

Using RNA Secondary Structure to Analyze Variations in the S (spike) Protein of SARS-CoV-2 and related Coronaviruses

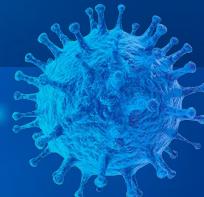
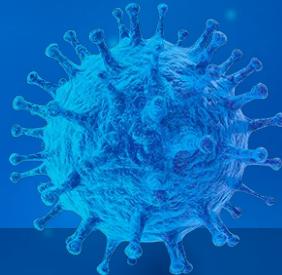


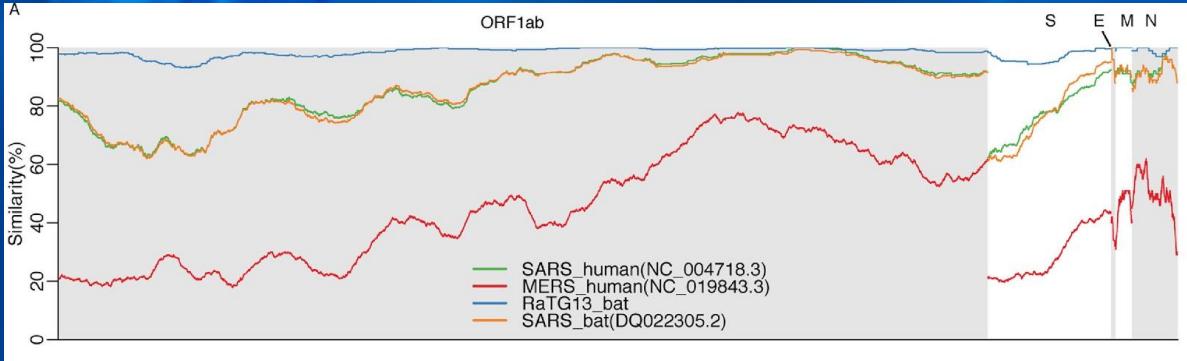
By Lawrence Chillrud



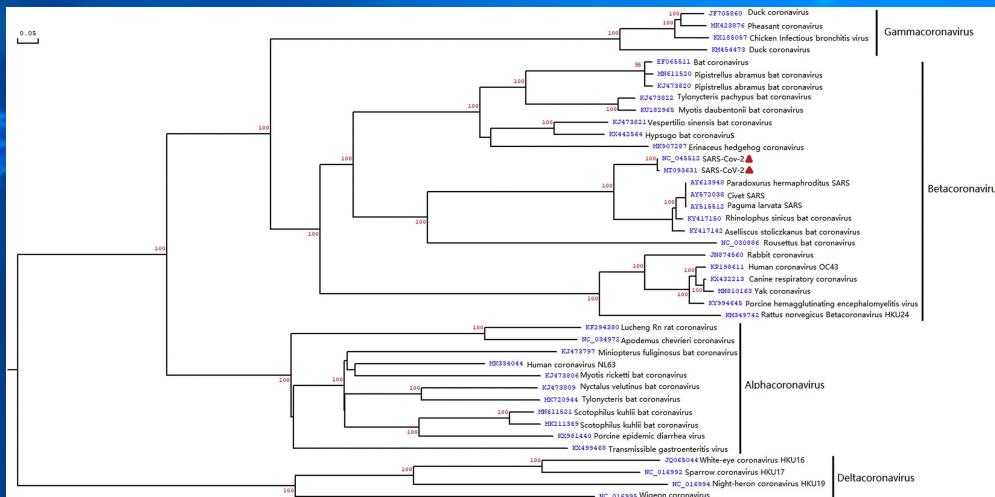
1. Introduction

RNA Secondary Structure, Coronaviruses, S (spike) Protein

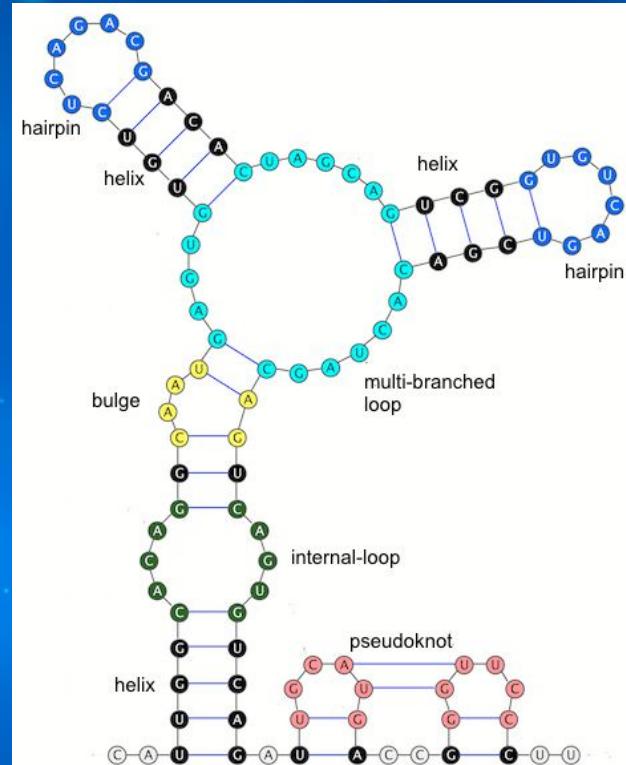




Reference: Wen, Feng, et al. "Identification of the Hyper-Variable Genomic Hotspot for the Novel Coronavirus SARS-CoV-2." *Journal of Infection*, 4 Mar. 2020, doi:10.1016/j.jinf.2020.02.027.



Reference: Li, Chun, et al. "Genetic Evolution Analysis of 2019 Novel Coronavirus and Coronavirus from Other Species." *Infection, Genetics and Evolution*, vol. 82, 10 Mar. 2020, p. 104285., doi:10.1016/j.meegid.2020.104285.

