# A comparative analysis for SARS-CoV-2

Göksel Mısırlı*

*School of Computing and Mathetmatics, Keele University*

E-mail: g.misirli@keele.ac.uk

## Abstract

COVID-19 has affected the world tremendously. It is critical that biological experiments and clinical designs are informed by computational approaches for time- and cost-effective solutions. Comparative analyses particularly can play a key role to reveal structural changes in proteins due to mutations, which can lead to behavioural changes, such as the increased binding of the SARS-CoV-2 surface glycoprotein to human ACE2 receptors. The aim of this report is to provide an easy to follow tutorial for biologists and others without delving into different bioinformatics tools. More complex analyses such as the use of large-scale computational methods can then be utilised. Starting with a SARS-CoV-2 genome sequence, the report shows visualising DNA sequence features, deriving amino acid sequences, and aligning different genomes to analyse mutations and differences. The report provides further insights into how the SARS-CoV-2 surface glycoprotein mutated for higher binding affinity to human ACE2 receptors, compared to the SARS-CoV protein, by integrating existing 3D protein models.

# Keywords

SAR-CoV-2, SARS-CoV, Comparative Genomics, Surface Glycoprotein, CLC Genomics Workbench

1

# Introduction

Severe acute respiratory syndrome (SARS)-related coronaviruses have previously caused two pandemics.[1] The recent SARS Coronavirus 2 (SARS2-CoV-2) outbreak has had unprecedented effects so far. It is essential to develop data integration mechanisms in order to gain insights using data that already exists. For example, genome-wide comparisons can be used to inform subsequent computational analyses which can potentially be used to search for drugs and to develop computational models.

Taxonomy classifications offer a systematic approach to link data from different organisms. The taxonomy id for SARS-CoV-2 is reported as 2697049 by the National Center for Biotechnology (NCBI) taxonomy browser[2] which also lists synonyms that can be used when searching for information in different databases (`https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?id=2697049`). Moreover, the taxonomy id 694009 is used to group all SARS-related coronavirus taxonomies. This list can especially be useful to search for related sequences in order to infer information via comparative genomics approaches. For example, the NCBI Virus[3] database can be queried with these taxonomy ids to retrieve SARS-CoV-2 related nucleotide and protein sequences.

The SARS-CoV-2 genome sequences can also be accessed directly from the NCBI's GenBank database (`https://www.ncbi.nlm.nih.gov/genbank/sars-cov-2-seqs`). Related information includes where a virus was isolated from. Some of these GenBank files include only sequences and some of them provide annotations about genome locations of important sequence features.

This tutorial initially uses two genome sequences: one for SARS-CoV-2 and another one for SARS-CoV. Regarding the former, the GenBank entry LC528232 was selected. LC528232 was created for a strain which was isolated in Japan. Regarding the latter, the GenBank entry AP006557, which was deposited in 2006, was selected.

