

Large-scale analysis of SARS-CoV-2 spike-glycoprotein mutants demonstrates the need for continuous screening of virus isolates

Barbara Schrörs¹, Ranganath Gudimella¹⁺, Thomas Bukur¹⁺, Thomas Rösler¹, Martin Löwer^{1*} & Ugur Sahin^{1,2*}

¹ TRON gGmbH - Translationale Onkologie an der Universitätsmedizin der Johannes Gutenberg-Universität Mainz gemeinnützige GmbH, Freiligrathstraße 12, 55131 Mainz, Germany

² BioNTech SE, An der Goldgrube 12, 55131 Mainz, Germany

+ contributed equally

* Correspondence:

Martin Löwer, martin.loewer@tron-mainz.de

Ugur Sahin, sahin@uni-mainz.de

Abstract

Due to the widespread of the COVID-19 pandemic, the SARS-CoV-2 genome is evolving in diverse human populations. Several studies already reported different strains and an increase in the mutation rate. Particularly, mutations in SARS-CoV-2 spike-glycoprotein are of great interest as it mediates infection in human and recently approved mRNA vaccines are designed to induce immune responses against it.

We analyzed 146,920 SARS-CoV-2 genome assemblies and 2,393 NGS datasets from GISAID, NCBI Virus and NCBI SRA archives focusing on non-synonymous mutations in the spike protein. Only around 13.6% of the samples contained the wild-type spike protein with no variation from the reference. Among the spike protein mutants, we confirmed a low mutation rate exhibiting less than 10 non-synonymous mutations in 99.98% of the analyzed sequences, but the mean and median number of spike protein mutations per sample increased over time. 2,592 distinct variants were found in total. The majority of the observed variants were recurrent, but only nine and 23 recurrent variants were found in at least 0.5% of the mutant genome assemblies and NGS samples, respectively. Further, we found high-confidence subclonal variants in about 15.1% of the NGS data sets with mutant spike protein, which might indicate co-infection with various SARS-CoV-2 strains and/or intra-host evolution. Lastly, some variants might have an effect on antibody binding or T-cell recognition.

These findings demonstrate the increasing importance of monitoring SARS-CoV-2 sequences for an early detection of variants that require adaptations in preventive and therapeutic strategies.

Introduction

Since the first report of the severe acute respiratory syndrome coronavirus-2 (SARS-CoV-2) outbreak (Zhu et al. 2020; Wu et al. 2020), it has transformed into a global pandemic infecting and threatening death for millions of people all over the globe. By January 20, 2021, the World Health Organization (WHO) reported 94,124,612 confirmed cases and 2,034,527 deaths caused by the SARS-CoV-2 outbreak (World Health Organization 2021). On verge of the approval of SARS-CoV-2 vaccines which are designed to invoke immune responses against the spike-glycoprotein (spike protein), it becomes necessary to track the mutations in spike protein and study their relevance for current and upcoming vaccines. Also the recently approved neutralizing antibody bamlanivimab targets the spike protein of SARS-CoV-2 (Mahase 2020).

Subunits of the spike protein are valuable targets for vaccine design as the protein is responsible for viral binding and entry to host cells (Tai et al. 2020; Shang et al. 2020). The spike protein consists of the N-terminal S1 and the C-terminal S2 subunits; the receptor-binding domain (RBD) in the S1 subunit binds to a receptor on the host cell surface and the S2 subunit fuses viral and host membranes (Li 2016). The receptor binding domain (RBD) of the SARS-CoV-2 spike protein recognizes human angiotensin-converting enzyme 2 (ACE2) as its entry receptor, similar to SARS-CoV (Zhou et al. 2020a). Interacting residues of the SARS-CoV-2 RBD with human ACE2 are highly conserved or share similar side chain properties with the SARS-CoV RBD (Lan et al. 2020). In addition, the SARS-CoV-2 RBD shows significantly higher binding affinity to ACE2 receptor compared to the SARS-CoV RBD. In order to repress the infection, blocking the RBD binding was effective in ACE2-expressing cells (Tai et al. 2020). Among the interacting sites in the SARS-CoV-2 RBD, particularly the amino acid residues L455, F486, Q493, S494, N501, and Y505 provide critical interactions with human ACE2 (Wan et al. 2020). These interacting residues vary due to natural selection in SARS-CoV-2 and other related coronaviruses (Tang et al. 2020). Similarly, worldwide SARS-CoV-2 genomic data shows ten RBD mutations which were caused due to natural selection by circulating among the human population (Ou et al. 2020). RBD mutations particularly at N501 may enhance the binding affinity between SARS-CoV-2 and human ACE2 significantly, improving viral infectivity and pathogenicity (Wan et al. 2020).

It is reported that continuous evolution of SARS-CoV-2 among the global population results into six major subtypes which involve the recurrent D614G mutation of the spike protein (Morais et al. 2020). Further, spread of such recurrent mutations within sub-populations might affect the severity of disease emergence and change the trajectory of the pandemic. Studies also report high intra-host diversity caused by low frequency subclonal mutations within a specific cohort (Kuipers et al. 2020). It is evident that changes in the SARS-CoV-2 genome over time might show new mutations which might influence the development efforts of of interventional strategies. The variability of epitopes of the RBD might hamper the development and use of neutralizing antibodies for cross-protective activities against mutant strains (Sun et al. 2020). Mutational variants of the spike protein might as well lead to escape variants with respect to pre-existing cross-reactive CD4+ T cell responses (Braun et al. 2020) or long-term protection from re-infection through T cell memory. Hence, there is a necessity of constant monitoring of the rapidly changing mutation rates in the spike protein in SARS-CoV-2, which could have significant impact on virus infection, transmissibility and pathogenicity in the current pandemic.

In this study, we gathered 147,413 genomic assemblies and 2,393 NGS sequencing datasets to detect non-synonymous spike protein mutations and infer their frequency within a given sample and the effect on potential antibody binding sites and known T cell epitopes.

Methods

SARS-CoV-2 assemblies

SARS-CoV-2 assemblies from human hosts were downloaded on October 2nd, 2020 from US National Center for Biotechnology Information (NCBI) Virus (protein sequences; Hatcher et al. 2017) and on October 2nd, 2020 from GISAID (nucleotide sequences; Elbe and Buckland-Merrett 2017). Pairwise alignments to the reference surface glycoprotein (NC_045512.2_cds_YP_009724390.1_3) were performed to extract the S gene sequences from GISAID samples using the R package Biostrings (version 2.52.0). Extracted sequences were translated with option `if.fuzzy.codon = "solve"`. Amino acid sequences of less than 100 length (440 samples) or premature stop codons (53 samples) were excluded from further analyses.

Non-synonymous variants were determined by pairwise alignment (Biostrings, version 2.52.0) of the protein sequences to the translated reference sequence.

NGS data processing

All available NGS data for SARS-CoV-2 was downloaded on October 14th, 2020 from the NCBI Sequence Read Archive (SRA) (<https://www.ncbi.nlm.nih.gov/genbank/sars-cov-2-seqs/>; Leinonen et al. 2011) and filtered for whole genome fastq data from Illumina instruments with a human sample background. Data were aligned to the reference MN908947.3 (Brister et al. 2015).