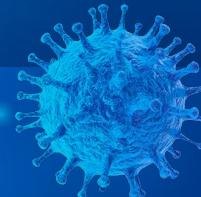
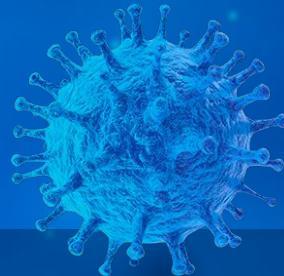


Reference: Mamuya, Adane & Merelli, Emanuela & Tesei, Luca. (2016). A Graph Grammar for Modelling RNA Folding. *Electronic Proceedings in Theoretical Computer Science*. 231. 31-41. 10.4204/EPTCS.231.3.



2. Methods

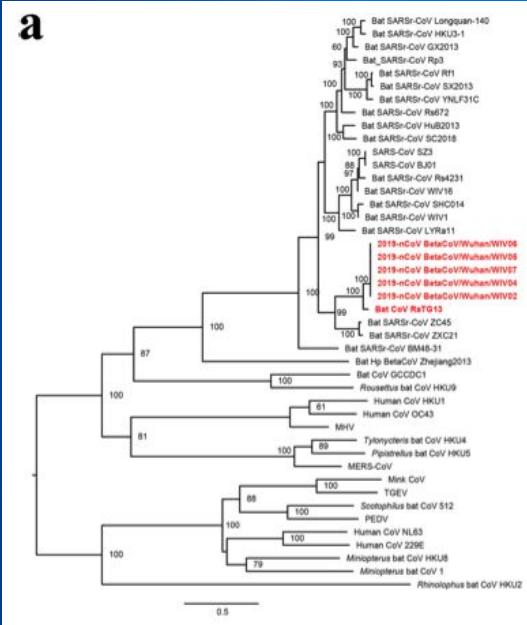
Sequence Selection, Alignment, Folding, Analysis



Interspecies Sequence Selection



a



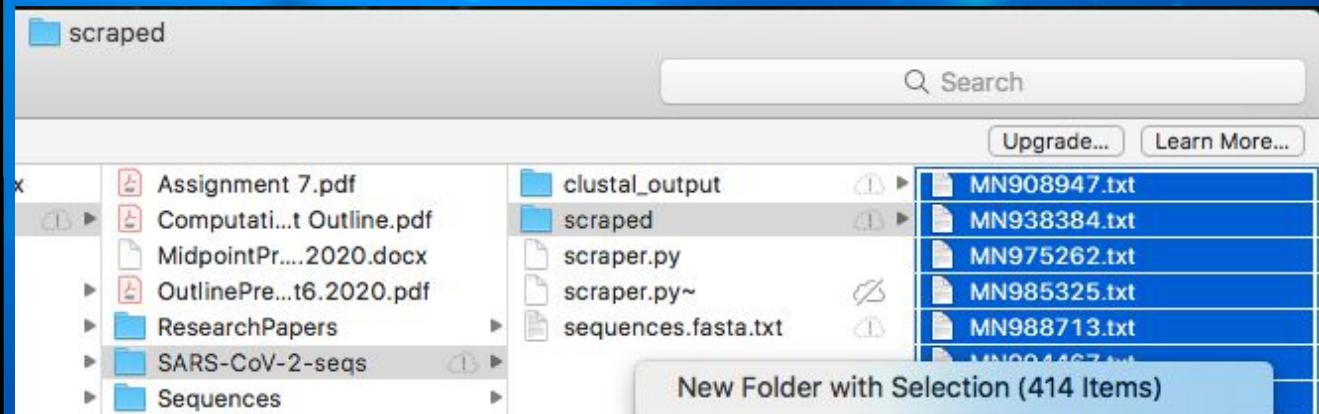
Reference: Zhou, Peng, et al. "A Pneumonia Outbreak Associated with a New Coronavirus of Probable Bat Origin." *Nature News*, Nature Publishing Group, 3 Feb 2020, www.nature.com/articles/s41586-020-1202-7.

