medRxiv preprint doi: https://doi.org/10.1101/2021.11.21.21266655; this version posted November 24, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity.

perpetuity. It is made available under a CC-BY-NC 4.0 International license .

Determining international spread of novel B.1.1.523 SARS-CoV-2 lineage

2

Lukas Zemaitis^{1, *+}, Gediminas Alzbutas²⁺, Dovydas Gecys¹, Andrey Komissarov³, Arnoldas Pautienius⁴ Rasa Ugenskiene⁵ Marius Sukys⁵ and Vaiva Lesauskaite¹

- ⁵ ¹Lithuanian University of Health Sciences, Institute of Cardiology, Laboratory of Molecular
- 6 Cardiology, Kaunas, LT-50162, Lithuania
- ⁷² Lithuanian University of Health Sciences, Institute for Digestive Research, Laboratory of
- 8 Translational Bioinformatics, Kaunas LT-50162, Lithuania
- 9 ³ Smorodintsev Research Institute of Influenza, Saint Petersburg, 197376, Russian Federation
- ⁴ Lithuanian University of Health Sciences, Institute of Microbiology and Virology, Kaunas LT-
- **11** 47181, Lithuania
- ⁵ Lithuanian University of Health Sciences Hospital Kaunas Clinics, Genetics and Molecular
- 13 Medicine Clinic, Kaunas LT- 50161, Lithuania
- 14 * lukas.zemaitis@lsmuni.lt
- 15 + these authors contributed equally to this work
- 16

17 ABSTRACT

- 18 Here we report the emergence of variant lineage B.1.1.523 that contains a set of mutations including
- 19 156_158del, E484K and S494P in Spike protein. E484K and S494P are known to significantly reduce
- 20 SARS-CoV-2 neutralization by convalescent and vaccinee sera and are considered as mutations of
- 21 concern. Lineage B.1.1.523 has presumably originated in Russian Federation and spread across
- 22 European countries with the peak of transmission in April May 2021. The B.1.1.523 lineage has
- 23 now been reported from 27 countries.

24 INTRODUCTION

25 The emergence of severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) in late 2019 led to the ongoing Coronavirus Disease 2019 (COVID -19), now a global pandemic with more 26 27 than 230 million cases of infection and about 5 million deaths worldwide ¹. On January 5, 2020, the first whole genome sequence of 2019-nCoV was completed by Wuhan Institute of Virology, China 28 29 Centre for Disease Control and Shanghai Public Health Clinical Centre of Fudan University². From this point on, genome sequencing has played an important role in vaccine development, 30 31 understanding viral evolution and epidemiological characteristics. In many countries, SARS-CoV-2 sequencing has been implemented at the national level as a tool for epidemiological management 32 33

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

A large dataset of SARS-CoV-2 genomes has been collected in the GISAID database, which now 34 contains more than 3.7 million sequenced genomes from around the world ⁴. As of May 31, 2021, 35 World Health Organisation (WHO) has proposed designations for global SARS-CoV-2 variants of 36 concern (VOC) and variants of interest (VOI) to be used alongside scientific nomenclature in 37 communications about variants to the public ⁵. This list includes the variants on the global list of 38 WHO VOC and VOI and will be updated as the list of WHO changes. There are currently three SARS-39 CoV-2 VOCs: Beta, Gama and Delta. The variant previously classified as Alpha (B.1.1.7) has been 40 reclassified as de-escalated due to the drastic decrease in prevalence in the EU/EEA⁶. 41

Through the routine analysis of National Lithuanian sequencing results from national 42 43 sequencing efforts coordinated by National Public Health Surveillance Laboratory, we have identified a novel SARS-CoV-2 variant classified as B.1 by PANGO but containing multiple S protein 44 mutations associated with effects on immunity (https://github.com/cov-lineages/pango-45 designation/issues/69), such as E484K; 156_158del; S494P. Preliminary phylogenetic analysis 46 indicated that this variant has a distinct viral lineage that may have originated in Russia. We reported 47 this variant to the PANGO curators and gave it the new phylum name B.1.1.523 48 (https://github.com/cov-lineages/pango-designation/issues/69). On 14-July-2021, WHO added this 49 variant to the list of variants under Monitoring section. By using bioinformatics tools, we performed 50 a detailed analysis of this lineage and disclosed our findings about this variant or the mutational 51 52 subgroups typical of this variant.

