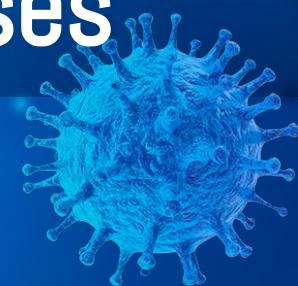


# Using RNA Secondary Structure to Analyze Variations in the S (spike) Protein of SARS-CoV-2 and related Coronaviruses

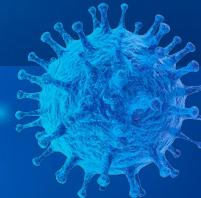
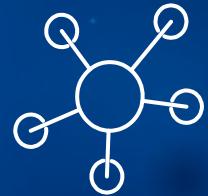
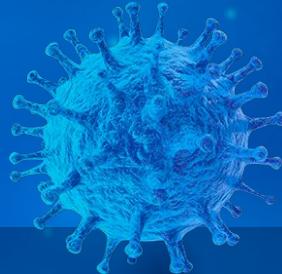


By Lawrence Chillrud

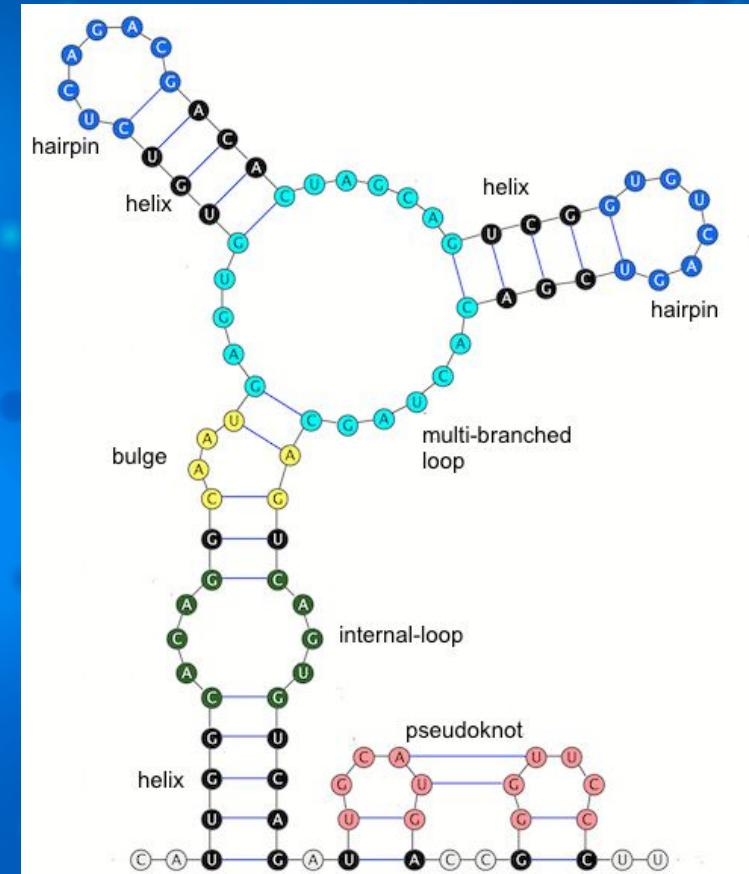
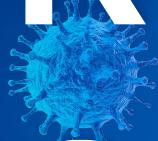


