

1 **INeo-Epp: A novel T-cell HLA class-I immunogenicity or** 2 **neoantigenic epitope prediction method based on sequence** 3 **related amino acid features**

4 Guangzhi Wang^{1,2}, Huihui Wan^{2,3}, Xingxing Jian^{2,4}, Yuyu Li¹, Jian Ouyang²,
5 XiaoxiuTan³, Yong Zhao^{1*}, Yong Lin^{3*}, Lu Xie^{1,2*}

6 ¹ College of Food Science and Technology, Shanghai Ocean University, Shanghai,
7 201306, China

8 ² Shanghai Center for Bioinformation Technology, Shanghai Academy of Science and
9 Technology, Shanghai, 201203, China

10 ³ School of Medical Instrument and Food Engineering, University of Shanghai for
11 Science and Technology, Shanghai, 200093, China

12 ⁴ Key Laboratory of Carcinogenesis and Cancer Invasion, Ministry of Education; Key
13 Laboratory of Carcinogenesis, National Health and Family Planning Commission,
14 Xiangya Hospital, Central South University, Changsha, 410008, China.

15 Correspondence should be addressed to Lu Xie; luxiex2017@outlook.com

16 **Abstract**

17 In silico T-cell epitope prediction plays an important role in immunization experimental
18 design and vaccine preparation. Currently, most epitope prediction research focuses on
19 peptide processing and presentation, e.g. proteasomal cleavage, transporter associated
20 with antigen processing (TAP) and major histocompatibility complex (MHC)
21 combination. To date, however, the mechanism for immunogenicity of epitopes remains
22 unclear. It is generally agreed upon that T-cell immunogenicity may be influenced by
23 the foreignness, accessibility, molecular weight, molecular structure, molecular
24 conformation, chemical properties and physical properties of target peptides to different
25 degrees. In this work, we tried to combine these factors. Firstly, we collected significant
26 experimental HLA-I T-cell immunogenic peptide data, as well as the potential
27 immunogenic amino acid properties. Several characteristics were extracted, including
28 amino acid physicochemical property of epitope sequence, peptide entropy, eluted
29 ligand likelihood percentile rank (EL rank(%)) score and frequency score for
30 immunogenic peptide. Subsequently, a random forest classifier for T cell immunogenic
31 HLA-I presenting antigen epitopes and neoantigens was constructed. The classification
32 results for the antigen epitopes outperformed the previous research (the optimal
33 AUC=0.81, external validation data set AUC=0.77). As mutational epitopes generated
34 by the coding region contain only the alterations of one or two amino acids, we assume
35 that these characteristics might also be applied to the classification of the endogenic
36 mutational neoepitopes also called ‘neoantigens’. Based on mutation information and
37 sequence related amino acid characteristics, a prediction model of neoantigen was
38 established as well (the optimal AUC=0.78). Further, an easy-to-use web-based tool
39 ‘INeo-Epp’ was developed (available at [http://www.biostatistics.online/INeo-](http://www.biostatistics.online/INeo-Epp/neoantigen.php)
40 [Epp/neoantigen.php](http://www.biostatistics.online/INeo-Epp/neoantigen.php)) for the prediction of human immunogenic antigen epitopes and
41 neoantigen epitopes.

42 Introduction

43 An antigen consists of several epitopes, which can be recognized either by B- or T-cells
44 and/or molecules of the host immune system. However, usually only a small number of
45 amino acid residues that comprise a specific epitope are necessary to elicit an immune
46 response [1]. The properties of these amino acid residues causing immunogenicity are
47 unknown. HLA-I antigen peptides are processed and presented as follows: a). cytosolic
48 and nuclear proteins are cleaved to short peptides by intracellular proteinases; b). some
49 are selectively transferred to endoplasmic reticulum (ER) by TAP transporter, and
50 subsequently are treated by endoplasmic reticulum aminopeptidase;c). antigen
51 presenting cells (APCs) present peptides containing 8-11 AA (amino acid) residues on
52 HLA class I molecules to CD8+ T cells [2]. Researchers can now simulate antigen
53 processing and presentation by computational methods to predict binding peptide-MHC
54 complexes (p-MHC). Several types of software systems have been developed,
55 including NetChop [3], NetCTL [4], NetMHCpan [5], MHCflurry [6]. However, the
56 binding to MHC molecules of most peptides is predicted, only 10%~15% of those have
57 been shown to be immunogenic [7-10]. For neoantigens the result was approximately
58 5% (range, 1%-20%) due to central immunotolerance [11, 12]. As a result, the cycle for
59 vaccine development and immunization research is extended. Here, we aim to develop
60 a T-cell HLA class-I immunogenicity prediction method to further identify real
61 epitopes/neoepitopes from p-MHC to shorten this cycle.

62 Many experimental human epitopes have been collected and summarized in the
63 immune epitope database (IEDB) [13], which makes it feasible to mathematically
64 predict human epitopes. However there still exist two limitations: i) a high level of
65 MHC polymorphism produces a severe challenge for T-cell epitope prediction. ii) there
66 is an extremely unequal distribution of data to compare epitopes and non-epitopes. It is
67 not conducive to analyze the potential deviation existing in TCR recognition owing to
68 the presentation of different HLA peptides. A general analysis of all HLA presented
69 peptides, ignoring the specific pattern of TCR recognition of individual HLA presented
70 peptides, may result in a lower predictive accuracy.

