1    **Diversity and Evolution of Computationally Predicted T Cell Epitopes against**

2    **Human Respiratory Syncytial Virus**

3    Short Title: Sequence-based characterization of RSV T-cell immune landscape

4

5    Jiani Chen [1][2][3][4], Swan Tan [1][3][4][5,] Vasanthi Avadhanula [6], Leonard Moise [3][7], Pedro A

6    Piedra [6], Anne S De Groot [3][7], Justin Bahl [1][2][3][4][5][8] *

7    [1] Center for Ecology of Infectious Diseases, University of Georgia, Athens, GA, USA

8    [2] Institute of Bioinformatics, University of Georgia, Athens, GA, USA

9    [3] Center for Vaccines and Immunology, University of Georgia, Athens, GA, USA

10    [4] Center for Influenza Disease and Emergence Response, University of Georgia, Athens, GA,

11    USA

12    [5] Department of Infectious Diseases, University of Georgia, Athens, GA, USA

13    [6] Department of Molecular Virology and Microbiology, Baylor College of Medicine, Houston,

14    TX, USA

15    [7] EpiVax Inc., Providence, RI, USA

16    [8] Department of Epidemiology and Biostatistics, University of Georgia, Athens, GA, USA

17    * Corresponding author email: justin.bahl@uga.edu

18    **Abstract**

19    Human respiratory syncytial virus (RSV) is a major cause of lower respiratory infection. Despite

20    more than 60 years of research, there is no licensed vaccine. While B cell response is a major

21    focus for vaccine design, the T cell epitope profile of RSV is also important for vaccine

22    development. Here, we computationally predicted putative T cell epitopes in the Fusion protein

23    (F) and Glycoprotein (G) of RSV wild circulating strains by predicting Major Histocompatibility

24    Complex (MHC) class I and class II binding affinity. We limited our inferences to conserved

25    epitopes in both F and G proteins that have been experimentally validated. We applied

26    multidimensional scaling (MDS) to construct T cell epitope landscapes to investigate the

27    diversity and evolution of T cell profiles across different RSV strains. We find the RSV strains

28    are clustered into three RSV-A groups and two RSV-B groups on this T epitope landscape.

29    These clusters represent divergent RSV strains with potentially different immunogenic profiles.

30    In addition, our results show a greater proportion of F protein T cell epitope content conservation

31    among recent epidemic strains, whereas the G protein T cell epitope content was decreased.

32    Importantly, our results suggest that RSV-A and RSV-B have different patterns of epitope drift

33    and replacement and that RSV-B vaccines may need more frequent updates. Our study provides

34    a novel framework to study RSV T cell epitope evolution. Understanding the patterns of T cell

35    epitope conservation and change may be valuable for vaccine design and assessment.

36 **Author Summary**

37 Lower respiratory infections caused by human respiratory syncytial virus (RSV) is a global

38 health challenge. B cell epitope immune response has been the major focus of RSV vaccine and

39 therapeutic development. However, T cell epitope induced immunity plays an important role in

40 the resolution of RSV infection. While RSV genetic diversity has been widely reported, few

41 studies focus on RSV T epitope diversity, which can influence vaccine effectiveness. Here, we

42 use computationally predicted T cell epitope profiles of circulating strains to characterize the

43 diversity and evolution of the T cell epitope of RSV A and B. We systematically evaluate the T

44 epitope profile of RSV F and G proteins. We provide a T cell epitope landscape visualization

45 that shows co-circulation of three RSV-A groups and two RSV-B groups, suggesting potentially

46 distinct T cell immunity. Furthermore, our study shows different levels of F and G protein T cell

47 epitope content conservation, which may be important to correlate with duration of vaccine

48 protection. This study provides a novel framework to study RSV T cell epitope evolution, infer

49 RSV T cell immunity at population levels and monitor RSV vaccine effectiveness.

50    **Introduction**

51    Human respiratory syncytial virus (RSV) is a negative-strand RNA virus that is classified in

52    the *Orthopneumovirus* genus of the family *Pneumoviridae*. It is a major cause of lower

53    respiratory disease in young infants, immunocompromised individuals, and elderly people,

54    resulting in annual epidemics worldwide [1]. The single-stranded RNA genome of RSV is

55    approximate 15.2 kb and encodes 11 viral proteins [2]. The Fusion (F) and Glycoprotein (G)

56    proteins are the two major surface proteins [3]. F protein is generally thought to be conserved

57    and therefore it is the focus of most current RSV vaccine designs. Although G protein is highly

58    variable, its contribution to disease pathogenesis and its role in the biology of infection suggest it

59    can also be an effective RSV vaccine antigen [4]. Despite the significant burden of RSV

60    infection worldwide, there is no licensed vaccine. The only approved intervention is passive

61    immuno-prophylaxis with palivizumab, which is achieved by administering the monoclonal

62    antibody (mAb) to a highly restricted group of infants under the age of 24 months and treatment

63    must be repeated monthly during the RSV season due to the relatively short half-life of the

64    antibody[5], [6]. Due to the high cost of monoclonal antibody treatments, this intervention is

65    limited to high-risk infants and is generally unavailable in developing countries. An RSV vaccine

66    is an urgent global healthcare priority, and it is likely that different strategies are needed for the

67    various high-risk groups.

68

69    A number of research teams have worked on the development of RSV vaccine since its isolation

70    and characterization in 1956 [7], [8]. However, vaccination with the formalin-inactivated, alum

71    precipitated RSV (FI-RSV) vaccine in RSV-naïve infants and young children, led to the

72    development of vaccine enhanced disease (VED) that hampered vaccine development for

73    decades to follow [9]. Many studies have been conducted to explain this undesirable outcome. It

74    is likely that formalin fixation led to a vaccine that mostly presented the post-fusion

75    conformation of RSV F protein, leading to an excess of non-neutralizing antibodies and immune

76    complex formation [10] [11] [12]. Other studies indicated that an impaired T cell response with

77    Th2 skewing [13], [14], as well as complement deposition in the lungs, contributed to enhanced

78    neutrophil recruitment [12]. Due to the recent breakthrough to structural constrain the F protein

79    in the pre-fusion conformation and the development of RSV rodent models, there has been a

80    surge in the number of RSV vaccine candidates undergoing clinical evaluation.

81

82    While most current RSV vaccination strategies focus on a B-cell-induced neutralization immune

83    response, T cell immunity also plays a major role in the resolution of virus infection and is

84    essential for RSV vaccine development [15], [16]. Once RSV infection of the lower airways is

85    established, CD8 T cells play an important part in viral clearance and CD4 helper T cells can

86    orchestrate cellular immune responses and stimulate B cells to produce antibodies. However,

87    Th2-biased responses have been associated with animal models of RSV VED, and measurement

88    of Th1 and Th2 responses are considered important to predict the safety of vaccine candidates

89    [12]. Therefore, induction of a balanced cell-mediated immune response through vaccination

90    would promote RSV clearance, but caution must be taken to avoid the potential for

91    immunopathology. Taken together, a closer examination of T cell immunity and the virus

92    sequences that induce T cell responses are needed for RSV vaccine development.

93

94    Human respiratory syncytial virus has a complex circulation pattern in the human population.

95    Within two antigenic groups, RSV-A and RSV-B, different genotypes can co-circulate within the

96    same community, while novel RSV genotypes with high genomic diversity may arise and

97    potentially replace the previously dominant genotypes [17]. In recent years, several unique

98    genetic modifications in RSV have been identified, including a 72-nucleotide (nt) duplication

99    (ON genotype) in RSV-A G gene and another with a 60-nt duplication (BA genotype) in RSV-B

100   at a similar region [18]. The observed RSV genetic diversity has raised a question about whether

101   it is necessary for an RSV vaccine to include several different strains to be effective. Most

102   current RSV vaccine developments are based on an RSV A2 laboratory strain, which is a

103   chimeric strain that belongs to subtype A [19]. While these treatments hold promise, there is the

104   possibility of viral strains developing escape mutations. For example, palivizumab-resistant

105   strains have been isolated from both RSV rodent models and human [20][16]. Monoclonal

106   antibody tests have demonstrated additional antigenic variability within RSV-A /B antigenic

107   groups and suggest that it may play a role in the ability of RSV to escape immune response and

108   established infections [21]. In addition, amino-acid variation at the T cell epitope level and the

109   emergence of novel T cell epitopes have been reported [22], but further studies are needed to

110   illustrate the effect of these variations on T cell recognition. T cell epitopes are sometimes cross-

111   reactive, which is defined as the recognition of two or more epitope peptide-MHC complexes by

112   the same T cell receptor and these cross-reactive epitopes are predicted to be cross-conserved

113   [23]. Prediction of cross-conservation is important for vaccine design because it would be useful

114   to predict protection against a pathogen with different lineages and identify escape variants.

115   Hence, characterizing T cell epitope profiles across different strains can be crucial for RSV

116   vaccine development.

117

118    In this study, we utilize immunoinformatic approaches that are implemented in the iVAX toolkit

119    [24] to predict T cell epitopes in RSV across different strains with a focus on the two major

120    surface proteins F and G. With the analysis of a comprehensive dataset, we evaluate the lineage-

121    specific T cell epitope profile of RSV. We also create sequence-based T cell epitope landscapes

122    based on epitope content comparison across different strains and further correlate RSV T cell

123    immunity change with virus evolution. The proportion of cross conserved T cell epitope content

124    between vaccine candidate strains that developed earlier and RSV circulating strains with

125    different isolated years and locations were also calculated. These analyses may aid in

126    understanding RSV T cell immunity across different strains and contribute to current vaccine