- ▶ 5 beta coronaviruses I am comparing across the beta coronavirus genus.
 - ▷ MN996532.1 (RaTG13)
 - ▷ NC_004718.3 (SARS)
 - ▷ NC_014470.1 (BM48-31 Bat CoV)
 - ▷ NC_019843.3 (MERS)
 - ▷ NC_045512.2 (SARS-CoV-2)
- ▶ This makes up the interspecies data set.

Intraspecies Sequence Selection

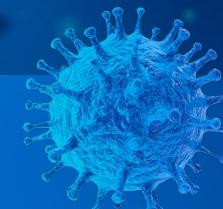


- ▶ 414 intraspecies sequences that capture the SARS-CoV-2 S gene.



Left: Screenshot of my python code to scrape genomic information from NCBI. Right: Screenshot of my folder containing the 414 scraped sequences encoding the S protein.

Sequence Alignment



Left: Screenshot of the files output by the Clustal Omega algorithm for multiple sequence alignment.

Bottom: Screenshot of Clustal Omega algorithm's webserver:
<https://www.ebi.ac.uk/Tools/msa/clustalo/>

The figure shows two windows of the ClustalW application. The left window displays a multiple sequence alignment of 12 entries (NC_019843.3, MN996532.1, NC_045512.2, NC_004718.3, NC_014470.1) against a reference sequence (NC_019843.3). The right window shows the resulting phylogenetic tree with the same set of sequences, rooted at a guide tree. A context menu is open over the tree, with options like 'phylogenetic_tree.phy.txt' and 'phylogenetic_tree.phy.tree'. The title bar of the right window indicates 'Percent Identity Matrix - created by Clustal2.1'.

Clustal Omega

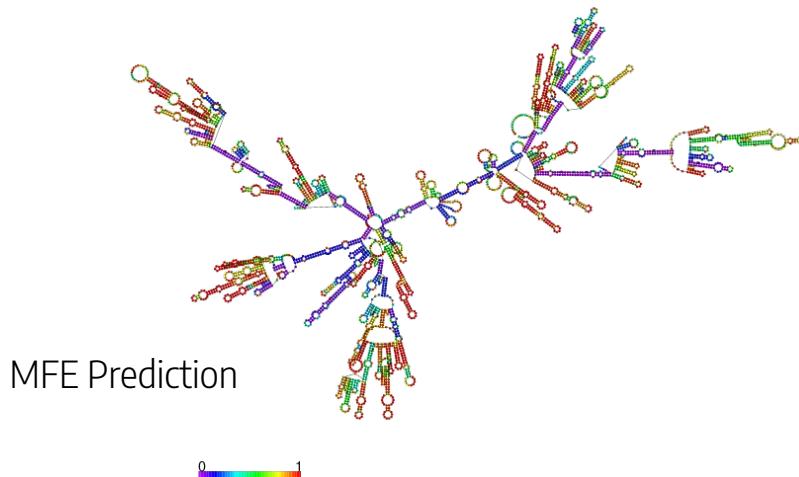
[Input form](#) | [Web services](#) | [Help & Documentation](#) | [Bioinformatics Tools FAQ](#)

Multiple Sequence Alignment

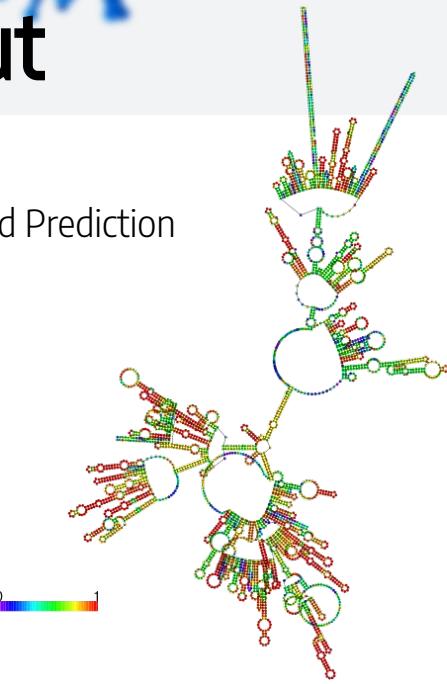
Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to align more sequences. For the alignment of two sequences please instead use our pairwise sequence alignment tools.

Important note: This tool can align up to 4000 sequences or a maximum file size of 4 MB.

