



**Environmental
Science**
Water Research & Technology

Prewhitening and normalization help detect a strong cross-correlation between daily wastewater SARS-CoV-2 RNA abundance and COVID-19 cases in a community

Journal:	<i>Environmental Science: Water Research & Technology</i>
Manuscript ID	EW-ART-12-2022-000951.R1
Article Type:	Paper

SCHOLARONE™
Manuscripts

Prewhitening and Normalization Help Detect a Strong Cross-Correlation Between Daily Wastewater SARS-CoV-2 RNA Abundance and COVID-19 Cases in a Community

Min Ki Jeon ¹, Bo Li ¹, Doris Yoong Wen Di ¹, and Tao Yan ^{1*}

¹ Department of Civil and Environmental Engineering, University of Hawaii at Manoa, Honolulu, HI 96822

*Corresponding author: Tao Yan, University of Hawaii at Manoa, Department of Civil and Environmental Engineering, 2540 Dole Street, 383 Holmes Hall, Honolulu, HI 96822. Phone: 808-956-6024. Fax: 808-956-5014. E-mail: taoyan@hawaii.edu

Abstract

Wastewater surveillance is a promising technology for real-time tracking and even early detection of COVID-19 infections in communities. Although correlation analysis between wastewater surveillance data and the daily clinical COVID-19 case numbers has been frequently conducted, the importance of stationarity of the time-series data has not been well addressed. In this study, we demonstrated that strong yet spurious correlation could arise from non-stationary time-series data in wastewater surveillance, and data prewhitening to remove trends by the first differences of y-values between two consecutive times helped to reveal distinct cross-correlation patterns between daily clinical case numbers and daily wastewater SARS-CoV-2 concentration during a lockdown period in 2020 in Honolulu, Hawaii. Normalization of wastewater SARS-CoV-2 concentration by the endogenous fecal viral markers in the same samples significantly improved the cross-correlation, and the best correlation was detected at a two-day lag of the daily clinical case numbers. The detection of a significant correlation between daily wastewater SARS-CoV-2

RNA abundance and clinical case numbers also suggests that disease burden fluctuation in the community should not be excluded as a contributor to the often observed weekly cyclic patterns of clinical cases.

Water impact

Wastewater surveillance represents an emerging water technology with significant human health benefits. The study demonstrated that non-stationary time-series data could lead to spurious correlation, highlighting the need for prewhitening. Normalization strategies could alleviate variations in sample collection and analyses, which is useful for detecting actual underlying relationships between wastewater surveillance data and clinical data.

1. Introduction

Since the outbreak of COVID-19 pandemic in late 2019 (1), wastewater surveillance has been explored as a new way to monitor the spread of SARS-CoV-2 in human communities. Many studies have shown the presence of SARS-CoV-2 viral particles or genomic RNA in bodily wastes, including feces (2), urine (3), and respiratory fluids (4, 5) in both symptomatic and asymptomatic patients. In particular, asymptomatic infections are now known to account for a large percentage of total COVID-19 infections (6, 7), and also shed SARS-CoV-2 virus in feces (8, 9). Since wastewater collects human wastes from all individuals in the wastewater service area and hence can provide comprehensive information on COVID-19 infection in the community. This enables a unique advantage of wastewater surveillance in that it can potentially capture the “actual” infection rates, including the asymptomatic or mildly symptomatic patients in the community who are less likely to seek clinical testing.

Wastewater surveillance may also be able to provide real-time tracking and even early detection of infectious in the communities. The sources of SARS-CoV-2 viral shedding to wastewater include mainly feces and partially saliva and sputum due to their high shedding probability and the possibility of entering the sewers (10). It is known that COVID-19-infected patients start to shed SARS-CoV-2 virus in feces during the incubation period between the infection and the symptom onset (11) and the peak of SARS-CoV-2 viral concentration is reported generally at the beginning of the symptom onset (12, 13). In addition, model fitting results from a meta-data analysis study using experimental findings from various clinical studies have estimated the highest viral concentration at 0.34 days after symptom onset (14).

Since data collected from both clinically confirmed COVID-19 cases and SARS-CoV-2 RNA abundance in wastewater are time series data, their relationships could be examined through time-series data analyses in particular cross-correlation. Removing the trend or seasonality of time-series data sets to achieve stationarity is an important prerequisite process to avoid spurious correlation (15). This transformation process is called “prewhitening” and it transforms time series data into stationary forms. One common prewhitening method is the single differencing of the time series data points by their first differences (16). However, many of the wastewater surveillance studies that correlated SARS-CoV-2 in wastewater and COVID-19 cases in the community did not prewhiten the time series data (17-20). As a result, the strong correlation coefficients observed could be attributed to trend or seasonality instead of the actual correlation of variation between the two types of data sets.

In this study, we re-analyzed previously collected time-series data on daily wastewater SARS-CoV-2 RNA abundance and clinical COVID-19 cases in a large metropolitan area to demonstrate the importance of prewhitening when conducting cross-correlation analysis. Both SARS-CoV-2 RNA concentration and its normalized abundance were subjected to time-series cross-correlation analysis to determine any presence of lags between the corresponding daily

clinical case numbers observed in the community. Also, normalization strategies were compared to distinguish the best improvement of correlation with the daily clinical case numbers.

2. Materials and Methods

2.1. Wastewater sampling, viral precipitation, RT process, and qPCR assays

Wastewater sampling, processing, and RT-qPCR quantification of SARS-CoV-2 RNA and several RNA viruses in wastewater were previously described in detail in Li *et al.* (21), which are briefly summarized in the following (Table S1 and S2). The two largest wastewater treatment plants (WWTP), Sand Island (SI) and Honouliuli (HO) in the City and County of Honolulu, were selected to collect wastewater samples to represent the wastewater of the community. Untreated primary influent wastewater samples were collected by daily flow-adjusted composite sampling from the SI and HO WWTPs from August 27th, 2020 to October 4th, 2020 (i.e. Day 0 to 38, each WWTP with $n = 39$). All daily collected wastewater samples were thoroughly mixed and subsequently centrifuged to collect suspended solids and supernatant, which were referred to as solid and liquid fractions of the wastewater samples, respectively. The exogenous process control bovine coronavirus (BCoV) (Zoetis; Kalamazoo, MI, USA) was spiked into the solid and liquid subsamples to detect inhibition and assess recovery (four batches of samples, $n = 46$) (Table S2). The liquid fractions were first treated by the polyethylene glycol (PEG) precipitation method (22) to concentrate and pellet viral particles in the liquid fractions. The precipitated pellets from the liquid fractions as well as the solid fractions were subjected to viral RNA extraction. All extracted viral RNA samples were reverse transcribed with random hexamer N6, and the produced cDNA samples were analyzed by qPCR assays targeting SARS-CoV-2 (the N1 and N2 assays (23) and the E gene assay (24)) and fecal RNA viral surrogates (F+ RNA coliphages Group II (G2) and Group III (G3) (25), and pepper mild mottle virus (PMMoV) (26)).

