# PSSMHCpan: a novel PSSM based software for predicting class I peptide-HLA binding affinity

Geng Liu[†1], Dongli Li[†1], Si Qiu[1, 2], Wenhui Li[1], Kun Ma[1], Jian Wang[1, 3], Huanming Yang[1, 3], Yong Hou[*1, 4], Bo Li[*1, 5]

[1] BGI-Shenzhen, Shenzhen 518083 China

[2] BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083, China

[3] James D. Watson Institute of Genome Sciences, Hangzhou 310058, China

[4] Department of Biology, University of Copenhagen, Denmark

[5] BGI-Forensics, Shenzhen 518083, China

†These authors contributed equally to this work as first authors.

*To whom correspondence should be addressed. Tel: (86) 18680679919; Fax: (86) 0755-25273620;

E-mail: libo@genomics.cn (B.L.). Correspondence may also be addressed. E-mail:

houyong@genomics.cn (Y.H.).

**Abstract**

**Background:** Predicting peptides binding affinity with human leukocyte antigen (HLA) is a crucial step in developing powerful antitumor vaccine for cancer immunotherapy. Currently available methods work reasonably well in predicting peptide binding affinity with HLA-A*0201, HLA-A*0101, and HLA-B*0702 in terms of sensitivity and specificity. However, it is unknown whether these methods can also predict well with other HLA alleles that are present in majority of human populations.

**Result:** Here we present a Position Score Specific Matrix (PSSM) based software called PSSMHCpan to accurately and efficiently predict peptide binding affinity with a broad coverage of HLA class I alleles. By analyzing 10 cross-validations on training database of 87 HLA alleles and an independent

1

25  dataset with NetMHC-4.0, NetMHCpan-3.0, PickPocket, and PSSMHCpan, we found that

26  PSSMHCpan is substantially better than the other three methods with accuracy ACC of 0.92 and

27  sensitivity of 0.87, as compared to 0.85, 0.85, 0.72 in 10 cross-validations and 0.73, 0.79, 0.75 in the

28  independent dataset evaluation. In addition, PSSMHCpan is more than 763 times faster than other three

29  methods to predict neoantigens from a breast tumor sample. Finally we built a neoantigen prediction

30  pipeline and identified 117,017 neoantigens from 467 cancer samples of diverse cancers from TCGA.

31  **Conclusion:** PSSMHCpan is superior to currently available methods in predicting peptide binding

32  affinity with a broad coverage of HLA class I alleles.

33  **Key words**: Antitumor vaccine, peptide-HLA binding affinity, PSSMHCpan, neoantigen.

34

35  **Background**

36  Cancer immunotherapy has been proved to be a promising strategy that enhances the strengths of the

37  immune system of cancer patients to fight cancer in recent years. This strategy exploits the fact that

38  surface of cancer cells have a variety of tumor antigens (i.e. peptides of 8-13 residues in lengths)

39  coming from various kinds of mutated proteins cleaved by the proteasomes intracellular. These

40  peptides are bound to HLA class I allelic specific molecules, forming peptide-HLA complexes which

41  are presented to T cell receptors (TCRs). If TCRs can recognize the complexes on the surface of cancer

42  cells, cytotoxic T lymphocytes (CTLs) will destroy cancer cells. Cancer cells are highly heterogeneous

43  in terms of morphological, phonotypical and genetic profiles. Cancer cells of different tumors and

44  within the same tumor could present hundreds of different types of peptides. The immune system of

45  cancer patients could only recognize small populations of cancer cells. In order to enhance the power of

46  the CTLs to recognize and eradicate as many cancer cells as possible, one strategy is to vaccinate

47  cancer patients with complex antitumor peptides. The first step to develop powerful antitumor vaccines

48  is to predict peptide binding affinity with HLA class I allele.

2

In order to predict peptide binding affinity with HLA class I allele, four types of methods have been developed, including structure based methods, machine learning based methods, PSSM based methods [1] and combined methods. The structure based methods predict peptide binding affinity by calculating the minimum free energy of peptide-HLA complex [2], which allows us to understand the peptide-HLA binding affinity at the structure level. However, the predicting speed of this types of methods is extremely slow, and inaccurate due to limited number of available crystal structures [3]. The machine learning based methods predict peptide binding affinity by learning a function that maps a given peptide to binding affinity based on available known bound peptides (binders). These methods can accurately predict peptides with specific HLA alleles of HLA-A*0201, HLA-A*0101, and HLA-B*0702 [4, 5]. Hence, they are widely used in many studies [6-8]. Thus far, many methods of machine learning have been developed, including support vector machine based method MHC2PRED [9], hidden markov model based method S-HMM [10], artificial neural network based method NetMHC [11, 12], and pan-specific method NetMHCpan [13-15]. However, machine learning methods cannot accurately predict peptide binding affinity with a broad range of HLA class I allelic coverage. Further, they are inefficient in predicting peptides from a large amount of sequencing data. The PSSM based methods predict peptide binding affinity by building a matrix from multiple peptides alignment results that represent the motif information (i.e. the binding anchor). These methods have a faster predicting speed because linear computational complexity of PSSM is much lower than nonlinear computational complexity of structure and machine learning based methods. Based on the mechanism of PSSM, several software have been developed such as PickPocket [16], SVMHC [17] and nHLAPred [18]. However the accuracy of current software is less than machine learning based methods [16].

