

1 Transmembrane helices are an overlooked and
2 evolutionarily conserved source of major
3 histocompatibility complex class I and II epitopes

4 Richèl J.C. Bilderbeek¹, Maxim Baranov¹, Geert van den
5 Bogaart¹, and Frans Bianchi¹

6 ¹GBB, University of Groningen, Groningen, The Netherlands

7 May 6, 2021

8

Abstract

9 Cytolytic T cell responses are predicted to be biased towards mem-
10 brane proteins. The peptide-binding grooves of most haplotypes of histo-
11 compatibility complex class I (MHC-I) are relatively hydrophobic, there-
12 fore peptide fragments derived from human transmembrane helices (TMHs)
13 are predicted to be presented more often as would be expected based
14 on their abundance in the proteome. However, the physiological reason
15 of why membrane proteins might be over-presented is unclear. In this
16 study, we show that the over-presentation of TMH-derived peptides is
17 general, as it is predicted for bacteria and viruses and for both MHC-
18 I and MHC-II. Moreover, we show that TMHs are evolutionarily more
19 conserved, because single nucleotide polymorphisms (SNPs) are present
20 relatively less frequently in TMH-coding chromosomal regions compared
21 to regions coding for extracellular and cytoplasmic protein regions. Thus,
22 our findings suggest that both cytolytic and helper T cells respond more to
23 membrane proteins, because these are evolutionary more conserved. We
24 speculate that TMHs therefore are less prone to escape mutations that
25 enable pathogens to evade T cell responses.

26 **Keywords:** antigen presentation, membrane proteins, bioinformatics, adap-
27 tive immunity, transmembrane domain, transmembrane helix, epitopes, T lym-
28 phocyte, MHC-I, MHC-II, evolutionary conservation

29 Abbreviations

Abbreviation	Full
MAP	Membrane-associated protein
TMH	Transmembrane helix
TMP	Transmembrane protein

30 1 Introduction

31 Our immune system fights diseases and infections from pathogens, such as fungi,
32 bacteria or viruses. An important part of the acquired immune response, that
33 develops specialized and more specific recognition of pathogens than the in-
34 nate immune response, are T cells which recognize peptides, called epitopes,
35 derived from antigenic proteins presented on Major Histocompatibility Com-
36 plexes (MHC) class I and II on the cell surface.

37 The MHC proteins are heterodimeric complexes encoded by the HLA (Hu-
38 man Leukocyte Antigens) genes. In humans, the peptide binding groove of
39 MHC-I is made by only the alpha subunit. There are three classical forms of
40 MHC-I, hallmarked by a highly polymorphic alpha chain called HLA-A, HLA-B
41 and HLA-C, that all present epitopes to cytolytic T cells. For MHC-II, both the
42 alpha and the beta chains contribute to the peptide binding groove. There are
43 three classical forms of MHC-II as well, called HLA-DR, HLA-DQ and HLA-DP,
44 that all present epitopes to helper T cells. Each MHC complex can present a
45 subset of all possible peptides. For example, HLA-A and HLA-B have no over-
46 lap in which epitopes they bind [1]. Moreover, the HLA genes of humans are
47 highly polymorphic, with hundreds to thousands of different alleles, and each
48 different HLA allele is called an MHC haplotype and presents a different subset
49 of peptides [2].

50 Humans mostly express two haplotypes per MHC form, one from the parental
51 and one from the maternal chromosome, and therefore an individual's immune
52 system detects only a fraction of all possible peptide fragments. However, at
53 the population level, the coverage of pathogenic peptides that are detected is
54 very high, because of the highly polymorphic MHC genes. It is therefore be-
55 lieved that MHC polymorphism improves immunity at the population level, as
56 mutations in a protein that disrupt a particular MHC presentation at the in-

57 individual level, so-called escape mutations, will not affect MHC presentation for
58 all haplotypes present in the population [3].

59 Many studies are aimed at identifying the repertoire of epitopes that are
60 presented in any MHC haplotype and determining which epitopes will result in
61 an immune response, as this will for instance aid the design of vaccines. These
62 studies have led to the development of prediction algorithms that allow for very
63 reliable *in silico* predictions of the binding affinities of peptides [4, 5, 6]. For
64 example, [6] found that, of the 432 peptides that were predicted to bind to
65 MHC, 86% were experimentally confirmed to do so.

66 Using these prediction algorithms, we recently predicted that peptides de-
67 rived from transmembrane helices (TMHs) will be more frequently presented
68 by MHC-I than expected based on their abundance [7]. Moreover, we showed
69 that some well-known immunodominant peptides stem from TMHs. This over-
70 presentation is attributed to the fact that the peptide-binding groove of most
71 MHC-I haplotypes is relatively hydrophobic, and therefore hydrophobic TMH-
72 derived peptides have a higher affinity to bind than their soluble hydrophilic
73 counterparts.

74 TMHs are hydrophobic as they need to span the hydrophobic lipid bilayer
75 of cellular membranes. They consist of an alpha helix of, on average, 23 amino
76 acids in length. TMHs can also be predicted with high accuracy from a protein
77 sequence by bioinformatics approaches [8, 9, 10, 11, 12, 13], for example, [11]
78 found that, from 184 transmembrane proteins (TMPs) with known topology,
79 80% of the TMH predictions replicated this finding.

80 TMHs are common structures in the proteins of humans and microbes. Dif-
81 ferent TMH prediction tools estimate that 15-39% of all proteins in the human
82 proteome contain at least one TMH [14]. However, the physiological reason
83 why peptides derived from TMHs would be presented more often than peptides

84 stemming from soluble (i.e., extracellular or cytoplasmic) protein regions is un-
85 known. We hypothesized that the presentation of TMH residues is evolutionary
86 selected for, because TMHs are less prone to undergo escape mutations. One
87 reason to expect such a reduced variability (and hence evolutionary conserva-
88 tion) in TMHs, is that these are restricted in their evolution by the functional
89 requirement to span a lipid bilayer. Due to this requirement, many of the amino
90 acids genuinely present in TMHs are limited to the ones with hydrophobic side
91 chains [15, 16]. Therefore, we speculated that the TMHs of pathogens might
92 have a lower chance to develop escape mutations, as many mutations will result
93 in a dysfunctional TMH and render the protein inactive.

94 This study had two objectives. First, we aimed to generalize our findings
95 by predicting the presentation of peptides from different kingdoms of life and
96 for both MHC-I and -II. From these *in silico* predictions, we conclude that
97 TMH-derived epitopes are presented more often than expected by chance, in a
98 human, viral and bacterial proteome, and for most haplotypes of both MHC-I
99 and II. We confirmed the presentation of TMH-derived peptides by re-analysis of
100 peptide elution studies. Second, we tested our hypothesis that TMHs are more
101 evolutionary conserved than soluble protein regions. Our analysis of human
102 single nucleotide polymorphisms (SNPs) showed that random point mutations
103 are indeed less likely to occur within TMHs. These findings strengthen the
104 emerging notion that TMHs are important for the T cell branch of the adaptive
105 immune system, and hence are of overlooked importance in vaccine development.

106 2 Methods

107 2.1 Predicting TMH epitopes

108 To predict how frequently epitopes overlapping with TMHs are presented, a
109 similar analysis strategy was applied as described in [7] for several haplotypes
110 of both MHC-I and MHC-II, and for a human, viral and bacterial proteome.
111 To summarize, for each proteome, all possible 9-mers (for MHC-I) or 14-mers
112 (MHC-II) were derived. For each of these peptides, we determined if it over-
113 lapped with a predicted TMH and if it was predicted to bind to each haplotype.

114 For MHC-I, 9-mers were used, as this is the length most frequently pre-
115 sented in MHC-I and was used in our earlier study [7]. For MHC-II, 14-
116 mers were used, as these are the most frequently occurring epitope length [17].
117 A human (UniProt ID UP000005640.9606), viral (SARS-CoV-2, UniProt ID
118 UP000464024) and bacterial (*Mycobacterium tuberculosis*, UniProt ID UP000001584)
119 reference proteome was used. TMHMM [8] was used to predict the topology of the
120 proteins within these proteomes. To predict the affinity of an epitope to a cer-
121 tain MHC haplotype, `EpitopePrediction` [7] for MHC-I and `MHCnuggets` [18]
122 for MHC-II was used. The 13 MHC-I haplotypes used in this study are the same
123 as used in the previous study [7]. For MHC-II, haplotypes were selected with a
124 phenotypic frequency of at least 14% in the human population [19], resulting in
125 21 MHC-II haplotypes.

126 In previous work, it was found that the over-presentation of TMH-derived
127 peptides can be explained from the hydrophobicity of the MHC-I binding cleft
128 [7]. Here, a similar analysis was applied, by correlating the percentage of pre-
129 dicted TMH-derived epitopes versus the mean hydrophobicity of all peptides.

130 This study differs in one important aspect from our previous work [7]. The
131 definition of a binder differs from [7]: in the current study, a peptide is called a
132 binder if, for a certain haplotype, any of its 9-mer or 14-mer peptides have an

133 IC50 value in the lowest 2% of all peptides within a *proteome* (see supplementary
134 Tables 4 and 5 for values), whereas the previous study defined a binder as having
135 an IC50 in the lowest 2% of the peptides within a *protein*. This revised definition
136 precludes bias of proteins that give rise to no or only very few MHC epitopes.
137 To verify that the results are similar, a side by side comparison was performed
138 shown in the supplementary materials.

139 **2.2 Peptide elution studies**

140 To obtain experimental evidence that epitopes derived from TMHs are pre-
141 sented in MHC, peptide elution studies for MHC-I [20] and MHC-II [17] were
142 reanalyzed. For each of the detected epitopes, its possible location(s) in a hu-
143 man reference proteome, with UniProt ID UP000005640.9606, was mapped.
144 For the epitopes that were present in the proteome exactly once, the topology
145 of the proteins in which these epitopes were located was predicted using both
146 TMHMM [8] and PureseqTM [13]. From this topology, we determined if the epitope
147 overlapped with a TMH.

148 The full analysis can be found at [https://github.com/richelbilderbeek/](https://github.com/richelbilderbeek/bbbq_article_issue_157)
149 [bbbq_article_issue_157](https://github.com/richelbilderbeek/bbbq_article_issue_157).

150 **2.2.1 Evolutionary conservation of TMHs**

151 To determine the evolutionary conservation of TMHs, human single nucleotide
152 polymorphisms (SNPs) were first collected that resulted in a single amino acid
153 substitution, and we then determined if this substitution occurred within a
154 predicted TMH or not.

155 As a data source, multiple NCBI (<https://www.ncbi.nlm.nih.gov/>) databases
156 were used: the *dbSNP* [21] database, which contains 650 million catalogued non-
157 redundant humane variations (called RefSNPs, <https://www.ncbi.nlm.nih.gov/>).

158 [gov/snp/docs/RefSNP_about/](https://www.ncbi.nlm.nih.gov/snp/docs/RefSNP_about/)), and the databases *gene* (for gene names, [22])
159 and *protein* (for proteins sequences, [23]).

160 The first query was a call to the *gene* database for the term 'membrane
161 protein' (in all fields) for the organism *Homo sapiens*. This resulted in 1,077
162 gene IDs (on December 2020). The next query was a call to the *gene* database
163 to obtain the gene names from the gene IDs. Per gene name, the *dbSNP* NCBI
164 database was queried for variations associated with the gene name. As the
165 NCBI API constrains its users to three calls per second (to assure fair use), we
166 had to limit the extent of our analysis.

167 The number of SNPs was limited to the first 250 variations per gene, resulting
168 in ≈ 61 k variations. Only variations that result in a SNP for a single amino acid
169 substitution were analyzed, resulting in ≈ 38 k SNPs. The exact amounts can be
170 found in the supplementary materials, Tables 9 and 10.

171 SNPs were picked based on ID number, which is linked to their discovery
172 date. To verify that these ID numbers are unrelated to SNP positions, the
173 relative positions of all analyzed SNPs in a protein were determined. This
174 analysis showed no positional bias of the SNPs, as shown in supplementary
175 figure 15.

176 Per SNP, the *protein* NCBI database was queried for the protein sequence.
177 For each protein sequence, the protein topology was determined using **Pureseq**™.
178 Using these predicted protein topologies, the SNPs were scored to be located
179 within or outside TMHs.

180 **3 Results**

181 **3.1 TMH-derived peptides are predicted to be over-presented** 182 **in MHC-I**

183 Figure 1A shows the predicted presentation of TMH-derived peptides in MHC-
184 I, for a human, viral and bacterial proteome. Per MHC-I haplotype, it shows
185 the percentage of binders that overlap with a TMH with at least one residue.
186 The horizontal line shows the expected percentage of TMH-derived epitopes
187 that would be presented, if TMH-derived epitopes would be presented just as
188 likely as epitopes derived from soluble regions. For 11 out of 13 MHC-I hap-
189 lotypes, TMH-derived epitopes are predicted to be presented more often than
190 the null expectation, for a human and bacterial proteome. For the viral pro-
191 teome, 12 out of 13 haplotypes present TMH-derived epitopes more often than
192 expected by chance. The extent of the over-presentation between the different
193 haplotypes is similar for the probed proteomes, which strengthens our previous
194 conclusion [7] that the hydrophobicity of the MHC-binding groove is the main
195 factor responsible for the predicted over-presentation of TMH-derived peptides.

196 **3.2 TMH-derived peptides are predicted to be over-presented** 197 **in MHC-II**

198 We next wondered if the over-representation of TMH-derived peptides would
199 also be present for MHC-II. Figure 1A shows the percentages of MHC-II epitopes
200 predicted to be overlapping with TMHs for our human, viral and bacterial
201 proteomes. We found that TMH-derived peptides are over-presented in all of
202 the 21 MHC-II haplotypes, for a human, bacterial and viral proteome, except
203 for HLA-DRB3*0101 in *M. tuberculosis*. See supplementary Table 8 for the exact
204 TMH and epitope counts.

