

**1 Viral genetic determinants of prolonged respiratory syncytial virus
2 infection among infants in a healthy term birth cohort.**

3 Dylan Lawless, PhD^{1*}, Christopher G. McKennan, PhD², Suman Das, PhD⁴, Thomas Junier, PhD³,
4 Zhi Ming Xu, MSc¹, Larry J Anderson, MD⁵, Tebeb Gebretsadik, MPH⁶, Meghan Shilts, MHS, MS⁷,
5 Emma Larkin, PhD⁸, Christian Rosas-Salazar, MD, MPH⁸, James D. Chappell, MD⁹, Jacques Fellay,
6 MD, PhD^{1,3,10}, Tina V. Hartert, MD, MPH^{7,9*}

***For correspondence:**

dylan.lawless@epfl.ch (DL);
tina.hartert@vumc.org (TVH)

7¹Global Health Institute, School of Life Sciences, École Polytechnique Fédérale de Lausanne,
8Lausanne, Switzerland, ²Department of Statistics, University of Pittsburgh, Pittsburgh, Pennsylvania,
9United States of America, ³Swiss Institute of Bioinformatics, Vital-IT Group, Switzerland, ⁴Division
10of Infectious Diseases, Department of Medicine, Vanderbilt University Medical Center, Nashville,
11Tennessee, United States of America, ⁵Department of Pediatrics, Emory University School of
12Medicine, Atlanta, Georgia, United States of America, ⁶Department of Biostatistics, Vanderbilt
13University Medical Center, Nashville, Tennessee, United States of America, ⁷Department of Medicine,
14Vanderbilt University Medical Center, Nashville, Tennessee, United States of America, ⁸Division of
15Allergy, Immunology, and Pulmonary and Critical Care Medicine, Department of Medicine, Vanderbilt
16University Medical Center, Nashville, Tennessee, United States of America, ⁹Department of Pediatrics,
17Vanderbilt University Medical Center, Nashville, Tennessee, United States of America, ¹⁰Precision
18Medicine Unit, Lausanne University Hospital and University of Lausanne, Lausanne, Switzerland

19 Abbreviations

20ALT (alternative); CI (confidence interval); GWAS (genome-wide association study); G (glycoprotein);
21H (hemagglutinin); HN (hemagglutinin-neuraminidase); IFN (interferon); IQR (interquartile range);
22INSPIRE (The INFant Susceptibility to Pulmonary Infections and Asthma Following RSV Exposure); LD
23(linkage disequilibrium); LRTI (lower respiratory tract infection); MAF (minor allele frequency); MSA
24(multiple sequence alignment); OR (odds ratio); PCR (polymerase chain reaction); PCA (Principal
25component analysis); REF (reference); RT (reverse transcription); SVD (singular value decomposition);
26SNP (single nucleotide polymorphism); VE (variance explained); MSA (multiple sequence alignment);
27RSV (respiratory syncytial virus).

28 Notice of Prior Presentation

29The results of the host genome wide association study analyses included in this manuscript were
30presented during the European Society of Human Genetics Conference in June 2020 in Berlin,
31Germany, which was held remotely (*Lawless et al., 2020*).

32 Ethics Statement for Human Subjects Research

33The protocol and informed consent documents were approved by the Institutional Review Board at
34Vanderbilt University Medical Center (#111299). One parent of each participant in the cohort study
35provided written informed consent for participation in this study. The informed consent document
36explained study procedures and use of data and biospecimens for future studies, including genetic
37studies.

³⁸ **Competing interests**

³⁹ All authors have completed a conflict of interest form (COI). There were no COI. Funding was
⁴⁰ supplied from National Institutes of Health and Swiss National Science Foundation. U19 AI 095227,
⁴¹ UG3/UH3 OD023282, UL1 TR002243, SNSF IZSEZ0_191968 (TVH), SNSF 310030L_197721 (JF), X01
⁴² HLG244 RS&G (EL).

43 **Abstract**

44 **Background:** Respiratory syncytial virus (RSV) is primarily associated with acute respiratory infection.
45 However, many RNA viruses can establish prolonged or persistent infection in some infected
46 individuals.

47 **Objectives:** To determine the impact of host genetics on first RSV infection and identify viral
48 genetic variants associated with prolonged infection.

49 **Methods:** In a population-based cohort study of healthy term infants, RSV infection was deter-
50 mined by biweekly surveillance for RSV and 1-year RSV serology. First, we tested the dependence of
51 first year RSV infection risk on the genotype at single nucleotide polymorphisms (SNPs) previously
52 shown to alter infant RSV lower respiratory tract infection or childhood asthma risk. Second, we
53 used RSV whole-genome sequencing to determine the relationship between viral amino acids (geno-
54 types) and prolonged infant RSV infection. Analyses were adjusted for viral and human population
55 structure and host features that alter infection risk.

56 **Results:** We found no evidence of host genetic risk for RSV infection. We identified two po-
57 tentially causal variants, p.E123K/D and p.P218T/S/L, in the RSV G protein that were associated
58 with prolonged infection after a Bonferroni correction for multiple testing. These variants were
59 associated exclusively with upper respiratory tract infection, and on average, milder clinical infection
60 compared with other circulating variants.

61 **Conclusions:** We identified a novel RSV viral variant associated with prolonged infection in
62 healthy infants and no evidence supporting host genetic susceptibility to RSV infection during
63 infancy. As the capacity of RSV for chronicity and its viral reservoir are not defined, these results are
64 important to understanding viral and host genetic determinants of chronic respiratory morbidity
65 due to early-life RSV infection along with sustained RSV endemicity.

66 **Introduction**

67 Human orthopneumovirus, formerly known (and frequently still referred to) as Respiratory syncytial
68 virus (RSV), results in significant global morbidity and mortality (*Hall et al., 2009*). By the age of two
69 to three years, nearly all children are infected with RSV at least once (*Glezen et al., 1986*). RSV is
70 a seasonal mucosal pathogen that primarily infects upper and lower respiratory tract epithelium,
71 although it has been recovered from non-airway sources (*Bokun et al., 2019; Cubie et al., 1997;*
72 *Nadal et al., 1990; O'Donnell et al., 1998; Rezaee et al., 2011; Rohwedder et al., 1998*). While RSV
73 is mainly associated with acute respiratory infection, many RNA viruses can establish prolonged or
74 persistent infection in some infected individuals (*Randall and Griffin, 2017*). Prolonged shedding of
75 RSV, especially in young infants and following first infection, has been demonstrated, with longer
76 average duration of viral shedding when polymerase chain reaction (PCR) is used to detect RSV
77 (*Munywoki et al., 2015*). While younger age and first infection are associated with protracted
78 infection (*Bagga et al., 2018; Glezen et al., 1986*), it is not known whether specific viral factors
79 contribute to prolonged RSV infection in infants. This is important, as prolonged infection may
80 contribute to enhanced transmission and developmental changes to the early life airway epithelium.
81 Further, the reservoir of RSV infection is not understood, and it is possible that some RSV strains
82 sustain a low levels of ongoing viral circulation in the community until seasonal or other influences
83 favor epidemic spread (*Okiro et al., 2010*).

84 The objectives of this study were therefore to determine if there exist host genetic risk alleles for

85 RSV infection and to identify viral genetic variation associated with prolonged infection. These moti-
86 vating questions are of fundamental interest in understanding viral and host genetic contributions
87 that may underlie the development of chronic respiratory morbidity due to RSV, including asthma.

88 **Methods**

89 **Study population**

90 The protocol and informed consent documents were approved by the Institutional Review Board at
91 Vanderbilt University Medical Center (#111299). One parent of each participant in the cohort study
92 provided written informed consent for participation in this study. The informed consent document
93 explained study procedures and use of data and biospecimens for future studies, including genetic
94 studies.

95 The study population is a longitudinal birth cohort, the INfant Susceptibility to Pulmonary
96 Infections and Asthma Following RSV Exposure (INSPIRE), specifically designed to capture the first
97 RSV infection in term healthy infants. Additional details of this birth cohort have been previously
98 published (*Larkin et al., 2015a*). Briefly, the cohort included 1949 term (≥ 37 weeks gestation),
99 non-low birth weight (≥ 2250 g, 5 lbs), otherwise healthy infants from a population-representative
100 sample of pediatric practices located in rural, suburban, and urban regions of the south-eastern
101 US during 2012-2014. Infants were born June through December; per study design, they were 6
102 months of age or less entering their first RSV season.

103 **Biweekly surveillance of RSV infection**

104 Infant (i.e., first year of life) RSV infection was ascertained through passive and active biweekly
105 surveillance during each infant's first RSV season and RSV serology (Table 1). If an infant met
106 pre-specified criteria for an acute respiratory infection, we conducted an in-person respiratory
107 illness visit at which time we administered a parental questionnaire, performed a physical exam,
108 collected a nasal wash, and completed a structured medical chart review for infants seen during an
109 unscheduled visit. RSV RNA in nasal samples was detected by reverse-transcription quantitative
110 PCR (*Larkin et al., 2015b*). We a priori defined the clinical entity of "prolonged" infection during
111 infancy as repeatedly meeting pre-specified criteria for an acute respiratory infection accompanied
112 by repeatedly positive RSV PCR separated by 15 or more days (Figure S1) (*Okiro et al., 2010*).

113 **Descriptive analyses**

114 Descriptive analyses of the cohort were conducted using R 4.0.5. Pearson or Wilcoxon tests were
115 used for comparing infants with and without prolonged RSV infection. The main descriptive features
116 are provided in Table 1.

117 **Host DNA collection and genotyping**

118 One-year blood samples were selected based on availability of DNA among a subset of children
119 with RSV infection and a random group of those without infection, and were genotyped with the
120 Multi-Ethnic Global Array microarray (Illumina, CA, United States) at the University of Washington
121 DNA Sequencing and Gene Analysis Center (Seattle, WA, United States).

