# Compressive Pangenomics using PanMANs

Prof. Yatish Turakhia

Department of ECE, UCSD

# Codebase and Preprint



Pangenome Mutation-Annotated Network

https://github.com/TurakhiaLab/panman

**Compressive Pangenomics Using Mutation-Annotated Networks**

Sumit Walia, Harsh Motwani, Kyle Smith, Russell Corbett-Detig, Yatish Turakhia

**doi:** https://doi.org/10.1101/2024.07.02.601807

This article is a preprint and has not been certified by peer review [what does this mean?].

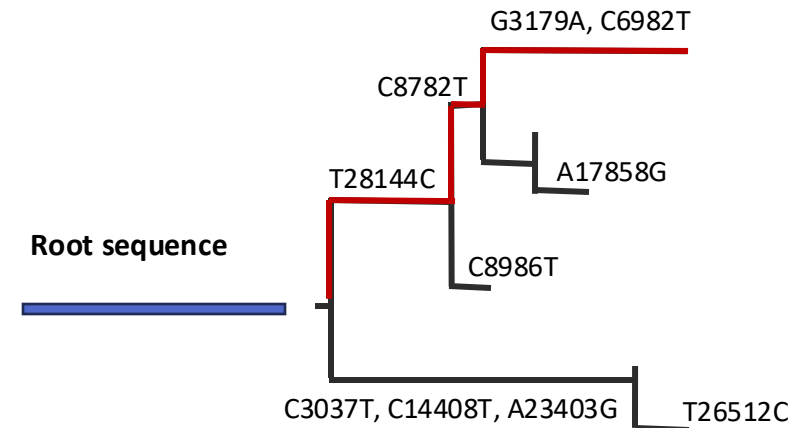https://www.biorxiv.org/content/10.1101/2024.07.02.601807v1
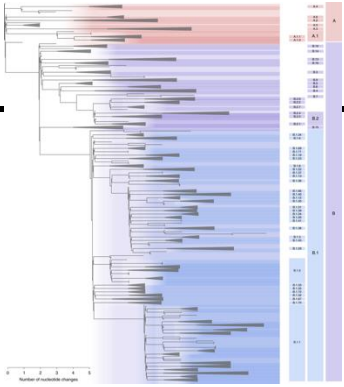
Sumit Walia

Harsh Motwani

# MAT: The data structure powering UShER

- **MAT**: mutation-annotated tree

- Stores:
  - **Tree topology** corresponding to the inferred phylogeny
  - A single **root sequence** (could be the reference genome)
  - **Mutations** inferred on each branch

- **Property:** sequence corresponding to every tip or internal node of the tree can be derived from the root sequence and the mutations on its path to the root
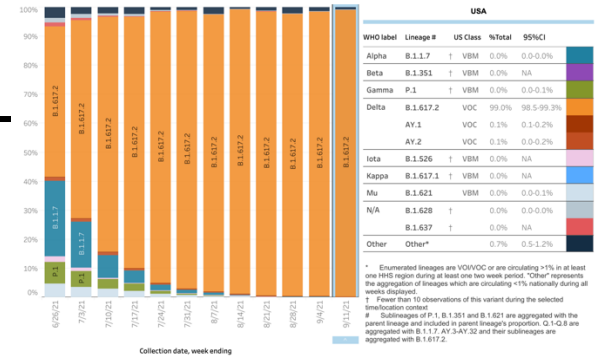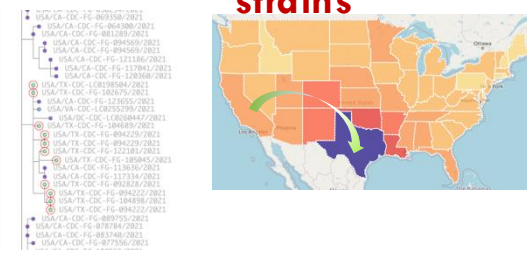
# Naming lineages



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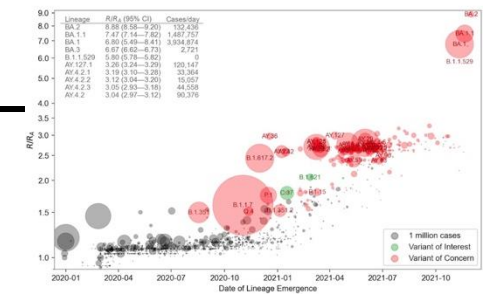
# Monitoring circulating lineages



(CDC.gov dashboard 2021)

# Identify newly-introduced strains
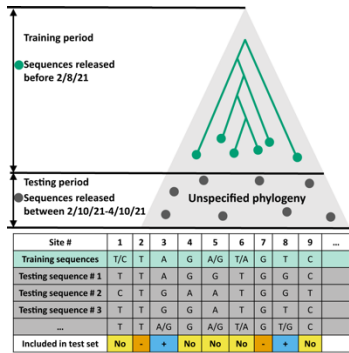


(McBroome et al., Virus Evol. 2022)

# Predicting fitness of a new strain



(Obermeyer et al., Science 2022)
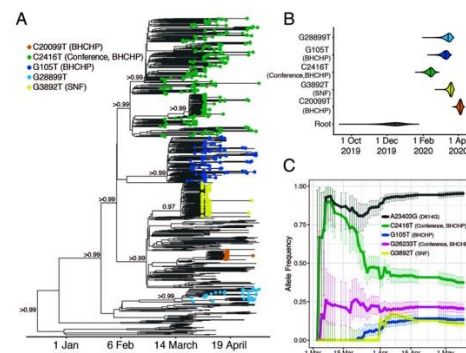


## UShER

# Predicting the next mutation



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

# Analyze outbreaks and superspreader events
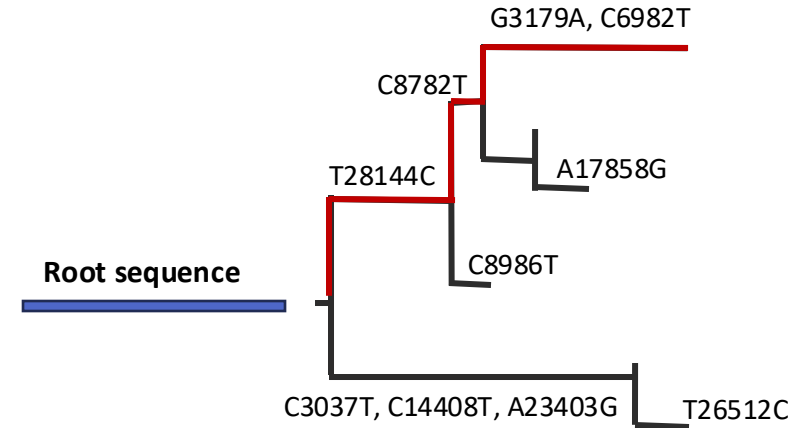


(Lemieux et al., Science 2021)

# Wastewater surveillance



(Karthikeyan et al., Nature 2022)

# Limitations of UShER-MAT

- **Reference-based**



S1 : AG**A**T**GC**T
S2 : **T**GCT**GC**T

⟶

```
0 1 2 3 4 5 6
AGCTATT
```

# Limitations of UShER-MAT

- **Reference-based**

- Only stores substitutions – **ignores indels**
  - Indels sometimes comprise lineage-defining mutations

