

# COVID-19 Spike Protein Phylogenetic trees

Ashvin Pidaparti, Bruno Bohn, Carrie Kistler, Kimberly Mandery

# Background and Summary of Previous Findings

D614G mutation

Spike specific

Increased prevalence of mutations with higher replication rates and fitness

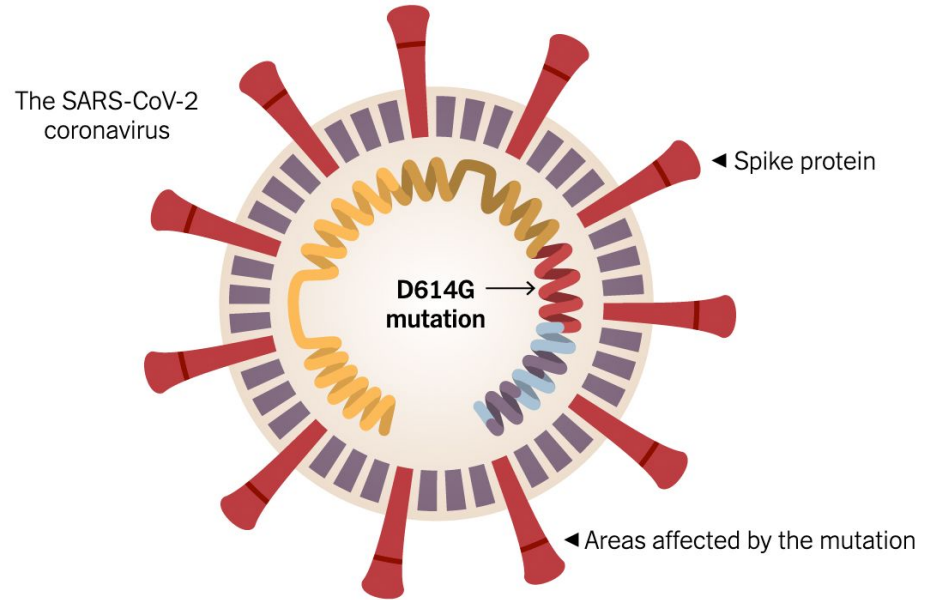


Image courtesy of the New York Times

# Sequences

Queried on **NCBI**

Spike protein (~1273 aa)

150 sequences:

- 15 states (10 each):

- 5 before July

- 5 July and after

Outgroup:

- 7 sequences from China,

Jan/Feb 2020

Refine Results

Reset

Virus

Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), taxid:2697049

Accession

Sequence Length

Min: 1250  
Max: 1300

Sequence Type

RefSeq Genome Completeness

Nucleotide Completeness

complete

Proteins

New!

surface glycoprotein

Provirus

Geographic Region

New!

USA

Host

Author

Isolation Source

Collection Date

From Jul 1, 2020  
To Dec 12, 2020

Selected Results: 0

Nucleotide (0)

Protein (3,912)

RefSeq Genome (0)

Select Columns

Expand Table

	Accession	Release Date	Species	Length	Nuc Completeness	Protein	Geo Location
	<a href="#">QPD95898</a>	2020-12-11	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Texas, B
	<a href="#">QPD95910</a>	2020-12-11	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Texas, B
	<a href="#">QPD95922</a>	2020-12-11	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Texas, B
	<a href="#">QPO15025</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15037</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15049</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15061</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15073</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15085</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15097</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15109</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15121</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15133</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15145</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15157</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco
	<a href="#">QPO15169</a>	2020-12-09	Severe acute respiratory s...	1273	complete	surface glycoprotein	USA: Minnesco

