

# 1 **Monitoring SARS-CoV-2 Populations in Wastewater by Amplicon Sequencing** 2 **and Using the Novel Program SAM Refiner**

3 Devon A. Gregory<sup>1</sup>, Chris G. Wieberg<sup>2</sup>, Jeff Wenzel<sup>3</sup>, Chung-Ho Lin<sup>4</sup> and Marc C. Johnson<sup>5,\*</sup>

4 <sup>1</sup> Dept. of Molecular Microbiology and Immunology; University of Missouri, Columbia, MO 65201, USA;  
5 gregoryde@missouri.edu

6 <sup>2</sup> Director, Water Protection Program, Missouri Department of Natural Resources

7 <sup>3</sup> Bureau Chief, Bureau of Environmental Epidemiology, Division of Community and Public Health, Missouri  
8 Department of Health and Senior Services

9 <sup>4</sup> Center for Agroforestry, University of Missouri-Columbia, USA; School of Natural Resources, University of  
10 Missouri-Columbia, USA

11 <sup>5</sup> Dept. of Molecular Microbiology and Immunology; University of Missouri, Columbia, MO 65201, USA;  
12 marcjohnson@missouri.edu

13 \* Correspondence: marcjohnson@missouri.edu

14 **Abstract:** Sequencing SARS-CoV-2 from wastewater has become a useful tool in monitoring the  
15 spread of variants. We use a novel computation workflow with SARS-CoV-2 amplicon sequencing  
16 in order to track wastewater populations of the virus. As part of this workflow, we developed a  
17 program for both variant reporting and removal of PCR generated chimeric sequences. With  
18 these methods, we are able to track viral population dynamics over time. We observe the  
19 emergence of the variants of concern B.1.1.7 and P.1, and their displacement of the D614G B.1  
20 variant.

21 **Keywords:** Coronavirus; Wastewater; Metagenomics; Molecular Epidemiology

22

## 23 **1. Introduction**

24 SARS-CoV-2 became pandemic and caused a world-wide health crisis starting in 2020 [1]. Full  
25 genome sequences of SARS-CoV-2 were rapidly made available within the first months of spread  
26 [2, 3]. Partial and whole genome sequencing of SARS-CoV-2 has been an important tool in  
27 monitoring transmission paths and the emergence of variant lineages. Most sequencing of SARS-  
28 CoV-2 has been done on clinical samples. However, early in the SARS-CoV-2 pandemic,  
29 wastewater began to be used to track community levels and spread of SARS-CoV-2 by RT-qPCR  
30 methods [4, 5]. Investigators have also used high throughput sequencing on wastewater samples  
31 to obtain full or partial SARS-CoV-2 genomic sequences which were used for metagenomic and  
32 epidemiologic analysis [6, 7, 8, 9, 10, 11, 12, 13]. Sequences identified in wastewater samples  
33 may reflect known lineages, as well as lineages not reported from clinical samples. Combinations  
34 of mutations not observed in clinical samples may represent new infections not yet picked up by  
35 clinical sampling or lineages that are under-represented in clinical samples. Approaches using  
36 wastewater are particularly relevant with the emergence of variant lineages that may vary from  
37 previous isolates in their fitness and/or disease.

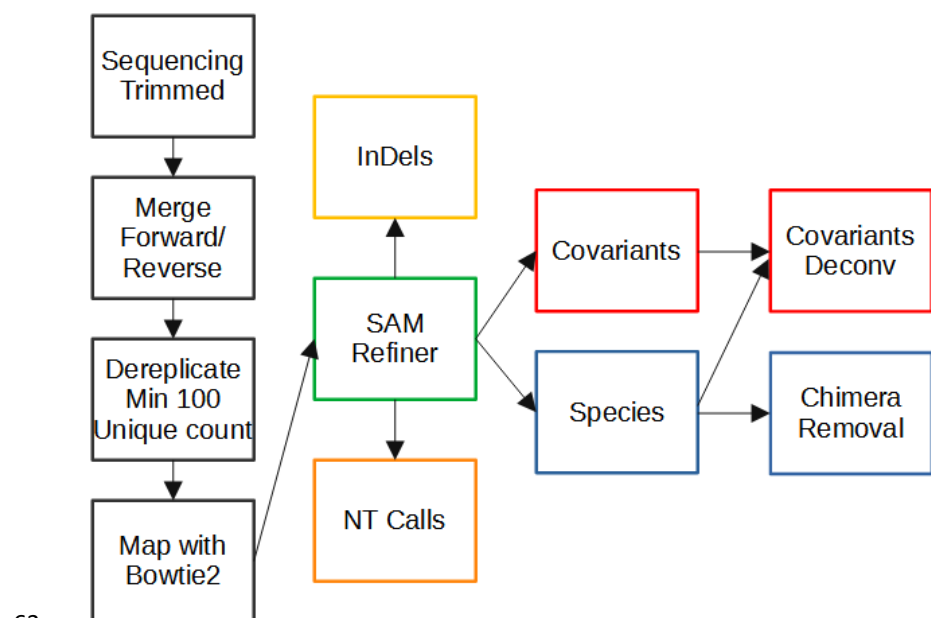
38

**NOTE:** This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

39 The state of Missouri has been monitoring wastewater with RT-qPCR to track the prevalence and  
40 spread of SARS-CoV-2  
41 (<https://storymaps.arcgis.com/stories/f7f5492486114da6b5d6fdc07f81aacf>). We sought to begin  
42 using the same samples for high throughput sequencing to track the presence and spread of  
43 known and previously unreported variant lineages. We were specifically interested in the spike  
44 gene and used primers to target 3 regions for amplification, the N-terminal domain (NTD), receptor  
45 binding domain (RBD) and the region of the S1 and S2 subunit split (S1S2). We chose these  
46 regions due to the numerous variations matching evolving lineages found in them and their  
47 significance in potential immune evasion [14]. While there are a number of high throughput  
48 sequencing technologies and methods, the sequence output is relatively standard; the processing  
49 and analysis of that sequence data is not. There are numerous programs and pipelines that can  
50 be used to obtain information from sequences and remove errors generated from PCR, such as  
51 single nucleotide polymorphisms (SNPs) and chimeric sequences. While many of these are  
52 quality approaches, we were unable to find a simple program or workflow with existing programs  
53 that provided easily human readable output that detailed variant lineages with the information we  
54 wanted and with sufficient removal of chimeric sequences. Specifically, we wished to include  
55 deletion and insertion events as well as SNPs and multiple nucleotide polymorphisms (MNP)s in  
56 our analysis and be able to view linked variances as single lineages. We also wished to be able to  
57 view downstream amino acid changes and have removal of chimeric sequences generated from  
58 PCR.

59

60 Here we detail the workflow (Fig1) we used to analyze high throughput sequencing data and the  
61 program we developed to provide a human readable, information dense output for viewing  
62 lineages.



