
TOTAL PREDICTED MHC-I EPITOPE LOAD IS INVERSELY ASSOCIATED WITH POPULATION MORTALITY FROM SARS-CoV-2

Eric A. Wilson
School of Molecular Sciences
Arizona State University

Gabrielle Hirneise
School of Life Sciences
Arizona State University

Abhishek Singharoy*
School of Molecular Sciences
Arizona State University

Karen S. Anderson*
Biodesign Institute
Arizona State University

December 15, 2020

ABSTRACT

Polymorphisms in MHC-I protein sequences across human populations significantly impacts viral peptide binding capacity and thus alters T cell immunity to infection. Consequently, allelic variants of the MHC-I protein have been found to be associated with patient outcome to various viral infections, including SARS-CoV. In the present study, we assess the relationship between observed SARS-CoV-2 population mortality and the predicted viral binding capacities of 52 common MHC-I alleles. Potential SARS-CoV-2 MHC-I peptides were identified using a consensus MHC-I binding and presentation prediction algorithm, called EnsembleMHC. Starting with nearly 3.5 million candidates, we resolved a few hundred highly probable MHC-I peptides. By weighing individual MHC allele-specific SARS-CoV-2 binding capacity with population frequency in 23 countries, we discover a strong inverse correlation between the predicted population SARS-CoV-2 peptide binding capacity and observed mortality rate. Our computations reveal that peptides derived from the structural proteins of the virus produces a stronger association with observed mortality rate, highlighting the importance of S, N, M, E proteins in driving productive immune responses. The correlation between epitope binding capacity and population mortality risk remains robust across a range of socioeconomic and epidemiological factors. A combination of binding capacity, number of deaths due to COPD complications, gender demographics, and the proportions of the population that were over the age of 65 and overweight offered the strongest determinant of at-risk populations. These results bring to light how molecular changes in the MHC-I proteins may affect population-level outcomes of viral infection.

Keywords SARS-CoV-2 · EnsembleMHC · MHC-I · risk-model · binding predictions · consensus models · population dynamics

*corresponding author

0 Introduction

1 In December 2019, the novel coronavirus, SARS-CoV-2 was
2 identified from a cluster of cases of pneumonia in Wuhan,
3 China^{1,2}. With over 73.1 million cases and over 1.6 million
4 deaths, the viral spread has been declared a global pandemic by
5 the World Health Organization³. Due to its high rate of trans-
6 mission and unpredictable severity, there is an immediate need
7 for information surrounding the adaptive immune response
8 towards SARS-CoV-2.

9 A robust T cell response is integral for the clearance of
10 coronaviruses, and generation of lasting immunity⁴. The poten-
11 tial role of T cells for coronavirus clearance has been supported
12 by the identification of immunogenic CD8⁺ T cell epitopes in
13 the S (Spike), N (Nucleocapsid), M (Membrane), and E (En-
14 velope) proteins⁵. Additionally, SARS-CoV specific CD8⁺ T
15 cells have been shown to provide long lasting immunity with
16 memory CD8⁺ T cells being detected up to 17 years post infec-
17 tion^{4,6,7}. The specifics of the T cell response to SARS-CoV-
18 2 peptides revealed a majority of the CD8⁺ T cell immune
19 response is targeted towards viral structural proteins (N, M,
20 S)⁸.

22 A successful CD8⁺ T cell response is contingent on the
23 efficient presentation of viral protein fragments by Major Histo-
24 compatibility Complex I (MHC-I) proteins. MHC-I molecules
25 bind and present peptides derived from endogenous proteins
26 on the cell surface for CD8⁺ T cell interrogation. The MHC-I
27 protein is highly polymorphic, with amino acid substitutions
28 within the peptide binding groove drastically altering the com-
29 position of presented peptides. Consequently, the influence
30 of MHC genotype to shape patient outcome has been well
31 studied in the context of viral infections⁹. For coronaviruses,
32 there have been several studies of MHC association with dis-
33 ease susceptibility. A study of a Taiwanese and Hong Kong
34 cohort of patients with SARS-CoV found that the MHC-I alle-
35 les HLA-B*07:03 and HLA-B*46:01 were linked to increased
36 susceptibility while HLA-Cw*15:02 was linked to increased
37 resistance¹⁰⁻¹². However, some of the reported associations
38 did not remain after statistical correction, and it is still un-
39 clear if MHC-outcome associations reported for SARS-CoV
40 are applicable to SARS-CoV-2^{13,14}. Recently, a comprehensive
41 prediction of SARS-CoV-2 MHC-I peptides indicated a relative
42 depletion of high affinity binding peptides for HLA-B*46:01,
43 hinting at a similar association profile in SARS-CoV-2¹⁵. More
44 importantly, it remains elusive if such a depletion of putative
45 high affinity peptides will impact patient outcome to SARS-
46 CoV-2 infections.

47 The lack of large scale genomic data linking individual
48 MHC genotype and outcome from SARS-CoV-2 infections

precludes a similar analysis as performed for SARS-CoV¹⁰⁻¹².
Therefore, we endeavored to assess the relationship between
the predicted SARS-CoV-2 binding capacity of a population
and the observed SARS-CoV-2 mortality rate. However, cur-
rent MHC-I prediction algorithms have been characterized by
a high false positive rate particularly when predicting peptides
that are naturally presented^{16,17}. To mitigate false positives and
identify the highest confidence SARS-CoV-2 MHC-I peptides,
we developed a consensus prediction algorithm, coined En-
sembleMHC, and predicted MHC-I peptides for a panel of 52
common MHC-I alleles¹⁸. This prediction workflow integrates
seven different algorithms that have been parameterized on
high-quality mass spectrometry data and provides a confidence
level for each identified peptide^{17,19-24}. The distribution of the
number of high-confidence peptides assigned to each allele was
used to assess a country-specific SARS-CoV-2 binding capacity,
called the EnsembleMHC population score, for 23 countries
(for selection criteria, please refer to the methods). This score
was derived by weighing the individual binding capacities of
the 52 MHC-I alleles by their endemic frequencies. We observe
a strong inverse correlation between the EnsembleMHC popu-
lation score and observed population SARS-CoV-2 mortality.
Furthermore, the correlation is shown to become stronger when
considering EnsembleMHC population scores based solely on
SARS-CoV-2 structural proteins, underlining their potential im-
portance in driving a robust immune response. Based on their
predicted binding affinity, expression, and sequence conserva-
tion in viral isolates, we identified 108 peptides derived from
SARS-CoV-2 structural proteins that are high-value targets for
CD8⁺ T cell vaccine development.

Results

**EnsembleMHC workflow offers more precise MHC-I pre-
sentation predictions than individual algorithms.** The ac-
curate assessment of differences in SARS-CoV-2 binding ca-
pacities across MHC-I allelic variants requires the isolation
of MHC-I peptides with a high probability of being presented.
EnsembleMHC provides the requisite precision through the use
of allele and algorithm-specific score thresholds and peptide
confidence assignment.

MHC-I alleles substantially vary in both peptide bind-
ing repertoire size and median binding affinity²⁵. The Ensem-
bleMHC workflow addresses this inter-allele variation by iden-
tifying peptides based on MHC allele and algorithm-specific
binding affinity thresholds. These thresholds were set by bench-
marking each of the seven component algorithms against 52
single MHC allele peptide data sets¹⁷. Each data set consists of
mass spectrometry-confirmed MHC-I peptides that have been
naturally presented by a model cell line expressing one of the 52

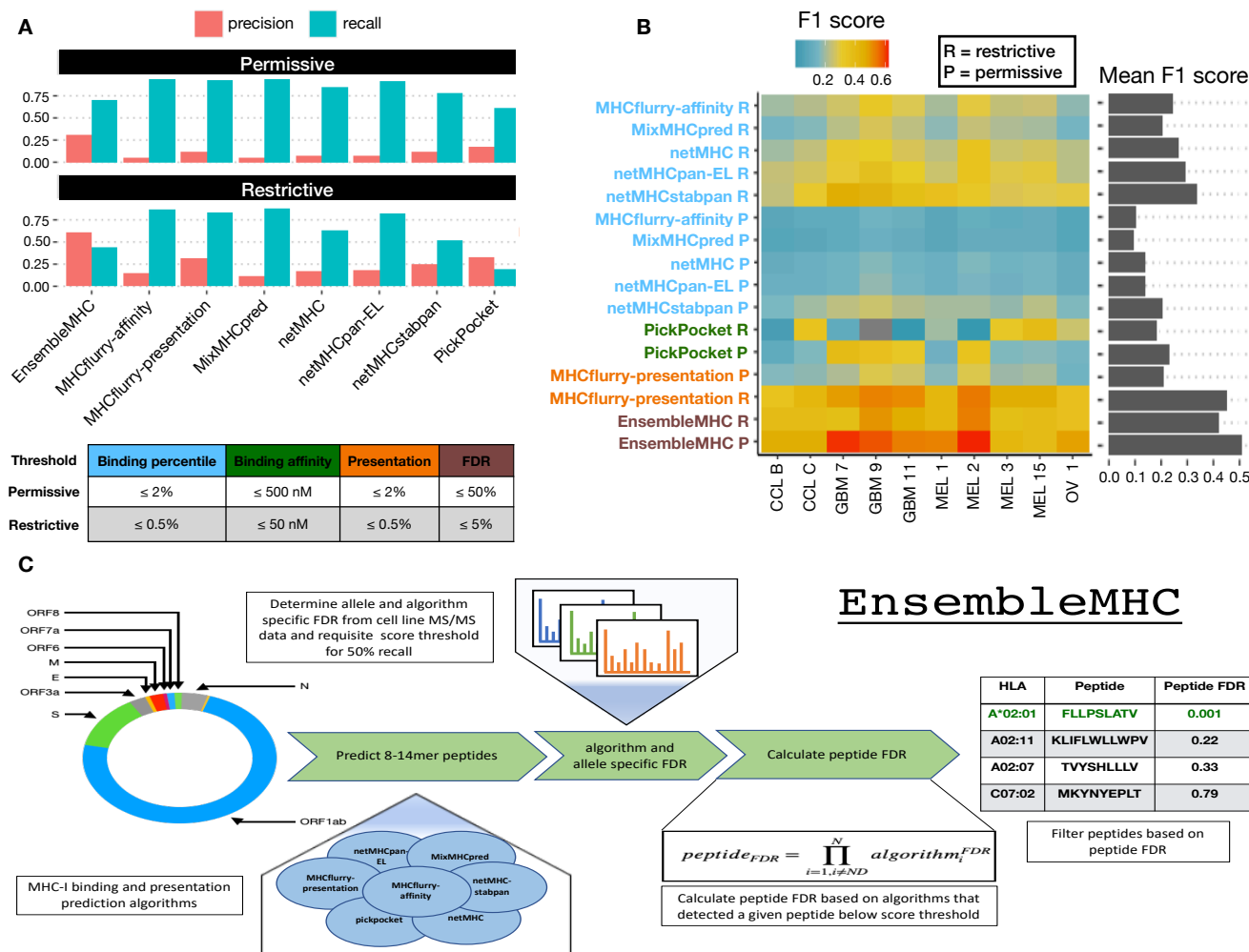


Figure 1: Application of the EnsembleMHC prediction algorithm. The EnsembleMHC prediction algorithm was used to recover MHC-I peptides from 10 tumor sample data sets. **A**, The average precision and recall for EnsembleMHC and each component algorithm was calculated across all 10 tumor samples. Peptide identification by each algorithm was based on commonly used restrictive (strong) or permissive (strong and weak) binding affinity thresholds (**inset table**). **B**, The F1 score of each algorithm was calculated for all tumor samples. Each algorithm is grouped into 1 of 4 categories: binding affinity represented by percentile score (blue), binding affinity represented by predicted peptide IC50 value (green), MHC-I presentation prediction (orange), and EnsembleMHC (brown). The heatmap colors indicate the value of the observed F1 score (color bar) for a given algorithm (y-axis) on a particular data set (x-axis). Warmer colors indicate higher F1 scores, and cooler colors indicate lower F1 scores. The average F1 score for each algorithm across all samples is shown in the marginal bar plot. **C**, The schematic for the application of the EnsembleMHC prediction algorithm to identify SARS-CoV-2 MHC-I peptides.

97 select MHC-I alleles. These experimentally validated peptides, 98 denoted target peptides, were supplemented with a 100-fold 99 excess of decoy peptides. Decoys were generated by randomly 100 sampling peptides that were not detected by mass spectrometry, 101 but were derived from the same protein sources as a detected 102 target peptide. Algorithm and allele-specific binding affinity 103 thresholds were then identified through the independent appli- 104 cation of each component algorithm to all MHC allele data sets.

For every data set and algorithm combination, the target and 105 decoy peptides were ranked by predicted binding affinity to 106 the MHC allele defined by that data set. Then, an algorithm- 107 specific binding affinity threshold was set to the minimum score 108 needed to isolate the highest affinity peptides commensurate 109 to 50% of the observed allele repertoire size (**methods, SI** 110 **A.1**). The observed allele repertoire size was defined as the 111 total number of target peptides within a given single MHC al- 112

113 lele data set. Therefore, if a data set had 1000 target peptides, 114 the top 500 highest affinity peptides would be selected, and 115 the algorithm-specific threshold would be set to the predicted 116 binding affinity of the 500th peptide. This parameterization 117 method resulted in the generation of a customized set of allele 118 and algorithm-specific binding affinity thresholds in which an 119 expected quantity of peptides can be recovered.

120 Consensus MHC-I prediction typically require a method 121 for combining outputs from each individual component algo- 122 rithm into a composite score. This composite score is then 123 used for peptide selection. EnsembleMHC identifies high- 124 confidence peptides based on filtering by a quantity called 125 $peptide^{FDR}$ (methods Eq. 1). During the identification of 126 allele and algorithm-specific binding affinity thresholds, the 127 empirical false detection rate (FDR) of each algorithm was cal- 128 culated. This calculation was based on the proportion of target 129 to decoy peptides isolated by the algorithm specific binding 130 affinity threshold. A $peptide^{FDR}$ is then assigned to each indi- 131 vidual peptide by taking the product of the empirical FDRs of 132 each algorithm that identified that peptide for the same MHC-I 133 allele. Analysis of the parameterization process revealed that 134 the overall performance of each included algorithms was compar- 135 able, and there was diversity in individual peptide calls by 136 each algorithm, supporting an integrated approach to peptide 137 confidence assessment (SI A.2). Peptide identification by En- 138 sembleMHC was performed by selecting all peptides with a 139 $peptide^{FDR}$ of less than or equal to 5%²⁶.

140 The efficacy of $peptide^{FDR}$ as a filtering metric was de- 141 termined through the prediction of naturally presented MHC-I 142 peptides derived from ten tumor samples¹⁷ (Figure 1). Similar 143 to the single MHC allele data sets, each tumor sample data set 144 consisted of mass spectrometry-detected target peptides and a 145 100-fold excess of decoy peptides. The relative performance of 146 EnsembleMHC was assessed via comparison with individual 147 component algorithms. Peptide identification by each algo- 148 rithm was based on a restrictive or permissive binding affinity 149 thresholds (Figure 1A (inset table)). For the component algo- 150 rithms, the permissive and restrictive thresholds correspond to 151 commonly used binding affinity cutoffs for the identification of 152 weak and strong binders, respectively²⁷. The performance of 153 each algorithm on the ten data sets was evaluated through the 154 calculation of the empirical precision, recall, and F1 score.

155 The average precision and recall of each algorithm across 156 all tumor samples demonstrated an inverse relationship (Figure 157 1A). In general, restrictive binding affinity thresholds produced 158 higher precision at the cost of poorer recall. When comparing 159 the precision of each algorithm at restrictive thresholds, En- 160 sembleMHC demonstrated a 3.4-fold improvement over the 161 median precision of individual component algorithms. Ensem- 162 bleMHC also produced the highest F1 score with an average 163 of 0.51 followed by mhcfurry-presentation with an F1 score

164 of 0.45, both of which are 1.5-2 fold higher than the rest of 165 the algorithms (Figure 1B). result was shown to be robust 166 across a range of $peptide^{FDR}$ cutoff thresholds (SI A.3) and 167 alternative performance metrics (SI A.4). Furthermore, Ensem- 168 bleMHC demonstrated the ability to more efficiently prioritize 169 peptides with experimentally established immunogenicity from 170 the Hepatitis-C genome polyprotein, the Dengue virus genome 171 polyprotein, and the HIV-1 POL-GAG protein (SI A.5). Taken 172 together, these results demonstrate the enhanced precision of 173 EnsembleMHC over individual component algorithms when 174 using common binding affinity thresholds.

175 In summary, the EnsembleMHC workflow offers two 176 desirable features. First, it determines allele-specific binding 177 affinity thresholds for each algorithm at which a known quantity 178 of peptides are expected to be successfully presented on the cell 179 surface. Second, it assigns a confidence level to each peptide 180 call made by each algorithm. Together, these traits enhance the 181 ability to identify MHC-I peptides with a high probability of 182 successful cell surface presentation.

183 EnsembleMHC was used to identify MHC-I peptides 184 for the SARS-CoV-2 virus (Figure 1C). The resulting identi- 185 fication of high-confidence SARS-CoV-2 peptides allows for 186 the characterization of alleles that are enriched or depleted 187 for predicted MHC-I peptides. The resulting distribution of 188 allele-specific SARS-CoV-2 binding capacities will then be 189 weighed by the normalized frequencies of the 52 alleles (SI A.6, 190 Methods Eq. 5-6) in 23 countries to determine the population- 191 specific SARS-CoV-2 binding capacity or EnsembleMHC pop- 192 ulation score (Methods Eq. 7). The potential impact of varying 193 population SARS-CoV-2 binding capacities on disease outcome 194 can then be assessed by correlating population SARS-CoV-2 195 mortality rates with EnsembleMHC population scores. Below, 196 we use EnsembleMHC population scores to stratify countries 197 based on their mortality risks.

198 **The MHC-I peptide-allele distribution for SARS-CoV-2** 199 **structural proteins is especially disproportionate.** MHC-I 200 peptides derived from the SARS-CoV-2 proteome were pre- 201 dicted and prioritized using EnsembleMHC. A total of 67,207 202 potential 8-14mer viral peptides were evaluated for each of the 203 considered MHC-I alleles. After filtering the pool of candi- 204 date peptides at the 5% $peptide^{FDR}$ threshold, the number of 205 potential peptides was reduced from 3.49 Million to 971 (658 206 unique peptides) (SI A.7, SI table B.1). Illustrated in Figure 207 2A, the viral peptide-MHC allele (or peptide-allele) distribu- 208 tion for high-confidence SARS-CoV-2 peptides was determined 209 by assigning the identified peptides to their predicted MHC-I 210 alleles. There was a median of 16 peptides per allele with a 211 maximum of 47 peptides (HLA-A*24:02), a minimum of 3 212 peptides (HLA-A*02:05), and an interquartile range (IQR) of 213 16 peptides. Quality assurance of the predicted peptides was 214 performed by computing the peptide length frequencies and

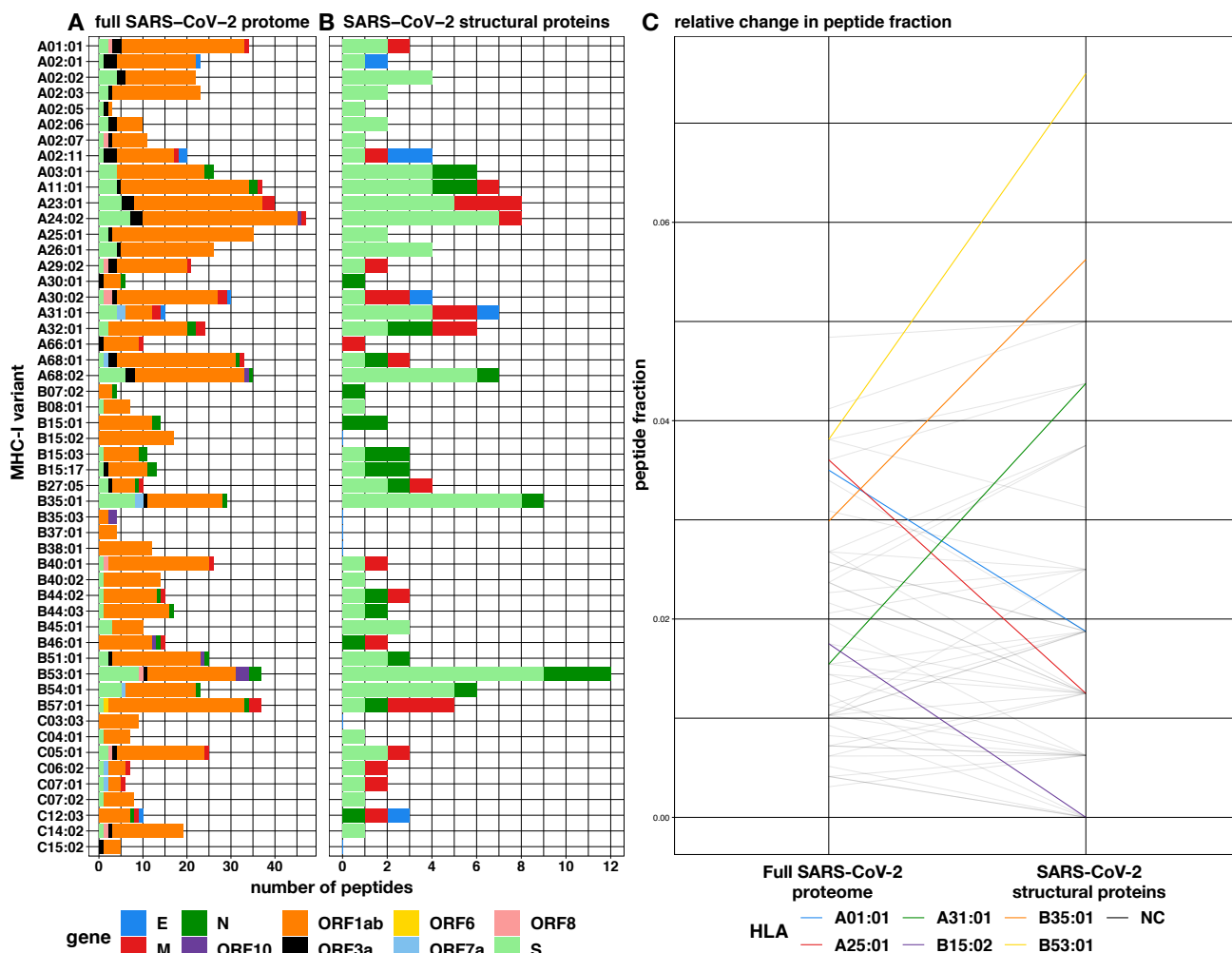


Figure 2: **Prediction of SARS-CoV-2 peptides across 52 common MHC-I alleles.** A-B, The EnsembleMHC workflow was used to predict MHC-I peptides for 52 alleles from the entire SARS-CoV-2 proteome or specifically SARS-CoV-2 structural proteins (envelope, spike, nucleocapsid, and membrane). C, The peptide fractions for both protein sets were calculated by dividing the number of peptides assigned to a given allele by the total number of identified peptides for that protein set. Each line indicates the change in peptide fraction observed by a given allele when comparing the viral peptide-MHC allele distribution for the full SARS-CoV-2 proteome or structural proteins. Alleles showing a change of greater than the median peptide fraction, $\tilde{X} = 0.015$, are highlighted in color. For the performance of EnsembleMHC at a range of different $peptide^{FDR}$ cutoff thresholds, refer to **SI A.3**

215 binding motifs. The predicted peptides were found to adhere to
 216 expected MHC-I peptide lengths²⁸ with 78% of the peptides being
 217 9 amino acids in length, 13% being 10 amino acids in length,
 218 and 8% of peptides accounting for the remaining lengths (**SI**
 219 **A.8**). Similarly, logo plots generated from predicted peptides
 220 were found to closely reflect reference peptide binding motifs
 221 for considered alleles²⁹ (**SI A.9**). Overall, the EnsembleMHC
 222 prediction platform demonstrated the ability to isolate a short

list of potential peptides which adhere to expected MHC-I peptide characteristics. 223
 224

The high expression, relative conservation, and reduced search space of SARS-CoV-2 structural proteins (S, E, M, and N) makes MHC-I binding peptides derived from these proteins high-value targets for CD8⁺ T cell-based vaccine development. **Figure 2B** describes the peptide-allele distribution for predicted MHC-I peptides originating from the four structural proteins. This analysis markedly reduces the number of con- 225
 226
 227
 228
 229
 230
 231

232 sidered peptides from 658 to 108 (**SI table B.1**). The median
233 number of predicted SARS-CoV-2 structural peptides assigned
234 to each MHC-I allele was found to be 2 with a maximum of
235 12 peptides (HLA-B*53:01), a minimum of 0 (HLA-B*15:02,
236 B*35:03,B*38:01,C*03:03,C*15:02), and a IQR of 3 peptides.
237 Analysis of the molecular source of the identified SARS-CoV-2
238 structural protein peptides revealed that they originate from en-
239 riched regions that are highly conserved (**SI A.10-A.11**). This
240 indicates that such peptides would be good candidates for tar-
241 geted therapies as they are unlikely to be disrupted by mutation,
242 and several peptides can be targeted using minimal stretches
243 of the source protein. Altogether, consideration of MHC-I
244 peptides derived only from SARS-CoV-2 structural proteins
245 reduces the number of potential peptides to a condensed set of
246 high-value targets that is amenable to experimental validation.

