# The Ensembl COVID-19 resource: ongoing integration of public SARS-CoV-2 data

**Nishadi H. De Silva** [ID], **Jyothish Bhai, Marc Chakiachvili, Bruno Contreras-Moreira,
Carla Cummins, Adam Frankish** [ID]**, Astrid Gall, Thiago Genez, Kevin L. Howe** [ID]**,
Sarah E. Hunt** [ID]**, Fergal J. Martin** [ID]**, Benjamin Moore, Denye Ogeh, Anne Parker,
Andrew Parton, Magali Ruffier** [ID]**, Manoj Pandian Sakthivel, Dan Sheppard, John Tate,
Anja Thormann, David Thybert, Stephen J. Trevanion, Andrea Winterbottom,
Daniel R. Zerbino** [ID]**, Robert D. Finn** [ID]**, Paul Flicek** [ID] **and Andrew D. Yates** [ID]*****

European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK

## ABSTRACT

**The COVID-19 pandemic has seen unprecedented use of SARS-CoV-2 genome sequencing for epidemiological tracking and identification of emerging variants. Understanding the potential impact of these variants on the infectivity of the virus and the efficacy of emerging therapeutics and vaccines has become a cornerstone of the fight against the disease. To support the maximal use of genomic information for SARS-CoV-2 research, we launched the Ensembl COVID-19 browser; the first virus to be encompassed within the Ensembl platform. This resource incorporates a new Ensembl gene set, multiple variant sets, and annotation from several relevant resources aligned to the reference SARS-CoV-2 assembly. Since the first release in May 2020, the content has been regularly updated using our new rapid release workflow, and tools such as the Ensembl Variant Effect Predictor have been integrated. The Ensembl COVID-19 browser is freely available at https://covid-19.ensembl.org.**

## INTRODUCTION

Over the past 20 years, multiple zoonotic respiratory diseases caused by coronaviruses have been identified. Examples include the SARS epidemic caused by severe acute respiratory syndrome coronavirus (SARS-CoV) in 2003 and the Middle East respiratory syndrome coronavirus (MERS-CoV) outbreak in 2012. Both belong to the *betacoronavirus* genus and are believed to have originated in bats with an intermediary animal host before transmission to humans (1).

Similarly, the SARS-CoV-2 virus responsible for the current COVID-19 pandemic is a *betacoronavirus,* with a 29 903-nucleotide positive-strand RNA genome encoding ∼30 known and hypothetical mature proteins. The first open reading frame (ORF), representing ∼67% of the entire genome, encodes 16 non-structural proteins (nsps). The remaining ORFs encode accessory proteins and four major structural proteins: spike surface glycoprotein (S), small envelope protein (E), matrix protein (M) and nucleocapsid protein (N).

Genomic sequencing has played a crucial role in understanding the mechanisms, spread and evolution of this virus, and the number of SARS-CoV-2 genomes sequenced has grown steadily. In the UK, for instance, close to 5% of all reported infections each week were being sequenced in January 2021 (COG-UK, January 2021: https://www.cogconsortium.uk/wp-content/uploads/2021/02/COG-UK-geo-coverage_2021--02-01_summary.pdf) but by April this figure had risen to 46% (https://www.cogconsortium.uk/wp-content/uploads/2021/04/COG-UK-geo-coverage_2021--04-05_summary.pdf). Established genomic resources can bring these data to new and existing user communities supporting efforts to combat the COVID-19 pandemic. Both the UCSC SARS-CoV-2 browser (2) and the WashU Virus browser (3) are examples of this, offering unique data sets (including immunological annotation and comparative alignments) and visualisations over the SARS-CoV-2 genome.

Ensembl (4,5) was launched to capture data from the Human Genome Project and has since developed into a large scale system for generating, integrating and disseminating reference genomes and annotation. The COVID-19 pandemic presented novel challenges related to presenting SARS-CoV-2 viral annotation within Ensembl. Meeting these, we launched the Ensembl COVID-19 browser

---

*To whom correspondence should be addressed. Tel: +44 1223 492538; Email: ayates@ebi.ac.uk

(https://covid-19.ensembl.org) in May 2020. In this manuscript, we detail its features and development that uses rapid Ensembl update cycles to react in the face of potential future outbreaks.