63

64 **Figure 1. Workflow of Amplicon Sequencing Analysis.** Computational processing of sequencing results prior to  
65 the use of SAM Refiner is seen in the black boxes. Paired end reads generated from an Illumina MiSeq were trimmed  
66 of low quality calls at the end of the reads. Paired end reads were then merged into single contiguous reads. Reads  
67 were then dereplicated to unique sequences with at least 100 counts while preserving the count information in the  
68 sequence IDs. Dereplicated sequences were then mapped to the sequence of the SARS-CoV-2 Spike ORF using

69 Bowtie2. SAM Refiner was then used to process the mapped reads to obtain information about the variant lineages  
70 observed, initially outputting 4 TSV files to report unique sequences, nt calls, indels and covariants. The unique  
71 sequences and covariants were further processed to remove chimeric PCR artifacts to produce covariant  
72 deconvolution and chimera removed outputs.

## 73 **2. Materials and Methods**

### 74 **2.1. Wastewater Collection**

75 Twenty-four hour composite samples were collected at wastewater treatment facilities (WWTF)  
76 and were maintained at 4°C until they were delivered to the analysis lab, generally within 24 hours  
77 of collection. Samples reported in this study were collected at the NPSD Interim Saline Creek  
78 Regional WWTF in Fenton, MO.

### 79 **2.2. RNA Extraction**

80 Wastewater samples were centrifuged at 3,000xg for 10 minutes and then filtered through a 0.22  
81 µM polyethersulfone membrane (Millipore). Approximately 37.5 mL of wastewater was mixed with  
82 12.5 mL solution containing 50% (w/vol) polyethylene glycol 8,000 and 1.2 M NaCl, mixed, and  
83 incubated at 4°C for at least 1 hr. Samples were then centrifuged at 12,000xg for 2h at 4°C.  
84 Supernatant was decanted and RNA was extracted from the remaining pellet (usually not visible)  
85 with the QIAamp Viral RNA Mini Kit (Qiagen) using the manufacturer's instructions. RNA was  
86 extracted in a final volume of 60 µL.

87

### 88 **2.3. Sequencing**

89 The primary RT-PCR (25 µl) was performed with 5 microliters of RNA extracted from wastewater  
90 samples with loci specific primers (0.5 µM each) shown in Table 1 using the Superscript IV One-  
91 Step RT-PCR System (Thermo Fisher). Primary RT-PCR amplification was performed as follows:  
92 25°C(2:00) + 50°C(20:00) + 95°C(2:00) + [95°C(0:15) + 55°C(0:30) + 72°C(1:00)] x 25 cycles.  
93 Secondary PCR (25 µl) was performed using 5 ul of the primary PCR as template with gene  
94 specific primers containing 5' adapter sequences (0.5 µM each), dNTPs (100 µM each) and Q5  
95 DNA polymerase (NEB). Secondary PCR amplification was performed as follows: 95°C(2:00) +  
96 [95°C(0:15) + 55°C(0:30) + 72°C(1:00)] x 20 cycles. A tertiary PCR (50 µl) was performed to add  
97 adapter sequences required for Illumina cluster generation with forward and reverse primers (0.2  
98 µM each), dNTPs (200 µM each), and Phusion High-Fidelity DNA Polymerase (1U). PCR  
99 amplification was performed as follows: 98°C(3:00) + [98°C(0:15) + 50°C(0:30) + 72°C(0:30)] x 7  
100 cycles +72°C(7:00). Amplified product (10 µl) from each PCR reaction is combined and  
101 thoroughly mixed to make a single pool. Pooled amplicons were purified by addition of Axygen  
102 AxyPrep MagPCR Clean-up beads in a 1.0 ratio to purify final amplicons. The final amplicon  
103 library pool was evaluated using the Agilent Fragment Analyzer automated electrophoresis  
104 system, quantified using the Qubit HS dsDNA assay (Invitrogen), and diluted according to  
105 Illumina's standard protocol. The Illumina MiSeq instrument was used to generate paired-end  
106 300 base pair length reads. Adapter sequences were trimmed from output sequences using  
107 cutadapt [15]. The raw and trimmed reads for the samples used in this report are available at  
108 [https://github.com/degregory/SR\\_manuscript/tree/master/Fenton\\_Data](https://github.com/degregory/SR_manuscript/tree/master/Fenton_Data).

109

Region	PCR	Orientation	Primer Sequences
RBD	Primary	forward	CTGCTTTACTAATGTCTATGCAGATTG
	Primary	reverse	TCCTGATAAAGAACAGCAACCT
	Secondary	forward	acactctttccctacacgacgctctccgatctGTGATGAAGTCAGACAAATCGC
	Secondary	reverse	gtgactggagttcagacgtgtgctctccgatctATGTCAAGAATCTCAAGTGTCTG
NTD	Primary	forward	GTGGTGTTTATTACCCGTGACAAAG
	Primary	reverse	GCTGTCCAACCTGAAGAAGA
	Secondary	forward	acactctttccctacacgacgctctccgatctCATTCAACTCAGGACTTGTCTT
	Secondary	reverse	gtgactggagttcagacgtgtgctctccgatctCCAATGGTTCTAAAGCCGAAA
S1S2	Primary	forward	GCCGGTAGCACACCTTGTA
	Primary	reverse	TGTGCAAAAACCTCTGGGTGT
	Secondary	forward	cactctttccctacacgacgctctccgatctCAGGCACAGGTGTTCTTACT
	Secondary	reverse	gtgactggagttcagacgtgtgctctccgatctGTCTTGGTCATAGACTGGTAG

110

111 **Table 1. PCR primers used to amplify Spike regions for MiSeq sequencing.** Upper case indicate SARS-CoV-2  
112 sequence. Lower case indicates adapter sequence.

### 113 3. Results

#### 114 3.1. Computational Pre-processing

115 Figure 1 illustrates the steps of our workflow. The two steps of our process after read trimming  
116 used the VSEARCH tool [16]. First, the trimmed paired reads were merged using vsearch –  
117 fastq\_merge with default parameters. Then merged reads were dereplicated using vsearch --  
118 derep\_fulllength with the arguments --minsize 100 and --sizeout. These arguments limit the output  
119 to unique sequences that occur at least 100 times and appends the sequence IDs with 'size=#',  
120 where # is the number of times that sequence occurred in the reads. The cutoff of 100 counts  
121 removes late stage PCR errors, leaving only sequences representing the original templates or  
122 errors that occurred in early cycles of the PCR. This removal makes further analysis simpler and  
123 faster. However, very low frequency original template sequences will also be removed by this cut  
124 off, so this step could be skipped to preserve such rare sequences. The resulting unique  
125 sequences were mapped to the sequence of SARS-CoV-2 (NCBI Reference Sequence:  
126 NC\_045512.2, [https://www.ncbi.nlm.nih.gov/nucleotide/NC\\_045512](https://www.ncbi.nlm.nih.gov/nucleotide/NC_045512)) spike ORF using Bowtie2 [17]  
127 with default parameters to generate standard SAM formatted files. Having SAM formatted files  
128 allows the use of the program we developed for amplicon sequencing results. All files associated  
129 with these steps for our analysis of the Fenton, MO sewershed in this manuscript can be accessed  
130 at [https://github.com/degregory/SR\\_manuscript/tree/master/Fenton\\_Data](https://github.com/degregory/SR_manuscript/tree/master/Fenton_Data).

