

PATHOGENETIC PERSPECTIVE OF MISSENSE MUTATIONS OF ORF3A PROTEIN OF SARS-CoV2

Sk. Sarif Hassan^{a,*}, Diksha Attrish^r^b, Shinjini Ghosh^{†c}, Pabitra Pal Choudhury^d, Bidyut Roy^e

^aDepartment of Mathematics, Pingla Thana Mahavidyalaya, Maligram 721140, India

^bDr. B. R. Ambedkar Centre For Biomedical Research (ACBR), University of Delhi (North Campus), Delhi 110007, India

^cDepartment of Biophysics, Molecular Biology and Bioinformatics, University of Calcutta, Kolkata 700009, West Bengal, India

^dApplied Statistics Unit, Indian Statistical Institute, Kolkata 700108, West Bengal, India

^eHuman Genetics Unit, Indian Statistical Institute, Kolkata 700108, West Bengal, India

Abstract

One of the most important proteins for COVID-19 pathogenesis in SARS-CoV2 is the ORF3a protein which is the largest accessory protein among others accessory proteins coded by coronavirus genome. The major roles of the protein include virulence, infectivity, ion channel activity, morphogenesis and virus release. The coronavirus, SARS-CoV2 is continuously evolving naturally and thereby the encoded proteins are also mutating rapidly. Therefore, critical study of mutations in ORF3a is certainly important from the pathogenetic perspective. Here, a sum of 175 various non-synonymous mutations in the ORF3a protein of SARS-CoV2 are identified and their corresponding effects in structural stability and functions of the protein ORF3a are studied. Broadly three different classes of mutations, such as neutral, disease and mixed (neutral and disease) type mutations were observed. Consecutive mutations in some ORF3a proteins are established based on timeline of detection of mutations. Considering the amino acid compositions over the ORF3a primary protein sequences, twenty clusters are detected based on K-means clustering method. Our findings on 175 novel mutations of ORF3a proteins will extend our knowledge of ORF3a, a vital accessory protein in SARS-CoV2, which would assist to enlighten on the pathogenicity of this life-threatening COVID-19.

Keywords: SARS-CoV2, ORF3a, COVID-19, Missense mutations, Shannon entropy and Genetic variations.

1. Introduction

. Severe Acute Respiratory Syndrome (SARS-CoV) emerged in 2002 infecting about 8000 people with a 10% mortality rate [1, 2]. Similarly, Middle East Respiratory Syndrome Coronavirus (MERS-CoV) emerged in 2012 with 2300 cases and a 35% mortality rate [3]. However, since the December 2019, another outbreak caused by a novel severe acute respiratory syndrome coronavirus 2 (SARS-CoV2) rapidly became a pandemic with the highest mortality rate of 3.4% within just 7 months; urging the World Health Organization to declare it as a Public Health Emergency of International Concern [4, 5, 6, 7]. It was found that SARS-CoV and SARS-CoV2 bear 79% of sequence identity [8, 9]. Similar to SARS-CoV, the ORF3a gene in SARS-CoV2 lies between the spike and envelope gene in virus genome [10]. Both the ORF3a protein of SARS-CoV and SARS-CoV2 contain a conserved cysteine residue which helps in protein-protein interaction [11, 12]. The RNA genome of SARS-CoV2 is about 30 kb in length and codes for 4 structural proteins, 16 non-structural proteins, and 6/7 accessory proteins [13, 14, 15, 16]. The structural proteins are known as Spike protein (S), Nucleocapsid protein (N), Membrane protein (M) and Envelope protein (E) [17].

*Corresponding author

Email addresses: sarimif@gmail.com (Sk. Sarif Hassan), dikshaattrish@gmail.com (Diksha Attrish^r), shinjinighosh2014@gmail.com (Shinjini Ghosh[†]), pabitrpalchoudhury@gmail.com (Pabitra Pal Choudhury), broy@isical.ac.in (Bidyut Roy)

Among the accessory proteins, our study is based on ORF3a, the largest accessory protein, and a unique membrane protein consisting of three transmembrane domains [18, 19]. SARS-CoV2 ORF3a is a 275 amino acid transmembrane protein that holds an N-terminal, three transmembrane helices followed by a cytosolic domain with multiple β -strands [20]. Functionally ORF3a proteins is divided into six domains [21]. Domain I contain N terminus signal peptide involved in subcellular localization of ORF3a protein [19]. Domain II contains a TRAF-3 binding motif (36-40 aa) through which it activates the NF- κ B and NLRP3 inflammasome by promoting TNF receptor-associated factor 3 (TRAF3)-mediated ubiquitination of apoptosis-associated speck-like protein containing a caspase recruitment domain (ASC) [21]. Domain III (93-133) is important for ion channel activity and has a Cysteine-rich domain which is associated with homodimerization of ORF3a protein which is very similar to SARS-CoV cysteine rich domain responsible for tetramerization (81-160) [22, 23]. Domain IV has a caveolin binding motif (141-149) which regulates viral uptake and trafficking of protein to the plasma membrane or intracellular membranes [24]. Domain V contains a tyrosine-based sorting motif $YXX\phi$ (160-163) which is responsible for Golgi to plasma membrane transport which in SARS-CoV is responsible for the surface expression [25]. Domain VI has an SGD motif (171-173) [23]. ORF3a has pro-apoptotic activity and membrane association is required for this activity. SARS-CoV2 ORF3a has relatively weaker proapoptotic activity and this property is probably contributing to asymptomatic infection and thus causing rapid transmission of the virus [26]. Therefore, ORF3a may become an important therapeutic target, and thus studying mutations in the ORF3a protein sequence becomes an important area in control of virus infection.

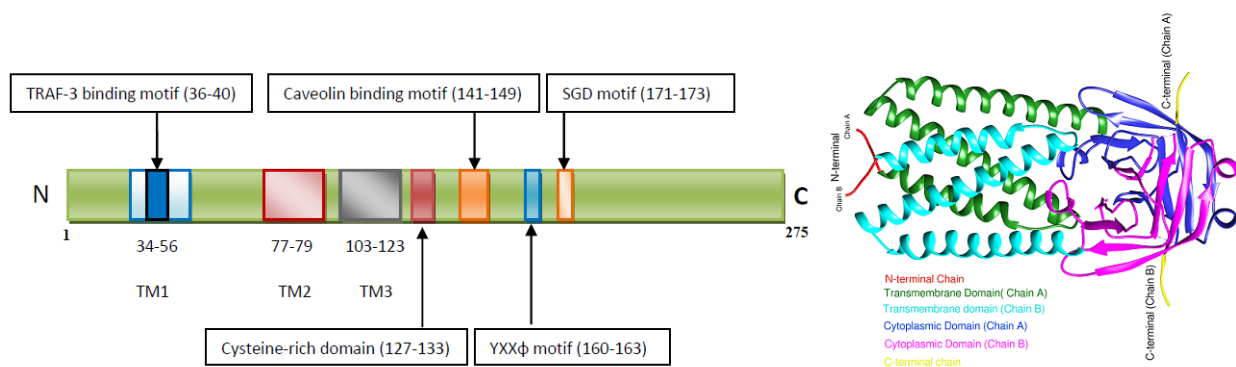


Figure 1: Schematic view of the domains in Primary and Tertiary structures of SARS-CoV2 ORF3a protein

In our present study, we found about 175 non-synonymous mutations in the ORF3a protein sequence. Among them, 32 are already reported previously [27, 23]. So, we accounted 143 new mutations along with the already existing ones. Mutations in the domain III alters the NF- κ B activation and NLRP3 inflammasome. Mutations in domain V were linked to the aggregation of the 3a protein in the Golgi apparatus [28]. Apart from these residues, mutations in 230(insertion of F), W131C, R134L, T151I, N152S and D155Y regions may contribute to a greater significance as they are poised to form a network of hydrophobic, polar and electrostatic interactions which mediate dimerization and tetramerization respectively [29]. To account for mutations of the ORF3a proteins of SARS-CoV2, we collected the SARS-CoV2 genome data from NCBI virus database, identified the mutations, predicted the effect of mutations based on chemical and structural properties. In addition, using the Meta-SNP and I-MUTANT web-servers, effect of the mutations in functions and structures are predicted [30, 31, 32]. We also performed K-means clustering of the distinct variants ORF3a proteins (available as on 27th July, 2020) in order to form twenty disjoint clusters based on the amino acid compositions embedded

in the proteins [33, 34]. In addition, Shannon entropy is employed to determine amount of disorderliness of the amino acids over the ORF3a proteins which amplify the wide distinct variations of ORF3a in the USA [35].

2. Data and Methods

This present study is based on available genome data of SARS-CoV2 from the NCBI virus database (<https://www.ncbi.nlm.nih.gov>).
45 Here we discuss about data followed by methods which are employed in this study.

2.1. Data

As on date 27th July, 2020, there were 7194 complete genomes of SARS-CoV2 available in the NCBI database and accordingly each genome contains one of the accessory proteins ORF3a and among them only 296 sequences are found to be distinct from each other. The amino acid sequences of ORF3a were exported in fasta format using file operations
50 through Matlab [36]. In this present study, we only concentrate on 296 ORF3a proteins which are listed in the Table-1 and Table-2. Note that, among these 296 sequences, three ORF3a proteins QKO00487 (India: Ahmedabad), QLA10225 (India: Vadodara) and QLA10069 (India: Surat) had the length 241, 253 and 257 respectively and were found to be truncated due to nonsense mutation at 242, 254 and 258 amino acid positions respectively. It is also note worthy that some (13.51%) of 296 ORF3a amino acid sequences contain ambiguous amino acids such as X , B and Z and so on. In
55 order to find mutations, we hereby consider the reference ORF3a protein as the ORF3a sequence (YP_009724391.1) of the SARS-CoV2 genome (NC_045512) from China: Wuhan [37].

2.2. Methods

Here in a nutshell, we present the methods used in this study as follows.

2.2.1. Frequency Probability of Amino Acids

60 A protein sequence of ORF3a is composed of twenty different amino acids with various frequencies. The probability of occurrence of each amino acid A_i is determined by the formula $\frac{f(A_i)}{l}$ where $f(A_i)$ denotes the frequency of occurrence of the amino acid A_i in the primary sequence ORF3a and l stands as the length of ORF3a protein [38]. Hence for each of the 296 ORF3a proteins, a twenty dimensional vector considering the frequency probability of twenty amino acids can be obtained. Based on these frequency probability vectors, a classification is performed using clustering technique.

65 2.2.2. K-means Clustering Algorithm

Clustering is one of the most widely used methods in vector-data analysis to develop an intuitive idea about closeness of data based on the structured feature vectors. By clustering we find homogeneous subclasses within the data such that data points in each cluster are as similar as possible according to a similarity measure such as euclidean-based distance. One of the most commonly used simple clustering techniques is the *K-means clustering* [33, 34].

70 **Algorithm:** K-means algorithm is an iterative algorithm that tries to form equivalence classes from the feature vectors into K (pre-defined) clusters where each data point belongs to only one cluster [33].

- Assign the number of desired clusters (K) (in the present study, $K = 20$).
- Finding centroids by first shuffling the dataset and then randomly selecting K data points for the centroids without replacement.

Table 1: List of accessions of the ORF3a protein, geo-location and respective data collection date

Accession	Geo-Location	Collection_Date	Accession	Geo-Location	Collection_Date	Accession	Geo-Location	Collection_Date
YP_00974391	China: Wuhan	2019-12	QLH00578	USA: CA	2020-04-23	QLA10225	India: Vadodara	2020-06-02
QLI49698	India: Himatnagar	2020-06-14	QLH01238	USA: CA	2020-04-21	QKY77929	USA: CA	2020-03-16
QLI50222	USA: New York, Rockland county	2020-06-26	QLH01250	USA: CA	2020-04-22	QKY59990	India: Surat	2020-06-11
QLI50282	USA: Wisconsin, Dane county	2020-06-26	QLH01298	USA: CA	2020-04-22	KX46204	USA	2020-05-11
QLI50414	USA: Wisconsin, Dane county	2020-06-28	QLH01334	USA: CA	2020-04-24	KX47995	Bangladesh: Rangpur	2020-06-07
QLI50570	USA: Wisconsin, Dane county	2020-06-27	QLH01382	USA: CA	2020-05-04	KX49024	Bangladesh	2020-05-23
QLI51038	USA: Wisconsin, Dane county	2020-06-30	QLH01502	USA: CA	2020-05-04	QKW8844	USA	2020-03-14
QLI51614	USA: Wisconsin, Ozaukee county	2020-03-19	QLF97736	Bangladesh	2020-06-17	QKW89480	USA	2020-03-25
QLI51746	USA: Wisconsin, Fond du Lac county	2020-03-31	QLF97772	Bangladesh	2020-06-18	QKY35400	USA: Washington, Yakima County	2020-04-15
QLI51782	USA: Wisconsin, Fond du Lac county	2020-03-31	QLF97844	Bangladesh	2020-06-18	QKY35688	USA: Washington, Yakima County	2020-04-13
QLI64290	USA: Arkansas, Little Rock	2020-04-01	QLF97952	India: Vadodara	2020-06-08	QKY36900	USA: Washington, Yakima County	2020-04-11
QLH64816	India: Modasa	2020-06-14	QLF98036	Bangladesh	2020-06-17	QKY37633	Australia: Victoria	2020-03-24
QLH93202	India: Surat	2020-06-13	QLF98045	Bangladesh	2020-06-17	QKY38005	Australia: Northern Territory	2020
QLH93429	Bangladesh: Jashore	2020-07-07	QLF98048	Bangladesh	2020-06-19	QKY38209	Australia: Victoria	2020-04-10
QLH93441	Bangladesh: Jashore	2020-07-07	QLF98201	India: Talod	2020-06-12	QKY38257	Australia: Victoria	2020-04-10
QLH93453	Bangladesh: Jashore	2020-07-07	QLF98261	India: Rajkot	2020-06-11	QKY38281	Australia: Victoria	2020-04-11
QLH55720	Bangladesh: Barishal	2020-07-06	QLF99991	India: Surat	2020-04-01	QKY38401	Australia: Victoria	2020-04-13
QLH55768	Bangladesh: Barishal	2020-07-06	QLF78310	USA: MD	2020-06-01	QKY38810	USA: Washington, Shohomish County	2020-04-18
QLH55816	Bangladesh: Barishal	2020-07-06	QLF80217	Poland	2020-03-13	QKY38894	USA: Washington, Yakima County	2020-05-03
QLH55840	Bangladesh: Barishal	2020-07-06	QLF95245	USA: Virginia	2020-03	QKY39324	USA: Washington, King County	2020-04-27
QLH56099	Saudi Arabia	2020-02-10	QLF95641	USA: Virginia	2020-03	QKY39588	USA: Washington, Shohomish County	2020-05-06
QLH56231	Saudi Arabia	2020-03-01	QLF95737	USA: Virginia	2020-03	QKY39840	USA: Washington, Yakima County	2020-05-06
QLH56255	Saudi Arabia	2020-03-01	QLF95773	USA: Virginia	2020-03	QKY40164	USA: Washington, Yakima County	2020-05-06
QLH56279	Bangladesh: Barishal	2020-07-06	QLE11150	Bangladesh	2020-06-18	QKY40440	USA: Washington, Yakima County	2020-05-06
QLH57751	USA: FL	2020-04-14	QLC91545	USA: Wisconsin, Dane County	2020-03-20	QKY40716	USA: Washington, Yakima County	2020-05-06
QLH57846	USA: FL	2020-04-14	QLC91617	USA: Wisconsin, Dane County	2020-03-19	QKY41592	USA: Washington, Yakima County	2020-04-22
QLH58037	USA: FL	2020-04-16	QLC91905	USA: Wisconsin, Dane County	2020-03-24	QKY41616	USA: Washington, Yakima County	2020-04-22
QLH58085	USA: FL	2020-04-16	QLC92097	USA: Wisconsin, Dane County	2020-03-31	QKY42204	USA: Washington, Benton County	2020-04-26
QLH58601	USA: FL	2020-05-14	QLC92421	USA: Wisconsin	2020-04-02	QKY42875	USA: Washington, Cowlitz County	2020-04-27
QLH58947	USA: FL	2020-06-02	QLC92553	USA: Wisconsin, Richland county	2020-04-08	QKY42947	USA: Washington, Cowlitz County	2020-04-29
QLH59007	USA: FL	2020-06-03	QLC92601	USA: Wisconsin, Dane County	2020-04-09	QKY26659	USA: Washington, Yakima County	2020-05
QLG75126	Bahrain	2020-06-22	QLC93129	USA: Wisconsin, Milwaukee county	2020-03-21	QKS98544	USA: Washington, King County	2020-03-04
QLG75678	Australia: Victoria	2020-06-01	QLC93357	USA: Wisconsin, Waukesha county	2020-03-24	QKS990192	USA: Washington, King County	2020-02-29
QLG75822	Australia: Victoria	2020-06-06	QLC94305	USA: Wisconsin, Milwaukee county	2020-04-13	QKU28463	USA: Washington, King County	2020-03-03
QLG75930	Australia: Victoria	2020-06-11	QLC94473	USA: Wisconsin, Milwaukee county	2020-04-14	QKU28847	USA: Washington, King County	2020-04-29
QLG75942	Australia: Victoria	2020-06-11	QLC94737	USA: Wisconsin, Milwaukee county	2020-03-24	QKU299039	USA: Washington, King County	2020-04-19
QLG76026	Australia: Northern Territory	2020	QLC46314	USA: FL	2020-04-03	QKU30570	USA: Washington, King County	2020-04-16
QLG76386	Australia: Victoria	2020-06-19	QLC46986	USA: FL	2020-04-21	QKU31182	USA: CA	2020-04-02
QLG76542	Australia: Victoria	2020-06-20	QLC47346	USA: FL	2020-05-03	QKU31266	USA: CA	2020-04-11
QLG97055	Italy	2020-04-04	QLB39261	USA	2020-04-06	QKU31638	USA: CA	2020-03-25
QLG97460	USA: Wisconsin, Dane county	2020-06-15	QLB39321	USA	2020-04-11	QKU31746	USA: CA	2020-03-20
QLG97484	USA: Wisconsin, Jackson county	2020-06-14	QLA47500	USA: Virginia	2020-05	QKU31806	USA: CA	2020-03-30
QLG98012	USA: CA	2020-06-01	QLA47776	USA: Virginia	2020-05	QKU31818	USA: CA	2020-03-30
QLG99677	USA: CA	2020-06-03	QKR84274	Egypt	2020-06-02	QKU32046	USA: CA	2020-05-01
QLG99737	USA: CA	2020-04-16	QKR84421	Egypt	2020-06-02	QKU32202	USA: CA	2020-03-30
QLH00026	USA: CA	2020-04-16	QKS66941	Egypt	2020-06-02	QKU32934	USA: CA	2020-03-24
QLH00290	USA: CA	2020-04-27	QLA100656	USA: Ak	2020-03-23	QKU32982	USA: CA	2020-03-26
QLH00362	USA: CA	2020-04-28	QLA100669	India: Surat	2020-06-11	QKU37034	Saudi Arabia: Jeddah	2020-03-15
QLH00362	USA: CA	2020-04-17	QLA10165	India: Kapsadvanj	2020-06-08	QKU37202	USA: CA	2020-04-18

Table 2: List of accessions of the ORF3a protein, geo-location and respective data collection date