Short-read whole genome sequencing data were aligned with `bwa` (version 0.7.17) `mem` (Li 2013). Output files in SAM format were sorted and converted to their binary form (BAM) using `SAMtools` (version 0.1.16) (Li et al. 2009). Variants were retrieved from the alignment files using `BCFtools` (version 1.9) `mpileup` (<http://samtools.github.io/bcftools/>) with the options to recalculate per-base alignment quality on the fly, disabling the maximum per-file depth, and retention of anomalous read pairs. Variants in gene `gp02` (i.e. S gene) were annotated using `SNPEff` (version 4.3t) “`ann`” (Cingolani et al. 2012).

Filtering subclonal variants

NGS variants were filtered with at least 30 reads coverage and a fraction of supporting reads of at least 0.1 and less than 0.95 to identify high-confidence sub-clonal mutations (Sashittal et al. 2020).

Calculation of solvent-accessible residues and corresponding solvent-accessible surface areas

Solvent-accessible residues of the spike protein were calculated using the rolling ball algorithm of the Swiss PDB Viewer (version 4.1.0; Guex and Peitsch 1997) with a parameter setting of $\geq 30\%$ accessible surface.

Solvent-accessible surface area (SASA) was calculated with tools from PyRosetta (version PyRosetta-4 2019) with default settings on reference `pdb-structure “6vxx”` for the spike protein (from PDB-Protein-Databank). SASA was calculated for every residue (in triplicates by the trimeric structure of the spike protein). The mutated structures were generated by introducing single mutations into the reference structure by tools from PyRosetta, too. This included merely a repacking of side-chains locally around the mutation side (with radius 3 Å), leaving the backbone unaltered.

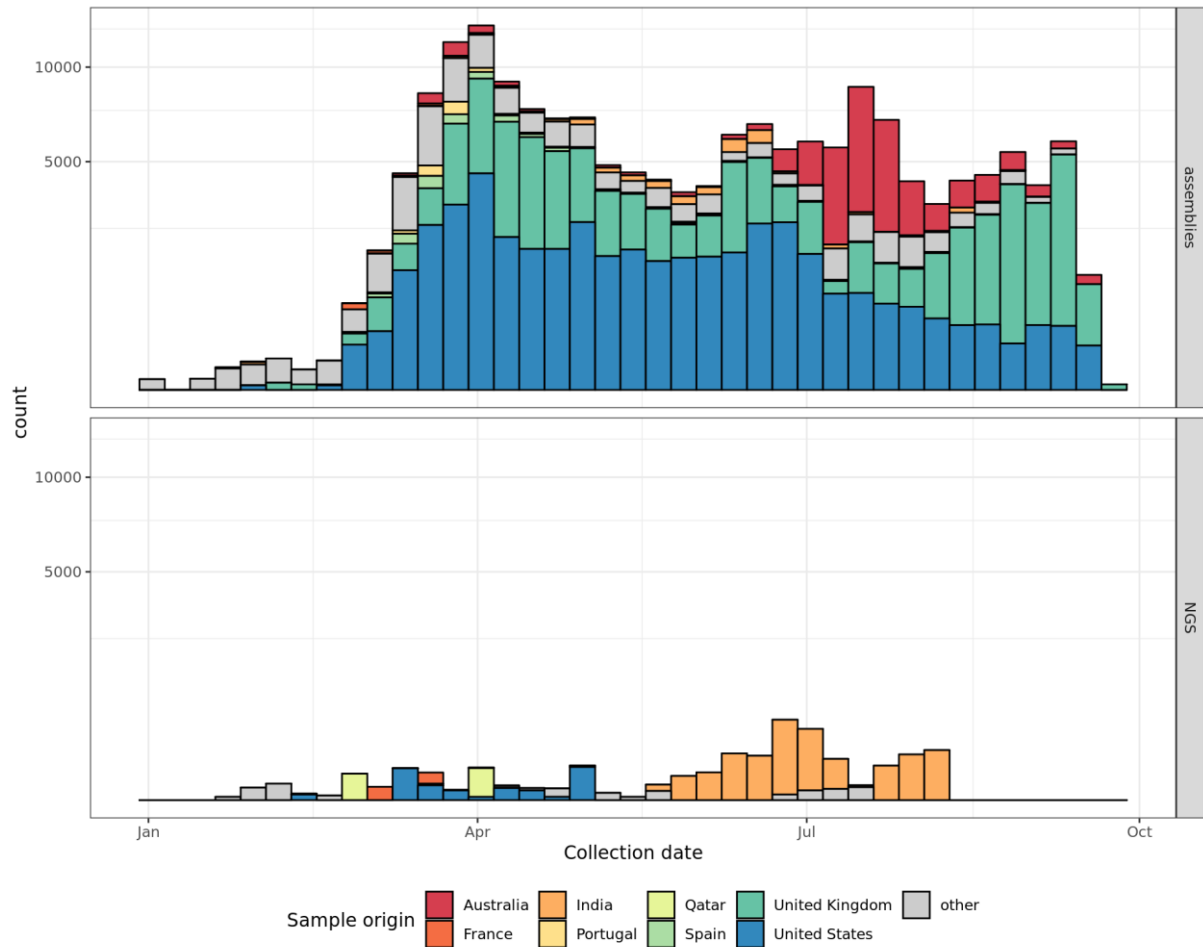
Published SARS-CoV-2 T-cell epitopes

SARS-CoV-2 antigens reported by Snyder et al. (2020) where downloaded from <https://clients.adaptivebiotech.com/pub/covid-2020> on 17NOV2020 (MIRA release 002.1).

Results

SARS-CoV-2 spike protein mutational profile from genome assemblies and NGS data

First, we determined the number of non-synonymous mutations in the spike protein per sample. Of the 146,920 analyzed genome assemblies and 2,393 NGS data sets (for geographic background, see Suppl. Figure S1), only 13.6% (20,248 samples) contained the WT spike protein (Figure 1A). Samples of mutant viruses exhibited only few mutations in the spike protein with less than ten mutations for all but 35 sequences. However, the mean and median number of mutations increased over time from December 2019 (mean: 0.14, median: 0) to September 2020 (mean: 1.98, median: 2, Figure 1B). Overall, we detected 2,592 distinct non-synonymous mutations in the spike protein (Supplementary Table S1).



