

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

Dylan Lawless, PhD¹, Christopher G. McKennan, PhD², Suman Das, PhD⁴, Thomas Junier, PhD³, Zhi Ming Xu, MSc¹, Larry J Anderson, MD⁵, Tebeb Gebretsadik, MPH⁶, Meghan Shilts, MHS, MS⁷, Emma Larkin, PhD⁸, Christian Rosas-Salazar, MD, MPH⁸, James D. Chappell, MD⁹, Jacques Fellay, MD, PhD^{1,3,10}, and Tina V. Hartert, MD, MPH^{7,9}

¹Global Health Institute, School of Life Sciences, École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland, ²Department of Statistics, University of Pittsburgh, Pittsburgh,

Pennsylvania, United States of America, ³Swiss Institute of Bioinformatics, Vital-IT Group, Switzerland, ⁴Division of Infectious Diseases, Department of Medicine, Vanderbilt University Medical Center, Nashville, Tennessee, United States of America, ⁵Department of Pediatrics, Emory University School of Medicine, Atlanta, Georgia, United States of America,

⁶Department of Biostatistics, Vanderbilt University Medical Center, Nashville, Tennessee, United States of America, ⁷Department of Medicine, Vanderbilt University Medical Center, Nashville, Tennessee, United States of America, ⁸Division of Allergy, Immunology, and Pulmonary and Critical Care Medicine, Department of Medicine, Vanderbilt University Medical Center, Nashville, Tennessee, United States of America, ⁹Department of Pediatrics, Vanderbilt University Medical Center, Nashville, Tennessee, United States of America,

¹⁰Precision Medicine Unit, Lausanne University Hospital and University of Lausanne, Lausanne, Switzerland

Abbreviations

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

Notice of Prior Presentation

The results of the host genome wide association study analyses included in this manuscript were presented during the European Society of Human Genetics Conference in June 2020 in

Berlin, Germany, which was held remotely [1].

Ethics Statement for Human Subjects Research

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

Competing interests

All authors have completed a conflict of interest form (COI). There were no COI. Funding was supplied from National Institutes of Health and Swiss National Science Foundation.

U19 AI 095227, UG3/UH3 OD023282, UL1 TR002243, SNSF IZSEZ0_191968 (TVH), SNSF 310030L_197721 (JF), X01 HLG244 RS&G (EL).

1 Abstract

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

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

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

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

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

2 Introduction

Human orthopneumovirus, formerly known (and frequently still referred to) as Respiratory syncytial virus (RSV), results in significant global morbidity and mortality [2]. By the age of two to three years, nearly all children are infected with RSV at least once [3]. RSV is a seasonal mucosal pathogen that primarily infects upper and lower respiratory tract epithelium, although it has been recovered from non-airway sources [4–9]. While RSV is mainly associated with acute respiratory infection, many RNA viruses can establish prolonged or persistent infection in some infected individuals [10]. Prolonged shedding of RSV, especially in young infants and following first infection, has been demonstrated, with longer average duration of viral shedding when polymerase chain reaction (PCR) is used to detect RSV [11]. While younger

age and first infection are associated with protracted infection [3; 12], it is not known whether specific viral factors contribute to prolonged RSV infection in infants. This is important, as prolonged infection may contribute to enhanced transmission and developmental changes to the early life airway epithelium. Further, the reservoir of RSV infection is not understood, and it is possible that some RSV strains sustain a low levels of ongoing viral circulation in the community until seasonal or other influences favor epidemic spread [13].

The objectives of this study were therefore to determine if there exist host genetic risk alleles for RSV infection and to identify viral genetic variation associated with prolonged infection. These motivating questions are of fundamental interest in understanding viral and host genetic contributions that may underlie the development of chronic respiratory morbidity due to RSV, including asthma.

3 Methods

3.1 Study population

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

The study population is a longitudinal birth cohort, the INFant Susceptibility to Pulmonary Infections and Asthma Following RSV Exposure (INSPIRE), specifically designed to capture the first RSV infection in term healthy infants. Additional details of this birth cohort have been previously published [14]. Briefly, the cohort included 1949 term (≥ 37 weeks gestation), non-low birth weight (≥ 2250 g, 5 lbs), otherwise healthy infants from a population-representative sample of pediatric practices located in rural, suburban, and urban regions of the south-eastern US during 2012-2014. Infants were born June through December; per study design, they were 6 months of age or less entering their first RSV season.

3.2 Biweekly surveillance of RSV infection

Infant (i.e., first year of life) RSV infection was ascertained through passive and active bi-weekly surveillance during each infant's first RSV season and RSV serology (Table 1). If an infant met pre-specified criteria for an acute respiratory infection, we conducted an in-person respiratory illness visit at which time we administered a parental questionnaire, performed a physical exam, collected a nasal wash, and completed a structured medical chart review for

infants seen during an unscheduled visit. RSV RNA in nasal samples was detected by reverse-transcription quantitative PCR [15]. We a priori defined the clinical entity of “prolonged” infection during infancy as repeatedly meeting pre-specified criteria for an acute respiratory infection accompanied by repeatedly positive RSV PCR separated by 15 or more days (Figure S1) [13].

3.3 Descriptive analyses

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

3.4 Host DNA collection and genotyping

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

3.5 Host genetic analyses of RSV infection in infancy

To determine whether host genetic factors associate with infant RSV infection risk, we examined single nucleotide polymorphisms (SNPs) previously shown to alter infant RSV lower respiratory tract infection (LRTI) or childhood asthma risk [16–18]. We also conducted a host GWAS to identify common variants associated with infant RSV infection, and examined narrow sense heritability to test for small cumulative effects. The GWAS was performed on 621 children with available DNA for the association between host genotype and RSV infection during infancy. Due to sample size constraints, we restricted our sub-analysis to the 54 host SNPs previously associated with RSV lower respiratory tract infection or childhood asthma [16–18]. We additionally evaluated the accumulation of small genetic effects that would go undetected in a GWAS by estimating the narrow sense heritability of RSV infection.

For GWAS analyses, the initial round of data quality control was performed on individual populations (self-reported as White, Black, and Hispanic) using PLINK version 1.9 [19]. Subjects with a missing genotype call rate above 5% were removed. The SNP minor allele frequency (MAF) threshold was set at > 0.01, 0.03, and 0.08 for White, Black, and Hispanic, respectively [20].

The groups were merged for a total of 1,086,830 variants and a genotyping rate of 0.78.

Subject independence was assessed using KING (<https://people.virginia.edu/~wc9c/KING/>) to prevent spurious associations. However, no probable relatives or duplicates were detected based on pairwise identify-by-state. We compared the genetic ancestry in cases to self-reported ethnicity to check for mislabelling. Reported and estimated sex was also examined for discrepancy. A second round of quality control on the combined dataset was conducted, which removed 74 samples due to genotype missingness and 399,991 variants with a genotyping rate ≤ 0.1 . Variants were checked for departure from Hardy-Weinberg equilibrium (HWE) ($P < 1e^{-6}$) to uncover features of selection, population admixture, cryptic relatedness, or genotyping error. This was only performed on controls to prevent removal of genuine genetic associations that can be associated with this measurement; 6,024 variants were removed. No variants had a MAF $MAF < 0.01$ after merging. SNP positions and identifiers were compared and updated according to dbNSFP4.0a (hg19) with 289 variants removed due to a missing coordinate and SNPs identifier [21]. This resulted in an analysis-ready dataset of 680,526 variants from 621 children (509 with and 112 without RSV infection in infancy), yielding a total genotyping rate of 0.98. No genomic inflation was evident with an estimated lambda (based on median chi-squared test) equal to 1. We then used genome-wide complex trait analysis (GCTA) software (<https://cnsgenomics.com/software/gcta/>) to calculate the genetic relationship matrix and performed principal component analysis (PCA) to account for population structure [20]. Genome-wide association analysis was performed using PLINK version 1.9 for logistic regression with multiple covariates consisting of the child's birth month, enrolment year (as a marker of RSV season), daycare attendance, presence of another child ≤ 6 years of age at home, sex, and 6 ancestry principal components (PCs) [19].

