

RESEARCH

Open Access



In depth analysis of Cyprus-specific mutations of SARS-CoV-2 strains using computational approaches

Anastasis Oulas^{1,2*}, Jan Richter³, Maria Zanti^{1,4,2}, Marios Tomazou^{1,5,2}, Kyriaki Michailidou^{4,2}, Kyroula Christodoulou^{5,2}, Christina Christodoulou^{3,2} and George M. Spyrou^{1,2}

Abstract

Background: This study aims to characterize SARS-CoV-2 mutations which are primarily prevalent in the Cypriot population. Moreover, using computational approaches, we assess whether these mutations are associated with changes in viral virulence.

Methods: We utilize genetic data from 144 sequences of SARS-CoV-2 strains from the Cypriot population obtained between March 2020 and January 2021, as well as all data available from GISAID. We combine this with countries' regional information, such as deaths and cases per million, as well as COVID-19-related public health austerity measure response times. Initial indications of selective advantage of Cyprus-specific mutations are obtained by mutation tracking analysis. This entails calculating specific mutation frequencies within the Cypriot population and comparing these with their prevalence world-wide throughout the course of the pandemic. We further make use of linear regression models to extrapolate additional information that may be missed through standard statistical analysis.

Results: We report a single mutation found in the *ORF1ab* gene (nucleotide position 18,440) that appears to be significantly enriched within the Cypriot population. The amino acid change is denoted as S6059F, which maps to the SARS-CoV-2 NSP14 protein. We further analyse this mutation using regression models to investigate possible associations with increased deaths and cases per million. Moreover, protein structure prediction tools show that the mutation infers a conformational change to the protein that significantly alters its structure when compared to the reference protein.

Conclusions: Investigating Cyprus-specific mutations for SARS-CoV-2 can lead to a better understanding of viral pathogenicity. Researching these mutations can generate potential links between viral-specific mutations and the unique genomics of the Cypriot population. This can not only lead to important findings from which to battle the pandemic on a national level, but also provide insights into viral virulence worldwide.

Keywords: SARS-CoV-2, Cyprus-specific mutations, Linear regression, Structural prediction

* Correspondence: anastasios@cing.ac.cy

¹Bioinformatics Department, Cyprus Institute of Neurology and Genetics, Nicosia, Cyprus

²The Cyprus School of Molecular Medicine, Nicosia, Cyprus

Full list of author information is available at the end of the article



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

Background

Cyprus' first case of COVID-19 was reported on March 9th 2020, with a gradual increase after that to formulate the first peak of the virus transmission seen in late March with all major Cypriot cities affected by the virus. This first wave was relatively mild with respect to reported daily new cases and deaths reaching a maximum of 58 cases on April 1st and 2 deaths. The second wave of virus transmissions hit Cyprus in mid-October and gradually increased to peak in December with a maximum of 907 and 8 reported daily new cases and deaths, respectively. This summary of the SARS-CoV-2 spread in the Cypriot population dictates that Cyprus is one of the least affected European countries during this pandemic (mostly with respect to deaths per million). This can be attributed to the rapid response time for austerity measures and the effective quarantine process for confirmed cases and their contacts. In addition, Cyprus is ranked 3rd in Europe as for the number of COVID tests performed per million (> 2,4 M tests).

SARS-CoV-2 is a viral quasispecies, therefore, humanity is facing a mutant cloud with trillions of different combinations [1]. Multiple national studies have been undertaken by a plethora of countries throughout the world, in order to generate and analyse high throughput sequencing country-specific data for SARS-CoV-2 strains [2–9]. Recently, the Cyprus Institute of Neurology and Genetics (CING) also published a study based on 144 NGS samples obtained from the Cypriot population representing the first documented genomic and epidemiological characterization of these samples [10]. In this study we expand on this work, by initially performing basic lineage analysis, as previously reported [10], to identify the dominant SARS-CoV-2 lineage(s) in Cyprus as well as a phylogenetic tree analysis to map the Cypriot strains against other strains present throughout the world. We further perform a more thorough, in depth genomic analysis of these samples by implementing variant-calling and displaying an overview of the genomic variation and reported frequencies in the Cypriot strains. We then focus on the 9 spike (S) protein mutations that delineate the UK B.1.1.7 lineage and how their founder effect in Cyprus has impacted the number of cases and deaths in the population. Furthermore, we identify Cypriot-specific/dominant mutations and investigate them using mutation tracking analysis in order to isolate their origin, determine their founder effect in Cyprus and also trace their overall prevalence and propagation in other countries. We capitalize on our previously published generalized linear models [11] to undertake virulence analysis on Cypriot-dominant mutations and show how their presence has affected the number of cases and deaths within the country. Finally, we perform structural modelling of the alternate versus

the reference mutation for selected mutations of interest and view their effects at the viral protein level.

Methods

Raw data analysis

The Burrows-Wheeler Aligner (BWA) [12], version: 0.7.15 was used to map the raw reads to Wuhan-Hu-1 (NCBI ID:NC_045512.2). Duplicate reads, which are likely to be the results of PCR bias, were marked using Picard (<http://broadinstitute.github.io/picard/>) version: 2.6.0. SAMtools [13], version: 0.1.19, was used for additional BAM/SAM file manipulations. The Genome Analysis Tool Kit (GATK) [14], version 3.6.0, Haplotype-Caller method was used for single nucleotide polymorphism (SNP) and insertion/deletion (indel) variant calling generating vcf files. All mutations were annotated according to the reference Wuhan-Hu-1 strain. These annotations include: genomic position, reference (REF) genotype, alternate (ALT) genotype, gene, encoded amino-acid with REF genotype, encoded amino-acid change (if not synonymous) with ALT genotype, amino-acid position according to main encoded viral protein (annotation file available as Supplementary Table S1). Finally, the GATK FastaAlternateReferenceMaker method was used for consensus sequence extraction from the vcf files.

Lineage assignment

The consensus sequences of all 144 Cypriot strains were uploaded to the Pangolin COVID-19 lineage assigner interface [15] (<https://pangolin.cog-uk.io/>). Further analysis of results and generation of visualizations was performed in R V.3.6.1 [16].

Phylogenetic analysis and comparison with other strains

Full genomic SARS-CoV-2 sequences of high sequencing resolution were obtained from GISAID (<https://www.gisaid.org/>, last accessed 15/01/21). The nextstrain [17] pipeline was downloaded locally and the commands for filtering, aligning and constructing phylogeny were used according to the nextstrain's best practices. MAFFT [18] was used to construct a multiple sequence alignment (MSA). Phylogeny was estimated using the RAXML [19] maximum likelihood algorithm for phylogenetic tree construction. Variant calling for the GISAID strains was achieved using the `snp_sites` tool available through github (<https://github.com/sanger-pathogens/snp-sites>). `Snpsites` was the only tool we could find that effectively calls variants directly from the fasta sequences, without the need for raw fastq files and bam files. A limitation of this tool is that it only calls for single nucleotide variants that are snp-like, therefore insertions and deletions are not obtained by this tool. We could therefore only track

7/9 mutations (excluding the 2 deletions) in the UK lineage tracking analysis.

Mutation tracking analysis

Relative frequencies across countries with at least one occurrence of the selected mutations of interest was visualized as bar plots across time (months). This provides an indication of the spread of the studied mutations across the general population. Analysis was performed using R (packages: dplyr, tidyr, ggplot2, ggtree, phytools, phangorn).

Identifying Cyprus-specific mutations

A generalized linear model was applied (see Formula 1 below), whereby the absence or, presence (values 0 or 1 respectively) of the S6059F mutation (*mut* variable) was assessed across samples from Cyprus, as well as the rest of the world (as defined by the *classes* variable, 1 denoting Cyprus and 0 for other countries).

$$\text{model} < -\text{glm}(\text{mut} \sim \text{classes}, \text{data} = \text{vcf}, \text{family} = \text{binomial}(\text{link} = \text{"logit"})) \quad (1)$$

An *odds* ratio analysis was also performed using a simple formula (see 2 below). Whereby, *cycounts* denotes the number of times the S6059F mutation was observed in Cyprus (0 for absence, 1 for presence) and similarly *othercounts* denotes the number of times the mutation was observed in other countries.

$$\text{OR} < -(\text{cycounts1} * \text{othercounts0}) / (\text{othercounts1} * \text{cycounts0}) \quad (2)$$

Structural analysis

Comparative structural modelling was carried out for unknown protein structures using the template-based web server, I-TASSER [20]. The accuracy of the method was assessed based on the I-TASSER template modeling score (TM-score), which indicates the structural similarity between model and templates. A TM-score higher than 0.5 indicates a model of correct topology, while a TM-score of less than 0.17 means a random topology [20]. I-TASSER was selected for protein structure modelling, since it outperformed other servers according to results from the 14th Community Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction (CASP14) (<https://zhanglab.ccmb.med.umich.edu/casp14/>, last accessed 23/03/2021). Mutagenesis was performed using the DynaMut suite (<http://biosig.unimelb.edu.au/dynamut/prediction>, last accessed 23/03/2021). The PyMOL software (v0.99) was used for the visualization of the protein molecules.

