



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Distinct patterns of within-host virus populations 2 between two subgroups of human respiratory syncytial 3 virus

Citation for published version:

Lin, G-L, Drysdale, SB, Snape, MD, O'Connor, DP, Brown, A, MacIntyre-Cockett, G, Mellado-Gomez, E, de Cesare, M, Bonsall, D, Ansari, MA, Öner, D, Aerssens, J, Butler, C, Bont, L, Openshaw, P, Martinon-Torres, F, Nair, H, Bowden, R, Golubchik, T & Pollard, AJ 2021, 'Distinct patterns of within-host virus populations 2 between two subgroups of human respiratory syncytial 3 virus', *Nature Communications*.
<https://doi.org/10.1038/s41467-021-25265-4>

Digital Object Identifier (DOI):

[10.1038/s41467-021-25265-4](https://doi.org/10.1038/s41467-021-25265-4)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Nature Communications

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 Distinct patterns of within-host virus populations
2 between two subgroups of human respiratory syncytial
3 virus

4 Gu-Lung Lin^{1,2,*}, Simon B Drysdale^{1,2,3}, Matthew D Snape^{1,2}, Daniel O'Connor^{1,2},
5 Anthony Brown⁴, George MacIntyre-Cockett⁵, Esther Mellado-Gomez⁵, Mariateresa de
6 Cesare⁵, David Bonsall^{5,6}, M Azim Ansari⁵, Deniz Öner⁷, Jeroen Aerssens⁷, Christopher
7 Butler⁸, Louis Bont^{9,10}, Peter Openshaw¹¹, Federico Martínón-Torres^{12,13}, Harish Nair¹⁴,
8 Rory Bowden^{5,15}, RESCEU Investigators^a, Tanya Golubchik⁶, and Andrew J Pollard^{1,2}

9 ¹Oxford Vaccine Group, Department of Paediatrics, University of Oxford, Oxford, UK.

10 ²NIHR Oxford Biomedical Research Centre, Oxford, UK.

11 ³Present address: Paediatric Infectious Diseases Research Group, Institute for Infection and Immunity, St
12 George's, University of London, London, UK.

13 ⁴Peter Medawar Building for Pathogen Research, University of Oxford, Oxford, UK.

14 ⁵Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK.

15 ⁶Big Data Institute, Nuffield Department of Medicine, University of Oxford, Oxford, UK.

16 ⁷Translational Biomarkers, Infectious Diseases Therapeutic Area, Janssen Pharmaceutica NV, Beerse,
17 Belgium.

18 ⁸Nuffield Department of Primary Care Health Sciences, University of Oxford, Oxford, UK.

19 ⁹Department of Pediatrics, Wilhelmina Children's Hospital, University Medical Center Utrecht, Utrecht,
20 Netherlands.

21 ¹⁰ReSViNET Foundation, Zeist, Netherlands.

22 ¹¹National Heart and Lung Institute, Imperial College London, London, UK.

23 ¹²Translational Pediatrics and Infectious Diseases, Hospital Clínico Universitario de Santiago de
24 Compostela, Santiago de Compostela, Spain.

25 ¹³Genetics, Vaccines, Infectious Diseases, and Pediatrics Research Group (GENVIP), Instituto de
26 Investigación Sanitaria de Santiago de Compostela, Santiago de Compostela, Spain.

27 ¹⁴Centre for Global Health, Usher Institute, Edinburgh Medical School, University of Edinburgh,

28 Edinburgh, UK.

29 ¹⁵Present address: Division of Advanced Technology and Biology, Walter and Eliza Hall Institute of
30 Medical Research, Melbourne, Australia.

31 ^aA full list of consortium members appears at the end of the paper.

32 These authors jointly supervised this work: Tanya Golubchik, Andrew J Pollard.

33 *E-mail: gu-lung.lin@paediatrics.ox.ac.uk

34 Abstract

35 Human respiratory syncytial virus (RSV) is a major cause of lower respiratory
36 tract infection in young children globally, but little is known about within-host
37 RSV diversity. Here, we characterised within-host RSV populations using deep-
38 sequencing data from 319 nasopharyngeal swabs collected during 2017–2020. RSV-B
39 had lower consensus diversity than RSV-A at the population level, while exhibiting
40 greater within-host diversity. Two RSV-B consensus sequences had an amino acid
41 alteration (K68N) in the fusion (F) protein, which has been associated with reduced
42 susceptibility to nirsevimab (MEDI8897), a novel RSV monoclonal antibody under
43 development. In addition, several minor variants were identified in the antigenic
44 sites of the F protein, one of which may confer resistance to palivizumab, the only
45 licensed RSV monoclonal antibody. The differences in within-host virus populations
46 emphasise the importance of monitoring for vaccine efficacy and may help to explain
47 the different prevalences of monoclonal antibody-escape mutants between the two
48 subgroups.

49 Introduction

50 Human respiratory syncytial virus (RSV) is the leading cause of lower respiratory tract
51 infection (LRTI) in young children, globally responsible for around 33 million episodes
52 of LRTI in children under 5 years of age annually with a disproportionately high burden
53 in infants younger than 1 year of age¹. Repeated infection is common throughout life²,
54 usually resulting in mild symptoms, but it can also cause serious disease in older (age
55 ≥ 65 years) or immunocompromised adults and people with chronic cardiopulmonary

56 disease³. Despite decades of effort, there is no efficacious antiviral for treatment or licensed
57 vaccine to prevent RSV infection, and thus the standard of care is supportive management
58 only. Palivizumab, an RSV-specific humanised monoclonal antibody, is the only available
59 immunoprophylactic agent. It requires multiple administrations over the RSV season and
60 is very expensive, so its use is limited to the highest-risk populations, namely infants born
61 preterm and those with congenital heart disease, chronic pulmonary disorders, or severe
62 combined immunodeficiency⁴.

63 RSV is a negative-sense single-stranded RNA virus with a genome containing 10 genes.
64 The F gene encodes the fusion (F) glycoprotein, which mediates the fusion of host cell and
65 viral membranes. The F protein is the main target for antibody-mediated neutralisation,
66 and has been the focus of the development of vaccines and monoclonal antibodies⁵. Through
67 the fusion process, the F protein changes from the prefusion to postfusion conformation.
68 Several antigenic sites (neutralising epitopes in particular) have been located on the surface
69 of the F protein. Antibodies exclusively targeting prefusion-specific antigenic sites (e.g.,
70 sites \emptyset and V) are more potent than those targeting sites that can be found in both
71 conformations (e.g., sites I, II, IV)⁶. Nirsevimab (MEDI8897), a recombinant human
72 monoclonal antibody currently in phase 3 clinical trials, exclusively targets antigenic site
73 \emptyset ⁷, and suptavumab (REGN2222), another prefusion-specific monoclonal antibody, binds
74 antigenic site V⁸. Palivizumab and its affinity-enhanced variant, motavizumab⁹, target
75 antigenic site II, and antibody 101F binds antigenic site IV¹⁰. Mutations in the antigenic
76 sites that confer resistance to monoclonal antibodies have been identified. For example,
77 mutants with N262S/Y, N268I, K272E/N/M/T/Q, or S275F/L in the F protein are less
78 susceptible to palivizumab¹¹⁻¹³, and nirsevimab has reduced neutralising activity against
79 mutants with N67I/N208Y, N208S/D, K68N/N201S, or K68N/N208S in the F protein⁷.

80 The G gene encodes the attachment (G) glycoprotein, a transmembrane protein
81 responsible for viral attachment. The extracellular portion (ectodomain) of the G protein
82 consists of two hypervariable mucin-like regions flanking a conserved central domain
83 (CCD)¹⁴. The CCD, containing antigenic sites γ 1 and γ 2, has been shown to be a target
84 for neutralising antibodies¹⁵ and is another focus of vaccine development^{16,17}. Outside the

85 CCD, the mucin-like regions also have multiple antigenic sites though less well-defined¹⁸.
86 The mucin-like region II (2nd hypervariable region) has been shown to have hypermutation
87 at the population level, and has thus been used widely in phylogenetic analyses¹⁹.

88 The two subgroups of RSV (A and B) co-circulate in epidemics, and both exhibit
89 rapid evolutionary dynamics²⁰. Molecular epidemiology and evolutionary dynamics of
90 RSV have been extensively studied at the consensus level; however, little is known about
91 virus populations in each infected individual (i.e., within-host or intrahost virus diversity).
92 Using high-throughput whole-genome sequencing, it is now possible to sequence viruses in
93 sufficient depth to obtain a complete picture of within-host populations. A previous study
94 showed that within-host RSV diversity increased in an immunocompromised infant with
95 persistent RSV infection following a haematopoietic stem cell transplant, and palivizumab
96 escape mutants emerged after multiple administrations of this drug²¹. Another study
97 demonstrated that RSV-A exhibited greater within-host virus diversity in experimentally
98 infected adults than in naturally infected infants²². However, these results were limited to
99 RSV-A infection and did not look at natural infections in adult populations. Analysing
100 within-host virus genetic diversity in infections that represent general seasonal epidemics
101 can aid understanding of the patterns of virus evolution and its driving forces, informing
102 the development of preventative and treatment measures.

103 In this study, we seek to characterise within-host RSV populations for the two subgroups,
104 RSV-A and RSV-B, using deep sequencing of samples collected from participants in three
105 prospective clinical studies. We find that RSV-B exhibits greater within-host diversity
106 than RSV-A, with two RSV-B consensus strains and one RSV-B minor variant likely
107 conferring resistance to nirsevimab or palivizumab. We also show that temporal changes of
108 intrahost viral populations follow stochastic patterns. Our work highlights the importance
109 of continued genetic surveillance of RSV to ensure the effectiveness of future RSV vaccines
110 and therapeutics.

