

# Global Phylogeny by Local Inferences

*Final Project for Phylogeny Inference & Applications*

Freie Universität Berlin: Winter Semester 2022-2023  
Meghan Kane  
Derek Shao

# Outline

- Introduction
- Motivation
- Methods
- Materials
- Results
- Discussion
- Contributions

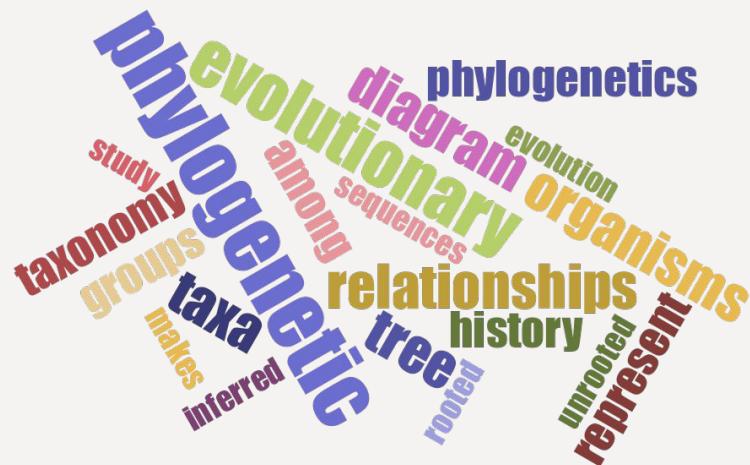
# Outline

- **Introduction**
- Motivation
- Methods
- Materials
- Results
- Discussion
- Contributions

# What is phylogeny?

The relationships between entities to reflect their evolutionary histories (De Bruyn et al., 2013)

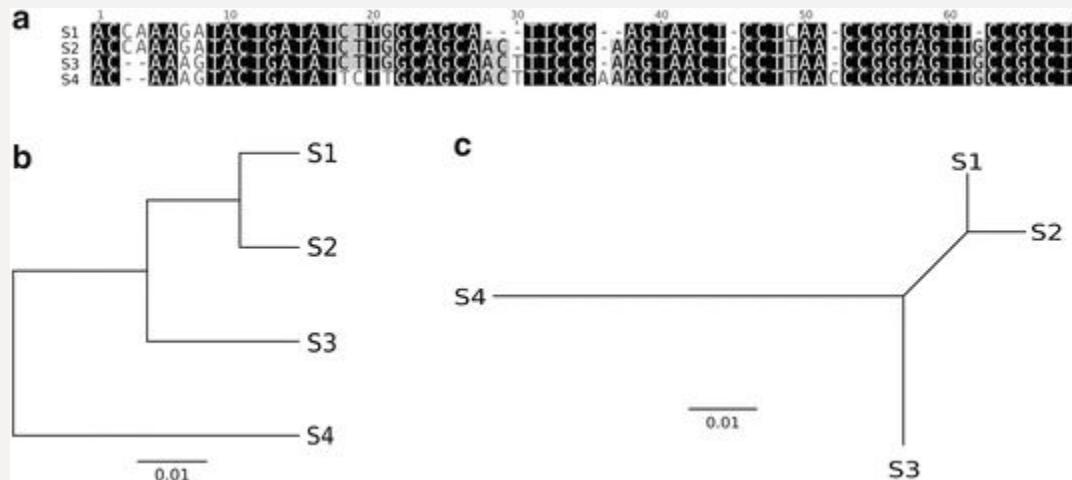
Key assumption: Measurable similarities between organisms are suggestive of their common evolutionary history.



Word cloud generated from <https://en.wikipedia.org/wiki/Phylogenetics>

# Phylogenetic tree

- Basic components: *vertex (node) & branch (edge)*
  - *Rooted* versus *Unrooted* trees
  - *Bifurcating* versus *Multifurcating* trees



# Comparison of methods in computational phylogenetics

	Attributes (+)	Notes (-)
Distance based (ape, ClustalOmega, EBI, phangorn, FastML)	<ul style="list-style-type: none"><li>Computational speed</li><li>Multiple distance metrics possible</li></ul>	<ul style="list-style-type: none"><li>Difficult to calculate distances for sequences with large divergence and alignment gaps</li><li>Less sophisticated than probabilistic methods</li></ul>
Maximum parsimony (PHYLIP, Rphylip, GCTree, phangorn, IgTree)	<ul style="list-style-type: none"><li>Intuitive algorithm</li><li>Clonal frequency incorporation (GCTree)</li><li>Polytomies and internal nodes (IgTree)</li></ul>	<ul style="list-style-type: none"><li>Ignores antibody specific properties (hotspots, transversions, transitions)</li><li>Long-branch attraction problem</li></ul>
Maximum likelihood (FastML, MEGA, IQ-TREE, dhami, IgPhyML)	<ul style="list-style-type: none"><li>Complex substitution models</li><li>Hotspot specific codon models (IgPhyML)</li></ul>	<ul style="list-style-type: none"><li>Computationally demanding</li><li>Sensitive to model misspecification</li></ul>
Bayesian (BEAST, Mr. Bayes, ImmuniTree)	<ul style="list-style-type: none"><li>Complex substitution models</li><li>Can produce rooted trees without explicit outgroup</li><li>Possible to incorporate biological knowledge with priors</li><li>Mutation rate returned in calendar time (BEAST)</li></ul>	<ul style="list-style-type: none"><li>Sensitive to model misspecification</li><li>Highest computational demands due to Markov chain Monte Carlo algorithm</li></ul>

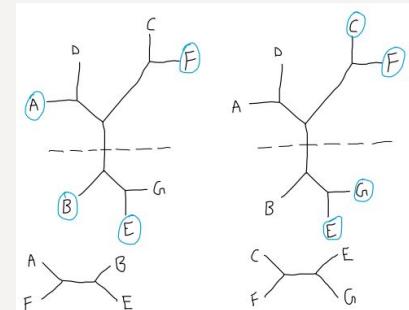
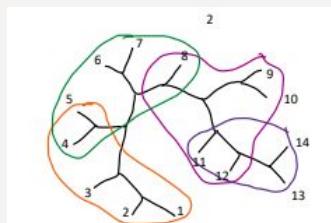
# Large-scale phylogenetic inference

Background:

- Exponentially increasing volume of DNA and protein sequence data
- Identifying the most parsimonious tree is NP-hard (Felsenstein, 2004)
- Finding maximum likelihood (ML) phylogenies is NP-hard (Roch, 2006; Chor & Tuller, 2006)

Common methods:

- Quartet-based method (Ranwez & Gascuel, 2001)
- Disc-covering method (Bayzid et al., 2014)

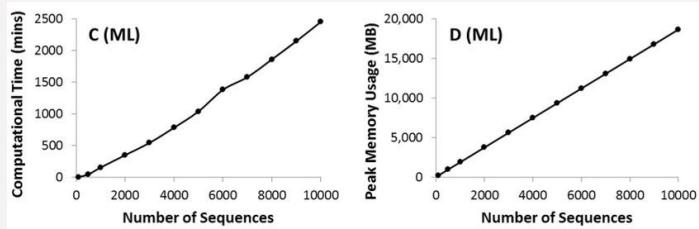


# Outline

- Introduction
- **Motivation**
- Methods
- Materials
- Results
- Discussion
- Contributions

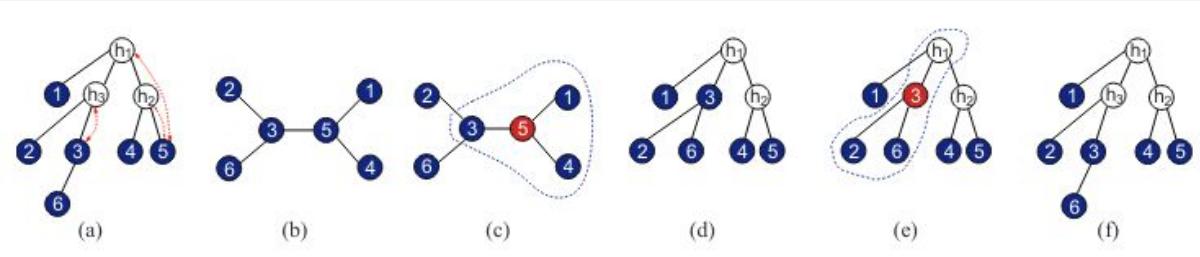
# Motivation

- Unacceptable computing time in MP and ML analyses when number of taxa increases



Kumar et al., 2016

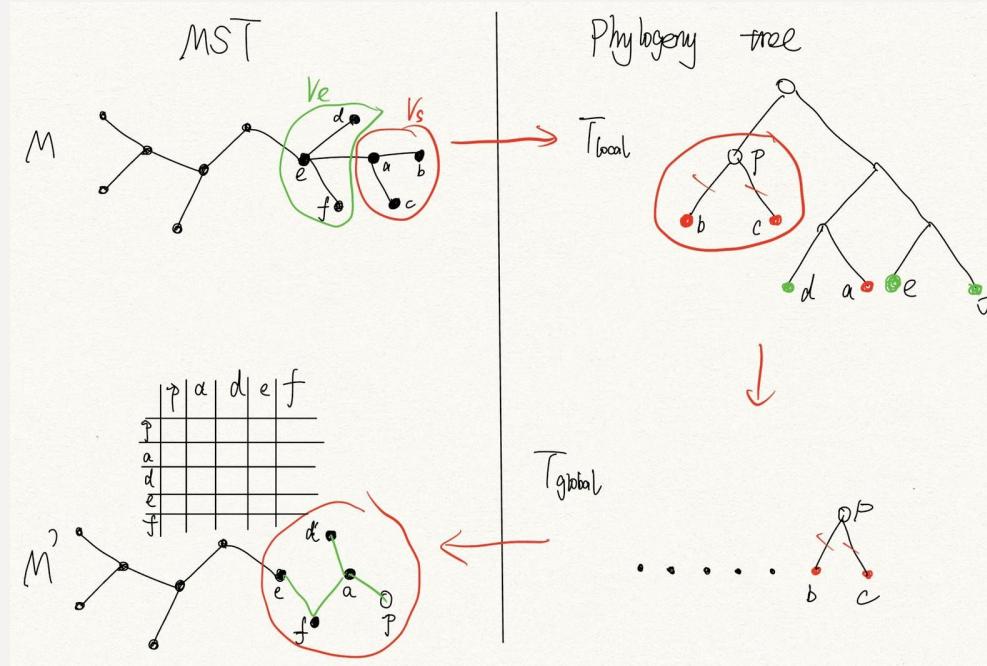
- Limitations of quartet-based and disc-covering methods
- Topological correspondence between unrooted phylogenetic tree and minimum spanning tree constructed using additive distances (Choi et al., 2010; Kalaghpati & Lengauer, 2017)



Choi et al., 2010

# Objective

An efficient and scalable algorithm for phylogenetic inference based on minimum spanning tree!



