

# TEINet: a deep learning framework for prediction of TCR-epitope binding specificity

Yuepeng Jiang<sup>1</sup>, Miaoze Huo<sup>1</sup>, and Shuai Cheng Li<sup>1\*</sup>

<sup>1</sup>Department of Computer Science, City University of Hong Kong, Hongkong, China.

\*Corresponding author. Email: [shuaicli@cityu.edu.hk](mailto:shuaicli@cityu.edu.hk)

## Abstract

The adaptive immune response to foreign antigens is initiated by T-cell receptor (TCR) recognition on the antigens. Recent experimental advances have enabled the generation of a large amount of TCR data and their cognate antigenic targets, allowing machine learning models to predict the binding specificity of TCRs. In this work, we present TEINet, a deep learning framework that utilizes transfer learning to address this prediction problem. TEINet employs two separately trained encoders to transform TCR and epitope sequences into numerical vectors, which are subsequently fed into a fully connected neural network to predict their binding specificities. A major challenge for binding specificity prediction is the lack of a unified approach to sample negative data. Here, we first assess the current negative sampling approaches comprehensively and suggest that the *Unified Epitope* is the most suitable one. Subsequently, we compare TEINet with three baseline methods and observe that TEINet achieves an AUROC of 0.760, which outperforms baseline methods by 6.4-26%. Furthermore, we investigate the impacts of the pretraining step and notice that excessive pretraining can adversely affect model performance. Our results and analysis show that TEINet can make an accurate prediction using only the TCR sequence (CDR3 $\beta$ ) and the epitope sequence, providing novel insights to understand the interactions between TCRs and epitopes. TEINet is available at <https://github.com/jiangdada1221/TEINet>.

## Introduction

T cells are critical for the adaptive immune system, providing protection against a wide range of pathogens. To recruit T cells in an immune response, the T cell receptors (TCRs) on their surface have to recognize a non-self-immunogenic peptide (epitope) presented in the context of major histocompatibility complex molecules (MHC). The generation of these protein receptors arises mainly from the quasirandom somatic V(D)J recombination process which theoretically can produce extremely high TCR diversity of  $10^{15}$ - $10^{20}$  in an individual, each with unique recognition capacity for antigens [1]. Understanding the mechanisms that govern the interaction between TCR and

32 peptide-MHC (pMHC) is considered an essential step toward personalized immunotherapy and the  
33 development of targeted vaccines.

34 Recent advancements in the high-throughput tetramer-associated T cell receptor sequencing  
35 technique [2] and other experimental approaches such as tetramer analysis [3] and T-scan [4] have  
36 enabled the generation of an increasing amount of data recording the binding of TCR and epitope.  
37 More and more interaction pairs are consistently being generated and stored in publicly available  
38 databases such as VDJdb [5], IEDB [6] and McPAS-TCR [7]. However, the available data are  
39 still scant compared to the theoretical TCR diversity. Further, the TCR-epitope paired data are  
40 imbalanced, as a single epitope is often linked by many TCRs. Both of them pose challenges to the  
41 development of *in silico* predictive methods.

42 Machine learning-based methods are able to capture the potential laws of TCR-epitope binding  
43 from a large amount of experimental data. With the help of advanced machine learning models,  
44 several computational methods have been proposed to assess the binding of a TCR and a pMHC  
45 (epitope). Previously, a branch of research focused on designing epitope-specific models with the  
46 aim of learning the pattern of TCRs binding to the same epitope. These models range from simple  
47 sequence alignment-based methods [8] to more complex machine learning models including random  
48 forest (e.g. TCRex [9]) and the Gaussian process classifier TCRGP [10]. However, they all share two  
49 downsides: each epitope needs a specific model trained separately; each model requires abundant  
50 training samples of epitope-specific TCRs, which are not always readily available.

51 To fulfill the need to predict the binding specificity of any TCR-epitope pair, previous studies  
52 have proposed generic models, which exploit the two-tower architecture to encode both the TCRs  
53 (CDR3 $\beta$ ) and the pMHCs (epitope) [11–17]. These generic models can fully capitalize on the cur-  
54 rently available paired data to unlock the binding patterns between TCRs and epitopes, and transfer  
55 the knowledge learned from paired samples of epitopes with sufficient binding TCRs to those with  
56 sparse linking TCRs. Current models have shown moderate predictive performance and demon-  
57 strated promising potential in understanding cancer progression, prognosis, and responsiveness to  
58 immunotherapy. For example, Dash *et al.* developed TCRdist [17] based on sequence similarity  
59 weighted distances; Moris *et al.* proposed a CNN-based model ImRex [11]; Weber *et al.* introduced  
60 TITAN that encodes epitopes at the atomic level with SMILES sequences using a pretrained deep  
61 learning model. Furthermore, Lu *et al.* presented pMTNet [14] that encodes TCRs and pMHCs by  
62 two respective pre-trained deep learning models and applied pMTNet to investigate tumor progres-  
63 sion and response to immunotherapy treatment. In particular, transfer learning is becoming a prior  
64 technique to develop advanced deep learning models for binding prediction since it helps leverage the  
65 knowledge from other pretraining tasks with abundant data and transfer it to the binding predic-  
66 tion task. For instance, TITAN, NetTCR [15], and pMTnet utilize pre-trained encoders. However,  
67 the impact of the pretraining step on the final performance of predicting TCR specificity remains  
68 undiscovered.

69 To train and evaluate supervised models, both positive and negative samples (TCRs and epitopes  
70 that do not interact with each other) are required. However, the public TCR-epitope interaction  
71 datasets only collect positive samples, which potentially poses a challenge in model training and  
72 evaluation. The method of generating negative samples based on the existing TCR and epitope

73 pairs directly affects the model performance. Currently, there are four major strategies for generating  
74 negative samples: (1) *Reference TCR* [15, 18]; (2) *Random TCR* [14]; (3) *Random Epitope* [12, 16,  
75 19]; (4) *Unified Epitope* [11, 20]. Different models might adopt different negative sampling strategies  
76 for this task, which make it difficult to fairly compare their performance. More importantly, which  
77 strategy leads to a better generalized model has not been explored and remains an open question.

78 In this work, we present TEINet for the prediction of the specificity of TCR binding, using the  
79 CDR3 $\beta$  chain of TCR and the epitope sequence within the pMHC complex. Following the concept of  
80 transfer learning, TEINet employs two separate pretrained encoders to convert TCRs and epitopes  
81 into numerical vectors, utilizing the architecture of recurrent neural networks to handle a variety of  
82 sequence lengths. We first contrast the four negative sampling strategies applied in the previous work  
83 to select the superior one. Next, we systematically validated TEINet using a large-scale TCR-epitope  
84 paired dataset and two independent validation datasets. The results demonstrated the enhancement  
85 in accuracy made over previous work. We also investigated the impact of the pretraining step on  
86 the final binding specificity prediction task. Overall, TEINet serves as a reliable computational tool  
87 for addressing the long-standing problem of predicting the TCR-epitope interaction.

