AI-Driven TCR Design: Leveraging Large Language Models for Personalized Cancer Immunotherapy

Sidan Yao^a, Yee Mun Lee^b, Yi Tian Png^b, <u>Manna Dai</u>^a, Jun Zhou^a, Tao Luo^a, Rick, Siow Mong Goh^a, Yong Liu^a, Chwee Ming Lim^b

^a Institute of High Performance Computing (IHPC) - Agency for Science, Technology and Research (A*STAR)

^b Singapore General Hospital

1. Introduction

T cell receptor-engineered T cell (TCR-T) therapy is a promising approach in cancer immunotherapy, enabling the selective targeting of tumor-associated antigens. However, its clinical efficacy is limited by key challenges, including the scarcity of truly tumorspecific antigens, which increases the risk of offtarget toxicity. Unlike CAR-T therapy, TCR-T relies on recognizing intracellular antigens presented by MHC, making antigen selection complex due to genetic variability in HLA alleles. Furthermore, the many-to-many binding nature of antigen-TCR interactions complicates the design of highly specific therapeutic TCRs, making it difficult to distinguish between functional and non-functional candidates.

A major hurdle in optimizing TCR-T therapy lies in computational modeling of TCR-antigen binding. Current models rely on manually curated datasets that are heavily biased toward Western populations, failing to capture the genetic diversity of Asian patients. This limits their predictive accuracy for non-Western populations. Additionally, the lack of experimentally validated negative binding data leads to overfitting, as models are trained primarily on positive interactions, reducing their generalizability.

To address these challenges, we propose an AIdriven framework that leverages large language models (LLMs) to learn antigen-TCR-HLA interactions. By treating antigen and TCR sequences as structured language data, our approach enables the generation of de novo TCRs with enhanced specificity and binding affinity. Integrating multimodal learning, our model incorporates protein structure analysis and physicochemical properties to refine predictions. Trained on a diverse dataset combining publicly available repositories and experimental data from Singapore General Hospital (SGH), this approach ensures better representation of underrepresented HLA alleles, particularly those prevalent in Asian populations. By reducing overfitting and improving generalizability, our framework aims to optimize personalized cancer immunotherapy, making TCR-T therapy safer and more effective across diverse patient populations.

2. Methodology

2.1 AI-Driven Antigen-TCR Binding Prediction

Our framework employs a fine-tuned large language model (LLM) to embed antigen sequences into high-dimensional vector spaces, enabling accurate binding affinity predictions. This approach surpasses traditional models that rely on manually curated peptide datasets, which are often incomplete and lack experimentally validated negative binding data.

Step 1: Data Preprocessing and Representation Learning

We aggregate antigen-TCR binding data from multiple sources, including publicly available databases (IEDB, NetMHCpan) and experimental datasets from Singapore General Hospital (SGH). Sequences are tokenized using a domain-specific tokenizer optimized for peptide sequences, ensuring efficient representation learning.

Step 2: Large Language Model Fine-Tuning

A pre-trained foundation model is fine-tuned on antigen-TCR sequence pairs, incorporating structural and biochemical constraints to enhance specificity. This process aligns embeddings with biologically relevant interactions, improving predictive accuracy.

Step 3: Binding Score Prediction using Domain Adaptation

To address the scarcity of experimentally validated negative binding data, we employ domain adaptation techniques to infer binding probabilities without artificial negative samples. This enhances generalizability beyond known antigen-TCR pairs, mitigating overfitting.

2.2 De Novo TCR Generation for Novel Antigens

A key feature of our framework is its ability to generate high-affinity, antigen-specific TCRs for previously unseen targets. This is achieved through a generative transformer-based decoder that synthesizes TCR sequences conditioned on antigen embeddings.

2.3 Clustering and Selection of High-Affinity TCRs

Deep clustering techniques categorize candidate TCRs based on predicted binding affinity scores. The top-ranked sequences undergo further validation via structure-based docking simulations before experimental testing, ensuring optimal therapeutic potential.



