

1 **Genetic Diversity and Evolutionary Convergence of Cryptic SARS-CoV-2 Lineages**
2 **Detected Via Wastewater Sequencing**

3
4 Devon A. Gregory¹, Monica Trujillo², Clayton Rushford¹, Anna Flury³, Sherin Kannoly⁴,
5 Kaung Myat San⁴, Dustin Lyfoung⁵, Roger W. Wiseman⁵, Karen Bromert⁶, Ming-Yi Zhou⁶,
6 Ellen Kesler⁶, Nathan Bivens⁶, Jay Hoskins⁷, Chung-Ho Lin⁸, David H. O'Connor⁵, Chris
7 Wieberg⁹, Jeff Wenzel¹⁰, Rose S. Kantor^{11#}, John J. Dennehy^{3,4#}, Marc C. Johnson^{1,#}

8
9 ¹Department of Molecular Microbiology and Immunology, University of Missouri-School
10 of Medicine, Columbia, MO, USA

11 ²Department of Biological Sciences and Geology, Queensborough Community College of
12 The City University of New York, Queens, NY, USA

13 ³Biology Doctoral Program, The Graduate Center of The City University of New York,
14 NYC, NY, USA

15 ⁴Biology Department, Queens College of The City University of New York, Queens, NY,
16 USA 11367

17 ⁵Department of Pathology and Laboratory Medicine, University of Wisconsin-Madison,
18 Madison, WI, USA 53706

19 ⁶Genomics Technology Core, University of Missouri, Columbia, MO, USA

20 ⁷Environmental Compliance Division, Engineering Department, Metropolitan St. Louis
21 Sewer District, St. Louis, MO, USA 63103

22 ⁸Center of Agroforestry, School of Natural Resources, University of Missouri, Columbia,
23 MO, USA

24 ⁹Water Protection Program, Missouri Department of Natural Resources, Jefferson City,
25 MO, USA

26 ¹⁰Bureau of Environmental Epidemiology, Division of Community and Public Health,
27 Missouri Department of Health and Senior Services, Jefferson City, MO, USA

28 ¹¹Department of Civil and Environmental Engineering, University of California, Berkeley,
29 663 Davis Hall, Berkeley, CA, USA 94720

30 #Corresponding Authors

31 Abstract

32 Wastewater-based epidemiology (WBE) is an effective way of tracking the appearance
33 and spread of SARS-COV-2 lineages through communities. Beginning in early 2021, we
34 implemented a targeted approach to amplify and sequence the receptor binding domain
35 (RBD) of SARS-COV-2 to characterize viral lineages present in sewersheds. Over the
36 course of 2021, we reproducibly detected multiple SARS-COV-2 RBD lineages that have
37 never been observed in patient samples in 9 sewersheds located in 3 states in the USA.
38 These cryptic lineages contained between 4 to 24 amino acid substitutions in the RBD
39 and were observed intermittently in the sewersheds in which they were found for as long
40 as 14 months. Many of the amino acid substitutions in these lineages occurred at residues
41 also mutated in the Omicron variant of concern (VOC), often with the same substitution.
42 One of the sewersheds contained a lineage that appeared to be derived from the Alpha
43 VOC, but the majority of the lineages appeared to be derived from pre-VOC SARS-COV-
44 2 lineages. Specifically, several of the cryptic lineages from New York City appeared to
45 be derived from a common ancestor that most likely diverged in early 2020. While the

46 source of these cryptic lineages has not been resolved, it seems increasingly likely that
47 they were derived from immunocompromised patients or animal reservoirs. Our findings
48 demonstrate that SARS-COV-2 genetic diversity is greater than what is commonly
49 observed through routine SARS-CoV-2 surveillance. Wastewater sampling may more
50 fully capture SARS-CoV-2 genetic diversity than patient sampling and could reveal new
51 VOCs before they emerge in the wider human population.

52

53 Author Summary

54 During the COVID-19 pandemic, wastewater-based epidemiology has become an
55 effective public health tool. Because many infected individuals shed SARS-CoV-2 in
56 feces, wastewater has been monitored to reveal infection trends in the sewersheds from
57 which the samples were derived. Here we report novel SARS-CoV-2 lineages in
58 wastewater samples obtained from 3 different states in the USA. These lineages
59 appeared in specific sewersheds intermittently over periods of up to 14 months, but
60 generally have not been detected beyond the sewersheds in which they were initially
61 found. Many of these lineages may have diverged in early 2020. Although these lineages
62 share considerable overlap with each other, they have never been observed in patients
63 anywhere in the world. While the wastewater lineages have similarities with lineages
64 observed in long-term infections of immunocompromised patients, animal reservoirs
65 cannot be ruled out as a potential source.

66

67 1. Introduction

68 SARS-CoV-2 is shed in feces of infected individuals [1,2], and SARS-CoV-2 RNA can be
69 extracted and quantified from community wastewater to provide estimates of SARS-CoV-
70 2 community prevalence [3,4]. This approach is especially powerful since it randomly
71 samples all community members and can detect viruses shed by individuals whose
72 infections are not recorded, such as asymptomatic individuals, those who abstain from
73 testing, or those who test at home [5,6]. Additionally, SARS-CoV-2 RNA isolated from
74 wastewater can be sequenced using high-throughput sequencing technologies to define
75 the composition of variants in the community [7–9].

76
77 The continuing evolution of SARS-CoV-2 [10] and the appearance of variants of concern
78 (VOC), such as the Omicron VOC [11], highlight the importance of maintaining a vigilant
79 watch for the emergence of unexpected, novel variants. The fact that the origins and early
80 spread of the Alpha and Omicron VOCs were not observed strongly motivates efforts to
81 detect and monitor novel variants [12]. However, whole genome sequencing of SARS-
82 CoV-2 RNA isolated from wastewater often suffers from low sequencing depth of
83 coverage in epidemiologically relevant areas of the genome, such as the Spike receptor
84 binding domain (RBD)[13–15]. Additionally, because wastewater may contain a mixture
85 of viral lineages and whole genome sequencing relies on sequencing small fragments of
86 the genome, computational strategies to identify variants with linked mutations often fail
87 to identify lineages present at low concentrations [16]. These features have made it

88 difficult to detect unexpected, novel variants from wastewater samples from whole
89 genome sequencing data.

90

91 To address these issues, we developed a “targeted” sequencing approach that amplifies
92 and sequences the Spike RBD of the SARS-CoV-2 genome as a single amplicon (Fig.
93 1A) [8,9]. Since the Spike RBD is relevant to SARS-CoV-2 infectivity, transmission, and
94 antibody-mediated neutralization [17–21], this approach ensures that the RBD receives
95 high sequencing coverage. Additionally, RBD sequencing enables linkage of
96 polymorphisms, forming short, phased haplotypes [16]. These phased haplotypes permit
97 easier lineage identification, even at low concentrations, if the targeted sequence(s) are
98 rich in lineage-defining polymorphisms [9].

99

100 Using our targeted sequencing approach, we identified and previously reported circulating
101 VOCs in different sewersheds around the United States [8,9]. Variant frequencies in these
102 sewersheds closely tracked VOCs frequency estimates from clinical sampling in the same
103 areas [8,9]. However, in some locations, we noted the presence of cryptic lineages not
104 observed in clinical samples anywhere in the world. Several of these lineages contained
105 amino acid substitutions rarely reported in global databases such as gisaid.org [22–24]
106 (e.g., N460K, Q493K, Q498Y, and N501S) [8]. Interestingly, polymorphisms in these
107 lineages show considerable overlap with the Omicron VOC, suggesting convergent
108 evolution due to similar selective pressures.

109

110 Here we describe an expanded set of cryptic lineages from multiple locations around the
111 United States. While each sewershed contains its own signature lineages and at least
112 some of the lineages appear to have diverged independently from one another, we
113 present evidence that some likely shared a common ancestor. Finally, we show evidence
114 of strong positive selection and rapid divergence of these lineages from ancestral SARS-
115 CoV-2.

116 2. Results

117 Beginning in early 2021, wastewater surveillance programs including RBD amplicon
118 sequencing (Fig 1A) were independently implemented in Missouri [9] and NYC [25]. A
119 similar strategy was subsequently adopted in California. While the vast majority of
120 sequences observed with this method matched to known lineages identified in patients,
121 reproducible lineages that did not match the known circulating lineages were also
122 detected. Herein, we refer to each RBD haplotype with a unique combination of amino
123 acid changes as a *lineage*, and combinations of lineages that all have specific amino acid
124 changes in common as *lineage classes*. Amino acid combinations identified that have not
125 been seen previously from patients are referred to as *cryptic* lineages. Here we describe
126 cryptic lineages detected from January 1, 2021 through March 15, 2022.

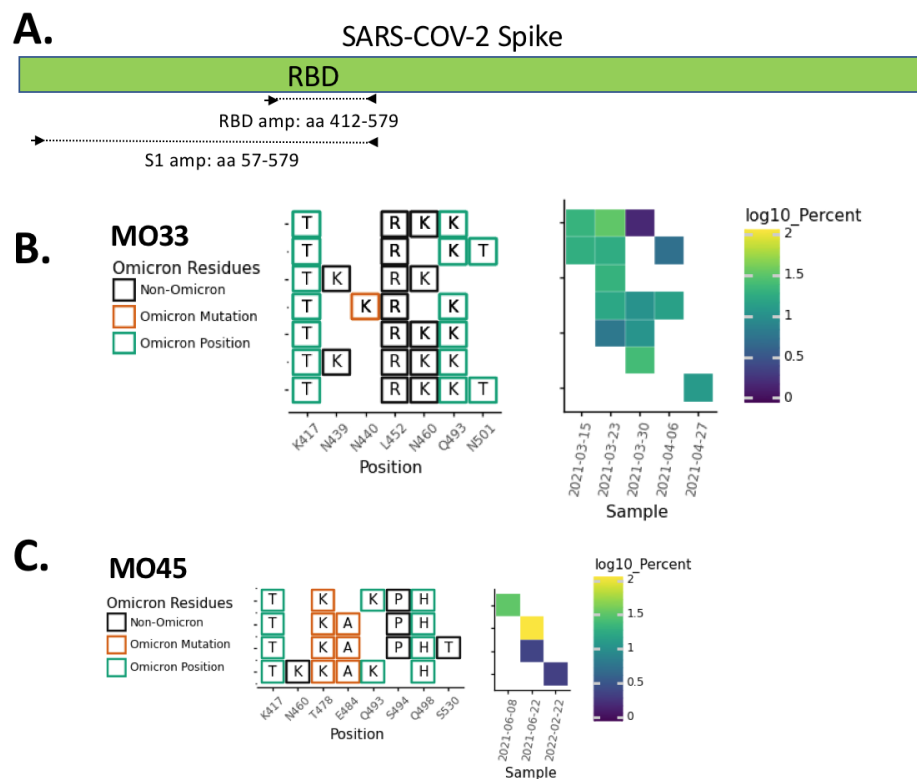
127 2.1 Lineage persistence and evolution over time

128 In total, cryptic lineages were observed in 9 sewersheds across 3 states (Table 1). Each
129 cryptic lineage class was generally unique to a sewershed. These lineages contained

130 between 4-24 non-synonymous substitutions, insertions, and deletions. In some cases,
131 lineages were detected for a short duration but with multiple similar co-occurring
132 sequences. For example, in Missouri sewershed MO33, a lineage class containing 4-5
133 RBD amino acid changes were consistently detected at low relative abundances from
134 March 15 to the end of April 2021 (Fig. 1B, Table 1). A total of 7 unique sequences were
135 spread across the 5 sampling events in this date range, and up to 5 unique sequences
136 co-occurred within a given sample.

137 Meanwhile, in other sewersheds, cryptic lineages were detected briefly, before
138 disappearing, and then reappearing many months later. For example, in Missouri
139 sewershed MO45, lineages were first detected in June 2021 and then were not seen
140 again until February 2022 (Fig. 1C, Table 1). The longest observed lineage class was in
141 sewershed NY3 where we previously reported a lineage class from January 2021 [8] that
142 was detected sporadically until March 2022 (Fig. 2A, Table 1). On average, cryptic
143 lineages lasted for around 6 months, such as the lineage class from NY14 which lasted
144 from May to October, 2021 (Fig. 2B).