## MATERIALS AND METHODS

### Reference assembly and a new gene annotation

The SARS-CoV-2 sequence represented in Ensembl (INSDC accession GCA_009858895.3, MN908947.3) is the viral RNA genome isolated from one of the first cases in Wuhan, China (6). It is widely used as the standard reference and has been incorporated into other resources including the UCSC SARS-CoV-2 genome browser. This assembly was imported from the European Nucleotide Archive (ENA) into an Ensembl database schema to facilitate analysis and display.

To enable the correct annotation of SARS-CoV-2, our gene annotation methods (7) were adapted to reflect the biology of the virus. To identify protein coding genes, we aligned SARS-CoV-2 proteins from RefSeq (8) to the genome using Exonerate (9). A challenge for annotation is that the first and largest ORF can result in either non-structural proteins nsp1-11 (ORF1a) or in nsp1-nsp10 and nsp12-nsp16 (ORF1ab) via an internal programmed translational frameshift (10). Exonerate handles this ribosomal slippage by inserting a gap in the alignment and thus allowing the annotation of the full ORF1ab locus. Our modified annotation methodology then removes the artificial gap to represent the slippage frameshift as an RNA edit and ensures a biologically accurate representation of the locus and product.

Our annotation approach was tested on 90 additional SARS-CoV-2 assemblies retrieved from the ENA. We assessed alignment coverage and percentage identity of the resultant gene translations to verify accuracy and consistency. In all cases, full length alignments were observed and average amino acid percentage identity across all genes in most assemblies were 99.9% or 100% (one assembly had 99.81% identity). These results demonstrate that our approach scales consistently to larger volumes of viral data.

In addition to our new gene annotation, we also provide the gene set submitted to the International Nucleotide Sequence Database Collaboration (INSDC) by the Shanghai Public Health Clinical Centre as a separate genome track as shown in Figure 1. The submitted gene annotation can be accessed via our 'Configure this page' option and under the 'Genes and transcripts' heading.

### Comparison of SARS-CoV-2 with 60 other *Orthocoronavirinae* genomes

We used Cactus (11) to align SARS-CoV-2 and 60 publicly available virus genomes from the *Orthocoronavirinae* subfamily resulting in 78% of the SARS-CoV-2 genome aligned with at least one other genome and 35% of the genome aligned with the complete set of *Orthocoronavirinae* genomes. The multiple sequence alignment gives evolutionary context for each region of the genome and is a powerful method to explore functionality such as gene annotation validation. For instance, multiple sequence alignments for

44 complete *Sarbecovirus* genomes have suggested a potentially novel alternate frame gene ORF3c and that ORF10, ORF9c, ORF3b and ORF3d are unlikely to be protein coding (12).

The alignment coverage (Figure 2) represents the number of genomes aligned to a given reference genomic position and is distributed heterogeneously across the SARS-CoV-2 genome. An immediate observation is that the central region of the genome (starting from ∼7.1 kb and ending at 21.3 kb), including a significant segment of the 3′ part of ORF1a, is highly shared across the *Orthocoronavirinae* subfamily. This indicates that the non-structural proteins encoded by this region (nsp3–nsp16) likely originate from the *Orthocoronavirinae* ancestral genome. Conversely, both ends of the SARS-CoV-2 genome have very low alignment coverage and are only shared with closely related viruses.

As a further demonstration of the utility of the alignment coverage, we focused in on the genomic region encoding for the SARS-CoV-2 spike protein. The spike protein has two subunits: S1, which binds to the host cell receptor angiotensin-converting enzyme 2 (ACE2), and S2, which is involved in membrane fusion. The region of the S ORF encoding for the S2 subunit of the spike protein clearly displays a high alignment coverage while the region encoding for the S1 subunit has large portions that are shared only by one other related genome as previously reported (13). This demonstrates our comparative methods—utilized on viruses for the first time—were capable of recapitulating this prior discovery.

