1 # TCRfinder: Improved TCR virtual screening for novel

2 # antigenic peptides with tailored language models

3

4 Yang Li[1,‡], Chaoting Zhang[2,‡], Xi Zhang[5], Yang Zhang[1,3,4,*]

5

6 *[1]Cancer Science Institute of Singapore, National University of Singapore, 117599,*

7 *Singapore.*

8 *[2]Key Laboratory of Carcinogenesis and Translational Research, Laboratory of*

9 *Biochemistry and Molecular Biology, Peking University Cancer Hospital and Institute,*

10 *100084, China.*

11 *[3]Department of Computer Science, School of Computing, National University of Singapore,*

12 *117417, Singapore.*

13 *[4]Department of Biochemistry, Yong Loo Lin School of Medicine, National University of*

14 *Singapore, 117596 Singapore.*

15 *[5]Center for AI and Computational Biology, Suzhou Institute of Systems Medicine, Chinese*

16 *Academy of Medical Sciences & Peking Union Medical College, Suzhou, 215123 China*

17

18 [‡] These authors contributed equally

19

20 *Correspondence should be addressed to Yang Zhang (email: zhang@zhanggroup.org)

21

22

23
24                                              **Abstract**
25
26      Accurate modeling of T-cell receptor (TCR) and peptide interactions is essential for
27      immunoreaction elucidation and T-cell-based immunotherapeutic developments. We
28      developed TCRfinder, a novel deep-learning architecture for TCR-peptide binding
29      prediction and virtual screening. Large-scale benchmark experiments demonstrated a
30      robust capability of TCRfinder in distinguishing interacting and non-interacting TCRs for
31      unseen peptides, with accuracy significantly beyond current state-of-the-art methods.
32      Furthermore, TCRfinder recognizes tumor neoantigen mutations from wild-type antigens
33      of given TCRs, with a success rate nearly 50% higher than the best of existing methods.
34      Detailed data analyses showed that the major advantage of TCRfinder lies in the specially
35      trained TCR and peptide language models tailored with iterative attention network
36      architecture, which can precisely reveal physical interaction patterns of cross-chain atoms
37      and substantially enhance the precision of TCR-peptide interaction predictions. The
38      open-source TCRfinder program can help facilitate large-scale deployment of high-quality
39      TCR and neoantigen virtual screening, offering exciting potential for personalized
40      TCR-based immunotherapies.
41
42      **Keywords:** TCR-peptide interaction, TCR and neoantigen virtual screening, language
43      model, attention network
44

**INTRODUCTION**

The T-cell receptor (TCR) is a protein complex found on the surface of T cells, playing a crucial role in the recognition of antigen peptides bound to the major histocompatibility complex (MHC) molecules. TCRs can recognize antigens, originating from viruses, bacteria, or somatically mutated genes (i.e., neoantigens), to activate cytotoxic T cells (CD8+ T cells) or helper T cells (CD4+ T cells), thereby cleaning up invading pathogens and tumor cells. As a result, the identification and characterization of TCRs that effectively interact with specific antigens are fundamental for understanding immune responses and designing targeted immunotherapies.

Previous studies primarily relied on experimental approaches, such as tetramer[1], TetTCR-seq[2] and T-scan[3], to identify antigen-specific TCRs. However, these methods are not only time-consuming and costly but also have a general low success rate[4]. Thus, there is an urgent demand for effective TCR screening based on computational methods, which would significantly reduce the time and cost of identifying antigen specific TCRs.

The initial computational methods, including TCRdist[5], DeepTCR[6], TCRGP[7], TCRex[8], and NetTCR[9], utilized various approaches such as CDR similarity-weighted distances, random forests[10], Gaussian process classification methods[11], and convolutional neural networks (CNNs) to distinguish between positive and negative TCRs for antigenic peptides. However, their utility is constrained, particularly when involving novel or unseen peptides that have not been encountered in the training datasets. Recent advancements in TCR screening methods aim to overcome the limitation of applicability to unseen peptides in methodology. For example, pMTnet[12] employs deep transfer learning and leverages pretrained autoencoders, to predict interactions between TCRs and peptides. ERGO[13,14] and DLpTCR[15] utilize CNNs, or ensembles of CNNs, to recognize TCR–antigen binding. More recent approaches, such as ATM-TCR[16], AttnTAP[17], and PanPep[18], incorporate deep learning models with attention mechanisms to predict these interactions. While these methods demonstrate the potential to identify TCR–peptide interactions for unseen peptides, their accuracy remains suboptimal[19], likely due to the constraints imposed by the limited availability of data.

In this work, we introduce TCRfinder, a novel approach that harnesses the power of language models[20] (LMs) to address the challenge posed by limited data in virtual TCR and antigen screenings. Specifically, we have trained two special language models, i.e., TCR LM and peptide LM, from scratch to enhance the representations of TCR and peptide sequences. We found that the two LMs exhibit an improved prediction power for different TCR-peptide interactions and achieve the best overall performance, when integrated with an iterative attention network architecture. TCRfinder represents a significant stride toward more reliable TCR and antigen virtual screening, showcasing the potential of language models in advancing our understanding of immune system responses to diverse antigens.

84

## RESULTS

86      The core of the TCRfinder is a deep learning model trained for accurately scoring the
87 interactions between the antigenic peptides and TCR sequences. Because the CDR3 region
88 of TCR β-chain is the key determinant of the interactions, our model is trained mainly on
89 the β-chain CDR3 regions of the TCR sequences. As outlined in **Figure 1A**, TCRfinder
90 first employs a Joint Embedder module (**Figure 1D**), which utilizes two pre-trained LMs
91 (**Figure 1B**) for peptide and TCRβ CDR3 sequences separately, followed by iterative
92 transformer block learnings starting with concatenated and split embeddings (**Figure 1C**).
93 The final TCR-peptide interaction scores are derived through a multi-layer perceptron
94 (MLP) network based on the concatenated representations of TCR and peptide
95 embeddings from the preceding transformer networks (**Figure 1**).

96

**Specifically trained language models enhance specificity of sequence representations**

98      Machine learning-based TCR virtual screening for unseen peptides is challenging
99 because, by definition, the interaction data between TCRs and the unseen peptides is
100 limited. To address the limitation, we trained two distinct language models tailored for
101 TCR (CDR3 region of the β chain) and peptide sequences, with the aim to incorporate
102 additional TCR and peptide sequence representations to better model the inherent
103 relationships between these sequences. It is worth noting that there are already large
104 language models designed for general protein sequences, e.g., ESM[21,22]. However, we
105 posited that TCR and peptide sequences exhibit unique sequence patterns that may need
106 specific model training.

107      In **Figure 2A**, we show a head-to-head comparison of perplexities of the ESM-2
108 model (t33_650M_UR50D) and the TCR LM trained from scratch by TCRfinder on 884
109 validation TCR sequences from VDJdb[23] that have a maximum sequence identity of 80%
110 to the training TCRs (see **Text S1** in **Supporting Information**, **SI**). Here, the (pseudo)
111 perplexity of a LM for a protein sequence is computed by

$$\text{Pseudo Perplexity} = \exp\left(\frac{1}{L}\sum_{l=1}^{L} -logp\left(x_l|x_{/l}\right)\right) \#(1)$$

112 where $L$ is the sequence length, $x_l$ the masked residue at position $l$, and $x_{/l}$ the context.
113 The value of the perplexity can serve as a metric to quantify the mean uncertainty of the
114 LM, reflecting how well the LM is able to predict the masked token in the sequence. A
115 lower perplexity indicates that the model is more certain and accurate in its predictions,
116 approaching the ideal scenario of one. In contrast, a higher perplexity suggests higher
117 uncertainty and a less effective model, with values approaching the number of unique
118 tokens, indicating randomness in predictions. Please note that for the evaluation of the
119 TCR LM, we specifically computed the perplexity of the 7 center amino acids of the CDR3

120    region, as these central regions are known to be more flexible.

