# Bacterial metatranscriptomes in wastewater can differentiate virally infected human populations

## (9/10 words)

Rodolfo A Salido[1,2,4], Cameron Martino[1,3,4], Smruthi Karthikeyan[1], Shi Huang[1,4], Gibraan Rahman[1,5], Antonio Gonzalez[1], Livia S. Zaramela[1], Kristen L Beck[6], Shrikant Bhute[4], Kalen Cantrell[4,5], Anna Paola Carrieri[7], Sawyer Farmer[1], Niina Haiminen[8], Greg Humphrey[1], Ho-Cheol Kim[6], Laxmi Parida[8], Alex Richter[4], Yoshiki Vázquez-Baeza[4], Karsten Zengler[1,2], Austin D. Swafford[4], Andrew Bartko[4], Rob Knight[1,2,4,5]

1. Department of Pediatrics, School of Medicine, University of California San Diego, La Jolla, CA, USA.
2. Department of Bioengineering, University of California San Diego, La Jolla, CA, USA.
3. Bioinformatics and Systems Biology Program, Jacobs School of Engineering, University of California San Diego, La Jolla, CA, USA.
4. Center for Microbiome Innovation, Jacobs School of Engineering, University of California San Diego, La Jolla, CA, USA.
5. Department of Computer Science and Engineering, Jacobs School of Engineering, University of California San Diego, La Jolla, CA, USA.
6. AI and Cognitive Software, IBM Research-Almaden, San Jose, CA, USA.
7. IBM Research Europe - Daresbury, UK.
8. IBM T.J Watson Research Center, Yorktown Heights, New York, USA.

Corresponding author: robknight@eng.ucsd.edu

**Abstract:**

Monitoring wastewater samples at building-level resolution screens large populations for SARS-CoV-2, prioritizing testing and isolation efforts. Here we perform untargeted metatranscriptomics on virally-enriched wastewater samples from 10 locations on the UC San Diego campus, demonstrating that resulting bacterial taxonomic and functional profiles discriminate SARS-CoV-2 status even without direct detection of viral transcripts. Our proof-of-principle reveals emergent threats through changes in the human microbiome, suggesting new approaches for untargeted wastewater-based epidemiology.

**Keywords:**

37 **Body:**

38 Our past work deploying a highly spatially resolved, high-throughput wastewater monitoring
39 system on a college campus (1) enabled collection and qPCR characterization of thousands of
40 wastewater samples, identifying 85% of SARS-CoV-2 clinical cases (2), and also enabling
41 genomic surveillance for emerging variants of concern by complete genome sequencing from
42 extracted RNA (3). Wastewater-based epidemiology (WBE) provides additional advantages in
43 that it is (i) non-invasive, (ii) cost-effective relative to individual clinical testing, (iii) does not
44 require individuals to consent to clinical testing that is often reported to public health agencies,
45 and (iv) can therefore benefit under-served populations (4-6). However, this WBE scheme is
46 currently limited to pathogen detection and characterization through targeted qPCR and
47 sequencing, and cannot detect agents of disease for which a screening test has not been
48 developed.

49

50 Here we describe an untargeted community/population level disease monitoring strategy using
51 metatranscriptomics, which leverages correlations in observable changes in wastewater
52 microbiomes with human microbiome disruptions associated with disease state. SARS-CoV-2,
53 like many pathogens, has been reported to cause systematic disruptions in the human gut
54 microbiome (7-9), which is the principal human microbial input to wastewater (10). We
55 employed this strategy to test whether information in the wastewater metatranscriptome could
56 discriminate SARS-CoV-2 positive from negative wastewater samples (assessed by qPCR) as a
57 proof-of-principle.

58

59 We present a high-throughput wastewater metatranscriptomics pipeline that lowers the
60 accessibility to an otherwise cost-prohibitive sequencing method at scale through
61 miniaturization, parallelization, and automation (11-12). (**Sup. Fig. S1**) Using this pipeline, we
62 generated metatranscriptomics sequencing data for 313 virally-enriched (VE) wastewater
63 samples collected from manholes servicing different residential buildings across a college
64 campus, including isolation housing buildings (Manhole IDs: C6M095-C6M098), from Nov 23
65 2020 to January 7 2021. Sequencing reads were demultiplexed, trimmed, and quality filtered
66 before being deposited in Qiita (13), where ribosomal reads were removed using SortMeRNA
67 (14) using default processing recommendations; non-ribosomal reads were aligned to genomes
68 or genes using Woltka (15) resulting in two different feature tables: taxonomic and functional
69 (details in Materials and Methods).