# 1. Introduction / Background

RNA Secondary Structure, Coronaviruses, S (spike) Protein

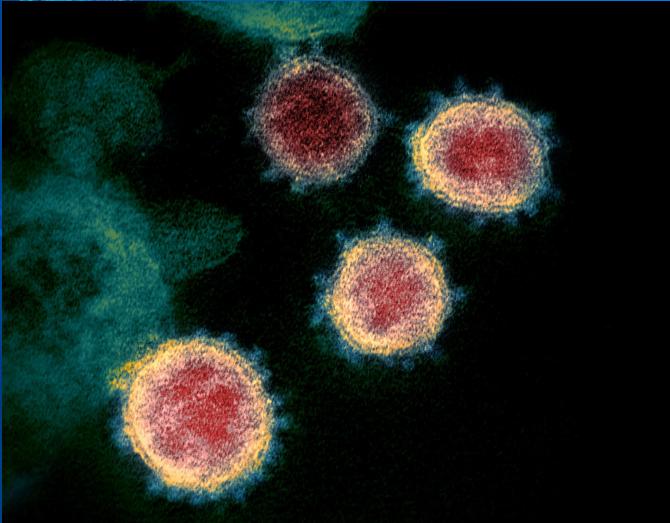


# RNA Secondary Structure



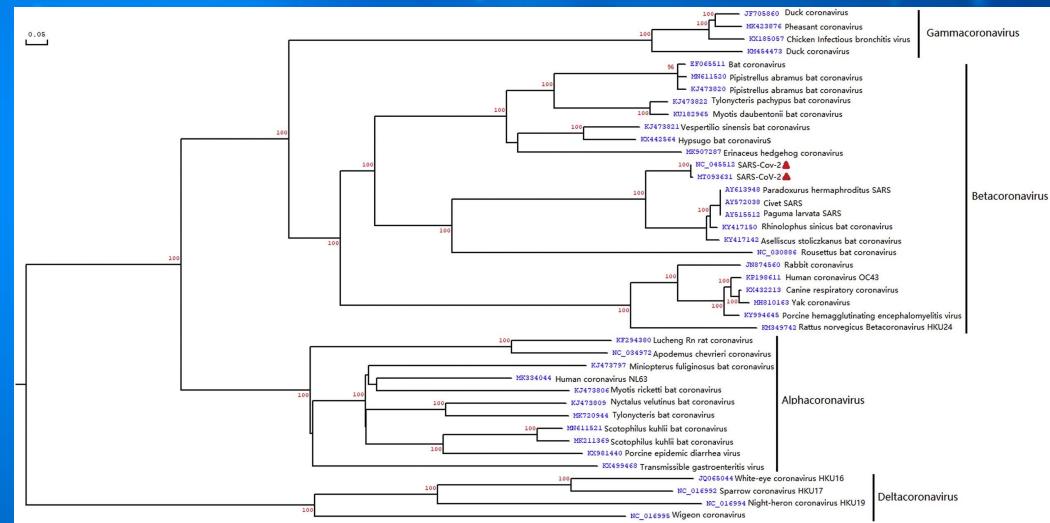
Reference: Mamuya, Adane & Merelli, Emanuela & Tesei, Luca. (2016). A Graph Grammar for Modelling RNA Folding. Electronic Proceedings in Theoretical Computer Science. 231. 31-41. 10.4204/EPTCS.231.3.