Accession	Geo_Location	Collection_Date	Accession	Geo_Location	Collection_Date	Accession	Geo_Location	Collection_Date	Accession	Geo_Location	Collection_Date
QKU37646	USA: CA	2020-04-02	QKG86518	USA	2020-04	QJR88822	Australia: Victoria	2020-03-20	QJR88822	Australia: Victoria	2020-03-20
QKU52834	USA: Washington,King County	2020-03-18	QKE61733	India: Rajkot	2020-04-28	QJR89110	Australia: Victoria	2020-03-22	QJR89110	Australia: Victoria	2020-03-22
QKU52870	USA: Washington,Shnohish County	2020-03-16	QKE44990	USA	2020-04	QJR89278	Australia: Victoria	2020-03-23	QJR89278	Australia: Victoria	2020-03-23
QKU53050	USA: Washington	2020-03-20	QKE45765	USA: CA	2020-04-26	QJR89362	Australia: Victoria	2020-03-23	QJR89362	Australia: Victoria	2020-03-23
QKU53650	USA: Washington,King County	2020-03-17	QKE45861	USA: CA	2020-04-30	QJR89446	Australia: Victoria	2020-03-24	QJR89446	Australia: Victoria	2020-03-24
QKU53854	USA: Washington,King County	2020-03-07	QKE45885	USA: CA	2020-04-30	QJR91282	Australia: Victoria	2020-03-26	QJR91282	Australia: Victoria	2020-03-26
QKV06224	USA: Washington,Yakima County	2020-04-02	QKE45933	USA: CA	2020-04-29	QJR91354	Australia: Victoria	2020-03-29	QJR91354	Australia: Victoria	2020-03-29
QKV06236	USA: Washington,Pierce County	2020-03-31	QKE10935	Czech Republic	2020-03-31	QJR95110	Australia: Victoria	2020-04-08	QJR95110	Australia: Victoria	2020-04-08
QKV07400	USA: Washington,Yakima County	2020-03-31	QJY78272	USA	2020-03-20	QJQ84173	USA: NEW ORLEANS, LA	2020-04-04	QJQ84173	USA: NEW ORLEANS, LA	2020-04-04
QKV07184	USA: Washington,King County	2020-03-31	QKC05357	USA	2020-03-11	QJQ38625	USA: CA	2020-04-22	QJQ38625	USA: CA	2020-04-22
QKV07340	USA: Washington,Yakima County	2020-04-02	QJY40110	USA	2020-03-17	QJQ39045	USA: MI	2020-03-13	QJQ39045	USA: MI	2020-03-13
QKV07400	USA: Washington,Yakima County	2020-03-26	QJY40506	India: Jmugadh	2020-05-09	QJQ39081	USA: MI	2020-03-16	QJQ39081	USA: MI	2020-03-16
QKV08048	USA: Washington,King County	2020-03-31	QJX68859	USA: Michigan	2020-03-16	QJQ39297	USA: MI	2020-03-18	QJQ39297	USA: MI	2020-03-18
QKS65597	USA: CA	2020-03-15	QJX70192	USA: Michigan	2020-03-30	QJQ39741	USA: MI	2020-03-25	QJQ39741	USA: MI	2020-03-25
QKS65621	USA: CA	2020-03-15	QJX70592	USA: Illinois	2020-04-14	QJQ107211	USA: VA	2020-04	QJQ107211	USA: VA	2020-04
QKS65777	USA: CA	2020-03-16	QJX45032	USA: CA	2020-03-23	QJ154123	USA: CA	2020-03-05	QJ154123	USA: CA	2020-03-05
QKS65849	USA: MA	2020-03-15	QJX45308	Poland	2020-04-11	QJ154254	USA: CA	2020-03-03	QJ154254	USA: CA	2020-03-03
QKS66041	USA: NJ	2020-03-14	QJW00412	India: Gandhinagar	2020-05-02	QJF75396	USA: Michigan	2020-03-20	QJF75396	USA: Michigan	2020-03-20
QKS66053	USA: NJ	2020-03-14	QJX44383	India: Ahmedabad	2020-04-29	QJF77147	USA: WA	2020-04-02	QJF77147	USA: WA	2020-04-02
QKS66305	USA: UT	2020-03-12	QJX44407	India: Ahmedabad	2020-04-29	QJF38451	USA: CA	2020-03-28	QJF38451	USA: CA	2020-03-28
QKS66737	USA: NY	2020-03-15	QJW69308	Germany: Bavaria	2020-03-23	QJD47203	USA: WA	2020-03-26	QJD47203	USA: WA	2020-03-26
QKS67001	USA	2020-04-09	QJU70306	USA: AK	2020-04-01	QJD47299	USA: WA	2020-03-28	QJD47299	USA: WA	2020-03-28
QKS67456	China	2020-01-23	QJV21807	USA: CA	2020-04-01	QJD47419	USA: WA	2020-04-05	QJD47419	USA: WA	2020-04-05
QJY78153	Egypt	2020-05-02	QJW28449	USA: VA	2020-04	QJD47539	USA: CT	2020-04-07	QJD47539	USA: CT	2020-04-07
QKQ63773	USA: Virginia	2020-04	QJW28665	USA: VA	2020-04	QJD47551	USA: CT	2020-04-06	QJD47551	USA: CT	2020-04-06
QKO25735	Bangladesh: Dhaka	2020-06-01	QJU11458	USA: FL	2020-03-06	QJD47849	Taiwan	2020-03-16	QJD47849	Taiwan	2020-03-16
QKO25747	Bangladesh: Dhaka	2020-06-01	QJT72327	France	2020-03-03	QJD47873	Taiwan	2020-03-18	QJD47873	Taiwan	2020-03-18
QKO00487	India: Ahmedabad	2020-05-27	QJT72387	France	2020-03	QJD47956	USA: WA	2020-03-10	QJD47956	USA: WA	2020-03-10
QKN19672	USA: Michigan	2020-04-26	QJT72471	France	2020-03	QJD48484	USA: WA	2020-03-13	QJD48484	USA: WA	2020-03-13
QKN20740	USA	2020-04-04	QJT72507	France	2020-03	QJD20838	Sri Lanka	2020-03-16	QJD20838	Sri Lanka	2020-03-16
QKN20812	USA	2020-04-03	QJT72951	France	2020-03	QJD23478	USA: NY	2020-03-18	QJD23478	USA: NY	2020-03-18
QKN20824	USA	2020-04-04	QJS53735	Greece: Athens	2020-03-12	QJD23730	USA: NY	2020-03-18	QJD23730	USA: NY	2020-03-18
QKM76547	Germany: Dusseldorf	2020-03-15	QJS53831	Greece: Athens	2020-03-13	QJD25758	USA: NY	2020-03-19	QJD25758	USA: NY	2020-03-19
QKM76907	Germany: Heinsberg	2020-02-28	QJS54023	Greece: Athens	2020-03-12	QJC19648	USA: WA	2020-03-31	QJC19648	USA: WA	2020-03-31
QKK12852	Bangladesh	2020-05-23	QJS54155	Greece: Athens	2020-03-08	QJC20380	USA: WA	2020-03-27	QJC20380	USA: WA	2020-03-27
QKK14612	USA	2020-05-11	QJS54191	Greece: Athens	2020-03-23	QJC20500	USA: WA	2020-03-30	QJC20500	USA: WA	2020-03-30
QKG87087	USA: Massachusetts	2020-04-01	QJS54383	Greece: Athens	2020-03-10	QJA17681	USA: PA	2020-03-07	QJA17681	USA: PA	2020-03-07
QKG87159	USA: Massachusetts	2020-04-02	QJS54923	USA: CA	2020-04-30	QJZ13336	USA	2020-03-23	QJZ13336	USA	2020-03-23
QKG87195	USA: Massachusetts	2020-03-27	QJS57052	USA: WA	2020-04-03	QJZ13838	USA	2020-03-22	QJZ13838	USA	2020-03-22
QKG87267	USA: Massachusetts	2020-04-01	QJS39520	Netherlands	2020-04-29	QJZ14498	USA: MA	2020-03-21	QJZ14498	USA: MA	2020-03-21
QKG88539	USA: Massachusetts	2020-04-02	QJS39568	Netherlands	2020-04-29	QJZ16438	USA: MA	2020-03-06	QJZ16438	USA: MA	2020-03-06
QKG90147	USA: Massachusetts	2020-04-01	QJS39616	Netherlands	2020-05-06	QJZ16548	Greece	2020-03-18	QJZ16548	Greece	2020-03-18
QKG90399	USA: Massachusetts	2020-03-21	QJR84550	USA: CA	2020-04-01	QJUI78768	Spain	2020-03-17	QJUI78768	Spain	2020-03-17
QKG90495	USA: Massachusetts	2020-03-26	QJR84790	USA: CA	2020-04-13	QIU81286	USA: WA	2020-03-13	QIU81286	USA: WA	2020-03-13
QKG90867	USA: Massachusetts	2020-03-25	QJR86050	Australia: Victoria	2020-03-15	QIS61075	USA: IL	2020-03-16	QIS61075	USA: IL	2020-03-16
QKG91107	USA: Massachusetts	2020-03-27	QJR87574	Australia: Victoria	2020-03-20	QIS61315	USA: WA	2020-03-18	QIS61315	USA: WA	2020-03-18
QKG64052	USA	2020-04	QJR87598	Australia: Victoria	2020-03-21	QIS30116	USA: San Francisco, CA	2020-03-16	QIS30116	USA: San Francisco, CA	2020-03-16
QKG81824	USA: Virginia	2020-04	QJR87730	Australia: Victoria	2020-03-21	QIH57239	USA	2020-02-25	QIH57239	USA	2020-02-25
QKG81932	USA: Virginia	2020-04	QJR88306	Australia: Victoria	2020-03-23	QHZ00380	South Korea	2020-01	QHZ00380	South Korea	2020-01
		2020-04	QJR88390	Australia: Victoria	2020-03-23						

- 75
- Keep iterating until there is no change to the centroids.
 - Find the sum of the squared distance between data points and all centroids.
 - Assign each data point to the closest cluster (centroid).
 - Compute the centroids for the clusters by taking the average of the all data points that belong to each cluster.

In this study we did clustering using Matlab by customizing the value of K and inputting the frequency of amino acid
80 compositions over the ORF3a proteins.

2.2.3. Amino Acid Conservation Shannon Entropy

How conserved/disordered the amino acids are, over ORF3a protein is addressed by the information theoretic measure known as 'Shannon entropy(SE)' which we deploy here to find out conservation entropy of each ORF3a protein. For each ORF3a protein, Shannon entropy of amino acid conservation over the amino acid sequence of ORF3a protein is computed
85 using the following formula [39]:

For a given amino acid sequence of ORF3a protein of length l , the conservation of amino acids is calculated as follows:

$$SE = - \sum_{i=1}^{20} p_{s_i} \log_{20}(p_{s_i})$$

where $p_{s_i} = \frac{k_i}{l}$; k_i represents the number of occurrences of an amino acid s_i in the given sequence.

In this study, SE describes the wide variety of 296 distinct ORF3a proteins collected from various countries across the world.

3. Results

90 All mutations, compared to Chinese Wuhan sequence, over the set of distinct ORF3a proteins are detected and consequently they have been classified based on their predicted effect as disease/neutral in important functions of ORF3a protein (Table 9). Also, some important known domains are identified for the observed mutations and accordingly the predicted effect of mutations in protein functions have been discussed. Further, consecutive mutations observed in ORF3a proteins according to the timelines of detection of various mutations for a subgroup of ORF3a proteins located
95 in Australia, Bangladesh, India, USA and so on is derived (Fig.7-11). Using a web-server ($i - MUTANT : http : //gpcr2.biocomp.unibo.it/cgi/predictors/I - Mutant3.0/I - Mutant3.0.cgi$) stability of ORF3a protein structures were predicted upon various mutations. At last, twenty clusters are formed using K-means clustering method based on frequency probability of amino acids of 296 ORF3a proteins. The wide variations of 296 ORF3a proteins are finally supported by the Shannon entropy (SE) and remarkably we found the most widest varieties of ORF3a proteins in virus detected in the
100 USA.

3.1. Mutations over the ORF3a protein of SARS-CoV2

Each of the ORF3a amino acid sequences (fasta formatted) are aligned with respect to the ORF3a protein (YP_009724391.1) from China-Wuhan using multiple sequence alignment tool (NCBI Blastp suite) and found the mutations and their associated positions were detected accordingly [40]. It is noted that a mutation from an amino acid A_1 to A_2 at a position p
105 is denoted by A_1pA_2 or $A_1(p)A_2$. The Fig.2 describes various mutations with their respective locations. The mutations

are found in the entire ORF3a sequence starting from the amino acid position 7 to 271. It is found that an amino acid at a fixed position mutates to two different amino acids. For examples, at 9th position of the reference ORF3a protein, the amino acid Threonine(T) maps to Isoleucine(I) and Lysine(K) in different ORF3a proteins. At the 18th position Glycine maps to three amino acids Valine, Serine and Cysteine. The amino acid Alanine(A) maps to Valine, Serine, Threonine and Aspartic acid at the 99th position.

From	I	F	T	T	A	V	V	T	Q	G	G	G	A	D	P	P	S	D	D	T	A	T	I	Q	Q	A	L	L	P	S	F	G
Position	7	8	9	10	13	13	14	17	18	18	18	23	23	25	25	26	27	27	32	33	34	35	38	38	39	41	41	42	42	43	44	
To	T	L	I	K	S	L	A	I	R	V	S	C	S	Y	L	S	L	H	Y	I	S	A	T	P	E	T	I	F	R	L	Y	V
From	W	L	V	V	V	A	L	L	A	V	V	F	Q	V	A	K	T	L	K	K	W	W	W	A	L	S	S	K	K	V	H	
Position	45	46	48	50	50	51	52	53	54	55	55	56	57	58	59	61	64	65	66	67	69	69	69	72	73	74	74	75	75	77	78	
To	L	F	F	I	A	S	I	F	S	F	G	C	H	H	L	V	N	I	F	N	N	L	R	C	S	F	P	F	R	E	F	Y
From	L	L	V	V	H	L	L	L	V	A	A	A	A	G	A	A	P	L	L	L	V	F	Q	S	I	F	I	G	M	R	L	W
Position	83	86	88	90	93	94	94	95	97	99	99	99	99	100	102	103	104	106	108	111	112	114	116	117	118	120	123	124	125	126	127	128
To	F	W	A	F	Y	P	F	F	A	V	S	T	D	V	V	S	S	F	F	S	F	C	H	L	V	F	V	V	I	S	I	L
From	T	E	W	W	W	W	R	R	S	L	A	D	C	T	N	C	Y	D	I	S	S	T	S	G	G	T	T	P	H	Q	G	T
Position	128	128	131	131	131	131	134	135	140	143	145	148	151	152	153	154	155	158	165	165	170	171	172	172	175	176	178	182	185	188	190	
To	A	L	R	C	S	L	L	C	P	F	S	Y	Y	I	S	Y	C	Y	T	L	I	S	L	V	C	I	S	Y	H	C	I	
From	W	W	E	S	G	G	V	V	D	Y	Q	Y	S	Q	L	S	T	G	R	T	V	V	D	D	E	P	E	Q	G	G	G	
Position	193	193	194	195	196	196	197	197	210	210	211	213	215	216	218	219	220	221	224	226	229	237	237	238	238	239	240	241	245	251	254	
To	R	C	Q	Y	R	L	I	I	Y	C	H	H	P	R	V	N	I	C	M	I	A	F	N	E	D	L	V	L	V	C	R	
From	V	V	N	N	V	M	M	P	P	I	Y	S	P	T																		
Position	255	256	257	257	259	260	260	262	262	263	264	265	267	271																		
To	L	I	Q	D	E	I	K	S	L	M	C	F	L	I																		

Reference ORF3a: YP_009724391.1

Figure 2: Mutations in the respective position in ORF3a protein sequence compared with reference Wuhan sequence YP_009724391.1. **Note:** From: existing amino acid in reference sequence; position: amino acid position in the sequence; To: mutated amino acid in studied sequence

Based on observed mutations, it is noticed that amino acids Alanine(A) and Tryptophan(W) are found to be most vulnerable to mutate to various amino acids. It is noted that the mutation of Tryptophan (W) at 131 position are found in the Cystine-rich domain (127-133).

Distinct mutations and its associated mutation of frequency are presented in Table-3. The most frequent mutation over the ORF3a is to be Q57H (Acidity: Neutral(Q) to Basic(weakly)(H)) with frequency 142. A pie chart accounting the frequency distribution of various mutations is shown in Fig.3. In addition to the list of mutations (Fig.2), two deletion and two insertion mutations were found in five different ORF3a proteins at various positions.

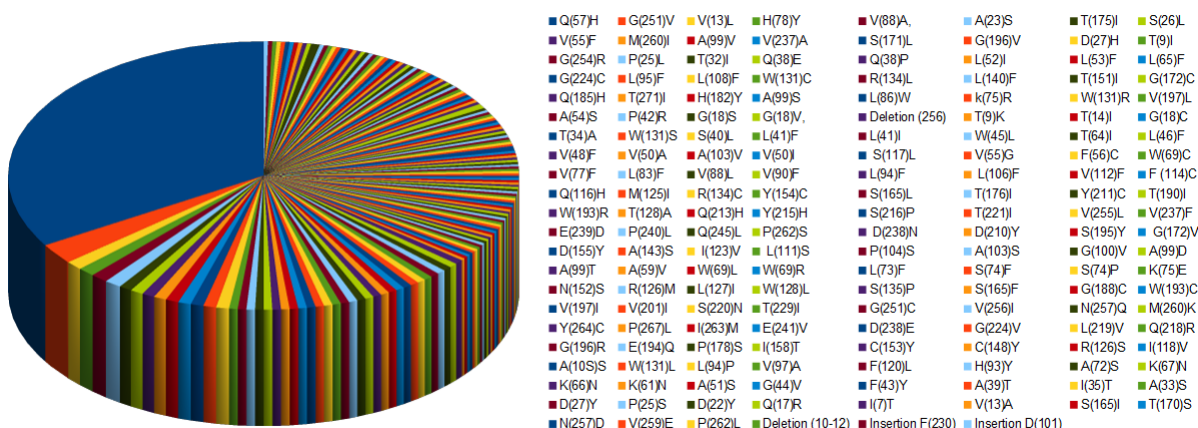


Figure 3: Pie chart of the frequency of distinct mutations

The details of mutations, in the 256 ORF3a unique proteins from viruses of 256 patients, in specific domain(s) and predicted effects of mutations viz. disease and neutral effects through the web-server Meta-SNP (<https://snps.biofold.org/meta-snp/>) are presented in the Tables-4, 5, 6, 7 & 8. Note that among 296 ORF3a proteins, 40 sequences possess only ambiguous mutations which we have neglected. A snapshot of predicted result (disease causing variant with reliability score 3) of the

Table 3: Distinct mutations across the ORF3a proteins and their respective frequency

Mutations	Q(57)H	G(251)V	A(23)S	H(75)Y	V(13)L	V(88)A	A(99)V	D(27)H	G(196)V	M(260)I	S(171)L	S(26)L
Frequency of Mutations	124	9	4	4	4	4	3	3	3	3	3	3
Mutations	T(175)I	V(237)A	V(55)F	A(54)S	A(99)S	D(155)Y	D(22)Y	Deletion (256)	G(I72)C	G(18)S	G(18)V	G(224)C
Frequency of Mutations	3	3	3	2	2	2	2	2	2	2	2	2
Mutations	G(254)R	H(182)Y	k(75)R	L(108)F	L(140)F	L(52)I	L(53)F	L(65)F	L(86)W	L(95)F	P(25)L	P(42)R
Frequency of Mutations	2	2	2	2	2	2	2	2	2	2	2	2
Mutations	Q(185)H	Q(38)E	Q(38)P	R(134)L	T(151)I	T(271)I	T(32)I	T(9)I	V(197)L	W(128)L	W(131)C	W(131)R
Frequency of Mutations	2	2	2	2	2	2	2	2	2	2	2	2
Mutations	D(238)N	G(172)V	I(123)V	L(106)F	L(111)S	S(117)L	A(103)S	A(103)V	A(105)S	A(143)S	A(33)S	A(39)T
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	A(51)S	A(59)V	A(72)S	A(99)D	A(99)T	C(148)Y	C(153)Y	D(210)Y	D(238)E	D(27)Y	Deletion IGTT(10-12)	E(194)Q
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	E(239)D	E(241)V	F(114)C	F(120)L	F(43)Y	F(56)C	G(100)V	G(18)C	G(188)C	G(196)R	G(224)V	G(251)C
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	G(44)V	H(93)Y	I(118)V	I(158)T	I(263)M	I(35)T	I(7)T	Insertion D(101)	Insertion F(230)	K(61)N	K(66)N	K(67)N
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	K(75)E	L(127)I	L(219)V	L(41)F	L(41)I	L(46)F	L(73)F	L(83)F	L(94)F	L(94)P	M(125)I	M(260)K
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	N(152)S	N(237)D	N(257)Q	P(104)S	P(178)S	P(240)L	P(25)S	P(262)L	P(262)S	P(267)L	Q(116)H	Q(17)R
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	Q(213)H	Q(218)R	Q(245)L	R(126)M	R(126)S	R(134)C	S(135)P	S(165)F	S(165)I	S(165)L	S(195)Y	S(216)P
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	S(220)N	S(40)L	S(74)F	S(74)P	T(128)A	T(14)I	T(170)S	T(176)I	T(190)I	T(221)I	T(229)I	T(34)A
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	T(64)I	T(9)K	V(112)F	V(13)A	V(197)I	V(201)I	V(237)F	V(255)L	V(256)I	V(259)E	V(48)F	V(50)A
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	V(50)I	V(55)G	V(77)F	V(88)L	V(90)F	V(97)A	W(131)L	W(131)S	W(193)C	W(193)R	W(46)L	W(59)A
Frequency of Mutations	1	1	1	1	1	1	1	1	1	1	1	1
Mutations	W(69)L	W(69)R	Y(154)C	Y(211)C	Y(215)H	Y(264)C						
Frequency of Mutations	1	1	1	1	1	1						

*124 is the frequency of the mutation Q to H occurred at the 57th position.

most frequent mutation Q57H is shown in Fig.4.