<

Page 1 of 20

>

# Sequences

Aligned using **Muscle**

\*tried T-Coffee but was ineffective\*

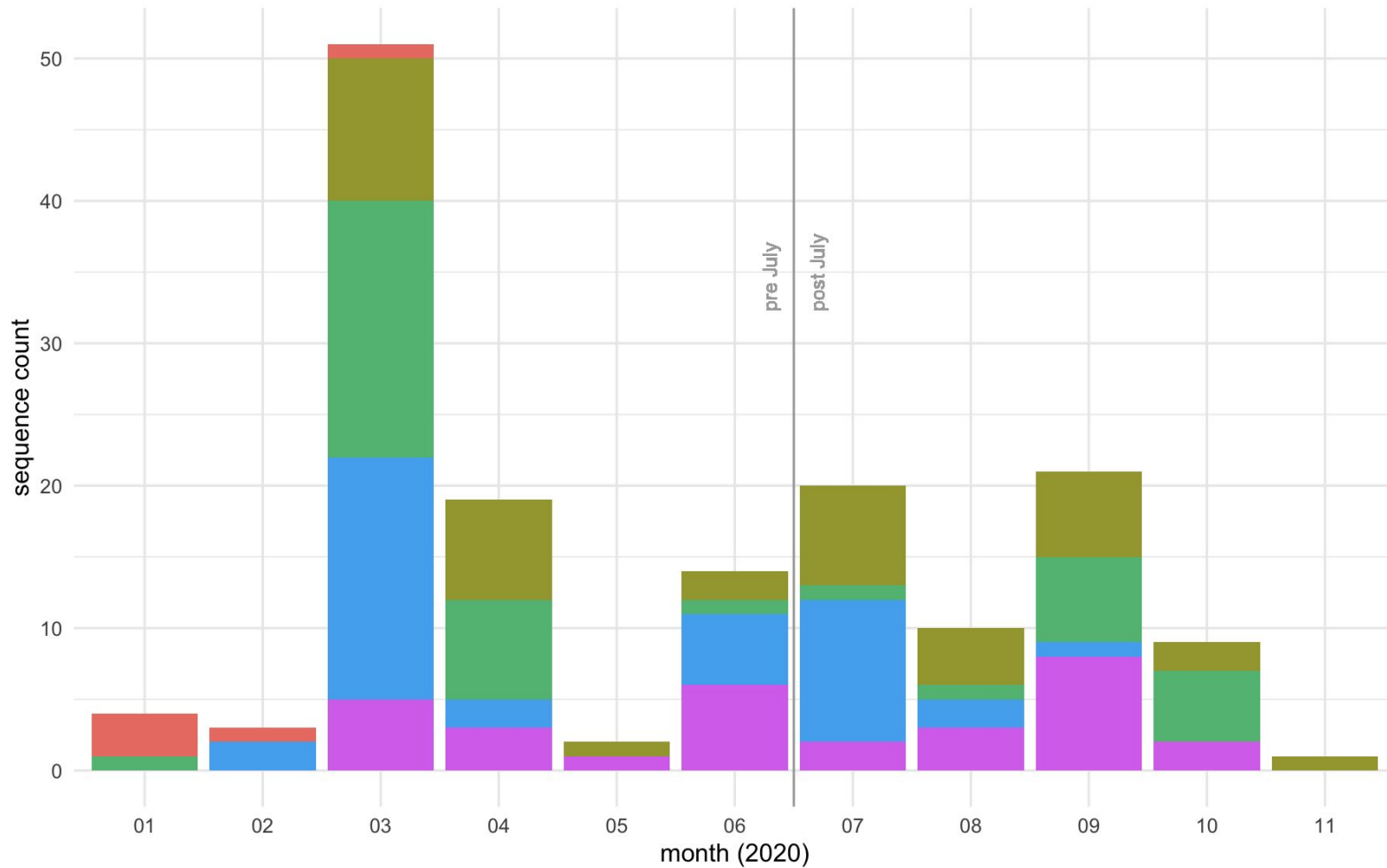
[Alignments](#)[Result Summary](#)[Phylogenetic Tree](#)[Results Viewers](#)[Submission Details](#)[Download Alignment File](#)[Show Colors](#)

```
>QPL23245.1 |surface glycoprotein [Severe acute respiratory syndrome coronavirus 2]|USA
MFVFLVLLPLVSSQCVNLTTRTQLPPAYTNSFTRGVYYPDKVFRSSVLHSTQDLFLPFFS
NVTWFWHAIHVSGTNGTKRFDNPVLPFNDGVYFASTEKSNIIRGWIFGTTLDSTQSLIIV
NNATNVVIKVCFFQFCNDPFLGVYHKNKSWMESEFRVYSSANNCTFEYVSQPFLMDLE
GKQGNFKNLREFVFKNIDGYFKIYSKHTPINLVRDLPQGFSALEPLVDLPIGINITRFQT
LLALHRSYLTTPGDSSSGWXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
XXXXXXXXXXXXXXXXXXXXPTESIVRFPNITNLCPFGEVFNATRFASVYAWNRRRISN
CVADYSVLVNSASFSTFKCYGVSPTKLNDLCFTNVYADSFVIRGDEVQRQIAPGQTGKIAD
YNYKLPDDFTGCVIAWNSNNLDSKVGGNYNLYRLFRKSNLKPFERDISTEIQAGSTPC
NGVEGFNCYFPLQSYGFQPTNGVGYPYRVVLSFELLHAPATVCGPKKSTNLVKNKCVN