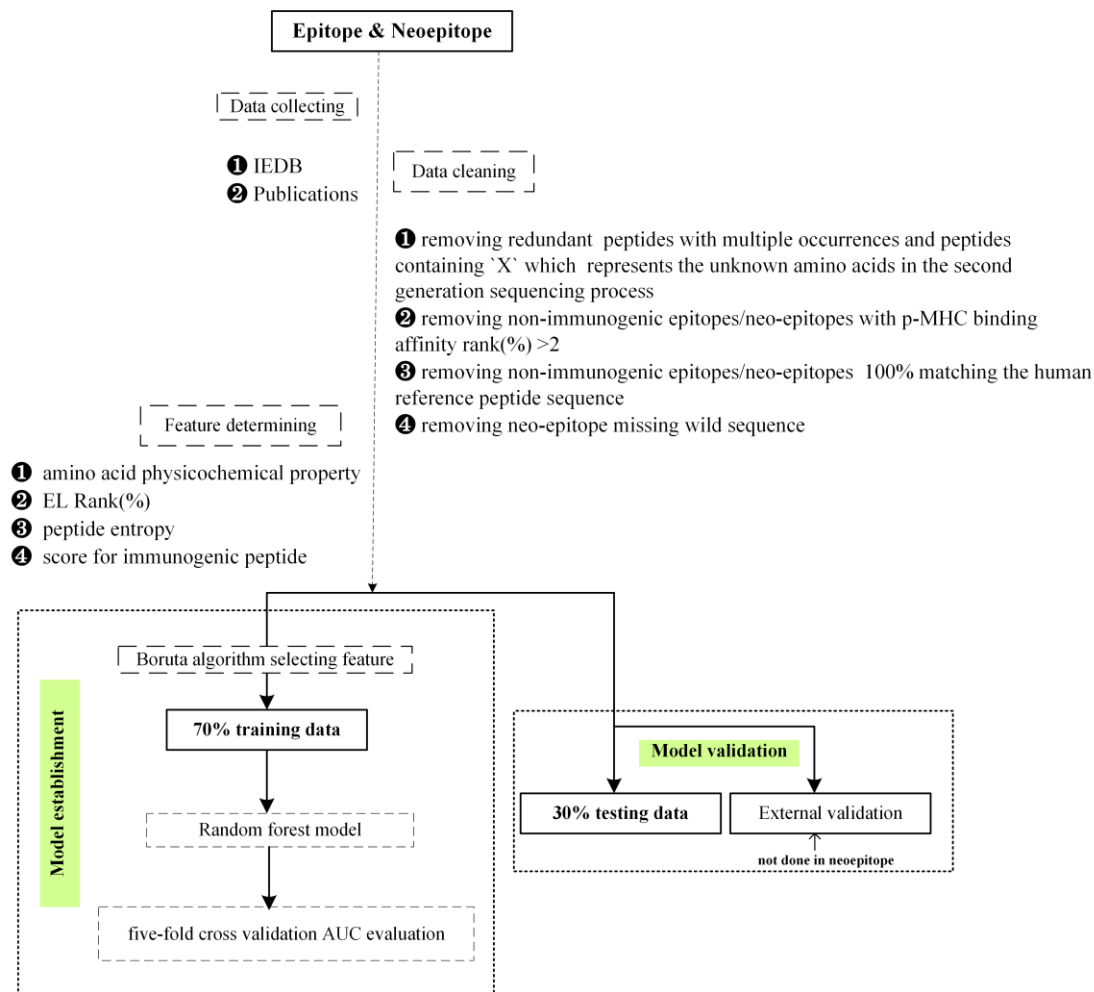
71 With the advances in HLA research, Sette *et al* [14] classified, for the first time,
72 overlapping peptide binding repertoires into nine major functional HLA supertypes (A1,
73 A2, A3, A24, B7, B27, B44, B58, B62). In 2008, John Sidney *et al* [15] made a further
74 refinement, in which over 80% of the 945 different HLA-A and -B alleles can be
75 assigned to the original nine supertypes. It has not been reported whether peptides
76 presented by different HLA alleles influence TCR recognition. Hence, we collected
77 experimental epitopes according to HLA alleles and assume that epitopes belonging to
78 the same HLA supertypes have similar properties.

79 Moreover, screening for endogenous mutational neoepitopes is one of the core steps
80 in tumor immunotherapy. In 2017, Ott PA *et al.* [16]and Sahin *et al* [17]. confirmed that
81 peptides and RNA vaccines made up of neoantigens in melanoma can stimulate and
82 proliferate CD8+ and CD4+ T cells. In addition, a recent research suggests that
83 including neoantigen vaccination not only can expand the existing specific T cells, but
84 also induce a wide range of novel T-cell specificity in cancer patients and enhance
85 tumor suppression[18]. Meanwhile, a tumor can be better controlled by the combination
86 therapy of neoantigen vaccine and programmed cell death protein 1 (PD-1)/PD1 ligand
87 1(PDL-1) therapy [19, 20]. Nevertheless, a considerable number of predicted candidate
88 p-MHC from somatic cell mutations may be false positive, which would fail to
89 stimulate TCR recognition and immune response. This is undoubtedly a challenge for
90 designing vaccines against neoantigens.

91 In our study, based on HLA-I T-cell peptides collected from experimentally
92 validated antigen epitopes and neoantigen epitopes, we aim to build a novel method to
93 further reduce the range of immunogenic epitopes screening based on predicted p-MHC.
94 Finally, a simple web-based tool, INeo-Epp (immunogenic epitope/neoepitope
95 prediction), was developed for prediction of human antigen and neoantigen epitopes.

96 Materials and Methods

97 The flow chart for 'INeo-Epp' prediction is shown as follows. (see Figure 1)



98

99

Figure 1: The flow chart for 'INeo-Epp' prediction

100 Construction of immunogenic and non-immunogenic epitopes

101 Peptides that can promote cytokine proliferation are considered to be immunogenic
102 epitopes. However, non-immunogenic epitopes may result for the following reasons: a)
103 p-MHC truly unrecognized by TCR; b) peptides not presented by MHC (quantitatively
104 expressed as rank(>) > 2, see rank(>) score (below: C24) for details); c) negative
105 selection/clonal presentation induced by excessive similarity to autologous
106 peptides[21]. In this work, to further study the recognition preferences of T cells,

107 peptides with >2 rank(%) were regarded as not in contact with TCR, and sequences
 108 100% matching the human reference peptides ([ftp://ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-97/fasta/homo_sapiens/pep/)
 109 [97/fasta/homo_sapiens/pep/](ftp://ftp.ensembl.org/pub/release-97/fasta/homo_sapiens/pep/)) were regarded as exhibiting immune tolerance. Hence,
 110 we removed these from the definition of non-immunogenic peptides.

111 Construction of data sets: epitopes, external validation epitopes and neoepitopes

112 Antigen epitope data were collected from IEDB (Linear epitope, Human, T cell assays,
 113 MHC class I, any disease were chosen). Data collection criteria: each HLA allele
 114 quantity >50 and frequency >0.5% (refer to allele frequency database [22]) (Table 1,
 115 check Table S1 for detailed information).

116 Table 1: Summary of IEDB epitope data

HLA supertype	IEDB HLA data	Number		HLA allele frequency Asian / Black / Caucasian	Motif view
		Negative	Positive		
A1	A01:01	811	103	0.154 / 0.046 / 0.164	1-2(ST)-3-4-5-6-7-8-9(Y)
	A26:01	83	19	0.041 / 0.014 / 0.030	1(DE)-2(ITV)-3-4-5-6-7-8-9(FMY)
A2	A02:01	1883	1580	0.049 / 0.123 / 0.275	1-2(LM)-3-4-5-6-7-8-9(ILV)-10(V)
	A11:01	196	174	0.139 / 0.014 / 0.060	1-2(IMSTV)-3-4-5-6-7-8-9(K)-10(K)
A3	A03:01	1400	169	0.063 / 0.083 / 0.139	1-2(ILMTV)-3-4-5-6-7-8-9(K)-10(K)
	A24:02	207	219	0.136 / 0.024 / 0.084	1-2(WY)-3-4-5-6-7-8-9(FIW)
A24	A23:01	1138	12	0.006 / 0.109 / 0.019	1-2(WY)-3-4-5-6-7-8-9-10(F)
	B7	B35:01	63	248	0.062 / 0.068 / 0.055
B7	B07:02	523	244	0.034 / 0.005 / 0.0143	1-2(p)-3-4-5-6-7-8-9(FLM)
	B51:01	13	51	0.074 / 0.021 / 0.047	1-2(P)-3-4-5-6-7-8-9(IV)
B8	B08:01	317	195	0.036 / 0.037 / 0.114	1-2-3-4-5(HKR)-6-7-8-9(FILMV)
B27	B27:05	100	86	0.008 / 0.008 / 0.037	1(RY)-2(R)-3(FMLWY)-4-5-6-7-8-9
B44	B37:01	1036	10	0.034 / 0.005 / 0.014	-
	B40:01	67	65	0.022 / 0.012 / 0.052	-
B44	B44:02	73	66	0.008 / 0.020 / 0.095	1-2(E)-3-4-5-6-7-8-9(FIWY)
	B58	B58:01	11	62	0.041 / 0.037 / 0.007
B62	B15:01	3	70	0.016 / 0.010 / 0.060	1-2(LMQ)-3-4-5-6-7-8-9(FY)
Total		7924	3373		
Remove negative rank(>2)		5123	3373		
Remove negative human 100% similar		4943	3373		