Vienna RNAfold: Example Output

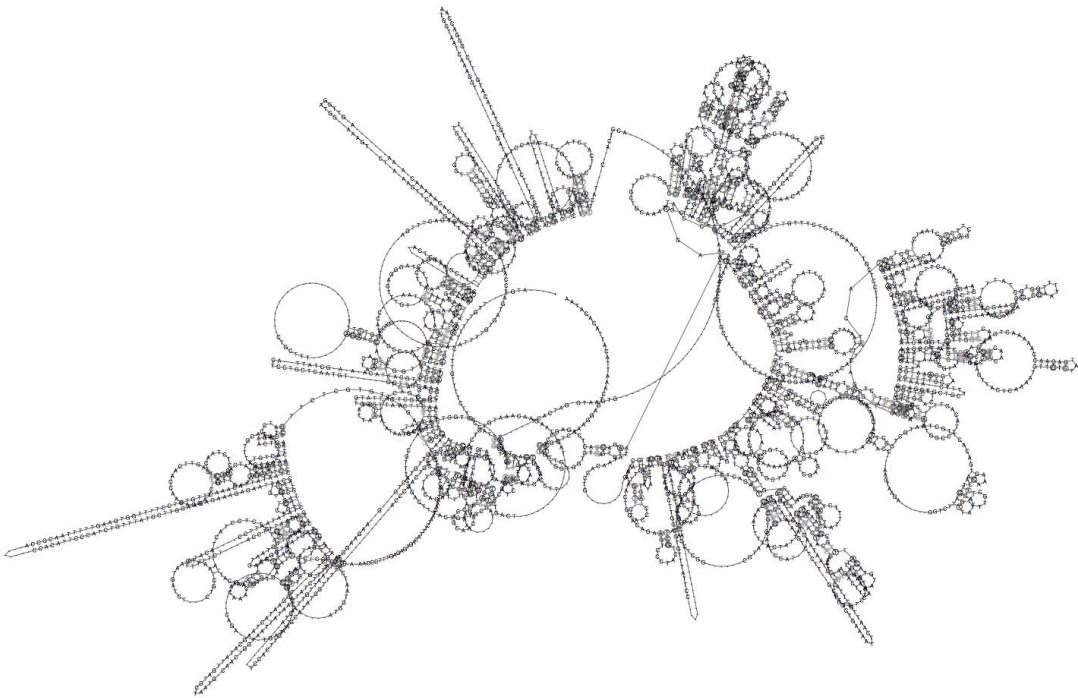


Centroid Prediction



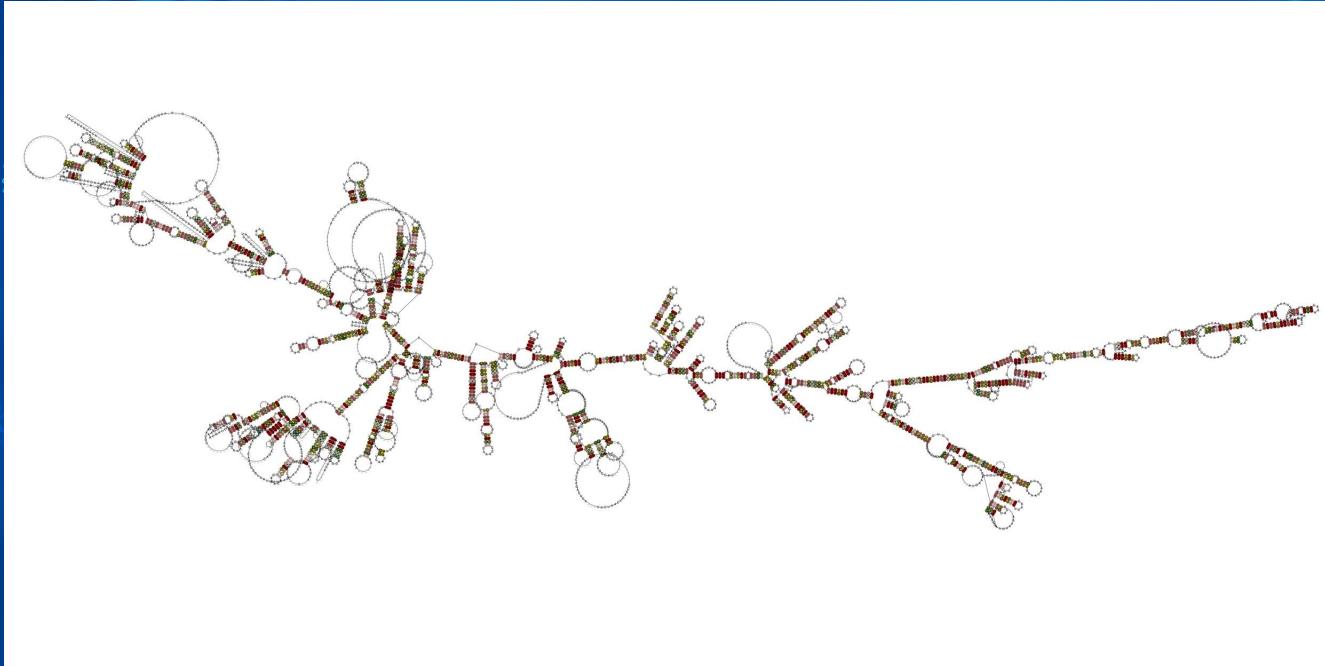
Downloaded output from Vienna RNAfold's webserver:
<http://rna.tbi.univie.ac.at/cgi-bin/RNAWebSuite/RNAfold.cgi>

