

Deep autoregressive generative models capture the intrinsics embedded in T cell receptor repertoires

Yuepeng Jiang¹ and Shuai Cheng Li^{1*}

¹Department of Computer Science, City University of Hong Kong, Kowloon Tong, Hong Kong

*shuaicli@cityu.edu.hk

ABSTRACT

T cell receptors (TCRs) play an essential role in the adaptive immune system. Probabilistic models for TCR repertoires can help decipher the underlying complex sequence patterns and provide novel insights into understanding the adaptive immune system. In this work, we develop TCRpeg, a deep autoregressive generative model to unravel the sequence patterns of TCR repertoires. TCRpeg outperforms state-of-the-art methods in estimating the probability distribution of a TCR repertoire, boosting the accuracy from 0.672 to 0.906 measured by the Pearson correlation coefficient. Furthermore, with promising performance in probability inference, TCRpeg improves on a range of TCR-related tasks: revealing TCR repertoire-level discrepancies, classifying antigen-specific TCRs, validating previously discovered TCR motifs, generating novel TCRs, and augmenting TCR data. Our results and analysis highlight the flexibility and capacity of TCRpeg to extract TCR sequence information, providing a novel approach to decipher the complex immunogenomic repertoires.

Introduction

The adaptive immune system consists of highly diverse B and T cells whose unique receptors can recognize enormous pathogens in vertebrates. The generation of these highly diverse receptors arises mainly from the genetic recombination of DNA segments from V, D, and J genes through V(D)J recombinations^{1,2}. T cells play an essential role in antiviral defense by selectively eliminating virus-infected cells³. Their ability to recognize specific short peptides; that is, peptide antigens bound to the major histocompatibility complex (MHC) molecules, is primarily determined by their unique receptor proteins^{4,5}. A receptor

1 contains an α polypeptide chain and an β polypeptide chain, both of which consist of two extracellular
2 domains: the variable (V) region and the constant (C) region⁶. The variable regions of the TCR α - and β -
3 chains both have three complementarity-determining regions (CDRs) that contribute to the specificity of
4 antigen recognition. Among these CDRs, the CDR3 region of the TCR β chain plays a pivotal role in the
5 recognition of the peptides presented by MHC. In contrast, the CDR1 and CDR2 regions contribute minor
6 effects to direct antigen recognition^{6,7}. Due to the importance of the highly diverse CDR3 region of the
7 TCR β chain in antigen recognition and data availability, this work focuses on deciphering the underlying
8 pattern of the CDR3 sequence.

9 Advancement in high-throughput sequencing techniques of the T cell receptor repertoire provides a
10 census of T cells found in blood or tissue samples^{8–11}. Large-scale sequencing data promote the investiga-
11 tion of the composition of immune repertoires, characterizing adaptive immune responses, and developing
12 descriptive models. The sampled repertoire of TCR serves as an indicator of the complete repertoire, re-
13 flecting the pathogenic history or the immune response to stimuli^{12–15}, with clinical applications including
14 cancer prediction and anticipation of immunotherapy. For example, Han *et al.* developed a statistical index
15 named TIR index based on TCR to predict the response and survival outcomes after immunotherapy¹⁶.
16 Beshnova *et al.* defined a cancer score for a given patient based on the predictive model trained on specific
17 TCR sequences that are assumed to be simply associated with cancers¹⁷.

18 Despite the success in predictive tasks associated with T cell repertoires, precise probabilistic distribu-
19 tion modeling is demanding to characterize the sequence information of a given repertoire sample, with
20 many potential applications such as estimating the relative ratio of CD4⁺ to CD8⁺^{18,19} and investigating
21 the differences in sequence characteristics between functional T cell subsets^{20,21}. Conventionally, model-
22 ing the sequence pattern behind a TCR repertoire is disentangled into two processes: generation (V(D)J
23 recombination)^{22,23} and selection^{19,24,25}. The ultimate probability assigned to a TCR sequence is the
24 product of the selection factor and the generation probability inferred from the selection process and the
25 generation process, respectively. However, the generation models learned from different individuals share
26 high mutual similarity^{22,25}, indicating that the selection process plays a central role in discriminating the
27 TCR repertoires sampled from different individuals. Therefore, instead of two-step disentanglement, we
28 can infer the probability of TCR sequences end-to-end.

1 In this work, we introduce a new software tool, TCRpeg, that utilizes deep learning techniques to learn
2 the underlying sequence patterns of TCR repertoires. Specifically, TCRpeg employs the architecture of
3 the deep autoregressive model with gated recurring units (GRU)²⁶ layers to characterize the repertoire
4 through the flexible and non-linear structure of deep neural networks. TCRpeg can infer the sequence
5 probability distribution with higher accuracy than other probabilistic models, boosting the performance
6 from 0.672 to 0.906 measured by the Pearson correlation coefficient. We applied the model to study type 1
7 diabetes (T1D) and found that TID samples have their iLN *CD8⁺* subrepertoires significantly diverged
8 from healthy controls with $D_{JS} \simeq 0.082$; while their pLN and spleen *CD8⁺* subrepertoires are similar to
9 healthy controls with $D_{JS} \simeq 0.018$ and $D_{JS} \simeq 0.025$, respectively. TCRpeg also provides high-quality
10 latent vector representations for TCR sequences. Based on these vector encodings of TCR sequences,
11 we build a fully connected neural network to classify the cancer-associated TCRs and SARS-CoV-2
12 epitope-specific TCRs, achieving 0.844 and 0.872 AUC, respectively; higher than DeepCAT¹⁷'s AUC
13 0.768 but slightly lower than TCRGP²⁷'s AUC 0.882. As a generative model, TCRpeg can generate new
14 TCR sequences, among which more than 50% of them share the same antigen specificity as the sequences
15 used in training according to the TCRMatch²⁸ software with a scoring threshold of 0.90, while the other
16 two generative models, TCRvae²⁹ and soNNia¹⁹, only achieve such a proportion of less than 40%. Further,
17 TCRpeg helps data augmentation; it shows a 7.4% accuracy gain in predicting cancer-associated TCRs
18 using the DeepCAT¹⁷ model.

19 Results

20 Autoregressive generative model for TCR sequences

21 Previously, the probabilistic sequence pattern of a TCR repertoire was modeled by either the disentangled
22 two processes of generation^{22,23} and selection^{19,24,25} (e.g., soNNia¹⁹) or the variational autoencoder
23 with convolutional neural networks (CNNs) as encoder and decoder (TCRvae²⁹). Although both models
24 achieved satisfactory performance, they lack the elegance to handle variable-length TCR sequence data.
25 The two types of models pad each sequence to a fixed length with an extra token representing the
26 padding positions. However, introducing the extra token could introduce noise to the original data and
27 partially conceal useful information from the diversity of sequence lengths, which is important for antigen

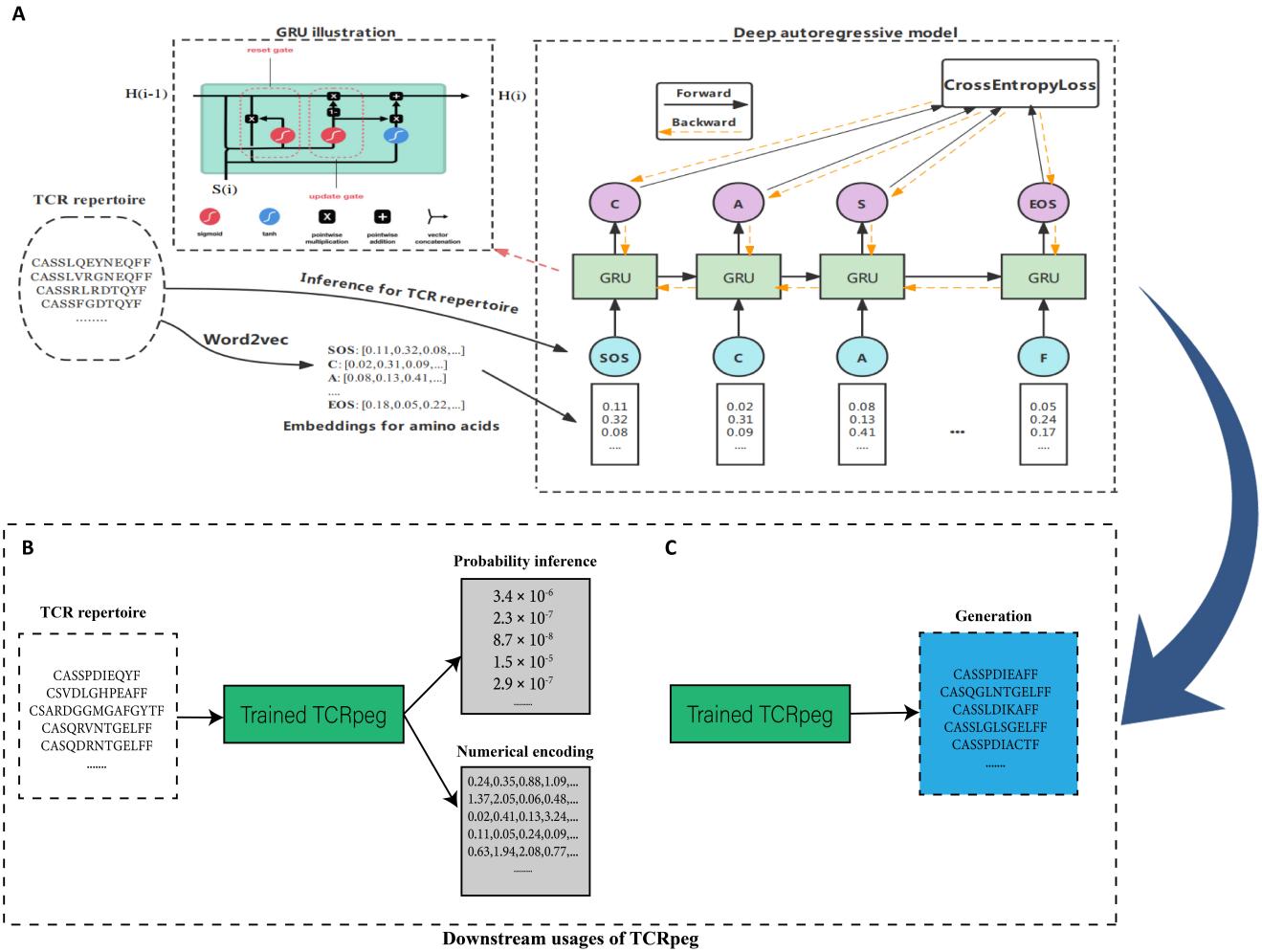


Figure 1. Workflow of TCRpeg to infer probabilistic patterns of immune receptor repertoires. **(A)** We have implemented a deep autoregressive network with GRU layers to process TCR sequences with different lengths to learn the hidden sequence pattern. The word2vec algorithm is first applied to the TCR repertoire to learn the numerical representations for each amino acid, regarding amino acids and TCR sequences as “Words” and “Sentences”. Then the TCR sequence is inputted to the deep autoregressive model sequentially. The model is updated by the gradient descent algorithm with the cross-entropy loss between the output logits and true labels. The trained TCRpeg model can be readily extended to downstream usages, including probability inference, encoding TCRs **(B)**, and generating similar new TCRs **(C)**. These functions and applications of TCRpeg are further elaborated in the Results section.

1 specificity^{30,31}.

2 In the past decade, deep learning models have achieved considerable success in handling sequential
 3 data^{26,32–35}. An autoregressive model processes the sequential data using observations from previous
 4 stages to infer the entry at the next time point. In the context of the TCR sequence, we can apply an
 5 autoregressive model to infer a residue using the amino acid subsequence proceeding from it. Therefore,
 6 we built TCRpeg, an autoregressive model that formulates the probability of a TCR sequence \mathbf{x} as $p(\mathbf{x}|\boldsymbol{\theta})$,
 7 where the parameters $\boldsymbol{\theta}$ capture the latent evolutionary patterns to generate \mathbf{x} . The probability density
 8 $p(\mathbf{x}|\boldsymbol{\theta})$ can be calculated by the product of probabilities conditioned on previous residues along a sequence

1 with length L through an autoregressive likelihood

$$p(\mathbf{x}|\boldsymbol{\theta}) = p(x_1|\boldsymbol{\theta}) \prod_{i=2}^L p(x_i|x_1, \dots, x_{i-1}; \boldsymbol{\theta}). \quad (1)$$

2 Figure 1 displays the TCRpeg workflow. We utilized gated recurrent units (GRUs)²⁶, commonly adopted
3 in recurrent neural networks, to model the autoregressive likelihood (Methods). Recurrent neural network
4 models might encounter gradient explosion for long peptide sequences^{36,37}. However, TCR sequences
5 contain mainly 12 to 17 residues (Supplementary S1). Thus, we can parameterize the generative process
6 with feed-forward GRU models that aggregate dependencies in sequences through the transmitting hidden
7 features controlled by the gate functions.

8 Training a GRU model requires vector representations for each amino acid. Instead of using one-hot
9 encodings or predefined characteristics of the analysis of principal components in biochemical features¹⁷,
10 we adopted the word2vec algorithm³⁸ to adaptively learn the embeddings for each amino acid from the TCR
11 sequencing data by treating an amino acid as a “Word” and each TCR sequence as a “Sentence” (Method).
12 Then, TCRpeg can be trained in a forward language modeling manner. To estimate the probability of a
13 given TCR sequence, we applied Eq.1 to the pre-trained TCRpeg. Details of the architecture of TCRpeg,
14 the training, and inferring processes are included in the Methods.

15 TCRpeg infers functional TCR repertoire probability distribution

16 First, we evaluated the probability distribution of the TCR sequences inferred by TCRpeg and compared
17 its accuracy with the other two probabilistic models, soNNia¹⁹ and TCRvae²⁹. To assess and compare
18 their performance, we constructed a universal TCR repertoire from a large cohort of 743 individuals from
19 Emerson *et al.*³⁹, following a similar data pre-processing strategy in Isacchini *et al.*¹⁹. Specifically, we
20 pooled the unique nucleotide sequences of TCRs from all individuals and constructed a universal TCR
21 repertoire. The universal repertoire was randomly divided into training and testing subrepertoires by a
22 50:50 split to ensure consistency with soNNia¹⁹ and TCRvae²⁹. Then we trained TCRpeg, soNNia, and
23 TCRvae on the training set (Methods).