# Limitations of UShER-MAT

- **Reference-based**

- Only stores substitutions – **ignores indels**
    - Indels sometimes comprise lineage-defining mutations

- Restricted to a single **tree topology** – cannot represent complex mutations (e.g., recombination or horizontal gene transfer) violating the vertical mode of evolution

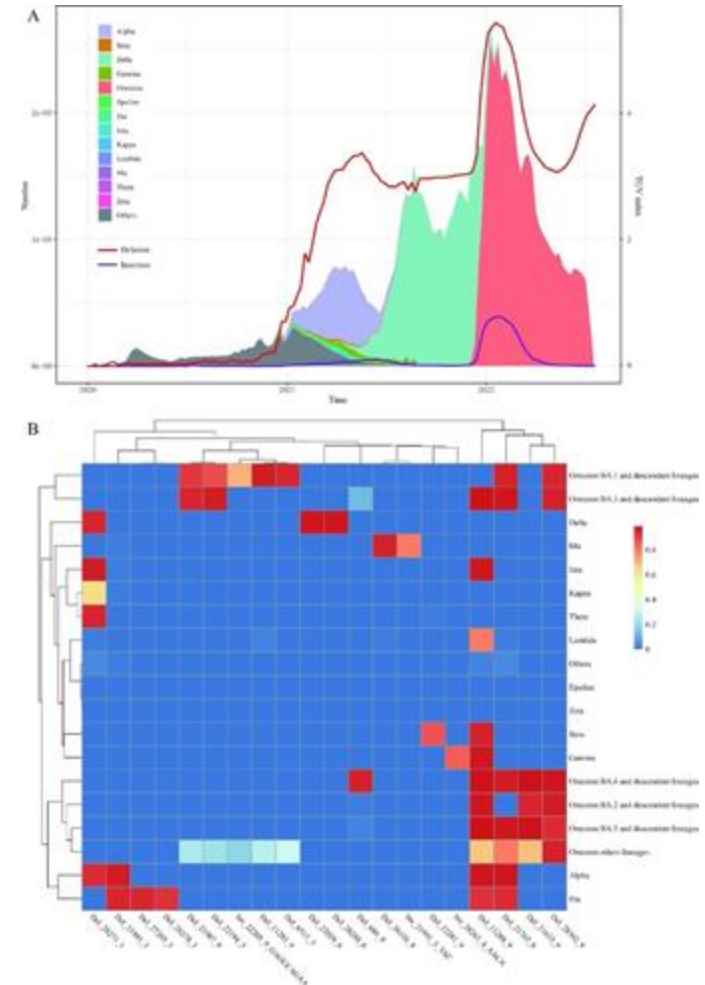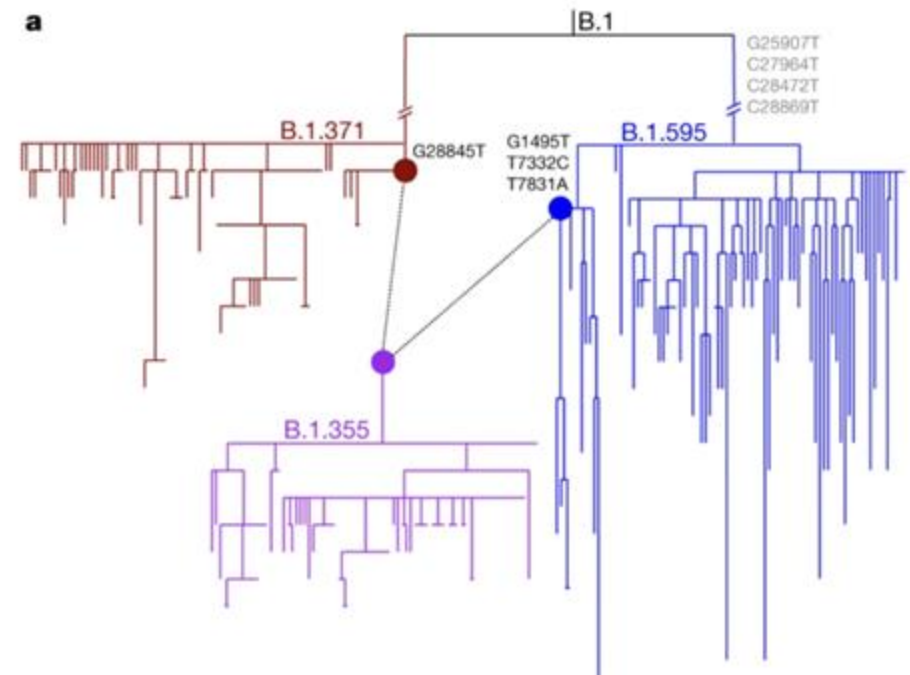# Summary of features in Pangenome formats

| | VG | GFA | GBZ | PanGraph | UShER-MAT | tskit |
|---|---|---|---|---|---|---|
| **Lossless Sequence Encoding** | ✓ | ✓ | ✓ | ✓ | | |
| **Genomic Variation / m-WGA** | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| **Phylogenetic Relationship** | | | | ✓ | ✓ | ✓ |
| **Single-nucleotide Substitutions** | | | | | ✓ | ✓ |
| **Small Indels** | | | | | | ✓ |
| **Structural Mutations** | | | | | | ✓ |
| **Complex Mutations** | | | | | | ✓ |

Mutations

# Summary of features in Pangenome formats

| | VG | GFA | GBZ | PanGraph | UShER-MAT | tskit | PanMAN (This work) |
|---|---|---|---|---|---|---|---|
| Lossless Sequence Encoding | ✓ | ✓ | ✓ | ✓ | | | ✓ |
| Genomic Variation / m-WGA | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Phylogenetic Relationship | | | | ✓ | ✓ | ✓ | ✓ |
| Single-nucleotide Substitutions | | | | | ✓ | ✓ | ✓ |
| Small Indels | | | | | | ✓ | ✓ |
| Structural Mutations | | | | | | ✓ | ✓ |
| Complex Mutations | | | | | | ✓ | ✓ |

