

1 **Adaptive advantage of deletion repair in the N-terminal domain of the**
2 **SARS-CoV-2 spike protein in variants of concern**

3 Miguel Álvarez-Herrera^a, Paula Ruiz-Rodríguez^a, Beatriz Navarro-
4 Domínguez^{ab}, Joao Zulaica^a, Brayan Grau^a, María Alma Bracho^{ac}, Manuel
5 Guerreiro^{de}, Cristóbal Aguilar-Gallardo^e, Fernando González-Candelas^{ac},
6 Iñaki Comas^{cf}, Ron Geller^a, Mireia Coscollá^{a*}

7 ^a*Institute for Integrative Systems Biology (I²SysBio, University of Valencia-CSIC),*
8 *FISABIO Joint Research Unit “Infection and Public Health”, Paterna, Spain;*
9 ^b*University of Granada, Department of Genetics, Granada, Spain;* ^c*CIBER in*
10 *Epidemiology and Public Health (CIBERESP), Madrid, Spain;* ^d*Haematology*
11 *Department, La Fe University and Polytechnic Hospital, Valencia, Spain;* ^e*La Fe*
12 *Health Research Institute (IIS-La Fe), Valencia, Spain;* ^f*Institute of Biomedicine of*
13 *Valencia (IBV-CSIC), Valencia, Spain*

14 * Email: mireia.coscolla@csic.es. Address: Institute for Integrative Systems Biology, C/
15 Catedrático Agustín Escardino 9, Building 4, 46980 Paterna, Spain

16 **Keywords**

17 SARS-CoV-2, COVID-19, spike, transmission, deletion, fusogenicity, antibody
18 neutralization, thermal stability

19 **Abstract**

20 Mutations within the N-terminal domain (NTD) of the spike (S) protein play a pivotal
21 role in the emergence of successful SARS-CoV-2 viral lineages. This study investigates
22 the influence on viral success of novel combinations of NTD lineage-defining mutations
23 found in the Alpha, Delta, and Omicron variants. We performed comparative genomics
24 of more than 10 million public SARS-CoV-2 samples to decipher the transmission
25 success of different combinations of NTD markers. Additionally, we characterized the
26 viral phenotype of such markers in a surrogate *in vitro* system. Alpha viruses bearing
27 repaired deletions S:ΔH69/V70 and S:ΔY144 in Alpha background were associated
28 with increased transmission relative to other combinations of NTD markers. After the
29 emergence of the Omicron BA.1 lineage, Alpha viruses harbouring both repaired
30 deletions still showed increased transmission compared to their BA.1 analogues.
31 Moreover, repaired deletions were more frequently observed among older individuals
32 infected with Alpha, but not with BA.1. *In vitro* biological characterization of Omicron
33 BA.1 spike deletion repair patterns also revealed substantial differences with Alpha. In
34 BA.1, S:ΔV143/Y145 repair enhanced fusogenicity and susceptibility to neutralization
35 by vaccinated individuals' sera. In contrast, the S:ΔH69/V70 repair did not significantly
36 alter these traits but reduced viral infectivity. Simultaneous repair of both deletions led
37 to lower fusogenicity. These findings highlight the intricate genotype-phenotype
38 landscape of the spike NTD in SARS-CoV-2, which impacts viral biology, transmission
39 efficiency, and susceptibility to neutralization. Overall, this study advances our
40 understanding of SARS-CoV-2 evolution, carrying implications for public health and
41 future research.

42

43 **Introduction**

44 As we navigate the constantly shifting landscape of the COVID-19 pandemic, the
45 remarkable potential of SARS-CoV-2 for genetic adaptation has taken centre stage.
46 Global turnover of SARS-CoV-2 lineages happened several times, with three variants of
47 concern (VOCs) displacing the previous predominant lineages just in 2021: first, Alpha
48 (B.1.1.7 and Q sublineages), then Delta (B.1.617.2 and AY sublineages), and lastly,
49 Omicron (B.1.1.529 and BA sublineages, among others) [1]. This successive
50 replacement of lineages across the pandemic suggests that the newer lineages had a
51 higher adaptive value than the previous ones [2,3], and thus, these three lineages are
52 supposed to carry mutations associated with higher transmissibility and/or immune
53 evasion.

54 The spike (S) protein—a type I transmembrane N-linked glycosylated protein of
55 150–200 kDa—is a hotspot for mutations with high adaptive value. Spike proteins are
56 located on the surface of SARS-CoV-2, and their main role is to mediate viral cell entry
57 [4]. The spike protein forms a homotrimer, which is cleaved post-transcriptionally into
58 two subunits: S1 and S2. The S1 consists of the amino or N-terminal domain (NTD) and
59 the receptor-binding domain (RBD), and it is responsible for binding to the host cell-
60 surface receptor, ACE2. The S2 subunit includes the trimeric core of the protein and is
61 responsible for membrane fusion [5,6]. Therefore, some amino acid changes in the S
62 protein may confer an advantage for transmission, considering its role in mediating viral
63 cell entry. Additionally, most antibodies target sites either in the NTD or RBD, and
64 therefore, mutations in these regions may enhance immune escape.

65 In Alpha lineages, deletions in the NTD (S:ΔH69/V70 and S:ΔY144) are
66 associated with antibody escape [7,8] and increased infectivity [9]. Acquisition of
67 deletions in the NTD of the spike glycoprotein during long-term infections of

68 immunocompromised patients has been reported and identified as an evolutionary
69 pattern defined by recurrent deletions that alter defined antibody epitopes [10,11].
70 Additionally, deletions may play a decisive role in SARS-CoV-2 adaptive evolution,
71 particularly on deletion-tolerant genome regions such as the S gene, as they can hardly
72 be corrected by the proofreading activity of its RNA-dependent RNA polymerase
73 (RdRP) [12–14]. Indeed, the NTD has been extensively impacted by deletions, which
74 have arisen multiple times in different variants, including Alpha, Delta, and Omicron.
75 For example, different deletions are observed in Delta (S:ΔE156/F157-R158G), Alpha
76 (S:ΔY144) and BA.1 (S:ΔV143/Y145) variants, all mapping to the same surface
77 indicating a convergent function [15,16]. All these mutations are within the NTD site
78 targeted by most anti-NTD neutralizing antibodies [8]. Despite not being found in BA.2,
79 its descendant lineage BJ.1 seems to have independently acquired S:ΔY144 (see cov-
80 lineages/pango-designation issues #915 and #922 on GitHub). Then, it was passed on to
81 XBB viruses through recombination with a BA.2-descended lineage, where it is
82 associated with increased immune escape regardless of the vaccination status or
83 infection history [17–19]. Subsequently, in January 2023, the newly emerging XBB.1.5
84 lineage was shown to display an increased receptor-binding affinity and infectivity with
85 respect to its parental lineage [20,21].

86 During the takeover of the Alpha variant by the Delta variant, we became
87 intrigued by the lack of overlap in NTD mutations between them. If mutations in the
88 NTD increased immune escape without compromising binding to ACE2, one might
89 expect that mutations in NTD have a cumulative (“the more the better”) effect.
90 However, that would not be the case if epistatic interactions between sites prevent the
91 fixation of particular mutations in different genomic and genetic backgrounds. Our
92 primary objective in this study is to investigate the effect of NTD deletion repair in

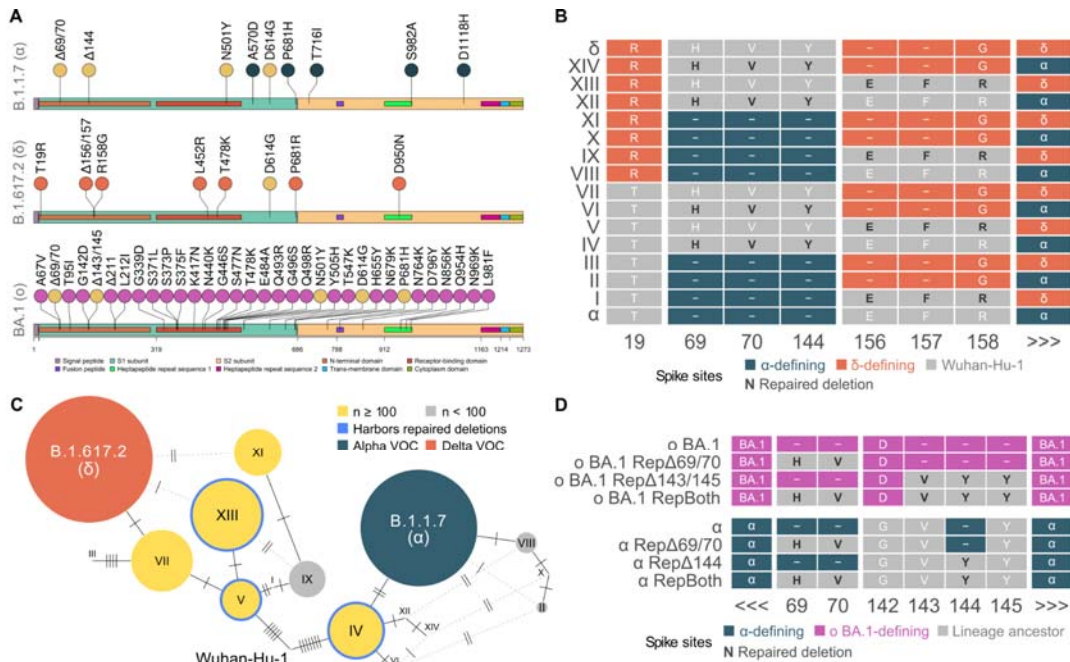
93 variants of concern. We focus on examining the differences in success between variants
94 of concern as a function of different mutational patterns in the NTD, including recurrent
95 deletion repair events occurring in distinct independent lineages. Our findings suggest a
96 non-linear effect of specific deletion repairs on viral phenotype, highlighting the
97 importance of examining the genomic context of SARS-CoV-2 mutations.

