



32 **ABSTRACT**

33

34 Wastewater-based epidemiology has emerged as a promising tool to monitor pathogens in a  
35 population, particularly when clinical diagnostic capacities become overwhelmed. During the  
36 ongoing COVID-19 pandemic caused by Severe Acute Respiratory Syndrome Coronavirus-2  
37 (SARS-CoV-2), several jurisdictions have tracked viral concentrations in wastewater to inform  
38 public health authorities. While some studies have also sequenced SARS-CoV-2 genomes  
39 from wastewater, there have been relatively few direct comparisons between viral genetic  
40 diversity in wastewater and matched clinical samples from the same region and time period.  
41 Here we report sequencing and inference of SARS-CoV-2 mutations and variant lineages  
42 (including variants of concern) in 936 wastewater samples and thousands of matched clinical  
43 sequences collected between March 2020 and July 2021 in the cities of Montreal, Quebec  
44 City, and Laval, representing almost half the population of the Canadian province of Quebec.  
45 We benchmarked our sequencing and variant-calling methods on known viral genome  
46 sequences to establish thresholds for inferring variants in wastewater with confidence. We  
47 found that variant frequency estimates in wastewater and clinical samples are correlated over  
48 time in each city, with similar dates of first detection. Across all variant lineages, wastewater  
49 detection is more concordant with targeted outbreak sequencing than with semi-random  
50 clinical swab sampling. Most variants were first observed in clinical and outbreak data due to  
51 higher sequencing rate. However, wastewater sequencing is highly efficient, detecting more  
52 variants for a given sampling effort. This shows the potential for wastewater sequencing to  
53 provide useful public health data, especially at places or times when sufficient clinical sampling  
54 is infrequent or infeasible.

## 55 INTRODUCTION

56 Quebec has been one of the Canadian provinces most affected by the COVID-19 pandemic,  
57 with 5,698 cases and 135 deaths per 100,000 inhabitants as of December 2021 (Public Health  
58 Agency of Canada 2021). After the first reported COVID-19 cases at the end of February 2020,  
59 the government implemented public health measures to control the spread of the virus,  
60 ranging from the closure of non-essential businesses to intra-provincial travel restrictions. The  
61 Coronavirus Sequencing in Quebec (CoVSeQ) consortium began sequencing SARS-CoV-2  
62 genomes from qPCR-positive nasal swabs sampled across the province, yielding insights into  
63 the early introduction events of the virus into the province, and their subsequent spread (Murall  
64 et al. 2021). As in many other places around the world, the Quebec Public Health lab (LSPQ)  
65 increased genomic surveillance of SARS-CoV-2 after the emergence of variants of concern  
66 (VOCs) at the end of 2020. Variants of interest (VOIs) are defined as SARS-CoV-2 lineages  
67 with mutations that have a potential impact on the clinical or epidemiological characteristics of  
68 the virus, while those with a demonstrated impact on disease transmission or clinical features  
69 are named VOCs (Institut national de santé publique du Québec 2021). For example, the  
70 Alpha variant (PANGO lineage B.1.1.7) had a significant transmission advantage relative to  
71 previously circulating lineages (Volz et al. 2021; Davies et al. 2021), which was then  
72 superseded by the even more transmissible Delta (B.1.617.2) variant (Campbell et al. 2021;  
73 Mlcochova et al. 2021). Although Quebec has the highest testing rate in Canada (491 tests  
74 performed per 100,000 habitants per week vs 335 on average in Canada as of December 18,  
75 2021) and a decreasing death rate since the beginning of 2021 (from 4.02 new deaths per  
76 100,000 habitants per week at the beginning of 2021 to 0.0 in mid August 2021), the recent  
77 increase in the case incidence rate (from 1.1 new cases per 100,000 habitants in mid July  
78 2021 to 3.3 in mid August 2021) and the threat of new variants with immune escape or further  
79 transmission advantages (Otto et al. 2021) calls for continued vigilance and improved  
80 surveillance.

81  
82 Wastewater (WW) surveillance has emerged as a method to track SARS-CoV-2 variants in a  
83 manner that is complementary to sequencing clinical samples. Wastewater-based  
84 epidemiology has the potential to variant lineages earlier than clinical sampling and provides  
85 valuable insights about the viral spread in the population (Bibby et al. 2021; Xiao et al. 2021).  
86 This is because wastewater sampling can catch asymptomatic cases that would otherwise not  
87 present for a nasal swab sample, but that do excrete viral RNA in stool (Bibby et al. 2021;  
88 Jones et al. 2020). Wastewater is a complex mixture of fragmented viral RNA, which can be  
89 difficult to sequence and confidently identify mutations. Recent studies in the United States  
90 have shown that the dominant SARS-CoV-2 variants observed in clinical samples are largely  
91 mirrored in wastewater sequences (Crits-Christoph et al. 2021; Baaijens et al. 2021).

92 Wastewater also has the potential to capture mutations and variant lineages not observed in  
93 human clinical samples, potentially including animal reservoirs (Smyth et al. 2021). However,  
94 more studies are required to assess the sensitivity and spatio-temporal resolution of  
95 wastewater sequencing in comparison to clinical sequencing.

