

# The emergence of highly fit SARS-CoV-2 variants accelerated by epistasis and caused by recombination

Michael R. Garvin<sup>1,2,\*</sup>, Erica T. Prates<sup>1,2+</sup>, Jonathon Romero<sup>3</sup>, Ashley Cliff<sup>3</sup>, Joao Gabriel Felipe Machado Gazolla<sup>1,2</sup>, Monica Pickholz<sup>4,5</sup>, Mirko Pavicic<sup>1,2</sup>, Daniel Jacobson<sup>1,2,\*</sup>

## Affiliations:

<sup>1</sup>Oak Ridge National Laboratory, Computational Systems Biology, Biosciences, Oak Ridge, TN; <sup>2</sup>National Virtual Biotechnology Laboratory, US Department of Energy; <sup>3</sup>The Bredesen Center for Interdisciplinary Research and Graduate Education, University of Tennessee Knoxville, Knoxville, TN; <sup>4</sup>Departamento de Física, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires, Buenos Aires, Argentina; <sup>5</sup>Instituto de Física de Buenos Aires (IFIBA), CONICET-Universidad de Buenos Aires, Buenos Aires, Argentina.

\*Correspondence: [garvinmr@ornl.gov](mailto:garvinmr@ornl.gov), [jacobsonda@ornl.gov](mailto:jacobsonda@ornl.gov)

+Contributed equally

**Key words:** SARS-CoV-2; recombination; haplotype; epistasis; variants of concern; COVID-19, syncytia

**Running title:** Highly fit SARS-CoV-2 generated from recombination

## Abstract

The SARS-CoV-2 pandemic has entered an alarming new phase with the emergence of the variants of concern (VOC), P.1, B.1.351, and B.1.1.7, in late 2020, and B.1.427, B.1.429, and B.1.617, in 2021. Substitutions in the spike glycoprotein (S), such as Asn<sup>501</sup>Tyr and Glu<sup>484</sup>Lys, are likely key in several VOC. However, Asn<sup>501</sup>Tyr circulated for months in earlier strains and Glu<sup>484</sup>Lys is not found in B.1.1.7, indicating that they do not fully explain those fast-spreading variants. Here we use a computational systems biology approach to process more than 900,000 SARS-CoV-2 genomes, map their spatiotemporal relationships, and identify lineage-defining mutations followed by structural analyses that reveal their critical attributes. Comparisons to earlier dominant mutations and protein structural analyses indicate that increased transmission is promoted by epistasis, i.e., the combination of functionally complementary mutations in S and in other regions of the SARS-CoV-2 proteome. We report that the VOC have in common mutations in non-S proteins involved in immune-antagonism and replication performance, such as the nonstructural proteins 6 and 13, suggesting convergent evolution of the virus. Critically, we propose that

recombination events among divergent coinfecting haplotypes greatly accelerates the emergence of VOC by bringing together cooperative mutations and explaining the remarkably high mutation load of B.1.1.7. Therefore, extensive community distribution of SARS-CoV-2 increases the probability of future recombination events, further accelerating the evolution of the virus. This study reinforces the need for a global response to stop COVID-19 and future pandemics.

*Notice: This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the US Department of Energy (DOE). The US government retains and the publisher, by accepting the article for publication, acknowledges that the US government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for US government purposes. DOE will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).*

*"Nothing in Biology Makes Sense Except in the Light of Evolution" -Theodosius Dobzhansky*

## 1. Introduction

In late 2020, three SARS-CoV-2 variants of concern, VOC; B.1.1.7, B.1.351, and P.1 (also called alpha, beta, and gamma respectively) rapidly became the predominant source of infections due to enhanced transmission rates and have since been linked to increased hospitalizations and mortalities (Alpert et al. 2021; Challen et al. 2021; Davies et al. n.d.; Faria et al. 2021; Funk et al. 2021; Sabino et al. 2021; Volz et al. n.d.; Washington et al. 2021). In early 2021, several new VOC appeared including, B.1.427 (epsilon), B.1.526 (iota), and B.1.617 (delta). B.1.617 is of immediate concern because it is responsible for the COVID-19 crisis that recently began in India (J. Singh et al. 2021), is causing the majority of new infections in the United Kingdom (UK), and the United States (USA), and has now been observed in more than 70 countries worldwide. Notably, several of these VOC have rapidly spread even in regions such as the UK that depend on robust sampling efforts for early detection. There is therefore a critical need to identify

accurate predictors and biological causes for the increased transmission of the next VOC, which will inevitably emerge if the viral spread is not globally restrained.

Although extensive efforts are underway to achieve these ends, integrating new findings is critical to unravel the multiple biomolecular and environmental factors influencing viral evolution. Toward a holistic understanding of VOC emergence, two major weaknesses need to be addressed: (1) currently, the mutations used to identify VOC and potentially explain the altered biology of the virus are predominantly focused on the changes observed in the spike glycoprotein (S) whereas those in other genomic regions are largely ignored, and (2) the molecular models used to reconstruct the evolutionary history of the virus employ phylogenetic trees that are useful for species-level but not population-based analyses and are unable to incorporate information and molecular events that are critical to understanding SARS-CoV-2 (Huson & Bryant 2006; Velasco 2013). Indeed, a recent study introduced novel methods to detect recombination in coronaviruses and found that their evolution shows very little resemblance to a tree structure (Müller et al. 2021).

For example, the Asn<sup>501</sup>Tyr substitution in S is likely key because it increases affinity for the host receptor, angiotensin-converting enzyme 2 (ACE2) (Liu et al. 2021), and is often used to identify the late 2020 VOC (Fratev n.d.; Luan et al. n.d.), but this mutation has been circulating widely at low frequency and only expanded seven months after being first detected. Similarly, the Glu<sup>484</sup>Lys substitution in S is often discussed in the context of P.1 and B.1.351 VOC and may allow escape from neutralizing antibodies (Greaney et al. n.d.; Starr et al. 2021), but is not found in B.1.1.7 and therefore does not explain the increased transmission of all three late 2020 VOC. These characteristics suggest that several mutations including those in S are being transmitted as

a linked set, i.e. a haplotype and their combined effects (i.e., epistasis) may be contributing to the rapid viral spread.

Widely used molecular evolutionary models based on phylogenetic trees are also problematic because the *algorithms* that are applied assume that mutations appearing in different SARS-CoV-2 haplotypes are due to repeated, independent mutations and the *scientific community* is interpreting this as evidence for the same; i.e., the logic is circular. Alternatively, these apparent repeated mutations may represent recombination, which is a common mechanism to accelerate evolution compared to single site mutations in positive strand RNA viruses such as SARS-CoV-2 (Bentley & Evans 2018; Simon-Loriere & Holmes 2011). Furthermore, phylogenetic trees are unable to incorporate important molecular events and metadata such as geospatial and temporal data that would be highly informative for detecting current and future VOC.

In contrast, median-joining networks (MJN) are an efficient and accurate means to analyze haploid genomic data at the population level (Bandelt et al. 1999) such as SARS-CoV-2, (Garvin, Prates, et al. 2020). Unlike independently segregating sites represented by phylogenetic trees, the unit of interest in an MJN is the haplotype, which more accurately reflects the biology of coronaviruses and enables the detection of important evolutionary events such as recombination. Furthermore, a network can be annotated with information including frequency, geospatial location, demographic, or clinical outcomes associated with a unique haplotype to create interpretable patterns of genome variation.

Here, we processed more than 900,000 SARS-CoV-2 genomes using a computational workflow that combines MJN and protein structural analysis (Garvin, Prates, et al. 2020; E. T. Prates et al. 2020) to identify critical attributes of the mutations that define these VOC and



provide substantial evidence that the genome-wide mutation load of the late 2020 VOC results from recombination between divergent strains. Using structural analysis and molecular dynamics simulations, we explore the individual effects of key mutations in S and other proteins of SARS-CoV-2 that are shared among different VOC. Additionally, we identify a signature of co-evolution between the residue 501 in S and ACE2, and propose the molecular basis for that. Overall, our results indicate that linked mutations in VOC generated from recombination act in an epistatic manner to enhance viral spread. This work emphasizes the important role of community spread in generating future VOC (Sheikh et al. 2021).

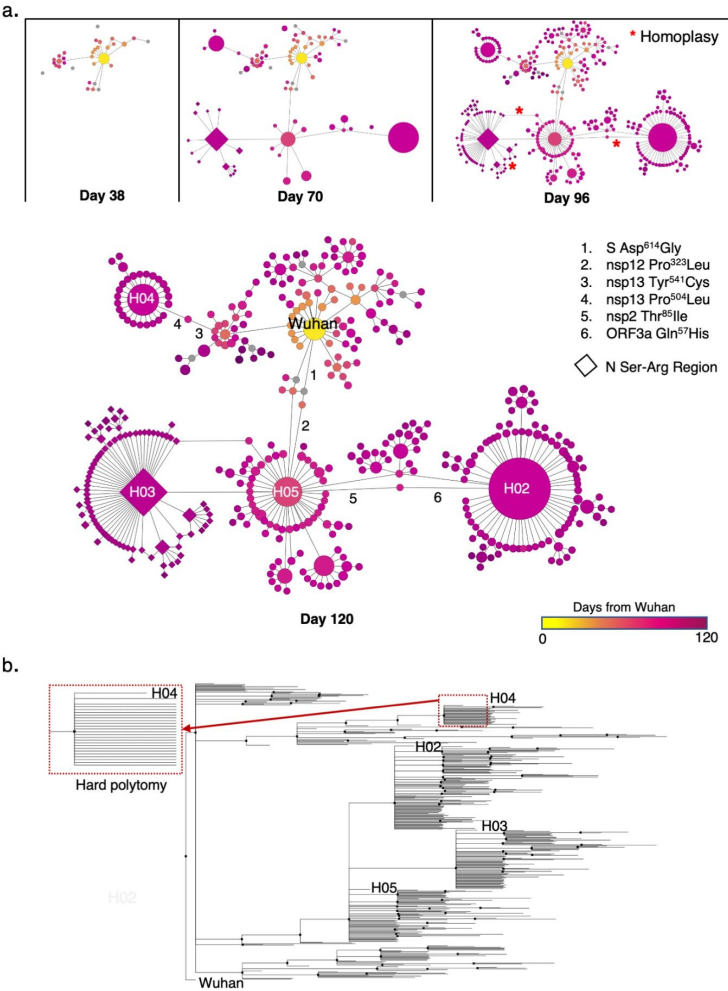
## 2. Results and Discussion

### *Network-based views of molecular evolution provide critical information that phylogenetic trees do not*

The COVID-19 pandemic is both an unprecedented tragedy and an opportunity to study molecular evolution given the abundant and global sampling of the mutational space of SARS-CoV-2. The MJN is a valuable method of integrating these data to understand viral evolution because the model assumes single mutational steps in which each node represents a haplotype and the edge between nodes is a mutation leading to a new one. Typically, a subsample of extant haplotypes for a taxon is obtained and unsampled, or extinct lineages are inferred. In contrast, SARS-CoV-2 sequence data repositories provide extensive sampling of haplotypes and collection dates (the calendar date of the sample). Given that in an MJN, the temporal distribution of haplotypes is inherent (the model assumes time-ordered sets of mutations), the mutational history of the virus can be traced as a genealogy that can incorporate both the relative

and absolute times of emergence of SARS-CoV-2 variants. Importantly, when the single-mutational step MJN model fails, it produces features such as loops or clusters of inferred haplotypes that can indicate biologically important processes such as recombination events, back mutations, or repeat mutations at a site that may be under positive selection.

In order to make a direct comparison, we generated a network and a phylogenetic tree of SARS-CoV-2 haplotypes that were identified from sequences sampled during the first four months of the pandemic and deposited into GISAID (A Global Initiative on Sharing Avian Flu Data, [gisaid.org](https://gisaid.org)) (Figure 1). Clearly, important metadata such as haplotype frequency, date of emergence, and mutations of interest are easily displayed on the network but are not on the phylogenetic tree. Likewise, at day 96, reticulations (i.e., homoplasy loops) begin to appear in the MJN, indicating reverse mutations to the ancestral states at specific sites or possibly recombination events that can be explored further. Another important feature identified when using networks, but is lost when using phylogenetic trees, is the presence of polytomies. So-called soft polytomies often indicate unsampled genomic information at the species level and hard polytomies are molecular events often found in rapidly expanding populations. For example, haplotype H04 in the MJN (Figure 1) represents a hard polytomy and indicates that a frequent variant is further undergoing multiple independent mutational events, but the phylogenetic tree is unable to convey this information. This example demonstrates the significant gain in information an MJN provides compared to a tree-based view of coronavirus evolution.



**Fig. 1. Comparison of a median-joining network (MJN) and phylogenetic tree generated with SARS-CoV-2 sequences sampled through April 2020.** a. MJN of SARS-CoV-2 haplotypes, 96, and 120 days. Node sizes in the MJN correspond to sample sizes for a given haplotype and node colors indicate the time of its first report relative to the putative origin of the pandemic in Wuhan. The most abundant haplotypes are named H02 - H05 and numerals 1 - 6 identify important mutations (Garvin, Prates, et al. 2020). Diamond shape nodes denote haplotypes that harbor a 3-bp mutation in the nucleocapsid gene (N) that is highly conserved and directly affects viral replication *in vitro* (Thorne et al. 2021; Tylor et al. 2009). b. The phylogenetic tree is unable to convey the same information. For example, rapidly expanding populations often display polytomies, i.e., single mutations from a common central haplotype. Those events are readily identified on the MJN but difficult to interpret on a tree because they are usually visualized as a multi-pronged fork (outlined in the dashed-line box) rather than a star pattern (compare H04 in (a) and (b)). These true biological processes also cause tree algorithms to perform poorly because they violate their assumptions, slowing convergence. Additionally, MJN are able to indicate reticulations (i.e., loops) that could denote recombination, reverse mutations, or other biologically important events whereas the forced bifurcation of phylogenetic tree algorithms is unable to display these. Reference sequence: NC\_045512, Wuhan, December 24, 2019.