With this report we aim to share a detailed analysis of discovered variant, evaluate its origin 53 as well as predict potential epidemiological impact and risks. 54

55

RESULTS 56

57

Mutation review of B.1.1.523 58

59

Several mutations in S region have been observed in B.1.1.523 variant, from which 60 156_158del, E484K, and S494P are considered as an attribute for VOCs (Fig. 1). According to 61 previously reported data, the 156_157del and G158R mutations in the Delta variant are matching to 62 the same surface as the 144 and 241–243 deletions in the Alpha and Beta (B.1.351) variants, 63 respectively. These altered residues are found in the NTD 'supersite' that is targeted by most anti-64 NTD neutralizing antibodies, thus providing a mode to dodge immune system ⁷. Moreover, E484K 65 mutation also contributes to SARS-CoV-2 immune system evasion. Several recent studies have 66 observed that E484K may significantly reduce convalescent serum neutralization ^{8,9}. Additionally, it 67 was observed that S494P mutation is related to 3-5-fold reduced SARS-CoV-2 neutralization in sera), 68 however, this mutation was not as potent at neutralization as E484K^{8,10}. With a combination of 69 156_158del, E484K, and S494P mutations, B1.1.523 lineage should remain on epidemiologists 70 71 watchlist as one of the most concerning SARS-CoV-2 lineages.

In addition to 156_158del, E484K, and S494P more than 70% of genomes attributed as 72 73 B.1.1.523 lineage possess F306L, D839V and T1027I in Spike and a set of substitutions in ORF1a 74 (NSP3:M84V, R1297I; NSP2: N269D; NSP5: V303I; NSP6: V84F; NSP10:T111I), ORF1b (NSP12:S229N,





93



medRxiv preprint doi: https://doi.org/10.1101/2021.11.21.21266655; this version posted November 24, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity.

perpetuity. It is made available under a CC-BY-NC 4.0 International license .

Figure 2. The overlap between the two data sets used for focused ML tree. The sequences were
either chosen based on Pango assignment or by the identity with a B.1.1.523 lineage sequence
(EPI_ISL_1590462). Most of the sequences, with high identity (> 0.993) to the Latvian B.1.1.523
lineage, were classified as belonging to the B.1.1.523. However, 21 (6%) sequences were not
assigned to the B.1.1.523.

- 101
- 102

Generation of ML revealed several interesting properties of B.1.1.523 (Suppl. Fig. 1). At the base of the lineages leading to the B.1.1.523 sequences having a full set of expected S protein mutations branches away in clusters of sequences having the triple S:156_158del deletion. The sequences having the additional substitutions at S:484 and S:494 positions emerge further in the evolution. However, here, we have no clear indications of that the mutations S:E484K, S:S494P are introduced sequentially to form the B.1.1.523 lineage.

109 We have observed that some sequences which originated from progenitors with full set of 110 expected substitutions at 484, 494, 156, 157, 158 S region positions, have underwent a reverse-type mutation to wild type variants. Such events are highly unlikely and usually are caused by erroneous 111 sequence assembly or could be an indication of low-guality data. Additionally, in order to discern 112 the cases where the mutations comprising the two regions of potential enhanced resistance to 113 immune response have been combined phylogeny analysis of all S protein unique variants was 114 performed. As it is depicted in Figure 3, there are at least three distinct cases where E484K and/or 115 S494P have been combined with the 156_158del. Phylogeny analysis of S protein indicated that 116 lineage assignments by Pango sometimes can be misleading. Furthermore, relying on plain 117 118 assignments could hide novel developments of SARS-CoV-2. For example, some sequences from Turkey that were assigned as lineage B.1.1.523 originated from distinct clusters based on S sequence 119 (Fig. 4A) and were evidently evolutionary distant from other B.1.1.523 lineage genomes. 120





Figure 3. Phylogeny based on S protein sequence. (A) The tree represents a maximum likelihood 124 tree based on all unique S protein sequences of the genomes deposited to the GISAID. The visible 125 subset of the tree matches lineages that lead to branches having 156_158del and E484K or S494P 126

mutations. The arrow "->" indicates haplotype transitions detected by comparison of parental 127 sequences with their offspring variants. The five-letter haplotype strings match 156, 157, 158, 484, 128 494 positions of the S protein with "." meaning the wild type. The green ovals and their number 129 indicate prominent transitions that are discussed in the main text. (B) The lower tree matches 130 maximum likelihood tree based on whole genomes of the cases visualised in the upper tree. The 131 132 black solid lines indicate a match of nodes for two B.1.1.523 sequences from Turkey in the phylogenies based on S protein and whole genome. 133

Transitions E1 (introduction of the triple deletion) and E2 (introduction of E484, S494P) (Fig. 134 135 3) indicate a pathway where the majority of B.1.1.523 lineage sequences took. As in the case of data 136 given in (Fig. 3) we do not detect sequential acquisition of E484K or S494P, and immediately next to the triple deletion we see the two aforementioned additional mutations. This could indicate a 137 potential recombination event, but it is most probably due to lack of insufficient data. 138

