trimming during consensus building  
89 ([https://github.com/W-L/ProblematicSites\\_SARS-CoV2](https://github.com/W-L/ProblematicSites_SARS-CoV2)). The A allele is one of the defining muta-

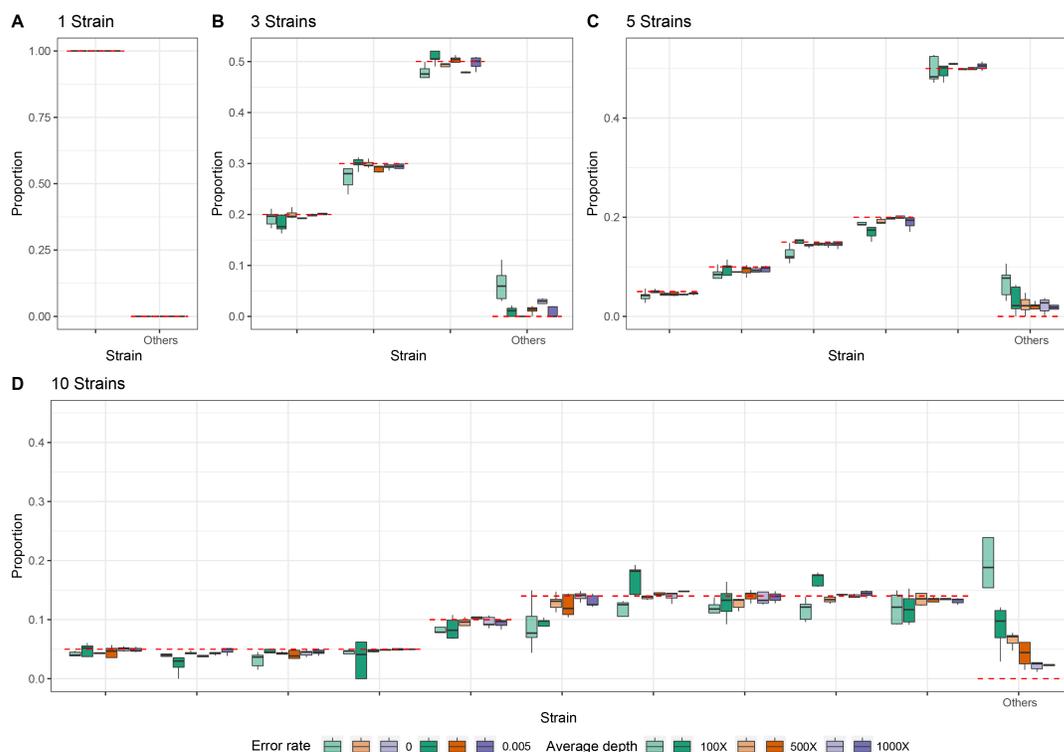


**Figure 1.** Proportion of misassigned bases along SARS-CoV-2 using the tree imputation method (A) and the common allele imputation method (B) against heterozygosity (C) using 3,117,131 SARS-CoV-2 genomes. Notice the difference in the scaling of the Y-axis of A and B.

90 tions of the delta strain and the apparent repeated re-emergence of the G allele within the delta  
 91 clade (Figure S1) is likely a consequence of this common sequencing error. Most other sites, in-  
 92 cluding the site with the highest heterozygosity, position 23,604 (Figure 1C), do not show a similar  
 93 pattern of homoplasy (Figure S2). This suggests that the sites with the highest apparent imputa-  
 94 tion error rate, might in fact have a much lower true imputation error; the *Tree imputation*  
 95 method may provide a more accurate assignment of alleles than the reported sequencing data for some  
 96 problematic sequencing sites.

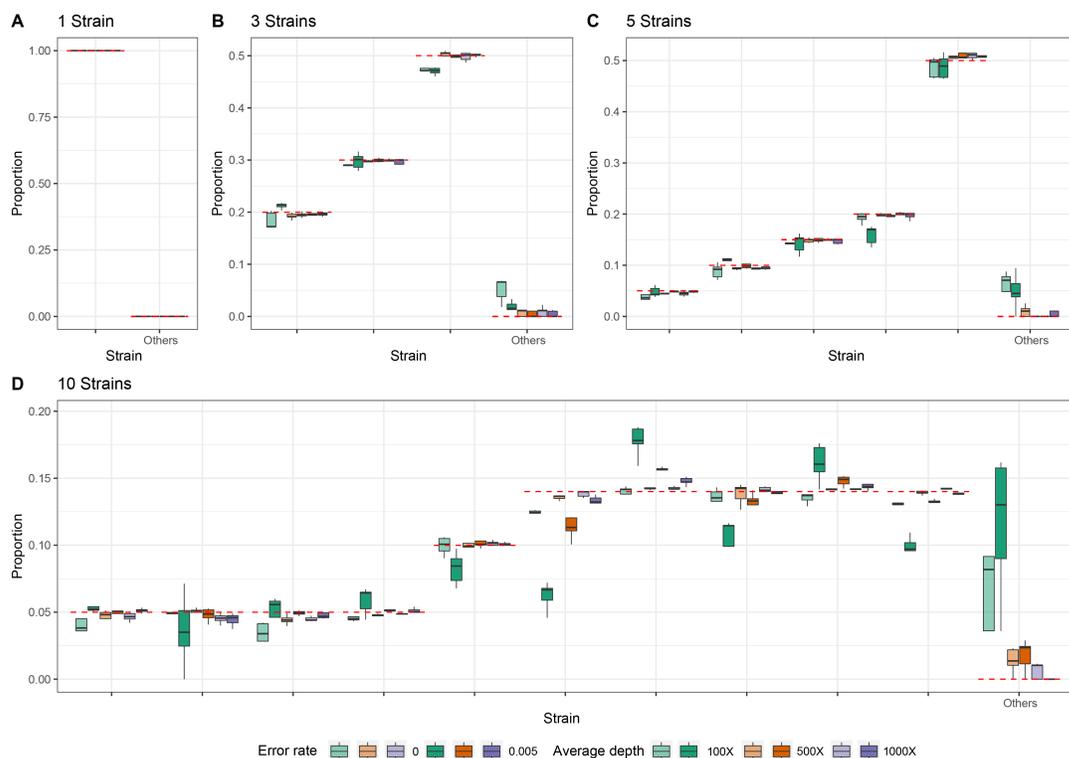
### 97 Simulations

98 In the Methods and Materials section, we describe an algorithm for estimating the proportion of  
 99 different SARS-CoV-2 strains in an environmental sample using maximum likelihood. To evaluate  
 100 the performance of the method, we simulate several sets of reads (single-end 300bp, paired-end  
 101 2x150bp, and paired-end 2x75bp) from 1, 3, 5, and 10 strains with an average depths of 100X, 500X,  
 102 1000X and a sequencing error rate of 0% and 0.5% (see Methods and Materials). We then apply  
 103 the method to these sets of reads using a database of 3,117,131 strains and report the estimated

