



OPEN ACCESS

EDITED BY

Jonathan S Duke-Cohan,
Dana–Farber Cancer Institute, United States

REVIEWED BY

Eun Young Choi,
Seoul National University, Republic of Korea
Santrupti Nerli,
Genentech, United States

*CORRESPONDENCE

H. Jabran Zahid

✉ hzahid@microsoft.com

Harlan Robins

✉ hrobins@adaptivebiotech.com

Jonathan M. Carlson

✉ carlson@microsoft.com

[†]Deceased

RECEIVED 31 March 2025

ACCEPTED 31 July 2025

PUBLISHED 27 August 2025

CITATION

Zahid HJ, Taniguchi R, Ebert P, Chow I-T,
Gooley C, Lv J, Pisani L, Rusnak M,
Elyanow R, Takamatsu H, Zhou W, Greissl J,
Robins H and Carlson JM (2025) Large-scale
statistical mapping of T-cell receptor β
sequences to human leukocyte antigens.
Front. Immunol. 16:1603730.
doi: 10.3389/fimmu.2025.1603730

COPYRIGHT

© 2025 Zahid, Taniguchi, Ebert, Chow, Gooley,
Lv, Pisani, Rusnak, Elyanow, Takamatsu, Zhou,
Greissl, Robins and Carlson. This is an open-
access article distributed under the terms of
the [Creative Commons Attribution License
\(CC BY\)](https://creativecommons.org/licenses/by/4.0/). 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.

Large-scale statistical mapping of T-cell receptor β sequences to human leukocyte antigens

H. Jabran Zahid^{1*}, Ruth Taniguchi², Peter Ebert^{2†}, I-Ting Chow²,
Chris Gooley¹, Jinpeng Lv¹, Lorenzo Pisani¹, Mikaela Rusnak²,
Rebecca Elyanow², Hiroyuki Takamatsu³, Wenyu Zhou²,
Julia Greissl¹, Harlan Robins^{2*} and Jonathan M. Carlson^{1*}

¹Microsoft Research, Redmond, WA, United States, ²Adaptive Biotechnologies, Seattle, WA, United States, ³Faculty of Transdisciplinary Sciences for Innovation, Institute of Transdisciplinary Sciences for Innovation/Department of Hematology, Kanazawa University, Kanazawa, Japan

Introduction: T-cell receptors (TCRs) interacting with peptides presented by human leukocyte antigens (HLAs) are the foundation of the adaptive immune system, but population-level analysis of TCR–HLA interactions is lacking.

Methods: We statistically associated approximately 10^6 public TCR β s to specific HLAs using TCR β repertoires sampled from 4,144 HLA-genotyped subjects. The TCR β s we associated were specific to unique HLA allotypes, not allelic groups, and to the paired α – β heterodimer of class II HLAs, though exceptions were observed.

Results: This specificity permitted highly accurate imputation of 248 class I and II HLAs from the TCR β repertoire. Notably, 45 HLA-DP and -DQ heterodimers lacked associated TCRs because they likely arise from non-functional trans-complementation. The public class I and II HLA-associated TCR β s we identified were primarily expressed on CD8⁺ and CD4⁺ memory T cells, respectively, which were responding to various common antigens.

Discussion: Our results recapitulate fundamental biology, provide insights into the functionality of HLAs, and demonstrate the power and potential of population-level TCR β repertoire sequencing.

KEYWORDS

human leukocyte antigen, T cell, immunology, T cell - HLA interaction, T cell repertoire

Introduction

The major histocompatibility complex (MHC) is a set of genes found in jawed vertebrates, which in humans encode for human leukocyte antigens (HLAs) (1). The primary function of HLAs is to present fragments of proteins (i.e., peptides or antigens) on the surface of cells for T cell recognition (2). TCRs on the surface of T cells interact with the peptides presented by HLAs (pHLA). The TCR–pHLA interaction is a key mechanism of

adaptive immunity and plays a central role in the immune system's response to infections, cancers, allergens and self-tissues targeted in autoimmunity and transplantation (3–10).

HLAs are both polygenic and polymorphic, allowing for a highly specific and fine-tuned adaptive immune response to diverse pathogens. The classical antigen-presenting MHC proteins—class I and class II HLAs—are found on the surface of all nucleated and professional antigen presenting cells, respectively. A large number of allelic variants have been identified across the six loci encoding class I (HLA-A, -B and -C) and class II (HLA-DP, -DQ and -DR) HLAs (11). Antigens presented to T cells bind to a single polymorphic α chain of class I HLAs (12, 13) and to the α and β chain of class II HLAs (14–16). Both the α and β chains are polymorphic for HLA-DP and -DQ whereas only the β chain is polymorphic for -DR. Furthermore, the β chain of HLA-DR may be encoded by four loci (*DRB1*, *DRB3*, *DRB4* and *DRB5*); all individuals have *DRB1* encoded on both instances of chromosome 6 and may additionally have one of *DRB3*, *DRB4* or *DRB5* on each chromosome 6.

HLA genes can be resolved to varying degrees by sequencing and several distinct naming systems can be found in the literature. Here we adopt the 2010 WHO HLA nomenclature (17). For each locus, the first two fields designate the allele group and specific protein (i.e. allotype), respectively. The third and fourth fields indicate synonymous substitutions in coding and non-coding regions, respectively. For class II HLAs, the α and β chains are encoded, sequenced and typed independently. Our sequencing of HLAs lacks information on parental haplotype (i.e., phasing). In the context of molecular epidemiology, identifying (or assuming) the resolution at which causal mechanisms (and thus, clinical associations) are likely operating remains challenging. HLAs in the same allelic group tend to present similar or identical peptides (18), share functional properties such as relative expression levels (19, 20), and serve as ligands for the same KIR receptors (21). Thus, many epidemiological and functional studies treat HLAs of the same allelic group interchangeably. However, structural (22), functional (23), and evolutionary evidence (24) suggest that, at least in some contexts, very similar HLA allotypes frequently interact with very different TCRs. A fundamental goal of this work is to establish the relationship between TCR specificity and HLA resolution and to explore the TCR specificity of class I and class II HLAs.

The human body maintains a diverse set of naive T cells where antigen specificity is determined by TCRs (25, 26). These T cells are selected such that their TCRs, which are generated via V(D)J recombination, interact with pHLAs in the thymus (27–30). Interaction with class I and class II pHLAs directs differentiation into CD8⁺ (cytotoxic T cell) (31, 32) and CD4⁺ (helper T cell) lineages (33), respectively. Antigen presentation by an HLA and subsequent TCR recognition in the appropriate immunological context triggers clonal expansion of naive T cells resulting in a large population of T cells expressing identical cognate TCRs (34). Clonal expansion of T cells with the same TCR greatly increases the chance of sampling these TCRs experimentally. As a result, subjects with matching HLAs and shared antigenic exposure have a

significantly higher likelihood of sharing subsets of TCRs compared to subjects with differing HLAs and/or antigenic exposure history (35–38). Here we leverage this aspect of T-cell biology to identify sets of public TCR β s that are over-represented in subjects sharing HLAs. We expect these sets to be enriched for HLA-restricted TCR β s specific to common antigens and we use them to probe the functional nature of HLAs.

The T-cell repertoire is a rich source of information for understanding adaptive immunity (39, 40). The vast majority of TCRs are heterodimers composed of an α and β chain which together encode pHLA specificity. Our data consists of T-cell repertoires of TCR β sequences; the paired α chain is unknown. While any given TCR β chain may randomly pair with many α chains, the memory T cell compartment appears to be dominated by β chains that pair with a single α chain (41). This observation results from the fact that TCR-pHLA binding only occurs with very specific TCR $\alpha\beta$ combinations and these specific TCRs are significantly more likely to be sampled in a repertoire due to clonal expansion driven by antigen recognition. Furthermore, if the response is to a common antigen, the TCR β may be observed in multiple subjects (35, 37).

