

QualD: Enabling Earlier Detection of Recently Emerged SARS-CoV-2 Variants of Concern in Wastewater

Nicolae Sapoval¹, Yunxi Liu¹, Esther G. Lou², Loren Hopkins^{3,4}, Katherine B Ensor⁴, Rebecca Schneider³, Lauren B Stadler², Todd J Treangen¹

¹ Department of Computer Science, Rice University, 6100 Main Street, Houston, TX 77005, USA

² Department of Civil and Environmental Engineering, Rice University, 6100 Main Street, Houston, TX 77005, USA

³ Houston Health Department, 8000 N. Stadium Dr., Houston, TX 77054

⁴ Department of Statistics, Rice University, 6100 Main Street, Houston, TX 77005, USA

Abstract

As clinical testing declines, wastewater monitoring can provide crucial surveillance on the emergence of SARS-CoV-2 variants of concern (VoC) in communities. Multiple recent studies support that wastewater-based SARS-CoV-2 detection of circulating VoC can precede clinical cases by up to two weeks. Furthermore, wastewater based epidemiology enables wide population-based screening and study of viral evolutionary dynamics. However, highly sensitive detection of emerging variants remains a complex task due to the pooled nature of environmental samples and genetic material degradation. In this paper we propose quasi-unique mutations for VoC identification, implemented in a novel bioinformatics tool (QualD) for VoC detection based on quasi-unique mutations. The benefits of QualD are three-fold: (i) provides up to 3 week earlier VoC detection compared to existing approaches, (ii) enables more sensitive VoC detection, which is shown to be tolerant of >50% mutation drop-out, and (iii) leverages all mutational signatures, including insertions & deletions.

Introduction

Wastewater monitoring is an invaluable tool for SARS-CoV-2 surveillance¹⁻⁸. Despite multiple recent successes in VoC monitoring and detection from wastewater sequencing data⁹⁻¹⁵, there are multiple challenges associated with the nature of the environmental data. Since wastewater represents a pooled sample of multiple hosts, it harbors a diversity of SARS-CoV-2 variants that are currently circulating in the population^{1,2,10,13}, including potentially previously unreported genotypes¹⁶. Furthermore, variant detection and phasing is further complicated by uneven genome coverage^{2,17,18} and environmental RNA degradation^{5,19,20} which render phased assembly extremely difficult²¹⁻²³. Despite these challenges, detection of VoCs in wastewater samples is important for monitoring the emergence and spread of variants and informing public health response^{4,5,11,24,25}. Current approaches for VoC detection in wastewater samples typically require sufficient depth and breadth of coverage of the variant genomes^{9,12}, and therefore depend on a large fraction of the sample representing the variant genotype¹⁰, hampering early detection. Furthermore, most of the current approaches discard insertion and deletion (indel) information and only rely on single nucleotide variants (SNVs) associated with the VoC^{9,12}. Finally, all approaches that rely on a database of previously collected SARS-CoV-2 genomes are biased by the contents of the database^{26,27}, which can lead to both false negative and false positive calls at the inference stage²⁸. This issue can be further amplified when the underlying database is not scrutinized for potential metadata errors^{29,30}.

To address these issues, we developed QualD: a computational pipeline for analyzing SARS-CoV-2 wastewater sequencing data and inferring presence of VoCs. We use empirical results on real Houston wastewater data and simulated data to validate the efficacy of our software, and compare it to Freyja⁹, another state-of-the-art tool for VoC detection. Our key goal is achieving sensitivity to newly emerging variants in scenarios where coverage breadth and depth can be uneven, and the VoC-associated genomes are present at low abundances. We also leverage the indel data that can be inferred from

the multiple sequence alignments to further improve the robustness of our VoC detection approach.

Results

Between February 23, 2021 and May 5th, 2022 we collected, processed, and analyzed 2,637 wastewater samples from the fifth-most populous metropolitan area in the US: Houston, Texas. Samples were collected weekly from 39 wastewater treatment plants (WWTPs, Supplementary Table 1, Supplementary Figure 1) distributed throughout the city of Houston and servicing more than 2 million Houston residents¹⁸. During the study period, the VoC detection signal clearly reflected the three major variants that affected Houston - Alpha, Delta, and Omicron (Figure 1A). QualD was able to detect the Delta VoC two weeks prior to the first sequenced clinical sample in Texas (marked by star in Figure 1B) and continued to provide detection signal for the four subsequent weeks after the first sequenced clinical sample (2021-04-05 to 2021-05-03).

In contrast, Freyja reliably picked up the Delta signal only once the VoC became more prevalent. Similarly for the Omicron VoC, QualD detected the presence of the variant in wastewater two weeks prior to the first clinical sample collection date, while Freyja required an additional week after the first clinical sample to detect Omicron presence.