98 **Materials and Methods**

99 *Data retrieval and preprocessing*

100 Screening Alpha and Delta combinations of NTD markers included 7,133,237
101 available SARS-CoV-2 sequences fetched from the GISAID [22] EpiCoV database on
102 January 16, 2022. The GISAID EPI_SET is available at DOI: 10.55876/gis8.230802fh.
103 We ran Nextalign CLI v1.9.0 [23] with default scoring settings to obtain aligned
104 genomes and peptides corresponding to each gene. Genomes from non-human hosts
105 were discarded. Then, we filtered out genomes containing more than 5% ambiguous
106 bases and 1,000 gaps, or at least one indeterminate position among the lineage-defining
107 sites of the spike protein. 6,589,393 genomes passed these filters. An additional filtering
108 step implemented in Nextclade CLI v1.7.0 [24] was used to flag and remove false
109 positive mixtures possibly caused by sample contamination or coinfections, considering
110 that we were exploring combinatorial variation during the time of co-existence of highly
111 successful variants. The same dataset was used to screen for patterns of deletion repair
112 in the Alpha variant, omitting the Nextclade filtering step as it was not deemed critical
113 when not searching for variant mixtures. We performed an analogous search of
114 deletion-missing sequences assigned to lineage BA.1 out of 11,334,504 samples
115 (10,353,158 after quality-filtering), fetched on June 15, 2022. The GISAID EPI_SET is
116 available at DOI: 10.55876/gis8.230801ex.

117 The classification of Alpha-Delta combinatorial haplotypes was based on the
 118 sequence correspondence of protein segments to the canonical variants. We defined
 119 three blocks within the N-terminus of the S protein that included all NTD lineage-
 120 defining mutations of the Alpha (S:ΔH69/V70 and S:ΔY144) and Delta (S:T19R and
 121 S:ΔE156/F157 + S:R158G) variants, each encompassing contiguous lineage-defining
 122 sites (S:16, S:69/70 + S:144 and S:156/158; see Figure 1B). We considered the
 123 remaining portion of the protein as background. Alpha and Delta backgrounds include
 124 the mutations in Figure 1A. Site S:142 was not considered in this analysis due to a
 125 known technical sequencing artifact in variant Delta [25]. See Supplemental Note 1 for
 126 further details about haplotype naming. This survey included sixteen haplotypes,
 127 including the canonical Alpha and Delta spike protein sequences and fourteen
 128 combinatorial haplotypes (I-XIV).



129 Figure 1. Combinatorial standpoint of the study. (A) Defining mutations of the B.1.1.7 (Alpha), B.1.617.2 (Delta) and
 130 BA.1 (Omicron) SARS-CoV-2 lineages on the spike protein. (B) Spike protein lineage-defining sites of the fourteen
 131 Alpha-Delta combinatorial haplotypes (I-XIV) and two variants of concern. (C) Haplotype network of the most
 132 abundant combinatorial haplotypes (see Table 1) and the B.1.1.7 and B.1.617.2 lineages, based on the haplotype
 133 group-level consensus of the S gene. Solid edges depict the minimum spanning tree of the graph. Each stroke on the
 134 edges stands for one nucleotide site that changes between the connected nodes. Node size represents the number of
 135 observations in log-scale. (D) Spike protein lineage-defining sites of samples that repair NTD deletions in common
 136 on Alpha and Omicron BA.1 background. Sites in grey are unaltered with respect to each lineage ancestral sequence.
 137

138 Table 1. Description of five Alpha-Delta combinatorial haplotypes with at least 100 observations. Values marked
139 with † were obtained with the dataset reduction strategy with CD-HIT.

Haplotype	Description	Observations	Transmission clusters	Minimum independent emergences
IV	Alpha + repaired S:ΔH69/V70 & S:ΔY144	736	72 78†	576 598†
V	Delta + S:R19T reversion & repaired S:ΔE156/F157-R158G	117	6	104
VII	Delta + S:R19T reversion	1,666	146	1,121
XI	Delta + S:ΔH69/V70 & S:ΔY144	304	19	267
XIII	Delta + repaired S:ΔE156/F157-R158G	4,781	431	2,880

140 Analogously, the survey of NTD deletion repair haplotypes in Alpha and
141 Omicron BA.1 backgrounds tracked the presence or absence of specific deletions in
142 sequences assigned to lineages B.1.1.7 (Alpha; deletions S:H69/V70 and S:ΔY144) and
143 BA.1 (Omicron; deletions S:H69/V70 and S:ΔV143/Y145), or any of their sublineages.
144 For more in-depth considerations about S:ΔV143/Y145 repair in lineage BA.1, see
145 Supplemental Note 2.

146 After data cleaning and homogenization of the retrieved sample metadata, we
147 were able to obtain the host age for approximately 50% of the analyzed samples of both
148 datasets. We used R v4.1.2 [26] along with *tidyverse* v1.3.1 [27] to conduct, manage
149 and visualize these analyses. Lollipop plots were generated using *trackViewer* v1.30.0
150 [28]. Haplotype networks with abundance weights of the S gene consensus sequences of
151 the combinatorial haplotypes were built using the randomized minimum spanning tree
152 (RMST) algorithm [29] implemented in *pegas* v1.1 [30]. We considered insertion or
153 deletion blocks as a single mutational event. Synthetic sequences were built to account
154 for unsampled haplotypes.

155 *Number of emergences and transmissions*

156 To infer the success of each combinatorial haplotype, we assessed the number of

157 minimum emergence events and whether these originated subsequent transmission
158 events. To overcome the challenging computational demands of working with global-
159 scale sequence datasets and, we devised a double approach, based on two distinct
160 optimality criteria. First, we placed genomes matching combinatorial haplotypes (I-
161 XIV) on a comprehensive phylogeny of public sequences up to the date of each survey
162 [31] under a maximum parsimony criterion, using UShER v0.5.3 [32]. Second, we used
163 an additional phylogenetic inference method that considered an evolutionary process
164 model while keeping the phylogenetic context of our datasets, to enhance the reliability
165 of our results. S gene clustering based on short word filtering was performed with CD-
166 HIT-EST, bundled in the CD-HIT suite v4.8.1 [33], with a word size of 8 and a
167 minimum sequence identity threshold equivalent to 5 amino acid changes in the spike
168 protein. For any given haplotype, sequences within clusters that included at least one of
169 the haplotype-assigned samples were selected as members of the reduced dataset for the
170 corresponding haplotype. Then, we inferred a whole-genome phylogeny with maximum
171 likelihood under a GTR model, with the Wuhan-Hu-1 sequence (NCBI RefSeq
172 accession no. NC_045512.2) as the outgroup, using IQ-TREE v2.1.2 COVID-edition
173 [34]. This approach was applied to haplotype IV (see Table 1), resulting in a reduced
174 phylogeny with 30,314 tips.

175 Quantification of transmission of Alpha and Delta mixtures was conducted
176 through an exhaustive breadth-first search, selecting clusters with at least 2 members
177 and at least 90% target samples, utilizing the implementation by Ruiz-Rodriguez et al.
178 [35]. The composition requirement was reduced from 100% to account for potential
179 sequencing errors and ambiguous sample placements on the phylogenetic tree. Due to
180 the increased dataset size and complexity, we quantified transmission of deletion repair
181 Alpha and BA.1 viruses using a reimplementaion of the algorithm that leverages

182 polytomies in a global-scale phylogeny for parallelization in an HPC environment,
183 based on *phylobase* v0.8.10 (<https://github.com/fmichonneau/phylobase>). The complete
184 pipeline is available as a Snakemake workflow at [https://github.com/PathoGenOmics-](https://github.com/PathoGenOmics-Lab/transcluster)
185 Lab/transcluster (v1.0.0). The minimum number of independent emergences for each
186 combinatorial haplotype was derived from the phylogenies by adding up the number of
187 transmission clusters to the number of emergences that were not transmitted. We
188 studied the location and time span of the transmission clusters as well (see
189 Supplemental Figures 1 and 2). Host age comparisons between transmitted and non-
190 transmitted samples included clusters with more than 2 members. We also conducted a
191 detailed analysis of the age distribution of transmission clusters associated with
192 combinatorial haplotype IV to exclude the influence of potential confounding factors
193 (see Supplemental Figure 3).

194 To adequately compare the relative transmission success of the most prevalent
195 haplotypes, we developed a method to estimate their transmission fitness for each
196 cluster. This estimation involved calculating the ratio between the cluster's size and the
197 number of sequences in GISAID that were collected between the first and last cases of
198 the transmission cluster, with 7 days as padding at both sides of these windows to
199 mitigate the effect of missing data and short cluster time spans. For clusters exhibiting
200 cross-border transmission (involving samples from different geographic locations), we
201 further divided the cluster time window into country-specific sub-windows. The
202 denominator of the transmission fitness estimate was then calculated as the sum of the
203 number of sequences in GISAID for each country-specific time window. This approach
204 allowed us to account for variations in sampling efforts and prevalence across different
205 time periods and geographic regions. Conducting the analyses without the 1-week
206 padding or without differentiating country-specific sub-windows yielded significantly

207 different estimates, but the overall differences between haplotypes did not change (see
208 Supplemental Figure 4). Differences in the distribution of the estimated transmission
209 fitness between haplotypes were evaluated using Wilcoxon rank-sum tests. Statistical
210 analyses were performed and visualized using R v4.1.2 [26] along with *tidyverse* v1.3.1
211 [27] and *ggpubr* v0.4.0 [36].