# Outline

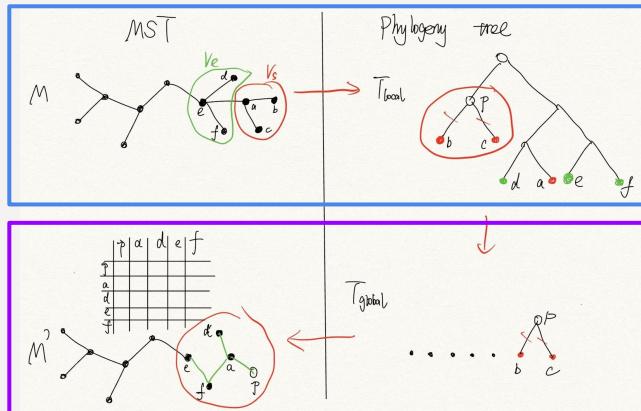
- Introduction
- Motivation
- **Methods**
- Materials
- Results
- Discussion
- Contributions

## Workflow

- Initiate global phylogeny tree
- Construct minimum spanning tree

While global tree is not connected:

- Construct a MST-induced MP local tree
- Update global tree & MST



**Input:** Aligned sequences  $X_L$  for each taxon in the taxon set  $L$ , and a subtree size threshold  $s$

Initialize a global unrooted phylogenetic tree  $T = (V_T, E_T)$  as  $V_T \leftarrow L$  and  $E_T \leftarrow \{\}$ ;

Compute a minimum spanning tree  $M = (V_M, E_M)$  using Hamming distances of each sequence pair in  $X_L$ ;

**while**  $|E_T| < |V_T| - 1$  **do**

**if**  $M$  contains a subtree  $\tau_s = (V_s, E_s)$  with more than  $s$  vertices, and  $|V_M| - |V_s| > s$  **then**

    Select a smallest subtree  $\tau_s$  of  $M$  containing more than  $s$  vertices;

    Select  $s$  vertices  $V_e$  of  $M$  that are visited using a BFS starting at the root of  $\tau_s$  such that  $V_e \cap V_s = \{\}$ ;

    Compute a ML phylogenetic tree  $t_\rho$  over vertices  $V_s \cup V_e$  by performing tree search using SEM-GM;

    Construct  $t$  by suppressing the root of  $t_\rho$ , and select the largest non-singleton subforest  $f_s$  of  $t$  that is induced by  $V_s$ ;

    Add edges in  $f_s$  to the global phylogenetic tree  $T$ ;

    Update  $M$  by removing leaves of subtrees, and adding the root of each subtree;

**else**

    Compute a ML phylogenetic tree  $t_\rho$  over all vertices in  $V_M$  using SEM-GM;

    Construct an unrooted phylogenetic tree  $t$  by suppressing the root in  $t_\rho$ , and add all edges in  $t$  to  $T$ ;

**end**

**end**

**Output:** A global unrooted phylogenetic tree  $T$

# Approach

1. Construct minimum-spanning tree (MST)
2. Define global phylogeny tree
3. Construct MST-induced local subtrees
4. Update MST

# Initiate global phylogeny tree

- Apply custom Tree class developed during semester
- Read sequence data from influenza\_98.fasta file
- **Assign all observed sequences to leaf nodes by setting out\_degree = 0**

So, the starting global tree has 98 vertices and 0 edges

# Construct MST

## I. Compute pairwise Jukes-Cantor distances between sequences

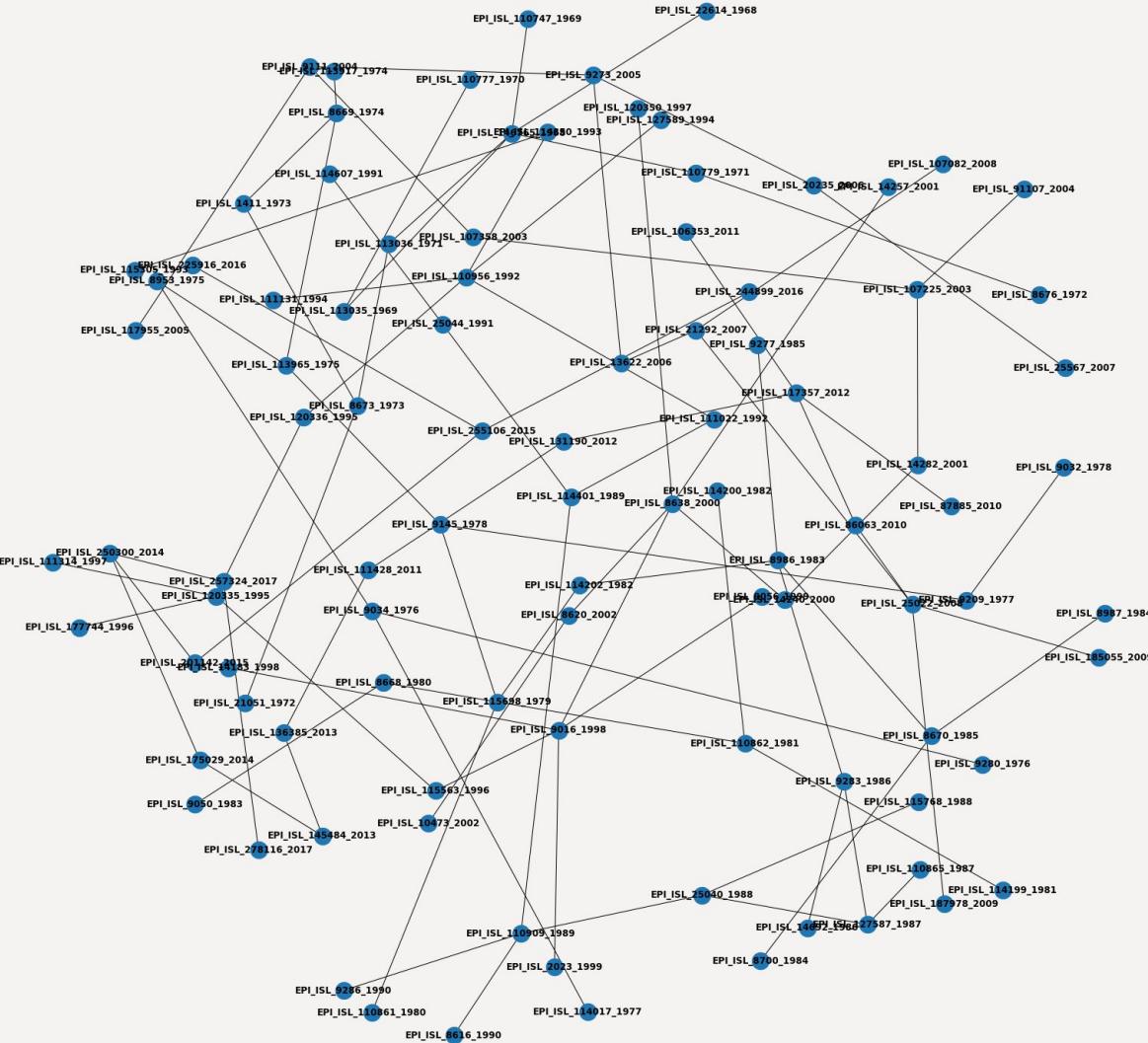
Estimate  $d = d_a + d_b$  given sequences  $X_a$  &  $X_b$  with Hamming distance  $\hat{p}$

$$\hat{d} = -\frac{3}{4} \ln \left( 1 - \frac{4\hat{p}}{3} \right)$$

- II. Create fully connected graph with distance matrix
- III. Compute MST in the connected graph
  - o networkx (Hagberg et al., 2008)

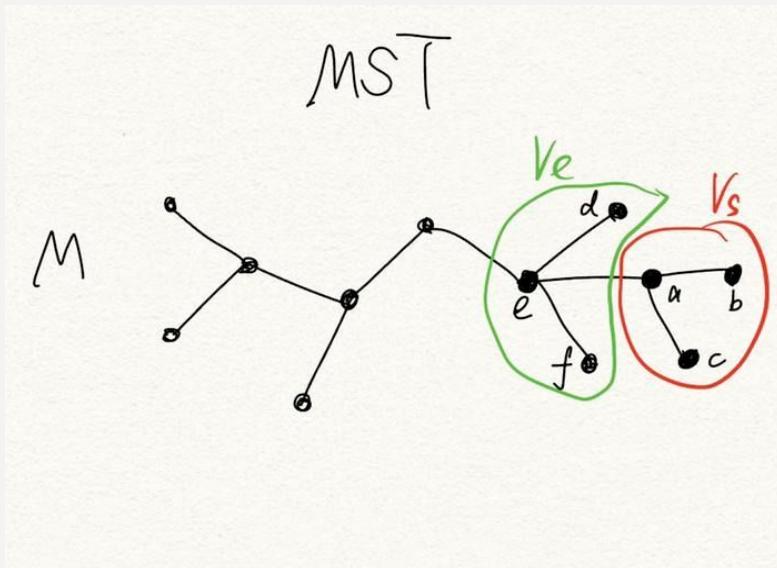
## Construct MST

influenza\_98.fasta



# Construct MST-induced MP local trees

- I. Convert networkx object to our custom Tree class
- II. Select  $V_s$  from MST s.t. a subtree is induced
  - *Threshold*: minimum number of vertices in  $V_s$  (subtree)



$$V_s = \{a, b, c\} \text{ at } Threshold = 3$$

## II. Select $V_s$ from MST s.t. a subtree is induced

**Algorithm 1** Select  $V_s$  from MST such that it induces a subtree

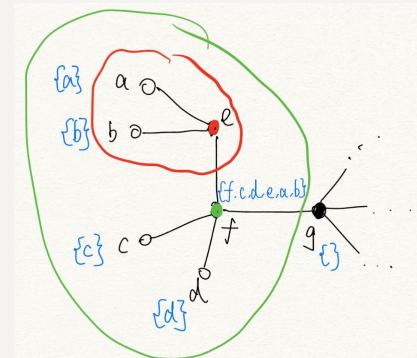
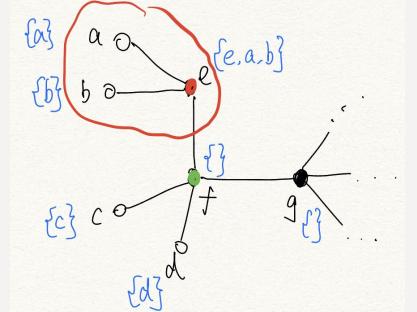
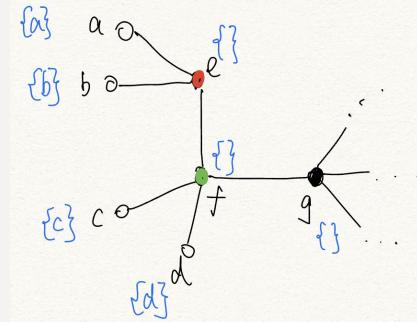
**Input:** MST nodes,  $Threshold$