Meta-SNP
Meta-predictor of disease causing variants

Mutation	PANTHER	PhD-SNP	SIFT	SNAP	Meta-SNP	RI
Q57H	NA -	Disease 0.673	Disease 0.000	Disease 0.730	Disease 0.637	3

Prediction:
Neutral: Neutral variants
Disease: Disease causing variants

Outputs: Value reported under each prediction
PANTHER: Between 0 and 1. (If >0.5 mutation is predicted Disease)
PhD-SNP: Between 0 and 1. (If >0.5 mutation is predicted Disease)
SIFT: Positive Value (If >0.05 mutation is predicted Neutral)
SNAP: Output normalized between 0 and 1 (If >0.5 mutation is predicted Disease)
Meta-SNP: Between 0 and 1. (If >0.5 mutation is predicted Disease)

RI: Reliability Index between 0 and 10.

Figure 4: A snapshot of the predicted effect of the frequently occurred mutation Q57H in ORF3a using Meta-SNP web-server

Based on the predicted type of mutations, all the 256 ORF3a proteins are classified into three classes which are presented in the Table 9. The three classes representing disease, neutral and mixture of disease as well as neutral mutations are constituted of protein IDs with respective geo-locations.

125

Table 4: protein IDs and respective mutations, geo-locations, total number of mutations in the protein, domains and predicted effect of the mutations

Protein ID	Country	Mutations	Total Mutations	Domain	Effect of mutation(s)	RI
QJD47419.1	USA	T(9)I	1	FD I	Disease (0.649)	3
QLH01250.1	USA	V(13)L	1	FD I	Neutral (0.119)	8
QLB39261.1	USA	T(14)I	1	FD I	Disease (0.650)	3
QJW69308.1	GERMANY	P(25)L	1		Neutral (0.125)	8
QKV38281.1	AUSTRALIA	S(26)L	1		Neutral (0.157)	7
QKS67456.1	CHINA	T(32)I	1		Disease (0.652)	3
QJS39568.1	Netherlands	T(34)A	1		Neutral (0.297)	4
QLH93429.1	Bangladesh	Q(38)E	1	FD II (TRAF3 binding domain)	Disease (0.631)	3
QLC46986.1	USA	Q(38)P	1	FD II (TRAF3 binding domain)	Disease (0.638)	3
QKE61733.1	India	L(41)F	1	FD II	Neutral (0.114)	8
QKV41616.1	USA	L(41)I	1	FD II	Neutral (0.266)	5
QJR88306.1	Australia	L(46)F	1	TransmembraneDomain I (FD II)	Neutral (0.114)	8
QLF97772.1	Bangladesh	V(48)F	1	TransmembraneDomain I (FD II)	Disease (0.717)	4
QLF95641.1	USA	Q(57)H	1	TransmembraneDomain I	Disease (0.637)	3
QJD23478.1	USA	V(50)A	1	TransmembraneDomain I	Disease (0.599)	2
QJR89110.1	AUSTRALIA	L(52)I	1	TransmembraneDomain I	Neutral (0.454)	1
QKG64052.1	USA	F(56)C	1	TransmembraneDomain I	Disease (0.673)	3
QLH58601.1	USA	Q(57)H	1	TransmembraneDomain I	Disease (0.637)	3
QKU53050.1	USA	Q(57)H	1	TransmembraneDomain I	* Disease (0.637)	3
QLA10225.1	India	Q(57)H	1	TransmembraneDomain I	Disease (0.637)	3
QJC20380.1	USA	Q(57)H	1	TransmembraneDomain I	Disease (0.637)	3
QKO25747.1	Bangladesh	W(69)L	1	TransmembraneDomain I	Disease (0.625)	3
QKX47995.1	Bangladesh	W(69)R	1		Disease (0.650)	3
QJT72387.1	France	L(73)F	1		Disease (0.623)	2
QLG75930.1	Australia	S(74)F	1		Neutral (0.478)	0
QKV38257.1	Australia	S(74)P	1		Disease (0.657)	3
QKQX49024.1	Bangladesh	K(75)E	1		Disease (0.649)	3
QKU37034.1	Saudi Arabia	W(88)A	1	FD III	Disease (0.636)	3
QKQ63773.1	USA	L(106)F	1	FD III	Disease (0.631)	3
QKU32202.1	USA	L(106)F	1	FD III	Disease (0.631)	3
QKV40716.1	USA	R(126)M	1	FD III	Disease (0.696)	4
QJZ16548.1	Greece	L(127)I	1	FD III (cysteine rich domain)	Neutral(0.447)	1
QKE45861.1	USA	W(128)L	1	FD III (cysteine rich domain)	Disease (0.675)	4
QJD47873.1	Taiwan	W(131)C	1	FD III (cysteine rich domain)	Disease (0.666)	3
QKV35688.1	USA	W(131)R	1	FD III (cysteine rich domain)	Disease (0.717)	4
QLC93357.1	USA	R(134)L	1	FD III	Disease(0.712)	4
QJL57239.2	USA	S(135)P	1	FD III	Disease(0.688)	3
QKU53854.1	USA	L(140)F	1	FD III	Disease(0.595)	2
QLF98261.1	India	T(151)I	1	FD III	Disease(0.624)	2
QKV07340.1	USA	S(165)F	1	FD VI (SGD motif)	Disease (0.614)	2
QLF80217.1	Brazil	S(171)L	1	FD VI (SGD motif)	Disease (0.602)	2
QLL50570.1	USA	G(172)C	1		Disease(0.646)	3
QLH59007.1	USA	T(175)I	1		Disease(0.728)	5
QKE10935.1	Bangladesh	W(188)C	1		Disease (0.668)	2
QLC92601.1	Czech Republic	W(193)C	1		Disease (0.600)	3
QKK14612.1	USA	V(197)I	1		Neutral (0.330)	3
QKU28463.1	USA	V(201)I	1		Disease (0.509)	0
QLF97844.1	Bangladesh	S(220)N	1		Neutral(0.255)	6
QKX46204.1	USA	T(229)I	1		Neutral (0.422)	1
QLH01382.1	USA	V(237)A	1		Disease(0.648)	3
					DiseaseE(0.583)	2

* Disease(0.637) denotes the effect of the mutation Q(57)H as 'disease' with the degree 0.637.

Table 5: protein IDs and respective mutations, geo-locations, total number of mutations, domain and predicted effect of the mutations

Protein ID	Country	Mutations	Total Mutations	Domain	Effect of mutation(s)	RI
QJY78272.1	USA	P(240)L	1		Disease(0.583)	2
QKU52834.1	USA	G(251)C	1		Disease(0.713)	4
QKU31182.1	USA	M(260)K	1		Disease(0.632)	3
QJX70592.1	USA	Y(264)C	1		Disease(0.651)	3
QLC92421.1	USA	P(267)L	1		Disease(0.525)	1
QKC05357.1	USA	T(271)I	1		Neutral(0.255)	5
QJD20838.1	Shri Lanka	I(263)M	1		Disease(0.510)	0
QKV07184.1	USA	G(254)R	1		Disease(0.728)	5
QLH57846.1	USA	G(251)V	1		Disease(0.770)	5
QJR89362.1	Australia	G(251)V	1		Disease(0.770)	5
QJS54191.1	Greece	E(241)V	1		Neutral(0.061)	9
QKV06236.1	USA	D(238)E	1		Neutral(0.244)	5
QJD47956.1	USA	G(224)V	1		Disease(0.686)	4
QJS39520.1	Netherlands	L(219)V	1		Neutral(0.137)	7
QKV42204.1	USA	Q(218)R	1		Disease(0.584)	2
QKG90867.1	USA	G(196)R	1		Disease(0.664)	3
QLG76026.1	Australia	G(196)V	1		Disease(0.687)	4
QLH56279.1	Bangladesh	E(194)Q	1		Neutral(0.140)	7
QJS39616.1	Netherlands	H(182)Y	1		Neutral(0.139)	7
QLG76386.1	Australia	P(178)S	1		Disease(0.565)	1
QLF97736.1	Bangladesh	G(172)V	1	FD VI (SGD motif)	Disease(0.646)	3
QLL51782.1	USA	I(158)T	1		Disease(0.734)	5
QLC92097.1	USA	D(155)Y	1		Disease(0.829)	7
QIS61315.1	USA	C(153)Y	1	FD VI	Disease(0.692)	4
QJF77147.1	USA	C(148)Y	1	FD VI	Disease(0.785)	6
QJE38451.1	USA	R(126)S	1	FD III	Disease(0.671)	3
QKS67001.1	USA	I(118)V	1	FD III	Neutral(0.063)	9
QLF78310.1	Poland	A(99)S	1	FD III	Disease(0.577)	2
QKO25735.1	Bangladesh	A(99)V	1	FD III	Disease(0.602)	2
QLF95773.1	USA	H(93)Y	1	FD III	Disease(0.649)	3
QLG75678.1	Australia	H(78)Y	1	FD III	Neutral(0.349)	3
QJZ14498.1	USA	A(72)S	1	TD2	Disease(0.580)	2
QKK12852.1	Bangladesh	K(67)N	1		Disease(0.551)	1
QKY59990.1	India	K(66)N	1		Neutral(0.031)	9
QJD23730.1	USA	K(61)N	1		Disease(0.622)	2
QLF98048.1	Bangladesh	A(54)S	1	TD1	Disease(0.613)	2
QLC94305.1	USA	A(39)T	1	FD II	Disease(0.648)	3
QLF95737.1	USA	Q(57)H	1	TD I	Disease(0.637)	3
QJT72951.1	France	A(33)S	1	TD I	Disease(0.578)	2
QKG81824.1	USA	D(27)H	1		Neutral(0.139)	7
QKU53650.1	USA	D(27)Y	1		Neutral(0.220)	6
QLH93202.1	India	A(23)S	1		Neutral(0.494)	0
QLH55768.1	Bangladesh	D(22)Y	1		Neutral(0.187)	6
QLH55720.1	Bangladesh	G(18)V	1		Neutral(0.036)	9
QKW88844.1	USA	Q(17)R	1		Neutral(0.139)	7
QKV07400.1	USA	I(7)T	1	FD I	Neutral(0.213)	6
QKW89480.1	USA	V(13)A	1	FD I	Neutral(0.175)	7
QLH01334.1	USA	V(13)L	1	FD I	Neutral(0.119)	8
QLH00290.1	USA	S(26)L	1		Neutral(0.157)	7
QKS66305.1	USA	Q(57)H	1	TD I	Disease(0.637)	3

Table 6: protein IDs and respective mutations, geo-locations, total number of mutations, domain and predicted effect of the mutations

Protein ID	Country	Mutations	Total Mutations	Domain	Effect of mutation(s)	RI
QKS65597.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QJS54923.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QJY78153.1	Egypt	Q(57)H	1	TD I	Disease (0.637)	3
QJQ39081.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QKV08048.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QKE45933.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QLJ51746.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QLG99773.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QLH00362.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QKU37646.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QKS65849.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QJC20500.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QKU32046.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QJU11458.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QJR87730.1	Australia	Q(57)H	1	TD I	Disease (0.637)	3
QKU31638.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QKU31746.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QJR89278.1	Australia	Q(57)H	1	TD I	Disease (0.637)	3
QJR89446.1	Australia	Q(57)H	1	TD I	Disease (0.637)	3
QLH01238.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QJQ39297.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QLB39321.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QLH01298.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QLG99737.1	USA	V(237)A	1	TD I	Disease (0.637)	2
QKV38401.1	Australia	V(259)E	1	DiseaseE(0.583)	Disease (0.637)	2
QJU78768.1	Spain	G(196)V	1	Disease (0.687)	Disease (0.687)	4
QKS89844.1	USA	P(262)L	1	Disease (0.687)	Disease (0.687)	4
QJD47539.1	USA	K(75)R	1	Disease (0.601)	Disease (0.601)	2
QJZ14498.1	USA	A(72)S	1	Disease (0.595)	Disease (0.595)	2
QJS54023.1	Greece	G(251)V	1	Disease (0.580)	Disease (0.580)	2
QKV35400.1	USA	W(131)R	1	Disease (0.770)	Disease (0.770)	5
QKU37202.1	USA	Q(57)H	1	FD III (cysteine rich domain)	Disease (0.717)	4
QIS30116.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QJR91354.1	Australia	Q(57)H	1	TD I	Disease (0.637)	3
QKN20740.1	USA	Q(57)H	1	TD I	Disease (0.637)	3
QIU81286.1	USA	F(8)L, deletion mutation(10-12)	1	FD I	Disease (0.637)	3
QLH55840.1	Bangladesh	Deletion (256)	1	FD I	Disease (0.642)	3
QJR84790.1	USA	Insertion F(230)	1	FD III	Neutral (0.055), Disease(0.637)	9,3
QLG97055.1	Italy	Insertion D(101)	1	TD I	Disease(0.649), Disease(0.637)	3,3
QLI50282.1	USA	G(18)S, Q(57)H	2	FD I, TD I	Disease(0.747), Disease(0.728)	5,5
QKS66053.1	USA	T(9)I, Q(57)H	2	FD I	Neutral (0.119), Neutral (0.142)	8, 7
QJC19648.1	USA	T(9)K, G(254)R	2	FD I, TD II	Neutral (0.119), Neutral (0.349)	8,3
QJR88390.1	AUSTRALIA	V(13)L, T(175)I	2	FD I, TD II	Disease (0.637), Neutral(0.134)	3,7
QJR88822.1	AUSTRALIA	V(13)L, H(78)Y	2	TD I,	Disease (0.637), Neutral(0.125)	3,8
QLI46290.1	USA	Q(57)H, G(18)C	2	TD I	Disease (0.637), Neutral(0.157)	3,7
QKV40164.1	USA	Q(57)H, P(25)L	2	TD I	Disease (0.637), Disease(0.652)	3,3
QLG97460.1	USA	Q(57)H, S(26)L	2	TD I	Disease (0.631), Neutral(0.349)	3,3
QJV21807.1	USA	Q(57)H, T(32)I	2	TD I	Disease(0.638), Disease(0.674)	3,3
QLF98036.1	Bangladesh	Q(38)E, H(78)Y	2	FD II (TRAF 3 binding domain), TDII	Disease(0.631), Neutral(0.349)	3,3
QKG81932.1	USA	Q(38)P, W(131)S	2	FD II (TRAF 3 binding domain); FDIII	Disease(0.638), Disease(0.674)	3,3

Table 7: protein IDs and respective mutations, geo-locations, total number of mutations, domain and predicted effect of the mutations

Protein ID	Country	Mutations	Total Mutations	Domain	Effect of mutation(s)	RI
QLI50414.1	USA	Q(57)H, S(40)L	2	TD II, FD II	Disease(0.637),Disease (0.628)	3,3
QLH93441.1	Bangladesh	W(45)L, T(64)I	2	FDII, TD I	Disease(0.664),Neutral(0.166)	3,7
QLF97952.1	India	V(50)I, A(103)V	2	TDI, FDIII	Disease(0.588),Neutral(0.139)	2,7
QKE45885.1	USA	Q(57)H, L(52)I	2	TDI, TDI	Disease (0.637),Neutral(0.454)	3,1
QKS65621.1	USA	Q(57)H, L(53)F	2	TDI, TDI	Disease (0.637),Disease(0.601)	3,2
QKU29039.1	USA	Q(57)H, V(55)F	2	TDI, TDI	Disease (0.637),Disease(0.702)	3,4
QKS66941.1	Egypt	V(55)F, S(117)L	2	TDI, FD III	Disease (0.702),Disease(0.623)	4,2
QLG76542.1	AUSTRALIA	Q(57)H, V(55)G	2	TDI, TDI	Disease(0.637),Disease(0.649)	3,3
QLA47500.1	USA	Q(57)H, L(65)F	2	TDI	Disease (0.637),Neutral(0.233)	3,5
QKN20812.1	USA	Q(57)H, W(69)C	2	TDI	Disease(0.637), Disease (0.642)	3,3
QIX44383.1	India	Q(57)H, V(77)F	2	TDI, TD II	Disease (0.637), Neutral (0.079)	3, 8
QLF95245.1	USA	Q(57)H, L(83)F	2	TDI, FD III	Disease(0.637), Disease(0.636)	3,3
QLI50222.1	USA	Q(57)H, V(88)L	2	TDI, FD III	Disease(0.637), Disease(0.665)	3,3
QLC94737.1	USA	Q(57)H, V(90)F	2	TDI, FD III	Disease(0.637),Disease(0.615)	3,2
QKG87087.1	USA	Q(57)H, L(94)F	2	TDI, FD III	Disease(0.637),Neutral(0.146)	3,7
QIZ13838.1	USA	Q(57)H, L(95)F	2	TDI, FD III	Disease(0.637),Disease(0.601)	3,2
QIQ84173.1	USA	Q(57)H, L(106)F	2	TDI, FD III	Disease (0.637),Disease(0.631)	3,3
QKG88539.1	USA	Q(57)H, L(108)F	2	TDI, FD III	Disease (0.637),Neutral(0.367)	3,3
QJY40110.1	USA	Q(57)H, V(112)F	2	TDI, FD III	Disease (0.637),Disease(0.621)	3,2
QJD47551.1	USA	Q(57)H, F(114)C	2	TDI, FD III	Disease (0.637),Disease(0.624)	3,2
QJD25758.1	USA	Q(57)H, Q(116)H	2	TDI, FD III	Disease (0.637),Disease(0.714)	3,4
QJD47849.1	Taiwan	Q(57)H, M(125)I	2	TDI, FD III	Disease (0.637),Disease(0.680)	3,4
QKU30570.1	USA	Q(57)H, W(131)C	2	TDI, FD III	Disease (0.637),Disease(0.666)	3,3
QKG90399.1	USA	Q(57)H, R(134)C	2	TDI, FD III	Disease (0.637),Disease(0.717)	3,4
QLF98201.1	India	Q(57)H, R(134)L	2	TDI, FD III	Disease (0.637),Disease(0.712)	3,4
QJR95110.1	AUSTRALIA	Q(57)H, L(140)F	2	TDI, FD III	Disease (0.637),Disease(0.595)	3,2
QIZ13336.1	USA	Q(57)H, T(151)I	2	TDI, FD III	Disease(0.637),Disease(0.624)	3,2
QJT72507.1	France	Q(57)H, Y(154)C	2	TDI, FD III	Disease(0.637),Disease(0.752)	3,5
QKV06224.1	USA	Q(57)H, S(165)L	2	TDI, TDI	Disease (0.637),Disease(0.592)	3,2
QLH58947.1	USA	Q(57)H, G(172)C	2	TDI, FD VI (SGD motif)	Disease(0.637),Disease(0.646)	3,3
QJ07211.1	USA	Q(57)H, T(176)I	2	TDI	Disease (0.637),Neutral(0.184)	3,6
QLH58085.1	USA	Q(57)H, Q(185)H	2	TDI	Disease (0.637),Disease(0.636)	3,3
QKO00487.1	India	Q(57)H, T(190)I	2	TDI	Disease (0.637),Neutral(0.118)	3,7
QKV39588.1	USA	Q(57)H, W(193)R	2	TDI	Disease (0.637),Neutral(0.067)	3,9
QKV38810.1	USA	Q(57)H, T(128)A	2	TDI	Disease (0.637),Disease(0.641)	3,3
QLC47346.1	USA	Q(57)H, Q(213)H	2	TDI, FD III	Disease (0.637),Disease(0.641)	3,3
QKG91107.1	USA	Q(57)H, Y(215)H	2	TDI	Disease (0.637),Neutral(0.139)	3,7
QLH56255.1	Saudi Arabia	Q(57)H, S(216)P	2	TDI	Disease (0.637),Disease(0.661)	3,3
QIQ39045.1	USA	Q(57)H, T(221)I	2	TDI	Disease(0.637),Disease(0.656)	3,3
QJU70306.1	USA	Q(57)H, G(224)C	2	TDI	Disease(0.637), Disease(0.693)	3, 4
QLG75126.1	Bahrain	Q(57)H, V(255)L	2	TDI	Disease (0.637),Disease(0.588)	3,2
QJT72327.1	France	Q(57)H, V(237)F	2	TDI	Disease (0.637),Disease(0.648)	3,3
QIZ16438.1	USA	Q(57)H, E(239)D	2	TDI	Disease (0.637),Neutral(0.051)	3,9
QLI51038.1	USA	Q(57)H, P(240)L	2	TDI	Disease(0.637),Disease(0.583)	3,2
QKG86518.1	USA	Q(57)H, Q(245)L	2	TDI	Disease (0.637),Disease(0.625)	3,3
QLG75942.1	Australia	Q(57)H, M(260)I	2	TDI	Disease(0.637), Disease (0.563)	3, 1

