

COVID-19 subject UPHS-0066

2021-05-05

The table below provides a summary of subject samples for which sequencing data is available.

The experiments column shows the number of sequencing experiments performed for each specimen.

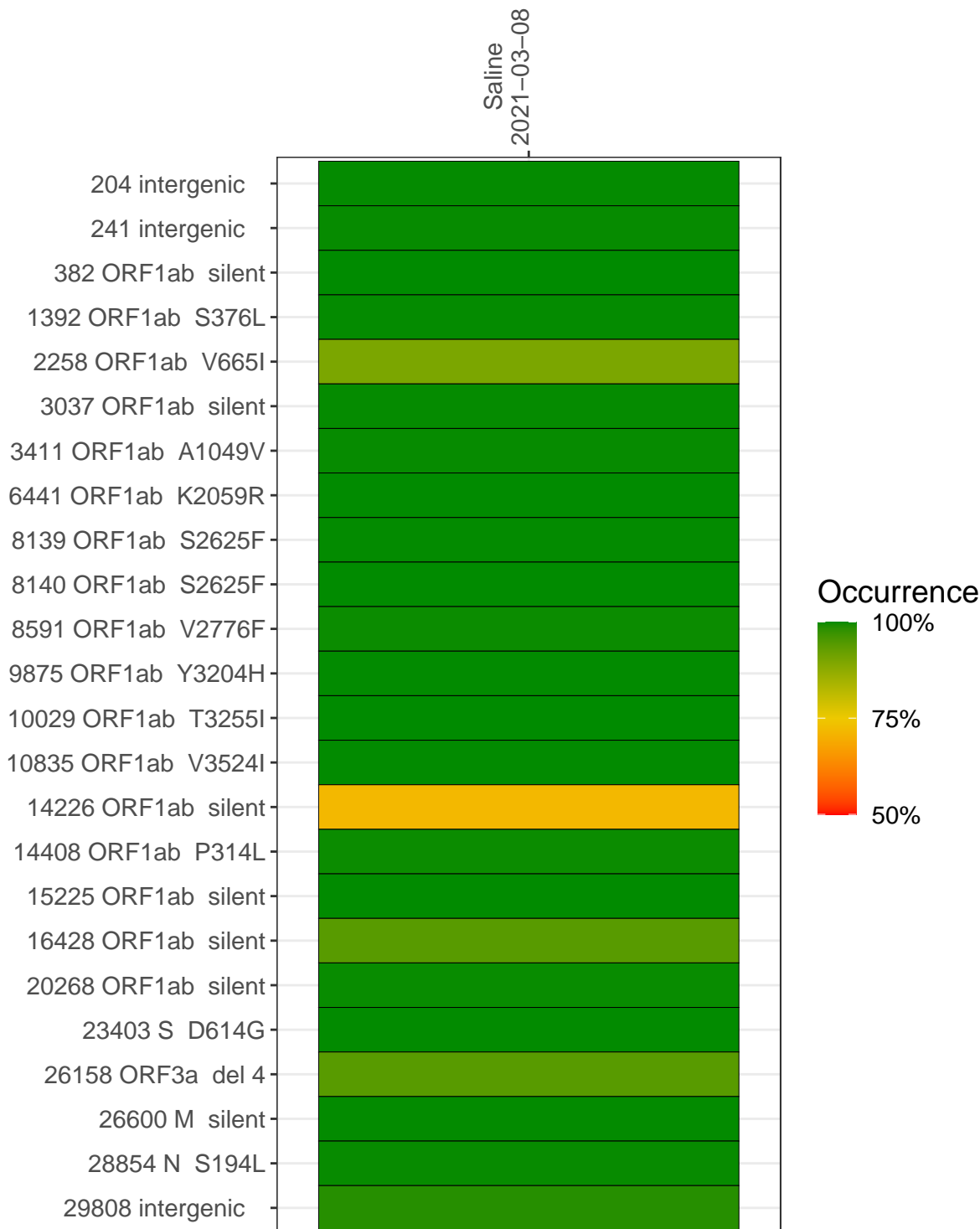
Experiment specific analyses are shown at the end of this report. Lineages are called with the Pangolin software tool (Rambaut et al 2020) for genomes with $> 90\%$ sequence coverage.