70

71 Samples obtained from each manhole have a distinct microbiome signature, likely a composite
72 of the individual microbiomes of the people contributing to each wastewater stream. Beta-
73 diversity analyses of both metatranscriptomic feature tables (taxonomic and functional)
74 measured by Aitchison distance and robust Aitchison principal component analysis (RPCA) (16)
75 reveal that wastewater samples cluster primarily by manhole source (manhole_id) (**Fig. 1A**),
76 with a stronger signal than SARS-CoV-2 detection status (**Fig. 1B**)(**Sup. Table ST1**).
77 Wastewater samples separate according to SARS-CoV-2 status based on these bacterial
78 profiles alone, but this signal is obscured in the RPCA ordination by the stronger manhole_id
79 clustering effect. Taxonomic features provide better separation by both SARS-CoV-2 status and
80 manhole_id than functional features (**Supp. Table ST1**), suggesting that microbial community

81  membership rather than current functional gene expression is more strongly affected by
82  infection.
83
84  To test whether the SARS-CoV-2 detection status-dependent microbiome signal can be
85  identified even against the stronger manhole_id clustering effect, we selected a subset of
86  samples for paired comparisons between SARS-CoV-2 positive and negative samples within
87  specific manholes across one week (selection process detailed in Materials and Methods). This
88  subset (squares, n=28 **Fig. 1A-B**) was analyzed by dimensionality reduction with compositional
89  tensor factorization (CTF) (17), which accounts for the intra-manhole sample correlation. The
90  resulting ordination shows that samples of the microbiome in any specific manhole undergo a
91  pronounced shift along one of the main principal components (PC1 for taxonomic, PC2 for
92  functional), when the subject population it services becomes infected with SARS-CoV-2 (**Fig.**
93  **1C-D**). Consequently, taxonomic features (genomes) that drive segregation along PC2 (**Fig.**
94  **1E**), or functional features (genes) along PC1 (**Fig. 1F**), can be positively or negatively
95  correlated with SARS-CoV-2 detection. Log-ratio analysis of the top and bottom ranked
96  taxonomic features as numerator and denominator respectively show a significant difference in
97  the means of the SARS-CoV-2 detection sample groupings (**Fig. 1G**). Similarly, a log-ratio of six
98  functional features positively and negatively ranked along PC2 also shows a significant
99  difference in the means of the SARS-CoV-2 detection sample groupings (**Fig. 1H**) (see
100 Materials and Methods).
101
102 The predictive power for wastewater SARS-CoV-2 status discrimination of the features selected
103 through CTF analysis was validated via log-ratios and random forest machine learning (RFML)
104 classification, using the remaining samples in this study (circles, **Fig. 1A-B**) plus an additional
105 validation set (total n=285, positive=179, negative=106, **Sup. Table ST2**). Log-ratios of selected
106 taxonomic and functional features showed a significant difference by SARS-CoV-2 detection
107 status across the validation sampleset, with function (*t*-test, $T$=-3.9 p=0.0001) (**Fig. 2A**) showing
108 a smaller effect than taxonomy (*t*-test, T=-8.8, p=1.3e-16) (**Fig. 2B**). Type II ANOVA of both log-
109 ratios shows that differences in sample means are larger across SARS-CoV-2 status groups
110 than manhole_id or sample_plate confounders (**Sup. Fig. S2**). The performances of the RFML
111 classification models were evaluated through average area under the curve of precision-recall
112 (AUC-PR) tests of stratified 5-fold cross validation classification tasks distinguishing samples'
113 SARS-CoV-2 status, manhole_id, and sample_plate. Lower dimensional feature tables from
114 feature selection show comparable SARS-CoV-2 status classification performance as full
115 feature tables for both data modalities (taxonomic and functional) (**Fig. 2C**), but reduced
116 classification performance when distinguishing confounding manhole_id (**Fig. 2D**) or
117 sample_plate (**Sup. Fig. S3**).
118
119 Our results demonstrate that wastewater metatranscriptomes can reveal traces of rare
120 pathogens through alterations of the microbiome of the afflicted individuals, which are eventually
121 reflected in the wastewater microbiome. When effects are confounded by site/population,
122 leveraging generalizable log-ratios separating positive/negative groupings across sites reduces
123 overfitting. This proof-of-principle justifies further research on high-throughput wastewater
124 metatranscriptome biomarker discovery for WBE; the untargeted nature of this data modality

