

Detecting SARS-CoV-2 lineages and mutational load in municipal wastewater; a use-case in the metropolitan area of Thessaloniki, Greece

Nikolaos Pechlivanis

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Maria Tsagiopoulou

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Maria Maniou

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Anastasis Togkousidis

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Evangelia Mouchtaropoulou

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Taxiarchis Chassalevris

School of Veterinary Medicine, Aristotle University of Thessaloniki

Serafeim Chaintoutis

School of Veterinary Medicine, Aristotle University of Thessaloniki

Maria Petala

Dept. of Civil Engineering, Aristotle University of Thessaloniki

Margaritis Kostoglou

Dept. of Chemistry, Aristotle University of Thessaloniki

Thodoris Karapantsios

Dept. of Chemistry, Aristotle University of Thessaloniki

Stamatia Laidou

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Elisavet Vlachonikola

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Anastasia Chatzidimitriou

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Agis Papadopoulos

EYATH S.A., Thessaloniki Water Supply and Sewerage Company S.A.

Nikolaos Papaioannou

School of Veterinary Medicine, Aristotle University of Thessaloniki

Chrysostomos Dovas

School of Veterinary Medicine, Aristotle University of Thessaloniki

Anagnostis Argiriou

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Fotis Psomopoulos (✉ fpsom@certh.gr)

Institute of Applied Biosciences, Centre of Research and Technology Hellas

Research Article

Keywords: SARS-CoV-2, next-generation sequencing (NGS), novel methodology, called lineagespot, entrepreneurs, traders

Posted Date: July 9th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-677811/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Scientific Reports on February 17th, 2022. See the published version at <https://doi.org/10.1038/s41598-022-06625-6>.

1 Detecting SARS-CoV-2 lineages and mutational load in municipal 2 wastewater; a use-case in the metropolitan area of Thessaloniki, Greece

3
4 Nikolaos Pechlivanis^{1,2}, Maria Tsagiopoulou¹, Maria Christina Maniou¹, Anastasis Togkousidis¹,
5 Evangelia Mouchtaropoulou¹, Taxiarchis Chassalevris³, Serafeim C. Chaintoutis³, Maria Petala⁴,
6 Margaritis Kostoglou⁵, Thodoris Karapantsios⁵, Stamatia Laidou^{1,2}, Elisavet Vlachonikola^{1,2}, Anastasia
7 Chatzidimitriou¹, Agis Papadopoulos⁶, Nikolaos Papaioannou³, Chrysostomos I. Dovas³, Anagnostis
8 Argiriou^{1,7}, Fotis Psomopoulos^{1,*}

9
10 ¹ Institute of Applied Biosciences, Centre of Research and Technology Hellas, Themi, 57001
11 Thessaloniki, Greece

12 ² Dept. of Genetics, Development and Molecular Biology, School of Biology, Aristotle University of
13 Thessaloniki, 54124 Thessaloniki, Greece

14 ³ School of Veterinary Medicine, Aristotle University of Thessaloniki, 54124 Thessaloniki, Greece

15 ⁴ Dept. of Civil Engineering, Aristotle University of Thessaloniki, Thessaloniki, 54 124, Greece

16 ⁵ Dept. of Chemistry, Aristotle University of Thessaloniki, 54 124 Thessaloniki, 54124, Greece

17 ⁶ EYATH S.A., Thessaloniki Water Supply and Sewerage Company S.A., Thessaloniki, 54636, Greece

18 ⁷ Department of Food Science and Nutrition, University of the Aegean, 81400 Myrina, Lemnos, Greece

19
20 *Correspondence to fpsom@certh.gr

22 Abstract

23 The SARS-CoV-2 pandemic represents an unprecedented global crisis necessitating novel approaches
24 for, amongst others, early detection of emerging variants relating to the evolution and spread of the
25 virus. Recently, the detection of SARS-CoV-2 RNA in wastewater has emerged as a useful tool to
26 monitor the prevalence of the virus in the community. Here, we propose a novel methodology, called
27 **lineagespot**, for the monitoring of variants and the detection of SARS-CoV-2 lineages in wastewater
28 samples using next-generation sequencing (NGS). Our proposed method was tested and evaluated
29 using NGS data produced by the sequencing of fourteen wastewater samples from the municipality of
30 Thessaloniki, Greece, covering a six month period. The results showed a clear identification of trends
31 in the presence of SARS-CoV-2 mutations in wastewater data. Lineagespot was able to record the
32 evolution and rapid domination of the B.1.1.7 lineage in the community, and allowed for a robust
33 inference between the variants evident through our approach and the variants observed in patients
34 from the same area time periods. Lineagespot is an open-source tool, implemented in R, and is freely
35 available on [GitHub](https://github.com).

36 Introduction

37 Nearly a year after the first report of SARS-CoV-2 in Wuhan, China, the virus has spread at an
38 unprecedented pace causing a global pandemic. As the main transmission process of the SARS-CoV-2

39 virus is through droplets and the contact between people, several testing strategies identify whether
40 a person is infected and, in cases of a positive sample, which the underlying virus variant is. However,
41 these methods are not easily scalable, especially in large urban areas, where a high number of
42 individuals have to be tested to assess virus and variant spread among the population. Interestingly,
43 the viral RNA can also be detected in wastewater, with SARS-CoV-2 RNA levels in wastewater
44 correlating with the COVID-19 epidemiology¹⁻³. Indeed, in the previous work of Petala et al³,
45 normalized viral copy levels in Thessaloniki wastewater agreed with the epidemiological conditions in
46 the city. Thessaloniki is the second largest city in Greece with around 1 million inhabitants. The city is
47 a chief gateway for entrepreneurs, traders, university students and tourists visiting Greece and, as
48 such, it was the place where the Greek patient “zero” appeared in March 2020 as well as the so-called
49 South African mutation for the first time. In other words, the city of Thessaloniki is a suitable place as
50 a case study for identifying new variants in the city wastewater before scattering to the rest of the
51 country.

52
53 The presence of SARS CoV-2 RNA in wastewater provides us with a unique opportunity, i.e., to identify
54 the most prevalent virus lineages through the analysis of the traces evident in wastewater samples.
55 So far, although there are few studies exploring the SARS-CoV-2 diversity in wastewater, it still remains
56 an open issue as there are no widely accepted methods that can sufficiently address this. The most
57 common used approaches involve the sequencing of the wastewater samples, and the consequent
58 application of low frequency variant analysis methods⁴ or metagenomic approaches^{5,6}. In either case,
59 the interpretation of the results focuses on the detection of specific variants⁴ or lineages² such as
60 B.1.1.7 and 501.V2, prevalent clades⁶ (19A, 20A and 20B) or new uncharacterized mutations⁶.

61
62 In this work we propose a novel methodology called *lineagespot*, implemented as a software tool that
63 can facilitate the detection of SARS-CoV-2 lineages in wastewater samples using next-generation
64 sequencing (NGS). The method is tested and validated across fourteen municipal wastewater samples
65 retrieved in Thessaloniki, Greece in fourteen different time periods, and correlated with the variants
66 and lineages observed in patients from the same area time points. Based on a variation of the Illumina
67 Arctic pipeline for the identification of mutations at low frequencies (< 0.01), and the lineage
68 assignments defined by Pangolin, this method identifies all SARS-CoV-2 mutations present in the
69 wastewater, and attempts to infer the potential distribution of the SARS-CoV-2 lineages. The
70 methodology is proven to be effective in detecting the mutational load in the wastewater, with the
71 inferred lineages being roughly aligned to the predominant lineages identified through targeted (and
72 therefore biased) patient-derived genotypes.

