



OPEN ACCESS

EDITED BY

Penelope Anne Morel,
University of Pittsburgh, United States

REVIEWED BY

Alok Vishnu Joglekar,
University of Pittsburgh, United States
Kevin Michalewicz,
Imperial College London, United Kingdom

*CORRESPONDENCE

Jason T. George
✉ jason.george@tamu.edu

†These authors have contributed
equally to this work and share
first authorship

RECEIVED 12 October 2024

ACCEPTED 24 December 2024

PUBLISHED 23 January 2025

CITATION

Teimouri H, Ghoreyshi ZS, Kolomeisky AB and
George JT (2025) Feature selection
enhances peptide binding predictions for
TCR-specific interactions.
Front. Immunol. 15:1510435.
doi: 10.3389/fimmu.2024.1510435

COPYRIGHT

© 2025 Teimouri, Ghoreyshi, Kolomeisky and
George. This is an open-access article
distributed under the terms of the [Creative
Commons Attribution License \(CC BY\)](#). The
use, distribution or reproduction in other
forums is permitted, provided the original
author(s) and the copyright owner(s) are
credited and that the original publication in
this journal is cited, in accordance with
accepted academic practice. No use,
distribution or reproduction is permitted
which does not comply with these terms.

Feature selection enhances peptide binding predictions for TCR-specific interactions

Hamid Teimouri^{1,2†}, Zahra S. Ghoreyshi^{2,3†},
Anatoly B. Kolomeisky^{1,2,4} and Jason T. George^{2,3,5,6*}

¹Department of Chemistry, Rice University, Houston, TX, United States, ²Center for Theoretical Biological Physics, Rice University, Houston, TX, United States, ³Department of Biomedical Engineering, Texas A&M University, College Station, TX, United States, ⁴Department of Chemical and Biomolecular Engineering, Rice University, Houston, TX, United States, ⁵Department of Hematopoietic Biology and Malignancy, MD Anderson Cancer Center, Houston, TX, United States, ⁶Department of Translational Medical Sciences, Texas A&M Health Science Center, Houston, TX, United States

Introduction: T-cell receptors (TCRs) play a critical role in the immune response by recognizing specific ligand peptides presented by major histocompatibility complex (MHC) molecules. Accurate prediction of peptide binding to TCRs is essential for advancing immunotherapy, vaccine design, and understanding mechanisms of autoimmune disorders.

Methods: This study presents a theoretical approach that explores the impact of feature selection techniques on enhancing the predictive accuracy of peptide binding models tailored for specific TCRs. To evaluate our approach across different TCR systems, we utilized a dataset that includes peptide libraries tested against three distinct murine TCRs. A broad range of physicochemical properties, including amino acid composition, dipeptide composition, and tripeptide features, were integrated into the machine learning-based feature selection framework to identify key properties contributing to binding affinity.

Results: Our analysis reveals that leveraging optimized feature subsets not only simplifies the model complexity but also enhances predictive performance, enabling more precise identification of TCR peptide interactions. The results of our feature selection method are consistent with findings from hybrid approaches that utilize both sequence and structural data as input as well as experimental data.

Discussion: Our theoretical approach highlights the role of feature selection in peptide-TCR interactions, providing a quantitative tool for uncovering the molecular mechanisms of the T-cell response and assisting in the design of more advanced targeted therapeutics.

KEYWORDS

immune response, feature selection, physicochemical properties, TCR-peptide interactions, binding affinity

1 Introduction

The host adaptive immune response, primarily driven by the activation of T-cells, orchestrates a precise and targeted defense by recognizing and responding to specific antigens (1, 2). T-cell receptors (TCRs) interact with peptide-major histocompatibility complex (MHC) through low-affinity, transient contacts, allowing them to identify the correct antigen while remaining sensitive to subtle molecular differences (3, 4). This low-affinity binding also allows for TCR cross-reactivity with diverse peptide sequences, broadening their recognition potential (2, 5). Accurately predicting peptide binding to specific TCRs is crucial for advancing immunotherapy and vaccine development, and for clarifying the underlying microscopic picture of immune response (5–8). However, this task remains very complex due to the immense variability of TCRs and peptides, in addition to the intricate nature of molecular mechanisms governing their binding affinities (9–12).

There are multiple experimental techniques available for investigating TCR-peptide interactions, including crystallography (13), surface plasmon resonance (14), and yeast display systems (5, 15–17). Recent advances include immunopeptidomics, which identifies naturally presented peptides on MHC complexes via mass spectrometry (18). Compared to other methods, yeast display offers a unique high-throughput advantage, allowing the screening of a large number of peptide variants simultaneously, which enables the rapid identification of high-affinity interactions. In a recent study, a yeast surface display system was developed to screen highly diverse libraries of peptides presented by MHC molecules, identifying those capable of binding specific TCRs (5). By coupling this approach with deep sequencing, the sequence diversity of peptides recognized by different TCRs was mapped, which helped to uncover critical binding motifs and interactions. After multiple rounds of selection, this dataset was refined to identify hundreds of high-affinity TCR-peptide interaction, which offers opportunities for identifying pertinent molecular features of TCR-peptide interactions by applying additional computational techniques.

The use of machine learning methods for predicting TCR-peptide interactions is a promising research direction that has the potential of studies that has the potential to overcome the limitations of traditional methods (19–24). Traditional experimental methods are pivotal but limited by being time-intensive, low-throughput, and reliant on structural data. In contrast, machine learning leverages sequence-based and physicochemical features, enabling scalable, high-throughput analyses and uncovering patterns in TCR binding specificity and cross-reactivity that are difficult to capture experimentally. By leveraging large datasets and incorporating structural, physicochemical, and sequence information, these models can learn the underlying principles that govern TCR specificity and binding affinity (25). Developing machine learning models to predict strong binder peptides for specific TCRs, however, involves several key challenges (25), including TCR cross-reactivity, wherein a single TCR can bind multiple peptides (25–27). This

property complicates the identification of true strong binders versus weaker ones, as a peptide that strongly binds to one TCR may bind weakly to another. One initial step for addressing this issue is to develop context-specific models to identify features that drive specificity in distinct functional scenarios, such as TCRs restricted to a common MHC allele that bind diverse peptide antigen sets. Additionally, quantification of such features in the context of the diversity to which strong binders are themselves identified, represents an important quantification of cross-reactivity within a given TCR system.

Recent computational frameworks (28–34), have made notable advancements in predicting TCR-peptide binding affinities. The Rapid Coarse-grained Epitope TCR (RACER) model is one particular example of a hybrid structural and sequence-based approach that uses a pairwise energy model trained on deep sequencing and crystallographic data to identify strong and weak TCR-peptide binders. RACER provides highly useful predictions by leveraging a sparse yet diverse set of experimentally determined TCR-peptide structures, enabling this model to generalize effectively across a wide range of cases. Furthermore, RACER and other similar predictive models aim to benefit from leveraging biophysical features to subsequently require a reduced number of positive (1) and negative (0) examples in training. These models are trained on a collection of TCR-peptide systems and can handle variations in either the TCR or peptide. In our study, RACER was employed to analyze datasets restricted to the same (mouse IE^K) MHC-II allele, where it was previously used to resolve strong and weak TCR binding profiles. However, these models do not comprehensively evaluate all of the variability within a given binding class when, for example, we have a fixed TCR and a large number of confirmed strong and weak binding peptide sequences corresponding to that single TCR. Incorporating an ML-based classifier could complement models like RACER by extracting more nuanced, context-specific features that confer binding specificity. This synergistic approach could improve predictive accuracy within specific binding classes, enhancing our understanding of TCR-peptide interactions. Additionally, this integration may reveal novel feature interactions that are critical for binding specificity, offering valuable insights for experimental validation.

Our study aims to apply machine learning techniques with feature selection to improve the accuracy of TCR-peptide specificity prediction to identify motifs that most highly resolve strong and weak binders given available large-scale binding datasets. By identifying features among a comprehensive set of physicochemical features that determine binding interactions, our model effectively distinguishes between strong binders and weak binders. We apply various feature extraction techniques and examine the robustness of each approach. To test our theoretical method, the analysis is applied to three distinct peptide pools. The model's ability to account for meaningful peptide variations that drive specificity is evaluated for each case, based on successful predictions of strong- and weak-binding TCR-peptide pairs.

2 Materials and methods

2.1 Dataset and data preprocessing

We employed a highly diverse set of peptide-MHC complexes derived from yeast-displayed peptide-MHC libraries, which includes three distinct types of murine TCRs: 2B4, 226, and 5cc7 (5). These TCRs were selected due to their distinct mechanisms of peptide recognition, which arise from variations in their structure, cross-reactivity, and interactions with MHC molecules (35). Specifically, the selected TCRs represent varying levels of cross-reactivity, with 2B4 exhibiting high specificity, 226 demonstrating broad cross-reactivity, and 5cc7 showing intermediate behavior. This diversity enabled us to assess the performance of our feature selection approach across different binding contexts. For each TCR, antigen libraries comprised of peptides having fixed length (13-amino acids) were subjected to multiple rounds of selection to enrich for TCR-binding peptides, which were subsequently analyzed using deep sequencing. The final dataset consisted of hundreds of unique peptide sequences, each characterized by specific TCR recognition motifs. In each dataset, each peptide sequence is assigned a “Post-selection” enrichment score following 5 rounds of affinity-driven selection. Highly enriched post-selection peptides are the ones that bind most strongly to the TCR in study. During each round of selection, the weaker binding peptides are gradually filtered out, and the frequency of peptides with stronger TCR affinities increases. Five rounds of affinity-based selection ultimately yield a list of sequences with their corresponding abundance, which indicates how many times that particular peptide was detected during sequencing and is proportional to their binding affinity.

