

# A SARS-CoV-2 Phylogenetic tree based on the Normalized Compression Distance

Presented by: Aloysio Galvão Lopes ([alloysio.galvao-lopes@ip-paris.fr](mailto:alloysio.galvao-lopes@ip-paris.fr))

Professor: Jean-Louis Dessalles

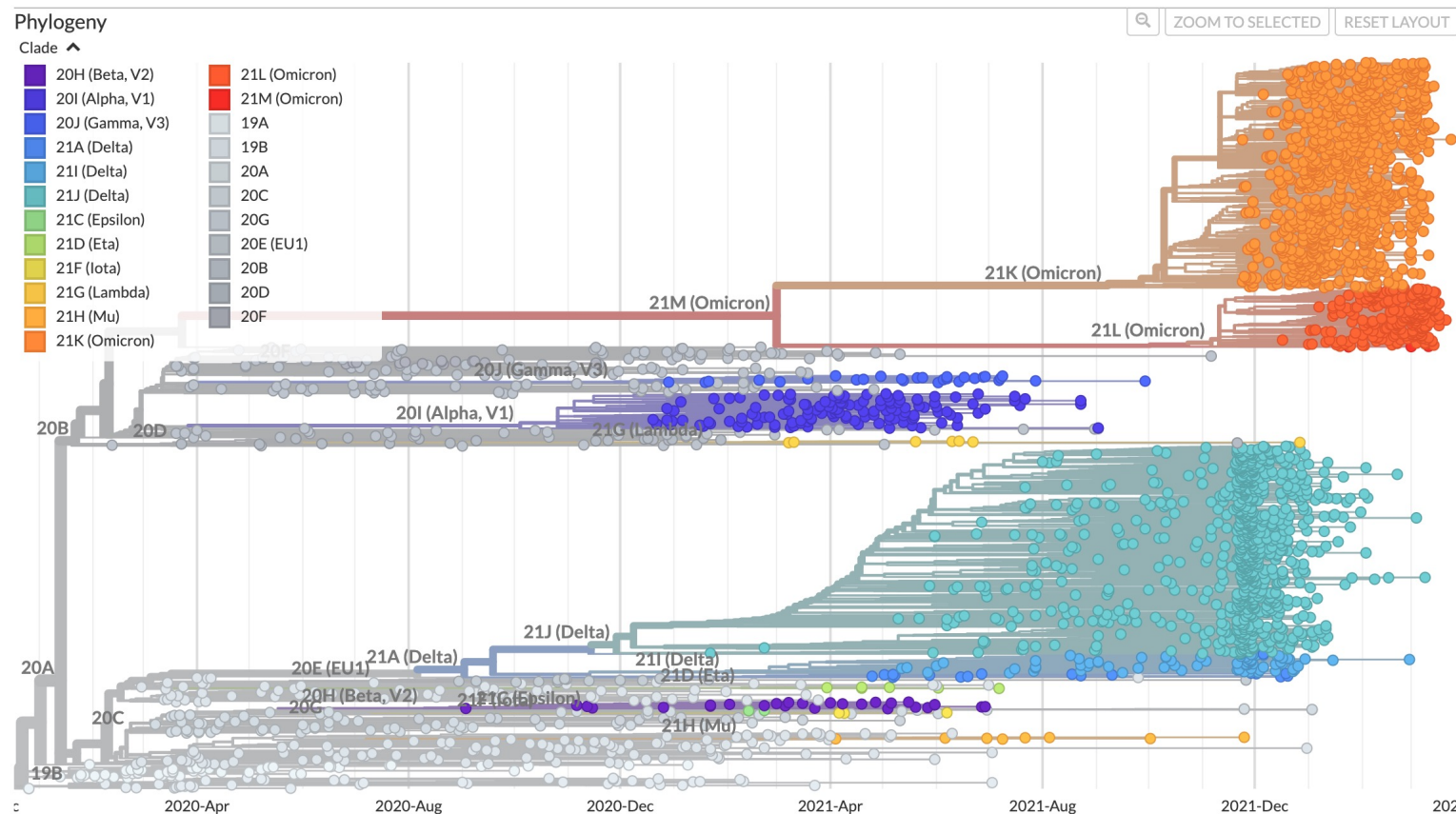
Key topics: Bioinformatics, Normalized Compression Distance, Algorithmic Information, Pangolin Classification, COVID-19

# Summary

- Introduction to SARS-CoV-2 Pango naming and lineage classification
- Approach
- Results and comparison with SOTA
- Conclusions

# Introduction

- How are SARS-CoV-2 lineages named?