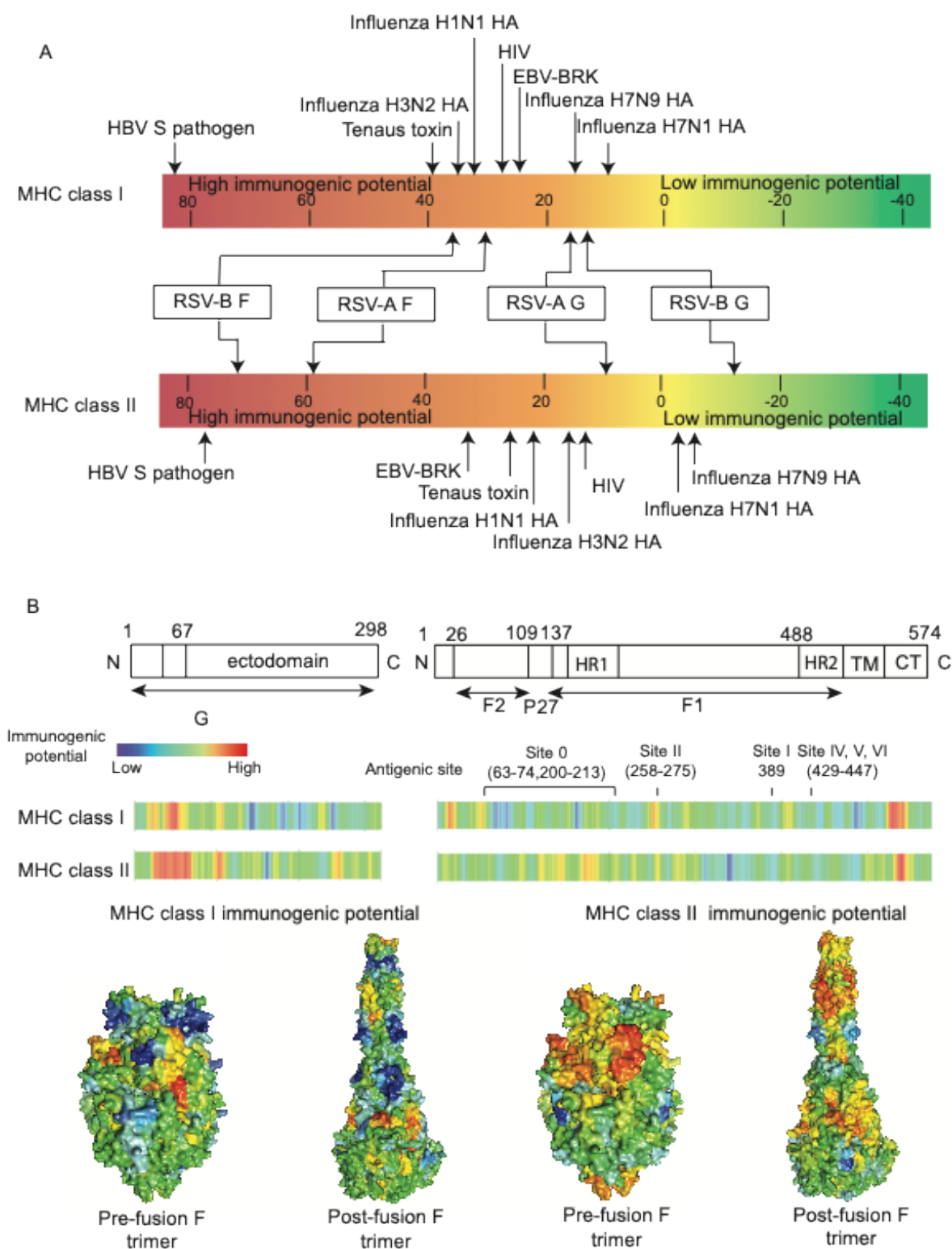127    design efforts.

128

129    **Results:**

130    **Distribution of T cell epitopes in RSV surface proteins**

131    We evaluated the T cell immunogenic potential across RSV surface proteins by scanning 9

132    residue regions to predict the binding probability to MHC class I and class II molecules (Figure

133    1). The epitope density of RSV surface proteins was evaluated using a normalized epitope

134    density score. F protein has an epitope density score greater than +20 for both the class I and

135    class II immunogenicity scale analysis, indicating significant immunogenic potential [24]. This

136    contrasts with lower G protein class I and class II epitope density protein scores for both

137    subtypes. The class I epitope density score of G protein was greater than +10 in both subtypes

138    but the class II density was lower than random expectation in the analysis of RSV-B (Figure 1A).

139    This result suggests that RSV surface proteins are likely to have the potential to stimulate T cells

140    that are required for protective immunity. We then investigated the distribution of T cell

7

141    immunogenicity across the proteins and found that there are regions with relatively high T cell

142    immunogenic potential (Figure 1B). The distribution of T cell immunogenicity of F protein was

143    mapped onto its protein structure and overlap between protein sequence regions with high T cell

144    immunity potential and the antibody neutralizing targets was observed at antigenic site Φ and

145    site II.

146

147 **Figure 1: T cell immunogenic potential for RSV surface proteins based on MHC binding**

148 **prediction**. (A) T cell immunogenic potential of RSV major surface proteins. T cell epitope

149 density scores for RSV major surface proteins and other pathogen proteins are labeled on a scale

150 bar. Low-scoring proteins are known to engender little to no immunogenicity while higher-

151 scoring proteins are known immunogens. Proteins scoring above +20 on this scale are considered

152 to have significant immunogenic potential. (B) Distribution of RSV T cell immunogenic

153 potential across F and G protein in RSV reference strain A2. Prefusion or post-fusion F protein

154 surface was colored by the relative immunogenetic potential at each location. Analyses are based

155 on the RSV-A reference sequence.

156

157 **Lineage specific T cell epitope profiles**

158 We then extended T cell epitope predictions from RSV representative strains to multiple wild-

159 circulating strains. The distribution and diversity of T cell epitopes across different strains are

160 illustrated in heatmaps with the corresponding time-scaled phylogenies (Supplementary Figure 1

161 and Supplementary Figure 2). Both F and G proteins contain epitopes that were conserved across

162 all RSV strains in almost 100% of sampled isolates, suggesting that they could serve as high-

163 quality T cell epitope candidates for vaccine design. In contrast, some epitopes were mutated in

164 selected strains, and those epitopes that only occurred in certain clades within the phylogeny

165 could be interpreted as clade-specific "fingerprints".

166

167 The G gene duplication events in RSV, which are unique gene signatures, can either shift the

168 position of epitopes or cause the emergence of novel epitopes. Two novel class I epitopes, (no.

169 31 and no. 40 in Supplementary Figure 2A), were found in RSV-A strains that contain G gene

170    duplication. In addition, an emergent class II epitope (no. 25 in Supplementary Figure 2A) was

171    identified in RSV-A sequences that contain G gene duplication, which was a shift from an

172    epitope (no. 24) that has been observed in other strains. From RSV-B strains that contain the G

173    gene duplication event, we also observed multiple lineage specific class I T cell epitopes, which

174    are caused by a 2-aa deletion (aa157 and aa158) in these strains instead of directly due to the 60-

175    nt duplication event. RSV-B G proteins that have the duplication event contain multiple novel

176    epitopes (no. 22, 23, 26, 28, 30, 37) but do not contain several epitopes (no. 24, 25, 27, 29, 31,

177    38) that are identified in other strains (Supplementary Figure 2B).

178

179    To further determine whether the T cell epitopes defined using EpiMatrix might be

180    immunogenic, we utilized the JanusMatrix [25] algorithm to identify the T cell epitopes that are

181    likely to be cross-conserved with human epitopes and thereby tolerated by the immune system.

182    Based on this analysis, 6.45% of putative class I epitopes and 1.12 % of putative class II epitopes

183    of RSV major surface proteins are cross-conserved with human proteome-derived epitopes at

184    TCR facing residues. As these peptides have similar HLA binding preferences that are contained

185    in human proteins (Supplementary Figure 3), they were therefore assumed not to be

186    immunogenic. After excluding the high-JanusMatrix score epitopes identified above, we were

187    able to identify T cell epitopes that were conserved in more than 60% of currently circulating

188    RSV strains. We searched the IEDB epitope database to determine if these epitopes were related

189    to experimentally validated RSV T cell epitopes or HLA ligands. The conserved RSV T cell

190    epitope sequences that may be important for future vaccine development are shown in Table 1

191    and Table 2.

11

192 **Table 1: Experimentally validated conserved MHC class I epitopes peptides in RSV major**

193 **surface proteins [a]**

194

| Subgroup | Protein | Epitope address | Epitope sequence [b] | Binding HLAs [c] | Conservation [d] | Number of human matches [e] | Epitope id in IEDB |
|---|---|---|---|---|---|---|---|
| RSV-A & RSV-B | F | 45-53 | **LSALRTGWY** | A0101 | 99.55%(A) & 74.24%(B) | 1 | 158982 |
| | | 140-148 | **FLLGVGSAI** | A0201 | 99.59%(A) & 97.98%(B) | 0 | 156869 |
| | | 250-258 | **YMLTNSELL** | A0201, A2402 | 99.59%(A) & 99.33%(B) | 0 | 156979 |
| | | 272-280 | **KLMSSNVQI** | A0201 | 66.64%(A) & 96.08%(B) | 3 | 156902 |
| | | 273-281 | **LMSSNVQIV** | A0201 | 66.56%(A) & 96.08%(B) | 1 | 156915 |
| | | 449-457 | **TVSVGNTLY** | A0101 | 99.75%(A) & 99.33%(B) | 0 | 97017 |
| RSV-A | F | 10-18 | AITTILAAV | A0201 | 84.69% | 3 | 156844 |
| | | 111-119 | LPRFMNYTL | B0702 | 91.18% | 0 | 158975 |
| | | 170-178 | ALLSTNKAV | A0201 | 99.67% | 2 | 156847 |
| | | 383-391 | NIDIFNPKY | A0101 | 95.86% | 0 | 159045 |
| | G | 25-33 | FISSCLYKL | A0201 | 99.26% | 0 | 158759 |
| | | 61-69 | FIASANHKV | A0201 | 82.08% | 0 | 158751 |
| RSV-B | F | 525-533 | IMITAIIIV | A0201 | 89.25% | 0 | 156892 |
| | | 540-548 | SLIAIGLLL | A0201 | 97.65% | 5 | 156960 |
| | G | 25-33 | VISSCLYKL | A0201 | 90.91% | 0 | 158759 |
| | | 61-69 | FIISANHKV | A0201 | 99.02% | 0 | 158751 |

195

196 a. This table contains putative MHC class I epitopes that have already been experimentally

197 validated in publications.

198 b. Epitopes sequences that are conserved in both RSV-A and RSV-B are in bold.

199 c. HLAs that have the top 1% binder scores in EpiMatrix for epitope sequence.

200 d. The conservation is evaluated by the presence of epitope peptides across all RSV-A or

201 RSV-B sequences that are publicly available (only epitope sequences with at least 60%

202 conservation are shown in the table).

203 e. Count of human epitopes found in the search database. JanusMatrix was used to search

204 human epitopes that are predicted to bind to the same allele as the RSV epitope and share

205 TCR facing contacts with the RSV epitope.

12

206

**Table 2: Experimentally validated conserved MHC class II epitopes peptides in RSV major surface proteins [a]**

| Subtype | Protein | Epitope address | Epitope sequence [b] | Conservation [c] | Number of human matches [d] | Epitope id in IEDB |
|---|---|---|---|---|---|---|
| RSV-A | F | 29 - 44 | TEEF**YQSTCSAVS**KGY | 98.53% | 3 | 956680 |
| | | 50 - 70 | TGW**YTSVITIELSNIK**ENKCN | 97.75% | 1 | 153700 |
| | | 167 - 192 | IKSALLSTNKAVVSLSNGVSVLTSKV | 93.14% | 4 | 545502 |
| | | 218 - 234 | ETVIEFQQKNNRLLEIT | 98.86% | 3 | 1087566 |
| | | 247 - 268 | VSTYMLTNSELLSLINDMPITN | 98.98% | 8 | 99471 |
| | | 288 - 310 | IMSIIKEEVLAYVVQLPLYGVID | 98.57% | 5 | 99334 |
| | | 399 - 418 | KTDVSSSV**ITSLGAIVS**CYG | 99.14% | 0 | 545603 |
| | | 453 - 470 | GNTLYYVNKQEGKSLYVK | 98.37% | 1 | 99691 |
| | | 492 - 510 | ISQVNEKI**NQSLAFIR**KSD | 80.32% | 1 | 153713 |
| | | 543 - 560 | AVG**LLLYCKARSTPV**TLS | 79.26% | 6 | 153641 |
| | G | 19 - 43 | TLNHLLFISSCLYKLNLKSIAQITL | 93.13% | 8 | 1087567 |
| RSV-B | F | 29 - 44 | TEE**FYQSTCSAVS**RGY | 99.78% | 3 | 956680 |
| | | 50 - 70 | TGW**YTSVITIELSNIK**ETKCN | 93.95% | 1 | 153700 |
| | | 192 - 218 | VLDLKNYINNQLLPIVNQQSCRISNIE | 83.43% | 4 | 153636 |
| | | 247 - 268 | LSTYMLTNSELLSLINDMPITN | 98.54% | 8 | 99471 |
| | | 399 - 418 | KTDISSSV**ITSLGAIVS**CYG | 98.88% | 0 | 545603 |
| | | 453 - 470 | GNTLYYVNKLEGKNLYVK | 98.77% | 0 | 99691 |
| | | 492 - 510 | ISQVNEKI**NQSLAFIR**RSD | 97.42% | 1 | 153713 |
| | | 543 - 560 | AIGL**LLYCKAKNTPV**TLS | 94.96% | 4 | 153641 |
| | G | 51 - 74 | STSLIIAAIIFIISANHKVTLTTV | 94.66% | 8 | 158751 |

209
210 a. This table contains putative MHC class II epitopes which share the identical binding groove sequence of the RSV class II epitopes that have already been experimentally validated in publications.