73 **Results**

74 **Comparison of variant calling methods**

75 The proposed methodology was initially applied on a selected wastewater sample (corresponding to
76 the 05-11 February 2021 time period), and for which three different variant callers were assessed: 1)
77 *freebayes*, 2) *mpileup* and 3) *GATK Mutect2* (cancer only mode). In terms of parameters, *freebayes*

78 was used with a low variant frequency parameter of 0.01, *mpileup* reported every position (either
 79 reference, or variant), and *GATK Mutect2* was used with the default parameters.

80

81 An example of the output produced by the methodology, regardless of the variant calling method, is
 82 shown in **Table 1**. In this table the overlap between the Pangolin's rules and the rules generated by
 83 the tool for the input dataset is captured for each lineage. In order to quantify the overlap, three basic
 84 metrics are produced; the overlap by considering pangolin's rules as a decision tree (*Tree Overlap*),
 85 the total overlap regardless of the rules order (*Total Overlap*), and the overlap for the rules that are
 86 satisfied only by the identified mutations (i.e., explicitly listed in the variants' file), and therefore
 87 excluding all rules based on the unmutated reference (*Total Overlap Var*). In addition to the previous
 88 metrics, the related ratio values are also calculated (*Tree ratio*, *Total Ratio*, *Total Ratio Var*). Finally,
 89 information regarding the read depth for each position (reference and variant) is also provided in the
 90 output file.

91

92 **Table 1:** Each row in the table corresponds to a single lineage rule defined by pangolin. The columns
 93 correspond to the different metrics captured, in order to perform the systematic evaluation.

Lineage	Rules	Total	Tree Overlap	Total Overlap	Total Overlap Var	Tree Ratio	Total Ratio	Total Ratio Var	Tree Avg AD	Total Avg AD	Total Avg DP	Total Sum AD	Total Sum DP	Avg DP	Total Run Reads	Avg AF
	26800															
B.1.177	=='C',	17	0	11	3	0	0.647	0.176	0	6.6	19	33	95	25.2	104214	0.347
	...															
	26800															
B.1.177	=='C',	11	0	7	2	0	0.636	0.181	0	7	23.66	21	71	25.2	104214	0.295
	...															
	26800															
B.1.1.7	!='C',	13	3	11	1	0.230	0.846	0.076	0	2	9	2	9	25.2	104214	0.222
	...															

94

95 Based on this detailed table, a second output is generated as a simplified summary. In this case, all
 96 rows for which the *Total Overlap Var* column was equal to 0 were removed, and therefore potential
 97 lineages that would be assigned based only on the unmutated reference (i.e., no actual mutations
 98 detected) are excluded from the analysis. The remaining rows were collapsed (**Table 2**) through a
 99 process in which the average values of the basic metrics were calculated for each lineage; i.e., the
 100 mean of the *Tree Ratio*, the *Total Ratio*, the *Total Ratio Var*, the *Tree Av AD*, the *Total Av AD*, and the
 101 *Avg AF* columns.

102

103 **Table 2:** Table shows the corresponding metrics to a unique lineage, after the merging process. The
 104 metrics can be consequently used to assess the presence of the particular variant in the dataset.

Lineage	Mean Tree Ratio	Mean Total Ratio	Mean Total Ratio Var	Mean Total Av AD	Mean AF
B.1.1.7	0.231	0.846	0.077	2	0.223
B.1.177	0.003	0.611	0.057	9.769	0.488

105

106 Depending on the variant caller tool used (*freebayes*, *mpileup* and *GATK Mutect2*), **lineagespot**
 107 generates a unique output. All outputs are compared pairwise, based on the decision tree rules (n_d),
 108 and the total number of rules satisfied (n_t). For each lineage, the absolute values of the differences
 109 between the two metrics (n_d , n_t) of the files are calculated. As an example, for lineage *A.1*, the output
 110 produced by using the *freebayes* variant caller tool in the first step of the methodology, returns $n_d =$
 111 2 and $n_t = 10$ rules satisfied, while the output of the *GATK Mutect2* tool returns $n_d = 1$ and $n_t = 4$
 112 rules satisfied. As a result, the two outputs exhibit a difference of $n_d = 1$ and $n_t = 6$ rules (**Table 3**).

113 In addition, the total number of lineages are shown in **Table 4**. In the same table, the maximum
 114 absolute n_d difference and the maximum n_t difference for each pair of files are calculated (**Figure 1A**).
 115 The latter is used for an overall comparison of the output files.

116 **Table 3:** Snapshot of the difference between the two metrics (n_d , n_t) across the output files coming
 117 from *freebayes* variant caller and *mpileup*.

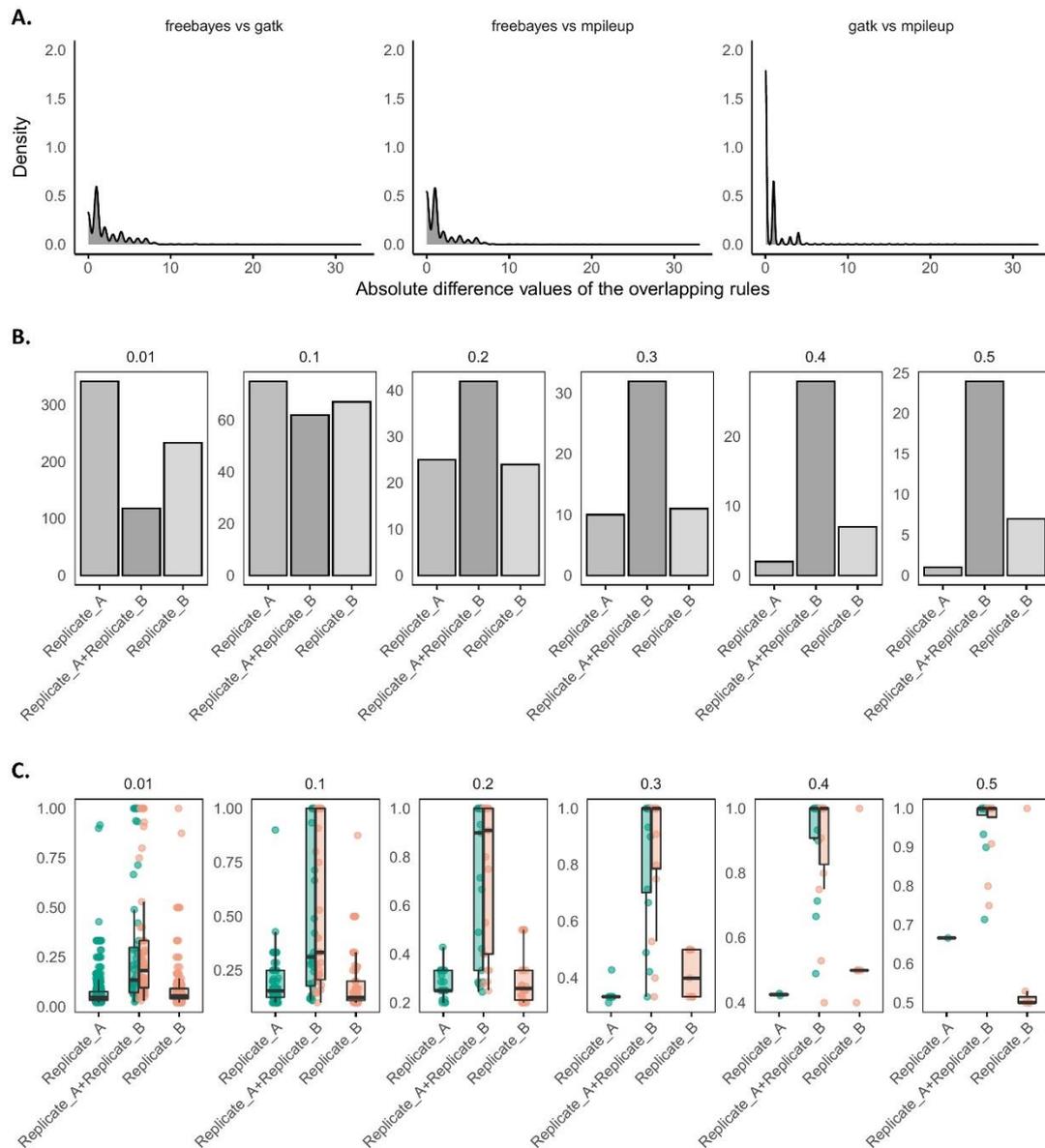
Lineage	Rules	N_t	$n_d^{\text{freebayes}}$	$n_t^{\text{freebayes}}$	n_d^{mpileup}	n_t^{mpileup}	$ n_d^{\text{freebayes}} - n_d^{\text{mpileup}} $	$ n_t^{\text{freebayes}} - n_t^{\text{mpileup}} $
B.1.1.7	26800!= 'C', ...	13	3	72	3	71	0	1
B.1.177	26800== 'C', ...	66	3	72	3	69	0	3
B.1.177	26800== 'C', ...	6	3	72	3	69	0	3
B.1.351	26800!= 'C', ...	38	3	162	3	160	0	2