As the multiple testing burden likely precluded identification of small genetic effects in our GWAS, we conducted an additional heritability analysis using the method described by Golan et al. [22] to estimate narrow-sense heritability of RSV infection during infancy on the latent liability scale (h_l^2), which, > 0 , would indicate an accumulation of small genetic effects. We estimated h_l^2 to be exactly 0, suggesting that, if present, infant RSV infection-related genetic signals are both small and sparse.

3.6 RSV whole-genome sequencing

RSV genome sequencing was performed on all specimens from subjects meeting illness criteria and with positive RSV PCR. Viral amino acid variants (genotype) of the F and G glycoprotein were tested for association with prolonged infection adjusting for host features associated with increased infection risk. We focused our analyses on the surface F (fusion) and G (attachment) proteins of RSV as they have been implicated in pathogenesis [23; 24], and both are targets for neutralizing antibodies during infection [25; 26]. Lastly, to determine if the variants of interest were enriched by selective pressure over time, we used public data from the past several decades to assess variant frequency over time.

RSV whole-genome sequencing of this study population has been previously described [27]. Briefly, RNA was extracted at J. Craig Venter Institute (JCVI) (<https://www.jcvi.org>) in Rockville, MD from nasal wash samples which were RSV PCR positive and collected during a respiratory illness visit triggered through biweekly surveillance of symptoms. Four forward reverse-transcription (RT) primers were designed and four sets of PCR primers were manually picked from primers designed across a consensus of complete RSV genome sequences using JCVI's automated primer design tool [28]. cDNA was generated from 4 µL undiluted RNA, using the pooled forward primers and SuperScript III Reverse Transcriptase (Thermo Fisher Scientific, Waltham, MA, USA). 100 ng of pooled DNA amplicons were sheared to create 400-bp libraries, which were pooled in equal volumes and cleaned with Ampure XP reagent (Beckman Coulter, Inc., Brea, CA, USA). Sequencing was performed on the Ion Torrent PGM using 316v2 or 318v2 chips (Thermo Fisher Scientific).

For samples requiring extra coverage, in addition to Ion Torrent sequencing, Illumina libraries were prepared using the Nextera DNA Sample Preparation Kit (Illumina, Inc., San Diego, CA, USA). Sequence reads were sorted by barcode, trimmed, and assembled de novo using CLC Bio's *clc_novo_assemble* program, and the resulting contigs were searched against custom, full-length RSV nucleotide databases to find the closest reference sequence. All sequence reads were then mapped to the selected reference RSV sequence using CLC Bio's *clc_ref_assemble_long* program [29]. Curated assemblies were validated and annotated with the viral annotation 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 (<https://www.ncbi.nlm.nih.gov/bioproject/225816>) [30] and PRJNA267583 (<https://www.ncbi.nlm.nih.gov/bioproject/267583>).

3.7 Viral sequence alignment

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

We computed a phylogenetic tree for each gene, as follows. NCBI-tools asn2fsa (<https://www.huge-man-linux.net/man1 asn2fsa.html>) was used to convert sequences to fasta format. Each sample consisted of 11 sequence segments (NS1, NS2, N, P, M, M2-1, M2-2, SH, G, F, and L) as shown in Figure S1. These were separated and repooled to create 11 single fasta files for each gene containing all 350 samples. Sequences were checked for at least 90% coverage of the corresponding gene to minimize loss of aligned positions when computing the phylogenetic tree. Each of the eleven resulting sets was aligned with MAFFT v7 (<https://mafft.cbrc.jp/alignment/software/>) [31], using default parameters. The sequence of the

orthologous gene from *Bovine orthopneumovirus* ([GenBank:NC_001989](#)) was added to each set as an outgroup.

IQ-Tree (<https://www.iqtree.org>) [32] was used with per-gene multiple sequence alignment (MSA) files based on amino acid sequence for estimating maximum-likelihood phylogenies using protein substitution model. Examining the sequences with an alignment viewer showed that a small number of sequences had frame-shift variants but these did not affect the regions included in our testing criteria.

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

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

A number of host features have been previously shown to influence infection susceptibility and were therefore included as covariates in our analysis [33]. Six samples were excluded due to insufficient covariate data, resulting in 344 test samples. Of these, 38 were from the same patients (prolonged infection) of which half (19) were included for association testing. Thus, the test set was comprised of single samples collected from 325 individuals.

3.8 Viral population structure

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

3.9 Viral variant association testing

Viral amino acids (genotype collapsed into REF/ALT) were tested for association with infection types (i.e., resolved and prolonged) including key covariates that alter infection risk. To reduce the multiple testing burden, proxy amino acid variants were identified by performing

clumping with ranking based on MAF and with a cut-off threshold of $r^2 \geq 0.8$ (Supplemental Figure S3). Since many variants within RSV coding genes have non-random association due to selection, analogous to linkage disequilibrium (LD) in human GWAS, we reduced the multiple testing burden by retaining proxy variants and removing those with $r^2 \geq 0.8$. Analysis was performed using logistic regression with the R stats (3.6.2) glm function. The model consisted of the binary response (prolonged infection Yes/No) and predictors viral genotype (REF/ALT amino acid, including multi-allelic non-REF collapsed into ALT), viral PCs 1-5, host sex, and host features that have been previously demonstrated as significantly associated with infection: self-reported race/ethnicity, daycare attendance, and living with siblings [33].

Environmental host covariates did not contribute significant effect in our model for candidate causal association. For this reason, in our main analysis, viral population structure was accounted for by the first five PCs. The Bonferroni correction for multiple testing was applied based on the number of variants tested. For the significant association found by proxy amino acid variants, the association model was repeated for all clumped variants to produce a LocusZoom-style Manhattan plot containing r^2 by color and p value statistics. R package stats was used for a range of analyses including glm for logistic regressions. R package MASS was used to analyse logistic regression model data. To test if the significantly associated variants were due to population structure, we re-estimated models using the subset of individuals infected with RSV strain B to confirm validity of combined analysis.

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

3.10 Public viral sequence data

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

3.11 Biological interpretation

As infant RSV infection stimulates an acute antiviral response and also results in decreased barrier function of the airway epithelium [34], we tested for association between host interferon (IFN) response and the amino acid (REF/ALT) identified as the viral variant associated with prolonged infection. A Wilcoxon test was performed to compare IFN- γ , and IFN- α , between RSV amino acid positions, with adjustment for the same covariates as in the main analysis. Protein structures were analysed with data sourced from RCSB PDB <https://www.rcsb.org>. Protein function and domains were assessed using UniProt (<https://www.uniprot.org>) for P03423 (GLYC_HRSVA) (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 modifications, motifs, and epitopes were assessed from the literature. Protein features were assessed using data from NCBI (https://www.ncbi.nlm.nih.gov/igp/NP_056862.1) and via sequence 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 according to available information on a broad range of biological and biochemical features, including native conformation (secondary, tertiary, and quaternary), domains and topology, disulfide bonds, glycosylation, interactions with other viral proteins and host-cell factors, proteolytic cleavage sites, normal patterns of intra-and/or extra-cellular distribution, and secretion status.