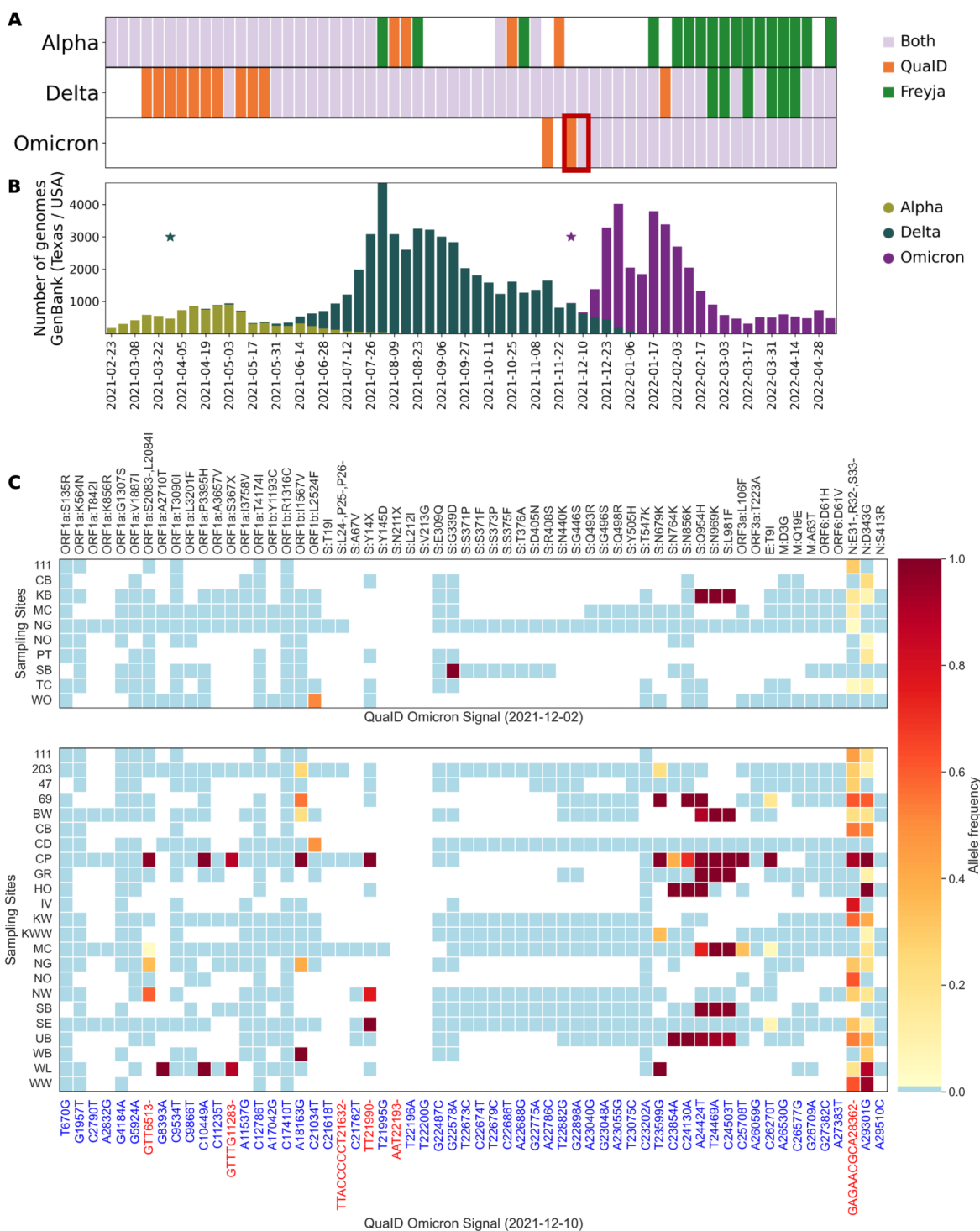


Figure 1. Detection of Alpha, Delta, and Omicron VoCs in Houston, TX wastewater. **A.** Early detection of the emerging variants of concern in Houston wastewater provided by QualID and Freyja pipelines. For Omicron and Delta variants, QualID provided earlier detection. Each week is presented as the aggregate signal from the 39 WWTPs with detections being reported if at least 2 WWTPs had any QualID signal or had any non-zero abundance of the VoCs reported by Freyja. **B.** Variant prevalence in the clinical data over the study period obtained from GenBank and restricted to Texas. Stars indicate the first occurrences of a Delta variant genome (yellow) and an Omicron variant genome (red). **C.** Heatmaps of WWTPs with detected Omicron variant quasi-unique mutations the week of December 2nd, 2021 (top) and December 10th, 2021 (bottom) in Houston. Blanks indicate lack of sequencing data, blue color indicates no mutation detected, and the gradient shows the allele frequency for detected mutations.

We further investigated the early detection of the Omicron variant in Houston wastewater by visualizing a heatmap of variant calls (Figure 1C) and examining the multiple sequence alignment (MSA) of SARS-CoV-2 Omicron variant genomes available on GISAID³¹ in early December 2021. We observed 50% (5 out of 10) of the samples with Omicron presence for the week of December 2nd, 2021 contained the 9bp deletion (N:DEL31/33), which is a stable mutation (95.1% prevalence among all Omicron genomes³²) for the Omicron variant (Figure 1C). Since the current version of Freyja relies on the UShER³³ phylogenetic tree for its designation of mutational signatures, no deletions are used in the inference process, highlighting one of the reasons for the delayed detection of the Omicron variant. In the subsequent week, December 10th, 2021, when both Freyja and QualD reported the presence of the Omicron variant in the wastewater, the N:DEL31/33 mutation was present in 16 of 23 sites with detections (Figure 1C), and for one of the samples with no deletion there was no coverage in the region flanking the deletion (Figure 1C, Sampling Site SB: Sims Bayou North).

