

1 HLA-dependent variation in SARS-CoV-2 CD8⁺ T cell cross- 2 reactivity with human coronaviruses

3
4 Paul R. Buckley^{1,2}, Chloe H. Lee^{1,2}, Mariana Pereira Pinho¹, Rosana Ottakandathil Babu^{1,2},
5 Jeongmin Woo^{1,2}, Agne Antanaviciute^{1,2}, Alison Simmons¹, Graham Ogg¹, Hashem
6 Koohy^{1,2,+}

7
8 ¹ MRC Human Immunology Unit, Medical Research Council (MRC) Human Immunology
9 Unit, MRC Weatherall Institute of Molecular Medicine (WIMM), John Radcliffe Hospital,
10 University of Oxford, Oxford, United Kingdom.

11
12 ² MRC WIMM Centre for Computational Biology, Medical Research Council (MRC)
13 Weatherall Institute of Molecular Medicine, John Radcliffe Hospital, University of Oxford,
14 Oxford, United Kingdom.

15
16 ⁺ Correspondence: hashem.koohy@rdm.ox.ac.uk

17 18 Abstract

19
20 Pre-existing T cell immunity to SARS-CoV-2 in individuals without prior exposure to SARS-
21 CoV-2 has been reported in several studies. While emerging evidence hints toward prior
22 exposure to common-cold human coronaviruses (HCoV), the extent of- and conditions for-
23 cross-protective immunity between SARS-CoV-2 and HCoVs remain open. Here, by
24 leveraging a comprehensive pool of publicly available functionally evaluated SARS-CoV-2
25 peptides, we report 126 immunogenic SARS-CoV-2 peptides with high sequence similarity to
26 285 MHC-presented target peptides from at least one of four HCoV, thus providing a map
27 describing the landscape of SARS-CoV-2 shared and private immunogenic peptides with
28 functionally validated T cell responses. Using this map, we show that while SARS-CoV-2
29 immunogenic peptides in general exhibit higher level of dissimilarity to both self-proteome
30 and -microbiomes, there exist several SARS-CoV-2 immunogenic peptides with high
31 similarity to various human protein coding genes, some of which have been reported to have
32 elevated expression in severe COVID-19 patients. We then combine our map with a SARS-
33 CoV-2-specific TCR repertoire data from COVID-19 patients and healthy controls and show
34 that whereas the public repertoire for the majority of convalescent patients are dominated by
35 TCRs cognate to private SARS-CoV-2 peptides, for a subset of patients, more than 50% of
36 their public repertoires that show reactivity to SARS-CoV-2, consist of TCRs cognate to shared
37 SARS-CoV-2-HCoV peptides. Further analyses suggest that the skewed distribution of TCRs
38 cognate to shared and private peptides in COVID-19 patients is likely to be HLA-dependent.
39 Finally, by utilising the global prevalence of HLA alleles, we provide 10 peptides with known
40 cognate TCRs that are conserved across SARS-CoV-2 and multiple human coronaviruses and
41 are predicted to be recognised by a high proportion of the global population. Overall, our work
42 indicates the potential for HCoV-SARS-CoV-2 reactive CD8⁺ T cells, which is likely
43 dependent on differences in HLA-coding genes among individuals. These findings may have
44 important implications for COVID-19 heterogeneity and vaccine-induced immune responses
45 as well as robustness of immunity to SARS-CoV-2 and its variants.

46 Introduction

47

48 After more than a year, the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)
49 pandemic remains a global health challenge and causes huge economic burden. SARS-CoV-2
50 virus gives rise to COVID-19 disease, which is characterised by a heterogenous clinical
51 outcome ranging from asymptomatic infection to severe acute respiratory distress and death.
52 The virus has proven to be dynamic, and the emergence of ‘variants of concern’ (e.g., the delta
53 variant) challenges the existing mitigation strategies including vaccine rollouts¹.

54
55 Although disease morbidity is associated with several factors including age, sex and aberrant
56 immune response; the mechanisms and factors underpinning the heterogeneity of disease are
57 incompletely understood². Furthermore, reports of differential immune responses following
58 vaccination have started to emerge, demonstrating prior SARS-CoV-2 infection can enhance
59 COVID-19 vaccine response compared with naïve individuals^{3,4}. Indeed, many questions
60 regarding the magnitude and robustness of immune response in disease, variants of concern
61 and/or COVID-19 vaccination in different individuals remain open despite the great recent
62 efforts.

63
64 Several studies^{5,6} have illustrated that the correlates of immunity to SARS-CoV-2 are
65 implicated by the presence of pre-existing immunological memory conferred from cross-
66 reactivity to other viruses. On the one hand such cross-reactivity could modulate disease
67 severity, vaccine response and/or protection against SARS-CoV-2 and its variants via presence
68 of antigen-specific memory T cells⁷. Conversely, cross-reactivity may provoke
69 immunopathology through mechanisms such as antibody-dependent enhancement of infection,
70 with the potential for virus-induced autoimmune disease in years to come⁸.

71
72 Coronavirus strains that infect humans belong to either alpha or beta genera. The
73 alphacoronaviruses contain HCoV-229E and -NL63 while the four lineages of
74 betacoronaviruses include HCoV -OC43 and -HKU1, SARS-CoV and -CoV-2, MERS-CoV
75 and other viruses only identified in bats. HCoV-OC43, -HKU1, -NL63 and -229E strains are
76 known to cause mild to moderate ‘common cold’ symptoms whereas MERS-, SARS-CoV-1
77 and -2 can cause severe respiratory tract disease and death. Previous natural and experimental
78 infection studies in humans suggest antibody cross-reactivity within- but minimal reactivity
79 between- endemic human alpha and beta coronaviruses. Unlike antibodies, T cell cross-
80 reactivity to SARS-CoV-2 appears to be more prevalent. Several recent studies have reported
81 existence of SARS-CoV-2-specific T cells in unexposed individuals⁹⁻¹⁵, although it appears
82 that T cell cross reactivity is more pronounced in CD4⁺ than CD8⁺ T cells in these subjects.

83
84 Recent studies have provided varying insights regarding the presence of pre-existing CD8⁺ T
85 cell immunity to SARS-CoV-2 conferred by HCoV. In an investigation into the
86 immunodominant SARS-CoV-2 SPR* epitope – associated with HLA-B*07:02 – Nguyen et
87 al¹⁶., found little evidence of cross-reactive exposure in pre-pandemic Australian samples. On
88 the other-hand, Francis et al¹⁵., found evidence of pre-existing memory CD8⁺ T cells in naïve
89 samples and have shown that HLA genotype conditions pre-existing CD8⁺ T cell memory to
90 SARS-CoV-2, and they suggest that unexposed individuals with specific HLA alleles (such as
91 HLA-B*07:02), may be more likely to possess cross-reactive memory T cells specific for the
92 SPR* SARS-CoV-2 epitope. These disparate results may stem from differences in regional
93 HLA allele frequencies and / or experimental methodology. Nevertheless, the extent to which
94 patients’ haplotypes and SARS-CoV-2-HCoV cross-reactivity - amongst other factors - are
95 linked to heterogeneous COVID-19 disease, robustness of immunity against SARS-CoV-2 and
96 its variants, and/or protection after vaccine-induced immune response, remains to be
97 elucidated.

98

99 In this study, we examined the evidence for SARS-CoV-2-specific T cell cross-reactivity with
100 common-cold HCoVs and identifying 126 immunogenic SARS-CoV-2 peptides that are highly
101 similar to 285 predicted HCoV pMHC. We additionally identified a set of SARS-CoV-2
102 peptides with high similarity to several human genes. We found that public TCR repertoires
103 reactive to SARS-CoV-2 in COVID-19 patients who carry specific HLA alleles primarily
104 recognise SARS-CoV-2 peptides with high similarity to HCoVs, suggesting that common-cold
105 HCoV cross-reactivity is variable and likely to be conditioned by HLA. It is plausible that
106 patients carrying these HLAs may exhibit more robust protection against SARS-CoV-2 and its
107 variants. We lastly identified a set of 10 peptides that are highly conserved across multiple
108 coronavirus strains, to serve not only as potential pan-coronavirus T cell targets, but we propose
109 are leading candidates as cross-reactive CD8⁺ T cell epitopes.

110

111 Results

112

113 Curation of functionally evaluated SARS-CoV-2 peptides

114

115 To investigate the potential for T cell cross-reactivity against SARS-CoV-2 as conferred by
116 common-cold HCoVs, we curated a comprehensive pool of SARS-CoV-2 class I and II
117 peptides from three datasets (see Methods), which have been functionally evaluated for CD4⁺
118 and CD8⁺ T cell responses (see Figure 1 for study overview). The data comprise 1799 and
119 1005 immunogenic and non-immunogenic SARS-CoV-2 *peptides* respectively (Fig 2A). Many
120 of these peptides were tested for T cell reactivity in the context of multiple HLA alleles and/or
121 by multiple assays (IFN γ , IL-5 production etc). Furthermore, some peptides are described by
122 qualitative labels corresponding to varying response magnitude (Positive-high and Positive-
123 low etc). Taking these combinations into account, we found 3979 and 2427 immunogenic and
124 non-immunogenic *complexes* (Fig 2B). For immunogenic complexes, the most common
125 lengths are 9-mers, followed by 15- and 10-mers (Fig 2C), and 36.0% are presented by class I
126 MHC, 32.9% by class II (Fig 2D) and for 31% MHC type is unknown (Fig S1A). For non-
127 immunogenic complexes, 36.1% are presented by class I, 26.4% by class II and for 37.51% the
128 MHC is unknown. At the gene level, HLA-allele specific information was available for 934
129 (56.5%) and 607 (42.2%) of immunogenic class I and II complexes respectively (Fig S1A).

130

131 Given the high proportion of missing MHC information, we employed netMHCpan 4.1 and
132 netMHCIIpan to predict presenting class I and class II alleles respectively for immunogenic
133 peptides (see Methods). Here, we were able to identify 98% of known MHC molecules,
134 providing confidence in predictions for unknown alleles (Fig S1B).

135

136 We next sought to examine whether HLAs exhibit preferences towards presenting peptides
137 from certain SARS-CoV-2 proteins. By employing a similar methodology to Karnaukhov et
138 al¹⁷, we gauged the enrichment and depletion of HLA ligands arising from these proteins (see
139 Methods). Indeed, we observed differential antigen presentation by HLAs e.g., HLA-C*07:02
140 appears to be the most consistently enriched in presenting 9mers from the examined proteins
141 (Fig 2E), while HLA-A*02:01 is enriched in presenting 9mers from ORFs but depleted for
142 10mers across most assessed proteins. This disparity may be due to a known preference of 9-
143 mers for HLA-A*02:01¹⁸. Furthermore, despite the prevalence of HLA A*02:01 in the global
144 population and in MHC presentation experiments, this allele appears to be depleted for
145 presenting ligands from SARS-CoV-2 proteins that have been the focus of intense experimental
146 work e.g., spike and nucleocapsid phosphoprotein.

147

148 These patterns of HLA preferences in presenting SARS-CoV-2 peptides appear to differ for 9-
149 and 10-mers. For example, whereas HLA-C*07:02 is enriched for presenting 9mers, this allele
150 appears to be a poor presenter of 10mers from each examined protein. It is unclear why
151 substantially fewer 10-mer HLA-C*07:02 ligands are predicted than 9mers, however it is
152 plausible that this allele may prefer 9mers, as appears to be the case with HLA-A*02:01, -
153 A*11:01 and -B*40:01¹⁸, or that this may be a SARS-CoV-2 specific effect.

154
155 Although it is of great interest to reveal the rate to which SARS-CoV-2 MHC-bound peptides
156 are immunogenic in humans¹⁹, it cannot be examined directly with existing data because not
157 all MHC-bound SARS-CoV-2 peptides have been evaluated for immunogenicity.
158 Nevertheless, we explored the pool of MHC-bound peptides in our dataset that have been
159 examined for a T cell response, to gauge the proportion that SARS-CoV-2 pMHC are
160 immunogenic. Overall, we observed low rates of immunogenic pMHC (Fig 2F), although
161 ligands of HLA-B*40:01 appear to be commonly immunogenic. Interestingly we observed that
162 HLA-C*07:02 does not present any 10-mers in our dataset. This apparent preference for 9-
163 mers is consistent with availability of HLA-C*07:02 ligands tested for T cell response in
164 humans from the IEDB, where there exist only 121 unique peptides, of which 73% are 9mers
165 and only 12% are 10mers. In summary, these data suggest length and source protein
166 preferences for HLA alleles presenting SARS-CoV-2 peptides and that HLA-B*40:01 SARS-
167 CoV-2 ligands are commonly immunogenic.

168

169 Identification of Shared and Private Immunogenic SARS-CoV-2 peptides

170

171 To discriminate SARS-CoV-2-HCoV shared (hereby referred to as ‘sCoV-2-HCoV’) peptides,
172 we compared immunogenic SARS-CoV-2 peptides to HCoV protein sequences. For this, we
173 define a metric that considers 1) sequence homology, 2) physicochemical similarities
174 (MatchScore²⁰) and 3) presentation status for which the source peptide from SARS-CoV-2 and
175 the target peptide from one of the HCOVs are required to be presented by the same HLA. A
176 source peptide is defined as shared if it fulfils all these three conditions otherwise is considered
177 as a private peptide (see Methods).

178

179 Using our metric, we identified 126 unique SARS-CoV-2 (immunogenic) peptides pointing to
180 285 highly similar peptides in HCOVs (Supplementary Data File 1). Hence, we provide a
181 comprehensive map of private and shared SARS-CoV-2 functionally evaluated immunogenic
182 peptides, and for sCoV-2-HCoV peptides, their matches from each HCoV.

183

184 Out of the HLAs tested (see Methods) 33 and 28 class I and II HLAs respectively were
185 predicted to present the target HCoV pMHCs (Fig3A). HLA-A*02:01 and HLA-B*27:05 were
186 the most and least common class I presenters respectively. For class II, DRB1-1501 and DRB5-
187 0101 were the most common presenters, while DRB1-0301 and DRB1-1303 were the least.
188 Most shared class I and II peptides were predicted to bind multiple HLA allelic variants (Fig
189 2SA). Compared with private peptides it appears that sCoV-2-HCoV peptides are presented by
190 less HLAs, although this was not significant (Fig S2C). Nevertheless, the range of predicted
191 alleles for these peptides suggests recognition in broad geographical and ethnic settings²¹.

192

193 For the 126 SARS-CoV-2 peptides with high similarity to HCoV, we also observed binding to
194 multiple HLAs (FigS2B). In addition, we found that 9mers comprise 54% of the 126 SARS-
195 CoV-2 peptides with high-similarity matches to HCoV, followed by 15mers (19%) and 10mers
196 (17.5%) (FigS2C). Consistent with previous reports²², the betacoronaviruses HKU1 and OC43
197 were most enriched in target matches (Fig3B), perhaps due to higher total sequence homology

198 among betacoronavirus strains²³. We next examined the extent to which immunogenic SARS-
199 CoV-2 peptides exhibit homology to *multiple* HCoV strains. Surprisingly, we found that 42
200 SARS-CoV-2 immunogenic peptides exhibit matches to at least three strains (Fig 3C).
201 However, we observed small clusters of peptides that only possess homology with one strain,
202 e.g. OC43 or HKU1. ORF1ab protein and spike surface glycoprotein produced the highest
203 quantity of shared SARS-CoV-2-HCoV peptides in both strains, and the protein regions from
204 which these peptides were found are similar in both HKU1 and OC43 (Fig S2E-H).

205

206 Of particular note about our map of shared and private peptides is that this map is subject to
207 thresholds that we used in our metric. The sequence homology threshold that was used here is
208 50% and most peptides had greater than or equal to 70% sequence homology (Fig S2E).
209 Although, more stringent sequence homology parameter will result a map containing fewer
210 shared peptides (Fig 3D), our main conclusions in this manuscript remain the same even with
211 sequence homology threshold of 70% (data are not shown).

212

213 Lastly, we compared the amino acid distribution between shared and private SARS-CoV-2
214 peptides for 9-mers, which is the most common peptide length in our dataset (Fig 3E). We
215 observed some moderate differences, e.g., increased prominence of Valine at position 9 within
216 shared peptides.

217

218 We have therefore identified a pool of 126 SARS-CoV-2 immunogenic peptides - that exhibit
219 high similarity to 291 peptides in HCoV strains - which are likely to be presented by an array
220 of class I and II HLA molecules. This array of presenting alleles suggests the potential for
221 broad global population coverage, which is explored later. We propose that this pool of
222 experimentally confirmed immunogenic SARS-CoV-2 peptides and their counterpart high
223 similarity matches be considered as potential targets for T cell cross-reactivity, therefore
224 warranting investigation into pre-existing immune memory from HCoV or a role in protection
225 from SARS-CoV-2 variants.

226

227 Identification of peptides with high similarity to self and self-microbiomes

228

229 To prevent aberrant T cell mediated inflammation and tissue damage, the immune system has
230 evolved several checkpoint mechanisms. These include thymus negative selection and
231 peripheral tolerance. Indeed, dissimilarity to self is increasingly recognised as a component of
232 peptide immunogenicity²⁵, which may assist in calibrating a balance between immunogenicity
233 and inflammatory pathogenesis.

234

235 To evaluate the extent to which dissimilarity to self and self-microbiomes contribute to SARS-
236 CoV-2 peptide immunogenicity, we took a similar approach and used our metric to compare
237 SARS-CoV-2 peptides to human self-proteome and microbiomes that include 457 gut and 50
238 airway microbiota. (see Methods). Here, for SARS-CoV-2 HLA class I presented 9- and 10-
239 mer peptides we observed that immunogenic SARS-CoV-2 peptides were significantly more
240 dissimilar to the human proteome than their non-immunogenic counterparts (Fig 4A, S3A).
241 Using this approach, we could not detect any significant difference between immunogenic and
242 none immunogenic class II peptides in their dissimilarity to self-proteome (Fig S3B).

