

Unraveling T Cell Receptor Specificity: An Integrated Approach to Sequence Analysis

Shalomi Fernandes

Department of Engineering and Mathematics
University of Bristol
Bristol, United Kingdom
mh23950@bristol.ac.uk

Jiadong Xu

Department of Engineering and Mathematics
University of Bristol
Bristol, United Kingdom
jjadong.xu.2021@bristol.ac.uk

Fangnan Wei

Department of Engineering and Mathematics
University of Bristol
Bristol, United Kingdom
tj23631@bristol.ac.uk

Jiahui Liu

Department of Engineering and Mathematics
University of Bristol
Bristol, United Kingdom
ye23356@bristol.ac.uk

Abstract—This study utilizes machine learning techniques to predict the specificity of T cell receptors (TCRs), which play a key role in the immune system’s response to pathogens and cancer cells. Our research focuses on the broad diversity of TCR sequences in the VDJdb dataset to understand their interactions with specific epitopes. We have developed a method that includes the preprocessing of TCR sequence data, we have used techniques such as Encoding, computational distance matrix, classification, dimensionality reduction, clustering techniques to process the VDJdb dataset, and we have also built an algorithm aimed at predicting antigen specificity. Our innovative computational approach is expected to overcome the limitations of current empirical approaches and significantly improve the effectiveness of immunotherapies and personalized medicine. This article presents the results of our research on the predictive accuracy of these methods and discusses potential improvements and future research directions.

I. INTRODUCTION

This study aims to predict the specificity of T cell receptors (TCRs) by machine learning techniques. TCRs are a key component of the immune system, which is able to recognize and respond to antigens presented by pathogens or cancerous cells through a highly variable molecular structure. This diversity is primarily generated through the recombination of variable (V), diversity (D), and junction (J) gene fragments, allowing TCRs to bind to a variety of peptides presented by major histocompatibility complex (MHC) molecules. However, due to the complexity and high variability of peptide-MHC-TCR interactions, the use of sequence data to predict the specificity of TCRs remains a challenge. Most of the existing traditional methods rely on empirical data, which often cannot accurately capture the nuances of antigen recognition. Therefore, this research aims to develop an innovative computational method that can more accurately simulate and predict TCR behavior through in-depth analysis of the vast TCR sequence data in the VDJdb database, which is expected to drive the development of personalized medicine and immunotherapy.

Predicting TCR specificity is critical for advancing immunotherapies and designing targeted therapies. The project involves preprocessing TCR sequence data, encoding, calculating distances, classification, clustering, dimensionality reduction techniques, and developing algorithms for predicting antigen specificity. We focus on extracting meaningful patterns from TCR sequences using dimensionality reduction and clustering techniques, with the aim of developing robust predictive models. Our approach aims not only to enhance the understanding of TCR-antigen interactions, but also to advance personalized medicine by facilitating the development of more effective immunotherapies. By combining computational methods with immunological insights, our research explores new areas of immunoinformatics, paving the way for breakthroughs in the understanding and therapeutic utilization of the immune system. This research has great potential to enhance the efficacy of personalized medicine and immunotherapy, and to address the limitations of current empirical approaches through innovative computational methods.

II. LITERATURE REVIEW

“GIANA allows computationally-efficient TCR clustering and multi-disease repertoire classification by isometric transformation” [1] introduces a TCR alignment algorithm GIANA (Geometric Isometry-based TCR AlignNment Algorithm) based on geometric isometric transformation, a new computational tool designed to improve the efficiency and accuracy of T cell receptors. GIANA not only improves calculation speed (600 times faster than TCRdist without sacrificing accuracy), but also facilitates fast queries of large reference queues. It successfully identified novel disease-associated TCRs and classified unseen samples in a variety of diseases, including cancer, infectious diseases, and autoimmune diseases.

"T cell receptor sequence clustering and antigen specificity" [2] describes clustering techniques for classifying TCR sequences based on their biological similarity and antigen specificity. Technological advances have made TCR sequencing data widely available, and there is increasing interest in understanding TCR-epitope interactions to predict disease outcome, track treatment efficacy, and stratify patients for treatment. This review discusses several sequencing-based methods, including sequence alignment, analysis of short TCR motifs, and the use of various computational tools to predict TCR binding based on structural and sequence data, highlighting the challenges of these methods, such as the need for large data set and the complexity of predicting protein interactions based on sequence data alone. The paper also explores future directions for improving TCR sequence analysis, showing that integrating more sophisticated machine learning models and structural data can improve the predictive accuracy of these methods.