139 An interesting case is with two highly diverged Turkey sequences that are classified by Pango 140 as B.1.1.523 lineage. In this case following the most likely scenario at first the E484K was acquired and only then the S494P and E156 R158del were acquired (Fig. 4). As in the case discussed before, 141 the sequential acquire steps are lacking from the data. Most probably the sequencing data was not 142 comprehensive enough to reveal the full picture on how the combination was formed. In addition, 143 these two highly diverged sequences indicate that there exist vast uncharted territories of COVID-144 19 evolution as we get only sparse sequencing data from Central Asia regions where the COVID-19 145 infection are barely controlled. An evident third case emergent combination of immune response 146 147 hindering mutations from distinct S protein regions are highlighted by E5 and E6 at Figure 3A. In this 148 case within the B.1.351 lineage at first the E494K mutation was introduced and then E156 R158del 149 followed. The presented data evidently shows that the immune response hindering mutations from RBD and NTD domains have been combined in one protein at least three times. In two of them 150 151 E484K occurred first, in one case the first one was the triple deletion. Most probably these are results of independent mutational events. 152

153



154

155

Figure 4. S protein region of B.1.1.523 sequences from Turkey. Turkey/HSGM-B11599/2021 and 156 Turkey/HSGM-B11931/2021 were classified by Pangolin as belonging to B.1.1.523 lineage. The top 157 figure is a snapshot from the Nextclade analysis. The top sequence is the reference sequence hCoV-158 19 that is used by GISAID. The asterisks ** indicate sequence from Turkey variants that have their 159 two most variable sequences swapped with corresponding regions from the refence sequence. The 160 two lower graphs indicate the sequence alignments from the two extremely variable fragments. The 161 order of sequences is the same as in the upper graph: reference sequence, three Turkey sequences 162 with swapped fragments, three original Turkey sequences. 163

164





166 Figure 5. Phylogeny based on S protein sequence with modified sequences from Turkey. The cladogram of maximum likelihood tree that includes the three pairs of Turkey sequences with their 167 168 variable regions either swapped with counterparts of reference sequences or being left original. Colour of the tips of leaves indicates if they are being classified as B.1.1.523. Tip labels indicates 169 170 their Pangolin assignment with colours indicating different lineages. The different tips connecting grey lines indicates a Turkey sequence. The asterisks ("**") in front of the sequence name labels 171

indicates the sequence variant where the variable regions were swapped with corresponding 172 regions from the reference sequences. Next to the sequence labels the haplotype at positions 156-173 158,484,494 with "." indicating the wild type and "-" - a gap. A Turkey sequences. 174

175 In Figure 5 the Turkey sequences are highly diverged from the wild type based on the region 176 variation of S protein. Based on the GISAID blast searches (performed at 2021 10 1 07) the first 177 variable region (bottom left, Fig. 5) from Turkey/HSGM-B11599/2021 was found in 7 sequences and 178 the second variable region (bottom right, Fig. 5) was found in 4 sequences. The sequences have been deposited in the GISAID during several submissions. Other genomic regions of the Turkey 179 180 sequences showed in the Fig. 4 do not contain such large SNPs clusters as the ones showed in the S protein. This indicates the changes in the placement of the three Turkey sequences upon switching 181 182 the two extremely deviated fragments with corresponding fragments with reference sequence.

183

B.1.1.523 spread worldwide 184

185

186 August 31st, 2021, 459 B.1.1.523 sequences spread across 27 countries have been Till published in the GISAID. According to analysis results, B1.1.523 has originated in Russian Federation 187 188 and spread across European countries (Fig 6). The sequenced clades peaked at week 25 and then subsided. In total, 95 transmission clusters have been identified. The peak of B.1.1.523 transmission 189 190 intensity was around April - May 2021. The most numerous transmission clusters were detected for with MRCA's originating from Germany and Russia. 191

192



193

Figure 6. The distribution of cases of the lineage B.1.1.523 across countries at different time 194 points. The "0" time point indicates the date of the earliest lineage sequence uploaded on the 195 GISAID database. Only sequences that have the typical set of S mutations were considered (E484K, 196

S494P, 156_158del). Only the cases which are corresponding to top 12 countries with the most 197 abundant detection rate are included in the underlying data. The top 12 countries correspond to 198 the 93 % of all cases. 199

- 200
- 201



203

204 Figure 7 B.1.1.523 Transmission clusters identified using the phylogeny tree. Sequences with 205 identity larger or equal to 99.3 % to the Latvian B.1.1.523 sequence (EPI ISL 1590462) were added to the analysis. Colours indicate sequences belonging to the same transmission cluster. Grey colour 206 207 marks sequences that were not assigned to any transmission colour (A). Colour encodes the country 208 of the most recent common ancestor (MRCA) of all sequences that constitute a cluster. Y axis represents size of a cluster and Y axis denotes the date inferred for the MRCA sequence (B). 209

