

## **The ongoing evolution of variants of concern and interest of SARS-CoV-2 in Brazil revealed by convergent indels in the amino (N)-terminal domain of the Spike protein**

Paola Cristina Resende<sup>1\*</sup>, Felipe G Naveca<sup>2\*</sup>, Roberto D. Lins<sup>3</sup>, Filipe Zimmer Dezordi<sup>4,5</sup>, Matheus V. F. Ferraz<sup>3,6</sup>, Emerson G. Moreira<sup>3,6</sup>, Danilo F. Coêlho<sup>3,6</sup>, Fernando Couto Motta<sup>1</sup>, Anna Carolina Dias Paixão<sup>1</sup>, Luciana Appolinario<sup>1</sup>, Renata Serrano Lopes<sup>1</sup>, Ana Carolina da Fonseca Mendonça<sup>1</sup>, Alice Sampaio Barreto da Rocha<sup>1</sup>, Valdinete Nascimento<sup>2</sup>, Victor Souza<sup>2</sup>, George Silva<sup>2</sup>, Fernanda Nascimento<sup>2</sup>, Lidio Gonçalves Lima Neto<sup>7</sup>, Fabiano Vieira da Silva<sup>7</sup>, Irina Riediger<sup>8</sup>, Maria do Carmo Debur<sup>8</sup>, Anderson Brandao Leite<sup>9</sup>, Tirza Mattos<sup>10</sup>, Cristiano Fernandes da Costa<sup>11</sup>, Felicidade Mota Pereira<sup>12</sup>, Cliomar Alves dos Santos<sup>13</sup>, Darcita Buerger Rovaris<sup>14</sup>, Sandra Bianchini Fernandes<sup>14</sup>, Adriano Abbud<sup>15</sup>, Claudio Sacchi<sup>15</sup>, Ricardo Khouri<sup>16</sup>, André Felipe Leal Bernardes<sup>17</sup>, Edson Delatorre<sup>18</sup>, Tiago Gräf<sup>19</sup>, Marilda Mendonça Siqueira<sup>1</sup>, Gonzalo Bello<sup>\*\*20</sup>, and Gabriel L Wallau<sup>\*\*4,5</sup> on behalf of Fiocruz COVID-19 Genomic Surveillance Network.

1. Laboratory of Respiratory Viruses and Measles (LVRS), Instituto Oswaldo Cruz, FIOCRUZ-Rio de Janeiro, Brazil.
2. Laboratório de Ecologia de Doenças Transmissíveis na Amazônia (EDTA), Instituto Leônidas e Maria Deane, FIOCRUZ-Amazonas, Brazil.
3. Department of Virology, Instituto Aggeu Magalhães, FIOCRUZ-Pernambuco, Brazil.
4. Departamento de Entomologia, Instituto Aggeu Magalhães, FIOCRUZ-Pernambuco, Brazil.
5. Núcleo de Bioinformática (NBI), Instituto Aggeu Magalhães FIOCRUZ-Pernambuco, Brazil.
6. Department of Fundamental Chemistry, Federal University of Pernambuco, Recife, Brazil
7. Laboratório Central de Saúde Pública do Estado do Maranhão (LACEN-MA), Brazil.
8. Laboratório Central de Saúde Pública do Estado do Paraná (LACEN-PR), Brazil.
9. Laboratório Central de Saúde Pública do Estado do Alagoas (LACEN-AL), Brazil.
10. Laboratório Central de Saúde Pública do Amazonas (LACEN-AM), Brazil.
11. Fundação de Vigilância em Saúde do Amazonas, Brazil.
12. Laboratório Central de Saúde Pública do Estado da Bahia (LACEN-BA), Brazil.
13. Laboratório Central de Saúde Pública do Estado de Sergipe (LACEN-SE), Aracajú, Sergipe, Brazil.
14. Laboratório Central de Saúde Pública do Estado de Santa Catarina (LACEN-SC), Florianópolis, Santa Catarina, Brazil.
15. Instituto Adolfo Lutz, São Paulo, São Paulo, Brazil.
16. Laboratório de Enfermidades Infecciosas Transmitidas por Vetores, Instituto Gonçalo Moniz, FIOCRUZ-Bahia, Salvador, Bahia, Brazil.
17. Laboratório Central de Saúde Pública do Estado de Minas Gerais (LACEN-MG).

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

18. Departamento de Biologia. Centro de Ciências Exatas, Naturais e da Saúde, Universidade Federal do Espírito Santo, Alegre, Brazil.
19. Plataforma de Vigilância Molecular, Instituto Gonçalo Moniz, FIOCRUZ-Bahia, Brazil.
20. Laboratório de AIDS e Imunologia Molecular, Instituto Oswaldo Cruz, FIOCRUZ-Rio de Janeiro, Brazil.

\*, \*\* These authors contributed equally to this work.

## Abstract

Mutations at both the receptor-binding domain (RBD) and the amino (N)-terminal domain (NTD) of the SARS-CoV-2 Spike (S) glycoprotein can alter its antigenicity and promote immune escape. We identified that SARS-CoV-2 lineages circulating in Brazil with mutations of concern in the RBD independently acquired convergent deletions and insertions in the NTD of the S protein, which altered the NTD antigenic-supersite and other predicted epitopes at this region. Importantly, we detected community transmission of four lineages bearing NTD indels: a P.1  $\Delta$ 69-70 lineage (which can impact several SARS-CoV-2 diagnostic protocols), a P.1  $\Delta$ 144 lineage, a P.1-like lineage carrying ins214ANRN, and the VOI N.10 derived from the B.1.1.33 lineage carrying three deletions ( $\Delta$ 141-144,  $\Delta$ 211 and  $\Delta$ 256-258). These findings support that the ongoing widespread transmission of SARS-CoV-2 in Brazil is generating new viral lineages that might be more resistant to antibody neutralization than parental variants of concern.

**Keywords:** COVID-19, pandemics, antibody escape, coronavirus, community transmission

## Introduction

Recurrent deletions in the amino (N)-terminal domain (NTD) of the spike (S) glycoprotein of SARS-CoV-2 have been identified during long-term infection of immunocompromised patients <sup>1-4</sup> as well as during extended human-to-human transmission <sup>3</sup>. Most of those deletions (90%) maintain the reading frame and cover four recurrent deletion regions (RDRs) within the NTD at positions 60-75 (RDR1), 139-146 (RDR2), 210-212 (RDR3), and 242-248 (RDR4) of the S protein <sup>3</sup>. The RDRs that occupy defined antibody epitopes within the NTD and RDR regions might alter antigenicity <sup>3</sup>. Interestingly, the RDRs overlap with four NTD Indel Regions (IR - IR-2 to IR-5) that are prone to gain or lose short nucleotide sequences during sarbecoviruses evolution both in animals and humans <sup>5,6</sup>.

Since late 2020, several more transmissible variants of concern (VOCs) and also variants of interest (VOI) with convergent mutations at the receptor-binding domain (RBD) of the S protein (particularly E484K and N501Y) arose independently in humans <sup>7,8</sup>. Some VOCs also displayed NTD deletions such as lineages B.1.1.7 (RDR2  $\Delta$ 144), B.1.351 (RDR4  $\Delta$ 242-244), and P.3 (RDR2  $\Delta$ 141-143) that were initially detected in the United Kingdom, South Africa, and the Philippines, respectively <sup>3</sup>. The VOCs B.1.1.7 and B.1.351 are resistant to neutralization by several anti-NTD monoclonal antibodies (mAbs) and NTD deletions at RDR2 and RDR4 are important for such phenotype <sup>9-14</sup>. Thus, NTD mutations and deletions represent an important mechanism of immune evasion and accelerate SARS-CoV-2 adaptive evolution in humans.