121         It is shown that the TCR LM trained in TCRfinder achieves an average perplexity of

122    14.289, significantly lower than the 18.321 by ESM-2, with a P-value of 4.98E-100 in

123    paired one-sided Student's t-test. Out of the 884 test cases, 78.96% (698 cases) exhibited a

124    lower perplexity with TCR LM, compared to only 21.04% (186 cases) for ESM-2. The

125    superiority of the TCR LM over ESM-2 is understandable because ESM-2 has been trained

126    on general protein sequences and is therefore less sensitive to the specific patterns present

127    in TCR sequences. These results also indicate that TCR sequence patterns significantly

128    diverge from general protein sequence patterns, highlighting the need for a customized

129    strategy when dealing with the specific patterns of TCR sequences.

130         In **Figure 2B**, we also present a perplexity comparison of ESM-2 and the peptide LM

131    trained in TCRfinder on 1000 validation peptides (see **Text S1**), where the specially trained

132    peptide LM demonstrates again a significantly lower perplexity (17.055) than ESM-2

133    (18.758) with a P-value of 6.54E-40 in a one-sided Student's t-test. Furthermore, the

134    peptide LM achieves a lower perplexity in 755 out of 1000 cases (75.50%), suggesting its

135    higher efficacy in capturing peptide sequence patterns compared to ESM-2.

136

137    **Impacts of language models on peptide specific TCR virtual screening**

138         One primary goal of training TCR and peptide LMs is to enhance the modeling

139    accuracy of TCR and peptide associations. To examine the impacts of these specially

140    trained LMs on TCR virtual screening, we present in **Figure 2C** a comparison of

141    ROCAUC achieved by TCRfinder with and without using LMs on 58 peptides randomly

142    selected from the validation dataset (**Text S2**). Here, ROCAUC stands for 'Areas under the

143    Receiver Operating Characteristic curve' computed in the classification experiment

144    distinguishing interacting from non-interacting TCRs for the given antigenic peptides. For

145    the case without LMs, TCRfinder models are trained with LM matrices replaced by the

146    query sequences only in the Joint Embedder module shown in **Figure 1D**.

147         It is shown that the incorporation of LMs significantly improves the screening

148    performance, with the average ROCAUC of full TCRfinder (0.711) being 5.2% higher than

149    that without using LMs (0.676). The ablation experiment shows that the major impact

150    comes from the TCR LM, as the ROCAUC of TCRfinder after removing TCR LM is 4.4%

151    lower than the full TCRfinder while removing peptide LM reduces ROCAUC only by

152    0.3%. The varied impact of LMs on the performance of TCR virtual screening might stem

153    from the lower specificity of the peptide language models, especially for the unseen

154    peptides in this validation dataset. This result is consistent with the data in **Figures 2A** and

155    **2B**, where the perplexity of the TCR LM (14.289) is considerably lower than that of the

156    peptide LM (17.055), indicating better sequence pattern recognition power of the TCR

157    over peptide LM.

158         In **Figure 2C**, we also list the average ROCAUC value (0.688) of TCRfinder when

159 replacing specially trained LMs by ESM-2 models, which is considerably lower than that
160 of original TCRfinder models (0.711). This result demonstrates again the importance of
161 using specially trained LMs over general protein LMs, even though the latter has been
162 trained in a much larger sequence space with significantly higher dimension of models[21,22]
163 (8.7 million short-length vs ~50 million full-length sequences with 6.1 million vs 651
164 million parameters for the TCRfinder and ESM-2 LMs, respectively).

165

166 **TCRfinder outperforms previous methods for TCR virtual screening**

167 To systematically examine the ability of TCRfinder in TCR virtual screening, we
168 collected a set of 38 nonredundant peptides from the VDJdb, each with 1-239 interacting
169 TCRβ CDR3 sequences unseen in any of the TCRfinder training datasets (see **Text S2**).
170 **Figure 3** summarizes TCR virtual screening results of TCRfinder, in control with six
171 state-of-the-art third-party methods, including ATM-TCR[16], PanPep[18], ERGO[13,14],
172 AttnTAP[17], pMTnet[12] and DLpTCR[15], which are all installed in our local computers and
173 implemented with default parameters.

174 In **Figure 3A**, it is observed that most methods exhibit limitations in effective TCR
175 recognitions for unseen peptides, where five control methods (PanPep, ERGO, AttnTAP,
176 pMTnet, and DLpTCR) generated average ROCAUC values of 0.534, 0.531, 0.504, 0.536,
177 and 0.401, respectively, marginally surpassing or even falling below a random guess
178 (ROCAUC = 0.5). Although ATM-TCR produces a considerably higher ROCAUC= 0.699
179 than other control methods, TCRfinder achieved the highest ROCAUC of 0.763, marking a
180 9.2% increase over ATM-TCR. **Figure 3B** provides a detailed head-to-head comparison of
181 peptide-wise ROCAUC between TCRfinder and control methods, where TCRfinder
182 achieves a higher ROCAUC in 28 cases compared to ATM-TCR, whereas ATM-TCR
183 accomplishes this in 10 cases. The numbers are 27/11, 32/6, 27/11, 29/9, and 32/6,
184 compared to PanPep, ERGO, AttnTAP, pMTnet and DLpTCR, respectively.

185 Notably, TCRfinder has a median ROCAUC of 0.912, with majority of test peptides
186 scoring between 0.8 and 1.0, representing a noteworthy 20.5% improvement over
187 ATM-TCR that has a median ROCAUC of 0.757. The superiority of TCRfinder in median
188 ROCAUC is particularly meaningful to experimentalists, as such a high median ROCAUC
189 suggests that with TCRfinder screening models, for half of the unseen peptides,
190 experimentalists only need to validate, on average, less than 10% (equals to 1-0.912) of the
191 sample pools to identify the true interacting TCRs. The result highlights the practical
192 usefulness of effective computational screening in helping significantly reduce the time
193 and cost of tumor specific TCRs recognitions and immunotherapies.

194 In **Figure 3C**, we extend our analysis to the comparisons of Enrichment Factors (EFs)
195 which is defined as

196

$$EF_{x\%} = \frac{N^{x\%}_{interact}/N^{x\%}_{select}}{N_{interact}/N_{select}} \#(2)$$

197  where $N_{interact}$ and $N_{select}$ are the total number of interacting and all TCRs in the
198  screening pool, respectively. $N^{x\%}_{interact}$ and $N^{x\%}_{select}$ are the numbers of true interacting
199  TCRs and the number of all candidates in the top-$x\%$ of the TCRs selected on the scores
200  predicted by the corresponding methods. Generally, a higher EF value indicates better
201  performance and $EF = 1$ indicates random selection. In this work, a wide range of
202  top-$x\%$ cut-offs (i.e., 1%, 5%, 10%, 15%, 20%, 25% and 30%) are used to compare the
203  performance under different scenarios.

