

# Generative Design of High-affinity T-cell Receptors by Progressive Learning with Structural Confidence

Xinyuan Zhu<sup>\*1</sup>, Qingshuo Jin<sup>\*2</sup>, Jiadong Lu<sup>2</sup>, Jun Wu<sup>2</sup>, Lixia Zhu<sup>3,4,5</sup>, Yuhua Shang<sup>4</sup>, Tengchuan Jin<sup>†3,4,5</sup>, and Fuli Feng<sup>†2</sup>

<sup>1</sup>School of Information Science and Technology, University of Science and Technology of China

<sup>2</sup>School of Artificial Intelligence and Data Science, University of Science and Technology of China

<sup>3</sup>Laboratory of Structural Immunology, State Key Laboratory of Immune Response and Immunotherapy, Division of Life Sciences and Medicine, University of Science and Technology of China

<sup>4</sup>Anhui Genebiol Biotech. Ltd

<sup>5</sup>Institute of Health and Medicine, Hefei Comprehensive National Science Center

## Supplementary Information

### Contents

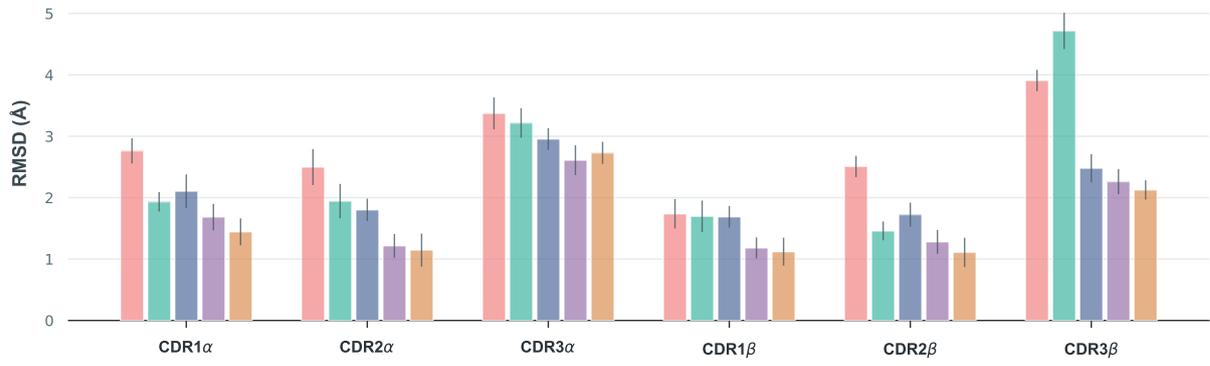
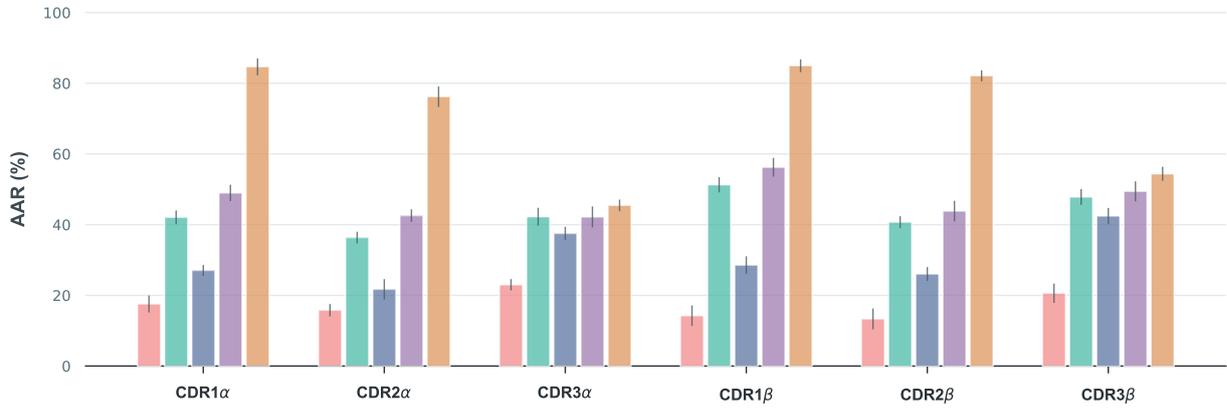
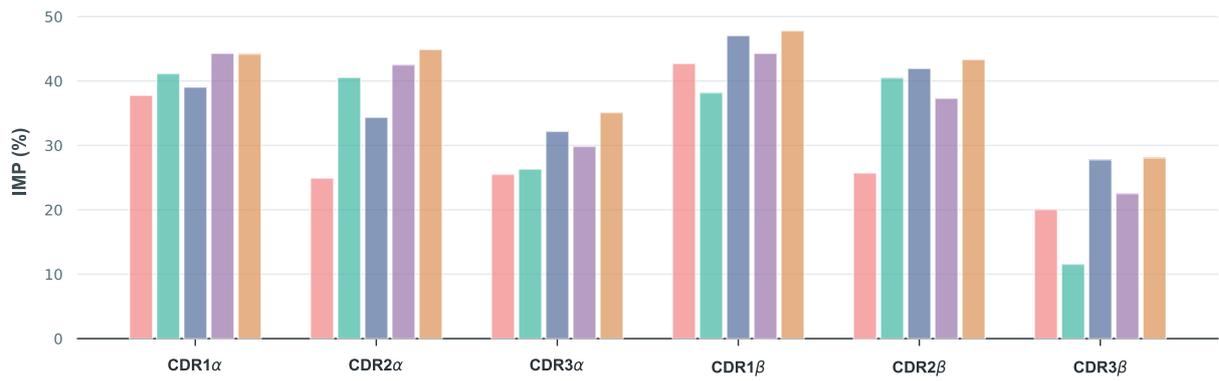
|          |  |          |
|----------|--|----------|
| <b>1</b> | <b>Supplementary Figures</b>                     | <b>1</b> |
|          | Supplementary Fig. 1 . . . . .                   | 1        |
|          | Supplementary Fig. 2 . . . . .                   | 1        |
|          | Supplementary Fig. 3 . . . . .                   | 1        |
|          | Supplementary Fig. 4 . . . . .                   | 1        |
| <b>2</b> | <b>Supplementary Tables</b>                      | <b>5</b> |
|          | Supplementary Table 3 . . . . .                  | 5        |
|          | Supplementary Table 4 . . . . .                  | 5        |
|          | Supplementary Table 5 . . . . .                  | 5        |
| <b>3</b> | <b>Supplementary Note</b>                        | <b>7</b> |
|          | 3.1 Details of numbering adaptation . . . . .    | 7        |
|          | 3.2 TCR full sequence reconstruction . . . . .   | 7        |
|          | 3.3 Surface plasmon resonance protocol . . . . . | 7        |

