

1 Wastewater sequencing uncovers early, cryptic SARS-CoV-2 variant transmission

2
3
4 Smruthi Karthikeyan^{1,#}, Joshua I Levy^{2,#}, Peter De Hoff^{3,9,21}, Greg Humphrey¹, Amanda
5 Birmingham⁴, Kristen Jepsen⁵, Sawyer Farmer¹, Helena M. Tubb¹, Tommy Valles¹, Caitlin E
6 Tribelhorn¹, Rebecca Tsai¹, Stefan Aigner³, Shashank Sathe³, Niema Moshiri⁶, Benjamin
7 Henson⁵, Adam M. Mark⁴, Abbas Hakim^{3,9,21}, Nathan A Baer³, Tom Barber³, Pedro Belda-
8 Ferre³, Marisol Chacón³, Willi Cheung^{3,9,21}, Evelyn S Cresini³, Emily R Eisner³, Alma L
9 Lastrella³, Elijah S Lawrence³, Clarisse A Marotz³, Toan T Ngo³, Tyler Ostrander³, Ashley
10 Plascencia³, Rodolfo A Salido³, Phoebe Seaver³, Elizabeth W Smoot³, Daniel McDonald¹,
11 Robert M Neuhard^{7,12}, Angela L Scioscia^{8,9}, Alysson M. Satterlund¹⁰, Elizabeth H Simmons¹¹,
12 Dismas B. Abelman¹², David Brenner¹², Judith C. Bruner¹², Anne Buckley¹², Michael Ellison¹²,
13 Jeffrey Gattas¹², Steven L. Gonias¹³, Matt Hale¹², Faith Hawkins¹², Lydia Ikeda¹², Hemlata
14 Jhaveri¹², Ted Johnson¹², Vince Kellen¹², Brendan Kremer¹², Gary Matthews¹², Ronald W.
15 McLawhon¹², Pierre Ouillet¹², Daniel Park¹², Allorah Pradenas¹², Sharon Reed¹², Lindsay
16 Riggs¹², Alison Sanders¹², Bradley Sollenberger¹², Angela Song^{7,12}, Benjamin White¹², Terri
17 Winbush¹², Christine M Aceves², Catelyn Anderson², Karthik Gangavarapu², Emory Hufbauer²,
18 Ezra Kurzban², Justin Lee², Nathaniel L Matteson², Edyth Parker², Sarah A Perkins², Karthik S
19 Ramesh², Refugio Robles-Sikisaka², Madison A Schwab², Emily Spencer², Shirlee Wohl², Laura
20 Nicholson², Ian H Mchardy², David P Dimmock¹⁵, Charlotte A Hobbs¹⁵, Omid Bakhtar¹⁶, Aaron
21 Harding¹⁶, Art Mendoza¹⁶, Alexandre Bolze¹⁷, David Becker¹⁷, Elizabeth T Cirulli¹⁷, Magnus
22 Isaksson¹⁷, Kelly M Schiabor Barrett¹⁷, Nicole L Washington¹⁷, John D Malone¹⁸, Ashleigh
23 Murphy Schafer¹⁸, Nikos Gurfield¹⁸, Sarah Stous¹⁸, Rebecca Fielding-Miller^{19,20}, Richard S.
24 Garfein¹⁹, Tommi Gaines²⁰, Cheryl Anderson¹⁹, Natasha K. Martin¹⁹, Robert Schooley¹⁹, Brett
25 Austin¹⁶, Duncan R. MacCannell²², Stephen F Kingsmore¹⁵, William Lee¹⁷, Seema Shah¹⁸, Eric
26 McDonald¹⁸, Alexander T. Yu²¹, Mark Zeller², Kathleen M Fisch^{4,9}, Christopher Longhurst^{1,23},
27 Patty Maysent²⁴, David Pride²⁵, Pradeep K. Khosla⁶, Louise C. Laurent^{3,9,26}, Gene W Yeo^{3,26,27},
28 Kristian G Andersen^{2,*}, Rob Knight^{1,6,28,*}

29
30 #equal contribution

31 *Senior author

32
33 ¹ Department of Pediatrics, University of California San Diego, La Jolla, CA, USA

34 ² Department of Immunology and Microbiology, The Scripps Research Institute, La Jolla, CA,
35 USA

36 ³ Expedited COVID Identification Environment (EXCITE) Laboratory, Department of Pediatrics,
37 University of California San Diego, La Jolla, CA, USA

38 ⁴ Center for Computational Biology and Bioinformatics, University of California San Diego, La
39 Jolla, CA, USA

40 ⁵ Institute for Genomic Medicine, University of California San Diego, La Jolla, CA, USA

41 ⁶ Department of Computer Science and Engineering, University of California San Diego, La
42 Jolla, CA, USA

43 ⁷ Operational Strategic Initiatives, University of California San Diego, La Jolla, CA, USA

44 ⁸ Student Health and Well-Being, University of California San Diego, La Jolla, CA, USA

45 ⁹ Department of Obstetrics, Gynecology, and Reproductive Sciences, University of California
46 San Diego, La Jolla, CA, USA

- 47 ¹⁰ Student Affairs, University of California San Diego, La Jolla, CA, USA
48 ¹¹ Academic Affairs, University of California San Diego, La Jolla, CA, USA
49 ¹² Return to Learn, University of California San Diego, La Jolla, CA, USA
50 ¹³ Department of Pathology, University of California San Diego, La Jolla, CA, USA
51 ¹⁴ Scripps Health, San Diego, La Jolla, CA, USA
52 ¹⁵ Rady Children's Institute for Genomic Medicine, San Diego, CA, USA
53 ¹⁶ Sharp Healthcare, San Diego, CA, USA
54 ¹⁷ Helix, San Mateo, CA, USA
55 ¹⁸ County of San Diego Health and Human Services Agency, San Diego, CA, USA
56 ¹⁹ Herbert Wertheim School of Public Health and Human Longevity Science, University of
57 California San Diego, La Jolla, CA, USA
58 ²⁰ Division of Infectious Disease and Global Public Health, University of California San Diego,
59 La Jolla, CA, USA
60 ²¹ COVID-19 Detection, Investigation, Surveillance, Clinical, and Outbreak Response, California
61 Department of Public Health, Richmond, CA, USA
62 ²² Office of Advanced Molecular Detection, Centers for Disease Control and Prevention, Atlanta,
63 GA, USA
64 ²³ Department of Biomedical Informatics, University of California, San Diego, La Jolla,
65 California, USA
66 ²⁴ Office of the UC San Diego Health CEO, University of California, San Diego
67 ²⁵ Departments of Pathology and Medicine, University of California, San Diego, La Jolla, CA
68 ²⁶ Sanford Consortium of Regenerative Medicine, University of California San Diego, La Jolla,
69 CA
70 ²⁷ Department of Cellular and Molecular Medicine, University of California San Diego, La Jolla,
71 CA
72 ²⁸ Department of Bioengineering, University of California San Diego, La Jolla, CA, USA

73

74 **Address correspondence to:**

75

76 Rob Knight
77 Department of Pediatrics
78 University of California San Diego
79 9500 Gilman Drive, MC 0763
80 La Jolla, CA 92093, USA
81 robknight@ucsd.edu
82 +1 858-246-1184

83

84

85 **Summary**

86

87 As SARS-CoV-2 continues to spread and evolve, detecting emerging variants early is critical for
88 public health interventions. Inferring lineage prevalence by clinical testing is infeasible at scale,
89 especially in areas with limited resources, participation, or testing/sequencing capacity, which
90 can also introduce biases. SARS-CoV-2 RNA concentration in wastewater successfully tracks
91 regional infection dynamics and provides less biased abundance estimates than clinical testing.
92 Tracking virus genomic sequences in wastewater would improve community prevalence

93 estimates and detect emerging variants. However, two factors limit wastewater-based genomic
94 surveillance: low-quality sequence data and inability to estimate relative lineage abundance in
95 mixed samples. Here, we resolve these critical issues to perform a high-resolution, 295-day
96 wastewater and clinical sequencing effort, in the controlled environment of a large university
97 campus and the broader context of the surrounding county. We develop and deploy improved
98 virus concentration protocols and deconvolution software that fully resolve multiple virus strains
99 from wastewater. We detect emerging variants of concern up to 14 days earlier in wastewater
100 samples, and identify multiple instances of virus spread not captured by clinical genomic
101 surveillance. Our study provides a scalable solution for wastewater genomic surveillance that
102 allows early detection of SARS-CoV-2 variants and identification of cryptic transmission.

103

104 **Introduction**

105

106 SARS-CoV-2 continues to evolve, producing diverse new lineages¹. Emerging variants of
107 concern (VOCs) and variants of interest (VOIs) demonstrate increased transmissibility, disease
108 severity, and/or immune escape². Timely and accurate quantification of local prevalence of
109 SARS-CoV-2 variants is thus essential for effective public health measures. However, existing
110 strategies for variant detection based on virus genome sequencing of biospecimens obtained from
111 clinical testing (“clinical genomic surveillance”) are expensive, inefficient, and have sampling
112 bias because of systemic healthcare disparities, particularly in poor and underserved
113 communities³⁻⁵.

114

115 In contrast, PCR-based wastewater surveillance of SARS-CoV-2 RNA is not subject to clinical
116 testing biases and can track temporal changes in overall SARS-CoV-2 prevalence in a region⁶⁻⁸,
117 but cannot identify epidemiological transmission links or monitor virus lineage prevalence,
118 which require genome sequence information. Virus genome sequencing from wastewater
119 (“wastewater genomic surveillance”) has the potential to cost-effectively capture community
120 virus spread^{9,10}, acting as a surrogate to clinical surveillance in elucidating lineage geospatial
121 distributions and track emerging SARS-CoV-2 variants (including new variants for which
122 targeted assays do not yet exist), and provide genome sequence data needed for transmission
123 network analysis and interpretation¹¹.

124

125 However, wastewater genomic surveillance is technically challenging¹⁰. Low viral loads, heavily
126 fragmented RNA, and PCR inhibitors in complex environmental samples lead to poor
127 sequencing coverage^{12,13}. Obtaining high quality sequences from samples with low viral load and
128 elevated levels of PCR inhibitors remains an outstanding technical challenge in implementation
129 of wastewater genomic surveillance at scale. Additionally, tools for SARS-CoV-2 lineage
130 classification, such as pangolin¹⁴ and UShER¹⁵, were designed for clinical samples containing a
131 single dominant variant, and cannot estimate relative abundances of multiple SARS-CoV-2
132 lineages in samples with virus mixtures such as wastewater.

133

134 Here, we report a high-resolution approach to study community virus transmission using
135 wastewater genomic surveillance, leveraging several technical advances in wastewater virus
136 concentration and nucleic acid sequencing, and a computational tool for resolving multiple
137 SARS-CoV-2 lineages in short-read sequence data from a mixed sample (lineage deconvolution).
138 We obtained near 95% genome coverage even for samples with low viral load, compared with

139 40% or below from previous studies¹¹⁻¹³, a key advance that allowed us to build a robust pipeline
140 to monitor virus lineage prevalence in community wastewater.

141
142 Because places of communal living, such as university campuses, are considered key sites for
143 virus spread and represent well-controlled and relatively isolated environments, they are ideal for
144 comparing the relative utility of clinical and wastewater genomic surveillance¹⁶. Accordingly, we
145 conducted a high-resolution, longitudinal wastewater genomic surveillance effort at the
146 University of California San Diego (UCSD) campus, in parallel with clinical genomic
147 surveillance from nasal swabs in the local community, from November 2020 to September 2021:
148 ten months that effectively capture the surges in the region caused by the three main VOCs (as
149 determined by US CDC) in the United States, Epsilon, Alpha and Delta¹. In more recent San
150 Diego-wide data collected from September 2021 to February 2022, we studied ongoing
151 transmission of the Delta variant and the rapid spread of the Omicron variant and its sublineages.

152
153 Our wastewater genomic surveillance approach identified VOCs up to 2 weeks prior to detection
154 through clinical genomic surveillance, even though a large proportion of clinical SARS-CoV-2
155 samples are sequenced in San Diego relative to other cities in the United States. In addition to
156 providing a detailed history of community virus spread, wastewater genomic surveillance also
157 identified multiple instances of cryptic community transmission not observed through clinical
158 genomic surveillance. Matching wastewater and clinical genome sequences provided
159 epidemiological information identifying specific transmission events. Our results demonstrate
160 the viability of wastewater genomic surveillance at scale, enabling early detection and tracking
161 of virus lineages and guiding clinical genomic surveillance efforts. This work informed public
162 health guidance and interventions on the UCSD campus as well as San Diego county in real
163 time, and our data and analyses were disseminated to both public health officials as well as the
164 general public via custom dashboards (see **Data Availability** for links).

165 166 **Results**

167
168 To directly compare wastewater genomic surveillance to clinical surveillance, we conducted a
169 large-scale SARS-CoV-2 genome sequencing study from wastewater samples collected daily
170 from 131 wastewater samplers covering 360 campus buildings, in many cases reaching single
171 building-level resolution. To identify epidemiological transmission links and monitor lineages in
172 the population, we sequenced all SARS-CoV-2 positive clinical and wastewater samples from
173 campus using a miniaturized tiled-amplicon sequencing approach. During this period of this
174 study, we collected and analyzed 21,383 wastewater samples: 19,944 wastewater samples from
175 the UCSD campus, and, for comparison, 1,475 wastewater samples from the greater San Diego
176 area, including the Point Loma wastewater treatment plant (the primary wastewater treatment
177 plant for the county with a catchment size of 2.3 million people) and 17 public schools spanning
178 four San Diego school districts¹⁷. We compared sequencing of 600 campus wastewater samples
179 to 759 genomes obtained from campus clinical swabs (46.2% of all positive tests on campus), all
180 processed by the CALM and EXCITE CLIA labs at UCSD. In addition, we compared 31,149
181 genomes obtained from clinical genomic surveillance of the greater San Diego community to
182 sequencing of 837 wastewater samples collected from San Diego county (including those from
183 the UCSD campus) during the same period.