247 Both the peptide-allele distributions, namely the ones de-
248 rived from the full SARS-CoV-2 proteome and those from the
249 structural proteins, were found to significantly deviate from an
250 even distribution of predicted peptides as apparent in **figure**
251 **2AB** and reflected in the Kolmogorov-Smirnov test p-values
252 (**SI A.12**, full proteome = $5.673e-07$ and structural proteins
253 = $1.45e-02$). These results support a potential allele-specific
254 hierarchy for SARS-CoV-2 peptide presentation.

255 To determine if the MHC-I binding capacity hierarchy
256 was consistent between the full SARS-CoV-2 proteome and
257 SARS-CoV-2 structural proteins, the relative changes in the ob-
258 served peptide fraction (number of peptides assigned to an allele
259 / total number of peptides) between the two protein sets was vi-
260 sualized (**Figure 2C**). Six alleles demonstrated changes greater
261 than the median peptide fraction ($\bar{X} = 0.015$) when comparing
262 the two protein sets. The greatest decrease in peptide fraction
263 was observed for A*25:01 (1.52 times the median peptide frac-
264 tion), and the greatest increase was seen with B*53:01 (2.38
265 times the median peptide fraction). Furthermore, the resulting
266 SARS-CoV-2 structural protein peptide-allele distribution was
267 found to be more variable than the distribution derived from the
268 full SARS-CoV-2 proteome with a quartile coefficient of disper-
269 sion of 0.6 compared to 0.44, respectively. This indicates that
270 peptides derived from SARS-CoV-2 structural proteins experi-
271 ence larger relative inter-allele binding capacity discrepancies
272 than peptides derived from the the full SARS-CoV-2 proteome.
273 Together, these results indicate a potential MHC-I binding ca-
274 pacity hierarchy that is more pronounced for SARS-CoV-2
275 structural proteins.

276 **Total population epitope load inversely correlates with re-**
277 **ported death rates from SARS-CoV-2.** The documented im-
278 portance of MHC-I peptides derived from SARS-CoV-2 struc-
279 tural proteins⁸, coupled with the observed MHC allele bind-
280 ing capacity hierarchy and the high immunogenicity rate of
281 SARS-CoV-2 structural protein MHC-I peptides identified by
282 EnsembleMHC (95% peptides tested *in vitro*, SI A.13), prompts

283 a potential relationship between MHC-I genotype and infec-
284 tion outcome. However, due to the absence of MHC genotype
285 data for SARS-CoV-2 patients, we assessed this relationship at
286 the population-level by correlating predicted country-specific
287 SARS-CoV-2 binding capacity (or EnsembleMHC population
288 score) with observed SARS-CoV-2 mortality.

289 EnsembleMHC population scores (EMP) were deter-
290 mined for 23 countries (**SI B.2**) by weighing the individual
291 binding capacities of 52 common MHC-I alleles by their nor-
292 malized endemic expression¹⁸ (**methods, SI A.6**). This results
293 in every country being assigned two separate EMP scores, one
294 calculated with respect to the 108 unique SARS-CoV-2 struc-
295 tural protein peptides (structural protein EMP) and the other
296 with respect to the 658 unique peptides derived from the full
297 SARS-CoV-2 proteome (full proteome EMP). The EMP score
298 corresponds to the average predicted SARS-CoV-2 binding ca-
299 pacity of a population. Therefore, individuals in a country with
300 a high EMP score would be expected, on average, to present
301 more SARS-CoV-2 peptides to CD8⁺ T cells than individuals
302 from a country with a low EMP score. The resulting EMP
303 scores were then correlated with observed SARS-CoV-2 mor-
304 tality (deaths per million) as a function of time. Temporal
305 variance in community spread within the cohort of countries
306 was corrected by truncating the SARS-CoV-2 mortality data
307 set for each country to start after a certain minimum death
308 threshold was met. For example, if the minimum death thresh-
309 old was 50, then day 0 would be when each country reported
310 at least 50 deaths. The number of countries included in each
311 correlation decreases as the number of days increases due to
312 discrepancies in the length of time that each country met a given
313 minimum death threshold (**SI table B.3**). Therefore, the cor-
314 relation between EMP score and SARS-CoV-2 mortality was
315 only estimated at time points where there were at least eight
316 countries. The eight country threshold was chosen because it is
317 the minimum sample size needed to maintain sufficient power
318 when detecting large effect sizes ($\rho > 0.85$). The strength of the
319 relationship between EMP score and SARS-CoV-2 mortality
320 was determined using Spearman's rank-order correlation (for
321 details concerning the choice of statistical tests, please refer
322 to the methods section). Accordingly, both EMP scores and
323 SARS-CoV-2 mortality data were converted into ascending
324 ranks with the lowest rank indicating the minimum value and
325 the highest rank indicating the maximum value. For instance, a
326 country with an EMP score rank of 1 and death per million rank
327 of 23 would have the lowest predicted SARS-CoV-2 binding
328 capacity and the highest level of SARS-CoV-2-related mortal-
329 ity. Using the described paradigm, the structural protein EMP
330 score and the full proteome EMP score were correlated with
331 SARS-CoV-2-related deaths per million for 23 countries.

332 Total predicted population SARS-CoV-2 binding capacity
333 exhibited a strong inverse correlation with observed deaths per

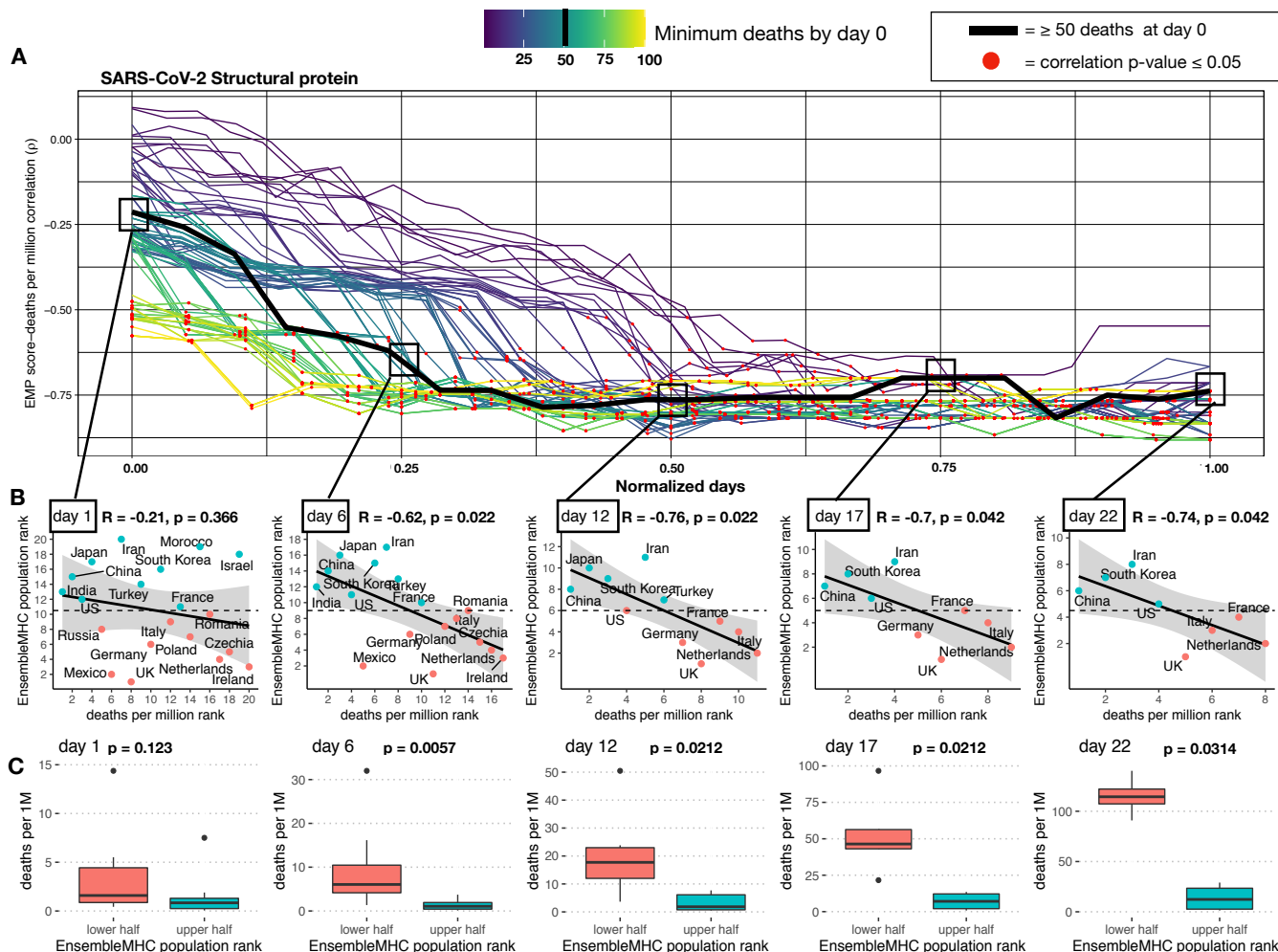


Figure 3: Predicted total epitope load within a population inversely correlates with mortality. **A**, SARS-CoV-2 structural protein-based EnsembleMHC population scores were assigned to 23 countries (SI B.2), and correlated with observed mortality rate (deaths per million). The correlation coefficient is presented as a function of time. Individual country mortality rate data were aligned by truncating each data set to start after a minimum threshold of deaths was observed in a given country (line color). The Spearman's rank correlation coefficient between structural protein EMP score and SARS-CoV-2 mortality rate was calculated at every day following day 0 for each of the minimum death thresholds. Due to the differing lengths of time series analysis at each minimum death threshold, the number of days were normalized to improve visualization. Thus, normalized day 0 represents the day when qualifying countries recorded at least the number of deaths indicated by the minimum death threshold, and normalized day 1 represents the final time point at which a correlation was measured. (For mapping between real and normalized days, see SI B.2). Correlations that were shown to be statistically significant (p -value ≤ 0.05) are indicated by a red point. **B**, The correlations between the structural protein EnsembleMHC population score (y-axis) and deaths per million (x-axis) were shown for countries meeting the 50 minimum deaths threshold at days 1, 6, 12, 17, and 22. Correlation coefficients and p-values were assigned using Spearman's rank correlation and the shaded region signifies the 95% confidence interval. Due to Spearman's rank correlation only considering data rank, Deaths per million and EnsembleMHC population score were converted to ascending rank values (low rank = low values, high rank = high values) to improve visualization of the measured relationship. Red points indicate a country that has an EnsembleMHC population rank less than the median EnsembleMHC population rank of all countries at that day, and blue points indicate a country with an EnsembleMHC population rank greater than the median EnsembleMHC population rank. **C**, The countries at each day were partitioned into an upper or lower half based on the median observed EnsembleMHC population rank. Therefore, countries with an EnsembleMHC population rank greater than the median group EnsembleMHC population score were assigned to the upper half (red), and the remaining countries were assigned to the lower half (blue). p-values were determined by Mann-Whitney U test. The presented box plots are in the style of Tukey (box defined by 25%, 50%, 75% quantiles, and whiskers $\pm 1.5 \times$ IQR). The increasing gap between the red and the blue box plots indicates a greater discrepancy in the number of deaths per million between the two groups. The p values in all figures were corrected using the Benjamini-Hochberg procedure³⁰ relative to the number of tests performed for each death threshold

334 million. This relationship was found to be true for correlations
335 based on the structural protein EMP (**Figure 3A**) and full pro-
336 teome EMP (**SI A.14**) scores with a mean effect size of -0.66
337 and -0.60, respectively. Significance testing of the correlations
338 produced by both EMP scores revealed that the majority of
339 reported correlations are statistically significant with 63% at-
340 taining a p-value of ≤ 0.05 . Correlations based on the structural
341 protein EMP score demonstrated a 24% higher proportion of
342 statistically significant correlations compared to the full pro-
343 teome EMP score (74% vs 51%). Furthermore, correlations
344 for EMP scores based on structural proteins produced narrower
345 95% confidence intervals (**SI A.15-A.16**, **SI table B.3**). Due
346 to relatively low statistical power of the obtained correlations
347 (**SI A.17**), the positive predictive value for each correlation
348 (**methods, Eq. 8**) was calculated. The resulting proportions
349 of correlations with a positive predictive value of $\geq 95\%$ were
350 similar to the observed significant p-value proportions with
351 62% of all measured correlations, 72% of structural protein
352 EMP score correlations, and 52% full proteome EMP score
353 correlations (**SI A.14**). The similar proportions of significant
354 p-values and PPVs supports that an overall true association is
355 being captured. Furthermore, analysis of similar sized peptide
356 sets sampled from the full SARS-CoV-2 proteome revealed that
357 the observed distinction between the correlations produced by
358 the two protein groups are unlikely to be due to differences in
359 peptide set sizes (**SI A.18**).

360 Finally, the reported correlations did not remain after ran-
361 domizing the allele assignment of predicted peptides prior to
362 *peptide*^{FD_R} filtering (**SI A.19**), through the use of any indi-
363 vidual algorithm (**SI A.20**). This indicates that the observed
364 relationship is contingent on the high-confidence peptide-allele
365 distribution produced by the EnsembleMHC prediction algo-
366 rithm. Altogether, these data demonstrate that the MHC-I allele
367 hierarchy characterized by EnsembleMHC is inversely asso-
368 ciated with SARS-CoV-2 population mortality, and that the
369 relationship becomes stronger when considering only the pre-
370 sentation of SARS-CoV-2 structural proteins.

371 The ability to use structural protein EMP score to identify
372 high and low risk populations was assessed using the median
373 minimum death threshold (50 deaths) at evenly spaced time
374 points (**Figure 3A, squares**). All correlations, with the excep-
375 tion of day 1, were found to be significant with an average
376 effect size of -0.71 (**Figure 3B**). Next, the countries at each
377 day were partitioned into a high or low group based based on
378 whether their assigned EMP score was higher or lower than
379 the median observed EMP score (**Figure 3C**). The resulting
380 grouping demonstrated a statistically significant difference in
381 the median deaths per million between countries with low struc-
382 tural protein EMP score and countries with high structural
383 protein EMP scores. Additionally, it was observed that deaths
384 per million increased much more rapidly in countries with low

385 structural protein EMP scores. Taken together, these results
386 indicate that structural protein EMP score may be useful for
387 assessing population risk from SARS-CoV-2 infections.

388 In summary, we make several important observations.
389 First, there is a strong inverse correlation between predicted
390 population SARS-CoV-2 binding capacity and observed deaths
391 per million. This finding suggests that outcome to SARS-CoV-
392 2 may be tied to total epitope load. Second, the correlation
393 between predicted epitope load and population mortality is
394 stronger for SARS-CoV-2 structural MHC-I peptides. This
395 suggests that CD8⁺ T cell-mediated immune response maybe
396 primarily driven by recognition of epitopes derived from these
397 proteins, a finding supported by recent T cell epitope mapping
398 of SARS-CoV-2⁸. Finally, the EnsembleMHC population score
399 can separate countries within the considered cohort into high
400 or low risk populations.

401 **Structural protein EMP score correlates better with popu-** 402 **lation outcome than identified individual risk factors.**

403 Recent large scale patient studies have identified several
404 socioeconomic and health-related factors associated with in-
405 creased risk of death from SARS-CoV-2 infections^{31,32}. To
406 delineate the relative importance of the structural protein EMP
407 score as a SARS-CoV-2 severity descriptor, 12 additional risk
408 factors were assessed for their ability to model population level
409 SARS-CoV-2 outcome in 21 countries (**SI B.4**).

410 Overall, the structural protein EMP scores produced a
411 significantly stronger association with population SARS-CoV-2
412 mortality compared to other 12 descriptors (**Figure 4A**). While
413 various effect size trends were observed, all additional covari-
414 ates failed to produce statistically significant correlations. To
415 determine if the modeling of SARS-CoV-2 mortality rate could
416 be improved by the combination of single socioeconomic or
417 health-related risk factors with structural protein EMP scores,
418 a set of linear models consisting of either a single risk factor
419 (single feature model) or that factor combined with structural
420 protein EMP scores (combination model) were generated for
421 every time point across each minimum death threshold (**meth-**
422 **ods**). Following model generation, the adjusted coefficient of
423 determination (R^2) and significance level of each individual
424 model was extracted and aggregated by dependent variable
425 (**SI A.21**). Single feature models were characterized by low
426 R^2 ($\tilde{x} = -0.0262$) while combination models showed signifi-
427 cant improvement ($\tilde{x} = .496$). Similarly, combination models
428 demonstrated a substantially higher proportion of statistical
429 significance (**SI A.21B**). To determine the set of features that
430 produce the best fitting model, all possible combinations of
431 explanatory factors (risk factors and structural protein EMP
432 score) were tested. Subsequently, the top ten performing mod-
433 els, ranked by adjusted R^2 value, were selected for analysis
434 (**Figure 4B**). The identified models were found to be largely sig-

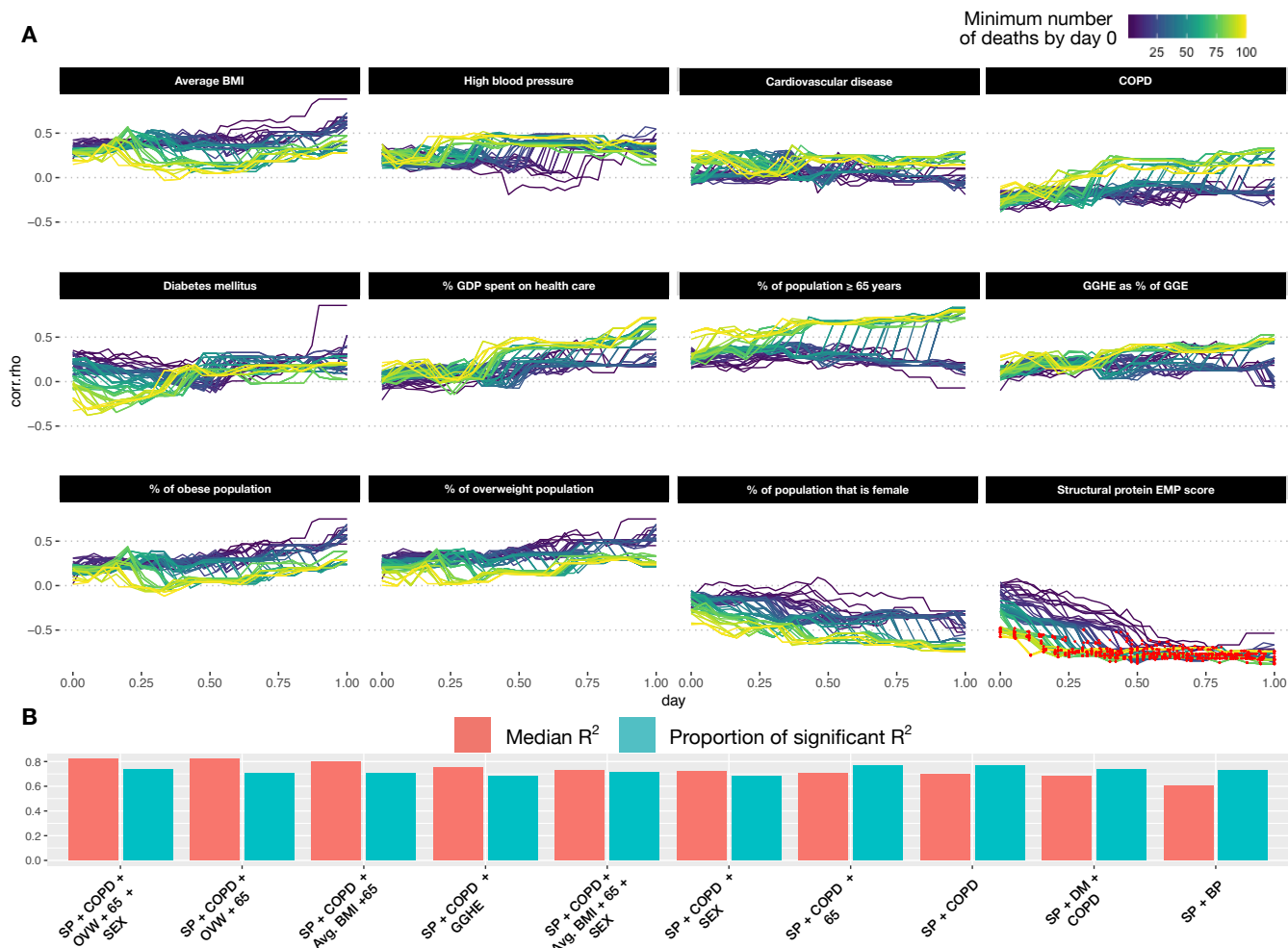


Figure 4: Analysis of other SARS-CoV-2 covariates with observed SARS-CoV-2 population mortality and development of an integrative model. **A**, 12 covariates associated with SARS-CoV-2 mortality on the individual patient level were assessed for correlation with population level mortality (SI table B.4). The correlation of each country-level covariate was determined at each time point after a minimum death threshold was met (line color). The x-axis represents the number of days (normalized) following when a minimum death threshold was met, and the y-axis indicates the observed effect size for that covariate at a given time point. Correlations achieving statistical significance are colored with a red dot. **B**, All possible combinations of covariates were used to fit a linear model. The top 10 models, ranked by median adjusted R^2 (red bars), were identified (**B**). The proportion of regressions performed by that model that were found to be statistically significant (F-test 0.05) are represented by the blue bars.

