# Tracking Recombination in the Evolution of SARS-CoV-2

Kristina Moen

UC Davis Math REU 2021

August 11, 2021

# Motivation - A Deadly and Mutating Virus

# Outline

# SARS-CoV-2 is a Single-Stranded RNA Virus



Spike Glycoprotein S
E Protein
Membrane Protein M
RNA genome
Envelope

Bases:

Adenine
Cytosine
Uracil
Guanine

~ 30,000
nucleotides

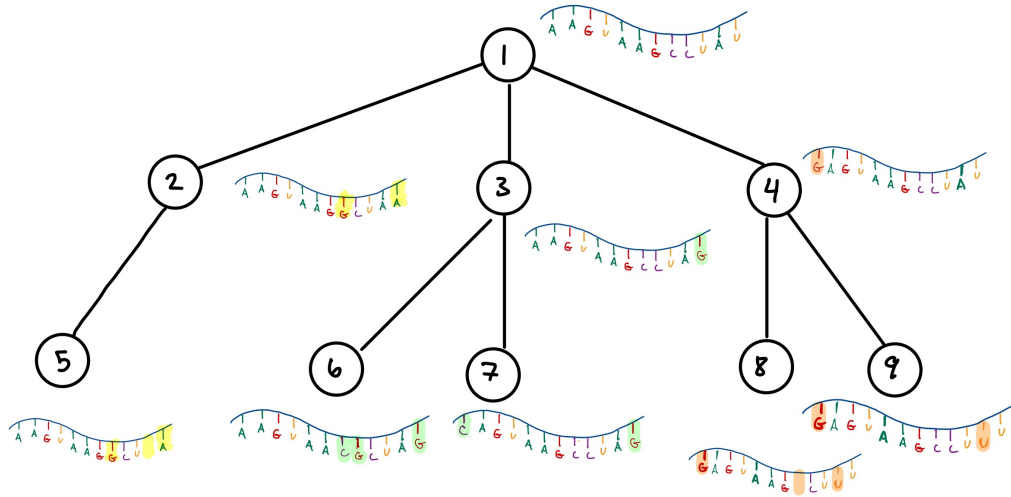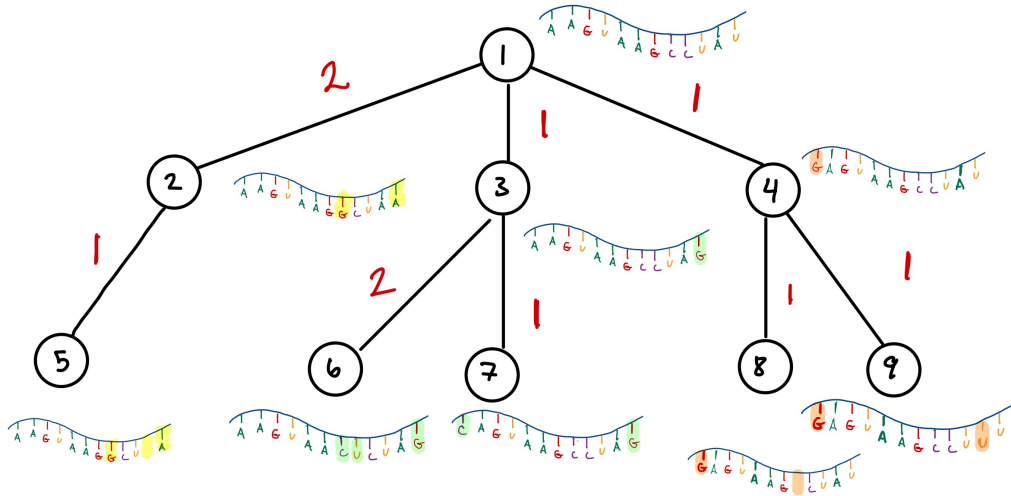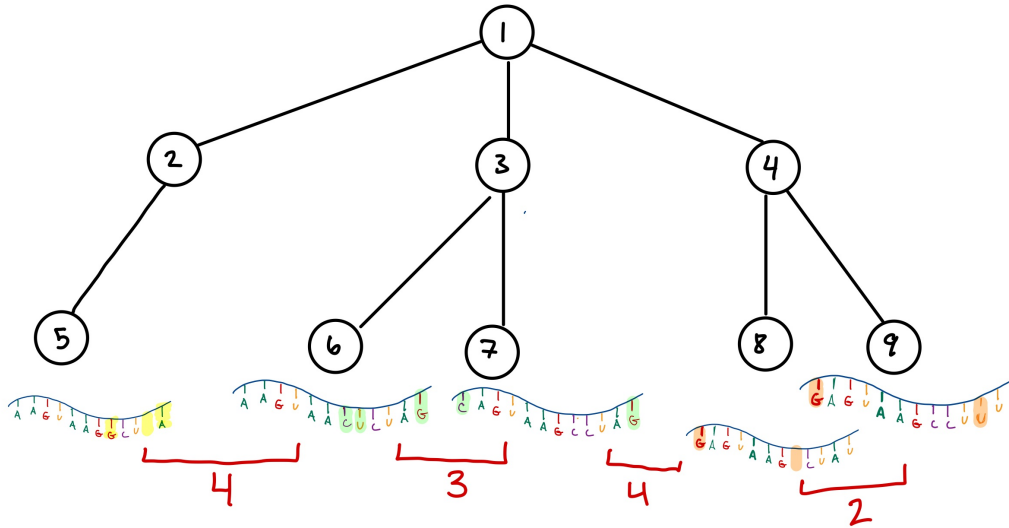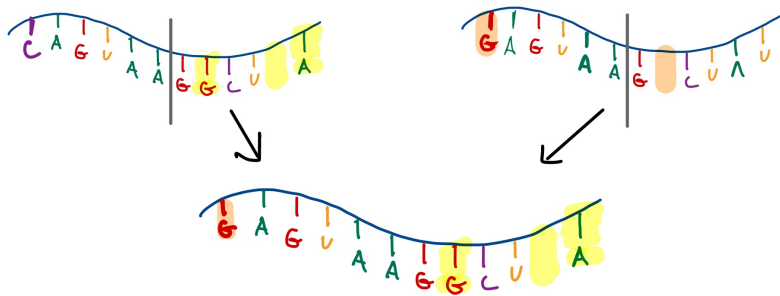$\{A, C, U, G\}$