204  Consistent with the ROCAUC results, ATM-TCR produces consistently higher EFs
205  than other control methods, where the P-values between TCRfinder and ATM-TCR,
206  PanPep, ERGO, AttnTAP, pMTnet and DLpTCR are 4.54E-03, 2.25E-05, 9.46E-08,
207  9.01E-06, 7.87E-06 and 3.32E-10, respectively, across all the top-$x\%$ cut-offs. For the
208  top-1% sample cutoff, for instance, TCRfinder achieves an average EF=27.01, which is
209  63.2% higher than ATM-TCR (16.53) and >4 times higher than that of the other three
210  methods (PanPep, AttnTAP, and pMTnet) that have achieved better than random EF (i.e.,
211  EF > 1.0).

212  **Figure 3D** further lists head-to-head comparisons of EF across multiple cut-offs
213  between TCRfinder and control methods. Here, EF values are scaled to [0.0, 1.0] by the
214  theoretical maximum value of the corresponding top-$x\%$ cut-offs, i.e., 100/$x$, for improved
215  visualization. Out of 266 indices (i.e., 38 peptides with 7 top-$x\%$ cut-offs each), TCRfinder
216  consistently achieves higher or equal scaled EF in 238, 224, 238, 235, 256, and 255 cases
217  compared to ATM-TCR, PanPep, ERGO, AttnTAP, pMTnet, and DLpTCR, respectively.
218  In contrast, the control methods generate better or equal performance in at most 172 cases.
219  The associated P-values, with a maximum of 9.20E-07, for the scaled EF, further validate
220  the statistical significance of TCRfinder's superior performance across a spectrum of
221  enrichment cut-offs.

222

223  **Case studies reveal effectiveness of unsupervised learning for cross-chain structure**
224  **interactions**
225  To further examine the strength and weakness of TCRfinder, as case studies we look
226  into details of two peptides of '**MMWDRGLGMM**' and '**EVDPIGHLY**', which
227  represent two typical examples standing at the upper-left and bottom-right corners of
228  ROCAUC-homology plot **Figure S1**.
229  First, '**MMWDRGLGMM**' is a synthetic epitope presented by HLA-A*02[24,25] with
230  only one interacting TCR identified in VDJdb. The CDR3 sequence '*CASSLSFGTEAFF*'
231  of this TCR is composed of the V and J genes of TRBV6-4 and TRBJ1-1, respectively. As

232    shown in **Figure S2A**, the most similar peptide in the training dataset is 'NMMWFQGQL',
233    which is another synthetic epitope presented by the same MHC, with a sequence identity of
234    40% to the test peptide. However, the two known interacting TCRs for 'NMMWFQGQL',
235    with β-chain CDR3 sequences of '*CASSRDTVNTEAFF*' and '*CASSRDFVSNEQYF*', have
236    completely different J genes (TRBJ1-1 and TRBJ2-7). Thus, the high ROCAUC of 0.985
237    (Figure **S1)** is not due to the simple homologous transfer from the training dataset.

238        The TCR-peptide interacting complex structure from PDB (PDBID: 6AMU) is
239    depicted in **Figure 4A**, with a detailed inter-chain Cα distance map presented in **Figure 4B**.
240    Interestingly, the central region of the TCR-peptide distance map, which corresponds to
241    physical inter-residue interactions, shows a similar pattern to the mean attention map
242    which was automatically learned from solely TCR-peptide sequence data by TCRfinder.
243    Here, the attention map between the TCRβ CDR3 region and peptide sequences was
244    extracted directly from the transformer block following the Joint Embedder and represents
245    the mean of attention weights from each head. Additionally, focusing on the TCRβ CDR3
246    region, **Figure 4C** illustrates the minimum distance between each residue of TCRβ CDR3
247    and the peptide residues, superposed with the accumulated attention weights
248    corresponding to each residue of TCRβ CDR3. A significant correlation (PCC of 0.549)
249    can be observed between the two signals, despite no structural information being included
250    during the TCRfinder training. These findings suggest that the attention mechanism in
251    TCRfinder captures meaningful information regarding the spatial binding interface of
252    residues in the TCRβ CDR3 region and their interactions with the peptide, an ability
253    critical to the enhanced TCR-peptide interaction predictions.

254        To further evaluate the effectiveness of TCR LM in this context, a basic
255    nearest-neighbor approach is implemented by encoding TCRs using both TCR LM and the
256    basic one-hot encoding. Candidate TCRs are searched by comparing the representations of
257    interacting TCRs of the template peptide in VDJdb and ranked according to the distance of
258    representative vectors between two TCRs. In this approach, the residue-wise
259    representations of TCRs can be simply averaged along the sequence axis into a
260    fixed-length vector. The ranks of interacting TCR for TCR LM encoding and the basic
261    one-hot encoding are 473 and 1999, respectively, as shown in **Figure 4D**. This
262    demonstrates that TCR LM learns more informative representations than just residue types.
263    Leveraging TCR LM, along with further supervised training, TCRfinder successfully
264    identifies the interacting TCR at the 75[th] trial from a pool of 5000 background TCRs,
265    constituting only 1/3 of the second-best ranking achieved by PanPep (243[rd]) (**Figure 4E**).
266    In contrast, other control methods exhibit rankings exceeding 800 and, in some cases, >
267    3500. These results strongly affirm the robust capability of the proposed TCRfinder to
268    effectively screen interacting TCRs, not only for unseen peptides but also for peptides that
269    are distinctly novel and challenging.

270        In the second example of '**EVDPIGHLY**', TCRfinder achieved a poor ROCAUC of

271    0.042, even though a highly homologous peptide (MEVDPIGHLY) exists in the training
272    dataset with 100% sequence identity (or 90% if normalized by the length of
273    'MEVDPIGHLY'), and both peptides belong to the MAGE-A3 epitope.

274    As shown in **Figure S2B**, the interacting TCR for the training peptide has a β-chain
275    CDR3 sequence of '*CSANPRTTLYEQYF*', formed by V and J genes of TRBV20-1 and
276    TRBJ2-7, respectively. In contrast, the interacting TCR for the test peptide has a distinct
277    β-chain CDR3 sequence of '*CASSFNMATGQYF*', with V and J genes of TRBV5-1 and
278    TRBJ2-7, respectively. The notable differences in TCRs despite the similar peptide
279    sequence make the virtual screening challenging. Furthermore, 84.14% of the randomly
280    selected background TCRs exhibited a higher sequence similarity to the positive TCR in
281    the training set. Consequently, TCRfinder failed to achieve a reasonable ranking for the
282    interacting TCR.

283    Notably, in VDJdb, the test and training peptides are presented by different MHCs, i.e.,
284    HLA-A*01 and HLA-B*18, respectively. As data accumulates, we aim to incorporate
285    more metadata (e.g., species and MHC types) and peripheral protein environments into
286    consideration for such cases when they are available.

287    Despite these challenges, it's noteworthy that for 8 peptides with a maximum
288    sequence identity exceeding 80%, 6 (75%) of them demonstrated high ROCAUC values of
289    at least 0.906 and an average rank of 202.23 out of 5000 candidates by TCRfinder (**Figure**
290    **S1**).

291

### Recognition of neoantigens against wildtype peptides

293    Neoantigens are newly formed antigen peptides generated by mutations in the DNA
294    of cancer cells. A critical step in T-cell based cancer therapeutics is the identification of
295    TCRs that can recognize neoantigens presented by MHC proteins[4,26]. However, since the
296    difference between a neoantigen and corresponding wild-type peptide often involves only
297    a single residue mutation, neoantigen-reactive TCRs might also interact the corresponding
298    wild-type peptide. Such interactions could lead to serious autoimmune disorders, creating
299    an urgent need to identify neoantigen-reactive TCRs that are less likely to interact with
300    wildtype peptides. But the high sequence similarity between between neoantigen and
301    wildtype peptides often renders the computational screening of neoantigens for given TCR
302    highly challenging.