117 The external antigen epitope validation set was collected from seven published
 118 independent human antigen studies [23-29], consisting of 577 non-immunogenic
 119 epitopes and 85 immunogenic epitopes (Table 2, S2 Table)

120 Table 2: External data included in validation set

Publication time	PMID	Author	non-epitopes	epitopes
2013	23580623	Weiskopf et al	477	42
2018	29397015	Hendrik Luxenburger et al	100	26
2018	30260541	Youchen Xia et al	-	1
2018	30487281	Hawa Vahed et al	-	4
2018	30518652	Atefeh Khakpoor et al	-	2
2018	30587531	Alina Huth et al	-	4
2018	30815394	Solomon Owusu Sekyere et al	-	6
Total			577	85
Remove negative with rank(>2) and HLA supertypes (not appeared in training set)			321	69

121 Here, we removed peptides for which HLA supertypes do not appear in training set,
 122 because we assume peptides belonging to the same HLA supertypes to have similar
 123 properties. In the external validation set, some peptides bind to rare HLA supertypes.
 124 Their characteristics were not included in the training set. Hence, these peptides in the
 125 external validation data might lead to a classification bias.

126 The neoantigens data were collected from 11 publications [19, 30-39] and IEDB
 127 mutational epitopes, and 13 published data sets collected by Anne-Mette B in one

128 publication [40] in 2017 (see Table 3, S3 Table for details) were also included.

129 Table 3: Neoepitopes data included in this study

Publication time	PMID	Author	Tumor Type	Non-immunogenic neo-epitopes	Immunogenic neo-epitopes	T-cell assay
2013-12	24323902	Darin A. W et al.	Ovarian Cancer	—	1	ELISPOT
2015-9	26359337	Eliezer M et al.	Melanoma	—	18	Clinical benefit
2015-11	26752676	Takahiro K et al.	Lung adenocarcinoma	—	4	—
2016-1	26901407	Alena Gros et al.	Melanoma	12	14	ELISPOT
2016-5	27198675	Erlend Strønen et al.	Melanoma	1134	16	CTL clone
2016-12	28405493	Annika Nelde et al.	Lymphoma	—	2	ELISPOT
2017-6	28619968	Xiuli Zhang et al.	Breast cancer	—	4	Flow cytometry
2017-10	29104575	Markus M et al.	Melanoma	10	16	—
2017-11	29187854	Anne-Mette B et al.	Polytype	1874	42	ELISPOT et al.
2017-11	29132146	Vinod P. B et al.	pancreatic	—	10	Flow Cytometry
2018-5	29720506	Tatsuo Matsuda et al.	Ovarian Cancer	—	3	ELISPOT
2018-12	29409514	Sonntag et al.	pancreatic ductal carcinoma	—	3	Flow Cytometry
2018-10	30357391	Randi Vita et al.	—	6	35	—
Total				3030	168	
Remove duplication				2837	164	
Remove negative rank(>2 and human 100% similar)				1697	164	

130 Construction of potential immunogenicity feature

131 **Characteristics calculation of peptides based on amino acid sequences.** The formula
 132 for calculating peptide characteristics is shown in (1). P_N , P_2 , P_C (N-terminal, position
 133 2, C-terminal as anchored sites by default) are considered to be embedded in HLA
 134 molecules and no contact with TCRs, therefore not evaluated.

$$135 \quad P_c = \left\{ \sum_{x \in Pos(P)}^{x \in (N,2,C)} P_{Ac} \right\} / (\text{len}(P) - 3) \quad (1)$$

136 P , peptide. c , characteristic. Where P_c represents characteristics of peptides. A , amino
 137 acid. N , N-terminal in a peptide. C , C-terminal in a peptide. Pos , amino acid position in
 138 peptide. Where P_{Ac} represents characteristics of amino acids in peptides.

139 **Frequency score for immunogenic peptide (C22).** Amino acid distribution frequency
 140 differences between immunogenicity and non-immunogenic peptides at TCR contact
 141 sites (excluding anchor sites) were considered as a feature (2).

$$142 \quad P_{score} = \sum_{x \in Pos(P)}^{x \in (N,2,C)} \{ P_{ie^+}(f'_A) - P_{ie^-}(f'_A) \} \quad (2)$$

143 P_{ie^+} , immunogenic peptides. P_{ie^-} , non-immunogenic peptides. f'_A , amino acid frequency
 144 in TCR contact position. Where $P_{ie^+}(f'_A)$ represents frequency of amino acids in
 145 immunogenic peptides at TCR contact sites.

146 **Calculating peptide entropy (C23).** Peptide entropy [41] was used as a feature (3).

$$147 \quad P_H = \left\{ - \sum_{x \in Pos(P)}^{x \in (N,2,C)} P_{f_A} * \log_2(P_{f_A}) \right\} / (\text{len}(P) - 3) \quad (3)$$

148 P_H , peptide entropy. f_A , amino acid frequency in human reference peptide sequence.
 149 Where P_{f_A} represents the frequency in human reference peptide sequence of amino
 150 acids in epitope peptides.

151 **Rank(%) score (C24).** HLA binding prediction were performed using netMHCpan4.0.
 152 rank(%) provides a robust filter for the identification of MHC-binding peptides, in
 153 which rank(%) was recommended as an evaluation standard, rank(%)<0.5 as strong

154 binders, $0.5 < \text{rank}(\%) < 2$ as weak binders, $\text{rank}(\%) > 2$ as no binders.

155 **Five-fold cross-validation, feature selection, random forests and ROC generation.**