**Output:** Subset  $V_s$  of MST nodes

```

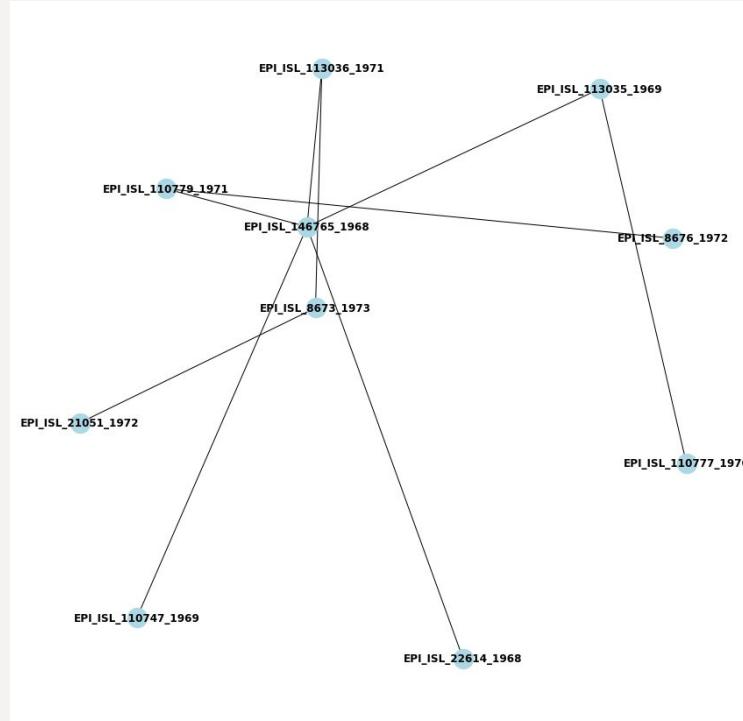
1: vertices_to_visit ← Assign all leaf nodes of MST
2:  $v \leftarrow$  Pop one random element from vertices_to_visit
3: while vertices_to_visit is not empty do
4:   if  $v$  is a leaf node then
5:      $v.\text{subtree} \leftarrow$  Assign  $v$  itself
6:     vertices_to_visit ← Append  $v.\text{neighbor}$  to the tail
7:   else if  $v$  is an internal node then
8:     Check the number of neighbors with empty subtree
9:     if Only one of its neighbors has an empty(unassigned) subtree then
10:        $v.\text{subtree} \leftarrow$  Assign the union of its neighbors' subtrees
11:     else
12:       vertices_to_visit ← Queue  $v$  for revisiting after updating its neighbors
13:     end if
14:   end if
15:   if Size of  $v.\text{subtree}$  is larger than  $Threshold$  then
16:     Return  $V_s = v.\text{subtree}$ 
17:   end if
18:    $v \leftarrow$  Pop the first element from vertices_to_visit

```



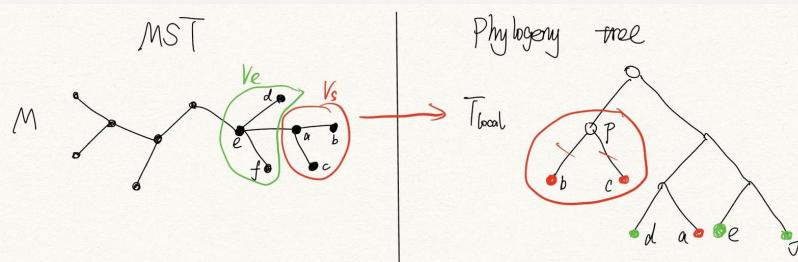
## II. Select $V_s$ from MST s.t. a subtree is induced

Example: Influenza\_98.fasta, threshold = 10



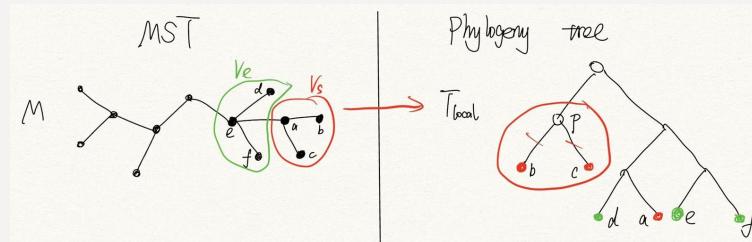
# Construct MST-induced MP local trees

- I. Convert networkx object to our custom Tree class
- II. Select  $V_s$  from MST s.t. a subtree is induced
  - Threshold: minimum number of vertices in  $V_s$  (subtree)
- III. Select  $V_e$  via BFS from root of  $V_s$ -induced subtree
- IV. Create one random staring tree topology from  $V_s \cup V_e$ 
  - ete3.populate - generates a random tree topology (Huerta-Cepas et al., 2016)
- V. Find local tree with maximum parsimony score
  - Bio.Phylo.TreeConstruction.ParsimonyTreeConstructor
- VI. Find minimum subtree unit induced by only  $V_s$  in local tree
  - Select  $V_s$  vertices with same parent (e.g., vertex  $b$  &  $c$ , and local parent vertex  $p$ )



# Construct MST-induced local subtrees (M)

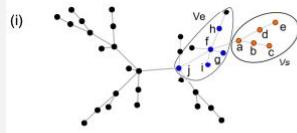
- I. Convert networkx object to our custom Tree class
- II. Select  $V_s$  from MST s.t. a subtree is induced
- III. Select  $V_e$  via BFS from root of induced subtree  $V_s$
- IV. Create random initial unrooted tree from  $V_s \cup V_e$ 
  - o ete3.populate - generates a random tree topology (Huerta-Cepas et al., 2016)
- V. Find local tree with maximum parsimony score
  - o Bio.Phylo.TreeConstruction.ParsimonyTreeConstructor
- VI. Find minimum subtree unit induced by only  $V_s$  in local tree
  - Select  $V_s$  vertices with same parent (e.g., vertices  $b$  &  $c$  and local parent  $p$ )



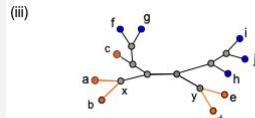
# Construct MST-induced local subtrees (M)

**TODO:** create & insert graph similar to Prabhav's (Fig 2) adjusted to our project

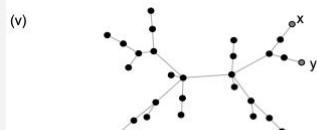
<https://www.biorxiv.org/content/10.1101/2020.06.30.180315v1.full.pdf>



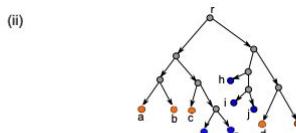
Select a vertex set  $V_s$  comprising  $s$  vertices such  $V_s$  induces a subtree in the MST. Start a breadth-first-search at the root of the subtree that is induced by  $V_s$  (vertex labeled  $a$ ), and stop if  $s$  vertices have been visited such that  $V_e$  and  $V_s$  have no vertices in common.



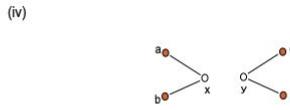
Suppress the root, and select all non-singleton subtrees that are induced by the vertices  $V_s$ . Edges of selected subtrees are highlighted in orange. Vertices  $x$  and  $y$  are the roots of the selected subtrees.



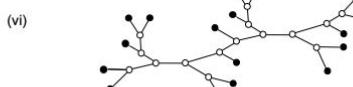
Update the MST by removing vertices that are leaves in  $F$ , and adding the roots of selected subtrees ( $x$  and  $y$ )



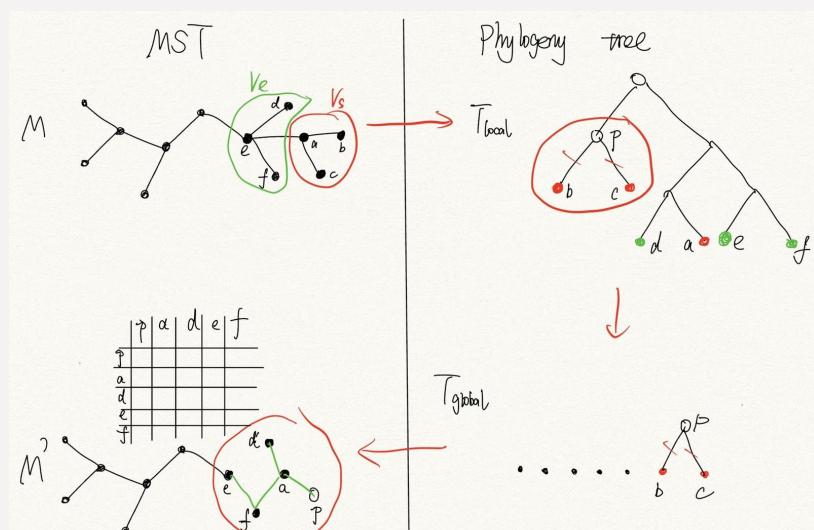
Fit a general Markov (GM) model on a rooted phylogenetic tree using SEM-GM.  $r$  indicates the root. Grey vertices represent the maximum a posteriori (MAP) estimate of ancestral sequences.



Update the global phylogenetic tree  $T$  by adding undirected edges in selected subtrees.



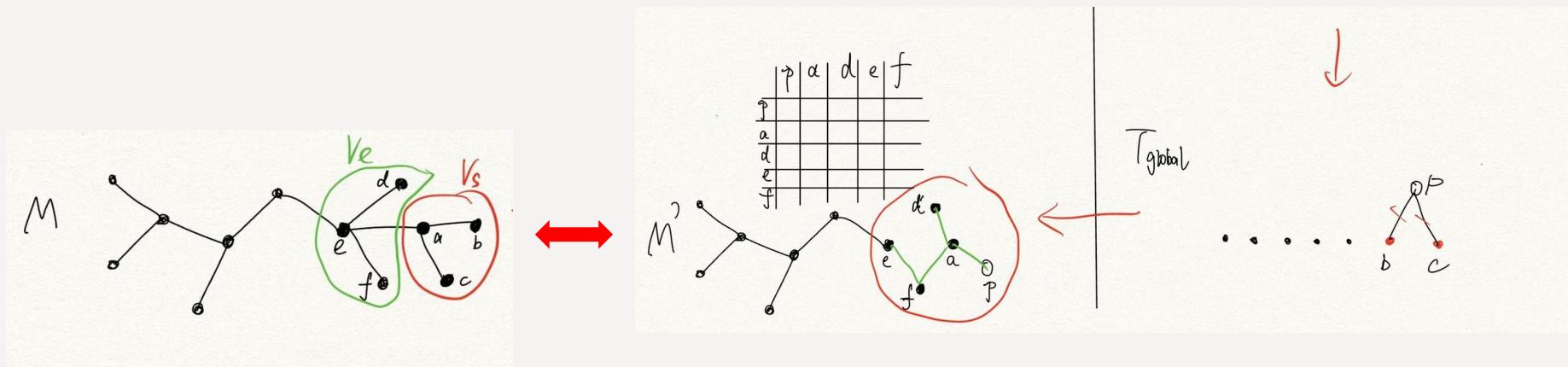
Iterate over steps (v) through (v), and stop if  $T$  is a connected graph. Root  $T$  at a vertex by fitting a GM model via SEM-GM such that the undirected topology is constrained to be  $T$ .



