

COVID-19 subject UPHS-0702

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.

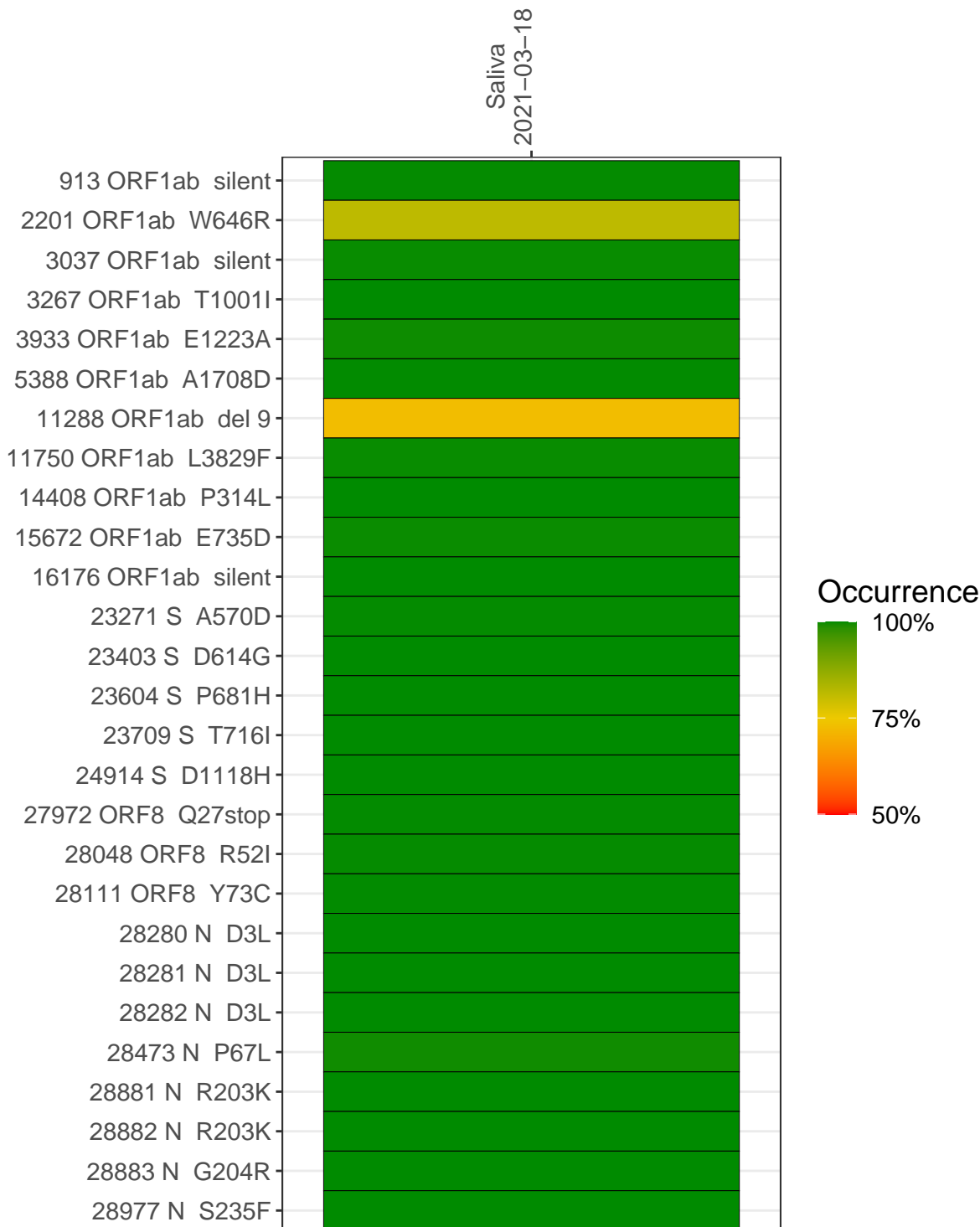
Experiment specific analyses are shown at the end of this report. Lineages are called with the Pangolin software tool (Rambaut et al 2020) for genomes with $> 90\%$ sequence coverage.

Table 1. Sample summary.

Experiment	Type	Genomes	Sample type	Sample date	Largest contig (KD)	Lineage	Reference read coverage	Reference read coverage (≥ 5 reads)
VSP1920-1	single experiment	NA	Saliva	2021-03-18	5.62	NA	90.7%	74.2%

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.



Saliva
2021-03-18

913 ORF1ab silent	2340
2201 ORF1ab W646R	735
3037 ORF1ab silent	1464
3267 ORF1ab T1001I	1277
3933 ORF1ab E1223A	1583
5388 ORF1ab A1708D	2237
11288 ORF1ab del 9	614
11750 ORF1ab L3829F	996
14408 ORF1ab P314L	1823
15672 ORF1ab E735D	1296
16176 ORF1ab silent	1756
23271 S A570D	1026
23403 S D614G	1112
23604 S P681H	2288
23709 S T716I	1525
24914 S D1118H	965
27972 ORF8 Q27stop	2812
28048 ORF8 R52I	1759
28111 ORF8 Y73C	1809
28280 N D3L	893
28281 N D3L	893
28282 N D3L	945
28473 N P67L	1215
28881 N R203K	277
28882 N R203K	275
28883 N G204R	275
28977 N S235F	327