Several SARS-CoV-2 variants with mutations in the RBD have been described in Brazil, including the VOC P.1 <sup>15</sup> and the VOIs P.2 <sup>16</sup>, N.9 <sup>17</sup> and N.10 <sup>18</sup>. With the exception of N.10, none of the other variants described in Brazil displayed indels in the NTD. Importantly, although the VOC P.1 displayed NTD mutations (L18F) that abrogate binding of some anti-NTD mAbs <sup>14</sup> and further showed reduced binding to RBD-directed antibodies, it is more susceptible to anti-NTD mAbs than other VOCs <sup>9-14,19</sup>. In this study, we characterized the emergence of RDR variants within VOC and VOIs circulating in Brazil that were genotyped by the Fiocruz COVID-19 Genomic Surveillance Network between November 2020 and February 2021.

## Results and Discussion

Our genomic survey identified 35 SARS-CoV-2 sequences from seven different Brazilian states that harbor a variable combination of mutations in the RBD (K417T,

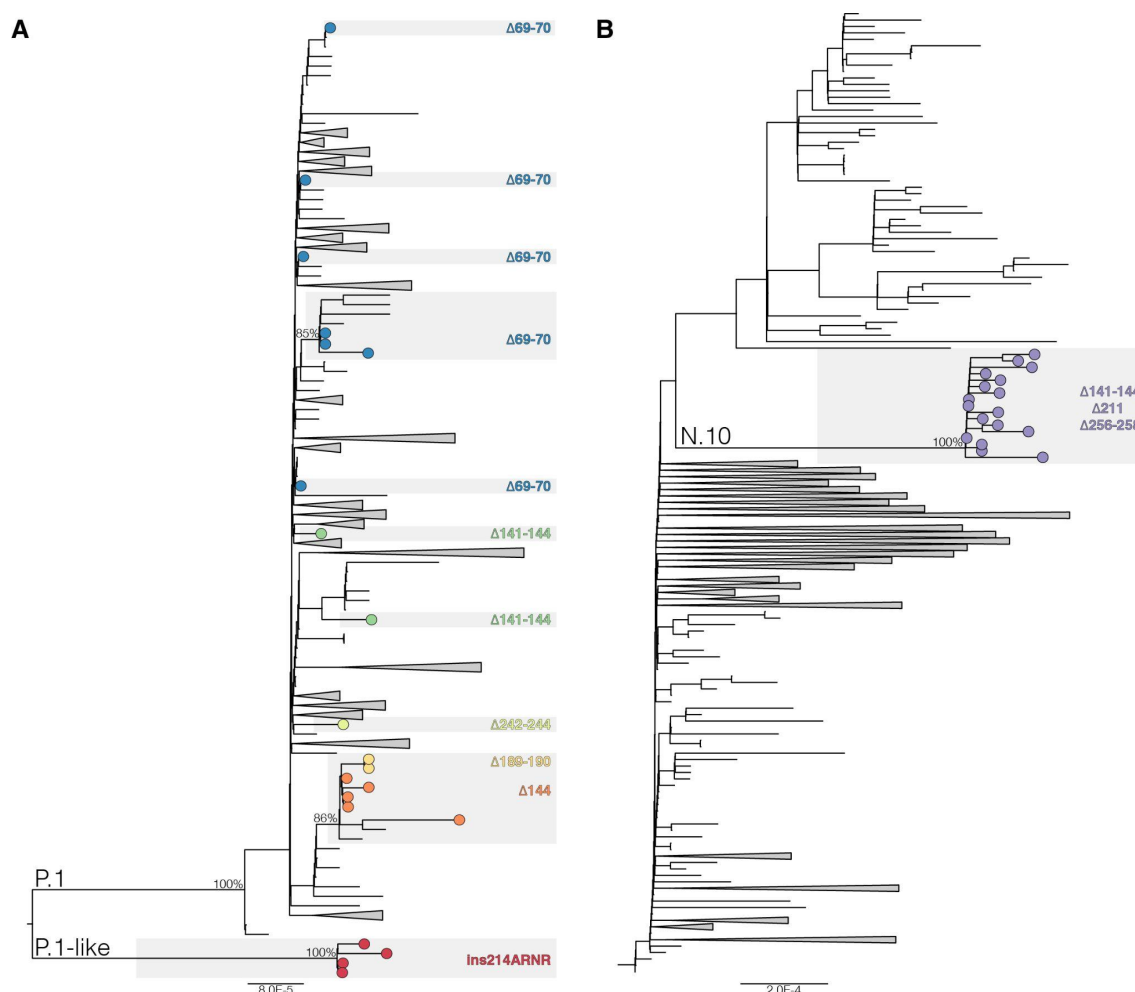
E484K, N501Y) and indels in the NTD region of the S protein. These genomes were classified within lineages N.10 (n = 16), P.1 (n = 14), P.2 (n = 1) and B.1.1.28 (P1-like, n = 4) (**Table 1**). Seven VOC P.1 sequences displayed deletion  $\Delta 69-70$  in the RDR1, three sequences (two VOC P.1 and one VOI P.2) displayed deletion  $\Delta 144$  in the RDR2, two P.1 sequences showed a four amino acid deletion  $\Delta 141-144$  in the RDR2, two P.1 sequences harbors a two amino acid deletion  $\Delta 189-190$ , and one P.1 sequence displayed a three amino acid deletion  $\Delta 242-244$  in the RDR4. We also detected four B.1.1.28 P.1-like genomes bearing an ins214ANRN insertion upstream to RDR 3 and sharing six out of 10 P.1 lineage-defining mutations in the Spike protein (L18F, P26S, D138Y, K417T, E484K, N501Y) as well as P.1 lineage-defining mutations in the NSP3 (K977Q), NS3 (S253P) and N (P80R) proteins <sup>20</sup>. The VOI N.10 displayed NTD indels  $\Delta 141-144$  at RDR2,  $\Delta 211$  at RDR3 and  $\Delta 256-258$  close to RDR4 <sup>18</sup>. Inspection of sequences available at EpiCoV database in the GISAID (<https://www.gisaid.org/>) at March 1st, 2021, retrieved three P.1 sequences from the Bahia state <sup>21</sup> and one B.1.1.28 sequence from the Amazonas state <sup>20</sup> with deletion  $\Delta 144$  (**Table 1**).