118

119 **Table 4:** Summary table of the three output files produced by *freebayes*, *mpileup* and *GATK mutect2*
 120 variant caller. The three files are compared in pairs.

Files	Number of differences	max absolute n_d dif- ference	max absolute n_t dif- ference
freebayes – gatk	3791	4	31
freebayes – mpileup	3140	4	33
gatk – mpileup	1571	1	31

121



122

123 **Figure 1:** Evolution of variants across different low frequency parameters. **A.** Density plot of the
 124 absolute n_t difference values between the output of the three variant calling tools use (pairwise
 125 comparisons). **B.** Number of reads for each replicate and for the common variants **C.** The
 126 corresponding allele frequency for each replicate and for the common variants

127 Based on the above comparison, we consider that the most productive and informative approach is
 128 to utilize *freebayes* as the variant calling tool. The rest of the results shown below, are based only on
 129 the *freebayes* tool output.

130 Evaluating lineage-specific mutations across time periods

131 Sensitivity

132 In order to investigate how specific variants, evolve over time, all variant files were merged into a
133 combined table in which all the detected nucleotide variants along with the corresponding amino acid
134 substitutions that have been identified are stored. Moreover, information about the read depth, the
135 allele frequency, and the overlapping gene is also provided for each variant. **Table 5** gives an example
136 of the overall information.

137

138 **Table 5:** Snapshot of table containing every mutation per sample along with the corresponding gene
139 and the amino acid variant.

CHROM	POS	REF	ALT	DP	AD alt	TYPE	Gene Name	HGVS	AF	sample
NC_045512.2	326	T	A	7	1	snp	ORF1ab	Leu21Ile	0.143	Sample A
NC_045512.2	378	T	C	10	1	snp	ORF1ab	Val38Ala	0.1	Sample A
NC_045512.2	408	A	T	10	1	snp	ORF1ab	Asp48Val	0.1	Sample B
NC_045512.2	433	T	C	10	2	snp	ORF1ab	Val56Val	0.2	Sample C
NC_045512.2	442	C	T	10	1	snp	ORF1ab	Gly59Gly	0.1	Sample C

140

141 Having a structured format, we first examined if the number of reads that were produced during PCR
142 amplification is affecting the number of variants that are identified in every sample. For this reason,
143 two replicates of the same biological specimen were produced. The first replicate (*Replicate A*)
144 contained 181,880 reads while the second replicate (*Replicate B*) contained 69,706 reads. The two
145 replicates were analyzed as described in the “*Raw data analysis*” Section and two VCF files were
146 produced which contained 401 and 293 variants respectively; of these, 59 variants were common in
147 both replicates, 342 variants were unique for Replicate A and 234 unique for Replicate B.

148

149 A Student's t-test was performed in order to compare the mean values of the allele frequency between
150 the two replicates and an F test to compare the variances. All tests resulted in p values higher than
151 the 0.05 threshold, meaning that no statistically significant difference was found between the two
152 replicates. On the contrary, when comparing two samples coming from different time periods there
153 was a significant difference on the variants' allele frequency (p values from Student's t-test and F test
154 were below 0.05).

155

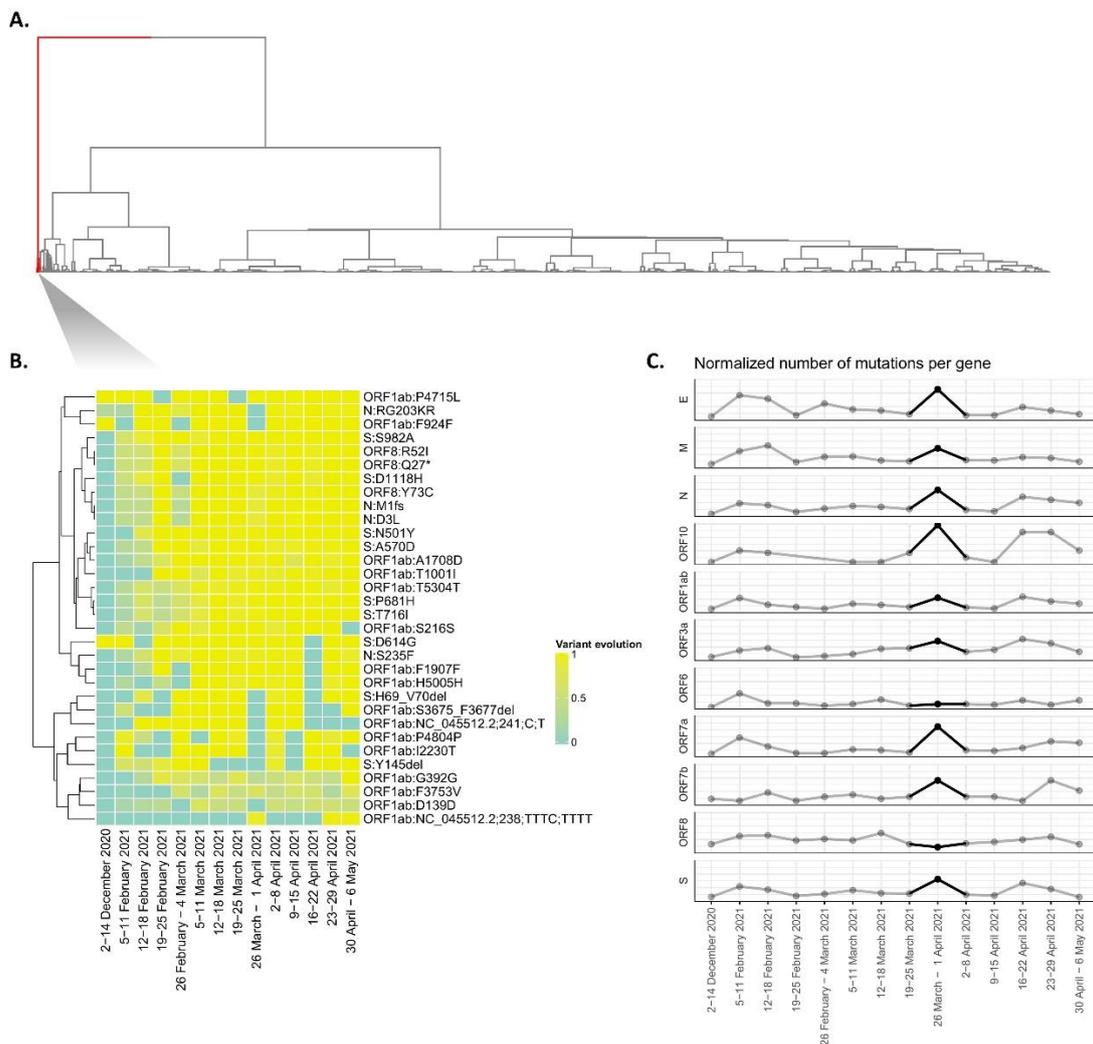
156 Moreover, for the same replicates, 6 different pairs of VCF files were generated in order to investigate
157 how does the number of variants located changes while increasing the low variant frequency
158 parameter of 0.01 that was used at *freebayes* tool during the raw data analysis. The thresholds that
159 were chosen were 0.01, 0.1, 0.2, 0.3, 0.4 and 0.5. In **Figure 1B**, the evolution of the number of variants
160 found is shown while **Figure 1C** gives the corresponding allele frequency.