---

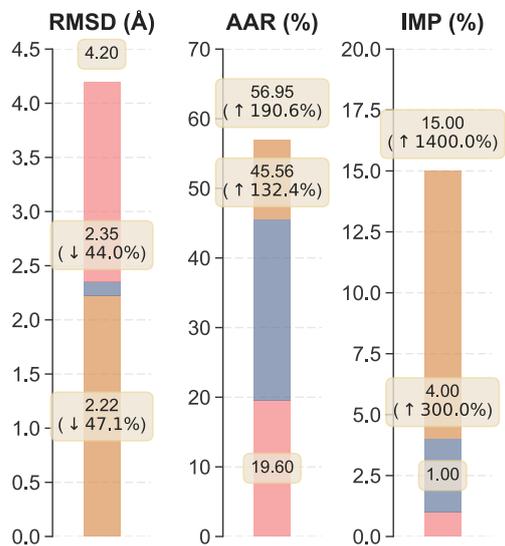
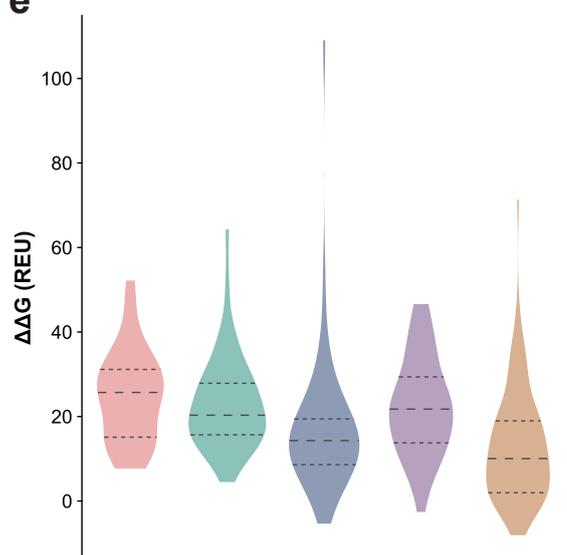
\*These authors contributed equally to this work.