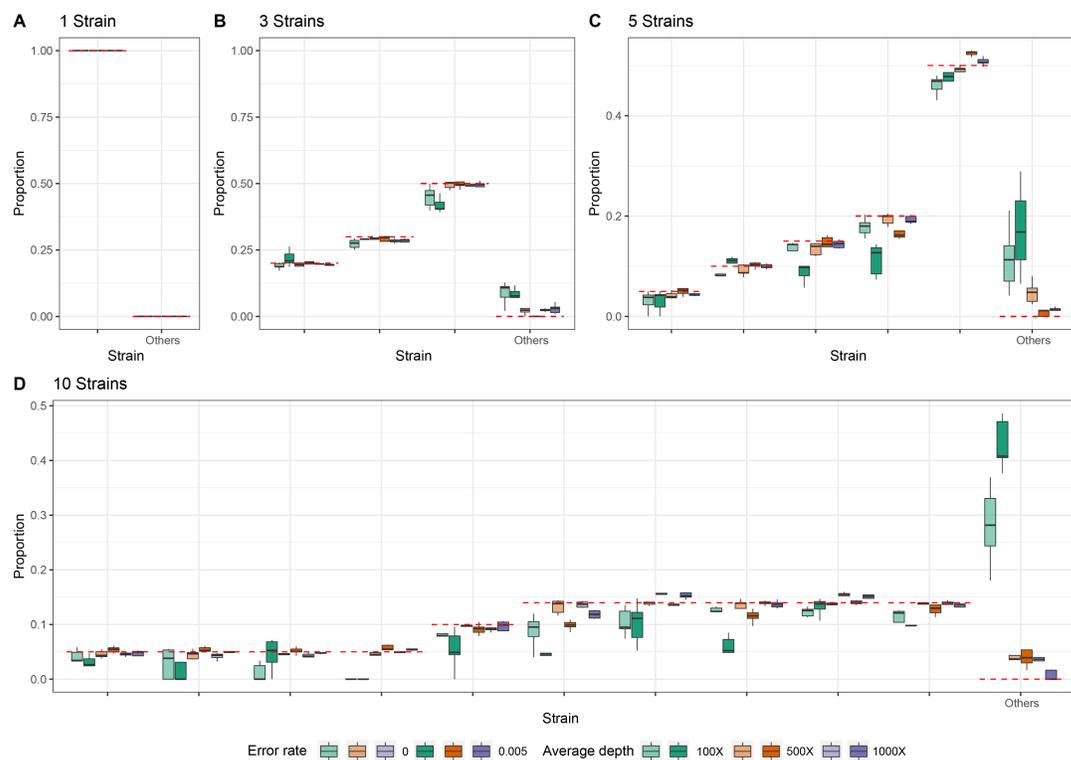


**Figure 2.** Estimated proportions for simulated 300 bp single-end reads with five replicates for when the sample truly contains 1 (A), 3 (B), 5 (C), or 10 (D) strains out of a total of 1,499,078 non-redundant candidate strains in the database. The red dashed lines indicate the true proportion of each strain. 'Other' indicates the sum of estimated proportions for all strains that are not truly represented in the sample.

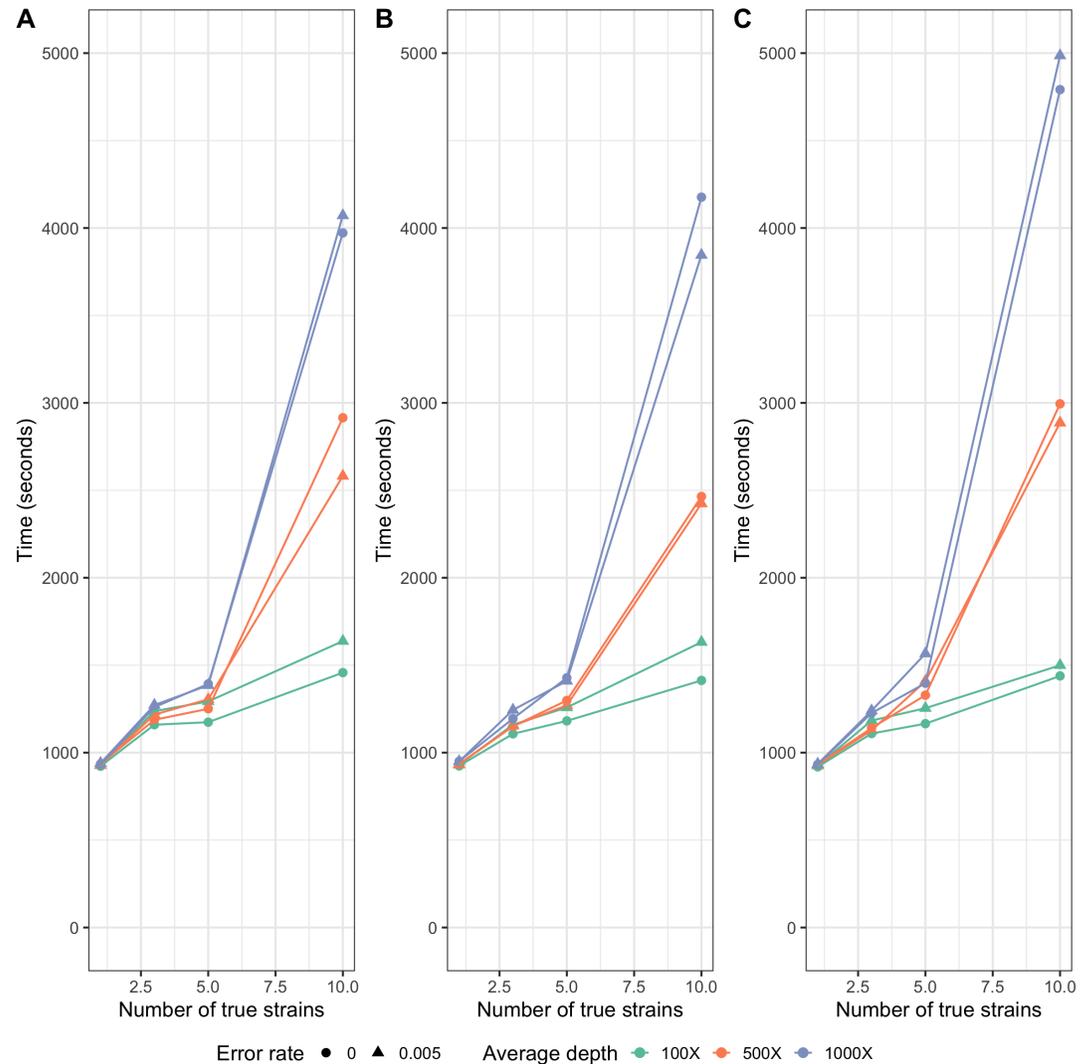
104 proportions of each candidate strains and compare them with the truth (Figure 2, 3 and 4).  
105 In most cases, the estimates are close to the true proportions, however, with a low coverage  
106 and high error rate, the proportions of the true strains will tend to be underestimated and strains  
107 that truly are not present will tend to be estimated as present in the sample. With one true strain in  
108 the sample, the proportion of this strain is always estimated to be 100%. For sufficiently high depth,  
109 e.g. 1000X corresponding to roughly a total of 30 Mb of data, the estimates of strain proportions  
110 are quite accurate, even when 10 strains are present and for strains with a proportion as low as 5%.  
111 There is similarly very little probability mass assigned to strains that are not truly in the sample. For  
112 example, for 150 bp paired-end reads with a +25 bp insert and 1000X average sequencing depth,  
113 the estimate of the cumulative average proportion of all strains not truly in the sample is 0.63%.  
114 The speed of the method is highly dependent on the number of true strains and the average  
115 depth (Figure 5), but for realistic sized data sets with a reference database of 3,117,131 strains, the  
116 typical computational time is between 15 minutes and two hours using a single core. This includes  
117 the initial time cost of ~10.5 minutes for reading the large panel of reference strains into mem-  
118 ory. There is no appreciable difference in speed between the different sequencing strategies used,  
119 except that paired-end 2x75bp sequences tends to take longer at higher average coverage. Simu-  
120 lations using the higher error rate (0.5%) are slower than simulations with no error. The average  
121 time for all sets of simulations with 5 or fewer true strains is <30 minutes for all coverages, while  
122 the average time for 10 true strains varies between ~24 to ~83 minutes depending on the average  
123 depth.  
124 In order to quantify the statistical evidence for the presence of a candidate strain in the sample,  
125 we propose a likelihood ratio test, LLR, formed by comparing the maximum likelihood value calcu-  
126 lated when the candidate strain is eliminated from the sample ( $p = 0$ ) to the maximum likelihood