# Genomic features

The GenBank files include annotations about sequence features, which can be visualised and analysed using various tools. In this report, the analyses were carried out using CLC Genomics Workbench 20.0.3. GenBank files were initially imported into this tool which was then used to analyse the annotated sequence features in more detail. For example, the surface glycoprotein denoted as 'S' is shown in Figure 1.
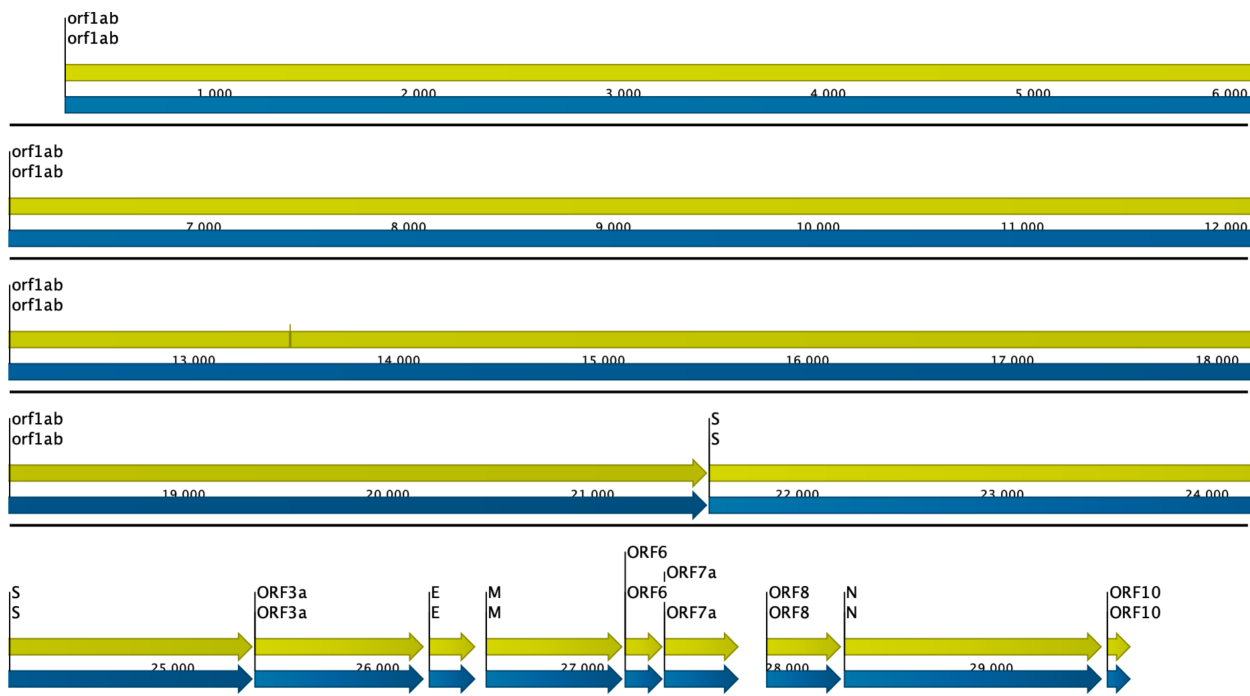


Figure 1: SARS-CoV-2 sequence features that are captured in GenBank files are shown.

In order understand the mutations in the surface glycoprotein, a new entry for the related coding sequence (CDS) was created in CLC Genomics Workbench. The corresponding CDS feature was selected and the nucleotides were copied to create this new entry. CLC Genomics Workbench was then used to display restriction sites and protein translation information (Figure 2).
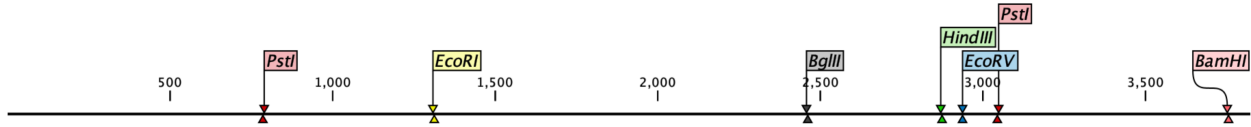
Figure 2: Restriction sites in the surface glycoprotein CDS, labelled as 'S'.

# Predicting secondary structures

Understanding secondary structure formations and where mutations occur can provide insights into the effects of these mutations in the SARS-CoV-2 surface glycoprotein CDS. Hence, a new entry for the corresponding amino acid sequences was created in CLC Genomics Workbench by using the '`Toolbox – Classical Sequence Analysis – Nucleotide Analysis – Translate to Protein`' option. Alpha-helix and beta-strand formations were predicted using the '`Toolbox – Classical Sequence Analysis – Protein Analysis – Predict Secondary Structure`' option. The inferred information was then incorporated as annotations into the entry for the amino acid sequences (Figure 3).
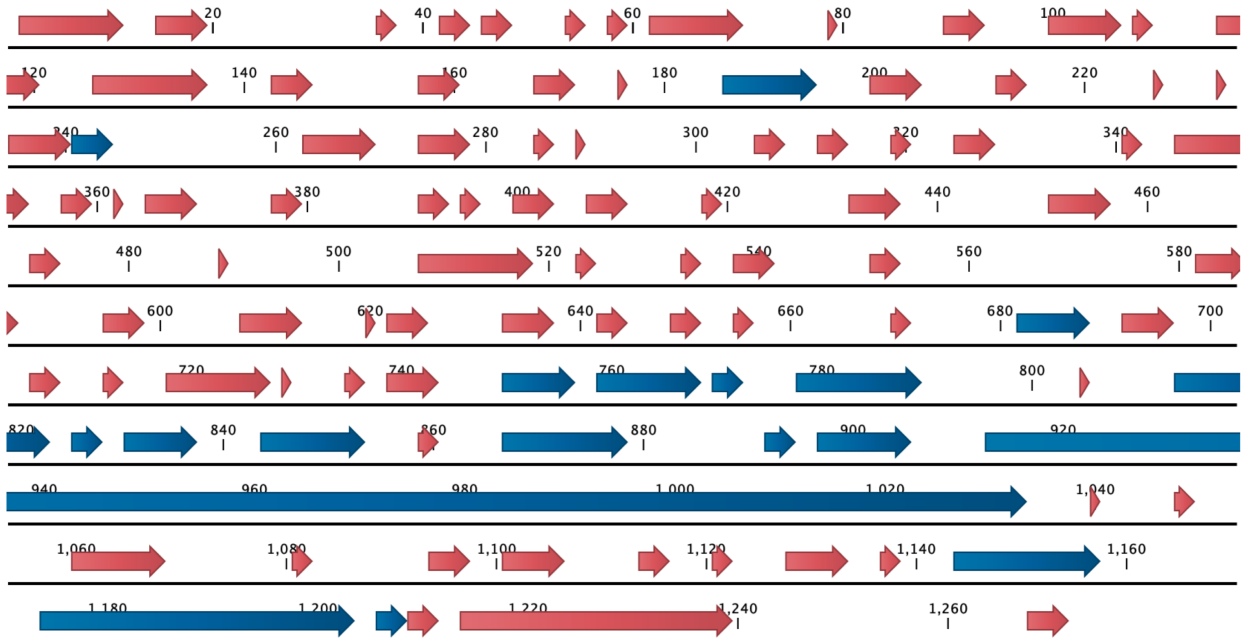