Mutations

**Inferred MSA, Phylogeny,** and **mutations** all in **one format!**
**PanMAN** is not just **information-rich** but also more **compact** and **scalable**
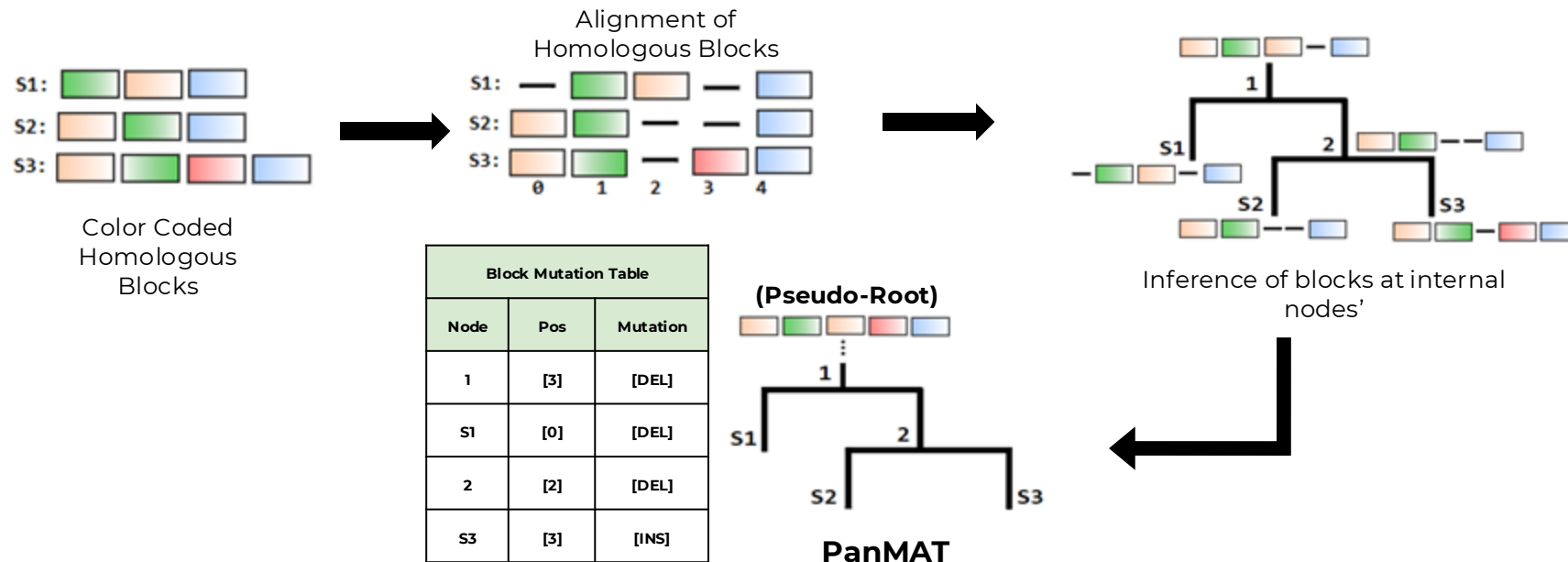
# PanMAT: Pangenome Mutation-Annotated Tree

- Incorporating insertions and deletions (indels) into a MAT
  - MSA defines the coordinate system
  - Gaps treated as special characters



Multiple Sequence Alignment

Phylogenetic Tree

Fitch / PastML

Inference of sequence at internal nodes

Root Sequence

PanMAT

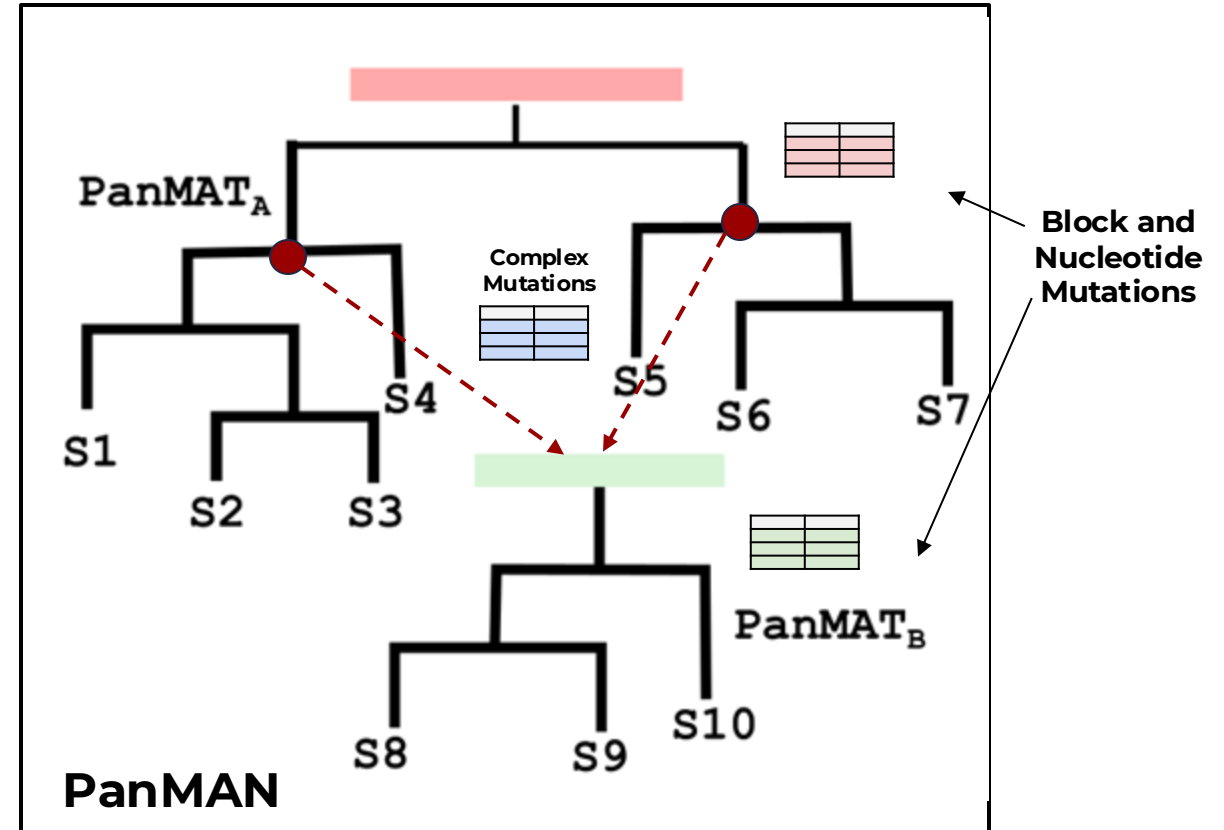| Nucleotide Mutation Table | | |
|---|---|---|
| Node | Pos | Allele |
| S1 | [0,3,4] | [-,G,C] |
| 2 | [6,9] | [A,-] |
| S2 | [8] | [-] |

# PanMAT: Pangenome Mutation-Annotated Tree

- Incorporating **structural changes** and **rearrangements**
    - Identify homologous blocks
    - MSA of homologous blocks
    - Block mutations are like substitutions to or from gaps



Color Coded Homologous Blocks

Alignment of Homologous Blocks

Inference of blocks at internal nodes'

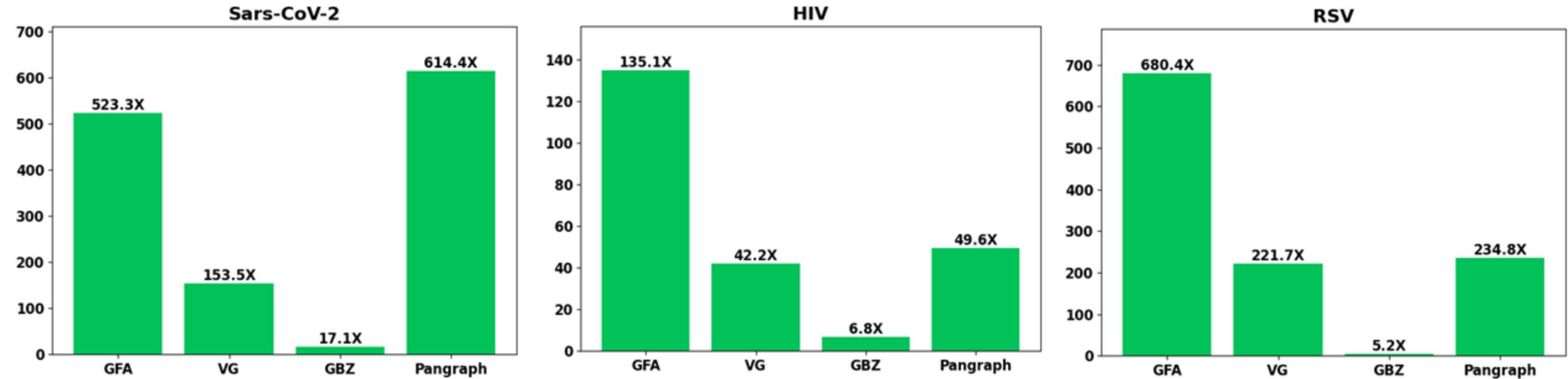| Block Mutation Table | | |
|---|---|---|
| **Node** | **Pos** | **Mutation** |
| 1 | [3] | [DEL] |
| S1 | [0] | [DEL] |
| 2 | [2] | [DEL] |
| S3 | [3] | [INS] |