213 b. Underlined sequences represent the nine-mer frames with the greatest potential to bind class II HLA. Epitope sequences that are in bold indicate sequences are predicted to bind class II HLA and are conserved in both RSV-A and RSV-B.

13

216      c.  Conservation is evaluated by the presence of epitope peptides across all RSV-A or RSV-

217          B sequences that are publicly available (Only epitope sequences with at least 60%

218          conservation are shown in the table).

219      d.  Count of human epitopes found in the search database. JanusMatrix was used to search

220          human epitopes that are predicted to bind to the same allele as the RSV epitope and share

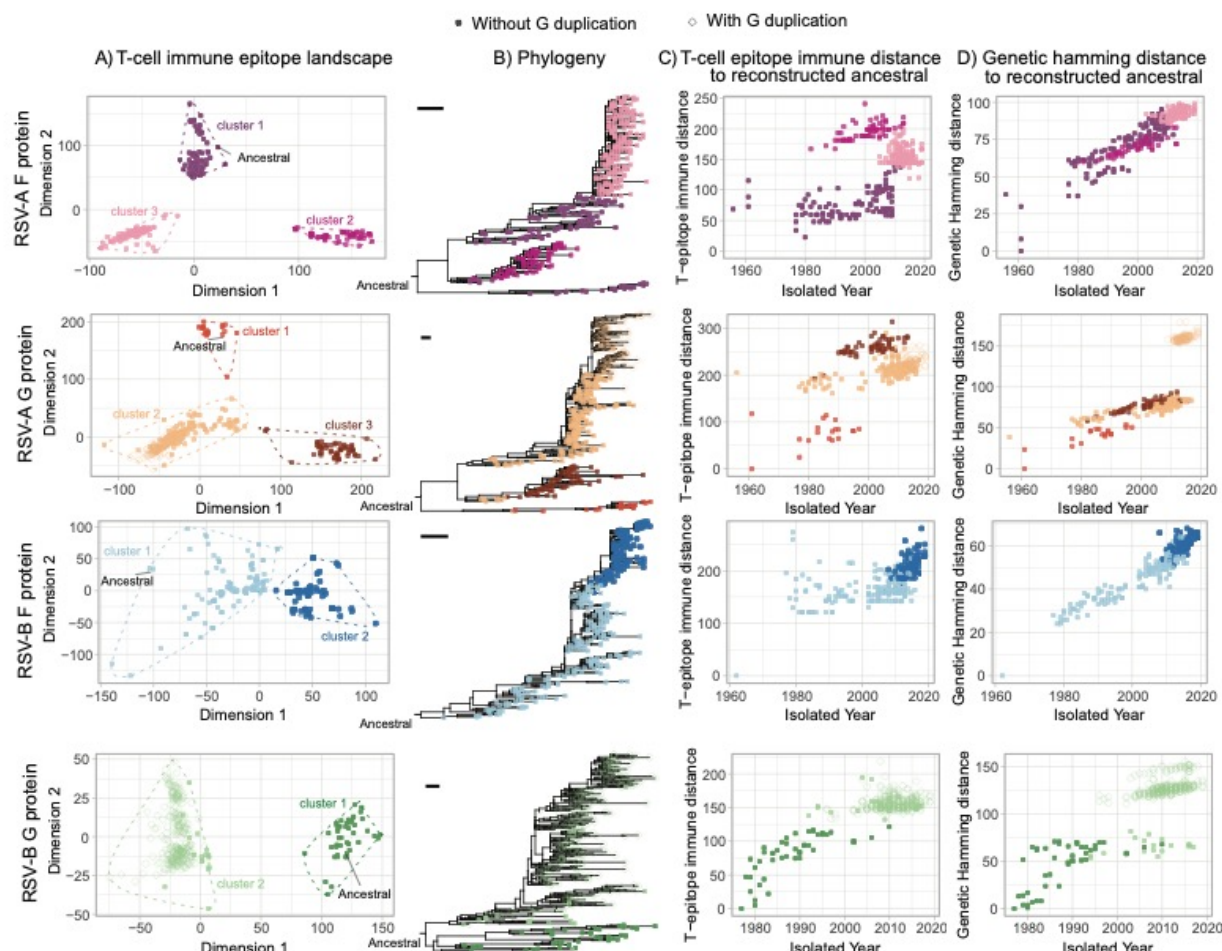221          TCR facing contacts with the RSV epitope.

222

**Predicted RSV T cell epitope landscapes**

To investigate the evolution of RSV on T cell immunity profiles, we develop a new approach to visualize the T immunity profile of multiple RSV strains on a landscape. We performed a T cell epitope content pairwise comparison between RSV strains, and T cell epitope distances between pair of RSV strains were then calculated using a T epitope distance algorithm. We then applied a multidimensional scaling (MDS) approach using these estimated pair-wise T epitope distances to map RSV strains to a landscape to characterize their T-cell immunity profile. We found both Class I and Class II T cell immunity profiles of F and G proteins of different RSV strains were clustered into groups on this T cell epitope landscapes (Supplementary Figure 4). Combining the Class I and Class II T-cell epitope biding profiles, RSV-A major surface protein isolates can be divided into three clusters and RSV-B major surface protein isolates can be divided into two clusters (Figure 2, Supplementary Figure 5). We observe that the G gene sequence isolates that contain 72-nt (RSV-A) or 60-nt (RSV-B) duplications clustered together with other sequences instead of forming isolated groups. To further investigate the T cell epitope diversity, we correlated this clustering pattern with the phylogenetic histories (Figure 2B). The phylogenetic tree topologies of the RSV-A F gene and G gene are similar. The F gene cluster 1 is paraphyletic, while cluster 2 and 3 are monophyletic. Cluster 1 is the closest to the ancestral sequence and mapping this group onto the phylogeny show that this cluster has a basal relationship with clusters 2 and 3 indicating that the phylogenetic divergence occurred prior to epitope drift. The RSV-B F and G gene genealogies are very different. In particular, the RSV-B F gene topologies is indicative of strong immune selection, similar to observed human influenza A virus or within host HIV phylogenies [26]. In contrast, the RSV-B G gene phylogeny shows the co-circulation of multiple lineages, though this could reflect the sequencing bias of G genes (Figure 2B). We

15

246    then calculated the T-cell epitope immune distance of each strain from a reconstructed ancestral

247    sequence (Figure 2C). These distances were then plotted against the year of isolation and colored

248    according to the cluster identified in Figure 2A. RSV-A shows that multiple predicted immune

249    phenotypes co-circulate and persist for long periods (>2 decades). Analysis of RSV-B shows a

250    turnover of the predicted immune phenotypes with short periods of co-circulation (<5 years) for

251    F and G protein T cell epitopes. The limited periods of co-circulation is again consistent with

252    phenotype patterns observed for viruses under strong immune selection (e.g H3N2 influenza A

253    virus) [27], [28]. In contrast, genetic distances from the reconstructed ancestral sequence plotted

254    against year of isolation show patterns typical of gradual genetic drift, except in the G gene

255    where a 72-nt and 60-nt insertion is present (Figure 2D). Taken together, these results suggest

256    that genetic and predicted T-cell epitope immune diversity are different and may be an important

257    factor to consider when evaluating RSV vaccine efficacy.

258

259    There are multiple methods available to predict T cell epitopes [29], which may result in

260    different reconstructed landscapes if there is a systematic bias in the prediction method. We used

261    the NetMHCpan method [30] to predict T cell epitopes and perform the same landscape

262    reconstruction using MHC class I binding predictions for RSV-A F protein. Our analysis showed

263    a consistent clustered pattern of RSV T epitope profile on the landscape regardless of T cell

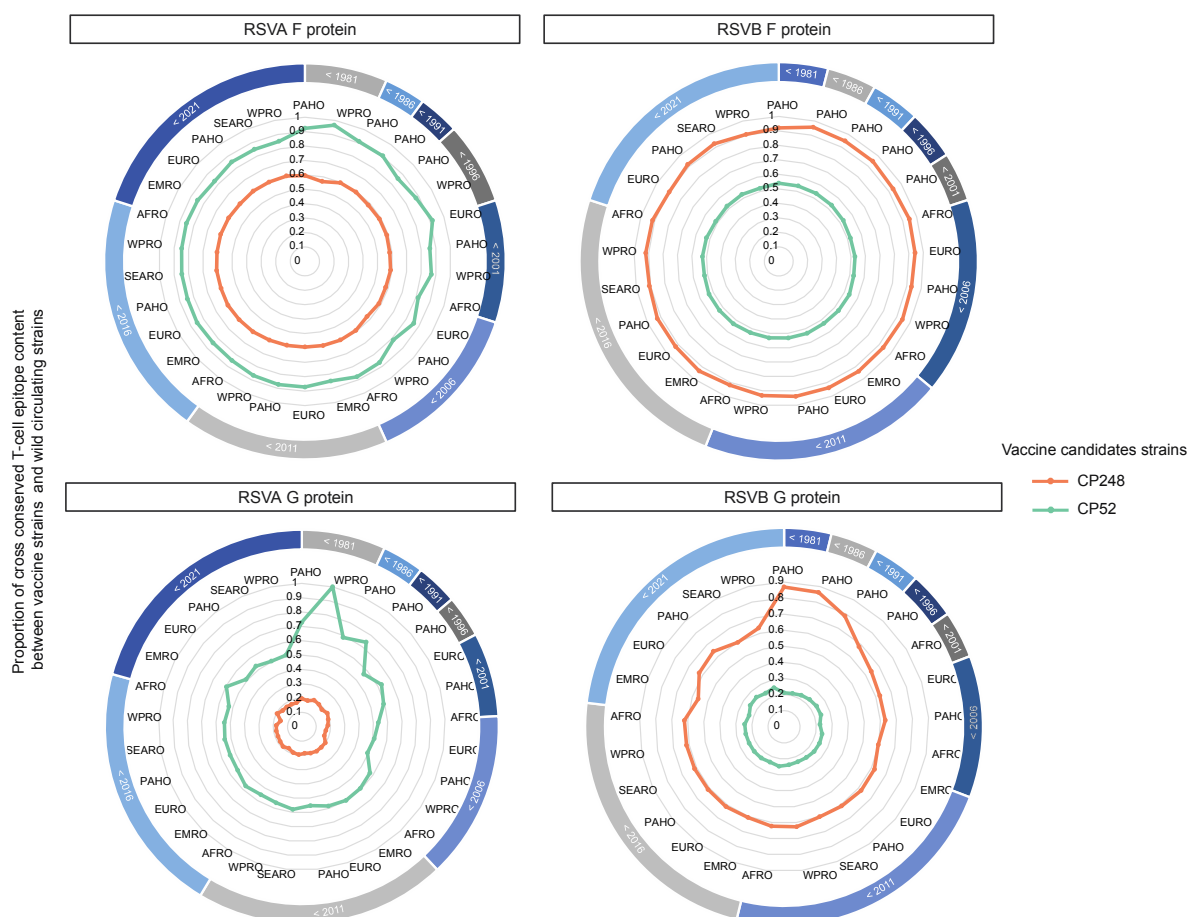264    epitope prediction method (Supplementary Figure 6).

