# Rapid SARS-CoV-2 Whole-genome Sequencing and Analysis for Informed Public Health Decision-making in the Netherlands Replication and Extension

## Presented by:

**De Ama, Syd**
**Dizon, Franz**
**Guaro, Sidney**

**CSC 771M**

**February 10, 2021**

# Presentation Outline

# Replication of Phylogenetic Analysis in Netherlands Paper

In this part of the report, we replicate the Phylogenetic Analysis part of Munnick's O., et al. (2020, Sep 26) paper. We skipped the Sequence Data Analysis part because the sequences are already available in the GISAID (https://www.gisaid.org/). Consequently, we only performed the Phylogenetic Analysis part. We also have variances in the tools used because of our computational unavailability.

## Data Acquisition

In December 2020, we collected all the sequences and Netherlands only sequences from GISAID (https://www.gisaid.org/) with a submission date of less than or equal to March 17, 2020. The sequences collected are 873 and 189 respectively. The 189 sequences from Netherlands validates our sequences as it has the same number as indicated in Munnink's O., et al. (2020, Sep 26) paper.

## Sequences Preprocessing

Similar to Munnink's O., et al. (2020, Sep 26), we filtered the 873 sequences collected from GISAID, which exclude sequences with > 10% NNNs, where the NNNs are known to be the sequence error reads. This is achieved by creating a small python script with the help of a biopython library developed by Cock PA., et al. (2009). The said python script is shown below.

```
from Bio import SeqIO
f_file = "inputFile.fasta"
count = 0
f_sequences = SeqIO.parse(open(f_file), 'fasta')

with open("outputFile.fasta","w") as f:
        for seq in f_sequences:
                seq_seq = seq.seq
                seq_length = len(seq_seq)
                if(seq_seq.count('N')/seq_length<.10):
                        SeqIO.write([seq],f,'fasta')
```

After applying the filter, there are 843 sequences retained in total. These sequences will then be used in Multiple Sequence Alignment (MSA).

## Multiple Sequence Alignment

The MSA on the collected sequences was performed using MAFFT (Multiple Alignment using Fast Fourier Transform) software. We used MAFFT instead of MUSCLE which has been used by Munnink's O., et al. (2020, Sep 26) due to time and memory constraints, MAFFT is a more lightweight program with multithreading support and memory saving option for computers without a huge amount of RAM (Katoh and Toh, 2010).

The command used are shown below:

```
mafft --retree 2 --thread -1 --inputorder "inputFile.fasta" > "outputFile.afa"
```

The parameter "--retree 2" is a default parameter which sets the multiple alignment strategy to use the FFT-NS-2 progressive method. FFT-NS-2 is stated to be capable of aligning a number of sequences up to 5000 on a standard desktop computer (Katoh, et al., 2010). The parameter "--thread -1" is used to automatically set the optimal number of cores to be used. The last parameter used for configuration is the "--inputorder" which is also a default parameter, it indicates that the order of of sequences in the output file, is the same as the order of sequences in the input file.

After the MSA is generated, we manually checked the sequences for discrepancies.

## Tree Construction

In this replication, the Maximum Likelihood (ML) tree was constructed using IQ-TREE (Nguyen, L. et al. 2014) for the aligned sequences output from MSA. We used the ultrafast bootstrap option with 1,000 replicates and GTR+F+I+64 model, the same with the paper of Munnick's O., et al. (2020, Sep 26).