**Figure 3.** Estimated proportions for simulated paired-end reads (2x150 bp with an insert size of +25 bp) with five replicates for when the sample truly contains 1 (A), 3 (B), 5 (C), or 10 (D) strains out of a total of 1,499,078 non-redundant candidate strains in the database. The red dashed lines indicate the true proportion of each strain. 'Other' indicates the sum of estimated proportions for all strains that are not truly represented in the sample.



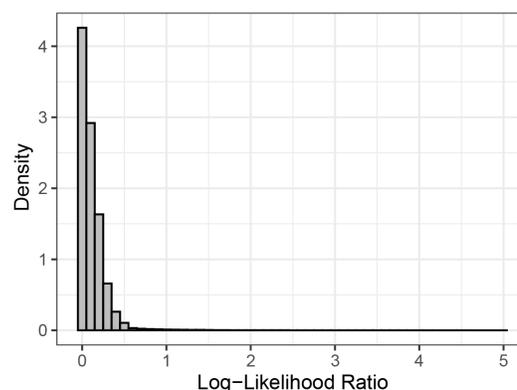
**Figure 4.** Estimated proportions for simulated paired-end reads (2x75 bp with an insert size of +25 bp) with five replicates for when the sample truly contains 1 (A), 3 (B), 5 (C), or 10 (D) strains out of a total of 1,499,078 non-redundant candidate strains in the database. The red dashed lines indicate the true proportion of each strain. 'Other' indicates the sum of estimated proportions for all strains that are not truly represented in the sample.



**Figure 5.** Average run times for single-end 300 bp (A), paired-end 2x150 bp (B), and paired-end 2x75 bp (C) read simulations using 100X, 500X, and 1000X average depth with an error rate of 0% and 0.5%. Each average run time reported is based on 5 replicates. Times were calculated using an AMD EPYC 7742 tetrahexaconta-core 2.25-3.40 GHz processor.

127 value calculated when allowing the strain to be present in the sample ( $p \geq 0$ ), where  $p$  is the pro-  
 128 portion of the strain in the sample (see *Methods and Materials*). Standard asymptotic theory for  
 129 the distribution of the likelihood ratio statistics does not apply to this situation for several reasons,  
 130 most importantly, a search is first made to find the strains that provide the largest increase in the  
 131 likelihood among many strains, and we only calculate the likelihood ratio for the strains with esti-  
 132 mates of  $p > 0$ . We, therefore, use simulations to evaluate the distribution of the likelihood ratio  
 133 test statistics under varying conditions. We simulated 1,000 data sets with different numbers of  
 134 true strains, coverage, read length and error rate and calculated the likelihood ratio for all strains  
 135 that were falsely inferred to be present in the sample (Figure 6). Since the frequency of  $LLR > 2$   
 136 and  $LLR > 4$  is about 0.001 and 0.0005, respectively, we recommend using 2 and 4 as thresholds  
 137 for strong and extremely strong evidence for presence of the strain in the sample.

138



**Figure 6.** Distribution of log-likelihood ratios from 1,000 data sets of simulated 100 ~ 300bp single-end reads. Those simulated data sets include 3 ~ 10 strains with proportions ranging from 5% to 50% and average depths ranging from 50X to 300X. The error rate varies from 0% to 0.5%.

