**NILE UNIVERSITY**

**Covid 19 PCR diagnostic kits**

A BMD301 Project Report

By

**Dina Yahia Zahran 202001218**

**Nada Osama Fikry 202002277**

**Dina Mohammed Sharwida 202002647**

**Alaa Hussein Mohamed 202001728**

**Fatema Gamal Soliman 202002674**

Submitted in partial fulfillment of the requirements.

for BMD301 Project

**January 30, 2023**

# Abstract

In our experiment, we used three different waves to obtain three sequences of the h-Cov19 for each wave. Waves 1, 2, and 3 occurred between 15 March 2020 and 30 June 2020, 1$^{st}$ September 2020 and 30 April 2021, and 1$^{st}$ July 2021 and 26 September 2021, respectively. The nine sequences were combined into a single file, and multiple sequence alignment was performed on them. We then created our tree to determine which sequences of the three waves are related with each other. NCBI website was used to create our primers as the following step. We selected six consecutive blocks, each of which has full stars at the end and took the first line of each one to represent the conserved region. Then, the conserved region was added to the NCBI Primer-BLAST tool. We received ten primers for the individual region, each with a forward and reverse sequence. To validate our primers, we extracted the forward and reverse sequences for each primer, gave them names, and entered those sequences with their names into a sequence manipulation suite tool. All ten of our primers that were validated had certain warnings, but the best one had only one warning which was that the temperature was over 58 degrees. To find the best conserved region with its best primers, we repeated the conserved region step three times, but each time the temperature warning appeared in the results.

**Table of Contents**

# 1. Introduction

The discipline of bioinformatics has been expanding globally recently. A branch of biology and computer science called bioinformatics deals with the collection, organization, analysis, and communication of biological data, most frequently DNA and amino acid sequences. Therefore, it is a branch of computational science that deals with the examination of biological molecule sequences [4]. Usually, it relates to proteins, DNA, RNA, or genes. In order to aid with biological concerns, the newly growing field of bioinformatics integrates biology, information science, and mathematics. The majority of bioinformatics researchers use Linux. Because it is unrestricted in terms of production, free, quick, secure, and does not slow down with time. Additionally, it includes tens of thousands of free products, does not require additional antivirus protection beyond routine updates, and has a vibrant community that is eager to assist you with any issues you may have. Because of this, Linux is required for genomics and a variety of other practical tasks. Numerous Linux distributions exist; they are referred to as "distros" since Ubuntu Linux is the most widely used. Linux also features a shell or terminal [4]. A shell is a software that takes user instructions, sends them to the OS for processing, and displays the results. Its key component is the Linux shell. Although some of its distributions include a GUI (graphical user interface), Linux primarily has a CLI (command line interface). Your computer's text user interface is called the Linux command line. A computer software designed to understand instructions is sometimes referred to by the names shell, terminal, console, command prompts, and several other names. Command lines are simple to use and may complete tasks in a few lines. They carry out actions in a substantial amount of code lines, unlike any other programming language. It offers the option of autocomplete. You may explore your computer's files and directories using the command line [4]. Nearly all of the operations that can be done

with a GUI can be done using the command line. However, many activities may be completed more quickly and more easily with automation and remote work.

The project's goal is to predict and compare covid 19 sequences across three waves in order to determine which wave has an effect on people and which has a negative impact on DNA sequences. To observe high coverage and complete sequence in genome, we compare sequences using multiple sequence alignment. In addition, to build multiple sequence alignments, we should examine several conserved regions to determine which sequences are related to each other in order to predict which waves have similar properties and therefore have the same effects on genes.

## 2. Methodology

First, we downloaded three sequences for each one of the three waves of hCov-19 from GISAID [2] based on the date of each one of them, where the first wave ranged from 15 March 2020 till 30 June 2020 [6]. While the second wave ranged from 1st September 2020 till 30 April 2021 [5]. Finally, the third wave ranged between 1st July 2021 till 26 September 2021 [1]. We specified the location of the virus to be Africa and we chose the host to be human. Also, we filtered the sequences by choosing only the ones that are complete and with high coverage so that there won't be gaps in the sequences chosen and to have them matched with each other.

*Figure 1: Wave 1*



*Figure 2: Wave 2*

*Figure 3: Wave 3*

## 2.1. Multiple Sequence Alignment

Then we installed muscle tool to perform multiple sequence alignment between nine sequences of the three waves.