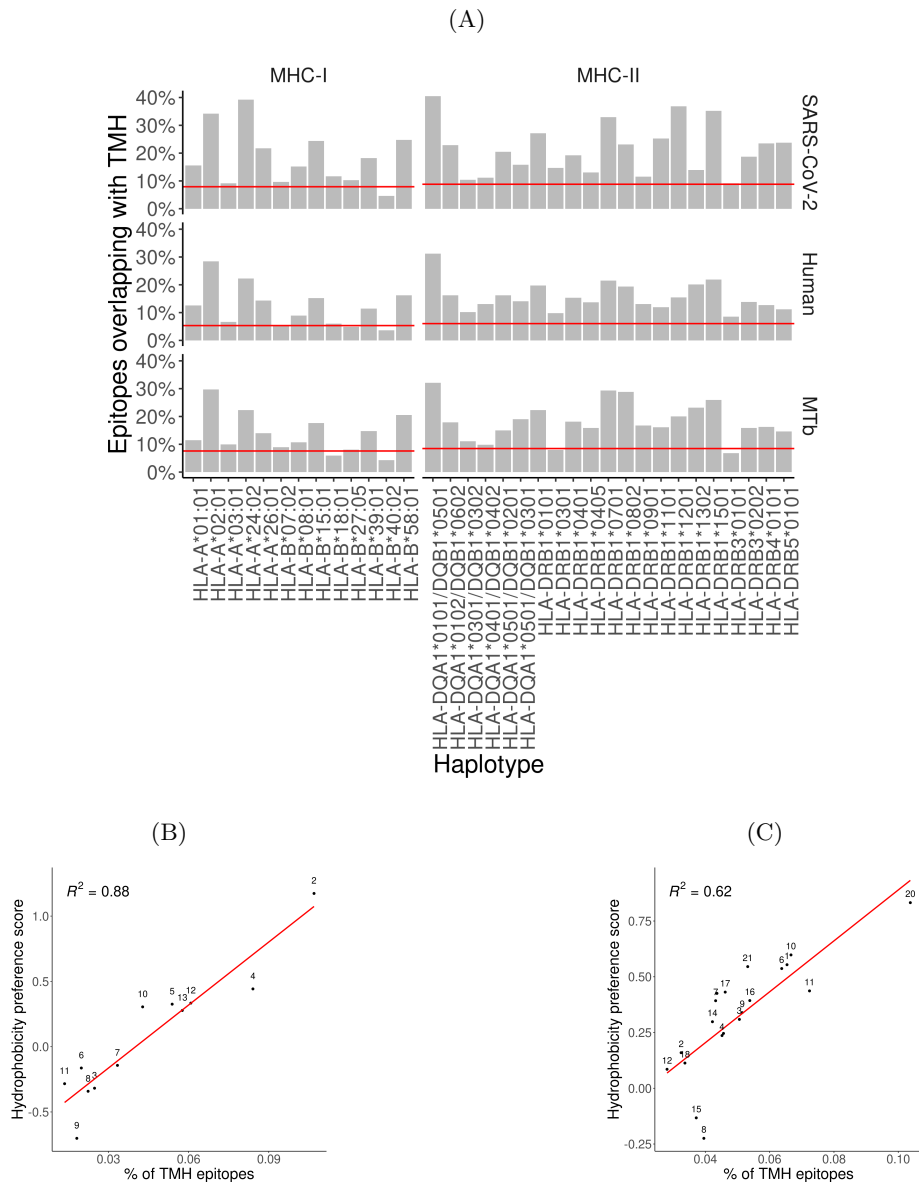


Figure 1: Over-presentation of TMH-derived epitopes on most MHC-I and -II haplotypes (A) The percentage of epitopes for MHC-I and -II haplotypes that are predicted to overlap with TMHs for the proteomes of SARS-CoV-2 (top row), human (middle row) and *M. tuberculosis* (bottom row). The red lines indicate the percentages as expected by chance. See supplementary Tables 7 and 8 for the exact TMH and epitope counts. **(B-C)** Correlation between the percentages of predicted TMH-derived epitopes and the hydrophobicity score of all predicted epitopes for MHC-I **(B)** and MHC-II haplotypes **(C)**. Red curve: linear regression analysis. Labels are shorthand for the HLA haplotypes, see the supplementary Table 6 for the names.

205 **3.3 The over-presentation of TMH-derived peptides is caused**
206 **by the hydrophobicity of the MHC peptide binding**
207 **groove**

208 For MHC-I, we previously showed that the over-presentation of TMH-derived
209 peptides is caused by the hydrophobicity of the peptide binding grooves [7]. Fig-
210 ures 1B and 1C show the extent of over-presentation of TMH-derived epitopes
211 as a function of the hydrophobicity preference score for the different haplo-
212 types. An assumed linear correlation explains 88% of the variability in MHC-
213 I. For MHC-II, 62% of the variability is explained by hydrophobicity. This
214 strengthens our previous finding [7] and indicates that TMH-derived peptides
215 are over-presented because the peptide binding grooves of most MHC-I and -II
216 haplotypes are relatively hydrophobic.

217 **3.4 Experimental validation of presentation of TMH-derived**
218 **peptides**

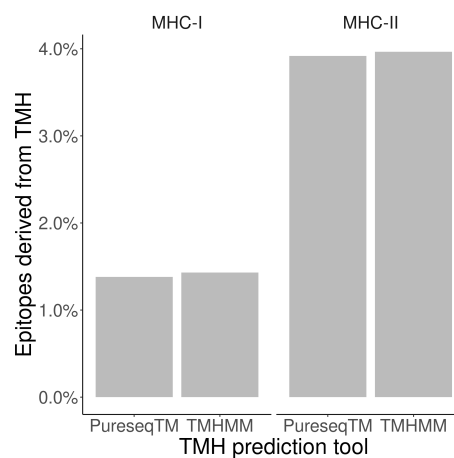


Figure 2: **Robust prediction that TMH epitopes are presented *in vivo*.** Bars show the percentage of peptides obtained from elution studies that is derived from a TMH, for MHC-I and -II, for two TMH prediction tools.

219 To obtain experimental confirmation that peptides stemming from TMHs are
220 presented in MHC-I and MHC-II, two peptide elution studies were reanalyzed:
221 For MHC-I, peptides presented *in vivo* by the (humane) haplotypes HLA-A
222 and B were sequenced [20], for MHC-II these were haplotypes DQ2.5, DQ2.2,
223 and DQ7.5 [17]. Figure 2 shows the percentages of epitopes derived from TMHs
224 found in the MHC-I and MHC-II elution studies, for the two topology prediction
225 tools TMHMM [8] and PureseqTM [13]. Regardless of the prediction tool, at least
226 100 epitopes were predicted to be derived from a TMH for each condition.
227 From these findings, it is robustly predicted that epitopes derived from TMHs
228 are presented in both MHC-I and MHC-II. See the supplementary Table 3 for
229 the exact values.

230 **3.5 Human TMHs are evolutionarily conserved**

231 We addressed the question whether there is an evolutionary advantage in pre-
232 senting TMHs. We determined the conservation of TMHs by comparing the
233 occurrences of SNPs located in TMHs or soluble protein regions for the genes
234 coding for membrane proteins. We obtained 911 unique gene names associated
235 with the phrase 'membrane protein', which are genes coding for both membrane-
236 associated proteins (MAPs, which have no TMH) and transmembrane proteins
237 (TMPs, which have at least one TMH). These genes are linked to 4,780 pro-
238 tein isoforms, of which 2,553 are predicted to be TMPs and 2,237 proteins are
239 predicted to be MAPs. We obtained 37,630 unique variations, of which 9,621
240 are SNPs that resulted in a straightforward amino acids substitution, of which
241 6,062 were located in predicted TMPs. See supplementary Tables 9 and 10 for
242 the detailed numbers and distributions of SNPs.

243 Per protein, we calculated two percentages: (1) the percentage of the total
244 protein predicted to be TMHs, and (2) the percentage of SNPs located within

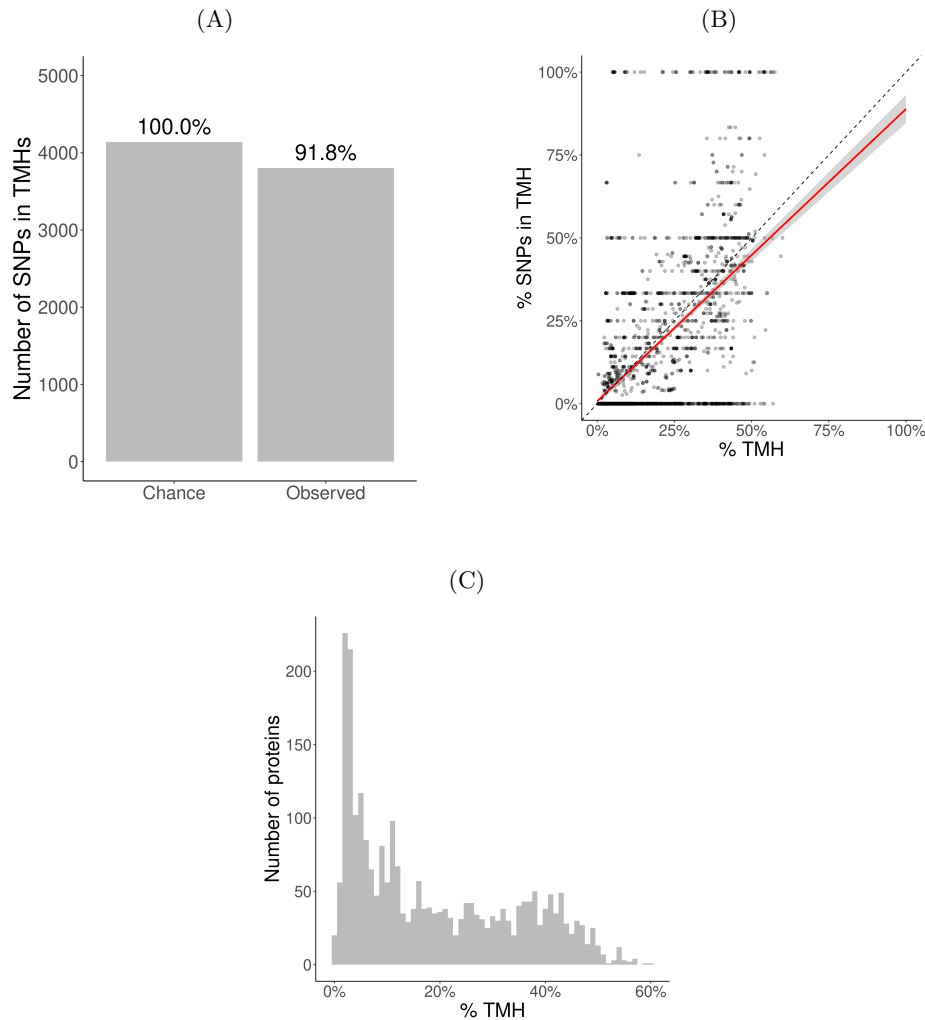


Figure 3: **Evolutionary conservation of human TMHs.** (A) The number of SNPs in TMHs as expected by chance (left bar) and found in the dbSNP database (right bar). Percentages show the relative conservation of SNPs in TMHs found. (B) Percentage of SNPs found in TMHs. Each point shows for one protein the predicted percentage of TMH (*x*-axis) and the observed occurrence of SNPs being located within a TMH (*y*-axis). The dashed diagonal line shows the line of equality (i.e., equal conservation of TMHs and soluble protein regions). The red line indicates a linear fit, the gray area its 95% confidence interval. (C) Distribution of the percentages of TMH in the TMPs used in this study.

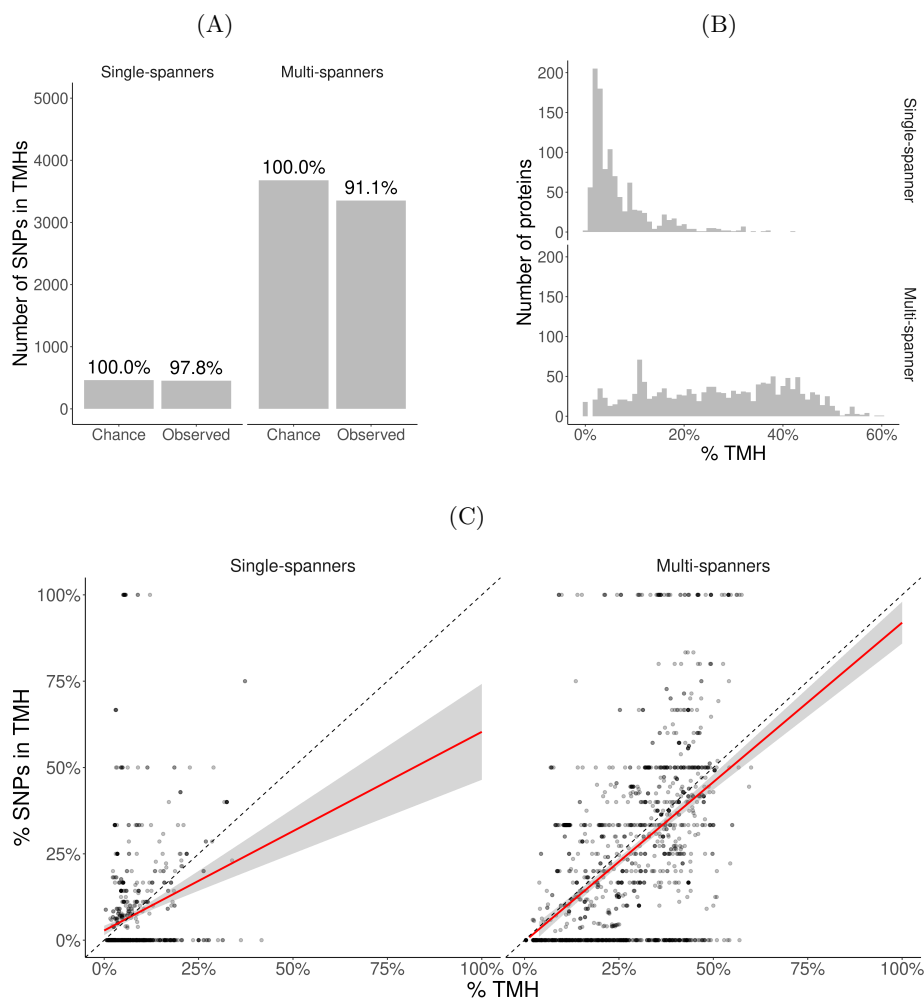


Figure 4: Membrane proteins with multiple TMHs are evolutionary more conserved than proteins with only a single TMH. (A) The number of SNPs in TMHs as expected by chance and observed in the dbSNP database, for TMPs with one TMH (single-spanners) and multiple TMHs (multi-spanners). Percentages show the relative conservation of SNPs in TMHs found. (B) Distribution of the proportion of amino acids residing in the plasma membrane. (C) Percentage of SNPs found in TMPs predicted to have only a single (left) or multiple (right) TMHs. Each point shows for one protein the predicted percentage of TMH (x -axis) and the observed occurrence of SNPs being located within a TMH (y -axis). The dashed diagonal lines show the line of equality (i.e., equal conservation of TMHs and soluble protein regions). The red line indicates a linear fit, the gray area its 95% confidence interval.