Relative solvent accessibility (RSA) was calculated using the Missense3D algorithm [21]. The DynaMut webserver [22], was used to visualize non-covalent molecular interactions, calculated by the Arpeggio algorithm [23]. Protein-protein complexes were constructed using the ClusPro (v2.0) [24] and HDOCK [25] algorithms and binding affinities and dissociation constants (Kd) were calculated using the PRODIGY webserver [26]. RNA-protein docking simulations were carried out using the HDOCK [27] and MPRDock algorithms. For RNA-protein docking simulations, as active residues we selected the active site residues of the SARS-CoV NSP14 protein, since the two proteins share a 99.1% sequence similarity [28]. Structural alignment was performed using the align tool of PyMOL and all-atom RMSD values were calculated without any outliers' rejection, with zero cycles of refinement. All docking simulations were performed in triplicates.

Results

Dominant lineages present throughout the COVID-19 pandemic in Cyprus – emphasis on UK lineage B.1.1.7

One of the most widely used international systems for detecting lineages that contribute most to active spread is the dynamic nomenclature system presented by Rambaut et al. [15]. Lineage analysis of the 144 Cypriot viral strain genomes showed 16 major lineages including strains originating from both A and B lineages, which are denoted as the root lineages of the phylogeny of SARS-CoV-2 (Fig. 1). Dominant lineages included: (i) B.1.258 (51.03% of sequenced strains), with most common countries of origin being UK, Denmark and Czech Republic and (ii) the UK lineage, B.1.1.7, having drawn attention by the recent outbreak in the UK with reported increased rates of viral transmission [29–31], which was detected with high prevalence within the Cypriot population (13.1% of analysed strains). This lineage is characterized by a series of 9 mutations (deletion HV69–70, deletion Y144, N501Y, A570D, P681H, T716I, S982A, D1118H, D614G) within the viral S protein. The D614G mutation has now become dominant in all strains but was included here for completion. Mutation tracking for 7/9 (see Methods for details) of these mutations reveals that the specific lineage was firstly reported in Cyprus on December 2020. This was the first reported case where the UK lineage happened to be sequenced and was derived from Cypriot citizens returning to Cyprus after visiting the UK. We acknowledge that there could have been other routes of importation from UK, or even other countries, around the same time period that were not catalogued. Moreover, the founder effect of this lineage within the Cypriot population appears to be associated with an increase in number of cases, a phenomenon also reported for the UK [31] (see Fig. 2).

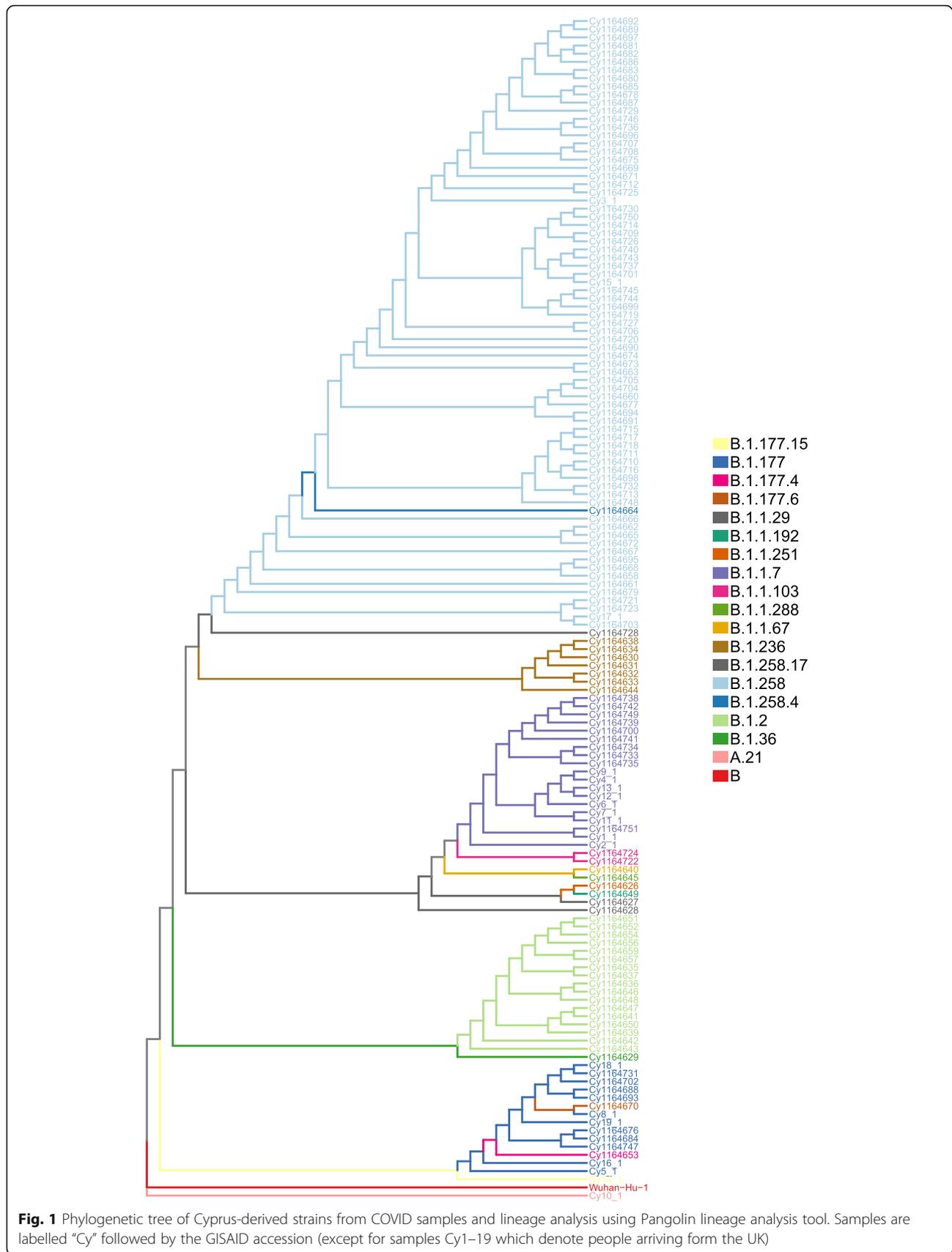
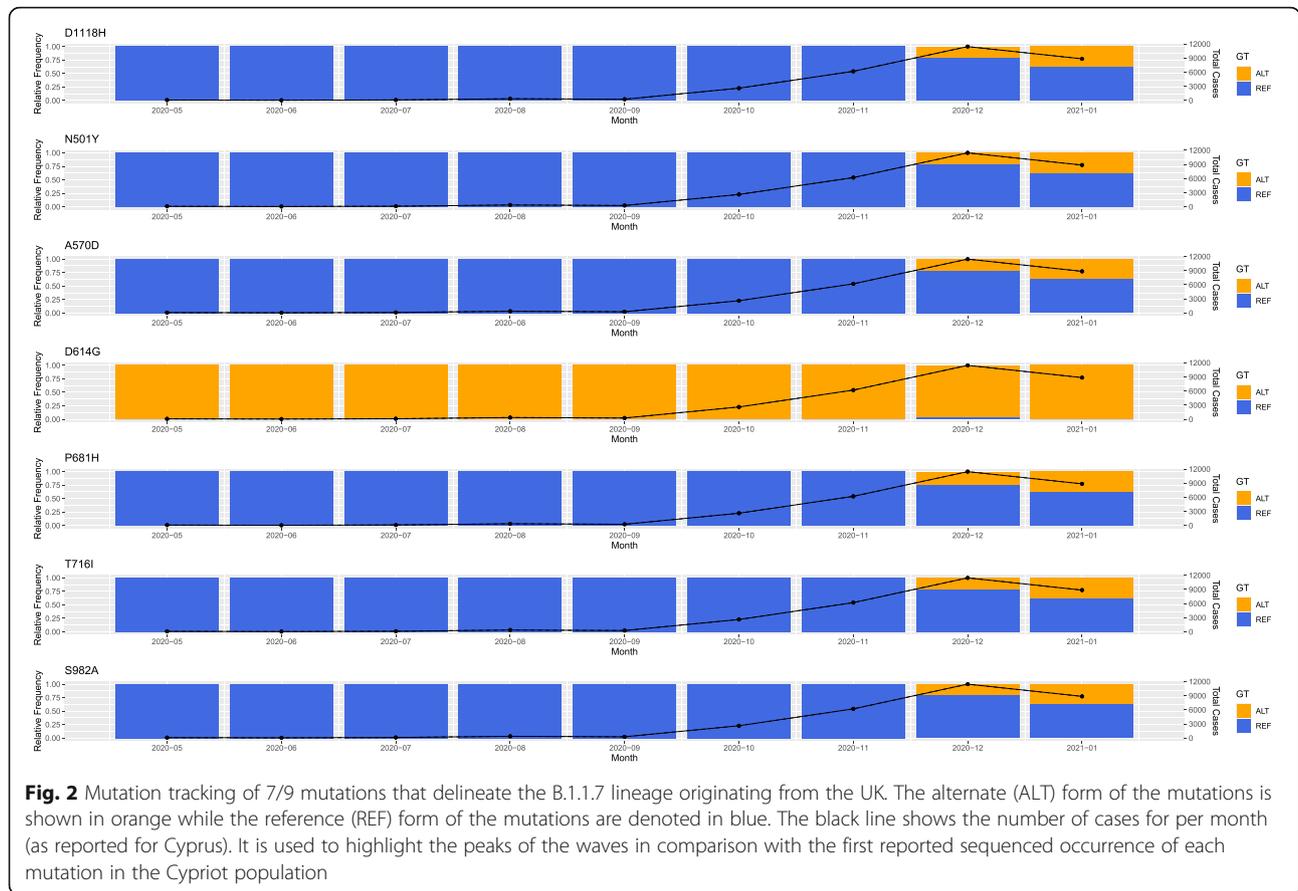