†Correspondence: jint@ustc.edu.cn; fengfl@ustc.edu.cn

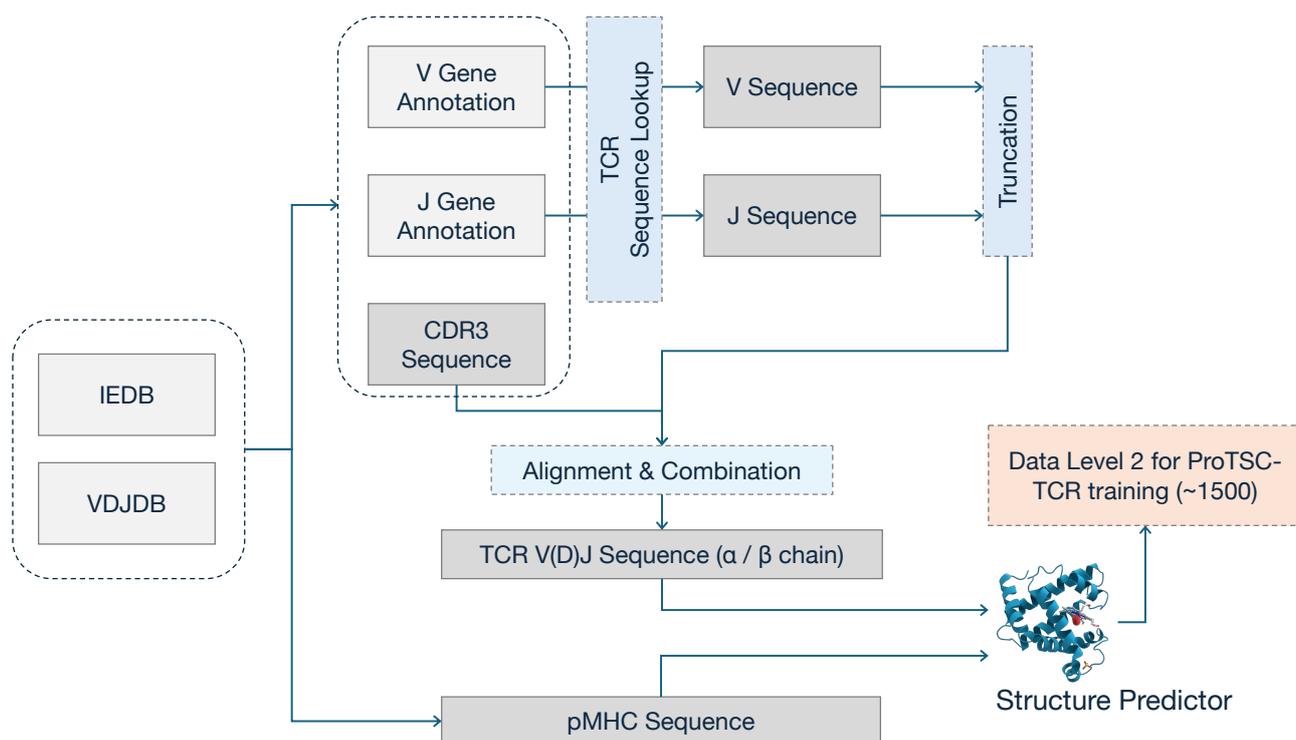


**a****b****c**

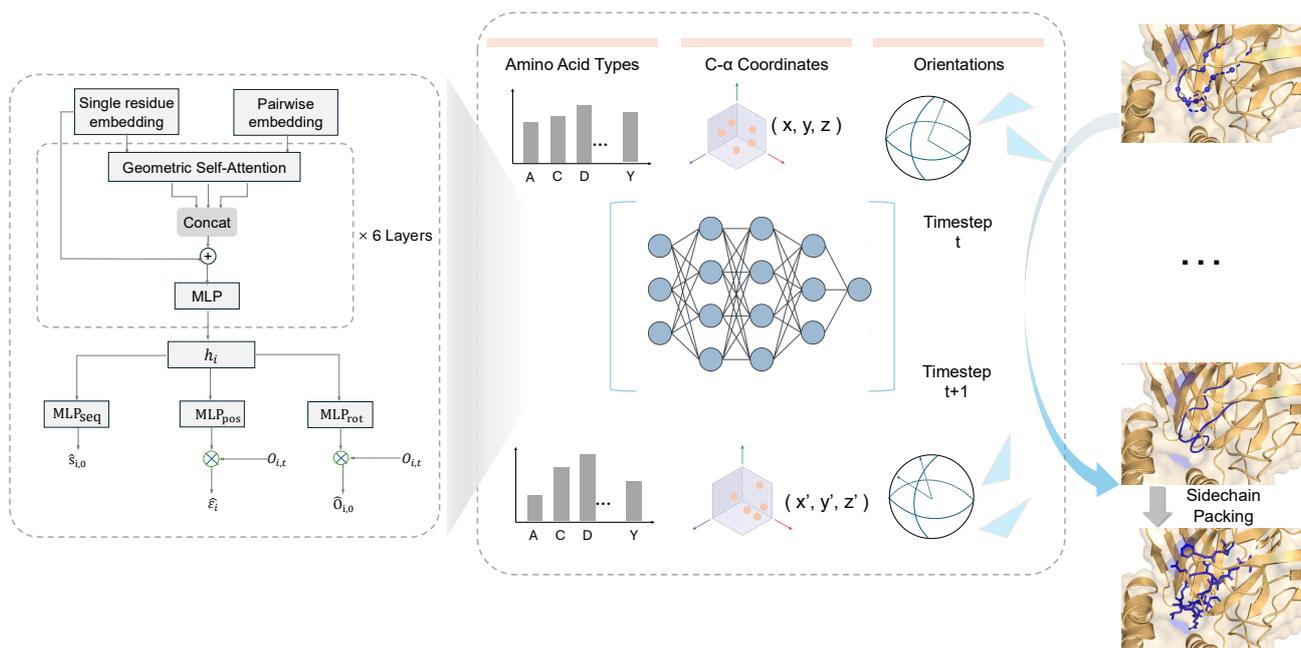
■ Level 1 Only   
 ■ Level 3 Only   
 ■ Level 1 + 2   
 ■ Level 1 + 3   
 ■ ProTSC-TCR

**d****e**

**Supplementary Fig. 2 (previous page) | Ablation across CDR regions and an additional case study under the *Top-100* protocol. a–c**, Performance across the six CDR regions (CDR1 $\alpha$ , CDR2 $\alpha$ , CDR3 $\alpha$ , CDR1 $\beta$ , CDR2 $\beta$ , CDR3 $\beta$ ) on the 35-target benchmark. Bars report backbone RMSD (**a**), AAR (**b**), and IMP (**c**) for variants: Level 1 Only, Level 3 Only, Level 1+2, Level 1+3, and the ProTSC-TCR. **d**, Stepwise improvements on an additional target (PDB ID 8SHI), highlighting how incorporating Data Level 2 and 3 contributes to superior RMSD, AAR, and IMP than Level 1 Only model. **e**, Distribution of Rosetta  $\Delta\Delta G$  over 100 independent designs for 8SHI, showing a progressive shift toward more favorable predicted binding energetics with ProTSC-TCR.



