# Pandemic-scale Phylogenetics
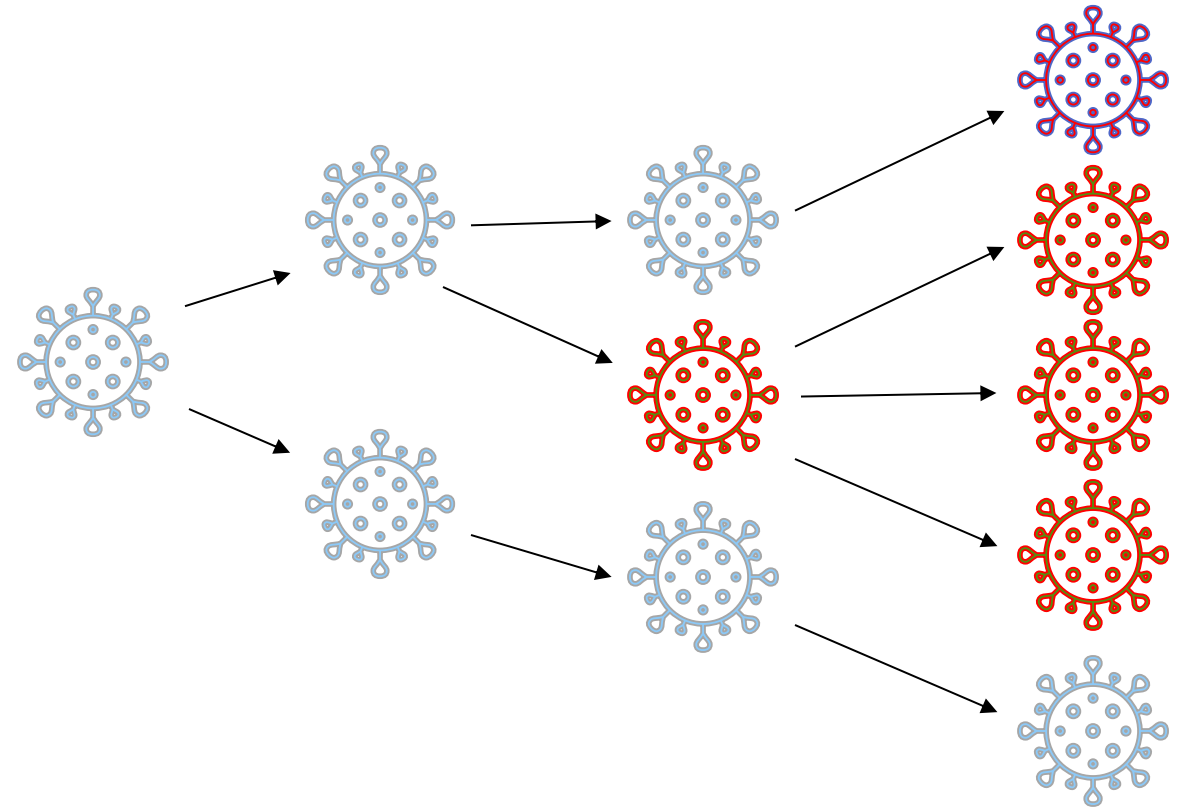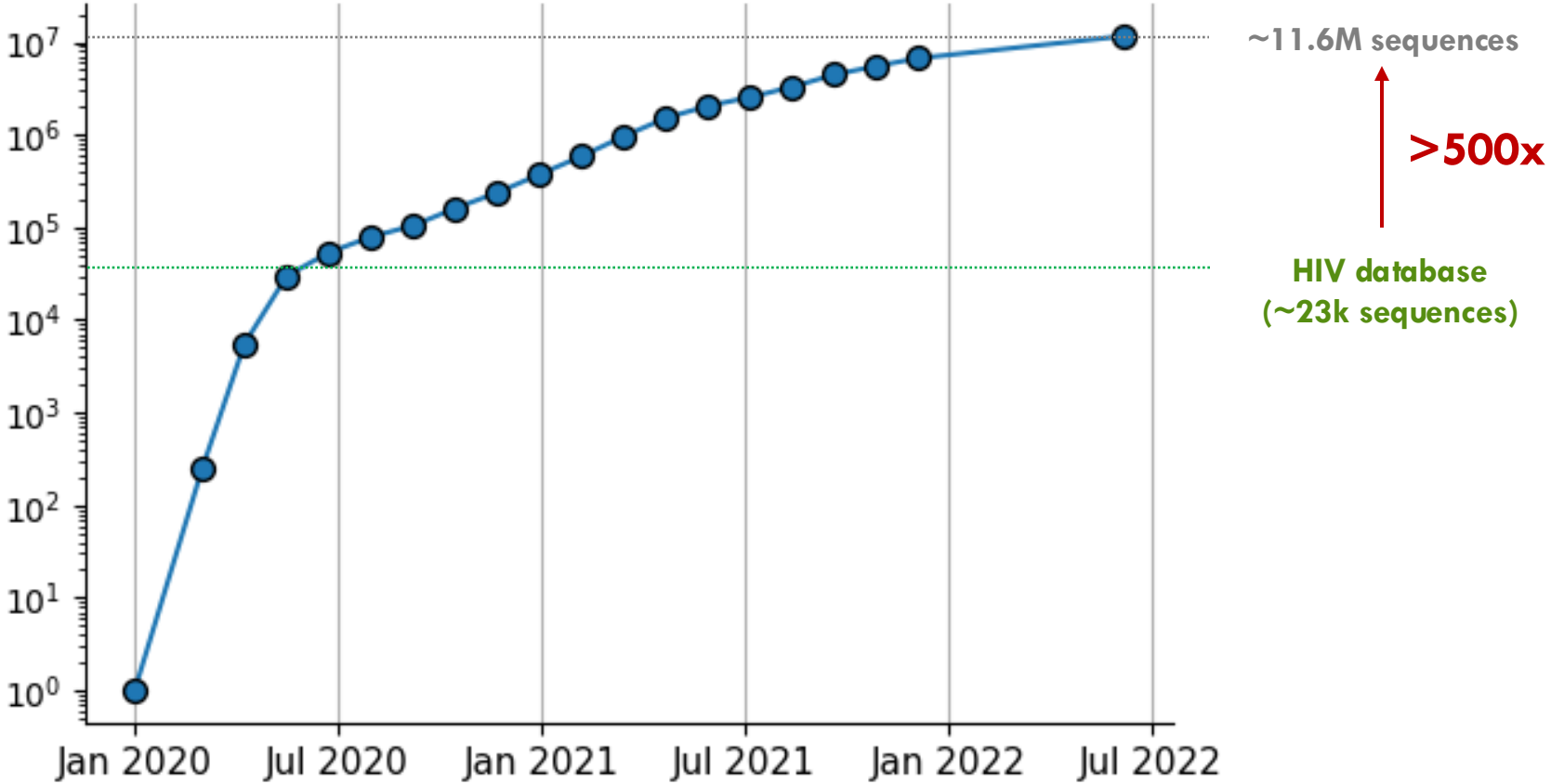
Yatish Turakhia

Assistant Professor, UC San Diego

# COVID-19 virus (SARS-CoV-2) is constantly mutating

- As the COVID-19 virus (SARS-COV-2) spreads, it **mutates**

- Certain mutations render the virus more **contagious**, **virulent** or capable of **evading** the vaccines and antibody-based therapies

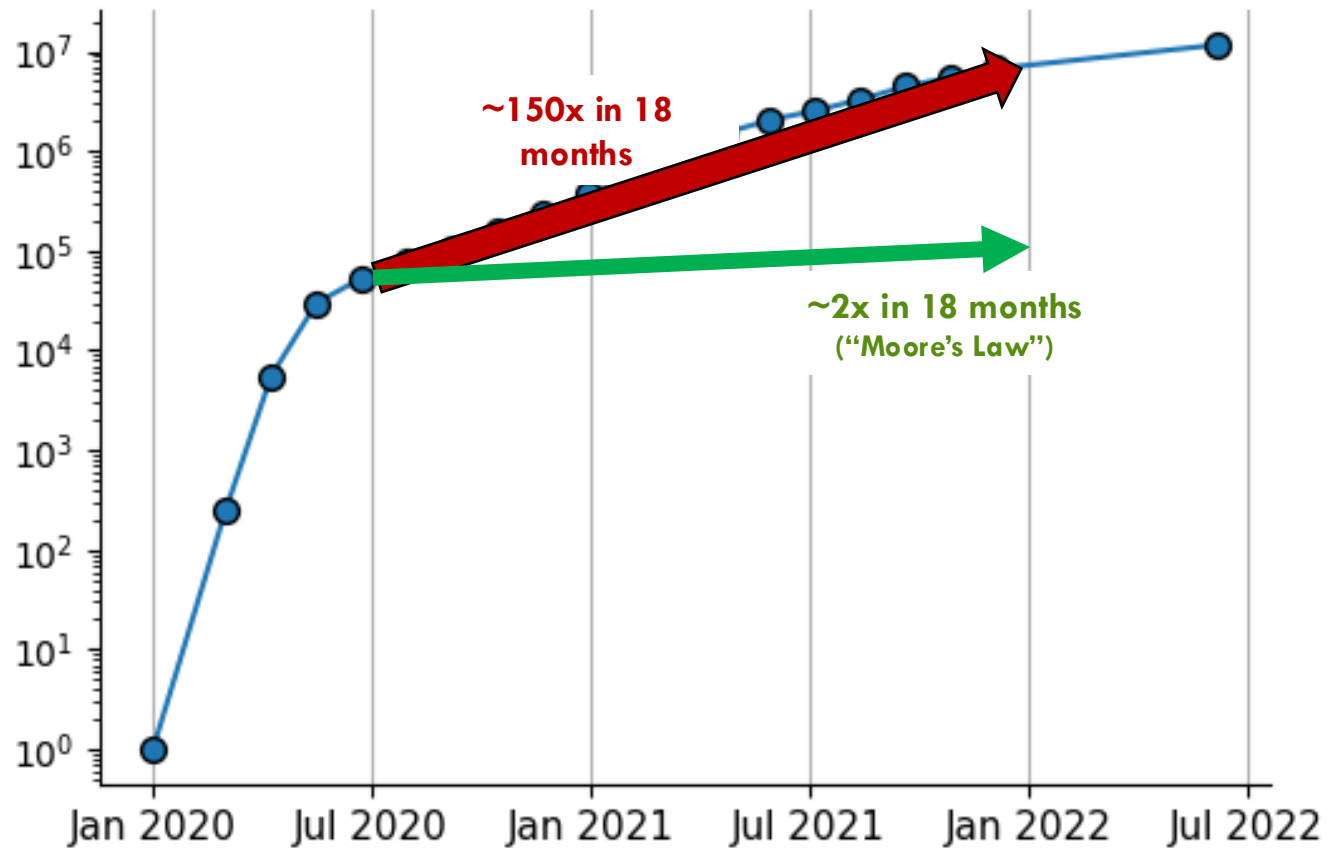- Genome sequencing helps **monitor** the **viral mutations** and the **evolutionary dynamics**

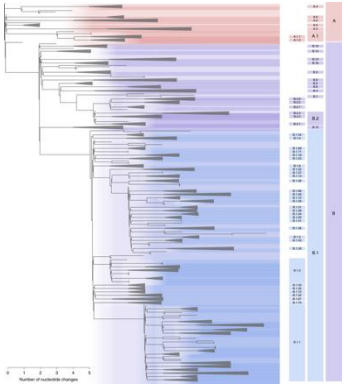# Number of Global SARS-CoV-2 Genome Sequences