96

## 97 **RESULTS AND DISCUSSION**

### 98 **Wastewater and clinical SARS-CoV-2 sampling and sequencing in Quebec**

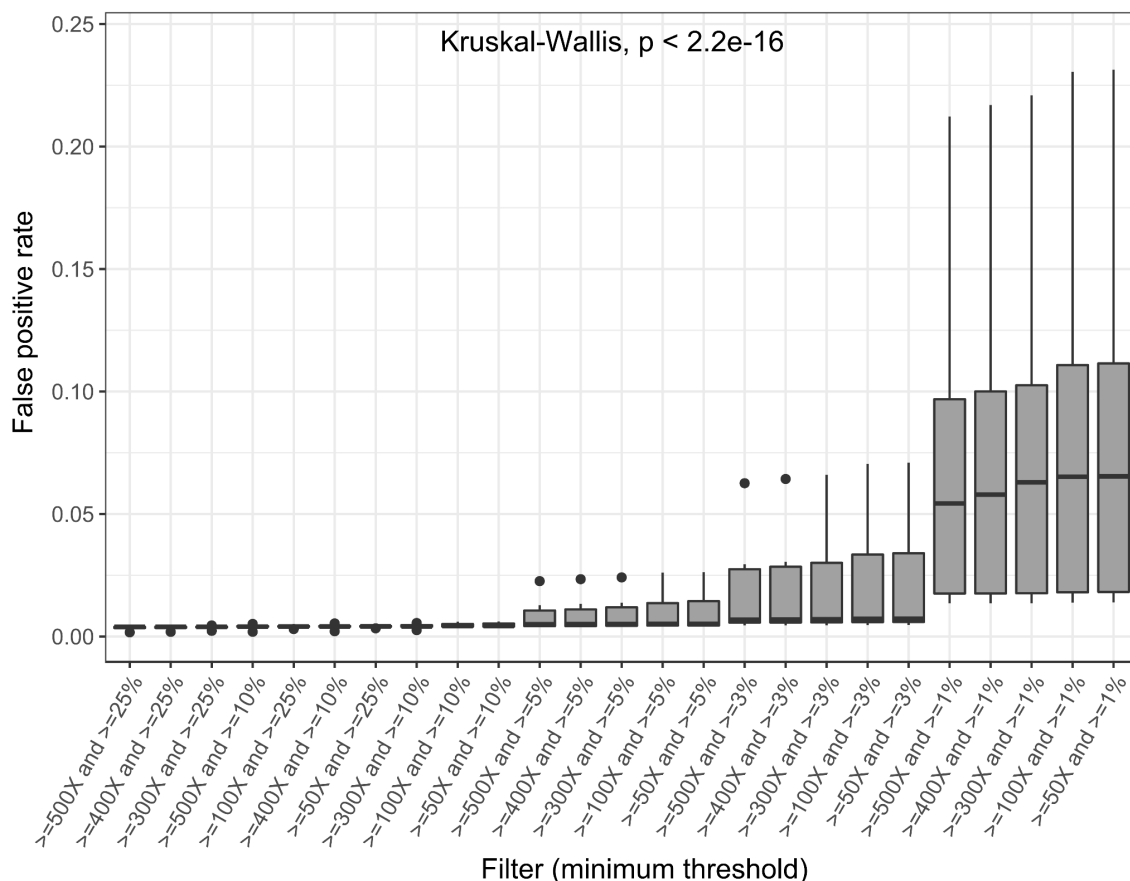
99 As a demonstration of WW-based VOCs and VOIs surveillance, we deep-sequenced SARS-  
100 CoV-2 genome-wide amplicons from 936 WW samples collected between March 2020 and  
101 July 2021 in the cities of Montreal, Quebec, and Laval, representing almost half of Quebec's  
102 population. The WW samples were collected from interceptors at the end of the WW network  
103 (36.1%), residential areas (29.9%), WW treatment facilities (21.2%), prisons (7.69%),  
104 industrial zones (4.90%) and long-term care and housing centers (0.21%). We compare the  
105 WW data to semi-random clinical samples from the same cities, which we refer to as the  
106 clinical dataset ( $N = 13,296$  sequences), and non-random priority sequencing of outbreaks  
107 and other samples of particular public health interest (e.g. suspected Alpha variant cases in  
108 early 2021, travel-related cases, etc.), which we refer to as 'outbreak' samples ( $N = 5,661$   
109 sequences).

110

### 111 **Establishing thresholds for calling single nucleotide variants (SNVs) with confidence**

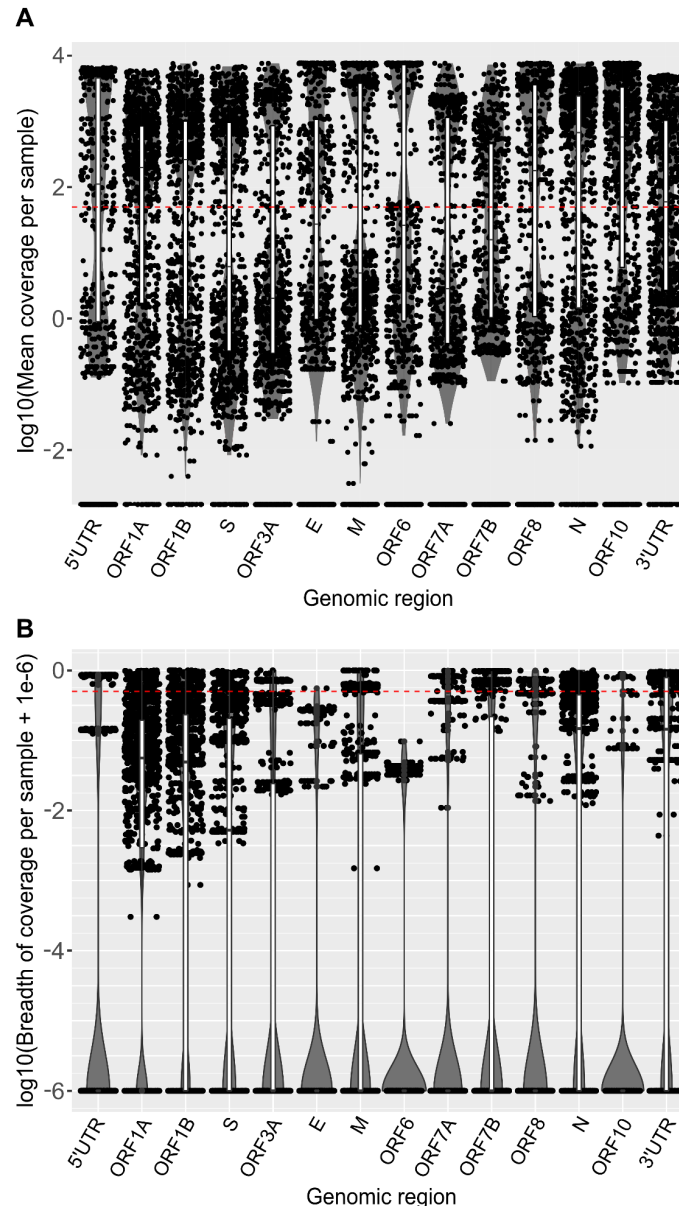
112 To study the genomic diversity of SARS-CoV-2 in WW samples, we first developed a  
113 bioinformatic pipeline to confidently identify single nucleotide variants (SNVs) while accounting  
114 for sequencing errors, then to infer VOIs, VOCs, and other variant lineages that contain  
115 combinations of signature SNVs (**Methods**). WW samples can contain a mixture of variants  
116 at different frequencies. True SNVs can be confused with errors that are introduced by sample  
117 processing, e.g. from errors introduced by the reverse transcriptase or the polymerase that  
118 are used for targeted SARS-CoV-2 amplification, or from read errors created by the  
119 sequencing platform. In this study, we chose illumina sequencing to minimize the per-read  
120 error rate (compared to Oxford Nanopore). To set appropriate thresholds for minimum SNV  
121 frequencies and coverage, we sequenced 14 SARS-CoV-2 positive control genomes from  
122 AccuGenomics, containing a set of known SNVs at 100% within-sample frequency, which we  
123 sequenced using the same Illumina amplicon strategy used for the WW samples. We  
124 assessed a range of filters for depth of coverage and minor allele frequency (MAF) required  
125 to call a SNV (**Figure 1**). Our SNV calling pipeline yields systematically low false positive rates  
126 (median  $< 7\%$  even at 1% MAF and effectively zero at 25% MAF) across a range of filters and  
127 expected frequency of the true positive SNVs (**Figure 1**). To be conservative for subsequent  
128 analyses we selected a minimum site coverage of 50X and a MAF of 25%. These filters