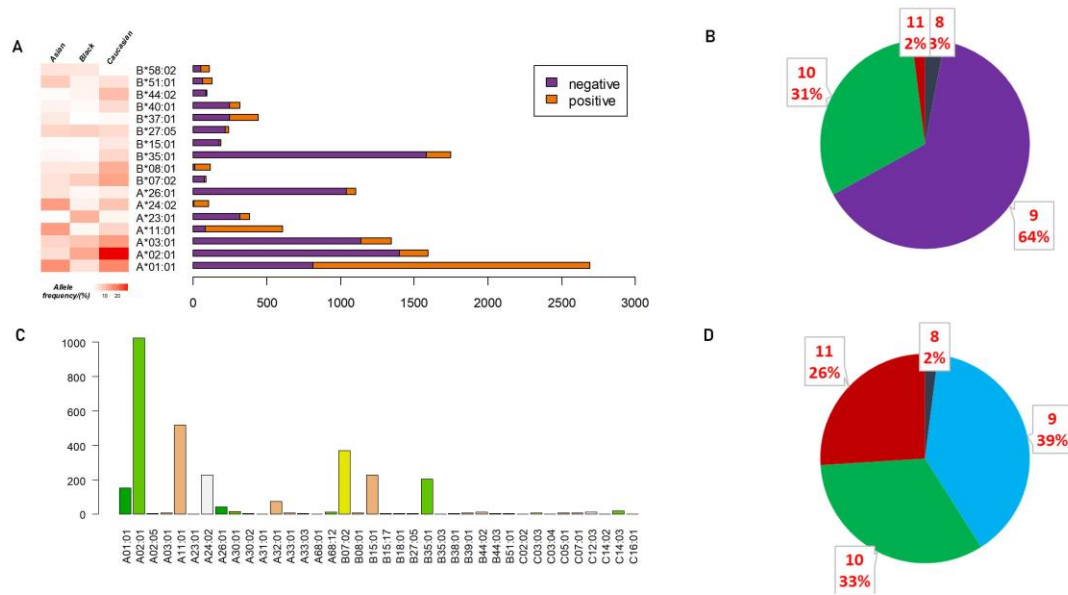
156 The 5-fold cross-validation was implemented in R using the package caret [42] (method
157 = "repeatedcv", number = 5, repeats = 3). The feature screening results were generated
158 in R using the package Boruta [43] (a novel random forest based feature selection
159 algorithm for finding all relevant variables, which provides unbiased and stable
160 selection of important and non-important attributes from an information system. It
161 iteratively removes the features which are proven by a statistical test to be less relevant
162 than random probes. It uses Z score (computed by dividing the average loss by its
163 standard deviation) as the importance measure and it takes into account the fluctuations
164 of the mean accuracy loss among trees in the forest). R package randomForest [44] was
165 used for training data (the R language machine learning package caret provides
166 automatic iteration selection of optimal parameters, mtry=15 for antigen epitope,
167 mtry=14 for neoantigen epitope, the remaining parameters use default values). R
168 package ROCR [45] was used for drawing ROC.

169 **Web tool implementation**

170 The front-end of Ineo-Epp was constructed via HTML/JavaScript/CSS. The back end
171 was written in PHP, connecting the web interface and Apache web server. A python
172 script was used for calculating peptide characteristics and extracting mutation
173 information. Models were built using R.

174 **Results**

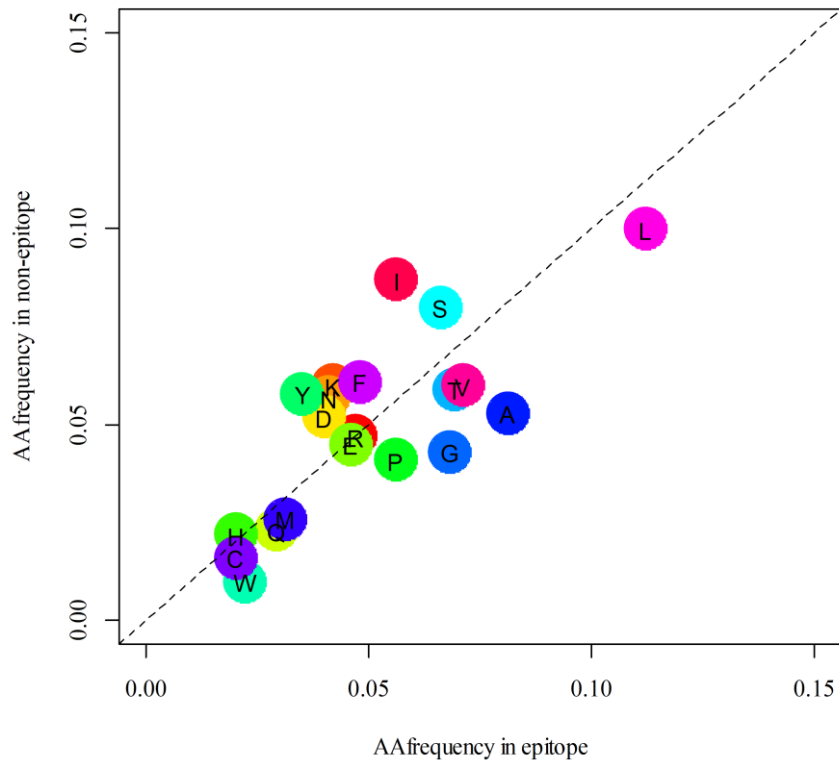
175 Ultimately, 11,297 validated epitopes and non-epitopes with the length of 8-11 amino
176 acids were collected from IEDB. T-cell responses included activation, cytotoxicity,
177 proliferation, IFN- γ release, TNF release, granzyme B release, IL-2 release, IL-10
178 release, etc. Seventeen different HLA alleles were collected (Fig 2A), and the detailed
179 antigen length distribution is shown in (Fig 2B). Additionally, we collected the
180 neoantigen data from 12 publications, including 2837 non-neoepitopes and 164
181 neoepitopes (Fig 2C), and the detailed neoantigen length distribution is shown in (Fig
182 2D).



183

184 Figure 2: Epitope/neoepitope peptides composition and amino acid lengths distribution. (a) Detailed data
 185 distribution of seventeen HLA alleles of antigen peptides and proportion of each HLA allele (positive
 186 and negative) epitopes and the corresponding HLA frequency in Asian, Black, Caucasian. (b) Proportion
 187 of antigen peptides of 8-11 AA lengths. (c) Data distribution of HLA alleles of neoantigen peptides. (d)
 188 Proportion of neoantigen peptides of 8-11 AA lengths.

189 The TCR contact position plays a crucial role in the analysis of immunogenicity,
 190 as TCRs might be more sensitive to some amino acids, the amino acids preference in
 191 antigen epitope peptide and antigen non-epitope peptide was further analyzed after
 192 excluding anchor sites (N-terminal, position 2, C-terminal) (Fig 3). We found that TCRs
 193 tend to identify hydrophobic amino acids. For example, 3/4 hydrophobic amino acids
 194 (L, W, P, A, V, M) occur more frequently in immunogenicity epitopes. Charged amino
 195 acids (*e.g.* D, K) are enriched in non-epitopes whereas the rest of charged amino acids
 196 (R, H, E) show no difference. Based on the result in figure 3, the amino acid distribution
 197 difference at the TCR contact sites was regarded by us as one of the immunogenicity
 198 features (*i.e.* Frequency score for immunogenic peptide (C22)).



199

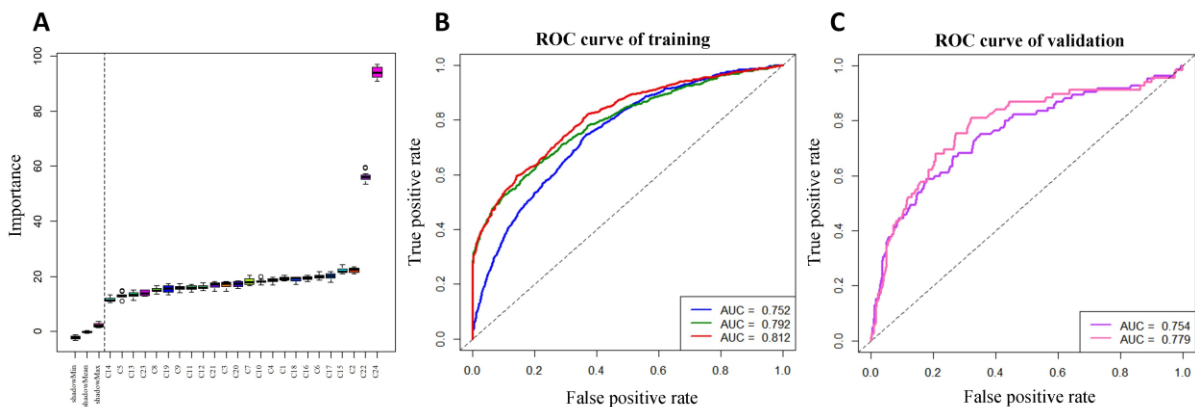
200 Figure 3: Antigen epitope amino acid distribution frequency in TCR contact site of epitopes and non-
201 epitopes. Frequency distribution of amino acids at TCR contact sites in antigen epitope and non-epitope
202 peptides, and the amino acids below the dotted line are preferred by the epitope.

