

COVID-19 subject SARS_CoV_93

2021-06-29

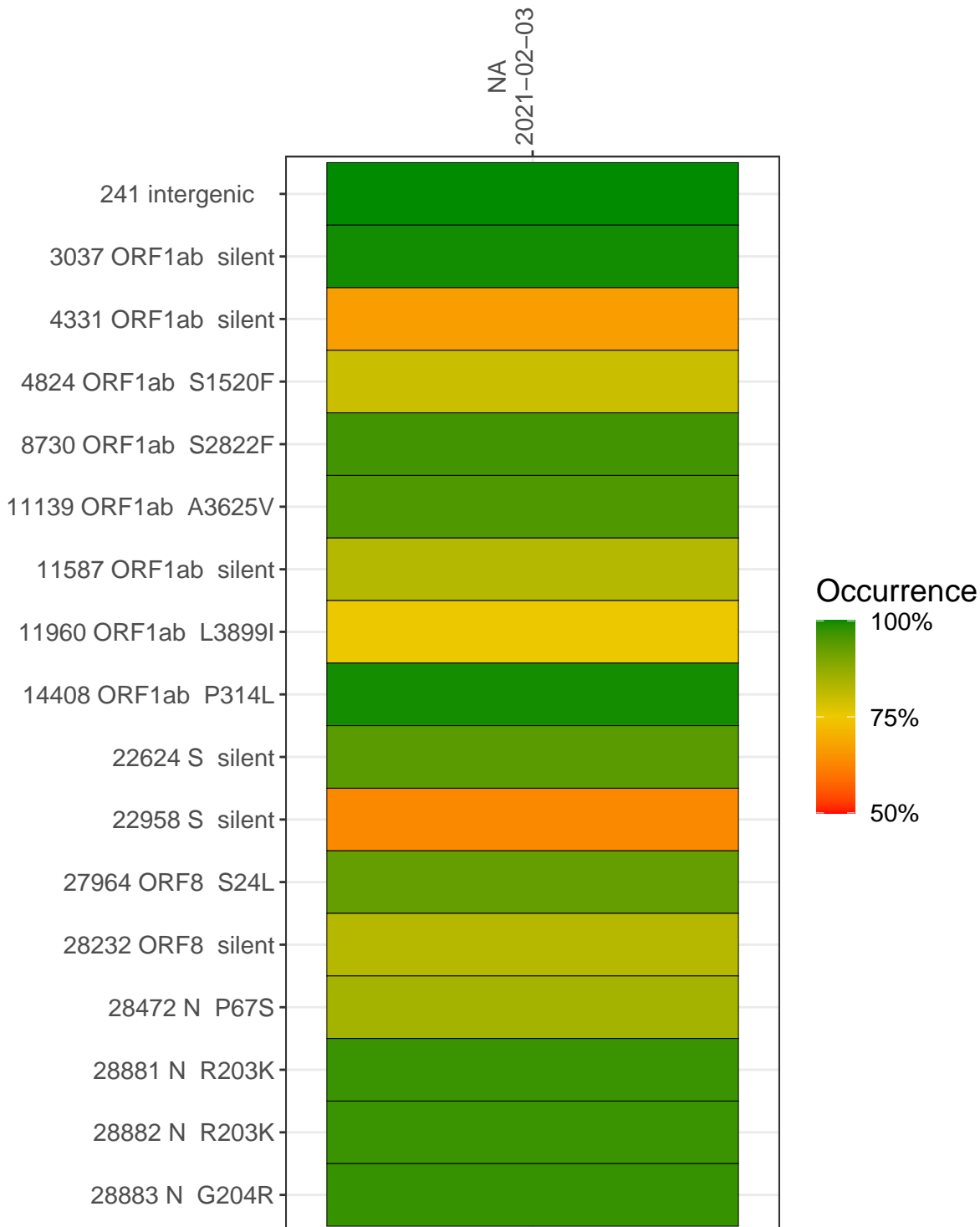
The table below provides a summary of subject samples for which sequencing data is available. The experiments column shows the number of sequencing experiments performed for each specimen. Experiment specific analyses are shown at the end of this report. Lineages are called with the Pangolin software 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)
VSP3015-1	single experiment	NA	NA	2021-02-03	11.72	B.1.1.161	99.7%	99.7%

Variants shared across samples

The heat map below shows how variants (reference genome `/home/common/SARS-CoV-2-Philadelphia/NC_045512`) 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.



NA
2021-02-03

241 intergenic	262
3037 ORF1ab silent	1677
4331 ORF1ab silent	555
4824 ORF1ab S1520F	1756
8730 ORF1ab S2822F	355
11139 ORF1ab A3625V	1441
11587 ORF1ab silent	3216
11960 ORF1ab L3899I	1670
14408 ORF1ab P314L	3836
22624 S silent	24510
22958 S silent	348
27964 ORF8 S24L	1193
28232 ORF8 silent	890
28472 N P67S	656
28881 N R203K	14841
28882 N R203K	14831
28883 N G204R	14833

