



High-throughput sequencing of SARS-CoV-2 in wastewater provides insights into circulating variants

Rafaela S. Fontenele^{a,b,1}, Simona Kraberger^{a,1}, James Hadfield^c, Erin M. Driver^d, Devin Bowes^d, LaRinda A. Holland^a, Temitope O.C. Faley^d, Sangeet Adhikari^{d,e}, Rahul Kumar^d, Rosa Inchausti^f, Wydale K. Holmes^f, Stephanie Deitrick^g, Philip Brown^h, Darrell Dutyⁱ, Ted Smith^j, Aruni Bhatnagar^j, Ray A. Yeager II^j, Rochelle H. Holm^j, Natalia Hoogesteijn von Reitzenstein^k, Elliott Wheeler^k, Kevin Dixon^k, Tim Constantine^k, Melissa A. Wilson^{b,1}, Efreem S. Lim^{a,b}, Xiaofang Jiang^m, Rolf U. Halden^{d,n}, Matthew Scotch^{d,o}, Arvind Varsani^{a,b,1,*}

^a The Biodesign Center for Fundamental and Applied Microbiomics, Arizona State University, 1001 S. McAllister Ave., Tempe, AZ 85281, USA

^b School of Life Sciences, Arizona State University, 427 East Tyler Mall, Tempe, AZ 85287, USA

^c Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA

^d Biodesign Center for Environmental Health Engineering, Biodesign Institute, Arizona State University, 1001 S. McAllister Ave., Tempe, AZ 85281, USA

^e School of Sustainable Engineering and the Built Environment, Arizona State University, Tempe, AZ USA

^f Strategic Management and Diversity Office, City of Tempe, 31 E Fifth Street, Tempe, AZ 85281, USA

^g Enterprise GIS & Data Analytics, Information Technology, 31 E Fifth Street, City of Tempe, Tempe, AZ 85281, USA

^h Municipal Utilities, City of Tempe, 31 E Fifth Street, Tempe, AZ 85281, USA

ⁱ Tempe Fire Medical Rescue, 31 E Fifth Street, City of Tempe, Tempe, AZ 85281, USA

^j Christina Lee Brown Envirome Institute, University of Louisville, 302 E. Muhammad Ali Blvd., Louisville, KY 40202, USA

^k Jacobs Engineering Group Inc., 1999 Bryan Street, Dallas, TX 75201, USA

^l Center for Evolution and Medicine, Arizona State University, 401 E. Tyler Mall, Tempe, AZ 85287, USA

^m National Library of Medicine, National Institute of Health, 8600 Rockville Pike, Bethesda, MD 20894, USA

ⁿ OneWaterOneHealth, Nonprofit Project of the Arizona State University Foundation, 1001 S. McAllister Ave., Tempe, AZ 85281, USA

^o College of Health Solutions, Arizona State University, 550 N. 3rd St, Phoenix, AZ 85004, USA

ARTICLE INFO

Keywords:

SARS-CoV-2
Wastewater
Surveillance
Wastewater-based epidemiology
High-throughput sequencing

ABSTRACT

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) likely emerged from a zoonotic spill-over event and has led to a global pandemic. The public health response has been predominantly informed by surveillance of symptomatic individuals and contact tracing, with quarantine, and other preventive measures have then been applied to mitigate further spread. Non-traditional methods of surveillance such as genomic epidemiology and wastewater-based epidemiology (WBE) have also been leveraged during this pandemic. Genomic epidemiology uses high-throughput sequencing of SARS-CoV-2 genomes to inform local and international transmission events, as well as the diversity of circulating variants. WBE uses wastewater to analyse community spread, as it is known that SARS-CoV-2 is shed through bodily excretions. Since both symptomatic and asymptomatic individuals contribute to wastewater inputs, we hypothesized that the resultant pooled sample of population-wide excreta can provide a more comprehensive picture of SARS-CoV-2 genomic diversity circulating in a community than clinical testing and sequencing alone. In this study, we analysed 91 wastewater samples from 11 states in the USA, where the majority of samples represent Maricopa County, Arizona (USA). With the objective of assessing the viral diversity at a population scale, we undertook a single-nucleotide variant (SNV) analysis on data from 52 samples with >90% SARS-CoV-2 genome coverage of sequence reads, and compared these SNVs with those detected in genomes sequenced from clinical patients. We identified 7973 SNVs, of which 548 were “novel” SNVs that had not yet been identified in the global clinical-derived data as of 17th June 2020 (the day after our last

* **Corresponding author at:** The Biodesign Center for Fundamental and Applied Microbiomics, Arizona State University, 1001 S. McAllister Ave., Tempe, Arizona, AZ 85281, USA

E-mail address: arvind.varsani@asu.edu (A. Varsani).

¹ Authors contributed equally to this work

<https://doi.org/10.1016/j.watres.2021.117710>

Received 25 January 2021; Received in revised form 15 September 2021; Accepted 22 September 2021

Available online 25 September 2021

0043-1354/© 2021 The Author(s).

Published by Elsevier Ltd.

This is an open access article under the CC BY-NC-ND license

(<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

wastewater sampling date). However, between 17th of June 2020 and 20th November 2020, almost half of the novel SNVs have since been detected in clinical-derived data. Using the combination of SNVs present in each sample, we identified the more probable lineages present in that sample and compared them to lineages observed in North America prior to our sampling dates. The wastewater-derived SARS-CoV-2 sequence data indicates there were more lineages circulating across the sampled communities than represented in the clinical-derived data. Principal coordinate analyses identified patterns in population structure based on genetic variation within the sequenced samples, with clear trends associated with increased diversity likely due to a higher number of infected individuals relative to the sampling dates. We demonstrate that genetic correlation analysis combined with SNVs analysis using wastewater sampling can provide a comprehensive snapshot of the SARS-CoV-2 genetic population structure circulating within a community, which might not be observed if relying solely on clinical cases.

1. Introduction

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) is the biggest pandemic since the 1918 H1N1 influenza A virus (Wang et al., 2020; Yan et al., 2020). The SARS-CoV-2 outbreak in humans likely emerged from a zoonotic transmission event(s), and was first recorded in December, 2019, in the City of Wuhan, China (Andersen et al., 2020; Boni et al., 2020; Zhang and Holmes, 2020). According to the Johns Hopkins Coronavirus Resource Center (Dong et al., 2020), there have been >95 million confirmed cases, resulting in more than 2 million deaths globally as of 18th January 2021. SARS-CoV-2 is a positive-sense single-stranded RNA virus in the family *Coronaviridae* (Gorbalenya et al., 2020) that can cause a range of symptoms in infected individuals including complications with breathing, dry cough, fever, and diarrhoea (Wang et al., 2020). However, the majority of individuals show little to no symptoms (Buitrago-Garcia et al., 2020; Byambasuren et al., 2020; Kimball et al., 2020; Syangtan et al., 2020).

Clinical testing of individuals for SARS-CoV-2 is the primary surveillance method for informing public health strategic interventions, and essential for implementing preventive measures, such as quarantine, to mitigate the spread of the virus. The most frequently used approach for clinical testing relies on the detection of genomic elements of SARS-CoV-2 by reverse transcription-quantitative polymerase chain reaction (RT-qPCR) based methods (CDC, 2020a; WHO 2020a)). The clinical analysis is now also being complemented with antibody-based assays (Adams et al., 2020; Becker et al., 2021; Bryant et al., 2020; CDC, 2020b; WHO 2020b)) that provide an indication of current or previous exposure to SARS-CoV-2.

High-throughput sequencing (HTS) technologies are being used to sequence the SARS-CoV-2 genome from a subset of the infected population globally using clinical samples. This has resulted in a large number of published genomes (Elbe and Buckland-Merrett, 2017; Shu and McCauley, 2017), and has provided insight into its origins, spread, and diversity via computational approaches in genomic epidemiology. Screening/testing of a large number of individuals for SARS-CoV-2 can be challenging particularly from a logistics perspective. Furthermore, in most countries it is largely the symptomatic population that is targeted for testing and therefore a large proportion of infected asymptomatic individuals may be missed. Nasopharyngeal swabs and saliva samples have been the principal sample types used for screening, however, SARS-CoV-2 has also been detected in other clinical specimens such as faeces, from both symptomatic and asymptomatic infected individuals (Chen et al., 2020; Jones et al., 2020; Park et al., 2020; Tang et al., 2020; Xing et al., 2020). Moreover, of late, wastewater samples have been utilized as a way to identify several pathogenic human viruses and, not surprisingly, it has gained attention for assessing population-level trends of SARS-CoV-2 infections.

Detection of SARS-CoV-2 in wastewater (untreated and treated) has been a focus of research, with feasibility highlighted in the review by Farkas et al. (2020) and with reported studies from locations including North America (D'Aoust et al., 2021; Nemudryi et al., 2020; Peccia et al., 2020; Wu et al., 2020), Europe (Balboa et al., 2021; Kocamemi et al., 2020; La Rosa et al., 2020; Medema et al., 2020; Randazzo et al., 2020;

Westhaus et al., 2021; Wurtzer et al., 2020), Asia (Kumar et al., 2020; Zhang et al., 2020) and Oceania (Ahmed et al., 2020). These studies used a range of sample concentration and viral RNA recovery approaches followed by RT-qPCR amplification to detect and determine the viral load. Two recent studies have sequenced the SARS-CoV-2 genomes recovered from wastewater (Crits-Christoph et al., 2021; Izquierdo-Lara et al., 2021).

Despite the promising success of these prior studies, it is still unclear how well wastewater-based epidemiology can identify the genetic diversity of SARS-CoV-2 in a given population and how this relates to known viral diversity of clinical cases. This is especially important as new lineages are being discovered. For example, the B.1.351 strain in the United Kingdom that contains single-nucleotide variants (SNVs) of potential biological significance such as N501Y (in the spike protein) (Rambaut et al., 2020b) and K417N, E484K and N501Y in South Africa (Tegally et al., 2021). To investigate the potential of using wastewater to gain insights into variants of SARS-CoV-2 circulating in the population, we used a tiling amplicon-based high-throughput sequencing approach to determine SARS-CoV-2 sequences (spanning the genome) in 91 wastewater samples collected from 11 states in the United States (USA) between 7th April 2020 and 16th June 2020. To further survey the viral diversity circulating within a community and to examine how these relate to sequences from clinical cases, we undertook SNV analysis and beta diversity analyses of SARS-CoV-2 sequences in 52 (>90% coverage) out of the 91 wastewater samples from 10 states. We focused specifically on spatial and temporal trends, and how they compare with clinically-derived data.

2. Material and methods

2.1. Sample collection and transport

Flow- or time-weighted, 24-hr composite samples of untreated wastewater were collected either from the headworks of the wastewater treatment plant, from within the wastewater collection system or at hospital facilities using high frequency automated samplers (Teledyne ISCO, USA) from locations across 11 states in the USA between 7th April 2020 and 16th June 2020 (Table 1, Fig. 1A, Sup Fig. 1). Most samplers had refrigeration capabilities or were supplied with an ice/dry ice blend to keep the interior collection vessel cool. During sample collection, wastewater was thoroughly mixed and transferred to high-density polyethylene sample bottles and placed on ice for transport. The samples were either hand delivered or shipped (next-day/2-day) in insulated shipping containers for subsequent processing and analysis.