Figure 3: The surface glycoprotein secondary structure predictions. Red arrows represent beta-strands and blue arrows represent alpha-helices.

4

## Sequence alignment

It is reported that the SARS-CoV-2 surface glycoprotein is optimised to bind to human ACE2 receptors.[4] In order to analyse the effects of these mutations further, the protein sequence can be aligned to previously known similar sequences. Figure 4 shows the alignment of the surface glycoprotein amino acid sequences from SARS-CoV-2 (LC528232) and SARS-CoV (AP006557). Although the alignment shows high similarities between the amino acid sequences, gaps and mutations also exist.
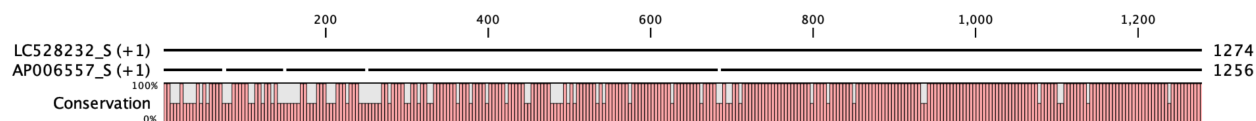


Figure 4: Aligned surface glycoprotein amino acid sequences from SARS-CoV-2 (LC528232) and SARS-CoV (AP006557).

A more detailed view of the alignment of the two protein sequences can be seen in Figure 5. Secondary structures are integrated into the view. The surface exposed regions are compared using different options such as the Kyte-Doolittle scale. Additional options can be configured from the 'Alignment Settings – Protein info' tab.
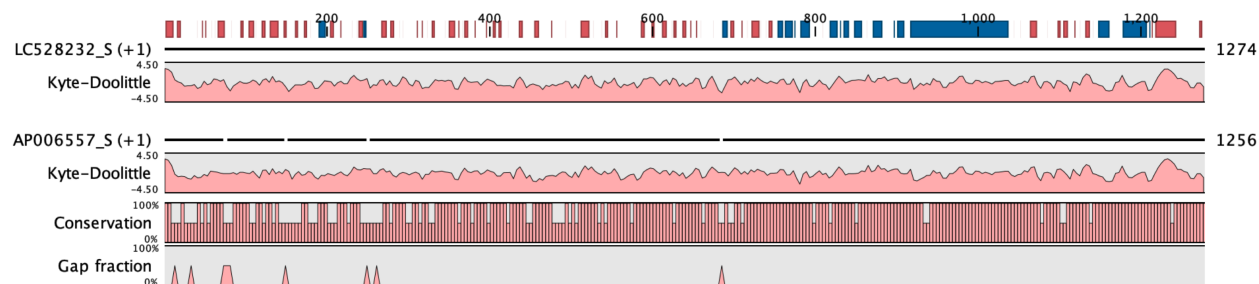


Figure 5: Aligned surface glycoprotein amino acid sequences from SARS-CoV-2 (LC528232) and SARS-CoV (AP006557). Red blocks at the top represent beta-strands and blue blocks represent alpha-helices.