129 maximized the F1 score, which combines information about the precision (i.e. the proportion  
 130 of SNV calls that are true positives) and the sensitivity (i.e. the proportion of true variants that  
 131 are identified). The F1 scores are poor when the expected frequency of the true positive SNVs  
 132 is 100% (**Figure S1A**) but improves when we allow the expected nucleotide to have a  
 133 frequency of at least 75% to allow for sequencing errors (**Figure S1B**). Note that very similar  
 134 F1 scores were obtained at 10% MAF and that our variant calling pipeline can detect low-  
 135 frequency mutations with high accuracy (**Figures 1 and S2**), but we proceeded with 25% to be  
 136 conservative. In addition to the AccuGenomics controls, we also sequenced mixtures of  
 137 two different SARS-CoV-2 viral cultures at known ratios, which we call “spike-in” samples. We  
 138 found that our SNV-calling pipeline identified the expected SNVs close to their expected  
 139 frequencies. Low variant frequencies (expected frequency < 5%) were inferred particularly  
 140 accurately, with more variation in the 25-50% range (**Figure S2**). Overall, these results  
 141 indicate that our SNV-calling filters are appropriate, and likely conservative, for identifying  
 142 SNVs in wastewater containing mixtures of viral genomes.  
 143



144  
 145 **Figure 1. False-positive SNV calling rates across different depth and frequency filters.**  
 146 In all cases, the expected variant frequency was 100% in AccuGenomics SARS-CoV-2  
 147 genome standards. The thresholds that define the filters are respectively the minimum  
 148 coverage (X) and the minimum minor SNV frequency. The Kruskal-Wallis test p-value (at the  
 149 top of the panel) indicates the significance of the differences across the different sets of filters.

150 **Sampling and sequencing SARS-CoV-2 from wastewater**  
151 Having established a reliable method for SNV calling, we applied it to a set of 936 WW samples  
152 collected between March 2020 and July 2021 in Montreal, Quebec City, and Laval (**Table S1**).  
153 Because SNV calling requires sufficient depth of coverage, we measured how coverage depth  
154 and breadth varied across the SARS-CoV-2 genome in our samples (**Figure 2**).



155  
156 **Figure 2 Sequencing coverage across the SARS-CoV-2 genome recovered from**  
157 **wastewater.** A) Depth of coverage, defined as the average number of sequencing reads  
158 covering a nucleotide site. Each point represents a unique wastewater sample. Each  
159 distribution is represented by a jitter plot (black points), a boxplot (black and white) and a violin  
160 plot (grey). The red dotted line indicates a coverage of 50 on a log10 scale. B) Breadth of  
161 coverage. The breadth of coverage represents the proportion of sites in each genomic region  
162 (gene or open reading frame; ORF) with a coverage of at least 50X. Boxplots are absent when  
163 75% of the samples have a breadth of coverage of 0 or, in other words, the 3 first quantiles  
164 of the distribution correspond to samples with a coverage <50 at the corresponding sites. The  
165 red dotted line indicates a breadth of coverage of 50% on a log10 scale.

166 The mean depth and breadth of coverage varied somewhat across genes, and was roughly  
167 bimodal, with samples tending to fall either close to 0-1X coverage or >1000X on average,  
168 with fewer samples in between (**Figure 2**). As expected, coverage correlated negatively with  
169 PCR Cycle threshold (linear regression adjusted  $R^2 = 19.2\%$ ; Permutational ANOVA  $p < 2e-$   
170  $4$ ,  $n = 5000$  permutations), as previously observed (Lythgoe et al. 2021; Baaijens et al. 2021).  
171 Due to this variation in coverage, not all nucleotide sites in the genome in all samples have  
172 sufficient depth to make SNV calls. However, given the substantial number of samples with  
173 >1000X coverage, we expect to obtain a reasonable number of high-confidence SNVs in our  
174 dataset. Using the  $\geq 50X$  coverage and 25% MAF filters described above, we called an  
175 average of 6.87 SNVs per sample (s.d. = 10.1) (**Figure S3**).

176

### 177 **Detecting SARS-CoV-2 lineages in wastewater using signature and marker mutations**

178 With these SNVs in hand, we sought to infer the presence of known variant lineages in WW  
179 samples. We defined “signature mutations” as SNVs with a minimum prevalence of 90%  
180 among the consensus sequences of a certain lineage (e.g. a particular VOC represented in a  
181 database). To find these signature mutations, we first calculated the prevalence of  
182 substitutions in thousands of publicly available consensus sequences collected during 2020  
183 and added data from CoV-Spectrum (Chen et al. 2021) to include under-represented lineages  
184 or lineages that emerged after 2020 (**Methods**). These variant lineages were named using  
185 the PANGO lineage designation scheme (Rambaut et al. 2020). We further defined “marker  
186 mutations” as signature mutations that are unique to a single variant lineage to the exclusion  
187 of others.

188 Using these criteria, we identified an average of 4.42 signature mutations per sample  
189 (s.d. = 6.94) and 2.01 marker mutations per sample (s.d. = 3.22) (**Figure S3**). To call a variant  
190 lineage as present in a WW sample, we required at least 3 signature mutations including at  
191 least one marker mutation (**Methods**), yielding an average of 0.79 variant lineages identified  
192 per sample (s.d. = 1.14; min.=0; max. = 6). This approach allowed us to track the spread of  
193 variants with a confidence score, i.e. the number of signature mutations supporting the  
194 presence of a lineage, which includes at least one of its marker mutations (**Figure S4 and**  
195 **Table S2**). As expected in this time period, the samples are dominated by the Alpha variant  
196 (B.1.1.7), with sporadic identification of other variants.

197 We also estimated the within-sample frequency of variant lineages using constrained  
198 linear models (**Methods**). These models infer the linear combination of known variant lineages  
199 (from the database of reference genomes) that best explain the frequencies of observed SNVs  
200 in the sample (**Figure S6**). We applied this approach to the 299 WW samples that have at  
201 least one lineage with at least 3 signature mutations, including at least one marker mutation.  
202 Among these samples, the regression model was significant for 250 samples (83.6%) and not

203 significant for the remaining 49 samples (16.4%). These variant frequency estimates (**Figure**  
204 **S5**) were generally concordant with the simpler approach relying on signature and marker  
205 mutations (**Figure S4**). Both approaches showed a predominance of the Alpha variant in all  
206 cities, with substantial numbers of B.1.160 in 2020 and early 2021, and A.2.5.X lineages  
207 previously described in Quebec clinical samples (Murall et al., 2021).