2.2. Wastewater sample processing and RNA extraction

Aliquots of 150 ml of each composite wastewater sample were filtered through a 0.45 µm polyethersulfone (PES) filter and then subsequently through a 0.2 µm (PES) filter. The filtrate was then concentrated using the Amicon® Ultra 15 Centrifugal Filter Units (MilliporeSigma, USA) by centrifuging at 4500 rpm for 15 min. For each sample, the process was repeated five times in total using two filter

Table 1

Summary of wastewater sample information. The collection date reflects influent from the previous day. Details of the location including state, city, and region of collection, and Ct value from the RT-qPCR SARS-CoV-2 detection assay targeting the E gene. The SARS-CoV2 genome percentage coverage based on the HTS for each sample is provided.

State	Location ID	Sampling date	Sample ID	Ct value	Mean coverage	Percentage coverage	Total reads
Arizona	G2	7-May-20	122	35.1	21.9801	37.91	8228
Arizona	G2	10-Jun-20	G3	32.2	82.9204	95.7246	30944
Arizona	Guadalupe	6-May-20	110	31.9	139.084	97.8267	51936
Arizona	Guadalupe	10-May-20	136	30.8	249.107	98.6426	93131
Arizona	Guadalupe	12-May-20	147	30.2	682.605	99.0555	254395
Arizona	Guadalupe	16-May-20	177	30.2	800.327	98.9946	298388
Arizona	Guadalupe	19-May-20	179	30.9	780.958	99.0217	291504
Arizona	Guadalupe	21-May-20	203	29.9	1496.09	99.1029	558227
Arizona	Guadalupe	26-May-20	227	30.6	563.257	98.9269	209969
Arizona	Guadalupe	30-May-20	253	28.9	1784.29	99.1097	665406
Arizona	Guadalupe	3-Jun-20	277	30.2	31.7447	71.6733	11859
Arizona	Guadalupe	5-Jun-20	303	30.6	18.0822	65.1061	6766
Arizona	Guadalupe	7-Jun-20	321	30.8	457.993	98.9269	170607
Arizona	Guadalupe	9-Jun-20	341	30.8	1111.99	98.998	414806
Arizona	Guadalupe	11-Jun-20	359	29.5	45.4204	83.8868	16957
Arizona	M1	27-Apr-20	80	32.7	20.8707	43.5666	7802
Arizona	M1	7-May-20	117	34.9	2.24021	7.66054	880
Arizona	M1	26-May-20	225	35.9	13.4329	37.7272	5035
Arizona	Rural	24-Oct-19	R19	NA	10.9956	1.29989	2698
Arizona	Rural	16-May-20	167	35.7	29.7984	54.0537	11099
Arizona	Rural	3-Jun-20	269	34.4	170.102	97.0279	63422
Arizona	Rural	6-Jun-20	305	33.3	87.2427	96.7435	32575
Arizona	Rural	9-Jun-20	338	33	81.784	97.1497	30496
Arizona	Rural	11-Jun-20	349	31.6	81.6799	96.0157	30520
Arizona	TP01	7-Apr-20	4	35	59.1029	66.643	22076
Arizona	TP01	8-Apr-20	3	37	0.646356	1.56054	255
Arizona	TP01	17-Apr-20	57	35	4.45655	15.1958	1667
Arizona	TP01	21-Apr-20	59	33	18.1784	39.5586	6761
Arizona	TP01	29-Apr-20	93	35	11.943	38.2418	4446
Arizona	TP01	12-May-20	137	34.7	47.4554	62.7061	17703
Arizona	TP01	26-May-20	220	35.5	35.8432	64.4122	13421
Arizona	TP01	2-Jun-20	260	33.6	586.011	99.0183	218520
Arizona	TP01	7-Jun-20	322	35.7	39.971	77.0048	14903
Arizona	TP01	9-Jun-20	348	31.5	339.292	98.9066	126569
Arizona	TP02	29-Apr-20	94	35	2.23134	7.12907	844
Arizona	TP02	12-May-20	138	35.8	5.71064	11.9055	2144
Arizona	TP02	30-May-20	247	35.1	52.7047	91.0226	19682
Arizona	TP02	2-Jun-20	261	32.6	106.321	96.0699	39581
Arizona	TP02	5-Jun-20	299	34	84.0252	96.3779	31348
Arizona	TP02	9-Jun-20	344	32.8	258.612	99.1165	96441
Arizona	TP03	6-Jun-20	312	34.5	130.712	97.2344	48711
Arizona	TP03	7-Jun-20	323	35.4	151.054	98.3514	56337
Arizona	TP04	28-May-20	274	34.5	34.992	71.6699	13061
Arizona	TP04	4-Jun-20	288	33	110.474	96.2053	41202
Arizona	TP04	5-Jun-20	129	32.7	31.8066	72.3368	11897
Arizona	TP04	6-Jun-20	314	34.7	191.419	98.8829	71379
Arizona	TP04	8-Jun-20	336	32.8	220.449	98.9371	82296
Arizona	TP05	25-Apr-20	69	31.2	15.223	41.1699	5678
Arizona	TP05	7-May-20	118	32.1	22.2285	50.7803	8291
Arizona	TP05	19-May-20	181	35.8	38.4298	59.3514	14304
Arizona	TP05	7-Jun-20	326	35.6	27.9763	66.1792	10443
Arizona	TP05	9-Jun-20	347	26.8	3735.92	99.1097	1510084
Arizona	TP05	11-Jun-20	358	31.5	37.94	75.9453	14211
Arizona	TP06	26-Apr-20	78	34.9	2.9937	5.98152	1127
Arizona	TP06	21-May-20	198	34.9	17.187	57.3034	6445
Arizona	TP06	28-May-20	234	34.7	784.976	98.998	292585
Arizona	TP06	3-Jun-20	271	33.3	61.7264	93.4159	23022
Arizona	TP06	5-Jun-20	296	32.8	92.836	97.3901	34617
Arizona	TP06	7-Jun-20	318	34.6	40.5103	90.8805	15096
Arizona	TP06	9-Jun-20	339	32.6	33.4383	86.1074	12474
Arizona	TP06	11-Jun-20	351	30.6	20.0344	65.7696	7472
Colorado	CO1	20-May-20	Jac_51	32.1	85.5393	93.1282	31953
Colorado	CO1	28-May-20	Jac_103	34	91.4798	96.124	34120
Georgia	GA1	14-May-20	Jac_33	29	68.8532	94.4078	25686
Idaho	ID1	18-May-20	Jac_56	34.7	88.5662	91.114	33005
Idaho	ID1	25-May-20	Jac_87	35.3	113.577	94.4416	42320
Illinois	IL1	19-May-20	Jac_45	33.3	79.0705	96.9331	29490
Illinois	IL1	1-Jun-20	Jac_106	33.1	54.4429	85.8332	20365
Illinois	IL2	7-May-20	Jac_12	33	71.8524	90.6875	26744
Illinois	IL2	1-Jun-20	Jac_127	31.9	77.5081	87.7357	28850
Kansas	KA1	20-May-20	Jac_58	33.2	91.4503	91.9265	34117
Kansas	KA1	27-May-20	Jac_96	31.7	364.619	98.9845	135932

(continued on next page)

units, and subsequently the concentrates were pooled per sample (from the two filter units). For each sample, a 200 μ l aliquot was used to extract total RNA using the RNeasy mini kit (Qiagen, USA).

2.3. SARS-CoV-2 RT-qPCR detection and high throughput sequencing of SARS-CoV-2 genome sequences

To determine the presence of SARS-CoV-2 in wastewater samples, the extracted RNA was used in a reverse transcription-quantitative PCR (RT-qPCR) assay targeting the E gene, as designed and validated by Corman et al. (2020) and cited by the WHO (WHO, 2020a). This probe-based assay was performed as per the specifications outlined in Corman et al. (2020) using the SuperScript III Platinum One-Step qRT-PCR Kit (Invitrogen, USA). This assay was validated and used by Holland et al. (2020) on SARS-CoV-2 clinical samples.

A total of 91 samples from 11 states in the USA (Fig. 1) were collected between 7th April 2020 and 16th June 2020 that tested positive, and one negative control sample collected in October 2019 in Tempe, Arizona (Table 1) were selected for sample processing and high-throughput SARS-CoV-2 amplicon sequencing. The SARS-CoV-2 RT-qPCR assay Ct values ranged from 26.8 to 36 for the 91 samples (Fig. 1). Total RNA (11 μ l) from each sample was used to generate cDNA using the Superscript® IV First-Strand Synthesis System (Thermo Fisher, USA). The manufacturer's protocol was followed, with one modification, the reverse transcription incubation step (50°C) was increased from 10 to 50 min. 10 μ l of cDNA from each sample was used to generate Illumina sequencing libraries (92 libraries in total) with the Swift Nomalase® Amplicon SARS CoV-2 Panel (SNAP) and these were subsequently normalized, pooled and sequenced at Psomagen (USA) on an Illumina HiSeq 2500 sequencer (2 \times 100 paired-end option on 1 lane in rapid mode).

2.4. Bioinformatics pipeline and analyses

The Illumina raw read sequences were aligned to the reference genome of SARS-CoV-2 (MN908947; RefSeq ID NC_045512.2) using the Burrows-Wheeler Alignment tool (BWA) MEM (Li and Durbin, 2009). The primers used for the tiling PCR-based amplification step were soft-clipped using iVar trim tool (Grubaugh et al., 2019) which also removed reads <30nts and reads that started outside of the primer region. Trimming with a sliding window of 4 for a minimum PHRED quality of 20 was performed as default by iVar. Primers that may have mismatches with the reference sequence were also evaluated and reads from those amplicons with varying primer binding efficiency were also removed as described by Grubaugh et al. (2019). The genome coverage (minimum quality of 20 and 10 \times coverage) and mean depth was

calculated for all samples. For the 52 samples with >90% genome coverage, variant calling was performed using iVar (Grubaugh et al., 2019) with minimum base quality of 20 and 20 \times coverage with no cut-off frequency since we have population-level sequence data. This approach was used because, unlike the case with a clinical sample from a single infected individual, wastewater contains material from a population that inhabits a particular region and therefore represents a collection of SARS-CoV-2 variants actively shed by infected individuals within the population. From the variants that were identified, only those with a p-value <0.05 in the Fisher's exact test implemented in iVar (tests if SNV frequency is higher than the mean error rate at the specific position) were maintained. Suggested masked sites due to biases shown by phylogenetic analysis or sequencing technology (De Maio et al., 2020) as of September 2020 were removed for downstream analyses. To identify the possible novel SNVs, the SNVs determined from the 52 wastewater samples with SARS-CoV-2 genome read coverage >90% were searched in all clinical data available in GISAID (Elbe and Buckland-Merrett, 2017; Shu and McCauley, 2017) at two time points (17th June 2020 and 20th November 2020). The 17th June 2020 download dataset includes clinical sample dates of 24th December 2019 to 11th June 2020, whereas 20th November 2020 download dataset includes those from 24th December 2019 to 16th November 2020. Variants that were not present in the GISAID deposited SARS-CoV-2 genomes were considered novel, however, to be more stringent, variants that were only present in one of the wastewater samples were removed from further analyses.