Base change

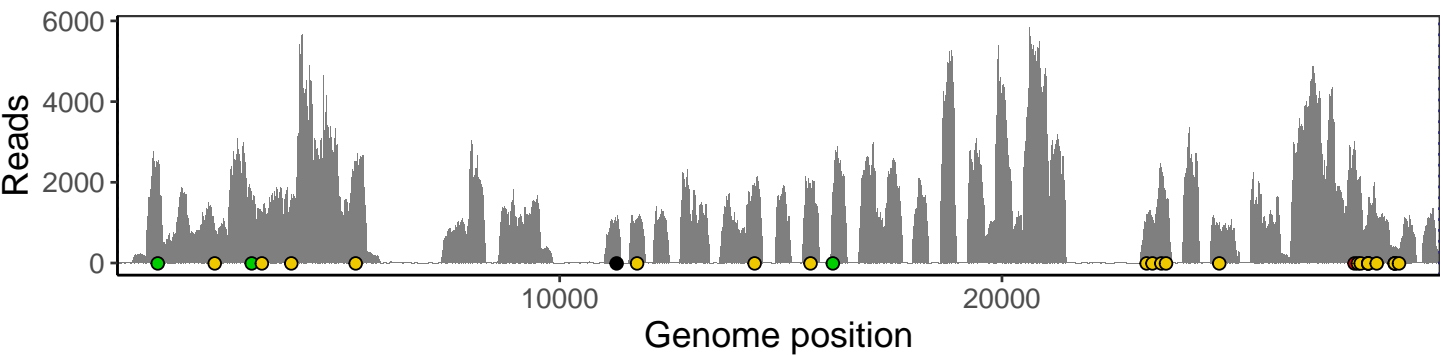
- Expected
- A
- T
- C
- G
- N
- Ins/Del
- No data

VSP1920-1

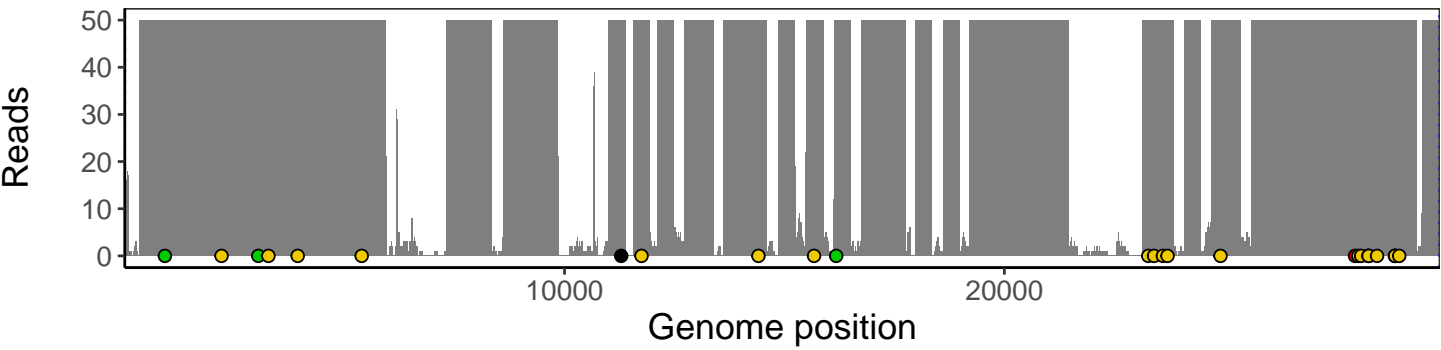
Analyses of individual experiments and composite results

VSP1920-1 | 2021-03-18 | Saliva | UPHS-0702 | genomes | single experiment

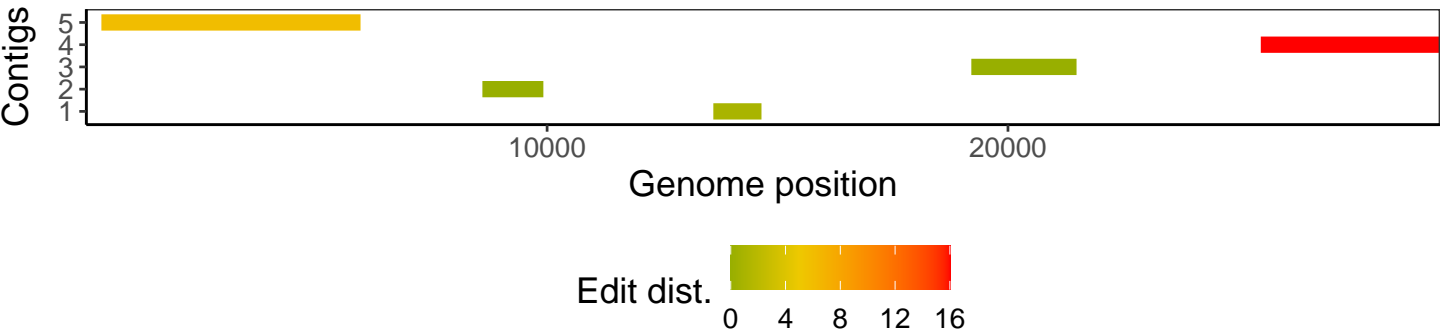
The plot below shows the number of reads covering each nucleotide position in the reference genome. Variants are shown as colored dots along the bottom of the plot and are color coded according by variant types: gray - transgenic, green - silent, gold - missense, red - nonsense, black - indel.



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