Wan and co-workers reported that mutations in the surface glycoprotein's five amino acids can play a critical role in binding to the human ACE2 receptors.[5] Anderson and co-workers provide their own insights, including another 6th critical amino acid. Here, we integrated secondary structure predictions and realigned the sequences in order to show how these five

mutations may affect the binding of the SARS-CoV-2 protein (Figure 6). The alignment of secondary structure predictions show the additions of three beta strands and the loss of one beta strand. These changes may affect the folding of the protein and hence its binding.
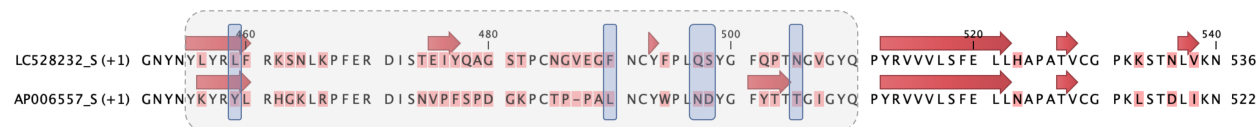


Figure 6: Five key mutations are shown using the blue boxes. These mutations are reported to change the binding affinity of the surface glycoprotein to the human ACE2 receptors. Comparisons are displayed for SARS-CoV-2 (LC528232) and SARS-CoV (AP006557).

Predictions for secondary structures may reveal some information about the binding of the surface glycoprotein. However, 3D models may help understanding the surfaces that are exposed and are likely to bind to other molecules. In order to understand the effects of the mutated region in Figure 6 (shown using the dashed box), an existing 3D model of the SARS-CoV protein was searched for. CLC Genomics Workbench's the 'Toolbox – Classical Sequence Analysis – Protein Analysis – Find and Model Structure' option was used to search for existing 3D models. The first ranked entry with the most 'Match identity' and 'Coverage' was selected. The Protein Data Bank (PDB)[6] identifier of this entry is '5X58 Prefusion Structure of SARS-CoV Spike Glycoprotein , Conformation 1'.[7] Using the CLC Genomics Workbench's 'Project Settings – Sequence tools – Show Sequence' option, the 5X58's 'Chain C' sequence was first displayed and it was then aligned to the SARS-CoV protein sequence (AP006557) using the 'Align to Existing Sequence' option. The first picture on the left in Figure 7 shows the structure of the chain. The SARS-CoV amino acid sequences from the area with the key mutations (shown using the dashed box in Figure 6) were highlighted using the sequence editor. The second picture on the right in Figure 7 shows the corresponding amino acid areas highlighted in the 3D structural view.