2.5. Support for lineages assigned by Phylogenetic Assignment of Named Global Outbreak Lineages (PANGOLIN)

Each environmental sample was compared against the SARS-CoV-2 genomes available in GISAID (Elbe and Buckland-Merrett, 2017; Shu and McCauley, 2017), an open-access genomic database, to collect a set of clinical genomes whose mutations were supported by the SNVs identified above. To reduce false positives, basal genomes, defined as those with 3 or fewer mutations relative to the reference (MN908947) were excluded. The set of genomes supported by each environmental sample SNV profile were grouped by lineages assigned by Phylogenetic Assignment of Named Global Outbreak Lineages (PANGOLIN) (Rambaut et al., 2020a) and lineages with fewer than 3 genomes were excluded to avoid any misannotations resulting in false positives. PANGOLIN is an online platform that assigns lineages to sequences (Rambaut et al., 2020a) and is updated as new metadata are submitted to GISAID. For each group of genomes (grouped per PANGOLIN), we then looked to see whether any genome was from North America and, if

Table 1 (continued)

State	Location ID	Sampling date	Sample ID	Ct value	Mean coverage	Percentage coverage	Total reads
Kentucky	S1	23-Apr-20	Lou_2	33.8	31.4104	70.7017	11723
Kentucky	S2	9-Jun-20	Lou_40	33.8	352.012	98.734	131165
Kentucky	S3	21-May-20	Lou_15	35.3	11.7138	36.1193	4379
Kentucky	S3	28-May-20	Lou_23	35.5	9.75725	33.6448	3640
Kentucky	S3	9-Jun-20	Lou_39	34.5	68.0629	87.6883	25339
Kentucky	S4	9-Jun-20	Lou_43	34.6	58.5395	92.2413	21876
Kentucky	S5	14-May-20	Lou_6	33.2	296.939	99.1233	110803
Kentucky	S5	9-Jun-20	Lou_38	31.4	393.77	99.0928	146800
Kentucky	S6	9-Jun-20	Lou_42	33.7	57.09	92.0856	21266
Kentucky	S7	23-Apr-20	Lou_3	33.2	63.1731	84.0764	23501
Kentucky	S8	21-May-20	Lou_13	34.8	148.323	98.5546	55410
Kentucky	S9	23-Apr-20	Lou_1	29.4	206.044	98.7577	76835
Massachusetts	MA1	27-May-20	Jac_89	32.8	89.2101	97.6236	33207
New Jersey	NJ1	3-May-20	Jac_04	31.2	62.1934	88.3518	23228
New Jersey	NJ1	11-May-20	Jac_30	32.6	1768.26	99.0759	658845
New Mexico	NM1	6-May-20	Jac_09	30.8	14.5232	42.8015	5435
New Mexico	NM1	13-May-20	Jac_31	33	127.887	98.1686	47610
New Mexico	NM1	21-May-20	Jac_69	34.3	139.456	94.8681	52042
New Mexico	NM1	27-May-20	Jac_90	34.1	223.602	98.321	83229
Oregon	OR1	27-May-20	Jac_92	34.7	9.50418	27.8291	3568

so, recorded the time between the genome's sampling date and the collection date of the environmental sample. Note that the set of genomes which we summarize as certain SARS-CoV-2 lineages assigned by PANGOLIN may be different for each environmental sample, and thus the time between clinical and environmental sampling dates depends on the particular SNV profile of the environmental sample. Given that linkage of SNVs is not possible via short read sequencing, support for mutation profiles observed in clinical genomes (and, correspondingly, PANGOLIN) does not guarantee that the lineages were present in the environmental sample.

2.6. Sample-based SARS-CoV-2 sequence distance calculation and ordination analysis

The 'genotype' of each sample was represented in a four-column matrix. In this matrix, each row corresponds to a position in the reference genome, and the value at each column is the frequency of occurrences for each nucleotide (A, C, G and T). At each genomic position, the Yue & Clayton measure of dissimilarity index (Yue and Clayton, 2005) on the nucleotide frequency of the compared samples was calculated. If

the nucleotide frequency at a position of a sample cannot be calculated due to zero depth, the Yue & Clayton measure of dissimilarity index at this position between this sample and any other sample compared is assumed to be zero. The sum of the Yue & Clayton dissimilarity (Yue and Clayton, 2005) of all genomic positions was used as a measure of distance between samples. The distance matrix was constructed by calculating pairwise distances of all samples and was subsequently used for principal coordinates analysis (PCoA) (Gower, 1966).

3. Results and discussion

3.1. Sample collection, processing, amplification and high-throughput sequencing of SARS-CoV-2 from wastewater samples

Sixty of our 91 samples (66%) were collected in Arizona (9 locations located in Maricopa County, Arizona Sup Fig. 1), 12 (13%) were collected from 9 locations in Louisville, Kentucky (Sup Fig. 1), and 19 (21%) were collected from other states, see Table 1 and Fig. 1 for details. The tiling PCR amplification enrichment process for the SARS-CoV-2 genome generated 341 amplicons covering ~99% of the genome

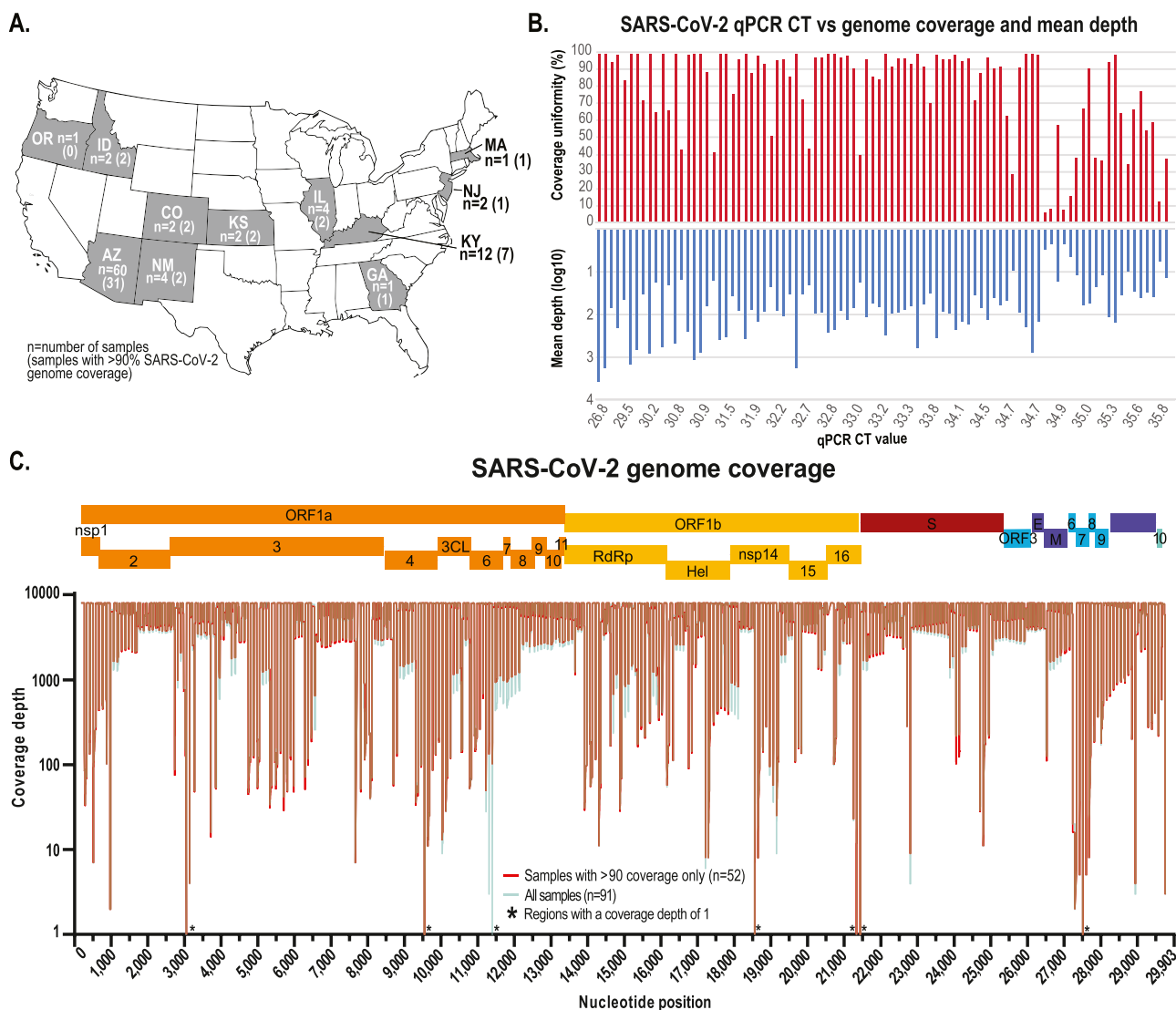


Fig. 1. A. Map of the United States of America with states where wastewater samples were collected for this study highlighted in grey. B. SARS-CoV-2 RT-qPCR Ct detection value for each sample and the corresponding SARS-CoV-2 genome coverage uniformity from the tiling amplicon-based HTS. C. SARS-CoV-2 genome coverage of the high-throughput sequencing of all the wastewater samples (cyan) and those with >90% coverage (red). * indicates that these sites have a coverage depth of 1.

albeit missing the 200 nts of 5' end and 162 nts from 3' end. The genome coverage calculated for all samples ranged between ~1.3% and ~99%. 52 of the 91 RT-qPCR positive samples showed >90% coverage (minimum quality of 20 and >10 reads per position) (Table 1). We note that there is no clear correlation between coverage and Ct values obtained using the RT-qPCR assay as some samples with lower Ct values appear to have low coverage and depth (Fig. 1). Nonetheless, the mean depth is lower with higher Ct values (Fig. 1). This has been shown in other wastewater-derived viral sequencing projects using an Illumina sequencing platforms via an amplification process (Izquierdo-Lara et al., 2021) and a capture approach (Crits-Christoph et al., 2021). This lack of correlation is not unexpected given the nature of wastewater, where dilution and degradation play a significant role, thereby this likely results in samples with differing levels of genomic RNA degradation. Furthermore, since the RT-qPCR assay only targets a specific small region of the genome, the Ct-value based quantification vary. Additionally, it is important to highlight that there are several variables attributed to the handling and transport process of the wastewater samples prior to concentration and RNA extraction. We acknowledge that we did not measure the recovery efficacy of SARS-CoV-2 in our extractions from wastewater (via a spiked surrogate). Recovery efficacy can help guide whether the majority of the SARS-CoV-2 sequence population in the sample has been "captured" for downstream analysis. To counter this, we use a conservative approach of only reporting and analysing the SNVs from the samples (n=52) for which we have >90% genome coverage. Certainly, there are SNVs in our samples that are likely not captured in our sequencing effort due to 1) preferential amplification of genomic regions due to sample quality; 2) SNVs in the non-coding regions not covered by the tiling amplicon approach; 3) sequence repeat regions that would yield low quality sequencing. Despite this, we were still able to identify 548 novel SARS-CoV-2 SNVs and therefore this conservative approach highlights the number of SNVs that are detectable in the excreted viral population with sequence coverage of >90%.