FNFNGLTGTGVLTESNKKFLPFQGFGRDIADTTDAVRDPQTLIELDITPCSFGGVSVITP
GTNTSNQVAVLYQGVNCTEVPVAIHADQLTPTRVYSTGSNVFQTRAGCLIGAEHVNNYSY
ECDIPIGAGICASYQTQTNSPRRARSVASQSI IAYTMSLGAENSVAYSNNIAIPTNFTI
SVTEILPVSMTKTSVDCTMYICGDSTECNLLQYGSFCTQLNRALTGIAVEQDKNTQE
VFAQVQKIYKTPPIKDFGGFNFSQLLPDPSKPSKRSFIEDLLFNKVTLADAGFIKQYGDC
LGDIAARDLICAQKFNGLTVLPLLLTDEMIAQYTSALLAGTITSGWTFGAGAALQIPFAM
QMAYRFNGIGVTQNVLYENQKLIANQFNSAIGKIQDSLSTASALGKLQDVVNQNAQALN
TLVKQLSSNFGAISSVLNDILSRDKVEAEVQIDRLITGRQLSLQTYVTQQLIRAAEIRA
SANLAATKMSECVLGQSKRVDFCGKGYHLMSPQSAHPGVVFLHVTVVPAQEKNFTTAPA
ICHDGKAHFPREGVFVSNGTHWVFVTRQNFYEPQIITTDNTFVSGNCDVVIGIVNNTVYDP
LQPELDSFKEELDKYFKNHTSPDVLGDISGINASVVNIQKEIDRLNEVAKNLNESLIDL
QELGKYEQYIKWPWYIWLGFIAGLIAIVMVTIMLCMTSCCSCCLKGCCSCGSCCKFDEDD
SEPVLGKGVKLHYT
```

```
>QPL23329.1 |surface glycoprotein [Severe acute respiratory syndrome coronavirus 2]|USA
MFVFLVLLPLVSSQCVNLTTRTQLPPAYTNSFTRGVYYPDKVFRSSVLHSTQDLFLPFFS
NVTWFWHAIHVSGTNGTKRFDNPVLPFNDGVYFASTEKSNIIRGWIFGTTLDSTQSLIIV
```