145

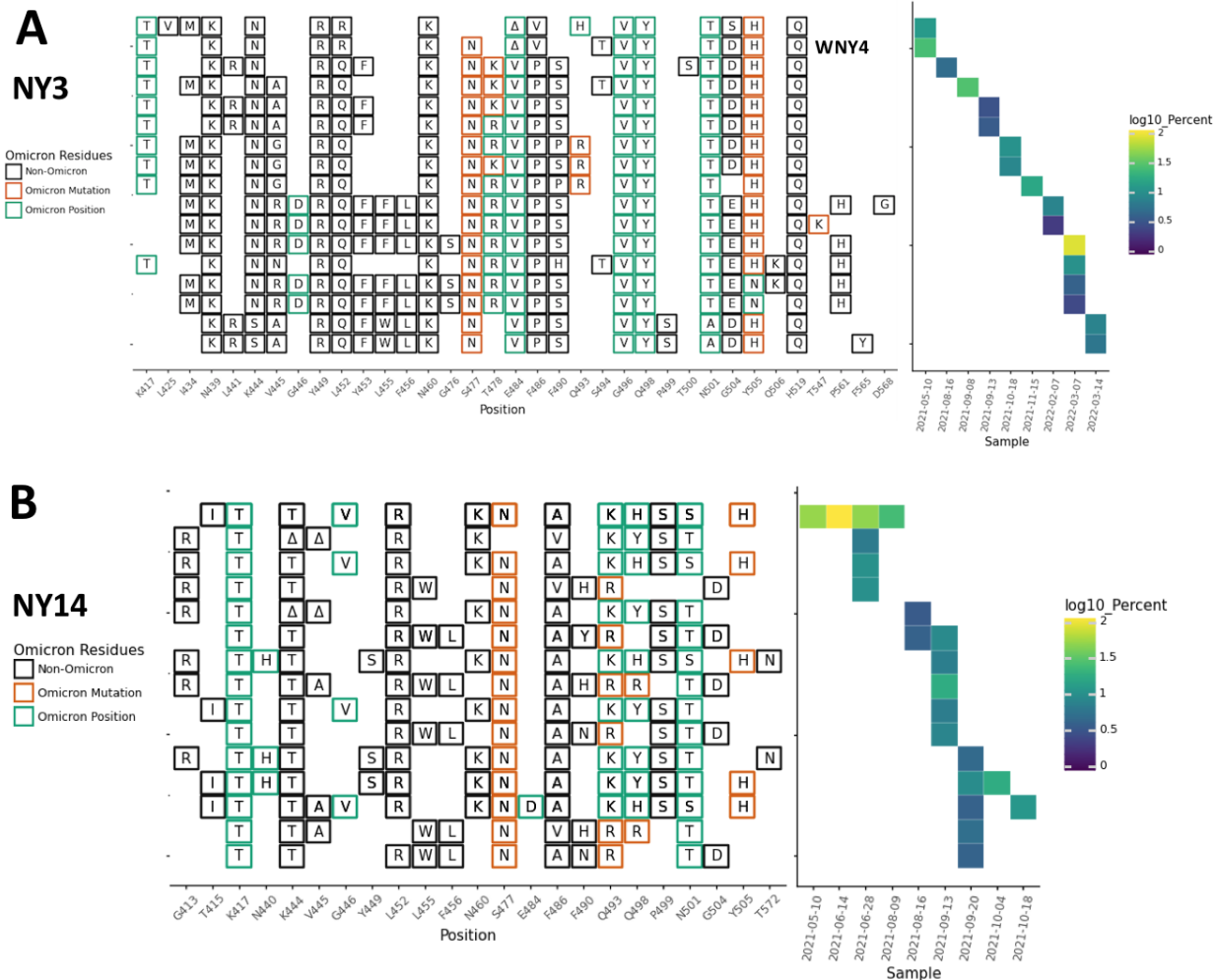


146

147 **Fig. 1** RBD amplification. A. Schematic of regions targeted by the RBD and S1 primer
 148 sets (see Methods for primer sequences). Overview of the SARS-COV-2 Spike RBD
 149 lineages identified in B. the MO33 sewershed and C. the MO45 sewershed. Each row
 150 represents a unique lineage and each column is an amino acid position in the Spike
 151 protein (left). Amino acid changes similar to (green boxes) or identical to (orange boxes)
 152 changes in Omicron (BA.1) are indicated. The heatmap (right) illustrates lineage (row)
 153 detection by date (column), colored by the log₁₀ percent relative abundance of that
 154 lineage.

155

156



157

158

159 **Fig. 2** NY3 and NY14 RBD amplifications. Overview of the SARS-COV-2 Spike RBD
 160 lineages identified from the A. NY3 and B. NY14 sewershed. Amino acid changes similar
 161 to or identical to changes in Omicron (BA.1) are indicated. The lineage previously referred
 162 to as WNY4 in [8] is indicated.

163 **Table 1**

Location	Date range when lineages appeared	Days within range	Number of samples	Number of RBD mutations
NY2	8/16/21-02/28/22	170	10	4-18
NY3	1/31/21 [8] -3/14/22	437	7	16-24
NY10	4/4/21-11/29/21	239	22	4-11
NY11	4/19/21-11/22/22	217	20	4-9
NY13	10/26/21-2/14/22	111	5	12-15
NY14	5/10/21-10/18/21	161	9	8-15
MO33	3/15/21-4/27/21	43	12	4-6
MO45	6/8/21-2/22/22	259	3	4-5
CA	11/4/21-12/21/21	47	3	16

164

165

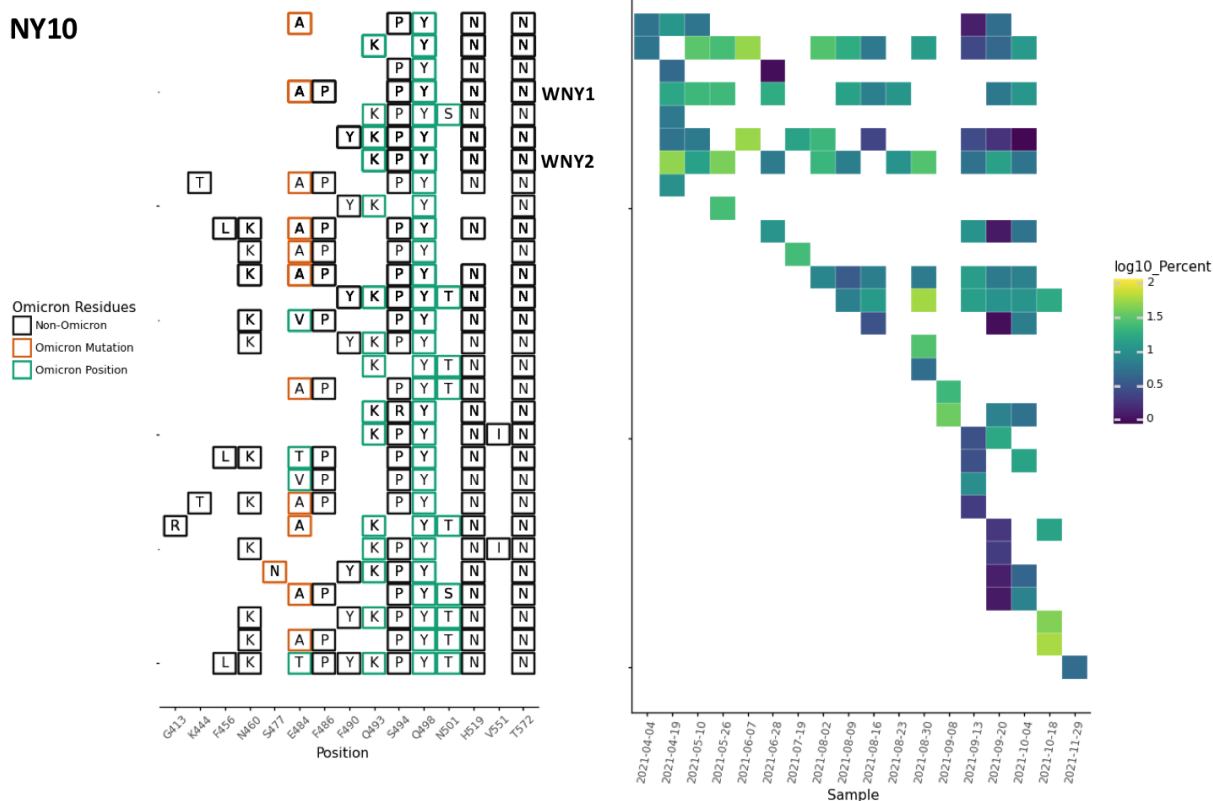
166 Each sewer shed had its own unique set of lineages, but these lineages were not static.

167 For instance, in NY10, the lineages first detected in April 2021 contained 4-5 RBD amino

168 acid changes, but by October and November the lineages contained 8-11 RBD amino

169 acid changes (Fig. 3, Table 1).

NY10

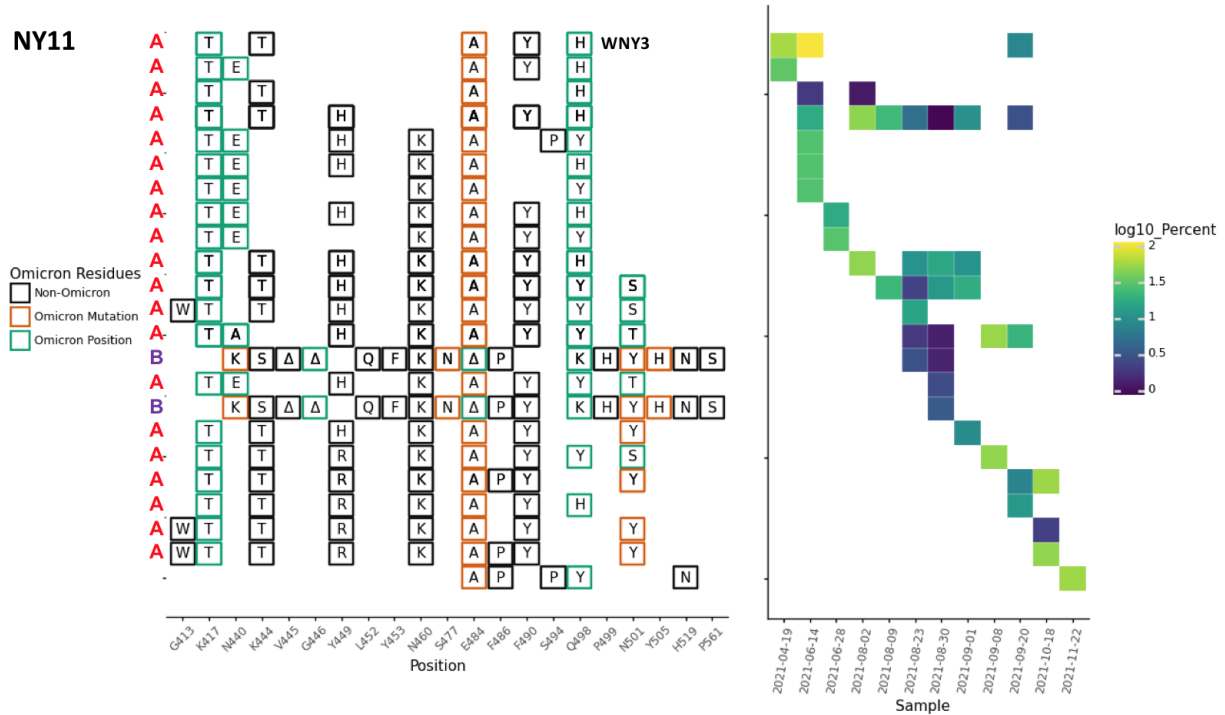


170

171 **Fig. 3** NY10 RBD amplifications. Overview of the SARS-COV-2 Spike RBD lineages
 172 identified from the NY10 sewershed. Amino acid changes similar to or identical to
 173 changes in Omicron (BA.1) are indicated. Two lineages previously referred to as WNY1
 174 and WNY2 in [8] are indicated.

175 In some cases, the sewersheds contained more than one lineage class. For instance, the
 176 NY11 sewershed contained several closely related lineages (class A) starting in April
 177 2021, but a new set of lineages (class B) were detected starting in August 2021. These

178 two classes were clearly distinct with very few amino acid changes in common (Fig. 4)

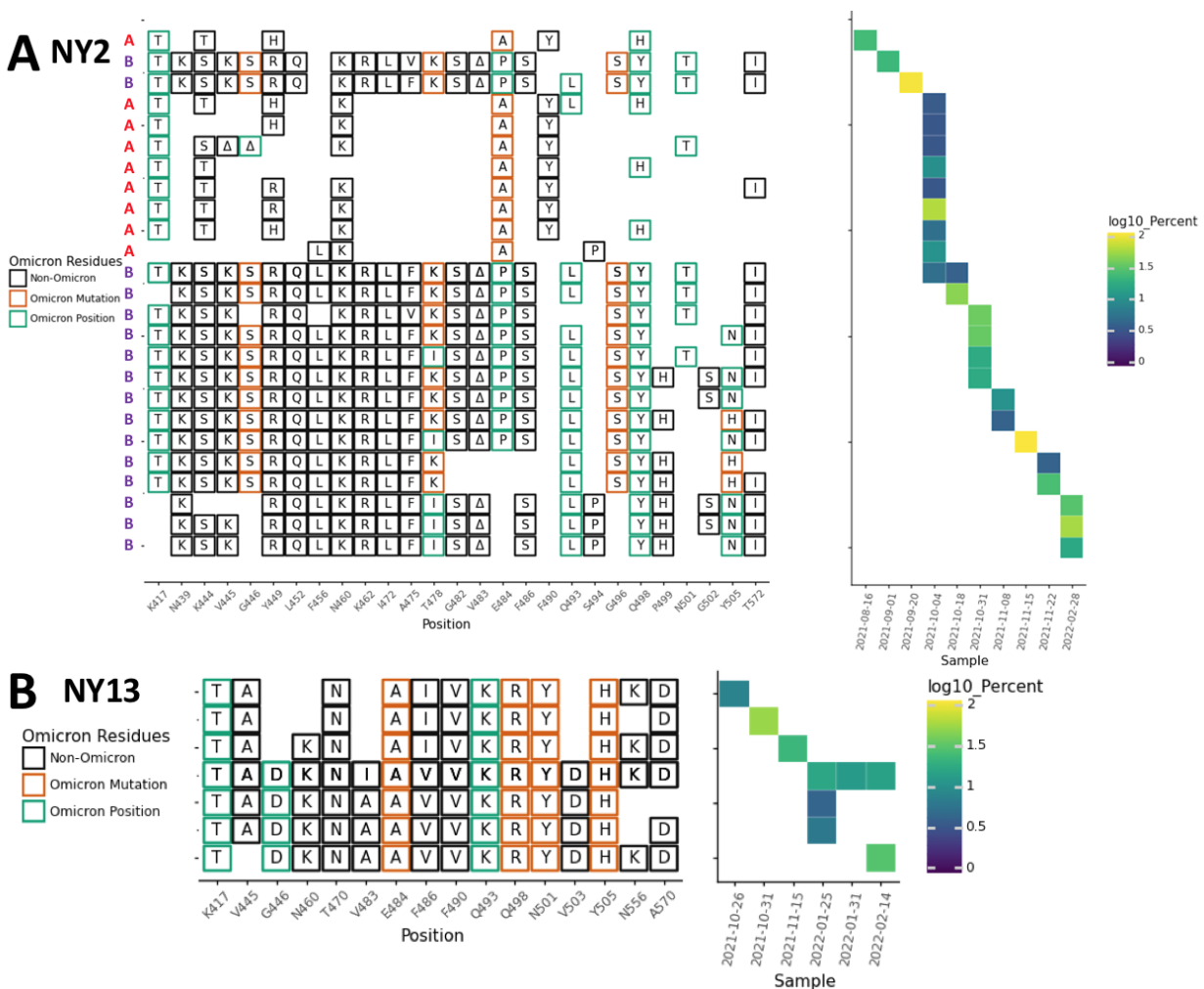


179

180 **Fig. 4** NY11 RBD amplifications. Overview of the SARS-COV-2 Spike RBD lineages
 181 identified from the NY11 sewershed. Lineages designated A and B belong to two lineages
 182 groups that appear unrelated. Amino acid changes similar to or identical to changes in
 183 Omicron (BA.1) are indicated. The lineage previously referred to as WNY3 in [8] is
 184 indicated.