Supplementary Figure S1: Number and origin of publicly available SARS-CoV-2 sequence data over time. The histogram shows the number of SARS-CoV-2 assembly sequences deposited at GISAID and NCBI Virus and NGS data deposited at SRA as of 02OCT2020. Color coding indicates the sample origin. Countries summarized as “other” include: Algeria, Andorra, Argentina, Aruba, Austria, Bahrain, Bangladesh, Belgium, Belize, Benin, Bosnia and Herzegovina, Botswana, Brazil, Bulgaria, Canada, Chile, China, Colombia, Congo [DRC], Costa Rica, Crimea, Croatia, Cuba, Cyprus, Czech Republic, Denmark, Dominican Republic, Ecuador, Egypt, Faroe Islands, Finland, Gabon, Gambia, Georgia, Germany, Ghana, Gibraltar, Greece, Guatemala, Hong Kong, Hungary, Iceland, Indonesia, Iran, Iraq, Ireland, Israel, Italy, Jamaica, Japan, Jordan, Kazakhstan, Kenya, Kuwait, Latvia, Lebanon, Lithuania, Luxembourg, Malaysia, Mali, Mexico, Moldova, Mongolia, Montenegro, Morocco, Myanmar, Netherlands, New Zealand, Nigeria, North Macedonia, Norway, Oman, Pakistan, Panama, Peru, Philippines, Poland, Puerto Rico, Reunion, Romania, Russia, Saudi Arabia, Senegal, Serbia, Sierra Leone, Singapore, Slovakia, Slovenia, South Africa, South Korea, Sri Lanka, Suriname, Sweden, Switzerland, Taiwan, Thailand, Tunisia, Turkey, Uganda, Ukraine, United Arab Emirates, Uruguay, Venezuela, Vietnam, Zambia and unknown.

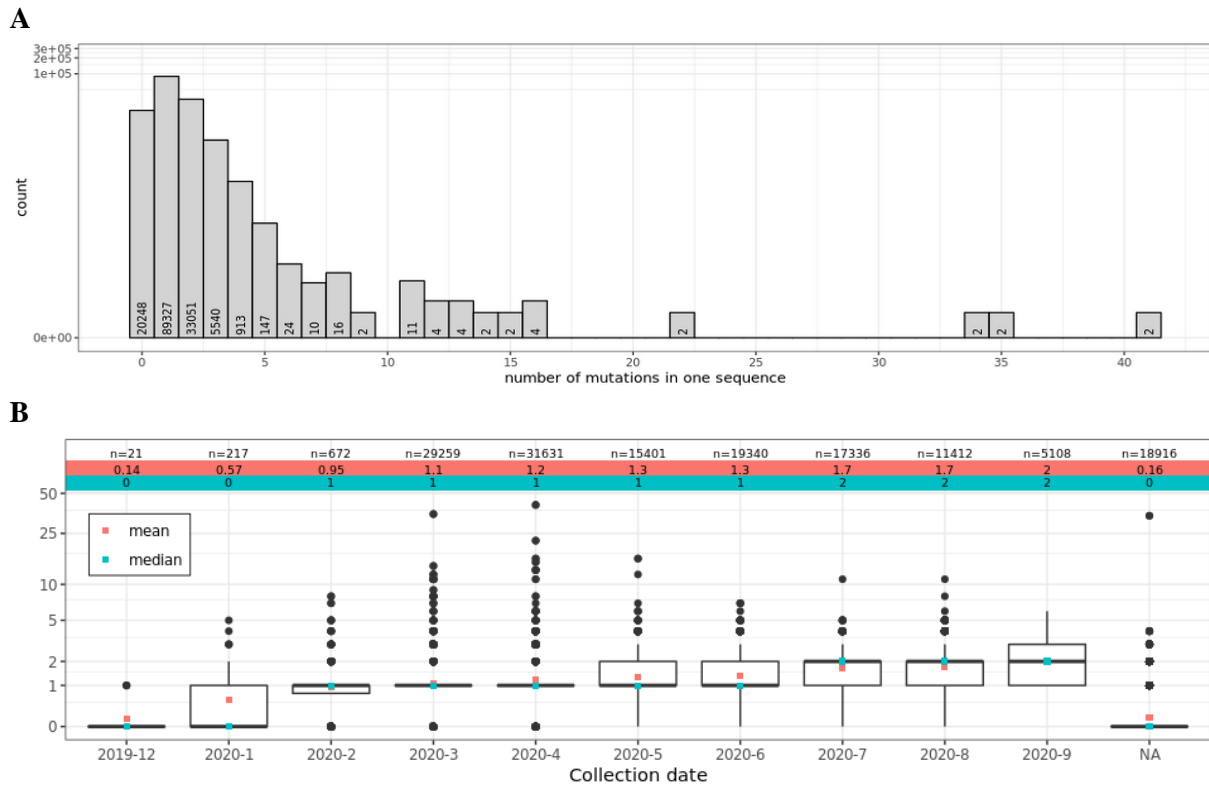


Figure 1: Most of the analyzed SARS-CoV-2 sequences differ from WT spike protein, but exhibit only few non-synonymous mutations. (A) The histogram shows the number of non-synonymous mutations in the spike protein detected in the analyzed samples. (B) The mean and median number of mutations per spike protein sequence increased over time.

Recurrent variants in SARS-CoV-2 spike protein

Most of the observed variants in the assembly and NGS data sets were recurrent (Figure 2A) and only 32.9% and 47.8% of the variants were singular events in the assembly and the NGS data, respectively. The recurrent variants were distributed throughout the whole spike protein (Figure 2B, C). Among the recurrent variants, nine and 23 mutations were found in at least 0.5% of the mutant assembly and NGS samples, respectively (labeled variants in Figure 2 B, C). The most common mutation was D614G in both the genome assemblies (124,179 samples) and the NGS data (1,792 samples) located outside the RBD (positions 319-529), followed by the RBD variants S477N in the assemblies (11,483 samples) and N440K in the NGS data (440 samples). In total, 339 distinct mutations (227 recurrent) were detected in the RBD in the assemblies out of which only two were common to more than 0.5% of the mutated assembly sequences (Figure 2A). For the NGS samples, 61 mutations in total (24 recurrent) were found in the RBD (Figure 2B) and again only two were detected in at least 0.5% of the mutant NGS samples. Overall, 246 mutations were commonly found in the assembly and NGS data (Figure 2C).