2.2. Data analysis

All data analyses used both the log-transformed SARS-CoV-2 concentration data determined by the three qPCR assays (i.e. log N1, log N2, and log E) and its abundances normalized by the three fecal viral indicators (i.e. log (N1/G2), log (N2/G2), log (E/G2), log (N1/G3), log (N2/G3), log (E/G3), log (N1/PMMoV), log (N2/PMMoV) and log (E/PMMoV)). Daily new COVID-19 case counts for Honolulu were sourced from the local COVID-19 dashboard by the Hawaii Emergency Management Agency. The study only used publicly available data at the population level, and thus required no IRB review. Cross-correlation was used to examine the time-lagged association between new clinical COVID-19 cases in the community and wastewater SARS-CoV-2 abundance. A prewhitening process was applied to all time-series data, including COVID-19 clinical case numbers and the wastewater SARS-CoV-2 wastewater abundance, to remove trends. The wastewater SARS-CoV-2 abundance data were prewhitened by log transformation followed by the first differences, and the clinical case data were prewhitened by the first differences. The original and the prewhitened data were tested for normality by using Shapiro-Wilk test (27). Mann-Kendall test (28, 29) was used for the assessment of trend significance before and after the prewhitening to verify the successful removal of trends.

Cross-correlation of original and prewhitened SARS-CoV-2 concentration and their normalized abundance in liquid or solid fractions and the daily new clinical COVID-19 cases were analyzed for the SI and HO WWTPs separately. The Cross-Correlation Functions (CCF) function in the R environment was used with a maximum lag of six days. A positive lag indicates that the SARS-CoV-2 concentration or normalized abundance was leading the clinical cases. Positive coefficients indicate a positive relationship between the SARS-CoV-2 concentration or normalized abundance and clinical cases.

Results of the cross-correlation analysis were visualized by heatmaps and boxplots. Additionally, correlation coefficients from the cross-correlation analysis were compared depending on normalization strategies by using p -values obtained from the pairwise t -test and were adjusted by the Benjamini and Hochberg correction (30) to determine which normalization strategy showed the best improvement. All statistical analyses and data visualization were conducted in R 4.2.1 (31) by using the packages *tidyverse* 1.3.1 (32), *ggpubr* 0.4.0 (33), *scales* 1.1.1 (34), and *rstatix* 0.7.0 (35).

3. Results

3.1. Importance of prewhitening on cross-correlation

Our previous study (21) detected downward trends for both daily clinical COVID-19 case numbers in the community and SARS-CoV-2 RNA abundances (both with or without normalization by fecal RNA viral markers) in the wastewater samples. The observed downward trends were the result of a public health lockdown implemented to counter the COVID-19 outbreak on the Island of Oahu. Therefore, the time-series data of raw wastewater SARS-CoV-2 RNA concentration and its normalized relative abundance, as well as the daily fluctuation of clinical case numbers all need to be prewhitened in order to be de-trended before cross-correlation analysis. The Shapiro-Wilk test showed that only 25 out of 48 of the original data ($p = 0.08 \pm 0.10$) were normally distributed, but all of the prewhitened data ($p = 0.39 \pm 0.24$) were normally distributed (Table S3). The Mann-Kendall test results confirmed the removal of trends from both SARS-CoV-2 abundance data from the wastewater and the daily clinical case data after the prewhitening (Table S3).

The prewhitened SI and HO WWTPs time-series SARS-CoV-2 concentration data were first compared with the prewhitened daily clinical case numbers, which showed only weak correlations (either positive or negative) (Figure 1). The only significant correlation was

observed at a two-day lag of the prewhitened daily clinical case numbers (x_{t+2}) behind the prewhitened daily wastewater SARS-CoV-2 concentrations of log N1 ($r = 0.38$, $p = 0.019$) in the liquid fraction of HO WWTP (Figure 1F). No significant correlations were found from any other lags from the two WWTPs with the prewhitened daily clinical case numbers, as indicated by the boxplots falling under the 95% confidence level (Figures 1BF).

Cross-correlation of the original non-stationary time-series SARS-CoV-2 data was also performed to illustrate the potential for spurious correlation (Figures 1CDGH). Both time-series concentration data (liquid and solid fractions) from SI (Figures 1CD) and HO (Figures 1GH) WWTPs showed all positive correlation coefficients and the majority of the correlation analyses showed statistically significant correlations with p -values less than 0.05 (SI: 21 out of 42 analyses; HO: 27 out of 42 analyses) with the original daily clinical case numbers (x_{t+h} , $h = \text{lag number}$). Because the normality assumption of cross-correlation analysis was not met, these high positive correlation coefficients are considered spurious and false positive.