Table 8: protein IDs and respective mutations, geo-locations, total number of mutations, geo-locations, domain and predicted effect of the mutations

Protein ID	Country	Mutations	Total Mutations	Domain	Effect of mutation(s)	RI
QKU28847.1	USA	Q(57)H, M(260)I	2	TDI	Disease (0.637), Disease (0.563)	3, 1
QLI49698.1	India	Q(57)H, T(271)I	2	TDI	Disease(0.637), Neutral (0.255)	3, 5
QKV37633.1	Australia	Q(57)H, P(262)S	2	TDI	Disease (0.637), Disease (0.601)	3, 2
QKG90495.1	USA	Q(57)H, D(238)N	2	TDI	Disease(0.637), Neutral(0.144)	3, 7
QLH58037.1	USA	Q(57)H, D(210)Y	2	TDI	Disease(0.637), Disease (0.610)	3, 2
QJX68859.1	USA	Q(57)H, S(195)Y	2	TDI	Disease (0.637), Disease (0.653)	3, 3
QKR84274.1	Egypt	Q(57)H, H(182)Y	2	TDI	Disease (0.637), Neutral(0.139)	3, 7
QKV38894.1	Egypt	Q(57)H, G(172)V	2	TDI, FD VI (SGD motif)	Disease (0.637), Disease (0.646)	3, 3
QJIS4155.1	Greece	Q(57)H, D(155)Y	2	TDI	Disease(0.637), Disease (0.829)	3, 7
QJX44407.1	India	Q(57)H, A(143)S	2	TDI, FD IV (Caveolin binding motif)	Disease (0.637), Disease (0.604)	3, 2
QKG87267.1	India	Q(57)H, I(123)V	2	TDI, FD III	Disease (0.637), Neutral (0.139)	3, 7
QJIS57052.1	USA	Q(57)H, L(111)S	2	TDI, FD III	Disease (0.637), Disease (0.636)	3, 3
QKG87195.1	USA	Q(57)H, P(104)S	2	TDI, FD III	Disease (0.637), Neutral (0.143)	3, 7
QLH93453.1	Bangladesh	Q(57)H, A(103)S	2	TDI, FD III	Disease (0.637), Neutral (0.448)	3, 1
QJIS61075.1	USA	Q(57)H, G(100)V	2	TDI, FD III	Disease (0.637), Disease (0.711)	3, 7
QJW28449.1	USA	Q(57)H, A(99)D	2	TDI, FD III	Disease (0.637), Disease (0.723)	3, 4
QLH56231.1	Saudi Arabia	Q(57)H, A(99)S	2	TDI, FD III	Disease (0.637), Disease (0.577)	3, 2
QLC91905.1	USA	Q(57)H, A(99)T	2	TDI, FD III	Disease (0.637), Disease (0.602)	3, 2
QJW72471.1	France	Q(57)H, A(99)V	2	TDI, FD III	Disease (0.637), Disease (0.602)	3, 2
QJW00412.1	India	Q(57)H, L(86)W	2	TDI, FD III	Disease (0.637), Disease (0.664)	3, 3
QKG88935.1	USA	Q(57)H, L(86)W	2	TDI, FD III	Disease (0.637), Disease(0.664)	3, 4
QK91545.1	USA	Q(57)H, H(78)Y	2	TDI, TD II	Disease (0.637), Neutral (0.349)	3, 3
QKV38005.1	Australia	Q(57)H, K(75)R	2	TDI, TD II	Disease (0.637), Disease (0.595)	3, 2
QKN20824.1	USA	Q(57)H, A(59)V	2	TDI, TDI	Disease (0.637), Disease (0.622)	3, 2
QKV38209.1	Australia	W(69)L, G(251)V	2	FD III	Disease (0.625), Disease (0.770)	3, 5
QLA09656.1	Australia	V(88)A, G(251)V	2	FD III	Disease(0.601), Neutral(0.189)	2, 6
QJD47203.1	USA	L(95)F, N(152)S	2	FD III	Disease (0.675), Disease (0.770)	2, 6
QHZ00380.1	South Korea	W(128)L, G(251)V	2	FD III	Disease (0.563), Disease (0.576)	4, 5
QLA10069.1	India	V(256)I, N(257)Q	2	FD III, FDIII (Cysteine rich domain)	Disease (0.770), Disease (0.563)	1, 2
QJIS3735.1	Greece	G(251)V, M(260)I	2	TD I, TD I	Disease (0.770), Disease (0.661)	5, 1
QLG98012.1	USA	A(103)S, W(131)L	2	TD I, TD I	Neutral (0.448), Disease (0.601)	1, 3
QLF98084.1	India	A(54)S, Q(57)H	2	TD I, TD I	Disease (0.613), Disease (0.637)	2, 3
QLH56099.1	Saudi Arabia	A(51)S, Q(57)H	2	TD I, TD I	Disease (0.600), Disease (0.637)	2, 3
QKV39324.1	USA	G(44)V, Q(57)H	2	FD II, TD I	Disease (0.628), Disease (0.637)	3, 3
QKU32982.1	USA	F(43)Y, Q(57)H	2	FD II, TD I	Disease (0.625), Disease (0.637)	3, 3
QLH64816.1	India	P(42)R, Q(57)H	2	FD II, TD I	Disease(0.615), Disease (0.637)	2, 3
QJAI7681.1	USA	P(42)R, Q(57)H	2	FD II, TD I	Disease(0.615), Disease (0.637)	2, 3
QJY40506.1	India	I(35)T, L(53)F	2	TD I	Disease (0.628), Disease (0.601)	3, 2
QLH57751.1	USA	D(27)H, Q(57)H	2	TD I	Neutral (0.139), Disease (0.637)	7, 3
QLC46314.1	USA	D(27)H, Q(57)H	2	TD I	Neutral (0.139), Disease (0.637)	7, 3
QKN19672.1	USA	P(25)S, T(175)I	2	TD I	Neutral (0.162), Disease (0.728)	7, 3
QLG75822.1	Australia	A(23)S, Q(57)H	2	TD I	Neutral (0.494), Disease (0.637)	7, 5
QLG97484.1	USA	D(22)Y, Q(57)H	2	TD I	Neutral(0.187), Disease (0.637)	0, 3
QLI50282.1	USA	G(18)S, Q(57)H	2	TD I	Neutral(0.055), Disease (0.637)	6, 3
QLA10165.1	India	G(18)V, Q(57)H	2	TD I	Neutral (0.036), Disease (0.637)	9, 3
QLI51614.1	USA	Q(57)H, V(197)L	2	TD I	Disease (0.637), Disease (0.509)	3, 0
QJD47299.1	USA	Q(57)H, S(165)I	2	TD I	Disease (0.637), Disease (0.605)	3, 2
QLF99991.1	USA	Q(57)H, T(170)S	2	TD I	Disease (0.637), Neutral(0.174)	3, 7
QJX70192.1	USA	Q(57)H, S(195)Y	2	TD I	Disease (0.637), Disease (0.653)	3, 3
QLE11150.1	Bangladesh	N(257)D, deletion(256)	2	TD I	Disease (0.590)	2
QJW28665.1	USA	Q(57)H, L(65)F, G(224)C	3	TD I	Disease (0.637), Neutral(0.233), Disease(0.693)	3, 5, 4
QKV26659.1	USA	Q(57)H, Q(185)H, Y(211)C	3	TD I	Disease (0.637), Disease(0.636), Disease(0.733)	3, 3, 5
QKG87159.1	USA	Q(57)H, A(99)V, V(237)A	3	TD I, FD III	Disease (0.637), Disease (0.602), Disease(0.583)	3, 2, 2
QKV42875.1	USA	V(88)A, S(171)L, G(251)V	3	FD III, FD VI (SGD motif)	Disease (0.636), Disease (0.602), Disease (0.770)	3, 2, 5
QKE44990.1	USA	L(94)P, V(97)A, F(120)L	3	FD III, FD III, FDIII	Disease(0.691), Neutral(0.157), Disease(0.641)	4, 7, 3
QLA47776.1	USA	Q(57)H, V(55), A(23)S	3	TD I, TD I	Disease (0.637), Disease(0.702)(3), Neutral (0.494)	3, 4, 0
QKV41592.1	USA	V(88)A, L(108)F, S(171)L, G(251)V	4	FD III, FDIII, FD VI (SGD motif)	Disease (0.636), Neutral(0.367), Disease(0.602), Disease (0.770)	3, 3, 2, 5

Table 9: ORF3a proteins possessed disease, neutral and mixed of neutral & disease type of predicted mutations

Disease		Disease		Disease		Neutral		Neutral & Disease	
Protein ID	Geo-location	Protein ID	Geo-location	Protein ID	Geo-location	Protein ID	Geo-location	Protein ID	Geo-location
QLG76542.1	Australia	QJD47849.1	Taiwan	QLC93357.1	USA	QJR88390.1	Australia	QLF98036.1	Bangladesh
QJR95110.1	Australia	QJD47873.1	Taiwan	QII57239.2	USA	QJR88822.1	Australia	QLH93441.1	Bangladesh
QLG75942.1	Australia	QKS66053.1	USA	QKU53854.1	USA	QKV38281.1	Australia	QLH93453.1	Bangladesh
QKV37633.1	Australia	QJC19648.1	USA	QKV07340.1	USA	QJR88306.1	Australia	QLF97952.1	India
QKV38005.1	Australia	QJV21807.1	USA	QLI50570.1	USA	QJR89110.1	Australia	QJX44383.1	India
QKV38209.1	Australia	QKG81932.1	USA	QLH59007.1	USA	QLG75930.1	Australia	QKO0487.1	India
QKV38257.1	Australia	QLI50414.1	USA	QKK14612.1	USA	QLG75678.1	Australia	QLI49698.1	India
QJR89362.1	Australia	QKS65621.1	USA	QKX46204.1	USA	QLF97844.1	Bangladesh	QLA50282.1	USA
QLG76026.1	Australia	QKU29039.1	USA	QLH01382.1	USA	QLH56279.1	Bangladesh	QLI46290.1	USA
QLG76386.1	Australia	QKN20812.1	USA	QJY78272.1	USA	QLH55768.1	Bangladesh	QKV40164.1	USA
QJR87730.1	Australia	QLF95245.1	USA	QKU52834.1	USA	QLH55720.1	Bangladesh	QLG97460.1	USA
QJR89278.1	Australia	QLI50222.1	USA	QKU31182.1	USA	QJW69308.1	Germany	QKE45885.1	USA
QJR89446.1	Australia	QLC94737.1	USA	QJX70592.1	USA	QIZ16548.1	Greece	QLA47500.1	USA
QKV38401.1	Australia	QKZ13838.1	USA	QLC92421.1	USA	QJS54191.1	Greece	QKG87087.1	USA
QJR91354.1	Australia	QJQ84173.1	USA	QKV07184.1	USA	QKE61733.1	India	QKG88539.1	USA
QLG75126.1	Baharain	QJY40110.1	USA	QLH57846.1	USA	QKY59990.1	India	QJH07211.1	USA
QLH93429.1	Bangladesh	QJD47551.1	USA	QJD47956.1	USA	QLH93202.1	India	QKV39588.1	USA
QLF97772.1	Bangladesh	QJQ47551.1	USA	QKV42204.1	USA	QJS39568.1	Netherlands	QKG91107.1	USA
QKO25747.1	Bangladesh	QKU30570.1	USA	QKG90867.1	USA	QJS39520.1	Netherlands	QIZ16438.1	USA
QKX47995.1	Bangladesh	QKG90399.1	USA	QLI51782.1	USA	QJS39616.1	Netherlands	QKG90495.1	USA
QKX49024.1	Bangladesh	QIZ13336.1	USA	QLC92097.1	USA	QLH01250.1	USA	QKR84274.1	USA
QLH55816.1	Bangladesh	QKV06224.1	USA	QIS61315.1	USA	QKV41616.1	USA	QKG87267.1	USA
QLF97736.1	Bangladesh	QLH58947.1	USA	QJF77147.1	USA	QLC92601.1	USA	QKG87195.1	USA
QKO25735.1	Bangladesh	QLH58085.1	USA	QJE38451.1	USA	QKU28463.1	USA	QLE91545.1	USA
QKK12852.1	Bangladesh	QKV38810.1	USA	QLF95773.1	USA	QKC05357.1	USA	QJD47203.1	USA
QLF98048.1	Bangladesh	QLC47346.1	USA	QIZ14498.1	USA	QKV06236.1	USA	QLG98012.1	USA
QLF80217.1	Brazil	QJQ39045.1	USA	QJD23730.1	USA	QKS67001.1	USA	QLH57751.1	USA
QKS67456.1	CHINA	QJU70306.1	USA	QLC94305.1	USA	QKG81824.1	USA	QLC46314.1	USA
QKE10935.1	Czech Republic	QLI51038.1	USA	QLF95737.1	USA	QKU53650.1	USA	QKN19672.1	USA
QKS66941.1	Egypt	QKG86518.1	USA	QKS66305.1	USA	QKW88844.1	USA	QKE44990.1	USA
QKV38894.1	Egypt	QKU28847.1	USA	QKS65597.1	USA	QKV07400.1	USA	QLA47776.1	USA
QJY78153.1	Egypt	QLH58037.1	USA	QJS54923.1	USA	QKW89480.1	USA	QKV41592.1	USA
QJT72507.1	France	QJX68859.1	USA	QJQ39081.1	USA	QLH01334.1	USA	QJW28665.1	USA
QJT72327.1	France	QJS57052.1	USA	QKV08048.1	USA	QLH00290.1	USA		
QJT72471.1	France	QIS61075.1	USA	QKE45933.1	USA				
QJT72387.1	France	QJW28449.1	USA	QLI51746.1	USA				
QJT72951.1	France	QLC91905.1	USA	QLG99773.1	USA				
QJS54155.1	Greece	QKG88935.1	USA	QLH00362.1	USA				
QJS53735.1	Greece	QKN20824.1	USA	QKU37646.1	USA				
QJS54023.1	Greece	QLA09656.1	USA	QKS65849.1	USA				
QLF98201.1	India	QKV39324.1	USA	QJC20500.1	USA				
QJX44407.1	India	QKU32982.1	USA	QKU32046.1	USA				
QJW00412.1	India	QJA17681.1	USA	QJU11458.1	USA				
QLA10069.1	India	QJD47419.1	USA	QKU31638.1	USA				
QLF98084.1	India	QLB39261.1	USA	QKU31746.1	USA				
QLH64816.1	India	QLC46986.1	USA	QLH01238.1	USA				
QJY40506.1	India	QLF95641.1	USA	QJQ39297.1	USA				
QLA10225.1	India	QJD23478.1	USA	QLB39321.1	USA				
QLF98261.1	India	QKG64052.1	USA	QLH01298.1	USA				
QLF78310.1	Poland	QLH58601.1	USA	QLG99737.1	USA				
QLH56255.1	Saudi Arabia	QKU53050.1	USA	QKS89844.1	USA				
QLH56231.1	Saudi Arabia	QJC20380.1	USA	QJD47539.1	USA				
QLH56099.1	Saudi Arabia	QKQ63773.1	USA	QIZ14498.1	USA				
QKU37034.1	Saudi Arabia	QKU32202.1	USA	QKV35400.1	USA				
QJD20838.1	Shri Lanka	QKV40716.1	USA	QKU37202.1	USA				
QHZ00380.1	South Korea	QKE45861.1	USA	QIS30116.1	USA				
QIU78768.1	Spain	QKV35688.1	USA	QKN20740.1	USA				
				QIU81286.1	USA				
				QKV26659.1	USA				
				QKG87159.1	USA				
				QKV42875.1	USA				

Almost 72% of the ORF3a proteins possess disease type of mutations whereas 14% (of which two mutations: 12%, three mutations: 1.5% and four mutations: 0.5%) and 14% of ORF3a proteins possess mixture type (i.e. both disease as well as neutral) and neutral types of mutations respectively (Fig.5).

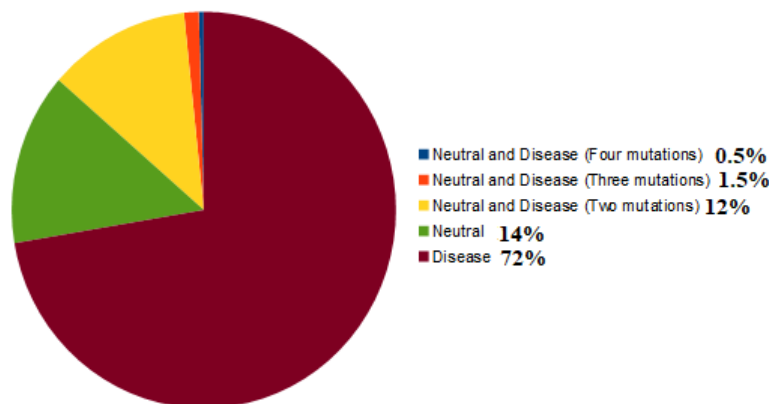


Figure 5: Percentage of disease, neutral and mix (neutral & disease) type of mutations over the ORF3a proteins

130 For each of the three types of mutations, we put the frequency and percentage of ORF3a proteins corresponding to each geo-locations as presented in the Table 10.