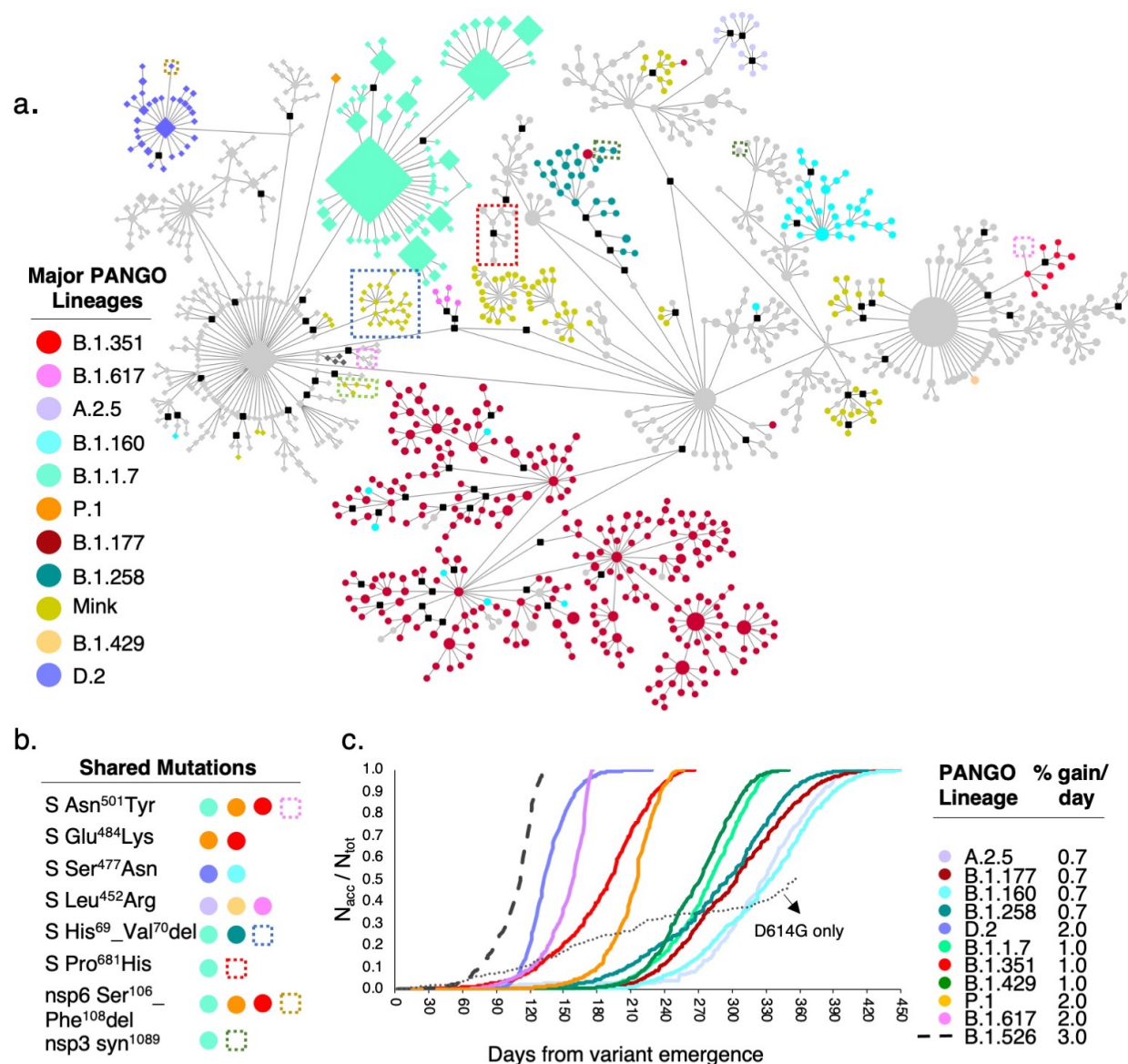
# *Networks identify lineage-defining mutations of Variants of Concern*

We processed more than 900,000 SARS-CoV-2 genomes from human and mink, built a MJN network using the 640,211 genomes that survived our quality control workflow, and annotated with PANGO lineages defined in the GISAID database (Figure 2, *see* Methods). This genealogy-based approach to molecular evolution identifies the mutations that define a given haplotype based on the edge between nodes. Here, they define the variants of concern and variants of interest (VOI) based on the edge that initiates their corresponding clusters of nodes (Table 1).

This approach also enables the identification of all the common features acquired in different VOC/VOI, which can elucidate the set of molecular features underlying their rapid spread. For example, the B.1.1.7, B.1.351, and P.1 variants, here referred to as late\_2020\_VOC, can all be defined by a triple amino acid deletion in the nonstructural protein 6 (nsp6; Ser<sup>106</sup>\_Phe<sup>108</sup>del) as well as Asn<sup>501</sup>Tyr in S. Notably this latter mutation has received considerable attention compared to the former (Figure 2a, Table 1), but both are likely key to the biology of these VOC. The MJN also reveals what appears to be a previous dominant but now rare variant (D.2) in Australia; an Ile<sup>120</sup>Phe mutation in nsp2 was followed immediately by S Ser<sup>477</sup>Asn, which seems to have led to its rapid expansion in April 2020, indicated by increased node size. In contrast, B.1.160 carries the S Ser<sup>477</sup>Asn mutation without nsp2 Ile<sup>120</sup>Phe, and did not demonstrate the same rapid expansion as D.2. This underscores the usefulness of the MJN approach as it is able to convey the sequence of timing of mutational events and the number of individuals carrying those haplotypes simultaneously.

In order to account for sampling bias (there are a disproportionate number of sequences contributed to GISAID by the UK and Australia as well as the high frequency of B.1.1.7 and D.2

197 in those two geographic locations), we plotted the number of daily samples of selected variants  
198 relative to their respective total number of cases to date (April, 2021) and compared the resulting  
199 slopes of the linear range of the curves (Figure 2c). The late\_2020\_VOC, early\_2021\_VOC, and  
200 early\_2021\_VOI (Table 1) display higher daily accumulation rates, between 1% and 3% of total  
201 observed cases per day, compared to other variants (e.g., B.1.177), which show less than 1%  
202 accumulation per day. Notably, the rapid increase in D.2 (2%), but not B.1.160 (0.7%) supports  
203 the MJN view of this as a likely VOC and the importance of epistasis (here S Ser<sup>477</sup>Asn and nsp2  
204 Ile<sup>120</sup>Phe). This analysis and the MJN confirm the importance of monitoring these variants  
205 closely and identify both S and non-S mutations that define the current and potential SARS-  
206 CoV-2 VOC (Table 1).



**Fig. 2. Median-joining network (MJN) of SARS-CoV-2 genomes.** **a.** MJN of haplotypes found in more than 30 individuals (N=640,211 sequences) using 2,128 variable sites. Colors identify PANGO lineages from GISAID. Diamond-shaped nodes correspond to haplotypes carrying a three base pair deletion in the nucleocapsid gene (N) at sites 28881-28883 (Arg<sup>203</sup>Lys, Gly<sup>203</sup>Arg). Black square nodes are inferred haplotypes, dashed-line box defines a subgroup of haplotypes within a lineage with a disjoint mutation that is also found in B.1.1.7. Several lineages show introgression from others (e.g., cyan nodes, B.1.160, into brick red, B.1.177). **b.** Several important mutations in S and non-S proteins appear in multiple lineages. For example, the B.1.1.7 variant carries four mutations that are in disjoint nodes: S Asn<sup>501</sup>Tyr, S Pro<sup>681</sup>His, a silent mutation in the codon for amino acid 1089 in nsp3, and the S His<sup>69</sup>\_Val<sup>70</sup>del that is also found in a clade of haplotypes from mink (blue dashed-line box in (a)). **c.** Accumulation rate for common GISAID lineages including VOC represented by the ratio between the accumulated number of reported sequences of a given lineage per day since the appearance of that haplotype ( $N_{acc}$ ) divided by the



corresponding total number ( $N_{\text{tot}}$ ) at the final sample date for this study. Colors of curves correspond to node colors in (a). All VOC display accumulation rates of at least 1% of the total for that variant per day. The remaining are less than 1% except for the VOI B.1.526 (not displayed in MJN) that is the highest with 3% per day, indicating further scrutiny of this variant is warranted. We also plotted the accumulation rate for lineages that carry the widely reported S Asp<sup>614</sup>Gly mutation but without the nsp12 Pro<sup>323</sup>Leu commonly found with it, supporting our previous hypothesis (Garvin, Prates, et al. 2020) that mutations in S alone are not responsible for the rapid transmission of these VOC/VOI but is a function of epistasis among S and non-S mutations. Reference sequence: NC\_045512, Wuhan, December 24, 2019.

**Table 1. Major lineages shown in the median-joining network and their defining mutations.** Center for disease control (CDC)-defined variants and their timing are listed under *Lineages* and discussed in the text. L-VOC denotes likely variants of concern, that is, those that we propose to have strong potential to become VOC. Non-VOC (N-VOC) are not identified by CDC as VOC. Potential epistatic non-S mutations lineage-defining mutations are listed for the VOC, L-VOC, and VOI. Sites in red font are discussed in the text.

Lineage	Class	Spike Mutation(s)	Likely non-S Epistatic Partner(s)	First major detection
B.1	Early_2020_VOC	D614G	nsp12 P323L	Germany
B.1.1.7 (alpha)	Late_2020_VOC	N501Y, del 69-70, P681H*, T716I, D1118H	nsp6 del 106-108, N L3D, N S235Y	United Kingdom
B.1.351 (beta)	Late_2020_VOC	N501Y, E484K, K417N	nsp6 del 106-108	South Africa
P.1 (gamma)	Late_2020_VOC	N501Y, E484K, K417N	nsp6 del 106-108	Brazil
B.1.427 (epsilon)	Early_2021_VOC	L452R, S13I	nsp13 D260Y	United States, California
B.1.429	Early_2021_VOC	L452R, W152C	nsp13 D260Y	United States, Washington
B.1.617 (delta)	Early_2021_VOC	L452R, E484Q, P681R*	N R203M, ORF7a V82G, ORF3a S26L	India
A.2.5	L-VOC	L452R, del 142-145	nsp1 L4P, nsp3 K839E, nsp4 P308Y	Panama
D.2	L-VOC	S477N	nsp2 I120F	Australia
B.1.160	N-VOC	S477N	na	Denmark
B.1.177	N-VOC	A222V	na	United Kingdom/Denmark
B.1.258	N-VOC	N434K, del 69-70	na	Denmark
B.1.526 (iota)	Early_2021_VOI	L5P, T95I, D253G	nsp6 del 106-108, nsp4 L438P, nsp13 Q88H	United States, New York

\* multibasic furin cleavage site

### *Recombination is the likely source for the rapidly expanding variants*

Haploid, clonally replicating organisms such as SARS-CoV-2 are predicted to eventually become extinct due to the accumulation of numerous slightly deleterious mutations over time, i.e., Muller's ratchet (Muller 1964). Recombination is not only a rescue from Muller's ratchet, it can also accelerate evolution by allowing for the union of advantageous mutations from divergent haplotypes (Bentley & Evans 2018). In SARS-CoV-2, recombination manifests as a template switch during replication when more than one haplotype is present in the host cell, i.e. the virus replisome stops processing a first RNA strand and switches to a second one from a different haplotype, producing a hybrid virus (Simon-Loriere & Holmes 2011). In fact, template

switching is a necessary step during the negative-strand synthesis of SARS-CoV-2 when the replisomes functions as an RNA-dependent RNA polymerase and pauses at transcription-regulatory motifs of the sub-genomic template to add the leader sequence from the 5' end of the genome (this “recombination” is not detected if only a single strain is present, i.e. there is no variation) (Kim et al. 2020),). Given this and the fact that recombination is a major mechanism of coronavirus evolution (Boni et al. 2020) it would be improbable for this process *not* to occur in the case of multiple strains infecting a cell (Gribble et al. 2021).

The late\_2020\_VOC exhibits large numbers of new mutations relative to any closely related sequence indicating rapid evolution of SARS-CoV-2 (Figure 3a). For example, the original node of B.1.1.7 differs from the most closely related node by 28 mutations. However, the majority of this total (15) corresponds to deletions that could be considered two single mutational events, as does a 3-bp change in N (28280-22883) since they occur in factors of three (a codon), maintaining the coding frame. By summing the two deletions and the full codon 3-bp change in N with the 10 remaining single site mutations, a conservative estimate would be 13 distinct mutational events leading to B.1.1.7. The plot of the accumulating mutations in the 640,211 haplotypes sampled to date reveals a linear growth of roughly 0.05 mutations per day (Figure 3b) and therefore, given this pace, it would be expected to take 260 days for these 13 mutational events to accumulate in a haplotype. For the B.1.1.7 to appear in October 2020 as reported, the genealogy would have to have been initiated in January 2020 and yet the nearest node harboring the S Asn<sup>501</sup>Tyr mutation was not sampled until June 2020 and no intermediate haplotypes have been identified to date.

Alternatively, it could be that the 13 mutational events occurred between June and October (122 days), but the probability of this is about one in 10<sup>15</sup> (Supplemental Methods).



Furthermore, all 28 differences between the Wuhan reference sequence and B.1.1.7 appeared in earlier haplotypes (Table S1) and therefore, if rapid evolution were the cause, it would require the extremely unlikely process of 28 independent, repeat mutations to the same nucleotide state. In order to test this, we plotted the population-level mutations per day (including repeat mutations at variable sites), which did not reveal any increase in mutation rate at the time of the B.1.1.7 and in fact displayed a *decrease* with the emergence of the late\_2020\_VOC (Figure 3c, Figure S1). Possible explanations are either a large increase in mutations in a small number of individuals over a short time period (that would have to occur on multiple continents to explain B.1.1.7, B.1.351, and P.1), or recombination between two or more divergent haplotypes carrying the VOC mutations (Figure 3d).