203 **Classification prediction model for antigen epitopes**

204 We constructed the features of peptides on the basis of the characteristics of amino acids
205 (see Materials and Methods section: Characteristics Calculation of peptides based on
206 amino acids). All amino acid characteristics were selected from Protscale [46] in
207 ExPASy (SIB bioinformatics resource portal). The 21 involved features are as follows:
208 Kyte–Doolittle numeric hydrophobicity scale (C1) [47], molecular weight (C2),
209 bulkiness (C3) [48], polarity (C4) [49], recognition factors (C5) [50], hydrophobicity
210 (C6) [51], retention coefficient in HPLC (C7) [52], ratio hetero end/side (C8)[49],
211 average flexibility (C9) [53], beta-sheet (C10) [54], alpha-helix (C11) [55], beta-turn
212 (C12) [55], relative mutability (C13) [56], number of codon(s) (C14), refractivity (C15)
213 [57], transmembrane tendency (C16) [58], accessible residues (%) (C17) [59], average
214 area buried (C18) [60], conformational parameter for coil (C19) [55], total beta-strand
215 (C20) [60], parallel beta-strand (C21) [61] (see Table S4 in detail). Also, frequency
216 score for immunogenic peptide (C22), peptide entropy (C23) and rank(%) (C24) were
217 also taken into consideration. Together, 24 immunogenic features were collected, and
218 all features were retained for antigen epitopes prediction after screening using the R
219 package Boruta. Compared with other characteristics, the frequency score for

220 immunogenic peptide and rank(%) have higher impact, suggesting they have more
221 significant influence on antigen epitopes classification (Figure 4A).

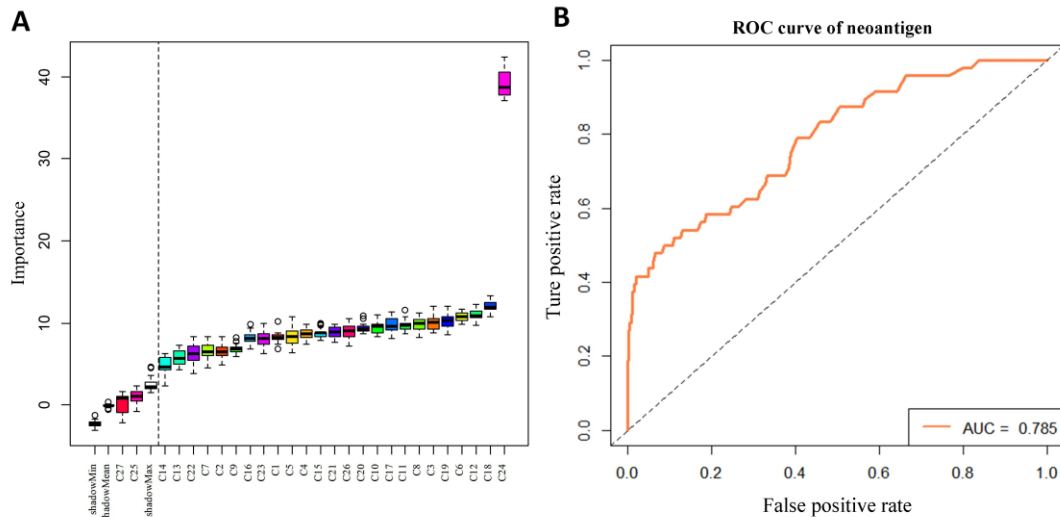
222 The receiver operator characteristic (ROC) curve of models are shown in Fig 4.
223 The five-fold cross validation AUC was 0.81 in the prediction model for antigen epitope
224 (line in red Fig 4B) and the externally validated (see table 2) AUC was 0.75 (line in
225 purple Fig 4C). Here, we tried to remove peptides for which HLA supertypes not
226 appearing in training set from the externally validated antigen data and, the AUC,
227 specificity, and sensitivity were increased to 0.78, 0.71, and 0.72, respectively. (line in
228 pink Fig4 C). This, to some extent, verifies our conjecture about TCR specific
229 recognition of different HLA alleles presenting peptides.
230



231 Figure 4: Feature selection in antigen epitopes and ROC curves of antigen epitopes classification.
232 (a)Peptide features: Twenty four features were screened and we defined the features on the right of the
233 dotted line as being effective. (b)Trained model: The line in blue represents antigen epitopes without
234 screening; the line in green represents selection with the deletion of rank(>2 non-epitope; and the line
235 in red represents selection with the deletion of the non-epitopes 100% matching human reference peptide
236 sequence. (c)External validation: The ROC curves for the external verification set, line in purple
237 represents modeling using antigen epitopes without filtering, the line in pink represents modeling using
238 antigen epitopes removing non-epitopes which rank(>2 and HLA for which supertypes not appearing
239 in training set.

240 Classification prediction model for neoantigen epitopes

241 Neoantigens derived from somatic mutations are different from the wild peptide
242 sequences. Therefore, some mutation-related characteristics were also taken into
243 account. For instance, difference in hydrophobicity before and after mutation (C25),
244 differential agretopicity index (DAI, C26) [62] and whether the mutation position was
245 anchored (C27). Finally, 27 features were selected for the neoantigen epitope prediction
246 model. However, only 25 neoantigen related features were retained after running Boruta,
247 because C25 and C27 were removed. Also, rank(%) showed a marked effect (Fig 5A).
248 in the five-fold cross-validation of the prediction model for neoantigen epitopes, AUC
249 was 0.78 (Fig 5B).



250

251 Figure 5: Feature selection in neoantigen epitopes and ROC curves of neoantigen epitopes classification.
252 (a) Twenty seven features were screened and the 25 features on the right of the dotted line were reserved
253 for modeling using a random forest algorithm. (b) ROC curves of neoantigen epitopes classification.

254 Web server for TCR epitope prediction

255 Based on these above-mentioned validated features, we established a web server for
256 TCR epitope prediction, named 'INeo-Epp'. This tool can be used to predict both
257 immunogenic antigen and neoantigen epitopes. For antigen, the nine main HLA
258 supertypes can be used. We recommend the peptides with the lengths of 8-12 residues,
259 but not less than 8. N-terminal, position 2, C-terminal were treated as anchored sites by
260 default. A predictive score value greater than 0.5 is considered as immunogenicity
261 (Positive-High), the score between 0.4-0.5 is considered as (Positive-Low), the score
262 less than 0.4 is considered as (Negative-High). It is critical to make sure that HLA-
263 subtype must match your peptides ($\text{rank}(\%) < 2$). Where HLA-subtypes mismatch, the
264 large deviation of $\text{rank}(\%)$ value may strongly influence the results. Additionally, the
265 neoantigen model requires providing wild type and mutated sequences at the same time
266 to extract mutation associated characteristics, and currently only immunogenicity
267 prediction for neoantigens of single amino acid mutations are supported. Users can
268 choose example options to test the INeo-Epp ([http://www.biostatistics.online/INeo-](http://www.biostatistics.online/INeo-Epp/neoantigen.php)
269 [Epp/neoantigen.php](http://www.biostatistics.online/INeo-Epp/neoantigen.php)).