243

244 Interestingly however, for peptides of both lengths 9 and 10, we identified several
245 immunogenic SARS-CoV-2 peptides with considerable sequence similarity to the human
246 proteome (Fig 4A-B, Table S1). For the top 10% of these peptides with highest similarity to
247 self, the mean amino acid conservation (the proportion of the amino acid sequence which is

248 exactly conserved) between these peptides and corresponding self matches is 72.1% with
249 8.33% standard deviation (see Supplementary Data File 2 for number of substitutions under
250 column ‘Hamming’). In general, T cells specific for these peptides should be subject to
251 negative selection otherwise it is plausible that aberrant immune responses may occur during
252 the course of the disease in the form of immunopathology or in the future in the form of
253 autoimmunity^{8,26,27}.

254
255 To investigate the potential association of these peptides in immunopathology further, we
256 predicted MHC presentation by a set of class I HLA alleles (see Methods) for the top 10% of
257 peptides most similar to the human proteome for 9mers and 10mers. We observed that these
258 peptides with high similarity to self are predicted to bind multiple HLAs (Fig 4C), and
259 interestingly, we found that in most cases, the SARS-CoV-2 immunogenic peptide and the
260 match from the human proteome are predicted to be presented by the same allele (Fig4C).

261
262 Next, we examined the list of genes with high sequence similarities to these SARS-CoV-2
263 immunogenic peptides (Table S1 & Supplementary Data File 2). Of particular interest, we
264 found e.g. CCL3 and CCL3L1 which are linked to cytokine storms and the expression of which
265 have been reported to be elevated in severe COVID-19 patients^{28–33} (Fig 4D, Supplementary
266 Data File: 1). We additionally observed CD163, similarly associated with severe COVID-19,
267 however the predicted presentation score of HLA-B*15:01 for peptides from CD163 with high
268 similarity to SARS-CoV-2 were slightly beyond the generally accepted ‘binding’ cutoff.
269 Interestingly, the SARS-CoV-2 peptides exhibiting sequence homology to CCL3 and CCL3L1
270 (and CD163) were private to SARS-CoV-2 (Fig 4D & Table 1) - which may increase the
271 likelihood of being involved in immunopathology after infection. Additionally, we observed
272 considerable amino acid conservation with matches from these genes, with 77.8% for 9mers
273 and 70% for 10mers (Table 1).

274
275 CCL3 and CCL3L1 are both ligands for CCR1 and CCR5. Interestingly, CCR1 variants are
276 linked to pulmonary macrophage infiltration in severe COVID-19³⁴ and inhibition of CCR5 in
277 critical COVID-19 patients has been associated with a decrease in plasma IL-6 and SARS-
278 CoV-2 RNA and an increase in CD8+ T cells³⁵. Additionally, intermediate monocytes which
279 constitutively express high levels of CCR5 have recently been suggested as playing a role in
280 post-acute sequelae of COVID-19³⁶ (often referred to as ‘long-COVID’). Of further interest,
281 we found SMPD4 and SLC1A4, which together with CCL3 and CCL3L1 are involved in the
282 response to TNF, which is part of the cytokine storm following COVID-19 disease.

283
284 By comparing SARS-CoV-2 peptides to human microbiomes, we observed subtle higher
285 dissimilarity of SARS-CoV-2 immunogenic peptides to the gut (Fig S3C) and airways (Fig
286 S3D) microbiomes, which may suggest a link between the diversity of both microbiota and
287 heterogeneity of the disease in populations, although this warrants further investigation.

288
289 Given the magnitude of the global pandemic and the widespread vaccination required to
290 combat it, future virus-induced autoimmune disease and immunopathology is of concern.
291 Overall, this analysis suggests dissimilarity of viral peptides to self-proteins as a correlate of
292 peptide immunogenicity. Furthermore, we present candidate genes and peptides with high
293 similarity to SARS-CoV-2 T cell targets, which we suggest as prime targets for further
294 investigations into their role in autoimmune disease and immunopathology following SARS-
295 CoV-2 infection and/or vaccination.

296

297 CD8+ T cell cross-reactivity and common-specificity within SARS-CoV-2

298

299 A valuable characteristic of our map of SARS-CoV-2 shared and private peptides, is that for
300 245 of these (out of 1279 class I immunogenic peptides), cognate TCRs at the beta chain
301 resolution are available in the IEDB. We therefore set out to map the TCR landscape through
302 a network approach to explore the potential for cross-reactivity among SARS-CoV-2 specific
303 CD8+ T cells, and their common-specificity. Here, to avoid overestimating connectivity, any
304 peptides of different lengths, which share starting positions in the SARS-CoV-2 proteome and
305 are recognised by identical sets of TCRs, are considered as one peptide.

306

307 Through a two-mode (bipartite) network-graph illustrating the connectivity of SARS-CoV-2
308 immunogenic peptides with their cognate TCRs, amongst a highly connected topography we
309 observed considerable connectivity for some sCoV-2-HCoV peptides e.g. “FLN..” (Fig S4A).
310 Exploring this further, we projected the bipartite network-graph into a one-mode graph where
311 nodes represent peptides and an edge between two nodes requires existence of a TCR
312 recognising both peptides (Fig S4B). The clustering around a small set of hubs suggests that
313 many experimentally assessed TCRs target a small set of SARS-CoV-2 peptides. Indeed, we
314 found that in this dataset, 80% of the TCRs are reported to recognise only 40 (16%) peptides,
315 of which 4 are sCoV-2-HCoV shared peptides and 36 are SARS-CoV-2 private (Fig S4C). This
316 dominant set of peptides may be due to experimental biases e.g., research may be heavily
317 biased toward several protein regions. However, this may also reflect a selection bias by SARS-
318 CoV-2 specific TCRs. In this regard, amongst these 80% of TCRs, we observed high usage of
319 V gene TRBV20-1³⁷ and J gene TRBJ2-1³⁸ (Fig S4D), that have been previously reported to
320 have implications in COVID-19 patients.

321

322 Similarly, we examined the extent of common specificity in SARS-CoV-2 specific T cells by
323 a one-mode graph in which nodes represent TCRs and an edge represents whether two nodes
324 (TCRs) recognize the same peptide (Fig S4E). Interestingly, this graph reveals a set of highly
325 connected hubs reflecting levels of common specificity, however there are many TCRs which
326 recognise only a single unique peptide. Comparing these two sets of TCRs, we did not observe
327 considerable differences in their CDR3 β sequences (Fig S4F-G), however we observed
328 differences in V and V-J gene usage (Fig S4H-J).

329

330 In summary, we employed peptides with known cognate TCRs in the IEDB database - although
331 limited in numbers – to explore SARS-CoV-2 CD8+ T cell cross-reactivity. Our network
332 approach demonstrates that SARS-CoV-2 CD8+ T cells can cross-react and exhibit common-
333 specificities.

334

335 Presence of public TCRs recognising sCoV-2-HCoV peptides in COVID-19 convalescents and 336 healthy subjects

337

338 We next integrated our map of SARS-CoV-2 shared and private peptides with a recently
339 published dataset known as ‘MIRA’³⁹ to track the patterns of public TCRs recognizing sCoV-
340 2-HCoV peptides in convalescents and/or healthy subjects. Here, Nolan et al., employed the
341 ‘Multiplex Identification of Antigen-Specific T cell receptors’ (MIRA) assay to identify
342 SARS-CoV-2 specific TCRs from PBMCs and naïve T cells. These data include more than
343 160k high confidence SARS-CoV-2-specific TCRs mapped to target peptides from 39 Healthy
344 controls (HC) (defined as unexposed to SARS-CoV-2) and 90 COVID-19 convalescent
345 patients. These data consist of 792 unique SARS-CoV-2 peptides, 54 of which are sCoV-2-
346 HCoV shared peptides.

347
348 To elucidate the landscape of public TCRs in HC and convalescent patients, we generated a
349 bipartite graph comprising all public TCRs (defined as CDR3b+V +J gene(s) present in at least
350 two subjects) cognate for SARS-CoV-2 private and sCoV-2-HCoV shared peptides (Fig 5A).
351 This graph revealed two clear hubs. In the first (green nodes), we observed that healthy subjects
352 were connected to public TCRs which recognise both sCoV-2-HCoV and SARS-CoV-2-
353 private peptides. In the second hub (red nodes) comprising convalescent patients, we observed
354 that generally their public TCR repertoires predominately recognise SARS-CoV-2-private
355 peptides. Indeed, it appears that cognate TCRs of sCoV-2-HCoV peptides are more pronounced
356 in HC (Fig S5A-Shared, wilcoxon $p=0.00029$) whereas cognate TCRs of SARS-CoV-2-private
357 peptides appear enriched in the convalescent cluster (Fig S5A-Private). Interestingly, we
358 observed a considerable number of TCRs recognising shared peptides which are common
359 between these two subject clusters, indicating that sCoV-2-HCoV-specific public TCRs are
360 present not only in COVID-19 patients but are also expanded from unexposed individuals (Fig
361 5A-B).

362
363 Given that in these healthy donors, the TCRs are generally from naïve CD8+ T cells which are
364 expanded and stimulated with SARS-CoV-2 peptide pools and analysed with the ‘MIRA’
365 assay, the presence of cognate TCRs recognising sCoV-2-HCoV peptides in HC as well as
366 COVID-19 patients may not necessarily translate into pre-existing T cell immunity. Rather, due
367 to the high similarity between the cognate SARS-CoV-2 antigens and (predicted) HCoV
368 presented peptides, we suggest it is plausible that these SARS-CoV-2 specific TCRs are cross-
369 reactive with HCoV peptides. Indeed, consistent with Francis et al.,¹⁵ who demonstrate pre-
370 existing memory CD8+ T cells to SPR* peptide in 80% of unexposed individuals, we found a
371 set of public TCRs – which are observed in both convalescent and unexposed individuals -
372 recognising this sCoV-2-HCoV peptide. In this light, we reveal candidate public TCRs and
373 corresponding SARS-CoV-2 peptides with high similarity to HCoVs, which should be
374 examined further for cross-reactive potential.

375
376 From these two bipartite graphs, we observed that healthy individuals respond to a balance of
377 SARS-CoV-2 private and sCoV-2-HCoV peptides, although it appears that infection *primarily*
378 dictates a dominant recognition of private SARS-CoV-2 peptides (Fig S5B). For convalescent
379 patients, we observed that public TCR repertoires of the majority (51/86) of patients are almost
380 entirely ($\geq 99\%$) occupied by TCRs recognising SARS-CoV-2 private peptides (Fig 5C).
381 However, in a subset of convalescent patients, public TCRs recognising sCoV-2-HCoV
382 peptides comprise a substantial fraction of the public repertoire. In fact, for 12 convalescent
383 patients, $>50\%$ of their public TCRs recognise sCoV-2-HCoV peptides.

384
385 Comparing these two groups of patients, we did not find evidence of a link toward biological
386 sex or age. To explore potential correlates, we first gathered the 12 patients whose public TCRs
387 most dominantly ($>50\%$) recognise sCoV-2-HCoV peptides (labelled *PubTCR-SharedEp*), and
388 then via sampling 12 patients 10 times from the set of 51 patients whose public TCRs almost
389 entirely recognise SARS-CoV-2 private peptides (labelled *PubTCR-Private*), we compared
390 HLA coding genes of these two groups. We observed that the *PubTCR-SharedEp* group is
391 statistically enriched for carrying HLA-B*07:02, HLA-C*07:02 and HLA A*03:01, whereas
392 the former group includes a broader set of HLAs among which HLA A*01:01 was more
393 pronounced (Fig 5D). The enrichment of HLA-B*07:02 in the *PubTCR-SharedEp* group is
394 consistent with recent work from Francis et al.¹⁵, and these data are in agreement with their
395 claim that CD8+ T Cell HCoV-CoV-2 cross-reactivity may be conditioned by HLA.

396

397 Employing these two groups and sampling a set of healthy patients (n=12), we reveal the set
398 of epitopes only recognised by public TCRs in these healthy patients, and those shared with
399 the convalescent *PubTCR-SharedEp* group (Fig S5C, Supplementary Data File 3 and 4).
400 Additionally, we reveal the peptides only observed in the *PubTCR-Private* convalescent group,
401 adding to previous insights that SARS-CoV-2 infection can provoke T cell responses to a novel
402 set of peptides compared to those expanded from unexposed patients⁷.

403
404 Recent work shows cross-reactive *private* TCRs from unexposed subject repertoires, capable
405 of recognising both the SARS-CoV-2 SPR* peptide and its LPR* homolog from HCoV OC43
406 and HKU1. By mapping out which SARS-CoV-2 peptides are recognised in which individuals
407 by private TCRs, we observed SPR* but also an additional set of sCoV-2-HCoV peptides
408 recognised in both healthy and convalescent patients (Fig S5D, Supplementary Data Files 5-
409 6). Lineburg et al.,⁴⁰ recently reported private TCRs in HLA-B*07:02⁺ unexposed individuals
410 which cross-react with both the SARS-CoV-2 SPR* peptide and the OC43/HKU1 homolog
411 LPR*, which indicates a level of pre-existing immunity. Of these TCRs, we found two (defined
412 as CDR3b, TRBV, TRBJ) which appear in two HLA-B*07:02⁺ unexposed individuals within
413 the MIRA dataset (Table S2). As these TCRs are now observed in two separate datasets, we
414 therefore propose these as public TCRs, capable – as identified by Lineburg et al., - of cross-
415 reacting with both SARS-CoV-2 SPR* and OC43/HKU1 LPR* peptides.

416
417 Taken together, we report existence of a set of CD8⁺ TCRs in both HC and COVID-19
418 convalescent patients that recognise SARS-CoV-2 peptides with high sequence similarity to a
419 pool of predicted HCoV pMHC. This high sequence similarity indicates cross-reactive
420 potential of these TCRs. Primarily however, we observed that COVID-19 patients develop
421 public TCR responses to private peptides – many of which are not observed in unexposed
422 individuals - indicating that any cross-reactive potential is limited. For the subset of COVID-
423 19 patients whose public TCRs are directed towards sCoV-2-HCoV peptides - and are observed
424 in HC - we found distinct HLA profiles. Therefore, in agreement with recent data from Francis
425 et al., we suggest that CD8⁺ T cell HCoV-CoV-2 cross-reactive potential is apparent, although
426 likely conditioned by patient HLA genotype. It is plausible that these patients may exhibit more
427 robust protection against SARS-CoV-2 and its variants.

428
429 Potential conserved coronavirus CD8⁺ T cell targets with broad population coverage
430

431 Given the emergence of new SARS-CoV-2 variants and concern over the theoretical capacity
432 of future mutants to evade current vaccine strategies¹, conserved CD8⁺ T cell targets across
433 multiple coronavirus strains with the potential to elicit T cell responses in a large percentage
434 of global populations are of interest. We therefore searched our peptide map for SARS-CoV-2
435 peptides with ‘high-similarity’ matches to *multiple* HCoVs, and with cognate TCRs in the
436 MIRA dataset. To select only the top ‘high-similarity’ SARS-CoV-2-HCoV matches for this
437 analysis, we applied a more stringent sequence homology threshold. Indeed, in addition to the
438 ‘MatchScore’ and peptide presentation criteria outlined previously (see Methods:
439 Discriminating shared and private SARS-CoV-2 peptides), we only retained matches with at
440 least 70% sequence conservation (i.e. allowing 30% amino acid substitution).

441
442 We found 86 peptides that match these criteria, 84 of which are recognised by TCRs in both
443 convalescent and HC (Fig 6A-B). We next focused on SARS-CoV-2 peptides with high
444 similarity matches in ≥ 3 HCoV strains (Table 2, Supplementary Data File 7). Of these SARS-
445 CoV-2-HCoV matches, the number of amino acid substitutions ranged between 0-3, with a
446 mean of 1.79 and standard deviation of 0.78. Additionally, while each of these peptides

447 exhibited a high similarity match to either MERS or SARS-CoV, the majority exhibited
448 homology with both of these viruses (Fig S6A). As well as high conservation across many
449 coronavirus strains, collectively these SARS-CoV-2 peptides are predicted to bind multiple
450 HLA alleles (Fig 6C), raising the possibility that this set of peptides may elicit T cell responses
451 in a substantial proportion of the global population.

452
453 We next sought to determine the extent in global and regional populations that these CD8+ T
454 cell targets may elicit T cell responses individually and accumulatively. We therefore used the
455 IEDB population coverage tool⁴¹, which employs global HLA allele prevalence data to predict
456 the percentage of individuals in a regional population to respond to a given epitope set. Starting
457 with each SARS-CoV-2 peptide and predicted HLAs individually, we find considerable
458 coverage of 55.32% for “LLLD*”, while “VQID*” exhibits the lowest predicted coverage of
459 7.09% (Fig 6D).

460
461 Similarly to a previous approach by Ahmedid et al⁴², we set out to predict the accumulated
462 global population coverage of the set. We found that 8 peptides collectively produce >90%
463 global coverage, while the entire set is predicted to elicit T cell responses in 92.93% of the
464 global population (Fig 6E). Regionally, Europe and North America exhibited the highest
465 predicted coverage (Fig 6F). Of note, Africa and Asia also exhibited high predicted coverage.
466 Central America (defined as Guatemala and Costa Rica) exhibited low coverage of 7%. It is
467 unclear why, and further investigation is necessary to produce a peptide set with high coverage
468 in these countries.

469
470 Overall, we identified a set of 10 SARS-CoV-2 immunogenic peptides, each highly conserved
471 across coronavirus strains, which collectively provide global population coverage of ~93%.
472 We believe that this is an encouraging insight in the search for pan-coronavirus T cell targets,
473 and additionally propose these as top candidates for cross-protective immunity.

474

475 Discussion

476

477 Our work demonstrates that T cells specific to SARS-CoV-2 peptides with high similarity to
478 HCoV pMHC can be expanded from naïve individuals, and that these cognate public TCRs are
479 also observed in a subset of recovered COVID-19 patients. This finding firstly suggests that
480 SARS-CoV-2-unexposed individuals could mount T cell responses to HCoVs that – due to
481 peptide similarity - could be cross-reactive with SARS-CoV-2 antigens. Furthermore, we
482 propose that while COVID-19 disease appears to primarily direct responses against SARS-
483 CoV-2-private peptides, patients with certain HLA alleles (e.g HLA-B*07:02, -C*07:02, -
484 A*03:01) may be more likely to possess sCoV-2-HCoV cross-reactive CD8+ T cells. It is
485 therefore plausible that SARS-CoV-2 naïve individuals with certain HLAs may be at lower
486 risk of severe disease – or experience augmented vaccine responses - if previously exposed to
487 endemic coronaviruses, however a direct link to pre-existing immunity requires further
488 investigation.