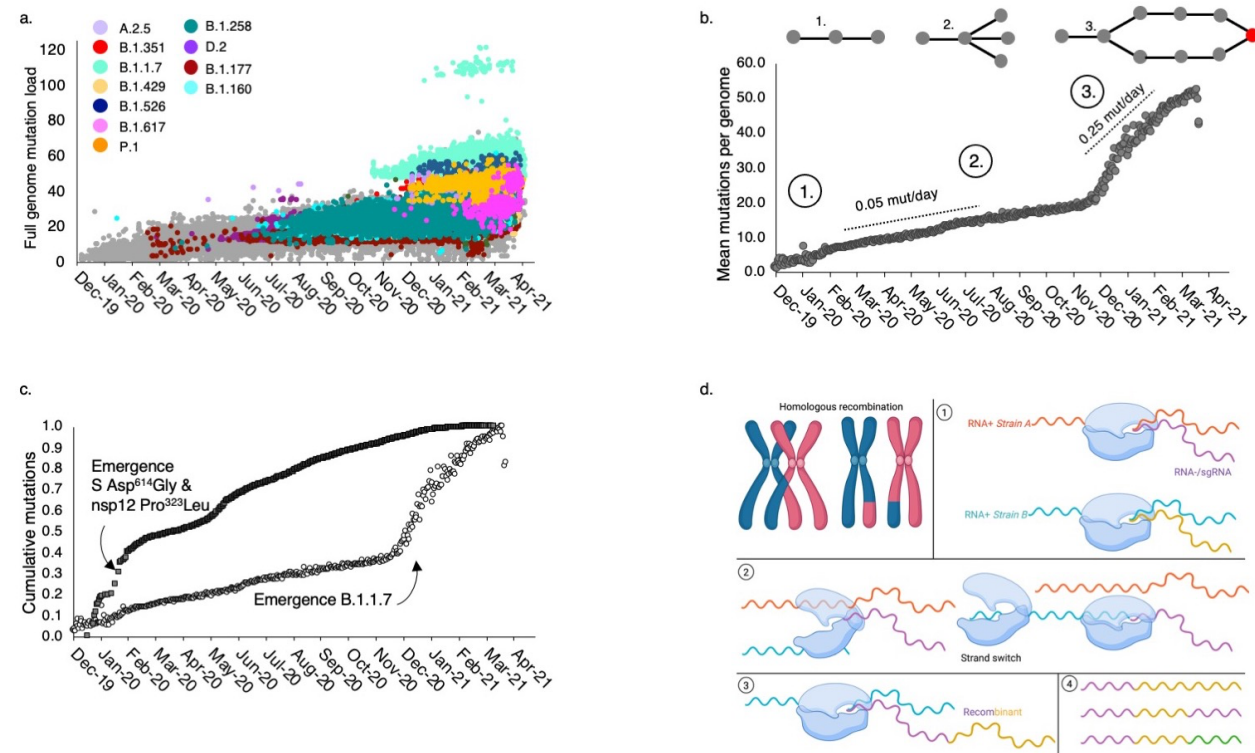
Recombination is the most parsimonious explanation given (1) the absence of a substantial increase in mutation rate at any time prior to the appearance of the VOC along with their increased mutation load, which can easily be explained by template-switching during replication (Simon-Loriere & Holmes 2011), (2) the widespread and early circulation of the majority of the mutations associated with them in other haplotypes and, (3) that several mutations appear disjointly across the MJN (Figure 2a). The first notable disjoint mutation is Ser<sup>477</sup>Asn in S that defines D.2 along with nsp2 Ile<sup>120</sup>Phe (Figure 2a), which then appears in B.1.160. Likewise, Asn<sup>501</sup>Tyr and Pro<sup>681</sup>His in S appear in divergent haplotypes, including one mink subgroup from Denmark and a basal node to B.1.351 (without the nsp6 deletion). It could be argued that those in S (Asn<sup>501</sup>Tyr, Ser<sup>477</sup>Asn, and Pro<sup>681</sup>His) are the result of multiple independent mutation events because they are under positive selection (Martin et al. 2021), but we also identified a mutation in nsp3 that is one of the lineage-defining mutations for B.1.1.7 and appears in disjoint nodes, but is unlikely to be under selection because it is synonymous.

Furthermore, mutation and selection are separate events, i.e., even if sites appearing in multiple lineages are under positive selection, their appearance in disjoint nodes still requires repeated mutations in the absence of recombination.

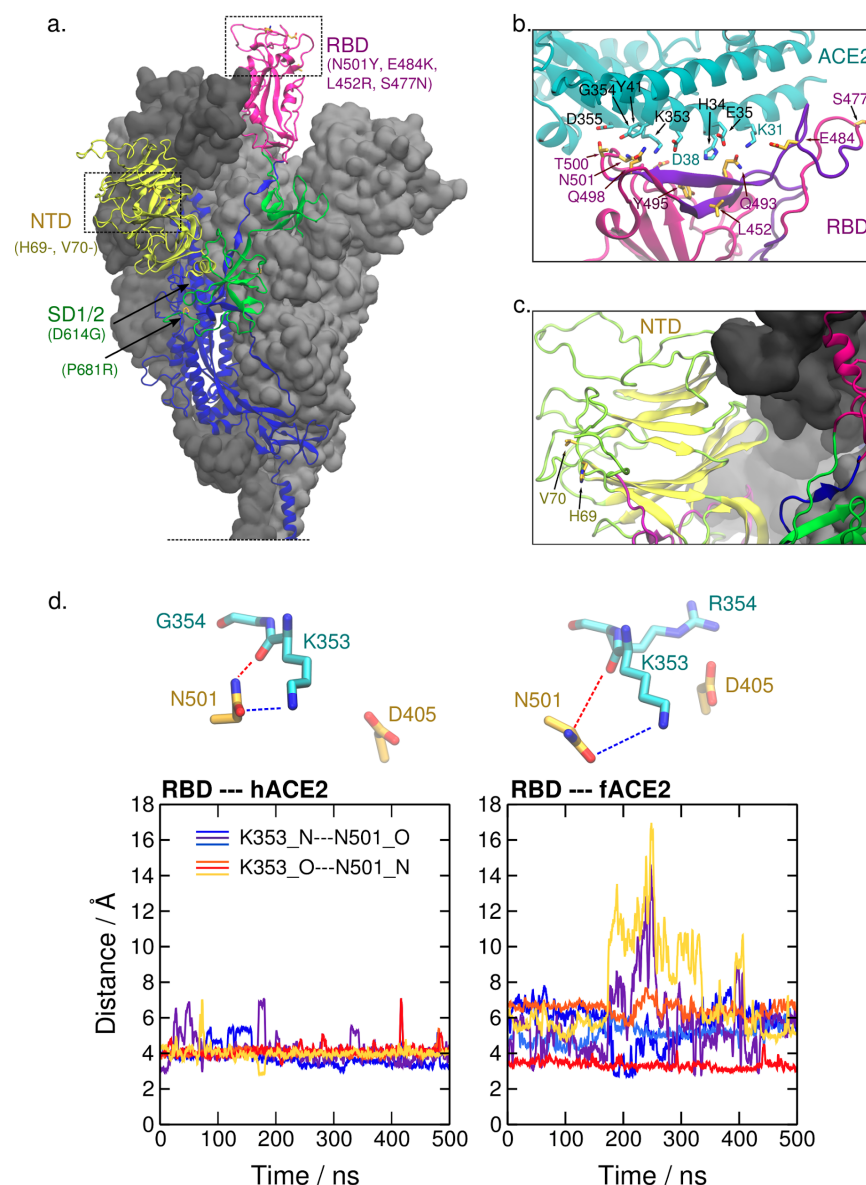
It should also be noted that recombination can generate a high number of false-positives when testing for signs of positive selection (Anisimova et al. 2003), and the complexity of coronavirus recombinants compared to those generated in diploid organisms through homologous chromosome crossovers (Figure 3d) makes that process difficult to detect. Therefore, analyses that test for positive selection based on multiple independent mutations at a site may, in fact, be false positives that result from recombination events. The majority of the mutations found in B.1.1.7 could be explained by the admixture and recombination among lineages and a random scan of 100 FASTQ files from B.1.1.7 available in the NCBI SRA database identified two co-infected individuals in further support of this hypothesis (Table S2). Large-scale analyses of these data with newer methods may enable the detection of recombinants (Müller et al. 2021).

A recent analysis of sequences in the U.K. that leveraged the geographical and temporal presence of specific sequences identified recombinants whose parents appeared to be from B.1.1.7 and B.1.1.177-derived strains (Jackson et al. n.d.). Although this demonstrates that recombination is occurring, it does not answer if B.1.1.7 itself arose through recombination. As noted above, the majority of B.1.1.7-defining mutations appeared much earlier in the pandemic outside the U.K. and therefore approaches that focus on a single geographic location will miss parental variants. The probability of proving B.1.1.7 originated from recombination using this method is further reduced if the parental lineages were less fit or poorly sampled and therefore rare. The present study demonstrates that detecting recombination events requires the use of data

on a global scale and the inclusion of rare variants.



**Fig. 3 Mutation rates and genomic mutation load of SARS-CoV-2.** **a.** A rapid increase in the number of mutations per individual genome is evident in the late\_2020\_VOC. The outliers of the B.1.1.7 lineages (mint green) are a subset of that lineage due to a single, 57 base pair deletion in ORF7a (amino acids 5-23). **b.** Mean mutation load per individual, based on 2,128 high-confidence sites by date. The SARS-CoV-2 virus accumulated an estimated 0.05 mutations per day until the appearance of B.1.1.7, when it increased five-fold. Circles with numbers denote three processes occurring at different timepoints: (1) emergence, (2) haplotype expansion, and (3) recombination of divergent lineages. **c.** A population-level analysis of new mutations per day over the same time period (dark squares) displays a declining rate of mutations with a slight increase around the emergence of D.2 in Australia but not an increase with the emergence of B.1.1.7 that could explain the rapid accumulation of mutations shown in (b) plotted as percent accumulation (unfilled circles). **d.** Recombination in a diploid organism results from the crossover of homologous chromosomes during meiosis. In RNA+ betacoronaviruses, recombination occurs when two or more strains (haplotypes) infect a single cell (1). The replisome dissociates (2) from one strand and switches to another, (3) generating a hybrid recombinant. The resulting chimera (4) can be as simple as a section of strain A fused to a section of strain B or more complex recombinants if strand switching occurs more than once or there are multiple strains per cell (green section). Reference sequence: NC\_045512, Wuhan, December 24, 2019.



**Fig. 4. Location of mutation sites of SARS-CoV-2 VOC on the structure of the spike glycoprotein.** **a.** Several mutations associated with dominant haplotypes are located in the receptor-binding domain (RBD, aa. 331-506), N-terminal domain (NTD, aa. 13-305), and subdomains 1 and 2 (SD1/2, aa. 528-685) of S. The structure of S in the prefusion conformation derived from PDB ID 6VSB (Wrapp et al. 2020) and completed *in silico* (Casalino et al. 2020) is shown. Glycosyl chains are not depicted and the S trimer is truncated at the connecting domain for visual clarity. The secondary structure framework of one protomer is represented and the neighboring protomers are shown as a gray surface. **b.** Mutation sites in the S RBD of SARS-CoV-2 VOC, such as 484, 452, 477, and 501 are located at or near the interface with ACE2. Notably, site 452 and 484 reside in an epitope that is a target of the adaptive immune response in humans (aa. 480-499, in violet) and site 501 is also located near it (B.-Z. Zhang et al. 2020). Dashed lines represent relevant polar interactions discussed here. PDB ID 6M17 was used (R. Yan et al. 2020). **c.** The sites 69 and 70 on the NTD, which are deleted in the VOC B.1.1.7, are also found near an epitope (aa. 21-45, in violet) (B.-Z. Zhang et al. 2020). **d.** Time progression of N---O distances between atoms of Asn<sup>501</sup> in RBD and Lys<sup>353</sup> in human and ferret ACE2 (hACE2 and fACE2, respectively) from the last 500 ns of the simulation runs. Colors in the plots correspond to the distances Lys<sup>353</sup>\_N---Asn<sup>501</sup>\_O (cold colors) and Lys<sup>353</sup>\_O---Asn<sup>501</sup>\_N (warm

colors) in three independent simulations of each system. These distances are represented in the upper part of the figure.

### ***The potential functional impact of key mutations in S and non-S proteins***

Given the results from the MJN analysis and our previous hypothesis (Garvin, Prates, et al. 2020) that the cooperative effects of mutations in S and non-S proteins (i.e., epistasis) define and are responsible for the increased transmission of prevalent SARS-CoV-2 variants (Lauring & Hodcroft 2021), we performed protein structural analyses and discuss below the functional effects of these individual and combined mutations in SARS-CoV-2 VOC. We analyze ten likely key mutation sites (red font, Table 1) in S and non-S proteins.

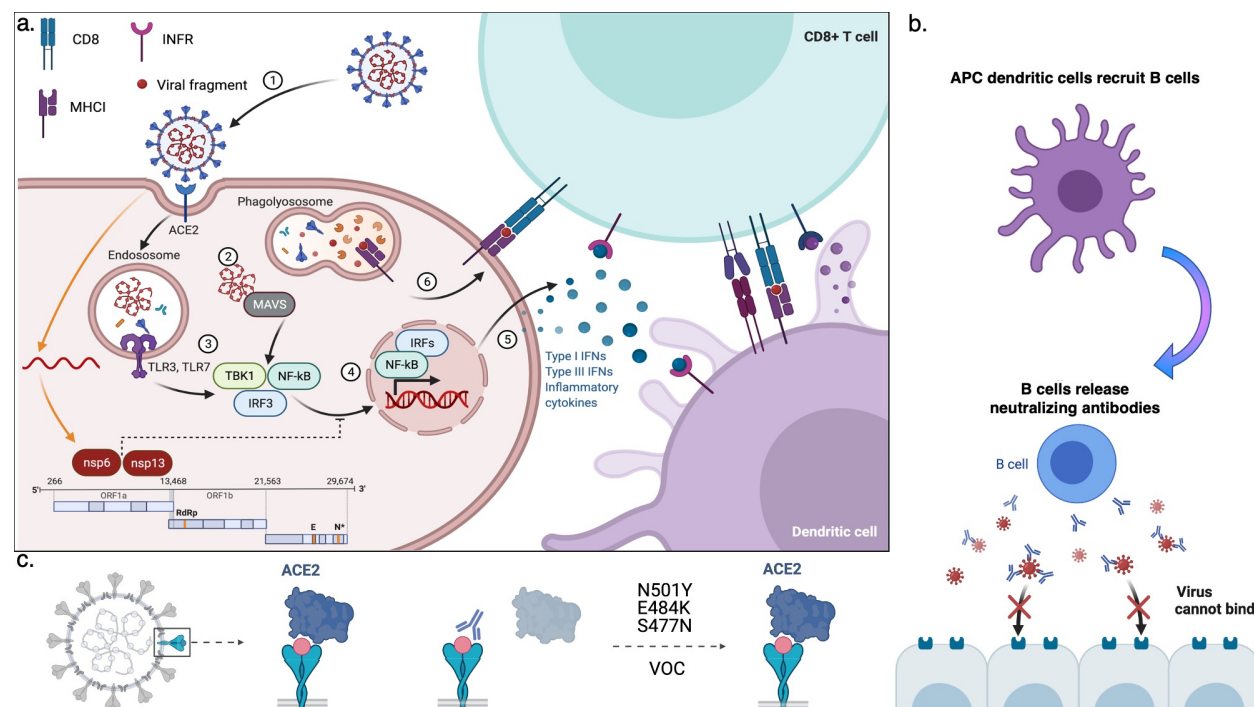
*S Asn<sup>501</sup>Tyr* - Located in the receptor-binding domain (RBD) of SARS-CoV-2 S, immunoprecipitation assays reveal that site 501 plays a major role in the affinity of the virus to the host receptor, ACE2 (Shang et al. 2020). Via structural analysis and extensive molecular dynamics simulations, Ali et al. highlighted the importance of the interactions with human ACE2 (hACE2) near the site 501 of the receptor-binding domain of S, particularly via a sustained hydrogen bond between RBD Asn<sup>498</sup> and hACE2 Lys<sup>353</sup> (Ali & Vijayan 2020). Deep mutational scanning of SARS-CoV-2 RBD reveals that the naturally occurring mutations at site 501, Asn<sup>501</sup>Tyr and Asn<sup>501</sup>Thr, lead to an increased affinity to hACE2 (Starr et al. n.d.). Additionally, this site is located near a linear B cell immunodominant site (B.-Z. Zhang et al. 2020), and therefore the mutation may allow SARS-CoV-2 variants to escape neutralizing antibodies (Figure 4, Figure 5). Indeed, neutralizing antibodies derived from vaccinations and natural infection have significantly reduced activity against pseudotyped viruses carrying this mutation (Wang et al. 2021).