# Coronaviruses and SARS-CoV-2



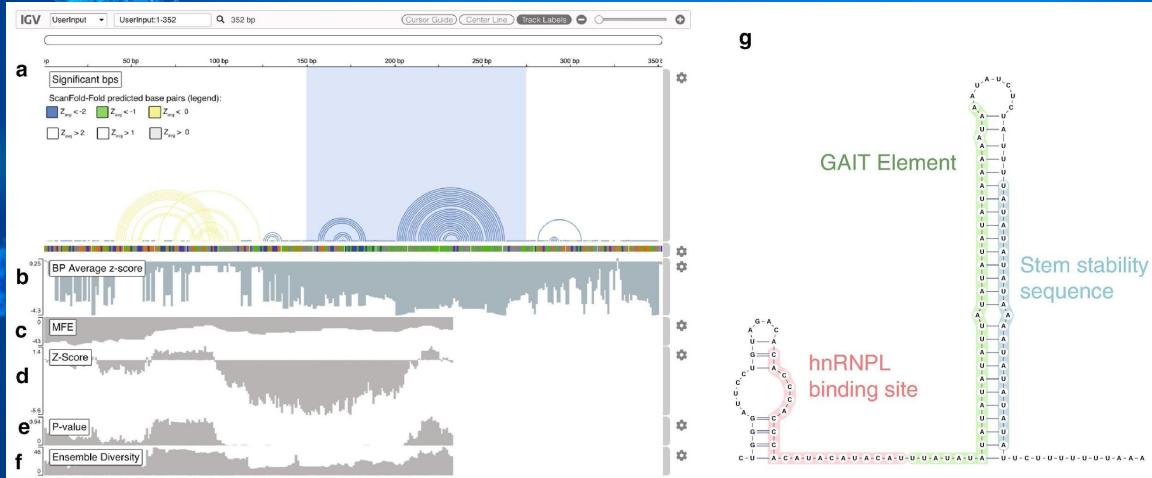
Reference: NIAID-RML.

Link: <https://www.flickr.com/photos/niad/49534865371/>



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.

# Secondary Structure Prediction Algorithms



Name	Description	Knots [Now 1]	Links	References
<b>CentroidFold</b>	Secondary structure prediction based on generalized centroid estimator	No	<a href="#">sourcecode</a> <a href="#">webserver</a>	[1]
<b>CentroidHomfold</b>	Secondary structure prediction by using homologous sequence information	No	<a href="#">sourcecode</a> <a href="#">webserver</a>	[2]
<b>ContextFold</b>	An RNA secondary structure prediction software based on feature-rich trained scoring models.	No	<a href="#">sourcecode</a> <a href="#">webserver</a>	[3]
<b>CONTRAFold</b>	Secondary structure prediction method based on conditional log-linear models (CLLMs), a flexible class of probabilistic models which generalize upon SCFGs by using discriminative training and feature-rich scoring.	No	<a href="#">sourcecode</a> <a href="#">webserver</a>	[4]
<b>Crumple</b>	Simple, cleanly written software to produce the full set of possible secondary structures for one sequence, given optional constraints.	No	<a href="#">sourcecode</a> <a href="#">webserver</a>	[5]
<b>Cy3Fold</b>	Secondary structure prediction method based on placement of helices allowing complex pseudoknots.	Yes	<a href="#">webserver</a>	[6]
<b>E2Efold</b>	A deep learning based method for efficiently predicting secondary structure by differentiating through a constrained optimization solver, without using dynamic programming.	Yes	<a href="#">link</a> <a href="#">sourcecode</a>	[7][8]
<b>GTFold</b>	Fast and scalable multicore code for predicting RNA secondary structure.	No	<a href="#">sourcecode</a>	[9]
<b>IPKnot</b>	Fast and accurate prediction of RNA secondary structures with pseudoknots using integer programming.	Yes	<a href="#">sourcecode</a> <a href="#">webserver</a>	[10]
<b>KineFold</b>	Folding kinetics of RNA sequences including pseudoknots by including an implementation of the partition function for knots.	Yes	<a href="#">linubinary</a> <a href="#">webserver</a>	[11][12]
<b>Mfold</b>	MFE (Minimum Free Energy) RNA structure prediction algorithm.	No	<a href="#">sourcecode</a> <a href="#">webserver</a>	[13]
<b>pKiss</b>	A dynamic programming algorithm for the prediction of a restricted class (H-type and kissing hairpins) of RNA pseudoknots.	Yes	<a href="#">sourcecode</a> <a href="#">webserver</a>	[14]
<b>Pknots</b>	A dynamic programming algorithm for optimal RNA pseudoknot prediction using the nearest neighbour energy model.	Yes	<a href="#">sourcecode</a>	[15]