Table 10: Frequency and percentage of ORF3a proteins located at various countries, having three type of mutations

Disease			Neutral			Neutral & Disease		
<i>Geo-location</i>	<i>Frequency</i>	<i>Percentage</i>	<i>Geo-location</i>	<i>Frequency</i>	<i>Percentage</i>	<i>Geo-location</i>	<i>Frequency</i>	<i>Percentage</i>
USA	116	66.30%	USA	14	41.20%	USA	26	78.80%
AUSTRALIA	15	8.60%	AUSTRALIA	7	20.60%	INDIA	4	12.10%
BANGLADESH	10	5.70%	BANGLADESH	4	11.80%	BANGLADESH	3	9.10%
INDIA	9	5.10%	NETHERLANDS	3	8.80%			
FRANCE	5	2.90%	INDIA	3	8.80%			
SAUDI ARABIA	4	2.30%	GREECE	2	5.90%			
EGYPT	3	1.70%	GERMANY	1	2.90%			
GREECE	3	1.70%						
TAIWAN	2	1.10%						
BAHARAIN	1	0.60%						
SOUTH KOREA	1	0.60%						
CHINA	1	0.60%						
BRAZIL	1	0.60%						
CZECH REPUBLIC	1	0.60%						
SHRI LANKA	1	0.60%						
POLAND	1	0.60%						
SPAIN	1	0.60%						

In USA, all three type of mutations over the ORF3a proteins are found to be dominant in percentage. In a Fig.6, the world maps are marked as per occurrence of three types of mutations in ORF3a variants.

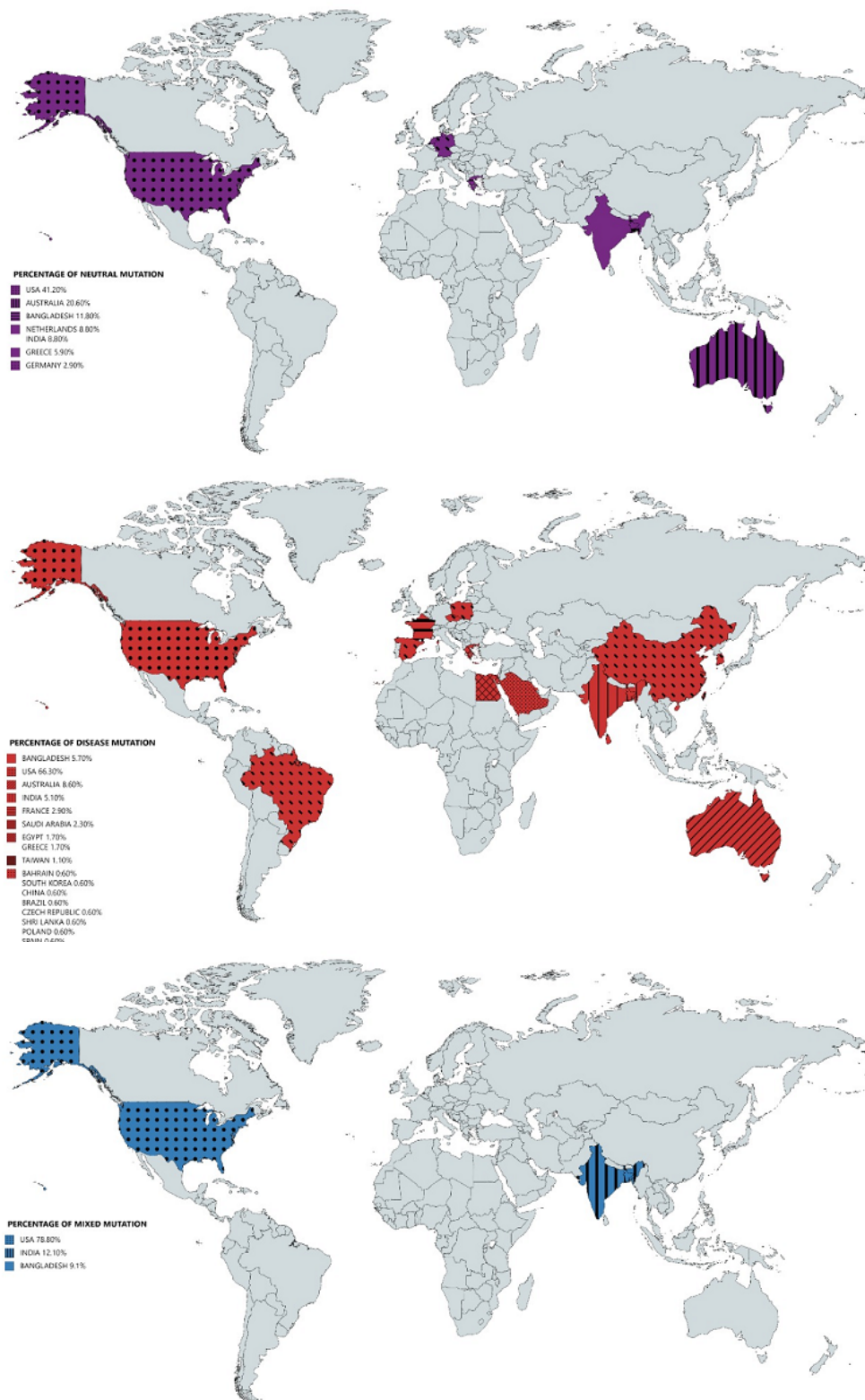


Figure 6: World maps of percentage of occurrence of neutral, disease and mixed type of mutations over the ORF3a proteins

Frequency of neutral mutation is 41.2% being the highest in the USA, according to prediction it shows that this mutation is neutral but still this mutation is supposed to be contributing to the weaker apoptotic activity of ORF3a and this weaker activity may be responsible for asymptomatic or relatively mildly symptomatic cases thus causing rapid transmission of the virus.

3.2. Possible Consecutive mutations over ORF3a proteins

Several ORF3a proteins (Tables 4-8) contain more than one mutations and maximally up to four mutations. It takes time for multiple mutations in a given ORF3a protein and relying on time-line and order occurrence of mutations several flow of consecutive mutations were derived. The predicted effects of these mutations on stability of the tertiary structure of the ORF3a proteins was determined in the flow of consecutive mutations (Table 11).

Table 11: ORF3a proteins with associated mutations and predicted effect in stability of the structures

Protein ID	Location	Mutation	Type of mutation	Effect on stability	* RI
QJR87730.1	Australia	Q(57)H	*P to P	Decrease	6
QKV38005.1	Australia	Q(57)H, K(75)R	P to P, P to P	Decrease, Increase	6, 3
QLG75822.1	Australia	Q(57)H, A(23)S	P to P, *NP to P	Decrease, Decrease	6, 8
QLG76542.1	Australia	Q(57)H, V(55)G	P to P, NP to NP	Decrease, Decrease	6,
QJR95110.1	Australia	Q(57)H, L(140)F	P to P, NP to NP	Decrease, Decrease	6, 9
QKU53050.1	USA	Q(57)H	P to P	Decrease	6
QKU30570.1	USA	Q(57)H, W(131)C	P to P, NP to P	Decrease, Decrease	6, 7
QIZ13838.1	USA	Q(57)H, L(95)F	P to P, NP to NP	Decrease, Decrease	6, 7
QKU289039.1	USA	Q(57)H, V(55)F	P to P, NP to NP	Decrease, Decrease	6, 9
QKU28847.1	USA	Q(57)H, M(260)I	P to P, NP to NP	Decrease, Decrease	6, 6
QKC88539.1	USA	Q(57)H, L(108)F	P to P, NP to NP	Decrease, Decrease	6, 7
QLI59282.1	USA	Q(57)H, G(18)S	P to P, NP to P	Decrease, Decrease	6, 8
QJU70306.1	USA	Q(57)H, G(224)C	P to P, NP to P	Decrease, Decrease	6, 3
QLA47776.1	USA	Q(57)H, V(55)F, A(23)S	P to P, NP to NP,	Decrease, Decrease, Decrease	6, 9, 8
QLH58085.1	USA	Q(57)H, Q(185)H	P to P, P to P	Decrease, Decrease	6, 3
QJW28665.1	USA	Q(57)H, G(224)C, L(65)F	P to P, NP to P, NP to NP	Decrease, Decrease, Decrease	6, 8, 7
QLA10225.1	Inida	Q(57)H	P to P	Decrease	6
QLF98201.1	Inida	Q(57)H, R(134)L	P to P, P to NP	Decrease, Decrease	6, 9
QLF98084.1	Inida	Q(57)H, A(54)S	P to P, NP to P	Decrease, Decrease	6, 8
QLH64816.1	Inida	Q(57)H, P(42)R	P to P, NP to P	Decrease, Decrease	6, 9
QLI49698.1	Inida	Q(57)H, T(271)I	P to P, P to NP	Decrease, Increase	6, 3
QLA10165.1	Inida	Q(57)H, G(18)V	P to P, NP to NP	Decrease, Decrease	6, 4
QLC46986.1	USA	Q38P	P to NP	Decreases	6
QKGS1932.1	USA	Q38P, W131S	P to NP, NP to P	Decreases, Decreases	6, 6
QKV07184.1	USA	G254R	NP to P	Decreases	7
QJC19648.1	USA	G254R, T9K	NP to P, P to P	Decreases, Decreases	7, 7
QKU53050.1	USA	Q57H	P to P	Decreases	6
QJ727471.1	FRANCE	Q57H, A99V	P to P, NP to NP	Decreases, Increases	6, 7
QKGS7159.1	USA	Q57H, A99V, V237A	P to P, NP to NP, NP to NP	Decreases, Increases, Decreases	6, 7, 9
QJ727507.1	FRANCE	Q57H, Y154C	P to P, P to NP	Decreases, Decreases	6, 5
QLG75678.1	AUSTRALIA	H78Y	P to NP	Increases	6
QJR88822.1	AUSTRALIA	H78Y, V13L	P to NP, NP to NP	Increases, Increases	6, 0
QLF98036.1	BANGLADESH	H78Y, Q38E	P to NP, P to P	Increases, Increases	6, 1
QJR89362.1	AUSTRALIA	G251V	NP to NP	Decreases	4
QKV38209.1	AUSTRALIA	G251V, W69L	NP to NP, NP to NP	Decreases, Decreases	4, 5
QJS54023.1	GREECE	G251V	NP to NP	Decreases	4
QJS53735.1	GREECE	G251V, M260I	NP to NP, NP to NP	Decreases, Decreases	4, 6
QLA09656.1	USA	G251V, V88A	NP to NP, NP to NP	Decreases, Decreases	4, 9
QKV42875.1	USA	G251V, V88A, S171L	NP to NP, NP to NP, P to NP	Decreases, Decreases, Increases	4, 9, 1
QKV41592.1	USA	G251V, V88A, S171L, L108F	NP to NP, NP to NP, P to NP, NP to NP	Decreases, Decreases, Increases, Decreases	4, 9, 1, 7

*Here P and NP stands for Polar, Non-Polar and RI : Reliability index

Flow of consecutive mutation-I: In the Australian region, it can be observed that the first mutation may have occurred in sequence QJR87730.1 with respect to the Wuhan sequence (YP_009724391.1) from Q to H at 57th position which is a disease type mutation and also this mutation is having the highest frequency which may indicate that it has an important role to play in infectivity part of the virus. As we move along the flow, six ORF3a sequences were considered based on the consecutive time scale of detection that was found to have 2nd mutation on the background of initial Q57H mutation with reference to Wuhan sequence (YP 009724391.1) (Fig.7).

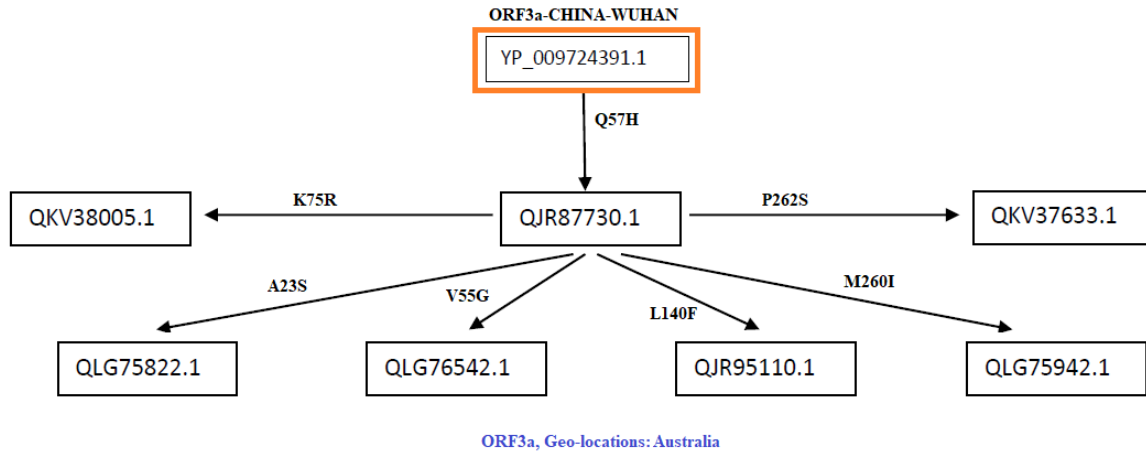


Figure 7: Flow of mutations in Australian ORF3a proteins

In this flow of mutation, six ORF3a proteins possess various mutations as follows:

- In QKV38005.1, there is a mutation K75R which was found to be a diseased type. We have to consider disease type mutation which may change the function of the protein.
- In QLG75822.1, there is a mutation A23S which was found to be a neutral type with no polarity change. So this is a synonymous mutation from the functionality perspective.
- In QLG76542.1, there is a mutation V55G which was found to be a diseased type, and and hydrophobicity changed to hydrophilicity. This indicates that there may be a functional importance of this mutation.
- In QJR95110.1, there is a mutation L140F which was found to be a diseased type with no polarity change. Since no polarity change is observed the type of amino acid remains same but the mutation effect becomes harmful for the host.
- In QLG75942.1, there is a mutation at M260I that was found to be a diseased type with no polarity change. This mutation may increase the virus virulence.
- In QKV37633.1, there is a mutation at P262S which was found to be a diseased type, and polarity changed from hydrophobic to hydrophilic. Consequently, it may account for change in structure of the protein.

Flow of consecutive mutation-II: The most frequent mutation Q57H occurred in the ORF3a protein QKV53050.1. In this network flow (Fig.8) there are other nine sequences which are considered based on the succeeding time scale that was found to have 2nd level mutations along with Q57H.

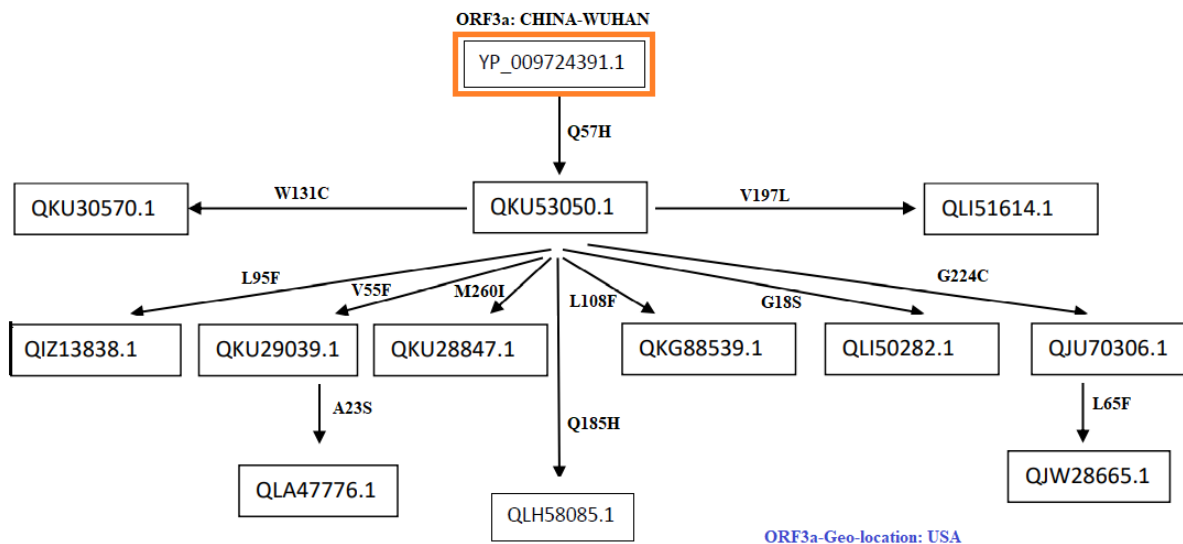


Figure 8: Flow of mutations in ORF3a proteins from the USA

- 165
- The ORF3a protein QKU30570.1 contains a mutation W131C which was found to be a diseased type and polarity changed from hydrophobic to hydrophilic. This mutation might affect the function of the ORF3a protein.
 - QIZ13838.1 possess a mutation L95F which was found to be a diseased type with no polarity change.
 - There is a mutation a V55F in QKU29039.1, which was found to be a diseased type with no polarity change. But the mutation may cause an increase in pathogenesis.
- 170
- In the protein QKU28847.1, a mutation M260I occurred which was found to be a diseased type with no polarity change and hence functional change of ORF3a can be expected.
 - In QLH58085.1, there is a mutation Q185H which was found to be a diseased type with no polarity change and so the structure of ORF3a protein may vary.
 - In QKG88539.1, there is a mutation at L108F which was found to be a neutral type with no polarity change. This mutation needs further investigation in order to confirm about its neutrality.
- 175
- In QLI50282.1, there is a mutation G18S which was found to be a neutral type, and polarity changed from hydrophobic to hydrophilic. Although this is a neutral mutation but the change in polarity may bear some significance in structural properties.
 - In QJU70306.1, there is a mutation at G224C which was found to be a diseased type polarity changed from hydrophobic to hydrophilic. This mutation may change the structure and functions of the protein.
- 180
- The ORF3a protein QLI51614.1 contains a mutation V197L which was found to be a diseased type with no polarity change.

In this network flow of mutations, it was also found sequences possessing 3rd level mutations which are described below:

- 185
- QLA47776.1: this sequence contains three mutations (Q57H, V55S, A23S), 3rd mutation is the neutral type, and polarity changed from hydrophobic to hydrophilic. Such mutations altogether may affect both structure and function of the protein.
 - QJW28665.1: this sequence contains three mutations (Q57H, G224C, L65F), 3rd mutation is the neutral type with no polarity change. The mutation L65F might not affect in virulence property of the SARS-CoV2.

190 **Flow of consecutive mutation-III:** In this case, network flow (Fig.9) of mutations is devised based on the ORF3a proteins of Indian origin. The sequence QLA10225.1 contains a mutation Q57H as usual. Further five ORF3a proteins are turned up in the network flow in the succeeding time scale of collection of samples. It was found that, all of them possess second mutation along with Q57H.