**Table 2. Surface exposed residues of ACE2 orthologues forming the region of contact with site 501 of SARS-CoV-2 S.** Relative to the human sequence, almost all these residues are either conserved (“|”) or replaced by a nearly equivalent amino acid in mouse, American mink, European mink, ferret, and pangolin. Notably, there is a nonconservative substitution of Gly<sup>354</sup> to a bulky positively charged amino acid in most species. Our structural analyses suggests that this substitution contributes to a putative host-dependent selective pressure at site 501 of SARS-CoV-2 S. Prevalent residues reported at this site are informed in order of frequency.

Species	Residues in ACE2									S 501
<i>Homo sapiens</i> (Human)	D38	Y41	Q42	L45	K353	G354	D355	R357	I358	Y, N
<i>Mus musculus</i> (House mouse)					H					Y, N
<i>Neovison vison</i> (American mink)	E					H				N, T
<i>Mustela lutreola</i> (European mink)	E					R				N, T
<i>Mustela putorius furo</i> (Ferret)	E					R				T, N
<i>Manis pentadactyla</i> (Pangolin)	E					H				N, T

Transmission between human and non-human hosts for SARS-CoV-2 provides further information on the evolutionary selectivity of site 501 in S. Repeated infection of mice with human SARS-CoV-2 resulted in the selection of a mouse-adapted strain carrying S Tyr<sup>501</sup> (Gu et al. 2020). It is possible that Asn<sup>501</sup>Tyr results in an additional stabilization of the RBD-ACE2 interaction via  $\pi$ -stacking of Tyr<sup>501</sup> with Tyr<sup>41</sup> in ACE2 (Figure 4a-b). In contrast, several introductions into farmed mink (*Neovison vison*), which caused a substantial increase in their mortality (Oude Munnink et al. 2021), have not led to the same selection. To date, reported sequences in GISAID of SARS-CoV-2 in this host carry either S Asn<sup>501</sup>, which is prevalent, or S Thr<sup>501</sup>, which appeared independently in mink farms (Table S3) (Oude Munnink et al. 2021). In ACE2 of these taxa, Tyr<sup>41</sup> is conserved, but near this site, a larger, positively charged amino acid, His<sup>354</sup>, replaces Gly<sup>354</sup>. Table 2 shows that the amino acids in the RBD 501-binding region of the ACE2 orthologues are conserved, except for Gly<sup>354</sup>, indicating that this site may play a key role in viral fitness.





**Fig. 5. Response to viral infection.** **a.** As part of the innate immune response, (Step 1) the SARS-CoV-2 virus is internalized into endosomes and degraded. (Step 2) viral RNA activates the mitochondrial antiviral innate immunity (MAVS) pathway and (Step 3) degraded proteins activate the toll receptor pathway (TLR3/TLR7), which result in the (Step 4) phosphorylation of TBK1 and translocation of NF-kB and IRF3 to the nucleus, where they regulate the transcription of immune genes including interferons (IFNs, Step 5). IFNs recruit CD8+ T cells that, (Step 6) recognize fragments of the virus on the cell surface via their class I major histocompatibility complex (MHC I) receptors and are activated by dendritic cells (antigen processing cells, or APC). If the virus bypasses innate immunity (orange arrows) nonstructural proteins (nsp6 and nsp13) block the IRF3 nuclear translocation. **b.** APCs recruit B lymphocytes and stimulate the production of antibodies that recognize SARS-CoV-2 S (whereas T cells recognize fragments of S bound to MHC I). **c.** The neutralizing antibodies block binding of the virus to the ACE2 receptor and can prevent re-infection but mutations in the receptor-binding domain (RBD), e.g., S Asn<sup>501</sup>Tyr, prevent binding of the antibodies and the virus is then able to bind the receptor again even if individuals experienced exposure to an earlier strain or were vaccinated. Created with BioRender.com.

Similarly, ferrets (*Mustela putorius furo*) and pangolins (*Manis pentadactyla*), relevant potential reservoirs of SARS-CoV-2, carry a large basic residue at site 354 (an arginine and histidine, respectively). Sawatzki et al. reported that the constant exposure of ferrets to infected humans did not result in natural transmission in a domestic setting, suggesting that ferret

infection may require improved viral fitness (Sawatzki et al. 2021). In agreement with that, Richard et al. (Richard et al. 2020) reported that the adaptive substitution Asn<sup>501</sup>Thr was detected in all experimentally infected ferrets in the laboratory. In order to further investigate the role of the Gly<sup>354</sup> versus Arg<sup>354</sup> in the adaptive mutation of site 501 in S RBD, we performed extensive molecular dynamics simulations of the truncated complexes of Asn<sup>501</sup>-carrying RBD of SARS-CoV-2 and ACE2, from human (hACE2) and ferret (fACE2). The simulations indicate that there is a remarkable difference in the interaction pattern between the two systems in the region surrounding site 501 of RBD. Firstly, we identified the main ACE2 contacts with Asn<sup>501</sup>, which were the same for both species, namely, Tyr<sup>41</sup>, Lys<sup>353</sup>, and Asp<sup>355</sup>, and we also show that the intensity of these contacts is lower in the simulations of fACE2 (Table S4).

To investigate further, we analyzed structural features in the interaction between ACE2 Lys<sup>353</sup> and RBD Asn<sup>501</sup>. Distances between polar atoms computed from the simulations indicate a weaker electrostatic interaction between this pair of residues in ferret compared to human (Figure 4d). This effect is accompanied by a conformational change of fACE2 Lys<sup>353</sup>. Figure S2 shows that, in ferret, the side chain of Lys<sup>353</sup> exhibits more stretched conformations, i.e., a higher population of the *trans* mode of the dihedral angle formed by the side chain carbon atoms. This conformational difference could be partially attributed to the electrostatic repulsion between the two consecutive bulky positively charged amino acids in ferrets, Lys<sup>353</sup> and Arg<sup>354</sup>. Additionally, the simulations suggest a correlation, in a competitive manner, between other interactions that these residues display with the RBD. For example, Figure S3 shows that the salt bridge fACE2\_Arg<sup>354</sup>---RBD\_Asp<sup>405</sup> and the HB interaction fACE2\_Lys<sup>353</sup>---RBD\_Tyr<sup>495</sup> (backbone) alternate in the simulations. This also suggests that the salt bridge formed by fACE2 Arg<sup>354</sup> drags Lys<sup>353</sup> apart from RBD Asn<sup>501</sup>, weakening the interaction between this pair of residues in ferrets



relative to humans.

Altogether, these analyses indicate that site 354 in ACE2 significantly influences the interactions with RBD in the region of site 501 and is likely playing a major role in the selectivity of the size and chemical properties of this residue in SARS-CoV-2. We propose that, in contrast to Tyr<sup>501</sup>, a smaller HB-interacting amino acid at site 501 of RBD, such as the threonine reported in farmed mink and ferrets, may ease the interactions on the region, e.g., the salt bridge between fACE2 Arg<sup>354</sup> and RBD Asp<sup>405</sup>. The differences in the region of ACE2 in contact with site 501 seem to have a key role for host adaptation and are worth further investigation as it may also reveal details of the origin of this zoonotic pandemic.

*S His<sup>69</sup>\_Val<sup>70</sup> deletion* - The His<sup>69</sup>\_Val<sup>70</sup> deletion (in B.1.1.7) is adjacent to a linear epitope at the N-terminal domain of S (Figure 4a,c) (B.-Z. Zhang et al. 2020), suggesting it too may improve fitness by reducing host antibody effectiveness.

*S Leu<sup>452</sup>Arg* - The Leu<sup>452</sup>Arg mutation in S is a core change in the early\_2021\_VOC (Table 1, Figure 4a-b). Although Leu<sup>452</sup> does not interact directly with ACE2, this mutation was shown to moderately increase infectivity in cell cultures and lung organoids using Leu<sup>452</sup>Arg-carrying pseudovirus (Deng et al. 2021). It is possible that the substitution of the leucine, hydrophobic, to arginine, a positively charged residue, creates a direct binding site with ACE2 via the electrostatic interaction with Glu<sup>35</sup>. However, in Starr et al., experiments with the isolated RBD expressed on the cell surface of yeast show that this mutation is associated with enhanced structural stability of RBD, while it only slightly improves ACE2-binding (Starr et al. n.d.). An alternative but not mutually exclusive hypothesis is that it causes a local conformational change

that impacts the complex dynamic interchange between interactions of RBD with the spike trimer itself and with the host receptor. Noteworthy, site 452 resides in a significant conformational epitope in RBD and Leu<sup>452</sup>Arg was shown to decrease binding to neutralizing antibodies (Figure 4b) (Deng et al. 2021; Li et al. 2021).

As noted in Deng et al., S Leu<sup>425</sup>Arg has been reported in rare variants starting in March 2020 from Denmark, i.e., several months before the surge of the VOC that carry this mutation (B.1.427, B.1.429, and B.1.617) (Deng et al. 2021). This indicates that the high transmissibility of the early\_2021\_VOC is not fully explained by the increased infectivity caused by Leu<sup>425</sup>Arg and combined mutations may be essential for the rapid spread. Besides the other mutations in the spike in these VOC, the substitution Asp<sup>260</sup>Tyr in the SARS-CoV-2 helicase (nsp13, below) is especially interesting, as it was identified in the MJN analysis as a defining mutation of both B.1.427 and B.1.429 variants.

*S Ser<sup>477</sup>Asn* - Variants carrying the S Ser<sup>477</sup>Asn mutation spread rapidly in Australia (Figure 1, Figure 3b). This site, located at the loop β4-5 of the RBD, is predicted not to establish persistent interactions with ACE2 (Ali & Vijayan 2020). However, deep scanning shows that this mutation is associated with a slight enhancement of ACE2-binding. Molecular dynamics simulations suggest that Ser<sup>477</sup>Asn affects the local flexibility of the RBD at the ACE2-binding interface, which could be underlying the highest binding affinity with ACE2 reported from potential mean force calculations (A. Singh et al. n.d.). Additionally, this site is located near an epitope and may alter antibody recognition and counteract the host immune response (Figure 4b).

*S Glu<sup>484</sup>Lys* - A recent computational study suggests that Glu<sup>484</sup> exhibits only intermittent interactions with Lys<sup>31</sup> in ACE2 (Ali & Vijayan 2020). Deep scanning shows that this mutation is associated with higher affinity to ACE2 (Starr et al. n.d.) and may be explained by its proximity to Glu<sup>75</sup> in ACE2, which would form a salt bridge with Lys<sup>484</sup>. Aside from the potential impact of Glu<sup>484</sup>Lys between virus-host cell interaction, this site is part of a linear B cell immunodominant site (B.-Z. Zhang et al. 2020) and this mutation was shown to impair antibody neutralization (Wang et al. 2021).

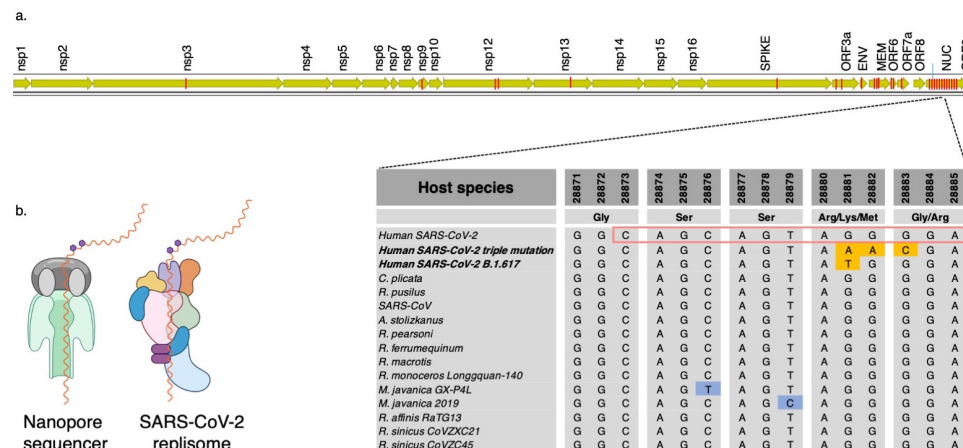
*S Pro<sup>681</sup>Arg and Pro<sup>681</sup>His* - These mutations in the multibasic furin cleavage site are particularly relevant given the importance of this region for cell-cell fusion (Hoffmann et al. 2020; Papa et al. n.d.). The presence of the multibasic motif of SARS-CoV-2 has shown to be essential to the formation of syncytia (i.e., multinucleate fused cells), and thus it is thought to be a key factor underlying pathogenicity and virulence differences between SARS-CoV-2 and other related betacoronaviruses. Hoffmann et al. recognized the importance of the furin cleavage site in SARS-CoV-2 and its biochemically basic signature and generated mutants to determine the effects of specific amino acids. Notably, they showed that pseudotyped virion particles bearing mutant SARS-CoV-2 S with additional basic residues in this region, including the substitution Pro<sup>681</sup>Arg (present in B.1.617), exhibits a remarkable increase in syncytium formation in lung cells *in vitro* (Hoffmann et al. 2020), which may explain the increased severity of the disease (Sheikh et al. 2021). Hoffman et al did not include a Pro<sup>681</sup>His change that is a defining mutation of B.1.1.7, and so it is not known if this too increases syncytium formation given that it is a basic amino acid, but should be the target of future studies.

*Nucleocapsid Arg<sup>203</sup>Met* - The main function of the nucleocapsid (N) protein in SARS-CoV-2 is to act as a scaffold for the viral genome and it is also the most antigenic protein produced by the virus (Dutta et al. 2020). In a previous study, we reported that the Ser-Arg-rich motif of this protein (a.a. 183-206), shown *in vitro* to be necessary for viral replication (Garvin, Prates, et al. 2020; Tylor et al. 2009), displays a high number of amino acid changes during the COVID-19 pandemic and is likely under positive selection. We propose that the RNA gene segment coding this particular subsequence may be linked to improved fitness of specific SARS-CoV-2 haplotypes including the rapidly spreading delta variant and is linked to epigenetic alterations (Figure 6a).

A recent deep transcriptome sequencing study used Oxford Nanopore™ technology to detect epigenetic modifications at 41 sites in the RNA genome that are associated with leader sequence addition to sub-genomic RNA transcripts, a recombination-like process of SARS-CoV-2 (Kim et al. 2020). Nanopore instruments can detect epigenetic modifications based on disruptions of the electrical current as the RNA molecule passes through the molecular pore (Rand et al. 2017; Simpson et al. 2017), which Kim et al propose is responsible for the pause that occurs before leader sequence addition. Twenty-five of the 41 modified sites reside in the N gene and the majority of the sites in this subset are found near the Ser-Arg-rich motif (Figure 6a).