265



266

**Figure 2: Predicted T cell epitope landscapes and genetic evolution of RSV surface proteins.** Epitope landscapes of RSV major surface proteins are built with MHC class I and class II epitope content comparison across different strains. Filled circles indicate F protein isolates or G protein isolates without duplication. Diamonds indicate G protein isolates with gene duplication. T cell immunity clusters are determined with *k-means* method and are used to color the sequenced isolates in the following panels. The corresponding Maximum Likelihood (ML) phylogenies are reconstructed and are rooted by mid-point. Scale bars indicate 0.005 nucleotide substitution per site. T cell epitope immune distance and genetic hamming distance from the estimated TMRCA are plotted against the isolated time of each sequence.

**Assessment of vaccine candidate strains with T cell epitope content**

To quantitatively evaluate whether it might be necessary to include multiple RSV strains to prepare an effective vaccine, two live attenuated RSV strains that are previously considered as vaccine candidates, CP248, a recombinant virus that belongs to subtype A, and CP52, which is a recombinant RSV-B strain, were included in our analysis and compared to wild-type strains using EpiCC. We calculated the average proportion of cross-conserved T cell epitope content between the selected vaccine strains and wild-circulating strains from different isolation years and WHO regional groups (Figure 3, Supplementary Figure 7). Different proportions of cross-conserved T cell epitope content against isolates from two different subtypes, A and B, were observed in both the F and G protein analyses. In the comparison of the vaccine strains and wild strains belonging to the same subtype, the proportion of cross-conserved T cell epitope in RSV F protein is relatively stable in different groups, all are higher than 78% for RSV-A and higher than 85% for RSV-B. In contrast, changes in the proportion of cross-conserved T cell epitopes were detected among groups within the same RSV subtype, especially in different temporal groups in the G protein analysis (Figure 3). Vaccine strain CP248 appears to have a relatively higher proportion of cross conserved T cell epitopes within G protein when compared to the RSV-A strains that were isolated before 1991 (> 70%) and a relatively lower degree of conservation against recently isolated strains. A similar decrease in T cell epitope conservation with time was identified for vaccine strain CP52 among circulating RSV-B strains.

18

**Figure 3: Evaluation of previously used RSV vaccine candidate strains with T cell epitope content of circulating strains**. RSV-A and RSV-B major surface protein sequences are subsampled and then grouped by isolation year and 6 isolated WHO regions. African Region (AFRO), Region of the Americas (PAHO), South-East Asia Region (SEARO), European Region (EURO), Eastern Mediterranean Region (EMRO) and Western Pacific Region (WPRO). The proportion of cross-conserved T cell epitope content between live attenuated strains (CP248 or CP52) and wild circulating strains are displayed as radar plots.

**Discussion**

19

305 Although both CD4 and CD8 T cells contribute to protection against RSV-induced disease

306 following primary infection [16][31], T cell epitopes have received limited attention in the RSV

307 research effort. We demonstrate RSV surface proteins appear to have significant potential to

308 drive T cell immunity using a computational approach, based on their T cell epitope density

309 scores as determined by MHC molecular binding prediction. The relatively high putative T-cell

310 epitope density might make F protein a good target for RSV vaccine. In addition to the analysis

311 of T cell epitope density and distribution in RSV major surface proteins, we also demonstrated

312 lineage-specific variations in T cell epitope content. Even though RSV F protein is believed to be

313 well conserved and G protein is reported to be highly variable, epitope mutations are observed

314 across different lineages within the F protein, and potential conserved T cell epitopes can still be

315 found in the highly variable G protein, suggesting that studying the lineage-specific T cell

316 epitopes in RSV can provide insight into the impact of immune selection on viral diversity and

317 persistence. While experimental validation is needed, this analysis highlights the importance of

318 understanding population-level epitope conservation as it may provide important insight into the

319 development of T cell epitope-driven vaccines against RSV infection.

320

321 A major focus of our work is the development of a sequence-based method to map the evolution

322 of T cell immunity across different stains. Following a previously pivotal work that used MDS

323 method to map the evolutionary adaptation of influenza A virus-induced by CD8 T cell using the

324 presence and absence of MHC class I epitopes [32], we constructed RSV T cell immunity

325 landscapes using immune distances that were generated by T cell epitope cross-conservation

326 analyses, which allows for easy visualization and intuitive understanding of the potential for T

327 cell immunity relationships among different strains. When comparing across strains, we found

328    that the T cell epitope content of RSV surface proteins from different strains can be clustered, as

329    has been observed for the antigenic relationship reported in other pathogens [49][50]. Our results

330    also demonstrate the correspondence between RSV T cell immunity clusters and their

331    corresponding phylogeny, with sequences in the same clade generally belonging to the same T

332    cell immune cluster. Importantly, we also observe different patterns of T-cell epitope evolution

333    of RSV wild strains compared with their genetic evolution, which highlights the importance of

334    characterizing T cell epitope changes in RSV.

335

336    We identified highly conserved RSV T cell epitopes in this study, some of which have already

337    been experimentally validated and published in the IEDB database. However, we also identified

338    several other conserved T cell epitopes that have not been previously described. These may be

339    valuable for vaccine design, although experimental validation will be needed. Furthermore, the

340    homology of selected RSV epitopes to human epitopes suggests that some predicted RSV T cell

341    epitopes might be tolerated by the human immune system, or could induce a harmful cross-

342    reactive immune response against human proteins when administered with an adjuvant [25].

343    Certain aspects of immunity to RSV were not addressed by this study. For example, neutralizing

344    antibody responses are currently considered to be the most important correlate of immunity.

345    While neutralizing antibodies would not directly be elicited by a T cell epitope-driven vaccine,

346    helper (CD4) T cell epitopes are required to generate high affinity, high specificity antibodies.

347    We also note that we have limited our focus on the two major RSV surface proteins in our

348    current analysis, but other RSV proteins like N, M, or M2-2 proteins might also contribute to

349    vaccine efficacy [36].

350

351    An effective vaccine against variable viruses should contain T cell epitopes that are highly

352    conserved among circulating strains [37]. Vaccine efficacy can be diminished if T cell epitopes

353    in a vaccine strain do not match when new strains of pathogens emerge. In this study, we used an

354    immunoinformatic-based approach to estimate cross-conserved T cell epitope contents between

355    two live attenuated vaccine candidate strains and RSV circulating wild strains. We found that

356    there was a low proportion of cross-conserved T cell epitope content with vaccine strains that

357    belonged to different antigenic groups, which indicates the risk of using a single-subtype strain

358    in RSV vaccines. In addition, we observed a lower proportion of cross-conserved G-protein T

359    cell epitope content between vaccine strains and recent circulating strains in the same antigenic

360    group, which suggests that including T cell epitopes from different strains in the same antigenic

361    group might also be important for RSV vaccine development. Although we did not observe a

362    significant change in cross-conserved T cell epitope content in F protein, we cannot rule out the

363    probability that variation of  F protein in the future could render a single-strain-based vaccine

364    less effective. Our current analysis is based on reduced datasets due to the heavy computational

365    capacity required to perform epitope content comparison. We constructed these representative

366    datasets by randomly subsampling the complete datasets according to geographical regions and

367    isolated years. Our findings may reflect the T cell epitope diversity of publicly available RSV

368    strains, however, additional RSV surveillance efforts may be required to get a full picture of the

369    T cell epitope variability of RSV.

370

371    The lack of experimental data might cause problems in epitope identification. Computational-

372    based T cell epitope landscapes have the potential for bias. But the observed clustered T cell

22

373   evolutionary pattern in RSV surface proteins provide valuable insights into virus evolution in the

374   aspects of T cell immunity and strain selection for vaccine design.

375

376   Overall, this study provides a focused analysis of T cell epitopes in RSV major surface proteins

377   using computational tools. We performed a comprehensive T cell epitope prediction for RSV

378   showing the immunological relationship of T cell epitopes in RSV surface proteins. This study

379   demonstrates that T cell epitope evolution may differ from genetic variation and provides a

380   framework for developing an integrated epitope-based RSV vaccine and evaluation methods that

381   could be used to optimize vaccination strategies.

382

383   **Materials and Methods**

384   **Dataset**

385   RSV GenBank records files were retrieved from NCBI's GenBank nucleotide database using the

386   search term "HRSVA" or "HRSVB" on June 22, 2020. F and G gene nucleotide sequences and

387   metadata including country of isolation and collection date were extracted using customized

388   python scripts. Genotype assignments were made with the program "LABEL", using a

389   customized RSV module [38] [39]. Countries of isolation were grouped into 6 WHO regions:

390   African Region, Region of the Americas, South-East Asia Region, European Region, Eastern

391   Mediterranean Region, and Western Pacific Region [40]. The following inclusion and exclusion

392   criteria were applied: (i) each sequence needed to have a known isolated geographic location and

393   isolated year, (ii) each sequence had to be at least 80% of the complete gene sequence in length,

394   (iii) identical sequences with the same isolate country were removed, and (iv) vaccine derivative

395   and recombinant sequences were removed. Using these criteria, comprehensive datasets of RSV

23

396     F and G genes were defined (RSV-A F gene = 1010, RSV-B F gene = 894, RSV-A G gene =

397     1488, RSV-B G gene = 1120). Nucleotide sequences from each dataset were aligned using

398     MAFFT.v7 [41] and were translated into amino acids using EMBOSS.v6.6.0 [42] for

399     immunoinformatic analyses. In addition, two artificial sequences, CP248 and CP52 (cold passage

400     live RSV strains that were previously evaluated as vaccine candidates, Accession No: U63644,

401     AF0132551 respectively ) were downloaded from the NCBI's GenBank nucleotide database

402     [43].