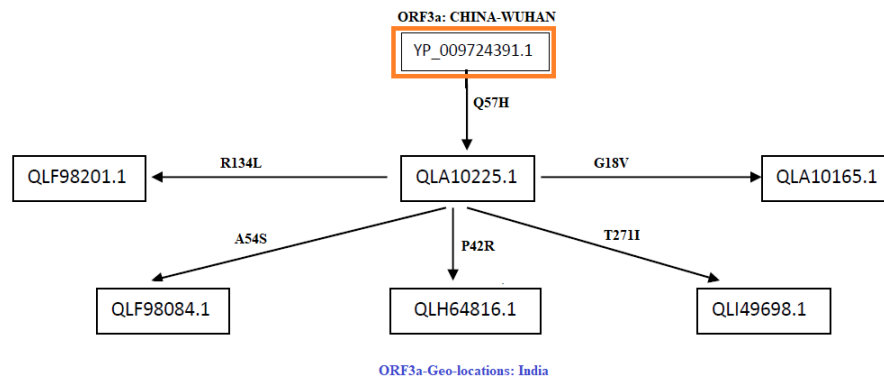


Figure 9: Flow of mutations in ORF3a proteins of Indian origin

- 195
- The mutation R134L in the ORF3a protein QLF98201.1, which was found to be a disease type and there was a polarity change from hydrophilic to hydrophobic. Here the change in mutations may lead to changes in tetramerization properties of the protein.
 - The protein QLF98084.1 possesses a mutation at A54S, which was found to be a disease type and the polarity changed from hydrophobic to hydrophilic and hence the structure of the protein is expected to be differed and accordingly the functions of the ORF3a protein would be affected.
 - QLH64816.1, there is a mutation at P42R which was found to be a disease type and there was a change in polarity from hydrophobic to hydrophilic and consequently the mutation may contribute to structural changes of the ORF3a protein.
 - The protein QLI49698.1 contains the mutation T271I which was found to be a neutral type and there was a change is polarity from hydrophilic to hydrophobic. Although the mutation is predicted to be neutral but the hydrophobicity is changed and hence alternation of functions of the proteins is anyway expected.
 - In ORF3a protein QLA10165.1, there is a mutation G18V which was found to be a neutral type of mutation and there is no change in polarity and consequently functions of the proteins would remain same.
- 200
- 205

Flow of consecutive mutation-IV: The sequence QLC46986.1 contains a mutation Q38P which is a disease mutation with the change in polarity from hydrophilic to hydrophobic which might cause a change in functions of the protein. The network flow of mutation id presented in the Fig.10.

210

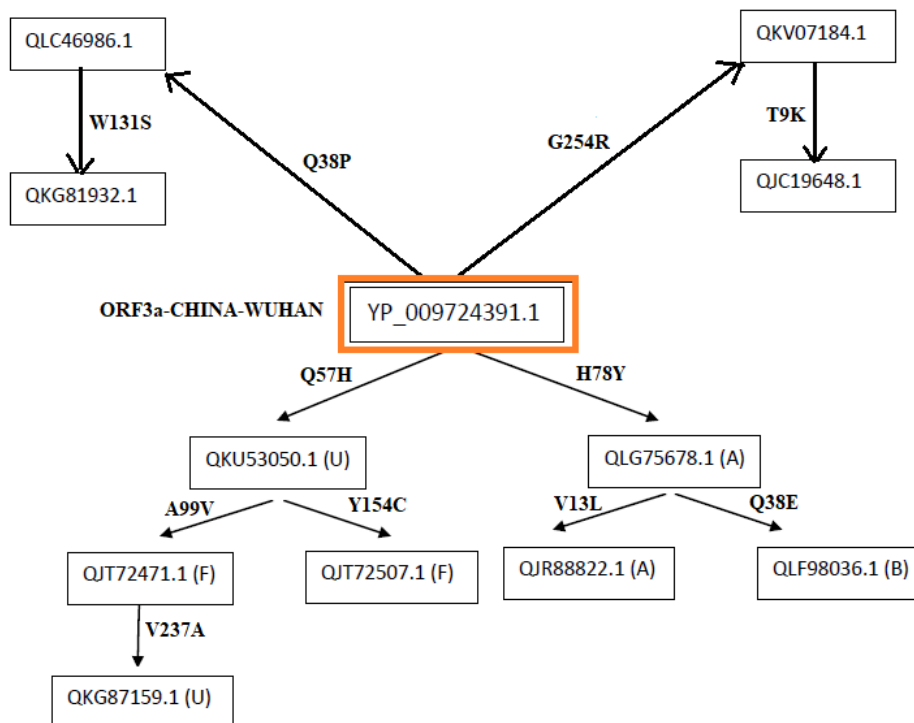


Figure 10: Network flow of mutations of ORF3a proteins considering from various geo-locations

A second level mutation along with Q38P occurred in QKG81932.1 sequence from W131S which is also a disease type mutation and polarity changed from hydrophobic to hydrophilic and so, it may change the structure of the protein. Also the ORF3a protein QKV07184.1 possesses G254R which changed the polarity from hydrophobic to hydrophilic and caused disease type mutation. On further analysis, the QJC19648.1 sequence was identified to have G254R along with T9K which is a disease mutation with no change in polarity. This is a mutation at the C-terminal region of protein so this mutation may effect the protein-protein interaction.

There is another sequence QKU53050.1(from USA) present in the work flow, which contains the usual mutation Q57H and a France based ORF3a sequence QJT72471.1 possessing a Q57H mutation along with A99V mutation which is a disease type mutation with no change in polarity. QJT72507.1 is another sequence of France origin, in which there is a mutation at Y154C along with Q57H mutation. Also in the QKG87159.1 sequence, another mutation apart from Q57H and A99V at position V237A which is a disease type with no change in polarity.

Another possible traffic of mutation was observed in which an Australian sequence QLG75678.1 had a mutation at 78th position from H to Y, a neutral mutation with no change in polarity which may be a virulence promoting factor. Another Australian sequence QJR88822.1 was identified in which H78Y mutation was observed with V13L which is a disease mutation with no change in polarity. So here we observed that along with a neutral mutation a disease mutation has occurred and it can be assumed that virus first evolved in terms of virulence then enhanced its functional activity. Although there is no change in polarity but it may affect the chemical properties. The sequence QLF98036.1 was another sequence from Bangladesh found to have H78Y mutation in addition to Q38E which is a disease mutation with no change in polarity. here also a disease mutation is observed along with neutral mutation again signifying the evolutionary importance of these mutations.

Flow of consecutive mutation-V: The network flow of mutations (Fig.11) with reference sequence of Wuhan's (ID

YP_009724391.1) is formed.

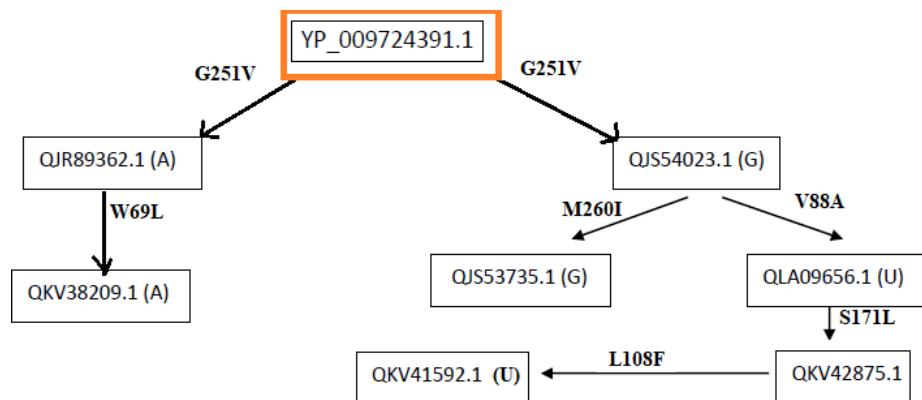


Figure 11: Network flow of mutations of ORF3a proteins considering from various geo-locations

Note: A: Australia, B: Bangladesh, F: France, G: Greece and U: USA.

The ORF3a protein QJR89362.1 possess a mutation G251V. It was found to be a disease type mutation and here no change in polarity is observed but it may have some significance as it is a disease causing mutation. From this originates another sequence in the flow whose sample collection date is ensuing to the previous one. This sequence (ID QKV38209.1) bears a mutation in W69L which is a disease mutation without any change of polarity that is both W and L are neutral. As this sequence has both the disease mutations, it indicates their functional importance.

In the second case, when the sequence (ID QJS54023.1) of geo-location Greece, is compared with the Wuhan sequence it bore the same mutation G251V. From here it is further divided into bi-flow according to geo-locations and all of them have the G251V mutation along with certain new:

1. The left one bears a sequence (ID QJS53735.1) of geo-location Greece which has a mutation M260I which is a disease type of mutation and has no change in polarity. Here, both the mutations are in the cytosolic domain indicating that these mutations are somehow important for the virus.
2. The right one is for the geo-location USA, which starts with the sequence (ID QLA09656.1) which has a mutation V88A. It is a disease type mutation with no change in polarity. So, it may be advantageous for virus in terms of functionality. Following there is another sequence (ID QKV42875.1) with respect to the time scale, bearing a mutation at S171L. This is a disease type mutation and there is a change in polarity from hydrophilic to hydrophobic. Since the polarity is changing which indicates that there is some effect on ionic and electrostatic interactions that may cause structural changes. Lastly, the sequence QKV41592.1 which bears a mutation at L108F which is a neutral mutation which has no change in polarity. This sequence has all disease mutations although no change in polarity is observed except for one mutation, so it signifies the order of occurrence of mutations allowing the virus to acquire new characteristics important for its survival.

In this study of mutation among many, we recognised five important mutations in the ORF3a proteins. While W131C, T151I, R134L and D155Y forms a network of hydrophobic, polar and electrostatic interactions which are important for the tetramerization process of ORF3a (the functional unit of ORF3a), F230 insertion is responsible for dimerization of ORF3a. We could see that all of the mutations have an effect of decrease in the stability apart from T151I which increases the stability of the protein. To get a better insight, we analysed for these mutations from a structural point of view:

Case-I: We collected the available structure of ORF3a (Protein ID: 6XDC) from Protein Data Bank(PDB), (leftmost figure shown in colour grey) in Fig.12

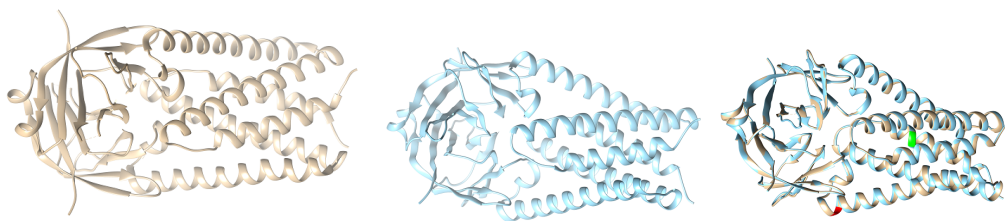


Figure 12: Structures of ORF3a (Reference coloured as grey in left), Structure of mutated ORF3a (coloured with blue in the middle) and Overlaid ORF3a (rightmost image)

260 Then we took the mutated sequence which contains the mutation W131C and performed homology modelling with the help of a web server called Swiss-model and built the corresponding structure of W131C (middle picture shown in blue) and finally we superimposed the structure of Wuhan (reference structure) with that of the modelled (right most picture) and checked for the corresponding differences with respect to structural change; labelling the mutated portions with colour green(Q57H) and red(W131C).

265 **Case-II:** In this case, we consider the mutated sequence which possesses the mutation T151Y and performed homology modelling and built the corresponding structure of T151Y (middle picture shown in blue) as shown in Fig.13.

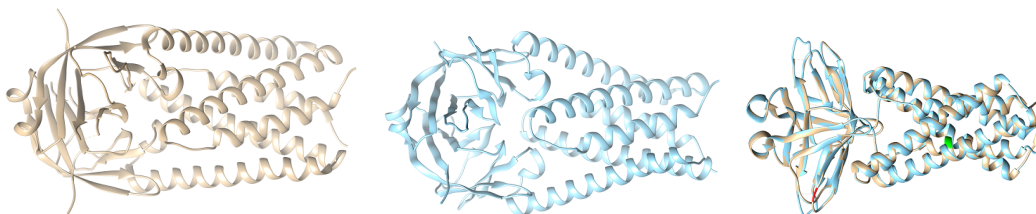


Figure 13: Structures of ORF3a (Reference coloured as grey in left), Structure of mutated ORF3a (coloured with blue in the middle) and Overlaid ORF3a (rightmost image)

Finally we overlaid the structure of Wuhan (reference structure) with that of the modelled (right most picture) and checked for the corresponding differences with respect to structural change; labelling the mutated portions with colour green(Q57H) and red(T151Y).

270 **Case-III:** With the available structure of ORF3a (Protein ID: 6XDC) from Protein Data Bank(PDB), (leftmost picture shown in colour grey) we took the mutated sequence of R134L and performed homology modelling and built the corresponding structure of R134L (middle picture shown in blue in Fig.14)

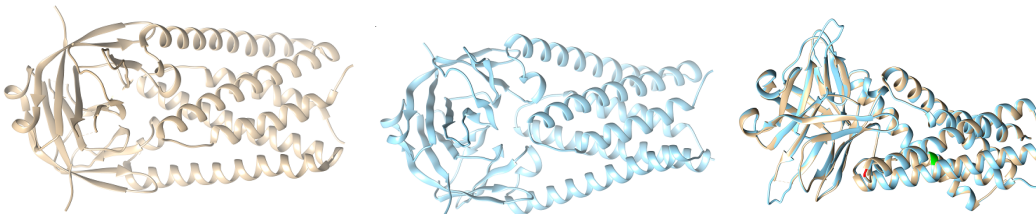


Figure 14: Structures of ORF3a (Reference coloured as grey in left), Structure of mutated ORF3a (coloured with blue in the middle) and Overlaid ORF3a (rightmost image)

Then we overlaid the structure of Wuhan (reference structure) with that of the modelled (right most picture) and checked for the corresponding differences with respect to structural change; labelling the mutated portions with colour green(Q57H) and red(R134L).
275

Case-IV: With the available structure of ORF3a (Protein ID: 6XDC) (leftmost picture shown in colour grey) and then we took the mutated sequence ORF3a considering the mutation D155Y and performed homology modelling and obtained the corresponding structure of D155Y (middle picture shown in blue in the Fig.15).

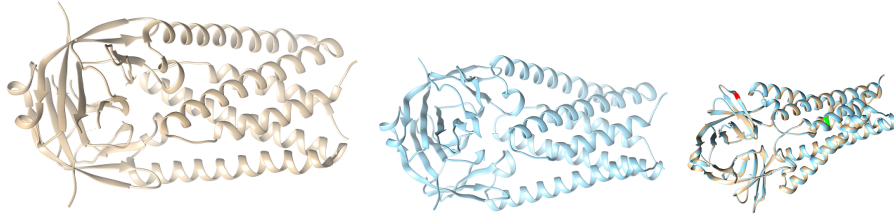


Figure 15: Structures of ORF3a (Reference coloured as grey in left), Structure of mutated ORF3a (coloured with blue in the middle) and Overlaid ORF3a (rightmost image)

We then overlaid the structure of Wuhan (reference structure) with that of the modelled (right most picture) and checked for the corresponding differences with respect to structural change; labelling the mutated portions with colour green(D155Y) and red(D155Y).
280

Case-V: Using the structure of the ORF3a (Protein ID: 6XDC) (leftmost picture shown in colour grey in Fig.16) by homology modelling the structure of the ORF3a protein which contains the insertion mutation F230 (middle picture shown in blue), is constructed.

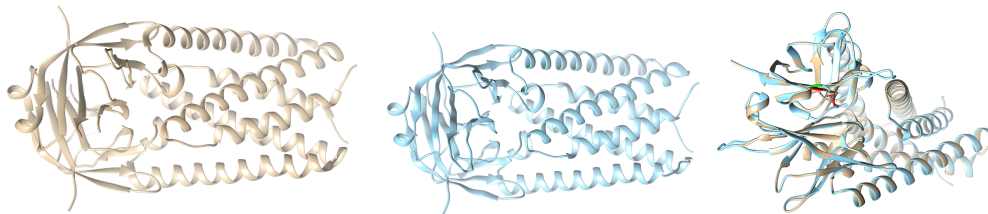


Figure 16: Structures of ORF3a (Reference coloured as grey in left), Structure of mutated ORF3a (coloured with blue in the middle) and Overlaid ORF3a (rightmost image)

Then we overlaid the structure of ORF3a based in Wuhan (reference structure) with that of the modelled (right most picture) and checked for the corresponding differences with respect to structural change; labelling the mutated portions with colour green(difference in structure) and red(inserted amino acid).
285

We did the above study and no significant change in protein structure was observed, we need a better soft-ware to find the difference between Wuhan sequence and mutated sequences.

290 3.3. Phylogeny and Clustering

We attempted to cluster each of the 296 ORF3a proteins into twenty disjoint clusters based on the probability distribution of amino acids using K-means clustering technique (Table 12). Note that, the number of clusters (twenty) is chosen optimally by heuristic method in such a manner that the clusters are separated from each other significantly. The frequency probability of each amino acids across all the 296 ORF3a proteins is available as a supplementary file-I. The three truncated ORF3a proteins (detected in Indian patients) are clustered in the cluster 11 as shown in Table 13.
295

Table 12: ORF3a proteins and corresponding cluster number based on amino acid distributions