## 88 Methods

### 89 Dataset

90 The CDR3 regions of TCR $\beta$  chains are located in the center of the paratope and are considered as the  
91 key determinant of specificity in antigen recognition [21]. Although CDR3- $\alpha$  and - $\beta$  synergistically  
92 drive TCR-epitope recognition [17, 22, 23], the current available databases still record mostly  $\beta$ chain  
93 paired samples. Thus, we restrict ourselves to CDR3 $\beta$  chain sequences in this study. Besides, with  
94 the aim of developing a general model that is suitable for most cases, we took the epitope sequence  
95 inside the pMHC complex as its representation. In order to construct a large and diverse dataset,  
96 we combined the data recorded in VDJdb database [5], McPAS database [7], and the data collected  
97 by Lu *et al.* [14] together.

98 The data from VDJdb was downloaded from its public website (<https://vdjdb.cdr3.net/>) on  
99 April 5, 2022. It consists of 89,321 curated pairs of CDR3  $\alpha/\beta$  sequences along with their binding  
100 epitopes and MHC classes, covering three species. We selected only human TCR sequences, removed  
101 duplicate cases, restricted only MHC class I entries, and only kept the CDR3 $\beta$  and epitope sequences  
102 whose lengths lie between 5-30 and 7-15 amino acids, respectively. After all these filterings, this  
103 dataset was reduced to 35,560 unique CDR3 $\beta$ -epitope pairs, among which 33,258 TCRs are assigned  
104 to 159 epitopes.

105 The McPAS-TCR dataset [7] originally contains 39,664 pairs ([http://friedmanlab.weizmann.  
106 ac.il/McPAS-TCR/](http://friedmanlab.weizmann.ac.il/McPAS-TCR/)) and Lu *et al.* collected a total of 32,607 paring data from a series of previous  
107 publications and four chromium single-cell immune profiling solution datasets. We performed the  
108 same preprocessing step on these two datasets and removed all TCR sequences with ambiguous  
109 amino acids (B, J, O, U, X). Then, these three datasets were merged together, followed by two  
110 additional filtering steps: removal of duplicate pairs and exclusion of epitopes with less than 10

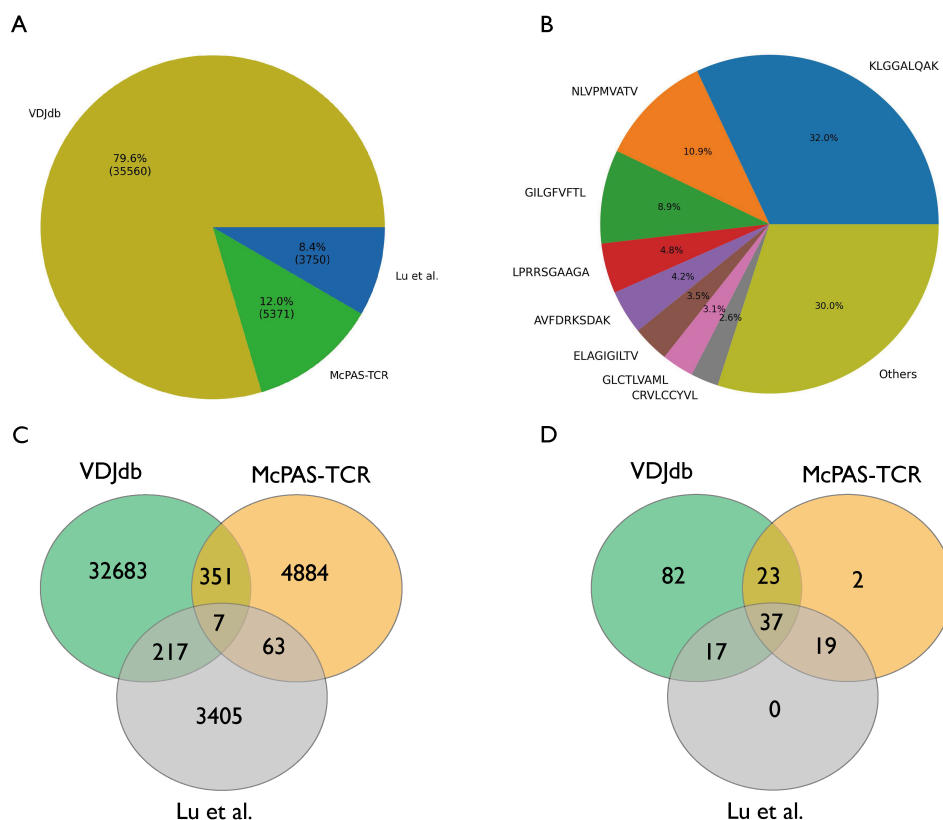


Figure 1: Overview of our constructed dataset. (A) The source of the paired samples in the dataset. (B) The number of the epitope-associated TCRs. Most TCRs are linked to a small group of epitopes. (C and D) Venn diagrams showing the number of (C) TCRs and (D) epitopes contributed by each source dataset.

111 associated TCR sequences since this merged dataset is highly imbalanced. At last, we constructed a  
 112 large dataset with 44,682 pairs of TCRs and epitopes, among which 41,610 TCRs are linked to 180  
 113 epitopes. An overview of the dataset is shown in Fig. 1.

## 114 Negative sampling strategies

115 Since the TCR-epitope dataset contains only positive samples, in order to train a generalized and  
 116 robust supervised model, the negative samples are required and should be generated via a biologically  
 117 and computationally plausible manner to serve as an unbiased estimate of the actual distribution  
 118 of non-binding pairs. For a positive sample  $d_i = (e_i, t_i) \in D = \{d_i\}_{i=1}^N$ , where  $e_i$  and  $t_i$  are the  
 119 interacting epitope and TCR for sample  $i$ , the corresponding negative samples are generated through  
 120 four major sampling strategies (Fig. 2A):

- 121 • *Reference TCR*. In this setting,  $e_i$  is combined with TCRs that are sampled uniformly from  
 122 the reference TCR dataset  $R = \{t_j\}$ . The negative samples for  $e_i$  are then represented as