111 **Results**

112 **Sample population**

113 We sequenced RSV from 858 nasopharyngeal swabs collected from 459 RSV-infected
114 patients in the United Kingdom, Spain, and the Netherlands during 2017–2020. Of these,
115 327 samples had sufficient viral load to generate more than 10,000 unique (deduplicated)
116 RSV reads. After removing five samples containing both RSV-A and RSV-B, 322 samples
117 were included in the within-host virus diversity analysis. Sequencing was carried out in four
118 batches, with 11, 113, 41, and 157 of the included samples from each batch respectively
119 (Supplementary Table 1). The 322 samples were collected from 267 different participants,
120 among which 34 participants had multiple samples (mean 2.6, range 2–5) collected on
121 different days (ranging from 1 to 8 days apart).

122 **Cumulative minor allele frequencies and minor variants**

123 Genomic positions with a read depth of less than 200 were excluded from the analysis.
124 Nearly 90% of the samples had >80% of the genome passing this threshold. Three samples
125 had a significantly high mean cumulative minor allele frequency (MAF) per sample: 0.52%
126 (from an RSV-A-infected infant; batch 4), 0.19% (from an RSV-B-infected adult; batch
127 2), and 0.17% (from an RSV-B-infected infant; batch 4). These samples presumably
128 represented a real or artefactual mixture of genetically distinct strains of the same RSV
129 subgroup and were thus excluded from the following analysis. The sources and sequencing
130 yields of the remaining 319 samples (collected from 264 participants) are shown in Table 1.

131 The median of the mean cumulative MAF per sample was 0.039% (range 0.025%–
132 0.068%) for the 319 samples. The distributions of the mean cumulative MAF per sample
133 were significantly different between samples from different sequencing batches (Supple-
134 mentary Fig. 1a), likely due to the differences in the ratio of duplicate read counts to
135 total RSV read counts (percent duplication rate) between batches (Supplementary Table
136 1). After adjusting for the observed batch effects (e.g., Supplementary Fig. 1b), RSV-B
137 samples had a higher mean cumulative MAF per sample than RSV-A samples (median of

138 the original data: 0.042% vs. 0.037%; multiple linear regression with batch and the number
139 of unique RSV reads as covariates, $P = 0.016$; Mann–Whitney U test on standardised
140 data, $P = 0.016$).

141 On average, each sample had 3.7 minor variants (range 0–30; defined as variants with
142 a frequency of $\geq 3\%$). 18.8% of the samples (60/319) did not have any minor variants. An
143 inverse correlation was noted between the number of unique RSV reads and the number
144 of minor variants ($r = -0.41$, $P = 4.2 \times 10^{-14}$; Supplementary Fig. 2), consistent with
145 a greater variance of MAF when the sampling fraction was small (i.e., few unique reads
146 were sequenced)²³. Variation rarely occurred at the same genomic position in different
147 samples. Among all minor variants found in this study, only 5.9% (57/972) were shared
148 by multiple samples (excluding 17 minor variants only shared by sequential samples from
149 the same participants), usually no more than five samples. However, there was one minor
150 variant shared by 59% (85/144) of the RSV-B samples, with a frequency between 3% and
151 11%. This minor variant had a G to A substitution at position 3403 of the L gene, causing
152 an amino acid alteration from glutamic acid to lysine at position 1135 (E1135K) of the
153 RNA-dependent RNA polymerase.

154 **Potential antigenic variants**

155 The sequences encoding the antigenic sites of the F protein were highly conserved at the
156 consensus level in this study. However, two RSV-B isolates from two infant participants,
157 both of whom had only one sample collected, had an A to T substitution at nucleotide
158 position 204 of the F gene. This substitution results in an amino acid alteration from lysine
159 to asparagine (K68N), which in a previous study was associated with a 4-fold reduction in
160 susceptibility to nirsevimab neutralisation in vitro⁷. No minor variant was found at this
161 position in these two samples.

162 The frequencies and distribution of all minor variants across the coding sequence of
163 the F gene are shown in Fig. 1a. There were one, eight, two, and three minor variants
164 identified in the antigenic sites \emptyset , II, IV, and V of the F protein, respectively (Table 2).
165 0%, 6.0% (6/100), and 1.6% (2/124) of the participants had potential antigenic variants

166 (i.e., minor variants encoding a nonsynonymous substitution in the antigenic sites) in the
167 2017–18, 2018–19, and 2019–20 RSV seasons, respectively. One of these minor variants
168 had two nucleotide substitutions with a frequency of $\geq 3\%$ in a single codon, encoding
169 an amino acid substitution from isoleucine to threonine at position 261 (I261T). Other
170 minor variants identified in the antigenic sites were from different samples. To date, none
171 of these variants have been reported to confer resistance to monoclonal antibodies.

172 We also looked at the frequencies and distribution of minor variants in the coding region
173 of the G gene (Fig. 1b). The median frequency of minor variants was significantly higher in
174 the G gene than in the F gene, either at potential antigenic sites (median: 9.3% vs. 4.6%;
175 Mann–Whitney U test, $P = 0.022$) or across the whole coding sequences (median: 8.3%
176 vs. 4.4%; Mann–Whitney U test, $P = 0.004$), consistent with previous studies identifying
177 the G gene as the most variable gene in the virus genome¹⁴. The median minor variant
178 frequency in the mucin-like region II of the G gene (13.7%) was greater than that in the
179 mucin-like region I (9.2%), which was greater than that in the CCD (4.0%). However,
180 these differences were not statistically significant (Kruskal–Wallis test, $P = 0.20$).

181 **Pairwise nucleotide diversity**

182 Within-host virus genetic diversity was estimated as pairwise nucleotide diversity (see
183 Methods). Pairwise nucleotide diversity did not correlate with the number of unique RSV
184 reads after adjusting for the batch effects (Supplementary Table 2 and Supplementary Fig.
185 3a), but was highly consistent with the mean cumulative MAF per sample ($r = 0.997$, P
186 $< 2.2 \times 10^{-16}$; Supplementary Fig. 3b). The median pairwise nucleotide diversity of the
187 whole dataset was 0.0007 (range 0.0005–0.0014). Gene-wise comparisons showed that the
188 L gene had significantly higher pairwise nucleotide diversity than the NS2, P, SH, and
189 G genes, but the other genes did not have significant differences in pairwise nucleotide
190 diversity between each other (Supplementary Fig. 4). These significant differences were by
191 definition due to the mean proportion of pairwise nucleotide differences at each genomic
192 position within the L gene instead of the length of the L gene.

193 RSV-B had greater pairwise nucleotide diversity than RSV-A (multiple linear regression,

194 $P = 0.044$, Supplementary Table 2; Fig. 2a), and older adults had a more diverse intrahost
195 RSV-B population than infants (multiple linear regression, $P = 0.0006$, Supplementary
196 Table 2; Fig. 2b). The subgroup difference was still significant if excluding adult samples
197 (Mann–Whitney U tests on standardised data, $P = 0.039$). The number of RSV reads and
198 the duration between symptom onset and sample collection were similar between both
199 RSV subgroups and between both age groups. Samples collected from different countries
200 or seasons or patients with different severity of RSV infections did not have significant
201 differences in pairwise nucleotide diversity (Supplementary Table 2).

202 Genetic distance

203 Within-host diversity levels between samples were compared using pairwise Manhattan
204 distances²⁴ at consensus-identical positions, where allele frequencies below the 3% threshold
205 were converted to 0. In contrast, consensus variations between samples were compared
206 using pairwise patristic distances, which are phylogenetic distances on RSV phylogenies
207 (Supplementary Fig. 5). To eliminate the batch effects, we only included pairwise distances
208 between samples in the second batch ($n = 112$; excluding one outlier). To reduce potential
209 bias from geographical and temporal differences, only pairwise distances between samples
210 from the same country and the same season were calculated.

211 Serial sample pairs had within-host diversity levels comparable to those of samples from
212 different participants (range: 0–3.34 vs. 0–5.03), despite having identical or nearly identical
213 consensus sequences, as indicated by their small patristic distances (range 2.0×10^{-6} –
214 7.5×10^{-5}). Excluding the serial sample pairs, RSV-B sample pairs had significantly
215 greater within-host diversity levels than RSV-A pairs (median: 1.24 vs. 0.86), whereas
216 the comparison of consensus sequences showed the opposite effect (Fig. 2c,d). Pairwise
217 patristic distances between RSV-A samples formed three clusters, corresponding to the
218 three main clades of the phylogenetic tree (Supplementary Fig. 5a). When using all allele
219 frequencies, including those below 3% MAF, to calculate Manhattan distances, RSV-B
220 sample pairs still had significantly greater pairwise Manhattan distances than RSV-A pairs
221 (median: 20.5 vs. 18.2, $P = 8.2 \times 10^{-58}$; Supplementary Fig. 6).

Temporal change of intrahost virus population

Putting all samples together, standardised pairwise nucleotide diversity did not have a significant temporal change within 7 days of symptom onset ($R^2 = 0.008$; $P = 0.122$). For the 34 participants with multiple samples collected daily during hospitalisation, pairwise nucleotide diversity was also evaluated in each set of serially collected samples, excluding those sequenced in different batches (Fig. 3). No significant trend was noted either in each participant or when combining all samples and adjusting for the batch effects. The only exception was the samples from GB-058, where pairwise nucleotide diversity increased by 0.000063 daily (95% confidence interval, 0.000046 to 0.000080; $P = 0.004$). This patient was a 19-day-old preterm neonate (gestational age of 33 weeks 6 days) with severe RSV infection requiring intensive care and mechanical ventilation.