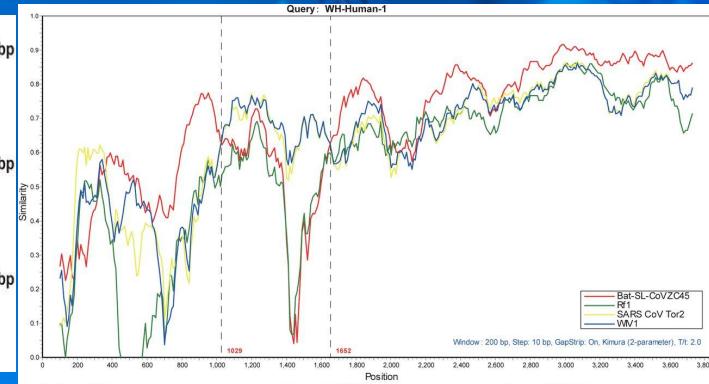
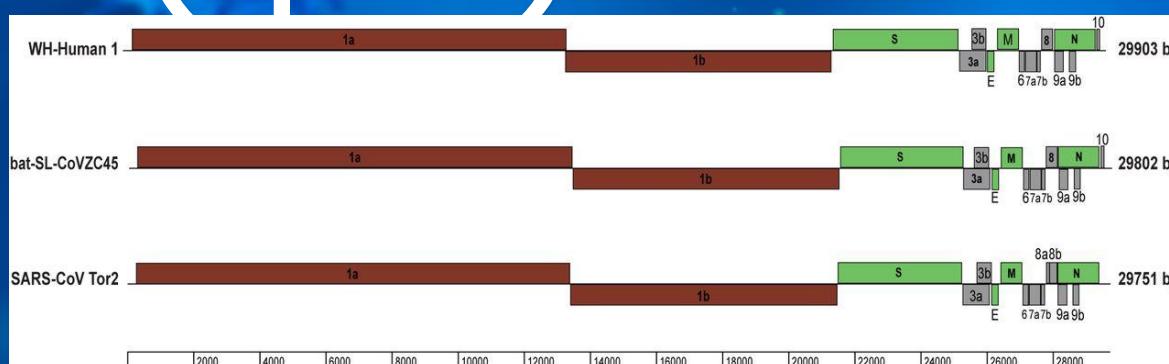
Reference:  
[https://en.wikipedia.org/wiki/List\\_of\\_RNA\\_structure\\_prediction\\_software](https://en.wikipedia.org/wiki/List_of_RNA_structure_prediction_software)



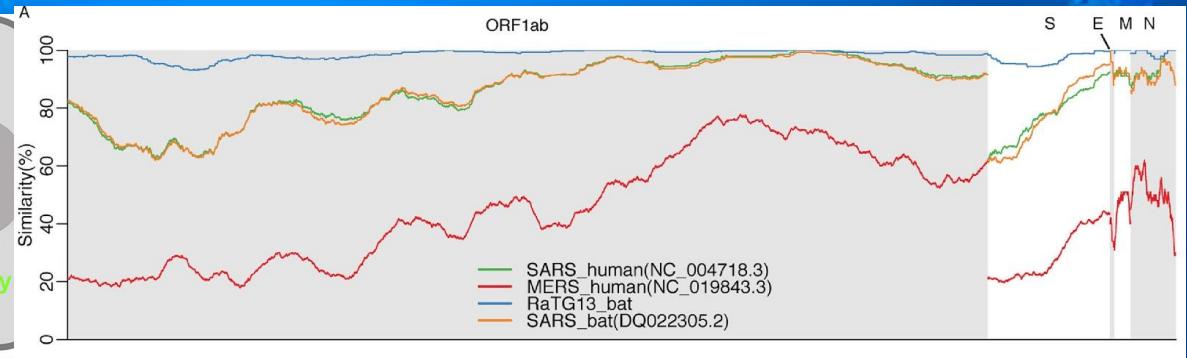
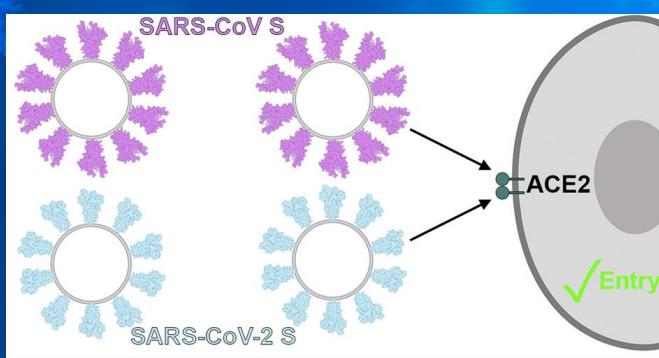
Reference: Andrews, Ryan J., et al. "Mapping the RNA Structural Landscape of Viral Genomes." *Methods*, 8 Nov. 2019, doi:10.1016/j.ymeth.2019.11.001.