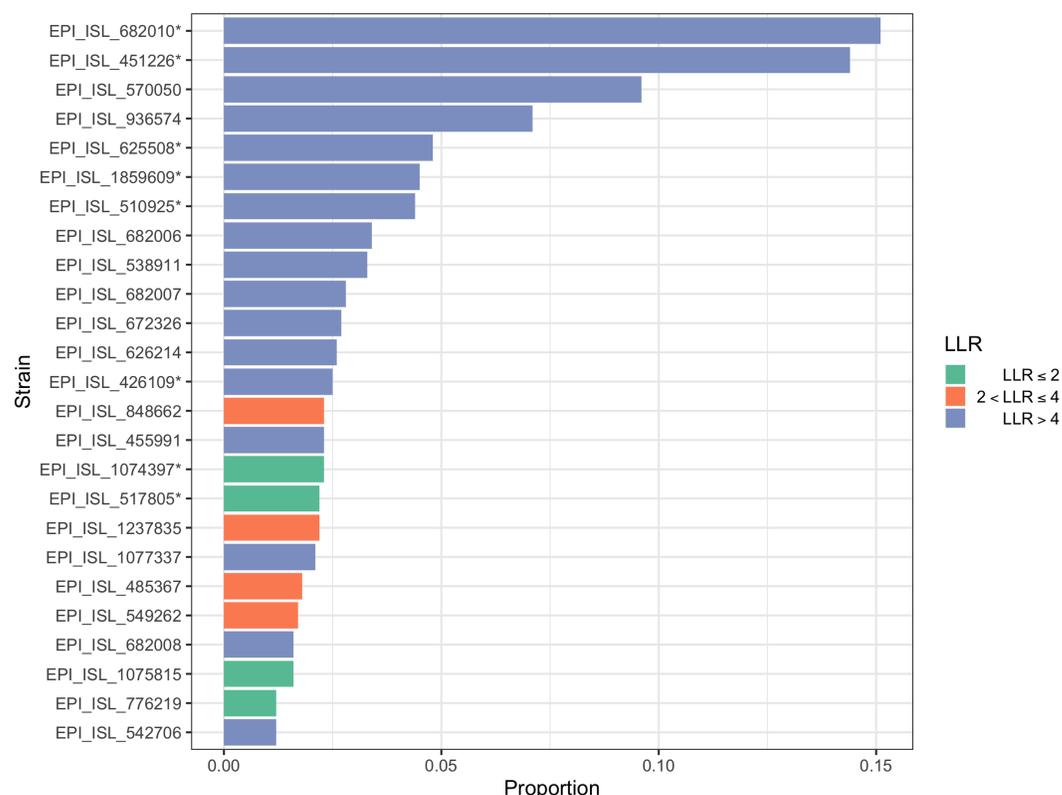
### 139 **Application to wastewater data from *Crits-Christoph et al. (2021)***

140 To apply our method to a published data set, we estimated the composition of SARS-CoV-2 lineages  
141 using wastewater shotgun sequencing data from *Crits-Christoph et al. (2021)* in Figure 7, which  
142 were all collected in the San Francisco Bay Area. Two out of the top ten strains were collected in  
143 Alameda county (EPI\_ISL\_625508, which is identical to EPI\_ISL\_625520, and EPI\_ISL\_672326), and  
144 the top five strains were all collected in North America.

### 145 **Discussion**

146 In order to allow for accurate inferences of strain composition, we first developed a new phylo-  
147 genetic method for data imputation for SARS-CoV-2 sequences. The method proved to be highly  
148 accurate with error rates comparable to, or lower, than typical sequencing error rates (Figure 1A). In  
149 fact, apparent wrongly inferred nucleotides may in many cases not be wrongly inferred but rather  
150 be inferences of the true allele, correcting a sequencing error in the reported sequence. Thus, sim-  
151 ilarly to imputation-based genotype calling in humans, this method could be used for correcting  
152 sequencing errors and incorporated formally into an algorithm of imputation-informed sequenc-  
153 ing where the quality scores from sequencing reads are combined with phylogenetically informed  
154 nucleotide probabilities to call nucleotides in each position. Computationally, this could be done  
155 simply by using the phylogenetic posterior probabilities of nucleotides as priors for genotype call-  
156 ing.

157 Our simulation results for the EM algorithm show that the new method can accurately estimate  
158 proportions of SARS-CoV-2 lineages in wastewater samples when up to 10 strains with frequencies  
159 as low as 5% are represented in the sample. Nonetheless, the estimated proportions for the true  
160 strains tend to be slightly lower than the actual proportions because the presence of other non-true  
161 strains is also estimated at low frequency. In order to have some probability for other non-true  
162 strains to be estimated, the true proportions for the true strains will naturally in average be slightly  
163 underestimated. In all sets of simulations of single-end 300bp reads (Figure 2), paired-end 2 x 75bp  
164 (Figure 3), and paired-end 2 x 150bp (Figure 4), the estimated proportions of the true strains tend  
165 to be more accurate as sequencing depth increases. When there are many strains (i.e., when there  
166 are 10 true strains) and sequencing depth is low (i.e., 100X), there is a high degree of noise in the  
167 data set. However, as the total sequencing depth increases, the estimates become progressively  
168 more accurate. We recommend that studies focused on identifying different strains of SARS-CoV-2  
169 in environmental samples aim to achieve an average depth of 1000X. Additionally, the method pre-  
170 sented here has only been evaluated for the estimation of proportions of strains with a frequency  
171 of 5% or larger. We recommend that strains identified in the sample at low frequencies are evalu-



**Figure 7.** Estimated proportions of the top 25 strains estimated from wastewater shotgun sequencing data from *Crits-Christoph et al. (2021)* and their log-likelihood ratios. Strains with an asterisk (\*) are identical with other strains. EPI\_ISL\_682010\* is identical to EPI\_ISL\_682025, EPI\_ISL\_1373628, EPI\_ISL\_1373632, and EPI\_ISL\_1373659. EPI\_ISL\_451226\* is identical to EPI\_ISL\_451227 and EPI\_ISL\_455983. EPI\_ISL\_625508\* is identical to EPI\_ISL\_625520, EPI\_ISL\_672318, EPI\_ISL\_672449, EPI\_ISL\_739003, EPI\_ISL\_739029, EPI\_ISL\_739135, EPI\_ISL\_739161, EPI\_ISL\_739207, and EPI\_ISL\_739286. EPI\_ISL\_1859609\* is identical to EPI\_ISL\_1859762. EPI\_ISL\_510925\* is identical to EPI\_ISL\_510926. EPI\_ISL\_426109\* is identical to EPI\_ISL\_486012, EPI\_ISL\_570168, EPI\_ISL\_570172, EPI\_ISL\_576500, and EPI\_ISL\_576501. EPI\_ISL\_1074397\* is identical to EPI\_ISL\_2190584. EPI\_ISL\_517805\* is identical to EPI\_ISL\_527398 and EPI\_ISL\_137362.

172 ated using the likelihood ratio test as they likely could be false positives.

173 Current strategies for monitoring community composition of SARS-CoV-2 strains include se-  
174 quencing a large number of clinical samples. As SARS-CoV-2 becomes endemic, tracking the rela-  
175 tive prevalence in local communities of different SARS-CoV-2 strains will be highly costly. Further-  
176 more, the use of clinical samples is associated with a lag from infection onset to hospitalization.  
177 Our results suggest an alternative strategy of monitoring using wastewater samples. Wastewater  
178 sequencing has already proved effective for tracking SARS-CoV-2 abundance ((*Korber et al., 2020*;  
179 *Rockett et al., 2020*)). With the computational framework developed here, it also promises to be-  
180 come an important cost-effective strategy for monitoring the local composition of different viral  
181 strains.