#### 131 3.2. SAM Refiner: SAM Processing

132 Our program, SAM Refiner, is currently a command line based python script and is available at  
133 [https://github.com/degregory/SAM\\_Refiner](https://github.com/degregory/SAM_Refiner) along with updated documentation. In order to run  
134 SAM Refiner, a python compiler or interpreter is needed  
135 (<https://docs.python.org/3/tutorial/interpreter.html>). Though only tested in a Linux environment, it  
136 should function with other common OSes. Figure 2 shows the command line usage for SAM  
137 Refiner. Standard SAM formatted files are the starting point for our program. These files are  
138 generated by many mapping programs, such as Bowtie2 [17] or BWA [18]. The default functions  
139 of SAM Refiner follow. Files with the extension .sam (case insensitive) in the working directory will  
140 be identified and processed. To process SAM files, SAM Refiner must be provided a FASTA  
141 formatted file for a reference sequence using the command line argument '-r reference.fasta',  
142 where the FASTA file contains the same sequence ID and sequence used to map the sequencing

143 reads in the SAM formatted file. If the IDs of the given reference and the reference of mapped  
144 sequences in the SAM file do not match, those sequences will be ignored. If the SAM formatted  
145 files were generated from dereplicated or collapsed sequences that still contain the unique read  
146 count, SAM Refiner can process the counts from certain formats. SAM Refiner will recognize the  
147 counts in sequence ids where the count is at the end of the id and denoted with a '=' or '-', i.e.  
148 'Seq1:1;counts=20' will be recognized as a sequence with 20 occurrences.

149

```
$ python SAM_Refiner.py -h
usage: SAM_Refiner.py [-h] [-r REF] [-S [SAM_FILES [SAM_FILES ...]]] [--use_count {0,1}]
  [--min_abundance1 MIN_ABUNDANCE1] [--min_abundance2 MIN_ABUNDANCE2] [--ntabund NTABUND]
  [--max_dist MAX_DIST] [--max_covar MAX_COVAR] [--Aareport {0,1}] [--AAcodonsMNP {0,1}]
  [--chim_in_abund CHIM_IN_ABUND] [--alpha ALPHA] [--foldab FOLDAB] [--redist {0,1}]
  [--max_cycles MAX_CYCLES] [--beta BETA] [--autopass AUTOPASS] [--colld COLLD] [--collect {0,1}]
  [--read {0,1}] [--nt_call {0,1}] [--ntvar {0,1}] [--indel {0,1}] [--seq {0,1}] [--covar {0,1}]
  [--pass_out {0,1}] [--chim_rm {0,1}] [--deconv {0,1}] [--wgs {0,1}]

process Sam files for variant information

optional arguments:
  -h, --help            show this help message and exit
  -r REF, --reference REF
                        reference fasta
  -S [SAM_FILES [SAM_FILES ...]], --Sam_files [SAM_FILES [SAM_FILES ...]]
                        optional .sam files, can use multiple files i.e. "-S Sample1.sam -S Sample2.sam" or "-S Sample1.sam Sample2.sam"
  --use_count {0,1}    Enable/Disable (1/0) use of counts in sequence IDs, default enabled (--use_count 1)
  --min_abundance1 MIN_ABUNDANCE1
                        Minimum observations required to be included in sample reports; >= 1 occurrence count; < 1 % observed (.1 = 10%), (default: .001)
  --min_abundance2 MIN_ABUNDANCE2
                        Minimum abundance required for variants to be included in collection reports; must be non-negative and < 1, % observed (.1 = 10%), (default: .01)
  --ntabund NTABUND    Minimum abundance relative to total reads required for a position to be reported in the nt call output; must be non-negative and < 1, % observed (.1 = 10%),
  (default: .001)
  --max_dist MAX_DIST  Maximum number of variances from the reference a sequence can have to be consider in covars processing (default: 40)
  --max_covar MAX_COVAR
                        Maximum number of variances from the reference to be reported in covars (default: 8)
  --Aareport {0,1}    Enable/Disable (1/0) amino acid reporting, default enabled (--Aareport 1)
  --AAcodonsMNP {0,1} Enable/Disable (1/0) reporting multiple nt changes in a single codon as one polymorphism, default enabled (--AAcodonsMNP 1), requires Aareport enabled
  --chim_in_abund CHIM_IN_ABUND
                        Minimum abundance a unique sequence must have to be considered in chimera removal / deconvolution (default: .001)
  --alpha ALPHA       Modifier for chim_rm chimera checking, default 1.2. Higher = more sensitive, more false chimeras removed; lower = less sensitive, fewer chimeras removed
  --foldab FOLDAB     Threshold for potential parent / chimera abundance ratio for chim_rm; default is 1.8
  --redist {0,1}      Enable/Disable (1/0) redistribution of chimera counts for chim_rm, default enabled (--redist 1)
  --max_cycles MAX_CYCLES
                        Max number of times chimera removal will be performed for chim_rm; default is 100
  --beta BETA         Modifier for covar pass checking, default 1. Higher = more sensitive, more failed checks; lower = less sensitive, fewer failed checks
  --autopass AUTOPASS threshold for a sequence to automatically pass the covar pass checking
  --colld COLLD       ID to prepend collections
  --collect {0,1}    Enable/Disable (1/0) collection step, default enabled (--collect 1)
  --nt_call {0,1}    Enable/Disable (1/0) nt_call output, default enabled (--nt_call 1)
  --indel {0,1}     Enable/Disable (1/0) indel output, default enabled (--indel 1)
  --seq {0,1}       Enable/Disable (1/0) unique seq output, default enabled (--seq 1)
  --covar {0,1}     Enable/Disable (1/0) covar output, default enabled (--covar 1)
  --pass_out {0,1}  Enable/Disable (1/0) covar_pass output, default disabled (--pass_out 0)
  --chim_rm {0,1}   Enable/Disable (1/0) chim_rm output, default enabled (--chim_rm 1)
  --deconv {0,1}    Enable/Disable (1/0) covar deconv, default enabled (--deconv 1)
```

150

151 **Figure 2. Command Line Usage of SAM Refiner.** The standard help output from SAM Refiner is shown. Syntax for  
152 the command line usage is seen, followed by details about potential arguments to modify program parameters.

153 For each SAM file, SAM Refiner initially outputs 4 tab separated values (TSV) files that can be  
154 read by any standard spreadsheet software. For a SAM file with the name Sample.sam, the  
155 outputs are named Sample\_unique\_seqs.tsv, Sample\_nt\_calls.tsv, Sample\_indels.tsv and  
156 Sample\_covars.tsv. Example outputs of each are provided in Supplemental Files 1, 2, 3, and 4,  
157 respectively ([https://github.com/degregory/SR\\_manuscript/tree/master/Supplementals](https://github.com/degregory/SR_manuscript/tree/master/Supplementals)). All reports  
158 are based on the FASTA reference relative to the SAM formatted file, so any errors made by the  
159 mapping or incongruence between the FASTA reference and the mapping reference will result in  
160 propagated errors. The reports also include the coded amino acids and their position in the coded  
161 peptide as if the reference is an in-frame coding sequence. If multiple nucleotides in a single  
162 codon differ from the reference, they will be reported together as a MNP with the associated amino  
163 acid change. Within the files, all of the sample specific outputs start with the name of the sample  
164 taken from the SAM file name followed in parenthesis by the count of reads mapped.