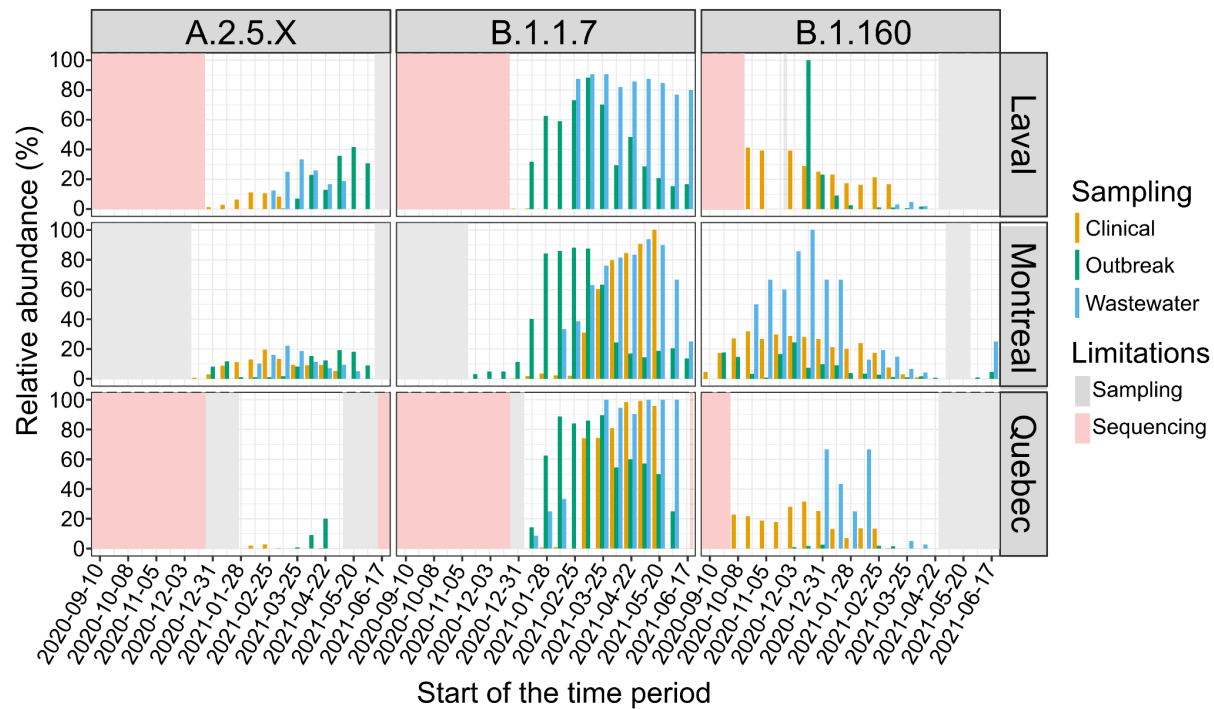
208 The absence of variant lineage detections can simply be due to lack of sampling, lack  
209 of sequencing depth, or a true absence. Sampling frequency is indeed variable over time and  
210 cities (**Figure S7**) and the absence of detections is mostly explained by the absence of WW  
211 samples (81.4% of time points with missing lineages; grey shading in **Figures S4 and S5**) and  
212 modestly by lack of sufficient sequencing depth (18.6%; transparent red shading in **Figures**  
213 **S4 and S5**). This suggests that increased sampling frequency, coupled with modest  
214 optimization of RNA extraction and sequencing protocols, could increase both the true positive  
215 and true negative rates.

216

### 217 **Comparisons between wastewater and nasal swab sequencing**

218 To assess how SARS-CoV-2 variant lineage detection in wastewater compared with  
219 sequencing of clinical samples in the same cities and time period, we considered the top three  
220 most prevalent lineages detected in our WW samples: Alpha (B.1.1.7), B.1.160, and A.2.5.X).  
221 For these analyses, we supplemented the 936 WW sequences with 13,296 clinical samples  
222 from a semi-random population sample of nasal swabs (of which 9,262 passed sequencing  
223 quality filters and were deposited in GISAID; **Table S3**) and 5,661 non-random 'outbreak'  
224 samples sequenced with high priority by the Quebec Public Health lab (of which 1,848 passed  
225 filters and were deposited in GISAID; **Table S4**; Methods). In each sample type, we defined a  
226 variant lineage's frequency relative to other lineages in a 2-week time window. The frequency  
227 estimates from WW and both types of clinical samples are significantly correlated (linear  
228 regression adjusted  $R^2 = 0.63$  with semi-random clinical data, and 0.34 with outbreak data;  
229 Permutational ANOVA both  $p < 2e-4$ , pooled over the three most prevalent variants;  $n = 5000$   
230 permutations). Note that this method ignores missing data, so sampling gaps do not contribute  
231 to the regression. In general, variants are detected first in clinical or outbreak data, and shortly  
232 thereafter in wastewater (**Figure 3**). The Alpha VOC (B.1.1.7) in particular tended to be  
233 detected first in outbreak samples, likely because suspected Alpha cases were prioritized for  
234 sequencing by the Quebec Public Health lab at the time. Notably, Alpha was detected in WW  
235 concurrently or shortly after outbreak samples, despite being detected much later (or not at  
236 all) in semi-random clinical sequences (**Figure 3**).





237

238 **Figure 3 Relative abundance of the most prevalent VOC/VOIs in Quebec in clinical,**  
 239 **outbreak and WW samples.** The relative abundance of each variant was estimated as the  
 240 percentage of samples in which the variant is present over a 2-week period. Gaps in the time  
 241 series can be explained by a lack of sampling (transparent grey), i.e. the absence of detections  
 242 of a particular lineage due to the absence of samples, or missing detections in the sequencing  
 243 data (transparent red), i.e. the absence of detections of a particular lineage although samples  
 244 were collected during the corresponding period. We only show time intervals starting in  
 245 September 2020 because sampling frequency was sparse before that (Figure S7).

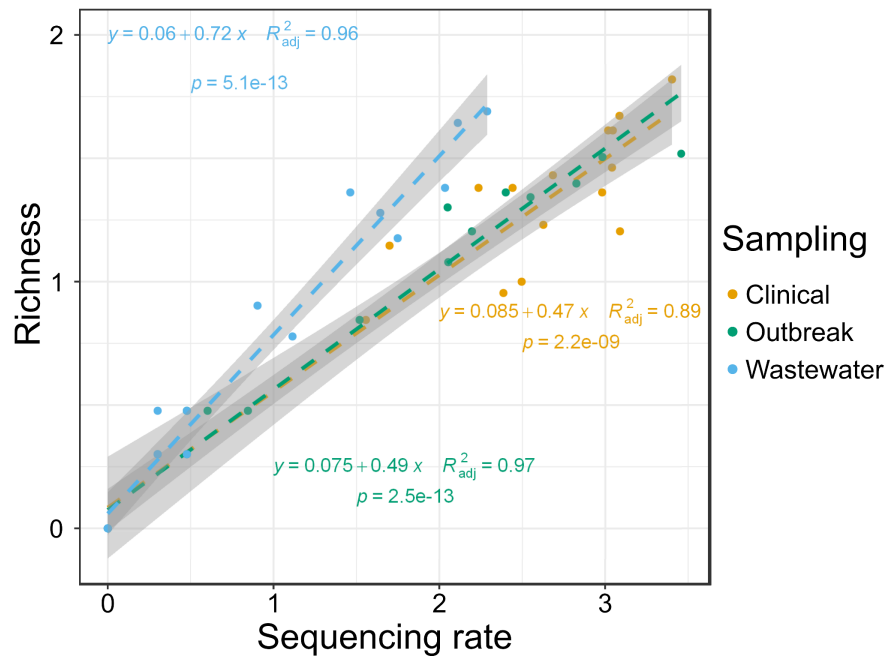
246

247 Next, we sought to generalize these results beyond the most prevalent variants and to  
 248 consider different possible time-lags between clinical and WW samples. For example, clinically  
 249 undetected asymptomatic or pre-symptomatic cases could be detected earlier with WW  
 250 sampling. Conversely, outbreak and clinical samplings might detect some lineages earlier than  
 251 WW when the testing rate is high and appropriately targeted. We considered time lags  
 252 between 0 and 8 days between collection of WW and clinical or outbreak samples, and  
 253 calculated concordance, defined simply as the detection of a variant in both sample types. We  
 254 found that concordance was maximized between WW and both clinical or outbreak samples  
 255 at a maximum time lag of around 7 days, regardless of the number of signature SNVs required  
 256 to define variants (**Figures S8 and S9**). The 7-day time interval is similar to the average pre-  
 257 symptomatic phase duration, i.e. 6.4 days (Backer, Klinkenberg, and Wallinga 2020). To  
 258 assess the significance of these concordance scores, we permuted (1000 times) the variant  
 259 detections across cities and time (days) and compared the concordance score after  
 260 permutations to the original concordance score to obtain a *p*-value. Using a 7-day maximum  
 261 time lag and defining variants with at least three signature SNVs, we estimated that 41.7% of  
 262 WW calls are concordant with semi-random clinical sampling data (permutation test *p*-value =

