

Uncanny similarity of unique inserts in the 2019-nCoV spike protein to HIV-1 gp120 and Gag

Prashant Pradhan^{\$1,2}, Ashutosh Kumar Pandey^{\$1}, Akhilesh Mishra^{\$1}, Parul Gupta¹, Praveen Kumar Tripathi¹, Manoj Balakrishnan Menon¹, James Gomes¹, Perumal Vivekanandan^{*1} and Bishwajit Kundu^{*1}

¹Kusuma School of biological sciences, Indian institute of technology, New Delhi-110016, India.

²Acharya Narendra Dev College, University of Delhi, New Delhi-110019, India

^{\$}Equal contribution

* Corresponding authors- email: bkundu@bioschool.iitd.ac.in

vperumal@bioschool.iitd.ac.in

Abstract:

We are currently witnessing a major epidemic caused by the 2019 novel coronavirus (2019-nCoV). The evolution of 2019-nCoV remains elusive. We found 4 insertions in the spike glycoprotein (S) which are unique to the 2019-nCoV and are not present in other coronaviruses. Importantly, amino acid residues in all the 4 inserts have identity or similarity to those in the HIV-1 gp120 or HIV-1 Gag. Interestingly, despite the inserts being discontinuous on the primary amino acid sequence, 3D-modelling of the 2019-nCoV suggests that they converge to constitute the receptor binding site. The finding of 4 unique inserts in the 2019-nCoV, all of which have identity /similarity to amino acid residues in key structural proteins of HIV-1 is unlikely to be fortuitous in nature. This work provides yet unknown insights on 2019-nCoV and sheds light on the evolution and pathogenicity of this virus with important implications for diagnosis of this virus.

Introduction

Coronaviruses (CoV) are single-stranded positive-sense RNA viruses that infect animals and humans. These are classified into 4 genera based on their host specificity: *Alphacoronavirus*, *Betacoronavirus*, *Deltacoronavirus* and *Gammacoronavirus* (Snijder et al., 2006). There are seven known types of CoVs that includes 229E and NL63 (Genus Alphacoronavirus), OC43, HKU1, MERS and SARS (Genus Betacoronavirus). While 229E, NL63, OC43, and HKU1 commonly infect humans, the SARS and MERS outbreak in 2002 and 2012 respectively occurred when the virus crossed-over from animals to humans causing significant mortality (J. Chan et al., n.d.; J. F. W. Chan et al., 2015). In December 2019, another outbreak of coronavirus was reported from Wuhan, China that also transmitted from animals to humans. This new virus has been temporarily termed as 2019-novel Coronavirus (2019-nCoV) by the World Health Organization (WHO) (J. F.-W. Chan et al., 2020; Zhu et al., 2020). While there are several hypotheses about the origin of 2019-nCoV, the source of this ongoing outbreak remains elusive.

The transmission patterns of 2019-nCoV is similar to patterns of transmission documented in the previous outbreaks including by bodily or aerosol contact with persons infected with the virus.

Cases of mild to severe illness, and death from the infection have been reported from Wuhan. This outbreak has spread rapidly distant nations including France, Australia and USA among others. The number of cases within and outside China are increasing steeply. Our current understanding is limited to the virus genome sequences and modest epidemiological and clinical data. Comprehensive analysis of the available 2019-nCoV sequences may provide important clues that may help advance our current understanding to manage the ongoing outbreak.

The spike glycoprotein (S) of coronavirus is cleaved into two subunits (S1 and S2). The S1 subunit helps in receptor binding and the S2 subunit facilitates membrane fusion (Bosch et al., 2003; Li, 2016). The spike glycoproteins of coronaviruses are important determinants of tissue tropism and host range. In addition the spike glycoproteins are critical targets for vaccine development (Du et al., 2013). For this reason, the spike proteins represent the most extensively studied among coronaviruses. We therefore sought to investigate the spike glycoprotein of the 2019-nCoV to understand its evolution, novel features sequence and structural features using computational tools.