Additionally, we applied our gene tree method (14) to group the protein coding genes into families and to predict orthologous and paralogous relationships between genes. These updates will be incorporated into the Ensembl COVID-19 resource by Q4 of 2021.

### Genetic variation data

Analysis of variants of viral genomes is important for understanding the spread of infection across different geographic regions. We display 6134 sequence variants for SARS-CoV-2 and show their regional frequency distributions alongside predicted molecular consequences calculated by the Ensembl Variant Effect Predictor (VEP) (15). The variants on our site are derived from overlapping sample sets produced by groups who used different analysis methods and a small collection of variants of special interest.

The first set of variants we integrated were from the Nextstrain project which creates phylogenetic trees for tracking pathogen evolution based on virus subsamples (16). We converted their SARS-CoV-2 data release from 08-04-2020 to VCF for integration into our system and display frequency distributions by country and Nextstrain-inferred clade.

The second variant set comes from the ENA team, who developed a LoFreq-based (17) analysis to call variants from SARS-CoV-2 sequence data sets submitted to their archives. LoFreq reports the proportion of each variant seen in a sample from an individual. For simplicity, we represent only the alleles seen in each sample and not the proportions estimated. Variants were called for each host sam-
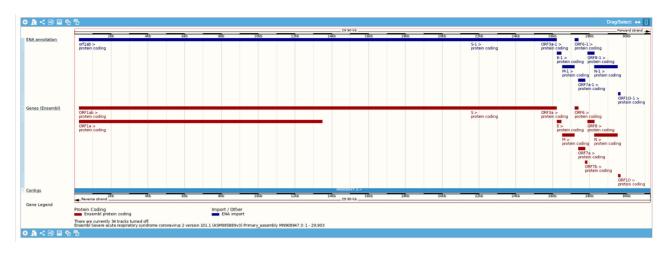
**Figure 1.** A comparison of the Ensembl gene set (displayed in red) and the gene set submitted to INSDC by the Shanghai Public Health Clinical Centre (displayed in blue) for the entire SARS-CoV-2 reference assembly. A notable difference between the two gene sets is the absence of ORF1a and ORF7b in the submitted gene set. Annotation tracks can be configured by clicking on the cog icon displayed in the top left of the figure.
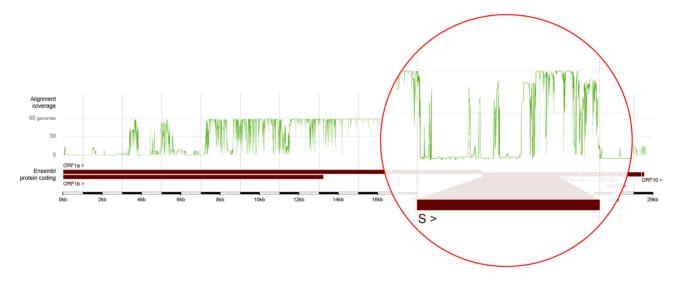


**Figure 2.** Alignment coverage across the SARS-CoV-2 reference genome based on a whole genome multiple sequence alignment with 60 other *Orthocoronavirinae* genomes. The green plot of alignment coverage shows that the central region of the genome is highly shared across the subfamily, while the ends are generally shared only with closely related viruses. The region encoding for the spike protein S has been enlarged within the red circle showing the difference between the low alignment coverage of the upstream S1 subunit (left) and the high coverage of the downstream S2 subunit (right). This demonstrates that our methods were able to reproduce the same observations made by other groups - that there is little conservation in S1 in the *Orthocoronavirinae* subfamily compared to S2.

ple individually and, to provide a more accurate estimation of the frequency of each allele across the entire sample set, it is assumed that sites at which a variant was not called in a sample match the reference genome used in the Ensembl COVID-19 browser. We currently display ENA's variant data from 17-08-2020 and have applied strict filters to reduce the proportion of lower confidence sites. Specifically, we have not included variants from sequence data sets with >40 calls and we have removed variants where no sample has a frequency of 20% or more for the non-reference allele and variants where all samples show strand bias.