Since the post-selection enrichment score gives a direct measure of how strongly and frequently each peptide binds to the TCR, we can classify peptides into strong binders (class 1) and weak binders (class 0) based on post-selection enrichment by setting a threshold value calculated from the RACER model. To calculate this threshold, a subset of 500 cases was selected from each dataset based on post-selection read counts. In this way, peptides with the highest read counts were selected from each dataset. This subset was chosen based on peptide quantities in Round 5, measured as ‘reads’ following the replication of yeast cells displaying the respective pMHCs. Peptides with the highest reads were selected, ensuring that all chosen peptides had high quantities, as these were likely to contain the majority of strong binders. Specifically, 140 instances with the highest reads were designated as strong binders for training, while the remaining 360 cases, which also had high but slightly lower reads, were used for testing. Peptides outside this subset were excluded because their reads were either zero or one, making them less relevant for separating strong and weak binders. For each strong binder, we generated a set of 1,000 decoy sequences by randomizing the peptide sequence and pairing it with the corresponding TCR structure. This approach created a comprehensive negative dataset to balance the strong binders, as previously established in (29). The inclusion of 1,000 decoys for each strong binder ensures robust statistical differentiation between

strong and weak binders, addressing challenges like global sparsity and enhancing the model’s predictive accuracy. This process resulted in a total of 140,000 decoy binders. The remaining cases were allocated for testing. We then applied RACER to compute thresholds that effectively separate the distributions of strong and weak binders for each TCR-pMHC case. The RACER model calculates binding energies by integrating high-throughput data from previously confirmed TCR-peptide interactions and crystal structures to train a residue-specific energy matrix. Key to this optimized energy matrix is in RACER’s training on confirmed strong and weak binding TCR-antigen pairs. We apply the RACER framework previously constructed in describing antigen-specific responses for 2B4, 226, and 5cc7 TCRs as demonstrated previously (28). This energy matrix, combined with available structural templates, is used to quantify TCR-peptide binding affinities. For our peptides of interest, we utilized crystal structures with PDB IDs 3QIB, 3QIU, and 4P2R corresponding to 2B4, 226, and 5cc7 TCRs, respectively.

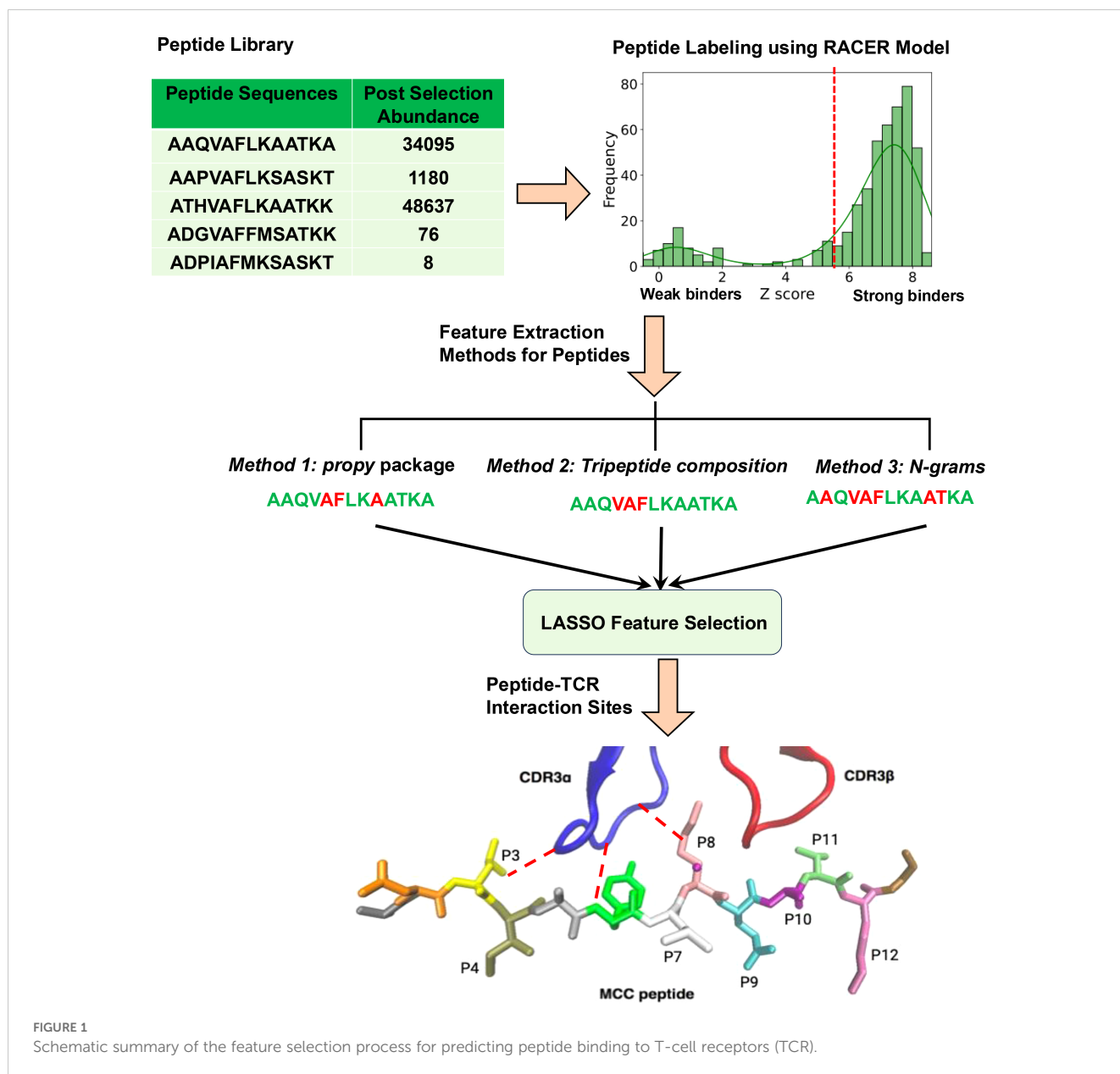
We predicted the binding energies and corresponding Z-scores for training and testing cases. To determine a partition threshold separating strong and weak binders, we analyzed histograms of the Z-scores for all 500 cases to identify the peaks of the distributions for both classes (Figure 1). Supplementary Figure S1 presents the histograms for peptide libraries targeting 2B4, 226, and 5cc7 TCRs. The threshold was defined as the midpoint between the peaks of the strong and weak binder distributions. Finally, we ranked all 500 cases in the dataset based on their Z-score values which are obtained by the RACER model, mapped each partition threshold to the corresponding post-selection enrichment score (5). The specific thresholds for each case are listed in Table 1. Using the calculated partition threshold values for each dataset, all peptides were partitioned into class 1 (strong binder) and class 0 (weak binder) cases. Since the dataset was highly imbalanced, with far more weak binders than strong binders, we performed controlled undersampling of the weak binders, which consisted of randomly selecting a subset of weak binders equal in size to the number of strong binders. This method helped mitigate the effects of class imbalance to improve the performance of the machine learning classifiers (36).

2.2 Extraction of physicochemical features for peptides

To capture the physicochemical properties of peptides that are crucial for their interaction with TCRs, we extracted multiple features using primary amino acid sequence information. This is a critical step for implementing machine learning models. To evaluate the robustness of our feature selection technique, we employed three different methods for extracting physicochemical features from the sequence data.

2.2.1 Propy package

First, we extracted a comprehensive set of physicochemical features from the amino acid sequence of each peptide using the



propy package (37). These features are broadly categorized into different groups, including charge, amino acid composition, dipeptide composition, autocorrelations, chemical composition features, and sequence order information. The physicochemical features generated using *propy* package have been utilized in a wide

TABLE 1 Summary of selected peptide datasets associated with each TCR.

TCR	Threshold	Peptides (0/1)
2B4	13	98/98
226	6	987/987
5c7	23	234/234

Data obtained from Ref (5). Partition thresholds used to distinguish strong into strong (1) and weak (0) binders after five rounds of affinity-based selection were obtained using the RACER model (28) (see text for details).

range of machine learning models, including classification of antimicrobial peptides (38, 39) and predicting protein-protein interactions (40).

Among the features extracted by *propy*, amino acid composition and dipeptide composition are particularly important for understanding the interactions between TCRs and various peptides, as they provide insights into how specific residues and their combinations influence binding affinity (25). For a peptide of L residues, amino acid composition, which represents the fraction of each amino acid type, reads as

$$f_i = \frac{N_i}{L}, \quad i = 1, 2, 3, \dots, 20 \quad (1)$$

where N_i is the number of amino acids of type i . Since there are 20 amino acids, the amino acid composition comprises 20 features among the *propy* features.

Similarly, the dipeptide composition represents the fraction of each possible dipeptide within the peptide, calculated as:

$$f_{ij} = \frac{N_{ij}}{L-1}, \quad i, j = 1, 2, 3, \dots, 20 \quad (2)$$

where N_{ij} is the number of dipeptides consisting of amino acids of type i, j . Consequently, the dipeptide composition contributes $20^2 = 400$ distinct features to the set of *propy* features.