Recently, in order to predict peptide-HLA binding affinity more accurately, scientists from several

3

71    groups combined different methods to develop new software including NetMHCcons [19], IEDB [20]

72    and HLaffy [21]. Although these combined methods indeed have shown a better predictive

73    performance as compared to individual methods, their predictive accuracy are still not satisfactory,

74    especially in clinical applications [22]. In order to develop more effective immunotherapy, it is

75    necessary to develop better software that can more accurately and efficiently predict peptide binding

76    affinity with a broad coverage of HLA class I alleles.

77      Here, we present a novel software called PSSMHCpan. We designed this software based on the

78    PSSM mechanism and using a more comprehensive training database containing 63,099 peptide-HLA

79    pairs to allele-specifically predict peptide binding affinity with HLA class I allele. In order to predict

80    peptide binding affinity with a broad coverage of HLA class I alleles, we induce a simple but powerful

81    pan-specific prediction approach based on the similarity of HLA protein sequences. We show that

82    PSSMHCpan can predict peptide binding affinity with a broad HLA class I allelic coverage of at least

83    87 types more accurately and efficiently than other available methods in 10 cross-validations and

84    independent dataset evaluation. Based on PSSMHCpan, we built a prediction pipeline to identify

85    neoantigens in 467 TCGA tumor samples across 10 types of cancers.

86

**Methods**

88    PSSM is represented as a motif of multiple sequence alignment result [23]. The basic principle of

89    PSSMHCpan is that peptides that bind to a specific HLA allele possess the motif information that can

90    be studied by PSSM. We propose the PSSMHCpan in two novel aspects. Firstly, we construct a

91    comprehensive training database to build allele-specific PSSMs for predicting peptide binding affinity

92    with characterized HLA class I allele (with binders in training database). Secondly, we use the

4

93    similarity of HLA sequences to induce a simple but powerful pan-specific prediction approach based

94    on our hypothesis below to predict peptide binding affinity with uncharacterized HLA class I allele

95    (without binder in training database). It is well known that peptides on the cell surface are bound to the

96    floor of the peptide-binding groove that is in the central region of the α1/α2 heterodimer (a molecule

97    composed of two non-identical subunits) of HLA protein sequences [24]. By analyzing the sequences

98    of HLA proteins, we noticed that HLA protein sequences are highly similar among different HLA

99    alleles (Figure 1), and that peptides bound to similar HLA alleles have similar binding affinity

100   according to predictive value of IC50. Thereby, we hypothesize that since different HLA protein

101   sequences are similar, the peptide binding affinity with different HLA alleles should be similar too.

102   Based on this hypothesis and the PSSM mechanism, we design the software PSSMHCpan as following

103   three steps: PSSM construction, allele-specific prediction, and pan-specific prediction. The flowchart of

104   PSSMHCpan is shown in Figure 2.

105

106   **PSSM construction**

107   We define PSSM as a matrix of M rows (Amino acid; M=20) and N columns (Length; N=8~25). Each

108   element $P_{ai}$ in the matrix is the likelihood of a given character (amino acid) at its position. We

109   calculate the element $P_{ai}$ through the following function,

$$P_{ai} = log\frac{F_{ai} + \omega}{BG_a}$$

111      Where $F_{ai}$ denotes the frequency of amino acid $a$ at position $i$; $BG_a$ denotes the background

112   frequency of amino acid $a$ from UniProt database [25]; and $\omega$ is a random value (ranging from 0 to 1)

113   generated from Dirichlet distribution [26].

114

5

**Allele-specific prediction**

To qualitatively predict peptide binding affinity with characterized HLA allele, we define a *binder_score* as a sum of the corresponding values of each amino acid of a given peptide at each position in the corresponding allele-specific PSSM.

$$binder\_score = \frac{\sum_{i=1}^{N} P_{ai}}{N}$$

We consider a peptide with *binder_score > 0* as a binder. The higher *binder_score* that a peptide has, the higher binding affinity this peptide would have.

We convert a binding affinity score (*binder_score*) into an IC50 value as follows:

$$IC50 = 50000^{Max-binder\_score}/_{Max-Min}$$

Where Max and Min denote the maximum and the minimum *binder_score*, respectively. We consider a peptide with *IC50 < 500nM* as a binder and a peptide with *IC50 < 50nM* as a strong binder.

**Pan-specific prediction**

Firstly, we construct a library of HLA similar weight (Button panel in Figure 2) that contains pairs of characterized and uncharacterized HLA alleles, and each pair has a weight value. We determine a pair of characterized and uncharacterized HLA alleles by using the BLOSUM62 [27] based BLAST alignment results of HLA protein sequences, and assign the alignment score as the weight value. We also extracted the nearest distance of HLA alleles from NetMHCpan-3.0 [15] as a pair of characterized and uncharacterized HLA alleles and assigned a constant as the weight value.

Secondly, we qualitatively predict the binding affinity of a given peptide with uncharacterized HLA allele with an *IC50_{un}* value which is calculated as below:

$$IC50_{un} = \frac{\sum_{i=1}^{S} w_i * IC50_i}{\sum_{i=1}^{S} w_i}$$

6

137 Where $S$ denotes the sum of characterized HLA alleles that pair up the specific uncharacterized

138 HLA allele according to the library of HLA similar weight. $w_i$ and $IC50_i$ denote the weight value

139 and the allele-specific prediction result of peptide binding affinity with HLA allele $i$. We also consider a

140 peptide with $IC50_{un} < 500nM$ as a binder, and a peptide with $IC50_{un} < 50nM$ as a strong binder.

141

**Data Description**

143 We collected our training database of HLA class I binders from the following resources: the Immune

144 Epitope Database and Analysis Resource (IEDB) [28], IEDB benchmark [29], SYFPEITHI [30],

145 MHCBN [31], and in-house experimental epitopes. After removing duplications, we obtained 64,677

146 peptide-HLA pairs that cover 162 HLA alleles (Table 1). We only selected HLA alleles that consist of

147 at least 10 binders with a fixed length. Finally, we built 241 PSSMs for allele-specific prediction of

148 peptide binding affinity with 123 HLA class I alleles (Additional file 1: Table S1).