4 Results

4.1 Cohort characteristics

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

4.2 Host genetic analyses

We explored whether RSV infection in infancy is a natural assignment (quasi-random) event and, unlike severity of early-life RSV infection [35], occurs independently of host genetics. For the candidate SNP analysis, we considered childhood asthma- and RSV LRTI-associated SNPs identified in Janssen et al. [16]; Pasanen et al. [17]; Pividori et al. [18]. The first is the largest childhood asthma GWAS to date, and, to our knowledge, the latter 2 represent the most comprehensive studies of RSV LRTI-associated SNPs. To further reduce the multiple testing burden, we only analysed SNPs with MAF ≥ 0.1 in at least one of the White, Black, or Hispanic ethnicity groups. Associations between genotype at the resulting 54 SNPs (50 childhood asthma- and 4 RSV LRTI-associated SNPs) and RSV infection in infancy in our data are given in Figure 1. The data are consistent with little to no effect of genotype at these SNPs on RSV infection in infancy.

We further investigated the possibility that the analysis was underpowered to identify associations with these SNPs by pooling information across SNPs to estimate the average genetic effect size. Our analysis in the supplement shows that the average genetic effect of these LRTI- and asthma-related SNPs on infant RSV infection is zero or trivially small (Supplemental 8).

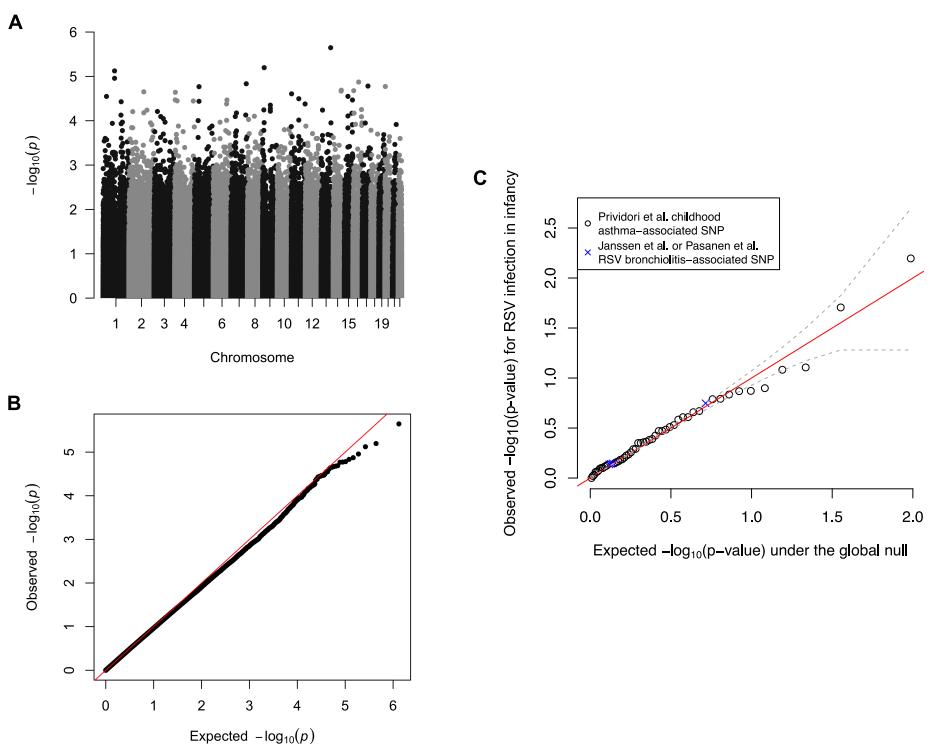


Figure 1: **Genetic analyses of RSV infection in infancy.** (A) The Manhattan plot shows no genome-wide significant associations (p value threshold of $5e^{-8}$). (B) The Q-Q plot demonstrates that the observed p values are congruent with those expected under the null hypothesis that RSV infection in 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.

4.3 Population structure

A summary of protein coding genes in RSV is illustrated in Figure 2 A. Our analysis focused on F and G protein. The phylogenetic tree based on multiple sequence alignment (MSA) of G protein amino acid sequences is shown in Figure 2 B. One obvious feature causing a separation in genetic diversity is G protein partial gene duplication, which has emerged in recent years within RSV-A strains [36]. RSV-B strains with an analogous duplication have existed for two decades, although the selection process leading to emergence and clinical implications have not been entirely defined.

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

4.4 Genetic invariance of prolonged infection

The duration of RSV shedding in Kenyan infants has been reported previously [13]. Based on these findings, infections separated by at least 15 days with symptoms were expected to be “new” infections [13]. Figure 2 D (panel [i]) summarizes every pairwise genetic distance between every viral sequence, where small distances indicate pairs with closely related sequences. Panels [ii] and [iii], which summarize the difference in sequence similarity distributions between viruses from the same host and different hosts, show that RSV sequences corresponding to initial and subsequent viral detections are nearly identical. These results support the conclusion that such cases are prolonged (i.e., failure to clear) infections rather than new infections.

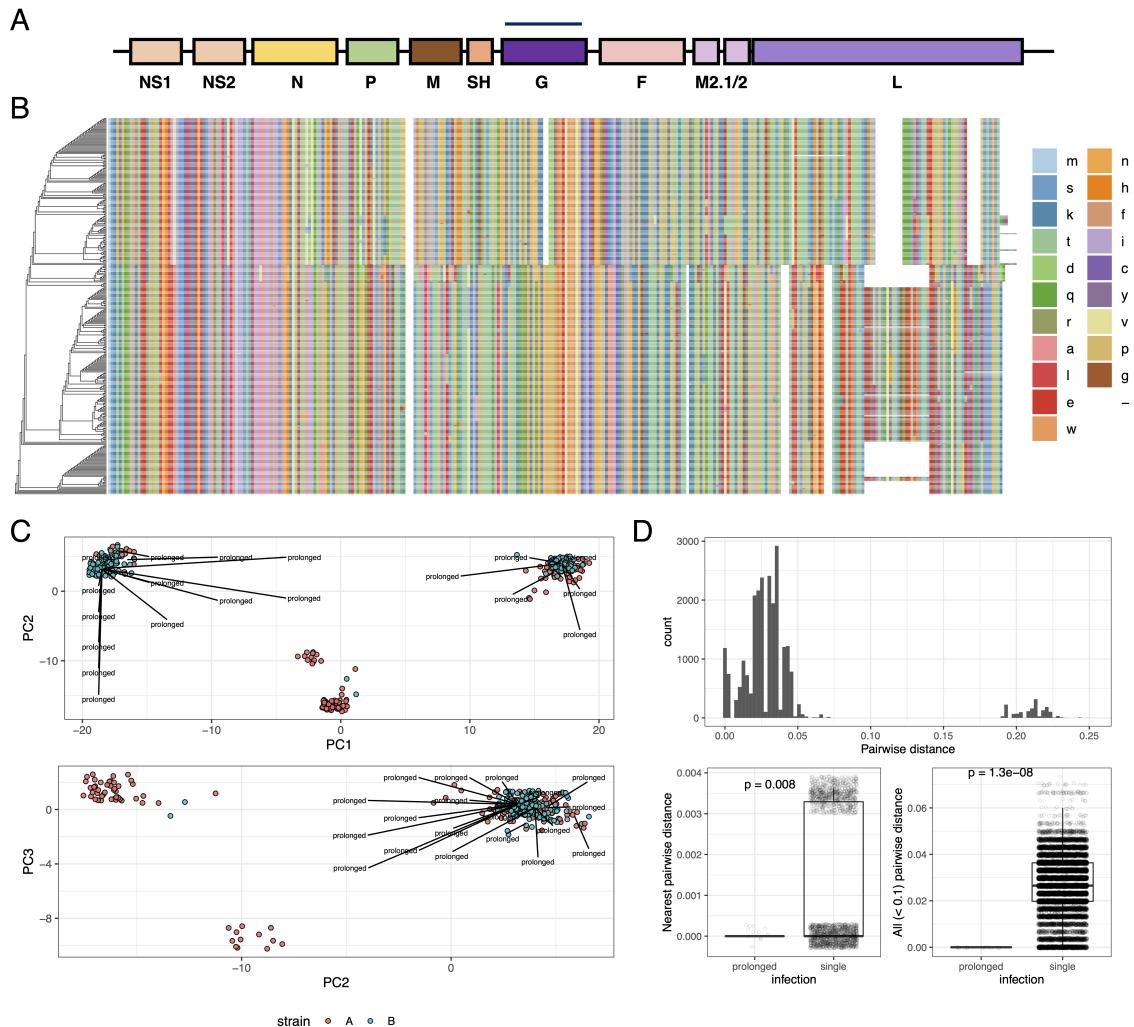


Figure 2: Viral population structure. (A) Linear map of the RSV genome. (B) Phylogenetic tree based on multiple sequence alignment MSA of G protein amino acid sequences. Color; amino acids. (C) Principal component (PC) analysis. PCs1-3 with labels indicating prolonged infections from different phylogenetic 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.