Methodology

Retrieval and alignment of nucleic acid and protein sequences

We retrieved all the available coronavirus sequences (n=55) from NCBI viral genome database (<https://www.ncbi.nlm.nih.gov/>) and we used the GISAID (Elbe & Buckland-Merrett, 2017)[<https://www.gisaid.org/>] to retrieve all available full-length sequences (n=28) of 2019-nCoV as on 27 Jan 2020. Multiple sequence alignment of all coronavirus genomes was performed by using MUSCLE software (Edgar, 2004) based on neighbour joining method. Out of 55 coronavirus genome 32 representative genomes of all category were used for phylogenetic tree development using MEGAX software (Kumar et al., 2018). The closest relative was found to be SARS CoV. The glycoprotein region of SARS CoV and 2019-nCoV were aligned and visualized using Multalin software (Corpet, 1988). The identified amino acid and nucleotide sequence were aligned with whole viral genome database using BLASTp and BLASTn. The conservation of the nucleotide and amino acid motifs in 28 clinical variants of 2019-nCoV genome were presented by performing multiple sequence alignment using MEGAX software. The three dimensional structure of 2019-nCoV glycoprotein was generated by using SWISS-MODEL online server (Biasini et al., 2014) and the structure was marked and visualized by using PyMol (DeLano, 2002).

Results

Uncanny similarity of novel inserts in the 2019-nCoV spike protein to HIV-1 gp120 and Gag

Our phylogenetic tree of full-length coronaviruses suggests that 2019-nCoV is closely related to SARS CoV [Fig1]. In addition, other recent studies have linked the 2019-nCoV to SARS CoV. We therefore compared the spike glycoprotein sequences of the 2019-nCoV to that of the SARS CoV (NCBI Accession number: AY390556.1). On careful examination of the sequence alignment we found that the 2019-nCoV spike glycoprotein contains 4 insertions [Fig.2]. To further investigate if these inserts are present in any other corona virus, we performed a multiple

sequence alignment of the spike glycoprotein amino acid sequences of all available coronaviruses (n=55) [refer Table S.File1] in NCBI refseq (ncbi.nlm.nih.gov) this includes one sequence of 2019-nCoV [Fig.S1]. We found that these 4 insertions [inserts 1, 2, 3 and 4] are unique to 2019-nCoV and are not present in other coronaviruses analyzed. Another group from China had documented three insertions comparing fewer spike glycoprotein sequences of coronaviruses. Another group from China had documented three insertions comparing fewer spike glycoprotein sequences of coronaviruses (Zhou et al., 2020).