185 In addition to amino acid changes, several of the lineages observed in these sewersheds
 186 contained amino acid deletions near positions 445 and 484. For instance, lineages from
 187 NY2 contained a 444-445 deletion, NY11 and NY14 contained 443-444 deletion, NY3
 188 and NY11 contained a deletion at position 484, and NY2 contained a deletion at position
 189 483 (Fig. 2B, 4, 5A).

190 Most cryptic lineages detected did not appear to be derived from any known VOCs. The
 191 one exception was a lineage class containing amino acid changes N501Y and A570D in
 192 NY13 that first appeared on September 26, 2021, which suggested derivation from the
 193 Alpha VOC (Fig. 5B; Table 1). The Alpha VOC had been the dominant lineage in NYC
 194 between April and June 2021, but by September 26, 2021, it had been supplanted by
 195 Delta VOC and was no longer being detected in NYC [26].



196
 197 **Fig. 5** NY2 and NY13 RBD amplifications. Overview of the SARS-COV-2 Spike RBD
 198 lineages identified from the A. NY3 and B. NY13 sewershed. Lineages designated A and

199 B belong to two lineages groups that appear unrelated. Amino acid changes similar to or
200 identical to changes in Omicron (BA.1) are indicated.

201 Overall, specific lineage classes persisted within, but did not spread beyond, their
202 individual sewersheds, with one notable exception. A cryptic lineage detected on August
203 16, 2021 in NYC sewershed NY2 precisely matched a lineage detected in sewershed
204 NY11 between June-September 2021 (Fig. 4, 5A). While this precise lineage was never
205 seen in NY2 again, several lineages with similar constellations of amino acid changes
206 appeared in NY2 after October 4, 2021. The NY11 and NY2 sewersheds do not border
207 each other, but are not separated by any bodies of water.

208 2.2 Rare and concerning amino acid changes are common in 209 cryptic lineages and are sometimes shared with Omicron

210 In November 2021, the Omicron VOC was first detected in South Africa. This VOC
211 contained eleven changes in the Spike protein between amino acids 410-510. Of these
212 eleven amino acid changes, four (K417T, S477N, T478K, and N501Y) were present in
213 previous VOCs. The remaining seven amino acid changes were rare prior to the Omicron
214 VOC. All seven of these new amino acid changes had been detected in at least one of
215 the wastewater lineages: N440K (MO33, NY11), G446S (NY2), E484A (MO45, NY10,
216 NY11, NY2, NY13, CA), Q493R (NY3, NY14), G496S (NY2), Q498R (NY14, NY13), and
217 Y505H (NY11, NY3, NY14, NY2, NY13, CA) (Figs. 1-6). None of the wastewater lineages
218 have combinations of amino acid changes consistent with having a common ancestor
219 with Omicron and most were initially detected prior to the emergence of Omicron.

220 However, these shared amino acid changes suggest that the cryptic lineages were under
221 selective pressures similar to those that shaped the Omicron lineage.

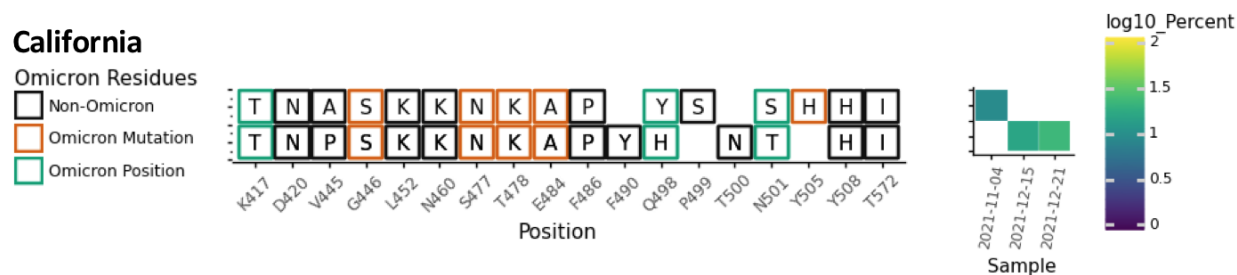
222 Although each sewershed with cryptic lineages had its own signature combinations of
223 amino acid changes, many of these changes were recurring among multiple sewersheds.
224 Some of the more striking examples are described below.

225 *N460K*. All nine of the sewersheds contained lineages with this change. Changes at this
226 position are known to lead to evasion of class I neutralizing antibodies [27,28]. However,
227 this amino acid change is very rare, appearing in less than 0.01% of sequences in GISAID
228 [22–24] submitted by March 15, 2022 (Table S1).

229 *K417T*. Eight of the nine sewersheds contained lineages with the amino acid change
230 *K417T*. Changes at this position are common and are known to participate in evasion
231 from class I neutralizing antibodies [27,28]. Although *K417T* was present in the Gamma
232 VOC, *K417N* is the more common amino acid change at this position. The *K417N* amino
233 acid change was not observed in any of the wastewater cryptic lineages.

234 *N501S/T*. The amino acid changes *N501S* and *N501T* were seen in four and seven of the
235 nine sewersheds, respectively. Changes at this position directly affect receptor binding
236 and can affect the binding of multiple classes of neutralizing antibodies [19,29,30].
237 Although mutations at this position are very common, the most common change by far is
238 *N501Y*, which was present in multiple VOCs. By contrast, *N501S* and *N501T* were
239 present in less than 0.01% and 0.1% of sequences in GISAID submitted by March 15,
240 2022 (Table S1).

241 Q498H/Y. Six of the nine sewersheds in this study contained lineages with the amino
 242 acid change Q498H or Q498Y. It should be noted that Q498Y differs from the Wuhan
 243 ancestral sequence by two nucleotide substitutions at the 498th codon (CAA→TAC).
 244 Q498H (CAA→CAC) is a necessary intermediary in this transition as TAA encodes a stop
 245 codon. In several cases both Q498H and Q498Y were seen in association with particular
 246 lineage classes including in NY2, 11, and 14 (Fig. 2B, 4, 5A) as well as a lineage class
 247 from California detected by the University of California, Berkeley wastewater monitoring
 248 laboratory (COVID-WEB) (Fig. 6). Changes at this position directly affect receptor binding
 249 and can affect the binding of multiple classes of neutralizing antibodies [19,29,30].
 250 Notably, Q498H and Q498Y have been associated with mouse adapted SARS-CoV-2
 251 lineages [31–33]. Both of these amino acid changes are very rare, appearing in less than
 252 0.01% of sequences in GISAID submitted by March 15, 2022. Prior to November 2021,
 253 Q498Y had never been seen in a patient sample (Table S1).



254
 255 **Fig. 6** Overview of the SARS-COV-2 Spike RBD lineages identified from the California
 256 sewersheds. Amino acid changes similar to or identical to changes in Omicron (BA.1) are
 257 indicated.

258 *E484A*. Six of the nine sewersheds contained lineages with the amino acid change
259 *E484A*. Changes at this position are known to participate in evasion from class II
260 neutralizing antibodies [27,28]. Prior to the emergence of Omicron in November 2021,
261 *E484A* was present in about 0.01% of sequences submitted to GISAID (Table S1).

262 *Q493K*. Five of the nine sewersheds contained lineages with the amino acid change
263 *Q493K*. Changes at this position directly affect receptor binding and can affect the binding
264 of multiple classes of neutralizing antibodies [19,27–30,34]. This amino acid change is
265 biophysically very similar to the *Q493R* mutation in Omicron. However, the *Q493K* amino
266 acid change is very rare in patient derived sequences, appearing in less than 0.01% of
267 sequences in GISAID submitted by March 15, 2022 (Table S1).

268 *Y505H*. Five of the nine sewersheds contained lineages with the amino acid change
269 *Y505H*. Prior to the emergence of Omicron in November 2021, *Y505H* was present in
270 about 0.01% of sequences submitted to GISAID (Table S1).

271 *K444T and K445A*. The amino acid changes *K444T* and *K445A* were each seen in four
272 of the nine sewersheds. Changes at these positions are known to participate in evasion
273 from class III neutralizing antibodies [28]. However, these amino acid changes are very
274 rare, each appearing in less than 0.01% of sequences in GISAID submitted by March 15,
275 2022 (Table S1).

276 Y449R. Three of the nine sewersheds contained lineages with the amino acid change
277 Y449R. This change is noteworthy because, as of March 15, 2022, no sequences with
278 this amino acid change had been submitted to GISAID (Table S1).

279 2.3 Long-read sequencing of S1 identifies substantial NTD 280 modifications and suggests high dN/dS ratio

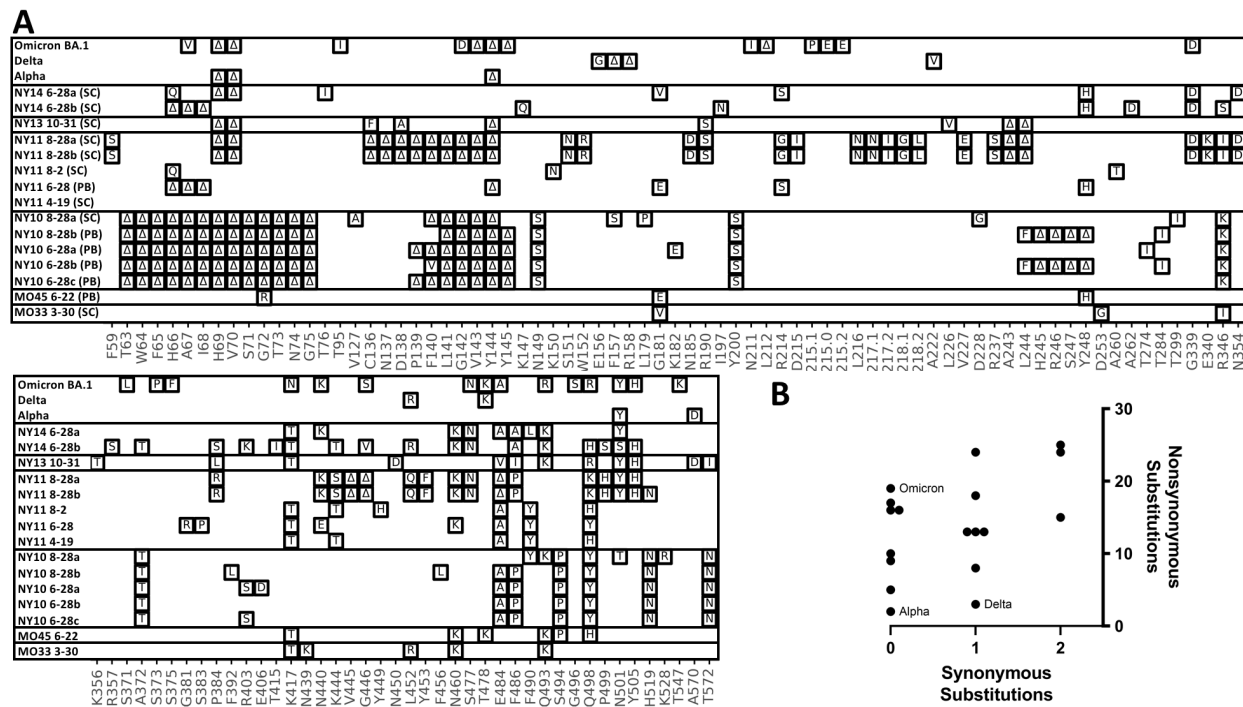
281 With each sample that contained novel cryptic lineages, attempts were made to amplify
282 a larger fragment of the S1 domain of Spike. Amplification of larger fragments from
283 wastewater is often inefficient, but sometimes can be achieved. To gain more information
284 about the S1 domain of Spike and independently confirm the authenticity of the RBD
285 lineages, we optimized a PCR strategy that amplifies 1.6 kb of the SARS-COV-2 Spike
286 encompassing amino acids 57-579. These fragments were then either subcloned and
287 sequenced or directly sequenced using Pacific Biosciences HiFi sequencing (Fig. 7A).