# Number of Global SARS-CoV-2 Genome Sequences



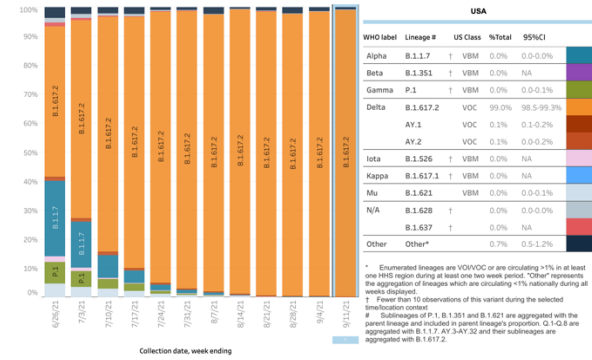~150x in 18 months

~75x higher semi-logarithmic slope compared to Moore's Law

~2x in 18 months ("Moore's Law")

# Naming lineages



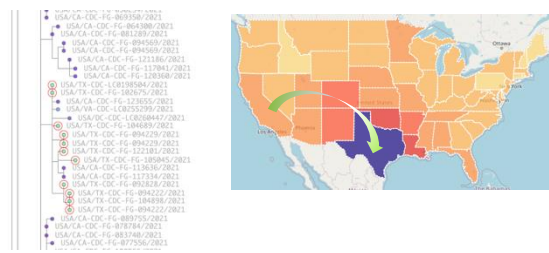(Rambaut et al., Nat. Microbiol. 2020)
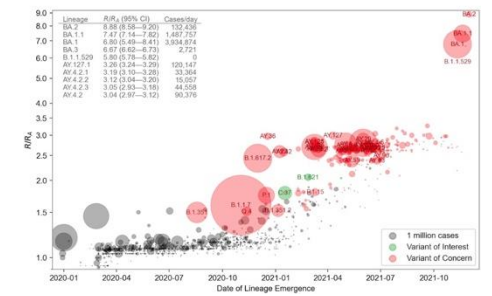
# Monitoring circulating lineages



(CDC.gov dashboard 2021)

# Identify newly-introduced strains

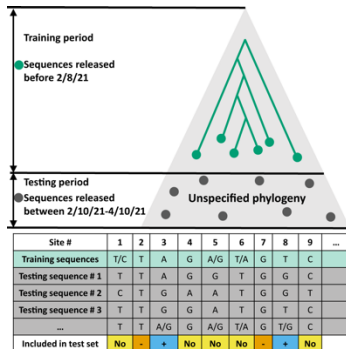

(McBroome et al., Virus Evol. 2022)

# Predicting fitness of a new strain



(Obermeyer et al., Science 2022)
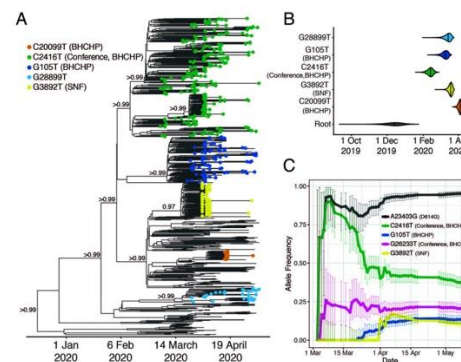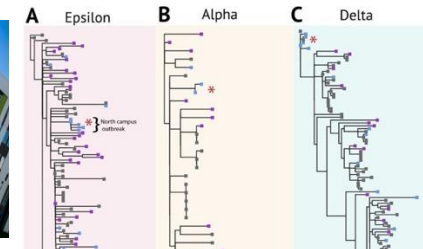
# UShER

# Predicting the next mutation



(Hallak et al., Nat. Comm Biol. 2022)

# Analyze outbreaks and superspreader events



(Lemieux et al., Science 2021)

# Wastewater surveillance



(Karthikeyan et al., Nature 2022)

# Overview of the UShER Package

- **UShER:** Phylogenetic placement
  - *Turakhia et al.,* **Nature Genetics** 2021

- **matOptimize:** Phylogenetic tree optimization
  - *Ye et al.,* **Bioinformatics** 2022

- **RIPPLES:** Find recombinant sequences using a phylogenomic approach
  - *Turakhia et al.,* bioRxiv 2021 (under revision, **Nature**)

- **matUtils:** Command-line tools for rapidly analyzing and interpreting SARS-CoV-2 mutation-annotated phylogenetic trees
  - *McBroome et al.,* Molecular Biology and Evolution (**MBE**) 2021

# Acknowledgments

Russell Corbett-Detig (UCSC)

Angie Hinrichs (UCSC)

*Idea*  *Implementation*  *Impact*

## UC Santa Cruz

- Bryan Thornlow
- Jakob McBroome
- Alexander Kramer
- Landen Gozashti
- Adriano Schneider
- Cade Mirchandani
- David Haussler

## UC San Diego

- Cheng Ye
- Sumit Walia
- Alireza M.
- Kyle Smith
- Kevin Liu
- Xuan Wang

- Carol Wang
- Devika Torvi
- Shoh Mollenkamp
- Arthur Lu

## ANU

- Robert Lanfear

## EBI/EMBL

- Nicola DeMaio
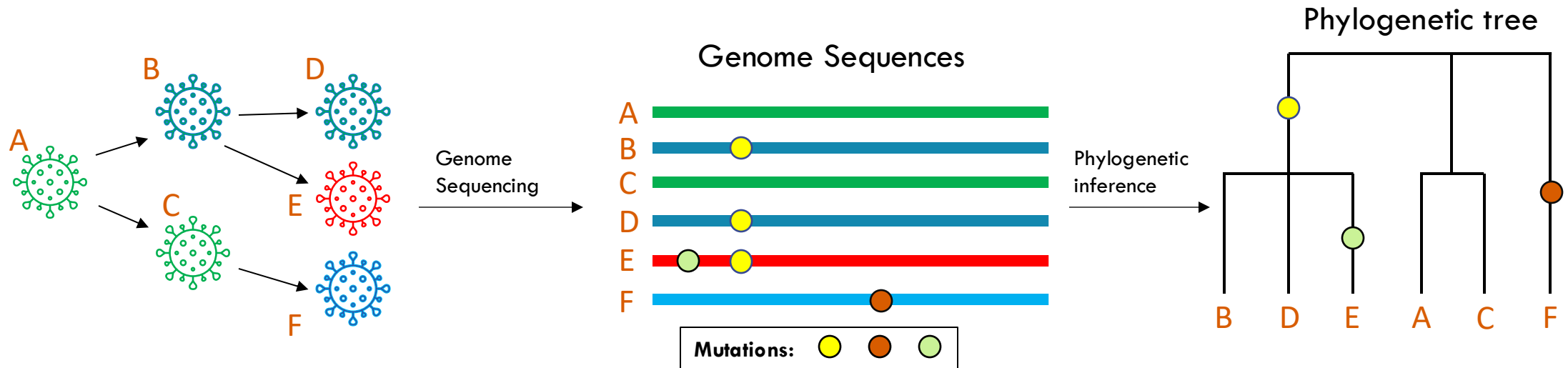- Nick Goldman

## Funding

- CDC
- NIH
- Schmidt Foundation
- UCOP Seed Funding for COVID-19

# Overview of the UShER Package

- **UShER: Phylogenetic placement**

- **matOptimize:** Phylogenetic tree optimization

- **RIPPLES:** Find recombinant sequences using a phylogenomic approach

- **matUtils:** Command-line tools for rapidly analyzing and interpreting SARS-CoV-2 mutation-annotated phylogenetic trees
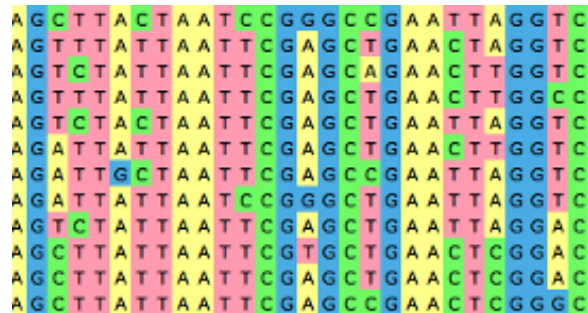
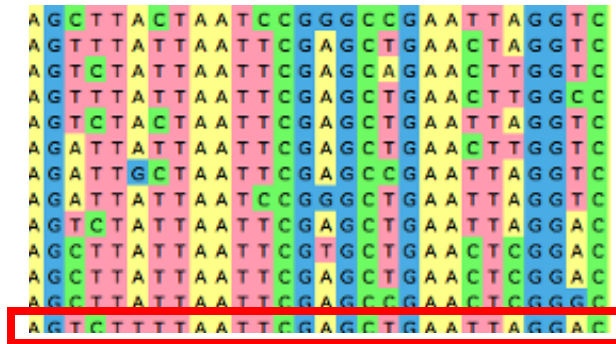# Phylogenetic analysis using genome sequence data

# SARS-CoV-2 phylogenetics with <u>new</u> sequences

**Approach 1: Re-infer** global phylogeny including the new sequences

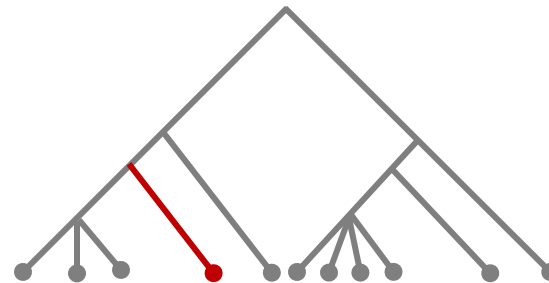1. Gather global sequences (GISAID etc.)

2. Add new sequences to MSA (MAFFT etc.)

Intractable for
SARS-CoV-2 scale!

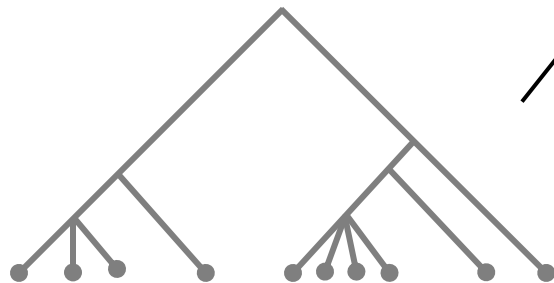3. De novo inference of phylogeny (IQ-Tree etc.)