Next, we converted the Fasta file including the aligned nine sequences to Clusterlw file to get sequences in blocks and extract the conserved region from them in order to design the primers suitable for this region.

```
dina@dina-VirtualBox:~/BMD301/PROJECT$ muscle -in waves.fasta -out alignment.cl
w -clw

MUSCLE v3.8.1551 by Robert C. Edgar

http://www.drive5.com/muscle
This software is donated to the public domain.
Please cite: Edgar, R.C. Nucleic Acids Res 32(5), 1792-97.

waves 9 seqs, lengths min 29175, max 29836, avg 29684
00:00:01      17 MB(1%)    Iter   1  100.00%  K-mer dist pass 1
00:00:01      17 MB(1%)    Iter   1  100.00%  K-mer dist pass 2
00:03:28   1116 MB(55%)    Iter   1  100.00%  Align node
00:03:28   1116 MB(55%)    Iter   1  100.00%  Root alignment
00:05:12   1116 MB(55%)    Iter   2  100.00%  Refine tree
00:05:12   1116 MB(55%)    Iter   2  100.00%  Root alignment
00:05:12   1116 MB(55%)    Iter   2  100.00%  Root alignment
00:11:56   1116 MB(55%)    Iter   3  100.00%  Refine biparts
00:18:38   1116 MB(55%)    Iter   4  100.00%  Refine biparts
```

## 2.2. Phylogenetic Tree

We also created the phylogenetic tree from the Clusterlw file to view matched and aligned sequences.

```
dina@dina-VirtualBox:~/BMD301/PROJECT$ muscle -in waves.fasta -out tree.clw -cl
w -tree1 tree.phy

MUSCLE v3.8.1551 by Robert C. Edgar

http://www.drive5.com/muscle
This software is donated to the public domain.
Please cite: Edgar, R.C. Nucleic Acids Res 32(5), 1792-97.

waves 9 seqs, lengths min 29175, max 29836, avg 29684
00:00:00      17 MB(1%)    Iter   1  100.00%  K-mer dist pass 1
00:00:00      17 MB(1%)    Iter   1  100.00%  K-mer dist pass 2
00:03:23   1116 MB(55%)    Iter   1  100.00%  Align node
00:03:23   1116 MB(55%)    Iter   1  100.00%  Root alignment
00:05:06   1116 MB(55%)    Iter   2  100.00%  Refine tree
00:05:06   1116 MB(55%)    Iter   2  100.00%  Root alignment
00:05:06   1116 MB(55%)    Iter   2  100.00%  Root alignment
00:11:50   1116 MB(55%)    Iter   3  100.00%  Refine biparts
00:19:01   1116 MB(55%)    Iter   4  100.00%  Refine biparts
dina@dina-VirtualBox:~/BMD301/PROJECT$
```

Phylogenetic tree (newick) viewer is a tool that was used in the project for visualizing the phylogenetic tree as a dendrogram to make the interpretation process of the tree easier.

## 2.3. Conserved Regions

We specified three different conserved regions from different consecutive blocks in Clusterlw file to see which one will have best primers. For the first conserved region, we chose six lines consecutively from the file to represent the conserved region of the multiple sequences and put it in a separated text file and selected from a random original sequence 20 bases before the first line chosen and 20 bases after that represent forward and reverse bases.



We performed the same steps for specifying the other two conserved regions and designing their own primers so that at the end, we can compare which resulted primers are the best.
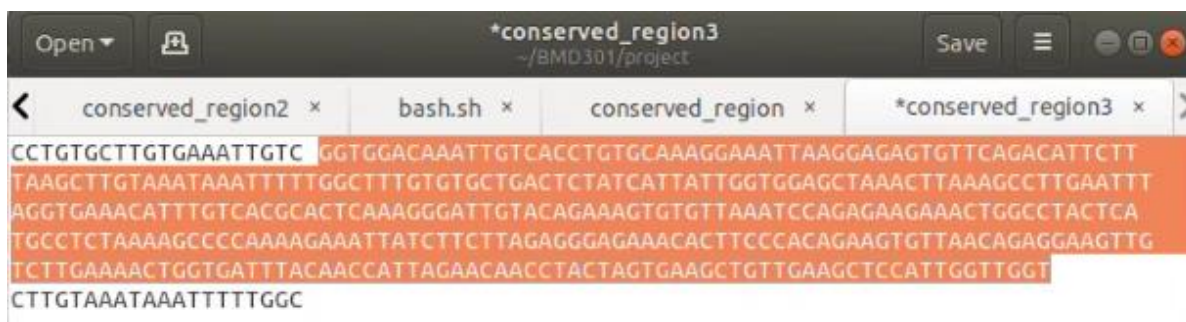
*Figure 5: Conserved Region 2*



*Figure 4: Conserved Region 3*

After specifying the first conserved region, we put it on NCBI Primer-BLAST [7] to design 10 primers for the specified conserved region.