(Pseudo-Root)

PanMAT

# PanMAN: Pangenome Mutation-Annotated Network

- **PanMAN:** Generalization of PanMAT to represent **complex mutations**

- One or more PanMATs are connected with network edges (red dotted lines)

- Network Edges stores complex mutations (blue table), i.e., Horizontal Gene Transfer (HGT) and Recombination
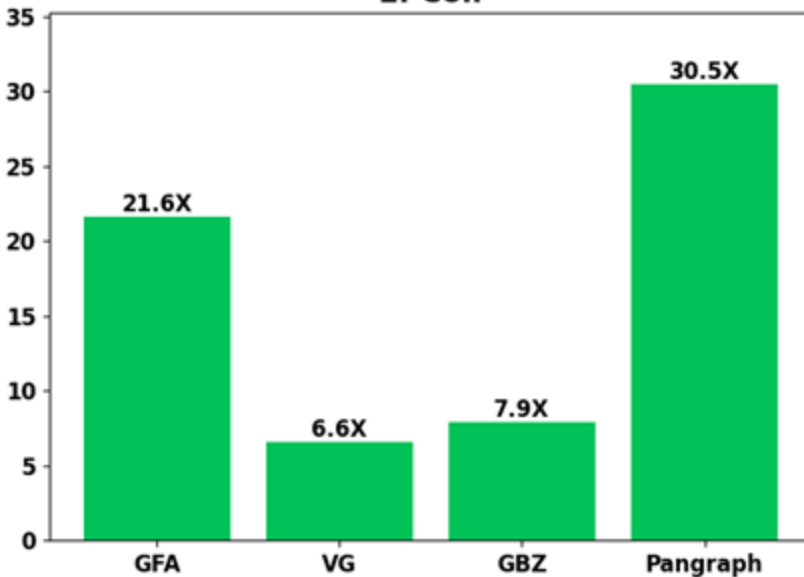
# PanMAN is the most compressive pangenomic format

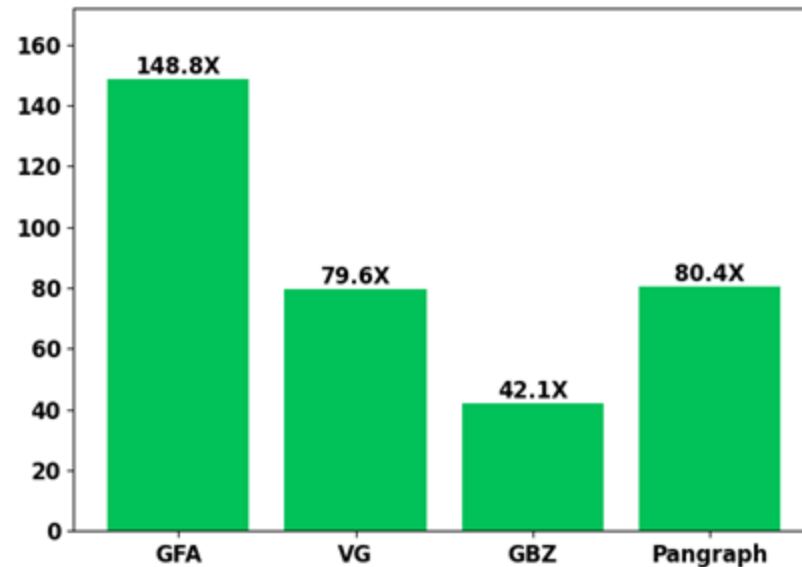**Compression achieved by PanMAN compared to other formats**