288 The S1 amplification from the MO33 and MO45 sewersheds contained the RBD amino
289 acid changes previously seen and each contained 3 additional amino acid changes
290 upstream from the region sequenced using the targeted amplicon strategy described
291 above (Fig. 7A). Many of the S1 amplifications from the NY10, NY11, NY13 and NY14
292 sewersheds contained numerous changes in S1 (Fig. 7A). In particular, many of the
293 sequences contained deletions near amino acid positions 63-75, 144, and 245-248. All
294 three of these areas are unstructured regions of the SARS-COV-2 spike where deletions
295 have been commonly observed in sequences obtained from patients [35]. Two distinct S1
296 sequences were detected from the NY14 sample collected on June 28, 2021.

297 Interestingly, the first sequence contained 13 amino acid changes which matched the
298 RBD sequences from the same sewershed. The second sequence did not match any
299 lineage that had been seen before, though it contained several mutations that were
300 commonly seen in other cryptic lineages (see section 2.2). This second sequence
301 presumably represented a unique lineage that had not been detected by routine
302 surveillance.

303 A single S1 sequence was obtained from the NY13 samples collected on October 31,
304 2021. This sequence generally matched the RBD sequence from the same date, but did
305 contain minor variations. Importantly, the S1 sequence contained deletions at positions
306 69-70 and 144, which, along with the amino acid changes N501Y and A570D, match the
307 changes found in the Alpha VOC lineage. This information is consistent with the NY13
308 lineages being derived from the Alpha VOC.

309 Comparing the number of non-synonymous to synonymous mutations in a sequence can
310 elucidate the strength of positive selection imposed on a sequence. The ratios of non-
311 synonymous and synonymous mutations in this region of S1 from the Alpha, Delta, and
312 Omicron VOCs (BA.1) were 19/0, 2/0, and 4/1, respectively. It was not possible to
313 calculate the formal dN/dS ratios since many of the sequences did not have synonymous
314 mutations in this region, so instead the numbers of non-synonymous and synonymous
315 mutations were plotted. The cryptic lineages contained 5 to 25 total non-synonymous
316 mutations and 0 to 2 total synonymous mutations (Fig. 7B).



317

318 **Fig. 7** S1 amplifications. A. Overview of the SARS-COV-2 Spike S1 lineages in the Alpha,
 319 Delta, Omicron VOCs and six of the sewersheds with cryptic lineages. S1 amplifications
 320 were sequenced by subcloning (SC) and Sanger sequencing, or were sequenced using
 321 a PacBio (PB) deep sequencing. B. Plot of the number of synonymous and non-
 322 synonymous changes in the S1 sequences shown.

323 **2.4 Cryptic lineages from NCBI suggest an early common ancestor**
 324 **for many of the NYC lineages**

325 In addition to RBD amplicon sequencing performed in our laboratories, we downloaded
 326 the 5609 SARS-CoV-2 wastewater fastq files from NCBI's Sequence Read Archive (SRA)
 327 that were publicly available on NCBI on January 21, 2022 (not including submissions from
 328 our own groups). We screened these sequences for cryptic lineages by searching for

329 recurring amino acid changes seen via RBD amplicon sequencing (K444T, Y449R,
330 N460K, E484A, F486P, Q493K, Q493R, Q498H, Q498Y, N501S, N501T, and Y505H)
331 (see above and Table S1), requiring at least two of these mutations with a depth of at
332 least 4 reads. This strategy identified samples from 15 sewersheds (Table 2). Four were
333 collected from unknown sewersheds in New Jersey and California in January 2021. The
334 other 11 were collected by the company Biobot from NYC between June and August
335 2021. All but one of the lineages closely matched the cryptic sequences that had been
336 observed via RBD amplicon sequencing from the same sewershed. The one exception
337 was SRR16038150, which contained 4 amino acid changes that had not been seen in
338 any of the previous sewershed samples in the same combination. The Biobot sequences
339 were 40-96% complete and appeared to contain 30-100% cryptic lineages based on the
340 frequency of mutation A23056C (Q498H/Y), a mutation shared with the lineages in all 11
341 sewershed samples from NYC. We speculate that the relative abundance of cryptic
342 lineages was high because, during this period, NYC experienced the lowest levels of
343 COVID-19 infections seen since the start of the pandemic. As a result, the sequences
344 that matched the known circulating lineage were at low abundance.

345 To compare the mutational profile among these different NYC samples, we first
346 determined all of the mutations that occurred in at least 3 of the 11 cryptic lineages. We
347 then produced a heat map to compare the frequency of each of these mutations from
348 wastewater samples with the mutations that were reported from New York patient
349 samples in June 2021 (Fig. 8). Surprisingly, the sewershed sequences often lacked two
350 of the four consensus sequences that define the B.1 PANGO lineage (GISAID G clades
351 or Nextstrain '20' clades) of SARS-COV-2 [36]. Almost all patient samples collected in

352 NYC during June 2020 contained the mutations C241T, C3037T, C14408T, and
353 A23403G. The cryptic lineages from NYC wastewater all appeared to contain the
354 mutations C3037T and A23403G, but possessed the ancestral sequences at positions
355 241 and 14408. In addition, there were two mutations in the S gene that were found in
356 nearly all of the cryptic lineages, A23056C (Q498H/Y) and C24044T (L828F). Both of
357 these mutations were found in less than 1% of patient samples. There were 3 additional
358 mutations outside of the S gene that were highly prevalent in most of the wastewater
359 samples, but essentially absent from patient samples: C25936G (Orf3 H182D), G25947C
360 (Orf3 Q185H), and T27322C (Orf6 S41P). While other mutations were detected
361 repeatedly within a sewershed, no other mutations spanned multiple sewersheds.

362 **Table 2. Cryptic lineage whole genome sequences from nationwide surveys.** ^an/a = not available; ND = none
 363 designated

SRA Accession	State	Submitter	Sample Date	Percent cryptic lineage	Genome coverage	Sewershed	PANGO assignment	RBD Changes
SRR17120725	CA	Aquavitas	2021-01-04	7%	27,403	n/a ^a	ND ^b	E484A/Q498H/H519N
SRR16638981	NJ	Aquavitas	2021-01-18	7%	28,185	n/a ^a	ND ^b	E484A/Q498H/H519N
SRR16542155	NJ	Aquavitas	2021-01-18	7%	27,295	n/a ^a	ND ^b	E484A/Q498H/H519N
SRR16362183	NJ	Aquavitas	2021-01-04	100%	15,217	n/a ^a	ND ^b	E484A/Q498H/H519N
SRR16038150	NY	Biobot Analytics	2021-08-17	79%	28,227	NY2	B.1.503	Y449P/E484A/F490Y/Q498H
SRR16038156	NY	Biobot Analytics	2021-08-09	92%	24,595	NY11	B.1.503	K417T/K444T/Y449H/N460K/E484A/F490Y/Q498H
SRR15706711	NY	Biobot Analytics	2021-08-09	100%	11,877	NY11	ND ^b	K417T/K444T/Y449H/N460K/E484A/F490Y/Q498H/A570D
SRR15384049	NY	Biobot Analytics	2021-07-12	99%	24,001	NY10	B.1	Q493K/Q498Y/H519N/T572N)
SRR15291305	NY	Biobot Analytics	2021-07-05	100%	22,316	NY11	P.1.15	K417T/K444T/Y449H/E484A/F490Y/Q498H

SRR15291304	NY	Biobot Analytics	2021-07-04	100%	28,634	NY10	B.1	Q493K/Q498/H519N/T572N
SRR15202285	NY	Biobot Analytics	2021-06-28	100%	12,209	NY2	ND ^p	K444S/V445K/G446V/Y449R/L452Q/N460K/K462R/S477N/T478E/T478R/ DEL483/E484P/F486I/F490P/G496S/Q498Y/P499S/N501T/Y505H/V511I
SRR15202284	NY	Biobot Analytics	2021-06-28	98%	16,281	NY14	ND ^p	K417T/K444S/DEL445- 6/L452R/N460K/S477D/F486V/Q493K/Q498Y/P499S/N501T
SRR15202279	NY	Biobot Analytics	2021-06-28	30%	21,974	NY11	B.1	N440K/K444S/DEL445- 6/L452Q/Y453F/N460K/S477N/D484/F486A/Q493K/Q498K/P499S/N501Y/ H519N
SRR15128983	NY	Biobot Analytics	2021-06-16	99%	21,152	NY11	A.29	K444T/Y449H/E484A/Y489Y/F490Y/Q498H
SRR15128978	NY	Biobot Analytics	2021-06-16	100%	15,593	NY10	ND ^p	E484A/F486P/S494/Q498Y/H519N

Changes in genome			Patient Seqs 6/2021	15291304	15384049	15128978	15291305	15706711	16038156	15128983	15202279	16038150	15202285	15202284
mutation	gene	AA change	New York	NY10	NY10	NY10	NY11	NY11	NY11	NY11	NY11	NY2	NY2	NY14
C00241T	5'UTR	-	99%	19%	0%		0%		63%	100%		63%	0%	
C14,408T	Orf1b	P314L	99%	26%	0%	0%	0%	0%	15%	0%	28%	45%		0%
C3,037T	Orf1a	silent	99%	100%	99%				100%			100%	100%	
A23,403G	S	D614G	100%	100%	100%	100%	100%	80%	100%	100%	99%	100%		
C1,059T	Orf1a	T265I	8%	100%	70%	100%		91%	67%	100%	48%	100%	0%	
A5,648C	Orf1a	K1795Q	10%	30%	99%		100%		99%	100%	0%	56%		
A23,056C	S	Q498X	0%	100%	99%	100%	100%	100%	92%	99%	0%	79%	0%	98%
C24,044T	S	L828F	0%	72%	99%				62%	64%	61%	85%		0%
G25,563T	Orf3a	Q57H	19%	100%	27%		53%	70%	81%	100%	99%	79%		39%
C25,936G	Orf3a	H182D	0%	75%	100%	18%	100%		70%	100%	30%	83%	0%	100%
G25,947C	Orf3a	Q185H	0%	63%	0%	0%	100%		75%	24%	11%	70%	0%	0%
T27,322C	Orf6	S41P	0%	42%	66%	100%	0%	90%	48%	0%	0%	41%		0%
C1,616A	Orf1a	L451I	0%	100%	61%	96%	0%	0%	0%	0%		0%	100%	
C3,267T	Orf1a	T1001I	32%	100%		100%			0%			0%		
G3,849T	Orf1a	S1195I	0%	81%	99%		0%	0%	0%			0%		0%
A4,178C	Orf1a	K1305Q	0%	60%	100%	0%		0%	0%	0%	1%	0%		
C5,178A	Orf1a	T1638N	0%	87%	53%		0%			0%	0%	0%	0%	
A6,328G	Orf1a	silent	0%	100%	99%	100%		0%	0%	0%	0%	0%		0%
T8,296C	Orf1a	silent	0%	98%	100%		0%			0%	0%	0%		0%
G22,599A	S	R346K	10%	100%	72%	51%	0%	0%	0%	0%		0%		0%
C23,039A	S	Q493K	0%	89%	96%	0%	0%	0%	0%	0%	28%	0%	0%	100%
C23,054T	S	Q498X	0%	100%	99%	100%	0%	0%	0%	1%	0%	0%	100%	99%
C23,117A	S	H519N	0%	100%	99%	99%	0%	0%	0%	1%	53%	0%	0%	0%
C23,277A	S	T572N	0%	100%	99%	52%	0%		0%	0%	0%	0%		
T23,406C	S	V615A	0%	52%	60%	50%	0%	0%	0%	0%	0%	0%		
G25,019A	S	D1153N	0%	75%	99%	100%	0%	0%	0%	0%	0%	0%	0%	0%
G25,116A	S	R1185H	0%	63%	100%	100%	0%		0%	0%	0%	0%	0%	52%
C28,887T	S	T205I	12%	89%	77%	71%	0%	0%	0%	0%	0%	9%	100%	54%
C4,113T	Orf1a	A1283V	0%	0%	0%	0%	100%	98%	74%	49%	99%	49%		0%
C4,230T	Orf1a	T1322I	0%	0%	0%	0%		96%	73%	50%	99%	0%		
T5,507G	Orf1a	L1748V	0%	0%			100%		89%	100%		66%		
A9,204G	Orf1a	D2980G	0%	41%	0%		100%		81%	100%	0%	100%		100%
G9,479T	Orf1a	G3072C	0%	0%			0%		96%	99%	99%	25%	0%	
C9,711T	Orf1a	S3149F	0%	0%	0%	0%	100%	50%	59%	100%	45%	81%	0%	0%
T9,982C	Orf1a	silent	0%	0%	0%		99%		50%		99%	100%	0%	
G11,670A	Orf1a	R3802H	0%	0%	0%	0%	91%	0%	92%		94%	76%	0%	0%
C11,916T	Orf1a	S3884L	0%	0%	0%	0%	100%	24%	100%	100%	42%	74%	0%	0%
G17,196A	Orf1b	silent	0%	0%	0%	0%	100%	70%	82%	100%	65%	56%	0%	0%
A17,496C	Orf1b	E1343D	0%	0%	0%		100%	89%	35%	100%		57%	0%	
T18,660C	Orf1b	silent	0%	0%	0%	0%	100%	87%	94%	93%	0%	63%	0%	0%
G22,340A	S	A260T	0%	0%		0%	93%	84%	100%			41%		0%
A22,893C	S	K444T	0%	0%	0%	0%	98%	100%	66%	11%	0%	62%	0%	0%
T22,907C	S	Y449H	0%	0%	0%	0%	99%	100%	86%	99%	0%	57%	100%	0%
A23,013C	S	E484A	0%	9%	3%	100%	100%	100%	99%	99%	0%	75%	100%	0%
C23,029T	S	silent	0%	0%	0%	0%	99%	20%	37%	98%	0%	0%	0%	99%
T23,031A	S	F490Y	0%	0%	0%	0%	99%	100%	100%	98%	0%	79%	4%	0%
C24,418T	S	silent	0%	0%	0%	0%	100%		72%	100%	0%	63%	0%	0%
A25,020C	S	D1153A	0%	0%	1%	0%	100%	6%	75%	100%	0%	72%	3%	14%
T25,570A	Orf3a	S60T	0%	0%	0%		0%	71%	54%	56%	1%	0%		0%
A27,330C	Orf6	silent	0%	0%	0%	0%	100%	94%	74%	99%	99%	66%		0%
T27,384C	Orf6	silent	1%	0%	0%	0%	100%	35%	41%	100%	68%	58%		0%
T27,907G	Orf8	V5G	0%	0%	0%		100%		65%		99%	100%		0%
C27,920T	Orf8	silent	0%	15%	0%		67%		27%		79%	61%		0%
T27,929A	Orf8	silent	0%	0%	0%		68%		28%		79%	61%		0%
A28,271T	UTR	-	1%	0%	1%	0%	49%		47%	0%	27%	56%	0%	0%
G29,540A	UTR	-	0%	0%	0%	0%	99%	84%	71%	88%	78%	65%	0%	0%