122 **Host genetic analyses of RSV infection in infancy**

123 To determine whether host genetic factors associate with infant RSV infection risk, we examined
124 single nucleotide polymorphisms (SNPs) previously shown to alter infant RSV lower respiratory tract
125 infection (LRTI) or childhood asthma risk (*Janssen et al., 2007; Pasanen et al., 2017; Pividori et al.,*
126 *2019*). We also conducted a host GWAS to identify common variants associated with infant RSV
127 infection, and examined narrow sense heritability to test for small cumulative effects. The GWAS
128 was performed on 621 children with available DNA for the association between host genotype
129 and RSV infection during infancy. Due to sample size constraints, we restricted our sub-analysis
130 to the 54 host SNPs previously associated with RSV lower respiratory tract infection or childhood
131 asthma (*Janssen et al., 2007; Pasanen et al., 2017; Pividori et al., 2019*). We additionally evaluated
132 the accumulation of small genetic effects that would go undetected in a GWAS by estimating the
133 narrow sense heritability of RSV infection.

134 For GWAS analyses, the initial round of data quality control was performed on individual popu-
135 lations (self-reported as White, Black, and Hispanic) using PLINK version 1.9 (*Purcell et al., 2007*).
136 Subjects with a missing genotype call rate above 5% were removed. The SNP minor allele frequency
137 (MAF) threshold was set at > 0.01, 0.03, and 0.08 for White, Black, and Hispanic, respectively (*Yang*
138 *et al., 2011*).

139 The groups were merged for a total of 1,086,830 variants and a genotyping rate of 0.78. Subject
140 independence was assessed using KING (<https://people.virginia.edu/~wc9c/KING/>) to prevent spuri-
141 ous associations. However, no probable relatives or duplicates were detected based on pairwise
142 identify-by-state. We compared the genetic ancestry in cases to self-reported ethnicity to check
143 for mislabelling. Reported and estimated sex was also examined for discrepancy. A second round
144 of quality control on the combined dataset was conducted, which removed 74 samples due to
145 genotype missingness and 399,991 variants with a genotyping rate < 0.1. Variants were checked for
146 departure from Hardy-Weinberg equilibrium (HWE) ($P < 1e^{-6}$) to uncover features of selection, pop-
147 ulation admixture, cryptic relatedness, or genotyping error. This was only performed on controls to
148 prevent removal of genuine genetic associations that can be associated with this measurement;
149 6,024 variants were removed. No variants had a MAF MAF < 0.01 after merging. SNP positions and
150 identifiers were compared and updated according to dbNSFP4.0a (hg19) with 289 variants removed
151 due to a missing coordinate and SNPs identifier (*Liu et al., 2016*). This resulted in an analysis-ready
152 dataset of 680,526 variants from 621 children (509 with and 112 without RSV infection in infancy),
153 yielding a total genotyping rate of 0.98. No genomic inflation was evident with an estimated lambda
154 (based on median chi-squared test) equal to 1. We then used genome-wide complex trait analysis
155 (GCTA) software (<https://cnsgenomics.com/software/gcta/>) to calculate the genetic relationship
156 matrix and performed principal component analysis (PCA) to account for population structure
157 (*Yang et al., 2011*). Genome-wide association analysis was performed using PLINK version 1.9 for
158 logistic regression with multiple covariates consisting of the child's birth month, enrolment year (as
159 a marker of RSV season), daycare attendance, presence of another child \leq 6 years of age at home,
160 sex, and 6 ancestry principal components (PCs) (*Purcell et al., 2007*).

161 As the multiple testing burden likely precluded identification of small genetic effects in our
162 GWAS, we conducted an additional heritability analysis using the method described by *Golan et al.*
163 *(2014)* to estimate narrow-sense heritability of RSV infection during infancy on the latent liability
164 scale (h_l^2), which, > 0 , would indicate an accumulation of small genetic effects. We estimated h_l^2 to

165 be exactly 0, suggesting that, if present, infant RSV infection-related genetic signals are both small
166 and sparse.

167 **RSV whole-genome sequencing**

168 RSV genome sequencing was performed on all specimens from subjects meeting illness criteria
169 and with positive RSV PCR. Viral amino acid variants (genotype) of the F and G glycoprotein were
170 tested for association with prolonged infection adjusting for host features associated with increased
171 infection risk. We focused our analyses on the surface F (fusion) and G (attachment) proteins of RSV
172 as they have been implicated in pathogenesis (*Boyoglu-Barnum et al., 2015; Bukreyev et al., 2012*),
173 and both are targets for neutralizing antibodies during infection (*Anderson et al., 1988; Ngwuta
et al., 2015*). Lastly, to determine if the variants of interest were enriched by selective pressure over
175 time, we used public data from the past several decades to assess variant frequency over time.

176 RSV whole-genome sequencing of this study population has been previously described (*Schobel
et al., 2016*). Briefly, RNA was extracted at J. Craig Venter Institute (JCVI) (<https://www.jcvi.org>)
177 in Rockville, MD from nasal wash samples which were RSV PCR positive and collected during a
178 respiratory illness visit triggered through biweekly surveillance of symptoms. Four forward reverse-
179 transcription (RT) primers were designed and four sets of PCR primers were manually picked from
180 primers designed across a consensus of complete RSV genome sequences using JCVI's automated
181 primer design tool (*Li et al., 2012*). cDNA was generated from 4 µL undiluted RNA, using the pooled
182 forward primers and SuperScript III Reverse Transcriptase (Thermo Fisher Scientific, Waltham, MA,
183 USA). 100 ng of pooled DNA amplicons were sheared to create 400-bp libraries, which were pooled
184 in equal volumes and cleaned with Ampure XP reagent (Beckman Coulter, Inc., Brea, CA, USA).
185 Sequencing was performed on the Ion Torrent PGM using 316v2 or 318v2 chips (Thermo Fisher
186 Scientific).

188 For samples requiring extra coverage, in addition to Ion Torrent sequencing, Illumina libraries
189 were prepared using the Nextera DNA Sample Preparation Kit (Illumina, Inc., San Diego, CA,
190 USA). Sequence reads were sorted by barcode, trimmed, and assembled de novo using CLC
191 Bio's *clc_novo_assemble* program, and the resulting contigs were searched against custom, full-
192 length RSV nucleotide databases to find the closest reference sequence. All sequence reads
193 were then mapped to the selected reference RSV sequence using CLC Bio's *clc_ref_assemble_long*
194 program (*Aarhus, 2016*). Curated assemblies were validated and annotated with the viral anno-
195 tation software called Viral Genome ORF Reader, VIGOR 3.0 (<https://sourceforge.net/projects/jcvi-vigor/files/>), before submission to GenBank as part of the Bioproject accession PRJNA225816
196 (<https://www.ncbi.nlm.nih.gov/bioproject/225816>) (*Wang et al., 2012*) and PRJNA267583 (<https://www.ncbi.nlm.nih.gov/bioproject/267583>).

199 **Viral sequence alignment**

200 The NCBI-tools Tbl2asn (<https://www.ncbi.nlm.nih.gov/genbank/tbl2asn2/>) was used in the creation
201 of sequence records for submission to GenBank (<https://www.ncbi.nlm.nih.gov/genbank/>). A total
202 of 350 viral sequences in .sqn file format were used for downstream analysis.

203 We computed a phylogenetic tree for each gene, as follows. NCBI-tools asn2fsa ([https://www.
204 huge-man-linux.net/man1/asn2fsa.html](https://www.huge-man-linux.net/man1/asn2fsa.html)) was used to convert sequences to fasta format. Each
205 sample consisted of 11 sequence segments (NS1, NS2, N, P, M, M2-1, M2-2, SH, G, F, and L) as
206 shown in Figure S1. These were separated and repooled to create 11 single fasta files for each gene

207 containing all 350 samples. Sequences were checked for at least 90% coverage of the corresponding
208 gene to minimize loss of aligned positions when computing the phylogenetic tree. Each of the
209 eleven resulting sets was aligned with MAFFT v7 (<https://mafft.cbrc.jp/alignment/software/>) (Katoch
210 and Standley, 2013), using default parameters. The sequence of the orthologous gene from *Bovine*
211 *orthopneumovirus* (GenBank:NC_001989) was added to each set as an outgroup.

212 IQ-Tree (<https://www.iqtree.org>) (Nguyen et al., 2015) was used with per-gene multiple sequence
213 alignment (MSA) files based on amino acid sequence for estimating maximum-likelihood phyloge-
214 nies using protein substitution model. Examining the sequences with an alignment viewer showed
215 that a small number of sequences had frame-shift variants but these did not affect the regions
216 included in our testing criteria.

217 Viral sequence data and clinical information were merged and cleaned with R. Clinical IDs
218 matching more than one viral sequence ID were used to re-identify samples from the same
219 individual as prolonged infections. Genetic variation was quantified in these samples, and for
220 subsequent analysis, only the first viral sequence was included for association testing. Antigenic
221 grouping of strain A and B had been completed previously and labels were included to annotate
222 each sample accordingly.

223 The cohort-specific variant frequency per position was calculated; residues were counted and
224 ranked by frequency with the most frequent residue defined as reference (REF) and alternative
225 (ALT) for variants. Positions with at least one ALT were checked for potential misalignment or
226 other sources of error. Variant positions were selected for association analysis, while non-variant
227 positions were ignored.

228 A number of host features have been previously shown to influence infection susceptibility and
229 were therefore included as covariates in our analysis (Rosas-Salazar et al., 2022). Six samples were
230 excluded due to insufficient covariate data, resulting in 344 test samples. Of these, 38 were from
231 the same patients (prolonged infection) of which half (19) were included for association testing.
232 Thus, the test set was comprised of single samples collected from 325 individuals.

233 **Viral population structure**

234 The genetic distances to nearest neighbors were computed based on phylogenetic trees generated
235 with MAFFT. PCA and singular value decomposition (SVD) were used in dimensionality reduction
236 for exploratory data analysis of viral phylogeny. The R package factoextra was used for PCA and to
237 visualise eigenvalues and variance. R package caret was used to analyse genetic correlations.