2.2.2 Tripeptide composition

Tripeptide composition represents the fraction of each possible tripeptide (formed by three consecutive amino acids) within a peptide sequence. The tripeptide composition is calculated as:

$$f_{ijk} = \frac{N_{ijk}}{L-2}, \quad i, j, k = 1, 2, 3, \dots, 20 \quad (3)$$

Where N_{ijk} is the number of tripeptides containing amino acids of type i, j, k . Tripeptide composition, which comprises $20^3 = 8000$ features, provides deeper insight into the peptide's structure by capturing the frequency of every unique combination of three consecutive amino acids. Since the *propy* package does not provide tripeptide composition features by default, we extracted these features separately. Tripeptide motifs were calculated by iterating through peptide sequences to count the occurrences of all 8000 possible tripeptides. These counts were subsequently normalized to generate relative frequencies, as detailed in Equation 3.

2.2.3 N-gram language model

A sequence of amino acids, whether forming a short peptide or a large protein, can be viewed as a text document, where amino acids function as the fundamental units, analogous to words (41). Text mining and natural language processing have been previously employed for bioinformatics applications such as protein clustering and classification, protein-protein interaction (PPI) prediction, protein folding analysis, and non-coding RNA identification (42, 43).

To analyze amino acid sequences using natural language processing methods, we can use the *N*-gram language model, which is a probabilistic model used to predict the next item in a sequence based on the preceding items. An *N*-gram is a contiguous sequence of *N* items from a given sequence of text. In our context, each amino acid represents an item (analogous to a word), and each *N*-gram represents a sequence of *N*-amino acids (analogous to a sentence). While *propy* can efficiently compute the frequency of single amino acids and dipeptides, the resulting dipeptide frequencies tend to be sparse, as many dipeptides may not appear in a given peptide. By incorporating common *N*-grams including unigrams (single amino acids), bigrams (pairs of amino acids), and trigrams (triplets), the model goes beyond mere composition analysis and captures the sequential order and local motifs within peptides (44). Moreover, the robustness of the overall predictive model can be enhanced by combining different types of amino acid composition, including unigrams, bigrams, and trigrams. This approach ensures that if one feature set fails to capture critical patterns, the other can compensate, leading to a more comprehensive and accurate analysis of TCR-peptide interactions.

Since all peptides are composed of 20 standard amino acids, the maximum vocabulary sizes for unigrams, bigrams, and trigrams are 20, $20^2 = 400$, and $20^3 = 8000$, respectively. This creates a fixed-size vocabulary that can be represented as a numerical feature vector, where each element corresponds to the frequency or presence of a specific *N*-gram in the sequence. The process of vectorizing a peptide sequence using the *N*-gram approach begins by breaking down each peptide into unigrams, bigrams, and trigrams, which serve as the fundamental building blocks of the sequence. Next, a complete vocabulary is composed of all possible *N*-grams that can occur within the sequence. Once the vocabulary is established, the sequence is vectorized by converting the frequency of each *N*-gram into a numeric vector. The resulting vectorized transformation, which has lower sparsity compared to features generated by *propy*, enables efficient processing of peptide sequences by machine learning algorithms.

2.3 Data normalization

For each peptide, the quantitative values of the physicochemical properties extracted from the methods described above have different numerical scales. It is important to initially re-scale all these values to fall between 0 and 1 so that every property is considered with a similar weight. To normalize this quantity to be in the range 0 and 1, we use the following re-scaling expression,

$$\hat{z} = \frac{(z - z_{min})}{(z_{max} - z_{min})}, \quad (4)$$

where z is the original value of the physicochemical property, z_{min} and z_{max} are limiting values for this property for all considered proteins, and \hat{z} is the normalized one that is specifically utilized in the analysis. It is important to note that to prevent leakage from the training set to the test set, we performed data normalization only after splitting the datasets into training and test sets.

2.4 Feature selection

In studying TCR-peptide interactions, our primary goal is to identify which specific physicochemical features – such as amino acid properties or sequence motifs – are most important for distinguishing between strong and weak binders. However, the extracted feature set often consists of high-dimensional data, meaning the number of features may exceed the available data, with some being irrelevant or highly correlated. Using all these features without selection can result in overfitting, where the model learns noise rather than meaningful correlations, reducing its predictive performance (38). To mitigate this issue, we employ LASSO (The Least Absolute Shrinkage and Selection Operator) techniques that mathematically assign zero weights to irrelevant or redundant features (38, 45). We note that features with non-zero weights are considered relevant, and the magnitude of these weights provides a measure of their relative importance in the predictive model. This property enables the identification of key features

contributing to TCR-peptide binding interactions. We note that while *propy* and N-gram features are less sparse, tripeptide-based features often include rare motifs that can introduce significant sparsity into the dataset. LASSO's L_1 -norm regularization effectively mitigates this sparsity by shrinking the coefficients of low-impact features, such as rare motifs, to zero. The natural exclusion of rare motifs simplifies the feature set without the need for explicit preprocessing steps, ensuring that the selected features are both robust and interpretable. Furthermore in datasets with complex feature sets, collinearity among motifs can reduce model interpretability and introduce redundancy. Again, LASSO's L_1 -norm regularization effectively addresses this by selecting one representative feature from groups of highly correlated features while shrinking the coefficients of others to zero. This property ensures that redundant features are excluded, simplifying the final model and enhancing interpretability. Additionally, because *propy* features focus on dipeptides and exclude tripeptides, the potential for correlated motifs is inherently reduced in our analysis. The overview of our feature selection procedure is presented in Figure 1. As part of the LASSO feature selection process, we optimized the regularization parameter λ using cross-validation. For each fold, the model was trained on the training set and evaluated on the held-out test set using mean squared error (MSE). The average MSE across all folds was computed for each λ (Supplementary Figure S6).

3 Results

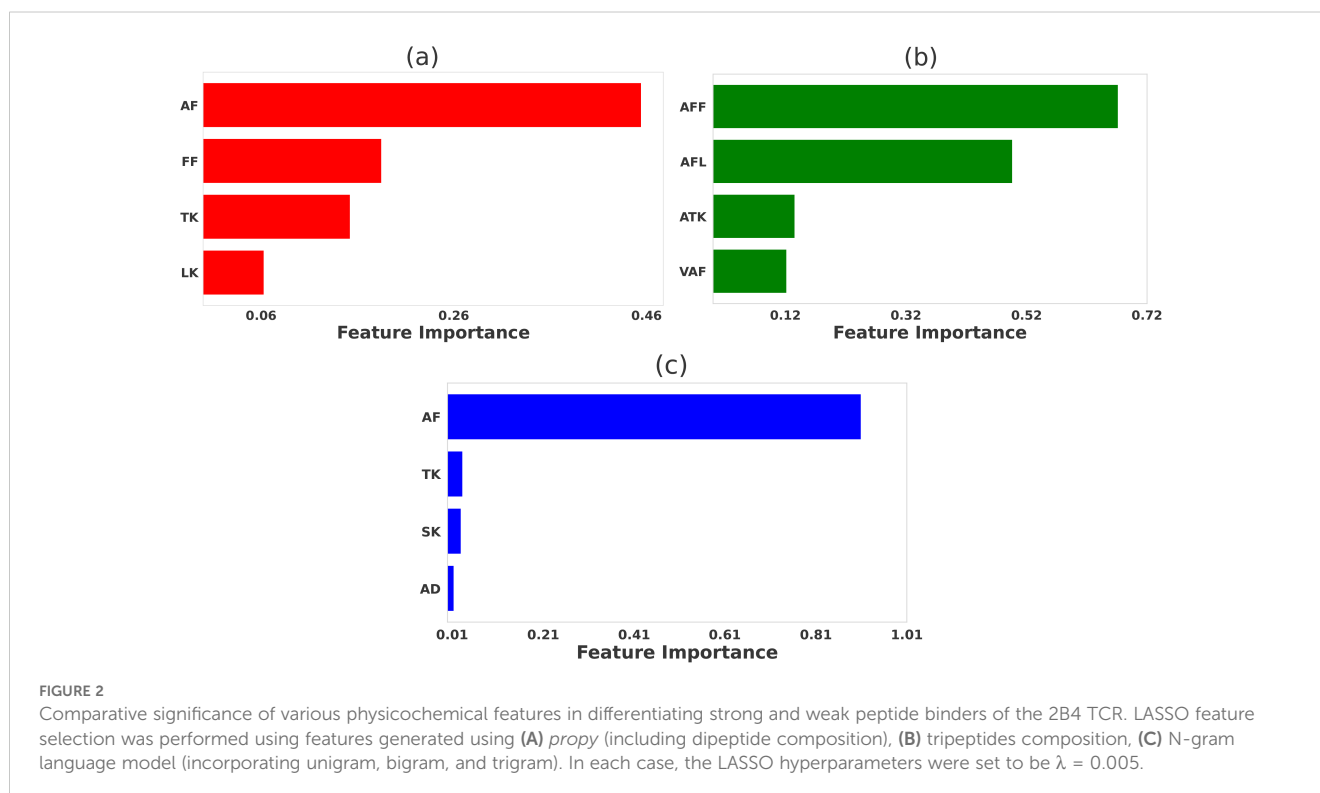
The relative significance of various physicochemical features distinguishing strong binders from weak binders in 2B4, 226, and

5cc7 peptide libraries are presented in Figures 2–4, respectively. We performed feature selection using three categories of properties: *propy* features, tripeptide composition, and N-gram language model (unigram, bigram, and trigram). This multi-faceted feature generation approach enabled us to extract key patterns and properties that significantly influence TCR binding behavior.

