1

2

# A Generalized Similarity Metric

# for Predicting Peptide Binding Affinity

5

6

Jacob Rodriguez[1,2], Siddharth Rath[1,2,3], Jonathan Francis-Landau[1,4], Yekta Demirci[5],

Burak Berk Ustundag[6], and Mehmet Sarikaya[1,2,3,7,8*]

9

[1] GEMSEC, Genetically Engineered Materials Science and Engineering Center, University of

Washington, Seattle, WA 98195, USA

[2] Department of Materials Science and Engineering, University of Washington, Seattle, WA

98195, USA

[3] Molecular Engineering and Science Institute, University of Washington, Seattle, WA 98195,

USA

[4] Department of Mathematics, University of Washington, Seattle, WA 98195, USA

[5] Department of Electrical and Electronics Engineering, Middle East Technical University

[6] Department of Computer and Informatics Engineering, Istanbul Technical University

[7] Department of Chemical Engineering, University of Washington, Seattle, WA 98195, USA

[8] Department of Oral Health Sciences, University of Washington, Seattle, WA 98195, USA

*Corresponding Author

E-mail: sarikaya@uw.edu

23

## Abstract

The ability to capture the relationship between similarity and functionality would enable the predictive design of peptide sequences for a wide range of implementations from developing new drugs to molecular scaffolds in tissue engineering and biomolecular building blocks in nanobiotechnology. Similarity matrices are widely used for detecting sequence homology but depend on the assumption that amino acid mutational frequencies reflected by each matrix are relevant to the system in which they are applied. Increasingly, neural networks and other statistical learning models solve problems related to functional prediction but avoid using known features to circumvent unconscious bias. We demonstrated an *iterative* alignment method that enhances predictive power of similarity matrices based a similarity metric, the Total Similarity Score. A generalized method is provided for application to amino acid sequences from inorganic and organic systems by benchmarking it on the debut quartz-binder set and 3 peptide-protein sets from the Immune Epitope Database. Pearson and Spearman Rank Correlations show that by treating the gapless Total Similarity Score as a predictor of relative binding affinity, prediction of test data has a 0.5-0.7 Pearson and Spearman Rank correlation. with respect to size of dataset. Since the benchmarks used herein are from a solid-binding peptide and a protein-peptide system, our proposed method could prove to be a highly effective general approach for establishing the predictive sequence-function relationships of among the peptides with different sequences and lengths in a wide range of biotechnology, nanomedicine and bioinformatics applications.

## Introduction and Background

The rapid development of target-specific drugs relies on the development of high-throughput and accurate methods of modelling molecular structures. The biology, pharmacology and bioengineering communities are interested in building widely applicable methods founded in predictive design of molecules that have specificity for biological targets, analytes and biomarkers [1-4]. Small peptides (7 to 40 amino acids) have high potential as both therapeutics [5-7] and high-performance molecular building blocks [8-10] due their diversity of binding affinity both quantitatively and specifically across 2D- and nano-materials.

Towards more accurate and fast predictions of affinity or conformation that would enable high-throughput drug and targeting peptide design, among some of the best performing methodologies are stochastic models such as NetMHCpan-4.0 [11], DeepMHC [12] and MHCflurry [13]. These methods use little or no prior information about the peptides to ensure only random walk identifies relevant patterns. By avoiding physiochemical properties published in the literature, these models are subject to inconsistent predictions between test peptide sets even for the same protein target. Alignment-free neural networks models have shown substantial success in predicting the binding affinity of the Immune Epitope Database (IEDB, www.iedb.org) datasets [12,14]. To avoid overfitting, they require hundreds of thousands of sequences and are not optimized for gaps in the binding domains [15,16].

The current state of the art in modelling tools, e.g., molecular dynamics (MD), molecular mechanics (MM), and Monte Carlo (MC) based methods, predict overall conformation from which binding energies may be calculated [9]. These approaches

70    utilize knowledge-based force fields [17,18] and energy minimization techniques to

71    sample the most probable structures [19]. Though solving conformational structures will

72    likely enable the most accurate predictions of peptide function, to date structural

73    information is avoided in models requiring large amounts of data. This is mostly due to

74    the large computational cost associated with calculating molecular structures of these

75    large molecules, which is a barrier to the development of both highly complex neural

76    networks and current MD/MC-based methods. The deeper networks rely less on learning

77    in space constrained by verified physiochemical trends and more on the number of

78    parameters and computational power. Less complex and more interpretable models

79    integrate known patterns while leaving space for optimization methods to learn unknown

80    patterns in the sequences.

81       Current alignment-based methods for high-throughput prediction functionality of

82    amino acid sequence information can be separated into two groups; pairwise [20-22] and

83    multiple sequence [16,23,24]. In general, pairwise alignment is ideal for shorter

84    sequences due to its higher computational cost per amino acid and is widely accepted to

85    be the optimal alignment [25]. Multiple sequence alignment is considered more

86    appropriate for longer sequences with suspected consensus domains. In both methods

87    Point Accepted Mutation (PAM) and Blocks Substitution Matrix (BLOSUM) matrices are

88    still the most widely used, and there are permutations of these matrices to serve more

89    specific tasks [17,26,27]. Overall, the limitations of PAM and BLOSUM provided

90    inspiration and guidance for generating matrices with increased accuracy based on larger

91    and more complete datasets [11,28-30]. Matrices such the PMBEC [27], have been

92    generated based on the two models that produce a minor increase in performance but

93    ultimately are vulnerable to the same factors as their predecessors [11]. In 2008, for

94    example, a miscalculation was discovered in the clustering protocol of the BLOSUM

95    matrix [31]. Despite extensive characterization of the mistake, BLOSUM is still the

96    standard for one of the largest alignment-capable databases available to date, BLAST

97    [11].

98        In contrast with PAM, BLOSUM, PMBEC [26] and the SAUSAGE Force Field

99    Matrix [16], the novel OCSimM and 8 property group-derived matrices (A-RMat) were

100   calculated from 527 physiochemical properties of amino acids [32,33]. AAindex is a vast

101   resource of high-quality amino acid properties collected from literature dating from the

102   mid-sixties to today [32,33]. Typically, either variable reduction methods (Principal

103   Component Analysis [6,34] or Factor Analysis [35]) or heuristic selection is performed to

104   shrink the huge dataset of over 550 amino acid properties to obtain an interpretable

105   solution. Variable reduction has significant advantages over a global analysis of

106   heuristically grouped properties because human error cannot influence the potential

107   relationships observed [35]. However, these methods still assume the relationship

108   between high-specificity peptides and low-specificity peptides is described by

109   physiochemical properties.

110       Previously, we have successfully used a matrix optimization method to a group of

111   peptides that were categorized as strong, weak or medium binders based on their binding

112   affinities to crystalline silica, quartz, using 40 sequences that were originally genetically

113   selected using M13 phage display peptide library [25]. The novel metric called the Total

114   Similarity Score (TSS$_{A-B}$) describes the average Global Alignment score of all peptides

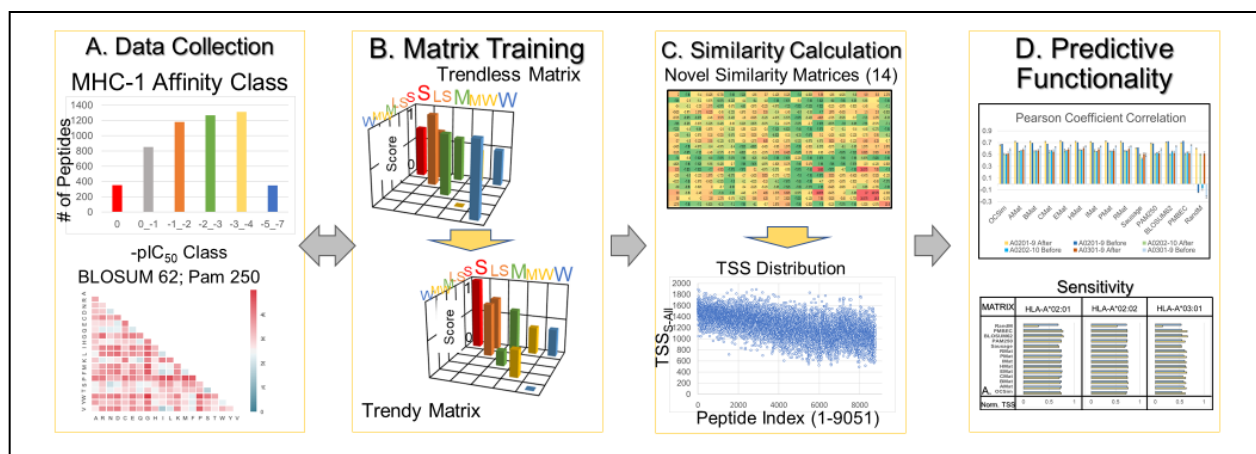115   from group-A to all of group-B [25]. The TSS score quantifies the similarity of a peptide to

116  a functional peptide set (i.e. affinity for a solid material). By keeping random changes to

117  a similarity matrix that increased the $TSS_{S-S}$ (TSS of strong binders with strong binders),

118  and decrease the $TSS_{S-W}$ (TSS of strong binders with weak binders), a similarity matrix

119  was obtained that could predict the semi-quantitative affinity of quartz-binding peptides

120  with 70-80% success. Despite its high predictive power, TSS has never been applied to

121  MHC data. Using the MHC data, here we demonstrate its implementation that strongly

122  suggests that TSS could be a predictive method for establishing sequence-function

123  relationships in a variety of large sequence-based data sets.