184

185

186 High-resolution spatial sampling reveals micro-scale community spread

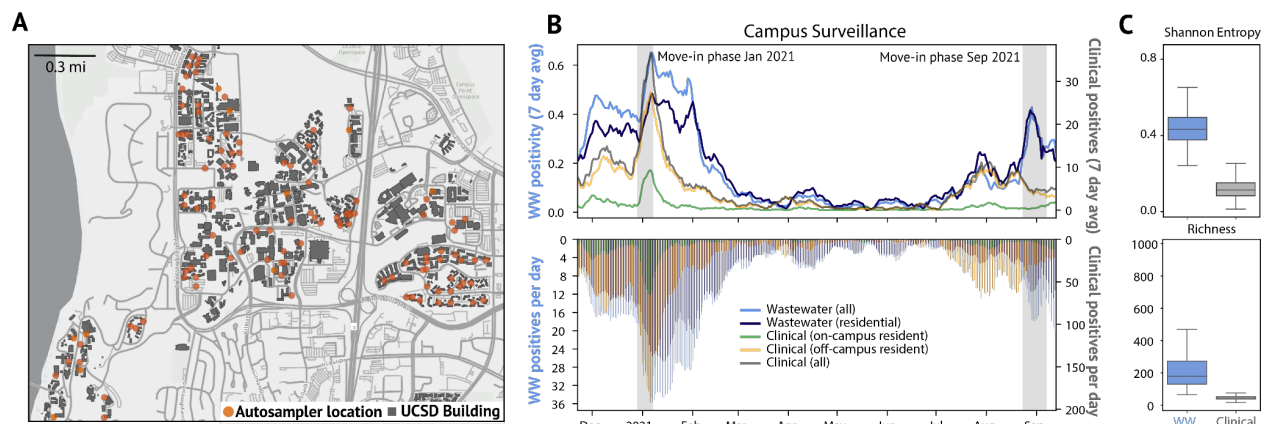
187

188 We implemented a GIS (geographic information system)-enabled building-level wastewater
189 surveillance system to cover 360 buildings on the UCSD campus (**Figure 1A**). During the period
190 of daily wastewater sampling, approximately 10,000 students lived on campus and 25,000
191 individuals were on campus on a daily basis. We found that wastewater test positivity correlated
192 strongly with the number of clinical positives (**Figure 1B** and **Extended Data Figure 1**),
193 showing that wastewater effectively captures the community infection dynamics based on total
194 viral load. This is also consistent with our past studies that showed SARS-CoV-2 RNA can be
195 detected ~85% of the time downstream from buildings containing individuals known to be
196 infected⁹.

197

198 Unlike qPCR-based mutant surveillance, genomic surveillance using full-length virus genomes
199 can detect which strains of SARS-CoV-2 are circulating in the population, and can identify
200 potential transmission links between infected individuals^{18,19}. While targeted qPCR mutant
201 panels have the ability to detect specific lineages in wastewater, they only target a small set of
202 mutations that must be known beforehand in addition to the development and validation time
203 before implementation. Furthermore, they cannot provide sub-lineage resolution (for instance,
204 BA.1 v. BA.2 sublineage of Omicron) and will fail altogether if a sublineage loses the specific
205 mutation targeted by the qPCR assay. To test the utility of wastewater genomic surveillance for
206 studying virus spread in the community, we obtained near complete virus genomes for
207 wastewater samples with cycle quantification (Cq) values as high as 38 (median genome
208 coverage: 96.49% [75.67% - 100.00%], **Extended Data Figure 2**). However, using two common
209 metrics of virus diversity, Shannon entropy (a measure of the uncertainty associated with
210 randomly sampling an allele) and richness (the number of single nucleotide variant, or SNV,
211 sites)²⁰, we found that SARS-CoV-2 genetic diversity is significantly greater in wastewater
212 samples than clinical samples (**Figure 1C**, Mann-Whitney U test, $p < 0.001$ for each, with effect
213 size $r = 0.99$, 0.97 for Shannon Entropy and Richness, respectively). This suggests that multiple
214 virus lineages, likely shed from different infected individuals, are often present in wastewater
215 samples.

216



217

218 **Figure 1: Campus sampling locations and SARS-CoV-2 testing statistics.** A. Geospatial
219 distribution of the 131 actively deployed wastewater autosamplers and the corresponding 360
220 university buildings on the campus sewer network. Building-specific data have been de-

221 identified in accordance with university reporting policies. B. Campus wastewater and diagnostic
222 testing statistics over the 295 day sampling period (WW = wastewater, positivity is the fraction
223 of WW samplers with a positive qPCR signal). C. Virus diversity in wastewater and clinical
224 samples: Boxplots of Shannon entropy (top) and richness (bottom) for each sample type.

225 226 **Sample deconvolution robustly recovers the abundance of SARS-CoV-2 lineages in mixed** 227 **samples**

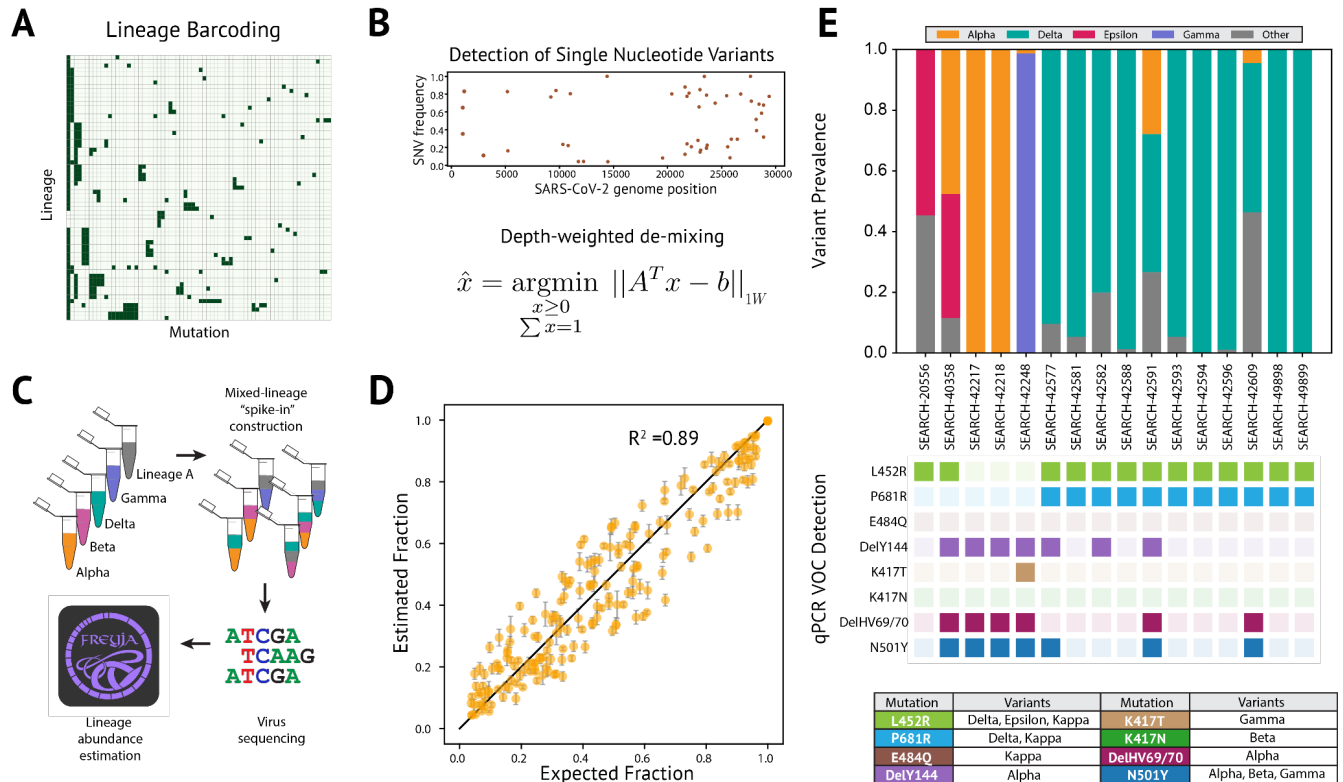
228
229 Wastewater systems aggregate stool, urine, and other biological waste products carrying viruses
230 from multiple infected individuals in the community in a single location, allowing for sampling
231 of virus mixtures that are representative of local lineage prevalence. However, existing methods
232 for determining virus lineage from sequencing are intended for non-mixed clinical samples and
233 can only be used to identify a single (dominant) lineage per sample.

234
235 To fully capture the virus diversity in community biospecimens, we developed Freyja, a tool to
236 estimate the relative abundance of virus lineages in a mixed sample. Freyja uses a “barcode”
237 library of lineage-defining mutations to represent each SARS-CoV-2 lineage in the global
238 phylogeny²¹(**Figure 2A**). To encode each sample, Freyja stores the SNV frequencies (proportion
239 of reads at a site that contain the SNV) for each of the lineage-defining mutations (**Figure 2B,**
240 **top**). Since SNV frequencies at positions with greater sequencing depth more accurately estimate
241 the true mutation frequency, Freyja recovers relative lineage abundance by solving a depth-
242 weighted least absolute deviation regression problem, a mixed sample analog of minimizing the
243 edit distance between sequences and a reference (**Figure 2B, bottom**). To ensure results are
244 meaningful, Freyja constrains the solution space such that each lineage abundance value is non-
245 negative, and overall lineage abundance sums to one. Importantly, Freyja performs site-specific
246 weighting to account for non-constant variance in measured SNV frequency across sites,
247 enabling prioritization of information at each site as a function of sequencing depth. Read depths
248 are log-transformed, providing robustness to common attributes of real sequencing data such as
249 heavily skewed read depth across amplicons.

250
251 To validate Freyja, we sequenced “spike-in” synthetic mixtures from five key SARS-CoV-2
252 lineages (Lineage A, Beta, Delta, Epsilon, and Gamma) at proportions ranging from 5% to 100%
253 in each sample, with between 1 and 5 different lineages per mixture (**Figure 2C**, and see **Table**
254 **1**). We found that Freyja robustly recovered the expected lineage abundances for all mixtures,
255 even for lineages at 5% abundance (**Figure 2D**, and see **Extended Data Figure 3** for lineage
256 specific predictions). To further validate Freyja, we used wastewater samples from the UCSD
257 isolation dorms as well as Point Loma wastewater treatment plant, collection sites likely to
258 contain mixed-lineage samples, to compare Freyja-detected lineages with qPCR testing for 8
259 mutations associated with different variants of concern (N501Y, DelHV69/70, DelY144, K417N,
260 K417T, E484Q, P681R and L452R, **Figure 2E**). We found that Freyja consistently identified the
261 same lineages as qPCR testing, but, as expected, also identified additional lineages with SNVs
262 not included in our qPCR panel that were known to be circulating in San Diego at the time of
263 collection. Combined, these results show that Freyja robustly estimates viral lineage abundance
264 from samples containing a mixture of lineages, including synthetic virus mixtures and field
265 wastewater collections.

266

267 To compare Freyja with other wastewater analysis pipelines, we tested the performance of other
 268 wastewater deconvolution methods including the method from Baaijens et al.¹², cojac²², and LCS
 269 ²³ using the spike-in mixtures (**Extended Data Figure 4**). We found that Freyja greatly
 270 outperforms other methods in terms of accuracy, false positive rate, and computational
 271 efficiency. The method from Baaijens et al. required greater than ten times more computation
 272 time per sample relative to Freyja (~13.2 minutes vs ~1.1 minutes per sample, respectively).
 273 Although cojac was fast, the small amplicon length used for the spike-in mixtures caused cojac
 274 to fail to identify most of the variants entirely, while LCS failed to return estimates within two
 275 days.
 276



277
 278
 279 **Figure 2: Sample deconvolution robustly recovers relative virus abundance.** A. Subset of
 280 lineage defining mutation “barcode” matrix. Each row represents one lineage (out of >1000
 281 lineages included in the UShER global phylogenetic tree), and individual nucleotide mutations
 282 are represented as columns. B. Single nucleotide variant frequencies obtained from iVar used for
 283 recovering relative abundance of each lineage. C. Schematic of the spike-in validation
 284 experiment. D. Depth-weighted de-mixing estimates of the virus abundance versus
 285 expected/known abundance. Details on lineage specific predictions are provided in **Extended**
 286 **Data Figure 3**. E. Comparison of wastewater sample deconvolution with VOC qPCR panel, with
 287 lookup table (bottom) showing amino acid mutations corresponding to each variant.
 288

289 Detection of early and cryptic community transmission in wastewater

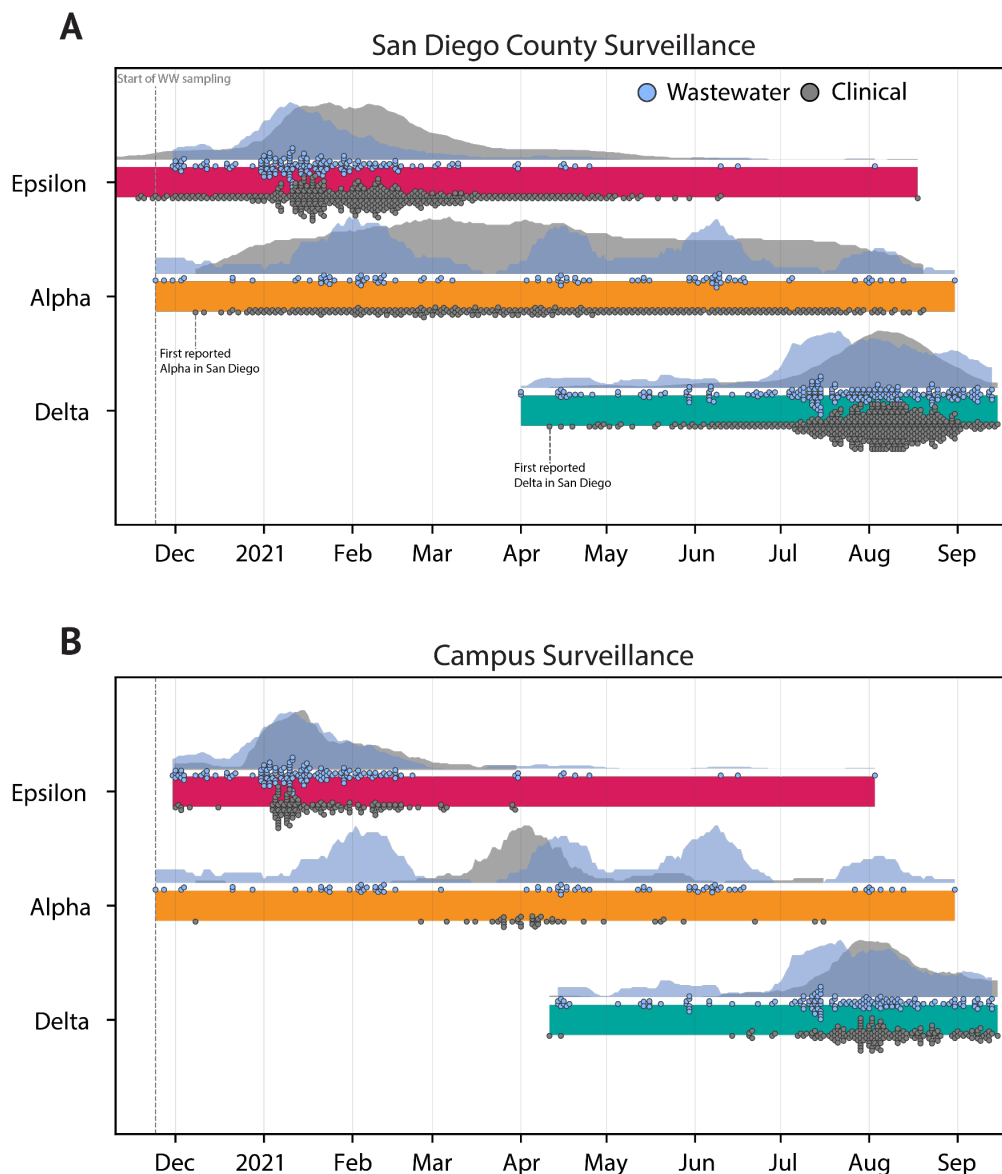
290
 291 SARS-CoV-2 RNA concentrations in wastewater have been shown to be an early indicator of
 292 rising COVID-19 community incidence^{9,24} (and see **Extended Data Figure 5A**), but whether

293 wastewater can be used to detect emerging variants, including VOCs and VOIs, prior to their
294 observation in clinical surveillance is unknown. To test if wastewater can enable early detection
295 of emerging lineages, we applied Freyja to our wastewater sequencing data and compared the
296 collection date of VOC positive samples from wastewater with the collection dates of samples
297 from clinical genomic surveillance (**Figure 3A**). With only 2.6% as many sequenced wastewater
298 samples as sequenced clinical samples, we detected the Alpha and Delta VOC lineages in
299 wastewater genomic surveillance up to 14 days prior to their first detection in genomic clinical
300 surveillance (Epsilon was circulating at the start of wastewater collection, and thus could not be
301 detected early). To further quantify our uncertainty in prevalence estimates, we used a fast
302 bootstrapping approach (**Extended Data Figure 6**) and found that the resampled distributions
303 did not include zero abundance. Since emerging VOC lineages may evade immune responses or
304 lessen the effectiveness of public health interventions¹⁸, this early detection provides additional
305 time to make necessary adjustments to existing countermeasures.

