Comparison of the complete genome sequences of COVID-14 and Zika virus using BioPython

Ahmad Hafezi', Zahra khamar'*

Department of Biology, Faculty of science, Ferdowsi University of Mashhad, Mashhad, Iran

*Corresponding author: <u>zahrakhamr@um.ac.ir</u>

Abstract

Background: Bioinformatics has always been one of the important tools in biological research. Increasing in the biological data, triggered some difficulties in the data analysis. To solve these problems BioPython developed, and created a better platform for researchers to conduct their study in a short way just by writing appropriate BioPython codes.

Methods: In this research, bioinformatics sites such as NCBI were first used so that the genome sequences of both Zika and Corona viruses were completely stored in the Festa file. The PyCharm program was then used to apply BioPython codes and facilitate the speed of data analysis. Finally, the sequences were analyzed with BioPython codes to determine the GC percentage, length and number of amino acids that each sequence encodes.

Results: The results showed that, although the zika virus has a shorter DNA sequence, the amino acids ratio encoded by the whole genome is almost the same as COVID-19 and the CG content of zika is more than coronavirus. After analyzing the sequence of both viruses, it was found that the three amino acids Thr, Ser and Leu are the three amino acids that both viruses had the most codons of their production in their genomes. In contrast, in both sequences, the amount of Met encoded by the entire genome is in fewer amounts.

Conclusions: In this research, the COVID-14 and Zika virus whole genomes are analyzed with BioPython and compared with each other. As a result, both viruses have historically been major threats to human health, so knowing enough and having the necessary information about their genomes will enable scientists to find the right drug or vaccine to deal with them in a shorter time. In this case, obtaining the desired information in the shortest possible time using programming can improve the speed of work and help to perform calculations in a shorter time. In the meantime, Python programming enables us to obtain valuable

https://bcs.cdsts.ir Page \



سومین کنگره توسعیی علمیی و فنیاوری دانشیجویان زیستشناسی و شیمی Srd Congress of Scientific and Technological Development of Biology and Chemistry Students

information about the desired strain in the shortest possible time in the next pandemic.
Keywords: BioPython, Data Analyzing, Bioinformatics, Sequence analyzing, COVID-19

https://bcs.cdsts.ir Page Y

1. Background

It was in 1991 when the first version of python published by Guido van Rossum (1). After that Brad Chapman and Jeff Chang started to utilize the python for biological research and they published the BioPython for the first time on $Y \cdot \cdot \cdot (Y)$. Python is a programming language that is extensively in academics and commercials and due to its user-friendly and its performance in all types of operating systems features, it has become the most well-known language programming in the world(*). BioPython is an open-source compilation of Python tools for natural computation, made by a universal group of developers(Υ, \mathcal{E}). **BioPython** has been significantly expanded during this century because using the appropriate modules and classes in biological analysis, enables us to solve complex bioinformatics problems. BioPython provides diverse types of tools to bioinformatics files and parse accessibility to online websites such as NCBI(o).

1.1. COVID-19

Severe acute respiratory syndrome coronavirus (SARS-CoV-Y) emerged in Y-19 in China. The coronaviridae family consist of four genera: alphacoronavirus,

betacoronavirus, gammacoronavirus and deltacoronavirus that can infect mammals and birds. It is known that alpha- and betacoronaviruses usually infect humans, whereas gamma- and deltacoronaviruses

infect birds. Coronavirus is a positive sense RNA that the length of RNA is approximately **.kb. The three types of proteins found in coronavirus include spike (s), membrane (M) and envelope (E) (Figure 1). The S protein contains two domains, s1 and sy, which are essential for entry into host cells. A cell surface peptidase called ACEY hydrolyzes angiotensin II and is naturally high-expressed in the epithelium of the lung and small intestine recognized as a receptor for coronavirus. after binding ACEY and s protein, s protein is subsequently cleaved and activated by another cell surface protease MPRSSY and it leads to membrane fusion (٦).

https://bcs.cdsts.ir Page **T**



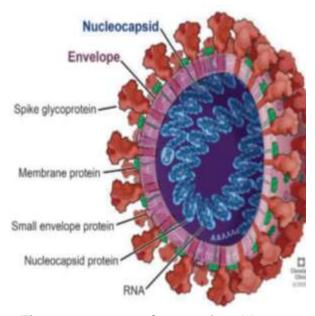


Figure \: structure of coronavirus (\v)

