# COVID-19 subject UPHS-0843

*2021-06-23*

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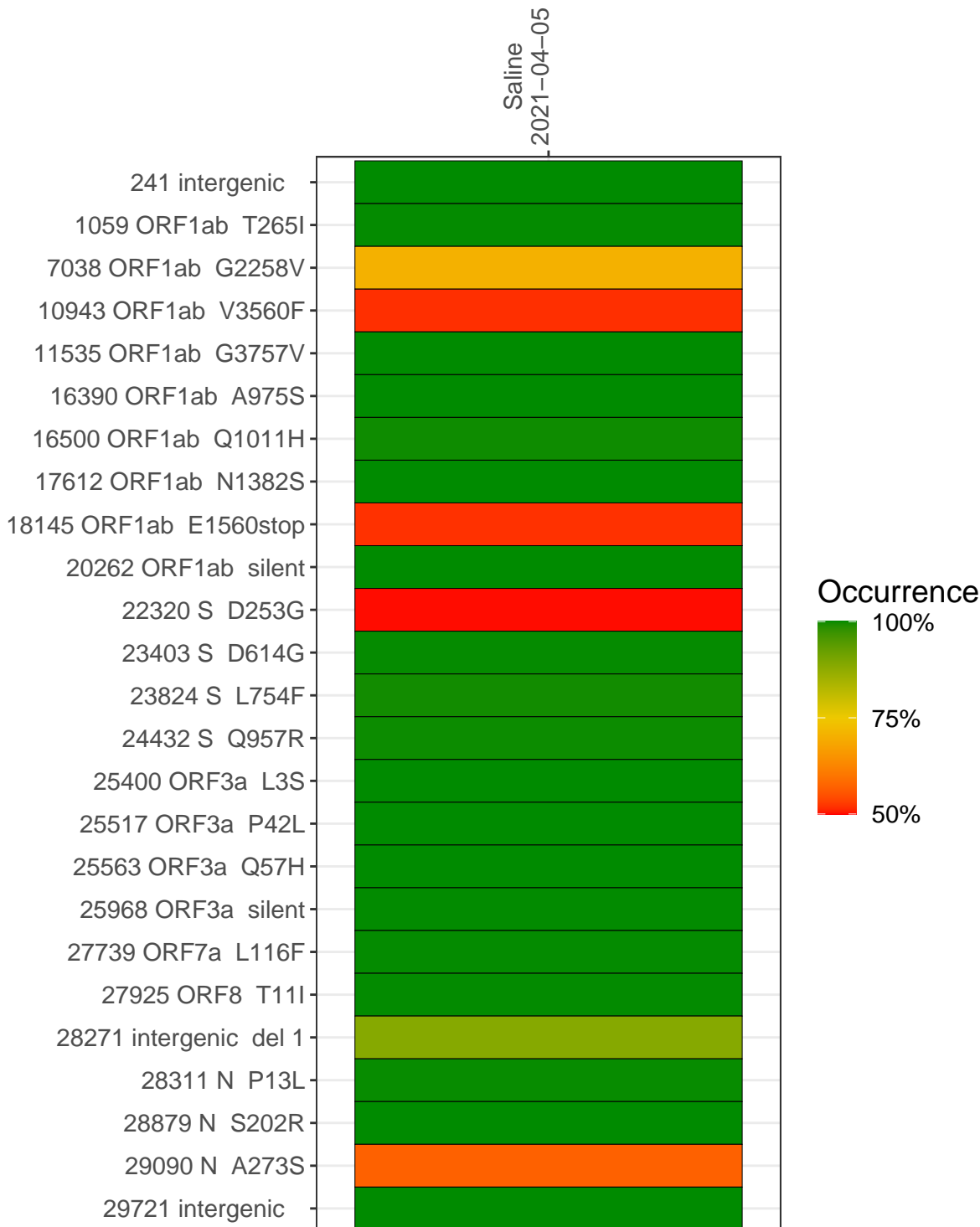 sequencing experiments performed for each specimen. Experiment specific analyses are shown at the end of this report. Lineages are called with the Pangolin software tool (Rambaut et al 2020) for genomes with > 90% sequence coverage.