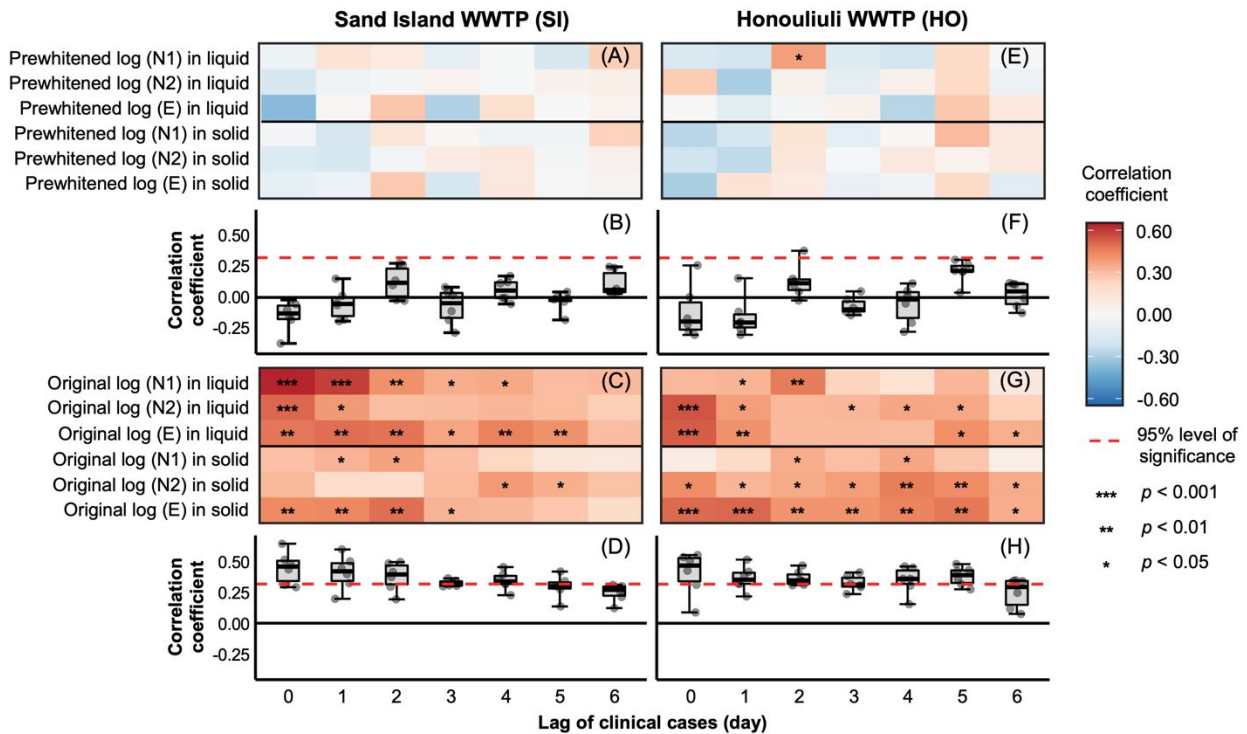


Figure 1. Cross-correlation between with and without prewhitening by the first differences of daily new clinical COVID-19 case numbers and measured SARS-CoV-2 RNA concentration (log 10 transformed) in wastewater samples from Sand Island (prewhitened: A, B; non-prewhitened: C, D) and Honouliuli (prewhitened: E, F; non-prewhitened: G, H). Red dashed lines represent a 95% level of significance and the p -value of the correlation less than 0.05 are displayed as asterisks. The middle, upper, and lower lines in the box of the boxplot represent the median, 25th, and 75th percentiles, respectively, and the whiskers represent the largest and smallest values outside of the interquartile range.

3.2. Impact of normalized abundance on cross-correlation

The concentrations of SARS-CoV-2 RNA measured from the samples of the two WWTPs are expected to be impacted by various processes during wastewater sampling and sample processing, including total fecal discharge in the area, sewer collection in the WWTPs, wastewater viral precipitation method, and molecular quantification steps. The resulting variations could be potentially mitigated by normalizing data to various endogenous fecal viral indicators (e.g. log N1/G2) (21). The normalized daily wastewater SARS-CoV-2 RNA data were

analyzed via cross-correlation with the prewhitened daily clinical case numbers. For wastewater samples from the SI WWTP, the normalization strategy produced significantly different cross-correlation patterns against different time lags (Figure 2) than those without normalization (Figure 1).

The most obvious improvement in the cross-correlation coefficient was observed at a two-day lag (x_{t+2} , Figures 2BDF), with a range of $r = -0.03-0.45$ (0.23 ± 0.13). The average cross-correlation coefficients were increased from 0.12 to 0.33, 0.19, and 0.16 for G2, G3, and PMMoV normalizations, respectively, which showed an average of 0.11 ± 0.09 increase than those without normalization (Figures 2BDF). Among all combinations, the best correlation coefficient ($r = 0.45$, $p = 0.004$) was observed between the log (E/G2) in the liquid fraction and the daily clinical case numbers. The normalization strategy showed a higher improvement of correlation coefficients in liquid fractions (0.14 improved) than in the solid fractions (0.07 improved) compared to the raw data.

At the two-day lag, statistically significant correlations were observed more frequently with the normalized SARS-CoV-2 abundance data than with the raw data. For example, at the SI WWTP, both log (N1/G2) and log (E/G2) showed statistically significant correlation coefficients in the liquid ($r = 0.38$ ($p = 0.017$) and $r = 0.45$ ($p = 0.004$), respectively) and solid fractions ($r = 0.32$ ($p = 0.048$) and $r = 0.38$ ($p = 0.018$), respectively). In contrast, there was no statistically significant correlation between clinical cases and raw wastewater SARS-CoV-2 concentration data at SI WWTP (Figure 1A). G3 normalization showed two statistically significant correlations from log (N1/G3) and log (E/G3) ($r = 0.33$ ($p = 0.043$) and $r = 0.33$ ($p = 0.044$), respectively) in the liquid fractions. PMMoV normalization showed only one statistically significant correlation coefficient from the log (E/PMMoV) in the solid fraction ($r = 0.33$, $p = 0.043$).

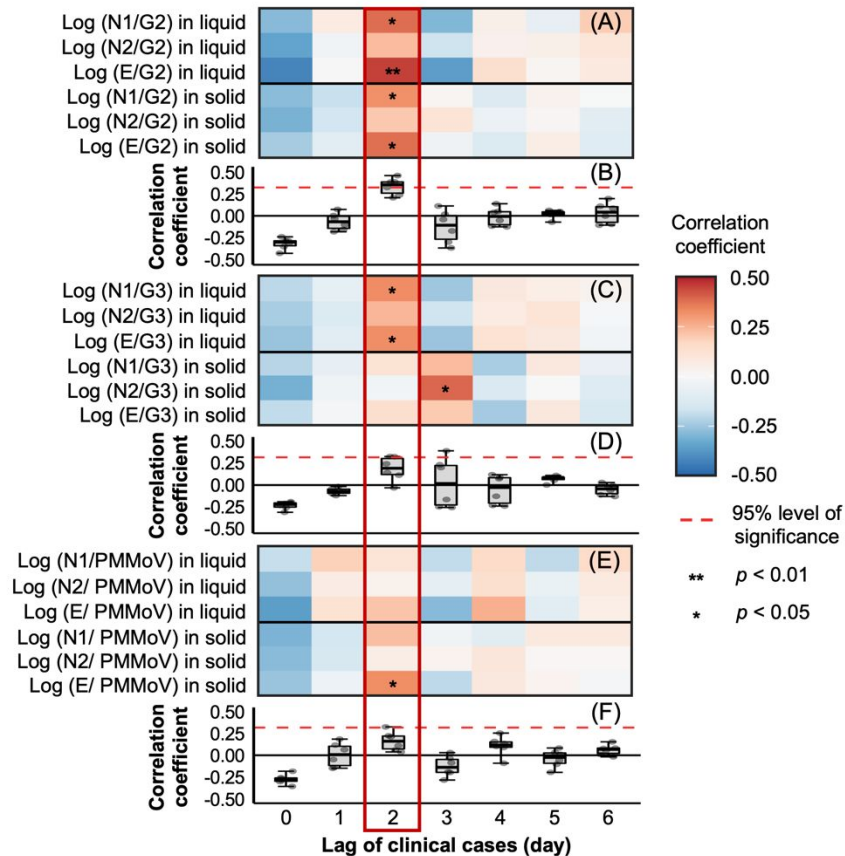


Figure 2. Cross-correlation between the prewhitened COVID-19 new case numbers and the prewhitened SARS-CoV-2 RNA normalized abundance in wastewater samples from the SI WWTP. The normalized abundance was calculated by dividing SARS-CoV-2 RNA abundance by F+ RNA coliphage Group II (A, B), Group III (C, D), and PMMoV (E, F). All normalized abundances were transformed into log forms. Red dashed lines represent a 95% level of significance and the p -value of the correlation less than 0.05 are displayed as asterisks. The middle, upper, and lower lines in the box of the boxplot represent the median, 25th, and 75th percentiles, respectively, and the whiskers represent the largest and smallest values outside of the interquartile range.