The changes in minor variants and variant frequencies in the serial samples were also evaluated at polymorphic sites where minor alleles were identified at more than three time points (Fig. 4). 79% of these minor variants had a nonsynonymous substitution. Only one minor variant with a G to A substitution at position 3403 of the L gene from participant NL-091, which was shared by 71 participants (85 samples), remained above the 3% threshold throughout the sampling period. This patient was a 42-day-old previously healthy infant with severe RSV infection requiring intensive care and mechanical ventilation. All other variants (including the aforementioned variant in other participants) were only detected either early, late, or intermittently during the course of sample collection.

Discussion

In this study, we sequenced 858 nasopharyngeal samples collected in three clinical studies during 2017–2020, and profiled within-host RSV populations from 319 samples. We demonstrated that RSV-B had greater within-host diversity than RSV-A, whereas RSV-A had greater consensus diversity than RSV-B. Two RSV-B isolates' consensus sequences had a mutation in the F protein (K68N), previously associated with reduced susceptibility to nirsevimab neutralisation. Several other minor variants were also identified in the

249 antigenic sites of the F protein. None of these variants have been reported before except for
250 S255N²⁵, whose susceptibility to monoclonal antibodies has not been examined. Stochastic
251 (random) patterns were found in the temporal changes of within-host virus diversity and
252 minor variants.

253 Low input genetic material (i.e., viral load) has been shown to reduce the sensitivity
254 and specificity of variant calling²⁶. In this study, we applied the quantitative methodology
255 of targeted metagenomics to library construction and used the number of unique RSV reads
256 as a proxy for viral load²⁷. The inclusion criterion of more than 10,000 unique RSV reads
257 corresponded with a viral load of approximately 2.4×10^6 copies/mL and above, sufficient
258 input levels for accurate minority variant calling²⁸. Given the large number of samples in
259 this study, batching was required for sequencing, resulting in variable percent duplication
260 rates and hence some batch effects on diversity metrics. We adopted two approaches to
261 account for the batch effects on the comparisons of mean cumulative MAF per sample
262 and pairwise nucleotide diversity: (i) including batch as a regression covariate and (ii)
263 standardising the values within each batch to z-scores (see Methods for details). Both
264 methods showed the same significant findings, making cross-batch comparisons robust. To
265 avoid any residual bias, for pairwise comparisons of genetic distances we used only samples
266 from the same batch (batch 2), which had very high percent duplication rates and similar
267 read counts for RSV-A and RSV-B (Table 1 and Supplementary Table 1), consistent with
268 capture saturation, and from which we could be confident of recovering the full range of
269 intrahost diversity.

270 The extent of intrahost virus diversity depends not only on the rate of virus evolution
271 (partly associated with the ability of proofreading for viral replication errors) but also on
272 the duration of infection. RNA viruses generally have a higher mutation rate than DNA
273 viruses²⁹, and are usually not able to correct the errors of viral replication, which DNA
274 viruses can³⁰. In our study, RSV had greater pairwise nucleotide diversity than has been
275 reported for influenza virus, another RNA virus causing acute respiratory infection (range
276 0.0005–0.0014 vs. 0–0.0002³¹). RSV intrahost diversity appears to be comparable with, or
277 slightly higher than, that of the DNA viruses in the family *Herpesviridae*, which cause

278 chronic infections³², but up to one to two orders of magnitude lower than that of persistent
279 RNA viruses (e.g., hepatitis C virus, human immunodeficiency virus) and persistent DNA
280 viruses (e.g., hepatitis B virus), which generally have pairwise nucleotide diversity above
281 0.005³².

282 Neutralisation escape mutants have been isolated in 0.7% of immunoprophylaxis-naïve
283 RSV-infected subjects¹³, 5–9% of RSV-breakthrough patients receiving palivizumab^{12,33},
284 and 8% of RSV-breakthrough cases receiving nirsevimab³⁴. In our study, isolates collected
285 from 0.8% (2/264) of the immunoprophylaxis-naïve participants were found to contain a
286 nirsevimab resistance-associated substitution at the consensus level. We also identified
287 an RSV-B minor variant with an amino acid change from serine to proline at position
288 275 (S275P) of the F protein. Other amino acid substitutions at this position have
289 demonstrated resistance to palivizumab (S275F/L)¹². Whether the mutation S275P also
290 alters the neutralising activity of palivizumab requires further investigation; however, all
291 three mutations at this position replaced a polar amino acid with a nonpolar one, which
292 may result in significant conformational or functional changes. It is important to identify
293 neutralisation escape mutants in immunoprophylaxis-naïve children in the era before RSV
294 monoclonal antibodies become extensively used. It indicates the circulation of escape
295 mutants in the community even though they generally have a selective disadvantage in
296 the absence of monoclonal antibodies¹³.

297 Our findings that RSV-B had greater pairwise nucleotide diversity and pairwise Man-
298 hattan distances than RSV-A both indicate that, at least in our dataset, RSV-B had a
299 more diverse intrahost virus population than RSV-A. These results do not correlate with
300 the duration between symptom onset and sample collection (Table 1), but are consist-
301 ent with previous studies on global RSV strains, which found that RSV-B has a higher
302 genome-wide evolutionary rate than RSV-A ($7.47\text{--}7.76 \times 10^{-4}$ substitutions/site/year
303 vs. $5.68\text{--}6.47 \times 10^{-4}$ substitutions/site/year)^{35,36}. This difference extends below the 3%
304 threshold for minority variant calling (Supplementary Fig. 6). On the basis of these
305 findings, we hypothesise that RSV-B is subject to greater immune pressure (e.g., by
306 innate immunity, neutralising antibodies, or T cell-mediated cytotoxicity) than RSV-A.

307 This hypothesis is in line with previous studies showing that intrahost RSV diversity
308 increased in response to an established immunity²¹, and that RSV-B has more amino acid
309 alterations³⁷, predicted O-glycosylation site changes³⁷, and indel mutations³⁶ in the G gene
310 than RSV-A, suggesting a stronger selective pressure acting on RSV-B than on RSV-A.

311 RSV-B exhibited higher within-host diversity in older adults than in infants in response
312 to different immune pressures between the two age groups. Of note, our dataset included
313 only eight adults, and this comparison was limited to seven adult samples and 137 infant
314 samples collected from those with RSV-B infection. Further studies enrolling more adults
315 would be of value to delineate the difference in within-host diversity between different age
316 groups. Furthermore, the temporal changes of pairwise nucleotide diversity and minor
317 variants were stochastic within each infected individual, suggesting the driving force of
318 evolutionary dynamics in global RSV populations is more likely from the selective pressure
319 imposed at the population level than within an individual host. Only samples that yielded
320 sufficient RSV reads were included in this study, so these temporal trends were confined
321 to samples collected over a short time frame (mostly within 5 days of symptom onset).
322 Nonetheless, a study on seasonal influenza virus also found limited evidence of positive
323 selection at the within-host evolutionary scale²⁴.

324 The greater within-host virus diversity observed in RSV-B than in RSV-A warrants
325 separate testing and close monitoring of the anti-RSV-B efficacy of vaccines and monoclonal
326 antibodies that are being developed. This is because the development of several RSV
327 vaccines in preclinical or clinical trials is based on the nucleotide sequences or structure of
328 RSV-A strains³⁸⁻⁴⁰. Some studies have also shown that RSV-B had more fixed mutations
329 in the antigenic sites of the F protein at the consensus level⁴¹, resulting in more variable
330 in vitro and clinical susceptibility to monoclonal antibodies than RSV-A. For example,
331 in a phase 2b trial of nirsevimab, the drug had reduced neutralising activity against two
332 RSV-B isolates collected from its recipients; one had a mutation of N208S and the other
333 had multiple mutations of I64T, K68E, I206M, and Q209R in the F protein³⁴. A phase
334 3 trial of another investigational RSV monoclonal antibody, suptavumab, failed to meet
335 its primary end point because all RSV-B strains identified in the trial carried two amino

336 acid changes in the F protein (L172Q and S173L), conferring resistance to the drug⁸. All
337 RSV-B samples in our study also harboured these two amino acid substitutions, except
338 for one that encoded isoleucine instead of leucine at position 173 (a nonpolar-to-nonpolar
339 substitution).

340 We excluded genomic positions where consensus bases were different from the calculation
341 of Manhattan distance, to ensure that between-host genetic distance would be driven by
342 differences in minor alleles rather than differences at the consensus level²⁴. We found that,
343 outside the consensus-different positions, serial samples from the same individual did not
344 have a shorter pairwise Manhattan distance than that of a randomly taken between-host
345 pair from the same country and season. This methodology change makes our results robust
346 to inter-host variation, in contrast to previous studies on influenza virus and RSV, where
347 distance metrics were largely driven by consensus differences^{42,43}.

348 Our findings suggest that RSV-B has a more diverse within-host population than
349 RSV-A, likely driven by selection pressure at the host-population level. This difference
350 between the two subgroups warrants close monitoring of vaccine efficacy and emergence of
351 neutralisation escape variants.

352 **Methods**

353 **Sample collection**

354 Nasopharyngeal swabs were collected from patients with respiratory symptoms under 1
355 year old or over 60 years old, from London and Oxford, United Kingdom, Santiago de
356 Compostela, Spain, and Utrecht, the Netherlands, during 2017–2020. These patients were
357 enrolled in three clinical studies of the REspiratory Syncytial virus Consortium in EUrope
358 project (RESCEU, ClinicalTrials.gov identifiers: NCT03627572⁴⁴, NCT03756766⁴⁵, and
359 NCT03621930⁴⁶), a European multicentre project investigating epidemiological, virological,
360 and immunological characteristics of RSV infection. None of these participants had received
361 any RSV monoclonal antibody or investigational vaccine. RSV infection was diagnosed
362 using molecular point-of-care testing on the AlereTM i RSV platform (Abbott, Illinois, US)

363 in infant participants and on the GeneXpert[®] influenza/RSV system (Cepheid, California,
364 US) in adult participants in a community setting, and using antigen and/or PCR tests
365 at a central laboratory in a hospital setting. A nasopharyngeal swab was collected from
366 each participant within 7 days of symptom onset, and daily swabs were also collected from
367 RSV-positive hospitalised infant participants where possible until hospital discharge. After
368 collection, swabs were immersed in M4RT[®] transport medium, aliquoted, and frozen at
369 -80°C until use.

