

COVID-19 subject HUP PH-0022

2021-05-05

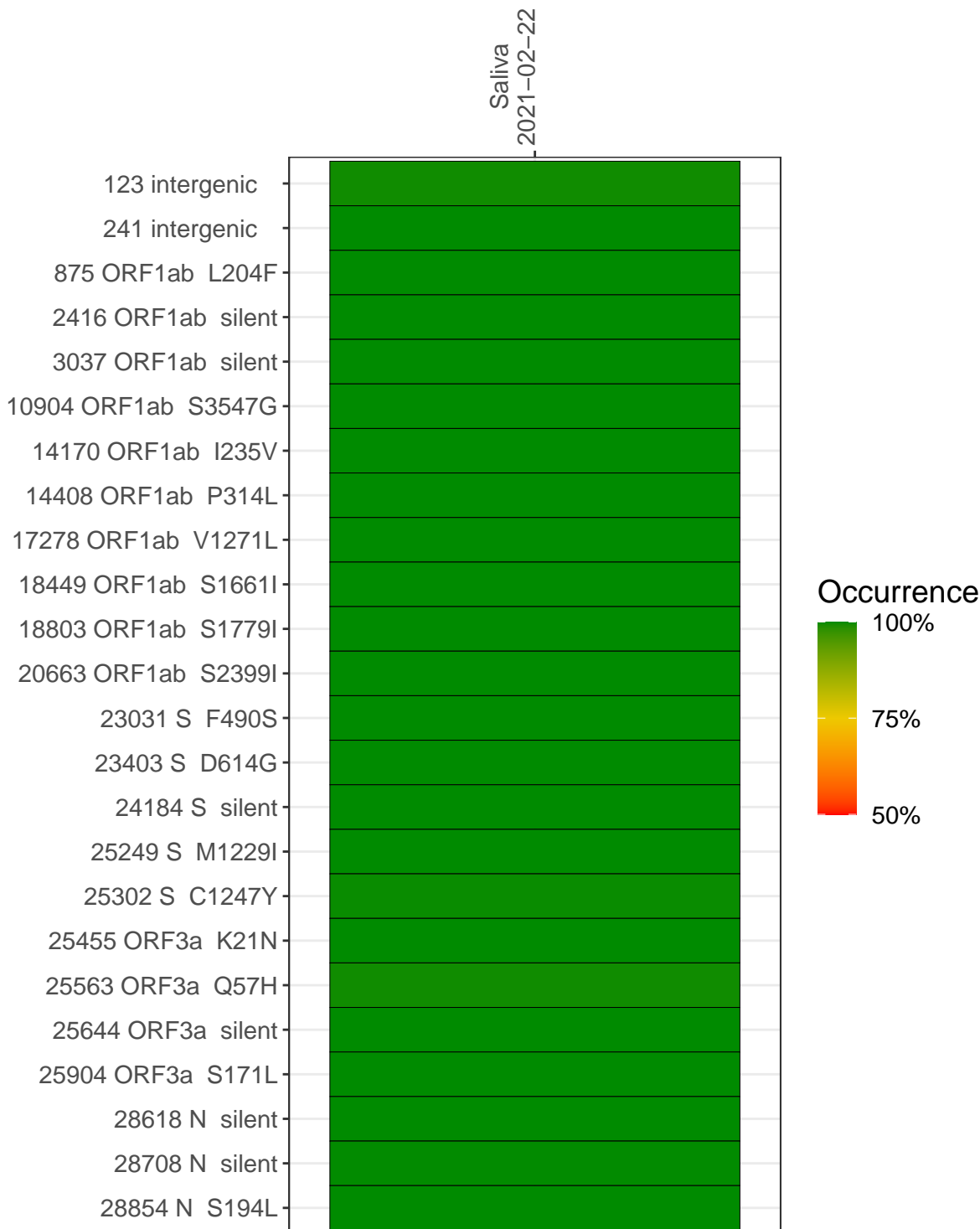
The table below provides a summary of subject samples for which sequencing data is available. The experiments column shows the number of sequencing experiments performed for each specimen. Experiment specific analyses are shown at the end of this report. Lineages are called with the Pangolin software 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)
VSP0866-1	single experiment	NA	Saliva	2021-02-22	13.31	B.1.480	99.8%	97.8%

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
2021-02-22

123 intergenic	241
241 intergenic	151
875 ORF1ab L204F	476
2416 ORF1ab silent	478
3037 ORF1ab silent	496
10904 ORF1ab S3547G	761
14170 ORF1ab I235V	875
14408 ORF1ab P314L	1115
17278 ORF1ab V1271L	1848
18449 ORF1ab S1661I	89
18803 ORF1ab S1779I	1788
20663 ORF1ab S2399I	1681
23031 S F490S	342
23403 S D614G	134
24184 S silent	766
25249 S M1229I	32
25302 S C1247Y	437
25455 ORF3a K21N	772
25563 ORF3a Q57H	1198
25644 ORF3a silent	331
25904 ORF3a S171L	52
28618 N silent	1597
28708 N silent	1880
28854 N S194L	12