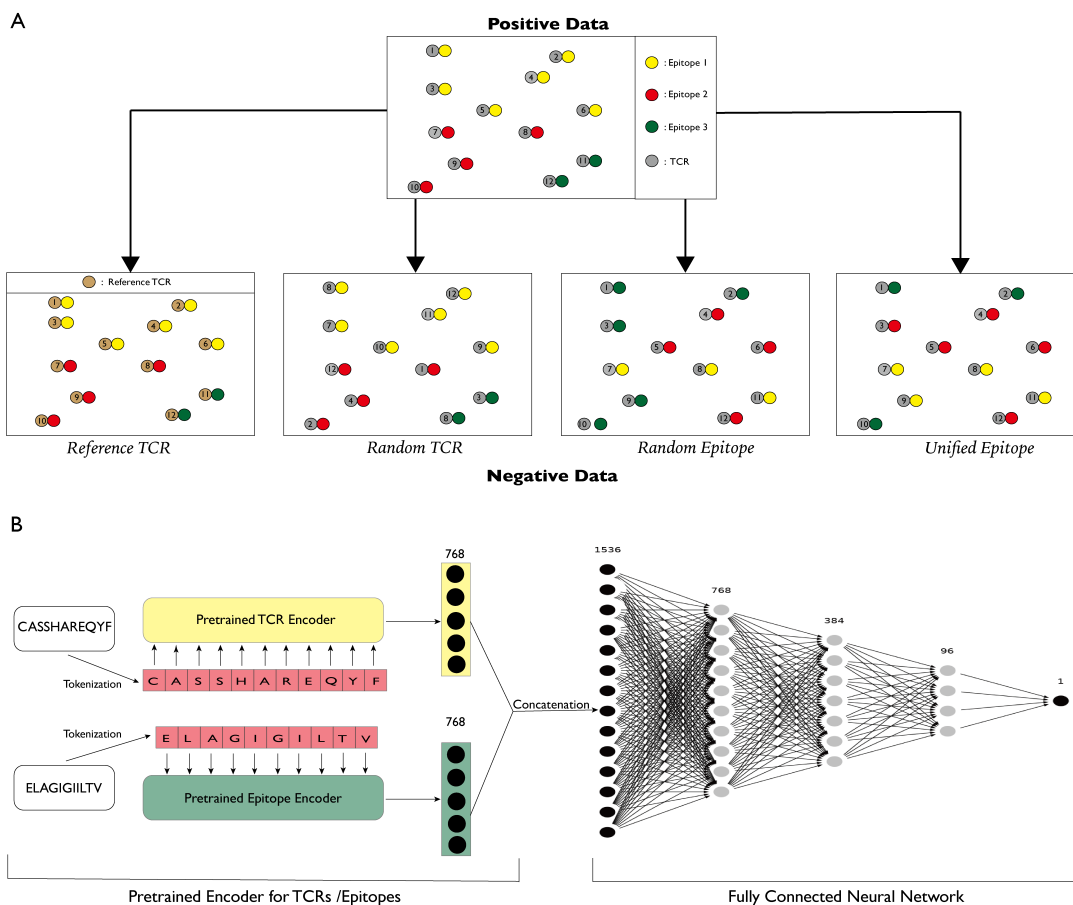


Figure 2: Illustration of each negative sampling strategy and the overall workflow of TEINet. (A) Sketch map of the four negative sampling strategies. In this example, there are in total 12 TCR-epitope binding pairs, with 3 different epitopes (depicted in yellow, red, and green) linking to 6, 4, and 2 TCRs, respectively. For *Reference TCR* strategy, the TCRs are randomly sampled from a reference TCR dataset inside which TCRs are considered unable to bind epitopes in the positive data. Here, we choose to generate the same number of pairs in negative data for demonstration. (B) General workflow of TEINet. TEINet is a two-stage deep learning model using transfer learning. At the first pretraining stage, two TCRpeg models are trained separately to learn the sequence pattern of TCRs and epitopes, and produce numerical encodings for them when the pretraining process is completed. At the next stage, encodings of TCRs and epitopes are concatenated together and output into a fully connected neural network to leverage the information from each part and make predictions accordingly.

123  $n_i = \{(e_j, t_j)\}_{j=1}^M$ , where  $t_j \in R$  and  $M$  is the number of negatives samples for a given positive  
 124 sample [15, 18]. This approach stands upon the assumption that TCRs from the reference  
 125 dataset are unlikely to bind epitopes in the positive dataset. The reference TCRs were obtained  
 126 from Montemurro *et al.* [15] where these TCRs had been exposed to all tested pMHC multimers  
 127 and no binding signals were detected.

128 • *Random TCR*. For this sampling approach, the negative TCRs for  $e_i$  are sampled uniformly

129 from the set of TCRs in the positive binding pairs while excluding its known true TCR binding  
130 partner(s) [14]. The negative samples are then represented as  $n_i = \{(e_i, t_k)\}_{k=1}^M$ , where  $t_k \in$   
131  $\{t_i\}$  and  $(e_i, t_k) \notin D$ .

132 • *Random Epitope*. In this strategy,  $t_i$  is combined with epitopes sampled uniformly from all  
133 epitopes without its true epitope binder(s) [12, 16, 19]. The sampled negative pairs are  $n_i =$   
134  $\{(e_j, t_i)\}_{j=1}^M$  with  $e_j \in \{e_i\}$  and  $(e_j, t_i) \notin D$ .

135 • *Unified Epitope*. Compared to *Random Epitope*, the only difference of *Unified Epitope* is that,  
136 the epitopes are sampled according to their frequency distributions in the positive dataset [11,  
137 20]; i.e.  $n_i = \{(e_j, t_i)\}_{j=1}^M$  and  $P_{pos}(e_j) \sim P_{neg}(e_j)$ . This strategy ensures that the frequencies  
138 of epitopes are unified in the negative data and positive data.

139 A systematical comparison between these four negative sampling strategies is an urgent need for  
140 benchmarking different models and guiding the development of accurate and generalized models in  
141 future works. To address this demand, we referred to the field of recommender system and selected  
142 three evaluation metrics that can be calculated without the attendance of negative samples.

143 **Precision@k and Recall@k**. These two metrics measure the exactness and completeness of  
144 the top k binding predictions for a given TCR. Assume that a TCR  $t_i$  in the test set  $\{(e_i, t_i)\}_{i=1}^N$   
145 can bind to a number of  $b_i$  epitopes (due to cross-reactivity), and a number of  $m_i$  true interacting  
146 pairs  $\{(e_j, t_j)\}_{j=1}^{m_i}$  lie in the top k predictions, then these two metrics are defined as follows:

$$Precision@k = \frac{1}{N} \sum_{i=1}^N \frac{m_i}{k} \quad (1)$$

147

$$Recall@k = \frac{1}{N} \sum_{i=1}^N \frac{m_i}{b_i} \quad (2)$$

148 where  $N$  is the total number of TCRs in the test set. A higher value of *Precision@k* indicates that  
149 the more true binding pairs can be found among the top k predicted pairs; And a higher value of  
150 *Recall@k* suggests a higher proportion of predicted binding pairs over all the true binding pairs.

151 **NDCG@k**. The previous two metrics overlook the order of the predictions since the ranking of  
152 the true predicted binding pairs does not affect the values of both metrics as long as they are in the  
153 top k predictions. The Normalized Discounted Cumulative Gain (NDCG) measures how relevant  
154 the predictions are and how good the ordering is, which is calculated by:

$$NDCG@k = \frac{DCG@k}{IDCG@k}, \quad (3)$$

155 where the definitions and formulas of DCG@k and IDCG@k are described in Supplementary Text S1.  
156 Overall, these three metrics are complementary to each other and help to determine the superiority  
157 of the four negative sampling strategies. A higher value of any of the three metrics indicates a better  
158 model performance.

## 159 Pretrained encoder

160 To numerically encode TCR and epitopes, we capitalized on the transfer learning technique. We have  
161 earlier proposed an autoencoder model, TCRpeg [24] that utilizes a recurrent neural network with  
162 GRU layers to characterize the TCR repertoires and demonstrated that it can produce high-quality  
163 vector encodings for TCR sequences. As an autoencoder model, TCRpeg is capable of capturing  
164 key features of sequence input via unsupervised learning of mapping between the latent space and  
165 sequence space, and more importantly, using TCRpeg for pretraining only needs plain amino acid  
166 sequences which are currently abundant. In addition, unlike the encoders in TITAN or pMTNet,  
167 TCRpeg can process sequences of arbitrary lengths without the need to pad them to a fixed length.  
168 Thus, we decided to employ two separate pretrained TCRpeg models as the encoders for TCRs and  
169 epitopes, respectively. A detailed description of TCRpeg is given in Supplementary Text S2.