Given that memory T cells undergo strong clonal selection for specific TCR $\alpha\beta$ pairs and that some TCRs recur across individuals who share both HLA alleles and exposure history, we hypothesize that public TCR β s associated with specific HLAs are memory T cells targeting the same antigens presented by multiple individuals who share the appropriate restricting HLAs and pathogenic exposure history. While TCRs are inherently cross-reactive and capable of recognizing many distinct pMHCs (42), any given individual encounters only a small subset of the total antigenic space. Moreover, only a fraction of exposures are sufficiently common to elicit reproducible, public TCR responses across individuals. We will show consistent associations between specific TCR β s and HLAs, reflecting both underlying biological specificity and the ability of our approach to identify reproducible, public signals. These associations suggest that the β chains are typically paired with compatible—albeit unknown— α chains that preserve specificity to the same pHLAs. Thus, we will demonstrate that for the public TCR β sequences that we statistically associate with specific HLAs, shared pHLA specificity can be inferred from TCR β alone. Recent independent work has also explored statistical associations between TCRs and HLA genotypes using public repertoire data (43, 44), further supporting the value of population-scale immunosequencing.

Here we use high-throughput genetic sequencing (45, 46) of the T-cell repertoires of 4,144 subjects with HLA genotypes measured from direct-sequencing to identify $\sim 10^6$ public TCR β s that are statistically associated with HLA allotypes. While observing any given TCR β in a repertoire may be rare, the TCR β repertoire of an individual expressing a given HLA will almost always contain many TCR β s that we associate with that HLA. The TCR β s we associate to HLAs provide a new window into understanding the interaction between TCR and HLAs. The public nature of these TCR β s and their robust HLA associations permit highly accurate imputation of

class I and II HLA allotypes solely from TCR β repertoires allowing us to probe functional characteristics of HLAs with respect to their TCR β interactions.

Results

Identification of HLA-associated TCR β s

Our data consist of the sequenced T-cell repertoires of 4,144 subjects with HLAs genotyped via next generation sequencing (NGS) (47) (see [Supplementary Figure S1](#) for demographic distributions). The median number of unique T cells sequenced from each individual is $\sim 227,000$ and 90% of subjects have counts between $\sim 74,000$ and $\sim 610,000$. The majority of our samples are taken from healthy adults residing in the United States; $\sim 5\%$ and $\sim 20\%$ are Lyme and Covid positive, respectively. For a given HLA, we separate subjects into cases and controls defined as those with and without the HLA, respectively. Subjects expressing an HLA that is in the same p-group¹ as the HLA of interest are excluded from the control group, as such HLAs have identical amino acid sequences in the peptide binding region and thus may share TCR specificity. HLA-DP and -DQ are treated as heterodimers, with cases and controls defined by α - β pairs. The α chain of HLA-DR is invariant and thus we treat these HLAs as monomers, similar to class I HLAs. We randomly select a fixed 80% and 20% of the samples for training and validation, respectively, and evaluate model generalization on an entirely independent cohort from Kanazawa, Japan.

We identify sets of public HLA-associated TCR β s using a statistical approach that enforces the assumption that any TCR β is associated with at most one HLA allotype (we test this assumption below). This assumption enables us to disentangle the effects of linkage disequilibrium (LD) among HLA loci (48) that would otherwise result in a large number of spurious HLA-TCR associations ([Supplementary Figure S2A](#)). We use exact matching of the TCR β V-gene, J-gene and CDR3 to identify sequences that are over-represented in subjects with a given HLA allotype. Thus, our association of TCR β s with HLAs is agnostic to the specific amino acid sequence, it solely relies on it being observed in multiple repertoires.

Our assumption that TCRs typically associate with a single HLA allotype is important for resolving LD but is also supported by current biological evidence. While TCRs may be theoretically capable of broad cross-reactivity (42), there is limited experimental evidence that public TCRs functionally recognize multiple unrelated pHLAs. Rather, evidence suggests that functional cross-reactivity may be rare, with TCRs typically exhibiting specificity for particular epitopes and HLAs (49–51). Nonetheless, the true extent of functional cross-reactivity in physiologically relevant settings remains an open question. Our sequence identification procedure is primarily powered to detect public TCR β s associated with a single HLA and has limited

sensitivity to detect sequences with broad cross-reactivity, if such sequences exist.

The identification procedure works as follows (see Methods for precise details): for each HLA allotype, we first create a set of candidate HLA-associated TCR β s using a one-sided Fisher's Exact Test (FET) to identify TCR β s over-represented in cases. We adopt a pre-specified fiducial p-value threshold, p^* . For each unique candidate TCR β , we fit an L1-regularized Logistic Regression (L1LR) model which predicts the presence of that TCR β in subjects given their HLAs (represented as a binary vector of indicator variables). We tune the L1 hyperparameter λ to be the smallest value for which exactly one HLA parameter is non-zero. In other words, we determine which single HLA allotype best predicts the observed distribution of a given TCR β in the repertoires of our training sample. We test all TCR β s with p-values $< p^*$ and retain only TCR β s which associate most strongly with the HLA being modeled. As our interest here is primarily in the characteristics of titsets of HLA-associated TCR β s, we set p^* to a permissive value of $p^* = 10^{-4}$ and use the hold-out repertoires for validation. Note that due to exclusion of p-group matched HLAs a small number of TCR β s are assigned to multiple HLAs due to variations in the training data ([Supplementary Figure S2B](#)).

We associate $\sim 10^6$ TCR β s to specific HLAs, for a median of 2,400 TCR β s per HLA with $\sim 70\%$ of HLAs having a total number of associated TCR β s in the range of 1,600–5,000. To maintain consistency with other work associating TCR β s with disease (37, 38, 52, 53), we refer to these HLA-associated TCR β s as *enhanced sequences* (ES).

TCR specificity

Most enhanced sequences are specific to HLA allotypes

To validate the HLA allotype specificity of ESs, we compare their abundance in HLA cases as compared to controls in our holdout set (see [Figures 1A, B](#)). Overall, we find clear separation of cases and controls in the holdout data across all functional HLAs ([Supplementary Figure S3–S8](#)), highlighting the specificity of these ESs at the HLA allotype level and to the heterodimer for class II HLAs. However, there are notable exceptions.

We identify one HLA class I allotype pair ([Figures 1A, B](#)) and 11 class II pairs ([Table 1](#)) that appear to completely share TCR β s despite differing in their two-field designation. We refer to these as “degenerate” HLAs. For each of these degenerate pairs, ESs specific to one HLA are equally distributed among individuals expressing either HLA ([Figures 1C, D](#)). Ten of the twelve degenerate HLAs we identify have amino acid differences in a single position not in the peptide presentation and TCR binding domain and thus are in the same p-group.

To further validate our assumption that most of these TCR β s are specific to HLA allotypes, not allelic groups (with the exception of noted degenerate pairs), we use the L1LR method to assign TCR β s to either the allelic group (1-field) or the allotype (2-field). We restrict our analysis to class I HLA groups observed in >200

¹ http://hla.alleles.org/alleles/p_groups.html.

