

# Wastewater sequencing reveals early cryptic SARS-CoV-2 variant transmission

<https://doi.org/10.1038/s41586-022-05049-6>

Received: 21 December 2021

Accepted: 29 June 2022

Published online: 7 July 2022

Open access



As SARS-CoV-2 continues to spread and evolve, detecting emerging variants early is critical for public health interventions. Inferring lineage prevalence by clinical testing is infeasible at scale, especially in areas with limited resources, participation, or testing and/or sequencing capacity, which can also introduce biases<sup>1–3</sup>. SARS-CoV-2 RNA concentration in wastewater successfully tracks regional infection dynamics and provides less biased abundance estimates than clinical testing<sup>4,5</sup>. Tracking virus genomic sequences in wastewater would improve community prevalence estimates and detect emerging variants. However, two factors limit wastewater-based genomic surveillance: low-quality sequence data and inability to estimate relative lineage abundance in mixed samples. Here we resolve these critical issues to perform a high-resolution, 295-day wastewater and clinical sequencing effort, in the controlled environment of a large university campus and the broader context of the surrounding county. We developed and deployed improved virus concentration protocols and deconvolution software that fully resolve multiple virus strains from wastewater. We detected emerging variants of concern up to 14 days earlier in wastewater samples, and identified multiple instances of virus spread not captured by clinical genomic surveillance. Our study provides a scalable solution for wastewater genomic surveillance that allows early detection of SARS-CoV-2 variants and identification of cryptic transmission.

SARS-CoV-2 continues to evolve, producing diverse new lineages<sup>6</sup>. Emerging variants of concern ( VOCs) and variants of interest ( VOIs) demonstrate increased transmissibility, disease severity and/or immune escape<sup>7</sup>. Timely and accurate quantification of local prevalence of SARS-CoV-2 variants is thus essential for effective public health measures. However, existing strategies for variant detection based on virus genome sequencing of biospecimens obtained from clinical testing ('clinical genomic surveillance') are expensive, inefficient and have sampling bias because of systemic healthcare disparities, particularly in poor and underserved communities<sup>1–3</sup>.

By contrast, PCR-based wastewater surveillance of SARS-CoV-2 RNA is not subject to clinical testing biases and can track temporal changes in overall SARS-CoV-2 prevalence in a region<sup>4,5,8</sup>, but cannot identify epidemiological transmission links or monitor virus lineage prevalence, which require genome sequence information. Virus genome sequencing from wastewater ('wastewater genomic surveillance') has the potential to cost-effectively capture community virus spread<sup>9,10</sup>, acting as a surrogate to clinical surveillance in elucidating lineage geo-spatial distributions and track emerging SARS-CoV-2 variants (including new variants for which targeted assays do not yet exist), and provide genome sequence data needed for transmission network analysis and interpretation<sup>11</sup>.

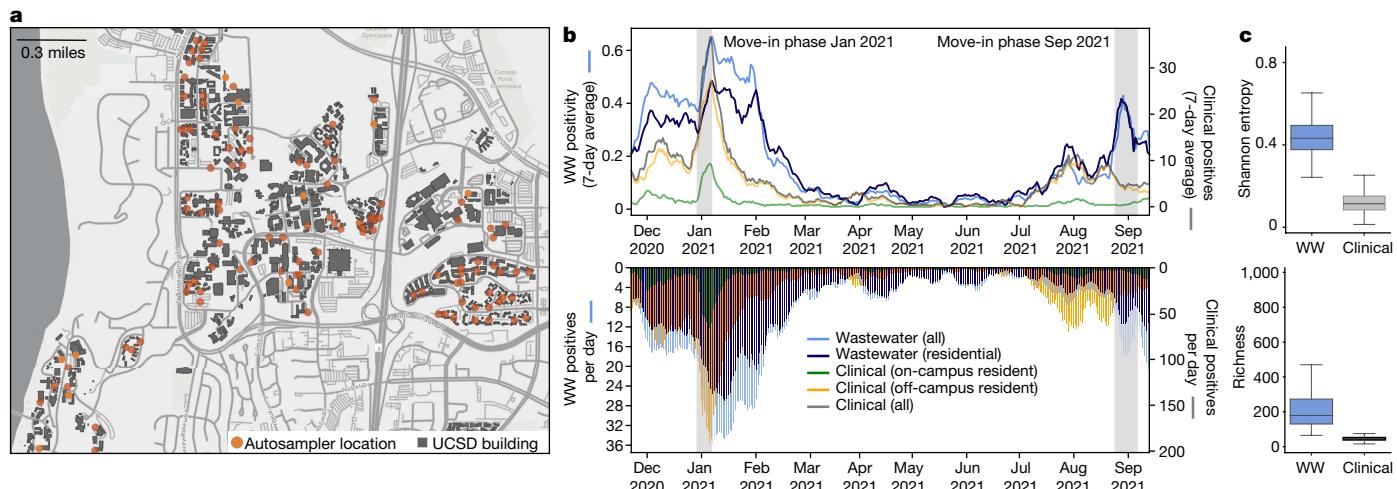
However, wastewater genomic surveillance is technically challenging<sup>10</sup>. Low viral loads, heavily fragmented RNA and PCR inhibitors in complex environmental samples lead to poor sequencing coverage<sup>12,13</sup>. Obtaining high-quality sequences from samples with low viral load and elevated levels of PCR inhibitors remains an outstanding technical challenge in implementation of wastewater genomic surveillance at scale. In addition, tools for SARS-CoV-2 lineage classification, such as

pangolin<sup>14</sup> and USHER<sup>15</sup>, were designed for clinical samples containing a single dominant variant, and cannot estimate relative abundances of multiple SARS-CoV-2 lineages in samples with virus mixtures such as wastewater.

Here we report a high-resolution approach to study community virus transmission using wastewater genomic surveillance, leveraging several technical advances in wastewater virus concentration and nucleic acid sequencing, and a computational tool for resolving multiple SARS-CoV-2 lineages in short-read sequence data from a mixed sample (lineage deconvolution). We obtained near 95% genome coverage even for samples with low viral load, compared with 40% or below from previous studies<sup>11–13</sup>, a key advance that allowed us to build a robust pipeline to monitor virus lineage prevalence in community wastewater.

Because places of communal living, such as university campuses, are considered key sites for virus spread and represent well-controlled and relatively isolated environments, they are ideal for comparing the relative utility of clinical and wastewater genomic surveillance<sup>16</sup>. Accordingly, we conducted a high-resolution, longitudinal wastewater genomic surveillance effort at the University of California San Diego (UCSD) campus, in parallel with clinical genomic surveillance from nasal swabs in the local community, from November 2020 to September 2021: 10 months that effectively capture the surges in the region caused by the three main VOCs (as determined by the US Centers of Disease Control and Prevention (CDC) in the United States: Epsilon, Alpha and Delta<sup>6</sup>). In more recent San Diego-wide data collected from September 2021 to February 2022, we studied ongoing transmission of the Delta variant and the rapid spread of the Omicron variant and its sublineages.

Our wastewater genomic surveillance approach identified VOCs up to 2 weeks before detection through clinical genomic surveillance,



**Fig. 1 | Campus sampling locations and SARS-CoV-2 testing statistics.**

**a**, Geospatial distribution of the 131 actively deployed wastewater autosamplers and the corresponding 360 university buildings on the campus sewer network. Building-specific data have been de-identified in accordance with university reporting policies. **b**, Campus wastewater (WW) and diagnostic testing statistics over the 295-day sampling period (positivity is the fraction of WW samplers with a positive qPCR signal). **c**, Virus diversity in wastewater and

clinical samples; boxplots of Shannon entropy (top) and richness (bottom) for each sample type ( $n = 153$  WW—a subset chosen to maximize sample independence; see Methods—and  $n = 5,888$  clinical). Box edges specify the first and third quartiles, the solid line indicates the median, and the whiskers delimit the maximum and minimum values. Map in **a** is the intellectual property of Esri and its licensors and are used herein under license. Copyright © 2022 Esri and its licensors. All rights reserved.

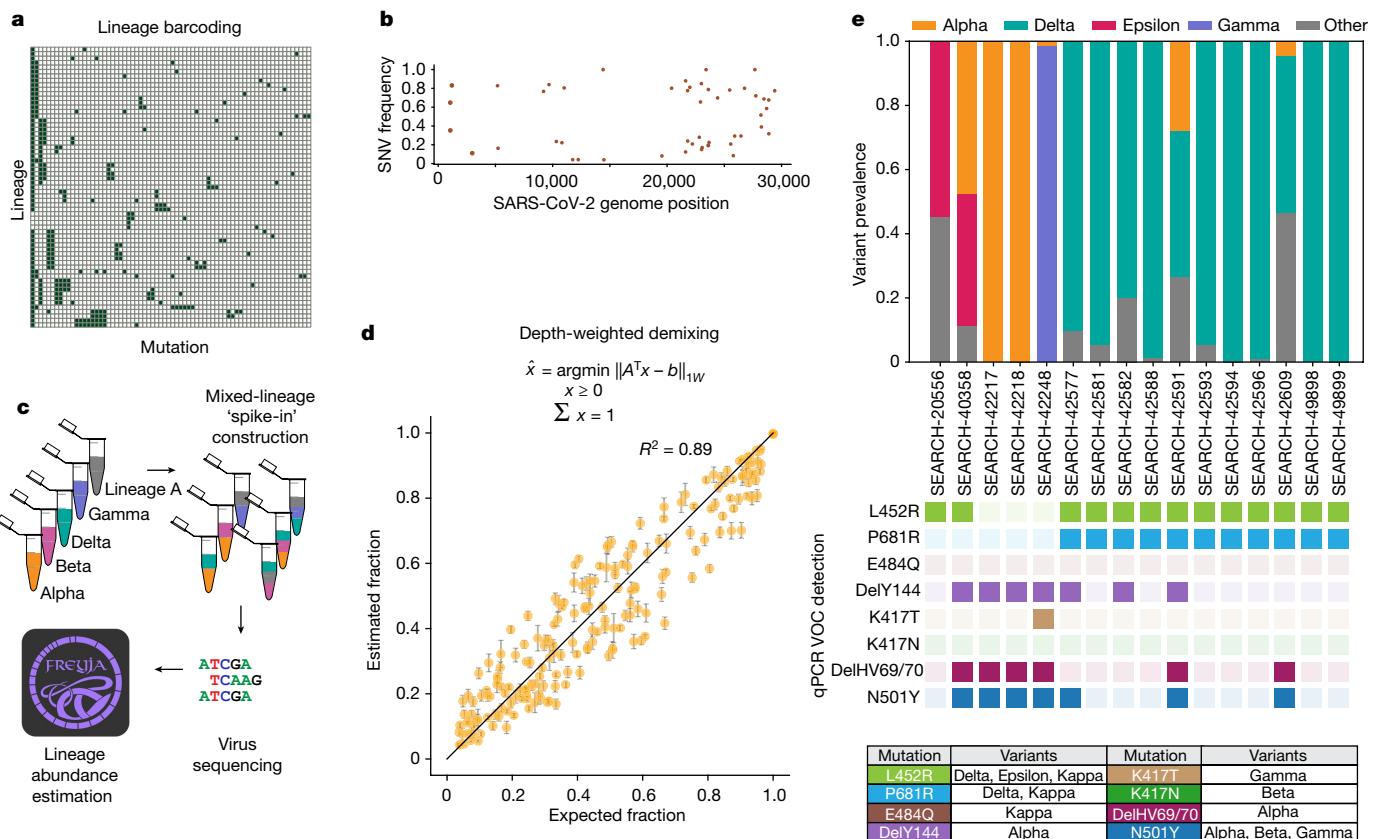
even though a large proportion of clinical SARS-CoV-2 samples is sequenced in San Diego relative to other cities in the United States. In addition to providing a detailed history of community virus spread, wastewater genomic surveillance also identified multiple instances of cryptic community transmission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

To directly compare wastewater genomic surveillance to clinical surveillance, we conducted a large-scale SARS-CoV-2 genome sequencing study from wastewater samples collected daily from 131 wastewater samplers covering 360 campus buildings, in many cases reaching single-building-level resolution. To identify epidemiological transmission links and monitor lineages in the population, we sequenced all clinical and wastewater samples positive for SARS-CoV-2 from campus using a miniaturized tiled-amplicon sequencing approach. During this period, we collected and analysed 21,419 wastewater samples: 19,944 wastewater samples from the UCSD campus, and, for comparison, 1,475 wastewater samples from the greater San Diego area, including the Point Loma wastewater treatment plant (the primary wastewater treatment plant for the county, with a catchment size of 2.3 million people) and 17 public schools spanning four San Diego school districts<sup>17</sup>. We compared sequencing of 600 campus wastewater samples to 759 genomes obtained from campus clinical swabs (46.2% of all positive tests on campus), all processed by the CALM and EXCITE CLIA laboratories at the UCSD. In addition, we compared 31,149 genomes obtained from clinical genomic surveillance of the greater San Diego community with sequencing of 837 wastewater samples collected from San Diego county (including those from the UCSD campus) during the same period.