435 nificant (average proportion of significant regressions = 72%)
 436 and produce strong fits to the data (average $R^2 = 0.7$).

437 Analysis of the dependent variables included in the top
 438 performing models revealed that all models included structural
 439 protein EMP scores followed by deaths per million due to
 440 complications from COPD (90% of models). The median model
 441 size included 3 features with a maximum of 5 features and a
 442 minimum of 2 features. The model producing the best fit (me-

443 dian $R^2 = 0.791$) consisted of structural protein EMP scores, 443
 444 gender demographics, number of deaths due to COPD compli- 444
 445 cations, the proportion of the population over the age of 65, 445
 446 and proportion of the population that is overweight (**Figure 4B**). 446
 447 All together, these results further indicate the robustness of the 447
 448 structural protein EMP score as a population level risk descrip- 448
 449 tor and identifies a potential candidate model for predicting 449
 450 pandemic severity. 450

451 Discussion

452 In the present study, we uncover evidence supporting an associ- 501
453 ation between population SARS-CoV-2 infection outcome and 502
454 MHC-I genotype. In line with related work highlighting the 503
455 relationship between total epitope load with HIV viral control³³, 504
456 we arrive at a working model that MHC-I alleles presenting 505
457 more unique SARS-CoV-2 epitopes will be associated with 506
458 lower mortality due to a higher number of potential T cell 507
459 targets. The SARS-CoV-2 binding capacities of 52 common 508
460 MHC-I alleles were assessed using the EnsembleMHC predic- 509
461 tion platform. These predictions identified 971 high-confidence 510
462 MHC-I peptides out of a candidate pool of nearly 3.5 million. In 511
463 agreement with other *in silico* studies^{15,34}, the assignment of the 512
464 predicted peptides to their respective MHC-I alleles revealed an 513
465 uneven distribution in the number of peptides attributed to each 514
466 allele. We discovered that the MHC-I peptide-allele distribution 515
467 originating from the full SARS-CoV-2 proteome undergoes a 516
468 notable rearrangement when considering only peptides derived 517
469 from viral structural proteins. The structural protein-specific 518
470 peptide-allele distribution produced a distinct hierarchy of al- 519
471 lele binding capacities. This finding has important clinical 520
472 implications as a majority of SARS-CoV-2 specific CD8⁺ T 521
473 cell response is directed towards SARS-CoV-2 structural pro- 522
474 teins⁸. Therefore, patients who express MHC-I alleles enriched 523
475 with a large potential repertoire of SARS-CoV-2 structural pro- 524
476 teins peptides may benefit from a broader CD8⁺ T cell immune 525
477 response. 526

478 The variations in SARS-CoV-2 peptide-allele distribu- 527
479 tions were analyzed at epidemiological scale to track its impact 528
480 on country-specific mortality. Each of the 23 countries were 529
481 assigned a population SARS-CoV-2 binding capacity (or En- 530
482 sembleMHC population score) based on the individual binding 531
483 capacities of the selected 52 MHC-I alleles weighted by their 532
484 endemic population frequencies. This hierarchization revealed 533
485 a strong inverse correlation between EnsembleMHC population 534
486 score and observed population mortality, indicating that popula- 535
487 tions enriched with high SARS-CoV-2 binding capacity MHC-I 536
488 alleles may be better protected. The correlation was shown to 537
489 be stronger when calculating the EnsembleMHC population 538
490 scores with respect to only structural proteins, reinforcing their 539
491 relevance to viral immunity. Finally, The molecular origin of 540
492 the 108 predicted peptides specific to SARS-CoV-2 structural 541
493 proteins revealed that they are derived from enriched regions 542
494 with a minimal predicted impact from amino acid sequence 543
495 polymorphisms. 544

496 The utility of structural protein EnsembleMHC popula- 545
497 tion scores was further supported by a multivariate analysis of 546
498 additional SARS-CoV-2 risk factors. These results emphasized 547
499 the relative robustness of structural protein EMP scores as a 548
500 population risk assessment tool. Furthermore, a linear model 549

501 based on the combination of structural protein EMP scores and 502
503 select population-level risk factors was identified a potential 504
505 candidate for a predictive model for pandemic severity. As 506
507 such, the incorporation of the structural protein EMP score in 508
509 more sophisticated models will likely improve epidemiological 510
511 modeling of pandemic severity. 512

513 In order to achieve the highest level of accuracy in MHC-I 514
515 predictions, the most up-to-date versions of each component 516
517 algorithm were used. However, this meant that several of the al- 518
519 gorithms (MHCflurry, netMHCpan-EL-4.0 and MixMHCpred) 520
521 were benchmarked against subsets of mass spectrometry data 522
523 that were used in the original training of these MHC-I predic- 524
525 tion models. While this could result in an unfair weight applied 526
527 to these algorithms in *peptide*^{FDR} calculation, the individual 528
529 FDRs of MHCflurry, netMHCpan-EL-4.0 and MixMHCpred 529
530 were comparable to algorithms without this advantage (**SI A.2**). 531
532 Furthermore, the peptide selection of SARS-CoV-2 peptides 532
533 was shown to be highly cooperative within EnsembleMHC (**SI** 533
534 **A.7**), and individual algorithms failed to replicate the strong 534
535 observed correlations between population binding capacity and 535
536 observed SARS-CoV-2 mortality (**SI A.20**). 536
537

538 In the future, the presented model could be applied to pre- 539
540 dict individual T cell capacity to mount a robust SARS-CoV-2 540
541 immune response. Evolutionary divergence of patient MHC-I 541
542 genotypes have shown to be predictive of response to immune 542
543 checkpoint therapy in cancer and HIV^{35,36}. However, confirma- 543
544 tion will require large data sets associating individual patient 544
545 MHC-I genotype and outcome. Additionally, future use of En- 545
546 sembleMHC to design personalized T cell vaccines will require 546
547 broad experimental validation of high scoring peptides, since 547
548 EnsembleMHC predicts MHC-I peptides with a high proba- 548
549 bility of antigen presentation as opposed to directly predicting 549
550 peptide immunogenicity. While previous work has determined 550
551 that a majority of successfully presented viral MHC-I peptides 551
552 are immunogenic³⁷, there is an expectation that some presented 552
553 SARS-CoV-2 MHC-I peptides will fail to produce an immune 553
554 response. 554

555 The current work assessed the relative importance of the 556
557 structural protein EMP score with respect to other population- 557
558 level risk factors (e.g. population incidence of risk-associated 558
559 commodities, healthcare infrastructure, age, sex), however, it 559
560 should be noted that the impacts these risk factors on patient 560
561 outcome are likely to vary significantly on a individual basis. 561
562 Furthermore, other genetic determinants of severity were not 562
563 considered³⁸. Therefore, a complete understanding of the rela- 563
564 tive importance of MHC genotype and SARS-CoV-2 presenta- 564
565 tion capacity on patient outcome will require the integration 565
566 individual patient genetic and clinical data. 566

567 The versatility of the proposed model will be improved 568
569 by the consideration of additional MHC-I alleles. To reduce 569
570

551 the presence of confounding factors, EnsembleMHC was pa- 595
552 rameterized on only a subset of common MHC-I alleles that 596
553 had strong existing experimental validation. While the selected 597
554 MHC-I alleles are among some of the most common, person- 598
555 alized risk assessment will require consideration of the full 599
556 patient MHC-I genotype. The continued mass spectrometry- 600
557 based characterization of MHC-I peptide binding motifs will 601
558 help in this regard. However, due to the large potential 602
559 sequence space of the MHC-I protein, extension of this model 603
560 will likely require inference of binding motifs based on MHC 604
561 variant clustering. 605

562 Acknowledgments

563 We would like to thank Drs. Diego Chowell, Matthew Scotch, 606
564 Sri Krishna, Shay Ferdosi, and Mr. John Vant, Mr. Ryan Boyd, 607
565 and Ms. Mollie Peters for critical feedback and discussion. 608
566 Finally, we would like to thank ASU Research computing for 609
567 allocating the computational resources. 610

568 References

- 569 [1] Zi Yue Zu et al. “Coronavirus disease 2019 (COVID- 612
570 19): a perspective from China”. In: *Radiology* (2020), 613
571 p. 200490. 614
- 572 [2] Qun Li et al. “Early transmission dynamics in Wuhan, 615
573 China, of novel coronavirus–infected pneumonia”. In: 616
574 *New England Journal of Medicine* (2020). 617
- 575 [3] Yan-Rong Guo et al. “The origin, transmission and clinical 618
576 therapies on coronavirus disease 2019 (COVID-19) 619
577 outbreak—an update on the status”. In: *Military Medical 620
578 Research* 7.1 (2020), pp. 1–10. 621
- 579 [4] Rudragouda Channappanavar, Jincun Zhao, and Stanley 622
580 Perlman. “T cell-mediated immune response to respira- 623
581 tory coronaviruses”. In: *Immunologic research* 59.1-3 624
582 (2014), pp. 118–128. 625
- 583 [5] Hsueh-Ling Janice Oh et al. “Understanding the T cell 626
584 immune response in SARS coronavirus infection”. In: 627
585 *Emerging microbes & infections* 1.1 (2012), pp. 1–6. 628
- 586 [6] Oi-Wing Ng et al. “Memory T cell responses target- 629
587 ing the SARS coronavirus persist up to 11 years post- 630
588 infection”. In: *Vaccine* 34.17 (2016), pp. 2008–2014. 631
- 589 [7] Nina Le Bert et al. “SARS-CoV-2-specific T cell immu- 632
590 nity in cases of COVID-19 and SARS, and uninfected 633
591 controls”. In: *Nature* 584.7821 (2020), pp. 457–462. 634
- 592 [8] Alba Grifoni et al. “Targets of T cell responses to SARS- 635
593 CoV-2 coronavirus in humans with COVID-19 disease 636
594 and unexposed individuals”. In: *Cell* (2020). 637
- [9] Vasiliki Matzaraki et al. “The MHC locus and genetic 638
susceptibility to autoimmune and infectious diseases”. 639
In: *Genome biology* 18.1 (2017), p. 76. 640
- [10] Marie Lin et al. “Association of HLA class I with severe 641
acute respiratory syndrome coronavirus infection”. In: 642
BMC Medical Genetics 4.1 (2003), p. 9. 643
- [11] Sheng-Fan Wang et al. “Human-leukocyte antigen class 644
I Cw 1502 and class II DR 0301 genotypes are associ-
ated with resistance to severe acute respiratory syndrome
(SARS) infection”. In: *Viral immunology* 24.5 (2011),
pp. 421–426.
- [12] Margaret HL Ng et al. “Association of human-leukocyte-
antigen class I (B* 0703) and class II (DRB1* 0301)
genotypes with susceptibility and resistance to the de-
velopment of severe acute respiratory syndrome”. In:
Journal of Infectious Diseases 190.3 (2004), pp. 515–
518.
- [13] MH Ng et al. “Immunogenetics in SARS: a case-control
study.” In: *Hong Kong medical journal= Xianggang yi
xue za zhi* 16.5 Suppl 4 (2010), p. 29.
- [14] Alicia Sanchez-Mazas. “HLA studies in the context
of coronavirus outbreaks”. In: *Swiss Medical Weekly*
150.1516 (2020).
- [15] Austin Nguyen et al. “Human leukocyte antigen suscep-
tibility map for SARS-CoV-2”. In: *Journal of Virology*
(2020).
- [16] Weilong Zhao and Xinwei Sher. “Systematically bench-
marking peptide-MHC binding predictors: From syn-
thetic to naturally processed epitopes”. In: *PLoS compu-
tational biology* 14.11 (2018).
- [17] Siranush Sarkizova et al. “A large peptidome dataset
improves HLA class I epitope prediction across most of
the human population”. In: *Nature Biotechnology* 38.2
(2020), pp. 199–209.
- [18] Faviel F González-Galarza et al. “Allele frequency net
2015 update: new features for HLA epitopes, KIR and
disease and HLA adverse drug reaction associations”. In:
Nucleic acids research 43.D1 (2015), pp. D784–D788.
- [19] Timothy J O’Donnell, Alex Rubinsteyn, and Uri Laser-
son. “MHCflurry 2.0: Improved Pan-Allele Prediction
of MHC Class I-Presented Peptides by Incorporating
Antigen Processing”. In: *Cell Systems* (2020).
- [20] Vanessa Jurtz et al. “NetMHCpan-4.0: improved peptide-
MHC class I interaction predictions integrating eluted
ligand and peptide binding affinity data”. In: *The Journal
of Immunology* 199.9 (2017), pp. 3360–3368.
- [21] Massimo Andreatta and Morten Nielsen. “Gapped se-
quence alignment using artificial neural networks: appli-
cation to the MHC class I system”. In: *Bioinformatics*
32.4 (2016), pp. 511–517.

- 645 [22] Michal Bassani-Sternberg et al. “Deciphering HLA-I 696
646 motifs across HLA peptidomes improves neo-antigen 697
647 predictions and identifies allosteric regulating HLA speci- 698
648 ficity”. In: *PLoS computational biology* 13.8 (2017), 699
649 e1005725. 700
- 650 [23] Hao Zhang, Ole Lund, and Morten Nielsen. “The Pick- 701
651 Pocket method for predicting binding specificities for 702
652 receptors based on receptor pocket similarities: applica- 703
653 tion to MHC-peptide binding”. In: *Bioinformatics* 25.10 704
654 (2009), pp. 1293–1299. 705
- 655 [24] Michael Rasmussen et al. “Pan-specific prediction of 706
656 peptide–MHC class I complex stability, a correlate of T 707
657 cell immunogenicity”. In: *The Journal of Immunology* 708
658 197.4 (2016), pp. 1517–1524. 709
- 659 [25] Sinu Paul et al. “HLA class I alleles are associated 710
660 with peptide-binding repertoires of different size, affinity, 711
661 and immunogenicity”. In: *The Journal of Immunology* 712
662 191.12 (2013), pp. 5831–5839. 713
- 663 [26] K Nichols. “False discovery rate procedures”. In: *Sta- 714
664 tistical Parametric Mapping*. Elsevier, 2007, pp. 246– 715
665 252. 716
- 666 [27] Morten Nielsen et al. “Immunoinformatics: Predicting 717
667 Peptide–MHC Binding”. In: *Annual Review of Biomed- 718
668 ical Data Science* 3 (2020). 719
- 669 [28] Thomas Trolle et al. “The length distribution of class I– 720
670 restricted T cell epitopes is determined by both peptide 721
671 supply and MHC allele–specific binding preference”. 722
672 In: *The Journal of Immunology* 196.4 (2016), pp. 1480– 723
673 1487. 724
- 674 [29] Nicolas Rapin et al. “The MHC motif viewer: a visualiza- 725
675 tion tool for MHC binding motifs”. In: *Current protocols 726
676 in immunology* 88.1 (2010), pp. 18–17. 727
- 677 [30] Yoav Benjamini and Yocef Hochberg. “Controlling the 728
678 false discovery rate: a practical and powerful approach 729
679 to multiple testing”. In: *Journal of the Royal statistical 730
680 society: series B (Methodological)* 57.1 (1995), pp. 289– 731
681 300. 732
- 682 [31] Elizabeth J Williamson et al. “OpenSAFELY: factors 733
683 associated with COVID-19 death in 17 million patients”. 734
684 In: *Nature* (2020), pp. 1–11. 735
- 685 [32] Simon de Lusignan et al. “Risk factors for SARS-CoV-2 736
686 among patients in the Oxford Royal College of General 737
687 Practitioners Research and Surveillance Centre primary 738
688 care network: a cross-sectional study”. In: *The Lancet 739
689 Infectious Diseases* (2020). 740
- 690 [33] Morgane Rolland et al. “Broad and Gag-biased HIV-1 741
691 epitope repertoires are associated with lower viral loads”. 742
692 In: *PloS one* 3.1 (2008). 743
- 693 [34] Katie M Campbell et al. “Prediction of SARS-CoV-2 744
694 epitopes across 9360 HLA class I alleles”. In: *bioRxiv* 745
695 (2020). 746
- 700 [35] Diego Chowell et al. “Evolutionary divergence of HLA 696
701 class I genotype impacts efficacy of cancer immunother- 697
702 apy”. In: *Nature medicine* 25.11 (2019), pp. 1715–1720. 698
- 703 [36] Jatin Arora et al. “HLA heterozygote advantage against 699
704 HIV-1 is driven by quantitative and qualitative differ- 700
705 ences in HLA allele-specific peptide presentation”. In: 701
706 *Molecular biology and evolution* 37.3 (2020), pp. 639– 702
707 650. 703
- 708 [37] Nathan P Croft et al. “Most viral peptides displayed by 704
709 class I MHC on infected cells are immunogenic”. In: 705
710 *Proceedings of the National Academy of Sciences* 116.8 706
711 (2019), pp. 3112–3117. 707
- 712 [38] Yanan Cao et al. “Comparative genetic analysis of the 708
713 novel coronavirus (2019-nCoV/SARS-CoV-2) receptor 709
714 ACE2 in different populations”. In: *Cell discovery* 6.1 710
715 (2020), pp. 1–4. 711
- 716 [39] Marek Prachar et al. “COVID-19 Vaccine Candidates: 712
717 Prediction and Validation of 174 SARS-CoV-2 Epi- 713
718 topes”. In: *bioRxiv* (2020). 714
- 719 [40] Syed Faraz Ahmed, Ahmed A Quadeer, and Matthew R 715
720 McKay. “COVIDep: A web-based platform for real-time 716
721 reporting of vaccine target recommendations for SARS- 717
722 CoV-2”. In: *Nature reviews microbiology* 15 (2020), 718
723 pp. 2141–2142. 719
- 724 [41] Helen M Berman et al. “The protein data bank”. In: *Nu- 720
725 cleic acids research* 28.1 (2000), pp. 235–242. 721
- 726 [42] Chengxin Zhang et al. “Protein structure and sequence 722
727 re-analysis of 2019-nCoV genome refutes snakes as its 723
728 intermediate host or the unique similarity between its 724
729 spike protein insertions and HIV-1”. In: *Journal of pro- 725
730 teome research* (2020). 726
- 731 [43] William Humphrey, Andrew Dalke, Klaus Schulten, et 727
732 al. “VMD: visual molecular dynamics”. In: *Journal of 728
733 molecular graphics* 14.1 (1996), pp. 33–38. 729
- 734 [44] Ensheng Dong, Hongru Du, and Lauren Gardner. “An 730
735 interactive web-based dashboard to track COVID-19 in 731
736 real time”. In: *The Lancet infectious diseases* (2020). 732
- 737 [45] Katherine S Button et al. “Power failure: why small sam- 733
738 ple size undermines the reliability of neuroscience”. In: 734
739 *Nature Reviews Neuroscience* 14.5 (2013), pp. 365–376. 735
- 740 [46] Andrew P Ferretti et al. “COVID-19 patients form mem- 736
741 ory CD8+ T cells that recognize a small set of shared 737
742 immunodominant epitopes in SARS-CoV-2”. In: (2020). 738
- 743 [47] Annika Nelde et al. “SARS-CoV-2-derived peptides de- 739
744 fine heterologous and COVID-19-induced T cell recog- 740
745 nition”. In: *Nature immunology* (2020), pp. 1–12. 741
- 746 [48] Thomas M Snyder et al. “Magnitude and dynamics of 742
747 the T-cell response to SARS-CoV-2 infection at both 743
748 individual and population levels”. In: *medRxiv* (2020). 744

- 745 [49] Ahmed A Quadeer, Syed Faraz Ahmed, and Matthew R
746 McKay. “Epitopes targeted by T cells in convalescent
747 COVID-19 patients”. In: *bioRxiv* (2020).

748 Methods

749 EnsembleMHC prediction workflow

750 **EnsembleMHC component binding and processing predic-**
751 **tion algorithms.** EnsembleMHC incorporates MHC-I binding
752 and processing predictions from 7 publicly available algorithms:
753 MHCflurry-affinity-1.6.0¹⁹, MHCflurry-presentation-1.6.0¹⁹,
754 netMHC-4.0²¹, netMHCpan-4.0-EL²⁰, netMHCstabpan-1.0²⁴,
755 PickPocket-1.1²³ and, MixMHCpred-2.0.2²². These algorithms
756 were chosen based on the criteria of providing a free academic
757 license, bash command line integration, and demonstrated ac-
758 curacy for predicting SARS-CoV-2 MHC-I peptides with ex-
759 perimentally validated binding stability³⁹.

760 Each of the selected algorithms cover components of
761 MHC-I binding and antigen processing that roughly fall into
762 two categories: ones based primarily on MHC-I binding affini-
763 ty predictions and others that model antigen presentation.
764 To this end, MHCflurry-affinity, netMHC, PickPocket, and
765 netMHCstabpan predict binding affinity based on quantitative
766 peptide binding affinity measurements. netMHCstabpan also
767 incorporates peptide-MHC stability measurements and Pick-
768 Pocket performs prediction based on binding pocket structural
769 extrapolation. To model the effects of antigen presentation,
770 MixMHCpred, netMHCpan-EL, and MHCflurry-presentation
771 are trained on naturally eluted MHC-I ligands. Additionally,
772 MHCflurry-presentation incorporates an antigen processing
773 term.

774 **Parameterization of EnsembleMHC using mass spectrom-**
775 **etry data.** EnsembleMHC is able to achieve high levels of
776 precision in peptide selection through the use of allele and
777 algorithm-specific binding affinity thresholds. These binding
778 affinity thresholds were identified through the parameterization
779 of each algorithm on high-quality mass spectrometry data sets¹⁷.
780 The mass spectrometry data sets used for algorithm parameteri-
781 zation were collected in the largest single laboratory MS-based
782 characterization of MHC-I peptides presented by single MHC
783 allele cell lines. These characteristics significantly reduces the
784 number of artifacts introduced by differences in peptide isola-
785 tion methods, mass spectrometry acquisition, and convolution
786 of peptides in multiallelic cell lines. An overview of the Ensem-
787 bleMHC parameterization is provided in supplemental figures
788 (SI A.1).

789 Fifty-two common MHC-I alleles were selected for pa-
790 rameterization based on the criteria that they were characterized
791 in *Sarkizova et al.* data sets and that all 7 component algorithms
792 could perform peptide binding affinity predictions for that allele.
793 Each target peptide (observed in the MS data set) was paired
794 with 100 length-matched randomly sampled decoy peptides
795 (not observed in the MS data set) derived from the same source

796 proteins. If a protein was less than 100 amino acids in length,
797 then every potential peptide from that protein was extracted. 797

798 Each of the seven algorithms were independently applied
799 to each of the 52 allele data sets. For each allele data set, the
800 minimum score threshold was determined for each algorithm
801 that recovered 50% of the allele repertoire size (the total num-
802 ber of target peptides observed in the MS data set for that allele).
803 Additionally, the expected accuracy of each algorithm was as-
804 sessed by calculating the observed false detection rate (the
805 fraction of identified peptides that were decoy peptides) using
806 the identified algorithm and allele specific scoring threshold.
807 The parameterization process was repeated 1000 times for each
808 allele through bootstrap sampling of half of the peptides in each
809 single MHC allele data set. The final FDR and score threshold
810 for each algorithm at each allele was determined by taking
811 the median value of both quantities reported during bootstrap
812 sampling.