There are several files that are being produced by the IQ-TREE, one of which has a file extension of .treefile which is what we used to visualize the tree. The said file has a file format of newick file, which is then imported to our custom python scripts which uses baltic library for visualizing the tree (Dudas, 2016).
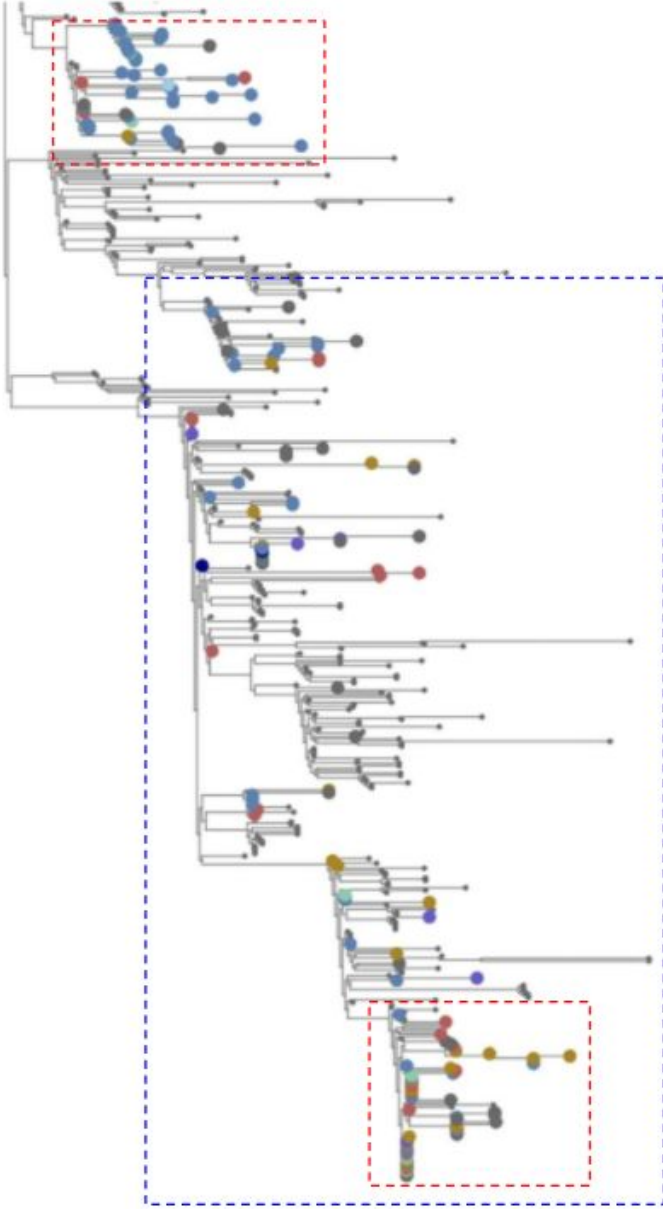
The command used are shown below:

```
.\iqtree2 -s netherlands_world_mafft.afa -bb 1000 -m GTR+F+I+64 -nt AUTO
```

The parameter -s specifies the name of the alignment file. The parameter -bb sets the number of replicates used for bootstrapping. The parameter -m specifies the model to be used. And lastly, the -nt parameter specifies the number of threads used.

# Discussion

| Munnink's, O., et al. (2020) phylo. tree | Our generated phylogenetic tree |
|:---:|:---:|
| Figure 1 (Whole tree shown on appendix A) | Figure 2 (Whole tree shown on appendix B) |



| | Drenthe |
|---|---|
| | Flevoland |
| | Gelderland |
| | Limburg |
| | Noord-Brabant |
| | Noord-Holland |
| | Overijssel |
| | Utrecht |
| | Zuid-Holland |

It can be observed that the phylogenetic tree for the Netherlands SARS-Cov-2 cases generated by Munnick's O., et al. (2020), shown in figure 1, was quite similar to the one we have created, which is shown in figure 2. The first major increase of cases in the Netherlands was likely due to multiple introductions of the virus, following the spring holidays where people travel to neighboring countries such as Northern Italy for its ski resorts. This caused several co-circulating viruses, resulting in cases appearing on different clusters in the phylogenetic tree, which can be seen on both figures 1 and 2 indicated by the blue-dashed box. Other cases which add to the diversity of the viruses, which further supports the evidence of the existence of the co-circulating viruses, can be observed on appendix A and B, outside of the green box. On the other-hand, clustering of cases which suggests an increasing number of local transmission, can also be detected on both phylogenetic trees. The upper-red box for both trees indicates the increasing number of local transmission within the province of Noord-Brabant, while the bottom-red box indicates the increasing number of inter-province transmission.