245 these predicted TMHs. Each percentage pair was plotted in figure 3B. The
246 proportion of SNPs found in TMHs varied from none (i.e. all SNPs were in
247 soluble regions) to all. To determine if SNPs were randomly distributed over the
248 protein, we performed a linear regression analysis, and added a 95% confidence
249 interval on this regression. This linear fit nearly goes through the origin and
250 has a slope below the line of equality, which shows that less SNPs are found in
251 TMHs than expected by chance.

252 We determined the probability to find the observed amount of SNPs in TMHs
253 by chance, i.e., when assuming SNPs occur just as likely in soluble domains as
254 in TMHs. We used a binomial Poisson distribution, where the number of trials
255 (n) equals the number of SNPs, which is 21,208. The probability of success
256 for the i th TMP (p_i), is the percentage of residues within a TMH per TMP.
257 These percentages are shown as a histogram in figure 3C. The expected number
258 of SNPs expected to be found in TMHs by chance equals $\sum p \approx 4,141$. As
259 we observed 3,803 SNPs in TMHs, we calculated the probability of having that
260 amount or less successes. We used the type I error cut-off value of $\alpha = 2.5\%$. The
261 chance to find, within TMHs, this amount or less SNPs equals $6.8208 \cdot 10^{-11}$. We
262 determined the relevance of this finding, by calculating how much less SNPs are
263 found in TMHs, when compared to soluble regions, which is the ratio between
264 the number of SNPs found in TMHs versus the number of SNPs as expected
265 by chance. In effect, per 1000 SNPs found in soluble protein domains, one finds
266 918 SNPs in TMHs, as depicted in figure 3A.

267 We split this analysis for TMPs containing only a single TMH (so-called
268 single-membrane spanners) and TMPs containing multiple TMHs (multi-membrane
269 spanners). We hypothesized that single-membrane spanners are less conserved
270 than multi-membrane spanners, because multi-membrane spanners might have
271 protein-protein interactions between their TMHs, for example to accommodate

272 active sites, and thus might have additional structural constraints. From the
273 split data, we did the same analysis as for the total TMPs. Figure 4C shows the
274 percentages of TMHs for individual proteins as a function of the percentage of
275 SNPs located in TMHs. For both single- and multi-spanners, a linear regression
276 shows that less SNPs are found in TMHs, than expected by chance.

277 We also determined the probability to find the observed amount of SNPs by
278 chance in single- and multi-spanners. For single-spanners, we found 452 SNPs
279 in TMH, where ≈ 462 were expected by chance. The chance to observe this or a
280 lower number by chance is 0.319. As this chance was higher than our $\alpha = 0.025$,
281 we consider this no significant effect. For the multi-spanners, we found 3,351
282 SNPs in TMH, where $\approx 3,678$ were expected by chance. The chance to observe
283 this or a lower number by chance is $8.315841 \cdot 10^{-12}$, which means this number
284 is significantly less as explained by variation.

285 Also, for single- and multi-spanners, we determined the relevance of this
286 finding by calculating how much less SNPs are found in TMHs when compared
287 to soluble regions, as depicted in figure 4A. In effect, per 1,000 SNPs found in
288 soluble protein domains, one finds 978 SNPs in TMHs of single-spanners and
289 911 SNPs in TMHs of multi-spanners.

290 4 Discussion

291 Epitope prediction is important to understand the immune system and for the
292 design of vaccines. In this study, we provide evidence that epitopes derived
293 from TMHs are a major but overlooked source of MHC epitopes. Our bioinfor-
294 matics predictions indicate that TMH-derived epitopes are presented to both
295 cytolytic and helper T cells more often than expected by chance, regardless of
296 the organism. Moreover, reanalysis of peptide elution studies confirmed the
297 presentation of TMH-derived epitopes. Finally, our SNP analysis shows that

298 TMHs are evolutionary more conserved than solvent-exposed protein regions.

299 **4.1 Mechanism of MHC presentation of TMH-derived epi-** 300 **topes**

301 Although our data show that TMH-derived epitopes are presented in MHC-I
302 and MHC-II, the molecular mechanisms of how integral membrane proteins are
303 processed for MHC presentation are largely unknown [7]. Most prominently, the
304 fundamental principles of how TMHs are extracted from their hydrophobic lipid
305 environments into the aqueous vacuolar lumen, and their prior or subsequent
306 proteolytic processing are unresolved.

307 A first possibility is that the extraction of TMPs from the membrane is
308 mediated by the ER-associated degradation (ERAD) machinery. For MHC class
309 I (MHC-I) antigen presentation of soluble proteins, the loading of the epitope
310 primarily occurs at the endoplasmic reticulum (ER). The chaperones tapasin
311 (TAPBP), ERp57 (PDIA3), and calreticulin (CALR) [24] first assemble and
312 stabilize the heavy and light chains of MHC-I. Later, this complex binds to the
313 transporter associated with antigen processing (TAP) leading to the formation of
314 the so-called peptide-loading complex (PLC). The PLC drives import of peptides
315 into the ER and mediates their subsequent loading into the peptide-binding
316 groove of MHC-I [25]. Membrane proteins first will have to be extracted from
317 the membrane before they become amenable to this MHC-I loading by the
318 PLC. In the ER, this process can be orchestrated by the ERAD machinery,
319 consisting of several chaperones that recognize TMPs, ubiquitinate them, and
320 extract them from the ER membrane into the cytosol (retrotranslocation) for
321 proteasomal degradation [26, 27]. Similar to the peptides generated from soluble
322 proteins, the TMP-derived peptides might then be re-imported by TAP into the
323 ER for MHC-I loading. This ERAD-driven antigen retrotranslocation might be

324 facilitated by lipid bodies (LBs) [28], since LBs can serve as cytosolic sites for
325 ubiquitination of ER-derived cargo [29].

326 A second possibility is that TMPs are proteolytically processed by intramem-
327 brane proteases that cleave TMHs while they are still membrane embedded.
328 Supporting this hypothesis is the well established notion that peptides gener-
329 ated by signal peptide peptidases (SPPs), an important class of intramembrane
330 proteases that cleave TMH-like signal sequences, are presented on a specialized
331 class of MHC-I called HLA-E [30]. The loading of peptides generated by SPP
332 onto MHC-I does not depend on the proteasome and TAP, possibly because
333 the peptides are directly released into the lumen of the ER [30]. However, this
334 mechanism would not explain how multispinner polytopic membrane proteins
335 can be processed for antigen presentation, because SPPs only cleave TMH-like
336 signal sequences at the N-terminus of a protein. Nevertheless, the presenta-
337 tion of peptides with a high hydrophobicity index was shown to be independent
338 of TAP as well [31], suggesting the TMH peptides might perhaps be released
339 directly in the ER lumen by other intramembrane proteases.

340 A third possibility is that peptide processing and MHC-loading occur in
341 multivesicular bodies (MVBs) [30]. TMPs can be routed from the plasma mem-
342 brane and other organelles by vesicular trafficking to endosomes. Eventually,
343 these TMPs can be sorted by the endosomal sorting complexes required for
344 transport (ESCRT) pathway into luminal invaginations that pinch off from the
345 limiting membrane and form intraluminal vesicles. This thus results in MVBs
346 where the membrane proteins destined for degradation are located in intralumi-
347 nal vesicles. Upon the fusion of MVBs with lysosomes, the entire intraluminal
348 vesicles including the TMPs are degraded [32]. Via this mechanism, TMPs
349 might well be processed for antigen presentation, particularly since the loading
350 of MHC class II molecules is well understood to occur in MVBs [33, 34, 35].

351 However, such processing of membrane proteins in MVBs for antigen presenta-
352 tion poses a problem, because complexes of HLA-DR with its antigen-loading
353 chaperon HLA-DM were only observed on intraluminal vesicles, but not on the
354 limiting membranes of MVBs [35], indicating that epitope loading of MHC-II
355 also occurs at intraluminal vesicles. This observation hence raises the question
356 how the intraluminal vesicles carrying the TMPs destined for antigen presen-
357 tation can be selectively degraded, while the intraluminal vesicles carrying the
358 MHC-II remain intact. A second problem is that phagosomes carrying inter-
359 nalized microbes lack intraluminal vesicles, and it is hence unclear how TMPs
360 from these microbes would be routed to MVBs for MHC-II loading [35].

361 Alternatively to the enzymatic degradation of lipids in MVBs by lipases
362 [36, 37], they might be oxidatively degraded by reactions with radical oxygen
363 species (ROS) produced by the NADPH oxidase NOX2 [38]. This oxidation can
364 result in a destabilization and disruption of membranes [38] and might thereby
365 lead to the extraction of TMPs. Due to the hydrophobic nature of TMHs,
366 however, the extracted proteins will likely aggregate and it is unclear how these
367 aggregates would be processed further for MHC loading.

368 **4.2 T cells recognize different protein regions than B cells**

369 An important implication from the over-presentation of TMH-derived epitopes
370 is that T cells will largely recognize different protein regions than B cells. Pre-
371 sentation of antigens by MHC-II is important for the activation of naive B cells
372 by helper T cells. For this activation, B cells first ingest antigen that is bound
373 to their B cell receptor, and subsequently present peptides derived from this
374 antigen in MHC-II to helper T cells. Following their activation by the T cells, B
375 cells mature into plasma cells and release antibodies which recognize the same
376 part of the antigen as the original B cell receptor. B cell receptors and antibod-

377 ies will thus recognize solvent-exposed regions of antigens that are accessible
378 for binding to the B cell receptor. However, the results from our study predict
379 that most MHC-II haplotypes present relatively hydrophobic peptides, which
380 are less likely to be solvent-exposed. It is unknown why B and T cells seem
381 to predominantly recognize different protein regions, but one possibility might
382 be that this lowers the chance of B cell mediated autoimmune diseases, because
383 auto-reactive B and T cells recognizing different parts of the same antigen would
384 need to be present for breakage of B cell tolerance.

385 **4.3 Evolutionary conservation of TMHs**

386 In general, one might expect that evolutionary selection results in an immune
387 system that is most attentive for protein regions that are essential for the sur-
388 vival, proliferation and/or virulence of pathogenic microbes, as these will be
389 most conserved. In SARS-CoV-2, for example, there is preliminary evidence
390 that the strongest selection pressure is upon residues that change its viru-
391 lence [39]. These regions, however, may only account for a small part of a
392 pathogen's proteome. Additionally, the structure and function of these essen-
393 tial regions might differ widely between different pathogenic proteins. Because
394 of this scarcity and variance in targets, one can imagine that it will be mostly
395 unfeasible to provide innate immune responses against such rare essential pro-
396 tein regions, as suggested in a study on influenza [40], where it was found that
397 the selection pressure exerted by the immune system was either weak or absent.

398 Evolutionary selection of pathogens by a host's immune system, however, is
399 likelier to occur for proteomic patterns that are general, over patterns that are
400 rare. While essential catalytic sites in a pathogenic proteome might be relatively
401 rare, TMHs are common and thus might be a more feasible target for evolution
402 to respond to. Indeed, we have found the signature of evolution when both

403 factors, that is, TMHs and catalytic sites are likely to co-occur, which is in TMPs
404 that span the membrane at least twice. In contrast to single-spanners, where
405 we found no significant evolutionary conservation, the TMHs of multi-spanners
406 are more evolutionary conserved than soluble protein regions. Likely, the TMHs
407 in many multi-spanners need to interact with each other for correct protein
408 structure and function and they might hence be more structurally constrained
409 compared to the TMHs of single-spanners. Thus, we speculate that the human
410 immune system is more attentive towards TMHs in multi-spanners, as these are
411 evolutionarily more conserved.

412 There have been more efforts to assess the conservation of TMHs, using
413 different methodologies. One such example is [41], in which aligned protein
414 sequence data was used. Also this study found that TMHs are evolutionarily
415 more conserved, as the mean amino acid substitution rate in TMHs is about ten
416 percent lower, which is a similar value as we found. Another example is a study
417 that estimated the conservation scores for TMHs and soluble regions based on
418 alignments of evolutionary related proteins, and also found that TMHs are more
419 conserved, with a conservation score that was 17% higher in TMHs [42]. Note
420 that the last study also found that mutations in human TMHs are likelier to
421 cause a disease, in line with our conclusion that TMHs are more conserved.

422 Together, from this study, two important conclusions can be drawn. First,
423 the MHC over-presentation of TMHs is likely a general feature and predicted
424 to occur for most haplotypes of both MHC-I and -II and for humans as well as
425 bacterial and viral pathogens. Second, TMHs are genuinely more evolutionary
426 conserved than soluble protein motifs, at least in the human proteome.

427 **5 Acknowledgments**

428 We thank the Center for Information Technology of the University of Gronin-
429 gen for its support and for providing access to the Peregrine high performance
430 computing cluster. FB is funded by a Veni grant from the Netherlands Orga-
431 nization for Scientific Research (016.Veni.192.026) and an Off-Road Grant from
432 the Dutch Medical Science Foundation (ZonMW 04510011910005). GvdB is
433 funded by a Young Investigator Grant from the Human Frontier Science Pro-
434 gram (HFSP; RGY0080/2018), and a Vidi grant from the Netherlands Orga-
435 nization for Scientific Research (NWO-ALW VIDI 864.14.001). GvdB has re-
436 ceived funding from the European Research Council (ERC) under the European
437 Union’s Horizon 2020 research and innovation programme (grant agreement No.
438 862137).