4.5 Variants in G glycoprotein significantly associated with prolonged infection

The consensus sequence within the cohort was assigned based on the major allele. Variants at the amino acid level were defined as either reference (REF) or alternative (ALT) and assessed for their association with prolonged infection. The model consisted of (A) the binary response (prolonged infection Yes/No), and (B) predictors: (1) viral genotype (REF/ALT amino acid), (2) viral PCs 1-5, (3) host sex, and host features that have been previously demonstrated as significantly associated with infection; (4) self-reported race/ethnicity, (5) child-care attendance, or living with another child \leq 6 years of age at home [37]. A significant genetic association was identified between prolonged infection and the lead variant after Bonferroni correction for multiple testing (threshold for number independent variants $< 0.05/23 = 0.002$), as shown in Figure 3 A, p value = 0.0006.

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

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

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

To determine whether p.E123K/D and p.P218T/S/L variant genotypes are novel and potentially influence viral fitness, we searched the public viral data repository of NCBI Human orthopneumovirus, taxid:11250, which contained data from 27 countries worldwide, sample collection dates from 1956 onward, and 1084 glycoprotein protein sequences after curation. The variants were present at a low and stable frequency, without obvious temporal enrichment (Supplemental Figure S4). Thus, while historical data reveal no positive selective advantage attached to p.E123K/D and p.P218T/S/L, longstanding circulation and linkage in prolonged

RSV infection suggest that these polymorphisms are present in the viral inoculum and do not arise through recurrent mutational events.

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

4.6 Functional interpretation

Cell-attachment proteins of paramyxoviruses (G protein in RSV) span the viral envelope and form spike-like projections from the virion surface. RSV G protein is a type II integral membrane protein consisting of 298 amino acid residues comprising N-terminal cytoplasmic (p.1-43), transmembrane helical (p.43-63), and extracellular (p.64-298) domains (Figure 3 D). RSV G protein ectodomain also exists in a soluble secreted form, p.66 – 298, which functions in immune evasion [38–40]. G protein interacts with the small hydrophobic (SH) protein [41] and, via the N-terminus, with matrix (M) [42] protein. It has also been reported to form homo-oligomers [43]. The variant amino acid positions associated with prolonged infection reside in a portion of the ectodomain of unassigned specific function and linearly non-contiguous with sequences that bind cell-surface heparan sulfate, which likely promotes RSV cell-attachment (p.187-198) [38–40]. In addition, these positions do not contribute to known neutralization epitopes on G protein. Information available in PDB was insufficient to infer effects of p.E123K/D and p.P218T/S/L on local or regional protein structure. The potential effect on glycosylation is indeterminate. Figure 3 D illustrates the position (dotted red lines) of these variants relative to summarised known functional features.

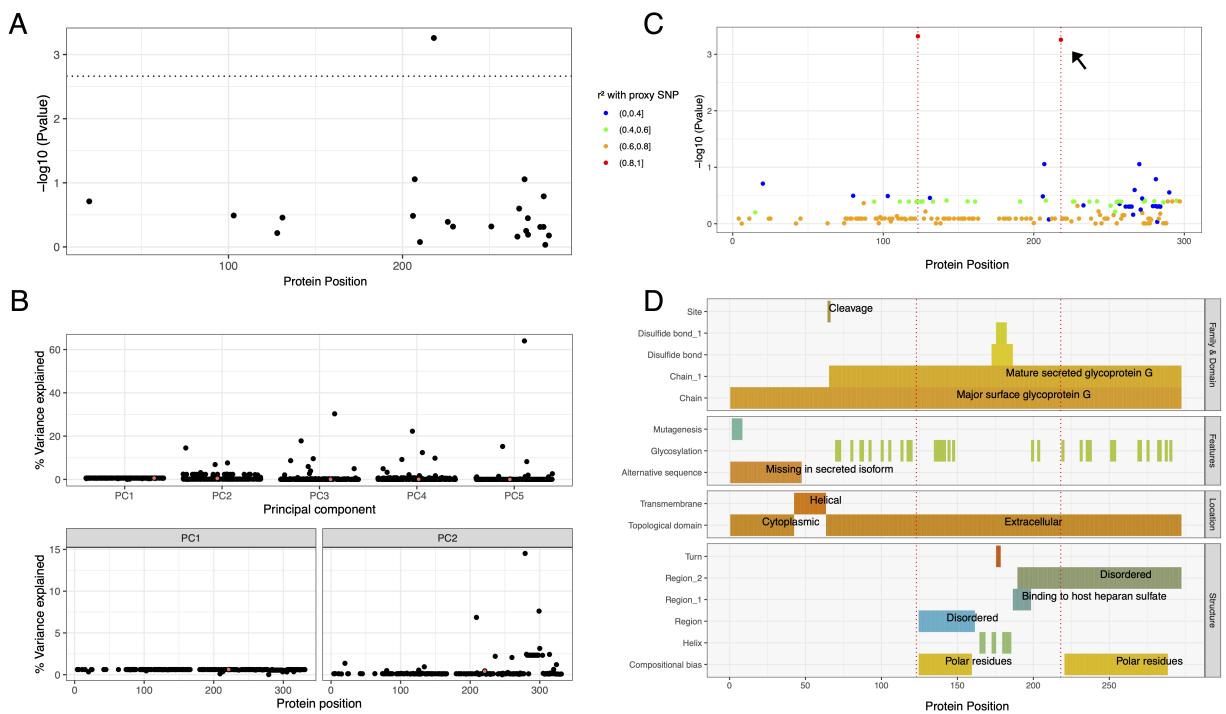


Figure 3: Viral genetic association with prolonged infection. (A) Amino acid association with prolonged infection after multiple testing correction (significant threshold shown by dotted line). (B) Variance explained (VE) within cohort. The effect of each variant on cohort structure is shown for PCs 1-2. 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.

4.7 Host response

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

5 Discussion

In this study of term healthy infants, we found no evidence of host genetic susceptibility to RSV infection during infancy. This allowed our analysis to focus on elucidation of viral drivers of prolonged infection. A significant viral genetic association in the RSV G protein, p.E123K/D and p.P218T/S/L, with prolonged infant RSV infection was identified. These variants were not associated with severe disease, and public data reveal their consistent presence at low frequencies over the past 30 years, without evidence of enrichment by positive selective pressure over time. The two variants we identified in G are correlated with non-random association, analogous to LD in the human diploid genome and therefore not likely random mutations, but instead co-inherited in the infecting inoculum. This suggests an evolutionary benefit and raises the question of why such variants have maintained a stable but low frequency in the human population for decades. These strains are a potential reservoir, emerging seasonally in response to immune, environmental, or other forces. Alternatively, the polymorphisms might recurrently arise de novo during infection of some individuals but are poorly transmissible because of suboptimal fitness. The possibility of viral mutational immune escape has been reported for infants who struggle to control primary RSV infections, allowing for prolonged viral replication and not previously described viral rebound [44].