124      The reliable prediction of peptide binding affinity has already led to ground-

125  breaking advances in oral health science and will continue to do so in areas requiring a

126  well-described soft interface between peptides and solid-state inorganic materials

127  [5,10,36]. Though affinity prediction is not the most descriptive or important

128  characterization of peptides, understanding the relationship among solid-binding peptides

129  [10] has led to many technologies such as sensors with high sensitivity, [5] assemblers in

130  nanotechnology, and tiny enzymes in biomineralization [37].

131  **Approach and Methodology**

132      Iterative Alignment (IA) creates a scoring matrix that provides scores correlating

133  with the positional composition of a peptide when compared to a weak and a strong

134  binding set. When a sequence of interest has high similarity to these strong binders and

135  low similarity to the weak binders, the sequence was given a higher $TSS_{Seq-S}$ and lower

136  $TSS_{Seq-W}$ (TSS of interesting sequence to weak binders). Training the similarity matrix

137  was done by increasing the differences in TSS to strong binders for two binding affinity

138  classes, strong and weak.

139    First, peptides were sorted by their affinity values shown in Fig 1A. The generated

140    trend was characterized by the positive correlation $TSS_{Seq-S}$ with binding affinity visualized

141    in the lower bar chart of Fig 1B. Once training was finished, we calculated the TSS to

142    strong binders a final time for all peptides in the set. The results from the trained randomly

143    initialized matrix (RandM) are shown in the scatterplot in Fig 1C. Next, several methods

144    are used to measure correlation of $TSS_{Seq-S}$ with the experimental affinity including

145    Pearson and Spearman Rank correlation, Root Mean Square Error and a binary

146    classification scheme (binder/nonbinder prediction). A sample of these results are

147    demonstrated in Fig 1D.



148    **Fig 1. Schematic of Iterative Alignment Procedure.** The iterative alignment procedure is executed

149    in four separate steps as show in the flow chart, that include: (A) Classification of MHC-I binding

150    peptides from the IEDB and the resultant matrices from AAindex; (B) Training the randomly

151    initialized matrix (RandM) which, before training, was uncapable to demonstrate the trend of

152    decreasing cross-similarity, but after training it becomes prominent indicating the successful

153    integration of the information; (C) Demonstration the total similarity score of the full allele set with

154    respect to the strongest determined binders of *HLA-A\*02:01* (TSS $_{HLA-A*02:01-S}$) for trained RandM.

155    Calculations were performed for all matrices before and after training; (D) Showcase correlation

156    and accuracy measurements (see details in the text and figures below).

157 **Data Collection**

158       Peptide sequences with affinity for HLA alleles were obtained from the Immune

159     Epitope Database (www.iedb.org), a common source of training and benchmark data for

160     predictive models of peptide function [14]. Quartz binders and the Quartz I matrix were

161     provided by GEMSEC at the MSE Department of the University of Washington [25]. The

162     Amino Acid Index (AAindex) is a large database of amino acid properties that were used

163     to calculate the cluster matrices (A-RMat) [32,33]. Within the site, similarity matrices

164     calculated by various studies are also provided, and it was from here that the SAUSAGE

165     force-field matrix was also chosen [17]. The PMBEC scoring-matrix [27] was included as

166     it was derived directly from binding affinity data from MHC-I. In general, the matrices

167     chosen are a diverse subsection of the types of information used to describe differences

168     between amino acids and therefore were an appropriate selection for yielding conclusions

169     about how the seed matrix would affect the overall result.

170 **Novel Matrix Calculation**

171       To explore the possibility that certain properties, e.g., hydrophobicity, electrical

172     properties, amino acid composition etc., may make better seed matrices, 9 similarity

173     matrices were calculated based on clusters optimized by Saha *et. al* [38]. After grouping

174     properties by alpha-helix or beta-sheet propensities, composition, electrical, hydrophobic,

175     and intrinsic characteristics, residue propensity, and physicochemical properties, we

176     performed Principal Component Analysis (PCA) on each group and all groups combined.

177     Using a Python library downloaded from scikit-learn.org [39], the principal components

178     were calculated which were most representative of the internal variation of property

179     subset. Because these principal components are orthogonal, Euclidean distance was the

180     most appropriate for calculating the actual similarity matrix. By calculating the difference

181     between the principal components of two amino acids, we were able to calculate nine (20

182     x 20) similarity matrices describing their quantitative physiochemical differences. These

183     matrices will be referred to for the rest of the work as AMat (Alpha-helix propensity), BMat

184     (Beta-sheet propensity), CMat (Composition), EMat (Electric), HMat (Hydrophobicity),

185     IMat (Intrinsic propensity), PMat (Physiochemical), RMat (Residue propensity) and

186     OCSim [Orthogonal Component Similarity matrix (all properties)].

187     **Code Implementation**

188        The newest version of the algorithm was written in Python, using a gapless scoring

189     method to calculate TSS scores. The gap calculation was excluded to rectify the issue

190     created by the changing gap position in each sequence. Per peptide-peptide scoring

191     operation (300 strong binders x 9000 peptides for *HLA-A\*02:01*), per iteration (5000 due

192     to randomly changing mutabilities) the gap is placed in one position. We suspected the

193     gap made recognizing the consistent amino acids between iterations difficult. The debut

194     implementation of the method [25] iteratively aligned less than 20 peptides per strong and

195     weak binding group. The IEDB dataset being substantially larger (i.e., over 9000 peptides

196     for the largest set) required the inclusion of more peptides per set in order to capture as

197     many of the features pertaining to binding affinity as possible.

198     **Designation of Affinity Classes**

199        The peptide sequences were first ordered by -$pIC_{50}$, and then segregated into

200     groups dependent on their affinity. For example, all peptides within the 3 chosen alleles

201     (*HLA-A\*02:01* [9-length], *HLA-A\*02:02* [10-length] and *HLA-A\*03:01* [9-length]) with a -

202     $pIC_{50}$ of 0 were named 'strong' (S) binders, creating 3 sets. The 'weak' (W) binders for

203   the 9-length and 10-length sets were those with a -pIC$_{50}$ of -5 to -7. From these, 80% of

204   a strong or weak peptide list was randomly chosen as training sets to obtain cross-

205   validation. To show the flexibility of the method, we chose several groups with differing

206   distributions to demonstrate the improvements are still achieved when only partial data is

207   available.

208   **Matrix Training**

209   To begin, two lists of peptide sequences (at least 6 in each) must be obtained, one

210   with higher 'internal similarity' and lower 'internal similarity'. Critically, peptides with high

211   binding affinity for the same material will also higher 'internal similarity' and those with low

212   affinity will have low 'internal similarity' [25]. Internal similarity refers to the sum of Global

213   Alignment (GA) scores of each peptide within a list to every other peptide within the same

214   list. Global Alignment is commonly referred to as the Needleman-Wunsch algorithm or

215   optimal alignment as it always obtains the optimum number and placement of gaps,

216   resulting in the most similar domains being recognized and aligned when they are

217   consensus [20]. It requires a similarity matrix to obtain scores between matches or

218   mismatches of amino acids, and many of these have been calculated throughout the

219   literature. For small peptides, it may not be the ideal alignment method considering their

220   short length makes scoring the entire sequence important.

221   While guaranteeing the optimal alignment, Global Alignment is computationally

222   very expensive and therefore impractical to apply to larger groups of sequences than

223   those used in previous work [25] (10 - 20 sequences per strong and weak group). The

224   updated method departs from the alignment methodology and scores peptides by their

225   positional composition only, which is essentially the same score without the gap

226  calculation. By greatly expanding the number of peptides used in the strong group, the

227  significance of GA is reduced due to a wider range of domain types and locations being

228  represented. In general, scoring with more peptides is just as beneficial as scoring a few

229  with GA. Global Alignment expands the number of sequences a peptide will have

230  consensus with; in a way making it appear as many peptides in the strong group.