Given this analysis, together with the epidemiological data (e.g. date of exposure, travel history, etc.), the Netherlands government decided to increase the movement restrictions and implement stricter preventive measures for the whole country to prevent further spread of SARS-CoV-2.

# Phylogenetic Analysis on Philippine SARS-CoV-2 Sequences Compared to World

The dissemination of SARS-CoV-2 in the world led the countries to control the transmission of diseases by implementing travel bans and community lockdowns. In the Philippines, with a number of foreign visitors and workers, the first two cases recorded were of a Chinese couple entering the country in January 2020. As of February 08, 2021, the COVID-19 cases in the Philippines reached up to 537,310. According to Tablizo, et al (2020), the sequenced Philippine isolates combined with GISAID data, can be divided into three main categories based on sampling dates and potential sources of infection: (1) Samples obtained in the early stages of the pandemic in January are closely related to isolates from Wuhan, China; (2) Samples from March, primarily related to the outbreak of the M/V Diamond Princess Cruise Ship; and (3) Samples clustered with European isolates in June. The clusters should still hold true in the case of the GISAID sequences since there are no recent updates in regards to the Philippine sequences.

In this part of the report, we extend Munnink's, O., et al. (2020, Sep 26) paper by utilizing the phylogenetic analysis methods stated in the study and applying it to the Philippine SARS-CoV-2 sequences. We will utilize the available sequences from the GISAID as it is the only known source which we can acquire the sequences publicly. In this extension, we present the MSA and tree construction  of Philippine SARS-CoV-2 genomes in order to analyze how the virus got into the Philippines or from the Philippines, how the virus spread.

## Data Acquisition

In January 2021, we investigated how to collect all sequences from GISAID. The option to download all sequences was not available during that time. The only functionally available was the option to download 10,000 sequences at a time from the interface thus we explored the option of web-scraping the GISAID. In mid-January 2021, we downloaded Philippine SARS-CoV-2 sequences for initial pipeline implementation and testing. In the third week of January 2021, we discovered an option that allows us to download all the sequences at once. This option involves communicating with the GISAID support team and acquiring their approval. Fortunately, we obtained the option and on February 03, 2021, the sequences were collected from GISAID.

The sequences collected were aligned and needed to be preprocess in order to perform the sampling selection. The total number of full alignment sequences is 422,925. The sequences are based on 456,943 submissions to EpiCoV. Both duplicate and low-quality sequences are filtered by GISAID where (>5% NNNNs) have been removed, using only complete sequences (length >29,000 bp).

## Sequences Preprocessing

The full alignment sequences are 422,925 in total. This will make the MSA and Tree Construction computationally expensive and burdensome on the professionals who analyze phylogenetic trees. To relieve the resources, we first removed the indels (-) in the sequences. Next, we filtered the sequences excluding sequences with > 10% NNNs. This is done in relation to Munnink's, O. et al. (2020, Sep 26) paper. Afterwards, we removed the Philippine sequences to avoid the duplication when adding the whole Philippine sequences.

After the filtration, we sampled the sequences using two methods. The first sampling method is the 10 sequences per country method. This sampling method assumes that the sequences are ordered in descending matter according to date. In the filtered sequences, the sequences order is ascending so we reversed the order and performed the sampling. In the 10 sequences per country sampling, the first 4 sequences, at maximum, from the first month are retained and followed by 4 sequences from the previous month, at maximum, and so on until the total number of sequences of a country is equal to 10. If a month has less than 4 sequences, it will proceed to the previous month (e.g. 2021-01: 1, 2020-12: 4, 2020-11: 4, 2020-10: 1). In this method, 1,195 samples are selected. The second sampling method is the 1 per month method. This sampling method retains the sequences one per month at maximum (e.g. 2021-01: 1, 2020-12: 1, 2020-11: 1, …, 2019-12: 1). In this method, 1,595 samples are selected.

Since the primary goal of this study is to analyze Philippine SARS-CoV-2 mutation and spread, after the sampled sequences are selected, the Philippines sequences are added to the sampled sequences. As a result, the total number of sequences on both sampling methods are 1,231 and 1,631 respectively.