The Maximum Likelihood (ML) phylogenetic analysis of lineage P.1 supports recurrent emergence of variants  $\Delta 141-144$  and  $\Delta 69-70$  and the monophyletic origin of variants  $\Delta 144$  and  $\Delta 189-190$  (**Fig. 1A**). Both P.1  $\Delta 141-144$  sequences recovered from patients from Amazonas and Rondônia states, all P.1  $\Delta 69-70$  sequences from Santa Catarina state and the P.1  $\Delta 242-244$  sequence from Sergipe state appeared as singletons intermixed among non-deleted P.1 sequences. The remaining P.1 variants with NTD deletions were distributed in two sub-clades that also include non-deleted P.1 sequences. One sub-clade (aLRT = 86%) was characterized by the mutations ORF1a:T951I and A18945G and comprises nine sequences: the five P.1  $\Delta 144$ , the two P.1  $\Delta 189-190$  and two P.1 from Amazonas and Goiás states. The other sub-clade (aLRT = 85%) was characterized by the synonymous mutations G29781A and T29834A and comprises seven sequences: the three P.1  $\Delta 69-70$  from São Paulo state plus four P.1 sequences from São Paulo, Amazonas and Tocantins states. The ML phylogenetic analyses further confirm that all sequences belonging to variants P.1-like ins214ANRN (**Fig. 1A**) and VOI N.10 (**Fig. 1B**) branched in highly supported (aLRT > 99%) monophyletic clades. These findings revealed that NTD deletions characteristic of VOCs B.1.1.7 ( $\Delta 69-70$  and  $\Delta 144$ ) and B.1.351 ( $\Delta 242-244$ ) occurred at multiple times during the evolution of lineage P.1 and also sporadically arose in lineages B.1.1.28, B.1.1.33 (N.10) and P.2.

**Table 1. SARS-CoV-2 Brazilian variants with indels at NTD of the Spike protein.**

Sample(s)	Lineage	NTD Indel	RBD	GISAID ID
AM-FIOCRUZ-20842572LS/2020*	B.1.1.28	Δ144	-	EPI_ISL_1068132
MG-FIOCRUZ-8180/2021	P.2	Δ144	E484K	EPI_ISL_1219137
SC-FIOCRUZ-13109/2021	P.1	Δ69-70	K417T	EPI_ISL_1533994
SC-FIOCRUZ-13111/2021			E484K	EPI_ISL_1533992
SC-FIOCRUZ-13113/2021			N501Y	EPI_ISL_1533996
SC-FIOCRUZ-13114/2021				EPI_ISL_1533993
SP-2075/2021				EPI_ISL_1498917
SP-2084/2021				EPI_ISL_1509639
SP-2088/2021				EPI_ISL_1509720
BA53/2021*	P.1	Δ144	K417T	EPI_ISL_1067729
BA54/2021*			E484K	EPI_ISL_1067733
BA55/2021*			N501Y	EPI_ISL_1067734
BA-FIOCRUZ-7029/2021*				EPI_ISL_1219136
AM-FIOCRUZ-21140861HC*				EPI_ISL_1533609
AL-FIOCRUZ-4795/2021*	P.1	Δ141-144	K417T	EPI_ISL_1219134
PR-FIOCRUZ-5273/2021**			E484K	EPI_ISL_1219133
			N501Y	
AL-FIOCRUZ-4786/2021*	P.1	Δ189-190	K417T	EPI_ISL_1219135
RS-FIOCRUZ-14243/2021			E484K	EPI_ISL_1534013
			N501Y	
SE-FIOCRUZ-10220/2021	P.1	Δ242-244	K417T	EPI_ISL_1534004
			E484K	
			N501Y	
MA-FIOCRUZ-6871/2021***	N.10	Δ141-144	V445A	EPI_ISL_1181371
		Δ211	E484K	
		Δ256-258		
AM-FIOCRUZ-20897269OP*	B.1.1.28	ins214ANRN	K417T	EPI_ISL_1068256
AM-FIOCRUZ-20897281WS*	(P.1-like)		E484K	EPI_ISL_1219132
AM-FIOCRUZ-21840593CL*			N501Y	EPI_ISL_1261122
PR-FIOCRUZ-5241/2021				EPI_ISL_1261123

\*Patient from Amazonas state or traveller returning from Amazonas state. \*\* Patient from Rondônia. \*\*\* Sequence representative of lineage N.10 (EPI\_ISL\_1181370, EPI\_ISL\_1465226, EPI\_ISL\_1465228, EPI\_ISL\_1465229, EPI\_ISL\_1465231, EPI\_ISL\_1465232, EPI\_ISL\_1465234, EPI\_ISL\_1465235, EPI\_ISL\_1465236, EPI\_ISL\_1465238, EPI\_ISL\_1465239, EPI\_ISL\_1465241, EPI\_ISL\_1465242, EPI\_ISL\_1465243, and EPI\_ISL\_1465245). Sequencing depth plots of the samples bearing indels are available in Supplementary Figure S1.



**Figure 1.** ML phylogenetic tree of whole-genome lineage P.1/P.1-like (A) and B.1.1.33 (B) Brazilian sequences showing the recurrent emergence of deletions at the NTD of the S protein. Tip circles representing the SARS-CoV-2 sequences with NTD indels are colored as indicated. The branch lengths are drawn to scale with the left bar indicating nucleotide substitutions per site. For visual clarity, some clades are collapsed into triangles.

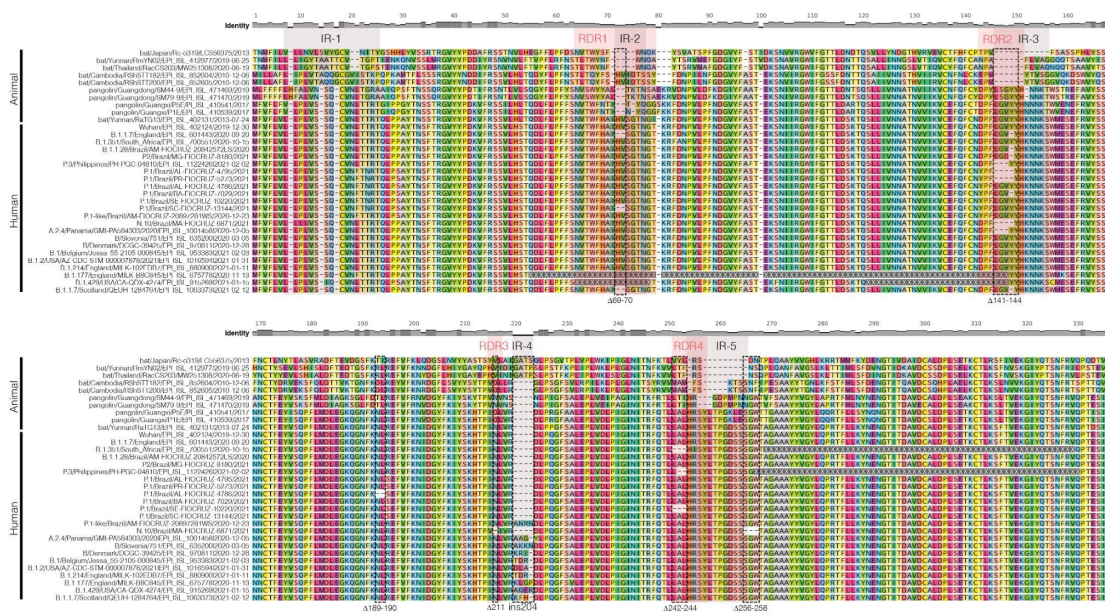
Most P.1  $\Delta 144$ ,  $\Delta 141-144$  and  $\Delta 189-190$  sequences were detected in the Amazonas or were recovered from individuals that were transferred from or that reported a travel history to the Amazonas state, like all P.1  $\Delta 144$  sequences from the Bahia state described previously<sup>21</sup> and in the present study (**Table 1**). Thus, those P.1 variants as well as the variant P.1-like ins214 probably emerged in the Amazonas state and some of them displayed low-level of community transmission. The P.1  $\Delta 69-70$  sequences, by contrast, were detected in autochthonous cases from Santa Catarina and Sao Paulo states that had no history of travel to the Amazonas. The phylogenetic clustering supports the independent origin of variant P.1  $\Delta 69-70$  in both Brazilian states and its local dissemination in São Paulo, but not in Santa Catarina. The lack of monophyletic clustering of P.1  $\Delta 69-70$  sequences from Santa Catarina, however, should

be interpreted with caution due to the paucity of synapomorphic mutations within diversity of lineage P.1. Other variants that also arose outside the Amazonas state were the P.1  $\Delta$ 242-244, the P.2  $\Delta$ 144 and the VOI N.10 detected in the states of Sergipe, Minas Gerais and Maranhão, respectively <sup>18</sup>.