"TCR meta-clonotypes for biomarker discovery with tcrdist3 enabled identification of public, HLA-restricted clusters of SARS-CoV-2 TCRs" [3] introduces a new framework to classify biochemically similar TCRs into "meta-clonotypes". Researchers used the newly developed open source software package tcrdist3 to analyze TCR data from COVID-19 patients to create and quantify meta-clonotypes, providing a more powerful tool for identifying disease-associated TCR patterns. Application of this method identified 1,831 public TCR metaclonotypes associated with SARS-CoV-2, and significant HLA restriction was observed. Meta-clonotypes were more frequently detected in the TCR repertoires of patients with specific HLA genotypes than exact amino acid matches.

III. METHODOLOGY

In this study, the data is preprocessed, encoded, distance calculation of TCR sequences is done followed by dimensionality reduction, clustering, classification, and analysis of predictions. The encoding methods used in this article is One-Hot, BLOSUM 62 and GIANA Encoding. The distance calculation methods used in this article are TCRDist, GIANA and Levenshtein. The dimensionality reduction methods used in this article are PCA, t-SNE and UMAP. The clustering methods used in this article are Agglomerative Hierarchical clustering and DBSCAN. The classification method used in this article are Logistic Regression, SVM, and RandomForestClassifier.

A. Pre-Processing

In this study, the data preparation process involved several key steps to ensure the cleanliness, relevance, and structured organization of the dataset. This included loading the TCR sequence data, followed by comprehensive data cleaning procedures to handle missing values and remove duplicates. Relevant features were carefully selected based on their significance to the analysis goals, while new features were

engineered to enhance the model's ability to discern meaningful patterns. Additionally, normalization or standardization techniques were applied to bring features to a comparable scale, and dimensionality reduction methods like PCA, t-SNE, or UMAP were employed to visualize complex datasets and streamline subsequent analyses. Overall, these steps were pivotal in readying the dataset for further exploration and modeling. These steps are discussed in detail in the Model Building Section.

B. Encoding Methods:

1) One-Hot Encoding:

One-hot encoding is a technique for representing categorical variables as binary vectors. In this encoding, the value of each variable is represented as a vector of length equal to the number of possible values, with only one element being 1 (representing the category to which the variable belongs) and all other elements being 0. In this work, we defined a function for one-hot encoding, which takes two parameters: the sequence to be encoded and a list of all possible amino acids. The function creates a two-dimensional array with the number of rows equal to the length of the sequence and the number of columns equal to the number of types of amino acids, it iterates through each amino acid in the sequence, and sets the value of the corresponding position to 1 if the amino acid is in a given list of amino acids, indicating that the amino acid appears, otherwise it is 0.

Using a one-hot representation to encode TCR sequences has several limitations in downstream analysis. One major limitation is the inability to capture the sequential nature of amino acids, resulting in a loss of important structural and functional information. Additionally, one-hot encoding can lead to high-dimensional sparse representations, especially for TCR sequences with variable lengths of CDR3 regions, making it challenging to analyze and compare sequences effectively. To overcome these limitations, a possible approach is to incorporate information about the CDR3 length distribution.

2) BLOSUM 62 ENCODING:

Since one-hot encoding represents amino acids as binary vectors without considering their biochemical properties or evolutionary relationships, it may not capture the nuanced differences between sequences effectively. In contrast, we utilized the BLOSUM62 (BLOCKS SUBSTITUTION MATRIX 62) for distance matrix calculation. BLOSUM62 assigns scores to pairs of amino acids based on their observed frequencies in related proteins, capturing both the similarities and differences between sequences in a biologically meaningful way. This approach provides a more accurate representation of sequence similarity and is better suited for our analysis compared to one-hot encoding.