Table 1. Sample summary.

| Experiment | Type | Genomes | Sample type | Sample date | Largest contig (KD) | Lineage | Reference read coverage | Reference read coverage (>= 5 reads) |
|---|---|---|---|---|---|---|---|---|
| VSP2057-2 | single experiment | NA | Saline | 2021-04-05 | 5.34 | NA | 76.0% | 75.1% |

# Variants shared across samples

The heat map below shows how variants (reference genome /home/common/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
2021−04−05

| Position | VSP2057−2 |
|---|---|
| 241 intergenic | 291 |
| 1059 ORF1ab T265I | 1007 |
| 7038 ORF1ab G2258V | 2499 |
| 10943 ORF1ab V3560F | 1327 |
| 11535 ORF1ab G3757V | 110 |
| 16390 ORF1ab A975S | 81 |
| 16500 ORF1ab Q1011H | 869 |
| 17612 ORF1ab N1382S | 567 |
| 18145 ORF1ab E1560stop | 1126 |
| 20262 ORF1ab silent | 715 |
| 22320 S D253G | 207 |
| 23403 S D614G | 1791 |
| 23824 S L754F | 414 |
| 24432 S Q957R | 661 |
| 25400 ORF3a L3S | 639 |
| 25517 ORF3a P42L | 500 |
| 25563 ORF3a Q57H | 716 |
| 25968 ORF3a silent | 2076 |
| 27739 ORF7a L116F | 3244 |
| 27925 ORF8 T11I | 3898 |
| 28271 intergenic del 1 | 686 |
| 28311 N P13L | 553 |
| 28879 N S202R | 185 |
| 29090 N A273S | 1010 |
| 29721 intergenic | 244 |

Base change
- Expected (gray)
- A (green)
- T (red)
- C (blue)
- G (yellow)
- N (purple)
- Ins/Del (black)
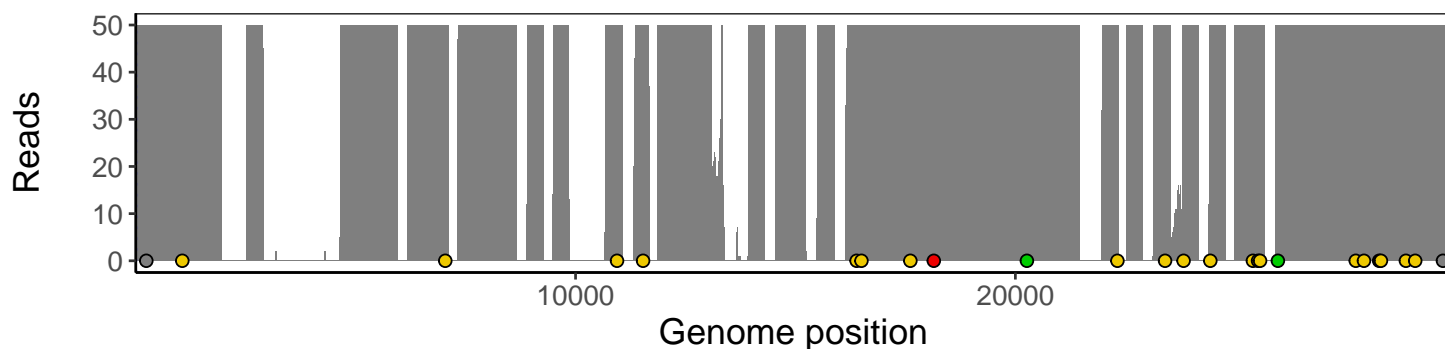- No data (light gray)

VSP2057−2

# Analyses of individual experiments and composite results

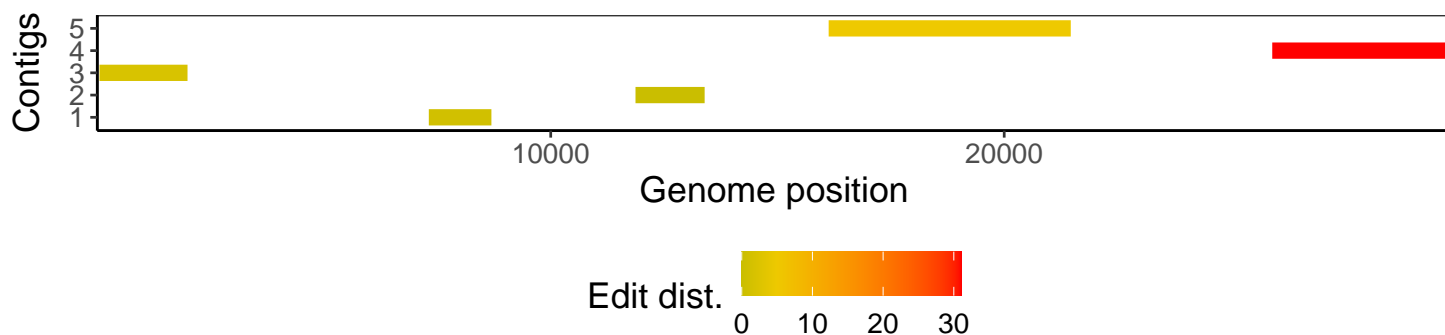**VSP2057-2 | 2021-04-05 | Saline | UPHS-0843 | genomes | single experiment**

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 | 3.1.3 |
| 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 |