439 **6 Data Accessibility**

440 All code is archived at <http://github.com/richelbilderbeek/someplace>,
441 with DOI <https://doi.org/12.3456/zenodo.1234567>.

442 **7 Authors’ contributions**

443 RJCB and FB conceived the idea for this research. MVB helped with the
444 proteome analysis of *M. tuberculosis*. RJCB wrote the code. RJCB, MB, GvdB
445 and FB wrote the article.

446 **References**

447 [1] Ole Lund, Morten Nielsen, Can Kesmir, Anders Gorm Petersen, Claus
448 Lundegaard, Peder Worning, Christina Sylvester-Hvid, Kasper Lamberth,

- 449 Gustav Røder, Sune Justesen, et al. Definition of supertypes for HLA
450 molecules using clustering of specificity matrices. *Immunogenetics*, 55(12):
451 797–810, 2004.
- 452 [2] Steven GE Marsh, ED Albert, WF Bodmer, RE Bontrop, B Dupont,
453 HA Erlich, M Fernández-Viña, DE Geraghty, R Holdsworth, CK Hurley,
454 et al. Nomenclature for factors of the HLA system, 2010. *Tissue antigens*,
455 75(4):291, 2010.
- 456 [3] Simone Sommer. The importance of immune gene variability (MHC) in evo-
457 lutionary ecology and conservation. *Frontiers in zoology*, 2(1):1–18, 2005.
- 458 [4] Mette Voldby Larsen, Alina Lelic, Robin Parsons, Morten Nielsen, Ilka
459 Hoof, Kasper Lamberth, Mark B Loeb, Søren Buus, Jonathan Bramson,
460 and Ole Lund. Identification of CD8+ T cell epitopes in the West Nile
461 virus polyprotein by reverse-immunology using NetCTL. *PloS one*, 5(9),
462 2010.
- 463 [5] Ingrid MM Schellens, Can Kesmir, Frank Miedema, Debbie van Baarle,
464 and José AM Borghans. An unanticipated lack of consensus cytotoxic T
465 lymphocyte epitopes in HIV-1 databases: the contribution of prediction
466 programs. *Aids*, 22(1):33–37, 2008.
- 467 [6] Sheila T Tang, Krista E van Meijgaarden, Nadia Caccamo, Giuliana Gug-
468 gino, Michèl R Klein, Pascale van Weeren, Fatima Kazi, Anette Stryhn,
469 Alexander Zaigler, Ugur Sahin, et al. Genome-based in silico identifica-
470 tion of new Mycobacterium tuberculosis antigens activating polyfunctional
471 CD8+ T cells in human tuberculosis. *The Journal of Immunology*, 186(2):
472 1068–1080, 2011.
- 473 [7] Frans Bianchi, Johannes Textor, and Geert van den Bogaart. Transmem-

- 474 brane helices are an overlooked source of Major Histocompatibility Com-
475 plex Class I epitopes. *Frontiers in immunology*, 8:1118, 2017.
- 476 [8] Anders Krogh, Björn Larsson, Gunnar Von Heijne, and Erik LL Sonnham-
477 mer. Predicting transmembrane protein topology with a hidden Markov
478 model: application to complete genomes. *Journal of molecular biology*, 305
479 (3):567–580, 2001.
- 480 [9] Lukas Käll, Anders Krogh, and Erik LL Sonnhammer. A combined trans-
481 membrane topology and signal peptide prediction method. *Journal of*
482 *molecular biology*, 338(5):1027–1036, 2004.
- 483 [10] Masafumi Arai, Hironori Mitsuke, Masami Ikeda, Jun-Xiong Xia, Takashi
484 Kikuchi, Masanobu Satake, and Toshio Shimizu. ConPred II: a consensus
485 prediction method for obtaining transmembrane topology models with high
486 reliability. *Nucleic acids research*, 32(suppl_2):W390–W393, 2004.
- 487 [11] David T Jones. Improving the accuracy of transmembrane protein topology
488 prediction using evolutionary information. *Bioinformatics*, 23(5):538–544,
489 2007.
- 490 [12] Martin Klammer, David N Messina, Thomas Schmitt, and Erik LL
491 Sonnhammer. MetaTM-a consensus method for transmembrane protein
492 topology prediction. *BMC bioinformatics*, 10(1):314, 2009.
- 493 [13] Qing Wang, Chongming Ni, Zhen Li, Xiufeng Li, Renmin Han, Feng Zhao,
494 Jinbo Xu, Xin Gao, and Sheng Wang. PureseqTM: efficient and accu-
495 rate prediction of transmembrane topology from amino acid sequence only.
496 *bioRxiv*, page 627307, 2019.
- 497 [14] Mamoun Ahram, Zoi I Litou, Ruihua Fang, and Ghaith Al-Tawallbeh.

- 498 Estimation of membrane proteins in the human proteome. *In silico biology*,
499 6(5):379–386, 2006.
- 500 [15] Tara Hessa, Nadja M Meindl-Beinker, Andreas Bernsel, Hyun Kim, Yoko
501 Sato, Mirjam Lerch-Bader, IngMarie Nilsson, Stephen H White, and Gun-
502 nar Von Heijne. Molecular code for transmembrane-helix recognition by
503 the sec61 translocon. *Nature*, 450(7172):1026–1030, 2007.
- 504 [16] DT Jones, WR Taylor, and JM Thornton. A model recognition approach
505 to the prediction of all-helical membrane protein structure and topology.
506 *Biochemistry*, 33(10):3038–3049, 1994.
- 507 [17] Elin Bergsgeng, Siri Dørum, Magnus Ø Arntzen, Morten Nielsen, Ståle
508 Nygård, Søren Buus, Gustavo A de Souza, and Ludvig M Sollid. Dif-
509 ferent binding motifs of the celiac disease-associated hla molecules DQ2.5,
510 DQ2.2, and DQ7.5 revealed by relative quantitative proteomics of endoge-
511 nous peptide repertoires. *Immunogenetics*, 67(2):73–84, 2015.
- 512 [18] Xiaoshan M Shao, Rohit Bhattacharya, Justin Huang, IK Ashok Sivaku-
513 mar, Collin Tokheim, Lily Zheng, Dylan Hirsch, Benjamin Kaminow, Ash-
514 ton Omdahl, Maria Bonsack, et al. High-throughput prediction of MHC
515 class I and II neoantigens with MHCnuggets. *Cancer Immunology Research*,
516 8(3):396–408, 2020.
- 517 [19] Jason Greenbaum, John Sidney, Jolan Chung, Christian Brander, Bjoern
518 Peters, and Alessandro Sette. Functional classification of class II human
519 leukocyte antigen (HLA) molecules reveals seven different supertypes and a
520 surprising degree of repertoire sharing across supertypes. *Immunogenetics*,
521 63(6):325–335, 2011.
- 522 [20] Ingrid MM Schellens, Ilka Hoof, Hugo D Meiring, Sanne NM Spijkers,
523 Martien CM Poelen, Kees van der Poel, Ana I Costa, Cecile ACM van

- 524 Els, Debbie van Baarle, Can Kesmir, et al. Comprehensive analysis of the
525 naturally processed peptide repertoire: differences between HLA-A and B
526 in the immunopeptidome. *PloS one*, 10(9):e0136417, 2015.
- 527 [21] Stephen T Sherry, M-H Ward, M Kholodov, J Baker, Lon Phan, Eliza-
528 beth M Smigielski, and Karl Sirotkin. dbSNP: the ncbi database of genetic
529 variation. *Nucleic acids research*, 29(1):308–311, 2001.
- 530 [22] Garth R Brown, Vichet Hem, Kenneth S Katz, Michael Ovetsky, Craig
531 Wallin, Olga Ermolaeva, Igor Tolstoy, Tatiana Tatusova, Kim D Pruitt,
532 Donna R Maglott, et al. Gene: a gene-centered information resource at
533 NCBI. *Nucleic acids research*, 43(D1):D36–D42, 2015.
- 534 [23] Eric W Sayers, Tanya Barrett, Dennis A Benson, Evan Bolton, Stephen H
535 Bryant, Kathi Canese, Vyacheslav Chetvernin, Deanna M Church, Michael
536 DiCuccio, Scott Federhen, et al. Database resources of the national center
537 for biotechnology information. *Nucleic acids research*, 39(suppl_1):D38–
538 D51, 2010.
- 539 [24] Kenneth L Rock, Eric Reits, and Jacques Neefjes. Present yourself! by mhc
540 class i and mhc class ii molecules. *Trends in immunology*, 37(11):724–737,
541 2016.
- 542 [25] Andreas Blees, Dovile Janulienė, Tommy Hofmann, Nicole Koller, Carla
543 Schmidt, Simon Trowitzsch, Arne Moeller, and Robert Tampé. Structure
544 of the human mhc-i peptide-loading complex. *Nature*, 551(7681):525–528,
545 2017.
- 546 [26] G Michael Preston and Jeffrey L Brodsky. The evolving role of ubiquitin
547 modification in endoplasmic reticulum-associated degradation. *Biochemical*
548 *Journal*, 474(4):445–469, 2017.

- 549 [27] Birgit Meusser, Christian Hirsch, Ernst Jarosch, and Thomas Sommer.
550 Erad: the long road to destruction. *Nature cell biology*, 7(8):766–772, 2005.
- 551 [28] Laurence Bougnères, Julie Helft, Sangeeta Tiwari, Pablo Vargas, Benny
552 Hung-Jun Chang, Lawrence Chan, Laura Campisi, Gregoire Lauvau,
553 Stephanie Hugues, Pradeep Kumar, et al. A role for lipid bodies in the
554 cross-presentation of phagocytosed antigens by mhc class i in dendritic
555 cells. *Immunity*, 31(2):232–244, 2009.
- 556 [29] Toyoshi Fujimoto and Yuki Ohsaki. The proteasomal and autophagic path-
557 ways converge on lipid droplets. *Autophagy*, 2(4):299–301, 2006.
- 558 [30] Cláudia C Oliveira and Thorbald van Hall. Alternative antigen processing
559 for mhc class i: multiple roads lead to rome. *Frontiers in immunology*, 6:
560 298, 2015.
- 561 [31] Georg Lautscham, Sabine Mayrhofer, Graham Taylor, Tracey Haigh, Ali-
562 son Leese, Alan Rickinson, and Neil Blake. Processing of a multiple mem-
563 brane spanning epstein-barr virus protein for cd8+ t cell recognition reveals
564 a proteasome-dependent, transporter associated with antigen processing–
565 independent pathway. *The Journal of experimental medicine*, 194(8):1053–
566 1068, 2001.
- 567 [32] Jean Gruenberg. Life in the lumen: the multivesicular endosome. *Traffic*,
568 21(1):76–93, 2020.
- 569 [33] Monique Kleijmeer, Georg Ramm, Danita Schuurhuis, Janice Griffith,
570 Maria Rescigno, Paola Ricciardi-Castagnoli, Alexander Y Rudensky, Ferry
571 Ossendorp, Cornelis JM Melief, Willem Stoorvogel, et al. Reorganization
572 of multivesicular bodies regulates mhc class ii antigen presentation by den-
573 dritic cells. *The Journal of cell biology*, 155(1):53–64, 2001.

- 574 [34] Peter J Peters, Jacques J Neefjes, Viola Oorschot, Hidde L Ploegh, and
575 Hans J Geuze. Segregation of mhc class ii molecules from mhc class i
576 molecules in the golgi complex for transport to lysosomal compartments.
577 *Nature*, 349(6311):669–676, 1991.
- 578 [35] Wilbert Zwart, Alexander Griekspoor, Coenraad Kuijl, Marije Marsman,
579 Jacco van Rheenen, Hans Janssen, Jero Calafat, Marieke van Ham, Lennert
580 Janssen, Marcel van Lith, et al. Spatial separation of hla-dm/hla-dr inter-
581 actions within miic and phagosome-induced immune escape. *Immunity*, 22
582 (2):221–233, 2005.
- 583 [36] Peter Sander, Katja Becker, and Michael Dal Molin. Lipase processing of
584 complex lipid antigens. *Cell chemical biology*, 23(9):1044–1046, 2016.
- 585 [37] Martine Gilleron, Marco Lepore, Emilie Layre, Diane Cala-De Paepe,
586 Naila Mebarek, James A Shayman, Stéphane Canaan, Lucia Mori, Frédéric
587 Carrière, Germain Puzo, et al. Lysosomal lipases plrp2 and lpla2 process
588 mycobacterial multi-acylated lipids and generate t cell stimulatory anti-
589 gens. *Cell chemical biology*, 23(9):1147–1156, 2016.
- 590 [38] Ilse Dingjan, Daniëlle RJ Verboogen, Laurent M Paardekooper, Natalia H
591 Revelo, Simone P Sittig, Linda J Visser, Gabriele Fischer Von Mollard,
592 Stefanie SV Henriët, Carl G Figdor, Martin Ter Beest, et al. Lipid perox-
593 idation causes endosomal antigen release for cross-presentation. *Scientific*
594 *reports*, 6(1):1–12, 2016.
- 595 [39] Lauro Velazquez-Salinas, Selene Zarate, Samantha Eberl, Douglas P
596 Gladue, Isabel Novella, and Manuel V Borca. Positive selection of ORF3a
597 and ORF8 genes drives the evolution of SARS-CoV-2 during the 2020
598 COVID-19 pandemic. *bioRxiv*, 2020.