210

211 Most of the cases transmission origin country (country inferred for a cluster MRCA sequence) 212 and the target country (country of a sequence indicated in the GISAID metadata) was Russia and as 213 we see in the Figure 6B majority of the transmission events happened within Russia. Evidently The 214 largest number of cases were transmission origin and target countries were different resembles 215 transmission from Russia to Germany. The data indicate one reverse type (Germany to Russia) 216 transmission. As of the date of writing this article, Germany can be considered as a reservoir of the 217 B.1.1.523 lineage. As for May 2021 the B.1.1.523 and Delta variants represented 0.002% and 0.054% respectively, while the data of August 2021 shows a steady increase in detected cases to 0.327% 218 219 and 32.358% respectively. The overall number of sequenced cases in Germany is becoming 220 overwhelmed by Delta variant, however the lineage B.1.1.523 lineage seems not to be fading away 221 (Fig. 8). Data shows that B.1.1.523 is able to steadily spread without interference with the Delta 222 variants and can be considered as a predominant background lineage of SARS-CoV-2. As indicated 223 in the Figure 6. the spread of the lineage (based on GISAID submissions) in Russia diminished while 224 it remains evidently circulating in Germany. It is unclear whether B.1.1.523 is completely out of

circulation in Russia or its frequency dropped below the sensitivity threshold of SARS-CoV-2 genomic 225

surveillance in this country. 226



227

Figure 8. Newly sequenced cases of B.1.1.523 and Delta lineages in Germany. Data number of new 228 229 cases per month.is based on GISAID metadata (August 31st, 2021). Numbers indicate the total number of sequences deposited in GISAID. 230

231

232 B.1.1.523 antibody escape

233 In Figure 9A the modelled free energy of complex separation is given as calculated for the 234 three sequence variants considering the NTD-antibody complex. As Rosetta documentation states¹¹ ΔG of complex separation (dG separated) shows the change in Rosetta energy when the interface 235 236 forming chains are separated, versus when they are complexed - therefore the lower the value the more energetically favourable is complex separation and less favourable is complex formation. In 237 238 other words - the higher the value - the less likely should be the neutralization by the antibody. Pairwise comparisons using Wilcoxon rank sum test with Benjamini-Hochberg correction for 239 240 multiple comparisons indicated that all three cases significantly differ from each other (p values are 241 given in Fig. 9A upper right corner).



243

Figure. 9. Escape effects of the d156 158 mutation (B.1.1.523) and del156 157&R158G mutations 244 (delta variant) based on NTD-directed neutralizing antibody 4-8 Fab (PDBID: 7LQV). Predicted 245 complex separation ΔG values for the mutants and the wild type complexes (A). The relationship 246 between structure deviation from the starting structure used for the docking with SnugDock and 247 248 the I sc score (B).

All calculated FoldX binding energies are given in the supplementary file FileS1.xlsx. At least 249 250 in four cases a significant synergy in the effect of the E484 and S494P mutations were observed: notably for the antibodies H11-H4 6ZH9, H11-D4 6YZ5, Sb45 7KGJ, Sb16 7KGK (Fig 10A). For the 251 most part the most significant antibody escape effect was noticed by the E484K mutation; however, 252 in some cases the S494P effect was also prominent. The largest effect by the S494P mutation was 253 observed for the E 7KN5¹² antibody structure (Fig. 10B): the S494P mutation increased the binding 254 $\Delta\Delta$ G by 74 % and the additional effect of E484 mutation was negligible (Fig. 10A). 255



	F 40 41/	۵۵G,%		AAAC %
	L404N	5494P	101 67%	
H11-D4 6YZ5	40.29% 37.09%	57.59%	95.02%	35.60% 4 cases of
Sb45_7KGJ	23.86%	29.36%	49.03%	19.68% synergy
Sb16_7KGK	14.40%	11.18%	24.93%	10.53%
E_7KN5	1.09%	73.26%	73.51%	0.26%
P2C-1A3_7CDJ	-0.04%	0.00%	0.18%	0.18%
CV30_6XE1	0.34%	0.15%	0.49%	0.15%
BD23_7BYR	0.00%	0.00%	0.10%	0.10%
CV07-250_6XKQ	0.00%	0.00%	0.00%	0.00%

256 Figure 10. Escape effects of the E484, S494P mutations and their combination. The ∆∆G values 257

258 indicates relative increase in binding energy compared to the wild type structure as inferred from the FoldX calculations. The $\Delta\Delta\Delta G$ indicates the minimum difference between the $\Delta\Delta G$ of the double 259 260 mutation and any of the two single point mutations. The large the value - the large the synergy (A). The structure of the antibody and receptor binding domain of the S protein complex which was 261 262 affected by the S494P mutation most significantly (PDB ID: 7KN5) (B).

