

Data-driven whole-genome clustering to detect geospatial, temporal, and functional trends in SARS-CoV-2 evolution

Jean Merlet¹, John Lagergren², Verónica G. Melesse Vergara², Mikaela Cashman³, Christopher Bradburne⁴, Raina Plowright⁵, Emily Gurley⁴, Wayne Joubert², Daniel Jacobson²

¹University of Tennessee, ²Oak Ridge National Laboratory, ³Berkeley National Laboratory, ⁴Johns Hopkins University, ⁵Cornell University

ORNL is managed by UT-Battelle LLC for the US Department of Energy



U.S. DEPARTMENT OF
ENERGY

Covid-19 pandemic

- More than 750 million cases
- Over 200,000 new cases / week

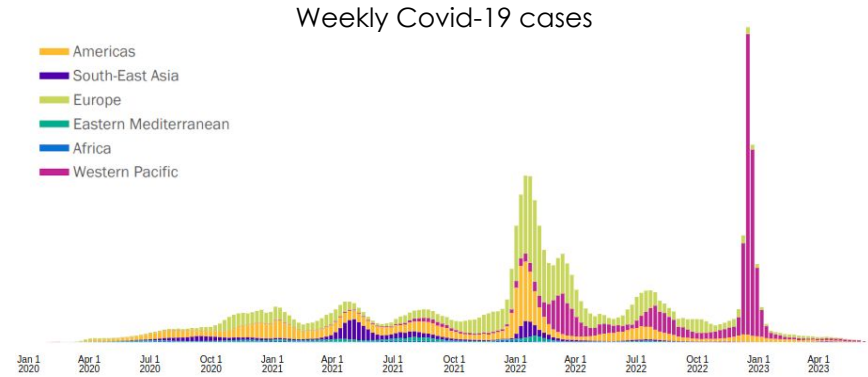
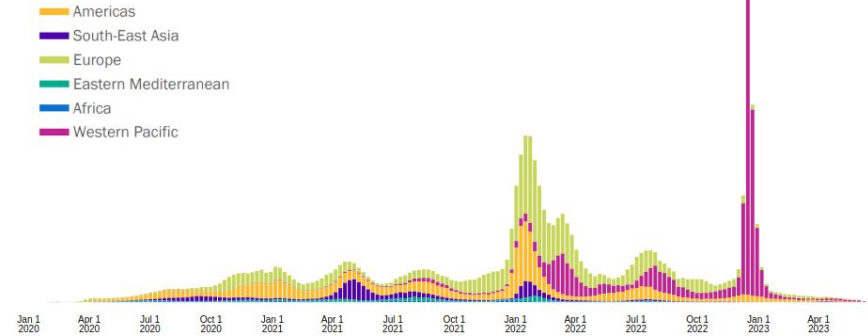


Image: WHO

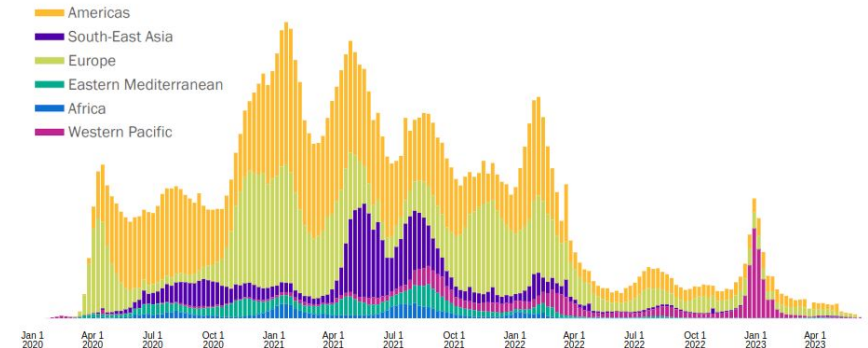
Covid-19 pandemic

- More than 750 million cases
- Over 200,000 new cases / week
- More than 7 million deaths
- Over 1,300 deaths / week

Weekly Covid-19 cases



Weekly Covid-19 deaths



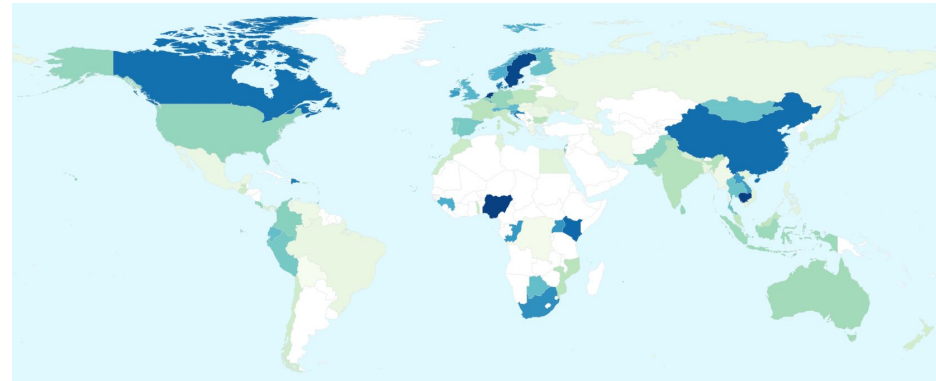
Images: WHO

Covid-19 pandemic

- More than 750 million cases
- Over 200,000 new cases / week
- More than 7 million deaths
- Over 1,300 deaths / week
- Nearly 16 million SARS-CoV-2 sequences on GISAID
- Over 17,000 uploads / month



Last 3 months of
sequence uploads



Percentage of COVID-19 cases shared via GISAID



Images: gisaid.org

SARS-CoV-2 genome sequences

- ~30,000 nucleotides
 - A, T, C, G, N, D



Image: astrochem.org

SARS-CoV-2 genome sequences

- ~30,000 nucleotides
 - A, T, C, G, N, D
- Known sequencing biases
 - Geospatial
 - Lag time
 - Selective sampling

Last 3 months of sequence uploads

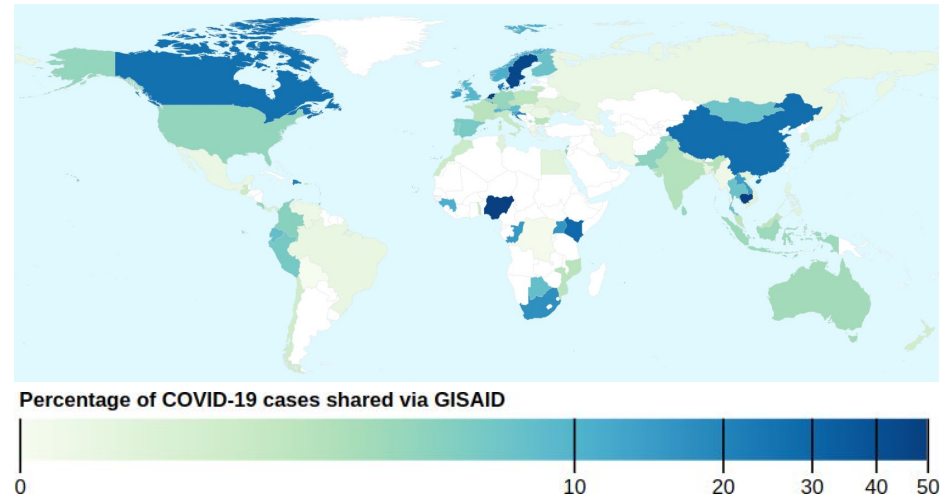
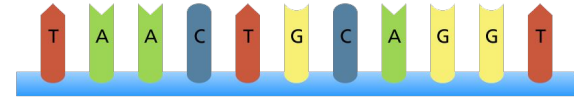


Image: gisaid.org

SARS-CoV-2 genome sequences

- ~30,000 nucleotides
 - A, T, C, G, N, D
- Known sequencing biases
 - Geospatial
 - Lag time
 - Selective sampling
- ~24,000 nucleotide positions have shown mutations

Original sequence



Point mutation

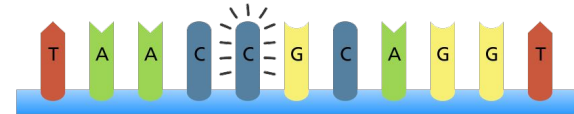


Image: yourgenome.org

Lineage classification

- **Phylogenetic Assignment of Named Global Outbreak**

Lineage classification

- **Phylogenetic Assignment of Named Global Outbreak**
- 3,053 lineages

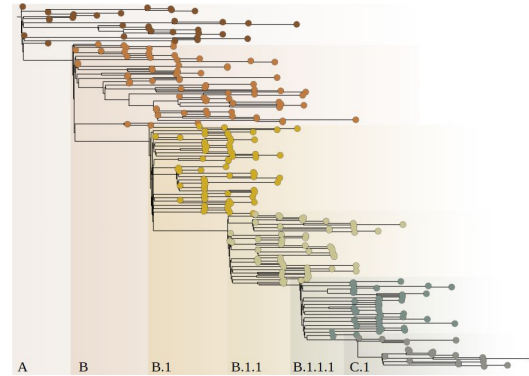


Image: pango.network

Lineage classification

- **Phylogenetic Assignment of Named Global Outbreak**
- 3,053 lineages
- Automated lineage assignment through pangolin