813 **Peptide confidence assessment.** Peptide confidence is as-
814 signed by calculating the $peptide^{FDR}$. This quantity is de-
815 fined as the product of the empirical FDRs of each individual
816 algorithm that detected a given peptide. The $peptide^{FDR}$ is
817 calculated using equation 1,

$$peptide^{FDR} = \prod_{i=1, i \neq ND}^N algorithm_i^{FDR} \quad (1)$$

818 , where N is the number of MHC-I binding and processing algo-
819 rithms, ND represents an algorithm that did not detect a given
820 peptide, and $algorithm^{FDR}$ represents the allele specific FDR
821 of the N th algorithm.

822 The $peptide^{FDR}$ represents the joint probability that all
823 MHC-I binding and processing algorithms that detected a partic-
824 ular peptide did so in error, and therefore returns a probability
825 of false detection. Unless otherwise stated, EnsembleMHC
826 selected peptides based on the criterion of a $peptide^{FDR} \leq$
827 5%.

828 Application of EnsembleMHC to tumor cell line data

829 **Tumor MHC-I peptide data sets.** Ten tumor samples were
830 obtained from the *Sarkizova et al.* data sets. Tumor samples
831 were selected for analysis if at least 50% of the expressed MHC-
832 I alleles for that sample were included in the 52 MHC-I alleles
833 supported by EnsembleMHC. For each data set, decoy peptides
834 were generated in a manner identical to the method used for
835 algorithm parameterization on single MHC allele data.

836 **Tumor MHC-I peptide identification.** Peptide identification
837 by each algorithm was based on restrictive or permissive bind-
838 ing affinities thresholds. These thresholds correspond to com-
839 monly used score cutoffs for the identification of strong binders
840 (restrictive) or all binders (permissive). These thresholds are

841 0.5% (percentile rank) or 50nM (IC50 value) for strong binders, 888
842 and 2% (percentile rank) or 500nM (IC50 value) for all binders. 889
843 Due to the lack of recommend score thresholds for MHCflurry- 890
844 presentation, the raw presentation score was converted to a 891
845 percentile score by histogramming the presentation scores pro- 892
846 duced by 100,000 randomly generated peptides. 893

847 **Application of EnsembleMHC for the prediction of** 848 **SARS-CoV-2 MHC-I peptides**

849 **SARS-CoV-2 reference sequence.** MHC-I peptide pre- 894
850 dictions for the SARS-CoV-2 proteome were performed 895
851 using the Wuhan-Hu-1(MN908947.3) reference sequence 896
852 (<https://www.ncbi.nlm.nih.gov/genbank/sars-cov-2-seqs/>). All 897
853 potential 8-14mer peptides (n= 67,207) were derived from the 898
854 open reading frames in the reported proteome, and each peptide 899
855 was evaluated by the EnsembleMHC workflow. 900

856 **SARS-CoV-2 polymorphism analysis and protein struc-** 905
857 **ture visualizations.** Polymorphism analysis of SARS-CoV-2 906
858 structural proteins were performed using 102,148 full length 907
859 protein sequences obtained from the COVDep database⁴⁰. 908
860 Solved structures for the E (5X29) and S (6VXX) proteins 909
861 (<http://www.rcsb.org/>)⁴¹ and predicted structures for the M and 910
862 N proteins⁴² were visualized using VMD⁴³. 911

863 **Application of EnsembleMHC to determine population** 864 **SARS-CoV-2 binding capacity**

865 The peptides identified by the EnsembleMHC workflow were 912
866 used to assess the SARS-CoV-2 population binding capacity by 913
867 weighing individual MHC allele SARS-CoV-2 binding capaci- 914
868 ties by regional expression (for a schematic representation see 915
869 SI A.6). 916

870 **Population-wide MHC-I frequency estimates by country.** 917
871 The selection of countries included in the EnsembleMHC popu- 918
872 lation binding capacity assessment was based on several criteria 919
873 regarding the underlying MHC-I allele data for that country (**SI** 920
874 **A.6**). The MHC-I allele frequency data used in our model was 921
875 obtained from the Allele Frequency Net Database (AFND)¹⁸,
876 and frequencies were aggregated by country. However, the
877 currently available population-based MHC-I frequency data
878 has specific limitations and variances, which we have addressed
879 as follows:

880
881 *Quality of MHC data within countries.* We define MHC- 922
882 typing breadth as the diversity of identified MHC-I alleles 923
883 within a given country, and its depth as the ability to accurately 924
884 achieve 4-digit MHC-I genotype resolution. High variability 925
885 was observed in both the MHC-I genotyping breadth and depth 926
886 (**SI A.6 inset**). Consequently, additional filter-measures were 927
887 introduced to capture potential sources of variance within the 928

analyzed cohort of countries. The thresholds for filtering the 888
country-wide MHC-I allele data were set based on meeting 889
two inclusion criteria: 1) MHC genotyping of at least 1000 890
individuals have been performed in that population, avoiding 891
skewing of allele frequencies due to small sample size. 2) 892
MHC-I allele frequency data for at least 51 of the 52 (95%) 893
MHC-I alleles for which the EnsembleMHC was parameterized 894
to predict, ensuring full power of the EnsembleMHC workflow. 895

Ethnic communities within countries. In instances where 896
the MHC-I allele frequencies would pertain to more than one 897
community, the reported frequencies were counted towards 898
both contributing groups. For example, the MHC-I frequency 899
data pertaining to the Chinese minority in Germany would 900
be factored into the population MHC-I frequencies for both 901
China and Germany. In doing so, this treatment resolves both 902
ancestral and demographic MHC-I allele frequencies. 903
904

Normalization of MHC allele frequency data. The focus of 905
this work was to uncover potential differences in SARS-CoV-2 906
MHC-I peptide presentation dynamics induced by the 52 se- 907
lected alleles within a population. Accordingly, the MHC-I 908
allele frequency data was carefully processed in order to main- 909
tain important differences in the expression of selected alleles, 910
while minimizing the effect of confounding variables. 911

The MHC-I allele frequency data for a given population 912
was first filtered to the 52 selected alleles. These allele fre- 913
quencies were then converted to the theoretical total number 914
of copies of that allele within the population (*allele count*) 915
following 916

$$allele\ count = allele_{freq} \times 2 \times n \quad (2)$$

, where $allele_{freq}$ is the observed allele frequency in a popula- 917
tion and n is the population sample size for which that allele 918
frequency was measured. The allele count is then normalized 919
with respect to the total allele count of selected 52 alleles within 920
that population using the following relationship 921

$$norm\ allele\ count_i = \frac{allele\ count_i}{\sum_{i=1}^{52} allele\ count_i} \quad (3)$$

, where i is one of the 52 selected alleles. This normaliza- 922
tion is required to overcome the potential bias towards *hidden* 923
alleles (alleles that are either not well characterized or not sup- 924
ported by EnsembleMHC) as would be seen using alternative 925
allele frequency accounting techniques (e.g. sample-weighted 926
mean of selected allele frequencies or normalization with re- 927
spect to all observed alleles within a population (**SI A.22**)). The 928
SARS-CoV-2 binding capacity of these *hidden alleles* cannot be 929
accurately determined using the EnsembleMHC workflow, and 930
therefore important potential relationships would be obscured. 931

932 **EnsembleMHC population score.** The predicted ability of a
933 given population to present SARS-CoV-2 derived peptides was
934 assessed by calculating the EnsembleMHC Population (EMP)
935 score. After the MHC-I allele frequency data filtering steps, 23
936 countries were included in the analysis. The calculation of the
937 EnsembleMHC population score is as follows

$$EMP\ score = \frac{\sum_{i=1}^{52} peptide_{frac} \times norm\ allele\ count_i}{N_{norm\ allele\ count \neq 0}} \quad (4)$$

938 , where *norm allele count* is the observed normal-
939 ized allele count for a given allele in a population,
940 $N_{norm\ allele\ count \neq 0}$ is the number of the 52 select alleles
941 detected in a given population (range 51-52 alleles), and
942 *peptide_{frac}* is the peptide fraction or the fraction of total pre-
943 dicted peptides expected to be presented by that allele within
944 the total set of predicted peptides with a $peptide^{FDR} \leq 5\%$.

945 **Death rate-presentation correlation.** The correlation be-
946 tween the EMP score and the observed deaths per million
947 within the cohort of selected countries was calculated as a func-
948 tion of time. SARS-Cov-2 data covering the time dependent
949 global evolution of the SARS-CoV-2 pandemic was obtained
950 from Johns Hopkins University Center for Systems Science
951 and Engineering⁴⁴ covering the time frame of January 22nd to
952 April 9th. The temporal variations in occurrence of community
953 spread observed in different countries were accounted for by
954 rescaling the time series data relative to when a certain mini-
955 mum death threshold was met in a country. This analysis was
956 performed for minimum death thresholds of 1-100 total deaths
957 by day 0, and correlations were calculated at each day sequen-
958 tially following day 0 until there were fewer than 8 countries
959 remaining at that time point. The upper-limit of 100-deaths was
960 chosen to ensure availability of death-rate data on at least 50%
961 of the countries for a minimum of 7 days starting following day
962 0. Additionally, a steep decline in average statistical power is
963 observed with day 1 death thresholds greater than 100 deaths
964 (**SI A.23**).

965 The time death correlation was computed using Spear-
966 man's rank correlation coefficient (two-sided). This method
967 was chosen due to the small sample size and non-normality
968 of the underlying data (**SI A.24**). The reported correlations
969 of EMP score and deaths per million using other correlation
970 methods can be seen in supplemental figure **SI A.25**.

971 The low statistical power for some of the obtained corre-
972 lations were addressed by calculating the Positive Predictive
973 Value (PPV) of all correlations using the following equation⁴⁵

$$PPV = \frac{1 - \beta \times R}{1 - \beta \times R + \alpha} \quad (5)$$

974 , where $1 - \beta$ is the statistical power of a given correlation,
975 R is the pre-study odds, and α is the significance level. A PPV

value of $\geq 95\%$ is analogous to a p value of ≤ 0.05 . Due to an
976 unknown pre-study odd (probability that probed effect is truly
977 non-null), R was set to 1 in the reported correlations. The pro-
978 portion of reported correlations with a PPV of 95% at different
979 R values can be seen in supplemental figure **SI A.17**. The sig-
980 nificance of partitioning high risk and low risk countries based
981 on median EMP score was determined using Mann-Whitney
982 U-test. Significance values were corrected for multiple tests
983 using the Benjamini-Hochberg procedure³⁰.
984

985 **Sub-sampling of peptides from the Full SARS-CoV-2 pro-**
986 **teome.** 108 unique peptides, derived from the Full SARS-
987 CoV-2 proteome and passing the 5% $peptide^{FDR}$ filter, were
988 randomly sampled. Then, the time series EMP score - death
989 per million correlation analysis used to generate **Figure 3** was
990 applied to each sampled peptide set. The sub-sampling proce-
991 dure was repeated for 1,000 iterations (**SI A.18A**). To quantita-
992 tively describe the similarity of the distributions, the Kullback-
993 Leibler divergence (KLD), a measure of divergence between
994 two probability distributions, was calculated for the correlation
995 distribution of each sub-sample iteration relative to either the
996 correlation distribution of the Full SARS-CoV-2 proteome or
997 SARS-CoV-2 structural proteins (**SI A.18B**).

998 **Analysis of additional SARS-CoV-2 risk factors**

999 **Additional SARS-CoV-2 risk factors.** Twelve poten-
1000 tial SARS-CoV-2 risk factors (**table B.4**) were selected
1001 for analysis. Country-specific data for each risk factor
1002 was obtained from the Global Health Observatory data
1003 repository provided by the World Health Organization
1004 (<https://apps.who.int/gho/data/node.main>). Countries were se-
1005 lected for analysis based on the criteria of having reported
1006 data in the WHO data sets and inclusion in the set of 23 coun-
1007 tries for which EnsembleMHC population scores were assigned
1008 (**table B.4A**). Data regarding the total number of noncommuni-
1009 cable disease-related deaths (Cardiovascular disease, Chronic
1010 obstructive pulmonary disease, and Diabetes mellitus) were
1011 converted to deaths per million.

1012 **Correlation of additional risk factors with observed deaths**
1013 **per million.** Correlation analysis of each additional factor was
1014 carried out in a similar manner to that of the EnsembleMHC
1015 population score. In short, Spearman's correlation coefficient
1016 between each individual factor and observed deaths per mil-
1017 lion was estimated as a function of time from when a specified
1018 minimum death threshold was met (**Figure 4**). The signifi-
1019 cance level was set to $p \leq 0.05$ and significant PPV was set to
1020 $PPV \geq 0.95$ (eq 8).

1021 **Linear models of SARS-CoV-2 mortality.** For the single and
1022 combination models, individual linear models were constructed
1023 for each considered death threshold as a function of time (simi-
1024 lar to the univariate correlation analysis). Each model consisted

1025 of 1 (a single socioeconomic or health-related risk factor) or 2 (a
1026 combination of 1 risk factor and structural protein EMP score)
1027 dependent variables and deaths per million as the independent
1028 variable. The adjusted R^2 value and statistical significance
1029 of the model (F-test) were then extracted from each individ-
1030 ual model and aggregated by dependent variable (**figure 4, SI**
1031 **A.21**).

1032 The best performing models were determined by assess-
1033 ing all possible combinations of factors including structural
1034 protein EMP score. This resulted in the consideration of 4,083
1035 different linear models. The top performing models were then
1036 selected by ranking each model by median adjusted R^2 .

1037 **Code and data availability.**

1038 All data analysis and statistical tests were performed using the
1039 R Statistical Computing Environment v.3.6.0 ([http://www.r-](http://www.r-project.org)
1040 [project.org](http://www.r-project.org)). Data sets and example code are avail-
1041 able at [https://github.com/eawilson-CompBio/EnsembleMHC-](https://github.com/eawilson-CompBio/EnsembleMHC-Covid.git)
1042 [Covid.git](https://github.com/eawilson-CompBio/EnsembleMHC-Covid.git)

A Supplemental figures: Total predicted MHC-I epitope load is inversely associated with population mortality from SARS-CoV-2

Eric Wilson, Gabrielle Herneise, Abhishek Singharoy, Karen S Anderson

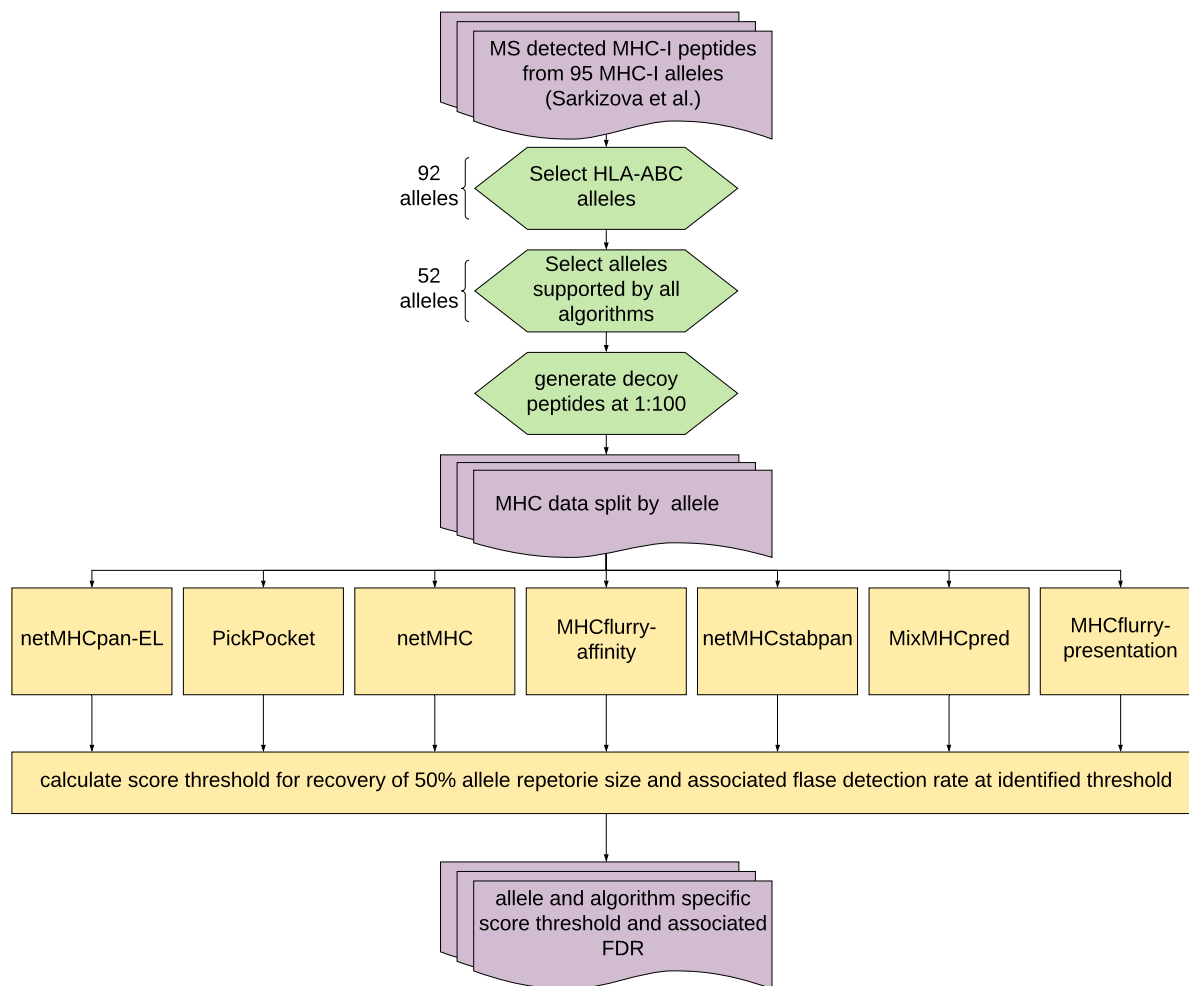


Figure A.1: EnsembleMHC Parameterization overview

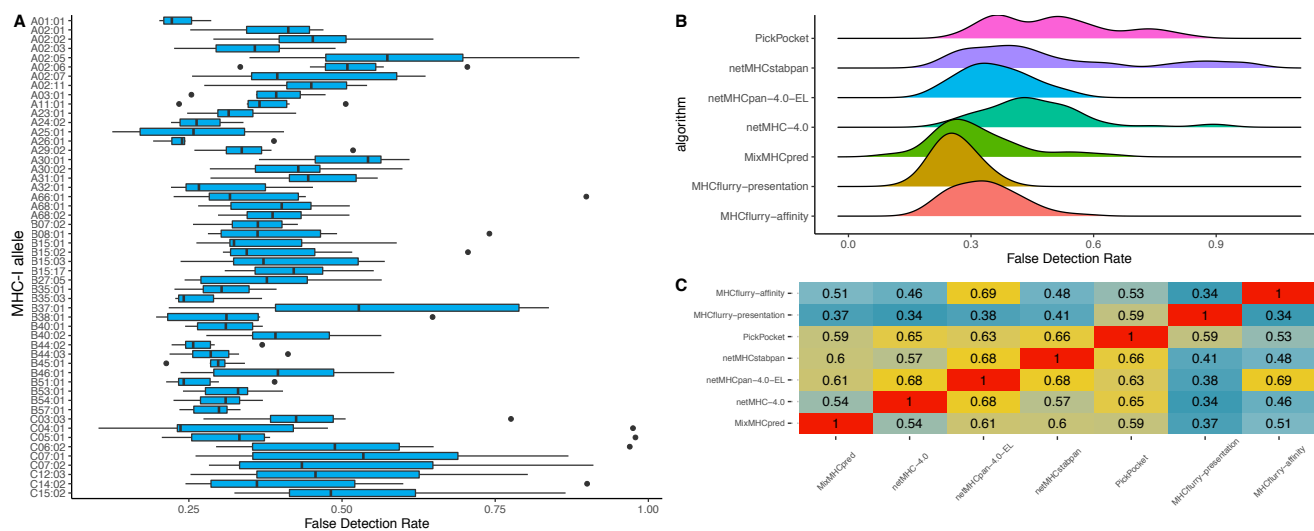


Figure A.2: EnsembleMHC prediction workflow. **A**, The EnsembleMHC score algorithm was parameterized using high quality mass spectrometry-detected MHC-I peptides paired with a 100-fold excess of randomly generated decoy peptides. Each bar represents the distribution of algorithm-specific false detection rates ($n = 7$) at that MHC allele. Each box plot is in the style of Tukey. **B**, a density plot of the observed FDRs for each algorithm across all alleles ($n = 52$). **C**, The correlation between individual peptide scores for each algorithm across all alleles was calculated using Pearson correlation. Warmer colors indicate a higher level of correlation while cooler colors indicate lower correlation.