**Supplementary Fig. 3 | Workflow for full TCR sequence reconstruction and construction of the Level 2 data.** To construct the dataset from binding records that typically lack full variable-domain sequences, a reconstruction pipeline was implemented. Germline V and J gene segments are retrieved based on annotated gene identifiers and truncated at conserved recombination motifs. These segments are then aligned and combined with the specific CDR3 sequence to reconstruct complete TCR  $\alpha$  and  $\beta$  V(D)J sequences. Finally, the reconstructed TCR sequences are paired with pMHC sequences and processed by a structure predictor (ColabFold-AF2) to generate 1,500 predicted TCR–pMHC complexes for ProTSC-TCR training.



**Supplementary Fig. 4 | Model architecture and generative inference workflow of ProTSC-TCR.** The model utilizes an SE(3)-equivariant architecture (left) that processes single-residue and pairwise geometric embeddings through invariant self-attention layers to produce context-aware hidden representations ( $h_i$ ). These representations are decoded by three modality-specific heads ( $MLP_{seq}$ ,  $MLP_{pos}$ , and  $MLP_{rot}$ ). The geometric predictions are computed in local frames and projected globally to ensure equivariance. The inference pipeline (right) demonstrates the iterative denoising of the joint state (sequence, C $\alpha$  coordinates, and orientations) from random noise, followed by side-chain packing to generate the final full-atom structure.

## 2 Supplementary Tables

**Supplementary Table 3 | Experimental SPR affinities and interface features for top-ranked KRAS-G12V TCR designs.** The wild-type reference complex (PDB ID 8I5D) and the top-10 energy-ranked ProTSC-TCR CDR3 $\beta$  designs. For each unique CDR3 $\beta$  sequence, we report the FoldX-predicted binding energy change ( $\Delta\Delta G$ ), the number of interfacial hydrogen bonds (#HB) and salt bridges (#SB), and the buried interface area ( $\text{\AA}^2$ ) computed from the structure generated by ProTSC-TCR, together with the SPR-measured dissociation constant ( $K_D$ ,  $\mu\text{M}$ ). Duplicate-ranked entries that share an identical CDR3 $\beta$  sequence are highlighted with matching colors, share the same measured  $K_D$ , and omit repeated structural statistics.

| ID                             | $\Delta\Delta G$ (FoldX) | #HB | #SB | Interface Area | KD ( $\mu\text{M}$ ) | CDR3 $\beta$            |
|--------------------------------|--------------------------|-----|-----|----------------|----------------------|-------------------------|
| <i>Reference (PDB ID 8I5D)</i> | –                        | 2   | 3   | 349.2          | 37.8                 | <b>CASSLEGTVEETLYF</b>  |
| <b>Design-1</b>                | -33.60                   | 4   | 3   | 474.7          | 22.9                 | <b>CASLLASAFGELYF</b>   |
| <b>Design-2</b>                | -30.57                   | 6   | 3   | 560.7          | 8.17                 | <b>CASSLRGGYFQPLHF</b>  |
| <b>Design-3</b>                | -28.79                   | –   | –   | –              | 22.9                 | <b>CASLLASAFGELYF</b>   |
| <b>Design-4</b>                | -27.23                   | 5   | 3   | 538.9          | 14.4                 | <b>CASSWWGDGGYYLYF</b>  |
| <b>Design-5</b>                | -24.21                   | 6   | 3   | 548.6          | 23.2                 | <b>CASSFSGLDGQPQTF</b>  |
| <b>Design-6</b>                | -23.86                   | 3   | 3   | 545.8          | 17.1                 | <b>CASSILGEPGNYVYF</b>  |
| <b>Design-7</b>                | -23.19                   | –   | –   | –              | 37.8                 | <b>CASSLEGTVEETLYF</b>  |
| <b>Design-8</b>                | -21.91                   | 7   | 4   | 555.1          | 3.54                 | <b>CASSLWLDMFYTLFYF</b> |
| <b>Design-9</b>                | -20.69                   | 3   | 2   | 413.4          | 8.67                 | <b>CASSDWFGSLDTLFF</b>  |
| <b>Design-10</b>               | -19.73                   | –   | –   | –              | 3.54                 | <b>CASSLWLDMFYTLFYF</b> |