ViennaRNA alidot



"The program `alidot` detects conserved secondary structure elements in relatively small sets of RNAs by combining multiple sequence alignments and secondary structure predictions. Both a (good) sequence alignment and predicted secondary structure predictions for each sequence in the alignment must be provided as inputs. `alidot` works either with predicted mfe structures, or with base pairing probability matrices. The basic idea behind `alidot` is to sort the individual base pairs by their *credibility* and to reduce the number of entries in the list by subsequent filtering steps until only those secondary structure elements are left that are consistently predicted."

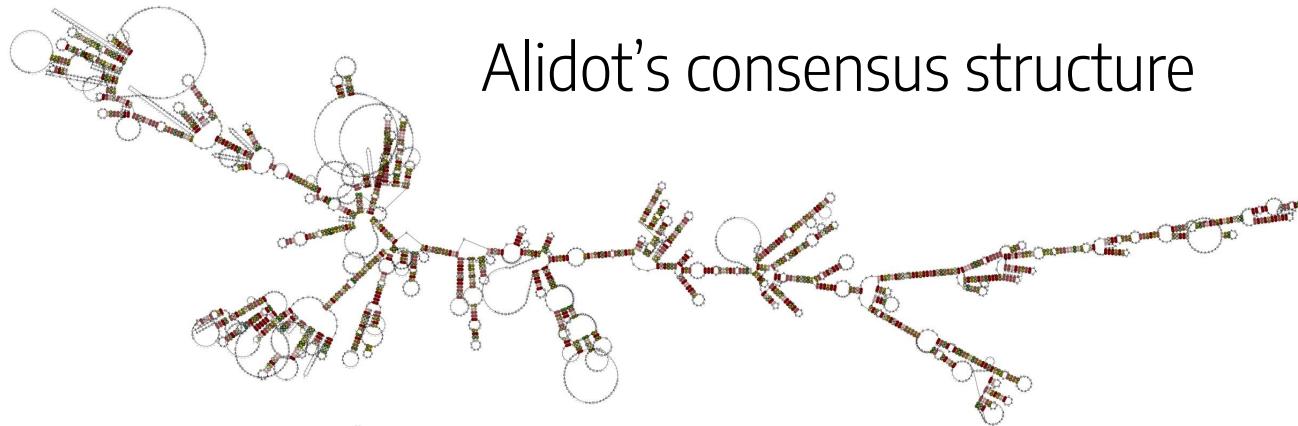
ViennaRNA alifold



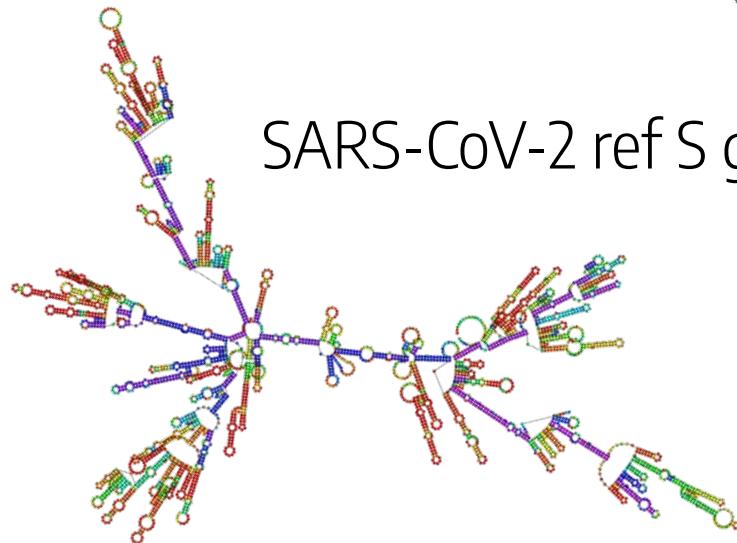
“The program `alifold` predict a consensus secondary structure of a set of aligned sequences..”