238 **Viral variant association testing**

239 Viral amino acids (genotype collapsed into REF/ALT) were tested for association with infection types
240 (i.e., resolved and prolonged) including key covariates that alter infection risk. To reduce the multiple
241 testing burden, proxy amino acid variants were identified by performing clumping with ranking
242 based on MAF and with a cut-off threshold of $r^2 \geq 0.8$ (Supplemental Figure S3). Since many variants
243 within RSV coding genes have non-random association due to selection, analogous to linkage
244 disequilibrium (LD) in human GWAS, we reduced the multiple testing burden by retaining proxy
245 variants and removing those with $r^2 \geq 0.8$. Analysis was performed using logistic regression with
246 the R stats (3.6.2) glm function. The model consisted of the binary response (prolonged infection
247 Yes/No) and predictors viral genotype (REF/ALT amino acid, including multi-allelic non-REF collapsed
248 into ALT), viral PCs 1-5, host sex, and host features that have been previously demonstrated as

249 significantly associated with infection: self-reported race/ethnicity, daycare attendance, and living
250 with siblings (*Rosas-Salazar et al., 2022*).

251 Environmental host covariates did not contribute significant effect in our model for candidate
252 causal association. For this reason, in our main analysis, viral population structure was accounted
253 for by the first five PCs. The Bonferroni correction for multiple testing was applied based on the
254 number of variants tested. For the significant association found by proxy amino acid variants, the
255 association model was repeated for all clumped variants to produce a LocusZoom-style Manhattan
256 plot containing r² by color and p value statistics. R package stats was used for a range of analyses
257 including glm for logistic regressions. R package MASS was used to analyse logistic regression
258 model data. To test if the significantly associated variants were due to population structure, we
259 re-estimated models using the subset of individuals infected with RSV strain B to confirm validity of
260 combined analysis.

261 The relatively small sample size of our cohort required analysis that targeted only genes which
262 were a priori likely to functionally contribute to the clinical phenotype. Therefore, our analysis
263 focused on the F and G glycoprotein.

264 **Public viral sequence data**

265 We gathered publicly available sequence data to further assess variants of interest. We used the
266 public viral data repository of NCBI (https://www.ncbi.nlm.nih.gov/labs/virus/vssi/#/virus?SeqType_=Nucleotide&VirusLineage_ss=Human%20orthopneumovirus,%20taxid:11250) to retrieve infor-
267 mation using search criteria that follow. Virus: Human orthopneumovirus (RSV), taxid:11250.
268 Proteins: attachment glycoprotein. Host: Homo (humans), taxid:9605. Collection dates: Jan 1,
269 1956 onward. Nucleotide and protein sequence data was collected, which contained data from 27
270 countries and 1084 glycoprotein protein sequences after curation. Sequence and meta data were
271 merged. Multiple sequence alignment was performed to find consensus relative positions for all
272 sequences. Regions of interest were then extracted and re-annotated with their correct amino acid
273 positions matching the reference sequence. Summary statistics were generated, including number
274 of samples, collection date, geo-location, variant frequency, and strain. for the specified amino acid
275 (Supplemental Figure S4).

277 **Biological interpretation**

278 As infant RSV infection stimulates an acute antiviral response and also results in decreased barrier
279 function of the airway epithelium (*Connelly et al., 2021*), we tested for association between host
280 interferon (IFN) response and the amino acid (REF/ALT) identified as the viral variant associated
281 with prolonged infection. A Wilcoxon test was performed to compare IFN- γ , and IFN- α , between
282 RSV amino acid positions, with adjustment for the same covariates as in the main analysis. Protein
283 structures were analysed with data sourced from RCSB PDB <https://www.rcsb.org>. Protein function
284 and domains were assessed using UniProt (<https://www.uniprot.org>) for P03423 (GLYC_HRSVA)
285 (strain A2) and O36633 (GLYC_HRSVB) (strain B1) in gff format (<https://www.uniprot.org/uniprot/P03423> and <https://www.uniprot.org/uniprot/O36633>, respectively). Interactions, post-translational
286 modifications, motifs, and epitopes were assessed from the literature. Protein features were
287 assessed using data from NCBI (https://www.ncbi.nlm.nih.gov/igp/NP_056862.1) and via sequence
288 viewer with O36633.1 human RSV B1, (<https://www.ncbi.nlm.nih.gov/projects/sviewer/?id=O36633.1>). Potential effects of amino acid variation on protein structure and function were considered

291 according to available information on a broad range of biological and biochemical features, including
292 native conformation (secondary, tertiary, and quaternary), domains and topology, disulfide bonds,
293 glycosylation, interactions with other viral proteins and host-cell factors, proteolytic cleavage sites,
294 normal patterns of intra-and/or extra-cellular distribution, and secretion status.

295 **Results**

296 **Cohort characteristics**

297 The INSPIRE cohort consisted of 1,949 enrolled infants among whom there were 2,093 in-person
298 respiratory illness visits completed during winter virus season, November – March, of each year
299 (Figure S1); the median (interquartile range [IQR]) number of in-person respiratory illness visits per
300 infant during this surveillance window was 1 (*Lawless et al., 2020; Hall et al., 2009*). There were
301 344 RSV PCR-positive samples from 325 individuals which were sequenced. Prolonged infection
302 was a priori defined as repeatedly meeting criteria for acute respiratory infection with RSV PCR
303 positive nasal samples \geq 15 days between testing. There were 19 infants who met the definition
304 of prolonged infection with available viral sequencing used to confirm clonality of original and
305 subsequent virus detections. The mean RSV CT value of first infections was 25.9 ± 7.1 , and second
306 detection was 31.6 ± 5.4 . The mean number of days between detections was 25 ± 25 days (Figure
307 S2). Table 1 lists the cohort characteristics of infants with prolonged RSV infection compared with
308 other RSV infection and the entire cohort.

		Prolonged RSV Infection N=19	RSV Infection N=342	Total N=1949
Illness	Age at first illness, months (median, IQR)	6 (4, 6)	4 (2, 5)	NA
	Respiratory severity score (median, IQR)	2.0 (1.2, 3.0)	3.0 (2.0, 4.0)	NA
RSV season	2012-13	68%	54%	44%
	2013-14	32%	46%	56%
Self reported Race	Non-Hispanic Black	37%	13%	18%
	Non-Hispanic White	63%	69%	65%
	Hispanic	0%	10%	9%
	Multi-race/ethnicity/other	0%	8%	8%
Sex	Female	53%	44%	48%
Second-hand smoke exposure	Yes	21%	23%	47%
Health insurance	Medicaid	68%	48%	54%
	Private	32%	51%	45%
	None/unknown	0%	1%	1%
Daycare/Sibling*	Yes	84%	78%	66%

Table 1. Characteristics of infants with prolonged RSV infection compared with other RSV infection and the entire cohort. Prolonged infection is defined as repeatedly RSV PCR-positive with ≥ 15 days between testing and meeting criteria for acute respiratory infection. *Presence of sibling or another child ≤ 6 years of age at home.

309 **Host genetic analyses**

310 We explored whether RSV infection in infancy is a natural assignment (quasi-random) event and,
 311 unlike severity of early-life RSV infection (*Larkin and Hartert, 2015*), occurs independently of host
 312 genetics. For the candidate SNP analysis, we considered childhood asthma- and RSV LRTI-associated
 313 SNPs identified in *Pividori et al. (2019)*; *Janssen et al. (2007)*; *Pasanen et al. (2017)*. The first is the
 314 largest childhood asthma GWAS to date, and, to our knowledge, the latter 2 represent the most
 315 comprehensive studies of RSV LRTI-associated SNPs. To further reduce the multiple testing burden,
 316 we only analysed SNPs with MAF ≥ 0.1 in at least one of the White, Black, or Hispanic ethnicity
 317 groups. Associations between genotype at the resulting 54 SNPs (50 childhood asthma- and 4 RSV
 318 LRTI-associated SNPs) and RSV infection in infancy in our data are given in Figure 1. The data are
 319 consistent with little to no effect of genotype at these SNPs on RSV infection in infancy.

320 We further investigated the possibility that the analysis was underpowered to identify associa-
 321 tions with these SNPs by pooling information across SNPs to estimate the average genetic effect
 322 size. Our analysis in the supplement shows that the average genetic effect of these LRTI- and
 323 asthma-related SNPs on infant RSV infection is zero or trivially small (Supplemental).

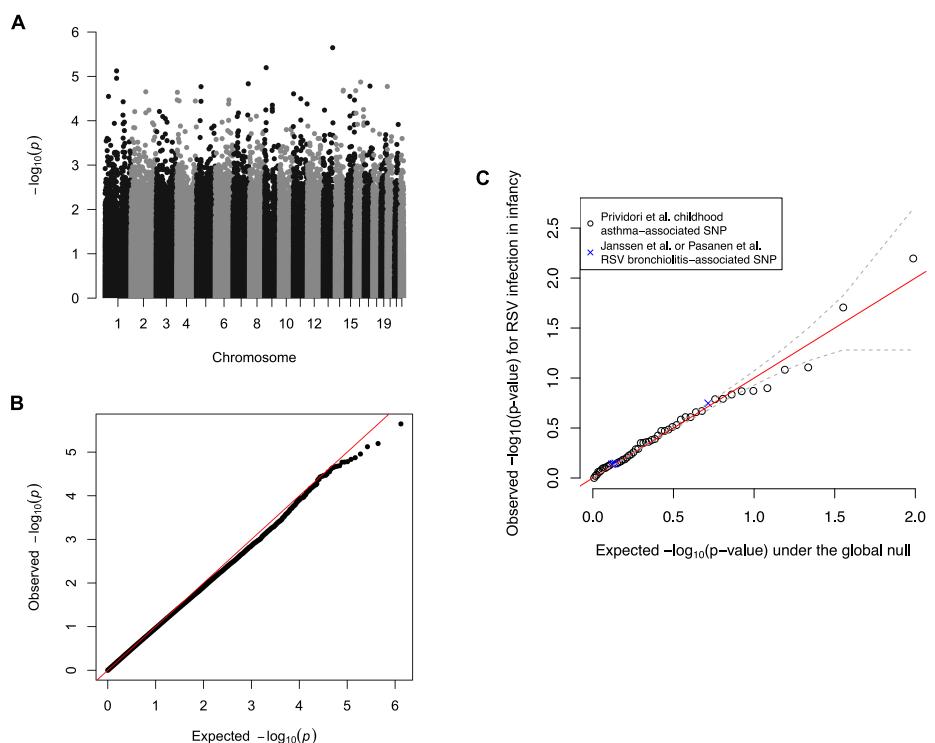


Figure 1. Genetic analyses of RSV infection in infancy. (A) The Manhattan plot shows no genome-wide significant associations (p value threshold of $5e^{-8}$). (B) The Q-Q plot demonstrates that the observed p values are congruent with those expected under the null hypothesis that RSV infection in infancy is independent of host genotype. (C) The association between the 54 selected childhood asthma- or RSV LRTI-associated SNPs and RSV infection in infancy in our data. The identity line is shown in red, and the dashed grey lines are \pm standard deviation around the expected $-\log_{10}(p$ value). RSV: respiratory syncytial virus; SNP: single nucleotide polymorphism.