306
307 To test if wastewater genomic surveillance can identify changes in the abundance of circulating
308 lineages, we compared VOC detection rates in clinical and wastewater sequencing over time. We
309 found that both wastewater and clinical genomic surveillance tracked changes in lineage
310 abundance, but increases in lineage detection frequency were generally observed first in
311 wastewater surveillance. For example, for the Epsilon variant, which was first detected in San
312 Diego in September of 2020, we observed increases in detection frequency in wastewater
313 approximately 5 days prior to the corresponding increase in clinical genomic surveillance data
314 (**Figure 3A**, see **Methods**). We noticed varying periods of ongoing lineage detection across
315 VOCs relative to clinical surveillance, possibly due to different virus shedding characteristics
316 across lineages²⁵. For Epsilon specifically, elevated sampling density on the UCSD campus
317 relative to elsewhere in the county early on in the experiment may have biased San Diego wide
318 detection trends towards campus trends, particularly during the end of the wave. We also
319 observed clear signatures of times with elevated travel, as seen in the pulsing of Alpha detections
320 in wastewater around the end of holidays and school breaks. During these periods as well as
321 other times of mass student arrival, students were mandated to test immediately upon arrival
322 before they moved into their respective on-campus housing. In late March of 2021 following the
323 university break, mandated clinical testing identified spread of the Alpha variant exclusively in
324 off-campus residents (see **Figure 1B**), suggesting that campus mitigation protocols kept the
325 Alpha outbreak from spreading on campus during this period.

326
327 To study the effectiveness of wastewater genomic surveillance at a smaller community scale, we
328 restricted our analysis to samples from the UCSD campus. We found that wastewater genomic
329 surveillance consistently identified the three major VOCs (Epsilon, Alpha, and Delta) throughout
330 their period of occurrence, despite detection gaps of one month or longer in clinical surveillance
331 that included regular asymptomatic testing, longer than the expected signal due to extended virus
332 shedding²⁶⁻²⁸ (**Figure 3B**). During these gaps, positive samples were collected from multiple
333 distinct locations, with most locations not repeated, suggesting that this continued detection in
334 wastewater was not simply due to extended shedding. From mid-December to late-March, the
335 Alpha variant was detected more than once per week on average in wastewater but was not
336 detected by clinical surveillance. Similarly, wastewater surveillance detected continued Delta
337 transmission from mid-April to mid-June, but no cases were identified by clinical surveillance.
338 This explains in part the long tails of wastewater positivity on campus relative to clinical

339 surveillance on campus (**Figure 1B**), in which we control for extended shedding by excluding
340 samples from campus isolation dorms (see **Methods** for details). The high wastewater positivity
341 level in February-March 2020 extends beyond the expected duration of extended shedding,
342 indicating that cryptic transmission likely played a significant role in campus virus spread during
343 this period.
344

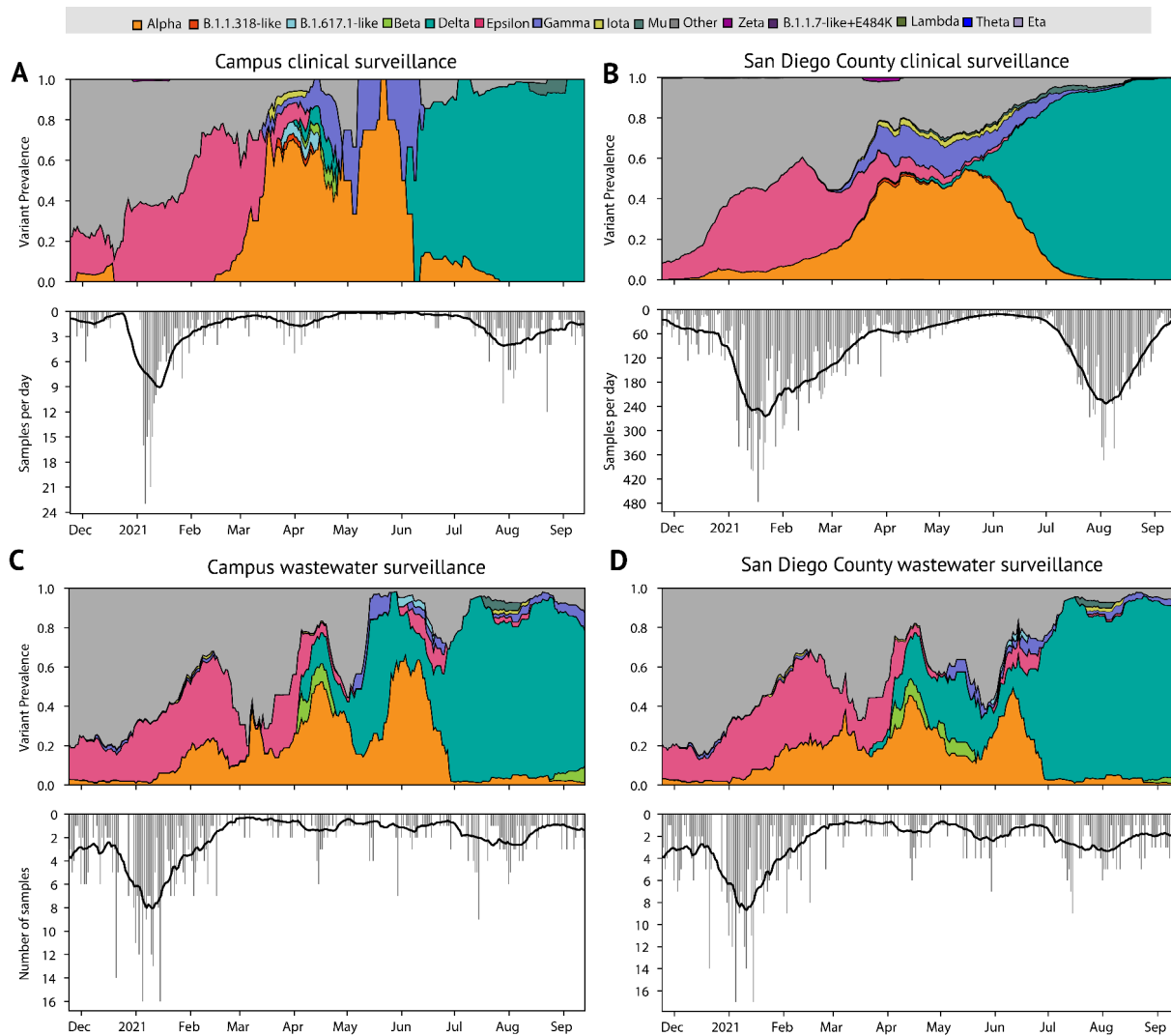


345
346
347 **Figure 3: Freyja recovers early and cryptic transmission of SARS-CoV-2 variants of**
348 **concern** A. Timeline and normalized epidemiological curves for VOC detection in both
349 wastewater and clinical sequences from San Diego County for the 3 major VOCs in circulation
350 during the sampling period. Both Alpha and Delta are detected first in wastewater before clinical
351 samples. Markers for clinical detections correspond to the ceiling of the daily detection count
352 divided by 30 (e.g. 1-30 samples= one marker, 31-60 = two markers), while wastewater markers
353 correspond to a single detection. B. Timeline and epidemiological curves for VOC detection in
354 the campus samples. Markers correspond to a single detection event for both clinical and

355 wastewater surveillance. All wastewater detections correspond to an estimated VOC prevalence
356 of at least 10%.

357
358 To study the effectiveness of wastewater surveillance in detecting and tracking other emerging
359 variants, we aggregated all wastewater sequencing data to estimate the temporal profile of
360 community lineage prevalence. We found that estimates of lineage abundance using wastewater
361 enable early identification of other VOCs/VOIs, even for lineages that are rarely observed in
362 clinical surveillance (**Figure 4**). For example, we detected the Mu (B.1.621) variant via
363 wastewater genomic surveillance on July 27th, nearly four weeks prior to its first detection
364 through clinical genomic surveillance on campus, on August 23rd (**Figure 4A,C**). However,
365 despite persistent Mu detection in campus wastewater throughout July and early August, we did
366 not detect the Mu variant in clinical or wastewater genomic surveillance on campus in
367 September, suggesting that local community transmission did not continue.

368
369 To test if Freyja continues to provide representative estimates of lineage prevalence for mixtures
370 containing closely related lineages, we analyzed the rise of the Delta variant (B.1.617.2) and its
371 sublineages (AY.*) in San Diego, from June-September 2021 (**Extended Data Figure 5B,C**). At
372 both the UCSD campus and the Point Loma wastewater treatment plant, we identified the rapid
373 emergence of B.1.617.2 and its sublineages (AY.*), along with low but persistent levels of the
374 P.1 (Gamma) variant. The relative abundances of each of the variants were within 2-fold of
375 prevalence estimates observed in clinical nasal swab data, suggesting that Freyja effectively
376 identifies prevalence even for closely related lineages, both at the university and county-scale.
377



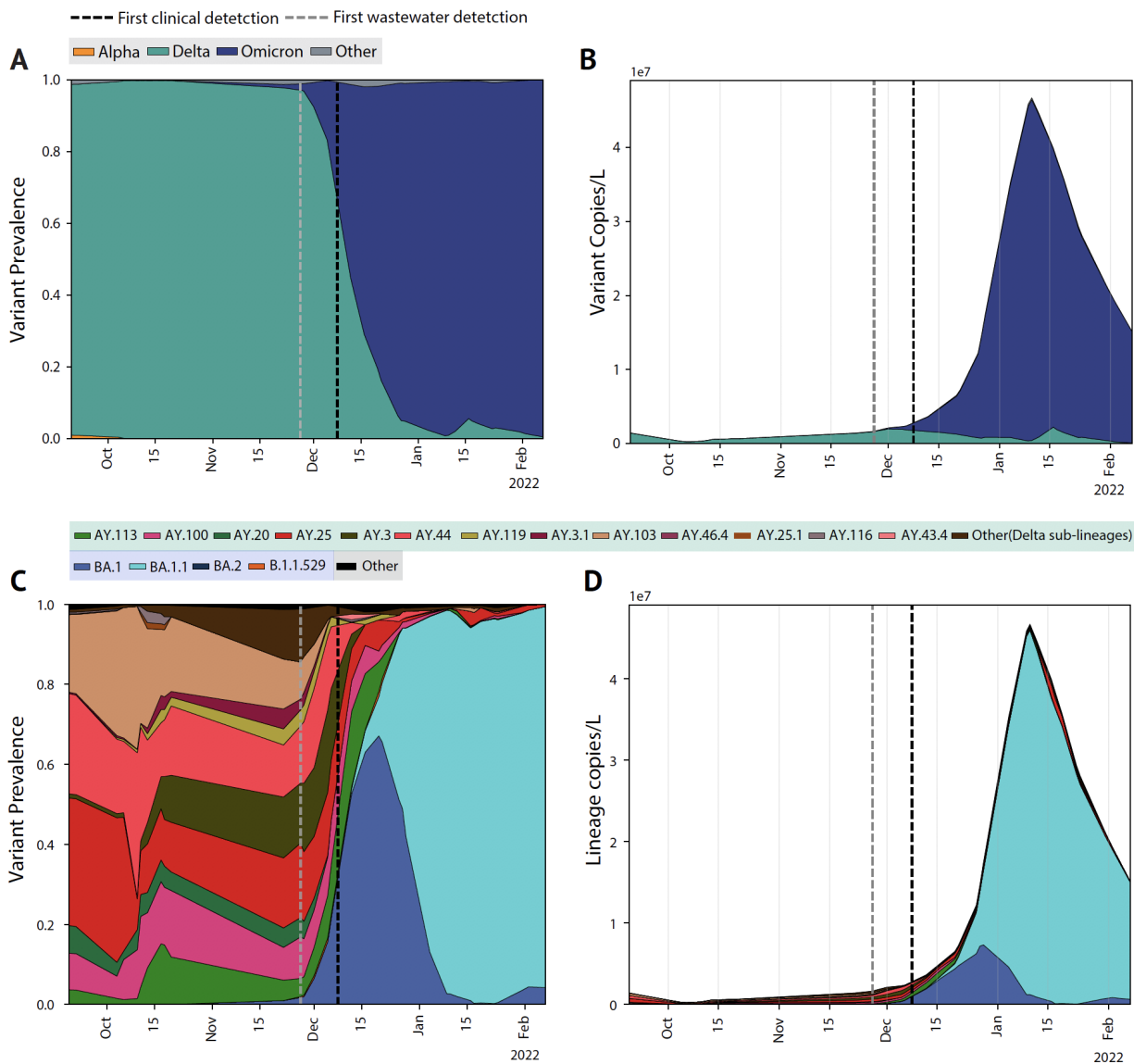
378
 379 **Figure 4: Deconvolution recovers a fine-grained estimate of virus population dynamics.** A.
 380 Prevalence of SARS-CoV-2 variants in UCSD clinical surveillance, and B. Variant prevalence in
 381 all clinical samples collected in San Diego County. C, D. Variant prevalence in wastewater at
 382 UCSD as well as the greater San Diego County (includes wastewater samples collected from
 383 Point Loma wastewater treatment plant as well as public schools in the San Diego districts).
 384 Further analysis of Point Loma wastewater samples is shown in **Extended Data Figure 5**. All
 385 curves show rolling average, window ± 10 days. “Other” contains all lineages not designated as
 386 VOCs. Bottom panels show number of sequenced samples per day.