# PanMAN is the most compressive pangenomic format
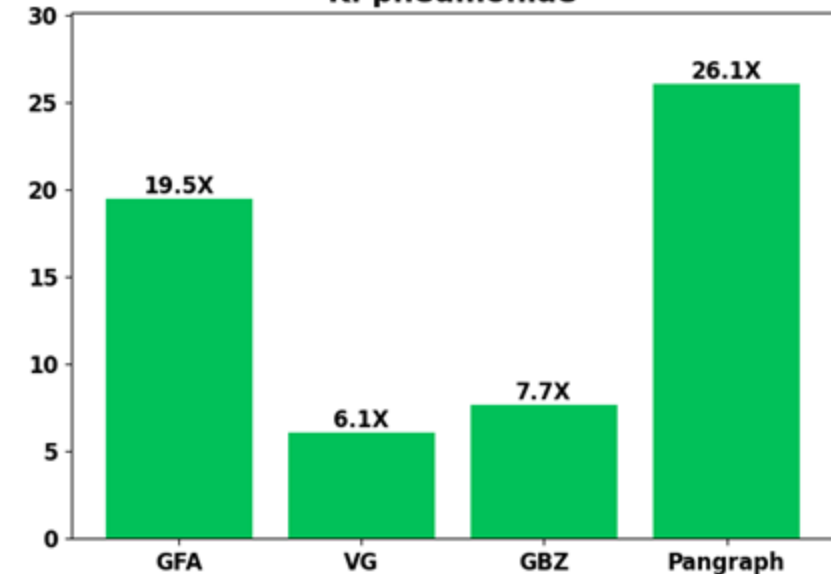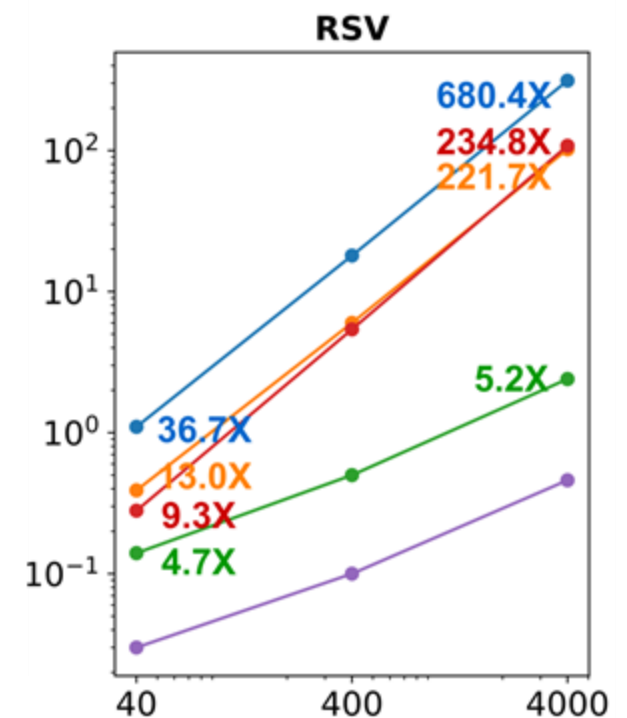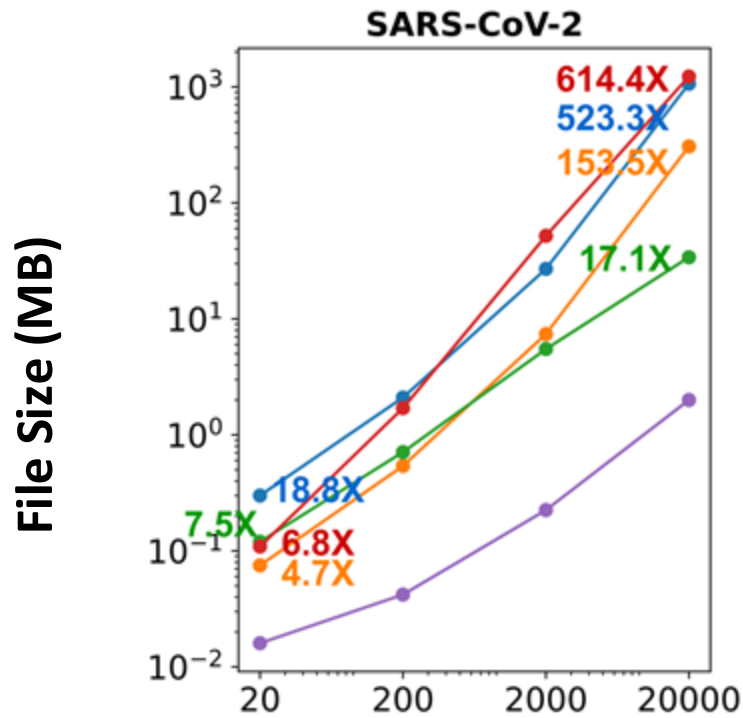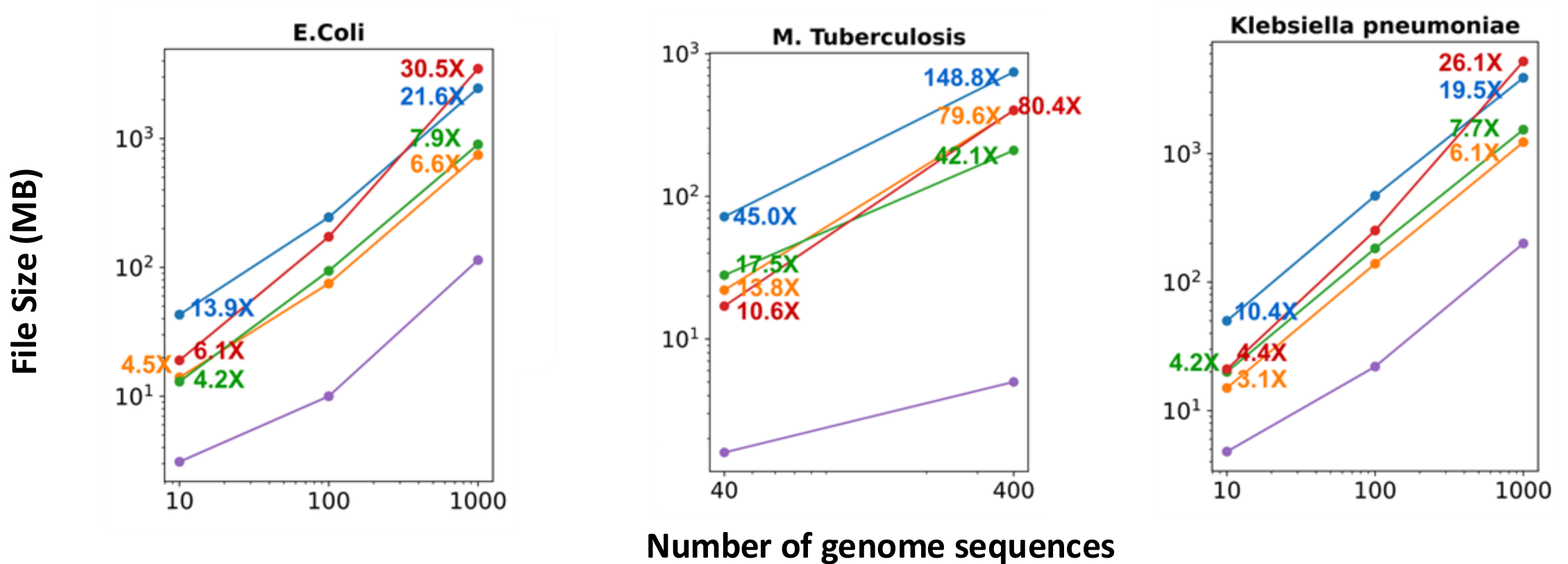
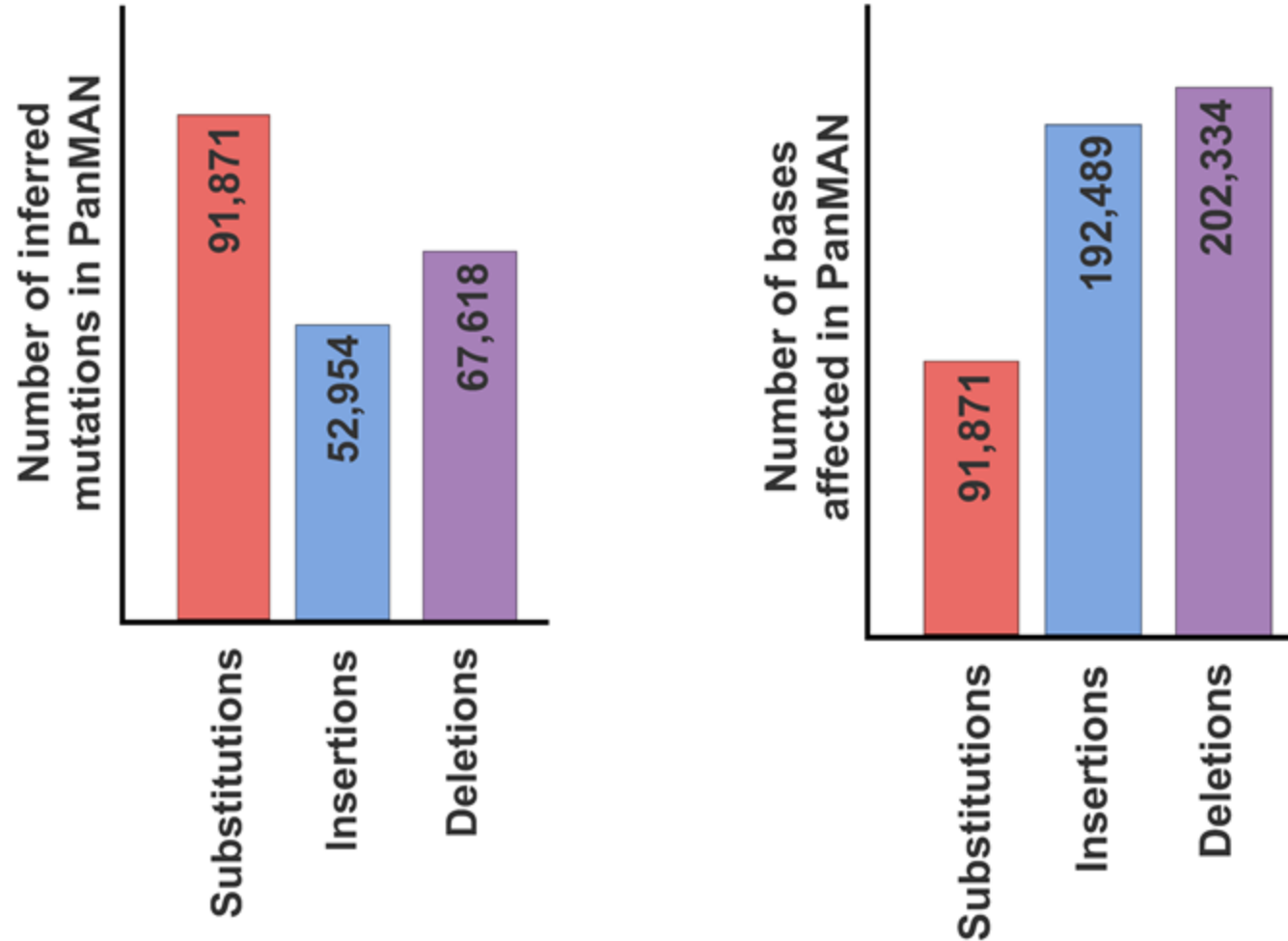**Compression achieved by PanMAN compared to other formats**