## 182 **Methods and Materials**

### 183 **SARS-CoV-2 Reference Database**

184 To build the SARS-CoV-2 reference database, a multiple sequence alignment (MSA) of 3,117,131  
185 SARS-CoV-2 genomes (msa\_2021-10-15.tar.xz) and the corresponding phylogenetic tree (GISAID-  
186 hCoV-19-phylogeny-2021-10-13.zip) was downloaded from GISAID ([www.gisaid.org](http://www.gisaid.org)) on October 16,  
187 2021. We pruned sequence EPI\_ISL\_4989640 from the tree since it was not present in the MSA. We  
188 use the function `collapse.singles` to collapse elbow nodes (i.e., nodes other than the root with  
189 two degrees) and `multi2di` to resolve multichotomies in the R `ape` package (*Paradis et al., 2004*).  
190 We impute missing data (i.e., every position in the MSA that did not contain an A, G, C, or T), using  
191 the phylogenetic tree. To do so, we first scale the branch lengths in terms of substitutions per site  
192 by dividing each reported branch length by the average sequence length (29618.5). For branch  
193 lengths that were reported to be 0, we define them to be 0.01 divided by the average sequence  
194 length. We impute missing nucleotides using the maximum of the posterior probability of each  
195 nucleotide in the leaf nodes under a standard Jukes and Cantor model (*Jukes et al., 1969*), using  
196 standard computational algorithms (*Yang, 2014*). In brief, because the model is time-reversible,  
197 the root can be placed in any particular node, and the fractional likelihoods (joint probabilities of a  
198 fraction of the data in the leaf nodes and the nucleotide state in the node) can be pulled recursively  
199 towards the node from both the child nodes and the parental node. The posterior probability in  
200 the leaf nodes of a nucleotide is calculated as the product of the stationary probability of the nu-  
201 cleotide multiplied by the fractional likelihood in the leaf node conditioned on the data in all other  
202 leaf nodes. This can be programmed so the calculation is linear in the number of leaf nodes using  
203 a single pre-order and a single post-order traversal of the tree that will calculate the posterior prob-  
204 abilities in all nodes. We note that other models than the Jukes and Cantor model could provide  
205 more accurate estimates, but at a computational cost.

206 Since calculating fractional likelihoods for the entire tree requires more RAM than was computa-  
207 tionally feasible for us (~72TB of RAM), we split the tree into partitions, and process each partition  
208 sequentially as follows:

209 Each internal node in the tree corresponds to a partition of leaf nodes into three sets. First, we iden-  
210 tify the node with the minimum variance in the number of elements among these three partitions,  
211 i.e. we find

$$\min_{n \in T} \left( \frac{(n_a - \frac{n_1+n_2+n_a}{3})^2 + (n_1 - \frac{n_1+n_2+n_a}{3})^2 + (n_2 - \frac{n_1+n_2+n_a}{3})^2}{3} \right) \quad (1)$$

212 where  $n$  is a node in the tree,  $T$  is the tree,  $n_1$  is the number of leaf nodes descending from the  
213 left child of  $n$ ,  $n_2$  is the number of leaf nodes descending from the right child, and  $n_a = N - n_1 - n_2$ ,  
214 where  $N$  is the total number of leaf nodes in the tree. We then split the tree into 3 subtrees by  
215 eliminating the identified node. We then iterate this procedure for the resulting subtrees until all  
216 trees contain at most 50,000 leaf nodes.

217 Using this partitioning procedure, we obtain 121 trees which we use to calculate the posterior  
218 probabilities at each site. After imputation, we trim the MSA to begin at the start of the Wuhan  
219 reference sequence (Wuhan-Hu-1), position 55 in the MSA, and we removed every position in the

220 MSA that contains a gap in Wuhan-Hu-1. After this trimming and imputation process, we save non-  
 221 informative invariant sites (856 sites), in order to reduce running time when eliminating unlikely  
 222 strains. We also remove all identical sequences, resulting in 1,499,078 non-redundant genomes.

### 223 Estimating the proportions of SARS-CoV-2 genomes

224 All sequencing reads are aligned to Wuhan-Hu-1 (NC\_045512.2) using bowtie2 (*Langmead and*  
 225 *Salzberg, 2012*) with the following command for single-end reads, `bowtie2 -all -f -x wuhCor1`  
 226 `-U`, and for paired-end reads, `bowtie2 -all -f -x wuhCor1 -1 -2`. For each read data set, we first  
 227 remove unlikely genomes from the candidate strain alignment by eliminating genomes with SNP  
 228 alleles that have an allele frequency in the read data less than a user-defined frequency thresh-  
 229 old. For the analyses in this data, that threshold was set to 0.01. This typically reduced the size of  
 230 the alignment to  $< 1,000$  relevant genomes. Using this reduced set of SARS-CoV-2 genomes, we  
 231 calculate a matrix of dimensions (number of reads) $\times$ (number of genomes) containing the number  
 232 of mismatches between each sequencing read and each genome,  $d = \{d_{ij}\}$ . For paired-end reads  
 233 with reads that overlap, we use the consensus nucleotide. If there is a conflict at any position in the  
 234 overlap of the paired-end reads, we omit this site. Based on the mismatch matrix,  $d$ , we first calcu-  
 235 late the probability of observing read  $j$  given that it comes from strain  $i$ , denoted as  $q_{ij}$ . Assuming  
 236 that the reads are independent (PCR clones removed) and a user-defined error rate  $\alpha$  (default =  
 237 0.005) at each nucleotide, this probability is given by

$$q_{ij} = \alpha^{d_{ij}} \times (1 - \alpha)^{n_j - d_{ij}}$$

238 where  $n_j$  is the length of read  $j$  and  $d_{ij}$  is the number of mismatches in read  $j$  given that it comes  
 239 from strain  $i$ . The log-likelihood is then given by

$$\log L(p_1, \dots, p_k) = \sum_{j=1}^n \log \sum_{i=1}^k q_{ij} p_i, \quad (2)$$

240 where  $p_i$  ( $i = 1, \dots, k$ ) is the proportion of strain  $i$ , i.e. the parameters we wish to estimate. We then  
 241 use the standard Expectation Maximization (EM) algorithm (*Dempster et al., 1977*) to maximize the  
 242 likelihood function with respect to these parameters 1:

---

#### Algorithm 1 EM algorithm for estimating the proportions of candidate strains

---

**Input:** The probability of observing read  $j$  given that it comes from strain  $i$ ,  $q_{ij}$ , for all  $i$  and  $j$ .