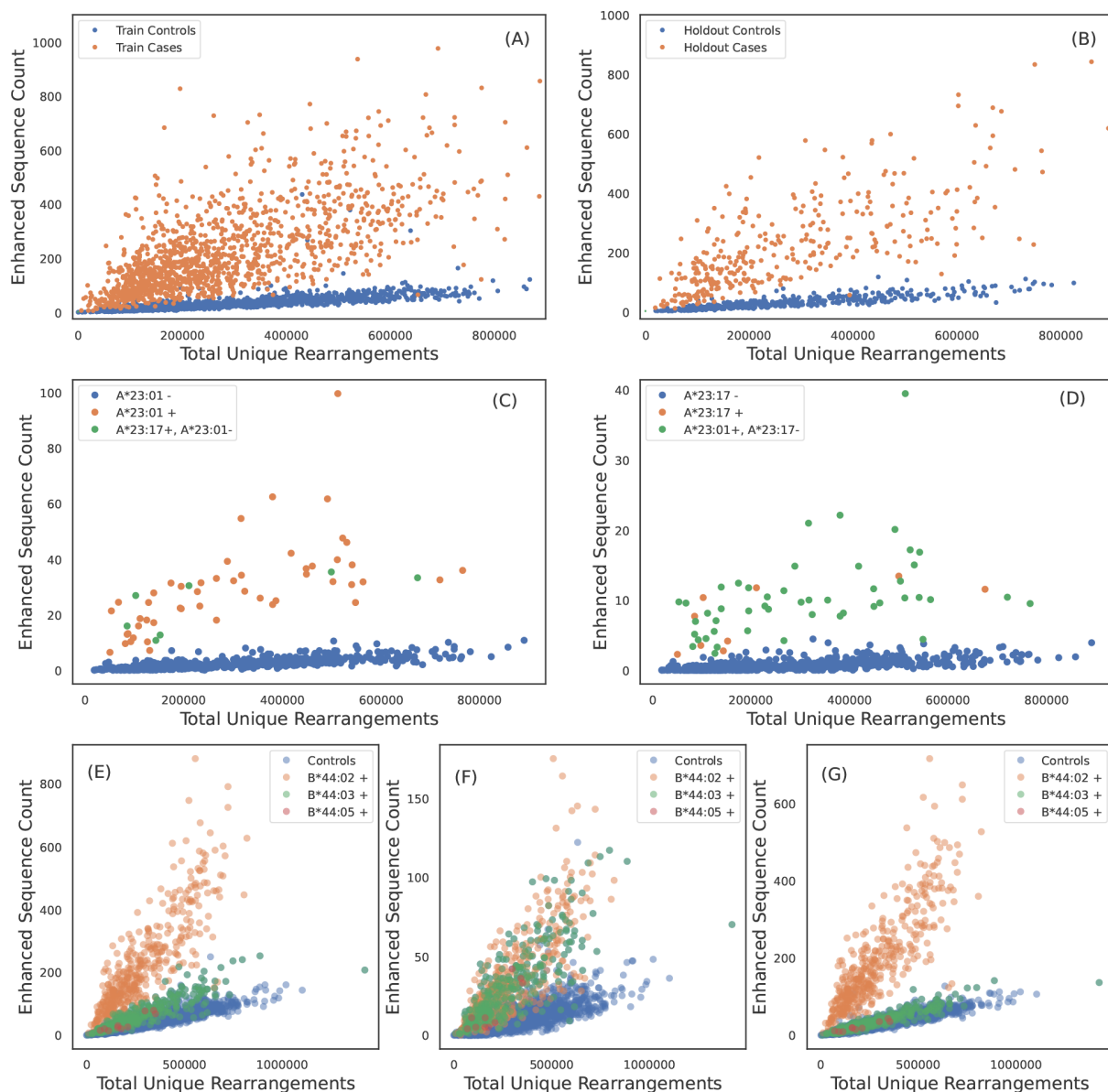


FIGURE 1

TCR β specificity and apparent TCR β sharing between some HLA allotypes. ESs for A*02:01 discriminate cases and controls in (A) train and (B) holdout samples. (C) A*23:01 ES counts observed in each sample as a function of sequencing depth. A*23:01 negative, A*23:17 positive subjects have counts consistent with A*23:01 positive subjects and vice-versa in (D). A*23:01 and A*23:17 appear to share the same TCR β specificity. (E) B*44:02 ES counts observed in each sample as a function of sequencing depth. ESs discriminate B*44:02 positive subjects from B*44:02 negative subjects. However, B*44:03 and B*44:05 positive subjects (green dots) have elevated counts as compared to controls (blue dots). (F) ES counts plotted against sequencing depth for the subset of ESs which associate more strongly with the B*44 group as compared to B*44:02 allotype (identified using the L1LR association method). The subset of ESs in (F) elevate all B*44 positive subjects equally suggesting the TCR β s in this ES subset are specific to the three allotypes in the group. (G) ES counts plotted against sequencing depth for the subset of ESs which associate only to B*44:02, i.e., this set excludes the ESs shown in (F). ESs plotted in (G) clearly separate B*44:02 positive subjects from B*44:02 negative subjects, including other allotypes in the B*44 group.

subjects with the most common allotype representing <70% of subjects. Among the six allelic groups tested, we find five have a negligible fraction (~1%) of ES associated to the group (A*30, A*33, A*68, B*15, B*35, C*07), indicating that the majority of HLA-associated TCR β s we identify are allotype specific.

We find one HLA group where a subset of TCR β s in the ES set appear to be specific to multiple allotypes in the group: B*44. We find that 10% of TCR β s originally identified as B*44:03-specific are assigned to the B*44 group (Figures 1E-G; similar conclusions are

reached when starting with the less prevalent B*44 allotypes). This set of B*44-specific TCR β s segregate all B*44 positive from negative individuals in the holdout (Figure 1F), while the remaining ~90% of TCR β s originally identified as B*44:03-specific separate B*44:03-expressing individuals from those who express B*44:02 or 44:05 (Figure 1G). These three B*44 allotypes differ in only two amino acids (residue 140 and/or 180), and B*44:02 and 44:03 are known to share a large fraction of their peptide repertoire and some of their TCR repertoire (54).

TABLE 1 List of degenerate HLAs.

Reported	Degenerate
A*23:01	A*23:17
DPA1*01:03+DPB1*03:01	DPA1*01:03+DPB1*104:01
DPA1*02:01+DPB1*17:01	DPA1*02:01+DPB1*131:01
DQA1*02:01+DQB1*02:02	DQA1*02:01+DQB1*02:01
DQA1*03:03+DQB1*03:01	DQA1*03:01+DQB1*03:01
DQA1*03:01+DQB1*03:02	DQA1*03:03+DQB1*03:02
DQA1*03:03+DQB1*02:02	DQA1*03:03+DQB1*02:01
DQA1*05:05+DQB1*03:01	DQA1*05:01+DQB1*03:01,
DQA1*05:05+DQB1*03:01	DQA1*05:05+DQB1*03:19
DQA1*05:01+DQB1*02:01	DQA1*05:05+DQB1*02:01
DRB4*01:03	DRB4*01:01
DRB5*02:02	DRB5*02:21

We are unable to distinguish between these sets of HLAs using TCR β based typing as they have identical TCR β specificities. When modeling degenerate HLAs, we only report the most commonly occurring HLA in the set.

Degenerate HLAs that completely share their TCR repertoire typically differ in one amino acid outside the binding domain. Similarly, B*44 has two polymorphic positions outside the binding domain and displays a high degree of sharing. On the other hand, groups in which we observe no TCR β sharing tend to have one or more amino acid differences in the binding domain. Thus, the degree of sharing we observe appears to be correlated to the number of differing residues and the position at which the differences occur. Taken together, these results show that many TCR β s (indeed, the vast majority of those identified here) are specific to distinct HLA allotypes, regardless of shared peptide repertoire or binding domain similarity. Thus, many characteristics of TCR-pHLA interactions differ among highly related HLA allotypes. As expected, no such specificity was observed among HLA allotypes that differ only in synonymous substitutions (ie, at 3- and 4- digit resolution; [Supplementary Figure S9](#)). We note that to mitigate the effects of linkage disequilibrium, our algorithm for identifying ESs assumes TCRs are specific to individual HLA allotypes (see methods). Thus, our methodology is not designed to identify an unbiased set of TCRs specific to multiple HLA allotypes. Our identification of a small fraction of ESs with specificity to multiple HLA allotypes warrants a systematic investigation which is beyond the scope of this work.

Most TCRs are specific to class II heterodimers, not subunits

Functional class II HLAs are stable heterodimers with both the α and β chain contacting the peptide. As such, we expect class II HLA-associated TCRs to be specific to the heterodimer and not the protein subunits (i.e. the α or β chains individually). To directly test this hypothesis, we use the LILR method to determine if a TCR β is more strongly enriched among individuals expressing both the α and β chains or individuals expressing only one or the other

subunit. Across 37 heterodimers, we find that ~ 146000 (70%), ~ 20500 (10%) and ~ 43000 (20%) of ESs are most strongly associated with the heterodimer, alpha and beta subunits, respectively. We note that not all 37 heterodimers exhibit single-chain specificity (see below). Thus, the vast majority of class II associated TCR β s appear to be specific to the combined $\alpha - \beta$ chains. This finding is bolstered by the fact that the ES sets we derive discriminate HLA-DP and -DQ heterodimers and not individual subunits ([Figure 2](#); see also [Supplementary Figures S6-7](#), which show ES distributions for all heterodimers).

We find that only a subset of ESs and HLAs exhibit exceptions to heterodimeric specificity. For example, some TCR β s appear to be associated to all heterodimers composed of the DQB1*05:01 subunit, indicating many of these TCR β s are associated with the subunit itself ([Figure 2](#)). DQB1*05:01 is the clearest example of TCR specificity to a subunit. However, we observe such subunit specificity across multiple subunits: DPB1*01:01, DQB1*02:01 (DQB1*02:02)², DQB1*03:01, DQB1*05:01, DQB1*06:03, DQA1*03:01 (DQA1*03:03) and DQA1*05:01 (DQA1*05:05). TCR specificity to subunits appears to be more common for the HLA β chain and to the HLA-DQ locus, though we observe single-chain specificity in both the α and β chains of HLA-DQ and a β chain of HLA-DP. We note that our identification is likely not exhaustive as many heterodimers lack enough diversity in one or both subunits to statistically associate TCRs independent of the heterodimer.