Furthermore, one specific epigenetic site is linked to two highly successful SARS-CoV-2 haplotypes. The first is a triple mutation at sites 28881-28883 (GGG to AAC, Arg<sup>203</sup>Lys) that is now found in nearly half of all sequences sampled across the globe (diamond nodes, Figure 1) and the second is Arg<sup>203</sup>Met, which is a defining mutation for the rapidly spreading B.1.617. Notably, this region of the genome is highly conserved across several hundred years of coronavirus evolution (Boni et al. 2020) (Figure 6a). Given that these epigenetic sites were

discovered because the RNA pauses as it passes across the pore of the molecular nanopore sequencer, one interesting hypothesis is that mutations at this region remove the epigenetic modification and speed the SARS-CoV-2 genome through the replisome (Figure 6b), increasing the production of virions, which is consistent with the more than 1000-fold higher virion count in those infected with B.1.617 (Lu et al. n.d.).

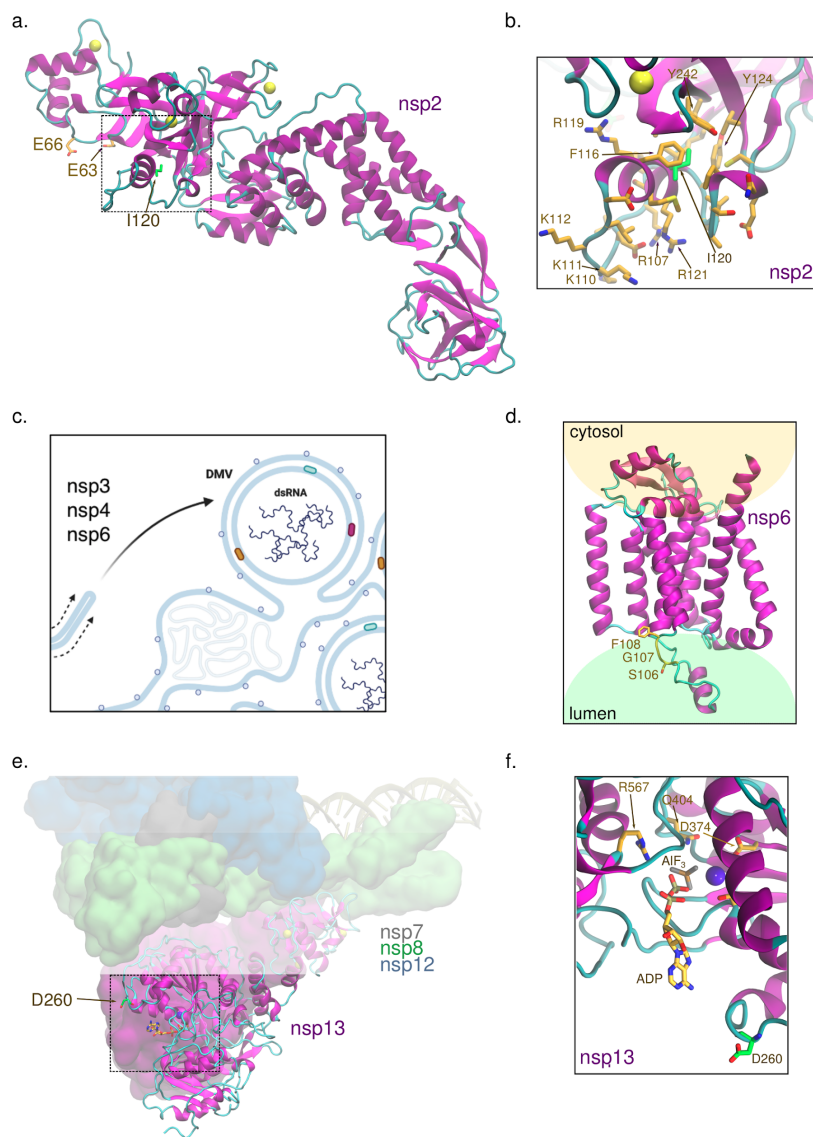


**Fig. 6. Modifications at the Ser-Arg-rich region of N may affect replication speed.** **a.** Location of 41 epigenetic sites reported in Kim et al. 2020 (red bars on SARS-CoV-2 genome). One of the sites in the nucleocapsid gene (nucleotides in red box of aligned sequences) is highly conserved across diverse host-defined coronaviruses. All bats and human coronavirus species from China are completely conserved at the epigenetic site 28881-28883, except for a 3-bp mutation in SARS-CoV-2 that occurred early in the pandemic and now corresponds to ~50% of all sequences globally (diamond nodes in Figure 1). **b.** Kim et al. proposed that *N*<sup>6</sup>-methyladenosine modification of the genome (purple hexagons), common in RNA viruses, caused the strand to pause while traversing the nanopore sequencing apparatus. We propose that loss of this site via mutations at site 203 in N may increase the replication rate of the RNA strand through the SARS-CoV-2 replisome. *Aselliscus stoliczkanus* - Stoliczka's trident bat, *Chaerephon plicata* - wrinkle-lipped free-tailed bat, *Rhinolophus pusillus* - least horseshoe bat, *R. pearsoni* - Pearson's horseshoe bat, *R. macrotis* - big-eared horseshoe bat, *R. ferrumequinum* - greater horseshoe bat, *R. monaceros* - Formosan lesser horseshoe bat, *R. affinis* intermediate horseshoe bat, *R. sinicus* Chinese rufous horseshoe bat, *R. mayalanis* - Mayalan horseshoe bat, *SARS* - Severe Acute Respiratory Syndrome, *Manis javanica* - Malayan pangolin. Created with BioRender.com.

*nsp2 Ile<sup>120</sup>Phe* - The main role of the nonstructural protein 2 (nsp2) in viral performance is not yet defined. Instead, this protein appears to be part of multiple interactions with host proteins involved in a range of processes including the regulation of mitochondrial respiratory function,

endosomal transport, and ribosome biogenesis (Verba et al. 2021). Very recently, deep learning-based methods of structure prediction and cryo-electron microscopy density were combined to provide the atomic model of nsp2 (PDB id 7MSW). In a preprint from Verba et al., structural information was used to localize the surfaces that are key for protein-protein interaction with nsp2 (Verba et al. 2021). From structural analysis and mass spectrometry experiments, the authors pose the interesting hypothesis that nsp2 interacts directly with ribosomal RNA via a highly conserved zinc ribbon motif to bring ribosomes close to the replication-transcription complexes.

Here we are particularly interested in the functional impact of the mutation Ile<sup>120</sup>Phe in nsp2 present in the D.2 variant. Site 120, identified in the nsp2 structure on Figure 7a, is a point of hydrophobic contact between a small helix, rich in positively charged residues, and a zinc binding site. The positively charged surface of the helix may be especially relevant for a putative interaction with the phosphate groups from ribosomal RNA. Normal mode analysis from DynaMut2 predicts that the substitution has a destabilizing effect in the protein structure (estimated  $\Delta\Delta G^{\text{stability}} = -2$  kcal/mol) (Rodrigues et al. 2021). Possibly, this could be caused by  $\pi$ -stacking interactions of the tyrosine with aromatic residues in the same helix that would disrupt the contacts anchoring it to the protein core (Figure 7b). Additionally, site 120 is spatially close to Glu<sup>63</sup> and Glu<sup>66</sup>, which were shown to be relevant for interactions with the endosomal/actin machinery via affinity purification mass spectrometry in HEK293T cells. Remarkably, upon mutation of these glutamates to lysines, there is increased interactions with proteins involved in ribosome biogenesis (Verba et al. 2021).



**Fig. 7. Location of mutations of prevalent SARS-CoV-2 variants on the structure of the nonstructural proteins nsp2, nsp6, and nsp13.** **a.** Site 120 in nsp2 is located in a small helix near a zinc-binding site and residues Glu<sup>63</sup> and Glu<sup>66</sup>, which play a role in the interaction with proteins involved in ribosome biogenesis and in the endosomal/actin machinery (Verba et al. 2021). PDB id 7MSW was used. **b.** Ile<sup>120</sup> forms some of the hydrophobic contacts that anchor the helix at the surface of nsp2, where this site resides, to the protein core. **c.** Nsp6 participates in generating double-membrane vesicles (DMV) for viral genome replication. Natural selection for the biological traits of viral entry and replication may explain the increased transmission of variants with adaptive mutations in both S and nsp6. DMVs isolate the viral genome from host cell attack to provide for efficient genome and sub-genome replication and generate virions. **d.** Sites 106-108 are predicted to be located at/near the protein region of nsp6 embedded in the endoplasmic reticulum lumen (structure generated by AlphaFold2 (Jumper et al. n.d.)). **e.** Nsp13 is the SARS-CoV-2 helicase and it is part of the replication complex. **f.** Asp<sup>260</sup> in nsp13 is mutated to tyrosine in B.1.427 and B.1.429 and it is located at the entrance of the NTP-binding site. PDB id 6XEZ was used (Chen et al. 2020).



*nsp6 Ser<sup>106</sup>\_Phe<sup>108</sup>deletion* - The nsp6 protein plays critical roles in viral replication and suppression of the host immune response (Figure 5a and Figure 7c) (Gupta et al. 2020). Along with nsp3 and nsp4, nsp6 is responsible for producing double-membrane vesicles from the endoplasmic reticulum (ER) to protect the viral RNA from host attack and increase replication efficiency (Figure 7c) (Santerre et al. 2020). The nsp6 Ser<sup>106</sup>\_Phe<sup>108</sup>del is predicted to be located at a loop in the interface between a transmembrane helix and the ER lumen based on a preliminary structural analysis of the model generated by the AlphaFold2 system (Figure 7d), and we hypothesize that the deletion may affect functional interactions of nsp6 with other proteins. In addition, in agreement with the enhanced suppression of innate immune response reported for B.1.1.7 (Thorne et al. 2021), changes in immune-antagonists, such as nsp6 Ser<sup>106</sup>\_Phe<sup>108</sup>del, may be key to prolonged viral shedding (Calistri et al. 2021).

*nsp13 Asp<sup>260</sup>Tyr* - The nonstructural protein 13 is a component of the viral replication-transcription complex, (nsp13; or SARS-CoV-2 helicase) and plays an essential role in unwinding the duplex oligonucleotides into single strands in a NTP-dependent manner (L. Yan et al. 2020). Hydrogen/deuterium exchange mass spectrometry demonstrates that the helicase and NTPase activities of SARS-CoV nsp13 are highly coordinated, and mutations at the NTPase active site impair both ATP hydrolysis and the unwinding process (Jia et al. 2019). Here we note that the substitution Asp<sup>260</sup>Tyr, present in B.1.427 and B.1.429, is located at the entrance of the NTPase active site and may favor  $\pi$ - $\pi$  stacking interactions with nucleobases (Figure 7e-f). Given that, at high ATP concentrations, SARS-CoV nsp13 exhibits increased helicase activity on duplex RNA (Jang et al. 2020), it is possible that, similarly, the putative optimization on NPT uptake in nsp13 Asp<sup>260</sup>Tyr favors RNA unwinding.



Additionally, nsp13 was shown to play an important role as an innate immune antagonist (Figure 5a). It contributes to the inhibition of the type I interferon response by directly binding to TBK1 and, with that, it impedes IRF3 phosphorylation (Guo et al. 2021). The dual role of nsp6 and nsp13 in immune suppression and viral replication may suggest a convergent evolution of SARS-CoV-2 manifested in most of the VOC, which carries either nsp6 Ser<sup>106</sup>\_Phe<sup>108</sup>del or nsp13 Asp<sup>260</sup>Tyr.

### 3. Concluding Remarks

From our thorough analysis of the spatiotemporal relationships of SARS-CoV-2 variants, we propose that the rapid increase of mutations in the late 2020 VOC is likely a consequence of the recombination of haplotypes carrying adaptive mutations in S and in non-S proteins that act cooperatively to enhance viral fitness. For example, as indicative of that, we call attention to five mutations that occur independently in disjoint clusters of our MJN, four of which (S Asn<sup>501</sup>Tyr, S His<sup>69</sup>\_Val<sup>70</sup>del, S Phe<sup>681</sup>His/Arg, and nsp6 Ser<sup>106</sup>\_Phe<sup>108</sup>del) are shared by different VOC, including B.1.1.7. Notably, S His<sup>69</sup>\_Val<sup>70</sup>del appeared in human and mink populations simultaneously in August 2020, prior to the emergence of B.1.1.7, indicating that mink should be further investigated as a possible component of a recombination event. In turn, our molecular dynamics simulations indicate that the molecular forces at site 501 in S and how they are altered upon mutation (S Asn<sup>501</sup>Tyr in B.1.1.7) are a key component to describe the history of transmission among other putative zoonotic reservoirs, such as farmed minks, ferrets, and pangolins.

The S Asp<sup>614</sup>Gly mutation has been shown to increase infectivity and is now predominant in the circulating virus (L. Zhang et al. 2020), and S Asn<sup>501</sup>Tyr is associated with higher

virulence (Gu et al. 2020). We show that the expansion of the strains carrying these mutations only occurred upon the additional substitutions in nsp12 Leu<sup>323</sup>Pro (Figure 2b) (Garvin, Prates, et al. 2020) and nsp6 Ser<sup>106</sup>\_Phe<sup>108</sup>del, respectively. A hypothesis consistent with these observations is that the changes in S enhance viral entry into the host cells but they do not easily transmit due to rapid suppression by a robust innate immune response. A secondary mutation is able to counteract the immune-driven suppression. In the case of S Asp<sup>614</sup>Gly, the nsp12 Leu<sup>323</sup>Pro may have increased the replication rate of the virus, which was supported by quantitative PCR from clinical samples with different viral strains in Korber et al. (Garvin, Prates, et al. 2020; Korber et al. 2020). However, the separate effects of S Asp<sup>614</sup>Gly and nsp12 Leu<sup>323</sup>Pro could not be described in the referred study because it did not include individuals infected with variants harboring only one of the mutations.

For the late 2020 VOC, nsp6 Ser<sup>106</sup>\_Phe<sup>108</sup>del may affect viral replication in DMVs or suppress the interferon-driven antiviral response (Xia et al. 2020). It is likely that other mutations also enhance viral mechanisms that impair the host immune response. For example, Thorne et al. recently showed that the B.1.1.7 VOC suppresses the innate immune response by host cells *in vitro* and attributed it to the increased transcription of the *orf9b* gene, nested within the gene coding the nucleocapsid protein. (Thorne et al. 2021), although they could not rule out the possibility that this was due to nsp6 Ser<sup>106</sup>\_Phe<sup>108</sup>del.