489

490 Indeed, our analysis indicates that after SARS-CoV-2 infection, a subset of individuals has
491 memory T cells that primarily recognize sCoV-2-HCoV peptides. In these convalescent
492 patients, it is unclear whether infection itself, and/or prior exposure to HCoVs are driving this
493 subset of individuals to select for these peptides. There is conflicting evidence surrounding the
494 existence of memory SARS-CoV-2 cross-reactive CD8+ T cells in unexposed
495 individuals^{15,16,40}, and a limitation of our work is that we could not to provide a direct link to
496 pre-existing immunity, because from healthy donors the MIRA dataset only evaluated

497 expanded naïve T cells and did not examine anti-viral efficacy of the responding T cells.
498 Indeed, although we cannot determine the cause or timeframe of this selection of sCoV-2-
499 HCoV peptides in this subset of individuals, the potential implications are interesting. It is
500 plausible that these patients may exhibit more robust protection against SARS-CoV-2 variants,
501 HCoVs or even future emerging coronavirus strains. Future work should explore any immunity
502 benefit of infection-induced cross-reactive T cell responses, and in addition, it will be
503 interesting to examine whether vaccination against SARS-CoV-2 can induce T cell memory
504 that is cross-reactive with SARS-CoV-2 variants and/or wider coronaviruses in such
505 individuals. Furthermore, by our identification of a set of 10 potentially cross-reactive peptides
506 with broad population coverage, it is possible that these peptides could be employed to test
507 which patients exhibit cross-reactive phenotypes e.g., after vaccination with relevant antigens.
508

509 More broadly, data are beginning to demonstrate distinct vaccine-induced responses linked to
510 differential patient exposure to SARS-CoV-2^{3,4}. In turn, it is possible that COVID-19 vaccine
511 boosted cross-reactive immune responses may influence vaccine-induced protection⁷. Indeed,
512 it will be important to explore whether COVID-19 vaccination can boost any infection-induced
513 cross-reactive T cell memory, and whether this affects robustness of protection from SARS-
514 CoV-2 variants or wider coronaviruses.

515
516 SARS-CoV-2 reactive CD8+ T cells have been associated with milder disease⁴³, and as
517 previously mentioned, conflicting evidence has recently emerged regarding the presence of
518 pre-existing CD8+ T cells in unexposed patients. Nguyen et al¹⁶., found that SARS-CoV-2
519 CD8+ T cells in Australian pre-pandemic samples, including those recognising the
520 immunodominant HLA-B*07:02-SPR* complex, predominately displayed a naïve phenotype,
521 indicating a lack of pre-existing memory conferred by HCoV. In contrast, Francis et al¹⁵., found
522 that ~80% of unexposed individuals carrying HLA-B*07:02 show a pre-existing CD8+ T cell
523 response to HLA-B*07:02-SPR*. Francis et al argue that these pre-existing memory pools are
524 likely induced by prior exposure to HCoV, and that only a subpopulation of individuals
525 carrying specific HLA would possess such memory T cells. Our work is consistent with a
526 subset of COVID patients enriched for carrying HLA-B*07:02, and observed that in these
527 patients, their public T cells respond primarily to shared peptides. Despite not providing a link
528 to memory vs naïve responses, we build upon existing work by proposing additional alleles
529 which may be carried by individuals who possess cross-reactive T cells, as well as those which
530 appear depleted or absent in these individuals. Few studies have examined associations
531 between HLA type and COVID disease or its severity^{15,44,45}. Nevertheless, the emerging
532 picture is indicating that HCoV-SARS-CoV-2 cross-reactivity is conditioned by HLA
533 genotype. Together, we provide a landscape of TCR-pMHC interactions (all TCR-pMHC
534 interactions used in the analyses are found in Supplementary Data File 8) which may be
535 involved in HCoV-SARS-CoV-2 cross-reactivity and provide a framework for further anti-
536 viral mechanistic studies.

537
538 Although our study provides a map of shared and private SARS-CoV-2 peptides to date and
539 offers the extent to which one may expect CD8+ T cells cross-reactivity between HCoVs and
540 SARS-CoV-2, a limitation is that for cross-reactivity insights, we had to limit ourselves only
541 on CD8+ T cells for which both peptides and their cognate TCRs information were available.
542 Additionally, our approach for identifying homologous sequences seems to work better for
543 MHC class I peptides that are considerably shorter in length than their class II counterparts.
544 With a more suitable metric for longer peptides, one may substantiate our insights for class II.
545

546 Our metric for discriminating shared and private peptides is based on three factors: 1) sequence
547 homology at 50%, 2) physicochemical similarity of 75% and 3) that both source and target
548 peptides must be presented by the same HLA. Of these three, 50% sequence homology may
549 seem too relaxed. In support of our use of this threshold we note that: a) factors 2 and 3 are
550 additionally applied to compensate for this, b) we have checked our results with 70% sequence
551 homology and observed that main conclusions are robust, 3) as this map is suggested for further
552 functional validation we favour minimizing false negatives at the cost of potential false
553 positives.

554
555 Through examining the potential for cross-reactivity between SARS-CoV-2 and HCoV strains,
556 we have predicted that a set of 10 highly conserved immunogenic peptides could mount CD8+
557 T cell responses in 99% of the global population. These peptides have been reported previously
558 in *in silico* and experimental work⁴⁶⁻⁵⁰ however to our knowledge their large accumulated
559 global population coverage has not yet been reported. Some of these peptides exhibit similar
560 population coverage although with different HLA profiles, therefore it may be possible to tailor
561 a smaller set of peptides to specific regions of interest (based on local HLA frequency), thus
562 maximising coverage with a minimal set of peptides. Our work firstly identifies these peptides
563 as top candidates for cross-reactivity. Secondly, we propose that their high conservation across
564 strains may be of interest as pan-coronavirus targets, to assist ongoing work in search of
565 mitigation strategies to reduce the threat from mutant variants or emerging coronaviruses⁵¹⁻⁵³.

566
567 A complex facet of severe COVID-19 disease and its diverse clinical manifestation is
568 immunopathogenesis. Indeed, exacerbated immune responses including cytokine storm are a
569 primary clinical characteristic in severe COVID-19 patients. Aberrant transcriptional
570 programming has been observed in response to SARS-CoV-2⁵⁴, characterised by a failure of
571 type-1 and -3 interferon responses and simultaneous high induction of chemoattractants. While
572 the growing evidence for pre-existing HCoV cross-reactive memory T cell responses may
573 simply translate into an immunity benefit in some patients, in concert with data from MERS
574 and SARS-CoV-1, there are considerable evidence that cross-reactive T and B cell responses
575 may on the other hand be involved in immunopathology with SARS-CoV-2.

576
577 Venkatakrisnan et al.,⁵⁵ identified peptides that are identical between SARS-CoV-2 and the
578 human proteome. Their work demonstrates that the genes giving rise to these peptides are
579 expressed in tissues implicated in COVID-19 pathogenesis. Our work expands their insights,
580 by identifying SARS-CoV-2 peptides that are experimentally confirmed to be immunogenic,
581 with high similarity to the human proteome. Consistent with their conclusions, we find
582 similarity of immunogenic SARS-CoV-2 peptides to human genes e.g., CCL3, CCL31 and
583 CD163. These insights are of particular interest given the elevated cytokine and chemokine
584 responses in severe COVID patients. More broadly, there is evidence that viral antigens that
585 are structurally similar to self-antigens can be involved in inducing autoimmunity via
586 molecular mimicry⁸. For these reasons, we propose these peptides as candidates which may
587 exhibit immunopathological or autoimmune associations.

588
589 In conclusion, we have employed an *in-silico* approach to examine the evidence surrounding
590 cross-reactive SARS-CoV-2 CD8+ T cell responses. We observed a set of SARS-CoV-2
591 candidates with high similarity to the human proteome and suggest investigation into whether
592 they provoke immunopathology. We have also provided evidence of CD8+ T cell cross-
593 reactivity, not only to an extent which indicates that naïve individuals could mount cross-
594 reactive responses to SARS-CoV-2 and common-cold coronaviruses, but we also found that
595 SARS-CoV-2 infection induces CD8+ T cell responses against peptides with high similarity to

596 HCoV in some COVID-19 patients. We build upon existing evidence that such cross-reactivity
597 is conditioned by presence of specific HLA alleles and envision that the insights presented here
598 are leveraged to explore whether these potentially cross-reactive T cells and cognate pMHCs
599 influence COVID-19 disease heterogeneity, vaccine- or infection-induced protection from
600 SARS-CoV-2 and its emerging variants of concern.

601

602 Acknowledgments:

603 We greatly acknowledge conversations and guidance from Dr Mikhail Shugay (ITM,
604 Moscow), and Dr Giorgio Napolitani (KCL, London).

605

606 Author contributions:

607 HK conceived, designed and supervised the project. PB performed computational analyses
608 with insights from CL and MPP and AA. HK and PB interpreted the results. AS, GO assisted
609 design, interpretation and supervision. ROB, JW commented on manuscript. HK, AS funded
610 the project. HK, PB, wrote the manuscript with contributions from CL, MPP and AA, AS and
611 GO.

612

613 Declaration of interest

614 The authors declare no competing interests.

615 **Methods**

616

617 **Data Processing and Analysis**

618

619 All data processing and analysis was performed using the R plugin for Pycharm 2020, in either R
620 4.0.3 or 4.0.1.

621

622 **Curating a pool of SARS-CoV-2 class I and II peptides**

623

624 Human immunogenic and non-immunogenic SARS-CoV-2 peptide data were gathered from the both
625 the IEDB and the Virus Pathogen Resource (VIPR) (accessed 11-02-2021). ‘T cell’ assay, ‘Human’
626 host and SARS-CoV-2 organism options were selected. If an observation was found in both datasets,
627 the one from the IEDB was retained. Protein names were cleaned and standardised where possible.
628 Immunogenic peptides not observed in either the IEDB or VIPR were also gathered from the ‘MIRA’
629 dataset which maps cognate TCRs and SARS-CoV-2 peptides.

630

631 **Retrieval of Coronavirus Proteome Sequences**

632

633 NCBI reference genomes were gathered for OC43

634 (<https://www.ncbi.nlm.nih.gov/nuccore/1578871709/>), HKU1

635 (https://www.ncbi.nlm.nih.gov/nuccore/NC_006577.2), 229E

636 (https://www.ncbi.nlm.nih.gov/nuccore/NC_002645.1), NL63

637 (<https://www.ncbi.nlm.nih.gov/nuccore/49169782/>), and SARS-CoV-2-Wuhan

638 (https://www.ncbi.nlm.nih.gov/nuccore/NC_045512.2).

639

640 **MHC Presentation Prediction**

641

642 Antigen presentation by MHC class I was predicted using NetMHCpan v4.1 against HLA-A*0101,
643 0201, 0301, 2402, HLA-B*0702, 4001, 0801, and HLA-C*0702, 0401, 0701 alleles. Antigen
644 presentation by MHC class II was predicted using netMHCIIpan against the most common sets of
645 alleles found in the IEDB, for which this model can make predictions to. The alleles are: DRB1-0101,
646 0102, 0301, 0401, 0402, 0402, 0404, 0701, 0801, 0901, 1001, 1101, 1104, 1201, 1202, 1301,1302,

647 1303, 1401, 1406, 1501, 1502, 1601, 1602, DRB3_0101, 0202 and DRB5_0101, 0102. Peptides with
648 a rank score ≤ 2.0 were classified as binders.

649

650 **HLA ligand enrichment analysis for SARS-CoV-2 proteins**

651

652 To provide reasonable statistical inference, we only examined proteins longer than 100 amino acids.

653 To compute enrichment or depletion, we followed the approach by Karnaukhov et al. First, we

654 predicted using netMHCpan v 4.1 the number of ligands N_i of length l from each SARS-CoV-2

655 protein i which adheres to the criteria. Probability of a HLA allele presenting a peptide was computed

656 as the average number of ligands per allele:

657

$$658 \quad p = \mu N_i / \mu L_i$$

659

660 where L_i is the corrected protein length (length of protein $- l$), and μ denotes the average over the

661 assessed SARS-CoV-2 proteins. It follows that the probability of observing a given number of ligands

662 from each SARS-CoV-2 protein is computed using the binomial distribution as:

663

$$664 \quad P(N_i) = P_{binom}(N_i | p, L_i)$$

665

666 The logs odds ratio (enrichment or depletion) is calculated as:

667

$$668 \quad \log\left(\frac{N_i}{pL_i}\right)$$

669

670 **Discriminating Shared and Private SARS-CoV-2 Peptides**

671

672 To compare a SARS-CoV-2 peptide a , of length N to a proteome of interest, all possible linear

673 peptides of length N were generated from said proteome. This can be thought of as scanning along the

674 proteome of interest with a step size of 1, generating all peptides of length N . The deriving protein

675 was recorded. Three metrics – which all must be satisfied - were used determine whether a peptide is

676 considered shared with HCoV or private to SARS-CoV-2. We below describe each metric, and then

677 explain the three thresholds which all must be achieved for a peptide to be classified as ‘shared.’.

678

679 Firstly, once all peptides from the proteome of interest of length N are generated, a similarity index

680 we call the ‘MatchScore’ is calculated for each pairwise comparison. This metric is charged with

681 assessing physicochemical similarity between two peptides of interest. For each SARS-CoV-2

682 peptide, the highest ‘MatchScore’ against each HCoV protein is retained and the rest are discarded.

683 To calculate the ‘MatchScore’, we employ the method designed by Bresciani et al²⁴. Briefly, for two

684 peptides a or b of length N , the similarity score is given as;

685

$$686 \quad MatchScore = \frac{bl(a, b)}{\sqrt{bl(a, a) \times bl(b, b)}}$$

687

688 where $bl(a, b)$ is the BLOSUM62 score for peptide a vs b , and $bl(a, a)$ is the BLOSUM62 score for

689 peptide a vs a , etc. BLOSUM62 local-global alignment scores (local or global would produce the

690 same score for a pairwise alignment of lengths N vs N) were computed using the pairwiseAlignment

691 function from the R package Biostrings, with high gap penalties (opening and extension of both 100).

692 The MatchScore function produces a score where 1 reflects an exact match, i.e no mismatches in two

693 sequences, and 0 reflects high dissimilarity.

694 Criteria 1: A shared peptide and its HCoV match must have a MatchScore of >0.75 .

695

696 The second metric is based on sequence homology between two sequences, essentially reflecting the

697 proportion of amino acid positions in the SARS-CoV-2 peptide, which are conserved in the HCoV

698 match. This is calculated as:

699

700

$$ProportionMismatched = \frac{HammingDistance}{Length}$$

701

702 where ‘*HammingDistance*’ is the hamming distance between two peptides of interest, which
703 calculates the number of different positions, and ‘*Length*’ is the length of the compared peptides.

704 Criteria 2: The *ProportionMismatched* between a shared peptide and its HCoV match must be < 0.5
705 (50%).

706

707 Naturally, the inverse of this is true, in that at least 50% amino acid conservation between a SARS-
708 CoV-2 peptide and a HCoV match must be observed for the peptide to be considered ‘shared’.

709

710 The third metric is based on predicted presentation by HLA of the SARS-CoV-2 peptide and its
711 HCoV match.

712 Criteria 3: Both the SARS-CoV-2 peptide and its HCoV match must be predicted to bind at least one
713 common HLA allele.

714

715 All three criteria must be satisfied for a SARS-CoV-2 peptide to be classified as a shared peptide and
716 also for a match from HCoV to be considered a homologous match. *doParallel* and *foreach* functions
717 were used to parallelise the processing.

718

719 **Sequence Similarity with the Human Proteome and Human Microbiomes**

720

721 Here, the same similarity criteria were employed as in the previous HCoV section. However, in
722 contrast with HCoV comparison, due to the size of the human proteomes and microbiomes, the best
723 match against the whole proteome is retained. *doParallel* and *foreach* functions were used to
724 parallelise the processing.

725

726 **Gathering human proteome sequence**

727

728 The human proteome was downloaded in fasta format from UniProt
729 <https://www.uniprot.org/proteomes/UP000005640>

730

731 **Gathering human microbiome sequences**

732

733 Human gut and airways microbiomes were downloaded from the HMP Data Analysis and
734 Coordination Center <http://www.hmpdacc.org/HMRGD>. The complete set of genomes were
735 downloaded in fasta format in ‘Protein multifasta (PEP) format’. For gut, the body site was specified
736 as ‘gastrointestinal tract’. 457 and 50 gut and airway microbiota were available respectively.

737

738 **Human gene sets with sequence similarity to SARS-CoV-2 immunogenic peptides**

739

740 The SARS-CoV-2 peptides of lengths 9 and 10 with a similarity score to the human proteome in the
741 top 10 percentile were gathered. Only predicted binders (see MHC presentation prediction) were
742 retained.

743

744 **CD8+ T cell cross-reactivity maps using IEDB receptor data**

745

746 The entire IEDB receptor data for SARS-CoV-2 peptides was downloaded. Bipartite graphs were
747 generated using *iGraph* and *Matrix* libraries in R. Bipartite graphs were projected into one-mode
748 graphs using the *bipartite_projection* function. All graphs were exported from *iGraph* into Cytoscape
749 v3.82 using the R function *createNetworkFromIgraph* from package *RCy3*. From Cytoscape,
750 ‘.graphml’ files were exported and opened with Gephi. Gephi was used to finalise the diagrams and
751 improve visual aesthetics. Either ‘ForceAtlas’ or ‘Fructerman-Reingold’ templates were used. Gravity
752 and repulsion parameters were altered to improve visual aesthetics.

753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805

CD8+ T cell CDR3 Kmer Enrichment

R Package *immunarch*⁵⁶ was used to compute Kmer (K=5 in this case) statistics for CDR3 sequences, and to visualise enrichment. See https://immunarch.com/articles/web_only/v9_kmers.html for full details.