270 Discussion

271 Due to the complexity of antigen presenting and TCR binding, the mechanism of TCR
272 recognition has not been clearly revealed. In 2013, J. A. Calis [63] developed a tool for
273 epitope identification for mice and humans (AUC = 0.68). Although mice and human
274 beings are highly homologous, the murine epitopes may very likely cause limitations
275 in identifying human epitopes. Inspired by J. A. Calis, our research here focused on
276 human beings' epitopes and has been conducted in a larger data set.

277 By analyzing epitope immunogenicity from the perspective of amino acid
278 molecular composition, we observed that TCRs do have a preference for hydrophobic
279 amino acid recognition. For short peptides presented by different HLA supertypes,
280 TCRs may have different identification patterns. The immunogenicity prediction based
281 on all HLA-presenting peptides may affect the accuracy of the prediction results. That
282 is, if the prediction could focus on specified HLA-presenting peptides the results may
283 improve. Therefore in our work we used HLA supertypes to improve the prediction of
284 HLA-presenting epitopes, including antigen epitopes and neoantigen epitopes, for a
285 better recognition by TCRs. At present, neoantigen epitopes that can be collected in
286 accordance with the standard for experimental verification are too few, the data of
287 positive and negative neoantigens are unbalanced, and there is not enough data to be
288 used for external verification set. In the future, we will continue to refine and expand
289 our training and verification datasets. Recently, Céline M. Laumont [64] demonstrated
290 that noncoding regions aberrantly expressing tumor-specific antigens (aeTSAs) may
291 represent ideal targets for cancer immunotherapy. These epitopes can also be studied in
292 the future. Increased epitope data may also help empower the prediction of potentially
293 immunogenic peptides or neopeptides.

294 **Conclusions**

295 Neoantigen prediction is the most important step at the start of preparation of
296 neoantigen vaccine. Bioinformatics methods can be used to extract tumor mutant
297 peptides and predict neoantigens. Most current strategies aimed at ended in presenting
298 peptides predictions and among the results of these predictions, probably only fewer
299 than 10 neoantigens might be clinically immunogenic and produce effective immune
300 response. It is time-consuming and costly to experimentally eliminate the false
301 positively predicted peptides. Our methods as developed in this study and the INeo-Epp
302 tool may help eliminate false positive antigen/neoantigen peptides, and greatly reduce
303 the amount of candidates to be verified by experiments. We believe that in the age of
304 biological systems data explosion, computational approaches are a good way to
305 enhance research efficiency and direct biological experiments. With the development
306 of machine learning and deep learning, we expect the prediction of epitope
307 immunogenicity will be continually improved.

308 In summary, this study provides a novel T-cell HLA class-I immunogenicity
309 prediction method from epitopes to neoantigens, and the INeo-Epp can be applied not
310 only to identify putative antigens, but also to identify putative neoantigens.

311 It needs to be stated here that we published the preprint [65] of this article in July
312 2019. This is a modified version.

313 **Data Availability**

314 The data used to support the findings of this study are included within the
315 supplementary information file(s).

316 **Competing of Interests**

317 The author(s) declare(s) that there is no conflict of interest regarding the publication of
318 this paper

319 **Funding Statement**

320 This work was funded by the National Natural Science Foundation of China (No.
321 31870829), Shanghai Municipal Health Commission, and Collaborative Innovation
322 Cluster Project (No. 2019CXJQ02). The funders had no role in study design, data
323 collection and analysis, decision to publish, or preparation of the manuscript.

324 **Acknowledgments**

325 We sincerely thank Drs. Menghuan Zhang, Hong Li and Qibing Leng for valuable
326 discussion. We also acknowledge Dr. Michael Liebman for his critical reading and
327 editing.

328 **Supplementary Material**

329 S1 Table **IEDB antigen epitopes summary**. Detailed description of 17 HLA molecules
330 collected from IEDB. (XLSX)
331 S2 Table **External validation antigen epitopes summary**. Epitope details of 7
332 publications. (XLSX)
333 S3 Table **Neoantigen epitopes summary**. Epitope details of 13 publications. (XLSX)
334 S4 Table **Summary of amino acid characteristics**. For all amino acid characteristics
335 (n=21) that are described in the ExpASY. (XLSX)

336 **References**

- 337 [1] D. V. Desai, and U. Kulkarni-Kale, “T-cell epitope prediction methods: an
338 overview,” *Methods Mol Biol*, vol. 1184, pp. 333-64, 2014.
- 339 [2] A. L. Goldberg, and K. L. Rock, “Proteolysis, proteasomes and antigen
340 presentation,” *Nature*, vol. 357, no. 6377, pp. 375-379, 1992.
- 341 [3] K. Can, A. K. Nussbaum, S. Hansjörg *et al.*, “Prediction of proteasome cleavage
342 motifs by neural networks,” *Protein Eng*, no. 4, pp. 4, 2002.
- 343 [4] M. V. Larsen, C. Lundegaard, K. Lamberth *et al.*, “An integrative approach to
344 CTL epitope prediction: A combined algorithm integrating MHC class I binding,
345 TAP transport efficiency, and proteasomal cleavage predictions,” *European
346 Journal of Immunology*, vol. 35, no. 8, pp. 2295-2303, 2005.
- 347 [5] V. Jurtz, S. Paul, M. Andreatta *et al.*, “NetMHCpan-4.0: Improved Peptide–
348 MHC Class I Interaction Predictions Integrating Eluted Ligand and Peptide
349 Binding Affinity Data,” *Journal of Immunology*, vol. 199, no. 9, pp. j11700893,
350 2017.
- 351 [6] T. J. O'Donnell, A. Rubinsteyn, M. Bonsack *et al.*, “MHCflurry: Open-Source
352 Class I MHC Binding Affinity Prediction,” *Cell Syst*, vol. 7, no. 1, pp. 129-
353 132.e4, Jul 25, 2018.