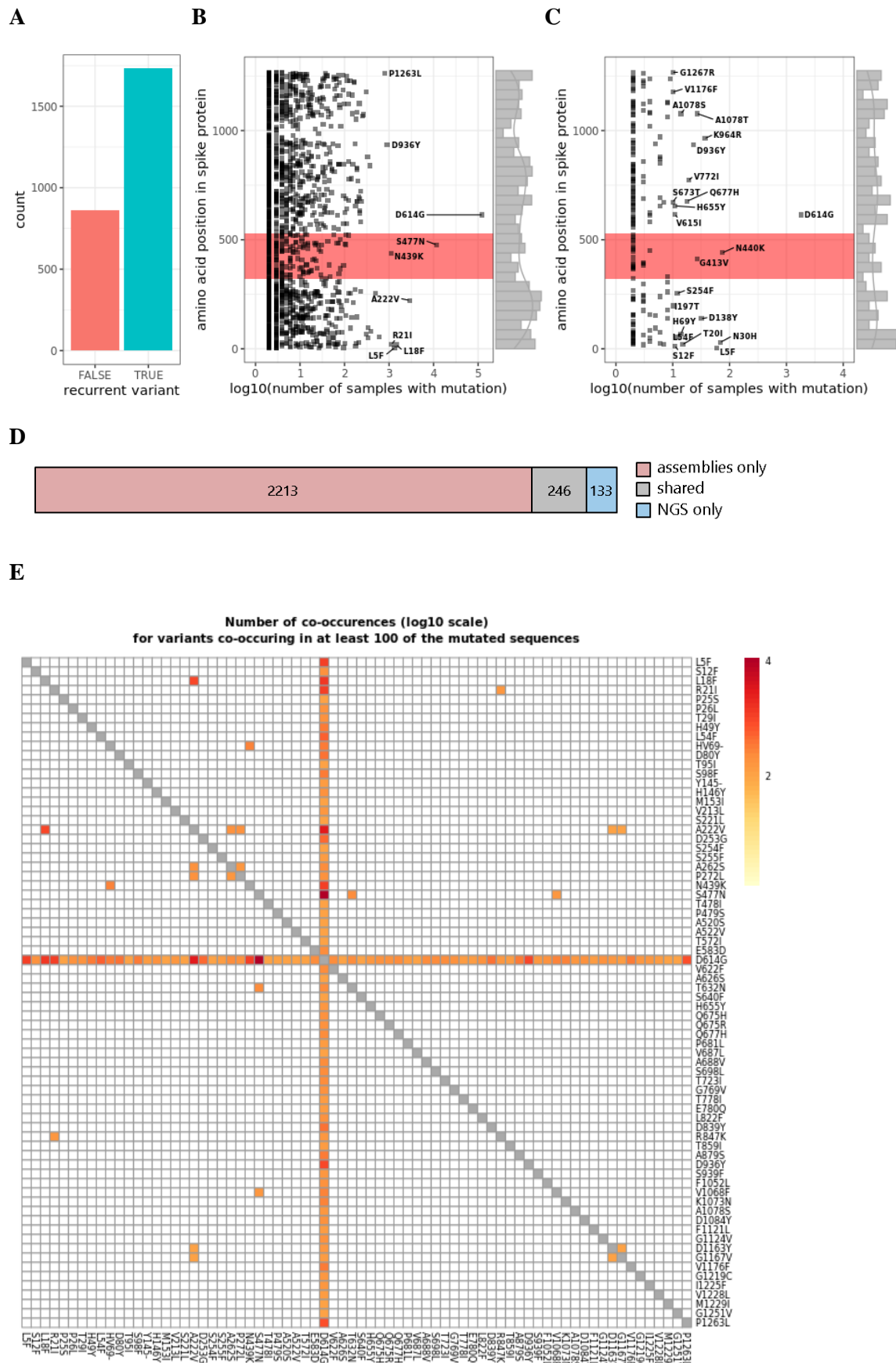


Figure 2: Recurrent variants are found throughout the whole spike protein. (A) Most of the detected variants were recurrent events occurring in at least two samples from the assembly or NGS data sets. (B, C) Each data point represents a distinct protein sequence mutation in the spike protein. The labels indicate the amino acid exchange for variants found in more

than 0.5% of the assemblies (B) or NGS samples (C). The RBD is highlighted in red. (D) 246 variants (grey) were detected both in the assemblies and the NGS data. (E) A subset of 72 variants co-occurred in at least 100 of the mutated spike protein sequences (assemblies and NGS data combined). For better visibility, co-occurrences in less than 100 samples were set to 0 (white tiles).

Furthermore, 72 (2.8%) of the detected variants co-occurred frequently in at least 100 of the mutated spike protein sequences when we combined assembly and NGS data (Figure 2D). Most prominent here, was the variant D614G which was found in combination with 1,385 other variants. The combination S477N/D614G was detected in 11,470 samples. These represented the above mentioned two most frequent variants in the assembly data. The most frequent co-occurring mutations not involving D614G were L18F/A222V.

Subclonal variants

In addition, we were interested in subclonal spike protein mutations (i.e. mutations with an observed variant frequency - as derived from the NGS reads - below 100%) which might either indicate co-infection with various SARS-CoV-2 strains and/or intra-host evolution of the virus. To this end, the fraction of variant supporting reads per sample of the detected mutations was determined. Most of the variants were observed with at least 95% of the reads supporting the respective variant nucleotide (Figure 3A, B). However, some mutations were only confirmed by a portion of the overlapping reads pointing to subclonal events. Filtering for a depth of at least 30 reads and a fraction of supporting reads between 0.1 and 0.95 (Sashittal et al. 2020) resulted in 363 mutations observed in 292 samples (i.e. 15.1% of the NGS data sets with mutant spike protein) that could be classified as high-confident subclonal (Figure 3B). Most of these subclonal events were recurrent variants (Figure 3C). Especially in the earlier samples, but also in some later cases, the fractions of supporting reads within the same sample differed notably.

