# COVID-19 subject 463

*2021-05-21*

The table below provides a summary of subject samples for which sequencing data is available. The experiments column shows the number of seque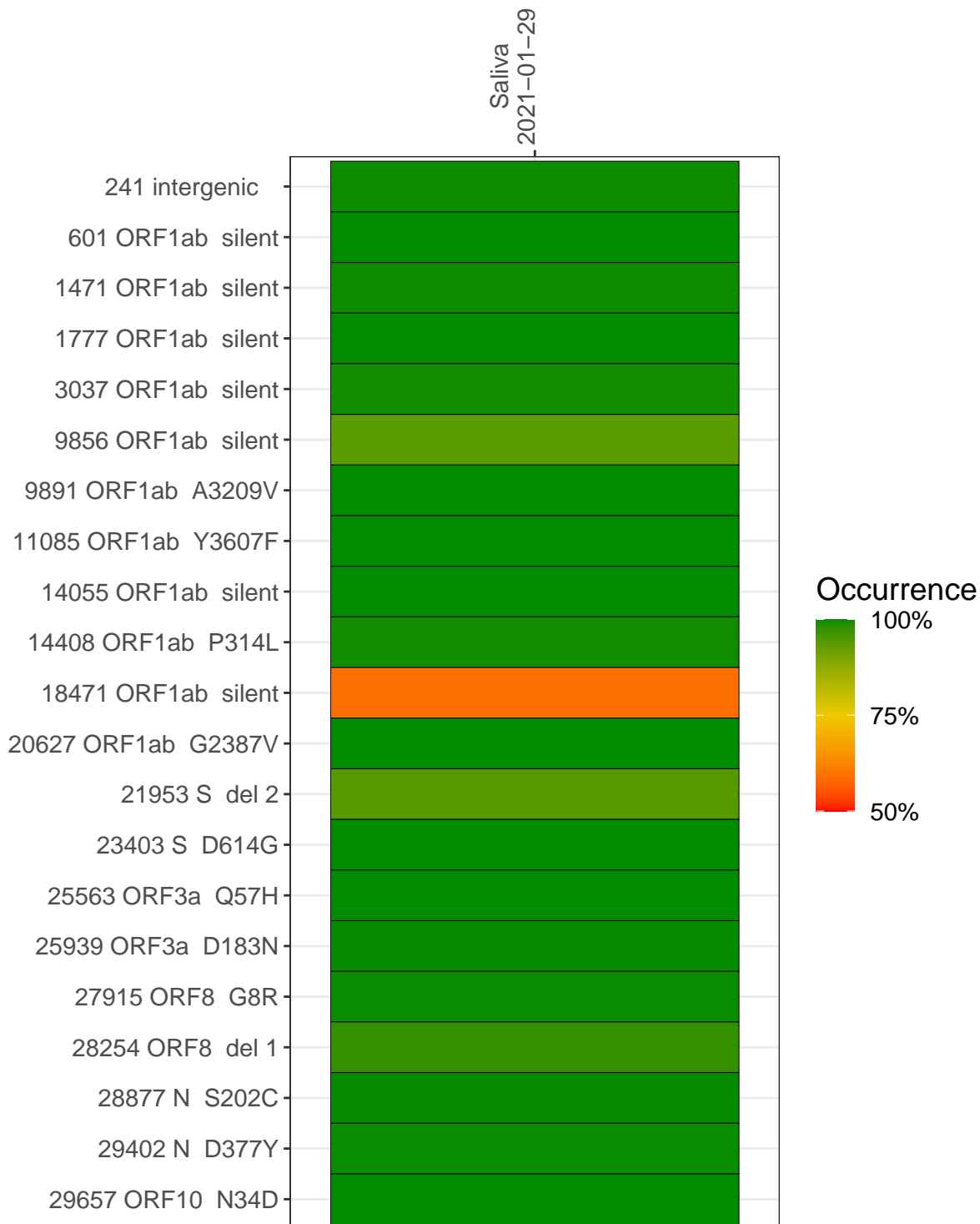ncing experiments performed for each specimen. Experiment specific analyses are shown at the end of this report. Lineages are called with the Pangolin software tool (Rambaut et al 2020) for genomes with > 90% sequence coverage.

Table 1. Sample summary.

| Experiment | Type | Genomes | Sample type | Sample date | Largest contig (KD) | Lineage | Reference read coverage | Reference read coverage (>= 5 reads) |
|---|---|---|---|---|---|---|---|---|
| VSP0789-1 | single experiment | NA | Saliva | 2021-01-29 | 11.16 | NA | 94.9% | 94.5% |

## 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
2021−01−29

| Position | Value |
|----------|-------|
| 241 intergenic | 292 |
| 601 ORF1ab  silent | 37 |
| 1471 ORF1ab  silent | 575 |
| 1777 ORF1ab  silent | 406 |
| 3037 ORF1ab  silent | 416 |
| 9856 ORF1ab  silent | 82 |
| 9891 ORF1ab  A3209V | 91 |
| 11085 ORF1ab  Y3607F | 120 |
| 14055 ORF1ab  silent | 297 |
| 14408 ORF1ab  P314L | 551 |
| 18471 ORF1ab  silent | 436 |
| 20627 ORF1ab  G2387V | 211 |
| 21953 S  del 2 | 74 |
| 23403 S  D614G | 627 |
| 25563 ORF3a  Q57H | 626 |
| 25939 ORF3a  D183N | 599 |
| 27915 ORF8  G8R | 1861 |
| 28254 ORF8  del 1 | 1033 |
| 28877 N  S202C | 674 |
| 29402 N  D377Y | 407 |
| 29657 ORF10  N34D | 401 |