303    The default TCRfinder model cannot be directly applied to neoantigen screening, as
304    the loss function focuses solely on the ranking between two TCRs for a given peptide (see
305    **Eq. 5** in **METHODS**). To address this, we developed an additional model specifically
306    focused on scoring interactions between diverse peptides for a given TCR as described in
307    '**Alternative model trained for TCR-based neoantigen screening**' in **METHODS**. To
308    quantitatively assess the capability of the model, we gathered a total of 64 independent
309    TCR-peptide pairs, including 32 experimentally validated TCR-neoantigen interactions

310    and 32 corresponding TCR-wildtype non-interacting pairs, also validated
311    experimentally[27-29]. Here, we used the data from VDJdb for training, where all pairs
312    sharing the same TCRs or peptides to the 64 test TCR-peptide pairs have been excluded
313    from the training set to prevent data contamination (**Text S2**).

314          In **Figure 5A**, we present a summary of success rate of TCRfinder on the 64
315    TCR-peptide pairs in control with other six third-party programs, where a success case
316    refers to that in which the program correctly prioritizes the TCR-neoantigen interaction
317    with a higher score than the TCR-wildtype peptide pairs. As shown in **Figure S3**, all the
318    control methods demonstrate relatively low success rates ranging from 0.375 (ERGO) to
319    0.563 (DLpTCR), indicating the challenging nature of this task. In contrast, TCRfinder
320    achieved a remarkable success rate of 0.844, due to the special training strategy. As shown
321    in **Figure 5B**, the P-value between the binding scores with neoantigen and wildtype
322    peptides by TCRfinder is 1.34E-07 in paired one-sided Student's t-test (**Figure 5B**),
323    highlighting its robustness in distinguishing TCRs with distinct interactions between
324    neoantigens and wild-type peptides.

325          The 32 TCR-neoantigen interactions involve a total of 12 unique neoantigens. **Figure
326    5C** lists the details of score differences, $\Delta S = S_{neoantigen} - S_{wildtype-peptide}$, on the
327    interactions grouped by neoantigens. It can be observed that TCRfinder correctly
328    recognizes 11 out of 12 (91.67%) neoantigens with positive $\Delta S$, failing only in one
329    neoantigen, 'EPLHTPTIM'. Such a high success rate at the neoantigen level is critical for
330    guiding the selection of TCRs with lower "toxicity" for personalized immunotherapies.

331

332    **DISCUSSION**
333          Accurate modeling of TCR and antigenic peptide interactions is essential for
334    elucidating immune responses and developing T-cell based immunotherapies. By
335    harnessing specifically trained TCR and peptide language models, we introduced
336    TCRfinder, a novel deep learning method for precise TCR-peptide binding interaction
337    prediction. Through a careful examination of various performance metrics, we found the
338    robust capability of TCRfinder in distinguishing interacting and non-interacting TCRs for
339    unseen peptides, with accuracy significantly beyond current state-of-the-art methods.
340    Furthermore, through rigorous training and validation, TCRfinder demonstrated a
341    remarkable ability to distinguish TCRs with distinct interactions between neoantigens and
342    wild-type peptides, with a success rate nearly 50% higher than the best existing method.
343    This capability is particularly important to minimize toxic autoimmune responses in T-cell
344    based immunotherapies.

345          Several unique designs have contributed to the superior performance of the
346    TCRfinder. First, the specifically trained TCR and peptide language models exhibit higher
347    specificity (or lower perplexity) than general large protein language models, such as
348    ESM-2, allowing better modeling of evolutionary patterns in diverse TCRs (particularly

349 the sensitive β-chain CDR3 regions) and peptide sequences. Second, as demonstrated by

350 the case studies, the specially designed Joint Embedder network, when integrated with the

351 LMs, can reveal meaningful information regarding the spatial binding interface of residues

352 within the TCRβ CDR3 region, highlighting its potential to bridge the gap between

353 sequence-based predictions and structural insights. Third, although both TCR and

354 neoantigen screens rely on the specific TCR-peptide interactions, the procedures involve

355 different sampling backgrounds (i.e., one on multiple TCRs for a given peptide and another

356 on highly similar peptide sequence for a given TCR). Therefore, separate models trained

357 with specially designed loss functions can further improve the accuracy and specificity of

358 the screening experiments.

359 Nevertheless, there are still substantial limitations in the TCRfinder model. Frist,

360 despite careful training, the peptide LM is still much less specific than the TCR β-chain

361 CDR3 region LM and thus has a less significant impact on TCRfinder's overall

362 performance. One reason could be the relatively small training dataset, where the trained

363 LM cannot accurately capture the diversity of TCR-bound peptide sequence patterns.

364 Second, the current TCRfinder model has difficulty in recognizing the closely homologous

365 peptides that bind with distinct TCRs. This limitation may also relate to the limited

366 sensitivity of the specially trained peptide LMs. Collecting larger datasets of

367 peptide/neoantigen sequences and incorporating peripheral features such as HMC and

368 TCR sequences in the LM training might help address the issues.

369 In summary, our findings highlight the significance of leveraging language models

370 in TCR virtual screening, offering a powerful approach to predict antigen-specific TCRs

371 for unseen antigenic peptide. Such advancement holds promise for applications in

372 individualized TCR-related therapies.

373

374 **METHODS**

375 TCRfinder is a deep learning pipeline designed for sequence based TCR and

376 neoantigen screenings. The core strategy of the pipeline is the development of accurate

377 interaction scoring models from the sequences of the β-chain TCR CDR3 region and target

378 peptides, which consists of four steps (**Figure 1A**). Firstly, two language models (TCR LM

379 and peptide LM) are pre-trained for the TCR and peptide sequences separately. A joint

380 embedder model is then learned from the concatenated representations of the TCR and

381 peptide LMs. Next, transformer blocks are implemented iteratively on both paired and

382 individual TCR and peptide representations to create new TCR and peptide embedding

383 matrixes. Finally, a TCR-peptide interaction scoring model is obtained from a Multi-Layer

384 Perceptron (MLP) block layer based on the new TCR and peptide embeddings. To address

385 issues of TCR and neoantigen recognitions, two models on peptide-based TCR screening

386 and TCR-based neoantigen screening are trained separately.

387

**TCR and peptide language model training**

388
389 The neural network architecture of TCR and peptide LMs is presented in **Figure 1B**.
390 Given a masked TCR sequence, the embedder layer is employed to extract hidden
391 representations in the forms of both sequential and pair matrices. The detailed flowchart of
392 the Embedder Module is shown in **Figure S4**, where the input is the one-hot encoding of
393 the masked TCR or peptide sequence, and the resulting output contains sequential and pair
394 representations. Specifically, the sequential representation is derived through a linear
395 transformation over the input features, added with 1-D positional encoding. The pair
396 representation is the outer sum of two separate hidden features, each obtained by individual
397 linear transformations. Additionally, the pair representation is enriched by the inclusion of
398 2-D positional encoding.

399 Next, the sequential and pair representations serve as inputs of a set of transformer
400 blocks, designed to model interactions between residues in the sequence (see **Figure 1C** as
401 well as the description in '**Transformer Block with pair modeling**' below). While the
402 transformer block generates both sequential and pair representations, we only utilize the
403 last sequential representation for predicting the residue type of the masked regions.