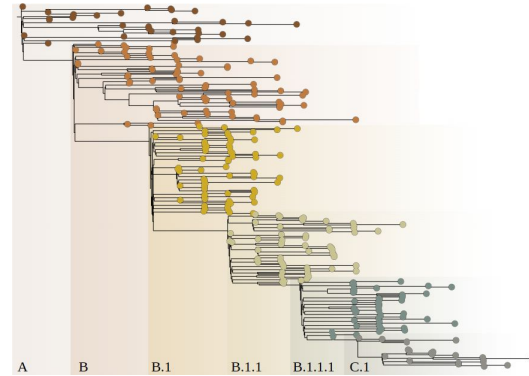


Image: pangolin.network

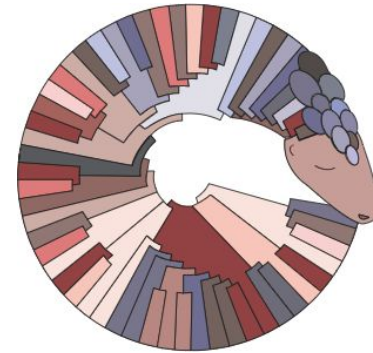


Image: cov-lineages.org

Lineage classification

- **Phylogenetic Assignment of Named Global Outbreak**
- 3,053 lineages
- Automated lineage assignment through pangolin
- Lineages are based on small, hand-curated sets of mutations

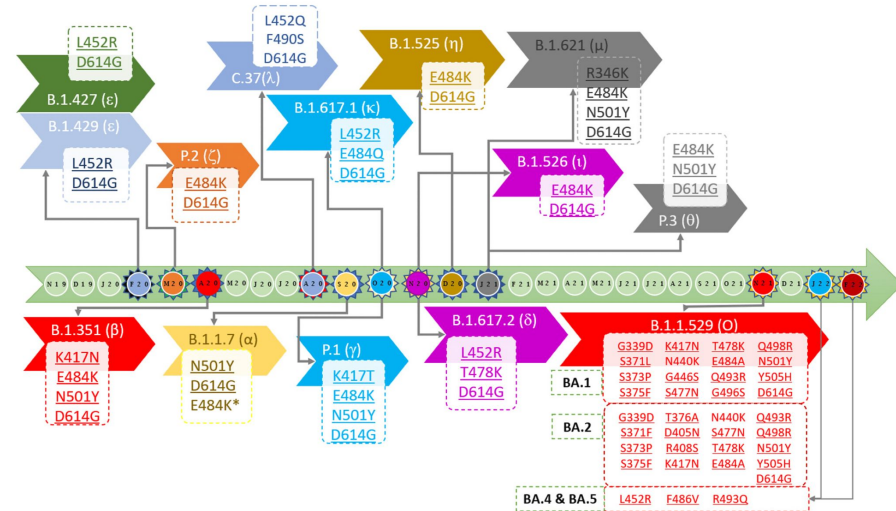


Image adapted from: Bhadane, R. and Salo-Ahen, O.M.H. 2022. High-Throughput Molecular Dynamics-Based Alchemical Free Energy Calculations for Predicting the Binding Free Energy Change Associated with the Selected Omicron Mutations in the Spike Receptor-Binding Domain of SARS-CoV-2. Biomedicines. MDPI AG.

Lineage classification

- **Phylogenetic Assignment of Named Global Outbreak**
- 3,053 lineages
- Automated lineage assignment through pangolin
- Lineages are based on small, hand-curated sets of mutations
- Need a data-driven, whole-genome model

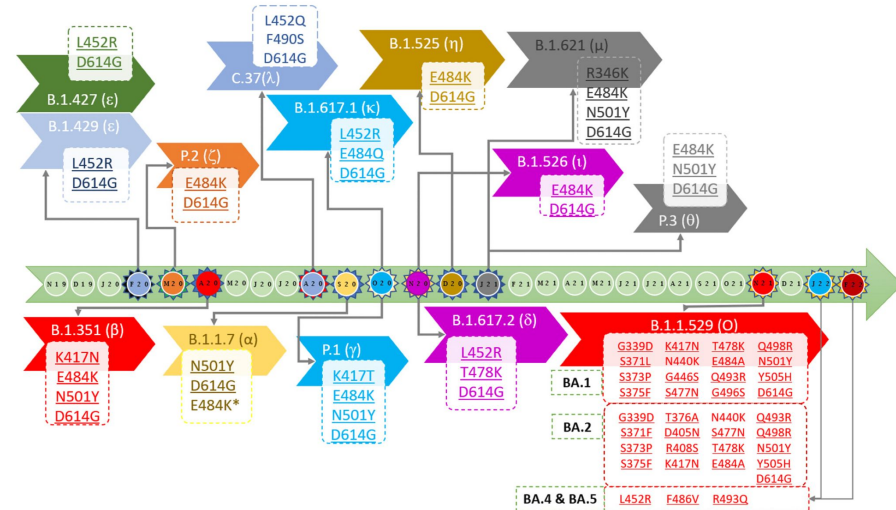


Image adapted from: Bhadane, R. and Salo-Ahen, O.M.H. 2022. High-Throughput Molecular Dynamics-Based Alchemical Free Energy Calculations for Predicting the Binding Free Energy Change Associated with the Selected Omicron Mutations in the Spike Receptor-Binding Domain of SARS-CoV-2. Biomedicines. MDPI AG.

Overview: data-driven strains

Millions of genomes

```
seq. 1: ATTAAG ...  
seq. 2: ATTAAG ...  
seq. 3: ATTAAG ...  
...
```

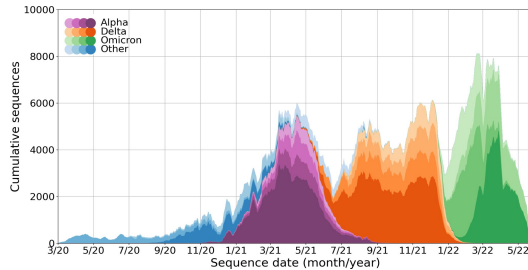
**Combinatorial
Metrics (CoMet)**

```
seq. 1: 0000101 ...  
seq. 2: 0000101 ...  
seq. 3: 0000101 ...  
...
```

**High Performance
Computing (HPC)**



Spatiotemporal trajectories
for data-driven lineages



**High Perf. Markov
Clustering (HipMCL)**

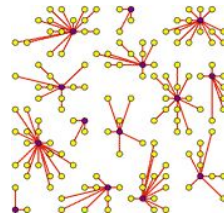


Image: micans.org/mcl

Genome-genome
similarity network

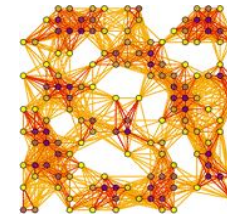


Image: micans.org/mcl

Overview: data-driven strains

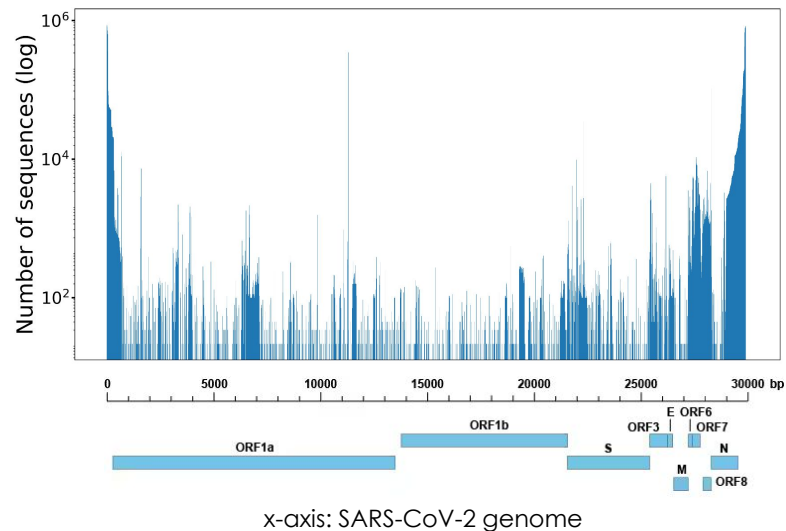
Millions of genomes

```
seq. 1: ATTAAAG ...  
seq. 2: ATTAAAG ...  
seq. 3: ATTAAAG ...  
...
```