170 To pretrain TCRpeg for encoding TCRs (TCRpeg-TCR), we fed TCRpeg with  $10^6$  TCR se-  
171 quences collected from Emerson *et al.* [25]. We set the feature size of TCRpeg to 768 and trained  
172 it for 20 epochs by minimizing the cross-entropy loss between the output soft-maxed logits and the  
173 one-hot encoded representation of the input sequences. For encodings of epitopes, we trained an-  
174 other TCRpeg model (TCRpeg-Epi) with the identical architecture of TCRpeg-TCR using 362,456  
175 unique epitope sequences collected from Mei *et al.* [26] with lengths ranging from 8 to 14 amino  
176 acids. Details on the pretraining process of TCRpeg are elaborated in Supplementary Text S3.

## 177 Model architecture

178 Figure 2B delineates an overview of the architecture of TEINet. Conceptually, the complex task  
179 of predicting the TCR-epitope interaction is decomposed into two steps to lower the difficulty level  
180 of the final prediction task. First, two encoding networks are pretrained so that the amino acid  
181 sequences of TCRs and epitopes could be represented by numerical vectors. Next, we concatenated  
182 these two vector encodings to form the final representations for TCR-epitope pairs. In the final step,  
183 we built a fully connected neural network (FCN) on top of these combined vector encodings to fuse  
184 the knowledge extracted from TCRs and epitopes. Specifically, the FCN consists of three hidden  
185 layers with 768, 384, and 96 neurons with the dropout [27] rate set to 0.15 to prevent overfitting.  
186 Before feature concatenation, we employed the layer normalization [28] to numerically stabilize each  
187 group of features. All neurons use the scaled exponential linear unit (SELU [29]) activation function,  
188 except for the output neuron which applies the sigmoid activation function.

## 189 Model training

190 TEINet was implemented in Python 3.6 and built on the deep learning framework PyTorch [30].  
191 TEINet was trained and evaluated under a 5-fold cross-validation procedure. Instead of inferring  
192 TEINet on a static dataset with negative pairs sampled prior to the training process, we adopted a  
193 dynamic sampling strategy: the negative examples are sampled on the fly at each training step using  
194 the sampling strategies described in the previous section. This dynamic sampling strategy demon-  
195 strates improved performance over static training (Supplementary Figure S1). For all experiments

196 in this work, the negative pairs were sampled 10 times more than positive pairs. TEINet optimized  
197 binary cross entropy loss with Adam algorithm [31] and an initial learning rate of  $1 \times 10^{-3}$ . The  
198 model was trained for 30 epochs with a batch size of 48. The learning rate was reduced at the 21st  
199 and 27th epoch by a factor of 0.1.

## 200 Results

### 201 Comparison of different negative sampling strategies

202 We trained TEINet with each negative sampling method and observed that they could achieve  
203 performance in different scales (Supplementary Figure S2). For example, using *Reference TCR*  
204 leads to an average AUROC (area under the receiver operating characteristic) of 0.797, whereas the  
205 performance achieves an AUROC of 0.934 under *Random Epitope*, which is unexpectedly high yet  
206 useless.

207 We first compared the three negative sampling methods: *Random TCR*, *Reference TCR*, and  
208 *Unified Epitope*. The negative data generated by these methods possess similar frequency distribu-  
209 tions of epitopes with those in the positive data. Table 1 shows the Precision, Recall, and NDCG of  
210 each schema using the TEINet. These results first demonstrated that *Random TCR* and *Reference*  
211 *TCR* obtained similar performance, indicating that sampling TCRs from the reference TCR pool or  
212 TCRs in positive data have a comparable effect on the model training. *Reference TCR* is slightly  
213 better than *Random TCR*, as TCRs drawn from another sequence pool constructed from healthy  
214 donors are less likely to interact with epitopes than shuffled TCRs from the positive data; i.e., *Ran-*  
215 *dom TCR* might produce more false negative pairs. *Unified Epitope* achieved superior performance  
216 among these three strategies by a large margin. It indicates that *Unified Epitope* can help develop a  
217 more robust and generalized model for the TCR-epitope interaction prediction task. We attributed  
218 its superior performance to the uniformity of the distribution of TCRs and epitopes across positive  
219 and negative data.

220 We next contrasted *Unified Epitope* with *Random Epitope*. It seems that *Random Epitope* is a  
221 perfect sampling strategy since it achieved extremely high values of Precision, Recall, and NDCG  
222 (Table 1), and achieved an average AUROC of 0.934. However, these high values are overesti-  
223 mated and misleading due to the inherent imbalance of the data. Note that the number of epitope-  
224 interacting TCRs follows an extreme long-tail distribution (Fig. 3A, 1B, and Supplementary Figure  
225 S3) that most TCRs (70%) are associated with the top 5% epitopes. As a result, *Random Epitope*  
226 would produce skewed negative data that most majority of epitopes were matched with far more  
227 negative TCRs than positive TCRs (Supplementary Figure S3). Trained with such a skewed dataset,  
228 TEINet was driven to make predictions based on the epitope sequences without the participation of  
229 TCRs, as discussed in Dens *et al.* [32]. That is, when the input pairs consist of frequent epitopes, the  
230 model tends to predict “1s”, and conversely, it is likely to predict “0s” when encountering pairs with  
231 infrequent epitopes. Thus, TEINet with *Random Epitope* obtains a misleading high performance:  
232 (1) for TCRs, TEINet often predicts high scores when they are linked to frequent epitopes, which  
233 results in high Precision, Recall, and NDCG since frequent epitopes appear in most paired samples;