387
 388 In more recent data from Point Loma wastewater treatment plant, we identified the Omicron
 389 variant (B.1.1.529 and descendants) at an abundance of near 1.7 % on November 27th, more than
 390 10 days prior to the first clinical detection in San Diego on December 8th (**Figure 5A-B**). To
 391 confirm these findings, we applied our VOC qPCR panel to the same samples and consistently
 392 detected two mutations associated with the Omicron variant (DelHV69/70 and N501Y) in
 393 samples detected after November 27th, while neither was detected in samples from earlier in
 394 November (**Extended Data Table 3**, P681R was included to confirm the presence of Delta).
 395

396 To visualize the dynamics of competition between the Delta and Omicron variants, we analyzed
397 wastewater collected at Point Loma from late September through early February. We found that
398 upon introduction to the community, Omicron rapidly rose to dominance and reached roughly
399 95% prevalence by December 26th. During the same period, the estimates for 95% Omicron
400 abundance in clinical samples tracked via S-gene target failures (SGTFs) was January 7th,
401 further suggesting wastewater genomic surveillance is a leading indicator of lineage dynamics
402 for emerging variants (**Figure 5A, Extended Data Figure 7**). To understand the magnitude of
403 lineage abundance, we scaled each sample by the measured virus RNA concentration of the
404 sample (**Figure 5B**). We observed that the absolute amount of circulating Delta variant remained
405 largely constant upon the introduction of Omicron, even as it appeared to decrease to a small
406 fraction of all viruses in the community.

407
408 To study the contribution of individual virus lineages to virus RNA concentration, we further
409 analyzed the growth dynamics of Delta and Omicron sub-lineages (**Figure 5C-D**). We found that
410 the many Delta lineages circulating in October and November were rapidly displaced by the
411 BA.1 Omicron lineage, which was soon after displaced by the BA.1.1 lineage, suggesting a
412 significant growth advantage over BA.1 and B.1.1.529. We did not observe significant levels of
413 any other Omicron sublineages.

414
415



416
 417 **Figure 5: Community wastewater enables early Omicron detection and reveals lineage**
 418 **dynamics.** A. Prevalence of SARS-CoV-2 VOCs in wastewater collected from the Point Loma
 419 wastewater treatment plant from late September 2021 to early February 2022. B. Estimated VOC
 420 concentrations, prevalence estimates scaled by normalized viral load in wastewater. C,D.
 421 Lineage-specific estimates of prevalence and concentration. All curves show an adaptive rolling
 422 average calculated using a local linear approximation (Savitzky-Golay filter) of virus copies/L,
 423 with window size ± 1 sampling date.

424
 425 **Wastewater identifies both known and unknown history of campus infections**

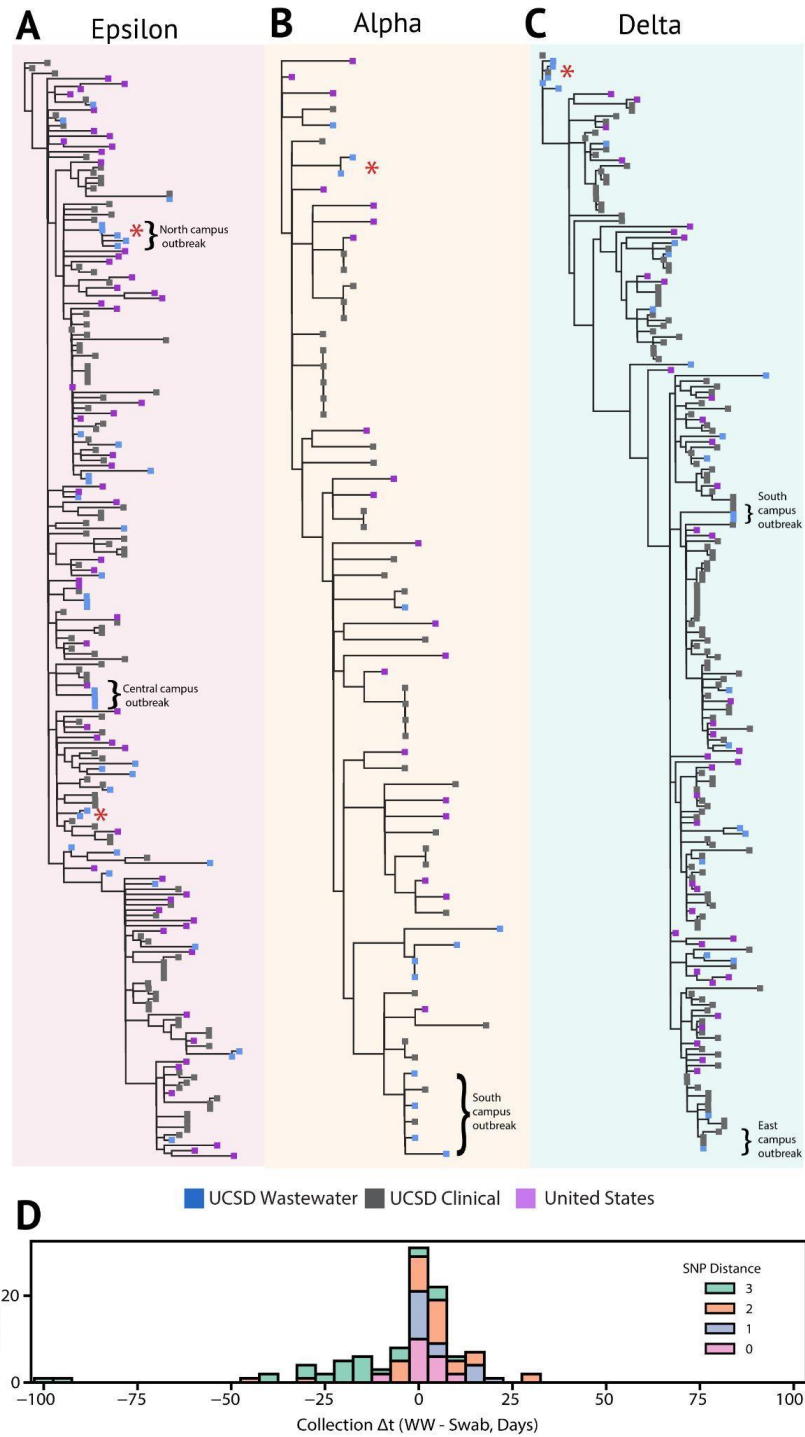
426
 427 Phylogenetic analysis of virus genomes can be used to identify fine-scale spatial and temporal
 428 transmission networks, but it is unknown if wastewater can be used to further refine possible
 429 sites of transmission, elucidate transmission networks (“who-infected-whom”), or identify
 430 specific infected individuals¹⁹. To investigate the scale, structure, and timing of SARS-CoV-2
 431 spread on campus, we reconstructed a maximum likelihood phylogenetic tree for each of the

432 major VOCs using all high-quality consensus genomes (see **Methods** for details) obtained from
433 the UCSD campus, as well as reference sequences for each lineage obtained elsewhere in the
434 United States (**Figure 6A-C**). In each tree, we identified many independent introductions, some
435 of which led to extended transmission on campus. The resulting virus diversity among the VOCs
436 present on campus enables ruling out of most transmission links and suggests campus virus
437 spread consisted of many separate, small outbreaks.

438
439 To analyze the spatial structure of virus spread, we identified collection sites for wastewater
440 sequences connected to transmission chains on campus, with building-specific resolution
441 (**Figure 6 A-C**). We observed multiple small, linked outbreaks clustered in nearby buildings.
442 Campus isolation protocol required students in congregate living to relocate to an isolation room
443 and linkages in the wastewater samples from buildings used for isolation reflected this co-
444 location. We also found multiple instances of successive exactly matching sequences from
445 wastewater collected from a single building, possibly due to continued viral shedding from the
446 same infected individuals from extended shedding in stool²⁶⁻²⁸ or a transmission chain in the
447 building leading to multiple infections by genetically identical viruses.

448
449 To study the temporal delay between clinical and wastewater lineage detection, we compared
450 collection times of sequences from campus wastewater that match sequences from campus
451 clinical surveillance (including non-VOC lineages). We found 20 exact sequence matches and
452 103 near-matches (SNP distance of 3 or less) but did not observe any overall bias towards earlier
453 or later detection in wastewater (**Figure 6D**), suggesting that on average, wastewater and clinical
454 genomic surveillance identify a similar timing of individual detection events. However, despite
455 current technical difficulties with isolating haplotypes from diverse virus mixtures, more than
456 half of the clinical-wastewater sequence pairs demonstrate earlier detection in wastewater or are
457 from the same date. Importantly, since detection is often delayed or missed by clinical
458 surveillance, detections occur first in wastewater (despite a loss of sequences due to limited
459 haplotype recovery), further suggesting that wastewater genomic surveillance can reveal the
460 presence of specific genome sequences prior to clinical surveillance.

461
462



463
 464 **Figure 6: Wastewater identifies clinically known and unknown virus transmission.** A-C.
 465 Maximum likelihood phylogenetic trees for each of the dominant variants of concern using high
 466 quality samples obtained at UCSD, as well as a representative set of sequences from the entire
 467 United States. Wastewater sequences from the same sampler that differ by 1 or fewer SNPs are
 468 denoted with a red asterisk. For all sequences, consensus bases were called at sites with >50%
 469 nucleotide frequency. Location information is provided for select outbreaks. D. Pairwise
 470 comparison of collection date for matching and near-matching wastewater and nasal swab

471 samples obtained at UCSD. Positive values indicate earlier collection in nasal swabs, and
472 negative values indicate earlier detection in wastewater.

473

474 **Discussion**

475

476 We show that improved virus concentration from wastewater, coupled with a method for
477 resolving multiple lineages from mixed samples, captures community virus lineage prevalence
478 and enables early detection of emerging variants, often before observation in clinical
479 surveillance. By sequencing both clinical and wastewater samples from the UCSD campus, we
480 detect VOCs persistently in wastewater even when their appearance in clinical samples is
481 intermittent. However, we also found occasions when rarer lineages, like B.1.1.318, were
482 detected in clinical samples but not in wastewater. This is not unexpected on campus since many
483 students living off-campus did not contribute to campus wastewater but were still clinically
484 tested as part of testing mandates and policies. In the larger San Diego community context, this
485 suggests that we may not be able to identify lineages circulating at low prevalence (< 1%) using
486 a single wastewater collection site. In addition, we note that clinical sequences identified from
487 the community may not be observable in the contributing catchment, as precise geolocation of all
488 clinical samples was not possible. On the other hand, we also observed rare lineages in
489 wastewater not seen in clinical samples from campus or the community. Since campus testing
490 mandates are unable to capture all cases (e.g. fully vaccinated individuals were not required to
491 test and not all community samples were sequenced), rare lineages can be missed.

492

493 The considerable benefits of wastewater surveillance may stem from biases in clinical testing,
494 including population testing availability and compliance, university quarantine policies, and
495 asymptomatic transmission, which may distort estimates of virus lineage prevalence from
496 clinical samples. Wastewater offers less biased and more consistent viral lineage prevalence
497 estimates, especially in areas with limited access and/or higher testing hesitancy rates, where
498 limited clinical surveillance can delay detection of emerging variants. Since it requires
499 considerably fewer samples, it is also more cost-effective than clinical testing, and could serve as
500 a long-term passive surveillance tool. This is particularly important for developing public health
501 interventions in low-resource and underserved communities, where widespread clinical genomic
502 surveillance for SARS-CoV-2 remains limited.

503

504 Wastewater is an information-dense resource for estimating the prevalence of specific viral
505 lineages, providing a community wide-snapshot not only of overall infection dynamics but of the
506 rise and fall of specific VOCs. Our method, Freyja, deconvolutes these information-rich mixtures
507 of virus lineages. For a large catchment area, such as San Diego's Point Loma wastewater
508 treatment plant, which covers over 2 million residents, even limited sampling may accurately
509 estimate lineage prevalence in the population and provide an early warning indicator of the rise
510 of new VOCs (as evidenced by the detection of Omicron at just over 1% abundance 11 days
511 ahead of the first local clinical observation). In addition, wastewater genomic surveillance with
512 building-level resolution provides a detailed description of the structure and dynamics of
513 community virus transmission, and can identify transmission links. It can be used to better direct
514 public health interventions, and can do so in real-time when combined with fast-turnaround
515 sequencing technologies. This high-resolution approach is of particular utility in community
516 gathering and transit sites, such as schools and airports, as well as sites with highly vulnerable

517 individuals, such as nursing homes and hospitals, where spatially resolved monitoring for
518 directing public health interventions is of great importance.

519
520 As SARS-CoV-2 continues to evolve, the risk of new VOCs remains high and there is a growing
521 need to identify these viruses ahead of their proliferation in the community. Accordingly,
522 development of technologies that are cost-effective, reduce biases, and provide leading rather
523 than trailing indicators of infection are essential to removing “blind spots” in our understanding
524 of local virus dynamics. Although technical issues have made wastewater sequencing difficult to
525 perform at scale, our key advances in virus concentration and sample deconvolution provide
526 evidence that this approach is now viable. Continued improvements to sequencing turnaround
527 speeds, lineage barcoding, and haplotype recovery from mixed samples will further accelerate
528 efforts to achieve earlier identification of emerging variants and improve the precision and
529 effectiveness of interventions.