263 0.23) and 85.7% of WW calls are concordant with outbreak sampling data (permutation test  
264  $p$ -value = 0.009). None of the tested time gaps yield a significant concordance between WW  
265 and semi-random clinical data (**Figure S8**). The stronger concordance between WW and  
266 outbreak sequences is surprising, as one would expect WW to better capture the same  
267 variants as the semi-random clinical samples. Although the reasons for this observation are  
268 unclear, it speaks to the unknown and potentially orthogonal biases implicit in each of the  
269 sampling schemes. Out of the 14 variant lineages detected both in WW and outbreak datasets,  
270 13 are also detected in the clinical dataset (**Figure S10**). This suggests that the weaker  
271 concordance between WW and clinical data relative to outbreak data is not explained by the  
272 identity of the lineage in question, but rather by discrepancies in the timing of detection.  
273 Indeed, 37 lineages were detected earlier by semi-random clinical sampling than WW or  
274 outbreak sampling, including A.2.5.X, B.1.1.7 (Alpha), B.1.160, B.1.526 (Iota) and R.1, while  
275 seven lineages were detected earlier in the outbreak dataset, including P.1 (Gamma),  
276 B.1.617.X (Kappa/Delta), B.1.1.519 and C.37 (Lambda). Only B.1.351 (Beta), B.1.621 (Mu)  
277 and B.1.214.2 were detected first in the WW data.

278 A major reason that more variant lineages are detected earlier in clinical and outbreak  
279 sequences is simply the sampling effort: there are about five times more available outbreak  
280 samples than WW samples in our dataset, and about ten times more semi-random clinical  
281 samples. As a result, there are more lineage detections in clinical data (**Figure S11A**), which  
282 is consistent with another study from the United States (Baaijens et al. 2021) and is explained  
283 by a higher monthly sequencing rate (**Figure S11B**). However, at an equal sequencing rate,  
284 WW is able to detect more variant lineages than clinical or outbreak sampling (**Figure 4**). This  
285 suggests that, for a given sequencing effort, WW surveillance could detect a higher diversity  
286 of variants. Even in our current limited sample of 936 WW sequences, we are able to infer the  
287 presence of certain variants that are not apparent in clinical sequences (**Figure S10**).

288



289  
290  
291  
292  
293  
294  
295

**Figure 4 Wastewater sampling detects more unique variant lineages at a given sequencing rate.** Monthly richness of detected lineages in function of sequencing rate for clinical (gold), outbreak (green) and wastewater (skyblue) samplings. We estimated the sequencing rate as  $\log_{10}(\text{number of sequences per month} + 1)$  and richness as  $\log_{10}(\text{number of variant lineages detected} + 1)$ .

## 296 CONCLUSION

297 Using a relatively conservative SNV calling pipeline and lineage detection method, we tracked  
298 the spread of SARS-CoV-2 lineages in wastewater between March 2020 and July 2021 across  
299 three cities in the province of Quebec. Although WW sequencing could be used to detect novel  
300 mutations or variant lineages not found in clinical samples (Smyth et al. 2021), here we  
301 focused on identifying known variants. Consistent with a similar study in the US, we found that  
302 sequencing coverage in WW was dependent on viral load, as measured by the qPCR cycle  
303 threshold (Ct) value (Baaijens et al. 2021). This suggests that WW sample concentration or  
304 optimization of RNA extraction could yield more sequence data and improve variant inference.  
305 We also found that the most prevalent variants were detected concordantly in WW and clinical  
306 samples over a 7-day time frame, but that variants were generally detected first in clinical  
307 samples, suggesting that clinical diagnostics efforts were effective at the time of sampling. We  
308 note that certain VOCs and VOIs, including B.1.351 (Beta), B.1.214.2, and B.1.621 (Mu) were  
309 detected in WW before clinical samples. Several other variants were detected only in WW but  
310 not clinical samples. These could be true positives that were undetected by clinical sampling,  
311 or false positives that could have arisen for a variety of reasons, including the threshold  
312 number of signature or marker SNVs required to identify a lineage. In other words, even if we  
313 trust the SNV calls, recurrent SNVs could potentially confuse one variant lineage for another.  
314 Outbreak samples sequenced as part of targeted public health investigations were particularly

315 concordant with WW samples, for reasons that remain unclear. Despite sources of error in  
316 sampling and sequencing, WW-based detection of variant lineages is generally aligned with  
317 clinical samples from the same city and time period. Our analyses pooled a variety of sample  
318 types at the city level. Future studies could examine the spatial scales (province, city,  
319 neighborhood, residence) at which wastewater sequencing is most concordant with clinical  
320 sampling.

321  
322 Importantly, WW samples are able to detect more variant lineages per sequencing run. This  
323 is intuitive, because WW includes a mixture of viruses, thereby often sampling multiple  
324 individuals at a time, depending on the wastewater catchment area. This mixture of viruses  
325 must then be inferred based on sequencing data. Here we used Illumina short-read data, and  
326 a SNV calling pipeline benchmarked on known standard SARS-CoV-2 genomes. We then  
327 inferred variant lineages as combinations of SNVs present in known variants, based on either  
328 a simple threshold of signature and marker mutations, or on a constrained linear model,  
329 yielding similar results. Long-read technology is potentially useful for resolving multiple SNVs  
330 on the same sequencing read, potentially providing more direct evidence for variant genomes  
331 and circumventing the need for inference based on unlinked SNV frequencies. However, RNA  
332 in wastewater is generally quite fragmented, which may impose an upper limit on the utility of  
333 such methods. Furthermore, the single read accuracy of Nanopore sequencing is to date not  
334 at par with the read accuracy offered by Illumina.

335  
336 While there is room for methodological improvements ranging from sample collection,  
337 sequencing, and computational inference of variants, the general concordance between  
338 wastewater and clinical sampling is encouraging. In contexts where clinical sampling is  
339 infrequent or infeasible, wastewater can provide a complementary window into VOC  
340 frequencies. The value of wastewater sequencing became exceedingly clear during the  
341 Omicron wave, which began in Canada in mid-December 2021 and is still ongoing in late  
342 January 2022. Anecdotally, we were able to detect Omicron with 11 signature mutations on  
343 December 4, 2021 in Montreal (data not shown). Together, our results suggest that  
344 wastewater sequencing can continue to provide a similar portrait of SARS-CoV-2 variant  
345 lineages, potentially with much less sampling and sequencing effort.