**Output:** The proportion of each candidate strain,  $p_i$ , for all  $i$ .

- 1: Initialize the proportions of each strain  $p_i(0), i = 1 \dots k$ , with uniform probabilities  $U(0, 1)$  and then re-scaled to 1.
  - 2: Compute the log-likelihood  $\ell_0 = \sum_{j=1}^n \log \sum_{i=1}^k q_{ij} p_i(0)$ ;
  - 3: **repeat**
  - 4:   Compute the proportion of each candidate strain at iteration  $t$  as  $p_i(t) = \frac{1}{n} \sum_{j=1}^n \frac{p_i(t-1)q_{ij}}{\sum_{i=1}^k p_i(t-1)q_{ij}}$ ;
  - 5:   Compute the log-likelihood at iteration  $t$  as  $\ell_t = \sum_{j=1}^n \log \sum_{i=1}^k q_{ij} p_i(t)$ ;
  - 6: **until**  $\ell_t - \ell_{t-1} < \epsilon$ , where  $\epsilon$  is a pre-defined stopping criterion.
- 

243 However, Algorithm 1 usually has a slow convergence rate, especially when the number of candi-  
 244 date strains  $k$  is large. Therefore, to accelerate the Algorithm 1, we use the SQUAREM algorithm  
 245 proposed by *Varadhan and Roland (2008)* with its implementation in the R package `turboEM` (*Bobb*  
 246 *and Varadhan, 2020*).

### 247 Determining unidentifiable strains

248 Note that if two stains have the same  $q_{ij}$ 's, say there exist  $i$  and  $i'$  such that  $q_{ij} = q_{i'j}$  for all  $j = 1, \dots, n$ ,  
 249 the log-likelihood (2) becomes

$$\log L = \sum_{j=1}^n \log \left[ \left( \sum_{r \in \{1, \dots, k\} \setminus \{i, i'\}} q_{rj} p_r \right) + q_{ij} (p_i + p_{i'}) \right]. \quad (3)$$

250 Therefore, as long as  $p_i + p_{i'}$  is fixed, (3) remains the same no matter what value  $p_i$  and  $p_{i'}$  take,  
251 making the model unidentifiable. To solve this problem, we gather strains with the same  $\{q_{ij}\}_{j=1}^n$   
252 into an unidentifiable group and estimate its overall proportion instead of the proportions of each  
253 strain in it.

### 254 **Quantifying the statistical evidence of the existence of each candidate strain**

255 To provide a measure of statistical support for the presence of strain  $i_0$ , i.e.  $p_{i_0} > 0$ , we remove strain  
256  $i_0$  from the candidate set of strains and re-run Algorithm 1 providing a new estimate  $\{\tilde{p}_i\}_{i=1}^k$  with  
257  $\tilde{p}_{i_0} = 0$ . Using (2), we can then calculate the difference in log likelihood before and after removing  
258 strain  $i_0$ , denoted as  $LLR_{i_0}$ . From our simulations (see Results), we recommend using  $LLR_{i_0} \geq 4$  as  
259 strong statistical evidence in favor of existence of strain  $i_0$  in the sample.

### 260 **Simulating missing data for imputation**

261 For every SARS-CoV-2 genome (out of a total of 3,117,131 genomes), we randomly remove 1% of  
262 nucleotides, and save the true nucleotide at each position that was removed. We then use the *Tree*  
263 *imputation* method and the *Common allele* method to impute the nucleotides that are missing.

### 264 **Simulating reads from SARS-CoV-2 genomes**

We choose 10 strains among 1,499,078 strains uniformly at random. Then, to simulate single-end  
reads from a strain, we choose a starting point uniformly at random and let it extend  $m_0$  bps, where  
 $m_0$  is the read length. For paired-end reads, we similarly choose a starting point at random and let  
it extend  $m_0$  bps. Then, starting from the end of this read, if the insert size is  $m_1$  is positive, we  
simulate the start of the reverse read  $m_1$  bps forward with length  $m_0$ ; if  $m_1$  is negative, we simulate  
the start of the reverse read  $m_1$  bps backwards. We then add sequencing errors independently  
with probability  $\alpha = 0.005$  at each site. Errors are induced by relabeling the nucleotide to any of the  
other three possible nucleotides with the following probabilities used in *Stephens et al. (2016)*:

	A	C	G	T
A	0	0.4918	0.3377	0.1705
C	0.5238	0	0.2661	0.2101
G	0.3754	0.2355	0	0.3890
T	0.2505	0.2552	0.4942	0

### 265 **Calculating time cost**

266 To calculate running time of the method we use `/usr/bin/time` on an AMD EPYC 7742 tetrahexaconta-  
267 core 2.25-3.40 GHz processor and report `real time` in the results (Figure 5). The running time that  
268 we calculate includes running the method from start (reading in the reference strains) to finish  
269 (reporting proportions) and includes the filtering step for eliminating unlikely strains. We report  
270 times that do not include calculating the log-likelihood ratio.

### 271 **Applying the method to wastewater data from *Crits-Christoph et al. (2021)***

272 Wastewater shotgun sequencing data from *Crits-Christoph et al. (2021)* was downloaded from  
273 NCBI BioProject ID PRJNA661613 (<https://www.ncbi.nlm.nih.gov/bioproject/?term=PRJNA661613>).  
274 All samples were pooled together and aligned against Wuhan-Hu-1 using BWA-MEM (*Li, 2013*) to  
275 identify SARS-CoV-2 reads.

### 276 **Data availability**

277 Simulations used in this manuscript can be downloaded at <https://doi.org/10.5281/zenodo.5838942>.  
278 The imputed MSA can be downloaded at <https://doi.org/10.5281/zenodo.5838946>. Identical strains  
279 are contained in the headers of the MSA separated by colons. Software for the method is available  
280 for download at [https://github.com/lpipes/SARS\\_CoV\\_2\\_wastewater\\_surveillance](https://github.com/lpipes/SARS_CoV_2_wastewater_surveillance).

## 281 **Competing interests**

282 We declare that we have no known competing financial interests or personal relationships that  
283 influenced this work.

## 284 **Acknowledgments**

285 We gratefully acknowledge all laboratories who submitted SARS-CoV-2 genome sequences to the  
286 GISAID EpiCoV database ([www.gisaid.org](http://www.gisaid.org)), which we used for the reference database for this method.  
287 We acknowledge Xiaoyi Gu for testing the software and for development of a website portal for the  
288 method, and Selina Kim for working on this project.