161

162 **Mutational load detection**

163 The proposed methodology was applied on the wastewater samples, across fourteen time periods, as
 164 show in Supplementary **Table 1**.

165
 166 To this end, **Table 5** was collapsed at the gene and amino acid substitution level, therefore reducing
 167 any data noise that is introduced by nucleotide mutations that correspond to the same amino acid
 168 change. In **Figure 2A** all variants were clustered accordingly based on Euclidean distance, while **Figure**
 169 **2B** highlights specific variants (cluster 1) that exhibit significant difference in behavior.
 170



171
 172 **Figure 2:** Unsupervised variant clustering was performed on a table containing all amino acid
 173 substitutions **A.** Hierarchical clustering shows the total variants using the Euclidean distance as a
 174 distance metric and ward.D as a clustering method. **B.** Hierarchical clustering based on the cluster 1
 175 of the Fig.2A. The heatmap shows the variant evolution across the different periods. **C.** Number of
 176 mutations per gene across the different periods. The values of the plot were normalized based on the
 177 length of each gene.

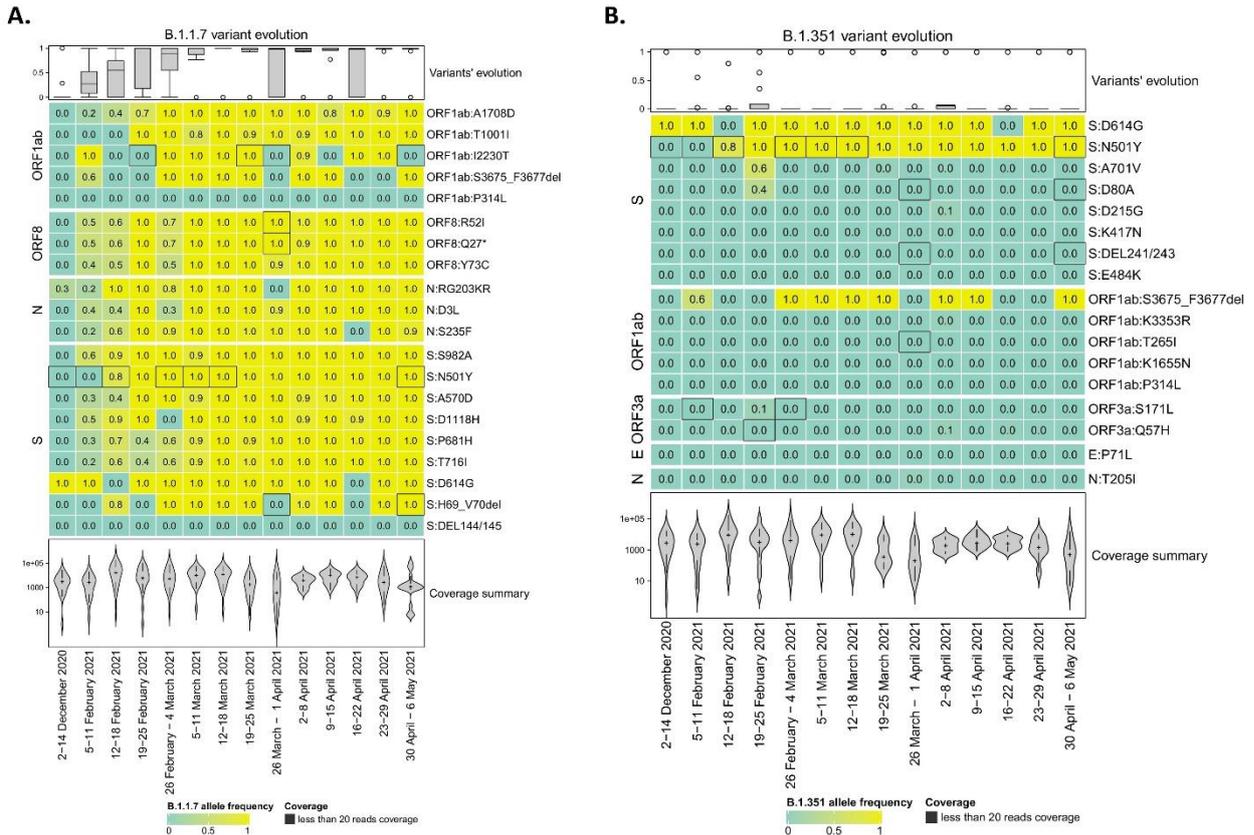
178 Moreover, we studied the number of amino acid variants per gene over each time period. In **Figure**
 179 **2C** the number of variants were clustered based on Euclidean distance, after all count values were
 180 normalized by each gene's length.

181

182 **Detection of Variants of Concern / Variants of Interest**

183 Furthermore, the evolution of lineage-specific variants over the eleven time periods was studied. To
 184 this end, two Variants of Concern were selected; the B.1.1.7 (UK) and the B.1.351 (South Africa)
 185 lineages which are shown in **Figure 3**.

186



187 **Figure 3:** Clustering amino acid substitutions for B.1.1.7 (UK lineage) and B.1.351 (South Africa lineage)
 188 strains. Heatmap displays the corresponding allele frequency (AF) of each period per amino acid
 189 variant. **A.** Evolution of B.1.1.7-detected mutations (UK Lineage) **B.** Evolution of B.1.351-detected
 190 mutations (South Africa lineage). Positions with low coverage (less than 20 reads) are depicted with
 191 dark gray color.

192

193 Moreover, we examined the prevalence rate of the two variants of concern by calculating the average
 194 allele frequency of their mutations. The results for the 14 time periods are presented in **Table 6** and
 195 show that it does not lead to 100% sum per time period. The latter is caused, due to overlapping
 196 mutations among the lineages and implies that the average value of all the mutations is not a reliable
 197 metric to characterize a specific lineage's presence.

198

200 **Table 6:** Allele frequency metrics computed for the comparison with the clinical data. To this end the
 201 average allele frequency of all mutations, the average allele frequency of the unique mutations and
 202 the minimum allele frequency were calculated for each time period.