BLOSUM 62 encoding is a protein sequence encoding method that is based on the score of the BLOSUM (Blocks

Substitution Matrix) substitution matrix. The BLOSUM 62 substitution matrix is derived from statistical analysis of a large number of known protein sequences to assess the similarity between two amino acids. In the BLOSUM 62 encoding, each amino acid is mapped to a vector whose length is equal to the number of rows of the substitution matrix. Each element in the vector represents the similarity score of that amino acid to the corresponding row in the substitution matrix.

3) GIANA ENCODING:

GIANA Encoding converts TCR sequences into numerical representations of fixed dimensions. GIANA first encodes the CDR3 sequence to numeric vectors using ordinal coding. The key computational innovation in GIANA encoding lies in the application of serial non-commuting linear transformations to these numeric vectors. This process leverages a unitary transformation matrix that is a member of a 6-order cyclic group, enabling the transformation of CDR3 sequence data into high-dimensional Euclidean space. This encoding scheme not only enhances computational efficiency by reducing the dimensionality and complexity of the data but also retains the biological interpretability necessary for effective TCR analysis.

C. Distance Calculation Methods:

1) *TCRdist*: We employed TCRdist, a specialized tool for computing pairwise distances between T-cell receptor sequences. We divided the dataset into six subsets: human alpha, human beta, combined human alpha-beta, mouse alpha, mouse beta, and combined mouse alpha-beta as shown in Fig.1. which shows the subdivision of data for TCR Distance Matrix Calculation. Each subset was processed to calculate the TCR distance matrix for both alpha and beta chains. The distance matrices were essential for quantifying the similarity between TCR sequences, which is a crucial step for further analysis such as clustering and dimensionality reduction.

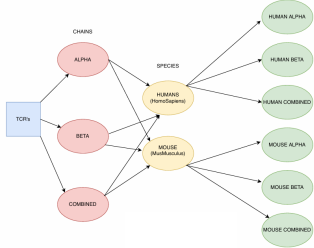


Fig. 1. Subdivision of TCR Distance Matrix Calculation

TCRdist [2] [3] is a toolkit for calculating TCR distances, which can be calculated based on their sequence characteristics (e.g., CDR3 sequences) and possibly other information (e.g., expression patterns of TCRs). For a given set of TCR sequences, the distances between all possible pairs are calculated, and the results are formed into a distance matrix, each element in the matrix represents the distance or similarity between the pairs, and once the distance matrix is calculated,

visualization tools can be used to show the similarity or difference between the TCRs.

2) *Levenshtein*:

Levenshtein distance is a measure of how similar two strings are. The Levenshtein distance, often referred to as edit distance, quantifies the minimum number of single-character edits (insertions, deletions, or substitutions) required to change one string into another. In the context of T-cell receptor (TCR) analysis, such as with TCRdist, Levenshtein distance is used to compute the similarities between TCR sequences by measuring how many changes are needed to convert one sequence into another. This approach can be especially useful in immunological studies where the similarity of TCR sequences can indicate clonal expansions or shared antigen specificity among different T cells. By constructing a distance matrix based on Levenshtein distances between all pairs of TCR sequences, we can perform clustering analysis, enabling the identification of TCR sequence groups that might respond to the same antigens, thereby providing valuable insights into immune system dynamics and pathogen recognition.

D. Dimensionality Reduction Method:

1) *UMAP*:

Uniform Manifold Approximation and Projection (UMAP) is a non-linear dimensionality reduction technique that is particularly efficient for high-dimensional data. It operates by constructing a high-dimensional graph representing the proximity of data points, which it then optimally projects onto a lower-dimensional space, making it particularly useful in bioinformatics for understanding T-cell receptor (TCR) complexities. In our work, UMAP can be employed to reduce the high-dimensional TCR data into a two-dimensional space for each of the TCR types: alpha chains, beta chains, and the combined alpha-beta chains. By plotting these reduced dimensions and coloring them based on antigen specificity, UMAP helps reveal underlying patterns and groupings that are not immediately apparent in higher-dimensional data. This visualization allows for a direct comparison of how TCR specificities cluster within each chain type. The differences in clustering patterns between the alpha, beta, and combined datasets can provide insights into the distinctiveness of their antigen recognition properties. By analyzing these plots, we can assess the consistency of TCR specificity across different chain types and potentially identify unique or shared characteristics among them.