# Update global tree & MST

## I. Update global phylogeny tree

- Add new local parent vertex ( $p$ ) & corresponding edges to global tree
- II. Delete selected  $V_s$  vertices & add new local parent vertices in MST
- III. Update sequence distance matrix with  $V_s \cup V_e$
- IV. Update the subgraph of  $V_s \cup V_e$  in MST

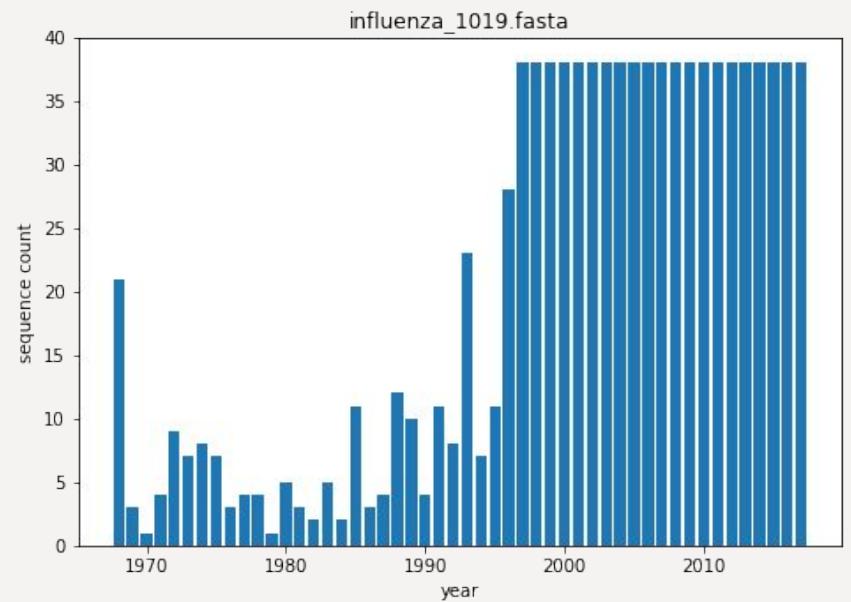
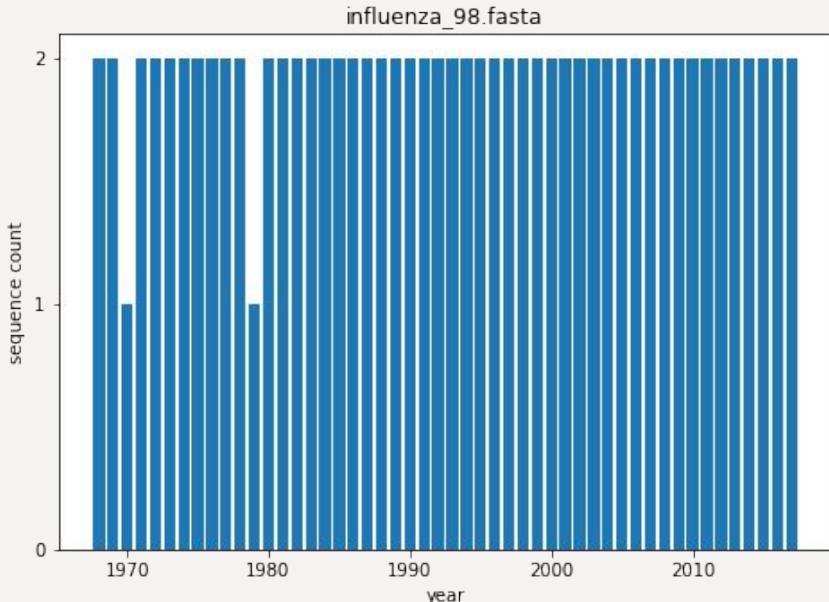


# Outline

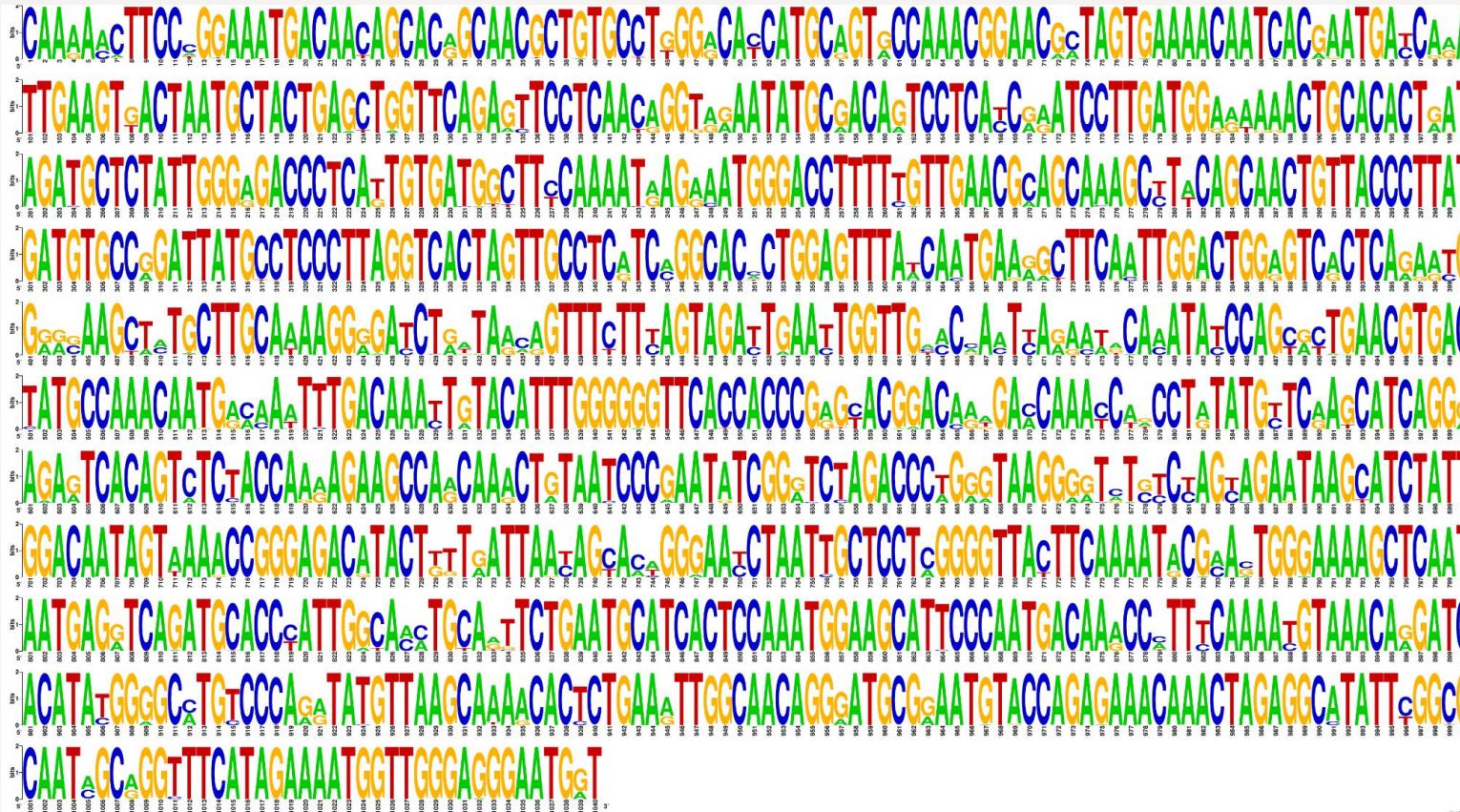
- Introduction
- Motivation
- Methods
- **Materials**
- Results
- Discussion
- Contributions

# Materials

- **2 fasta files:** influenza\_98.fasta & influenza\_1019.fasta
- **Sequence length:** 1040 nucleotides
- **Time period:** 1968 - 2017, distribution of data different between datasets



# Nucleotide distribution in Influenza\_98.fasta

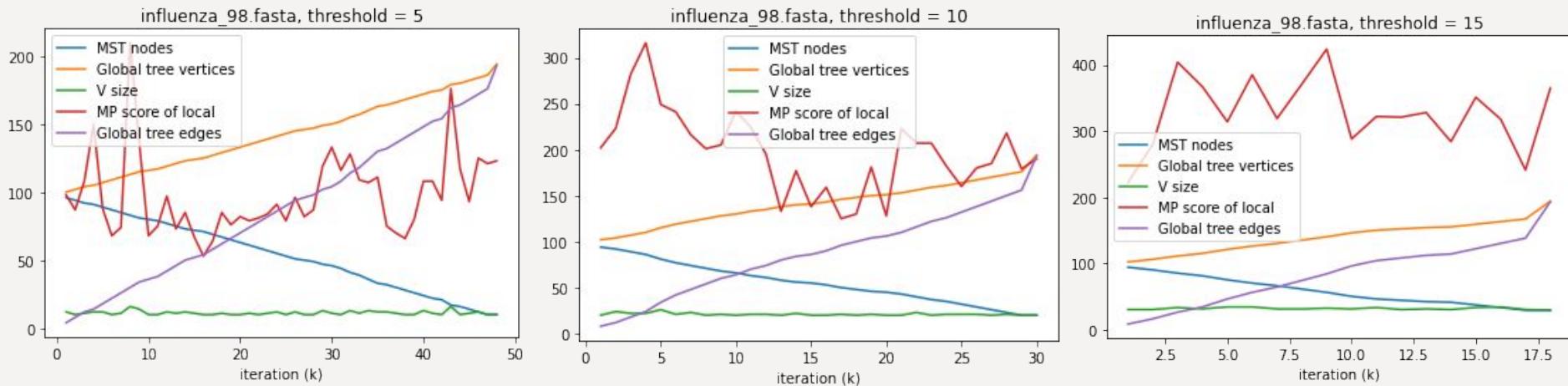


Crooks et al., 2004

# Outline

- Introduction
- Motivation
- Methods
- Materials
- **Results**
- Discussion
- Contributions

# Method Performance: influenza\_98



**Runtime:** larger tree space to search when threshold is increased

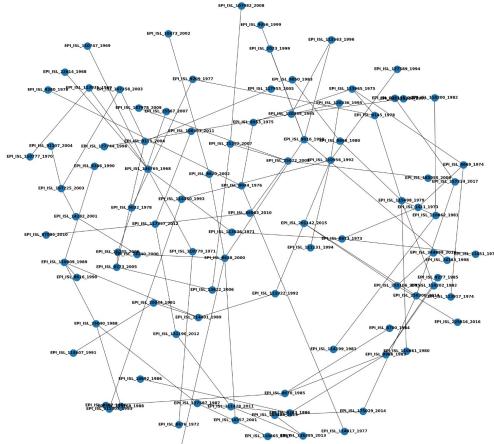
**96 s**  
(1 min 36 s)

**637 s**  
(10 min 37 s)

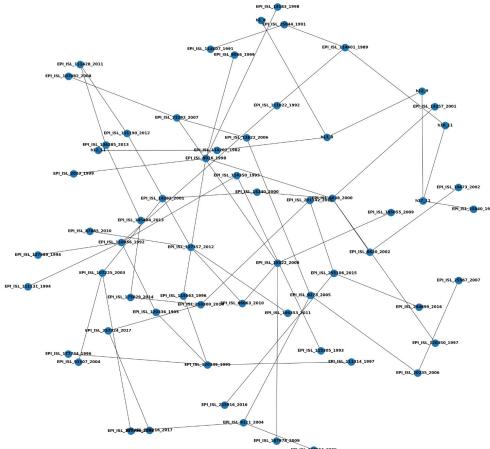
**1571 s**  
(26 min 11 s)