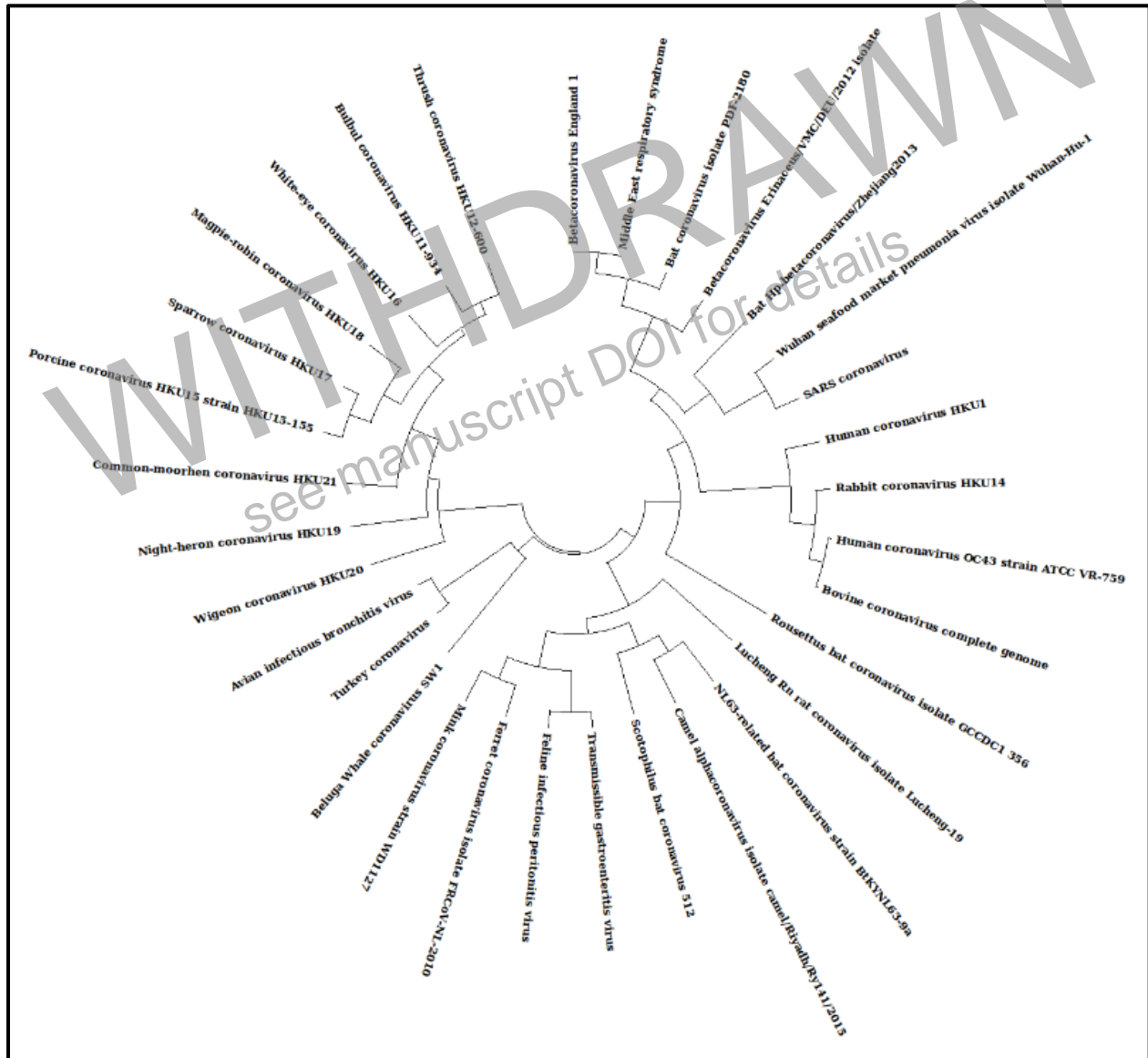


Figure 1: Maximum likelihood genealogy show the evolution of 2019- nCoV: The evolutionary history was inferred by using the Maximum Likelihood method and JTT matrix-based model. The tree with the highest log likelihood (12458.88) is shown. Initial tree(s) for the heuristic search were obtained automatically by applying Neighbor-Join and BioNJ algorithms to a matrix of pairwise distances estimated using a JTT model, and then selecting the topology with superior log likelihood

value. This analysis involved 5 amino acid sequences. There were a total of 1387 positions in the final dataset. Evolutionary analyses were conducted in MEGA X.



Figure 2: Multiple sequence alignment between spike proteins of 2019-nCoV and SARS. The sequences of spike proteins of 2019-nCoV (Wuhan-HU-1, Accession NC_045512) and of SARS CoV (GZ02, Accession AY390556) were aligned using MultiAlin software. The sites of difference are highlighted in boxes.

We then analyzed all available full-length sequences (n=28) of 2019-nCoV in GISAID (Elbe & Buckland-Merrett, 2017) as on January 27, 2020 for the presence of these inserts. As most of these sequences are not annotated, we compared the nucleotide sequences of the spike glycoprotein of all available 2019-nCoV sequences using BLASTp. Interestingly, all the 4 insertions were absolutely (100%) conserved in all the available 2019-nCoV sequences analyzed [Fig.S2, Fig.S3].

We then translated the aligned genome and found that these inserts are present in all Wuhan 2019-nCoV viruses except the 2019-nCoV virus of Bat as a host [Fig.S4]. Intrigued by the 4 highly conserved inserts unique to 2019-nCoV we wanted to understand their origin. For this purpose, we used the 2019-nCoV local alignment with each insert as query against all virus genomes and considered hits with 100% sequence coverage. Surprisingly, each of the four inserts aligned with short segments of the Human immunodeficiency Virus-1 (HIV-1) proteins. The amino acid positions of the inserts in 2019-nCoV and the corresponding residues in HIV-1 gp120 and HIV-1 Gag are shown in Table 1. The first 3 inserts (insert 1,2 and 3) aligned to short segments of amino acid residues in HIV-1 gp120. The insert 4 aligned to HIV-1 Gag. The insert 1 (6 amino acid residues) and insert 2 (6 amino acid residues) in the spike glycoprotein of 2019-nCoV are 100% identical to the residues mapped to HIV-1 gp120. The insert 3 (12 amino acid residues) in 2019-nCoV maps to HIV-1 gp120 with gaps [see Table 1]. The insert 4 (8 amino acid residues) maps to HIV-1 Gag with gaps.