3.2. Wastewater-derived SARS-CoV-2 sequence analyses

A total of 7973 SNVs were detected for the 52 analysed samples with >90% genome coverage after quality control steps from which the number of detected SNVs per sample ranged from 24 to 793 (Sup. Table 1, Sup. Table 2, Fig. 2A). As expected, mean depth is correlated

with the number of SNVs detected in each sample (Fig. 2B), the regression analysis indicates the trend.

To determine unique variants within the 52 wastewater-derived SARS-CoV-2 sequences, SNVs counted in more than one sample at each site were removed and this resulted in 5680 unique SNVs identified across the genome. Of these, 4372 are non-synonymous and 1308 are synonymous substitutions (Sup. Table 2). Additionally, 246 result in nonsense mutations and 64 are in non-coding regions. We highlight that SNV A23403G responsible for the spike protein substitution D614G that is frequently seen in clinical data, although it has not thus far been shown to be under strong positive selection (Volz et al., 2021), was present in all 52 wastewater-derived SARS-CoV-2 sequences. From one sample (sample #147, Tempe, Arizona), a new variant A23403T was also identified that results in a D614V substitution in the spike protein, but at very low frequency (Sup. Table 1).

3.3. Comparative analysis of SARS-CoV-2 SNVs in clinical and wastewater-derived samples during the collection period

The wastewater-derived SARS-CoV-2 SNVs were compared with substitutions that have been detected in clinical-derived sequences available in public databases. The first aim was to identify possible "novel" SNVs present in the analysed wastewater samples that had not yet been identified in any of the sequences available in GISAID (Elbe and Buckland-Merrett, 2017; Shu and McCauley, 2017) from clinical samples globally. To accomplish this, we initially undertook an analysis to identify all the detected SNVs in the clinical data available from GISAID up until the 17th June 2020 (subsequent to the last day of wastewater sampling in this study - 16th June 2020) which consisted of 45,836 SARS-CoV-2 genome sequences. A total of 548 novel SNVs (Sup. Table 3) were identified in the 52 wastewater samples collectively, of these 469 were non-synonymous (not including nonsense mutations) and 79 were synonymous substitutions (Fig. 3). Since we evaluated all variants regardless of frequency, some locations (as expected) had more than one possible variant and are illustrated in Fig. 3 and outlined in Sup Table 1. These 548 SNVs are distributed along the SARS-CoV-2 genome with three of those located in non-coding regions. The vast majority of "novel" SNVs were detected in up to 8 of the wastewater samples analysed. The exceptions are four non-synonymous mutations, three on the ORF1ab and one in the N gene that are present in >8 samples (Fig. 3 and Sup Table 1). It is important to highlight that not all the novel SNVs may

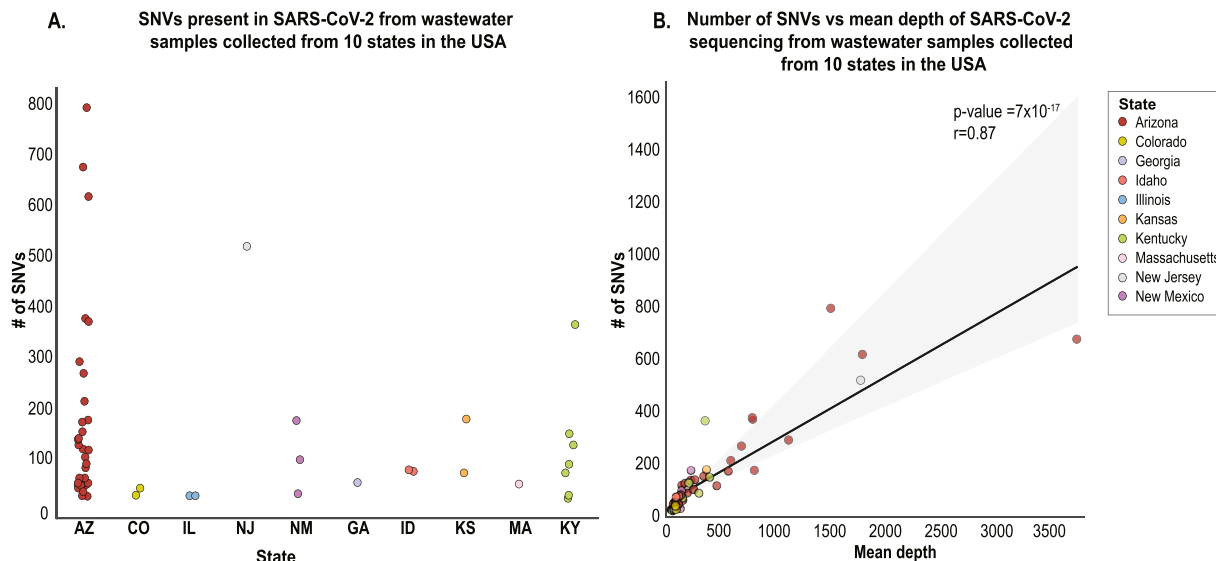


Fig. 2. A. Number of single nucleotide variants (SNV) per sample across 10 states (each state is represented by a different colour). B. Regression analysis, with 95% confidence interval, of the number of wastewater-derived SARS-CoV-2 SNVs detected versus the mean depth for each of the 52 samples with >90% coverage that were analysed. The colour code indicates the states in which the samples were collected.

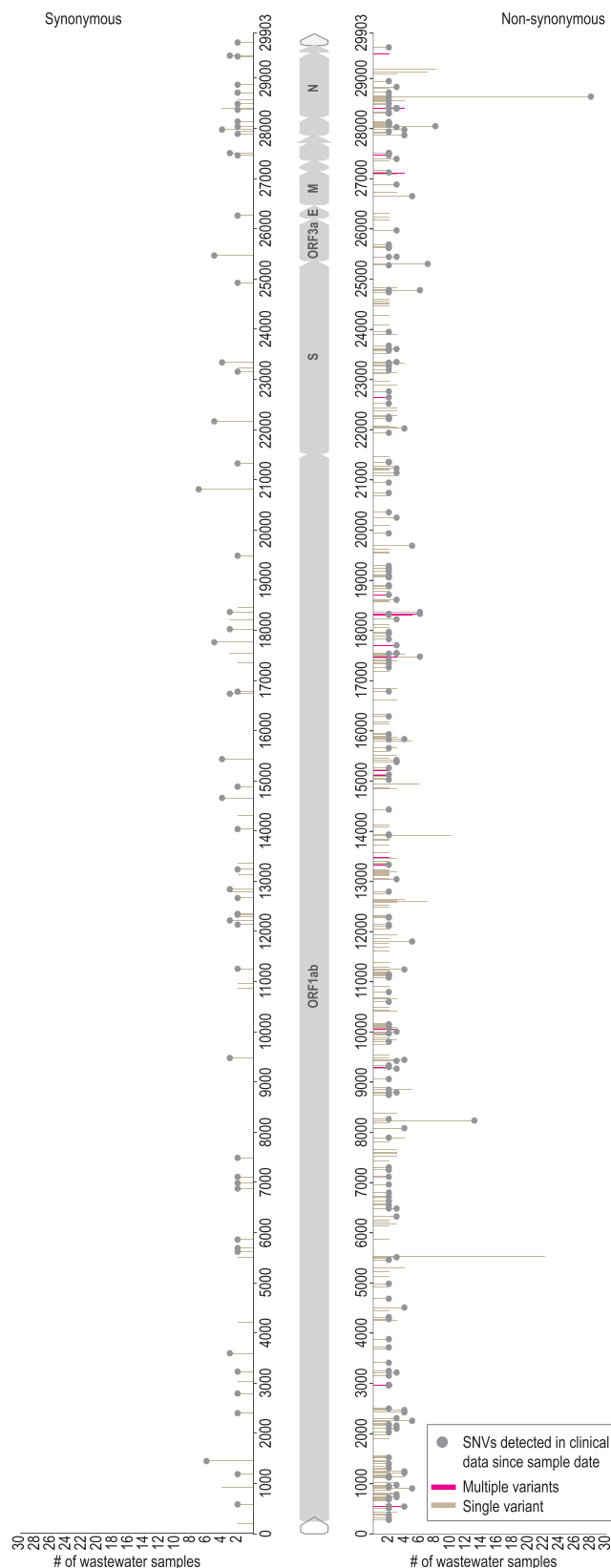


Fig. 3. Novel SARS-CoV-2 SNVs (i.e. not yet detected in clinical-derived samples as of 17th June 2020) identified in the 52 wastewater samples analysed. On the y-axis are the number of samples containing the SNV and on the x-axis is the relative position of SNV in the SARS-CoV-2 genome. Positions with multiple variants are marked in red and those marked with grey circles represent the SNVs that have been detected up until 20th November 2020 in clinical samples.

be associated with a viral lineage circulating in infected individuals, but it is highly likely that a portion of those SNVs are associated with viral genomes that have not yet been sampled and/or deposited in public sequence databases.

3.4. Identification of SARS-CoV-2 SNVs in wastewater samples in clinical-derived samples post-collection period

To determine how many SNVs have been identified post wastewater sample collection (16th June 2020), a second SNV comparison was performed with all the clinical-derived sequence data available as of 20th November 2020 (203,741 SARS-CoV-2 genomes available at GISAID). Based on the analysis of samples during the collection period, SNVs that were not detected in the clinical-derived sequence data up until 17th June 2020 were considered as novel SNVs. From the 548 SNVs considered as novel from the wastewater-derived samples, 263 SNVs were subsequently identified in clinical-derived samples in the period of 17th June - 20th November 2020 (Sup Table 1, Fig. 3). From those, 126 (~47%) SNVs were identified in the USA which provides good support that novel identified SNVs in wastewater samples are associated with circulating SARS-CoV-2 lineages in infected individuals (Sup. Table 3). However, none of those 126 SNVs detected in the wastewater samples were detected in clinical-derived sequence data from the same USA states. Most of the novel SNVs we report were identified in wastewater derived SARS-CoV-2 sequences from Arizona, however, only 54 patient-derived genomes had been submitted to GISAID between 17th June 2020 and 20th November 2020 and those had collection dates ranging from 12th March 2020 to 27th June 2020. Before 17th June 2020 only 86 patient-derived genomes from Arizona were available in GISAID with collection dates ranging from 22nd January 2020 to 2nd April 2020. The remainder of the novel SNVs (137/263) were identified in clinical-derived samples from a variety of countries (Sup. Table 3) suggesting that sequences containing those SNVs were likely circulating in the population in the USA but had not been sampled in a clinical setting or made available in public databases. Two hundred eighty-five SNVs identified in the wastewater-derived samples with the last sampling date of 16th June 2020 had not been identified in clinical-derived SARS-CoV-2 sequences deposited between then and 20th November 2020.