For the HO WWTP, the largest correlation coefficients from the cross-correlation between the normalized SARS-CoV-2 abundance and the daily clinical case numbers were also observed at a two-day lag (x_{t+2} , Figures 3BD), which is similar to the results of SI WWTP. The average cross-correlation coefficients were increased from 0.13 to 0.21 and 0.24 for G2 and G3 normalizations, respectively, which showed an average of 0.10 ± 0.02 increase than those without normalization (Figures 3BD). The best correlation coefficient among all three

normalization methods was observed between the log (N1/G2) in the liquid fraction and the daily clinical case numbers ($r = 0.35$, $p = 0.029$). Similar to the SI WWTP, the normalization strategy showed a larger improvement of correlation coefficients in liquid fractions (0.07 improved) than in the solid fractions (0.03 improved) compared to the raw data.

At the two-day lag, all normalized forms of log N1 in the liquid fractions (log (N1/G2): $r = 0.35$, $p = 0.029$; log (N1/G3): $r = 0.35$, $p = 0.030$; log (N1/PMMoV): $r = 0.33$, $p = 0.043$) showed statistically significant correlations between the daily clinical cases and the normalized SARS-CoV-2 RNA abundance from HO WWTP (Figures 3ACE). Although no statistically significant correlation was observed from solid fractions from any normalization strategies ($r = -0.05-0.27$, 0.16 ± 0.11), all correlation coefficients observed from the solid fractions were increased after normalization by G2 and G3.

The correlation coefficients of all genes from both liquid and solid fractions of HO WWTP were decreased when they were normalized by PMMoV with average correlation coefficients decreased from 0.13 to 0.09 at the two-day lag (Figure 3E). Interestingly, PMMoV normalized SARS-CoV-2 RNA abundances in the liquid fraction showed large correlation coefficients ($r = 0.32 \pm 0.06$) at a five-day lag of daily clinical case numbers (Figure 3F) showing a statistically significant correlation with log (E/PMMoV) ($r = 0.39$, $p = 0.014$).

The better correlation exhibited by the SI WWTP than the HO WWTP could be due to the SI WWTP treating a much larger fraction of the City's wastewater (58%) than the HO WWTP (24%), and hence is expected to be more representative of the community disease burden.

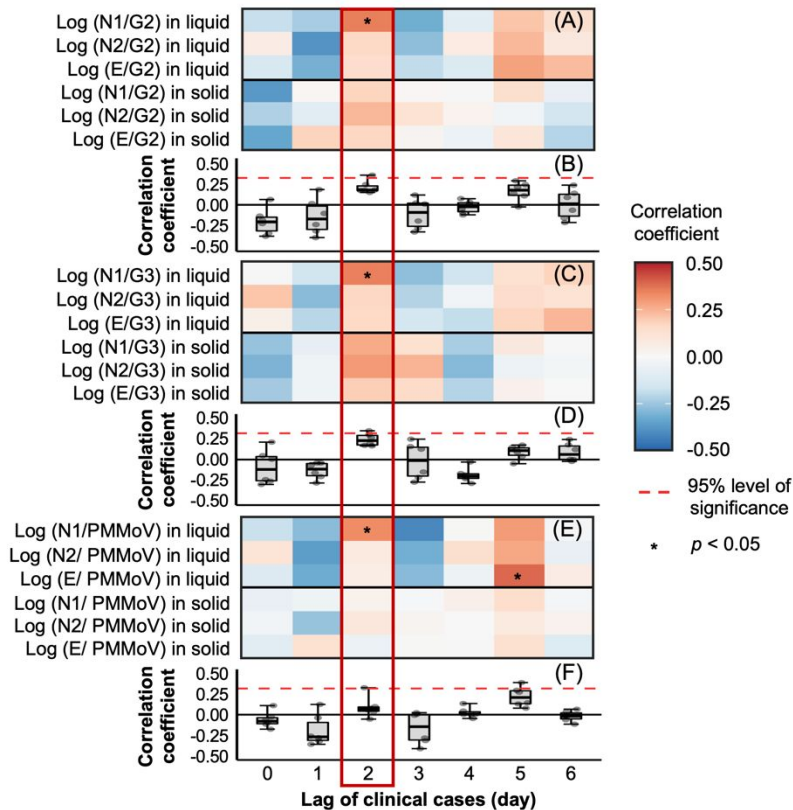


Figure 3. Cross-correlation between the prewhitened COVID-19 new case numbers and the prewhitened SARS-CoV-2 RNA normalized abundance in wastewater samples from the HO WWTP. The normalized abundance was calculated by dividing SARS-CoV-2 RNA abundance by F+ RNA coliphage Group II (A, B), Group III (C, D), and PMMoV (E, F). All normalized abundances were transformed into log forms. Red dashed lines represent a 95% level of significance and the p -value of the correlation less than 0.05 are displayed as asterisks. The middle, upper, and lower lines in the box of the boxplot represent the median, 25th, and 75th percentiles, respectively, and the whiskers represent the largest and smallest values outside of the interquartile range.

3.3. Comparison of different normalization strategies for cross-correlation analysis

The correlation coefficients between the daily clinical case numbers and SARS-CoV-2 RNA abundance at a two-day lag of clinical cases were compared with respect to the normalization strategies by using a pairwise t -test (Figure 4). Among the three different endogenous fecal viral RNA controls used, only normalizing the data with G2 ($r = 0.33 \pm 0.10$, $p = 0.002$) showed a statistically significant improvement of correlation coefficients in comparison

to that using the raw data ($r = 0.12 \pm 0.13$). While G3 ($r = 0.19 \pm 0.14$, $p = 0.396$) and PMMoV ($r = 0.16 \pm 0.11$, $p = 0.198$) mildly improved the correlation (Figure 4A), but not statistically.

Furthermore, the G2 normalization showed a significantly larger correlation coefficient value than both G3 ($p = 0.04$) and PMMoV ($p = 0.014$) normalization methods.