A similar analysis was also carried out for the SARS-CoV-2 Chain C. The '6W41' entry,[8] including the details about the crystal structure of the SARS-CoV-2 receptor binding domain, was downloaded from the Protein Data Bank. A collection of related-entries can be found

```
                                    420                              440                                      460                                      480                                  500
                                     |                                |                                        |                                        |                                  |
(5X58 – SPIKE GLYCOPROTEIN) C MGCV  LAWNTRNIDA TSTGNYNYKY RYLRHGKLRP FERDISNVPF SPDGKPCTPP ALNCYWPLND YGFYTTTGIG YQPYRVVVLS 479
            AP006557_S(+1) MGCV  LAWNTRNIDA TSTGNYNYKY RYLRHGKLRP FERDISNVPF SPDGKPCTPP ALNCYWPLND YGFYTTTGIG YQPYRVVVLS 500
```
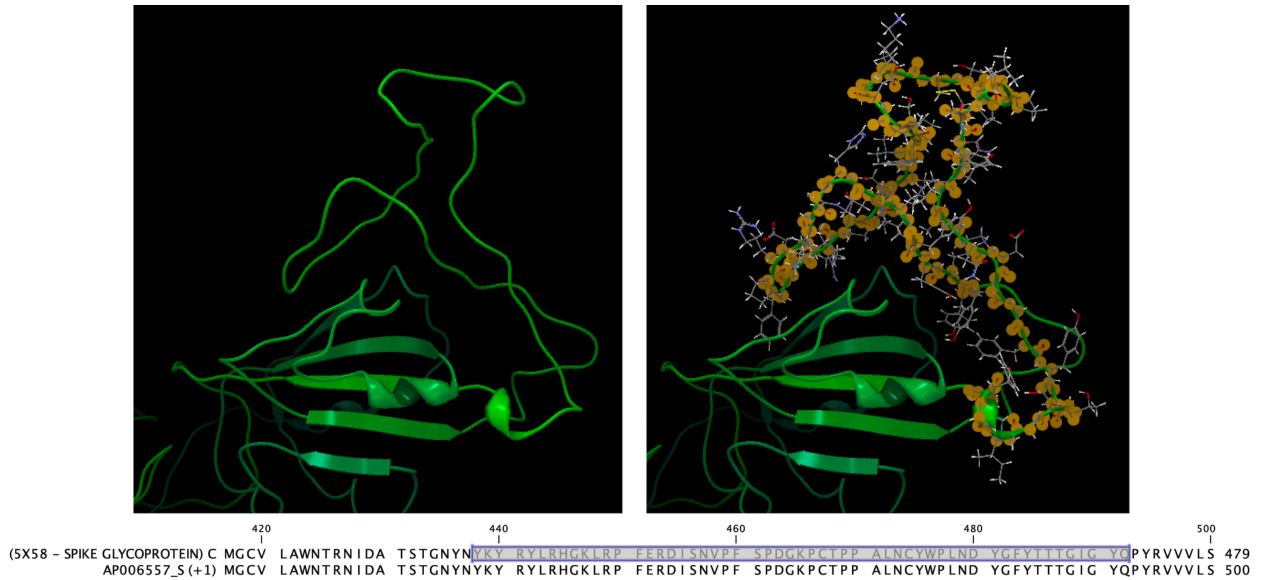
Figure 7: The binding region in SARS-CoV. The left picture shows the SARS-CoV Chain C, which binds to the human ACE2 receptors. The right picture shows the area corresponding to the highlighted sequence below, which represent the mutation area shown using the dashed box in Figure 6.

at the European Bioinformatics Institute's COVID-19 page.[9] The sequence from the '6W41' model was then aligned to the SARS-CoV-2 sequence (LC528232). Compared to the SARS-CoV secondary structures, both the CLC Genomics Workbench predictions and the SARS-CoV-2 '6W41' model reveal additional beta strands in the mutated region (shown using the dashed box in Figure 6) of the surface glycoprotein. The 3D view of the binding region is shown in Figure 8. The second picture on the right in Figure 8 shows the corresponding amino acid areas highlighted in the 3D structural view.