324 **Population structure**

325 A summary of protein coding genes in RSV is illustrated in Figure 2 A. Our analysis focused on F
326 and G protein. The phylogenetic tree based on multiple sequence alignment (MSA) of G protein
327 amino acid sequences is shown in Figure 2 B. One obvious feature causing a separation in genetic
328 diversity is G protein partial gene duplication, which has emerged in recent years within RSV-A
329 strains (*Eshaghi et al., 2012*). RSV-B strains with an analogous duplication have existed for two
330 decades, although the selection process leading to emergence and clinical implications have not
331 been entirely defined.

332 PCA was used for reducing the dimensionality of sequence data, where PC1 accounted for
333 95.19% of cumulative variance, and variance attributed to other PCs was roughly uniformly dis-
334 tributed (Figure 2 C). We observed prolonged infections by viruses from different phylogenetic
335 clades, rather than one specific clade (Figure 2 C), indicating that these results are not confounded
336 by latent clade membership.

337 **Genetic invariance of prolonged infection**

338 The duration of RSV shedding in Kenyan infants has been reported previously (*Okiro et al., 2010*).
339 Based on these findings, infections separated by at least 15 days with symptoms were expected
340 to be "new" infections (*Okiro et al., 2010*). Figure 2 D (panel [i]) summarizes every pairwise genetic
341 distance between every viral sequence, where small distances indicate pairs with closely related
342 sequences. Panels [ii] and [iii], which summarize the difference in sequence similarity distributions
343 between viruses from the same host and different hosts, show that RSV sequences corresponding
344 to initial and subsequent viral detections are nearly identical. These results support the conclusion
345 that such cases are prolonged (i.e., failure to clear) infections rather than new infections.

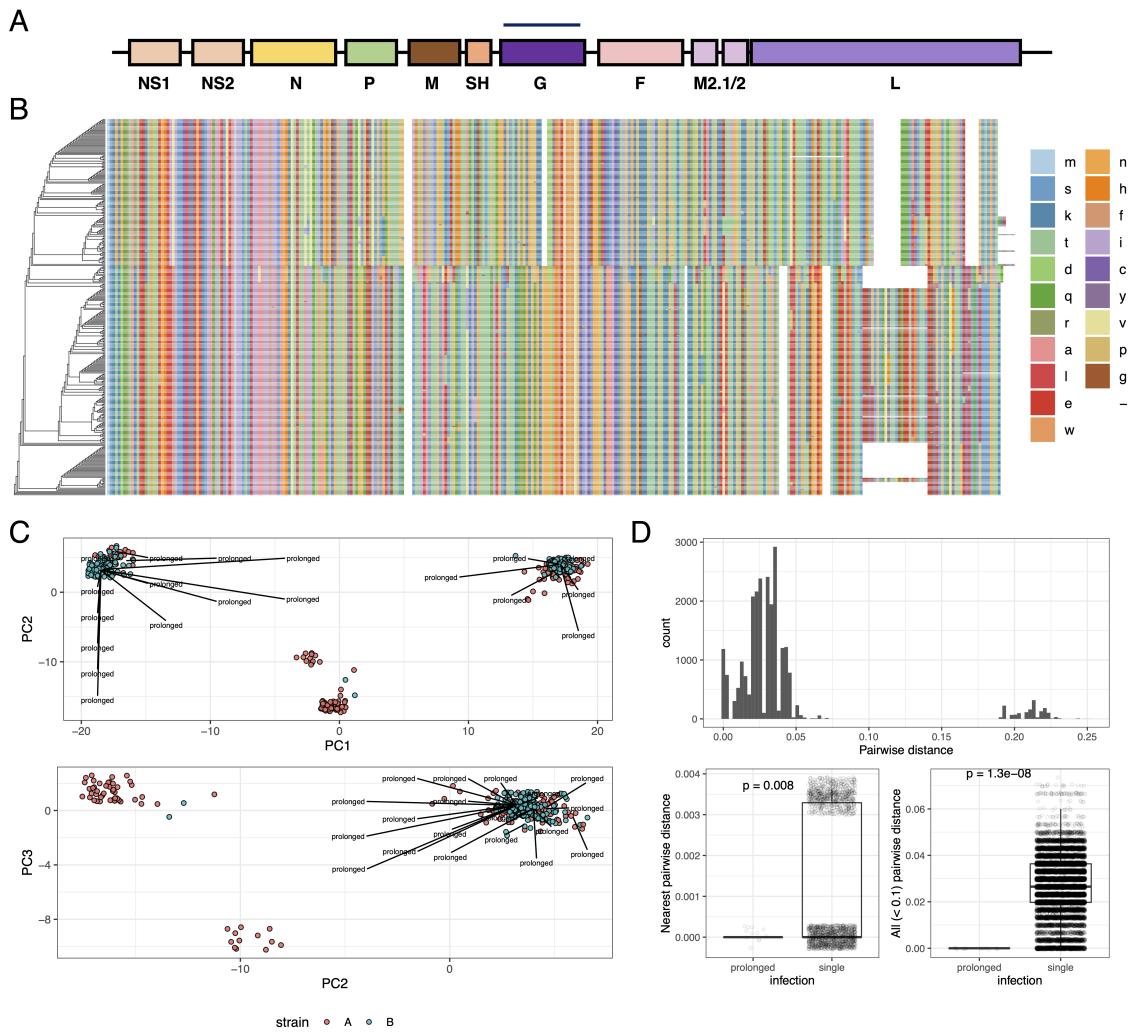


Figure 2. Viral population structure. (A) Linear map of the RSV genome. (B) Phylogenetic tree based on multiple sequence alignment MSA of G protein amino acid sequences. Color; amino acids. (C) Principal component (PC) analysis. PCs 1-3 with labels indicating prolonged infections from different phylogenetic clades. (D) Panel [i] summarises every pairwise genetic distance between every viral sequence. Genetic invariance in prolonged infections separated by at least 15 days compared to other genetic variation within the most closely related sequences (panel [ii]) and within all possible closely related pairs (panel [iii]). Jitter applied for visualisation.

346 **Variants in G glycoprotein significantly associated with prolonged infection**

347 The consensus sequence within the cohort was assigned based on the major allele. Variants at the
348 amino acid level were defined as either reference (REF) or alternative (ALT) and assessed for their
349 association with prolonged infection. The model consisted of (A) the binary response (prolonged
350 infection Yes/No), and (B) predictors: (1) viral genotype (REF/ALT amino acid), (2) viral PCs 1-5, (3)
351 host sex, and host features that have been previously demonstrated as significantly associated with
352 infection; (4) self-reported race/ethnicity, (5) child-care attendance, or living with another child \leq 6
353 years of age at home (*Hall et al., 1976*). A significant genetic association was identified between
354 prolonged infection and the lead variant after Bonferroni correction for multiple testing (threshold
355 for number independent variants $< 0.05/23 = 0.002$), as shown in Figure 3 A, p value = 0.0006.

356 To determine whether this association was simply due to population stratification between
357 strains A and B, a subset analysis was performed using independently assessed clinical laboratory
358 strain labels for A and B. The same direction of effect indicated that the association was not a
359 false positive, although the significantly smaller sample size prevented the sub-analysis result from
360 crossing the significance threshold.

361 To assess the possibility of a false positive association due to population structure within our
362 cohort, we assessed the magnitude of variance explained (VE) at every amino acid position. Figure 3
363 B (panel [i]) shows the variance explained by each amino acid in PCs1-5. The cumulative proportion
364 of variance for PCs 1-5 was 99.5% (PC1 = 95%, PC2 = 3%). The values are illustrated according to
365 protein position in panels [ii-iii]. The lead association variant had 0.603% VE for PC1 and 0.458%
366 VE for PC2, a negligible effect that precludes spurious association by allele frequency between
367 populations.

368 After identifying a significant viral genetic association with prolonged infection, we quantified
369 the correlation of variants with the lead proxy. Clumping was performed with ranking based on
370 MAF and with a cut-off threshold of $r^2 \geq 0.8$ (Supplemental Figure S3). The association model was
371 repeated for all variants, defining protein p.E123K/D and p.P218T/S/L as candidate causal variants
372 associated with prolonged infection as shown in Figure 3 C. No other variants were correlated with
373 this outcome.

374 To determine whether p.E123K/D and p.P218T/S/L variant genotypes are novel and potentially
375 influence viral fitness, we searched the public viral data repository of NCBI Human orthopneu-
376 movirus, taxid:11250, which contained data from 27 countries worldwide, sample collection dates
377 from 1956 onward, and 1084 glycoprotein protein sequences after curation. The variants were
378 present at a low and stable frequency, without obvious temporal enrichment (Supplemental Figure
379 S4). Thus, while historical data reveal no positive selective advantage attached to p.E123K/D and
380 p.P218T/S/L, longstanding circulation and linkage in prolonged RSV infection suggest that these
381 polymorphisms are present in the viral inoculum and do not arise through recurrent mutational
382 events.

383 Due to multiple testing correction according to our statistical analysis plan, an association also
384 originally identified in F protein was rejected and therefore omitted from further discussion. For
385 posterity, the variant position was p.N116S (relative to strain A GenBank: AMN91253.1).