24 We evaluated the three models, each to estimate a probability distribution $P_{infer}(\mathbf{x})$ for the test set;
25 TCRpeg shows high accuracy with substantial improvement over soNNia and TCRvae, yet requires lower

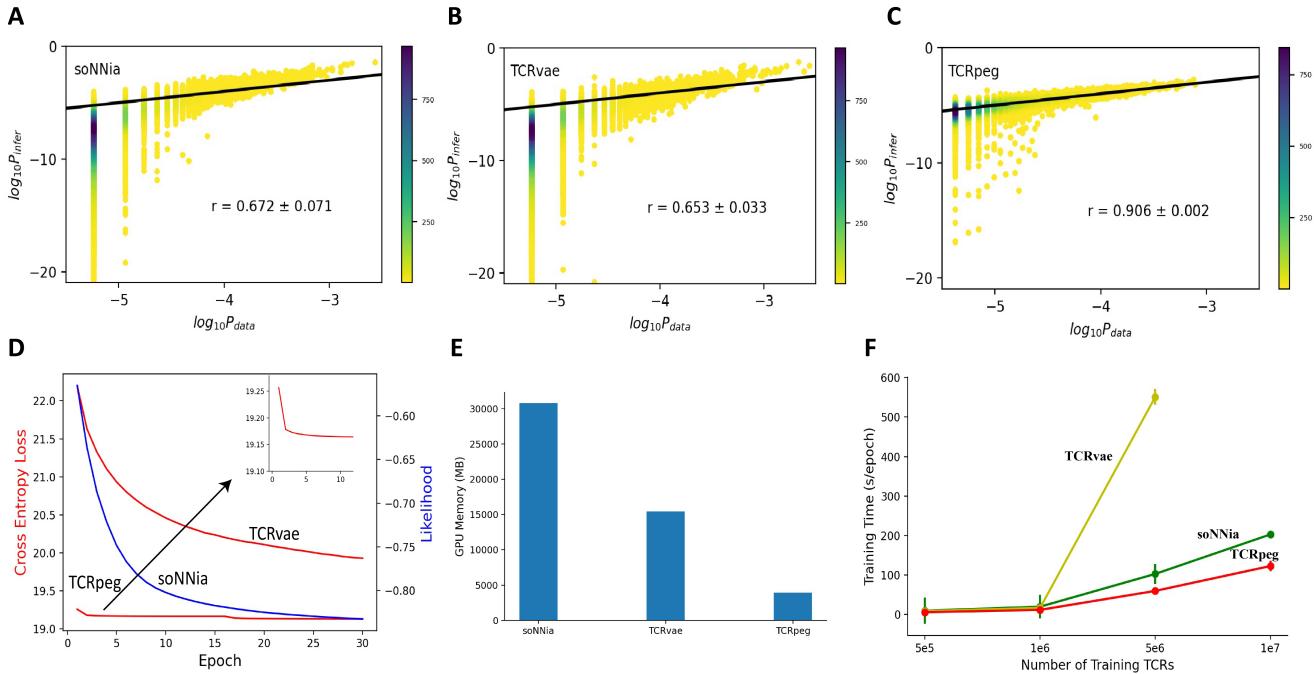


Figure 2. Performance of TCRpeg compared to the other two baseline methods soNNia and TCRvae. **(A-C)** Scatterplots of observed frequency P_{data} vs. estimated probability $P_{inferred}$ for **(A)** soNNia, **(B)** TCRvae and **(C)** TCRpeg models trained on the large TCR pool combining 743 individuals from Emerson *et al.*³⁹, along with the corresponding Pearson correlation coefficient r . The color indicates the number of sequences. **(D-F)** Comparison of soNNia, TCRvae, and TCRpeg model from practical aspects. Experiments are conducted under the same settings (learning rate and batch size) on a single Nvidia Tesla V100 GPU card with maximum 32 gigabytes memory. **(D)** The training curves for these three models. The soNNia model uses the likelihood as the model objective function (shown in the blue curve), while TCRvae and TCRpeg model minimize the cross-entropy loss (shown in the red curves). TCRpeg only needs less than ten epochs to converge, while the other two take around 30 epochs to converge. **(E)** The bar plot shows the GPU memory required to train each model. TCRpeg is more hardware-friendly. **(F)** The training speed of each model. TCRpeg takes less time to complete one training epoch compared to soNNia and TCRvae.

1 resources to train. Prediction accuracy can be quantified through the Pearson correlation coefficient r
 2 between the inferred and true probability distributions, i.e., $P_{inferred}(\mathbf{x})$ and $P_{data}(\mathbf{x})$, on test set (Methods).
 3 TCRpeg achieved $r \simeq 0.906$; however, soNNia and TCRvae obtained $r \simeq 0.672$ and $r \simeq 0.653$, respectively
 4 (Fig. 2A-2C). TCRpeg also performs stably and robustly when training on a small proportion of training
 5 data consisting of only 2×10^5 TCR sequences (Supplementary S2). In addition to the substantial
 6 accuracy improvement, TCRpeg converges faster and costs significantly less GPU memory (Fig. 2D-2F).
 7 It converges within five epochs, whereas the other two methods require around 30 epochs. Moreover,
 8 SoNNia and TCRae consume six times and three times more memory than TRCpeg.

9 TCRpeg gains new insights into TCR subrepertoires from T1D donors

10 The learned probability distribution can help profile the TCR subrepertoires. Here, we were interested in
 11 learning the cell-type level discrepancy and exploring the tissue level differences since T cells migrate
 12 and reside in different tissues and are influenced by different tissue environments. We chose Type 1

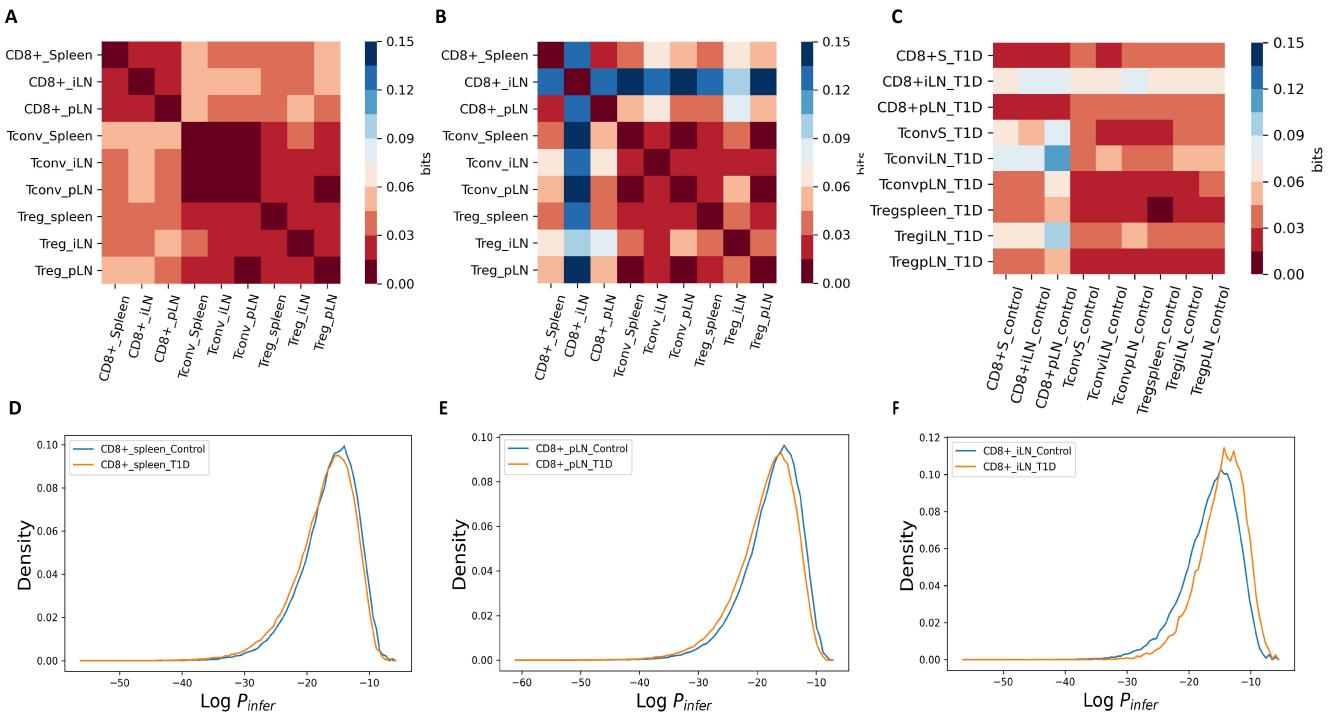


Figure 3. Jensen-Shannon divergences between TCR subrepertoires at the cell type and tissue level. (A and B) Jensen-Shannon divergences (D_{JS}) computed from TCRpeg trained in different subrepertoires of (A) healthy controls and (B) donors with T1D. (C) Jensen-Shannon divergences for subrepertoire pairs r^i from control and r^j from T1D donors. It reveals the divergences of the TCR subrepertoires of different cell types and tissue between healthy controls and T1D donors. (D - F) Distributions of the logarithmic probability of CD8⁺ TCR subrepertoires from control and T1D donors according to the TCRpeg model. (D), (E) and (F) show the log-probability distribution of CD8⁺ TCR subrepertoires from spleen, pLN, and iLN, respectively. CD8⁺ TCRs within iLN demonstrate significant divergence between control and T1D donors.

1 diabetes (T1D) as a case study. TID is an autoimmune disease that originates from insulin deficiency
 2 and causes hyperglycemia. The adaptive immune repertoire plays a critical role in type 1 diabetes (T1D)
 3 pathogenesis^{40–43}. Therefore, exploring the homogeneity and heterogeneity among the TCR subrepertoires
 4 in a probabilistic aspect might be illuminating to understand the immune response to T1D. We pooled
 5 TCRs with unique nucleotide sequences from nine control (healthy) individuals and seventeen donors
 6 with T1D from Seay *et al.*²¹. These TCR sequences were sorted into three cell types (CD4⁺ conventional
 7 T cells [Tconvs], CD4⁺ regulatory T cells [Tregs], and CD8⁺ T cells) and collected from three tissues
 8 (pancreatic draining lymph nodes [pLNs], mesenteric or inguinal “irrelevant” lymph nodes [iLNs], and
 9 spleen); that is, we have nine classes of the subrepertoires. We applied TCRpeg to infer the probability
 10 distribution of each subrepertoire and quantified the difference between these distributions using the
 11 Jensen-Shannon divergence D_{JS} (Methods).

12 First, in healthy donors, the subrepertoires belonging to the same cell type are more conserved across
 13 different tissues. The same cell type across different tissues shows a lower TCR subrepertoire divergence,

1 with an average Jensen-Shannon divergence as $D_{JS} \simeq 0.014$ bits (Fig. 3A). However, the divergence is
2 high between CD8⁺ and CD4⁺ TCR subrepertoires with the average $D_{JS} \simeq 0.041$ bits. Tconv and Treg
3 within the class of CD4⁺ cells demonstrate moderate similarities, with average $D_{JS} \simeq 0.024$ bits. Notably,
4 these observations confirm the results from Isacchini *et al.*¹⁹, where larger divergence between CD8⁺ and
5 CD4⁺ TCR subrepertoires and lower difference between Tconv and Treg TCR subrepertoires were shown.

6 Second, the CD8⁺ TCR subrepertoires show higher divergences across the T1D tissues (Fig. 3B). The
7 iLN CD8⁺ TCR subrepertoire differs significantly from those in spleen or pLN, with average $D_{JS} \simeq 0.128$
8 bits. In contrast, subrepertoires from spleen and pLN display a lower average $D_{JS} \simeq 0.016$ bits. This
9 observation indicates an active immune response of CD8⁺ associated with T1D. Further, we notice that
10 the scale of D_{JS} in T1D donors (0 to 0.15 bits, Fig. 3B) is much larger than that in healthy donors (0 to
11 0.06 bits, Fig. 3A), implying a further differentiation of T cells in response to T1D. This finding is in
12 accord with the results from Iria *et al.*⁴⁴, who conclude that the T1D TCR sequences show larger diversity
13 through the observation of a higher number of unique clonotypes in T1D donors.

14 Third, T1D donors have more divergent TCR subrepertoires than healthy controls, and their iLN TCR
15 subrepertoires diverged largely from healthy controls. Figure 3C shows the Jensen-Shannon divergence
16 between the healthy and T1D TCR subrepertoires. Between healthy controls and T1D donors, the spleen
17 and pLN CD8⁺ TCR subrepertoires show lower divergences (Fig. 3C-E), with an average $D_{JS} \simeq 0.025$
18 and $D_{JS} \simeq 0.018$, respectively; however, the iLN CD8⁺ subrepertoires from T1D donors demonstrate a
19 large discrepancy compared to healthy controls with $D_{JS} \simeq 0.082$ (Fig. 3C and 3F). These results imply a
20 pathogenic role of CD8⁺ from iLN in adaptive immunity for T1D that the T1D abnormal TCRs in iLN
21 might mistakenly result in the destruction of insulin-producing cells.

22 **Classification of cancer-associated TCRs and SARS-CoV-2 epitope-specific TCRs**

23 TCRpeg yields vector embeddings for TCRs sequences. Compared to the predefined or manually designed
24 encoding method for TCR sequences, TCRpeg provides a learnable way to encode TCR sequences into
25 vector representations. The update and reset gates of the GRU layers are learned during the training
26 process to determine how much of the previous information stored in the hidden features needs to be
27 passed along or abandoned²⁶ (Fig. 1A). Therefore, the hidden features of the GRU layers at the last

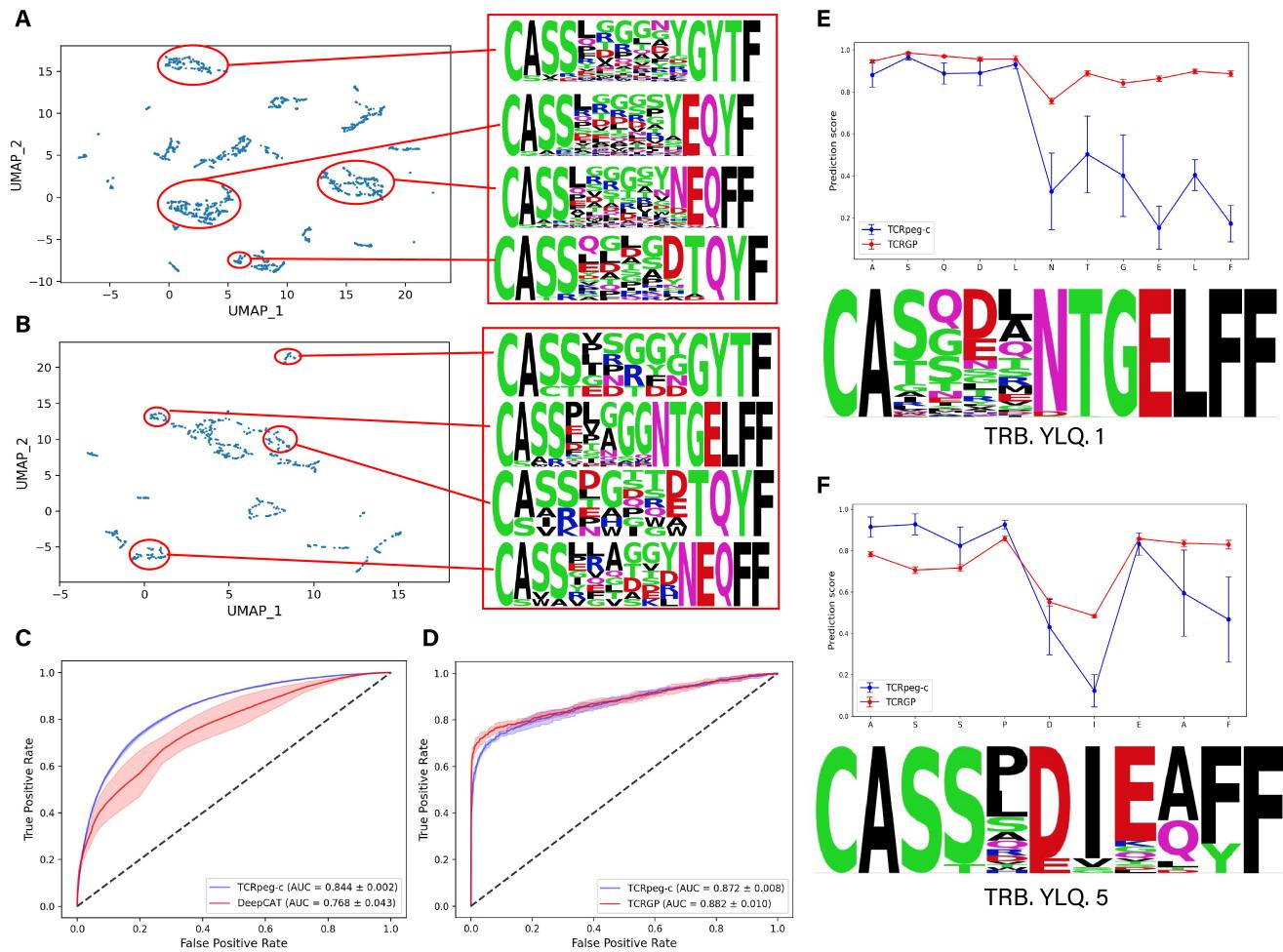


Figure 4. 2D illustration of TCRpeg-based encodings and predictive performance for downstream classification tasks. **(A and B)** 2D projection map of encodings obtained from TCRpeg trained on **(A)** caTCRs and **(B)** specific TCRs of the SARS-CoV-2 epitope (YLQPRTFLL). More projecting results can be found in Supplementary S3. **(C and D)** ROC curves for tasks of **(C)** predicting caTCR and **(D)** SARS-CoV-2 epitope YLQPRTFLL. **(E and F)** Sensitivity analysis through amino acid substitutions used TCRpeg-c and TCRGP for two previously identified TCR motifs. For each position other than the two ends, we changed the amino acid at that position to the four other most frequent AAs and used these two models to score the modified sequences. TCRpeg-c is more sensitive than TCRGP to substitutions of amino acids inside the motifs.