- 599 [40] Alvin X Han, Sebastian Maurer-Stroh, and Colin A Russell. Individual
600 immune selection pressure has limited impact on seasonal influenza virus
601 evolution. *Nature ecology & evolution*, 3(2):302–311, 2019.
- 602 [41] Timothy J Stevens and Isaiah T Arkin. Substitution rates in α -helical
603 transmembrane proteins. *Protein Science*, 10(12):2507–2517, 2001.
- 604 [42] Amit Oberai, Nathan H Joh, Frank K Pettit, and James U Bowie. Struc-
605 tural imperatives impose diverse evolutionary constraints on helical mem-
606 brane proteins. *Proceedings of the National Academy of Sciences*, 106(42):
607 17747–17750, 2009.
- 608 [43] Morten Nielsen, Claus Lundegaard, Thomas Blicher, Bjoern Peters,
609 Alessandro Sette, Sune Justesen, Søren Buus, and Ole Lund. Quantita-
610 tive predictions of peptide binding to any HLA-DR molecule of known
611 sequence: NetMHCIIpan. *PLoS computational biology*, 4(7), 2008.
- 612 [44] Edita Karosiene, Michael Rasmussen, Thomas Blicher, Ole Lund, Søren
613 Buus, and Morten Nielsen. NetMHCIIpan-3.0, a common pan-specific
614 MHC class II prediction method including all three human MHC class
615 II isotypes, HLA-DR, HLA-DP and HLA-DQ. *Immunogenetics*, 65(10):
616 711–724, 2013.
- 617 [45] Richèl J C Bilderbeek. tmhmm, 2019. [https://github.com/
618 richelbilderbeek/tmhmm](https://github.com/richelbilderbeek/tmhmm) [Accessed: 2019-03-08].
- 619 [46] Richèl J C Bilderbeek. pureseqtmr, 2020. [https://github.com/
620 richelbilderbeek/pureseqtmr](https://github.com/richelbilderbeek/pureseqtmr) [Accessed: 2020-05-19].
- 621 [47] Richèl J C Bilderbeek. netmhc2pan, 2019. [https://github.com/
622 richelbilderbeek/netmhc2pan](https://github.com/richelbilderbeek/netmhc2pan) [Accessed: 2019-03-08].

- 623 [48] Richèl J C Bilderbeek. bbbq, 2020. [https://github.com/](https://github.com/richelbilderbeek/bbbq)
624 [richelbilderbeek/bbbq](https://github.com/richelbilderbeek/bbbq) [Accessed: 2020-09-02].
- 625 [49] Steffen Möller, Michael DR Croning, and Rolf Apweiler. Evaluation of
626 methods for the prediction of membrane spanning regions. *Bioinformatics*,
627 17(7):646–653, 2001.
- 628 [50] Claus Lundegaard, Ole Lund, and Morten Nielsen. Prediction of epitopes
629 using neural network based methods. *Journal of immunological methods*,
630 374(1-2):26–34, 2011.
- 631 [51] Morten Nielsen, Claus Lundegaard, Peder Worning, Sanne Lise Lauemøller,
632 Kasper Lamberth, Søren Buus, Søren Brunak, and Ole Lund. Reliable
633 prediction of T-cell epitopes using neural networks with novel sequence
634 representations. *Protein Science*, 12(5):1007–1017, 2003.
- 635 [52] Edita Karosiene, Claus Lundegaard, Ole Lund, and Morten Nielsen.
636 NetMHCcons: a consensus method for the major histocompatibility com-
637 plex class I predictions. *Immunogenetics*, 64(3):177–186, 2012.
- 638 [53] Morten Nielsen, Claus Lundegaard, Peder Worning, Christina Sylvester
639 Hvid, Kasper Lamberth, Søren Buus, Søren Brunak, and Ole Lund. Im-
640 proved prediction of MHC class I and class II epitopes using a novel Gibbs
641 sampling approach. *Bioinformatics*, 20(9):1388–1397, 2004.
- 642 [54] David J. Winter. rentrez: an R package for the NCBI eUtils API. *The R*
643 *Journal*, 9:520–526, 2017.
- 644 [55] Richèl J C Bilderbeek. sprentrez, 2021. [https://github.com/](https://github.com/richelbilderbeek/sprentrez)
645 [richelbilderbeek/sprentrez](https://github.com/richelbilderbeek/sprentrez) [Accessed: 2021-02-09].
- 646 [56] Lucia Musumeci, Jonathan W Arthur, Florence SG Cheung, Ashraf
647 Hoque, Scott Lippman, and Juergen KV Reichardt. Single nucleotide dif-

648 ferences (SNDs) in the dbSNP database may lead to errors in genotyping
649 and haplotyping studies. *Human mutation*, 31(1):67–73, 2010.

650 [57] Ryan Hunt, Zuben E Sauna, Suresh V Ambudkar, Michael M Gottesman,
651 and Chava Kimchi-Sarfaty. Silent (synonymous) SNPs: should we care
652 about them? *Single nucleotide polymorphisms*, pages 23–39, 2009.

653 A Supplementary materials

654 A.1 Differences with Bianchi et al., 2017

655 A part of this study does the same analysis as Bianchi et al., 2017. mainly
656 concern the use of different software and a different definition of what an MHC
657 binder is.

658 The earlier study defined a peptide an MHC binder if *within the protein* in
659 which it was found, is was among the peptides with the 2% lowest IC50 val-
660 ues. This can be seen at [https://github.com/richelbilderbeek/bianchi_](https://github.com/richelbilderbeek/bianchi_et_al_2017/blob/master/predict-binders.R)
661 [et_al_2017/blob/master/predict-binders.R](https://github.com/richelbilderbeek/bianchi_et_al_2017/blob/master/predict-binders.R), where the binders are written
662 to file.

663 However, in this study, an MHC binder is defined as a peptide within a
664 *proteome* in which it is found, that is among the peptides with the 2% lowest
665 IC50 values. Subsection A.8 shows the IC50 values for a binder per haplotype.
666 We believe that our revised definition is more correct, as it overcomes bias from
667 proteins with very low numbers of peptides and/or MHC-predicted binders.

668 Our previous study used the TMHMM web server to predict TMHs. The
669 desktop version of TMHMM, however, gives an error message on the 25 seleno-
670 proteins found in the human reference proteome. For the sake of reproducible
671 research, we used the desktop version (as we can call it from scripts) and, due
672 to this, we removed the selenoproteins from this analysis.

673 To verify if the previous and the current method give rise to notable dif-
674 ference, we show a side-by-side comparison in figures 5A and 5B. The figures
675 that haplotypes that over-present or under-present TMH-derived epitopes, do
676 so in both studies. The extent to which TMH-derived epitopes are presented,
677 however, is more extreme in our current setup.

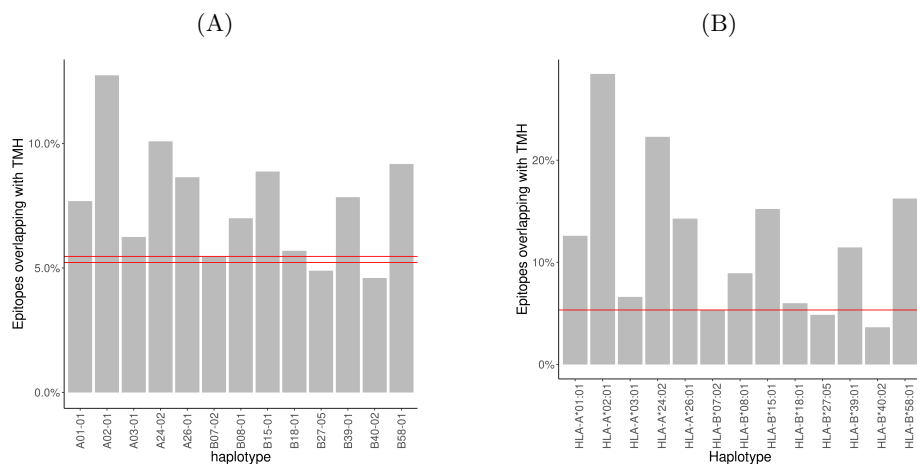


Figure 5: **(A)** Results for [7]. Red lines denotes the coincidence interval. **(B)** Results for this study. Red line denotes the percentage as expected by chance.

Goal	Tool	Reference
Predict topology	TMHMM	[8]
Predict topology	PureseqTM	[13]
Predict epitopes MHC-I	epitope-prediction	[7]
Predict epitopes MHC-II	NetMHCIIpan	[43, 44]
Call TMHMM from R	tmhmm	[45]
Call PureseqTM from R	pureseqtmr	[46]
Call NetMHCIIpan from R	netmhc2pan	[47]
Combine all	bbbq	[48]

Table 1: Overview of all software used in this research.

678 A.2 Prediction software used

679 For this research, we needed software to predict protein topology, as well as the
680 MHC-I and MHC-II binding affinities of epitopes. We selected our software, by
681 searching the scientific literature to identify the most recent free and open source
682 (FOSS) prediction software. This was done by searching for papers that (1) cite
683 older prediction software, and (2) present a novel method to make predictions.
684 As a starting point, per type of prediction software, a review paper was used
685 ([49] for protein topology, [50] for MHC-I binding affinities and [51] for MHC-II

686 binding affinities).

687 There are multiple computational tools developed to predict which parts of
688 a protein forms a TMH. In 2001, multiple of such prediction tools have been
689 compared [49], of which TMHMM [8] turned out to be the most accurate, as
690 is used in the previous study [7]. However, TMHMM has a restrictive software
691 license and is nearly two decades old. Therefore, PureseqTM [13], was also used
692 in this study, which has been more recently developed and has a free software
693 license.

694 For MHC-I, there are multiple computational tools developed to predict epi-
695 topes. According to [50], at that time, NetMHCcons [52] gave the best predic-
696 tions. We used the same tool as used in our earlier study, `epitope-prediction`
697 [7],

698 Also for MHC-II, there are multiple computational tools developed to pre-
699 dict epitopes, such as using a trained neural network [51] or a Gibbs sam-
700 pling approach [53]. According to [50], in 2011, from a set of multiple tools,
701 NetMHCIIpan [43, 44] made the most accurate predictions. The most recent
702 FOSS tool available now appears to be MHCnuggets [18], which can do both
703 MHC-I and MHC-II predictions. As we already use `epitope-prediction` [7]
704 for MHC-I predictions, we use MHCnuggets only for MHC-II predictions.

705 To retrieve the data from the NCBI databases the `rentrez` R package [54]
706 was used that calls the NCBI website's API. To provide for a stable user expe-
707 rience for all users, this API limits the user to 3 calls per second. Additionally,
708 the API splits the result of a bigger query into multiple pages, each of which
709 needs one API call. We wrote the `sprentrez` package [55] to provide for bigger
710 queries of multiple (and delayed) API calls.

711 **A.3 Prediction software written**

712 The R programming language is used for the complete experiment, including the
713 analysis. The complete experiment is bundled in the 'bbq' R package, which
714 is dependent on 'tmhmm', 'pureseqtmr', 'epitope-prediction' and 'mhc-nuggetsr'
715 as described below.

716 The R package 'tmhmm' was developed to do the similar topology predic-
717 tions as our earlier study (that used 'TMHMM'), yet in an automated way.
718 'TMHMM' has a restrictive software license [8] and allows a user to download a
719 pre-compiled executable after confirmation that he/she is in academia. The R
720 package respects this restriction and allows the user to install and use TMHMM
721 from within R, as done in this study. 'tmhmm' has been submitted to and is
722 accepted by the Comprehensive R Archive Network (CRAN).

723 To be able to call, from R, the TMH prediction software 'PureseqTM' [13],
724 which is written in C, the package 'pureseqtmr' has been developed. 'pureseq-
725 qtmr' allows to install 'PureseqTM' and use most of its features. 'pureseqtmr'
726 has been submitted to and is accepted by CRAN.

727 MHCnuggets is a free and open-source Python package to predict epitope
728 affinity for many MHC-I and MHC-II variants [18]. The R package 'mhc-
729 nuggetsr' allows one to install and use MHCnuggets from within R. Also 'mhc-
730 nuggetsr' has been submitted to and is accepted by CRAN.

731 To reproduce the full experiment presented in this paper, the functions
732 needed are bundled in the 'bbq' R package. This package is too specific to
733 be submitted to CRAN.

Table 2: Percentage of spots and spots that overlap with a TMH

target	mhc_class	n_spots	n_spots_tmh	f_tmh
covid	1	14207	1124	7.91
covid	2	14137	1245	8.81
human	1	11220940	598391	5.33
human	2	11118448	672273	6.05
myco	1	1299707	98613	7.59
myco	2	1279742	108419	8.47

734 **A.4 Prediction of percentage of epitopes overlapping with** 735 **a TMH**

736 Supplementary Table 2 shows an overview of the findings, where a target speci-
737 fies the source of the proteome, where `covid` denotes SARS-CoV-2 and `myco` de-
738 notes *Mycobacterium tuberculosis*. `mhc_class` denotes the MHC class, `n_spots`
739 the number of possible 9-mers (for MHC-I) or 14-mers (for MHC-II) possible.
740 `n_spots_tmh` the number of epitopes that overlapped with a TMH that were
741 binders. `f_tmh` the percentage of peptides that had at least 1 residue overlapping
742 with a TMH.

743 **A.5 Minor methods**

744 These are details that are removed from the 'Methods' section.

745 PureseqTM does not predict the topology of proteins that have less than
746 three amino acids. The TRDD1 ('T cell receptor delta diversity 1') protein,
747 however, is two amino acids long. The R package `pureseqtmr`, however, predicts
748 that mono- and di-peptides are cytosolic.

749 **A.6 Minor discussion**

750 These are details that are removed from the 'Discussion' section.

751 In this experiment we predicted epitopes that overlap with TMHs from a

752 human, bacterial and viral proteome, would these proteins be expressed in a
753 human host. Bacteria, however have different cell membranes and cell walls,
754 hence different structural requirements for a TMH. Both topology prediction
755 tools were trained to recognize human TMHs, thus we cannot be sure that
756 the transmembrane regions predicted in bacterial proteins are actually part of a
757 TMH. For the purpose of this study, we assume the error in topology predictions
758 to be unbiased way towards topology. In other words: that a bacterial TMH is
759 incorrectly predicted to be absent just as often as it is incorrectly predicted to
760 be present elsewhere.

761 Regarding the evolutionary conservation of TMHs using SNPs, again, it is
762 estimated that approximately ten percent of SNPs is a false positive that result
763 from the methods to determine a SNP. One example is that sequence variations
764 are incorrectly detected due to highly similar duplicated sequences [56]. We
765 assume that these duplications occur as often in TMHs as in regions around
766 these, hence we expect this not to affect our results.