370 Severity of an RSV infection was defined using the ReSVinet scale⁴⁷ in infants. This scale
371 accounts for several clinical variables, including feeding intolerance, medical intervention,
372 respiratory difficulty, respiratory frequency, apnoea, general condition, and fever. The
373 score ranges from 0 to 20; a score of 0–7 was defined as mild, a score of 8–13 as moderate,
374 and a score of 14–20 as severe. In older adults, those who did not require any treatment or
375 medical attendance were defined as having a mild disease, those requiring hospitalisation
376 were defined as having a severe disease, and the rest were defined as having a moderate
377 RSV disease.

378 These clinical studies were conducted in accordance with the provisions of the De-
379 clarations of Helsinki and were approved by the relevant ethics committees at each site,
380 including the University of Oxford, the Health Research Authority (IRAS IDs: 224156 and
381 231136), the NHS National Research Ethics Service Oxfordshire Committee A (reference
382 number: 15/SC/0335), the South Central and Hampshire A Research Ethics Committee
383 (reference number: 17/SC/0522), and the London—Central Research Ethics Committee
384 (reference number: 17/LO/1210) in the UK; Hospital Clínico Universitario de Santiago
385 de Compostela, and Comité de Ética de la Investigación de Santiago-Lugo (reference
386 number: 2017/395) in Spain; the Medical Ethical Committee, University Medical Center
387 Utrecht (reference number: 17/563), and the Ethical Review Authority (reference number:
388 NL60910.041.17) in the Netherlands. All adult participants and the parents or guardians
389 of all infant participants provided written, informed consent.

390 **Nucleic acid isolation and whole-genome sequencing**

391 All RSV-positive samples were selected for whole-genome sequencing. Nucleic acid isolation,
392 library construction, and sequencing were performed in four different batches. To minimise
393 the risk of RNA degradation, nucleic acid was extracted locally from primary samples,
394 and the extractions were scheduled as close as practical to the time of sequencing.

395 Total nucleic acid extraction was carried out using the NucliSENS® easyMAG® system
396 (BioMérieux, Marcy-l'Étoile, France), following the manufacturer's instructions. 500 µL of
397 each sample was used to get 25 µL eluate in the first and fourth batches, and 35 µL in the
398 second and third batches.

399 Sequencing libraries were constructed using the methodology of targeted metagen-
400 omics²⁷, a modification of the veSEQ-HIV protocol⁴⁸. A 12-µL aliquot of each nucleic
401 acid sample was first concentrated to 3 µL with RNAClean XP magnetic beads (Beck-
402 man Coulter, California, United States). Dual-indexed libraries for Illumina sequencing
403 were then constructed using the SMARTer Stranded Total RNA-Seq Kit v2 - Pico In-
404 put Mammalian (Takara Bio USA, California, United States), where first-strand reverse
405 transcription was primed with tagged random hexamers and double-stranded cDNA was
406 synthesised with sets of i5 and i7 index primers, as previously described elsewhere⁴⁹.
407 These gave unique dual indexing (UDI) for the samples, thus minimising the risk of index
408 misassignment during sequencing. After 12 cycles of PCR amplification of the cDNA,
409 10 µL of each library was pooled and purified using AMPure XP (Beckman Coulter). A
410 750-ng aliquot was taken from the pool and captured using a predesigned SureSelect
411 RNA Target Enrichment multi-pathogen probe set (Agilent, California, United States).
412 This probe set (each 120 nucleotides long) targeted more than 100 pathogenic bacteria
413 and viruses, including both RSV-A and RSV-B⁵⁰. 16 cycles of PCR were performed for
414 post-capture amplification, and the final product was purified by AMPure XP.

415 Sequencing was performed on the Illumina MiSeq platform (Illumina, California, US)
416 with the MiSeq Reagent Kit v3 (600-cycle) for the first and third batches, generating
417 265-bp and 300-bp paired-end reads respectively. The second and fourth batches were
418 sequenced on the Illumina NovaSeq 6000 system with the NovaSeq 6000 SP Reagent Kit

419 v1.5 (300-cycle), generating 151-bp paired-end reads.

420 **Genome reconstruction**

421 The first six bases of read 1 and the first three bases of read 2 were clipped off to
422 remove random hexamer primers and the SMARTer adapter sequences, respectively. An
423 extra three bases at the 5' end of MiSeq-generated read 2 were also cut off as they
424 had reduced quality. Trimmomatic (v0.39)⁵¹ was then used to trimmed off adapter
425 sequences and low-quality bases with a Phred score below 20 (option: Adapters:2:10:7:1:true
426 LEADING:20 TRAILING:20 SLIDINGWINDOW:4:20 MINLEN:50). De novo assembly
427 of the trimmed reads was carried out using both IVA (v1.0.8)⁵² and SPAdes (v3.14.1)⁵³,
428 in each case selecting the contig sequences with a higher N50 for genome reconstruction
429 using shiver⁵⁴. Internally, BLASTN (v2.7.1+)⁵⁵ was used for read and contig classification,
430 MAFFT (v7.471)⁵⁶ was used for sequence alignment, and Bowtie 2 (v2.4.1)⁵⁷ was used
431 for read alignment (option: --very-sensitive-local). A minimum base quality of 35 and
432 mapping quality of 30 were required for a base or an alignment to be counted as mapped.
433 Mapped RSV reads were deduplicated with Picard MarkDuplicates (v2.18.14, <https://broadinstitute.github.io/picard/>). Pre-deduplicated per-position mapped read counts,
434 generated by shiver, were used for downstream within-host virus diversity analysis.

436 **Within-host virus diversity analysis**

437 Only samples generating more than 10,000 unique (i.e., deduplicated) RSV reads and
438 containing a single subgroup of RSV were included in within-host virus genetic diversity
439 analysis. We have previously shown that RSV viral load highly correlates with the number
440 of unique RSV reads generated by this sequencing method²⁷, consistent with high quality
441 RNA being recovered in a quantitative way. Ten thousand unique RSV reads correspond
442 to a viral load of approximately 2.4×10^6 copies/mL. Allele frequencies were calculated at
443 each genomic position, excluding those supported by fewer than 200 reads. The choice of
444 this cut-off was based on a predefined criterion that 90% of the included samples had at
445 least 80% of the genome fulfilling this cut-off (Supplementary Fig. 7). Cumulative minor

446 allele frequency (MAF) was defined as 1 minus major allele frequency, and polymorphic
 447 sites were those with a cumulative MAF of $\geq 3\%$. Mean cumulative MAF per sample was
 448 calculated as the sum of cumulative MAF at each genomic position divided by the total
 449 number of positions. Minor variants, or intrahost single nucleotide variants, were defined
 450 as variants with an allele frequency of $\geq 3\%$ and $< 50\%$.

Intrahost virus diversity was estimated as pairwise nucleotide diversity (π)⁵⁸. The proportion of pairwise nucleotide differences (D) at each genomic position was calculated as

$$D_i = \frac{A_i \times C_i + A_i \times G_i + A_i \times T_i + C_i \times G_i + C_i \times T_i + G_i \times T_i}{(N_i^2 - N_i)/2} \quad (1)$$

where A_i , C_i , G_i , and T_i represent the copy number of allele A, C, G, and T, respectively, and N_i is the total count of the four alleles (i.e., depth of coverage) at a given locus i , so $N_i = A_i + C_i + G_i + T_i$. Loci with a total count of less than 200 were excluded. Pairwise nucleotide diversity across a genome (π) was then calculated as

$$\pi = \sum_{i=1}^L \frac{D_i}{L} \quad (2)$$

451 where L is the number of genomic positions with a read depth of at least $200\times$.

Manhattan (L1-norm) distance was used to compare within-host diversity levels between samples, calculated as

$$d_i(\mathbf{p}, \mathbf{q}) = \sum_{k=1}^4 |\mathbf{p}_k - \mathbf{q}_k| \quad (3)$$

$$M = \sum_{i=1}^N d_i \times \frac{S}{N} \quad (4)$$

452 where d_i is the distance between two samples at a given locus i with vectors \mathbf{p} and \mathbf{q}
 453 containing relative frequencies of four possible alleles (i.e., A, C, G, and T), M is the
 454 Manhattan distance between the coding sequences of two samples, N is the number of
 455 coding sequence positions where both samples have the same consensus base and a read
 456 depth of at least $200\times$, and S is the total length of the coding sequence. To remove
 457 potential background noise in Manhattan distance calculations, allele frequencies of $< 3\%$
 458 were changed to 0, and those of $> 97\%$ were changed to 100%.

459 Nucleotide positions were numbered from the first base of the coding sequence of each
460 gene according to the NCBI reference sequences with the accession numbers of NC_038235
461 and NC_001781 for RSV-A and RSV-B, respectively. Amino acid positions were numbered
462 from the first methionine of each protein according to the same NCBI reference sequences.