1 sequence position store summative information of the TCR sequences with different lengths; and these
 2 feature vectors provide an embedding for the TCR sequences.
 To illustrate the embedding of the TCR sequence, we first collected cancer-associated TCR (caTCR)
 from Beshnova *et al.*¹⁷ ($N \sim 43,000$) and SARS-CoV-2 epitope (YLQPRTFLL) specific TCRs from VDJdb
 database⁴⁵ ($N=683$). We trained TCRpeg on these two datasets separately and obtained the respective
 numerical TCR embedding vectors. The UMAP dimensionality reduction⁴⁶ was applied to project these
 vectors onto 2D space (Fig. 4A, Fig. 4B and Supplementary S3), showing that TCRs with a similar
 pattern (motif) tend to be clustered. It implies that the encodings could be helpful for antigen-specific TCR
 clustering. To further demonstrate the utility of TCRpeg-based encodings, we evaluated the classification

1 performance on caTCRs and SARS-CoV-2-epitope-specific TCRs using a fully connected neural network
2 (FCN), taking these vector encodings as inputs. We refer to this network as “TCRpeg-c” (Methods and
3 Supplementary S4). To collect negative (or control) samples for the epitope-specific TCR dataset, we
4 randomly sampled ten times more negative data than positive data from the universal TCR repertoire
5 mentioned above. We selected the CNN-based model - DeepCAT, developed in Beshnova *et al.* to compare
6 with in the caTCR prediction task, adopting the five-fold cross-validation procedure. In this prediction
7 task, we observed an improvement in accuracy and predictive stability for TCRpeg-c with an average AUC
8 $\simeq 0.844$ compared to DeepCAT with an average AUC $\simeq 0.768$ of caTCRs (Fig. 4C).

9 In the more challenging epitope-specific TCR prediction task with scant data, TCRpeg-c still demon-
10 strated competitive performance with AUC $\simeq 0.872$ compared to the baseline method TCRGP²⁷ with
11 AUC $\simeq 0.882$ (Fig. 4D). However, the TCRGP model is sophisticated, and it is designed specifically for
12 the TCR-epitope mapping problem with low data size, combining multiple techniques including alignment
13 of TCR sequences, Gaussian process (GP) and variational inference.

14 TCRpeg-c finds TCR motifs through perturbation analysis. TCR motifs are important and instructive
15 in determining their specificity to antigens⁴⁷. Previously, motif discovery for TCR repertoires was mainly
16 accomplished either by exploration of similarities between TCRs such as the TCRNET method^{48–50}
17 or investigation of frequency enhancement of k-mers for TCRs⁴⁷. Here, we used predictive models to
18 test the sensitivity of previously identified TCR motifs for specific TCRs of the SARS-CoV-2 epitope
19 YLQPRTFLL (Methods). We observed correspondence between previously identified TCR motifs and
20 sensitive residues according to scores predicted by TCRpeg-c, indicating the importance of TCR motifs
21 for epitope binding (Fig. 4E and 4F). However, although TCRGP achieves high predictive performance, it
22 lacks the ability to detect sensitive residues (Fig. 4E and 4F). We attribute its insensitivity to the necessity
23 of padding TCR sequences to a fixed length, which could lower the degree of variation caused by amino
24 acid substitution.

25 **Generating more TCR sequences with potentially the same specificity**

26 A good generative model could be beneficial for adoptive transfer of TCR engineered T cells (TCR-T)
27 that has been applied to treat viral infections such as hepatitis B and C^{51,52}, cancer immunotherapy^{53,54},

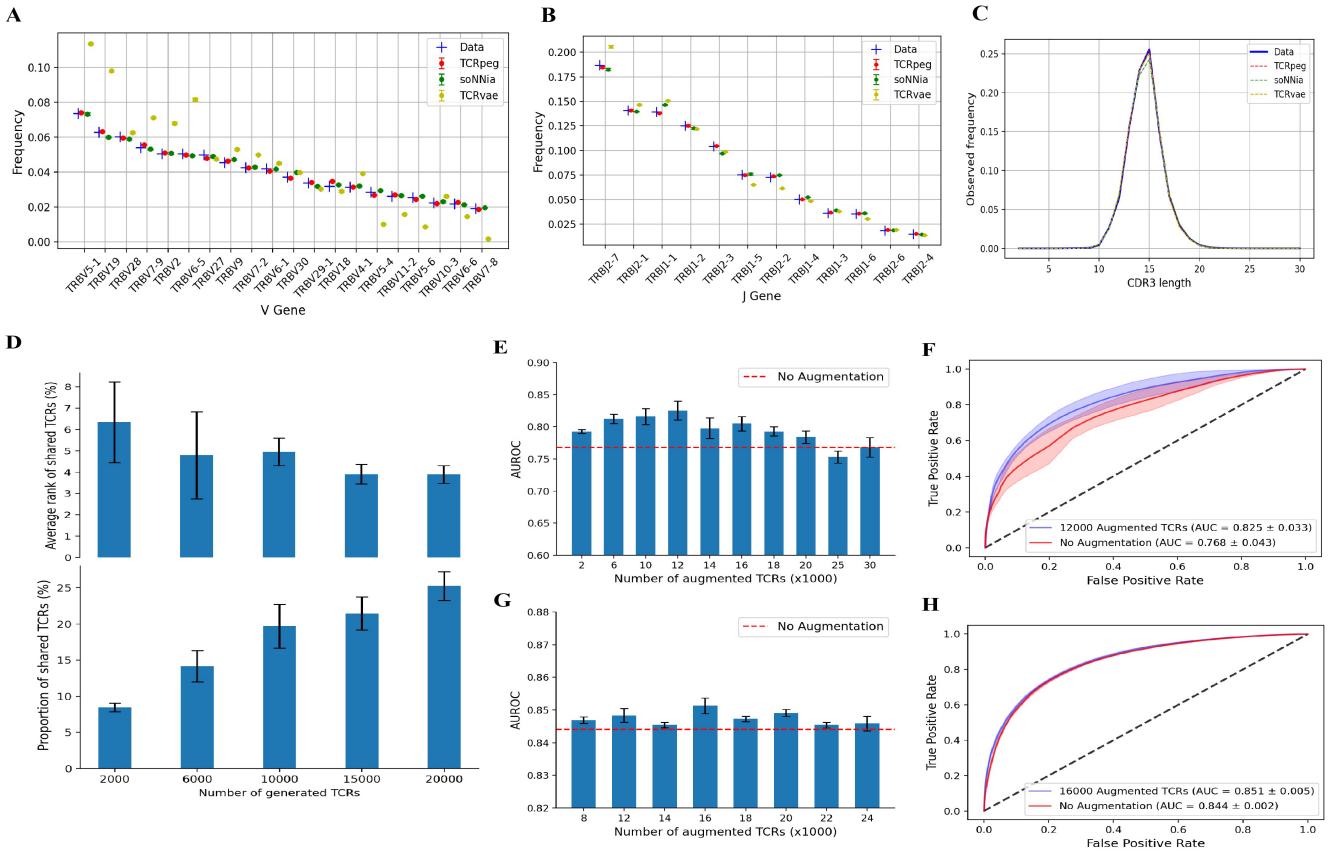


Figure 5. Characteristics of the TCR sequences generated by the three generative models. **(A-C)** Comparison of the statistical distributions of the generated sequences with the real data with respect to **(A)** V gene usage, **(B)** J gene usage and **(C)** length distribution. In **(A)**, only the top 20 frequent V genes are listed. We include the figure of full V gene usage and the distributions of amino acids in Supplementary S7 and S6. **(D)** The proportion of the TCR sequences in the test set that also appears in the generated TCRs (bottom panel) and the average probability rank of those shared TCRs among the generated TCRs (top panel).

With more TCRs being generated, more of them can be found in the test set. **(E - H)** Performance in the task of predicting catTCRs by applying the TCR-specific data augmentation technique. **(D and F)** The AUC scores with a different number of augmented TCR sequences when using the **(E)** DeepCAT model and **(G)** TCRpeg-c. **(F and H)** ROC curves for the DeepCAT model **(E)** and TCRpeg-c **(G)** with the best number of augmented TCRs.

1 and autoimmune disease therapy⁵⁵ through *in silico* generation of similar TCR sequences guiding the in
 2 vitro TCR design. We extended TCRpeg to be generative through a simple sampling strategy (Methods).
 3 We first aimed to systematically evaluate and compare the generation ability of TCRpeg with the
 4 baseline methods, soNNia and TCRvae, in terms of the statistical properties between the generated
 5 TCR sequences and real sequences. Specifically, we investigated the distributions of sequence lengths,
 6 positions of amino acids, V gene and J gene usages. We observed strong agreement between the probability
 7 distributions of *in silico* and real TCR repertoires for both the TCRpeg and soNNia models. For the position
 8 distributions of each amino acid, the TCRpeg- and soNNia- generated sequences successfully fitted the
 9 original statistics with an average Pearson correlation coefficient $r \simeq 1.0$ and $r \simeq 0.999$, respectively,
 10 compared to TCRvae with $r \simeq 0.982$ (Supplementary S6). For V and J gene usages, the TCRpeg

1 and soNNia models still outperform TCRvae, achieving average $r \simeq 0.999$ and $r \simeq 0.998$ compared to
2 $r \simeq 0.949$ of the TCRvae model for V gene usage distribution (Fig. 5A and Supplementary S7), and
3 $r \simeq 1.0$ and $r \simeq 0.997$ over $r \simeq 0.993$ for J gene usage distribution (Fig. 5B). For the length distribution,
4 these three models all achieved highly accurate performance with $r \simeq 1.0, 0.998, 1.0$ for TCRpeg, soNNia,
5 and TCRvae, respectively (Fig. 5C). The generation performance of TCRpeg is stable and accurate even
6 when trained on a small subset of TCRs (Supplementary S8 and S9). These results together highlight that
7 TCRpeg is reliable for summarizing a TCR repertoire, and consequently, generating new sequences in
8 recovering the real statistical distributions.

9 A reliable generative model should be able to produce new TCR sequences with “hidden similarity” to
10 real TCR data aside from statistical similarity. Here, we were interested to determine whether the generated
11 TCR sequences possess the same epitope specificity with the data used in training TCRpeg. To verify this,
12 we retrained TCRpeg on the training set of the TCRs specific to the epitope YLQPRTFLL and utilized
13 it to generate new sequences accordingly. We first noticed that some of the TCRs in the test set could
14 also be found in the generated data set (Fig. 5D), which shows the generative power of TCRpeg given
15 the wide potential diversity of TCR sequences. To take a closer look at these generated TCR sequences,
16 we observed that those TCRs that were also found in test set possessed high generation probabilities
17 (averagely ranked < 10% among generated sequences, Fig. 5D). Finally, we utilized the TCRMatch²⁸
18 software to further validate the hidden similarity of the generated TCR sequences and observed that
19 50 – 60% of them possess the same epitope specificity as the TCR sequences used in training according
20 to a scoring threshold of 0.9 (Supplementary S10). On the contrary, although the soNNia and TCRvae
21 models achieve comparable performance with respect to the statistical similarities, only less than 40% of
22 the generated sequences possess the same epitope specificity determined by the same scoring threshold
23 (Supplementary S10). Overall, our results indicate that TCRpeg can generate new TCR sequences with
24 statistical and possible hidden similarities to the TCRs used for training.

25 Augmenting TCR sequencing data

26 Data augmentation techniques are used ubiquitously in machine learning tasks to increase the generality
27 of data by adding similar samples generated by either slightly modifying the original data or synthesizing

similar data. They act as regularizers to alleviate the issue of overfitting and improve the generalization capacity of machine learning models, especially when applied to computer vision tasks⁵⁶ or natural language processing tasks⁵⁷. Adopting the data augmentation techniques here should improve the classification of TCR sequences. TCR sequences might abolish their epitope specificity by amino acid substitutions, especially when they happen inside contact motifs^{47,58}; therefore, directly performing amino acid substitutions, insertions, or deletions on TCR sequences cannot work as data augmentation. However, with the strong generative ability, TCRpeg may generate similar TCR sequences, and serve as a computational tool for TCR-specific data augmentation.

To analyze the feasibility of TCR-specific data augmentation, we evaluated and compared the predictive performance of classifying caTCRs with and without data augmentation while keeping all other training settings unchanged. For the DeepCAT model, we observe a large performance gain with up to 0.057 higher AUC when applying data augmentation technique (Fig. 5E and 5F). For the TCRpeg-c model, we still find accuracy enhancement in the AUC value from 0.844 to 0.851 with data augmentation (Fig. 5G and 5H). Besides, the AUPRC (area under the precision-recall curve) also increases and the test loss decreases, which is a positive sign of mitigation of overfitting (Supplementary S11). To further validate the utility of our TCRpeg-based augmentation technique, we performed classification for the influenza epitope GILGFVFTL and EBV GLCTLVAML epitope-specific TCRs with 3,406 and 962 positive samples, respectively, using the TCRex model⁵⁹. Without changing any training settings, we observed up to 2.1% and 21.4% accuracy enhancement for these two TCR datasets (Supplementary S12).

Discussion

An accurate probabilistic model for large-scale TCR sequencing data is a cornerstone for a better understanding of functional TCR repertoire. Previous works have developed selection models soNia²⁵, soNNia¹⁹, and the VAE-based model TCRvae²⁹ to characterize the distribution of productive TCR sequences. However, they are all intrinsically unable to capture the information behind the length variation. In this work, we introduced TCRpeg, an autoregressive deep learning model that utilizes a recurrent neural network with GRU layers to characterize the TCR repertoires. Unlike soNia, soNNia, and TCRvae which need to pad every TCR sequence to the same length, TCRpeg can process TCR sequences with any lengths.

