

A resource for SARS-Cov2  
mutation tracking and analysis:  
**The LANL/GISAID COVID pipeline**

Bette Korber

Will Fischer

Los Alamos National Laboratory  
Los Alamos, New Mexico, USA

5 November 2020

powered by



# Unprecedented *outbreak*: Unprecedented *data*

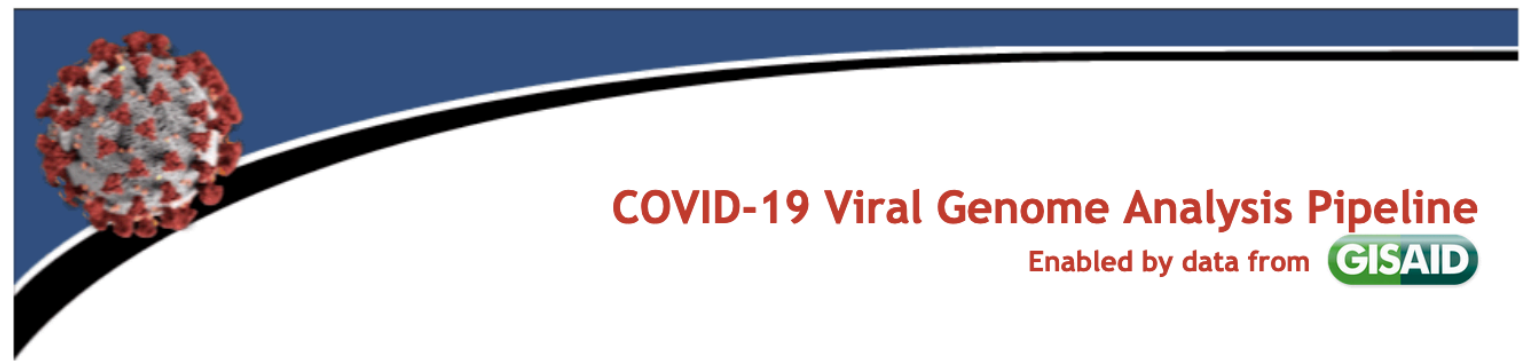
## High case numbers meet high sequencing capacity and bioinformatic infrastructure

- The GISAID database (initially tasked with archiving influenza sequences) collects and distributes SARS-Cov2 sequence data (January 2020).
- The Korber team at the LANL HIV database begins analyzing SARS-Cov2 data, processing the raw sequences for mutational analysis (February 2020).
- LANL goal: forward-looking threat assessment of mutational escape from vaccines/therapeutics



# LANL COVID-19/ SARS-CoV2 website: <https://cov.lanl.gov>

- Leverages infrastructure of LANL HIV-DB
- Provides analytical tools and links to computed analyses
- Emphasis on geographical distribution and temporal evolution of non-synonymous mutations in coding regions, especially the Spike protein
- Sequence data and analytical results updated daily.



[Analytical resources](#)

[Tools](#)

[Home](#)

This website provides analyses and tools for exploring accruing mutations in hCoV-19 (SARS-CoV-2) geographically and over time, with an emphasis on the Spike protein, using data from GISAID.

The SARS-CoV-2 sequence data used for these analyses was updated from GISAID on Nov 1, 2020

The analyses provided are based on a trimmed full length SARS-CoV-2 alignment containing 98,946 sequences:

[sequence names and ID numbers used for full-length analyses](#),

or on a Spike alignment containing 135,322 sequences:

[sequence names and ID numbers used for spike-only analyses](#).

The details of the analyses are described in:

**Tracking changes in SARS-CoV-2 Spike: evidence that D614G increases infectivity of the COVID-19 virus.**

Korber B, Fischer WM, Gnanakaran S, Yoon H, Theiler J, Abfalterer W, Hengartner N, Giorgi EE, Bhattacharya T, Foley B, Hastie KM, Parker MD, Partridge DG, Evans CM, Freeman TM, de Silva TI\*, McDanal C, Perez LG, Tang H, Moon-Walker A, Whelan SP, LaBranche CC, Saphire EO, and Montefiori DC.

\*on behalf of the Sheffield COVID-19 Genomics Group

Cell, June 2020

[DOI:10.1016/j.cell.2020.06.043](https://doi.org/10.1016/j.cell.2020.06.043)

## News

Nov 3, 2020

We have released an updated [summary of SARS-CoV-2 variation](#).

Sep 3, 2020

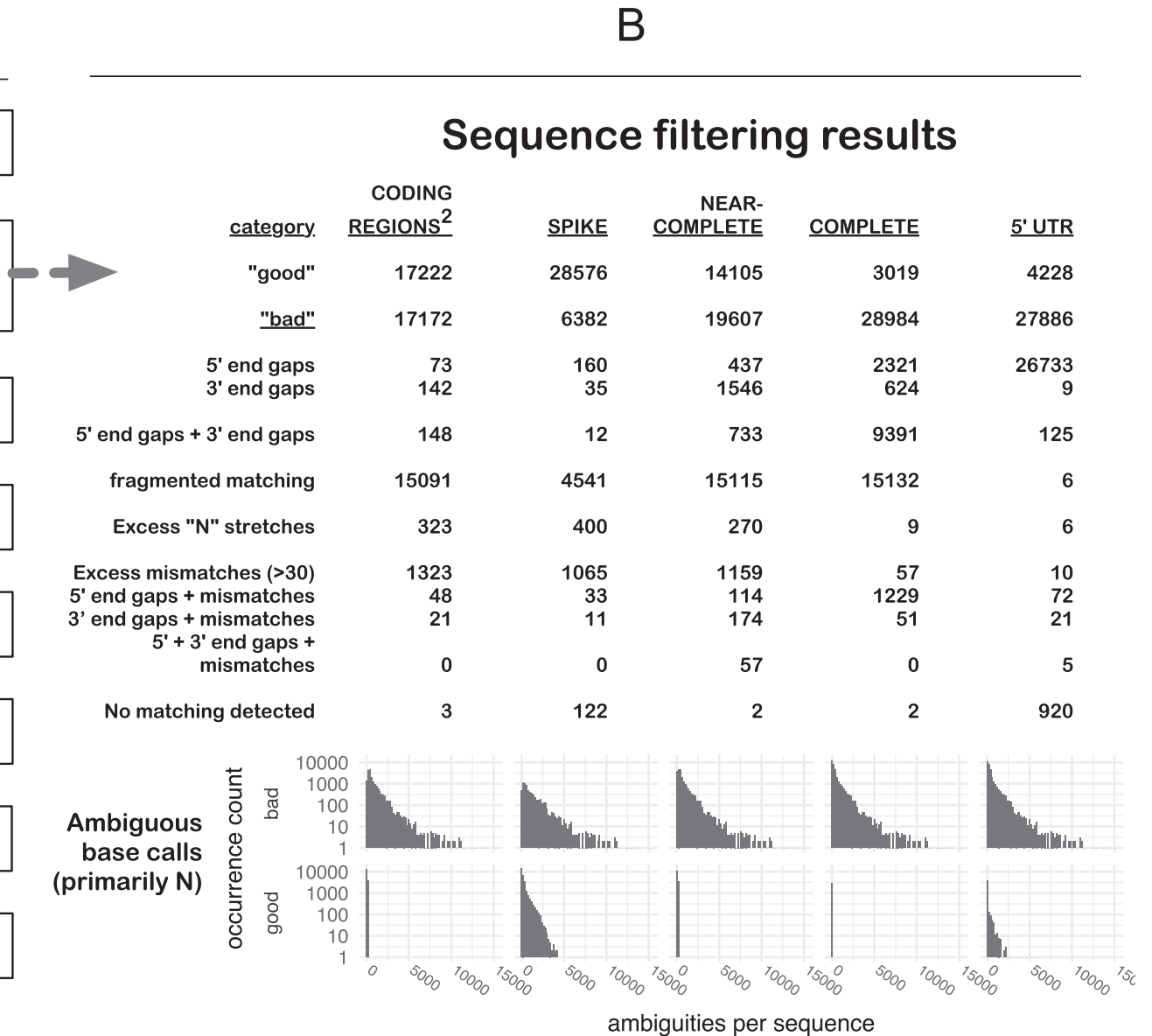
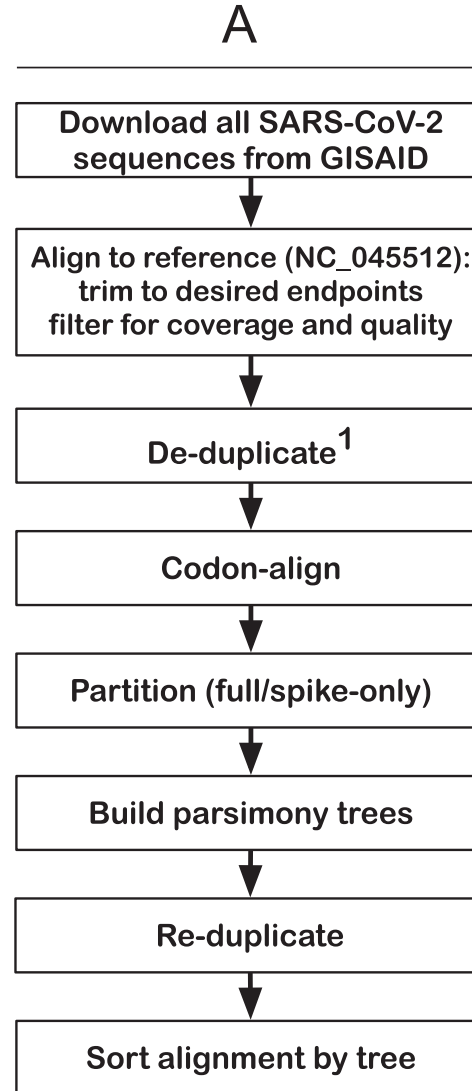
We have released a [summary of SARS-CoV-2 variation](#), with a focus on Spike mutations that might impact antibodies, and considerations for vaccine reagents.