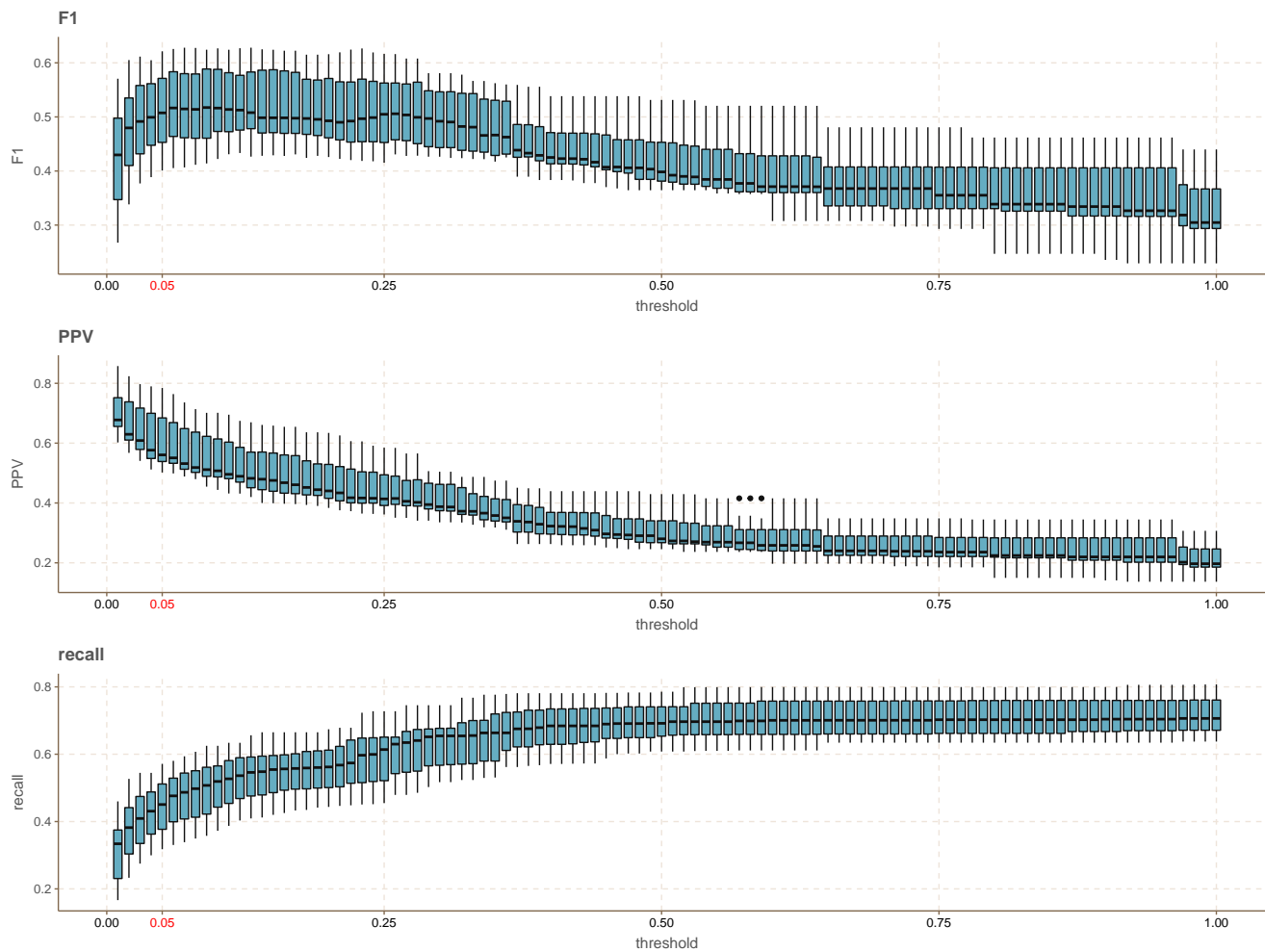


Figure A.3: **The effect of different peptide FDR threshold levels.** The effect of different $peptide^{FDR}$ cutoff thresholds on the results reported in figure 1 was evaluated for a range of 0.01-1. The $peptide^{FDR}$ selected for use in this study is highlighted in red.

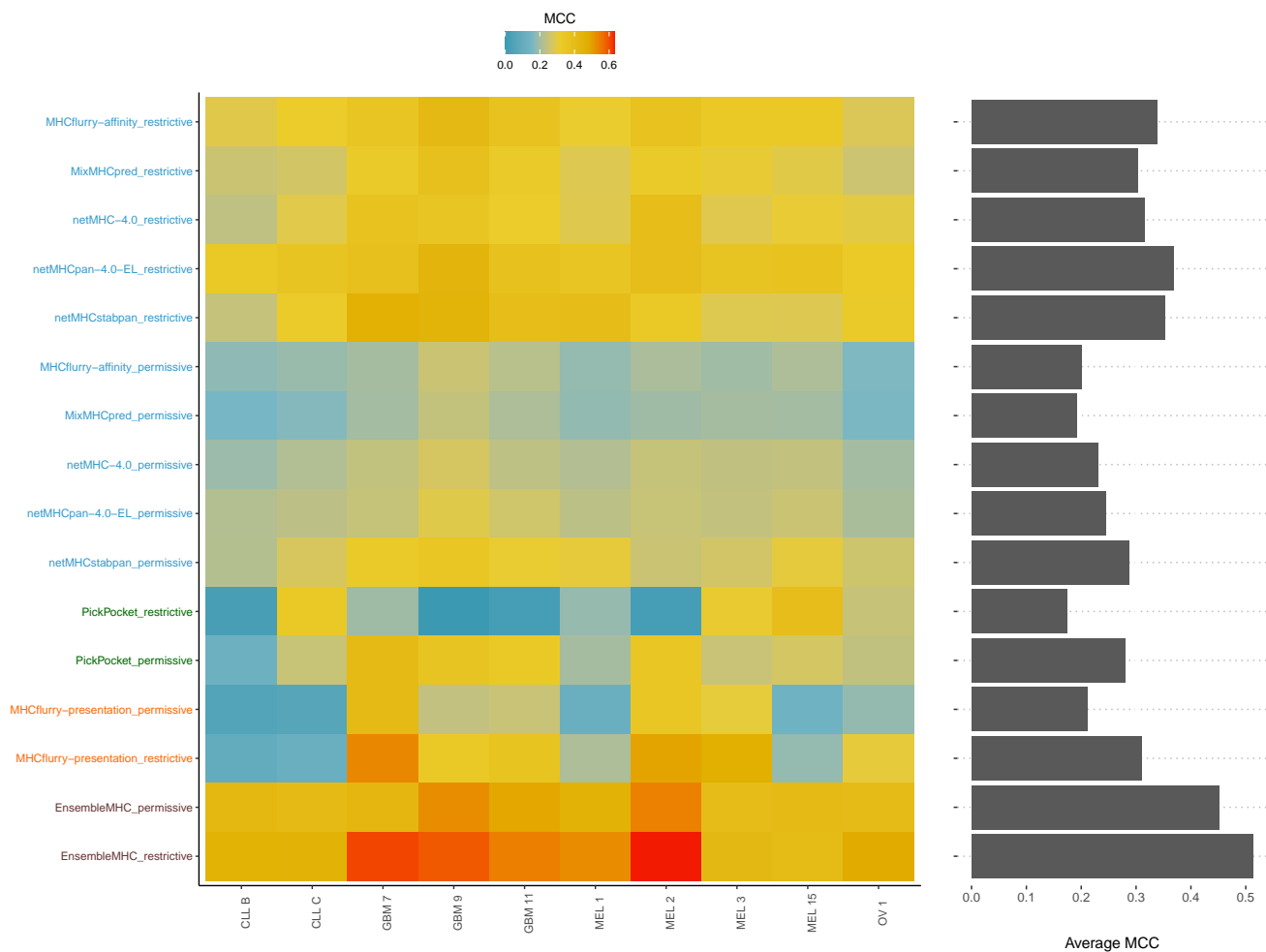


Figure A.4: **Evaluation of individual algorithms using Matthew's correlation coefficient.** As an alternative to F1 score (Figure 1B), Matthew's correlation coefficient was calculated for each algorithm. Warm colors indicate higher MCC while cooler colors indicate lower MCC. The average MCC for each algorithm is represented by the margin bar plot on the right.

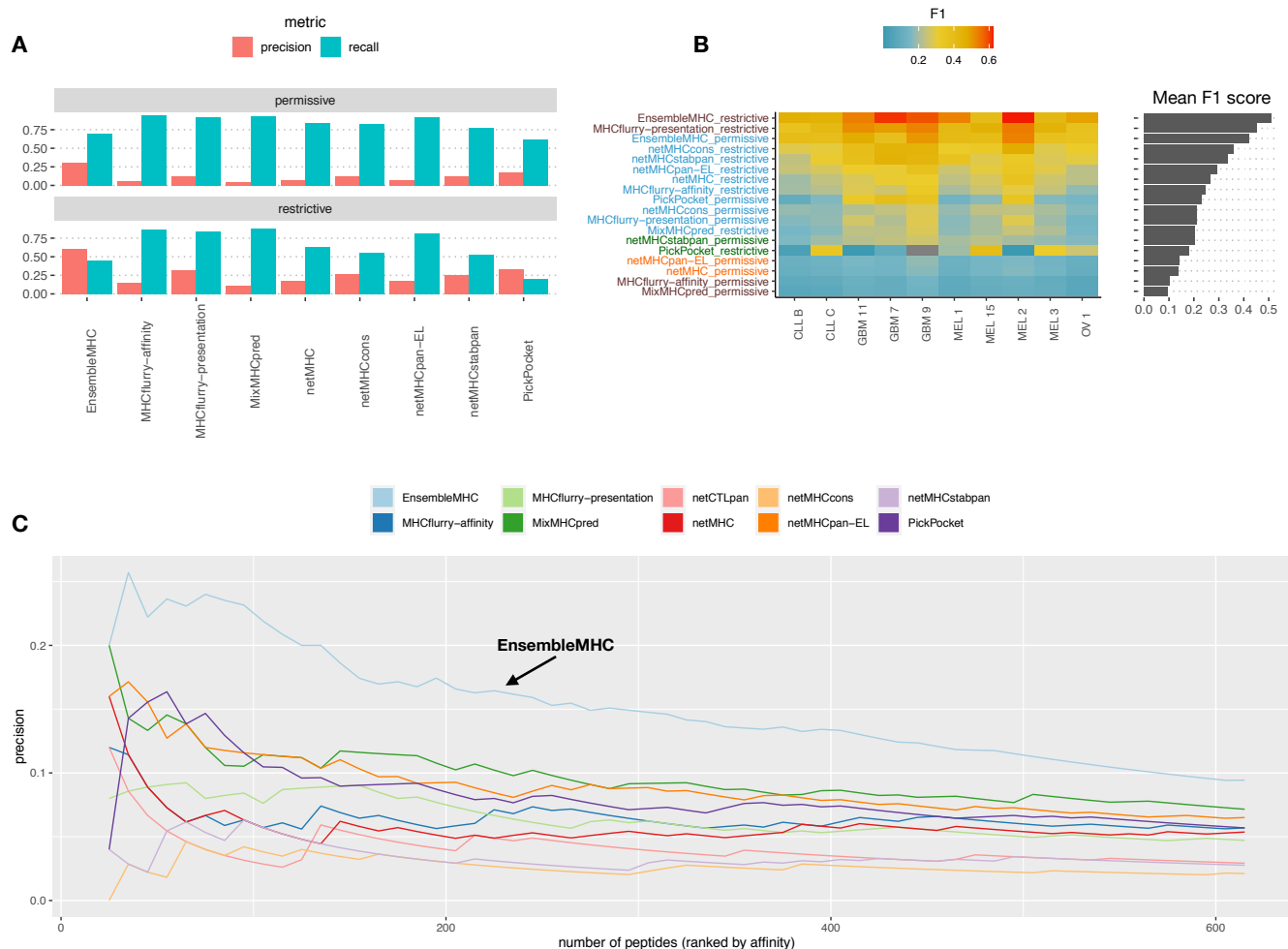


Figure A.5: **Prediction of viral immunogenic epitopes across ensemble-based algorithms.** A-B, The results in reported **figure 1 (A-B)** were compared across ensemble-based MHC-I prediction algorithms. C, The positive predictive value (precision) of the algorithms to select immunogenic MHC-I peptides was assessed across Hepatitis-C genome polyprotein (P26664), Dengue virus genome polyprotein (P14340), and the HIV-1 POL-GAG protein (P03369). All potential 8 – 14mer peptides were extracted from each protein and the resulting peptides were checked against the Immune Epitope database to identify peptides with experimentally validated immunogenicity. The result of this analysis was the generation of a data set comprised of 616 experimentally validated immunogenic peptides and 54,663 putative non-immunogenic peptides. The performance of each algorithm was then assessed by calculating the precision when selecting n number of top scoring peptides as determined by a given algorithm. Precision was calculated for each algorithm across a range of $n = 25$ to $n = 615$.

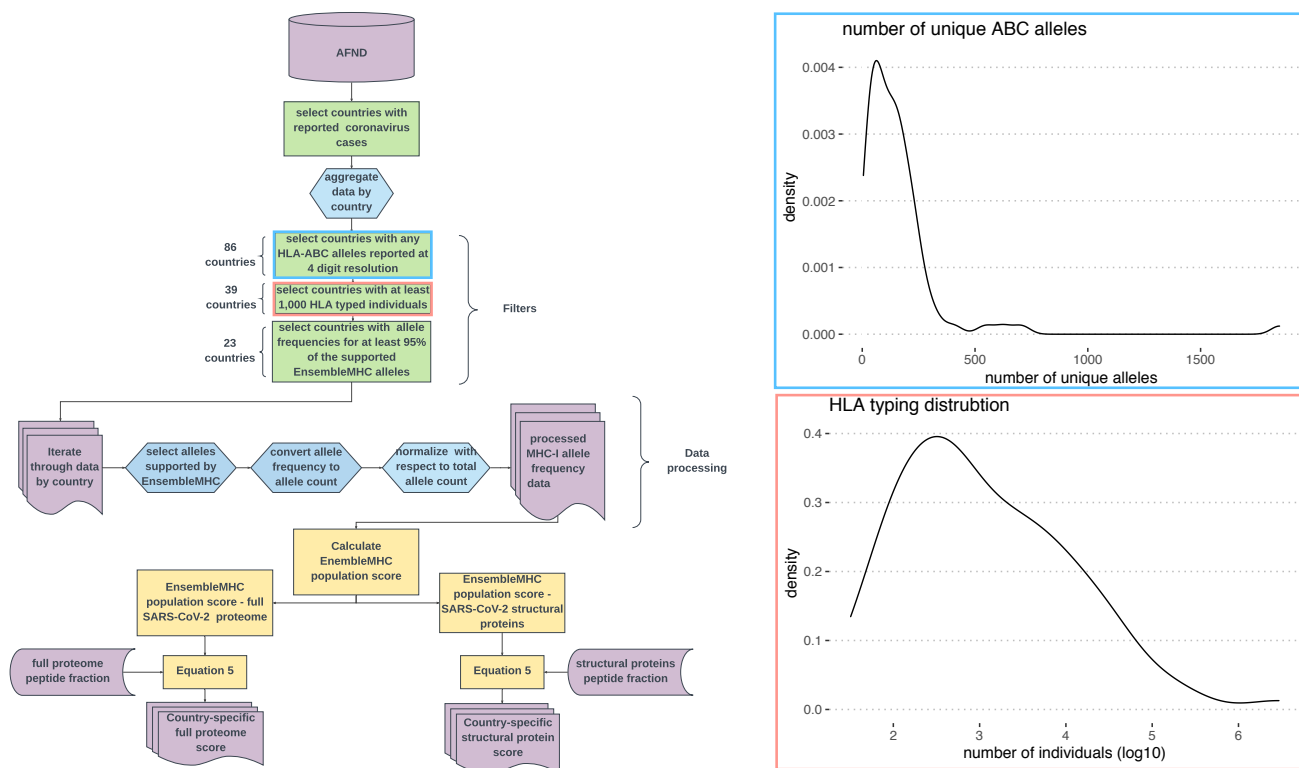


Figure A.6: **Data processing EnsembleMHC population score calculation.** The overview of the data processing steps for the global MHC-I allele frequency data and its application in the calculation the EnsembleMHC population score with respect to the full SARS-CoV-2 proteome and SARS-CoV-2 structural proteins. (inset plots), The blue inset plot illustrates MHC-typing breadth and depth variation by showing the distribution of the total number of MHC-I alleles reported at 4-digit resolution in 86 countries. The red inset plot shows the distribution of the number of MHC-genotyped individuals in the set of countries with at least 1 reported coronavirus case. **AFND = Allele Frequency Net Database**

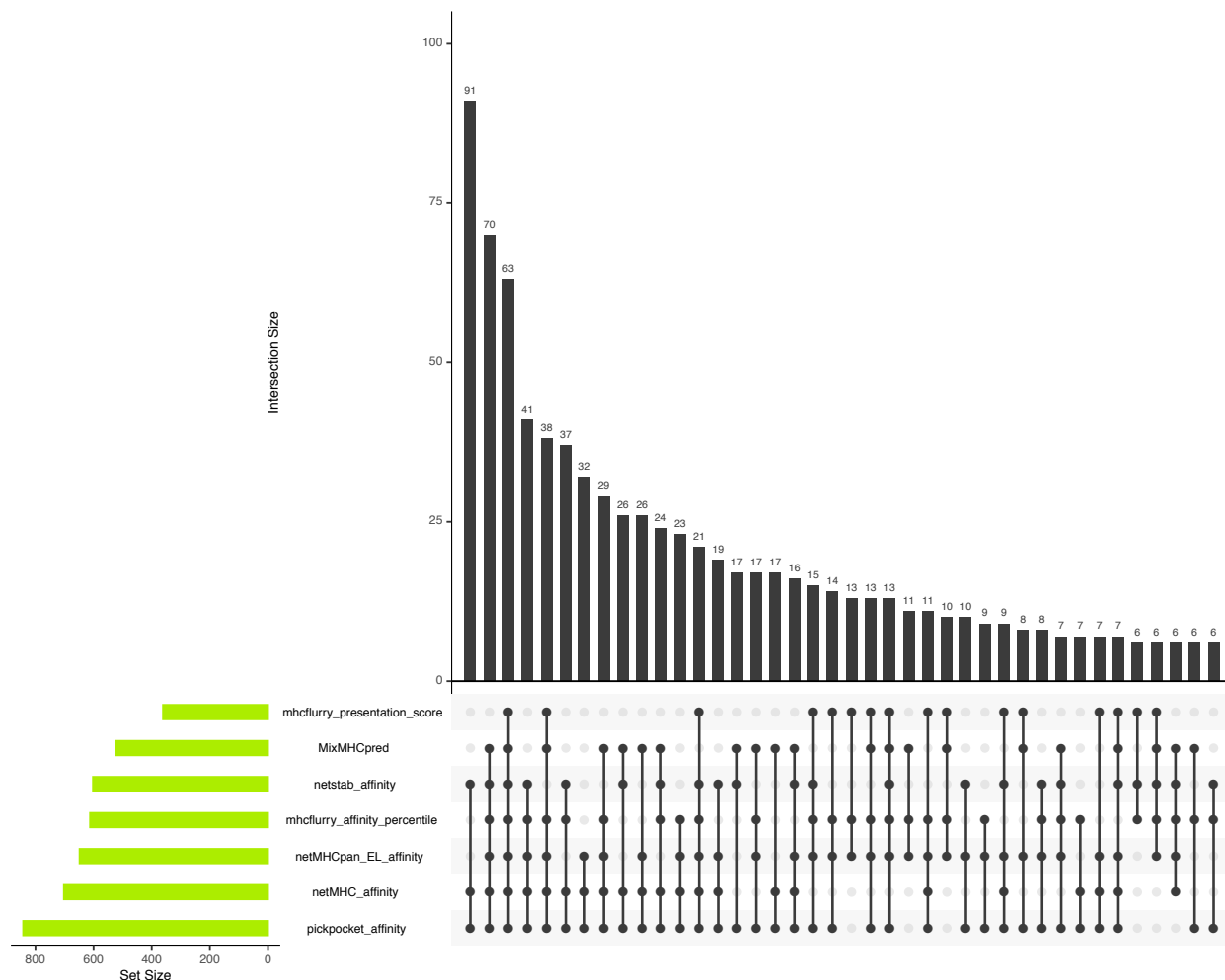


Figure A.7: **Contribution of each EnsembleMHC component algorithm to peptide selection for predicted SARS-CoV-2 peptides.** The UpSet plot shows the contribution of each individual component algorithm to the 658 unique SARS-CoV-2 peptides identified by EnsembleMHC. The top bar plot indicates the number of unique peptides identified by the combination of algorithms shown by the points and segments located under each bar. The bar plot on the left-hand side of the plot indicates the total number of peptides identified by each algorithm.

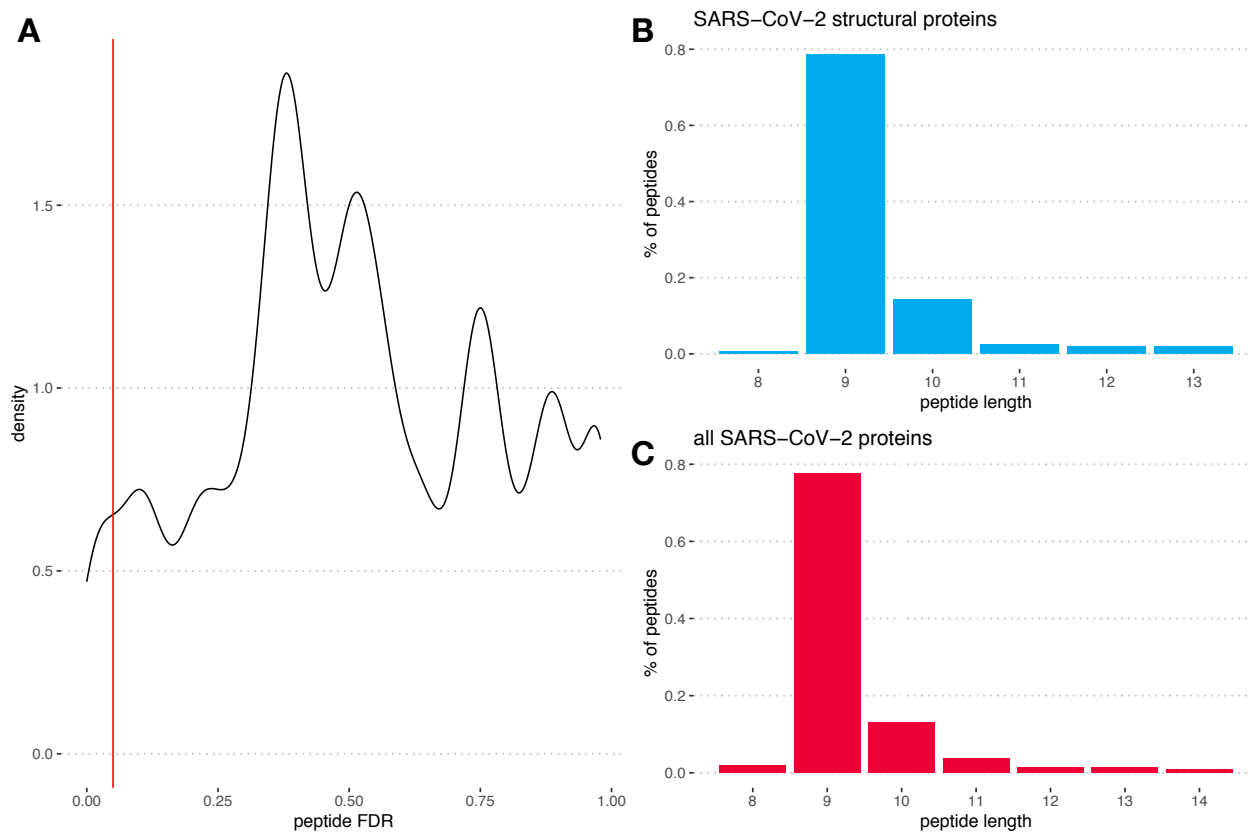


Figure A.8: **EnsembleMHC peptide^{FDR} and length distributions of predicted SARS-CoV-2 MHC-I peptides.** **A**, The distribution of the $peptide^{FDR}$ for 9,712 SARS-CoV-2 peptides that fell with the score threshold of at least one component algorithm. The red line indicates an $peptide^{FDR}$ level of $\leq 5\%$. **B**, The length distribution of the 108 high-confidence peptides identified from SARS-CoV-2 structural proteins. **C**, The length distribution of the 658 high-confidence peptides identified from full SARS-CoV-2 proteome.