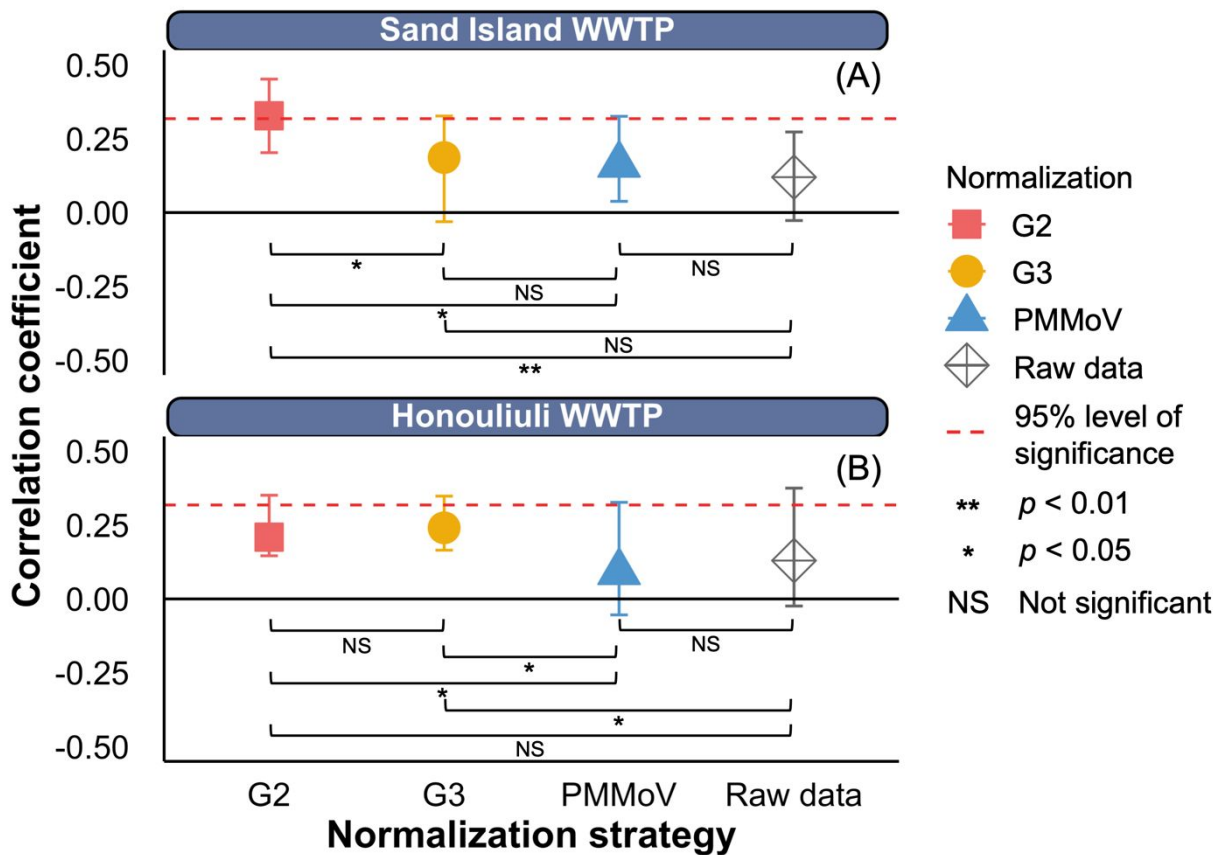


Figure 4. Cross-correlation between both prewhitened daily new clinical COVID-19 case numbers in Honolulu County and normalized SARS-CoV-2 RNA concentration (log 10 transformed) by F+ RNA coliphage Group II, Group III, and PMMoV in wastewater samples from Sand Island (A) and Honouliuli (B) at a two-day lag of clinical cases. Red dashed lines represent a 95% level of significance and the whiskers represent the largest and the smallest values.

For the HO WWTP, the G3 ($r = 0.24 \pm 0.08$, $p = 0.032$) normalization method was the only strategy that significantly improved the correlation coefficients (Figure 4B). Although G2 ($r = 0.21 \pm 0.08$, $p = 0.058$) did not statistically improve the correlations with the raw data, it increased the mean values of correlation coefficients ($r = 0.13 \pm 0.14$). However, PMMoV ($r = 0.09 \pm 0.13$, $p = 0.32$) normalization even decreased the mean values of the correlation

coefficients than the raw data. Moreover, G2 ($p = 0.028$) and G3 ($p = 0.028$) normalizations improved the correlation when compared with PMMoV normalization.

The overall results indicate that the normalization of SARS-CoV-2 RNA abundance improved the cross-correlations with the daily clinical case numbers. G2 normalization showed the largest improvement of cross-correlations in SI WWTP samples. On the contrary, the G3 normalization was the largest improvement of the cross-correlations in HO WWTP samples. When considering the liquid fractions only, the G2 normalization of log N1 from SI ($r = 0.38$, $p = 0.017$) and HO ($r = 0.35$, $p = 0.029$) WWTPs showed both significant correlations with the daily clinical case numbers.

4. Discussion

In our previous study (21), we observed simultaneous downward trends between SARS-CoV-2 RNA abundances (both with and without normalization by fecal viral markers) in wastewater samples from the SI and HO WWTPs and the daily clinical COVID-19 case numbers during a COVID-19 public health lockdown. This is congruent with previous observations where increases in wastewater SARS-CoV-2 RNA concentration corresponded with rapidly expanding COVID-19 outbreaks (17, 36-39). The fine-scale temporal dynamics afforded by daily sampling also detected significant intra-day fluctuation of the wastewater SARS-CoV-2 RNA abundance, even within the same weeks. Many factors could have contributed to the observed intra-day fluctuation, including errors in wastewater sampling and sample analysis, variations in viral shedding by infected individuals, and daily fluctuations in disease burden in the community. Since similar trends were detected in the two replicate WWTPs, over multiple weeks, and regardless of normalization strategies, the former two (i.e., sampling and analysis errors and variation in viral shedding) are unlikely to entirely explain the observations.

Cross-correlation analysis of the time-series data of wastewater SARS-CoV-2 abundance and clinical case numbers could be used to infer potential association and determine if daily fluctuations in disease burden in the community contributed to the observed intra-day fluctuation. Many wastewater surveillance studies have compared wastewater SARS-CoV-2 abundance with clinical case data in the community through correlation analysis (17-20). However, few previous studies have conducted prewhitening treatment to achieve stationarity of the time-series data. Stationarity in time-series data indicates consistency of the distribution (mean and variance) over time (40), and non-stationary time-series data can often lead to spurious outcomes in correlation analyses. This was clearly demonstrated when we observed significant spurious cross-correlations with the original data that contain trends (Figure 1 CDGH). Similar phenomenon may explain reports of wastewater surveillance showing strong correlations to clinical cases, especially when the studies were conducted during a period when COVID-19 clinical cases were continuously increasing or decreasing (36, 41, 42). Therefore, prewhitening the data for wastewater surveillance to meet the stationarity requirement of cross-correlation analysis must be practiced in order to identify the actual association between data sets.

After prewhitening the data, the cross-correlation coefficients decreased significantly, and the overall patterns with respect to the time lag also changed drastically (Figure 1). For example, at zero-time lag, cross-correlation using the original data detected the best positive correlation, whereas the prewhitened data actually detected some negative correlations. The low levels of correlation detected are not entirely unexpected, considering the extraordinary complexity involved in collecting the wastewater SARS-CoV-2 RNA data and resulting variations. Many factors, including varying fecal discharge by infected individuals, dilution and fluctuation during transportation in sanitary sewers, and wastewater sample collection and processing, could have contributed to the variations in SARS-CoV-2 RNA in wastewater samples. The molecular quantification processes could also introduce additional variations to the results; for example,

RNA recovery during sample extraction could have varying efficiencies, and subsequent reverse transcription and qPCR quantification could introduce additional biases.