Based on a limited set of solved structures, a conserved binding pattern has been proposed for class II TCR-pHLA interactions such that TCR α contacts the α helix of the HLA β chain and TCR β contacts the α helix of the HLA α chain ([55](#)). Thus, the various specificity patterns we observe may reflect TCR interactions strongly mediated by peptides and not the direct interactions between TCRs and HLAs which may be conserved. Future analyses based on larger sets of solved or perhaps *in silico* generated structures of class II TCR-pHLA complexes will be informative for exploring the structural basis of the various specificity patterns we identify.

TCR breadth is proportional to zygosity

HLA homozygosity has been epidemiologically linked to poor clinical prognosis in the context of both chronic HIV infection ([56](#)) and cancer checkpoint-inhibitor immunotherapy ([57](#)), possibly due to the reduced size of the HLA-restricted peptide repertoire available for T-cell recognition. This reduced size of the antigen repertoire, coupled with higher relative surface concentration of pHLAs associated with the homozygous protein, may directly impact the TCR repertoire by increasing the probability of clonal expansion of T cells expressing cognate TCRs. Consistent with this hypothesis, we find that the distribution of ES counts is elevated among homozygous individuals ([Figures 3A-F](#)). Across all HLAs,

2 Degenerate subunit pairs are indicated by one of the degenerate subunits shown in parenthesis.

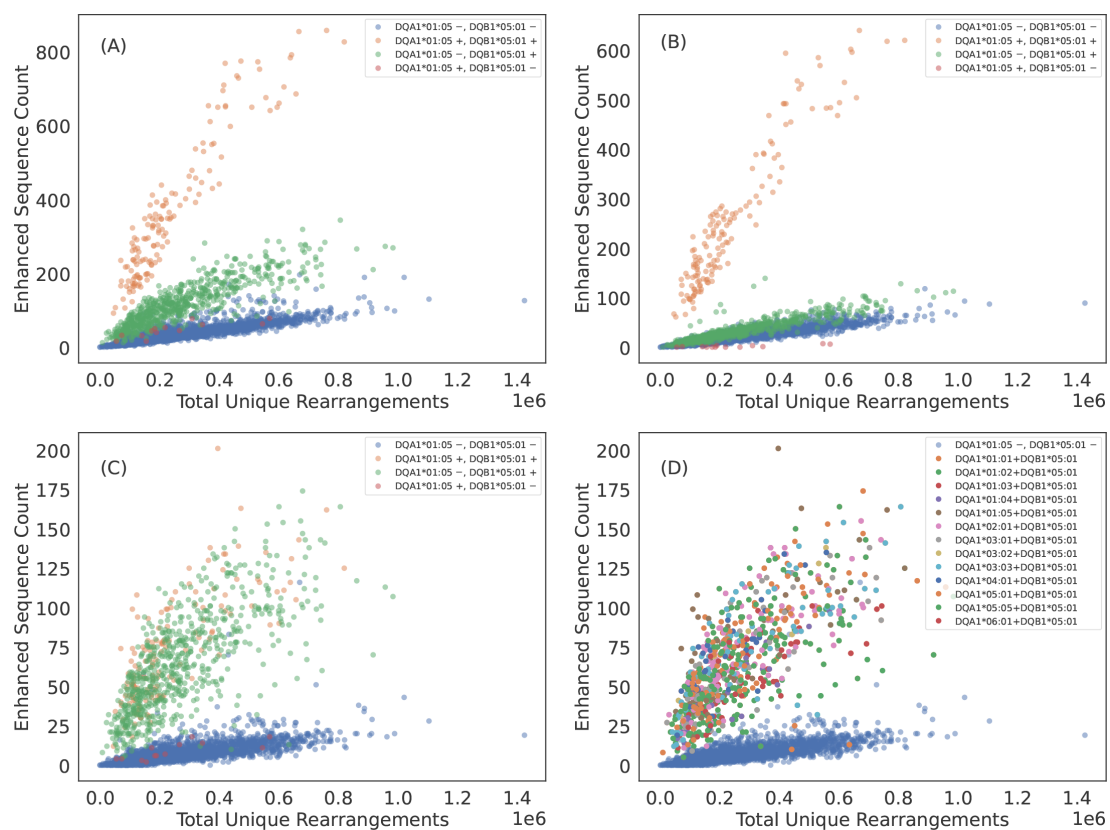


FIGURE 2

Some TCR β s appear to be specific to the class II HLA subunit rather than heterodimer. (A) DQA1*01:05+DQB1*05:01 ES counts observed in each sample as a function of sequencing depth. The ESs separate DQA1*01:05+DQB1*05:01 positive subjects from DQA1*01:05+DQB1*05:01 negative subjects. However, subjects who have the DQB subunit appear elevated above the control population of subjects with neither subunit. (B) ES count plotted against sequencing depth for the subset of ESs which associate most strongly with DQA1*01:05+DQB1*05:01 heterodimer via the L1LR method. A majority of TCR β s which make up the ES set for DQA1*01:05+DQB1*05:01 appear to be specific to the heterodimer. (C) ES count plotted against sequencing depth for the subset of ESs which associate most strongly to subunit DQB1*05:01 via the L1LR method. A small fraction of TCR β s which make up the ES set for DQA1*01:05+DQB1*05:01 appear to be specific only to the β chain subunit. (D) Same as (C) but color-coding subjects with DQB1*05:01 by their various α chain pairings. This subset of TCR β s appear to be specific to all possible heterodimeric combinations which include DQB1*05:01, thus suggesting specificity solely to the β chain subunit. Results are shown for train sample and are consistent with holdout sample.

the distribution of ESs is (on average) about one standard deviation higher for homozygous as compared to heterozygous individuals (Figures 3G, H). Notably, the ES distribution is an additional standard deviation higher among individuals homozygous at both the α and β locus for HLA-DP or -DQ (“double homozygous”, Figure 3H).

While these results are consistent with increased relative antigen abundance increasing clonal expansion of associated T cells, an alternative hypothesis is that a decrease in the relative surface expression or decreased antigenic diversity in homozygous subjects results in less crowding out by other HLAs or TCRs, respectively. This crowding hypothesis implies an increased breadth across multiple loci within a given class for homozygous subjects. We test this hypothesis by examining whether a homozygous subject at one locus has higher breadth at another class-matched loci. For example, if HLA and/or TCRs crowd each other, we expect that a subject homozygous for HLA-A will also have a higher average breadth in their HLA-B and/or HLA-C response due to lower diversity at the class I loci. We do not find

evidence of such an effect (see Supplementary Figure S10) and thus conclude that crowding out by HLAs and/or TCRs may not be the primary driver of increased breadth unless it is restricted to a particular locus.

Taken together, these results suggest homozygosity at a particular locus increases the breadth of the T-cell response against peptides presented by that HLA, possibly through increased surface expression and antigen presentation.

Imputing HLA genotype from TCR repertoires

The clear separation of ES counts in HLA cases versus controls implies that HLA allotypes can be easily imputed from HLA-associated TCR β s alone. To this end, we fit a simple logistic regression model for each HLA allotype observed in at least 30 training samples (representing $\sim 1\%$ expression frequency), predicting whether an individual expresses that HLA as a