Data pre-processing

- Challenges
 - Low-quality sequences

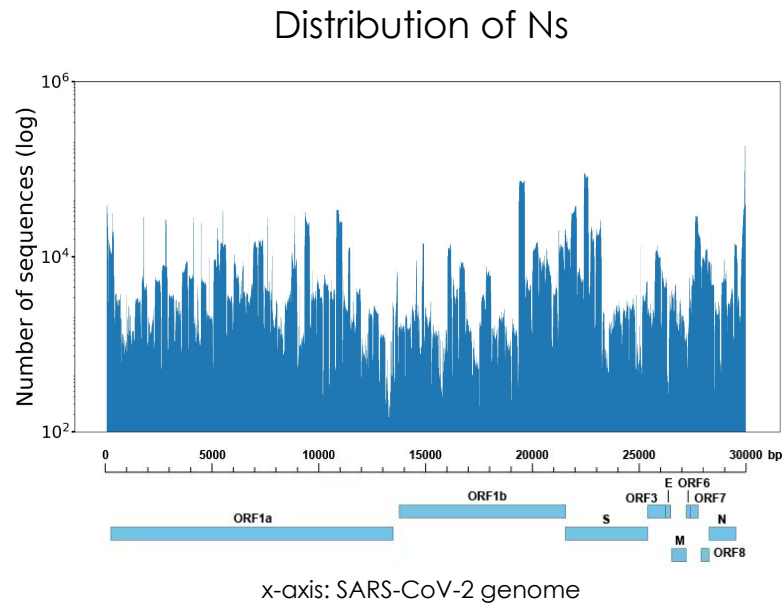
Distribution of deletions (Ds)



Genome adapted from novusbio.com

Data pre-processing

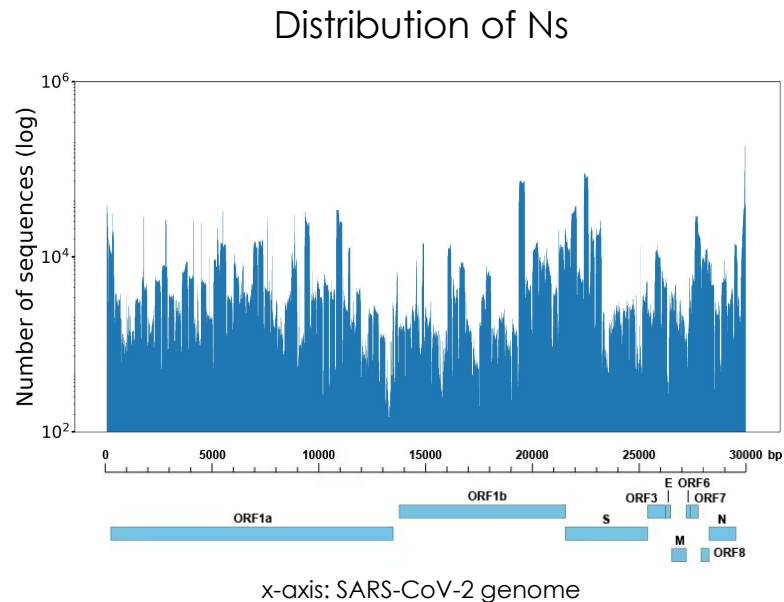
- Challenges
 - Low-quality sequences
 - Sequencing error



Genome adapted from novusbio.com

Data pre-processing

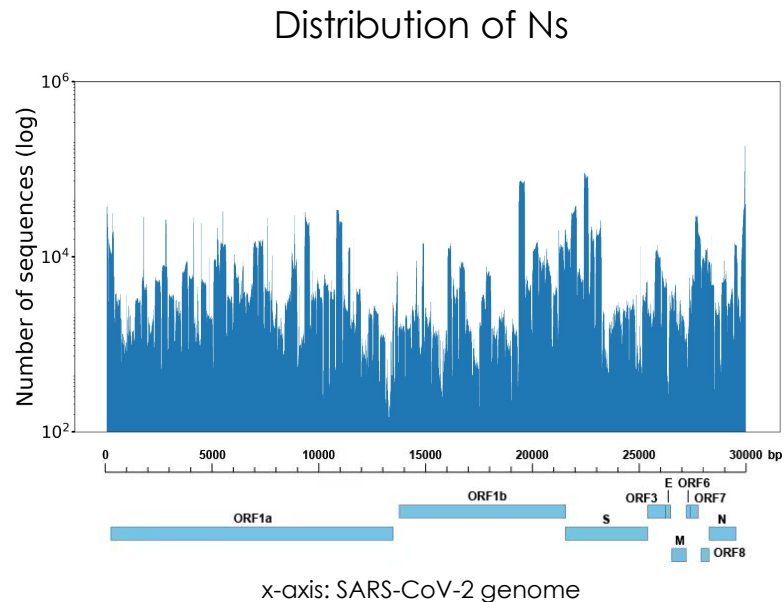
- Challenges
 - Low-quality sequences
 - Sequencing error
 - Incomplete metadata



Genome adapted from novusbio.com

Data pre-processing

- Challenges
 - Low-quality sequences
 - Sequencing error
 - Incomplete metadata
- 11.0 -> 7.7 million sequences

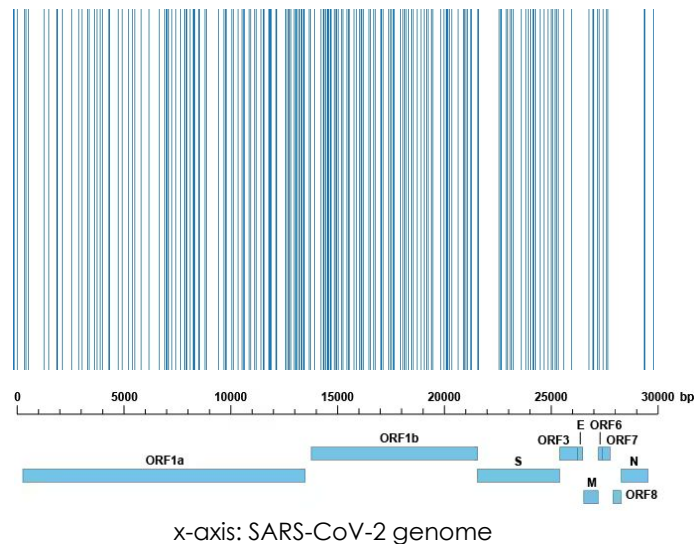


Genome adapted from novusbio.com

Data pre-processing

- Challenges
 - Low-quality sequences
 - Sequencing error
 - Incomplete metadata
- 11.0 -> 7.7 million sequences
- Detect significant mutation positions

Distribution of Ns

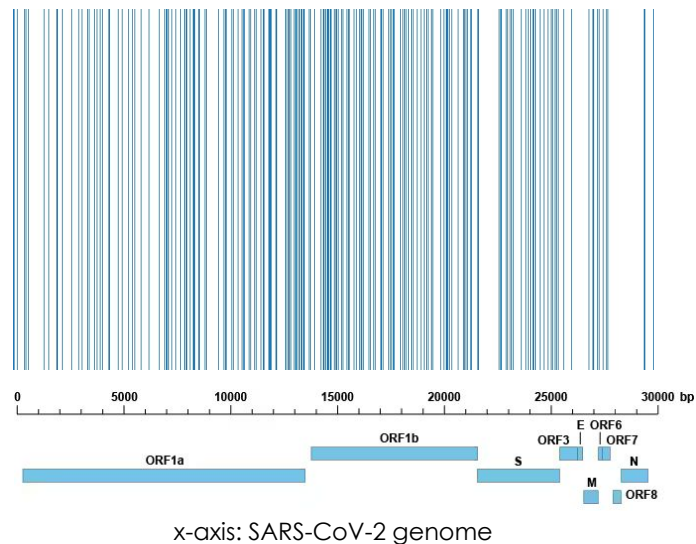


Genome adapted from novusbio.com

Data pre-processing

- Challenges
 - Low-quality sequences
 - Sequencing error
 - Incomplete metadata
- 11.0 -> 7.7 million sequences
- Detect significant mutation positions
- 29.9 -> 21.4 thousand nucleotides