While SARS-CoV-2 variants harboring NTD deletions at RDR2 and RDR4 have emerged in many different lineages globally, insertions in the S protein are rare events. Our search of SARS-CoV-2 sequences available at EpiCoV database in the GISAID (<https://www.gisaid.org/>) on March 1st retrieved only 146 SARS-CoV-2 sequences of lineages A.2.4 (n = 52), B (n = 3), B.1 (n = 7), B.1.1.7 (n = 1), B.1.177 (n = 1), B.1.2 (n = 1), B.1.214 (n = 80) and B.1.429 (n = 1) that displayed an insert motif of three to four amino acids (AKKN, KLGB, AQER, AAG, KFH, KRI, and TDR) in position 214 (**Appendix Table 1**). Most ins214 motifs were unique, except for the ins214TDR that arose independently in lineages B.1 and B.1.214. With the only exception of one lineage B sequence sampled in March 2020, all SARS-CoV-2 ins214 variants were only detected since November 2020, and its frequency increased in 2021 mainly due to the recent dissemination of lineage A.2.4 ins214AAG in Central and North America and of lineage B.1.214 ins214TDR in Europe.

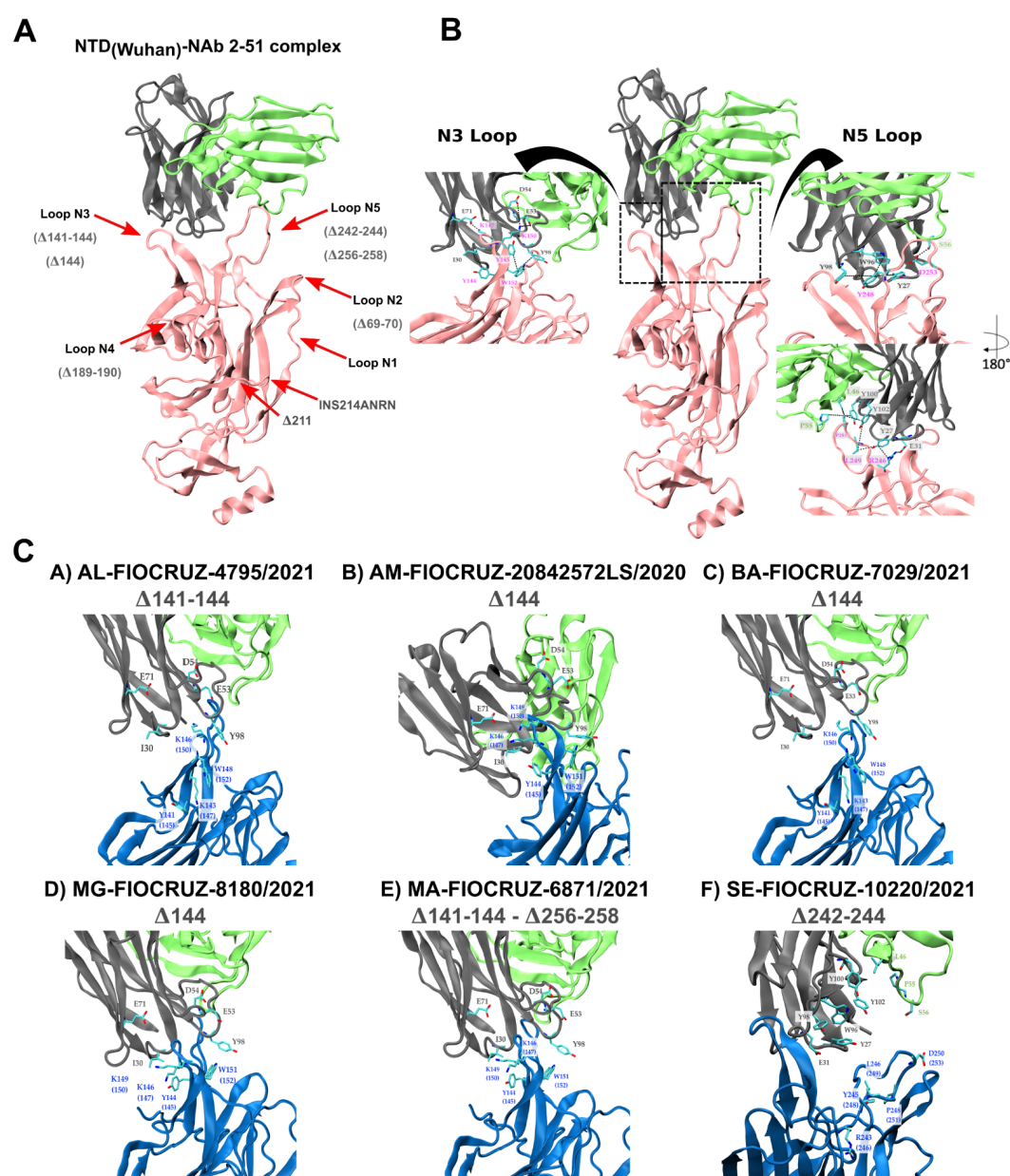
To better understand the evolutionary context of NTD indels, we aligned the S protein of representative sequences of SARS-CoV-2 lineages with NTD indels and SARS-CoV-2-related coronavirus (SC2r-CoV) lineages from bats and pangolins <sup>22</sup>. Inspection of the alignment confirms that most NTD indels detected in the SARS-CoV-2 lineages occur within IR previously defined in sarbecovirus (**Fig. 2**). The  $\Delta$ 141-144 occurs in the IR-3 located in the central part of the NTD, where some bats SC2r-CoV also have deletions. The  $\Delta$ 211 and ins214 occurs near the IR-4 where some bat SC2r-CoV from China (RmYN02, ins214GATP), Thailand (RacCS203, ins214GATP), and Japan (Rc-o319, ins214GATS) displayed a four amino acids insertion. Despite amino acid motifs at ins214 are very different across SARS-CoV-2 and SC2r-CoV lineages, the insertion size (3-4 amino acids) was conserved. Deletions  $\Delta$ 242-244 and  $\Delta$ 256-258 occur immediately upstream and downstream to IR-5, respectively, where some bat and pangolin SC2r-CoV lineages also displayed deletions. Thus, NTD regions that are prone to gain indels during viral transmission among animals are the same as those detected during transmissions in humans.





**Figure 2.** Amino acid alignment of Sarbecovirus NTD Spike region up to amino acid 335 including representative sequences of SARS-CoV-2 lineages harboring indels in the NTD and SARS-CoV-2-related coronavirus (SC2r-CoV) from bats and pangolins. IRs and RDRs positions (gray and red shaded areas, respectively) are approximations due to the high genetic variability in these alignment positions. Dotted rectangles highlight the indels identified in this study. The relative identity level estimated for each position of the alignment is displayed at the top.