# SARS-CoV-2 phylogenetics with <u>new</u> sequences

~~**Approach 1: Re-infer** global phylogeny including the new sequences~~

**Approach 2: Place** new sequences on an existing phylogeny



(Placement tool)

| Tool | Time to place 1K samples on 100K tree | Memory required |
|---|---|---|
| IQ-TREE 2 | 6h 9m | 120.2GB |

(using e2-highmem-16 instance)

Ultrafast Sample Placement on Existing Trees

https://github.com/yatisht/usher

**Placing 1K samples on 100K tree**

1439x

1304x

# Parallelizing over multiple CPU instances



Devika Torvi, UCSD **Bioinformatics** undergrad

Kyle Smith, UCSD **Bioinformatics** undergrad

13

# Scaling analysis for placing 100K samples on a 1M-sample tree

**Strong Scaling**

**Weak Scaling**



| UShER | | |
|---|---|---|
| **vCPU** | **Samples placed** | **Time** |
| 64 | 6.25K | 26m 48s |
| 128 | 12.5K | 28m 22s |
| 256 | 25K | 30m 41s |
| 512 | 50K | 33m 36s |
| 1024 | 100K | 37m 07s |

# Why is UShER so fast?

For SARS-CoV-2 relative to the highly popular software, **IQ-TREE,** that has amassed **>10K citations**

# What makes UShER so fast?

1. Choice of **algorithm**: maximum parsimony over maximum likelihood

2. Efficient **data structure**: mutation-annotated tree (MAT)

3. **Pre-processing** for sequential placement

4. Efficiently **parallelizing** the placement step

# 1. Choice of Algorithm: MP over ML (10-100x speedup)

## Maximum Likelihood (ML)



$$L(\theta;D) = \prod_{i=1}^{L} \sum_{x \in X(D)} \pi_{x_r^i} \prod_{(p,c)\in E} P(x_c^i | x_p^i, t_e, \theta) \quad P(D|\theta) = \prod_{1}^{L} P(D^i|\theta)$$

## Maximum Parsimony (MP)



Internal node $u$ with children $v$ and $w$:

$$S_u(x) = \min_y (S_v(y) + W_{xy}) + \min_y (S_w(y) + W_{xy})$$

# MP and ML trees practically the same for SARS-CoV-2



Bryan Thornlow,
UCSC -> ROME Therapeutics

# MP trees are also easier to analyze & interpret!

# 2. Efficient data structure: Mutation-annotated tree (MAT)



| Node | List of Mutations |
|------|-------------------|
| 1 | [G1449U, C7869U, G3179A] |
| 2 | [C9977A] |
| S1 | [C5005U] |
| S2 | [ ] |
| S3 | [ ] |
| S4 | [ ] |
| S5 | [A2869G] |
| S6 | [A6693G] |

174,000 bytes $\longrightarrow$ 186 bytes

"evolutionary compression"

# 3. Pre-processing for sequential placement: ~50x speedup



~89 GB for 3M-seqs

Parallel Fitch-Sankoff (~18 min for 3M-tree)

| Node | List of Mutations |
|------|-------------------|
| 1 | [G1449U, C7869U, G3179A] |
| 2 | [C9977A] |
| S1 | [C5005U] |
| S2 | [ ] |
| S3 | [ ] |
| S4 | [ ] |
| S5 | [A2869G] |
| S6 | [A6693G] |

Place new sequence **S7**

| Node | List of Mutations |
|------|-------------------|
| 1 | [G1449U] |
| 3 | [C7869U, G3179A] |
| 2 | [C9977A] |
| S1 | [C5005U] |
| S2 | [ ] |
| S3 | [ ] |
| S4 | [ ] |
| S5 | [A2869G] |
| **S7** | **[C9977A]** |
| S6 | [A6693G] |

Store as **pre-processed** protobuf file

~200 MB for 3M-tree

Load MAT File (~**20 sec** for 3M-tree)

**~50x speedup**

| Node | List of Mutations |
|------|-------------------|
| 1 | [G1449U, C7869U, G3179A] |
| 2 | [C9977A] |
| S1 | [C5005U] |
| S2 | [ ] |
| S3 | [ ] |
| S4 | [ ] |
| S5 | [A2869G] |
| S6 | [A6693G] |

Place new sequence **S7**

| Node | List of Mutations |
|------|-------------------|
| 1 | [G1449U] |
| 3 | [C7869U, G3179A] |
| 2 | [C9977A] |
| S1 | [C5005U] |
| S2 | [ ] |
| S3 | [ ] |
| S4 | [ ] |
| S5 | [A2869G] |
| **S7** | **[C9977A]** |
| S6 | [A6693G] |

Output MAT

# 4. Efficiently parallelizing the placement step



| Node | List of Mutations |
|------|-------------------|
| 1 | [G1449U, C7869U, G3179A] |
| 2 | [C9977A] |
| S1 | [C5005U] |
| S2 | [ ] |
| S3 | [ ] |
| S4 | [ ] |
| S5 | [A2869G] |
| S6 | [A6693G] |

**S7**      **[G1449U, C9977A]**

| Node | List of Mutations |
|------|-------------------|
| 1 | [G1449U] |
| 3 | [C7869U, G3179A] |
| 2 | [C9977A] |
| S1 | [C5005U] |
| S2 | [ ] |
| S3 | [ ] |
| S4 | [ ] |
| S5 | [A2869G] |
| **S7** | **[C9977A]** |
| S6 | [A6693G] |

# But … (greedy) sequential placement can lead to <u>suboptimal</u> trees occasionally

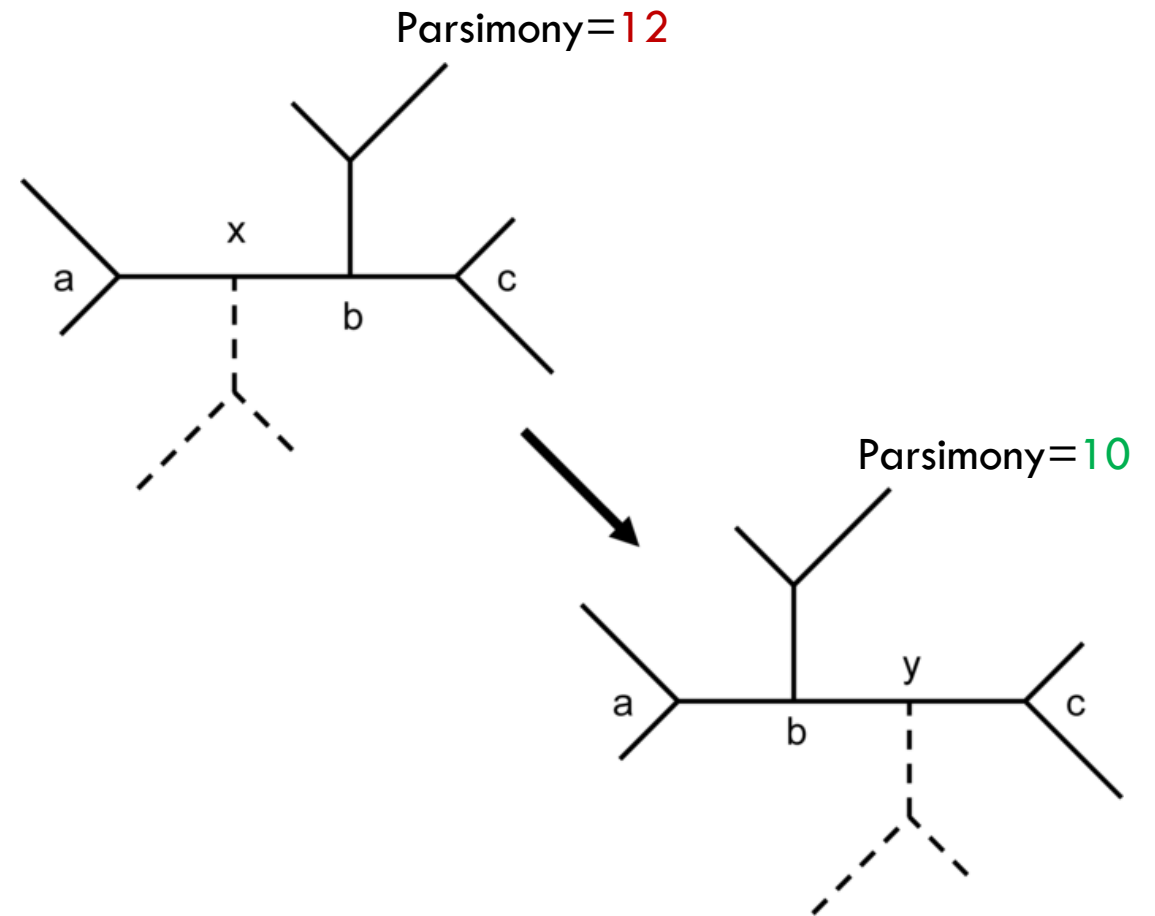And suboptimalities could accumulate!

# Overview of the UShER Package

- UShER: Phylogenetic placement

- **matOptimize: Phylogenetic tree optimization**

- **RIPPLES:** Find recombinant sequences using a phylogenomic approach

- **matUtils:** Command-line tools for rapidly analyzing and interpreting SARS-CoV-2 mutation-annotated phylogenetic trees
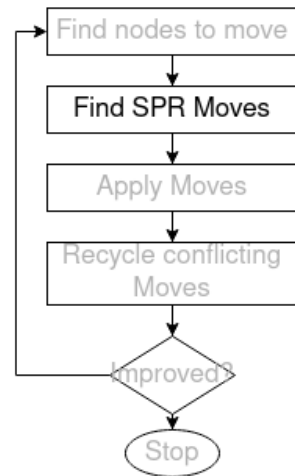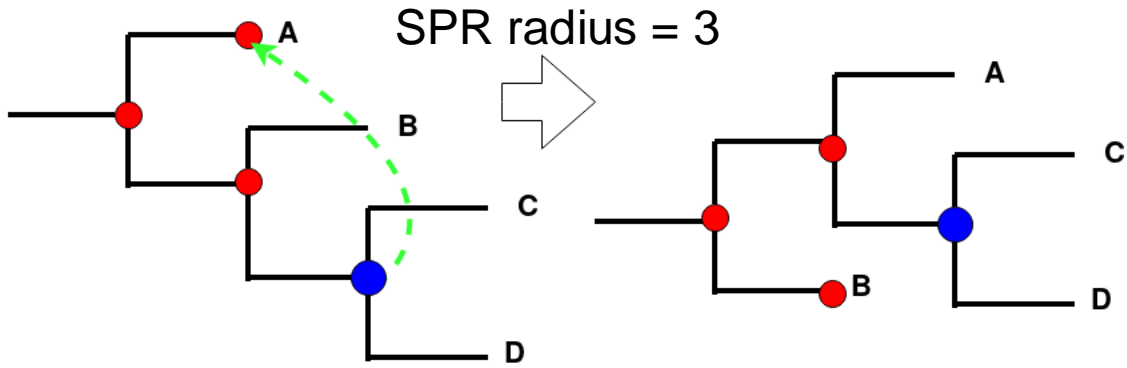
# Tree optimization programs can help ameliorate suboptimal placements

- Use tree re-arrangement (NNI, SPR, TBR etc.)

- Three step process:
  - (Split tree into subtrees/sectors)
  - Identify **profitable** tree rearrangement moves in each sector
  - Apply **non-conflicting** profitable moves in each sector, with a tie-breaking strategy for conflicts
  - (Merge sectors to optimized full tree)

- Some programs maintain several tree copies (applying different moves in each copy)

Parsimony=12

Parsimony=10

# matOptimize (ours) outperforms TNT for SARS-CoV-2

SPR radius = 3



>23x

0.84%

0.75%

Parsimony score improvement

Time (h)

Find nodes to move

Find SPR Moves

Apply Moves

Recycle conflicting Moves

Improved?