212 ***Biological characterization of BA.1 deletion repair***

213 Combinations of deletion repairs of S:ΔH69/V70 and S:ΔV143/ΔY145 were introduced
214 into a pCG1 plasmid encoding a codon-optimized BA.1 spike protein [37] by site-
215 directed mutagenesis. All the constructs were verified by Sanger sequencing.
216 Pseudotyped vesicular stomatitis virus (VSV) encoding a GFP reporter gene and
217 carrying the different spike proteins was produced as previously reported [38]. To
218 assess the effects on virus production, pseudotyped VSV carrying each construct were
219 produced independently three times. The resulting viruses were then titrated by
220 infecting Vero E6 cells (kindly provided by Dr. Luis Enjuanes; CNB, Spain) or Vero
221 E6-TMPRSS2 cells (JCRB Cell bank catalogue code: JCRB1819) for 16 hours,
222 followed by quantification of GFP-expressing infected cells using a live cell microscope
223 (Incucyte SX5; Sartorius) to obtain the number of focus forming units (FFU) per
224 millilitre. To assess thermal stability, 500 FFU of these pseudotyped viruses were
225 incubated for 15 min at a range of temperature in a thermal cycler (30.4, 31.4, 33.0,
226 35.2, 38.2, 44.8, 47.0, 48.6 and 49.6°C; Biometra T one Gradient, Analytik Jena) and
227 the surviving virus was used to infect VeroE6-TMPRSS2 cells for 16 h. The GFP signal
228 in each well was then determined using a live-cell microscope (Incucyte SX5,
229 Sartorius). The average GFP signal observed in mock-infected wells was subtracted
230 from all infected wells, followed by standardization of the GFP signal to the average
231 GFP signal from wells incubated at 30.4 °C. Finally, we fitted a three-parameter log-

232 logistic function to the data using the *drc* v3.0-1 R package (*LL.3* function) and
233 calculated the temperature resulting in 50% reduction in virus infection (*ED* function).
234 To assess the effects on neutralization by polyclonal sera, we used six sera from
235 convalescent patients from the first COVID-19 wave in Spain and six sera from
236 individuals that had been administered two doses of the BioNTech-Pfizer Comirnaty
237 COVID-19 vaccine. The neutralization capacity of the sera was obtained as previously
238 described on VeroE6-TMPRSS2 cells [37]. We used a previously described flow
239 cytometry assay based on the use of polyclonal sera to examine surface expression [39].
240 Briefly, HEK293T cells were transfected with the different S mutants using the calcium
241 chloride method. 24 h later, cells were detached using PBS with 1 mM EDTA, washed,
242 and incubated on ice with different polyclonal sera (three from convalescent patients
243 from the 1st COVID-19 wave in Spain and one from individuals that had been
244 administered two doses of the BioNTech-Pfizer Comirnaty COVID-19 vaccine) at a
245 1:300 dilution in PBS containing 0.5 % BSA and 2 mM EDTA for 30 min. Next, cells
246 were washed three times with PBS, stained with anti-IgG Alexa Fluor 647 (Thermo
247 Fisher Scientific) at a 1:400 dilution, and analyzed by flow cytometry similarly treated
248 un-transfected controls to set the threshold for positive cells. For cell-cell fusion assays
249 we used a split Venus fluorescent protein system [40]. Briefly, HEK293T cells were
250 grown overnight in 24 well plates (1.5×10^5 cells/well) using DMEM supplemented
251 with 10 % FBS. After 24 hours, cells were transfected using Lipofectamine 2000
252 (Invitrogen) with 0.5 μ g of either a 1:1 mixture of the S plasmids and a Jun-Nt Venus
253 fragment (Addgene 22012) plasmid or a mixture of hACE2 plasmid (kindly provided by
254 Dr. Markus Hoffman, German Primate Center, Goettingen/Germany) [41] and the Fos-
255 Ct Venus fragment (Addgene 22013). After 24 h, cells were counted, and the S-
256 transfected cells were mixed at a 1:1 ratio with ACE2-transfected cells and seeded in 96

257 well plates (3 x 10⁴ cells/well) in 100 μ L of media. Cells transfected with the Wuhan-
258 Hu-1 served as a positive control, while cells transfected with hACE2 and Jun-Nt Venus
259 were used as negative control. We obtained the GFP Integrated Intensity
260 (GCU· μ m²/image) in each condition using a live-cell imaging platform (Incucyte SX5,
261 Sartorius) at 24 hours post mixing and standardized to the signal obtained from the
262 positive control (Wuhan-Hu-1 spike protein). All experiments were performed at least
263 three times in triplicates.

264 Statistical analyses were performed and visualized using R v4.1.2 [26] along
265 with *tidyverse* v1.3.1 [27] and *ggpubr* v0.4.0 [36] to facilitate the analysis and enhance
266 the visualization of the results. Comparisons were conducted utilizing t-tests (unpaired
267 for all assays and paired for neutralization and surface expression, as we used the same
268 sets of polyclonal sera in each experiment) after verifying that the data met the
269 assumptions of normality using a Shapiro-Wilk test. To determine fold-change values,
270 the ratio of group average values was calculated.

271 **Results**

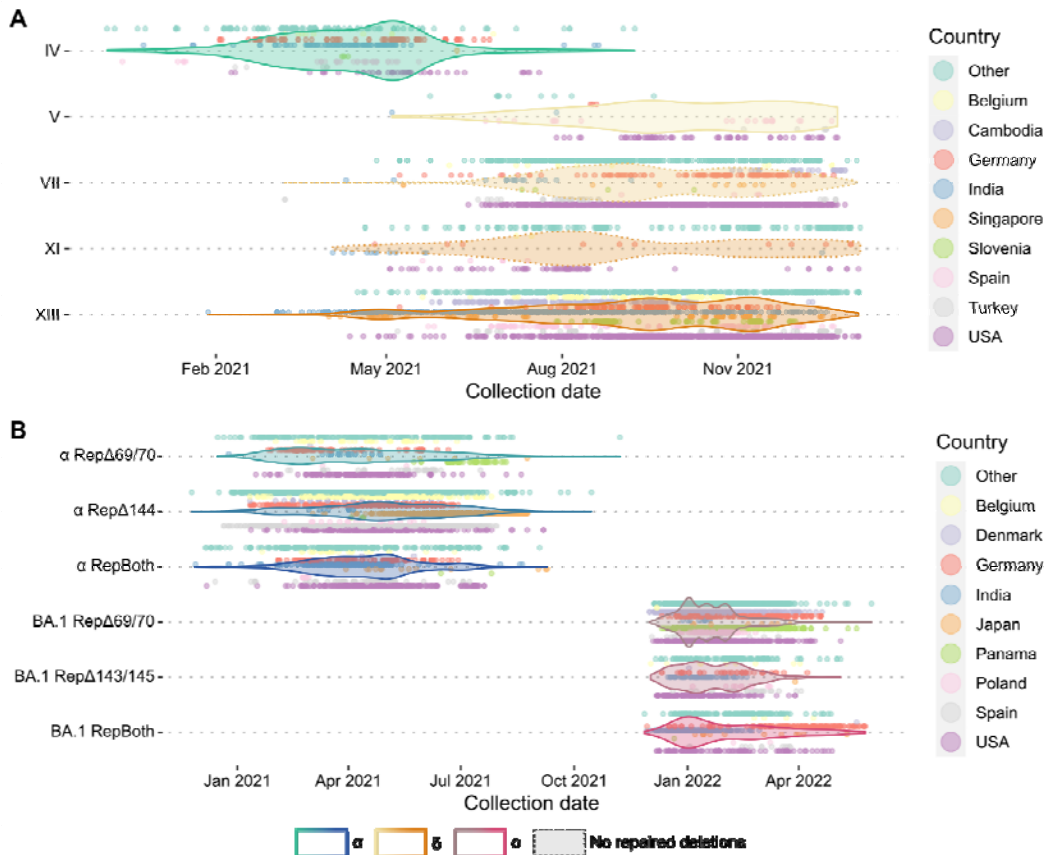
272 *Favoured and forbidden mixtures of NTD marker combinations in the Alpha* 273 *and Delta variants*

274 In the first part of this work, we focused on examining the differences in success
275 between the Alpha and Delta variants in the NTD of the spike protein. This domain is
276 characterized by recurrent deletions occurring in distinct independent lineages
277 [8,14,15,17,42,43]. We identified 7,706 samples that carry a mixture of Alpha and
278 Delta-defining mutations in this region, out of the 7.13 million sequences submitted to
279 the GISAID's EpiCoV repository as of January 2022 (Supplemental Table 1). By
280 comparing these two variants, we aimed to gain insights into potential variations in their

281 adaptive characteristics and overall performance in this specific region. In total, 14
282 combinatorial haplotypes were surveyed (termed I-XIV, Figure 1B). Differences in the
283 number of observations were apparent, with the noticeable absence ($n = 0$) of
284 haplotypes I, X, XII and XIV, and low prevalence ($n < 100$) of haplotypes II, III, VI,
285 VIII, and IX. These haplotypes were omitted from further analysis, as they were not
286 considered to be epidemiologically relevant. On the contrary, haplotypes IV, V, VII, XI,
287 and XIII were continuously sampled throughout the time window of the variant
288 takeover (Figure 1C and Figure 2A), with a remarkably high prevalence compared to
289 the rest (Table 1). The haplotype network showed two well-separated haplotype groups
290 around the two main variants of concern connected by the reference genome resembling
291 a distance-based unrooted phylogeny with added alternative connections, i.e. mutational
292 jumps (Figure 1C). Interestingly, three out of five of these combinatorial haplotypes
293 (IV, V and XIII) bore repaired NTD deletions (Figure 1B).