Gathering clinical and TCR repertoire data for COVID-19 patients and healthy subjects

The COVID-19 MIRA dataset (>160k high-confidence SARS-CoV-2-specific TCRs) was downloaded from <https://clients.adaptivebiotech.com/pub/covid-2020> with corresponding sample metadata. These data contain TCR repertoire data mapped to SARS-CoV-2 epitopes from 5 patient cohorts, including COVID convalescent patients and healthy subjects with no known exposure to SARS-CoV-2. Only convalescent patients and healthy subjects were used in the analysis due to low numbers of subjects for other cohorts.
<https://clients.adaptivebiotech.com/pub/covid-2020>

Networks of TCRs recognising Shared and/or Private Peptides

A public TCR is defined as a CDR3 sequence and V and J gene which is observed in more than one patient in the MIRA dataset. All graphs were first generated using *iGraph* in R, exported to Cytoscape using the *createNetworkFromIgraph* function in the *RCy3* package. From cytoscape, all graphs were exported as .graphml files and read into Gephi. In Gephi, either ‘ForceAtlas’ and ‘Fruchterman-Reingold’ templates were used. In almost all cases, gravity and repulsion parameters were adjusted to improve visual aesthetics.

Estimating population coverage of SARS-CoV-2 peptides with high conservation to three or more HCoV

We followed the approach by Ahmedid et al⁴². Population coverage is an estimate of the proportion of individuals in a given population that may mount a T cell response against a peptide. Population coverage is predicted based on HLA alleles for each immunogenic peptide as predicted by netMHCpan 4.1, leading to individual population coverage of a peptide. To predict accumulated coverage, we began with the peptide with the highest individual coverage “FVDG*”, and incrementally added a peptide and predicted accumulated coverage. The population coverage of a set of peptides (i.e accumulated coverage), is defined as the proportion of individuals able to mount a T cell response to at least one peptide in the set. Python code for the IEDB tool to compute the population coverage was downloaded from <http://tools.iedb.org/population/download> on 24-11-20.

Reference:

1. Krause, P. R. *et al.* SARS-CoV-2 Variants and Vaccines. (2021).
2. Doshi, P. Covid-19: Do many people have pre-existing immunity? *BMJ* **370**, (2020).
3. Reynolds, C. J. *et al.* Prior SARS-CoV-2 infection rescues B and T cell responses to variants after first vaccine dose. *Science* (80-.). eabh1282 (2021)
doi:10.1126/science.abh1282.
4. Goel, R. R. *et al.* Distinct antibody and memory B cell responses in SARSCoV-2 naïve and recovered individuals following mRNA vaccination. *Sci. Immunol.* **6**, 1–19 (2021).
5. Beretta, A., Cranage, M. & Zipeto, D. Is Cross-Reactive Immunity Triggering COVID-19 Immunopathogenesis? *Front. Immunol.* **11**, 1–9 (2020).
6. Huang, A. T. *et al.* A systematic review of antibody mediated immunity to coronaviruses: kinetics, correlates of protection, and association with severity. *Nat.*

- 806 *Commun.* **11**, 1–16 (2020).
- 807 7. Grifoni, A. *et al.* SARS-CoV-2 Human T cell Epitopes: adaptive immune response
808 against COVID-19. *Cell Host Microbe* (2021) doi:10.1016/j.chom.2021.05.010.
- 809 8. Smatti, M. K. *et al.* Viruses and Autoimmunity: A Review on the Potential Interaction
810 and Molecular Mechanisms. *Viruses* **11**, 762 (2019).
- 811 9. Braun, J. *et al.* Presence of SARS-CoV-2-reactive T cells in COVID-19 patients and
812 healthy donors. *medRxiv* (2020) doi:10.1101/2020.04.17.20061440.
- 813 10. Grifoni, A. *et al.* Targets of T Cell Responses to SARS-CoV-2 Coronavirus in Humans
814 with COVID-19 Disease and Unexposed Individuals. *Cell* **181**, 1489–1501.e15 (2020).
- 815 11. Le Bert, N. *et al.* SARS-CoV-2-specific T cell immunity in cases of COVID-19 and
816 SARS, and uninfected controls. *Nature* **584**, 457–462 (2020).
- 817 12. Mateus, J. *et al.* Selective and cross-reactive SARS-CoV-2 T cell epitopes in
818 unexposed humans. *Science (80-.)*. **370**, (2020).
- 819 13. Peng, Y. *et al.* Broad and strong memory CD4+ and CD8+ T cells induced by SARS-
820 CoV-2 in UK convalescent individuals following COVID-19. *Nat. Immunol.* **21**,
821 1336–1345 (2020).
- 822 14. Weiskopf, D. *et al.* Phenotype of SARS-CoV-2-specific T-cells in COVID-19 patients
823 with acute respiratory distress syndrome. *medRxiv* 1–29 (2020)
824 doi:10.1101/2020.04.11.20062349.
- 825 15. Francis, J. M. *et al.* Allelic variation in Class I HLA determines pre-existing memory
826 responses to SARS-CoV-2 that shape the CD8 + T cell repertoire upon viral exposure
827 Collection and Processing Team. doi:10.1101/2021.04.29.441258.
- 828 16. Nguyen, T. H. O. *et al.* CD8+ T cells specific for an immunodominant SARS-CoV-2
829 nucleocapsid epitope display high naïve precursor frequency and T cell receptor
830 promiscuity. *Immunity* (2021) doi:10.1016/j.immuni.2021.04.009.
- 831 17. Karnaukhov, V. *et al.* HLA binding of self-peptides is biased towards proteins with
832 specific molecular functions. *bioRxiv* 2021.02.16.431395 (2021).
- 833 18. Trolle, T. *et al.* The Length Distribution of Class I–Restricted T Cell Epitopes Is
834 Determined by Both Peptide Supply and MHC Allele–Specific Binding Preference. *J.*
835 *Immunol.* **196**, 1480–1487 (2016).
- 836 19. Croft, N. P. *et al.* Most viral peptides displayed by class I MHC on infected cells are
837 immunogenic. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 3112–3117 (2019).
- 838 20. Bresciani, A. *et al.* T-cell recognition is shaped by epitope sequence conservation in
839 the host proteome and microbiome. *Immunology* **148**, 34–39 (2016).
- 840 21. Tarke, A. *et al.* Comprehensive analysis of T cell immunodominance and
841 immunoprevalence of SARS-CoV-2 epitopes in COVID-19 cases. *Cell Reports Med.*
842 **2**, 100204 (2021).
- 843 22. Tan, H.-X. *et al.* Adaptive immunity to human coronaviruses is widespread but low in
844 magnitude 1 2. *medRxiv* 2021.01.24.21250074 (2021)
845 doi:10.1101/2021.01.24.21250074.
- 846 23. Lee, C. H. *et al.* Potential CD8+ T Cell Cross-Reactivity Against SARS-CoV-2
847 Conferred by Other Coronavirus Strains. *Front. Immunol.* **11**, 2878 (2020).
- 848 24. Bresciani, A. *et al.* T-cell recognition is shaped by epitope sequence conservation in
849 the host proteome and microbiome. *Immunology* **148**, 34–39 (2016).
- 850 25. Wells, D. K. *et al.* Key Parameters of Tumor Epitope Immunogenicity Revealed
851 Through a Consortium Approach Improve Neoantigen Prediction. *Cell* **0**, (2020).
- 852 26. Wraith, D. C., Goldman, M. & Lambert, P. H. Vaccination and autoimmune disease:
853 What is the evidence? *Lancet* vol. 362 1659–1666 (2003).
- 854 27. Boekel, L. *et al.* Perspective of patients with autoimmune diseases on COVID-19
855 vaccination. *The Lancet Rheumatology* vol. 3 e241–e243 (2021).

- 856 28. Cao, X. COVID-19: immunopathology and its implications for therapy. *Nature*
857 *Reviews Immunology* vol. 20 269–270 (2020).
- 858 29. Gómez-Rial, J. *et al.* Increased Serum Levels of sCD14 and sCD163 Indicate a
859 Preponderant Role for Monocytes in COVID-19 Immunopathology. *Front. Immunol.*
860 **11**, 1–8 (2020).
- 861 30. Liu, X. *et al.* Single-cell analysis reveals macrophage-driven T cell dysfunction in
862 severe COVID-19 patients. *medRxiv* (2020) doi:10.1101/2020.05.23.20100024.
- 863 31. Qin, C. *et al.* Dysregulation of immune response in patients with coronavirus 2019
864 (COVID-19) in Wuhan, China. *Clin. Infect. Dis.* **71**, 762–768 (2020).
- 865 32. Wang, J., Jiang, M., Chen, X. & Montaner, L. J. Cytokine storm and leukocyte
866 changes in mild versus severe SARS-CoV-2 infection: Review of 3939 COVID-19
867 patients in China and emerging pathogenesis and therapy concepts. *J. Leukoc. Biol.*
868 **108**, 17–41 (2020).
- 869 33. Szabo, P. A. *et al.* Longitudinal profiling of respiratory and systemic immune
870 responses reveals myeloid cell-driven lung inflammation in severe COVID-19.
871 *Immunity* **54**, 797–814.e6 (2021).
- 872 34. Stikker, B. Severe COVID-19 associated variants linked to chemokine receptor gene
873 control in monocytes and macrophages Authors: *bioRxiv Prepr. Serv. Biol.* (2021).
- 874 35. Patterson, B. K. *et al.* CCR5 inhibition in critical COVID-19 patients decreases
875 inflammatory cytokines, increases CD8 T-cells, and decreases SARS-CoV2 RNA in
876 plasma by day 14. *Int. J. Infect. Dis.* **103**, 25–32 (2021).
- 877 36. Patterson, B. K. *et al.* Persistence of SARS CoV-2 S1 Protein in CD16 + Monocytes in
878 Post- Acute Sequelae of COVID-19 (PASC) Up to 15 Months Post-Infection. (2021).
- 879 37. Wang, P. *et al.* Comprehensive analysis of TCR repertoire in COVID-19 using single
880 cell sequencing. *Genomics* **113**, 456–462 (2021).
- 881 38. Zhang, J. Y. *et al.* Single-cell landscape of immunological responses in patients with
882 COVID-19. *Nat. Immunol.* **21**, 1107–1118 (2020).
- 883 39. Nolan, S. *et al.* A large-scale database of T-cell receptor beta (TCR β) sequences and
884 binding associations from natural and synthetic exposure to SARS-CoV-2. *Res. Sq.* 1–
885 28 (2020) doi:10.21203/rs.3.rs-51964/v1.
- 886 40. Lineburg, K. E. *et al.* CD8+ T cells specific for an immunodominant SARS-CoV-2
887 nucleocapsid epitope cross-react with selective seasonal coronaviruses. *Immunity* **54**,
888 1055-1065.e5 (2021).
- 889 41. Bui, H. H. *et al.* Predicting population coverage of T-cell epitope-based diagnostics
890 and vaccines. *BMC Bioinformatics* **7**, 1–5 (2006).
- 891 42. Ahmedid, S. F., Quadeer, A. A., Barton, J. P. & McKay, M. R. Cross-serotypically
892 conserved epitope recommendations for a universal T cell-based dengue vaccine. *PLoS*
893 *Negl. Trop. Dis.* **14**, 1–22 (2020).
- 894 43. Sekine, T. *et al.* Robust T Cell Immunity in Convalescent Individuals with
895 Asymptomatic or Mild COVID-19. *Cell* **183**, 158-168.e14 (2020).
- 896 44. Shkurnikov, M. *et al.* Association of HLA Class I Genotypes With Severity of
897 Coronavirus Disease-19. *Front. Immunol.* **12**, (2021).
- 898 45. Anzurez, A. *et al.* Association of HLA- DRB1 *09:01 with severe COVID -19 . *Hla*
899 1–6 (2021) doi:10.1111/tan.14256.
- 900 46. Qiao, R., Tran, N. H., Shan, B., Ghodsi, A. & Li, M. Personalized workflow to identify
901 optimal T-cell epitopes for peptide-based vaccines against COVID-19. *arXiv* (2020).
- 902 47. Slathia, P. S. & Sharma, P. Prediction of T and B cell epitopes in the proteome of
903 SARS-CoV-2 for potential use in diagnostics and vaccine design. *ChemRxiv* (2020)
904 doi:10.26434/chemrxiv.12116943.v1.
- 905 48. Pacholczyk, M. & Rieske, P. In silico studies suggest T-cell cross-reactivity between

- 906 SARS-CoV-2 and less dangerous coronaviruses. 1–12.
- 907 49. Dijkstra, J. M. & Hashimoto, K. Expected immune recognition of COVID-19 virus by
908 memory from earlier infections with common coronaviruses in a large part of the
909 world population. *F1000Research* **9**, (2020).
- 910 50. Gangaev, A. & Kvistborg, P. 863 Identification and characterization of an
911 immunodominant SARS-CoV-2-specific CD8 T cell response. *J. Immunother. Cancer*
912 **8**, A916–A916 (2020).
- 913 51. A. Almofti, Y., Ali Abd-elrahman, K., Abd Elgadir Gassmallah, S. & Ahmed Salih,
914 M. Multi Epitopes Vaccine Prediction against Severe Acute Respiratory Syndrome
915 (SARS) Coronavirus Using Immunoinformatics Approaches. *Am. J. Microbiol. Res.* **6**,
916 94–114 (2018).
- 917 52. Li, M. *et al.* Rational Design of a Pan-Coronavirus Vaccine Based on Conserved CTL
918 Epitopes. *Viruses* **13**, 1–11 (2021).
- 919 53. Koff, W. C. & Berkley, S. F. A universal coronavirus vaccine. *Science (80-.)*. **371**,
920 759 (2021).
- 921 54. Blanco-Melo, D. *et al.* Imbalanced Host Response to SARS-CoV-2 Drives
922 Development of COVID-19. *Cell* **181**, 1036-1045.e9 (2020).
- 923 55. Venkatakrishnan, A. J. *et al.* Benchmarking evolutionary tinkering underlying human–
924 viral molecular mimicry shows multiple host pulmonary–arterial peptides mimicked
925 by SARS-CoV-2. *Cell Death Discov.* **6**, (2020).
- 926 56. Nazarov, V., immunarch.bot & Rumynskiy, E. immunomind/immunarch: 0.6.5: Basic
927 single-cell support. (2020) doi:10.5281/ZENODO.3893991.
- 928

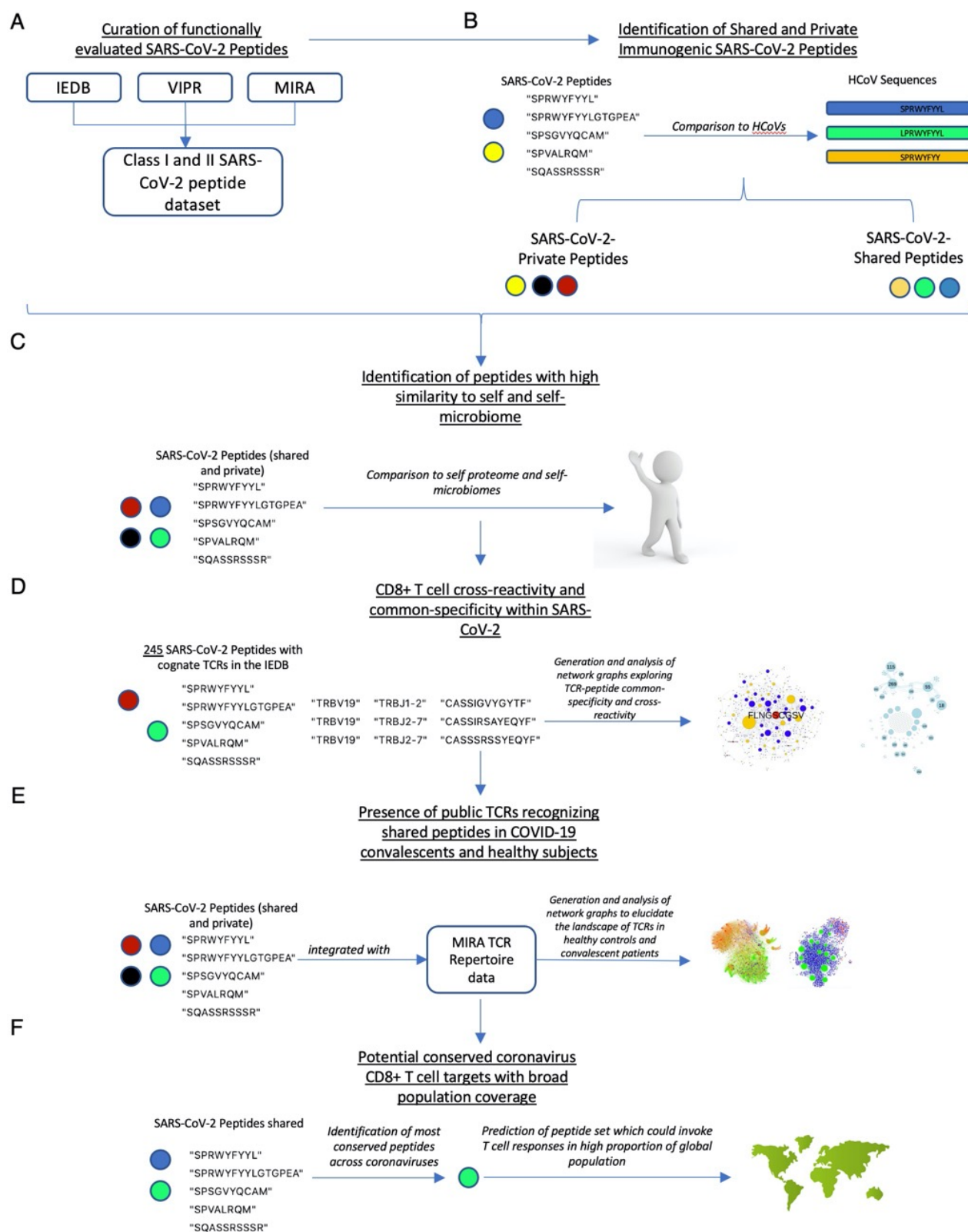


Figure 1: Overview of the study. A) Functionally evaluated SARS-CoV-2 peptides are gathered from three online repositories. Data are cleaned and integrated, thus curating a comprehensive pool of SARS-CoV-2 class I and II peptides. Exploratory data analysis followed. B) Next, each immunogenic SARS-CoV-2 peptide of length 1 is compared to each possible linear peptide of length 1 from four common-cold causing human coronavirus strains. Based on similarity criteria and following confirmation that the target hit from HCoV is predicted to bind HLA, peptides are classified as 'shared' SARS-CoV-2-HCoV peptides. Those which do not adhere to the criteria are classified as SARS-CoV-2 private. C) The entire set of SARS-CoV-2 shared and private, immunogenic and non-immunogenic peptides is compared to the human proteome and gut and airways microbiomes. D) 245 peptides from our SARS-CoV-2 peptide dataset have known cognate TCRs in the IEDB. These peptide-TCR associations were examined to explore the extent of cross-reactivity and common-specificity within SARS-CoV-2. E) Both shared and private immunogenic SARS-CoV-2 peptides are integrated with the COVID-19 MIRA TCR repertoire dataset and employed to examine the presence of TCRs recognizing shared/private SARS-CoV-2 peptides in health and/or disease. F) The entire set of 126 shared SARS-CoV-2 peptides is searched for those most highly conserved across coronaviruses. This resulted in 17 peptides, of which we used to predict global and regional population coverage, given predicted HLA alleles.