Source: <https://nextstrain.org/ncov/gisaid/global>

# Introduction

- The Pango nomenclature system  
A, A.1, B.1, B.1.1

# Introduction

- The Pango nomenclature system

A, A.1, B.1, B.1.1, B.1.529 (WHO name: Omicron)

More info: <https://www.pango.network/>

# Introduction

- The Pango nomenclature system

A, A.1, B.1, B.1.1, B.1.529 (WHO name: Omicron)

More info: <https://www.pango.network/>

- Pangolin uses a ML model to assign lineages, this process needs sequence alignment, they're aligned based on A



Sources: <https://cov-lineages.org/> and [https://en.wikipedia.org/wiki/Sequence\\_alignment](https://en.wikipedia.org/wiki/Sequence_alignment)

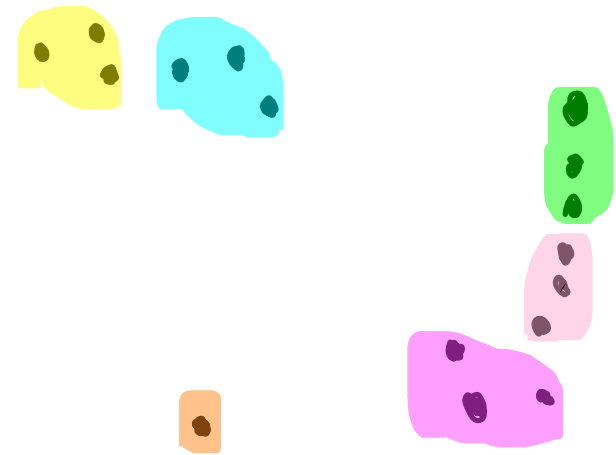
# Approach

- All sequences taken from [NCBI](#), a little over 52 GB of genome data. Total 883,020 genomes
- Sampled max 2 genomes per day totaling 1247 genomes
- Distance matrix computed in parallel using Python's [zlib](#)

# Approach

- The distance used is a slightly changed version of the Normalized Compression distance (NCD)
- Samples are clustered using hierarchical clustering with complete linkage to make for vizualisation and computation purposes

$$NCD(x, y) = \frac{\frac{Z(xy) + Z(yx)}{2} - \min(Z(x), Z(y))}{\max(Z(x), Z(y))}$$