404 To train the TCR LM, we collected all downloadable TCR sequences from TCRdb[30],
405 comprising a total of 7,259,306 TCR β-chain CDR3 sequences, while the validation set for
406 the TCR LM contains 884 TCRs included in VDJdb[23], which have the maximum sequence
407 identity of 80% to the training TCRs. In the case of the peptide LM, we collected all unique
408 linear peptides from IEDB[31] and randomly selected 1000 peptides for evaluation, with the
409 remainder allocated for training, where a maximum sequence identity cutoff of 80% is
410 maintained between the evaluation and training sets for peptide LM (see **Text S1** in **SI**).

411 For both TCR and peptide LMs, we set the number of blocks as 4 and 7, respectively.
412 The dimensions for sequential and pair representations in the TCR LM are set as 128 and
413 96, while for the peptide LM, they are set as 256 and 96, respectively. Two masking
414 policies are employed during training: (1) Randomly masking 10% of the residues; (2)
415 Randomly masking a continuous region with a length of approximately 10% of the
416 sequence length.

417 Both TCR and peptide LMs were trained using the Adam optimizer[32] with default
418 parameters for approximately 10 epochs. The training process was guided by minimizing
419 the cross-entropy loss of the masked regions given a sequence, as defined by:

$$L_{CE} = -\sum_{m \in M} \log p(x_m) \tag{3}$$

421 where $p(x_m)$ is the softmax probability for the $m$th masked residue and $M$ refers to the
422 masked residue set.

423

**Joint embedding of TCR and peptide sequences**

425 Built on the pre-trained LMs, TCRfinder employs the Joint Embedder module to
426 transform inputs into diverse representations. As illustrated in **Figure 1D**, the TCR and

427  peptide sequences are initially encoded using one-hot encoding. Each sequence undergoes
428  a linear transformation, accompanied by 1-D positional encoding injection (denoted as
429  'pos' in the architecture), resulting in hidden sequential representations. Simultaneously,
430  the one-hot encoding features are input into their respective LMs to derive sequential and
431  pair representations for TCR and peptide, respectively. Note that the representations from
432  LMs will be projected by the linear layers to the desired dimensions. The output sequential
433  representations are obtained by element-wise additions of (1) the linear layer with
434  positional encoding and (2) the corresponding language model's output after projection.
435  The pair representations after the projection will also be considered as the outputs of the
436  Joint Embedder.

437      On the other hand, the sequential representations of TCR and peptide will be fed into
438  two linear layers individually. The resultant sequential representations are concatenated
439  separately to create two joint sequential representations. A subsequent Outer Sum
440  operation is applied to these joint sequential representations, resulting in the joint pair
441  representation. The joint pair representation will have the 2-D positional encoded (denoted
442  as 'relpos' in the **Figure 1D** architecture). The dimension sizes for sequential and pair
443  representation are uniformly set to 64 and 48, respectively.

444

445  **Transformer Block with pair modeling**
446      The joint sequence and pair representations from the Joint Embedder will be
447  considered as the input of a single transformer block to capture interaction patterns. This
448  simplified design aligns with the philosophy that a compact deep learning model can
449  effectively learn intricate interaction patterns when transferring strong prior information
450  from, for example, pretrained language models. The architecture of the transformer block
451  is an extension inspired by the Evoformer in AlphaFold2[33], as detailed in **Figure 1C**.
452      The training process begins with sequential and pair representations being through a
453  sequence row-wise gated self-attention layer. Notably, the pair representation is
454  transformed and added as the bias term of the attention map, facilitating the effective
455  incorporation of relationships between residues in the output sequential representation. A
456  sequence transition layer, comprising two linear layers, follows the sequence self-attention.
457  This transition layer first expands the dimension and then projects it back to the original
458  channel size. The resultant sequence representation is subsequently transformed into a pair
459  representation through an outer product mean (OPM) block. The 2-D representation from
460  the OPM block undergoes a sequential progression through a series of blocks, comprising
461  (1) a triangle multiplicative update block using outgoing edges, (2) a triangle multiplicative
462  update block using incoming edges, (3) a triangle self-attention block around the starting
463  node, (4) a triangle self-attention block around the ending node, and (5) a pair transition
464  block. It is important to highlight that the sequence and pair blocks are stacked residually
465  for efficient and stable training.

466    The attention layer is configured with 8 heads, with each head having a channel size
467 of 16. In the sequence transition layer, the expand factor is set to 2, meaning that the
468 sequential representation is initially expanded to twice its original dimension. Within the
469 OPM Module, the sequential representation is first projected into two representations with
470 a channel size of 12, followed by the outer product operation to obtain a pair representation
471 with a channel size of 12×12=144. This 144-D pair representation is then projected down
472 to 64. For subsequent modules applied to this pair representation, the number of heads and
473 the channel size of each head are also set to 8 and 16, respectively.

474

**Asymmetric modeling of TCR and peptide representations**

476    The Joint Embedder and subsequent transformer block described above play crucial
477 roles in integrating information across both TCR and peptide representations. Following
478 this, the joint representations are spitted into TCR and peptide representations, each
479 directed into separate branches of transformer blocks. This establishes an asymmetric deep
480 learning structure for TCRs and peptides, characterized by the absence of parameters
481 sharing between the two transformer blocks and further emphasized by a distinct layer
482 configuration, i.e., $N_1$=2 layers for TCRs and $N_2$=1 layer for peptides (**Figure 1A**).

483    The asymmetry of the network is designed based on the empirical observation that the
484 number of unique TCRs significantly exceeds that of unique peptides in the training set.
485 Although alternative strategies, such as employing different hidden dimensions, are also
486 plausible for achieving asymmetry, we opt for simplicity by differing in the block number
487 and found that the latter design achieves slightly better performance compared to other
488 selections.

489

**Training of TCR-peptide interaction model with triplet loss function design**

491    Following the asymmetric modeling of TCR and peptide representations, the output
492 TCR and peptide sequence representations are concatenated and go through an average
493 pooling operation along the sequence axis, resulting a fixed-length vector. The final scores,
494 which can quantify the interactions between the input TCR and peptide sequences, are then
495 obtained through the final MLP module, which consists of 2 hidden layers with 126 and 64
496 neurons, respectively. (**Figure 1A**).

497    We have formulated a loss function based on the triplet loss for the training of
498 TCRfinder. In each training step, a triplet of sequences is randomly sampled from the
499 training set. This triplet consists of a peptide sequence (considered as the anchor), a TCR
500 sequence that positively interacts with the peptide (denoted as $TCR_p$), and a negative TCR
501 sequence ($TCR_n$). The TCRfinder model then computes the interaction scores between the
502 anchor peptide and the two TCRs, i.e., $s(peptide, TCR_p)$ and $s(peptide, TCR_n)$
503 respectively, which are the target of interaction models. The loss function for each triplet is
504 given by:

$$L_{triplet} = \max\left(0, m + s(peptide, TCR_n) - s(peptide, TCR_p)\right) \#(4)$$

505    where *m* is the hyper-parameter that describes the margin and was set to 1.0 when training

506    TCRfinder. Although the structure of the loss function in Eq. (3) resembles the previous

507    triplet loss network protocol[34,35], the learning in those triplet loss networks utilizes a given

508    form of distance between samples, whereas in TCRfinder, the interacting score between

509    difference sequences is unknown and trained as the target of the deep learning models.