Via focused protein structural analysis, we identify other mutations shared among different VOC that reside in key locations of proteins involved in viral replication and/or in suppressing the innate immune suppression, such as nsp13, suggesting a convergent evolution of SARS-CoV-2. This emphasizes the importance of tracking mutations in a genome-wide manner as a strategy to avoid the emergence of future VOC. For example, an earlier dominant variant in

Australia (D.2) that carried the mutations Ser<sup>477</sup>Asn in S and Ile<sup>120</sup>Phe in nsp6 was successfully restrained. However, we note that variants harboring only the S Ser<sup>477</sup>Asn substitution are currently circulating in several European countries (Figure 2, Table S5) and may only need to recombine with a variant carrying an advantageous complementary mutation to become the next VOC.

A second and equally significant outcome from recombination-driven haplotypes is the generation of variants that allow escape from neutralizing antibodies produced by an adaptive immune response (Garvin, T Prates, et al. 2020) (Figure 5c). As a case in point, the resurgence of COVID-19 in Manaus, Brazil, in January 2021, where seroprevalence was above 75% in October 2020, is due to immune escape of new SARS-CoV-2 lineages (Sabino et al. 2021). Broad disease prevalence and community spread of COVID-19 increase the probability that divergent haplotypes may come in contact, thereby dramatically accelerating the evolution and transmission of the virus. This emphasizes that regions with low sequence surveillance can be viral breeding grounds for the next SARS-CoV-2 VOC.

## 4. Methods

### *Sequence data pre-processing*

We downloaded SARS-CoV-2 sequences in FASTA format and corresponding metadata from GISAID and processed as we have reported previously (Garvin, Prates, et al. 2020; E. Prates et al. 2021). To ensure that deletions were accounted for, full genome sequences were aligned with MAFFT (Katoh et al. 2002) to the established reference genome (accession NC\_045512), uploaded into CLC Genomics Workbench, and trimmed to the start and stop codons (nsp1 start

site and ORF10 stop codon). Aligned sequences in tab-delimited format were imported into R to count the number of variable accessions at each of the 29,409 sites.

Variable sites were determined with all sequences downloaded up through the end of January, 2021. In order to reduce false-positive mutation sites (those that were due to technical error), we selected sites that were variable in 25 or more individuals (0.01%) compared to the reference (all 25 were required to be the same state: A, G, T, C, or -). We further pruned these by removing sites in which 20% or more of the accessions harbored an unknown character state (“N”), leaving 2,128 variable sites for downstream analyses. After removing sequences with an “N” at any of these sites, we retained 280,409 individuals. Prior to submission, we updated the number of sequences through April 19, 2021, keeping the same 2128 variable sites, which allowed us to capture the most up-to date metadata and produced 640,211 for analysis. We kept haplotypes that occurred in more than 35 individuals to remove rare or artifact-derived haplotypes. For the comparison of median-joining networks and phylogenetic trees, we used sequences from the pandemic sampled through the end of April, 2020. We used variable sites found in more than ten individuals and haplotypes found in five or more individuals as we had in previous work (Garvin, Prates, et al. 2020). This produced 410 unique haplotypes based on 467 variable sites.

# *Median-joining network (MJN)*

Haplotypes were coded in NEXUS format and uploaded to PopArt (Leigh & Bryant 2015). An MJN was produced with the epsilon parameter set to 0. The networks were exported as a table and visualized in Cytoscape (Shannon et al. 2003) with corresponding metadata. The date of emergence of each haplotype was defined by the sample date subtracted from the report date for

the Wuhan reference sequence (December 24, 2019) and then one day was added to remove zeros. For samples that only reported the month but no day, we recorded the day as the 15th of that month. We excluded samples with no sampling date.

### *Phylogenetic tree*

We used the program MrBayes to generate a phylogenetic tree (Ronquist & Huelsenbeck 2003). Parameters were set to *Nucmodel=4by4*, *Nst=6*, *Code=Universal*, and *Rates=Invgamma*. We performed 5,000,000 mcmc generations, which produced a stable standard deviation of split frequencies of 0.014. A consensus tree was generated using the 50% majority rule and visualized using FigTree v1.4.4 (<http://tree.bio.ed.ac.uk/software/figtree/>).

### *Estimation of genome mutation load*

We estimated the mutation load using two data sets. First, we used the 640,211 sequences based on 2,128 variable sites used for the MJN because these represent high-confidence mutations. For each of the 640,211 accessions, we counted the number of differences of the 2,128 variable sites compared to the reference genome (accession NC\_045512) and recorded the day of emergence. The mutational load for all accessions for a given day was then averaged and this was plotted across time. For the second estimate of mutation rate, we used all variable sites across the full genome (29,409 sites) to include rare variants and removed all sequences with at least one ambiguous site, leaving 584,119 accessions.

For the population-level estimate of mutation accumulation, we applied the filters used to identify the 2,128 variable sites that were used for the MJN for all sequences up through April

19, 2021. We did not include new mutations because the B.1.1.7 VOC and its downstream haplotypes had become the predominant variants globally at that time and, consequently, much early information of the molecular evolution is lost when applying frequency filters on the entire GISAID database. This is exacerbated with the MJN approach because the software algorithm used to generate the network is computationally intractable with greater than 1,000 haplotypes and therefore future efforts will either need to ignore early molecular events or use new methods that can handle the large datasets and any recombination events that occur (an alternative approach would be to now use the alpha or delta variant as the reference sequence because they are now the predominant strains globally).

For calculations of population-level mutation accumulation, it is possible (and necessary) to include all sequences to determine if mutation or recombination are the cause of the high mutation load seen in the late 2020 VOC. After applying the frequency and haplotype filters, we retained 5,011 variable sites that define 12,282 unique haplotypes for further analysis. Mutations to five possible states (A, G, T, C, and -) were counted at each site on the first date that they appeared and their appearance at later dates were excluded. Multiple mutations at a site to different states were counted with this method.

For lineage-specific mutation curves, we extracted all sequences based on their PANGO lineage listed in the metadata from GISAID that also had a sample data and plotted the cumulative number over time, where time is represented by days from first appearance. To estimate the rate of accumulation, we calculated the slope for the linear portion of each of the curves.

## *Probability of mutation accumulation*

To calculate the chance of accumulating several mutations in a certain period, the probability density function for a normal distribution is used:

$$PDF(x) = \exp(-(x - \mu)^2 / 2\sigma^2) / \sqrt{2\pi * \sigma^2},$$

where  $\mu$  is the expected number of mutations for that date,  $x$  is the measured value, and  $\sigma$  is the standard deviation of error calculated from the data shown in Fig. 1b, considering the difference between the actual and predicted number of mutations. The expected value of mutations  $\mu$  for a given time period is computed from the estimated rate of mutations per day (Figure 3, 0.05). The period of interest to our discussion (June-October 2020) corresponds to 122 days, for which, the integral of  $PDF(x=13)$  gives the probability of  $1 \times 10^{-15}$  to accumulate 13 mutational events.

## *Screen for coinfecting individuals with UK B.1.1.7*

We extracted 25 samples from the Sequence Read Archive at NCBI for each of the months of October, November, December, and January listed as variant B.1.1.7 from the UK (Table S2) for a total of 100 samples to check for coinfection. The reads were mapped to the NC\_045512 Wuhan reference using CLC Genomics Workbench using the default parameters except for length fraction and similarity fraction were set to 0.9. Three sites specific to UK B.1.1.7 were analyzed for possible heterozygosity. Of the 100 we sampled, two appeared to be cases of coinfection. This supports the hypothesis that the large expansion in overall mutations seen in UK B.1.1.7 are likely due to recombination. In addition, it also supports the case that coinfection is occurring at a baseline sufficient to allow for occasional recombination.

## *Protein structure analysis*

VMD was used to visualize the protein structures and analyze the potential functional effects of mutations (Humphrey et al. 1996). Figure 3 was created using Inkscape (<https://inkscape.org/>) and Gimp 2.8 (<https://www.gimp.org>) ('The GIMP Development Team. (2019). GIMP. Retrieved from <https://www.gimp.org>' n.d.).

## *Molecular dynamics simulations*

Molecular dynamics (MD) simulations were used to study interactions between SARS-CoV-2 RBD and ACE2 from ferret and human. Three independent extensive MD simulations were performed for each species using GROMACS 2020 package (Lindahl et al. 2020) and the CHARMM36 force field for protein and glycans (Guvench et al. 2011; J. Huang & MacKerell 2013). Each simulation ran up to 800 ns, being the last 500 ns used for analysis. PDB id 6M17 was used to build the ACE2-RBD complexes. Given the high sequence identity between human and ferret ACE2 (83%), we performed local modeling of the non-conserved amino acid residues in ferret ACE2 using the human homolog as the template, via RosettaRemodel (P.-S. Huang et al. 2011).

The inputs for simulations were generated using CHARMM-GUI (Jo et al. 2008). Counterions were added for electroneutrality (0.1 M NaCl). The complexes were surrounded by TIP3P water molecules to form a layer of at least 10 Å relative to the box borders (Jorgensen et al. 1983). Simulations were performed using the NPT ensemble. The temperature was maintained at 310 K with the Nosé–Hoover thermostat using a time constant of 1.0 ps (Evans & Holian 1985). The pressure was maintained at 1 bar with the isotropic Parrinello–Rahman barostat using a compressibility of  $4.5 \times 10^{-5} \text{ bar}^{-1}$  and a time constant of 1.0 ps in a rectangular



simulation box (Parrinello & Rahman 1981). The particle mesh Ewald method was used for the treatment of periodic electrostatic interactions with a cutoff distance of 1.2 nm (Darden et al. 1993). The Lennard–Jones potential was smoothed over the cutoff range of 1.0–1.2 nm by using the force-based switching function. Only atoms in the Verlet pair list with a cutoff range reassigned every 20 steps were considered. The LINCS algorithm was used to constrain all bonds involving hydrogen atoms to allow the use of a 2 fs time step (Hess et al. 1997). The suggested protocol for nonbonded interactions with the CHARMM36 force field when used in the GROMACS suite was followed.

The Hbonds plugin in VMD was used to identify hydrogen bond interactions along the simulations (Humphrey et al. 1996). The geometric criteria adopted are a cutoff of 3.5 Å for donor-acceptor distance and 30° for acceptor-donor-H angle. The Timeline plugin was used to count contacts formed by a given amino acid residue. We defined the distance of 4 Å between any atom pairs as the cutoff for contact.

## 5. Data Access

All SARS-CoV-2 sequences used in this study are available from the public repositories Genome Initiative on Sharing Avian Influenza Data (GISAID, [gisaid.org](https://gisaid.org)), the National Center for Biotechnology Information (NCBI, <https://www.ncbi.nlm.nih.gov/sars-cov-2/>) and the COVID-19 Genomics UK Consortium (COG, <https://www.sanger.ac.uk/collaboration/covid-19-genomics-uk-cog-uk-consortium/>)

## 6. Acknowledgments

The viral evolution research was funded by the Laboratory Directed Research and Development Program of Oak Ridge National Laboratory, managed by UT-Battelle, LLC for the US Department of Energy (LOIS:10074) and the structural implication work was funded via the DOE Office of Science through the National Virtual Biotechnology Laboratory (NVBL), a consortium of DOE national laboratories focused on the response to COVID-19, with funding provided by the Coronavirus CARES Act. This work was also funded by the United States Government. This research used resources of the Oak Ridge Leadership Computing Facility (OLCF) and the Compute and Data Environment for Science (CADES) at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725. Figures generated with Biorender and VMD. We gratefully acknowledge the Originating laboratories responsible for obtaining the viral specimens and the Submitting laboratories where genetic sequence data were generated and shared via the GISAID Initiative, on which this research is based.

## *Author Contribution*

**MR Garvin:** Conceptualization, Data curation, Funding acquisition, Formal Analysis, Investigation, Methodology, Visualization, Writing – original draft, Writing – review & editing.  
**ET Prates:** Formal Analysis, Investigation, Visualization, Writing - original draft, Writing - review & editing  
**J Romero:** Conceptualization, Formal Analysis, Investigation, Methodology, Software, Writing – original draft, Writing – review & editing.  
**A Cliff:** Methodology, Software, Writing – review & editing  
**JGFM Gazolla:** Software, Formal Analysis, Investigation, Data Curation, Visualization, Writing - Review and Editing.  
**M Pickholz:** Investigation, Visualization, Writing – original draft, Writing – review & editing  
**M Pavicic:** Investigation, Writing – original draft, Writing – review & editing  
**DA Jacobson:** Conceptualization, Funding acquisition, Formal Analysis, Investigation, Project

administration, Supervision, Resources, Writing – original draft, Writing – review & editing

## 7. References

- Ali, A., & Vijayan, R. (2020). ‘Dynamics of the ACE2-SARS-CoV-2/SARS-CoV spike protein interface reveal unique mechanisms’, *Scientific reports*, 10/1: 14214.
- Alpert, T., Brito, A. F., Lasek-Nesselquist, E., Rothman, J., Valesano, A. L., MacKay, M. J., Petrone, M. E., et al. (2021). ‘Early introductions and transmission of SARS-CoV-2 variant B.1.1.7 in the United States’, *Cell*, 184/10: 2595–604.e13.
- Anisimova, M., Nielsen, R., & Yang, Z. (2003). ‘Effect of recombination on the accuracy of the likelihood method for detecting positive selection at amino acid sites’, *Genetics*, 164/3: 1229–36.
- Bandelt, H. J., Forster, P., & Röhl, A. (1999). ‘Median-joining networks for inferring intraspecific phylogenies’, *Molecular biology and evolution*, 16/1: 37–48.
- Bentley, K., & Evans, D. J. (2018). ‘Mechanisms and consequences of positive-strand RNA virus recombination’, *The Journal of general virology*, 99/10: 1345–56.
- Boni, M. F., Lemey, P., Jiang, X., Lam, T. T.-Y., Perry, B. W., Castoe, T. A., Rambaut, A., et al. (2020). ‘Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible for the COVID-19 pandemic’, *Nature microbiology*, 5/11: 1408–17.
- Calistri, P., Amato, L., Puglia, I., Cito, F., Di Giuseppe, A., Danzetta, M. L., Morelli, D., et al. (2021). ‘Infection sustained by lineage B.1.1.7 of SARS-CoV-2 is characterised by longer persistence and higher viral RNA loads in nasopharyngeal swabs’, *International journal of infectious diseases: IJID: official publication of the International Society for Infectious Diseases*, 105: 753–5.
- Casalino, L., Gaieb, Z., Goldsmith, J. A., Hjorth, C. K., Dommer, A. C., Harbison, A. M., Fogarty, C. A., et al. (2020). ‘Beyond Shielding: The Roles of Glycans in the SARS-CoV-2 Spike Protein’, *ACS central science*, 6/10: 1722–34.
- Challen, R., Brooks-Pollock, E., Read, J. M., Dyson, L., Tsaneva-Atanasova, K., & Danon, L. (2021). ‘Risk of mortality in patients infected with SARS-CoV-2 variant of concern 202012/1: matched cohort study’, *BMJ*, 372: n579.
- Chen, J., Malone, B., Llewellyn, E., Grasso, M., Shelton, P. M. M., Olinares, P. D. B., Maruthi, K., et al. (2020). ‘Structural Basis for Helicase-Polymerase Coupling in the SARS-CoV-2 Replication-Transcription Complex’, *Cell*, 182/6: 1560–73.e13.
- Darden, T., York, D., & Pedersen, L. (1993). ‘Particle mesh Ewald: An N·log(N) method for Ewald sums in large systems’, *The Journal of chemical physics*, 98/12: 10089–92. American Institute of Physics.
- Davies, N. G., Jarvis, C. I., John Edmunds, W., Jewell, N. P., Diaz-Ordaz, K., Keogh, R. H., & CMMID COVID-19 Working Group. (n.d.). ‘Increased mortality in community-tested cases of SARS-CoV-2 lineage B.1.1.7’. DOI: 10.1101/2021.02.01.21250959
- Deng, X., Garcia-Knight, M. A., Khalid, M. M., Servellita, V., Wang, C., Morris, M. K., Sotomayor-González, A., et al. (2021). ‘Transmission, infectivity, and neutralization of a spike L452R SARS-CoV-2 variant’, *Cell*, 184/13: 3426–37.e8.
- Dutta, N. K., Mazumdar, K., & Gordy, J. T. (2020). ‘The Nucleocapsid Protein of SARS-CoV-2: a Target for Vaccine Development’. *Journal of Virology*. DOI: 10.1128/jvi.00647-20
- Evans, D. J., & Holian, B. L. (1985). ‘The Nose–Hoover thermostat’, *The Journal of chemical*

*physics*, 83/8: 4069–74. American Institute of Physics.