## Uncovering microscale community spread

We implemented a geographical information system (GIS)-enabled building-level wastewater surveillance system to cover 360 buildings on the UCSD campus (Fig. 1a). During the period of daily wastewater sampling, approximately 10,000 students lived on campus and 25,000 individuals were on campus on a daily basis. We found that wastewater test positivity correlated strongly with the number of clinical positives (Fig. 1b and Extended Data Fig. 1), showing that wastewater effectively captures the community infection dynamics based on total viral load. This is also consistent with our past studies that showed that SARS-CoV-2 RNA can be detected approximately 85% of the time downstream from buildings containing individuals known to be infected<sup>9</sup>.

Unlike quantitative PCR (qPCR)-based mutant surveillance, genomic surveillance using full-length virus genomes can detect which strains of SARS-CoV-2 are circulating in the population, and can identify potential transmission links between infected individuals<sup>18,19</sup>. Although targeted qPCR mutant panels have the ability to detect specific lineages in wastewater, they only target a small set of mutations that must be known beforehand and require development and validation time before implementation. To test the utility of wastewater genomic surveillance for studying virus spread in the community, we obtained near-complete virus genomes for wastewater samples with cycle quantification values as high as 38 (median genome coverage of 96.49% (range 75.67–100.00%); Extended Data Fig. 2). However, using two common metrics of virus diversity—Shannon entropy (a measure of the uncertainty associated with randomly sampling an allele) and richness (the number of single-nucleotide variant (SNV) sites)<sup>20</sup>—we found that the genetic diversity of SARS-CoV-2 is significantly greater in wastewater samples than clinical samples (Fig. 1c; Mann–Whitney *U*-test,  $P < 0.001$  for each, with effect size  $r = 0.99, 0.97$  for Shannon entropy and richness, respectively). This suggests that multiple virus lineages, probably shed from different infected individuals, are often present in wastewater samples, whereas clinical samples generally contain a single virus lineage shed from one individual.



**Fig. 2 | Sample deconvolution robustly recovers relative virus abundance.**

**a**, Subset of lineage defining mutation ‘barcode’ matrix. Each row represents one lineage (out of more than 1,000 lineages included in the UShER global phylogenetic tree), and individual nucleotide mutations are represented as columns. **b**, Single-nucleotide variant (SNV) frequencies obtained from iVar used for recovering relative abundance of each lineage. **c**, Schematic of the

spike-in validation experiment. **d**, Depth-weighted demixing estimates of the virus abundance versus expected or known abundance. Details on lineage-specific predictions are provided in Extended Data Fig. 3. Error bars indicate s.d. of estimates across mixture replicates. **e**, Comparison of wastewater sample deconvolution with VOC qPCR panel, with lookup table (bottom) showing amino acid mutations corresponding to each variant.

## Sample deconvolution robustly recovers the abundance of SARS-CoV-2 lineages in mixed samples

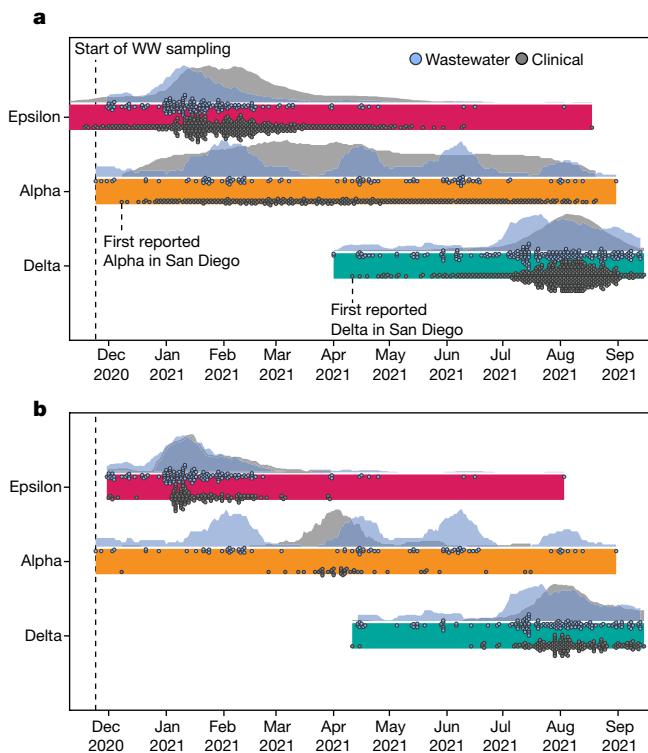
Wastewater systems aggregate stool, urine and other biological waste products carrying viruses from multiple infected individuals in the community in a single location, allowing for sampling of virus mixtures that are representative of local lineage prevalence. However, existing methods for determining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

To fully capture the virus diversity in community biospecimens, we developed Freyja, a tool to estimate the relative abundance of virus lineages in a mixed sample. Freyja uses a ‘barcode’ library of lineage-defining mutations to represent each SARS-CoV-2 lineage in the global phylogeny<sup>21</sup> (Fig. 2a). To encode each sample, Freyja stores the SNV frequencies (the proportion of reads at a site that contain the SNV) for each of the lineage-defining mutations (Fig. 2b, top). As SNV frequencies at positions with greater sequencing depth more accurately estimate the true mutation frequency, Freyja recovers relative lineage abundance by solving a depth-weighted least absolute deviation regression problem, a mixed sample analogue of minimizing the edit distance between sequences and a reference (Fig. 2b, bottom). To ensure results are meaningful, Freyja constrains the solution space such that each lineage abundance value is non-negative and overall lineage abundance sums to one. Freyja performs site-specific weighting to account for non-constant variance in measured SNV frequency across sites, enabling prioritization of information at each site as a

function of sequencing depth. Read depths were log-transformed, providing robustness to common attributes of real sequencing data such as heavily skewed read depth across amplicons.

To validate Freyja, we sequenced ‘spike-in’ synthetic mixtures from five key SARS-CoV-2 lineages (lineage A, Beta, Delta, Epsilon and Gamma) at proportions ranging from 5% to 100% in each sample, with between one and five different lineages per mixture (Fig. 2c and Extended Data Table 1). We found that Freyja robustly recovered the expected lineage abundances for all mixtures, even for lineages at 5% abundance (Fig. 2d and see Extended Data Fig. 3 for lineage-specific predictions). To further validate Freyja, we used wastewater samples from the UCSD isolation dorms as well as Point Loma wastewater treatment plant, collection sites likely to contain mixed-lineage samples, to compare Freyja-detected lineages with qPCR testing for eight mutations associated with different VOCs (N501Y, DelHV69/70, DeY144, K417N, K417T, E484Q, P681R and L452R; Fig. 2e). We found that Freyja consistently identified the same lineages as qPCR testing, but, as expected, also identified additional lineages with SNVs not included in our qPCR panel that were known to be circulating in San Diego at the time of collection. Combined, these results show that Freyja robustly estimates viral lineage abundance from samples containing a mixture of lineages, including synthetic virus mixtures and field wastewater collections.

To compare Freyja with other wastewater analysis pipelines, we tested the performance of other wastewater deconvolution methods including a method from Baaijens et al.<sup>12</sup>, cojac<sup>22</sup> and Lineage deComposition for SARS-CoV-2 (LCS)<sup>23</sup> using the spike-in mixtures (Extended



**Fig. 3 | Freyja recovers early and cryptic transmission of SARS-CoV-2 variants of concern.** **a**, Timeline and normalized epidemiological curves for VOC detection in both wastewater and clinical sequences from San Diego County (includes wastewater samples collected from Point Loma wastewater treatment plant, UCSD, as well as public schools in the San Diego districts) for the three major VOCs in circulation during the sampling period ( $n = 475$  wastewater,  $n = 22,504$  clinical). Both Alpha and Delta variants are detected first in wastewater before clinical samples. Markers for clinical detections correspond to the ceiling of the daily detection count divided by 30 (for example, 1–30 samples = one marker, 31–60 = two markers), whereas wastewater markers correspond to a single detection. **b**, Timeline and epidemiological curves for VOC detection in the campus samples ( $n = 364$  wastewater,  $n = 333$  clinical). Markers correspond to a single detection event for both clinical and wastewater surveillance. All wastewater detections correspond to an estimated VOC prevalence of at least 10%.

Data Fig. 4). We found that Freyja greatly outperformed other methods in terms of accuracy, false-positive rate and computational efficiency. The method from Baaijens et al.<sup>12</sup> required more than ten times more computation time per sample relative to Freyja (approximately 13.2 min versus approximately 1.1 min per sample, respectively). Although cojac was fast, the small amplicon length used for the spike-in mixtures resulted in cojac failing to identify most of the variants entirely, whereas LCS failed to return estimates within 2 days.

### Detection of community transmission in wastewater

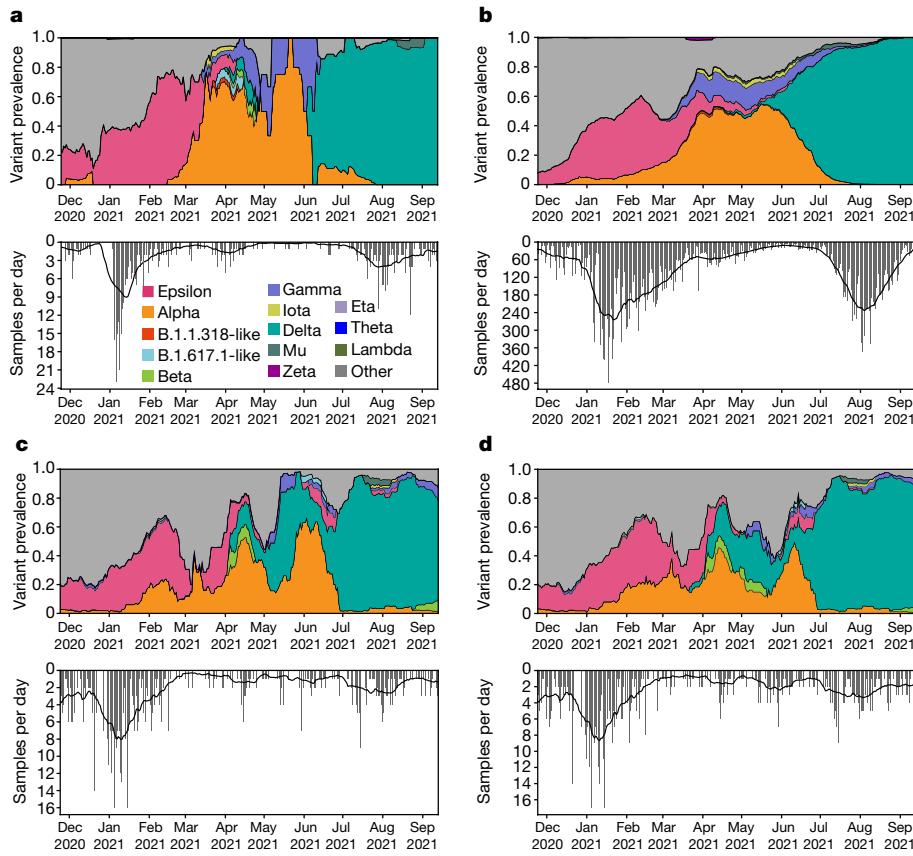
SARS-CoV-2 RNA concentrations in wastewater have been shown to be an early indicator of rising COVID-19 community incidence<sup>9,24</sup> (Extended Data Fig. 5A), but whether wastewater can be used to detect emerging variants, including VOCs and VOIs, before their observation in clinical surveillance is unknown. To test whether wastewater can enable early detection of emerging lineages, we applied Freyja to our wastewater sequencing data and compared the collection date of VOC-positive samples from wastewater with the collection dates of samples from clinical genomic surveillance (Fig. 3a). With only 2.6% as many sequenced wastewater samples as sequenced clinical samples, we detected the Alpha and Delta VOC lineages in wastewater

genomic surveillance up to 14 days before their first detection in genomic clinical surveillance (Epsilon was circulating at the start of wastewater collection, and thus could not be detected early). To further quantify our uncertainty in prevalence estimates, we used a fast bootstrapping approach (Extended Data Fig. 6) and found that the resampled distributions did not include zero abundance. As emerging VOC lineages may evade immune responses or lessen the effectiveness of public health interventions<sup>18</sup>, this early detection provides additional time to make necessary adjustments to existing countermeasures.

To test whether wastewater genomic surveillance can identify changes in the abundance of circulating lineages, we compared VOC detection rates in clinical and wastewater sequencing over time. We found that both wastewater and clinical genomic surveillance tracked changes in lineage abundance, but increases in lineage detection frequency were generally observed first in wastewater surveillance. For example, for the Epsilon variant, which was first detected in San Diego in September 2020, we observed increases in detection frequency in wastewater approximately 5 days before the corresponding increase in clinical genomic surveillance data (Fig. 3a; see Methods). We noticed varying periods of ongoing lineage detection across VOCs relative to clinical surveillance, possibly due to different virus shedding characteristics across lineages<sup>25</sup>. For Epsilon specifically, elevated sampling density on the UCSD campus relative to elsewhere in the county early on in the experiment may have biased San Diego-wide detection trends towards campus trends, particularly during the end of the wave. We also observed clear signatures of times with increased travel, as seen in the pulsing of Alpha detections in wastewater around the end of holidays and school breaks. During these periods as well as other times of mass student arrival, students were mandated to test immediately on arrival before they moved into their respective on-campus housing. In late March 2021 following the university break, mandated clinical testing identified the spread of the Alpha variant exclusively in off-campus residents (Fig. 1b), suggesting that campus mitigation protocols kept the Alpha outbreak from spreading on campus during this period.