165

166 The Sample\_unique\_seqs.tsv file (Sup. 1) lists the unique sequence reads mapped in the SAM file  
167 using a variance notation to list the variations from the reference along with occurrence count and  
168 abundance. For example, using the previously mentioned SARS-CoV-2 spike ORF as the  
169 reference sequence, a sequence read that matches the reference except for having a T at position



170 1501 instead of the reference A would be reported simply as '1501A(N501Y)'. The abundance  
171 reported uses a decimal notation, so 0.2 represents 20%. Unique sequences that have an  
172 abundance below 0.001 are not reported.

173

174 The Sample\_nt\_calls.tsv file (Sup. 2) has a line for each nt position covered in at least 0.1% of the  
175 reads. Based on the reference sequence, each line first reports the nt position, the reference nt,  
176 the amino acid position, and the reference amino acid residue. The line then reports the number  
177 of calls for each base and for deletions at that position, followed by the most abundant (primary)  
178 call and its counts and abundance. If that primary nt is different from the reference, the amino  
179 acids encoded by the primary nt sequence and by the reference sequence with only that nt  
180 changed are reported. Further, if the second (secondary) and third (tertiary) most abundance nts  
181 are above .1% of the total read counts, those nts, their counts and abundances, and their  
182 associated amino acid changes are also reported.

183

184 The Sample\_indels.tsv (Sup. 3) file lists each insertion or deletion found in the mapping along with  
185 its occurrence count and abundance. Reported insertions have the format of 'position-  
186 insertNT(s)', so an insertion between nt positions 54 and 55 of the sequence 'GCA' will be  
187 reported as '55-insertGCA'. Reported deletions have the format 'start Position-end positionDel',  
188 so a deletion of the nts at positions 61 through 64 would be reported as '61-64Del'. Amino acid  
189 changes are reported if the indel maintains the reading frame. If there are no indels in the reads,  
190 no indel report will be generated.

191

192 Finally, the Sample\_covars.tsv (Sup. 4) file lists all observed single variances and variance  
193 combinations relative to the reference sequence. The number and abundance of sequence reads  
194 containing each covariant (covar) are reported regardless of whether any of those reads have  
195 other variations or not. As an example of this processing, the sequence '1212G(G404G)  
196 1501T(N501Y) 1709A(A570D)' with 100 counts would have the covariants of '1212G(G404G)',  
197 '1501T(N501Y)', '1709A(A570D)', '1212G(G404G) 1501T(N501Y)', '1212G(G404G)  
198 1709A(A570D)', '1501T(N501Y) 1709A(A570D)' and '1212G(G404G) 1501T(N501Y)  
199 1709A(A570D)', and contribute 100 counts to each. Because unique sequences that fall below  
200 the .1% reporting cutoff can still contribute to covariants, there may be variances in the reported  
201 covariants that aren't seen in the unique sequence output. Any sequence with more than 40  
202 variances from the reference are ignored. While all sequences with 40 or fewer variances are  
203 analyzed, only combinations of 8 or less variances are reported.

204

205 Once the above outputs are generated from each SAM file found, SAM Refiner will collect  
206 information from each sample and report them in a single file for the covars and unique\_seqs  
207 reports (Collected\_Covariances.tsv and Collected\_Unique\_Seqs.tsv). These collections have a  
208 threshold of 1% occurrence for reporting.

209

210 Many options are available as command line arguments to change parameters of the SAM  
211 processing of SAM Refiner (Fig. 2). There are no strictly required command line arguments,  
212 though the -r argument is required for the SAM processing. Omitting the reference will cause SAM  
213 Refiner to skip the SAM processing and only perform the collections and chimera removal (see  
214 below), which require per-existing outputs. The other input option is the '-S' argument which  
215 provides SAM Refiner with SAM files to process instead of allowing it to search the working  
216 directory. The use of dereplicated/collapsed counts in the SAM files can be disabled with '--  
217 use\_counts 0'. There are also options for the outputs. All outputs can be separately suppressed  
218 with the arguments '--seq 0', '--nt\_call 0', '--indel 0', '--covar 0' and '--collect 0'. The collections file  
219 names can be prepended with a string specified by the argument '--colID'. To change the  
220 reporting threshold for the sample and collected outputs, arguments '--min\_abundance1' and '--  
221 min\_abundance2' are used respectively. For '--min\_abundance1', despite its name, the value can  
222 be used to either set a minimal abundance threshold or a minimal count threshold. Values of 1 or  
223 greater will set a count threshold, while those less than 1 will set an abundance threshold. Only  
224 an abundance threshold is available for '--min\_abundance2'. All amino acid information in the  
225 reports can be suppressed with the argument '--AAreport 0'. This disabling is recommended if the  
226 reference doesn't primarily provide an in frame coding sequence. Users can also have all nt  
227 changes processed independently, even if they are in the same codon with --AAcodonasMNP 0'.  
228 Using '--ntabund' will change the required mapped coverage threshold for reporting a position in  
229 the nt\_calls output. Finally, '--max\_dist' and '--max\_covar' allow changes to the covar processing  
230 and reporting. Sequences with more variations than the amount specified by '--max\_dist' are not  
231 included in the covar analysis. The maximum number of variances reported in a combination can  
232 be set with '--max\_covar'. As an example, if '--max\_covar 2' were used for Sup. 4, then '1216-  
233 1216Del 1501T(N501Y) 1709A(A570D)', '1212G(G404G) 1501T(N501Y) 1709A(A570D)' and  
234 '1217-1217Del 1501T(N501Y) 1709A(A570D)' would not be reported.

235

236 Using the SAM files generated from the sequencing data of the Fenton sewershed, we ran SAM  
237 Refiner with the SARS-CoV-2 (NCBI Reference Sequence: NC\_045512.2) spike ORF sequence  
238 as a reference, the same as was used with the Bowtie2 mapping. The resulting outputs can be  
239 accessed at [https://github.com/degregory/SR\\_manuscript/tree/master/Fenton\\_Data](https://github.com/degregory/SR_manuscript/tree/master/Fenton_Data). These  
240 outputs allow us to see the variant lineages present at different dates in this sewer shed.  
241 However, as can be seen in Sup. 1, many of the sequences reported appear to be chimeric  
242 sequences arising from template jumping. While these outputs can still be used for further  
243 analysis, removing chimeric sequences makes such analysis easier. SAM Refiner also has  
244 methods to remove such chimeric sequences.