Faria, N. R., Mellan, T. A., Whittaker, C., Claro, I. M., Candido, D. da S., Mishra, S., Crispim, M. A. E., et al. (2021). ‘Genomics and epidemiology of the P.1 SARS-CoV-2 lineage in Manaus, Brazil’, *Science*, 372/6544: 815–21.

Fravev, F. (n.d.). ‘The N501Y and K417N mutations in the spike protein of SARS-CoV-2 alter the interactions with both hACE2 and human derived antibody: A Free energy of perturbation study’. DOI: 10.1101/2020.12.23.424283

Funk, T., Pharris, A., Spiteri, G., Bundle, N., Melidou, A., Carr, M., Gonzalez, G., et al. (2021). ‘Characteristics of SARS-CoV-2 variants of concern B.1.1.7, B.1.351 or P.1: data from seven EU/EEA countries, weeks 38/2020 to 10/2021’, *Euro surveillance: bulletin Europeen sur les maladies transmissibles = European communicable disease bulletin*, 26/16. DOI: 10.2807/1560-7917.ES.2021.26.16.2100348

Garvin, M. R., Prates, E. T., Pavicic, M., Jones, P., Amos, B. K., Geiger, A., Shah, M., et al. (2020). ‘Potentially adaptive SARS-CoV-2 mutations discovered with novel spatiotemporal and explainable-AI models’, *Genome biology*, in press.

Garvin, M. R., T Prates, E., Pavicic, M., Jones, P., Amos, B. K., Geiger, A., Shah, M. B., et al. (2020). ‘Potentially adaptive SARS-CoV-2 mutations discovered with novel spatiotemporal and explainable AI models’, *Genome biology*, 21/1: 304.

Greaney, A. J., Starr, T. N., Gilchuk, P., Zost, S. J., Binshtein, E., Loes, A. N., Hilton, S. K., et al. (n.d.). ‘Complete mapping of mutations to the SARS-CoV-2 spike receptor-binding domain that escape antibody recognition’. DOI: 10.1101/2020.09.10.292078

Gribble, J., Stevens, L. J., Agostini, M. L., Anderson-Daniels, J., Chappell, J. D., Lu, X., Pruijssers, A. J., et al. (2021). ‘The coronavirus proofreading exoribonuclease mediates extensive viral recombination’, *PLoS pathogens*, 17/1: e1009226.

Gu, H., Chen, Q., Yang, G., He, L., Fan, H., Deng, Y.-Q., Wang, Y., et al. (2020). ‘Adaptation of SARS-CoV-2 in BALB/c mice for testing vaccine efficacy’, *Science*, 369/6511: 1603–7.

Guo, G., Gao, M., Gao, X., Zhu, B., Huang, J., Luo, K., Zhang, Y., et al. (2021). ‘SARS-CoV-2 non-structural protein 13 (nsp13) hijacks host deubiquitinase USP13 and counteracts host antiviral immune response’, *Signal transduction and targeted therapy*, 6/1: 119.

Gupta, R., Charron, J., Stenger, C. L., Painter, J., Steward, H., Cook, T. W., Faber, W., et al. (2020). ‘SARS-CoV-2 (COVID-19) structural and evolutionary dynamicome: Insights into functional evolution and human genomics’, *The Journal of biological chemistry*, 295/33: 11742–53.

Guvench, O., Mallajosyula, S. S., Raman, E. P., Hatcher, E., Vanommeslaeghe, K., Foster, T. J., Jamison, F. W., 2nd, et al. (2011). ‘CHARMM additive all-atom force field for carbohydrate derivatives and its utility in polysaccharide and carbohydrate-protein modeling’, *Journal of chemical theory and computation*, 7/10: 3162–80.

Hess, B., Bekker, H., Berendsen, H. J. C., & Fraaije, J. G. E. M. (1997). ‘LINCS: A linear constraint solver for molecular simulations’, *Journal of computational chemistry*, 18/12: 1463–72. Wiley.

Hoffmann, M., Kleine-Weber, H., & Pöhlmann, S. (2020). ‘A Multibasic Cleavage Site in the Spike Protein of SARS-CoV-2 Is Essential for Infection of Human Lung Cells’, *Molecular cell*, 78/4: 779–84.e5.

Huang, J., & MacKerell, A. D., Jr. (2013). ‘CHARMM36 all-atom additive protein force field: validation based on comparison to NMR data’, *Journal of computational chemistry*, 34/25: 2135–45.

Huang, P.-S., Ban, Y.-E. A., Richter, F., Andre, I., Vernon, R., Schief, W. R., & Baker, D. (2011). 'RosettaRemodel: a generalized framework for flexible backbone protein design', *PloS one*, 6/8: e24109.

Humphrey, W., Dalke, A., & Schulten, K. (1996). 'VMD: visual molecular dynamics', *Journal of molecular graphics*, 14/1: 33–8, 27–8.

Huson, D. H., & Bryant, D. (2006). 'Application of phylogenetic networks in evolutionary studies', *Molecular biology and evolution*, 23/2: 254–67.

Jackson, B., Boni, M. F., Bull, M. J., Collier, A., Colquhoun, R. M., Darby, A., Haldenby, S., et al. (n.d.). 'Generation and transmission of inter-lineage recombinants in the SARS-CoV-2 pandemic'. DOI: 10.1101/2021.06.18.21258689

Jang, K.-J., Jeong, S., Kang, D. Y., Sp, N., Yang, Y. M., & Kim, D.-E. (2020). 'A high ATP concentration enhances the cooperative translocation of the SARS coronavirus helicase nsP13 in the unwinding of duplex RNA', *Scientific reports*, 10/1: 4481.

Jia, Z., Yan, L., Ren, Z., Wu, L., Wang, J., Guo, J., Zheng, L., et al. (2019). 'Delicate structural coordination of the Severe Acute Respiratory Syndrome coronavirus Nsp13 upon ATP hydrolysis', *Nucleic acids research*, 47/12: 6538–50.

Jorgensen, W. L., Chandrasekhar, J., Madura, J. D., Impey, R. W., & Klein, M. L. (1983). 'Comparison of simple potential functions for simulating liquid water', *The Journal of chemical physics*, 79/2: 926–35. American Institute of Physics.

Jo, S., Kim, T., Iyer, V. G., & Im, W. (2008). 'CHARMM-GUI: a web-based graphical user interface for CHARMM', *Journal of computational chemistry*, 29/11: 1859–65.

Jumper, J., Tunyasuvunakool, K., Kohli, P., Hassabis, D., & AlphaFold Team. (n.d.). 'Computational predictions of protein structures associated with COVID-19'. *Deep Mind*. Retrieved from <<https://deepmind.com/research/open-source/computational-predictions-of-protein-structures-associated-with-COVID-19>>

Katoh, K., Misawa, K., Kuma, K.-I., & Miyata, T. (2002). 'MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform', *Nucleic acids research*, 30/14: 3059–66.

Kim, D., Lee, J.-Y., Yang, J.-S., Kim, J. W., Kim, V. N., & Chang, H. (2020). 'The Architecture of SARS-CoV-2 Transcriptome', *Cell*, 181/4: 914–21.e10.

Korber, B., Fischer, W. M., Gnanakaran, S., Yoon, H., Theiler, J., Abfalterer, W., Hengartner, N., et al. (2020). 'Tracking Changes in SARS-CoV-2 Spike: Evidence that D614G Increases Infectivity of the COVID-19 Virus', *Cell*, 182/4: 812–27.e19.

Lauring, A. S., & Hodcroft, E. B. (2021). 'Genetic Variants of SARS-CoV-2-What Do They Mean?', *JAMA: the journal of the American Medical Association*, 325/6: 529–31.

Leigh, J. W., & Bryant, D. (2015). 'popart: full-feature software for haplotype network construction', *Methods in ecology and evolution / British Ecological Society*, 6/9: 1110–6. John Wiley & Sons, Ltd.

Lindahl, Abraham, Hess, & Spoel, V. der. (2020). *GROMACS 2020 Source code*. DOI: 10.5281/zenodo.3562495

Liu, Y., Liu, J., Plante, K. S., Plante, J. A., Xie, X., Zhang, X., Ku, Z., et al. (2021). 'The N501Y spike substitution enhances SARS-CoV-2 transmission', *bioRxiv : the preprint server for biology*. DOI: 10.1101/2021.03.08.434499

Li, Y., Ma, M.-L., Lei, Q., Wang, F., Hong, W., Lai, D.-Y., Hou, H., et al. (2021). 'Linear epitope landscape of the SARS-CoV-2 Spike protein constructed from 1,051 COVID-19 patients', *Cell reports*, 34/13: 108915.



- 1009 Luan, B., Wang, H., & Huynh, T. (n.d.). ‘Molecular Mechanism of the N501Y Mutation for
- 1010 Enhanced Binding between SARS-CoV-2’s Spike Protein and Human ACE2 Receptor’.
- 1011 DOI: 10.1101/2021.01.04.425316
- 1012 Lu, J., Li, B., Deng, A., Li, K., Hu, Y., Li, Z., Xiong, Q., et al. (n.d.). ‘Viral infection and
- 1013 transmission in a large, well-traced outbreak caused by the SARS-CoV-2 Delta variant’.
- 1014 DOI: 10.21203/rs.3.rs-738164/v1
- 1015 Martin, D. P., Weaver, S., Tegally, H., San, E. J., Shank, S. D., Wilkinson, E., Giandhari, J., et
- 1016 al. (2021). ‘The emergence and ongoing convergent evolution of the N501Y lineages
- 1017 coincides with a major global shift in the SARS-CoV-2 selective landscape’, *medRxiv : the*
- 1018 *preprint server for health sciences*. DOI: 10.1101/2021.02.23.21252268
- 1019 Muller, H. J. (1964). ‘THE RELATION OF RECOMBINATION TO MUTATIONAL
- 1020 ADVANCE’, *Mutation research*, 106: 2–9.
- 1021 Müller, N. F., Kistler, K. E., & Bedford, T. (2021). ‘Recombination patterns in coronaviruses’,
- 1022 *bioRxiv : the preprint server for biology*. DOI: 10.1101/2021.04.28.441806
- 1023 Oude Munnink, B. B., Sikkema, R. S., Nieuwenhuijse, D. F., Molenaar, R. J., Munger, E.,
- 1024 Molenkamp, R., van der Spek, A., et al. (2021). ‘Transmission of SARS-CoV-2 on mink
- 1025 farms between humans and mink and back to humans’, *Science*, 371/6525: 172–7.
- 1026 Papa, G., Mallery, D. L., Albecka, A., Welch, L., Cattin-Ortolá, J., Luptak, J., Paul, D., et al.
- 1027 (n.d.). ‘Furin cleavage of SARS-CoV-2 Spike promotes but is not essential for infection and
- 1028 cell-cell fusion’. DOI: 10.1101/2020.08.13.243303
- 1029 Parrinello, M., & Rahman, A. (1981). ‘Polymorphic transitions in single crystals: A new
- 1030 molecular dynamics method’, *Journal of applied physics*, 52/12: 7182–90. American
- 1031 Institute of Physics.
- 1032 Prates, E., Garvin, M., Jones, P., Miller, J. I., Kyle, S., Cliff, A., Gazolla, J. G. F. M., et al.
- 1033 (2021). ‘Antiviral Strategies Against SARS-CoV-2 – For a Bioinformatics Approach’. Hann
- 1034 J. J., Bintou A., & Keng C. (eds) *SARS-CoV-2 Methods and Protocols*. Springer.
- 1035 Prates, E. T., Garvin, M. R., Pavicic, M., Jones, P., Shah, M., Demerdash, O., Amos, B. K., et al.
- 1036 (2020). ‘Potential pathogenicity determinants identified from structural proteomics of
- 1037 SARS-CoV and SARS-CoV-2’, *Molecular biology and evolution*. DOI:
- 1038 10.1093/molbev/msaa231
- 1039 Rand, A. C., Jain, M., Eizenga, J. M., Musselman-Brown, A., Olsen, H. E., Akeson, M., & Paten,
- 1040 B. (2017). ‘Mapping DNA methylation with high-throughput nanopore sequencing’, *Nature*
- 1041 *methods*, 14/4: 411–3.
- 1042 Richard, M., Kok, A., de Meulder, D., Bestebroer, T. M., Lamers, M. M., Okba, N. M. A.,
- 1043 Fentener van Vlissingen, M., et al. (2020). ‘SARS-CoV-2 is transmitted via contact and via
- 1044 the air between ferrets’, *Nature communications*, 11/1: 3496.
- 1045 Rodrigues, C. H. M., Pires, D. E. V., & Ascher, D. B. (2021). ‘DynaMut2: Assessing changes in
- 1046 stability and flexibility upon single and multiple point missense mutations’, *Protein science: a publication of the Protein Society*, 30/1: 60–9.
- 1047 Ronquist, F., & Huelsenbeck, J. P. (2003). ‘MrBayes 3: Bayesian phylogenetic inference under
- 1048 mixed models’, *Bioinformatics*, 19/12: 1572–4.
- 1049 Sabino, E. C., Buss, L. F., Carvalho, M. P. S., Prete, C. A., Jr, Crispim, M. A. E., Fraiji, N. A.,
- 1050 Pereira, R. H. M., et al. (2021). ‘Resurgence of COVID-19 in Manaus, Brazil, despite high
- 1051 seroprevalence’, *The Lancet*. DOI: 10.1016/S0140-6736(21)00183-5
- 1052 Santerre, M., Arjona, S. P., Allen, C. N., Shcherbik, N., & Sawaya, B. E. (2020). ‘Why do
- 1053 SARS-CoV-2 NSPs rush to the ER?’, *Journal of neurology*. DOI: 10.1007/s00415-020-
- 1054



10197-8

Sawatzki, K., Hill, N. J., Puryear, W. B., Foss, A. D., Stone, J. J., & Runstadler, J. A. (2021). 'Host barriers to SARS-CoV-2 demonstrated by ferrets in a high-exposure domestic setting', *Proceedings of the National Academy of Sciences of the United States of America*, 118/18. DOI: 10.1073/pnas.2025601118

Shang, J., Ye, G., Shi, K., Wan, Y., Luo, C., Aihara, H., Geng, Q., et al. (2020). 'Structural basis of receptor recognition by SARS-CoV-2', *Nature*, 581/7807: 221–4.