E. Clustering Methods

1) *DBSCAN*:

The DBSCAN (Density-Based Spatial Clustering of Applications with Noise) algorithm is an effective density-based clustering method that identifies clusters as high-density areas separated by low-density areas. This approach is particularly advantageous in immunology for clustering T-cell receptor (TCR) sequences, as it does not require pre-specification of the number of clusters and can handle noise

and outliers effectively. DBSCAN operates by identifying 'core' samples that have a minimum number of neighbors within a given radius and marking as 'border' points those that are within the radius but do not meet the core criteria. This methodology can adeptly group TCRs that exhibit high similarity in sequence and functionality, allowing researchers to discern patterns of immune response and identify unique or shared antigenic encounters across different T cells.

For assessing the clustering quality of TCRs based on their specificity, metrics such as Purity Fraction, Purity Retention, and Normalized Mutual Information (NMI) are employed. Purity Fraction measures the homogeneity of the clusters, where a higher purity indicates that more sequences within a cluster recognize the same antigen. Purity Retention quantifies how much of the initial purity (before clustering) is retained in the clusters, reflecting the algorithm's ability to preserve natural grouping in the data. NMI, a normalization of the Mutual Information score that adjusts for chance, evaluates the overall clustering performance by measuring how much information about the true class labels is gained by the clustering, providing insights into the effectiveness of the clustering in grouping TCRs by their antigenic specificity. Together, these metrics provide a comprehensive assessment of how well the TCRs are clustered in relation to their specificity, guiding improvements in clustering strategies and understanding of immune diversity.

2) *Agglomerative Hierarchical Clustering*: Agglomerative Hierarchical Clustering is a bottom-up clustering algorithm that gradually merges data points into clusters, and represents the similarity relationship between data points through a tree-like structure. This method is particularly suited for biological data like T-cell receptor (TCR) sequences, where it can provide detailed insights into the relationship and similarity between various sequences. Agglomerative Hierarchical Clustering takes the distance between data points as a measure of similarity, and then gradually merges the data points into different clusters based on preset parameters (number of clusters, distance metric, and linking method).

The process results in a dendrogram that illustrates the sequence of cluster mergers and can be cut at various levels to obtain a specific number of clusters. This method allows researchers to analyze the hierarchical structure of TCR data, potentially revealing new insights into how TCR diversity relates to immune response functions. Similar to the DBSCAN Clustering, In the context of TCR clustering, metrics such as Purity Fraction, Purity Retention, and Normalized Mutual Information (NMI) are used for evaluating the effectiveness of agglomerative hierarchical clustering.

F. Model Building (Classification Models):

1) *Data Preparation*: The dataset contained several columns that were not necessary for building the predictive

models. By retaining key features the dataset is refined to include only those attributes that provide essential information about the TCRs and their interactions with antigens. Epitopes with fewer than 10 occurrences are considered insufficient for reliable modeling and were thus filtered out. This focuses the dataset on more common epitopes, which improves the model's ability to learn relevant patterns and make accurate predictions. We Calculated the length of each CDR3 sequence and retaining those within a practical range (10 to 20 amino acids) as depicted by the box plot shown in fig 2, as very short or very long sequences might represent sequencing errors or unusual variations that could skew the model training.

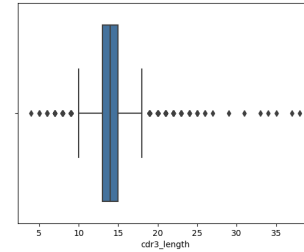


Fig. 2. Distribution of cdr3 sequence lengths

We counted the occurrences of each species and filtered to include only those species with sufficient representation (over 1000 occurrences), thus focusing on more common contexts in the training data. Further, we filtered the dataset to exclude rows where complex.id is 0. This step removed entries that did not meet certain criteria essential for our analysis, such as incomplete data entries or placeholders that do not correspond to actual TCR sequences. By doing this, we ensured that the dataset consists only of meaningful and relevant entries, enhancing the quality of our training data.