Epitope mapping showed that neutralizing antibodies are primarily directed against the RBD and NTD of the S protein<sup>9,23-26</sup>. Some of the RBD mutations (K417T and E484K) detected in the VOCs and VOIs circulating in Brazil have been associated with increased resistance to neutralization by mAbs, or polyclonal sera from convalescent and vaccinated subjects<sup>27-31</sup>. The RDR2 and RDR4 are located in the N3 (residues 141 to 156) and N5 (residues 246 to 260) loops that composes the NTD antigenic-supersite<sup>32,33</sup> and deletions at those RDRs are also an essential mechanism for SARS-CoV-2 immune evasion of anti-NTD Abs<sup>3,9,10,34,14,3,9,10,34</sup>. To further visualize the potential impact of NTD deletions on immune recognition, we performed a modeling analysis of the binding interface between the NTD region and the NTD-directed neutralizing antibody (NAb) 2-51 derived from a convalescent donor<sup>23,33</sup>. The NAb 2-51 interacts with the wildtype NTD antigenic-supersite (EPI\_ISL\_402124) through several contacts with loops N3 and N5, with a predominance of hydrophobic contacts and dispersion interactions in N5 and saline interactions in N3 (**Fig. 3A and B**).



**Figure 3.** Representation of the Spike NTD 3D structure of wild type (pink) and NTD deleted variants (colored in blue) complexed to the NAb 2-51 heavy (gray) and light (green) chains. **A)** Relative position of the five NTD loops (red arrows) and the NTD deletions detected in our sample. **B)** Native interactions of mAb NAb 2-51 with N3 (left close-up) and N5 (right close-up) loops on the 3D structure of the wild type Spike NTD antigenic supersite. The N5 loop representation is also rotated 180° around its z-axis. **C)** Potential interactions of mAb NAb 2-51 with N3 and N5 loops on the 3D structure of the Spike NTD of N3 and N5 deleted variants. Residues making contact in the interface are depicted in the licorice representation, with carbon atoms in cyan, nitrogen atoms in blue and oxygen atoms in red. The dotted lines indicate the interacting residues-pair.

Our analyses corroborate that deletions at RDR2/IR-3 ( $\Delta 144$ ,  $\Delta 141-144$ ) and RDR4/IR-5 ( $\Delta 242-244$ ,  $\Delta 256-258$ ) detected in Brazilian sequences impact the N3 and N5 loops' size and conformation, disrupting the native contacts and reducing the interacting hydrophobic surface accessible area, mainly due to the loss of the hydrophobic pocket (**Figure 3C**). Indels around the N3/N5 loops resulted in a significant loss of interactions (both electrostatic and hydrophobic) that can dramatically impact the binding free energy, and therefore the binding affinity, between those NTD deletion variants and the NAb 2-51. Variant P.1  $\Delta 242-244$  displayed the largest loss of interactions, followed by variants N10, P.1  $\Delta 141-144$ , P.1  $\Delta 144$ , P.2  $\Delta 144$ , and B.1.1.28  $\Delta 144$  (**Table 2**). The NTD indels  $\Delta 69-70$ ,  $\Delta 189-190$ ,  $\Delta 211$  and ins214ANRN did not affect the NTD antigenic-supersite (**Figures 3A**), but they occur at other loops that comprise putative epitope regions covering residues 64-83, 168/173-188 and 209-216 (**Appendix Table 2**) and leads to conformational changes (**Supplementary Figure S2**) which might affect Ab binding outside the NTD antigenic-supersite. These findings suggest that NTD indels detected here probably abrogate the binding of NAb directed against the antigenic-supersite and other epitopes.

Several studies of SARS-CoV-2 evolution *in vitro* and *ex vivo* also support that NTD indels here observed in Brazilian SARS-CoV-2 VOC and VOI represent a mechanism of ongoing adaptive evolution to escape from dominant neutralizing antibodies directed against the NTD. *In vitro* co-incubation of SARS-CoV-2 with highly neutralizing plasma from COVID-19 convalescent patient, has revealed an incremental resistance to neutralization followed by the stepwise acquisition of indels at N3/N5 loops<sup>35</sup>. SARS-CoV-2 challenge in hamsters previously treated with anti-NTD mAbs revealed selection of two escape mutants harboring NTD deletions  $\Delta 143-144$  and  $\Delta 141-144$ <sup>14</sup>. Studies of intra-host SARS-CoV-2 evolution in immuno-compromised hosts revealed the emergence of viral variants with NTD deletions at RDR1 ( $\Delta 69-70$ ), RDR2 ( $\Delta 144$  and  $\Delta 141-144$ ) and RDR4 ( $\Delta 243-244$ ) following therapy with convalescent plasma<sup>1,3,4,36,37</sup>. Another study revealed the emergence of several Spike gene mutations, including inframe deletions  $\Delta 141-143$ ,  $\Delta 141-144$ ,  $\Delta 145$  and  $\Delta 211-212$ , during persistent SARS-CoV-2 infection in two individuals with partial humoral immunity<sup>38</sup>. Finally, a recent longitudinal analysis of intra-host SARS-CoV-2 evolution during acute infection in one immunocompetent individual revealed the emergence of virus haplotypes bearing deletions  $\Delta 144$  and  $\Delta 141-144$  in the NTD following the development of autologous anti-NTD specific antibodies<sup>39</sup>.

**Table 2. Impact of indels on the binding between SARS-CoV-2 NTDs and NAb 2-51, expressed as loss of putative interactions.**

Variant	$\Delta$ H-bond	$\Delta$ Salt-bridge	$\Delta$ pi-stacking	$\Delta$ Hydrophobic SASA [ $\text{\AA}^2$ ]	Native Contacts Lost
B.1.1.28 $\Delta$ 144	-2	-3	-1	-1	K147-E71 K150-E53 K150-D54 Y145-Y98
P.2 $\Delta$ 144	-2	-3	-1	-104	K147-E71 K150-E53 K150-D54 Y145-Y98
P.1 $\Delta$ 144	-2	-3	-1	-111	K147-E71 K150-E53 K150-D54 Y145-Y98
P.1 $\Delta$ 141-144	-2	-3	-1	-313	K147-E71 K150-E53 K150-D54 Y145-Y98
N.10 $\Delta$ 141-144 $\Delta$ 256-258	-3	-3	-1	-439	Y147-E71 K150-E53 K150-D54 Y145-Y98 D253-S56 P251-P55 P251-L46 P251-Y100
P.1 $\Delta$ 242-244	-3	-1	-2	-746	Y248-Y27 Y248-W96 Y248-Y98 L249-Y27 R246-E31 R246-Y27 D253-S56

Recent genomic findings showed a sudden landscape change in SARS-CoV-2 evolution since October 2020, coinciding with the independent emergence of VOCs carrying multiple convergent amino acid replacements at the RBD of the S protein <sup>40</sup>. One hypothesis is that such a major selection pressure shift on the virus genome is driven by the increasing worldwide human population immunity acquired from natural SARS-CoV-2 infection that might also select for convergent deletions at NTD. Our findings suggest that P.1, P.2 and N.10 variants with NTD indels here detected might have evolved to escape from NAb against NTD and could be even more resistant to



neutralization than the parental viruses. Notably, the sequential acquisition of RBD and NTD mutations observed in the VOC P.1 recapitulates the evolution pattern of the VOC B.1.351 that first acquired RBD mutations E484K and N501Y and sometime later the NTD deletion  $\Delta 242-244$ <sup>7</sup>. The detection of P.1 genomes with convergent NTD deletions with VOCs B.1.1.7 ( $\Delta 69-70$ ,  $\Delta 144$ ) and B.1.351 ( $\Delta 242-244$ ) bring caution about the specificity of published real-time RT-PCR protocols to distinguish different VOCs in Brazil and also alert against use the failure to detect the S gene (due to mutation  $\Delta 69-70$ ) by certain tests, known as S gene target dropout<sup>41,42</sup>, as a definitive proof of circulation of the VOC B.1.1.7 in Brazil.