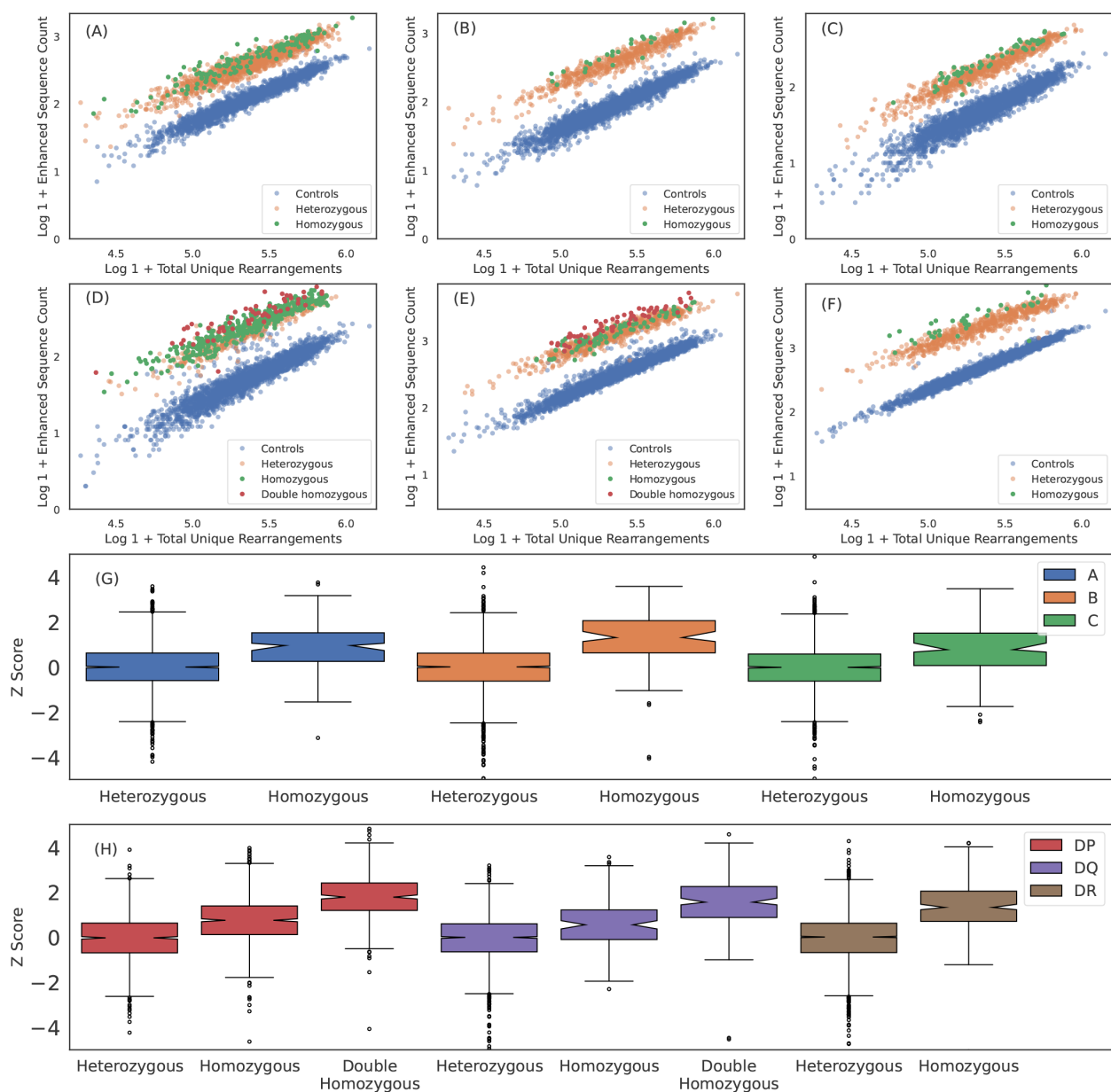


FIGURE 3

Breadth of T cell response is proportional to HLA zygosity. (A) A*02:01, (B) B*07:02, (C) C*04:01, (D) DPA1*01:03+DPB1*02:01, (E) DQA1*01:02+DOB1*06:02 and (F) DRB1*07:01 ESs observed in each sample plotted against total number of unique rearrangements. HLA negative, heterozygous positive subjects are shown in blue, orange and green, respectively. For HLA-DP and DQ, double homozygous subjects are shown in red. The breadth of the ES response appears to be correlated with homozygosity across all six loci. We quantify the increased breadth resulting from homozygosity by fitting the mean and standard deviation of the ES counts in heterozygous cases for each HLA as a function of sequencing depth and then calculating the z-score for all subjects and all well-represented HLAs. We aggregate z-score distributions per-loci for (G) class I and (H) class II HLAs. Results are shown for train sample and are consistent with holdout sample.

function of observed ES and total-unique-rearrangement (log) counts (see Figures 4A, B for representative examples, Supplementary Figures S3-S8 for all HLAs tested). The number of subjects with each HLA allotype based on HLA genotyping and the model performance for each imputed allotype are provided in Table 2 (class I) and Table 3 (class II); model performance is also shown in Figure 4C.

Over all the imputation accuracy is extremely high with area under the receiver operating characteristic curve (AUC-ROC)

scores of ≥ 0.9 for all but 3 of the 120 HLA-A, -B, and -DR allotypes modeled. This accuracy highlights the specificity of HLA ESs, indicating HLAs at these loci can be accurately imputed from immunosequencing alone. Furthermore, this accuracy confirms that HLA β alone is typically sufficient to identify shared antigen specificity of public T cells. Among class I HLAs, HLA-C allotypes are a notable outlier with relatively lower classification performance even among models with a significant number of positive training examples (Figure 4D). We hypothesize that this reduced

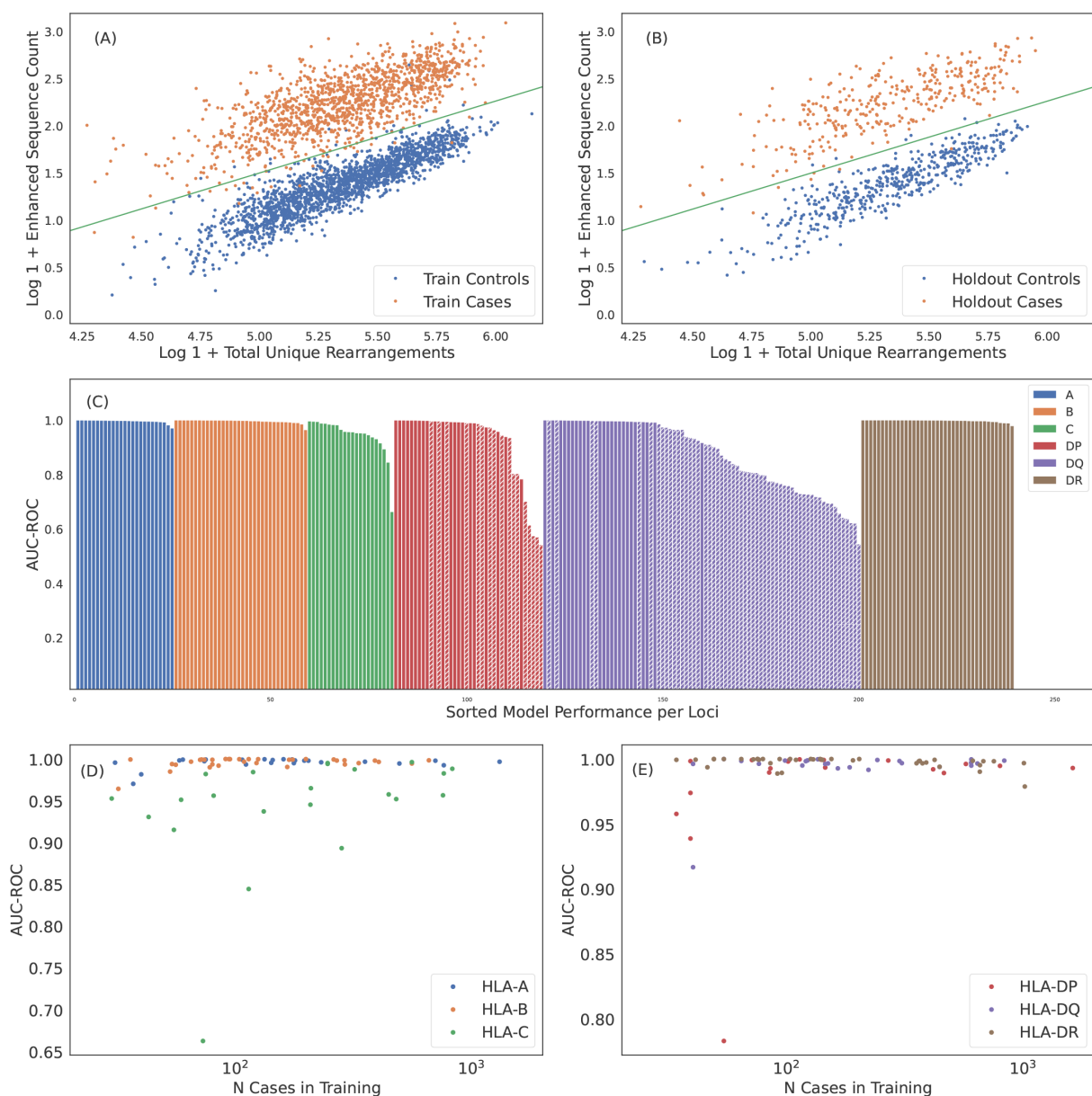


FIGURE 4

Robust predictions of hundreds of HLAs solely based on public T cells observed in the repertoire. ES counts for A*02:01 cases and controls plotted against total number of unique rearrangements for (A) train and (B) holdout samples. The green line indicates the call threshold. (C) AUC-ROC of HLA models across each loci sorted by performance. The hashed bars for HLA-DP and -DQ indicate heterodimer resulting from subunits combinations in linkage equilibrium suggesting trans-complementation. We show below that many of these HLAs are likely non-functional. Performance of (D) class I and (E) class II HLA models as a function of the number of training samples color-coded by loci. Performance correlates with the number of training samples, as expected. HLAs shown by hashed bars in (C) are excluded in (E). In (C-E) we show performance derived from 5-fold cross validation (CV) of the train sample. CV performance is consistent with holdout performance.

performance is due in part to the $\sim 10\times$ -lower surface expression of HLA-C compared to allotypes expressed by other class I genes (58).