## Multiple Sequence Alignment

The MSA on the collected sequences was performed using MAFFT software. We used MAFFT for similar reasons stated in the replication. We also used the similar command stated in the replication and manually checked the sequences for discrepancies (See [10_per_country_msa_visualization.pdf](#) and [1_per_country_msa_visualization.pdf](#) files for the visualization of MSA).

## Tree Construction

In the paper extension, the ML trees were constructed using IQ-TREE (Nguyen, L. et al. 2014) for both the multiple sequence aligned outputs. We used the ultrafast bootstrap option with 1,000 replicates for both the aligned sampled sequences. For the 1 per month sampling, the best model is GTR+F+R3 and for the 10 per country, the best model is TIM+F+R3.

The command used are shown below:

```
.\iqtree2 -s algined_sampled_sequences.afa -bb 1000 -nt AUTO
```

Between the two phylogenetic trees, we select the 10 per country tree in which we will perform our analysis on. We chose the 10 per country tree because it provides a more general clustering on different countries, which may further help us determine the origins of the sequence residing in the Philippines by looking at which cluster it branched out. Similar to the replication part, the output of the IQTREE was visualized using our custom python script that uses the baltic library (Dudas, 2016).

## Discussion



Generated phylogenetic tree for Philippines analysis
Figure 3 (Whole tree shown on appendix C)

Philippines
Other countries

Phylogenetic tree shown on figure 3 is the one that will be used for the analysis of SARS-Cov-2 cases in the Philippines. It can be observed that there exists two forming clusters, highlighted by the red-dashed box and purple-dashed box, which suggests that an ongoing local transmission has occurred in the country. There are also other cases that are located on various parts of the tree, highlighted by the yellow-dashed box, which suggests that there may be new sequences that have been introduced from other countries, and may introduce further local transmission of cases if no action is taken.
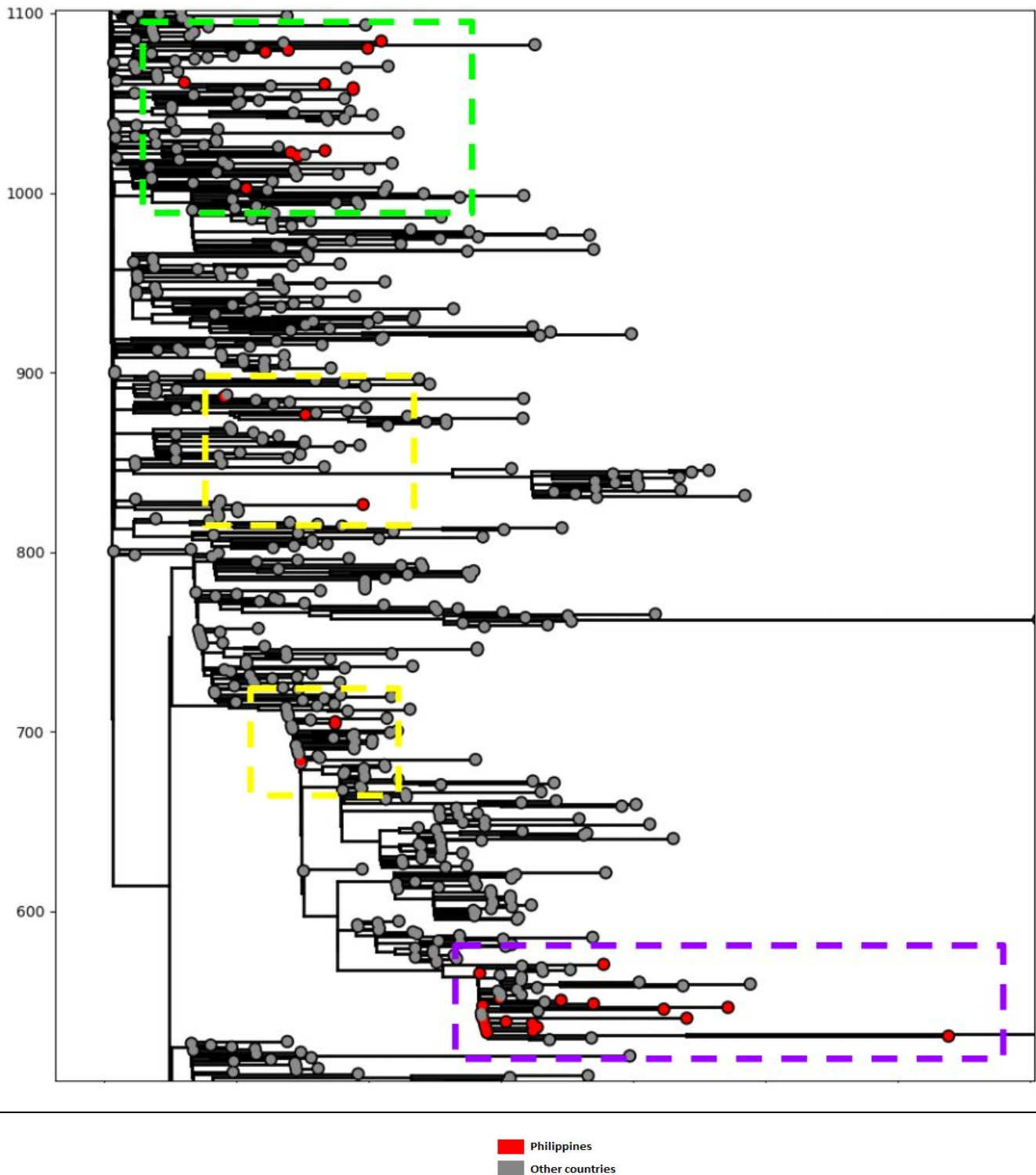