## 289 **Funding**

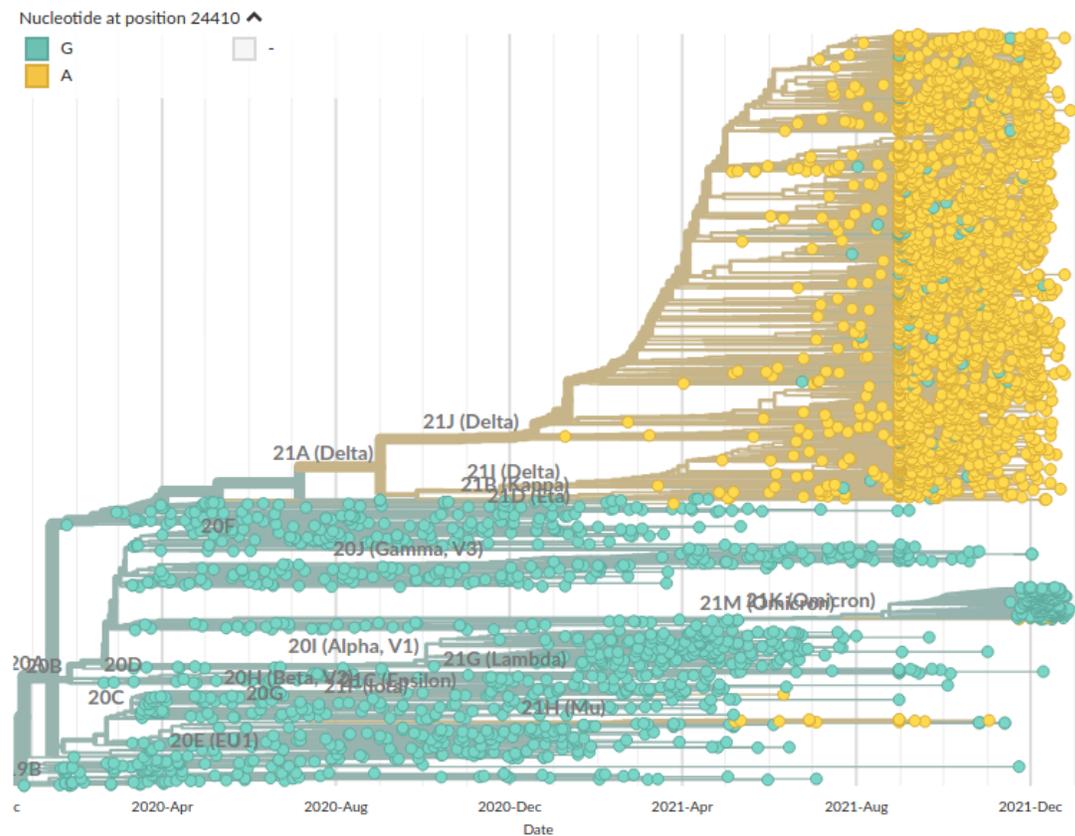
290 This work used the Extreme Science and Engineering Discovery Environment (XSEDE) Bridges-2 sys-  
291 tem at the Pittsburgh Supercomputing Center through allocation BIO180028 and was supported  
292 by NIH grant 1R01GM138634-01.

## 293 **References**

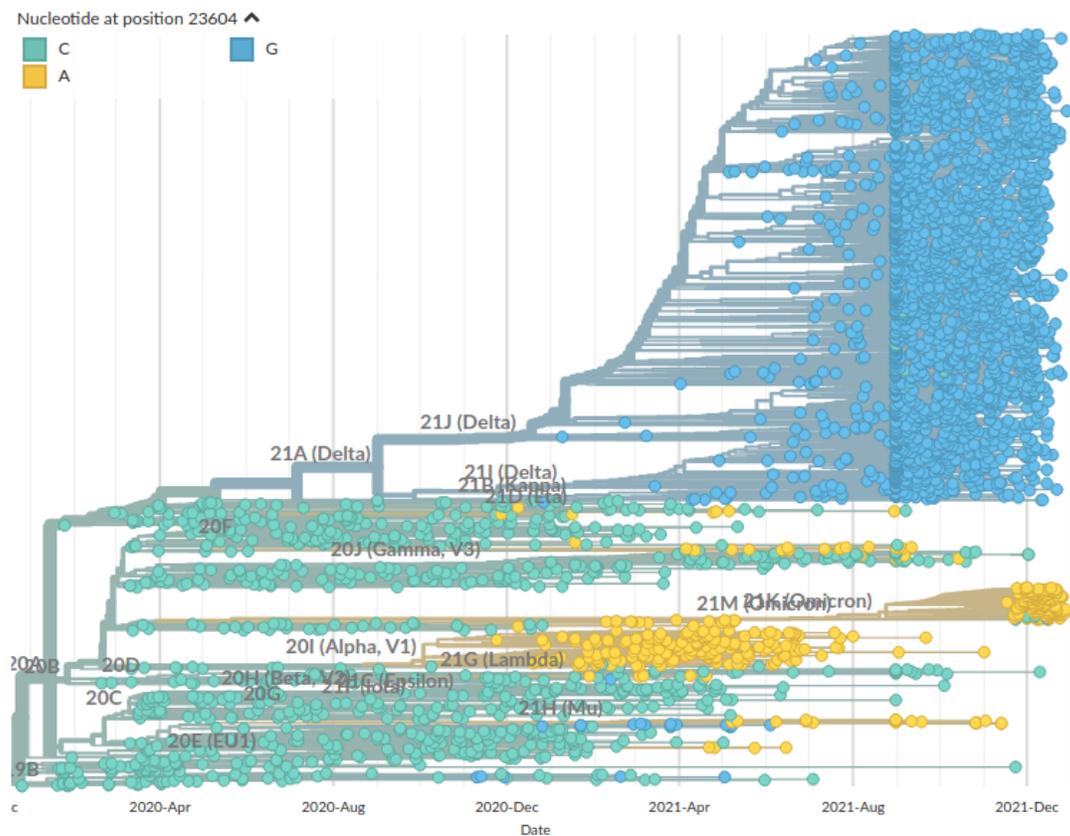
- 294 **Ahmed W**, Angel N, Edson J, Bibby K, Bivins A, O'Brien JW, Choi PM, Kitajima M, Simpson SL, Li J, et al. First con-  
295 firmed detection of SARS-CoV-2 in untreated wastewater in Australia: A proof of concept for the wastewater  
296 surveillance of COVID-19 in the community. *Science of The Total Environment*. 2020; p. 138764.
- 297 **Bobb JF**, Varadhan R. turboEM: A Suite of Convergence Acceleration Schemes for EM, MM and Other Fixed-Point  
298 Algorithms; 2020, <https://CRAN.R-project.org/package=turboEM>, r package version 2020.1.
- 299 **Crits-Christoph A**, Kantor RS, Olm MR, Whitney ON, Al-Shayeb B, Lou YC, Flamholz A, Kennedy LC, Greenwald  
300 H, Hinkle A, et al. Genome sequencing of sewage detects regionally prevalent SARS-CoV-2 variants. *MBio*.  
301 2021; 12(1):e02703–20.
- 302 **Dempster AP**, Laird NM, Rubin DB. Maximum likelihood from incomplete data via the EM algorithm. *Journal*  
303 *of the Royal Statistical Society: Series B (Methodological)*. 1977; 39(1):1–22.
- 304 **Faria NR**, Morales Claro I, Candido D, Moyses Franco L, Andrade PS, Coletti TM, Silva CA, Sales FC, Manuli  
305 ER, Aguiar RS, et al., Genomic characterisation of an emergent SARS-CoV-2 lineage in Manaus: preliminary  
306 findings. *Virological*; 2021.
- 307 **Farkas K**, Hillary LS, Malham SK, McDonald JE, Jones DL. Wastewater and public health: the potential of wastew-  
308 ater surveillance for monitoring COVID-19. *Current Opinion in Environmental Science & Health*. 2020; 17:14–  
309 20.
- 310 **Jukes TH**, Cantor CR, et al. Evolution of protein molecules. *Mammalian protein metabolism*. 1969; 3:21–132.
- 311 **Korber B**, Fischer WM, Gnanakaran S, Yoon H, Theiler J, Abfalterer W, Hengartner N, Giorgi EE, Bhattacharya  
312 T, Foley B, et al. Tracking changes in SARS-CoV-2 Spike: evidence that D614G increases infectivity of the  
313 COVID-19 virus. *Cell*. 2020; 182(4):812–827.
- 314 **Kumar M**, Patel AK, Shah AV, Raval J, Rajpara N, Joshi M, Joshi CG. First proof of the capability of wastewater  
315 surveillance for COVID-19 in India through detection of genetic material of SARS-CoV-2. *Science of The Total*  
316 *Environment*. 2020; 746:141326.
- 317 **Langmead B**, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature methods*. 2012; 9(4):357–359.
- 318 **Li H**. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv preprint*  
319 *arXiv:13033997*. 2013; .
- 320 **Marchini J**, Howie B. Genotype imputation for genome-wide association studies. *Nature Reviews Genetics*.  
321 2010; 11(7):499–511.
- 322 **Medema G**, Heijnen L, Elsinga G, Italiaander R, Brouwer A. Presence of SARS-Coronavirus-2 RNA in sewage  
323 and correlation with reported COVID-19 prevalence in the early stage of the epidemic in the Netherlands.  
324 *Environmental Science & Technology Letters*. 2020; 7(7):511–516.