## 2.4. Automated Bash Script

At the end, we automated the multiple sequence alignment, phylogenetic tree and primer design parts using bash script, where the user is asked to enter the specific file, put a name for the output file and choose the process he/she wants to perform on the sequence, in order to make the process more generic and not specific for a certain file. This was done using "nano bashscript.sh" and writing the commands for performing such processes inside it.

```
  GNU nano 6.2                        bashfile.sh *
echo "Enter the file: "
read file
echo -n "1-MSA     2-Tree    3-primer  :  "
read number
echo "Write the name of the output file: "
read fileoutput
if [[ $number == 1 ]]
then
   muscle -in $file -out $fileoutput.clw -clw
elif [[$number == 2]]
then
   muscle -in $file -out $fileoutput.clw -clw -tree1 $fileoutput.phy
else
  primer3_core  $file  > $fileoutput.txt
  cat $fileoutput.txt
fi
```

# 3. Results

## 3.1. Phylogenetic Tree

The following figure represents the phylogenetic tree produced from the nine sequences of the three different waves. It was found that there are two different sequences from two different waves are similar to each other which are the two sequences from Libya in wave one and two, where both of them are complete and of high coverage. The rest of the similar sequences are in the same waves but might be from different countries.

The output of the Phylogenetic tree (newick) viewer tool [8] is represented in the following figure.



## 3.2. Multiple Sequence Alignment

We use command 'cat' to visualize all aligned sequences in file alignment.clw to extract the conserved regions for designing the primers.
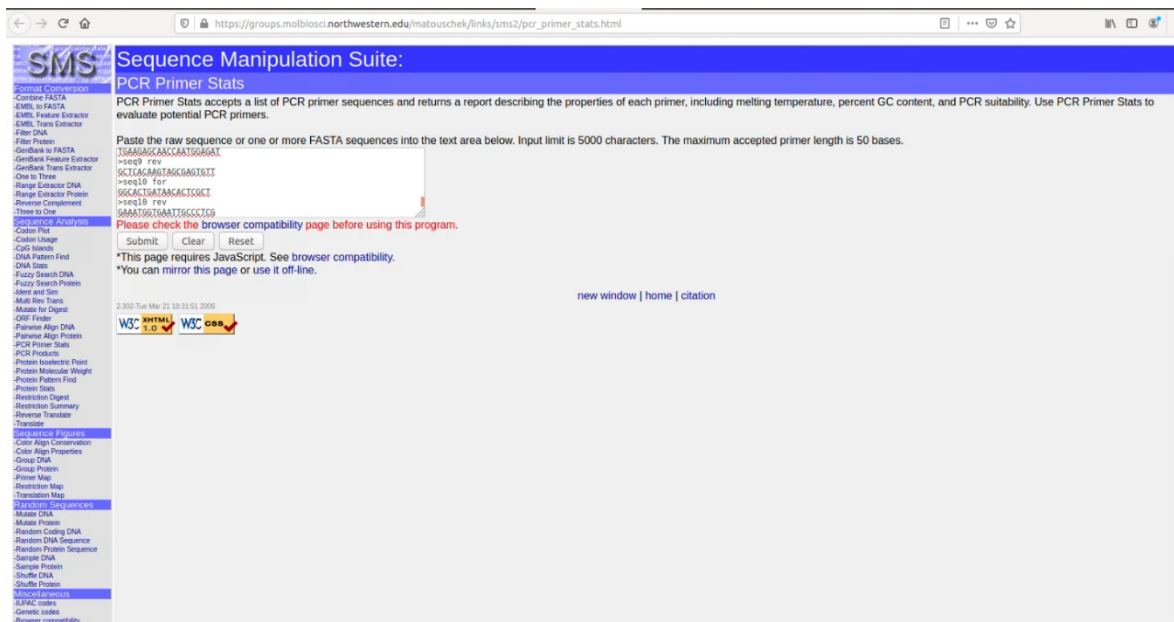
## 3.3. Conserved Regions

According to the three conserved regions selected from the aligned sequence file, we designed 10 primers for each region.

We then used Sequence Manipulation Suite tool (PSR Primer State) to validate the 10

primers of each conserved region and compare between them.