149 **Table 1** Summary of training database.

| Database | IEDB | IEDB benchmark | SYFPEITHI | MHCBN | Combined | Training database |
|---|---|---|---|---|---|---|
| **HLA alleles** | 166 | 95 | 109 | 103 | 162 | 123 |
| **Binders** | 54,272 | 40,930 | 3,329 | 4,070 | 64,677 | 63,099 |

150 We collected 64 uncharacterized HLA class I alleles that cannot be predicted with NetMHC-4.0 but

151 can be predicted with NetMHCpan-3.0. We extracted 2064 binders that bind to the 64 uncharacterized

152 HLA alleles from our training database as a dataset for pan-specific evaluation.

153 To construct a library of HLA weight similarity, we collected 657,397 pairs of characterized and

154 uncharacterized HLA class I alleles from 13,957 HLA protein sequences in IMGT/HLA (Release

155 3.23.0) [32], and 2800 pairs from the nearest distance of HLA alleles in NetMHCpan-3.0, respectively.

7

156  After removing duplications, we retained 657,930 pairs for pan-specific prediction of peptide binding

157  affinity with 4,778 HLA class I alleles (Additional file 1: Table S1).

158   We also collected an independent dataset of binders from the Peptide Database of Cancer Immunity

159  [33]. Then we selected 285 binders that cover 38 HLA alleles of HLA-A, HLA-B, HLA-C, including

160  35 from tumor antigens resulting from mutations, 91 from shared tumor-specific antigens, 63 from

161  differentiation antigens and 96 from antigens overexpressed in tumors. After removing duplications, we

162  retained 273 binders for validation.

163   To detect pan-cancer neoantigens, we obtained somatic mutations of 467 TCGA cancer samples

164  across 10 cancer types (Table 2) from GDC data portal (https://gdc-portal.nci.nih.gov/) and the RSEM

165  gene expression data of these tumors and their corresponding normal samples from FireBrowse

166  (http://firebrowse.org/). We also obtained the tumor RNASeq aligned bam files from dbGAP.

167

168  **Table 2 Summary of 467 cancer samples from TCGA cohort.**

| Cancer type | Patient # | Cancer type | Patient # |
|---|---|---|---|
| BLCA | 19 | LIHC | 47 |
| BRCA | 93 | LUAD | 57 |
| COAD | 16 | PRAD | 43 |
| HNSC | 39 | STAD | 28 |
| KIRC | 67 | THCA | 58 |

169

170  **Analyses**

171  **Evaluation of peptide binding affinity prediction with a broad HLA class I allelic coverage**

8

172   In order to evaluate the allele-specific prediction accuracy of PSSMHCpan with a broad HLA class I

173   allelic coverage, we performed 10 cross-validations on training data of 87 HLA class I alleles that

174   contain at least 12 binders. We generated non-binders randomly with the same number of binders, and

175   performed allele-specific prediction of peptide-HLA binding affinity using our PSSMHCpan, and the

176   two well-known and currently considered as the best software for peptide-HLA binding affinity

177   prediction NetMHC-4.0 and NetMHCpan-3.0 [4], and with the latest reported PSSM based software

178   PickPocket, respectively. We found that the performance of the four software appeared similar in terms

179   of the average area under receiver operating characteristic curve (AUC) with the HLA alleles of

180   HLA-A*0101, HLA-A*0201, and HLA-B*0702 (Additional file 1: Table S2). However, in terms of the

181   prediction accuracy ACC (ACC $= \frac{TP+TN}{TP+FP+TN+FN}$, where TP, FP, TN and FN, represent true-positive,

182   false-positive, true-negative and false-negative) under the cutoff at 500nM, PSSMHCpan is larger than

183   NetMHC-4.0, NetMHCpan-3.0 and PickPocket (Table 3), suggesting that the PSSMHCpan delivers

184   more accurate than the other three software in predicting peptide binding affinity with the HLA alleles

185   of HLA-A*0101, HLA-A*0201, and HLA-B*0702 at 500nM. We also noticed that although the overall

186   AUC of PSSMHCpan is slightly larger than that of any of the software with the rest HLA class I alleles

187   (ranging from 1% to 2%; Figure 3a), the ACC of PSSMHCpan is much larger than those of other three

188   software (ranging from 7% to 20%). By comparing the ACC of each HLA allele with a fixed peptide

189   length among the four software, we found that the median ACC of PSSMHCpan is significantly larger

190   than other three software ($P <0.01$, *Paired T test*; Figure 3b).

191

192   **Table 3** Assessments (ACC values) of four software to predict peptide binding affinity with three HLA

193   alleles.

| | A*0101 9mer | A*0201 9mer | B*0702 9mer | A*0101 10mer | A*0201 10mer | B*0702 10mer |
|---|---|---|---|---|---|---|
| **PSSMHCpan** | **0.96** | **0.88** | **0.91** | **0.96** | **0.92** | **0.96** |
| **NetMHC-4.0** | 0.86 | 0.86 | 0.87 | 0.86 | 0.88 | 0.90 |
| **NetMHCpan-3.0** | 0.85 | 0.86 | 0.87 | 0.83 | 0.87 | 0.88 |
| **PickPocket** | 0.65 | **0.88** | 0.85 | 0.53 | 0.89 | 0.81 |

194

195     Considering a one-time 10 cross-validation of randomly selection and non-binders construction

196     might produce biased results, we repeated another five times of 10 cross-validations, and found that

197     (Table 4) the standard deviations (SD) of AUCs are $\leq$ 0.0005, indicating no bias in the 10

198     cross-validation.