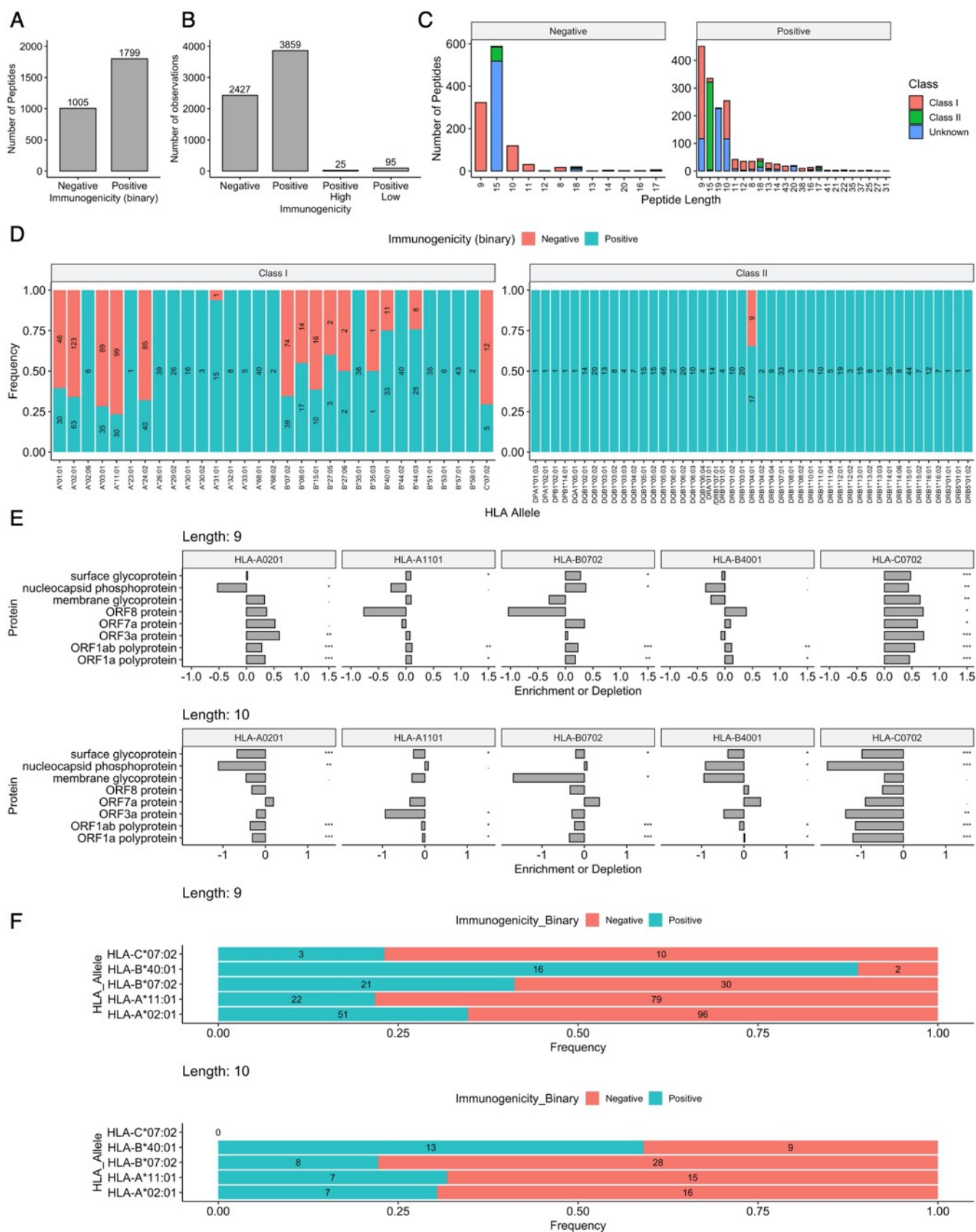


Figure 2: A comprehensive pool of functionally validated SARS-CoV-2 peptides. Barplots showing A) The number of SARS-CoV-2 peptides deemed 'positive' or 'negative'. A 'negative' label reflects that for a peptide, there are only negative (nonimmunogenic) qualitative observations while 'positive' reflects at least one immunogenic observation. B) The number of total pMHC complexes observed in the dataset, including all assay and HLA combinations for each peptide. C) The distribution of lengths of complexes amongst our SARS-CoV-2 dataset. Left plot shows nonimmunogenic 'Negative' peptides. Right plot shows immunogenic or 'Positive' peptides. MHC class is colour coded. D) The frequencies of immunogenic or non-immunogenic class I or II observations, for peptides where specific HLA allele information is available. Numeric labels show the number of peptides in each group. E) The log odds ratio of observed and expected number of presented peptides of lengths 9 and 10 by common HLA alleles for SARS-CoV-2 proteins > 100 amino acids in length. Significance calculated using binomial distribution. F) The frequency of immunogenic and nonimmunogenic peptides as presented by common class I HLA alleles arising from SARS-CoV-2. Numeric labels show the number of observations per immunogenicity status.

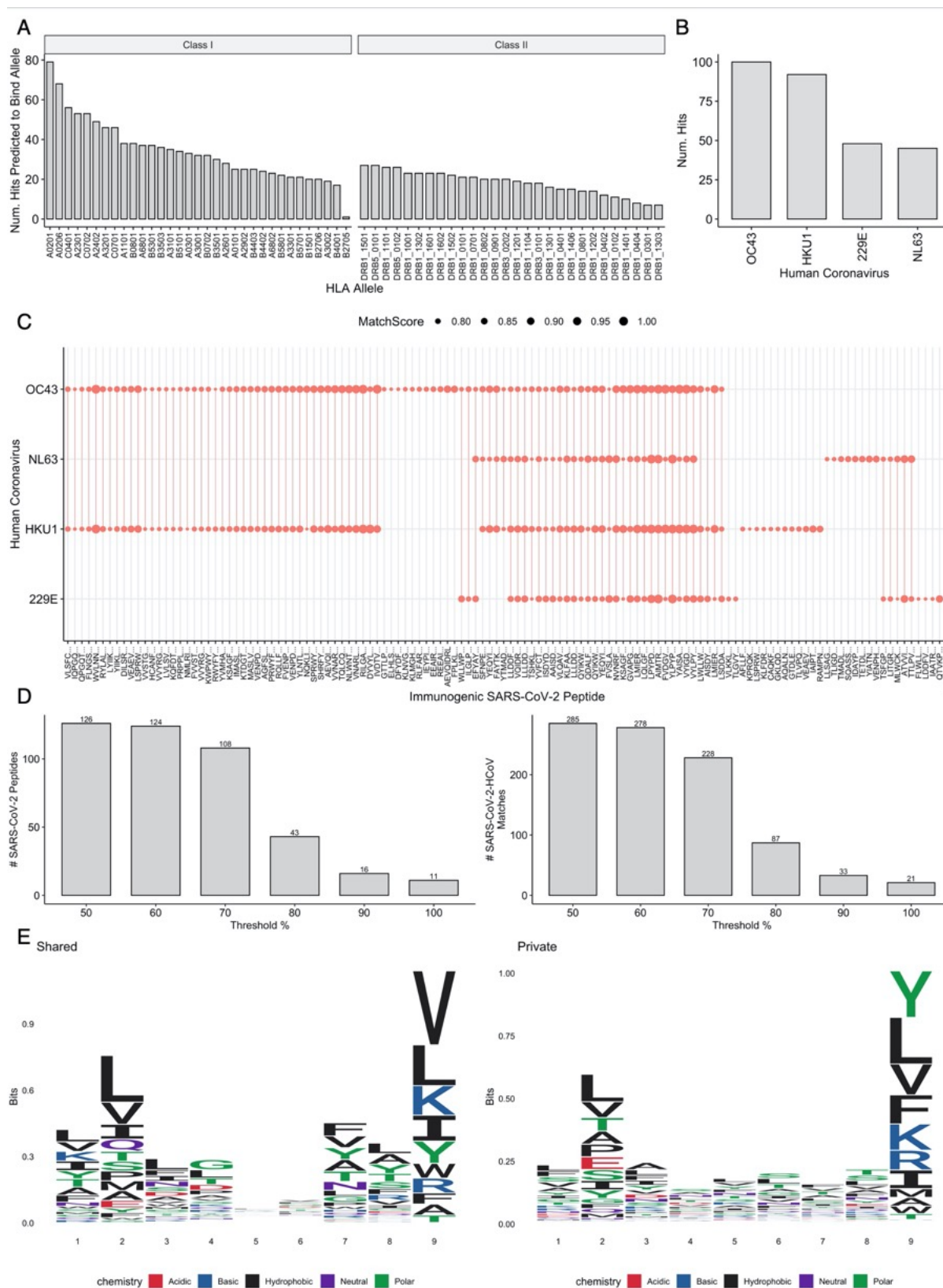


Figure 3: A set of peptides from human coronavirus strains with high similarity to immunogenic SARS-CoV-2 peptides. A) A barplot showing the number of high similarity matches predicted to bind a set of common HLA class I and class II alleles. B) A barplot showing the number of unique high similarity matches derived from each human common-cold-causing coronavirus. Each hit is defined as a unique observation with MatchScore > 0.75, between an immunogenic SARS-CoV-2 peptide with length x, and a stretch of length x from one viral protein. C) A dot and line plot showing each SARS-CoV-2 peptide and to which common-cold-causing coronavirus it exhibits a high similarity match. The size of each point reflects the MatchScore, i.e the primary similarity metric. D) Barplots showing the number of unique SARS-CoV-2 peptides (left) and SARS-CoV-2-HCoV matches (right), at different thresholds of the sequence homology metric, i.e the % of the amino acids that must be conserved between the SARS-CoV-2 peptide and its HCoV match. E) Sequence logo plots comparing amino acid usage amongst SARS-CoV-2 shared and private 9-mer peptide sequences.

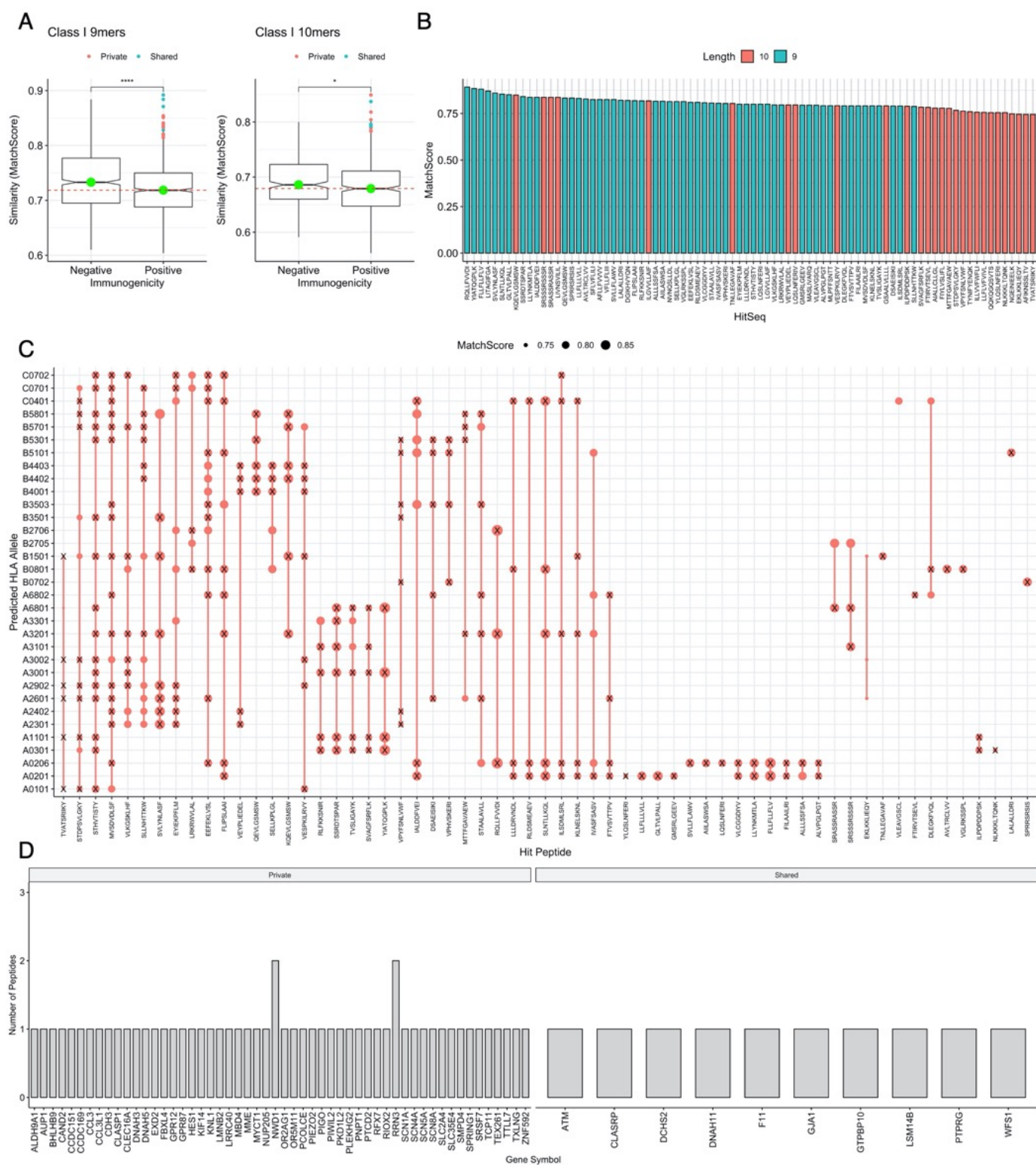


Figure 4: A pool of immunogenic SARS-CoV-2 peptides with high similarity to human genes. A) Notched boxplots showing the similarity (evaluated by the MatchScore) of nonimmunogenic and immunogenic SARS-CoV-2 peptides with sequences derived from the human proteome of lengths 9 (immunogenic $n = 906$, nonimmunogenic $n=734$ peptides) and 10 (immunogenic $n=394$ and nonimmunogenic $n=394$ peptides). The green dot shows the median of each group, red line represents the median of the immunogenic group. B) A barplot showing the MatchScores of hits with high similarity to the human proteome, colour-coded by protein length. C) A dotplot showing the predicted HLAs of candidate peptides with high similarity to SARS-CoV-2 proteome. The size of the point reflects the MatchScore. An 'X' shows where the SARS-CoV-2 derived peptide and the match are predicted to bind the same allele. D) A bar chart showing the genes from which the high similarity match peptides arise in the human proteome, separated by the "Private" or "Shared" epitope status of the SARS-CoV-2 epitope.