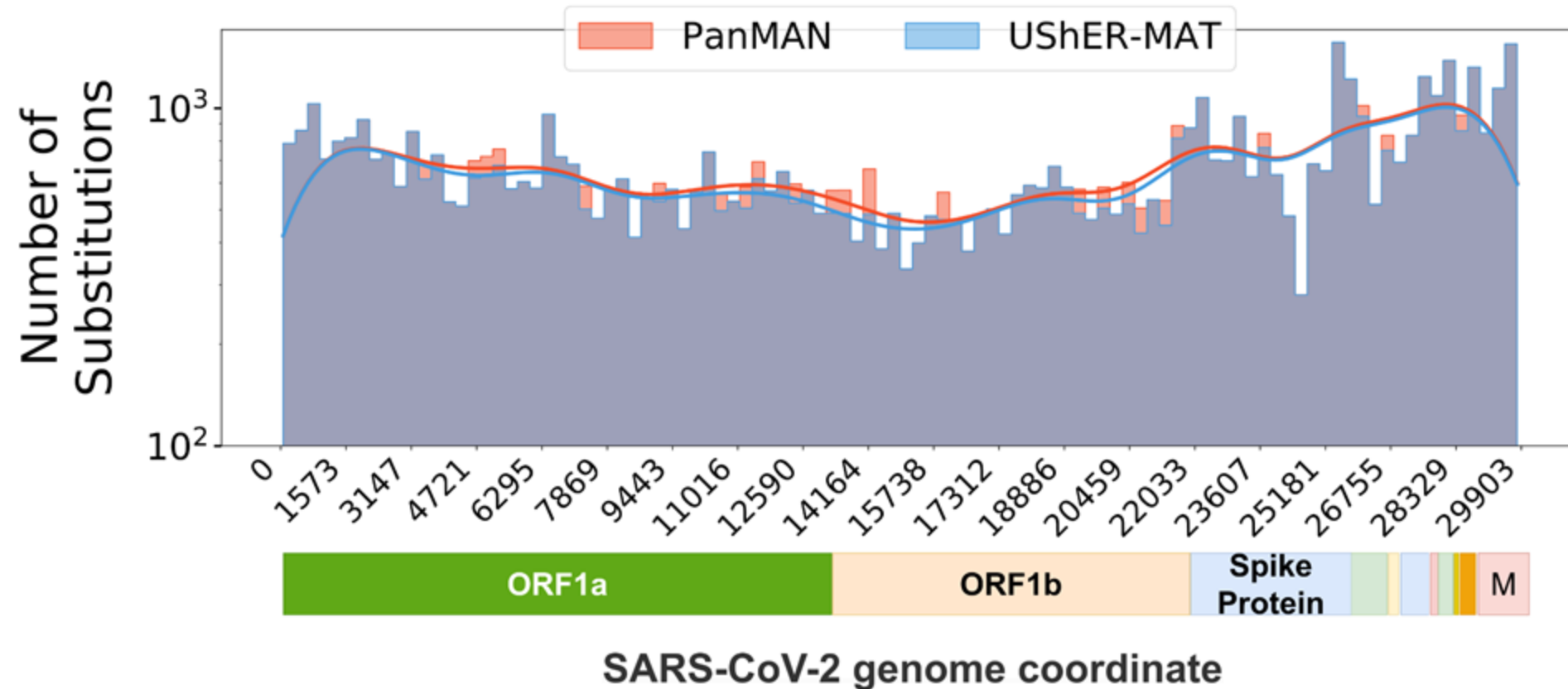# PanMAN scales well relative to other formats