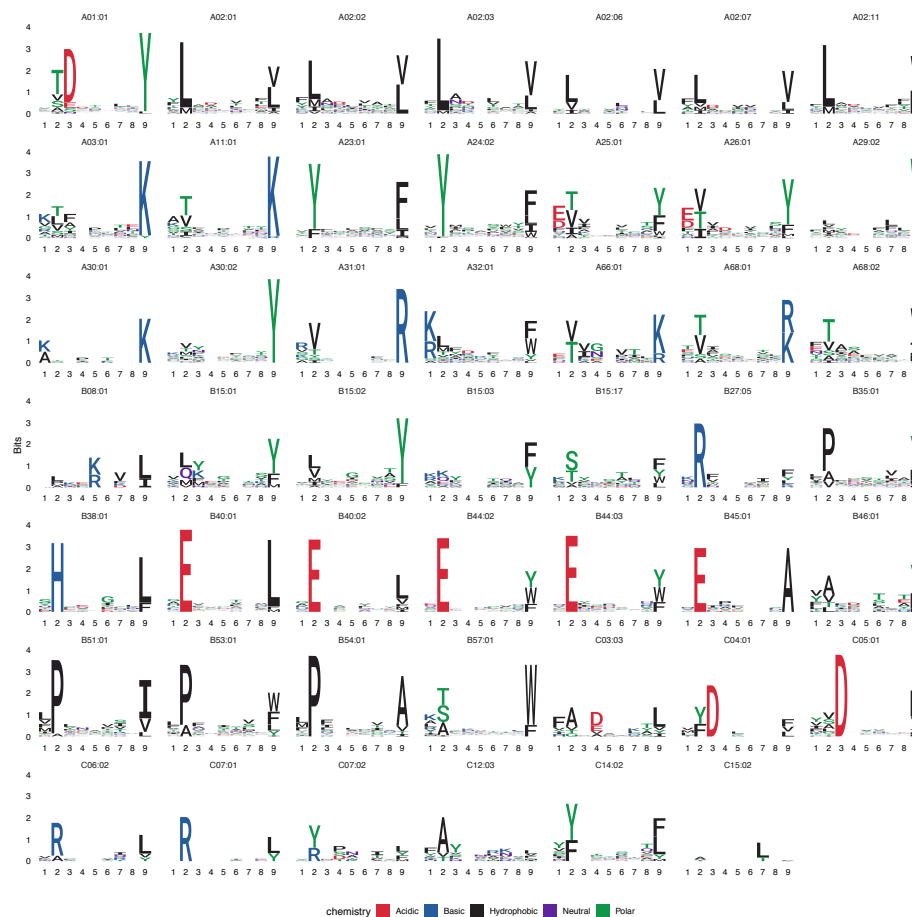


Figure A.9: Logo plots for the identified peptides from the SARS proteome. Logo plots were generated for MHC alleles with at least 5 peptides identified by EnsembleMHC prediction. Peptides shorter than 9 amino acids had random amino acid inserted into a non-anchor position while peptides longer than 9 amino acids had a random non-anchor position deleted. Large amino acid character height indicates a high frequency of that amino acid at that position. Amino acids are colored residue type.

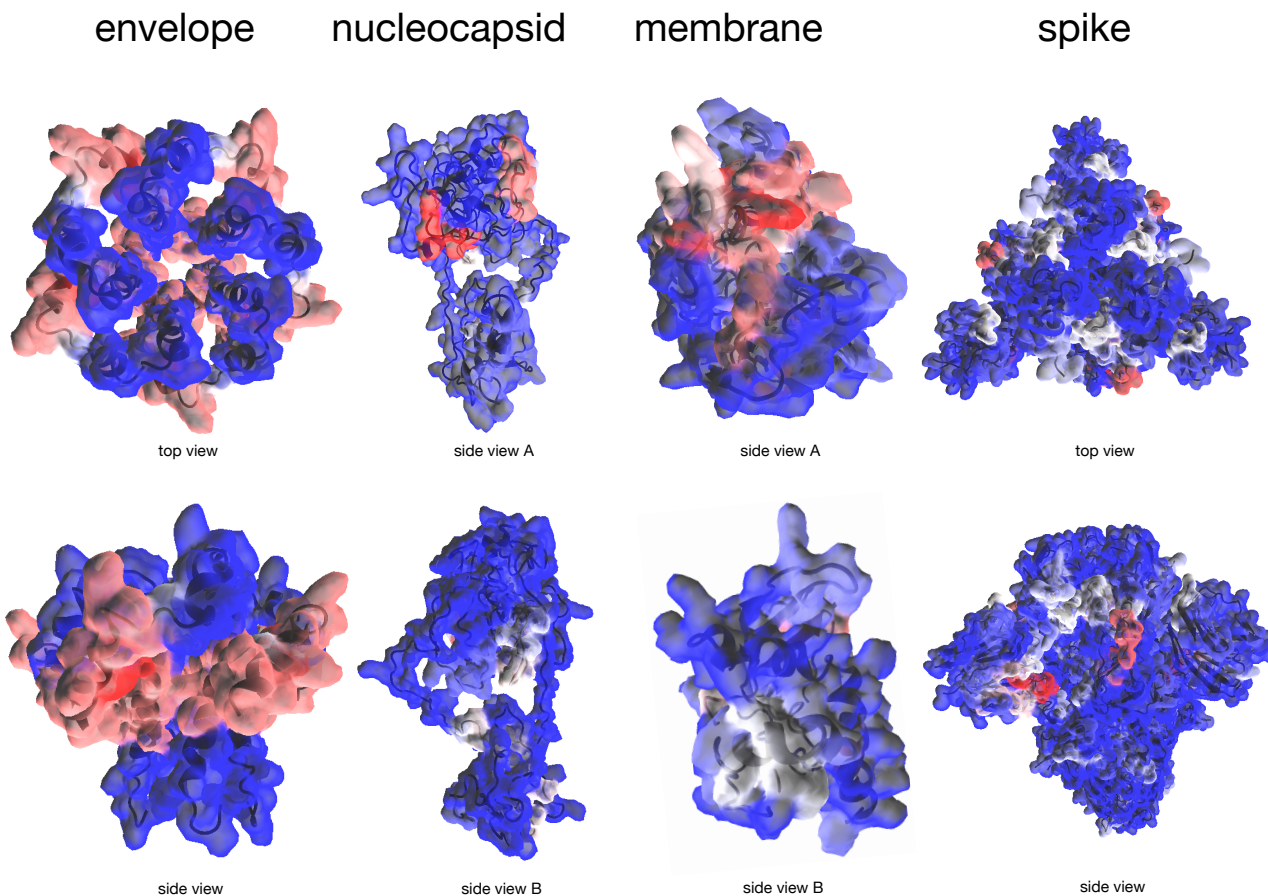


Figure A.10: **Molecular origin of predicted SARS-CoV-2 structural protein MHC-I peptides.** The predicted SARS-CoV-2 structural protein MHC-I peptides were mapped onto the solved structures for the envelope and spike proteins, and the predicted structures for the nucleocapsid and membrane proteins. Red highlighted regions indicate an enrichment of predicted peptides while blue regions indicate a depletion of predicted peptides.

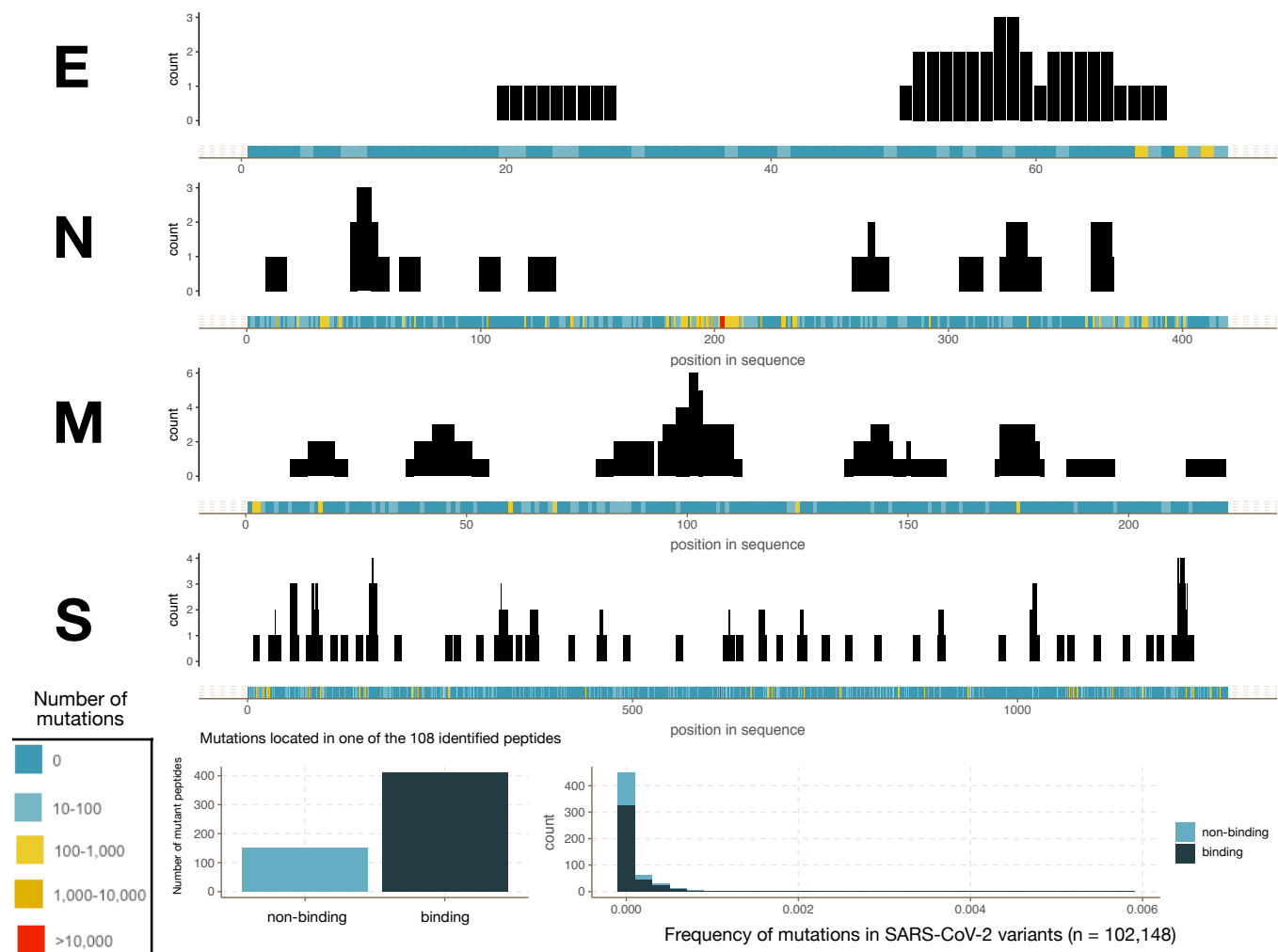


Figure A.11: **Impact of SARS-CoV-2 sequence polymorphism on predicted SARS-CoV-2 structural MHC-I peptides.** Top 4 panels, The incidence of sequence position mutations (colored bar) and the frequency of each amino acid position in one of the 108 SARS-CoV-2 structural protein peptides (black bars) by protein were calculated for 102,148 SARS-CoV-2 sequence variants. Lower left panel, all potential mutations arising in an EnsembleMHC-predicted MHC-I peptides were evaluated for changes in binding affinity ($peptide^{FDR} > 0.05$). Lower right panel, The overall frequency of mutations impacting EnsembleMHC-predicted peptides with light blue indicating deleterious mutations, and dark blue indicating neutral mutations.

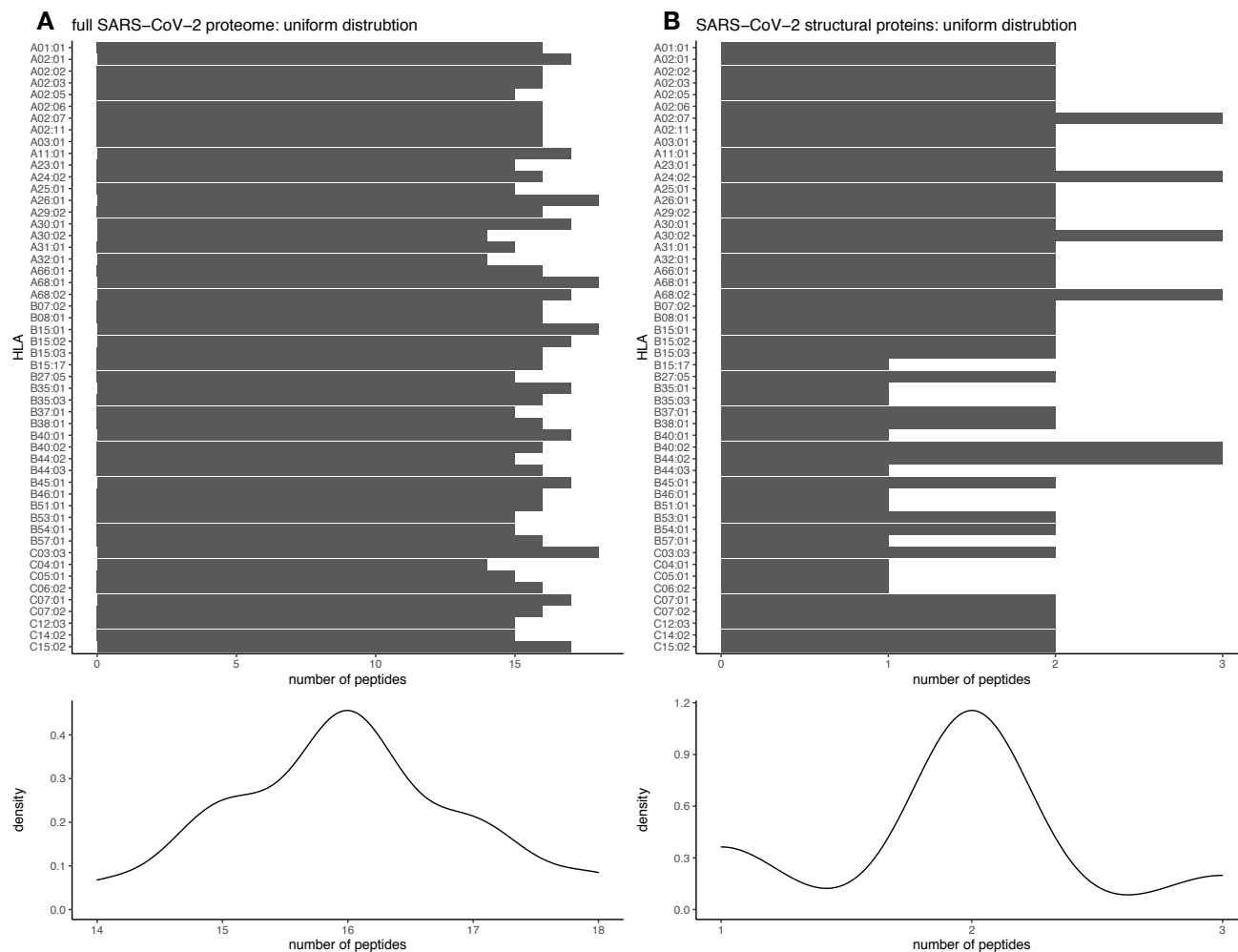


Figure A.12: Simulated even peptide-allele distribution. The resulting SARS-CoV-2 peptide-MHC allele distribution was compared to an even distribution by means of the Kolmogorov-smirnov test. **A**, An even peptide-allele distribution for the full SARS-CoV-2 proteome was simulated by sampling from a discrete symmetrical distribution centered around the median number of peptides assigned to individual alleles ($\bar{X} = 16$). **B**, An even peptide-allele distribution for the SARS-CoV-2 structural proteins was simulated by sampling from a discrete symmetrical distribution centered around the median number of peptides assigned to individual alleles ($\bar{X} = 2$).

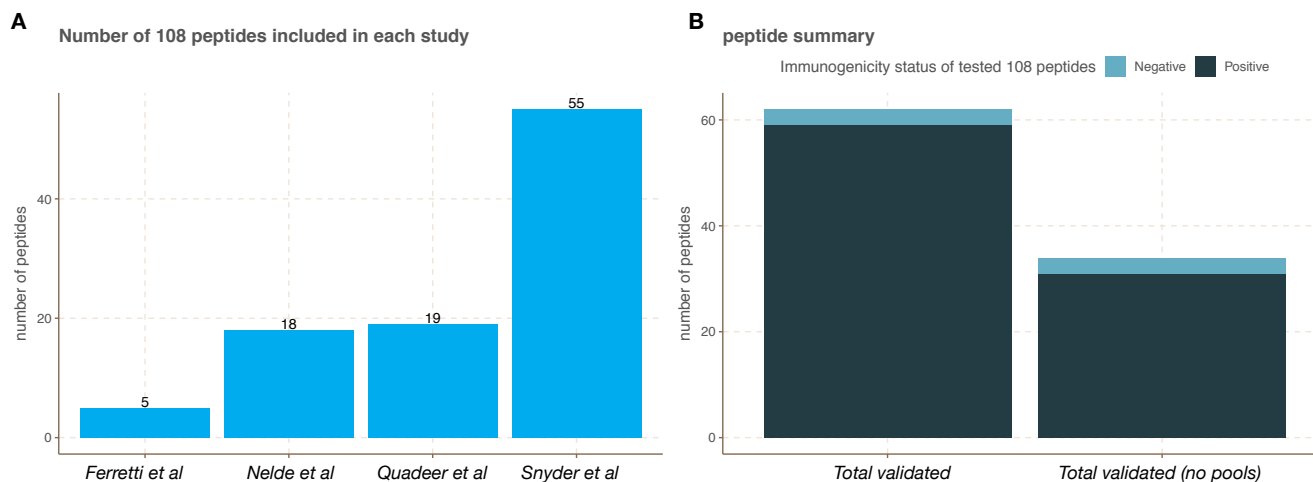


Figure A.13: **External experimental validation of the 108 high confidence SARS-Cov-2 structural protein peptides.** Experimentally validated immunogenic peptides derived from SARS-CoV-2 structural proteins were obtained from 4 independent studies (Ferretti et al⁴⁶, Nelde et al⁴⁷, and Snyder et al⁴⁸) and 1 meta-study (Quadeer et al⁴⁹). These peptides were then assessed for overlap with the 108 SARS-CoV-2 peptides identified by EnsembleMHC (*108 peptides*). **A**, The total number of peptides from the *108 peptides* set that were included for testing in each study. **B**, The summary of immunogenicity status of *108 peptides* across all studies. These summaries were split into two groups. *Total validated* indicates the total number of experimentally validate 108 peptides while *total validated (no pools)* indicates the number of experimentally validated peptides excluding those only tested in peptide pools. This distinction was made due to the potential of peptide pools to obscure which tested peptide is truly responsible for the observed immune response. Overall, 57% of the predicted 108 structural protein peptides were tested with 95% of tested peptides producing an immune response.

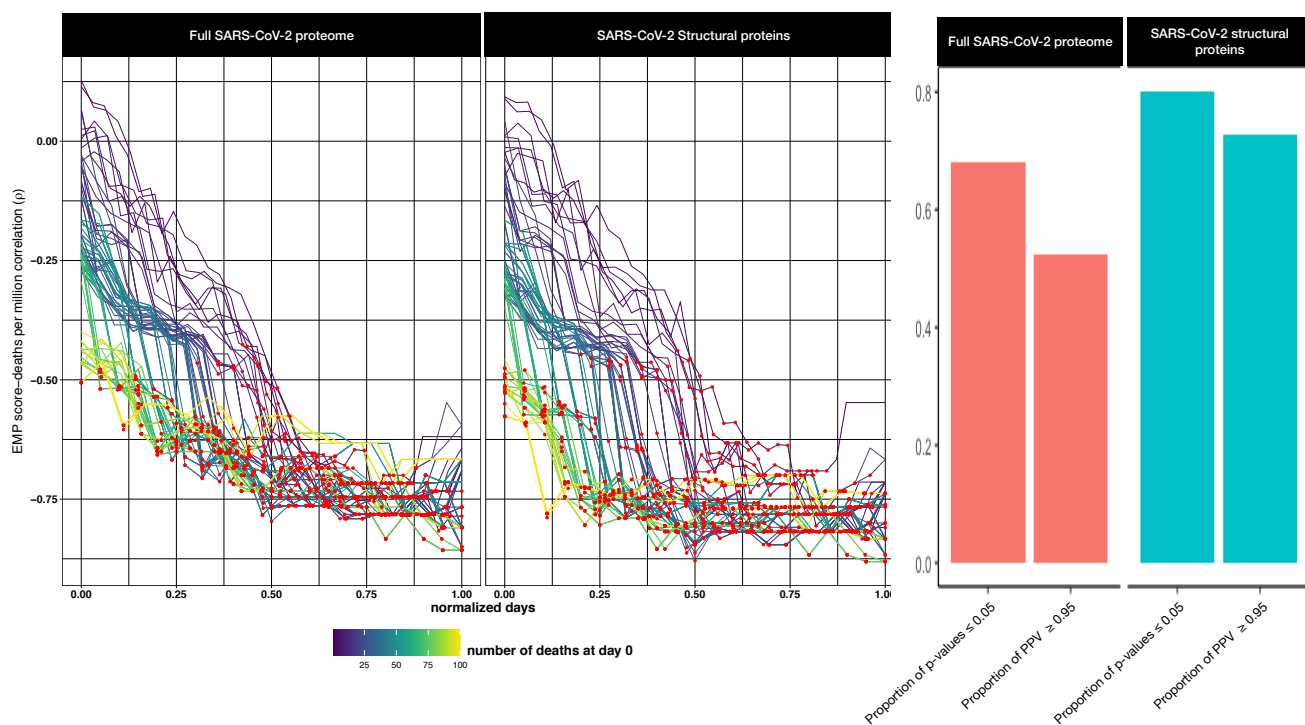


Figure A.14: **Correlation of EMP score based on full SARS-CoV-2 proteome or SARS-CoV-2 structural proteins with observed deaths per million.** **A**, The correlations between EnsembleMHC population score based on the full SARS-CoV-2 proteome (**left**) or EnsembleMHC population score based on SARS-CoV-2 structural proteins (**right**). **B**, The difference in the proportions of significant p-values and PPV between the full SARS-CoV-2 proteome (**left**) and SARS-CoV-2 structural proteins (**right**) (not corrected for multiple testing).