463 **Phylogeny reconstruction**

464 Maximum likelihood phylogenies of consensus coding sequences, supported by at least
465 two unique (deduplicated) RSV reads, were estimated using RAxML (v8.2.12)⁵⁹ with
466 the general time-reversible nucleotide substitution model and gamma-distributed rate
467 heterogeneity. Bootstrapping with 1,000 replicates was used to assess the robustness of
468 tree topologies. Pairwise patristic distances were calculated from the maximum-likelihood
469 trees using the cophenetic function of the R package ape (v5.4-1)⁶⁰. Phylogenetic trees
470 were visualised using the R package ggtree (v2.2.4)⁶¹.

471 **Statistical analysis**

472 Continuous variables were summarised using mean, median, maximum, and minimum.
473 All comparisons of continuous variables between groups were conducted by two-tailed
474 Mann–Whitney U tests (two groups) or Kruskal–Wallis tests (three groups). Post hoc
475 application of the Benjamini–Hochberg procedure was used to control false discovery
476 rates for multiple testing. Chi-square tests with Yates’ continuity correction were used
477 for contingency analysis; Fisher’s exact tests were performed when the expected value
478 of a cell was less than 5. Logistic regression was employed to model a binary dependent
479 variable while adjusting for a covariate. Two-tailed Pearson correlation analysis was used
480 to evaluate the relationship between two variables. Temporal changes of a variable were
481 determined by ordinary least-squares linear regression. Two approaches were applied
482 to account for batch effects on the comparisons of diversity metrics: (i) including batch
483 as a regression covariate (e.g., regression of pairwise nucleotide diversity on sampling
484 country, sampling season, RSV subgroup, RSV read count, participant age group, disease
485 severity, and ‘batch’ as in Supplementary Table 2); and (ii) standardising the values within

486 each batch to z-scores, that is, to a mean of zero and a standard deviation of 1 (e.g.,
487 Mann–Whitney U test on z-score standardised pairwise nucleotide diversity as in Fig. 2).
488 Missing data were imputed using the `aregImpute` function, implemented in the R package
489 `Hmisc` (v4.5-0)⁶². All statistical analyses were performed using R (v4.0.2)⁶³. P values or
490 adjusted P values of less than 0.05 were considered to indicate statistical significance.

491 **Data availability**

492 The sequencing read data generated in this study have been deposited in the European
493 Nucleotide Archive under study accession PRJEB34042 ([https://www.ebi.ac.uk/ena/
494 data/view/PRJEB34042](https://www.ebi.ac.uk/ena/data/view/PRJEB34042)). The RSV genomic sequences generated in this study have
495 been deposited in GenBank under accession numbers LR699315, LR699726, LR699734,
496 LR699736–LR699744, and MZ515551–MZ516143. The RSV reference sequences used
497 in this study are available in GenBank under accession numbers NC_038235 ([https:
498 //www.ncbi.nlm.nih.gov/nuccore/NC_038235](https://www.ncbi.nlm.nih.gov/nuccore/NC_038235)) and NC_001781 ([https://www.ncbi.nlm.nih.
499 gov/nuccore/NC_001781](https://www.ncbi.nlm.nih.gov/nuccore/NC_001781)). The associated sample and de-identified clinical information
500 used in this study is provided in Supplementary Data 1.

501 **References**

- 502 1. Shi, T. *et al.* Global, regional, and national disease burden estimates of acute lower
503 respiratory infections due to respiratory syncytial virus in young children in 2015: a
504 systematic review and modelling study. *Lancet* **390**, 946–958 (2017).
- 505 2. Varga, S. M. & Braciale, T. J. The adaptive immune response to respiratory syncytial
506 virus. *Curr. Top. Microbiol. Immunol.* **372**, 155–171 (2013).
- 507 3. Falsey, A. R., Hennessey, P. A., Formica, M. A., Cox, C. & Walsh, E. E. Respiratory
508 syncytial virus infection in elderly and high-risk adults. *N. Engl. J. Med.* **352**, 1749–
509 1759 (2005).

- 510 4. American Academy of Pediatrics. Updated guidance for palivizumab prophylaxis
511 among infants and young children at increased risk of hospitalization for respiratory
512 syncytial virus infection. *Pediatrics* **134**, e620–e638 (2014).
- 513 5. Ruckwardt, T. J., Morabito, K. M. & Graham, B. S. Immunological lessons from
514 respiratory syncytial virus vaccine development. *Immunity* **51**, 429–442 (2019).
- 515 6. McLellan, J. S. Neutralizing epitopes on the respiratory syncytial virus fusion gly-
516 coprotein. *Curr. Opin. Virol.* **11**, 70–75 (2015).
- 517 7. Zhu, Q. *et al.* Prevalence and significance of substitutions in the fusion protein of
518 respiratory syncytial virus resulting in neutralization escape from antibody medi8897.
519 *J. Infect. Dis.* **218**, 572–580 (2018).
- 520 8. Simoes, E. A. F. *et al.* Suptavumab for the prevention of medically attended respiratory
521 syncytial virus infection in preterm infants. *Clin. Infect. Dis.* (2020).
- 522 9. Wu, H. *et al.* Development of motavizumab, an ultra-potent antibody for the prevention
523 of respiratory syncytial virus infection in the upper and lower respiratory tract. *J.*
524 *Mol. Biol.* **368**, 652–665 (2007).
- 525 10. Wu, S. J. *et al.* Characterization of the epitope for anti-human respiratory syn-
526 cytial virus f protein monoclonal antibody 101f using synthetic peptides and genetic
527 approaches. *J. Gen. Virol.* **88**, 2719–2723 (2007).
- 528 11. Zhao, X., Chen, F. P., Megaw, A. G. & Sullender, W. M. Variable resistance to
529 palivizumab in cotton rats by respiratory syncytial virus mutants. *J. Infect. Dis.* **190**,
530 1941–1946 (2004).
- 531 12. Zhu, Q. *et al.* Analysis of respiratory syncytial virus preclinical and clinical variants
532 resistant to neutralization by monoclonal antibodies palivizumab and/or motavizumab.
533 *J. Infect. Dis.* **203**, 674–682 (2011).

- 534 13. Zhu, Q. *et al.* Natural polymorphisms and resistance-associated mutations in the
535 fusion protein of respiratory syncytial virus (rsv): effects on rsv susceptibility to
536 palivizumab. *J. Infect. Dis.* **205**, 635–638 (2012).
- 537 14. Battles, M. B. & McLellan, J. S. Respiratory syncytial virus entry and how to block
538 it. *Nat. Rev. Microbiol.* **17**, 233–245 (2019).
- 539 15. Fedechkin, S. O., George, N. L., Wolff, J. T., Kauvar, L. M. & DuBois, R. M. Structures
540 of respiratory syncytial virus g antigen bound to broadly neutralizing antibodies. *Sci.*
541 *Immunol.* **3**, eaar3534 (2018).
- 542 16. Power, U. F. *et al.* Safety and immunogenicity of a novel recombinant subunit
543 respiratory syncytial virus vaccine (bbg2na) in healthy young adults. *J. Infect. Dis.*
544 **184**, 1456–1460 (2001).
- 545 17. Choi, Y. *et al.* Antibodies to the central conserved region of respiratory syncytial
546 virus (rsv) g protein block rsv g protein cx3c-cx3cr1 binding and cross-neutralize rsv a
547 and b strains. *Viral Immunol.* **25**, 193–203 (2012).
- 548 18. Lee, J., Klenow, L., Coyle, E. M., Golding, H. & Khurana, S. Protective antigenic
549 sites in respiratory syncytial virus g attachment protein outside the central conserved
550 and cysteine noose domains. *PLoS Pathog.* **14**, e1007262 (2018).
- 551 19. Eshaghi, A. *et al.* Genetic variability of human respiratory syncytial virus a strains
552 circulating in ontario: a novel genotype with a 72 nucleotide g gene duplication. *PLoS*
553 *One* **7**, e32807 (2012).
- 554 20. Peret, T. C., Hall, C. B., Schnabel, K. C., Golub, J. A. & Anderson, L. J. Circulation
555 patterns of genetically distinct group a and b strains of human respiratory syncytial
556 virus in a community. *J. Gen. Virol.* **79**, 2221–2229 (1998).
- 557 21. Grad, Y. H. *et al.* Within-host whole-genome deep sequencing and diversity analysis
558 of human respiratory syncytial virus infection reveals dynamics of genomic diversity
559 in the absence and presence of immune pressure. *J. Virol.* **88**, 7286–7293 (2014).

- 560 22. Lau, J. W. *et al.* Deep sequencing of rsv from an adult challenge study and from
561 naturally infected infants reveals heterogeneous diversification dynamics. *Virology*
562 **510**, 289–296 (2017).
- 563 23. Lythgoe, K. A. *et al.* Sars-cov-2 within-host diversity and transmission. *Science* **372**,
564 eabg0821 (2021).
- 565 24. McCrone, J. T. *et al.* Stochastic processes constrain the within and between host
566 evolution of influenza virus. *Elife* **7**, e35962 (2018).
- 567 25. Tabor, D. E. *et al.* Global molecular epidemiology of respiratory syncytial virus from
568 the 2017-2018 inform-rsv study. *J Clin Microbiol* **59**, e01828–20 (2020).
- 569 26. McCrone, J. T. & Luring, A. S. Measurements of intrahost viral diversity are
570 extremely sensitive to systematic errors in variant calling. *J. Virol.* **90**, 6884–6895
571 (2016).
- 572 27. Lin, G. L. *et al.* Simultaneous viral whole-genome sequencing and differential expression
573 profiling in respiratory syncytial virus infection of infants. *J. Infect. Dis.* **222**, S666–
574 S671 (2020).
- 575 28. Xue, K. S., Moncla, L. H., Bedford, T. & Bloom, J. D. Within-host evolution of
576 human influenza virus. *Trends. Microbiol.* **26**, 781–793 (2018).
- 577 29. Duffy, S., Shackelton, L. A. & Holmes, E. C. Rates of evolutionary change in viruses:
578 patterns and determinants. *Nat. Rev. Genet.* **9**, 267–276 (2008).
- 579 30. Sanjuan, R. & Domingo-Calap, P. Mechanisms of viral mutation. *Cell. Mol. Life Sci.*
580 **73**, 4433–4448 (2016).
- 581 31. Valesano, A. L. *et al.* Influenza b viruses exhibit lower within-host diversity than
582 influenza a viruses in human hosts. *J. Virol.* **94**, e01710–19 (2020).
- 583 32. Cudini, J. *et al.* Human cytomegalovirus haplotype reconstruction reveals high diversity
584 due to superinfection and evidence of within-host recombination. *Proc. Natl. Acad.*
585 *Sci. U. S. A.* **116**, 5693–5698 (2019).