Fig. 1 Phylogenetic tree of Cyprus-derived strains from COVID samples and lineage analysis using Pangolin lineage analysis tool. Samples are labelled "Cy" followed by the GISAID accession (except for samples Cy1–19 which denote people arriving from the UK)



Cyprus specific mutations

Cyprus-specific mutations were obtained by counting the frequency of all single nucleotide mutations present within the vcf file generated from our 144 sequenced samples. In order to obtain a manageable number of prevalent mutation within the Cypriot population an arbitrary threshold was set that filters out for most rare mutations in the dataset. We selected for mutations that appeared in at least 30% of the samples and thus identified 18 mutations that met this criterion. These included mutations that were very cosmopolitan (e.g. D614G mutation) but also included mutations that appeared primarily in the Cypriot population according to our sequenced samples. We performed an odds ratio (OR) analysis to investigate which mutations appear to be more frequent in the Cypriot population vs. strains obtained from the rest of the world. We also applied a generalized linear model (glm) with a logistic regression using the *logit* function in order to obtain a model for each of the 18 selected mutations. The model assumes a value of 1 if the mutation is present and 0 if not, through this way trying to see how well this fits the ideal situation of mutations only occurring in the Cypriot samples vs. the rest of the world (for details see

Materials and methods). *P-values* were generated to show how well the glm performs under this scenario. We specifically highlight the S6059F mutation, located in the NSP14 protein of the *ORF1ab* gene, because it obtained the highest OR statistic (6921.17) and was also deemed most significant according to our glm model (*p-value* = 6.94E-180) (see Table 1). We performed mutation tracking analysis to further investigate the prevalence of this specific mutation throughout the course of the pandemic, both in Cyprus as well as worldwide. The founder effect appeared to have occurred in Russia, as shown by timestamps for the strains as well as phylogenetic tree analysis for the specific strains with the S6059F mutation (see Fig. 3). This event did not lead to the establishment of this mutation as the dominant form within the Russian population. Similar results appear for other countries with at least one reported strain with this mutation. This is in contrast to Cyprus where the alternate S6059F mutation clearly appears to be replacing the reference form (see Fig. 4).

The NSP14 protein is bifunctional and contains two domains: a 3'-to-5' exonuclease (ExoN) and a guanine-N7-methyltransferase (N7-MTase). It is presumed that the N7-MTase domain supports mRNA

Table 1 Cyprus specific mutations ordered by *p*-value for best fit according to our glm model

REF ^a	ALT ^a	GENE ^a	NT.POS ^a	REF.AA ^a	ALT.AA ^a	AA.POS ^a	CY counts1 ¹	CY counts0 ²	Other counts1 ³	Other counts0 ⁴	P-value ⁵	OR ⁶
C	T	ORF1ab	18,440	S	F	6059	60	84	25	242,241	6.94E-180	6921.17
G	T	ORF1ab	11,557	E	D	3764	69	75	928	241,338	1.03E-124	239.26
C	T	E	26,313	F	F	23	69	75	1016	241,250	4.25E-122	218.45
C	T	ORF1ab	20,451	N	N	6729	76	68	2983	239,283	5.26E-104	89.65
T	C	M	26,972	R	R	150	77	67	3988	238,278	1.86E-96	68.67
T	C	S	24,910	T	T	1116	76	68	3967	238,299	8.45E-95	67.14
G	A	ORF1ab	15,598	V	I	5112	76	68	3973	238,293	9.45E-95	67.03
G	T	ORF1ab	12,988	M	I	4241	76	68	3976	238,290	1.00E-94	66.98
G	T	ORF1ab	18,028	A	S	5922	76	68	4020	238,246	2.27E-94	66.24
C	A	ORF7b	27,800	A	A	15	76	68	4431	237,835	3.13E-91	59.99
T	C	ORF1ab	7767	I	T	2501	77	67	5213	237,053	1.06E-87	52.26
C	T	ORF1ab	8047	Y	Y	2594	76	68	5199	237,067	4.44E-86	50.96
C	T	ORF1ab	17,104	H	Y	5614	76	68	5486	236,780	2.38E-84	48.24
C	A	S	22,879	N	K	439	76	68	5501	236,765	2.92E-84	48.10
A	G	ORF1ab	20,268	L	L	6668	83	61	16,967	225,299	4.21E-57	18.07
C	T	ORF1ab	14,408	P	L	4715	143	1	225,633	16,633	2.08E-04	10.54
C	T	ORF1ab	3037	F	F	924	143	1	225,640	16,626	2.08E-04	10.54
A	G	S	23,403	D	G	614	143	1	225,843	16,423	2.36E-04	10.40

^aREF reference nucleotide, ALT alternate nucleotide, GENE gene name, NT.POS nucleotide position in gene, REF.AA reference amino acid, ALT.AA alternate amino acid, AA.POS amino acid position in protein

¹CY counts1 number of times ALT form of the S6059F mutation (value 1) was found in Cyprus strains

²CY counts0 number of times REF form of the S6059F mutation (value 0) was found in Cyprus strains

³Other counts1 number of times ALT form of the S6059F mutation (value 1) was found strains from other countries

⁴Other counts0 number of times REF form of the S6059F mutation (value 0) was found strains from other countries

⁵GLM *p*-value

⁶Odds ratio (OR) statistic

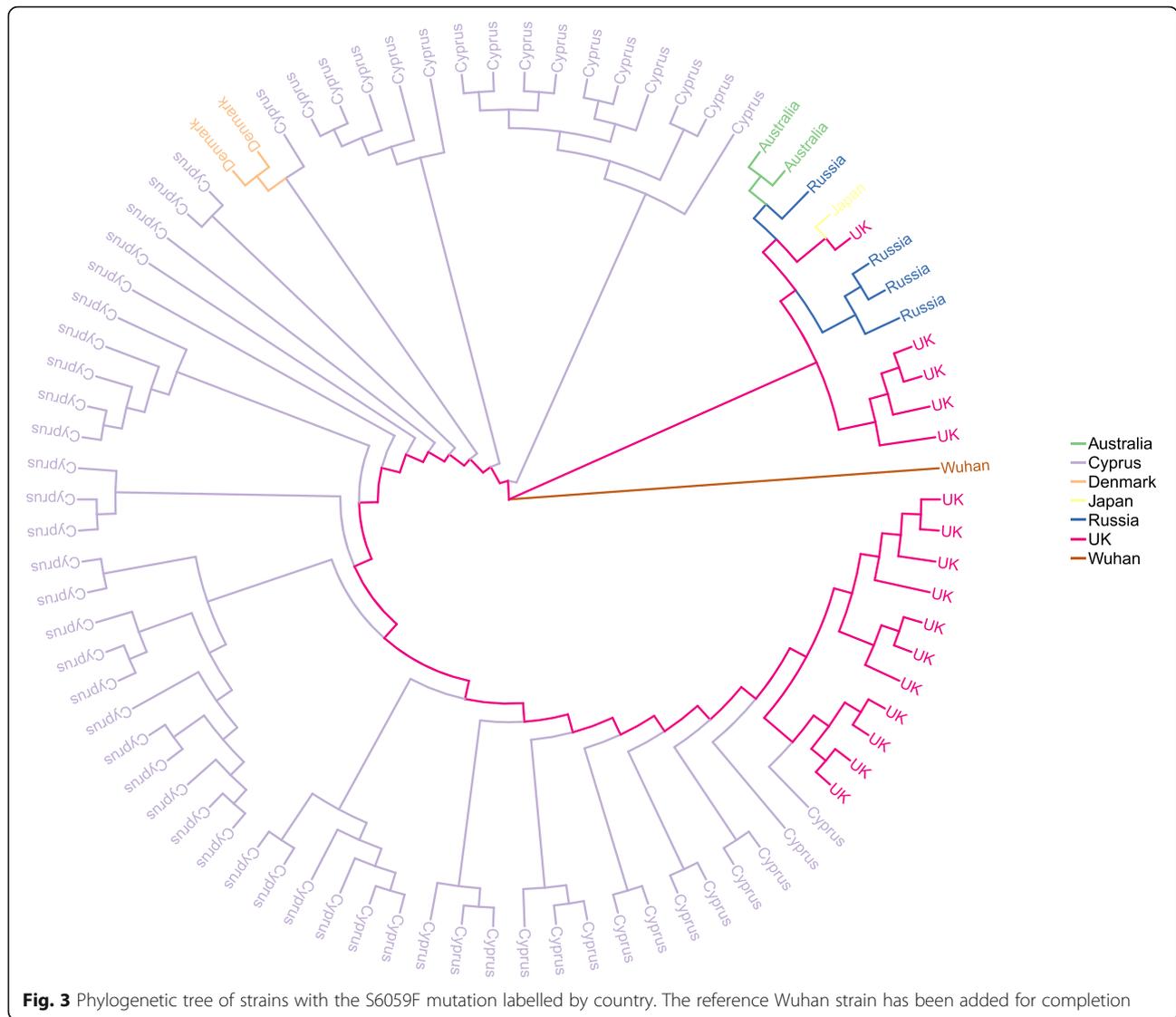
capping, while the ExoN domain is believed to mediate proofreading during genome replication [32]. Previous studies have shown that ExoN knockout mutants or severe acute respiratory syndrome coronavirus (SARS-CoV) exhibit a dysfunctional yet viable hypermutation phenotype both in cell culture as well as animal models [33–35] while a SARS-CoV-2 ExoN knockout mutant was found to be unable to replicate [32].