This model first uses a combination of Label Encoding and Binary Encoding to encode features like 'gene', 'species', 'epitope', 'mhc.a', 'mhc.b', 'v.segm' and 'j.segm' and 'mhc.class'. Next, using a function `GIANA_encoder_pd`, the CDR3 sequence is encoded into a vector representation, first traversing each amino acid sequence, for each sequence, removing the C and J segments at the beginning, and the F segment at the end, encoding the sequence using a predefined M6 matrix. For each row, the encoded CDR3 sequence and the normalized eigenvalues are first concatenated, and then they are combined into a single array. Finally, the resulting array contains the encoded CDR3 sequence and all the normalized eigenvalues, so that each row represents a complete sample eigenvector.

2) *Data Selection*: We made a decision to focus exclusively on Human Species (*Homo Sapiens*) for model building and predictions due to the significantly larger volume of data available—51,535 entries for humans compared to just 4,173 for mice (*Mus Musculus*). This stark disparity allows for a more robust and generalizable model for humans,

where the richer dataset ensures more reliable predictive performance. Additionally, focusing on human data aligns with the primary goal of applying findings directly to human medicine, avoiding the risks of undertraining and overfitting associated with the limited mouse data. This approach not only makes efficient use of the extensive human data but also enhances the clinical relevance of the predictions.

We subdivided the dataset to obtain the alpha and the beta chains and performed model evaluations. Each chain potentially interacts differently with antigens, influencing the specificity and strength of immune responses. By analyzing them separately, it's possible to isolate and understand the unique contributions of each chain to antigen recognition. This approach allows for a more nuanced analysis of TCR behavior, providing insights that could be obscured in a combined model. Moreover, separate evaluations enable the identification of chain-specific patterns, which are crucial for developing targeted therapeutic strategies and enhancing the accuracy of predictive immunological models.

3) *Logistic Regression (Baseline Model)*: The baseline model for predicting T-cell receptor (TCR) specificity utilizes logistic regression, a choice driven by the model's robustness and simplicity for binary classification tasks. During training, the dataset was split into training and testing sets to validate the model's performance on unseen data. The primary evaluation metric selected was the F1 score, due to the imbalanced nature of the dataset where accuracy alone could be misleading. This was particularly important given the variability in epitope representation within the dataset.

4) *SVM Model*: In this project, the Support Vector Machine (SVM) classifier was employed as another method for building predictive models to understand antigen specificity from T-cell receptor (TCR) sequences. The SVM was chosen for its effectiveness in high-dimensional spaces, as it excels in finding the optimal hyperplane that maximizes the margin between different classes. Specifically, the linear kernel was utilized to handle the linearly separable aspects of the data. The training process involved splitting the dataset into training and testing sets to both train the SVM and validate its performance on unseen data. This method provided a systematic approach to model antigen recognition patterns in an attempt to predict responses based on the molecular features of the TCR sequences.

5) *Random Forest Classifier Model*: The Random Forest classifier demonstrates a marked improvement over the logistic regression and SVM models primarily due to its ensemble learning approach, which integrates multiple decision trees to produce a more stable and accurate prediction by averaging the results. This method effectively reduces the variance and overfitting problems often seen in single decision tree models, leading to enhanced generalizability across diverse datasets. Specifically, the Random Forest model exhibited

higher weighted average F1 scores, indicating superior performance in balancing precision and recall across various classes within the imbalanced dataset.

G. Model Tuning

Hyperparameter tuning of the Random Forest model significantly enhanced its performance for both Alpha and Beta chains, illustrating the importance of optimizing model parameters in handling complex immunological datasets. The tuning was conducted using RandomizedSearchCV, where parameters such as the number of estimators, maximum features, maximum depth, minimum samples split, minimum samples leaf, and the use of bootstrap were varied across a predefined grid. This process was repeated across 10 iterations with 3-fold cross-validation to maximize robustness and effectiveness.

These traits make Random Forest particularly effective in handling the complex patterns and intricacies of antigen specificity prediction, where robustness against diverse and skewed data distributions is crucial. This enhanced capability to deliver more reliable and consistent predictions across a wider range of T-cell receptor specificities underlines why Random Forest is a preferable choice over the simpler logistic regression and the linear-bound SVM in this context.

H. Figures and Tables

a) *Positioning Figures and Tables*: Place figures and tables at the top and bottom of columns. Avoid placing them in the middle of columns. Large figures and tables may span across both columns. Figure captions should be below the figures; table heads should appear above the tables. Insert figures and tables after they are cited in the text. Use the abbreviation

TABLE I
TABLE TYPE STYLES

Table Head	Table Column Head		
	Table column subhead	Subhead	Subhead
copy	More table copy ^a		

^aSample of a Table footnote.