1 Such capability can eliminate the noise introduced by adding an extra “amino acid” for padding and take
2 advantage of the information behind the variance in lengths.

3 We first demonstrated that TCRpeg can improve the statistical characterization of TCR repertoires in
4 a large cohort of individuals³⁹ compared to soNNia and TCRvae by a large margin, which implies that
5 TCRpeg can better learn the TCR sequence pattern. We attribute the superior performance of TCRpeg
6 to its ability to process TCRs with different lengths and its transmission of hidden features that properly
7 store the previous information. In particular, TCRpeg takes less iterations to converge and requires lower
8 computation resources. These results indicate the advantages of using an autoregressive model that is
9 capable of processing TCR sequences with different lengths to describe large-scale TCR sequencing data
10 from a probabilistic perspective.

11 Using the statistical inference power of TCRpeg, we explored the differences and similarities between
12 functional TCR subrepertoires collected from different cell types, tissues, or donor status (healthy or T1D)
13 at the repertoire level. We discovered that TCR subrepertoires belonging to families with more closely
14 related developmental paths (i.e., Tconvs and Tregs) possess high statistical similarities. They both show
15 large differences with CD8⁺ T cells that diverged earlier in T cell maturation. In response to T1D disease,
16 we found that different subsets of T cells further adapt themselves, especially for CD8⁺ T cells within
17 the “irrelevant” lymph nodes (iLN) that diverge largely from other subrepertoires of TCR. The prominent
18 changes in the TCR subrepertoires inside iLN suggest a possible etiopathogenesis of T1D that these
19 abnormal TCRs might mistakenly destroy insulin-producing cells. Kurrer and his colleagues partially
20 support our hypothesis; they observed the initial onset of diabetes in healthy mice on day eight after T
21 cells collected from the spleen and mesenteric lymph nodes of overtly diabetic mice⁶⁰. Previous work also
22 identified significant changes in cellular composition or T cell quantities for T1D or other donors with
23 autoimmune diseases^{61–63}. More studies are needed to investigate the immune response at the tissue level
24 to better understand the pathogenesis of T1D or other autoimmune diseases and develop more effective
25 therapies.

26 On the basis of the architecture of TCRpeg, we can obtain helpful vector representations of TCR
27 sequences from the trained TCRpeg model, which is not provided by soNNia or TCRvae. Compared to
28 other predefined or hand-designed encoding methods for TCR sequences, TCRpeg provides a learnable

1 way to encode TCR sequences by updating functional gates inside GRU layers²⁶. We observed that
2 TCRpeg-based TCR encodings could reflect the degrees of similarities between TCR sequences that
3 sequences with a similar pattern (motifs) tend to cluster together (Fig. 4A and 4B). This suggests a
4 potential application of antigen-specific TCR clustering, since shared TCR motifs indicate the same
5 antigen specificity.

6 To examine the performance of TCRpeg-based encodings in a predictive manner, we assessed the
7 classification performance of caTCRs and YLQPRTFLL epitope-specific TCRs using a fully connected
8 neural network taking these vector encodings as input (TCRpeg-c). For the caTCR prediction task, we
9 chose the DeepCAT model developed by Beshnova *et al.* as the baseline method. We observed a significant
10 improvement in accuracy and predictive stability for TCRpeg-c compared to DeepCAT in the prediction
11 of caTCRs (Fig. 4C). With such high precision, TCRpeg-c could facilitate cancer detection through the
12 process introduced in Beshnova *et al.*. In recent years, multiple machine learning methods have been
13 developed to predict the epitope specificity of TCRs, such as TCRex⁵⁹, DeepTCR⁵⁸, and TCRGP²⁷. All of
14 these methods have explored the problem in slightly different settings and compared with each other. In the
15 more challenging classification task of predicting SARS-CoV-2 epitope (YLQPRTFLL)-specific TCRs, we
16 compared TCRpeg-c to a representative of the above group of machine learning models, TCRGP, which is
17 a combination of multiple functional modules including TCR alignment, Gaussian process, and variational
18 inference. TCRpeg-c demonstrated competitive performance in this task compared to TCRGP (Fig. 4D).
19 In particular, TCRpeg-c is sensitive to substituting for an amino acid primarily when it occurs inside the
20 TCR motifs, while TCRGP is insensitive to that (Fig. 4E and 5F). This finding indicates that TCRpeg-c
21 can be used for motif validation and help with TCR engineering for immunotherapies⁶⁴. In addition, this
22 perturbation analysis might reveal *de novo* motifs that have not yet been discovered using nonpredictive
23 methods (Supplementary S5). The comparable accuracy performances in the above two classification
24 challenges validate the advantage of TCRpeg-based encodings, which can be further concatenated with
25 epitope features to facilitate the unseen epitope-TCR interaction prediction task⁶⁵.

26 One direct application of TCRpeg is to generate new TCR sequences with characteristics similar to
27 those of natural sequences. We first compared the generation capability of TCRpeg with soNNia and
28 TCRvae with respect to the statistical distributions on the large universal TCR pool we have constructed.

1 We showed that TCRpeg-generated TCR sequences had the closest amino acid distributions, length
2 distribution, and V/J gene usages to the real sequences compared to the other baseline methods. Next,
3 we found that some TCRs in the test set could also be found in the generated dataset, and those shared
4 TCRs have high generation probabilities among the generated dataset (Fig. 5D). These results imply that
5 newly generated TCR sequences with high probabilities might share the same epitope specificity with the
6 data used in training, providing a potential way to meet the demand for more data. We further applied the
7 TCRMatch²⁸ software to validate this implication and show that 50 – 60% of the generated TCRs share
8 the same epitope specificity as the TCRs used for training. On the contrary, less than 40% of the TCRs
9 generated using TCRvae or soNNia share the same specificity (Supplementary S10). The generative power
10 of TCRpeg can also be used to design similar TCRs to facilitate immunotherapy for T cell transfer^{53–55}.

11 Data augmentation is a ubiquitous technique used to increase the performance of machine learning
12 models, especially in computer vision systems⁵⁶. Given that more and more machine learning models
13 have been developed for TCR-related tasks and the acquisition of more data is costly and time consuming,
14 which restricts the development of highly accurate machine learning models, we developed and validated
15 the TCR-specific data augmentation technique empowered by TCRpeg to relieve such restriction. For the
16 caTCR classification task, we observed a notable improvement with data augmentation (Fig. 5E and 5F).
17 In addition, we further validated the utility of data augmentation using another machine learning model -
18 TCRex⁵⁹ in the prediction tasks of GILGFVFTL and GLCTLVAML specific TCRs and again observed an
19 improvement in accuracy (Supplementary S12). However, in the SARS-CoV-2 specific TCR recognition
20 task, data augmentation failed to boost the model performance. When learning from such a small data size,
21 TCRpeg tends to generate highly similar TCRs with those in the training set and thus provides limited
22 additional information to the predictive model, which might result in more severe overfitting. Nevertheless,
23 TCRpeg-based data augmentation is a free option for boosting model performance without any extra cost.

24 In summary, we have introduced a new holistic software tool TCRpeg for probability inference,
25 encoding TCRs, and generating new TCR sequences with both statistical and hidden similarities. These
26 functions can be applied to different tasks and bring us new insights for understanding the complex
27 genomic concepts hidden behind TCR repertoires.

¹ Methods

² Data Description

³ The data sets used in this work are classified into three groups to evaluate the performance of TCRpeg. We
⁴ filter out TCRs with lengths greater than 30 or not starting with a cysteine in all data sets. We also verified
⁵ sequences that are written as V gene, CDR3 sequence, J gene and removed sequences with unknown
⁶ genes.

⁷ 1. To quantify the precision of the inference of TCRpeg along with the other two baseline methods, we
⁸ used the TCR repertoires sampled from a large cohort, including 743 individuals from Emerson
⁹ *et al.*³⁹ We pooled the unique nucleotide sequences of receptors from all individuals and built a
¹⁰ universal TCR pool that contains around 10^9 sequences in total. The multiplicity of an amino acid
¹¹ sequence in this universal TCR pool indicates the number of independent recombination events that
¹² led to that receptor. We randomly and equally split the TCR pool into a training set and a test set.

¹³ 2. To characterize the differences between the TCR subrepertoires of functional cell types collected
¹⁴ from different tissues, we pooled unique TCRs from 9 control donors and 17 T1D donors from Seay
¹⁵ *et al.*²¹ at a tissue level. These TCR sequences were sorted into three cell types and collected from
¹⁶ three tissues. Thus, for each donor status (healthy or T1D), we have nine groups of TCRs.

¹⁷ 3. To evaluate the performance of TCRpeg-c in classification tasks, we first collected cancer-associated
¹⁸ TCRs (caTCRs) from Beshnova *et al.*. Briefly, Beshnova and his colleagues collected TCR sequences
¹⁹ from approximately 4,200 recorded samples downloaded from The Cancer Genome Atlas (TCGA)
²⁰ and excluded those sequences that are also found in healthy donors. The remaining around 43000
²¹ TCR sequences are assumed to be cancer-associated TCRs (caTCRs). We extracted the SARS-CoV-
²² 2 epitope (YLQPRTFLL) specific TCRs from VDJdb⁴⁵ database (positive TCRs N = 683, extracted
²³ on 24 January 2022). We then randomly sampled ten times more negative data than positive data
²⁴ from the universal TCR pool constructed previously to serve as the control TCRs.

¹ TCRpeg and TCRpeg-c

² The illustrations of TCRpeg and TCRpeg-c are shown in Fig. 1A and Supplementary S4. To enable the
³ training of TCRpeg, we first trained the word2vec³⁸ model on 1×10^6 TCR sequences randomly sampled
⁴ from the pooled universal repertoire aforementioned to obtain the numerical embeddings for each amino
⁵ acid, regarding the amino acid as the “words” and the TCR sequences as the “sentences”. Specifically,
⁶ we adopted the skip-gram architecture with the window size and embedding size set to 2 and 32 and
⁷ trained it for 20 epochs. For the TCRpeg model, the GRU modules have three layers with the size of the
⁸ hidden feature set to 64. We trained TCRpeg using the Adam⁶⁶ optimizer for 20 epochs to minimize the
⁹ cross-entropy loss between the soft-maxed logits and the one-hot encoded representation of the discrete
¹⁰ categorical outputs of the network. The probability of a given TCR sequence $P_{infer}(\mathbf{x})$ is estimated using
¹¹ Equation 1. Specifically, we input the given TCR sequence to TCRpeg and obtain the corresponding
¹² output probability distribution of the amino acid at the next time step. Thus, $P_{infer}(\mathbf{x})$ is the multiplication
¹³ of the probabilities of amino acids at each time step.

¹⁴ For the TCRpeg-c model, the size of the hidden feature is increased to 512 to better capture the hidden
¹⁵ sequence features for classification tasks. On top of the pre-trained TCRpeg, the fully connected neural
¹⁶ network contains two hidden layers with 384 and 96 neurons, followed by the ReLU activation function.
¹⁷ In the task of predicting caTCRs, we trained TCRpeg-c for 30 epochs to minimize the loss of cross-entropy
¹⁸ between the output logits and true labels, with dropout operations ($p=0.2$) to reduce the issue of overfitting.
¹⁹ In classifying epitope-specific SARS-CoV-2 TCRs, we trained TCRpeg-c for 20 epochs with a dropout
²⁰ rate set to 0.4. In both above-mentioned classification tasks, the TCRpeg was trained on the respective
²¹ training set to provide the numerical embeddings for TCRs. The trained TCRpeg-c can be used to find
²² TCR motifs through perturbation analysis. Specifically, we permuted each position of the TCR sequences
²³ except for the first and last positions, with four other amino acids that most likely appeared at that position
²⁴ according to the amino acid frequency at that position. We adopted this strategy to avoid skewed permuted
²⁵ sequences containing amino acids at some positions with nearly zero probabilities. We then applied the
²⁶ trained TCRpeg-c to score each permuted sequence to determine residues that are sensitive to changes.

1 Quantifying the accuracy of probability inference.

2 To evaluate the precision of probability inference, we compared the estimated probabilities $P_{infer}(\mathbf{x})$ to
3 the observed frequencies $P_{data}(\mathbf{x})$ of the test set. The accuracy can be quantified by Pearson's correlation
4 coefficient r between $P_{infer}(\mathbf{x})$ and $P_{data}(\mathbf{x})$. A higher value of r indicates a better model. The calculation
5 of $P_{infer}(\mathbf{x})$ for TCRpeg is described in the previous section using Eq.1. For the two baseline methods
6 TCRvae and soNNia, we compute $P_{infer}(\mathbf{x})$ by:

$$P_{infer}(\mathbf{x}) = \sum_{v,j} P_{infer}(\mathbf{x}, v, j), \quad (2)$$

7 which sums the V and J genes along with the TCR sequence \mathbf{x} . Finally, we normalize the inferred
8 probabilities $P_{infer}(\mathbf{x})$ and consider them as the approximation of the real probability distribution.

9 Quantifying of difference between TCR subrepertoires

10 We used the Jensen-Shannon divergence $D_{JS}(r^i, r^j)$ to characterize the difference between two TCR
11 subrepertoires r^i and r^j :

$$D_{JS}(r^i, r^j) = \frac{1}{2} D_{KL}(P_{infer}^i, M) + \frac{1}{2} D_{KL}(P_{infer}^j, M), \quad (3)$$

12 where $M = \frac{1}{2}(P_{infer}^i + P_{infer}^j)$ and D_{KL} represent the Kullback-Leibler divergence. To characterize the
13 differences between the TCR subrepertoires of functional cell types collected from different tissues, we
14 first trained TCRpeg on each tissue-level TCR subset for 20 epochs with hidden size and the number of
15 layers set to 128 and 3, respectively. Then we applied Eq.3 to calculate the JS divergences between each
16 pair of those TCR subrepertoires.

17 Using TCRpeg to generate TCR sequences

18 We adopted a simple sampling method to generate new TCR sequences using TCRpeg. Specifically, we
19 first input the start token (“<SOS>”) to the TCRpeg and then randomly sampled the amino acid for the
20 next position from the output probability distribution (computed using the Softmax operation). Following
21 the same procedure, at each time step, we randomly sampled the amino acid for that time step according to

1 the probability distribution defined by the predicted scores and input it to the next time step to obtain the
2 following amino acids. This stochastic generation procedure can be described by the formula stated below:

$$AA_t = P(AA|AA_{t-1:0}; \boldsymbol{\theta}), \quad (4)$$

3 where AA_0 stands for the start token and $\boldsymbol{\theta}$ represents the TCRpeg parameters. The generation process
4 stops when the special stop token (“<EOS>”) is generated. To allow the ability to infer the corresponding
5 V and J gene along with the TCR sequence, we extended TCRpeg and formulated the probability of a
6 given TCR sequence \mathbf{x} with specific V and J genes as:

$$p(\mathbf{x}, V, J | \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3) = p(x_1 | \boldsymbol{\theta}_1) \prod_{i=2}^L p(x_i | x_1, \dots, x_{i-1}; \boldsymbol{\theta}_1) p(V | \mathbf{x}; \boldsymbol{\theta}_2) p(J | \mathbf{x}; \boldsymbol{\theta}_3), \quad (5)$$

7 where $p(V | \mathbf{x}; \boldsymbol{\theta}_2)$ and $p(J | \mathbf{x}; \boldsymbol{\theta}_3)$ are the probabilities conditioning on the TCR sequence \mathbf{x} ; $\boldsymbol{\theta}_2$ and $\boldsymbol{\theta}_3$ are
8 parameterized by two respective single-layer fully connected neural networks. The TCRpeg, soNNia and
9 TCRvae models were inferred from the universal TCR repertoire aforementioned, and then we applied
10 them to generate new TCR sequences along with V and J genes.