- 354 [7] M. Wang, K. Lamberth, M. Harndahl *et al.*, “CTL epitopes for influenza A
355 including the H5N1 bird flu; genome-, pathogen-, and HLA-wide screening,”
356 *Vaccine*, vol. 25, no. 15, pp. 0-2831,2007.
- 357 [8] C. L. Perez, M. V. Larsen, R. Gustafsson *et al.*, “Broadly Immunogenic HLA
358 Class I Supertype-Restricted Elite CTL Epitopes Recognized in a Diverse
359 Population Infected with Different HIV-1 Subtypes,” *Journal of Immunology*,
360 vol. 180, no. 7, pp. 5092-5100,2008.
- 361 [9] C. Lundegaard, I. Hoof, O. Lund *et al.*, “State of the art and challenges in
362 sequence based T-cell epitope prediction,” *Immunome Research*, vol. 6 Suppl 2,
363 no. Suppl 2, pp. S3, 2010.
- 364 [10] J. L. Sanchez-Trincado, G.-P. Marta, and R. P. A., “Fundamentals and Methods
365 for T- and B-Cell Epitope Prediction,” *Journal of Immunology Research*, vol.
366 pp. 1-14,2017.
- 367 [11] E. G. Phimister, and V. N. Kristensen, “The Antigenicity of the Tumor Cell —
368 Context Matters,” *New England Journal of Medicine*, vol. 376, no. 5, pp. 491-
369 493,2017.
- 370 [12] K. Kiyotani, H. T. Chan, and Y. Nakamura, “Immunopharmacogenomics
371 towards personalized cancer immunotherapy targeting neoantigens,” *Cancer Sci*,
372 vol. 109, no. 3, pp. 542-549, Mar, 2018.
- 373 [13] V. Randi, J. A. Overton, J. A. Greenbaum *et al.*, “The immune epitope database
374 (IEDB) 3.0,” *Nucleic Acids Research*, no. D1, pp. D1, 2014.
- 375 [14] A. Sette, and J. Sidney, “Nine major HLA class I supertypes account for the vast
376 preponderance of HLA-A and -B polymorphism,” *Immunogenetics*, vol. 50, no.
377 3-4, pp. 201-12, Nov, 1999.
- 378 [15] J. Sidney, B. Peters, N. Frahm *et al.*, “HLA class I supertypes: a revised and
379 updated classification,” vol. 9, no. 1, pp. 1-0,2008.
- 380 [16] “An immunogenic personal neoantigen vaccine for patients with melanoma.”
- 381 [17] “Personalized RNA mutanome vaccines mobilize poly-specific therapeutic
382 immunity against cancer,” *Nature*, vol. 547, no. 7662, pp. 222-226,2017.
- 383 [18] Z. Hu, P. A. Ott, and C. J. Wu, “Towards personalized, tumour-specific,
384 therapeutic vaccines for cancer,” *Nat Rev Immunol*, vol. 18, no. 3, pp. 168-182,
385 Mar, 2018.
- 386 [19] E. M. Van Allen, D. Miao, B. Schilling *et al.*, “Genomic correlates of response
387 to CTLA-4 blockade in metastatic melanoma,” *Science*, vol. 350, no. 6257, pp.
388 207-211, 2015.
- 389 [20] M. Efremova, F. Finotello, D. Rieder *et al.*, “Neoantigens Generated by
390 Individual Mutations and Their Role in Cancer Immunity and Immunotherapy,”
391 *Front Immunol*, vol. 8, pp. 1679, 2017.
- 392 [21] L. Klein, M. Hinterberger, G. Wirnsberger *et al.*, “Antigen presentation in the
393 thymus for positive selection and central tolerance induction,” *Nature reviews*.
394 *Immunology*, vol. 9, no. 12, pp. 833-844,2009.
- 395 [22] F. F. Gonzalez-Galarza, A. McCabe, E. J. Melo Dos Santos *et al.*, “Allele

- 396 Frequency Net Database,” *Methods Mol Biol*, vol. 1802, pp. 49-62, 2018.
- 397 [23] D. Weiskopf, M. A. Angelo, E. L. D. Azeredo *et al.*, “Comprehensive analysis
398 of dengue virus-specific responses supports an HLA-linked protective role for
399 CD8(+) T cells,” *Proc Natl Acad Sci U S A*, vol. 110, no. 22, pp. E2046-E2053,
400 2013.
- 401 [24] H. Luxenburger, F. Grass, J. Baermann *et al.*, “Differential virus-specific CD8(+)
402 T-cell epitope repertoire in hepatitis C virus genotype 1 versus 4,” *J Viral Hepat*,
403 vol. 25, no. 7, pp. 779-790, Jul, 2018.
- 404 [25] Y. Xia, W. Pan, X. Ke *et al.*, “Differential escape of HCV from CD8+ T cell
405 selection pressure between China and Germany depends on the presenting HLA
406 class I molecule,” *Journal of Viral Hepatitis*, vol. 26, no. 1, pp. 73-82, 2019.
- 407 [26] H. Vahed, A. Agrawal, R. Srivastava *et al.*, “Unique Type I Interferon,
408 Expansion/Survival Cytokines, and JAK/STAT Gene Signatures of
409 Multifunctional Herpes Simplex Virus-Specific Effector Memory CD8 T Cells
410 Are Associated with Asymptomatic Herpes in Humans,” *Journal of Virology*,
411 vol. 93, no. 4, pp. e01882-18, 2019.
- 412 [27] A. Khakpoor, Y. Ni, A. Chen *et al.*, “Spatiotemporal Differences in Presentation
413 of CD8 T Cell Epitopes during Hepatitis B Virus Infection,” *J Virol*, vol. 93, no.
414 4, Feb 15, 2019.
- 415 [28] A. Huth, X. Liang, S. Krebs *et al.*, “Antigen-Specific TCR Signatures of
416 Cytomegalovirus Infection,” *J Immunol*, vol. 202, no. 3, pp. 979-990, Feb 1,
417 2019.
- 418 [29] S. O. Sekyere, B. Schlevogt, F. Mettke *et al.*, “HCC immune surveillance and
419 antiviral therapy of hepatitis C virus infection,” *Liver cancer*, vol. 8, no. 1, pp.
420 41-65, 2019.
- 421 [30] D. A. Wick, J. R. Webb, J. S. Nielsen *et al.*, “Surveillance of the Tumor
422 Mutanome by T Cells during Progression from Primary to Recurrent Ovarian
423 Cancer,” *Clinical Cancer Research*, vol. 20, no. 5, 2013.
- 424 [31] T. Karasaki, K. Nagayama, M. Kawashima *et al.*, “Identification of Individual
425 Cancer-Specific Somatic Mutations for Neoantigen-Based Immunotherapy of
426 Lung Cancer,” *Journal of Thoracic Oncology Official Publication of the
427 International Association for the Study of Lung Cancer*, vol. 11, no. 3, pp. 324-
428 333, 2015.
- 429 [32] A. Gros, M. R. Parkhurst, E. Tran *et al.*, “Prospective identification of
430 neoantigen-specific lymphocytes in the peripheral blood of melanoma patients,”
431 *Nature Medicine*, vol. 22, no. 4, pp. 433-438, 2016.
- 432 [33] E. Strønen, M. Toebe, S. Kelderman *et al.*, “Targeting of cancer neoantigens
433 with donor-derived T cell receptor repertoires,” *Science*, vol. 352, no. 6291, pp.
434 1337-1341, 2016.
- 435 [34] A. Nelde, J. S. Walz, D. J. Kowalewski *et al.*, “HLA class I-restricted MYD88
436 L265P-derived peptides as specific targets for lymphoma immunotherapy,”
437 *OncImmunology*, vol. 6, no. 3, Mar 4, 2017.