To further examine sensitivity of the QualD and Freyja to degradation of the sequencing data, we constructed a simulated data experiment in which we varied the fraction of SNVs dropped out from the variant calling results. In the real wastewater sequencing data, 37.7% of all samples had less than 25% of the SNVs associated with the Omicron VoC via UShER barcodes covered by at least one read (Supplementary Figure 2B), and 24.4% of samples had less than 10% of all Omicron-associated SNVs with at least one read. Thus, we constructed three simulation scenarios with each retaining 10%, 25%, or 50% of all SNV calls at random. Our results show that due to the inclusion of deletion information in the inference process, QualD remained sensitive even when only 10% of all SNV calls were retained, while Freyja required at least 50% of the calls to be included to reliably detect the VoC presence. In particular, when only 10% of all SNV calls were retained, QualD still detected the presence of Delta and Omicron VoCs reliably, and Alpha and Gamma VoCs sparsely, while Freyja failed to estimate the abundance of any of the VoCs (Alpha, Delta, Gamma, and Omicron) present in the simulated samples (Figure 2A, Supplementary Figures 3A-6A). Furthermore, when 25% of all SNVs were retained, QualD identified the present VoCs in the majority of the simulated samples, while Freyja provided sparse detection in the samples dominated by a single VoC (Figure 2B). Finally, when 50% of all SNVs are retained, Freyja detected most of the VoCs present in the samples, and in several instances recovered the correct relative abundance. However, even in this scenario 8 Omicron dominated samples failed to be correctly identified by Freyja, while QualD correctly inferred the presence of the VoC. Additionally, we observe that the stability of the coverage for the N:DEL31/33 is further empirically supported by our data, which indicated that among all samples more than 61% have at least 10 reads that cover the bases immediately flanking the deletion (Supplementary Figure 2C).

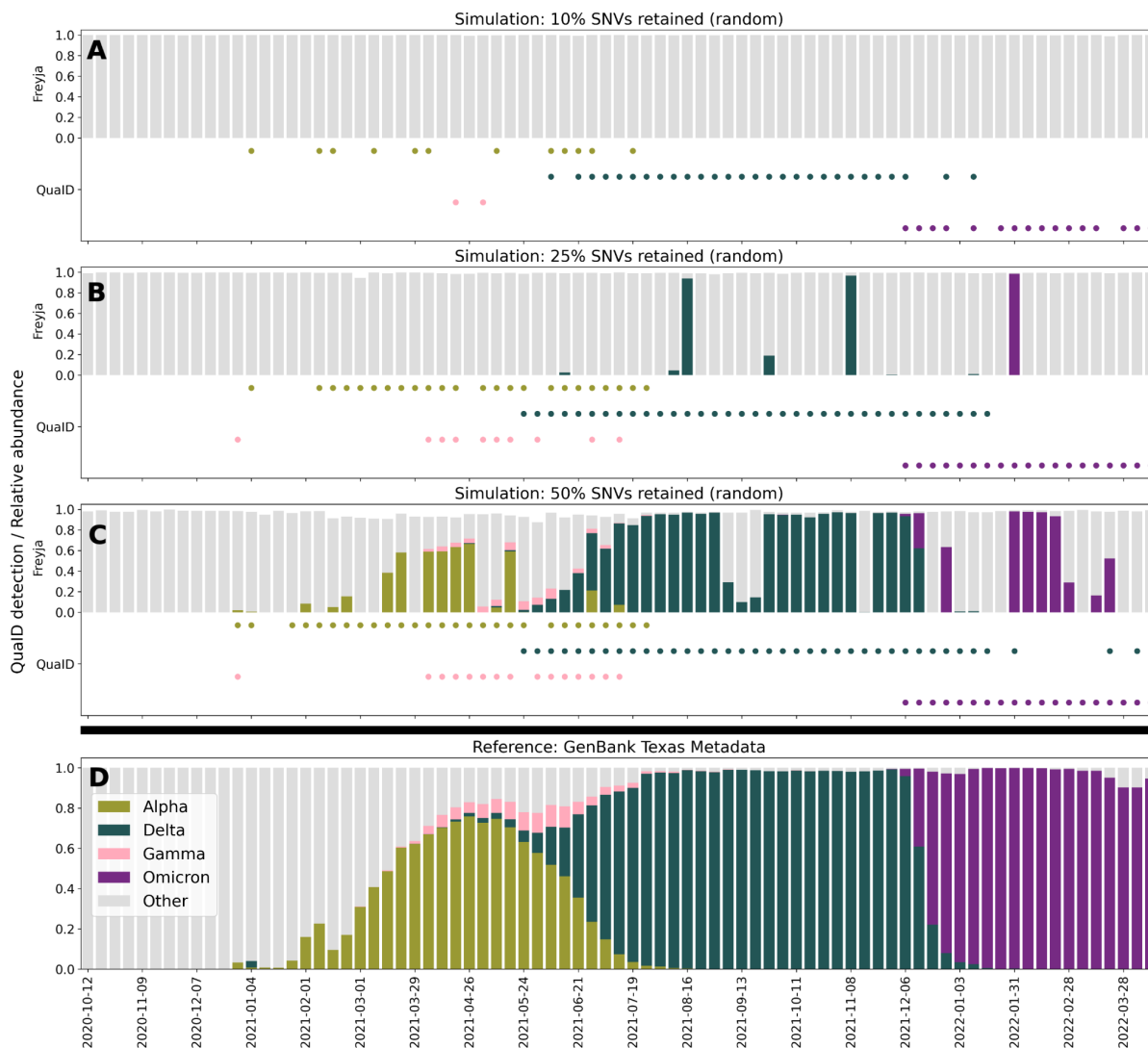


Figure 2. Detection of VoCs in simulated data at various levels of SNV dropout. A. Freyja relative abundance estimates and QualD detection signal on simulated data from GenBank (USA/TX) with 10% of all SNVs retained at random. Freyja is unable to detect any of the four (Alpha, Delta, Gamma, Omicron) VoCs. **B.** Freyja relative abundance estimates and QualD detection signal on simulated data from GenBank (USA/TX) with 25% of all SNVs retained at random. Freyja sparsely detects major VoCs (Delta, Omicron). QualD detections are less sparse for all VoCs. **C.** Freyja relative abundance estimates and QualD detection signal on simulated data from GenBank (USA/TX) with 50% of all SNVs retained at random. **D.** Metadata from GenBank (USA/TX) showing the fraction of genomes belonging to different VoCs for a given week. In this simulated experiment the fractions shown correspond to true relative abundances in the simulated mixture.