In order to investigate for proofreading dysfunctionality in the strains with the S6059F mutation within the Cypriot population, we compared mutation counts for all strains with the alternative (ALT) form of this mutation to the strains with the reference (REF) form. The 60 strains with the ALT S6059F mutation all belonged to the B.1.258 lineage and showed a higher mean (25.2) for the number of mutations when compared to the REF S6059F strains (23.4). As the data does not follow a normal distribution (according to Shapiro-Wilks normality test - *p*-value = 2.486e-05) we use the Wilcoxon test to assess whether the distributions are different. The results show that the distributions are significantly different (Wilcoxon test - *p*-value 0.00025 – see Fig. 5). We

propose that this difference in mutation counts within the strains with different forms of the S6059F mutation (ALT vs. REF) can be attributed to increased hypermutability as a consequence of the amino acid change in the NSP14 protein.

Virulence analysis of Cyprus-specific SARS-COV-2 mutations

Capitalizing on previously published generalized linear models (glms) that provide a measure of association between specific SARS-CoV-2 mutations and increased or decreased viral cases or deaths [11], we proceeded to perform virulence analysis on the 18 mutations found to be prevalent in the Cypriot population. We focused on the S6059F mutation due to its exclusive propagation within Cyprus compared to the rest of the world. Our glms are uniquely designed to incorporate mutation frequency, austerity measure response time and mutation occurrence information as predictor variables and cases or deaths per million as response variables into a single model. This allows for a fit that determines whether a

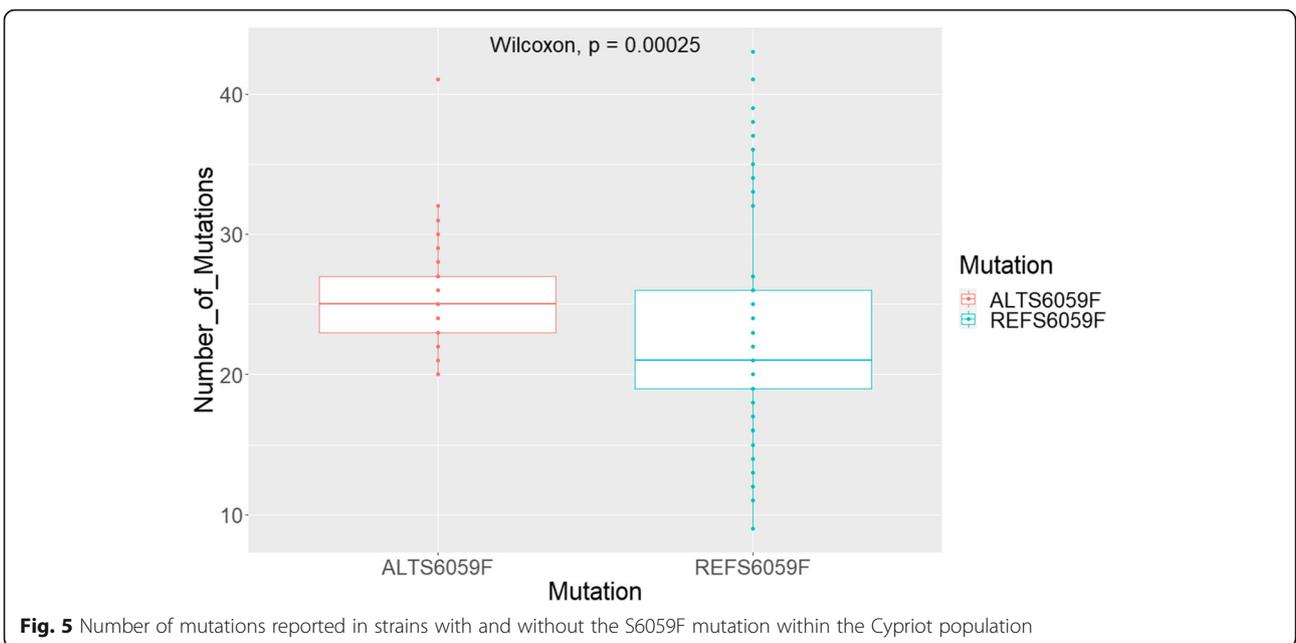
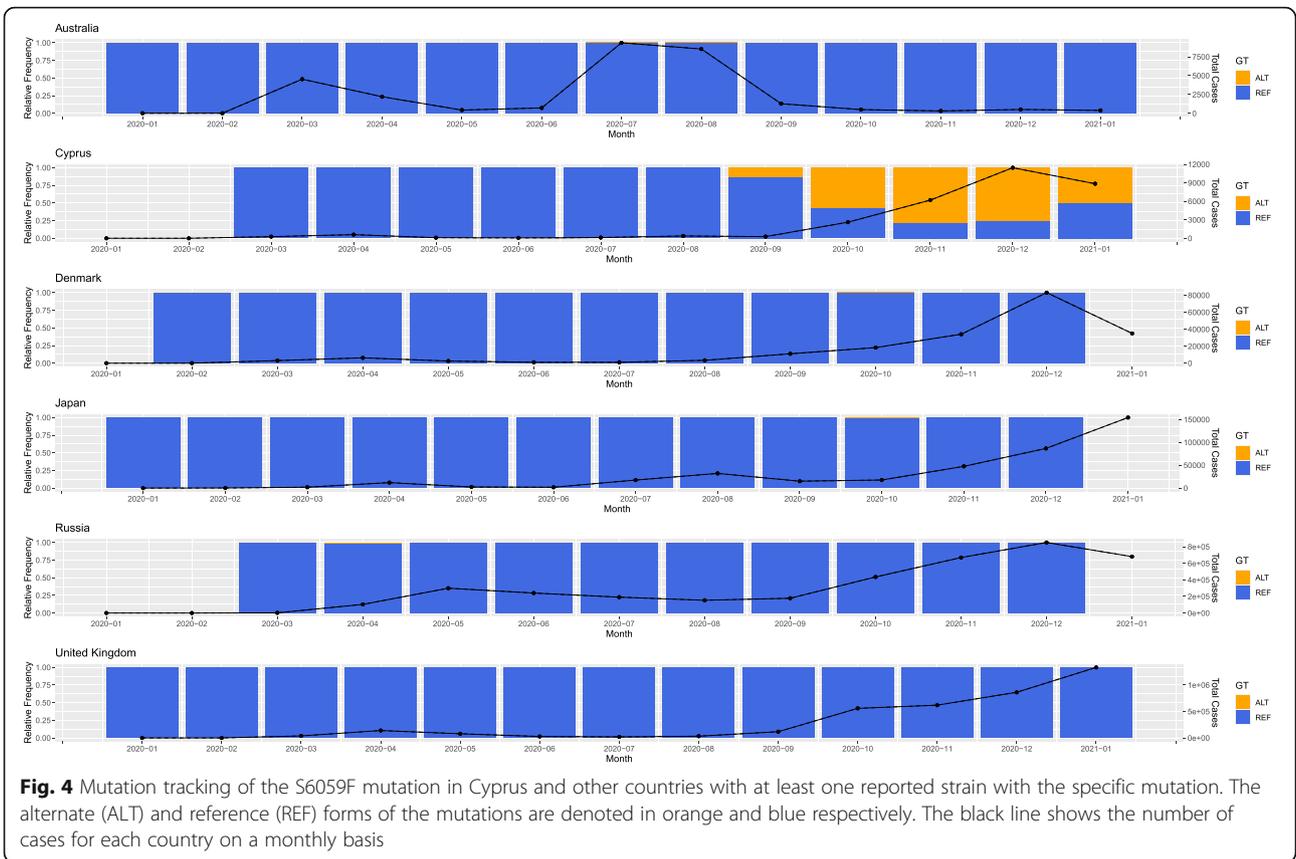


given mutation is positively or negatively correlated with number of cases and deaths per country.

We focused on geographic regional data reporting deaths and transmissions according to the [Worldometers.info](https://www.worldometers.info) website (last accessed 08/03/21). Populations harbouring the mutation show a higher mean of deaths per million (1363) compared to the converse (1206). Statistical analysis of the populations with and without the S6059F mutation shows substantial evidence that the two distributions are significantly different (Wilcoxon test – p -value $2.2e-16$) (see Fig. 6A). Results obtained from using the cases per million parameter to assess the different distributions, show a reversed pattern with samples with the mutation exhibiting a lower mean of cases (49224) compared to the reference samples (64630) (see Fig. 6B). Statistical analysis further provides evidence that the two

distributions are significantly different (Wilcoxon test – p -value $2.2e-16$). We also included response time separation in the box plots which shows how the mutation segregates across countries that responded differently to the COVID austerity measures (see supplementary Fig. S1A, B).