[See more](#)

# Sequence Processing Pipeline

Quality control and processing

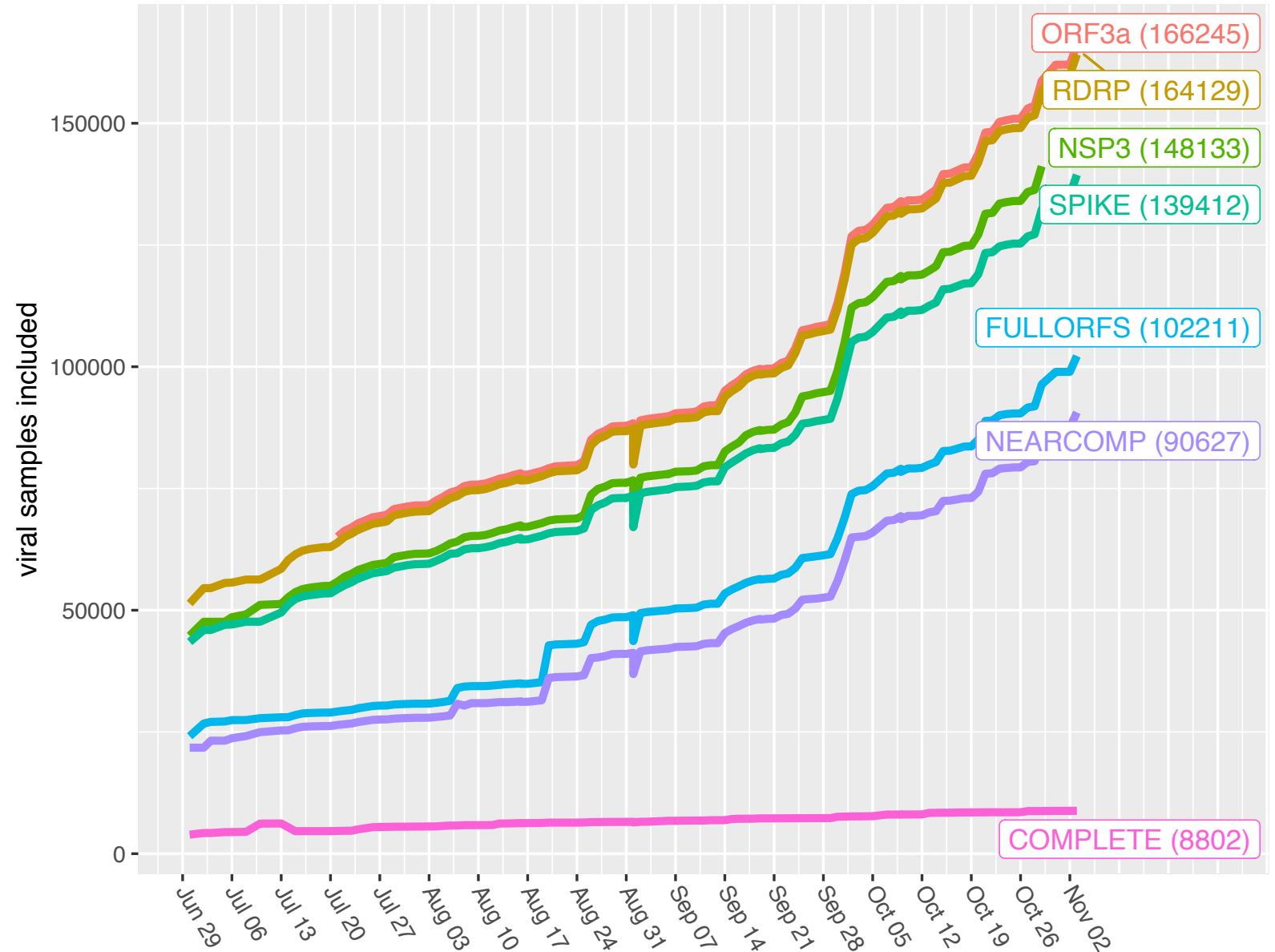
- Data received at LANL in custom JSON feed
- Processing:
  - mapping to reference strain
  - quality filtering
  - alignment
  - phylogenetic inference



# Continued growth in sequence submission

- Sequence counts tripled through July-October 2020
- LANL partitions sequences into regions of interest, assembling aligned datasets containing only high-quality bases throughout a given region
  - FULLORFS: start codon of first reading frame (ORF1ab) to the stop codon of the last reading frame (ORF10): 29,409 nt
  - NEARCOMP: includes the most-commonly-used primer sites: 29,782 nt
  - SPIKE and other coding regions
  - Sequences are excluded from a given partition if they have many uncalled bases, long indels, or extensive mismatches

Sequence counts in LANL partitions of GISAID sequences: growth over time (to 2020-11-03)



The O, S, L, and V clade are rarely sampled after June 1, through the summer of 2020

G is Spike D614G:

G has 3 sub-lineages

GR,  
GH,  
and GV (Spike A222V)

GR is the most frequently sampled, but is very common in the UK which is highly sampled.

GH is also frequently sampled, and is common in the US.

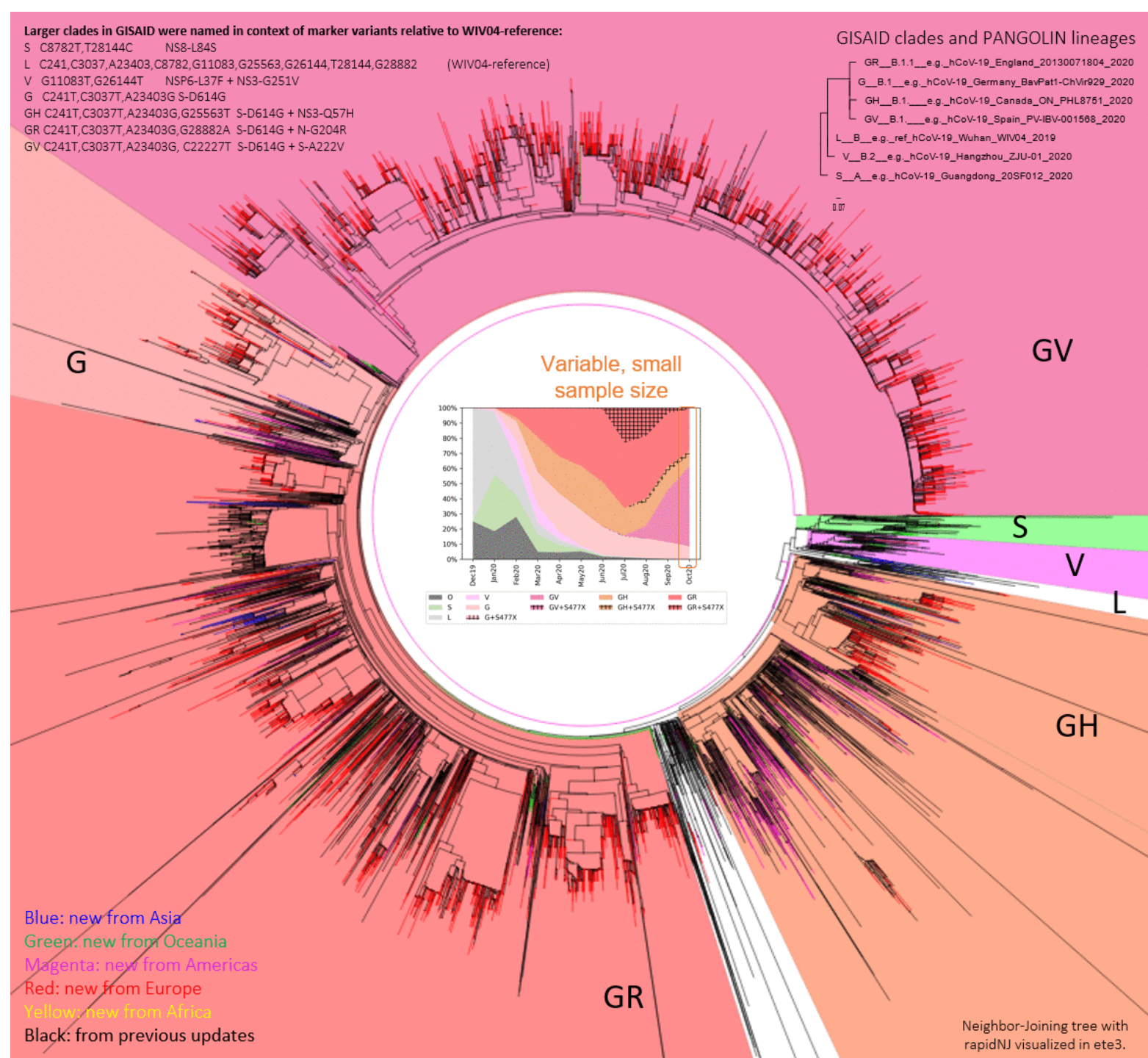
L is complex, may include recombinants?

Larger clades in GISAID were named in context of marker variants relative to WIV04-reference:

S C8782T,T28144C NS8-L84S  
 L C241,C3037,A23403,C8782,G11083,G25563,G26144,T28144,G28882 (WIV04-reference)  
 V G11083T,G26144T NSP6-L37F + NS3-G251V  
 G C241T,C3037T,A23403G S-D614G  
 GH C241T,C3037T,A23403G,G25563T S-D614G + NS3-Q57H  
 GR C241T,C3037T,A23403G,G28882A S-D614G + N-G204R  
 GV C241T,C3037T,A23403G, C22227T S-D614G + S-A222V

GISAID clades and PANGOLIN lineages

GR\_B.1.1\_e.g.\_hCoV-19\_England\_20130071804\_2020  
 G\_B.1.1\_e.g.\_hCoV-19\_Germany\_BavPat1-ChVir029\_2020  
 GH\_B.1.1\_e.g.\_hCoV-19\_Canada\_ON\_PHL8751\_2020  
 GV\_B.1.1\_e.g.\_hCoV-19\_Spain\_FV-IBV-001568\_2020  
 L\_B\_e.g.\_ref\_hCoV-19\_Wuhan\_WIV04\_2019  
 V\_B.2\_e.g.\_hCoV-19\_Hangzhou\_ZJU-01\_2020  
 S\_A\_e.g.\_hCoV-19\_Guangdong\_20SF012\_2020



Full genome tree derived from all outbreak sequences 2020-10-30

Notable changes:

155,961 full genomes (+7,632) (excluding low coverage, out of 167,272 entries)

Updated clades:  
 S clade 6,495 (+39)  
 L clade 4,236 (+1)  
 V clade 5,343 (+21)  
 G clade [#S477X] 28,633 [105] (+769 [+3])  
 GR clade [#S477X] 60,971 [8,868] (+2,130 [+33])  
 GH clade [#S477X] 34,922 [1,098] (+932 [+204])  
 GV clade [#S477X] 11,431 [5] (+3,719 [+0])  
 Other clades 3,930 (+21)

Blue: new from Asia  
 Green: new from Oceania  
 Magenta: new from Americas  
 Red: new from Europe  
 Yellow: new from Africa  
 Black: from previous updates

We gratefully acknowledge the Authors from Originating and Submitting laboratories of sequence data on which the analysis is based.