# PanMAN scales well relative to other formats

# Exploration of SARS-CoV-2 mutational and evolutionary landscape using PanMAN

# Exploration of SARS-CoV-2 mutational and evolutionary landscape using PanMAN

# Exploration of SARS-CoV-2 mutational and evolutionary landscape using PanMAN
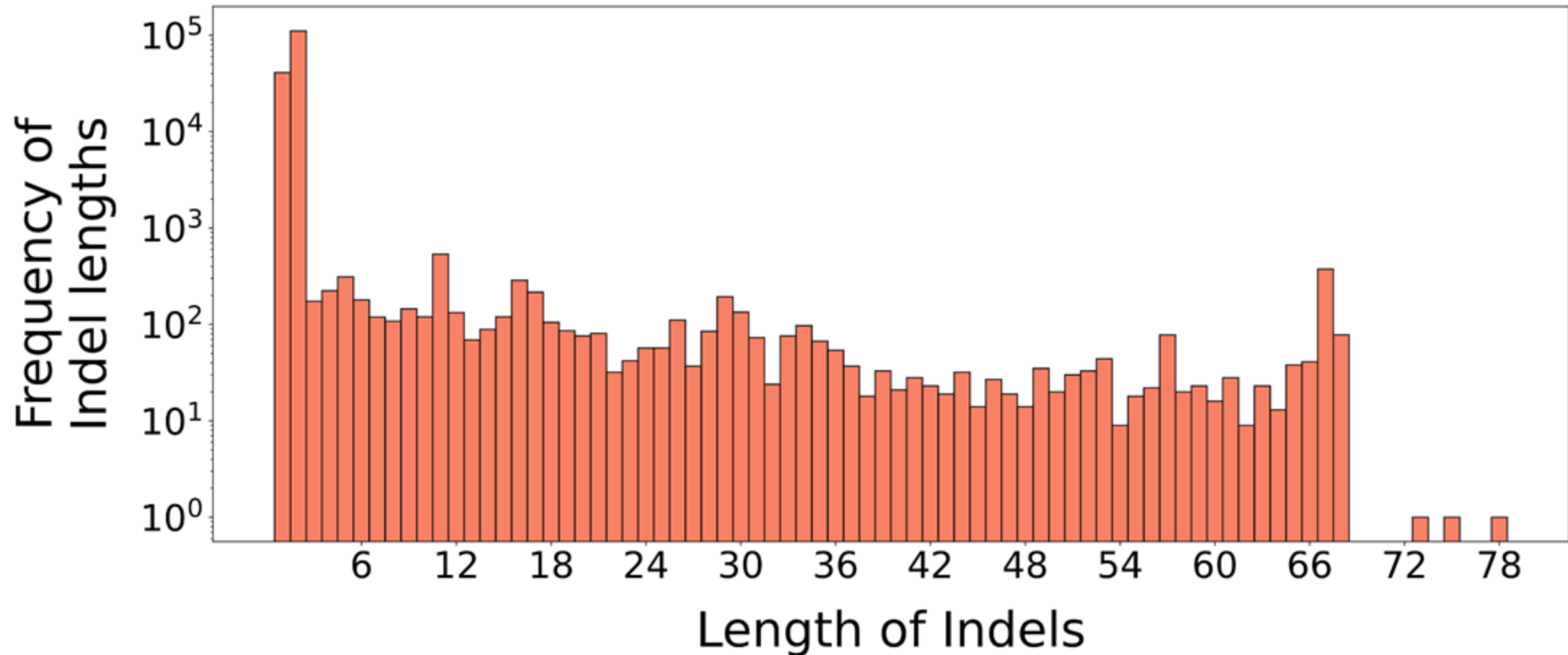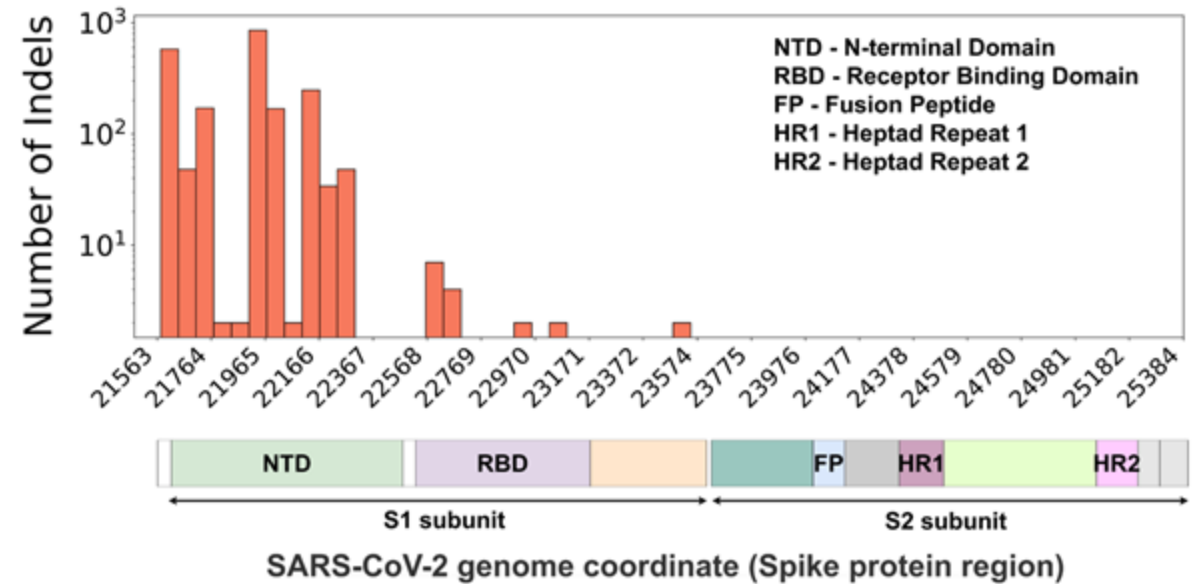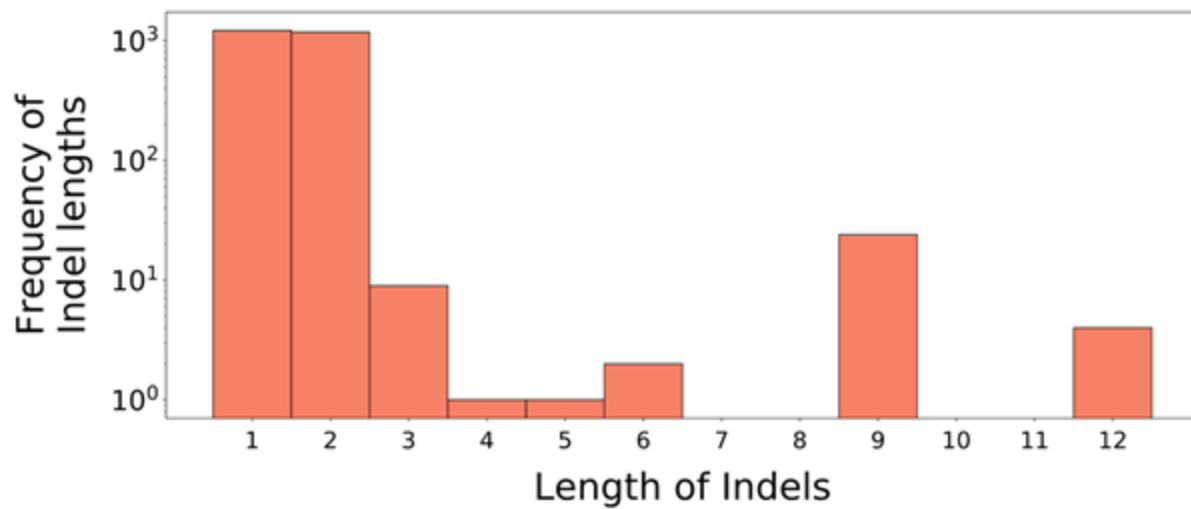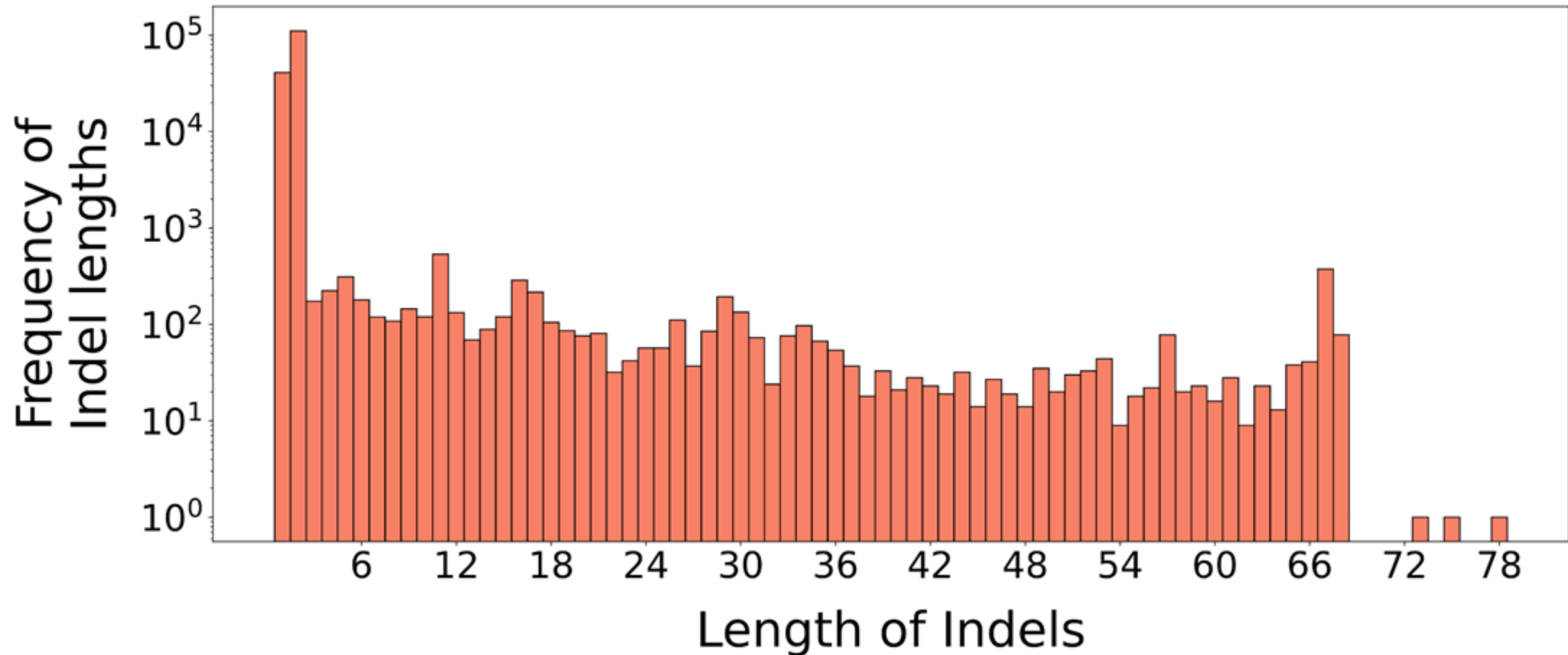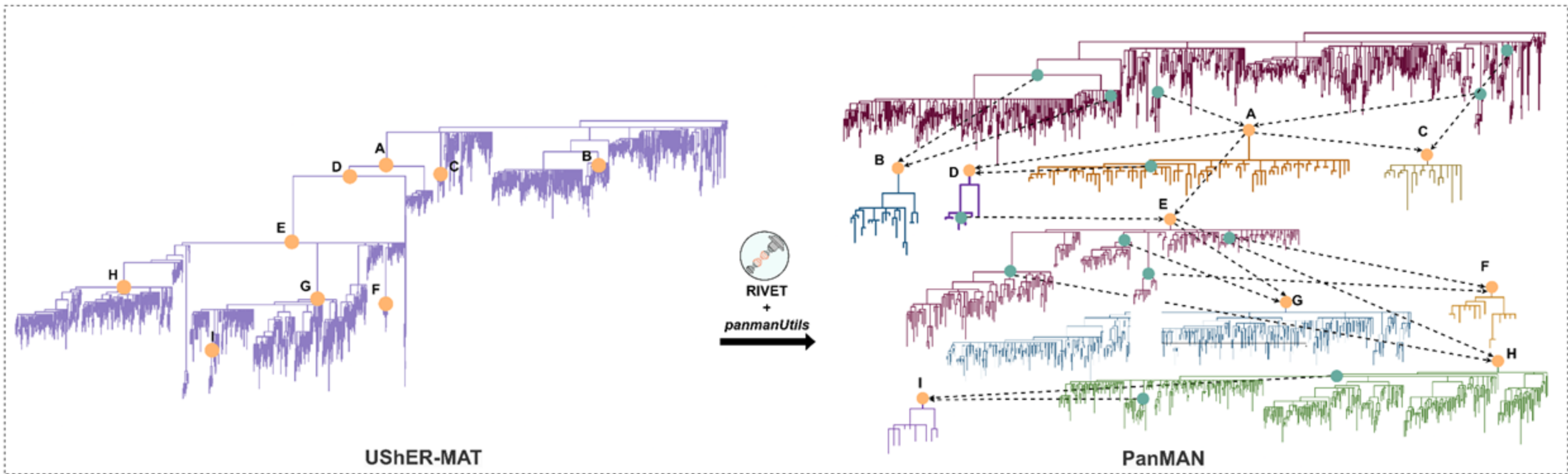
# Exploration of SARS-CoV-2 mutational and evolutionary landscape using PanMAN