231  However, the domains being aligned and the values scoring the alignments are different

232  from one iteration to another, resulting in a lack of consistent scoring between sequence

233  domains. Therefore, we justify the departure from GA as both a necessity and a benefit

234  to ensure the method runs within a practical time constraint.
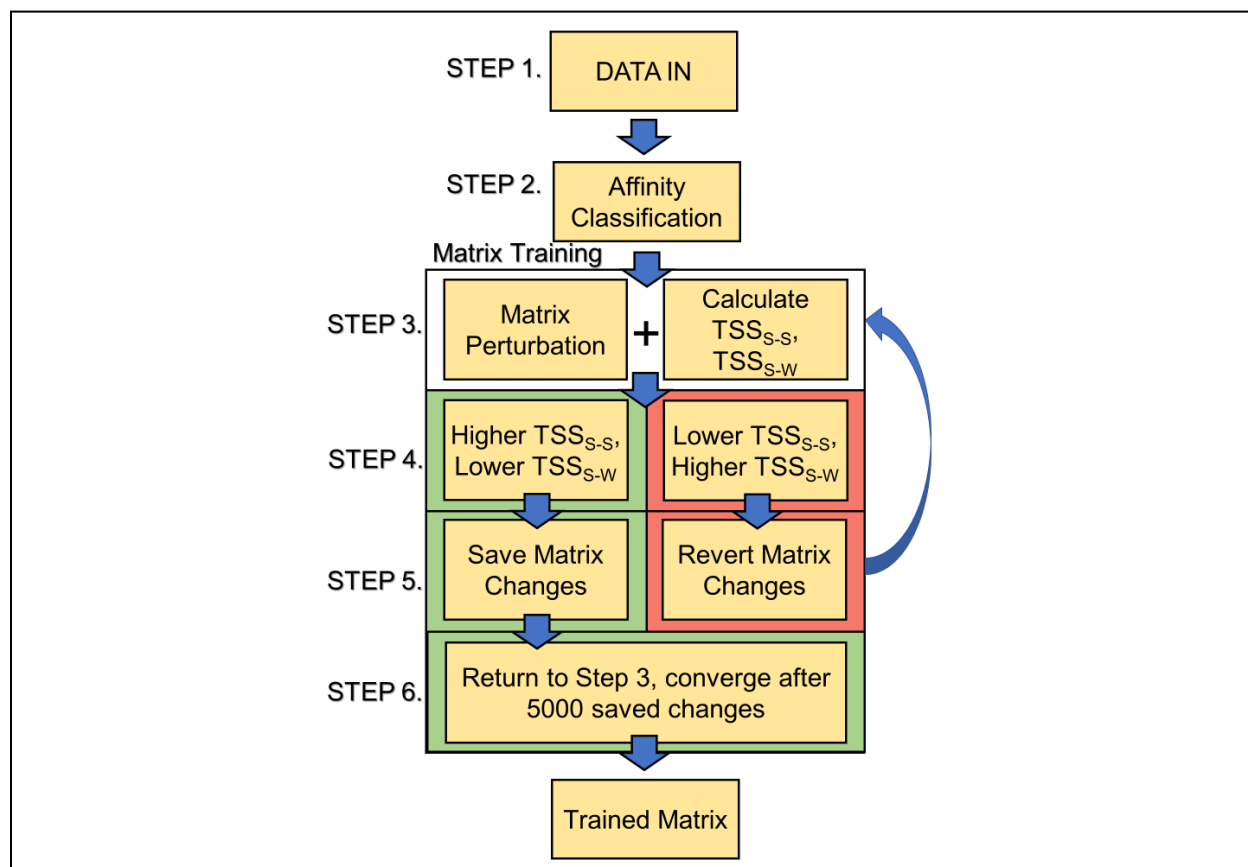
235      The procedure for one iteration can be described in 6 steps (see Fig 2). After the

236  affinity classes have been designated (Fig 2, Step 2), a seed similarity matrix is used to

237  calculate TSS$_{S-S}$ and TSS$_{S-W}$ (same as internal but to separate group of peptides)

238  similarity for each peptide (Fig 2, Step 3). External similarity is calculated by aligning the

239  strong binders to each in the low internal similarity group. Within each list, the average is

240  found and form the cost functions for IA, the Total Similarity Score Strong-Strong (TSS$_{S-}$

241  $_{S}$)   and   Total   Similarity   Score   Strong-Weak   (TSS$_{S-W}$),   respectively.

242  Mathematically, the expression for general TSS calculation is given by Equation (1) as

243  $$TSS_{A-B}\left[|A|\,_{xa}^{ya} - |B|\,_{xb}^{yb}\right] = 1/[xa \times (xb - \delta_{AB})] \times \sum_{i=1,j=1}^{xa,xb} PSS_{ij}(1 - \delta_{ij}\delta_{AB})$$

244                                                                                            (1)

245  where, TSS$_{A-B}$ is the Total Similarity Score (TSS) between peptide sets A and B, PSS$_{ij}$ is

246  the pairwise similarity score (PSS) between sequences i and j of sets A and B

247  respectively, *xa* and *xb* are the total number of sequences in sets A and B, and *δ* is the

248  Kronecker delta function ($\delta_{ij}$ = 1 if i = j, otherwise  $\delta_{ij}$ = 0).

249        After the values of $TSS_{S-S}$ and $TSS_{S-W}$ have been calculated and saved for the first

250    time, the similarity matrix is perturbed by making random changes (1-20) to the matrix

251    values by either adding 1 or subtracting 1 (Fig 2, Step 3). Using the new matrix, $TSS_{S-S}$

252    and $TSS_{S-W}$ are calculated again and compared with the previous TSS (Fig 2, Step 4). A

253    change to the matrix is considered beneficial if $TSS_{S-S,NEW}$ is greater than $TSS_{S-S,OLD}$ and

254    $TSS_{S-W,NEW}$ is less than $TSS_{S-W,OLD}$. Beneficial changes are saved for the next round (Fig

255    2, Step 5). If the change is not beneficial, then the previous matrix (before mutation) is

256    perturbed again and the process repeats (Fig 2, Step 5). The algorithm could continue

257    indefinitely but we considered the matrix converged when over 5,000 iterations occurred

258    without a beneficial change (Fig 2, Step 6).



259    **Fig 2. Schematics of matrix training procedure.** Peptides were first downloaded and classified

260    by their affinity. The similarity matrix is perturbed randomly and then TSS scores are calculated.

261    Depending on the outcome, changes to the matrix were either saved or discarded. The matrix

262    was considered 'converged' after 5000 beneficial changes total, or 5000 negative changes in a

263    row, occur.

264    **Benchmark with the Previous Work**

265    To prove the updated methodology was up to par with the original implementation

266    of the procedure, we obtained the Quartz I matrix and silica binding peptides used by

267    Oren *et. al* [25] The same procedure was followed by mutating PAM250 and training on

268    the same strong and weak groups. After training, IA converged on a matrix capable of

269    predicting binding affinity with similar accuracy to the debut implementation [25]. Using a

270    Pearson correlation of the external similarity to affinity of any silica binding peptide to the

271    group of strong binders designated by [25], we calculated a 51% correlation with our

272    matrix. Previous work obtained a 46% correlation with Quartz I, demonstrating the

273    equivalent capabilities of the updated method. P-values for these correlations were less
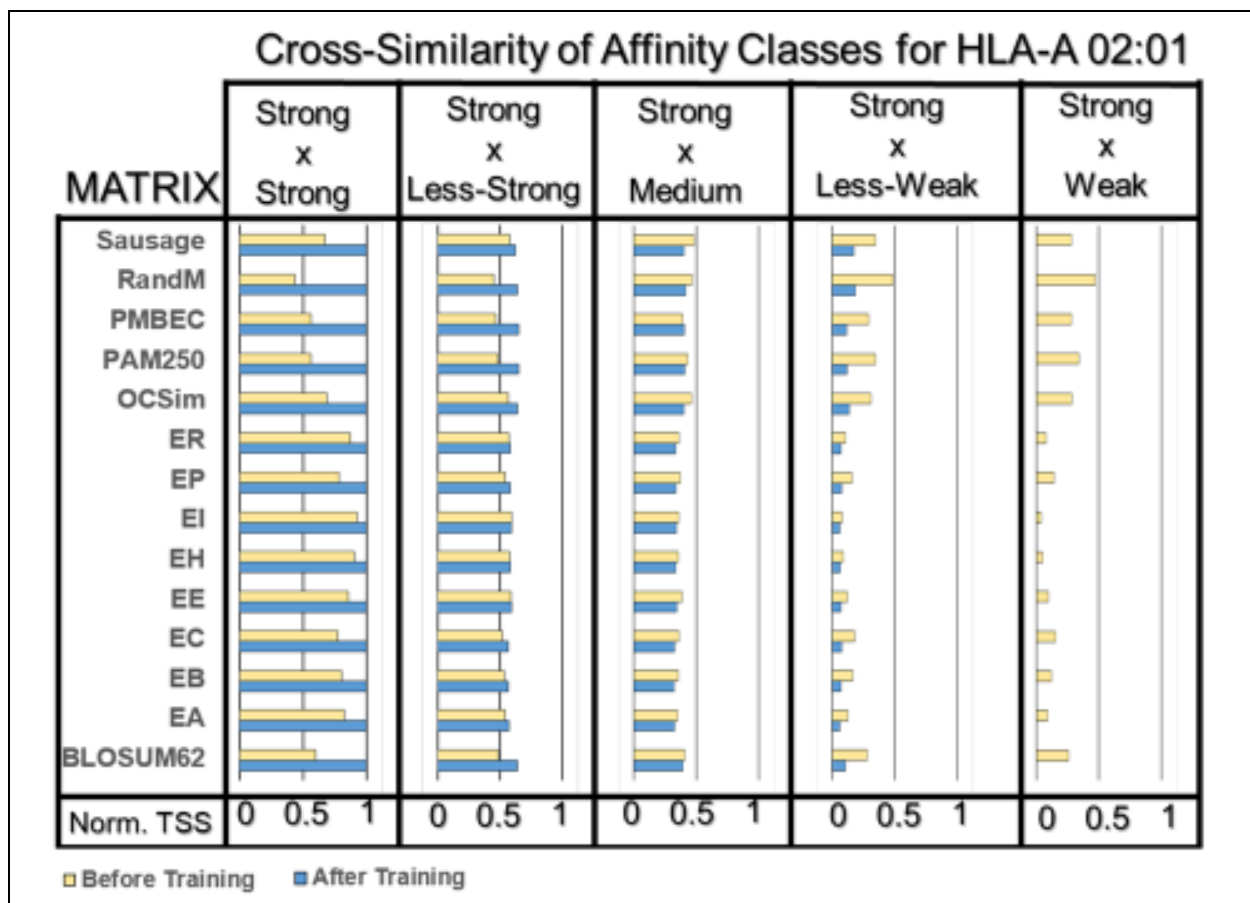
274    than 0.0005.