530

531 **Methods**

532

533 **Wastewater sampling**

534

535 *High-resolution spatial sampling at the campus level*

536 131 wastewater autosamplers collecting 24h time-weighted composites were deployed across
537 manholes or sewer cleanouts of 360 campus buildings. GIS (geographic information systems)
538 informed analyses as well as agent-based network modeling of SARS-CoV-2 transmission on the
539 UCSD campus enabled identification of most optimal locations for wastewater sampling. During
540 the pilot phase (November 23-Dec 29th 2020), 68 samplers were prioritized to cover 239
541 residential buildings identified as the highest risk areas for large outbreaks on campus as a part of
542 an observational study of wastewater monitoring in high-density buildings²⁹. This was based on
543 preliminary dynamic modeling which showed the largest potential outbreaks to occur within the
544 largest residential buildings⁹. In addition to the observational study of wastewater monitoring in
545 these high-density buildings, a cluster randomized study was also performed concurrently. This
546 included a randomized modified version of a stepped wedge crossover design, in which there
547 was random assignment of manholes for wastewater sampling. Clusters of manholes associated
548 with residential buildings were randomized to receive wastewater monitors at one of two-time
549 steps to evaluate the impact of wastewater monitoring on outbreak size in the associated
550 buildings. During the same time period, all students in these residences were mandated to
551 undergo weekly diagnostic testing which was used to validate the utility of building-level
552 wastewater monitoring. Furthermore, on-campus residences were initially focused due to the
553 relatively static nature of the population which enabled a more robust cross-validation of the
554 sensitivity and efficacy of the wastewater surveillance. The coverage of wastewater surveillance
555 was then increased to cover the rest of the campus buildings (including non-residential buildings
556 on campus) from January 2021. Four of the deployed wastewater samplers covered the
557 designated isolation and quarantine buildings on campus.

558

559 Wastewater composites were collected from the 131 samplers every day for the on-campus
560 residence buildings and Monday through Friday for the nonresidential campus buildings. 19,944
561 wastewater samples were collected and analyzed for the presence of SARS-CoV-2 RNA via RT-
562 qPCR between November 23rd 2020 and September 20th 2021. During this time, 9700 students

563 lived in campus residences and 25,000 worked on campus on a daily basis. Between October
564 2020 to January 1st 2021, all on-campus residents were mandated to test on a bi-weekly (once
565 every 2 weeks) basis and on a weekly basis from January 2nd 2021 (start of the Winter term).
566 However, fully vaccinated individuals were not mandated to test on a regular basis. Campus
567 protocols required SARS-CoV-2 positive students living in congregate housing to relocate to
568 designated isolation housing. Accordingly, our analysis of wastewater positivity (**Figure 1B**) did
569 not include isolation housing samplers, in order to control - as best as possible, a small number
570 of students in non-congregate housing spaces were allowed to isolate “in-place”, for example -
571 for possible repeat detection due to extended shedding from infected individuals. Automated,
572 localized wastewater-triggered notifications were sent to the residents/employees of buildings
573 associated with a positive wastewater signal which further led to a surge in testing uptake rates
574 by 2 to 40-fold in the associated buildings.

575

576 *Wastewater sampling at the county level*

577 24h flow-weighted composites were collected thrice a week from the main pump station for the
578 Point Loma wastewater treatment plant, the primary treatment plant serving the greater San
579 Diego county with a catchment size of approximately 2.3 million. 132 wastewater samples were
580 collected between February 24th 2021 to February 7th, 2022.

581

582 **Wastewater sample processing and viral genome sequencing**

583

584 *Sample processing*

585 SARS-CoV-2 RNA was concentrated from 10ml of raw sewage and processed as described
586 elsewhere⁷. In brief, the viral RNA was concentrated using an automated affinity capture
587 magnetic hydrogel particle (Ceres Nanosciences Inc., USA) based concentration method after
588 which the nucleic acid was extracted and sample eluted in 50uL of elution buffer. The extracted
589 RNA was then screened for SARS-CoV-2 RNA via real-time RT-qPCR for 3 gene targets (N1,
590 N2 and E-gene). PMMoV (pepper mild mottle virus) was also screened to adjust for changes in
591 load. Positive wastewater samples were sequenced within 1-2 weeks of collection, comparable to
592 the delay for clinical samples. To cross-validate the ability of the deconvolution tool in reliably
593 resolving mixtures of strains in wastewater, the wastewater samples from the county as well as
594 the ones from the isolation dorms on campus (where multiple infected individuals were isolating)
595 were also run through a PCR panel targeting 8 mutations associated with the strains designated
596 as VOCs. The mutations screened for in wastewater using RT-qPCR included N501Y,
597 DelHV69/70, DelY144, K417N, K417T, E484Q, P681R and L452R (Promega Corp. Cat#
598 CS3174B02).

599

600 *Miniaturized wastewater SARS-CoV-2 amplicon sequencing*

601 The Swift Normalase® Amplicon Panels (SNAP) kit (PN: SN-5X296 (core) COVG1V2-96
602 (amplicon primers), Integrated DNA Technologies, Coralville, IA) was used on RNA from
603 wastewater samples that were positive for SARS-CoV-2 RNA to prepare the multiplex NGS
604 amplicon libraries and indexed using the SN91384 series of dual indexing oligos, yielding up to
605 1536 index pairs per pool. A miniaturized version of the protocol was used with the following
606 modifications: the Superscript IV VILO (Thermo Fisher, Carlsbad, CA) cDNA synthesis
607 reaction was scaled down to ~1/12 the normal reaction volume with 0.333uL of enzyme mix and
608 1.333uL of RNA being used. The multiplex amplicon amplification and Ampure XP bead

609 purification steps were scaled down ~1/6 the normal reaction volume. The Index adapter PCR
610 reaction and Ampure XP bead purification steps were scaled down to ~2/13 the normal reaction
611 volume. The final library resuspension volume was 29uL. 1uL of each library was pooled for an
612 initial shallow NGS run on a MiSeq (Illumina, San Diego, CA) using a Nano flow cell. This
613 equal volume pool was used to estimate the differential volumes required for similar read depths
614 across samples using a NovaSeq SP or S4 flow cell (Illumina, San Diego, CA). Between 5uL and
615 0.2uL of library material, depending on the data provided from the MiSeq Nano run, was
616 pipetted into a single pool for the NovaSeq run. Transfer volumes were capped at 5uL to reduce
617 pipetting time and because these types of “high volume” samples typically contained a higher
618 proportion of likely adapter dimers that inhibit flow cell performance for all samples. A
619 Dragonfly Discovery (SPT Labtech, UK) was used to dispense reaction master mixes or water
620 depending on the step. A BlueWasher (BlueCatBio, MA) was used for high throughput
621 centrifugal 384-well plate washing during the AmpureXP bead reaction cleanup steps. An IKA
622 MS3 Control linear plate mixer (IKA Works Inc, Wilmington, NC) set to 2600 RPM for 5’ was
623 used to resuspend the AmpureXP beads during the rehydration steps. A Mosquito Genomics HV
624 16 channel robotic liquid handler (SPT Labtech, UK) was used to dispense the RNA, the reaction
625 master mixes, and prepare the equal volume pools for the initial MiSeq Nano (Illumina, San
626 Diego, CA) balancing runs. A Mosquito X1 single channel “hit picker” robotic liquid handler
627 (SPT Labtech, UK) was used for the final library balancing for the NovaSeq (Illumina, San
628 Diego, CA) NGS lanes.

629
630 Sequencing data were analyzed using the C-VIEW (COVID-19 Viral Epidemiology Workflow)
631 platform for initial QC and SARS-CoV-2 lineage assignment and phylogenetics. In brief,
632 sequencing reads are aligned with minimap2³⁰, and primer sequences trimming and quality
633 filtering is applied using the iVar trim method²⁰. Sequencing depth and single nucleotide variant
634 (SNV) calls are obtained using samtools mpileup³¹ and the iVar variants method²⁰.

635
636 Controls were included at all stages of sample processing (viral concentration, extraction, qPCR
637 and sequencing) to assess potential inhibition and cross-contamination. Most of the sample
638 processing steps were performed by liquid handling robots for consistency and to minimize
639 human error. Replicates were included for all wastewater samples. If any of the controls failed or
640 indicated cross-contamination, the entire batch was rerun. The clinical samples and wastewater
641 samples were processed separately for sequencing due to significant differences in viral load
642 between the two sample types.

643 644 **Virus diversity**

645 As reported previously²⁰, virus SNVs were used to characterize the populations derived from
646 wastewater and clinical samples. Richness was defined as the total number of SNV sites, and
647 mean Shannon entropy $H(p)$ was defined as

$$648 \quad H(p) = \frac{1}{N} \sum_{i=1}^N -p_i \log_2(p_i) - (1 - p_i) \log_2(1 - p_i).$$

649 where p_i is the SNV frequency of at the i -th site, of N total sites. For statistical testing, a Mann-
650 Whitney U test was performed using all wastewater samples that were not sampled from the
651 same source within a 10 day period in order to ensure independence across samples, as well as
652 all clinical samples. Effect size was calculated using the rank-biserial correlation,
653

654 $r = 2U/(n_{WW}n_{CS}) - 1$ where U is the Mann-Whitney test statistic and n_{WW} and n_{CS} are the
655 numbers of wastewater and clinical samples, respectively.

656

657 **Wastewater sample deconvolution**

658

659 To infer relative abundance within a wastewater sample, we used a “barcode” matrix containing
660 the lineage defining mutations for each known virus lineage,

661

$$A = \begin{bmatrix} a_{1,1} & \dots & a_{1,N} \\ \vdots & \ddots & \vdots \\ a_{M,1} & \dots & a_{M,N} \end{bmatrix}$$

662

663

664 where $a_{i,j}$ denotes the i -th lineage, at mutation j . Lineage defining mutations were obtained from
665 the USHER global phylogenetic tree using the matUtils package¹⁵. Similarly, we let b and d
666 encode the frequency of each mutation and the corresponding sequencing depth (using the log-
667 transform $d_i = \log_2(\text{depth}_i + 1)$) to adjust for large differences in depth across amplicons, which
668 we use to control for heteroskedasticity and down-weight the importance of sites with little or no
669 sequencing depth),

670

$$b = \begin{bmatrix} b_1 \\ \vdots \\ b_N \end{bmatrix}, d = \begin{bmatrix} d_1 \\ \vdots \\ d_N \end{bmatrix}.$$

671

672

673 We were then able to write this as a constrained (weighted) least absolute deviations problem

$$\hat{x} = \underset{x \geq 0}{\operatorname{argmin}} \sum_{x=1} \|A^T x - b\|_{1W}, \quad \text{where} \quad \|\mu\|_{1W} = \sum_{i=1}^N d_i |\mu_i|$$

674

675 which yields the “demixing” vector $\hat{x} = [\hat{x}_1 \dots \hat{x}_M]$ that specifies the relative abundances of
676 each of the known haplotypes. Analysis was only performed on samples with greater than 70%
677 coverage, with the exception of March samples from UCSD for which all samples with greater
678 than 50% coverage were used. Constrained minimization was performed in Python using the
679 cvxpy convex optimization package^{32,33}. Mapping of lineages to variant WHO lineages (VOCs,
680 VUMs, etc.) was performed using curated lineage data from outbreak.info¹. We note that the
681 Epsilon variant received different maximum escalation levels at CDC and WHO, which assigned
682 VOC and VOI status, respectively. Since the Epsilon variant was widespread in California and
683 much of the United States, we use the more “local” CDC designation.

684

685 **Fast-bootstrapping method**

686 Bootstrapping was performed at the nucleotide level by resampling each site based on a
687 multinomial distribution of read depth across all sites, where the event probabilities are
688 determined by the fraction of the total sample reads found at each site, followed by a secondary
689 resampling at each site according to a multinomial distribution (i.e. binomial when there was
690 only one SNV at a site), where event probabilities were determined by the frequencies of each
691 base at the site, and the number of trials is given by the sequencing depth. 1000 resamplings and
692 demixings were performed for all samples.

693

694 **Spike-in mixture experiment**

695 RNA was isolated from supernatants of a mammalian cell culture infected with one of five
696 strains of SARS-CoV-2. (A, B.1.1.7, B.1.351, P.1, or B.1.617.2).

697

698 *RNA concentration standardization*

699 Virus concentration was quantified by the UCSD EXCITE COVID testing laboratory using the
700 Thermo COVID-19 Test kit (PN:A47814, Thermo Scientific Corporation, Carlsbad, CA). The
701 median Cq values (N-gene, Orflab, & S-gene (where applicable)) was calculated and used to
702 determine how much the RNA needed to be diluted with water to reach a Cq value of 23. A post
703 dilution RT-qPCR reaction was performed and used to calculate the final dilution of the more
704 concentrated samples to the new target value of Cq 23.296. The number of freeze thaw cycles
705 between RNA samples was kept the same.

706

707 *Virus Mixing*

708 RNA standardized in the prior section was used to make a volumetric mixing array (final volume
709 10uL) using a Mosquito X1 HV robotic liquid handler (SPT Labtech, UK). Pairwise mixes of
710 5:95, 10:90, 20:80, 60:40, and 50:50 were made for each virus lineage and in both directions.
711 Equal mixes (20%) for each of the five test strains were made. 25% mixes and 33% mixes were
712 made for a subset of possible combinations and controls of 100:0 were prepared. See **Extended**
713 **Data Table 1** for complete array. Corrected estimates of the fraction of each virus lineage based
714 were performed using the final measured Cq values for each pure virus lineage sample to control
715 for issues encountered during the dilution step (repeat Cq measurements had a coefficient of
716 variation of 0.007, **Extended Data Table 2**). Across all 95 mixtures, we observed a coefficient
717 of variation of 0.016. Since initial virus concentrations are controlled for using measured Cq
718 values, we expect remaining lineage specific bias (see **Extended Data Figure 3**) is likely due to
719 experimental inconsistencies encountered during mixture creation.

720

721 **Deconvolution method performance comparison**

722 A subset of the spike-in mixtures (1 of each type, for a total of 95 mixtures) was used to compare
723 Freyja, cojac (using VOC definitions from the public cojac github repository; Lineage A and
724 Epsilon definitions were created manually), the Kallisto-based method from Baaijens et al. 2021,
725 and LCS. Kallisto was run using 10 cores (with no bootstrapping), and LCS was run using 16
726 cores, both on an Intel Xeon processor (2.2GHz). LCS was run for 48 hours, but failed to
727 complete. Timing was performed using the “time” command, and included all steps after
728 alignment, trimming, and sorting. Times correspond to total CPU time.

729

730

731 **Estimation of delay in detection frequency**

732 Estimation of the lag time between epidemiological curves for wastewater and clinical
733 surveillance of the Epsilon variant in San Diego was performed by identifying the shift with
734 maximal cross-correlation. All time points leading up to the time of initial peak in detection
735 frequency were included for both wastewater and clinical data.

736

737 **Phylogenetic analyses**

738 Reconstruction of maximum likelihood trees was performed on all SARS-CoV-2 VOC genomes
739 with 10x (10 reads or greater per site) genome coverage >95% and quality score >20 obtained
740 from UCSD campus sampling, using IQtree³⁴. This analysis included 150 (112 clinical, 38
741 wastewater) Epsilon, 49 (37 clinical, 12 wastewater) Alpha, and 160 (136 clinical, 24
742 wastewater) Delta lineage genomes from UCSD, in addition to 60 Epsilon, 20 Alpha, and 39
743 Delta randomly selected genomes from elsewhere in the United States. We used iVar²⁰ to
744 identify consensus sequences for all San Diego samples. Bases were only included in the
745 sequence if there was a consensus base at the site (>50% nucleotide frequency). We also masked
746 known homoplastic sites prior to tree reconstruction³⁵. Analysis of temporal comparison was
747 performed on 608 samples (443 clinical, 165 wastewater, all lineages were included) with 10x
748 genome coverage >95% and quality score >20 from UCSD. Sample collection SNP distances
749 were calculated without considering ambiguous bases and gaps.