3.1 Features selection for 2B4 data

Our feature selection method for 2B4 data yields different but in many aspects overlapping selected features that contribute to strong binder peptides. Among features generated from the *propy* tool, the most important dipeptide compositions such as 'AF', 'FF', 'TK', and 'LK' likely represent amino acid pairs that significantly enhance peptide stability or affinity to the TCR [see Figure 2A]. Similar motifs are predicted when tripeptide compositions are utilized in the feature selection method, as shown in Figure 2B. Specifically, the tripeptide motif 'AFF' can be broken down into two dipeptides 'AF' and 'FF', both of which are captured by the *propy* method. The N-gram method also yields similar results, although selected features do not fully overlap with the results of other methods [see Figure 2B].

This observation is strongly supported by the experimental data, which highlights the amino acid preferences at key TCR contact positions (P3, P5, and P8) during peptide-MHC interactions (5). Notably, positions like P3 show a clear preference for aromatic residues such as phenylalanine (F) and tyrosine (Y), aligning with the dipeptides 'AF' and 'FF', and the tripeptide 'AFF', identified in our study. The overlap between motifs



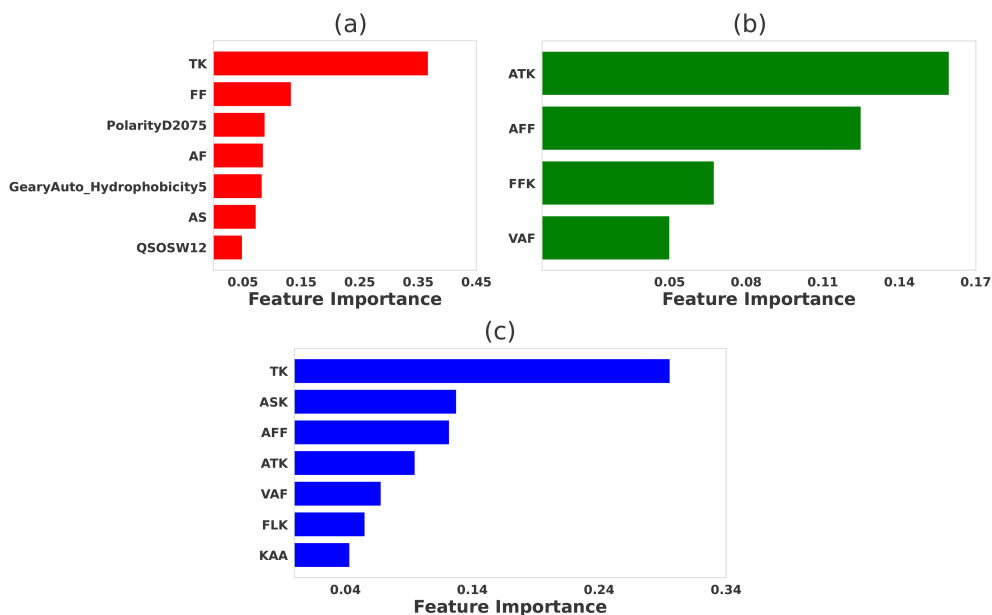


FIGURE 3 Comparative significance of various physicochemical features in differentiating strong binder and weak binder for peptides targeting 226 TCR. LASSO feature selection was performed using features generated using (A) propy (including dipeptide composition), (B) tripeptides composition, (C) N-gram language model (incorporating unigram, bigram, and trigram). For LASSO the hyperparameters were set to be $\lambda = 0.015$, $\lambda = 0.015$, and $\lambda = 0.01$, respectively.

such as dipeptides ‘AF’, ‘FF’ and tripeptides ‘AFF’, ‘AFL’ across independent feature selection methods could be a positive indicator of their relevance. Since feature selection for propy features and tripeptides was performed independently, the consistent identification of these motifs across models provides strong

evidence of their importance for distinguishing strong binders. This consistency is also reflected in the N-gram results, where motifs such as ‘AF’ dominate, highlighting the key structural patterns that underlie strong TCR-peptide interactions. To further validate the significance of motifs identified through feature

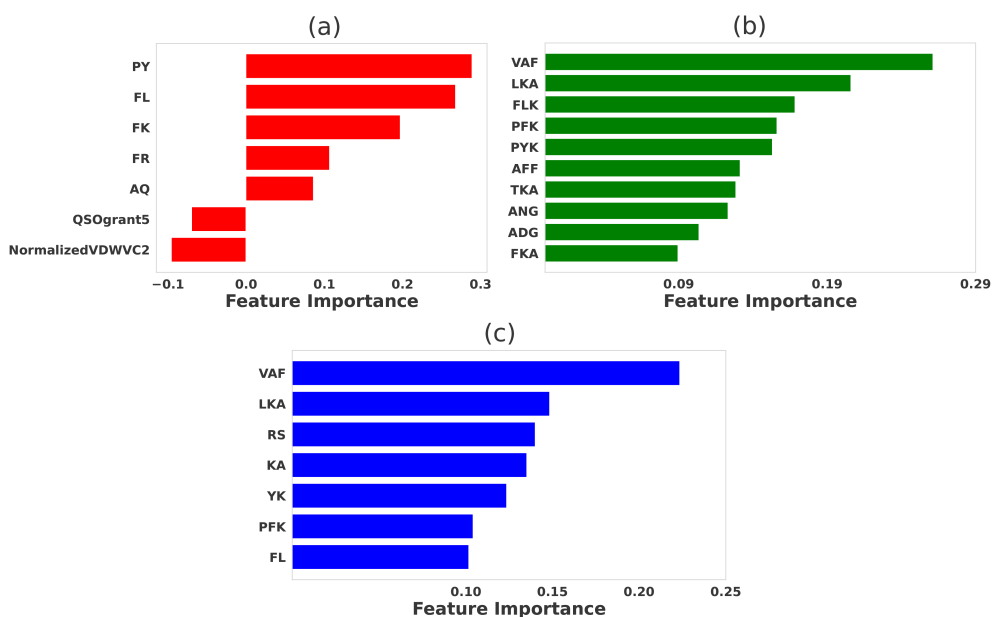


FIGURE 4 Comparative significance of various physicochemical features in differentiating strong and weak peptide binders of 5cc7 TCR. LASSO feature selection was performed using features generated using (A) propy (including dipeptide composition), (B) tripeptides composition, (C) N-gram language model (incorporating unigram, bigram, and trigram). In each case, the LASSO hyperparameters were set to be $\lambda = 0.01$ for all cases.

selection, we generated sequence logos for strong binders (Supplementary Figure S7), which provide a visual representation of motif enrichment by highlighting the frequency and conservation of amino acid patterns within the dataset.

3.2 Features selection for 226 data

For the 226 TCR dataset, our feature selection method also identified several key physicochemical features that distinguish strong binders from weak binders, as presented in Figure 3. Among the features generated by the *propy* tool [see Figure 3A], dipeptide compositions such as ‘TK’, ‘FF’, ‘AF’ emerged as the most important contributors to TCR-peptide affinity. These motifs likely play critical roles in stabilizing TCR-peptide interactions by complementing specific amino acid residues on the TCR. Structural analysis of the 226 TCR-pMHC (PDB ID:3QIU) reveals that the ‘TK’ motif contributes to electrostatic interactions, and ‘FF’ and ‘AF’ reinforce hydrophobic interactions that enhance the stability of the peptide-MHC complex. Definitions of other selected *propy* features, including PolarityD2075, GearyAuto Hydrophobicity5, and QSOSW12 are presented in Supplementary Table S1 in the Supplementary Information.

Specifically, in the ‘TK’ motif, the weakly acidic threonine residue (‘T’ at P8) can interact via hydrogen bonding with asparagine on the TCR’s CDR3 β loop. This interaction is identified through the contact map generated from the 3QIU crystal structure with a maximum distance ($r_{\max} = 8.5\text{\AA}$) (Supplementary Figures S5C, D). Moreover, for the ‘FF’ and ‘AF’ motifs, which do not exhibit these features in the original peptide, the RACER-derived pairwise amino acid energy matrix (Supplementary Figure S4B) predicts a favorable interaction between phenylalanine (‘F’) and alanine (‘A’). Similarly, alanine is predicted to favorably interact with proline (‘P’), methionine (‘M’), and phenylalanine (‘F’). These findings suggest that the ‘F’ residues in the peptide engage in favorable interactions with ‘A’ and other hydrophobic residues on the TCR. These hydrophobic interactions contribute to the stability of the TCR-peptide complex, underscoring the importance of the ‘FF’ and ‘AF’ motifs in facilitating binding through hydrophobic contacts.

Tripeptides such as ‘ATK’, ‘AFF’, and ‘FFK’ were identified as highly significant for distinguishing strong binders (see Figure 3B). These tripeptides can be deconstructed into dipeptides like ‘TK’, ‘FF’, and ‘AF’, which contain residues that are also prominent in the dipeptide analysis. The similarity between the tripeptide and dipeptide results reinforces the importance of these specific motifs, suggesting that key residues such as lysine (‘K’) and phenylalanine (‘F’) play central roles in facilitating hydrogen bonding and electrostatic interactions with the TCR’s complementarity-determining regions (CDRs). For example, lysine (‘K’) in the ‘ATK’ and ‘FFK’ motifs likely contributes to salt bridge formation, enhancing electrostatic interactions between the peptide and TCR. Similarly, phenylalanine (‘F’) in the ‘AFF’ and ‘FFK’ motifs strengthens binding through hydrophobic

interactions, which help stabilize the TCR-peptide complex within the MHC groove.