Applying our previously published glm model [11], on the 18 Cyprus-specific mutations allows for the analysis of these less studied mutations and how they appear to correlate with death and transmission rates. According to the fitting plot of our glm models, the unique S6059F mutation appears to be positively correlated with deaths per million ($R = 0.61$) and neutral with respect to cases ($R = 0.15$) per million (see supplementary Fig. S2). However, it should be noted that this mutation did not show significance according to the model obtained p -values ($p > 0.05$). This can most probably be attributed to its



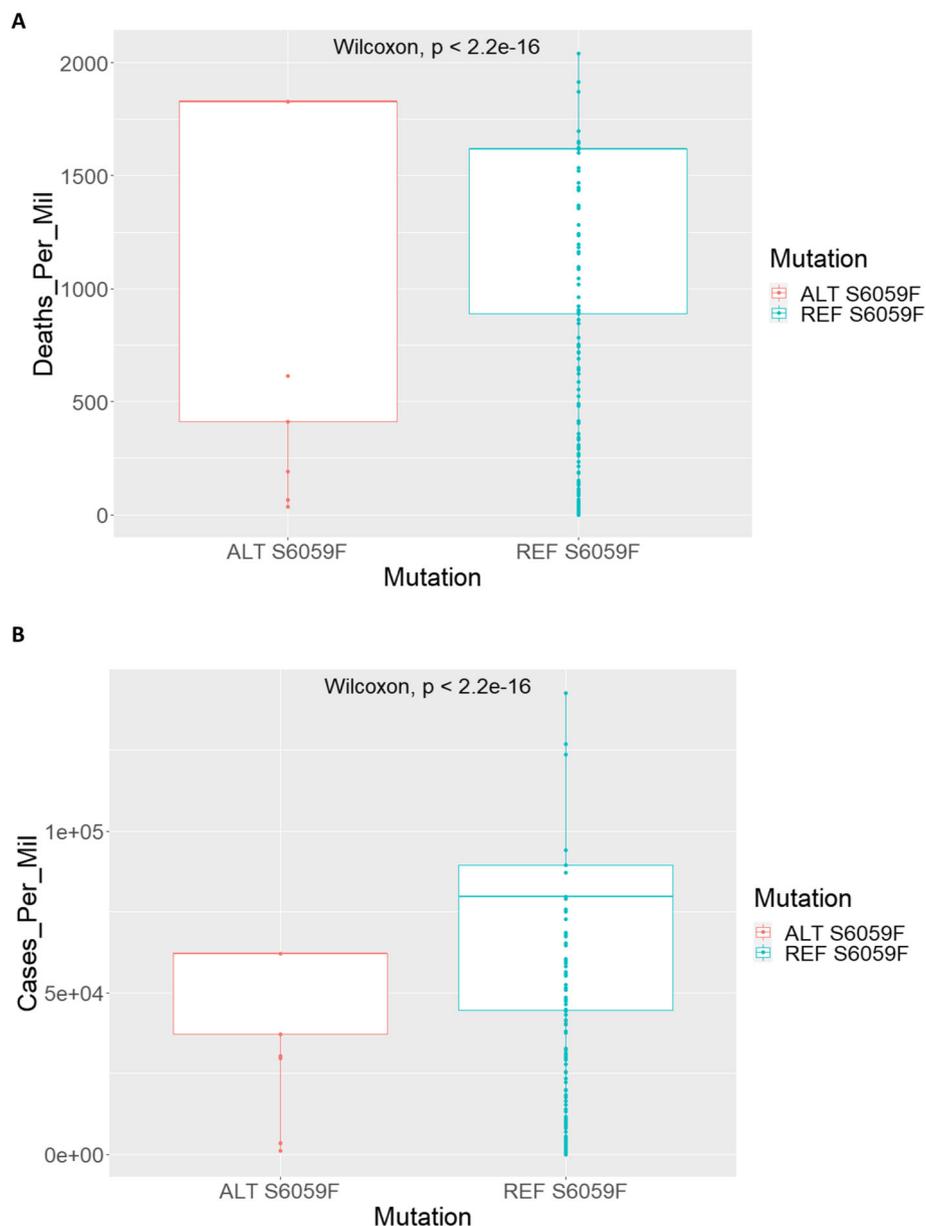


Fig. 6 Boxplot distributions for S6059F mutation. **A** Deaths per million for countries with the ALT and REF forms of the S6059F. **B** Cases per million for countries with the ALT and REF forms of the S6059F mutation. Note that each data point can represent more than one country

low frequency of occurrence across different countries besides Cyprus.

Structural analysis

In order to investigate potential reasons why the S6059F mutation appears to be associated with increased mutagenesis, we turn to structural prediction and docking tools. The Ser6059 residue is located on the surface of the NSP14 protein encoded by the *ORF1ab* gene, a bi-functional protein with an N-terminal exonuclease (ExoN) domain with a proofreading function and a C-

terminal N7-guanine methyltransferase (N7-MTase) domain, implicated in the methylation of viral RNA cap structures. Both domains are crucial for the maintenance and stability of the viral RNA by lowering its sensitivity to RNA mutagens and evading its degradation from the host immune response [28]. The QHD43415 structure (Estimated TM-score = 0.87) of the NSP14 protein was retrieved from the I-TASSER repository containing the 3D structural models of all proteins encoded by the genome of SARS-CoV-2. The S6059F residue is located in the NSP14 domain directly involved in the physical

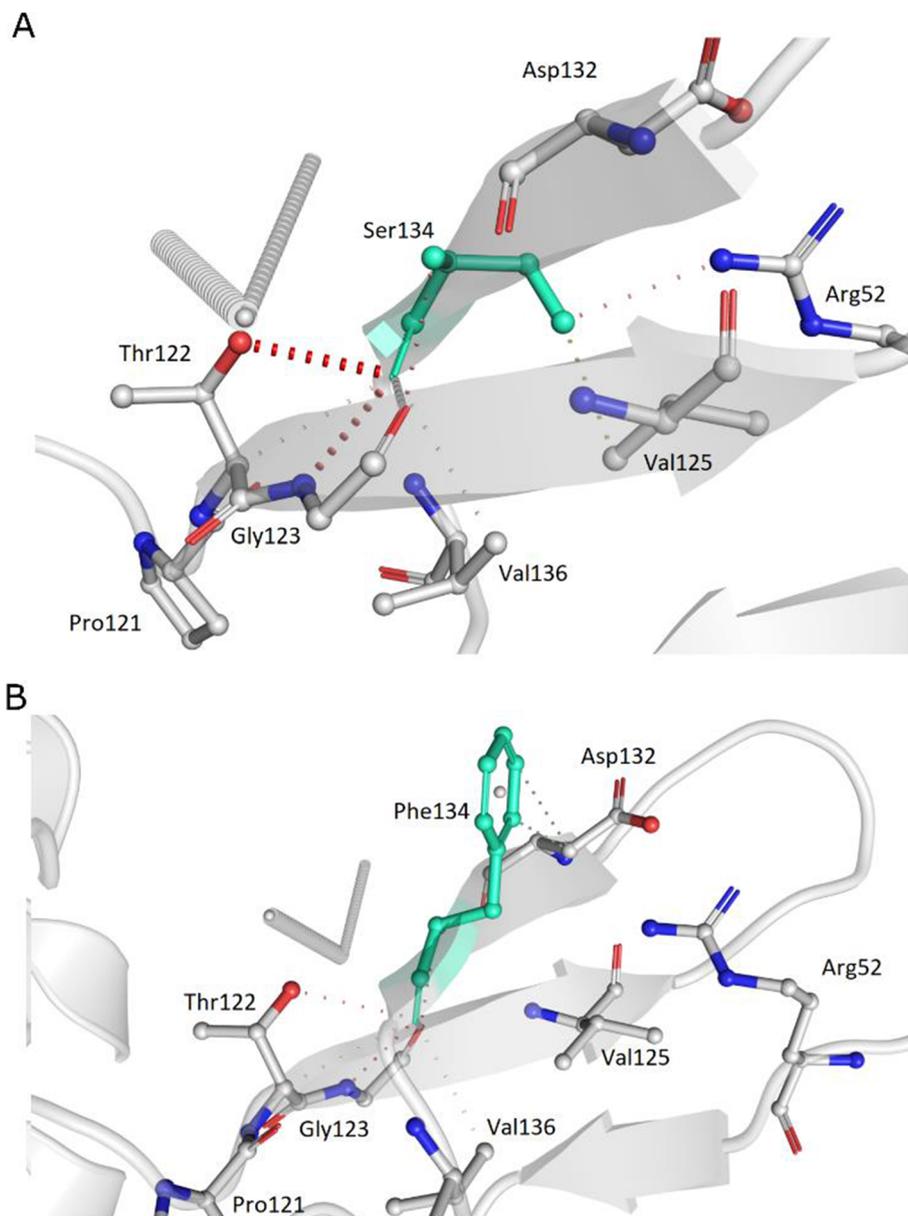
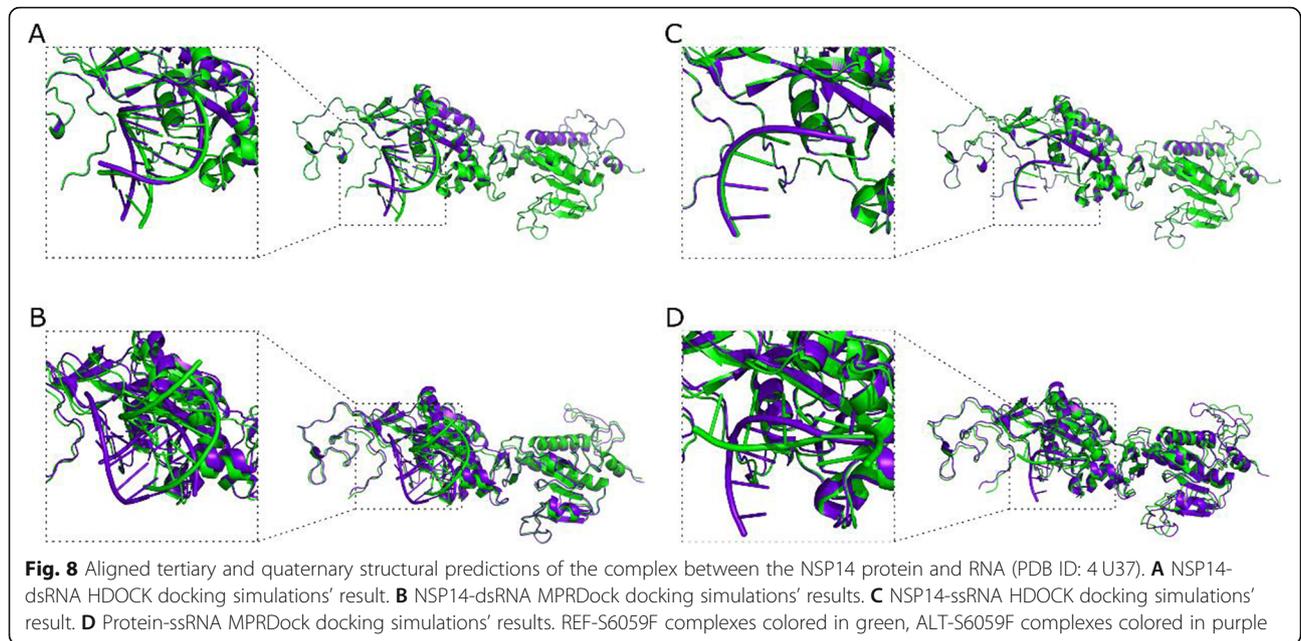