199 **Table 4** The AUC and SD values in 5 times 10 cross-validations.

| Time | PSSMHCpan | NetMHC-4.0 | NetMHCpan-3.0 | PickPocket |
|------|-----------|------------|---------------|------------|
| 1 | 0.9693 | 0.9623 | 0.965 | 0.9494 |
| 2 | 0.9703 | 0.9633 | 0.9661 | 0.9507 |
| 3 | 0.9703 | 0.9633 | 0.9661 | 0.9506 |
| 4 | 0.9699 | 0.9632 | 0.966 | 0.9505 |
| 5 | 0.9699 | 0.9633 | 0.9657 | 0.9506 |
| SD | 0.0004 | 0.0004 | 0.0005 | 0.0005 |

200

201     To evaluate our pan-specific prediction, we retrained PSSMs without binders from the dataset for

202     pan-specific evaluation. And then we predicted binders from the dataset for pan-specific evaluation and

203     2,064 randomly constructed non-binders by PSSMHCpan. Although the AUC of PSSMHCpan (0.93) is

204     slightly lower than those of NetMHCpan-3.0 and PickPocket (0.96; Figure 3c; Additional file 1: Table

205     S3), the ACC of PSSMHCpan (0.86) is much larger than those two software (0.75 and 0.73). By

206     comparing the allele-specific prediction and pan-specific prediction of 3,408 correctly predicted

207     peptides from the dataset for pan-specific evaluation, we found a high correlation between

208     allele-specific and pan-specific prediction (Pearson' rho=0.89, $P$<0.01; Figure 3d), suggesting that our

209     PSSMHCpan can quantitatively predict peptide-HLA binding affinity with profound accuracy.

210     We compared the performance of our PSSMHCpan with the latest software HLaffy developed by

211     Mukherjee et al (2016) using the same peptides from the MHCBN. We removed all the binders from

212     MHCBN in our training database and retrained our PSSMs with the rest of binders. Because the

213     number of non-binders is much smaller than that of the binders in MHCBN, we only used binders to

214     evaluate and calculated the prediction accuracy by sensitivity (Sen $= \frac{TP}{TP+FP}$). We found that our

215     PSSMHCpan correctly detected 1309 binders, while HLaffy correctly detected 1179 binders (Table 5).

216 **Table 5** Assessments of PSSMHCpan and HLaffy. The prediction of HLaffy was performed on

217     webserver (http://proline.biochem.iisc.ernet.in/HLaffy/).

| Allele | PSSMHCpan | HLaffy |
|--------|-----------|--------|
| HLA-A*0201 | **100.00%** | 91.99% |
| HLA-A*0203 | **100.00%** | 93.22% |

10

| | | |
|---|---|---|
| HLA-A*0206 | **100.00%** | 93.44% |
| HLA-A*0301 | **100.00%** | 83.93% |
| HLA-A*1101 | **100.00%** | 96.00% |
| HLA-A*2402 | **100.00%** | 76.60% |
| HLA-A*3301 | **100.00%** | 83.33% |
| HLA-A*6801 | **100.00%** | 94.12% |
| HLA-A*6802 | **95.45%** | 72.73% |
| HLA-B*0702 | **100.00%** | 87.88% |
| HLA-B*3501 | **99.16%** | 89.08% |
| HLA-B*5301 | **100.00%** | 91.84% |
| HLA-B*5401 | **100.00%** | 88.10% |
| All | **99.85%** | 89.93% |

218

**Evaluation of peptide binding affinity prediction with an independent dataset**

220 Considering cross validation might overestimate prediction accuracy, we reevaluated PSSMHCpan,

221 NetMHC-4.0, NetMHCpan-3.0 and PickPocket with an independent dataset containing 273

222 non-duplicated experimental binders from the Peptide Database of Cancer Immunity. If a peptide binds

223 to any 4-digital HLA allele that belong to the given 2-digital HLA allele with a predicting binding

224 affinity IC50 less than 500nM, we considered as binder. Totally, 245 of 273 (90%) binders were

225 identified with the four software. Of the 245 binders identified, PSSMHCpan, NetMHC-4.0,

226 NetMHCpan-3.0 and PickPocket identified 237, 199, 216, and 204, respectively (Figure 4; Additional

227 file 1: Table S4), again indicating that PSSMHCpan can predict more binders than either NetMHC-4.0,

228 NetMHCpan-3.0, or PickPocket can.

229

**Evaluation of the software efficiency**

231 As whole genome sequencing (WGS) and whole exome sequencing (WES) of cancer genome data are

232 rapidly increasing, there is an urgent need to develop software that can quickly identify neoantigens

233 from cancer genome data. To compare the efficiency of PSSMHCpan, NetMHC-4.0, NetMHCpan-3.0

11

234    and PickPocket (Table 6), we first calculated the predicting speed of 10-cross validation on training

235    database with 87 HLA class I alleles and found that PSSMHCpan is much faster than other three. We

236    then used each software to independently predict binding affinity of the same set of 661,263 peptides

237    generated from a breast tumor sample containing 3062 somatic mutations and 6 HLA class I alleles. We

238    found that it took about 6 seconds for PSSMHCpan to complete the analysis. In contrast, NetMHC-4.0,

239    took 3.61 hours, NetMHCpan-3.0 took 28.63 hours, and PickPocket took 1.34 hours to complete the

240    analysis. In general, PSSMHCpan are not only more accuracy but also faster than other methods.