The RSV variants associated with prolonged infection in our cohort, G p.E123K/D and p.P218T/S/L, lie in the extracellular region, and there are no known mechanistic features that directly overlap, although it is possible that variant positions approximate sequences that bind a putative viral receptor, heparan sulfate [39], in the G protein three-dimensional structure. G protein amino acid positions 123 and 218 are not part of known antibody neutralization epi-

topes or CD8+ cytotoxic T-cell epitopes (Figure 3 D). In addition to heparan sulfate, interactions between viral G protein and CX3CR1, the receptor for the CX3C chemokine fractalkine, have been reported to modulate the immune response and facilitate infection [38–40; 45–47]. Furthermore, the mature secreted isoform of G protein (p.66-298) is thought to facilitate viral antibody evasion by acting as an antigen decoy and modifying the activity of leukocytes bearing Fc-gamma receptors [48]. Our findings raise the interesting prospect that G protein variants associated with prolonged infection alter a key interaction at the immune interface between pathogen and host.

Although this study was not designed to define mechanisms underlying the association of G protein variants with prolonged infection, these sequence changes might dampen antiviral immune responses and thereby delay viral clearance. Although we observed differences in the acute antiviral response between subjects with resolved and prolonged infection, specifically increased levels of types 1 and 2 IFN in nasal secretions, we could not make causal inference about variant sequences because of confounding by co-linearity of these polymorphisms with RSV antigenic group. Results of nasal cytokine analysis are nevertheless consistent with a contemplated role for altered immune responses in extended infections by G protein variant strains. It is also possible that strains harbouring G protein p.E123K/D and p.P218T/S/L variants are cleared more slowly and foster an immune environment of low-level chronic stimulation or exhaustion. We previously demonstrated that infants infected with RSV in their first year of life have dampened subsequent antiviral immune responses in early childhood [49] as well as changes in airway epithelial cell metabolism [34].

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

we cannot exclude modulatory effects of maternal antibody, which we did not measure.

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

6 Links

6.1 Software

R v4.1.0 was used for data preparation and analysis <http://www.r-project.org>.

R package *caret* was used for analysis: genetic correlations.

R package *dplyr* was used for data curation.

R package *factoextra* was used for analysis: PCA, and to visualise eigenvalues and variance.

R package *ggplot2* was used for data visualisation.

R package *MASS* was used to analysis: logistic regression model data.

R package *stats* was used for analysis: including *glm* for logistic regressions.

R package *stringr* was used for data curation.

R package *tidyR* was used for data curation.

asn2fsa <https://www.huge-man-linux.net/man1 asn2fsa.html>

clc_novo_assemble qiagenbioinformatics.com

Clustal Omega <https://www.ebi.ac.uk/Tools/msa/clustalo/>

dbNSFP (database) <http://database.liulab.science/dbNSFP> [21]

GCTA <https://cnsgenomics.com/software/gcta/> [20]

GenBank <https://www.ncbi.nlm.nih.gov/genbank/>

IQ-Tree <https://www.iqtree.org/> [32]

KING <https://people.virginia.edu/~wc9c/KING/> [50]

MAFFT <https://mafft.cbrc.jp/alignment/software/> [31]

NextAlign <https://github.com/nextstrain/nextclade>

PLINK <http://zzz.bwh.harvard.edu/plink/> [19]
Tbl2asn <https://www.ncbi.nlm.nih.gov/genbank/tbl2asn2/>
Viral Genome ORF Reader, VIGOR 3.0 <https://sourceforge.net/projects/jcvi-vigor/files/>
RCSB PDB <https://www.rcsb.org>
UniProt <https://www.uniprot.org>

6.2 Data sources

Dataset <https://www.ncbi.nlm.nih.gov/bioproject/267583>.
Dataset <https://www.ncbi.nlm.nih.gov/bioproject/225816>.
J. Craig Venter Institute <https://www.jcvi.org>.
GenBank:NC_001989 *Bovine orthopneumovirus*, complete genome https://www.ncbi.nlm.nih.gov/nuccore/NC_001989.
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.
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.
Reference data for all public NCBI Virus <https://www.ncbi.nlm.nih.gov/labs/virus/vssi/> for species: *Human orthopneumovirus*; genus: *Orthopneumovirus*; family: *Pneumoviridae*.
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.
Reference https://www.ncbi.nlm.nih.gov/protein/NP_056862.1
GCF_002815475.1 (release 2018-08-19) Nucleotide Accessions: NC_038235.1, protein: Y_009518856.1
Reference https://www.ncbi.nlm.nih.gov/protein/YP_009518856.1
GCF_000855545.1 (release 2015-02-12) Nucleotide Accessions: NC_001781.1, protein: NP_056862.1 (strain B1).

7 Code availability

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

References

- [1] D Lawless, C Rosas-Salazar, T Gebretsadik, K Turi, B Snyder, P Wu, J Fellay, and T Hartert. Genome-wide association study of susceptibility to respiratory syncytial virus infection during infancy. In *European Journal of Human Genetics*, volume 28, pages 319–319. Springer Nature Campus, 4 Crinan St, London, N1 9XW, England, 2020.
- [2] C. B. Hall, G. A. Weinberg, M. K. Iwane, A. K. Blumkin, K. M. Edwards, M. A. Staats, P. Auinger, M. R. Griffin, K. A. Poehling, D. Erdman, C. G. Grijalva, Y. Zhu, and P. Szilagyi. The burden of respiratory syncytial virus infection in young children. *N Engl J Med*, 360(6):588–98, February 2009. ISSN 0028-4793 (Print) 0028-4793. doi: 10.1056/NEJMoa0804877. Edition: 2009/02/07.
- [3] W Paul Glezen, Larry H Taber, Arthur L Frank, and Julius A Kasel. Risk of primary infection and reinfection with respiratory syncytial virus. *American journal of diseases of children*, 140(6):543–546, 1986.
- [4] V. Bokun, J. J. Moore, R. Moore, C. C. Smallcombe, T. J. Harford, F. Rezaee, F. Esper, and G. Piedimonte. Respiratory syncytial virus exhibits differential tropism for distinct human placental cell types with Hofbauer cells acting as a permissive reservoir for infection. *PLoS One*, 14(12):e0225767, 2019. ISSN 1932-6203. doi: 10.1371/journal.pone.0225767. Edition: 2019/12/04.
- [5] H. A. Cubie, L. A. Duncan, L. A. Marshall, and N. M. Smith. Detection of respiratory syncytial virus nucleic acid in archival postmortem tissue from infants. *Pediatr Pathol Lab Med*, 17(6):927–38, November 1997. ISSN 1077-1042 (Print) 1077-1042. Edition: 1997/11/14.
- [6] D. Nadal, W. Wunderli, O. Meurmann, J. Briner, and J. Hirsig. Isolation of respiratory syncytial virus from liver tissue and extrahepatic biliary atresia material. *Scand J Infect Dis*, 22(1):91–3, 1990. ISSN 0036-5548 (Print) 0036-5548. doi: 10.3109/00365549009023125. Edition: 1990/01/01.
- [7] D. R. O'Donnell, M. J. McGarvey, J. M. Tully, I. M. Balfour-Lynn, and P. J. Openshaw. Respiratory syncytial virus RNA in cells from the peripheral blood during acute infection. *J Pediatr*, 133(2):272–4, August 1998. ISSN 0022-3476 (Print) 0022-3476. doi: 10.1016/s0022-3476(98)70234-3. Edition: 1998/08/26.
- [8] F. Rezaee, L. F. Gibson, D. Piktel, S. Othumpangat, and G. Piedimonte. Respiratory syncytial virus infection in human bone marrow stromal cells. *Am J Respir Cell Mol Biol*, 45(2):277–86, August 2011. ISSN 1044-1549 (Print) 1044-1549. doi: 10.1165/rcmb.2010-0121OC. Edition: 2010/10/26.
- [9] A. Rohwedder, O. Keminer, J. Forster, K. Schneider, E. Schneider, and H. Werchau. Detection of respiratory syncytial virus RNA in blood of neonates by polymerase chain