CLC Genomic Workbench was then used to align the two SARS-CoV-2 and SARS-CoV 3D structures using the 'Project Settings – Structure tools – Align Protein Structure' option. Figure 9 shows the alignment using red and blue colours for SARS-CoV-2 and SARS-CoV respectively. The third picture on the right shows the region that aligns with the ACE2 receptors during binding.[10] The middle picture shows the four key amino acid sequences that mutated in SARS-CoV-2. These key mutations affect the binding affinity.[10]
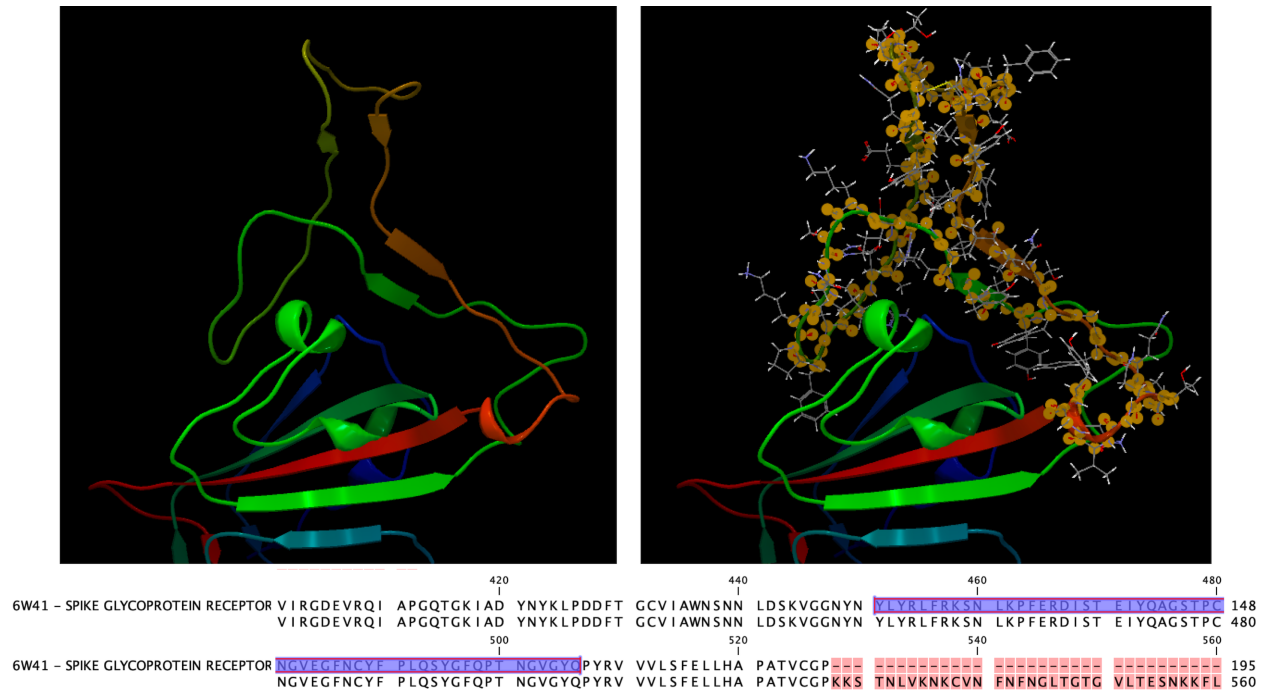
Figure 8: The binding region in SARS-CoV-2. The sequences from 6W41 and LC528232 were aligned. The left picture shows the SARS-CoV-2 Chain C, which binds to the human ACE2 receptors. The right picture shows the area corresponding to the highlighted sequence below, which represents the mutation area shown using the dashed box in Figure 6.
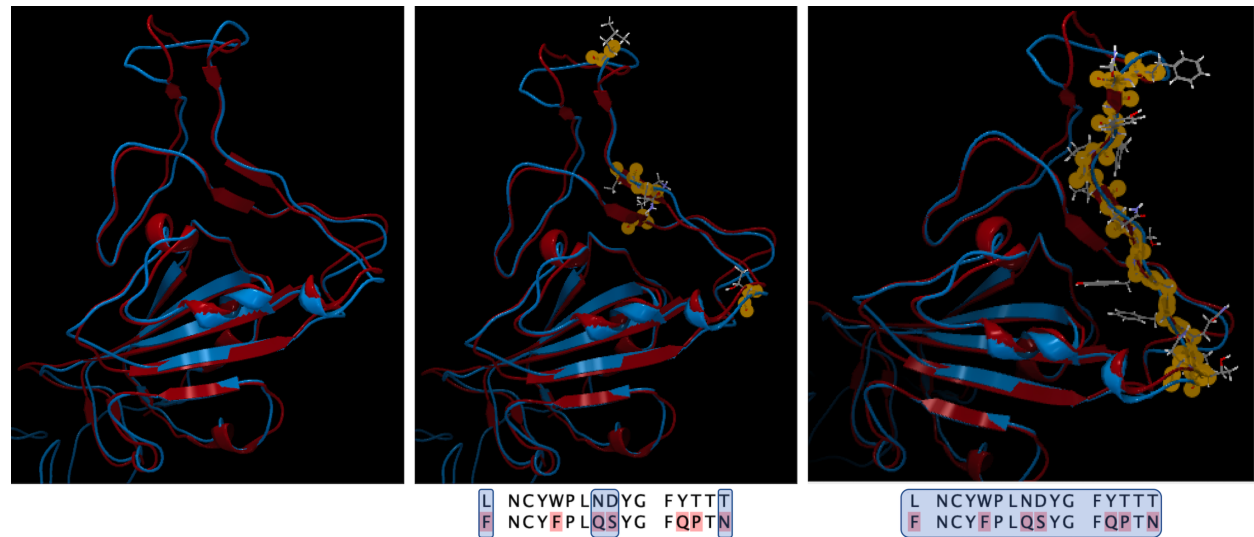


Figure 9: SARS-CoV-2 (in red colour) and SARS-CoV (in blue colour) 3D models of the binding regions to ACE2 sites are aligned. The middle picture shows the locations of the four key mutations in SARS-CoV-2. The third picture on the right shows the ACE2 binding area. SARS-CoV-2 sequences are shown at the bottom and the highlighted sequences are displayed using blue boxes.