Screenshot of the consensus secondary structures predicted by ViennaRNA’s `alifold` program when given the interspecies dataset

Alidot's consensus structure



SARS-CoV-2 ref S gene



My Own Algorithm



CLUSTAL 0(1.2.4) multiple sequence alignment

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

GATACACTGAGTTCTAC-TGATGTTCTGGTAAACC-----TACAGAAA
-ATGT-----TTGTTTT-C-TTGTGTTATTGCCACTAG-----TTCTAG-TCA
-ATGT-----TTGTTTT-C-TTGTGTTATTGCCACTAG-----TCTCTAG-TCA
-ATGT-----TTATGTTCTATTATTCCTTACTGCTACTAGGTTGTA
GAAATT-TTGGCTTTCTCTGGCTTCTGGCTTCTGCAACGCTCAAGTGCGAAG

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

```

ACGTTGATGAGGGCAGATTCTTAA---GTCGCTGATGTTATTGAGGTGATA
          TGTAACTCTA---CAACTAGAACCTAGGTT---A
          TGTAACTCTTA---CAACCGAGAACCTCAATT---A
ACCTTGAC---CGGTGCAACCATCTTGTATGDTGTCAGC---T
          TGTTACACTAT---CTAAATAAAAGTC---C---A

```

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

ACAGACTTCTTGTAAAACTTGCCCTAGGCCATTGATTTCTAAGGCTGACGGT
TC—C—TGCATACACAACCTTCATA—CCCCGGTGTCTTACCCGTACAAA
CC—C—TGCATACACTAACTTCTCA—CACGTGGTTTATACACTCTGACAAA
TAATA—TCACTAACATCTTCATCA—TGGAGGGGTTTTACTATCTGTGAA
TAAGCT—TCACTAGACCTCTTCTGA—GGAGGGGTTTTTATTATTTGTGAC

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

```

TAT-ATACCCCTCAAGGCCGTCATACTTACATACTACATTACTAAGGTCTTT
TTTCAAGATCTTCAGTTTACA-TTAACTACAGG-----ATTGGTCT
TTTCAGATCCTCA-----GTTTACA-TTCAACTCAGG-----ACTGGTCT
TTTATGAGATCCTCA-----ACCTTACA-TTAACTACAGG-----ATTGGTCT
TTTAACTGCTTCAGTGTGTCACACTGG-----CCATTTC

```

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

CTTATCAGGGAGCACATTGTGATATGTATGTTACTCTGCAGGACATGCTACAGGCC
CCCTTCTT---TCCTCAATGTGACCCTGGTCCA---TGCTATACATGTTCCAGGGAC
CTTCTCT---TTTCAATGTTACTGGTCCA---TGCTATACATGTTCTGGGAC
CCTATT---ATTCTAATGTTACAGGGTTCA---TACTATT
CTTTTTA---ATACTAACCTTACTGGTATT---GACTTTAAA-----GTC

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

CTCCA---CAAAGTTGTTGTAAGCTAATTCTCAGGCGCTAAACAGTTGCTAA
ATGGTATTAAAAGGT----TTGATAAC--CCAGTTCTGCATTCAAGGA
TGTTGACTAAAGGGT----TTGATAAC--CTGTCTTCACATTAAATGA
---AACATACGT-----TTGCAAC--CTCTGCATACCTTTAAAGGA
ATGGTAAGCAGGAGGATTATTGATAAT-----CCAAACATTACATTGGTA

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

```

GGTTTGTCTGCGTATAGGAGCAGTCGCAATTCCACTG-
GGTC- -TATTTTGCTTCCACTGAGAAGTCATAATAAAAGAGGATGGATTT
GTGTT- -TATTTTGCTTCCACTGAGAAGTCATAATAAAAGAGGCTGGATTT
GTATT- -TATTTTGCTGCCAACAGAGAACATTAAGTTTCTGGCTGGTTGGTTT
GTGTT- -TATTTGGTCAACCGAGAAAATTAATGTTTCTGGGTTGGATTT

```

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

-----GCACTGTTATTAGGCCATCAGCGG
GTACTACCTTAGATTGAAAGCCAGCTCTACTATTGTTAATAACGCTACTAATG
GTACTACTTTAGATTGAAAGCCAGCTCTACTATTGTTAATAACGCTACTAATG
GTTCTTACCATGACAAACAGTCAGTCGGTATTATTAACAACTTACTAATG
TGTCGACATAGAACACACAACTCTGCTGTTCTTTAAATGTCACACAT

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

C T A C G A A A A A T T T A C C T G C T T T T A T G C T G G G T T C T C A G T T G G T A A T T T C C A G A
— T T A T T A A A G T C T G T G A A T T C A —
— T T T A A A G T C T G T G A A T T C A —
— T T A C G A G C A T G T A A C T T T G A —
— T T A G A T G T G T G A A C T T T A A —