The results show that a proportion of those novel SNVs are present in lineages circulating in the USA but we acknowledge that not all of these novel SNVs are necessarily fixed in SARS-CoV-2 lineages that are actively being transmitted nor is it possible to determine if any of these SNVs are linked within lineages. Nonetheless, the identification of the novel SNVs clearly demonstrates the relevance of wastewater-derived SARS-CoV-2 sequence analysis which can provide valuable information on SNVs that have not yet been captured using clinical-derived approaches. Wastewater-derived sequence analysis provides information at a population scale and can allow for rapid detection of clinically relevant / important SNVs. Our SNV analysis shows that there is no particular region of the genome that is a SNV hotspot and this mirrors what is observed in clinical-derived samples (see updates on SNVs analysis of clinical data at <https://nextstrain.org/ncov/global>).

3.5. Determination of putative lineages of SARS-CoV-2 in wastewater-derived sequences

Given that wastewater harbours a collective population of SARS-CoV-2 and therefore likely many variants, it is not ideal to determine consensus sequences and consensus sequences-based phylogeny. Therefore, our first approach was to evaluate which clades in the global phylogeny of clinical-derived sequences are supported by the SNVs present in each sample based on the SARS-CoV-2 lineages assigned by PANGOLIN (Rambaut et al., 2020a). The represented SARS-CoV-2 lineages for each wastewater sample that are supported are shown in Fig. 4. We determined the time frames for which these lineages were first detected in North American clinical-derived sequences relative to the

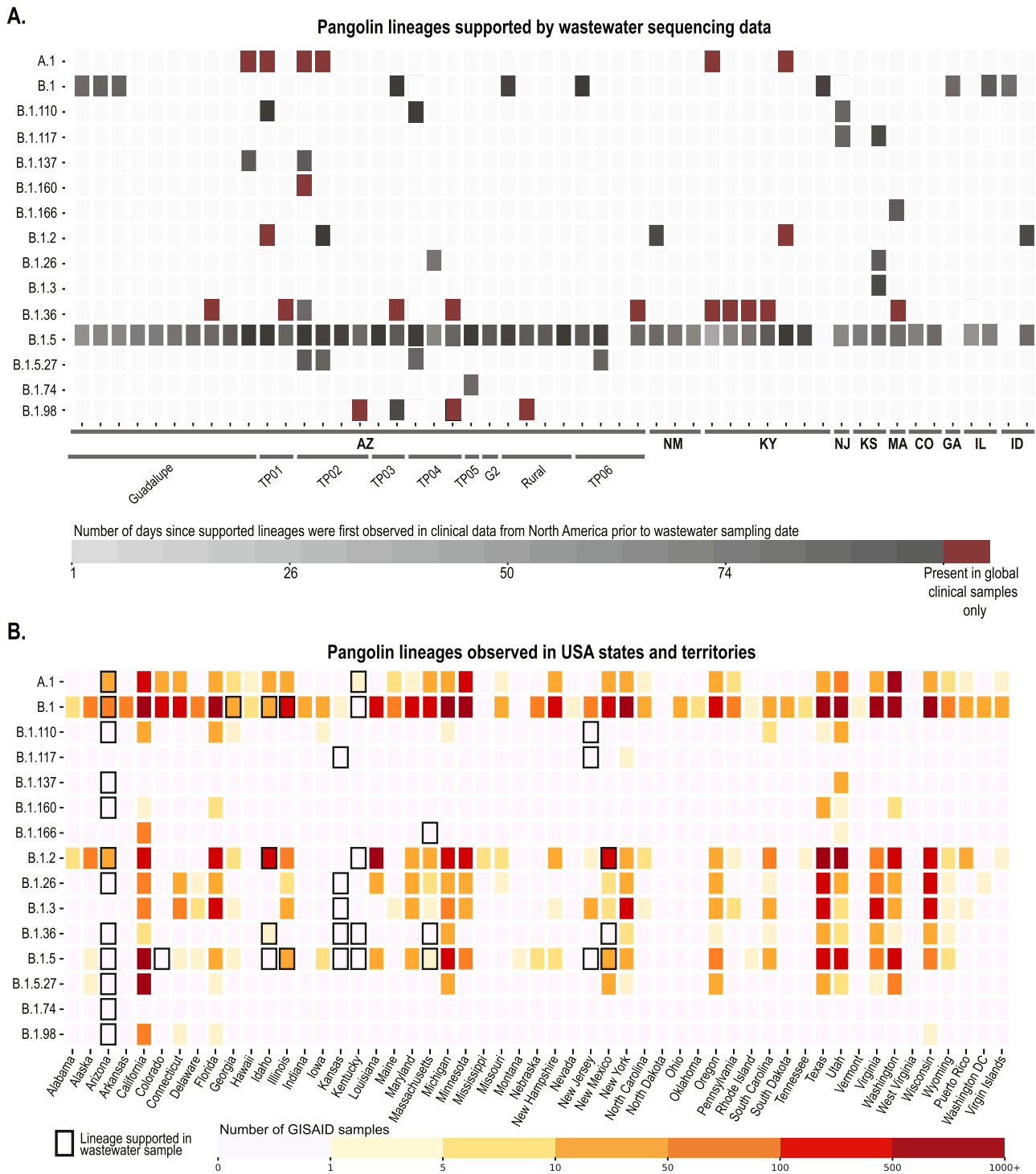


Fig. 4. Publicly available genomes from clinically derived data deposited in GISAID, grouped by PANGOLIN, whose mutations were consistent with those observed in wastewater samples. **A.** Heatmap showing the number of days between sample collection and when supported lineages were first observed in clinical data. Each wastewater sample (52 samples across 10 states) contained support for different clinical samples which are grouped here by PANGOLIN, some of which have only been observed outside North America (indicated as “global only”). **B.** Clinical genomes reported in USA states and territories which were assigned to PANGOLIN supported by at least one environmental sample. Black borders indicate lineages supported in environmental samples from the respective location.

date each wastewater sample was collected (Fig. 4A).

We also undertook a comprehensive analysis of all the lineages detected in each state in the USA up to November 2020 that were supported by at least one environmental sample, this included the number of clinical-derived SARS-CoV-2 genomes sequenced in each lineage (Fig. 4B). This approach helps to determine whether wastewater-based surveillance for SARS-CoV-2 can provide valuable insights on putative

circulating lineages in the wastewater contributing population. Although there are several limitations to the analysis of wastewater-derived SARS-CoV-2 sequences, our analysis of SNV-based supported lineages revealed some interesting findings. From the 52 analysed wastewater samples, 15 SARS-CoV-2 lineages assigned by PANGOLIN (Rambaut et al., 2020a) were supported, with lineage B.1.5 being the most prominent for the wastewater-derived sequences. The B.1.5

lineage has been identified in clinical samples in 27 USA states. Our wastewater-derived sequence data suggests that B.1.5 may also be present in 6 additional states in the USA (Arizona, Colorado, Idaho, Kansas, Kentucky and New Jersey). In 17 of the 52 wastewater samples, there were up to two supported SARS-CoV-2 lineages that had not been detected in North American clinical samples, during the period of our wastewater collection, as of 17th June 2020 (Fig. 4). These 17 samples were from the states of Arizona, Kentucky and Massachusetts (Fig. 4B). In wastewater-derived sequences from Arizona, which represents the greatest proportion of samples, the observed circulating lineages based on clinical-derived sequences are well represented (Ladner et al., 2020),

with an additional nine possible circulating lineages identified.

Although wastewater-based SARS-CoV-2 sequence analysis does not provide the same level of genome confidence (and thus lineage assignment) as those from clinical samples, the wastewater-derived data can be used to identify possible circulating lineages and assess the diversity of SARS-CoV-2. We would like to emphasize that despite us identifying supported lineages based on SNVs analysis, without verification of full genomes using long read sequencing technologies it is not possible to confirm all the specific lineages present in the wastewater. Nevertheless, it is apparent that valuable population-level variant information on SARS-CoV-2 can be gleaned from wastewater sampling, including