294 Table 2. Description of six Alpha and Omicron BA.1 haplotypes with repaired NTD deletions in sites S:69/70 and
295 S:144.

Haplotype	Description	Observations	Transmission clusters	Minimum independent emergences
Alpha Rep Δ 69/70	Alpha + repaired S: Δ H69/V70	722	0	722
Alpha Rep Δ 144	Alpha + repaired S: Δ Y144	4,571	332	2,823
Alpha RepBoth	Alpha + repaired S: Δ H69/V70 & S: Δ Y144	1,297	119	982
BA.1 Rep Δ 69/70	Omicron BA.1 + repaired S: Δ H69/V70	5,166	198	3,969
BA.1 Rep Δ 143/145	Omicron BA.1 + repaired S: Δ V143/Y145	637	27	537
BA.1 RepBoth	Omicron BA.1 + repaired S: Δ H69/V70 & S: Δ V143/Y145	737	63	453

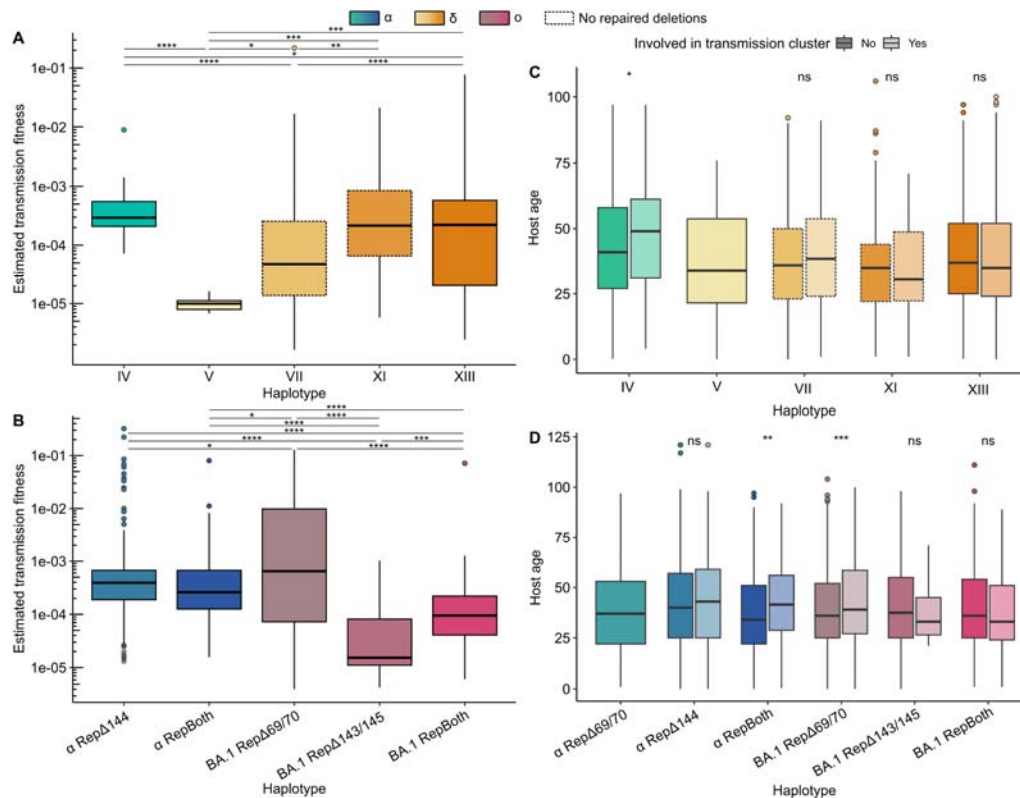


296
297
298
299
300
301

Figure 2. Sample collection timelines broken down by sampling location. The nine countries with the highest number of observations in each dataset are displayed for each panel. Each point represents a different sample assigned to the corresponding haplotype. (A) Observations of each Alpha and Delta NTD combinatorial haplotype with more than one hundred observations. (B) Observations of each deletion repair variant on Alpha and BA.1 (Omicron) background (see Table 2).

302 Nevertheless, estimates of prevalence can be distorted by several factors,
303 including the geographical location and temporal specificity associated with both
304 sequencing efforts and the distribution of lineages throughout the pandemic.
305 Transmission, on the other hand, is usually considered as a proxy of viral fitness,
306 because it is related to its basic reproductive number, reflecting the ability of the virus
307 to replicate, persist and spread within hosts and in the population [44–47]. Therefore,
308 we estimated the transmission of each haplotype by interrogating worldwide SARS-
309 CoV-2 sequence data. This enabled the group-wise quantification of the minimum
310 number of emergence events. We observed limited cross-border transmission, as 89 %
311 of all transmission clusters were contained in a single country. The average within-
312 cluster collection date window was 30 ± 44 days (Supplemental Figure 1). We found

313 vast differences in the number of clusters among haplotypes (see Table 1). In fact, we
 314 observed that this indicator, as well as the number of emergences, tended to increase
 315 linearly with the number of observations (both $R^2 = 1.0$, with $p = 2.2 \cdot 10^{-16}$ and $p =$
 316 $4.7 \cdot 10^{-16}$, respectively). To mitigate this effect and better evaluate differences in
 317 transmission of the most frequent haplotypes, we estimated their transmission fitness by
 318 adjusting their cluster sizes to account for variation in sampling effort and prevalence
 319 across different time periods and geographic regions. We found that the two haplotypes
 320 with the highest median estimated transmission fitness bore repaired deletions, both on
 321 Delta (haplotype XIII) and Alpha (haplotype IV) genomic backgrounds (Figure 3A).



322
 323 Figure 3. Differences in the estimated transmission fitness (A, B) and host age (C, D) between five Alpha and Delta
 324 combinatorial haplotype groups (A, C), and all six Alpha and Omicron BA.1 deletion repair haplotypes (B, D).
 325 Colour gradients indicate the variant representing each haplotype background. The differences between and within
 326 haplotypes were assessed using Wilcoxon rank-sum tests (ns: $p > 0.05$; *: $p \leq 0.05$; **: $p \leq 0.01$; ***: $p \leq 0.001$;
 327 ****: $p \leq 0.0001$). Only significant differences in panels A and B are displayed for clarity. Note that the haplotype
 328 corresponding to the Alpha variant with repaired S: Δ H69/V70 is missing from panel B due to its lack of
 329 transmission.

330 We sought to identify possible adaptive drivers among the clinical variables
 331 associated with these samples. We did not detect any significant differences related to

332 host sex. However, we found an association of host age with deletion repair in Alpha
333 background: viruses assigned to haplotype IV infected older hosts (average 43 ± 21
334 years old) compared to the rest of NTD combinations (average difference of 6 ± 1 years;
335 all $p < 0.005$, Wilcoxon rank-sum test), except V ($p = 0.085$, Wilcoxon rank-sum test).
336 These differences are presented in Supplemental Figure 3A. We then assessed whether
337 our results could be biased by the epidemiology and demography of the Alpha and
338 Delta variants, which led to distinct spread patterns among different population groups.
339 We confirmed that host age in haplotype IV was also higher when compared to samples
340 that were collected in the same time window but not assigned to haplotype IV, or any of
341 the remaining combinatorial haplotypes (both $p = 1.6 \cdot 10^{-6}$, Wilcoxon rank-sum test;
342 Supplemental Figure 3B). Thus, sampling bias is unlikely to be the primary driver of the
343 observed age differences between combinatorial haplotypes. Additionally, the
344 associated host age was higher in samples that were involved in transmission events
345 (Figure 3B). This points to the non-essentiality of S: Δ H69/V70 and S: Δ Y144 for the
346 infection of older individuals—who are often immunocompromised—in the Alpha
347 background. However, it could be argued that certain outbreaks could skew the age
348 distribution—for instance, if several large outbreaks occurred in elderly care facilities.
349 To control for this factor, we analysed the within-location cluster size distribution at a
350 regional level and found that hosts were generally older, but there were no dominant
351 transmission events (Supplemental Figure 3C and 3D). Based on these findings, we find
352 no evidence to suggest that age effects were driven by a few specific outbreaks.

353 To summarize, haplotypes featuring repaired deletions in the NTD exhibit an
354 increased number of observations and transmission capabilities among all combinatorial
355 possibilities, except for the established variants of concern. Furthermore, our findings
356 show a pronounced difference in distribution and fitness between combinatorial

357 variation of Alpha and Delta, emphasizing the significance of genetic context in the
358 evaluation of genotype-phenotype relationships.