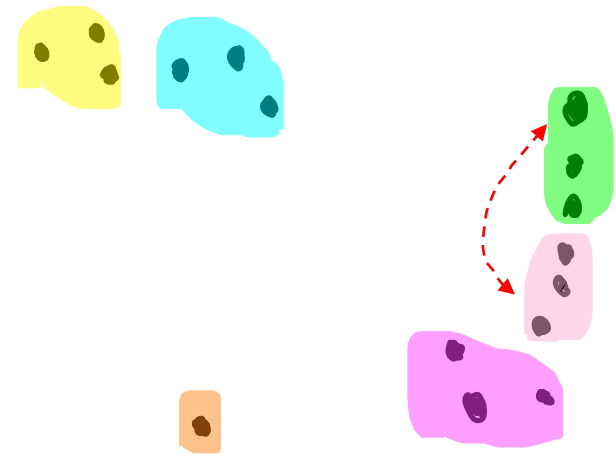




# Approach

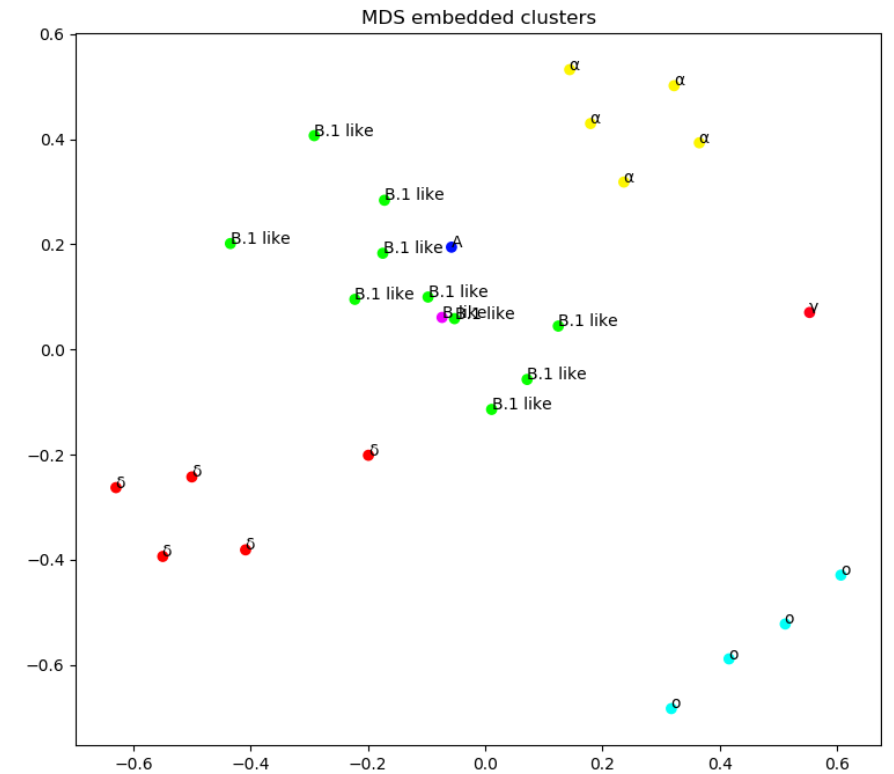
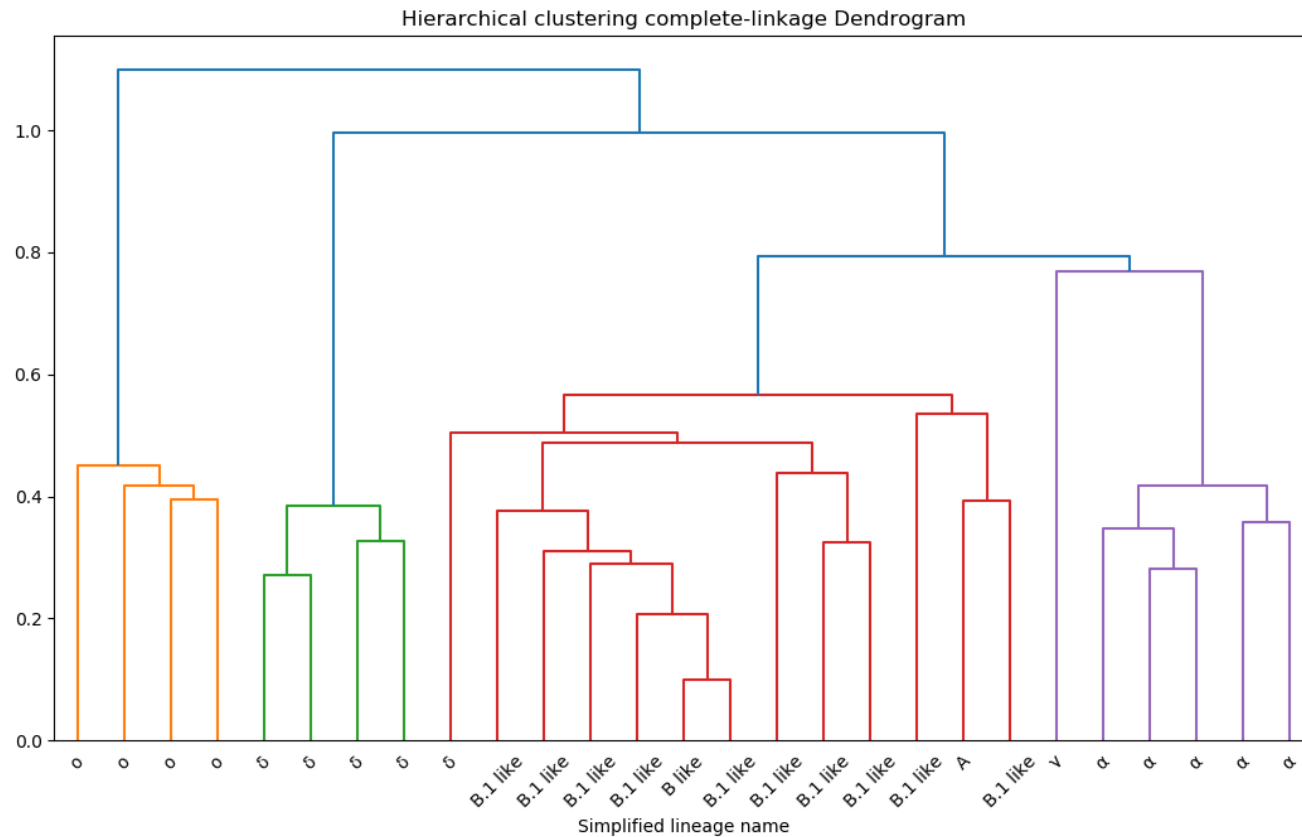
- The distance used is a slightly changed version of the Normalized Compression distance (NCD)
- Samples are clustered using hierarchical clustering with complete linkage to make for vizualisation and computation purposes
- 27 cluters were generated and used for the phylogenetic tree