Table 1. Sample summary.

Experiment	Type	Genomes	Sample type	Sample date	Largest contig (KD)	Lineage	Reference read coverage	Reference read coverage (≥ 5 reads)
VSP0998-1	single experiment	NA	Saline	2021-03-08	29.88	B.1.234	99.9%	99.7%

Variants shared across samples

The heat map below shows how variants (reference genome /home/everett/projects/SARS-CoV-2-Philadelphia/Wuhan-Hu-1) are shared across subject samples where the percent variance is colored. Variants are called if a variant position is covered by 5 or more reads, the alternative base is found in > 50% of read pairs and the variant yields a PHRED score > 20. Gray tiles denote positions where the variant was not the major variant or no variants were found. The relative base compositions of each experiment used to calculate tiles are shown in the following plot where the total number of position reads are shown atop of each plot.



Saline
2021-03-08

204 intergenic	4915
241 intergenic	3331
382 ORF1ab silent	3482
1392 ORF1ab S376L	4210
2258 ORF1ab V665I	4126
3037 ORF1ab silent	5881
3411 ORF1ab A1049V	7271
6441 ORF1ab K2059R	13045
8139 ORF1ab S2625F	7597
8140 ORF1ab S2625F	7563
8591 ORF1ab V2776F	4627
9875 ORF1ab Y3204H	3387
10029 ORF1ab T3255I	3414
10835 ORF1ab V3524I	9949
14226 ORF1ab silent	9858
14408 ORF1ab P314L	5929
15225 ORF1ab silent	9845
16428 ORF1ab silent	15480
20268 ORF1ab silent	2514
23403 S D614G	11392
26158 ORF3a del 4	8883
26600 M silent	6019
28854 N S194L	1265
29808 intergenic	140