359 ***Common patterns of NTD deletion repair confer different degrees of viral***
360 ***success on Alpha and Omicron BA.1 backgrounds***

361 The Omicron BA.1 lineage emerged after November 2021 and rapidly outcompeted the
362 Delta variant. The BA.1 lineage is defined by deletions S: Δ H69/V70 and
363 S: Δ V143/Y145, which map to Alpha-defining deletions (see Figure 1A) and are known
364 to also confer adaptive advantages in Alpha viruses [14,42]. Building upon our prior
365 findings about the simultaneous repair of these two deletions in the Alpha variant
366 (haplotype IV), we interrogated the possibility of similar repairs occurring in the
367 Omicron BA.1 genetic background by investigating different patterns of deletion repair
368 of epidemiological significance and the drivers behind their emergence. We performed
369 an analogous global survey followed by a phylogenetic estimation of viral success of
370 individual and dual deletion repairs in Alpha and Omicron BA.1 backgrounds. Due to
371 the increased number of BA.1-defining markers compared to previous hegemonic
372 variants of concern (see Figure 1A), we based the survey solely on the presence or
373 absence (i.e. presence of the ancestral state) of the specific deletions of interest (Figure
374 1D). We identified 13,130 samples that carried repaired NTD deletions in Alpha or
375 Omicron BA.1 backgrounds as of June 2022 (Supplemental Table 2).

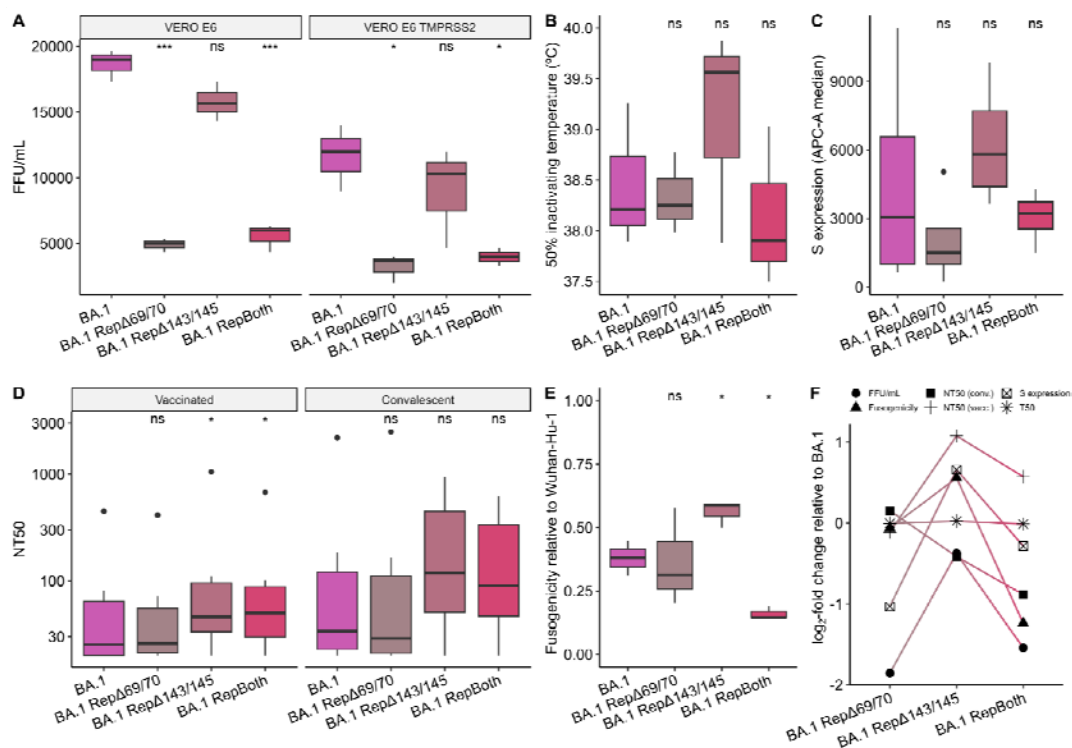
376 In terms of the number of observations, there was a clear disparity in repair
377 patterns between Alpha and Omicron BA.1 samples. The most frequently observed
378 group consisted of samples exhibiting the repaired S: Δ H69/V70 in BA.1 background,
379 followed by S: Δ Y144 repair in Alpha, while the remaining repair patterns were
380 significantly less frequently observed (Table 2). The overall number of observations of
381 Omicron BA.1 viruses with double repair was predictably lower than that of the

382 predominant parental lineage, as the sampling window was nearly two-thirds narrower
383 compared to that of other groups (Figure 2), but similar to that of combinatorial
384 haplotypes in our previous analysis. We measured an average within-cluster collection
385 date window of 25 ± 31 days, with 87 % of clusters showing within-border transmission
386 (Supplemental Figure 2). The BA.1 haplotype bearing the S: Δ H69/V70 single repair
387 had the highest median transmission fitness, while Alpha with the S: Δ H69/V70 single
388 repair was not transmitted at all (Figure 3C). Incidentally, repair of S: Δ H69/V70 had no
389 significant impact on transmission fitness when co-occurring with repaired S: Δ Y144 in
390 Alpha background. In turn, BA.1 with both repairs had a lower median transmission
391 fitness than Alpha with both repairs. This was also the case of BA.1 with the
392 S: Δ V143/Y145 single repair compared with the analogous Alpha with the S: Δ Y144
393 single repair.

394 We further investigated the potential role of host age as an adaptive determinant
395 of NTD deletion repair in Alpha and BA.1 background. We observed a significant
396 difference in host age for the Alpha haplotype with both repairs, with transmitted
397 samples exhibiting an elevated age compared to non-transmitted samples. This finding
398 paralleled our previous observations for haplotype IV, which is genetically analogous
399 (see Figures 1A and 1D for reference). However, this pattern did not hold for the
400 Omicron BA.1 haplotype with both deletion repairs. In contrast, the single repair of
401 S: Δ H69/V70 in BA.1 did show such an association (Figure 3D). In short, we observed a
402 lack of correspondence in the population effect of NTD deletion repairs in both
403 backgrounds. This fact, combined with the well-described effect of deletions in Alpha
404 background (reviewed later in the Discussion section), suggested that the genetic
405 background in which these deletions emerge is likely to have a differential functional
406 impact.

407 ***NTD deletion repair in Omicron BA.1 background has a non-accumulative***
 408 ***effect on viral phenotype***

409 To investigate whether certain viral characteristics exist that act as drivers of deletion
 410 repair in Omicron, we analysed the effect of deletion repair patterns in the Omicron
 411 BA.1 using pseudotyped VSV bearing BA.1 spike proteins with repaired S: Δ H69/V70,
 412 S: Δ V143/Y145, or both. Specifically, we examined the efficiency of virus production,
 413 thermal stability, surface expression, susceptibility to antibody neutralization, and
 414 fusogenicity. To assess virus production, VSV was pseudotyped with all spike
 415 constructs under identical conditions at the same time, and the amount of virus
 416 produced titrated on Vero E6 cells (Figure 4A).



417 Figure 4. Experimental assessment of the effect of repairing NTD deletions in Omicron BA.1 background. (A)
 418 Comparison of viral titres obtained for each spike variant using pseudotyped VSV on the indicated cell line. (B)
 419 Comparison of the temperature resulting in 50% inactivation of pseudotyped VSV carrying the indicated S protein.
 420 (C) Surface expression of the spike protein in transfected HEK293T cells quantified using flow cytometry. (D)
 421 Reciprocal 50% neutralization titres (NT50) of sera from convalescent and vaccinated individuals. (E) Relative
 422 ability of each spike variant to drive cell-cell fusion relative to the Wuhan-Hu-1 spike. (F) Summary of the log₂-
 423 transformed fold change for each spike variant relative to that of BA.1 in terms of virus production, fusogenicity,
 424 NT50, surface expression and 50% inactivating temperature. Subfigures A-E represent the median and interquartile
 425 range of at least three replicates. Differences were assessed using Wilcoxon rank-sum tests (ns: $p > 0.05$; *: $p \leq 0.05$;
 426 **: $p \leq 0.01$; ***: $p \leq 0.001$).

428 Virus production in Vero E6 cells was not affected by S: Δ V143/Y145 repair (0.85-fold;
429 $p = 0.06$) but was significantly reduced upon repairing S: Δ H69/V70 (0.26-fold; $p =$
430 $6.6 \cdot 10^{-4}$) and the repair of both deletions (0.30-fold; $p = 1.6 \cdot 10^{-4}$). A similar effect was
431 observed when the titre of the virus was evaluated in cells expressing the TMPRSS2 co-
432 receptor, indicating a substantial negative effect of S: Δ H69/V70 repair by itself on virus
433 production. We assessed whether this was the result of reduced thermal stability of the
434 spike proteins by obtaining the temperature resulting in a 50% reduction of virus titre
435 (i.e. 50% inactivation temperature; Figure 4B). No significant differences were
436 observed, suggesting stability was not the driver of these differences.

437 Next, we examined whether the effects on virus production in Vero E6 cells stemmed
438 from altered cell expression of the different constructs by flow cytometry, using
439 polyclonal sera from four individuals (Figure 4C). We did not detect differences in the
440 median cell surface expression of the spike protein between the canonical BA.1 protein
441 and any deletion repair haplotype. However, the repair of S: Δ V143/Y145 resulted in
442 higher expression than S: Δ H69/V70 (2.2-fold; $p = 0.0030$), and than the double repair
443 (2.1-fold; $p = 0.032$), resembling the effect observed with virus production.

444 Then, we questioned whether deletion repair affected neutralization by polyclonal sera
445 from convalescent donors from the first epidemic wave in Spain and those dually
446 vaccinated with the Comirnaty mRNA vaccine ($n = 6$ each; Figure 4D). Significant
447 increases in susceptibility to neutralization in sera from vaccinated individuals against
448 viruses with repaired S: Δ V143/Y145 (2.1-fold; $p = 0.011$) or both S: Δ H69/V70 and
449 S: Δ V143/Y145 repaired (1.5-fold; $p = 0.020$) were observed. A similar trend was
450 observed with convalescent sera, but statistical significance was not reached ($p = 0.19$
451 and $p = 0.42$, respectively) due to higher intra sample variability. Thus, these results
452 indicate that increased neutralization is driven by the repair of S: Δ V143/Y145 in BA.1

453 background.