### 245 **3.3. SAM Refiner: Chimera Removal**

246 PCR amplification can introduce sequence errors that obscure the original template sequences.  
247 Of most concern are the introduction of false SNPs and chimeric reads. Most PCR introduced  
248 SNPs can be removed from analysis by the use of an abundance threshold such as is the default  
249 for SAM Refiner, or as was used in our pre-processing dereplication step. There are also  
250 numerous programs that can be used to attempt to remove such errors. Chimeric sequences are  
251 generally more difficult to remove. Many programs exist for this task; however, we were unable to  
252 find any that provided satisfying results for our amplicon sequencing. We developed two

253 algorithms for SAM Refiner in order to remove chimeric errors arising from PCR template jumping  
 254 from the SAM processing outputs. They are redundant in their function but use different methods,  
 255 allowing for increased confidence in results by crosschecking between the two methods.

256

Variant Sequence	Counts	Abundance
1450A(E484K) 1709A(A570D)	1478	0.006
Potential Chimera Parent Pairs		
Left Parent : Abundance	Right parent : Abundance	Multiplied Abundance
1450A(E484K) : 0.07	1501T(N501Y) / 1709A(A570D) : 0.486	0.03402 (73%)
1450A(E484K) : 0.07	1709A(A570D) : 0.097	0.00679 (14%)
1450A(E484K) : 0.07	1450A(E484K) 1501T(N501Y) / 1709A(A570D) : 0.033	0.00231 (5%)
1450A(E484K) / 1501T(N501Y) 1709A(A570D) : 0.033	1709A(A570D) : 0.097	0.003201 (7%)
1450A(E484K) / 1501T(N501Y) : 0.006	1709A(A570D) : 0.097	0.000582 (1%)
	Total:	0.046903 (100%)
Query (actual) Abundance	< Multiplied Parent (expected) Abundance	
0.006	< 0.046903 × 1.2	
1450A(E484K) 1709A(A570D) flagged as chimera, counts redistributed		

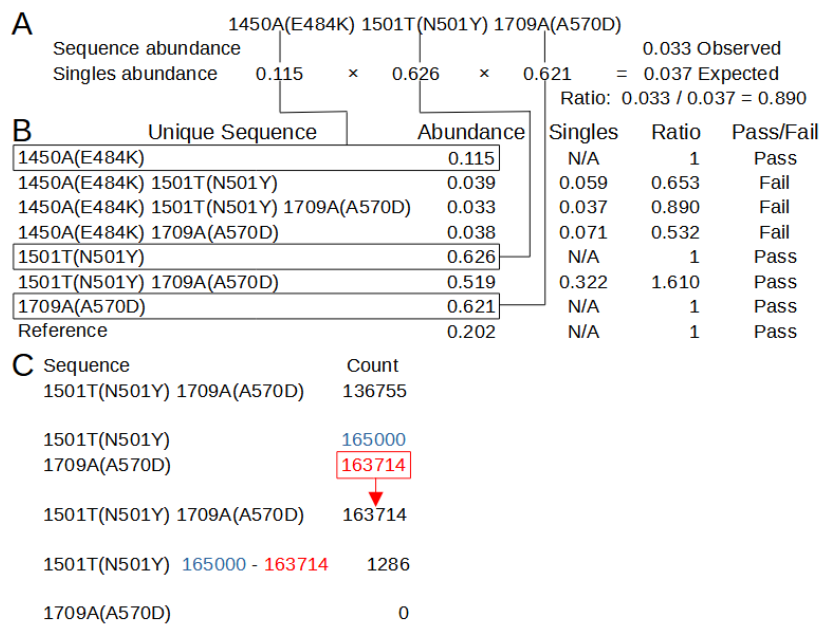
257

258 **Figure 3. First method of detection and removal of chimeras, Chimera Removed.** Using the sequences shown in  
 259 Sup. 5, the query of the least abundant sequence is shown. Potential parents whose recombination could result in the  
 260 query sequence are found. The abundances of each potential pair are multiplied. The sum the multiplied abundances  
 261 of the pairs (expected) is then compared to the abundance of the query sequence (actual) to determine if the query  
 262 sequence is a chimera. If the actual abundance is greater or equal to 1.2 times the expected abundance, the  
 263 sequence is considered non-chimeric.

264 The algorithms to remove chimeric sequences rely on the unique sequence and covariant files  
 265 generated by the SAM processing. The first algorithm, chimera removed (chim rm), goes through  
 266 the individual unique sequences, starting with the lowest abundance, and determines if the  
 267 sequences are chimeric. Figure 3 shows a simplified example of how the determination is made  
 268 on the lowest abundant sequence of an example unique sequence output (Sup. 5). The more  
 269 detailed and exact method is as follows. The sequence being considered as a potential chimera is  
 270 broken up into all possible dimeric halves. Each pair is then compared to all the other sequences  
 271 to detect potential parents. A sequence is flagged as a potential parent if its abundance is greater  
 272 than or equal to the abundance of the potential chimera multiplied by 1.8 (foldab) and there is at  
 273 least one other sequence that would be a matched parent to the complimentary dimeric half.  
 274 When a pair of dimeric halves have potential parents, the abundances of parent pairs are  
 275 multiplied. The products from each potential parent pairings are summed as an expected  
 276 abundance value and compared to the observed abundance of the potential chimera. If the  
 277 potential chimeras abundance is less than that of the expected value multiplied by 1.2 (alpha), that  
 278 sequence is flagged as a chimera and removed. The counts attributed to that flagged chimeric  
 279 sequence are then redistributed to the parents sequences based on the relative expected  
 280 contribution to recombination. Once this process has been done for all the sequences, it is  
 281 repeated until no more sequences are flagged as chimeric or 100 chimera removal cycles have  
 282 completed. The results of this algorithm that have a recalculated abundance of 0.001 or greater  
 283 are output in a new file (Sup. 6 Example\_a1.2f1.8rd1\_chim\_rm.tsv). The added string represent  
 284 values of the parameters used for the processing (alpha, foldab and redist, see below for more  
 285 information on the parameters).

286





287

288 **Figure 4. Second chimera removal method in SAM Refiner, Covariant Deconvolution.** A. Calculations of the  
 289 singles / expected abundance and abundance ratio for one of the unique sequences from Sup 5 and the abundances  
 290 from Sup 7. Lines connect the singles and their abundance to the same in B. B. Calculations for determining if a  
 291 unique sequence passes the initial check. Sequences pass when they have an Abundance/Singles ratio of 1 or  
 292 greater. C. Passed sequences are processed in order of greatest ratio to least. Counts of the sequence are set to the  
 293 counts of the least abundant single variant, and that count is then removed from all single variants in that sequence.