Stop

matOptimize
TNT

939GB

21 x

42.9GB

Memory (GB)
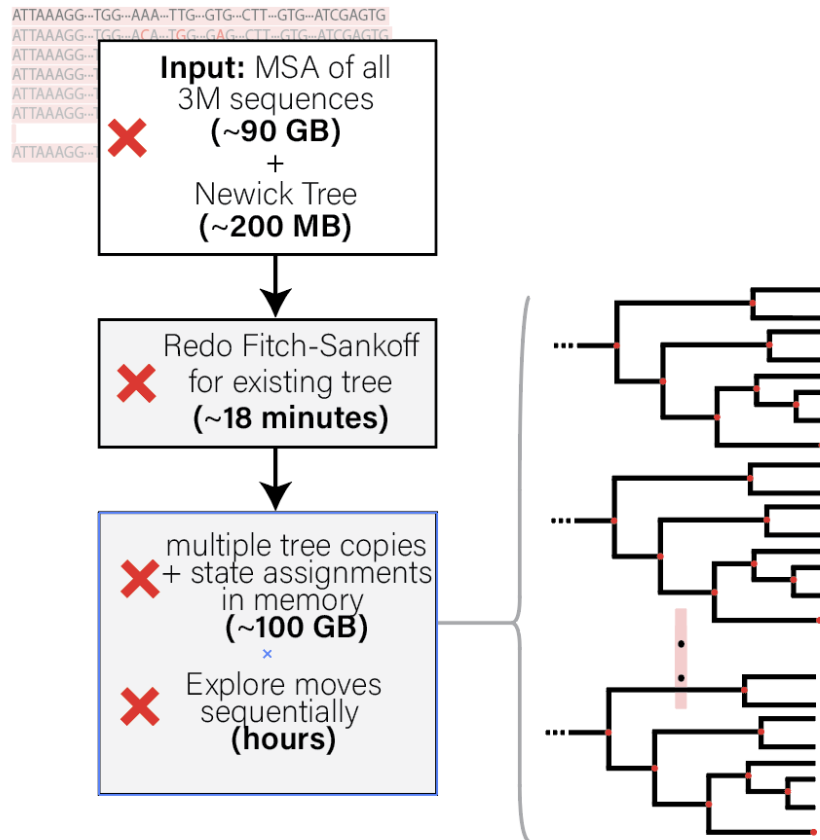
TNT    matOptimize

(Cheng Ye, UCSD ECE undergrad)                                                    26

# Innovative optimizations in matOptimize

## Previous approaches



**Input:** MSA of all 3M sequences **(~90 GB)** + Newick Tree **(~200 MB)** ❌

Redo Fitch-Sankoff for existing tree **(~18 minutes)** ❌

multiple tree copies + state assignments in memory **(~100 GB)** ❌

Explore moves sequentially **(hours)** ❌

## Ours (matOptimize)



**Input:** MAT **(~200 MB)** ✓

MAT loading **(~22 seconds)** ✓

Single MAT copy in memory **(~200 MB)** ✓

Paralellize profitable move search **(minutes)** ✓
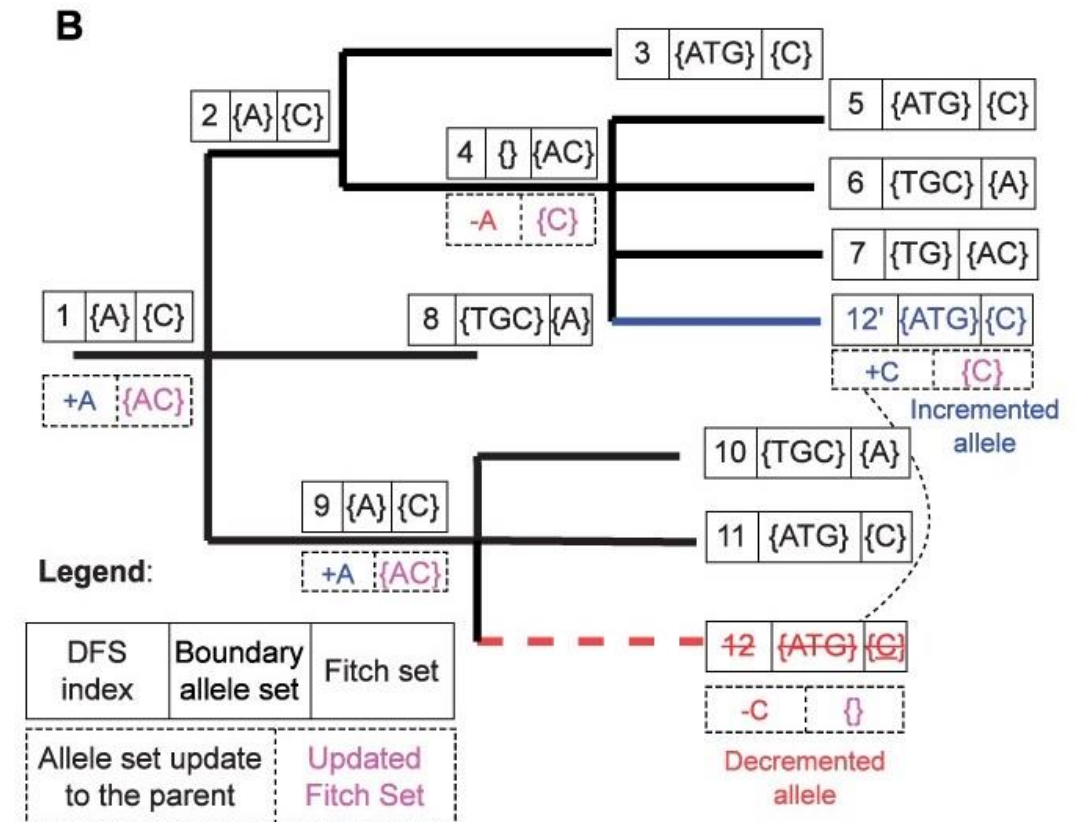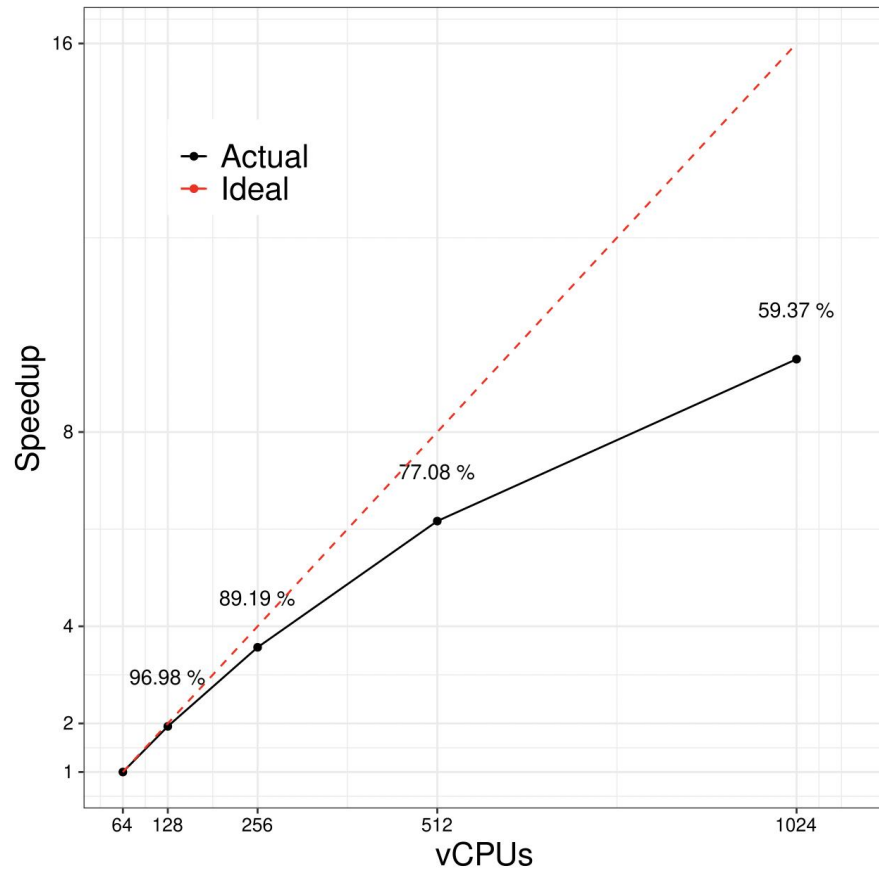
# Innovative optimizations in matOptimize

- More space-efficient and optimization-friendly **MAT format**

- Separate profitable move search and application phase to achieve **high parallelism**
  - Supports multi-node parallelism with MPI

- Modified **Gladstein's incremental update** method to calculate change in parsimony score resulting from a move

- Novel **search space pruning**



(Ye et al. Bioinformatics 2022)

# Multi-node scaling of matOptimize

## Strong Scaling



## Weak Scaling

| vCPU | Source nodes explored | Time |
|------|----------------------|------|
| 64 | 39789 | 10m 45s |
| 128 | 79577 | 11m 54s |
| 256 | 159154 | 11m 51s |
| 512 | 318308 | 11m 58s |
| 1024 | 636616 | 11m 30s |

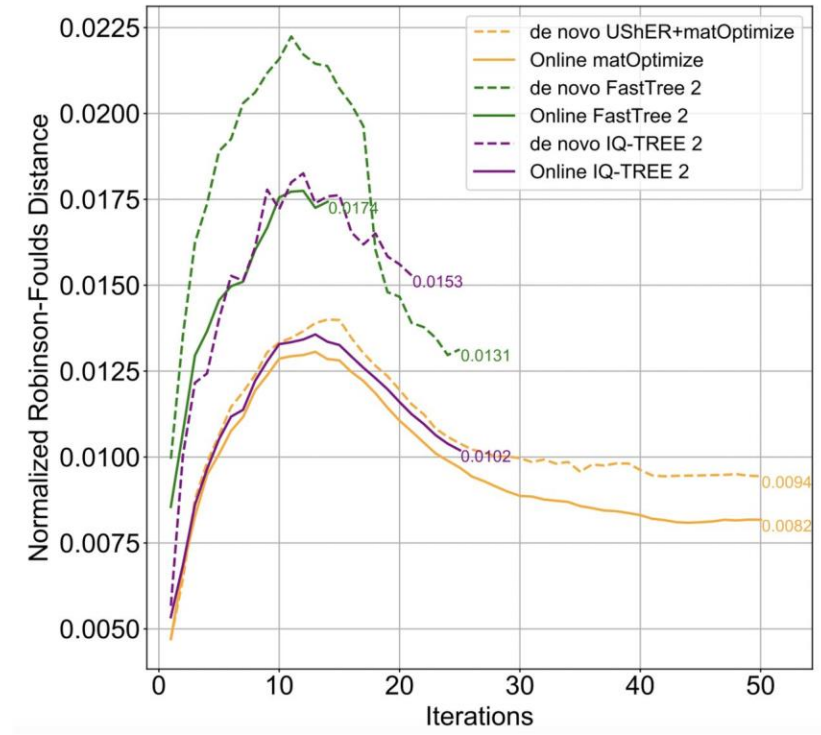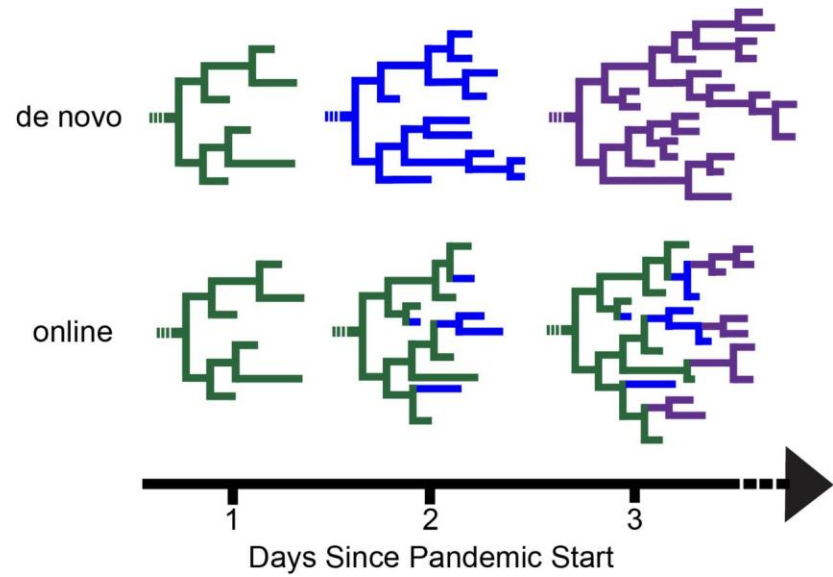# matOptimize: a parallel tree optimization method enables online phylogenetics for SARS-CoV-2 🔓