263 DISCUSSION

264 We have identified a new SARS-CoV-2 virus lineage with multiple mutations associated with immune escape and reported this to Pango at 5'th of May [https://github.com/cov-lineages/pango-265 designation/issues/69], which mandates the new lineage designation B.1.1.523. This Lineage was 266 first determined in March 2021 and at the time of writing this article, the total amount of cases has 267 reached 598 over 32 countries 12. It is likely that the rapid increase in circulation of Delta variant 268 269 could have diminished the rise of B.1.1.523 lineage, however, the spread of the novel SARS-CoV-2 270 lineage not only has not ceased, but even has started to rise.

271 Currently, a vast growth of B.1.1.523 can now be observed in Germany. Interestingly, the 272 transmission of this lineage has diminished in Russia, where it was most expected to rise. This can 273 be explained by different diagnostic strategy approaches in Russian Federation, where the testing 274 is performed on non-randomly selected sources in the country. Alternatively, this could be explained by the steep rise of Delta variant in Russia, which started a month earlier than in Europe. E.g. in mid 275 June we had >80% delta, while in Germany the same frequency of Delta was observed only in mid 276 277 July.

278 The B.1.1.523 lineage possesses three or more mutations that characterizes SARS-CoV-2 VOCs, including S:D156-158 deletion, S:E484K and S:S494P. D156-158 deletion at β-hairpin antigenic 279 supersite, that is located at the same region typical for the Delta variant (E156G and 157-158del)¹³. 280 E484K mutation has been detected in Beta variant (B.1.351) and VUM Zeta (B.1.1.28). The mutation 281 is in the genomic region coding SARS-CoV-2 spike protein, and it appears to have a significant impact 282 on the body's immune response and possibly, vaccine efficacy. On February 1st, Public Health 283

England (PHE) announced that the Covid-19 Genomics (COG-UK) consortium had identified this 284 same E484K mutation in 11 samples carrying the UK variant B.1.1.7 (sometimes called the Kent 285 variant), after analysing 214 159 sequences^{14.} 286

An in vitro study for SARS-CoV-2 spike protein mutations that are responsible for antibody 287 evasiveness has identified that S494P mutation 15 reduce SARS-CoV-2 neutralization by 3-5-fold in 288 some convalescent sera. However, this mutation was not as potent at neutralization as E484K 16. 289 The results show that S494P mutation increases the spike protein stability. Also, applying docking 290 by HADDOCK displayed higher binding affinity to hACE2 for mutant spike than wild type possibly 291 292 due to the increased β-strand and turn secondary structures which increases surface accessibly 293 surface area (SASA) and chance of interaction. Currently deposited sequences in GISAID do not 294 support hypothesis that S:156 158del were combined with S:E484K and S:S494P during a recombination event; rather it looks like that initially the triple deletion was introduced and then 295 followed addition of S:E484K and S:S494P. 296

We have showed by molecular modelling that in at least one case of antibody the triple 297 deletion del156-158 could decrease interaction. The combination with other immune escape 298 enhancing mutation at RBD this could result in highly resistant variant to immunity that was formed 299 by the initial virus variants. Delta variant also poses sequence changes at the S protein residues 156-300 158 that can induce immune escape and recombination with the B.1.1.523 variant or de novo 301 introduction of the N484K and S494P mutations could make the Delta variant even more dangerous. 302 The case with Turkey sequences is controversial. The sequences at the S region significantly deviated 303 fragments with many "private" mutations. These sequences have been seen across several 304 submissions; hence, these might be not artificial. If these sequences are not artifacts, then the 305 306 Turkey sequences classified by Pangolin as belonging to B.1.1.523 lineage, resulted after a recombination event between a highly diverged Turkey variant and a typical B.1.1.523 lineage 307 sequence. Turkey has been an extensive place of the virus spread with limited control and 308 reportability ¹⁷; therefore, highly diverged variants could have evolved. 309

The results indicate that Pangolin classification should not be taken with granted. Out of the 310 two Turkey sequences that were assigned to the B.1.1.523 lineage, only one had characteristic to 311 the lineage SNP's at the S protein, the other one (Turkey/HSGM-B11931/2021) has mutations 312 characteristic for the Delta variant (double deletion and G residue at 156-158 region). 313

314

CONCLUSIONS 315

316 Presence and spread of SARS-CoV-2 B.1.1.523 lineage is evident regardless of the rapid spread of the delta variant. This variant needs to be carefully observed and studied to keep a look 317 318 out for new mutations, that may cause even more harm in the Covid-19 pandemic. It is also important to monitor other SARS-CoV-2 variants to keep track if this or similar mutations occur 319 320 spontaneously or by recombination.