Although, the 4 inserts represent discontinuous short stretches of amino acids in spike glycoprotein of 2019-nCoV, the fact that all three of them share amino acid identity or similarity with HIV-1 gp120 and HIV-1 Gag (among all annotated virus proteins) suggests that this is not a random fortuitous finding. In other words, one may sporadically expect a fortuitous match for a stretch of 6-12 contiguous amino acid residues in an unrelated protein. However, it is unlikely that all 4 inserts in the 2019-nCoV spike glycoprotein fortuitously match with 2 key structural proteins of an unrelated virus (HIV-1).

The amino acid residues of inserts 1, 2 and 3 of 2019-nCoV spike glycoprotein that mapped to HIV-1 were a part of the V4, V5 and V1 domains respectively in gp120 [Table 1]. Since the 2019-nCoV inserts mapped to variable regions of HIV-1, they were not ubiquitous in HIV-1 gp120, but were limited to selected sequences of HIV-1 [refer S.File1] primarily from Asia and Africa.

The HIV-1 Gag protein enables interaction of virus with negatively charged host surface (Murakami, 2008) and a high positive charge on the Gag protein is a key feature for the host-virus interaction. On analyzing the pI values for each of the 4 inserts in 2019-nCoV and the corresponding stretches of amino acid residues from HIV-1 proteins we found that a) the pI values were very similar for each pair analyzed b) most of these pI values were 10 ± 2 [Refer Table 1]. Of note, despite the gaps in inserts 3 and 4 the pI values were comparable. This uniformity in the pI values for all the 4 inserts merits further investigation.

As none of these 4 inserts are present in any other coronavirus, the genomic region encoding these inserts represent ideal candidates for designing primers that can distinguish 2019-nCoV from other coronaviruses.

Motifs	Virus Glycoprotein	Motif Alignment	HIV protein and Variable region	HIV Genome Source Country/ subtype	Number of Polar Residues	Total Charge	pI Value
Insert 1	2019- nCoV (GP) HIV1(GP120)	71 76 TNGTKR TNGTKR 404 409	gp120- V4	Thailand */ CRF01_ AE	5 5	2 2	11 11
Insert 2	2019- nCoV (GP) HIV1(GP120)	145 150 HKNNKS HKNNKS 462 467	gp120- V5	Kenya* G	6 6	2 2	10 10
Insert 3	2019- nCoV (GP) HIV1(GP120)	245 256 RSYL----TPGDSSSG RTYLFNETRGNSSSG 136 150	gp120- V1	India*/C	8 10	2 1	10.84 8.75
Insert 4	2019- nCoV (Poly P) HIV1(gag)	676 684 QTNS-----PRRA QTNSSILMQRSNFKG PRRA 366 384	Gag	India*/C	6 12	2 4	12.00 12.30

Table 1: Aligned sequences of 2019-nCoV and gp120 protein of HIV-1 with their positions in primary sequence of protein. All the inserts have a high density of positively charged residues. The deleted fragments in insert 3 and 4 increase the positive charge to surface area ratio. *please see Supp. Table 1 for accession numbers

The novel inserts are part of the receptor binding site of 2019-nCoV

To get structural insights and to understand the role of these insertions in 2019-nCoV glycoprotein, we modelled its structure based on available structure of SARS spike glycoprotein (PDB: 6ACD.1.A). The comparison of the modelled structure reveals that although inserts 1,2 and 3 are at non-contiguous locations in the protein primary sequence, they fold to constitute the part of glycoprotein binding site that recognizes the host receptor (Kirchdoerfer et al., 2016) (Figure 4). The insert 1 corresponds to the NTD (N-terminal domain) and the inserts 2 and 3 correspond to the CTD (C-terminal domain) of the S1 subunit in the 2019-nCoV spike glycoprotein. The insert 4 is at the junction of the SD1 (sub domain 1) and SD2 (sub domain 2) of the S1 subunit (Ou et al., 2017). We speculate, that these insertions provide additional flexibility to the glycoprotein binding site by forming a hydrophilic loop in the protein structure that may facilitate or enhance virus-host interactions.