Distribution of Ns

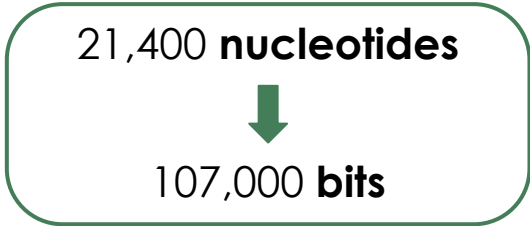
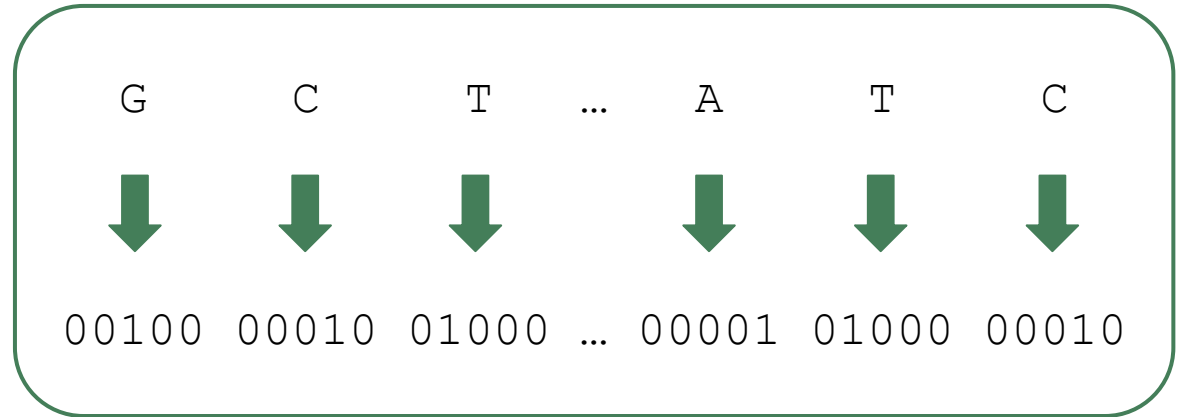


Genome adapted from novusbio.com

SARS-CoV-2 genome vector binarization

Binarization scheme

A: 00001
C: 00010
G: 00100
T: 01000
D: 10000
N: 00000



Duo similarity metric

- Compares binary vectors

	G		C		T	
00	1	00	00010	01	0000	...
00	1	00	00001	01	0000	...
	G		A		T	

Duo similarity metric

- Compares binary vectors
- 2-way
 - i and j binary vectors
 - N total vectors
 - $q = 2$ scaling constant
 - D_{ij} comparison frequency
 - f_i vector i binary frequency

	G		C		T	
00	1	00	000	10	01	0000 ...
00	1	00	0000	01	01	0000 ...
	G		A		T	

$$\text{Duo}_{ij} = 5D_{ij} \left(1 - \frac{f_i}{q}\right) \left(1 - \frac{f_j}{q}\right)$$

$$D_{ij} = \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{i_n = 1, j_n = 1\}$$

Duo similarity metric

- Compares binary vectors
- 2-way
 - i and j binary vectors
 - N total vectors
 - $q = \frac{2}{3}$ scaling constant
 - D_{ij} comparison frequency
 - f_i^{ij} vector i binary frequency
- Choice of $q = 0$ converts to Sørensen-Dice index

	G		C		T	
00	1	00	000	10	01	000 ...
00	1	00	0000	1	01	000 ...
	G		A		T	

$$SD_{ij} = 5D_{ij}$$

$$D_{ij} = \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{i_n = 1, j_n = 1\}$$

Duo similarity metric

- Compares binary vectors
- 2-way
 - i and j binary vectors
 - N total vectors
 - $q = \frac{2}{3}$ scaling constant
 - D_{ij} comparison frequency
 - f_i^{ij} vector i binary frequency
- Choice of $q = 0$ converts to Sørensen-Dice index
- Also have 3-way capability

	G		C		T	
00	1	00	000	10	01	000 ...
00	1	00	0000	1	01	000 ...
	G		A		T	

$$SD_{ij} = 5D_{ij}$$

$$D_{ij} = \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{i_n = 1, j_n = 1\}$$

Overview: data-driven strains

Millions of genomes

```
seq. 1: ATTAAG ...  
seq. 2: ATTAAG ...  
seq. 3: ATTAAG ...  
...
```



**Combinatorial
Metrics (CoMet)**

```
seq. 1: 0000101 ...  
seq. 2: 0000101 ...  
seq. 3: 0000101 ...  
...
```



**High Performance
Computing (HPC)**



Genome-genome
similarity network

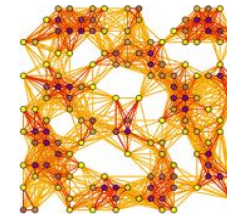
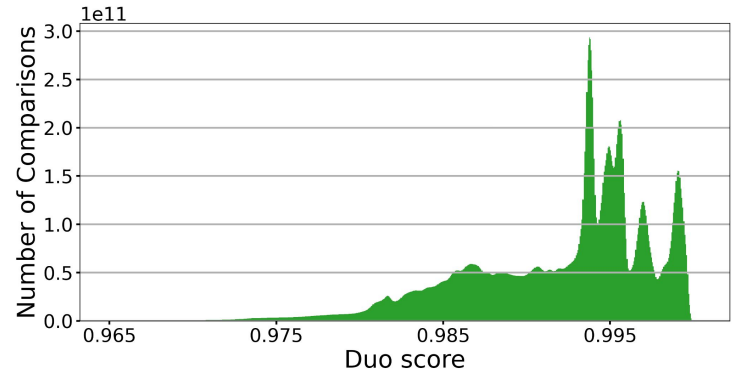


Image: micans.org/mcl



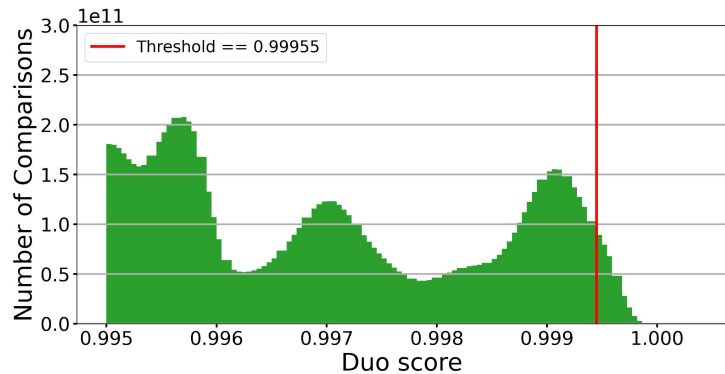
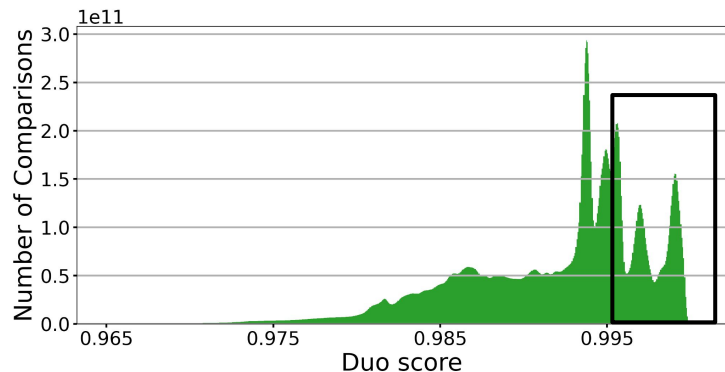
CoMet: Network generation

- 7.7 million binarized vectors
- Vectors each 107,000 bits
- 2.98×10^{13} vector comparisons



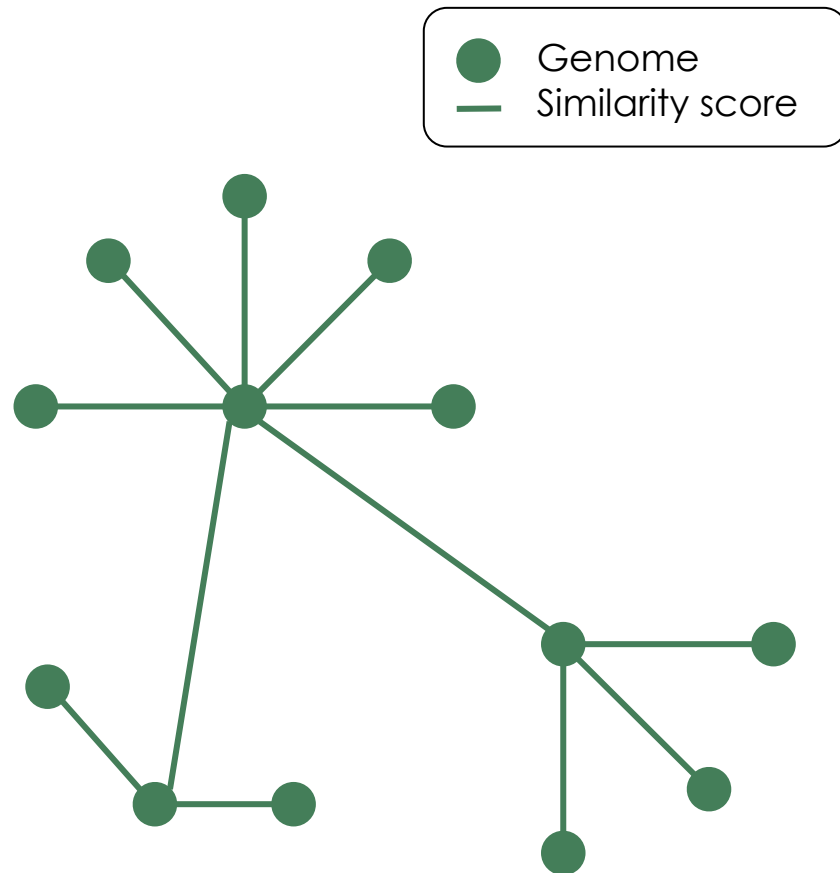
CoMet: Network generation

- 7.7 million binarized vectors
- Vectors each 107,000 bits
- 2.98×10^{13} vector comparisons
- Top 0.1% most similar vectors
 - 5.3 million nodes
 - 778 trillion edges
 - 147 thousand edges / node