Some sites are annotated as a further guide to quality. For example, variants seen in more than one sample in either set have an evidence status of 'Multiple observations' and variants at sites recommended for masking by De Maio *et*

*al.* (18) have a flag of 'Suspect reference location'. Variants can be displayed as separate tracks in the genome browser: those from ENA, those from Nextstrain and those observed in more than one sample in either project as shown in Figure 3.

In December 2020, we also incorporated a set of variants which were reported as a tracking priority by the COVID-19 Genomics UK Consortium at the time (COG-UK, https://www.cogconsortium.uk/). This included 17 variants from the rapidly spreading B.1.1.7 strain (https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563) and four variants from the mink associated strain. The D614G, A222V and N439K mutations associated with
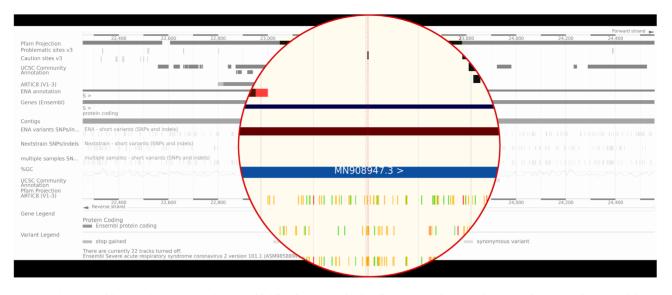
**Figure 3.** The browser with several tracks turned on and highlighting a substitution flagged early in UCSC's community annotation at position 23403 (D614G) in the S spike glycoprotein gene. Due to the prompt nature of community driven annotation, this was available on our browser as soon as the annotation was published as a pre-print. It is labelled as a common missense mutation in SARS-CoV-2 with a notably high difference in resulting isoelectric point (D→G). Studies have shown this missense mutation in the spike protein is predominantly observed in Europe (26); patterns that were also observed in the variation data we host when first imported.

an effect on transmissibility, a fast growing lineage and increased binding affinity to the ACE2 receptor (19–21), respectively, were also included. We extracted the gene and protein change information from the reports and used the Ensembl VEP to map these descriptions to genomic coordinates and create a VCF file, which was then loaded into the Ensembl database with associated phenotype information. We use the NCBI SPDI (22) convention to name all variants for consistency while also enabling searching using the popularly used names, such as 'N501Y'. While the other data sets are limited to single nucleotide substitutions, the COG-UK set also includes deletions and can be viewed as a separate 'COG-UK priority mutations' track alongside the gene annotations.

### Alignment of data from other resources

To enrich the SARS-CoV-2 genome annotation we aligned and integrated data from several external repositories in a similar manner to other genomes available in Ensembl.

Specifically, we aligned Rfam (23) covariance models using their COVID-19 release 14.2 (http://rfam.xfam.org/covid-19) to highlight conserved non-coding RNA structures which are responsible for various stages of the viral life cycle. These include the frame shifting stimulation element and the pseudoknot necessary for the genome replication of SARS-CoV-2 (24). We also provide cross references to proteins from RefSeq, UniProt (25) and the INSDC; functional annotation from the Gene Ontology Consortium; and annotation of protein domains using InterProScan. These additional annotations are accessible via our region views and the gene and transcript tabs. We also created a genome browser track projecting the protein-domain annotations onto the genome to facilitate a genome-oriented view of the gene products including the non-structural cleavage products of ORF1a/ORF1ab.

The browser also displays community annotation of sites and regions using results co-ordinated by the UCSC genome browser. Additions to this annotation resource are open to all and done via a publicly available spreadsheet hosted by UCSC (http://bit.ly/cov2annots), the data from which is integrated periodically into the Ensembl browser. This is achieved via specialised code that uses Git workflows to convert the annotations into BigBed files that can be visualised on a variety of genome browsers (https://github.com/Ensembl/sarscov2-annotation).