386 **Functional interpretation**

387 Cell-attachment proteins of paramyxoviruses (G protein in RSV) span the viral envelope and form
 388 spike-like projections from the virion surface. RSV G protein is a type II integral membrane protein
 389 consisting of 298 amino acid residues comprising N-terminal cytoplasmic (p.1-43), transmembrane
 390 helical (p.43-63), and extracellular (p.64-298) domains (Figure 3 D). RSV G protein ectodomain also
 391 exists in a soluble secreted form, p.66 – 298, which functions in immune evasion (*Levine et al., 1987; Feldman et al., 1999, 2000*). G protein interacts with the small hydrophobic (SH) protein
 392 (*Rixon et al., 2005*) and, via the N-terminus, with matrix (M) (*Ghildyal et al., 2005*) protein. It has
 393 also been reported to form homo-oligomers (*Collins and Mottet, 1992*). The variant amino acid
 394 positions associated with prolonged infection reside in a portion of the ectodomain of unassigned
 395 specific function and linearly non-contiguous with sequences that bind cell-surface heparan sulfate,
 396 which likely promotes RSV cell-attachment (p.187-198) (*Levine et al., 1987; Feldman et al., 1999, 2000*). In addition, these positions do not contribute to known neutralization epitopes on G protein.
 397 Information available in PDB was insufficient to infer effects of p.E123K/D and p.P218T/S/L on local
 398 or regional protein structure. The potential effect on glycosylation is indeterminate. Figure 3 D
 399 illustrates the position (dotted red lines) of these variants relative to summarised known functional
 400 features.

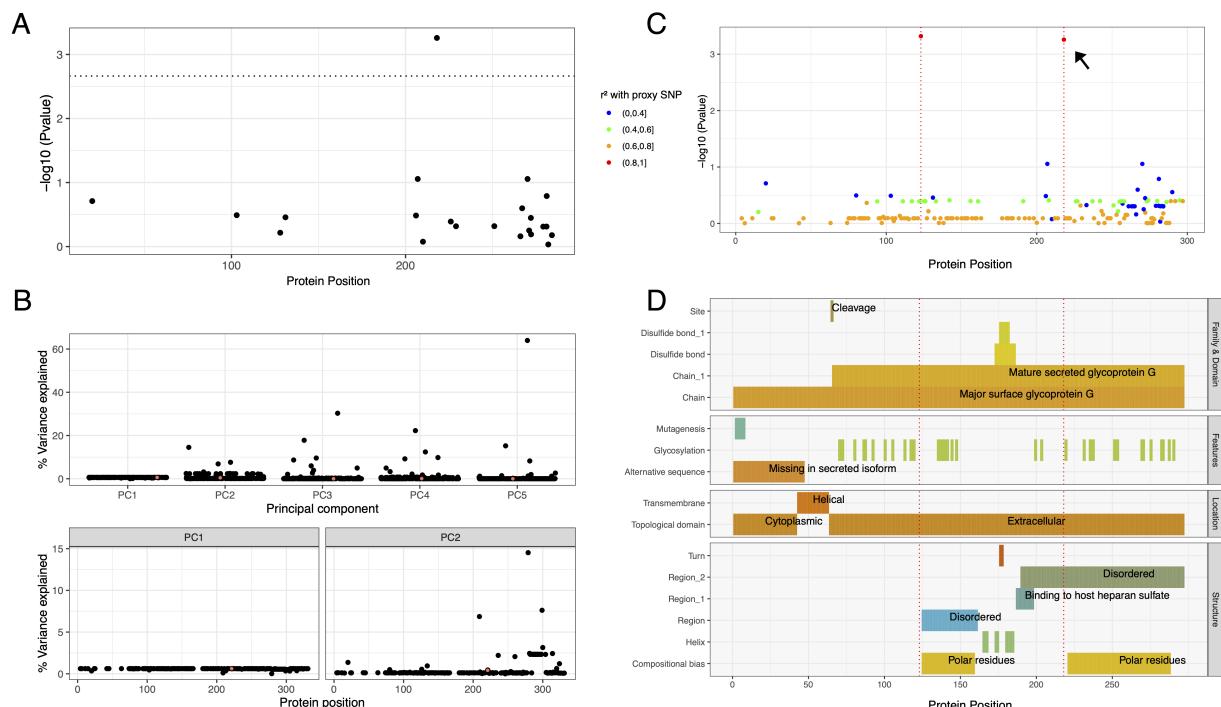


Figure 3. Viral genetic association with prolonged infection. (A) Amino acid association with prolonged infection after multiple testing correction (significant threshold shown by dotted line). (B) Variance explained (VE) within cohort. The effect of each variant on cohort structure is shown for PCs1-2. The small % VE for a significantly associated lead variant supports a true positive. (C) Variants in strong correlation were clumped for association testing using proxies for $r^2 \geq 0.8$. One significant association was identified (shown in A); the r^2 values for all other variants show a single highly correlated variant with the lead proxy (red), identifying p.E123K/D and p.P218T/S/L. (D) Evidence for biological interpretation for every amino acid position is summarised. Dotted red lines indicate the positions at p.123, p.218.

403 **Host response**

404 Prolonged infections associated with G protein variants p.E123K/D and p.P218T/S/L were on average
405 less severe compared with other circulating variants, and all were limited to the upper respiratory
406 tract (Table 1). Therefore, we analysed nasal wash samples collected during acute RSV infection for
407 a panel of cytokines involved in antiviral immune responses and observed differential IFN α and IFN γ
408 levels segregating according to viral antigenic group—A or B. Both cytokines were elevated in group
409 B infections compared to group A. The groups A and B median (lower-and upper-quartile) values
410 were 9.5 (3-22.5) and 12.6 (4.1-25.8), respectively, for IFN α and 3.6 (1-7) and 4 (2-7.4), respectively, for
411 IFN γ (group A, n = 149; group B, n = 103). As prolonged infections with p.E123K/D and p.P218T/S/L
412 genotypes were exclusively group B, the dichotomous relationship of IFN α and IFN γ levels to
413 antigenic group precluded evaluation of G protein variants as independent predictors of IFN α and
414 IFN γ production.

415 **Discussion**

416 In this study of term healthy infants, we found no evidence of host genetic susceptibility to RSV
417 infection during infancy. This allowed our analysis to focus on elucidation of viral drivers of
418 prolonged infection. A significant viral genetic association in the RSV G protein, p.E123K/D and
419 p.P218T/S/L, with prolonged infant RSV infection was identified. These variants were not associated
420 with severe disease, and public data reveal their consistent presence at low frequencies over the
421 past 30 years, without evidence of enrichment by positive selective pressure over time. The two
422 variants we identified in G are correlated with non-random association, analogous to LD in the
423 human diploid genome and therefore not likely random mutations, but instead co-inherited in
424 the infecting inoculum. This suggests an evolutionary benefit and raises the question of why such
425 variants have maintained a stable but low frequency in the human population for decades. These
426 strains are a potential reservoir, emerging seasonally in response to immune, environmental, or
427 other forces. Alternatively, the polymorphisms might recurrently arise de novo during infection
428 of some individuals but are poorly transmissible because of suboptimal fitness. The possibility of
429 viral mutational immune escape has been reported for infants who struggle to control primary RSV
430 infections, allowing for prolonged viral replication and not previously described viral rebound (**Brint**
431 *et al., 2017*).

432 The RSV variants associated with prolonged infection in our cohort, G p.E123K/D and p.P218T/S/L,
433 lie in the extracellular region, and there are no known mechanistic features that directly overlap,
434 although it is possible that variant positions approximate sequences that bind a putative viral
435 receptor, heparan sulfate (**Feldman et al., 1999**), in the G protein three-dimensional structure. G
436 protein amino acid positions 123 and 218 are not part of known antibody neutralization epitopes
437 or CD8+ cytotoxic T-cell epitopes (Figure 3 D). In addition to heparan sulfate, interactions between
438 viral G protein and CX3CR1, the receptor for the CX3C chemokine fractalkine, have been reported
439 to modulate the immune response and facilitate infection (**Levine et al., 1987; Feldman et al.,**
440 **1999, 2000; Johnson et al., 2015; Tripp et al., 2001; Jeong et al., 2015**). Furthermore, the mature
441 secreted isoform of G protein (p.66-298) is thought to facilitate viral antibody evasion by acting as
442 an antigen decoy and modifying the activity of leukocytes bearing Fc-gamma receptors (**Bukreyev**
443 *et al., 2008*). Our findings raise the interesting prospect that G protein variants associated with
444 prolonged infection alter a key interaction at the immune interface between pathogen and host.

445 Although this study was not designed to define mechanisms underlying the association of G
446 protein variants with prolonged infection, these sequence changes might dampen antiviral immune
447 responses and thereby delay viral clearance. Although we observed differences in the acute
448 antiviral response between subjects with resolved and prolonged infection, specifically increased
449 levels of types 1 and 2 IFN in nasal secretions, we could not make causal inference about variant
450 sequences because of confounding by co-linearity of these polymorphisms with RSV antigenic
451 group. Results of nasal cytokine analysis are nevertheless consistent with a contemplated role for
452 altered immune responses in extended infections by G protein variant strains. It is also possible
453 that strains harbouring G protein p.E123K/D and p.P218T/S/L variants are cleared more slowly
454 and foster an immune environment of low-level chronic stimulation or exhaustion. We previously
455 demonstrated that infants infected with RSV in their first year of life have damped subsequent
456 antiviral immune responses in early childhood (*Chirkova et al., 2022*) as well as changes in airway
457 epithelial cell metabolism (*Connelly et al., 2021*).

458 While this study has a number of significant strengths, including one of few population-based
459 surveillance studies of first RSV infections during infancy among term healthy infants, our findings
460 are also subject to some limitations. First, this study was not designed with the primary intention to
461 examine infection duration, and additional sampling following initial RSV infection was triggered
462 by a repeat acute respiratory illness. Asymptomatic prolonged infections would therefore not
463 have been captured. Second, our study cohort was small, necessitating focus on viral surface
464 glycoproteins, F and G, due to their variability and importance in host immunity. A larger cohort
465 with serial sampling would be required to diminish the impact of co-linearity of viral genotypes
466 with antigenic groups and to perform informative viral whole genome analysis. Genome-wide
467 information might elucidate other determinants of prolonged infection or pathogen fitness that
468 mediate and/or modulate effects of phenotype-driving variations. Third, again due to small sample
469 size, we could only investigate host genetic risk for infection, not prolonged infection. While we have
470 not specifically assessed subjects for rare monogenic variants that may underlie immunodeficiency,
471 our enrolment criteria included only infants who were term and otherwise healthy. While we
472 performed an interaction analysis for the outcome of host asthma, host genetics, and pathogen
473 genetics and found no significant interaction, our sample size is unlikely sufficient to exclude such
474 an interaction. Lastly, while we do not expect a role for immune memory in these first-in-life RSV
475 infections, we cannot exclude modulatory effects of maternal antibody, which we did not measure.