Fig. 7 Tertiary and quaternary structural predictions of the NSP14 protein. **A** Detailed molecular structural conformational changes of NSP14 showing the REF-S6059F. **B** Same molecular analysis for the structure with ALT-S6059F. Hydrogen bonds in red, weak hydrogen bonds in orange, halogen bonds in blue, ionic interactions in yellow, aromatic contacts in light-blue, hydrophobic contacts in green, carbonyl interactions in pink, VdW interactions in grey

interaction with NSP10 (residues 1–76 and 119–145), which enhances the ExoN activity by more than 35-fold [28]. The 6059 residue is close to the active site of the enzyme (D90, E92, E191 and D273) [28], and both amino acids are exposed to the surface with RSA values 24.6 and 51.7% for the uncharged Serine (S) and Phenylalanine (F) residues, respectively. Further, in depth molecular analysis reveals disruption of hydrogen bonds and hampering and shifts in other inter-molecular

interactions in the ALT-S6059F compared to the REF-S6059F structure of the protein (see Fig. 7A, B).

The experimentally-solved three-dimensional protein structure of the SARS-CoV-2 NSP10 protein was retrieved from the protein databank (PDB ID: 6W75). Protein-protein docking analysis revealed changes on the assembly of NSP14-NSP10 complexes ($\text{RMSD } 5.841 \pm 5.631 \text{ \AA}$), while the ALT-S6059F-NSP10 complex exhibited a minor increase in complex affinity ($\Delta G - 17.5 \pm$



0.141 kcal/mol and K_d $1.5E-13 \pm 2.828E-14$ M) compared to the REF-S6059F-NSP10 complex ($\Delta G - 16.55 \pm 0.495$ kcal/mol and K_d $8.35E-13 \pm 6.576E-13$ M).

Investigation of the NSP14-RNA interaction was also carried out, using the native 7mer-dsRNA (PDB ID: 4 U37) and the derived 7-mer-ssRNA. Upon RNA-protein docking using the HDOCK and MPRDock algorithms,

structural changes are evident as demonstrated upon structural alignment for both dsRNA (mean RMSD for both HDOCK and MPRDock 4.255 ± 0.222 Å) (Fig. 8A, B) and ssRNA complexes (mean RMSD for both HDOCK and MPRDock 1.946 ± 2.708 Å) (Fig. 8C, D).

Investigation of the NSP10-NSP14-RNA interaction was also carried out. Upon RNA-protein docking using

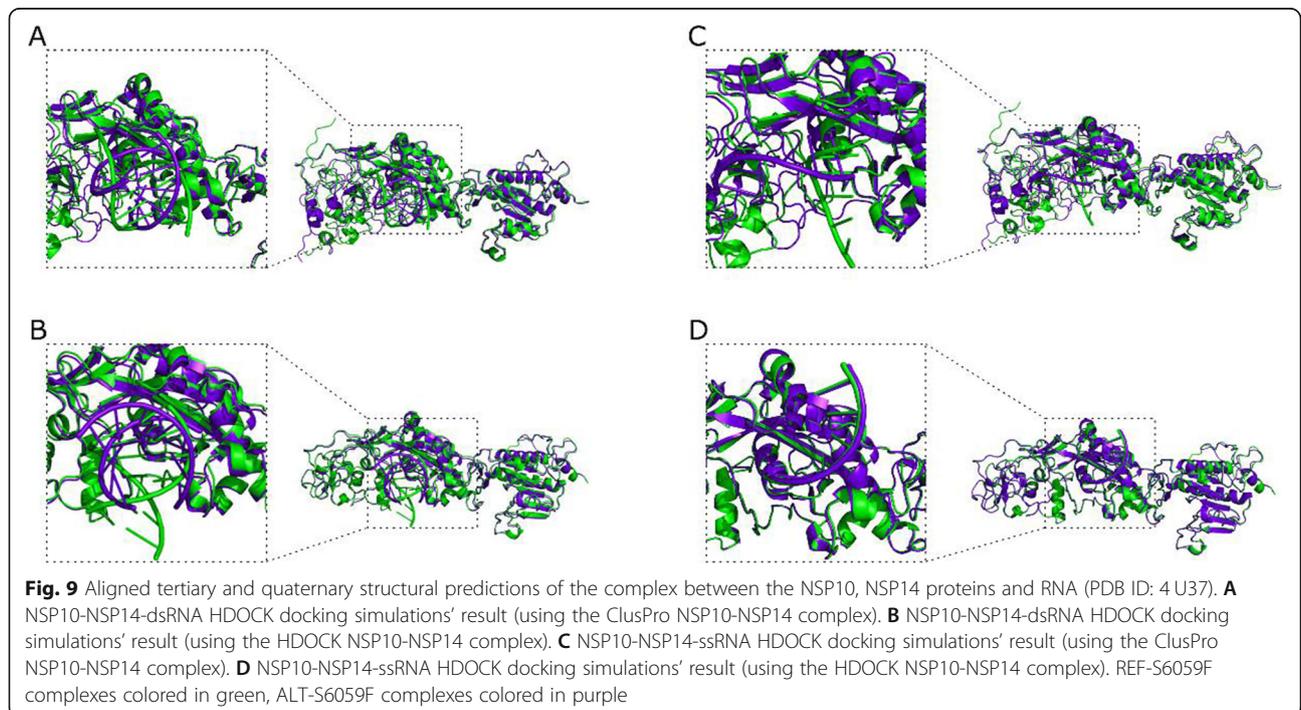


Table 2 HDOCK free energy and RMSD scores for REF- and ALT- NSP14-RNA and NSP14-NSP10-RNA complexes

	Mean docking score (kcal/mol)	Mean ligand rmsd (Å)
REF-S6059F-ssRNA	-210.759	147.084
ALT-S6059F-ssRNA	-210.208	147.022
REF-S6059F-dsRNA	-207.585	144.734
ALT-S6059F-dsRNA	-206.933	144.414
Nsp14-nsp10-REF-S6059F-ssRNA	-190.768, -197.32 ^a	80.506, 66.387
Nsp14-nsp10-ALT-S6059F-ssRNA	-195.228, -207.88	106.924, 67.26
Nsp14-nsp10-REF-S6059F-dsRNA	-188.021, -205.339	75.54, 69.589
Nsp14-nsp10-ALT-S6059F-dsRNA	-192.135, -215.547	108.507, 69.443

the HDOCK algorithm, structural changes were observed upon structural alignment for both dsRNA (mean RMSD 6.946 ± 4.232 Å) (Fig. 9A, B) and ssRNA complexes (mean RMSD 5.871 ± 5.708 Å) (Fig. 9C, D).

Since the NSP14 protein is known to be involved in the maintenance and stability of the viral RNA by lowering its sensitivity to RNA mutagens and evading its degradation from the host immune response [28], the large structural changes observed upon mutagenesis at the protein and RNA-protein complex levels, could ultimately lead to a reduced enzyme activity and could be the functional aetiology for observing a greater mutagenesis rate.