- 586 33. Papenburg, J. *et al.* Molecular evolution of respiratory syncytial virus fusion gene,
587 canada, 2006-2010. *Emerg. Infect. Dis.* **18**, 120–124 (2012).
- 588 34. Griffin, M. P. *et al.* Single-dose nirsevimab for prevention of rsv in preterm infants. *N.*
589 *Engl. J. Med.* **383**, 415–425 (2020).
- 590 35. Tan, L. *et al.* The comparative genomics of human respiratory syncytial virus subgroups
591 a and b: genetic variability and molecular evolutionary dynamics. *J Virol* **87**, 8213–
592 8226 (2013).
- 593 36. Schobel, S. A. *et al.* Respiratory syncytial virus whole-genome sequencing identifies
594 convergent evolution of sequence duplication in the c-terminus of the g gene. *Sci Rep*
595 **6**, 26311 (2016).
- 596 37. Matheson, J. W. *et al.* Distinct patterns of evolution between respiratory syncytial
597 virus subgroups a and b from new zealand isolates collected over thirty-seven years. *J*
598 *Med Virol* **78**, 1354–1364 (2006).
- 599 38. Smith, G. *et al.* Respiratory syncytial virus fusion glycoprotein expressed in insect
600 cells form protein nanoparticles that induce protective immunity in cotton rats. *PLoS*
601 *One* **7**, e50852 (2012).
- 602 39. Pierantoni, A. *et al.* Mucosal delivery of a vectored rsv vaccine is safe and elicits
603 protective immunity in rodents and nonhuman primates. *Mol. Ther. Methods Clin.*
604 *Dev.* **2**, 15018 (2015).
- 605 40. Crank, M. C. *et al.* A proof of concept for structure-based vaccine design targeting
606 rsv in humans. *Science* **365**, 505–509 (2019).
- 607 41. Bin, L. *et al.* Emergence of new antigenic epitopes in the glycoproteins of human
608 respiratory syncytial virus collected from a us surveillance study, 2015-17. *Sci. Rep.* **9**,
609 3898 (2019).
- 610 42. Poon, L. L. *et al.* Quantifying influenza virus diversity and transmission in humans.
611 *Nat. Genet.* **48**, 195–200 (2016).

- 612 43. Githinji, G. *et al.* Assessing the utility of minority variant composition in elucidating
613 rsv transmission pathways. *bioRxiv* 411512 (2018). URL [https://www.biorxiv.org/
614 content/10.1101/411512v1.full](https://www.biorxiv.org/content/10.1101/411512v1.full).
- 615 44. Wildenbeest, J. G. *et al.* Respiratory syncytial virus consortium in europe (resceu)
616 birth cohort study: defining the burden of infant respiratory syncytial virus disease in
617 europe. *J. Infect. Dis.* **222**, S606–S612 (2020).
- 618 45. Jefferies, K. *et al.* Presumed risk factors and biomarkers for severe respiratory syncytial
619 virus disease and related sequelae: protocol for an observational multicenter, case-
620 control study from the respiratory syncytial virus consortium in europe (resceu). *J.*
621 *Infect. Dis.* **222**, S658–S665 (2020).
- 622 46. Korsten, K. *et al.* Burden of respiratory syncytial virus infection in community-dwelling
623 older adults in europe (resceu): an international prospective cohort study. *Eur. Respir.*
624 *J.* **57**, 2002688 (2021).
- 625 47. Justicia-Grande, A. J. *et al.* Development and validation of a new clinical scale for
626 infants with acute respiratory infection: the resvinet scale. *PLoS One* **11**, e0157665
627 (2016).
- 628 48. Bonsall, D. *et al.* A comprehensive genomics solution for hiv surveillance and clinical
629 monitoring in low income settings. *J. Clin. Microbiol.* **58**, e00382–20 (2020).
- 630 49. Faircloth, B. C. & Glenn, T. C. Not all sequence tags are created equal: designing and
631 validating sequence identification tags robust to indels. *PLoS One* **7**, e42543 (2012).
- 632 50. Goh, C. *et al.* Targeted metagenomic sequencing enhances the identification of
633 pathogens associated with acute infection. *bioRxiv* 716902 (2019). URL [https:
634 //www.biorxiv.org/content/10.1101/716902v1.full](https://www.biorxiv.org/content/10.1101/716902v1.full).
- 635 51. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for illumina
636 sequence data. *Bioinformatics* **30**, 2114–2120 (2014).

- 637 52. Hunt, M. *et al.* Iva: accurate de novo assembly of rna virus genomes. *Bioinformatics*
638 **31**, 2374–2376 (2015).
- 639 53. Bankevich, A. *et al.* Spades: a new genome assembly algorithm and its applications
640 to single-cell sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).
- 641 54. Wymant, C. *et al.* Easy and accurate reconstruction of whole hiv genomes from
642 short-read sequence data with shiver. *Virus Evol.* **4**, vey007 (2018).
- 643 55. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
644 alignment search tool. *J Mol Biol* **215**, 403–10 (1990).
- 645 56. Katoh, K. & Standley, D. M. Mafft multiple sequence alignment software version 7:
646 improvements in performance and usability. *Mol Biol Evol* **30**, 772–80 (2013).
- 647 57. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with bowtie 2. *Nat.*
648 *Methods* **9**, 357–359 (2012).
- 649 58. Nelson, C. W. & Hughes, A. L. Within-host nucleotide diversity of virus populations:
650 insights from next-generation sequencing. *Infect. Genet. Evol.* **30**, 1–7 (2015).
- 651 59. Stamatakis, A. Raxml version 8: a tool for phylogenetic analysis and post-analysis of
652 large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
- 653 60. Paradis, E. & Schliep, K. ape 5.0: an environment for modern phylogenetics and
654 evolutionary analyses in r. *Bioinformatics* **35**, 526–528 (2019).
- 655 61. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T. ggtree: an r package for
656 visualization and annotation of phylogenetic trees with their covariates and other
657 associated data. *Methods Ecol. Evol.* **8**, 28–36 (2017).
- 658 62. Harrell Jr, F. E., with contributions from Charles Dupont & many others. Hmisc:
659 Harrell miscellaneous (2021). URL <https://CRAN.R-project.org/package=Hmisc>. R
660 package version 4.5-0.
- 661 63. R Core Team. *R: A Language and Environment for Statistical Computing* (R Founda-
662 tion for Statistical Computing, Vienna, 2018).

663 **Acknowledgements**

664 This work was supported by the National Institute for Health Research (NIHR) Oxford Bio-
665 medical Research Centre, the NIHR Thames Valley and South Midlands Clinical Research
666 Network, the British Research Council, and the REspiratory Syncytial virus Consortium in
667 EUrope (RESCEU) project. RESCEU has received funding from the Innovative Medicines
668 Initiative 2 Joint Undertaking (grant number 116019). This Joint Undertaking receives
669 support from the European Union Horizon 2020 Research and Innovation Program and
670 European Federation of Pharmaceutical Industries and Associations.

671 **Author contributions**

672 G.-L. L., T. G., and A. J. P. conceived and designed the work. G.-L. L., S. B. D., M. D. S.,
673 D. Ö., J. A., C. B., L. B., P. O., F. M.-T., H. N., and A. J. P. conducted and supervised
674 the clinical studies. M. A. A. designed the probe set that was used for capture. M. d. C.,
675 D. B., and R. B. designed the sequencing protocol. G.-L. L., A. B., G. M.-C., E. M.-G.,
676 and M. d. C. performed the experiments. G.-L. L., T. G., D. O’C., and A. J. P. analysed
677 and interpreted the data. G.-L. L. drafted the manuscript, and T. G., D. O’C., and A.
678 J. P. substantively revised it. T. G. and A. J. P. supervised the work. All authors have
679 approved the submitted version and agreed to submit the manuscript.

680 **Competing interests**

681 S. B. D has been an investigator for clinical trials of vaccines and antimicrobials for
682 pharmaceutical companies including AstraZeneca, Merck, and Janssen, and sits on an RSV
683 advisory board for Sanofi Pasteur. D. Ö. and J. A. are employees of Janssen Pharmaceutica
684 NV. F. M.-T. has received honoraria from GSK, Pfizer Inc., Sanofi Pasteur, MSD, Seqirus,
685 and Janssen for taking part in advisory boards and expert meetings and for acting as a
686 speaker in congresses outside the scope of the submitted work. F. M.-T. has also acted as
687 principal investigator in randomised controlled trials of the above-mentioned companies as

688 well as Ablynx, Regeneron, Roche, Abbott, Novavax, and MedImmune, with honoraria
689 paid to his institution. F. M.-T. receives support for his research activities from the
690 Instituto de Salud Carlos III (Proyecto de Investigación en Salud, Acción Estratégica en
691 Salud): Fondo de Investigación Sanitaria (FIS;PI1601569/PI1901090) del plan nacional de
692 I+D+I and ‘fondos FEDER’. A. J. P. is a National Institute for Health Research (NIHR)
693 Senior Investigator with funding from the British Research Council. The remaining authors
694 declare no competing interests. The views expressed in this article are those of the authors
695 and may not be understood or quoted as being made on behalf of or reflecting the position
696 of the organizations with which the authors are employed/affiliated.