454 Finally, as cell to cell spread via fusion of the plasma membrane could potentially
455 reduce exposure to neutralizing antibodies, we examined whether deletion repair could
456 alter the ability of the spike protein to fuse cells (Figure 4E). Interestingly, the repair of
457 S:ΔV143/Y145 increased cell fusion relative to the BA.1 spike protein (1.5-fold; $p =$
458 0.026), while the repair of both S:ΔH69/V70 and S:ΔV143/Y145 led to a decrease of
459 more than 50% in the average fusogenicity (0.42-fold; $p = 0.021$). The presence of
460 S:ΔH69/V70 by itself did not seem to play a role in cell-cell fusion in a BA.1
461 background.

462 **Discussion**

463 Deletions in the SARS-CoV-2 genome have a significant impact on viral adaptation and
464 fitness, often surpassing the effects of single nucleotide variants (SNVs)
465 [9,14,16,42,48]. In fact, deletions in the NTD region of the spike protein are fixed in
466 prominent viral variants. Thus, our understanding of the repair patterns of deletion is
467 crucial for elucidating the genetic and phenotypic characteristics of the variants of
468 concern. In this work, we demonstrate that repairing these deletions can alter viral
469 characteristics and potentially influence the success of specific viral haplotypes. Our
470 findings provide valuable insights into the genetic and phenotypic characteristics of
471 these variants, shedding light on the factors driving their emergence and transmission
472 dynamics.

473 We first examined the global distribution of SARS-CoV-2 samples carrying
474 different combinatorial patterns of spike NTD marker segments in the Alpha and Delta
475 variants. Upon examination of the genetic variations among the most prevalent
476 combinatorial haplotypes, it became apparent that only a single haplotype is assigned to

477 the Alpha variant: the one bearing repaired S: Δ H69/V70 and S: Δ Y144 (IV, n = 736).
478 The fact that only one combinatorial haplotype was derived from the Alpha lineage —
479 compared to all the rest derived from Delta— could imply that Alpha viruses were more
480 constrained regarding NTD mutation interactions, although the difference in prevalence
481 and time frames between variants may have also played a role. Indeed, Alpha bearing
482 S: Δ E156/F157-R158G (X) and Delta bearing S: Δ H69/V70 and S: Δ Y144 (XI) share an
483 NTD that contains all Alpha and Delta lineage-defining changes. However, the former
484 has never been detected (n = 0), while the latter is the fourth most abundant
485 combinatorial haplotype (n = 304). The fact that an Alpha spike with S: Δ E156/F157-
486 R158G has never been recorded might be suggestive of an epistatic incompatibility with
487 infection success. It might also be related to these sites mapping to the same b-hairpin in
488 the NTD supersite as S:144 [49], which is typically deleted in Alpha. This also ties to
489 previous suggestions of host and variant-specific structural restrictions being imposed
490 strictly by the length of NTD loops [9] and emphasizes the significance of genetic
491 context in genotype-phenotype relationships.

492 Furthermore, we investigated clinical variables associated with these samples,
493 and found an association between host age and deletion repair in the Alpha background,
494 with Alpha viruses with repaired S: Δ H69/V70 and S: Δ Y144 being more prevalent
495 among older individuals. This suggests a non-essential role of certain deletions in
496 elderly hosts. The association between deletion repair and age was further supported by
497 the higher age in samples involved in transmission events. Besides, the association was
498 not driven by just a few outbreaks. Focusing on these deletions, both S: Δ H69/V70 and
499 S: Δ Y144 map to the NTD antigenic supersite [8,49] and have been recurrently
500 identified in immunocompromised individuals with chronic infections before
501 widespread vaccination against COVID-19 [7,10,11,14]. Earlier in the pandemic, Meng

502 et al. [42] reported the independent emergence of S:ΔH69/V70 after the occurrence of
503 infectivity-impairing amino acid changes that, in turn, could promote immune escape or
504 stronger receptor-binding affinity. Interestingly, this study also showed that
505 S:ΔH69/V70 compensated for the deleterious effect of these mutations in a surrogate
506 system by increasing viral infectivity and cell-cell fusogenicity in Alpha background.
507 Previous research has pointed out that variation in the NTD can act as a fine-tuning
508 vehicle for accommodating diverse pressures [9]. During the initial emergence of
509 S:ΔH69/V70 itself, the lack of vaccination limited selective pressures [50]. In this
510 perspective, repair of S:ΔH69/V70 and S:ΔY144 might be able to stay as a permissive
511 mutation in immunocompromised patients with a minimally constrained adaptive
512 landscape. This, in turn, could facilitate the further acquisition of otherwise deleterious
513 mutations.

514 In the second part of our study, we focused on the emergence of repaired
515 deletions in the Omicron BA.1 lineage compared to the Alpha variant, finding striking
516 differences. We detected a higher number of cases with repaired deletions in lineage
517 BA.1 than in Alpha. Interestingly, BA.1 with repaired S:ΔH69/V70 exhibited the
518 highest transmission fitness, while Alpha with the same repair did not transmit to any
519 extent. Additionally, we did not detect any age association with repaired deletion in
520 Omicron background like we detected in Alpha. This could be due to several factors,
521 such as the influence of viral genetic context or the different immune status of the
522 population when Alpha and Omicron BA.1 were circulating.

523 Survey efforts are expected to be biased by geographic and temporal differences
524 in sequencing efforts and lineage prevalence, and thus the phylogenetic analysis of
525 transmission and emergence may not capture all transmission events or account for
526 other epidemiological factors that affect viral success. However, phylogenetic

527 estimation of transmission and emergence rates effectively enabled us to overcome
528 these biases and better understand viral spread dynamics. The consistency between our
529 two methodologies for estimating transmission (operating under maximum likelihood
530 and maximum parsimony) reassures our inferences about the genetic relationships and
531 evolutionary dynamics within our dataset. In summary, our exhaustive approach
532 revealed a wide variation in the number of observations of these haplotypes,
533 highlighting their differences in their success and persistence over time. Overall, these
534 results demonstrate the distinct effect on viral success of common NTD markers
535 depending on their genomic background and lineage, highlighting once more the critical
536 importance of the genetic context when describing the genotype-phenotype relationship
537 of mutations.

538 To gain further insights into the viral characteristics driving deletion repair in
539 Omicron BA.1 background, we conducted an array of *in vitro* experiments using spike
540 proteins bearing different deletion repair patterns. We assessed virus production,
541 thermal stability, surface expression, susceptibility to neutralization, and fusogenicity of
542 each spike haplotype. The similarity in cell surface expression of the engineered spike
543 protein compared to the canonical spike rejected the possibility of other viral
544 phenotypes in our surrogate system being dependent solely on this factor. Our results
545 show that deletion repair patterns have an impact on virus production and infectivity,
546 with repair of S: Δ H69/V70 leading to reduced virus titres by itself. This viral phenotype
547 parallels that of the Alpha spike protein, with earlier studies reporting that repair of
548 S: Δ H69/V70 in this background results in a decrease in infectivity and spike
549 incorporation into virions [42]. We did not observe differences in thermal stability with
550 the BA.1 protein, suggesting that the effect of S: Δ H69/V70 and S: Δ V143/Y145 repair
551 does not affect full-protein stability.

552 However, deletion repair did affect sera neutralization. In our BA.1 background,
553 repaired S: Δ V143/Y145 resulted in an increase in sensitivity to neutralization by
554 polyclonal sera obtained from vaccinated individuals, in line with previous findings of
555 S: Δ Y144 facilitating escape from NTD neutralizing antibodies in the Alpha background
556 [51]. The same was not observed for sera from convalescent, "first-wave" patients. In
557 turn, repair of S: Δ H69/V70 did not affect neutralization. This paralleled prior research
558 on Alpha and other earlier variants showing that neither S: Δ H69/V70 nor its repair
559 influenced the spike sensitivity to NTD neutralizing antibodies [9,42,51]. These results
560 point to the variant-independent influence of NTD variability on immune effects.
561 Incidentally, the association of higher host age with viral transmission of S: Δ H69/V70
562 repair in BA.1 background may be attributed to a lack immune selection at the S:69/70
563 sites following vaccination efforts. This is coherent with deletion repair being able to
564 emerge in elderly patients with reduced adaptive constraints, potentially facilitating
565 functional adaptation during prolonged infections.

566 Research has shown that cell-cell fusion and subsequent formation of syncytia
567 can be mediated by viral membrane glycoproteins [52,53] such as the spike protein
568 found in SARS-CoV-2. In fact, this is a key feature of SARS-CoV-2 infection
569 [41,54,55]. In line with prior research, our results point to a lower fusogenicity of the
570 unaltered BA.1 spike protein compared with Wuhan-Hu-1, which has been attributed to
571 its inefficient spike cleavage and unfavoured TMPRSS2-mediated cell entry [42].
572 Interestingly, we found that deletion repair patterns exerted a non-accumulative
573 influence on the fusogenicity of the BA.1 spike protein. Proteins with repaired
574 S: Δ V143/Y145 alone promoted higher fusogenicity, while S: Δ H69/V70 alone did not
575 have a significant impact. However, repair of both deletions led to a significant decrease
576 in fusogenicity. In contrast, the spike protein of the Alpha variant has been shown to

577 exhibit an increased potential for cell-to-cell fusion. Nevertheless, upon S:ΔH69/V70
578 repair, the fusogenicity is reduced to levels comparable to the Wuhan-Hu-1 spike [42].