Furthermore, feature selection based on the N-gram language model [see Figure 3C] revealed a strong overlap with the amino acid patterns identified in both the dipeptide and tripeptide analyses. The most important features included ‘TK’, ‘ASK’, ‘AFF’, and ‘ATK’, which closely correspond to amino acids found at critical TCR contact points. This consistency across different feature selection methods reinforces the importance of these motifs in contributing to strong TCR-peptide interactions. The recurrence of ‘TK’, in particular, highlights the role of lysine in driving strong electrostatic interactions, while ‘AFF’ and ‘FFK’ emphasize the contribution of hydrophobic residues like phenylalanine (F) in maintaining binding affinity and structural stability.

3.3 Features selection for 5cc7 data

Our feature selection approach, trained using features generated from *propy* tool, predicts that ‘PY’, ‘FL’, ‘FK’, and ‘FR’, are crucial in determining peptide stability and TCR binding affinity [see Figure 4A]. These dipeptides are likely involved in stabilizing peptide-MHC interactions, with ‘FL’ and ‘FK’ contributing to hydrophobic and polar interactions, respectively. Hydrophobic residues such as phenylalanine (F) are known to form important nonpolar contacts that help to stabilize the peptide in the TCR binding groove, enhancing binding affinity. Definitions of the acronyms for other selected *propy* features — QSOgrant5 and NormalizedVDWVC2 — are presented in Supplementary Table S1 of the Supplementary Information.

When the tripeptide compositions are considered [see Figure 4B], motifs like ‘VAF’, ‘LKA’, and ‘FLK’ emerged as highly significant. These tripeptides suggest a combination of hydrophobic, polar, and charged interactions, which enhance the binding potential by promoting stable contacts of different natures with the TCR. For instance, ‘LKA’ features a combination of leucine (‘L’) and alanine (‘A’), hydrophobic residues, and lysine (‘K’), a positively charged residue, both of which are known to interact favorably with the TCR’s (PDB ID: 4P2R) binding pocket through hydrophobic and electrostatic interactions. Notably, the ‘LKA’ motif does not appear in the original peptide contact map generated with a maximum distance ($r_{\max} = 8.5\text{\AA}$) (Supplementary Figures S5E, F). However, analysis of the energy matrix (Supplementary Figure S4C) reveals that leucine (‘L’) has high affinity for proline (‘P’), alanine (‘A’), and phenylalanine (‘F’), suggesting potential hydrophobic interactions with these residues on the TCR. Similarly, lysine (‘K’) shows high affinity for tryptophan (‘W’), alanine (‘A’), phenylalanine (‘F’), and leucine (‘L’), indicating possible favorable interactions with these residues. Furthermore, alanine (‘A’), due to its small side chain, provides structural flexibility, allowing optimal positioning of neighboring residues for interaction. These findings suggest that the ‘LKA’ motif may enhance TCR-peptide interactions through hydrophobic and electrostatic interactions, as indicated by the energy matrix analysis, even though these interactions are not apparent in the contact map.

The N-gram language model (Figure 4C) further emphasized the importance of these motifs by identifying similar patterns. Features like 'VAF', 'LKA', 'RS', and 'KA' were among the most important for distinguishing strong binders, reflecting the same key interactions seen with the dipeptide and tripeptide compositions. The prevalence of hydrophobic residues such as valine (V), phenylalanine (F), and leucine (L) in these motifs highlight the critical role in stabilizing the peptide-MHC-TCR complex.

3.4 Prediction of strong vs weak binders using selected features

After selecting the essential physicochemical features from each peptide dataset, we aim to leverage these attributes to accurately and reliably predict strong and weak binders for each TCR type through logistic regression models. Table 2 summarizes key performance metrics, including Accuracy, Recall, F1 Score, Matthews Correlation Coefficient (MCC), and AUC (Area Under the ROC Curve), averaged over 10 cross-validation sets with an 80/20 train/test split for each fold. We employed StratifiedShuffleSplit for cross-validation to ensure that the equal class distribution achieved through undersampling was preserved in both training and test sets across all folds. By maintaining this balance, the model was evaluated on data representative of the balanced dataset used for training. While StratifiedShuffleSplit does not guarantee non-overlapping test sets across folds, each fold ensures that the test data is unseen during training for that specific fold.

For the 2B4 dataset, all three selected feature categories (*propy*, tripeptide composition, and N-gram) performed exceptionally well, with predictive accuracy ranging from 0.94 to 0.96 and an AUC reaching up to 0.98. These findings indicate that for this dataset, the selected features were highly informative, resulting in predictive models that perform well in identifying strong binders with high precision. The strong performance of the models for the 2B4 dataset can be attributed to the lower cross-reactivity of 2B4 i.e. it binds to a narrower range of peptides compared to more flexible TCRs,

making the binding interactions easier to model and predict. Furthermore, a smaller dataset (98 peptides) with clear sequence patterns provides the machine-learning models with less variability to account for, resulting in higher accuracy and AUC values.

In contrast, the 226 datasets demonstrated somewhat lower overall performance across all feature methods. Accuracy and AUC values were notably lower, with *propy* yielding the highest performance at 0.77 accuracy and 0.78 AUC, while the tripeptide and N-gram methods scored marginally lower. The relatively low MCC values (0.55 for *propy* and below 0.50 for others) suggest that the model's predictions are less consistent for this dataset. This result is likely due to the existence of more complex or less distinguishable features between strong and weak binders. The relatively lower performance of the models for the 226 dataset can be attributed to several factors related to the biological properties of the 226 TCR and the complexity of its dataset. The 226 TCR is known for its high degree of cross-reactivity (35), meaning it can recognize and bind to a much wider range of peptide sequences than more specific TCRs like 2B4. This broad recognition profile introduces greater variability in the peptide sequences classified as binders and non-binders, making it harder for machine learning models to identify clear patterns that distinguish strong from weak binders. Thus, the larger size of the 226 dataset, which includes 987 peptides, increases the diversity of peptide sequences.

For the 5cc7 dataset, however, performance is intermediate, with accuracy values ranging from 0.82 to 0.85 and an AUC as high as 0.87 for the tripeptide method. Here, the MCC values indicate that the models were relatively effective, with the N-gram method achieving the highest MCC (0.79), suggesting that it provided a more balanced prediction between strong and weak binders compared to the other methods. The F1 scores consistently reflect solid performance in identifying true strong binders, particularly with the tripeptide method ($F_1 = 0.86$). The moderate performance of the models for the 5cc7 dataset can be explained by the balance between specificity and cross-reactivity in the 5cc7 TCR and the size of the dataset. Unlike the highly specific 2B4 TCR or the highly

TABLE 2 Performance comparison of feature selection methods for three TCR datasets.

TCR Data	Feature Category	Accuracy	Recall	F1 Score	MCC	AUC
2B4	<i>propy</i>	0.94	0.94	0.93	0.87	0.96
	tripeptide	0.96	0.96	0.96	0.92	0.97
	N-gram	0.96	0.96	0.96	0.93	0.98
226	<i>propy</i>	0.77	0.77	0.79	0.55	0.78
	tripeptide	0.74	0.74	0.76	0.49	0.66
	N-gram	0.73	0.73	0.76	0.48	0.7
5cc7	<i>propy</i>	0.82	0.82	0.83	0.65	0.83
	tripeptide	0.85	0.85	0.86	0.71	0.87
	N-gram	0.85	0.85	0.85	0.79	0.84

Metrics include Accuracy, Recall, Matthews Correlation Coefficient (MCC), F1 Score, and AUC for trained baseline models (Logistic Regression) using selected features from LASSO. Values reflect the average across 10 cross-validation sets, with an 80/20 train/test split for each fold.

cross-reactive 226 TCR, 5cc7 exhibits an intermediate level of specificity. It binds to a moderately diverse set of peptides, leading to less sequence variability than 226 but more than 2B4.

It is important to highlight that while our datasets include an equal number of strong and weak binders, the overall peptide data are highly imbalanced in favor of weak binders over strong ones. Quantitatively, for a peptide of L residues, the total number of possible peptide sequences are 20^L , and an overwhelming majority of these sequences are weak binders. This discrepancy presents significant challenges in accurately predicting peptide specificity. However, despite these challenges, the close alignment between our $F1$ score and recall metrics (Table 2) indicates that the model achieves balanced performance in handling false positives (FP) and false negatives (FN). The balance between Recall and $F1$ score is especially critical in this context, where accurately identifying strong binders is essential, but misclassifying weak binders as strong could lead to a false sense of antigen coverage by a particular TCR, which can significantly affect T-cell based immunotherapeutic strategies. The fact that both metrics are comparable across datasets and feature selection methods indicates that the models are balanced in their sensitivity and specificity and are robustly selecting relevant features to resolve strong and weak binders (40).

3.5 Comparison with the RACER model

Our feature selection-based approach for predicting TCR-peptide binding, like several before it (31, 32, 34, 46–48), is purely sequence-based, relying on the selection of key features derived from amino acid sequences. By focusing on the sequence characteristics of peptides, we identified key dipeptide and tripeptide motifs that are enriched in strong binders. These selected motifs, without relying on detailed structural information, were essential for distinguishing between high- and low-affinity peptide antigens for specific TCR. In contrast, the RACER model adopts a hybrid approach by combining sequence data with structural insights to predict TCR-peptide binding affinities (22, 28, 29, 49, 50). RACER utilizes a pairwise energy framework, integrating residue-specific energy matrices derived from high-throughput data on experimentally confirmed TCR-peptide interactions, along with crystal structures of these complexes. The structural templates provided by crystal data allow RACER to quantify the binding energy with greater precision by modeling the physical interactions at atomic resolution. After identifying motifs enriched in strong binders, we then aimed to apply RACER's pairwise energy framework to test our sequence-based approach. This allowed us to pinpoint the specific positions within the peptide sequence where these motifs are predicted to have the most significant impact on binding energy.