275    **Application to MHC Data**

276    To test whether the modified methodology would perform on organic materials, we

277    needed a set of peptides with affinity for a biological target. The IEDB provides high quality

278    sequence data including binding affinities for multiple Major-Histocompatibility

279    Complexes which provided a perfect opportunity to test performance [14]. By designating

280    peptides with -pIC$_{50}$ (negative logarithm of IC$_{50}$) of 0 as strong-binders peptides and weak-

281    binders having -pIC$_{50}$ of -5 to -7 (Fig 2, Step 2) from three alleles (*HLA-A*02:01*, *HLA-*

282    *A*03:01*, and *HLA-A*02:02*), we optimized 14 similarity matrices capable of ranking

283    peptides by their binding affinities via their total similarity to strong binders. Matrices were

284    optimized by iteratively perturbing a seed similarity matrix and keeping those changes

285    which ultimately increased the self-similarity of the strong binders and cross-similarity of

286    the strong with weak binders.
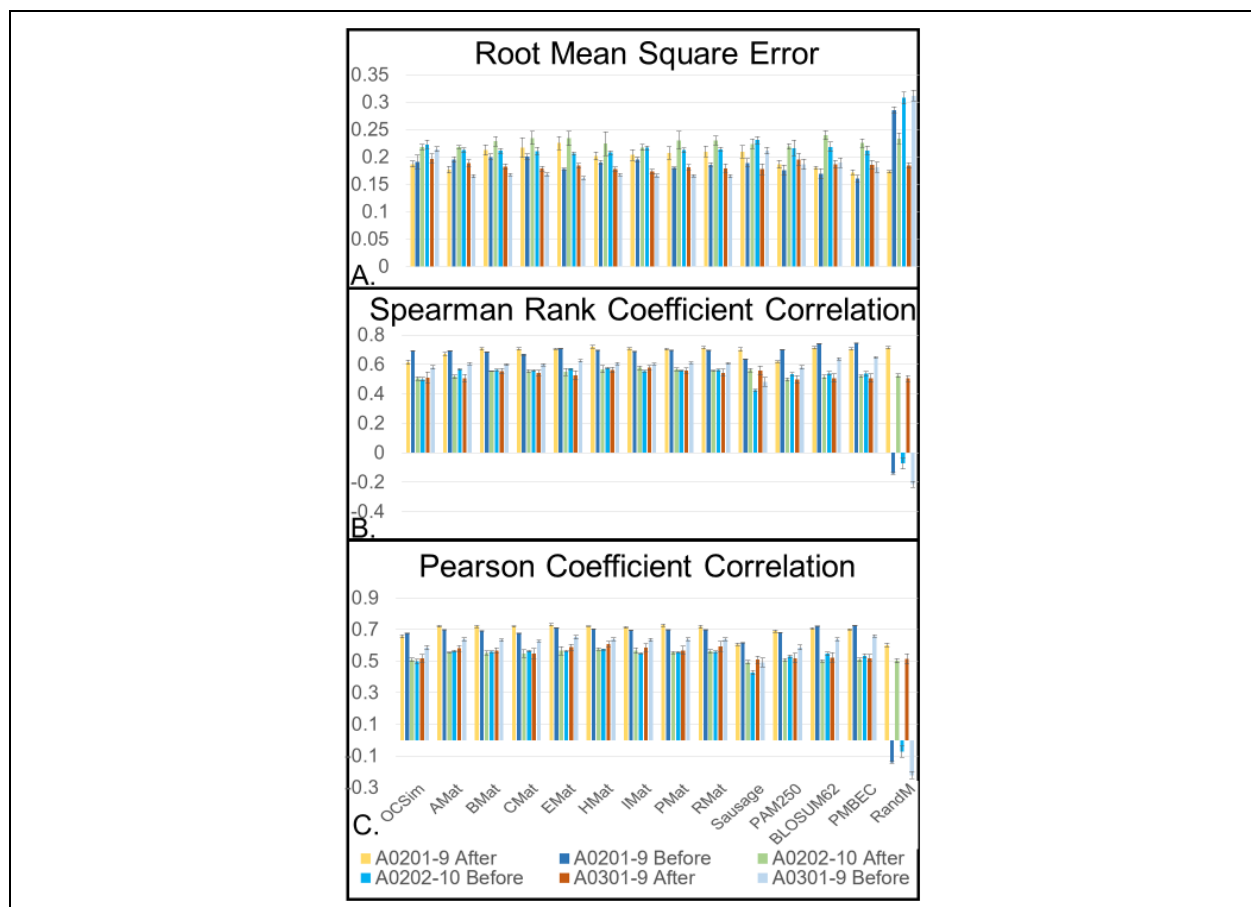
## Results and Discussions

288        **Cross-Similarity Analysis.** Fig 3 shows the cross-similarity results of 5 subsets

289    of peptides deemed Strong (S; -pIC$_{50}$:0), Less Strong (LS; -pIC$_{50}$:-1 to -2), Medium (M; -

290    pIC$_{50}$:-2 to -3), Medium Weak (MW; -pIC$_{50}$:-3 to -4) and Weak (W; -pIC$_{50}$:-5 to -7) based

291    on their binding affinity to alleles MHC-I *HLA-A*02:01* for all matrices before and after

292    training. Each set of bars per matrix was normalized by the largest value of both before

293    and after results. In addition, these bars are the results of 5 average TSS subsets (80%

294    randomly chosen from each affinity class). Previous work showed the TSS of a peptide

295    with high similarity to the peptides that are strong binders of a solid-state material

296    indicates that the peptide in question likely also has strong binding capability [25].

297    Therefore, the average TSS of peptides with an affinity for a protein should decrease with

298    their experimental affinity. Fig 3 shows that before training (yellow bars) the trend is

299    somewhat present but not very defined (Strong x Weak is comparable to Strong x Less-

300    Weak) but after training (blue bars) the trend is very pronounced. For each matrix and

301    across all three alleles (see S1 and S2 Figs) we observe average TSS when grouped by

302    affinity class to strong binders correlated with experimental affinity. Most notably, the

303    randomly initialized matrix RandM despite having no initial correlation was able to show

304    the trend as definitively as the others after matrix training.

**Fig 3. External similarity results for matrices before and after training with *HLA-A 02:01* binders.** Five subsets of peptides were created from the full list from each allele. Blue bars represent after training and yellow before training. The TSS for each group to strong binders ($TSS_{S\text{-}LS,M,MW,W}$) was calculated in addition to each group to itself ($TSS_{S\text{-}S}$, $TSS_{LS\text{-}LS}$, $TSS_{M\text{-}M}$, $TSS_{MW\text{-}MW}$, $TSS_{W\text{-}W}$). The y-axis for each bar chart denotes the matrix, the x-axis is the normalized $TSS_{S\text{-}S,LS,M,MW,W}$ values. The results show, especially in RandM's case, that we can improve similarity matrices to predict a trend correlated to binding affinity. This trend is characterized by decreasing $TSS_{S\text{-}S,LS,M,MW,W}$ correlating with decreasing binding affinity.