750

751 **Code availability**

752 Freyja is hosted publicly on github (<https://github.com/andersen-lab/Freyja>) and is available
753 under a BSD-2-Clause License. Freyja is accessible as a package via bioconda
754 (<https://bioconda.github.io/recipes/freyja/README.html>) in container form via dockerhub
755 (<https://hub.docker.com/r/andersenlabapps/freyja>). COVID-19 Viral Epidemiology Workflow
756 (C-VIEW) is available at <https://github.com/ucsd-ccbb/C-VIEW> as an open-source, end-to-end
757 workflow for viral epidemiology focused on SARS-CoV-2 lineage assignment and
758 phylogenetics.

759

760 **Data Availability**

761 All raw wastewater sequencing data is available via the NCBI Sequence Read Archive under the
762 BioProject ID PRJNA819090. Consensus sequences from clinical and wastewater surveillance
763 are all available on GISAID. Spike-in sequencing data is available via google cloud
764 (https://console.cloud.google.com/storage/browser/search-reference_data). The UCSD campus
765 dashboard can be accessed at <https://returntolearn.ucsd.edu/dashboard/>. The county wastewater
766 data from Point Loma are available through the public dashboard that can be accessed at
767 <https://searchcovid.info/dashboards/wastewater-surveillance/>. The SEARCH genomic
768 surveillance dashboard is available at <https://searchcovid.info/dashboards/sequencing-statistics/>.

769

770 **Acknowledgements**

771 We thank Laralyn Asato and the Microbiology Lab at the SD Public Utilities Department for
772 providing us with county wastewater samples. We thank UC San Diego's Return to Learn (RTL)
773 program for funding the campus-wide wastewater surveillance efforts. We also thank Jason
774 Kayne, Rich Cota, Jesus Ortiz, and the Facilities management team (FM) at UCSD, Joseph
775 Mayer from the Center for Aerosol Impacts on Chemistry of the Environment (CAICE) and
776 Luke Arnold of the Campus Research Machine Shop (CRMS) for assistance with the installation
777 and operation of the autosamplers; Robbie Jacobs, Shawn Knepple and their team at UCSD
778 Logistics for assisting with our daily sampling efforts; Brett Pollak and the UCSD Information
779 Technology Services team for assisting with the daily notifications; Office of Academic Affairs
780 for contact tracing and targeted campus messaging assistance; Jana Severson, Patrick Hochstein,
781 the UCSD HDH team, and the UCSD Environmental Health and Safety (EHS) personnel; Jack
782 Gilbert and the Microbiome Sample Processing Core at UCSD for access to qPCR equipment.
783 We also thank the CDC SPHERES consortium, SEARCH (San Diego Epidemiology and

784 Research for COVID Health) Alliance, and members of the Andersen lab for discussion and help
785 with logistics. We thank the healthcare workers, frontline workers, and patients who made the
786 collection of this SARS-CoV-2 dataset possible and all those who made genomic data available
787 for analysis via GISAID.

788

789 **Funding**

790 This work has been funded by CDC BAA contracts 75D30121P10258 (Helix) and
791 75D30120C09795 (G.W.Y., R.K., L.C.L., and K.G.A.), NIH NIAID 3U19AI135995-03S2
792 (K.G.A.), U19AI135995 (K.G.A.), U01AI151812 (K.G.A.), NIH NCATS UL1TR002550
793 (K.G.A.), the Conrad Prebys Foundation (K.G.A.), NIH 5T32AI007244-38 (J.I.L.), NIH Pioneer
794 Grant 1DP1AT010885 (R.K), NSF RAPID 2029069 (R.K.), San Diego County Health and
795 Human Services Agency (R.F.M), NIH S10OD026929 (K.J.). The findings and conclusions in
796 this report are those of the authors and do not necessarily represent the official position of the
797 Centers for Disease Control and Prevention. Use of trade names is for identification only and
798 does not imply endorsement by the Centers for Disease Control and Prevention.

799

800 **Ethics declarations**

801 The University of California San Diego Institutional Review Boards (IRB) provided human
802 subject protection oversight of the of the data obtained by the EXCITE lab for the campus
803 clinical samples (IRB approval #210699, #200477). All necessary patient/participant consent has
804 been obtained and the appropriate institutional forms have been archived, and any sample
805 identifiers included were de-identified. The wastewater component of this project was discussed
806 with our Institutional Review Board, and was not deemed to be human subject research, as it did
807 not record personally identifiable information.

808

809 **Author Contributions**

810 Conceptualization: RK, KGA

811 Methodology: SK, JIL, NKM, PDH, AK, SS, KMF, AB, LCL, GWY, KGA, RK

812 Software: JIL, DM, NM, KMF, AB, BH, SS, KG, NLM, KSR, CMA, EH, AMM

813 Formal Analysis: SK, JIL

814 Investigation: SK, JIL, NM, SF, HMT, TV, CET, RT, NAB, TB, MC, WC, ESC, ERE, AH, GH,
815 ALL, EL, TTN, TO, AP, RAS, PS, PBF, EWS, SA, PDH, CAM, LCL, GWY, CA, EK, MAS,
816 SAP, JL, EP, MZ, ES, RFK, TG, RG, KGA, RK

817 Resources: CA, NKM, RMN, RS, EHS, AMS, SFK, DPD, CAH, AM, SS, BA, SS, NG, JDM,
818 EM, IAM, AH, OB, AM, AB, KMSB, ETC, NLW, WL, MI, DB, LN, SW, MZ, RRS, RFM, TG,
819 RG, DBA, DB, JCB, AB, ME, JG, SLG, MH, FKH, LI, HJ, TJ, VK, BK, LR, CAH, GCM, PM,
820 RM, PO, DP, AP, AMS, BS, AS, BW, TW, SR, PKK ATY, DMC

821 Data curation: SK, JIL, PDH, GH, SF, HMT, CET, RT, TV, PDH, AB, NM, AMM, KMF

822 Writing – Original Draft: SK, JIL, KGA, RK

823 Writing – Review and editing: all authors

824 Visualization: SK, JIL

825 Supervision: RMN, NKM, RS, ALS, EHS, AMS, PDH, LCL, DP, GWY, KGA, RK

826 Project administration: RMN, NKM, RS, ALS, EHS, AMS, PDH, LCL, GWY, KGA, RK

827 Funding acquisition: RK, KGA

References

1. Julia L. Mullen, Ginger Tsueng, Alaa Abdel Latif, Manar Alkuzweny, Marco Cano, Emily Haag, Jerry Zhou, Mark Zeller, Emory Hufbauer, Nate Matteson, Kristian G. Andersen, Chunlei Wu, Andrew I. Su, Karthik Gangavarapu, Laura D. Hughes, and the Center for Viral Systems Biology. outbreak.info. *outbreak.info* <https://outbreak.info/> (2021).
2. Harvey, W. T. *et al.* SARS-CoV-2 variants, spike mutations and immune escape. *Nat. Rev. Microbiol.* **19**, 409–424 (2021).
3. Reitsma, M. B. *et al.* Racial/Ethnic Disparities In COVID-19 Exposure Risk, Testing, And Cases At The Subcounty Level In California. *Health Aff.* **40**, 870–878 (2021).
4. Lieberman-Cribbin, W., Tuminello, S., Flores, R. M. & Taioli, E. Disparities in COVID-19 Testing and Positivity in New York City. *Am. J. Prev. Med.* **59**, 326–332 (2020).
5. Brito, A. F. *et al.* Global disparities in SARS-CoV-2 genomic surveillance. *medRxiv* (2021) doi:10.1101/2021.08.21.21262393.
6. Hata, A., Hara-Yamamura, H., Meuchi, Y., Imai, S. & Honda, R. Detection of SARS-CoV-2 in wastewater in Japan during a COVID-19 outbreak. *Sci. Total Environ.* **758**, 143578 (2021).
7. Karthikeyan, S. *et al.* High-Throughput Wastewater SARS-CoV-2 Detection Enables Forecasting of Community Infection Dynamics in San Diego County. *mSystems* **6**, (2021).
8. Randazzo, W. *et al.* SARS-CoV-2 RNA in wastewater anticipated COVID-19 occurrence in a low prevalence area. *Water Res.* **181**, 115942 (2020).
9. Karthikeyan, S. *et al.* Rapid, Large-Scale Wastewater Surveillance and Automated Reporting System Enable Early Detection of Nearly 85% of COVID-19 Cases on a University Campus. *mSystems* **6**, e0079321 (2021).
10. Mercer, T. R. & Salit, M. Testing at scale during the COVID-19 pandemic. *Nat. Rev. Genet.* **22**, 415–426 (2021).
11. Crits-Christoph, A. *et al.* Genome Sequencing of Sewage Detects Regionally Prevalent SARS-CoV-2 Variants. *MBio* **12**, (2021).

12. Baaijens, J. A. *et al.* Variant abundance estimation for SARS-CoV-2 in wastewater using RNA-Seq quantification. *medRxiv* (2021) doi:10.1101/2021.08.31.21262938.
13. Amman, F. *et al.* National-scale surveillance of emerging SARS-CoV-2 variants in wastewater. *medRxiv* 2022.01.14.21267633 (2022).
14. Rambaut, A. *et al.* A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol* **5**, 1403–1407 (2020).
15. Turakhia, Y. *et al.* Ultrafast Sample placement on Existing tRees (UShER) enables real-time phylogenetics for the SARS-CoV-2 pandemic. *Nat. Genet.* **53**, 809–816 (2021).
16. Walke, H. T., Honein, M. A. & Redfield, R. R. Preventing and Responding to COVID-19 on College Campuses. *JAMA* **324**, 1727–1728 (2020).
17. Fielding-Miller, R. K. *et al.* Wastewater and surface monitoring to detect COVID-19 in elementary school settings: The Safer at School Early Alert project. (2021) doi:10.1101/2021.10.19.21265226.
18. Ladner, J. T., Grubaugh, N. D., Pybus, O. G. & Andersen, K. G. Precision epidemiology for infectious disease control. *Nat. Med.* **25**, 206–211 (2019).
19. Grubaugh, N. D. *et al.* Tracking virus outbreaks in the twenty-first century. *Nat Microbiol* **4**, 10–19 (2019).
20. Grubaugh, N. D. *et al.* An amplicon-based sequencing framework for accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biol.* **20**, 8 (2019).
21. McBroome, J. *et al.* A daily-updated database and tools for comprehensive SARS-CoV-2 mutation-annotated trees. *Mol. Biol. Evol.* (2021) doi:10.1093/molbev/msab264.
22. Jahn, K. *et al.* Detection and surveillance of SARS-CoV-2 genomic variants in wastewater. *bioRxiv* (2021) doi:10.1101/2021.01.08.21249379.
23. Valieris, R. *et al.* A mixture model for determining SARS-Cov-2 variant composition in pooled samples. *Bioinformatics* (2022) doi:10.1093/bioinformatics/btac047.
24. Peccia, J. *et al.* Measurement of SARS-CoV-2 RNA in wastewater tracks community infection dynamics. *Nat. Biotechnol.* **38**, 1164–1167 (2020).

25. Singanayagam, A. *et al.* Community transmission and viral load kinetics of the SARS-CoV-2 delta (B.1.617.2) variant in vaccinated and unvaccinated individuals in the UK: a prospective, longitudinal, cohort study. *Lancet Infect. Dis.* **22**, 183–195 (2022).
26. Cevik, M. *et al.* SARS-CoV-2, SARS-CoV, and MERS-CoV viral load dynamics, duration of viral shedding, and infectiousness: a systematic review and meta-analysis. *Lancet Microbe* **2**, e13–e22 (2021).
27. Wu, Y. *et al.* Prolonged presence of SARS-CoV-2 viral RNA in faecal samples. *Lancet Gastroenterol Hepatol* **5**, 434–435 (2020).
28. Xu, Y. *et al.* Characteristics of pediatric SARS-CoV-2 infection and potential evidence for persistent fecal viral shedding. *Nat. Med.* **26**, 502–505 (2020).
29. Goyal, R., Hotchkiss, J., Schooley, R. T., De Gruttola, V. & Martin, N. K. Evaluation of SARS-CoV-2 transmission mitigation strategies on a university campus using an agent-based network model. *Clin. Infect. Dis.* (2021) doi:10.1093/cid/ciab037.
30. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
31. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
32. Diamond, S. & Boyd, S. CVXPY: A Python-Embedded Modeling Language for Convex Optimization. *J. Mach. Learn. Res.* **17**, (2016).
33. Agrawal, A., Verschueren, R., Diamond, S. & Boyd, S. A rewriting system for convex optimization problems. *Journal of Control and Decision* **5**, 42–60 (2018).
34. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
35. Issues with SARS-CoV-2 sequencing data. <https://virological.org/t/issues-with-sars-cov-2-sequencing-data/473> (2020).