11 Data availability

12 All data analyzed in this work can be found in the original publications that collected the data^{17, 17, 39, 45, 61},
13 and we include the preprocessed data at <https://github.com/jiangdada1221/TCRpeg>.

14 Code availability

15 TCRpeg was written in Python using the deep learning library Pytorch⁶⁷ and is available as a python pack-
16 age. Source code, use-case tutorials, and documentations can be found at <https://github.com/jiangdada1221/TCRpeg>.
17 Users can install directly from Github or PyPI via pip.

1 References

- 2 1. Susumu Tonegawa. Somatic generation of antibody diversity. *Nature*, 302(5909):575–581, 1983.
- 3 2. Mark M Davis and Pamela J Bjorkman. T-cell antigen receptor genes and t-cell recognition. *Nature*,
- 4 334(6181):395–402, 1988.
- 5 3. Srinika Ranasinghe, Pedro A Lamothe, Damien Z Soghoian, Samuel W Kazer, Michael B Cole,
- 6 Alex K Shalek, Nir Yosef, R Brad Jones, Faith Donaghey, Chioma Nwonu, et al. Antiviral cd8+ t cells
- 7 restricted by human leukocyte antigen class ii exist during natural hiv infection and exhibit clonal
- 8 expansion. *Immunity*, 45(4):917–930, 2016.
- 9 4. P Anton Van Der Merwe and Omer Dushek. Mechanisms for t cell receptor triggering. *Nature*
- 10 *Reviews Immunology*, 11(1):47–55, 2011.
- 11 5. Philippa Marrack and John Kappler. The t cell receptor. *Science*, 238(4830):1073–1079, 1987.
- 12 6. Jannie Borst, Heinz Jacobs, and Gaby Brouns. Composition and function of t-cell receptor and b-cell
- 13 receptor complexes on precursor lymphocytes. *Current opinion in immunology*, 8(2):181–190, 1996.
- 14 7. Nishant K Singh, Timothy P Riley, Sarah Catherine B Baker, Tyler Borrman, Zhiping Weng, and
- 15 Brian M Baker. Emerging concepts in tcr specificity: rationalizing and (maybe) predicting outcomes.
- 16 *The Journal of Immunology*, 199(7):2203–2213, 2017.
- 17 8. XL Hou, L Wang, YL Ding, Q Xie, and HY Diao. Current status and recent advances of next
- 18 generation sequencing techniques in immunological repertoire. *Genes & Immunity*, 17(3):153–164,
- 19 2016.
- 21 9. Sebastian Zeissig, Elisa Rosati, C Marie Dowds, Konrad Aden, Johannes Bethge, Berenice Schulte,
- 22 Wei Hung Pan, Neha Mishra, Maaz Zuhayra, Marlies Marx, et al. Vedolizumab is associated with
- 23 changes in innate rather than adaptive immunity in patients with inflammatory bowel disease. *Gut*,
- 24 68(1):25–39, 2019.
- 25 10. Jonathan R McDaniel, Brandon J DeKosky, Hidetaka Tanno, Andrew D Ellington, and George
- 26 Georgiou. Ultra-high-throughput sequencing of the immune receptor repertoire from millions of
- lymphocytes. *Nature protocols*, 11(3):429–442, 2016.

- 1 **11.** Maria A Turchaninova, Olga V Britanova, Dmitriy A Bolotin, Mikhail Shugay, Ekaterina V Putintseva,
2 Dmitriy B Staroverov, George Sharonov, Dmitriy Shcherbo, Ivan V Zvyagin, Ilgar Z Mamedov, et al.
3 Pairing of t-cell receptor chains via emulsion pcr. *European journal of immunology*, 43(9):2507–2515,
4 2013.
- 5 **12.** Pradyot Dash, Andrew J Fiore-Gartland, Tomer Hertz, George C Wang, Shalini Sharma, Aisha
6 Souquette, Jeremy Chase Crawford, E Bridie Clemens, Thi HO Nguyen, Katherine Kedzierska,
7 et al. Quantifiable predictive features define epitope-specific t cell receptor repertoires. *Nature*,
8 547(7661):89–93, 2017.
- 9 **13.** David Masopust and Louis J Picker. Hidden memories: frontline memory t cells and early pathogen
10 interception. *The Journal of Immunology*, 188(12):5811–5817, 2012.
- 11 **14.** Nina Le Bert, Anthony T Tan, Kamini Kunasegaran, Christine YL Tham, Morteza Hafezi, Adeline
12 Chia, Melissa Hui Yen Chng, Meiyin Lin, Nicole Tan, Martin Linster, et al. Sars-cov-2-specific t cell
13 immunity in cases of covid-19 and sars, and uninfected controls. *Nature*, 584(7821):457–462, 2020.
- 14 **15.** Thomas M Snyder, Rachel M Gittelman, Mark Klinger, Damon H May, Edward J Osborne, Ruth
15 Taniguchi, H Jabran Zahid, Ian M Kaplan, Jennifer N Dines, Matthew N Noakes, et al. Magnitude
16 and dynamics of the t-cell response to sars-cov-2 infection at both individual and population levels.
17 *MedRxiv*, 2020.
- 18 **16.** Jiefei Han, Ruofei Yu, Jianchun Duan, Jin Li, Wei Zhao, Guoshuang Feng, Hua Bai, Yuqi Wang,
19 Xue Zhang, Rui Wan, et al. Weighting tumor-specific tcr repertoires as a classifier to stratify the
20 immunotherapy delivery in non–small cell lung cancers. *Science Advances*, 7(21):eabd6971, 2021.
- 21 **17.** Daria Beshnova, Jianfeng Ye, Oreoluwa Onabolu, Benjamin Moon, Wenxin Zheng, Yang-Xin Fu,
22 James Brugarolas, Jayanthi Lea, and Bo Li. De novo prediction of cancer-associated t cell receptors
23 for noninvasive cancer detection. *Science translational medicine*, 12(557), 2020.
- 24 **18.** Ryan Emerson, Anna Sherwood, Cindy Desmarais, Sachin Malhotra, Deborah Phippard, and Harlan
25 Robins. Estimating the ratio of cd4+ to cd8+ t cells using high-throughput sequence data. *Journal of*
26 *immunological methods*, 391(1-2):14–21, 2013.

- 1 **19.** Giulio Isacchini, Aleksandra M Walczak, Thierry Mora, and Armita Nourmohammad. Deep generative
2 selection models of t and b cell receptor repertoires with sonnia. *Proceedings of the National Academy*
3 *of Sciences*, 118(14), 2021.
- 4 **20.** Jason A Carter, Jonathan B Preall, Kristina Grigaityte, Stephen J Goldfless, Eric Jeffery, Adrian W
5 Briggs, Francois Vigneault, and Gurinder S Atwal. Single t cell sequencing demonstrates the functional
6 role of $\alpha\beta$ tcr pairing in cell lineage and antigen specificity. *Frontiers in immunology*, 10:1516, 2019.
- 7 **21.** Howard R Seay, Erik Yusko, Stephanie J Rothweiler, Lin Zhang, Amanda L Posgai, Martha Campbell-
8 Thompson, Marissa Vignali, Ryan O Emerson, John S Kaddis, Dave Ko, et al. Tissue distribution and
9 clonal diversity of the t and b cell repertoire in type 1 diabetes. *JCI insight*, 1(20), 2016.
- 10 **22.** Anand Murugan, Thierry Mora, Aleksandra M Walczak, and Curtis G Callan. Statistical inference of
11 the generation probability of t-cell receptors from sequence repertoires. *Proceedings of the National*
12 *Academy of Sciences*, 109(40):16161–16166, 2012.
- 13 **23.** Quentin Marcou, Thierry Mora, and Aleksandra M Walczak. High-throughput immune repertoire
14 analysis with igor. *Nature communications*, 9(1):1–10, 2018.
- 15 **24.** Yuval Elhanati, Anand Murugan, Curtis G Callan, Thierry Mora, and Aleksandra M Walczak. Quantifi-
16 fying selection in immune receptor repertoires. *Proceedings of the National Academy of Sciences*,
17 111(27):9875–9880, 2014.
- 18 **25.** Zachary Sethna, Giulio Isacchini, Thomas Dupic, Thierry Mora, Aleksandra M Walczak, and Yuval
19 Elhanati. Population variability in the generation and thymic selection of t-cell repertoires. *arXiv*
20 *preprint arXiv:2001.02843*, 2020.
- 21 **26.** Junyoung Chung, Caglar Gülcöhre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of
22 gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014.
- 23 **27.** Emmi Jokinen, Jani Huuhtanen, Satu Mustjoki, Markus Heinonen, and Harri Lähdesmäki. Determin-
24 ing epitope specificity of t cell receptors with tcrgp. *BioRxiv*, page 542332, 2019.
- 25 **28.** William D Chronister, Austin Crinklaw, Swapnil Mahajan, Randi Vita, Zeynep Koşaloğlu-Yalçın,
26 Zhen Yan, Jason A Greenbaum, Leon E Jessen, Morten Nielsen, Scott Christley, et al. Tcrmatch:

- 1 Predicting t-cell receptor specificity based on sequence similarity to previously characterized receptors.
- 2 *Frontiers in immunology*, 12:673, 2021.
- 3 **29.** Kristian Davidsen, Branden J Olson, William S DeWitt III, Jean Feng, Elias Harkins, Philip Bradley,
4 and Frederick A Matsen IV. Deep generative models for t cell receptor protein sequences. *Elife*,
5 8:e46935, 2019.
- 6 **30.** K Christopher Garcia and Erin J Adams. How the t cell receptor sees antigen—a structural view. *Cell*,
7 122(3):333–336, 2005.
- 8 **31.** Kai W Wucherpfennig, Etienne Gagnon, Melissa J Call, Eric S Huseby, and Matthew E Call. Structural
9 biology of the t-cell receptor: insights into receptor assembly, ligand recognition, and initiation of
10 signaling. *Cold Spring Harbor perspectives in biology*, 2(4):a005140, 2010.
- 11 **32.** Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–
12 1780, 1997.
- 13 **33.** Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
14 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 15 **34.** Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
16 Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio,
17 H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information
18 Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 19 **35.** James Douglas Hamilton. *Time series analysis*. Princeton university press, 2020.
- 20 **36.** Razvan Pascanu, Tomás Mikolov, and Yoshua Bengio. Understanding the exploding gradient problem.
21 *CoRR*, abs/1211.5063, 2012.
- 22 **37.** Jung-Eun Shin, Adam J Riesselman, Aaron W Kollasch, Conor McMahon, Elana Simon, Chris Sander,
23 Aashish Manglik, Andrew C Kruse, and Debora S Marks. Protein design and variant prediction using
24 autoregressive generative models. *Nature communications*, 12(1):1–11, 2021.
- 25 **38.** Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representa-
26 tions in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

- 1 **39.** Ryan O Emerson, William S DeWitt, Marissa Vignali, Jenna Gravley, Joyce K Hu, Edward J Osborne,
2 Cindy Desmarais, Mark Klinger, Christopher S Carlson, John A Hansen, et al. Immunosequencing
3 identifies signatures of cytomegalovirus exposure history and hla-mediated effects on the t cell
4 repertoire. *Nature genetics*, 49(5):659–665, 2017.
- 5 **40.** Mark A Atkinson, George S Eisenbarth, and Aaron W Michels. Type 1 diabetes. *The Lancet*,
6 383(9911):69–82, 2014.
- 7 **41.** Bart O Roep and Mark Peakman. Diabetogenic t lymphocytes in human type 1 diabetes. *Current*
8 *opinion in immunology*, 23(6):746–753, 2011.
- 9 **42.** Greig P Lennon, Maria Bettini, Amanda R Burton, Erica Vincent, Paula Y Arnold, Pere Santamaria,
10 and Dario AA Vignali. T cell islet accumulation in type 1 diabetes is a tightly regulated, cell-
11 autonomous event. *Immunity*, 31(4):643–653, 2009.
- 12 **43.** Bart O Roep and Mark Peakman. Antigen targets of type 1 diabetes autoimmunity. *Cold Spring*
13 *Harbor perspectives in medicine*, 2(4):a007781, 2012.
- 14 **44.** Iria Gomez-Tourino, Yogesh Kamra, Roman Baptista, Anna Lorenc, and Mark Peakman. T cell
15 receptor β -chains display abnormal shortening and repertoire sharing in type 1 diabetes. *Nature*
16 *communications*, 8(1):1–15, 2017.
- 17 **45.** Mikhail Shugay, Dmitriy V Bagaev, Ivan V Zvyagin, Renske M Vroomans, Jeremy Chase Crawford,
18 Garry Dolton, Ekaterina A Komech, Anastasiya L Sycheva, Anna E Koneva, Evgeniy S Egorov, et al.
19 Vdjdb: a curated database of t-cell receptor sequences with known antigen specificity. *Nucleic acids*
20 *research*, 46(D1):D419–D427, 2018.
- 21 **46.** Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and
22 projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 23 **47.** Jacob Glanville, Huang Huang, Allison Nau, Olivia Hatton, Lisa E Wagar, Florian Rubelt, Xuhuai
24 Ji, Arnold Han, Sheri M Krams, Christina Pettus, et al. Identifying specificity groups in the t cell
25 receptor repertoire. *Nature*, 547(7661):94–98, 2017.

- 1 **48.** Mikhail V Pogorelyy and Mikhail Shugay. A framework for annotation of antigen specificities in
2 high-throughput t-cell repertoire sequencing studies. *Frontiers in immunology*, 10:2159, 2019.
- 3 **49.** Paul-Gydeon Ritvo, Ahmed Saadawi, Pierre Barennes, Valentin Quiniou, Wahiba Chaara, Karim
4 El Soufi, Benjamin Bonnet, Adrien Six, Mikhail Shugay, Encarnita Mariotti-Ferrandiz, et al. High-
5 resolution repertoire analysis reveals a major bystander activation of tfh and tfr cells. *Proceedings of
6 the National Academy of Sciences*, 115(38):9604–9609, 2018.
- 7 **50.** Dmitry V Bagaev, Renske MA Vroomans, Jerome Samir, Ulrik Stervbo, Cristina Rius, Garry Dolton,
8 Alexander Greenshields-Watson, Meriem Attaf, Evgeny S Egorov, Ivan V Zvyagin, et al. Vdjdb
9 in 2019: database extension, new analysis infrastructure and a t-cell receptor motif compendium.
10 *Nucleic Acids Research*, 48(D1):D1057–D1062, 2020.
- 11 **51.** Janine Kah, Sarene Koh, Tassilo Volz, Erica Ceccarello, Lena Allweiss, Marc Lütgehetmann, Antonio
12 Bertoletti, Maura Dandri, et al. Lymphocytes transiently expressing virus-specific t cell receptors
13 reduce hepatitis b virus infection. *The Journal of clinical investigation*, 127(8):3177–3188, 2017.
- 14 **52.** Anangi Balasiddaiah, Haleh Davanian, Soo Aleman, Anna Pasetto, Lars Frelin, Matti Sällberg, Volker
15 Lohmann, Sarene Koh, Antonio Bertoletti, and Margaret Chen. Hepatitis c virus-specific t cell receptor
16 mrna-engineered human t cells: impact of antigen specificity on functional properties. *Journal of
17 virology*, 91(9):e00010–17, 2017.
- 18 **53.** Steven A Rosenberg, Nicholas P Restifo, James C Yang, Richard A Morgan, and Mark E Dudley.
19 Adoptive cell transfer: a clinical path to effective cancer immunotherapy. *Nature Reviews Cancer*,
20 8(4):299–308, 2008.
- 21 **54.** Steven A Rosenberg and Nicholas P Restifo. Adoptive cell transfer as personalized immunotherapy
22 for human cancer. *Science*, 348(6230):62–68, 2015.
- 23 **55.** James L Riley, Carl H June, and Bruce R Blazar. Human t regulatory cell therapy: take a billion or so
24 and call me in the morning. *Immunity*, 30(5):656–665, 2009.
- 25 **56.** Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning.
26 *Journal of Big Data*, 6(1):1–48, 2019.