Base change

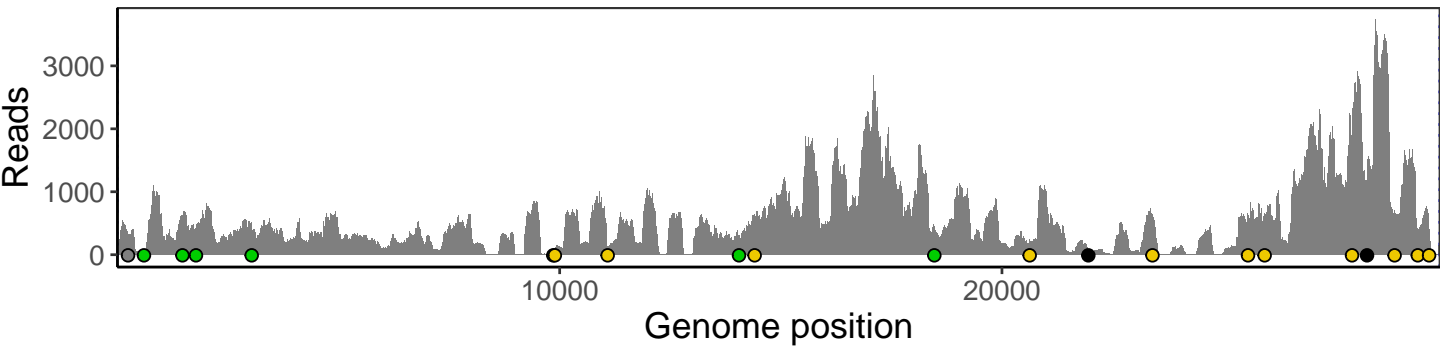- Expected
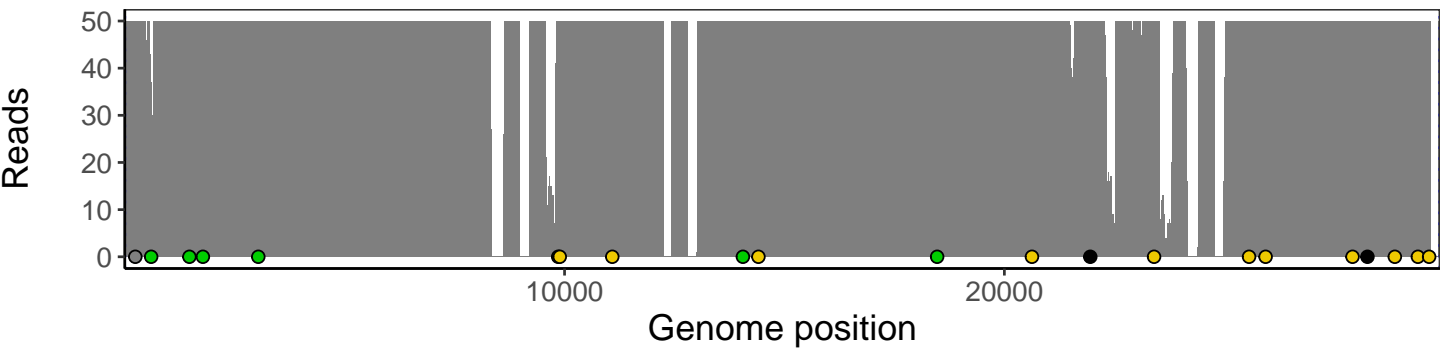- A
- T
- C
- G
- N
- Ins/Del
- No data

VSP0789−1

3

# Analyses of individual experiments and composite results

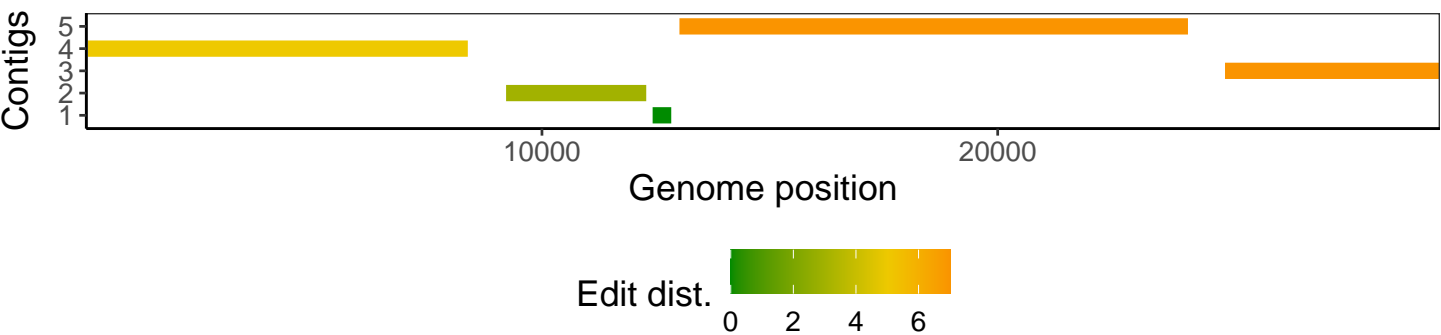**VSP0789-1 | 2021-01-29 | Saliva | 463s | genomes | single experiment**

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 htslib 1.10 |
| bcftools | 1.10.2-34-g1a12af0-dirty Using htslib 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 | 3.3.3 |
| 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 |