CoMet: Network generation

- 7.7 million binarized vectors
- Vectors each 107,000 bits
- 2.98×10^{13} vector comparisons
- Top 0.1% most similar vectors
 - 5.3 million nodes
 - 778 trillion edges
 - 147 thousand edges / node
- Genome-to-genome network



CoMet

- Deployed on various HPC systems
 - **Summit** (#5 TOP500 list)
 - **Frontier** (#1 TOP500 list)
 - JUWELS Booster
 - Perlmutter
- Awarded Gordon Bell in 2018



Images: olcf.ornl.gov, fz-juelich.de, nersc.gov

CoMet

- Deployed on various HPC systems
 - **Summit** (#5 TOP500 list)
 - **Frontier** (#1 TOP500 list)
 - JUWELS Booster
 - Perlmutter
- Awarded Gordon Bell in 2018
- Application agnostic



Images: olcf.ornl.gov, fz-juelich.de, nersc.gov

CoMet

- Deployed on various HPC systems
 - **Summit** (#5 TOP500 list)
 - **Frontier** (#1 TOP500 list)
 - JUWELS Booster
 - Perlmutter
- Awarded Gordon Bell in 2018
- Application agnostic
- Record-setting computations

2.34 ExaFLOPs



6.6 ExaFLOPs



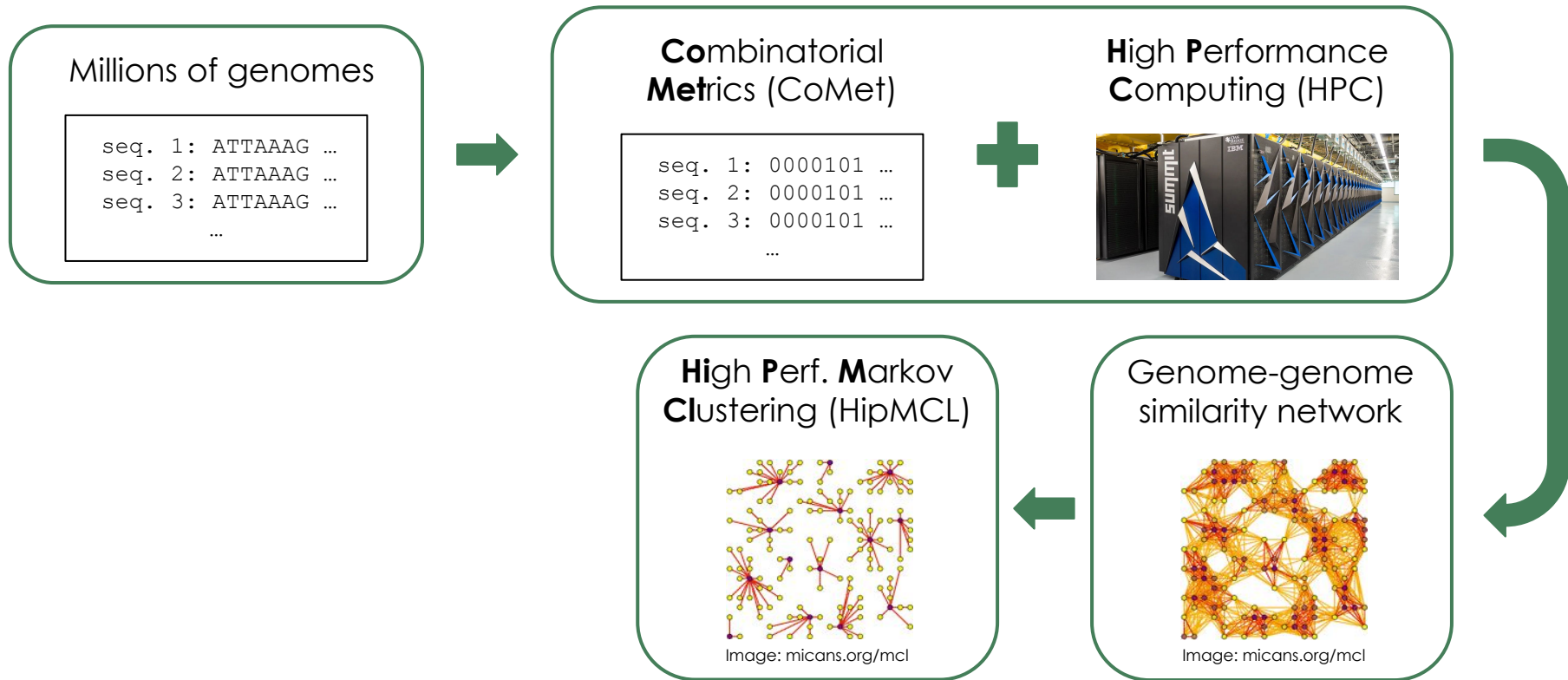
9.37 ExaFLOPs



TBD ExaFLOPs

Images: olcf.ornl.gov, fz-juelich.de, nersc.gov

Overview: data-driven strains



Markov clustering with HipMCL*

- Unsupervised algorithm

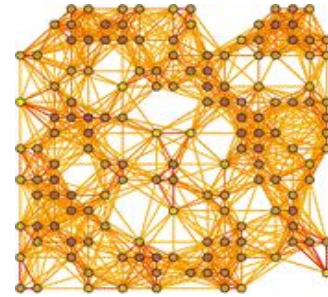
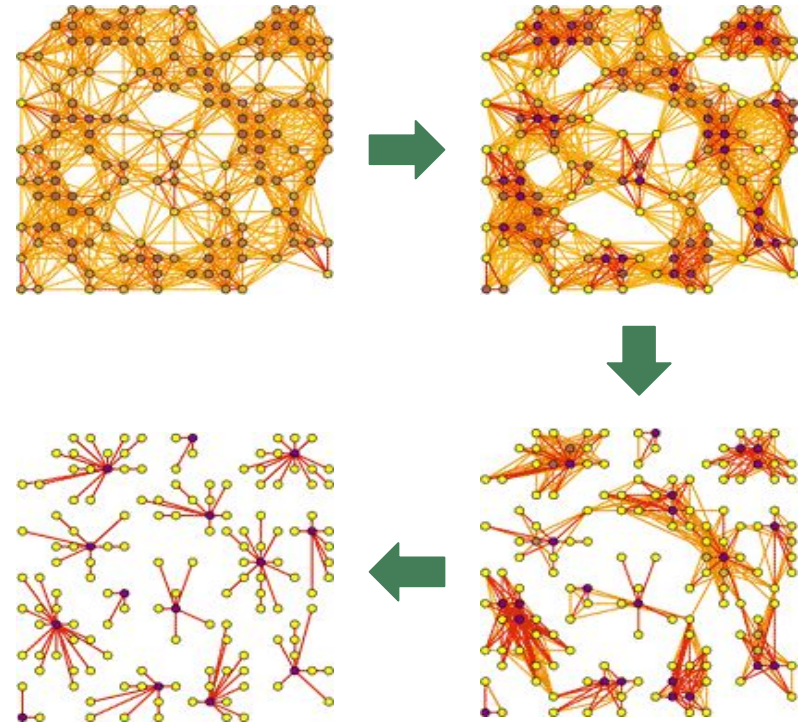


Image: micans.org/mcl

* Azad, A. et al. 2018. HipMCL: a high-performance parallel implementation of the Markov clustering algorithm for large-scale networks. Nucleic Acids Research. Oxford University Press (OUP).