366 **Fig. 8** Polymorphisms from wastewater genomes. Shown are all mutations present in at
367 least three of the whole genome sequences from NYC listed in Table 2 and their
368 corresponding amino acid changes. First column lists the prevalence of each mutation
369 among all patients samples collected in June 2021 from New York. Each other column
370 lists the prevalence of each mutation in each of the genome sequences.

371 To confirm that some of the cryptic lineages lacked the B.1 lineage consensus mutations,
372 we designed primers to amplify and sequence the C14408 region of SARS-CoV-2 RNA
373 isolated from wastewater. Indeed, samples from NY11 and NY10 that had a high
374 prevalence of cryptic lineages were found to contain sequences that lacked C14408T
375 (Fig. S1). However, when samples were amplified from the NY13 sewershed when the
376 cryptic lineages there were present, we observed only the modern C14408T, as would be
377 expected if the NY13 lineage were derived from the Alpha VOC. In addition, we performed
378 whole genome sequencing on a March 30, 2021 sample from MO33 when the cryptic
379 lineages were highly prevalent and did not detect any sequence that lacked C241T or
380 C14408T, suggesting the cryptic lineages in this sewershed diverged after the emergence
381 of the B.1 lineage (Fig. S2). Finally, we also analyzed the sequences from NCBI that
382 contained the cryptic lineages from NJ and CA and did not find any sequences lacking
383 C241T or C14408T. Thus, the lineages lacking C241T and C14408T appear to be limited
384 to a subset of the cryptic lineages from NYC. It would appear that a SARS-CoV-2 lineage
385 bearing mutations C3037T and A23403G, but possessing the ancestral genotype at
386 positions 241 and 14408, was the direct ancestor of most of the cryptic lineages found in
387 NYC.

388 3. Discussion

389 Our results point to the evolution of numerous SARS-CoV-2 lineages under positive
390 immune selection whose source/host remains unknown.

391 3.1 Relatedness of and origin of cryptic lineages

392 We previously detected cryptic lineages via targeted amplicon sequencing [8], but lacked
393 information about their derivation. Here, from comparison of the sewersheds for which
394 whole genome sequencing is available, it is clear that the cryptic lineages from
395 wastewater are not all derived from a common ancestor. The NY13 lineage appeared to
396 be derived from the Alpha VOC. If this is true, the NY13 lineage most likely branched off
397 from Alpha sometime in early to mid-2021 when that variant was common in NYC.
398 However, many lineages from the NY10, NY11, NY2, and NY14 sewersheds in New York
399 appear to share a common ancestor that branched off from a pre-B.1 lineage.
400 Additionally, we often observed swarms of related sequences that co-occurred within a
401 sewershed on a single date, and accumulated new mutations over time, suggesting
402 continued diversification from a single origin within each sewershed.

403 3.2 Comparison with the Omicron VOC

404 The Omicron VOC and the wastewater lineages appear to have been subjected to high
405 positive selection. While prior VOCs had 3 or fewer amino acid changes in the amplified
406 region of the RBD, the Omicron VOC (BA.1) contained 11 and the cryptic lineages from

407 wastewater averaged over 10. By comparison, a cluster of SARS-COV-2 sequences that
408 appear to have circulated in white-tailed deer for over a year accumulated only 2 amino
409 acid changes in this region [37]. Of the nonsynonymous RBD mutations in Omicron, four
410 were in at least one prior VOC: K417N, S477N, T478K, and N501Y. The other seven
411 were relatively rare; N440K was present in 0.2% of sequences and the other six were
412 each present in less than 0.1% of sequences in GISAID prior to November 1, 2021. All of
413 the rare Omicron changes were observed in at least one of the cryptic wastewater
414 lineages. Collectively, this suggests that the wastewater lineages and the Omicron VOC
415 likely arose under similar selective pressures. The high dN/dS ratios found in cryptic
416 lineages and in Omicron suggest that these selective pressures must be exceptionally
417 strong.

418 3.3 Source of Lineages

419 In spite of detailed tracking and cataloging of the cryptic lineages, the question where
420 they are coming from remains unanswered. The most parsimonious explanations are 1)
421 undetected spread within the human population, 2) prolonged shedding by individuals,
422 most likely immunosuppressed, or 3) spread in animal reservoirs.

423 Undetected spread in the population appears unlikely. While the sequencing rate for US
424 patient samples is not 100%, it is high enough that population-level spread of cryptic
425 lineages would not be missed. Alternatively, as it is known that SARS-CoV-2 can replicate
426 in gastrointestinal sites [38,39], the lack of detection of cryptic lineages by clinical
427 sequencing could be explained by the potential adaptation of some SARS-CoV-2 to

428 replicate exclusively in the gastrointestinal tract [1,38]. Nonetheless, even if replication of
429 these lineages were occurring outside of the nasopharyngeal region, this could not
430 explain why cryptic lineages generally remain geographically constrained.

431 The simplest explanation for the appearance of cryptic lineages in wastewater is that they
432 are shed by immunosuppressed patients with persistent infections. Indeed, the vast
433 majority of amino acid changes in the RBD of the Omicron VOC and the cryptic lineages
434 confer resistance to neutralizing antibodies. In particular, substitutions at positions 417,
435 440, 460, 484, 493, 498 and 501 have all been well documented to lead to immune
436 evasion [17,27,34,40–42]. Additionally, RBD changes K417T, N440K, N460K, E484A,
437 Q493K, and N501Y have all been observed in persistent infections of
438 immunocompromised patients [43,44]. Given the repeated appearance of these
439 mutations in diverse sewersheds, the majority of the selective pressure on the cryptic
440 lineages is almost certainly immune pressure. The counterargument to this explanation
441 is the sheer volume of viral shedding required to account for the wastewater signal. Many
442 of the sewersheds process 50-100 million gallons of wastewater per day. Reliable
443 amplification of a sequence from wastewater generally requires that the sequence is
444 present at least 10,000 copies per liter. Therefore, detection of a specific virus lineage in
445 such a sewershed would seem to require several trillion virus particles to be deposited
446 each day. If this signal were derived from a single infected patient or even a small group
447 of patients, those patients would have to shed exponentially more viruses than typical
448 COVID-19 patients.

449 The final explanation for the cryptic lineages in wastewater is that they are shed into
450 wastewater by an animal host population. Previously, we determined through rRNA
451 analysis of several NYC sewersheds that the major non-human mammals that contribute
452 to the wastewater are cats, rats, and dogs [8]. Of these three, rats were the only species
453 that seemed to be a plausible candidate. Indeed, we also showed that the cryptic lineages
454 from the sewersheds had the ability to utilize rat and mouse ACE2 [8]. However, one of
455 the sewersheds with the most consistent signal in 2021 was NY10, which had little to no
456 rat rRNA. In addition, it is not clear why circulation in an immune competent animal such
457 as a dog or a rat would result in a more rapid selection of immune escape mutations than
458 circulation in humans, yet the cryptic lineages display accumulation of many times more
459 immune escape changes than seen in viruses circulating in the human population.

460 3.4 The importance of wastewater sequencing methodology for 461 identification of novel variants

462 In order to provide information regarding the appearance and spread of SARS-CoV-2
463 variants in communities, next generation sequencing technologies have been applied to
464 sequence SARS-CoV-2 genetic material obtained from sewersheds around the world
465 [45–47]. Commonly, SARS-CoV-2 RNA extracted from wastewater is amplified using
466 SARS-CoV-2 specific primers that cover the entire genome [48–50]. Bioinformatic
467 pipelines are employed to identify circulating SARS-CoV-2 variants [16,51]. In general,
468 the presence and abundance of variants in wastewater corresponds to data obtained from
469 clinical sequencing [45,46]. However, to our knowledge, there have been no other reports
470 of cryptic lineages detected in wastewater that were not also observed in clinical

471 sequence data. A major issue with generating whole genome sequence data from nucleic
472 acid isolated from wastewater is sequence dropout over diagnostically important regions
473 of the genome [48,52,53]. In some cases, diagnostically important regions of the genome
474 that accumulate many mutations, such as the Spike RBD, receive little to no sequence
475 coverage, making variant attribution difficult. Since wastewater contains a mixture of virus
476 lineages and whole genome sequencing relies on sequencing of small genome
477 fragments, mutations appearing on different reads cannot be linked together. Indeed,
478 some variant identification pipelines map reads to reference genomes to estimate the
479 probability that mutations are found in the same genome [16]. Such strategies would not
480 be able to detect variants containing unique constellations of mutations. Detecting novel
481 variants that are present at low relative abundances may be better achieved by targeted
482 amplicon sequencing, such as the strategy we present here.

483 **Summary**

484 Over the past 15 months, cryptic SARS-CoV-2 lineages never seen in human patients
485 have appeared in community wastewater in several locations across the USA [8]. These
486 lineages have persisted, intermittently, often as swarms of closely related haplotypes that
487 acquired additional amino acid changes over time, for up to 14 months. Evidence
488 suggests that some of the lineages may have arisen during the initial phases of the
489 pandemic in early to mid-2020. Significantly, these lineages often contained amino acid
490 changes that have rarely or never appeared in contemporaneous variants, at least until
491 the appearance of the Omicron VOC. Many of these amino acid changes are associated
492 with evasion of antibody-mediated neutralization. Collectively, nonsynonymous

493 substitutions in these lineages overwhelmingly outnumbered synonymous substitutions,
494 indicating that these lineages have undergone exceptionally strong positive selection.

495

496 Three hypotheses for the origins of these lineages have been proposed: 1) undetected
497 transmission, 2) long-term infections of immunocompromised patients and 3) possible
498 animal reservoirs. Although immunosuppressed populations are the simplest explanation,
499 it is difficult to reconcile the magnitude of the signal with individual patients being the
500 source. Regardless of the origins and dynamics of cryptic variant shedding, our results
501 highlight the ability of wastewater-based epidemiology to more completely monitor SARS-
502 CoV-2 transmission and genetic diversity than can patient based sampling, at scale and
503 at a greatly reduced cost. Given that multiple VOCs may have gone undetected until
504 suddenly appearing, highly mutated, in apparently single evolutionary leaps [12], it is
505 crucial to the early detection of the next variant of concern that novel SARS-CoV-2
506 genotypes are monitored for evidence of significant expansion. Importantly, patient
507 sampling efforts, despite occurring with an intensity not seen in any prior epidemic, were
508 unable to identify intermediary forms of most VOCs. Monitoring of wastewater, particularly
509 using a targeted sequencing approach, likely provides the best avenue for detecting
510 developing VOCs.

511 4. Materials and Methods

512 Wastewater sample processing and RNA extraction

513 24-hr composite samples of wastewater were collected weekly from the inflow at each of
514 the wastewater treatment plans.

515 NYC: Samples were processed on the day they were collected and RNA was isolated
516 according to our previously published protocol [6]. Briefly, 250 mL from a 24-hr
517 composite wastewater sample from each WWTP were centrifuged at 5,000 x g for 10
518 min at 4°C to pellet solids. A 40 mL aliquot from the centrifuged samples was passed
519 through a 0.22 µM filter (Millipore). To each corresponding filtrate, 0.9 g sodium chloride
520 and 4.0 g PEG 8000 (Fisher Scientific) were added. The tubes were kept at 4°C for 24
521 hrs and then centrifuged at 12,000 x g for 120 minutes at 4 °C to pellet the precipitate.
522 The pellet was resuspended in 1.5 mL TRIzol (Fisher Scientific), and RNA was purified
523 according to the manufacturer's instructions.