125 makes it flexible enough to monitor multiple diseases at the population scale (through traditional
126 direct detection of known sequences from pathogens, but also by leveraging microbiome
127 perturbations as a proxy), and is superior to metagenomic monitoring because it encompasses
128 all living organisms and viruses(18). One of the limitations of the proposed strategy is the
129 narrow stability of the samples' RNA molecules. However, our methods don't claim to
130 comprehensively characterize the wastewater metatranscriptome and instead focus on the fact
131 that changes in the observable bacterial metatranscriptome are sufficient to discriminate the
132 wastewater's viral status, with SARS-CoV-2 detection status serving as a relevant case study.
133 Although key features of the bacterial metatranscriptome discriminate SARS-CoV-2 detection,
134 further work is needed to determine how broadly this phenomenon generalizes to other
135 pathogens. Lastly, our methodology allows automated high-throughput metatranscriptomics
136 processing, applicable to many biospecimen types, and could have considerable impact beyond
137 WBE.

138

**139 Acknowledgments**

144

**145 Conflict of Interest:**

146 A.D.S. is currently Chief Technology Officer of InterOme, Inc. a digital health company which
147 offers wastewater testing and monitoring of pathogens including SARS-CoV-2 among its
148 services

149

150    **References:**

151    1. Karthikeyan, S. *et al.* High-Throughput Wastewater SARS-CoV-2 Detection Enables
152       Forecasting of Community Infection Dynamics in San Diego County. *mSystems* **6**,
153       (2021).

154    2. Karthikeyan, S. *et al.* Rapid, Large-Scale Wastewater Surveillance and Automated
155       Reporting System Enable Early Detection of Nearly 85% of COVID-19 Cases on a
156       University Campus. *mSystems* **6**, 793–814 (2021).

157    3. Karthikeyan, S. *et al.* Wastewater sequencing uncovers early, cryptic SARS-CoV-2
158       variant transmission. *medRxiv* 2021.12.21.21268143 (2021).
159       doi:10.1101/2021.12.21.21268143

160    4. Fielding-Miller, R. *et al.* Wastewater and surface monitoring to detect COVID-19 in
161       elementary school settings: The Safer at School Early Alert project. *medRxiv*
162       2021.10.19.21265226 (2021). doi:10.1101/2021.10.19.21265226

163    5. Reitsma, M. B. *et al.* Racial/ethnic disparities in covid-19 exposure risk, testing, and
164       cases at the subcounty level in California. *Health Aff.* **40**, 870–878 (2021).

165    6. Lieberman-Cribbin, W., Tuminello, S., Flores, R. M. & Taioli, E. Disparities in COVID-
166       19 Testing and Positivity in New York City. *Am. J. Prev. Med.* **59**, 326–332 (2020).

167    7. Wu, Y. *et al.* Altered oral and gut microbiota and its association with SARS-CoV-2 viral
168       load in COVID-19 patients during hospitalization. *npj Biofilms Microbiomes* **7**, 61
169       (2021).

170    8. Xu, R. *et al.* Temporal association between human upper respiratory and gut bacterial
171       microbiomes during the course of COVID-19 in adults. *Commun. Biol. 2021 41* **4**,
172       1–11 (2021).

173    9. Gu, S. *et al.* Alterations of the Gut Microbiota in Patients With Coronavirus Disease
174       2019 or H1N1 Influenza. *Clin. Infect. Dis.* **71**, 2669–2678 (2020).

175    10. Newton, R. J. *et al.* Sewage reflects the microbiomes of human populations. *MBio* **6**,
176        (2015).

177    11. Sboner, A., Mu, X. J., Greenbaum, D., Auerbach, R. K. & Gerstein, M. B. The real
178        cost of sequencing: Higher than you think! *Genome Biol.* **12**, 1–10 (2011).

179    12. Mayday, M. Y., Khan, L. M., Chow, E. D., Zinter, M. S. & DeRisi, J. L. Miniaturization
180        and optimization of 384-well compatible RNA sequencing library preparation. *PLoS*
181        *One* 14, e0206194 (2019).

182  13. Gonzalez, A. *et al.* Qiita: rapid, web-enabled microbiome meta-analysis. *Nat. Methods*
183      **15**, 796–798 (2018).

184  14. Kopylova, E., Noé, L. & Touzet, H. SortMeRNA: fast and accurate filtering of
185      ribosomal RNAs in metatranscriptomic data. *Bioinformatics* **28**, 3211–3217 (2012).