Markov clustering with HipMCL

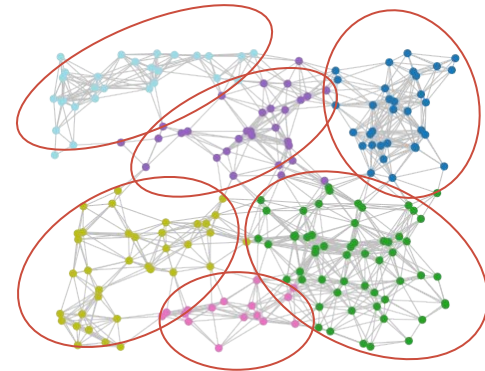
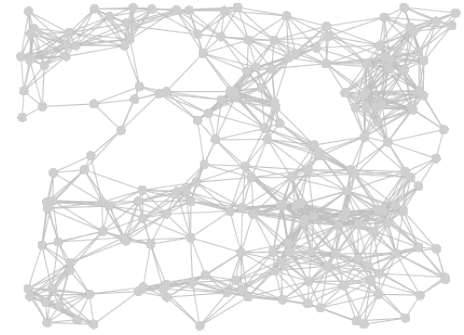
- Unsupervised algorithm
- Random walk
 - expansion (flow spreads)
 - inflation (flow recedes)



Images: micans.org/mcl

Markov clustering with HipMCL

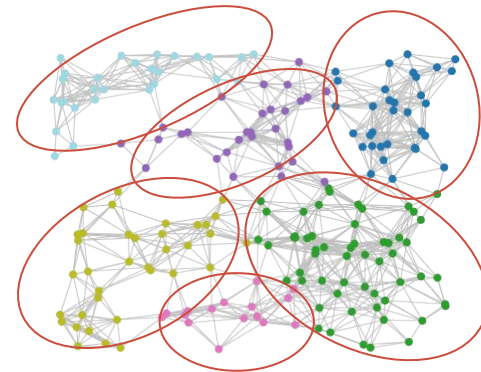
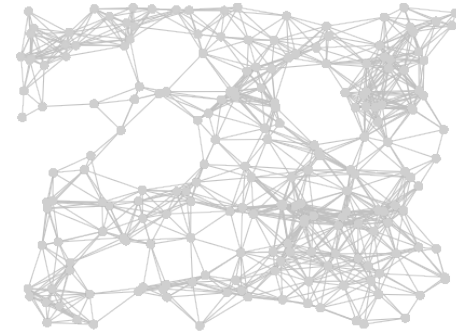
- Unsupervised algorithm
- Random walk
 - expansion (flow spreads)
 - inflation (flow recedes)
- Network \rightarrow clusters



Images: github.com/GuyAllard/markov_clustering

Markov clustering with HipMCL

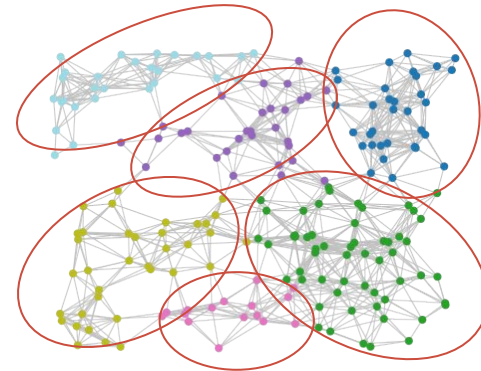
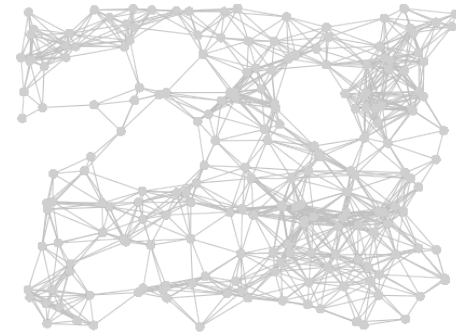
- Unsupervised algorithm
- Random walk
 - expansion (flow spreads)
 - inflation (flow recedes)
- Network -> clusters
- 22,000 SARS-CoV-2 strains



Images: github.com/GuyAllard/markov_clustering

Markov clustering with HipMCL

- Unsupervised algorithm
- Random walk
 - expansion (flow spreads)
 - inflation (flow recedes)
- Network -> clusters
- 22,000 SARS-CoV-2 strains
- Cluster number driven by graph topology



Images: github.com/GuyAllard/markov_clustering

Overview: data-driven strains

Millions of genomes

```
seq. 1: ATTAAG ...  
seq. 2: ATTAAG ...  
seq. 3: ATTAAG ...  
...
```

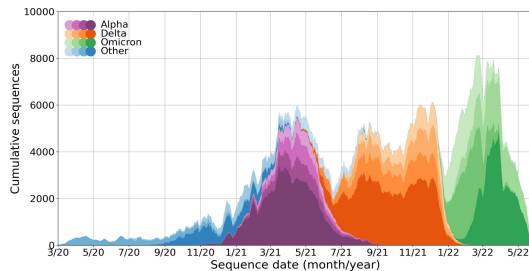
**Combinatorial
Metrics (CoMet)**

```
seq. 1: 0000101 ...  
seq. 2: 0000101 ...  
seq. 3: 0000101 ...  
...
```

**High Performance
Computing (HPC)**



Spatiotemporal trajectories
for data-driven lineages



**High Perf. Markov
Clustering (HipMCL)**

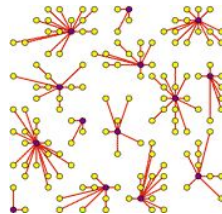


Image: micans.org/mcl

Genome-genome
similarity network

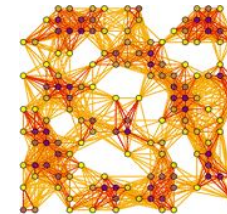
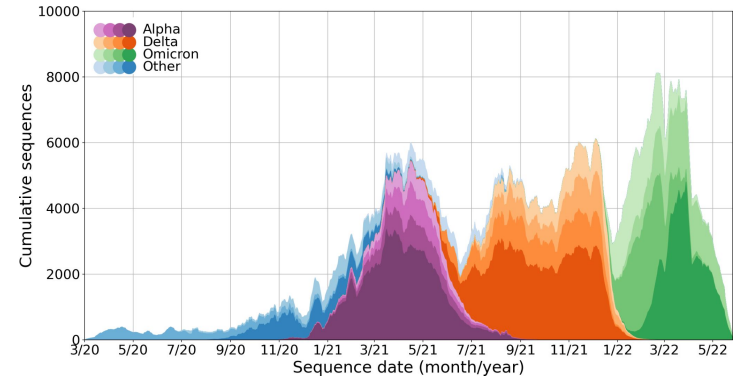


Image: micans.org/mcl

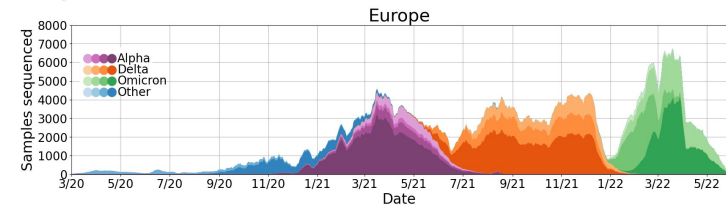
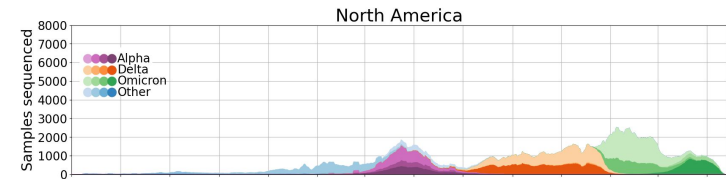
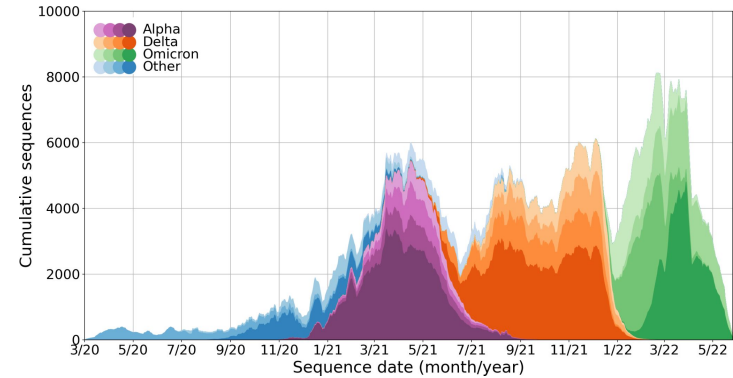
Our SARS-CoV-2 strains and WHO variants

- Mapped clusters to WHO variant
- Four largest clusters per WHO variant
- Longitudinal view