To study the effectiveness of wastewater genomic surveillance at a smaller community scale, we restricted our analysis to samples from the UCSD campus. We found that wastewater genomic surveillance consistently identified the three major VOCs (Epsilon, Alpha and Delta) throughout their period of occurrence, despite detection gaps of 1 month or longer in clinical surveillance that included regular asymptomatic testing, longer than the expected signal due to extended virus shedding<sup>26–28</sup> (Fig. 3b). During these gaps, positive samples were collected from multiple distinct locations, with most locations not repeated, suggesting that this continued detection in wastewater was not simply due to extended shedding. From mid-December to late March, the Alpha variant was detected more than once per week on average in wastewater but was not detected by clinical surveillance. Similarly, wastewater surveillance detected continued Delta transmission from mid-April to mid-June, but no cases were identified by clinical surveillance. This explains, in part, the long tails of wastewater positivity on campus relative to clinical surveillance on campus (Fig. 1b), in which we controlled for extended shedding by excluding samples from campus isolation dorms (see Methods for details). The high wastewater positivity level in February–March 2020 extends beyond the expected duration of extended shedding, indicating that cryptic transmission probably had a substantial role in the spread of the virus on campus during this period.

To study the effectiveness of wastewater surveillance in detecting and tracking other emerging variants, we aggregated all wastewater sequencing data to estimate the temporal profile of community lineage prevalence. We found that estimates of lineage abundance using wastewater enable early identification of other VOCs or VOIs, even for lineages that are rarely observed in clinical surveillance (Fig. 4). For



**Fig. 4 | Deconvolution recovers a fine-grained estimate of virus population dynamics.** **a,b**, Prevalence of SARS-CoV-2 variants in UCSD clinical surveillance (**a**) and variant prevalence in all clinical samples collected in San Diego County (**b**). **c,d**, Variant prevalence in wastewater at UCSD (**c**) and the greater San Diego

County (**d**). Further analysis of Point Loma wastewater samples is shown in Extended Data Fig. 5. All curves show the rolling average, with a window of  $\pm 10$  days. ‘Other’ contains all lineages not designated as VOCs. The bottom panels show the number of sequenced samples per day.

example, we detected the Mu (B.1.621) variant via wastewater genomic surveillance on 27 July, nearly 4 weeks before its first detection through clinical genomic surveillance on campus on 23 August (Fig. 4a,c). However, despite persistent Mu detection in campus wastewater throughout July and early August, we did not detect the Mu variant in clinical or wastewater genomic surveillance on campus in September, suggesting that local community transmission did not continue.

To test whether Freyja continues to provide representative estimates of lineage prevalence for mixtures containing closely related lineages, we analysed the rise of the Delta variant (B.1.617.2) and its sublineages (AY<sup>+</sup>) in San Diego, from June to September 2021 (Extended Data Fig. 5B,C). At both the UCSD campus and the Point Loma wastewater treatment plant, we identified the rapid emergence of B.1.617.2 and its sublineages (AY<sup>+</sup>), along with low but persistent levels of the P.1 (Gamma) variant. The relative abundances of each of the variants were within twofold of prevalence estimates observed in clinical nasal swab data, suggesting that Freyja effectively identifies prevalence even for closely related lineages, at both the university and the county scale.

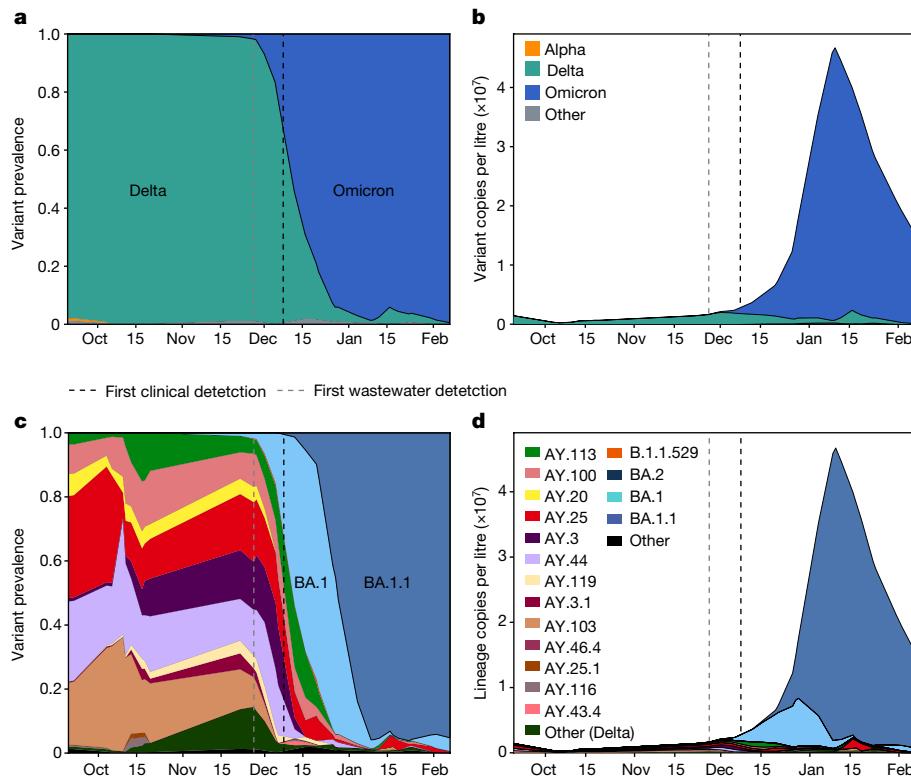
From more recent data from Point Loma wastewater treatment plant, we identified the Omicron variant (B.1.1.529 and descendants) at an abundance of near 1.7% on 27 November, more than 10 days before the first clinical detection in San Diego on 8 December (Fig. 5a,b). To confirm these findings, we applied our VOC qPCR panel to the same samples and consistently detected two mutations associated with the Omicron variant (DelHV69/70 and N501Y) in samples detected after 27 November, whereas neither was detected in samples from earlier in November (Extended Data Table 3; P681R was included to confirm the presence of Delta).

To visualize the dynamics of competition between the Delta and Omicron variants, we analysed wastewater collected at Point Loma from late September through to early February. We found that upon introduction to the community, Omicron rapidly rose to dominance and reached roughly 95% prevalence by 26 December. During the same period, the estimates for 95% Omicron abundance in clinical samples tracked via S-gene target failures was 7 January, further suggesting that wastewater genomic surveillance is a leading indicator of lineage dynamics for emerging variants (Fig. 5a and Extended Data Fig. 6E). To understand the magnitude of lineage abundance, we scaled each sample by the measured virus RNA concentration of the sample (Fig. 5b). We observed that the absolute amount of circulating Delta variant remained largely constant upon the introduction of Omicron, even as it appeared to decrease to a small fraction of all viruses in the community.

To study the contribution of individual virus lineages to virus RNA concentration, we further analysed the growth dynamics of Delta and Omicron sublineages (Fig. 5c,d). We found that the many Delta lineages circulating in October and November were rapidly displaced by the BA.1 Omicron lineage, which was soon after displaced by the BA.1.1 lineage, suggesting a growth advantage over BA.1 and B.1.1.529. We did not observe significant levels of any other Omicron sublineages.

## Wastewater reveals history of campus infections

Phylogenetic analysis of virus genomes can be used to identify fine-scale spatial and temporal transmission networks, but it is not known whether wastewater can be used to further refine possible sites of transmission, elucidate transmission networks ('who infected whom') or identify



**Fig. 5 | Community wastewater enables early Omicron detection and reveals lineage dynamics.** **a**, Prevalence of SARS-CoV-2 VOCs in wastewater collected from the Point Loma wastewater treatment plant from late September 2021 to early February 2022. **b**, Estimated VOC concentrations; prevalence estimates

were scaled by normalized viral load in wastewater. **c,d**, Lineage-specific estimates of prevalence (**c**) and concentration (**d**). All curves show an adaptive rolling average calculated using a local linear approximation (Savitzky–Golay filter) of virus copies per litre, with a window size of  $\pm 1$  sampling date.

specific infected individuals<sup>19</sup>. To investigate the scale, structure and timing of SARS-CoV-2 spread on campus, we reconstructed a maximum likelihood phylogenetic tree for each of the major VOCs using all high-quality consensus genomes (see Methods for details) obtained from the UCSD campus, as well as reference sequences for each lineage obtained elsewhere in the United States (Fig. 6a–c). In each tree, we identified many independent introductions, some of which led to extended transmission on campus. The resulting virus diversity among the VOCs present on campus enables ruling out of most transmission links and suggests that the spread of virus on campus comprised many separate, small outbreaks.

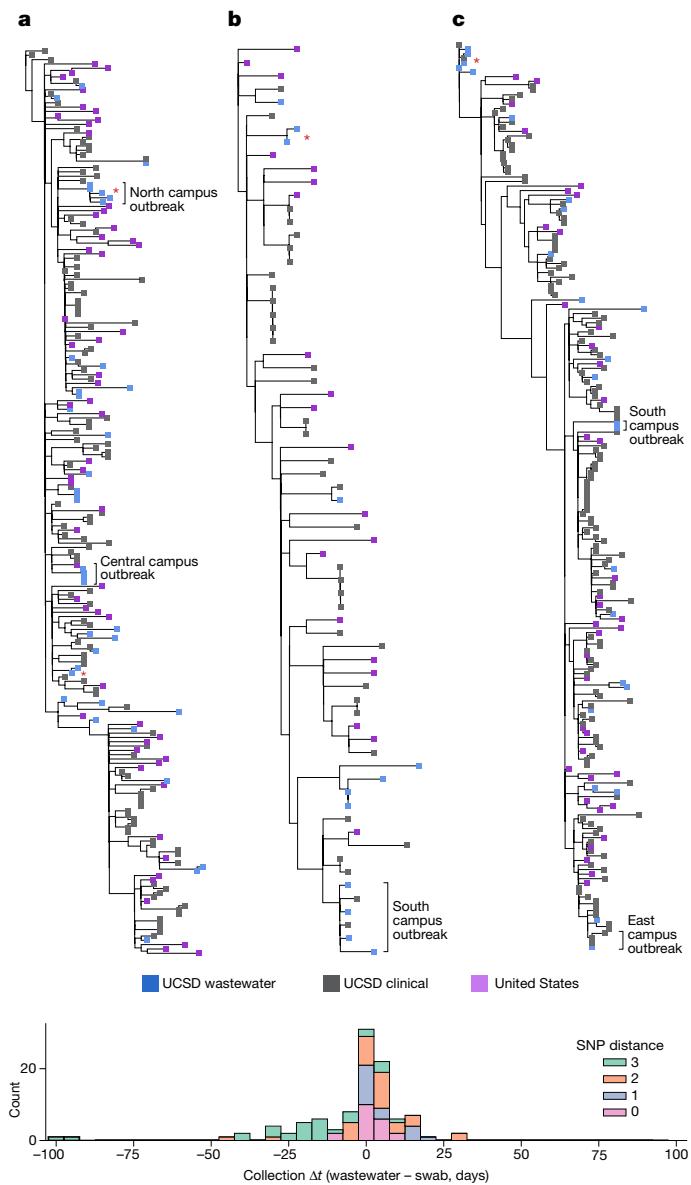
To analyse the spatial structure of virus spread, we identified collection sites for wastewater sequences connected to transmission chains on campus, with building-specific resolution (Fig. 6a–c and see detailed example in Extended Data Fig. 7). We observed multiple small, linked outbreaks clustered in nearby buildings. Campus isolation protocol required students in congregate living to relocate to an isolation dorm, and linkages in the wastewater samples from buildings used for isolation reflected this co-location. We also found multiple instances of successive exactly matching sequences from wastewater collected from a single building, possibly due to continued viral shedding from the same infected individuals from extended shedding in stool<sup>26–28</sup> or a transmission chain in the building leading to multiple infections by genetically identical viruses.

To study the temporal delay between clinical and wastewater lineage detection, we compared collection times of sequences from campus wastewater that match sequences from campus clinical surveillance (including non-VOC lineages). We found 20 exact sequence matches and 103 near-matches (SNP distance of three or less), but did not observe any overall bias towards earlier or later detection in

wastewater (Fig. 6d), suggesting that, on average, wastewater and clinical genomic surveillance identify a similar timing of individual detection events. However, despite current technical difficulties with isolating haplotypes from diverse virus mixtures, more than half of the clinical–wastewater sequence pairs demonstrate earlier detection in wastewater or are from the same date. As detection is often delayed or missed by clinical surveillance, detections occur first in wastewater (despite a loss of sequences due to limited haplotype recovery), further suggesting that wastewater genomic surveillance can reveal the presence of specific genome sequences before clinical surveillance.