241    **Table 6** The predicting speed (CPU time) of the four software. The fastest ones were marked in bold.

| Methods | 10-cross validations | Breast tumour neoantigens prediction |
|---|---|---|
| **PSSMHCpan** | **18.40s** | **6.34s** |
| **NetMHC-4.0** | 1056.83s | 13001.57s |
| **NetMHCpan-3.0** | 5371.16s | 103060.24s |
| **PickPocket** | 282.83s | 4839.63s |

242    CPU time was measured by second (s).

243

244    **Pan-cancer neoantigens**

245    To identify neoantigens that can be used as candidate markers to develop antitumor vaccine, we

246    develop a neoantigen prediction pipeline to determine what types of mutated peptides in cancer cells

247    could be brought to the cell surface by HLAs based on somatic small mutations (SSMs). In order to

248    maximize prediction accuracy, we include PSSMHCpan, NetMHC-4.0, NetMHCpan-3.0 and

249    PickPocket into our pipeline to detect neoantigens in TCGA tumor samples as following (Figure 5a).

250    We first annotate missense SSMs including single nucleotide variants (SNVs), insertions and deletions

12

251 (InDels) with ANNOVAR [34] to create a list of tumor-specific peptides (8-13) with an in-house script.

252 After HLA alleles are predicted with Seq2HLA [35], we predict neoantigens with PSSMHCpan,

253 NetMHC-4.0, NetMHCpan-3.0 and PickPocket, respectively. Finally, we select a list of neoantigens

254 that meet the following conditions: 1) Predicting as binders (IC50<500nM) by at least 2 software and

255 taking the median value of IC50 as final result; 2) The IC50 value of a given SNV-derived neoantigen

256 must be smaller than that of its corresponding wile type (WT) peptide [36]. Using this pipeline, we

257 analyzed the neoantigens across 10 cancer types from TCGA cohort.

258 Totally we identified 117,017 neoantigens from 467 TCGA cancer samples. We calculated the

259 number of neoantigens per SSM in different types of cancer and observed that STAD, PRAD and

260 BRCA had the highest neoantigens with 2.54, 1.52 and 1.43 per SNV, respectively (Figure 5b), whereas

261 the highest neoantigens per InDel were 2.76, 2.59 and 2.34 in PRAD, STAD and KIRC, respectively

262 (Figure 5c). We also compared the neoantigen loads (number of neoantigens per sample) across 10

263 cancer types and found that STAD, COAD and BLCA tumors had the highest neoantigen loads with

264 median values of 302, 182 and 163, while the THCA tumors had a lowest median neoantigen load of 30

265 (Figure 5d).

266 On average we identified 251 neoantigens in each tumor. We then investigated whether the

267 expression level of HLA class I would be increased in cancer cells to bind neoantigens. Indeed, by

268 looking at the mRNA expression in 467 TCGA tumor samples and their paired normal tissues, we

269 found that the expression of HLA class I was markedly elevated in most tumors (Figure 5e). Since the

270 amount of neoantigens differs substantially among different tumors, we examined whether the number

271 of neoantigens was correlated with HLA class I expression level in each tumor. However, we did not

272 find a correlation between the number of neoantigens and the HLA class I expression levels in tumors

273    (Pearson' rho=-0.05, *P*=0.33).

274

275    **Discussion**

276    Designing antitumor vaccine requires predicting peptide-HLA binding affinity with high accuracy. In

277    this article, we have presented a novel software PSSMHCpan that allows us to predict peptide binding

278    affinity with a broad coverage of HLA class I alleles. By comparing our PSSMHCpan with the most

279    popular machine learning based methods NetMHC-4.0, NetMHCpan-3.0 and the most recently

280    published PSSM based method PickPocket, we demonstrated that overall our PSSMHCpan is

281    substantially better than the other three in predicting peptide-HLA binding affinity, in terms of accuracy

282    and efficiency.

283    In recent years, PSSM based methods to predict peptide-HLA binding affinity were gradually

284    replaced by machine learning based methods that are believed to have reliable accuracy and larger data

285    prediction capability [3]. However, by comparing our PSSMHCpan with machine learning based

286    methods NetMHC-4.0 and NetMHCpan-3.0, we show that our PSSMHCpan exhibits a higher

287    predicting accuracy than NetMHC-4.0 and NetMHCpan-3.0, respectively. In terms of data prediction

288    capability, PSSMHCpan can allele-specifically and pan-specifically predict peptides that bind to 241

289    and 4778 HLA class I alleles, while NetMHC-4.0 and NetMHCpan-3.0 can only predict 89 and 2924

290    HLA class I alleles, respectively. Furthermore, the PSSMHCpan displays much higher prediction

291    efficiency as compared to NetMHC-4.0 and NetMHCpan-3.0 (Table 6).

292    We noticed that the size of training database appeared to directly affect the prediction accuracy. A

293    larger training database could improve the prediction accuracy of PSSMHCpan. For instance, the

294    PSSMHCpan prediction accuracy ACC in predicting 9mer peptides bind to HLA-A*0101 and

14

295    HLA-B*5703 are 0.96 and 0.70. We found that in our training database, there are 813 binders for

296    HLA-A*0101 and only 25 binders for HLA-B*5703, respectively. We believed that in order to improve

297    the prediction accuracy, it is necessary to increase the size of training database.

298    Based on the evaluation results (Figure 4), we recognized that none of the available software is

299    perfect and that in order to maximize the prediction accuracy, it is necessary to use multiple software.