Discussion

Wastewater monitoring of the SARS-CoV-2 variant emergence and spread offers unique benefits based on the early detection of the variant arrival prior to the clinical data^{3,34,35,15,25}, and broad surveillance coverage of the population^{20,34,36,37}. QualD offers a highly sensitive pipeline for VoC detection using wastewater SARS-CoV-2 sequencing data. In comparison to one of the leading tools for analyzing SARS-CoV-2 wastewater sequencing data for VoC detection, QualD demonstrated superior sensitivity. This is particularly important given that the underlying sample quality varies and the depth and breadth of coverage of amplicon sequencing data can vary widely across samples¹⁸. Furthermore, the ability to leverage indel information in the inference process makes QualD overall more robust than approaches that rely solely on SNVs^{9,15}. Methods that currently rely on phylogenetic placement or utilize phylogenetic trees inferred by tools like USHER³³ will likely continue to lack support for indel based detection.

While full phylogenetic placement remains the gold standard for strain level analyses of mixed population samples³⁸, challenges posed by the data quality of SARS-CoV-2 wastewater samples limit applicability of such approaches. Thus, QualD is designed as an early detection tool that does not perform full phylogenetic placement of reads. Additionally, since our main goal was high sensitivity to emerging variants in scenarios where the underlying mutational signal is low, QualD treats each observed mutation as an independent event, and hence is not in its current form suited to perform relative abundance estimates.

Since QualD relies on the definition of the quasi-unique mutations for the VoCs, the selection of thresholds for mutation inclusion into the quasi-unique set for a given lineage and for the exclusion of a mutation from this list if found in another lineage affects the sensitivity and specificity of our method. Namely, setting a higher inclusion threshold (i.e. requiring that more of the target lineage genomes contain the mutation) and lower exclusion threshold (i.e. allowing a lower mutational prevalence in other lineages to exclude mutation from the list) will increase specificity, but decrease sensitivity of the method. Our choice of the 50% as the threshold for both inclusion and exclusion is motivated by the prior work for bacterial identification in the field of metagenomics^{39,40}. Furthermore, given that the current genomic database contains more than 10 million SARS-CoV-2 sequences with major variants contributing millions of sequences, setting restrictive thresholds (e.g. requiring that all genomes classified as a given lineage contain the mutation, i.e. setting inclusion threshold to 100%; or alternatively requiring that no genomes outside the target lineage contain the mutation) can result in empty quasi-unique sets for certain lineages.

We envision QualD to be one of several tools routinely employed in wastewater monitoring efforts. For example, QualD could be used in parallel with Freyja to achieve high sensitivity for detecting emerging variants, and relative abundance estimates of the dominant circulating variants. Furthermore, future work on extending the framework of QualD and other tools to other pathogens that can be detected in the wastewater can enable sensitive and continuous environmental monitoring^{41,42} beyond the COVID-19 pandemic. Finally, given the multitude of technical challenges posed by the inherent variability and quality of wastewater sequencing data, we believe that establishing extensive sets of simulated and synthetic datasets that emulate challenges in variant calling in wastewater samples are required to further expand our understanding of how RNA degradation, sample preparation and storage techniques, and sequencing protocols affect the downstream data and analyses.

Methods

Wastewater sample collection, RNA extraction, and sequencing

Houston Water collected and provided weekly 24-hour time-weighted composite influent (raw wastewater) samples from 39 wastewater treatment plants (WWTPs) in Houston covering a service area of approximately 580 miles² and serving over 2.3 million people. In total, 2,637 samples were analyzed. SARS-CoV-2 was concentrated in wastewater samples using an electronegative filtration method as previously described⁴³. We followed the same RNA extraction, library preparation, and sequencing protocols, as described in prior work¹⁸. Details on WWTP sample sites, and methods regarding sample collection procedures, and quantification of SARS-CoV-2 in wastewater samples can also be found in our previous publication¹⁸. Estimates of the viral load are provided by the Houston Health Department following the same methodology as outlined in the SARS-CoV-2 Wastewater Monitoring Dashboard⁴⁴. Raw sequencing reads for the samples used in this study can be found in SRA under BioProject accession: PRJNA796340.

Amplicon sequencing data processing

We processed the MiSeq paired-end data through a standard sequence of steps consisting of quality control report generation (FastQC⁴⁵, default parameters), quality and adapter trimming (BBduk⁴⁶, quality trimming both ends of the read with threshold 15, and

trimming standard PhiX adapter sequences), read mapping (BWA MEM⁴⁷, default parameters), and primer site soft clipping (iVar⁴⁸, ARTIC v3⁴⁹ primer scheme, minimum quality threshold 15). The summary overview of the whole processing pipeline is presented in Figure 3.