Shannon, P., Markiel, A., Ozier, O., Baliga, N. S., Wang, J. T., Ramage, D., Amin, N., et al. (2003). 'Cytoscape: a software environment for integrated models of biomolecular interaction networks', *Genome research*, 13/11: 2498–504.

Sheikh, A., McMenamin, J., Taylor, B., & Robertson, C. (2021). 'SARS-CoV-2 Delta VOC in Scotland: demographics, risk of hospital admission, and vaccine effectiveness'. *The Lancet*. DOI: 10.1016/s0140-6736(21)01358-1

Simon-Loriere, E., & Holmes, E. C. (2011). 'Why do RNA viruses recombine?', *Nature reviews. Microbiology*, 9/8: 617–26.

Simpson, J. T., Workman, R. E., Zuzarte, P. C., David, M., Dursi, L. J., & Timp, W. (2017). 'Detecting DNA cytosine methylation using nanopore sequencing', *Nature methods*, 14/4: 407–10.

Singh, A., Steinkellner, G., Köchl, K., Gruber, K., & Gruber, C. C. (n.d.). 'Serine 477 plays a crucial role in the interaction of the SARS-CoV-2 spike protein with the human receptor ACE2'. DOI: 10.21203/rs.3.rs-106969/v1

Singh, J., Rahman, S. A., Ehtesham, N. Z., Hira, S., & Hasnain, S. E. (2021). 'SARS-CoV-2 variants of concern are emerging in India', *Nature medicine*. DOI: 10.1038/s41591-021-01397-4

Starr, T. N., Greaney, A. J., Addetia, A., Hannon, W. W., Choudhary, M. C., Dingens, A. S., Li, J. Z., et al. (2021). 'Prospective mapping of viral mutations that escape antibodies used to treat COVID-19', *Science*. DOI: 10.1126/science.abf9302

Starr, T. N., Greaney, A. J., Hilton, S. K., Crawford, K. H. D., Navarro, M. J., Bowen, J. E., Alejandra Tortorici, M., et al. (n.d.). 'Deep mutational scanning of SARS-CoV-2 receptor binding domain reveals constraints on folding and ACE2 binding'. DOI: 10.1101/2020.06.17.157982

'The GIMP Development Team. (2019). GIMP. Retrieved from <https://www.gimp.org>'. (n.d.). <https://www.gimp.org>.

Thorne, L. G., Bouhaddou, M., Reuschl, A.-K., Zuliani-Alvarez, L., Polacco, B., Pelin, A., Batra, J., et al. (2021). 'Evolution of enhanced innate immune evasion by the SARS-CoV-2 B.1.1.7 UK variant', *bioRxiv : the preprint server for biology*. DOI: 10.1101/2021.06.06.446826

Tylor, S., Andonov, A., Cutts, T., Cao, J., Grudesky, E., Van Domselaar, G., Li, X., et al. (2009). 'The SR-rich motif in SARS-CoV nucleocapsid protein is important for virus replication'. *Canadian Journal of Microbiology*. DOI: 10.1139/w08-139

Velasco, J. D. (2013). 'Phylogeny as population history'. *Philosophy and Theory in Biology*. DOI: 10.3998/ptb.6959004.0005.002

Verba, K., Gupta, M., Azumaya, C., Moritz, M., Pourmal, S., Diallo, A., Merz, G., et al. (2021). 'CryoEM and AI reveal a structure of SARS-CoV-2 Nsp2, a multifunctional protein involved in key host processes', *Research square*. DOI: 10.21203/rs.3.rs-515215/v1

Volz, E., Mishra, S., Chand, M., Barrett, J. C., Johnson, R., Geidelberg, L., Hinsley, W. R., et al.

(n.d.). ‘Transmission of SARS-CoV-2 Lineage B.1.1.7 in England: Insights from linking epidemiological and genetic data’, DOI: 10.1101/2020.12.30.20249034

Wang, Z., Schmidt, F., Weisblum, Y., Muecksch, F., Barnes, C. O., Finkin, S., Schaefer-Babajew, D., et al. (2021). ‘mRNA vaccine-elicited antibodies to SARS-CoV-2 and circulating variants’, *Nature*. DOI: 10.1038/s41586-021-03324-6

Washington, N. L., Gangavarapu, K., Zeller, M., Bolze, A., Cirulli, E. T., Schiabor Barrett, K. M., Larsen, B. B., et al. (2021). ‘Emergence and rapid transmission of SARS-CoV-2 B.1.1.7 in the United States’, *Cell*, 184/10: 2587–94.e7.

Wrapp, D., Wang, N., Corbett, K. S., Goldsmith, J. A., Hsieh, C.-L., Abiona, O., Graham, B. S., et al. (2020). ‘Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation’, *Science*, 367/6483: 1260–3.

Xia, H., Cao, Z., Xie, X., Zhang, X., Chen, J. Y.-C., Wang, H., Menachery, V. D., et al. (2020). ‘Evasion of Type I Interferon by SARS-CoV-2’, *Cell reports*, 33/1: 108234.

Yan, L., Zhang, Y., Ge, J., Zheng, L., Gao, Y., Wang, T., Jia, Z., et al. (2020). ‘Architecture of a SARS-CoV-2 mini replication and transcription complex’, *Nature communications*, 11/1: 5874.

Yan, R., Zhang, Y., Li, Y., Xia, L., Guo, Y., & Zhou, Q. (2020). ‘Structural basis for the recognition of SARS-CoV-2 by full-length human ACE2’, *Science*, 367/6485: 1444–8.

Zhang, B.-Z., Hu, Y.-F., Chen, L.-L., Yau, T., Tong, Y.-G., Hu, J.-C., Cai, J.-P., et al. (2020). ‘Mining of epitopes on spike protein of SARS-CoV-2 from COVID-19 patients’, *Cell research*, 30/8: 702–4.

Zhang, L., Jackson, C. B., Mou, H., Ojha, A., Peng, H., Quinlan, B. D., Rangarajan, E. S., et al. (2020). ‘SARS-CoV-2 spike-protein D614G mutation increases virion spike density and infectivity’, *Nature communications*, 11/1: 6013.

# **Supplementary Material**

## **The emergence of highly fit SARS-CoV-2 variants accelerated by recombination**

Michael R. Garvin<sup>1,2+\*</sup>, Erica T. Prates<sup>1,2+</sup>, Jonathon Romero<sup>3</sup>, Ashley Cliff<sup>3</sup>, Joao Gabriel Felipe Machado Gazolla<sup>1,2</sup>, Monica Pickholz<sup>4,5</sup>, Mirko Pavicic<sup>1,2</sup>, Daniel Jacobson<sup>1,2,\*</sup>

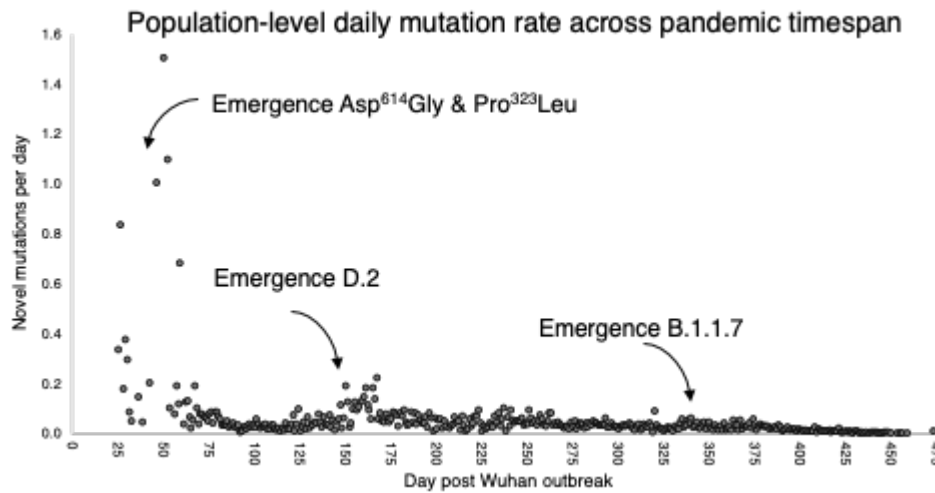
### **Affiliations:**

<sup>1</sup>Oak Ridge National Laboratory, Computational Systems Biology, Biosciences, Oak Ridge, TN; <sup>2</sup>National Virtual Biotechnology Laboratory, US Department of Energy; <sup>3</sup>The Bredesen Center for Interdisciplinary Research and Graduate Education, University of Tennessee Knoxville, Knoxville, TN; <sup>4</sup>Departamento de Física, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires, Buenos Aires, Argentina; <sup>5</sup> Instituto de Física de Buenos Aires (IFIBA), CONICET-Universidad de Buenos Aires, Buenos Aires, Argentina.

**\*Correspondence:** [garvinmr@ornl.gov](mailto:garvinmr@ornl.gov), [jacobsonda@ornl.gov](mailto:jacobsonda@ornl.gov)

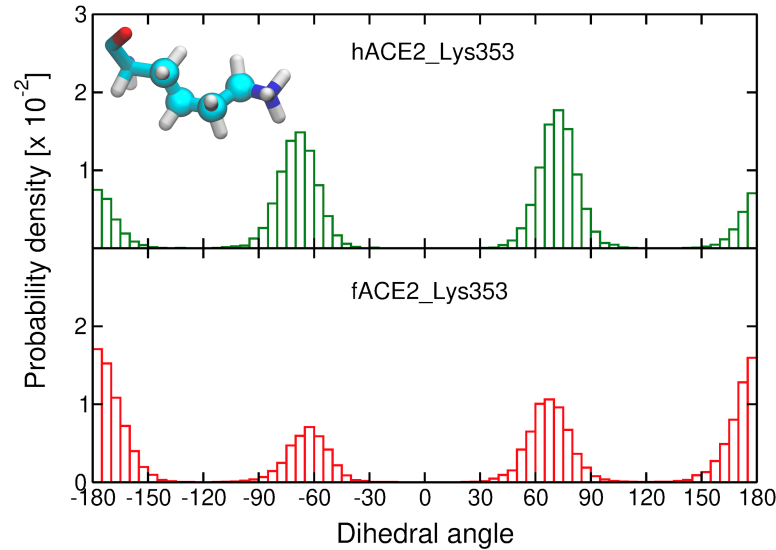
**+Contributed equally**

# 1. Supplementary Figures



**Fig. S1.**

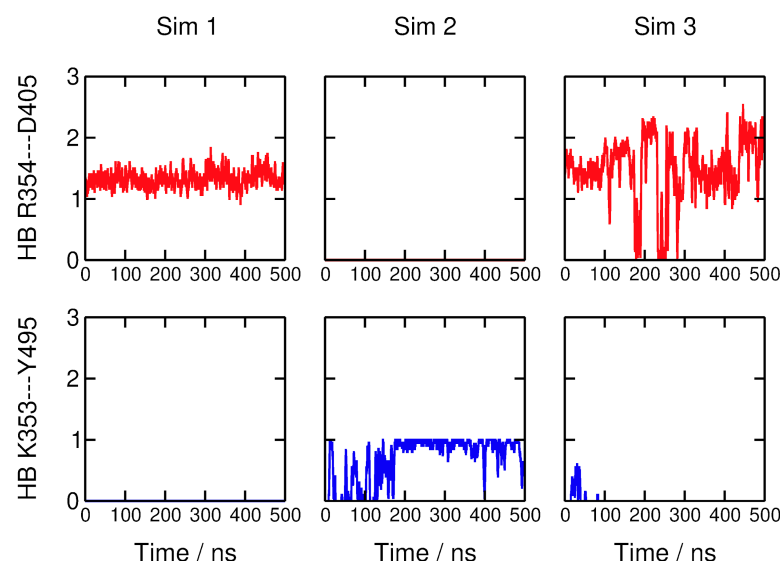
**Population level mutation rate over the course of the pandemic.** Number of novel mutations sampled across the globe for each day are plotted against time (days from the Wuhan outbreak). Emergence of major VOC are provided for context and show small increases in the number of new mutations but there is an overall decrease across time, even accounting for multiple mutations at a site to different nucleotide states and deletions.



**Fig. S2.**

**The probability density of the conformations of Lys<sup>353</sup> in human and ferret ACE2 in the simulations.**

Histograms of the distribution of a dihedral angle of the Lys<sup>353</sup> side chain carbon atoms in human ACE2 (hACE2, upper figure) and ferret ACE2 (fACE2, lower figure) in complex with the SARS-CoV-2 S receptor-binding domain. The atoms forming the selected dihedral are depicted as spheres in the molecular representation of Lys<sup>353</sup>. Three independent simulations are considered for the calculation of the histograms. Dihedral angles near  $\pm 180^\circ$  correspond to a more stretched conformation (i.e., *trans*).



**Fig. S3.**  
**Competing hydrogen bond interactions formed between positively charged amino acid residues in ferret ACE2 (fACE2) and the SARS-CoV-2 S receptor-binding domain.** Time evolution of the number of hydrogen bonds (HB) that fACE2 Arg3<sup>54</sup> and Lys<sup>353</sup> form with Asp<sup>405</sup> and Tyr<sup>495</sup> from the SARS-CoV-2 S receptor-binding domain. The columns correspond to the three simulation replicas. The geometric criteria adopted for hydrogen bonds are a cutoff of 3.0 Å for donor-acceptor distance and 20° for acceptor-donor-H angle.



## 2. Supplementary Tables

**Table S4.**

**Average number of contacts formed between Asn<sup>501</sup> in the receptor-binding domains of SARS-CoV-2 S and residues in ACE2 from human (hACE2) and ferret (fACE2).** A distance of 4 Å between any atom pairs was defined as the cut-off for contact statistics.

ACE2 residue	hACE2	fACE2
<b>Tyr<sup>41</sup></b>	0.96 ± 0.02	0.80 ± 0.03
<b>Lys<sup>353</sup></b>	0.99 ± 0.01	0.90 ± 0.01
<b>Asp<sup>353</sup></b>	0.98 ± 0.01	0.70 ± 0.04

s