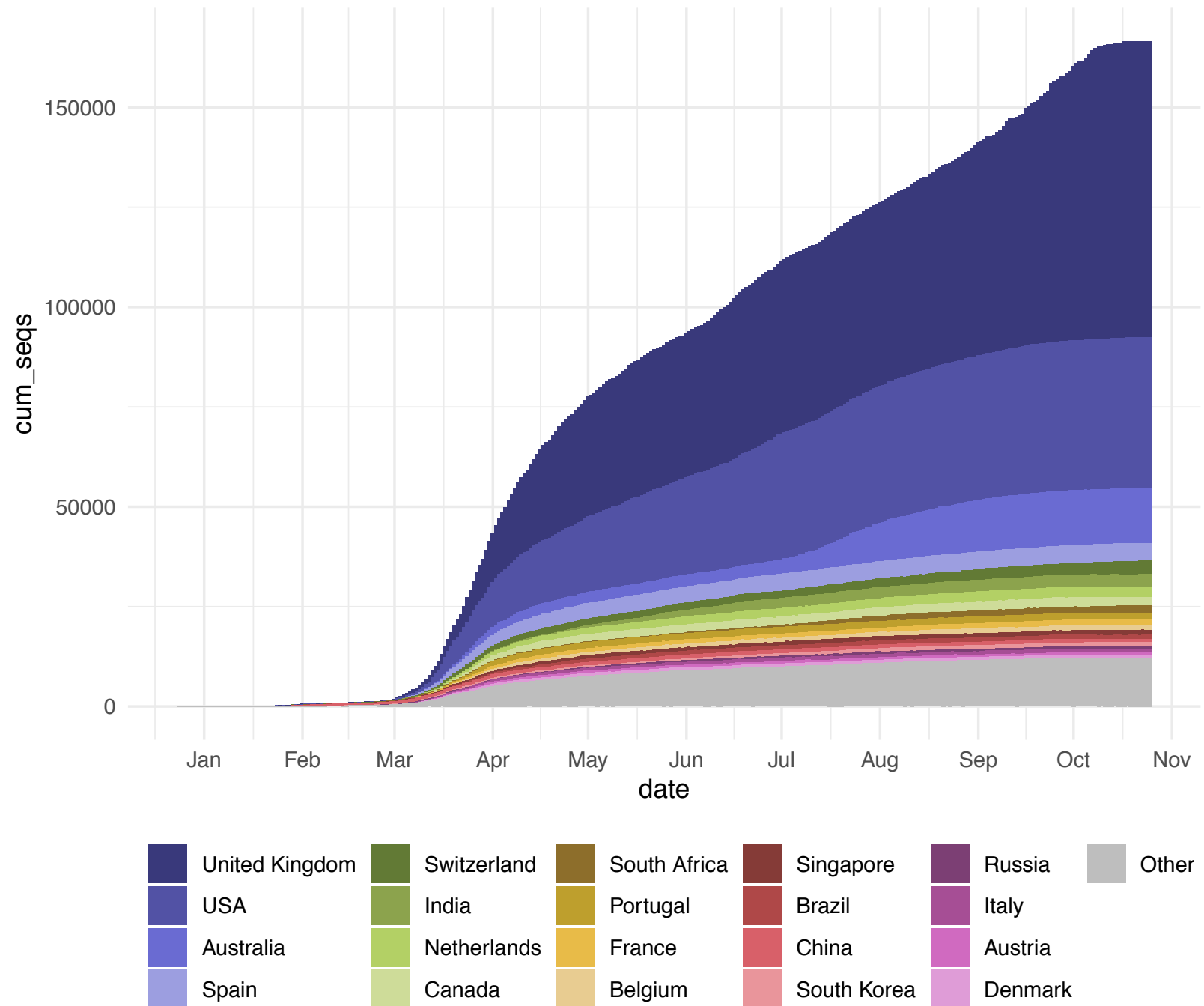


Neighbor-Joining tree with rapidNJ visualized in ete3.

by BII/GIS, A\*STAR Singapore

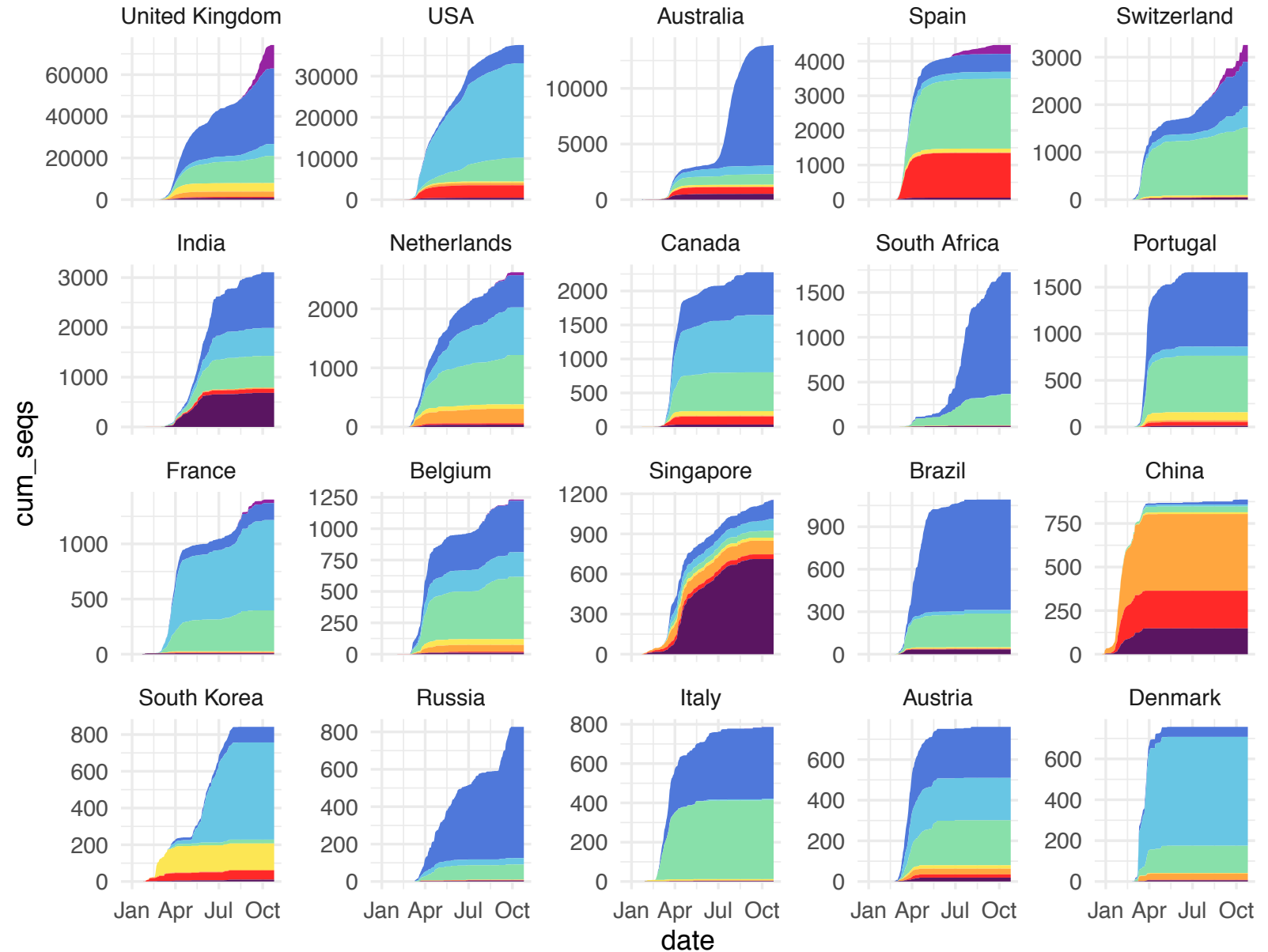
# Sequencing effort is highly biased

- Geographically:
  - The top 20 countries account for 92% of all sequences submitted.
  - 78% of GISAID SARS-Cov2 sequence data come from the top 4 countries (U.K., USA, Australia, Spain); >44% from the U.K. alone.
- Temporally
  - many countries employed concentrated bursts of sequencing effort, with extensive sampling gaps



# Clade distribution varies between countries

- Cumulative sequence counts (vertical change shows added sequences)
- In most countries, early clades (O,S,L,V) are replaced by G clade and its descendants
- Examples:
  - arrival and expansion of GR clade in Australia in July 2020
  - appearance of GV in mid-August in the U.K., France, and Switzerland; earlier in Spain
  - continued detection of the original "O" form in Singapore
  - "frozen" early frequency distributions where few later sequences reported (e.g., China, Portugal, Denmark)

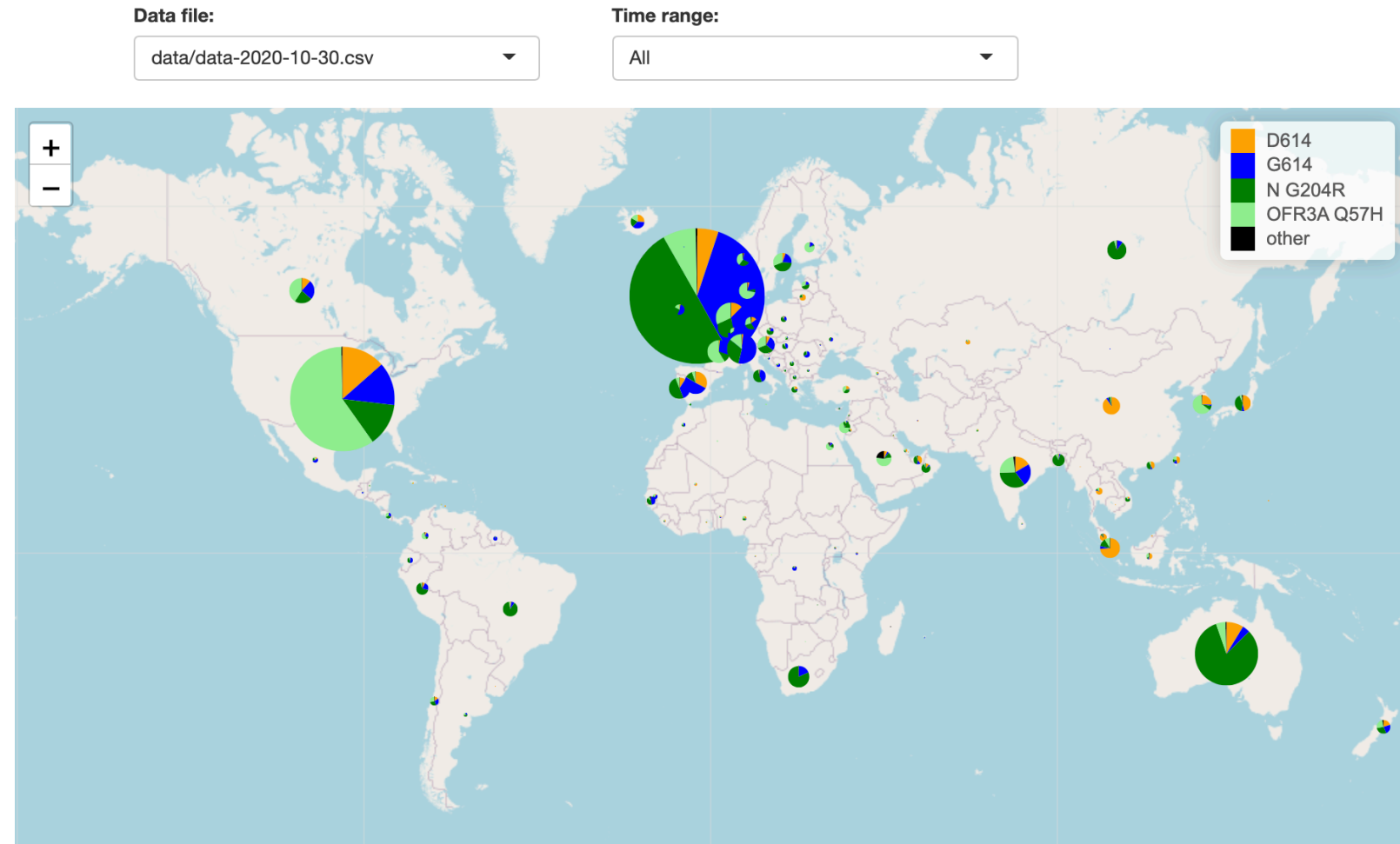




# Mutational clade distributions

- Distributions sized by total sequence count (not case count!)
- Two-week intervals from March 1, 2020
- After March, "G" and descendent clades predominate worldwide