## Discussion

We show that improved virus concentration from wastewater, coupled with a method for resolving multiple lineages from mixed samples, captures community virus lineage prevalence and enables early detection of emerging variants, often before observation in clinical surveillance. By sequencing both clinical and wastewater samples from the UCSD campus, we detected VOCs persistently in wastewater even when their appearance in clinical samples is intermittent. However, we also found occasions when rarer lineages, such as B.1.1.318, were detected in clinical samples but not in wastewater. This is not unexpected on campus as many students living off-campus did not contribute to campus wastewater but were still clinically tested as part of testing mandates and policies. In the larger San Diego community context, this suggests that we may not be able to identify lineages circulating at low prevalence (less than 1%) using a single wastewater collection site. In addition, we note that clinical sequences identified from the community may not be observable in the contributing catchment, as precise geolocation of



**Fig. 6 | Wastewater identifies clinically known and unknown virus transmission.** **a–c**, Maximum likelihood phylogenetic trees for each of the dominant VOCs (Epsilon (a), Alpha (b) and Delta (c)) using high-quality samples obtained at UCSD, as well as a representative set of sequences from the entire United States. Wastewater sequences from the same sampler that differ by one or fewer SNPs are denoted with a red asterisk. For all sequences, consensus bases were called at sites with more than 50% nucleotide frequency. Location information is provided for select outbreaks. **d**, Pairwise comparison of collection date for matching and near-matching wastewater and nasal swab samples obtained at UCSD. Positive values indicate earlier collection in nasal swabs and negative values indicate earlier detection in wastewater.

all clinical samples was not possible. Conversely, we also observed rare lineages in wastewater not seen in clinical samples from the campus or the community. As campus testing mandates are unable to capture all cases (for example, fully vaccinated individuals were not required to test and not all community samples were sequenced), rare lineages can be missed.

The considerable benefits of wastewater surveillance may stem from biases in clinical testing, including population testing availability and compliance, university quarantine policies and asymptomatic transmission, which may distort estimates of virus lineage prevalence from

clinical samples. Wastewater offers less biased and more consistent viral lineage prevalence estimates, especially in areas with limited access and/or higher testing hesitancy rates, in which limited clinical surveillance can delay detection of emerging variants. As it requires considerably fewer samples, it is also more cost-effective than clinical testing and could serve as a long-term passive surveillance tool. This is particularly important for developing public health interventions in low-resource and underserved communities, where widespread clinical genomic surveillance for SARS-CoV-2 remains limited.

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

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

## Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41586-022-05049-6>.

1. 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).
2. 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).
3. Brito, A. F. et al. Global disparities in SARS-CoV-2 genomic surveillance. Preprint at medRxiv <https://doi.org/10.1101/2021.08.21.21262393> (2021).
4. Karthikeyan, S. et al. High-throughput wastewater SARS-CoV-2 detection enables forecasting of community infection dynamics in San Diego County. *mSystems* **6**, e00045-21 (2021).
5. Randazzo, W. et al. SARS-CoV-2 RNA in wastewater anticipated COVID-19 occurrence in a low prevalence area. *Water Res.* **181**, 115942 (2020).
6. Mullen, J. L., the Center for Viral Systems Biology et al. *outbreak.info* <https://outbreak.info/> (2021).
7. Harvey, W. T. et al. SARS-CoV-2 variants, spike mutations and immune escape. *Nat. Rev. Microbiol.* **19**, 409–424 (2021).
8. 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).

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**, e02703-20 (2021).
12. Baaijens, J. A. et al. Variant abundance estimation for SARS-CoV-2 in wastewater using RNA-seq quantification. Preprint at medRxiv <https://doi.org/10.1101/2021.08.31.21262938> (2021).
13. Amman, F. et al. Viral variant-resolved wastewater surveillance of SARS-CoV-2 at national scale. *Nat. Biotechnol.* <https://doi.org/10.1038/s41587-022-01387-y> (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. Preprint at medRxiv <https://doi.org/10.1101/2021.10.19.21265226> (2021).
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.* **38**, 5819–5824 (2021).
22. Jahn, K. et al. Early detection and surveillance of SARS-CoV-2 genomic variants in wastewater using COJAC. *Nat. Microbiol.* **7**, 1151–1160 (2022).
23. Valeris, R. et al. A mixture model for determining SARS-CoV-2 variant composition in pooled samples. *Bioinformatics* **38**, 1809–1815 (2022).
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).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022

**Smruthi Karthikeyan**<sup>1,29</sup>, **Joshua I. Levy**<sup>2,29</sup>, **Peter De Hoff**<sup>3,4,5</sup>,  
**Greg Humphrey**<sup>1</sup>, **Amanda Birmingham**<sup>6</sup>, **Kristen Jepsen**<sup>7</sup>,  
**Sawyer Farmer**<sup>1</sup>, **Helena M. Tubb**<sup>1</sup>, **Tommy Valles**<sup>1</sup>,  
**Caitlin E. Tribelhorn**<sup>1</sup>, **Rebecca Tsai**<sup>1</sup>, **Stefan Aigner**<sup>3</sup>,  
**Shashank Sathe**<sup>3</sup>, **Niema Moshiri**<sup>8</sup>, **Benjamin Henson**<sup>3</sup>,  
**Adam M. Mark**<sup>6</sup>, **Abbas Hakim**<sup>3,4,5</sup>, **Nathan A. Baer**<sup>3</sup>, **Tom Barber**<sup>3</sup>,  
**Pedro Belda-Ferre**<sup>3</sup>, **Marisol Chacón**<sup>3</sup>, **Willi Cheung**<sup>3,4,5</sup>,  
**Evelyn S. Cresini**<sup>3</sup>, **Emily R. Eisner**<sup>3</sup>, **Alma L. Lastrella**<sup>3</sup>

**Elijah S. Lawrence**<sup>3</sup>, **Clarisse A. Marotz**<sup>3</sup>, **Toan T. Ngo**<sup>3</sup>,  
**Tyler Ostrander**<sup>3</sup>, **Ashley Plascencia**<sup>3</sup>, **Rodolfo A. Salido**<sup>3</sup>,  
**Phoebe Seaver**<sup>3</sup>, **Elizabeth W. Smoot**<sup>3</sup>, **Daniel McDonald**<sup>1</sup>,  
**Robert M. Neuhard**<sup>9,10</sup>, **Angela L. Scioscia**<sup>4,11</sup>, **Alysson M. Satterlund**<sup>12</sup>,  
**Elizabeth H. Simmons**<sup>13</sup>, **Dismas B. Abelman**<sup>10</sup>, **David Brenner**<sup>10</sup>,  
**Judith C. Bruner**<sup>10</sup>, **Anne Buckley**<sup>10</sup>, **Michael Ellison**<sup>10</sup>, **Jeffrey Gattas**<sup>10</sup>,  
**Steven L. Gonias**<sup>14</sup>, **Matt Hale**<sup>10</sup>, **Faith Hawkins**<sup>10</sup>, **Lydia Ikeda**<sup>10</sup>,  
**Hemlata Jhaveri**<sup>10</sup>, **Ted Johnson**<sup>10</sup>, **Vince Kellen**<sup>10</sup>, **Brendan Kremer**<sup>10</sup>,  
**Gary Matthews**<sup>10</sup>, **Ronald W. McLawhon**<sup>10</sup>, **Pierre Ouillet**<sup>10</sup>,  
**Daniel Park**<sup>10</sup>, **Allorah Pradenas**<sup>10</sup>, **Sharon Reed**<sup>10</sup>, **Lindsay Riggs**<sup>10</sup>,  
**Alison Sanders**<sup>10</sup>, **Bradley Sollenberger**<sup>10</sup>, **Angela Song**<sup>9,10</sup>,  
**Benjamin White**<sup>10</sup>, **Terri Winbush**<sup>10</sup>, **Christine M. Aceves**<sup>2</sup>,  
**Catelyn Anderson**<sup>2</sup>, **Karthik Gangavarapu**<sup>2</sup>, **Emory Hufbauer**<sup>2</sup>,  
**Ezra Kurzban**<sup>2</sup>, **Justin Lee**<sup>2</sup>, **Nathaniel L. Matteson**<sup>2</sup>, **Edyth Parker**<sup>2</sup>,  
**Sarah A. Perkins**<sup>2</sup>, **Karthik S. Ramesh**<sup>2</sup>, **Refugio Robles-Sikisaka**<sup>2</sup>,  
**Madison A. Schwab**<sup>2</sup>, **Emily Spencer**<sup>2</sup>, **Shirlee Wohl**<sup>2</sup>,  
**Laura Nicholson**<sup>15</sup>, **Ian H. McHardy**<sup>15</sup>, **David P. Dimmock**<sup>16</sup>,  
**Charlotte A. Hobbs**<sup>16</sup>, **Omid Bakhtar**<sup>17</sup>, **Aaron Harding**<sup>17</sup>,  
**Art Mendoza**<sup>17</sup>, **Alexandre Bolze**<sup>18</sup>, **David Becker**<sup>18</sup>,  
**Elizabeth T. Cirulli**<sup>18</sup>, **Magnus Isaksson**<sup>18</sup>, **Kelly M. Schiabor Barrett**<sup>18</sup>,  
**Nicole L. Washington**<sup>18</sup>, **John D. Malone**<sup>19</sup>, **Ashleigh Murphy Schafer**<sup>19</sup>,  
**Nikos Gurfield**<sup>19</sup>, **Sarah Stous**<sup>19</sup>, **Rebecca Fielding-Miller**<sup>20,21</sup>,  
**Richard S. Garfein**<sup>20</sup>, **Tommi Gaines**<sup>21</sup>, **Cheryl Anderson**<sup>20</sup>,  
**Natasha K. Martin**<sup>21</sup>, **Robert Schooley**<sup>21</sup>, **Brett Austin**<sup>17</sup>,  
**Duncan R. MacCannell**<sup>22</sup>, **Stephen F. Kingsmore**<sup>16</sup>, **William Lee**<sup>18</sup>,  
**Seema Shah**<sup>19</sup>, **Eric McDonald**<sup>19</sup>, **Alexander T. Yu**<sup>5</sup>, **Mark Zeller**<sup>2</sup>,  
**Kathleen M. Fisch**<sup>4,6</sup>, **Christopher Longhurst**<sup>1,23</sup>, **Patty Maysent**<sup>24</sup>,  
**David Pride**<sup>14,25</sup>, **Pradeep K. Khosla**<sup>8</sup>, **Louise C. Laurent**<sup>3,4,26</sup>,  
**Gene W. Yeo**<sup>3,26,27</sup>, **Kristian G. Andersen**<sup>2,30</sup> & **Rob Knight**<sup>1,8,28,30</sup>✉

<sup>1</sup>Department of Pediatrics, University of California San Diego, La Jolla, CA, USA. <sup>2</sup>Department of Immunology and Microbiology, The Scripps Research Institute, La Jolla, CA, USA.

<sup>3</sup>Expedited COVID Identification Environment (EXCITE) Laboratory, Department of Pediatrics, University of California San Diego, La Jolla, CA, USA. <sup>4</sup>Department of Obstetrics, Gynecology, and Reproductive Sciences, University of California San Diego, La Jolla, CA, USA. <sup>5</sup>COVID-19 Detection, Investigation, Surveillance, Clinical, and Outbreak Response, California

Department of Public Health, Richmond, CA, USA. <sup>6</sup>Center for Computational Biology and Bioinformatics, University of California San Diego, La Jolla, CA, USA. <sup>7</sup>Institute for Genomic Medicine, University of California San Diego, La Jolla, CA, USA. <sup>8</sup>Department of Computer Science and Engineering, University of California San Diego, La Jolla, CA, USA. <sup>9</sup>Operational Strategic Initiatives, University of California San Diego, La Jolla, CA, USA. <sup>10</sup>Return to Learn, University of California San Diego, La Jolla, CA, USA. <sup>11</sup>Student Health and Well-Being, University of California San Diego, La Jolla, CA, USA. <sup>12</sup>Student Affairs, University of California San Diego, La Jolla, CA, USA. <sup>13</sup>Academic Affairs, University of California San Diego, La Jolla, CA, USA. <sup>14</sup>Department of Pathology, University of California San Diego, La Jolla, CA, USA.

<sup>15</sup>Scripps Health, San Diego, La Jolla, CA, USA. <sup>16</sup>Rady Children's Institute for Genomic Medicine, San Diego, CA, USA. <sup>17</sup>Sharp Healthcare, San Diego, CA, USA. <sup>18</sup>Helix, San Mateo, CA, USA. <sup>19</sup>County of San Diego Health and Human Services Agency, San Diego, CA, USA.

<sup>20</sup>Herbert Wertheim School of Public Health and Human Longevity Science, University of California San Diego, La Jolla, CA, USA. <sup>21</sup>Division of Infectious Disease and Global Public Health, University of California San Diego, La Jolla, CA, USA. <sup>22</sup>Office of Advanced Molecular Detection, Centers for Disease Control and Prevention, Atlanta, GA, USA. <sup>23</sup>Department of Biomedical Informatics, University of California San Diego, La Jolla, CA, USA. <sup>24</sup>Office of the UC San Diego Health CEO, University of California San Diego, La Jolla, CA, USA. <sup>25</sup>Department of Medicine, University of California San Diego, La Jolla, CA, USA. <sup>26</sup>Sanford Consortium of Regenerative Medicine, University of California San Diego, La Jolla, CA, USA. <sup>27</sup>Department of Cellular and Molecular Medicine, University of California San Diego, La Jolla, CA, USA.