# Method Performance: MST plots

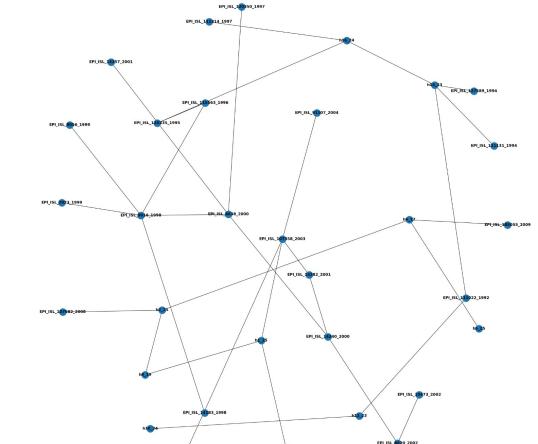
*start*



*checkpoint #1*



*checkpoint #2*

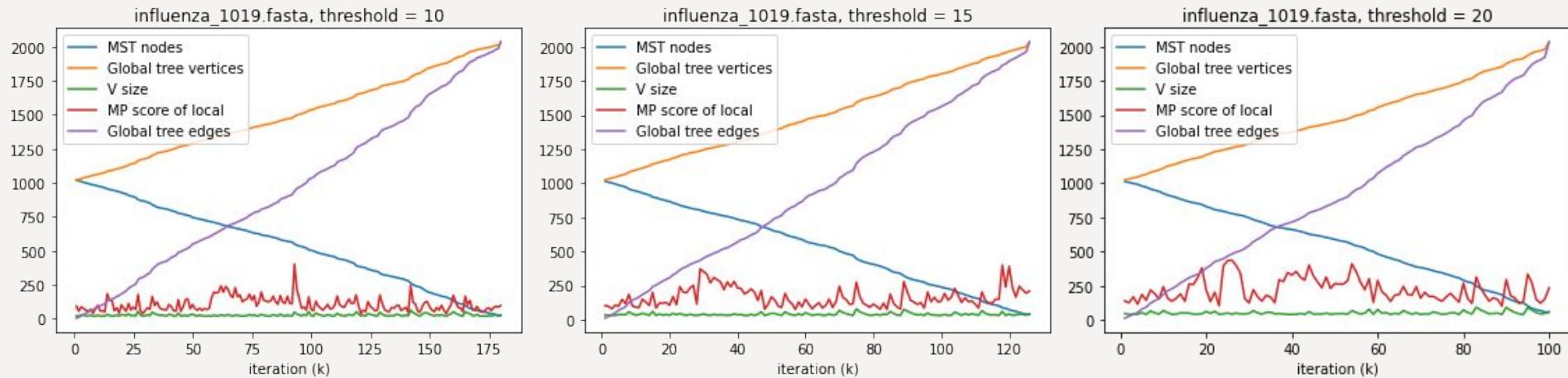


# nodes: 98

63

32

# Method Performance: influenza\_1019



Runtime:

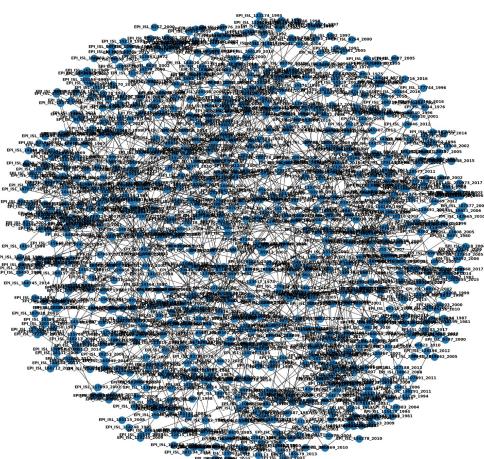
**3048 s**  
(50 min 48 s)

**5212 s**  
(86 min 52 s)

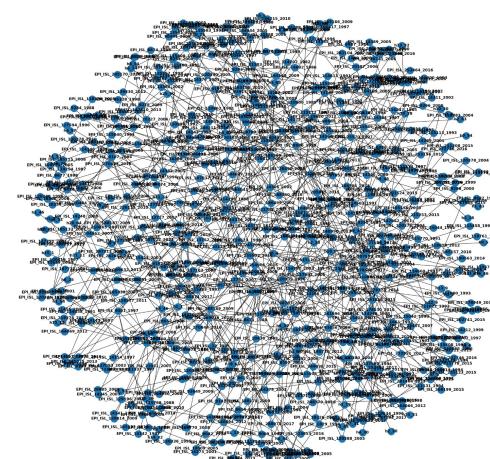
**11958 s**  
(199 min 18 s)

## Method Performance: MST plots

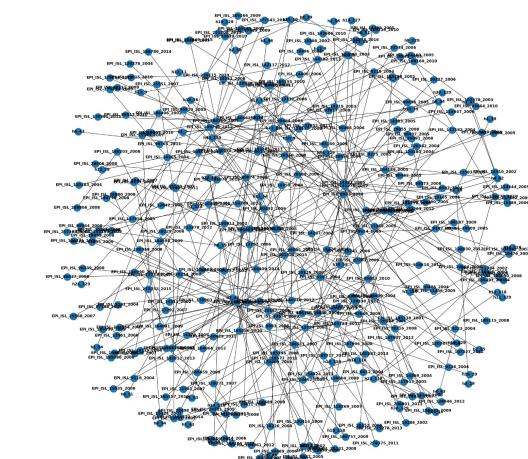
*start*



*checkpoint #1*



*checkpoint #2*



# nodes: 1019

676

335

# Runtime

<b>Dataset</b>	<b>Threshold</b>	<b>Runtime (our method)</b>	<b>Runtime (RAxML benchmark)</b>
influenza_98	5	96 s (1 min 36 s)	66 s (1 min 6 s)
	10	637 s (6 min 37 s)	
	15	1571 s (26 min 11 s)	
influenza_1019	10	3048 s (50 min 48 s)	6152 s (102 min 32 s)
	15	5212 s (86 min 52 s)	
	20	11958 s (199 min 18 s)	

# OLD (skip)

## Results vs Benchmark: Robinson-Foulds distance (RF)

$$\text{normalized RF} = \frac{|S_1 \cup S_2| - |S_1 \cap S_2|}{|S_1 \cup S_2|}$$

Dataset	Threshold	$ S_1 \cup S_2 $	$ S_1 \cap S_2 $	RF distance	RF max	Normalized RF distance
influenza_98	5	246	142	104	190	0.423
	10	249	139	110	190	0.442
	15	259	129	130	190	0.502
influenza_1019	10	2875	1197	1678	2032	0.584
	15	2881	1991	1690	2032	0.587
	20	2891	1181	1710	2032	0.591

# Results vs Benchmark: Robinson-Foulds distance (RF)

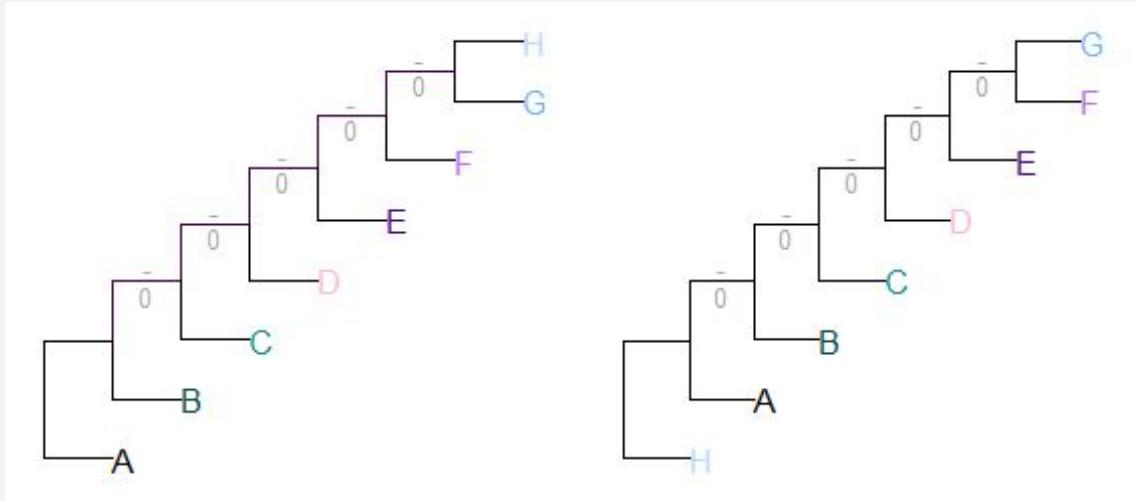
$$normalized\ RF = \frac{|S_1 \cup S_2| - |S_1 \cap S_2|}{max\ RF}$$

<b>Dataset</b>	<b>Threshold</b>	$ S_1 \cup S_2 $	$ S_1 \cap S_2 $	<b>RF distance</b>	<b>RF max</b>	<b>Normalized RF distance</b>
<b>influenza_98</b>	5	246	142	104	190	0.548
	10	249	139	110	190	0.579
	15	259	129	130	190	0.684
<b>influenza_1019</b>	10	2875	1197	1678	2032	0.826
	15	2881	1991	1690	2032	0.832
	20	2891	1181	1710	2032	0.842

# Drawbacks of using RF distance

- Maximum RF distance occurs after moving a single leaf (H)

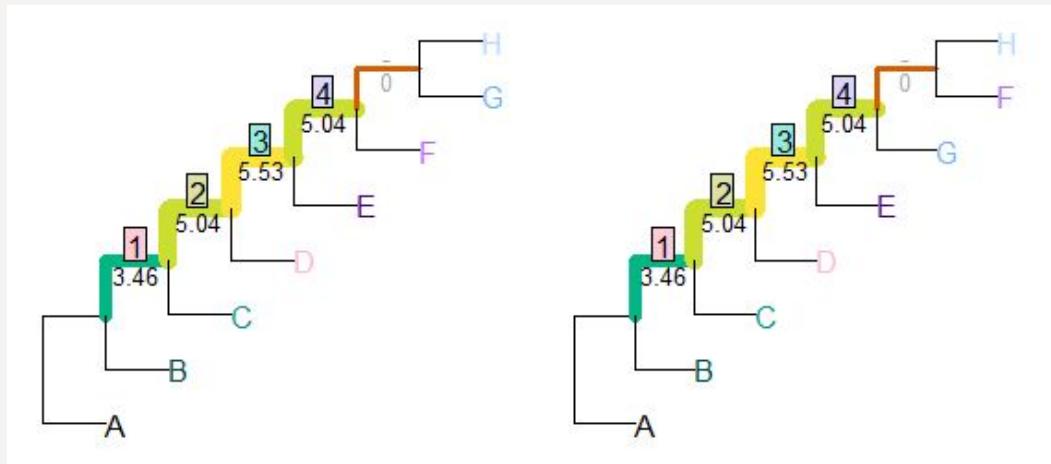
$$\text{normalized RF} = \frac{|S_1 \cup S_2| - |S_1 \cap S_2|}{\max RF} = \frac{19 - 9}{10} = 1$$



# Quantifying Information in Split

- Inspired by Claude Shannon's information theory
- Information contained in some splits is more instructive than others

$$\text{information content ( bits)} = - \log_2(\text{probability})$$



Normalized RF = 1  
Normalized info RF = 0.153

Information corrected RF distance (info RF) =  $22.54 - (3.46 + 5.04 + 5.53 + 5.04) = 3.46$  bits  
Normalized info RF =  $3.46/22.54 = 0.153$