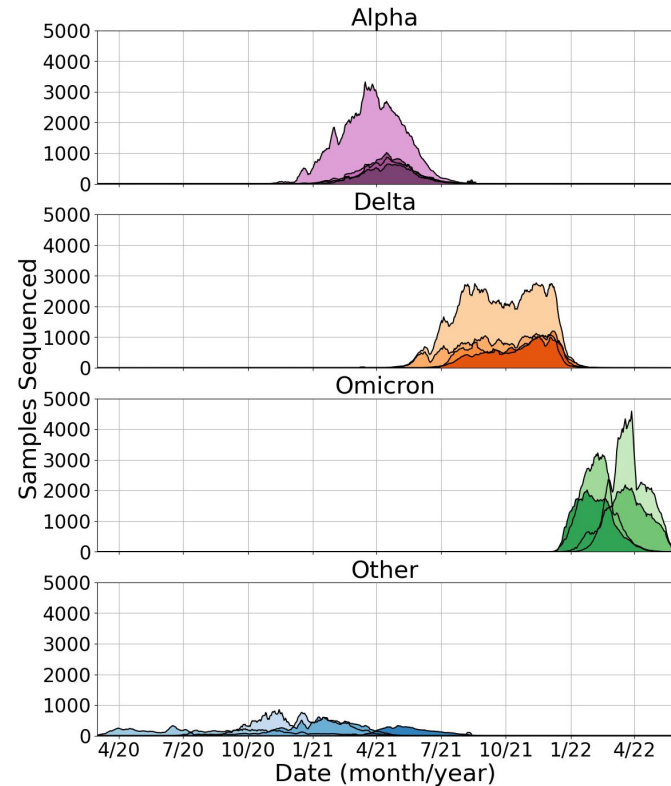
Our SARS-CoV-2 strains and WHO variants

- Mapped clusters to WHO variant
- Four largest clusters per WHO variant
- Longitudinal view
 - By continent



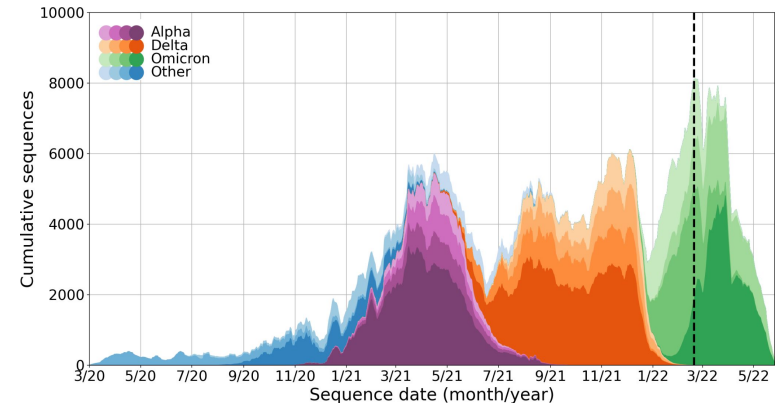
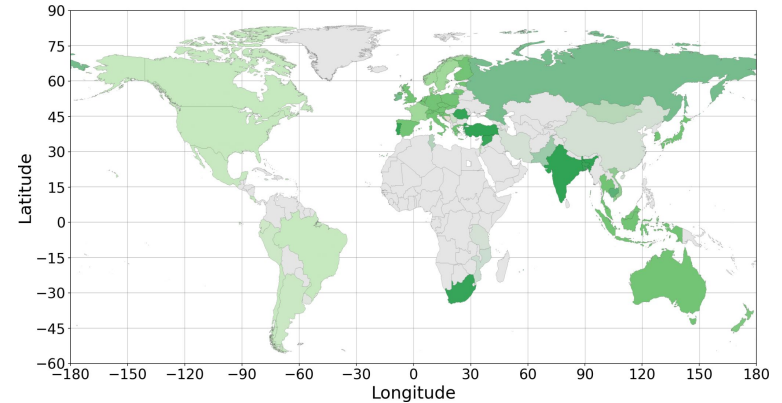
Our SARS-CoV-2 strains and WHO variants

- Mapped clusters to WHO variant
- Four largest clusters per WHO variant
- Longitudinal view
 - By continent
 - By WHO variant



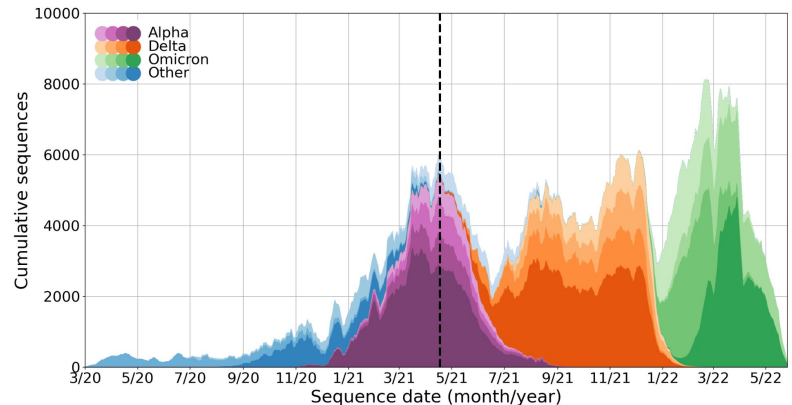
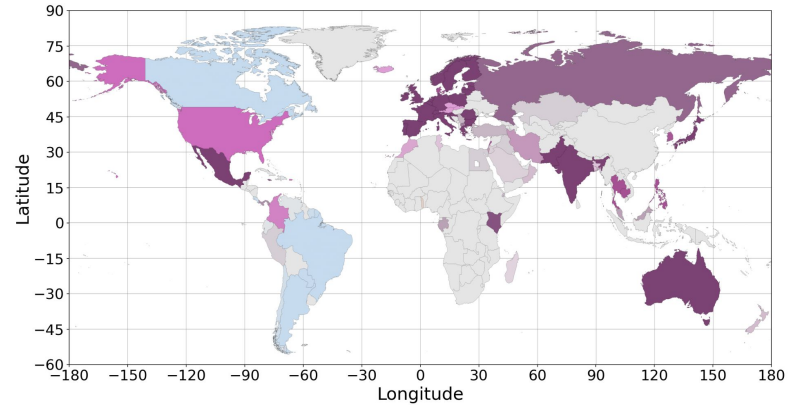
Our SARS-CoV-2 strains and WHO variants

- Mapped clusters to WHO variant
- Four largest clusters per WHO variant
- Longitudinal view
 - By continent
 - By WHO variant
- Geospatial view



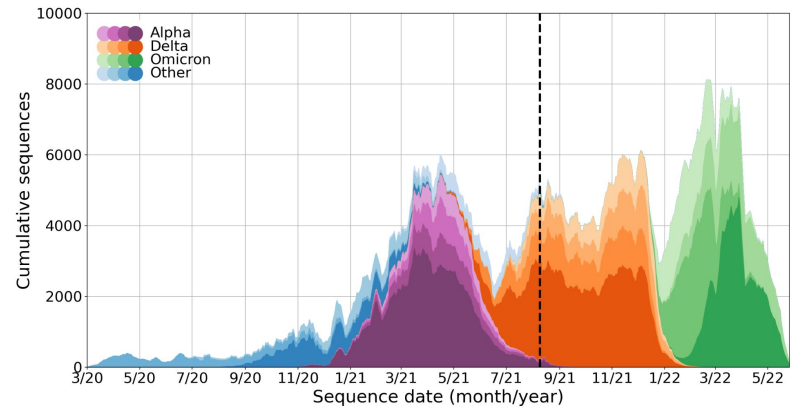
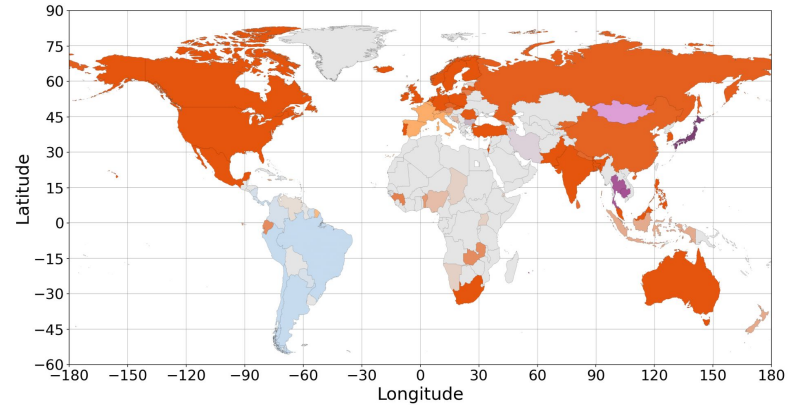
Geospatial asynchronicity - Alpha

- Geospatial - longitudinal view
- 10 largest clusters (dominant strains with most samples)
- Comprised x % of all sequences
- Shared country borders, different strains