# Discussion

In this report, insights from the sequence analysis of SARS-CoV-2 were presented. Genome-wide comparisons between SARS-CoV-2 and SARS-CoV further shed lights into understanding the potential effects of key mutations in the surface glycoprotein amino acid sequences. Our analysis is inline with the current reports.[10,11] Although, some of the mutations may play an important role in binding to the human ACE2 receptors, there are several mutations that may have strengthened the binding affinity of the surface glycoprotein. Both the predictions and the 3D models reveal additional beta strands and hydrogen bonds. These additional mutations may cause changes in structural conformations and cause a higher binding affinity.

This report has been prepared in a tutorial style using CLC Genomics Workbench. The approach can be adopted by biologists and others to have initial understanding of SARS-CoV-2 mutations, and their relationships to SARS-CoV and other closely related viruses. Initial findings can then be explored further by using other specialised tools and approaches.

Computational methods are especially promising. Machine learning, artificial intelligence, data integration and mining, visualisation, computational and mathematical modelling for key biochemical interactions and disease control mechanisms can usefully be applied to provide solutions in a cost- and time- effective manner. We hope that this report is useful to those who wish to understand essential information about SARS-CoV-2 for subsequent analyses.

# Acknowledgements

# References

(1) Zhou, P.; Yang, X.-L.; Wang, X.-G.; Hu, B.; Zhang, L.; Zhang, W.; Si, H.-R.; Zhu, Y.; Li, B.; Huang, C.-L.; Chen, H.-D.; Chen, J.; Luo, Y.; Guo, H.; Jiang, R.-D.; Liu, M.-Q.; Chen, Y.; Shen, X.-R.; Wang, X.; Zheng, X.-S.; Zhao, K.; Chen, Q.-J.; Deng, F.; Liu, L.-L.; Yan, B.; Zhan, F.-X.; Wang, Y.-Y.; Xiao, G.-F.; Shi, Z.-L. A pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* **2020**, *579*, 270–273.

(2) The National Center for Biotechnology, Taxonomy Browser. `https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi`.

(3) The National Center for Biotechnology, Virus Database. `https://www.ncbi.nlm.nih.gov/labs/virus/vssi`.

(4) Andersen, K. G.; Rambaut, A.; Lipkin, W. I.; Holmes, E. C.; Garry, R. F. The proximal origin of SARS-CoV-2. *Nature Medicine* **2020**,

(5) Wan, Y.; Shang, J.; Graham, R.; Baric, R. S.; Li, F. Receptor Recognition by the Novel Coronavirus from Wuhan: an Analysis Based on Decade-Long Structural Studies of SARS Coronavirus. *Journal of Virology* **2020**, *94*.

(6) Berman, H.; Henrick, K.; Nakamura, H. Announcing the worldwide Protein Data Bank. *Nature Structural & Molecular Biology* **2003**, *10*, 980.

(7) Yuan, Y.; Cao, D.; Zhang, Y.; Ma, J.; Qi, J.; Wang, Q.; Lu, G.; Wu, Y.; Yan, J.; Shi, Y.; Zhang, X.; Gao, G. F. Cryo-EM structures of MERS-CoV and SARS-CoV spike glycoproteins reveal the dynamic receptor binding domains. *Nature Communications* **2017**, *8*, 15092.

(8) Protein Data Bank, Structure of 2019-nCoV chimeric receptor-binding domain complexed with its receptor human ACE2. `http://www.rcsb.org/structure/6VW1`.

(9) The European Bioinformatics Institute, COVID-19 Outbreak. `https://www.ebi.ac.uk/ena/pathogens/covid-19`.

(10) Shang, J.; Ye, G.; Shi, K.; Wan, Y.; Luo, C.; Aihara, H.; Geng, Q.; Auerbach, A.; Li, F. Structural basis of receptor recognition by SARS-CoV-2. *Nature* **2020**,

(11) Wu, F.; Zhao, S.; Yu, B.; Chen, Y.-M.; Wang, W.; Song, Z.-G.; Hu, Y.; Tao, Z.-W.; Tian, J.-H.; Pei, Y.-Y.; Yuan, M.-L.; Zhang, Y.-L.; Dai, F.-H.; Liu, Y.; Wang, Q.-M.; Zheng, J.-J.; Xu, L.; Holmes, E. C.; Zhang, Y.-Z. A new coronavirus associated with human respiratory disease in China. *Nature* **2020**, *579*, 265–269.