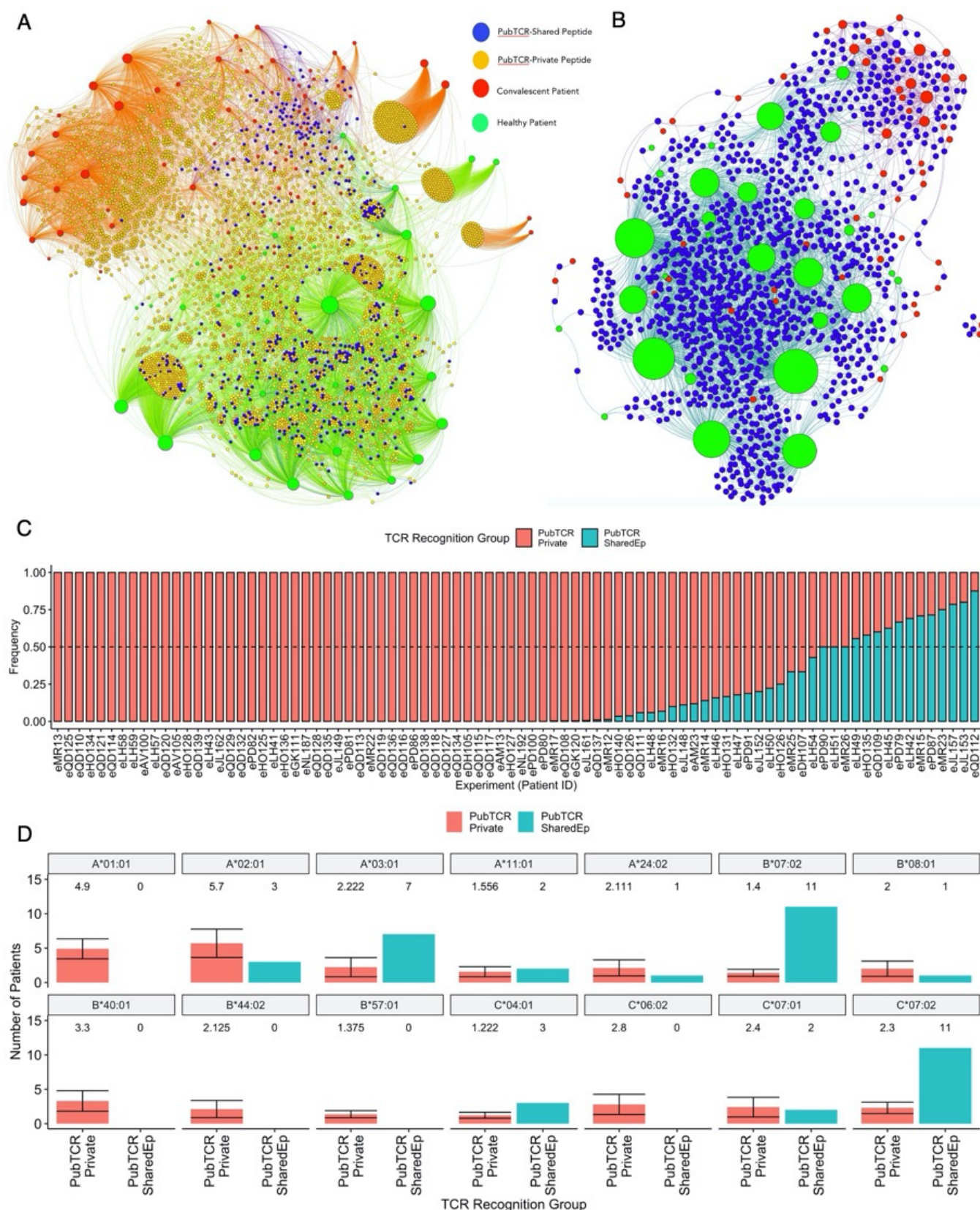


Figure 5: A landscape of T cell responses against SARS-CoV-2-HCoV Shared and SARS-CoV-2-private peptides in healthy or COVID-19 convalescent individuals: A) A bipartite network graph showing SARS-CoV-2-specific public TCRs which recognize shared or private SARS-CoV-2 peptides in healthy or convalescent patients. TCRs are colour-coded by whether they recognize only shared peptides (blue), only private peptides (orange) or both (yellow). COVID-19 convalescent patients are labelled red while healthy controls are labelled green. Node size reflects degree of connectivity, i.e. the quantity of an individual's TCRs which are shared with other patients. B) A bipartite network graph showing SARS-CoV-2 public TCRs that recognize SARS-CoV-2-HCoV shared peptides. Patient node size reflects the quantity of their TCRs which are shared with another patient. Healthy patients are labelled green, COVID-19 convalescent are labelled red, and (public) TCRs are labelled blue. C) A barplot showing the frequency that each convalescent patient's public TCRs recognize SARS-CoV-2-private (red) or SARS-CoV-2-HCoV-shared (blue) peptides. Patients with identical frequencies are ordered by the number of TCRs. D) Barplots showing the quantities of COVID-19 convalescent patients who carry 14 class I HLA alleles of interest. Patients are grouped by whether their public TCRs predominately recognize "Private" (PubTCR_Private, n=12, sampled 10 times) or "Shared" (PubTCR_SharedEp, n=12) peptides. For the PubTCR-Private group, 12 patients were sampled 10 times and the number of patients carrying alleles was measured. The mean number of patients carrying each allele and the error are visualized. For the PubTCR-SharedEp group, the data contain only 12 patients of interest, thus the number of patients carrying each allele is measured and visualised.

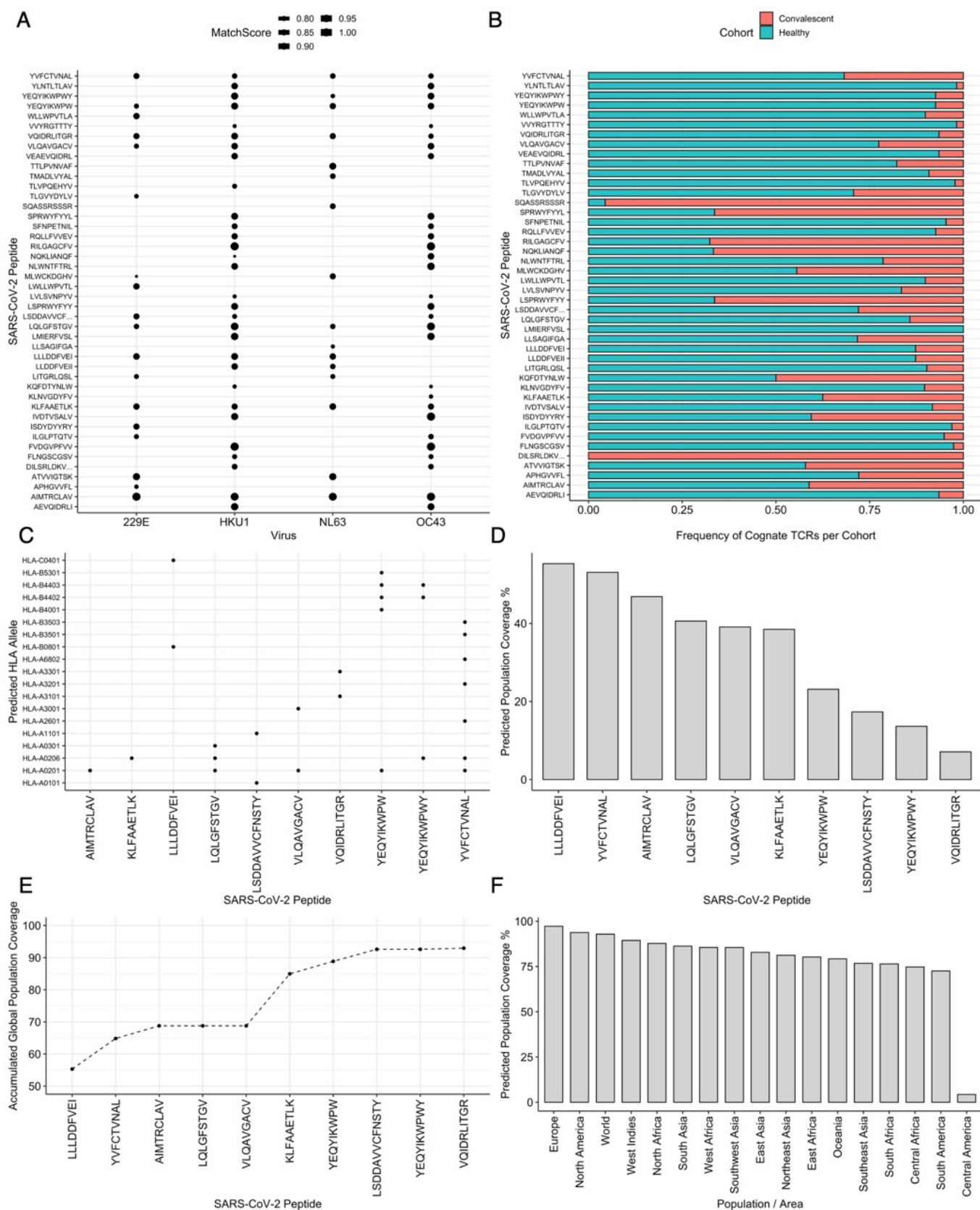


Figure 6: T cell epitopes with known cognate TCRs which are conserved across multiple HCoV and SARS-CoV-2 exhibit broad population coverage: A) A dot plot showing SARS-CoV-2 peptides with high similarity to more than one HCoV that are recognized by TCRs in the MIRA dataset. Size of the dot represents the MatchScore. B) The frequency of cognate TCRs which recognize these peptides from the COVID-19 convalescent or healthy cohorts. C) The HLA alleles predicted to present SARS-CoV-2 peptides with high similarity matches to 3 or 4 HCoV strains. D) Global population coverage as calculated by the 'IEDB population coverage tool' for each individual SARS-CoV-2 peptide with high similarity matches to 3 or 4 HCoV strains. E) Accumulated global population coverage predicted by the IEDB population coverage tool. F) Regional population coverage for the entire set of 10 SARS-CoV-2 peptides with matches to 3 or 4 HCoV.

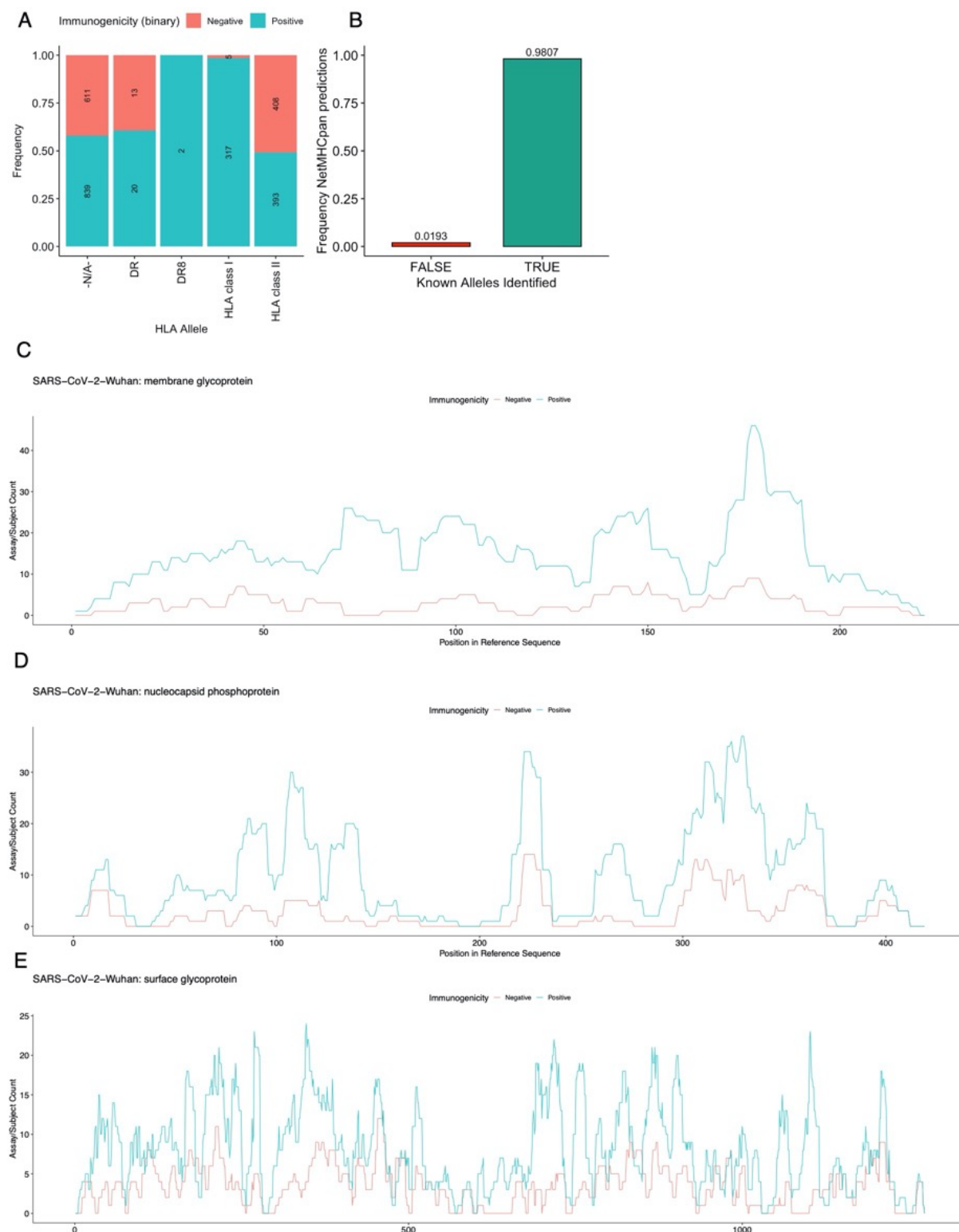
Human Match	MatchScore	SARS-CoV-2 Peptide	HLA Allele	Shared/Private	Protein	Length	Proportion Conserved
STAALAVLL	0.806	STAALGVLM	HLA-A*26:01	Private	>sp P10147 CCL3_HUMAN C-C motif chemokine 3 OS=Homo sapiens OX=9606 GN=CCL3 PE=1 SV=1	9	0.778
STAALAVLL	0.806	STAALGVLM	HLA-A*26:01	Private	>sp P16619 CL3L1_HUMAN C-C motif chemokine 3-like 1 OS=Homo sapiens OX=9606 GN=CCL3L1 PE=1 SV=1	9	0.778
ILGVVLLAIF	0.818	IVGVALLAVF	HLA class I	Private	>sp Q86VB7 C163A_HUMAN Scavenger receptor cysteine-rich type 1 protein M130 OS=Homo sapiens OX=9606 GN=CD163 PE=1 SV=2	10	0.7
ILGVVLLAIF	0.818	IVGVALLAVF	HLA class I	Private	>tr F5GZZ9 F5GZZ9_HUMAN Scavenger receptor cysteine-rich type 1 protein M130 OS=Homo sapiens OX=9606 GN=CD163 PE=1 SV=1	10	0.7
ILGVVLLAIF	0.818	IVGVALLAVF	HLA class I	Private	>tr H0YFM0 H0YFM0_HUMAN Scavenger receptor cysteine-rich type 1 protein M130 (Fragment) OS=Homo sapiens OX=9606 GN=CD163 PE=1 SV=1	10	0.7
ILGVVLLAIF	0.818	IVGVALLAVF	HLA class I	Private	>tr C9JHR8 C9JHR8_HUMAN Scavenger receptor cysteine-rich type 1 protein M130 OS=Homo sapiens OX=9606 GN=CD163 PE=1 SV=1	10	0.7

Table 1: SARS-CoV-2 peptides with high similarity to the human proteome, from genes reported to be involved with severe COVID-19. Note, peptides derived from CD163 were not predicted to bind HLA.

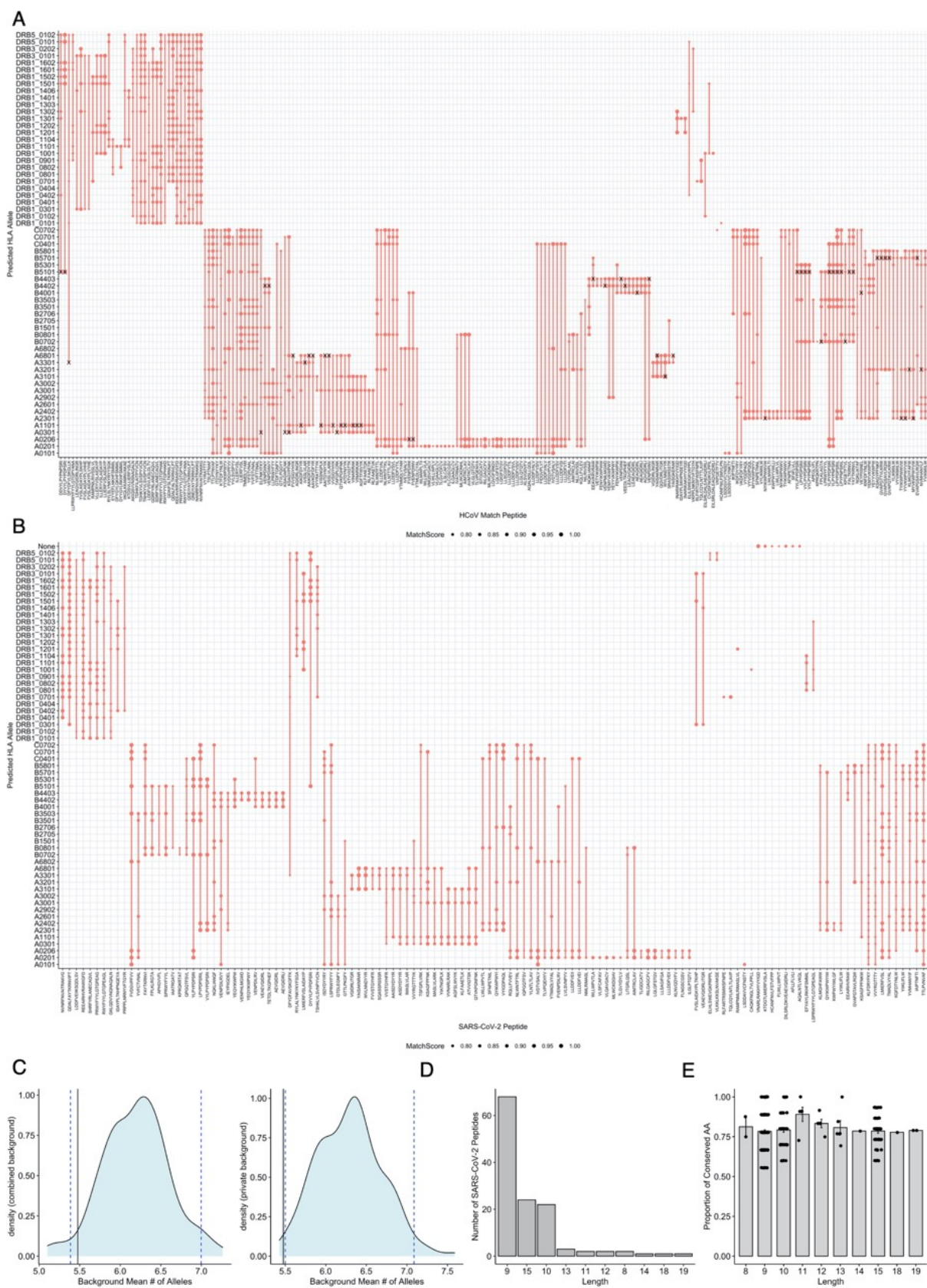
SARS-CoV-2 Peptide	Virus	Protein	MatchScore
AIMTRCLAV	229E,OC43,HKU1,NL63,SARS-CoV,MERS_CoV	replicase polyprotein 1ab,ORF1ab polyprotein,ORF1ab polyprotein,replicase polyprotein 1ab,ORF1ab polyprotein,1AB polyprotein	1,1,0.977,0.977,1,0.977
KLFAAETLK	NL63,229E,HKU1,OC43,SARS-CoV,MERS_CoV	replicase polyprotein 1ab,replicase polyprotein 1ab,ORF1ab polyprotein,ORF1ab polyprotein,ORF1ab polyprotein,1AB polyprotein	0.881,0.857,0.847,0.847,1,0.929
LLDDFVEI	HKU1,229E,NL63,OC43,SARS-CoV	ORF1ab polyprotein,replicase polyprotein 1ab,replicase polyprotein 1ab,ORF1ab polyprotein,ORF1ab polyprotein	0.894,0.871,0.86,0.777,1
LQLGFSTGV	OC43,HKU1,229E,NL63,MERS_CoV,SARS-CoV	ORF1ab polyprotein,ORF1ab polyprotein,replicase polyprotein 1ab,replicase polyprotein 1ab,1AB polyprotein,ORF1ab polyprotein	0.977,0.955,0.809,0.809,1,1
LSDDAVVCFNSTY	229E,HKU1,OC43,SARS-CoV,MERS_CoV	replicase polyprotein 1ab,ORF1ab polyprotein,ORF1ab polyprotein,ORF1ab polyprotein,1AB polyprotein	0.843,0.789,0.789,0.872,0.789
VLQAVGACV	HKU1,OC43,229E,NL63,SARS-CoV,MERS_CoV	ORF1ab polyprotein,ORF1ab polyprotein,replicase polyprotein 1ab,replicase polyprotein 1ab,ORF1ab polyprotein,1AB polyprotein	0.876,0.876,0.795,0.773,1,0.832
VQIDRLITGR	HKU1,229E,NL63,OC43,SARS-CoV,MERS_CoV	spike protein (all)	0.887,0.845,0.845,0.804,1,0.804
YEQYIKWPW	HKU1,OC43,NL63,229E,SARS-CoV	spike protein (all)	0.903,0.873,0.855,0.794,1
YEQYIKWPWY	HKU1,OC43,NL63,SARS-CoV	spike protein (all)	0.913,0.886,0.775,1
YVFTVNAL	229E,NL63,HKU1,OC43,SARS-CoV	replicase polyprotein 1ab,replicase polyprotein 1ab,ORF1ab polyprotein,ORF1ab polyprotein,ORF1ab polyprotein	0.84,0.818,0.809,0.809,1