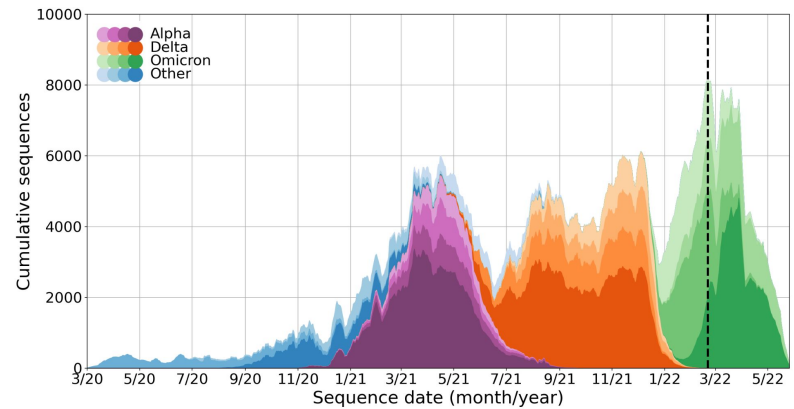
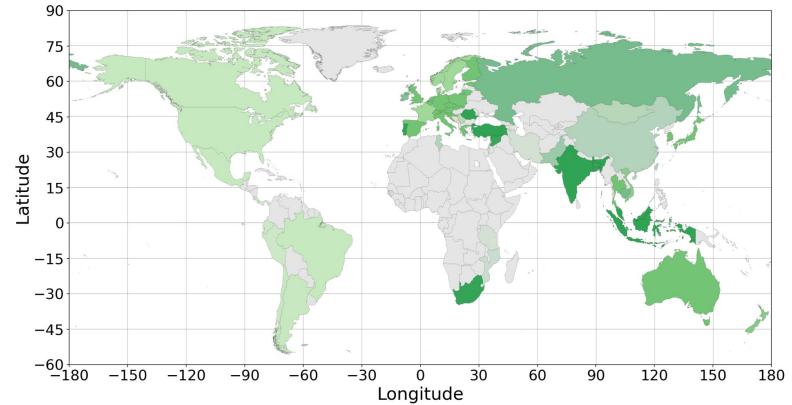
Geospatial asynchronicity - Delta

- Geospatial - longitudinal view
- 10 largest clusters (dominant strains with most samples)
- Comprised x % of all sequences
- Shared country borders, different strains



Geospatial asynchronicity - Omicron

- Geospatial - longitudinal view
- 10 largest clusters (dominant strains with most samples)
- Comprised x % of all sequences
- Shared country borders, different strains

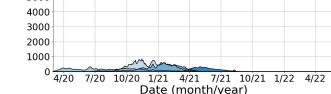
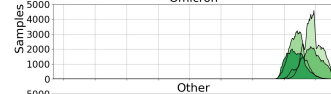
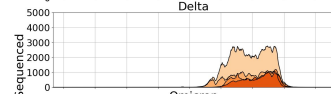
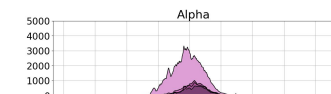
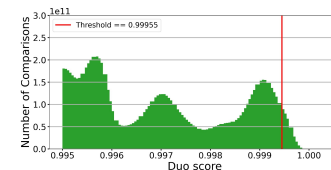
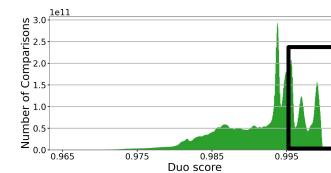
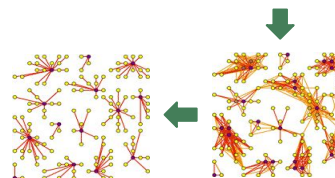
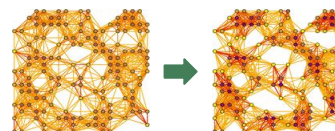
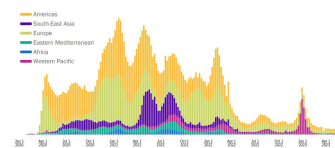
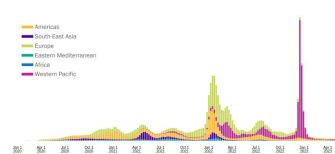


Geospatial asynchronicity - early pandemic

- Early pandemic view on Frontier
- 10 largest clusters (dominant strains with most samples)
- Shared country borders, different strains

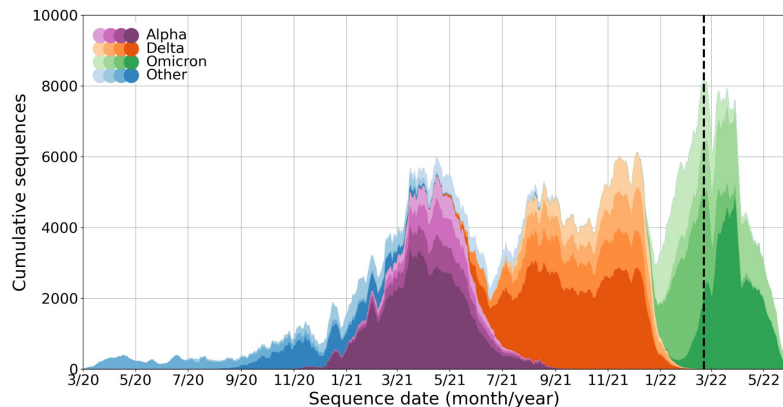
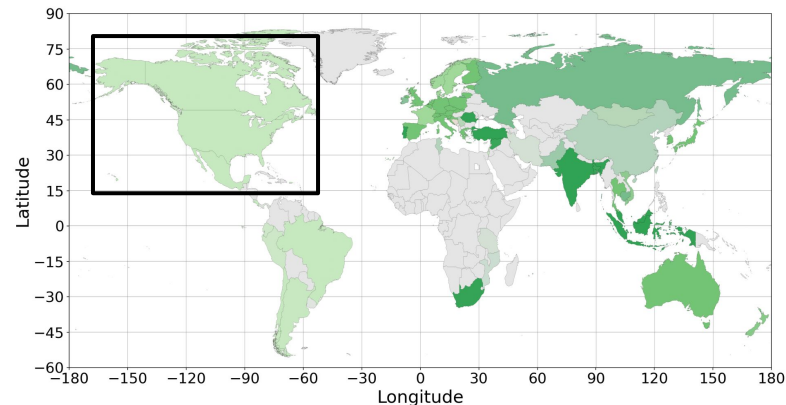
Summary

- All-against-all, whole-genome vector comparison
- > 22,000 SARS-CoV-2 strains
- Geospatial and temporal trajectories for each strain
- Spatially asynchronicity of dominant strains
- Possibly different pattern of epidemiological risk



Future work

- County or city geo-resolution
- Triplets of vectors (3-way)
- Assimilation of new samples
- Recombination analyses
 - Phylogenetic trees
 - exclude recombination
- Relationship to other variables
 - Environmental
 - Demographic
 - Mortality rate



Acknowledgements

Oak Ridge National Lab

Dan Jacobson

John Lagergren

Verónica G. Melesse Vergara

Wayne Joubert



Johns Hopkins University

Christopher Bradburne

Emily Gurley

Berkeley National Lab

Mikaela Cashman



Cornell University

Raina Plowright

Funding

DE-AC02-05CH11231

DE-AC02-06CH11357

DE-AC05-00OR22725

DE-AC02-98CH10886

EF-2133763

DOE INCITE

Questions?

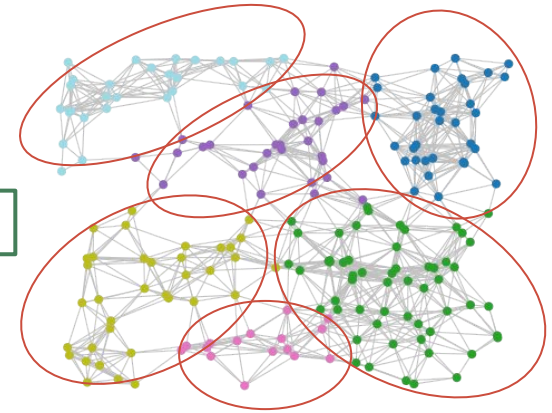
merletji@ornl.gov

jacobsonda@ornl.gov

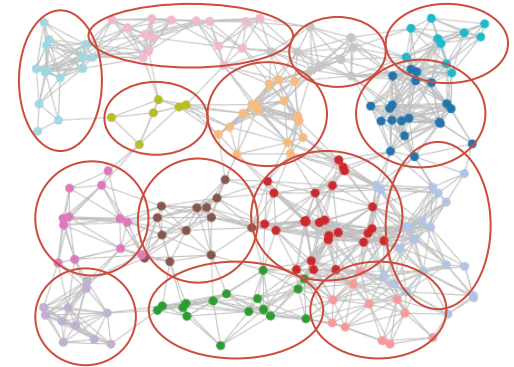
Markov clustering with HipMCL

- Unsupervised algorithm
- Random walk
 - expansion (flow spreads)
 - inflation (flow recedes)
- Network \rightarrow clusters
- 22,000 SARS-CoV-2 strains
- Cluster number affected by
 - graph topology
 - inflation value

Low inflation



High inflation



Images: github.com/GuyAllard/markov_clustering