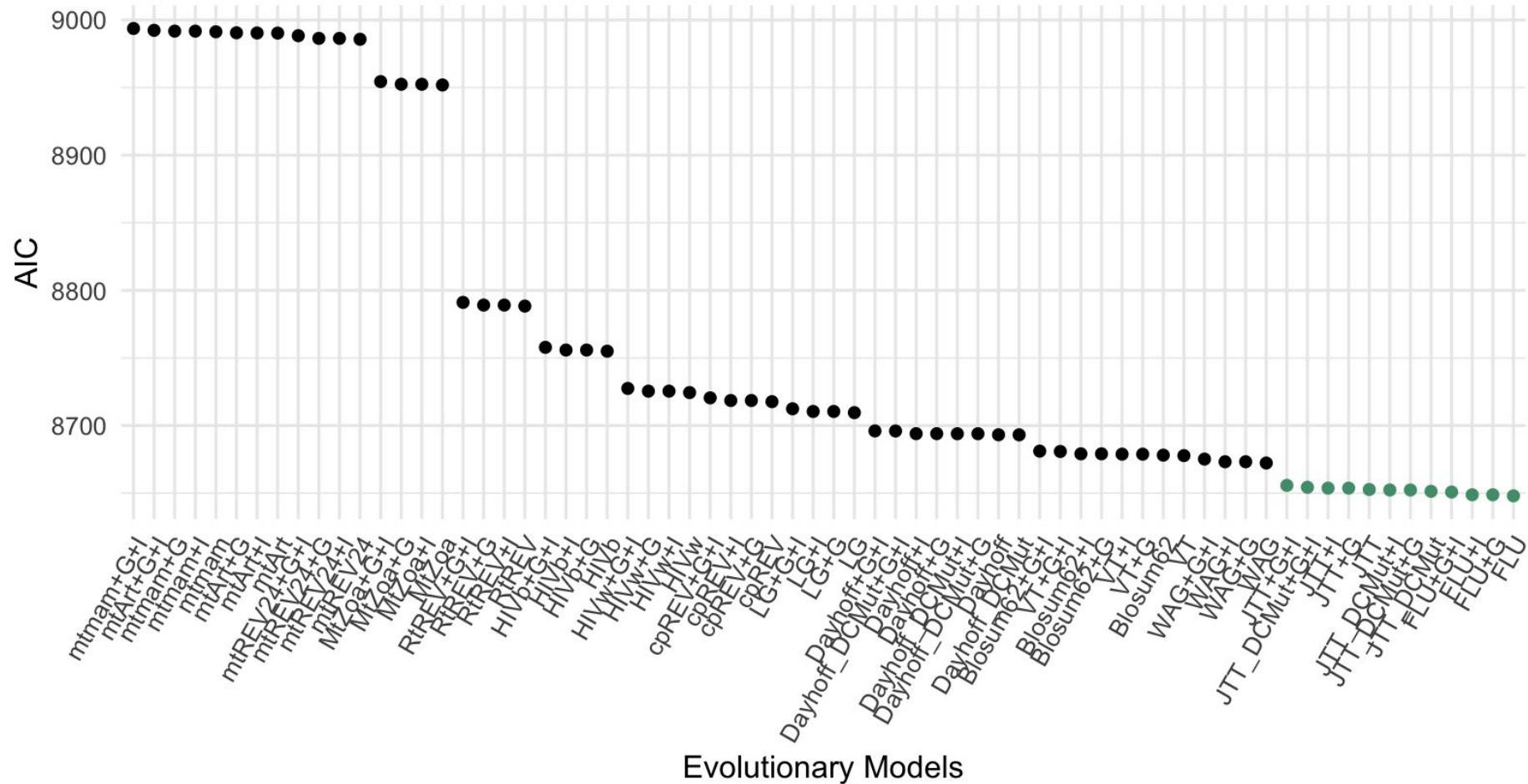
# COVID spike protein

Sequence counts by US regions **East**, **Midwest**, **South**, and **West**, including outgroup **China**, split into time intervals



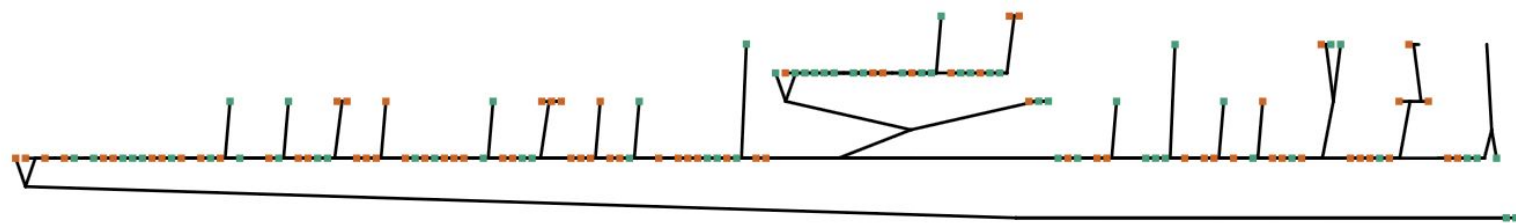
## AIC measures

for finding **best evolutionary models** to be used in distance matrix calculation