Dynamic Programming Algorithms, Nearest Neighbor Classifiers, NLP approaches involving CFGs, Other approaches

# S(spike) Protein



Reference: Wu, Fan, et al. "Complete Genome Characterisation of a Novel Coronavirus Associated with Severe Human Respiratory Disease in Wuhan, China." *BioRxiv*, 2020, doi:10.1101/2020.01.24.919183.

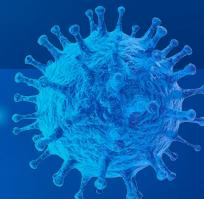
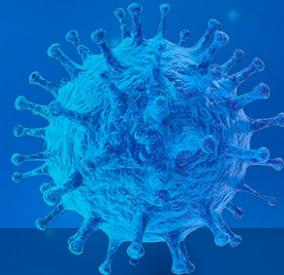


Reference: Walls, Alexandra C., et al. "Structure, Function, and Antigenicity of the SARS-CoV-2 Spike Glycoprotein." *Cell*, 9 Mar. 2020, doi:10.1016/j.cell.2020.02.058.

Reference: Wen, Feng, et al. "Identification of the Hyper-Variabile Genomic Hotspot for the Novel Coronavirus SARS-CoV-2." *Journal of Infect*

## 2. Task / Questions

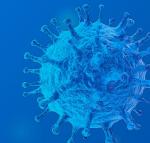
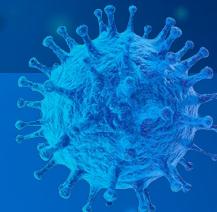
Can we link secondary structure to variation and even selection?