Extended Data:

Extended Data Table 1: Platemap of spike-in mixtures used for method validation

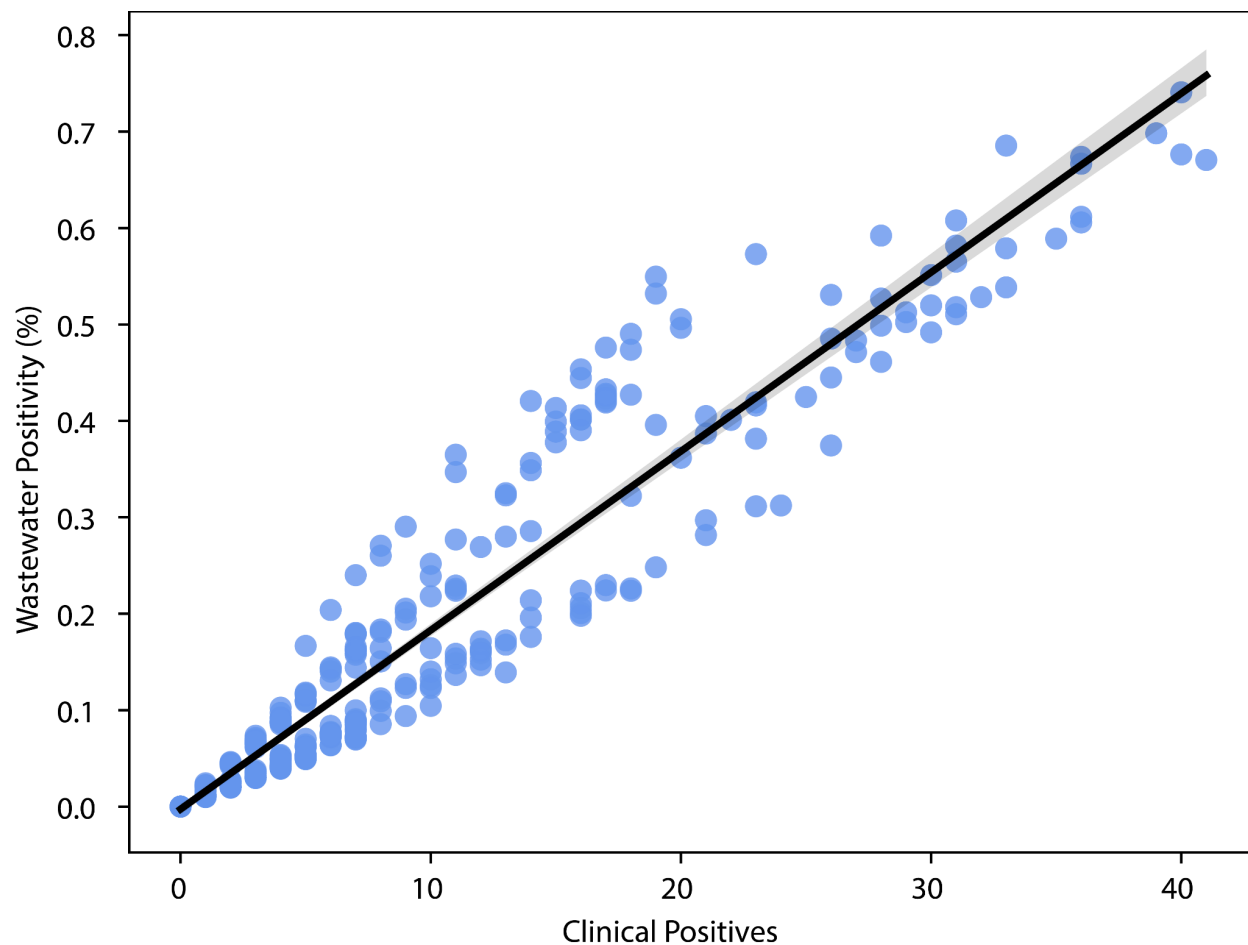
	1	2	3	4	5	6
A	5% Delta: 95% A	10% Delta: 90% A	20% Delta: 80% A	40% Delta: 60% A	50% Delta: 50% A	100% A
B	5% Delta: 95% Beta	10% Delta: 90% Beta	20% Delta: 80% Beta	40% Delta: 60% Beta	50% Delta: 50% Beta	100% Delta
C	5% Delta: 95% Gamma	10% Delta: 90% Gamma	20% Delta: 80% Gamma	40% Delta: 60% Gamma	50% Delta: 50% Gamma	100% Beta
D	5% Delta: 95% Alpha	10% Delta: 90% Alpha	20% Delta: 80% Alpha	40% Delta: 60% Alpha	50% Delta: 50% Alpha	100% Gamma
E	5% Beta: 95% A	10% Beta: 90% A	20% Beta: 80% A	40% Beta: 60% A	50% Beta: 50% A	100% Alpha
F	5% Beta: 95% Delta	10% Beta: 90% Delta	20% Beta: 80% Delta	40% Beta: 60% Delta	50% Beta: 50% Delta	20% A: 20% Delta: 20% Beta: 20% Gamma: 20% Alpha
G	5% Beta: 95% Gamma	10% Beta: 90% Gamma	20% Beta: 80% Gamma	40% Beta: 60% Gamma	50% Beta: 50% Gamma	25% Delta: 25% Beta : 25% Gamma: 25% Alpha
H	5% Beta: 95% Alpha	10% Beta: 90% Alpha	20% Beta: 80% Alpha	40% Beta: 60% Alpha	50% Beta: 50% Alpha	25% Delta: 25% Beta: 25% Gamma: 25% A
I	5% Gamma: 95% A	10% Gamma: 90% A	20% Gamma: 80% A	40% Gamma: 60% A	50% Gamma: 50% A	25% Delta: 25% Beta: 25% A: 25% Alpha
J	5% Gamma: 95% Delta	10% Gamma: 90% Delta	20% Gamma: 80% Delta	40% Gamma: 60% Delta	50% Gamma: 50% Delta	25% Delta: 25% A: 25% Gamma: 25% Alpha
K	5% Gamma: 95% Beta	10% Gamma: 90% Beta	20% Gamma: 80% Beta	40% Gamma: 60% Beta	50% Gamma: 50% Beta	25% A: 25% Beta: 25% Gamma: 25% Alpha
L	5% Gamma: 95% Alpha	10% Gamma: 90% Alpha	20% Gamma: 80% Alpha	40% Gamma: 60% Alpha	50% Gamma: 50% Alpha	33% Delta: 33% Beta: 33% Gamma
M	5% Alpha: 95% A	10% Alpha: 90% A	20% Alpha: 80% A	40% Alpha: 60% A	50% Alpha: 50% A	33% Delta: 33% Beta: 33% Alpha
N	5% Alpha: 95% Delta	10% Alpha: 90% Delta	20% Alpha: 80% Delta	40% Alpha: 60% Delta	50% Alpha: 50% Delta	33% Delta: 33% Alpha: 33% Gamma
O	5% Alpha: 95% Beta	10% Alpha: 90% Beta	20% Alpha: 80% Beta	40% Alpha: 60% Beta	50% Alpha: 50% Beta	33% Alpha: 33% Beta: 33% Gamma
P	5% Alpha: 95% Gamma	10% Alpha: 90% Gamma	20% Alpha: 80% Gamma	40% Alpha: 60% Gamma	50% Alpha: 50% Gamma	Neg

Extended Data Table 2: Consistency of Lineage A Cq values across repeated measurements

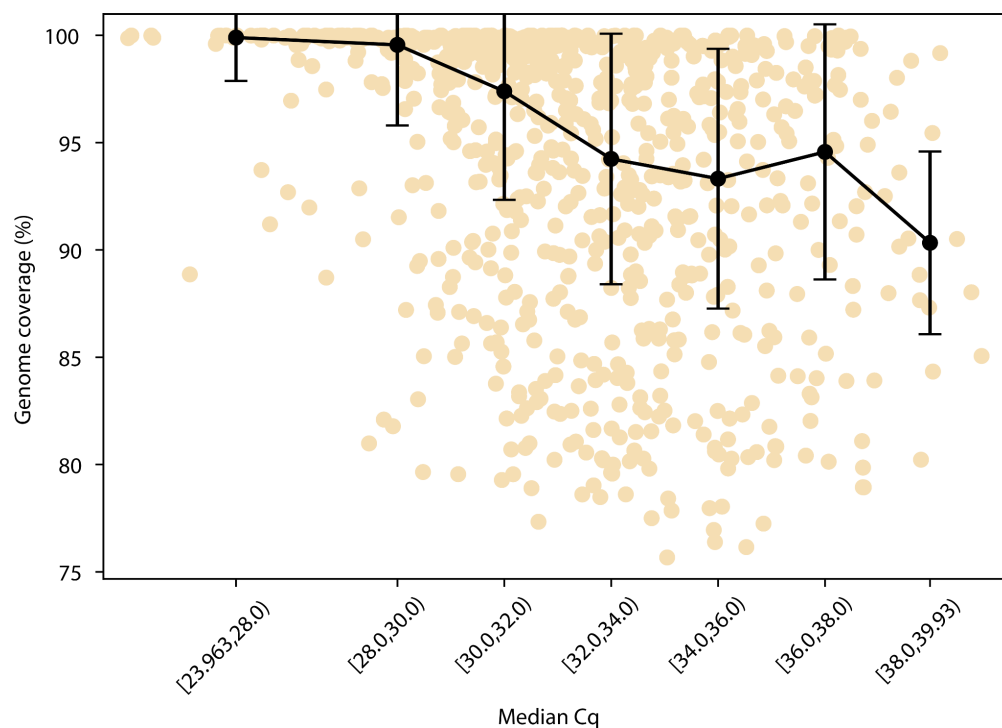
Replicate Number	Cq N Gene	Cq Orflab	Cq S Gene	Cq RNaseP	Average
1	31.228	30.807	30.045	29.581	30.693
2	30.783	29.77	29.546	29.49	30.033
3	31.201	30.622	29.733	29.745	30.519
4	30.621	30.953	29.578	28.925	30.384
5	31.188	30.073	29.366	28.745	30.209
6	30.604	29.788	29.829	28.797	30.074
7	30.308	30.335	29.573	29.149	30.072
8	30.738	30.36	29.711	28.79	30.269
9	31.144	29.97	30.045	28.79	30.386
10	31.122	30.822	29.566	29.671	30.503
11	31.825	29.763	29.833	29.134	30.474
12	31.434	30.18	29.773	29.133	30.462
13	31.209	29.793	29.402	29.559	30.135
14	30.641	30.181	29.816	29.833	30.213
15	30.744	29.371	29.695	29.257	29.937
16	30.396	29.728	29.441	28.428	29.855
17	30.957	29.449	29.913	28.242	30.107
18	30.791	30.113	29.601	29.277	30.169
19	31.561	29.839	29.943	29.06	30.448
20	31.434	29.711	29.568	28.864	30.238

Extended Data Table 3: Omicron surveillance at Point Loma Wastewater Treatment Plant

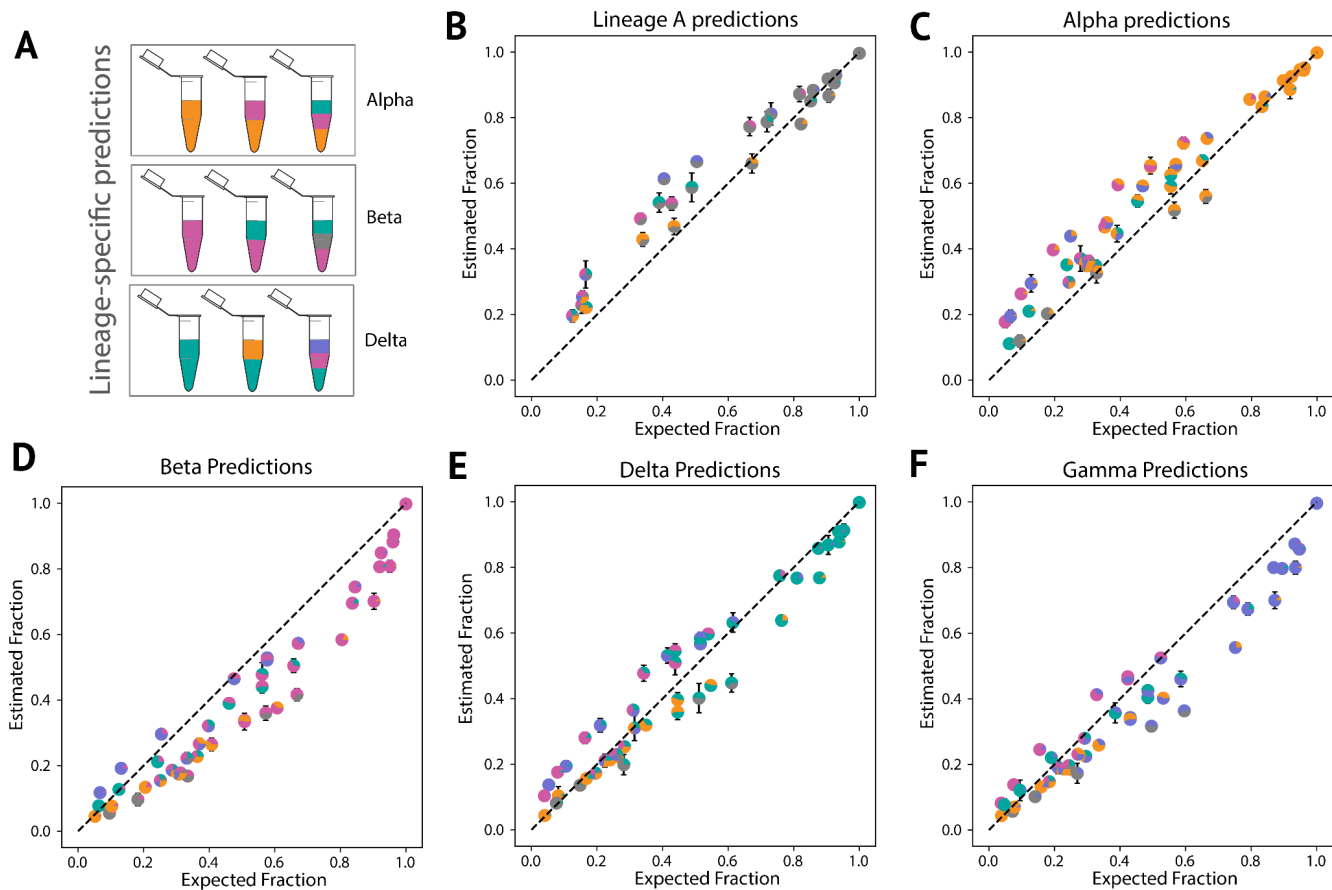
Collection Date	Avg. Estimated Omicron Abundance (%)	qPCR Detection		
		DelHV69/70	N501Y	P681R
10/04/21	0			X
10/06/21	0			X
10/10/21	0			X
10/11/21	0			X
10/13/21	0			X
10/17/21	0			X
10/18/21	0			X
10/20/21	0			X
11/12/21	0			X
11/22/21	0			X
11/27/21	1.726	X	X	X
11/28/21	1.967	X	X	X
12/1/21	2.439	X	X	X
12/5/21	17.11	X	X	X
12/6/21	19.764	X	X	X
12/12/21	50.65	X	X	X
12/16/21	67.14	X	X	X
12/20/21	79.135	X	X	X
12/21/21	80.567	X	X	X



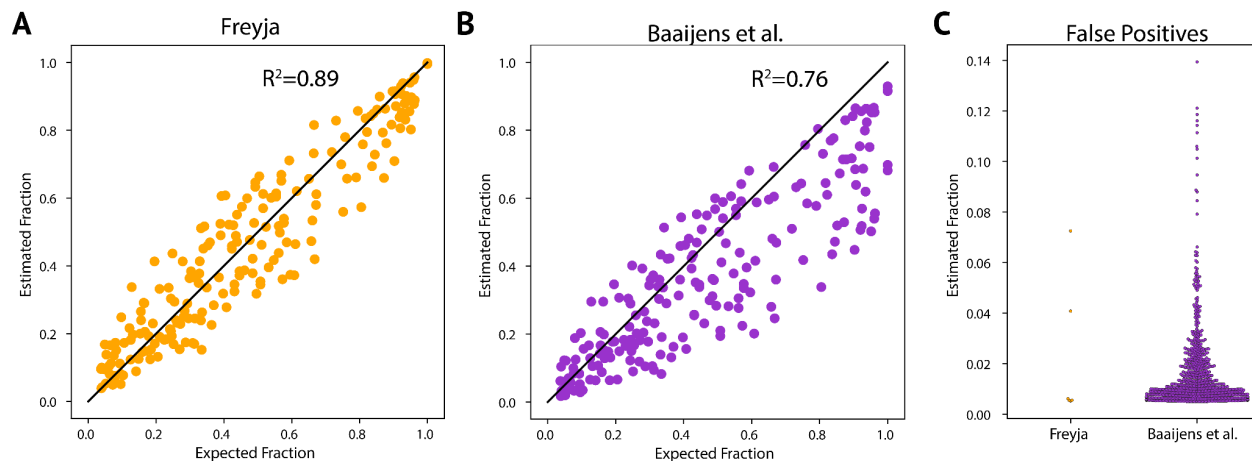
Extended Data Figure 1: Relationship of daily UCSD campus wastewater sampler positivity and campus clinical positives. Black line indicates the linear fit to the data, with bootstrap 95% confidence interval shown in gray.