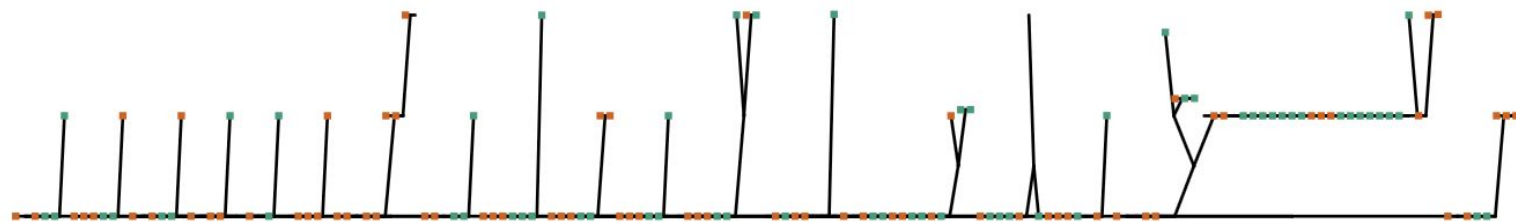


Model selection based on time intervals **Pre July** and **Post July**

**FastME**



**NJ**

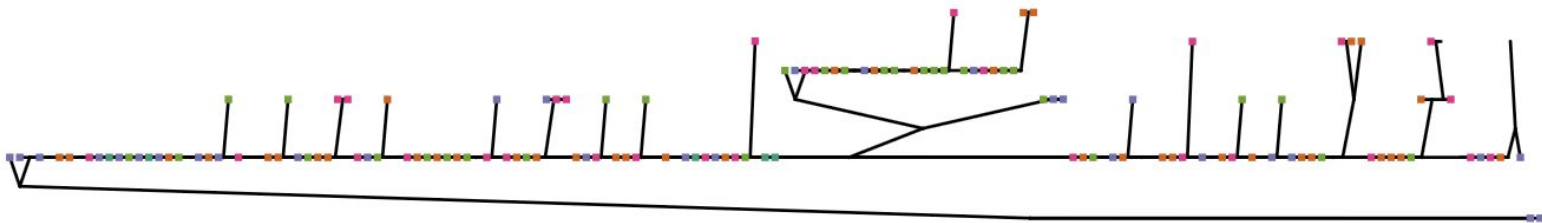


**UPGMA**

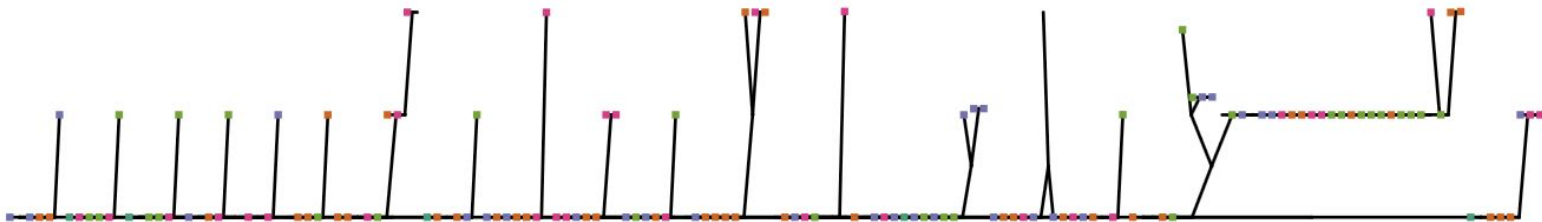


Model selection based off of US region East Midwest South and West with outgroup China

FastME



NJ



UPGMA





# Conclusions

**No pattern** in clustering was observed **by time period or US region**.

UPGMA suggested all **outgroup samples** (China, early 2020) to be **clustered with a large group of other samples**.

This result **suggests genetic variance largely occurred within the US**, without geographic limits.

A **larger sample size** and more **refined sample selection** methods **are warranted** to confirm findings.