- 321
- 322

METHODS 323

Collection of SARS-CoV-2 sequences and initial data processing 324

Sequences used for the analyses were downloaded from GISAID¹⁸ as for date of August 31st, 325 2021. metadata were extracted 326 Fasta files and using ncov-ingest tool [https://github.com/nextstrain/ncov-ingest (cloned at 2021 04 19)]. Lineages for all downloaded 327 sequences were assigned using pangolin 3.1.11 (pangoLEARN 2021-08-24 and pango-designation 328 v1.2.66)¹⁹ [https://github.com/cov-lineages/pangolin]. 329

General sequence quality evaluation, extraction and alignment of S protein sequences, 330 variant calling was performed using Nextclade 1.3.0²⁰ [https://github.com/nextstrain/nextclade]. 331

332

333 Transmission cluster analysis

334 In order to elucidate potential origin of the lineage and transmission cluster, a phylogeny analysis of full genomes representing a small subset of GISAID has been done. The chosen sequences 335 for analysis composed from the union of two sets of sequences: (i) sequences that were assigned 336 B.1.1.523 lineage by pangolin, (ii) sequences that were at least 99.3 % identical to the Latvian 337 B.1.1.523 sequence EPI ISL 1590462 and number of matched residues makes up equal or more 338 339 than 95% of the reference sequence. The reference sequence was chosen as it was closest to the 340 one of the first this lineage sequences sequenced at Lithuania but with smaller gaps regions. The 341 alignment against the GISAID sequences was conducted using minimap2 2.20-r1061 [https://github.com/lh3/minimap2/]. The limits were chosen arbitrary after several tries looking for 342 343 cut-offs resulting in a set of sequences that includes majority of sequences from the lineage and some more diverged ones that are classified as belonging to other lineages by pangolin. The 344 sequences with quality control overall status being bad or having more than 1000 bps missing (as 345 indicated by Nextclade analysis) were discarded. The maximum likelihood tree was calculated using 346 347 a modified version of Nextstrain workflow [https://github.com/nextstrain/zika (cloned at 2021 05 01)]. The tree was build using IQ-TREE with 2.1.2 General time reversible model with unequal rates 348 and unequal base frequencies were used²¹ allowing for a proportion of invariable sites together with 349 discrete Gamma model²². Ultrafast bootstrap with 1000 replicates was used. The maximum 350 351 likelihood emergence time and origin of country for inner nodes were calculated by treetime 0.8.1 352 [https://libraries.io/pypi/phylo-treetime] as described by the aforementioned workflow. The set of 353 sequences collected as described above were clustered into transmission clusters using Phydelity 354 v2.0 [https://github.com/alvinxhan/Phydelity]²³.

355 S protein phylogeny analysis

The S protein-based phylogeny was based on the S protein sequences extracted from GISAID 356 by the Nextclade and aligned to the reference COVID-19 sequence. Sequences shorter than 1175 357 residues or having more than one stop codon or having any number of undetermined residues were 358 discarded. Sequences were further clustered into identical sequence clusters using CD-HIT v4.8.1 359 (command line option "-c 1.0". The sequences representing all high-quality S protein variants were 360 361 used for maximum likelihood tree calculation VeryFastTree 3.0.1 by [DOI:

10.1093/bioinformatics/btaa582] using LG substitution model. The sequence alignment that was 362 used to construct the tree was composed from the alignment produced by the Nextclade leaving 363 only the representative sequences of the clusters and CAT approximation with 20 rate categories. 364 Ig -gamma". Also command line flags that should increase calculations accuracy were added: "-spr 365 4 -mlacc 2 -slownni -double-precision ". The tree was re-rooted using sequence matching to the 366 EPI_ISL_402124 as an out-group using gotree 0.4.1 [https://anaconda.org/bioconda/gotree/files]. 367

Ancestral states for inner nodes for 156, 157, 158, 484, 494 positions of protein S were 368 inferred using command line version GRASP-suite²⁵. This tool for inference was chosen due to high 369 370 speed and, most importantly, ability to handle insertion and deletions. The same substitution model 371 that used for phylogeny tree was also used in this case (LG). The resulting tree was analysed detected potential changes in the haplotypes was done using a custom script written in julia 1.6 372 exploiting capabilities of the NewickTree library [https://github.com/arzwa/NewickTree.jl]. The tree 373 was trimmed to keep only those inner nodes that leads to leaves containing the triple deletion at 374 146-148 positions and either E484K, S494P and visualised using ggtree 3.0.426. Additionally, a set 375 of full genome sequences has been composed matching the leaves of the aforementioned trimmed 376 tree and corresponding maximum likelihood tree was calculated using the aforementioned 377 Nextstrain workflow. 378