In summary, these findings suggest that SARS-CoV-2 VOC and VOI are continuously adapting and evolving in Brazil through acquisition of Spike NTD indels. Some variants like P.1  $\Delta 69-70$ , P.1  $\Delta 144$  and P.1-like ins214ANRN might represent newly emergent VOC/VOI and its communitary dissemination, as well as that of VOI N.10, requires careful monitoring. These findings highlight the urgent need to address the SARS-CoV-2 vaccines' efficacy towards emergent SARS-CoV-2 variants carrying both RBD and NTD mutations and deletions of concern and the risk of ongoing uncontrolled community transmission of SARS-CoV-2 in Brazil for the generation of more transmissible variants. The recurrent emergence of NTD ins214 variants in different SARS-CoV-2 lineages circulating in the Americas and Europe since November 2020 and its impact on vaccine efficacy also deserves further attention.

## Material and Methods

### SARS-CoV-2 and SARS-CoV-2-related coronavirus (SC2r-CoV) sequences

Our genomic survey of SARS-CoV-2 positive samples sequenced by the Fiocruz COVID-19 Genomic Surveillance Network between 12th March 2020 and 28th February 2021 identified 11 sequences with mutations of concern in the RBD and indels in the NTD (**Appendix Table 1**). The SARS-CoV-2 genomes were recovered using Illumina sequencing protocols as previously described<sup>43,44</sup>. The FASTQ reads obtained were imported into the CLC Genomics Workbench version 20.0.4 (Qiagen A/S, Denmark), trimmed, and mapped against the reference sequence EPI\_ISL\_402124 available in EpiCoV database in the GISAID (<https://www.gisaid.org/>). The alignment was refined using the InDels and Structural Variants module. Additionally, the same reads were imported in a different pipeline<sup>45</sup> based on Bowtie2 and bcftools<sup>46</sup> mapping and consensus generation allowing us to further confirm the indels supported by

paired-end reads, removing putative indels with less than 10x of sequencing depth and with mapping read quality score below to 10 for all samples sequenced in this study. BAM files were used as input to generate sequencing coverage plots onto indels using the Karyoploter R package <sup>47</sup>. Sequences were combined with SARS-CoV-2 and SC2r-CoV from bats and pangolins available in the EpiCoV database in GISAID by 1st March 2021 (**Appendix Table 1**). This study was approved by the FIOCRUZ-IOC (68118417.6.0000.5248 and CAAE 32333120.4.0000.5190) and the Amazonas State University Ethics Committee (CAAE: 25430719.6.0000.5016) and the Brazilian Ministry of the Environment (MMA) A1767C3.

### **Maximum Likelihood Phylogenetic Analyses**

SARS-COV-2 sequences here obtained were aligned with high quality (<1% of N) and complete (>29 kb) lineages B.1.1.28, P.1, P2 and B.1.1.33 sequences that were available in EpiCoV database in the GISAID (<https://www.gisaid.org/>) at March 1st, 2021 and subjected to maximum-likelihood (ML) phylogenetic analysis using IQ-TREE v2.1.2 <sup>48</sup>. The S amino acid sequences from selected SARS-CoV-2 and SC2r-CoV lineages available in the EpiCoV database were also aligned using Clustal W <sup>49</sup> adjusted by visual inspection.

### **Structural Modeling**

The resolved crystallographic structure of SARS-CoV-2 NTD protein bound to the neutralizing antibody 2-51 was retrieved from the Protein Databank (PDB) under the accession code 7L2C <sup>33</sup>. Missing residues of the chain A, corresponding to the NTD coordinates, were modeled using the user template mode of the Swiss-Model webserver (<https://swissmodel.expasy.org/>) <sup>50</sup> and was used as starting structure for the NTD wildtype. This structure was then used as a template to model the NTD variants using the Swiss-Model webserver. The modeled structures of the NTDs variants were superimposed onto the coordinates of the PDB ID 7L2C to visualize the differences between the NTD-antibody binding interfaces. Image rendering was carried out using Visual Molecular Dynamics (VMD) software <sup>51</sup>. The NTD-antibody complexes were geometry optimized using a maximum of 5,000 steps or until it reached a convergence value of 0.001 REU (Rosetta energy units) using the limited-memory BroydenFletcher-Goldfarb-Shanno algorithm, complying with the Armijo-Goldstein condition, as implemented in the Rosetta suite of software 3.12 <sup>52</sup>. Geometry optimization was accomplished through the atomistic Rosetta energy function 2015 (REF15), while preserving backbone torsion angles. Protein-protein interface analyses

were performed using the Protein Interactions Calculator (PIC) webserver (<http://pic.mbu.iisc.ernet.in/>)<sup>53</sup>, the ‘Protein interfaces, surfaces and assemblies’ service (PISA) at the European Bioinformatics Institute (<https://www.ebi.ac.uk/pdbe/pisa/pistart.html>)<sup>54</sup> and the InterfaceAnalyzer protocol of the Rosetta package interfaced with the RosettaScripts scripting language<sup>55</sup>. For the interfaceAnalyzer, the maximum SASA that is allowed for an atom to be defined as buried is 0.01 Å<sup>2</sup>, with a SASA probe radius of 1.2 Å.

### Epitope prediction

Epitopes in the NTD region were predicted by the ElliPro Antibody Epitope Prediction server<sup>56</sup>. NTD are shown as predicted linear epitopes when using PDB accession codes 6VXX<sup>57</sup> and 6VSB<sup>58</sup>, (structural coordinates corresponding to the entire S protein), along with a minimum score of 0.9, *i.e.*, a highly strict criterion.

### Acknowledgements

The authors wish to thank all the health care workers and scientists who have worked hard to deal with this pandemic threat, the GISAID team, and all the EpiCoV database's submitters, GISAID acknowledgment table containing sequences used in this study are attached to this post (**Appendix Table 3**). We also appreciate the support of the Fiocruz COVID-19 Genomic Surveillance Network (<http://www.genomahcov.fiocruz.br/>) members, the Respiratory Viruses Genomic Surveillance Network of the General Laboratory Coordination (CGLab), Brazilian Ministry of Health (MoH), Brazilian Central Laboratory States (LACENs) and the Amazonas surveillance teams for the partnership in the viral surveillance in Brazil. Financial support was provided by FAPEAM (PCTI-EmergeSaude/AM call 005/2020 and Rede Genômica de Vigilância em Saúde - REGESAM); Ministério da Ciência, Tecnologia, Inovações e Comunicações/Conselho Nacional de Desenvolvimento Científico e Tecnológico - CNPq/Ministério da Saúde - MS/FNDCT/SCTIE/Decit (grants 402457/2020-9 and 403276/2020-9); Inova Fiocruz/Fundação Oswaldo Cruz (Grants VPPCB-007-FIO-18-2-30 and VPPCB-005-FIO-20-2-87) and INCT-FCx (465259/2014-6). Computer allocation was partly granted by the Brazilian National Scientific Computing Center (LNCC). FGN, GLW, RDL and GB are supported by the CNPq through their productivity research fellowships (306146/2017-7, 303902/2019-1,

425997/2018-9 and 302317/2017-1 respectively). G.B. is also funded by the Fundação Carlos Chagas Filho de Amparo à Pesquisa do Estado do Rio de Janeiro – FAPERJ (Grant number E-26/202.896/2018).