Fig. 1: This figure illustrates the AI-driven framework for TCR generation and antigen-TCR binding prediction. The process begins with a customized tokenizer, which encodes antigen (X) and TCR (Y) sequences. The Antigen Encoder converts antigen data into a meaningful antigen embedding, while the TCR Encoder processes TCR sequences into a high-dimensional TCR embedding. The model optimizes binding affinity by minimizing Embedding Distance Loss and ensures generated TCRs maintain functional diversity by incorporating TCR Distance Loss. The TCR Decoder then generates new TCR sequences with improved specificity. This AI-driven approach accelerates therapeutic TCR discovery, enhancing personalized cancer immunotherapy.

3. Results and Preliminary Validation

Our preliminary results demonstrate that the proposed LLM-based framework surpasses conventional HLA-peptide binding predictors in both accuracy and generalization, offering a more robust solution for antigen-TCR interaction modeling. Key findings highlight significant improvements in predictive performance, novel TCR generation, and error minimization:

Enhanced Binding Prediction Accuracy: Benchmarking against established tools such as NetMHCpan and MHCflurry reveals a 15% increase in positive predictive value when evaluated on an independent dataset enriched with Asian-specific HLA alleles. This improvement underscores the model's ability to better accommodate diverse HLA profiles, addressing a key limitation in existing methods that are biased toward Western populations. Moreover, the LLM-driven embedding of antigen sequences captures intricate biochemical and structural interactions that traditional models often overlook, leading to higher precision in identifying true binding events and minimizing misclassification of weak binders.

Effective TCR Generation for Novel Antigens: A core strength of our framework lies in its generative capability. Unlike existing approaches that rely on predefined antigen-TCR pairs, our model successfully synthesizes TCR sequences for antigens absent from the training data. Notably, 80% of the generated TCRs exhibit predicted binding affinities comparable to experimentally validated sequences, demonstrating the model's ability to infer functionally relevant interactions. This capability is critical for applications in personalized cancer immunotherapy, where patient-specific neoantigens often lack known TCR counterparts. Additionally, structure-based docking simulations confirm that the generated TCRs adopt conformations conducive to high-affinity binding, reinforcing the model's practical utility.

Reduction in False Positives: Conventional models suffer from high false positive rates due to artificial negative sampling. Our domain adaptation strategy infers binding probabilities without explicitly labeled non-binding pairs, reducing false positives and enhancing reliability. Improved specificity minimizes erroneous candidate selection, crucial for preventing severe off-target effects in therapy.

These results collectively highlight the potential of our AI-driven approach to transform antigen-TCR modeling, paving the way for more accurate, scalable, and clinically relevant immunotherapy strategies. Future work will focus on further refining the generative model, expanding its applicability across diverse HLA backgrounds, and integrating real-world experimental validation to solidify its translational impact.

4. Discussion and Impact

The ability to computationally generate highly specific TCRs has profound implications for cancer immunotherapy. Our model enables:

Personalized Immunotherapy: By incorporating diverse HLA alleles, the method is particularly suited for developing TCR-T therapies tailored to Asian cancer patients. Scalability and Efficiency: The AI-driven framework reduces reliance on wet-lab screening, accelerating the identification of therapeutic TCRs while lowering costs. Generalization to Other Diseases: Beyond cancer, this approach can be extended to infectious diseases and autoimmune disorders requiring precise antigen recognition.

5. Conclusion and Future Work

This study presents an AI-driven framework for antigen-TCR binding prediction and de novo TCR generation, addressing critical challenges in personalized cancer immunotherapy. By leveraging large language models and domain adaptation techniques, we have improved the specificity and scalability of TCR discovery. Our approach offers a computationally efficient and biologically relevant solution to accelerate the development of TCR-based therapies, particularly for underrepresented populations with diverse HLA alleles.

By integrating binding prediction with generative modeling, we have demonstrated that AI can not only classify antigen-TCR interactions but also generate high-affinity TCR candidates tailored to specific tumor antigens. The results indicate that this model significantly reduces false positives, improves binding affinity predictions, and enhances the diversity of generated TCRs.