1.1 Zika virus

In 1967, the Zika virus was first isolated in a rhesus monkey (λ). Seven years later, the first human infection emerged in Nigeria (9). Zika virus was isolated from some mosquitoes in Africa when researchers studied yellow fever (λ ,1.). The first official report of zika fever occurred in $\Upsilon \cdot \cdot \cdot \Upsilon$ (11) and a considerable epidemic emerged in French Polynesia in the South Pacific in $\Upsilon \cdot \cdot \Upsilon$ and $\Upsilon \cdot \cdot \cdot \Upsilon$ (17). It was evaluated that approximately $\Upsilon \cdot \cdot \cdot \cdot \Upsilon$ symptomatic infections were recorded (1 Υ ,1 Υ). This epidemic outbreak in $\Upsilon \cdot \cdot \Upsilon \in \Upsilon$ in New

Caledonia and initially disseminated in the Americas (Brazil in March)(10).

Y. Material and Methods

To have a better platform for running the BioPython codes, the PyCharm was downloaded the https://www.jetbrains.com/pycharm/ site website. Afterwards, the whole genome of respiratory severe acute syndrome coronavirus 7 and Zika virus were download from https://www.ncbi.nlm.nih.gov/ website with $ID = NC_{\cdot} \epsilon \circ 017.7$ and $KX \wedge 97 \wedge 00.1$ respectively. These FASTA-sequence files were transferred to the PyCharm directory and used the below codes to analyze these sequences (\7):

>>>from Bio.SeqRecord import SeqRecord

>>>from Bio.SeqIO import parse

>>>from Bio.Seq import Seq

>>>COVID_seq =

open("FASTA_sequence.fasta")

>>>detail = parse(COVID_seq, "fasta")

>>>for i in detail:

>>> print("ID: %s " % i.id)

>>> print("Name: %s " % i.name)

https://bcs.cdsts.ir Page £

سومین کنگره توسعی علمی و فنیاوری دانشیجویان زیستشناسی و شیمی Srd Congress of Scientific and Technological Development of Biology and Chemistry Students

```
>>> print("Description: %s " % i.description)
```

>>> print("Sequence Data: %s " % i.seq)

This code separates the ID, Name, Description and the main sequence in the FASTA file. By falling into different parts of the FASTA component, it becomes easy to parse and analyze the main sequence. To find out the GC percentage in this sequence it is necessary to import the GC class from Bio.SeqUtils module.

```
>>>from Bio.SeqUtils import GC
```

>>> main_seq = []

>>>for i in detail:

>>> main_seq = i.seq

>>>print(GC(main_seq))

To transcription this sequence, the following code was used. After that to translate this sequence the translate() should be applied. Also, the length of both RNA and amino acids sequences will be appeared by applying the len() command:

```
>>>from Bio.Seq import Seq
>>>from Bio.Seq import transcribe
>>>COVID_seq =
open("FASTA_sequence.fasta")
>>>detail = parse(COVID_seq, "fasta")
>>>RNA = []
>>> main_seq = []
>>> for i in detail:
>>> main_seq = i.seq
```

```
>>> RNA = transcribe(main_seq)
>>>print(RNA)
>>>print(len(RNA))
>>>print(RNA.translate())
>>>print(len(RNA.translate()))
```

To indicate the standard codons table that uses for translation the CodonTable must be imported from Bio.Data:

```
>>>from Bio.Data import CodonTable
>>>codon_table =
CodonTable.ambiguous_dna_by_name['Stan dard']
>>>print(codon_table)
```

Next, the number of each amino acids were calculated by writing below codes:

>>>from Bio.Data import IUPACData

>>>amino acids =

IKGLYLPR*QTNQLSISCRSVL*TNFKIC VAVTRLHA*CTHAV*LITNYCR*QDTS NSSIFCRLL...MCKINFSSAIPM*F**LLR RMTKKKKKKKKKK

```
>>>type_amino =
IUPACData.protein_letters
>>>count_amino =
{}.fromkeys(type_amino, .)
>>>stop_codon = {}.fromkeys('stop_codon', .)
>>>for char in amino_acids:
>>> if char in type_amino:
>>> count_amino[char] += .
```

https://bcs.cdsts.ir Page 4

سومین کنگره توسعی علمی و فنیاوری دانشیجویان زیستشناسی و شیمی Signature Congress of Scientific and Technological Development of Biology and Chemistry Students