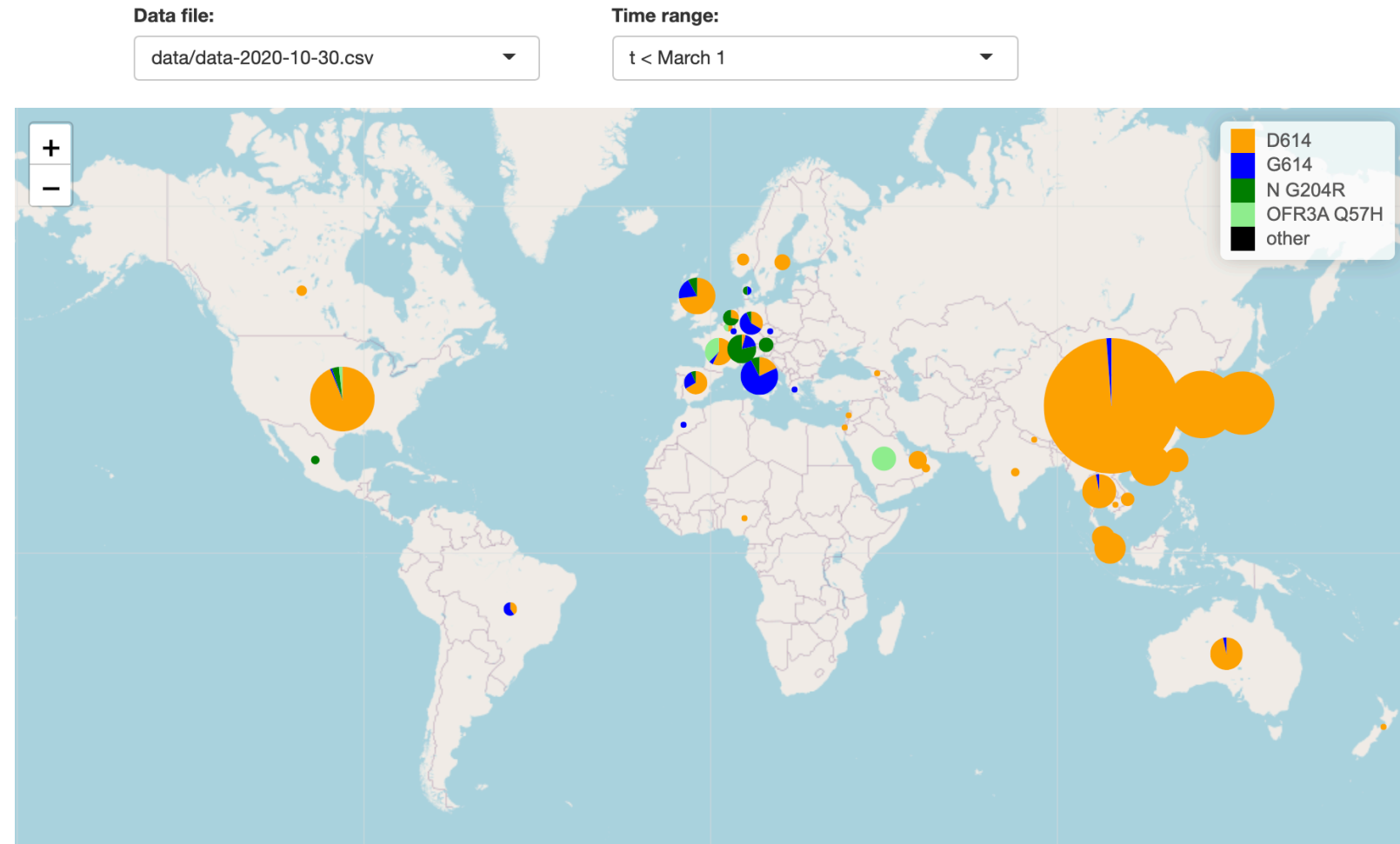
## Distribution of D614, G614, N G204R and OFR3A Q57H



# Mutational clade distributions

- Distributions sized by total sequence count (not case count!)
- Two-week intervals from March 1, 2020
- After March, "G" and descendent clades predominate worldwide

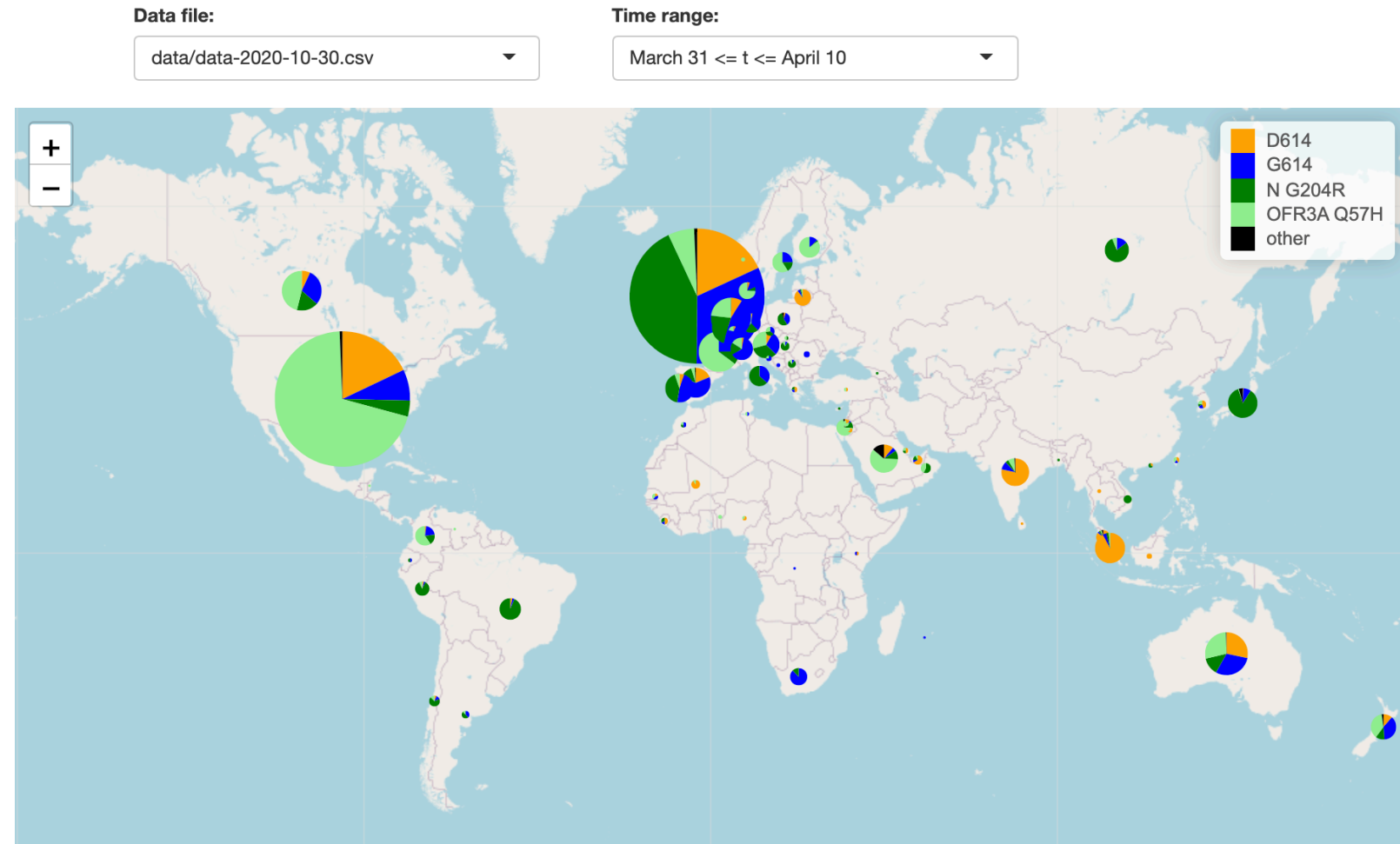
## Distribution of D614, G614, N G204R and OFR3A Q57H



# Mutational clade distributions

- Distributions sized by total sequence count (not case count!)
- Two-week intervals from March 1, 2020
- After March, "G" and descendent clades predominate worldwide

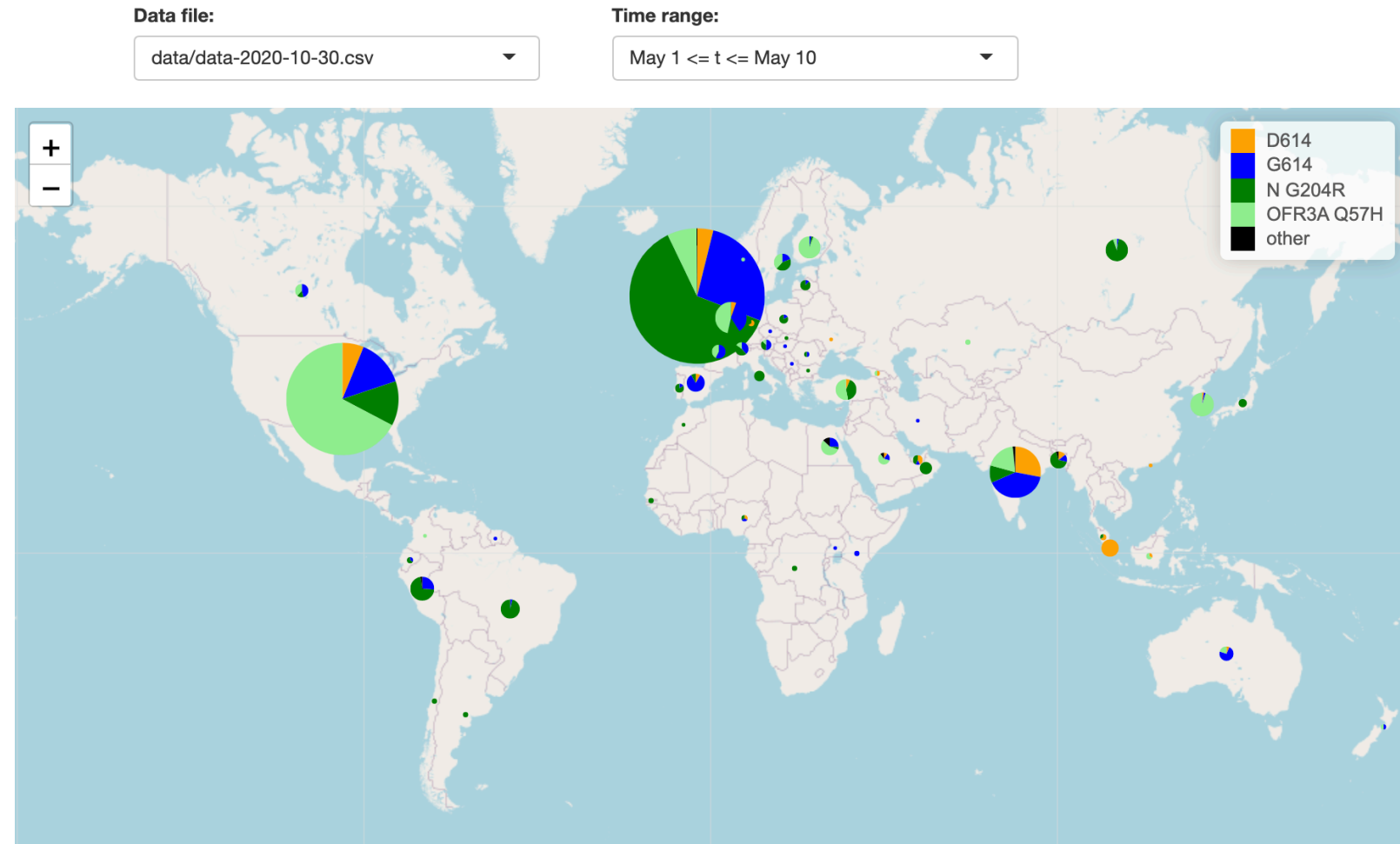
## Distribution of D614, G614, N G204R and OFR3A Q57H