To determine the positions of the selected features, we performed *in silico* mutation in all two-adjacent (selected dipeptide motifs) and three-adjacent (selected tripeptide motifs) amino acids at every peptide amino acid position containing the selected features. We then used RACER to estimate the binding energy for each mutant

peptide. The binding energies for all possible mutant peptides are plotted for selected dipeptides (Supplementary Figure S2) and tripeptides (Supplementary Figure S3). We then compared the binding energy of each oligopeptide motif located at each position to the baseline binding energy of the wild-type (WT) TCR-peptide (WT given by the black dashed line in Supplementary Figures S2, S3). If a mutated TCR-peptide showed increased binding energy (above the dashed line), it indicates that the mutation enhanced the binding affinity above that of the WT (strong) binding peptide, suggesting that the underlying importance of the selected dipeptide at that specific position. On the other hand, mutation may also result in significantly lower predicted affinities. This decrease indicates that certain substitutions disrupt key interactions necessary for strong binding, effectively identifying sequences that function as weak binders. By recognizing these sequences, we not only validate the specificity of our selected motifs but also enhance our understanding of the structural and sequence determinants that diminish binding affinity. This dual identification of both strong and weak binders underscores the robustness of our approach in mapping the landscape of TCR-peptide interactions.

For example, Supplementary Figure S2A shows that for 2B4, dipeptide 'AF' at positions (2, 3), (3, 4), and (5, 6) increased binding affinity, with a particularly significant increase at positions (5, 6). Although all three TCRs retain a WT-like TCR recognition motif, each TCR exhibits some variation in positional preference (Figure 5). For instance, whereas 2B4 can recognize Lysine at position P8 [P5 in (5)], 5cc7 accommodates Leucine and Arginine at P8. The P6 [P3 in (5)] TCR contact position showed the least variance across all three TCRs, with either Phenylalanine or Valine being required for 2B4 and 5cc7, and Phenylalanine, Lysine, or Arginine being required for 226. As previously reported (5), 226 demonstrated a greater degree of cross-reactivity, being able to recognize 897 unique peptide sequences. The larger number of peptides recognized was largely due to a higher tolerance for substitutions at TCR-neutral and MHC-contacting residues, such as position P9 (Figure 5B).

Combining the predictions identified in our feature selection framework with RACER-predicted position-specific information provides an opportunity to construct heatmaps (Figure 5) enriched in beneficial dipeptide compositions that maximally resolve strong and weak binding peptides. These results can be directly compared to those from the original work by Birnbaum et al. (5), which provided a similar description of binding motifs acquired experimentally. Notably, while their work identified *single amino acid* hotspots indicative of strong binders by analyzing the abundance of amino acids in strong-binding peptides, our approach focuses on *dipeptide motifs*, identifying them based on their enrichment in strong binders relative to non-binders.

This methodological difference is evident when considering anchor residues like P4 and P12, which are restricted to isoleucine, leucine, valine, and lysine, respectively. While lysine is ubiquitous among all strong binders, it is similarly present in weak binders and thus does not emerge as a distinguishing amino acid at P12 in our analysis (Figure 5). This underscores one way in which conserved residues might mask a model's discriminatory power

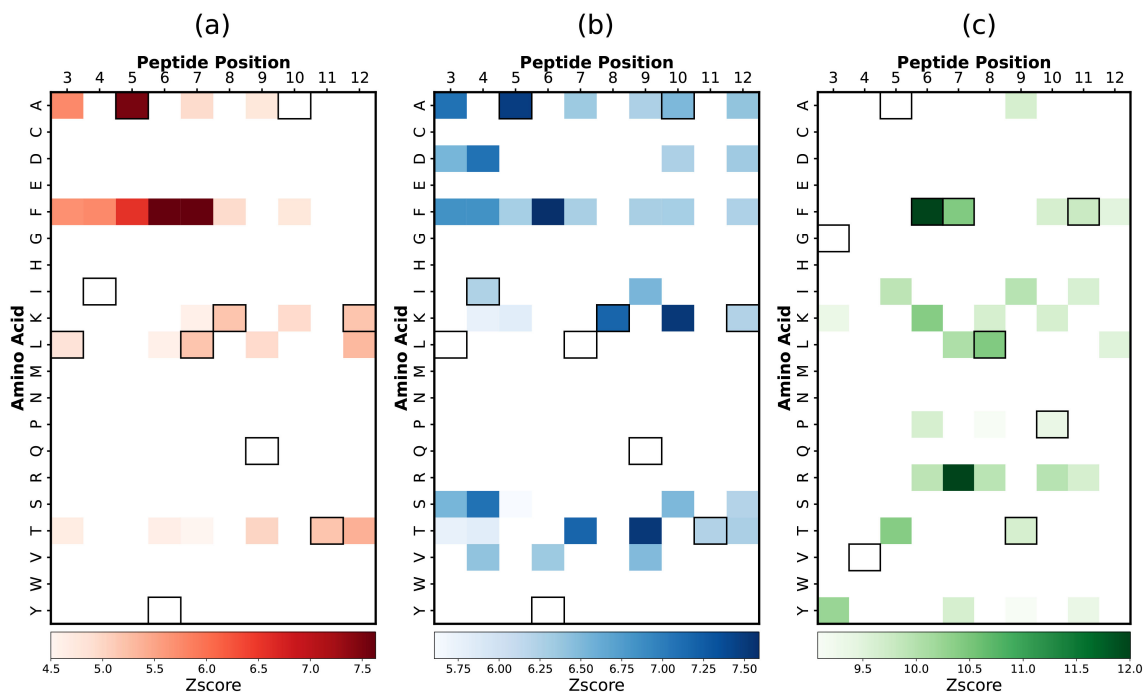


FIGURE 5

Predicted heatmaps using RACER model for (A) 2B4, (B) 226, and (C) 5cc7 peptide libraries. The sequence for the peptide is represented via outlined boxes.

were it to identify those features as important for strong binding. In addition to this, in the experimental data from (5), heatmaps are generated after three rounds of selection, whereas in our approach, we generate heatmaps after five rounds of selection.

Despite these challenges, the key motifs identified through our sequence-based feature selection are corroborated by RACER's binding energy predictions. Particularly in the case of the 2B4 TCR, where we predict the enrichment of specific motifs such as the 'AFF' motif. RACER's energy calculations confirm that these motifs contribute significantly to binding affinity (as indicated by larger interaction values with complementary amino acids 'WSQ' in the 2B4 TCR CDR3 β domain, and 'RA' and 'G' in the CDR3 α domain), aligning with experimental observations. The 2B4 TCR seems well-characterized by a single important motif, as evidenced by the peptide position curves showing a single sharp peak corresponding to a small number of features (Supplementary Figures S2A, S3A). In contrast, the 226 and 5cc7 TCRs display different binding characteristics, highlighting the unique specificity of each TCR. For 5cc7, we observe lower intensity compared to 2B4, with a wide and smooth peak across a larger number of features (Supplementary Figures S2C, S3C). For 226, we observe sharp peaks at different positions for a large number of features, which aligns well with the high cross-reactivity we previously mentioned for 226 (Supplementary Figures S2B, S3B). In both 2B4 and 226, we observe the importance of features containing phenylalanine ('F') in the first part of the peptide across positions P3 to P7. This is because the peaks for dipeptide and tripeptide motifs containing 'F'

are prominent in these positions. When considering dipeptide and tripeptide motifs, 'F' appears in all positions P3–P7, whereas in single amino acid analysis, 'F' does not appear in all positions since only the best location is selected. This indicates that 'F' can be a very important feature when combined with other amino acids in motifs. Together, these findings underscore the power of our model in identifying critical dipeptide and tripeptide motifs, which are more informative than single amino acid motifs, thereby enhancing predictive performance and providing deeper insights into TCR-peptide interactions.

4 Discussion

The interaction between T-cell receptors (TCRs) and peptide-MHC complexes is a critical component of the adaptive immune system, enabling T-cells to detect and respond to specific antigens. This process, however, is complicated by the TCR cross-reactivity, where a single TCR can recognize multiple peptide sequences. Understanding cross-reactivity is important as many TCRs are known to confer coverage across distinct peptide systems. Cross-reactivity also presents a major challenge for reliably predicting TCR specificity, which is required for optimal vaccine design and selection of T-cell-based immunotherapy. In this study, we leveraged machine learning techniques with refined feature selection to improve the accuracy and generalizability of TCR-peptide interaction predictions. Our findings show that focusing on

specific physicochemical features significantly enhances the model's ability to distinguish between strong and weak binders, offering new insights into the molecular mechanism of TCR recognition.