379

380 Analysis of potential recombination events

381 The focused set of sequences was composed based on the sequences used for the S protein 382 phylogeny. The set further narrowed to the sequences containing either the triple deletion at 156-383 158 positions or a combination of E484K and S494P. The detection of recombination events at DNA level was done using PoSeiDon workflow27 [https://github.com/hoelzer/poseidon cloned at 2021] 384 09 28] that runs GARD program to identify recombination events²⁸. The genomic regions matching 385 to the S protein were used. As the workflow is limited to 100 sequences the initial set 110 was 386 387 clustered to 100 sequences using cd-hit-est program from the CD-HIT package 4.8.1 using identity 388 cut-off equal to 0.9998658. The detection of recombination effects at protein level was done using detREC tool29 [https://github.com/gianfeng2/detREC program (cloned at 2021-10-01)]. The 389 recombination detection was conducted using either genomic or protein sequence corresponding 390 391 to the S protein. The full genome scale recombination analysis was done for the same set of 110 sequences using 3SEQ program^{28.} 392

393 Antibody escape effect estimation

394 Two S protein sequence variants 156_158del (B.1.1.523 lineage) and 156_157del & R158G (delta variant) were modelled using Rosetta package (2021.16.61629 bundle)³¹. The escape effect 395 was evaluated based on the structure PDBID: 7LQV containing a NTD domain targeting binding 396 antibody. The S protein structure from the "A" chain with matching antibody was used. The S protein 397 398 residues from 283 position to the "N" terminus were discarded. The PDB structures for fragment 399 library generation and other required databases for Rosetta were downloaded at 2021 06 02. The S structure models were created using a comparative modelling approach. The antibody structure 400 401 including side chains was kept constrained using CoordinateConstraintGenerator. For the relaxation

part "InterfaceRelax2019" script was used. The lowest energy model was chosen from 200 models. 402 The antibody docking was done using the SnugDock program³² from the Rosetta package, using 403 these non-default flags: "-auto generate h3 kink constraint - h3 loop csts hr -h3 filter false -404 docking centroid inner cycles=20 -docking centroid outer cycles=2 ". for each model docking 405 was run 300 times and the top 175 structures based on the interface score were chosen for analysis. 406 Antibody escape effect of E484K and S494P mutations were analysed using FoldX 5 407 program³³. The antibody structures targeting the receptor binding domain of the protein S were 408 downloaded from the CoV3D³⁴ as available at 2021 05 01. The mutations into the complexes were 409 introduced by consecutively applying FoldX command: RepairPDB, BuildModel, AnalyseComplex. 410 The residues at 484 and 494 positions were modelled either to the residue matching the mutation 411 412 or to the wild type residue and changes in free energy upon binding were evaluated.

413

414 **COMPETING INTERESTS**

The authors declare no competing interests. 415

medRxiv preprint doi: https://doi.org/10.1101/2021.11.21.21266655; this version posted November 24, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity.

perpetuity. It is made available under a CC-BY-NC 4.0 International license .