Model performance among HLA-DP and -DQ heterodimers is substantially more variable than performance at the other loci. We find that 8 of the 30 HLA-DP and 37 of 81 of HLA-DQ fail to achieve AUC-ROC scores of > 0.9 even among heterodimers expressed in a high number of individuals in our training population (Figure 4E). Notably, HLA-DQ and -DP are the only

loci we study with polymorphic α chains, suggesting that heterodimer incompatibility may explain the lack of associated TCRs and the corresponding inability to impute expression of these heterodimers. We explore this issue in the next section.

We assess the generalizability of our imputation models and rule out overfitting or data leakage by evaluating model performance on an entirely independent, out-of-distribution cohort of 136 individuals from Kanazawa, Japan. Subjects in this cohort have HLAs genotyped by direct sequencing and

none are included in our training or internal holdout sets. This population differs both genetically and geographically from the predominantly U.S.-based training cohort, thus providing a stringent test of the model's generalizability to out-of-distribution data. For this analysis, we use a set of 135 HLAs across all six loci that yield models with cross-validation precision > 0.9 , recall > 0.8 and ≥ 30 positive training cases.

We apply the 135 HLA models to the Japanese cohort, resulting in an average of 8.5 imputations per individual, compared to 11.0 per individual in our internal validation cohort. This difference is expected and reflects the shift in HLA allele frequencies between the Japanese and U.S.-based cohorts. Despite these differences, we observe comparable model performance for the subset of overlapping HLAs. The overall F1 score³ for the Japanese cohort, aggregated across all imputations, is statistically consistent ($< 3\sigma$) to that observed in the internal validation cohort (0.925 ± 0.006 and 0.940 ± 0.002 , respectively). These results confirm that model performance generalizes across genetic and geographic backgrounds and that the high accuracy we observe is not a consequence of data leakage, but rather reflects a strong HLA-specific signal in TCR repertoires.

Poor-performing class II models are trans-complemented, non-functional HLAs

HLAs are inherited as a haplotype, such that one full set is inherited on a single chromosome from each parent (59). The α and β chains of HLA-DP and -DQ⁴ pair after synthesis, yielding a phenotype of up to four unique heterodimers in each individual: two formed in *cis*, where both subunits are encoded on the same chromosome, and two in *trans*, where the subunits are encoded on opposite chromosomes. Given the high degree of polymorphism observed in HLA-DP and -DQ subunits, it is perhaps unsurprising that some pairs of α and β chains do not form stable heterodimers (60–62). Based on structural and sequence analysis of HLA-DQ, Tollefsen et al. (62) propose specific group pairings that likely form stable heterodimers (though they note a small number of exceptions to the pairing rules).

In the context of the present study, the proposed existence of incompatible (and thus non-functional) α and β chains implies two testable hypotheses: (1) that co-inheritance of incompatible pairs on the same chromosome will be under strong negative selection (62); and (2) that incompatible pairs are unable to elicit a T-cell response and thus will not be associated with any public TCR β s.

The first hypothesis implies that co-expression of incompatible pairs almost always results from *trans*-complementation; as such, incompatible pairs will be in linkage equilibrium since they are not

co-inherited. Conversely, *cis*-complementation should necessarily result in functional pairs thus all pairs forming from subunits in linkage disequilibrium should be functional. For two subunits α and β , with respective expression frequencies f_α and f_β and co-expression frequency $f_{\alpha+\beta}$, linkage equilibrium results in an expected co-expression frequency of approximately $E_{LE}[f_{\alpha+\beta}] = 2f_\alpha f_\beta$ (63). Overall, we find that a majority (75 of 119) HLA-DP and -DQ $\alpha + \beta$ pairs appear to be in linkage equilibrium, i.e., not co-inherited (Figure 5A; only subunits comprising heterodimers observed in at least 30 individuals are considered). Notably, all of the pairs that Tollefsen et al. (62) propose to be incompatible are in equilibrium (Figure 5B).

Given the tight linkage between genes encoding subunits, any pair expressed in *cis* in at least one individual is likely to be in LD in a large cohort. Consistent with this hypothesis, every individual heterozygous at both the α and β locus has at least two of four α/β pairs in LD; conversely, no individual should have more than two α/β pairs in equilibrium. We confirm that no subject in our cohort has more than two heterodimers formed from subunits in linkage equilibrium, as expected. Taken together, we conclude that LD is a strong proxy for *cis*-complementarity, and that, given strong selection pressure, incompatible pairs are only expressed in *trans*.

If incompatible pairs are truly non-functional (with respect to antigen presentation and T-cell recognition), then there should not be any TCR β s that are specific to such pairs. As an example, we are unable to identify distinguishing TCR β s for DQA1*01:02+DQB1*03:01, which both violates the Tollefsen et al. (62) pairing rules and is in linkage equilibrium (Figure 5C). Moreover, all of our high-frequency, poor-performing HLA-DP and -DQ imputation models are for heterodimers forming from subunits that are in linkage equilibrium, and thus are almost certainly expressed in *trans* (Figure 5D; see also Figure 3C). Moreover, almost all pairs that violate the Tollefsen et al. (62) pairing rules have lower-than-expected imputation performance (Figure 5D).

Notably, while heterodimer incompatibility implies both linkage equilibrium and poor model performance, not all heterodimers forming from subunits in linkage equilibrium result in poor performing models. Thus, *trans*-complementation may yield heterodimers which can drive a T-cell response resulting in identifiable TCR β s and a high-accuracy imputation model (Figure 5D).

Using these results on model performance and gene linkage, we can extend the Tollefsen et al. (62) pairing rules to the HLA-DP locus: DPA1*02 appears to form unstable heterodimers when paired with DPB1*02 and DPB1*04; DPA1*01 is apparently unrestricted, forming stable heterodimers with subunits from all DPB1 groups in our sample.

HLA associated sequences are memory T cells responding to common antigens

The T-cell repertoire consists of a mixture of naive and memory T cells. In principle, HLA-restricted antigen presentation will bias both thymic selection and clonal expansion, and thus HLA-specific signatures may exist in both compartments. However, our statistical approach to identifying HLA-specific public TCR β s likely favors identification of TCRs from memory T-cells.

³ F1 score is the harmonic mean of precision and recall and a useful metric of classification accuracy when there is large class imbalance such as the one we have for HLA imputation.

⁴ We note that HLA-DR heterodimers are not subject to such pairing because the α chain is nearly invariable and thus these heterodimers behave similarly to class I HLAs.