# OLD (skip)

## Beyond Classic RF distance: generalized RF distance

$$\text{information content ( bits)} = - \log_2(\text{probability})$$

Dataset	Threshold	RF	Info RF	Shared Phylogenetic Info	Different Phylogenetic Info	Mutual Clustering Info
influenza_98	5	0.547	0.554	0.723	0.277	0.727
	10	0.579	0.622	0.700	0.304	0.701
	15	0.684	0.764	0.614	0.386	0.634

\*TreeDist tool wouldn't return results for influenza\_1019 dataset

# Metrics Beyond Classic RF distance

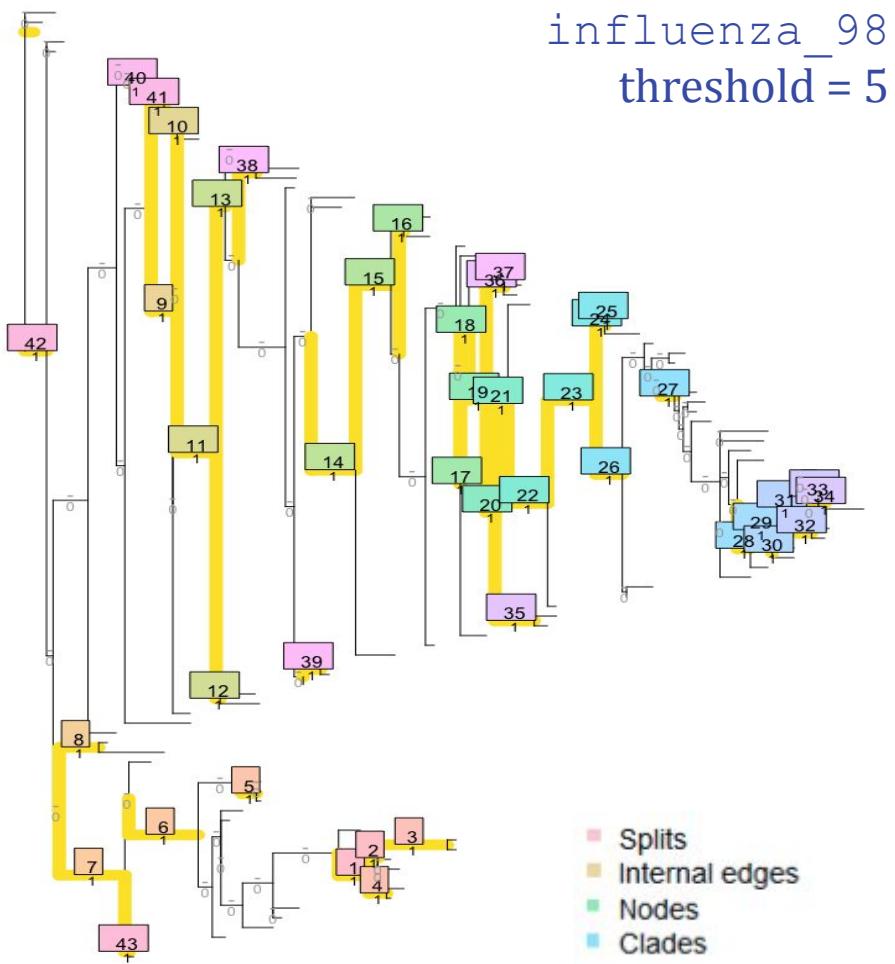
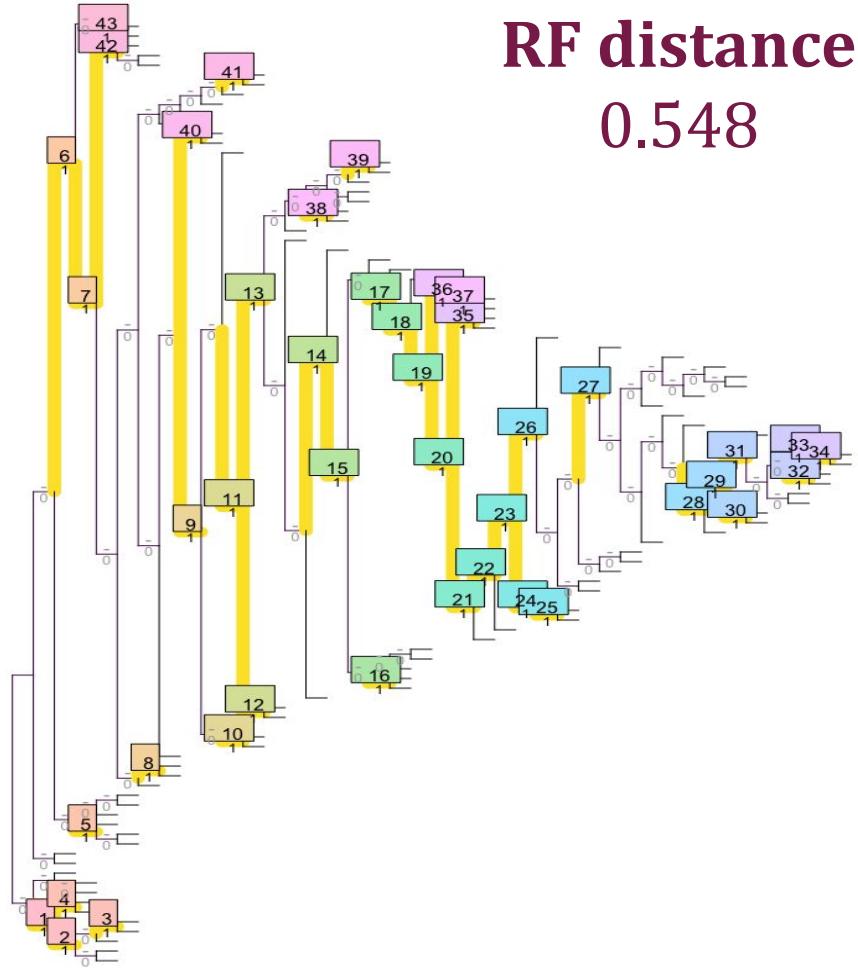
$$\text{information content ( bits)} = - \log_2(\text{probability})$$

Dataset	Threshold	RF	Info RF	Shared Phylogenetic Info	Different Phylogenetic Info
influenza_98	5	0.547	0.554	0.723	0.277
	10	0.579	0.622	0.700	0.304
	15	0.684	0.764	0.614	0.386

\*TreeDist tool wouldn't return results for influenza\_1019 dataset

# RF distance

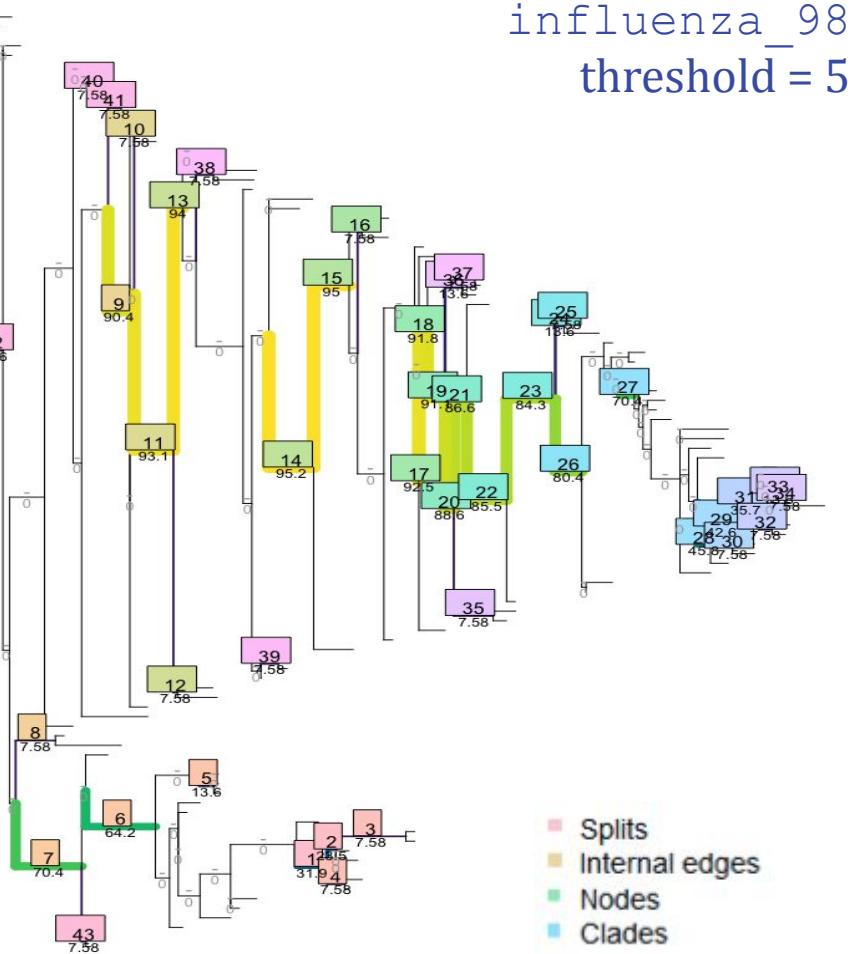
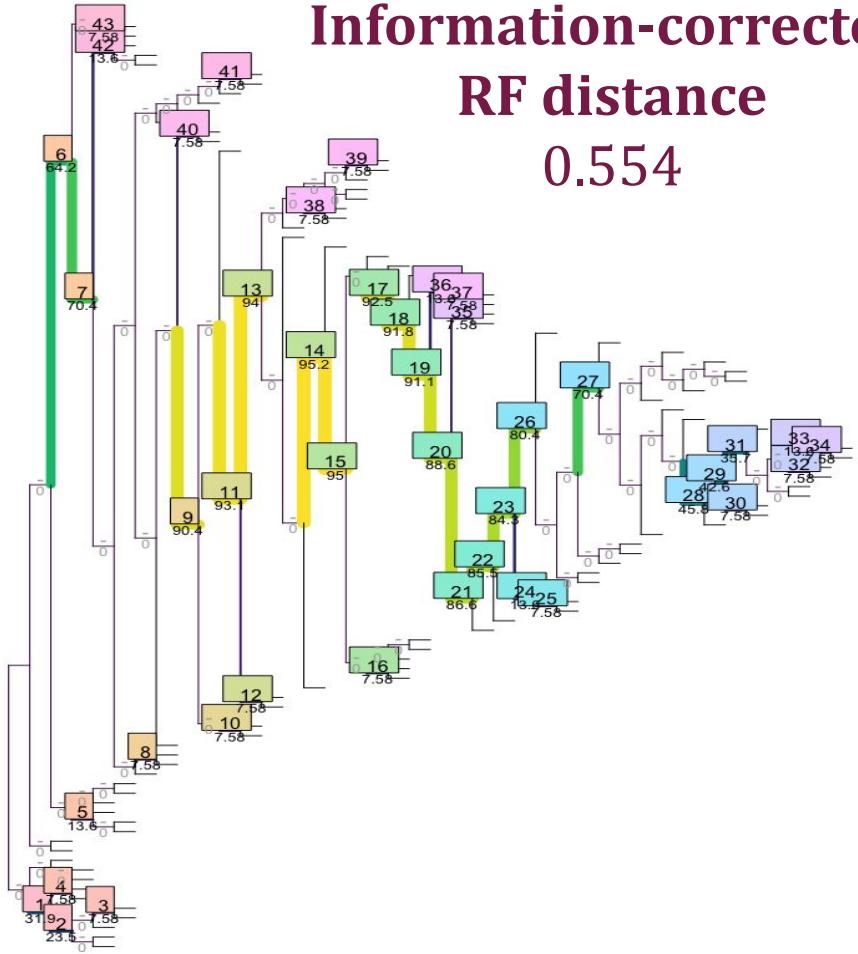
0.548