524 MO: Samples were processed as previously described [9]. Briefly, wastewater samples
525 were centrifuged at 3000×g for 10 min and then filtered through a 0.22 µM
526 polyethersulfone membrane (Millipore, Burlington, MA, USA). Approximately 37.5 mL of
527 wastewater was mixed with 12.5 mL solution containing 50% (w/vol) polyethylene glycol
528 8000 and 1.2 M NaCl, mixed, and incubated at 4 °C for at least 1 h. Samples were then
529 centrifuged at 12,000×g for 2 h at 4°C. Supernatant was decanted and RNA was extracted
530 from the remaining pellet (usually not visible) with the QIAamp Viral RNA Mini Kit (Qiagen,
531 Germantown, MD, USA) using the manufacturer's instructions. RNA was extracted in a
532 final volume of 60 µL.

533 CA: Samples were processed as previously described [54]. Briefly, 40 mLs of influent
534 was mixed with 9.35g NaCl and 400 uL of 1M Tris pH 7.2, 100mM EDTA. Solution was

535 filtered through a 5-um PVDF filter and 40 mLs of 70% EtNY11 was added. Mixture was
536 passed through a silica spin column. Columns were washed with 5 mL of wash buffer 1
537 (1.5 M NaCl, 10 mM Tris pH 7.2, 20% EtNY11), and then 10 mL of wash buffer 2 (100
538 mM NaCl, 10 mM Tris pH 7.2, 80% EtNY11). RNA was eluted with 200 ul of ZymoPURE
539 elution buffer.

540 **Targeted PCR: MiSeq sequencing**

541 The primary RBD RT-PCR was performed using the Superscript IV One-Step RT-PCR
542 System (Thermo Fisher Scientific,12594100). Primary RT-PCR amplification was
543 performed as follows: 25 °C (2:00) + 50 °C (20:00) + 95 °C (2:00) + [95 °C (0:15) + 55 °C
544 (0:30) + 72 °C (1:00)] × 25 cycles using the MiSeq primary PCR primers
545 CTGCTTTACTAATGTCTATGCAGATTC and NCCTGATAAAGAACAGCAACCT.
546 Secondary PCR (25 µL) was performed on RBD amplifications using 5 µL of the primary
547 PCR as template with MiSeq nested gene specific primers containing 5' adapter
548 sequences (0.5 µM each)
549 acactcttccctacacgacgctctccgatctGTRATGAAGTCAGMCAAATYGC and
550 gtgactggagttcagacgtgtgctctccgatctATGTCAAGAATCTCAAGTGTCTG, dNTPs (100 µM
551 each) (New England Biolabs, N0447L) and Q5 DNA polymerase (New England Biolabs,
552 M0541S). Secondary PCR amplification was performed as follows: 95 °C (2:00) + [95 °C
553 (0:15) + 55 °C (0:30) + 72 °C (1:00)] × 20 cycles. A tertiary PCR (50 µL) was performed
554 to add adapter sequences required for Illumina cluster generation with forward and
555 reverse primers (0.2 µM each), dNTPs (200 µM each) (New England Biolabs, N0447L)
556 and Phusion High-Fidelity or (KAPA HiFi for CA samples) DNA Polymerase (1U) (New

557 England Biolabs, M0530L). PCR amplification was performed as follows: 98 °C (3:00) +
558 [98 °C (0:15) + 50 °C (0:30) + 72 °C (0:30)] × 7 cycles +72 °C (7:00). Amplified product
559 (10 µl) from each PCR reaction is combined and thoroughly mixed to make a single
560 pool. Pooled amplicons were purified by addition of Axygen AxyPrep MagPCR Clean-up
561 beads (Axygen, MAG-PCR-CL-50) or in a 1.0 ratio to purify final amplicons. The final
562 amplicon library pool was evaluated using the Agilent Fragment Analyzer automated
563 electrophoresis system, quantified using the Qubit HS dsDNA assay (Invitrogen), and
564 diluted according to Illumina's standard protocol. The Illumina MiSeq instrument was used
565 to generate paired-end 300 base pair reads. Adapter sequences were trimmed from
566 output sequences using Cutadapt.

567 **Long PCR and subcloning.**

568 The long RBD RT-PCR was performed using the Superscript IV One-Step RT-PCR
569 System (Thermo Fisher Scientific, 12594100). Primary long RT-PCR amplification was
570 performed as follows: 25 °C (2:00) + 50 °C (20:00) + 95 °C (2:00) + [95 °C (0:15) + 55 °C
571 (0:30) + 72 °C (1:30)] × 25 cycles using primary primers
572 CCCTGCATACACTAATTCTTTCAC and TCCTGATAAAGAACAGCAACCT. Secondary
573 PCR (25 µL) was performed on RBD amplifications using 5 µL of the primary PCR as
574 template with nested primers (0.5 µM each) CATTCAACTCAGGACTTGTTCTT and
575 ATGTCAAGAATCTCAAGTGTCTG, dNTPs (100 µM each) (New England Biolabs,
576 N0447L) and Q5 High-Fidelity DNA Polymerase (New England Biolabs, M0491L).
577 Secondary PCR amplification was performed as follows: 95 °C (2:00) + [95 °C (0:15) +
578 55 °C (0:30) + 72 °C (1:30)] × 20 cycles.

579 Positive amplifications were visualized in an agarose gel stained with ethidium bromide,
580 excised, and purified with a NucleoSpin Gel and PCR Clean-up Kit (Macherey-Nagel,
581 74609.250). Gel purified DNA was subcloned using a Zero Blunt TOPO PCR Cloning Kit
582 (Invitrogen, K2800-20SC). Individual colonies were transferred to capped test tubes
583 containing 10 ml of 2X YT broth (ThermoFisher, BP9743-5). Test tubes were incubated
584 at 37°C and shook at 250 rpm for 24 hours. The resulting *E. Coli* colonies were centrifuged
585 for 10 minutes at 5000 xg and the supernatant was decanted. Plasmid DNA was extracted
586 from the pellet using a GeneJet Plasmid Miniprep Kit (ThermoFisher, K0503). The
587 concentration of plasmid DNA extracts was measured using a NanoDrop One
588 (ThermoFisher, ND-ONE-W).

589 **Subcloning**

590

591 The 1.6kb S1 fragment was cloned into the pMiniT vector using the NEB PCR Cloning Kit
592 (NEB #E1202) protocol. Briefly, 5 ul of the RT-PCRRed fragment was ligated to the
593 linearized pMiniT 2.0 Vector and transformed into NEB 10-beta Competent *E. coli* (NEB
594 #C3019). Transformed cells were outgrown in 950ul of NEB 10-beta/Stable Outgrowth
595 Medium at 37°C and shaken at 200 rpm for 1 hour. Dilutions of 1:10, 1:100, and 1:1000
596 were plated onto 100 ug/mL LB-ampicillin using glass beads. Plates were incubated
597 overnight at 37°C. Following this, single colonies were swatched onto 100ug/mL LB-
598 ampicillin plates and pools of 5 colonies were lysed in 20ul H₂O for 5 minutes at 95°C. 5
599 ul of lysed template was added to the following PCR mix: 10 ul Q5 High-Fidelity 2X Master
600 Mix (NEB #M0492S) + 0.8 ul 5uM CATTCAACTCAGGACTTGGTTCTT forward primer
601 (2402F) + 0.8 ul 5uM ATGTCAAGAATCTCAAGTGTCTG (2376R) reverse primer+ 5 ul

602 H₂O. Thermocycling conditions were as follows: 95°C (2:00) + [95°C (00:15) + 55°C
603 (00:30) + 72°C (3:00)] x 40 cycles + 72°C (1:00). PCR products were run on a 0.8%
604 agarose gel and successfully cloned colonies were determined after ethidium bromide
605 staining. Positive pools were then sent for Genewiz Sanger sequencing with the following
606 primer pairs: 2402F (5uM), 2376R (5uM), and NEB Cloning Kit Cloning Analysis Forward
607 Primer (100uM), Cloning Analysis Reverse Primer (100uM). The following additional
608 primer pairs were included as needed: TGCGAATAATTGCACTTTTGA and
609 TGCTACCGGCCTGATAGATT, GGACCTTGAAGGAAAACAGG and
610 TGCTACCGGCCTGATAGATT, and ATCTCCCTCAGGGTTTTTCG and
611 CCATTACAAGGTGTGCTACCG. Pools with mixed Sanger sequencing signals were
612 resequenced as individual colonies picked from the swatch plates and DNA was isolated
613 using the QIAprep Spin Miniprep Kit (#27106).

614

615 **PacBio sequencing**

616 A nested RT-PCR protocol was used to generate 1.6kb Spike amplicons from wastewater
617 RNAs for PacBio sequencing. The primary RT-PCR amplification was performed with the
618 SuperscriptTM IV One-Step RT-PCR System (Invitrogen) and the same thermal cycling
619 program as described above for MiSeq amplicons. These inter Spike gene-specific primer
620 sequences (5'-[BC10ab]-ATTCAACTCAGGACTTGTTCTT and 5'-[BC10xy]-
621 ATGTCAAGAATCTCAAGTGTCTG) were tagged directly on their 5' ends with standard
622 16 bp PacBio barcode sequences (Supplemental Table Z) and used with asymmetric
623 barcode combinations that allow large numbers of samples to be pooled prior to
624 sequencing. The following thermal cycling profile was used for nested PCR: 98 °C (2 min)

625 + [98 °C (10 sec) + 55 °C(10 sec) + 72 °C (1 min)] x 20 cycles + 72 °C (5 min). The
626 resulting PCR amplicons were then subjected to three rounds of purification with AMPure
627 XP beads (Beckman Coulter Life Sciences) in a ratio of 0.7:1 beads to PCR. Purified
628 amplicons were quantified using a Qubit™ dsDNA HS kit (ThermoFisher Scientific) and
629 pooled prior to PacBio library preparation.

630 After ligation of SMRTbell adaptors according to the manufacturer's protocol, sequencing
631 was completed on a PacBio Sequel II instrument (PacBio, Menlo Park, CA USA) in the
632 Genomic Sequencing Laboratory at the Centers of Disease Control in Atlanta, GA, USA.
633 Raw sequence data was processed using the SMRT Link v10.2 command line toolset
634 ([Software downloads - PacBio](#)). Circular consensus sequences were demultiplexed
635 based on the asymmetric barcode combinations and subjected to PB Amplicon Analysis
636 to obtain high-quality consensus sequences and search for minor sequence variants.

637 **Bioinformatics**

638 MiSeq and PacBio processing

639 Sequencing reads were processed as previously described. Briefly, VSEARCH tools were
640 used to merge paired reads and dereplicate sequences [55]. Dereplicated sequences
641 from RBD amplicons were mapped to the reference sequence of SARS-CoV-2
642 (NC_045512.2) spike ORF using Minimap2 [56]. Mapped amplicon sequences were then
643 processed with SAM Refiner using the same spike sequence as a reference and the
644 command line parameters "--Alpha 1.8 --foldab 0.6" [9].

645 The covariant deconvolution outputs were used to generate the haplotype plots in figures
646 1-7. MiSeq sequences from those outputs were collected by sewershed and multiple runs
647 of the same sample averaged. The collected MiSeq sequences were processed to
648 remove the Alpha, Beta, Gamma, Delta and Omicron VOC lineage sequences and
649 remove individual polymorphisms that appeared only once in sewersheds that had
650 unknown lineages on more than 3 dates. Sequences that appeared in only one sample
651 and had an abundance less than .02 were discarded. The resulting sequences were then
652 plotted along with their date and abundance. The PacBio sequences were similarly
653 collected to generate the haplotype plot in Fig. 7. The covariant outputs, their collection
654 and variant lineage processing outputs, and the scripts for processing them are available
655 upon request. The MO raw sequence reads are available in NCBI's SRA under the
656 BioProject accession PRJNA748354. The NY raw sequence reads are available in
657 NCBI's SRA under the BioProject accession PRJNA715712.

658 **NCBI SRA screening**

659 Raw reads were downloaded and then processed similar to MiSeq sequencing except
660 the reads were mapped to the entire SARS-CoV-2 genome and SAM Refiner was run
661 with the parameters '--wgs 1 --collect 0 --indel 0 --covar 0 --min_count 1 --
662 min_samp_abund 0 --min_col_abund 0 --ntabund 0 --ntcover 1'. Unique sequence
663 outputs from SAM Refiner were then screened for specific amino acid changes with a
664 custom script (available upon request). The nt call outputs of samples of interest were
665 used to determine other variations in the genomes sequenced.

666 **14408 sequencing**

667 The long RBD RT-PCR was performed using the Superscript IV One-Step RT-PCR
668 System (Thermo Fisher Scientific, 12594100). Primary long RT-PCR amplification was
669 performed as follows: 25 °C (2:00) + 50 °C (20:00) + 95 °C (2:00) + [95 °C (0:15) + 55 °C
670 (0:30) + 72 °C (1:30)] × 25 cycles using primary primers ATACAAACCACGCCAGGTAG
671 and AACCTTAGACACAGCAAAGT. Secondary PCR (25 µL) was performed on RBD
672 amplifications using 5 µL of the primary PCR as template with nested primers (0.5 µM
673 each)
674 AACTCTTTCCCTACACGACGCTCTTCCGATCTGGTAGTGGAGTTCCTGTTGTAG
675 and
676 GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTAGCACGTAGTGCGTTTATCT,
677 dNTPs (100 µM each) (New England Biolabs, N0447L) and Q5 High-Fidelity DNA
678 Polymerase (New England Biolabs, M0491L). Secondary PCR amplification was
679 performed as follows: 95 °C (2:00) + [95 °C (0:15) + 55 °C (0:30) + 72 °C (1:30)] × 20
680 cycles.