Figure Labels: Use 8 point Times New Roman for Figure labels. Use words rather than symbols or abbreviations when writing Figure axis labels to avoid confusing the reader. As an example, write the quantity "Magnetization", or "Magnetization, M", not just "M". If including units in the label, present them within parentheses. Do not label axes only with units. In the example, write "Magnetization (A/m)" or "Magnetization {A[m(1)]}", not just "A/m". Do not label axes with a ratio of quantities and units. For example, write "Temperature (K)", not "Temperature/K".

IV. RESULTS AND DISCUSSION

Reporting on the experiments with discussion on insights. Technical challenges are to be discussed here too.

pre-tuning scores, indicating a more precise balance between recall and precision in predicting various classes. Similarly, the Beta chain also showed improved performance with the same weighted average F1 score of 0.91. These enhancements suggest that the tuned Random Forest model is more adept at managing the skewed distributions and diversity of the dataset, leading to more reliable predictions. This improvement is particularly significant in a clinical setting, where accurate and dependable model predictions are crucial for developing targeted immunotherapies based on T-cell receptor specificities.

V. FURTHER WORK AND IMPROVEMENT

A. Short Term Improvements:

1) *Feature Engineering and Optimization:* With more time, further refinement of the feature set could be beneficial. This could involve exploring more complex features derived from the existing data, such as higher-order interactions between features or more sophisticated transformation of the CDR3 lengths.

2) *Model Tuning and Regularization Techniques:* While initial models have been established, additional tuning of hyperparameters and the application of regularization techniques could enhance model performance and prevent overfitting.

3) *Cross-Validation and Robustness Checks:* Implementing a more rigorous cross-validation framework to ensure the models are robust and generalizable across different data splits, which would help in assessing the stability of the predictions.

4) *Incremental Data Integration:* We can start incorporating incremental learning capabilities so that the model can continuously learn and improve from new data without the need for retraining from scratch.

B. Long Term Extensions:

1) *Artificial Intelligence and Deep Learning Models:* In the long run, integrating deep learning approaches like convolutional neural networks (CNNs) or recurrent neural networks (RNNs) could significantly enhance the ability to model complex patterns in antigen specificity more effectively. This could also include exploring transformer models that have shown success in other sequence-based tasks in bioinformatics.

2) *Automated Pipeline for Real-Time Analysis:* Building an automated pipeline that integrates with existing healthcare systems to provide real-time analysis and predictions. This system could use the latest models to provide immediate insights into patient data, helping clinicians make faster and more informed decisions.

3) *Expansion to other Species or Conditions:* Broadening the scope of the model to include other species or additional immunological conditions could vastly increase its utility. This would involve collecting and integrating diverse datasets, possibly leading to a more universally applicable tool.

4) *Strategic Partnerships and Collaborative Research:* Establishing partnerships with research institutions and other companies in the biotech sector could lead to shared innovations and enhancements in predictive modeling. Collaborative

projects could also provide access to proprietary datasets and specialized knowledge.

VI. CONCLUSION

A brief summary of the key insights in your report.

REFERENCES

- [1] H. Zhang, X. Zhan, and B. Li, "GIANA allows computationally-efficient TCR clustering and multi-disease repertoire classification by isometric transformation," *Nature Communications*, vol. 12, no. 1, pp. 4699, 2021.
- [2] M. Vujovic, K. F. Degen, F. I. Marin, et al., "T cell receptor sequence clustering and antigen specificity," *Computational and Structural Biotechnology Journal*, vol. 18, pp. 2166-2173, 2020.
- [3] K. Mayer-Blackwell, S. Schattgen, L. Cohen-Lavi, et al., "TCR meta-clonotypes for biomarker discovery with tcrdist3 enabled identification of public, HLA-restricted clusters of SARS-CoV-2 TCRs," *Elife*, vol. 10, p. e68605, 2021.

APPENDIX

The document up to this section should be no more than 8 pages. The appendix section is optional. You can include additional material here, but it will not be marked.