Variant calling

We obtained two sets of variant calls for each sample: one with iVar⁴⁸ (minimum quality 20, minimum allele frequency 0) and the other with LoFreq⁵⁰ (after adjusting quality scores for indel calling with the `\lofreq indelqual --dindel`` call, variant calling parameters are set to default). Both variant callers were configured to output all variant calls regardless of the allele frequency. We then used custom Python code to perform a variant call merge-and-filter operation which retained only those variant calls that were supported by both variant callers and had an allele frequency equal or above the user-defined threshold (default: 0.02) according to at least one of the two variant callers (while allele frequency estimates are typically close between the two variant callers differences of <0.01 occur).

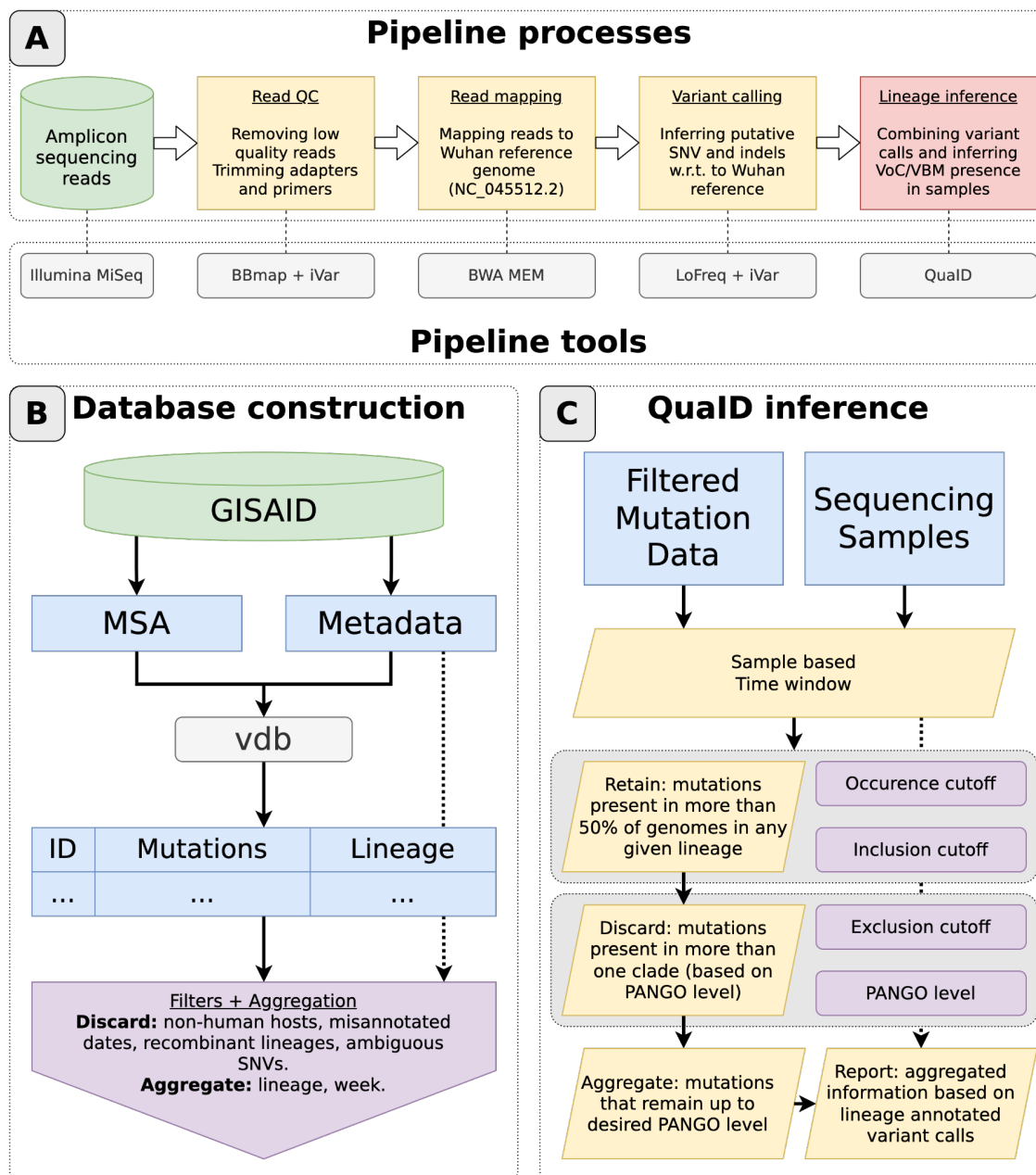


Figure 3. Overview of QualID computational approach. **A.** Overview of the complete sequencing data processing pipeline employed in the analysis of the Houston wastewater SARS-CoV-2 sequencing data. **B.** Schematics overview of the mutation database construction and sanitation used by QualID. MSA and metadata obtained from GISAID are first pre-processed with vdb to extract mutations, and then custom Python code is used to perform additional filtering and aggregation procedures. **C.** QualID VoC/VBM inference process overview. Parameters that affect described subroutines (yellow parallelograms) are provided in the purple rounded rectangles.

Sequence database sanitation

Prior to the subsequent analysis we used metadata obtained from GISAID website to filter out sequences that were marked as incomplete or that had an associated host other than *Homo Sapiens*. Additionally, VoCs with a large amount of clinical sequencing data available (Alpha, Delta, Omicron) are more prone to human error in the metadata entries. Thus, we implemented a filter that removed any genomes: annotated as Alpha with submission date prior to September 3rd, 2020, annotated as Delta with submission date prior to March 1st, 2021, and annotated as Omicron with submission date prior to September 1st, 2021 (first detection dates based on cov-lineages.org VoC reports). Finally, we excluded all recombinant PANGO lineages⁵¹ (X*) from the analysis.