- 325 **Medema G**, Heijnen L, Elsinga G, Italiaander R, Brouwer A. Presence of SARS-Coronavirus-2 RNA in Sewage  
326 and Correlation with Reported COVID-19 Prevalence in the Early Stage of the Epidemic in The Netherlands.  
327 *Environmental Science & Technology Letters*. 2020; 7(7):511–516. <https://doi.org/10.1021/acs.estlett.0c00357>,  
328 doi: 10.1021/acs.estlett.0c00357.
- 329 **Paradis E**, Claude J, Strimmer K. APE: Analyses of Phylogenetics and Evolution in R language. *Bioinformatics*.  
330 2004 01; 20(2):289–290. <https://doi.org/10.1093/bioinformatics/btg412>, doi: 10.1093/bioinformatics/btg412.
- 331 **Rockett RJ**, Arnott A, Lam C, Sadsad R, Timms V, Gray KA, Eden JS, Chang S, Gall M, Draper J, et al. Revealing  
332 COVID-19 transmission in Australia by SARS-CoV-2 genome sequencing and agent-based modeling. *Nature*  
333 *medicine*. 2020; 26(9):1398–1404.
- 334 **Shu Y**, McCauley J. GISAID: Global initiative on sharing all influenza data—from vision to reality. *Eurosurveillance*.  
335 2017; 22(13):30494.
- 336 **Stephens ZD**, Hudson ME, Mainzer LS, Taschuk M, Weber MR, Iyer RK. Simulating Next-Generation Sequencing  
337 Datasets from Empirical Mutation and Sequencing Models. *PLOS ONE*. 2016 11; 11(11):1–18. [https://doi.org/](https://doi.org/10.1371/journal.pone.0167047)  
338 [10.1371/journal.pone.0167047](https://doi.org/10.1371/journal.pone.0167047), doi: 10.1371/journal.pone.0167047.
- 339 **Tang JW**, Tambyah PA, Hui DS. Emergence of a new SARS-CoV-2 variant in the UK. *Journal of Infection*. 2020; .
- 340 **Varadhan R**, Roland C. Simple and globally convergent methods for accelerating the convergence of any EM  
341 algorithm. *Scandinavian Journal of Statistics*. 2008; 35(2):335–353.
- 342 **Volz E**, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, Hinsley WR, Laydon DJ, Dabrera G, O’Toole Á,  
343 et al. Transmission of SARS-CoV-2 Lineage B. 1.1. 7 in England: Insights from linking epidemiological and  
344 genetic data. *medRxiv*. 2021; p. 2020–12.
- 345 **Weber S**, Ramirez CM, Weiser B, Burger H, Doerfler W. SARS-CoV-2 Worldwide Replication Drives Rapid Rise  
346 and Selection of Mutations across the Viral Genome: A Time-Course Study Potential Challenge for Vaccines  
347 and Therapies. *medRxiv*. 2021; .
- 348 **Wu F**, Zhang J, Xiao A, Gu X, Lee WL, Armas F, Kauffman K, Hanage W, Matus M, Ghaeli N, Endo N, Duvallet C,  
349 Poyet M, Moniz K, Washburne AD, Erickson TB, Chai PR, Thompson J, Alm EJ. SARS-CoV-2 Titers in Wastewater  
350 Are Higher than Expected from Clinically Confirmed Cases. *mSystems*. 2020; 5(4). [https://msystems.asm.org/](https://msystems.asm.org/content/5/4/e00614-20)  
351 [content/5/4/e00614-20](https://msystems.asm.org/content/5/4/e00614-20), doi: 10.1128/mSystems.00614-20.
- 352 **Yang Z**. *Molecular evolution: a statistical approach*. Oxford University Press; 2014.
- 353 **Zhang W**, Davis B, Chen SS, Martinez JS, Plummer JT, Vail E. Emergence of a novel SARS-CoV-2 strain in Southern  
354 California, USA. *medRxiv*. 2021; .

355 **Supplementary Material**



**Figure S1.** Screenshot of SARS-CoV-2 phylogeny from nextstrain.org for nucleotide position 24,410 taken on January 5, 2022.



**Figure S2.** Screenshot of SARS-CoV-2 phylogeny from nextstrain.org for nucleotide position 23,604 taken on January 5, 2022.