**Correlations with experimental affinity.** In the previous work, binding affinity was predicted by placing peptides into semi-quantitative groups of strong, medium and weak by their total similarity score to the strong binding peptide sequences of quartz [25]. The trend of decreasing $TSS_{Seq\text{-}S}$ was correlated with experimental affinity by using $TSS_{Seq\text{-}S}$

317   as a threshold to determine whether a peptide would fall into an affinity class (binary

318   classification) [25]. Though significant predictability (70-80%) was obtained using the

319   semi-quantitative scoring method, it falls short of the trend prediction needed to be

320   comparable with MHCFlurry, NetMHC and DeepMHC [12,13,37]. To enable more direct

321   comparisons the Pearson correlation coefficient (linear, Fig 4C) and Spearman rank

322   correlation coefficient (nonlinear, Fig 4B) were calculated, which can determine whether

323   the predicted binding affinity trend (TSS to strong binders) matches the experimental

324   binding affinity trend. In addition, a classifier scheme is included that can recognize

325   whether a peptide is a strong or weak binder by the magnitude of its $TSS_{Seq-S}$. Further, a

326   root mean square error (RMSE) is calculated from the normalized trend of TSS and RMSE

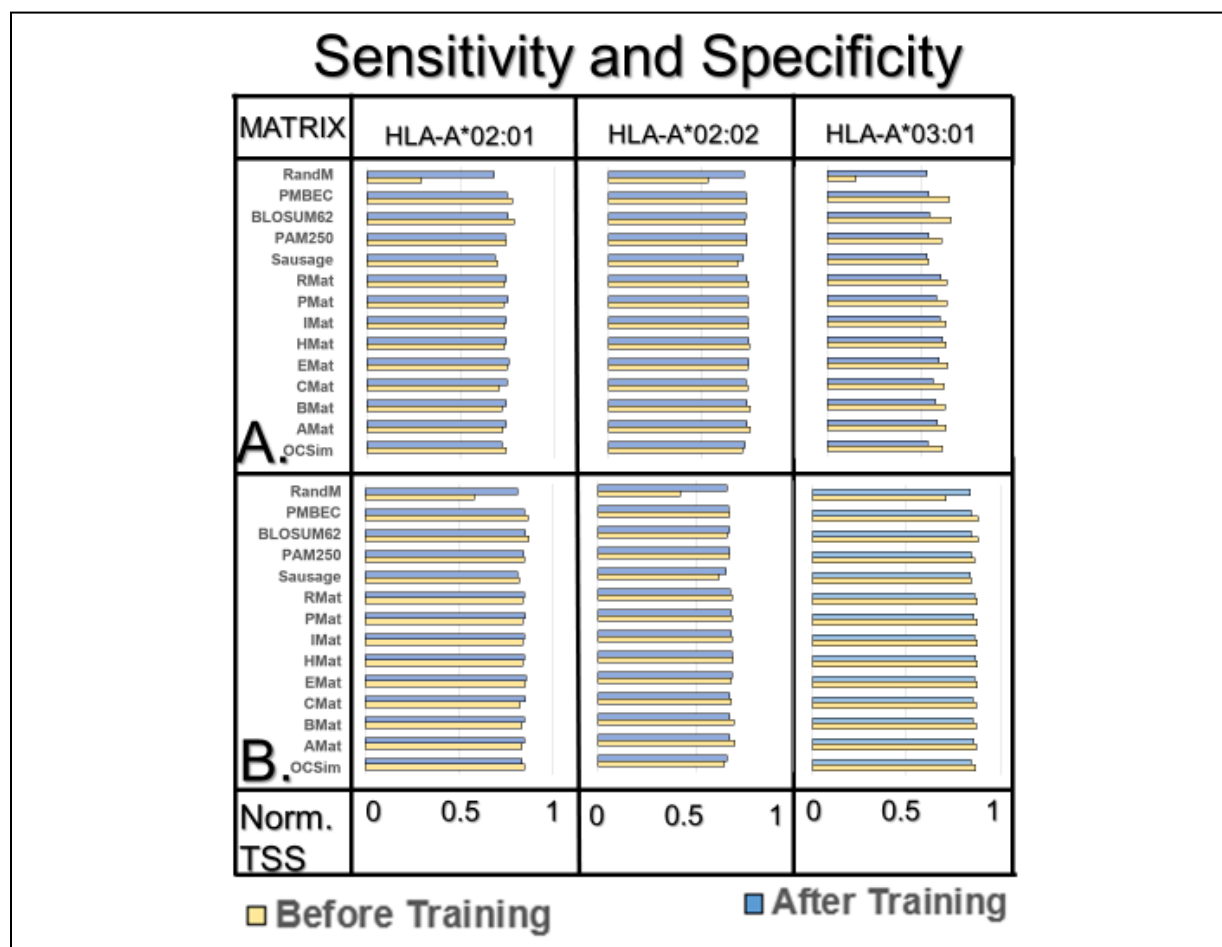327   to get an idea of close the TSS scores are to the experimental affinity (Fig 4A).

328    **Fig 4. RMSE (A), Spearman Rank (B), and Pearson (C) correlations for TSS trends**

329    **calculated using trained and untrained matrices.** The $TSS_{Seq-S}$ of the *HLA-A\*02:01* list was

330    calculated by aligning each peptide with the top binders of the allele and correlating the list of

331    values with the list of experimentally determined binding affinities using linear (Pearson, Fig 4C)

332    and nonlinear (Spearman Rank, Fig 4B) methods. RMSE (Fig 4A) is calculated by obtaining the

333    root mean square of the difference between the normalized (0-1) $-pIC_{50}$ and $TSS_{Seq-S}$. Error bars

334    are 1 standard deviation from the average of each set. The data shows the method can improve

335    literature and calculated matrices, most significantly that a trained randomly initialized matrix

336    (RandM) is more reflective of mutability information in the MHC-I context than all literature and

337    calculated matrices before training.

338        $TSS_{Seq-S}$ were then correlated with those of the experimental affinity. The values

339    of $TSS_{Seq-S}$ served as the predicted binding affinity ranking and was correlated with the

340    experimentally determined binding affinity using Pearson and Spearman Rank functions.

341    Fig 4 shows the score of each correlation for trained matrices and untrained matrices.

342    The error bars are one standard deviation from the average of these scores. All p-values

343    were less than 0.0005, except for in the case of RandM. Considering the substantially

344    less amount of data used (~350 peptide sequences for *HLA-A\*02:01*) compared with

345    DeepMHC and NetMHC (80% of the full set;~7200 sequences for *HLA-A\*02:01*), the

346    range of 0.5-0.7 is significant and is reflective of mutability information being captured. In

347    addition, the RMSE scores show that in general one TSS score is insufficient to describe

348    the exact binding affinity. While it is clear from the improvement in Pearson and Spearman

349    Rank correlation that these matrices are capturing some similarity information using the

350    method, no matrix alone can produce a TSS ranking exactly correlating with the rest of

351    the set. The integration of several TSS rankings into a single score could prove to be a

352    relevant predictor if they are capturing diverse similarity information unique to their matrix

353    values.

354        **Binary Classification.** Sensitivity and specificity were also recorded as a measure

355    of binary prediction accuracy, shown in Fig 5A and Fig 5B respectively. A

356    binder/nonbinder classification was performed via observing peptides conserved as

357    binders through the magnitude of their $TSS_{S\text{-}Seq}$. The sequences having greater than 500

358    $IC_{50}$ [12] were considered binders. Therefore, peptides given a predicted ranking above

359    the $TSS_{S\text{-}Seq}$ threshold correlating with the 500 $IC_{50}$ [12] bar were considered predicted

360    binders. True positives and negatives, and false positive and negatives were calculated

361    by observing which predicted binders were also in the actual binder group.

362  **Fig 5. Results of classification analysis.** Sensitivity (A) and specificity (B) measures

363  calculated from the results of binary classification of binding. In general, these results show the

364  training function did not improve the predictive ability of any matrix besides RandM, providing

365  evidence that TSS is a relevant predictor while noting the training operation is an ineffective

366  application of TSS.

367  Across all the matrices a similar specificity/sensitivity was observed before and

368  after training. This indicates the cost function did not improve the ability of the

369  calculated/literature matrices to classify peptides based on TSS values. RandM showed

370  marked improvement across all the alleles but yields lower accuracies than other

371  matrices. This demonstrates that information can be integrated into a similarity matrix up

372  to a limit. In general, the prediction metrics show that the separation of $TSS_{S-S}$ and $TSS_{S-W}$

373  $_W$ may not be the appropriate cost function to improve a predictive model. However,

374  $TSS_{Seq-S}$ is a highly relevant predictor of affinity. Though the model was trained on only

375  the dominant features of the peptide set represented by strong binders, the affinity trend

376  was generally conserved by $TSS_{S-Seq}$ scoring.