We have also integrated Oxford Nanopore sequencing primers (version 3) made available by the ARTIC network (https://artic.network/ncov-2019) to assist in sequencing the virus. Though mainly focused on the Oxford Nanopore MinION sequencer, some aspects of the protocol may be generalised to other sequencing platforms. The complete list of primers included is available on GitHub (https://github.com/artic-network/artic-ncov2019/blob/master/primer_schemes/nCoV-2019/V3/nCoV-2019.tsv).

Finally, we provide tracks to visualise problematic and caution sites, which result from common systematic errors associated with laboratory protocols and have been observed in submitted sequences (18). Inclusion of these can adversely influence phylogenetic and evolutionary inference. Visualising these in the browser alongside the locations of primers and other community derived annotations helps determine how best to proceed with analysis of each these sites.

### Integration and dissemination

The Ensembl COVID-19 resource features a newly designed landing page, which prioritises key views and data to help direct researchers into relevant sections of the site. To support expeditious data release, we have not made potentially time-consuming virus-specific modifications to our existing

web codebase—such as showing a single nucleic acid strand and removing all mentions of exons—because we felt the data could be effectively understood without these changes. However, we have altered the vocabulary wherever possible and are reviewing feedback as we receive it.

Our COVID-19 resource is also integrated into the European COVID-19 Data Portal hosted by EMBL-EBI (https://www.covid19dataportal.org/) (27). The portal enables searches across the multiple research outputs on COVID-19 including viral and human sequences; relevant biochemical pathways, interactions, complexes, targets and compounds; protein and expression data; and literature. Our hosted filtered variants from ENA flow directly to this COVID-19 portal.

We have engaged our existing and new user communities using our blog and social media accounts to announce the release and updates to the Ensembl COVID-19 resource. We also highlighted the changes made to our gene annotation method to ensure the complete set of ORFs as these have been overlooked by other annotation tools.

## RESULTS AND DISCUSSION

The swift spread of COVID-19 has highlighted the necessity for data resources to be prepared for rapid adaptation to developing outbreaks. Our development and release of the Ensembl COVID-19 resource leveraged our experience integrating thousands of genomes into the Ensembl infrastructure and supporting hundreds of thousands of active users. The Ensembl COVID-19 browser provides a unique view on SARS-CoV-2 using our gene annotation method and variation data processed to focus on the highest confidence variants. Additionally, the Ensembl VEP and haplotype views enable the consequences of the variants to be assessed within the context of specific strains and geographical locations. The data is made accessible via the widely used Ensembl platform making it immediately familiar to a large userbase who may be able to repurpose existing software and browser knowledge to support their work during the pandemic and beyond.

When the COVID-19 pandemic hit, we had been working for several months to develop Ensembl Rapid Release (https://rapid.ensembl.org) to distribute annotated genomes within days of their annotation being completed. This experience proved useful in bringing the COVID-19 site to public release quickly. We have also demonstrated the flexibility of the Ensembl infrastructure and its value as a platform for research and discovery. All of our pipelines and schemas worked seamlessly, even though Ensembl was not designed to support RNA genomes or used with viral genomes. The adapted gene annotation method, for instance, produced consistent annotation with ribosomal slippage correctly modelled and can be reused in the future. Similarly, our gene tree and alignment methods have been applied to the viral data with only minimal changes to parameters. We will continue to regularly update the site as new data emerges such that the Ensembl COVID-19 resource supports research into understanding the genomic evolution of this virus, identifying hotspots of genomic variation and enabling the rational design of future therapeu-

tics, vaccines and policies well beyond the end of the current pandemic.

## DATA AVAILABILITY

The new Ensembl COVID-19 resource is available without restrictions at https://covid-19.ensembl.org. The reference genome assembly for SARS-CoV-2 with the accession GCA_009858895.3 was obtained from the European Nucleotide Archive (https://www.ebi.ac.uk/ena/browser/view/GCA_009858895.3). Ensembl code is available under an Apache 2.0 license at https://github.com/Ensembl.