<sup>28</sup>Department of Bioengineering, University of California San Diego, La Jolla, CA, USA. <sup>29</sup>These authors contributed equally: Smruthi Karthikeyan, Joshua I. Levy. <sup>30</sup>These authors jointly supervised this work: Kristian G. Andersen, Rob Knight. ✉e-mail: robknight@ucsd.edu

## Methods

### Ethics declarations

The UCSD Institutional Review Board (IRB) provided human subject protection oversight of the data obtained by the EXCITE laboratory for the campus clinical samples (IRB approval nos. #210699 and #200477). All necessary patient and/or participant consent was obtained and the appropriate institutional forms have been archived, and any sample identifiers included were de-identified. The wastewater component of this project was discussed with our IRB, and was not deemed to be human subject research as it did not record personally identifiable information.

### Wastewater sampling

**High-resolution spatial sampling at the campus level.** One hundred and thirty-one wastewater autosamplers collecting 24-h time-weighted composites were deployed across manholes or sewer cleanouts of 360 campus buildings. GIS-informed analyses as well as agent-based network modelling of SARS-CoV-2 transmission on the UCSD campus enabled identification of most optimal locations for wastewater sampling. During the pilot phase (23 November to 29 December 2020), 68 samplers were prioritized to cover 239 residential buildings identified as the highest risk areas for large outbreaks on campus as a part of an observational study of wastewater monitoring in high-density buildings<sup>29</sup>. This was on the basis of preliminary dynamic modelling, which showed the largest potential outbreaks to occur within the largest residential buildings<sup>9</sup>. In addition to the observational study of wastewater monitoring in these high-density buildings, a cluster randomized study was also performed concurrently. This included a randomized modified version of a stepped wedge crossover design, in which there was random assignment of manholes for wastewater sampling. Clusters of manholes associated with residential buildings were randomized to receive wastewater monitors at one of two-time steps to evaluate the effect of wastewater monitoring on outbreak size in the associated buildings. During the same time period, all students in these residences were mandated to undergo weekly diagnostic testing, which was used to validate the utility of building-level wastewater monitoring. Furthermore, on-campus residences were initially focused due to the relatively static nature of the population, which enabled a more robust cross-validation of the sensitivity and efficacy of the wastewater surveillance. The coverage of wastewater surveillance was then increased to cover the rest of the campus buildings (including non-residential buildings on campus) from January 2021. Four of the deployed wastewater samplers covered the designated isolation and quarantine buildings on campus.

Wastewater composites were collected from the 131 samplers every day for the on-campus residence buildings and Monday through to Friday for the non-residential campus buildings. Wastewater samples ( $n = 19,944$ ) were collected and analysed for the presence of SARS-CoV-2 RNA via qPCR with reverse transcription (RT-qPCR) between 23 November 2020 and 20 September 2021. During this time, 9,700 students lived in campus residences and 25,000 worked on campus on a daily basis. Between October 2020 to 1 January 2021, all on-campus residents were mandated to test on a biweekly (once every 2 weeks) basis and on a weekly basis from 2 January 2021 (start of the winter term). However, fully vaccinated individuals were not mandated to test on a regular basis. Campus protocols required students positive for SARS-CoV-2 living in congregate housing to relocate to designated isolation housing. Accordingly, our analysis of wastewater positivity (Fig. 1b) did not include isolation housing samplers, to control—as best as possible, as a small number of students in non-congregate housing spaces were allowed to isolate ‘in place’, for example—for possible repeat detection due to extended shedding from infected individuals. Automated, localized wastewater-triggered notifications were sent to the residents and/or employees of buildings associated with a positive wastewater signal, which further led to a surge in testing uptake rates by 2- to 40-fold in the associated buildings.

**Wastewater sampling at the county level.** Twenty-four-hour flow-weighted composites were collected three times a week from the main pump station for the Point Loma wastewater treatment plant, which is the primary treatment plant serving the greater San Diego county with a catchment size of approximately 2.3 million. Wastewater samples ( $n = 132$ ) were collected between 24 February 2021 to 7 February 2022.

### Wastewater sample processing and viral genome sequencing

**Sample processing.** SARS-CoV-2 RNA was concentrated from 10 ml of raw sewage and processed as described elsewhere<sup>4</sup>. In brief, the viral RNA was concentrated using an automated affinity capture magnetic hydrogel particle (Ceres Nanosciences Inc.) based concentration method after which the nucleic acid was extracted and sample eluted in 50  $\mu$ l of elution buffer. The extracted RNA was then screened for SARS-CoV-2 RNA via real-time RT–qPCR for three gene targets (N1, N2 and E-gene). PMMoV (pepper mild mottle virus) was also screened to adjust for changes in load. Positive wastewater samples were sequenced within 1–2 weeks of collection, comparable to the delay for clinical samples. To cross-validate the ability of the deconvolution tool in reliably resolving mixtures of strains in wastewater, the wastewater samples from the county as well as the samples from the isolation dorms on campus (where multiple infected individuals were isolating) were also run through a PCR panel targeting eight mutations associated with the strains designated as VOCs. The mutations screened for in wastewater using RT–qPCR included N501Y, DelHV69/70, DelY144, K417N, K417T, E484Q, P681R and L452R (CS3174B02, Promega Corp.).

**Miniaturized wastewater SARS-CoV-2 amplicon sequencing.** The **Swift Normalase Amplicon Panels (SNAP) kit (PN: SN-SX296 (core))** COV-G1V2-96 (amplicon primers), Integrated DNA Technologies was used on RNA from wastewater samples that were positive for SARS-CoV-2 RNA to prepare the multiplex next-generation sequencing (NGS) amplicon libraries and indexed using the SN91384 series of dual indexing oligos, yielding up to 1,536 index pairs per pool. A miniaturized version of the protocol was used with the following modifications: the Super-script IV VILO (Thermo Fisher) cDNA synthesis reaction was scaled down to approximately one-twelfth the normal reaction volume with 0.333  $\mu$ l of enzyme mix and 1.333  $\mu$ l of RNA being used. The multiplex amplicon amplification and Ampure XP bead purification steps were scaled down approximately one-sixth the normal reaction volume. The index adapter PCR and Ampure XP bead purification steps were scaled down to approximately two-thirteenths the normal reaction volume. The final library resuspension volume was 29  $\mu$ l. Of each library, 1  $\mu$ l was pooled for an initial shallow NGS run on a MiSeq (Illumina) using a Nano flow cell. This equal volume pool was used to estimate the differential volumes required for similar read depths across samples using a NovaSeq SP or S4 flow cell (Illumina). Between 5  $\mu$ l and 0.2  $\mu$ l of library material, depending on the data provided from the MiSeq Nano run, was pipetted into a single pool for the NovaSeq run. Transfer volumes were capped at 5  $\mu$ l to reduce pipetting time and because these types of ‘high-volume’ samples typically contained a higher proportion of likely adapter dimers that inhibit flow cell performance for all samples. A Dragonfly Discovery (SPT Labtech) was used to dispense reaction master mixes or water depending on the step. A BlueWasher (BlueCatBio) was used for high-throughput centrifugal 384-well plate washing during the AmpureXP bead reaction cleanup steps. An IKA MS3 Control linear plate mixer (IKA Works Inc.) set to 2,600 r.p.m. for 5 min was used to resuspend the AmpureXP beads during the rehydration steps. A Mosquito Genomics HV16 channel robotic liquid handler (SPT Labtech) was used to dispense the RNA, the reaction master mixes and prepare the equal volume pools for the initial MiSeq Nano (Illumina) balancing runs. A Mosquito X1 single-channel ‘hit picker’ robotic liquid handler (SPT Labtech) was used for the final library balancing for the NovaSeq (Illumina) NGS lanes.

# Article

Sequencing data were analysed using the C-VIEW (COVID-19 Viral Epidemiology Workflow) platform for initial quality control and SARS-CoV-2 lineage assignment and phylogenetics. In brief, sequencing reads were aligned with minimap2 (ref.<sup>30</sup>), and primer sequence trimming and quality filtering were applied using the iVar trim method<sup>20</sup>. Sequencing depth and SNV calls were obtained using samtools mpileup<sup>31</sup> and the iVar variants method<sup>20</sup>.

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

## Virus diversity

As reported previously<sup>20</sup>, virus SNVs were used to characterize the populations derived from wastewater and clinical samples. Richness was defined as the total number of SNV sites, and mean Shannon entropy  $H(p)$  was defined as:

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

where  $p_i$  is the SNV frequency of at the  $i$ -th site, of  $N$  total sites. For statistical testing, a Mann–Whitney  $U$ -test was performed using all wastewater samples that were not sampled from the same source within a 10-day period to maximize independence across samples, as well as all clinical samples. Effect size was calculated using the rank-biserial correlation,  $r = \frac{2U}{n_{ww}n_{cs}} - 1$ , where  $U$  is the Mann–Whitney test statistic and  $n_{ww}$  and  $n_{cs}$  are the numbers of wastewater and clinical samples, respectively.

## Wastewater sample deconvolution

To infer relative abundance within a wastewater sample, we used a ‘barcode’ matrix containing the lineage-defining mutations for each known virus lineage:

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

where  $a_{ij}$  denotes the  $i$ -th lineage at mutation  $j$ . Lineage-defining mutations were obtained from the USHERglobal phylogenetic tree using the matUtils package<sup>15</sup>. Similarly, we let  $b$  and  $d$  encode the frequency of each mutation and the corresponding sequencing depth (using the log-transform  $d_i = \log_2(\text{depth}_i + 1)$  to adjust for large differences in depth across amplicons, which we use to control for heteroskedasticity and down-weight the importance of sites with little or no sequencing depth):

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

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

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

which yields the ‘demixing’ vector  $\hat{x} = [\hat{x}_1 \dots \hat{x}_M]$  that specifies the relative abundances of each of the known haplotypes. Analysis was only

performed on samples with greater than 70% coverage, with the exception of March samples from the UCSD for which all samples with greater than 50% coverage were used. Constrained minimization was performed in Python using the cvxpy convex optimization package<sup>32,33</sup>. Mapping of lineages to variant WHO lineages (VOCs, VOIs and so on) was performed using curated lineage data from outbreak.info<sup>6</sup>. We note that the Epsilon variant received different maximum escalation levels at the US CDC and WHO, which assigned VOC and VOI status, respectively. As the Epsilon variant was widespread in California and much of the United States, we used the more ‘local’ CDC designation.

## Fast bootstrapping method

Bootstrapping was performed at the nucleotide level by resampling each site based on a multinomial distribution of read depth across all sites, where the event probabilities were determined by the fraction of the total sample reads found at each site, followed by a secondary resampling at each site according to a multinomial distribution (that is, binomial when there was only one SNV at a site), where event probabilities were determined by the frequencies of each base at the site, and the number of trials is given by the sequencing depth. Resamplings and demixings ( $n = 1,000$ ) were performed for all samples.

## Spike-in mixture experiment

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

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

**Virus mixing.** RNA standardized in the previous section was used to make a volumetric mixing array (final volume of 10  $\mu$ l) using a Mosquito X1 HV robotic liquid handler (SPT Labtech). Pairwise mixes of 5:95, 10:90, 20:80, 60:40 and 50:50 were made for each virus lineage and in both directions. Equal mixes (20%) for each of the five test strains were made. 25% mixes and 33% mixes were made for a subset of possible combinations and controls of 100:0 were prepared. See Extended Data Table 1 for complete array. Corrected estimates of the fraction of each virus lineage were performed using the final measured Cq values for each pure virus lineage sample to control for issues encountered during the dilution step (repeat Cq measurements had a coefficient of variation of 0.007; Extended Data Table 2). Across all 95 mixtures, we observed a coefficient of variation of 0.016. As initial virus concentrations were controlled for using measured Cq values, we expect that remaining lineage-specific bias (Extended Data Fig. 3) is probably due to experimental inconsistencies encountered during mixture creation.

## Deconvolution method performance comparison

A subset of the spike-in mixtures (one of each type, for a total of 95 mixtures) was used to compare Freyja<sup>34</sup>, cojac (using VOC definitions from the public cojac GitHub repository; lineage A and Epsilon definitions were created manually), the Kallisto-based method from Baaijens et al.<sup>12</sup> and LCS. Kallisto was run using 10 cores (with no bootstrapping), and LCS was run using 16 cores, both on an Intel Xeon processor (2.2 GHz). LCS was run for 48 h, but failed to complete. Timing was performed using the ‘time’ command, and included all steps after alignment, trimming and sorting. Times correspond to total CPU time.

## Estimation of delay in detection frequency

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

## Phylogenetic analyses

**Reconstruction of maximum likelihood trees was performed on all SARS-CoV-2 VOC genomes with 10 $\times$  (10 reads or greater per site) genome coverage of more than 95% and quality score of more than 20 obtained from UCSD campus sampling, using IQtree<sup>35</sup>. This analysis included 150 (112 clinical and 38 wastewater) Epsilon, 49 (37 clinical and 12 wastewater) Alpha, and 160 (136 clinical and 24 wastewater) Delta lineage genomes from the UCSD, in addition to 60 Epsilon, 20 Alpha and 39 Delta randomly selected genomes from elsewhere in the United States. We used iVar<sup>20</sup> to identify consensus sequences for all San Diego samples.** Bases were only included in the sequence if there was a consensus base at the site (more than 50% nucleotide frequency). We also masked known homoplasic sites before tree reconstruction<sup>36</sup>. Analysis of temporal comparison was performed on 608 samples (443 clinical and 165 wastewater, all lineages were included) with 10 $\times$  genome coverage of more than 95% and quality score of 20 from the UCSD. Sample collection SNP distances were calculated without considering ambiguous bases and gaps.