346

## 347 **ACKNOWLEDGEMENTS**

348 We are particularly grateful to Carole Fleury, Suzanne Boulet, Luc Carrière from the City of  
349 Montreal, David Kaiser from Montreal Public Health, Mario Gagné from the City of Laval and  
350 Diarra Zeinab from the Centre de technologies de l'eau for their assistance in obtaining  
351 samples. We also thank Frédéric Cloutier, Denis Dufour and François Proulx from Quebec  
352 City for assistance with sampling and logistics support.

353 This study was supported by the Canadian Institutes for Health Research (CIHR) operating  
354 grant to the Coronavirus Variants Rapid Response Network (CoVaRR-Net) to JR and BJS and  
355 the CFI project “High throughput SARS-CoV-2 genome sequencing at the McGill Genome  
356 Center” to JR. DF was supported by the Fonds de Recherche du Québec (FRQ) and a McGill  
357 MI4 Emergency COVID-19 Research Fund grant. SD was supported by FRQ and The Trottier  
358 and Molson Foundations. Data analyses were enabled by compute and storage resources  
359 provided by Compute Canada and Calcul Québec. We would like to acknowledge the team at  
360 the Canadian Centre for Computational Genomics for their contributions to the development  
361 and testing of the pipeline used for these analyses.

362

### 363 **ETHICS STATEMENT**

364 This investigation was carried out in accordance with the legal mandate granted to public  
365 health authorities by the Public Health Act (LRQ, chapter S-2.2. Article 1;  
366 <http://legisquebec.gouv.qc.ca/fr/ShowDoc/cs/S-2.2>) as part of a public health intervention. All  
367 data was treated confidentially and analysed without nominal identification in accordance with  
368 the Policy on Information Protection and Security (PO-04-2014) of the National Public Health  
369 Institute of Quebec (INSPQ). The INSPQ therefore deemed this study exempt from ethical  
370 oversight according to provincial legislation.

371

### 372 **DATA & CODE AVAILABILITY**

373 Raw wastewater sequencing data is available in the NCBI SRA database under BioProject  
374 ID PRJNA788395 (<http://www.ncbi.nlm.nih.gov/bioproject/788395>). Viral genomes from  
375 clinical samples are available in GISAID under IDs listed in Tables S3 and S4. The SNV  
376 calling and post-variant-calling pipeline is available at:  
377 [https://github.com/arnaud00013/Wastewater\\_surveillance\\_pipeline](https://github.com/arnaud00013/Wastewater_surveillance_pipeline).

## 378 **METHODS**

### 379 **Sample collection and sequencing**

380 To perform the genomic surveillance of SARS-CoV-2, we collected WW samples at the  
381 institutional, district, and municipal scale from Montreal, Quebec City, and Laval. The samples  
382 were collected by composite sampling, grab sampling, and passive sampling. The composite  
383 samples were collected with autosamplers, which collected wastewater every 10 minutes, over  
384 a 24, 48, or 72 hour time period. The passive samples were collected through 2 absorbent  
385 materials, Q-tips and negatively charged membranes, Mixed Cellulose Ester (MCE) filters,  
386 which were housed in torpedoes (Schang et al. 2021). The torpedoes were also deployed over  
387 a 24, 48, or 72 hour time period. For time series analyses, the samples were pooled from the  
388 different scales (institution, district) and from the different sampling methods (composite, grab,  
389 and passive) for each date; thus the analyses were conducted at the municipal scale.  
390 Wastewater samples, grab or composite, were additionally concentrated by filtration. First, in  
391 50 ml of wastewater, the pH was adjusted to between 3.5-4.5, and MgCl was added to a final  
392 concentration of 25 mM. Then the samples were filtered through a 0.45 µm MCE filter. All  
393 samples were processed within at most 72 hours; the MCE filters and Q-tips were stored at -  
394 80°C. RNA was then extracted using the Allprep Powerviral DNA/RNA kit. The protocol was  
395 followed according to the manufacturer, with the exception of the lysis step, where a final  
396 concentration of 10% Beta-Mercaptoethanol was used in the lysis buffer and incubation time  
397 was raised to 30 minutes at 55°C. After extraction, RNA samples were submitted to the McGill  
398 Genome Center for reverse transcription followed by targeted SARS-CoV-2 amplification  
399 using the ARTIC V3 primer scheme ([https://artic.network/resources/ncov/ncov-amplicon-](https://artic.network/resources/ncov/ncov-amplicon-v3.pdf)  
400 [v3.pdf](https://artic.network/resources/ncov/ncov-amplicon-v3.pdf)). Samples were purified and a Nextera DNA Flex library preparation was performed for  
401 Illumina paired-end amplicon sequencing (PE150) on a NovaSeq instrument at the McGill  
402 Genome Center. The detailed protocol can be accessed at:  
403 [dx.doi.org/10.17504/protocols.io.by6xpzfn](https://doi.org/10.17504/protocols.io.by6xpzfn).

404

### 405 **SNV calling and quality control**

406 For each sample, we performed quality control of the raw reads using fastp (v.0.20.0, read  
407 length >=70, Phred Score >20 and cut\_tail option (S. Chen et al. 2018). We then aligned reads  
408 to the reference genome (MN908947.3 or NC\_045512.2) using bwa (v.0.7.17) (Li 2013) and  
409 performed a coverage analysis with samtools (v.1.10) depth (Li 2011). In total, 936 WW  
410 samples were collected and sequenced with Illumina. For each of these samples, we  
411 performed SNV calling from the read mapping using samtools (v.1.10) mpileup and varscan  
412 (v.2.4.1) pileup2snp (Li 2011; Koboldt et al. 2012). To select coverage and SNV frequency  
413 filters for SNV calling, we analyzed 14 positive controls from AccuGenomics for which we  
414 know the list of expected mutations (<https://accugenomics.com/accukit-sars-cov-2/>). These

415 controls allow us to measure a limit of detection and to evaluate the background error due to  
416 library preparation, sequencing, or other sources of noise. Using these controls, we can then  
417 determine the variant calling filters that allows to minimize the background error, i.e. the false  
418 positive and false negative rates, and maximizes the accuracy of SNV calling, i.e. F1 score.  
419 We calculated this score using the precision, i.e. the proportion of SNV calls (positives) that  
420 are true positives, and the recall/sensitivity, i.e. the proportion of expected variants that have  
421 been called (Powers 2020). The SNV calling and post-variant-calling pipeline is available at  
422 [https://github.com/arnaud00013/Wastewater\\_surveillance\\_pipeline](https://github.com/arnaud00013/Wastewater_surveillance_pipeline).