510        In addition, inspired by the loss function utilized in pMTnet[12], we introduce a new

511    regularization term to Eq. (3) to regulate the scales of the scores. The ultimate form of the

512    loss function is expressed as:

$$L = L_{triplet} + \rho\sqrt{s(peptide, TCR_n)^2 + s\left(peptide, TCR_p\right)^2} \#(5)$$

513    where $\rho$ denotes the regularization factor and sets to 0.03 in the training of TCRfinder.

514        The training of TCRfinder's prediction model involves the utilization of the Adam

515    optimization method and the loss function defined in Eq (5), implemented with PyTorch. In

516    each sampling step, one interacting and 10 non-interacting TCR sequences are randomly

517    selected for a given peptide sequence. This sampling strategy has been observed to yield

518    slightly improved performance compared to sampling across all possible pairs. The

519    rationale behind this approach is the uneven distribution of interacting TCRs for each

520    peptide. Directly sampling over all pairs may lead to overfitting on peptides with a larger

521    number of interacting TCRs, and the adopted strategy aims to mitigate such potential

522    biases.

523        The loss function for this specific sampling strategy is calculated as the average of ten

524    triplet losses in a sampling step. In our formulation, we treat the sampling of one peptide

525    and its interacting TCRs (including 10 negative TCRs) as a single sample, and one batch

526    contains a fix number of 16 samples. Gradient accumulation has been utilized to reduce

527    GPU memory consumption. Each of the 5 TCRfinder sub models has been trained for

528    10,000 batches and the early stopping mechanism is implemented to identify and select the

529    optimal models.

530

531    **Alternative model trained for TCR-based neoantigen screening**

532        To address the issue of TCR-based neoantigen screening, we have trained a separate

533    model for recognizing mutant neoantigens from wildtype peptides under the same TCRs.

534    For this, we have made straightforward adjustments to the sampling and loss function in

535    the training process of the TCRfinder model, to enable peptide ranking given a TCR on the

536    same dataset.

537        Each sample in the training set contains one TCR, one interacting peptide, and 10

538    non-interacting peptides randomly selected from the database. Consequently, the loss

539    function for peptide ranking is formulated as:

$$L = \max\left(0, m + s(TCR, peptide_n) - s(TCR, peptide_p)\right)$$

540

$$+\rho\sqrt{s(TCR, peptide_n)^2 + s(TCR, peptide_p)^2} \#(6)$$

541  where $s(TCR, peptide_p)$ and $s(TCR, peptide_n)$ are the predicted scores between the
542  TCR and interacting and non-interacting peptides, respectively. We trained a total of 5
543  models over 5,000 batches, implementing an early stopping mechanism to identify and
544  select optimal models. The final score is the average score of the 5 models.
545

546  **DATA AVAILABILITY**
547  The datasets collected and used in this work are available at
548  https://zhanggroup.org/TCRfinder/dataset. We show structures of 6AMU obtained by
549  four-digit accession codes in the PDB repository (https://www.rcsb.org/). TCR and peptide
550  pairs for training were collected from VDJdb (https://vdjdb.cdr3.net/). The TCR and
551  peptide sequences for training language models were downloaded from TCRdb
552  (https://guolab.wchscu.cn/TCRdb//#/download) and IEDB (https://www.iedb.org/)
553  respectively.
554

555  **CODE AVAILABILITY**
556  The online server and standalone package of TCRfinder are freely available at
557  https://zhanggroup.org/TCRfinder and https://zhanggroup.org/TCRfinder/download,
558  respectively. Data were analyzed using Numpy v.1.20.3
559  (https://github.com/numpy/numpy), SciPy v.1.7.1 (https://www.scipy.org/), and Matplotlib
560  v.3.4.3 (https://github.com/matplotlib/matplotlib). Structures were visualized by Pymol
561  v.2.3.0 (https://github.com/schrodinger/pymol-open-source).
562

572  **AUTHOR CONTRIBUTIONS**
573  Y.Z. conceived the project and designed the experiments; Y.L. developed methods
574  and designed and performed experiments; C.Z. participated in discussion and helped
575  design the experiments; X.Z. constructed the online server. Y.L. wrote the initial
576  manuscript; all authors proofread and approved the final manuscript.
577

578 **COMPETING INTERESTS**
579     The authors declare no competing interests.
580
581

## Reference

1. Altman, J. D. *et al.* Phenotypic Analysis of Antigen-Specific T Lymphocytes. *Science* **274**, 94-96, doi:10.1126/science.274.5284.94 (1996).

2. Zhang, S.-Q. *et al.* High-throughput determination of the antigen specificities of T cell receptors in single cells. *Nature Biotechnology* **36**, 1156-1159, doi:10.1038/nbt.4282 (2018).

3. Kula, T. *et al.* T-Scan: A Genome-wide Method for the Systematic Discovery of T Cell Epitopes. *Cell* **178**, 1016-1028.e1013, doi:10.1016/j.cell.2019.07.009 (2019).

4. Baulu, E., Gardet, C., Chuvin, N. & Depil, S. TCR-engineered T cell therapy in solid tumors: State of the art and perspectives. *Sci Adv* **9**, eadf3700, doi:10.1126/sciadv.adf3700 (2023).

5. Dash, P. *et al.* Quantifiable predictive features define epitope-specific T cell receptor repertoires. *Nature* **547**, 89-93, doi:10.1038/nature22383 (2017).

6. Sidhom, J.-W., Larman, H. B., Pardoll, D. M. & Baras, A. S. DeepTCR is a deep learning framework for revealing sequence concepts within T-cell repertoires. *Nature Communications* **12**, 1605, doi:10.1038/s41467-021-21879-w (2021).

7. Jokinen, E., Huuhtanen, J., Mustjoki, S., Heinonen, M. & Lähdesmäki, H. Predicting recognition between T cell receptors and epitopes with TCRGP. *PLOS Computational Biology* **17**, e1008814, doi:10.1371/journal.pcbi.1008814 (2021).

8. Gielis, S. *et al.* Detection of Enriched T Cell Epitope Specificity in Full T Cell Receptor Sequence Repertoires. *Frontiers in Immunology* **10**, doi:10.3389/fimmu.2019.02820 (2019).

9. Montemurro, A. *et al.* NetTCR-2.0 enables accurate prediction of TCR-peptide binding by using paired TCRα and β sequence data. *Communications Biology* **4**, 1060, doi:10.1038/s42003-021-02610-3 (2021).

10. Breiman, L. Random Forests. *Machine Learning* **45**, 5-32, doi:10.1023/A:1010933404324 (2001).

11. Rasmussen, C. E. in *Advanced Lectures on Machine Learning: ML Summer Schools 2003, Canberra, Australia, February 2 - 14, 2003, Tübingen, Germany, August 4 - 16, 2003, Revised Lectures* (eds Olivier Bousquet, Ulrike von Luxburg, & Gunnar Rätsch) 63-71 (Springer Berlin Heidelberg, 2004).

12. Lu, T. *et al.* Deep learning-based prediction of the T cell receptor–antigen binding specificity. *Nature Machine Intelligence* **3**, 864-875, doi:10.1038/s42256-021-00383-2 (2021).

13. Springer, I., Besser, H., Tickotsky-Moskovitz, N., Dvorkin, S. & Louzoun, Y. Prediction of Specific TCR-Peptide Binding From Large Dictionaries of TCR-Peptide Pairs. *Frontiers in Immunology* **11**, doi:10.3389/fimmu.2020.01803 (2020).