403

404     **Phylogenetic inference**

405     The nucleotide sequences of  RSV major surface proteins were used to reconstruct the

406     maximum-likelihood (ML) phylogeny of RSV using RAxML.v8 with GTR+GAMMA

407     substitution model [44]. The best-scoring ML tree was automatically generated from five runs by

408     RAxML. Time-scaled phylogenies were further reconstructed with the best scoring ML trees

409     using the program "Timetree" [45]. The phylogenies are visualized in the R package "ggtree"

410     [46].

411

412     **T cell epitope prediction**

413     RSV major surface protein sequences were scored for binding potential against a globally

414     representative panel of Human Leukocyte Antigen (HLA) class I and class II alleles using the

415     EpiMatrix algorithm. This algorithm as well as the ClustiMer, JanusMatrix, and EpiCC

416     algorithms discussed below are part of the iVAX toolkit developed by EpiVax, which is

417     available for use under a license or through academic collaborations [24].

418

24

419    Evaluation of class I epitopes was made based on predictions for four HLA-A and two HLA-B

420    supertype alleles: A*01:01, A*02:01, A*03:01, A*24:02, B*07:02, B*44:03. Class II epitopes

421    were identified for nine HLA-DR supertype alleles: DRB1*01:01, DRB1*03:01, DRB1*04:01,

422    DRB1*07:01, DRB1*08:01, DRB1*09:01, DRB1*11:01, DRB1*13:01, and DRB1*15:01.

423    EpiMatrix parsed 9-mer sequence frames (each one overlapping the previous one by one amino

424    acid) from the antigen sequence and assigned a score for each nine-mer/allele pair on a

425    normalized Z distribution. Nine-mer sequences that had Z-scores of at least 1.64 are considered

426    to be in the top 5% of any randomly generated set of 9-mer sequences and to have a high

427    likelihood of binding to HLA molecules and being presented to T cells. Sequences that score

428    above 2.32 on the Z-scale (top 1%) are extremely likely to bind to a particular HLA allele and to

429    be immunogenic. For this analysis, HLA-class I restricted 9-mer sequences that had top 1%

430    binder scores to at least one HLA class I supertype allele were considered to be putative class I

431    epitopes [24].

432    To identify putative class II epitopes, we used an algorithm called ClustiMer [24] to screen

433    EpiMatrix scoring results for the nine class II alleles.  ClustiMer identifies contiguous regions of

434    15–30 amino acids that have a high density of MHC class II binding potential. Epitope density

435    within a cluster is reported as an EpiMatrix Cluster Score, where scores of 10 and above are

436    likely to be recognized in the context of multiple class II alleles and to be high-quality class II

437    epitopes.

438    We also applied analysis of human homology to this study. The JanusMatrix algorithm was used

439    to assess the potential cross-conservation of T cell epitopes with epitopes restricted by the same

440    HLA alleles in the human proteome [25]. Briefly, JanusMatrix scans each identified epitope and

441    examines the shared T cell receptor (TCR) contacts with class II epitopes present in the human

25

442     proteome, to compute a JanusMatrix Human Homology Score. As defined in retrospective

443     studies, foreign class I epitopes that score greater than 2 and class II epitopes that score greater

444     than 5 may be less immunogenic due to T cell tolerance.

445

446     **Protein-level T cell immunogenic potential evaluation**

447     RSV reference sequences (RSV-A: NC_038235, RSV-B: NC_001781) were downloaded from

448     the NCBI RefSeq database and were used to evaluate the protein-level immunogenic potential of

449     RSV major surface proteins. The protein-level immunogenic potential as represented by the

450     EpiMatrix-defined T cell epitope density score was computed by summing the top 5% binder

451     scores across HLA alleles and normalizing for a 1000-amino acid protein length. Zero on this

452     scale is set to indicate the average number of top 5% binders that would be observed in 10,000

453     random protein sequences with natural amino acid frequencies. Proteins scoring above +20 have

454     been observed to have the significant immunogenic potential [47]. Fully human proteins

455     generally score lower than zero on the EpiMatrix immunogenicity scale.

456     To investigate the distribution of T cell immunogenic potential across RSV protein sequence

457     regions, we summed up the binding scores of HLA alleles for each nine-mer frame, to get a

458     frame-specific immunogenic potential score and standardized this score to a relative scale. The

459     relative immunogenic potential across protein structure was represented by a color scale and the

460     visualization of F protein structure was built with PyMOL Molecular Graphics System, Version

461     2.0 (Schrödinger, LLC). Protein data bank (PDB) files 5UDE [48] and 3RRR [49] were used for

462     the pre-fusion and post-fusion forms.

463

464     **Subsampling strategy**

26

465    Considering the heavy computational load that would be required to evaluate all available RSV

466    sequences and to correct the overrepresentation of recently sampled strains, the comparative

467    analysis for T cell epitope content was conducted with datasets in which overrepresented groups

468    were reduced. A maximum of five sequences of each isolation year from different WHO region

469    groups were subsampled randomly from the original datasets (RSV-A F gene = 402, RSV-B F

470    gene = 319, RSV-A G gene = 390, RSV-B G gene = 359).

471

472    **T cell epitope content comparison**

473    The Epitope Content Comparison (EpiCC) algorithm, which is implemented in iVAX was used

474    to evaluate pairwise T cell epitope cross-conservation potential within each subsampled dataset

475    [50]. Briefly, T cell epitope cross-conservation was evaluated by the binding likelihood of

476    epitopes from different antigens with identical T cell receptor-facing residues (TCR$f$), which are

477    predicted to bind to the same MHC allele. Two epitope sequences are assumed to be potentially

478    cross-conserved if they have identical TCR$f$ (position 4, 5, 6, 7, 8 for class I epitopes binding

479    core and 2, 3, 5, 7, 8 for class II epitopes binding core) regardless of differences on their

480    MHC-facing amino acids. To simplify the analysis, the binding of 9-mer epitopes within protein

481    sequences are assumed to be mutually exclusive and uniform.

482

483    Therefore, T cell epitope immune distance ($D$) between two wild circulating strains ($w_1$ and $w_2$)

484    can be defined as the sum of Z-scaled binding probabilities of paired epitopes that are unable to

485    induce cross-reactivity (non-cross conserved epitopes) within two protein sequences using

486    equations (1.1 and 1.2). $d$ is the T cell immunity distance between a pair of epitopes, $i$ and $j$ are

27

487    the non-cross conserved T cell epitopes from two protein sequences, $a$ is a class I or class II

488    allele, $p$ is the predicted binding probability against allele $a$, A is a set of alleles.

489    $$d(i,j) = p(i)a + p(j)a \#(1.1)$$

490    $$D(w_1,w_2) = \sum_{i \in w_1 j \in w_2} \sum_{a \in A} d(i,j) \#(1.2)$$

491

492    To evaluate the T cell epitope immune distance generated by the EpiCC algorithm, we further

493    adapt the equations (1.1 and 1.2) to re-calculate T cell epitope immune distance with customized

494    Python scripts using  MHC binding prediction results that are generated from publicly available

495    T cell epitope prediction tool,  netMHCpan EL 4.1 methods in the Immune Epitope Database

496    (IEDB)  [51]. Eigenvalues of each sequence that were calculated from the pairwise distance

497    matrix with "RSpectra" package were used to statistically examine the correlation of the epitope

498    distances that are computed from the two methods, and Pearson correlation test was used to test

499    the correlation hypothesis.

500

501    The cross-conservation of vaccine strains against circulating RSV can be evaluated by the cross-

502    conservation of the epitopes within vaccine strains ($v$) and wild circulating strains ($w$). T cell

503    cross-conservation between two epitopes can be represented by a joint probability estimation and

504    therefore T cell cross-conservation between two sequences can be represented by summing T

505    cell cross-conservation of the paired T-epitopes within two protein sequences. The proportion of

506    T cell cross-conservation between the vaccine and circulating strains ($P$) with a set of alleles (A)

507    can be represented as the equations (2.1 and 2.2), where $p$ is the predicted binding probability in

508    EpiMatrix, $i$ and $j$ are the cross conserved T cell epitopes, $a$ is a class I or class II allele.

509    $$S(i,j) = p(i)a * p(j)a \#(2.1)$$

28

510
$$P(v,w) = \frac{\sum_{i \in v, j \in w} \sum_{a \in A} S(i,j)}{\sum_{i \in v, j \in v} \sum_{a \in A} S(i,j)} \#(2.2)$$

511

512 **Dimension reduction**

513 The equation to calculate T cell epitope immune distance was applied iteratively to the

514 subsampled dataset and therefore the pairwise T cell epitope immune distances are structured

515 into an *n x n* square-distance matrix. Given that each protein is described by a relative distance to

516 the rest of n-1 proteins, the data must be dimensionally reduced to be graphed. Classic (metric)

517 multidimensional scaling (MDS) can be used to preserve the distances between a set of

518 observations in a way that allows the distances to be represented in a two-dimensional space.

519 MDS was performed as previously described by Gower [52]. The MDS method first constructs

520 an n-dimensional Euclidean space using the distance matrix in which all distances are conserved,

521 and then principal component analysis is performed. MDS and Goodness-of-fit (GOF) [53] were

522 carried out using the *cmdscale* package in R [52]. K-means clustering was performed using

523 the *kmeans* function in base R. Due to the lack of previous characterizations of RSV T cell

524 immunity clusters, the number of T cell immunity groups was determined using the optimized

525 within-cluster sum of square (wss) with Elbows method [54].

526

527 **Calculation of genetic hamming distance**

528 Genetic hamming distance, which is defined as the number of bases by which two nucleotide

529 sequences differ, was calculated by comparing the number of different bases between each

530 sequence in the subsampled datasets. The reconstructed most recent common ancestor (TMRCA)

531 sequences for each dataset (subsampled F and G protein sequences of subtype A and subtype B,

29

532    respectively) were estimated using the program "Treetime" and were used as root in our analysis

533    [45].

534

542

543    **Competing interests:** A.S.DeG. and L.M are both paid employees of EpiVax. Some of the

544    epitope prediction tools used in this study were developed by EpiVax.

545

546    **Data and materials availability:** Accession number to RSV sequence in the paper are available

547    in supplementary materials. Code to generate T epitope landscapes are deposited in GitHub

548    https://github.com/JianiC/RSV_Epitope.

549

550    **References**

551    [1]    J. E. Crowe Jr and J. V Williams, "Paramyxoviruses: Respiratory Syncytial Virus and

552          Human Metapneumovirus," *Viral Infect. Humans Epidemiol. Control*, pp. 601–627, Feb.

553          2014, doi: 10.1007/978-1-4899-7448-8_26.

554    [2]    W.-J. Lee, Y. Kim, D.-W. Kim, H. S. Lee, H. Y. Lee, and K. Kim, "Complete Genome

555        Sequence of Human Respiratory Syncytial Virus Genotype A with a 72-Nucleotide

556        Duplication in the Attachment Protein G Gene," *J. Virol.*, vol. 86, no. 24, pp. 13810 LP –

557        13811, Dec. 2012, doi: 10.1128/JVI.02571-12.