300    We then included PSSMHCpan, NetMHC-4.0, NetMHCpan-3.0 and PickPocket to build a neoantigen

301    prediction pipeline that allowed us to detect 117,017 neoantigens in 467 TCGA tumor samples across

302    10 types of cancer. We believe that in order to provide actionable neoantigens that can be used in

303    cancer immunotherapy, it requires more efforts to validate the function and immunogenicity of the

304    predicted neoantigens experimentally.

305    In conclusion, our PSSMHCpan can predict peptide binding affinity with a broad coverage of HLA

306    class I alleles more accurately and efficiently compared with currently most popular peptide binding

307    affinity prediction software. Our PSSMHCpan can not only help develop personalized antitumor

308    vaccines, but also has great potentials in other aspects of cancer immunotherapy including designing

309    dendritic cell (DC) vaccines, inducing DC-CTL, TCR-T, and assessing the PD-1/CTLA4 prognosis.

310

311    **Availability and requirements**

312    ● Project name: PSSMHCpan

313    ● Project home page: https://github.com/BGI2016/PSSMHCpan

314    ● Operating system: Platform independent

315    ● Programming language: Perl

316    ● Other requirements: ActivePerl 5.8

15

317 • License: OSI

318

**Availability of supporting data and materials**

320 The supporting data from this study will be hosted in the additional files and PSSMHCpan home page.

321

**Additional file**

323 Additional file 1: Supplementary tables for supporting the analysis part

324 Table S1 is the list of HLA class I alleles for allele-specific and pan-specific prediction. Table S2 is

325 10-cross validations results of alleles-specific prediction. Table S3 is the pan-specific prediction results.

326 Table S4 is prediction results the independent dataset.

327

**Competing interests**

329 The authors declare no competing financial interests.

330

**Authors' contributions**

332 G. L., D. L, B. L. Y. H, J. W. and H. Y. conceived of study and designed the project. G. L. and D. L.

333 performed software development, computational analyses and prepared figures. S. Q., W. L. performed

334 pan-cancer neoantigen analysis. G. L., B. L. and K. M. wrote the manuscript.

335

**Acknowledgements**

342

16

343 **Reference**

344 1. Liao WW, Arthur JW. Predicting peptide binding to Major Histocompatibility Complex molecules.
345 Autoimmunity reviews. 2011;10(8):469-73. doi:10.1016/j.autrev.2011.02.003.

346 2. Schueler-Furman O, Altuvia Y, Sette A, Margalit H. Structure-based prediction of binding peptides to
347 MHC class I molecules: application to a broad range of MHC alleles. Protein science : a publication of
348 the Protein Society. 2000;9(9):1838-46. doi:10.1110/ps.9.9.1838.

349 3. Luo H, Ye H, Ng HW, Shi L, Tong W, Mendrick DL et al. Machine Learning Methods for Predicting
350 HLA-Peptide Binding Activity. Bioinformatics and biology insights. 2015;9(Suppl 3):21-9.
351 doi:10.4137/BBI.S29466.

352 4. Zhang GL, Ansari HR, Bradley P, Cawley GC, Hertz T, Hu X et al. Machine learning competition in
353 immunology - Prediction of HLA class I binding peptides. Journal of immunological methods.
354 2011;374(1-2):1-4. doi:10.1016/j.jim.2011.09.010.

355 5. Nielsen M, Lundegaard C, Worning P, Lauemoller SL, Lamberth K, Buus S et al. Reliable prediction of
356 T-cell epitopes using neural networks with novel sequence representations. Protein science : a
357 publication of the Protein Society. 2003;12(5):1007-17. doi:10.1110/ps.0239403.

358 6. Carreno BM, Magrini V, Becker-Hapak M, Kaabinejadian S, Hundal J, Petti AA et al. Cancer
359 immunotherapy. A dendritic cell vaccine increases the breadth and diversity of melanoma
360 neoantigen-specific T cells. Science. 2015;348(6236):803-8. doi:10.1126/science.aaa3828.

361 7. Yadav M, Jhunjhunwala S, Phung QT, Lupardus P, Tanguay J, Bumbaca S et al. Predicting
362 immunogenic tumour mutations by combining mass spectrometry and exome sequencing. Nature.
363 2014;515(7528):572-6. doi:10.1038/nature14001.

364 8. Walter S, Weinschenk T, Stenzl A, Zdrojowy R, Pluzanska A, Szczylik C et al. Multipeptide immune
365 response to cancer vaccine IMA901 after single-dose cyclophosphamide associates with longer patient
366 survival. Nature medicine. 2012;18(8):1254-61. doi:10.1038/nm.2883.

367 9. Lata S, Bhasin M, Raghava GP. Application of machine learning techniques in predicting MHC
368 binders. Methods in molecular biology. 2007;409:201-15. doi:10.1007/978-1-60327-118-9_14.

369 10. Noguchi H, Kato R, Hanai T, Matsubara Y, Honda H, Brusic V et al. Hidden Markov model-based
370 prediction of antigenic peptides that interact with MHC class II molecules. Journal of bioscience and
371 bioengineering. 2002;94(3):264-70.

372 11. Lundegaard C, Lund O, Nielsen M. Prediction of epitopes using neural network based methods.
373 Journal of immunological methods. 2011;374(1-2):26-34. doi:10.1016/j.jim.2010.10.011.

374 12. Andreatta M, Nielsen M. Gapped sequence alignment using artificial neural networks: application
375 to the MHC class I system. Bioinformatics. 2016;32(4):511-7. doi:10.1093/bioinformatics/btv639.