681 **Whole genome sequencing**

682 Whole genome sequencing of the SARS-CoV-2 genome from the MO33 sewershed was
683 performed using the NEBNext ARTIC SARS-CoV-2 Library Prep Kit (Illumina). Amplicons
684 were sequenced on an Illumina MiSeq instrument. Output sequences were analyzed
685 using the program SAM Refiner [57].

686 **Author Contributions** <https://casrai.org/credit/>

687 Conceptualization: DAG, JJD, MCJ, RSK; Data Curation: DAG, MT, CR, AF, SK, KMS,
688 DL, RWW, DB, KL, SN, KB, M-YZ, EK, NB, JL, JH, C-HL, DHO, RSK, JJD, MCJ; Formal

689 Analysis: DAG, DHO, JL, JJD, MCJ, RSK; Funding Acquisition DHO, JW, RSK, JJD, MCJ:
690 JJD; Investigation: DAG, MT, CR, AF, SK, KMS, DL, RWW, DB, KL, SN, KB, M-YZ, EK,
691 NB, JL, JH, C-HL, DHO, RSK, JJD, MCJ; Methodology: DAG, JJD, DHO; Project
692 Administration: DHO, JL, NB, CW, JW, JJD, MCJ, RSK; Resources: DHO, JJD, MCJ,
693 RSK; Software: DAG, RSK, DHO; Supervision: JL, NB, JW, JJD, MCJ, RSK;
694 Visualization: DAG, JJD, MCJ, RSK; Writing – Original Draft Preparation: DAG, JJD,
695 MCJ, RSK; Writing – Review & Editing: DAG, MT, CW, JW, DHO, RSK, JJD, MCJ

696

697 **Acknowledgments**

698 The authors thank Benjamin Martin-Rambo, Dhvani Batra, Kristine Lacek, Sarah Nobles,
699 and Justin Lee at the Centers for Disease Control and Prevention Genomic Sequencing
700 Lab for assistance with PacBio sequencing. We also thank Thomas Peacock for valuable
701 advice and feedback during the preparation of this manuscript. Thanks to Kristen Cheung,
702 Anna Gao, Nanami Kubota, and Shyanon Rai for experimental assistance.

703 **Funding**

704 This project has been funded in part with federal funds from the NIDA/NIH
705 (<https://www.nida.nih.gov/>) under contract numbers 1U01DA053893-01 to JW and MCJ. This
706 work was supported by grants from the New York City Department of Environmental Protection
707 (<https://www1.nyc.gov>) to JJD. This work was supported by financial support through Rockefeller
708 Regional Accelerator for Genomic Surveillance (<https://www.rockefellerfoundation.org>,133
709 AAJ4558), Wisconsin Department of Health Services Epidemiology and Laboratory Capacity
710 funds (<https://www.dhs.wisconsin.gov>, 144 AAJ8216) to DHO. The work was supported by funds

711 from the California Department of Health (<https://www.dhcs.ca.gov/>). The funders had no role in
712 study design, data collection and analysis, decision to publish, or preparation of the manuscript.

713

714 **Data Availability**

715 Data Availability: Raw data are available under NCBI's SRA under the BioProject
716 accession PRJNA748354. The NY raw sequence reads are available in NCBI's SRA
717 under the BioProject accession PRJNA715712.

718 **Competing Interest**

719 The authors declare no competing interests.

720 **References**

- 721 1. Cheung KS, Hung IFN, Chan PPY, Lung KC, Tso E, Liu R, et al. Gastrointestinal
722 Manifestations of SARS-CoV-2 Infection and Virus Load in Fecal Samples From a Hong
723 Kong Cohort: Systematic Review and Meta-analysis. *Gastroenterology*. 2020;159: 81–95.
724 doi:10.1053/j.gastro.2020.03.065
- 725 2. Parasa S, Desai M, Thoguluva Chandrasekar V, Patel HK, Kennedy KF, Roesch T, et al.
726 Prevalence of Gastrointestinal Symptoms and Fecal Viral Shedding in Patients With
727 Coronavirus Disease 2019: A Systematic Review and Meta-analysis. *JAMA Netw Open*.
728 2020;3: e2011335. doi:10.1001/jamanetworkopen.2020.11335
- 729 3. Ahmed W, Tschärke B, Bertsch PM, Bibby K, Bivins A, Choi P, et al. SARS-CoV-2 RNA

- 730 monitoring in wastewater as a potential early warning system for COVID-19 transmission in
731 the community: A temporal case study. *Sci Total Environ.* 2021;761: 144216.
732 doi:10.1016/j.scitotenv.2020.144216
- 733 4. Gonzalez R, Curtis K, Bivins A, Bibby K, Weir MH, Yetka K, et al. COVID-19 surveillance in
734 Southeastern Virginia using wastewater-based epidemiology. *Water Res.* 2020;186:
735 116296. doi:10.1016/j.watres.2020.116296
- 736 5. Hoar C, Chauvin F, Clare A, McGibbon H, Castro E, Patinella S, et al. Monitoring SARS-
737 CoV-2 in wastewater during New York City's second wave of COVID-19: Sewershed-level
738 trends and relationships to publicly available clinical testing data. *medRxiv.* 2022;
739 2022.02.08.22270666. doi:10.1101/2022.02.08.22270666
- 740 6. Trujillo M, Cheung K, Gao A, Hoxie I, Kannoly S, Kubota N, et al. Protocol for Safe,
741 Affordable, and Reproducible Isolation and Quantitation of SARS-CoV-2 RNA from
742 Wastewater. *medRxiv.* 2021; 2021.02.16.21251787. doi:10.1101/2021.02.16.21251787
- 743 7. Kirby AE, Welsh RM, Marsh ZA, Yu AT, Vugia DJ, Boehm AB, et al. Notes from the Field:
744 Early Evidence of the SARS-CoV-2 B.1.1.529 (Omicron) Variant in Community Wastewater
745 - United States, November-December 2021. *MMWR Morb Mortal Wkly Rep.* 2022;71: 103–
746 105. doi:10.15585/mmwr.mm7103a5
- 747 8. Smyth DS, Trujillo M, Gregory DA, Cheung K, Gao A, Graham M, et al. Tracking cryptic
748 SARS-CoV-2 lineages detected in NYC wastewater. *Nat Commun.* 2022;13: 635.
749 doi:10.1038/s41467-022-28246-3
- 750 9. Gregory DA, Wieberg CG, Wenzel J, Lin C-H, Johnson MC. Monitoring SARS-CoV-2
751 Populations in Wastewater by Amplicon Sequencing and Using the Novel Program SAM
752 Refiner. *Viruses.* 2021;13. doi:10.3390/v13081647
- 753 10. Martin DP, Weaver S, Tegally H, San JE, Shank SD, Wilkinson E, et al. The emergence and
754 ongoing convergent evolution of the SARS-CoV-2 N501Y lineages. *Cell.* 2021/09/07 ed.
755 2021;184: 5189-5200.e7. doi:10.1016/j.cell.2021.09.003

- 756 11. Callaway E. BEYONDOMICRON: WHAT'S NEXT FOR SARS-COV-2 EVOLUTION.
757 NATURE. 2021;600: 204–207.
- 758 12. Hill V, Du Plessis L, Peacock TP, Aggarwal D, Colquhoun R, Carabelli AM, et al. The origins
759 and molecular evolution of SARS-CoV-2 lineage B.1.1.7 in the UK. bioRxiv. 2022;
760 2022.03.08.481609. doi:10.1101/2022.03.08.481609
- 761 13. Swift CL, Isanovic M, Correa Velez KE, Norman RS. Community-level SARS-CoV-2
762 sequence diversity revealed by wastewater sampling. Sci Total Environ. 2021/08/18 ed.
763 2021;801: 149691–149691. doi:10.1016/j.scitotenv.2021.149691
- 764 14. Herold M, d'Hérouël AF, May P, Delogu F, Wienecke-Baldacchino A, Tapp J, et al. Genome
765 Sequencing of SARS-CoV-2 Allows Monitoring of Variants of Concern through Wastewater.
766 Water. 2021;13. doi:10.3390/w13213018
- 767 15. Fontenele RS, Kraberger S, Hadfield J, Driver EM, Bowes D, Holland LA, et al. High-
768 throughput sequencing of SARS-CoV-2 in wastewater provides insights into circulating
769 variants. medRxiv. 2021; 2021.01.22.21250320. doi:10.1101/2021.01.22.21250320
- 770 16. Baaijens JA, Zulli A, Ott IM, Petrone ME, Alpert T, Fauver JR, et al. Variant abundance
771 estimation for SARS-CoV-2 in wastewater using RNA-Seq quantification. medRxiv. 2021;
772 2021.08.31.21262938. doi:10.1101/2021.08.31.21262938
- 773 17. Greaney AJ, Starr TN, Barnes CO, Weisblum Y, Schmidt F, Caskey M, et al. Mapping
774 mutations to the SARS-CoV-2 RBD that escape binding by different classes of antibodies.
775 Nat Commun. 2021;12: 4196. doi:10.1038/s41467-021-24435-8
- 776 18. Harvey WT, Carabelli AM, Jackson B, Gupta RK, Thomson EC, Harrison EM, et al. SARS-
777 CoV-2 variants, spike mutations and immune escape. Nat Rev Microbiol. 2021;19: 409–424.
778 doi:10.1038/s41579-021-00573-0
- 779 19. Shang J, Ye G, Shi K, Wan Y, Luo C, Aihara H, et al. Structural basis of receptor recognition
780 by SARS-CoV-2. Nature. 2020;581: 221–224. doi:10.1038/s41586-020-2179-y
- 781 20. Liu H, Wei P, Kappler JW, Marrack P, Zhang G. SARS-CoV-2 Variants of Concern and

- 782 Variants of Interest Receptor Binding Domain Mutations and Virus Infectivity. *Front*
783 *Immunol.* 2022;13. Available: <https://www.frontiersin.org/article/10.3389/fimmu.2022.825256>
- 784 21. Li Q, Wu J, Nie J, Zhang L, Hao H, Liu S, et al. The Impact of Mutations in SARS-CoV-2
785 Spike on Viral Infectivity and Antigenicity. *Cell.* 2020/07/17 ed. 2020;182: 1284-1294.e9.
786 doi:10.1016/j.cell.2020.07.012
- 787 22. Khare S, Gurry C, Freitas L, Schultz MB, Bach G, Diallo A, et al. GISAID's Role in
788 Pandemic Response. *China CDC Wkly.* 2021;3: 1049–1051. doi:10.46234/ccdcw2021.255
- 789 23. Elbe S, Buckland-Merrett G. Data, disease and diplomacy: GISAID's innovative contribution
790 to global health. *Glob Chall Hoboken NJ.* 2017;1: 33–46. doi:10.1002/gch2.1018
- 791 24. Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data - from vision to
792 reality. *Euro Surveill Bull Eur Sur Mal Transm Eur Commun Dis Bull.* 2017;22: 30494.
793 doi:10.2807/1560-7917.ES.2017.22.13.30494
- 794 25. Smyth DS, Trujillo M, Cheung K, Gao A, Hoxie I, Kannoly S, et al. Detection of Mutations
795 Associated with Variants of Concern Via High Throughput Sequencing of SARS-CoV-2
796 Isolated from NYC Wastewater. *medRxiv.* 2021; 2021.03.21.21253978.
797 doi:10.1101/2021.03.21.21253978
- 798 26. [nychealth/coronavirus-data](https://github.com/nychealth/coronavirus-data). NYC Department of Health and Mental Hygiene; Available:
799 <https://github.com/nychealth/coronavirus-data/blob/master/variants/variant-epi-data.csv>
- 800 27. Starr TN, Greaney AJ, Dingens AS, Bloom JD. Complete map of SARS-CoV-2 RBD
801 mutations that escape the monoclonal antibody LY-CoV555 and its cocktail with LY-
802 CoV016. *Cell Rep Med.* 2021;2: 100255. doi:10.1016/j.xcrm.2021.100255
- 803 28. Starr TN, Greaney AJ, Addetia A, Hannon WW, Choudhary MC, Dingens AS, et al.
804 Prospective mapping of viral mutations that escape antibodies used to treat COVID-19.
805 *Science.* 2021;371: 850–854. doi:10.1126/science.abf9302
- 806 29. Starr TN, Greaney AJ, Hilton SK, Ellis D, Crawford KHD, Dingens AS, et al. Deep
807 Mutational Scanning of SARS-CoV-2 Receptor Binding Domain Reveals Constraints on