Table 2: Highly conserved CD8+ T cell peptides across SARS-CoV-2 and HCoV strains, with high population coverage

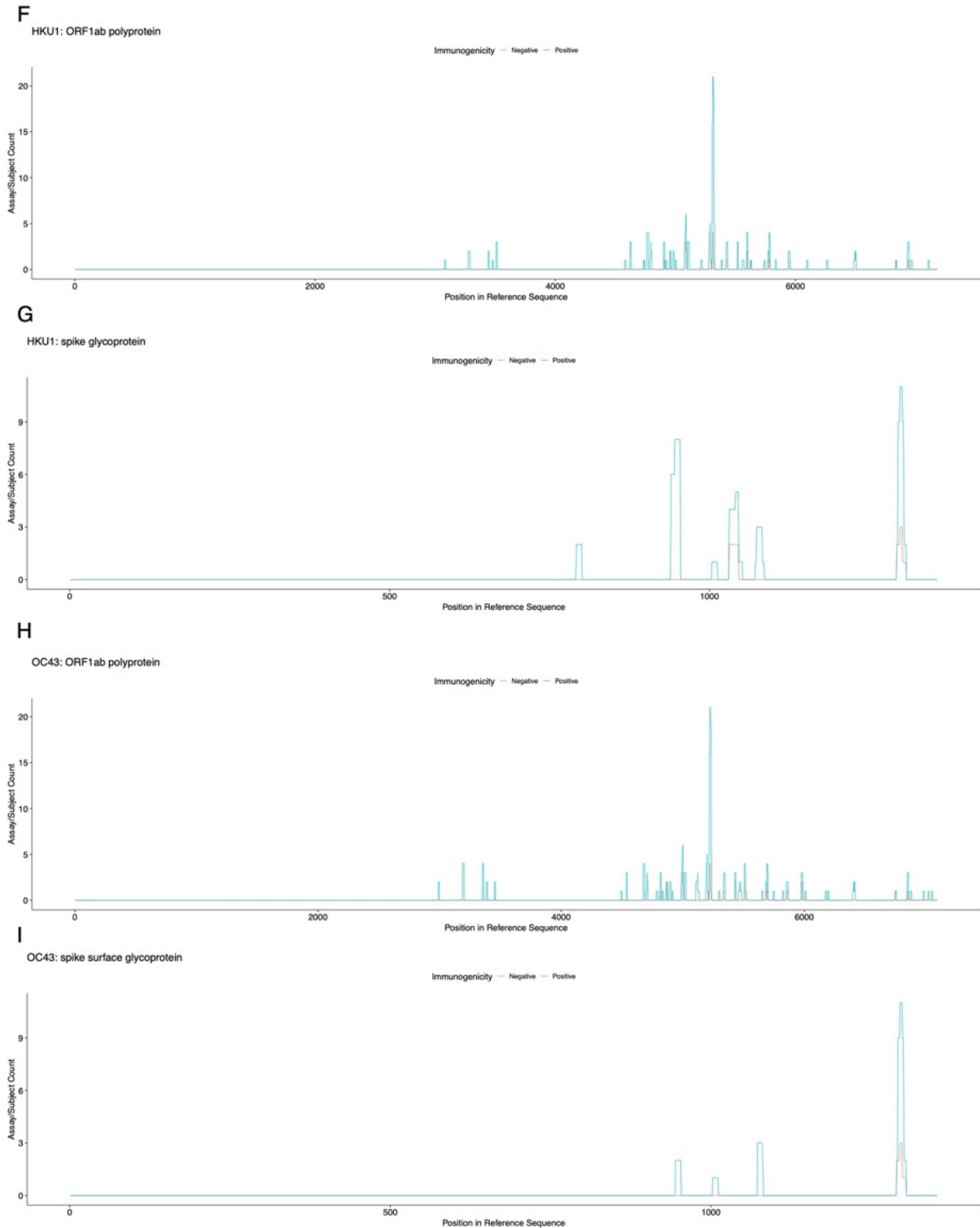
Supplementary



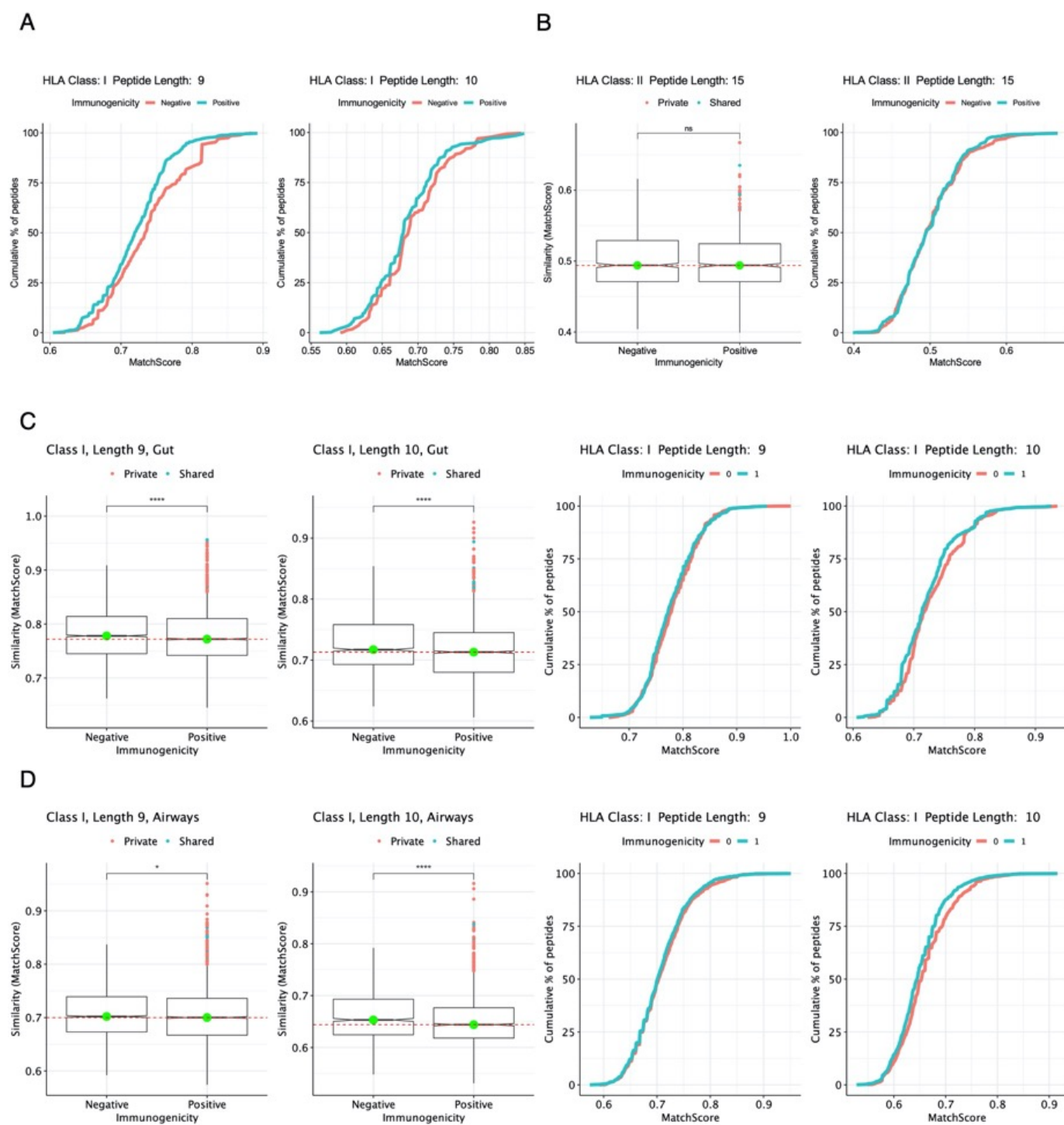
Supplementary Figure 1: A-B) Barplots showing A) the frequency of immunogenic epitopes amongst HLA labels that do not designate a specific allele. B) the frequency of antigen presentation predictions where the correct allele for the SARS-CoV-2 peptide (where known) was identified. C-E) Line plots showing the immunogenic regions of the most immunodominant SARS-CoV-2 proteins, C) membrane glycoprotein, D) nucleocapsid phosphoprotein, E) ‘spike’ surface glycoprotein.



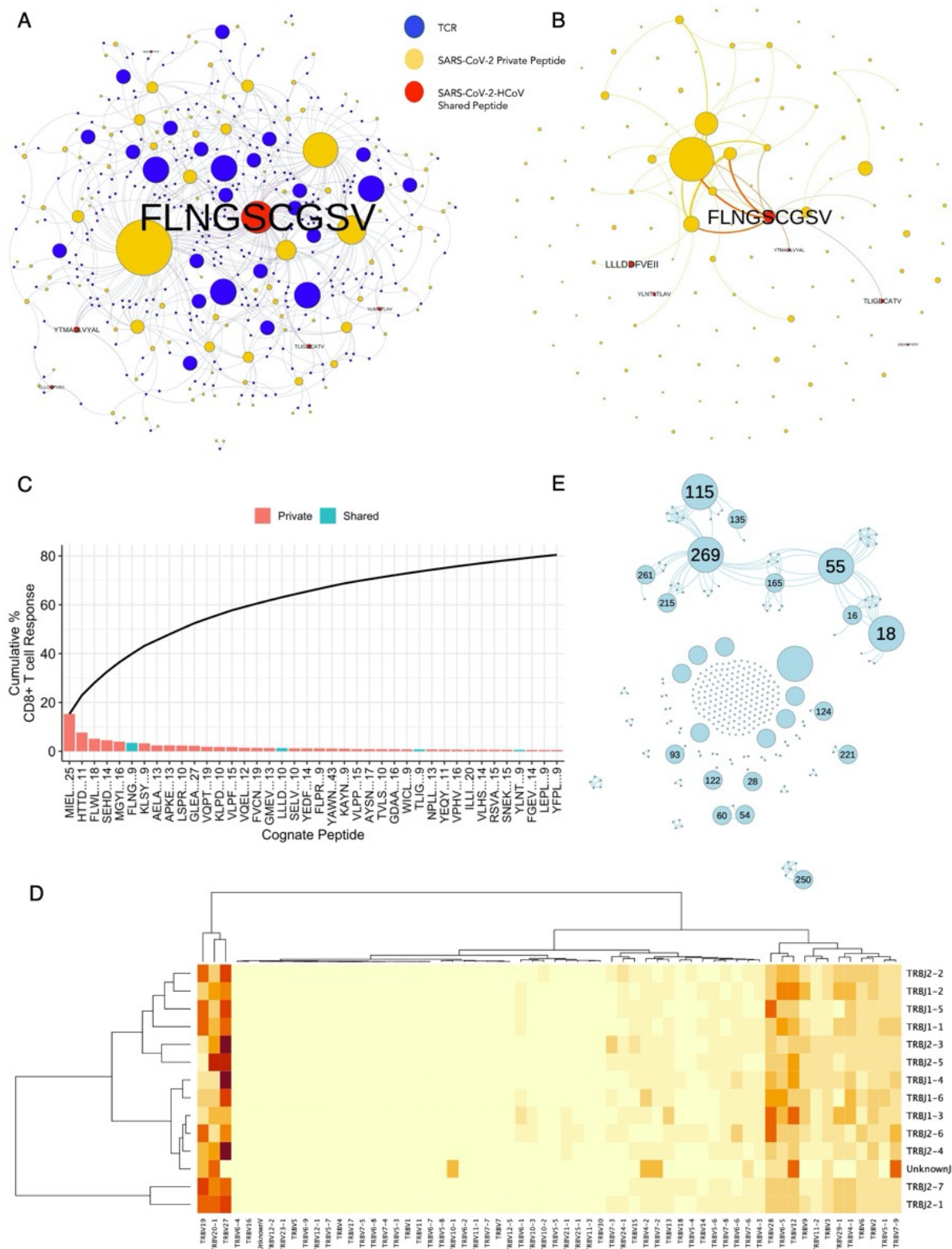
Supplementary Figure 2A-E, continued overleaf.



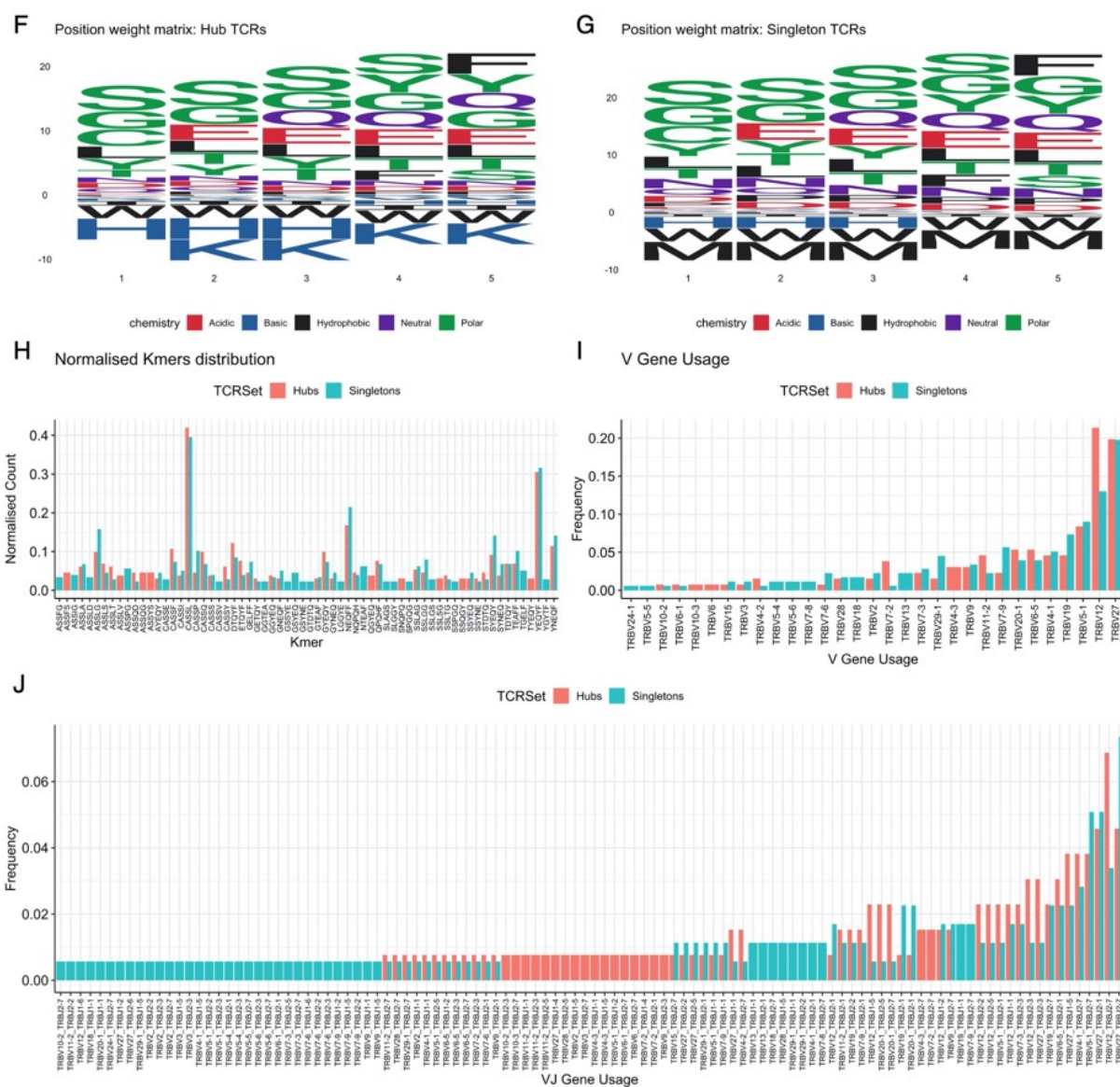
Supplementary Figure 2: A) A dot plot showing predicted HLA alleles for peptides from HCoV with high-similarity to immunogenic SARS-CoV-2 peptides. Size of the point reflects the similarity score to the corresponding SARS-CoV-2 peptide. Lines link predicted alleles. An 'X' indicates - where possible - if there is experimental evidence that the corresponding SARS-CoV-2 peptide binds the predicted HLA allele for the HCoV match. B) a dot plot showing predicted HLA alleles for the 129 shared SARS-CoV-2 immunogenic peptides. C) Density plots showing the mean number of alleles for which shared peptides are predicted to bind (solid black line), compared with a background distribution generated by randomly shuffling the dataset using both private and shared peptides (left) or private peptides only (right). Dashed lines show ± 2 standard deviations from the mean. D) Bar plot showing the distribution of lengths for 129 shared SARS-CoV-2 immunogenic peptides. E) Bar plot showing the distribution of amino acid conservation between the 285 matches and the 126 SARS-CoV-2 shared peptides. F-I Line plots showing the HCoV protein regions which produce immunogenic HCoV-CoV-2 shared peptides from F) HKU1 ORF1ab polyprotein, G) HKU1 spike glycoprotein, H) OC43 ORF1ab polyprotein, I) OC43 spike glycoprotein.