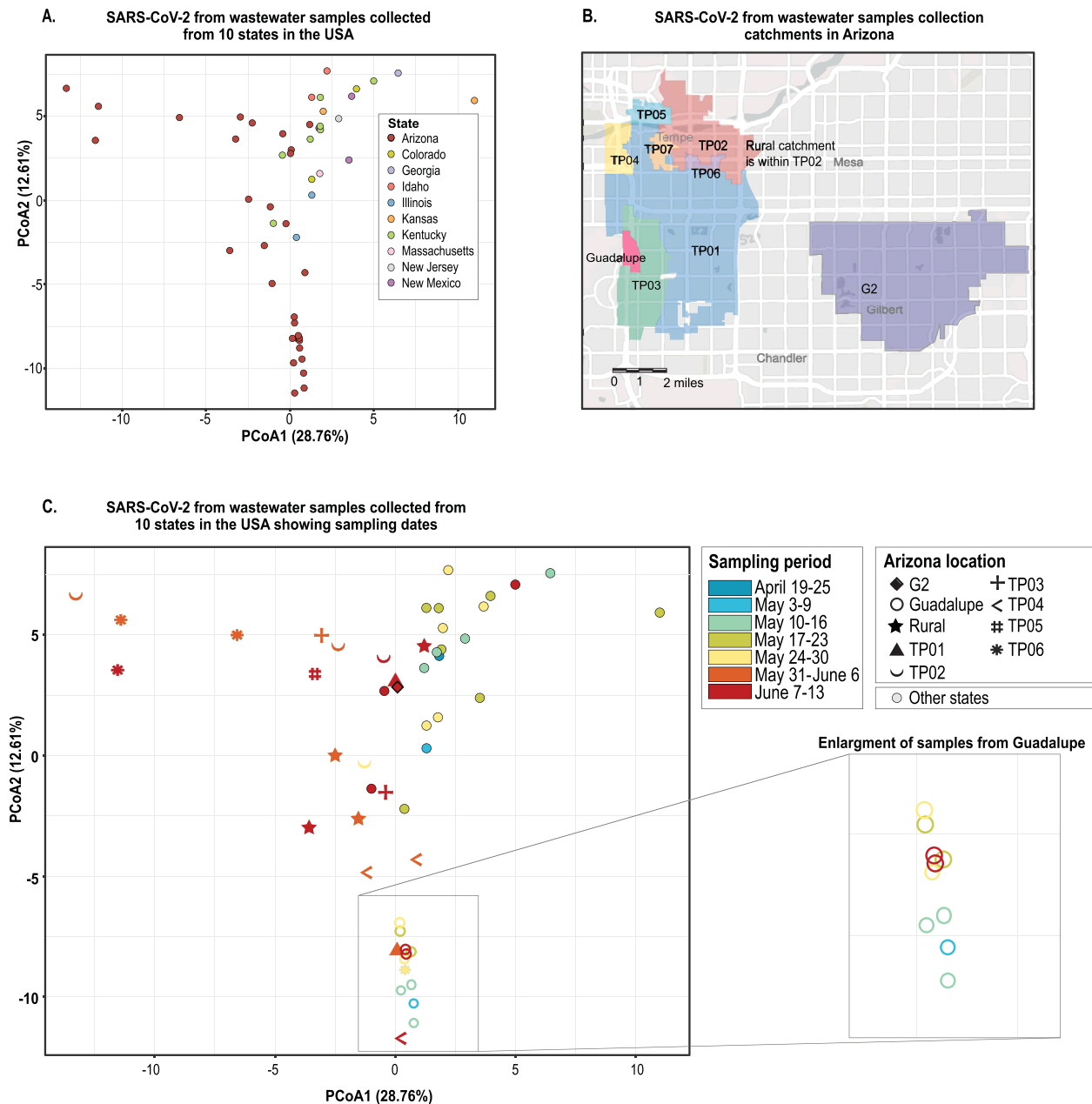


Fig. 5. Principal coordinate analysis (PCoA) of SARS-CoV-2 sequence data derived from wastewater samples. **A.** Distribution of sequences from samples collected in ten states (each represented by a different colour) in the USA showing pairwise distance based on genomic composition between viral populations present in each sample. **B.** Sampling catchments in Tempe, Guadalupe and Gilbert, Arizona. **C.** Spatial and temporal (shown by the colour gradient) representation of samples taken from the sample locations across ten USA states, focusing on Arizona, between April-June 2020 with pairwise distance based on genomic composition between viral populations present in each sample. Enlargement of samples collected from Guadalupe shown as example of location site where the viral population did not appear to diversify greatly over time.

significant sequence data that are potentially missed in clinical-derived sequence data where genomes are sequenced from predominantly infected individuals who might represent a small percentage of those shedding virus in a community.

3.6. Principal coordinates analysis (PCoA) analysis of nucleotide frequencies

In Fig. 5A, we show our PCoA analysis results using nucleotide frequencies to evaluate the viral population diversity within and between samples. SARS-CoV-2 sequences in the samples from the ten states were overall highly diverse, and those with two or more samples from the same state tend to cluster closer together (Fig. 5A and C). The main exceptions are those from Kansas (20th May 2020 and 27th May 2020) and Colorado (20th May 2020 and 28th May 2020) that do not cluster together, both were collected a week apart, and the locations have an estimated human population size of ~25,900 and ~8,300, respectively. Additionally, the Arizona wastewater SARS-CoV-2 sequences are broadly distributed in the PCoA plot which is likely a consequence of the large number of samples collected over a three-month period across several sites within Maricopa County, Arizona (Tempe sites, Guadalupe and Gilbert) (Fig. 5B and C). In comparison to those in the Arizona wastewater samples, the SARS-CoV-2 sequences in samples from Louisville (Kentucky) are much more tightly clustered in the PCoA plot despite sampling from several locations in the city over a two-month period (Fig. 5A). Despite the large number of samples collected in Arizona compared to Kentucky, and the other states, if seven individual samples were to be randomly picked from each location over the same period as those from Kentucky the SARS-CoV-2 genetic distance between them would still be apparently higher for Arizona. We hypothesize that one contributing factor to the differences in viral diversity present in these two areas i.e. Maricopa County Arizona and Louisville (Kentucky), is that, Tempe (the region where the majority of the samples were collected) is home to one of the largest universities in the USA, Maricopa County is the 4th most populous county in the USA with ~4.4 million inhabitants (Maricopa County 2020) and a major travel hub with an international airport.

The highest number of samples collected within a state both temporally and spatially for this study was in Arizona. In Arizona, we note that the wastewater-derived SARS-CoV-2 sequences in samples from the same locations do not necessarily cluster together in the PCoA plot (Fig. 5A and C). Nonetheless, there are clear shifts in the SARS-CoV-2 sequence variants in each location over time (Fig. 5B and C). This is most evident for the Town of Guadalupe (Arizona) given the sampling effort here, where the SARS-CoV-2 sequences in the samples collected in early May 2020 cluster with lower distance but we can see a clear shift in the viral population starting late May 2020 through to early June (Fig. 5C) which coincides with stay at home lockdown being lifted on 15th May 2020. It is important to highlight that the Town of Guadalupe (Arizona) has a small resident community (~6,500) from where wastewater was collected. Moreover, SARS-CoV-2 sequences in the samples from the same location at closer timepoints are often more likely to be similar, yet there are exceptions such as the samples from site TP04 (Tempe, Arizona) that have no resident population (Fig. 5B and C). The shift in SARS-CoV-2 sequence diversity in locations such as TP04 (Tempe, Arizona) over time may be due to new infections given the transient population.

Increases in SARS-CoV-2 viral RNA in wastewater have been correlated to an increase in the number of cases locally (Medema et al., 2020). Observing a shift in the SARS-CoV-2 population diversity through wastewater analysis with time provides insights into corresponding dynamics of increased infection in the community. For example, in Tempe, the number of recorded cases nearly doubled in June 2020. When analysing wastewater-derived SARS-CoV-2 sequence data and correlating it with human dynamics, business districts in the cities will certainly see the activity of transient community members and this will

likely reflect in sequence diversity data.

4. Conclusion

The SARS-CoV-2 pandemic response has relied mostly on clinical-based epidemiology as surveillance for informed response to mitigate viral transmission. However, there are certain limitations to clinical-derived epidemiology such as the number of patient samples that can be analysed based on resources, as well as a bias towards sampling predominantly symptomatic patients. Wastewater-based analyses has been shown to be a useful approach for monitoring of genomic levels of SARS-CoV-2 and community-level surveillance. Further, HTS of SARS-CoV-2 in wastewater samples could provide a population-level analysis of circulating lineages and complement surveillance based on clinical-derived sequences.

In this study, we analyse HTS data of wastewater-derived SARS-CoV-2 sequences to determine SNVs, putative circulating lineages and population structure at a spatial and temporal scale. We were able to recover near full-length genome coverage from ~55% of the analysed samples which demonstrates that wastewater can provide useful genomic data for epidemiology despite high level of variability of handling and processing of samples, as well as viral RNA degradation. In addition, by identifying SNVs in SARS-CoV-2 sequences from each wastewater sample, we were able to determine likely PANGOLIN lineages, some of which were not known to be circulating in the USA as of 20th November 2020. In conjunction with diversity analyses using distance matrices, we show trends in viral populations which can help monitor the shifts in the SARS-CoV-2 sequences within regions.

This study supports the use of wastewater sampling as a tool suitable for analysing the genomics of ongoing outbreaks of infectious diseases, such as SARS-CoV-2. As demonstrated here, HTS of RNA from wastewater can provide novel information on SNVs and lineages which, when coupled with that derived from clinical data, can help identify new emerging variants/lineages of clinical importance within a population. The study results indicating a shift in the SARS-CoV-2 sequence variation in wastewater from each location over time shows the ongoing need for such approaches. As a collective, the approaches we have outlined in this study can be used within a public health setting as an early warning tool to inform infectious disease mitigation measures, especially in situations where obtaining clinical-derived sequences is difficult.

Sequence data

Sequences are deposited in NCBI's SRA under the project number PRJNA662596; SRA # SRR12618464 - SRR12618554 and SRR13289969.

Declaration of Competing Interest

R.U.H. and E. M. D. are cofounders of AquaVitas, LLC, 9260 E. Raintree, Ste 130, Scottsdale, AZ 85260, USA, an ASU startup company providing commercial services in wastewater-based epidemiology. The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

The research reported in this publication was supported by the National Library of Medicine of the National Institutes of Health under Award Number U01LM013129 to RUH, MS and AV. The work of XJ was supported by the Intramural Research Program of the National Library of Medicine, National Institutes of Health. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. The work in Louisville, KY, was supported in part by grants from the James Graham Brown

Foundation and the Owsley Brown II Family Foundation. The authors would like to thank William Mancini (Enterprise GIS & Data Analyst, City of Tempe, USA) for providing the base map for Tempe.

Supplementary materials

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.watres.2021.117710.