- reaction. *J Med Virol*, 54(4):320–7, April 1998. ISSN 0146-6615 (Print) 0146-6615. doi: 10.1002/(sici)1096-9071(199804)54:4(320::aid-jmv13)3.0.co;2-j. Edition: 1998/04/29.
- [10] Richard E Randall and Diane E Griffin. Within host rna virus persistence: mechanisms and consequences. *Current opinion in virology*, 23:35–42, 2017.
- [11] P. K. Munywoki, D. C. Koech, C. N. Agoti, N. Kibirige, J. Kipkoech, P. A. Cane, G. F. Medley, and D. J. Nokes. Influence of age, severity of infection, and co-infection on the duration of respiratory syncytial virus (RSV) shedding. *Epidemiol Infect*, 143(4):804–12, March 2015. ISSN 0950-2688 (Print) 0950-2688. doi: 10.1017/s0950268814001393. Edition: 2014/06/06.
- [12] Bindya Bagga, L Harrison, P Roddam, and JP DeVincenzo. Unrecognized prolonged viral replication in the pathogenesis of human rsv infection. *Journal of Clinical Virology*, 106: 1–6, 2018.
- [13] Emelda A Okiro, Lisa J White, Mwanajuma Ngama, Patricia A Cane, Graham F Medley, and D James Nokes. Duration of shedding of respiratory syncytial virus in a community study of kenyan children. *BMC infectious diseases*, 10(1):1–7, 2010.
- [14] E. K. Larkin, T. Gebretsadik, M. L. Moore, L. J. Anderson, W. D. Dupont, J. D. Chappell, P. A. Minton, R. S. Peebles, Jr., P. E. Moore, R. S. Valet, D. H. Arnold, C. Rosas-Salazar, S. R. Das, F. P. Polack, and T. V. Hartert. Objectives, design and enrollment results from the Infant Susceptibility to Pulmonary Infections and Asthma Following RSV Exposure Study (INSPIRE). *BMC Pulm Med*, 15:45, April 2015. ISSN 1471-2466. doi: 10.1186/s12890-015-0040-0. Edition: 2015/05/30.
- [15] Emma K Larkin, Tebeb Gebretsadik, Martin L Moore, Larry J Anderson, William D Dupont, James D Chappell, Patricia A Minton, R Stokes Peebles, Paul E Moore, Robert S Valet, et al. Objectives, design and enrollment results from the infant susceptibility to pulmonary infections and asthma following rsv exposure study (inspire). *BMC pulmonary medicine*, 15(1):1–12, 2015.
- [16] Riny Janssen, Louis Bont, Christine LE Siezen, Hennie M Hodemaekers, Marieke J Ermers, Gerda Doornbos, Ruben van't Slot, Ciska Wijmenga, Jelle J Goeman, Jan LL Kimpen, et al. Genetic susceptibility to respiratory syncytial virus bronchiolitis is predominantly associated with innate immune genes. *Journal of Infectious Diseases*, 196(6): 826–834, 2007.
- [17] Anu Pasanen, Minna K Karjalainen, Louis Bont, Eija Piippo-Savolainen, Marja Ruotsalainen, Emma Goksör, Kuldeep Kumawat, Hennie Hodemaekers, Kirsi Nuolivirta, Tuomas Jartti, et al. Genome-wide association study of polymorphisms predisposing to bronchiolitis. *Scientific reports*, 7(1):1–9, 2017.

- [18] Milton Pividori, Nathan Schoettler, Dan L Nicolae, Carole Ober, and Hae Kyung Im. Shared and distinct genetic risk factors for childhood-onset and adult-onset asthma: genome-wide and transcriptome-wide studies. *The Lancet Respiratory Medicine*, 7(6): 509–522, 2019.
- [19] Shaun Purcell, Benjamin Neale, Kathe Todd-Brown, Lori Thomas, Manuel AR Ferreira, David Bender, Julian Maller, Pamela Sklar, Paul IW De Bakker, Mark J Daly, et al. Plink: a tool set for whole-genome association and population-based linkage analyses. *The American journal of human genetics*, 81(3):559–575, 2007.
- [20] Jian Yang, S Hong Lee, Michael E Goddard, and Peter M Visscher. Gcta: a tool for genome-wide complex trait analysis. *The American Journal of Human Genetics*, 88(1): 76–82, 2011.
- [21] Xiaoming Liu, Chunlei Wu, Chang Li, and Eric Boerwinkle. dbnsfp v3. 0: A one-stop database of functional predictions and annotations for human nonsynonymous and splice-site snvs. *Human mutation*, 37(3):235–241, 2016.
- [22] David Golan, Eric S Lander, and Saharon Rosset. Measuring missing heritability: inferring the contribution of common variants. *Proceedings of the National Academy of Sciences*, 111(49):E5272–E5281, 2014.
- [23] Seyhan Boyoglu-Barnum, Sean O Todd, Tatiana Chirkova, Thomas R Barnum, Kelsey A Gaston, Lia M Haynes, Ralph A Tripp, Martin L Moore, and Larry J Anderson. An anti-g protein monoclonal antibody treats rsv disease more effectively than an anti-f monoclonal antibody in balb/c mice. *Virology*, 483:117–125, 2015.
- [24] Alexander Bukreyev, Lijuan Yang, and Peter L Collins. The secreted g protein of human respiratory syncytial virus antagonizes antibody-mediated restriction of replication involving macrophages and complement. *Journal of virology*, 86(19):10880–10884, 2012.
- [25] Larry J Anderson, P Bingham, and JC Hierholzer. Neutralization of respiratory syncytial virus by individual and mixtures of f and g protein monoclonal antibodies. *Journal of virology*, 62(11):4232–4238, 1988.
- [26] JO Ngwuta, M Chen, K Modjarrad, MG Joyce, M Kanekiyo, A Kumar, HM Yassine, SM Moin, AM Killikelly, GY Chuang, et al. Prefusion f-specific antibodies determine the magnitude of rsv neutralizing activity in human sera. *sci transl med* 7: 309ra162, 2015.
- [27] S. A. Schobel, K. M. Stucker, M. L. Moore, L. J. Anderson, E. K. Larkin, J. Shankar, J. Bera, V. Puri, M. H. Shilts, C. Rosas-Salazar, R. A. Halpin, N. Fedorova, S. Shrivastava, T. B. Stockwell, R. S. Peebles, T. V. Hartert, and S. R. Das. Respiratory Syncytial Virus whole-genome sequencing identifies convergent evolution of sequence duplication in the C-terminus of the G gene. *Sci Rep*, 6:26311, May 2016. ISSN 2045-2322. doi: 10.1038/srep26311. Edition: 2016/05/24.