767 In our evolutionary experiment, we removed variations that were synony-
768 mous mutations (i.e. resulted in the same amino acid, from a different genetic
769 code) from our analysis. There is evidence, however, that these synonymous mu-
770 tations do have an effect and may even be evolutionary selected for [57]. As the
771 possible effect of synonymous mutations is ignored by our topology prediction
772 software, we do so as well.

773 **A.7 Elution studies**

774 **A.8 IC50 values of binders per haplotype**

775 Per target proteome (i.e. human, SARS-CoV-2, *M tuberculosis*), we collected all
776 9-mers (for MHC-I) and 14-mers (for MHC-II), after removing the selenoproteins
777 and proteins that are shorter than the epitope length. From these epitopes, per

MHC class	Tool	n
I	PureseqTM	1.38% (109/7897)
I	TMHMM	1.43% (113/7897)
II	PureseqTM	3.92% (498/12712)
II	TMHMM	3.96% (504/12712)

Table 3: Percentage of epitopes derived from a TMH found in the two elution studies, for the two different kind of topology prediction tools. The values between braces show the the number of epitopes that were predicted to overlapping with a TMH per all epitopes that could be uniquely mapped to the representative human reference proteome.

Table 4: IC50 values (in nM) per haplotype below which a peptide is considered a binder. percentage used: 2

haplotype	covid	human	myco
HLA-A*01:01	1470.5912	2545.9537	2812.1714
HLA-A*02:01	118.9596	218.7274	186.7565
HLA-A*03:01	537.0144	804.7455	1544.1073
HLA-A*24:02	984.8147	1590.0623	1971.8258
HLA-A*26:01	1095.2591	1771.6924	1526.1101
HLA-B*07:02	1215.7734	705.6514	435.5361
HLA-B*08:01	886.5661	883.0951	1023.2213
HLA-B*18:01	921.4157	1063.2215	1319.0445
HLA-B*27:05	1186.0963	689.8815	475.6130
HLA-B*39:01	437.3506	484.3843	399.3873
HLA-B*40:02	585.6308	541.2392	600.1688
HLA-B*58:01	435.4693	591.0526	538.9063
HLA-B*15:01	281.9129	440.6541	482.8369

778 MHC haplotype, we predicted the IC50 (in nM) using epitope-prediction
779 (for MHC-I) and MHCnuggets (for MHC-II). Here, we show the IC50 value per
780 haplotype that is used to determine if a peptide binds to the haplotype’s MHC
781 for MHC-I (see supplementary Table 4) and MHC-II (see supplementary Table
782 5).

783 A.9 Presentation of TMH-derived epitopes

784 Supplementary Table 6 shows the shorthand notation for the HLA haplotypes.

785 Supplementary Tables 7 and 8 show the exact number of binders, binders

Table 5: IC50 values (in nM) per haplotype below which a peptide is considered a binder. percentage used: 2

haplotype	covid	human	myco
HLA-DRB1*0101	7.3896	9.72	9.9600
HLA-DRB1*0301	121.8420	198.40	164.4900
HLA-DRB1*0401	59.8780	74.92	84.3112
HLA-DRB1*0405	46.2324	51.88	66.7100
HLA-DRB1*0701	17.7464	22.40	28.1700
HLA-DRB1*0802	99.7592	137.16	67.9900
HLA-DRB1*0901	42.3464	53.52	41.5400
HLA-DRB1*1101	35.9988	39.01	48.9200
HLA-DRB1*1201	194.4408	248.72	289.7300
HLA-DRB1*1302	21.1084	40.59	35.4100
HLA-DRB1*1501	32.6196	40.69	46.6700
HLA-DRB3*0101	175.2984	298.94	218.7300
HLA-DRB3*0202	176.8168	291.95	405.8724
HLA-DRB4*0101	47.6384	51.04	62.7800
HLA-DRB5*0101	32.8872	43.52	60.2312
HLA-DQA1*0501/DQB1*0201	193.1108	209.89	174.2124
HLA-DQA1*0501/DQB1*0301	51.2028	43.47	20.3200
HLA-DQA1*0301/DQB1*0302	361.8180	365.96	296.4712
HLA-DQA1*0401/DQB1*0402	214.1932	242.68	199.8912
HLA-DQA1*0101/DQB1*0501	550.4488	674.95	930.9612
HLA-DQA1*0102/DQB1*0602	157.4480	174.82	114.3512

786 that overlap with TMHs and the percentage of binders that overlap with TMHs,
787 as visualized by figure 1A.

788 **A.10 Relative presentation of TMH-derived epitopes**

789 To compare the over-presentation of TMH-derived epitopes between the differ-
790 ent proteomes, we normalized this percentages in such a way that 1.0 is the
791 percentage of TMH-derived epitopes that would be expected by chance. Figure
792 6 and 7 show these normalized values for the MHC-I and MHC-II haplotypes
793 respectively.

794 To determine the additional over-presentation of TMH-derived epitopes in
795 MHC-II (as compared to MHC-I), we normalized the data to enable a side-
796 by-side comparison. The percentage of TMH-derived epitopes presented was
797 normalized to the expected percentage of TMH-derived epitopes, where 1.0
798 denotes that the percentage of presented TMH-derived epitopes matches the
799 values as expected by chance. The normalized values per haplotype are shown
800 in figure 8. To compare the TMH-derived over-presentation per MHC class, we
801 grouped the normalized values per haplotype, and plot the mean and standard
802 error, as shown in figure 9.

803 **A.11 Evolutionary conservation**

804 See supplementary Tables 9 and 10 for an overview of all amounts. In supple-
805 mentary Table 9 there are multiple instances where the amounts are expected
806 to add up, yet don't, as one SNP can work on multiple isoforms. For example,
807 there are 9,621 unique SNPs found in all proteins, of which 4,219 around found
808 in MAPs and 6,026 in TMPs. Apparently, 624 SNPs work on a set of isoforms
809 that contains both MAPs and TMPs.

810 Figure 10 shows the distribution of the number of SNPs per gene name, at

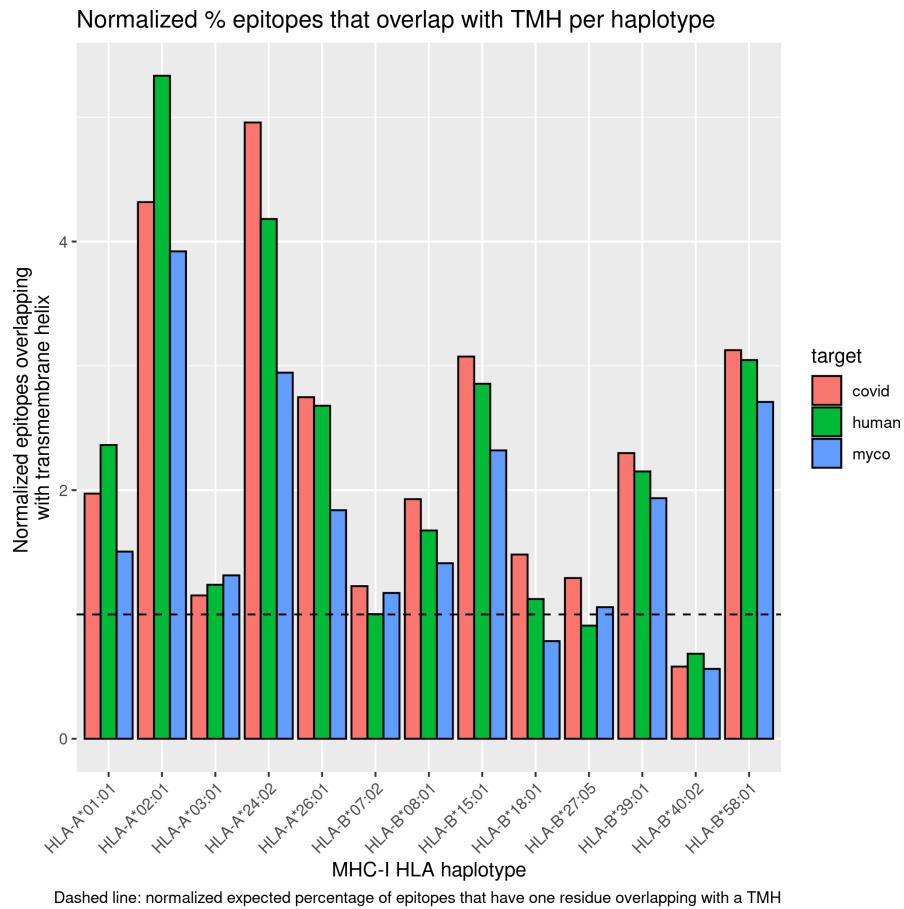


Figure 6: Normalized proportion of MHC-I epitopes overlapping with TMHs for human, viral and bacterial proteomes. Legend: covid = SARS-CoV-2, human = *Homo sapiens*, myco = *Mycobacterium tuberculosis*

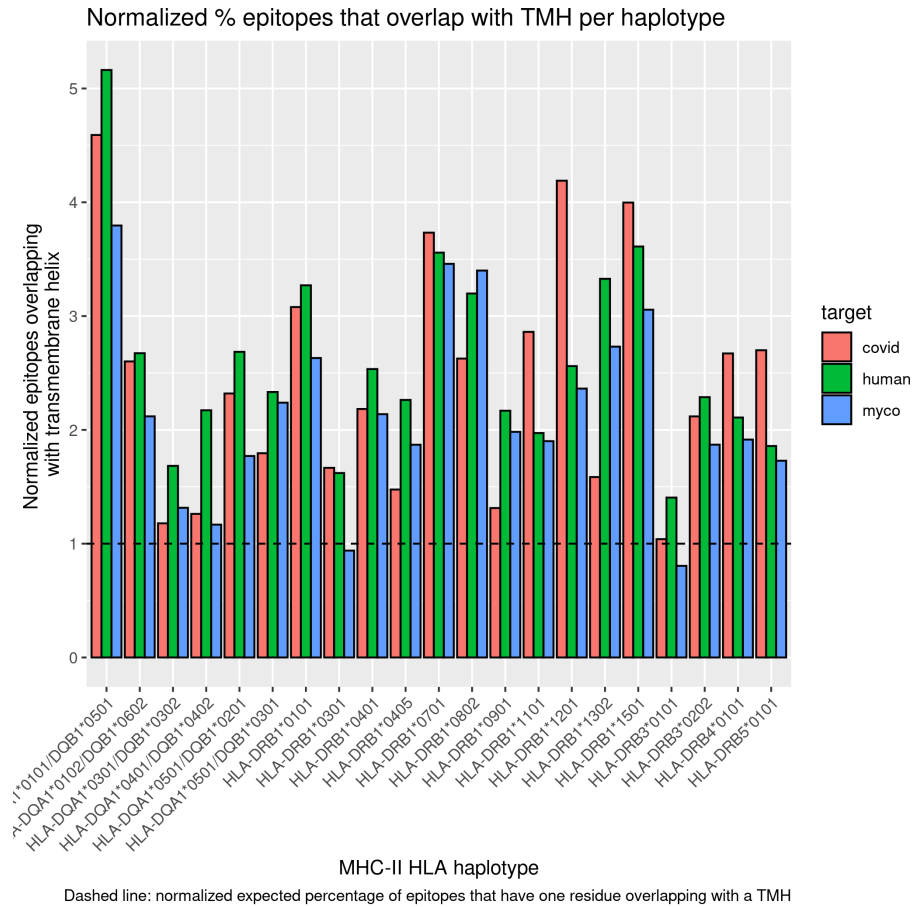


Figure 7: Normalized proportion of MHC-II epitopes overlapping with TMHs for human, viral and bacterial proteomes. Legend: covid = SARS-CoV-2, human = *Homo sapiens*, myco = *Mycobacterium tuberculosis*

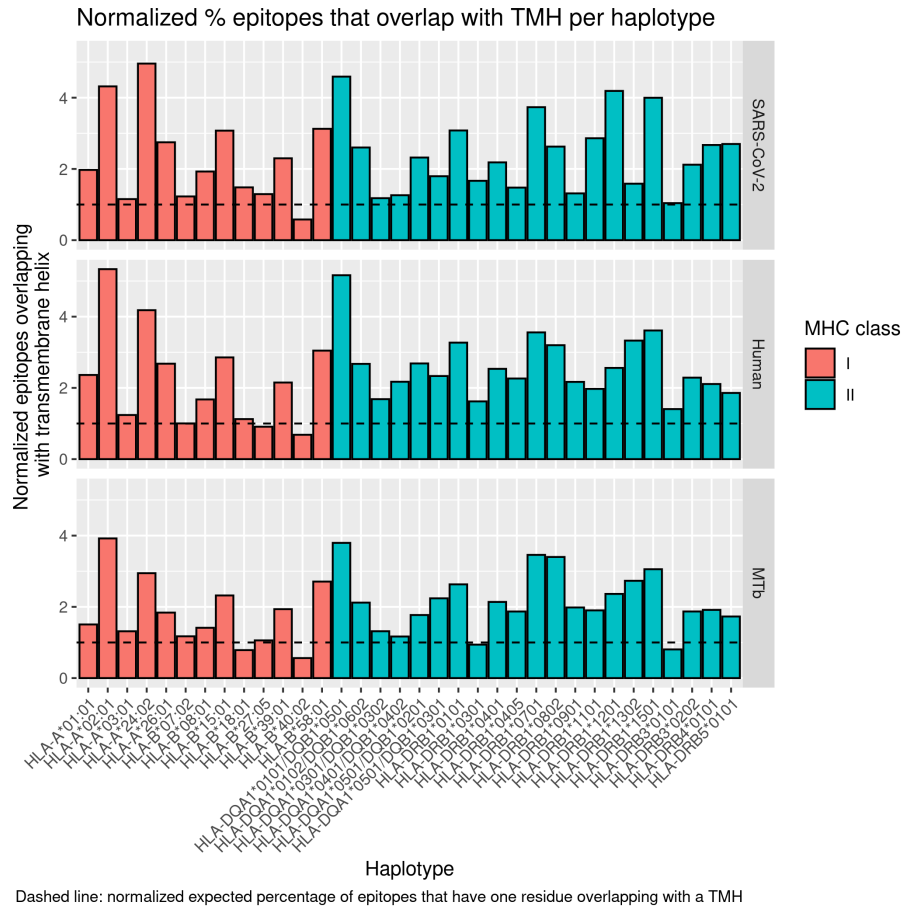


Figure 8: Normalized proportion of MHC-I and MHC-II epitopes overlapping with TMHs, for the different haplotypes and proteomes