Cheng Ye, Bryan Thornlow, Angie Hinrichs, Alexander Kramer, Cade Mirchandani,

Devika Torvi, Robert Lanfear, Russell Corbett-Detig, Yatish Turakhia ✉

# Online matOptimize produces phylogenies most similar to ground truth on simulated SARS2 data



Bryan Thornlow

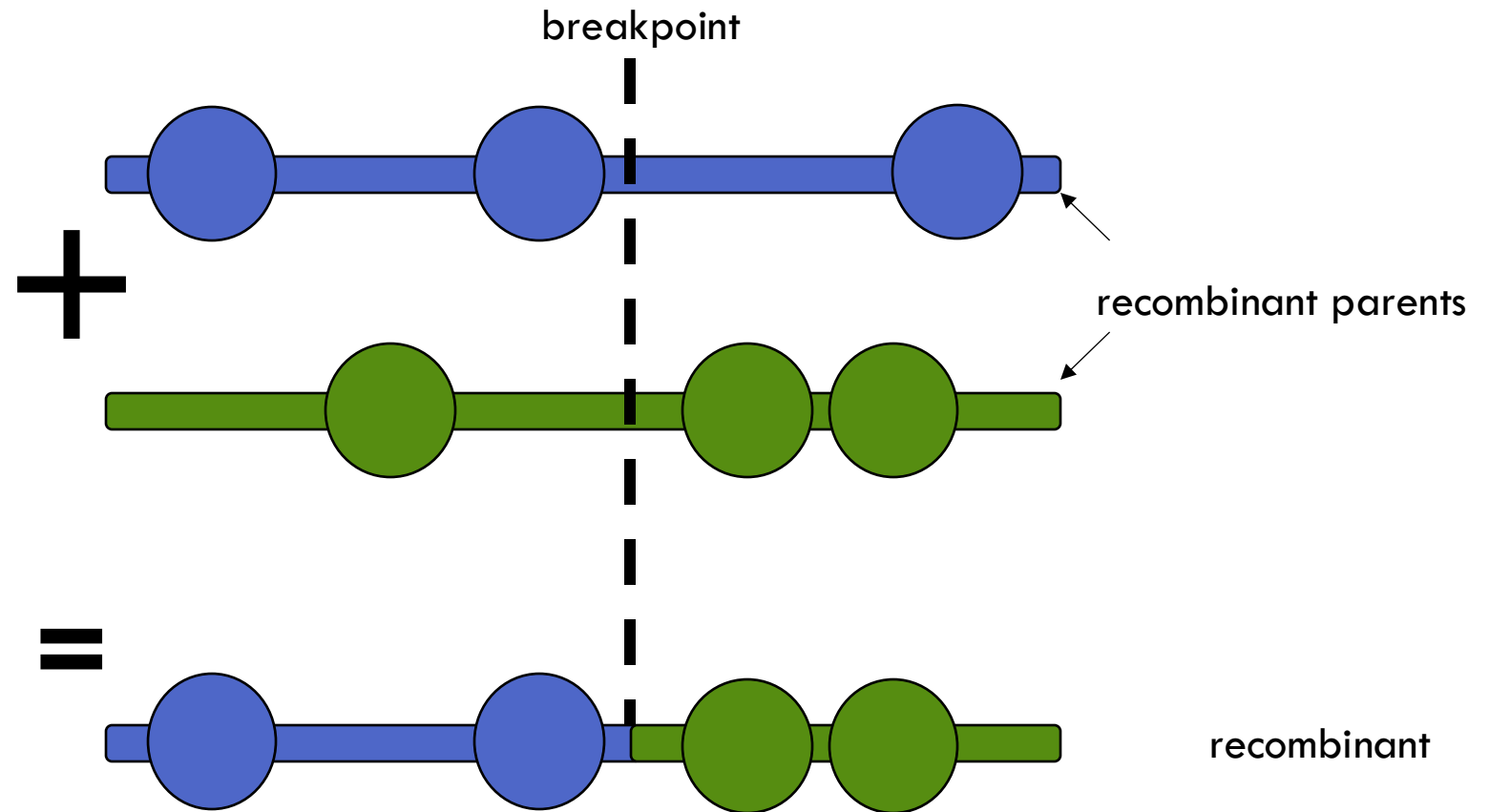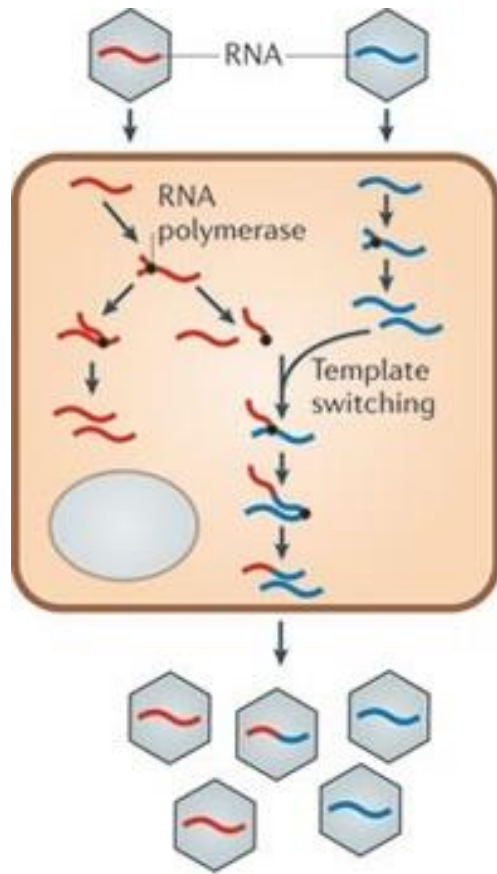# There is another spooky mechanism through which the virus evolves ...

It's called Recombination
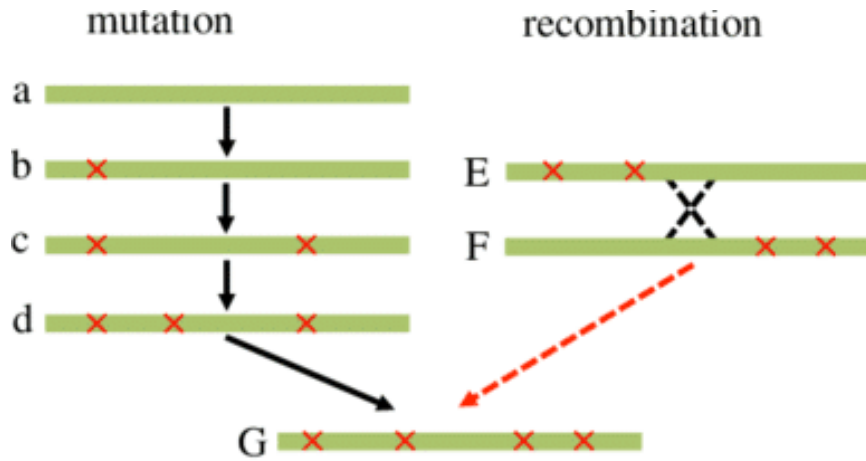
# Overview of the UShER Package

- UShER: Phylogenetic placement

- matOptimize: Phylogenetic tree optimization

- **RIPPLES: Find recombinant sequences using a phylogenomic approach**

- **matUtils:** Command-line tools for rapidly analyzing and interpreting SARS-CoV-2 mutation-annotated phylogenetic trees
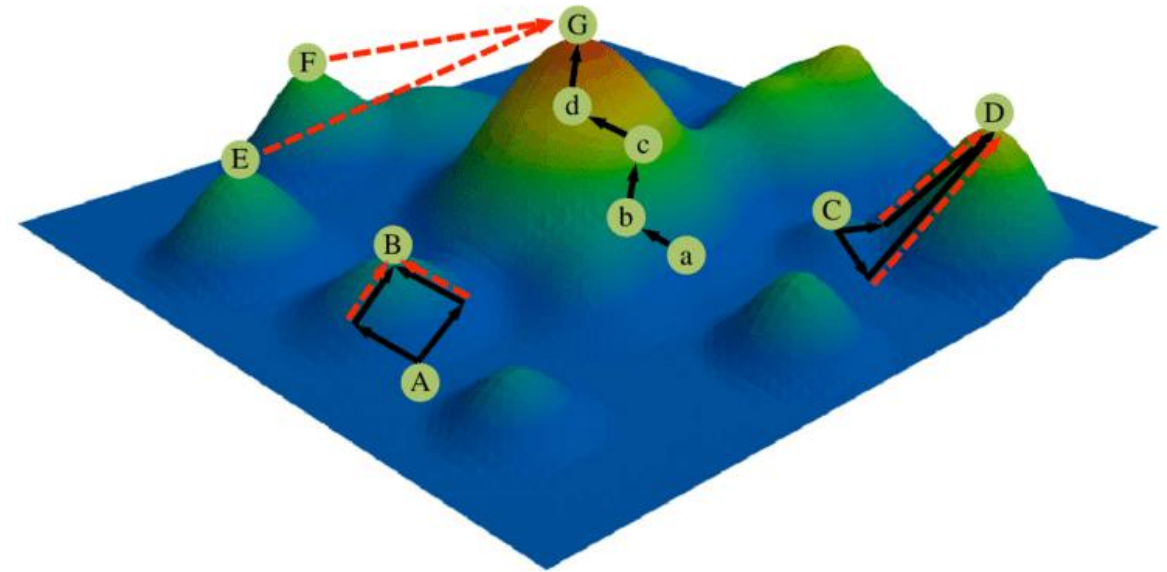
# Viral Recombination

# Recombination may lead to drastic jumps in fitness!

**Sequence evolution by single mutations v/s recombination**

**Fitness landscape**

# Deltacron & some other recombinants are real!

## What is the hybrid 'deltacron' variant of the coronavirus?

Scientists have detected a handful of cases of the delta-omicron hybrid but say it's unlikely to cause a new surge.