$$NCD(x, y) = \frac{\frac{Z(xy) + Z(yx)}{2} - \min(Z(x), Z(y))}{\max(Z(x), Z(y))}$$

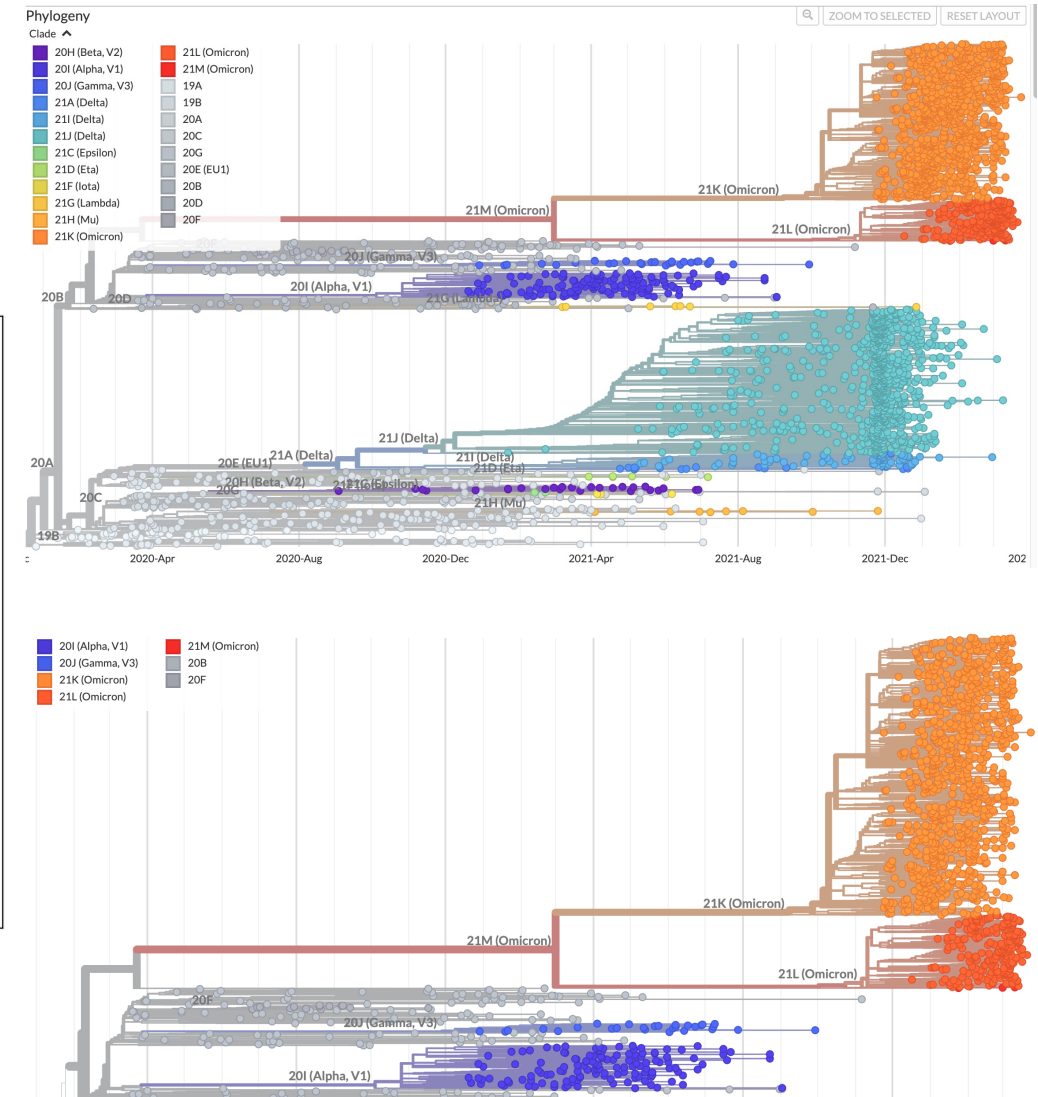
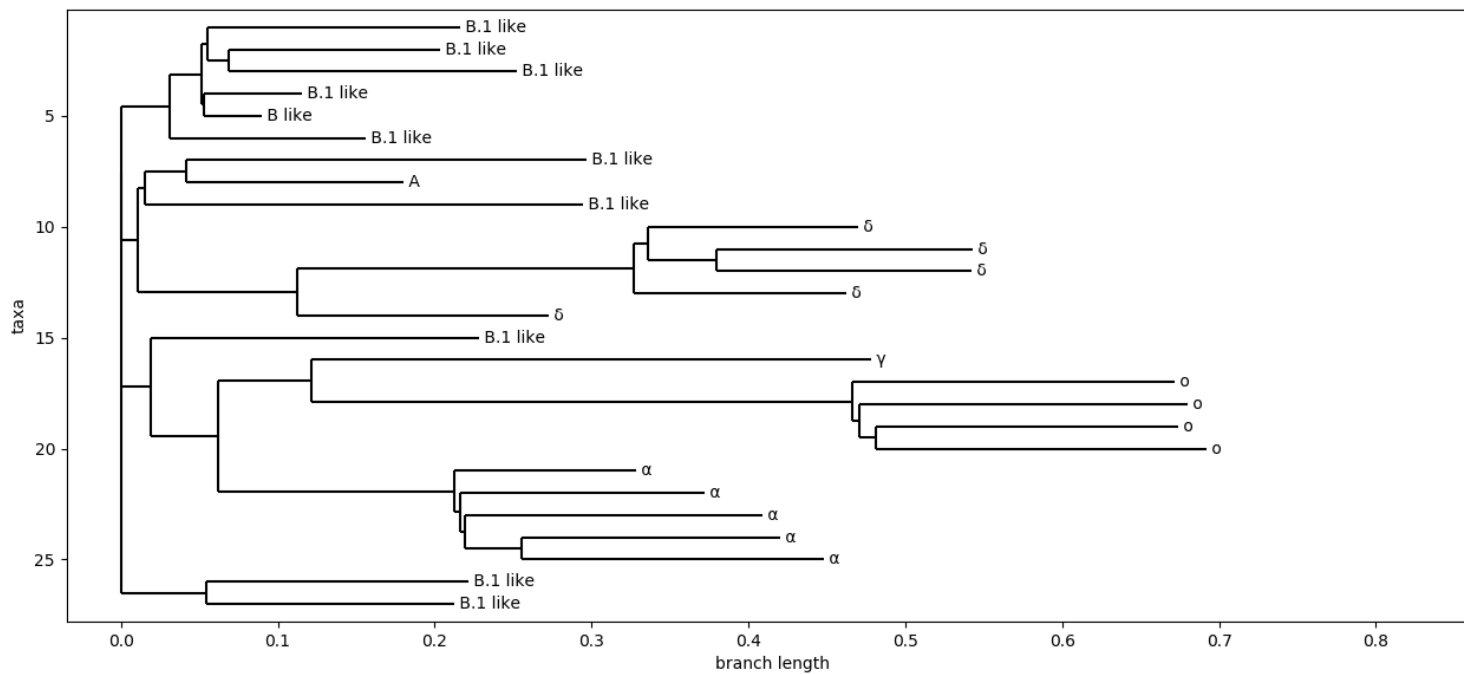


# Results

Variant	Start date	Approx onset date	End date
$\delta$	2021-07-14	2021-05-XX	2022-03-04
$\alpha$	2021-01-05	2020-11-XX	2021-11-24
$\sigma$	2021-12-31	2021-11-XX	2022-03-21
$\gamma$	2021-05-05	2021-02-XX	2022-03-18



# Results



# Conclusions

## Summary

- Using the normalized compression distance (NCD), I've computed distances between a subsample of all covid cases
- I've created agglomerative clusters with complete linkage
- I've used those clusters to build the final tree

## Takeaways

- The general genealogy of the virus has been successfully found
- AIT methods don't require sequence alignment
- The NCD can be difficult to interpret, as well as we lose some information when working only with distances