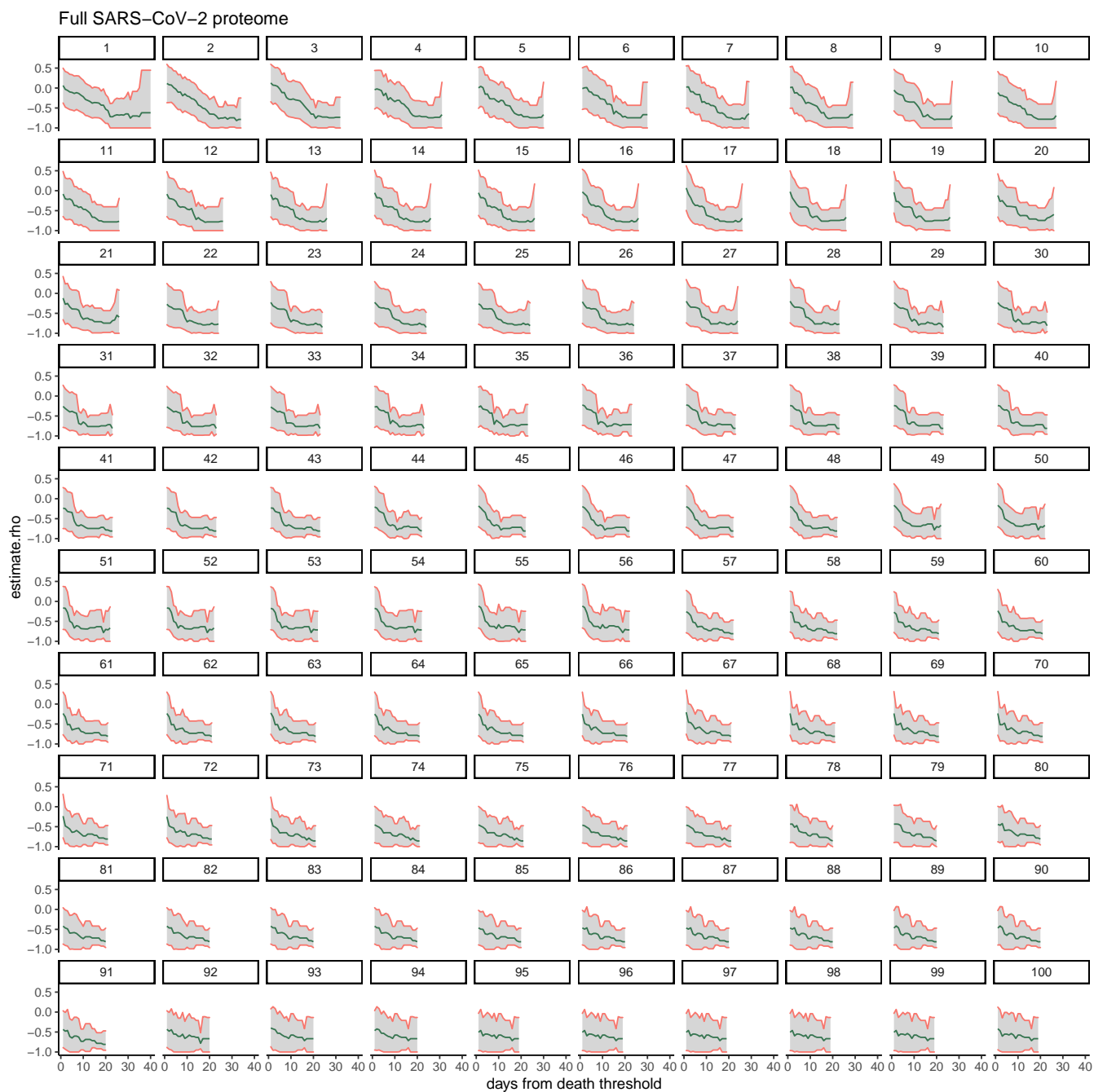


Figure A.15: **95% Confidence interval for the correlations between the EnsembleMHC score based on the full SARS-CoV-2 proteome and observed deaths per million.** Each individual plot shows the 95% confidence interval (grey region) for the correlations between EMP scores based on the full SARS-CoV-2 proteome and observed deaths per million (blue line) for all starting minimum death thresholds (indicated by number above plot).

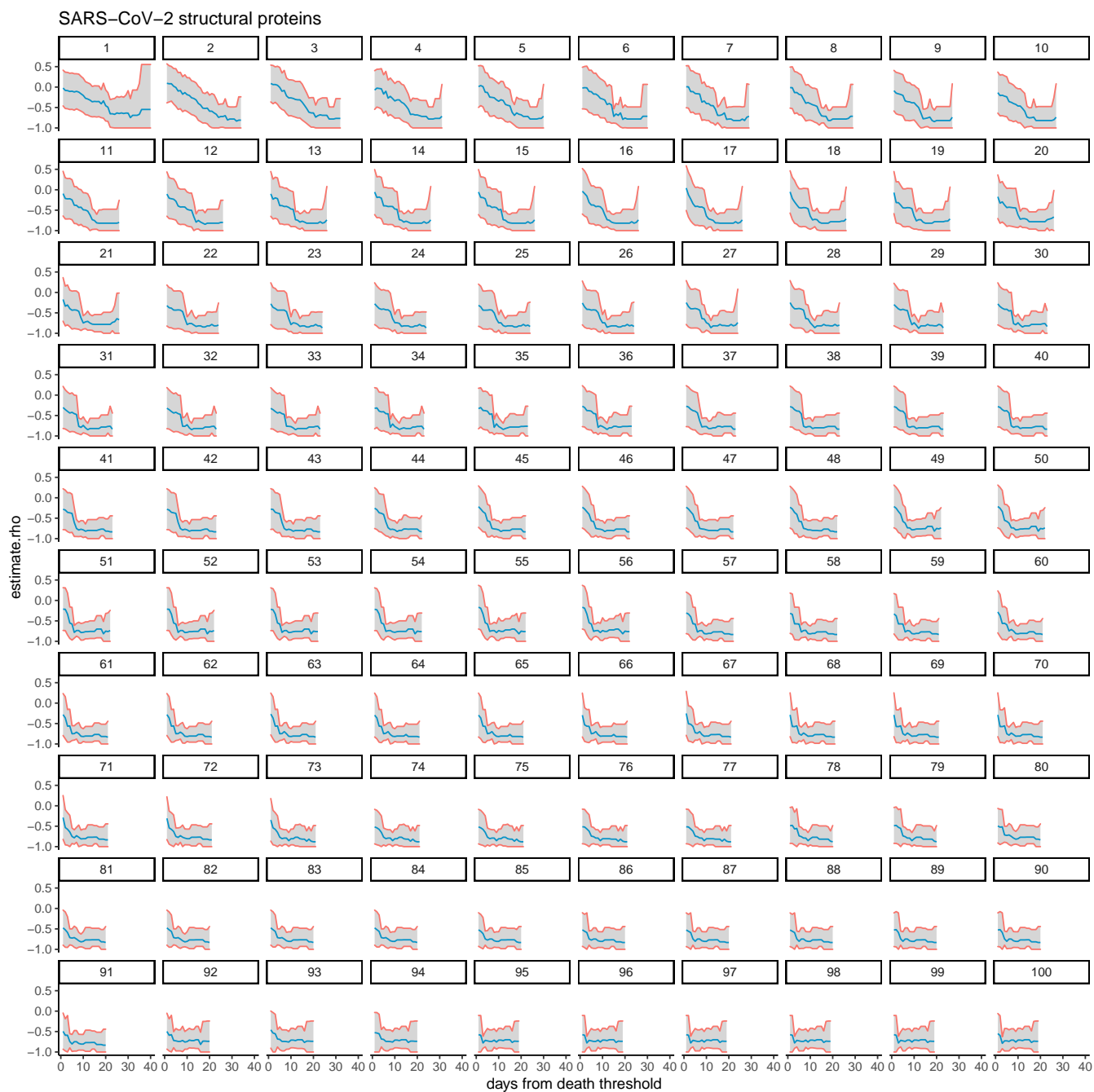


Figure A.16: **95% Confidence interval for the correlations between the EnsembleMHC score based on SARS-CoV-2 structural proteins and observed deaths per million.** Each individual plot shows the 95% confidence interval (grey region) for the correlations between EMP scores based on SARS-CoV-2 structural proteins and observed deaths per million (blue line) for all starting minimum death thresholds (indicated by number above plot).

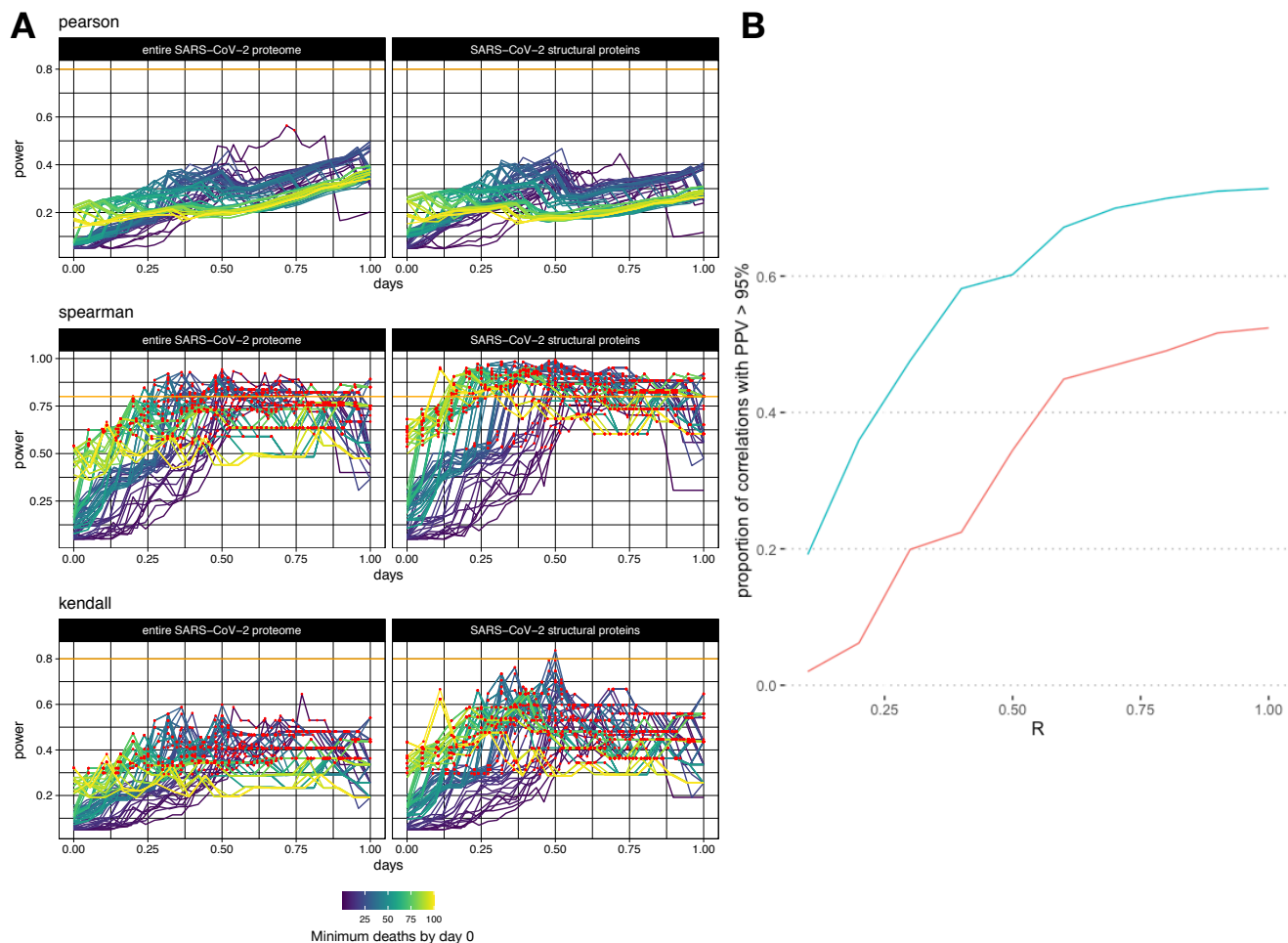


Figure A.17: **Analysis of statistical power for different correlation methods.** **A**, The statistical power of each reported correlation between EnsembleMHC population score with respect to the full SARS-CoV-2 proteome (left column) or specifically structural proteins (right column) and deaths per million were calculated at each day starting from the day a country passed a particular minimum death threshold using Pearson's correlation (**top**), Spearman's rho (**middle**), and Kendall's tau (**bottom**). The days from each start point were normalized, and correlations that were shown to be statistically significant are colored with a red point. The orange line indicates a power threshold of 80%. **B**, The line plot on the right shows the proportion of points achieving a significant PPV at different thresholds for pre-study odds (R) for the spearman correlation carried out in section 3. The blue line represents the proportion of significant correlations for EnsembleMHC score based on SARS-CoV-2 structural proteins while the red line represents the same correlations with an EnsembleMHC population score based on the full SARS-CoV-2 proteome.

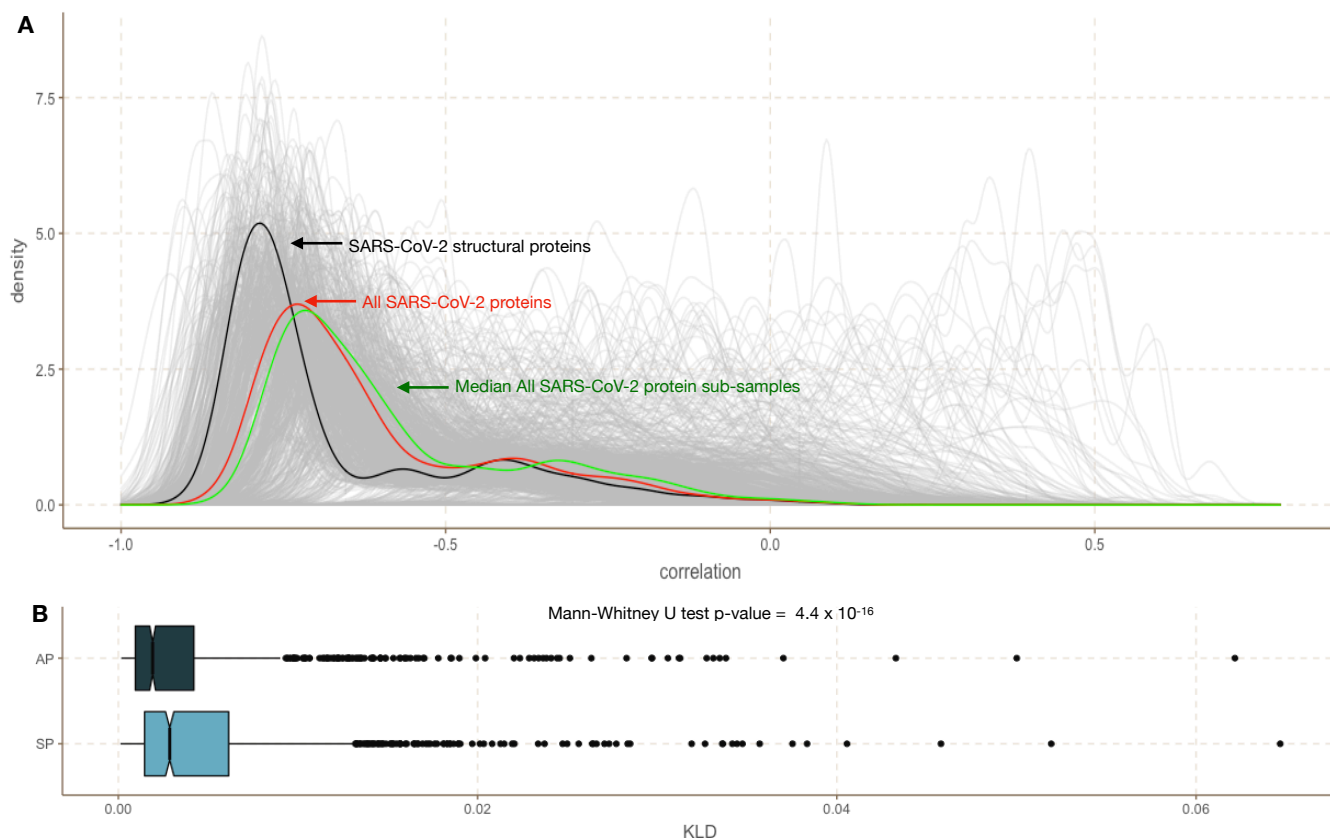


Figure A.18: **The effect of sub-sampling EnsembleMHC-identified MHC-I peptides derived from the full SARS-CoV-2 proteome on EMP score - deaths per million correlation distribution.** The robustness of the observed distinction between the EMP score - deaths per million correlation distributions between SARS-CoV-2 structural proteins and all SARS-CoV-2 proteins was assessed by performing sub-sampling of non-structural SARS-CoV-2 MHC-I peptides. **A**, 1,000 sub-sampling iterations were performed by randomly selecting 108 peptides from the full SARS-CoV-2 proteome that passed the 5% $peptide^{FDR}$ filter. The correlation between the population EMP score produced by each sub-sampled set of peptides and observed deaths per million were plotted (grey lines). The correlation distribution observed for identified SARS-CoV-2 structural protein peptides (black line), all SARS-CoV-2 proteins (red line), and the median correlation distribution across all subsampling iterations (green line) were plotted for comparison. **B**, Kullback-Leibler divergence was calculated for the correlation distribution of each down sample iteration relative to either the correlation distribution of the all peptide group (AP) or the structural peptide group (SP).

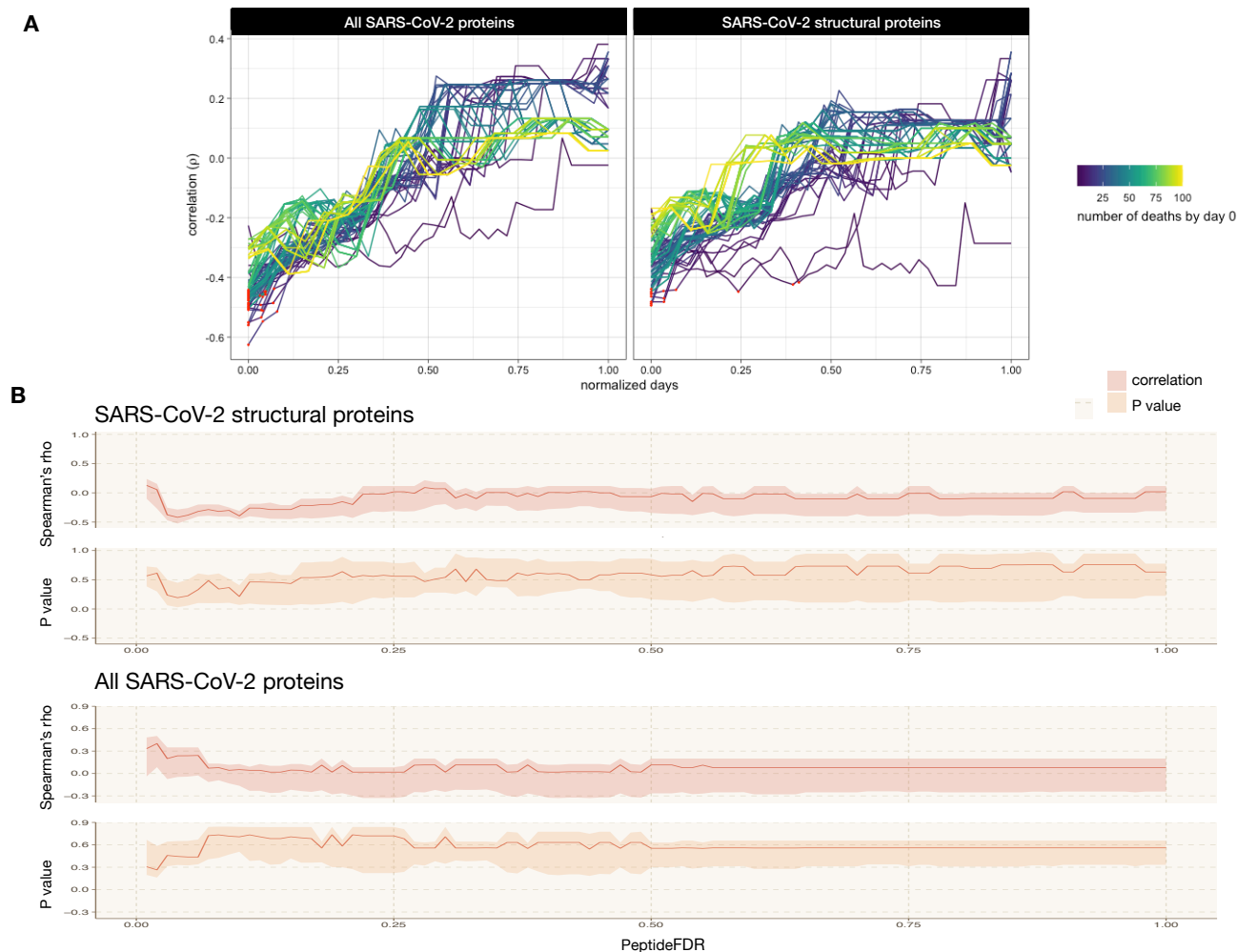


Figure A.19: **The effect of MHC shuffling on the correlation of EMP score and deaths per million.** **A**, The MHC-I allele assessment of peptides that passed an individual algorithm binding affinity thresholds were shuffled prior to $peptide^{FDR}$ filtering. The red points indicate correlations with a p-value $\leq 5\%$. **B**, The impact of varying $peptide^{FDR}$ cutoff threshold on the shuffled MHC data set. For each $peptide^{FDR}$ cutoff threshold (x-axis), the upper bound of the shaded region indicates the 75th percentile, the lower bound indicates the 25th percentile, and the solid line indicates the median.

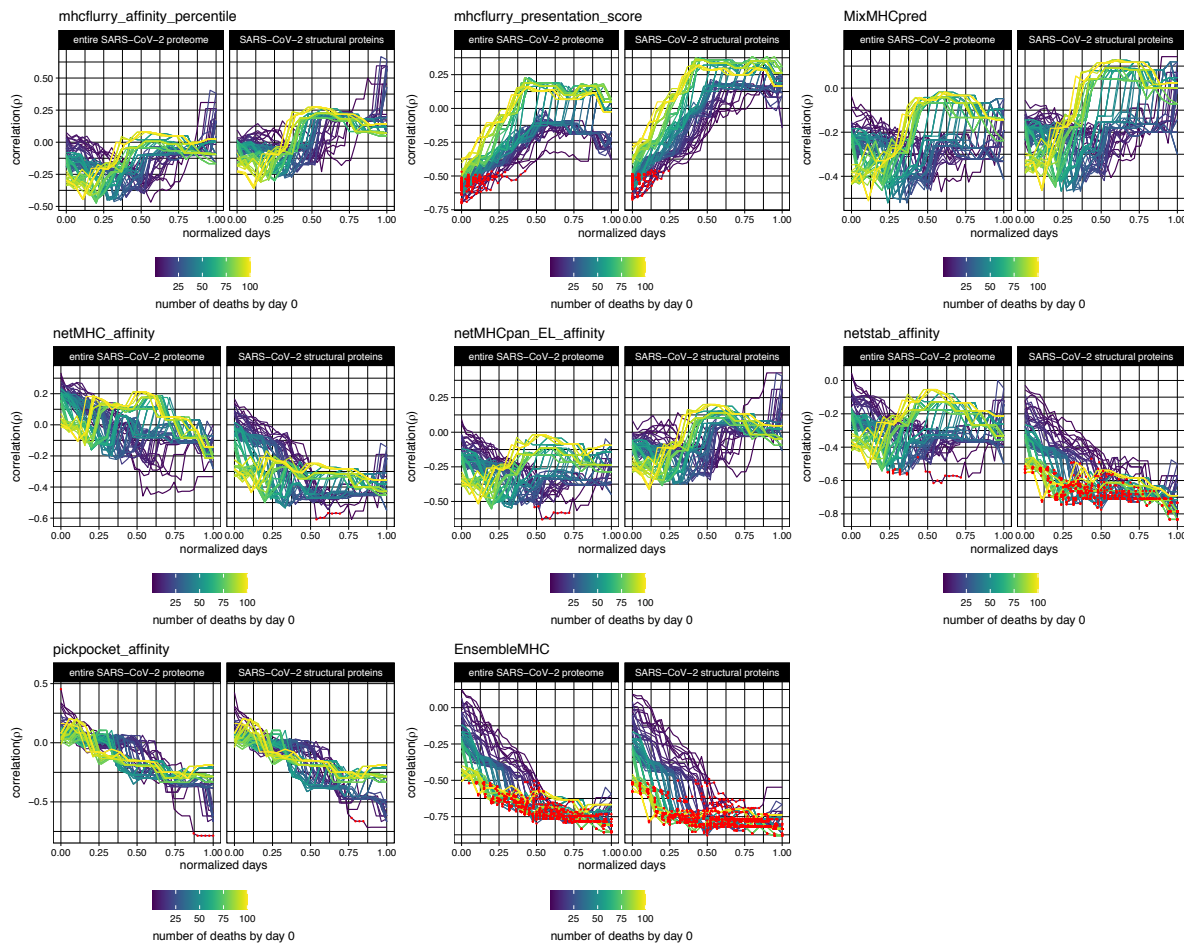


Figure A.20: **Individual algorithms are unable to fully recreate the correlation with population mortality reported by EnsembleMHC.** Population SARS-CoV-2 binding capacities using only single algorithms were correlated to observed deaths per million. For each algorithm, the population SARS-CoV-2 binding capacity was calculated from the resulting viral peptide-MHC allele distribution using restrictive MHC-I binding affinity cutoffs ($\leq 0.5\%$ for binding percentile scores, top 0.5% MHCflurry presentation score, and $\leq 50nm$ for PickPocket). Red points indicate a PPV $\geq 95\%$.