## REFERENCES

1. Wu,A., Peng,Y., Huang,B., Ding,X., Wang,X., Niu,P., Meng,J., Zhu,Z., Zhang,Z., Wang,J. *et al.* (2020) Genome composition and divergence of the novel coronavirus (2019-nCoV) originating in China. *Cell Host Microbe*, **27**, 325–328.
2. Fernandes,J.D., Hinrichs,A.S., Clawson,H., Gonzalez,J.N., Lee,B.T., Nassar,L.R., Raney,B.J., Rosenbloom,K.R., Nerli,S., Rao,A.A. *et al.* (2020) The UCSC SARS-CoV-2 Genome Browser. *Nat. Genet.*, **52**, 991–998.
3. Flynn,J.A., Purushotham,D., Choudhary,M.N.K., Zhuo,X., Fan,C., Matt,G., Li,D. and Wang,T. (2020) Exploring the coronavirus pandemic with the WashU Virus Genome Browser. *Nat. Genet.*, **52**, 986–991.
4. Howe,K.L., Contreras-Moreira,B., De Silva,N., Maslen,G., Akanni,W., Allen,J., Alvarez-Jarreta,J., Barba,M., Bolser,D.M., Cambell,L. *et al.* (2020) Ensembl Genomes 2020—enabling non-vertebrate genomic research. *Nucleic. Acids. Res.*, **48**, D689–D695.
5. Howe,K.L., Achuthan,P., Allen,J., Allen,J., Alvarez-Jarreta,J., Amode,M.R., Armean,I.M., Azov,A.G., Bennett,R., Bhai,J. *et al.* (2021) Ensembl 2021. *Nucleic Acids Res.*, **49**, D884–D891.
6. Wu,F., Zhao,S., Yu,B., Chen,Y.-M., Wang,W., Song,Z.-G., Hu,Y., Tao,Z.-W., Tian,J.-H., Pei,Y.-Y. *et al.* (2020) A new coronavirus associated with human respiratory disease in China. *Nature*, **579**, 265–269.
7. Aken,B.L., Ayling,S., Barrell,D., Clarke,L., Curwen,V., Fairley,S., Fernandez Banet,J., Billis,K., García Girón,C., Hourlier,T. *et al.* (2016) The Ensembl gene annotation system. *Database*, **2016**, baw093.
8. O'Leary,N.A., Wright,M.W., Brister,J.R., Ciufo,S., Haddad,D., McVeigh,R., Rajput,B., Robbertse,B., Smith-White,B., Ako-Adjei,D. *et al.* (2016) Reference sequence (RefSeq) database at NCBI: current

status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.*, **44**, D733–D745.

9. Slater,G.S.C. and Birney,E. (2005) Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics*, **6**, 31.

10. Chen,Y., Liu,Q. and Guo,D. (2020) Emerging coronaviruses: genome structure, replication, and pathogenesis. *J. Med. Virol.*, **92**, 418–423.

11. Armstrong,J., Hickey,G., Diekhans,M., Fiddes,I.T., Novak,A.M., Deran,A., Fang,Q., Xie,D., Feng,S., Stiller,J. *et al.* (2020) Progressive cactus is a multiple-genome aligner for the thousand-genome era. *Nature*, **587**, 246–251.

12. Jungreis,I., Sealfon,R. and Kellis,M. (2021) SARS-CoV-2 gene content and COVID-19 mutation impact by comparing 44 Sarbecovirus genomes. *Nat. Commun.*, **12**, 2642.

13. Lokman,S.M., Rasheduzzaman,M., Salauddin,A., Barua,R., Tanzina,A.Y., Rumi,M.H., Hossain,M.I., Siddiki,A.M.A.M.Z., Mannan,A. and Hasan,M.M. (2020) Exploring the genomic and proteomic variations of SARS-CoV-2 spike glycoprotein: a computational biology approach. *Infect. Genet. Evol.*, **84**, 104389.

14. Herrero,J., Muffato,M., Beal,K., Fitzgerald,S., Gordon,L., Pignatelli,M., Vilella,A.J., Searle,S.M.J., Amode,R., Brent,S. *et al.* (2016) Ensembl comparative genomics resources. *Database*, **2016**, bav096.

15. McLaren,W., Gil,L., Hunt,S.E., Riat,H.S., Ritchie,G.R.S., Thormann,A., Flicek,P. and Cunningham,F. (2016) The Ensembl variant effect predictor. *Genome Biol.*, **17**, 122.