- 1 **57.** Edward Ma. Nlp augmentation. <https://github.com/makcedward/nlpaug>, 2019.
- 2 **58.** John-William Sidhom, H Benjamin Larman, Drew M Pardoll, and Alexander S Baras. Deeptcr
3 is a deep learning framework for revealing sequence concepts within t-cell repertoires. *Nature*
4 *communications*, 12(1):1–12, 2021.
- 5 **59.** Sofie Gielis, Pieter Moris, Wout Bittremieux, Nicolas De Neuter, Benson Ogunjimi, Kris Laukens,
6 and Pieter Meysman. Detection of enriched t cell epitope specificity in full t cell receptor sequence
7 repertoires. *Frontiers in immunology*, 10:2820, 2019.
- 8 **60.** Michael O Kurrer, Syamasundar V Pakala, Holly L Hanson, and Jonathan D Katz. β cell apoptosis in t
9 cell-mediated autoimmune diabetes. *Proceedings of the National Academy of Sciences*, 94(1):213–218,
10 1997.
- 11 **61.** LGM Van Baarsen, MJH de Hair, TH Ramwadhoebe, IJ AJ Zijlstra, M Maas, DM Gerlag, and
12 PP Tak. The cellular composition of lymph nodes in the earliest phase of inflammatory arthritis.
13 *Annals of the rheumatic diseases*, 72(8):1420–1424, 2013.
- 14 **62.** H Chakir, DE Lefebvre, H Wang, E Caraher, and FW Scott. Wheat protein-induced proinflammatory t
15 helper 1 bias in mesenteric lymph nodes of young diabetes-prone rats. *Diabetologia*, 48(8):1576–1584,
16 2005.
- 17 **63.** Jennie HM Yang, Leena Khatri, Marius Mickunas, Evangelia Williams, Danijela Tatovic, Mohammad
18 Alhadj Ali, Philippa Young, Penelope Moyle, Vishal Sahni, Ryan Wang, et al. Phenotypic analysis
19 of human lymph nodes in subjects with new-onset type 1 diabetes and healthy individuals by flow
20 cytometry. *Frontiers in immunology*, 10:2547, 2019.
- 21 **64.** Hsueh-Ling Janice Oh, Adeline Chia, Cynthia Xin Lei Chang, Hoe Nam Leong, Khoon Lin Ling,
22 Gijsbert M Grotzbreg, Adam J Gehring, Yee Joo Tan, and Antonio Bertoletti. Engineering t cells
23 specific for a dominant severe acute respiratory syndrome coronavirus cd8 t cell epitope. *Journal of*
24 *virology*, 85(20):10464–10471, 2011.
- 25 **65.** Pieter Moris, Joey De Pauw, Anna Postovskaya, Sofie Gielis, Nicolas De Neuter, Wout Bittremieux,
26 Benson Ogunjimi, Kris Laukens, and Pieter Meysman. Current challenges for unseen-epitope

1 tcr interaction prediction and a new perspective derived from image classification. *Briefings in*
2 *Bioinformatics*, 22(4):bbaa318, 2021.

3 **66.** Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
4 *arXiv:1412.6980*, 2014.

5 **67.** Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,
6 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch.
7 2017.

1 Supplementary

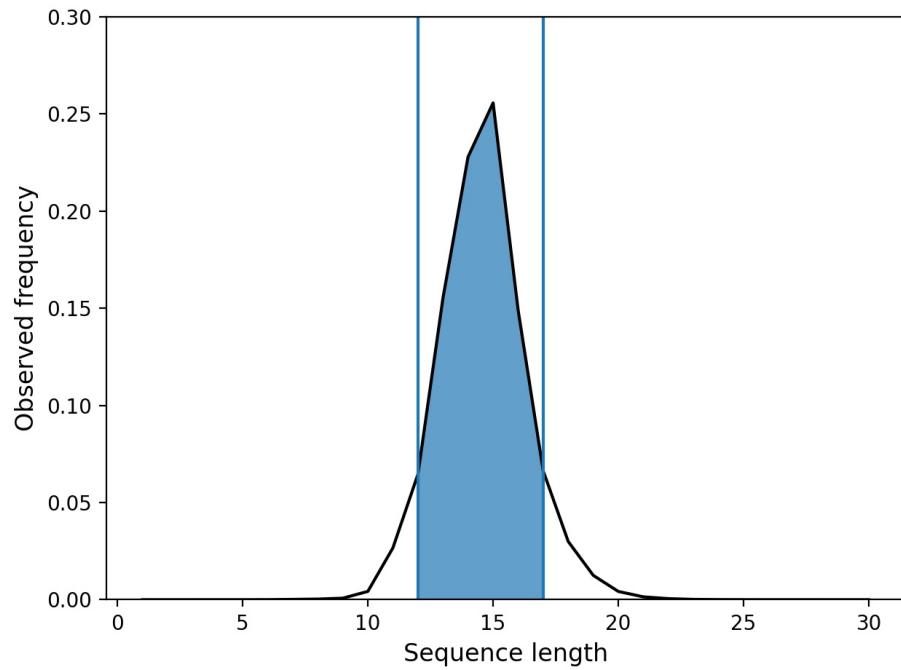


Figure S1. Length distribution of TCR sequences in our constructed universal TCR pool. The lengths of TCRs mainly range from 12 to 17, which take up 84% of the total TCRs (shown in blue shadow). This range is suitable for inferring a deep autoregressive model without causing the issue of gradient explosion or gradient vanishing.

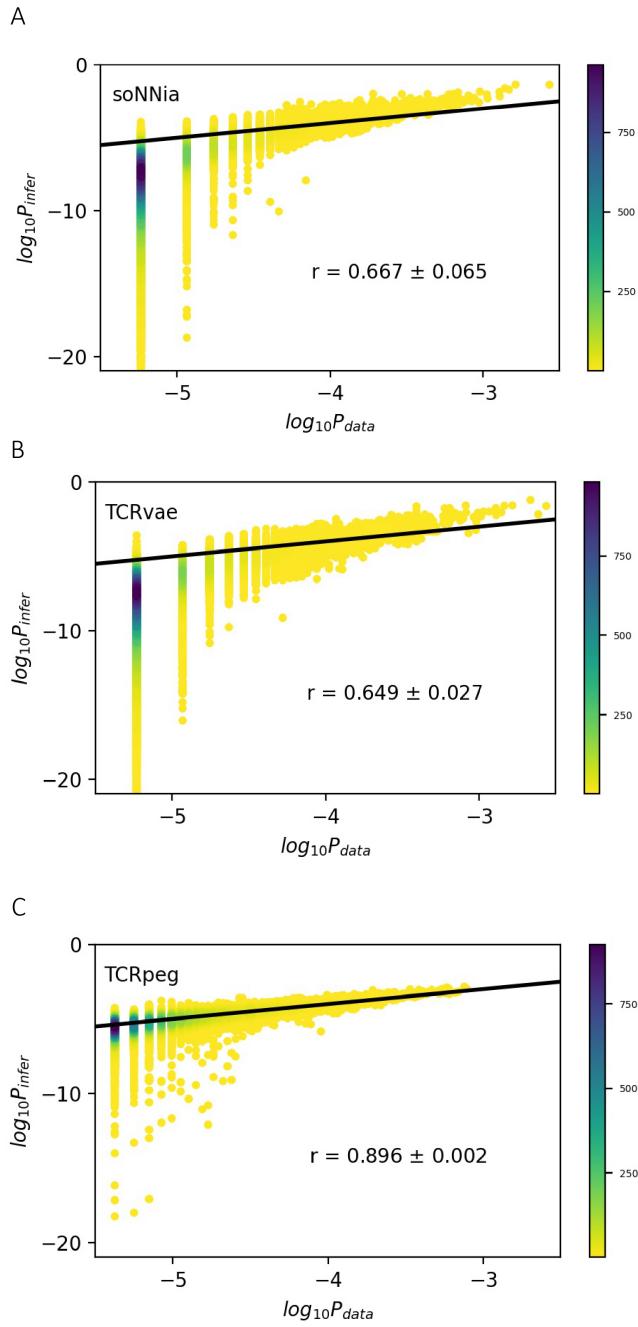


Figure S2. Inference performance of the soNNia, TCRvae, and TCRpeg models in a small proportion of the large universal TCR pool. To assess the stability of these generative models, we re-trained each of them on 200 thousand TCR sequences randomly sampled from the training set. We evaluated them on the test set under the same training settings. The results show that all these models are stable in probability inference, and TCRpeg still surpasses the other two baseline models by large margins.

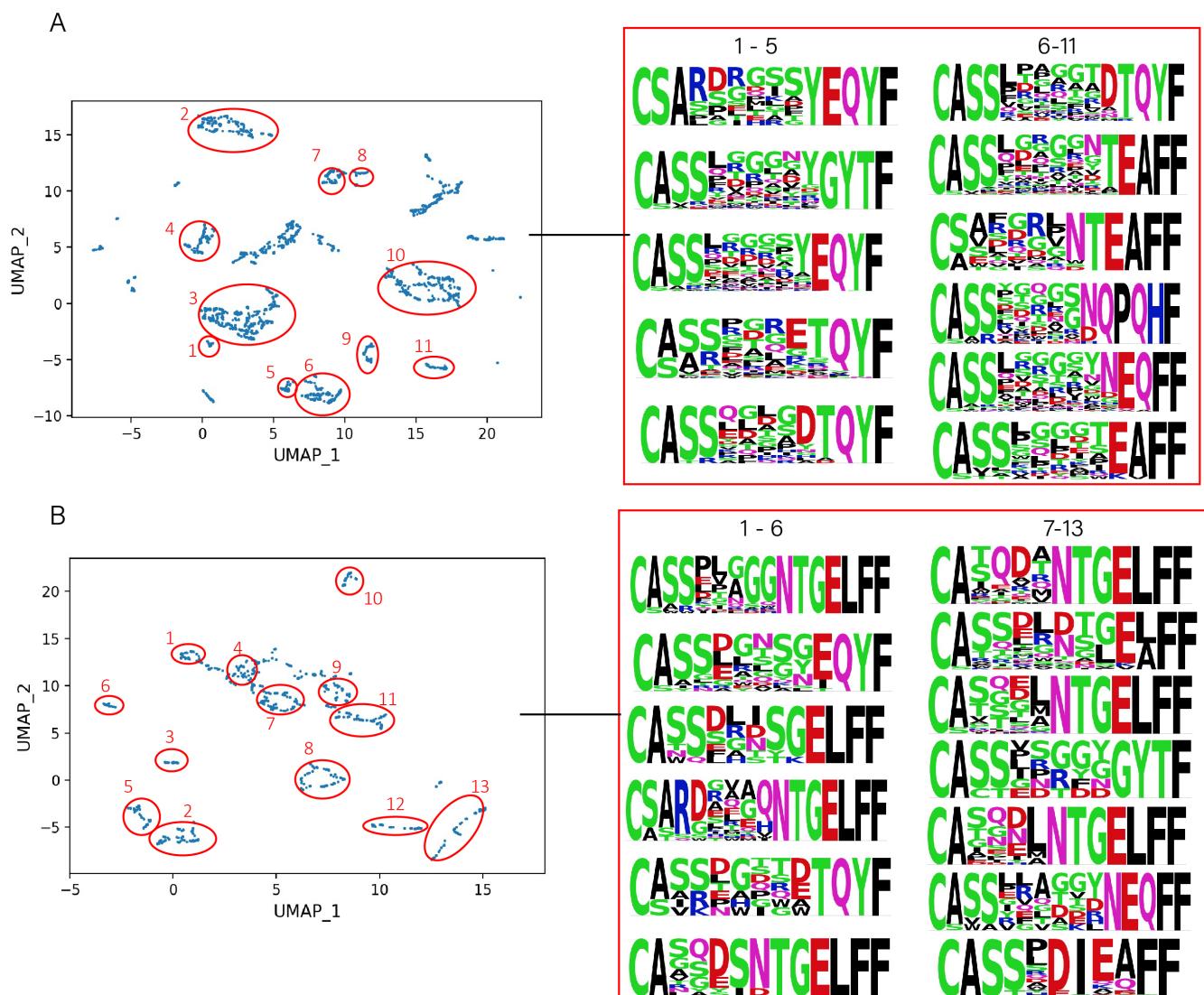


Figure S3. 2D projection map of TCRpeg-based encodings of (A) caTCRs and (B) YLQPRTFLL specific TCRs. We demonstrate more TCR patterns corresponding to different clusters in the projection map that are not shown in Fig. 4 due to the size limit.

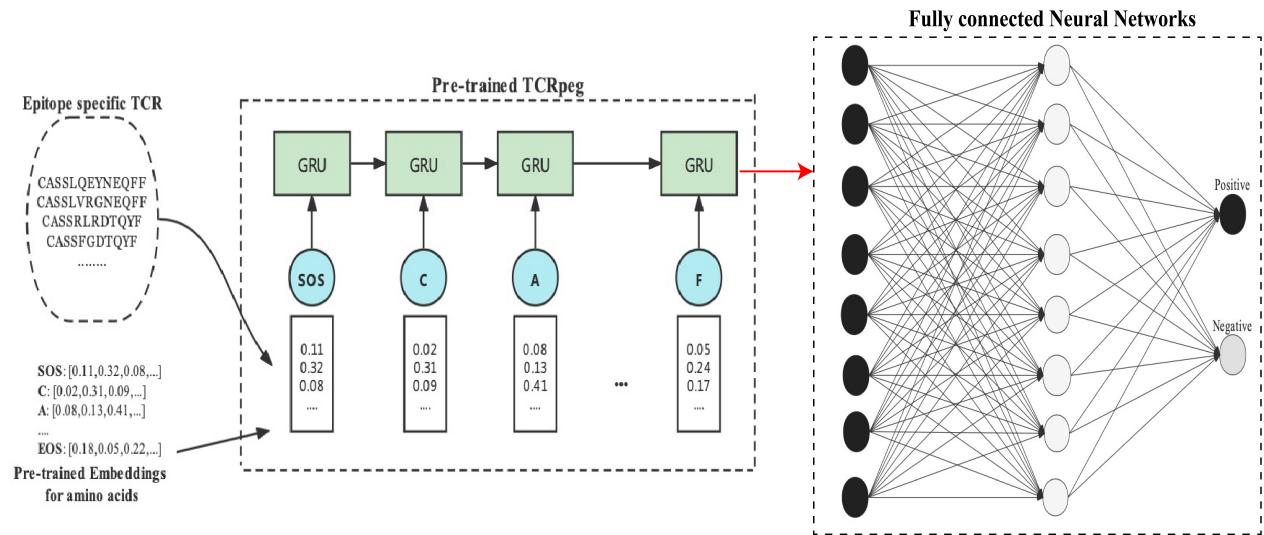


Figure S4. The illustration of the architecture of TCRpeg-c. First, we trained the TCRpeg on epitope-specific TCR sequences (or other specific TCRs) to obtain the numerical encodings for each TCR. Then these encodings were inputted into a fully connected neural network for classification purposes.