The incorporation of normalization strategies and the resulting relative abundance of SARS-CoV-2 RNA in the wastewater samples led to the identification of a two-day lag showing the best correlation (Figures 2 and 3). Given the complex and multi-step process required for quantifying SARS-CoV-2 RNA in wastewater, the selection of endogenous viral RNA control for global normalization is also important to reduce the variations during the analysis and enable statistical comparison. Amongst the three fecal RNA viruses as endogenous controls tested in this study, normalization by G2 provided the most significant improvement in correlation between wastewater SARS-CoV-2 RNA abundance and clinical new cases. Cole *et al.* found that G2 had the highest proportion among the total F+ RNA groups in WWTP samples (51.9%) and G2 was found more in human-impacted wastes than G3 (43). This supports our results of G2 normalization of SARS-CoV-2 RNA abundance having higher correlation coefficients compared to G3. On the other hand, PMMoV only provided marginal improvement in correlation. This difference could be attributed to their respective sources in human feces where G2 and G3 are inherently linked with fecal coliforms while PMMoV is subjected to dietary variation in pepper consumption. Some previous wastewater surveillance studies that used PMMoV for the SARS-CoV-2 abundance normalization also reported that the PMMoV did not improve the correlation with the clinical cases (44-46). Other biomarkers and chemical indicators for population normalization are also recently considered, such as paraxanthine (47), cross-assembly phage (48, 49), Human RNase P (50), total nitrogen, and phosphate (51). Therefore, more studies related to improving the normalization methods in the wastewater surveillance field with the consideration of stationarity are required to efficiently reduce the variations and accurately quantify SARS-CoV-2 abundance.

The significant cross-correlation between the normalized abundance of daily wastewater SARS-CoV-2 and new clinical cases in the community is highly intriguing. Since the COVID-19 pandemic, weekly intra-day oscillations in new clinical case numbers have been widely observed in communities across the globe (52). One school of thought is that these weekly intra-day oscillations are primarily a reflection of diagnostic and reporting biases (53, 54), while a competing theory is that this could be caused by actual disease transmission dynamics caused by weekly behavior patterns (55-57). The strong correlation between the intra-day fluctuation and weekly oscillation of wastewater SARS-CoV-2 RNA abundance and clinical case numbers that we observed in this study suggest that the observed weekly oscillation of clinical cases may be indeed a true reflection of disease burden in the community, in addition to contributions from clinical sampling and reporting biases and errors.

Since the average turnaround for clinical testing during the study period was approximately one day, with the assumption of one day lag between symptom onset and specimen collection, the observed two-day lag in cross-correlation analysis indicates that the wastewater SARS-CoV-2 viral RNA abundance may be synchronizing with symptom development of new COVID-19 cases in the community. Studies at the early stage of the pandemic, which likely experienced clinical testing delays, have reported the detection of the SARS-CoV-2 viral RNA in wastewater about one week ahead of reported clinical cases in the communities (17, 37), while a study reported wastewater sludge SARS-CoV-2 concentration leading the specimen collection by 0-2 days (23). The apparent synchronous correspondence supports the possibility of using wastewater for early detection of viral transmission in communities, as viral shedding can start 3-5 days before and peaks around symptom onset (12) and asymptomatic infections and often lead to symptomatic infections and represent a significant portion of the overall community disease burden (58).

In this study, both the solid and liquid fractions of the same wastewater samples were analyzed separately, and SARS-CoV-2 normalized abundance data in both fractions showed similar cross-correlation patterns (Figures 2 and 3). In the previous study (21), the solid fractions contained higher per mass concentration of SARS-CoV-2 RNA than the liquid fractions, while the normalized abundances between the two fractions were quite similar. The normalized abundance of SARS-CoV-2 RNA in liquid fractions exhibited slightly stronger correlations with the clinical COVID-19 case numbers in the community than the normalized abundance of SARS-CoV-2 in the solid fraction. This could be attributed to the more complex matrix effects in the solid fractions than the liquid fractions, as indicated in our previous study, where lower recovery and higher variation of the spiked BCoV as exogenous control were observed in the solid fractions than in the liquid fractions (21).

It is important to note that all three quantification assays showed similar cross-correlation patterns between normalized SARS-CoV-2 RNA abundance in wastewater and clinical case numbers in the community. While the SARS-CoV-2 RNA genome contains a single copy of N and E genes, our previously published study (21) and many other studies (37, 59) have shown that different assays usually generate different abundance levels, indicating that the molecular quantification processes have variations. Nevertheless, all three assays were able to reveal strong correlations at a two-day lag between normalized SARS-CoV-2 RNA abundance in the wastewater and community disease burden, with the N1 gene assay providing the highest correlation coefficients in both tested WWTPs (Figures 2 and 3). Even though the E gene showed strong positive cross-correlation patterns at a two-day lag from both SI and HO WWTPs, it is considered the least specific PCR target for SARS-CoV-2 detection due to homologous sequence similarities with other coronaviruses together with recurrent mutations (60). Our previous paper (21) showed that the N2 gene assay had the least sensitivity, consequently, more wastewater surveillance studies for SARS-CoV-2 RNA detection are using the N1 assay and showed higher

positivity rates compared to the E gene assay (61, 62). With the reasons above, we would recommend the N1 gene assay for future wastewater surveillance for SARS-CoV-2 RNA detection.

5. Conclusions

This study demonstrated the importance of prewhitening to remove the trends of the daily fluctuation of wastewater surveillance data and the clinical case numbers before cross-correlation analysis of the time-series data sets to avoid spurious correlations. We also observed that normalization strategies to account for various variations in the process are helpful in improving cross-correlation analysis. Amongst the various normalization strategies, SARS-CoV-2 RNA abundances normalized with F+RNA coliphage G2 provided the best correlation coefficients in this study. We observed that the N1 assay was showing the best correlation, while N2 showed less sensitivity and the E gene has been reported to be less specific. Although there were significant inherent variations in the data due to the complexity of the samples, we could draw a conclusion based on multiple gene markers and multiple normalization strategies that the daily clinical case numbers were two days behind the SARS-CoV-2 RNA detection in wastewater. This supports the notion that wastewater surveillance has the potential to provide earlier detection of SARS-CoV-2 signal than clinical diagnosis. Most interestingly, the strong cross-correlation between the intra-day fluctuation and weekly oscillation of wastewater SARS-CoV-2 RNA abundance and clinical cases suggest that the observed weekly oscillation of clinical cases is likely caused (at least partially) by disease burden fluctuation in the community in addition to clinical sampling and reporting biases and errors, which requires further research to delineate their respective contributions.

Conflicts of interest

There are no conflicts to declare.

Acknowledgments

This material is based upon work supported by the National Science Foundation under Grant No. CBET-2027059.