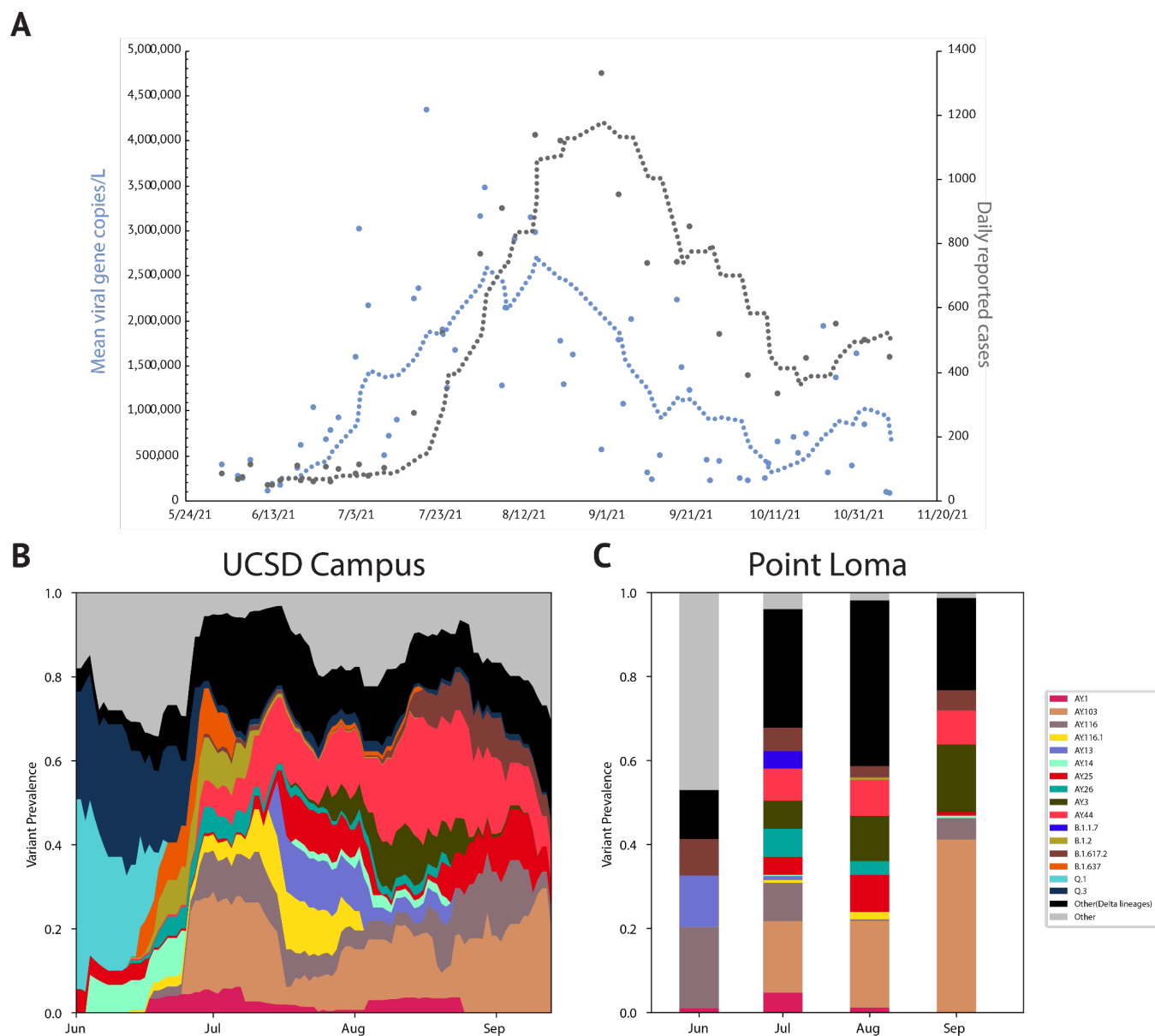
Extended Data Figure 2: Relationship between genome coverage and cycle quantification values. 10x genome coverage (fraction of sites with 10 reads or greater) remains high, even for Cq values of nearly 38. Points indicate median value in each bin, while error bars indicate the median absolute deviation.



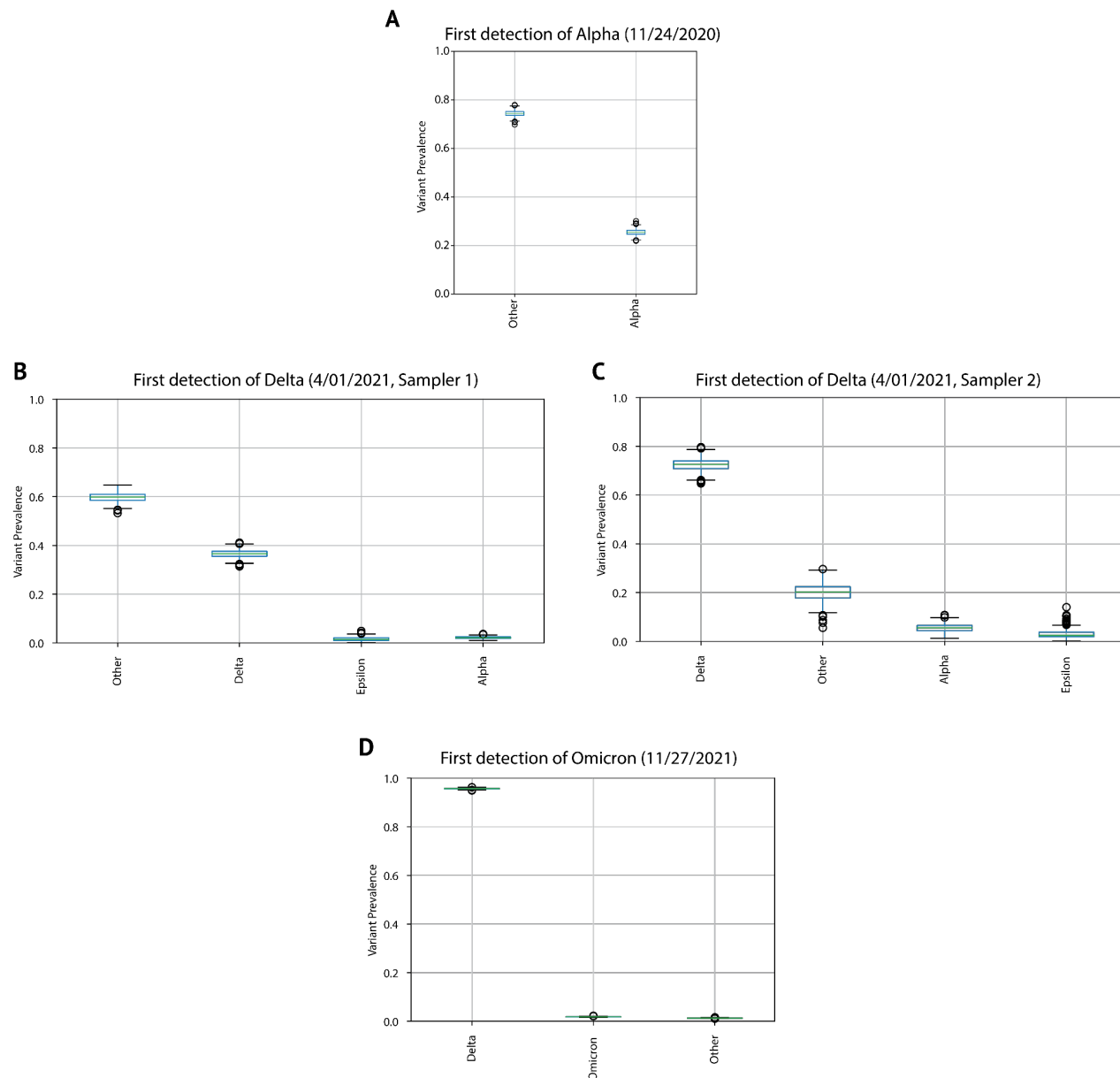
Extended Data Figure 3: Lineage-specific prediction of variant abundance in spike-in validation samples. A. Schematic of “spike-in” sample design. B-F. Lineage specific prediction. Proportions of each lineage in the sample are shown as a pie chart marker (Grey = Lineage A, Orange = Alpha, Pink = Beta, Turquoise = Delta, and Purple = Gamma) with error bars indicating the standard deviation from the mean, across four replicates.



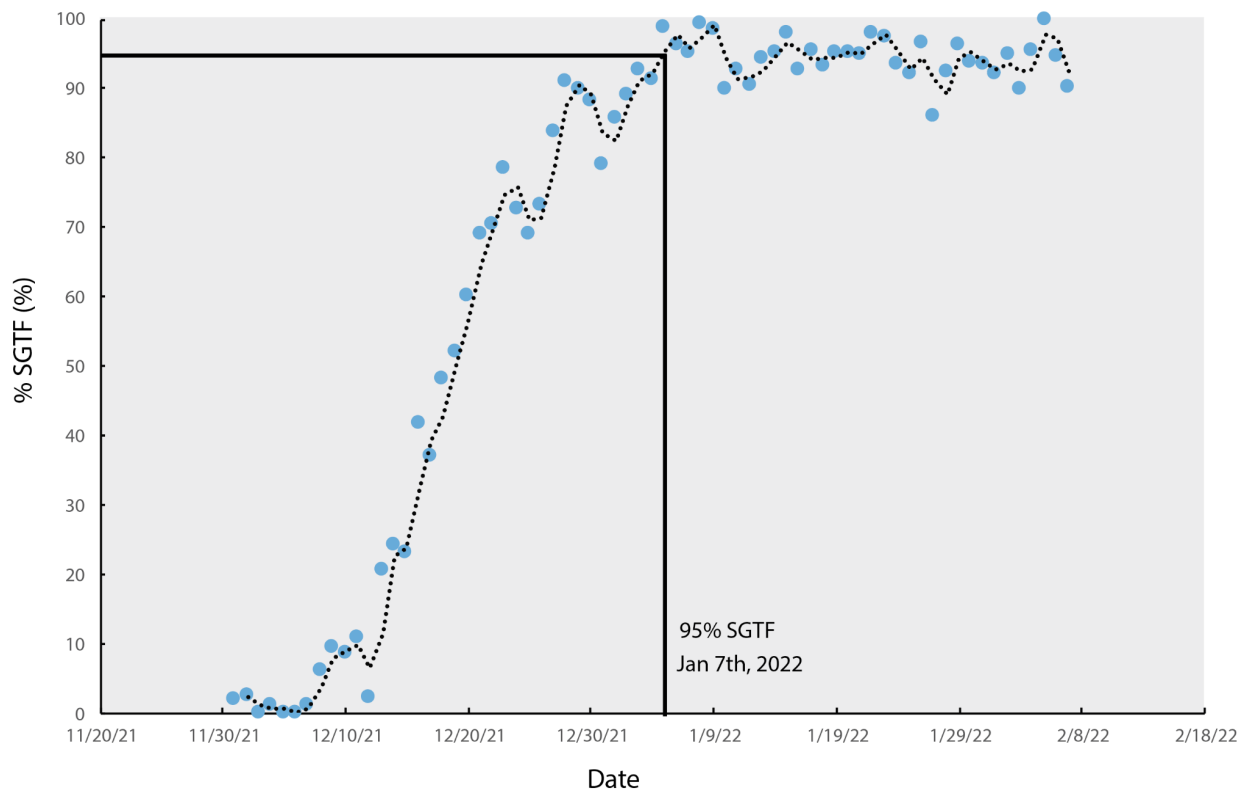
Extended Data Figure 4: Freyja more accurately estimates virus abundance, with fewer false positives. A-B. Estimated vs expected fraction of each lineage in the mixture. The Kallisto-based approach from Baaijens et. al shows a wider range of estimates for each known mix fraction, and generally underestimates the fraction. C. False positives with abundance greater than 0.5%.



Extended Data Figure 5: The rise of the Delta variant during Summer 2021. A. Mean SARS-CoV-2 viral gene copies/L of raw sewage (blue) collected from the Point Loma Wastewater Treatment Plant and caseload (gray) reported by the county during the same period. SARS-CoV-2 concentrations were normalized by PMMoV (pepper mild mottle virus) concentration to adjust for load changes. B. Lineage distribution in UCSD campus wastewater. C. Monthly lineage averages for wastewater collected at Point Loma Wastewater Treatment Plant during the Delta surge (N= 5, 20, 25, 7)



Extended Data Figure 6: Quantification of deconvolution uncertainty in first detection of VOCs. A-D. Bootstrap distributions of Freyja abundance estimates obtained by resampling read data from each sample corresponding to the first detection of that VOC in San Diego. Two samplers were found to contain Delta on the same day. First detections were also confirmed using a VOC qPCR panel, as shown in Figure 2 and Extended Data Table 3. 95% Confidence intervals for variant prevalence for each first detection event: A. Alpha: (0.232, 0.278), B. Delta: (0.336, 0.397), C. Delta: (0.676, 0.772), D. Omicron: (0.017, 0.021).



Extended Data Figure 7: Estimated proportion of Omicron sequences in clinical data.

Omicron estimates tracked via S-gene target failure, SGTF (characteristic of Omicron lineage BA.1 and its descendants) qPCR assays for clinical samples in San Diego between November 27th, 2021-February 7th, 2022. First detection of Omicron through clinical genomic sequencing in San Diego was December 8th. Dotted line shows a rolling average with a window size of seven days.