Although we chose different conserved regions throughout the multiple sequence

alignment of the nine sequences, the primers designed for each one of them had the same results

in all the criteria they are evaluated upon, where the best primer in each had only warning in the

temperature part saying that "Tm is greater than 58", while the rest of the criteria such as GC

clamp, self-annealing, and hairpin formation were good without any warnings.

## 3.4. Automated Bash Script

In the bash script part, if the user chose the option of performing multiple sequence

alignment on the file he/she desires, the output will be as the following resulting in a "file.afa".

```
dina@dina-VirtualBox:~/BMD301/PROJECT$ bash bashfile.sh
Enter the file:
wave2.fasta
1-MSA     2-Tree    3-primer   :  1
Write the name of the output file:
msa_trial

MUSCLE v3.8.1551 by Robert C. Edgar

http://www.drive5.com/muscle
This software is donated to the public domain.
Please cite: Edgar, R.C. Nucleic Acids Res 32(5), 1792-97.

wave2 3 seqs, lengths min 29175, max 29789, avg 29553
00:00:00      16 MB(1%)   Iter   1   100.00%  K-mer dist pass 1
00:00:00      16 MB(1%)   Iter   1   100.00%  K-mer dist pass 2
00:00:00      41 MB(2%)   Iter   1    50.00%  Align node
```

On the other hand, if option two was chosen, which is making a phylogenetic tree for the

specified sequence, the output will be the following figure resulting in a "file.phy".

```
dina@dina-VirtualBox:~/BMD301/PROJECT$ bash bashfile.sh
Enter the file:
wave2.fasta
1-MSA     2-Tree    3-primer   :  2
Write the name of the output file:
tree_trial

MUSCLE v3.8.1551 by Robert C. Edgar

http://www.drive5.com/muscle
This software is donated to the public domain.
Please cite: Edgar, R.C. Nucleic Acids Res 32(5), 1792-97.

wave2 3 seqs, lengths min 29175, max 29789, avg 29553
00:00:00      16 MB(1%)   Iter   1   100.00%  K-mer dist pass 1
00:00:00      16 MB(1%)   Iter   1   100.00%  K-mer dist pass 2
^C:00:00      41 MB(2%)   Iter   1    50.00%  Align node
```

Finally, if option three, which is designing the primer for the conserved region, is chosen then the output will be a number of primers designed for the region including the forward and reverse strands of each one of them as well as other parameters for each primer.



# 4. Conclusion

After we made alignment, we concluded that there are multiple sequences similar to each other in the same wave but in different countries. In addition, when we specified the conserved regions from the multiple sequence alignment, it was also observed that the conserved region is considered as similar sequences in same wave but in different countries. Depending on that, when we made multiple primers for different conserved regions, this gave the same result, so we conducted that the alignment was done correctly.

# 5. References

[1] "Africa Faces Steepest COVID-19 Surge Yet." WHO | Regional Office for Africa,

www.afro.who.int/news/africa-faces-steepest-covid-19-surge-yet. Accessed 30 Jan. 2023.

[2] Epicov.org, 2023, www.epicov.org/epi3/frontend#4c9d24.

[3] https://www.bioinformatics.org/sms2/

[4] National Human Genome Research Institute. "Bioinformatics." Genome.gov, 6 Sept. 2022,

www.genome.gov/genetics-glossary/Bioinformatics.

[5] "New COVID-19 Variants Fuelling Africa's Second Wave." WHO | Regional Office for

Africa, www.afro.who.int/news/new-covid-19-variants-fuelling-africas-second-wave.

[6] "PCR Primer Stats." Groups.molbiosci.northwestern.edu,

groups.molbiosci.northwestern.edu/matouschek/links/sms2/pcr_primer_stats.html.

Accessed 30 Jan. 2023.

[7] "Primer Designing Tool." Nih.gov, 2019, www.ncbi.nlm.nih.gov/tools/primer-blast/

[8] "Tree Viewer - Online Visualization of Phylogenetic Trees (Newick) and Alignments."

Etetoolkit.org, etetoolkit.org/treeview/.

# 6. Team Members Contributions

- **Fatema Gamal Soliman:** Downloaded the sequences from different waves and performed multiple sequence alignment.

- **Dina Mohammed Sharwida:** Made the phylogenetic tree for the nine sequences downloaded and analyzed its results.

- **Dina Yahia Zahran:** Extracted conserved regions from multiple sequence alignment file and designed primers for such regions.

- **Alaa Hussein Mohamed:** Performed Primer validation among different conserved regions and compared between their results.

- **Nada Osama Fikry:** Automated bash script for multiple sequence alignment, phylogenetic tree and primer designing.