Mutation database construction

We used the pre-generated MSA file from GISAID to extract all mutations using vdb⁵² in nucleotide mode with ambiguous bases included. We then trimmed the resulting list of mutations using the vdb trim command. Finally, we linked the resulting mutation list with the metadata based on the genome accession IDs, and the resulting data were aggregated by week and lineage through custom Python code. Additionally, any SNVs that resulted in an ambiguous base call (e.g. N, W, S, etc.) were removed from the database (summary view provided in Supplementary Figure 2B). The resulting data were used as the mutation tables to calculate prevalence of mutations in PANGO lineages⁵¹ over a user-defined time window (default: 4 weeks).

Quasi-unique mutations

For each lineage and mutation combination, the prevalence of the mutation occurring in the corresponding lineage's genomes was calculated and then converted to a fraction of all genomes assigned to the lineage. Mutations that appeared in more than 50% of all genomes for a single lineage (i.e. not appearing in any other lineage at 50% or more) were considered quasi-unique for that lineage. The above choice of inclusion (what fraction of genomes in the lineage must have the mutations) and exclusion (what fraction of genomes in any other lineage precludes the mutation from being selected) corresponds to the definition of a consensus genome but can be modified to arbitrary values by the end user. Setting stricter thresholds (requiring more of the target lineage genomes to have a mutation) will lead to smaller sets of quasi-unique mutations of high confidence, trading of sensitivity for specificity. Furthermore, since often we want to report detections at a higher level (e.g. any Omicron sub-lineage as opposed to a specific leaf node like BA.2.1) when determining which genomes are used for the exclusion rule, all the genomes that come from the same sub-clade at a fixed level (default: 4) in PANGO hierarchy are omitted from the exclusion check. Thus, mutations common to BA.1 and BA.2 can still be considered as quasi-unique for the Omicron VoC. Note that since vdb reports out deletions and we only filter out ambiguous SNVs, a quasi-unique mutation can be a deletion. Additionally, in order to reduce potential noise from extremely rare lineages, we omit any lineages which have less than a user-defined count of genomes (default: 2) within the designated time window. An overview of these processes is presented in Figure 3C. Finally, for each quasi-unique mutation QualD estimated its predictive power as the posterior probability of observing a particular lineage given the observed mutation. Formally, for a lineage of interest l and the quasi-unique mutation m we computed $P(l|m)$ using Bayes' theorem $P(l|m) = \frac{P(m|l)P(l)}{P(m)}$. We let $P(m)$ to be the ratio of the number of genomes with the mutation m observed to the total number of genomes in the database. Next, we let $P(l)$ to be the ratio of the number of genomes belonging to the lineage l and the total number of genomes. Finally, we let $P(m|l)$ to be the fraction of genomes in the lineage l containing the mutation m . While we did not provide any filtering based on the estimated predictive power of the quasi-unique mutations, these probabilities can be used in the downstream analyses to improve the interpretations of the detection signal provided by QualD.

Mutational signature aggregation

Since the PANGO lineage hierarchy continuously expands potentially introducing new sub-levels for any lineage, it is useful to aggregate quasi-unique mutations into sets that correspond to a node at a fixed level of the hierarchy. For example, Omicron variant is defined as any descendant of B.1.1.529 PANGO lineage, and thus Omicron corresponds to level 4 in the hierarchy. When aggregating quasi-unique mutational signatures up to a given level, we took the union of all descendant lineage quasi-unique mutation sets. Note that the aggregation step always uses the same level of the hierarchy as the exclusion step of the quasi-unique mutation set construction.

Variant of concern detection

Given a wastewater-based sequencing sample collected on a given date D , we constructed the corresponding sets of quasi-unique mutations in the time-window prior to and including weeks up to date D (in case when there is no database information for the week(s) immediately preceding the target date D , the last available time-window was used). Then we merged the filtered set of variant calls for the sample with the quasi-unique set of mutations with the key set to the nucleotide change. We also filtered out any SNVs from the sample that result in synonymous mutations. Once the combined data is obtained, we report for each sample the total combined allele frequency and total count of observed quasi-unique mutations, as well as the total possible number of quasi-unique mutations for the variants of interest at the desired level. Additionally, we reported what percentage of the quasi-unique mutation sites had coverage (with deletions being evaluated based on the genomic positions flanking the deletion) in order to distinguish between the “no detection” and “no coverage” scenarios.

Data availability

Raw sequencing data used in this study is available on SRA under BioProject accession PRJNA796340. Software developed in this manuscript and used to generate the results is available at <https://gitlab.com/treangenlab/quaid>.

Acknowledgements

The authors thank the GISAID contributors who provided the SARS-CoV-2 assemblies. We thank Roger Sealy, Pamela Brown, Ryker Penn, and Yanlai Lai (Houston Health Department) for their assistance in sample collection and sequencing. The authors also would like to thank Lauren Bauhs, Russell Carlson-Stadler, Madeline Wolken, Kyle Palmer, Whitney Rich (Rice University) for their assistance in sample collection, processing, and analysis. This work was supported by the Houston Health Department. Funding sources for sequencing by the Houston Health Department were CDC ELC Enhanced Detection, CDC ELC Enhanced Detection Expansion, and CDC ELC Advanced Molecular Detection. N.S., Y.L. and T.J.T. were supported in part by the C3.ai DTI, Centers for Disease Control (CDC) contract 75D30121C11180 and P01-AI152999 NIH award. E.G.L. and L.B.S. were supported in part by the National Science Foundation (CBET 2029025), and seed funds from Rice University. K.B.E. was supported in part by National Institute of Environmental Health Sciences, R01ES028819.