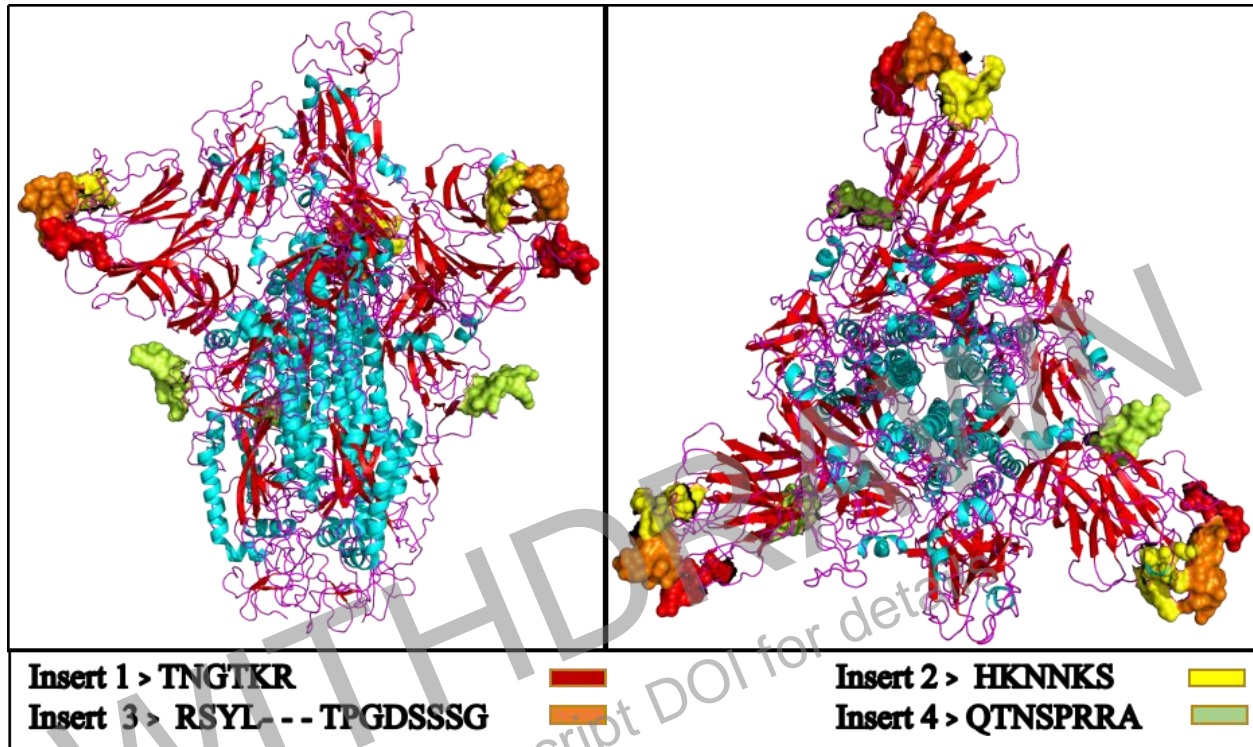


Figure 3. Modelled homo-trimer spike glycoprotein of 2019-nCoV virus. The inserts from HIV envelop protein are shown with colored beads, present at the binding site of the protein.

Evolutionary Analysis of 2019-nCoV

It has been speculated that 2019-nCoV is a variant of Coronavirus derived from an animal source which got transmitted to humans. Considering the change of specificity for host, we decided to study the sequences of spike glycoprotein (S protein) of the virus. S proteins are surface proteins that help the virus in host recognition and attachment. Thus, a change in these proteins can be reflected as a change of host specificity of the virus. To know the alterations in S protein gene of 2019-nCoV and its consequences in structural re-arrangements we performed *in-silico* analysis of 2019-nCoV with respect to all other viruses. A multiple sequence alignment between the S protein amino acid sequences of 2019-nCoV, Bat-SARS-Like, SARS-GZ02 and MERS revealed that S protein has evolved with closest significant diversity from the SARS-GZ02 (Figure 1).

Insertions in Spike protein region of 2019-nCoV

Since the S protein of 2019-nCoV shares closest ancestry with SARS GZ02, the sequence coding for spike proteins of these two viruses were compared using MultiAlin software. We found four new insertions in the protein of 2019-nCoV- “GTNGTKR” (IS1), “HKNNKS” (IS2), “GDSSSG” (IS3) and “QTNSPRRA” (IS4) (Figure 2). To our surprise, these sequence insertions were not only absent in S protein of SARS but were also not observed in any other member of the *Coronaviridae* family (Supplementary figure). This is startling as it is quite unlikely for a virus to have acquired such unique insertions naturally in a short duration of time.