423 In addition to the AccuGenomics controls, we also sequenced mixtures of two different  
424 SARS-CoV-2 viral cultures at known ratios, which we called “spike-in” samples. In these  
425 samples, a viral culture extract was spiked into a SARS-CoV-2 positive control sample at a  
426 concentration of 1%, 2%, 5%, 10%, 20% or 50%. Both samples were obtained from the first  
427 wave of the pandemic and differed by 11 mutations . We performed variant calling on the  
428 sequenced spiked-in samples based on the same method as the WW samples, except for the  
429 VAF filter, which we removed to be able to detect expected low-frequency mutations. Because  
430 each Spike-sample is a mixture of variants and are thus not clonal, the expected variant  
431 frequency (VAF) is calculated based on the virus concentration and initial VAF in the samples.  
432  $Expected\_VAF = (Concentration\_positive\_control * VAF\_in\_positive\_ctrl) +$   
433  $(Concentration\_ViralCulture * VAF\_in\_L00241026\_ViralCulture)$  where  
434  $Concentration\_positive\_control$  would be 1%, 2%, 5%, 10%, 20% and 50%, and  
435  $Concentration\_ViralCulture = 1 - Concentration\_positive\_control$ .

436

### 437 **Detection of lineages of interest in WW samples and estimation of their within-sample** 438 **frequency**

439 To infer the presence of SARS-CoV-2 lineages in WW samples, we used the number of  
440 lineage signature mutations, i.e. the number of mutations that have a minimum prevalence of  
441 90% among the consensus sequences of the lineage. To find these signature mutations, we  
442 first calculated the prevalence of substitutions in thousands of publicly available consensus  
443 sequences collected during 2020 and added data from CoV-Spectrum about under-  
444 represented lineage in the database or lineages that emerged during 2021 (C. Chen et al.  
445 2021). The database contains 755 lineages and is available at  
446 [https://github.com/arnaud00013/Wastewater\\_surveillance\\_pipeline/blob/main/](https://github.com/arnaud00013/Wastewater_surveillance_pipeline/blob/main/Post_variant_calling_analysis/Database_all_mutations_prevalence_in_SC2_lineages_consensus_sequences_as_of_2021_07_08.json)  
447 [Post\\_variant\\_calling\\_analysis/Database\\_all\\_mutations\\_prevalence\\_in\\_SC2\\_lineages\\_conse](https://github.com/arnaud00013/Wastewater_surveillance_pipeline/blob/main/Post_variant_calling_analysis/Database_all_mutations_prevalence_in_SC2_lineages_consensus_sequences_as_of_2021_07_08.json)  
448 [nsus\\_sequences\\_as\\_of\\_2021\\_07\\_08.json](https://github.com/arnaud00013/Wastewater_surveillance_pipeline/blob/main/Post_variant_calling_analysis/Database_all_mutations_prevalence_in_SC2_lineages_consensus_sequences_as_of_2021_07_08.json). At least 3 signature mutations including a marker  
449 mutation, i.e. a substitution that have a very high prevalence ( $\geq 90\%$ ) only among consensus  
450 sequences from a certain lineage, are required to call a lineage present in a WW sample. We  
451 selected this filter arbitrarily as a compromise between confidence of detection and

452 sensitivity/stringency. We also made sure to evaluate the effect of selecting different filters on  
453 the concordance between WW and clinical sampling.

454 For the calculation of within-sample frequency of SARS-CoV-2 lineages, we used a  
455 linear model to fit the lineages' signature mutations data to the within-sample mutation  
456 frequency data. The rationale of the approach is that the frequency of lineage signature  
457 mutations within a sample is a linear combination of the frequency of the lineages and the  
458 prevalence of the mutations in the consensus sequences of these lineages. To make sure that  
459 linear regression converged to a good solution, we needed to apply certain constraints (**Figure**  
460 **S6**). Thus, we implemented this analysis using the "ConsReg()" function from the R (R  
461 Development Core Team 2011) package ConsReg (v.0.1.0), which allows us to perform linear  
462 regressions under constraints for regression coefficients. We explored the space of solutions  
463 using a Grid Search (coefficients initial values = 0.1,0.5 or 0.9) and a Monte Carlo Markov  
464 Chain optimizer.

465

## 466 REFERENCES

467 Baaijens, Jasmijn A., Alessandro Zulli, Isabel M. Ott, Mary E. Petrone, Tara Alpert,  
468 Joseph R. Fauver, Chaney C. Kalinich, et al. 2021. "Variant Abundance Estimation for  
469 SARS-CoV-2 in Wastewater Using RNA-Seq Quantification." *medRxiv*, September,  
470 2021.08.31.21262938.

471  
472 Backer, Jantien A., Don Klinkenberg, and Jacco Wallinga. 2020. "Incubation Period of  
473 2019 Novel Coronavirus (2019-nCoV) Infections among Travellers from Wuhan, China,  
474 20-28 January 2020." *Euro Surveillance* 25 (5). [https://doi.org/10.2807/1560-](https://doi.org/10.2807/1560-7917.ES.2020.25.5.2000062)  
475 [7917.ES.2020.25.5.2000062](https://doi.org/10.2807/1560-7917.ES.2020.25.5.2000062).

476  
477 Bibby, Kyle, Aaron Bivins, Zhenyu Wu, and Devin North. 2021. "Making Waves:  
478 Plausible Lead Time for Wastewater Based Epidemiology as an Early Warning System  
479 for COVID-19." *Water Research* 202 (September): 117438.

480  
481 Campbell, Finlay, Brett Archer, Henry Laurenson-Schafer, Yuka Jinnai, Franck Konings,  
482 Neale Batra, Boris Pavlin, et al. 2021. "Increased Transmissibility and Global Spread of  
483 SARS-CoV-2 Variants of Concern as at June 2021." *Euro Surveillance* 26 (24).  
484 <https://doi.org/10.2807/1560-7917.ES.2021.26.24.2100509>.

485  
486 Chen, Chaoran, Sarah Nadeau, Michael Yared, Philippe Voinov, Ning Xie, Cornelius  
487 Roemer, and Tanja Stadler. 2021. "CoV-Spectrum: Analysis of Globally Shared SARS-  
488 CoV-2 Data to Identify and Characterize New Variants." *Bioinformatics* , December.  
489 <https://doi.org/10.1093/bioinformatics/btab856>.

490  
491 Chen, Shifu, Yanqing Zhou, Yaru Chen, and Jia Gu. 2018. "Fastp: An Ultra-Fast All-in-  
492 One FASTQ Preprocessor." *Bioinformatics* 34 (17): i884–90.

493  
494 Crits-Christoph, Alexander, Rose S. Kantor, Matthew R. Olm, Oscar N. Whitney, Basem  
495 Al-Shayeb, Yue C. Lou, Avi Flamholz, et al. 2021. "Genome Sequencing of Sewage  
496 Detects Regionally Prevalent SARS-CoV-2 Variants." *mBio* 12 (1): e02703–20.