186  15. Zhu, Q. *et al.* OGUs enable effective, phylogeny-aware analysis of even shallow
187      metagenome community structures. *bioRxiv* 2021.04.04.438427 (2021).
188      doi:10.1101/2021.04.04.438427

189  16. Martino, C. *et al.* A Novel Sparse Compositional Technique Reveals Microbial
190      Perturbations. *mSystems* **4**, e00016-19 (2019).

191  17. Martino, C. *et al.* Context-aware dimensionality reduction deconvolutes gut microbial
192      community dynamics. *Nat. Biotechnol. 2020 392* **39**, 165–168 (2020).

193  18. Sims, N. & Kasprzyk-Hordern, B. Future perspectives of wastewater-based
194      epidemiology: Monitoring infectious disease spread and resistance to the
195      community level. *Environ. Int.* **139**, 105689 (2020).

196

197
198

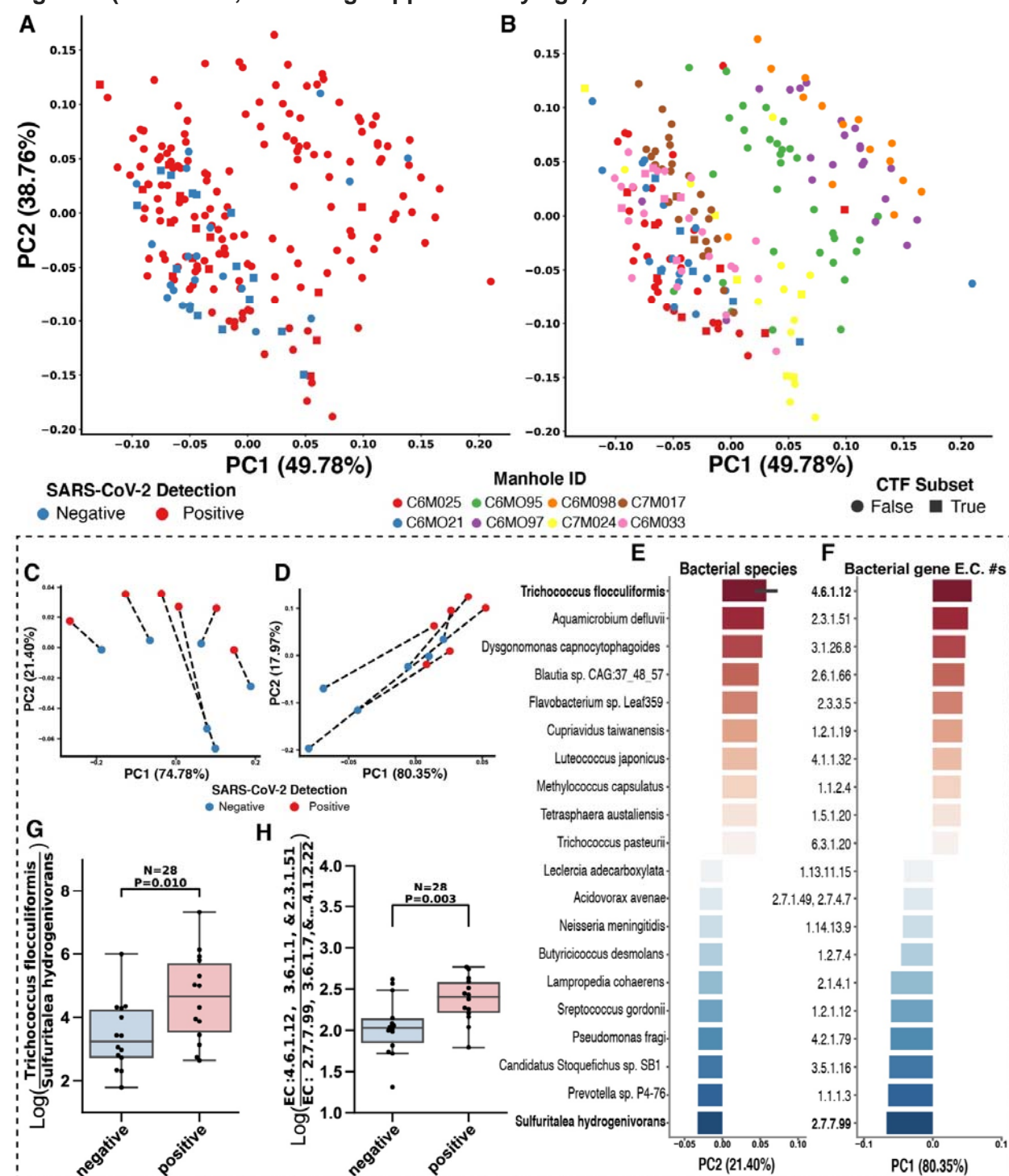Figures: (223 words, excluding supplementary figs)