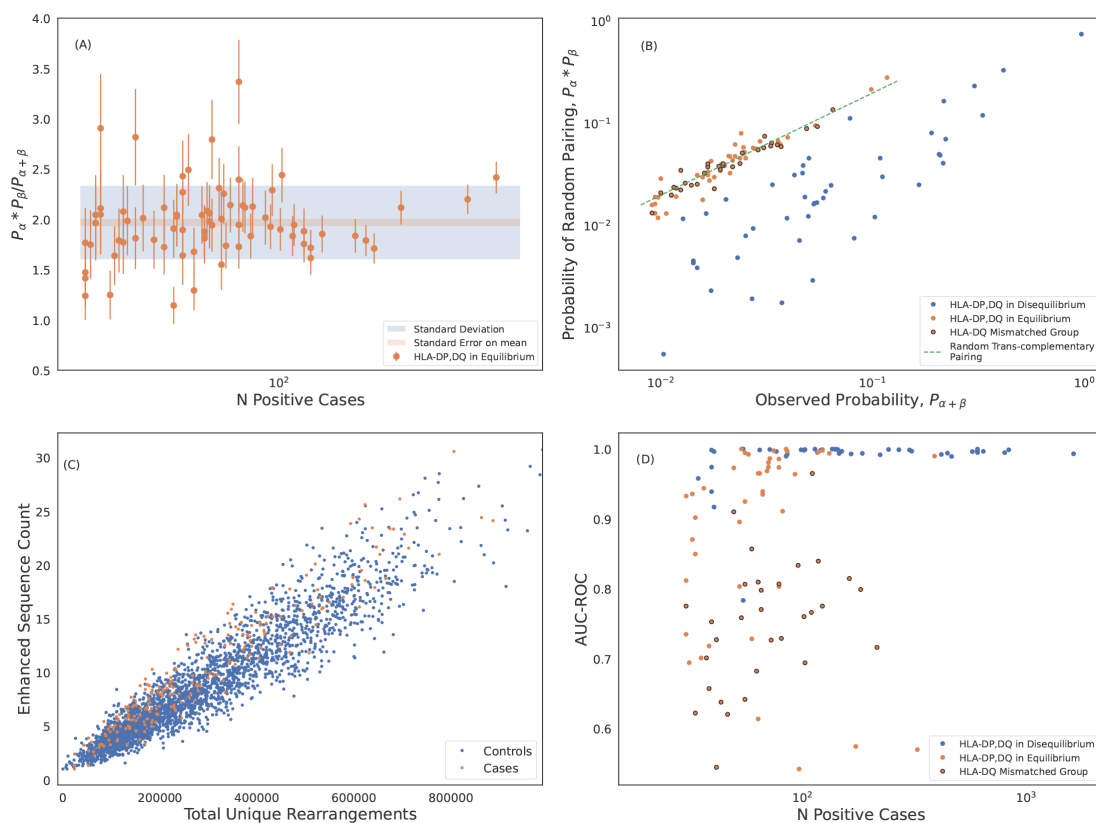


FIGURE 5

Trans-complemented heterodimers may not form stable HLAs. **(A)** Statistical analysis of the marginal-to-joint probability ratio of heterodimers forming from subunits in linkage equilibrium. We identify these heterodimers as those which differ $> 5\sigma$ from ratio of 2 expected for genes in linkage equilibrium. We measure an average probability ratio of 1.97 ± 0.04 . Error bars are derived from propagating poisson uncertainties. **(B)** Expected probability of randomly pairing α and β chains of DP and DQ heterodimers plotted against the observed joint probabilities. Probabilities are calculated from the normalized inverse frequencies. The probability of randomly pairing is calculated as the product of the observed marginal probabilities of the subunits. The dashed green line is the expected correlation for random trans-complementary pairing. We identify heterodimers formed from subunits in linkage equilibrium (shown in orange) as in **(A)**. DQ heterodimers formed from pairing of mismatched groups as defined by (62) are shown with black circles. These mismatched group heterodimers cluster around the dashed green line indicating random trans-complementary pairing. **(C)** DQA1*01:02+DQB1*03:01 is an example of an HLA where we are unable to identify any ESs that separate cases and controls. **(D)** Model performance of HLA-DP and -DQ models including HLAs forming from subunits in linkage equilibrium. Here we show performance derived from 5-fold cross validation (CV) of the train sample which is consistent with holdout performance.

We investigate characteristics of HLA associated sequences by comparing the distribution of ESs observed in the memory and naive compartments sequenced from 45 individuals who were not included in the original study. For each subject, we sequence five separate repertoires: the memory and naive compartments of CD8⁺ and CD4⁺ T cells, respectively, and the unsorted repertoire. As sequence-based typing was unavailable for these individuals, we treat the imputed HLAs from the the unsorted repertoire as ground truth (limiting to 90 models with >0.95 precision and recall). For each repertoire and each possible HLA, we compute the weighted breadth of the HLA-specific ESs observed in the repertoire. Across all 45 unsorted repertoires, the weighted breadth of both class I and class II ESs is substantially higher for the HLAs an individual expresses compared to those they do not express (Figure 6A; compare HLA-positive to -negative).

Within the sorted compartments, a striking pattern emerges: in the naive compartments, there is little difference in the breadth of ES specific for an individual's expressed HLAs compared to

background (Figures 6C, E), while in the memory compartments, ESs specific to the individuals' expressed HLAs have far higher breadth (Figures 6B, D). Moreover, within the memory compartment, CD8⁺ cells are closely linked to high relative breadth of class I HLA ESs, while CD4⁺ cells are closely linked to high relative breadth of class II HLA ESs. This result provides further confirmation that ESs are correctly mapped to HLAs despite the challenges of HLA LD.

The centrality of the memory compartment in driving our HLA imputation signal raises several additional hypotheses. The first is that clonal expansion increases the likelihood of detection for ESs. This increased likelihood implies that, while ESs are frequently observed in individuals without the associated HLA, they will tend to be at higher repertoire frequency among individuals who do express the HLA. Indeed, we observe a notable increase in the distribution of clonal frequency among cases compared controls (Figure 7A).

The second hypothesis is that clonal expansion results from antigen exposure, which will also result in polyclonal expansion of T

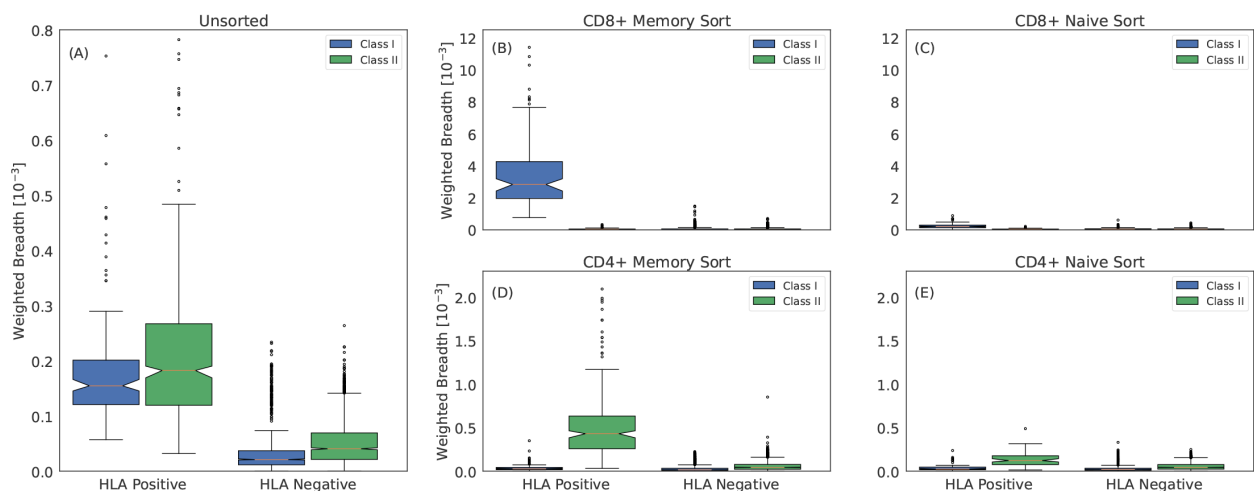


FIGURE 6

Class I and class II associated ESs are TCR β s from CD8 $^{+}$ and CD4 $^{+}$ memory T cells, respectively. (A) Weighted breadth of 90 HLA associated ESs measured in 45 subjects with imputed HLAs. The breadth is sorted by whether the subject has the HLA or not and is then aggregated across all subjects and HLAs. To generate a comparable breadth across HLAs, which have varying number of ESs, we measure the median breadth in controls for all 90 HLAs, which we rescale to a mean of 1. We then normalize the breadth of any given HLA by this value. Breadth measured in (B) CD8 $^{+}$ memory sorted, (C) CD8 $^{+}$ naive sorted, (D) CD4 $^{+}$ memory sorted and (E) CD4 $^{+}$ naive sorted repertoires. We measure slightly elevated breadth in the CD8 $^{+}$ and CD4 $^{+}$ naive sorted repertoires for class I and class II HLAs, respectively. This elevation may be due to surface markers not perfectly discriminating naive and memory compartments or to a weak HLA specific signal due to the HLA interactions required for maintaining homeostatic equilibrium of naive T cells (75).

cells that express distinct TCRs that all respond to the same antigen. An extreme form of polyclonal expansion is *convergent recombination*, in which multiple, distinct TCR DNA sequences encode identical amino acid sequences. Consistent with this hypothesis, the per-ES convergent recombination rates are higher when ESs are observed in cases as compared to controls (Figure 7B; conditional on the ES being present in the repertoire).