476 In summary, we identified a novel RSV viral variant associated with prolonged infection in
477 healthy infants, but no evidence of host genetic susceptibility to infant RSV infection. Understanding
478 host and viral mechanisms that contribute to prolonged infection will be important in crafting
479 strategies to control the short and long-term impact of RSV infection. The identification of RSV
480 variants associated with prolonged infection might also improve vaccine design, particularly if
481 these variants stimulate robust immunity or, in contrast, escape the immune response or induce
482 immunopathologic conditions. The growing availability of large genomic and functional data sources
483 provides opportunities for advancing our understanding of the pathogenesis of infant RSV infection,
484 defining the contribution of viral genetic variants to acute and chronic disease, and informing
485 the development of effective vaccines. As neither the capacity of RSV for prolonged infection in
486 immunocompetent hosts nor a viral reservoir has been delineated, these results are of fundamental
487 interest in understanding viral and host genetic contributions that may promote prolonged infection

488 influence development of chronic respiratory morbidity.

489 **Links**

490 **Software**

491 R v4.1.0 was used for data preparation and analysis <http://www.r-project.org>.
492 R package *caret* was used for analysis: genetic correlations.
493 R package *dplyr* was used for data curation.
494 R package *factoextra* was used for analysis: PCA, and to visualise eigenvalues and variance.
495 R package *ggplot2* was used for data visualisation.
496 R package *MASS* was used to analysis: logistic regression model data.
497 R package *stats* was used for analysis: including *glm* for logistic regressions.
498 R package *stringr* was used for data curation.
499 R package *tidyR* was used for data curation.
500 asn2fsa <https://www.huge-man-linux.net/man1/asn2fsa.html>
501 clc_novo_assemble qiagenbioinformatics.com
502 Clustal Omega <https://www.ebi.ac.uk/Tools/msa/clustalo/>
503 dbNSFP (database) <http://database.liulab.science/dbNSFP> (*Liu et al., 2016*)
504 GCTA <https://cnsgenomics.com/software/gcta/> (*Yang et al., 2011*)
505 GenBank <https://www.ncbi.nlm.nih.gov/genbank/>
506 IQ-Tree <https://www.iqtree.org/> (*Nguyen et al., 2015*)
507 KING <https://people.virginia.edu/~wc9c/KING/> (*Manichaikul et al., 2010*)
508 MAFFT <https://mafft.cbrc.jp/alignment/software/> (*Katoh and Standley, 2013*)
509 NextAlign <https://github.com/nextstrain/nextclade>
510 PLINK <http://zzz.bwh.harvard.edu/plink/> (*Purcell et al., 2007*)
511 Tbl2asn <https://www.ncbi.nlm.nih.gov/genbank/tbl2asn2/>
512 Viral Genome ORF Reader, VIGOR 3.0 <https://sourceforge.net/projects/jcvi-vigor/files/>
513 RCSB PDB <https://www.rcsb.org>
514 UniProt <https://www.uniprot.org>

515 **Data sources**

516 Dataset <https://www.ncbi.nlm.nih.gov/bioproject/267583>.
517 Dataset <https://www.ncbi.nlm.nih.gov/bioproject/225816>.
518 J. Craig Venter Institute <https://www.jcvi.org>.
519 GenBank:NC_001989 *Bovine orthopneumovirus*, complete genome https://www.ncbi.nlm.nih.gov/nucleotide/NC_001989.
520 Reference data <https://www.ncbi.nlm.nih.gov/gene/?term=1489824>. G attachment glycoprotein [*Human orthopneumovirus*]; ID: 1489824; Location: NC_001781.1 (4675..5600); Aliases: HRSVgp07.
521 Reference data <https://www.ncbi.nlm.nih.gov/gene/?term=37607642>. G attachment glycoprotein [*Human orthopneumovirus*]; ID: 37607642; Location: NC_038235.1 (4673..5595); Aliases: DZD21_gp07.
522 Reference data for all public NCBI Virus <https://www.ncbi.nlm.nih.gov/labs/virus/vssi/> for species:
523 *Human orthopneumovirus*; genus: *Orthopneumovirus*; family: *Pneumoviridae*.

529 Reference data https://www.ncbi.nlm.nih.gov/labs/virus/vssi/#/virus?SeqType_s=Nucleotide&VirusLineage_ss=Human%20orthopneumovirus,%20taxid:11250 - contains sequence data for Virus Lineage ss=Human orthopneumovirus, taxid:11250 nucleotide: 26'965, protein: 53'804, RefSeq Genomes: 2.
530
531
532
533 Reference https://www.ncbi.nlm.nih.gov/protein/NP_056862.1
534 GCF_002815475.1 (release 2018-08-19) Nucleotide Accessions: NC_038235.1, protein: Y_009518856.1
535 Reference https://www.ncbi.nlm.nih.gov/protein/YP_009518856.1
536 GCF_000855545.1 (release 2015-02-12) Nucleotide Accessions: NC_001781.1, protein: NP_056862.1
537 (strain B1).

538 **Code availability**

539 Analysis code is available at <https://github.com/DylanLawless/inspire2022lawless.github.io>.

540 **References**

- 541 Aarhus Q. White paper on de novo assembly in CLC Assembly Cell 4.0. digitalinsights. 2016 Jun; p. 14. <https://digitalinsights.qiagen.com/files/whitepapers/whitepaper-denovo-assembly.pdf>, place: Denmark Publisher: Qiagen.
- 544 Anderson LJ, Bingham P, Hierholzer J. Neutralization of respiratory syncytial virus by individual and mixtures of F and G protein monoclonal antibodies. *Journal of virology*. 1988; 62(11):4232-4238.
- 546 Bagga B, Harrison L, Roddam P, DeVincenzo J. Unrecognized prolonged viral replication in the pathogenesis of human RSV infection. *Journal of Clinical Virology*. 2018; 106:1-6.
- 548 Bokun V, Moore JJ, Moore R, Smallcombe CC, Harford TJ, Rezaee F, Esper F, Piedimonte G. Respiratory syncytial virus exhibits differential tropism for distinct human placental cell types with Hofbauer cells acting as a permissive reservoir for infection. *PLoS One*. 2019; 14(12):e0225767. doi: [10.1371/journal.pone.0225767](https://doi.org/10.1371/journal.pone.0225767), edition: 2019/12/04.
- 552 Boyoglu-Barnum S, Todd SO, Chirkova T, Barnum TR, Gaston KA, Haynes LM, Tripp RA, Moore ML, Anderson LJ. An anti-G protein monoclonal antibody treats RSV disease more effectively than an anti-F monoclonal antibody in BALB/c mice. *Virology*. 2015; 483:117-125.
- 555 Brint ME, Hughes JM, Shah A, Miller CR, Harrison LG, Meals EA, Blanch J, Thompson CR, Cormier SA, DeVincenzo JP. Prolonged viral replication and longitudinal viral dynamic differences among respiratory syncytial virus infected infants. *Pediatric research*. 2017; 82(5):872-880.
- 558 Bukreyev A, Yang L, Collins PL. The secreted G protein of human respiratory syncytial virus antagonizes antibody-mediated restriction of replication involving macrophages and complement. *Journal of virology*. 2012; 86(19):10880-10884.
- 561 Bukreyev A, Yang L, Fricke J, Cheng L, Ward JM, Murphy BR, Collins PL. The secreted form of respiratory syncytial virus G glycoprotein helps the virus evade antibody-mediated restriction of replication by acting as an antigen decoy and through effects on Fc receptor-bearing leukocytes. *Journal of virology*. 2008; 82(24):12191-12204.
- 564 Chirkova T, Rosas-Salazar C, Gebretsadik T, Jadhao SJ, Chappell JD, Peebles Jr RS, Dupont WD, Newcomb DC, Berdnikovs S, Gergen PJ, et al. Effect of Infant RSV Infection on Memory T Cell Responses at Age 2-3 Years. *Frontiers in immunology*. 2022; 13:826666.
- 567 Collins PL, Mottet G. Oligomerization and post-translational processing of glycoprotein G of human respiratory syncytial virus: altered O-glycosylation in the presence of brefeldin A. *Journal of General Virology*. 1992; 73(4):849-863.