**Figure 1: Microbial community composition changes can be observed in SARS-CoV-2 positive vs. negative wastewater samples**. Robust principal component analysis (RPCA) of wastewater samples colored by SARS-CoV-2 detection status (**A**) and manhole source (**B**). A subset of samples (squares) was selected for pairwise comparisons of SARS-CoV-2 positive and negative wastewater microbiomes within a manhole and a week using compositional tensor

factorization (CTF) on taxonomic (genomes, **C**) and functional (genes, **D)** features. Results shown in the dashed box are exclusive to this subset of samples. Important bacterial genomes (**E**) and genes (**F**) identified from CTF show significant differences between positive and negative sample groupings by log-ratios of top and bottom ranked features respectively (**G-H**). Error bar on the x-axis of the ranked features plot represents the standard error in the PC2 loadings across strains within the same species. The log-ratio boxplot elements are defined as follows: the centerline is the median of the distribution, box limits represent upper and lower quartiles, whiskers span 1.5x of the interquartile range, and points represent all data points.
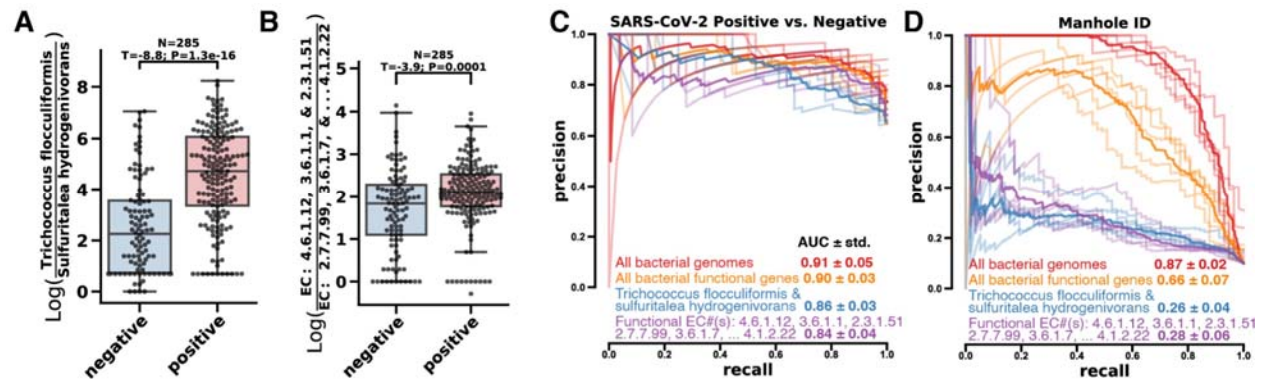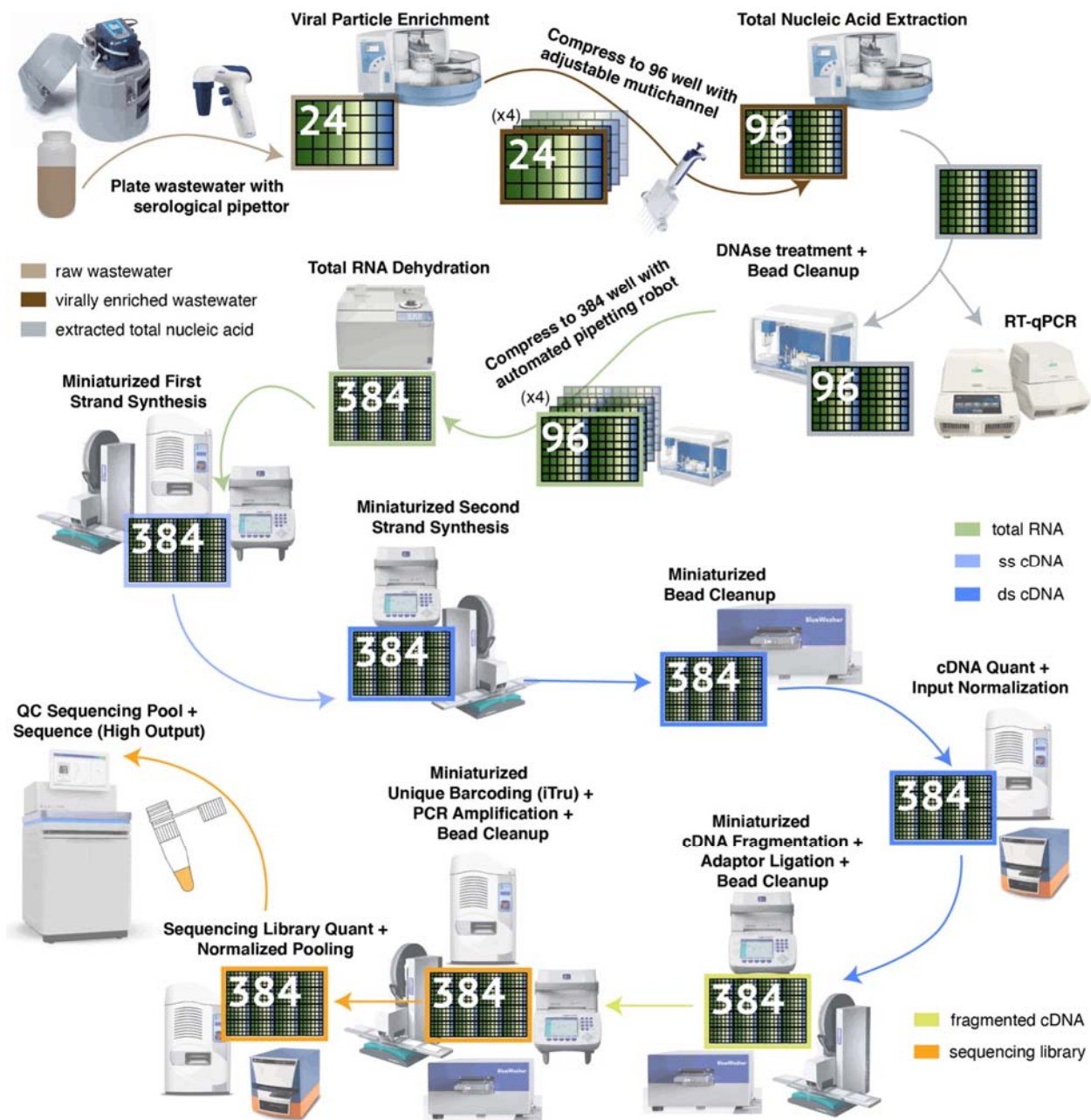