**Supplementary Table 4 | Experimental SPR affinities and interface features for top-ranked SARS-CoV-2 Spike-Y453F TCR designs.** The wild-type reference complex (PDB ID 8YE4) and the top-10 energy-ranked ProTSC-TCR CDR3 $\beta$  designs. For each unique CDR3 $\beta$  sequence, we report the FoldX-predicted binding energy change ( $\Delta\Delta G$ ), the number of interfacial hydrogen bonds (#HB) and salt bridges (#SB), and the buried interface area ( $\text{\AA}^2$ ) computed from the structure generated by ProTSC-TCR, together with the SPR-measured dissociation constant ( $K_D$ ,  $\mu\text{M}$ ). Duplicate-ranked entries that share an identical CDR3 $\beta$  sequence are highlighted with matching colors, share the same measured  $K_D$ , and omit repeated structural statistics.

| ID                             | $\Delta\Delta G$ (FoldX) | #HB      | #SB      | Interface Area | $K_D$ ( $\mu\text{M}$ ) | CDR3 $\beta$         |
|--------------------------------|--------------------------|----------|----------|----------------|-------------------------|----------------------|
| <i>Reference (PDB ID 8YE4)</i> | \                        | 0        | 0        | 155            | 6.37                    | CASSETGGYEQYF        |
| Design-1                       | -7.55                    | 3        | 7        | 252.7          | 16.9                    | CASSDEPFDEQYF        |
| Design-2                       | -3.79                    | 4        | 0        | 376            | 20.3                    | CASSLLGPGEQYF        |
| Design-3                       | -3.19                    | 2        | 3        | 269.9          | 12.1                    | CASSDDASNEQYF        |
| Design-4                       | -2.73                    | 3        | 4        | 251.7          | 9.65                    | CASSDGAGDEQYF        |
| <b>Design-5</b>                | <b>-1.43</b>             | <b>4</b> | <b>5</b> | <b>229.9</b>   | <b>4.71</b>             | <b>CASSDDAGDEQFF</b> |
| <b>Design-6</b>                | <b>0.32</b>              | <b>1</b> | <b>1</b> | <b>244.4</b>   | <b>3.97</b>             | <b>CASSDGRGNEQFF</b> |
| Design-7                       | 0.78                     | 0        | 0        | 277            | 8.86                    | CASSNGPYGEQYF        |
| Design-8                       | 0.88                     | 4        | 2        | 262            | 13.1                    | CASSYRLGDEQYF        |
| Design-9                       | 1.18                     | 5        | 2        | 220.2          | 11.1                    | CASSDGWPDEQFF        |
| Design-10                      | 1.26                     | 3        | 2        | 214.9          | 10.3                    | CASSYRLGDEQYF        |

**Supplementary Table 5 | Distribution of MHC Class and Alleles in the 35-target benchmark.**

| MHC Category | Allele/Locus   | Count     |
|--------------|----------------|-----------|
| HLA Class I  | A              | 15        |
|              | B              | 9         |
|              | C              | 1         |
|              | E              | 2         |
| HLA Class II | DP             | 2         |
|              | DQ             | 2         |
| Mouse MHC    | H-2 (Class I)  | 2         |
|              | I-A (Class II) | 2         |
| <b>Total</b> |                | <b>35</b> |