>>>for i in amino_acids:

>>>stop_codon=amino_acids.count('*')

>>>print(count_amino)

>>>print(IUPACData.protein_letters_\text{"to\)

>>>print('stop_codon:',stop_codon)

All these steps were done for Zika virus genome. These codes were run in the PyCharm Console and then the results were checked.

r. Results

In the first, the Codon Table was shown to understand what three types of basses in mRNA create a specific type of amino acids. In this table, four types of basses (T, C, A and G) are listed and the type of amino acid that will be created is shown beside of each the three basses with a special letter (Figure

r). It is helpful to have an abbreviation of the name of twenty types of amino acids. Therefore, by using the IUPACData.protein_letters_rto\,

all of the amino acids were shown with their abbreviation letters:

{'Ala': 'A', 'Cys': 'C', 'Asp': 'D', 'Glu': 'E', 'Phe': 'F', 'Gly': 'G', 'His': 'H', 'Ile': 'I', 'Lys': 'K', 'Leu':

'L', 'Met': 'M', 'Asn': 'N', 'Pro': 'P', 'Gln': 'Q', 'Arg': 'R', 'Ser': 'S', 'Thr': 'T', 'Val': 'V', 'Trp': 'W', 'Tyr': 'Y'}.

https://bcs.cdsts.ir Page 1



```
C | A | G
T | TTT F | TCT S | TAT Y | TGT C
                               T
T | TTC F | TCC S | TAC Y | TGC C
T | TTA L | TCA S | TAA Stop | TGA Stop | A
T | TTG L(s) | TCG S | TAG Stop | TGG W | G
  C | CTT L | CCT P | CAT H | CGT R | T
C | CTC L | CCC P | CAC H | CGC R | C
C | CTA L | CCA P | CAA Q | CGA R | A
C | CTG L(s) | CCG P | CAG Q | CGG R | G
  A | ATT I | ACT T | AAT N | AGT S | T
A | ATC I | ACC T | AAC N | AGC S | C
A | ATA I | ACA T | AAA K | AGA R | A
A | ATG M(s) | ACG T | AAG K | AGG R | G
  G GTT V GCT A GAT D GGT G T
G | GTC V | GCC A | GAC D | GGC G | C
G | GTA V | GCA A | GAA E
                       GGA G A
G | GTG V | GCG A | GAG E | GGG G | G
```

Figure 7: The standard genetic codons demonstrated in the PyCharm console (17)

r.1. Covid 19 sequence analysis

In the first written code, all information which is available in the FASTA file, will be shown separately in the PyCharm console. ID and name indicate the specific identification which exclusively belongs to the COVID-19 sequence. The main information indicated in the description part gives complete data about the sequences. Finally, the sequence is imported into a separate part.

ID: NC . 20017.7

Name: NC_ · ٤00 \ Y. Y

Sequence Data:

ATTAAAGGTTTATACCTTCCCAGGTA ACAAACCAACCAACTTTCGATCTCTT GTAGATCTGTTCTCTAAACGAACTTT AAAATCTGTGTGGCTGTCACTCGGCT GCATGCTTAGTGCATTTCCTATTCCTT AGGGAACGTGGTTGACCTACACAGGT

https://bcs.cdsts.ir Page V

GC content was evaluated by using BioPython according to the main sequence. For COVID-19, GC% = TV.4VTVVATO.4V1&A was estimated. BioPython used the transcribe module to replace the U

nucleotide with the T nucleotide in the main sequence. And the whole sequence was completely translated. Also, the len() was used to estimates the length of the COVID-14 genome.

Length = 799.% (bp)

After transcribing, all of the amino acids that are encoded by the whole genome were determined and its length was measured as well. In addition, the number of each amino acid was revealed beside the number of stop codons (*) in the PyCharm console (Figure *):

IKGLYLPR*QTNQLSISCRSVL*TNFKIC VAVTRLHA*CTHAV*LITNYCR*QDTS NSSIFCRLL...MCKINFSSAIPM*F**LLR RMTKKKKKKKKKK

Length= 997V amino acids

stop_codon: ٧٧٤

{'A': ٣٧0, 'C': ٦٣0, 'D': ٢٩٠, 'E': ٢٧٠, 'F': 09٣,

'G': ٣٩٤, 'H': ٣٣٢, 'I': ٤٣٦, 'K': ٤١٣, 'L': ٨٨٦,

'M': 11V, 'N': EVT, 'P': T9T, 'Q': TTO, 'R': OOA,

'S': AI., 'T': 7/9, 'V': 0EA, 'W': 774, 'Y': 0.0}.