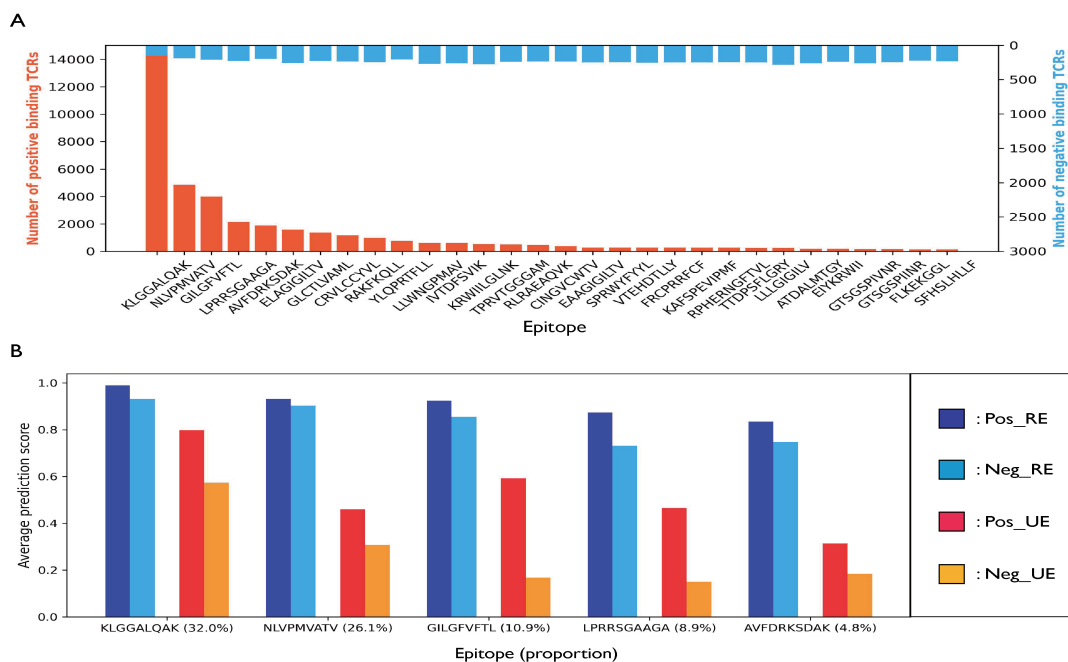


Figure 3: Distribution of the number of epitope-associated TCRs and the average prediction scores for them. (A) Distribution of the number of positive and negative TCRs sampled by *Random Epitope* for the 30 most abundant epitopes. Given that the epitopes are sampled randomly for each TCR and the epitope-associated TCRs follow an extreme long-tail distribution, there are far fewer negative samples than positive samples for abundant epitopes, whereas for most epitopes, there are far more negative samples than positive samples. (B) The average prediction scores for the positive and negative pairs of the top 5 most abundant epitopes. “Pos” and “Neg” stand for positive and negative samples; “RE” and “UE” represent *Random Epitope* and *Unified Epitope*. We observed that for both positive and negative pairs of abundant epitopes, *Random Epitope* will produce high predictive scores. Such a problem is greatly relieved by *Unified Epitope*.

234 (2) for epitopes, TEINet tends to predict high scores for pairs with frequent epitopes that possess  
 235 abundant positive binding TCRs and sparse negative binding TCRs, and low scores for pairs with  
 236 rare epitopes that are linked to abundant negative TCRs and sparse positive binding TCRs, which  
 237 leads to high AUROC. Indeed, TEINet with *Random Epitope* obtained high prediction scores for  
 238 both positive and negative pairs with frequent epitopes (Fig. 3B). For instance, it outputs an aver-  
 239 age prediction score of 0.99 and 0.93 for respective positive and negative pairs of the most frequent  
 240 epitope (KLGALQAK). As a result, those negative pairs will be classified as false positives in the  
 241 generic performance evaluation. Moreover, due to the long-tail distribution of the epitope-associated  
 242 TCRs, *Random Epitope* will generate far fewer negative pairs than positive pairs, so that those false  
 243 positives have minor impact on the generic performance evaluation, resulting in a misleading high  
 244 AUROC.

245 Overall, our results and analysis indicate that *Unified Epitope* is more appropriate for negative  
 246 sampling in the TCR-epitope prediction task, which is further supported in the evaluation on in-  
 247 dependent datasets (see the following section). To eliminate potential model bias introduced by

248 TEINet, we performed the same experiments using the ImRex model and obtained similar results  
 249 (Supplementary Table S1). In the remaining experiments, *Unified Epitope* is selected as the default  
 250 strategy.

Table 1: The Precision, Recall, and NDCG of each negative sampling method.

Method	Precision@3	Recall@3	NDCG@3	Precision@10	Recall@10	NDCG@10
<i>Reference TCR</i>	0.093±0.002	0.275±0.007	0.255±0.006	0.036±0.001	0.356±0.009	0.284±0.007
<i>Random TCR</i>	0.085±0.002	0.251±0.004	0.226±0.004	0.033±0.001	0.322±0.003	0.252±0.003
<i>Unified Epitope</i>	0.129±0.004	0.380±0.012	0.334±0.011	0.052±0.001	0.506±0.014	0.380±0.012
<i>Random Epitope</i>	0.192±0.002	0.567±0.006	0.484±0.006	0.081±0.001	0.788±0.002	0.565±0.005

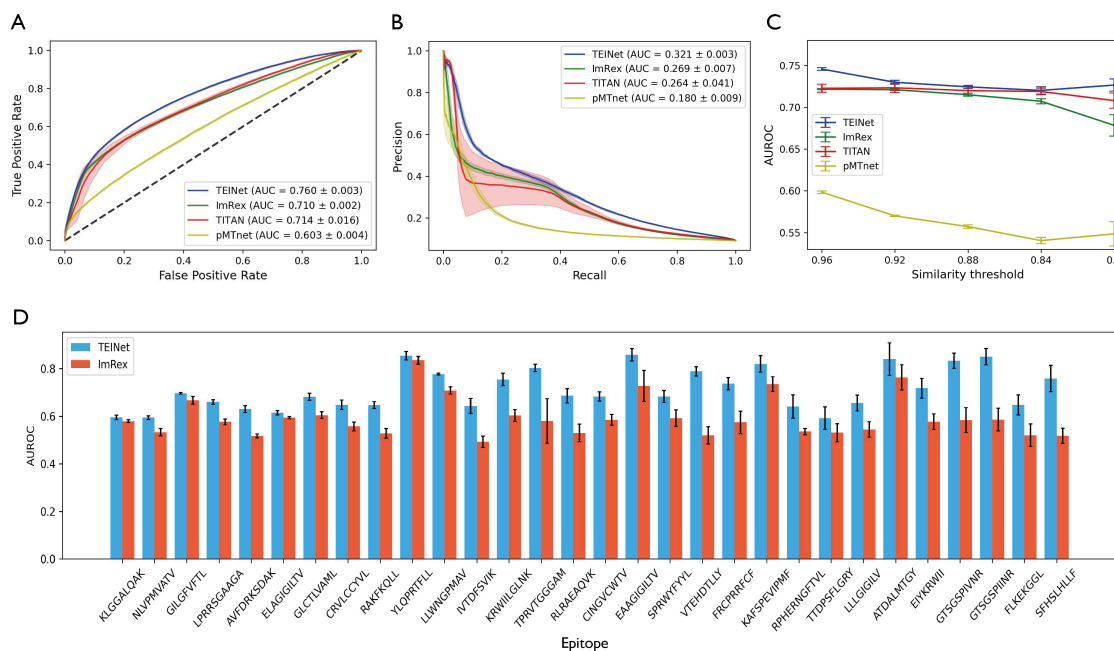


Figure 4: Performance of TEINet and the three baseline models. (A) The receiver operator characteristic (ROC) curves for each model. The area under the ROC curve (AUROC) values are shown in the legend. (B) The precision-recall (PRC) curves for each model. The area under the PRC curve (AUPRC) values are shown in the legend. (C) The AUROC for each model according to different similarity thresholds for filtering the test set. (D) The per-epitope AUROC performance for the top 30 most abundant epitopes. TEINet outperforms ImRex in these epitopes.

## 251 Performance of TEINet

252 To assess the performance of TEINet, we compared it with three existing approaches: ImRex [11],  
 253 TITAN [12] and pMTNet [14]. ImRex encodes TCRs and epitopes based on their physicochemical  
 254 properties and utilizes a CNN (convolutional neural network) to process the combined encodings.

255 Similar to our proposed TEINet, TITAN and pMTNet both make use of the pretrained encoders.  
256 The pMTNet additionally incorporates the information of the MHC allele associated with the epitope  
257 to make the prediction.