References

- Adams, E.R., Ainsworth, M., Anand, R., Andersson, M.I., Auckland, K., Baillie, J.K., Barnes, E., Beer, S., Bell, J.I., Berry, T., Bibi, S., Carroll, M., Chinnakannan, S.K., Clutterbuck, E., Cornall, R.J., Crook, D.W., de Silva, T., Dejnirattisai, W., Dingle, K. E., Dold, C., Espinosa, A., Eyre, D.W., Farmer, H., Fernandez Mendoza, M., Georgiou, D., Hoosdally, S.J., Hunter, A., Jefferey, K., Kelly, D.F., Klenerman, P., Knight, J., Knowles, C., Kwok, A.J., Leuschner, U., Levin, R., Liu, C., Lopez-Camacho, C., Martinez, J., Matthews, P.C., McGivern, H., Mentzer, A.J., Milton, J., Mongkolsapaya, J., Moore, S.C., Oliveira, M.S., Pereira, F., Perez, E., Peto, T., Ploeg, R.J., Pollard, A., Prince, T., Roberts, D.J., Rudkin, J.K., Sanchez, V., Sreaton, G.R., Semple, M.G., Slon-Campos, J., Skelly, D.T., Smith, E.N., Sobrinodiaz, A., Staves, J., Stuart, D.I., Supasa, P., Surik, T., Thraves, H., Tsang, P., Turtle, L., Walker, A.S., Wang, B., Washington, C., Watkins, N., Whitehouse, J., National, C.T.S.A.P., 2020. Antibody testing for COVID-19: a report from the national COVID scientific advisory panel. *Wellcome Open Res.* 5 (139), 139 <https://doi.org/10.12688/wellcomeopenres.15927.1>.
- Ahmed, W., Angel, N., Edson, J., Bibby, K., Bivins, A., O'Brien, J.W., Choi, P.M., Kitajima, M., Simpson, S.L., Li, J., Tschärke, B., Verhagen, R., Smith, W.J.M., Zaugg, J., Dierens, L., Hugenholz, P., Thomas, K.V., Mueller, J.F., 2020. 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. *Sci. Total Environ.* 728, 138764 <https://doi.org/10.1016/j.scitotenv.2020.138764>.
- Andersen, K.G., Rambaut, A., Lipkin, W.I., Holmes, E.C., Garry, R.F., 2020. The proximal origin of SARS-CoV-2. *Nat. Med.* 26 (4), 450–452 <https://doi.org/10.1038/s41591-020-0820-9>.
- Balboa, S., Mauricio-Iglesias, M., Rodriguez, S., Martinez-Lamas, L., Vassallo, F.J., Regueiro, B., Lema, J.M., 2021. The fate of SARS-CoV-2 in WWTPs points out the sludge line as a suitable spot for detection of COVID-19. *Sci. Total Environ.* 772, 145268 <https://doi.org/10.1016/j.scitotenv.2021.145268>.
- Becker, M., Strengert, M., Junker, D., Kaiser, P.D., Kerrinnes, T., Traenkler, B., Dinter, H., Haring, J., Ghozzi, S., Zeck, A., Weise, F., Peter, A., Horber, S., Fink, S., Ruoff, F., Dulovic, A., Bakchoul, T., Baillet, A., Lohse, S., Cornberg, M., Illig, T., Gottlieb, J., Smola, S., Karch, A., Berger, K., Rammensee, H.G., Schenke-Layland, K., Nelde, A., Marklin, M., Heitmann, J.S., Walz, J.S., Templin, M., Joos, T.O., Rothbauer, U., Krause, G., Schneiderhan-Marra, N., 2021. Exploring beyond clinical routine SARS-CoV-2 serology using MultiCoV-Ab to evaluate endemic coronavirus cross-reactivity. *Nat. Commun.* 12 (1), 1152 <https://doi.org/10.1038/s41467-021-20973-3>.
- Boni, M.F., Lemey, P., Jiang, X., Lam, T.T., Perry, B.W., Castoe, T.A., Rambaut, A., Robertson, D.L., 2020. Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible for the COVID-19 pandemic. *Nat. Microbiol.* 5 (11), 1408–1417 <https://doi.org/10.1038/s41564-020-0771-4>.
- Bryant, J.E., Azman, A.S., Ferrari, M.J., Arnold, B.F., Boni, M.F., Boum, Y., Hayford, K., Luquero, F.J., Mina, M.J., Rodriguez-Barraguer, I., Wu, J.T., Wade, D., Vernet, G., Leung, D.T., 2020. Serology for SARS-CoV-2: apprehensions, opportunities, and the path forward. *Sci. Immunol.* 5 (47) <https://doi.org/10.1126/sciimmunol.abc6347>.
- Buitrago-Garcia, D., Egli-Gany, D., Counotte, M.J., Hossmann, S., Imeri, H., Ipekci, A.M., Salanti, G., Low, N., 2020. Occurrence and transmission potential of asymptomatic and presymptomatic SARS-CoV-2 infections: a living systematic review and meta-analysis. *PLoS Med.* 17 (9), e1003346 <https://doi.org/10.1371/journal.pmed.1003346>.
- Byambasuren, O., Cardona, M., Bell, K., Clark, J., McLaws, M.-L., Glasziou, P., 2020. Estimating the extent of asymptomatic COVID-19 and its potential for community transmission: Systematic review and meta-analysis. *Off. J. Assoc. Med. Microbiol. Infect. Dis. Can.* 5 (4), 223–234 <https://doi.org/10.3138/jammi-2020-0030>.
- CDC 2020a Centers for disease control and prevention - CDC diagnostic tests for COVID-19. <https://www.cdc.gov/coronavirus/2019-ncov/lab/testing.html>.
- CDC 2020b Centers for disease control and prevention - serology testing for COVID-19 at CDC. <https://www.cdc.gov/coronavirus/2019-ncov/lab/serology-testing.html>.
- Chen, Y., Chen, L., Deng, Q., Zhang, G., Wu, K., Ni, L., Yang, Y., Liu, B., Wang, W., Wei, C., Yang, J., Ye, G., Cheng, Z., 2020. The presence of SARS-CoV-2 RNA in the feces of COVID-19 patients. *J. Med. Virol.* 92 (7), 833–840 <https://doi.org/10.1002/jmv.25825>.
- Corman, V.M., Landt, O., Kaiser, M., Molenkamp, R., Meijer, A., Chu, D.K., Bleicker, T., Brunink, S., Schneider, J., Schmidt, M.L., Mulders, D.G., Haagmans, B.L., van der Veer, B., van den Brink, S., Wijsman, L., Goderski, G., Romette, J.L., Ellis, J., Zambon, M., Peiris, M., Goossens, H., Reusken, C., Koopmans, M.P., Drosten, C., 2020. Detection of 2019 novel coronavirus (2019-nCoV) by real-time RT-PCR. *Euro Surveill.* 25 (3) <https://doi.org/10.2807/1560-7917.ES.2020.25.3.2000045>.
- Crits-Christoph, A., Kantor, R.S., Olm, M.R., Whitney, O.N., Al-Shayeb, B., Lou, Y.C., Flamholz, A., Kennedy, L.C., Greenwald, H., Hinkle, A., Hetzel, J., Spitzer, S., Koble, J., Tan, A., Hyde, F., Schroth, G., Kuersten, S., Banfield, J.F., Nelson, K.L., 2021. Genome sequencing of sewage detects regionally prevalent SARS-CoV-2 variants. *mBio* 12 (1) <https://doi.org/10.1128/mBio.02703-20>.
- D'Aoust, P.M., Mercier, E., Montpetit, D., Jia, J.J., Alexandrov, I., Neault, N., Baig, A.T., Mayne, J., Zhang, X., Alain, T., Langlois, M.A., Servos, M.R., MacKenzie, M., Figeys, D., MacKenzie, A.E., Graber, T.E., Delatolla, R., 2021. Quantitative analysis of SARS-CoV-2 RNA from wastewater solids in communities with low COVID-19 incidence and prevalence. *Water Res.* 188, 116560 <https://doi.org/10.1016/j.watres.2020.116560>.
- De Maio, N., Walker, C., Borges, R., Weilguny, L., Slodkowitz, G. and Goldman, N. 2020. Issues with SARS-CoV-2 sequencing data. <https://virological.org/t/issues-with-sars-cov-2-sequencing-data/473/1>.
- Dong, E., Du, H., Gardner, L., 2020. An interactive web-based dashboard to track COVID-19 in real time. *Lancet Infect. Dis.* 20 (5), 533–534 [https://doi.org/10.1016/S1473-3099\(20\)30120-1](https://doi.org/10.1016/S1473-3099(20)30120-1).
- Elbe, S., Buckland-Merrett, G., 2017. Data, disease and diplomacy: GISAID's innovative contribution to global health. *Glob. Chall.* 1 (1), 33–46 <https://doi.org/10.1002/gch2.1018>.
- Farkas, K., Hillary, L.S., Malham, S.K., McDonald, J.E., Jones, D.L., 2020. Wastewater and public health: the potential of wastewater surveillance for monitoring COVID-19. *Curr. Opin. Environ. Sci. Health* 17, 14–20 <https://doi.org/10.1016/j.coesh.2020.06.001>.
- Gorbalenya, A.E., Baker, S.C., Baric, R.S., de Groot, R.J., Drosten, C., Gulyaeva, A.A., Haagmans, B.L., Lauber, C., Leontovich, A.M., Neuman, B.W., Penzar, D., Perlman, S., Poon, L.L.M., Samborskiy, D.V., Sidorov, I.A., Sola, I., Ziebuhr, J., Coronaviridae Study Group of the International Committee on Taxonomy of, V., 2020. The species severe acute respiratory syndrome-related coronavirus: classifying 2019-nCoV and naming it SARS-CoV-2. *Nat. Microbiol.* 5 (4), 536–544 <https://doi.org/10.1038/s41564-020-0695-z>.
- Gower, J.C., 1966. Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika* 53 (3/4), 325 <https://doi.org/10.2307/2333639>.
- Grubaugh, N.D., Gangavarapu, K., Quick, J., Matteson, N.L., De Jesus, J.G., Main, B.J., Tan, A.L., Paul, L.M., Brackney, D.E., Grewal, S., Gurfield, N., Van Rompay, K.K.A., Isern, S., Michael, S.F., Coffey, L.L., Loman, N.J., Andersen, K.G., 2019. An amplicon-based sequencing framework for accurately measuring intrahost virus diversity using PrimalSeq and iVar. *Genome Biol.* 20 (1), 8 <https://doi.org/10.1186/s13059-018-1618-7>.
- Holland, L.A., Kaelin, E.A., Maqsood, R., Estifanos, B., Wu, L.I., Varsani, A., Halden, R.U., Hogue, B.G., Scotch, M., Lim, E.S., 2020. An 81-nucleotide deletion in SARS-CoV-2 ORF7a identified from sentinel surveillance in Arizona (January to March 2020). *J. Virol.* 94 (14) <https://doi.org/10.1128/JVI.00711-20>.
- Izquierdo-Lara, R., Elsinga, G., Heijnen, L., Oude Munnink, B.B., Schapendonk, C.M.E., Nieuwenhuijse, D., Kon, M., Lu, L., Aarestrup, F.M., Lyckett, S., Medema, G., Koopmans, M.P.G., de Graaf, M., 2021. Monitoring SARS-CoV-2 circulation and diversity through community wastewater sequencing, the Netherlands and Belgium. *Emerg. Infect. Dis.* 27 (5) <https://doi.org/https://wwwnc.cdc.gov/eid/article/27/5/20-4410.article>.
- Jones, D.L., Baluja, M.Q., Graham, D.W., Corbishley, A., McDonald, J.E., Malham, S.K., Hillary, L.S., Connor, T.R., Gaze, W.H., Moura, I.B., Wilcox, M.H., Farkas, K., 2020. 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. *Sci. Total Environ.* 749, 141364 <https://doi.org/10.1016/j.scitotenv.2020.141364>.
- Kimball, A., Hatfield, K.M., Arons, M., James, A., Taylor, J., Spicer, K., Bardossy, A.C., Oakley, L.P., Tanwar, S., Chisty, Z., Bell, J.M., Methner, M., Harney, J., Jacobs, J.R., Carlson, C.M., McLaughlin, H.P., Stone, N., Clark, S., Brostrom-Smith, C., Page, L.C., Kay, M., Lewis, J., Russell, D., Hiatt, B., Gant, J., Duchin, J.S., Clark, T.A., Honein, M. A., Reddy, S.C., Jernigan, J.A., Public Health, S., King, C., Team, C.C.-I., 2020. Asymptomatic and presymptomatic SARS-CoV-2 infections in residents of a long-term care skilled nursing facility - king county, Washington, March 2020. *MMWR Morb. Mortal. Wkly. Rep.* 69 (13), 377–381 <https://doi.org/10.15585/mmwr.mm6913e1>.
- Kocamei, B.A., Kurt, H., Sait, A., Sarac, F., Saatci, A.M. and Pakdemirli, B. 2020. SARS-CoV-2 detection in istanbul wastewater treatment plant sludges. <https://doi.org/10.1101/2020.05.12.20099358>.
- Kumar, M., Patel, A.K., Shah, A.V., Raval, J., Rajpara, N., Joshi, M., Joshi, C.G., 2020. First proof of the capability of wastewater surveillance for COVID-19 in India through detection of genetic material of SARS-CoV-2. *Sci. Total Environ.* 746, 141326 <https://doi.org/10.1016/j.scitotenv.2020.141326>.
- La Rosa, G., Iaconelli, M., Mancini, P., Bonanno Ferraro, G., Veneri, C., Bonadonna, L., Lucentini, L., Suffredini, E., 2020. First detection of SARS-CoV-2 in untreated wastewaters in Italy. *Sci. Total Environ.* 736, 139652 <https://doi.org/10.1016/j.scitotenv.2020.139652>.
- Ladner, J.T., Larsen, B.B., Bowers, J.R., Hepp, C.M., Bolyen, E., Folkerts, M., Sheridan, K., Pfeiffer, A., Yaglom, H., Lemmer, D., Sahl, J.W., Kaelin, E.A., Maqsood, R., Bokulich, N.A., Quirk, G., Watts, T.D., Komatsu, K.K., Waddell, V., Lim, E.S., Caporaso, J.G., Engelthaler, D.M., Worobey, M., Keim, P., 2020. An early pandemic analysis of SARS-CoV-2 population structure and dynamics in Arizona. *mBio* 11 (5) <https://doi.org/10.1128/mBio.02107-20>.
- Li, H., Durbin, R., 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25 (14), 1754–1760 <https://doi.org/10.1093/bioinformatics/btp324>.
- Maricopa County 2020 Maricopa County, AZ. <https://www.maricopa.gov/>.
- Medema, G., Heijnen, L., Elsinga, G., Italiaander, R., Brouwer, A., 2020. 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. *Environ. Sci. Technol. Lett.* 7 (7), 511–516 <https://doi.org/10.1021/acs.estlett.0c00357>.
- Nemudryi, A., Nemudraia, A., Wiegand, T., Surya, K., Buyukyoruk, M., Vanderwood, K. K., Wilkinson, R., Wiedenheft, B., 2020. Temporal detection and phylogenetic