# Point Mutations and Phylogenetic Tree

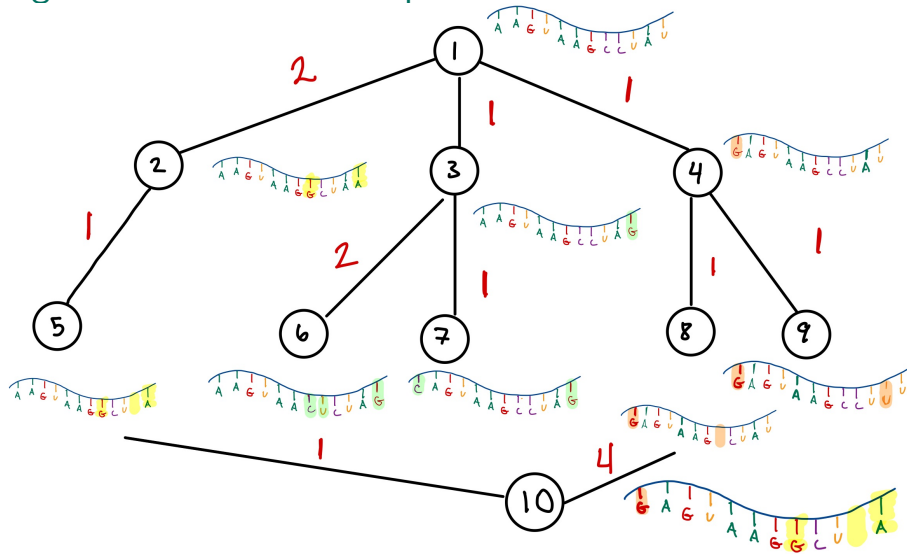# Hamming Distance Along Edges of Phylogenetic Tree - Point Mutations Add Up Through Vertical Evolution

# Basic Problem of Phylogenetics - Inferring Evolutionary History

# Horizontal Evolution - Recombination is an Exchange of Genetic Information
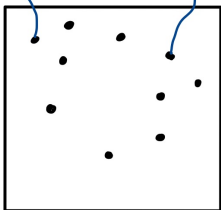
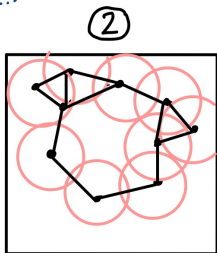# Phylogenetic Trees Fail to Capture Recombination



[1]Chan, J., et al., 2013. Topology of viral evolution. PNAS.

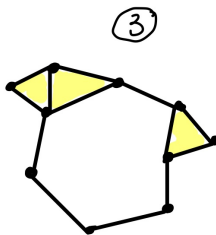# Topological Data Analysis as a Tool for Finding 1-Dimensional Holes