Protein ID	Cluster No	Protein ID	Cluster No	Protein ID	Cluster No	Protein ID	Cluster No	Protein ID	Cluster No
QLL46290.1 USA	1	QLH93441.1 Bangladesh	5	QKS66305.1 USA	6	QKX49024.1 Bangladesh	15		
QKGS1932.1 USA	1	QLF97772.1 Bangladesh	5	QKS65907.1 USA	6	QKV06236.1 USA	15		
QKN20812.1 USA	1	QLF97952.1 India	5	QKS88935.1 USA	6	QLX99737.1 USA	15		
QJD47551.1 USA	1	QKS66941.1 Egypt	5	QKS66041.1 USA	6	QIX45308.1 Poland	15		
QJD25758.1 USA	1	QKQ664052.1 USA	5	QJF54254.1 USA	6	QKS66053.1 USA	16		
QKU30570.1 USA	1	QKQ25747.1 Bangladesh	5	QKE61733.1 India	7	QJV21807.1 USA	16		
QKG990399.1 USA	1	QKX47995.1 Bangladesh	5	QKV41616.1 USA	7	QLF95641.1 USA	16		
QLH5894.1 USA	1	QLG75930.1 Bangladesh	5	QJR88306.1 Australia	7	QKE45885.1 USA	16		
QLH5894.1 USA	1	QKV38257.1 Australia	5	QJR89110.1 Australia	7	QKS65621.1 USA	16		
QKV26659.1 USA	1	QKV41592.1 USA	5	QJT72387.1 France	7	QKU29039.1 USA	16		
QLH58085.1 USA	1	QKV42875.1 USA	5	QJD47203.1 USA	7	QLG76542.1 Australia	16		
QLC47346.1 USA	1	QLA09656.1 USA	5	QKQ63773.1 USA	7	QJW28665.1 USA	16		
QKG91107.1 USA	1	QLG97055.1 Italy	5	QKU32202.1 USA	7	QLA47500.1 USA	16		
QJL72507.1 France	1	QKV40716.1 USA	5	QIZ16548.1 Greece	7	QJX44383.1 India	16		
QLH49698.1 USA	1	QKE45861.1 USA	5	QKU53854.1 USA	7	QLF9245.1 USA	16		
QJY75153.1 Egypt	1	QJD47873.1 Taiwan	5	QKU31806.1 USA	7	QLC94737.1 USA	16		
QKV08048.1 USA	1	QKV35688.1 USA	5	QJX45032.1 USA	7	QKQ87087.1 USA	16		
QKE45933.1 USA	1	QLC93357.1 USA	5	QJX45032.1 USA	7	QIZ13838.1 USA	16		
QKE99773.1 USA	1	QLH59007.1 USA	5	QJR91282.1 Australia	7	QJQ84173.1 USA	16		
QKV38894.1 USA	1	QJH7239.2 USA	5	QJR87598.1 Australia	7	QKG88539.1 USA	16		
QJX44407.1 India	1	QLF98261.1 India	5	QJH38451.1 USA	7	QJY40110.1 USA	16		
QKC290500.1 USA	1	QLF80217.1 Brazil	5	QJF39741.1 USA	7	QJX47849.1 Taiwan	16		
QKGS7267.1 USA	1	QLL50570.1 USA	5	QLF78310.1 Poland	7	QJR95110.1 Australia	16		
QJ557052.1 USA	1	QLH5907.1 USA	5	QJH9282.1 USA	7	QIZ13336.1 USA	16		
QLH93453.1 Bangladesh	1	QLEH5816.1 Bangladesh	5	QJH47539.1 USA	7	QKU53050.1 USA	16		
QJUI1458.1 USA	1	QKE10935.1 Czech Republic	5	QJZ14498.1 USA	7	QJI07211.1 USA	16		
QJ61075.1 USA	1	QLC92601.1 USA	5	QJ553831.1 Greece	7	QKV38810.1 USA	16		
QJR87730.1 Australia	1	QKK14612.1 USA	5	QJ554023.1 Greece	7	QJO39045.1 USA	16		
QJW28449.1 USA	1	QKU28463.1 USA	5	QLF98048.1 Bangladesh	7	QJT72327.1 France	16		
QLH56231.1 Saudi Arabia	1	QLF97844.1 Bangladesh	5	QJH48484.1 USA	7	QJX290380.1 USA	16		
QKS91905.1 USA	1	QKX46204.1 USA	5	QJY40506.1 India	7	QLG75942.1 Australia	16		
QKGS7159.1 USA	1	QJH84790.1 USA	5	QJH4498.1 USA	7	QKU28847.1 USA	16		
QKT2471.1 France	1	QJY78272.1 USA	5	QJZ14498.1 USA	7	QLI51746.1 USA	16		
QKU31638.1 USA	1	QJY52834.1 USA	5	QJH8822.1 Australia	8	QLH00362.1 USA	16		
QLH01238.1 USA	1	QLH1150.1 Bangladesh	5	QJR98036.1 Bangladesh	8	QKW31746.1 USA	16		
QKN20824.1 USA	1	QLH55840.1 Bangladesh	5	QJS39616.1 Netherlands	8	QJW00412.1 India	16		
QLF98084.1 India	1	QKU31182.1 USA	5	QJS61315.1 USA	8	QJR89278.1 Australia	16		
QLH56099.1 Saudi Arabia	1	QJX70592.1 USA	5	QJF77147.1 USA	8	QJR89446.1 Australia	16		
QKV39324.1 USA	1	QJC92421.1 USA	5	QLF95773.1 USA	8	QLH57751.1 USA	16		
QJ8300116.1 France	1	YJ029724391.1 China	5	QLG75678.1 Australia	9	QLA4776.1 USA	16		
QLC46314.1 USA	1	QJC05357.1 USA	5	QJ54923.1 USA	9	QLG97460.1 USA	16		
QLG75822.1 Australia	1	QJD20838.1 Sri Lanka	5	QJH47299.1 USA	9	QLI50414.1 USA	17		
QLU50282.1 USA	1	QKV36900.1 USA	5	QJQ38625.1 USA	9	QLI50222.1 USA	17		
QLA10165.1 India	2	QKV07184.1 USA	5	QKG90147.1 USA	9	QLF98201.1 India	17		
QKG90495.1 USA	2	QKM76547.1 Germany	5	QKV42947.1 USA	10	QKV06224.1 USA	17		
QLH58037.1 USA	2	QJF75396.1 USA	5	QKQ00487.1 —truncated India	11	QLH58601.1 USA	17		
QKR84274.1 Egypt	2	QKV7929.1 Bangladesh	5	QLA10225.1 —truncated India	11	QLG56255.1 Saudi Arabia	17		
QLC91545.1 Greece	2	QLG76386.1 Australia	5	QLA10069.1 —truncated India	11	QLG75126.1 Bahrain	17		
QLC92097.1 USA	2	QKY77929.1 USA	5	QKX89480.1 USA	12	QKG86518.1 USA	17		
QKU32982.1 USA	2	QJG99677.1 USA	5	QJH23478.1 USA	12	QJX68859.1 USA	17		
QKGS1824.1 USA	2	QKU31266.1 USA	5	QKJ32046.1 USA	12	QKU37646.1 USA	17		
QLK53650.1 USA	2	QLI51782.1 USA	5	QLH01382.1 USA	12	QLI51614.1 USA	17		
QLG97484.1 USA	2	QLC91617.1 USA	5	QKV06920.1 USA	12	QLJ399297.1 USA	17		
QLH565768.1 Bangladesh	3	QLC98012.1 USA	5	QLH01298.1 USA	12	QJQ339321.1 USA	17		
QJX70192.1 USA	3	QLC94473.1 USA	5	QJF54123.1 USA	12	QKX8442.1 Egypt	17		
QKS65777.1 USA	3	QLC900578.1 USA	5	QKE44990.1 USA	12	QJR91354.1 Australia	17		
QKV35400.1 USA	4	QKK12852.1 Bangladesh	5	QKS89844.1 USA	12	QKS54383.1 Greece	18		
QKV40440.1 USA	4	QKE45765.1 USA	5	QKX38209.1 Australia	13	QJ554383.1 Greece	18		
QJD47419.1 USA	5	QLH00026.1 USA	5	QJZ000380.1 South Korea	13	QKV38401.1 Australia	18		
QJL19648.1 USA	5	QJH52870.1 India	5	QJ53735.1 Greece	13	QKU32934.1 USA	18		
QJR88390.1 Australia	5	QJW59990.1 India	5	QLH57846.1 USA	13	QKS66737.1 USA	19		
QLH01250.1 USA	5	QJW52870.1 USA	5	QJH98362.1 Australia	13	QKV40164.1 USA	19		
QLH01334.1 USA	5	QJW52870.1 USA	5	QJS54191.1 Greece	13	QKV39588.1 USA	20		
QLH01334.1 USA	5	QJC92553.1 USA	5	QJH47956.1 USA	13	QLI51038.1 USA	20		
QLW69308.1 Germany	5	QJR87574.1 Australia	5	QLG76026.1 Australia	13	QJQ37633.1 Australia	20		
QKV38281.1 Australia	5	QKU31818.1 USA	5	QLF97736.1 Bangladesh	13	QKX39081.1 USA	20		
QLH00290.1 USA	5	QJR86050.1 Australia	5	QLH01250.1 USA	13	QKV38005.1 Australia	20		
QKS67456.1 China	5	QJC94305.1 USA	5	QJH78768.1 Spain	13	QKV42204.1 USA	20		
QJ839568.1 Netherlands	5	QKN19672.1 USA	5	QJH78768.1 Spain	13	QLH64816.1 India	20		
QLC46986.1 USA	5	QKS90192.1 USA	5	QKS67001.1 USA	13	QJA17681.1 USA	20		
		QJH1286.1 USA	5	QLH55720.1 Bangladesh	13	QLF95737.1 USA	20		
		QKV07400.1 USA	5	QJR84550.1 USA	14	QKN20740.1 USA	20		
				QLH93429.1 Bangladesh	15	QKW88844.1 USA	20		

Table 13: Clusters and its frequencies

Cluster	Frequency	Cluster	Frequency
1	47	11	3
2	11	12	10
3	3	13	13
4	1	14	2
5	86	15	5
6	5	16	36
7	28	17	16
8	7	18	3
9	4	19	2
10	1	20	13

The largest cluster 5 contains 53 ORF3a proteins of the USA patients including other 33 from various geo-locations as shown in Table 12. It is found that the ORF3a variants of the USA belong to each of the clusters except the cluster 11 which contained only three truncated proteins belong. This observation confirms the diversity of ORF3a isolates from the USA. It has been seen that the clusters 4, 6, 9 and 10 contains only the ORF3a proteins which are isolated from USA patients. It has been observed that the ORF3a proteins belonging to the clusters 4 and 10 do not possess any mutations and clearly these two ORF3a sequence contain ambiguous amino acids otherwise they would not have appeared as distinct variants.

Based on the hierarchical clustering method, single linkage dendrogram was obtained using the distance matrix of the clusters formed by the K-means clustering method over the 296 ORF3a proteins. This dendrogram (Fig.17) depicts the nearness of the clusters which are formed.

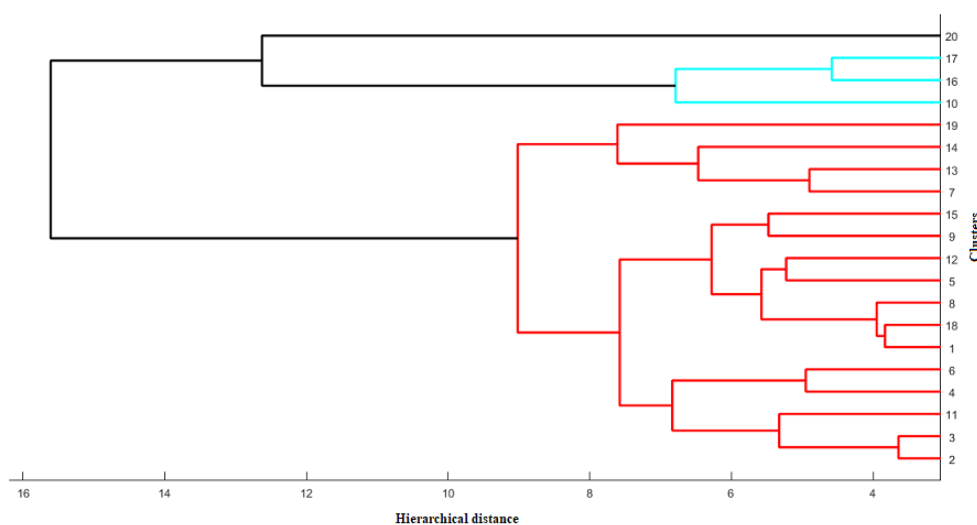


Figure 17: Dendrogram of the twenty clusters

The most nearest pair of clusters are (2, 3), (4, 6), (1, 18), (5, 12), (9, 15), (7, 13) and (16, 17) as observed from the dendrogram (Fig.17).

3.4. Variability of ORF3a Isolates

The variations among the ORF3a proteins based on the disorderly character of the amino acids over the proteins were determined using Shannon entropy (SE). For each sequence, SE is determined according to the formula stated in the method 2.2.3 and shown in Table 14.

Table 14: Shannon entropy of amino acid conservations of the 296 ORF3a distinct variants across the world

Protein ID	Geo-location	SE	Protein ID	Geo-location	SE	Protein ID	Geo-location	SE
QJR88390.1	Australia	0.957	QJD20838.1	Sri Lanka	0.959	QJD47203.1	USA	0.957
QJR88822.1	Australia	0.956	QJD47849.1	Taiwan	0.955	QKQ63773.1	USA	0.958
QKV38281.1	Australia	0.957	QJD47873.1	Taiwan	0.957	QKU32202.1	USA	0.958
QJR88306.1	Australia	0.958	QKS66053.1	USA	0.958	QKV40716.1	USA	0.958
QJR89110.1	Australia	0.958	QJD47419.1	USA	0.958	QKE45861.1	USA	0.955
QLG76542.1	Australia	0.958	QJC19648.1	USA	0.959	QKV35688.1	USA	0.957
QJR95110.1	Australia	0.958	QKW89480.1	USA	0.958	QLC93357.1	USA	0.955
QLG75942.1	Australia	0.955	QLH01250.1	USA	0.957	QLI57239.2	USA	0.958
QKV37633.1	Australia	0.957	QLH01334.1	USA	0.957	QKU53854.1	USA	0.958
QJR87730.1	Australia	0.957	QLB39261.1	USA	0.958	QKU31806.1	USA	0.958
QJR89278.1	Australia	0.959	QLI46290.1	USA	0.958	QKS66041.1	USA	0.960
QJR89446.1	Australia	0.958	QKV40164.1	USA	0.956	QKV07340.1	USA	0.958
QKV38005.1	Australia	0.958	QLG97460.1	USA	0.957	QLI50570.1	USA	0.958
QKV38209.1	Australia	0.955	QLH00290.1	USA	0.957	QLH59007.1	USA	0.958
QLG75930.1	Australia	0.958	QJV21807.1	USA	0.958	QLC92601.1	USA	0.958
QKV38257.1	Australia	0.958	QKG81932.1	USA	0.955	QKK14612.1	USA	0.957
QJR89362.1	Australia	0.958	QLC46986.1	USA	0.957	QKU28463.1	USA	0.958
QLG76026.1	Australia	0.957	QLI50414.1	USA	0.957	QKX46204.1	USA	0.958
QLG76386.1	Australia	0.957	QKV41616.1	USA	0.958	QJR84790.1	USA	0.958
QKV38401.1	Australia	0.960	QLF95641.1	USA	0.958	QLH01382.1	USA	0.958
QJR87598.1	Australia	0.958	QJD23478.1	USA	0.958	QJY78272.1	USA	0.956
QLG75678.1	Australia	0.957	QKE45885.1	USA	0.958	QKU52834.1	USA	0.958
QJR91282.1	Australia	0.959	QKS65621.1	USA	0.958	QKU31182.1	USA	0.956
QJR87574.1	Australia	0.957	QKU29039.1	USA	0.958	QJX70592.1	USA	0.959
QJR86050.1	Australia	0.957	QKG64052.1	USA	0.958	QLC92421.1	USA	0.956
QJR91354.1	Australia	0.958	QJW28665.1	USA	0.959	QKC05357.1	USA	0.958
QLG75822.1	Australia	0.957	QLA47500.1	USA	0.958	QJX45032.1	USA	0.958
QLG75126.1	Bahrain	0.957	QKN20812.1	USA	0.957	QKV36900.1	USA	0.958
QLH93429.1	Bangladesh	0.957	QLF95245.1	USA	0.958	BKVL07184.1	USA	0.958
QLF98036.1	Bangladesh	0.956	QLI50222.1	USA	0.957	QLH57846.1	USA	0.957
QLH93441.1	Bangladesh	0.956	QLC94737.1	USA	0.958	QLH01502.1	USA	0.957
QLF97772.1	Bangladesh	0.958	QKG87087.1	USA	0.958	QKV06236.1	USA	0.958
QLH93453.1	Bangladesh	0.957	QIZ13838.1	USA	0.958	QJF75396.1	USA	0.957
QKQ25747.1	Bangladesh	0.955	QJQ84173.1	USA	0.958	QKV06920.1	USA	0.957
QKX47955.1	Bangladesh	0.957	QKG88539.1	USA	0.958	QJD47956.1	USA	0.957
QKX49024.1	Bangladesh	0.957	QJY40110.1	USA	0.958	QKV42204.1	USA	0.958
QLH55816.1	Bangladesh	0.958	QJD47551.1	USA	0.958	QJQ38625.1	USA	0.959
QLF97844.1	Bangladesh	0.959	QJL25758.1	USA	0.957	QKG90867.1	USA	0.958
QLE11150.1	Bangladesh	0.957	QKU30570.1	USA	0.957	QKY77929.1	USA	0.958
QLH55840.1	Bangladesh	0.958	QKG90399.1	USA	0.957	QKV40440.1	USA	0.954
QLH56279.1	Bangladesh	0.958	QIZ13336.1	USA	0.958	QLH01298.1	USA	0.958
QLF97736.1	Bangladesh	0.957	QKS66305.1	USA	0.959	QLG99737.1	USA	0.958
QKQ25735.1	Bangladesh	0.957	QKS65597.1	USA	0.960	QLG99677.1	USA	0.957
QKK12852.1	Bangladesh	0.958	QKV06224.1	USA	0.957	QKU31266.1	USA	0.957
QLF98048.1	Bangladesh	0.957	QLH58601.1	USA	0.957	QLI51782.1	USA	0.957
QLH55768.1	Bangladesh	0.957	QLH58947.1	USA	0.958	QLC92097.1	USA	0.957
QLH55720.1	Bangladesh	0.957	QKU53050.1	USA	0.958	QIS61315.1	USA	0.956
QLF80217.1	Brazil	0.957	QJH07211.1	USA	0.958	QJF77147.1	USA	0.956
QKS67456.1	China	0.958	QKV26659.1	USA	0.958	QKG90147.1	USA	0.945
YP_009724391.1	China	0.958	QLH58085.1	USA	0.957	QJE38451.1	USA	0.956
QKE10935.1	Czech Republic	0.957	QKV39588.1	USA	0.957	QI54254.1	USA	0.943
QKS66941.1	Egypt	0.958	QKV38810.1	USA	0.958	QKS67001.1	USA	0.957
QJY78153.1	Egypt	0.957	QJS54923.1	USA	0.959	QLC91617.1	USA	0.958
QKR84274.1	Egypt	0.957	QLC47346.1	USA	0.957	QJ154123.1	USA	0.961
QKR84421.1	Egypt	0.957	QKG91107.1	USA	0.958	QJQ39741.1	USA	0.958
QJT72507.1	France	0.959	QJQ39045.1	USA	0.958	QJR84550.1	USA	0.960
QJT72327.1	France	0.958	QJU70306.1	USA	0.958	QLG98012.1	USA	0.955
QJT72471.1	France	0.957	QIZ16438.1	USA	0.957	QKE44990.1	USA	0.958
QJT72387.1	France	0.958	QLI51038.1	USA	0.956	QKU52934.1	USA	0.959
QJT72951.1	France	0.957	QK66518.1	USA	0.956	QLF95773.1	USA	0.957
QJW69308.1	Germany	0.956	QJC20380.1	USA	0.958	QLC94473.1	USA	0.958
QKM76547.1	Germany	0.957	QKU28847.1	USA	0.958	QLH00578.1	USA	0.958
QKM76907.1	Germany	0.958	QJQ39081.1	USA	0.957	QLC93129.1	USA	0.958
QJS54155.1	Greece	0.957	QKG90495.1	USA	0.958	QKS66737.1	USA	0.968
QIZ16548.1	Greece	0.958	QLH58037.1	USA	0.957	QKS65777.1	USA	0.965
QJS53735.1	Greece	0.955	QJX68859.1	USA	0.958	QKS8044.1	USA	0.956
QJS54383.1	Greece	0.958	QKV08048.1	USA	0.957	QJD47539.1	USA	0.957
QJS54191.1	Greece	0.956	QKE45933.1	USA	0.957	QIZ14498.1	USA	0.957
QJS53831.1	Greece	0.955	QLI51746.1	USA	0.958	QKE45765.1	USA	0.957
QJS54023.1	Greece	0.954	QLG99773.1	USA	0.957	QLH00026.1	USA	0.957
QKQ00487.1 truncated	India	0.957	QLH00362.1	USA	0.958	QJD23730.1	USA	0.958
QLA10225.1 truncated	India	0.957	QKV38894.1	USA	0.957	QKU52870.1	USA	0.957
QLA10069.1 truncated	India	0.957	QKV37646.1	USA	0.958	QLC92553.1	USA	0.957
QKE61733.1	India	0.958	QKS665849.1	USA	0.959	QKU31818.1	USA	0.957
QLF97952.1	India	0.957	QLI51614.1	USA	0.957	QKV35400.1	USA	0.954
QJX44383.1	India	0.958	QJD47299.1	USA	0.951	QKV42947.1	USA	0.952
QLF98201.1	India	0.955	QJC20500.1	USA	0.954	QKU37202.1	USA	0.956
QLI449698.1	India	0.958	QKGS7267.1	USA	0.957	QKV39324.1	USA	0.957
QJX44407.1	India	0.957	QKU32046.1	USA	0.956	QIS30116.1	USA	0.957
QJW00412.1	India	0.959	QJS57052.1	USA	0.958	QKU32982.1	USA	0.957
QLF98261.1	India	0.958	QKGS7195.1	USA	0.957	QJA17681.1	USA	0.957
QKV59990.1	India	0.958	QJU11458.1	USA	0.957	QLC94305.1	USA	0.957
QLF98084.1	India	0.957	QIS61075.1	USA	0.957	QJD48484.1	USA	0.958
QLH64816.1	India	0.958	QJW28449.1	USA	0.957	QLF95737.1	USA	0.958
QJY40506.1	India	0.958	QLC91905.1	USA	0.957	QKN20740.1	USA	0.958
QLH93202.1	India	0.957	QKG87159.1	USA	0.958	QLH57751.1	USA	0.959
QLA10165.1	India	0.957	QKU31638.1	USA	0.957	QLC46314.1	USA	0.958
QLG97055.1	Italy	0.958	QKU31746.1	USA	0.958	QKG81824.1	USA	0.958
QJS39568.1	Netherlands	0.958	QLF99991.1	USA	0.959	QKU53650.1	USA	0.957
QJS39520.1	Netherlands	0.958	QJX70192.1	USA	0.961	QKN19672.1	USA	0.957
QJS39616.1	Netherlands	0.957	QKGS88935.1	USA	0.963	QLA47776.1	USA	0.957
QLF78310.1	Poland	0.957	QLC91545.1	USA	0.957	QLG97484.1	USA	0.957
QJX45308.1	Poland	0.957	QLH01238.1	USA	0.957	QLI50282.1	USA	0.957
QLH56255.1	Saudi Arabia	0.958	QKN20824.1	USA	0.957	QKW88844.1	USA	0.958
QLH56231.1	Saudi Arabia	0.957	QJQ39297.1	USA	0.958	QKS90192.1	USA	0.958
QKU37034.1	Saudi Arabia	0.958	QLB39321.1	USA	0.958	QKV39840.1	USA	0.960
QLH56099.1	Saudi Arabia	0.957	QKV41592.1	USA	0.958	QIU81286.1	USA	0.957
QHZ00380.1	South Korea	0.955	QKV42875.1	USA	0.957	QKV07400.1	USA	0.957
QIU78768.1	Spain	0.956	QLA09656.1	USA	0.958			

The SE of all the ORF3a proteins is bounded by the global minima 0.943 and global maxima 0.968 which are indeed the same as the minima and maxima of the ORF3a proteins which belongs to the USA (Table 15). Clearly, the amount of disorderliness of the amino acids over the ORF3a proteins is extremely high.