# Mutational clade distributions

- Distributions sized by total sequence count (not case count!)
- Two-week intervals from March 1, 2020
- After March, "G" and descendent clades predominate worldwide

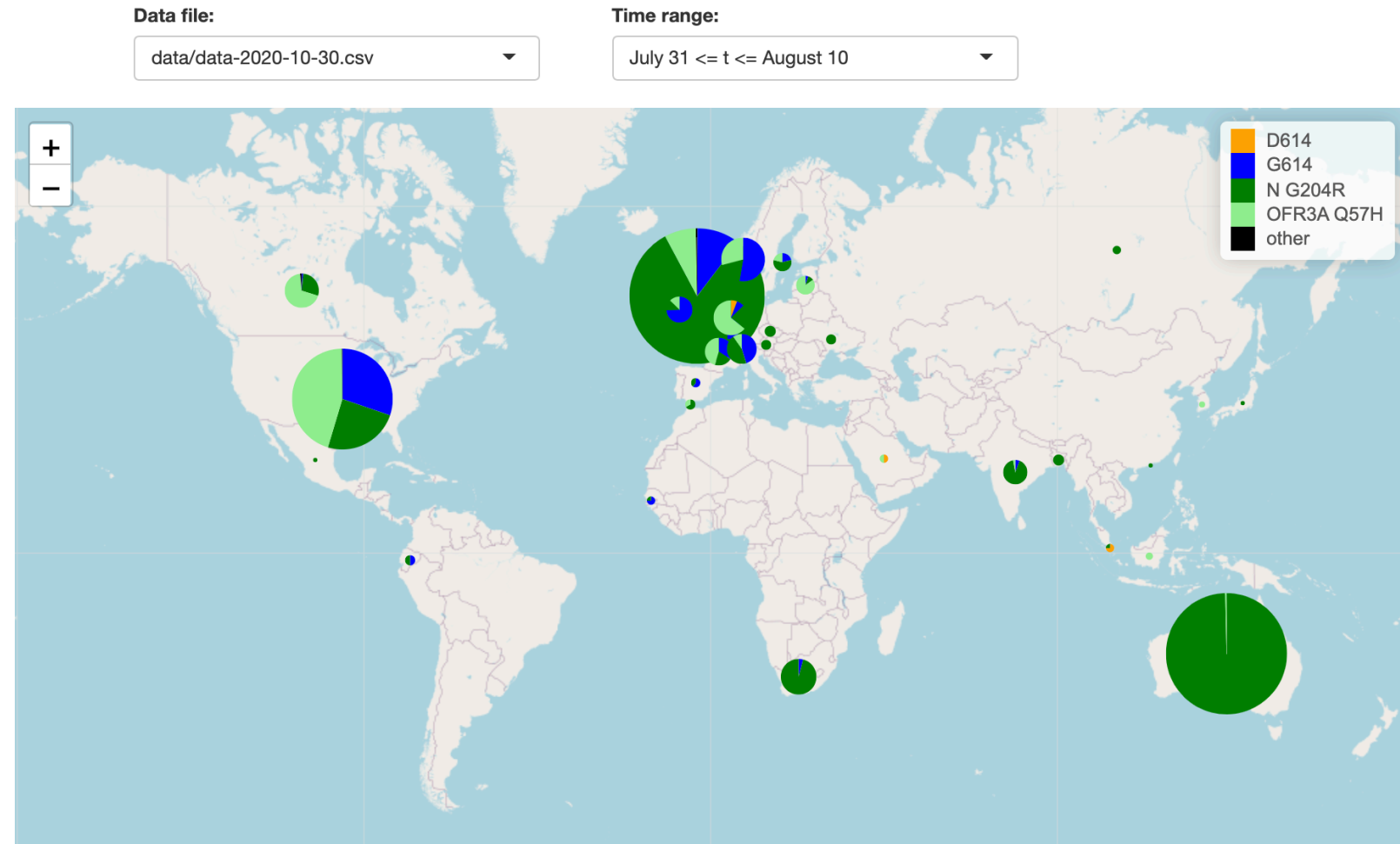
## Distribution of D614, G614, N G204R and OFR3A Q57H



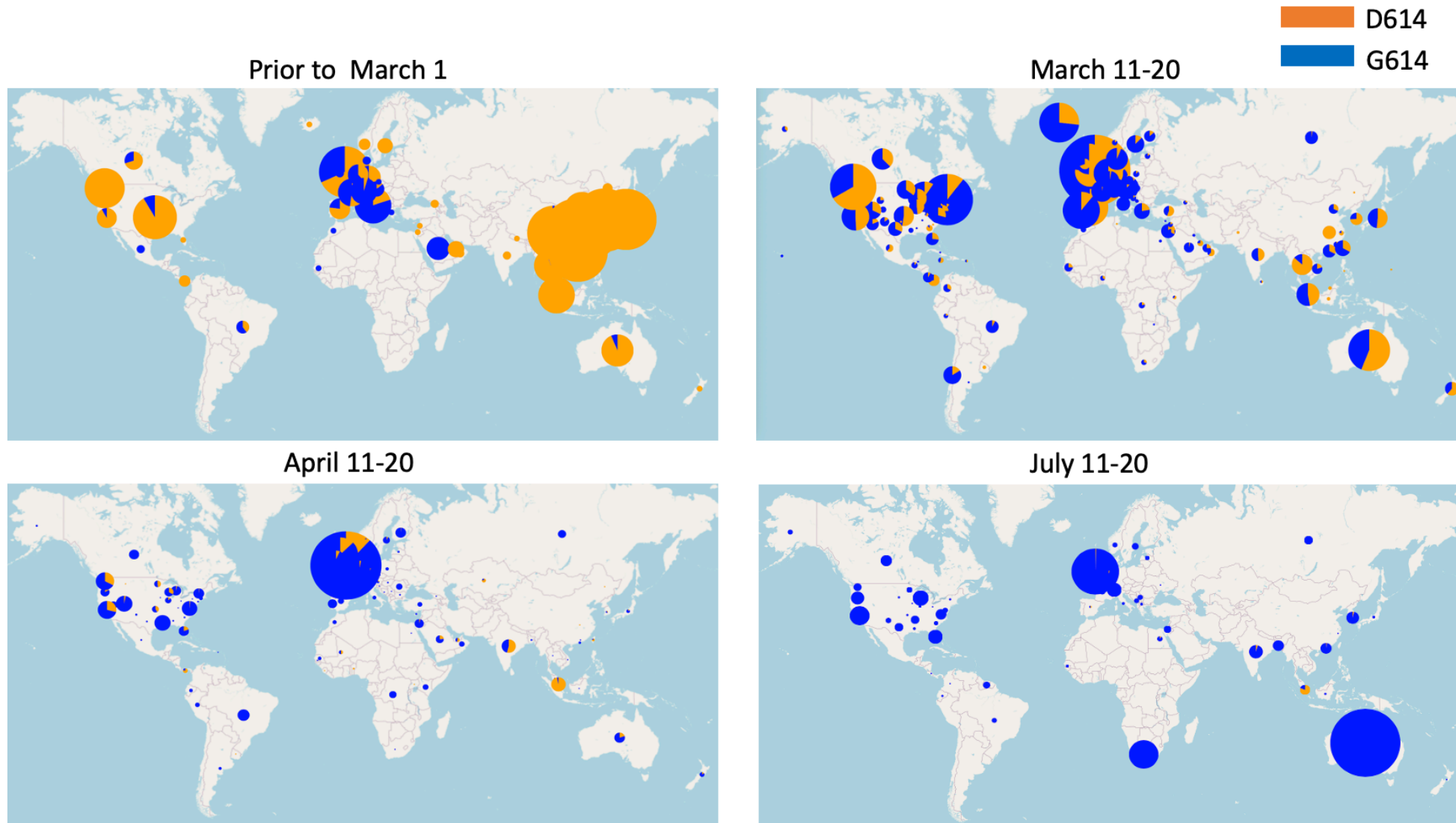
# Mutational clade distributions

- Distributions sized by total sequence count (not case count!)
- Two-week intervals from March 1, 2020
- After March, "G" and descendent clades predominate worldwide

## Distribution of D614, G614, N G204R and OFR3A Q57H



# A rapid change in frequency of the G614 form was apparent early in 2020



Geographic/temporal mutation tracking in the COVID-19 pandemic:  
<https://cov.lanl.gov>

# The D614G mutation

SARS-CoV-2 Spike is the protein that mediates virus entry into cells and it's the prime target of COVID vaccines.

A single mutation in Spike, D614 to G614, has become the dominant form of the virus in the pandemic.