14. Springer, I., Tickotsky, N. & Louzoun, Y. Contribution of T Cell Receptor Alpha and Beta CDR3, MHC Typing, V and J Genes to Peptide Binding Prediction. *Frontiers in Immunology* **12**, doi:10.3389/fimmu.2021.664514 (2021).

15. Xu, Z. *et al.* DLpTCR: an ensemble deep learning framework for predicting immunogenic peptide recognized by T cell receptor. *Briefings in Bioinformatics* **22**, bbab335, doi:10.1093/bib/bbab335 (2021).

621  16  Cai, M., Bang, S., Zhang, P. & Lee, H. ATM-TCR: TCR-Epitope Binding Affinity Prediction
622      Using a Multi-Head Self-Attention Model. *Frontiers in Immunology* **13**,
623      doi:10.3389/fimmu.2022.893247 (2022).
624  17  Xu, Y. *et al.* AttnTAP: A Dual-input Framework Incorporating the Attention Mechanism for
625      Accurately Predicting TCR-peptide Binding. *Frontiers in Genetics* **13**,
626      doi:10.3389/fgene.2022.942491 (2022).
627  18  Gao, Y. *et al.* Pan-Peptide Meta Learning for T-cell receptor–antigen binding recognition.
628      *Nature Machine Intelligence* **5**, 236-249, doi:10.1038/s42256-023-00619-3 (2023).
629  19  Grazioli, F. *et al.* On TCR binding predictors failing to generalize to unseen peptides.
630      *Frontiers in Immunology* **13**, doi:10.3389/fimmu.2022.1014256 (2022).
631  20  Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. in *North American Chapter of the
632      Association for Computational Linguistics.*
633  21  Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language
634      model. *Science* **379**, 1123-1130, doi:10.1126/science.ade2574 (2023).
635  22  Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised learning
636      to 250 million protein sequences. *Proceedings of the National Academy of Sciences* **118**,
637      e2016239118, doi:10.1073/pnas.2016239118 (2021).
638  23  Shugay, M. *et al.* VDJdb: a curated database of T-cell receptor sequences with known
639      antigen specificity. *Nucleic Acids Research* **46**, D419-D427, doi:10.1093/nar/gkx760 (2018).
640  24  Hellman, L. M. *et al.* Improving T Cell Receptor On-Target Specificity via Structure-Guided
641      Design. *Molecular Therapy* **27**, 300-313, doi:10.1016/j.ymthe.2018.12.010 (2019).
642  25  Riley, T. P. *et al.* T cell receptor cross-reactivity expanded by dramatic peptide–MHC
643      adaptability. *Nature Chemical Biology* **14**, 934-942, doi:10.1038/s41589-018-0130-4
644      (2018).
645  26  Lang, F., Schrors, B., Lower, M., Tureci, O. & Sahin, U. Identification of neoantigens for
646      individualized therapeutic cancer vaccines. *Nat Rev Drug Discov* **21**, 261-282,
647      doi:10.1038/s41573-021-00387-y (2022).
648  27  Tran, E. *et al.* T-Cell Transfer Therapy Targeting Mutant KRAS in Cancer. *New England
649      Journal of Medicine* **375**, 2255-2262, doi:10.1056/NEJMoa1609279 (2016).
650  28  Lu, Y.-C. *et al.* An Efficient Single-Cell RNA-Seq Approach to Identify
651      Neoantigen-Specific T Cell Receptors. *Molecular Therapy* **26**, 379-389,
652      doi:10.1016/j.ymthe.2017.10.018 (2018).
653  29  Arnaud, M. *et al.* Sensitive identification of neoantigens and cognate TCRs in human solid
654      tumors. *Nature Biotechnology* **40**, 656-660, doi:10.1038/s41587-021-01072-6 (2022).
655  30  Chen, S.-Y., Yue, T., Lei, Q. & Guo, A.-Y. TCRdb: a comprehensive database for T-cell
656      receptor sequences with powerful search function. *Nucleic Acids Research* **49**, D468-D474,
657      doi:10.1093/nar/gkaa796 (2021).
658  31  Vita, R. *et al.* The Immune Epitope Database (IEDB): 2018 update. *Nucleic Acids Res* **47**,
659      D339-D343, doi:10.1093/nar/gky1006 (2019).

660    32    Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *CoRR* **abs/1412.6980**
661          (2014).
662    33    Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,
663          583-589, doi:10.1038/s41586-021-03819-2 (2021).
664    34    Schroff, F., Kalenichenko, D. & Philbin, J. Facenet: A unified embedding for face
665          recognition and clustering. *Proceedings of the IEEE conference on computer vision and*
666          *pattern recognition*, 815-823 (2015).
667    35    Zhu, Y. H. *et al.* TripletGO: Integrating Transcript Expression Profiles with Protein
668          Homology Inferences for Gene Function Prediction. *Genomics Proteomics Bioinformatics*
669          **20**, 1013-1027, doi:10.1016/j.gpb.2022.03.001 (2022).
670
671

672 **FIGURES**

673



674

675 **Figure 1. Flowchart of TCRfinder for sequence based TCR and peptide interaction**

676 **predictions using deep learning networks.** (A) Overview of the TCRfinder pipeline. (B)

677 Detailed architecture of TCR and peptide language models. (C) Details of transformer

678 block with pair modeling, utilized in both TCRfinder model and TCR and peptide language

679 models. (D) Detailed architecture of Joint Embedder module to encode TCR and peptide

680 sequences into sequential and pair representations.

681

**Figure 2. Comparative Analysis of Language Models in TCRfinder.** (A) Head-to-head perplexity comparison of ESM-2 and TCR LM models on TCR β-chain CDR3 sequences. (B) Head-to-head perplexity comparison of ESM-2 and specially trained peptide LM models on peptide sequences. (C) ROCAUC of TCR screening experiments on the validation dataset with and without using language models.