Base change

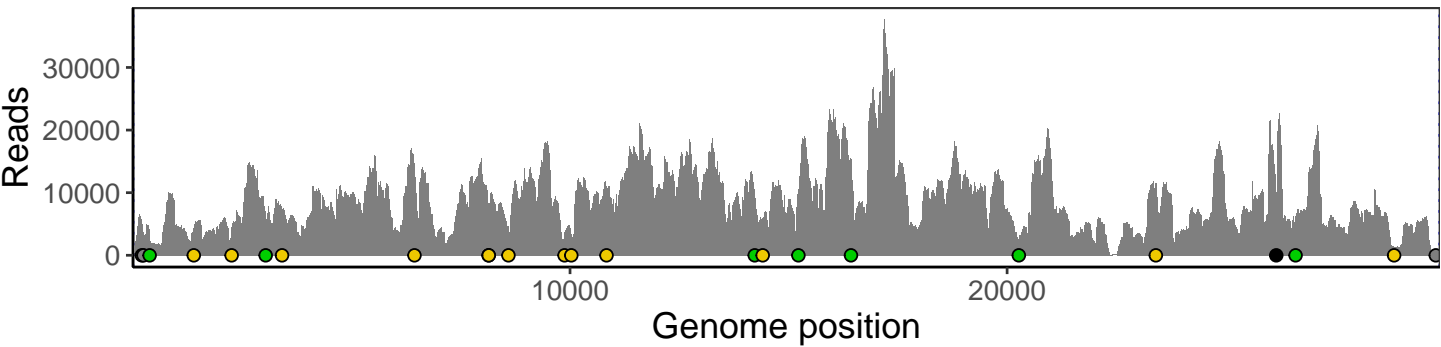


VSP0998-1

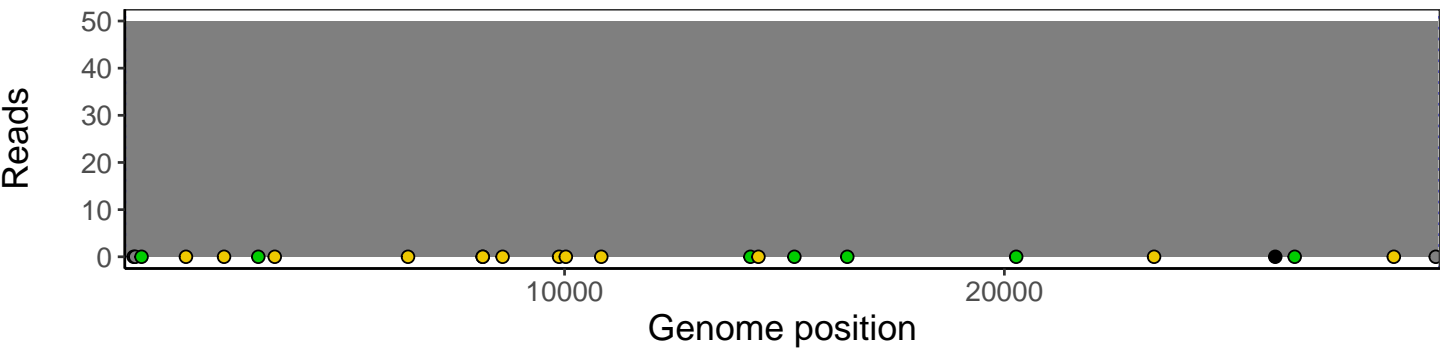
Analyses of individual experiments and composite results

VSP0998-1 | 2021-03-08 | Saline | UPHS-0066 | genomes | single experiment

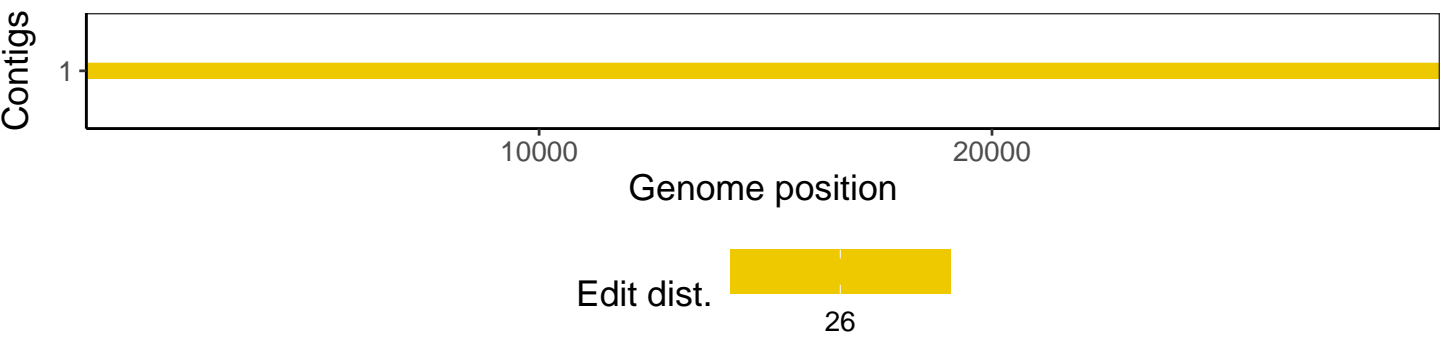
The plot below shows the number of reads covering each nucleotide position in the reference genome. Variants are shown as colored dots along the bottom of the plot and are color coded according to variant types: gray - transgenic, green - silent, gold - missense, red - nonsense, black - indel.



Excerpt from plot above focusing on reads coverage from 0 to 50 NT.



The longest five assembled contigs are shown below colored by their edit distance to the reference genome.



Software environment

Software/R package	Version
R	3.4.0
bwa	0.7.17-r1198-dirty
samtools	1.10 Using htlib 1.10
bcftools	1.10.2-34-g1a12af0-dirty Using htlib 1.10.2-57-gf58a6f3
pangolin	2.3.8
genbankr	1.4.0
optparse	1.6.0
forcats	0.3.0
stringr	1.4.0
dplyr	0.8.1
purrr	0.2.5
readr	1.1.1
tidyr	0.8.1
tibble	2.1.2
ggplot2	3.0.0
tidyverse	1.2.1
ShortRead	1.34.2
GenomicAlignments	1.12.2
SummarizedExperiment	1.6.5
DelayedArray	0.2.7
matrixStats	0.54.0
Biobase	2.36.2
Rsamtools	1.28.0
GenomicRanges	1.28.6
GenomeInfoDb	1.12.3
Biostrings	2.44.2
XVector	0.16.0
IRanges	2.10.5
S4Vectors	0.14.7
BiocParallel	1.10.1
BiocGenerics	0.22.1