- 438 [35] X. Zhang, S. Kim, J. Hundal *et al.*, “Breast Cancer Neoantigens Can Induce
439 CD8 T-Cell Responses and Antitumor Immunity,” *Cancer Immunology*
440 *Research*, vol. 5, no. 7, pp. 516-523, 2017.
- 441 [36] M. Markus, G. David, C. George *et al.*, “‘Hotspots’ of Antigen Presentation
442 Revealed by Human Leukocyte Antigen Ligandomics for Neoantigen
443 Prioritization,” *Front Immunol*, vol. 8, pp. 1367,2017
- 444 [37] V. P. Balachandran, M. Łuksza, J. N. Zhao *et al.*, “Identification of unique
445 neoantigen qualities in long-term survivors of pancreatic cancer,” *Nature*, vol.
446 551, no. 7681, pp. 512-516,2017.
- 447 [38] T. Matsuda, M. Leisegang, J.-H. Park *et al.*, “Induction of Neoantigen-Specific
448 Cytotoxic T Cells and Construction of T-cell Receptor-Engineered T Cells for
449 Ovarian Cancer,” *Clinical cancer research : an official journal of the American*
450 *Association for Cancer Research*, vol. 24, no. 21, pp. 5357-5367, 2018.
- 451 [39] K. Sonntag, H. Hashimoto, M. Eyrih *et al.*, "Immune monitoring and TCR
452 sequencing of CD4 T cells in a long term responsive patient with metastasized
453 pancreatic ductal carcinoma treated with individualized, neopeptide-derived
454 multi-peptide vaccines: a case report," *Journal of translational medicine*,
455 16,2018.
- 456 [40] A.-M. Bjerregaard, M. Nielsen, V. Jurtz *et al.*, "An Analysis of Natural T Cell
457 Responses to Predicted Tumor Neoepitopes," *Frontiers in immunology*, 8, 2017.
- 458 [41] C. E. Shannon, “A Mathematical Theory of Communication,” *Bell System*
459 *Technical Journal*, vol. 27, 1948.
- 460 [42] M. Kuhn, “Building Predictive Models in R Using the caret Package,” *Journal*
461 *of Statistical Software*, 2008.
- 462 [43] M. B. Kursa, and W. R. Rudnicki, “Feature Selection with the Boruta Package,”
463 *Journal of Statistical Software*, vol. 036, 2010.
- 464 [44] A. Liaw, and M. Wiener, “Classification and Regression by randomForest,” *R*
465 *News*, vol. 23, no. 23, 2002.
- 466 [45] T. Sing, O. Sander, N. Beerenwinkel *et al.*, “ROCR: visualizing classifier
467 performance in R,” *Bioinformatics (Oxford, England)*, vol. 21, no. 20, pp. 3940-
468 3941, 2005.
- 469 [46] Walker, and M. J., “The proteomics protocols handbook,” *Biochemistry*, vol. 71,
470 no. 6, pp. 696-696, 2006.
- 471 [47] J. Kyte, and R. F. Doolittle, “A simple method for displaying the hydrophobic
472 character of a protein,” vol. 157, no. 1, pp. 105-132,1982.
- 473 [48] J. M. Zimmerman, N. Eliezer, and R. Simha, “The characterization of amino
474 acid sequences in proteins by statistical methods,” *Journal of theoretical biology*,
475 vol. 21, no. 2, pp. 170-201,1968.
- 476 [49] Grantham, and R., “Amino Acid Difference Formula to Help Explain Protein
477 Evolution,” *Science*, vol. 185, no. 4154, pp. 862-864,1974.
- 478 [50] Fraga, and Serafin, “Theoretical prediction of protein antigenic determinants
479 from amino acid sequences,” *Canadian Journal of Chemistry*, vol. 60, no. 20,

- 480 pp. 2606-2610,1982.
- 481 [51] R. M. Sweet, and D. Eisenberg, "Correlation of sequence hydrophobicities
482 measures similarity in three-dimensional protein structure," *Journal of*
483 *molecular biology*, vol. 171, no. 4, pp. 479-488,1983.
- 484 [52] Meek, and L. J., "Prediction of peptide retention times in high-pressure liquid
485 chromatography on the basis of amino acid composition," *Proceedings of the*
486 *National Academy of Sciences of the United States of America*, vol. 77, no. 3,
487 pp. 1632-1636,1980.
- 488 [53] G. D. Rose, A. R. Geselowitz, G. J. Lesser *et al.*, "Hydrophobicity of amino acid
489 residues in globular proteins," *Science (New York, N.Y.)*, vol. 229, no. 4716, pp.
490 834-838, 1985.
- 491 [54] P. Y. Chou, and G. D. Fasman, "Prediction of the secondary structure of proteins
492 from their amino acid sequence," *Advances in enzymology and related areas of*
493 *molecular biology*, vol. 47, pp. 45-148, 1978, 1978.
- 494 [55] G. Deléage, and B. Roux, "An algorithm for protein secondary structure
495 prediction based on class prediction," *Protein engineering*, vol. 1, no. 4, pp.
496 289-294, 1987 Aug-Sep, 1987.
- 497 [56] A. Burger, "Atlas of Protein Sequence and Structure 1969," *Journal of*
498 *Medicinal Chemistry*, vol. 13, no. 2, pp. 337-337, 1970.
- 499 [57] D. D. Jones, "Amino acid properties and side-chain orientation in proteins: A
500 cross correlation approach," *Journal of Theoretical Biology*, vol. 50, no. 1, pp.
501 167-183,1975.
- 502 [58] G. Zhao, and E. London, "Strong Correlation Between Statistical
503 Transmembrane Tendency and Experimental Hydrophobicity Scales for
504 Identification of Transmembrane Helices," *Journal of Membrane Biology*, vol.
505 229, no. 3, pp. p.165-168,2009.
- 506 [59] J. Janin, "Surface and inside volumes in globular proteins," *Nature*, vol. 277,
507 no. 5696, pp. 491-492, 1979.
- 508 [60] J. R. Green, M. J. Korenberg, R. David *et al.*, "Recognition of Adenosine
509 Triphosphate Binding Sites Using Parallel Cascade System Identification,"
510 *Annals of Biomedical Engineering*, vol. 31, no. 4, pp. 462-470,2003.
- 511 [61] S. Lifson, and C. Sander, "Antiparallel and parallel beta-strands differ in amino
512 acid residue preferences," *Nature*, vol. 282, no. 5734, pp. 109-111, 1979.
- 513 [62] F. Duan, J. Duitama, S. Al Seesi *et al.*, "Genomic and bioinformatic profiling of
514 mutational neoepitopes reveals new rules to predict anticancer immunogenicity,"
515 *J Exp Med*, vol. 211, no. 11, pp. 2231-48, Oct 20, 2014.
- 516 [63] J. J. A. Calis, M. Maybeno, J. A. Greenbaum *et al.*, "Properties of MHC class I
517 presented peptides that enhance immunogenicity," *PLoS computational biology*,
518 vol. 9, no. 10, pp. e1003266,, 2013.
- 519 [64] C. M. Laumont, K. Vincent, L. Hesnard *et al.*, "Noncoding regions are the main
520 source of targetable tumor-specific antigens," *Sci Transl Med*, vol. 10, no. 470,
521 Dec 5, 2018.

- 522 [65] G. Wang, H. Wan, X. Jian *et al.*, "INeo-Epp: T-cell HLA class I immunogenic
523 or neoantigenic epitope prediction via random forest algorithm based on
524 sequence related amino acid features," bioRxiv, 2019.
525