## Statistics and reproducibility

Experiments for retrieving sequences from samples reported in Fig. 5 and Extended Data Fig. 5C were run twice along with positive (spike-in controls of known SARS-CoV-2 lineages derived from mammalian cells as well as heat-inactivated SARS-CoV-2 viral particles in wastewater) and negative controls. Experiments were repeated twice for a batch of 207 wastewater samples. All attempts at replication were successful. For spike-in data reported in Fig. 2 and Extended Data Fig. 3, extraction and RT-qPCR for spike-ins of lineage A from clinical samples were repeated with 20 replicates to check for overall assay variability (reported in Extended Data Table 2).

## Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

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

## Code availability

Freyja is hosted publicly on GitHub (<https://github.com/andersen-lab/Freyja>) and is available under a BSD-2-Clause license (<https://doi.org/10.5281/zenodo.6585067>, all versions). Freyja is accessible as a package via bioconda (<https://bioconda.github.io/recipes/freyja/README.html>) in container form via dockerhub (<https://hub.docker.com/r/andersenlabapps/freyja>). COVID-19 Viral Epidemiology Workflow

(C-VIEW) is available at <https://github.com/ucsd-ccb/C-VIEW> as an open-source, end-to-end workflow for viral epidemiology focused on SARS-CoV-2 lineage assignment and phylogenetics. C-VIEW uses minimap2 (v2.17), samtools(v1.11), iVar(v1.3.1) and pangolin (varying versions).

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.* **73**, 1735–1741 (2021).
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**, 83 (2016).
33. Agrawal, A., Verschueren, R., Diamond, S. & Boyd, S. A rewriting system for convex optimization problems. *J. Control Decision* **5**, 42–60 (2018).
34. Levy, J. I., McDonald, D., Tomkins-Tinch, C. & Pettit, R. A. III. andersen-lab/Freyja: 1.3.7. Zenodo <https://doi.org/10.5281/zenodo.6585068> (2022).
35. 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).
36. Issues with SARS-CoV-2 sequencing data. *Virological* <https://virological.org/t/issues-with-sars-cov-2-sequencing-data/473> (2020).

**Acknowledgements** We thank L. Asato and the Microbiology Laboratory at the SD Public Utilities Department for providing us with county wastewater samples; UCSD's Return to Learn (RTL) program for funding the campus-wide wastewater surveillance efforts; J. Kayne, R. Cota, J. Ortiz and the Facilities management team at the UCSD; J. Mayer from the Center for Aerosol Impacts on Chemistry of the Environment (CAICE); L. Arnold of the Campus Research Machine Shop (CRMS) for assistance with the installation and operation of the autosamplers; R. Jacobs, S. Knepple and their team at UCSD Logistics for assisting with our daily sampling efforts; B. Pollak and the UCSD Information Technology Services team for assisting with the daily notifications; the Office of Academic Affairs for contact tracing and targeted campus messaging assistance; J. Severson, P. Hochstein, the UCSD HDH team and the UCSD Environmental Health and Safety personnel; J. Gilbert and the Microbiome Sample Processing Core at UCSD for access to qPCR equipment; the CDC SPHERES consortium, SEARCH (San Diego Epidemiology and Research for COVID Health) Alliance and members of the Andersen laboratory for discussion and help with logistics; and the healthcare workers, frontline workers and patients who made the collection of this SARS-CoV-2 dataset possible and all those who made genomic data available for analysis via GISAID. This work has been funded by CDC BAA contracts 75D30121P10258 (Helix) and 75D30120C09795 (to G.W.Y., R.K., L.C.L. and K.G.A.), NIH NIAID 3U19AI135995-03S2 (to K.G.A.), U19AI135995 (to K.G.A.), U01AI151812 (to K.G.A.), NIH NCATS UL1TR002550 (to K.G.A.), the Conrad Prebys Foundation (to K.G.A.), NIH 5T32AI007244-38 (to J.I.L.), NIH Pioneer grant 1DP1AT010885 (to R.K.), NSF RAPID 2029069 (to R.K.), San Diego County Health and Human Services Agency (to R.F.-M.) and NIH S10OD026929 (to K.J.). The findings and conclusions in this report are those of the authors and do not necessarily represent the official position of the CDC and California Department of Public Health or the California Health and Human Services Agency. Use of trade names is for identification only and does not imply endorsement by the CDC.

**Author contributions** R.K. and K.G.A. conceptualized the study. S.K., J.I.L., N.K.M., P.D.H., A. Birmingham, S. Sathe, K.M.F., L.C.L., G.W.Y., K.G.A. and R.K. contributed the methodology. J.I.L., D.M., N.M., K.M.F., A. Birmingham, K.J., B.H., S. Sathe, K.G., N.L.M., K.S.R., C.M.A., E.H. and A.M.M. provided software. S.K. and J.I.L. performed formal analysis. S.K., J.I.L., N.M., S.F., H.M.T., T.V., C.E.T., R.T., N.A.B., T.B., M.C., W.C., E.S.C., E.R.E., A. Hakim, G.H., A.L.L., E.S.L., T.T.N., T.O., A. Plascencia, R.A.S., P.S., P.B.-F., E.W.S., S.A., P.D.H., C.A.M., L.C.L., G.W.Y., Cheryl Anderson, E.K., M.A.S., S.A.P., J.L., E.P., M.Z., E.S., R.F.-M., T.G., R.S.G., K.G.A. and R.K. conducted the investigation. Catelyn Anderson, N.K.M., R.M.N., R.S., E.H.S., A.M. Schafer, S.F.K., D.P.D., C.A.H., A.M., S. Strous, B.A., S. Shah, N.G., J.D.M., E.M., I.H.M., A. Harding, O.B., A.M., A. Bolze, K.M.S.B., E.T.C., N.L.W., W.L., M.I., D. Becker, L.N., S.W., M.Z., R.R.-S., R.F.-M., T.G., R.S.G., D.B.A., D. Brenner, J.C.B., A. Buckley, M.E., J.G., S.L.G., M.H., F.H., L.I., H.J., T.J., V.K., B.K., L.R., C.A.H., G.M., P.M., R.F.-M., P.O. Park, A. Pradens, A.M.S., B.A., S. Song, B.W., T.W., S.R., P.K.K., A.T.Y., D.R.M., F.H., G.M., R.W.M. and C.L. provided resources. S.K., J.I.L., P.D.H., G.H., S.F., H.M.T., C.E.T., R.T., T.V., A. Birmingham, N.M., A.M.M. and K.M.F. curated data. S.K., J.I.L., K.G.A. and R.K. wrote the original draft of the manuscript. All authors reviewed and edited the manuscript. S.K. and J.I.L. performed visualization. R.M.N., N.K.M., R.S., A.L.S., E.H.S., A.M. Satterlund, P.D.H., L.C.L., D. Pride, G.W.Y., K.G.A. and R.K. supervised the study. R.M.N., N.K.M., R.S., A.L.S., E.H.S., A.M. Satterlund, P.D.H., L.C.L., G.W.Y., K.G.A. and R.K. undertook project administration. R.K. and K.G.A. acquired funding.

**Competing interests** A. Bolze, D. Becker, E.T.C., M.I., K.M.S.B., N.L.W. and W.L. are employees of Helix. K.G.A. has received consulting fees for advising on SARS-CoV-2, variants and the COVID-19 pandemic.

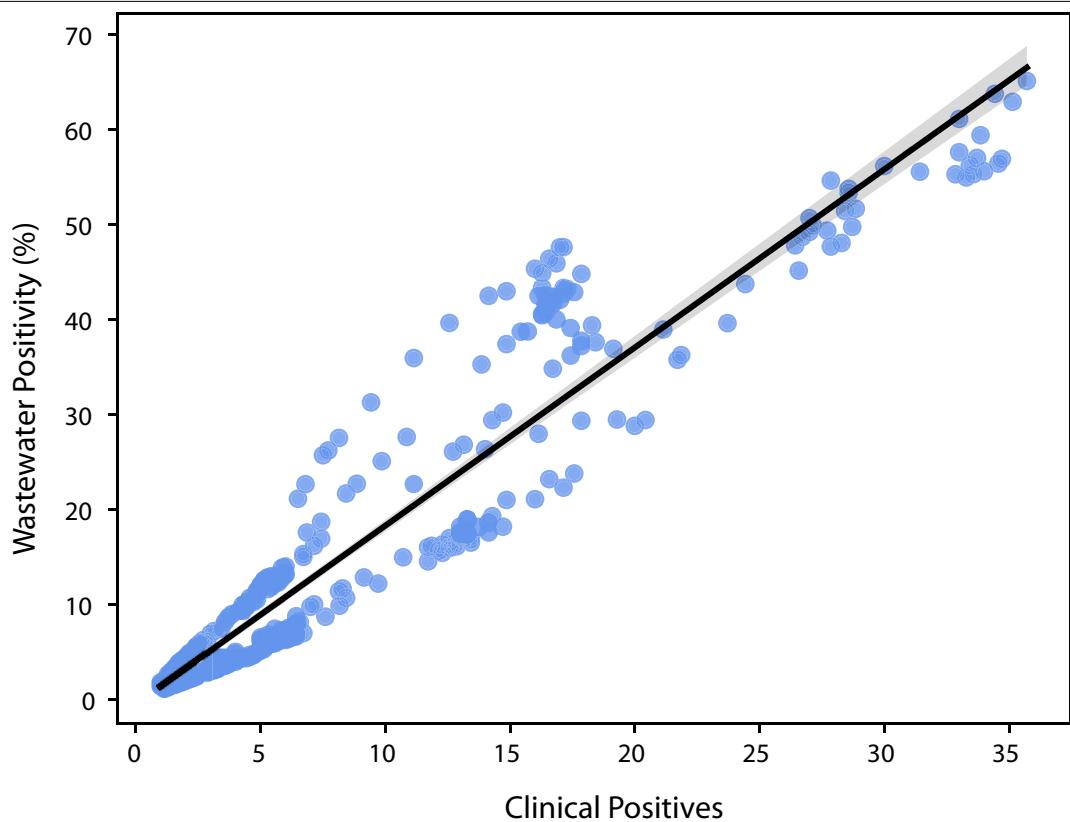
## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41586-022-05049-6>.

**Correspondence and requests for materials** should be addressed to Rob Knight.

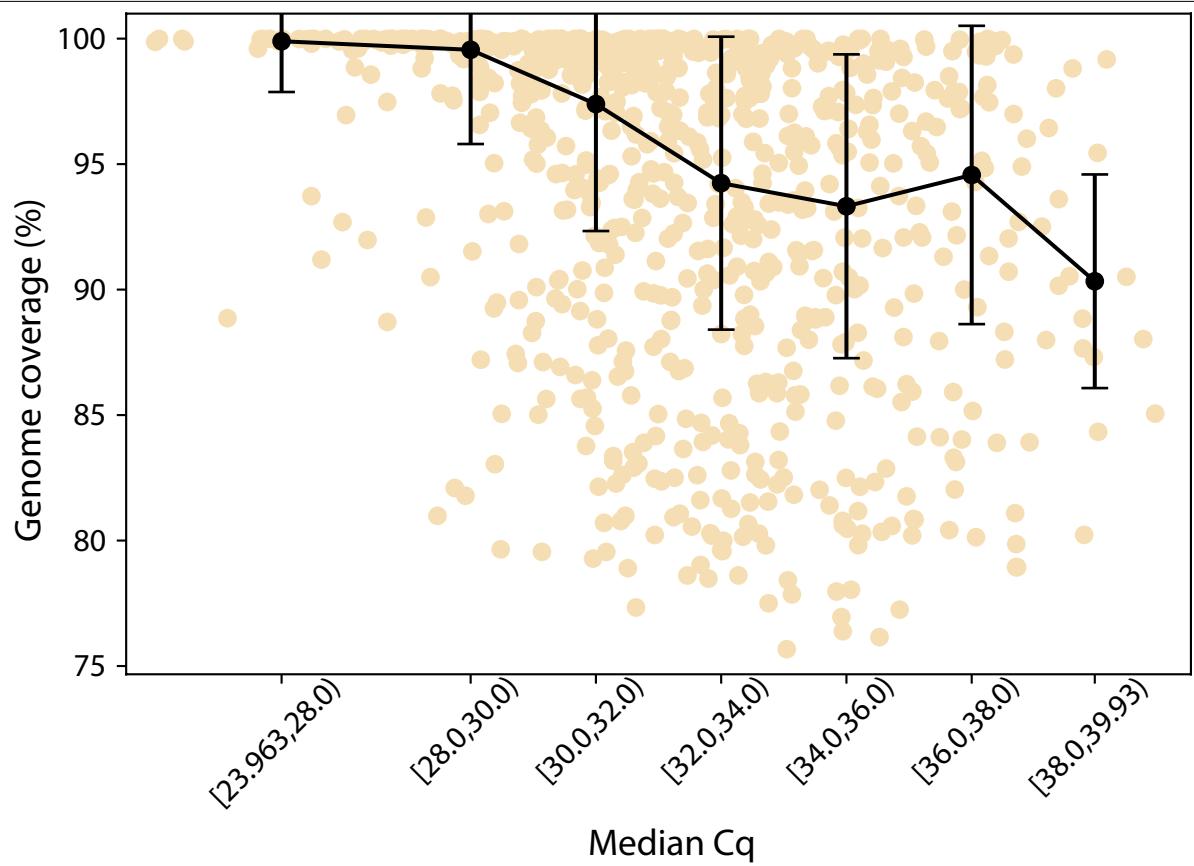
**Peer review information** *Nature* thanks Timothy Julian and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>.