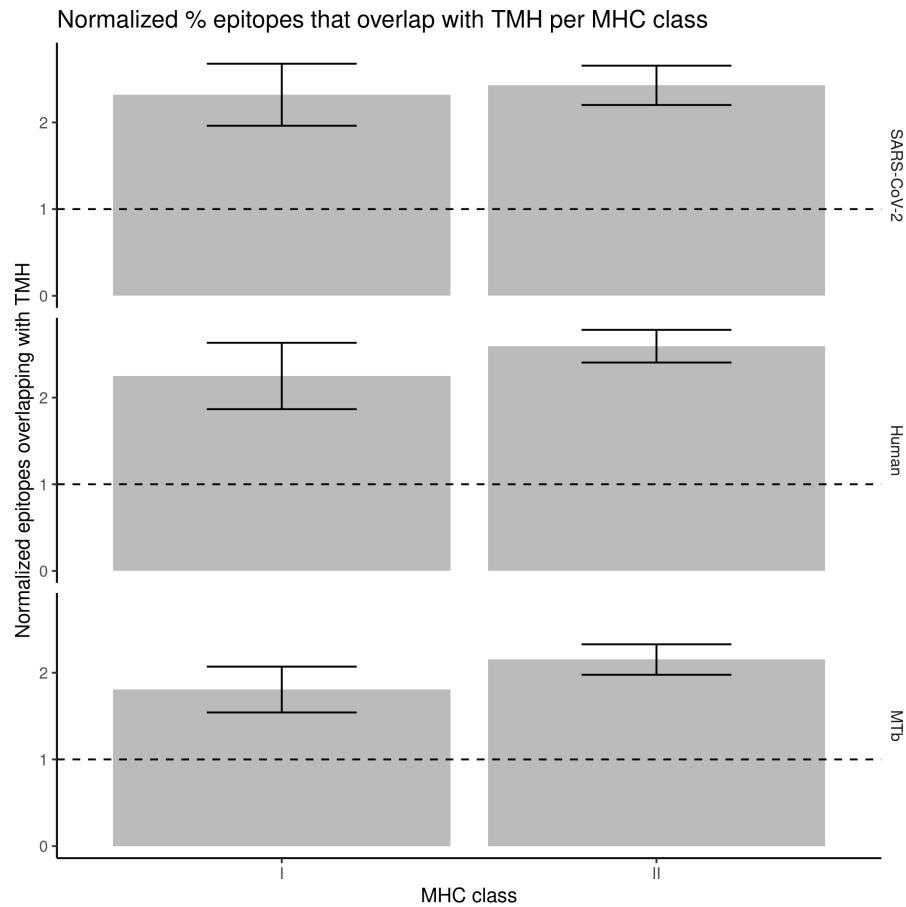


Figure 9: Normalized proportion of MHC-I and MHC-II epitopes overlapping with TMHs, for the different MHC classes and proteomes. Error bars denote the standard error.

811 the date we started the experiment, at December 14th 2020.

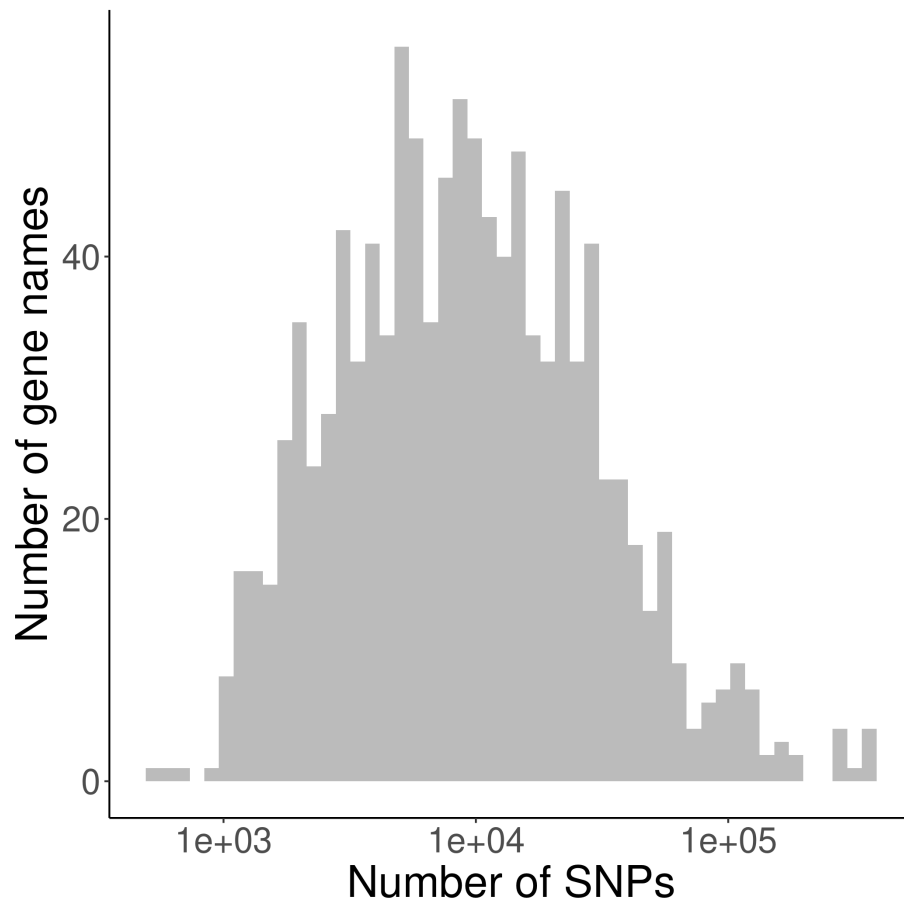


Figure 10: Distribution of the number of SNPs per gene name in the NCBI database.

812 To verify if SNPs were sampled uniformly over proteins, we show the dis-
813 tribution of the relative position in figure 15. We find no clear evidence of a
814 bias.

815 Supplementary Table 11 shows the statistics for all SNPs, where supplemen-
816 tary Tables 12 and 13 show the statistics for only single-spanners and multi-
817 spanners respectively.

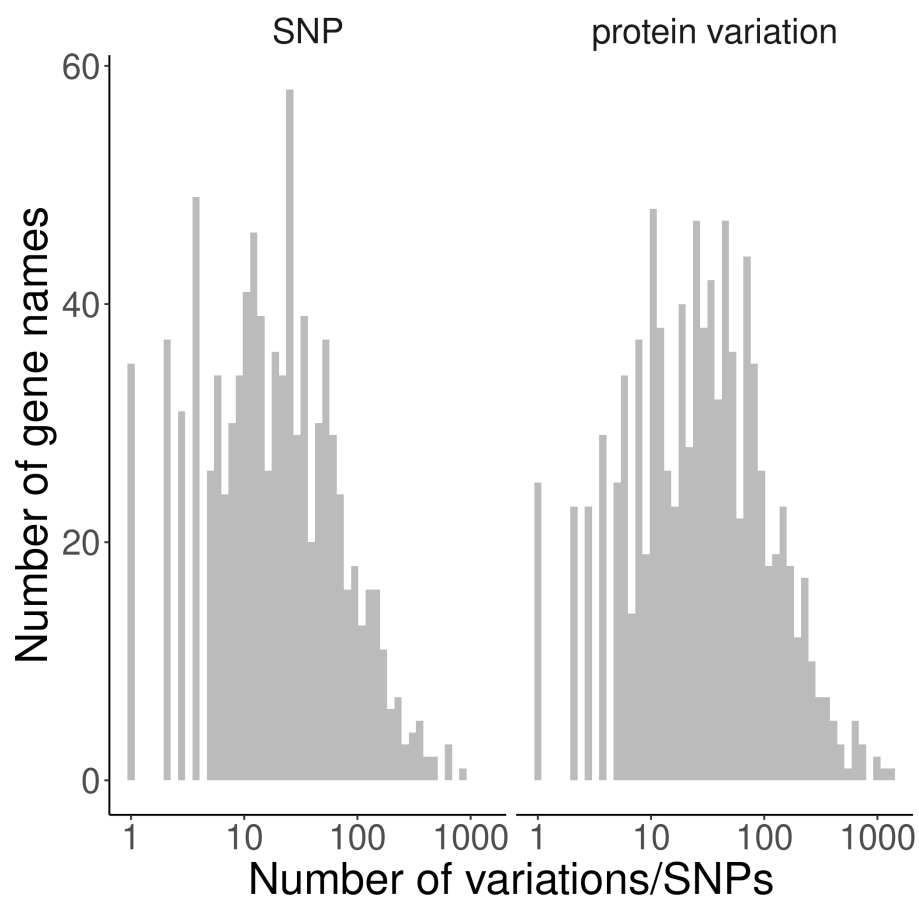


Figure 11: Distribution of the number of protein variations and SNPs per gene name processed.

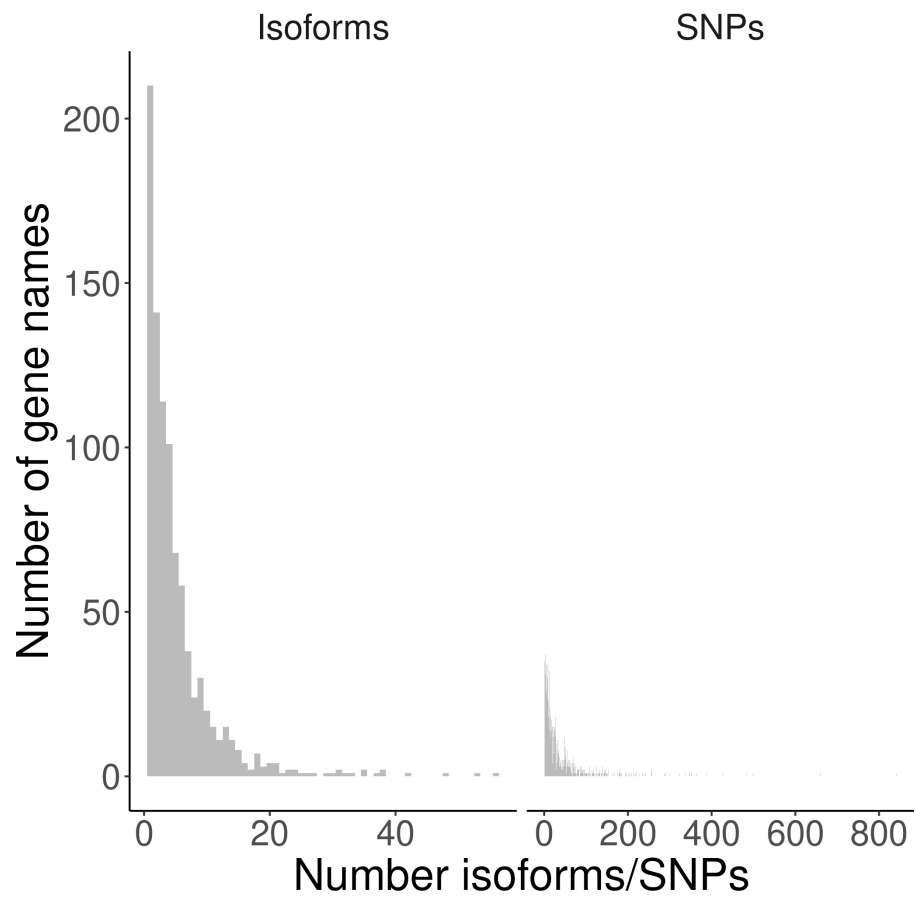


Figure 12: Histogram of the number of proteins found per gene name. Most often, a gene name is associated with one proteins.

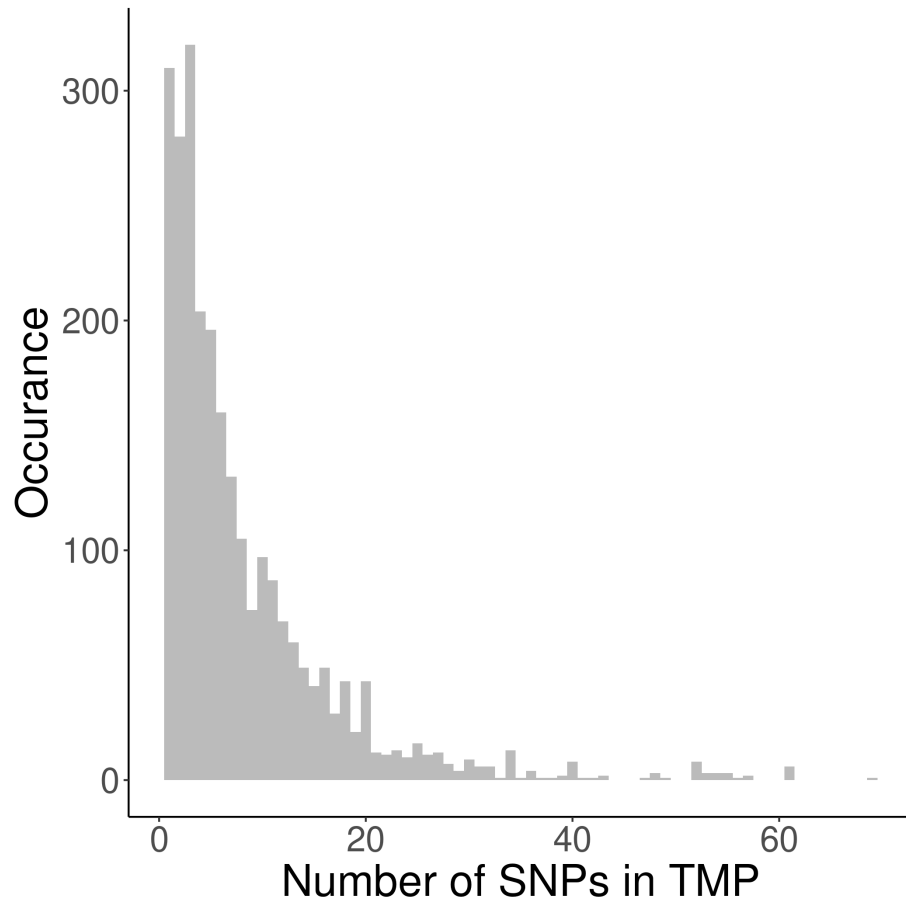


Figure 13: Histogram of the number of SNPs per trans-membrane protein. Dashed vertical line: average number of SNPs per TMP