**Figure 2: Key bacterial features identified in small paired subset show significant differences in larger validation dataset and provide RFML the ability to accurately predict SARS-CoV-2 status but not manhole source in wastewater.** Log-ratios of important features, taxonomic (**A**) and functional (**B**), identified by CTF significantly separate wastewater samples by SARS-CoV-2 detection status in the remaining samples not included in the CTF subset. The log-ratio boxplot elements are defined as follows: the centerline is the median of the distribution, box limits represent upper and lower quartiles, whiskers span 1.5x of the interquartile range, and points represent all data points. **C**) Random forest machine learning 5-fold cross-validation shows high precision-recall of samples with positive SARS-CoV-2 detection status from taxonomic and functional tables with all features or a few selected features. **D**) Feature selection reduces Manhole ID classification performance while retaining SARS-CoV-2 discrimination, suggesting a reduction of overfitting. The translucent precision-recall curve traces of each feature table reflect all 5-fold cross-validation results while the bold trace represents the average.

**Supplementary Figure S1: High Throughput pipeline for Virally Enriched (VE) wastewater metatranscriptomics.** Flow diagram of metatranscriptomic data generation from VE wastewater samples, from auto-sampler to sequencer. Key robotic instrumentation and tools are depicted alongside each step. The flow diagram is color coded according to the different stages of sample processing. The high throughput pipeline increases sample processing parallelization through incremental compression of samples from 24-well plates to 384-well plates. Significant per sample cost savings are achieved through miniaturization of molecular reactions in 384-well format, for which specialized low volume liquid handling infrastructure is needed.