NC_019843.3_21456-25517
MN996532.1_21545-25354
NC_045512.2_21563-25384
NC_004718.3_21492-25259
NC_014470.1_21391-25170

GTAATTA—ATGGGGCGCTTCAATCATACTCTAGTTCTTTGCCGATGGATGTGGC
ATCCATTATTTGGTGTTTATTACCAACAAAACCAA—AAAGTTGGATGGA
ATCATTTTTGGGTTTTATTACCAACAAAACCAA—AAAGTTGGATGGA
ACCTCTTCTTGGTTTAACTAACCA—TTGGGTAC—
ATCAATTTGGTGTCAATA—TTGG—A

```

531 (((((.....))))),..,((.....))),.....,((.....))),.....,((.....))),.....,
461 (((((.....))))),.....,((.....))),.....,((.....))),.....,((.....))),.....,((.....)),
461 (((((.....))))),.....,((.....))),.....,((.....))),.....,((.....))),.....,((.....)),
488 (((((.....))))),.....,((.....))),.....,((.....))),.....,((.....))),.....,((.....)),
489 (((((.....))))),.....,((.....))),.....,((.....))),.....,((.....))),.....,((.....))),.....,

```

My Own Algorithm



	1	2	3	4	5	6	7
SC2	A	U	G	-	C	A	U
Bat	A	U	G	C	C	A	U
SARS	A	C	A	U	U	C	U
BM48-31	A	C	A	G	U	C	U
MERS	A	C	U	A	G	C	U

(((.)))

	1	2	3	4	5	6	7
SC2	X	X	X	X	X	X	X
Bat	X	X	X	Y	X	X	X
SARS	X	Y	Y	Z	Y	Y	X
BM48-31	X	Y	Y	P	Y	Y	X
MERS	X	Y	Z	Q	Z	Y	X

(((.)))

```

> 274-292 log likelihood: -37.62211326680391random seq of same size log likelihood: -52.7956961917466
UUCGUUCCACUGAGAACUC
(((.....))).).

> 471-516 log likelihood: -119.79646528234582random seq of same size log likelihood: -112.13864671127861
AGAGUUUAUUCUAGUGCGAAUAAUUGCACUUUUGAAU AUGUCUCUC
(((.....(((((.....))))....))))....))).

> 1267-1286 log likelihood: -36.05579777198838random seq of same size log likelihood: -45.16881916213479
AUAAAUAACCAGAUAGAUUU
..(((.....))).).

> 1294-1304 log likelihood: -26.6124989581944random seq of same size log likelihood: -28.890622866778113
GCGUUAUGCU
((.....))..

> 1571-1592 log likelihood: -45.9933061656108random seq of same size log likelihood: -57.10043842636637
UUGGGGACCUAAAAGUCUACU
..(((.....))).).

> 1822-2203 log likelihood: -49.70515948468804random seq of same size log likelihood: -44.66402909690834
CAGUGCUAUGGCCAACAUUC
..(((.....))).).

> 2244-2264 log likelihood: -55.533570155296665random seq of same size log likelihood: -33.862962813153295
UGCAGCAUCUUUUGUUGCAA
(((.....))).).

> 2276-2297 log likelihood: -49.93569809155589random seq of same size log likelihood: -48.72767996143369
UUGUACACAAUUAACCGUGCU
..(((.....))).).

> 2452-2471 log likelihood: -47.222710506602894random seq of same size log likelihood: -56.07092404684964
UUGAAGAUCAUCUUUCAAC
((.....))).).

> 2740-2758 log likelihood: -25.877639661851013random seq of same size log likelihood: -43.556302151640764
AUGUUCUCUJAUGAGAACCA
..(((.....))).).

> 3045-3065 log likelihood: -39.01146005367711random seq of same size log likelihood: -52.74933696093571
GCAGAAAUCAGACUUUCGCU
((.....))).).

> 3165-3179 log likelihood: -25.621606268582298random seq of same size log likelihood: -41.41538476623717
GCACCUCAUGGUGUA
((.....))).).

> 3435-3455 log likelihood: -44.41372625555621random seq of same size log likelihood: -39.82257756574254
GACUCAUCAAGGAGGAGUUA
((.....))).).

> 3578-3593 log likelihood: -23.686253578019425random seq of same size log likelihood: -45.7479335967179
AAUAGAAUCUCUACU
..(((.....))).).

> 3676-3688 log likelihood: -25.95255005869852random seq of same size log likelihood: -31.210463338852193
CCAUAGUAAUGGU
((.....))).).

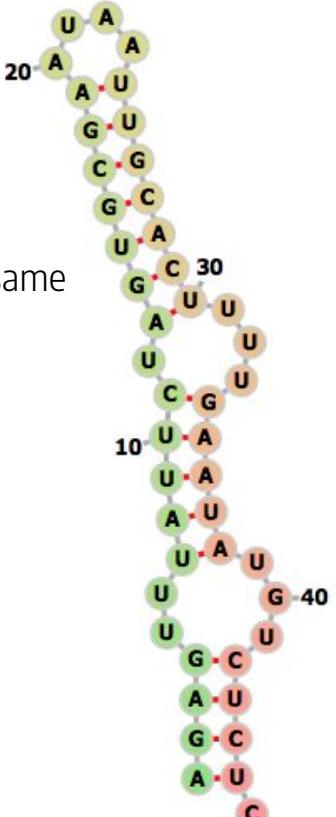
```



Largest potentially conserved secondary structure across interspecies dataset found at positions 471-516.