Table 15: Maxima and minima of SEs across geo-locations

<i>Geo-location</i>	<i>Min</i>	<i>Max</i>	<i>Range</i>
Australia	0.955	0.96	0.005
India	0.955	0.959	0.004
USA	0.943	0.968	0.025
Bangladesh	0.955	0.959	0.004

315 The range of SE of the ORF3a proteins of SARS-CoV2 collected from USA is comparatively more than others and it ensures the wide variety of distinct ORF3a in USA patients. The SEs of 296 ORF3a proteins are plotted (Blue line) in the Fig.18. We found various non-smooth peaks and those are clearly the SEs of the ORF3a proteins of the USA patients and that is reconfirmed in the SE plot (Red line) of the ORF3a proteins of the USA.

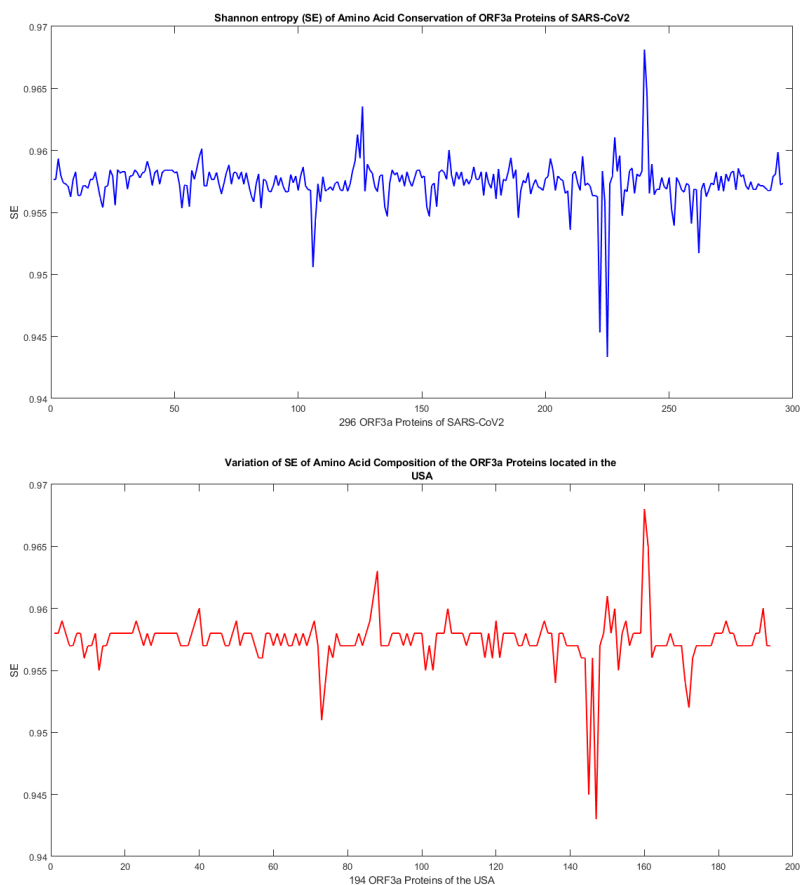


Figure 18: SE of amino acid compositions of ORF3a proteins

4. Discussions

320 A total of 175 distinct mutations across the distinct 256 ORF3a proteins of SARS-CoV2 are detected and further analyzed. Among all the mutations, 32 mutations were already reported [27, 23]. It was reported that in SARS-CoV, there exists an intensive interchain disulphide bonds with that of the spike protein with the help of the cysteine residues present in ORF3a protein. SARS-CoV2 ORF3a, contains a similar functional region (Domain III: C133) which is found to be conserved, as we did not find any mutation in this region. So, it can be assumed that these cysteine domains

325 perform a similar function as in SARS-CoV and is functionally important for virulence. In SARS-CoV, it was reported that tetramerization of the ORF3a protein is an important step for the ion channel formation which further increased the infectivity of the virus. From this study we found mutations W131C, T151I, R134L and D155Y which may facilitate the tetramerization process in SARS-CoV2 and thereby assisting the ion channel formation and favouring the virus with its infectivity. Similar to that of SARS-CoV, it is also responsible for apoptosis mediated by TRAF-3 (Domain III). We found 330 two mutations in this region Q38E and Q38P which may enhance the effect of apoptosis but further studies are required. Caveolin-binding domain is responsible for viral uptake of the host cell and its translocation to various endomembrane organelle. We have also isolated mutations in this zone (C148Y and A143S) which may enhance the viral uptake by the host, thereby increasing the infectivity rate. However, it is noteworthy that in $YXX\phi$ motif domain, no mutation is observed so far and consequently this domain is conserved. In seven ORF3a variants from the USA, two mutations are 335 found in SGD domain (S171L & G172C), however the function of this SGD domain is unknown.

We characterized the mutations into three types: Neutral, Disease and Mixed. Among these three mutations we found that disease mutations are highly prevalent with a percentage of 66% in the geo-location of USA, indicating disease-causing character of the virus getting intensified and thus posing threat to mankind. Simultaneously, we have the mixed type mutation occurring with a rate of 79% in the geo-location of USA. Mixed type had both disease and neutral occurring 340 together. Although, neutral mutations are there in mixed type but frequency of disease mutation is high, again pointing towards the viral advantage over host. In France although the infectivity rate was very high, but disease (2.9%) mutation rate was low compared to USA; where we find the maximum variety of mutation as shown with Shannon entropy in this study. So, we can suggest that the possible wide variety of mutations in USA is due to the high rate of travel 345 within USA and from outside USA, while in France there might be within-country transmission which resulted in less frequent mutations. We also checked the mortality rate of USA (3.3%), France (13.4%) and India (2.1%) and from the results we found that France has the highest mortality rate than USA followed by India. So, consequently we can draw a conclusion that since France has only disease type mutation unlike that of USA and India in which all three type of mutations are present. This may prove that the presence of only disease type of mutation in a sequence may pose more danger to mankind than a sequence containing either mixed type or neutral type of mutations. Next, we analysed 350 consecutive mutations within a protein sequence on the basis of chronological order of the time-line of sample collection from COVID-19 infected patients.

We further went on to analyse the mutations responsible for tetramerization and dimerization with respect to structure and found that there were no significant structural changes observed by homology modelling method. So, other method should be used to detect the effect of mutations on the 3D structure of the protein and results need to be experimentally 355 validated. Finally, twenty clusters are formed from 296 distinct variants of ORF3a of SARS-CoV2 based on the amino acid compositions of the proteins. It also shows wide variety of compositions of ORF3a variants in the USA which is further quantitatively supported by the SE. This study of comprehensive 175 novel mutations would help in understanding the pathogenetic contribution of the ORF3a proteins. This understanding is an important aspect in devising vaccine for COVID-19.

360 **Author Contributions**

SH conceived the problem. DA, SG and SH examined the mutations. All the authors analysed the data and result. SH wrote the initial draft which was checked and edited by all other authors to generate the final version.

Conflict of Interests

The authors do not have any conflicts of interest to declare.

365 Acknowledgement

[‡]Ms. Diksha Attrish and [†]Ms. Shinjini Ghosh are Interns under the supervision of Dr. Sk. Sarif Hassan through Virtual Internship with Science Leader (VISL) Programme, 2020. Authors thank to the Virtual Internship with Science Leader (VISL) program for their supports.

References

- 370 [1] J. Guarner, Three emerging coronaviruses in two decades: the story of sars, mers, and now covid-19 (2020).
- [2] Y. Huang, et al., The sars epidemic and its aftermath in china: a political perspective, in: Learning from SARS: preparing for the next disease outbreak: workshop summary, National Academies Press, 2004, pp. 116–136.
- [3] A. M. Al-Osail, M. J. Al-Wazzah, The history and epidemiology of middle east respiratory syndrome corona virus, Multidisciplinary respiratory medicine 12 (1) (2017) 20.
- 375 [4] A. Perrella, N. Carannante, M. Berretta, M. Rinaldi, N. Maturo, L. Rinaldi, Editorial—novel coronavirus 2019 (sars-cov2): a global emergency that needs new approaches, Eur Rev Med Pharmacol 24 (2020) 2162–2164.
- [5] J. M. Hintze, C. W. Fitzgerald, B. Lang, P. Lennon, J. B. Kinsella, Mortality risk in post-operative head and neck cancer patients during the sars-cov2 pandemic: early experiences, European Archives of Oto-Rhino-Laryngology (2020) 1–4.
- 380 [6] G. Fiorino, M. Allocca, F. Furfaro, D. Gilardi, A. Zilli, S. Radice, A. Spinelli, S. Danese, Inflammatory bowel disease care in the covid-19 pandemic era: the humanitas, milan, experience, Journal of Crohn's and Colitis (2020).
- [7] H. Harapan, N. Itoh, A. Yufika, W. Winardi, S. Keam, H. Te, D. Megawati, Z. Hayati, A. L. Wagner, M. Mudatsir, Coronavirus disease 2019 (covid-19): A literature review, Journal of Infection and Public Health (2020).
- [8] N. Van Doremalen, T. Bushmaker, D. H. Morris, M. G. Holbrook, A. Gamble, B. N. Williamson, A. Tamin, J. L. Harcourt, N. J. Thornburg, S. I. Gerber, et al., Aerosol and surface stability of sars-cov-2 as compared with sars-cov-1, New England Journal of Medicine 382 (16) (2020) 1564–1567.
- 385 [9] K. G. Andersen, A. Rambaut, W. I. Lipkin, E. C. Holmes, R. F. Garry, The proximal origin of sars-cov-2, Nature medicine 26 (4) (2020) 450–452.
- [10] P. T. Law, C.-H. Wong, T. C. Au, C.-P. Chuck, S.-K. Kong, P. K. Chan, K.-F. To, A. W. Lo, J. Y. Chan, Y.-K. Suen, et al., The 3a protein of severe acute respiratory syndrome-associated coronavirus induces apoptosis in vero e6 cells, Journal of general virology 86 (7) (2005) 1921–1930.
- 390 [11] J. L. Meitzler, S. Hinde, B. Bánfi, W. M. Nauseef, P. R. O. de Montellano, Conserved cysteine residues provide a protein-protein interaction surface in dual oxidase (duox) proteins, Journal of Biological Chemistry 288 (10) (2013) 7147–7157.

- 395 [12] J. To, J. Torres, Beyond channel activity: protein-protein interactions involving viroporins, in: *Virus Protein and Nucleoprotein Complexes*, Springer, 2018, pp. 329–377.
- [13] X. Tang, C. Wu, X. Li, Y. Song, X. Yao, X. Wu, Y. Duan, H. Zhang, Y. Wang, Z. Qian, et al., On the origin and continuing evolution of sars-cov-2, *National Science Review* (2020).
- [14] Z. Shen, Y. Xiao, L. Kang, W. Ma, L. Shi, L. Zhang, Z. Zhou, J. Yang, J. Zhong, D. Yang, et al., Genomic diversity of sars-cov-2 in coronavirus disease 2019 patients, *Clinical Infectious Diseases* (2020).
400
- [15] T. Phan, Genetic diversity and evolution of sars-cov-2, *Infection, genetics and evolution* 81 (2020) 104260.
- [16] Y.-Z. Zhang, E. C. Holmes, A genomic perspective on the origin and emergence of sars-cov-2, *Cell* (2020).
- [17] U. J. Buchholz, A. Bukreyev, L. Yang, E. W. Lamirande, B. R. Murphy, K. Subbarao, P. L. Collins, Contributions of the structural proteins of severe acute respiratory syndrome coronavirus to protective immunity, *Proceedings of the National Academy of Sciences* 101 (26) (2004) 9804–9809.
405
- [18] Y. Gao, L. Yan, Y. Huang, F. Liu, Y. Zhao, L. Cao, T. Wang, Q. Sun, Z. Ming, L. Zhang, et al., Structure of the rna-dependent rna polymerase from covid-19 virus, *Science* 368 (6492) (2020) 779–782.
- [19] W. Lu, K. Xu, B. Sun, Sars accessory proteins orf3a and 9b and their functional analysis, in: *Molecular Biology of the SARS-Coronavirus*, Springer, 2010, pp. 167–175.
- 410 [20] W. Lu, B.-J. Zheng, K. Xu, W. Schwarz, L. Du, C. K. Wong, J. Chen, S. Duan, V. Deubel, B. Sun, Severe acute respiratory syndrome-associated coronavirus 3a protein forms an ion channel and modulates virus release, *Proceedings of the National Academy of Sciences* 103 (33) (2006) 12540–12545.
- [21] K.-L. Siu, K.-S. Yuen, C. Castano-Rodriguez, Z.-W. Ye, M.-L. Yeung, S.-Y. Fung, S. Yuan, C.-P. Chan, K.-Y. Yuen, L. Enjuanes, et al., Severe acute respiratory syndrome coronavirus orf3a protein activates the nlrp3 inflammasome by promoting traf3-dependent ubiquitination of asc, *The FASEB Journal* 33 (8) (2019) 8865–8877.
415
- [22] K. Wang, S. Xie, B. Sun, Viral proteins function as ion channels, *Biochimica et Biophysica Acta (BBA)-Biomembranes* 1808 (2) (2011) 510–515.
- [23] E. Issa, G. Merhi, B. Panossian, T. Salloum, S. Tokajian, Sars-cov-2 and orf3a: Nonsynonymous mutations, functional domains, and viral pathogenesis, *Msystems* 5 (3) (2020).
- 420 [24] K. Padhan, C. Tanwar, A. Hussain, P. Y. Hui, M. Y. Lee, C. Y. Cheung, J. S. M. Peiris, S. Jameel, Severe acute respiratory syndrome coronavirus orf3a protein interacts with caveolin, *Journal of General Virology* 88 (11) (2007) 3067–3077.
- [25] R. Minakshi, K. Padhan, The yxx ϕ motif within the severe acute respiratory syndrome coronavirus (sars-cov) 3a protein is crucial for its intracellular transport, *Virology journal* 11 (1) (2014) 75.
- 425 [26] Y. Ren, T. Shu, D. Wu, J. Mu, C. Wang, M. Huang, Y. Han, X.-Y. Zhang, W. Zhou, Y. Qiu, et al., The orf3a protein of sars-cov-2 induces apoptosis in cells, *Cellular & molecular immunology* (2020) 1–3.

- [27] S. S. Hassan, P. P. Choudhury, P. Basu, S. S. Jana, Molecular conservation and differential mutation on orf3a gene in indian sars-cov2 genomes, *Genomics* (2020).
- [28] E. Smirnova, A. E. Firth, W. A. Miller, D. Scheidecker, V. Brault, C. Reinbold, A. M. Rakotondrafara, B. Y.-W. Chung, V. Ziegler-Graff, Discovery of a small non-aug-initiated orf in polioviruses and luteoviruses that is required for long-distance movement, *PLoS Pathog* 11 (5) (2015) e1004868.
- [29] D. M. Kern, B. Sorum, C. M. Hoel, S. Sridharan, J. P. Remis, D. B. Toso, S. G. Brohawn, Cryo-em structure of the sars-cov-2 3a ion channel in lipid nanodiscs, *BioRxiv* (2020).
- [30] E. Capriotti, R. B. Altman, Y. Bromberg, Collective judgment predicts disease-associated single nucleotide variants, *BMC genomics* 14 (S3) (2013) S2.
- [31] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, I. H. Witten, The weka data mining software: an update, *ACM SIGKDD explorations newsletter* 11 (1) (2009) 10–18.
- [32] E. Capriotti, P. Fariselli, R. Casadio, I-mutant2. 0: predicting stability changes upon mutation from the protein sequence or structure, *Nucleic acids research* 33 (suppl_2) (2005) W306–W310.
- [33] A. Likas, N. Vlassis, J. J. Verbeek, The global k-means clustering algorithm, *Pattern recognition* 36 (2) (2003) 451–461.
- [34] W. Zhong, G. Altun, R. Harrison, P. C. Tai, Y. Pan, Improved k-means clustering algorithm for exploring local protein sequence motifs representing common structural property, *IEEE transactions on Nanobioscience* 4 (3) (2005) 255–265.
- [35] B. J. Strait, T. G. Dewey, The shannon information entropy of protein sequences, *Biophysical journal* 71 (1) (1996) 148–155.
- [36] The Mathworks, Inc., Natick, Massachusetts, MATLAB version 9.3.0.713579 (R2020a) (2020).
- [37] X. Wang, Q. Zhou, Y. He, L. Liu, X. Ma, X. Wei, N. Jiang, L. Liang, Y. Zheng, L. Ma, et al., Nosocomial outbreak of covid-19 pneumonia in wuhan, china, *European Respiratory Journal* 55 (6) (2020).
- [38] D. J. Brooks, J. R. Fresco, A. M. Lesk, M. Singh, Evolution of amino acid frequencies in proteins over deep time: inferred order of introduction of amino acids into the genetic code, *Molecular Biology and Evolution* 19 (10) (2002) 1645–1655.
- [39] F. Johansson, H. Toh, Relative von neumann entropy for evaluating amino acid conservation, *Journal of bioinformatics and computational biology* 8 (05) (2010) 809–823.
- [40] F. Madeira, Y. M. Park, J. Lee, N. Buso, T. Gur, N. Madhusoodanan, P. Basutkar, A. R. Tivey, S. C. Potter, R. D. Finn, et al., The embl-ebi search and sequence analysis tools apis in 2019, *Nucleic acids research* 47 (W1) (2019) W636–W641.