Mean free energy values were also reported for the 10 best models attained in order to observe for differences in stability of the NSP14-RNA and the NSP14-NSP10-RNA complexes. HDOCK free energy values are only reported as there were some missing values from the MPRDock algorithm. Results show minor differences in stability between the REF- and ALT-NSP14-RNA complexes for both ds- and ss-RNA molecules examined here (see Table 2). More evident differences in complex stability as denoted by free energy scores are observed for the NSP14-NSP10-RNA complexes, with the ALT forms of the complex attaining overall lower free energy values, indicating a higher affinity for RNA binding for the ALT forms of this complex (see Table 2).

^aTwo values are shown for rows displaying the Nsp14-nsp10 complex and RNA binding results. The first value denotes scores from obtaining the protein complex using ClusPro while the second one using HDOCK

Discussion

Investigating mutations for SARS-CoV-2 that are unique to a specific country can not only lead to interesting results for the underlying country but also provide insights into viral virulence worldwide. Countries like Cyprus with a small isolated population can be invaluable in assessing viral transmission and death rates, as well as potentially predicting future trends of the virus in the

rest of the world. This work was based on the analysis of a small number ($N = 144$) of viral strains from the Cypriot population. To this end we initially perform basic phylogenetic and lineage analysis as previously reported for these samples [10]. We extend previous work by selecting for 18 Cyprus-specific mutations exhibiting a high frequency of occurrence that is unique to Cyprus compared to the rest of the world. We highlight a single mutation with the highest frequency of occurrence in Cyprus alone and further analyse this by mutation tracking and regressions analysis (glms) in order to obtain a greater understanding of the nature of this mutation. We show that this mutation causes an amino acid change (S6059F) on the NSP14 exoribonuclease of the *ORF1ab* protein. We furthermore assess whether this mutation affects the molecular functionality of the NSP14 protein, which is known to be implicated in viral mutation proofreading during replication. We provide evidence that it may actually allude to a dysfunctional NSP14 protein that causes hypermutability in the strains with this mutation. This is further supported by structural modelling of the NSP14 protein with and without the S6059F mutation, which clearly points to a different structural conformation of the ALT vs. the REF form of the protein. This structural variation is exhibited by the NSP14 protein alone as well as in complex with NSP10 and with both ds- and ss-RNA molecules. Moreover, free energy values report greater stability in the ALT-NSP14-NSP10-RNA complex. The mutation-generated alteration in RNA or DNA affinity has also been investigated with other viral proteins that are involved in RNA/DNA processing during viral replication. A recent example comes from the Herpes Simplex Virus UL42 protein which binds DNA and plays an essential role in viral DNA replication by acting as the polymerase accessory subunit. It has been shown that engineered viruses expressing mutant forms of the UL42 proteins that increase its affinity for DNA binding, exhibited increased mutation frequencies and elevated ratios of virion DNA copies [36]. These results suggest the

SARS-CoV-2 S6059F-generated increased affinity for RNA seen in our structural models, may be the cause of the hypermutation also observed in our data for the Cypriot strains with the ALT form of the S6059F mutation. As previously reported, certain viruses may have evolved so that their RNA or DNA binding proteins, implicated in viral replication, neither bind DNA too tightly nor too weakly to optimize virus production and replication [36]. In addition, studies from the first SARS-CoV epidemic has implicated the ExoN activity within NSP14 in a controversial mechanism for how coronaviruses (CoVs) regulate replication fidelity. This new model of CoV replication regulation is characterized by increased proofreading capabilities. At the time, this was contradictory to the widely accepted paradigm that proofreading is predominantly of low-fidelity in RNA viruses. This raises important issues on how the current SARS-CoV-2 RNA-dependent RNA polymerases (RdRps), including NSP14, maintain a delicate balance between viral quasispecies generation and the accumulation of too many deleterious mutations to conserve viral replication [37].

It is important to stress that our proposed mechanism for the S6059F mutation alluding to a dysfunctional NSP14 protein, which in turn may cause hypermutability, requires additional in-vitro experimental verification before it can be fully validated.

Recent research has shown that host genetic factors can affect susceptibility to COVID-19. Notable examples include: 1) DNA polymorphisms in ACE2 and TMPRSS2, which are two key host factors of SARS-CoV-2 [38], 2) variants in HLA genes related to type I interferon immunity that might predispose patients toward life-threatening COVID-19 pneumonia [39]. A current study has linked NSP14 with mechanisms used by SARS-CoV-2 to evade host antiviral responses, specifically those exhibited by type I interferon (IFN-I) [40]. Moreover, inborn errors of immunity specifically of IFN-I, have been shown to affect the risk of severe COVID-19 in patient cohorts [41]. These findings allude to another potential role for the NSP14 mutation found uniquely amongst the Cypriot population. One that might not be related to factors of the SARS-CoV-2 per se but, to the unique genomics of the Cypriots. Concluding evidence from studies by Hsu et al [40] and Zhang et al [41], reveal that the NSP14 evasive action through translational shutdown of IFN-I, can be affected by the actual genetic composition of the underlining population. Since Cyprus is an island with a relatively isolated population, risk factors of severe COVID-19 are most probably affected by the unique genomics of the Cypriots and a likely candidate is the IFN-I gene. It remains to be seen if genetic variants of the IFN-I gene

within the Cypriot population could be the delineating factors that link the NSP14 S6059F mutation to COVID-19 severity.

To complicate things even further, recent research on drug discovery for SARS-CoV-2 has associated structural and function implications of RdRps (NSP12 and NSP14) mutations in relation to resistance against the drug Remdesivir [42]. This is yet another factor that should be taken under consideration for the Cypriot population, given the unique NSP14 S6059F mutation found in high frequency amongst the Cypriot strains.

Conclusions

We propose a theory whereby the increased mutability caused by decreased proofreading as a consequence of the S6059F mutation, causes a greater diversity of viral strains or quasispecies resident within the same host. This can potentially impact the infection patterns descending from this host by altering transmission rates and possibly death rates. We outline two potential scenarios where this could be of significant impact for downstream viral infection trends. If the “main” nested strain is of high virulence, viral replication with increased mutability will allude to a greater range of viral quasispecies from which to infect downstream hosts, thus less chances of the host transmitting the strain with high virulence. If, on the other hand, the nested strain is one of low virulence, increased mutability could potentially lead to the development and consequent transmission of more virulent strains in downstream infections. One very important aspect that ought to be discussed and investigated further, is the potential impact of such viable mutations that cause hypermutability on vaccination efficacy. It can be argued that infections caused from viral strains with such proofreading dysfunction can act as generators of diverse viral strains that allude current vaccination attempts at a faster rate compared to strains with more regular proofreading functionality.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12863-021-01007-9>.

Additional file 1.

Additional file 2.

Acknowledgements

Not applicable.

Authors' contributions

Conceptualization, A. O, J.R, C. C, and G.M.S.; methodology, A. O, M. Z, M.T; software, A. O, M. Z, M.T; validation, A. O, M. Z, M.T; formal analysis, A. O, M. Z, M.T; investigation, A. O, M. Z, M.T; re-sources, J. R, C.C; data curation, J. R, C.C; writing—original draft preparation, A. O, M. Z, M. T, J.R, K. M, K. C, C. C, and G.M.S.; writing—review and editing, A. O, M. Z, M. T, J.R, K. M, K. C, C. C, and G.M.S.; visualization, A. O, M. Z, M.T; supervision, K. M, C.C, and G.M.S.;

project administration, A. O, J.R, C. C, and G.M.S. All authors have read and agreed to the published version of the manuscript.

Funding

Not applicable.

Availability of data and materials

The sequences used and/or analysed during the current study are available from the GISAID (<https://www.gisaid.org/>) EpiCov database under the accession numbers: EPI_ISL_1164626–EPI_ISL_1164751 and EPI_ISL_463741–EPI_ISL_463748.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Bioinformatics Department, Cyprus Institute of Neurology and Genetics, Nicosia, Cyprus. ²The Cyprus School of Molecular Medicine, Nicosia, Cyprus. ³Molecular Virology Department, Cyprus Institute of Neurology and Genetics, Nicosia, Cyprus. ⁴Biostatistics Unit, Cyprus Institute of Neurology and Genetics, Nicosia, Cyprus. ⁵Neurogenetics Department, Cyprus Institute of Neurology and Genetics, Nicosia, Cyprus.