The number of peptides in our final dataset reflected the varying specificity and cross-reactivity of the TCRs, which in turn explained the differences in model performance. The 2B4 dataset, with only 98 peptides, highlighted the high specificity of the 2B4 TCR, which lead to clearer binding patterns and superior model performance. In contrast, the 226 dataset, which includes 987 peptides, highlighted the TCR's greater cross-reactivity, thereby making binding patterns more complex and harder to capture with our feature selection methods, which resulted in lower performance metrics. The 5cc7 dataset, with 234 peptides, demonstrated intermediate specificity and moderate cross-reactivity, aligning with its intermediate model performance. This variation in dataset sizes reflects the inherent biological properties of each TCR, with more specific TCRs resulting in smaller datasets and higher model performance. Undersampling is a common technique when dealing with class imbalance in machine learning, and it can potentially impact both feature selection and model performance. In our study, we resorted to undersampling to balance the number of weak binders (class 0) with the number strong binders (class 1). This may have resulted in the loss of valuable information or features associated with weak binders. However, as noted earlier, the closeness of F1 score and recall demonstrated that undersampling did not affect the balance between false positive (FP) and false negatives (FN). These metrics, along with additional performance measures such as Accuracy, MCC, and AUC, confirmed the robustness and reliability of our feature selection approach despite the use of undersampling.

We employed the LASSO feature selection method to extract meaningful physicochemical properties from the peptide sequences to identify key features that contribute to binding affinity, including amino acid composition, dipeptide frequency, and tripeptide motifs. Among these, dipeptide compositions and tripeptide compositions emerged as particularly important, consistently ranking among the most predictive for distinguishing strong from weak binders across the different TCR datasets. This finding suggests that the arrangement of amino acids in short peptide sequences plays a crucial role in TCR recognition, and our optimized feature set provides a robust foundation for understanding TCR-peptide interactions and highlighting the importance of tripeptide features.

Identified important tripeptides in TCR-peptide binding can be further understood by examining the molecular interactions between peptide residues and the CDR3 α and CDR3 β loops of the TCR, as illustrated in [Figure 1](#). It is known that a single amino acid in the peptide can simultaneously interact with residues from both the CDR3 α and CDR3 β regions (35). For example, if we consider a symbolic tripeptide sequence like 'LTP' the first residue, 'L', may form contacts with both CDR3 α and CDR3 β , providing a dual interaction site. In contrast, the second and third residues, 'T' and 'P', may predominantly interact with only CDR3 β . This picture highlights how specific residues within a tripeptide can influence the binding strength by creating multiple interaction points, making tripeptides like 'LTP' particularly important for determining binding affinity. The ability of certain tripeptides to establish

multiple points of contact contributes to the overall specificity and affinity of TCR recognition.

MHC molecules, are known to play a critical role in immune recognition by influencing peptide presentation to T-cell receptors (51). In this study, however, the datasets utilized specific MHC alleles [e.g., I-Ek in murine TCR experiments, as described in (5)], effectively eliminating the influence of MHC polymorphism. This approach allowed us to isolate the effects of peptide features on TCR binding. Future studies could expand on this work by incorporating data from multiple MHC alleles to explore how MHC polymorphism shapes TCR-peptide interactions and the generalizability of predictive models.

An improved quantitative understanding of the features that derive TCR specificity remains a significant obstacle in the fields of immunotherapy and vaccine design (9, 52). By identifying key dipeptide and tripeptide motifs predictive of TCR-peptide binding, our results contributes to a deeper understanding of the sequence-level determinants of TCR specificity. These insights could be leveraged to design peptides with optimized binding profiles, enhancing immune responses in therapeutic contexts. For example, the ability to predict TCR-peptide interactions could aid in developing personalized TCR therapies, where T-cells are designed (e.g., CAR T-cells) or identified (e.g., adoptive cell therapy) to recognize tumor-associated antigens. Similarly, in vaccine design, these motifs could help identify or engineer peptides that elicit strong and targeted immune responses, improving the efficacy of peptide-based vaccines.

While our sequence-based approach successfully identifies key dipeptide and tripeptide motifs enriched in strong binders, it has certain limitations. Our purely sequence-driven model may miss rare or unconventional motifs and struggle in cases of extreme TCR cross-reactivity. Additionally, our findings are derived from a relatively limited set of TCR-peptide interactions, which may limit the generalizability of the identified motifs across all TCRs, particularly those with unique binding preferences. Moreover, certain TCRs may prioritize interactions with MHC residues over peptides, a factor that our current model does not fully address. To overcome these limitations, future work will explore hybrid models that integrate structural insights, allowing for more accurate predictions of TCR-peptide dynamics. Despite these challenges, this approach is able to extract meaningful motifs for resolving TCR specificity based on TCR and peptide primary sequences. Future work will be directed at using these learned features to train a classification model for identifying strong binding pairs from a variety of possible TCR and peptide test sequences. Moreover, advanced deep learning approaches, such as attention mechanisms and transformer architectures, can further investigate TCR-peptide binding specificity.

Data availability statement

Publicly available datasets were analyzed in this study. This data can be found here: https://github.com/XingchengLin/RACER/tree/main/raw_data, https://github.com/TAMUGeorgeGroup/Feature_Selection_TCR-Specific_Interaction.git.

Author contributions

HT: Writing – original draft, Writing – review & editing, Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Software, Supervision, Validation, Visualization. ZG: Writing – original draft, Writing – review & editing, Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Software, Supervision, Validation, Visualization. AK: Writing – original draft, Writing – review & editing, Conceptualization, Formal analysis, Funding acquisition, Investigation, Methodology, Project administration, Resources, Supervision. JG: Writing – original draft, Writing – review & editing, Conceptualization, Formal analysis, Funding acquisition, Investigation, Methodology, Project administration, Resources, Supervision.

Funding

The author(s) declare financial support was received for the research, authorship, and/or publication of this article. AK was supported by the Welch Foundation (C-1559), and by the Center for Theoretical Biological Physics sponsored by the NSF (PHY-2019745). JG was supported by the Cancer Prevention Research Institute of Texas (RR210080) and the National Institute of General Medical Sciences of the NIH (R35GM155458). JG is a CPRIT Scholar in Cancer Research.

References

- Ca J. The immune system in health and disease. (New York: Garland Pub.) (2001). Available online at: <http://www.garlandscience.com>.
- Birnbaum ME, Dong S, Garcia KC. Diversity-oriented approaches for interrogating t-cell receptor repertoire, ligand recognition, and function. *Immunol Rev.* (2012) 250:82–101. doi: 10.1111/imr.2012.250.issue-1
- Rudolph MG, Stanfield RL, Wilson IA. How tcrs bind mhc, peptides, and coreceptors. *Annu Rev Immunol.* (2006) 24:419–66. doi: 10.1146/annurev.immunol.23.021704.115658
- Garcia KC, Adams EJ. How the t cell receptor sees antigen—a structural view. *Cell.* (2005) 122:333–6. doi: 10.1016/j.cell.2005.07.015
- Birnbaum ME, Mendoza JL, Sethi DK, Dong S, Glanville J, Dobbins J, et al. Deconstructing the peptide-mhc specificity of t cell recognition. *Cell.* (2014) 157:1073–87. doi: 10.1016/j.cell.2014.03.047
- Morse MA, Gwin WR III, Mitchell DA. Vaccine therapies for cancer: then and now. *Target Oncol.* (2021) 16:121–52. doi: 10.1007/s11523-020-00788-w
- Yee C. Adoptive t cell therapy: addressing challenges in cancer immunotherapy. *J Trans Med.* (2005) 3:1–8. doi: 10.1186/1479-5876-3-17
- Harkioliaki M, Holmes SL, Svendsen P, Gregersen JW, Jensen LT, McMahon R, et al. T cell-mediated autoimmune disease due to low-affinity crossreactivity to common microbial peptides. *Immunity.* (2009) 30:348–57. doi: 10.1016/j.immuni.2009.01.009
- Buhrman JD, Jordan KR, Munson DJ, Moore BL, Kappler JW, Slansky JE. Improving antigenic peptide vaccines for cancer immunotherapy using a dominant tumor-specific t cell receptor. *J Biol Chem.* (2013) 288:33213–25. doi: 10.1074/jbc.M113.509554
- Løset GÅ, Berntzen G, Frigstad T, Pollmann S, Gunnarsen KS, Sandlie I. Phage display engineered t cell receptors as tools for the study of tumor peptide-mhc interactions. *Front Oncol.* (2015) 4:378. doi: 10.3389/fonc.2014.00378
- Gee MH, Han A, Lofgren SM, Beausang JF, Mendoza JL, Birnbaum ME, et al. Antigen identification for orphan t cell receptors expressed on tumor-infiltrating lymphocytes. *Cell.* (2018) 172:549–63. doi: 10.1016/j.cell.2017.11.043
- Springer I, Besser H, Tickotsky-Moskovitz N, Dvorkin S, Louzoun Y. Prediction of specific tcr-peptide binding from large dictionaries of tcr-peptide pairs. *Front Immunol.* (2020) 11:1803. doi: 10.3389/fimmu.2020.01803

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fimmu.2024.1510435/full#supplementary-material>