## Conclusions and Future Work

378  The predicted correlation range of 0.5-0.7 determined by Pearson and Spearman

379  Rank of the similarity matrix methodology demonstrates similarity matrices can predict

380  functionality (i.e. solid substrate binding specificity) of peptides using the Total Similarity

381  Score. Previous work provided definitive evidence concluding the average similarity score

382  (TSS) of a peptide towards strong binding peptides of an inorganic solid material is

383  positively correlated with the binding affinity of that peptide. Using the Total Similarity

384  Score, we modified a computational method and applied it to a substantially larger dataset

385    to demonstrate that across organic and inorganic materials the metric applies. Though

386    we use substantially lower training data than other methods, similarity matrices were

387    obtained that recognize the dominant features of the strongest binding peptides, which in

388    turn describe those of the weaker binders. Therefore, the strongest binders of the full set

389    can adequately describe the behavior of the remaining peptides Though the training

390    method is insufficient to produce a trend capable of ranking affinity with comparable

391    accuracy to other MHC predictors, we postulate that based on the diversity of the matrices

392    trained that they are capturing different subsections of the total similarity information.

393    Therefore, integrating the trends of multiple matrices into a single score would produce

394    comparable accuracy even when trained on substantially less data. In this work, we show

395    that we can capture similarity information using different matrices and that TSS to strong

396    binders is a relevant predictor of affinity in both organic and inorganic systems.

397        To uncover the relationship between $TSS_{Seq-S}$ and the experimentally measured

398    affinity, the future work would involve integrating the TSS score with recent statistical

399    learning techniques. If the matrix cannot be optimized, then the value of $TSS_{Seq-S}$ may not

400    be the highest achievable even if the sequence is a strong binder. The sequences with

401    amino acids in similar positions to the strong binding group will, however, tend to give the

402    same average score. Therefore, if the goal is to predict the similarity of sequences based

403    on their positional composition, conserving the common score range will also retain their

404    sequence information. An additional problem may also arise when considering the

405    diversity of the strong binding group. If a given peptide is a strong binder having a

406    completely unique sequence compared to those of the other strong-binding peptides, it

407    will have a low $TSS_{Seq-S}$. TSS scoring assumes that weak and medium binders are

408 mutations of stronger binders. Future methods will capitalize on the information hidden

409 within weak/medium binders and use it to describe the full strong binding space. The full

410 results, gapless Iterative Alignment Python program for calculating similarity matrices,

411 and all the data used to train the matrices are located online on GitHub

412 (https://github.com/Sarikaya-Lab-GEMSEC/Iterative-Alignment-Gapless).

## Acknowledgements

414 We appreciate the data sets and computational facilities provided GEMSEC labs at the

415 University of Washington.

416

417 We declare there were no conflicts of interest for this work.

## References

419 1. Georgoulia PS, Glykos NM. Molecular simulation of peptides coming of age:

420 Accurate prediction of folding, dynamics and structures. Archives of Biochemistry

421 and Biophysics. 2019;664:76-88.

422 2. Minami S, Sawada K, Chikenji G. MICAN: a protein structure alignment

423 algorithm that can handle Multiple-chains, Inverse alignments, Cα only models,

424 Alternative alignments, and Non-sequential alignments. BMC Bioinformatics.

425 2013;14(1):24.

426 3. Muthukrishnan S, Puri M. Harnessing the evolutionary information on oxygen

427 binding proteins through Support Vector Machines based modules. BMC

428 Research Notes. 2018;11(1).

429    4. Nakai K, Kidera A, Kanehisa M. Cluster analysis of amino acid indices for

430    prediction of protein structure and function. "Protein Engineering, Design and

431    Selection". 1988;2(2):93-100.

432    5. Khatayevich D, Page T, Gresswell C, Hayamizu Y, Grady W, and Sarikaya M,

433    Selective Detection of Target Proteins by Peptide-Enabled Graphene Biosensor,

434    Small, 2014; 10(8): 1505-1513.

435    6. Hayamizu Y, So CR, Dag S, Page TS, Starkebaum D, Sarikaya M. Bioelectronic

436    interfaces by spontaneously organized peptides on 2D atomic single layer

437    materials. Scientific Reports. 2016 Sep 22;6:33778.

438    7. Fosgerau K, Hoffmann T. Peptide therapeutics: current status and future

439    directions. Drug Discovery Today. 2015;20(1):122-128.

440    8. Barbosa AJM, Oliveira AR, Roque ACA. Protein- and Peptide-Based Biosensors

441    in    Artificial    Olfaction.    *Trends    Biotechnol.*    2018;36(12):1244–1258.

442    doi:10.1016/j.tibtech.2018.07.004

443    9. Hughes ZE, Walsh TR. What makes a good graphene-binding peptide?

444    Adsorption of amino acids and peptides at aqueous graphene interfaces. Journal

445    of Materials Chemistry B. 2015;3(16):3211-21.

446    10. Sarikaya M, Tamerler C, Jen AK-Y, Schulten K, Baneyx F. Molecular

447    biomimetics: nanotechnology through biology. Nature Materials. 2003;2(9):577–

448    85.

449    11. Jurtz V, Paul S, Andreatta M, Marcatili P, Peters B, Nielsen M. NetMHCpan-

450    4.0: Improved Peptide–MHC Class I Interaction Predictions Integrating Eluted

Ligand and Peptide Binding Affinity Data. The Journal of Immunology. 2017;199(9):3360-8.

12. Hu J, Liu Z. DeepMHC: Deep Convolutional Neural Networks for High-performance peptide-MHC Binding Affinity Prediction. Cold Spring Harbor Laboratory; 2017.

13. Odonnell TJ, Rubinsteyn A, Bonsack M, Riemer AB, Laserson U, Hammerbacher J. MHCflurry: Open-Source Class I MHC Binding Affinity Prediction. Cell Systems. 2018;7(1).

14. Vita R, Mahajan S, Overton JA, Dhanda SK, Martini S, Cantrell JR, et al. The Immune Epitope Database (IEDB): 2018 update. Nucleic Acids Research. 2018;47(D1).

15. Wang S, Ma J, Peng J, Xu J. Protein structure alignment beyond spatial proximity. Scientific Reports. 2013;3(1).

16. Wright ES. DECIPHER: harnessing local sequence context to improve protein multiple sequence alignment. BMC Bioinformatics. 2015;16(1).

17. Dosztanyi Z, Torda AE. Amino acid similarity matrices based on force fields. Bioinformatics. 2001;17(8):686-99.

18. Dasetty S, Barrows JK, Sarupria S. Adsorption of Amino Acids on Graphene: Assessment of Current Force Fields. American Chemical Society (ACS); 2019.

19. Alford RF, Leaver-Fay A, Jeliazkov JR, O'Meara MJ, DiMaio FP, Park H, et al. The Rosetta all-atom energy function for macromolecular modeling and design. J Chem Theory Comput. 2017;13(6):3031–48.

473    20. Needleman SB, Wunsch CD. A general method applicable to the search for

474    similarities in the amino acid sequence of two proteins. Journal of Molecular

475    Biology. 1970;48(3):443-53.

476    21. Jones DT, Taylor WR, Thornton JM. The rapid generation of mutation data

477    matrices from protein sequences. Bioinformatics. 1992;8(3):275-82.

478    22. Brown P, Pullan W, Yang Y, Zhou Y. Fast and accurate non-sequential protein

479    structure alignment using a new asymmetric linear sum assignment heuristic.

480    Bioinformatics. 2015;32(3):370-7.

481    23. Deorowicz S, Debudaj-Grabysz A, Gudyś A. FAMSA: Fast and accurate

482    multiple sequence alignment of huge protein families. Scientific Reports.

483    2016;6(1).

484    24. Thompson JD, Higgins DG, Gibson TJ. CLUSTAL W: improving the sensitivity

485    of progressive multiple sequence alignment through sequence weighting, position-

486    specific gap penalties and weight matrix choice. Nucleic Acids Research.

487    1994;22(22):4673–80.

488    25. Oren EE, Tamerler C, Sahin D, Hnilova M, Seker UOS, Sarikaya M, et al. A

489    novel knowledge-based approach to design inorganic-binding peptides.

490    Bioinformatics. 2007;23(21):2816-22.

491    26. Keul F, Hess M, Goesele M, Hamacher K. PFASUM: a substitution matrix from

492    Pfam structural alignments. BMC Bioinformatics. 2017;18(1).

493    27. Kim Y, Sidney J, Pinilla C, Sette A, Peters B. Derivation of an amino acid

494    similarity matrix for peptide:MHC binding and its application as a Bayesian prior.

495    BMC Bioinformatics. 2009;10(1):394.

496  28. Andreatta M, Nielsen M. Gapped sequence alignment using artificial neural

497  networks: application to the MHC class I system. Bioinformatics. 2015;32(4):511-

498  7

499  29. Han Y, Kim D. Deep convolutional neural networks for pan-specific peptide-

500  MHC class I binding prediction. BMC Bioinformatics. 2017;18(1).

501  30. Liang G, Yang L, Chen Z, Mei H, Shu M, Li Z. A set of new amino acid

502  descriptors applied in prediction of MHC class I binding peptides. European

503  Journal of Medicinal Chemistry. 2009;44(3):1144-54.

504  31. Hess M, Keul F, Goesele M, Hamacher K. Addressing inaccuracies in

505  BLOSUM computation improves homology search performance. BMC

506  Bioinformatics. 2016;17(1).

507  32. Kawashima S. AAindex: Amino Acid index database. Nucleic Acids Research.