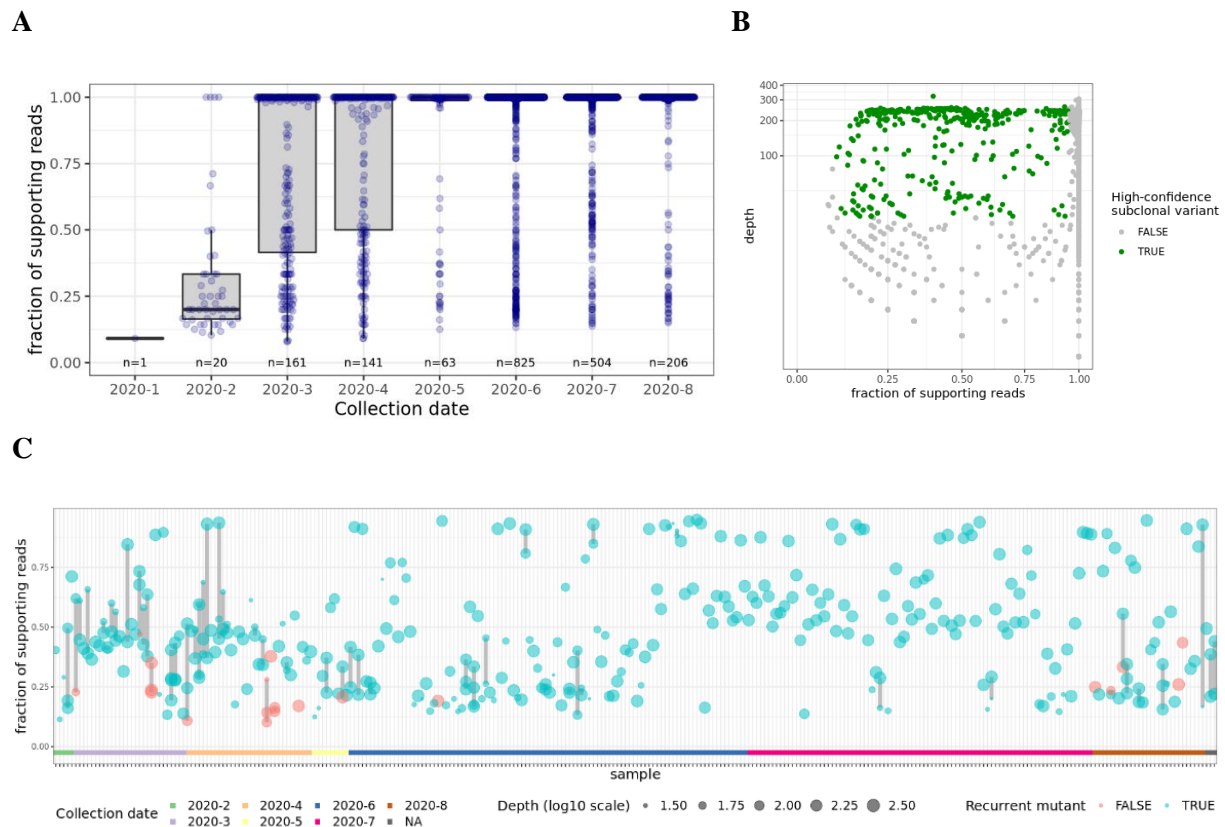


Figure 3: Variant frequencies of spike protein mutants indicate presence of multiple SARS-CoV-2 mutants in some samples. (A) The boxplot shows the distributions of the fraction of supporting reads of the mutations found in the NGS data. The numbers of underlying samples are indicated (n). Most of the observed variants have a variant allele frequency of ≥ 0.95

and can be accounted as clonal. (B) Filtering for high-confidence subclonal variants (green) with sequencing depth ≥ 30 reads and fractions of supporting reads between 0.1 and 0.95. (C) Sample-wise depiction of high-confidence subclonal events. Some of the observed subclonal variants were recurrent (blue) and only few were individual (red). The samples were ordered by collection date (see also color bar at the bottom of the plot) and point sizes indicate sequencing depth (log10 scale). Subclonal variants of the same sample are linked with grey lines. The fraction of supporting reads of variants found in the same sample differed notably in some cases.

Effect of detected spike protein variants on potential antibody and T cell target sites

Next, we investigated whether the observed spike protein variants were relevant in the context of antibody binding or T cell recognition. In order to be visible for antibodies, a mutation has to hit a residue on the surface of the trimeric spike protein complex. 432 (16,7%) of 2,592 unique variants affected surface residues. For the 20 most frequent among these occurring in at least 50 samples, the change in SASA from wild type to mutation at the mutated residue position was investigated (Figure 4A). The SASA changed for all but one (H245Y) of the variants which might influence the accessibility of neutralizing antibodies. Furthermore, 2,544 (98.1%) of the 2,592 distinct variants hit at least one CD8+ or CD4+ T-cell epitope (Figure 4B) when compared to the T-cell epitopes reported by Snyder et al. 2020 no matter if they were recurrent or individual events.

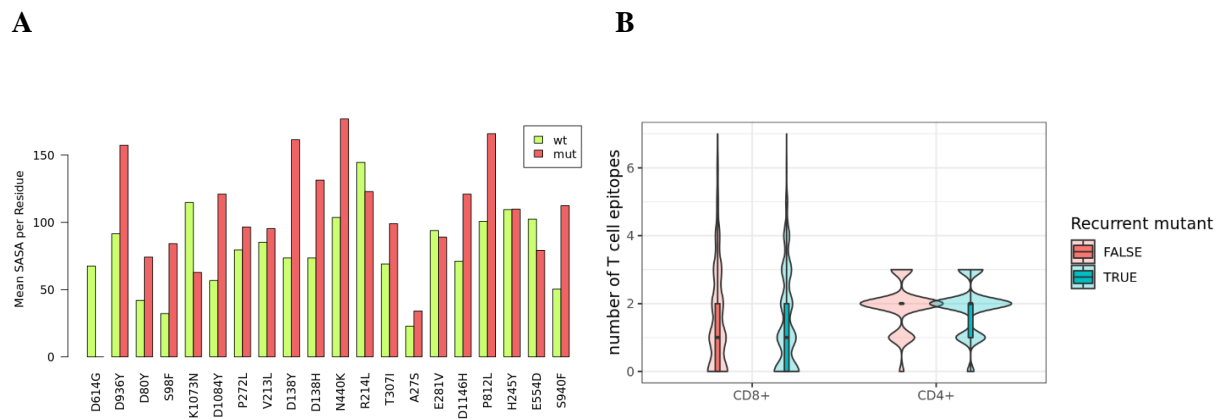


Figure 4: Variants affect antibody and T cell target sites. (A) Solvent-accessible surface area (SASA) values compared between wild type (wt) and mutation residue (mut) for surface variants occurring in at least 50 samples. The values are taken as the mean of the three replicated residues (3-meric structure of spike protein). Each time a new spike-protein structure has been generated by mutating the respective residue. The backbone of the mutated structure has not been re-modelled. The change in surface value is mainly due to change of amino acid and calculated optimal side-chain conformation. (B) The number of published T-cell epitopes (presented by MHC I or MHC II) that are affected by spike protein variants occurring in at least 50 analyzed samples is depicted. Most of the variants hit at least one epitope.

Discussion

Our study sheds light on non-synonymous variants in the spike protein of SARS-CoV-2 in a large cohort of samples from all over the world. While most analyzed sequences vary from the reference sample from Wuhan, China, our analysis of almost 150,000 assembly and NGS samples shows an overall low mutation burden in the SARS-CoV-2 spike protein across different host populations (Figure 1). However, the mean and median number of variants per sample increased over time. Coronaviruses (CoV) have fewer mutations compared to any other RNA virus due to its inherent 3' to 5' exonuclease activity (Minskaia et al. 2006). This suggests that the SARS-CoV-2 genome is genetically stable and the vast majority of mutations have no phenotypic effect such as virus transmissibility and virulence (Grubaugh et al. 2020b; Ruffell 2020). However, mutations of critical

residues in the RBD of the spike protein might increase the virus transmission ability by enhancing the interaction (Korber et al. 2020a). Furthermore, vaccines or treatments targeting the spike protein might become less efficient, if the number of variants in the spike protein increases further.

We identified a subset of mutations from the assembly and NGS data that are recurrent variants in the spike protein. van Dorp et al. (2020a) have already reported such recurrent variants in SARS-CoV-2 evolution, which is a likely phenomenon of positive selection signifying the adaption of SARS-CoV-2 in human hosts. Furthermore, most recurrent variants show no evidence in increase of viral transmission and are likely induced by host immunity through RNA editing mechanisms (van Dorp et al. 2020b). However, some variants might significantly influence SARS-CoV-2 transmission and infectivity. Among such variants, the non-synonymous D614G mutation has become most prevalent among several populations. We identified around 84.4% of the samples with a D614G variant, which supports a previous theory of an increasing frequency of the D614G variant in the global pandemic (Korber et al. 2020a; Korber et al. 2020b). Studies show evidence that the D614G variant is associated with high levels of viral RNA in COVID-19 patients, suggesting a role of D614G mutations in enhancing the viral infectivity in patients (Korber et al. 2020b; Hu et al. 2020; Ozono et al. 2020; Plante et al. 2020). In contrast to these findings, it remains unclear whether the D614G variant makes the infections more severe or may impact vaccine design (Grubaugh et al. 2020a), as the viral load does not correlate with disease severity and the variant is not in the RBD of the spike protein, which interacts with the human ACE2 protein.