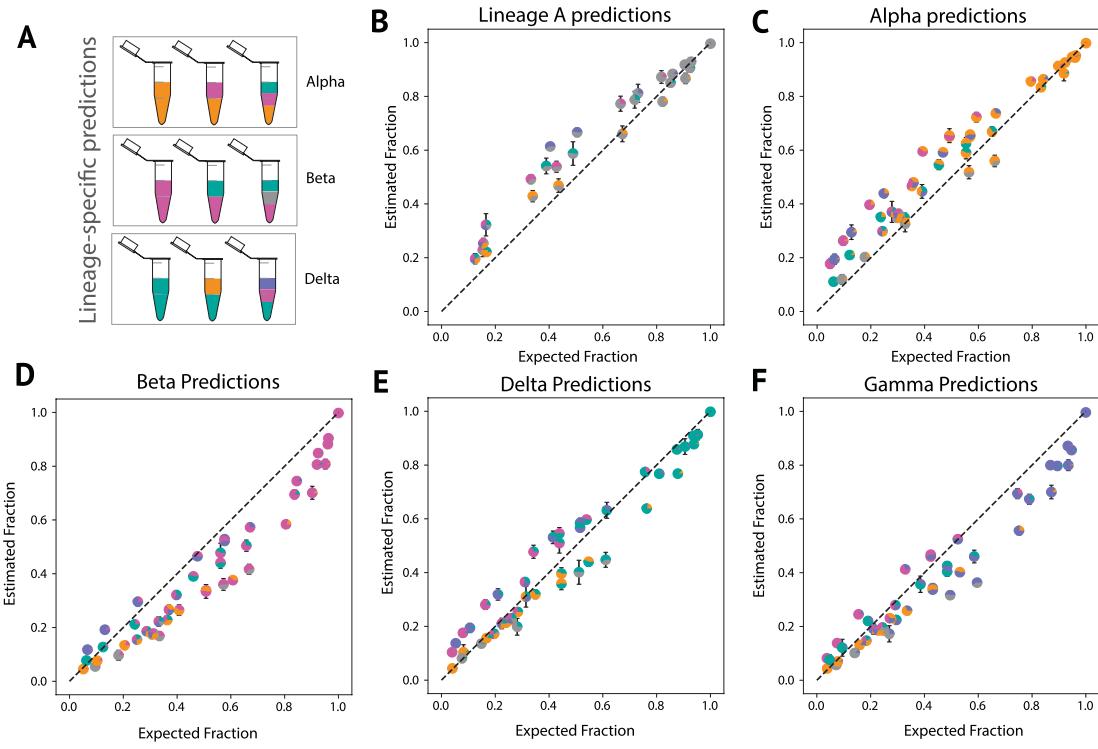


**Extended Data Fig. 1 | Relationship of daily UCSD campus wastewater sampler positivity and campus clinical positives.** Black line indicates the linear regression fit (slope = 1.88 %/clinical positive, intercept = -0.45%) to the

data ( $n = 321$ ), with bootstrap 95% confidence interval (resampled 1000 times with replacement) shown in gray (median slope = 1.88%/clinical positive, intercept = -0.47%).

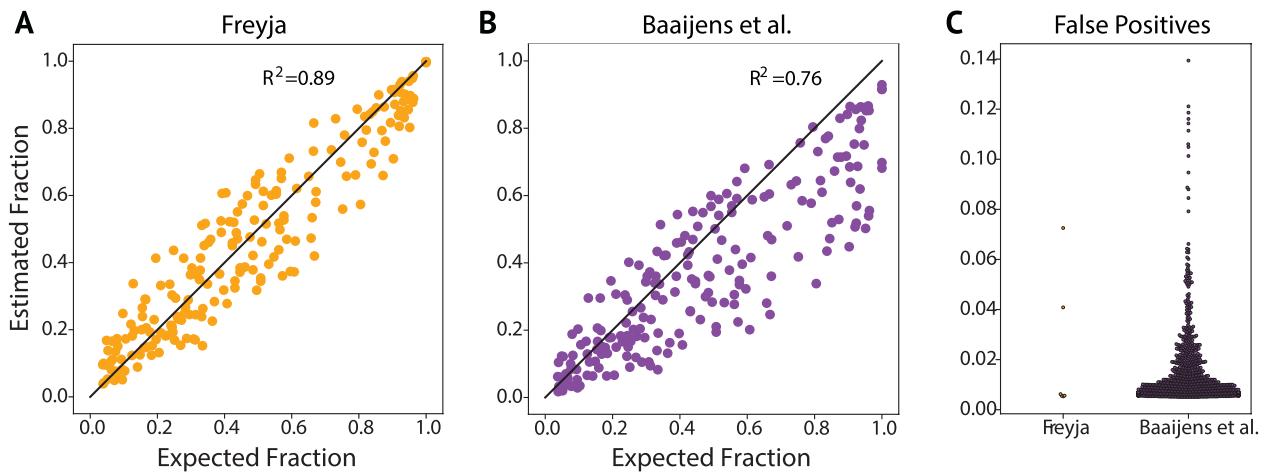


**Extended Data Fig. 2 | Relationship between genome coverage and cycle quantification values.** 10x genome coverage (fraction of sites with 10 reads or greater) remains high, even for Cq values of nearly 38 ( $n = 786$ ). Points indicate median value in each bin, while error bars indicate the median absolute deviation.



**Extended Data Fig. 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 ( $n = 380$ , four samples per mixture type).

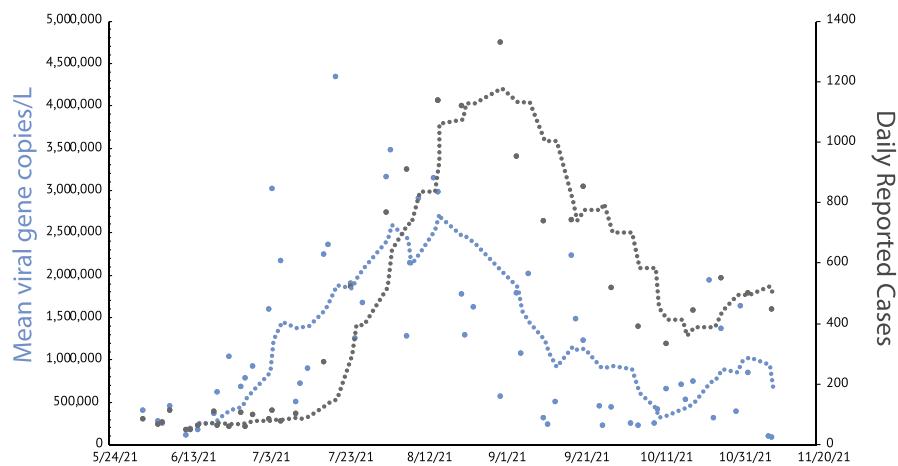


**Extended Data Fig. 4 | Freyja more accurately estimates virus abundance, with fewer false positives.** A-B. Estimated vs expected fraction of each lineage in the mixture ( $n = 95$ , one sample per mixture type). The Kallisto-based

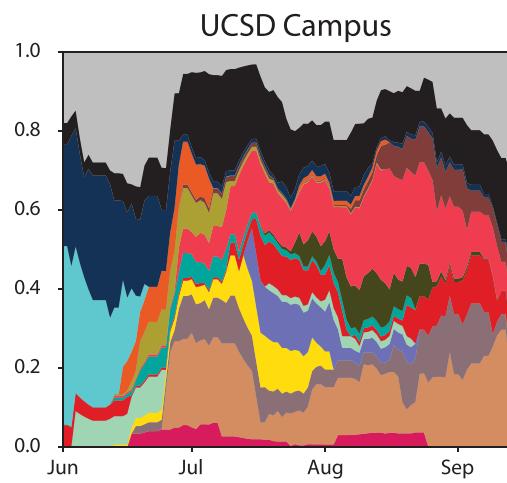
approach from Baaijens et. al shows a wider range of estimates for each known mix fraction, and generally underestimates the fraction. C. False positives with abundance greater than 0.5%.

# Article

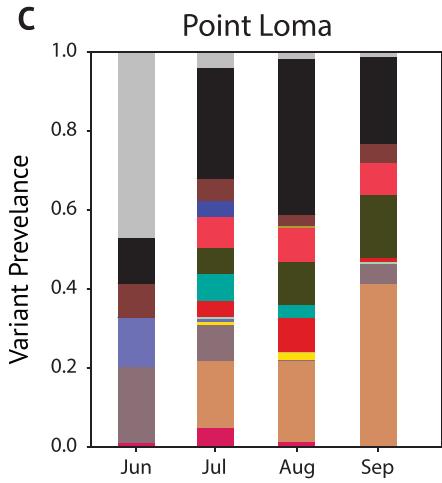
A



B



C

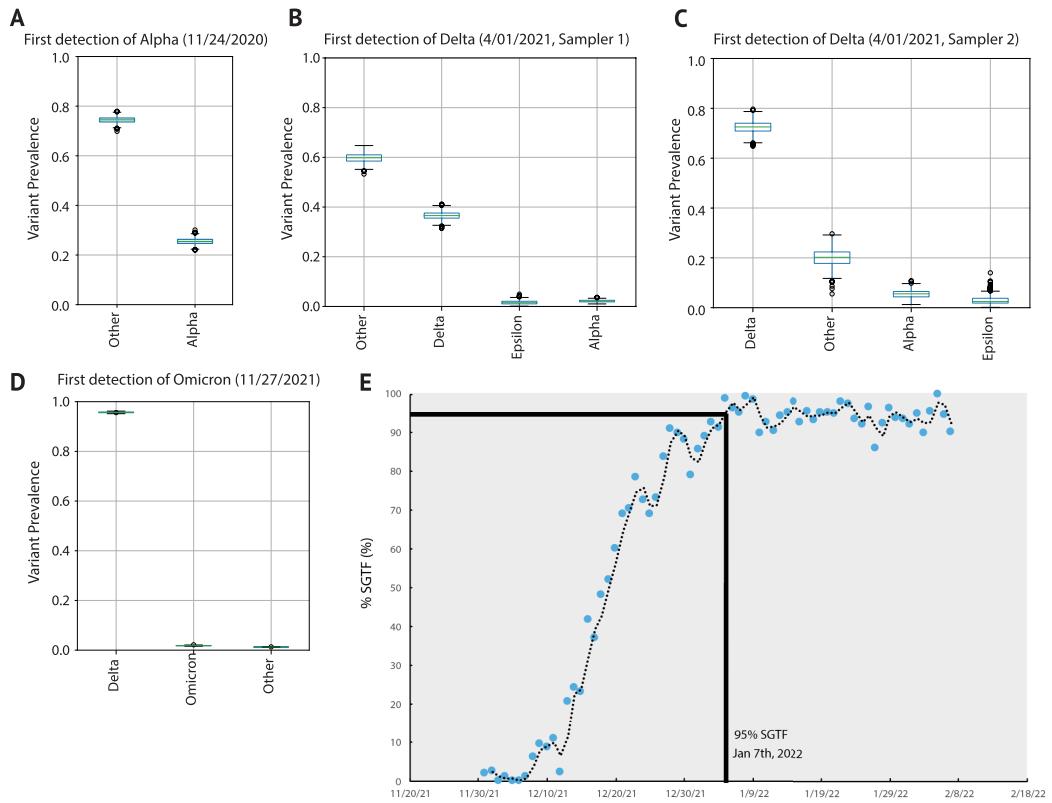


## Extended Data Fig. 5 | The rise of the Delta variant during Summer 2021.

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.