	Average allele frequency of the mutation		Average allele frequency of the unique mutation		Minimum allele frequency of the present (non-zero) unique mutation	
	B.1.1.7	B.1.351	B.1.1.7	B.1.351	B.1.1.7	B.1.351
2 – 14 Dec. 2020	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
5 – 11 Feb. 2021	32.71%	0.00%	35.15%	0.00%	16.67%	0.00%
12 – 18 Feb. 2021	51.92%	9.05%	49.20%	0.21%	39.53%	1.47%
19 – 25 Feb. 2021	79.78%	22.15%	85.84%	14.20%	58.06%	35.48%
26 Feb. – 4 Mar. 2021	78.18%	11.11%	77.51%	0.00%	31.88%	0.00%
5 – 11 Mar. 2021	95.54%	11.11%	96.28%	0.00%	76.71%	0.00%
12 – 18 Mar. 2021	94.12%	11.11%	92.31%	0.00%	100.0%	0.00%
19 – 25 Mar. 2021	92.58%	11.41%	91.14%	0.56%	92.75%	3.92%
26 Mar. – 1 Apr. 2021	74.09%	11.51%	74.27%	0.63%	89.49%	4.44%
2 – 8 Apr. 2021	96.89%	12.30%	95.98%	1.56%	92.14%	2.03%
9 – 15 Apr. 2021	86.20%	11.10%	81.99%	0.00%	77.15%	0.00%
16 - 22 Apr. 2021	74.32%	10.97%	74.50%	0.00%	89.22%	0.00%
23 - 29 Apr. 2021	81.60%	11.11%	83.81%	0.00%	93.41%	0.00%
30 Apr. - 6 May 2021	93.54%	11.11%	91.55%	0.00%	93.85%	0.00%

203
 204 Further in the analysis, we attempted to examine the minimal lineage support per time period. In this
 205 step we calculated the average allele frequency of the mutations unique to each variant of concern
 206 and produced a table which is more comparable with the clinical data metric. In this context, the
 207 minimum allele frequency of the present (greater than zero) mutations indicates the level of our
 208 confidence. The results of the two metrics are shown in **Table 6**.

209 **Assessing lineage assignment**

210 Qualitative assessment between the major lineages found in the targeted clinical samples, retrieved
 211 from the ENA study ID: PRJEB44141 (ERP128154), and across the fourteen time periods. The clinical
 212 samples exhibit the lineage distribution shown in **Table 7**.

213
 214 **Table 7:** Lineage distribution of samples across all periods. The prevalent strain is depicted with bold
 215 font

Period	No. of patients / period	Lineages Detected	Percentage (%)
2 – 14 Dec. 2020	3	B.1.177	66.7
		B.1.1.189	33.3
		B.1.1.7	43.9
5 – 11 Feb. 2021	57	B.1.177	54.4
		B.1.351	1.75

		B.1.1.7	53.1
12 – 18 Feb. 2021	32	B.1.177	37.5
		B.1.351	9.38
19 – 25 Feb. 2021	21	B.1.1.7	81
		B.1.177	19
		B.1	0.935
		B.1.1.7	83.2
26 Feb. – 4 Mar. 2021	214	B.1.177	11.7
		B.1.258	2.8
		B.1.351	1.4
5 – 11 Mar. 2021	13	B.1.1.7	100
12 – 18 Mar. 2021	22	B.1.1.7	100
19 – 25 Mar. 2021	5	B.1.1.7	100
26 Mar. – 1 Apr. 2021	3	B.1.1.7	100
2 – 8 Apr. 2021	3	B.1.1.7	100
9 – 15 Apr. 2021	35	B.1.1.7	100
16 – 22 Apr. 2021	32	B.1.1	3.12
		B.1.1.7	96.9
23 – 29 Apr. 2021	66	B.1.1.318	1.52
		B.1.1.7	98.5
30 Apr. – 6 May 2021	45	B.1.1.7	100

216

217 Based on the clinical results, as derived from targeted and non-randomized sampling of the
 218 Thessaloniki area, we can perform a direct comparison (**Figure 4**) between the level of presence in a
 219 particular variant of concern (B.1.1.7) in clinical and wastewater data.

220

B.1.1.7	0.022619	0.634453	0.051619	29.74789	0.482472	12 – 18 Feb. 2021
B.1.177	0	0.576864	0.046912	14.04043	0.349998	
B.1.1.7	0.029478	0.633467	0.059038	16.58363	0.505483	19 – 25 Feb. 2021
B.1.177	0	0.563693	0.02373	17.74751	0.508294	
B.1.1.7	0.02381	0.628576	0.051367	17.47607	0.487571	26 Feb. – 4 Mar. 2021
B.1.177	0	0.571146	0.009646	11.87439	0.337707	
B.1.1.7	0.022619	0.630203	0.054135	38.854	0.514036	5 – 11 Mar. 2021
B.1.177	0	0.558882	0.022893	30.16625	0.512106	
B.1.1.7	0.022619	0.627289	0.071746	24.92428	0.557327	12 – 18 Mar. 2021
B.1.177	0	0.563743	0.018941	18.61508	0.446698	
B.1.1.7	0.027891	0.633067	0.05223	469.058	0.516416	19 – 25 Mar. 2021
B.1.177	0	0.568151	0.025152	762.8515	0.47739	
B.1.1.7	0.041931	0.645114	0.122677	1487.607	0.52834	26 Mar. – 1 Apr. 2021
B.1.177	0	0.559649	0.047282	1920.774	0.55107	
B.1.1.7	0.05285	0.611235	0.111091	938.4528	0.515389	2 – 8 Apr. 2021
B.1.177	0	0.5658	0.020531	1172.69	0.527678	
B.1.1.7	0.022619	0.64304	0.085195	1750.538	0.400779	9 – 15 Apr. 2021
B.1.177	0	0.557766	0.01395	1538.826	0.445925	
B.1.1.7	0.257666003	0.724088929	0.067766428	28.4705103	0.455354099	16 – 22 Apr. 2021
B.1.177	0	0.549734606	0.028445651	43.86612546	0.563500191	
B.1.1.7	0.132953798	0.688224115	0.066440254	33.76142848	0.495087239	23 – 29 Apr. 2021
B.1.177	0.000575705	0.557364992	0.036746317	47.34227275	0.51233052	
B.1.177	0	0.550166357	0.02576803	24.96695208	0.429484878	30 Apr. – 6 May 2021
B.1.1.7	0	0.644013622	0.026942757	43.19821767	0.645142425	

235 Discussion

236 Analyzing wastewater, i.e. used water that goes through the drainage system to a treatment facility,
237 is a way that researchers and surveillance systems can track pathogens, such as SARS-CoV-2, or
238 biomarkers that are excreted in urine or feces. Monitoring effluents could be a reliable and more
239 effective tool to estimate SARS-CoV-2 spread compared to sampling and testing the population,
240 because wastewater surveillance can account for those who have not been tested and have only mild
241 or no symptoms. Moreover, an effective and reliable methodology able to detect viral load and SARS-
242 CoV-2 variants from municipal wastewater samples could drastically, or at least help, decrease the
243 cost of virus variant detection in the general population based on whole genome sequencing, since
244 only a few samples must be processed and analyzed.

245
246 In this manuscript we present and validate a methodology named *lineagespot*, making use of NGS
247 data, able to detect lineages and mutational load of SARS-CoV2. The methodology aims to aid the
248 epidemiological system for the monitoring of COVID-19 pandemic in urban areas.

249
250 The method has been tested in different time point samples taken from the main Municipal
251 Wastewater Treatment Plant of Thessaloniki - Greece, where effluents from approx. 750,000
252 inhabitants are collected. The *lineagespot* method demonstrated to be sensitive enough to identify
253 and quantify differences in the mutational load, across various time points and was capable of
254 recording the evolution of the B.1.1.7 lineage in the community. Moreover, the quantitative data
255 obtained using *lineagespot* are in accordance with the trends of well-known mutations (such as
256 Asp614Gly) in the same period with the overall epidemiological status of the municipal area. The
257 application of *lineagespot* in such complex samples, like those from Wastewater Treatment Plants,
258 was able to assign lineages and in agreement with the trend of the major lineages detected in the area
259 of Thessaloniki, in the fourteen time points by whole genome sequencing of samples from the general
260 population.