497



- 498 Davies, Nicholas G., Sam Abbott, Rosanna C. Barnard, Christopher I. Jarvis, Adam J.  
499 Kucharski, James D. Munday, Carl A. B. Pearson, et al. 2021. "Estimated  
500 Transmissibility and Impact of SARS-CoV-2 Lineage B.1.1.7 in England." *Science*,  
501 March. <https://doi.org/10.1126/science.abg3055>.  
502
- 503 Institut national de santé publique du Québec. 2021. "Définitions Pour La Vigie  
504 Sanitaire Des Variants Du SRAS-CoV-2 et Classification Des Lignées Détectées Au  
505 Québec." 2021. [https://www.inspq.qc.ca/sites/default/files/publications/3138-definition-  
506 vigie-sanitaire-variants-sras-cov-2.pdf](https://www.inspq.qc.ca/sites/default/files/publications/3138-definition-vigie-sanitaire-variants-sras-cov-2.pdf).  
507
- 508 Jones, David L., Marcos Quintela Baluja, David W. Graham, Alexander Corbishley,  
509 James E. McDonald, Shelagh K. Malham, Luke S. Hillary, et al. 2020. "Shedding of  
510 SARS-CoV-2 in Feces and Urine and Its Potential Role in Person-to-Person  
511 Transmission and the Environment-Based Spread of COVID-19." *Science of the Total  
512 Environment* 749 (December): 141364.  
513
- 514 Koboldt, Daniel C., Qunyuan Zhang, David E. Larson, Dong Shen, Michael D. McLellan,  
515 Ling Lin, Christopher A. Miller, Elaine R. Mardis, Li Ding, and Richard K. Wilson. 2012.  
516 "VarScan 2: Somatic Mutation and Copy Number Alteration Discovery in Cancer by  
517 Exome Sequencing." *Genome Research* 22 (3): 568–76.  
518
- 519 Li, Heng. 2011. "A Statistical Framework for SNP Calling, Mutation Discovery,  
520 Association Mapping and Population Genetical Parameter Estimation from Sequencing  
521 Data." *Bioinformatics* 27 (21): 2987–93.  
522
- 523 ———. 2013. "Aligning Sequence Reads, Clone Sequences and Assembly Contigs  
524 with BWA-MEM." *arXiv [q-bio.GN]*. arXiv. <http://arxiv.org/abs/1303.3997>.  
525
- 526 Lythgoe, Katrina A., Matthew Hall, Luca Ferretti, Mariateresa de Cesare, George  
527 MacIntyre-Cockett, Amy Trebes, Monique Andersson, et al. 2021. "SARS-CoV-2 within-  
528 Host Diversity and Transmission." *Science*, March.  
529 <https://doi.org/10.1126/science.abg0821>.  
530
- 531 Mlcochova, Petra, Steven A. Kemp, Mahesh Shanker Dhar, Guido Papa, Bo Meng,  
532 Isabella A. T. M. Ferreira, Rawlings Datir, et al. 2021. "SARS-CoV-2 B.1.617.2 Delta  
533 Variant Replication and Immune Evasion." *Nature*, September.  
534 <https://doi.org/10.1038/s41586-021-03944-y>.  
535
- 536 Murall, Carmen Lía, Eric Fournier, Jose Hector Galvez, Arnaud N'Guessan, Sarah J.  
537 Reiling, Pierre-Olivier Quirion, Sana Naderi, et al. 2021. "A Small Number of Early  
538 Introductions Seeded Widespread Transmission of SARS-CoV-2 in Québec, Canada."  
539 *Genome Medicine* 13 (1): 169.  
540
- 541 Murall, Carmen Lía, Fatima Mostefai, Jean-Christophe Grenier, Raphaël Poujol, Julie  
542 Hussin, Sandrine Moreira, B. Jesse Shapiro, CoVSeQ consortium. 2021. "Recent  
543 Evolution and International Transmission of SARS-CoV-2 Clade 19B (Pango A  
544 Lineages)." June 2, 2021. [https://virological.org/t/recent-evolution-and-international-  
545 transmission-of-sars-cov-2-clade-19b-pango-a-lineages/711](https://virological.org/t/recent-evolution-and-international-transmission-of-sars-cov-2-clade-19b-pango-a-lineages/711).  
546
- 547 Otto, Sarah P., Troy Day, Julien Arino, Caroline Colijn, Jonathan Dushoff, Michael Li,  
548 Samir Mechai, et al. 2021. "The Origins and Potential Future of SARS-CoV-2 Variants  
549 of Concern in the Evolving COVID-19 Pandemic." *Current Biology*.  
550 <https://doi.org/10.1016/j.cub.2021.06.049>.  
551
- 552 Powers, David M. W. 2020. "Evaluation: From Precision, Recall and F-Measure to

- 553 ROC, Informedness, Markedness and Correlation.” *arXiv [cs.LG]*. arXiv.  
554 <http://arxiv.org/abs/2010.16061>.  
555  
556 Public Health Agency of Canada. 2021. “COVID-19 Daily Epidemiology Update.” 2021.  
557 [https://health-infobase.canada.ca/covid-19/epidemiological-summary-covid-19-](https://health-infobase.canada.ca/covid-19/epidemiological-summary-covid-19-cases.html?stat=rate&measure=deaths&map=pt#a2)  
558 [cases.html?stat=rate&measure=deaths&map=pt#a2](https://health-infobase.canada.ca/covid-19/epidemiological-summary-covid-19-cases.html?stat=rate&measure=deaths&map=pt#a2).  
559  
560 Rambaut, Andrew, Edward C. Holmes, Áine O’Toole, Verity Hill, John T. McCrone,  
561 Christopher Ruis, Louis du Plessis, and Oliver G. Pybus. 2020. “A Dynamic  
562 Nomenclature Proposal for SARS-CoV-2 Lineages to Assist Genomic Epidemiology.”  
563 *Nature Microbiology* 5 (11): 1403–7.  
564  
565 R Development Core Team. 2011. “R: A Language and Environment for Statistical  
566 Computing.” *R Foundation for Statistical Computing*.  
567  
568 Schang, Christelle, Nicolas D. Crosbie, Monica Nolan, Rachael Poon, Miao Wang,  
569 Aaron Jex, Nijoy John, et al. 2021. “Passive Sampling of SARS-CoV-2 for Wastewater  
570 Surveillance.” *Environmental Science & Technology* 55 (15): 10432–41.  
571  
572 Smyth, Davida S., Monica Trujillo, Devon A. Gregory, Kristen Cheung, Anna Gao,  
573 Maddie Graham, Yue Guan, et al. 2021. “Tracking Cryptic SARS-CoV-2 Lineages  
574 Detected in NYC Wastewater.” medRxiv. <https://doi.org/10.1101/2021.07.26.21261142>.  
575  
576 Volz, Erik, Swapnil Mishra, Meera Chand, Jeffrey C. Barrett, Robert Johnson, Lily  
577 Geidelberg, Wes R. Hinsley, et al. 2021. “Assessing Transmissibility of SARS-CoV-2  
578 Lineage B.1.1.7 in England.” *Nature* 593 (7858): 266–69.  
579  
580 Xiao, Amy, Fuqing Wu, Mary Bushman, Jianbo Zhang, Maxim Imakaev, Peter R. Chai,  
581 Claire Duvallet, et al. 2021. “Metrics to Relate COVID-19 Wastewater Data to Clinical  
582 Testing Dynamics.” *medRxiv*. June. <https://doi.org/10.1101/2021.06.10.21258580>.