258 Figure 4A and 4B show the AUROC and AUPRC of TEINet as well as the three baseline mod-  
259 els. TEINet outperforms the baseline methods with an AUROC of 0.760 and an AUPRC of 0.321,  
260 while the second best comparative model ImRex has an AUROC of 0.714 and an AUPRC of 0.269.  
261 Moreover, we calculated the Precision, Recall, and NDCG of ImRex and still observed superior per-  
262 formance of TEINet (Supplementary Table S1). With learnable encoders that possess the capability  
263 of processing sequences in any length, TEINet can better extract sequence information and conse-  
264 quently make more accurate predictions. To investigate whether the superiority of TEINet retains  
265 when the similarity of TCRs between training and evaluation datasets decreases, we filtered out pairs  
266 in the test set with specific TCRs according to the Levenstein similarity thresholds (Supplementary  
267 Text S4). Figure 4C demonstrates the corresponding performance of each model under different  
268 similarity thresholds. Again, TEINet outperforms other baseline models. Further, to resolve the  
269 concern that pairs consisting of frequent epitopes would dominate the effects on performance, we  
270 report the per-epitope AUROC derived by evaluating on paired data for one specific epitope in  
271 Fig. 4D. We found no explicit correlation between the AUROC and the number of training samples,  
272 indicating that the complexities of the binding pattern for each epitope are different. Similar results  
273 were also found in Moris *et al.* [11]. Besides, TEINet is still superior, with the ImRex lagging behind  
274 for most epitopes (Fig. 4D).

## 275 **Impact of pretraining**

276 Transfer learning is becoming an integral part of the design of deep learning models for the prediction  
277 of TCR binding specificity. Recently developed models tend to employ pretrained encoders to  
278 transform amino acid sequences into vector representations [12, 14, 15, 19]. An analysis of the  
279 impact of the pretraining step is in demand to provide a better understanding of the pretrained  
280 encoders.

281 First, without the pretraining step, the performance of TEINet dropped significantly with an  
282 AUROC of 0.675, which demonstrated the necessity of the pretraining step. Next, we explored the  
283 influence of the TCR and epitope encoders singly and simultaneously (Fig. 5A-C). It is clear that the  
284 pretraining of TCRs greatly enhanced the model performance (Fig. 5A), whereas the pretraining of  
285 epitopes only brought about slight and unstable improvement (Fig. 5B). Given that the diversity of  
286 TCRs (41,610 unique samples) is much higher than that of epitopes (180 unique samples), pretraining  
287 of TCRs enables them to be distributed separably in the feature space, which is more important  
288 for making a prediction. Further, these two encoders improved the performance synergistically and  
289 achieved the best performance (Fig. 5C). Utilizing both pretrained encoders enhanced the AUROC  
290 by around 0.01 than using the TCR encoder alone. Notably, we observed that when the pretraining  
291 of the TCR encoder exceeded a certain epoch, the final performance dropped (Fig. 5A and C).  
292 Thus, the degree of the pretraining needs to be tuned carefully; otherwise, the model might have  
293 the problem of overfitting.

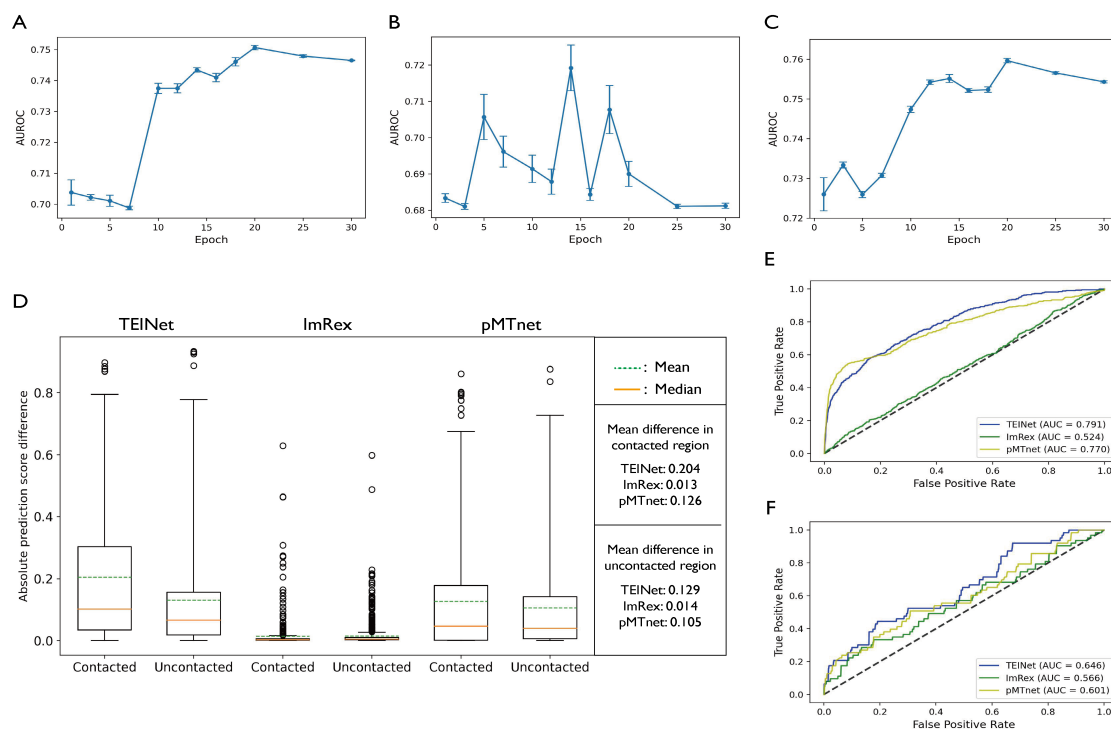


Figure 5: Investigation of the impact of the pretraining stage and further validations of TEINet. (A-C) The AUROC values with the encoders trained for different epochs. (A) Only pretrain the encoder for TCRs. (B) Only pretrain the encoder for epitopes. (C) Pretrain the encoders for TCRs and epitopes. (D) The absolute difference of prediction scores for each model between the contacted and uncontacted residues. Using TEINet, residues with direct contacts are more likely to induce larger changes in the predicted binding strength than non-contact residues. (E and F) The ROC curves for each model in the two independent test sets: (E) TBAdb and (F) PDB. The corresponding AUROC values are shown in the legend.

## 294 Structural analysis

295 Perturbation (mutational) analysis can be used to detect the important amino acid residues for  
 296 the model prediction [14, 24, 33]. We grouped TCR residues by whether or not they formed any  
 297 direct contact with any residue of epitopes within  $5\text{\AA}$  and assumed that substitutions inside the  
 298 contact region would lead to dramatic changes in the predicted binding score. To analyze the  
 299 effects of predictive models on the contact/non-contact region, we collected 105 solved TCR-epitope  
 300 interacting complex structures from the public RCSB Protein Data Bank (PDB) database [34] as  
 301 the ground truth data. We performed the alanine scanning technique in biophysics studies [35]  
 302 on the TCRs in the PDB database using the predictive models. Figure 5D illustrates the average  
 303 score difference for each model inside the contact and non-contact region. We observed that for  
 304 TEINet, the contact residues were more likely to induce larger drops in predicted TCR-epitope  
 305 binding strength than non-contact residues, which supports our assumption.