### What is deltacron?

Recombinants can emerge when a cell is infected with two different strains of a virus at the same time – in this case, the delta variant and the omicron variant. As the viruses invade the cell and replicate, they can, in rare cases, swap parts of their genome and pick up mutations from each other.
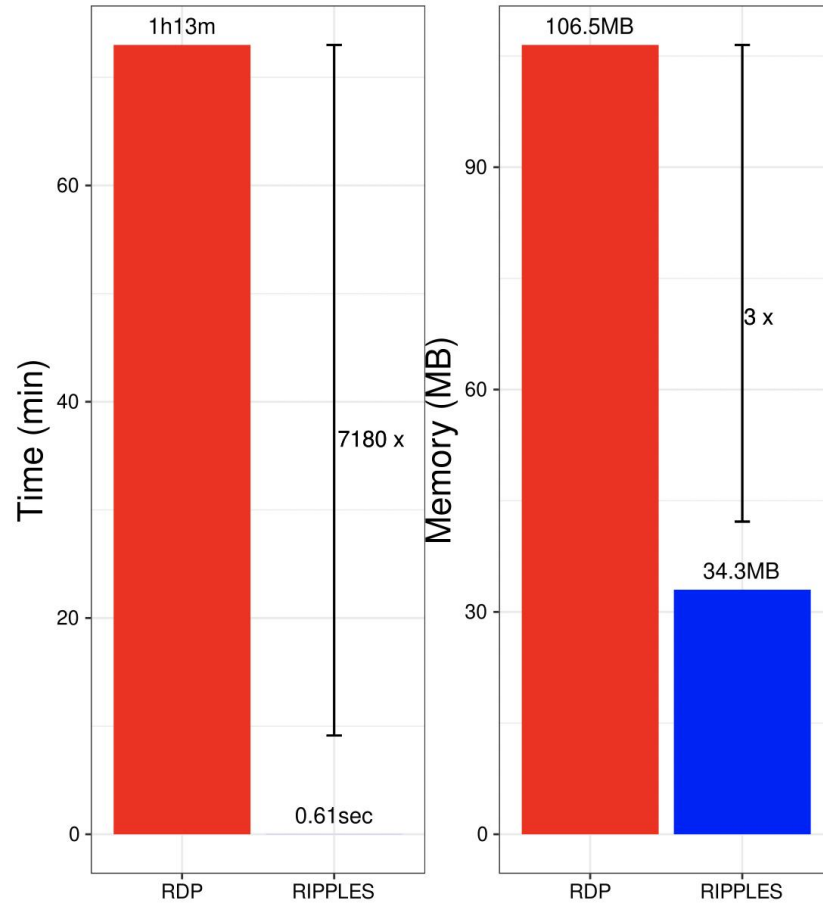
HEALTHCARE • CORONAVIRUS
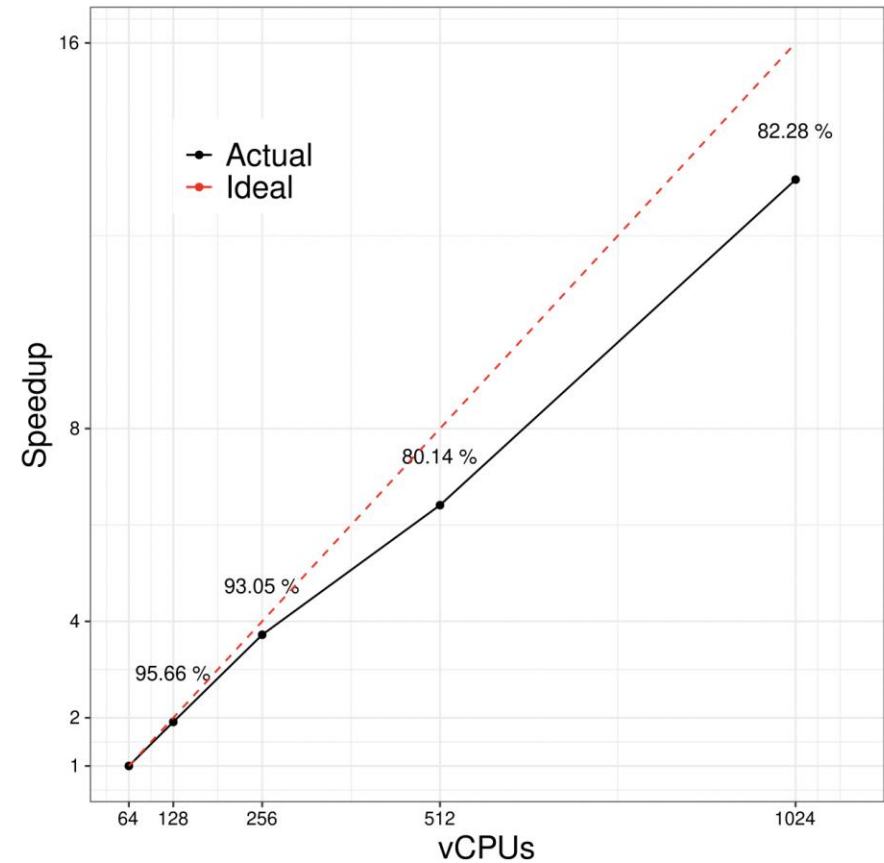
## An Omicron-Omicron Recombinant—BA.4

**William A. Haseltine** Contributor ⓘ

Follow

# RIPPLES orders of magnitude faster using phylogenomic insights



Analyzed on 1K SARS-CoV-2 seqs



~50 min for 1M-sample tree

(Turakhia et al., bioRxiv 2021)

# Innovative optimizations in RIPPLES



**Previous approaches**

**Ours (RIPPLES)**

# RIPPLES discovers >600 recombination events!

- **>600 unique** SARS-CoV-2 recombination events discovered from a 1.6M sample phylogeny!
  - ~2.7% sequences have a **detectable** recombinant ancestry

- This is the **largest recombination study** to our knowledge

- With our latest RIPPLES software (`ripples-fast`), we can infer recombinants from a **~10M SARS-CoV-2 mutation-annotated tree in ~2 hours!**



Cheng Ye

# Recombination breakpoints are elevated in the SARS-CoV-2 Spike protein region

# Overview of the UShER Package

- UShER: Phylogenetic placement

- matOptimize: Phylogenetic tree optimization

- RIPPLES: Find recombinant sequences using a phylogenomic approach
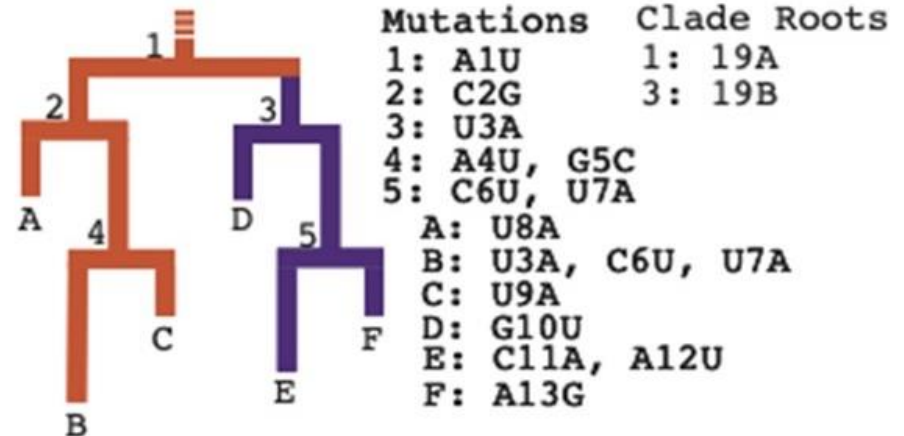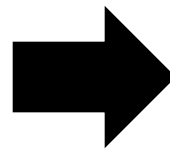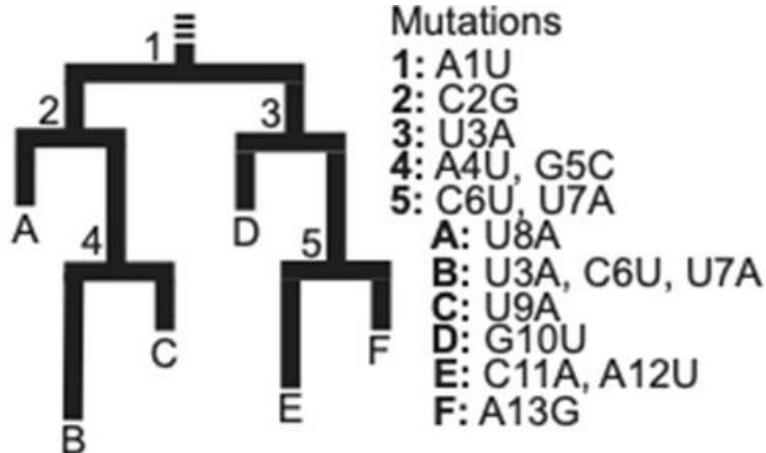
- **matUtils: Command-line tools for rapidly analyzing and interpreting SARS-CoV-2 mutation-annotated phylogenetic trees**

# matUtils: Toolkit to rapidly query and interpret MATs

| Subcommand | Arguments | Description |
|---|---|---|
| **annotate** | **--clade-to-nid (1: 19A, 3: 19B)** | **Assigns clades to nodes** |
| extract | --clade 19B | Extract subtree based on clade, mutation, branch length and other conditions |
| uncertainty | --find-epps | Output sample placement uncertainty metrics, for e.g., number of equally parsimonious placements |
| introduce | --population-samples (D: USA, E: USA, F:USA) | Identify internal nodes corresponding to introduction of one or more infection clusters in a geographic region of interest |
| summary | N/A | Output basic statistics of a MAT |



(McBroome et al., MBE 2021)

# matUtils: Toolkit to rapidly query and interpret MATs

| Subcommand | Arguments | Description |
|---|---|---|
| annotate | --clade-to-nid (1: 19A, 3: 19B) | Assigns clades to nodes |
| **extract** | **--clade 19B** | **Extract subtree based on clade, mutation, branch length and other conditions** |
| uncertainty | --find-epps | Output sample placement uncertainty metrics, for e.g., number of equally parsimonious placements |
| introduce | --population-samples (D: USA, E: USA, F:USA) | Identify internal nodes corresponding to introduction of one or more infection clusters in a geographic region of interest |
| summary | N/A | Output basic statistics of a MAT |



(McBroome et al., MBE 2021)

# matUtils: Toolkit to rapidly query and interpret MATs

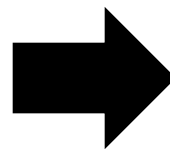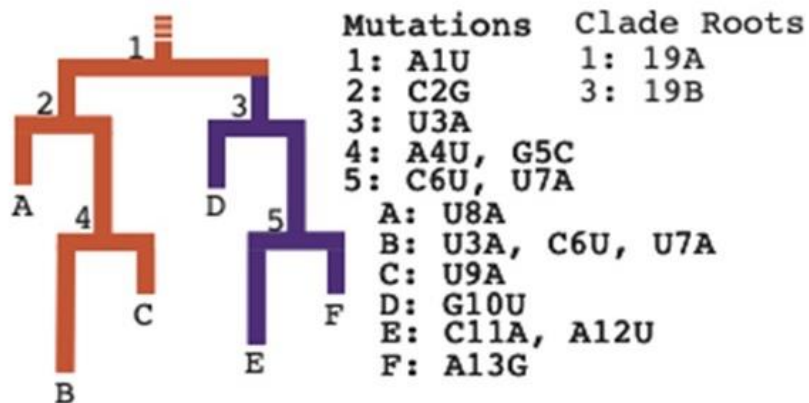| Subcommand | Arguments | Description |
|---|---|---|
| annotate | --clade-to-nid (1: 19A, 3: 19B) | Assigns clades to nodes |
| extract | --clade 19B | Extract subtree based on clade, mutation, branch length and other conditions |
| **uncertainty** | **--find-epps** | **Output sample placement uncertainty metrics, for e.g., number of equally parsimonious placements** |
| introduce | --population-samples (D: USA, E: USA, F:USA) | Identify internal nodes corresponding to introduction of one or more infection clusters in a geographic region of interest |
| summary | N/A | Output basic statistics of a MAT |