The RBD of the spike protein is a potential target for neutralizing antibodies and the variants in these regions might influence the infectivity and pathogenicity. We have identified high frequency variants in the RBD region from the assembly data, i.e. S477N, N439K, N440K and G413V (Figure 2B, C). S477N occurs frequently almost similar to the D614G variant and studies show that S477N has potential to affect the RBD stability and strengthen the binding with the human ACE2 protein (Starr et al. 2020; Singh et al. 2020). In our study, S477N was most frequently co-occurring with D614G (Figure 2D). This combination was estimated to spread more rapidly than the D614G mutant alone (He and Wong). Other RBD variants such as N439K and N440K also show enhanced binding affinity to the human ACE2 receptor and result in immune escape from a panel of neutralizing monoclonal antibodies (Zhou et al. 2020b; Thomson et al. 2020; Weisblum et al. 2020). Antibody-resistant RBD variants might affect the therapeutic potential of neutralizing monoclonal antibodies by escaping through disruption of epitopes.

However, a significant portion of the detected variants represent individual events based on what could be deduced from the available data. This indicates the necessity to further collect SARS-CoV-2 isolates and monitor newly occurring variants. Here, the combination of assembly data (which appeared to be available in a timelier manner) and NGS samples (which also contain information on the clonality of the observed variants but which might be deposited with some delay) provide a valuable resource.

Further, we identified subclonal variants with a fraction of supporting reads between 0.1 and 0.95 at a sequencing depth of more than 30 reads in 15.1% of the NGS samples with mutant spike protein (Figure 3). Subclonal variants are indicative of within-host viral diversity leading to transmission of multiple strains (Sashittal et al. 2020). Low frequency variants could have been part of parallel evolution, where the same mutation rises to detectable frequencies in different lineages and it is observed as part of SARS-CoV-2 virus adaptation (Wright et al. 2020). Further, recurrent mutations might point to co-infection with multiple strains. Sample-specific variants in turn might rather indicate that the mutation occurred after infection within the host. This viral diversity within the host might prevent complete clearance after treatment and thus might lead to the development of resistant strains. Also, subclonal variants should be considered for vaccine design as these might represent the next generation of the virus.

The analyzed data sets also showed that a notable portion of the individual and recurrent mutations in the spike protein (98.1%) overlap with at least one known T-cell epitope. They also may change the solvent-accessible area and thus antibody binding when they involve surface residues of the trimeric

spike protein complex as shown for the 20 most frequent solvent-accessible mutations. While we had no information on the HLA-restriction of the published T-cell epitopes, the influence on CD8 T cell epitope generation by different HLA alleles was investigated for the three common mutations L5F, D614G and G1124V (Guo and Guo 2020). These mutations were predicted to result in epitope gains, losses or higher or lower HLA binding affinities. Greaney et al. (2021) presented a system to map mutations in the SARS-CoV-2 RBD that escape antibody binding. However, there is no overlap with our exemplary analysis on SASA changes. In agreement with the increase of the SASA of the mutation N440K, the binding affinity of this mutant to antibody REGN10933_REGN10987 is strengthened (Chen et al.). All these findings demonstrate that SARS-CoV-2 mutants need to be set in the context of immune recognition to evaluate their implications for the global spreading of the pandemic and future preventive or therapeutic approaches in a timely manner.

Conclusion and Outlook

Human infections with SARS-CoV-2 are spreading globally since the beginning of 2020, necessitating preventive or therapeutic strategies and first steps towards an end to this pandemic were done with the approval of the first mRNA vaccines against SARS-CoV-2. Here, we show different types of variants (recurrent vs. individual, clonal vs. subclonal, hitting T-cell or antibody target sites vs. not-hitting) that can be incorporated in global efforts to sustainably prevent or treat infections. The underlying computational strategy might serve as a template for a platform to constantly analyze globally available sequencing data. In combination with a web-based platform to administer the results, this could help guiding global vaccine design efforts to overcome the threats of this pandemic.

The importance of our approach is underlined by the recently emerging UK lineage B.1.1.7 of SARS-CoV-2 (Rambaut et al. 2020), which is characterized by the accumulation of 17 variants; eight of those are located in the S protein. This lineage has a higher transmissibility compared to other lineages (Volz et al. 2021). The occurrence of this lineage questioned the efficacy of current vaccines, but first results showed that it at least unlikely will escape BNT162b-induced protection (Muik et al. 2021). Interestingly, the individual variants can be traced back to samples from March (P681H, T716I) and April (Y114del, N501Y, A570D) of 2020. It needs to be mentioned that the available data, although representing a large cohort, might not reflect the real distribution of the circulating variants as mostly samples of specific interest will be sequenced. International sequencing efforts, combined data analysis and prediction of variant impact will be important tools for the future in order to ensure an early detection of such genomic variants of concern.

Conflict of Interest

Author U.S. is co-founder, shareholder and CEO at BioNTech SE. The remaining authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Acknowledgments

We thank Pablo Riesgo Ferreira and Patrick Sorn for critical discussions. We gratefully acknowledge the authors from the originating laboratories responsible for obtaining the specimens, as well as the submitting laboratories where the sequence data were generated and shared via GISAID, NCBI Virus or the NCBI SRA, on which this research is based.