16. Hadfield,J., Megill,C., Bell,S.M., Huddleston,J., Potter,B., Callender,C., Sagulenko,P., Bedford,T. and Neher,R.A. (2018) Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics*, **34**, 4121–4123.

17. Wilm,A., Aw,P.P.K., Bertrand,D., Yeo,G.H.T., Ong,S.H., Wong,C.H., Khor,C.C., Petric,R., Hibberd,M.L. and Nagarajan,N. (2012) LoFreq: a sequence-quality aware, ultra-sensitive variant caller for uncovering cell-population heterogeneity from high-throughput sequencing datasets. *Nucleic Acids Res.*, **40**, 11189–11201.

18. De Maio,N., Walker,C., Borges,R., Weilguny,L., Slodkowicz,G. and Goldman,N. (2020) In: *Issues with SARS-CoV-2 sequencing data*.

19. Volz,E., Hill,V., McCrone,J.T., Price,A., Jorgensen,D., O'Toole,Á., Southgate,J., Johnson,R., Jackson,B., Nascimento,F.F. *et al.* (2021) Evaluating the Effects of SARS-CoV-2 Spike Mutation D614G on Transmissibility and Pathogenicity. *Cell*, **184**, 64–75.

20. Bartolini,B., Rueca,M., Gruber,C., Messina,F., Giombini,E., Ippolito,G., Capobianchi,M. and Di Caro,A. (2020) The newly introduced SARS-CoV-2 variant A222V is rapidly spreading in Lazio region, Italy. medRxiv doi: https://doi.org/10.1101/2020.11.28.20237016, 30 November 2020, preprint: not peer reviewed.

21. Thomson,E., Rosen,L., Shepherd,J., Spreafico,R., da Silva Filipe,A., Wojcechowskyj,J., Davis,C., Piccoli,L., Pascall,D., Dillen,J. *et al.* (2020) The circulating SARS-CoV-2 spike variant N439K maintains fitness while evading antibody-mediated immunity. *Cell*, **184**, 1171–1187.

22. Holmes,J.B., Moyer,E., Phan,L., Maglott,D. and Kattman,B. (2020) SPDI: data model for variants and applications at NCBI. *Bioinformatics*, **36**, 1902–1907.

23. Kalvari,I., Nawrocki,E.P., Ontiveros-Palacios,N., Argasinska,J., Lamkiewicz,K., Marz,M., Griffiths-Jones,S., Toffano-Nioche,C., Gautheret,D., Weinberg,Z. *et al.* (2021) Rfam 14: expanded coverage of metagenomic, viral and microRNA families. *Nucleic Acids Res.*, **49**, D192–D200.

24. Williams,G.D., Chang,R.-Y. and Brian,D.A. (1999) A phylogenetically conserved hairpin-type 3′ untranslated region pseudoknot functions in coronavirus RNA replication. *J. Virol.*, **73**, 8349–8355.

25. Bateman,A., Martin,M.-J., Orchard,S., Magrane,M., Agivetova,R., Ahmad,S., Alpi,E., Bowler-Barnett,E.H., Britto,R., Bursteinas,B. *et al.* (2021) UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res.*, **49**, D480–D489.

26. Isabel,S., Graña-Miraglia,L., Gutierrez,J.M., Bundalovic-Torma,C., Groves,H.E., Isabel,M.R., Eshaghi,A., Patel,S.N., Gubbay,J.B., Poutanen,T. *et al.* (2020) Evolutionary and structural analyses of SARS-CoV-2 D614G spike protein mutation now documented worldwide. *Sci. Rep.*, **10**, 14031.

27. Harrison,P.W., Lopez,R., Rahman,N., Allen,S.G., Aslam,R., Buso,N., Cummins,C., Fathy,Y., Felix,E., Glont,M. *et al.* (2021) The COVID-19 Data Portal: accelerating SARS-CoV-2 and COVID-19 research through rapid open access data sharing. *Nucleic Acids Res.*, **49**, W619–W623.