Insertions share similarity to HIV

The insertions were observed to be present in all the genomic sequences of 2019-nCoV virus available from the recent clinical isolates (Supplementary Figure 1). To know the source of these insertions in 2019-nCoV a local alignment was done with BLASTp using these insertions as query with all virus genome. Unexpectedly, all the insertions got aligned with Human immunodeficiency Virus-1 (HIV-1). Further analysis revealed that aligned sequences of HIV-1 with 2019-nCoV were derived from surface glycoprotein gp120 (amino acid sequence positions: 404-409, 462-467, 136-150) and from Gag protein (366-384 amino acid) (Table 1). Gag protein of HIV is involved in host membrane binding, packaging of the virus and for the formation of virus-like particles. Gp120 plays crucial role in recognizing the host cell by binding to the primary receptor CD4. This binding induces structural rearrangements in GP120, creating a high affinity binding site for a chemokine co-receptor like CXCR4 and/or CCR5.

Discussion

The current outbreak of 2019-nCoV warrants a thorough investigation and understanding of its ability to infect human beings. Keeping in mind that there has been a clear change in the preference of host from previous coronaviruses to this virus, we studied the change in spike protein between 2019-nCoV and other viruses. We found four new insertions in the S protein of 2019-nCoV when compared to its nearest relative, SARS CoV. The genome sequence from the recent 28 clinical isolates showed that the sequence coding for these insertions are conserved amongst all these isolates. This indicates that these insertions have been preferably acquired by the 2019-nCoV, providing it with additional survival and infectivity advantage. Delving deeper we found that these insertions were similar to HIV-1. Our results highlight an astonishing relation between the gp120 and Gag protein of HIV, with 2019-nCoV spike glycoprotein. These proteins are critical for the viruses to identify and latch on to their host cells and for viral assembly (Beniac et al., 2006). Since surface proteins are responsible for host tropism, changes in these proteins imply a change in host specificity of the virus. According to reports from China, there has been a gain of host specificity in case 2019-nCoV as the virus was originally known to infect animals and not humans but after the mutations, it has gained tropism to humans as well.

Moving ahead, 3D modelling of the protein structure displayed that these insertions are present at the binding site of 2019-nCoV. Due to the presence of gp120 motifs in 2019-nCoV spike glycoprotein at its binding domain, we propose that these motif insertions could have provided an enhanced affinity towards host cell receptors. Further, this structural change might have also increased the range of host cells that 2019-nCoV can infect. To the best of our knowledge, the function of these motifs is still not clear in HIV and need to be explored. The exchange of genetic material among the viruses is well known and such critical exchange highlights the risk and the need to investigate the relations between seemingly unrelated virus families.

Conclusions

Our analysis of the spike glycoprotein of 2019-nCoV revealed several interesting findings: First, we identified 4 unique inserts in the 2019-nCoV spike glycoprotein that are not present in any other coronavirus reported till date. To our surprise, all the 4 inserts in the 2019-nCoV mapped to

short segments of amino acids in the HIV-1 gp120 and Gag among all annotated virus proteins in the NCBI database. This uncanny similarity of novel inserts in the 2019-nCoV spike protein to HIV-1 gp120 and Gag is unlikely to be fortuitous. Further, 3D modelling suggests that at least 3 of the unique inserts which are non-contiguous in the primary protein sequence of the 2019-nCoV spike glycoprotein converge to constitute the key components of the receptor binding site. Of note, all the 4 inserts have pI values of around 10 that may facilitate virus-host interactions. Taken together, our findings suggest unconventional evolution of 2019-nCoV that warrants further investigation. Our work highlights novel evolutionary aspects of the 2019-nCoV and has implications on the pathogenesis and diagnosis of this virus.