Author contributions

N.S. has conducted data analyses, designed and implemented software methods, designed and implemented the experiments. Y.L. has tested and implemented software methods. Y.L., E.G.L., and R.S. have conducted data collection and analyses. L.H., K.B.E., L.B.S., T.J.T. have overseen the data collection and analyses, designed and reviewed the experiments. All authors have participated in writing and reviewing the manuscript.

Competing interests

Authors declare no competing interests.

References

1. Herold, M. *et al.* Genome Sequencing of SARS-CoV-2 Allows Monitoring of Variants of Concern through Wastewater. *Water* **13**, 3018 (2021).
2. Fontenele, R. S. *et al.* High-throughput sequencing of SARS-CoV-2 in wastewater provides insights into circulating variants. *Water Research* **205**, 117710 (2021).
3. Randazzo, W. *et al.* SARS-CoV-2 RNA in wastewater anticipated COVID-19 occurrence in a low prevalence area. *Water Research* **181**, 115942 (2020).

4. McClary-Gutierrez, J. S. *et al.* SARS-CoV-2 Wastewater Surveillance for Public Health Action. *EID* **27**, (2021).
5. Polo, D. *et al.* Making waves: Wastewater-based epidemiology for COVID-19 – approaches and challenges for surveillance and prediction. *Water Research* **186**, 116404 (2020).
6. Wu, F. *et al.* SARS-CoV-2 Titers in Wastewater Are Higher than Expected from Clinically Confirmed Cases. *mSystems* **5**, e00614-20 (2020).
7. Peccia, J. *et al.* Measurement of SARS-CoV-2 RNA in wastewater tracks community infection dynamics. *Nat Biotechnol* **38**, 1164–1167 (2020).
8. Kitajima, M. *et al.* SARS-CoV-2 in wastewater: State of the knowledge and research needs. *Science of The Total Environment* **739**, 139076 (2020).
9. Karthikeyan, S. *et al.* Wastewater sequencing uncovers early, cryptic SARS-CoV-2 variant transmission. *Nature* (2022) doi:10.1038/s41586-022-05049-6.
10. Baaijens, J. A. *et al.* Variant abundance estimation for SARS-CoV-2 in wastewater using RNA-Seq quantification. 2021.08.31.21262938 Preprint at <https://doi.org/10.1101/2021.08.31.21262938> (2021).
11. Kirby, A. E. Notes from the Field: Early Evidence of the SARS-CoV-2 B.1.1.529 (Omicron) Variant in Community Wastewater — United States, November–December 2021. *MMWR Morb Mortal Wkly Rep* **71**, (2022).
12. Ellmen, I. *et al.* *Alcov: Estimating Variant of Concern Abundance from SARS-CoV-2 Wastewater Sequencing Data*. 2021.06.03.21258306 <https://www.medrxiv.org/content/10.1101/2021.06.03.21258306v1> (2021) doi:10.1101/2021.06.03.21258306.
13. Crits-Christoph, A. *et al.* Genome Sequencing of Sewage Detects Regionally Prevalent SARS-CoV-2 Variants. *mBio* **12**, e02703-20 (2021).
14. Amman, F. *et al.* Viral variant-resolved wastewater surveillance of SARS-CoV-2 at national scale. *Nat Biotechnol* **1–9** (2022) doi:10.1038/s41587-022-01387-y.
15. Jahn, K. *et al.* Early detection and surveillance of SARS-CoV-2 genomic variants in wastewater using COJAC. *Nat Microbiol* **1–10** (2022) doi:10.1038/s41564-022-01185-x.

16. Smyth, D. S. *et al.* Tracking cryptic SARS-CoV-2 lineages detected in NYC wastewater. *Nat Commun* **13**, 635 (2022).
17. Pérez-Cataluña, A. *et al.* Spatial and temporal distribution of SARS-CoV-2 diversity circulating in wastewater. *Water Research* **211**, 118007 (2022).
18. Lou, E. G. *et al.* Direct comparison of RT-ddPCR and targeted amplicon sequencing for SARS-CoV-2 mutation monitoring in wastewater. *Science of The Total Environment* **833**, 155059 (2022).
19. McCall, C. *et al.* Modeling SARS-CoV-2 RNA degradation in small and large sewersheds. *Environmental Science: Water Research & Technology* **8**, 290–300 (2022).
20. Wu, F. *et al.* SARS-CoV-2 RNA concentrations in wastewater foreshadow dynamics and clinical presentation of new COVID-19 cases. *Science of The Total Environment* **805**, 150121 (2022).
21. Baaijens, J. A., Stougie, L. & Schönhuth, A. Strain-aware assembly of genomes from mixed samples using flow variation graphs. 645721 Preprint at <https://doi.org/10.1101/645721> (2020).
22. Chiara, M. *et al.* Next generation sequencing of SARS-CoV-2 genomes: challenges, applications and opportunities. *Briefings in Bioinformatics* **22**, 616–630 (2021).
23. Izquierdo-Lara, R. *et al.* Monitoring SARS-CoV-2 Circulation and Diversity through Community Wastewater Sequencing, the Netherlands and Belgium - Volume 27, Number 5—May 2021 - Emerging Infectious Diseases journal - CDC. *Emerging Infectious Diseases* **27**, (2021).
24. Graham, K. E. *et al.* SARS-CoV-2 RNA in Wastewater Settled Solids Is Associated with COVID-19 Cases in a Large Urban Sewershed. *Environ. Sci. Technol.* **55**, 488–498 (2021).
25. Lamba, S. *et al.* SARS-CoV-2 infection dynamics and genomic surveillance reveals early variant transmission in urban wastewater. 2022.07.14.22277616 Preprint at <https://doi.org/10.1101/2022.07.14.22277616> (2022).