Base change

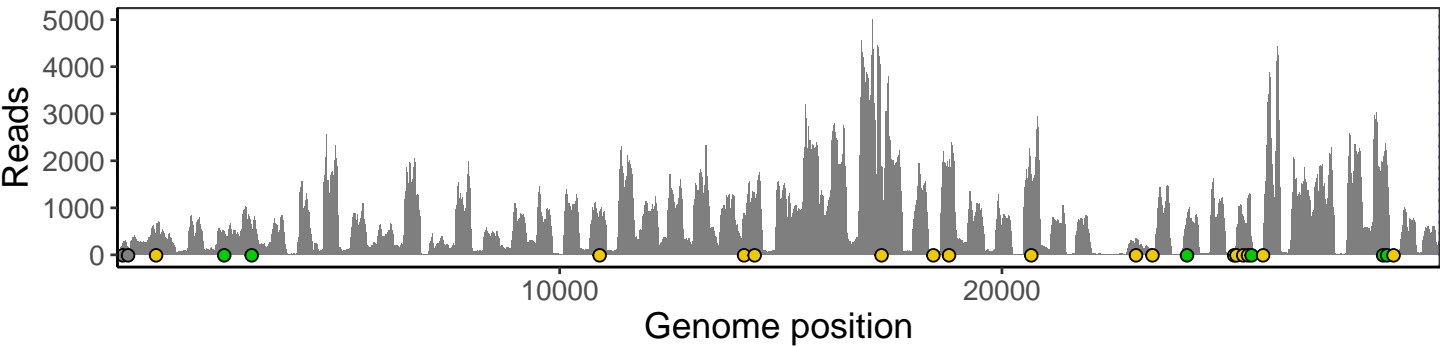


VSP0866-1

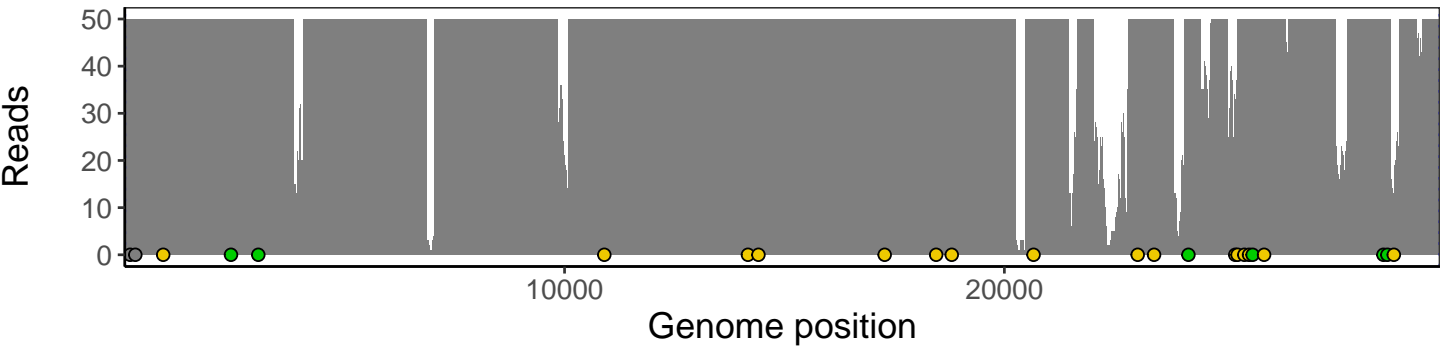
Analyses of individual experiments and composite results

VSP0866-1 | 2021-02-22 | Saliva | HUP-PH-0022 | genomes | single experiment

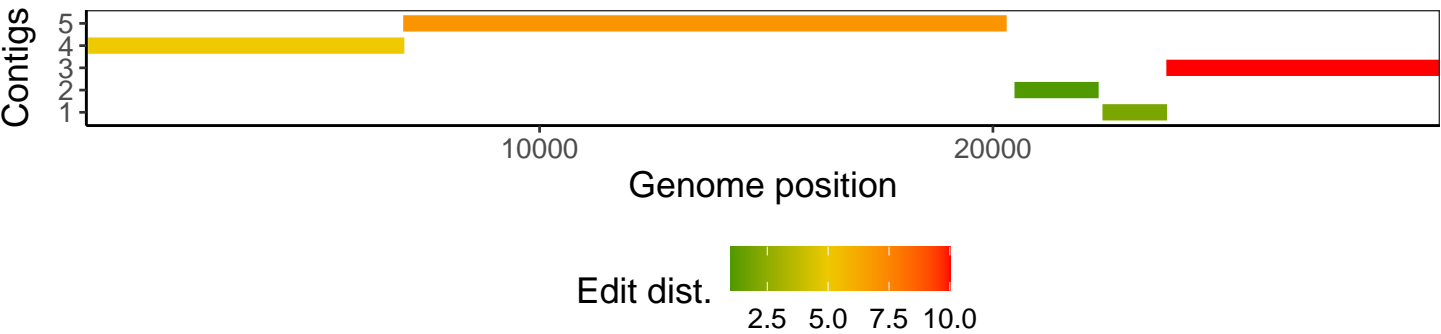
The plot below shows the number of reads covering each nucleotide position in the reference genome. Variants are shown as colored dots along the bottom of the plot and are color coded according by variant types: gray - transgenic, green - silent, gold - missense, red - nonsense, black - indel.



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