References

- Beniac, D. R., Andonov, A., Grudeski, E., & Booth, T. F. (2006). Architecture of the SARS coronavirus prefusion spike. *Nature Structural and Molecular Biology*, 13(8), 751–752. <https://doi.org/10.1038/nsmb1123>
- Biasini, M., Bienert, S., Waterhouse, A., Arnold, K., Studer, G., Schmidt, T., Kiefer, F., Cassarino, T. G., Bertoni, M., Bordoli, L., & Schwede, T. (2014). SWISS-MODEL: Modelling protein tertiary and quaternary structure using evolutionary information. *Nucleic Acids Research*. <https://doi.org/10.1093/nar/gku340>
- Bosch, B. J., van der Zee, R., de Haan, C. A. M., & Rottier, P. J. M. (2003). The Coronavirus Spike Protein Is a Class I Virus Fusion Protein: Structural and Functional Characterization of the Fusion Core Complex. *Journal of Virology*, 77(16), 8801–8811. <https://doi.org/10.1128/jvi.77.16.8801-8811.2003>
- Chan, J. F.-W., Kok, K.-H., Zhu, Z., Chu, H., To, K. K.-W., Yuan, S., & Yuen, K.-Y. (2020). Genomic characterization of the 2019 novel human-pathogenic coronavirus isolated from a patient with atypical pneumonia after visiting Wuhan. *Emerging Microbes & Infections*, 9(1), 221–236. <https://doi.org/10.1080/22221751.2020.1719902>
- Chan, J. F. W., Lau, S. K. P., To, K. K. W., Cheng, V. C. C., Woo, P. C. Y., & Yuen, K.-Y. (2015). *Middle East Respiratory Syndrome Coronavirus: Another Zoonotic Betacoronavirus Causing SARS-Like Disease*. <https://doi.org/10.1128/CMR.00102-14>
- Chan, J., To, K., Tse, H., Jin, D., microbiology, K. Y.-T. in, & 2013, undefined. (n.d.). Interspecies transmission and emergence of novel viruses: lessons from bats and birds. *Elsevier*.
- Corpet, F. (1988). Multiple sequence alignment with hierarchical clustering. *Nucleic Acids Research*. <https://doi.org/10.1093/nar/16.22.10881>
- DeLano, W. L. (2002). The PyMOL Molecular Graphics System, Version 1.1. *Schrödinger LLC*. <https://doi.org/10.1038/hr.2014.17>
- Du, L., Zhao, G., Kou, Z., Ma, C., Sun, S., Poon, V. K. M., Lu, L., Wang, L., Debnath, A. K., Zheng, B.-J., Zhou, Y., & Jiang, S. (2013). Identification of a Receptor-Binding Domain in the S Protein of the Novel Human Coronavirus Middle East Respiratory Syndrome Coronavirus as an Essential Target for Vaccine Development. *Journal of Virology*, 87(17), 9939–9942. <https://doi.org/10.1128/jvi.01048-13>