242
243
244
245
246
247

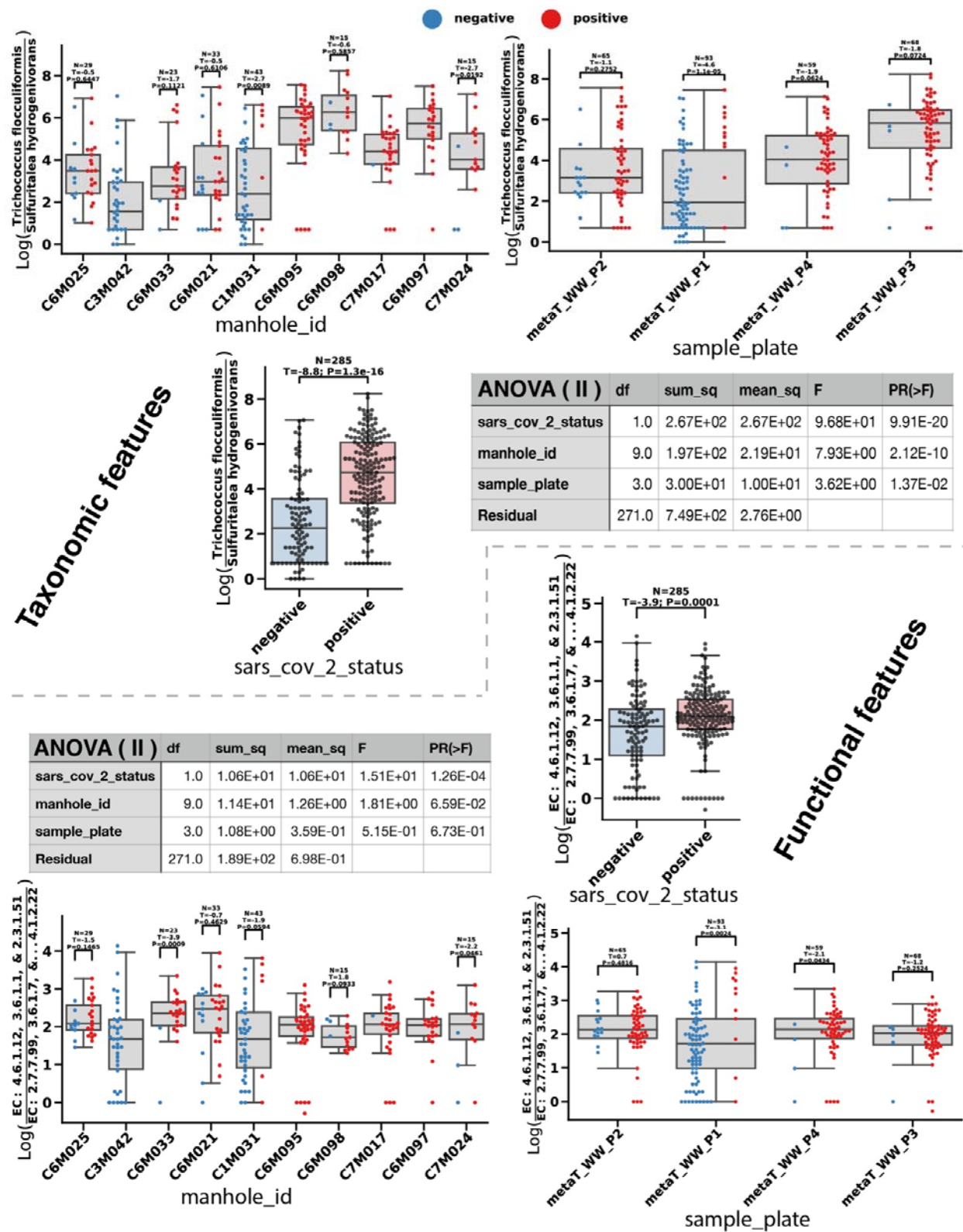|  |  |  | PERMANOVA: F-stat. | PERMANOVA: p-value |
|---|---|---|---|---|
| **across-all** | **taxonomic** | manhole_id | 23.9008 | 0.0002 |
|  |  | time_encoded | 0.8869 | 0.7321 |
|  |  | sars_cov_2_status | 8.8129 | 0.0002 |
|  |  | sample_plate | 21.2303 | 0.0002 |
|  | **functional** | manhole_id | 11.9542 | 0.0002 |
|  |  | time_encoded | 0.9860 | 0.5055 |
|  |  | sars_cov_2_status | 4.0365 | 0.0180 |
|  |  | sample_plate | 9.1532 | 0.0002 |

248
249 **Supplementary Table ST1: PERMANOVA results on RPCA distance matrix show stronger**
250 **manhole of origin effect than SARS-CoV-2 status.** An analysis of variance of the Aitchison
251 distance between wastewater samples shows that manhole of origin has the strongest effect
252 size, followed by sample processing plate, and SARS-CoV-2 status. Samples from different
253 manholes were not uniformly distributed across sample processing plates, confounding the
254 effect sizes for both independent variables.