- 808 Folding and ACE2 Binding. *Cell*. 2020;182: 1295-1310.e20. doi:10.1016/j.cell.2020.08.012
- 809 30. Lan J, Ge J, Yu J, Shan S, Zhou H, Fan S, et al. Structure of the SARS-CoV-2 spike
810 receptor-binding domain bound to the ACE2 receptor. *Nature*. 2020;581: 215–220.
811 doi:10.1038/s41586-020-2180-5
- 812 31. Dinno KH, Leist SR, Schäfer A, Edwards CE, Martinez DR, Montgomery SA, et al. A
813 mouse-adapted model of SARS-CoV-2 to test COVID-19 countermeasures. *Nature*.
814 2020;586: 560–566. doi:10.1038/s41586-020-2708-8
- 815 32. Wang J, Shuai L, Wang C, Liu R, He X, Zhang X, et al. Mouse-adapted SARS-CoV-2
816 replicates efficiently in the upper and lower respiratory tract of BALB/c and C57BL/6J mice.
817 *Protein Cell*. 2020;11: 776–782. doi:10.1007/s13238-020-00767-x
- 818 33. Gawish R, Starkl P, Pimenov L, Hladik A, Lakovits K, Oberndorfer F, et al. ACE2 is the
819 critical in vivo receptor for SARS-CoV-2 in a novel COVID-19 mouse model with TNF- and
820 IFN γ -driven immunopathology. *eLife*. 2022;11: e74623. doi:10.7554/eLife.74623
- 821 34. Weisblum Y, Schmidt F, Zhang F, DaSilva J, Poston D, Lorenzi JC, et al. Escape from
822 neutralizing antibodies by SARS-CoV-2 spike protein variants. *eLife*. 2020;9: e61312.
823 doi:10.7554/eLife.61312
- 824 35. McCarthy Kevin R., Rennick Linda J., Nambulli Sham, Robinson-McCarthy Lindsey R., Bain
825 William G., Haidar Ghady, et al. Recurrent deletions in the SARS-CoV-2 spike glycoprotein
826 drive antibody escape. *Science*. 2021;371: 1139–1142. doi:10.1126/science.abf6950
- 827 36. Korber B, Fischer WM, Gnanakaran S, Yoon H, Theiler J, Abfalterer W, et al. Tracking
828 Changes in SARS-CoV-2 Spike: Evidence that D614G Increases Infectivity of the COVID-19
829 Virus. *Cell*. 2020/07/03 ed. 2020;182: 812-827.e19. doi:10.1016/j.cell.2020.06.043
- 830 37. Pickering B, Lung O, Maguire F, Kruczkiewicz P, Kotwa JD, Buchanan T, et al. Highly
831 divergent white-tailed deer SARS-CoV-2 with potential deer-to-human transmission.
832 *Microbiology*; 2022 Feb. doi:10.1101/2022.02.22.481551
- 833 38. Natarajan A, Zlitni S, Brooks EF, Vance SE, Dahlen A, Hedlin H, et al. Gastrointestinal

- 834 symptoms and fecal shedding of SARS-CoV-2 RNA suggest prolonged gastrointestinal
835 infection. *Med.* 2022; S2666634022001672. doi:10.1016/j.medj.2022.04.001
- 836 39. Zollner A, Koch R, Jukic A, Pfister A, Meyer M, Rössler A, et al. Postacute COVID-19 is
837 Characterized by Gut Viral Antigen Persistence in Inflammatory Bowel Diseases.
838 *Gastroenterology.* 2022; S0016508522004504. doi:10.1053/j.gastro.2022.04.037
- 839 40. Greaney AJ, Starr TN, Barnes CO, Weisblum Y, Schmidt F, Caskey M, et al. Mutational
840 escape from the polyclonal antibody response to SARS-CoV-2 infection is largely shaped by
841 a single class of antibodies. *Microbiology;* 2021 Mar. doi:10.1101/2021.03.17.435863
- 842 41. Greaney AJ, Loes AN, Crawford KHD, Starr TN, Malone KD, Chu HY, et al. Comprehensive
843 mapping of mutations in the SARS-CoV-2 receptor-binding domain that affect recognition by
844 polyclonal human plasma antibodies. *Cell Host Microbe.* 2021;29: 463-476.e6.
845 doi:10.1016/j.chom.2021.02.003
- 846 42. Liu Z, VanBlargan LA, Bloyet L-M, Rothlauf PW, Chen RE, Stumpf S, et al. Identification of
847 SARS-CoV-2 spike mutations that attenuate monoclonal and serum antibody neutralization.
848 *Cell Host Microbe.* 2021;29: 477-488.e4. doi:10.1016/j.chom.2021.01.014
- 849 43. Coronavirus Antiviral & Resistance Database. Stanford University; Available:
850 <https://covdb.stanford.edu/search-drdb>
- 851 44. Wilkinson SA, Richter A, Casey A, Osman H, Mirza JD, Stockton J, et al. Recurrent SARS-
852 CoV-2 Mutations in Immunodeficient Patients. *medRxiv.* 2022; 2022.03.02.22271697.
853 doi:10.1101/2022.03.02.22271697
- 854 45. Crits-Christoph A, Kantor RS, Olm MR, Whitney ON, Al-Shayeb B, Lou YC, et al. Genome
855 Sequencing of Sewage Detects Regionally Prevalent SARS-CoV-2 Variants. Pettigrew MM,
856 editor. *mBio.* 2021;12: e02703-20. doi:10.1128/mBio.02703-20
- 857 46. Fontenele RS, Kraberger S, Hadfield J, Driver EM, Bowes D, Holland LA, et al. High-
858 throughput sequencing of SARS-CoV-2 in wastewater provides insights into circulating
859 variants. *Water Res.* 2021;205: 117710. doi:10.1016/j.watres.2021.117710

- 860 47. Izquierdo-Lara R, Elsinga G, Heijnen L, Munnink BBO, Schapendonk CME, Nieuwenhuijse
861 D, et al. Monitoring SARS-CoV-2 Circulation and Diversity through Community Wastewater
862 Sequencing, the Netherlands and Belgium. *Emerg Infect Dis.* 2021;27: 1405–1415.
863 doi:10.3201/eid2705.204410
- 864 48. Cotten M, Lule Bugembe D, Kaleebu P, V T Phan M. Alternate primers for whole-genome
865 SARS-CoV-2 sequencing. *Virus Evol.* 2021;7: veab006–veab006. doi:10.1093/ve/veab006
- 866 49. Xiao M, Liu X, Ji J, Li M, Li J, Yang L, et al. Multiple approaches for massively parallel
867 sequencing of SARS-CoV-2 genomes directly from clinical samples. *Genome Med.*
868 2020;12: 57. doi:10.1186/s13073-020-00751-4
- 869 50. Addetia Amin, Lin Michelle J., Peddu Vikas, Roychoudhury Pavitra, Jerome Keith R.,
870 Greninger Alexander L., et al. Sensitive Recovery of Complete SARS-CoV-2 Genomes from
871 Clinical Samples by Use of Swift Biosciences' SARS-CoV-2 Multiplex Amplicon Sequencing
872 Panel. *J Clin Microbiol.* 59: e02226-20. doi:10.1128/JCM.02226-20
- 873 51. Dezordi FZ, Neto AM da S, Campos T de L, Jeronimo PMC, Aksenon CF, Almeida SP, et
874 al. ViralFlow: A Versatile Automated Workflow for SARS-CoV-2 Genome Assembly, Lineage
875 Assignment, Mutations and Intra-host Variant Detection. *Viruses.* 2022;14: 217.
876 doi:10.3390/v14020217
- 877 52. Van Poelvoorde LAE, Delcourt T, Coucke W, Herman P, De Keersmaecker SCJ, Saelens X,
878 et al. Strategy and Performance Evaluation of Low-Frequency Variant Calling for SARS-
879 CoV-2 Using Targeted Deep Illumina Sequencing. *Front Microbiol.* 2021;12. Available:
880 <https://www.frontiersin.org/article/10.3389/fmicb.2021.747458>
- 881 53. Lin X, Glier M, Kuchinski K, Ross-Van Mierlo T, McVea D, Tyson JR, et al. Assessing
882 Multiplex Tiling PCR Sequencing Approaches for Detecting Genomic Variants of SARS-
883 CoV-2 in Municipal Wastewater. *mSystems.* 2021/10/19 ed. 2021;6: e0106821–e0106821.
884 doi:10.1128/mSystems.01068-21
- 885 54. N Whitney O, Al-Shayeb B, Crits-Cristoph A, Chaplin M, Fan V, Greenwald H, et al. V.4 -

886 Direct wastewater RNA capture and purification via the Ca^{2+} -Sewage, Salt, Silica and
 887 SARS-CoV-2 (4S) Ca^{2+} ; method v4. 2020 Nov. doi:10.17504/protocols.io.bpdfmi3n
 888 55. Rognes T, Flouri T, Nichols B, Quince C, Mahé F. VSEARCH: a versatile open source tool
 889 for metagenomics. PeerJ. 2016;4: e2584. doi:10.7717/peerj.2584
 890 56. Li H. Minimap2: pairwise alignment for nucleotide sequences. Birol I, editor. Bioinformatics.
 891 2018;34: 3094–3100. doi:10.1093/bioinformatics/bty191
 892 57. Gregory DA, Wieberg CG, Wenzel J, Lin C-H, Johnson MC. Monitoring SARS-CoV-2
 893 Populations in Wastewater by Amplicon Sequencing and Using the Novel Program SAM
 894 Refiner. Viruses. 2021;13: 1647. doi:10.3390/v13081647
 895

896 **Supporting information captions**

897 **Supplemental Table 1. Prevalence in GISAID of common substitutions found in**
 898 **cryptic lineages.** Data reflects the number of sequences from humans deposited into
 899 GISAID by the indicated dates [22].

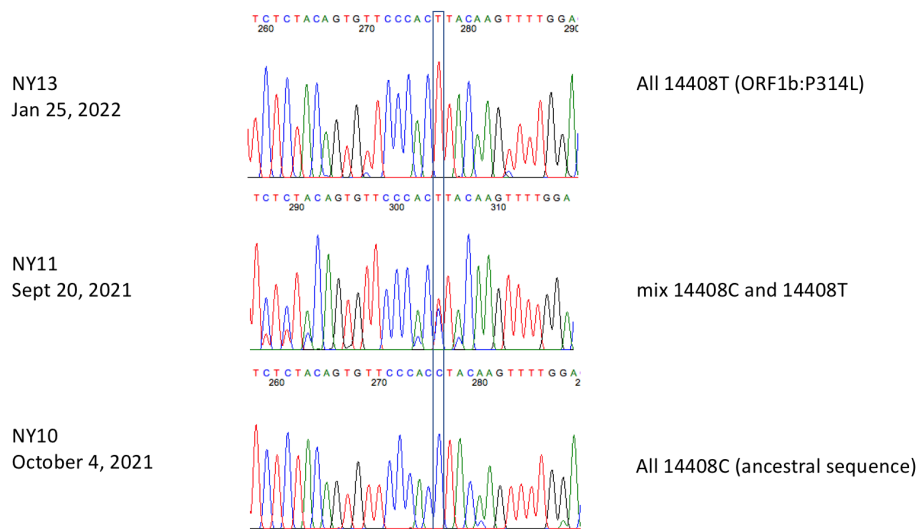
Substitution	Location	Global Prevalence in Humans 11/1/21	Global Prevalence 3/15/22
Total Sequences		4,824,812	9,349,201
G413R	NY10, NY14	147	239
K417T	MO33, MO45, NY2, NY3, NY11, NY13, NY14, CA	108,537	119,127
N439K	MO33, NY2, NY3	37,227	40,274
N440K	MO33, NY11	9,154	1,652,112

K444Δ	NY14	14	68
K444S	NY2, NY3, NY11	21	26
K444T	NY2, NY10, NY11, NY14	58	217
V445Δ	NY2, NY11, NY14	24	84
V445A	CA, NY3, NY13, NY14	298	557
G446S	CA, NY2	472	1,336,425
G446Δ	NY2, NY11	19	82
G446D	NY3, NY13	75	224
Y449H	NY2, NY11	821	1,094
Y449R	NY2, NY3, NY11	0	0
L452Q	NY2, NY3, NY11	10,498	12,137
L452R	MO33, NY3, NY14	2,315,718	4,324,990
Y453F	NY3, NY11	1,320	1,497
L455W	NY3, NY14	5	21
F456L	NY2, NY3, NY10, NY14	259	736
N460K	NY2, CA, MO33, MO45, NY3, NY10, NY11, NY13, NY14	76	242
S477N	CA, NY3, NY10, NY11, NY14	71,960	2,164,897

T478K	CA, MO45, NY2, NY3	2,249,016	6,340,479
V483Δ	NY2	49	932
E484Δ	NY3, NY11	31	912
E484A	CA, MO45, NY2, NY10, NY11, NY13	551	2,087,453
E484P	CA, NY2	0	73
E484V	NY3, NY10	104	1,610
F486P	CA, NY3, NY10, NY11	2	3
F486V	NY3, NY13, NY14	4	34
F490H	NY3, NY14	1	2
F490Y	CA, NY2, NY10, NY11, NY14	120	163
Q493R	NY3, NY14	261	2,083,669
Q493K	MO33, MO45, NY10, NY13, NY14	152	835
S494P	MO45, NY2, NY10, NY11	12,916	15,009
Q498H	CA, MO45, NY2, NY11, NY14	36	57
Q498R	NY13	91	2,007,408
Q498Y	NY2, CA, NY3, NY10, NY11,	0	13

	NY14		
P499H	NY2, NY11	77	118
P499S	CA, NY3, NY14	216	353
N501S	CA, NY10, NY11, NY14	166	663
N501T	CA, MO33, NY2, NY3, NY10, NY11, NY14	4,742	5,639
N501Y	NY11, NY13	1,325,387	3,389,688
G504D	NY3, NY14	190	580
Y505H	CA, NY2, NY3, NY11, NY13, NY14	133	2,013,881
Y505N	NY2, NY3	2	5
H519N	NY10, NY11	13	31
T572I	CA, NY2	16,610	26,948
T572N	NY10, NY14	148	298

900



901

902 **Supplemental Figure 1.** Sequence of nt 14408 from NYC wastewater.

903

904