References

- Braun, Julian; Loyal, Lucie; Frentsch, Marco; Wendisch, Daniel; Georg, Philipp; Kurth, Florian et al. (2020): Presence of SARS-CoV-2 reactive T cells in COVID-19 patients and healthy donors.
- Brister, J. Rodney; Ako-Adjei, Danso; Bao, Yiming; Blinkova, Olga (2015): NCBI viral genomes resource. In *Nucleic acids research* 43 (Database issue), D571-7. DOI: 10.1093/nar/gku1207.
- Chen, Jiahui; Gao, Kaifu; Wang, Rui; Wei, Guowei: Prediction and mitigation of mutation threats to COVID-19 vaccines and antibody therapies. Available online at <http://arxiv.org/pdf/2010.06357v1>.
- Cingolani, Pablo; Platts, Adrian; Le Wang, Lily; Coon, Melissa; Nguyen, Tung; Wang, Luan et al. (2012): A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff. SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. In *Fly* 6 (2), pp. 80–92. DOI: 10.4161/fly.19695.
- Elbe, Stefan; Buckland-Merrett, Gemma (2017): Data, disease and diplomacy. GISAID's innovative contribution to global health. In *Global challenges (Hoboken, NJ)* 1 (1), pp. 33–46. DOI: 10.1002/gch2.1018.
- Greaney, Allison J.; Starr, Tyler N.; Gilchuk, Pavlo; Zost, Seth J.; Binshtein, Elad; Loes, Andrea N. et al. (2021): Complete Mapping of Mutations to the SARS-CoV-2 Spike Receptor-Binding Domain that Escape Antibody Recognition. In *Cell host & microbe* 29 (1), 44-57.e9. DOI: 10.1016/j.chom.2020.11.007.
- Grubaugh, Nathan D.; Hanage, William P.; Rasmussen, Angela L. (2020a): Making Sense of Mutation. What D614G Means for the COVID-19 Pandemic Remains Unclear. In *Cell*. DOI: 10.1016/j.cell.2020.06.040.
- Grubaugh, Nathan D.; Petrone, Mary E.; Holmes, Edward C. (2020b): We shouldn't worry when a virus mutates during disease outbreaks. In *Nature microbiology* 5 (4), pp. 529–530. DOI: 10.1038/s41564-020-0690-4.
- Guex, N.; Peitsch, M. C. (1997): SWISS-MODEL and the Swiss-PdbViewer. An environment for comparative protein modeling. In *Electrophoresis* 18 (15), pp. 2714–2723. DOI: 10.1002/elps.1150181505.
- Guo, Elisa; Guo, Hailong (2020): CD8 T cell epitope generation toward the continually mutating SARS-CoV-2 spike protein in genetically diverse human population. Implications for disease control and prevention. In *PloS one* 15 (12), e0239566. DOI: 10.1371/journal.pone.0239566.
- Hatcher, Eneida L.; Zhdanov, Sergey A.; Bao, Yiming; Blinkova, Olga; Nawrocki, Eric P.; Ostapchuck, Yuri et al. (2017): Virus Variation Resource - improved response to emergent viral outbreaks. In *Nucleic acids research* 45 (D1), D482-D490. DOI: 10.1093/nar/gkw1065.
- He, Shiyu; Wong, Samuel W. K.: Statistical challenges in the analysis of sequence and structure data for the COVID-19 spike protein. Available online at <http://arxiv.org/pdf/2101.02304v1>.
- Hu, Jie; He, Chang Long; Gao, Qingzhu; Zhang, Gui Ji; Cao, Xiao Xia; Long, Quan Xin et al. (2020): The D614G mutation of SARS-CoV-2 spike protein enhances viral infectivity. In *bioRxiv : the preprint server for biology*. Available online at <https://doi.org/10.1101/2020.06.20.161323>.
- Korber, Bette; Fischer, Will; Gnanakaran, S. Gnana; Yoon, Heyjin; Theiler, James; Abfalterer, Werner et al. (2020a): Spike mutation pipeline reveals the emergence of a more transmissible form of SARS-CoV-2. In *bioRxiv : the preprint server for biology*. Available online at <https://doi.org/10.1101/2020.04.29.069054>.
- Korber, Bette; Fischer, Will M.; Gnanakaran, Sandrasegaram; Yoon, Hyejin; Theiler, James; Abfalterer, Werner et al. (2020b): Tracking Changes in SARS-CoV-2 Spike. Evidence that D614G Increases Infectivity of the COVID-19 Virus. In *Cell*. DOI: 10.1016/j.cell.2020.06.043.

Kuipers, Jack; Batavia, Aashil A.; Jablonski, Kim Philipp; Bayer, Fritz; Borgsmüller, Nico; Dondi, Arthur et al. (2020): Within-patient genetic diversity of SARS-CoV-2. In *bioRxiv : the preprint server for biology*.

Lan, Jun; Ge, Jiwan; Yu, Jinfang; Shan, Sisi; Zhou, Huan; Fan, Shilong et al. (2020): Crystal structure of the 2019-nCoV spike receptor-binding domain bound with the ACE2 receptor. In *bioRxiv : the preprint server for biology*.

Leinonen, Rasko; Sugawara, Hideaki; Shumway, Martin (2011): The sequence read archive. In *Nucleic acids research* 39 (Database issue), D19-21. DOI: 10.1093/nar/gkq1019.

Li, Fang (2016): Structure, Function, and Evolution of Coronavirus Spike Proteins. In *Annual review of virology* 3 (1), pp. 237–261. DOI: 10.1146/annurev-virology-110615-042301.

Li, Heng (2013): Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. Available online at <http://arxiv.org/pdf/1303.3997v2>.

Li, Heng; Handsaker, Bob; Wysoker, Alec; Fennell, Tim; Ruan, Jue; Homer, Nils et al. (2009): The Sequence Alignment/Map format and SAMtools. In *Bioinformatics (Oxford, England)* 25 (16), pp. 2078–2079. DOI: 10.1093/bioinformatics/btp352.

Mahase, Elisabeth (2020): Covid-19. FDA authorises neutralising antibody bamlanivimab for non-admitted patients. In *BMJ (Clinical research ed.)* 371, m4362. DOI: 10.1136/bmj.m4362.

Minskaia, Ekaterina; Hertzog, Tobias; Gorbalenya, Alexander E.; Campanacci, Valérie; Cambillau, Christian; Canard, Bruno; Ziebuhr, John (2006): Discovery of an RNA virus 3'-5' exoribonuclease that is critically involved in coronavirus RNA synthesis. In *Proceedings of the National Academy of Sciences of the United States of America* 103 (13), pp. 5108–5113. DOI: 10.1073/pnas.0508200103.

Morais, Ivair José; Polveiro, Richard Costa; Souza, Gabriel Medeiros; Bortolin, Daniel Inserra; Sasaki, Flávio Tetsuo; Lima, Alison Talis Martins (2020): The global population of SARS-CoV-2 is composed of six major subtypes. In *Scientific reports* 10 (1), p. 18289. DOI: 10.1038/s41598-020-74050-8.

Muik, Alexander; Wallisch, Ann-Kathrin; Sängler, Bianca; Swanson, Kena A.; Mühl, Julia; Chen, Wei et al. (2021): Neutralization of SARS-CoV-2 lineage B.1.1.7 pseudovirus by BNT162b2 vaccine-elicited human sera.

Ou, Junxian; Zhou, Zhonghua; Zhang, Jing; Lan, Wendong; Zhao, Shan; Wu, Jianguo et al. (2020): RBD mutations from circulating SARS-CoV-2 strains enhance the structural stability and human ACE2 affinity of the spike protein. In *bioRxiv : the preprint server for biology*.

Ozono, Seiya; Zhang, Yanzhao; Ode, Hirotaka; Seng, Tan Toong; Imai, Kazuo; Miyoshi, Kazuyasu et al. (2020): Naturally mutated spike proteins of SARS-CoV-2 variants show differential levels of cell entry. In *bioRxiv : the preprint server for biology*. Available online at <https://doi.org/10.1101/2020.06.15.151779>.

Plante, Jessica A.; Liu, Yang; Liu, Jianying; Xia, Hongjie; Johnson, Bryan A.; Lokugamage, Kumari G. et al. (2020): Spike mutation D614G alters SARS-CoV-2 fitness. In *Nature*. DOI: 10.1038/s41586-020-2895-3.