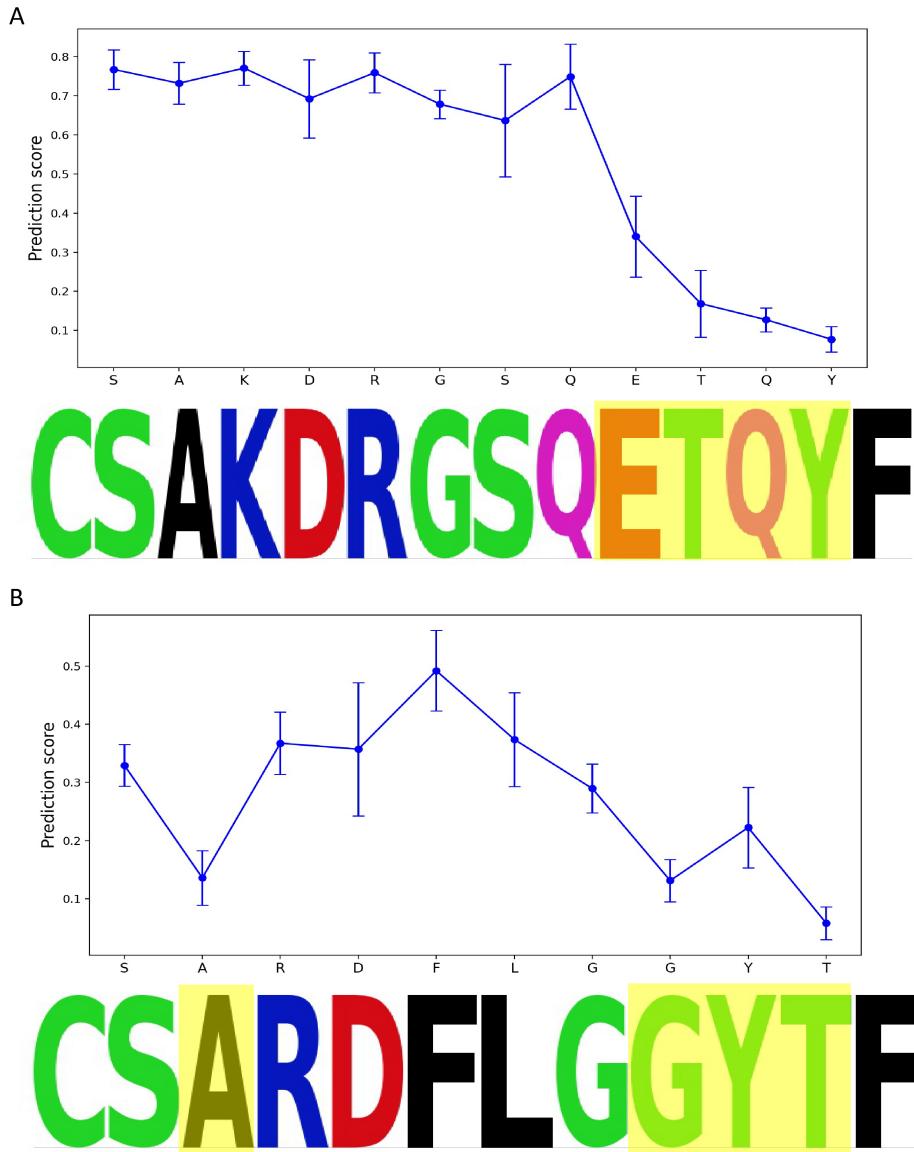


Figure S5. Perturbation analysis to discover *de novo* motifs. In addition to validating previously identified motifs, we performed a perturbation analysis on some TCRs to find potential motifs. Blue curves indicate the prediction scores from TCRpeg-c for the permuted TCR sequences. We observe dropoffs in the region of “ETQY” within the TCR sequence “CSAKDRGSQETQYF” (**A**) and “GYT” as well as “A” within “CSARDFLGGYTTF” (**B**). We wish that when more YLQPRTFLL-specific TCRs are available, these two motifs could be identified through the TCR similarity network.

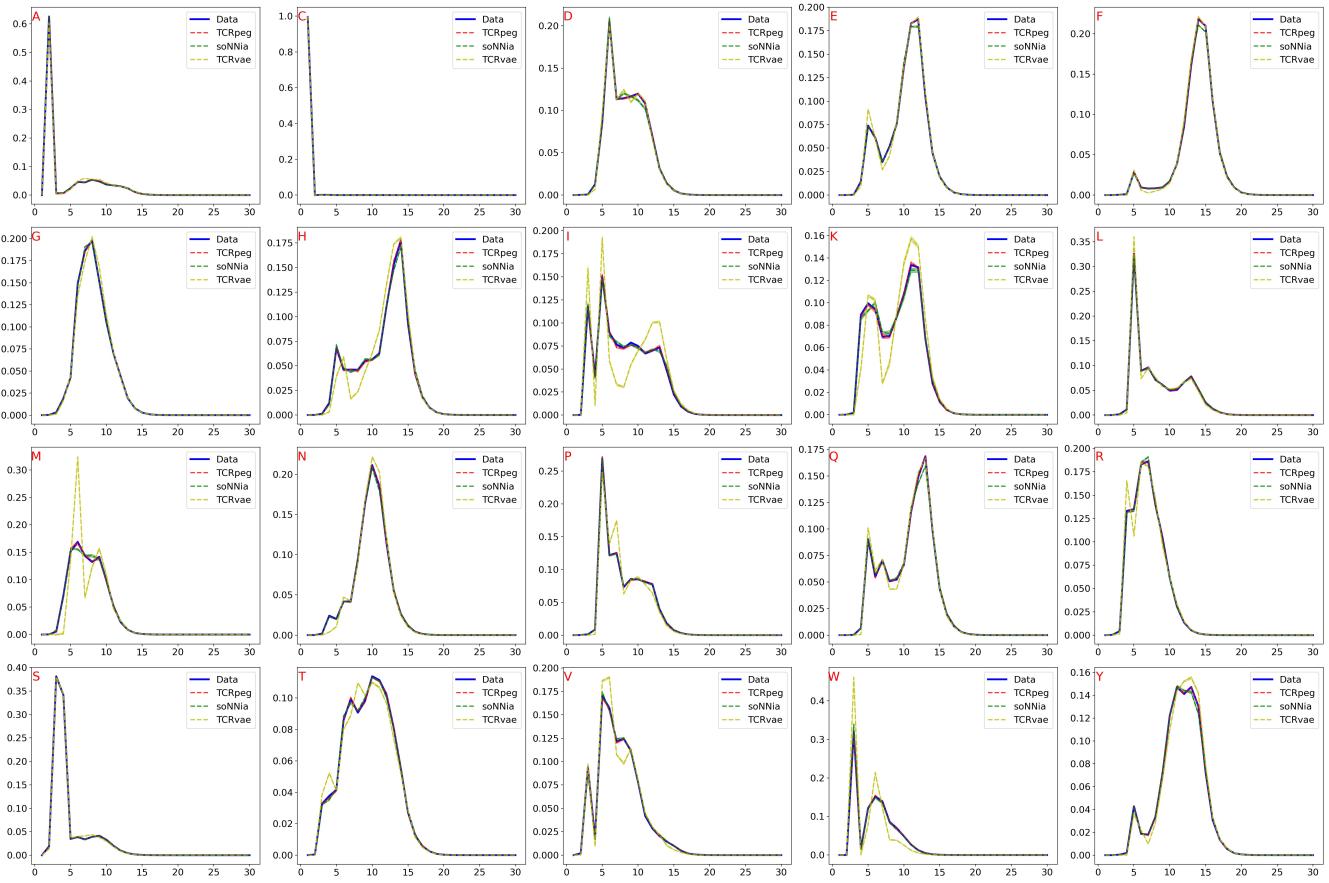


Figure S6. The position distributions of amino acids for the generated TCR sequences compared to TCRs from the universal TCR pool. Each subplot shows the position distribution of a particular amino acid. The x-axis represents the TCR sequence's position index (from left to right), and the y-axis shows the amino acid frequency appearing at each position index. The TCRvae model performs the worst in this task. The generated TCR sequences using the soNNia and TCRpeg models possess a close distribution to the observed data, and TCRpeg performs slightly better than soNNia.

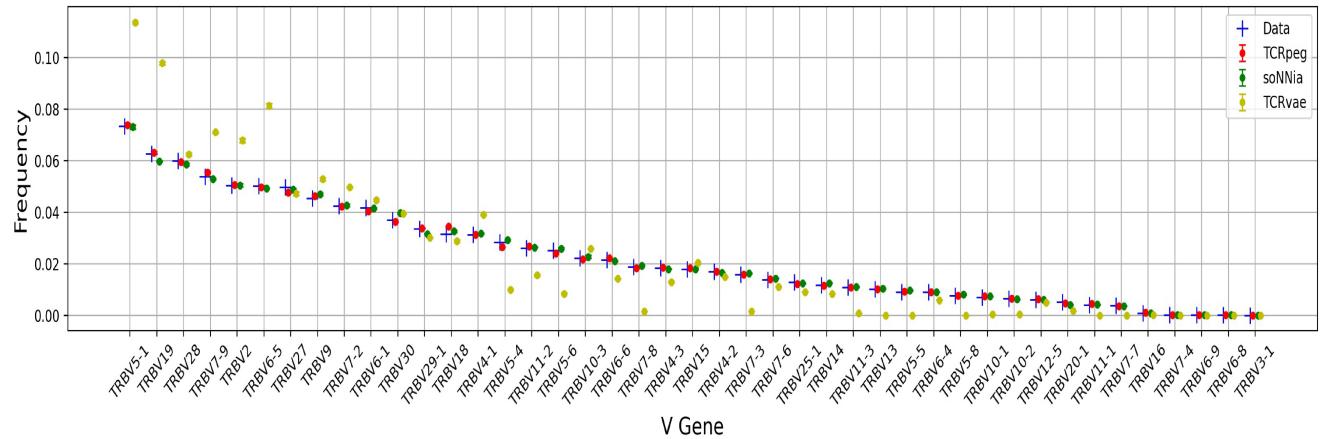


Figure S7. The full version of Fig. 5A, showing the statistical distributions of all 43 V genes from the generated repertoires using the three generative models.

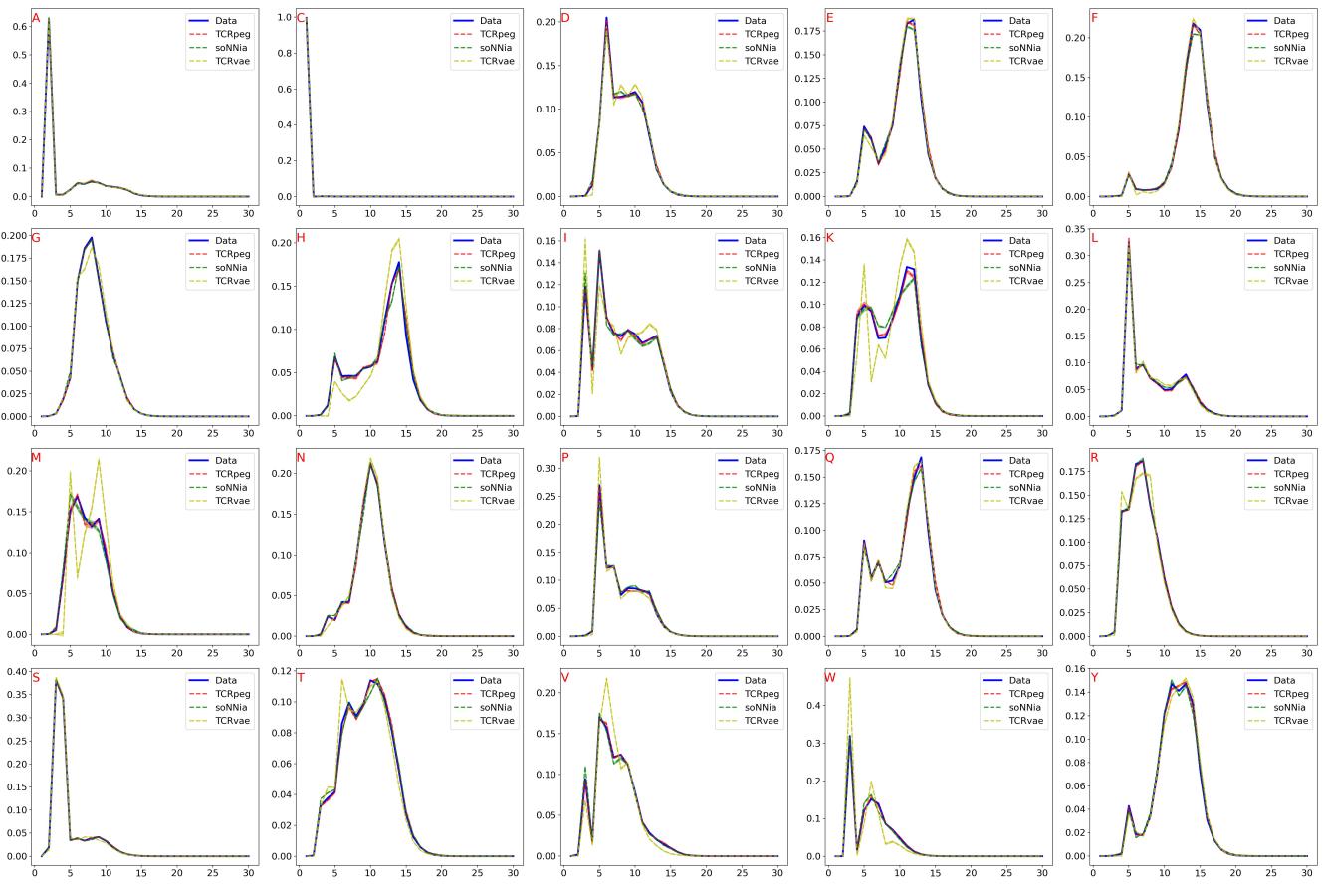
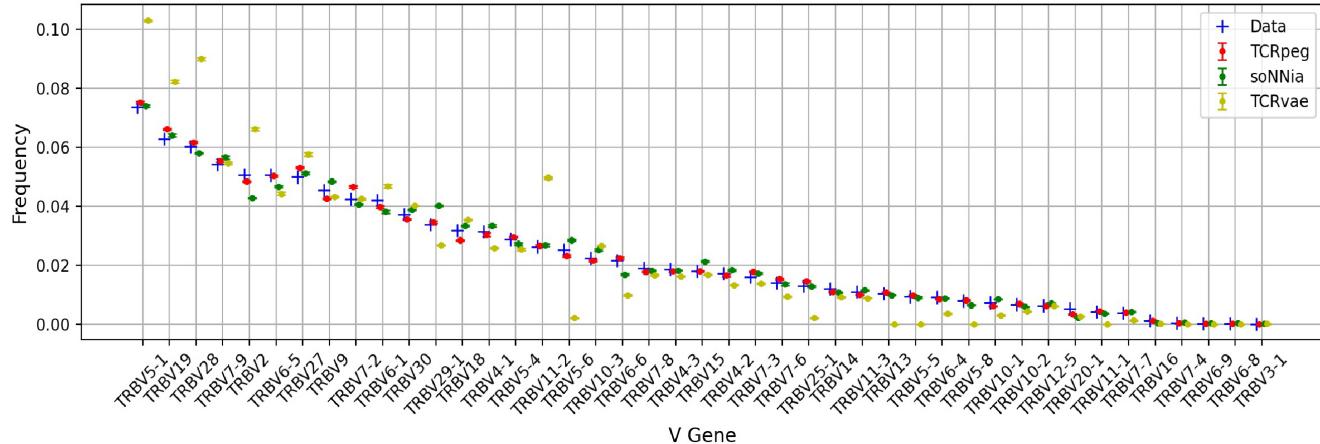
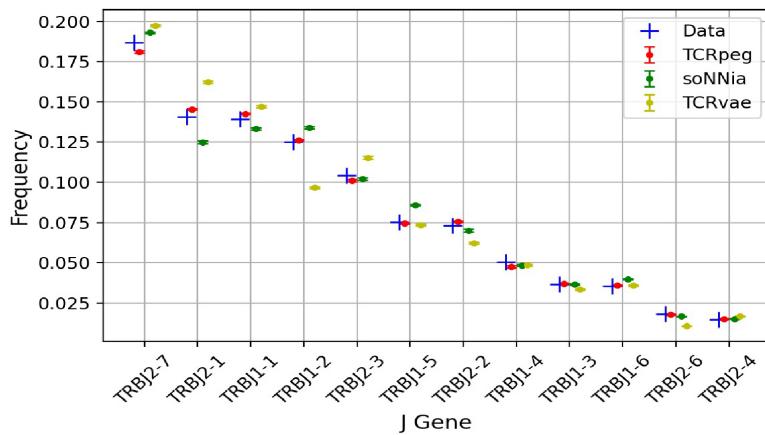