255

256

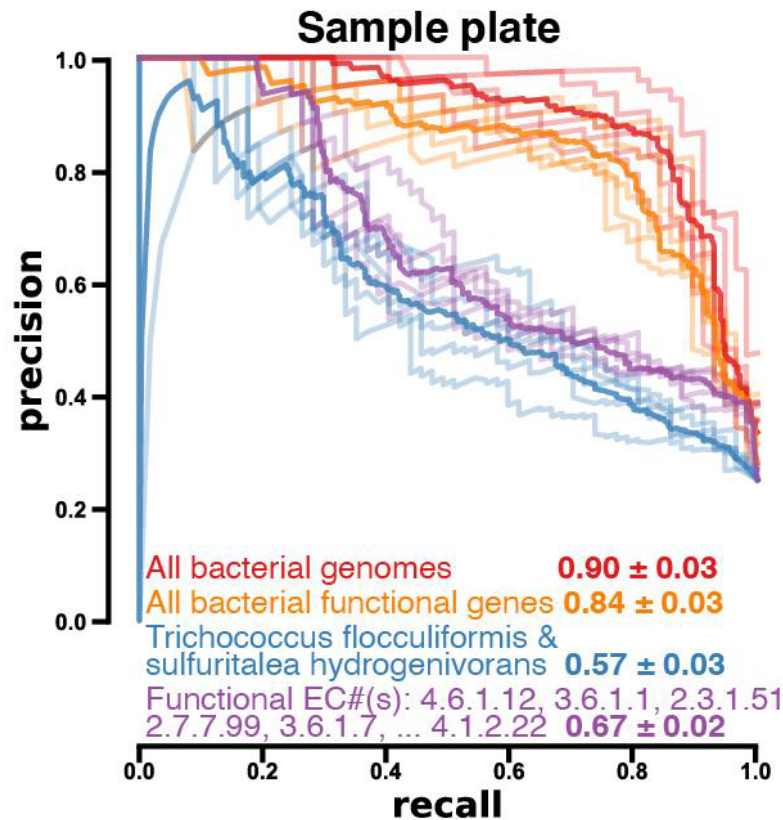| sample_plate | manhole_id | sars_cov_2_status | samples |
|---|---|---|---|
| | C1Mo31 | negative | 37 |
| | | positive | 6 |
| Sample_Plate_1 | C3Mo42 | negative | 38 |
| | C6Mo21 | negative | 7 |
| | | positive | 5 |
| | C6Mo21 | negative | 5 |
| | | positive | 16 |
| Sample_Plate_2 | C6Mo25 | negative | 10 |
| | | positive | 19 |
| | C6Mo95 | positive | 15 |
| | C6Mo33 | negative | 2 |
| | | positive | 7 |
| | C6Mo95 | positive | 22 |
| Sample_Plate_3 | C6Mo97 | positive | 22 |
| | C6Mo98 | negative | 3 |
| | | positive | 12 |
| | C6Mo33 | positive | 14 |
| | C7Mo17 | negative | 1 |
| Sample_Plate_4 | | positive | 29 |
| | C7Mo24 | negative | 3 |
| | | positive | 12 |
| | TOTAL | | 285 |

257

258 **Supplemental Table ST2: Description of validation dataset for Random Forest Machine**
259 **Learning (RFML).** Distribution of samples across different groupings relevant to the observed
260 variance in the unsupervised learning analysis. Sample plate 1 was added, as an additional
261 validation set, to the RFML analyses. The validation dataset excludes the subset of samples
262 selected for the CTF analysis (n=28).

263

**Supplementary Figure S2: Analysis of variance (ANOVA) of both log-ratios show that SARS-CoV-2 status has the strongest effect size.** Boxplots with overlaid swarmplots show the distribution of selected log-ratios for both taxonomic and functional feature tables, grouped

268    by relevant sample metadata. The log-ratio boxplot elements are defined as follows: the
269    centerline is the median of the distribution, box limits represent upper and lower quartiles,
270    whiskers span 1.5x of the interquartile range, and points represent all data points. Results from
271    ANOVA (type II) analyses are shown as tables for each feature modality. Statistical tests results
272    (Student's $t$-test) between SARS-CoV-2 status subgroupings (negative=blue / positive=red) in
273    manhole_id and sample_plate plots are also shown, evidencing that the log-ratios generalize
274    and perform better at discriminating SARS-CoV-2 status across all samples than within specific
275    manholes.
276

**Sample plate**

277

**Supplementary Figure S3:** Random forest machine learning 5-fold cross-validation shows a decrease in precision-recall of samples' processing plate (sample plate) from feature selection of taxonomic and functional feature tables in comparison to full feature tables, suggesting a reduction of overfitting on a possible technical confounder. The translucent precision-recall curve traces of each feature table reflect all 5-fold cross-validation results while the bold trace represents the average.

284