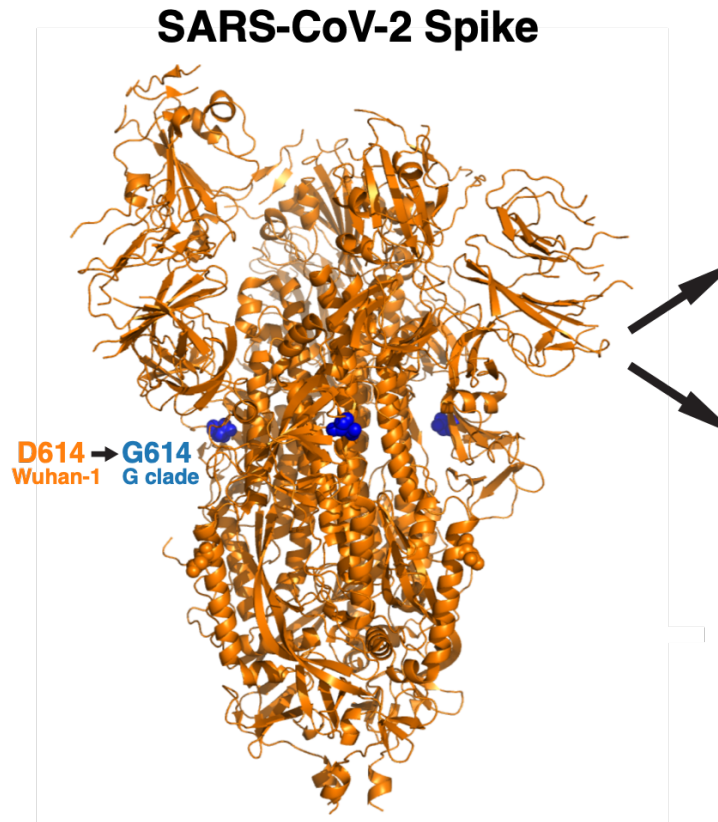
The G614 form is more infectious.

The G614 form is associated with higher levels of viral RNA in the upper respiratory tract of infected people.

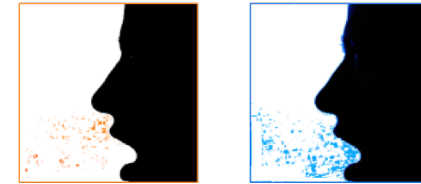
The G614 was *not* associated with increased hospitalization.

*Tracking Changes in SARS-CoV-2 Spike: Evidence that D614G Increases Infectivity of the COVID-19 Virus.*

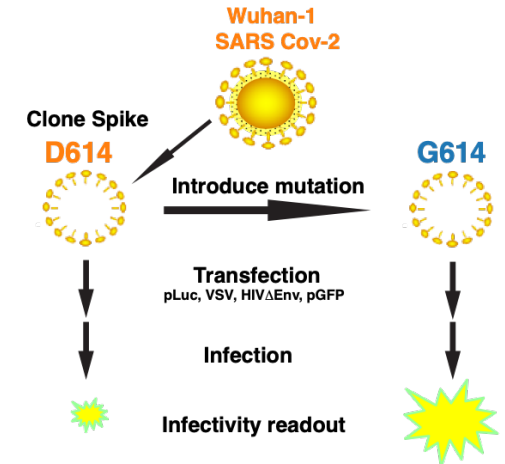
Korber B, Cell. 2020 Aug 20;182(4):812-827.e19.



G614 is associated with higher viral loads in patients indicated by lower RT PCR cycle thresholds for detection

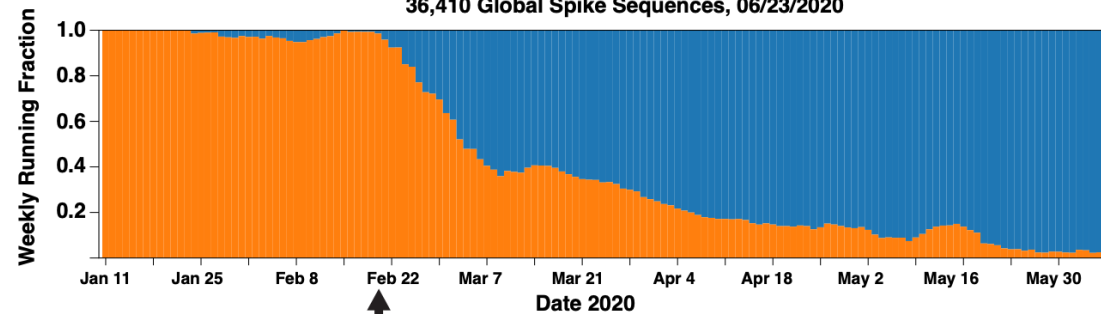


G614 is associated with higher infectious titers of spike pseudotyped virus



Global Transition from D614 to G614 variants

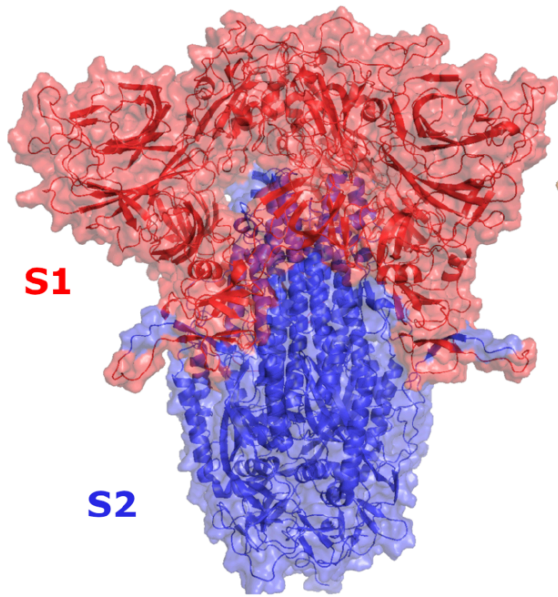
36,410 Global Spike Sequences, 06/23/2020



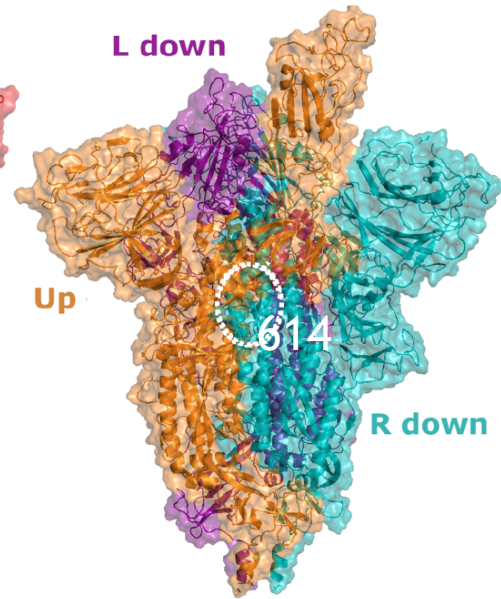
G614 emerges in Europe

**G614** is likely to be more infectious because it favors an open “one-up” conformation that makes its ACE2 receptor binding site more accessible.

A) S1/S2 in all down



B) Protomers in 1up

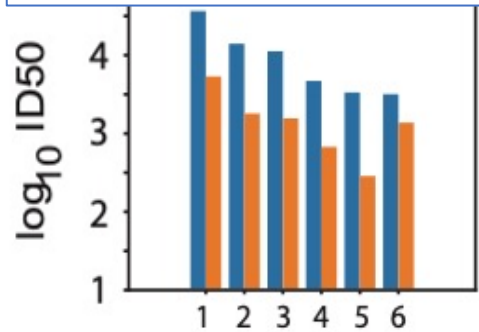


*The SARS-CoV-2 Spike Variant D614G Favors an Open Conformational State*  
 Mansbach et al. submitted and under review,  
 also in: bioRxiv. 2020 doi: 10.1101/2020.07.26.219741.  
 LANL/Duke collaboration

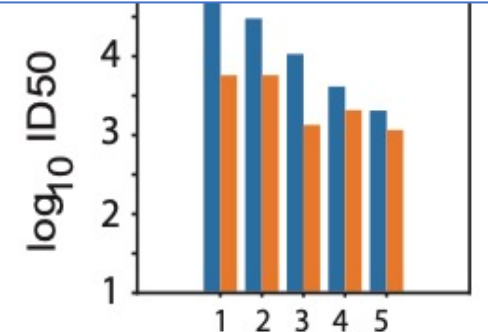
We were concerned that the more infectious **G614** mutated virus might be more resistant to antibodies.

Instead, it is even *more* sensitive to vaccine-induced antibodies, COVID-19 convalescent sera, and Spike antibodies. This is also likely a result of favoring the 1-up conformation.

**G614** Spike is more sensitive to sera from Spike vaccinated monkeys than **D614**



**G614** Spike is more sensitive to sera from Spike vaccinated people than **D614**



*D614G Spike Mutation Increases SARS CoV-2 Susceptibility to Neutralization.*

Weissman et al. in press, Cell Host and Microbe, Oct. 2020

Also in: medRxiv: doi.org/10.1101/2020.07.22.20159905

Duke/U. Penn./LANL collaboration



# Sites of particular interest in Spike, Nov. 2, 2020

- The original viruses that carried **D614** are very rarely now sampled globally. Among the **D614G** G clade viruses, GH and GR clades have emerged. The GR clade is now globally the most common, and is tending to increase in frequency relative to the G and GH clade.
  - The defining amino acid substitutions for the GR clade are outside of Spike
- The Spike **S477N** mutation virus started to dominant the sample in Australia in the summer, particularly in sequences from Victoria, and is now 6.3% of the global GISAID sample. It is increasing in the UK and several countries in Europe as well.
  - It is always found in a D614G context, but arises independently within the G, GR and GH clades
- The Spike **A222V** mutation has overtaken S477N as the second most common mutation in GISAID it is now at 7.9% of the global sample.
  - It is always found in a D614G context
  - It is most common found in the UK, but in increasing now in several countries.
- The Spike **N439K** mutation is becoming more common in the UK. It is currently present in 0.78% of the global sample. It is of particular interest because it is in the most common mutation in the RBM
  - The UK is by far the most heavily sampled country in GISAID, so this biases the global sample. It is not very common outside the UK at the point.
  - It is always found in a D614G context.
  - Embedded in several known neutralizing antibody epitope
    - Starr TN, et al. ... Bloom JD. Cell. 2020 Aug 11;S0092-8674(20)31003-5. doi: 10.1016/j.cell.2020.08.012. PMID: 32841599
  - It can confer escape and is a selected resistance mutation of Nab C135
    - Weisblum Y et al. ...Bieniasz PD. bioRxiv. 2020 Jul 22:2020.07.21.214759. doi: 10.1101/2020.07.21.214759. Preprint. PMID: 32743579