## 306 Evaluation on independent datasets

307 To further compare the predictive performance of each model, we collected two independent test  
308 sets. We selected the TBADB [36] dataset, which includes 439 binding pairs on 414 unique TCRs  
309 and 42 epitopes as our first independent test set; The 105 interacting pairs extracted from the PDB  
310 database aforementioned were selected as the second independent test set. As before, the same  
311 filtering procedure was applied to them. Figure 5E and 5F show the performance of each model  
312 on the independent test sets. Again, TEINet achieved superior performance over the other baseline  
313 methods. Note that for the PDB dataset, TEINet obtained a lower AUROC value of 0.646. We  
314 attributed it to the small overlap of epitopes as there is only 1 epitope in the PDB dataset that also  
315 appears in the training data. Moreover, given that the PDB dataset is an approximately balanced  
316 dataset with each epitope binding with 1 or 2 TCRs, the *Random Epitope* and *Unified Epitope* will  
317 generate similar negative data, which enables us to compare these two strategies by the AUROC  
318 value. Thus, we trained two TEINets each using *Random Epitope* or *Unified Epitope* during the  
319 training process and then evaluated them on the PDB database constructed with *Random Epitope*.  
320 We observed that TEINet trained with *Random Epitope* obtained an AUROC of 0.572, which was  
321 surpassed by *Unified Epitope* by a large margin with an AUROC of 0.644 (Supplementary Figure  
322 S4). This finding further supports the advantage of *Unified Epitope*.

## 323 1 Discussion

324 The prediction of TCR specificity to epitope has been a challenging problem. The immense search-  
325 ing space of immune receptors, lack of curated training samples, and absence of negative samples  
326 remain issues for algorithm development. In recent years, public databases have been accumulat-  
327 ing an enormous amount of TCR-epitope interacting data. Benefiting from the enrichment of data  
328 enrichment of available data, it is possible to develop accurate deep learning models to tackle the  
329 challenging task of TCR-epitope interaction prediction.

330 In this work, we have proposed TEINet, a new deep learning model for predicting the TCR  
331 binding specificity. TEINet only requires the CDR3 $\beta$  chain of the TCR and epitope sequence of  
332 the pMHC complex to make the prediction. Though the CDR3 $\alpha$  chain and the MHC allele are  
333 shown to be beneficial in this task [10, 13–15, 23], the paired data is still rare compared to single-  
334 chain data, which limits the generalizability of the pair-chain model. We leave the exploration of  
335 both CDR3 chains and MHC alleles in future work. TEINet employed the TCRpeg [24], a deep  
336 autoregressive model, to extract the sequence information of TCRs and epitopes and transform  
337 them into numerical vector space. The TCRpeg was pretrained in a self-supervised learning manner  
338 on large-scale sequence data to learn a more general pattern to encode TCRs/epitopes. TEINet  
339 then combined the encodings of TCRs and epitopes and used a fully-connected neural network to  
340 make the final prediction, leveraging the knowledge from TCRs and epitopes.

341 To train and evaluate a supervised model, negative samples are required. However, currently  
342 there is no unified method for negative sampling, which poses a challenge for comparing different  
343 models. For example, *Random TCR* was applied in pMTNet [14]; *Reference TCR* was applied in

344 NetTCR [15, 18]; *Random Epitope* was employed in TITAN [12]; *Unified Epitope* was employed  
345 in ImRex [11]. We thus proposed three metrics, Precision, Recall, and NDCG that are unrelated  
346 to negative samples to compare different sampling strategies. We manifested that *Unified Epitope*  
347 is the winner among these four sampling schemas for the development of a more accurate model,  
348 given that it achieved superior Precision, Recall, and NDCG among the first three schemas and  
349 that *Random Epitope* breaks the uniformity between positive and negative data, which leads to  
350 misleading performance. Thus, we recommend *Unified Epitope* as the default negative sampling  
351 method in future works.

352 To showcase the predictive strength of TEINet, we compared TEINet with another three pub-  
353 lished deep learning models: ImRex [11], TITAN [12], and pMTNet [14]. We performed the 5-fold  
354 cross-validation procedure on our constructed dataset which consists of 44,682 interacting pairs. We  
355 observed that TEINet achieved an AUROC of 0.760 and an AUPRC of 0.321 and outperformed other  
356 comparative models with the best AUROC of 0.714 and AUPRC of 0.269. Further, we also evaluated  
357 and compared these models on two additional independent test sets. Again, TEINet surpassed other  
358 baseline models.

359 The usage of the transfer learning technique has become a trend in the design of deep learning  
360 models for the TCR-epitope binding prediction task. Instead of using the physicochemical properties  
361 of amino acid sequences to construct the features of TCRs and epitopes, many recently published  
362 models capitalized on the pretrained encoders that leveraged the knowledge learned from other  
363 tasks with abundant data [12, 14, 15, 19]. However, the impact of the pretraining step on the final  
364 prediction accuracy remains unknown, which could potentially hinder the exploitability of pretrained  
365 encoders. Here, we disentangled the effect from each encoder (Fig. 5A-5C). We first observed that  
366 the pretraining of the TCR encoder improved the TEINet by a much larger margin than that of  
367 the epitope encoder, which could be explained by the vast diversity of TCRs. More importantly,  
368 we found that excessive pretraining might harm the performance, so that the degree of pretraining  
369 needs to be tuned carefully.

370 At last, we analyzed whether the prediction from TEINet can reveal the structural information  
371 of the interacting complex. We grouped residues of TCRs that form any contact with epitope within  
372  $5\text{\AA}$  into the contact region. Contact residues should be more important than non-contact residues  
373 in forming the interaction between TCRs and epitopes [37]. Indeed, larger drops of predicted scores  
374 were observed inside the contact region than non-contact region using TEINet.

375 In summary, we have designed TEINet to predict the interaction between TCRs and their epi-  
376 tope targets. Our results demonstrate that TEINet achieved superior performance over three other  
377 comparative models only by using the information of CDR3 $\beta$  chains and epitope sequences. We  
378 also compared different negative sampling strategies and suggested that *Unified Epitope* is more  
379 appropriate for the development of a generalized model. We expected that with enhanced accuracy  
380 in predicting the potential immune response of T-cells to epitopes, TEINet could be beneficial for  
381 the *in silico* design and implementation of immunotherapy in the era of personalized medicine.

## Acknowledgments

We thank all contributors to VDJdb, McPAS-TCR, and other TCR specificity datasets for making their data publicly available.