- Splits
- Internal edges
- Nodes
- Clades

# Information-corrected RF distance

0.554

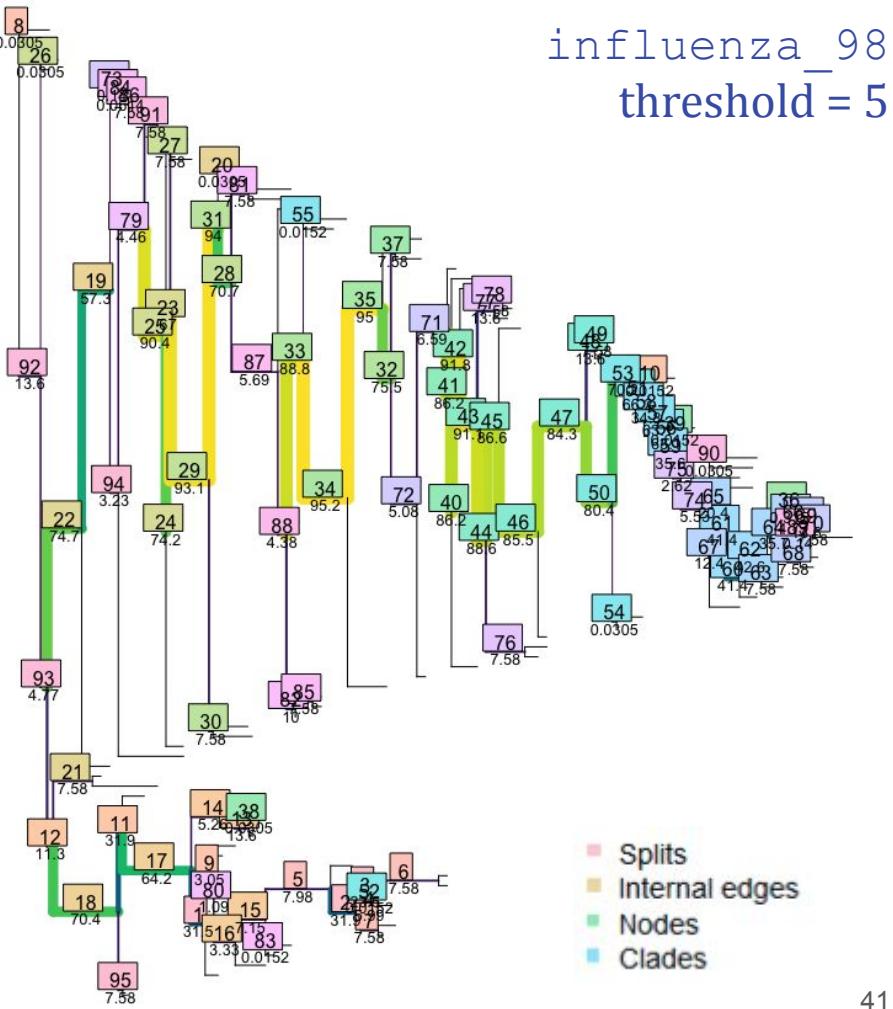
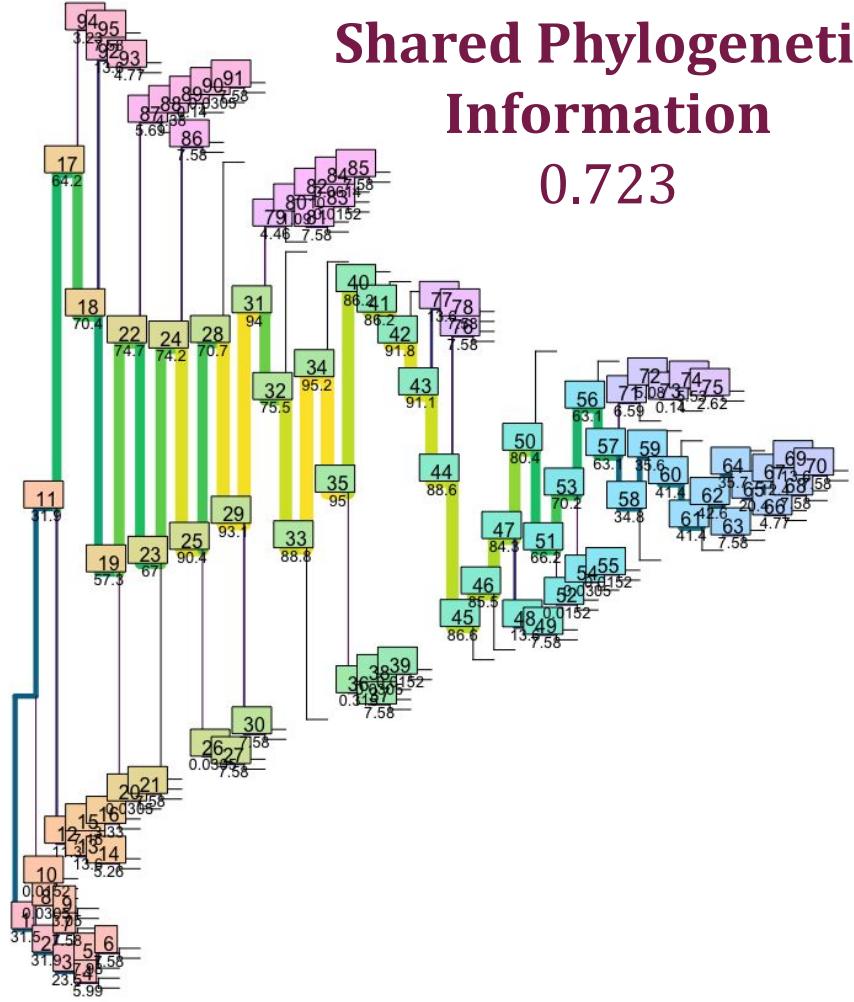


influenza\_98  
threshold = 5

- Splits
- Internal edges
- Nodes
- Clades

# Shared Phylogenetic Information

0.723

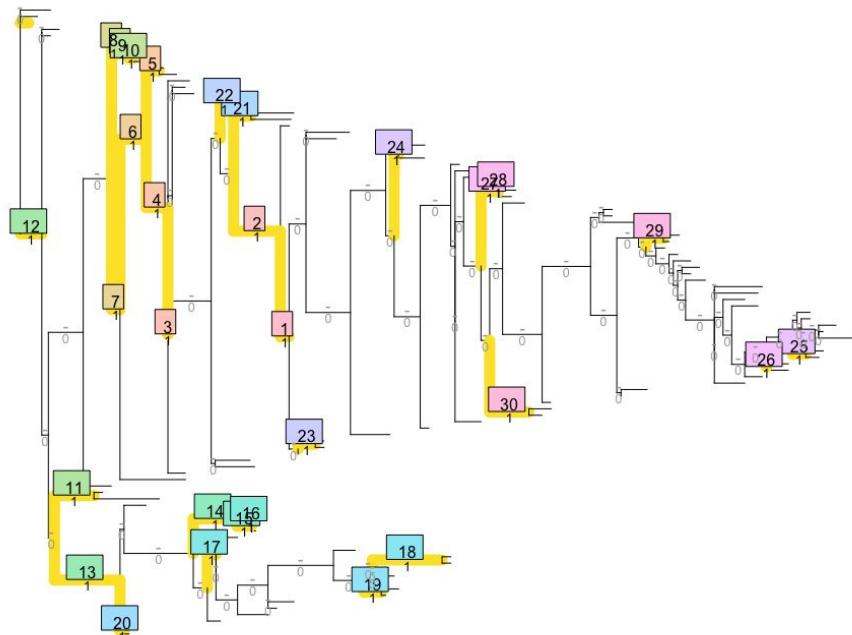
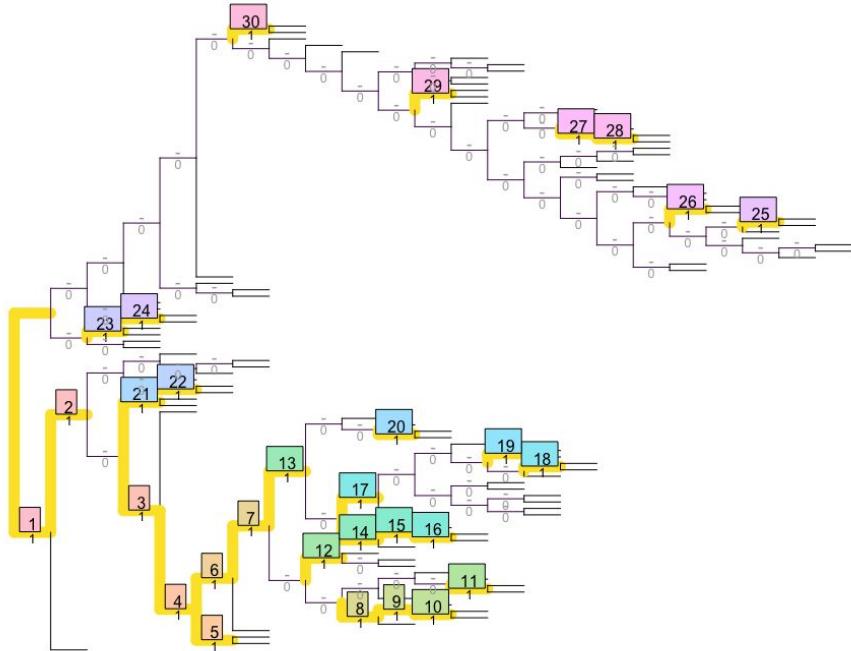