# matUtils: Toolkit to rapidly query and interpret MATs

| Subcommand | Arguments | Description |
|---|---|---|
| annotate | --clade-to-nid (1: 19A, 3: 19B) | Assigns clades to nodes |
| extract | --clade 19B | Extract subtree based on clade, mutation, branch length and other conditions |
| uncertainty | --find-epps | Output sample placement uncertainty metrics, for e.g., number of equally parsimonious placements |
| **introduce** | **--population-samples (D: USA, E: USA, F:USA)** | **Identify internal nodes corresponding to introduction of one or more infection clusters in a geographic region of interest** |
| summary | N/A | Output basic statistics of a MAT |

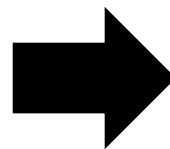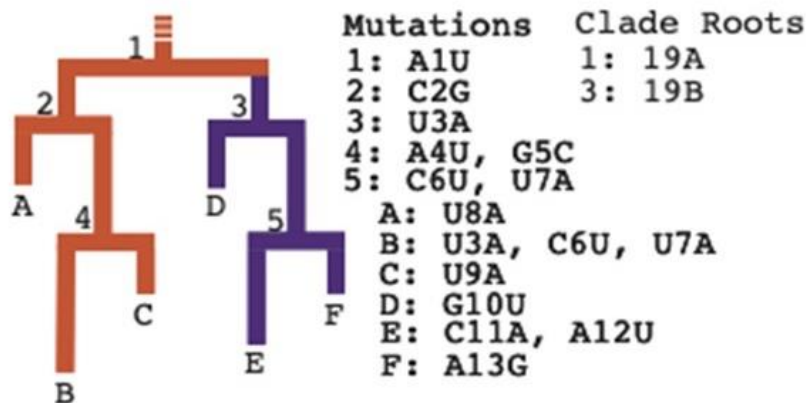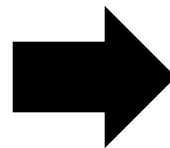# matUtils: Toolkit to rapidly query and interpret MATs

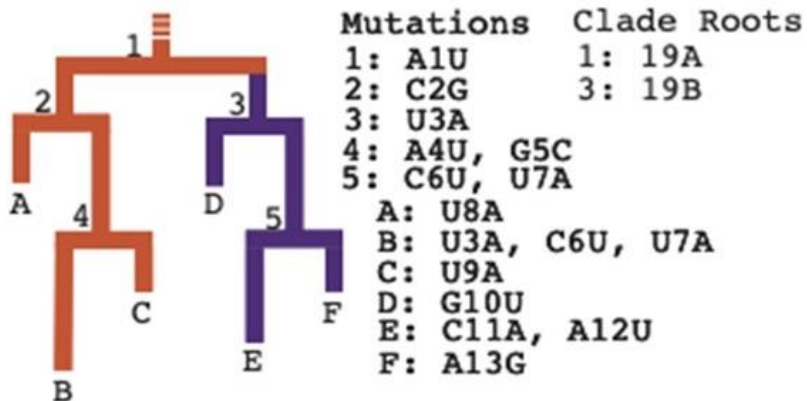| Subcommand | Arguments | Description |
|---|---|---|
| annotate | --clade-to-nid (1: 19A, 3: 19B) | Assigns clades to nodes |
| extract | --clade 19B | Extract subtree based on clade, mutation, branch length and other conditions |
| uncertainty | --find-epps | Output sample placement uncertainty metrics, for e.g., number of equally parsimonious placements |
| introduce | --population-samples (D: USA, E: USA, F:USA) | Identify internal nodes corresponding to introduction of one or more infection clusters in a geographic region of interest |
| **summary** | **N/A** | **Output basic statistics of a MAT** |

# matUtils is fast!

On **~1M-sample** SARS-CoV-2 phylogeny:

- ~5 sec to compute **summary** statistics

- ~5 sec to **extract** a subtree of specified samples

- ~15 sec to **extract** mutation paths from root to every sample in the tree

Jacob McBroome, UCSC

- ~9 sec to **resolve** all polytomies
  - Takes ~37 min using `ape`

- ~1 min to identify **introductions**

# Real-world impact

# UShER added to UCSC Genome Browser



Angie Hinrichs, UCSC

Rapid cross-referencing with 50+ molecular and structural biology tracks

# YouTube Tutorials on the UCSC SARS-CoV-2 Genome Browser



**SARS-CoV-2 Genome Browser**

5 videos • 65 views • Last updated on Mar 22, 2022

A fast-paced introduction to using the SARS-CoV-2 Genome Browser to look at Genome Annotation Tracks and the UShER web interface to make phylogenetic trees of your samples and compare with millions of SARS-CoV-2 sequences

**UCSC Genome Browser**

SUBSCRIBE

1. **Introduction to the UCSC SARS-CoV-2 Genome Browser**
   UCSC Genome Browser
   5:03

2. **1A - Uploading a sequence file to the UCSC SARS-CoV-2 Genome Browser**
   UCSC Genome Browser
   6:23

3. **1B - Using UShER custom tracks with the UCSC SARS-CoV-2 Genome Browser**
   UCSC Genome Browser
   3:42

4. **2A - Exploring additional data tracks in the UCSC SARS-CoV-2 Genome Browser**
   UCSC Genome Browser
   6:26

5. **2B - Uploading custom data tracks into the UCSC SARS-CoV-2 Genome Browser**
   UCSC Genome Browser
   4:46

https://www.youtube.com/playlist?list=PL5G1tMPaNuswZM4zAkmS6F7q1G-7D-Pl-

# Largest SARS-CoV-2 phylogenies

**Covered by UCSC news and
Santa Cruz Tech Beat**

**Taxonium view of >10M SARS-CoV-2 sample phylogeny**



UCSC's Million-COVID-Genome Tree Could be a First

ngie Hinrichs, Senior Software Architect, UCSC Genome Browser

Yatish Turakhia, UCSC Postdoc scholar, incoming Assistant Professor, UCSD

Russell Corbett-Detig, UCSC Assistant Professor, Biomolecular Engineering

Solving a computational puzzle, a UCSC team created a dynamic evolutionary tree to enable real-time genomic contact tracing

SANTA CRUZ, CA – April 13, 2020 – Early in the pandemic, UCSC knew they wanted to help researchers tracking the virus. During the 2013 Ebola outbreak, the seasoned Browser team had used their coding skills to build a virus browser. Since Ebola, a new era of fast, cheap sequencing has created a mountain of genomic data, changing the research landscape. Traditional display code just wouldn't keep pace with this novel coronavirus.

Theo Sanderson,
Francis Crick Institute

# UShER default engine in Pangolin SARS-CoV-2 lineage assigner

## pangolin v3.0

aineniamh released this on May 27 · 110 commits to master since this release

### Release notes: Major release

- pangolin 3.0 comes with additional functionality and requires an environment update as extra dependencies now include UShER and scorpio.
- Requires pangoLEARN data >= 2021-05-27

### Lineage assignment updates

- PANGO assignment uses a sequence hash from all currently designated sequences to assign lineages.
- PLEARN (pangoLEARN) assignment using a machine-learning model to assign the most likely lineage. Current model is decision tree model.
- PUSHER (pangolin-UShER, pangUSHlin, pangUShER) assignment uses fast parsimony placement of a query sequence into a protobuf file based on currently designated sequences and infers most pasimonious lineage based on this placement (thanks to @AngieHinrichs and the rest of the UShER team).
- Default --max-ambig value changed to 0.3.

## pangolin v4.0

aineniamh released this Apr 01, 2022    · 51 commits to master since this release    ⬡ v4.0    ⊶ 20eb73e ✓

Compare ▾

### Release notes

pangolin has had a big code overhaul recently, which should help with maintainability going forward, but there are some main changes the user will be concerned with that I wanted to flag here before the release:

- Notably, the default mode is shifting from pangoLEARN to UShER. If you run large amounts of sequences through pangolin routinely you should be aware this update will impact the speed of pangolin for large amounts of data and you may want to consider parallelisation, using the optional usher assignment cache file (accessed with --add-assignment-cache and --use-assignment-cache flags) or using the --analysis-mode pangoLEARN flag.
- The pangoLEARN model being trained is a random forest rather than a decision tree, so the confidence scores reflect the assignment probability from the random forest model now rather than the number of suitable categories as is the case with the decision tree model.
- Changes to dependencies: We're rationalising the pangoLEARN repository and the file accessed from pango-designation into a single repository called pangolin-data, so pangoLEARN and pango-designation are no longer needed as dependencies.
- Changes to versioning: pangolin-data will have the same version number as the pango-designation tag as the lineages version in UShER protobuf file and the pangoLEARN model, giving a less convoluted versioning system than has previously been the case.

Thanks to:
- Áine O'Toole
- Emily Scher
- Rachel Colquhoun
- Andrew Rambaut

Angie Hinrichs, UCSC

# UShER included in the CDC COVID-19 Genomic Epidemiology Toolkit



## Part 3: Implementation
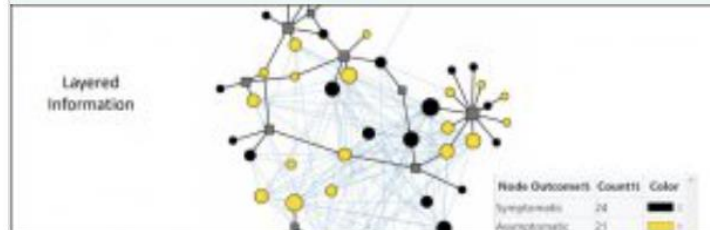
**Module 3.1**

**Getting started with Nextstrain**

Introducing Nextstrain, an interactive tool for visualizing phylogenetic trees
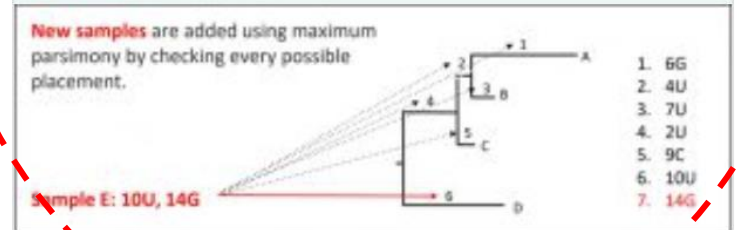
**Module 3.2**