- assessment of SARS-CoV-2 in municipal wastewater. *Cell Rep Med* 1 (6), 1000098 <https://doi.org/10.1016/j.xcrm.2020.100098>.
- Park, S.K., Lee, C.W., Park, D.I., Woo, H.Y., Cheong, H.S., Shin, H.C., Ahn, K., Kwon, M. J., Joo, E.J., 2020. Detection of SARS-CoV-2 in fecal samples from patients with asymptomatic and mild COVID-19 in Korea. *Clin. Gastroenterol. Hepatol.* 10.1016/j.cgh.2020.06.005. <https://doi.org/10.1016/j.cgh.2020.06.005>.
- Peccia, J., Zulli, A., Brackney, D.E., Grubaugh, N.D., Kaplan, E.H., Casanovas-Massana, A., Ko, A.I., Malik, A.A., Wang, D., Wang, M., Warren, J.L., Weinberger, D. M., Arnold, W., Omer, S.B., 2020. Measurement of SARS-CoV-2 RNA in wastewater tracks community infection dynamics. *Nat. Biotechnol.* 38 (10), 1164–1167 <https://doi.org/10.1038/s41587-020-0684-z>.
- Rambaut, A., Holmes, E.C., O'Toole, A., Hill, V., McCrone, J.T., Ruis, C., du Plessis, L., Pybus, O.G., 2020a. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat. Microbiol.* 5 (11), 1403–1407 <https://doi.org/10.1038/s41564-020-0770-5>.
- Rambaut, A., Loman, N., Pybus, O., Barclay, W., Barrett, J., Carabelli, A., Connor, T., Peacock, T., Robertson, D.L. and Volz, E. 2020b Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in the UK defined by a novel set of spike mutations. <https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563>.
- Randazzo, W., Truchado, P., Cuevas-Ferrando, E., Simon, P., Allende, A., Sanchez, G., 2020. SARS-CoV-2 RNA in wastewater anticipated COVID-19 occurrence in a low prevalence area. *Water Res.* 181, 115942 <https://doi.org/10.1016/j.watres.2020.115942>.
- Shu, Y., McCauley, J., 2017. GISAID: Global initiative on sharing all influenza data - from vision to reality. *Euro Surveill.* 22 (13) <https://doi.org/10.2807/1560-7917.ES.2017.22.13.30494>.
- Syangtan, G., Bista, S., Dawadi, P., Rayamajhee, B., Shrestha, L.B., Tuladhar, R. and Joshi, D.R. 2020. Asymptomatic SARS-CoV-2 carriers: a systematic review and meta-analysis. *front public health* 8, 587374. <https://doi.org/10.3389/fpubh.2020.587374>.
- Tang, A., Tong, Z.D., Wang, H.L., Dai, Y.X., Li, K.F., Liu, J.N., Wu, W.J., Yuan, C., Yu, M. L., Li, P., Yan, J.B., 2020. Detection of novel coronavirus by RT-PCR in stool specimen from asymptomatic child, China. *Emerg. Infect. Dis.* 26 (6), 1337–1339 <https://doi.org/10.3201/eid2606.200301>.
- Tegally, H., Wilkinson, E., Giovanetti, M., Iranzadeh, A., Fonseca, V., Giandhari, J., Doolabh, D., Pillay, S., San, E.J., Msomi, N., Mlisana, K., von Gottberg, A., Walaza, S., Allam, M., Ismail, A., Mohale, T., Glass, A.J., Engelbrecht, S., Van Zyl, G., Preiser, W., Petruccione, F., Sigal, A., Hardie, D., Marais, G., Hsiao, N.Y., Korsman, S., Davies, M.A., Tyers, L., Mudau, I., York, D., Maslo, C., Goedhals, D., Abrahams, S., Laguda-Akingba, O., Alisoltani-Dehkordi, A., Godzik, A., Wibmer, C. K., Sewell, B.T., Lourenco, J., Alcantara, L.C.J., Kosakovsky Pond, S.L., Weaver, S., Martin, D., Lessells, R.J., Bhiman, J.N., Williamson, C., de Oliveira, T., 2021. Detection of a SARS-CoV-2 variant of concern in South Africa. *Nature* 592 (7854), 438–443 <https://doi.org/10.1038/s41586-021-03402-9>.
- Volz, E., Hill, V., McCrone, J.T., Price, A., Jorgensen, D., O'Toole, A., Southgate, J., Johnson, R., Jackson, B., Nascimento, F.F., Rey, S.M., Nicholls, S.M., Colquhoun, R. M., da Silva Filipe, A., Shepherd, J., Pascall, D.J., Shah, R., Jesudason, N., Li, K., Jarrett, R., Pacchiarini, N., Bull, M., Geidelberg, L., Siveroni, I., Consortium, C.-U., Goodfellow, I., Loman, N.J., Pybus, O.G., Robertson, D.L., Thomson, E.C., Rambaut, A., Connor, T.R., 2021. Evaluating the effects of SARS-CoV-2 spike mutation D614G on transmissibility and pathogenicity. *Cell* 184 (1), 64–75 e11 <https://doi.org/10.1016/j.cell.2020.11.020>.
- Wang, H., Li, X., Li, T., Zhang, S., Wang, L., Wu, X., Liu, J., 2020. The genetic sequence, origin, and diagnosis of SARS-CoV-2. *Eur. J. Clin. Microbiol. Infect. Dis.* 39 (9), 1629–1635 <https://doi.org/10.1007/s10096-020-03899-4>.
- Westhaus, S., Weber, F.A., Schiwy, S., Linnemann, V., Brinkmann, M., Wiedera, M., Greve, C., Janke, A., Hollert, H., Wintgens, T., Ciesek, S., 2021. Detection of SARS-CoV-2 in raw and treated wastewater in Germany - suitability for COVID-19 surveillance and potential transmission risks. *Sci. Total Environ.* 751, 141750 <https://doi.org/10.1016/j.scitotenv.2020.141750>.
- WHO 2020a World health organisation - SARS-CoV-2 assay. https://www.int/docs/default-source/coronaviruse/whoinhouseassays.pdf?sfvrsn=de3a76aa_2.
- WHO 2020b World health organisation - serology in the context of COVID-19. <https://www.who.int/emergencies/diseases/novel-coronavirus-2019/serology-in-the-context-of-covid-19>.
- Wu, F., Zhang, J., Xiao, A., Gu, X., Lee, W.L., Armas, F., Kauffman, K., Hanage, W., Matus, M., Ghaeli, N., Endo, N., Duvallet, C., Poyet, M., Moniz, K., Washburne, A.D., Erickson, T.B., Chai, P.R., Thompson, J., Alm, E.J., 2020. SARS-CoV-2 titers in wastewater are higher than expected from clinically confirmed cases. *mSystems* 5 (4), 00614–00620 <https://doi.org/10.1128/mSystems.00614-20>.
- Wurtzer, S., Marechal, V., Mouchel, J.M., Maday, Y., Teyssou, R., Richard, E., Almayrac, J.L. and Moulin, L. 2020. Evaluation of lockdown effect on SARS-CoV-2 dynamics through viral genome quantification in waste water, Greater Paris, France, 5 March to 23 April 2020. *Euro Surveill* 25(50), 2000776. <https://doi.org/10.2807/1560-7917.ES.2020.25.50.2000776>.
- Xing, Y.H., Ni, W., Wu, Q., Li, W.J., Li, G.J., Wang, W.D., Tong, J.N., Song, X.F., Wing-Kin Wong, G., Xing, Q.S., 2020. Prolonged viral shedding in feces of pediatric patients with coronavirus disease 2019. *J. Microbiol. Immunol. Infect.* 53 (3), 473–480 <https://doi.org/10.1016/j.jmii.2020.03.021>.
- Yan, Y., Shin, W.I., Pang, Y.X., Meng, Y., Lai, J., You, C., Zhao, H., Lester, E., Wu, T., Pang, C.H., 2020. The first 75 days of novel coronavirus (SARS-CoV-2) outbreak: recent advances, prevention, and treatment. *Int. J. Environ. Res. Public Health* 17 (7) <https://doi.org/10.3390/ijerph17072323>.
- Yue, J.C., Clayton, M.K., 2005. A Similarity measure based on species proportions. *Commun. Stat. - Theory Methods* 34 (11), 2123–2131 <https://doi.org/10.1080/sta-200066418>.
- Zhang, D., Ling, H., Huang, X., Li, J., Li, W., Yi, C., Zhang, T., Jiang, Y., He, Y., Deng, S., Zhang, X., Wang, X., Liu, Y., Li, G., Qu, J., 2020. Potential spreading risks and disinfection challenges of medical wastewater by the presence of Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) viral RNA in septic tanks of Fangcang Hospital. *Sci. Total Environ.* 741, 140445 <https://doi.org/10.1016/j.scitotenv.2020.140445>.
- Zhang, Y.Z., Holmes, E.C., 2020. A Genomic perspective on the origin and emergence of SARS-CoV-2. *Cell* 181 (2), 223–227 <https://doi.org/10.1016/j.cell.2020.03.035>.