https://bcs.cdsts.ir Page A



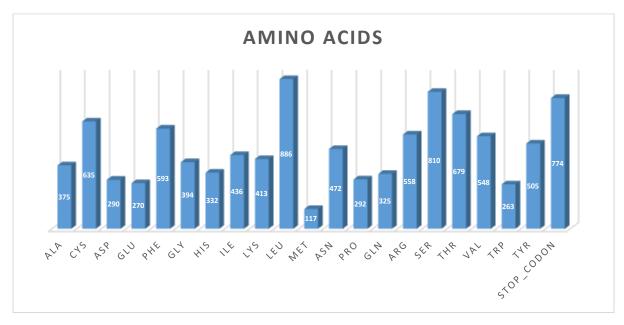


Figure r: The number of each amino acid and stop codons transcribed from the COVID-19 whole genome

r.r. Zika virus sequence analysis

ID, Name, description and Sequence data in the FASTA file of the Zika virus fell into different sections and the main sequence was analyzed more by BioPython.

ID: KXA9TA00.1

Name: KXA9TA00.1

Description: KX \ 9\% \ 2ika virus strain

Zika virus/Homosapiens/VEN/UF-۲/۲۰۱٦,

complete genome

Sequence Data:

AGTTGTTACTGTTGCTGACTCAGACT GCGACAGTTCGAGTTT...GGGGAAATC CATGGGTCTT Transcribe:

AGUUGUUACUGUUGCUGACUCAGAC UGCGACAGUUCGAGUUU...GGGGAAA UCCAUGGGUCUU

Length = $1 \cdot A \cdot A$ (bp)

%GC = 01.717.7017797071

Translate:

SCYCC*LRLRQFEFEAKASNSINRFYFG FGNESFWS*KTQKRNPEDSGLSIC...QK RD*WLEETPRKTQNSILTLGKTRDSMS FHHAGRQAQIAE*RRPVWGNPWV.

Length = ٣٦٠٢

stop_codon: ۲۱۳

https://bcs.cdsts.ir Page 9

{'A': \n, 'C': \rv, 'D': \r, 'E': \n, 'F': \r, 'G': \r, 'H': \equiv , 'I': \l, 'K': \q, 'L':

ΥΥΥ, 'M': Α٩, 'N': ٥٥, 'P': ΥξΨ, 'Q': ١٩٦, 'R':
ΥΛξ, 'S': ΨΥο, 'T': Υ٦٠, 'V': ١ΥΨ, 'W': ١٧Λ,
'Y': ΨΥ} (Figure ξ).

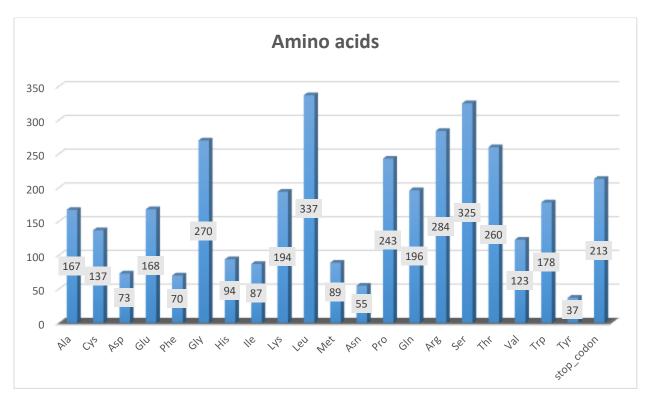


Figure £: The number of each amino acid transcribed from the Zika virus whole genome

By analyzing both genomes, the result showed that Thr, Ser and Leu are three of the most common amino acids in both sequences. Although Tyr and Asn have fewer codons in the Zika whole genome, they are two of the more typical amino acids in the COVID-19 genome. Moreover, in both sequences, the amount of Met that is encoded by the whole genome is in fewer

amounts. Another common feature in these results is related to the stop codons. As it is indicated, the number of stop codons in both genomes is in a high amount. Even though the length of the Zika genome is approximately three times shorter than that of the COVID-14 genome, the ratio of amino acids encoded by the complete Zika genome is almost identical to the ratio of amino acids

https://bcs.cdsts.ir Page 1 ·

encoded by the COVID-19 complete genome.

predict the strains that can trigger the infection in advance. In this way, BioPython has become an open-source programming application that has tremendous potential for bioinformatics researchers.