261
262 Overall, the method developed herein was proven superior compared to other methodologies (Sanger
263 sequencing)^{4,7}, since it was more informative and sensitive enough to detect mutations with low
264 frequency and able to assign with good approximation the correct lineage present in the municipality.

265 Methods

266 Sampling and isolation

267 Wastewater samples were collected from the entrance of the main Municipal Wastewater Treatment
268 Plant of the city which accommodates sewerage of about 750,000 inhabitants. Wastewater entering
269 this plant refers exclusively to citizens from urban districts of the city. Typical values of certain
270 physicochemical properties of wastewater samples tested in this study are displayed in **Table 9**. These
271 properties demonstrate, among others, the existence of suspended solids, dissolved organic matter,
272 dissolved oxygen and salinity that may have strong impact on viral adsorption and decay because of

273 oxidation and increased metabolic activity of microorganisms in wastewater. The residence time of
 274 wastewater until the entrance of the Plant is between 2 and 7 hours (depending on the area), which
 275 is more than enough to allow viral adsorption and decay. Identification of mutations may be hindered
 276 by viral adsorption and decay and for this reason the present effort is particularly significant.

277

278 **Table 9:** Main quality characteristics of wastewater samples

Parameter/ Sample pe- riod	pH	Electri- cal Conduc- tivity (S/cm)	Total Sus- pended Solids (mg/L)	BOD ₅ (mg/L)	COD (mg/L)	Dis- solved Organic Carbon (mg/L)	UV absorp- tion at 254 nm (1/cm)	Total Nitro- gen (mg/L)	Ammo- nium Nitrogen (mg/L)	Total Phos- phorus (mg/L)	Cop- ies/μl
02-14 Dec. 2021	7.5	8.5	620	385	960	35	0.35	62	28.5	11.5	36
05-11 Feb. 2021	7.8	9.6	930	525	1,250	49	0.4	76	33	11.5	68
12-18 Feb. 2021	7.8	4.6	1,200	650	1,570	44	0.45	95	38	12	53
19-25 Feb. 2021	7.9	3.5	1,225	684	1,610	56	0.49	95	38.2	15.2	82
26 Feb. – 4 Mar. 2021	7.8	2.9	1,225	535	1,383	53.5	0.47	78.5	36.7	12.4	179
5-11 Mar. 2021	7.6	2.8	1,017	540	1,373	50.2	0.47	71.7	36.8	12.1	102
12-18 Mar. 2021	7.8	4.5	852	580	1,285	66.3	0.48	76.4	37.4	11.7	277
19-25 Mar. 2021	7.6	4.1	926	582	1,467	60.8	0.48	79.3	39.6	11.6	467
26 Mar. – 3 Apr. 2021	7.6	3.4	1,095	660	1,708	52.1	0.44	85.9	41.4	12.3	494
2-08 Apr. 2021	7.6	4.2	1,054	667	1,537	52	0.48	88.1	40.5	13.2	498
09-15 Apr. 2021	7.7	4.1	1,025	579	1,464	55.6	0.5	80	32.1	11.4	505

279

280 Sampling and handling of the wastewater samples were performed according to Petala et al³. Briefly,
 281 samples were obtained using a refrigerated autosampler (6712 Teledyne ISCO) programmed to deliver
 282 a 24-hours composite sample by mixing consecutive half-hour samples. Samples were transported to
 283 the lab on ice and were processed immediately. Three 50-mL aliquots of each untreated wastewater
 284 sample were subjected to centrifugation at 4,000 × g for 30 min. Afterwards, a composite sample was
 285 obtained from supernatants and pH was adjusted to 4 using 2 M HCl solution. Then, three aliquots of
 286 40 mL each, were filtered through respective 0.45-μm pore-size, 47-mm diameter electronegative
 287 membranes (HAWP04700; Merck Millipore, Ireland). Each membrane filter was rolled into a Falcon™
 288 15-mL conical centrifuge tube with the top side facing inward, and was subjected to RNA extraction.

289 **RNA extraction and SARS-CoV-2 quantification**

290 Each electronegative membrane was subjected to RNA extraction process based on a phenol-
291 chloroform-method⁸ coupled with magnetic bead binding. The following reagents were added
292 sequentially: a) 900 μ L of guanidinium isothiocyanate-based “Lysis buffer I” [5 M guanidinium
293 isothiocyanate, 25 mM EDTA, 25 mM sodium citrate (pH 7.0), 25 mM phosphate buffer (pH 6.6)]
294 containing 1% N-Lauroylsarcosine, 2% Triton X-100, 2% CTAB and 2% PVP, b) 18 μ L β -mercaptoethanol,
295 c) 300 μ L H₂O. Tubes mixed thoroughly by inversion and were stored at 4 °C for 10-30 min.
296 Subsequently, 1.2 mL of “Lysis buffer II” [prepared by mixing 152.5 gr guanidinium hydrochloride,
297 31.25 mL of 2 M acetate buffer (pH 3.8) and water-saturated phenol stabilized (pH 4), at a final volume
298 of 500 ml] was added, followed by incubation on a horizontal rotator (150 rpm, 10 min, RT). The liquid
299 phase was transferred into a 2-mL microcentrifuge tube, was clarified by centrifugation (21,000 \times *g*, 5
300 min, 4 °C) and 1.6 mL of the supernatant were transferred to a new 2-mL tube, wherein 200 μ L
301 chloroform-isoamyl alcohol (24:1) were added, followed by vigorous shaking for 30 s, incubation (-20
302 °C, 30 min) and centrifugation (21,000 \times *g*, 10 min, 4 °C). The upper aqueous phase (800 μ L) was
303 transferred and mixed with 667 μ L isopropanol and 20 μ L of magnetic beads (IDEXX Water DNA/RNA
304 Magnetic Bead Kit; IDEXX Laboratories Inc., Westbrook, ME, USA), followed by incubation on a
305 horizontal rotator (150 rpm, 15 min, RT). The beads were washed according to the manufacturer’s
306 protocol. RNA was eluted in 100 μ L buffer, and eluates were subjected to filtration, using the
307 OneStep™ PCR Inhibitor Removal Kit (Zymo Research Corporation, Irvine, CA, USA) and were stored
308 at -80 °C. Extracted RNAs originating from 12 processed electronegative membranes and spanning 6
309 different days of sampling were pooled (1.1 mL total RNA extract) and mixed with 2.2 mL binding
310 buffer containing isopropanol (IDEXX Water DNA/RNA Magnetic Bead Kit). Half of the mixture (1.65
311 mL) was incubated with 20 μ L of magnetic beads on a horizontal rotator (150 rpm, 15 min, RT). After
312 the magnetic separation of beads, the supernatant was removed and the procedure was repeated by
313 adding the remaining 1.65 mL of the mixture. The beads were washed according to the manufacturer’s
314 protocol and RNA was eluted in 60 μ L buffer.

315

316 Concentrated RNAs were subjected to real-time RT-PCR testing for SARS-CoV-2 quantification, utilizing
317 the N2 protocol proposed by the Centers for Disease Control and Prevention (CDC) for the diagnosis
318 of COVID-19 in humans (CDC, 2020). The assay was performed on a CFX96 Touch™ Real-Time PCR
319 Detection System (Bio-Rad Laboratories, Hercules, CA, USA). Calibration curves were generated using
320 the synthetic single-stranded RNA standard “EURM-019” (Joint Research Centre, European
321 Commission) and SARS-CoV-2 viral loads were expressed as genome copies per μ L of RNA extract.