508  2000;28(1):374-.

509  33. Kawashima S, Pokarowski P, Pokarowska M, Kolinski A, Katayama T,

510  Kanehisa M. AAindex: amino acid index database, progress report 2008. Nucleic

511  Acids Research. 2007;36(Database):D202-D5.

512  34. Hemmateenejad B, Miri R, Elyasi M. A segmented principal component

513  analysis—regression approach to QSAR study of peptides. Journal of Theoretical

514  Biology. 2012;305:37-44.

515  35. Atchley W, Zhao J, Fernandes A, Druke T. Solving the protein sequence metric

516  problem. Proceedings of the National Academy of Sciences. 2005;102(18):6395-

517  6400.

518    36. Doytchinova IA, Walshe V, Borrow P, Flower DR. Towards the chemometric

519    dissection of peptide – HLA-A*0201 binding affinity: comparison of local and global

520    QSAR models. Journal of Computer-Aided Molecular Design. 2005;19(3):203-12.

521    37. Dogan S, Fong H, Yucesoy DT, Cousin T, Gresswell C, Dag S, Huang G,

522    Sarikaya M. Biomimetic tooth repair: amelogenin-derived peptide enables in vitro

523    remineralization of human enamel. ACS Biomaterials Science & Engineering.

524    2018 Mar 9;4(5):1788-96.

525    38. Saha I, Maulik U, Bandyopadhyay S, Plewczynski D. Fuzzy clustering of

526    physicochemical and biochemical properties of amino Acids. Amino Acids.

527    2011;43(2):583-594.

528    39. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O et. al.

529    Scikit-learn: Machine learning in Python. Journal of machine learning research.

530    2011;12(Oct):2825-30.

531

532

533

534

535

536

537

538

539

540

541

542

## Supporting Information

**S1 Fig. Cross-similarity results for the HLA-02:02 allele.** Each bar chart shows the average normalized TSS of the "Strong" affinity class with itself and each other class. The decreasing trend similarity of "Strong" peptides with those of decreasing affinity demonstrates the successful optimization of each matrix for the HLA-A*02:02 allele.
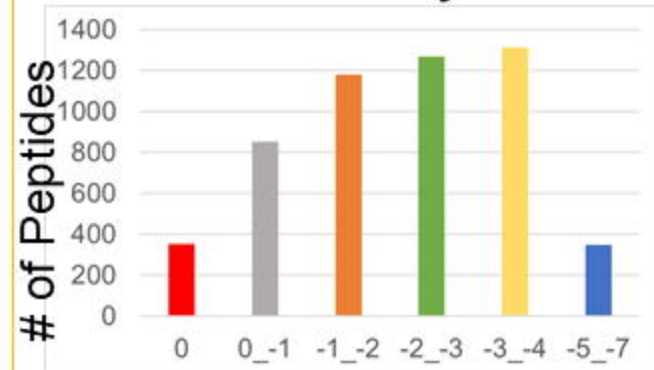
**S2 Fig. Cross-similarity results for the HLA-03:01 allele.** Each bar chart shows the average normalized TSS of the "Strong" affinity class with itself and each other class. The decreasing trend similarity of "Strong" peptides with those of decreasing affinity demonstrates the successful optimization of each matrix for the HLA-A*03:01 allele.

# Sensitivity and Specificity

| MATRIX | HLA-A*02:01 | HLA-A*02:02 | HLA-A*03:01 |
|---|---|---|---|

**A.**

| | HLA-A*02:01 | HLA-A*02:02 | HLA-A*03:01 |
|---|---|---|---|
| RandM | | | |
| PMBEC | | | |
| BLOSUM62 | | | |
| PAM250 | | | |
| Sausage | | | |
| RMat | | | |
| PMat | | | |
| IMat | | | |
| HMat | | | |
| EMat | | | |
| CMat | | | |
| AMat | | | |
| OCSim | | | |

**B.**

| | HLA-A*02:01 | HLA-A*02:02 | HLA-A*03:01 |
|---|---|---|---|
| RandM | | | |
| PMBEC | | | |
| BLOSUM62 | | | |
| PAM250 | | | |
| Sausage | | | |
| RMat | | | |
| PMat | | | |
| IMat | | | |
| HMat | | | |
| EMat | | | |
| CMat | | | |
| BMat | | | |
| AMat | | | |
| OCSim | | | |

| Norm. TSS | 0    0.5    1 | 0    0.5    1 | 0    0.5    1 |
|---|---|---|---|

■ Before Training          ■ After Training

A. Data Collection

MHC-1 Affinity Class

-pIC$_{50}$ Class

BLOSUM 62; Pam 250

B. Matrix Training

Trendless Matrix

Trendy Matrix

C. Similarity Calculation

Novel Similarity Matrices (14)

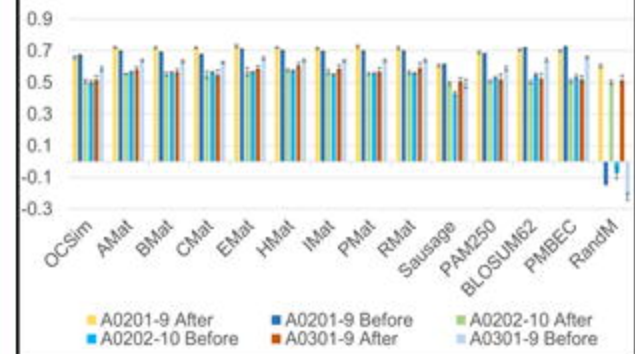TSS Distribution

D. Predictive Functionality

Pearson Coefficient Correlation

A0201-9 After    A0201-9 Before    A0202-10 After
A0202-10 Before    A0301-9 After    A0301-9 Before

Sensitivity

STEP 1. DATA IN

STEP 2. Affinity Classification

Matrix Training

STEP 3. Matrix Perturbation + Calculate $TSS_{S-S}$, $TSS_{S-W}$

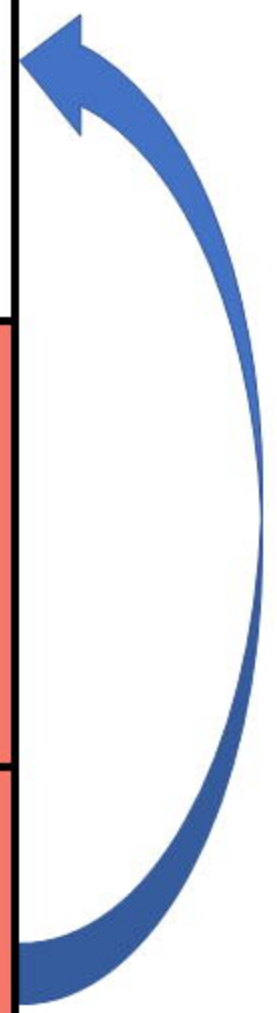STEP 4. Higher $TSS_{S-S}$, Lower $TSS_{S-W}$ | Lower $TSS_{S-S}$, Higher $TSS_{S-W}$
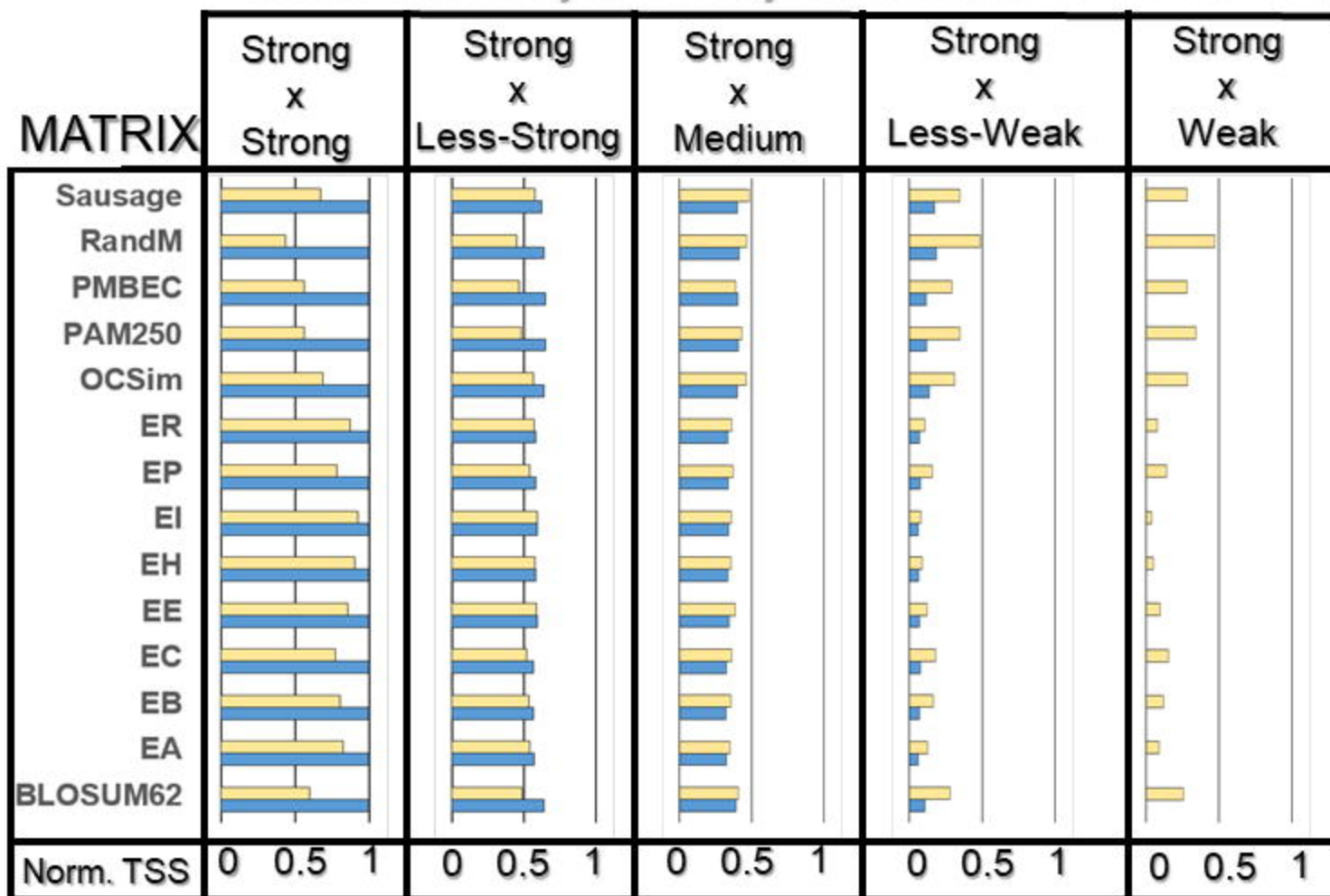
STEP 5. Save Matrix Changes | Revert Matrix Changes