References

1. Zhou P, Yang X-L, Wang X-G, Hu B, Zhang L, Zhang W, et al. A pneumonia outbreak associated with a new coronavirus of probable bat origin. *nature*. 2020;579(7798):270-3.
2. Chen Y, Chen L, Deng Q, Zhang G, Wu K, Ni L, et al. The presence of SARS-CoV-2 RNA in the feces of COVID-19 patients. *Journal of medical virology*. 2020;92(7):833-40.
3. Kashi AH, De la Rosette J, Amini E, Abdi H, Fallah-Karkan M, Vaezjalali M. Urinary viral shedding of COVID-19 and its clinical associations: a systematic review and meta-analysis of observational studies. *Medrxiv*. 2020.
4. To KK-W, Tsang OT-Y, Leung W-S, Tam AR, Wu T-C, Lung DC, et al. Temporal profiles of viral load in posterior oropharyngeal saliva samples and serum antibody responses during infection by SARS-CoV-2: an observational cohort study. *The Lancet infectious diseases*. 2020;20(5):565-74.
5. Wyllie AL, Fournier J, Casanovas-Massana A, Campbell M, Tokuyama M, Vijayakumar P, et al. Saliva or nasopharyngeal swab specimens for detection of SARS-CoV-2. *New England Journal of Medicine*. 2020;383(13):1283-6.
6. Long Q-X, Tang X-J, Shi Q-L, Li Q, Deng H-J, Yuan J, et al. Clinical and immunological assessment of asymptomatic SARS-CoV-2 infections. *Nature medicine*. 2020;26(8):1200-4.
7. Buitrago-Garcia D, Egli-Gany D, Counotte MJ, Hossmann S, Imeri H, Ipekci AM, et al. Occurrence and transmission potential of asymptomatic and presymptomatic SARS-CoV-2 infections: A living systematic review and meta-analysis. *PLoS medicine*. 2020;17(9):e1003346.
8. Schmitz BW, Innes GK, Prasek SM, Betancourt WQ, Stark ER, Foster AR, et al. Enumerating asymptomatic COVID-19 cases and estimating SARS-CoV-2 fecal shedding rates via wastewater-based epidemiology. *Science of The Total Environment*. 2021;801:149794.
9. Park S-k, Lee C-W, Park D-I, Woo H-Y, Cheong HS, Shin HC, et al. Detection of SARS-CoV-2 in fecal samples from patients with asymptomatic and mild COVID-19 in Korea. *Clinical Gastroenterology and Hepatology*. 2021;19(7):1387-94. e2.
10. Li X, Kulandaivelu J, Guo Y, Zhang S, Shi J, O'Brien J, et al. SARS-CoV-2 shedding sources in wastewater and implications for wastewater-based epidemiology. *Journal of hazardous materials*. 2022;432:128667.
11. He X, Lau EH, Wu P, Deng X, Wang J, Hao X, et al. Temporal dynamics in viral shedding and transmissibility of COVID-19. *Nature medicine*. 2020;26(5):672-5.

12. Jones DL, Baluja MQ, Graham DW, Corbishley A, McDonald JE, Malham SK, et al. Shedding of SARS-CoV-2 in feces and urine and its potential role in person-to-person transmission and the environment-based spread of COVID-19. *Science of the Total Environment*. 2020;749:141364.
13. Cevik M, Tate M, Lloyd O, Maraolo AE, Schafers J, Ho A. SARS-CoV-2, SARS-CoV, and MERS-CoV viral load dynamics, duration of viral shedding, and infectiousness: a systematic review and meta-analysis. *The lancet microbe*. 2021;2(1):e13-e22.
14. Miura F, Kitajima M, Omori R. Duration of SARS-CoV-2 viral shedding in faeces as a parameter for wastewater-based epidemiology: Re-analysis of patient data using a shedding dynamics model. *Science of The Total Environment*. 2021;769:144549.
15. Cryer JD, Chan K-S. *Time series analysis: with applications in R*: Springer; 2008.
16. Bayazit M, Önöz B. To prewhiten or not to prewhiten in trend analysis? *Hydrological Sciences Journal*. 2007;52(4):611-24.
17. Peccia J, Zulli A, Brackney DE, Grubaugh ND, Kaplan EH, Casanovas-Massana A, et al. Measurement of SARS-CoV-2 RNA in wastewater tracks community infection dynamics. *Nature biotechnology*. 2020;38(10):1164-7.
18. Aberi P, Arabzadeh R, Insam H, Markt R, Mayr M, Kreuzinger N, et al. Quest for optimal regression models in SARS-CoV-2 wastewater based epidemiology. *International journal of environmental research and public health*. 2021;18(20):10778.
19. Lara-Jacobo LR, Islam G, Desaulniers J-P, Kirkwood AE, Simmons DB. Detection of SARS-CoV-2 Proteins in Wastewater Samples by Mass Spectrometry. *Environmental science & technology*. 2022;56(8):5062-70.
20. Li L, Mazurowski L, Dewan A, Carine M, Haak L, Guarin TC, et al. Longitudinal monitoring of SARS-CoV-2 in wastewater using viral genetic markers and the estimation of unconfirmed COVID-19 cases. *Science of The Total Environment*. 2022;817:152958.
21. Li B, Di DYW, Saingam P, Jeon MK, Yan T. Fine-scale temporal dynamics of SARS-CoV-2 RNA abundance in wastewater during a COVID-19 lockdown. *Water Research*. 2021;197:117093.
22. Hjelmsø MH, Hellmér M, Fernandez-Cassi X, Timoneda N, Lukjancenko O, Seidel M, et al. Evaluation of methods for the concentration and extraction of viruses from sewage in the context of metagenomic sequencing. *PloS one*. 2017;12(1):e0170199.
23. Lu X, Wang L, Sakthivel SK, Whitaker B, Murray J, Kamili S, et al. US CDC real-time reverse transcription PCR panel for detection of severe acute respiratory syndrome coronavirus 2. *Emerging infectious diseases*. 2020;26(8):1654.
24. Corman VM, Landt O, Kaiser M, Molenkamp R, Meijer A, Chu DK, et al. Detection of 2019 novel coronavirus (2019-nCoV) by real-time RT-PCR. *Eurosurveillance*. 2020;25(3):2000045.

25. Friedman SD, Cooper EM, Calci KR, Genthner FJ. Design and assessment of a real time reverse transcription-PCR method to genotype single-stranded RNA male-specific coliphages (Family Leviviridae). *Journal of virological methods*. 2011;173(2):196-202.
26. Rosario K, Symonds EM, Sinigalliano C, Stewart J, Breitbart M. Pepper mild mottle virus as an indicator of fecal pollution. *Applied and environmental microbiology*. 2009;75(22):7261-7.
27. Shapiro SS, Wilk MB. An analysis of variance test for normality (complete samples). *Biometrika*. 1965;52(3/4):591-611.
28. Mann HB. Nonparametric tests against trend. *Econometrica: Journal of the econometric society*. 1945:245-59.
29. Kendall MG. Rank correlation methods. 1948.
30. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)*. 1995;57(1):289-300.
31. Team RC. A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing; 2022. 2022.
32. Wickham H, Averick M, Bryan J, Chang W, McGowan LDA, François R, et al. Welcome to the Tidyverse. *Journal of open source software*. 2019;4(43):1686.
33. Kassambara A. ggpubr: 'ggplot2' Based Publication Ready Plots. R package version 0.4.0. Computer software] <https://cran-r-project.org/web/packages/ggpubr/index.html>. 2020.
34. Wickham H, Seidel D. scales: Scale functions for visualization. R package version. 2020;1(1):678.
35. Kassambara A. rstatix: Pipe-friendly framework for basic statistical tests. R package version 0.6.0. 2020.
36. Róka E, Khayer B, Kis Z, Kovács LB, Schuler E, Magyar N, et al. Ahead of the second wave: early warning for COVID-19 by wastewater surveillance in Hungary. *Science of The Total Environment*. 2021;786:147398.
37. Medema G, Heijnen L, Elsinga G, Italiaander R, Brouwer A. Presence of SARS-Coronavirus-2 RNA in sewage and correlation with reported COVID-19 prevalence in the early stage of the epidemic in the Netherlands. *Environmental Science & Technology Letters*. 2020;7(7):511-6.
38. Randazzo W, Truchado P, Cuevas-Ferrando E, Simón P, Allende A, Sánchez G. SARS-CoV-2 RNA in wastewater anticipated COVID-19 occurrence in a low prevalence area. *Water research*. 2020;181:115942.
39. Ahmed W, Angel N, Edson J, Bibby K, Bivins A, O'Brien JW, et al. First confirmed detection of SARS-CoV-2 in untreated wastewater in Australia: a proof of concept for the wastewater surveillance of COVID-19 in the community. *Science of the Total Environment*. 2020;728:138764.