416 **REFERENCES**

417 World Health Organisation. WHO Coronavirus (COVID-19) Dashboard. WHO Coronavirus 1. 418 (COVID-19) Dashboard (2021). 419 F, W. et al. A new coronavirus associated with human respiratory disease in China. Nature 2. 420 **579**, 265–269 (2020). 421 European Centre for Disease Prevention and Control. How ECDC collects and processes 3. 422 COVID-19 data. https://www.ecdc.europa.eu/en/covid-19/data-collection (2021). 423 4. S, E. & G, B.-M. Data, disease and diplomacy: GISAID's innovative contribution to global 424 health. Global challenges (Hoboken, NJ) 1, 33–46 (2017). 425 Organisation. 5. World Health Tracking SARS-CoV-2 variants. 426 https://www.who.int/en/activities/tracking-SARS-CoV-2-variants/ (2021). 427 6. European Centre for Disease Prevention and Control. SARS-CoV-2 variants of concern as of 428 30 September 2021. (2021). 429 7. Planas, D. et al. Reduced sensitivity of SARS-CoV-2 variant Delta to antibody neutralization. 430 Nature 2021 596:7871 596, 276-280 (2021). 431 Collier, D. A. et al. SARS-CoV-2 B.1.1.7 escape from mRNA vaccine-elicited neutralizing 8. 432 antibodies. medRxiv 2021.01.19.21249840 (2021) doi:10.1101/2021.01.19.21249840. 433 9. AJ, G. et al. Comprehensive mapping of mutations in the SARS-CoV-2 receptor-binding 434 domain that affect recognition by polyclonal human plasma antibodies. Cell host & microbe 29, 463-435 476.e6 (2021). 436 Lasek-Nesselquist, E., Pata, J., Schneider, E. & George, K. St. A tale of three SARS-CoV-2 10. 437 variants with independently acquired P681H mutations in New York State. medRxiv 438 2021.03.10.21253285 (2021) doi:10.1101/2021.03.10.21253285. 439 11. Lewis, S., Stranges, B. P. & Adolf-Bryfogle, J. Interface Analyzer Documentation. 440 https://www.rosettacommons.org/docs/latest/application documentation/analysis/interface-441 analyzer (2016). 442 12. Koenig, P. A. et al. Structure-guided multivalent nanobodies block SARS-CoV-2 infection and 443 suppress mutational escape. Science **371**, (2021). 444 13. Latif, A. A. *et al.* B.1.1.523 Lineage Report. https://outbreak.info/situation-445 reports?pango=B.1.1.523 (2021). 446 14. McCallum, M. et al. Molecular basis of immune evasion by the delta and kappa SARS-CoV-2 447 variants. bioRxiv 2021.08.11.455956 (2021) doi:10.1101/2021.08.11.455956. 448 15. Wise, J. Covid-19: The E484K mutation and the risks it poses. BMJ 372, n359 (2021). 449 16. Z, L. et al. Identification of SARS-CoV-2 spike mutations that attenuate monoclonal and serum 450 antibody neutralization. Cell host & microbe 29, 477-488.e4 (2021). 451 17. AJ, G. et al. Complete Mapping of Mutations to the SARS-CoV-2 Spike Receptor-Binding 452 Domain that Escape Antibody Recognition. *Cell host & microbe* **29**, 44-57.e9 (2021). S, K. & A, K. Under-reporting of COVID-19 cases in Turkey. The International journal of health 453 18. 454 planning and management **35**, 1009–1013 (2020). 455 Y, S. & J, M. GISAID: Global initiative on sharing all influenza data - from vision to reality. Euro 19. 456 surveillance : bulletin Europeen sur les maladies transmissibles = European communicable disease 457 bulletin 22, (2017). 458 20. O'Toole, Á. et al. Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and 459 B.1.351/501Y-V2. Wellcome Open Research 2021 6:121 6, 121 (2021). 460 21. Hadfield, J. et al. Nextstrain: real-time tracking of pathogen evolution. Bioinformatics 34, 461 4121-4123 (2018). 462 Tavaré, S. Some probabilistic and statistical problems in the analysis of DNA sequences. (The 22. 463 American Mathematical Society, 1986). 464 23. X, G., YX, F. & WH, L. Maximum likelihood estimation of the heterogeneity of substitution 465 rate among nucleotide sites. *Molecular biology and evolution* **12**, 546–557 (1995).

466	24.	AX, H., E, P., S, MS. & CA, R. Inferring putative transmission clusters with Phydelity. Virus				
467	evolut	evolution 5 , (2019).				
468	25.	SQ, L. & O, G. An improved general amino acid replacement matrix. Molecular biology and				
469	evolut	olution 25 , 1307–1320 (2008).				
470	26.	Foley, G. et al. Identifying and engineering ancient variants of enzymes using Graphical				
471	Repre	presentation of Ancestral Sequence Predictions (GRASP). <i>bioRxiv</i> 2019.12.30.891457 (2020)				
472	doi:10	oi:10.1101/2019.12.30.891457.				
473	27.	Yu, G. Using ggtree to Visualize Data on Tree-Like Structures. Current Protocols in				
474	Bioinf	<i>ioinformatics</i> 69 , e96 (2020).				
475	28.	M, H. & M, M. PoSeiDon: a Nextflow pipeline for the detection of evolutionary recombination				
476	event	nts and positive selection. Bioinformatics (Oxford, England) 37, 1018–1020 (2021).				
477	29.	SL, K. P., D, P., MB, G., CH, W. & SD, F. GARD: a genetic algorithm for recombination detection.				
478	Bioinf	informatics (Oxford, England) 22 , 3096–3098 (2006).				
479	30.	Feng, Q. et al. A scalable method for identifying recombinants from unaligned sequences.				
480	bioRx	R <i>xiv</i> 2020.11.18.389262 (2020) doi:10.1101/2020.11.18.389262.				
481	31.	HM, L., O, R. & MF, B. Improved Algorithmic Complexity for the 3SEQ Recombination				
482	Detec	Detection Algorithm. <i>Molecular biology and evolution</i> 35 , 247–251 (2018).				
483	32.	Leman, J. K. et al. Macromolecular modeling and design in Rosetta: recent methods and				
484	frame	meworks. <i>Nature Methods 2020 17:7</i> 17 , 665–680 (2020).				
485	33.	Sircar, A. & Gray, J. J. SnugDock: Paratope Structural Optimization during Antibody-Antigen				
486	Docki	ng Compensates for Errors in Antibody Homology Models. PLOS Computational Biology 6,				
487	e1000	21000644 (2010).				
488	34.	Schymkowitz, J. et al. The FoldX web server: an online force field. Nucleic acids research 33,				
489	(2005).				
490	35.	Gowthaman, R. et al. CoV3D: a database of high resolution coronavirus protein structures.				
491	Nucle	Nucleic acids research 49 , D282–D287 (2021).				
492						