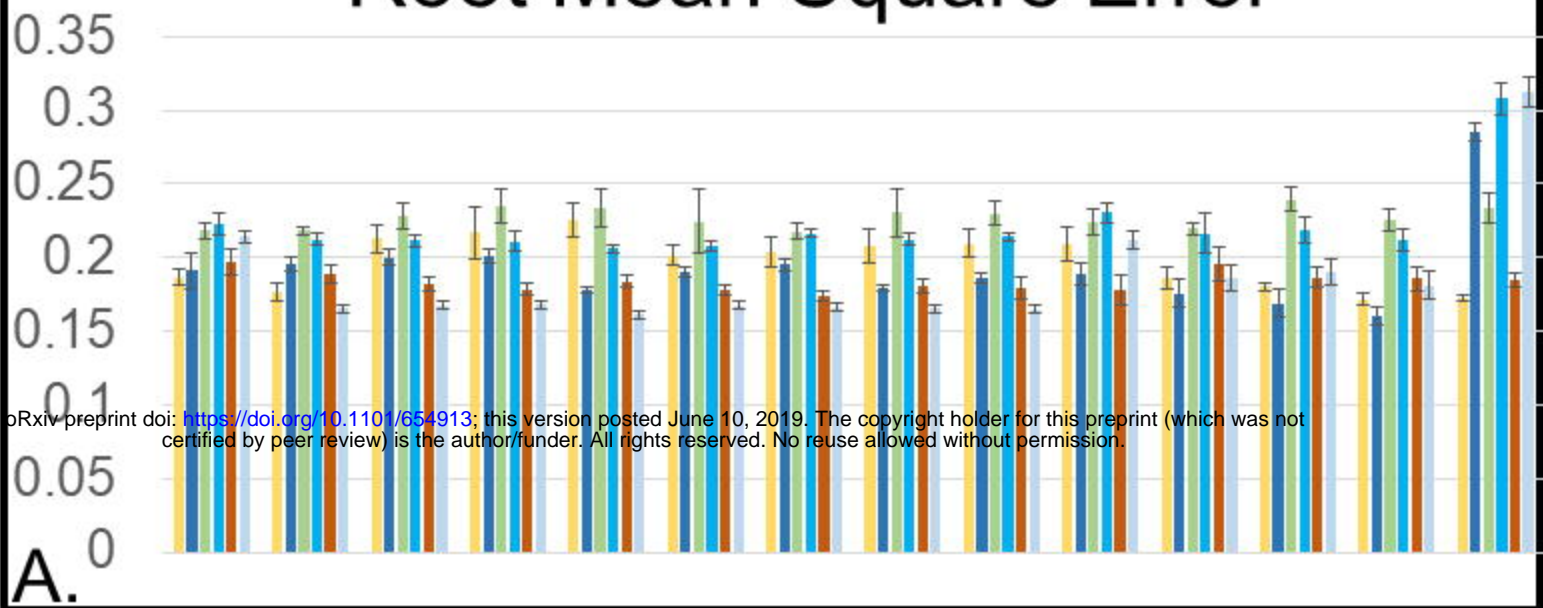STEP 6. Return to Step 3, converge after 5000 saved changes

Trained Matrix

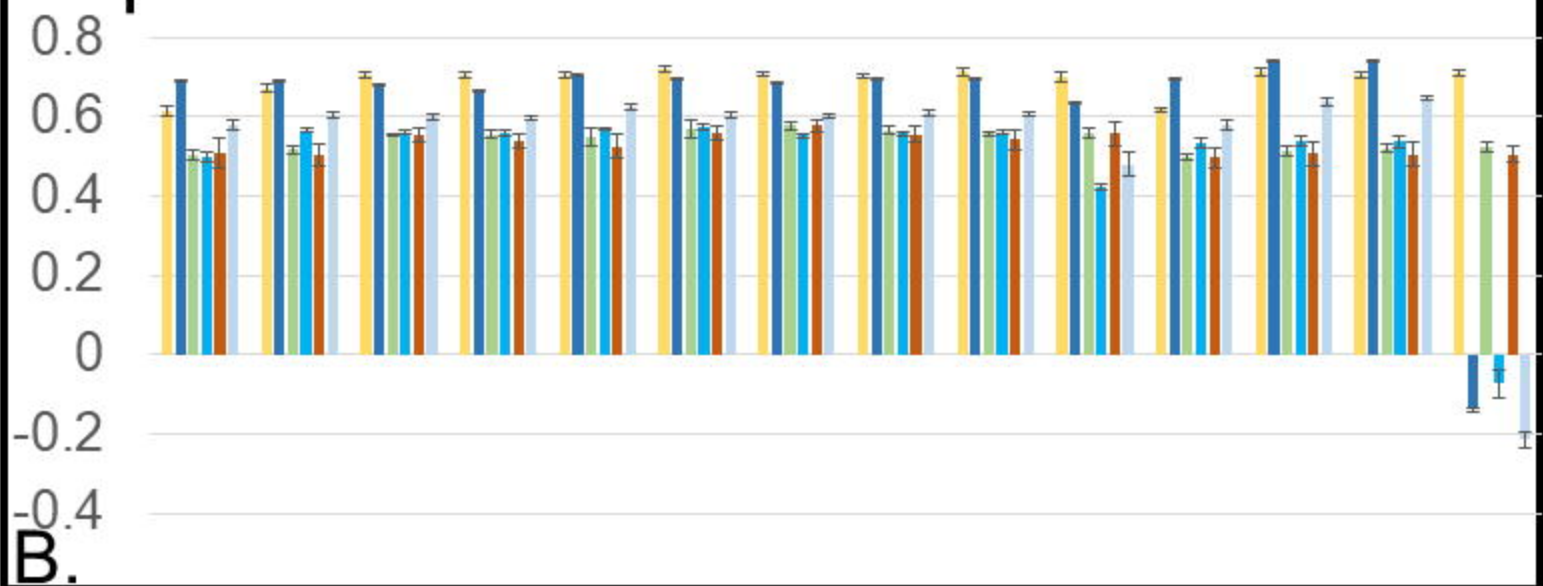Cross-Similarity of Affinity Classes for HLA-A 02:01

# Root Mean Square Error

**A.**

# Spearman Rank Coefficient Correlation

**B.**

# Pearson Coefficient Correlation

OCSim, AMat, BMat, CMat, EMat, HMat, IMat, PMat, RMat, Sausage, PAM250, BLOSUM62, PMBEC, RandM

**C.**

- ■ A0201-9 After
- ■ A0201-9 Before
- ■ A0202-10 After
- ■ A0202-10 Before
- ■ A0301-9 After
- ■ A0301-9 Before