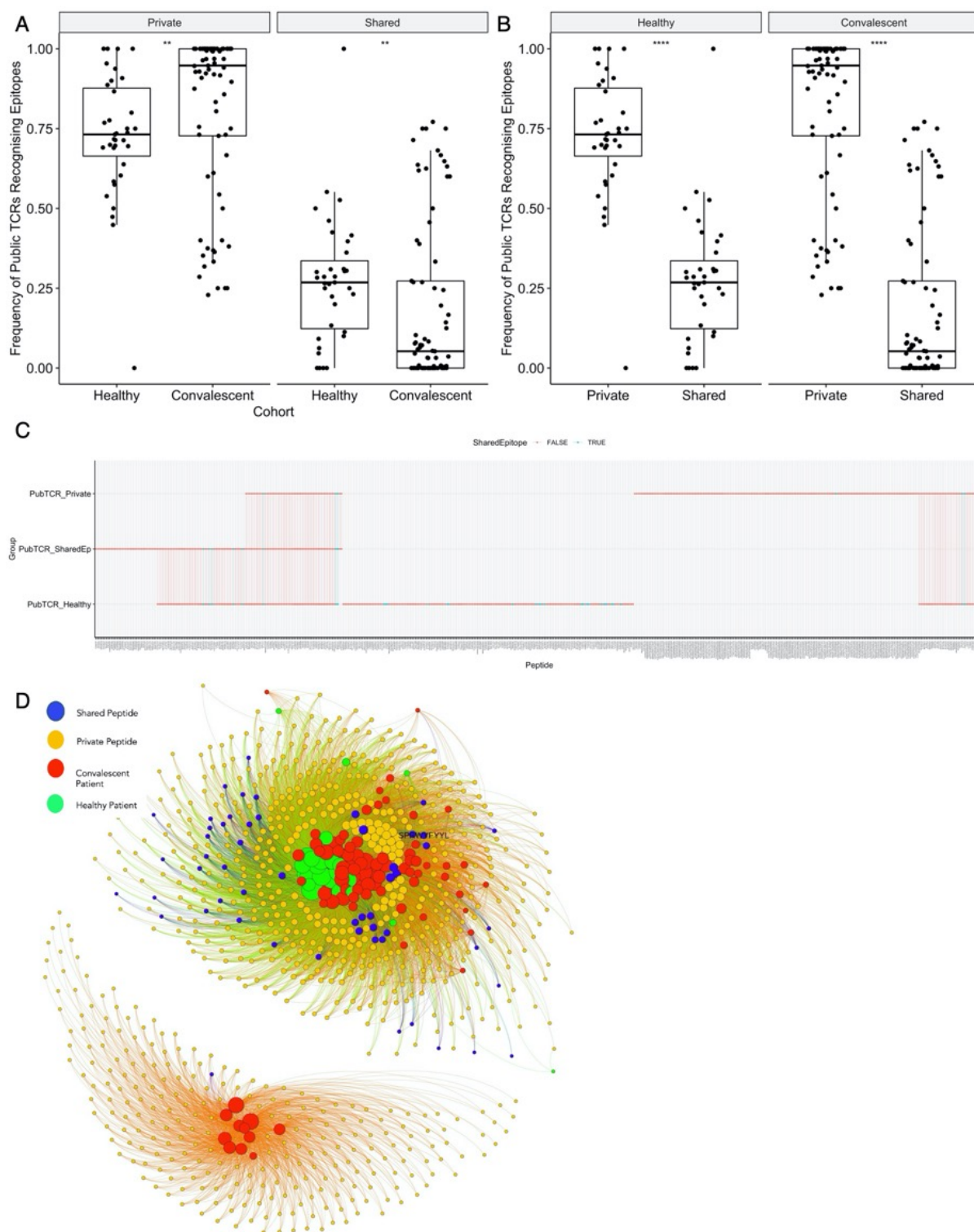
Supplementary Figure 3: A) Empirical cumulative distribution plots for HLA class I peptides of length 9 and 10, showing the cumulative % of peptides to reach a specific MatchScore, color labelled by immunogenicity status. B) A notched boxplot and empirical cumulative distribution plot showing the similarity -as evaluated by the MatchScore- of nonimmunogenic and immunogenic class II SARS-CoV-2 peptides with sequences derived from the human proteome of length 15 (immunogenic n=955, nonimmunogenic n=953 peptides). C-D) Notched boxplots and empirical cumulative distribution plots showing the similarity of nonimmunogenic and immunogenic SARS-CoV-2 peptides with sequences derived from gut C) and airway microbiomes D).



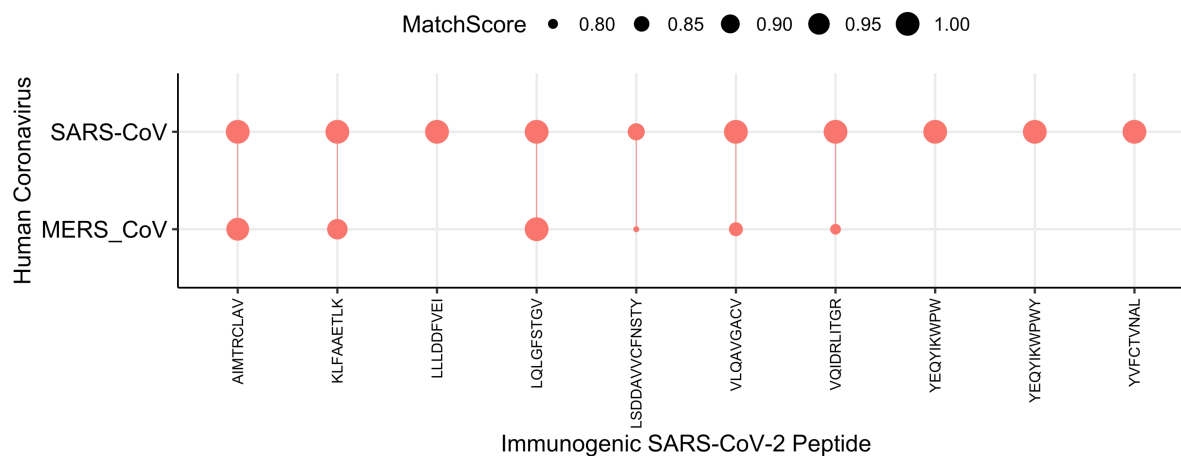
Supplementary Figure 4 A-D, continued overleaf.



Supplementary Figure 4: CD8⁺ T cell cross-reactivity against SARS-CoV-2. A) A bipartite network graph depicting the interactions with SARS-CoV-2 immunogenic peptides (SARS-CoV-2-shared are colored red, -private are coloured yellow) and their cognate TCRs (blue). Node size represents the degree of connectivity. B) A one-mode network graph showing shared and private SARS-CoV-2 peptides. An edge between a peptide node demonstrates that a peptide is recognized by the same TCR. Node size reflects the degree of connectivity. C) A barplot showing the top 80% of the cumulative SARS-CoV-2-specific CD8⁺ T cell response and the peptides recognized, as per the cognate TCR dataset from the IEDB. D) A heatmap showing common V and J gene usage for SARS-CoV-2 specific TCR β sequences. E) A one-mode network graph showing the common specificity of SARS-CoV-2 specific TCRs. Each node is a TCR, and an edge reflects whether two TCRs recognize the same peptide. Node size reflects the number of peptides recognized by a TCR. F) A sequence logo plot visualizing the position weight matrix for CDR3b 5mers in the IEDB dataset, for “Hub TCRs”, those with considerable common-specificity. G) A sequence logo plot visualizing the position weight matrix for CDR3b 5mers in the IEDB dataset, for “Singletons”, those recognizing only one unique SARS-CoV-2 peptide H) A barplot contrasting the Kmer distribution for “Hub” and “Singleton” TCRs. The count is normalized by the number of TCRs in each group (Hub or Singletons). I) A barplot contrasting the V gene usage of “Hub” and “Singleton” TCRs. Y axis shows the frequency to which the V gene is used in each group. J) A barplot contrasting the V-J gene usage of “Hub” and “Singleton” TCRs. Y axis shows the frequency to which the V-J gene combination is used in each group.



Supplementary Figure 5: A) Boxplot showing for each patient, the frequency of public TCRs recognizing SARS-CoV-2 private or SARS-CoV-2-HCoV shared peptides, grouped by Private/Shared epitopes and contrasting patient cohort. B) Boxplot showing for each patient, the frequency of public TCRs recognizing SARS-CoV-2 private or SARS-CoV-2-HCoV shared peptides, grouped by patient cohort and contrasting Private/Shared epitopes. C) Line graph showing the peptides recognized by public TCRs in different patient cohorts: Healthy patients (PubTCR_Healthy) (n=12), COVID convalescent patients whose public TCRs predominately recognize Private (PubTCR_Private) or Shared (PubTCR_SharedEp) peptides (both n=12). D) A bipartite graph showing the common-specificity of private TCRs recognizing SARS-CoV-2 shared and private peptides. An edge between a patient and a peptide is observed if a patient possesses a private TCR recognizing that peptide.



Supplementary Figure 6: A) A dot and line plot showing each SARS-CoV-2 peptide on the x-axis. A dot shows a high similarity match to MERS or SARS-CoV. The size of each point reflects the MatchScore, i.e the similarity metric.

SARS-CoV-2 Peptide	Qualitative_Measure	HLA_Allele	HitSeq	MatchScore	SharedEpitope	Length	SYMBOL	HitSeq_Binder	AASeq_Similarity
RQLLFVVEV	Positive	HLA class I	RQLLFVVDI	0.892	TRUE	9	GTPBP10	TRUE	0.778
YIATNGPLK	Positive	HLA-A*11:01	YIATQGPKL	0.884	TRUE	9	PTPRG	TRUE	0.889
LLSAGIFGA	Positive	HLA class I	LITAGIFGA	0.871	TRUE	9	SLC12A3	FALSE	0.778
ALNTLVKQL	Positive-Low	HLA-A*02:01	SLNTLLKQL	0.854	FALSE	9	CCDC169	TRUE	0.778
GLTVLPPLL	Positive,Positive	HLA-A*02:01,HLA-A*02:01	GLTVLPALL	0.851	FALSE	9	SLC2A4	TRUE	0.889
SSRGTSPAR	Positive	HLA class I	SSRDTSPAR	0.841	FALSE	9	CLASP1	TRUE	0.889
LLFNKVTLA	Positive,Positive	HLA class I,HLA-A*02:01	LLYNKMTLA	0.837	FALSE	9	MME	TRUE	0.778
QEILGTVSW	Positive,Positive	HLA-B*44:03,HLA-B*44:03	QEVLSGMSW	0.833	FALSE	9	NWD1	TRUE	0.667
SPRRARSVA	Positive	HLA-B*07:02	SPRRRSIS	0.833	FALSE	9	SRSF7	TRUE	0.667
AIMTRCLAV	Positive	HLA class I	AVLTRCLVV	0.828	TRUE	9	LSM14B	TRUE	0.667
LALLLDRL	Positive,Positive	HLA-A*02:01,HLA-A*02:01	LALALLDRI	0.821	FALSE	9	NUP205	TRUE	0.778
FLLPSLATV	Positive,Positive,Positive	HLA-A*02:01,HLA-A*02:01,HLA-A*02:01	FLIPSLAAI	0.819	FALSE	9	OR2AG1	TRUE	0.667
RLFRKSNLK	Positive,Positive,Positive	HLA class I,HLA-A*31:01,HLA-A*03:01	RLFKSNIR	0.818	FALSE	9	GPR87	TRUE	0.667
NLNESLIDL	Positive,Positive-Low,Positive	HLA class I,HLA-A*02:01,HLA-A*02:01	NVNQSLDL	0.814	FALSE	9	FTHL17	FALSE	0.667
SELLTPLGI	Positive	HLA-B*40:01	SELLKPLGL	0.814	FALSE	9	MBD4	TRUE	0.778
EAFEKMVSL	Positive-Low	HLA-B*08:01	EEFEKLVSL	0.81	FALSE	9	BHLHB9	TRUE	0.778
RLDKVEAEV	Positive	HLA class I	RLDSMEAEV	0.81	FALSE	9	CCDC151	TRUE	0.778
STAALGVLM	Positive	HLA-A*26:01	STAALAVLL	0.806	FALSE	9	CCL3	TRUE	0.778
STAALGVLM	Positive	HLA-A*26:01	STAALAVLL	0.806	FALSE	9	CCL3L1	TRUE	0.778
VPHISRQRL	Positive	HLA-B*07:02	VPHVSKERI	0.804	FALSE	9	KNL1	TRUE	0.556
EYVSQPFLM	Positive	HLA-A*24:02	EYIEKPFM	0.8	FALSE	9	TTL7	TRUE	0.667
LLDRLNQL	Positive,Positive-Low,Positive-Low,Positive,Positive	HLA-A*02:01,HLA-A*02:01,HLA-A*02:01,HLA-A*02:01,HLA-A*02:01,HLA-A*02:01	LLDRVNDL	0.8	FALSE	9	DNAH5	TRUE	0.778
STNVTIATY	Positive,Positive	HLA-A*01:01,HLA-A*32:01	STHVTISTY	0.8	FALSE	9	RIOX2	TRUE	0.778
VLKGVKLHY	Positive	HLA-A*29:02	VLKSKLHF	0.796	FALSE	9	KIF14	TRUE	0.778
GMSRIGMEV	Positive-Low,Positive,Positive,Positive	HLA-A*02:01,HLA-A*02:01,HLA-A*02:01,HLA-A*02:01	GMSRLGEEV	0.795	FALSE	9	EXD2	TRUE	0.778
MASLVLARK	Positive	HLA-A*68:01	MASLIVARQ	0.795	TRUE	9	SLC24A3	FALSE	0.667
MASLVLARK	Positive	HLA-A*68:01	MASLIVARQ	0.795	TRUE	9	SLC24A4	FALSE	0.667
VLQAVGACV	Positive	HLA class I	VLEAVGACL	0.795	TRUE	9	ATM	TRUE	0.667
KQEILGTVSW	Positive,Positive	HLA-B*44:02,HLA-B*44:03	KQEVLSGMSW	0.849	FALSE	10	NWD1	TRUE	0.7
SQASSRSSR	Positive	HLA class I	SRSSRSR	0.837	TRUE	10	CLASRP	TRUE	0.8
SQASSRSSR	Positive	HLA class I	SRASSRASSR	0.837	TRUE	10	GJA1	TRUE	0.8
IVGVALLAVF	Positive	HLA class I	ILGVLLAIF	0.818	FALSE	10	CD163	FALSE	0.7
TNVLEGSVAY	Positive	HLA-B*35:01	TNLEGAFAV	0.804	FALSE	10	AUP1	TRUE	0.7
IEYPIIGDEL	Positive	HLA-B*40:01	VEYPLIEDEL	0.796	TRUE	10	DNAH11	TRUE	0.7
VENPDILRVY	Positive	HLA-B*44:02	VESPKILRVY	0.792	TRUE	10	F11	TRUE	0.8
ILPDPKPSK	Positive	HLA class I	ILPDPDDPSK	0.789	FALSE	10	ZNF592	TRUE	0.8
KVAGFAKFLK	Positive	HLA-A*11:01	SVAGFSRFLK	0.784	FALSE	10	FBXL4	TRUE	0.7
FTISVTTEIL	Positive	HLA class I	FTIRVTSEVL	0.783	FALSE	10	PNPT1	TRUE	0.7

Table S1: Peptides identified with high similarity to human proteome. HitSeq shows the match from the self proteome. SharedEpitope shows whether the peptide is a sCoV-2-HCoV shared peptide or not. SYMBOL shows the gene from which the hit peptide is derived. HitSeq_Binder shows whether the hit peptide is predicted to bind an HLA allele. AASeq_Similarity shows the proportion of amino acids conserved between the SARS-CoV-2 peptide and the human proteome match.

Experiment	Cohort	Age	Gender	HLA-B	TCR	Peptide
eAV91	Healthy (No known exposure)	31	M	B*07:02	CASSELPGPPGEQYF+TCRBV0201+TCRBJ0207	LSPRWYFYI
eAV91	Healthy (No known exposure)	31	M	B*07:02	CASSELPGPPGEQYF+TCRBV0201+TCRBJ0207	SPRWYFYI
eXL31	Healthy (No known exposure)	28	M	B*07:02	CASTLAGGPYNEQFF+TCRBV0501+TCRBJ0201	LSPRWYFYI
eXL31	Healthy (No known exposure)	28	M	B*07:02	CASTLAGGPYNEQFF+TCRBV0501+TCRBJ0201	SPRWYFYI

Table S2: Two previously reported private TCRs are identified in additional HLA-B*07:02+ individuals at beta chain resolution, indicating these are cross-reactive public TCRs.

Supplementary Data File description

- 1: Data file containing each of the 126 SARS-CoV-2 peptides which map to 285 targets from HCoV.
- 2: Data file containing the full set of SARS-CoV-2 peptides which have high similarity to the human proteome.
- 3: Data file containing a summary of the peptides recognized by public TCRs in the PubTCR-SharedEp and PubTCR-Private groups, supplemented with those recognized in a sampled set of healthy patients. File contains a 1 or a 0 demonstrating whether a public TCR in each group of patients recognizes the peptide.
- 4: Data file containing the information from data file 3, but also includes information pertaining to whether key class I HLA alleles are observed in a patient with a public TCR recognizing each peptide.
- 5: Data file containing cohort information regarding the peptides most commonly recognized by private TCRs in the MIRA dataset.
- 6: Data file reporting the private TCRs which recognize the sCoV-2-HCoV shared peptides shown in data file 5.
- 7: Data file containing information from Figure 6, exhibiting the peptides with hits to ≥ 3 HCoV strains. File provides detailed information regarding SARS-CoV-2 peptide and the corresponding hit to HCoV.
- 8: Data file contains every TCR from the IEDB and the MIRA datasets which we used in this analysis, mapped to the recognized SARS-CoV-2 peptide.