Example: tracking  
variation in site Spike 477

# Summary of Spike Mutations: Spreadsheet

- A unannotated version of both of the spreadsheet tables that summarize variation are provided with a daily updates from GISAID at [cov.lanl.gov](https://cov.lanl.gov)
- The tab labeled “Spike Variation” contains a row for each position in Spike that includes:
  - The number of each variant, the entropy of each site, the local entropy of each 10 amino acid stretch
  - Spikes sites with > 0.3% variation (or within 4 Å of ACES, 0.1%) in GISAID are highlighted in red
  - Sites and local linear regions that have relatively high entropy are highlighted in yellow
  - Sites are annotated with Spike regions and mAb features are informed by f the following sources:
    - Starr TN, et al. ... Bloom JD. Cell. 2020 Aug 11:S0092-8674(20)31003-5. doi: 10.1016/j.cell.2020.08.012. PMID: 32841599  
[Deep Mutational Scanning of SARS-CoV-2 Receptor Binding Domain Reveals Constraints on Folding and ACE2 Binding](#)  
Annotation is based on: [https://jbloombio.github.io/SARS-CoV-2-RBD\\_DMS/](https://jbloombio.github.io/SARS-CoV-2-RBD_DMS/)
    - Weisblum Y et al. ...Bieniasz PD. bioRxiv. 2020 Jul 22:2020.07.21.214759. doi: 10.1101/2020.07.21.214759. Preprint. PMID: 32743579  
[Escape from neutralizing antibodies by SARS-CoV-2 spike protein variants](#)
    - Barnes CO, et al. ... Bjorkman PJ. Cell. 2020 Aug 20;182(4):828-842.e16. doi: 10.1016/j.cell.2020.06.025 .PMID: 32645326  
[Structures of Human Antibodies Bound to SARS-CoV-2 Spike Reveal Common Epitopes and Recurrent Features of Antibodies](#)
- The tab labeled “Common Spike mutations”, summarizes amino acids that on Nov 2. had >0.3% variants
  - We excluded D614G because otherwise it overwhelms the information as it is globally dominant now
  - For other varying sites we provide counts, codons, amino acid variants, and geographic regions
  - A version of this table is provided with a daily update is available at [cov.lanl.gov](https://cov.lanl.gov)

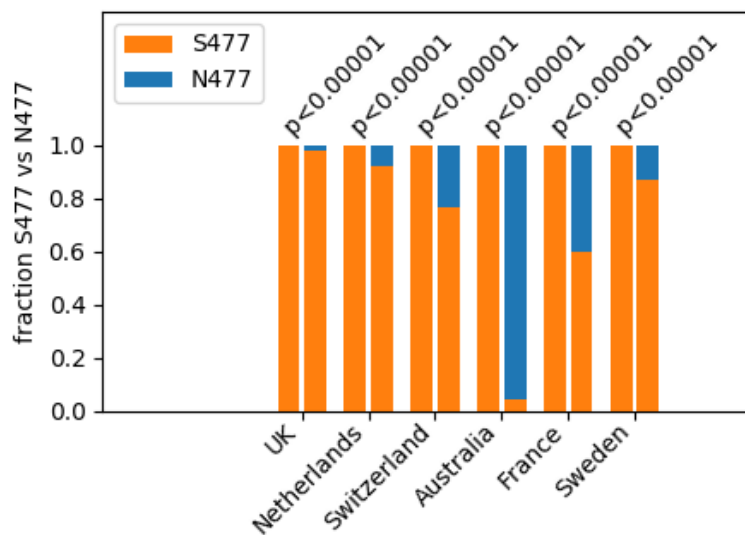


# Steps:

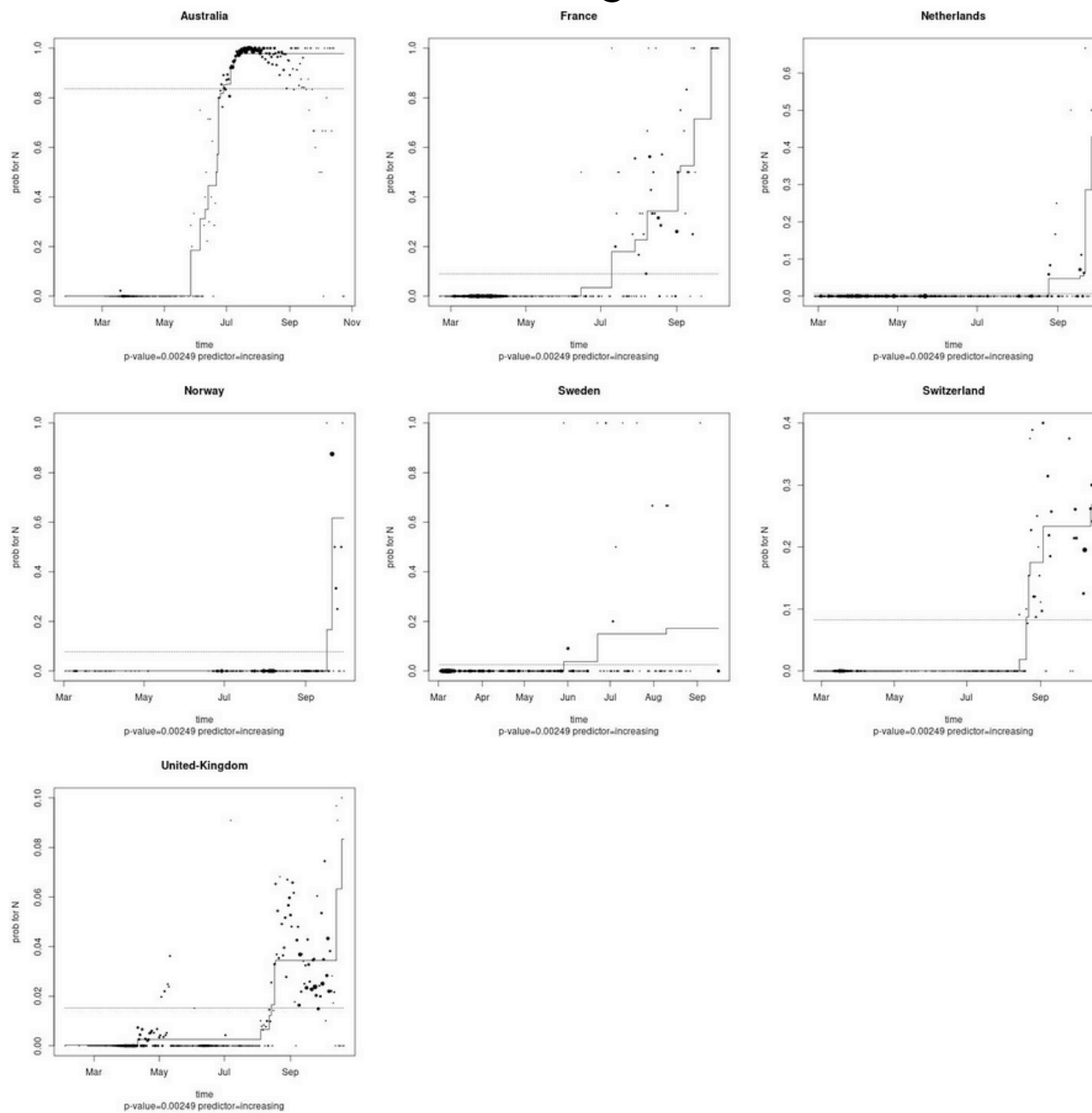
- 1) Tables of variation in Spike
  - 1) Summaries of sites that have  $>0.3\%$  variation, updated daily: under 477, you can see mutations in this site are by far the most common in Australia
  - 2) Variation in all positions in Spike are be tracked
- 2) Tracking places where variation is accruing:
  - 1) New mutations are accruing in a G614 background, so exclude D614 the original form
  - 2) Using Isotonic Regression
  - 3) Using Relative Frequency Change By Geographical Region
- 3) Tracking mutations: Follow variation in a position over time in a given place: e.g. 477 Australia
- 4) Analyze Align: e.g. How often are mutations in 614, 222, 439, 477 tracking together?

The mutation S477N has recently increased to 6.3% of GISAID, and is beginning to significantly increase in a number of countries

## Relative Frequency Change By Geographical Region



## Isotonic Regression

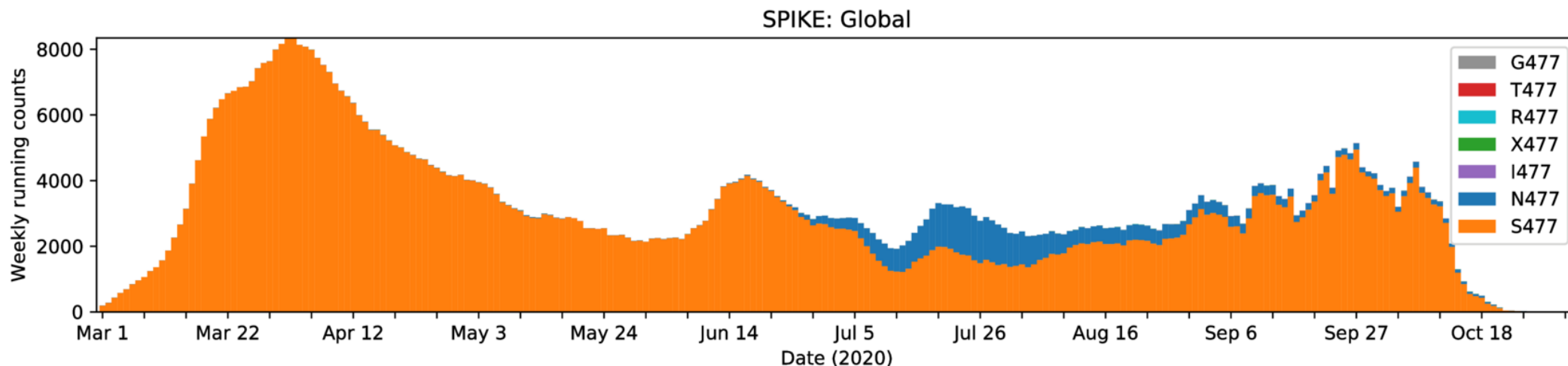


# Tracking mutations:

## The frequency **S477N** in the global sample is dominated by Australia

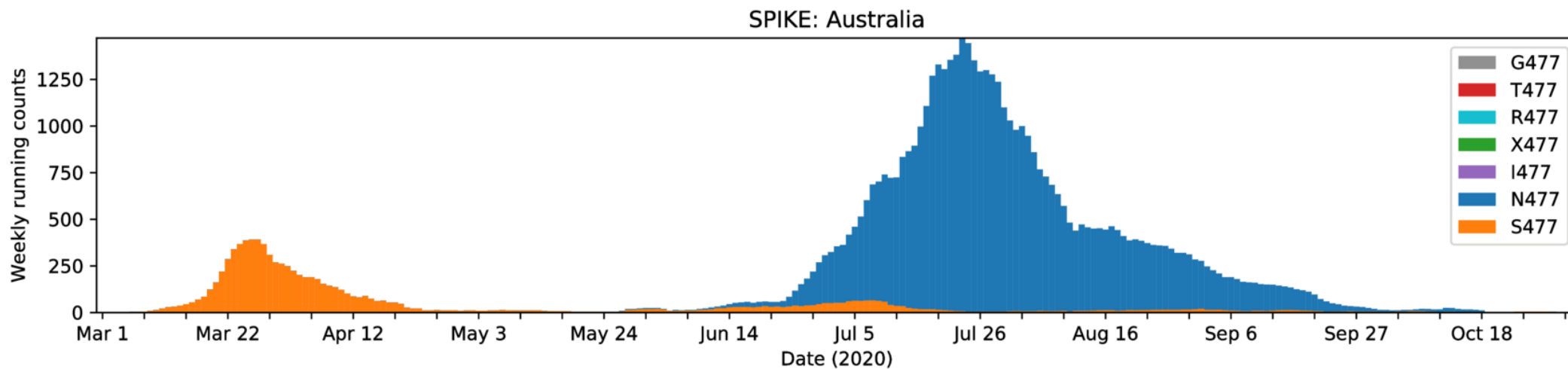
**SPIKE 477 Global: 116454 good entries: 7% S477N total, 6% from Australia, 1% from other places**