Figure 4. Zoom-in of green-dashed box from figure 3

In the zoom-in view of the green-dashed box shown on figure 4, it suggests that cases in the Philippines that can be found from this cluster may have originated from Europe, and it also indicates that a possible local transmission on the country with this type of strain have already occurred. This is because most of the cases that can be found in this subtree, originated from Europe.
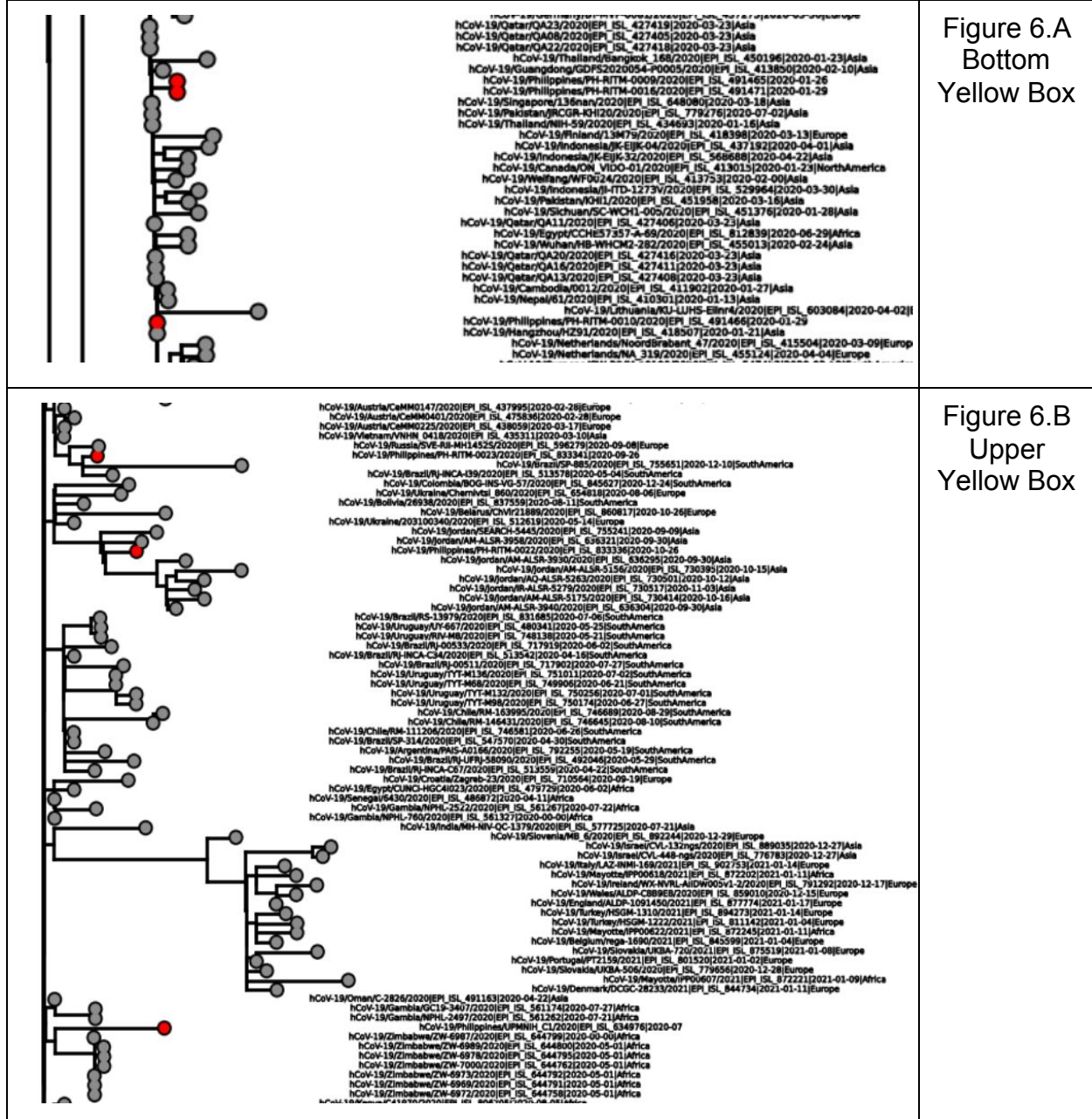
Figure 5. Zoom-in of purple-dashed box from figure 3