558   [3]   J. S. McLellan, W. C. Ray, and M. E. Peeples, "Structure and function of respiratory

559        syncytial virus surface glycoproteins," *Curr. Top. Microbiol. Immunol.*, vol. 372, pp. 83–

560        104, 2013, doi: 10.1007/978-3-642-38919-1_4.

561   [4]   J. Lee, L. Klenow, E. M. Coyle, H. Golding, and S. Khurana, "Protective antigenic sites in

562        respiratory syncytial virus G attachment protein outside the central conserved and cysteine

563        noose domains," *PLoS Pathog.*, vol. 14, no. 8, pp. e1007262–e1007262, Aug. 2018, doi:

564        10.1371/journal.ppat.1007262.

565   [5]   "Updated guidance for palivizumab prophylaxis among infants and young children at

566        increased risk of hospitalization for respiratory syncytial virus infection.," *Pediatrics*, vol.

567        134, no. 2, pp. 415–420, Aug. 2014, doi: 10.1542/peds.2014-1665.

568   [6]   "Palivizumab, a humanized respiratory syncytial virus monoclonal antibody, reduces

569        hospitalization from respiratory syncytial virus infection in high-risk infants. The IMpact-

570        RSV Study Group.," *Pediatrics*, vol. 102, no. 3 Pt 1, pp. 531–537, Sep. 1998.

571   [7]   T. F. Schwarz *et al.*, "Three dose levels of a maternal respiratory syncytial virus vaccine

572        candidate are  well tolerated and immunogenic in a randomized trial in non-pregnant

573        women.," *J. Infect. Dis.*, Jun. 2021, doi: 10.1093/infdis/jiab317.

574   [8]   C. Biagi *et al.*, "Current State and Challenges in Developing Respiratory Syncytial Virus

575        Vaccines," *Vaccines*, vol. 8, no. 4, p. 672, Nov. 2020, doi: 10.3390/vaccines8040672.

576   [9]   H. W. Kim *et al.*, "Respiratory syncytial virus disease in infants despite prior

577        administration of  antigenic inactivated vaccine.," *Am. J. Epidemiol.*, vol. 89, no. 4, pp.

578        422–434, Apr. 1969, doi: 10.1093/oxfordjournals.aje.a120955.

579  [10]  B. R. Murphy, A. V Sotnikov, L. A. Lawrence, S. M. Banks, and G. A. Prince, "Enhanced

580        pulmonary histopathology is observed in cotton rats immunized with formalin-inactivated

581        respiratory syncytial virus (RSV) or purified F glycoprotein and challenged with RSV 3-6

582        months after immunization.," *Vaccine*, vol. 8, no. 5, pp. 497–502, Oct. 1990, doi:

583        10.1016/0264-410x(90)90253-i.

584  [11]  A. M. Killikelly, M. Kanekiyo, and B. S. Graham, "Pre-fusion F is absent on the surface

585        of formalin-inactivated respiratory syncytial virus," *Sci. Rep.*, vol. 6, p. 34108, Sep. 2016,

586        doi: 10.1038/srep34108.

587  [12]  F. P. Polack *et al.*, "A role for immune complexes in enhanced respiratory syncytial virus

588        disease.," *J. Exp. Med.*, vol. 196, no. 6, pp. 859–865, Sep. 2002, doi:

589        10.1084/jem.20020781.

590  [13]  M. Connors, N. A. Giese, A. B. Kulkarni, C. Y. Firestone, H. C. 3rd Morse, and B. R.

591        Murphy, "Enhanced pulmonary histopathology induced by respiratory syncytial virus

592        (RSV) challenge of formalin-inactivated RSV-immunized BALB/c mice is abrogated by

593        depletion of interleukin-4 (IL-4) and IL-10.," *J. Virol.*, vol. 68, no. 8, pp. 5321–5325, Aug.

594        1994, doi: 10.1128/JVI.68.8.5321-5325.1994.

595  [14]  M. E. Waris, C. Tsou, D. D. Erdman, S. R. Zaki, and L. J. Anderson, "Respiratory synctial

596        virus infection in BALB/c mice previously immunized with formalin-inactivated virus

597        induces enhanced pulmonary inflammatory response with a predominant Th2-like

598        cytokine pattern.," *J. Virol.*, vol. 70, no. 5, pp. 2852–2860, May 1996, doi:

599        10.1128/JVI.70.5.2852-2860.1996.

600  [15]  N. I. Mazur *et al.*, "The respiratory syncytial virus vaccine landscape: lessons from the

601    graveyard and  promising candidates.," *Lancet. Infect. Dis.*, vol. 18, no. 10, pp. e295–

602    e311, Oct. 2018, doi: 10.1016/S1473-3099(18)30292-5.

603  [16]  B. N. Blunck, W. Rezende, and P. A. Piedra, "Profile of respiratory syncytial virus

604    prefusogenic fusion protein nanoparticle  vaccine.," *Expert Rev. Vaccines*, vol. 20, no. 4,

605    pp. 351–364, Apr. 2021, doi: 10.1080/14760584.2021.1903877.

606  [17]  X. Liang *et al.*, "Gradual replacement of all previously circulating respiratory syncytial

607    virus A strain with the novel ON1 genotype in Lanzhou from 2010 to 2017," *Medicine*

608    *(Baltimore).*, vol. 98, no. 19, 2019.

609  [18]  A. Ahmed *et al.*, "Co-Circulation of 72bp Duplication Group A and 60bp Duplication

610    Group B Respiratory Syncytial Virus (RSV) Strains in Riyadh, Saudi Arabia during

611    2014," *PLoS One*, vol. 11, no. 11, pp. e0166145–e0166145, Nov. 2016, doi:

612    10.1371/journal.pone.0166145.

613  [19]  D. Tian *et al.*, "Structural basis of respiratory syncytial virus subtype-dependent

614    neutralization by an antibody targeting the fusion glycoprotein," *Nat. Commun.*, vol. 8, no.

615    1, p. 1877, 2017, doi: 10.1038/s41467-017-01858-w.

616  [20]  K. Hashimoto and M. Hosoya, "Neutralizing epitopes of RSV and palivizumab resistance

617    in Japan," *Fukushima J. Med. Sci.*, vol. 63, no. 3, pp. 127–134, Dec. 2017, doi:

618    10.5387/fms.2017-09.

619  [21]  W. M. Sullender, "Respiratory syncytial virus genetic and antigenic diversity," *Clin.*

620    *Microbiol. Rev.*, vol. 13, no. 1, pp. 1–15, Jan. 2000, doi: 10.1128/CMR.13.1.1.

621  [22]  X. Chen *et al.*, "Genetic variations in the fusion protein of respiratory syncytial virus

622    isolated from children hospitalized with community-acquired pneumonia in China," *Sci.*

623    *Rep.*, vol. 8, no. 1, p. 4491, 2018, doi: 10.1038/s41598-018-22826-4.

624  [23]  G. Petrova, A. Ferrante, and J. Gorski, "Cross-reactivity of T cells and its role in the

625        immune system," *Crit. Rev. Immunol.*, vol. 32, no. 4, pp. 349–372, 2012, doi:

626        10.1615/critrevimmunol.v32.i4.50.

627  [24]  A. S. De Groot *et al.*, "Better Epitope Discovery, Precision Immune Engineering, and

628        Accelerated Vaccine Design Using Immunoinformatics Tools   ," *Frontiers in*

629        *Immunology* , vol. 11. p. 442, 2020.

630  [25]  L. He, A. S. De Groot, A. H. Gutierrez, W. D. Martin, L. Moise, and C. Bailey-Kellogg,

631        "Integrated assessment of predicted MHC binding and cross-conservation with self

632        reveals patterns of viral camouflage.," *BMC Bioinformatics*, vol. 15 Suppl 4, no. Suppl 4,

633        p. S1, 2014, doi: 10.1186/1471-2105-15-S4-S1.

634  [26]  B. T. Grenfell, "Unifying the Epidemiological and Evolutionary Dynamics of Pathogens,"

635        *Science (80-. ).*, vol. 303, no. 5656, pp. 327–332, Jan. 2004, doi:

636        10.1126/science.1090727.

637  [27]  S. D. J. *et al.*, "Mapping the Antigenic and Genetic Evolution of Influenza Virus," *Science*

638        *(80-. ).*, vol. 305, no. 5682, pp. 371–376, Jul. 2004, doi: 10.1126/science.1097211.

639  [28]  T. Bedford *et al.*, "Integrating influenza antigenic dynamics with molecular evolution,"

640        *Elife*, vol. 3, Feb. 2014, doi: 10.7554/eLife.01914.

641  [29]  B. Korber, M. LaBute, and K. Yusim, "Immunoinformatics comes of age," *PLoS Comput.*

642        *Biol.*, vol. 2, no. 6, pp. e71–e71, Jun. 2006, doi: 10.1371/journal.pcbi.0020071.

643  [30]  B. Reynisson, B. Alvarez, S. Paul, B. Peters, and M. Nielsen, "NetMHCpan-4.1 and

644        NetMHCIIpan-4.0: improved predictions of MHC antigen presentation by concurrent

645        motif deconvolution and integration of MS MHC eluted ligand data," *Nucleic Acids Res.*,

646        vol. 48, no. W1, pp. W449–W454, Jul. 2020, doi: 10.1093/nar/gkaa379.

34

647 [31]  C. D. Russell, S. A. Unger, M. Walton, and J. Schwarze, "The Human Immune Response

648       to Respiratory Syncytial Virus Infection.," *Clin. Microbiol. Rev.*, vol. 30, no. 2, pp. 481–

649       502, Apr. 2017, doi: 10.1128/CMR.00090-16.

650 [32]  R. G. Woolthuis, C. H. van Dorp, C. Keşmir, R. J. de Boer, and M. van Boven, "Long-

651       term adaptation of the influenza A virus by escaping cytotoxic T-cell recognition," *Sci.*

652       *Rep.*, vol. 6, no. 1, p. 33334, 2016, doi: 10.1038/srep33334.

653 [33]  L. C. Katzelnick *et al.*, "Dengue viruses cluster antigenically but not as discrete

654       serotypes," *Science*, vol. 349, no. 6254, pp. 1338–1343, Sep. 2015, doi:

655       10.1126/science.aac5017.

656 [34]  "Correction: Antigenic cartography of immune responses to Plasmodium falciparum

657       erythrocyte membrane protein 1 (PfEMP1)," *PLOS Pathog.*, vol. 15, no. 8, p. e1008018,

658       Aug. 2019.

659 [35]  D. J. Smith *et al.*, "Mapping the antigenic and genetic evolution of influenza virus,"

660       *Science (80-. ).*, vol. 305, no. 5682, pp. 371–376, Jul. 2004, doi: 10.1126/science.1097211.

661 [36]  J. Liu, T. J. Ruckwardt, M. Chen, T. R. Johnson, and B. S. Graham, "Characterization of

662       respiratory syncytial virus M- and M2-specific CD4 T cells in a murine model," *J. Virol.*,

663       vol. 83, no. 10, pp. 4934–4941, May 2009, doi: 10.1128/JVI.02140-08.