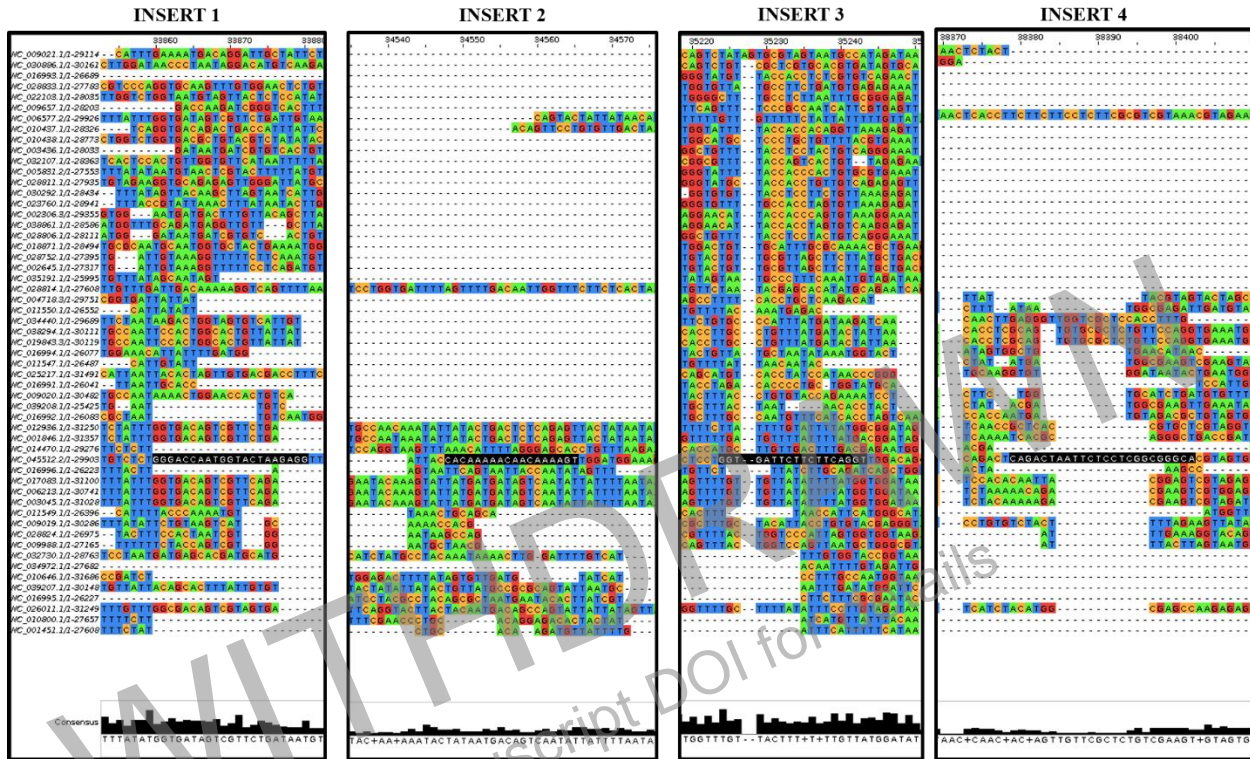
- Edgar, R. C. (2004). MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Research*. <https://doi.org/10.1093/nar/gkh340>
- Elbe, S., & Buckland-Merrett, G. (2017). Data, disease and diplomacy: GISAID's innovative contribution to global health. *Global Challenges*. <https://doi.org/10.1002/gch2.1018>
- Kirchdoerfer, R. N., Cottrell, C. A., Wang, N., Pallesen, J., Yassine, H. M., Turner, H. L., Corbett, K. S., Graham, B. S., McLellan, J. S., & Ward, A. B. (2016). Pre-fusion structure of a human coronavirus spike protein. *Nature*. <https://doi.org/10.1038/nature17200>
- Kumar, S., Stecher, G., Li, M., Knyaz, C., & Tamura, K. (2018). MEGA X: Molecular evolutionary genetics analysis across computing platforms. *Molecular Biology and Evolution*. <https://doi.org/10.1093/molbev/msy096>
- Li, F. (2016). Structure, Function, and Evolution of Coronavirus Spike Proteins. *Annual Review of Virology*, 3(1), 237–261. <https://doi.org/10.1146/annurev-virology-110615-042301>
- Murakami, T. (2008). Roles of the interactions between Env and Gag proteins in the HIV-1 replication cycle. *Microbiology and Immunology*, 52(5), 287–295. <https://doi.org/10.1111/j.1348-0421.2008.00008.x>
- Ou, X., Guan, H., Qin, B., Mu, Z., Wojdyla, J. A., Wang, M., Dominguez, S. R., Qian, Z., & Cui, S. (2017). Crystal structure of the receptor binding domain of the spike glycoprotein of human betacoronavirus HKU1. *Nature Communications*. <https://doi.org/10.1038/ncomms15216>
- Snijder, E. J., van der Meer, Y., Zevenhoven-Dobbe, J., Onderwater, J. J. M., van der Meulen, J., Koerten, H. K., & Mommaas, A. M. (2006). Ultrastructure and origin of membrane vesicles associated with the severe acute respiratory syndrome coronavirus replication complex. *Journal of Virology*, 80(12), 5927–5940. <https://doi.org/10.1128/JVI.02501-05>
- 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., ... Shi, Z.-L. (2020). Discovery of a novel coronavirus associated with the recent pneumonia outbreak in humans and its potential bat origin. *BioRxiv*. <https://doi.org/10.1101/2020.01.22.914952>
- Zhu, N., Zhang, D., Wang, W., Li, X., Yang, B., Song, J., Zhao, X., Huang, B., Shi, W., Lu, R., Niu, P., Zhan, F., Ma, X., Wang, D., Xu, W., Wu, G., Gao, G. F., & Tan, W. (2020). A Novel Coronavirus from Patients with Pneumonia in China, 2019. *New England Journal of Medicine*, NEJMoa2001017. <https://doi.org/10.1056/NEJMoa2001017>



Fig.S1 Multiple sequence alignment of glycoprotein of *coronaviridae* family, representing all the four inserts.



Fig.S3 Phylogenetic tree of 28 clinical isolates genome of 2019-nCoV including one from bat as a host.



Supplementary Fig 4. Genome alignment of Coronaviridae family. Highlighted black sequences are the inserts represented here