- Garboczi DN, Ghosh P, Uetz U, Fan QR, Biddison WE, Wiley DC. Structure of the complex between human t-cell receptor, viral peptide and hla-a2. *Nature.* (1996) 384:134–41. doi: 10.1038/384134a0
- Margulies DH, Plaksin D, Khilko S, Jelonek MT. Studying interactions involving the t-cell antigen receptor by surface plasmon resonance. *Curr Opin Immunol.* (1996) 8:262–70. doi: 10.1016/S0952-7915(96)80066-5
- Boder ET, Bill JR, Niels AW, Marrack PC, Kappler JW. Yeast surface display of a noncovalent mhc class ii heterodimer complexed with antigenic peptide. *Biotechnol Bioeng.* (2005) 92:485–91. doi: 10.1002/(ISSN)1097-0290
- Jiang W, Boder ET. High-throughput engineering and analysis of peptide binding to class ii mhc. *Proc Natl Acad Sci.* (2010) 107:13258–63. doi: 10.1073/pnas.1006344107
- Starwalt SE, Masteller EL, Bluestone JA, Kranz DM. Directed evolution of a single-chain class ii mhc product by yeast display. *Protein Eng.* (2003) 16:147–56. doi: 10.1093/proeng/gzg018
- Purcell AW, Ramarathinam SH, Ternette N. Mass spectrometry-based identification of mhc-bound peptides for immunopeptidomics. *Nat Protoc.* (2019) 14:1687–707. doi: 10.1038/s41596-019-0133-y
- Peng X, Lei Y, Feng P, Jia L, Ma J, Zhao D, et al. Characterizing the interaction conformation between t-cell receptors and epitopes with deep learning. *Nat Mach Intell.* (2023) 5:1–13. doi: 10.1038/s42256-023-00634-4
- Hudson D, Fernandes RA, Basham M, Ogg G, Koohy H. Can we predict t cell specificity with digital biology and machine learning? *Nat Rev Immunol.* (2023) 23:511–21. doi: 10.1038/s41577-023-00835-3
- Lee CH, Huh J, Buckley PR, Jang M, Pinho MP, Fernandes RA, et al. A robust deep learning workflow to predict cd8+ t-cell epitopes. *Genome Med.* (2023) 15:70. doi: 10.1186/s13073-023-01225-z
- Ghoreyshi ZS, George JT. Quantitative approaches for decoding the specificity of the human t cell repertoire. *Front Immunol.* (2023) 14:1228873. doi: 10.3389/fimmu.2023.1228873
- Zhang W, Hawkins PG, He J, Gupta NT, Liu J, Choonoo G, et al. A framework for highly multiplexed dextramer mapping and prediction of t cell receptor sequences to antigen specificity. *Sci Adv.* (2021) 7:eabf5835. doi: 10.1126/sciadv.abf5835

24. Jurtz VI, Jessen LE, Bentzen AK, Jespersen MC, Mahajan S, Vita R, et al. Netter: sequence-based prediction of tcr binding to peptide-mhc complexes using convolutional neural networks. *BioRxiv*. (2018), 433706. doi: 10.1101/433706
25. Lee CH, Salio M, Napolitani G, Ogg G, Simmons A, Koohy H. Predicting cross-reactivity and antigen specificity of t cell receptors. *Front Immunol*. (2020) 11:565096. doi: 10.3389/fimmu.2020.565096
26. Petrova G, Ferrante A, Gorski J. Cross-reactivity of t cells and its role in the immune system. *Crit ReviewsTM Immunol*. (2012) 32:349–72. doi: 10.1615/CritRevImmunol.v32.i4.50
27. Antunes DA, Rigo MM, Freitas MV, Mendes MF, Sinigaglia M, Lizée G, et al. Interpreting t-cell cross-reactivity through structure: implications for tcr-based cancer immunotherapy. *Front Immunol*. (2017) 8:1210. doi: 10.3389/fimmu.2017.01210
28. Lin X, George JT, Schafer NP, Ng Chau K, Birnbaum ME, Clementi C, et al. Rapid assessment of t-cell receptor specificity of the immune repertoire. *Nat Comput Sci*. (2021) 1:362–73. doi: 10.1038/s43588-021-00076-1
29. Wang A, Lin X, Chau KN, Onuchic JN, Levine H, George JT. Racer-m leverages structural features for sparse t cell specificity prediction. *Sci Adv*. (2024) 10:ead0161. doi: 10.1126/sciadv.adl0161
30. Ghoreyshi ZS, Teimouri H, Kolomeisky AB, George JT. Integration of kinetic data into affinity-based models for improved t cell specificity prediction. *Biophys J*. (2024) 123(23):4115–22. doi: 10.1101/2024.06.17.599469
31. Meynard-Piganeau B, Feinauer C, Weigt M, Walczak AM, Mora T. Tulip: A transformer-based unsupervised language model for interacting peptides and t cell receptors that generalizes to unseen epitopes. *Proc Natl Acad Sci*. (2024) 121:e2316401121. doi: 10.1073/pnas.2316401121
32. Dash P, Fiore-Gartland AJ, Hertz T, Wang GC, Sharma S, Souquette A, et al. Quantifiable predictive features define epitope-specific t cell receptor repertoires. *Nature*. (2017) 547:89–93. doi: 10.1038/nature22383
33. Glanville J, Huang H, Nau A, Hatton O, Wagar LE, Rubelt F, et al. Identifying specificity groups in the t cell receptor repertoire. *Nature*. (2017) 547:94–8. doi: 10.1038/nature22976
34. Montemurro A, Schuster V, Povlsen HR, Bentzen AK, Jurtz V, Chronister WD, et al. Netcr-2.0 enables accurate prediction of tcr-peptide binding by using paired tcr α and β sequence data. *Commun Biol*. (2021) 4:1060. doi: 10.1038/s42003-021-02610-3
35. Newell EW, Ely LK, Kruse AC, Reay PA, Rodriguez SN, Lin AE, et al. Structural basis of specificity and cross-reactivity in t cell receptors specific for cytochrome c-i-ek. *J Immunol*. (2011) 186:5823–32. doi: 10.4049/jimmunol.1100197
36. He H, Garcia EA. Learning from imbalanced data. *IEEE Trans Knowledge Data Eng*. (2009) 21:1263–84. doi: 10.1109/TKDE.2008.239
37. Cao DS, Xu QS, Liang YZ. propy: a tool to generate various modes of chou's pseaac. *Bioinformatics*. (2013) 29:960–2. doi: 10.1093/bioinformatics/btt072
38. Teimouri H, Medvedeva A, Kolomeisky AB. Bacteria-specific feature selection for enhanced antimicrobial peptide activity predictions using machine-learning methods. *J Chem Inf Model*. (2023) 63:1723–33. doi: 10.1021/acs.jcim.2c01551
39. Lee EY, Lee MW, Fulan BM, Ferguson AL, Wong GC. What can machine learning do for antimicrobial peptides, and what can antimicrobial peptides do for machine learning? *Interface Focus*. (2017) 7:20160153. doi: 10.1098/rsfs.2016.0153
40. Teimouri H, Medvedeva A, Kolomeisky AB. Unraveling the role of physicochemical differences in predicting protein-protein interactions. *J Chem Phys*. (2024) 161:045102-1-45102-11. doi: 10.1063/5.0219501
41. Cavnar WB, Trenkle JM, et al. (1994). in: *Proceedings of SDAIR-94, 3rd annual symposium on document analysis and information retrieval*, Ann Arbor, Michigan, Vol. 161175. p. 14.
42. Zeng Z, Shi H, Wu Y, Hong Z. Survey of natural language processing techniques in bioinformatics. *Comput Math Methods Med*. (2015) 2015:674296. doi: 10.1155/2015/674296
43. Ding H, Lin H, Chen W, Li ZQ, Guo FB, Huang J, et al. Prediction of protein structural classes based on feature selection technique. *Interdiscip Sci: Comput Life Sci*. (2014) 6:235–40. doi: 10.1007/s12539-013-0205-6
44. Islam SA, Heil BJ, Kearney CM, Baker EJ. Protein classification using modified n-grams and skip-grams. *Bioinformatics*. (2018) 34:1481–7. doi: 10.1093/bioinformatics/btx823
45. Tibshirani R. Regression shrinkage and selection via the lasso. *J R Stat Soc Ser B: Stat Method*. (1996) 58:267–88. doi: 10.1111/j.2517-6161.1996.tb02080.x
46. Detours V, Mehr R, Perelson AS. A quantitative theory of affinity-driven t cell repertoire selection. *J Theor Biol*. (1999) 200:389–403. doi: 10.1006/jtbi.1999.1003
47. Košmrlj A, Jha AK, Huseby ES, Kardar M, Chakraborty AK. How the thymus designs antigen-specific and self-tolerant t cell receptor sequences. *Proc Natl Acad Sci*. (2008) 105:16671–6. doi: 10.1073/pnas.08080811105
48. George JT, Kessler DA, Levine H. Effects of thymic selection on t cell recognition of foreign and tumor antigenic peptides. *Proc Natl Acad Sci*. (2017) 114:E7875–81. doi: 10.1073/pnas.1708573114
49. Ng Chau K, George JT, Onuchic JN, Lin X, Levine H. Contact map dependence of a t-cell receptor binding repertoire. *Phys Rev E*. (2022) 106:014406. doi: 10.1103/PhysRevE.106.014406
50. Bradley P. Structure-based prediction of t cell receptor: peptide-mhc interactions. *eLife*. (2023) 12:e82813. doi: 10.7554/eLife.82813
51. Radwan J, Babik W, Kaufman J, Lenz TL, Winternitz J. Advances in the evolutionary understanding of mhc polymorphism. *Trends Genet*. (2020) 36:298–311. doi: 10.1016/j.tig.2020.01.008
52. McMahan RH, McWilliams JA, Jordan KR, Dow SW, Wilson DB, Slansky JE, et al. Relating tcr-peptide-mhc affinity to immunogenicity for the design of tumor vaccines. *J Clin Invest*. (2006) 116:2543–51. doi: 10.1172/JCI26936