- 570 Connolly AR, Jeong BM, Coden ME, Cao JY, Chirkova T, Rosas-Salazar C, Cephus JY, Anderson LJ, Newcomb DC,
571 Hartert TV, et al. Metabolic Reprogramming of Nasal Airway Epithelial Cells Following Infant Respiratory
572 Syncytial Virus Infection. *Viruses*. 2021; 13(10):2055.
- 573 Cubie HA, Duncan LA, Marshall LA, Smith NM. Detection of respiratory syncytial virus nucleic acid in archival
574 postmortem tissue from infants. *Pediatr Pathol Lab Med*. 1997 Nov; 17(6):927–38. Edition: 1997/11/14.
- 575 Eshaghi A, Duvvuri VR, Lai R, Nadarajah JT, Li A, Patel SN, Low DE, Gubbay JB. Genetic variability of human
576 respiratory syncytial virus A strains circulating in Ontario: a novel genotype with a 72 nucleotide G gene
577 duplication. *PLoS one*. 2012; 7(3):e32807.
- 578 Feldman SA, Audet S, Beeler JA. The fusion glycoprotein of human respiratory syncytial virus facilitates
579 virus attachment and infectivity via an interaction with cellular heparan sulfate. *Journal of Virology*. 2000;
580 74(14):6442–6447.
- 581 Feldman SA, Hendry RM, Beeler JA. Identification of a linear heparin binding domain for human respiratory
582 syncytial virus attachment glycoprotein G. *Journal of virology*. 1999; 73(8):6610–6617.
- 583 Ghildyal R, Li D, Peroulis I, Shields B, Bardin PG, Teng MN, Collins PL, Meanger J, Mills J. Interaction between
584 the respiratory syncytial virus G glycoprotein cytoplasmic domain and the matrix protein. *Journal of General
585 Virology*. 2005; 86(7):1879–1884.
- 586 Glezen WP, Taber LH, Frank AL, Kasel JA. Risk of primary infection and reinfection with respiratory syncytial
587 virus. *American journal of diseases of children*. 1986; 140(6):543–546.
- 588 Golan D, Lander ES, Rosset S. Measuring missing heritability: inferring the contribution of common variants.
589 *Proceedings of the National Academy of Sciences*. 2014; 111(49):E5272–E5281.
- 590 Hall CB, Weinberg GA, Iwane MK, Blumkin AK, Edwards KM, Staat MA, Auinger P, Griffin MR, Poehling KA, Erdman
591 D, Grijalva CG, Zhu Y, Szilagyi P. The burden of respiratory syncytial virus infection in young children. *N Engl J
592 Med*. 2009 Feb; 360(6):588–98. doi: 10.1056/NEJMoa0804877, edition: 2009/02/07.
- 593 Hall CB, Geiman JM, Biggar R, Kotok DI, Hogan PM, Douglas Jr RG. Respiratory syncytial virus infections within
594 families. *New England journal of medicine*. 1976; 294(8):414–419.
- 595 Janssen R, Bont L, Siezen CL, Hodemaekers HM, Ermers MJ, Doornbos G, Slot Rv, Wijmenga C, Goeman JJ,
596 Kimpen JL, et al. Genetic susceptibility to respiratory syncytial virus bronchiolitis is predominantly associated
597 with innate immune genes. *Journal of Infectious Diseases*. 2007; 196(6):826–834.
- 598 Jeong KI, Piepenhagen PA, Kishko M, DiNapoli JM, Groppo RP, Zhang L, Almond J, Kleanthous H, Delagrave
599 S, Parrington M. CX3CR1 is expressed in differentiated human ciliated airway cells and co-localizes with
600 respiratory syncytial virus on cilia in a G protein-dependent manner. *PLoS one*. 2015; 10(6):e0130517.
- 601 Johnson SM, McNally BA, Ioannidis I, Flano E, Teng MN, Oomens AG, Walsh EE, Peebles ME. Respiratory
602 syncytial virus uses CX3CR1 as a receptor on primary human airway epithelial cultures. *PLoS pathogens*.
603 2015; 11(12):e1005318.
- 604 Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance
605 and usability. *Molecular biology and evolution*. 2013; 30(4):772–780.
- 606 Larkin EK, Gebretsadik T, Moore ML, Anderson LJ, Dupont WD, Chappell JD, Minton PA, Peebles RS Jr, Moore PE,
607 Valet RS, Arnold DH, Rosas-Salazar C, Das SR, Polack FP, Hartert TV. Objectives, design and enrollment results
608 from the Infant Susceptibility to Pulmonary Infections and Asthma Following RSV Exposure Study (INSPIRE).
609 *BMC Pulm Med*. 2015 Apr; 15:45. doi: 10.1186/s12890-015-0040-0, edition: 2015/05/30.
- 610 Larkin EK, Gebretsadik T, Moore ML, Anderson LJ, Dupont WD, Chappell JD, Minton PA, Peebles RS, Moore PE,
611 Valet RS, et al. Objectives, design and enrollment results from the infant susceptibility to pulmonary infections
612 and asthma following RSV Exposure Study (INSPIRE). *BMC pulmonary medicine*. 2015; 15(1):1–12.

- 613 Larkin EK, Hartert TV. Genes associated with RSV lower respiratory tract infection and asthma: the application
614 of genetic epidemiological methods to understand causality. Future virology. 2015; 10(7):883–897.
- 615 Lawless D, Rosas-Salazar C, Gebretsadik T, Turi K, Snyder B, Wu P, Fellay J, Hartert T. Genome-wide association
616 study of susceptibility to respiratory syncytial virus infection during infancy. In: *European Journal of Human*
617 *Genetics*, vol. 28 Springer Nature Campus, 4 Crinan St, London, N1 9XW, England; 2020. p. 319–319.
- 618 Levine S, Klaiber-Franco R, Paradiso P. Demonstration that glycoprotein G is the attachment protein of respira-
619 tory syncytial virus. Journal of General Virology. 1987; 68(9):2521–2524.
- 620 Li K, Shrivastava S, Brownley A, Katzel D, Bera J, Nguyen AT, Thovarai V, Halpin R, Stockwell TB. Automated
621 degenerate PCR primer design for high-throughput sequencing improves efficiency of viral sequencing. Virol
622 J. 2012 Nov; 9:261. doi: 10.1186/1743-422x-9-261, edition: 2012/11/08.
- 623 Liu X, Wu C, Li C, Boerwinkle E. dbNSFP v3. 0: A one-stop database of functional predictions and annotations
624 for human nonsynonymous and splice-site SNVs. Human mutation. 2016; 37(3):235–241.
- 625 Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen WM. Robust relationship inference in genome-wide
626 association studies. Bioinformatics. 2010; 26(22):2867–2873. Number: 22 ISBN: 1460-2059 Publisher: Oxford
627 University Press.
- 628 Munywoki PK, Koech DC, Agoti CN, Kibirige N, Kipkoech J, Cane PA, Medley GF, Nokes DJ. Influence of age,
629 severity of infection, and co-infection on the duration of respiratory syncytial virus (RSV) shedding. Epidemiol
630 Infect. 2015 Mar; 143(4):804–12. doi: 10.1017/s0950268814001393, edition: 2014/06/06.
- 631 Nadal D, Wunderli W, Meurmann O, Briner J, Hirsig J. Isolation of respiratory syncytial virus from liver tissue and
632 extrahepatic biliary atresia material. Scand J Infect Dis. 1990; 22(1):91–3. doi: 10.3109/00365549009023125,
633 edition: 1990/01/01.
- 634 Nguyen LT, Schmidt HA, Von Haeseler A, Minh BQ. IQ-TREE: a fast and effective stochastic algorithm for
635 estimating maximum-likelihood phylogenies. Molecular biology and evolution. 2015; 32(1):268–274.
- 636 Ngwuta JO, Chen M, Modjarrad K, Joyce MG, Kanekiyo M, Kumar A, Yassine HM, Moin SM, Killikelly AM, Chuang
637 GY, et al. Prefusion F-specific antibodies determine the magnitude of RSV neutralizing activity in human sera.
638 Science translational medicine. 2015; 7(309):309ra162–309ra162.
- 639 O'Donnell DR, McGarvey MJ, Tully JM, Balfour-Lynn IM, Openshaw PJ. Respiratory syncytial virus RNA in cells
640 from the peripheral blood during acute infection. J Pediatr. 1998 Aug; 133(2):272–4. doi: 10.1016/s0022-
641 3476(98)70234-3, edition: 1998/08/26.
- 642 Okiro EA, White LJ, Ngama M, Cane PA, Medley GF, Nokes DJ. Duration of shedding of respiratory syncytial virus
643 in a community study of Kenyan children. BMC infectious diseases. 2010; 10(1):1–7.
- 644 Pasanen A, Karjalainen MK, Bont L, Piippo-Savolainen E, Ruotsalainen M, Goksör E, Kumawat K, Hodemaekers
645 H, Nuolivirta K, Jartti T, et al. Genome-wide association study of polymorphisms predisposing to bronchiolitis.
646 Scientific reports. 2017; 7(1):1–9.
- 647 Pividori M, Schoettler N, Nicolae DL, Ober C, Im HK. Shared and distinct genetic risk factors for childhood-onset
648 and adult-onset asthma: genome-wide and transcriptome-wide studies. The Lancet Respiratory Medicine.
649 2019; 7(6):509–522.
- 650 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J, Sklar P, De Bakker PI, Daly MJ, et al.
651 PLINK: a tool set for whole-genome association and population-based linkage analyses. The American journal
652 of human genetics. 2007; 81(3):559–575.
- 653 Randall RE, Griffin DE. Within host RNA virus persistence: mechanisms and consequences. Current opinion in
654 virology. 2017; 23:35–42.

- 655 **Rezaee F**, Gibson LF, Piktel D, Othumpangat S, Piedimonte G. Respiratory syncytial virus infection in human bone
656 marrow stromal cells. *Am J Respir Cell Mol Biol.* 2011 Aug; 45(2):277–86. doi: 10.1165/rcmb.2010-0121OC, edition: 2010/10/26.
- 658 **Rixon HM**, Brown G, Murray J, Sugrue R. The respiratory syncytial virus small hydrophobic protein is phosphory-
659 lated via a mitogen-activated protein kinase p38-dependent tyrosine kinase activity during virus infection.
660 *Journal of General Virology.* 2005; 86(2):375–384.
- 661 **Rohwedder A**, Kemerer O, Forster J, Schneider K, Schneider E, Werchau H. Detection of respiratory syncytial
662 virus RNA in blood of neonates by polymerase chain reaction. *J Med Virol.* 1998 Apr; 54(4):320–7. doi:
663 10.1002/(sici)1096-9071(199804)54:4<320::aid-jmv13>3.0.co;2-j, edition: 1998/04/29.
- 664 **Rosas-Salazar C**, Tang ZZ, Shilts MH, Turi KN, Hong Q, Wiggins DA, Lynch CE, Gebretsadik T, Chappell JD,
665 Peebles Jr RS, et al. Upper respiratory tract bacterial-immune interactions during respiratory syncytial virus
666 infection in infancy. *Journal of Allergy and Clinical Immunology.* 2022; 149(3):966–976.
- 667 **Schobel SA**, Stucker KM, Moore ML, Anderson LJ, Larkin EK, Shankar J, Bera J, Puri V, Shilts MH, Rosas-Salazar C,
668 Halpin RA, Fedorova N, Shrivastava S, Stockwell TB, Peebles RS, Hartert TV, Das SR. Respiratory Syncytial Virus
669 whole-genome sequencing identifies convergent evolution of sequence duplication in the C-terminus of the
670 G gene. *Sci Rep.* 2016 May; 6:26311. doi: 10.1038/srep26311, edition: 2016/05/24.
- 671 **Tripp RA**, Jones LP, Haynes LM, Zheng H, Murphy PM, Anderson LJ. CX3C chemokine mimicry by respiratory
672 syncytial virus G glycoprotein. *Nature immunology.* 2001; 2(8):732–738.
- 673 **Wang S**, Sundaram JP, Stockwell TB. VIGOR extended to annotate genomes for additional 12 different viruses.
674 *Nucleic Acids Res.* 2012 Jul; 40(Web Server issue):W186–92. doi: 10.1093/nar/gks528, edition: 2012/06/07.
- 675 **Yang J**, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait analysis. *The American
676 Journal of Human Genetics.* 2011; 88(1):76–82.