## 3 Supplementary Note

### 3.1 Details of numbering adaptation

---

**Algorithm 1:** Unified Receptor Numbering Adaptation (Chothia Scheme)

---

```
Input: Raw amino acid sequence  $S$ 
Input: Topology Constant Table  $\mathcal{T}_{topo}$  (contains boundary definitions for Heavy/Light topologies)
Output: Standardized Chain object  $\mathcal{C}_{final}$  with unified numbering
1  $(\mathcal{A}, \tau) \leftarrow \text{RunHMMAlignment}(S, \text{'chothia'})$ ;
2 if  $\tau \in \{IgH, TCR\beta, TCR\delta\}$  then
3    $\mathcal{B} \leftarrow \text{Lookup}(\mathcal{T}_{topo}, \text{Topology}_{heavy})$ ; ; // Retrieves {CDR1, CDR2, CDR3} bounds for Heavy-like
   chains
4 else
5    $\mathcal{B} \leftarrow \text{Lookup}(\mathcal{T}_{topo}, \text{Topology}_{light})$ ; ; // Retrieves {CDR1, CDR2, CDR3} bounds for Light-like
   chains
6 end
7  $\mathcal{C}_{final} \leftarrow \emptyset$ ;
8 foreach residue tuple  $(n, i, aa) \in \mathcal{A}$  do
9    $\rho \leftarrow \text{DetermineRegion}(n, \mathcal{B})$ ;
10   $p \leftarrow \text{CreatePosition}(n, i, \tau, \rho)$ ;
11   $\mathcal{C}_{final}.append((p, aa))$ ;
12 end
13 return  $\mathcal{C}_{final}$ ;
```

---

### 3.2 TCR full sequence reconstruction

To generate full T cell receptor (TCR) variable region sequences from available V(D)J gene annotations, we employed a standardized assembly pipeline. Germline amino acid sequences corresponding to V and J gene annotations were retrieved from the IMGT/Gene-DB repository. These sequences were precisely truncated at conserved recombination signal motifs and subsequently concatenated with the experimentally determined CDR3 hypervariable junctions to assemble the complete variable domains of both TCR  $\alpha$  and  $\beta$  chains. This reconstruction strategy maintains high data fidelity by directly querying up-to-date reference databases and explicitly omits the constant regions, thereby focusing the input features on the variable domains critical for discriminating specific binding interactions.

### 3.3 Surface plasmon resonance protocol

**Instrumentation and reagents** SPR assays were conducted on a Biacore 8K system (Cytiva) using Series S Sensor Chip CM5. The running buffer used throughout the experiments for immobilization, binding, and dissociation contained 20 mM HEPES (pH 7.5), 150 mM NaCl, and 0.05% Tween-20.

**Ligand Immobilization** The pMHC ligands (specifically pHLA-A\*11:01 for KRAS-G12V targets and pHLA-A\*24:02 for SARS-CoV-2 targets) were immobilized using standard amine coupling chemistry. The CM5 chip surface was first activated with a 1:1 mixture of 0.4 M 1-ethyl-3-(3-dimethylaminopropyl)carbodiimide hydrochloride (EDC) and 0.1 M N-hydroxysuccinimide (NHS). The ligands were diluted to a final concentration of 20  $\mu\text{g}/\text{mL}$  in 10 mM sodium acetate buffer (pH 4.5) and injected over the activated surface at a flow rate of 10  $\mu\text{L}/\text{min}$  until a response level of approximately 2000 RU was achieved. Residual activated groups were blocked by injecting 1 M ethanolamine-HCl (pH 8.5).

**Kinetic analysis** Binding kinetics were measured using a multi-cycle kinetics approach. TCR analytes were diluted in the running buffer to generate a series of five concentrations using a 2-fold serial dilution factor. Samples were injected from low to high concentration. For each cycle, the association phase was monitored for 120 s, followed by a dissociation phase of 600 s. Between cycles, the sensor surface was regenerated with a brief injection of 50 mM NaOH to remove any bound TCRs.

**Data evaluation** The resulting sensorgrams were analyzed using the Biacore Insight Evaluation Software. Equilibrium dissociation constants ( $K_D$ ) were derived by fitting the data to a 1:1 Langmuir binding model based on the multi-cycle kinetics data.