579 Regardless, our surrogate system might not fully reflect the complex interactions
580 of the spike protein with host cells and the immune system *in vivo*, and other regions of
581 the spike protein or the viral genome might also contribute to SARS-CoV-2 fitness and
582 adaptation. Besides, despite our attempts to account for demographic differences
583 between variants, our analysis could be affected by the non-random nature of the
584 sequences available in GISAID. The GISAID SARS-CoV-2 database is a remarkable
585 initiative and tool for monitoring, which has greatly advanced the knowledge of the
586 pandemic. Still, it also has significant limitations in terms of the origin and distribution
587 of the sequences that are submitted, and it cannot be assumed to reflect unbiased
588 genome surveillance. We also acknowledge that our suggestion of epistatic
589 incompatibility and genetic constraints is not conclusive, since we have limited our
590 analysis to a specific region of the S gene to narrow the number of possible
591 combinatorial haplotypes. A more comprehensive analysis could consider the whole
592 gene, or even the entire genome, thus requiring more computational and experimental
593 resources and the availability of sufficiently curated sequence data.

594 Our study highlights the importance of deletion repair in the spike protein NTD
595 in SARS-CoV-2 and its impact on viral fitness, transmission, and clinical characteristics
596 depending on the genetic context. We have provided novel insights into NTD marker
597 combinations in the Alpha and Delta variants, the repair patterns of NTD deletions in
598 the Alpha and Omicron BA.1 backgrounds, and the phenotypic consequences of NTD
599 deletion repair in the Omicron BA.1 background. There have been efforts to
600 characterize the variability of the S protein NTD through full swap assays, comparing
601 the ancestral (Wuhan-Hu-1) spike with that of Alpha and Omicron BA.1 [9], Delta and

602 Omicron BA.1 and BA.2 [42] and even with more distant virus like SARS-CoV [56].
603 However, to our best knowledge, no previous study has performed an exhaustive
604 combinatorial approach to characterize each separate marker in each genomic context.
605 By comparing the outcomes of deletion repair in different variants, we aim to gain
606 insights into the specific adaptive characteristics and genomic context that influence the
607 genotype-phenotype relationships. Our findings reveal that repair of specific deletions
608 in different genetic backgrounds can be driven by distinct phenotypic traits, such as
609 enhanced viral transmission, altered host age distribution, and changes in viral
610 characteristics, including fusogenicity, infectivity and susceptibility to neutralization.
611 Understanding these genotype-phenotype relationships can provide valuable insights
612 into the evolutionary dynamics and adaptation of SARS-CoV-2 variants, aiding in the
613 development of effective control strategies and therapeutic interventions. Future studies
614 building upon these findings will contribute to the development of effective strategies to
615 monitor and mitigate the impact of NTD deletions in emerging SARS-CoV-2 variants.

616 **Acknowledgements**

617 We thank Dr. Óscar González Recio (INIA-CSIC, Spain) for providing sequencing data
618 from a representative BA.1 sample with S:142G and repaired NTD deletions (GISAID
619 accession code: EPI_ISL_9805648), Dr. Luis Enjuanes (CNB, Spain) for providing
620 Vero E6 cells, and Dr. Markus Hoffman (German Primate Center,
621 Goettingen/Germany) for providing the hACE2 plasmid. We also thank Francisco José
622 Martínez Martínez (IBV-CSIC, Spain) for comments on phylogeny visualization and
623 analysis of transmission clusters. The computations were performed on the HPC cluster
624 Garnatxa at the Institute for Integrative Systems Biology (I²SysBio). The I²SysBio is a
625 is a joint collaborative research institute involving the University of Valencia (UV) and
626 the Spanish National Research Council (CSIC). We gratefully acknowledge all data

627 contributors, i.e., the Authors and their Originating laboratories responsible for
628 obtaining the specimens, and their Submitting laboratories for generating the genetic
629 sequence and metadata and sharing via the GISAID Initiative, on which this research is
630 based.

631 **Funding details**

632 This research work was funded by the European Commission – NextGenerationEU
633 (Regulation EU 2020/2094), through CSIC's Global Health Platform (PTI+ Salud
634 Global) to MC, RG, IC and FGC. MAH is supported by the Generalitat Valenciana and
635 the European Social Fund “ESF Investing in your future” through grant
636 CIACIF/2022/333. This work was also a part of projects CNS2022-135116 (MC) and
637 CNS2022-135100 (RG) funded by MCIN/AEI/10.13039/501100011033 and the
638 European Union NextGenerationEU/PRTR.

639 **Ethics statement**

640 Sera samples for the biological characterization of BA.1 deletion repair were obtained
641 from the La Fe University and Polytechnic Hospital of Valencia and were collected
642 after informed written consent had been obtained, with approval by the ethical
643 committee and institutional review board (registration number 2020-123-1).

644 **Contributions**

645 MAH, BND, PRR and MC conceived the theoretical framework. MAH and BND
646 devised the initial idea. MAH and MC conceived the second part of the study. MAH
647 implemented and performed the data retrieval, data curation and computations. JZ, BG,
648 and RG designed the experiments. JZ and BG performed the experiments. MAH carried
649 out the statistical analyses. MG and CAG contributed biological samples. MAB did
650 project management. MAH drafted the manuscript with support from BDN, PRR, RG

651 and MC. FGC aided in interpreting the results. MAH, PRR, BND, RG, FGC, IC and
652 MC discussed the results and commented on the manuscript. MC supervised the project.

653 **Declaration of interest statement**

654 The authors report there are no competing interests to declare.

655 **References**

- 656 [1] Rambaut A, Holmes EC, O’Toole Á, et al. A dynamic nomenclature proposal
657 for SARS-CoV-2 lineages to assist genomic epidemiology. *Nat Microbiol.*
658 2020;5:1403–1407.
- 659 [2] Chen Z, Chong KC, Wong MCS, et al. A global analysis of replacement of
660 genetic variants of SARS-CoV-2 in association with containment capacity and
661 changes in disease severity. *Clin Microbiol Infect.* 2021;27:750–757.
- 662 [3] da Silva Francisco Junior R, Lamarca AP, de Almeida LGP, et al. Turnover of
663 SARS-CoV-2 Lineages Shaped the Pandemic and Enabled the Emergence of
664 New Variants in the State of Rio de Janeiro, Brazil. *Viruses.* 2021;13:2013.
- 665 [4] Letko M, Marzi A, Munster V. Functional assessment of cell entry and receptor
666 usage for SARS-CoV-2 and other lineage B betacoronaviruses. *Nat Microbiol.*
667 2020;5:562–569.
- 668 [5] Harvey WT, Carabelli AM, Jackson B, et al. SARS-CoV-2 variants, spike
669 mutations and immune escape. *Nat Rev Microbiol.* 2021;19:409–424.
- 670 [6] V’kovski P, Kratzel A, Steiner S, et al. Coronavirus biology and replication:
671 implications for SARS-CoV-2. *Nat Rev Microbiol.* 2021;19:155–170.
- 672 [7] Kemp SA, Collier DA, Datir RP, et al. SARS-CoV-2 evolution during treatment
673 of chronic infection. *Nature.* 2021;592:277–282.
- 674 [8] McCallum M, De Marco A, Lempp FA, et al. N-terminal domain antigenic
675 mapping reveals a site of vulnerability for SARS-CoV-2. *Cell.* 2021;184:2332-
676 2347.e16.
- 677 [9] Cantoni D, Murray MJ, Kalemera MD, et al. Evolutionary remodelling of N-
678 terminal domain loops fine-tunes SARS-CoV-2 spike. *EMBO Rep.*
679 2022;n/a:e54322.

- 680 [10] Avanzato VA, Matson MJ, Seifert SN, et al. Case Study: Prolonged Infectious
681 SARS-CoV-2 Shedding from an Asymptomatic Immunocompromised
682 Individual with Cancer. *Cell*. 2020;183:1901-1912.e9.
- 683 [11] Choi B, Choudhary MC, Regan J, et al. Persistence and Evolution of SARS-
684 CoV-2 in an Immunocompromised Host. *N Engl J Med*. 2020;383:2291–2293.
- 685 [12] Minskaia E, Hertzog T, Gorbalenya AE, et al. Discovery of an RNA virus 3'→5'
686 exoribonuclease that is critically involved in coronavirus RNA synthesis. *Proc*
687 *Natl Acad Sci U S A*. 2006;103:5108–5113.
- 688 [13] Denison MR, Graham RL, Donaldson EF, et al. Coronaviruses: an RNA
689 proofreading machine regulates replication fidelity and diversity. *RNA Biol*.
690 2011;8:270–279.
- 691 [14] McCarthy KR, Rennick LJ, Nambulli S, et al. Recurrent deletions in the SARS-
692 CoV-2 spike glycoprotein drive antibody escape. *Science*. 2021;371:1139–1142.
- 693 [15] Planas D, Veyer D, Baidaliuk A, et al. Reduced sensitivity of SARS-CoV-2
694 variant Delta to antibody neutralization. *Nature*. 2021;596:276–280.
- 695 [16] Mishra T, Dalavi R, Joshi G, et al. SARS-CoV-2 spike E156G/Δ157-158
696 mutations contribute to increased infectivity and immune escape. *Life Sci*
697 *Alliance*. 2022;5:e202201415.
- 698 [17] Tamura T, Ito J, Uriu K, et al. Virological characteristics of the SARS-CoV-2
699 XBB variant derived from recombination of two Omicron subvariants. *Nat*
700 *Commun*. 2023;14:2800.
- 701 [18] Zhang X, Chen L-L, Ip JD, et al. Omicron sublineage recombinant XBB evades
702 neutralising antibodies in recipients of BNT162b2 or CoronaVac vaccines.
703 *Lancet Microbe*. 2022;0.
- 704 [19] Mykytyn AZ, Rosu ME, Kok A, et al. Antigenic mapping of emerging SARS-
705 CoV-2 omicron variants BM.1.1.1, BQ.1.1, and XBB.1. *Lancet Microbe*.
706 2023;0.
- 707 [20] Uriu K, Ito J, Zahradnik J, et al. Enhanced transmissibility, infectivity, and
708 immune resistance of the SARS-CoV-2 omicron XBB.1.5 variant. *Lancet Infect*
709 *Dis*. 2023;23:280–281.
- 710 [21] Yue C, Song W, Wang L, et al. ACE2 binding and antibody evasion in enhanced
711 transmissibility of XBB.1.5. *Lancet Infect Dis*. 2023;23:278–280.
- 712 [22] Khare S, Gurry C, Freitas L, et al. GISAID's Role in Pandemic Response. *China*
713 *CDC Wkly*. 2021;3:1049–1051.