294 The second algorithm, covariant deconvolution (covar deconv), is a two-step process. Figure 4  
 295 shows these processes using example outputs in Sup. 5 and 7. The first step determines if a  
 296 sequence is likely to be a true or chimeric sequence by obtaining the ratio of the frequency of a  
 297 given covariant sequence relative to an expected abundance of that covariant sequence assuming  
 298 random recombination of its individual polymorphisms. The expected abundance is obtained by  
 299 multiplying the abundances of each individual variance that is present in that covariant sequence.  
 300 For instance, in a sample where '1501T(N501Y)' has an abundance of 0.32 and '1709A(A570D)'  
 301 has an abundance of 0.35, the expected abundance of the covariant '1501T(N501Y)  
 302 1709A(A570D)' would be 0.112 [0.32 × 0.35]. If the ratio of the observed abundance to the  
 303 expected abundance is equal to or greater than 1 (beta), that covariant passes the check and is  
 304 sent to the second step. Any sequence that has an abundance of 0.3 or greater is automatically  
 305 passed. If such a sequence has an observed/expected ratio less than 1, it will be assigned a ratio  
 306 of 1. The second step processes the passed sequences in order of greatest observed/expected  
 307 ratio to least. If multiple sequences have the same ratio, they are processed in order of greatest to  
 308 least distance from the reference. Sequences that automatically pass the first step are processed  
 309 after the other sequences and in order of least abundant to greatest. Sequences are assigned a  
 310 new occurrence count based on their constituent individual variances. For the sequence being  
 311 processed, the count for the least abundant individual variance is assigned to the sequence and  
 312 constituent variances making up the sequence have their count reduced by the amount of the  
 313 least variance. This reduction means the individual variance that had the least counts is assigned  
 314 0 counts, so any sequence not yet processed in which that variance is present is functionally  
 315 removed. This process is repeated until all sequences have been reassessed or removed. The  
 316 final results with an abundance of 0.001 or greater are reported in a new file (Sup. 8  
 317 Example\_covar\_deconv.tsv).

318

319 As before, the results from individual samples are collected and reported for entries above 1%  
320 occurrence. A number of command line arguments will also influence the chimera removal  
321 algorithms. Both chimera removal algorithms run by default, but either or both can be disabled (--  
322 chim\_rm 0 and --covar\_deconv 0). The collections are again disabled with '--collect 0'. An  
323 additional output of the covariants that passed the first step of the second algorithm can be  
324 generated with '--pass\_out 1'(Sup. 9). The outputs are constrained as before by a minimum  
325 abundance with command line arguments '--min\_abundance1' and '--min\_abundance2'.  
326 Collection file names are also prepended with '--collID'. The only input parameter that can be  
327 changed by command line argument is the abundance of sequences or covariants that will be  
328 considered in the algorithms. By default, only entries from the inputs that have a 0.001 abundance  
329 or greater are processed. This threshold can be changed with '--chim\_in\_abund'.

330

331 Four parameters can be altered for the first algorithm. The abundance ratio that is used as a  
332 threshold for selecting potential parents of potential chimera can be set with '--foldab'. Larger  
333 values will generally reduce the pool of sequences that will be considered as potential parents,  
334 thus potentially reducing the total expected abundance obtained from parent pairs and number of  
335 sequences flagged as chimeric. In the most simple theoretical model of PCR chimera generation,  
336 two parents generate one chimera. The parents have at least twice the abundance of the chimera  
337 as they would exist and have been amplified prior to the chimera. The reality of chimera  
338 generation can be much more complex, as many sequences may generate identical chimeras  
339 multiple times. If a sample has little chimera generation, a --foldab value close to 2, such as the  
340 default of 1.8, should be sufficient to remove chimeras without also removing non-chimeric  
341 sequences in error. However, the more chimera generation observed, the more the --foldab value  
342 needs to be reduced to accurately remove all chimeric sequences, even to 0 to not exclude any  
343 sequence from being considered a potential parent (though it will likely be vary rare for such a  
344 value to be necessary). Lower values, however, will also increase the likelihood of a sequence  
345 being flagged as a chimera in error. Users may need to empirically determine the best value for  
346 their samples.

347

348 The multiplier for the parental summed abundance for determining if a sequence is a chimera can  
349 be set with '--alpha'. Larger values will generally result in a greater number of sequences flagged  
350 as chimeric. As with --foldab, the optimal value for --alpha will depend on the extent of chimera  
351 generation in the samples being processed, with a value near 1 for minimal chimera generation  
352 (such as the default 1.2) and 2 or even higher for rampant chimera generation. Once again, the  
353 later would also increase the likelihood of sequences being flagged as chimeric in error.

354

355 Redistribution of the counts from the chimera to the parent sequences can be disabled with '--  
356 redist 0'. Redistribution is meant to give an estimate of the counts and abundances that would  
357 have been observed without chimera generation, which users may wish to forgo. The maximum  
358 number of chimera removal cycles can be change by '--max\_cycles', ei '--max\_cycles 2' will only  
359 allow two iterations of the chimera removal. Multiple removal cycles allows chimeras to be found

360 based on new counts and abundances resulting from previous cycles, increasing the likelihood  
361 chimeras are removed from a sample.

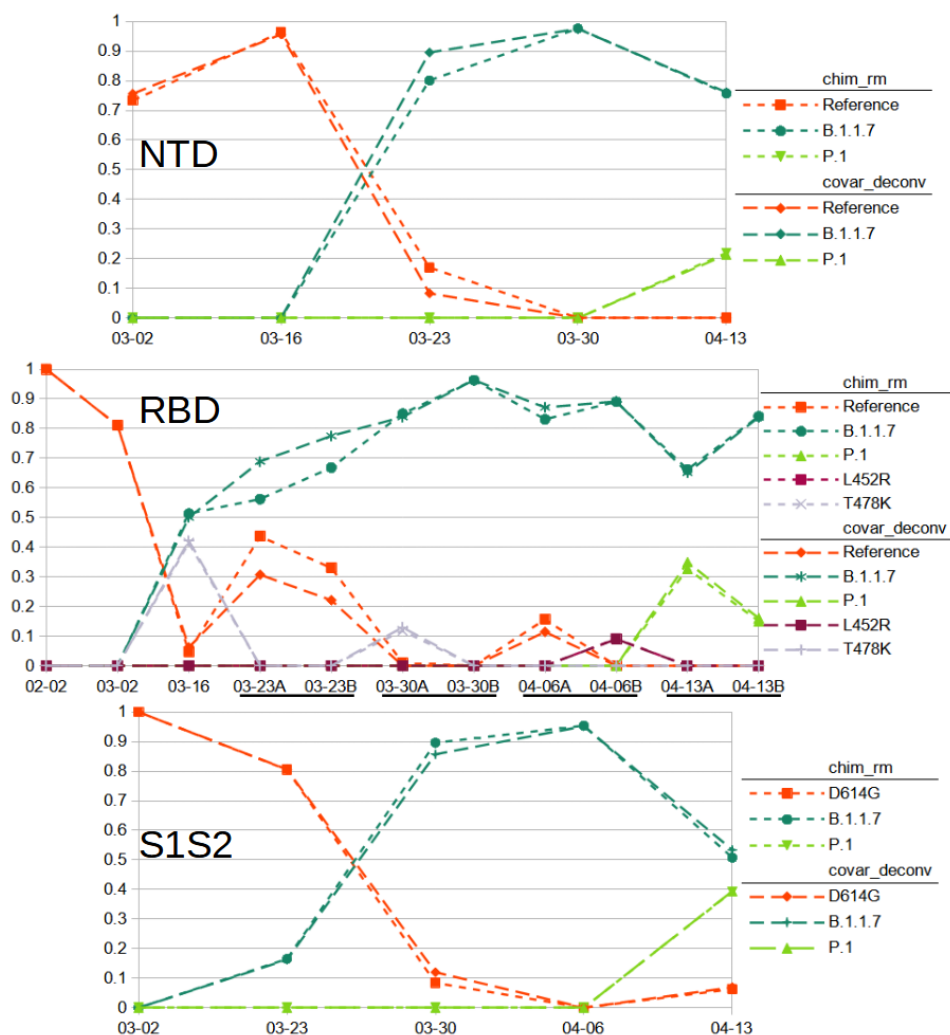
362

363 The second algorithm has two parameter that can be changed. The ratio threshold at which a  
364 covariant will be passed to the second step can be altered with '--beta'. The abundance at which  
365 a covariant will automatically be passed can be changed with '--autopass'.