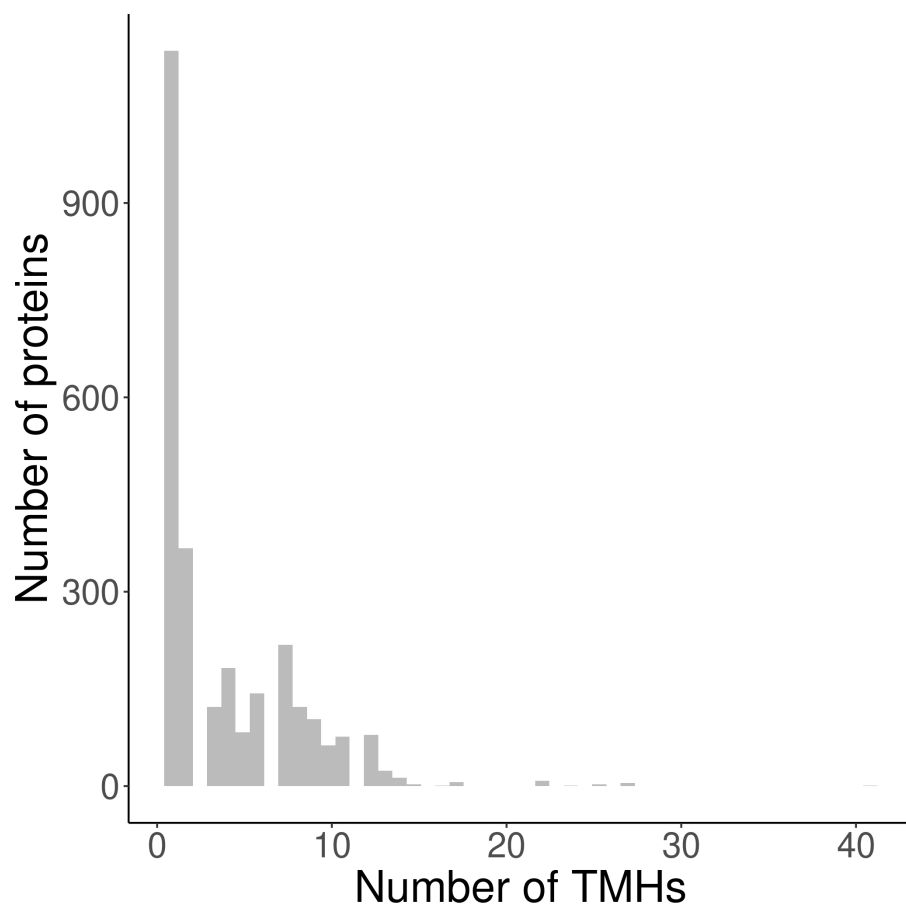


Figure 14: Histogram of the number of TMHs predicted per protein, for the trans-membrane proteins used.

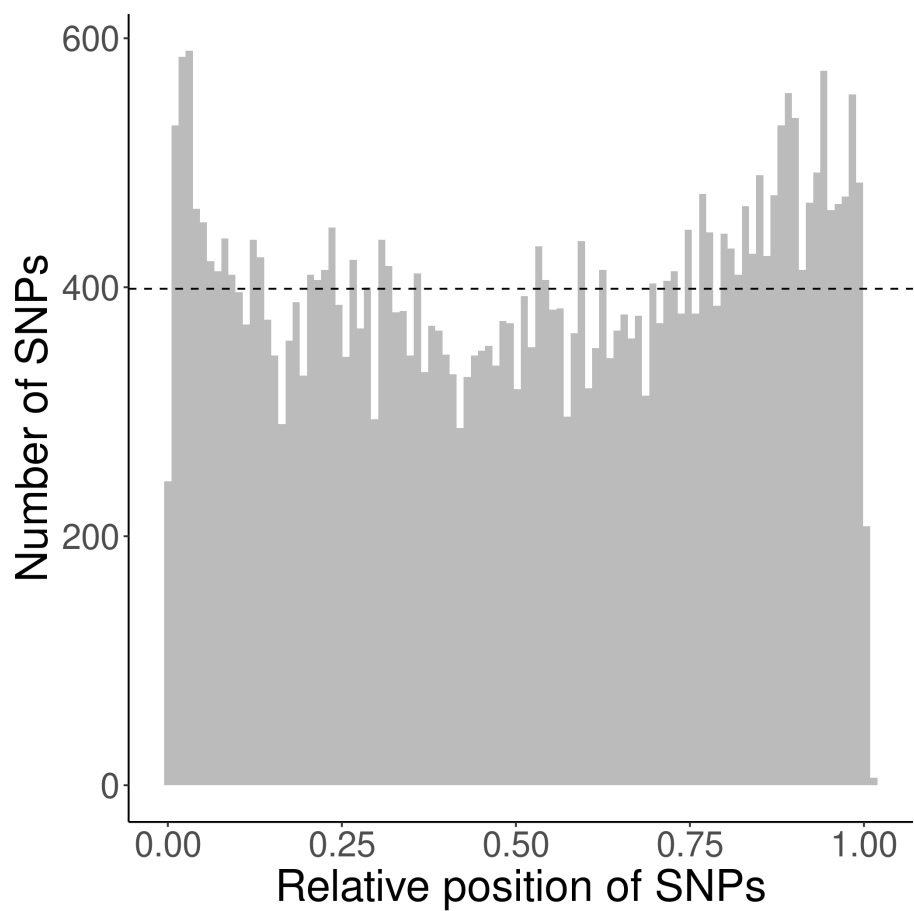


Figure 15: Distribution of the relative position of the SNPs used, where a relative position of zero denotes the first amino acid at the N-terminus, where a relative position of one indicates the last residue at the C-terminus.

index	haplotype_name
1	HLA-A*01:01
2	HLA-A*02:01
3	HLA-A*03:01
4	HLA-A*24:02
5	HLA-A*26:01
6	HLA-B*07:02
7	HLA-B*08:01
8	HLA-B*18:01
9	HLA-B*27:05
10	HLA-B*39:01
11	HLA-B*40:02
12	HLA-B*58:01
13	HLA-B*15:01
1	HLA-DRB1*0101
2	HLA-DRB1*0301
3	HLA-DRB1*0401
4	HLA-DRB1*0405
5	HLA-DRB1*0701
6	HLA-DRB1*0802
7	HLA-DRB1*0901
8	HLA-DRB1*1101
9	HLA-DRB1*1201
10	HLA-DRB1*1302
11	HLA-DRB1*1501
12	HLA-DRB3*0101
13	HLA-DRB3*0202
14	HLA-DRB4*0101
15	HLA-DRB5*0101
16	HLA-DQA1*0501/DQB1*0201
17	HLA-DQA1*0501/DQB1*0301
18	HLA-DQA1*0301/DQB1*0302
19	HLA-DQA1*0401/DQB1*0402
20	HLA-DQA1*0101/DQB1*0501
21	HLA-DQA1*0102/DQB1*0602

Table 6: Abbreviations of the haplotype names

Table 7: Percentage of MHC-I 9-mers overlapping with TMH. Values in brackets show the number of binders that have at least one residue overlapping with a TMH (first value) as well as the number of binders (second value). percentage used: 2

haplotype	covid	human	myco
HLA-A*01:01	15.603 (44/282)	12.600 (28377/225209)	11.424 (2947/25797)
HLA-A*02:01	34.155 (97/284)	28.441 (63994/225003)	29.749 (7646/25702)
HLA-A*03:01	9.122 (27/296)	6.606 (14851/224796)	9.972 (2565/25721)
HLA-A*24:02	39.223 (111/283)	22.297 (50313/225648)	22.346 (5752/25741)
HLA-A*26:01	21.739 (65/299)	14.287 (32232/225598)	13.950 (3598/25793)
HLA-B*07:02	9.712 (27/278)	5.347 (11893/222429)	8.899 (2291/25744)
HLA-B*08:01	15.248 (43/282)	8.935 (19981/223616)	10.714 (2750/25667)
HLA-B*15:01	24.324 (72/296)	15.228 (34498/226542)	17.600 (4547/25835)
HLA-B*18:01	11.724 (34/290)	5.993 (13409/223745)	5.960 (1536/25773)
HLA-B*27:05	10.227 (27/264)	4.854 (10882/224178)	8.031 (2063/25688)
HLA-B*39:01	18.182 (50/275)	11.468 (25621/223419)	14.682 (3787/25793)
HLA-B*40:02	4.594 (13/283)	3.647 (8147/223408)	4.264 (1097/25729)
HLA-B*58:01	24.731 (69/279)	16.245 (36409/224119)	20.558 (5292/25742)

Table 8: Percentage of MHC-II 14-mers overlapping with TMH. Values in brackets show the number of binders that have at least one residue overlapping with a TMH (first value) as well as the number of binders (second value). percentage used: 2

haplotype	covid	human	myco
HLA-DQA1*0101/DQB1*0501	40.433 (112/277)	31.214 (69752/223464)	32.158 (8187/25459)
HLA-DQA1*0102/DQB1*0602	22.910 (74/323)	16.167 (35753/221147)	17.950 (4608/25671)
HLA-DQA1*0301/DQB1*0302	10.381 (30/289)	10.179 (22623/222248)	11.144 (2842/25502)
HLA-DQA1*0401/DQB1*0402	11.111 (32/288)	13.135 (29319/223219)	9.890 (2524/25522)
HLA-DQA1*0501/DQB1*0201	20.430 (57/279)	16.240 (36186/222820)	14.999 (3823/25489)
HLA-DQA1*0501/DQB1*0301	15.808 (46/291)	14.106 (31046/220089)	18.969 (4878/25715)
HLA-DRB1*0101	27.119 (80/295)	19.774 (43968/222349)	22.293 (5692/25533)
HLA-DRB1*0301	14.676 (43/293)	9.801 (21831/222752)	7.956 (2025/25451)
HLA-DRB1*0401	19.231 (55/286)	15.325 (34011/221930)	18.113 (4641/25623)
HLA-DRB1*0405	12.996 (36/277)	13.684 (30380/222012)	15.837 (4036/25484)
HLA-DRB1*0701	32.877 (96/292)	21.512 (47856/222465)	29.304 (7471/25495)
HLA-DRB1*0802	23.132 (65/281)	19.339 (42859/221623)	28.805 (7358/25544)
HLA-DRB1*0901	11.565 (34/294)	13.111 (29043/221520)	16.798 (4301/25605)
HLA-DRB1*1101	25.197 (64/254)	11.924 (26582/222928)	16.103 (4101/25467)
HLA-DRB1*1201	36.897 (107/290)	15.482 (34596/223464)	20.018 (5098/25467)
HLA-DRB1*1302	13.962 (37/265)	20.121 (44798/222646)	23.141 (5935/25647)
HLA-DRB1*1501	35.206 (94/267)	21.836 (48671/222893)	25.891 (6584/25430)
HLA-DRB3*0101	9.158 (25/273)	8.496 (18884/222274)	6.819 (1740/25517)
HLA-DRB3*0202	18.657 (50/268)	13.832 (30687/221859)	15.843 (4059/25620)
HLA-DRB4*0101	23.529 (68/289)	12.749 (28376/222568)	16.221 (4131/25467)
HLA-DRB5*0101	23.776 (68/286)	11.235 (24993/222464)	14.648 (3732/25478)

Table 9: Amounts. raw = all variations, including DNA variations. all_proteins = all proteins. map = membrane associated protein. tmp = transmembrane protein. in_tmh = in transmembrane helix of TMP. in_sol = in soluble region of TMP.

what	raw	all_proteins	map	tmp	in_tmh	in_sol
Number of variations	60931	37831	16623	21208	3803	17405
Number of unique variations	60544	37630	16606	21024	3789	17235
Number of unique SNPs	NA	9621	4219	6026	1140	4936
Number of unique gene names	953	911	457	605	325	590
Number of unique protein names	5163	4780	2227	2553	1280	2467
Percentage TMH	NA	10	0	19	26	18

Table 10: Amounts. `single_in_tmh` = in transmembrane helix of single-spanner. `single_in_sol` = in soluble region of single-spanner. `multi_in_tmh` = in transmembrane helix of multi-spanner. `multi_in_sol` = in soluble region of multi-spanner.

what	single_in_tmh	single_in_sol	multi_in_tmh	multi_in_sol
Number of variations	452	7734	3351	9671
Number of unique variations	451	7733	3338	9502
Number of unique SNPs	160	2393	994	2762
Number of unique gene names	96	282	243	344
Number of unique protein names	304	1032	976	1435
Percentage TMH	11	5	35	26

Table 11: Statistics for all TMPs. `p` = p value. `n` = number of SNPs. `n_success` = number of SNPs found in TMHs (dashed blue line). `E(n_success)` = expected number of SNPs to be found in TMHs (dashed red line).

parameter	value
<code>p</code>	6.820823e-11
<code>n</code>	21208
<code>n_success</code>	3803
<code>E(n_success)</code>	4140.56

Table 12: Statistics for the single-spanners. `p` = p value. `n` = number of SNPs in single-spanners. `n_success` = number of SNPs found in TMHs of single-spanners (dashed blue line). `E(n_success)` = expected number of SNPs to be found in TMHs of single-spanners (dashed red line).

parameter	value
<code>p</code>	0.3189532
<code>n</code>	8186
<code>n_success</code>	452
<code>E(n_success)</code>	462.1535

Table 13: Statistics for the multi-spanners. `p` = p value. `n` = number of SNPs in multi-spanners. `n_success` = number of SNPs found in TMHs of multi-spanners (dashed blue line). `E(n_success)` = expected number of SNPs to be found in TMHs of multi-spanners (dashed red line).

parameter	value
<code>p</code>	8.315841e-12
<code>n</code>	13022
<code>n_success</code>	3351
<code>E(n_success)</code>	3678.406