In the zoom-in view of the purple-dashed box shown on figure 5, it suggests that cases in the Philippines that can be found from this cluster may have originated from the neighboring countries like Malaysia, Singapore, Brunei, and others. Furthermore, since the cases here are observably denser, it suggests that a relatively large amount of local transmision took place with respect to the green-dashed box. This may be one of the reasons why the Philippine government decided to implement lockdowns and cancellation of flights.

Figure 6. Zoom-ins for yellow-dashed boxes from figure 3

| | |
|---|---|
|  | Figure 6.A Bottom Yellow Box |
|  | Figure 6.B Upper Yellow Box |

In the zoom-in views of the yellow-dashed boxes shown on figure 6, it suggests that there are also different strains of viruses that are circulating in the Philippines. One possible reason that we can think for this is due to an increased number of Overseas-Filipino-Workers (OFW) that are going back home in the country, specially the neighboring cases that can be found under the subtree originated from countries like Qatar, Jordan, and Singapore, which are known countries to have a number of OFWs.
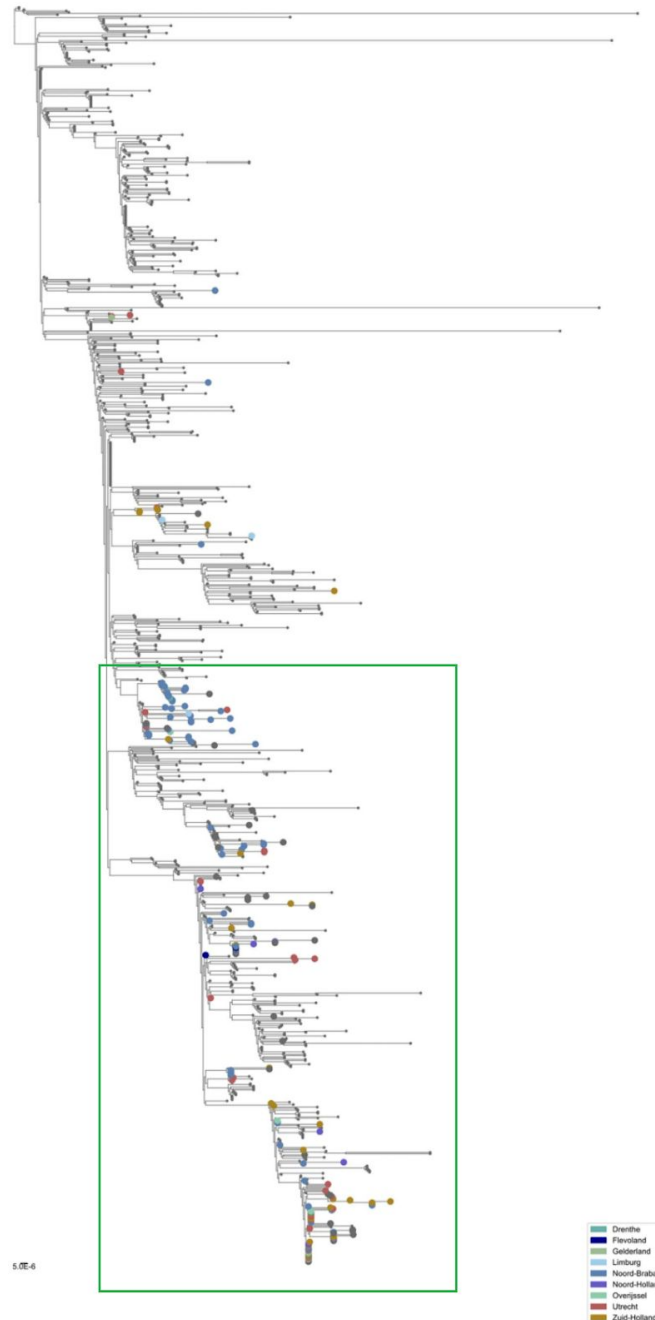
# References

— Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., & Minh, B. Q. (2014, November 3). *IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies*. OUP Academic. https://academic.oup.com/mbe/article/32/1/268/2925592?login=true.

— Munnink, B. B. O., Nieuwenhuijse, D. F., Stein, M., O'Toole, Á., Haverkate, M., Mollers, M., … Koopmans, M. (2020, July 16). *Rapid SARS-CoV-2 whole-genome sequencing and analysis for informed public health decision-making in the Netherlands*. Nature News. https://www.nature.com/articles/s41591-020-0997-y.

— Tablizo, F. A., Lapid, C. M., Maralit, B. A., Yap, J. M. C., Destura, R. V., Alejandria, M. A., ... & Saloma, C. P. (2020). ANALYSIS OF SARS-COV-2 GENOME SEQUENCES FROM THE PHILIPPINES: GENETIC SURVEILLANCE AND TRANSMISSION DYNAMICS. medRxiv.

— Katoh, Misawa, Kuma, Miyata 2002 (*Nucleic Acids Res.* **30**:3059-3066)
MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform.
(describes the FFT-NS-1, FFT-NS-2 and FFT-NS-i strategies)

— Katoh, Toh 2010 (*Bioinformatics* **26**:1899-1900)
Parallelization of the MAFFT multiple sequence alignment program.
(describes the multithread version)

— Dudas, G. (2016). *baltic.* Retrieved from Github: https://github.com/evogytis/baltic

— Cock PA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T, Kauff F, Wilczynski B and de Hoon MJL (2009) Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25, 1422-1423

— ETE Toolkit. (n.d.). Analysis and Visualization of (phylogenetic) trees. Accessed February 09, 2021, from http://etetoolkit.org/