- [28] K. Li, S. Shrivastava, A. Brownley, D. Katzel, J. Bera, A. T. Nguyen, V. Thovarai, R. Halpin, and T. B. Stockwell. Automated degenerate PCR primer design for high-throughput sequencing improves efficiency of viral sequencing. *Virol J*, 9:261, November 2012. ISSN 1743-422x. doi: 10.1186/1743-422x-9-261. Edition: 2012/11/08.
- [29] QIAGEN Aarhus. White paper on de novo assembly in CLC Assembly Cell 4.0. *digitalinsights*, page 14, June 2016. URL <https://digitalinsights.qiagen.com/files/whitepapers/whitepaper-denovo-assembly.pdf>. Place: Denmark Publisher: Qiagen.
- [30] S. Wang, J. P. Sundaram, and T. B. Stockwell. VIGOR extended to annotate genomes for additional 12 different viruses. *Nucleic Acids Res*, 40(Web Server issue):W186–92, July 2012. ISSN 0305-1048 (Print) 0305-1048. doi: 10.1093/nar/gks528. Edition: 2012/06/07.
- [31] Kazutaka Katoh and Daron M Standley. Mafft multiple sequence alignment software version 7: improvements in performance and usability. *Molecular biology and evolution*, 30(4):772–780, 2013.
- [32] Lam-Tung Nguyen, Heiko A Schmidt, Arndt Von Haeseler, and Bui Quang Minh. Iq-tree: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular biology and evolution*, 32(1):268–274, 2015.
- [33] Christian Rosas-Salazar, Zheng-Zheng Tang, Meghan H Shilts, Kedir N Turi, Qilin Hong, Derek A Wiggins, Christian E Lynch, Tebeb Gebretsadik, James D Chappell, R Stokes Peebles Jr, et al. Upper respiratory tract bacterial-immune interactions during respiratory syncytial virus infection in infancy. *Journal of Allergy and Clinical Immunology*, 149(3):966–976, 2022.
- [34] Andrew R Connelly, Brian M Jeong, Mackenzie E Coden, Jacob Y Cao, Tatiana Chirkova, Christian Rosas-Salazar, Jacqueline-Yvonne Cephus, Larry J Anderson, Dawn C Newcomb, Tina V Hartert, et al. Metabolic reprogramming of nasal airway epithelial cells following infant respiratory syncytial virus infection. *Viruses*, 13(10):2055, 2021.
- [35] Emma K Larkin and Tina V Hartert. Genes associated with rsv lower respiratory tract infection and asthma: the application of genetic epidemiological methods to understand causality. *Future virology*, 10(7):883–897, 2015.
- [36] AliReza Eshaghi, Venkata R Duvvuri, Rachel Lai, Jeya T Nadarajah, Aimin Li, Samir N Patel, Donald E Low, and Jonathan B Gubbay. Genetic variability of human respiratory syncytial virus a strains circulating in ontario: a novel genotype with a 72 nucleotide g gene duplication. *PloS one*, 7(3):e32807, 2012.
- [37] Caroline Breese Hall, Joyce M Geiman, Robert Biggar, David I Kotok, Patricia M Hogan, and R Gordon Douglas Jr. Respiratory syncytial virus infections within families. *New England journal of medicine*, 294(8):414–419, 1976.

- [38] S Levine, R Klaiber-Franco, and PR Paradiso. Demonstration that glycoprotein g is the attachment protein of respiratory syncytial virus. *Journal of General Virology*, 68(9): 2521–2524, 1987.
- [39] Steven A Feldman, R Michael Hendry, and Judy A Beeler. Identification of a linear heparin binding domain for human respiratory syncytial virus attachment glycoprotein g. *Journal of virology*, 73(8):6610–6617, 1999.
- [40] Steven A Feldman, Susette Audet, and Judy A Beeler. The fusion glycoprotein of human respiratory syncytial virus facilitates virus attachment and infectivity via an interaction with cellular heparan sulfate. *Journal of Virology*, 74(14):6442–6447, 2000.
- [41] HW McL Rixon, G Brown, JT Murray, and RJ Sugrue. The respiratory syncytial virus small hydrophobic protein is phosphorylated via a mitogen-activated protein kinase p38-dependent tyrosine kinase activity during virus infection. *Journal of General Virology*, 86 (2):375–384, 2005.
- [42] Reena Ghildyal, Dongsheng Li, Irene Peroulis, Benjamin Shields, Phillip G Bardin, Michael N Teng, Peter L Collins, Jayesh Meanger, and John Mills. Interaction between the respiratory syncytial virus g glycoprotein cytoplasmic domain and the matrix protein. *Journal of General Virology*, 86(7):1879–1884, 2005.
- [43] Peter L Collins and Geneviève Mottet. 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*, 73(4):849–863, 1992.
- [44] Monica E Brint, Joshua M Hughes, Aditya Shah, Chelsea R Miller, Lisa G Harrison, Elizabeth A Meals, Jacqueline Blanch, Charlotte R Thompson, Stephania A Cormier, and John P DeVincenzo. Prolonged viral replication and longitudinal viral dynamic differences among respiratory syncytial virus infected infants. *Pediatric research*, 82(5): 872–880, 2017.
- [45] Sara M Johnson, Beth A McNally, Ioannis Ioannidis, Emilio Flano, Michael N Teng, Antonius G Oomens, Edward E Walsh, and Mark E Peebles. Respiratory syncytial virus uses cx3cr1 as a receptor on primary human airway epithelial cultures. *PLoS pathogens*, 11(12):e1005318, 2015.
- [46] Ralph A Tripp, Les P Jones, Lia M Haynes, HaoQiang Zheng, Philip M Murphy, and Larry J Anderson. Cx3c chemokine mimicry by respiratory syncytial virus g glycoprotein. *Nature immunology*, 2(8):732–738, 2001.
- [47] Kwang-II Jeong, Peter A Piepenhagen, Michael Kishko, Joshua M DiNapoli, Rachel P Groppo, Linong Zhang, Jeffrey Almond, Harry Kleanthous, Simon Delagrange, and Mark Parrington. Cx3cr1 is expressed in differentiated human ciliated airway cells and co-localizes with respiratory syncytial virus on cilia in a g protein-dependent manner. *PloS one*, 10(6):e0130517, 2015.

- [48] Alexander Bukreyev, Lijuan Yang, Jens Fricke, Lily Cheng, Jerrold M Ward, Brian R Murphy, and Peter L Collins. 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*, 82(24):12191–12204, 2008.
- [49] Tatiana Chirkova, Christian Rosas-Salazar, Tebeb Gebretsadik, Samadhan J Jadhao, James D Chappell, R Stokes Peebles Jr, William D Dupont, Dawn C Newcomb, Sergejs Berdnikovs, Peter J Gergen, et al. Effect of infant rsv infection on memory t cell responses at age 2-3 years. *Frontiers in immunology*, 13:826666, 2022.
- [50] Ani Manichaikul, Josyf C. Mychaleckyj, Stephen S. Rich, Kathy Daly, Michèle Sale, and Wei-Min Chen. Robust relationship inference in genome-wide association studies. *Bioinformatics*, 26(22):2867–2873, 2010. Number: 22 ISBN: 1460-2059 Publisher: Oxford University Press.

8 Supplemental

8.1 Supplemental host genetic analyses

We further investigated the possibility that the analysis was underpowered to identify associations with reported childhood asthma- and RSV LRTI-associated SNPs [16–18]. This was done by pooling information across SNPs to estimate the average genetic effect size. In brief, we computed a z-score for each SNP, where the average (across SNPs) squared. As \bar{G} is an average of $p = 54$ approximately independent statistics, it is approximately $N(n\mu^2 + 1, 2/p)$, where $n = 621$ is the sample size and μ^2 is a function of the average squared genetic effect on RSV infection in infancy. Using the genetic effect estimates from Janssen et al. [16]; Pasanen et al. [17]; Pividori et al. [18], we calculated that we would have 80% power to reject the global null hypothesis of no genetic effect at any of these SNPs (i.e., $\mu^2 = 0$) if, on average across the 54 SNPs, the genetic effect on RSV infection in infancy was at least 61% as large as those estimated in the aforementioned 3 studies. The z-score \bar{G} is proportional to the average squared genetic effect on RSV infection in infancy. We found $\bar{G}=1.00$ in our data, which corresponds to a p value of 0.50. This result indicates that the genetic effect on RSV infection in infancy is zero or small at SNPs likely to be associated with RSV infection a priori.