Received: 22 July 2021 Accepted: 26 October 2021

Published online: 13 November 2021

References

- Andino R, Domingo E. Viral quasispecies. *Virology*. 2015;479–480:46.
- Zrelovs N, Ustinova M, Silamikelis I, Birzniece L, Megnis K, Rovite V, et al. First report on the Latvian SARS-CoV-2 isolate genetic diversity. *Front Med*. 2021; 8:241. <https://doi.org/10.3389/fmed.2021.626000>.
- Geoghegan JL, Douglas J, Ren X, Storey M, Hadfield J, Silander OK, et al. The power and limitations of genomics to track COVID-19 outbreaks: a case study from New Zealand. *Emerging Infectious Diseases*. 2020.
- Elizondo V, Harkins GW, Mabvakure B, Smidt S, Zappile P, Marier C, et al. SARS-CoV-2 genomic characterization and clinical manifestation of the COVID-19 outbreak in Uruguay. *Emerg Microbes Infect*. 2021;10(1):51–65. <https://doi.org/10.1080/22221751.2020.1863747>.
- Kozlovskaya L, Piniava A, Ignatyev G, Selivanov A, Shishova A, Kovpak A, et al. Isolation and phylogenetic analysis of SARS-CoV-2 variants collected in Russia during the COVID-19 outbreak. *Int J Infect Dis*. 2020;99:40–6. <https://doi.org/10.1016/j.ijid.2020.07.024>.
- Taboada B, Vazquez-Perez JA, Muñoz-Medina JE, Ramos-Cervantes P, Escalera-Zamudio M, Boukadida C, et al. Genomic analysis of early SARS-CoV-2 variants introduced in Mexico. *J Virol*. 2020;94(18). <https://doi.org/10.1128/JVI.01056-20>.
- Zhang W, Govindavari JP, Davis BD, Chen SS, Kim JT, Song J, et al. Analysis of genomic characteristics and transmission routes of patients with confirmed SARS-CoV-2 in Southern California during the early stage of the US COVID-19 pandemic. *JAMA Netw Open*. 2020;3(10):e2024191. <https://doi.org/10.1001/jamanetworkopen.2020.24191>.
- Gómez-Carballa A, Bello X, Pardo-Seco J, Del Molino MLP, Martínón-Torres F, Salas A. Phylogeography of SARS-CoV-2 pandemic in Spain: a story of multiple introductions, micro-geographic stratification, founder effects, and super-spreaders. *Zool Res*. 2020;41(6):605–20. <https://doi.org/10.24272/j.issn.2095-8137.2020.217>.
- Sekizuka T, Itokawa K, Hashino M, Kawano-Sugaya T, Tanaka R, Yatsu K, et al. A genome epidemiological study of SARS-CoV-2 introduction into Japan. *mSphere*. 2020;5:e00786-20.
- Richter J, Fanis P, Tryfonos C, Koptides D, Krashias G, Bashiards S, et al. Molecular epidemiology of SARS-CoV-2 in Cyprus. *PLoS One*. 2021;16:e0248792. <https://doi.org/10.1371/journal.pone.0248792>.
- Oulas A, Zanti M, Tomazou M, Zachariou M, Minadakis G, Bourdakou MM, et al. Generalized linear models provide a measure of virulence for specific mutations in SARS-cov-2 strains. *PLoS One*. 2021;16(1):e0238665. <https://doi.org/10.1371/journal.pone.0238665>.
- Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol*. 2009;10:R25. <https://doi.org/10.1186/gb-2009-10-3-r25>.
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The sequence alignment/map format and SAMtools. *Bioinformatics*. 2009;25:2078–9. <https://doi.org/10.1093/bioinformatics/btp352>.
- McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al. The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res*. 2010;20(9):1297–303. <https://doi.org/10.1101/gr.107524.110>.
- Rambaut A, Holmes EC, O'Toole Á, Hill V, McCrone JT, Ruis C, et al. A dynamic nomenclature proposal for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol*. 2020;5(11):1403–7. <https://doi.org/10.1038/s41564-020-0770-5>.
- RC T. R: a language and environment for statistical computing. Vienna: Austria R Found Stat Comput; 2013.
- Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al. NextStrain: real-time tracking of pathogen evolution. *Bioinformatics*. 2018;34:4121–3.
- Katoh K. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res*. 2002;30(14):3059–66. <https://doi.org/10.1093/nar/gkf436>.
- Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*. 2014;30(9):1312–3. <https://doi.org/10.1093/bioinformatics/btu033>.
- Yang J, Zhang Y. Protein structure and function prediction using I-TASSER. *Curr Protoc Bioinformatics*. 2015;52(1):5.8.1–5.8.15. <https://doi.org/10.1002/0471250953.bi0508s52>.
- Ittisoponpisan S, Islam SA, Khanna T, Alhuzimi E, David A, Sternberg MJE. Can predicted protein 3D structures provide reliable insights into whether missense variants are disease associated? *J Mol Biol*. 2019;431:2197–212.
- Rodrigues CH, Pires DE, Ascher DB, RenéRen I, Rachou R, Oswaldo CF. DynaMut: predicting the impact of mutations on protein conformation, flexibility and stability. *Nucleic Acids Res*. 2018;46(W1):W350–5. <https://doi.org/10.1093/nar/gky300>.
- Jubb HC, Higuero AP, Ochoa-Montaña B, Pitt WR, Ascher DB, Blundell TL. Arpeggio: a web server for calculating and visualising interatomic interactions in protein structures. *J Mol Biol*. 2017;429(3):365–71. <https://doi.org/10.1016/j.jmb.2016.12.004>.
- Vajda S, Yueh C, Beglov D, Bohnuud T, Mottarella SE, Xia B, et al. New additions to the ClusPro server motivated by CAPRI. *Proteins Struct Funct Bioinforma*. 2017;85:435–44.
- Yan Y, Tao H, He J, Huang SY. The HDock server for integrated protein–protein docking. *Nat Protoc*. 2020;15:1829–52.
- Xue LC, Rodrigues JP, Kastritis PL, Bonvin AM, Vangone A. PRODIGY: a web server for predicting the binding affinity of protein-protein complexes. *Bioinformatics*. 2016;32:3676–8.
- He J, Tao H, Huang SY. Protein-ensemble-RNA docking by efficient consideration of protein flexibility through homology models. *Bioinformatics*. 2019;35:4994–5002.
- Arya R, Kumari S, Pandey B, Mistry H, Bihani SC, Das A, et al. Structural insights into SARS-CoV-2 proteins. *J Mol Biol*. 2021;433(2):166725. <https://doi.org/10.1016/j.jmb.2020.11.024>.
- Rambaut A, Loman N, Pybus O, Barclay W, Barrett J, Carabelli A, et al. Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in the UK defined by a novel set of spike mutations. *VirologicalOrg*. 2020.
- Davies NG, Abbott S, Barnard RC, Jarvis CI, Kucharski AJ, Munday J, et al. Estimated transmissibility and severity of novel SARS-CoV-2 Variant of Concern 202012/01 in England. *Science*. 2020.
- Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Hopkins S, et al. Transmission of SARS-CoV-2 lineage B.1.1.7 in England: insights from linking epidemiological and genetic data. *medRxiv*. 2021;2020.12.30.20249034.
- Ogando NS, Zevenhoven-Dobbe JC, van der Meer Y, Bredenbeek PJ, Posthuma CC, Snijder EJ. The enzymatic activity of the nsp14 exoribonuclease is critical for replication of MERS-CoV and SARS-CoV-2. *J Virol*. 2020;94(23). <https://doi.org/10.1128/JVI.01246-20>.

33. Eckerle LD, Lu X, Sperry SM, Choi L, Denison MR. High fidelity of murine hepatitis virus replication is decreased in nsp14 exoribonuclease mutants. *J Virol.* 2007;81(22):12135–44. <https://doi.org/10.1128/JVI.01296-07>.
34. Eckerle LD, Becker MM, Halpin RA, Li K, Venter E, Lu X, et al. Infidelity of SARS-CoV Nsp14-exonuclease mutant virus replication is revealed by complete genome sequencing. *PLoS Pathog.* 2010;6:1–15.
35. Graham RL, Becker MM, Eckerle LD, Bolles M, Denison MR, Baric RS. A live, impaired-fidelity coronavirus vaccine protects in an aged, immunocompromised mouse model of lethal disease. *Nat Med.* 2012;18: 1820–6.
36. Jiang C, Komazin-Meredith G, Tian W, Coen DM, Hwang CBC. Mutations that increase DNA binding by the processivity factor of herpes simplex virus affect virus production and DNA replication fidelity. *J Virol.* 2009.
37. Smith EC, Denison MR. Coronaviruses as DNA wannabes: a new model for the regulation of RNA virus replication fidelity. *PLoS Pathog.* 2013;9(12): e1003760. <https://doi.org/10.1371/journal.ppat.1003760>.
38. Hou Y, Zhao J, Martin W, Kallianpur A, Chung MK, Jehi L, et al. New insights into genetic susceptibility of COVID-19: an ACE2 and TMPRSS2 polymorphism analysis. *BMC Med.* 2020;18(1):216. <https://doi.org/10.1186/s12916-020-01673-z>.
39. Secolin R, de Araujo TK, Gonsales MC, Rocha CS, Naslavsky M, De Marco L, et al. Genetic variability in COVID-19-related genes in the Brazilian population. *Hum Genome Var.* 2021;8(1):15. <https://doi.org/10.1038/s41439-021-00146-w>.
40. Hsu JCC, Laurent-Rolle M, Pawlak JB, Wilen CB, Cresswell P. Translational shutdown and evasion of the innate immune response by SARS-CoV-2 NSP14 protein. *Proc Natl Acad Sci U S A.* 2021;118(24):e2101161118. <https://doi.org/10.1073/pnas.2101161118>.
41. Zhang Q, Liu Z, Moncada-Velez M, Chen J, Ogishi M, Bigio B, et al. Inborn errors of type I IFN immunity in patients with life-threatening COVID-19. *Science* (80-). 2020.
42. Shannon A, Le NTT, Selisko B, Eydoux C, Alvarez K, Guillemot JC, et al. Remdesivir and SARS-CoV-2: structural requirements at both nsp12 RdRp and nsp14 exonuclease active-sites. *Antivir Res.* 2020;178:104793. <https://doi.org/10.1016/j.antiviral.2020.104793>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