# Appendices

A. Munnink's, O., et al. (2020) whole generated phylogenetic tree using Maximum likelihood for sequences gathered on March 22, 2020
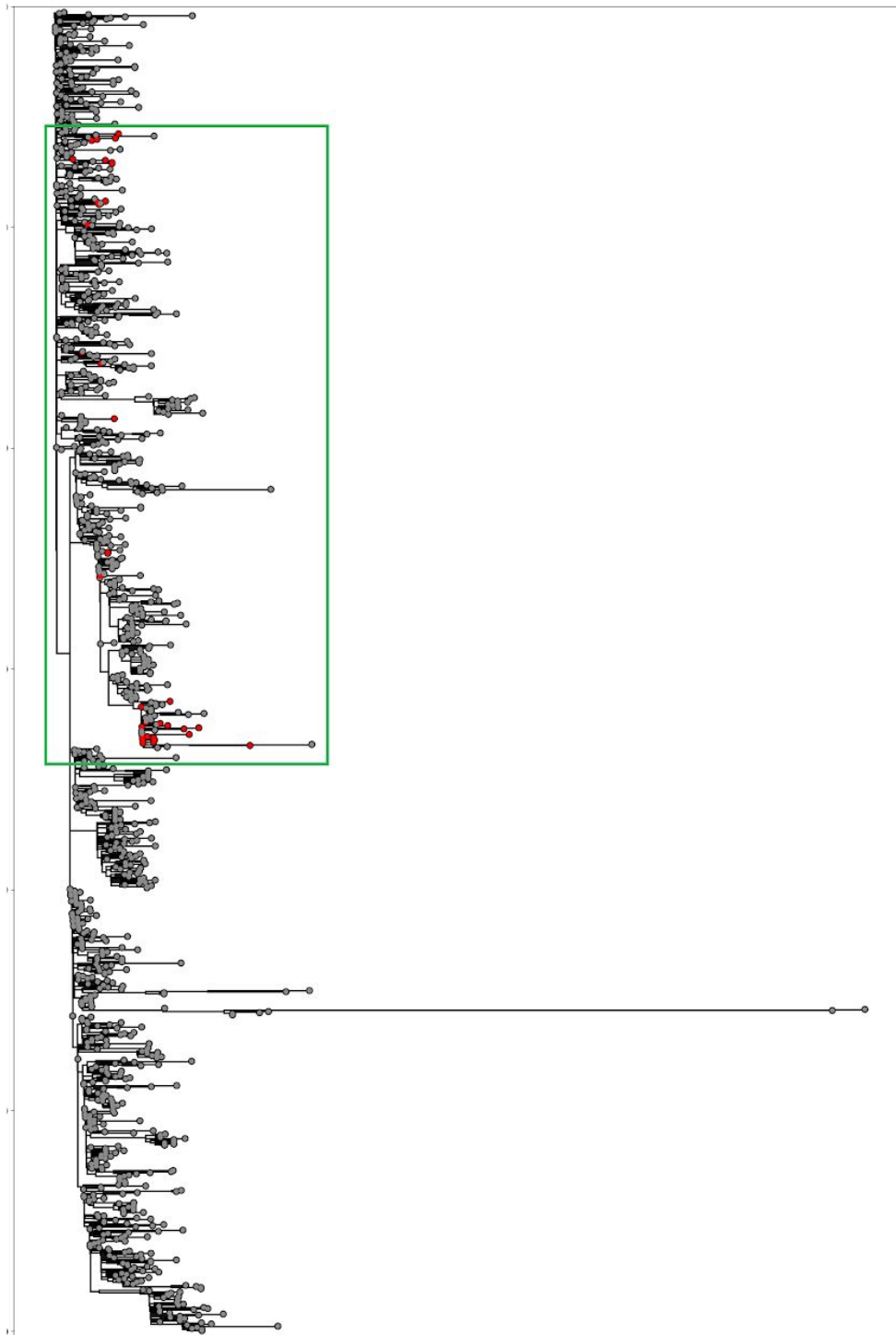


* Green box is what is shown on figure 1.

B. Our generated phylogenetic tree for sequences of Covid-19
   cases from March 22, 2020 and earlier.



* Green box is what is shown on figure 2.
* Extremely long branches are cut, since it would collapse the view of the tree.

C. Phylogenetic tree using Maximum likelihood for sequences
gathered 10 per country and 36 for Philippines



* The red points are the Philippine sequences and the gray points are sequences from other countries.
* Green box is what is shown on figure 3.