# Exploration of SARS-CoV-2 mutational and evolutionary landscape using PanMAN

# Exploration of SARS-CoV-2 mutational and evolutionary landscape using PanMAN

| Pango Designation (WHO labels) | Mutation Type | Mutated Characters | Mutated Position | Mutated Length | Represented in PanMAN? |
|---|---|---|---|---|---|
| BA.1 (Omicron) | Insertion | GAGCCAGAA | 22205 | 9 | Yes |
| | Deletion | N/A | 11283 | 9 | Yes |
| | Deletion | N/A | 6513 | 3 | Yes |
| | Deletion | N/A | 21765 | 6 | Yes* |
| | Deletion | N/A | 21987 | 9 | Yes* |
| | Deletion | N/A | 22194 | 3 | Yes |
| BA.2 (Omicron) | Deletion | N/A | 11288 | 9 | Yes* |
| | Deletion | N/A | 21633 | 9 | Yes |
| | Deletion | N/A | 28362 | 9 | Yes* |
| P.1 (Gamma) | Deletion | N/A | 11288 | 9 | Yes |
| | Insertion | AACA | 28263 | 4 | Yes |
| B.1.617.2 (Delta) | Deletion | N/A | 22029 | 6 | Yes |
| | Deletion | N/A | 28271 | 1 | Yes* |
| | Deletion | N/A | 28248 | 6 | Yes |
| B.1.1.7 (Alpha) | Deletion | N/A | 11288 | 9 | Yes |
| | Deletion | N/A | 21765 | 6 | Yes |
| | Deletion | N/A | 21991 | 3 | Yes |

UShER-MAT

RIVET
+
*panmanUtils*

PanMAN

# PanMANs using likelihood

- Ancestral sequences in PanMAN can be inferred by a a variety of techniques:
  - Parsimony, e.g. Fitch algorithm
  - Likelihood, e.g. PastML, MPPA

- Appears to have a noticeable impact on the file sizes

| Dataset | # sequences | File size (MB) | | Ratio (LK/parsimony) |
|---|---|---|---|---|
| | | PanMAT-Parsimony | PanMAT-LK | |
| Sars-CoV2 | 20 | 0.019 | 0.019 | 1 |
| | 200 | 0.083 | 0.1 | 1.2 |
| | 2000 | 0.68 | 0.91 | 1.3 |
| | 20000 | 4.8 | 6.1 | 1.3 |
| RSV | 50 | 0.047 | 0.17 | 3.6 |
| | 500 | 0.137 | 0.65 | 4.7 |
| | 5000 | 1.1 | 4.3 | 3.9 |
| TB | 40 | 1.9 | 9.3 | 4.9 |
| | 400 | 5.1 | 39.7 | 7.8 |

# PanMAN Utility for Common Bioinformatic Analyses