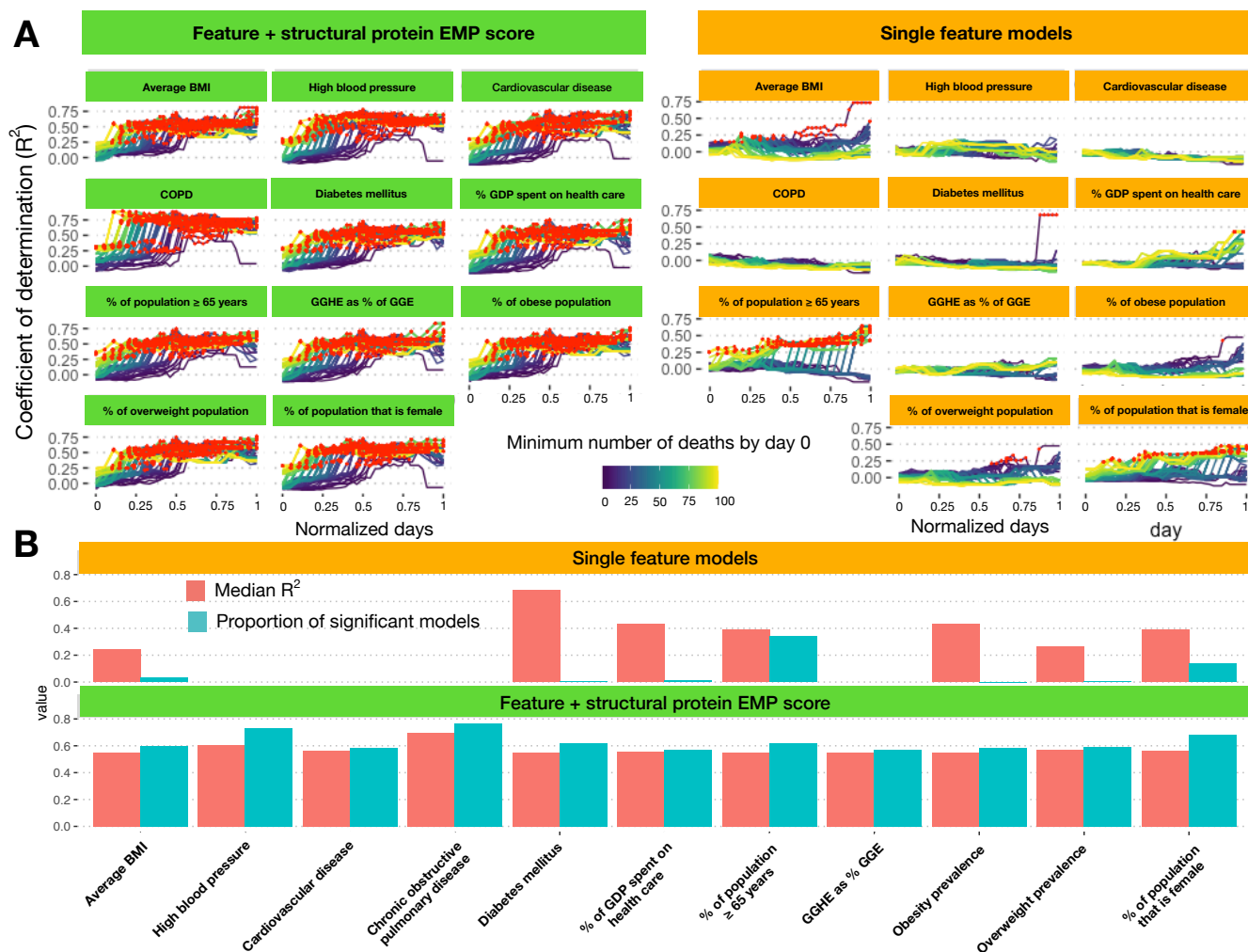


Figure A.21: **Addition of structural protein EMP score significantly improves linear model fit to observed deaths per million.** **A**, Linear models were constructed using either a single risk factor (yellow) or a combination of a risk factor and structural protein EMP scores (green). The x-axis indicates the number of normalized days from when a minimum death threshold was met (line color), and the y-axis indicates the observed adjusted R^2 value. **B**, A summary of results obtained from single feature linear models (top panel, yellow) or the combination models (bottom panel, green). The red bars indicate the median R^2 value achieved by that model and the blue bars indicate the proportion of regressions that were found to be significant ($F\text{-test} \leq 0.05$).

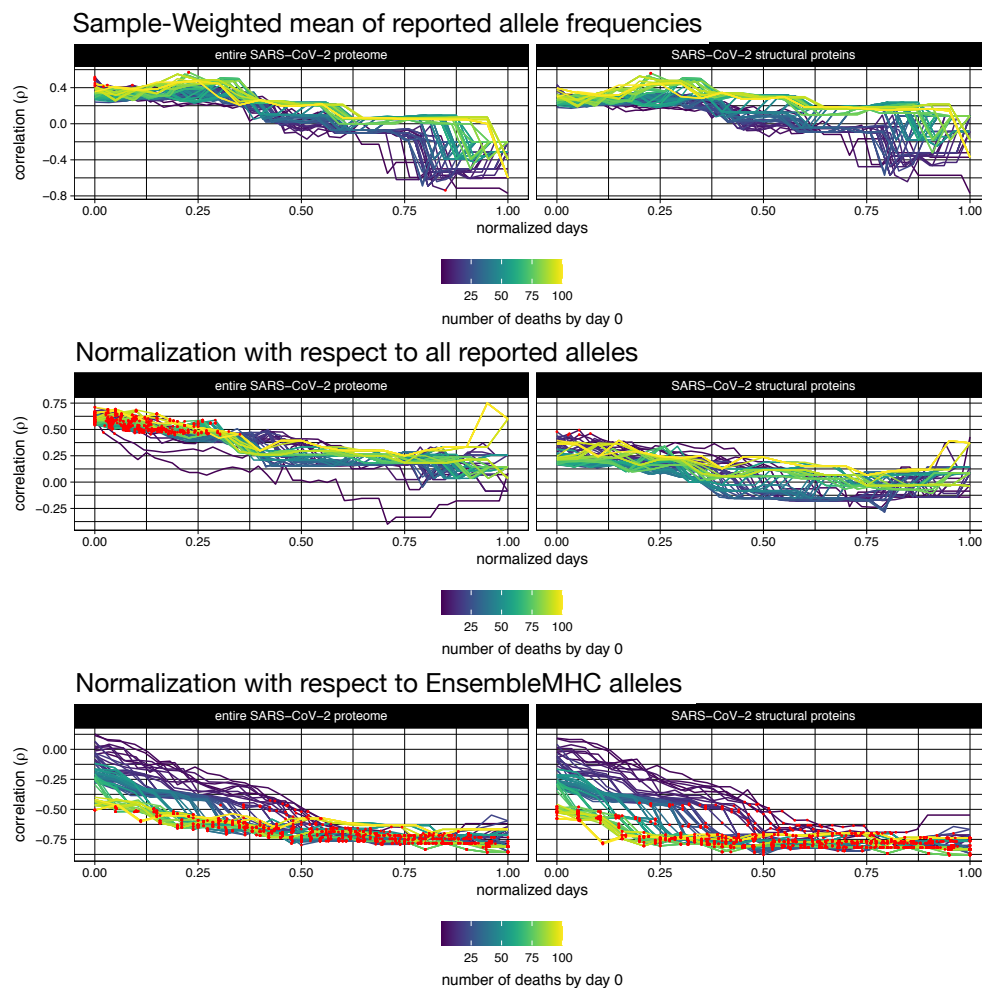


Figure A.22: **EnsembleMHC population score and deaths per million correlation using different allele frequency accounting methods.** The effect of different allele frequency normalization techniques on the reported correlations between SARS-CoV-2 mortality and EMP scores based on the full SARS-CoV-2 proteome (left column) or SARS-CoV-2 structural proteins (right column). **Top panel**, The aggregation of allele frequencies within a particular country by taking the sample-weighted mean of reported frequencies for the 52 selected MHC-I alleles. **Middle panel**, Normalizing allele count with respect to all detected alleles in a given population. **Bottom panel**, Normalizing allele count with respect to only the 52 select alleles.

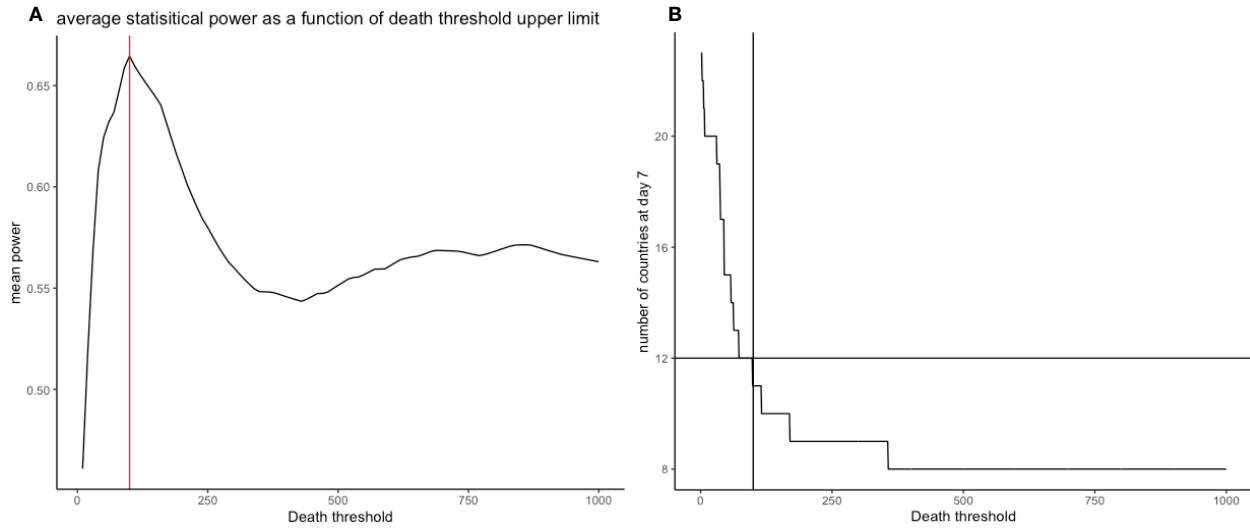


Figure A.23: **Justification of upper limit for death threshold.** **A**, The mean statistical power of all resulting correlations between EnsembleMHC population scores and observed deaths per million at different minimum reported death thresholds. The red line indicates a minimum death threshold of 100 deaths by day 0, the selected upper limit for analysis. **B**, The number of countries remaining at day seven using different minimum death thresholds. The cross bar indicates that there would be more than half of the considered countries remaining at day 7 when using the 100 minimum death threshold.

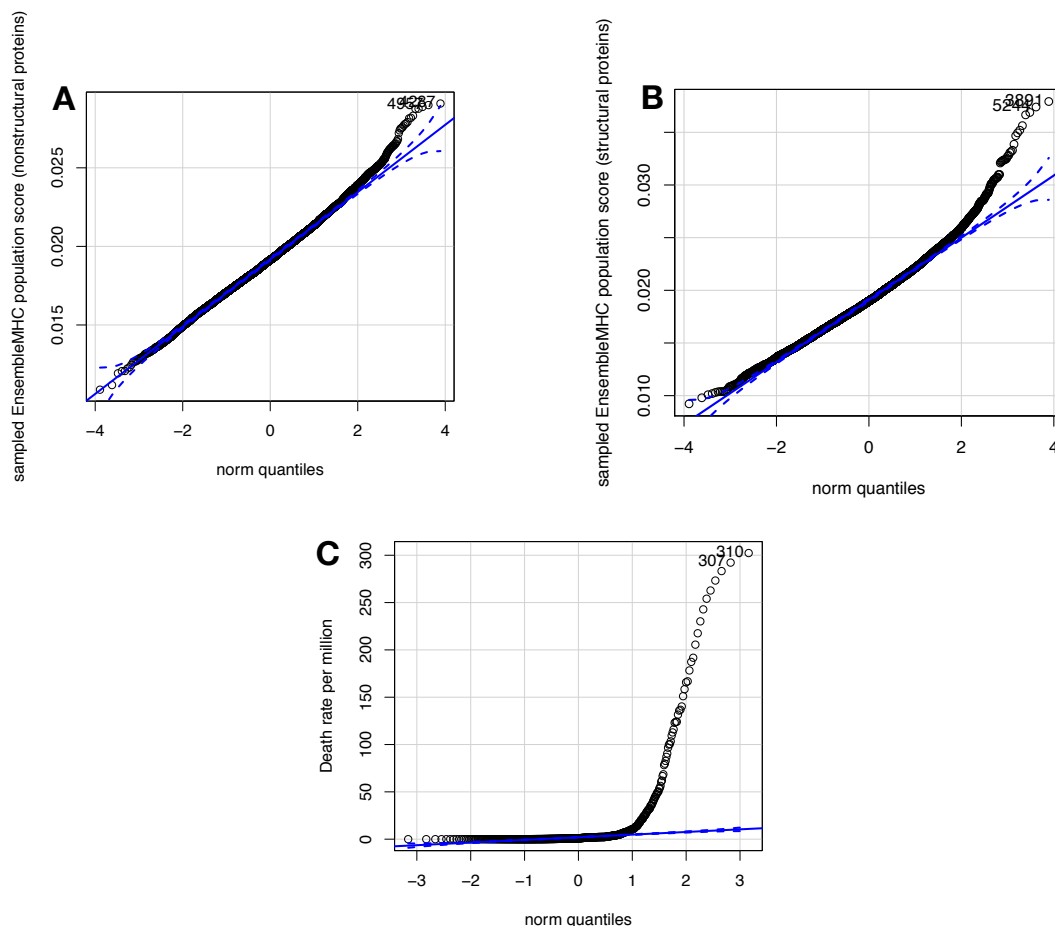


Figure A.24: **Justification of non-parametric correlation analysis.** The use of non-parametric correlation analysis, namely spearman's rho, is justified by the non-normality of the underlying data. EnsembleMHC population scores based on the full SARS-CoV-2 proteome and SARS-CoV-2 structural proteins were calculated for 10,000 simulated countries. Allele frequencies for simulated countries were generated by randomly sampling an observed allele frequency for each of the 52 alleles and re-normalizing to ensure the sum of allele frequencies were equal to one. **A**, The Q-Q plot for the simulated EnsembleMHC population score distribution based on the full SARS-CoV-2 proteome. **B**, The Q-Q plot for the simulated EnsembleMHC population score distribution based on SARS-CoV-2 structural proteins. **C**, The Q-Q plot for all reported deaths per Million. All three quantities show a considerable level of positive skewing, indicating non-normality.

population EnsembleMHC score and death rate correlation

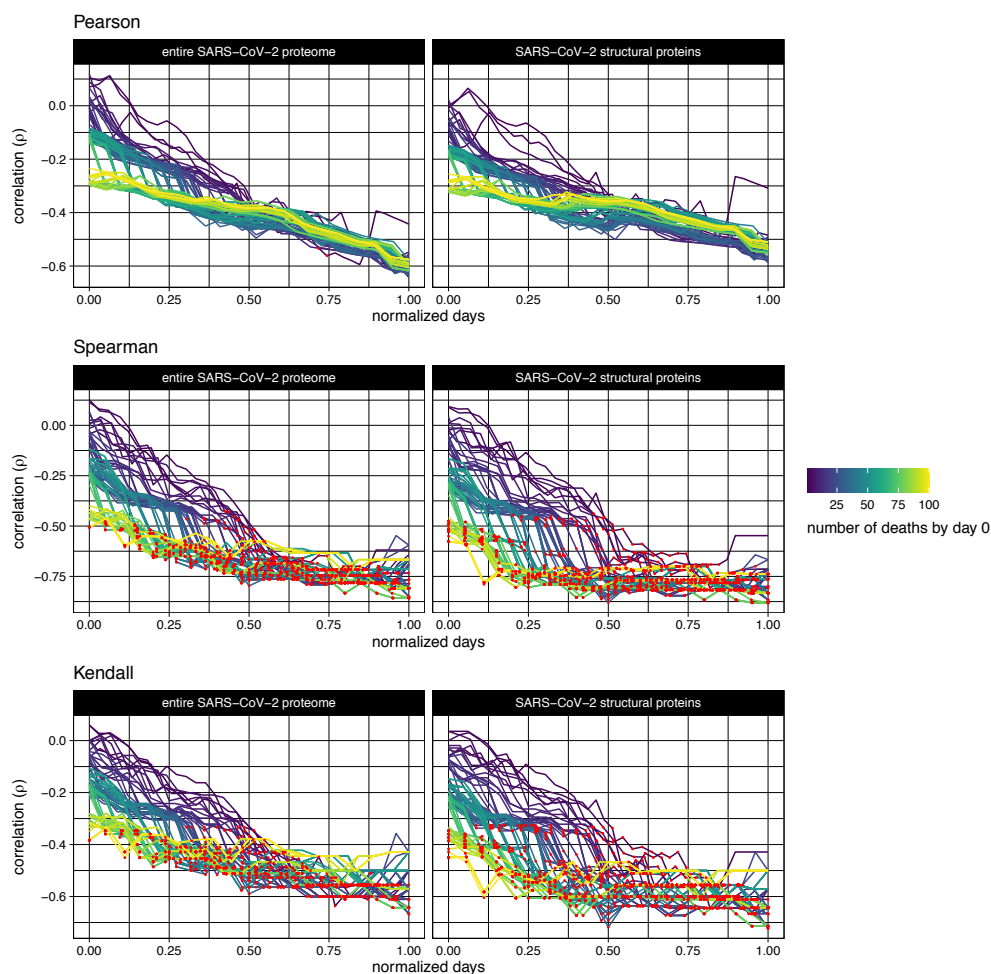


Figure A.25: **The effect of different correlation methods on the relationship between EnsembleMHC score and deaths per million.** The correlation between EnsembleMHC population score with respect to all SARS-CoV-2 proteins (**left column**) or SARS-CoV-2 structural proteins (**right columns**) and deaths per million using Pearson's r (**top**), Spearman's rho (**middle**), and Kendall's tau (**bottom**). Correlations that were shown to be statistically significant are colored with a red point.

B tables

Figure B.1: **MHC-I peptide identified by Ensemble MHC**. All identified peptides with a $peptide^{FDR} \leq 0.05$ (not pictured due to size).

A		B				
Countries		Minimum death threshold	Normalized day: 0.25	Normalized day: 0.5	Normalized day: 0.75	Normalized day: 1
China		5	8	15	23	30
Japan		10	7	14	20	27
South Korea		15	7	13	20	26
Taiwan		20	7	13	20	26
US		25	7	12	18	24
Hong Kong		30	7	12	17	23
France		35	7	12	17	23
Germany		40	7	12	17	23
India		45	6	11	17	22
Italy		50	6	11	17	22
Russia		55	6	11	17	22
UK		60	6	11	16	21
Iran		65	6	11	16	21
Israel		70	6	11	16	21
Croatia		75	6	11	16	21
Romania		80	6	11	15	20
Netherlands		85	6	11	15	20
Mexico		90	6	11	15	20
Ireland		95	5	10	14	19
Czechia		95	5	10	14	19
Morocco		100	5	10	14	19

non-normalized days

Figure B.2: Description of the additional socioeconomic factors selected for analysis. A, The 23 countries for which SARS-CoV-2 population binding capacities were calculated. B, The mapping of normalized days to real days for normalized day quartiles (0.25, , 0.5, 0.75, 1) at select minimum death thresholds.

Figure B.3: **EMP score correlation data.** All correlation data pertaining to the correlations between EMP score and deaths per million. This includes rho estimate, 95% CI, non-normalized days, and sample size for each correlation (not pictured due to size).

A		B		
Countries		Factor	Abbreviation	Description
China		% of population ≥ 65 years	65	Percentage of the population that is 65 years of age or older (2020).
Japan		Average BMI	Avg. BMI	The age-standardized average population body mass index (2016).
South Korea		Cardiovascular disease	CD	The deaths per million due to cardiovascular disease (2016).
US		Chronic obstructive pulmonary disease	COPD	The deaths per million due to complications from chronic obstructive pulmonary disease (2016).
France		Diabetes mellitus	DM	The deaths per million due to complications from diabetes mellitus (2016).
Germany		High blood pressure	BP	The age-standardized percentage of the population with a systolic blood pressure ≥ 140 or diastolic blood pressure ≥ 90 (2015).
India		Obesity prevalence	OBS	The age-standardized percentage of the population with a BMI ≥ 30 (2016).
Italy		Overweight prevalence	OVW	The age-standardized percentage of the population with a BMI ≥ 25 (2016).
Russia		Structural protein EMP score	SP	The SARS-CoV-2 structural protein presentation score.
UK		% of GDP spent on health care	GDP	Current health expenditure (CHE) as percentage of gross domestic product (2017).
Iran		% of total gov. expenditure on health care	GGHE	General government expenditure on health as a percentage of total (2014).
Israel		% of population that is female	SEX	The proportion of the total population that is female (2020).
Croatia				
Romania				
Netherlands				
Mexico				
Ireland				
Czechia				
Morocco				

Figure B.4: **Socioeconomic and health-related risk factors.** **A**, 21 countries were selected for analysis based on the existence of data in the Global Health Observatory data repository and inclusion in the 23 country set used for EMP score analysis. **B**, Descriptions and abbreviations for the selected risk factors. Each factor is labeled with the year that the data was collected. In every case, the most recent data was selected for analysis.