322 **Library preparation and sequencing**

323 The targeted sequencing method was applied by preparing 400nt amplicons using the ARTIC v3
324 protocol developed by Wellcome Sanger Institute⁹, with some modifications. First, cDNA synthesis
325 was prepared from 10 μ L of RNA using Super Script II reverse transcriptase (Invitrogen - Thermo Fisher
326 Scientific, USA) and 50 ng/ μ L of random primers according to the protocol guidelines. For subsequent
327 cDNA amplification, 2.5 μ L of the generated cDNA was used instead of 6 μ L, using ARTIC PCR primer
328 pools (v3). Finally, the NEBNext adaptor (New England Biolabs, US, #7335) was used in the ligation
329 reaction, diluted with adaptor dilution buffer at 10 μ M final concentration. All purification steps were

330 performed according to the ARTIC protocol. The samples were paired-end sequenced on a MiSeq
331 platform (Illumina, USA) with a read length of 2×300 bp.

332 **Raw data analysis**

333 The initial phase of the bioinformatics analysis is to produce an alignment of the sequencing reads,
334 while maintaining extremely strict criteria, in order to remove any potential contaminants and/or
335 sequencing errors. The first step is the adaptor removal process, where any adaptors were removed
336 from the raw *fastq* sequences, with the cleaned reads mapped to the SARS-CoV-2 reference genome
337 (Wuhan variant, NC_045512), using minimap2 tool¹⁰ with a minimal chaining score (matching bases
338 minus log gap penalty) equal to 40. From this process, only the paired-end sequences were retained,
339 while any other (unmatched, multiple mappings, etc.) were removed. In the next step, two different
340 computational workflows were employed, corresponding to the two different sequencing protocols
341 employed. For the first seven samples and the last three (2 December 2020 – 18 March 2021 and 16
342 April 2021 – 06 May 2021), the primer sequences are excluded using the *ivar* tool¹¹, setting a minimum
343 of 200 length in nucleotides for a read to be retained after trimming, and a minimum threshold for
344 sliding window of 15 quality to pass (width of sliding window equal to 4). The final sequences are then
345 remapped to the same reference genome (minimal chaining score equal to 40). For the four samples
346 between 19 March 2021 and 15 April 2021 primer trimming and remapping to reference genome was
347 not applied, as the updated protocol used did not necessitate this step. Finally, in order to be able to
348 detect low frequency variants, the *freebayes* variant caller was used with a low variant frequency
349 parameter of 0.01. Ultimately, all identified mutations were annotated using the *SnpEff* tool¹² and the
350 NC_045512.2 (version 5.0) database.

351 **Downstream analysis of lineages detection**

352 In order to identify and assign different SARS-CoV-2 lineages based on the mutations detected from a
353 single wastewater sample, we implemented the proposed methodology in a tool named *lineagespot*.
354 The tool accepts as input a VCF file, which contains all mutations identified in the sample, along with
355 the reference SARS-CoV-2 genome file, and a file containing all lineage-assignment rules, as retrieved
356 from the *pangolin* tool¹³ repository. After analyzing all inputs, a tab-delimited file (TSV file) is produced
357 containing the most probable lineages that have been found. **Figure 5** shows an overview of the tool's
358 functionalities, which can be described in 2 phases:

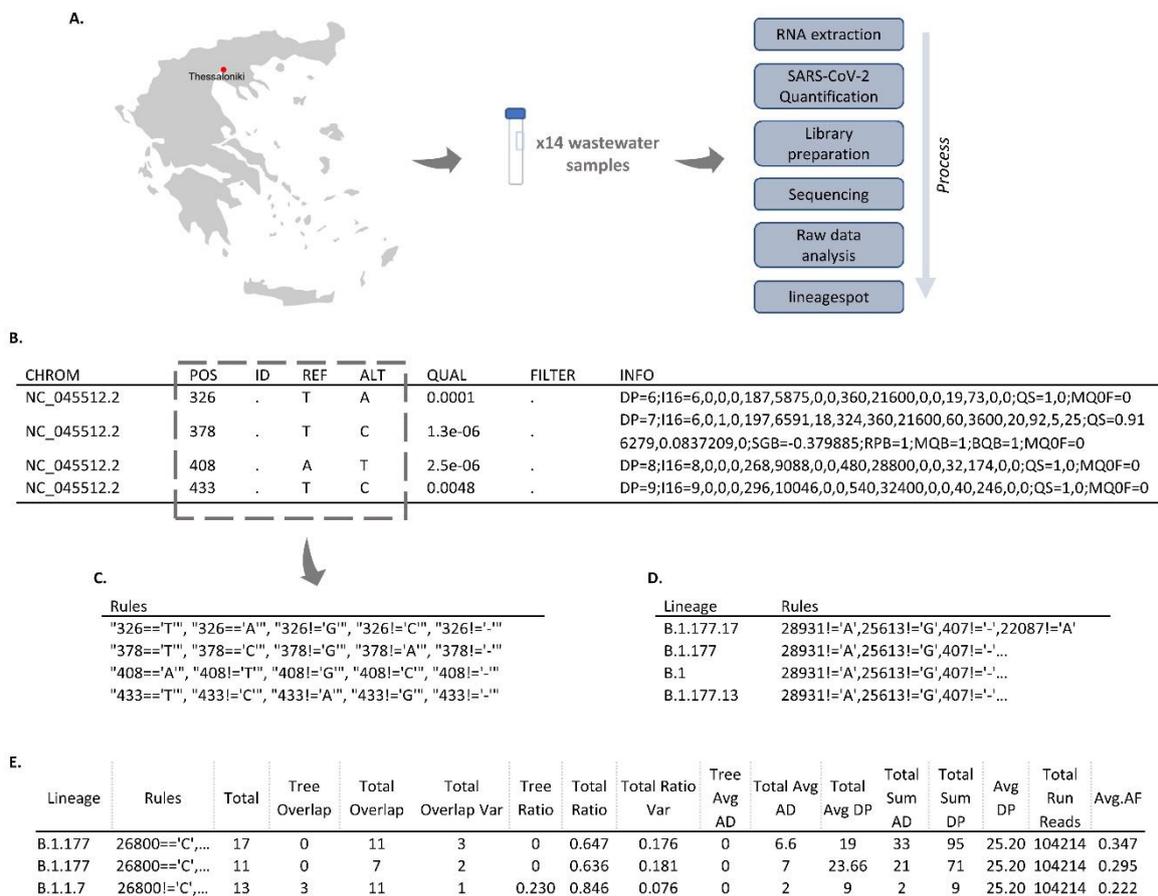
359

360 *i) Creating rules from variants*

361 In this phase all rules that can be derived from the VCF file are constructed. Initially, a vector of all
362 genome positions is created, for which each position is set to be equal to the reference genome
363 nucleotides. Then, the VCF file is read and a set of new rules is formed by setting each position of the
364 file with the reported variant or multiple variants (in case there is more than one reported variant at
365 the same position). It should be noted that positions that have been detected with more than one
366 variant, should include all of them at the VCF's ALT column in a comma-delimited format. Most of the
367 variant caller tools (*freebayes*, *GATK*, etc.) are doing this by default. Finally, positions with reference
368 read depth equal to zero are removed from the first vector. The remaining two vectors are merged
369 into one.