376 13. Hoof I, Peters B, Sidney J, Pedersen LE, Sette A, Lund O et al. NetMHCpan, a method for MHC class
377 I binding prediction beyond humans. Immunogenetics. 2009;61(1):1-13.
378 doi:10.1007/s00251-008-0341-z.

379 14. Nielsen M, Lundegaard C, Blicher T, Lamberth K, Harndahl M, Justesen S et al. NetMHCpan, a
380 method for quantitative predictions of peptide binding to any HLA-A and -B locus protein of known
381 sequence. PloS one. 2007;2(8):e796. doi:10.1371/journal.pone.0000796.

382 15. Nielsen M, Andreatta M. NetMHCpan-3.0; improved prediction of binding to MHC class I
383 molecules integrating information from multiple receptor and peptide length datasets. Genome
384 medicine. 2016;8(1):33. doi:10.1186/s13073-016-0288-x.

385 16. Zhang H, Lund O, Nielsen M. The PickPocket method for predicting binding specificities for
386 receptors based on receptor pocket similarities: application to MHC-peptide binding. Bioinformatics.

17

387    2009;25(10):1293-9. doi:10.1093/bioinformatics/btp137.

388    17. Donnes P, Kohlbacher O. SVMHC: a server for prediction of MHC-binding peptides. Nucleic acids
389    research. 2006;34(Web Server issue):W194-7. doi:10.1093/nar/gkl284.

390    18. Bhasin M, Raghava GP. A hybrid approach for predicting promiscuous MHC class I restricted T cell
391    epitopes. Journal of biosciences. 2007;32(1):31-42.

392    19. Karosiene E, Lundegaard C, Lund O, Nielsen M. NetMHCcons: a consensus method for the major
393    histocompatibility    complex    class    I    predictions.    Immunogenetics.    2012;64(3):177-86.
394    doi:10.1007/s00251-011-0579-8.

395    20. Trolle T, Metushi IG, Greenbaum JA, Kim Y, Sidney J, Lund O et al. Automated benchmarking of
396    peptide-MHC    class    I    binding    predictions.    Bioinformatics.    2015;31(13):2174-81.
397    doi:10.1093/bioinformatics/btv123.

398    21. Mukherjee S, Bhattacharyya C, Chandra N. HLaffy: estimating peptide affinities for Class-1 HLA
399    molecules    by    learning    position-specific    pair    potentials.    Bioinformatics.    2016.
400    doi:10.1093/bioinformatics/btw156.

401    22. Backert L, Kohlbacher O. Immunoinformatics and epitope prediction in the age of genomic
402    medicine. Genome medicine. 2015;7(1):119. doi:10.1186/s13073-015-0245-0.

403    23. Xia X. Position weight matrix, gibbs sampler, and the associated significance tests in motif
404    characterization and prediction. Scientifica. 2012;2012:917540. doi:10.6064/2012/917540.

405    24. Toh H, Savoie CJ, Kamikawaji N, Muta S, Sasazuki T, Kuhara S. Changes at the floor of the
406    peptide-binding groove induce a strong preference for proline at position 3 of the bound peptide:
407    molecular    dynamics    simulations    of    HLA-A*0217.    Biopolymers.    2000;54(5):318-27.
408    doi:10.1002/1097-0282(20001015)54:5<318::AID-BIP30>3.0.CO;2-T.

409    25. Apweiler R, Bairoch A, Wu CH, Barker WC, Boeckmann B, Ferro S et al. UniProt: the Universal
410    Protein    knowledgebase.    Nucleic    acids    research.    2004;32(Database    issue):D115-9.
411    doi:10.1093/nar/gkh131.

412    26. Altschul SF, Gertz EM, Agarwala R, Schaffer AA, Yu YK. PSI-BLAST pseudocounts and the minimum
413    description length principle. Nucleic acids research. 2009;37(3):815-24. doi:10.1093/nar/gkn981.

414    27. Styczynski MP, Jensen KL, Rigoutsos I, Stephanopoulos G. BLOSUM62 miscalculations improve
415    search performance. Nature biotechnology. 2008;26(3):274-5. doi:10.1038/nbt0308-274.

416    28. Vita R, Overton JA, Greenbaum JA, Ponomarenko J, Clark JD, Cantrell JR et al. The immune epitope
417    database    (IEDB)    3.0.    Nucleic    acids    research.    2015;43(Database    issue):D405-12.
418    doi:10.1093/nar/gku938.

419    29. Kim Y, Sidney J, Buus S, Sette A, Nielsen M, Peters B. Dataset size and composition impact the
420    reliability of performance benchmarks for peptide-MHC binding predictions. BMC bioinformatics.
421    2014;15:241. doi:10.1186/1471-2105-15-241.

422    30. Schuler MM, Nastke MD, Stevanovikc S. SYFPEITHI: database for searching and T-cell epitope
423    prediction. Methods in molecular biology. 2007;409:75-93.

424    31. Bhasin M, Singh H, Raghava GP. MHCBN: a comprehensive database of MHC binding and
425    non-binding peptides. Bioinformatics. 2003;19(5):665-6.

426    32. Robinson J, Soormally AR, Hayhurst JD, Marsh SG. The IPD-IMGT/HLA Database - New
427    developments    in    reporting    HLA    variation.    Human    immunology.    2016.
428    doi:10.1016/j.humimm.2016.01.020.

429    33. Vigneron N, Stroobant V, Van den Eynde BJ, van der Bruggen P. Database of T cell-defined human
430    tumor antigens: the 2013 update. Cancer immunity. 2013;13:15.