## References

- 1 Avanzato VA, Matson MJ, Seifert SN *et al.* Case Study: Prolonged Infectious SARS-CoV-2 Shedding from an Asymptomatic Immunocompromised Individual with Cancer. *Cell* 2020; **183**: 1901-1912.e9.
- 2 Choi B, Choudhary MC, Regan J *et al.* Persistence and Evolution of SARS-CoV-2 in an Immunocompromised Host. *N Engl J Med* 2020; **383**: 2291–2293.
- 3 McCarthy KR, Rennick LJ, Nambulli S *et al.* Recurrent deletions in the SARS-CoV-2 spike glycoprotein drive antibody escape. *Science* 2021; **6950**: 6–6.
- 4 Kemp SA, Collier DA, Datir RP *et al.* SARS-CoV-2 evolution during treatment of chronic infection. *Nature* 2021; : 1–10.
- 5 Spike protein mutations in novel SARS-CoV-2 ‘variants of concern’ commonly occur in or near indels. Virological. 2021.<https://virological.org/t/spike-protein-mutations-in-novel-sars-cov-2-variants-of-concern-commonly-occur-in-or-near-indels/605> (accessed 14 Mar2021).
- 6 Spike protein sequences of Cambodian, Thai and Japanese bat sarbecoviruses provide insights into the natural evolution of the Receptor Binding Domain and S1/S2 cleavage site. Virological. 2021.<https://virological.org/t/spike-protein-sequences-of-cambodian-thai-and-japanese-bat-sarbecoviruses-provide-insights-into-the-natural-evolution-of-the-receptor-binding-domain-and-s1-s2-cleavage-site/622> (accessed 14 Mar2021).
- 7 Tegally H, Wilkinson E, Giovanetti M *et al.* Emergence of a SARS-CoV-2 variant of concern with mutations in spike glycoprotein. *Nature* 2021; : 1–8.
- 8 Preliminary genomic characterisation of an emergent SARS-CoV-2 lineage in the UK defined by a novel set of spike mutations - SARS-CoV-2 coronavirus / nCoV-2019 Genomic Epidemiology. Virological. 2020.<https://virological.org/t/preliminary-genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-the-uk-defined-by-a-novel-set-of-spike-mutations/563> (accessed 14 Mar2021).
- 9 Wang R, Zhang Q, Ge J *et al.* Spike mutations in SARS-CoV-2 variants confer resistance to antibody neutralization. *bioRxiv* 2021; : 2021.03.09.434497.
- 10 Wang P, Nair MS, Liu L *et al.* Antibody Resistance of SARS-CoV-2 Variants B.1.351 and B.1.1.7. *bioRxiv* 2021; : 2021.01.25.428137.
- 11 Collier DA, De Marco A, Ferreira IATM *et al.* Sensitivity of SARS-CoV-2 B.1.1.7 to mRNA vaccine-elicited antibodies. *Nature* 2021; : 1–8.
- 12 Gobeil S, Janowska K, McDowell S *et al.* Effect of natural mutations of SARS-CoV-2 on spike structure, conformation and antigenicity. *bioRxiv* 2021; : 2021.03.11.435037.
- 13 Wang P, Wang M, Yu J *et al.* Increased Resistance of SARS-CoV-2 Variant P.1 to Antibody Neutralization. *bioRxiv* 2021; : 2021.03.01.433466.
- 14 McCallum M, Marco AD, Lempp FA *et al.* N-terminal domain antigenic mapping



- reveals a site of vulnerability for SARS-CoV-2. *Cell* 2021; **0**. doi:10.1016/j.cell.2021.03.028.
- 15 Genomic characterisation of an emergent SARS-CoV-2 lineage in Manaus: preliminary findings - SARS-CoV-2 coronavirus / nCoV-2019 Genomic Epidemiology. Virological. 2021.<https://virological.org/t/genomic-characterisation-of-an-emergent-sars-cov-2-lineage-in-manaus-preliminary-findings/586> (accessed 14 Mar2021).
- 16 Voloch CM, da Silva Francisco R, de Almeida LGP *et al*. Genomic characterization of a novel SARS-CoV-2 lineage from Rio de Janeiro, Brazil. *J Virol* 2021. doi:10.1128/JVI.00119-21.
- 17 Resende PC, Gräf T, Paixão ACD *et al*. A potential SARS-CoV-2 variant of interest (VOI) harboring mutation E484K in the Spike protein was identified within lineage B.1.1.33 circulating in Brazil. *bioRxiv* 2021; : 2021.03.12.434969.
- 18 Identification of a new B.1.1.33 SARS-CoV-2 Variant of Interest (VOI) circulating in Brazil with mutation E484K and multiple deletions in the amino (N)-terminal domain of the Spike protein - SARS-CoV-2 coronavirus / nCoV-2019 Genomic Epidemiology. Virological. 2021.<https://virological.org/t/identification-of-a-new-b-1-1-33-sars-cov-2-variant-of-interest-voi-circulating-in-brazil-with-mutation-e484k-and-multiple-deletions-in-the-amino-n-terminal-domain-of-the-spike-protein/675> (accessed 7 Apr2021).
- 19 Dejnirattisai W, Zhou D, Supasa P *et al*. Antibody evasion by the P.1 strain of SARS-CoV-2. *Cell* 2021; **0**. doi:10.1016/j.cell.2021.03.055.
- 20 Naveca, Felipe, Nascimento, Valdinete, Souza, Victor *et al*. COVID-19 epidemic in the Brazilian state of Amazonas was driven by long-term persistence of endemic SARS-CoV-2 lineages and the recent emergence of the new Variant of Concern P.1. *Res Sq Prepr Platf Makes Res Commun Faster Fairer More Useful* 2021. doi:10.21203/rs.3.rs-275494/v1.
- 21 Tosta S, Giovanetti M, Nardy VB *et al*. Early genomic detection of SARS-CoV-2 P.1 variant in Northeast Brazil. *medRxiv* 2021; : 2021.02.25.21252490.
- 22 Wacharapluesadee S, Tan CW, Maneeorn P *et al*. Evidence for SARS-CoV-2 related coronaviruses circulating in bats and pangolins in Southeast Asia. *Nat Commun* 2021; **12**: 972.
- 23 Liu L, Wang P, Nair MS *et al*. Potent neutralizing antibodies against multiple epitopes on SARS-CoV-2 spike. *Nature* 2020; **584**: 450–456.
- 24 Voss WN, Hou YJ, Johnson NV *et al*. Prevalent, protective, and convergent IgG recognition of SARS-CoV-2 non-RBD spike epitopes in COVID-19 convalescent plasma. *bioRxiv* 2020; : 2020.12.20.423708.
- 25 Piccoli L, Park Y-J, Tortorici MA *et al*. Mapping Neutralizing and Immunodominant Sites on the SARS-CoV-2 Spike Receptor-Binding Domain by Structure-Guided High-Resolution Serology. *Cell* 2020; **183**: 1024-1042.e21.
- 26 Barnes CO, West AP, Huey-Tubman KE *et al*. Structures of Human Antibodies Bound to SARS-CoV-2 Spike Reveal Common Epitopes and Recurrent Features of Antibodies. *Cell* 2020; **182**: 828-842.e16.
- 27 Greaney AJ, Loes AN, Crawford KHD *et al*. Comprehensive mapping of mutations in the SARS-CoV-2 receptor-binding domain that affect recognition by polyclonal human plasma antibodies. *Cell Host Microbe* 2021; **29**: 463-476.e6.
- 28 Hoffmann M, Arora P, Groß R *et al*. SARS-CoV-2 variants B.1.351 and B.1.1.248:

- Escape from therapeutic antibodies and antibodies induced by infection and vaccination. *bioRxiv* 2021; : 2021.02.11.430787.
- 29 Baum A, Fulton BO, Wloga E *et al.* Antibody cocktail to SARS-CoV-2 spike protein prevents rapid mutational escape seen with individual antibodies. *Science* 2020; **369**: 1014–1018.
  - 30 Nelson G, Buzko O, Spilman P, Niazi K, Rabizadeh S, Soon-Shiong P. Molecular dynamic simulation reveals E484K mutation enhances spike RBD-ACE2 affinity and the combination of E484K, K417N and N501Y mutations (501Y.V2 variant) induces conformational change greater than N501Y mutant alone, potentially resulting in an escape mutant. *bioRxiv* 2021; : 2021.01.13.426558.
  - 31 Ferraz M, Moreira E, Coêlho DF, Wallau G, Lins R. SARS-CoV-2 VOCs Immune Evasion from Previously Elicited Neutralizing Antibodies Is Mainly Driven by Lower Cross-Reactivity Due to Spike RBD Electrostatic Surface Changes. 2021. doi:10.26434/chemrxiv.14343743.v1.
  - 32 Chi X, Yan R, Zhang J *et al.* A neutralizing human antibody binds to the N-terminal domain of the Spike protein of SARS-CoV-2. *Science* 2020; **369**: 650–655.
  - 33 Cerutti G, Guo Y, Zhou T *et al.* Potent SARS-CoV-2 neutralizing antibodies directed against spike N-terminal domain target a single supersite. *Cell Host Microbe* 2021. doi:10.1016/j.chom.2021.03.005.
  - 34 Wibmer CK, Ayres F, Hermanus T *et al.* SARS-CoV-2 501Y.V2 escapes neutralization by South African COVID-19 donor plasma. *Nat Med* 2021; : 1–4.
  - 35 Andreano E, Piccini G, Licastro D *et al.* SARS-CoV-2 escape in vitro from a highly neutralizing COVID-19 convalescent plasma. *bioRxiv* 2020; : 2020.12.28.424451.
  - 36 Khatamzas E, Rehn A, Muenchhoff M *et al.* Emergence of multiple SARS-CoV-2 mutations in an immunocompromised host. *medRxiv* 2021; : 2021.01.10.20248871.
  - 37 Chen L, Zody MC, Mediavilla JR *et al.* Emergence of multiple SARS-CoV-2 antibody escape variants in an immunocompromised host undergoing convalescent plasma treatment. *medRxiv* 2021; : 2021.04.08.21254791.
  - 38 Truong TT, Ryutov A, Pandey U *et al.* Persistent SARS-CoV-2 infection and increasing viral variants in children and young adults with impaired humoral immunity. *medRxiv* 2021; : 2021.02.27.21252099.
  - 39 Ko SH, Mokhtari EB, Mudvari P *et al.* High-throughput, single-copy sequencing reveals SARS-CoV-2 spike variants coincident with mounting humoral immunity during acute COVID-19. *PLOS Pathog* 2021; **17**: e1009431.
  - 40 Martin DP, Weaver S, Tegally H *et al.* The emergence and ongoing convergent evolution of the N501Y lineages coincides with a major global shift in the SARS-CoV-2 selective landscape. *medRxiv* 2021; : 2021.02.23.21252268.
  - 41 Korukluoglu G, Kolukirik M, Bayrakdar F *et al.* 40 minutes RT-qPCR Assay for Screening Spike N501Y and HV69-70del Mutations. *bioRxiv* 2021; : 2021.01.26.428302.
  - 42 Bal A, Destras G, Gaymard A *et al.* Two-step strategy for the identification of SARS-CoV-2 variant of concern 202012/01 and other variants with spike deletion H69–V70, France, August to December 2020. *Eurosurveillance* 2021; **26**: 2100008.
  - 43 Nascimento VA do, Corado A de LG, Nascimento FO do *et al.* Genomic and phylogenetic characterisation of an imported case of SARS-CoV-2 in Amazonas State, Brazil. *Mem Inst Oswaldo Cruz* 2020; **115**. doi:10.1590/0074-02760200310.
  - 44 Resende PC, Motta FC, Roy S *et al.* SARS-CoV-2 genomes recovered by long

- amplicon tiling multiplex approach using nanopore sequencing and applicable to other sequencing platforms. *bioRxiv* 2020; : 2020.04.30.069039.
- 45 Paiva MHS, Guedes DRD, Docena C *et al.* Multiple Introductions Followed by Ongoing Community Spread of SARS-CoV-2 at One of the Largest Metropolitan Areas of Northeast Brazil. *Viruses* 2020; **12**: 1414.
  - 46 Li H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinforma Oxf Engl* 2011; **27**: 2987–2993.
  - 47 Gel B, Serra E. karyoploteR: an R/Bioconductor package to plot customizable genomes displaying arbitrary data. *Bioinformatics* 2017; **33**: 3088–3090.
  - 48 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol Biol Evol* 2015; **32**: 268–274.
  - 49 Larkin M a, Blackshields G, Brown NP *et al.* Clustal W and Clustal X version 2.0. *Bioinforma Oxf Engl* 2007; **23**: 2947–8.
  - 50 Waterhouse A, Bertoni M, Bienert S *et al.* SWISS-MODEL: homology modelling of protein structures and complexes. *Nucleic Acids Res* 2018; **46**: W296–W303.
  - 51 Humphrey W, Dalke A, Schulten K. VMD: Visual molecular dynamics. *J Mol Graph* 1996; **14**: 33–38.
  - 52 Leaver-Fay A, Tyka M, Lewis SM *et al.* ROSETTA3: an object-oriented software suite for the simulation and design of macromolecules. *Methods Enzymol* 2011; **487**: 545–574.
  - 53 Tina KG, Bhadra R, Srinivasan N. PIC: Protein Interactions Calculator. *Nucleic Acids Res* 2007; **35**: W473–476.
  - 54 Krissinel E, Henrick K. Inference of Macromolecular Assemblies from Crystalline State. *J Mol Biol* 2007; **372**: 774–797.
  - 55 Fleishman SJ, Leaver-Fay A, Corn JE *et al.* RosettaScripts: A Scripting Language Interface to the Rosetta Macromolecular Modeling Suite. *PLOS ONE* 2011; **6**: e20161.
  - 56 Ponomarenko J, Bui H-H, Li W *et al.* ElliPro: a new structure-based tool for the prediction of antibody epitopes. *BMC Bioinformatics* 2008; **9**: 514.
  - 57 Wrapp D, Wang N, Corbett KS *et al.* Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation. *Science* 2020; **367**: 1260–1263.
  - 58 Walls AC, Park Y-J, Tortorici MA, Wall A, McGuire AT, Veesler D. Structure, Function, and Antigenicity of the SARS-CoV-2 Spike Glycoprotein. *Cell* 2020; **181**: 281–292.e6.