370
371
372
373
374
375
376
377

In addition to finding all positions that need to be set equal to the base that has been allocated, four more rules are added for each genome position. These rules contain all bases *not* equal to the nucleotide of the original rule. For example, if position 5388 is equal to base 'A' (representing rule 5388=='A'), then four new rules are added containing all bases not equal to 'A', e.g., 5388!='T', 5388!='G', 5388!='C', 5388!='-' (where the '-' symbol stands for a gap in the referred sequence). Finally, all rules are merged into a single vector representing this particular lineage.



378
379
380
381
382

Figure 5: Snapshots of the intermediate steps. **A.** A summary plot showing the overall process from the sampling to LineageSpot **B.** A VCF file produced by the chosen variant caller **C.** The rules as they are generated by the VCF file. **D.** Pangolin's decision rules. **E.** A tab-delimited file as produced by *lineagespot*.

383
384
385
386
387
388
389
390

ii) Comparing with pangolin rules

The second phase aims to compare the rules derived from the VCF file with the assignment rules provided by the pangolin tool. Specifically, all decision rules are read from the input pangolin file, and for each lineage, the related rules are compared with the final rule vector.

Three metrics are computed and stored in the output file; the total number of rules leading to the related lineage (N_r), the number of rules satisfied by the created rule vector, considering pangolin's

391 rules as a decision tree (n_d), and the total number of rules satisfied (n_t). Also, the related ratio values
392 are being computed, giving a satisfaction percentage of each lineage:

393

394

$$R_d = \frac{n_d}{N_r}, R_t = \frac{n_t}{N_r}$$

395 Underlying assumptions of the method

396 It should be noted that the methodology relies on the following assumption. Given a group of reads
397 that satisfy a rule A of lineage L, and another group of reads that satisfy rule B from the same lineage
398 L, then the lineage L is incorrectly assigned. As an example, suppose that a group of reads satisfy only
399 the first two rules from lineage B.1.177.17 (28931!= 'A', 25613!= 'G'), and another group of reads satisfy
400 the next two rules from the same lineage (407!= '-', 22087!= 'A'). Based on the method description
401 above, lineage B.1.177.17 will be marked as an identified lineage, even though none of the reads
402 satisfy all of the lineage's rules.

403

404 In order to mitigate this risk, we are taking into consideration a number of different indicators, that
405 reflect the number of total rules satisfied, the number of rules that are satisfied based on the detected
406 mutations, and the overall number of reads that support both reference and allele for each of the
407 rules.

408

409 Source of lineage-specific mutations

410 For the detection of lineage-specific mutations, three different data sources were used; pangolin, VEO
411 and *outbreak.info*. An example of the differences of the three sources is provided in **Supplementary**
412 **Figure 1**, which compares the detection of B.1.1.7 using **A** outbreak.info, **B** pangolin and **C** VEO data.

413 References

- 414 1. Nemudryi, A. *et al.* Temporal detection and phylogenetic assessment of SARS-CoV-2 in municipal
415 wastewater. <http://medrxiv.org/lookup/doi/10.1101/2020.04.15.20066746> (2020)
416 doi:10.1101/2020.04.15.20066746.
- 417 2. Jahn, K. *et al.* Detection of SARS-CoV-2 variants in Switzerland by genomic analysis of
418 wastewater samples. <http://medrxiv.org/lookup/doi/10.1101/2021.01.08.21249379> (2021)
419 doi:10.1101/2021.01.08.21249379.
- 420 3. Petala, M. *et al.* A physicochemical model for rationalizing SARS-CoV-2 concentration in sewage.
421 Case study: The city of Thessaloniki in Greece. *Science of The Total Environment* **755**, 142855
422 (2021).

- 423 4. Martin, J. *et al.* Tracking SARS-CoV-2 in Sewage: Evidence of Changes in Virus Variant
424 Predominance during COVID-19 Pandemic. *Viruses* **12**, 1144 (2020).
- 425 5. Crits-Christoph, A. *et al.* Genome Sequencing of Sewage Detects Regionally Prevalent SARS-CoV-
426 2 Variants. *mBio* **12**, e02703-20, /mbio/12/1/mBio.02703-20.atom (2021).
- 427 6. Izquierdo-Lara, R. *et al.* Monitoring SARS-CoV-2 circulation and diversity through community
428 wastewater sequencing. <http://medrxiv.org/lookup/doi/10.1101/2020.09.21.20198838> (2020)
429 doi:10.1101/2020.09.21.20198838.
- 430 7. Daughton, C. G. Wastewater surveillance for population-wide Covid-19: The present and future.
431 *Science of The Total Environment* **736**, 139631 (2020).
- 432 8. Chaintoutis, S. C., Papadopoulou, E., Melidou, A., Papa, A. & Dovas, C. I. A PCR-based NGS
433 protocol for whole genome sequencing of West Nile virus lineage 2 directly from biological
434 specimens. *Molecular and Cellular Probes* **46**, 101412 (2019).
- 435 9. Pipelines, D. *et al.* COVID-19 ARTIC v3 Illumina library construction and sequencing protocol v4
436 (protocols.io.bgxjxkn). doi:10.17504/protocols.io.bgxjxkn.
- 437 10. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100
438 (2018).
- 439 11. Grubaugh, N. D. *et al.* An amplicon-based sequencing framework for accurately measuring
440 intrahost virus diversity using PrimalSeq and iVar. *Genome Biol* **20**, 8 (2019).
- 441 12. Cingolani, P. *et al.* Using *Drosophila melanogaster* as a Model for Genotoxic Chemical Mutational
442 Studies with a New Program, SnpSift. *Front. Gene.* **3**, (2012).
- 443 13. pangolin; *Phylogenetic Assignment of Named Global Outbreak LINEages*.

444 Acknowledgements / CRediT author statement

445 This work was supported by the “Greece vs Corona: Flagship Action to address the SARS-CoV-2 crisis.
446 *Epidemiological study in Greece through extensive testing for virus and antibody detection, viral*
447 *genome sequencing and genetic analysis of patients”* project, which is funded by the General
448 Secretariat for Research and Innovation, under the Public Investments Program (PIP). Support was
449 also received by the Region of Central Macedonia through the project on “*Epidemiological status of*
450 *COVID-19 disease based on viral loading monitoring in wastewater”*, as well as through the “SARS-

451 *CoV-2 RNA monitoring in untreated wastewater*” Initiative of the Hellenic National Public Health
452 Organization. Finally, this work was supported by *ELIXIR*, the research infrastructure for life-science
453 data.

454

455 NPe, MT, MCM and AT developed the *lineagespot* code and performed the downstream analysis. NPe
456 and MT performed the analysis of the raw NGS data. EM, SL and EV did the library preparation. CID
457 and AC provided the data and reviewed the submitted version. CID, TC, SCC, MP and MK participated
458 in experimental investigation. CID, AP, NPa, FP and TK participated in project conceptualization,
459 management and funding. FP and AA designed, supervised the study, and reviewed the manuscript.
460 All authors contributed to the article and approved the submitted version.

461 Supplementary Material

462 **S1: Code and data**

463 The implemented code that produces the results of this paper, starting from the VCF files, is available
464 on the GitHub repository: <https://github.com/BiodataAnalysisGroup/lineagespot>.

465

466 All raw FASTQ files are deposited on ENA: Project IDs **PRJEB44141** (for *patient samples*) and
467 **PRJEB44548** (for *wastewater samples*).

468

469 **Supplementary Figure legend**

470

471 **Supplementary Figure 1:** Detected mutations of B.1.1.7 (UK Lineage) using data provided by **A.**
472 outbreak.info, **B.** pangolin and **C.** VEO.

473

474 **Supplementary Table caption**

475

476 **Supplementary Table 1:** Wastewater samples used for the SARS-CoV-2 variant analysis.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplTable1.pdf](#)
- [SupplementaryFigure1.jpg](#)