## References

- [1] D. J. Laydon, C. R. Bangham, and B. Asquith, “Estimating t-cell repertoire diversity: Limitations of classical estimators and a new approach,” *Philosophical Transactions of the Royal Society B: Biological Sciences*, vol. 370, no. 1675, p. 20140291, 2015.
- [2] S.-Q. Zhang *et al.*, “High-throughput determination of the antigen specificities of t cell receptors in single cells,” *Nature biotechnology*, vol. 36, no. 12, pp. 1156–1159, 2018.
- [3] J. D. Altman *et al.*, “Phenotypic analysis of antigen-specific t lymphocytes,” *Science*, vol. 274, no. 5284, pp. 94–96, 1996.
- [4] T. Kula *et al.*, “T-scan: A genome-wide method for the systematic discovery of t cell epitopes,” *Cell*, vol. 178, no. 4, pp. 1016–1028, 2019.
- [5] M. Shugay *et al.*, “Vdjdb: A curated database of t-cell receptor sequences with known antigen specificity,” *Nucleic acids research*, vol. 46, no. D1, pp. D419–D427, 2018.
- [6] R. Vita *et al.*, “The immune epitope database (iedb): 2018 update,” *Nucleic acids research*, vol. 47, no. D1, pp. D339–D343, 2019.
- [7] N. Tickotsky, T. Sagiv, J. Prilusky, E. Shifrut, and N. Friedman, “Mcpas-tcr: A manually curated catalogue of pathology-associated t cell receptor sequences,” *Bioinformatics*, vol. 33, no. 18, pp. 2924–2929, 2017.
- [8] W. D. Chronister *et al.*, “Tcrmatch: Predicting t-cell receptor specificity based on sequence similarity to previously characterized receptors,” *Frontiers in immunology*, vol. 12, p. 640725, 2021.
- [9] S. Gielis *et al.*, “Detection of enriched t cell epitope specificity in full t cell receptor sequence repertoires,” *Frontiers in immunology*, vol. 10, p. 2820, 2019.
- [10] E. Jokinen, J. Huuhtanen, S. Mustjoki, M. Heinonen, and H. Lähdesmäki, “Predicting recognition between t cell receptors and epitopes with tcrgp,” *PLoS computational biology*, vol. 17, no. 3, e1008814, 2021.
- [11] P. Moris *et al.*, “Current challenges for unseen-epitope tcr interaction prediction and a new perspective derived from image classification,” *Briefings in Bioinformatics*, vol. 22, no. 4, bbaa318, 2021.
- [12] A. Weber, J. Born, and M. Rodriguez Martinez, “Titan: T-cell receptor specificity prediction with bimodal attention networks,” *Bioinformatics*, vol. 37, no. Supplement\_1, pp. i237–i244, 2021.

- 416 [13] W. Zhang *et al.*, “A framework for highly multiplexed dextramer mapping and prediction of t  
417 cell receptor sequences to antigen specificity,” *Science Advances*, vol. 7, no. 20, eabf5835, 2021.
- 418 [14] T. Lu *et al.*, “Deep learning-based prediction of the t cell receptor–antigen binding specificity,”  
419 *Nature Machine Intelligence*, vol. 3, no. 10, pp. 864–875, 2021.
- 420 [15] A. Montemurro *et al.*, “Nettcr-2.0 enables accurate prediction of tcr-peptide binding by using  
421 paired tcr $\alpha$  and  $\beta$  sequence data,” *Communications biology*, vol. 4, no. 1, pp. 1–13, 2021.
- 422 [16] I. Springer, H. Besser, N. Tickotsky-Moskovitz, S. Dvorkin, and Y. Louzoun, “Prediction of  
423 specific tcr-peptide binding from large dictionaries of tcr-peptide pairs,” *Frontiers in immunol-*  
424 *ogy*, p. 1803, 2020.
- 425 [17] P. Dash *et al.*, “Quantifiable predictive features define epitope-specific t cell receptor reper-  
426 toires,” *Nature*, vol. 547, no. 7661, pp. 89–93, 2017.
- 427 [18] V. I. Jurtz *et al.*, “Nettcr: Sequence-based prediction of tcr binding to peptide-mhc complexes  
428 using convolutional neural networks,” *BioRxiv*, p. 433 706, 2018.
- 429 [19] Y. Fang, X. Liu, and H. Liu, “Attention-aware contrastive learning for predicting t cell receptor-  
430 antigen binding specificity,” *bioRxiv*, 2022.
- 431 [20] M. Cai, S. Bang, and H. Lee, “Tcr-epitope binding affinity prediction using multi-head self  
432 attention model,”
- 433 [21] X. Hou *et al.*, “Analysis of the repertoire features of tcr beta chain cdr3 in human by high-  
434 throughput sequencing,” *Cellular Physiology and Biochemistry*, vol. 39, no. 2, pp. 651–667,  
435 2016.
- 436 [22] E. Lanzarotti, P. Marcatili, and M. Nielsen, “T-cell receptor cognate target prediction based on  
437 paired  $\alpha$  and  $\beta$  chain sequence and structural cdr loop similarities,” *Frontiers in immunology*,  
438 vol. 10, p. 2080, 2019.
- 439 [23] I. Springer, N. Tickotsky, and Y. Louzoun, “Contribution of t cell receptor alpha and beta  
440 cdr3, mhc typing, v and j genes to peptide binding prediction,” *Frontiers in immunology*,  
441 vol. 12, p. 664 514, 2021.
- 442 [24] Y. Jiang and S. C. Li, “Deep autoregressive generative models capture the intrinsics embedded  
443 in t-cell receptor repertoires,” *bioRxiv*, 2022.
- 444 [25] R. O. Emerson *et al.*, “Immunosequencing identifies signatures of cytomegalovirus exposure  
445 history and hla-mediated effects on the t cell repertoire,” *Nature genetics*, vol. 49, no. 5,  
446 pp. 659–665, 2017.
- 447 [26] S. Mei *et al.*, “Anthem: A user customised tool for fast and accurate prediction of binding be-  
448 tween peptides and hla class i molecules,” *Briefings in Bioinformatics*, vol. 22, no. 5, bbaa415,  
449 2021.
- 450 [27] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A  
451 simple way to prevent neural networks from overfitting,” *The journal of machine learning*  
452 *research*, vol. 15, no. 1, pp. 1929–1958, 2014.



- 453 [28] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*,  
454 2016.
- 455 [29] G. Klambauer, T. Unterthiner, A. Mayr, and S. Hochreiter, “Self-normalizing neural net-  
456 works,” *Advances in neural information processing systems*, vol. 30, 2017.
- 457 [30] A. Paszke *et al.*, “Pytorch: An imperative style, high-performance deep learning library,”  
458 *Advances in neural information processing systems*, vol. 32, 2019.
- 459 [31] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint*  
460 *arXiv:1412.6980*, 2014.
- 461 [32] C. Dens, W. Bittremieux, F. Affaticati, K. Laukens, and P. Meysman, “Interpretable deep  
462 learning to uncover the molecular binding patterns determining tcr–epitope interactions,”  
463 *bioRxiv*, 2022.
- 464 [33] J.-W. Sidhom, H. B. Larman, D. M. Pardoll, and A. S. Baras, “Deeptcr is a deep learning  
465 framework for revealing sequence concepts within t-cell repertoires,” *Nature communications*,  
466 vol. 12, no. 1, pp. 1–12, 2021.
- 467 [34] J. L. Sussman *et al.*, “Protein data bank (pdb): Database of three-dimensional structural  
468 information of biological macromolecules,” *Acta Crystallographica Section D: Biological Crys-*  
469 *tallography*, vol. 54, no. 6, pp. 1078–1084, 1998.
- 470 [35] G. A. Weiss, C. K. Watanabe, A. Zhong, A. Goddard, and S. S. Sidhu, “Rapid mapping of  
471 protein functional epitopes by combinatorial alanine scanning,” *Proceedings of the National*  
472 *Academy of Sciences*, vol. 97, no. 16, pp. 8950–8954, 2000.
- 473 [36] W. Zhang *et al.*, “Pird: Pan immune repertoire database,” *Bioinformatics*, vol. 36, no. 3,  
474 pp. 897–903, 2020.
- 475 [37] D. Chowell *et al.*, “Tcr contact residue hydrophobicity is a hallmark of immunogenic cd8+ t  
476 cell epitopes,” *Proceedings of the National Academy of Sciences*, vol. 112, no. 14, E1754–E1762,  
477 2015.