Rambaut, Andrew; Loman, Nick; Pybus, Oliver; Barclay, Wendy; Barrett, Jeff; Carabelli, Alesandro et al. (2020): Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in the UK defined by a novel set of spike mutations. COVID-19 Genomics Consortium UK (CoG-UK). Available online at <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>, checked on 1/20/2021.

Ruffell, Daniela (2020): Coronavirus SARS-CoV-2. Filtering fact from fiction in the infodemic: Q&A with virologist Professor Urs Greber. In *FEBS letters* 594 (7), pp. 1127–1131. DOI: 10.1002/1873-3468.13784.

Sashittal, Palash; Luo, Yunan; Peng, Jian; El-Kebir, Mohammed (2020): Characterization of SARS-CoV-2 viral diversity within and across hosts. In *bioRxiv : the preprint server for biology*. Available online at <https://doi.org/10.1101/2020.05.07.083410>.

Shang, Weiling; Yang, Yi; Rao, Yifan; Rao, Xiancai (2020): The outbreak of SARS-CoV-2 pneumonia calls for viral vaccines. In *NPJ vaccines* 5, p. 18. DOI: 10.1038/s41541-020-0170-0.

Singh, Amit; Steinkellner, Georg; Köchl, Katharina; Gruber, Karl; Gruber, Christian C. (2020): Serine 477 plays a crucial role in the interaction of the SARS-CoV-2 spike protein with the human receptor ACE2.

Snyder, Thomas M.; Gittelman, Rachel M.; Klinger, Mark; May, Damon H.; Osborne, Edward J.; Taniguchi, Ruth et al. (2020): Magnitude and Dynamics of the T-Cell Response to SARS-CoV-2 Infection at Both Individual and Population Levels. In *medRxiv : the preprint server for health sciences*. DOI: 10.1101/2020.07.31.20165647.

Starr, Tyler N.; Greaney, Allison J.; Hilton, Sarah K.; Ellis, Daniel; Crawford, Katharine H. D.; Dingens, Adam S. et al. (2020): Deep Mutational Scanning of SARS-CoV-2 Receptor Binding Domain Reveals Constraints on Folding and ACE2 Binding. In *Cell* 182 (5), 1295–1310.e20. DOI: 10.1016/j.cell.2020.08.012.

Sun, C.; Chen, L.; Yang, J.; Luo, C.; Zhang, Y.; Li, J. et al. (2020): SARS-CoV-2 and SARS-CoV spike-RBD structure and receptor binding comparison and potential implications on neutralizing antibody and vaccine development. *bioRxiv*. Published online February 20, 2020. In *Google Scholar*.

Tai, Wanbo; He, Lei; Zhang, Xiujuan; Pu, Jing; Voronin, Denis; Jiang, Shibo et al. (2020): Characterization of the receptor-binding domain (RBD) of 2019 novel coronavirus. Implication for development of RBD protein as a viral attachment inhibitor and vaccine. In *Cellular & molecular immunology* 17 (6), pp. 613–620. DOI: 10.1038/s41423-020-0400-4.

Tang, Xiaolu; Wu, Changcheng; Li, Xiang; Song, Yuhe; Yao, Xinmin; Wu, Xinkai et al. (2020): On the origin and continuing evolution of SARS-CoV-2. In *National Science Review* 7 (6), pp. 1012–1023. DOI: 10.1093/nsr/nwaa036.

Thomson, Emma; Rosen, Laura; Shepherd, James; Spreafico, Roberto; da Silva Filipe, Ana; Wojcechowskyj, Jason et al. (2020): The circulating SARS-CoV-2 spike variant N439K maintains fitness while evading antibody-mediated immunity: *bioRxiv*.

van Dorp, Lucy; Acman, Mislav; Richard, Damien; Shaw, Liam P.; Ford, Charlotte E.; Ormond, Louise et al. (2020a): Emergence of genomic diversity and recurrent mutations in SARS-CoV-2. In *Infection, genetics and evolution : journal of molecular epidemiology and evolutionary genetics in infectious diseases* 83, p. 104351. DOI: 10.1016/j.meegid.2020.104351.

van Dorp, Lucy; Richard, Damien; Tan, Cedric C. S.; Shaw, Liam P.; Acman, Mislav; Balloux, François (2020b): No evidence for increased transmissibility from recurrent mutations in SARS-CoV-2. In *bioRxiv : the preprint server for biology*.

Volz, Erik; Mishra, Swapnil; Chand, Meera; Barrett, Jeffrey C.; Johnson, Robert; Geidelberg, Lily et al. (2021): Transmission of SARS-CoV-2 Lineage B.1.1.7 in England. Insights from linking epidemiological and genetic data. In *medRxiv : the preprint server for health sciences*. DOI: 10.1101/2020.12.30.20249034.

Wan, Yushun; Shang, Jian; Graham, Rachel; Baric, Ralph S.; Li, Fang (2020): Receptor Recognition by the Novel Coronavirus from Wuhan. An Analysis Based on Decade-Long Structural Studies of SARS Coronavirus. In *Journal of virology* 94 (7). DOI: 10.1128/JVI.00127-20.

Weisblum, Yiska; Schmidt, Fabian; Zhang, Fengwen; DaSilva, Justin; Poston, Daniel; Lorenzi, Julio Cc et al. (2020): Escape from neutralizing antibodies by SARS-CoV-2 spike protein variants. In *eLife* 9. DOI: 10.7554/eLife.61312.

World Health Organization (2021): WHO Coronavirus Disease (COVID-19) Dashboard. Data last updated: 2021/1/19, 6:44pm CET. Available online at <https://covid19.who.int/>.

Wright, Erik Scott; Lakdawala, Seema S.; Cooper, Vaughn S. (2020): SARS-CoV-2 genome evolution exposes early human adaptations. In *bioRxiv : the preprint server for biology*. Available online at <https://doi.org/10.1101/2020.05.26.117069>.

Wu, Fan; Zhao, Su; Yu, Bin; Chen, Yan-Mei; Wang, Wen; Song, Zhi-Gang et al. (2020): A new coronavirus associated with human respiratory disease in China. In *Nature* 579 (7798), pp. 265–269. DOI: 10.1038/s41586-020-2008-3.

Zhou, Peng; Yang, Xing-Lou; Wang, Xian-Guang; Hu, Ben; Zhang, Lei; Zhang, Wei et al. (2020a): A pneumonia outbreak associated with a new coronavirus of probable bat origin. In *Nature* 579 (7798), pp. 270–273. DOI: 10.1038/s41586-020-2012-7.

Zhou, Wenyang; Xu, Chang; Wang, Pingping; Luo, Meng; Xu, Zhaochun; Cheng, Rui et al. (2020b): N439K variant in spike protein may alter the infection efficiency and antigenicity of SARS-CoV-2 based on molecular dynamics simulation: bioRxiv.

Zhu, Na; Zhang, Dingyu; Wang, Wenling; Li, Xingwang; Yang, Bo; Song, Jingdong et al. (2020): A Novel Coronavirus from Patients with Pneumonia in China, 2019. In *The New England journal of medicine* 382 (8), pp. 727–733. DOI: 10.1056/NEJMoa2001017.