Base change

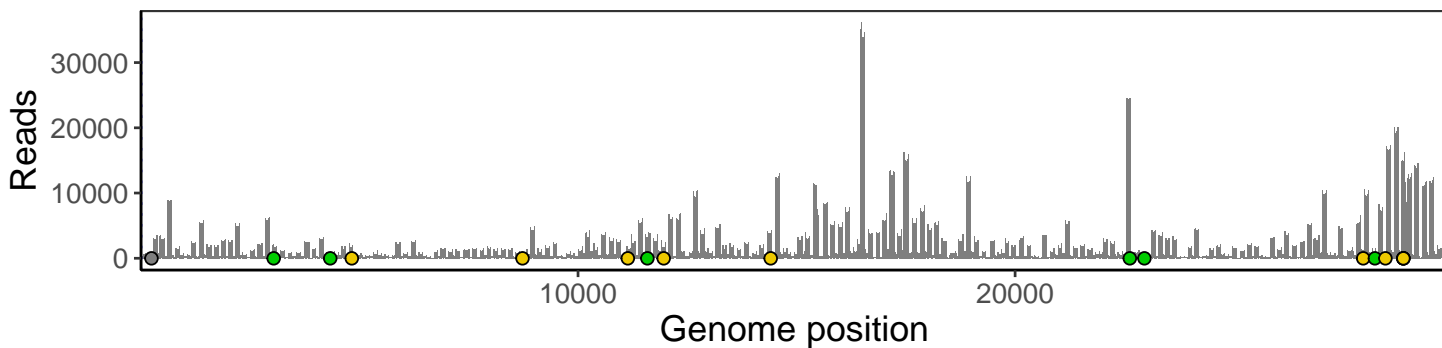
- Expected
- A
- T
- C
- G
- N
- Ins/Del
- No data

VSP3015-1

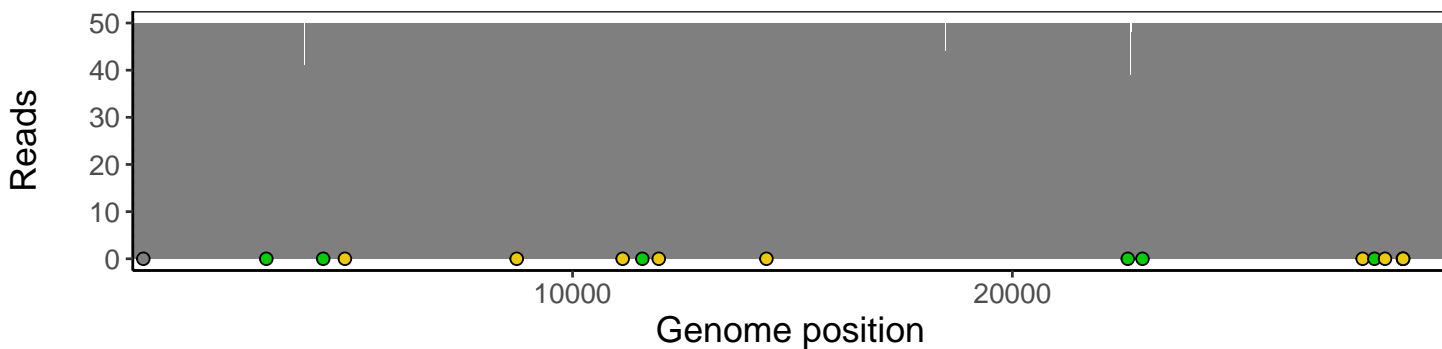
Analyses of individual experiments and composite results

VSP3015-1 | 2021-02-03 | NA | SARS_CoV_93 | genomes | single experiment

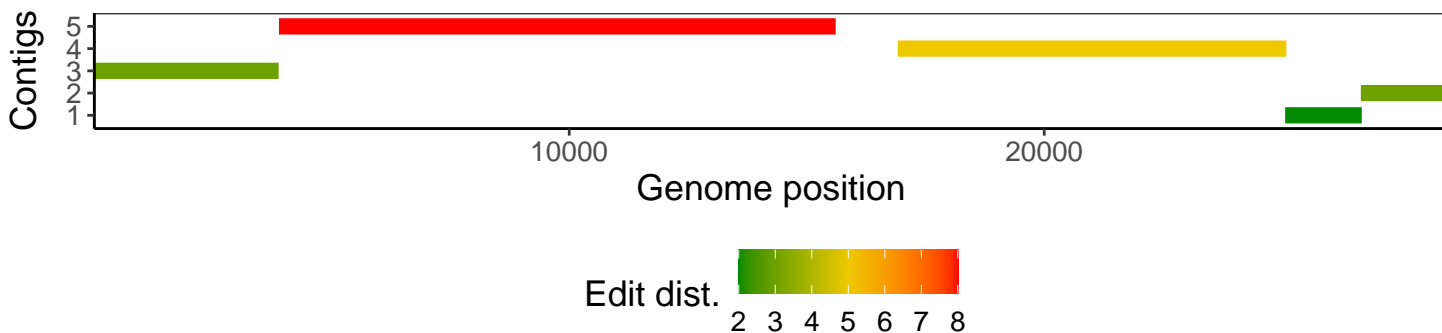
The plot below shows the number of reads covering each nucleotide position in the reference genome. Variants are shown as colored dots along the bottom of the plot and are color coded according by variant types: gray - transgenic, green - silent, gold - missense, red - nonsense, black - indel.



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	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