677 **Supplemental**

678 **Supplemental host genetic analyses**

679 We further investigated the possibility that the analysis was underpowered to identify associations
680 with reported childhood asthma- and RSV LRTI-associated SNPs (**Pividori et al., 2019; Janssen et al.,
681 2007; Pasanen et al., 2017**). This was done by pooling information across SNPs to estimate the
682 average genetic effect size. In brief, we computed a z-score for each SNP, where the average (across
683 SNPs) squared. As \bar{G} is an average of $p = 54$ approximately independent statistics, it is approximately
684 $N(n\mu^2 + 1, 2/p)$, where $n = 621$ is the sample size and μ^2 is a function of the average squared genetic
685 effect on RSV infection in infancy. Using the genetic effect estimates from **Pividori et al. (2019);
686 Janssen et al. (2007); Pasanen et al. (2017)**, we calculated that we would have 80% power to reject
687 the global null hypothesis of no genetic effect at any of these SNPs (i.e., $\mu^2 = 0$) if, on average
688 across the 54 SNPs, the genetic effect on RSV infection in infancy was at least 61% as large as those
689 estimated in the aforementioned 3 studies. The z-score \bar{G} is proportional to the average squared
690 genetic effect on RSV infection in infancy. We found $\bar{G}=1.00$ in our data, which corresponds to a p
691 value of 0.50. This result indicates that the genetic effect on RSV infection in infancy is zero or small
692 at SNPs likely to be associated with RSV infection a priori.

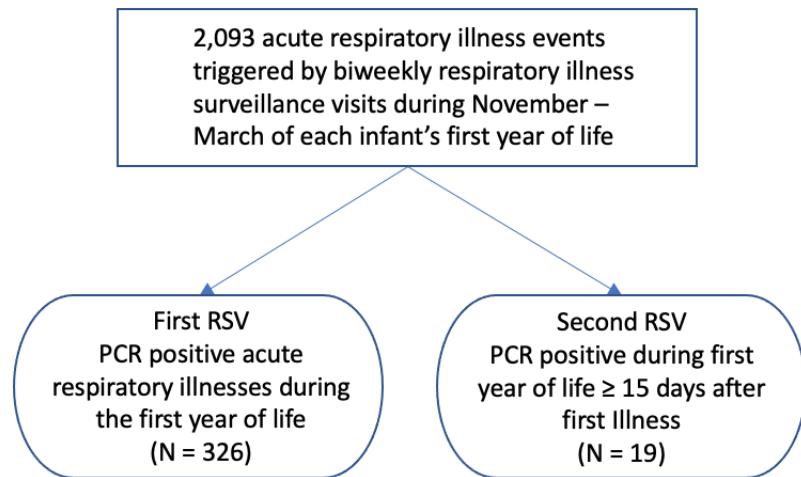


Figure S1. Supplemental: Infant Susceptibility to Pulmonary Infections and Asthma Following RSV Exposure (INSPIRE). The study population is a longitudinal birth cohort specifically designed to capture the first RSV infection in term healthy infants. Prolonged infection was a priori defined as repeatedly meeting criteria for acute respiratory infection with RSV PCR positive nasal samples ≥ 15 days between testing.

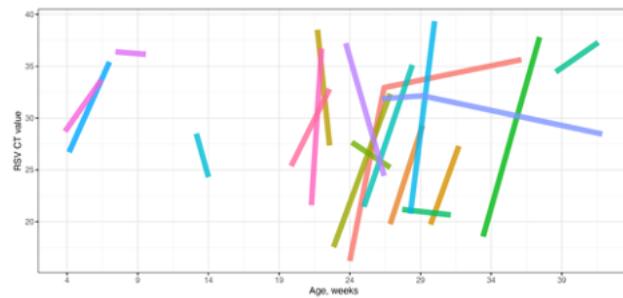


Figure S2. Supplemental: Infant RSV prolonged infections. Each line represents an infant in the study, and line start and end correspond to clinical respiratory illness sampling timepoints. CT values are inversely related to viral RNA abundance.

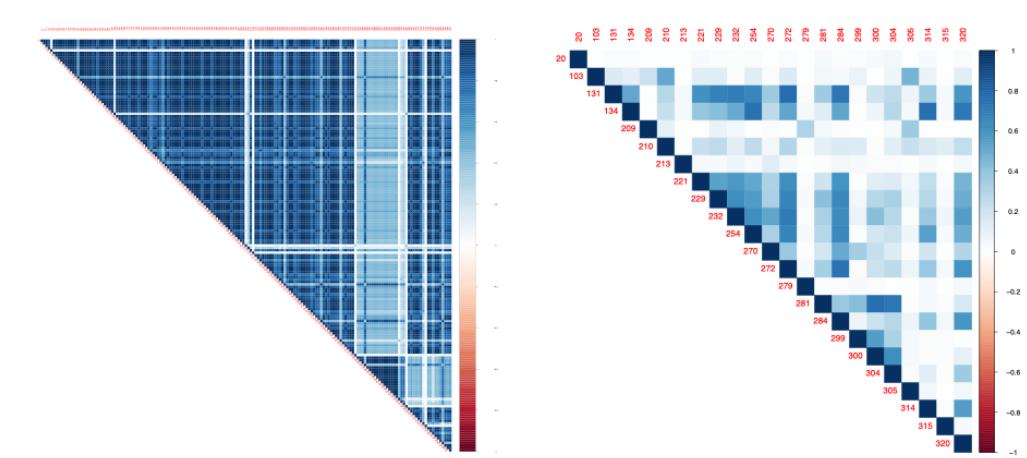


Figure S3. Supplemental: Variant clumping for reduction in association testing. [Left] Correlation between all positions. [Right] Correlation between proxy variants were clumped to remove $r^2 \geq 0.8$. Values indicate relative amino acid positions within MSA. r^2 indicated by color.

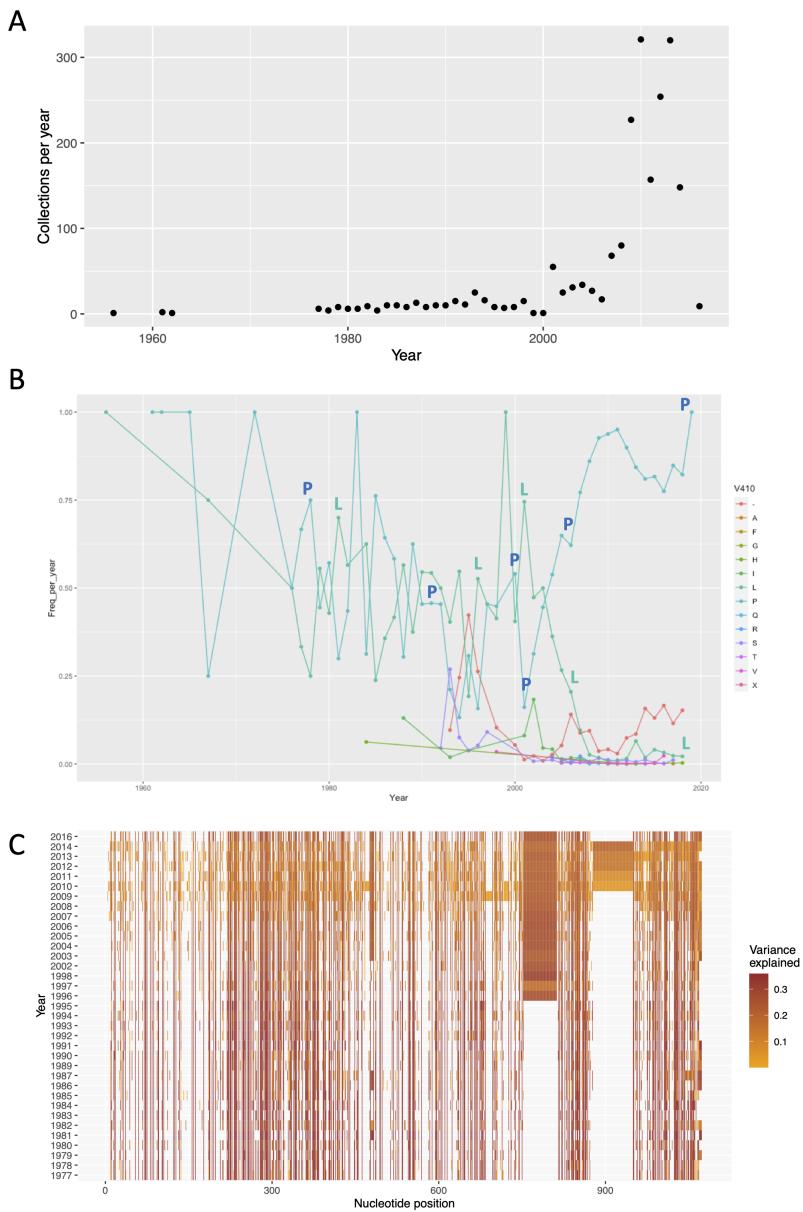


Figure S4. Supplemental: Publicly available RSV sequence data for > 60 years. (A) Global sample collection per year. (B) Variant associated with prolonged infection tracked in public data. The lead proxy SNP, p.P218T/S/L is illustrated here (relative amino acid positive 410 in MSA). The major alleles (proline, leucine) are seen for group A/B, with minor alleles (serine, threonine) generally at low frequency <10%. (C) % variance explained per year for all G protein amino acid variants from 1977-2016 (years with very low coverage removed).