366

367 The chimera removal methods of SAM Refiner were also used on the Fenton sewershed  
368 sequencing data. Due to the relatively high amount of chimeric sequences in our samples, we  
369 used the command line arguments '--foldab=0.6 --alpha=2.2'. The outputs generated for the  
370 Fenton sewershed from 2-2-21 to 4-13-21 can be accessed at  
371 [https://github.com/degregory/SR\\_manuscript/tree/master/Fenton\\_Data](https://github.com/degregory/SR_manuscript/tree/master/Fenton_Data). The two different chimera  
372 removal methods showed good concordance, validating each as being a viable method. Duplicate  
373 RT-PCR preparation and sequencing of the same wastewater sample also generally provided  
374 similar results, though less consistently (Fig. 5. Compare A and B RBD amplicon preparations).  
375 These differences were more pronounced with covariants with relatively low abundance, such as  
376 is seen with 3-30 RBD samples, where one detects T478K and the other does not (Fig. 5). These  
377 differences illustrate the stochastic nature of RT-PCR amplification.

378



379

380 **Figure 5.** Relative Abundance of Reference and Variant SARS-CoV-2 Sequences Observed in Fenton, MO  
 381 sewershed from February to March. Results from sequencing of Spike amplicons of the NTD, RBD and S1S2 junction  
 382 regions are shown. Lines of short dashes connect values obtained by the chimera removed method, lines of long  
 383 dashes connect values obtained by the covariant deconvolution method. All amplicons show a population shift from  
 384 the reference with D614G to B.1.1.7 sequences with the appearance of P.1 sequences at the last time point.  
 385 Additionally, known common polymorphisms T478K and L452R were observed from the RBD amplicons. RT-PCR for  
 386 the RBD amplicon was performed in duplicate for some samples.

387 We used the chimera removed and covariant deconvolution outputs to assign sequences to known  
 388 variant lineages or the reference (Sup. 10, 11 & 12) based on variances present. Variances that  
 389 only appeared in one sequencing run and did not appear frequently in GSIAD  
 390 (<https://www.gisaid.org/>) were considered likely PCR error and not taken into account for  
 391 sequence assignment. Based on the assignments, we were able to observe the changes to virus  
 392 populations in the sewershed over time (Fig. 5). We classified the sequences found from the NTD  
 393 amplicon as matching reference sequence, lineage B.1.1.7 with '203-208Del 429-431Del' or  
 394 lineage P.1 with '412T(D138Y) 570T(R190S)' (Sup 10). Sequences from the RBD amplicon  
 395 matched reference sequence, lineages B.1.1.7 with '1501T(N501Y) 1709A(A570D)' or P.1 with  
 396 '1250C(K417T) 1450A(E484K) 1501T(N501Y)', or had the single variations of T478K or L452R  
 397 (Sup 11). T478K and L452R each have lineage associations. However, no other variances are  
 398 associated with these in the RBD amplicons, nor were any variances present in the other  
 399 amplicons that would indicate the presence of any associated lineages. While these SNPs could  
 400 be the result of PCR error, it is more likely the associated lineages exist in the sewershed, but due

401 to stochastic effects the other associated variances in the other amplicons were not detected.  
402 They could have also arisen in the reference background. As we can not assign them to known  
403 variant lineages with any certainty, we assigned them to their own category. Sequences from the  
404 S1S2 amplicon matched lineage B.1.1.7 with '1841G(D614G) 2042A(P681H) 2147T(T716I)',  
405 lineage P.1 with '1841G(D614G) 1963T(H655Y) 2063T(A688V)' or only had the now ubiquitous  
406 D614G variation (Sup 12). The 03-23 S1S2 sample had a sequence '1841G(D614G)  
407 2037G(N679K) 2063T(A688V)'. While A688V is associated with P.1, it does not appear in that  
408 context here. As that is the only sample where those covariant sequences were observed and the  
409 variances are not frequently reported in GISAID (outside of P.1 for A688V), we assigned it to the  
410 reference category. Looking at all samples over time and the three amplicon regions in concert,  
411 we can conclude that the SARS-CoV-2 population of this sewershed changed from almost  
412 exclusively having only the D614G variation to mainly the B.1.1.7 lineage, with the introduction of  
413 P.1 early in April. This general method is now being used to track SARS-CoV-2 variants in many  
414 Missouri sewersheds (<https://storymaps.arcgis.com/stories/f7f5492486114da6b5d6fdc07f81aacf>).

### 415 **3.4. SAM Refiner: Limitations and Future Development**

416 While the outputs of SAM Refiner can be very informative, the program has some limitations,  
417 some of which may be overcome in future development. Currently the greatest limitation is the  
418 need for users to be familiar with command line usage. We hope to develop a graphical user  
419 interface version of these programs to overcome this user hurdle in the future. We also intend to  
420 develop SAM Refiner to be available from widely used functional collections, such as BioConda  
421 (<https://bioconda.github.io/>) and Galaxy (<https://usegalaxy.org/>).

422

423 Though SAM Refiner can be used on sequencing not based on amplicons, its usefulness will be  
424 more limited, as the relative abundance of sequences and covariants will be calculated based on  
425 total reads and not positional coverage. Future development may include modes for whole  
426 genome sequencing or multiple amplicons, even the ability to use multiple sequences for a  
427 reference.

428

429 The accuracies of the chimera removal algorithms will vary greatly depending on the parameters  
430 used and the sample they are being run on. Due to the stochastic nature of chimera generation  
431 and amplification during PCR and the possible complexity of the original template sequences,  
432 samples will sometimes be refractory to chimera removal algorithms. This problem is faced by all  
433 programs designed for this purpose. The ability to modify parameters in the algorithms and having  
434 two algorithms with different approaches to the chimera removal improves the accuracy the user  
435 can achieve with this software. Some samples will, however, always fail to be processed  
436 accurately by one or both methods.