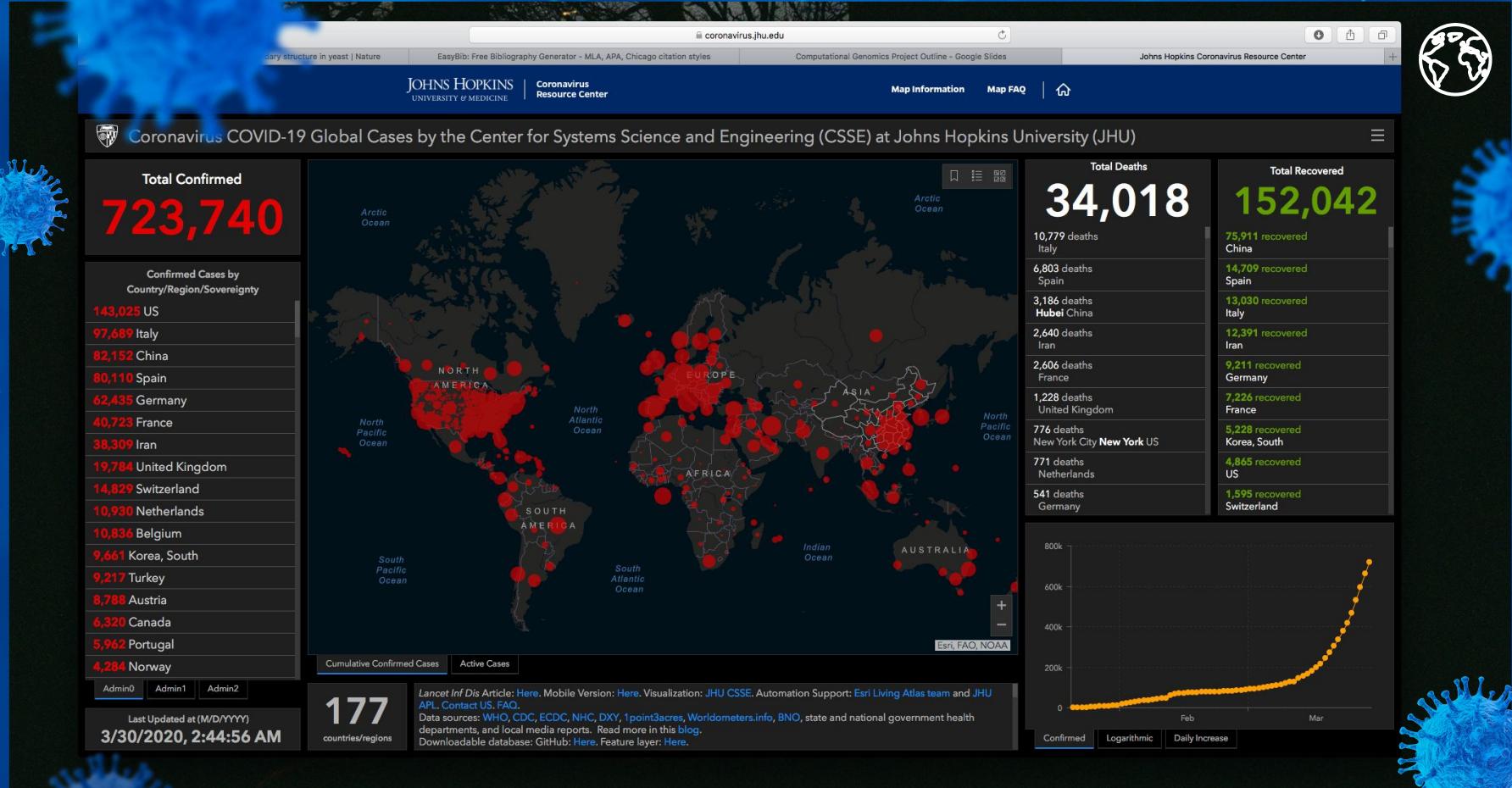


## 2. Task / Questions



- ▶ Can we link mutation rates / the mutations we see in the spike protein regions across the coronavirus genome with what we see in the corresponding secondary structures?
- ▶ What about across the SARS-CoV-2 genome?
- ▶ From any links we can establish between variation in the primary and secondary structures, what do they tell us about selection processes?
- ▶ Do structural insights help at all with developing a vaccine to exploit potential structural weaknesses?





# 3. Project Plan



Determine + Gather Data, Determine + Run Software, Analyse Findings

# Project Plan



## Step 1

Figure out what the relevant RNA sequences I want to analyze are, then gather them. Both within SARS-CoV-2 genome as well as outside of it, i.e. in the broader coronavirus genome.

## Step 2

Determine appropriate secondary structure prediction software + mutation detection / seq. alignment software and apply it to my collection of sequences gathered in Step 1.

## Step 3

Write code to organize and analyse the output from Step 2 in the context of my task / questions.



# Rough Timeline



Determine and  
Gather all Relevant  
Sequences

By 4/12

Run Software on  
Seqs.

By 4/18

Synthesize it all +  
write up findings

By 4/30

Determine  
Relevant Software

By 4/15

By 4/26

Write + run code to  
organize + analyse  
output

# Questions?