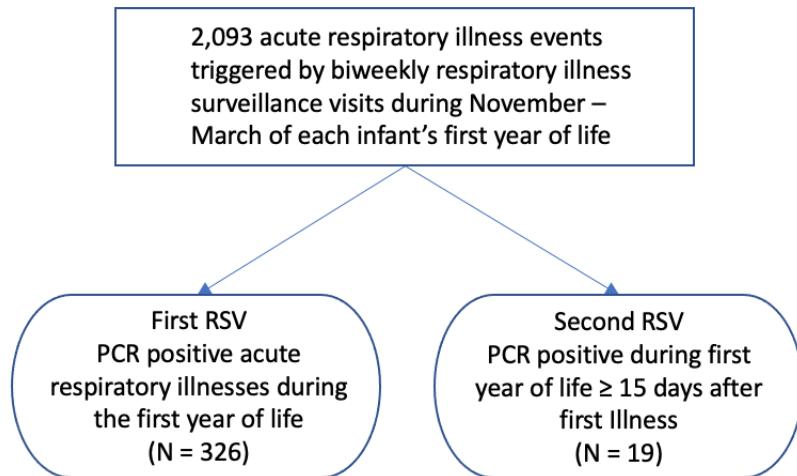


Figure S1: Supplemental: INFant Susceptibility to Pulmonary Infections and Asthma Following RSV Exposure (INSPIRE). The study population is a longitudinal birth cohort specifically designed to capture the first RSV infection in term healthy infants. Prolonged infection was a priori defined as repeatedly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

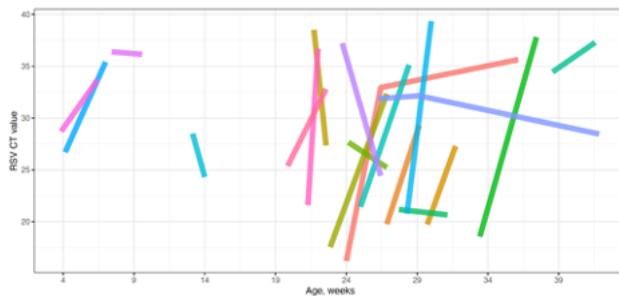


Figure S2: Supplemental: Infant RSV prolonged infections. Each line represents an infant in the study, and line start and end correspond to clinical respiratory illness sampling timepoints. CT values are inversely related to viral RNA abundance.

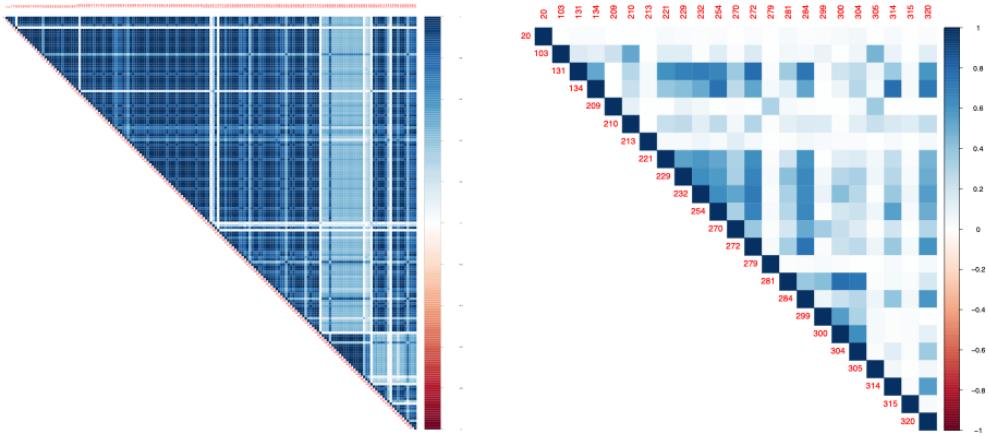


Figure S3: Supplemental: Variant clumping for reduction in association testing.
 [Left] Correlation between all positions. [Right] Correlation between proxy variants were clumped to remove $r^2 \geq 0.8$. Values indicate relative amino acid positions within MSA. r^2 indicated by color.

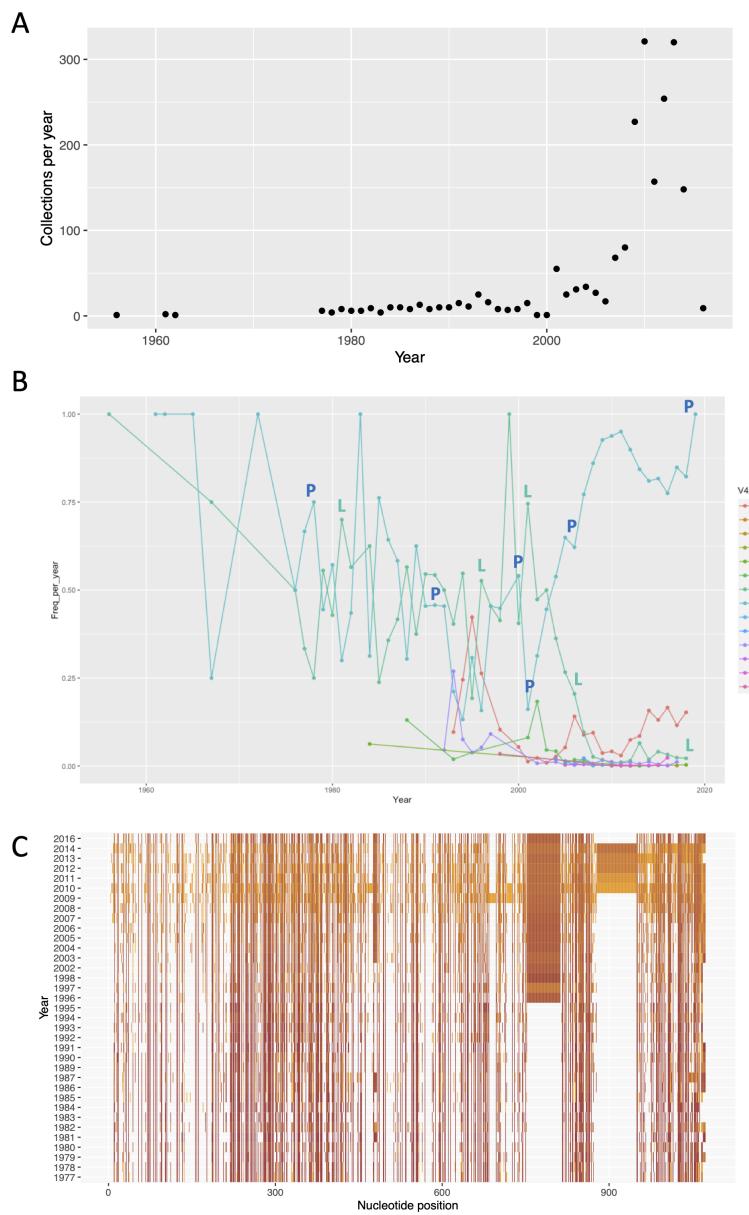


Figure S4: Supplemental: Publicly available RSV sequence data for > 60 years.
 (A) Global sample collection per year. (B) Variant associated with prolonged infection tracked in public data. The lead proxy SNP, p.P218T/S/L is illustrated here (relative amino acid positive 410 in MSA). The major alleles (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).