437

438 **Supplementary Materials:** The following are available online at  
439 [https://github.com/degregory/SR\\_manuscript/tree/master/Supplementals](https://github.com/degregory/SR_manuscript/tree/master/Supplementals), Sup. 1 Example of SAM  
440 Refiner's Output for Reporting Unique Sequences, Sup. 2 Example of SAM Refiner's Output for  
441 Reporting Positional NT Calls, Sup. 3 Example of SAM Refiner's Output for Reporting Insertions  
442 and Deletions, Sup. 4 Example of SAM Refiner's Output for Reporting Covariance, Sup. 5 Sample



443 Unique Sequences Output With Chimeric Sequences, Sup. 6 Sample Output of Sequences of  
444 SAM Refiner's Chimeras Removed, Sup. 7 Sample Covariance Output With Chimeric Sequences,  
445 Sup. 8 Sample Passed Sequences Output from the First Part of SAM Refiner's Covariant  
446 Deconvolution Method, Sup. 9 Sample Output of Sequences by SAM Refiner's Covariant  
447 Deconvolution Method, Sup. 10 Assignment of NTD Covariant Sequences to Variants and  
448 Lineages, Sup. 11 Assignment of RBD Covariant Sequences to Variants and Lineages, Sup. 12  
449 Assignment of S1S2 Covariant Sequences to Variants and Lineages

450 **Author Contributions:** Conceptualization, MJ.; methodology, MJ; software, DG; validation, DG  
451 and MJ.; formal analysis, DG and MJ; investigation, DG and MJ; resources, MJ, JW; data curation,  
452 DG, CW, CL and MJ; writing—original draft preparation, DG; writing—review and editing, DG, JW,  
453 CW and MJ; visualization, DG; supervision, MJ; project administration, CW, JW, and MJ; funding  
454 acquisition, JW and MJ. All authors have read and agreed to the published version of the  
455 manuscript

456 **Funding:** Funding for the project was administered by the Missouri Department of Health and  
457 Senior Services (DHSS). This project was supported by funding from the Centers for Disease  
458 control and the National Institutes of Health grant U01DA053893-01

459 **Data Availability Statement:** Raw and processed data can be accessed at  
460 [https://github.com/degregory/SR\\_manuscript](https://github.com/degregory/SR_manuscript)

461 **Acknowledgments:** We would like to acknowledge Christopher Bottoms for assistance in  
462 software development, and the University of Missouri DNA Core for assistance in developing deep  
463 sequencing protocols.

## 464 References

- 465 1. Ghebreyesus, Tedros Adhanom. (2020) Speech.  
466 <https://www.who.int/director-general/speeches/detail/who-director-general-s-opening-remarks-at-the-media-briefing-on-covid-19---11-march-2020>  
467
- 468 2. Zhu N., Zhang D., Wang W., Li X., Yang B., et al. (2020). A novel coronavirus from patients  
469 with pneumonia in China, 2019. *N. Engl. J. Med.* 382 727–733
- 470 3. Wu F., Zhao S., Yu B., Chen Y. M., Wang W., et al. (2020). A new coronavirus associated with  
471 human respiratory disease in China. *Nature* 579 265–269
- 472 4. Ahmed W, Angel N, Edson J, Bibby K, Bivins A, et al. (2020 ) First confirmed detection of  
473 SARS-CoV-2 in untreated wastewater in Australia: A proof of concept for the wastewater  
474 surveillance of COVID-19 in the community. *Sci Total Environ.* 2020 Aug 1;728:138764
- 475 5. Medema G., Heijnen L., Elsinga G., Italiaander R., and Brouwer A. (2020) Presence of SARS-  
476 Coronavirus-2 RNA in Sewage and Correlation with Reported COVID-19 Prevalence in the  
477 Early Stage of the Epidemic in The Netherlands. *Environmental Science & Technology Letters*  
478 2020 7 (7), 511-516
- 479 6. Nemudryi A, Nemudraia A, Wiegand T, Surya K, Buyukyoruk M, et al. (2020 )Temporal  
480 Detection and Phylogenetic Assessment of SARS-CoV-2 in Municipal Wastewater. *Cell Rep*  
481 *Med.* 2020 Sep 22;1(6):100098

- 482 7. Martin J, Klapsa D, Wilton T, Zambon M, Bentley E et al. (2020) Tracking SARS-CoV-2 in  
483 Sewage: Evidence of Changes in Virus Variant Predominance during COVID-19 Pandemic.  
484 *Viruses*. 2020 Oct 9;12(10):1144.
- 485 8. Ul-Rahman A, Shabbir MAB, Aziz MW, Yaqub S, Mehmood A, et al. (2020) A comparative  
486 phylogenomic analysis of SARS-CoV-2 strains reported from non-human mammalian species  
487 and environmental samples. *Mol Biol Rep*. 2020 Nov;47(11):9207-9217.
- 488 9. Crits-Christoph A, Kantor RS, Olm MR, Whitney ON, Al-Shayeb B, et al. (2021) Genome  
489 Sequencing of Sewage Detects Regionally Prevalent SARS-CoV-2 Variants. *mBio*. 2021 Jan  
490 19;12(1):e02703-20.
- 491 10. Izquierdo-Lara R, Elsinga G, Heijnen L, Munnink BBO, Schapendonk CME, et al. (2021)  
492 Monitoring SARS-CoV-2 Circulation and Diversity through Community Wastewater  
493 Sequencing, the Netherlands and Belgium. *Emerg Infect Dis*. 2021 May;27(5):1405-1415
- 494 11. La Rosa G, Mancini P, Bonanno Ferraro G, Veneri C, et al. (2021) Rapid screening for SARS-  
495 CoV-2 variants of concern in clinical and environmental samples using nested RT-PCR assays  
496 targeting key mutations of the spike protein. *Water Res*. 2021 Jun 1;197:117104.
- 497 12. Smyth D, Trujillo M, Cheung K, Gao A, Hoxie I, et al. (2021) Detection of Mutations Associated  
498 with Variants of Concern Via High Throughput Sequencing of SARS-CoV-2 Isolated from NYC  
499 Wastewater. *medRxiv [Preprint]*. 2021
- 500 13. Fontenele RS, Kraberger S, Hadfield J, Driver EM, Bowes D et al. (2021) High-throughput  
501 sequencing of SARS-CoV-2 in wastewater provides insights into circulating variants. *medRxiv*  
502 *[Preprint]*. 2021 Jan 25:2021.01.22.21250320.
- 503 14. Weisblum Y, Schmidt F, Zhang F, DaSilva J, Poston D, et al. Escape from neutralizing  
504 antibodies by SARS-CoV-2 spike protein variants. *Elife*. 2020 Oct 28;9:e61312. doi:  
505 10.7554/eLife.61312.
- 506 15. Martin, Marcel. (2011) Cutadapt removes adapter sequences from high-throughput  
507 sequencing reads. *EMBnet.journal*, [S.I.], v. 17, n. 1, p. pp. 10-12, may 2011.
- 508 16. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. (2016) VSEARCH: a versatile open source  
509 tool for metagenomics. *PeerJ* 4:e2584.
- 510 17. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nature Methods*. 2012  
511 Mar 4;9(4):357-9.
- 512 18. Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler  
513 Transform. *Bioinformatics*, 25:1754-60.