40. Kwiatkowski D, Phillips PC, Schmidt P, Shin Y. Testing the null hypothesis of stationarity against the alternative of a unit root: How sure are we that economic time series have a unit root? *Journal of econometrics*. 1992;54(1-3):159-78.
41. Sangsanont J, Rattanakul S, Kongprajug A, Chyerochana N, Sresung M, Sriporatana N, et al. SARS-CoV-2 RNA surveillance in large to small centralized wastewater treatment plants preceding the third COVID-19 resurgence in Bangkok, Thailand. *Science of the Total Environment*. 2022;809:151169.
42. Weidhaas J, Aanderud ZT, Roper DK, VanDerslice J, Gaddis EB, Ostermiller J, et al. Correlation of SARS-CoV-2 RNA in wastewater with COVID-19 disease burden in sewersheds. *Science of The Total Environment*. 2021;775:145790.
43. Cole D, Long SC, Sobsey MD. Evaluation of F+ RNA and DNA coliphages as source-specific indicators of fecal contamination in surface waters. *Applied and environmental microbiology*. 2003;69(11):6507-14.
44. Ai Y, Davis A, Jones D, Lemeshow S, Tu H, He F, et al. Wastewater SARS-CoV-2 monitoring as a community-level COVID-19 trend tracker and variants in Ohio, United States. *Science of The Total Environment*. 2021;801:149757.
45. Feng S, Roguet A, McClary-Gutierrez JS, Newton RJ, Kloczko N, Meiman JG, et al. Evaluation of sampling, analysis, and normalization methods for SARS-CoV-2 concentrations in wastewater to assess COVID-19 burdens in Wisconsin communities. *ACS ES&T Water*. 2021;1(8):1955-65.
46. Greenwald HD, Kennedy LC, Hinkle A, Whitney ON, Fan VB, Crits-Christoph A, et al. Tools for interpretation of wastewater SARS-CoV-2 temporal and spatial trends demonstrated with data collected in the San Francisco Bay Area. *Water Research X*. 2021;12:100111.
47. Hsu S-Y, Bayati M, Li C, Hsieh H-Y, Belenchia A, Klutts J, et al. Biomarkers Selection for Population Normalization in SARS-CoV-2 Wastewater-based Epidemiology. *Water research*. 2022:118985.
48. Bivins A, Crank K, Greaves J, North D, Wu Z, Bibby K. Cross-assembly phage and pepper mild mottle virus as viral water quality monitoring tools—potential, research gaps, and way forward. *Current Opinion in Environmental Science & Health*. 2020;16:54-61.
49. Green H, Wilder M, Collins M, Fenty A, Gentile K, Kmush BL, et al. Quantification of SARS-CoV-2 and cross-assembly phage (crAssphage) from wastewater to monitor coronavirus transmission within communities. *MedRxiv*. 2020.
50. El-Malah SS, Saththasivam J, Jabbar KA, Arun K, Gomez TA, Ahmed AA, et al. Application of human RNase P normalization for the realistic estimation of SARS-CoV-2 viral load in wastewater: A perspective from Qatar wastewater surveillance. *Environmental Technology & Innovation*. 2022;27:102775.
51. Isaksson F, Lundy L, Hedström A, Székely AJ, Mohamed N. Evaluating the use of alternative normalization approaches on sars-cov-2 concentrations in wastewater: Experiences from two catchments in northern sweden. *Environments*. 2022;9(3):39.

52. Bukhari Q, Jameel Y, Massaro JM, D'Agostino RB, Khan S. Periodic oscillations in daily reported infections and deaths for coronavirus disease 2019. *JAMA network open*. 2020;3(8):e2017521-e.
53. Bragato PL. Assessment of the weekly fluctuations of the Covid-19 cases in Italy and worldwide. *Preprints*. 2020.
54. Bergman A, Sella Y, Agre P, Casadevall A. Oscillations in US COVID-19 incidence and mortality data reflect diagnostic and reporting factors. *Msystems*. 2020;5(4):e00544-20.
55. Cecil WT. COVID-19: daily fluctuations, a weekly cycle, and a negative trend. *The American Journal of Managed Care*. 2020;26(7):284-5.
56. Derakhshan M, Ansarian HR, Ghomshei M. Temporal variations in COVID-19: an epidemiological discussion with a practical application. *Journal of International Medical Research*. 2021;49(8):03000605211033208.
57. Ricon-Becker I, Tarrasch R, Blinder P, Ben-Eliyahu S. A seven-day cycle in COVID-19 infection and mortality rates: Are inter-generational social interactions on the weekends killing susceptible people. *medRxiv*. 2020.
58. Oran DP, Topol EJ. Prevalence of asymptomatic SARS-CoV-2 infection: a narrative review. *Annals of internal medicine*. 2020;173(5):362-7.
59. Nalla AK, Casto AM, Huang M-LW, Perchetti GA, Sampoleo R, Shrestha L, et al. Comparative performance of SARS-CoV-2 detection assays using seven different primer-probe sets and one assay kit. *Journal of clinical microbiology*. 2020;58(6):e00557-20.
60. Artesi M, Bontems S, Göbbels P, Franckh M, Maes P, Boreux R, et al. A recurrent mutation at position 26340 of SARS-CoV-2 is associated with failure of the E gene quantitative reverse transcription-PCR utilized in a commercial dual-target diagnostic assay. *Journal of clinical microbiology*. 2020;58(10):e01598-20.
61. Muenchhoff M, Mairhofer H, Nitschko H, Grzimek-Koschewa N, Hoffmann D, Berger A, et al. Multicentre comparison of quantitative PCR-based assays to detect SARS-CoV-2, Germany, March 2020. *Eurosurveillance*. 2020;25(24):2001057.
62. Pérez-Cataluña A, Cuevas-Ferrando E, Randazzo W, Falcó I, Allende A, Sánchez G. Comparing analytical methods to detect SARS-CoV-2 in wastewater. *Science of the Total Environment*. 2021;758:143870.