**Getting started with MicrobeTrace**

Introducing MicrobeTrace, an interactive tool for transmission network analysis

**Module 3.3**

**Real-time phylogenetics with UShER**

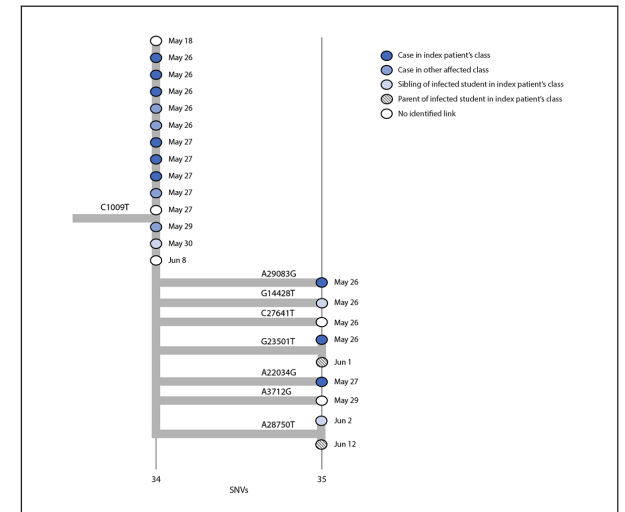Introducing UShER, a web portal for fast calculation of phylogenetic trees

# Delta outbreak analysis at a Calif. elementary school used UShER



The New York Times

## How the Delta Variant Infiltrated an Elementary School Classroom

A detailed study in California found that the variant easily spread from an unvaccinated teacher to children and, in a few cases,

# UShER has been used for outbreak analysis in several parts of the world (including LMIC)

Repeated transmission of SARS-CoV-2 in an overcrowded Irish emergency department elucidated by whole genome sequencing

Daniel Hare [1,2,3], Carolyn Meaney [1], James Powell [1,4], Barbara Slevin [5], Breda O' Brien [5], Lorraine Power [1], Nuala H. O' Connell [1,2,4], Cillian F. De Gascun [3], Colum P. Dunne [2,4], Patrick J. Stapleton [1,2]

Emergence of a novel SARS-CoV-2 Pango lineage B.1.1.526 in West Bengal, India

Rakesh Sarkar [a], Ritubrita Saha [a], Pratik Mallick [b], Ranjana Sharma [a], Amandeep Kaur [c], Shanta Dutta [a], Mamta Chawla-Sarkar [a]

## Assessment of Inter-Laboratory Differences in SARS-CoV-2 Consensus Genome Assemblies between Public Health Laboratories in Australia

by Charles S. P. Foster [1,2,*], Sacha Stelzer-Braid [1,2], Ira W. Deveson [3,4], Rowena A. Bull [2,5], Malinna Yeang [1,2], Jane-Phan Au [1,2], Mariana Ruiz Silva [1,2], Sebastiaan J. van Hal [6,7], Rebecca J. Rockett [8,9], Vitali Sintchenko [8,9,10,11], Ki Wook Kim [1,12] and William D. Rawlinson [1,2,12,13]

## Imported SARS-CoV-2 Variants of Concern Drove Spread of Infections across Kenya during the Second Year of the Pandemic

by Carolyne Nasimiyu [1,2,†], Damaris Matoke-Muhia [3,†], Gilbert K. Rono [4,†], Eric Osoro [1,2], Daniel O. Ouso [4], J. Milkah Mwangi [3], Nicholas Mwikwabe [3], Kelvin Thiong'o [3], Jeanette Dawa [1], Isaac Ngere [1], John Gachohi [1,5], Samuel Kariuki [3], Evans Amukoye [3], Marianne Mureithi [6], Philip Ngere [7], Patrick Amoth [7], Ian Were [7], Lyndah Makayotto [7], Vishvanath Nene [4], Edward O. Abworo [4], + Show full author list

Short communication

The dynamic change of SARS-CoV-2 variants in Sierra Leone

Lei Lin [a,1], Juling Zhang [b,1], James Rogers [c], Allan Campbell [d], Jianjun Zhao [a], Doris Harding [d], Foday Sahr [c], Yongjian Liu [a], Isata Wurie [c]

## The rise and spread of the SARS-CoV-2 AY.122 lineage in Russia

Galya V Klink, Ksenia R Safina, Elena Nabieva, Nikita Shvyrev, Sofya Garushyants, Evgeniia Alekseeva, Andrey B Komissarov, Daria M Danilenko, Andrei A Pochtovyi, Elizaveta V Divisenko, Lyudmila A Vasilchenko, Elena V Shidlovskaya, Nadezhda A Kuznetsova, The Coronavirus Russian Genetics Initiative (CoRGI) Consortium, Anna S Speranskaya, Andrei E Samoilov, Alexey D Neverov, Anfisa V Popova, Gennady G Fedonin, The CRIE Consortium, Vasiliy G Akimkin, Dmitry Lioznov, Vladimir A Gushchin, Vladimir Shchur, Georgii A Bazykin
    Author Notes

## Three SARS-CoV-2 recombinants identified in Brazilian children

Luciane Sussuchi da Silva
    Dasa

# First flagging of Omicron (B.1.1.529) variant

thomasppeacock commented on Nov 23, 2021 · edited by chrisruis · · ·

**New proposed lineage**
By Tom Peacock

**Description**
**Sub-lineage of:** B.1.1
Earliest Sequence: 2021-11-11
Latest Sequence: 2021-11-13

Countries circulating: Botswana (3 genomes), Hong Kong ex S. Africa (1 genome, partial)

Description:
**Conserved Spike mutations -** A67V, Δ69-70, T95I, G142D/Δ143-145, Δ211/L212I, ins214EPE, G339D, S371L, S373P, S375F, K417N, N440K, G446S, S477N, T478K, E484A, Q493R, G496S, Q498R, N501Y, Y505H, T547K, D614G, H655Y, N679K, P681H, N764K, D796Y, N856K, Q954H, N969K, L981F

**Conserved non-Spike mutations -** NSP3 – K38R, V1069I, Δ1265/L1266I, A1892T; NSP4 – T492I; NSP5 – P132H; NSP6 – Δ105-107, A189V; NSP12 – P323L; NSP14 – I42V; E – T9I; M – D3G, Q19E, A63T; N – P13L, Δ31-33, R203K, G204R

Currently only 4 sequences so would recommend monitoring for now. Export to Asia implies this might be more widespread than sequences alone would imply. Also the extremely long branch length and incredibly high amount of spike mutations suggest this could be of real concern (predicted escape from most known monoclonal antibodies)

**Genomes:**
EPI_ISL_6590608 (partial RBD Sanger sequencing from Hong Kong)
EPI_ISL_6640916
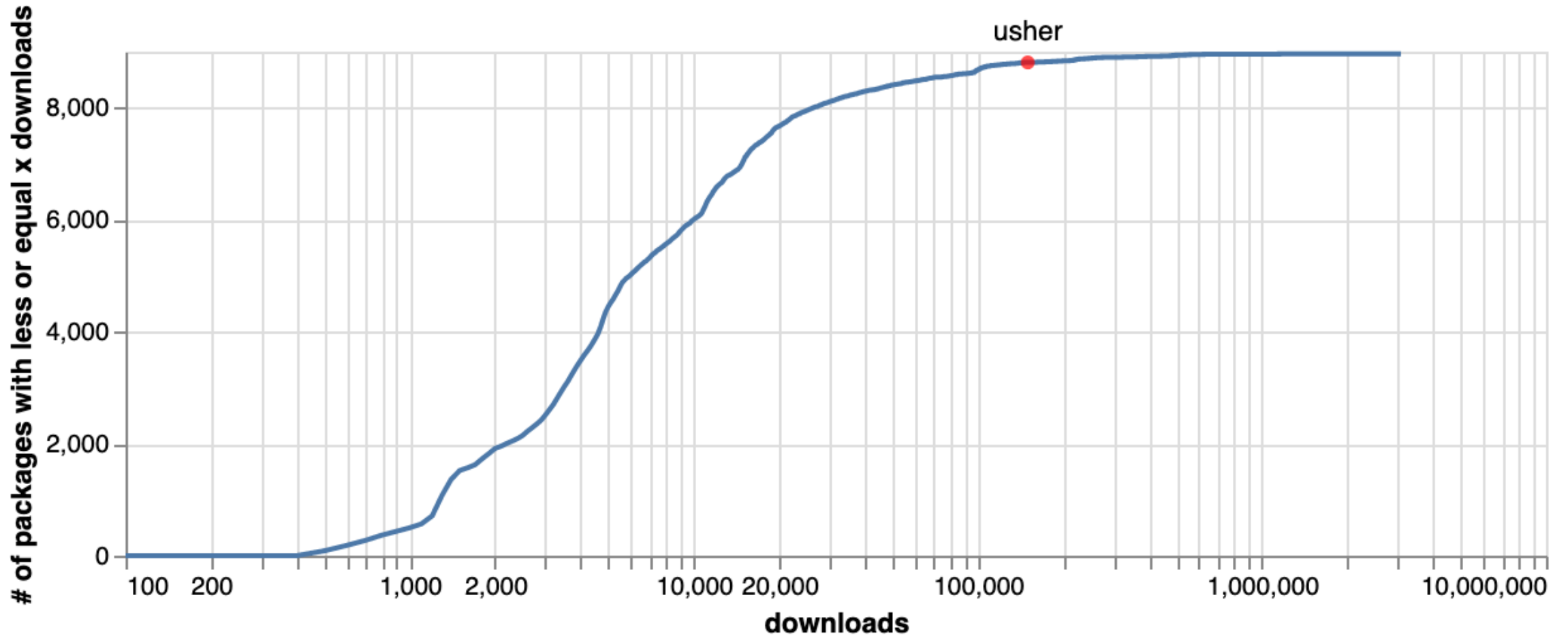EPI_ISL_6640919
EPI_ISL_6640917

Tom Peacock, Imperial College

**UShER-based phylogenetic analysis with the first Omicron sequences in red**

# 150K+ downloads on bioconda

# Press coverage

**Medical Xpress**

MAY 10, 2021

**New tools enable rapid analysis of coronavirus sequences and tracking of variants**

**TheScientist**
EXPLORING LIFE, INSPIRING INNOVATION

**Plenty of Evidence for Recombination in SARS-CoV-2**

Different variants of the virus behind the COVID-19 pandemic are swapping chunks of genetic material, but it's not yet clear what implications that may have for public health.

**genomeweb**

**Nature Papers Present Nautilus Genome, Tool to Analyze Single-Cell Data, More**

May 13, 2021

UC SANTA CRUZ    **The team behind a tree of 10 million Covid sequences**

June 21, 2022
By Rose Miyatsu

**NEWSCENTER**

**Medical Xpress**

JUNE 23, 2022

**New phylogenetic tool can handle the SARS-CoV-2 data load**

by Kiran Kumar and Katherine Connor, University of California - San Diego

**WIRED**

**How the Delta variant took over**

The Delta variant accounts for 90 per cent of new Covid-19 cases in the UK. Scientists fear its global spread is going unchecked

**NEWS MEDICAL LIFE SCIENCES**

**New phylogenomic platform identifies increased recombination rates in SARS-CoV-2**

UC SANTA CRUZ

Home / 2022 / April / Genomics Institute tool becomes primary method to identify lineages of COVID-19 worldwide

**Genomics Institute tool becomes primary method to identify lineages of COVID-19 worldwide**

**NEWSCENTER**

Widespread use of the "UShER" tool will enable public health officials to more accurately identify and track the virus's variants

April 04, 2022
By Emily Cerf

# Community response

# TURAKHIA LAB

# We are hiring at all levels!

- Computational Genomics / Bioinformatics
- High-performance computing (HPC)
- GPU / FPGA computing
- Computer Architecture
- Hardware-Software Co-design
- VLSI Design





Webpage: https://turakhia.eng.ucsd.edu/
Email: yturakhia@ucsd.edu