26. Berry, I. M. *et al.* High confidence identification of intra-host single nucleotide variants for person-to-person influenza transmission tracking in congregate settings. 2021.07.01.450528 Preprint at <https://doi.org/10.1101/2021.07.01.450528> (2021).
27. Nasko, D. J., Koren, S., Phillippy, A. M. & Treangen, T. J. RefSeq database growth influences the accuracy of k-mer-based lowest common ancestor species identification. *Genome Biology* **19**, 165 (2018).
28. Limited genomic reconstruction of SARS-CoV-2 transmission history within local epidemiological clusters | Virus Evolution | Oxford Academic. <https://academic.oup.com/ve/article/8/1/veac008/6523092?login=true>.
29. Quiñones, M. *et al.* “METAGENOTE: a simplified web platform for metadata annotation of genomic samples and streamlined submission to NCBI’s sequence read archive”. *BMC Bioinformatics* **21**, 378 (2020).
30. Schmedes, S. E., King, J. L. & Budowle, B. Correcting Inconsistencies and Errors in Bacterial Genome Metadata Using an Automated Curation Tool in Excel (AutoCurE). *Frontiers in Bioengineering and Biotechnology* **3**, (2015).
31. Elbe, S. & Buckland-Merrett, G. Data, disease and diplomacy: GISAID’s innovative contribution to global health. *Global Challenges* **1**, 33–46 (2017).
32. Gangavarapu, K. *et al.* Outbreak.info genomic reports: scalable and dynamic surveillance of SARS-CoV-2 variants and mutations. 2022.01.27.22269965 Preprint at <https://doi.org/10.1101/2022.01.27.22269965> (2022).
33. 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).
34. Lastra, A. *et al.* SARS-CoV-2 detection in wastewater as an early warning indicator for COVID-19 pandemic. Madrid region case study. *Environmental Research* **203**, 111852 (2022).
35. Sutton, M. *et al.* Detection of SARS-CoV-2 B.1.351 (Beta) Variant through Wastewater Surveillance before Case Detection in a Community, Oregon, USA. *Emerging Infectious Diseases* **28**, (2022).
36. Maida, C. M. *et al.* Wastewater-based epidemiology for early warning of SARS-

- COV-2 circulation: A pilot study conducted in Sicily, Italy. *International Journal of Hygiene and Environmental Health* **242**, 113948 (2022).
37. Larsen, D. A., Green, H., Collins, M. B. & Kmush, B. L. Wastewater monitoring, surveillance and epidemiology: a review of terminology for a common understanding. *FEMS Microbes* **2**, xtab011 (2021).
 38. Czech, L., Stamatakis, A., Dunthorn, M. & Barbera, P. Metagenomic Analysis Using Phylogenetic Placement—A Review of the First Decade. *Frontiers in Bioinformatics* **2**, (2022).
 39. Segata, N. *et al.* Metagenomic microbial community profiling using unique clade-specific marker genes. *Nat Methods* **9**, 811–814 (2012).
 40. Beghini, F. *et al.* Integrating taxonomic, functional, and strain-level profiling of diverse microbial communities with bioBakery 3. *eLife* **10**, e65088 (2021).
 41. Sims, N. & Kasprzyk-Hordern, B. Future perspectives of wastewater-based epidemiology: Monitoring infectious disease spread and resistance to the community level. *Environment International* **139**, 105689 (2020).
 42. Lee, W. L. *et al.* Monitoring human arboviral diseases through wastewater surveillance: Challenges, progress and future opportunities. *Water Research* 118904 (2022) doi:10.1016/j.watres.2022.118904.
 43. LaTurner, Z. W. *et al.* Evaluating recovery, cost, and throughput of different concentration methods for SARS-CoV-2 wastewater-based epidemiology. *Water Research* **197**, 117043 (2021).
 44. City of Houston SARS-CoV-2 Wastewater Monitoring Dashboard.
<https://covidwwtp.spatialstudieslab.org/>.
 45. Andrews, S. FastQC: A Quality Control Tool for High Throughput Sequence Data.
<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/> (2010).
 46. Bushnell, B. BBMap. *SourceForge* <https://sourceforge.net/projects/bbmap/>.
 47. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv:1303.3997 [q-bio]* (2013).
 48. Grubaugh, N. D. *et al.* An amplicon-based sequencing framework for accurately

measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biology* **20**, 8 (2019).

49. Tyson, J. R. *et al.* Improvements to the ARTIC multiplex PCR method for SARS-CoV-2 genome sequencing using nanopore. 2020.09.04.283077 Preprint at <https://doi.org/10.1101/2020.09.04.283077> (2020).
50. Wilm, A. *et al.* LoFreq: a sequence-quality aware, ultra-sensitive variant caller for uncovering cell-population heterogeneity from high-throughput sequencing datasets. *Nucleic Acids Res* **40**, 11189–11201 (2012).
51. Rambaut, A. *et al.* A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol* **5**, 1403–1407 (2020).
52. West, A. P. *et al.* Detection and characterization of the SARS-CoV-2 lineage B.1.526 in New York. *bioRxiv* 2021.02.14.431043 (2021) doi:10.1101/2021.02.14.431043.