689



690

**Figure 3. The comparison of TCRfinder with the control methods in TCR screening on 38 unseen peptides.** (A) Distribution of ROCAUC with the central mark indicating the mean. (B) Head-to-head ROCAUC comparison. (C) Enrichment Factors (EFs) versus percentage of sampling. (D) Head-to-head comparisons of EFs that are scaled to [0.0, 1.0] by the theoretical maximum EF value at corresponding top-$x$% cut-offs.
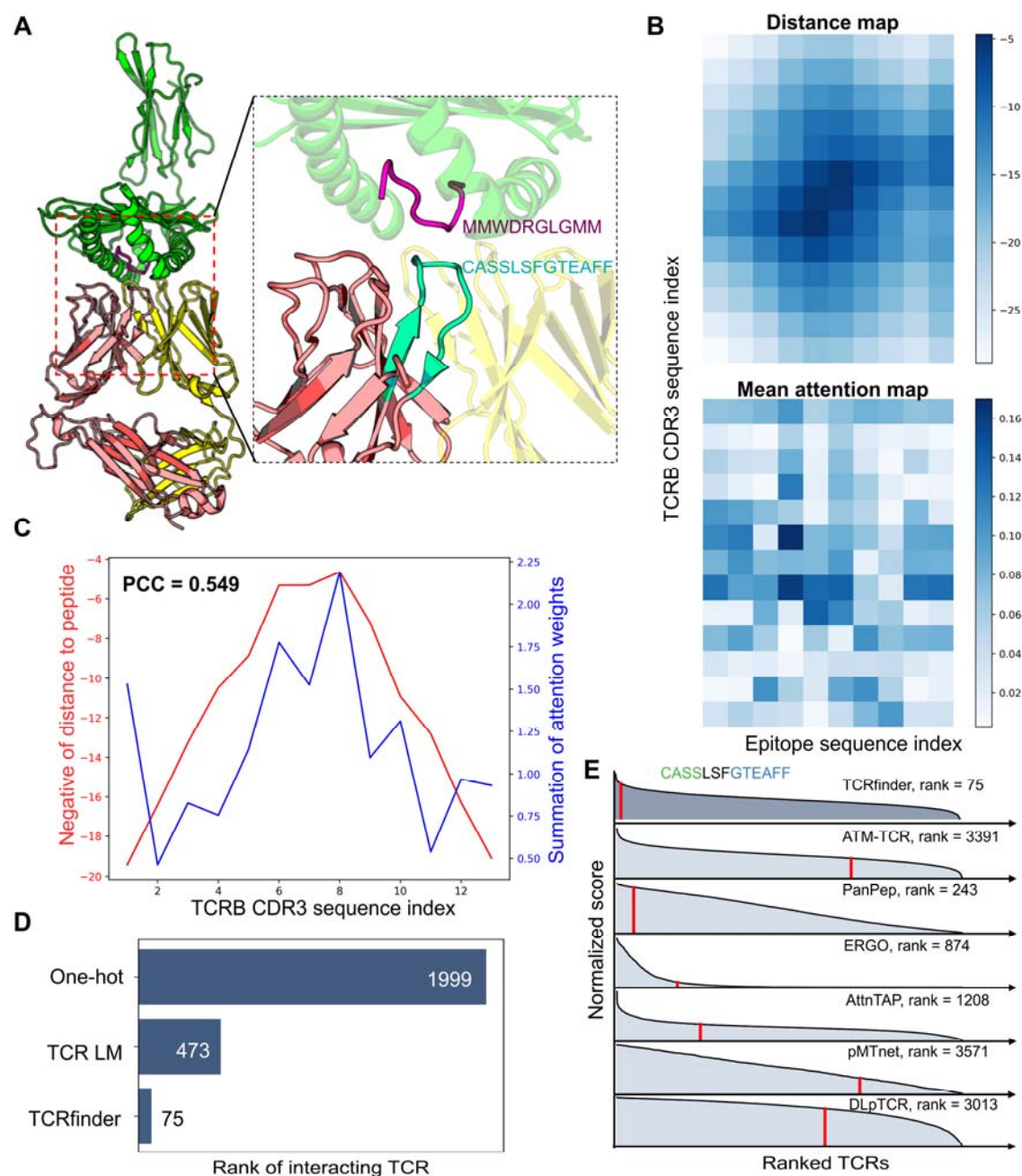
696

697

**Figure 4. Case study of TCR screening for peptide 'MMWDRGLGMM'.** (A) Experimental structure of TCR-peptide interacting complex with TCRβ CDR3 region and the peptide highlighted in cyan and purple, respectively. (B) Distance map between interchain Cα atoms, attention map between TCRβ CDR3 region and peptide sequences. (C) Residue wise minimal distance and attention weights to the peptide. (D) Comparison of TCR LM encoding and one-hot encoding in a nearest-neighbor approach for the studied case. (E) identification of the interacting TCR for the case peptide by TCRfinder and other methods.
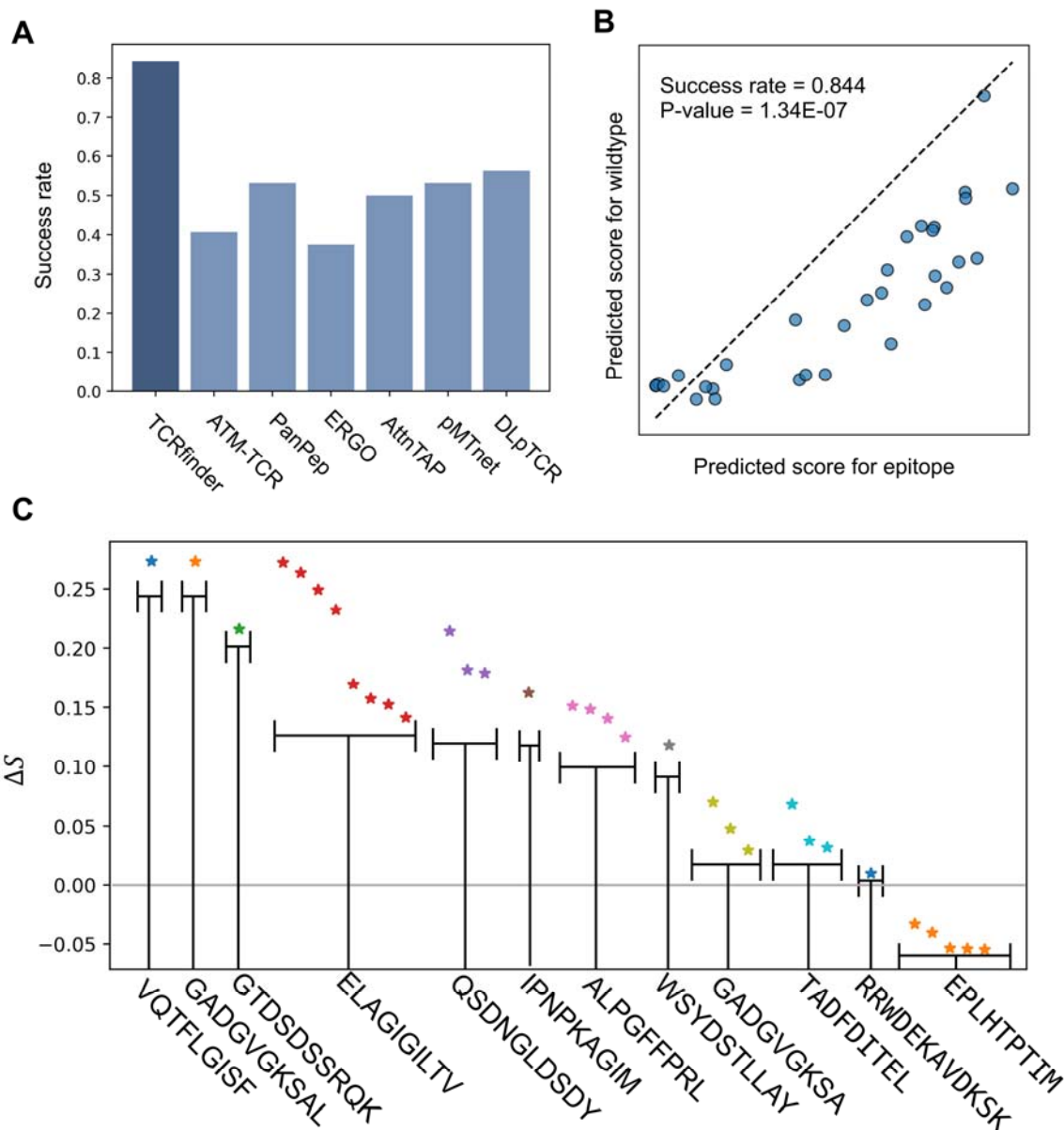
706

**Figure 5. Test results of neoantigen recognition on 64 TCR-peptide pairs.** (A) Success rate of different programs in correctly recognizing TCR-neoantigen interactions. (B) Predicted interaction score for TCR-neoantigen versus that for TCR-wildtype peptide by TCRfinder. (C) Score difference between TCR-neoantigen and TCR-wildtype pairs by TCRfinder as grouped by distinct neoantigens.