18

431  34. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from
432  high-throughput sequencing data. Nucleic acids research. 2010;38(16):e164. doi:10.1093/nar/gkq603.
433  35. Boegel S, Lower M, Schafer M, Bukur T, de Graaf J, Boisguerin V et al. HLA typing from RNA-Seq
434  sequence reads. Genome medicine. 2012;4(12):102. doi:10.1186/gm403.
435  36. Hundal J, Carreno BM, Petti AA, Linette GP, Griffith OL, Mardis ER et al. pVAC-Seq: A
436  genome-guided in silico approach to identifying tumor neoantigens. Genome medicine. 2016;8(1):11.
437  doi:10.1186/s13073-016-0264-5.

438

439  **FIGURE LEGENDS**

440  **Figure 1** Heat map of HLA protein sequence similarity. The larger the Z-Score, the more similar of the

441  pair HLA protein sequences. It showed high similarity between different types of HLA alleles within

442  the same gene locus.

443  **Figure 2** Method of PSSMHCpan. The three mainly steps are shown in grey background.

444  **Figure 3** Evaluation on broad HLA allelic coverage. (a) The allele-specific prediction evaluation

445  results showed by ROC curse of PSSMHCpan, NetMHC-4.0, NetMHCpan-3.0 and PickPocket. This

446  result was except 9mer and 10mer of HLA-A*0101, HLA-A*0201 and HLA-B*0702. The ACC,

447  sensitivity and specificity at cutoff of 500nM were also shown. (b) The boxplot of individual ACC of

448  particular HLA allele with fixed peptide length. Comparison between PSSMHCpan and other three

449  methods were performed by using paired T test. "*" denotes $P<0.05$ and "**" denotes $P<0.01$. (c) The

450  evaluation results showed by ROC curse of PSSMHCpan in pan-specific prediction, NetMHCpan-3.0

451  and PickPocket. The ACC, sensitivity and specificity at cutoff of 500nM were also shown. (d)

452  Correlation analysis of peptide-HLA binding affinity result of IC50 value in log2 between

453  allele-specific prediction and pan-specific prediction.

454  **Figure 4** The evaluation result of the independent dataset. We denoted IC50<500nM was positive

455  prediction.

456  **Figure 5** Pan-cancer neoantigens. (a) The flow-char of neoantigen prediction pipeline. Software with

19

457  parameters using in the pipeline are shown in dashed procedure. (b) The distribution of neoantigens

458  generated from each SNV across diverse cancers. (c) The distribution of neoantigens generated from

459  each InDel across diverse cancers. (d) The distribution of neoantigen loads across 10 cancer types. The

460  cancer types are sorted by median value of neoantigen loads. (e) The expression of HLA class I in

461  tumor and corresponding normal samples.

Figure

**a**



ACC$^{500nM}_{PSSMHCpan}$=0.92

ACC$^{500nM}_{NetMHC}$=0.85

ACC$^{500nM}_{NetMHCpan}$=0.85

ACC$^{500nM}_{PickPocket}$=0.72

PSSMHCpan:0.97
NetMHC:0.96
NetMHCpan:0.97
PickPocket:0.95

**b**



**c**



ACC$^{500nM}_{PSSMHCpan}$=0.87

ACC$^{500nM}_{NetMHCpan}$=0.75

ACC$^{500nM}_{PickPocket}$=0.73

PSSMHCpan:0.93
NetMHCpan:0.96
PickPocket:0.96
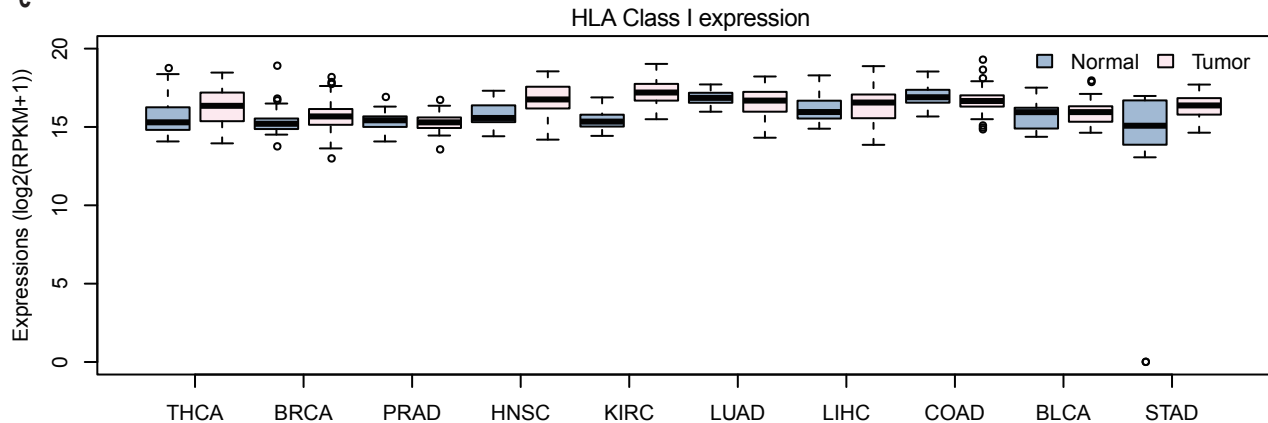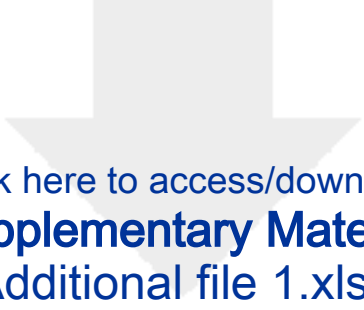
**d**

Pearson's rho = 0.89
P−value < 0.01

Figure

Click here to access/download
**Supplementary Material**
Additional file 1.xlsx