S477: 108094, N477: 8240, I477: 54, X477: 48, R477: 16, T477: 2, G477: 0



**SPIKE 477 Australia: 8391 good entries. S477N is 83.6% of the Australian sample.**

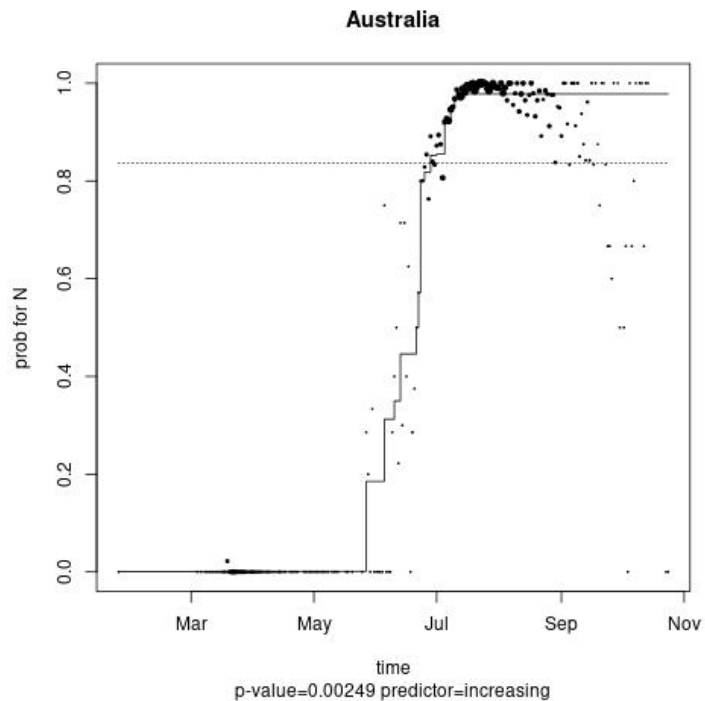
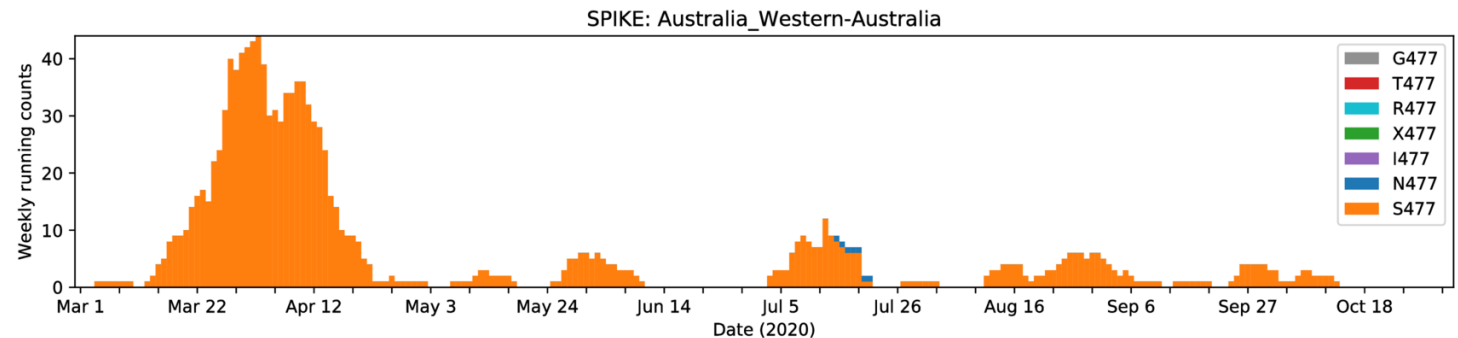
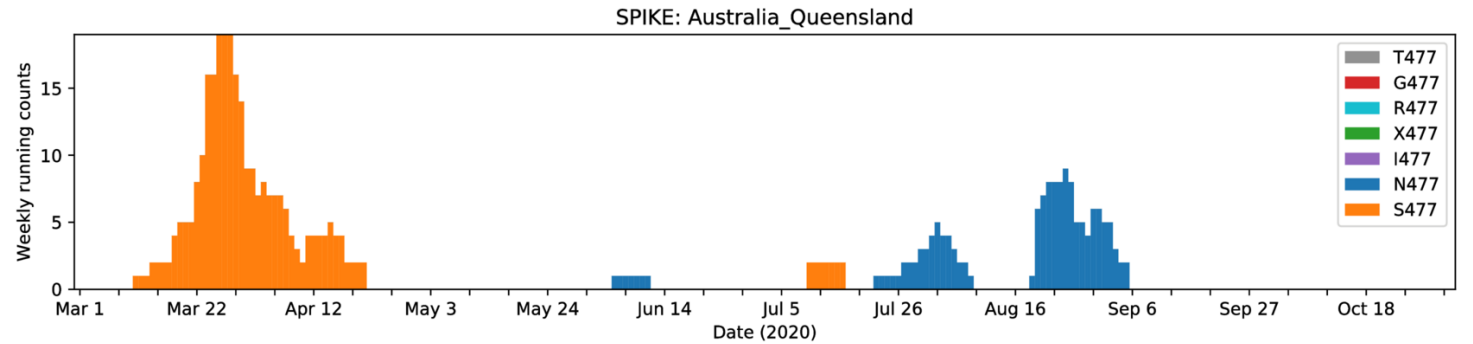
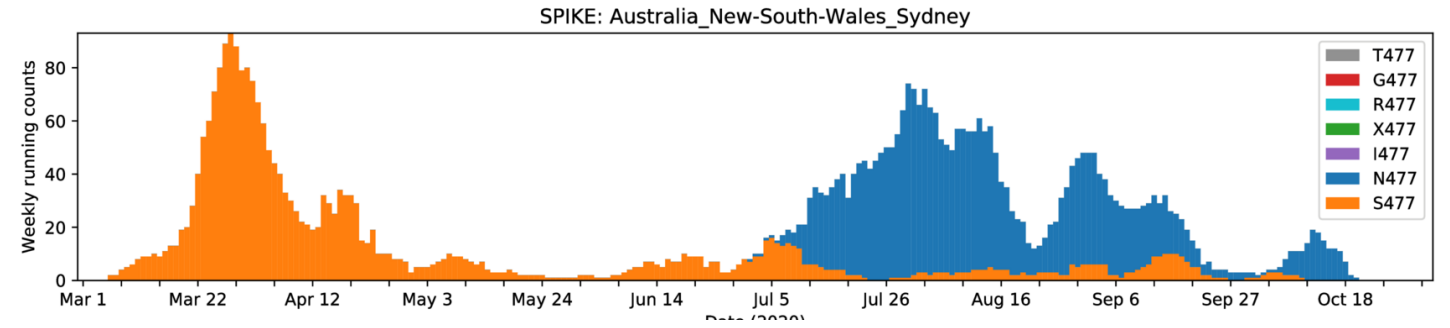
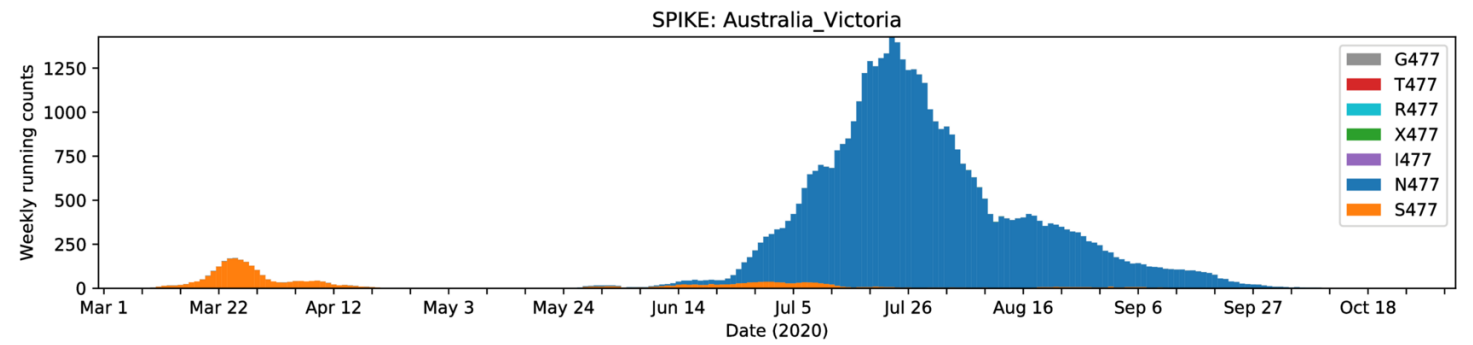
S477: 1373, N477: 7017, I477: 0, X477: 1, R477: 0, G477: 0, T477: 0



# Australia 477

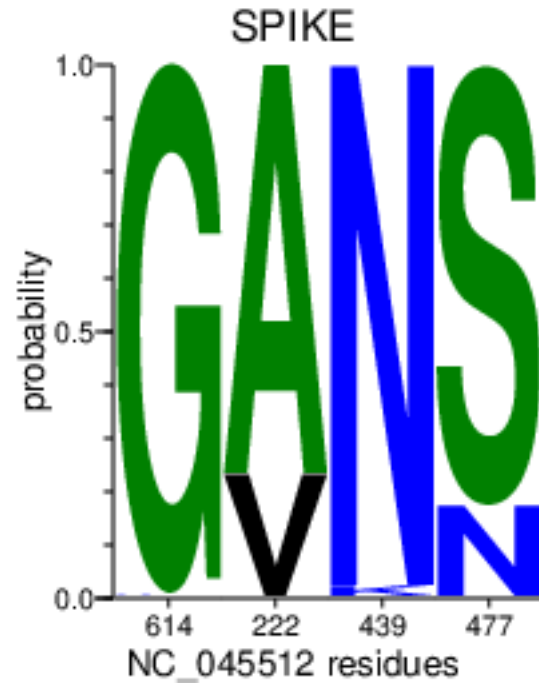
First became common in GISAID due to over 6000 sequences from Victoria sampled and provided to GISAID during the late summer. Victoria had completely switched to S477N.

Most Australian samples reflect this switch (Sydney and Queensland), but not all (Western Australia)





# Analyze Align 7/1/2020 - 11/3/2020



Form	N muts	Count	Percent
<b>GANS</b>	<b>0</b>	<b>25964</b>	<b>56.85</b>
<b>-V--</b>	<b>1</b>	<b>10510</b>	<b>23.01</b>
<b>---N</b>	<b>1</b>	<b>7749</b>	<b>16.97</b>
<b>--K-</b>	<b>1</b>	<b>977</b>	<b>2.14</b>
<b>D---</b>	<b>1</b>	<b>202</b>	<b>0.44</b>
<b>---I</b>	<b>1</b>	<b>41</b>	<b>0.09</b>
<b>-V-N</b>	<b>2</b>	<b>16</b>	<b>0.04</b>

121 were found in Singapore

Combinations are very rare