664 [37]  C. Viboud *et al.*, "Beyond clinical trials: Evolutionary and epidemiological considerations

665       for development of a universal influenza vaccine," *PLOS Pathog.*, vol. 16, no. 9, p.

666       e1008583, Sep. 2020.

667 [38]  J. Chen *et al.*, "Novel and extendable genotyping system for Human Respiratory Syncytial

668       Virus based on whole-genome sequence analysis," *Authorea Prepr.*, Jun. 2021, doi:

669       10.22541/AU.162251650.01202619/V1.

670 [39] S. S. Shepard, C. T. Davis, J. Bahl, P. Rivailler, I. A. York, and R. O. Donis, "LABEL:

671 Fast and Accurate Lineage Assignment with Assessment of H5N1 and H9N2 Influenza A

672 Hemagglutinins," *PLoS One*, vol. 9, no. 1, p. e86921, Jan. 2014.

673 [40] "WHO | World Health Statistics 2011," *WHO*, 2011.

674 [41] J. Rozewicki, S. Li, K. M. Amada, D. M. Standley, and K. Katoh, "MAFFT-DASH:

675 integrated protein sequence and structural alignment," *Nucleic Acids Res.*, vol. 47, no.

676 W1, pp. W5–W10, Jul. 2019, doi: 10.1093/nar/gkz342.

677 [42] P. Rice, I. Longden, and A. Bleasby, "EMBOSS: The European Molecular Biology Open

678 Software Suite," *Trends Genet.*, vol. 16, no. 6, pp. 276–277, Jun. 2000, doi:

679 10.1016/S0168-9525(00)02024-2.

680 [43] S. S. Whitehead *et al.*, "Replacement of the F and G Proteins of Respiratory Syncytial

681 Virus (RSV) Subgroup A with Those of Subgroup B Generates Chimeric Live Attenuated

682 RSV Subgroup B Vaccine Candidates," *J. Virol.*, vol. 73, no. 12, pp. 9773 LP – 9780,

683 Dec. 1999, doi: 10.1128/JVI.73.12.9773-9780.1999.

684 [44] A. Stamatakis, "RAxML version 8: a tool for phylogenetic analysis and post-analysis of

685 large phylogenies," *Bioinforma. Appl.*, vol. 30, no. 9, pp. 1312–1313, 2014, doi:

686 10.1093/bioinformatics/btu033.

687 [45] P. Sagulenko, V. Puller, and R. A. Neher, "TreeTime: Maximum-likelihood

688 phylodynamic analysis," *Virus Evol.*, vol. 4, no. 1, Jan. 2018, doi: 10.1093/VE/VEX042.

689 [46] G. Yu, D. K. Smith, H. Zhu, Y. Guan, and T. T. Y. Lam, "ggtree: an r package for

690 visualization and annotation of phylogenetic trees with their covariates and other

691 associated data," *Methods Ecol. Evol.*, vol. 8, no. 1, pp. 28–36, Jan. 2017, doi:

692 10.1111/2041-210X.12628.

693   [47]  L. Moise *et al.*, "iVAX: An integrated toolkit for the selection and optimization of

694        antigens and the design of epitope-driven vaccines," *Hum. Vaccin. Immunother.*, vol. 11,

695        no. 9, pp. 2312–2321, 2015, doi: 10.1080/21645515.2015.1061159.

696   [48]  Q. Zhu *et al.*, "A highly potent extended half-life antibody as a potential RSV vaccine

697        surrogate for all infants," *Sci. Transl. Med.*, vol. 9, no. 388, p. eaaj1928, May 2017, doi:

698        10.1126/scitranslmed.aaj1928.

699   [49]  J. S. McLellan, Y. Yang, B. S. Graham, and P. D. Kwong, "Structure of respiratory

700        syncytial virus fusion glycoprotein in the postfusion conformation reveals preservation of

701        neutralizing epitopes," *J. Virol.*, vol. 85, no. 15, pp. 7788–7796, Aug. 2011, doi:

702        10.1128/JVI.00555-11.

703   [50]  A. H. Gutiérrez *et al.*, "T-cell epitope content comparison (EpiCC) of swine H1 influenza

704        A virus hemagglutinin," *Influenza Other Respi. Viruses*, vol. 11, no. 6, pp. 531–542, Nov.

705        2017, doi: 10.1111/irv.12513.

706   [51]  E. Karosiene, C. Lundegaard, O. Lund, and M. Nielsen, "NetMHCcons: a consensus

707        method for the major histocompatibility complex class I predictions," *Immunogenetics*,

708        vol. 64, no. 3, pp. 177–186, 2012, doi: 10.1007/s00251-011-0579-8.

709   [52]  J. C. GOWER, "Some distance properties of latent root and vector methods used in

710        multivariate analysis," *Biometrika*, vol. 53, no. 3–4, pp. 325–338, Dec. 1966, doi:

711        10.1093/biomet/53.3-4.325.

712   [53]  J. Sánchez, "MARDIA, K. V., J. T. KENT, J. M. BIBBY: Multivariate Analysis.

713        Academic Press, London-New York-Toronto-Sydney-San Francisco 1979. xv, 518 pp., $

714        61.00," *Biometrical J.*, vol. 24, no. 5, p. 502, Jan. 1982, doi:

715        https://doi.org/10.1002/bimj.4710240520.

716 [54] M. Charrad, N. Ghazzali, V. Boiteau, and A. Niknafs, "NbClust: An R Package for

717 Determining the Relevant Number of Clusters in a Data Set," *J. Stat. Software; Vol 1,*

718 *Issue 6* , 2014, doi: 10.18637/jss.v061.i06.

719

720 **Supporting information**

721 **S1 Fig. Distribution and diversity of T cell epitopes in RSV F protein.** The tree panel on the

722 left is a time-scaled phylogeny build with RSV-A **(A)** or RSV-B **(B)** F gene nucleotide

723 sequences using the ML approach. Determined genotypes are labeled on the right with black

724 bars. Each color column on the right side represents the presence of an MHC class I or class II

725 epitope. Only the epitopes that are present in more than 1% of sampled isolates are displayed.

726 The column color indicates different numbers of epitope sequences at the same location.

727

728 **S2 Fig. Distribution and diversity of T cell epitopes in RSV G protein.** The tree panel on the

729 left is a time-scaled phylogeny build with RSV-A **(A)** or RSV-B **(B)** G gene nucleotide

730 sequences using the ML approach. The clades that contain novel 72-nt or 60-nt duplication at the

731 second hypervariable region of G gene were highlighted in red. Determined genotypes are

732 labeled on the right with black bars. Each color column on the right side represents the presence

733 of an MHC class I or class II epitope. Only the epitopes that are present in more than 1% of

734 sampled isolates were displayed. The column color indicates different numbers of epitope

735 sequences at the same location.

736

737 **S3 Fig. Distribution of JanusMatrix Human Homology score for putative RSV MHC class I**

738 **and class II epitopes.** The cross-reactive potential of identified putative T cell epitopes and

739    human host was represented with a JanusMatrix Human Homology score. 6.45% identified

740    putative class I epitopes and 1.12% class II epitopes are cross-conserved on the TCR face with

741    human epitopes.

742

743    **S4 Fig. Predicted T cell epitope landscapes of RSV surface proteins.** RSV T cell epitope

744    landscapes were built with sequenced-based MHC class I epitope binding prediction (left), MHC

745    class II epitope binding prediction (middle) or combining class I and class II epitope biding

746    prediction (right). Sequences are colored by the epitope cluster determined by epitope landscapes

747    built with combining Class I and Class II epitope prediction

748

749    **S5 Fig. Total within sum of squares (*wss*) using *k-means* algorithm.** Totals within sum of

750    squares in epitope topographies were calculated after clustering into k (from 1 to 10) groups with

751    *k-means*. The optimal number of clusters is determined to be 3 in the analysis of RSV-A F and G

752    proteins and is determined to be 2 in the analysis of RSV-B F and G proteins using the Elbow

753    method.

754

755    **S6 Fig. Validation of T cell epitope distance estimation using the IEDB analysis resource**.

756    Validation is performed with MHC class I epitope binding prediction of RSV-A F protein. **(A)**

757    Heatmaps for pairwise MHC class I epitope distance estimated in iVAX toolkits or calculated

758    with custom python scripts using MHC class I molecule binding prediction that is implemented

759    in IEDB. **(B)** Eigenvalues for each sequence are calculated from pairwise distance matrices using

760    "RSpectra" package in R. The Pearson correlation test significantly supports a non-zero

761    correlation between T cell epitope distance estimated with EpiCC and T cell epitope distance

39

762    estimated with IEDB. **(C)** T cell epitope topographies are built with pairwise epitope distances

763    estimated from EpiCC or IEDB. Both methods resulted in a similar cluster pattern for the CD8 T

764    cell epitope profile of RSV-A F protein.

765

766    **S7 Fig. Evaluation of RSV vaccine candidate strains with T cell epitope content in different**

767    **WHO regions**. RSV-A and RSV-B major surface protein sequences were grouped by isolation

768    year and 6 isolated WHO regions, African Region (AFRO), Region of the Americas (PAHO),

769    South-East Asia Region (SEARO), European Region (EURO), Eastern Mediterranean Region

770    (EMRO) and Western Pacific Region (WPRO). Each year group was labeled by the latest

771    isolated year of sequences after the previous group label. The proportion of cross-conserved T

772    cell epitope content between vaccine strains (CP248 or CP52) and wild circulating strains in

773    different year groups was represented by bar graphs.