- 714 [23] Hadfield J, Megill C, Bell SM, et al. Nextstrain: real-time tracking of pathogen
715 evolution. *Bioinformatics*. 2018;34:4121–4123.
- 716 [24] Aksamentov I, Roemer C, Hodcroft EB, et al. Nextclade: clade assignment,
717 mutation calling and quality control for viral genomes. *J Open Source Softw*.
718 2021;6:3773.
- 719 [25] Sanderson T, Barrett JC. Variation at Spike position 142 in SARS-CoV-2 Delta
720 genomes is a technical artifact caused by dropout of a sequencing amplicon.
721 *Wellcome Open Res*. 2021;6:305.
- 722 [26] R Core Team. R: A Language and Environment for Statistical Computing
723 [Internet]. Vienna, Austria: R Foundation for Statistical Computing; 2021.
724 Available from: <https://www.R-project.org/>.
- 725 [27] Wickham H, Averick M, Bryan J, et al. Welcome to the tidyverse. *J Open*
726 *Source Softw*. 2019;4:1686.
- 727 [28] Ou J, Zhu LJ. trackViewer: a Bioconductor package for interactive and
728 integrative visualization of multi-omics data. *Nat Methods*. 2019;16:453–454.
- 729 [29] Paradis E. Analysis of haplotype networks: The randomized minimum spanning
730 tree method. *Methods Ecol Evol*. 2018;9:1308–1317.
- 731 [30] Paradis E. pegas: an R package for population genetics with an integrated–
732 modular approach. *Bioinformatics*. 2010;26:419–420.
- 733 [31] McBroom J, Thornlow B, Hinrichs AS, et al. A Daily-Updated Database and
734 Tools for Comprehensive SARS-CoV-2 Mutation-Annotated Trees. *Mol Biol*
735 *Evol*. 2021;38:5819–5824.
- 736 [32] Turakhia Y, Thornlow B, Hinrichs AS, et al. Ultrafast Sample placement on
737 Existing tRees (USHER) enables real-time phylogenetics for the SARS-CoV-2
738 pandemic. *Nat Genet*. 2021;53:809–816.
- 739 [33] Fu L, Niu B, Zhu Z, et al. CD-HIT: accelerated for clustering the next-
740 generation sequencing data. *Bioinformatics*. 2012;28:3150–3152.
- 741 [34] Minh BQ, Schmidt HA, Chernomor O, et al. IQ-TREE 2: New Models and
742 Efficient Methods for Phylogenetic Inference in the Genomic Era. Teeling E,
743 editor. *Mol Biol Evol*. 2020;37:1530–1534.
- 744 [35] Ruiz-Rodríguez P, Francés-Gómez C, Chiner-Oms Á, et al. Evolutionary and
745 Phenotypic Characterization of Two Spike Mutations in European Lineage 20E
746 of SARS-CoV-2. *mBio*. 2021;12:e02315-21.

- 747 [36] Kassambara A. ggpubr: “ggplot2” Based Publication Ready Plots. 2023.
748 Available from: <https://rpkgs.datanovia.com/ggpubr/>.
- 749 [37] Giménez E, Albert E, Zulaica J, et al. Severe Acute Respiratory Syndrome
750 Coronavirus 2 Adaptive Immunity in Nursing Home Residents Following a
751 Third Dose of the Comirnaty Coronavirus Disease 2019 Vaccine. *Clin Infect*
752 *Dis.* 2022;75:e865–e868.
- 753 [38] Gozalbo-Rovira R, Gimenez E, Latorre V, et al. SARS-CoV-2 antibodies, serum
754 inflammatory biomarkers and clinical severity of hospitalized COVID-19
755 patients. *J Clin Virol.* 2020;131:104611.
- 756 [39] Grzelak L, Temmam S, Planchais C, et al. A comparison of four serological
757 assays for detecting anti-SARS-CoV-2 antibodies in human serum samples from
758 different populations. *Sci Transl Med.* 2020;12:eabc3103.
- 759 [40] García-Murria MJ, Expósito-Domínguez N, Duart G, et al. A Bimolecular
760 Multicellular Complementation System for the Detection of Syncytium
761 Formation: A New Methodology for the Identification of Nipah Virus Entry
762 Inhibitors. *Viruses.* 2019;11:229.
- 763 [41] Hoffmann M, Kleine-Weber H, Schroeder S, et al. SARS-CoV-2 Cell Entry
764 Depends on ACE2 and TMPRSS2 and Is Blocked by a Clinically Proven
765 Protease Inhibitor. *Cell.* 2020;181:271-280.e8.
- 766 [42] Meng B, Kemp SA, Papa G, et al. Recurrent emergence of SARS-CoV-2 spike
767 deletion H69/V70 and its role in the Alpha variant B.1.1.7. *Cell Rep.*
768 2021;35:109292.
- 769 [43] McCallum M, Walls AC, Sprouse KR, et al. Molecular basis of immune evasion
770 by the Delta and Kappa SARS-CoV-2 variants. *Science.* 2021;374:1621–1626.
- 771 [44] Domingo E. Chapter 7 - Long-Term Virus Evolution in Nature. In: Domingo E,
772 editor. *Virus Popul.* Boston: Academic Press; 2016. p. 227–262.
- 773 [45] Achaiah NC, Subbarajasetty SB, Shetty RM. R0 and Re of COVID-19: Can We
774 Predict When the Pandemic Outbreak will be Contained? *Indian J Crit Care Med*
775 *Peer-Rev Off Publ Indian Soc Crit Care Med.* 2020;24:1125–1127.
- 776 [46] Zhao L, Wymant C, Blanquart F, et al. Phylogenetic estimation of the viral
777 fitness landscape of HIV-1 set-point viral load. *Virus Evol.* 2022;8:veac022.
- 778 [47] Chaudhry MRA. Chapter 5 - Coronavirus infection outbreak: comparison with
779 other viral infection outbreak. In: Qureshi AI, Saeed O, Syed U, editors.
780 *Coronavirus Dis.* Academic Press; 2022. p. 47–57.

- 781 [48] Venkatakrishnan AJ, Anand P, Lenehan PJ, et al. Expanding repertoire of
782 SARS-CoV-2 deletion mutations contributes to evolution of highly transmissible
783 variants. *Sci Rep.* 2023;13:257.
- 784 [49] Cerutti G, Guo Y, Zhou T, et al. Potent SARS-CoV-2 neutralizing antibodies
785 directed against spike N-terminal domain target a single supersite. *Cell Host*
786 *Microbe.* 2021;29:819-833.e7.
- 787 [50] MacLean OA, Lytras S, Weaver S, et al. Natural selection in the evolution of
788 SARS-CoV-2 in bats created a generalist virus and highly capable human
789 pathogen. *PLoS Biol.* 2021;19:e3001115.
- 790 [51] Graham C, Seow J, Huettner I, et al. Neutralization potency of monoclonal
791 antibodies recognizing dominant and subdominant epitopes on SARS-CoV-2
792 Spike is impacted by the B.1.1.7 variant. *Immunity.* 2021;54:1276-1289.e6.
- 793 [52] Li W, Moore MJ, Vasilieva N, et al. Angiotensin-converting enzyme 2 is a
794 functional receptor for the SARS coronavirus. *Nature.* 2003;426:450–454.
- 795 [53] Tseng C-TK, Tseng J, Perrone L, et al. Apical Entry and Release of Severe
796 Acute Respiratory Syndrome-Associated Coronavirus in Polarized Calu-3 Lung
797 Epithelial Cells. *J Virol.* 2005;79:9470–9479.
- 798 [54] Buchrieser J, Dufloo J, Hubert M, et al. Syncytia formation by SARS-CoV-2-
799 infected cells. *EMBO J.* 2020;39:e106267.
- 800 [55] Bussani R, Schneider E, Zentilin L, et al. Persistence of viral RNA, pneumocyte
801 syncytia and thrombosis are hallmarks of advanced COVID-19 pathology.
802 *EBioMedicine.* 2020;61:103104.
- 803 [56] Qing E, Kicmal T, Kumar B, et al. Dynamics of SARS-CoV-2 Spike Proteins in
804 Cell Entry: Control Elements in the Amino-Terminal Domains. *mBio.*
805 2021;12:10.1128/mbio.01590-21.