697 **RESCEU Investigators**

698 Harry Campbell¹⁴, Steve Cunningham¹⁴, Debby Bogaert^{9,16}, Philippe Beutels¹⁷, Joanne
699 Wildenbeest⁹, Elizabeth Clutterbuck¹, Joseph McGinley¹, Ryan Thwaites¹¹, Dexter
700 Wiseman¹¹, Alberto Gómez-Carballa¹³, Carmen Rodriguez-Tenreiro¹³, Irene Rivero-Calle¹³,
701 Ana Dacosta-Urbiet¹³, Terho Heikkinen¹⁸, Adam Meijer¹⁹, Thea Kølsten Fischer²⁰, Maarten
702 van den Berge²¹, Carlo Giaquinto²², Michael Abram²³, Philip Dormitzer²⁴, Sonia Stoszek²⁵,
703 Scott Gallichan²⁶, Brian Rosen²⁷, Eva Molero²⁸, Nuria Machin²⁸, Martina Spadetto²⁸.

704 ¹⁶Queen’s Medical Research Institute, University of Edinburgh, Edinburgh, UK. ¹⁷Centre
705 for Health Economics Research and Modelling Infectious Diseases, Vaccine and Infectious
706 Disease Institute, University of Antwerp, Antwerp, Belgium. ¹⁸Department of Pediatrics,
707 University of Turku, Turku University Hospital, Turku, Finland. ¹⁹National Institute for
708 Public Health and the Environment, Bilthoven, Netherlands. ²⁰Statens Serum Institut,
709 Copenhagen, Denmark. ²¹Department of Pulmonary Diseases, University of Groningen,
710 University Medical Center Groningen, Groningen, Netherlands. ²²PENTA Foundation,
711 Padua, Italy. ²³AstraZeneca, Gaithersburg, MD, US. ²⁴Pfizer, Pearl River, NY, US.
712 ²⁵GlaxoSmithKline, Potomac, MD, US. ²⁶Sanofi Pasteur, Toronto, Ontario, Canada.
713 ²⁷Novavax, Potomac, MD, US. ²⁸Team-It Research, Barcelona, Spain.

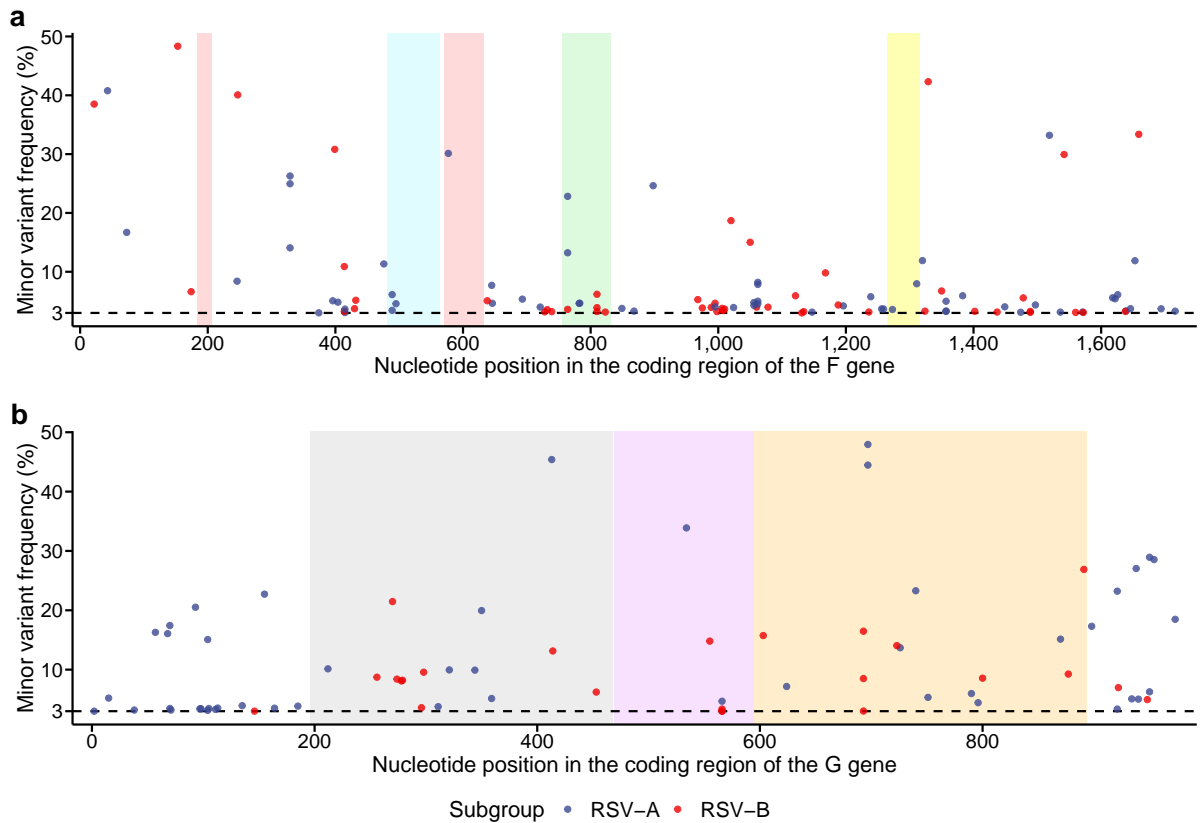


Fig. 1 Minor variants in the coding region of the **F** and **G** genes among 175 **RSV-A** and 144 **RSV-B** samples. **a** **F** gene. Shaded regions represent known antigenic sites (neutralising epitopes in particular): red, prefusion-specific antigenic site \emptyset (target for nirsevimab); green, site II (target for palivizumab and motavizumab); yellow, site IV (target for 101F); and blue, prefusion-specific site V (target for suptavumab). **b** **G** gene. The purple region represents the conserved central domain (target for 3D3 and 2D10), flanked by highly variable mucin-like regions I (grey) and II (orange). Nirsevimab, palivizumab, motavizumab, 101F, suptavumab, 3D3, and 2D10 are RSV-specific monoclonal antibodies. Each dot denotes a minor variant, coloured by subgroup. Black dashed line represents minor allele frequency of 3%, used to define a minor variant. Positions are numbered from the first base of the coding sequence of each gene according to the NCBI reference sequence (accession number NC_038235).

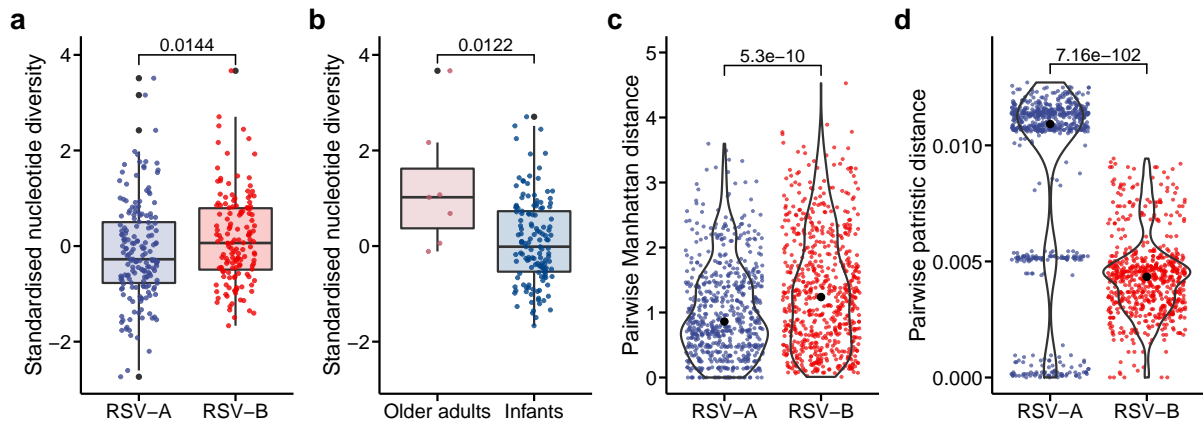


Fig. 2 Z-score standardised pairwise nucleotide diversity and pairwise genetic distances. **a** Comparison of standardised pairwise nucleotide diversity between 175 RSV-A and 144 RSV-B samples. **b** Comparison of standardised pairwise nucleotide diversity of RSV-B between 7 adult samples and 137 infant samples. RSV-A isolates were excluded from this comparison because only one adult had RSV-A infection. **c** Comparison of pairwise Manhattan distances. **d** Comparison of pairwise patristic distances. Only pairwise distances between samples from the second sequencing batch, the same country, the same season, and different participants were included in **c** and **d** (650 RSV-A pairs and 656 RSV-B pairs). Each dot represents an individual sample in **a** and **b**, and a sample pair in **c** and **d**. Two-tailed Mann–Whitney U tests were used to evaluate the significance of the differences. P values are shown above the plots. For **a** and **b**, the centre line of each box denotes the median; box limits, the first and third quartiles; whiskers, the highest and lowest values within 1.5 times the interquartile range from the box limits; and outlying points, outliers. For **c** and **d**, the violin plots summarise the distribution of the data, and the black dots denote the median value of each group.