influenza\_98  
threshold = 5

- Splits
- Internal edges
- Nodes
- Clades

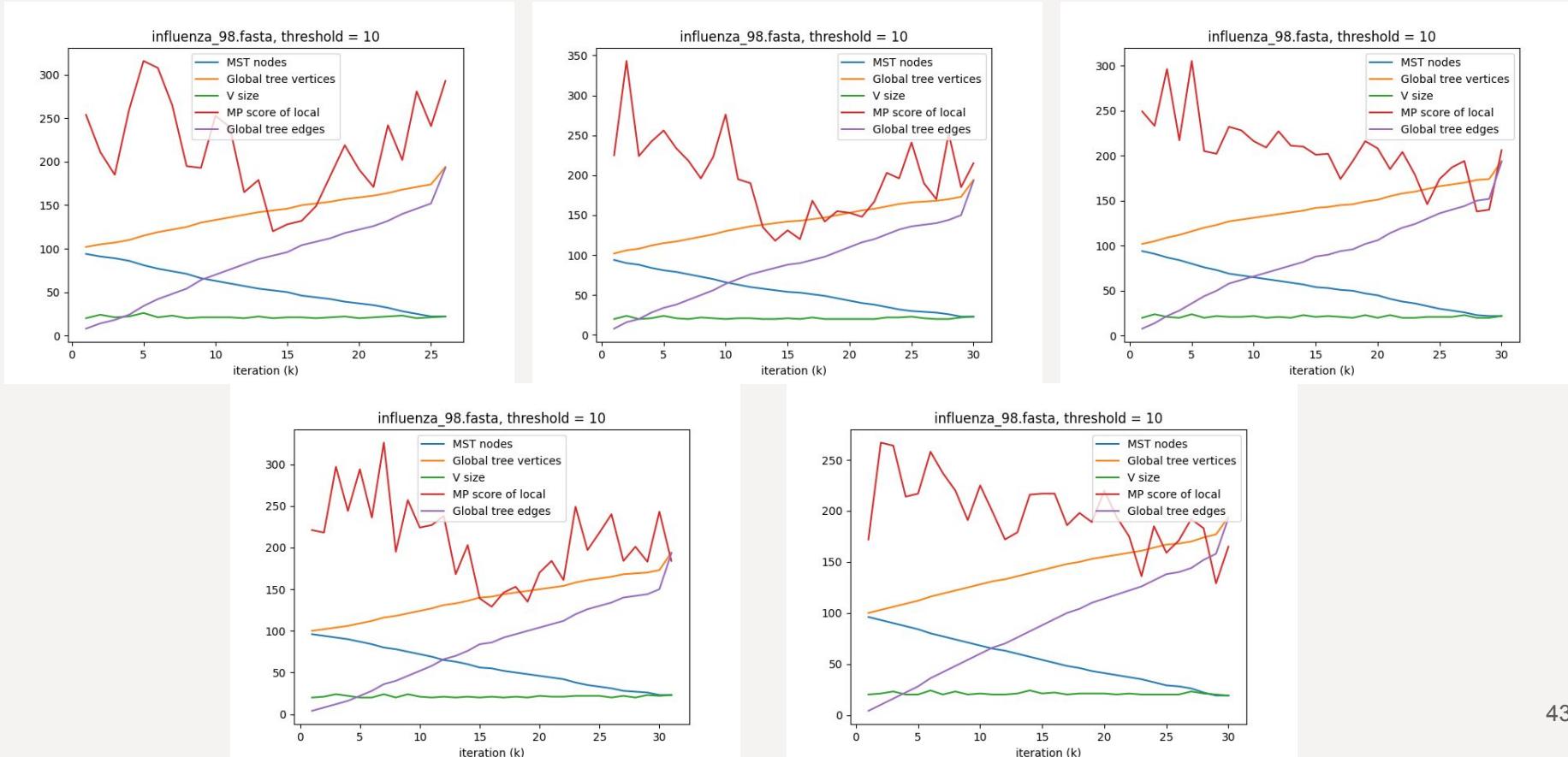
■ Splits  
■ Internal edges  
■ Nodes  
■ Clades

# RF distance: influenza\_98, threshold = 15



# Method Performance: influenza\_98, threshold = 10 (5 runs)

- Run 5 times with 5 random seeds for more robust statistics



Runtime: influenza\_98 , threshold = 10 (5 runs)

**Runtime (average) = 475 s (7 min 55 s)**

Dataset	Run	Runtime
influenza_98 threshold = 10	1	348 s (5 min 48 s)
	2	497 s (8 min 17 s)
	3	425 s (7 min 5 s)
	4	624 s (10 min 24 s)
	5	482 s (8 min 2 s)
	average	475 s (7 min 55 s)

**RF distance:** influenza\_98 , threshold = 10 (5 runs)

**normalized RF distance (average) = 0.473**

$$\text{normalized RF} = \frac{|S_1 \cup S_2| - |S_1 \cap S_2|}{\max RF}$$

Dataset	Run #	$ S_1 \cup S_2 $	$ S_1 \cap S_2 $	RF distance	RF max	Normalized RF distance
influenza_98 threshold = 10	1	263	125	138	190	0.525
	2	256	132	124	190	0.484
	3	251	137	114	190	0.454
	4	247	141	106	190	0.429
	5	254	134	120	190	0.472

## Normalized RF distance: influenza\_98 , threshold = 10

Pairwise: 5 runs with different random seeds

	Seed 1	Seed 2	Seed 3	Seed 4	Seed 5
Seed 1	0	0.552	0.574	0.558	0.579
Seed 2		0	0.536	0.574	0.547
Seed 3			0	0.541	0.579
Seed 4				0	0.530
Seed 5					0

## **Info RF distance: influenza\_98 , threshold = 10**

**Pairwise:** 5 runs with different random seeds

	<b>Seed 1</b>	<b>Seed 2</b>	<b>Seed 3</b>	<b>Seed 4</b>	<b>Seed 5</b>
<b>Seed 1</b>	0	0.903	0.909	0.859	0.880
<b>Seed 2</b>		0	0.846	0.927	0.858
<b>Seed 3</b>			0	0.808	0.830
<b>Seed 4</b>				0	0.837
<b>Seed 5</b>					0

## Shared Phylogenetic Info: influenza\_98, threshold = 10

**Pairwise:** 5 runs with different random seeds

	Seed 1	Seed 2	Seed 3	Seed 4	Seed 5
Seed 1	1	0.567	0.503	0.541	0.516
Seed 2		1	0.544	0.558	0.583
Seed 3			1	0.645	0.604
Seed 4				1	0.635
Seed 5					1

# Beyond Classic RF distance

Dataset	Threshold	RF	Info RF	Shared Phylogenetic Info	Different Phylogenetic Info	Mutual Clustering Info
influenza_98	5	0.547	0.547	0.723	0.277	0.727
	10	0.579	0.622	0.700	0.304	0.701
	15	0.684	0.764	0.614	0.386	0.634
influenza_1019	10					
	15					
	20					

# Method implementation on GitHub

<https://github.com/DerekSHAOZH/Global Phylogeny by Local Inferences>

The screenshot shows a GitHub repository page with the following details:

- Repository Name:** DerekSHAOZH / Global\_Phylology\_by\_Local\_Inferences (Public)
- Branches:** main (selected), 2 branches, 0 tags
- Code Protection:** Your main branch isn't protected. A button to "Protect this branch" is visible.
- Commits:** A list of 54 commits from the main branch, showing changes to various files like README.md, .gitkeep, and Python scripts (infer\_global.py, infer\_local.py, main.py, mst.py, rf\_distance.py, tree.py).
- Readme:** Readme file present.
- Statistics:** 1 star, 1 watching, 0 forks.
- Releases:** No releases published. A link to "Create a new release" is available.
- Packages:** No packages published. A link to "Publish your first package" is available.
- Contributors:** 2 contributors listed: DerekSHAOZH and megaphone (Meghan Kane).

# Outline

- Introduction
- Motivation
- Methods
- Materials
- Results
- **Discussion**
- Acknowledgement

# Discussion

## Summary of results interpretation

- Lower thresholds gave shorter runtime and better recall (accuracy)
- Inconsistency between global trees was possibly caused by local optima of MP
- Global phylogenetic analysis was prone to erroneous local inferences
- Other metrics beyond RF distance gave us other insights into tree similarity

## Further Research

- Improve accuracy
- Multiple starting tree topologies to avoid local optima (20 in RAxML by default)

# Outline

- Introduction
- Motivation
- Methods
- Materials
- Results
- Discussion
- **Contributions**

# Contributions

Method Design - Prabhav, Derek & Meghan

Method Implementation:

- Tree.py - Derek & Meghan
- Mst.py - Meghan
- Infer\_local.py - Derek
- Infer\_global.py - Derek & Meghan

Program Execution - Derek

- Our program: req. 50GB RAM
- RAxML benchmark: Apple M1, 8 cores, 16 GB RAM, GTR model, ML

Data Analyses - Meghan

# References

- Bayzid, M. S., Hunt, T., & Warnow, T. (2014). Disk covering methods improve phylogenomic analyses. *BMC Genomics*, 15(S6). <https://doi.org/10.1186/1471-2164-15-s6-s7>
- Böcker, Sebastian, et al. "The Generalized Robinson-Foulds Metric." *Lecture Notes in Computer Science*, 2013, pp. 156–169., [https://doi.org/10.1007/978-3-642-40453-5\\_13](https://doi.org/10.1007/978-3-642-40453-5_13).
- Choi, M. J., Tan, V. Y. F., Anandkumar, A., & Willsky, A. S. (2010, September 14). *Learning latent tree graphical models*. arXiv.org. Retrieved March 7, 2023, from <https://doi.org/10.48550/arXiv.1009.2722>
- Chor, B., & Tuller, T. (2006). Finding a maximum likelihood tree is hard. *Journal of the ACM*, 53(5), 722–744. <https://doi.org/10.1145/1183907.1183909>
- Crooks, G. E., Hon, G., Chandonia, J.-M., & Brenner, S. E. (2004). WebLogo: A sequence logo generator. *Genome Research*, 14(6), 1188–1190. <https://doi.org/10.1101/gr.849004>
- De Bruyn, A., Martin, D. P., & Lefevre, P. (2013). Phylogenetic Reconstruction Methods: An overview. *Methods in Molecular Biology*, 257–277. [https://doi.org/10.1007/978-1-62703-767-9\\_13](https://doi.org/10.1007/978-1-62703-767-9_13)
- Hagberg , A. A., Schult , D. A., & Swart , P. J. (2008). Exploring Network Structure, Dynamics, and Function using NetworkX. *Proceedings of the 7th Python in Science Conference (SciPy 2008)*.
- Huerta-Cepas, J., Serra, F., & Bork, P. (2016). ETE 3: Reconstruction, analysis, and visualization of Phylogenomic Data. *Molecular Biology and Evolution*, 33(6), 1635–1638. <https://doi.org/10.1093/molbev/msw046v>

# References

- Kalaghatgi, P. (2020). Phylogeny inference under the general markov model using MST-backbone.  
<https://doi.org/10.1101/2020.06.30.180315>
- Kalaghatgi, P., & Lengauer, T. (2017). Computing phylogenetic trees using topologically related minimumspanning trees. *Journal of Graph Algorithms and Applications*, 21(6), 1003–1025. <https://doi.org/10.7155/jgaa.00447>
- Kumar, S., Stecher, G., & Tamura, K. (2016). Mega7: Molecular evolutionary genetics analysis version 7.0 for bigger datasets. *Molecular Biology and Evolution*, 33(7), 1870–1874. <https://doi.org/10.1093/molbev/msw054>
- Ranwez, V., & Gascuel, O. (2001). Quartet-based phylogenetic inference: Improvements and limits. *Molecular Biology and Evolution*, 18(6), 1103–1116. <https://doi.org/10.1093/oxfordjournals.molbev.a003881>
- Roch, S. (2006). A short proof that phylogenetic tree reconstruction by maximum likelihood is hard. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 3(1), 92–94. <https://doi.org/10.1109/tcbb.2006.4>
- Smith, Martin R. *Comparing Splits Using Information Theory*,  
<https://cran.r-project.org/web/packages/TreeDist/vignettes/information.html>.
- Yermanos, A. D., Dounas, A. K., Stadler, T., Oxenius, A., & Reddy, S. T. (2018). Tracing antibody repertoire evolution by systems phylogeny. *Frontiers in Immunology*, 9. <https://doi.org/10.3389/fimmu.2018.02149>