774

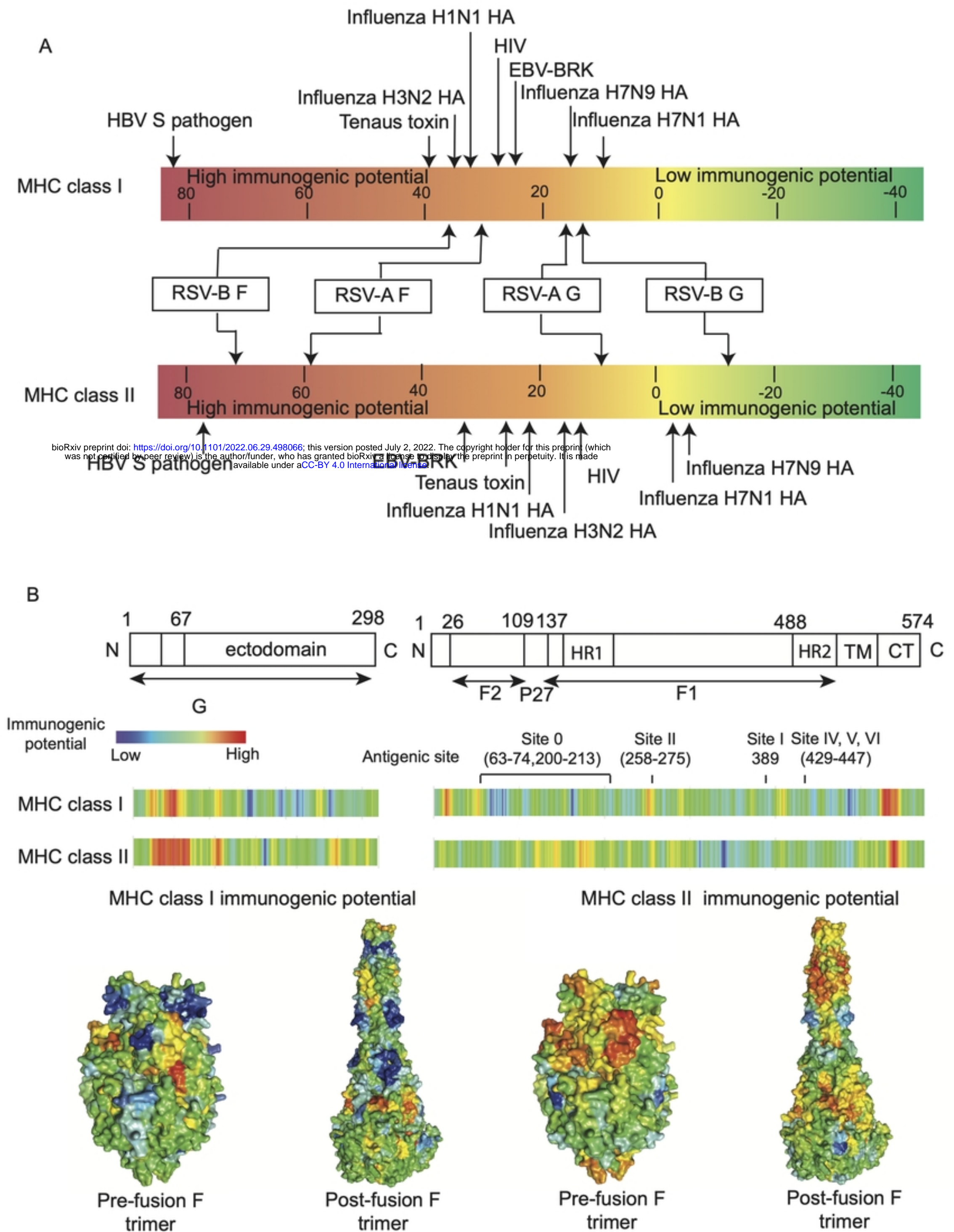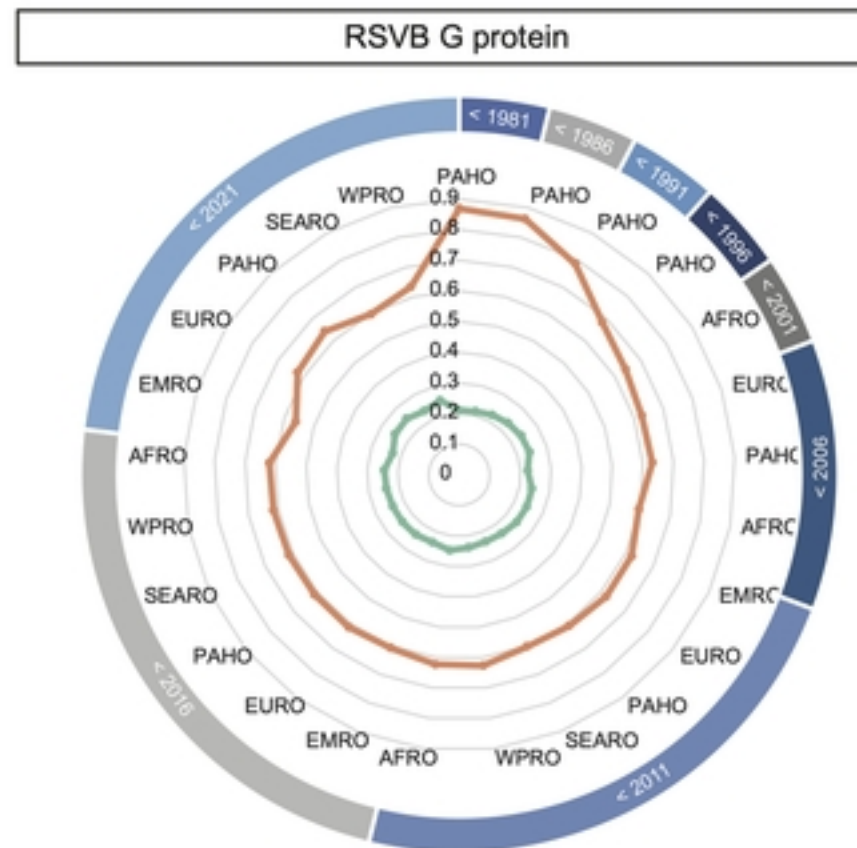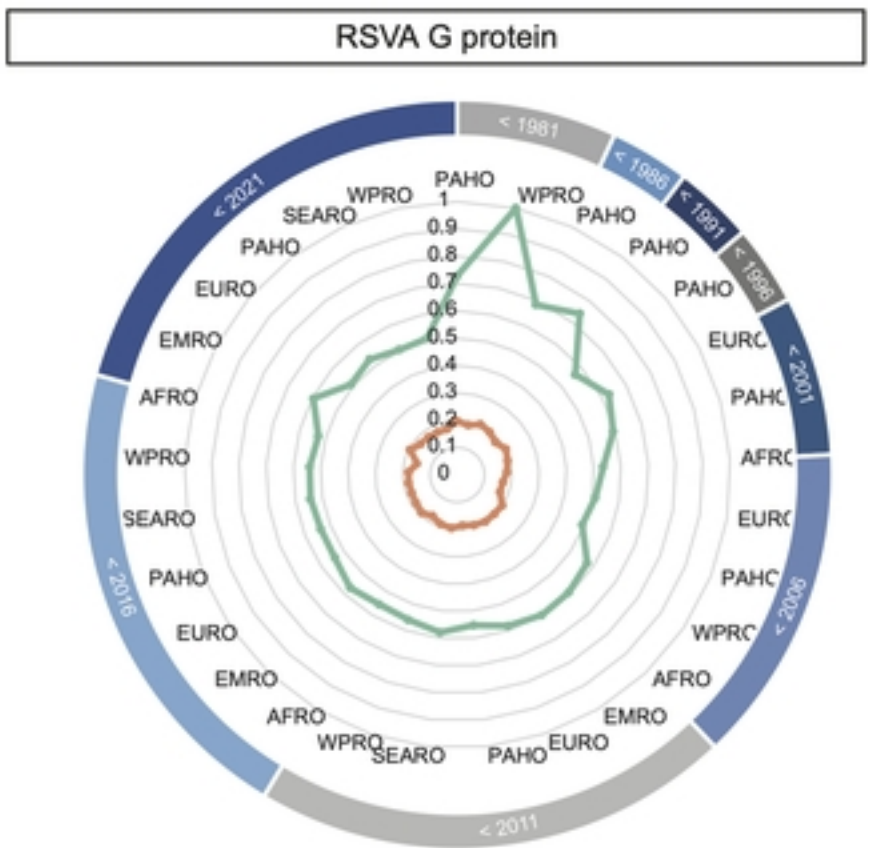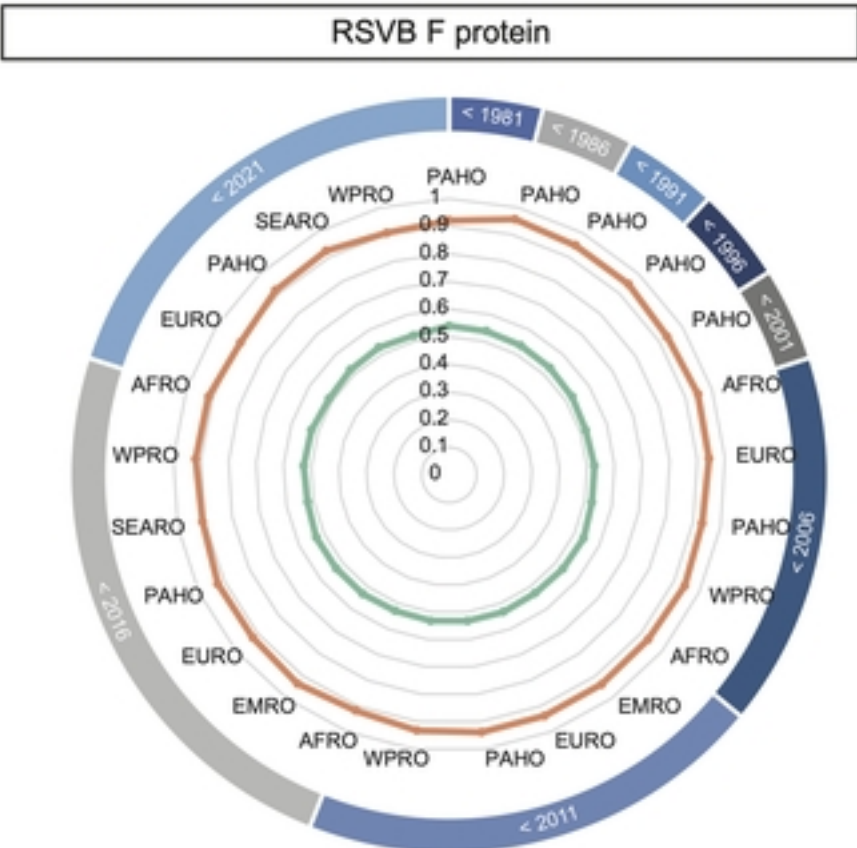775    **S1 File. Accession number to RSV sequence that are used in this study.**

A

Influenza H1N1 HA
HIV
EBV-BRK
Influenza H3N2 HA       Influenza H7N9 HA
Tenaus toxin            Influenza H7N1 HA

HBV S pathogen

**MHC class I**

High immunogenic potential          Low immunogenic potential

80    60    40    20    0    -20    -40

RSV-B F    RSV-A F    RSV-A G    RSV-B G

**MHC class II**

80    60    40    20    0    -20    -40

High immunogenic potential          Low immunogenic potential

HBV S pathogen                  RK

Tenaus toxin
Influenza H1N1 HA                HIV          Influenza H7N9 HA
Influenza H3N2 HA                             Influenza H7N1 HA

B

1    67                    298          1    26    109 137                            488          574

N  [    ectodomain    ] C          N [            HR1                              HR2  TM  CT ] C

G                                   F2  P27              F1

Immunogenic
potential
Low ──── High

Antigenic site    Site 0                Site II        Site I    Site IV, V, VI
                  (63-74,200-213)       (258-275)      389       (429-447)

MHC class I

MHC class II

MHC class I immunogenic potential          MHC class II immunogenic potential

Pre-fusion F       Post-fusion F          Pre-fusion F       Post-fusion F
trimer             trimer                 trimer             trimer

Figure1

Figure3

- Without G duplication ◇ With G duplication

A) T-cell immune epitope landscape
B) Phylogeny
C) T-cell epitope immune distance to reconstructed ancestral
D) Genetic hamming distance to reconstructed ancestral

Figure2