Figure S8. Position distributions of amino acids for the generated TCR sequences compared to TCRs from the universal TCR pool. We randomly sampled 200 thousand TCRs from the TCR pool and trained each model on this subset. TCRpeg still achieved the best performance with $r \simeq 0.999$ compared to $r \simeq 0.998$ for soNNia and $r \simeq 0.982$ for TCRvae.

A



B



C

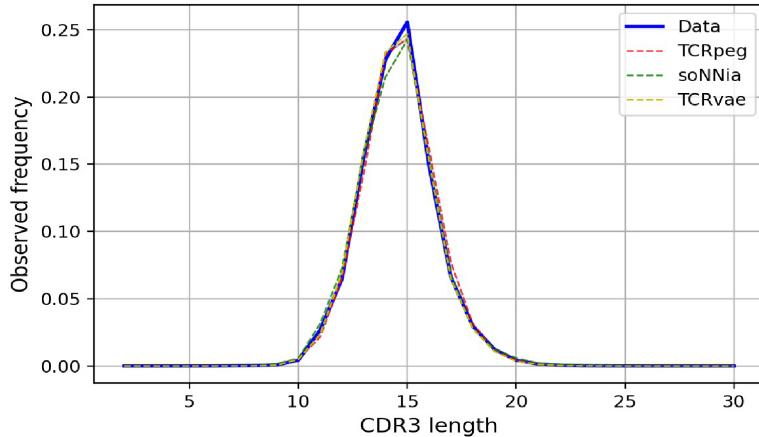


Figure S9. Comparisons of (A) V gene, (B) J gene, and (C) length distributions between generated and real sequences when the generative models are inferred from a small subset of the universal TCR pool. Here, we randomly sampled 200 thousand TCRs from the TCR pool and trained each model on this subset. TCRpeg still outperformed the other two models in V and J gene usage distribution with $r \simeq 0.997, 0.998$ compared to $r \simeq 0.993, 0.992$ for soNNia and $r \simeq 0.951, 0.980$ for TCRvae. For length distribution, TCRvae achieved the best performance with $r \simeq 0.992$ compared to $r \simeq 0.990$ for TCRpeg and $r \simeq 0.982$ for soNNia.

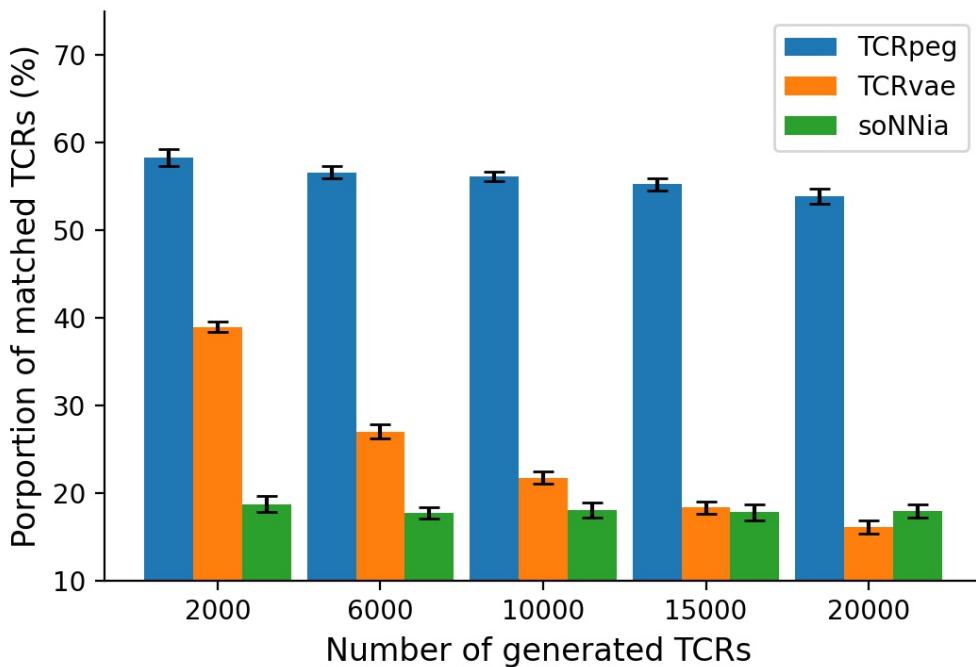


Figure S10. Investigation of epitope specificity for TCRpeg-generated TCR sequences using the TCRMatch software. TCRMatch is a software package for predicting TCR specificity based on sequence similarity to previously characterized TCRs. In this experiment, we set the similarity cutoff to 0.9 and assume that those generated TCRs with similarity scores higher than 0.9 to any TCR in the inferring repertoire share the same epitope specificity. Specifically, we inferred the three generative models on the YLQPRTFLL-specific TCRs and used the trained models to generate new TCRs. We observed that more than 50% different numbers of generated TCRs possess epitope specificity to YLQPRTFLL using TCRpeg. These results support our claim that the TCRpeg-generated TCRs may share a hidden similarity to the TCRs used to train the model. In contrast, only less than 40% proportion of the soNNia- and TCRvae- generated sequences possess epitope specificity.

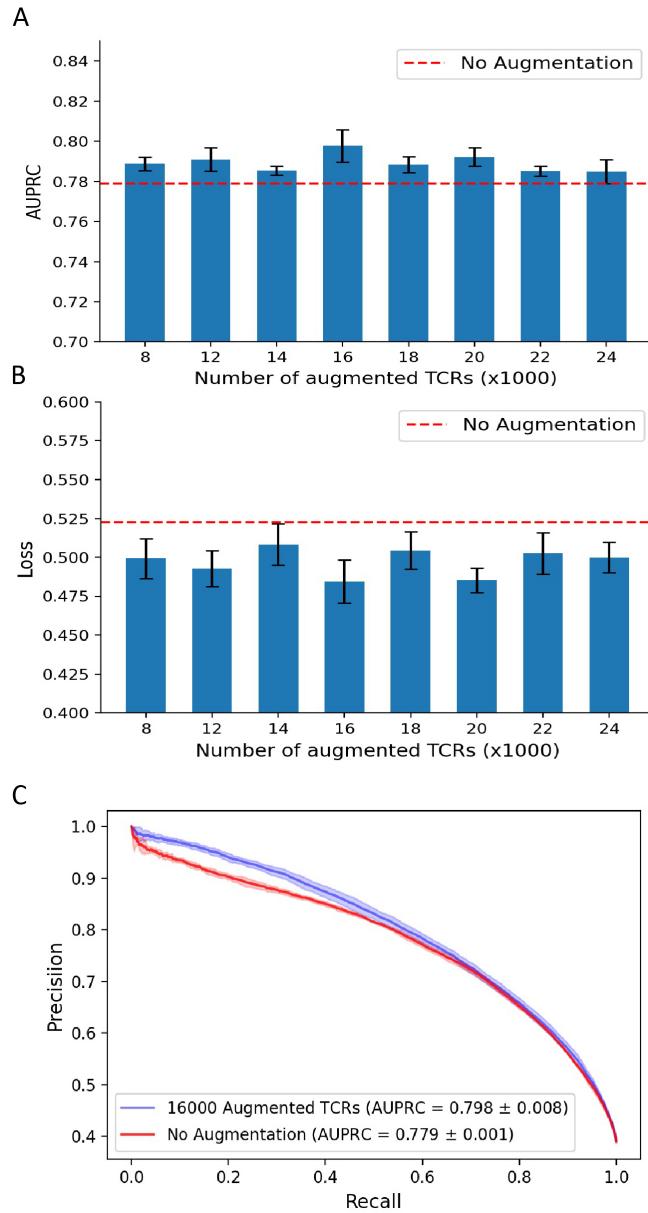


Figure S11. Effects of using the TCR-specific data augmentation technique for classifying caTCRs. **(A)** The area under the precision-recall curve (AUPRC) values when data augmentation is applied using the TCRpeg-c model. Within a large range of the number of generated TCR sequences, we observe an enhancement of AUPRC value. **(B)** The binary cross-entropy loss (BCE loss) that is evaluated on the test set. When data augmentation is applied, test loss decreases, which is a positive sign of alleviation of the overfitting problem. **(C)** The precision-recall curve when the number of augmented TCRs is 16,000 where AUC, AUPRC, and the reduction of test loss achieve the highest.

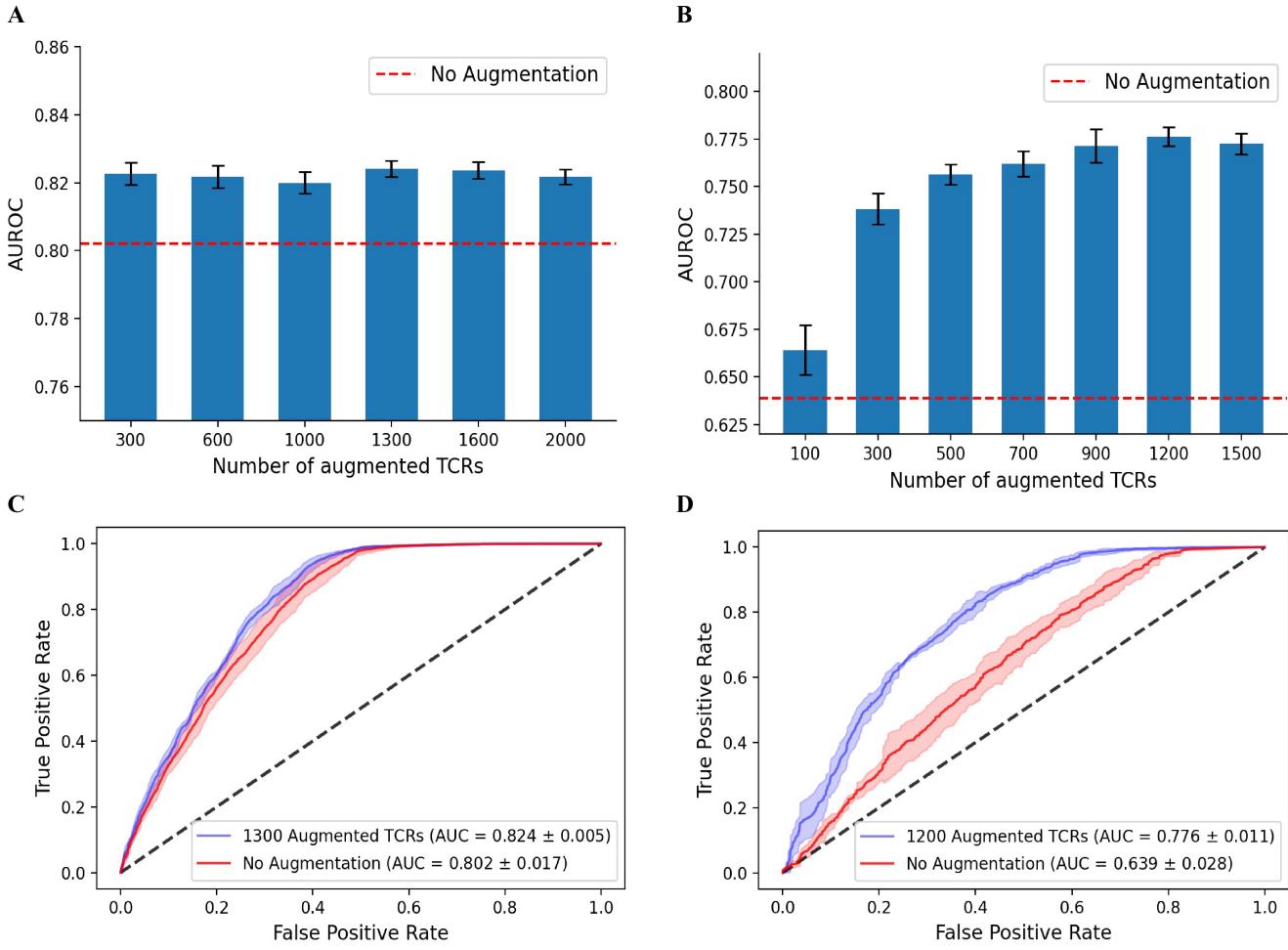


Figure S12. Validation of the utility of the TCRpeg-based data augmentation technique using the TCReX model. We collected epitope-specific TCRs for GILGFVFTL and GLCTLVAML from the VDJdb database (downloaded on Mar 21, 2022) with 3406 and 962 positive samples. We randomly sampled the same number of TCRs from the universal TCR pool as the control samples. We trained TCReX in the default setting in all experiments. The TCRpeg model with hidden size and number of layers set to 256 and 3 was trained for 30 epochs and then used to generate new TCRs as augmentation. We show the AUC values when applying data augmentation with a different number of augmented TCRs to the classification of (A) GILGFVFTL and (B) epitope-specific TCRs of GLCTLVAML. The receiver operating characteristic curves for the GILGFVFTL and GLCTLVAML with the highest performance boost are shown in (C) and (D) respectively. We observe that applying data augmentation can bring about 2.1% and 21.4% AUC enhancement.