① Embed sequences in a metric space

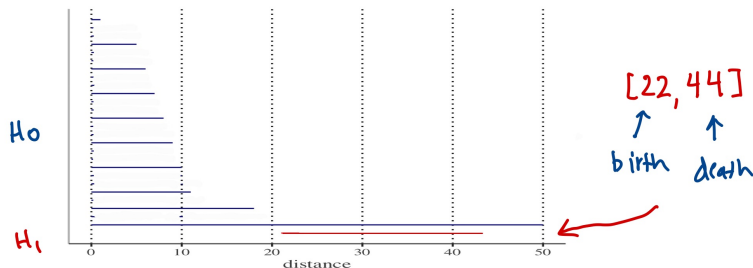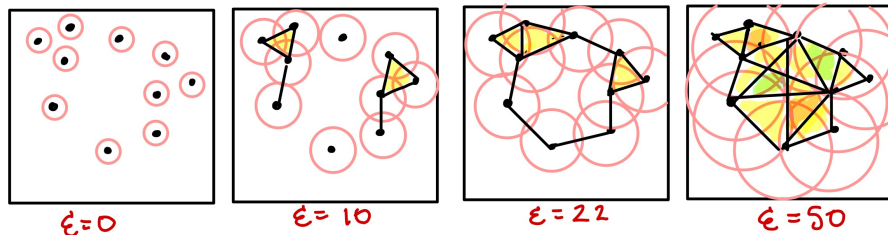② Draw balls around points, and connect points with intersecting balls
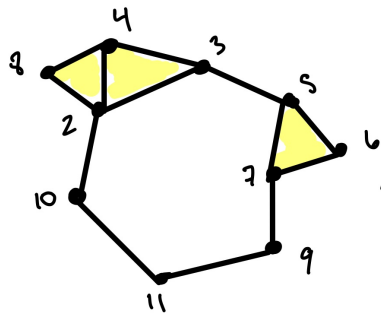
③ Construct simplicial complex

④ Compute homology groups to find connected components and 1-dimensional holes

# Persistent Homology Can Be Visualized as a Barcode



$\epsilon = 0$    $\epsilon = 10$    $\epsilon = 22$    $\epsilon = 50$

$H_0$

$H_1$

distance

[22, 44]
↑      ↑
birth  death

2

---

[2]Ghrist, R., 2007. Barcodes: The persistent topology of data. Bulletin of AMS.

$H_1$ generator:

Nodes and edges that bound a 1-dimensional hole at birth

← Possible generator:

[3,5]
[5,7]
[7,9]
[9,11]
[11,10]
[10,2]
[2,3]
[3,5]

What can topological data analysis tell us about the evolution of SARS-CoV-2?

# Collecting SARS-CoV-2 Sequences

[3] www.gisaid.org

# Data is Downloaded as Text File



FASTA
file

# Align Sequences Before Computing Distance

Mega X: Molecular Evolutionary Genetic Analysis

4

# TDA Results Were Inconclusive



Contra Costa.
100 sequences

San Francisco - 129
sequences

5

What do TDA barcode lengths and generators reveal about the underlying data?

# Simulation Method

- Start with randomly generated 1,000-length sequence

- Choose population size, number of generations, mutation rate, survivors per generation

- 1 homologous recombination event per simulation cut at 50%

6

18

# Example of Simulation

# Recombination Cycle and TDA Results

$H_1$ Birth: 21
Death: 23

Generator: [28:32], [24:32]
[8:24], [8:28]

No hole formed

# An Aside: How a Recombinant Can Be Genetically Closer to a Common Ancestor than Its Direct Ancestor

# Another Example of Missed Recombination

# Hypothesis to Explain Missed Homologous Recombinations

Homologous recombination is missed when the recombinant sequence has closer genetic distance to a common ancestor than to one or more of its direct ancestors.

# TDA Results of Simulation



Generator

Birth: 33
Death: 36

# Where Are the Missing Ancestors In the Generator?



Birth: 33

Death: 36

# Conjecture of Relationship Between Recombinant and Other Sequences in Generator



$$i \in \mathbb{Z}/n\mathbb{Z}$$

$$R_i = \frac{d(V_i, V_{i-1}) + d(V_i, V_{i+1})}{d(V_{i-1}, V_{i+1})}$$

R-score of $V_i$

$$\frac{R_i + R_{i-1} + R_{i+1}}{3}$$

C l l R s i l ly recombinant

# Further Examples of Generator Conjecture

# Cycling Back - Looking at the California Data



San Francisco - H1 Cycle - Birth: 62, Death: 85, Length: 23

San Francisco
Feb - April 2021

R
(likely)

▶ Develop Generator Conjecture - Run simulations with different types of recombinations, metrics, complexes.

▶ Analyze SARS-CoV-2 data with respect to H1 generators - Compare suspected recombinants point-by-point with other sequences in the generator.

## THANK YOU!

Dr. Javier Arsuaga, Dr. Mariel Vazquez, Sofia Jakovcevic, Nathan Solomon, Michael Keith, Emil Geisler, Georgina Gonzalez, Arsuaga Vazquez Lab, Greg Kuperberg and the UC Davis REU

# References

Allman, E., Rhodes, J., 2016. Lecture notes: the mathematics of phylogenetics, IAS/Park City Mathematics Institute and University of Alaska Fairbanks.

Chan, J., Carlsson, G. and Rabadan, R., 2013. Topology of viral evolution. Proceedings of the National Academy of Sciences, 110(46), pp.18566-18571.

Ghrist, R., 2007. Barcodes: The persistent topology of data. Bulletin of the American Mathematical Society, 45(01), pp.61-76.

Jordá, Teresa Díaz, 2020. 'Characterization of Horizontal Evolution of RNA Viruses Using Topological Data Analysis', Bachelor's Thesis, University of California, Davis and Universitat Politècnica de València.