46 nt long.

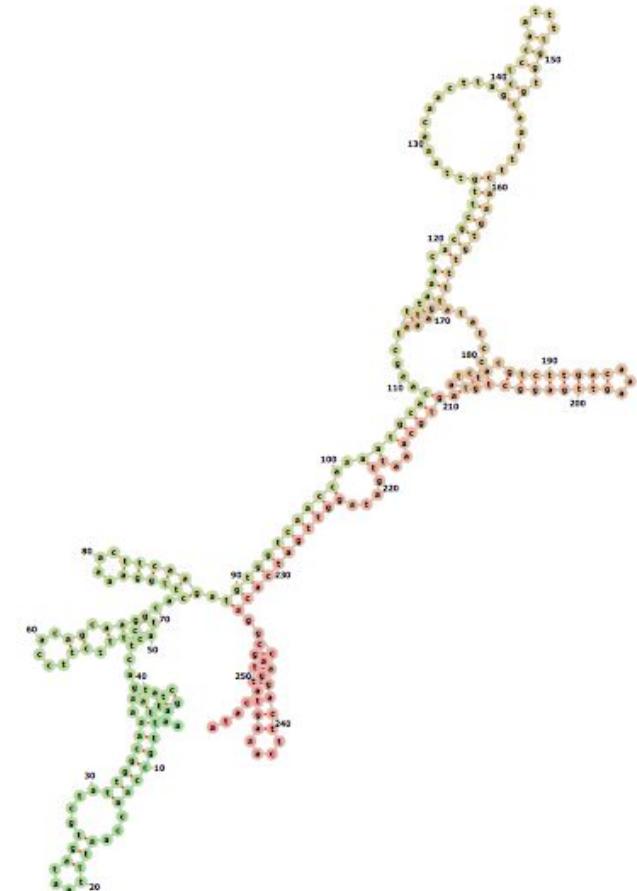
Less “likely” than a random, uniformly distributed seq. of same size.



```
#  
#  
# Percent Identity Matrix – created by Clustal2.1  
#  
#
```

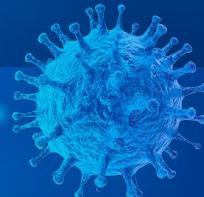
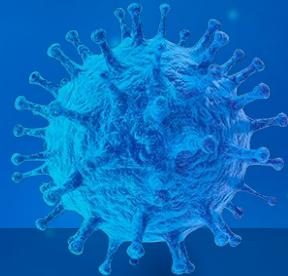
1: NC_019843.3_21456–25517	100.00	53.01	52.64	51.69	51.79
2: MN996532.1_21545–25354	53.01	100.00	93.15	73.89	70.24
3: NC_045512.2_21563–25384	52.64	93.15	100.00	73.97	70.40
4: NC_004718.3_21492–25259	51.69	73.89	73.97	100.00	71.79
5: NC_014470.1_21391–25170	51.79	70.24	70.40	71.79	100.00



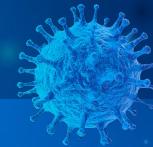


3. Results

Difficulty getting past the bugs...



Predictions...



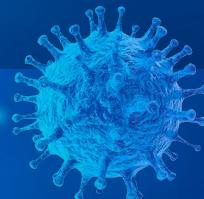
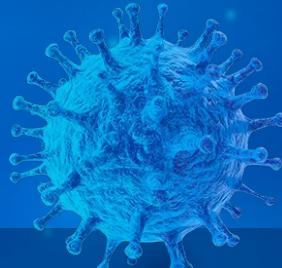
When only looking at SARS-CoV-2 and RaTG13, found 54 potential conserved structures, the longest one (that is more likely than expected) is 214 nucleotides long... (32/54 “likely”)

As we add in species to look at, we get less potential conserved structures found, they’re usually shorter, and usually they aren’t as likely as in the lower dimensional cases...

When looking across all 5 species, found 15 potential conserved structures, the longest one is 45 nucleotides long, although it’s not as likely as we would expect. (10/15 “likely”)

4. Conclusions

Summary, Limitations, Prospective Future Work



Summary, Limitations, Future Work



Picked out potential conserved secondary structures but with limited confidence – bugs really dragged this project down along with limited resources.

Loads of assumptions were made that decrease the accuracy of the project – e.g., collapsation of the quintets, only utilizing the SARS-CoV-2 secondary structure, abbreviated application of the MSA. Notion of “likely” is fuzzy as well; statistical significance analysis needed.

Future work would involve getting rid of the limitations discussed, making the methods more valid! Also a need to take into account thermodynamic data as other analyses do...

Questions?