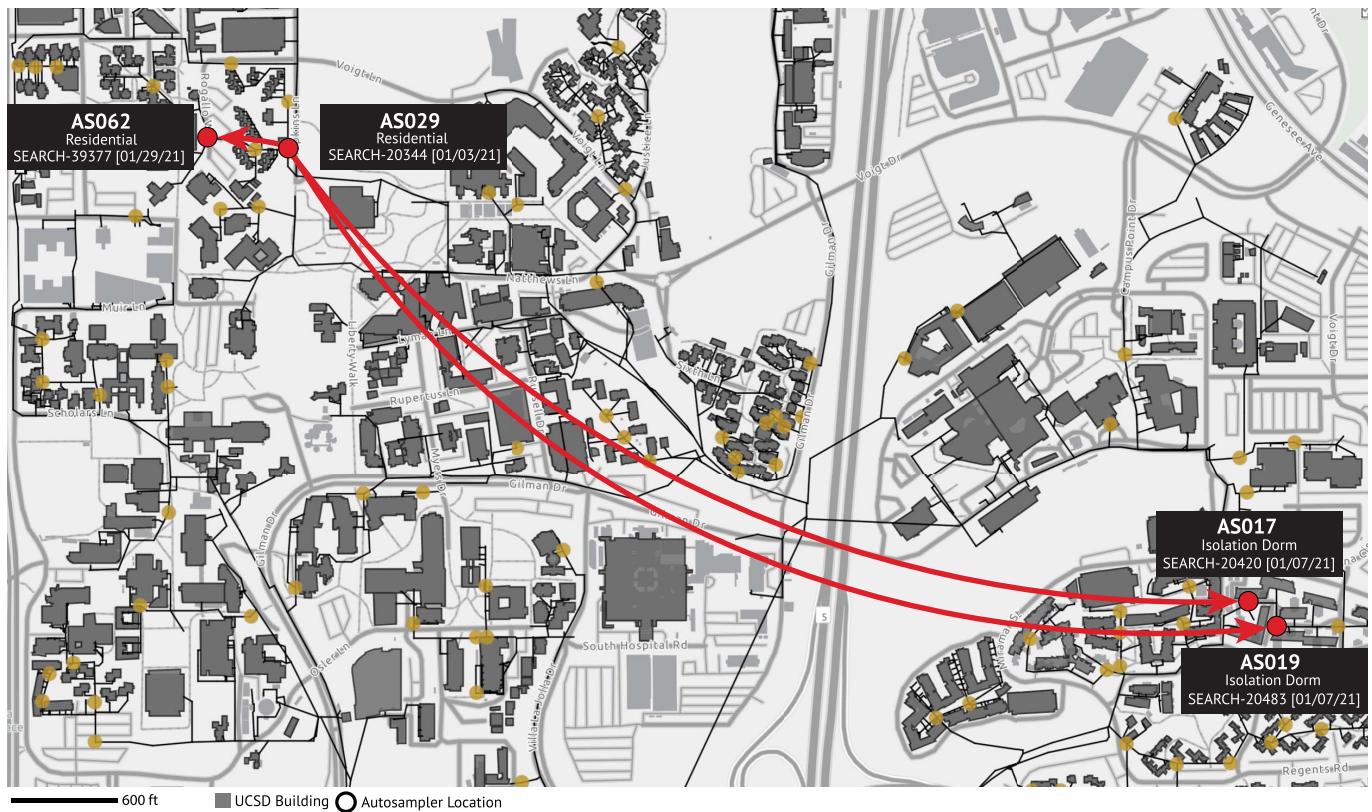
B. Lineage distribution in UCSD campus wastewater. C. Monthly lineage averages for wastewater collected at Point Loma Wastewater Treatment Plant during the Delta surge (N = 5, 20, 25, 7).



**Extended Data Fig. 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 1000 times with replacement. In all boxplots, box edges specify the first and third quartiles, solid line indicates the median, and whiskers delimit the maximum and minimum values within 1.5 times the inter-quartile range (IQR) of box edges. Outliers are denoted with individual markers. Two samplers were found to contain Delta on the same day. First detections were also confirmed using a VOC qPCR panel, as shown in Fig. 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). **E. 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.

# Article



**Extended Data Fig. 7 | Temporal and spatial dynamics of an Epsilon outbreak at UCSD.** After initial detection on January 3rd 2021, infected individuals were transferred to isolation housing where they continued to shed virus. At the end of January, a matching virus was detected in a residence nearby

the original site of detection. All four samples have perfectly matching virus genomes. Maps are the intellectual property of Esri and its licensors and are used herein under license. Copyright © 2022 Esri and its licensors. All rights reserved.

**Extended Data Table 1 | Plate map of spike-in mixtures used for method validation**

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

# Article

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

Replicate Number	Cq N Gene	Cq Orf1ab	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

Corresponding author(s): Rob Knight

Last updated by author(s): May 30, 2022

## Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided  
*Only common tests should be described solely by name; describe more complex techniques in the Methods section.*
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g.  $F$ ,  $t$ ,  $r$ ) with confidence intervals, effect sizes, degrees of freedom and  $P$  value noted  
*Give  $P$  values as exact values whenever suitable.*
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's  $d$ , Pearson's  $r$ ), indicating how they were calculated

*Our web collection on [statistics for biologists](#) contains articles on many of the points above.*

### Software and code

Policy information about [availability of computer code](#)

Data collection	All data collection code is publicly available in our COVID-19 Viral Epidemiology Workflow (C-VIEW) at <a href="https://github.com/ucsd-ccb/C-VIEW">https://github.com/ucsd-ccb/C-VIEW</a> as an open-source, end-to-end workflow for viral epidemiology focused on SARS-CoV-2 lineage assignment and phylogenetics. C-VIEW uses minimap2 (v2.17), samtools(v1.11), iVar(v1.3.1), and pangolin (varying versions).
Data analysis	All data analysis is performed using Freyja (v.1.3.7) as well as custom Python 3 scripts. Freyja is hosted publicly on github ( <a href="https://github.com/andersen-lab/Freyja">https://github.com/andersen-lab/Freyja</a> ) and is available under a BSD-2-Clause License (doi: 10.5281/zenodo.6585067, version 1.3.7). Freyja is accessible as a package via bioconda ( <a href="https://bioconda.github.io/recipes/freyja/README.html">https://bioconda.github.io/recipes/freyja/README.html</a> ) in container form via dockerhub ( <a href="https://hub.docker.com/r/andersenlabapps/freyja">https://hub.docker.com/r/andersenlabapps/freyja</a> ).

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

## Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

All raw wastewater sequencing data is available via the NCBI Sequence Read Archive under the BioProject ID PRJNA819090. Spike-in sequencing data is available via google cloud ([https://console.cloud.google.com/storage/browser/search-reference\\_data](https://console.cloud.google.com/storage/browser/search-reference_data)). The UCSD campus dashboard can be accessed at <https://returntoteach.ucsd.edu/dashboard/>. The county wastewater data from Point Loma are available through the public dashboard that can be accessed at <https://searchcovid.info/dashboards/wastewater-surveillance/>. The SEARCH genomic surveillance dashboard is available at <https://searchcovid.info/dashboards/sequencing-statistics/>. Consensus sequences from clinical and wastewater surveillance are all available on GISAID. Further details are provided here: <https://github.com/andersen-lab/HCoV-19-Genomics>.

## Human research participants

Policy information about [studies involving human research participants](#) and [Sex and Gender in Research](#).

Reporting on sex and gender

[Sex/gender-based data were not considered in this study design.](#)

Population characteristics

No covariate relevant population characteristics (age, gender etc.) was used in this study

Recruitment

Participants were recruited on the basis that they had a positive COVID-19 test at UCSD, as approved under a waiver of consent framework, or on the basis that they volunteered for an associated study. Both of these pathways for recruitment were IRB-approved (see below).

Ethics oversight

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

Note that full information on the approval of the study protocol must also be provided in the manuscript.

## Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences       Behavioural & social sciences       Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://nature.com/documents/nr-reporting-summary-flat.pdf)

## Ecological, evolutionary & environmental sciences study design

All studies must disclose on these points even when the disclosure is negative.

Study description

To compare the relative utility of wastewater genomic surveillance in enabling early detection of SARS-CoV-2 VOCs and transmission dynamics, a high-resolution 295-day wastewater and clinical sequencing effort was performed in the controlled environment of a large university campus and the broader context of the surrounding county.

Research sample

Between November 2020–September 2021, 21,383 wastewater samples were collected from UCSD campus and the surrounding county and analyzed for presence of SARS-CoV-2 viral RNA. The campus wastewater samples covered 343 campus buildings representing a catchment of approximately 25,000 students and staff (capturing all on-campus residents). The wastewater treatment plant at Point Loma covers a catchment of 2.2 million residents in the greater San Diego county. The samples from the public school districts cover a combined population of 10,000 students and staff. Samples positive for the RNA as determined via qPCR were sequenced using a miniaturized tiled amplicon approach. During the same period of study, clinical samples were also sequenced from the UCSD campus as well the surrounding county. Sequencing of 600 campus wastewater samples were compared to 759 genomes obtained from campus clinical swabs (processed by the CALM and EXCITE CLIA labs at UCSD). 31,149 genomes obtained from clinical genomic surveillance of the greater San Diego community were compared to 837 wastewater sample sequences collected from San Diego county.

Sampling strategy

Sampling sites were identified via campus GIS, which provides mapping and flow direction of interconnected sewer lines and locations of manholes where samplers could be placed. In the UCSD campus, wastewater autosamplers were first prioritized to the most proximal manhole to buildings with residential populations greater than 150. This decision was made based on agent-based

network modeling of SARS-CoV-2 transmission on the UCSD campus, indicating that the highest risk areas for large outbreaks on campus were buildings containing the largest residential populations (details provided in Goyal et al 2021, see References). 131 wastewater autosamplers were placed optimally to cover 343 campus buildings thereby capturing the majority of the campus population. For the San Diego county samples, the main wastewater treatment plant was chosen as the sampling site since it captured a majority (2.2 million) of the area's residents. For the public school districts, pilot sites were selected from ZIP Codes with COVID-19 rates above the county median and with high levels of social vulnerability according to the California Healthy Places Index. The samplers for the campus and school districts were placed at manholes or sewer cleanouts such that they captured the waste from all the schools/building's residents.

Sampling procedures are described in detail at: [dx.doi.org/10.17504/protocols.io.bshvnb66](https://doi.org/10.17504/protocols.io.bshvnb66).

#### Data collection

Sample data logging was streamlined via integration with the campus GIS (geographic information system) server. The spatially enabled sewer network and subsequent trace of samplers to buildings were stored in and performed by ArcGIS Pro 2.7 (Esri). Details on sample collection and data integration are also provided at: <https://doi.org/10.1128/mSystems.00793-21>. The unique autosampler barcode and the sample bottle barcodes were scanned by the sample collection staff using the ArcGIS Survey123 mobile app (ESRI) which enabled automatic data integration into the ArcGIS Online environment for trace analysis.

#### Timing and spatial scale

During the 10 month study period, wastewater samples were collected on a daily basis from the UCSD campus from 131 sampling sites covering nearly 350 campus buildings. Wastewater samples were collected 5 days a week from the public schools in the various San Diego school districts and 3 times a week from the Point Loma wastewater treatment plant (the primary plant serving the greater San Diego area). This is part of an ongoing study and samples continue to be collected from these sites regularly as of May 2022. The spatial scale of the campus data is shown in Fig. 1A of the current manuscript. The spatial scale of the county samples are provided in Karthikeyan et al., 2021 (doi: 10.1128/mSystems.00045-21) and Fielding-Miller et al., 2021 (doi: /10.1101/2021.10.19.21265226)

#### Data exclusions

All wastewater sequences used had greater than 70% coverage, with the exception of March samples from UCSD for which all samples with greater than 50% coverage were used due to low sample numbers during that period. No data were excluded from the study unless they failed to meet the quality threshold specified above (70% coverage of the SARS-CoV-2 genome with no evidence of cross-contamination as well as the positive and negative controls passing QC for the specific run). 4.7% of the sequenced samples failed to meet the QC threshold.

#### Reproducibility

Controls were included at all stages of sample processing (viral concentration, extraction, qPCR and sequencing) to assess potential inhibition and cross contamination. Most of the sample processing steps were performed by liquid handling robots to minimize human error and replicates included. If any of the controls failed or indicated cross-contamination, the entire batch was rerun. The clinical samples and wastewater samples were processed separately for sequencing due to significant differences in viral load between the two sample types. Due to the sample heterogeneity in complex environmental matrices such as wastewater, controls are vital in aiding of data interpretation. With every sequencing run spike-in controls of a known SARS-CoV-2 lineage (Lineage A) was chosen as a positive control and a no template extraction blank was used as a negative control to assess cross-contamination during the library generation and the sequencing stage. If significant SARS-CoV-2 mapping reads were found passing QC in the negative controls, the entire batch was rerun from the library generation step. Experiments for retrieving sequences from samples reported in Fig. 5 and Extended Data Figure 5C were run twice along with positive (spike-in controls of known SARS-CoV-2 lineages derived from mammalian cells as well as heat-inactivated SARS-CoV-2 viral particles in wastewater) and negative controls. Experiments were repeated twice for a batch of 207 wastewater samples. All attempts at replication were successful. For spike-in data reported in Fig 2 and Extended Data Fig. 3, extraction and RT-qPCR for spike-ins of Lineage A from clinical samples were repeated with 20 replicates to check for overall assay variability (reported in Extended Data Table 2). Detailed wastewater sample processing steps are also provided here: [dx.doi.org/10.17504/protocols.io.bshvnb66](https://doi.org/10.17504/protocols.io.bshvnb66) as well as in the Methods section of the manuscript

#### Randomization

All sequences/samples that met our threshold for quality were used in the study. Randomization was not required since we did not perform experiments to evaluate the effects of specific treatments or groups.

#### Blinding

Blinding was not applicable to this study as we did not perform experiments involving specific treatments or groups and all clinical data used were de-identified from the point of receipt.

Did the study involve field work?  Yes  No

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

**Materials & experimental systems**

n/a	Involved in the study
<input checked="" type="checkbox"/>	Antibodies
<input checked="" type="checkbox"/>	Eukaryotic cell lines
<input checked="" type="checkbox"/>	Palaeontology and archaeology
<input checked="" type="checkbox"/>	Animals and other organisms
<input checked="" type="checkbox"/>	Clinical data
<input checked="" type="checkbox"/>	Dual use research of concern

**Methods**

n/a	Involved in the study
<input checked="" type="checkbox"/>	ChIP-seq
<input checked="" type="checkbox"/>	Flow cytometry
<input checked="" type="checkbox"/>	MRI-based neuroimaging