٤. Conclusions

COVID-19 and Zika are two viruses which triggered a large pandemic in the world. Therefore, it is necessary to analyze their genomes to extract the information that can help us to find an appropriate view to be prepared for the next pandemic. Today, BioPython has become one of the state-ofthe-art language programming which can parse and analyze bioinformatics files in a short time. In this regard, Kukreja and Kumari conducted research to analyze the data related to brain cancer with BioPython. In the first step, they used the KEGG website to find the interaction between Glioma molecules. After the selection of related molecules, the GC content and the length of these genes were assessed (o).

Overall, by parsing the sequence, the disease prediction will be done easily. For example, it is possible to analyze what happens if a nucleotide changes in a specific part of the genome with this awareness, some diseases can be treated early in diagnosis. Furthermore, by analyzing the infection strains genome it would be possible to

References

- Van Rossum G. A brief timeline of Python. Hist Python. Y., 9;
- Y. Chapman B, Chang J. Biopython: Python tools for computational biology. ACM Sigbio Newsl.

- Oliphant TE. Python for scientific computing. Comput Sci Eng.
 Y··V; ٩(٣): Y·-Y·.
- ¿. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely available Python tools for computational molecular biology and bioinformatics. Bioinformatics. Y: ٩:٢٥(١١):١٤٢٢.
- •. Kukreja V, Kumari U. Data Analysis of Brain Cancer with Biopython.
- The molecular virology of coronaviruses. J Biol Chem.
- V. Bergmann CC, Silverman RH. COVID-19: Coronavirus replication, pathogenesis, and therapeutic strategies. Cleve Clin J Med.

https://bcs.cdsts.ir Page \\\

سومین کنگره توسعی علمی و فنیاوری دانشیجویان زیستشناسی و شیمی Signature Congress of Scientific and Technological Development of Biology and Chemistry Students

۲۰۲۰;۸۷(٦):۳۲۱<u>–</u>۷.

- A. Dick GWA, Kitchen SF, Haddow AJ. Zika virus (I). Isolations and serological specificity. Trans R Soc Trop Med Hyg. 1907; £7(0):0.9-7.
- Macnamara FN. Zika virus: a report on three cases of human infection during an epidemic of jaundice in Nigeria.

Trans R Soc Trop Med Hyg. $190\xi:\xi\lambda(7):179=\xi\circ$.

- Marchette NJ, Garcia R, Rudnick A. Isolation of Zika virus from Aedes aegypti mosquitoes in Malaysia. Am J Trop Med Hyg. ۱۹۲۹; ۱۸(۳).
- Duffy MR, Chen TH, Hancock WT, Powers AM, Kool JL, Lanciotti RS, et al. Zika virus outbreak on Yap Island, federated states of Micronesia. N Engl J Med. ۲۰۰۹; ۳٦٠ (٢٤): ٢٥٣٦— ٤٣.
- Y. Cao-Lormeau VM, Roche C, Teissier A, Robin E, Berry AL, Mallet HP, et al. Zika virus, French polynesia, South pacific, Y.Y. Emerg Infect Dis. Y.Y: Y.(1):Y.Ao.
- Musso D, Nilles EJ, Cao-Lormeau VM. Rapid spread of emerging Zika virus in the Pacific area. Clin Microbiol Infect. Y · \(\xi \xi; \Y \cdot(\Y) \: O \(\quad \cdot \cdot \)
- Musso D, Gubler DJ. Zika virus. Clin Microbiol Rev. ۲۰۱٦;۲۹(۳): ۱۹۸۷-۱۶۰۶.
- Campos GS, Bandeira AC, Sardi SI. Zika virus outbreak, bahia, brazil. Emerg Infect Dis. Y. 10; Y. 10. 11.

NT. Biopython - Quick Guide [Internet]. Available from: https://www.tutorialspoint.com/biopython/biopython_quick_guide.htm

https://bcs.cdsts.ir Page \Y