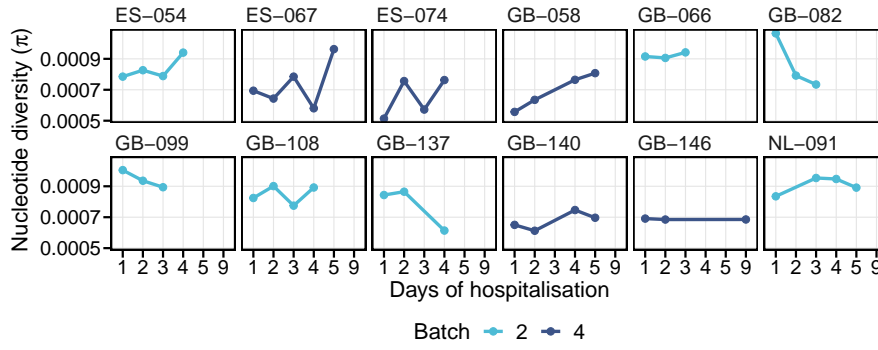


Fig. 3 Temporal change of pairwise nucleotide diversity. Pairwise nucleotide diversity of serial samples collected at more than two time points and sequenced in the same batch are shown here. Three participants whose samples were sequenced in different batches and 19 participants who had only two samples collected are not shown. Each panel is labelled with the participant ID.

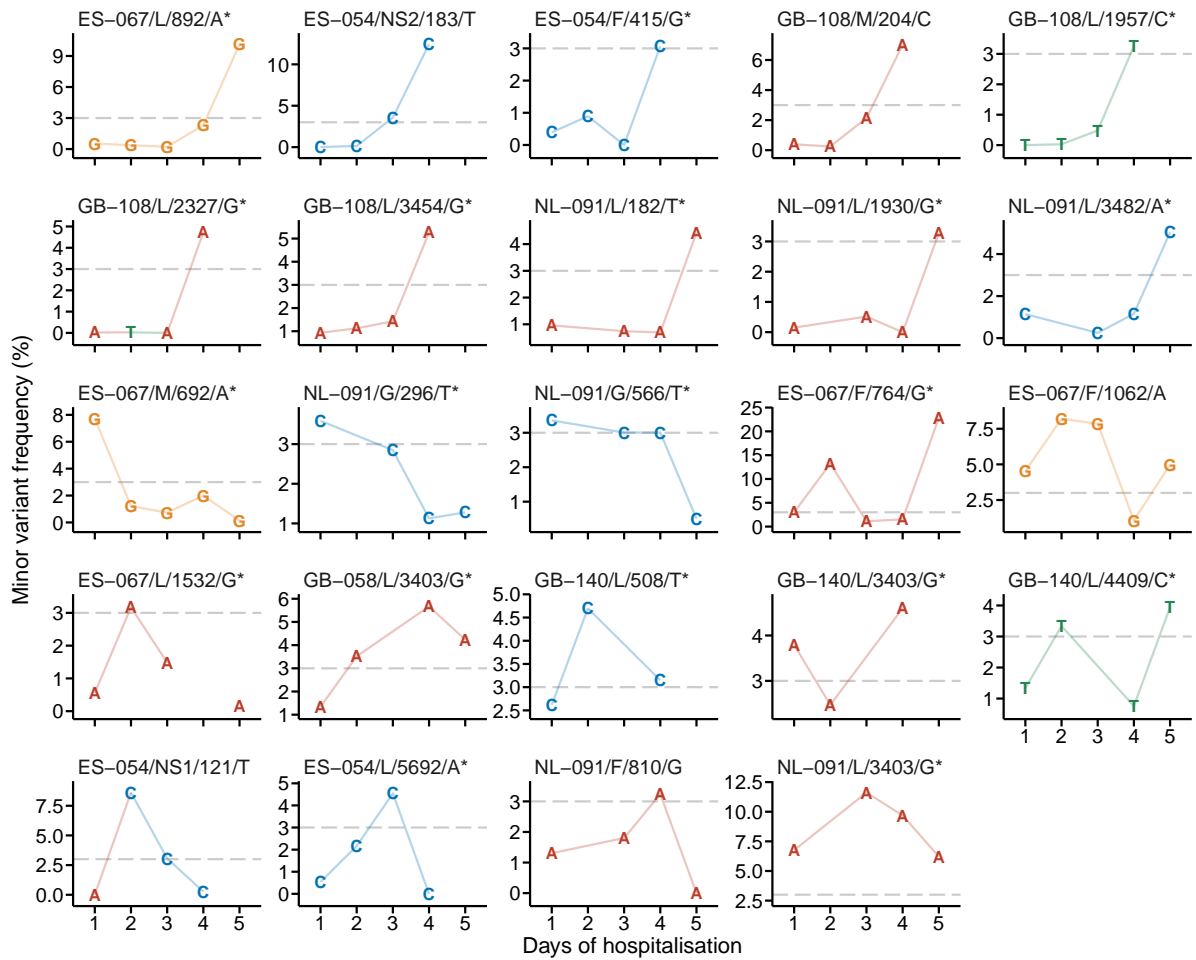


Fig. 4 Temporal change in minor alleles. Minor alleles and allele frequencies are shown at polymorphic sites within the coding sequence of the serial samples, where minor alleles were detected at ≥ 3 time points. The grey dashed lines represent the 3% threshold, which defines a minor variant. Panels are labelled with the participant ID, followed by the gene name, the nucleotide position, and the consensus base. Asterisks denote nonsynonymous substitutions. Letters in the plots denote minor allele bases. Panels are ordered by the trend of the change: increased, decreased, and fluctuated. Positions are numbered from the first base of the coding sequence of each gene according to the NCBI reference sequences with the accession numbers of NC_038235 and NC_001781 for RSV-A and RSV-B, respectively.

Table 1 Characteristics of RSV samples by subgroup.

	RSV-A (N = 175)	RSV-B (N = 144)	P Value ^a
Host number			0.12 ^b
Infants	141	115	
Older adults	1	7	
Host age, median (range)			
Infants (month) ^c	4.5 (0.5–11.6)	4.3 (0.2–11.7)	0.72
Older adults (year)	69	75 (72–78)	0.19
Sample source			0.45
United Kingdom	74	64	
Netherlands	58	53	
Spain	43	27	
Sampling season			2.9×10^{-5}
2017–18	14	33	
2018–19	65	63	
2019–20	96	48	
Days between symptom onset and sample collection, median (range) ^d	4 (1–11)	4 (1–9)	0.11
Number of unique RSV read pairs (\log_{10}), median (range)	4.6 (4.0–5.8)	4.7 (4.0–5.9)	0.22
Batch 1	4.9 (4.1–5.5)	5.3 (4.4–5.6)	0.50
Batch 2	4.6 (4.0–5.6)	4.6 (4.0–5.9)	0.98
Batch 3	4.4 (4.0–4.8)	4.5 (4.0–5.5)	0.18
Batch 4	4.7 (4.0–5.8)	4.9 (4.0–5.6)	0.12
Minimum genome coverage (%)	99.9	100	0.37
Average depth of coverage, median (range)	3,372 (696–7,897)	3,650 (525–7,930)	0.41
Batch 1	2,940 (696–6,823)	4,975 (1,295–7,601)	0.63
Batch 2	3,561 (1,092–7,452)	3,469 (1,091–7,930)	0.68
Batch 3	2,045 (803–3,224)	2,258 (525–7,157)	0.37
Batch 4	3,736 (847–7,897)	4,505 (719–7,798)	0.23

^a Unless otherwise specified, chi-square tests with Yates' continuity correction or Fisher's exact tests were used for contingency analysis, and two-tailed Mann–Whitney U tests were used to compare numeric variables between subgroups.

^b Logistic regression was used to adjust for sampling season. Samples were collected from older adults only in 2017–18 and 2018–19 RSV seasons, when RSV-B was the predominant circulating subgroup.

^c One infant with RSV-B infection had missing information on age.

^d Six infants with RSV-A infection and five infants with RSV-B infection had missing information on date of symptom onset.

Table 2 Characteristics of minor variants within the antigenic sites of the fusion protein.

Nucleotide position ^a	Codon change	Amino acid change ^b	Antigenic site	Subgroup/Country/Season/ Minor allele frequency (%) ^c
489	GAA:GAt	E163D	V	A/GB/2018–19/3.4 A/GB/2018–19/6.1
495	AAC:AAt	N165	V	A/GB/2018–19/4.6
577	CCA:tCA ^d	P193S	∅	A/GB/2018–19/30.1
764	AGT:AaT	S255N	II	A/ES/2019–20/13.2 ^e A/ES/2019–20/22.8 ^e B/GB/2018–19/3.6
782	ATC:Act	I261T	II	A/ES/2018–19/4.6 ^f
783	ATC:Act	I261T	II	A/ES/2018–19/4.7 ^f
810	CAG:CAa	Q270	II	B/NL/2017–18/3.2 B/NL/2017–18/6.2 ^g B/NL/2018–19/3.9
823	TCA:cCA	S275P	II	B/GB/2018–19/3.1
1273	TCA:cCA	S425P	IV	A/GB/2019–20/3.6
1311	AAC:AAt	N437	IV	A/NL/2019–20/8.0

^a Positions are numbered from the first base of the coding sequence of the F gene according to the NCBI reference sequence (accession number NC_038235).

^b Positions are numbered from the first methionine of the fusion protein according to the NCBI reference sequence (accession number NC_038235).

^c GB denotes the United Kingdom; ES, Spain; and NL, the Netherlands.

^d 55.7% (98/176) of the RSV-A samples had a consensus base of T, and all RSV-B samples had a consensus base of T at this position.

^e These two variants were found in samples collected from the same participant on day 2 (13.2%) and day 5 (22.8%) of hospitalisation, respectively. Samples collected from this participant on other days (days 1, 3, and 4) did not have variants with a frequency of $\geq 3\%$ at this position.

^f These two were co-occurring mutations, identified in the same minor variant.

^g Except for this variant, which was in a sample from an adult participant, other minor variants were identified in infant samples.