The publicity of HLA ESs combined with their apparent antigen-specific clonal expansion suggests that these sequences are responding to antigens which are common in the human population (64). We test this hypothesis in the context of SARS-CoV-2 exposure. Because of its novel nature, SARS-CoV-2 is especially well-suited for this task as we are able to confidently assign Covid-19 negativity to samples collected before 2020. In our training sample for deriving HLA restricted sequences, 694 of the 4,144 subjects (~20%) are Covid-19 positive based on PCR labels; the remaining samples are collected before 2020 and thus Covid-19 negative. Because SARS-CoV-2 exposure is relatively high in our training sample, we expect that some of the HLA-specific ESs we identify are SARS-CoV-2 specific.

We identify 6866 SARS-CoV-2-specific ESs using a set of 1523 SARS-CoV-2 PCR-positive samples that have no overlap with the typed HLA training samples and 4386 controls (1008 controls overlap with the typed HLA samples). A high proportion (80% of the most confident SARS-CoV-2 specific ESs) are also HLA-specific ESs (Figure 7C, blue line). This overlap in ESs results from the fact that ~20% of our typed HLA training samples are SARS-CoV-2 positive, i.e., SARS-CoV-2 specific ESs are predictive of HLAs as well. In contrast, if we define an alternative set of HLA-ESs using only 3,450 repertoires which were sampled prior to 2020, we find

very little overlap (Figure 7C, orange line). The limited overlap we do observe may be due to cross-reactivity to homologous epitopes from other coronaviruses and/or false positive sequences in the SARS-CoV-2/HLA ES set.

The intersection of SARS-CoV-2 specific ES set with the and HLA specific ES sets yields 1,880 TCR β s (using a threshold of $p^* = p < 10^{-4}$) that are associated with a particular HLA in the context of SARS-CoV-2 infection. Thus, these ESs are almost certainly specific to commonly targeted SARS-CoV-2 T-cell epitopes. To confirm this specificity, we impute HLAs for samples in the SARS-CoV-2 training data⁵ and then compute the fraction of HLA-associated SARS-CoV-2 ES observed in their repertoire. For “HLA +” and “HLA -” subjects we only count sequences if the subject has or does not have the HLA which the SARS-CoV-2 ES is associated with, respectively. Notably, both class-I and class-II associated TCR β s are far more commonly observed in individuals with the restricting HLA and known SARS-CoV-2 infection (Figure 7D), thus confirming the HLA- and pathogen-specificity of these TCR β s and further demonstrating that while TCR β s may be cross-reactive to many distinct antigens, the likelihood of cross-reactivity *in vivo* is low (65).

We conclude that, taken together, these results demonstrate that the majority of HLA-specific ESs are the result of clonal expansion of T cells responding to common antigens, thus establishing an immunological foundation for HLA imputation from TCR β repertoires.

⁵ These data do not have sequencing based HLA typing.

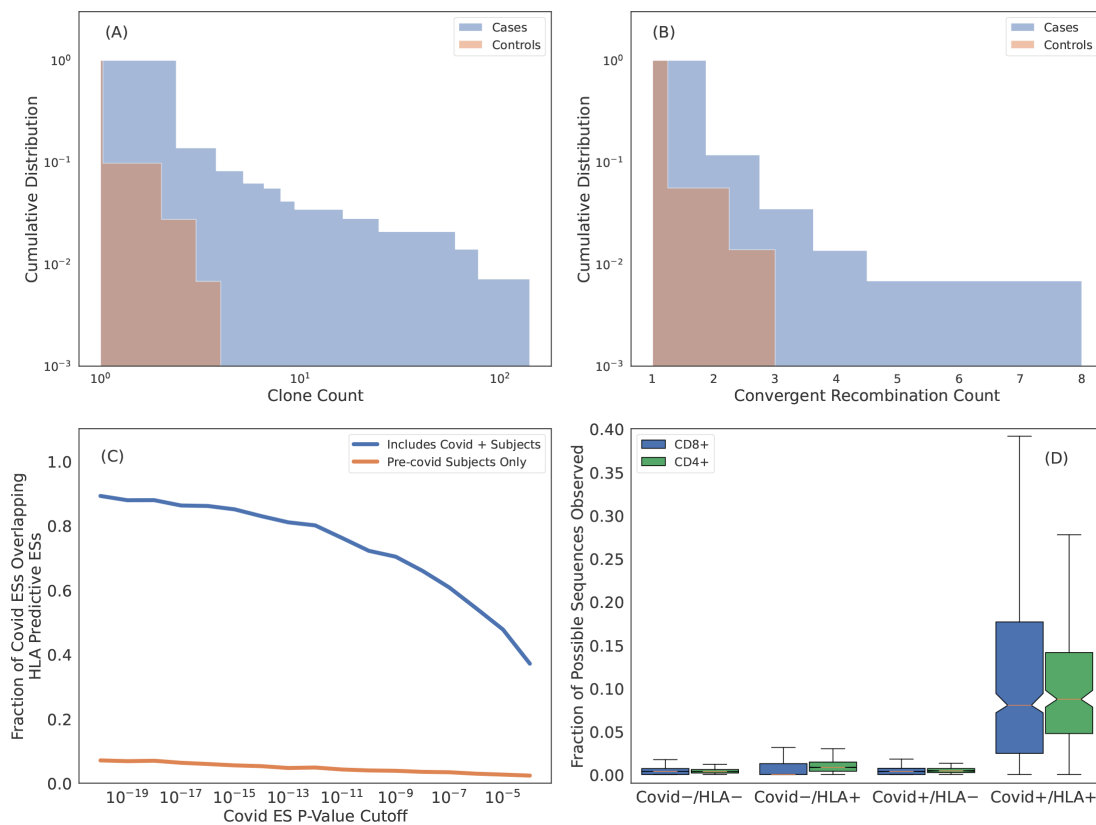


FIGURE 7

Many HLA ESs are T cells responding to common pathogens. **(A)** Cumulative distribution of clone count per unique rearrangement as measured in cases (blue) and controls (orange). ESs are generally more expanded when they are observed in subjects with the restricting HLA. **(B)** Cumulative distribution of the number of unique rearrangements mapping to an ES, i.e., convergent recombination count. ESs show more convergent recombination when observed in cases (blue) as compared to controls (orange). **(C)** Intersection of SARS-CoV-2 specific ESs derived via a FET on samples with PCR labels and HLA ES sets. Blue curves are HLA ES sets derived from samples which include Covid-19 positive subjects and orange curves are from samples with all Covid-19 positive subjects removed. **(D)** Fraction of SARS-CoV-2 ESs observed in subjects relative to the total fraction possible given the ES HLA association determined from intersecting the SARS-CoV-2 ESs with the HLA ES sets. Here we impute HLAs using our models. For “HLA +” we only count sequences associated with HLAs inferred for the subject and for “HLA -” we only count sequences associated with HLAs that are not inferred for the subject.

Discussion

The T-cell repertoire of any individual consists of $\sim 10^8$ unique TCRs out of an estimated $\sim 10^{16}$ possibilities (66, 67), of which only $10^5 - 10^6$ TCRs may be practically sampled from any single repertoire using current techniques. Given the enormous diversity of possible TCRs, naively, little overlap in the TCR repertoire of different subjects may be expected. However, several factors significantly increase the likelihood of observing public TCRs: (1) the probability distribution of TCRs generated via V(D)J recombination is non-uniform and spans ~ 25 orders of magnitude (68), such that higher generation probability TCRs are more commonly observed in the naive repertoire; (2) antigen-experienced T cells clonally expand and are thus more likely to be observed in a repertoire; and (3) the immune response is focused on only a few immunogenic epitopes per HLA out of the many possible derived from any given antigenic exposure, an effect called *immunodominance* (69). As a consequence, the likelihood of observing the same TCRs in the repertoires of multiple subjects

with shared antigenic exposure and appropriate restricting HLA is significantly higher than naively expected. Studies have shown that public TCRs can be identified that permit sensitive and specific diagnosis of individuals with past SARS-CoV-2 infection (37) and Lyme disease (38) as well as to determine who is seropositive for Cytomegalovirus (35). Given that these antigens are presented by HLAs, it is perhaps unsurprising that these public TCRs are also specific to the HLA context (35, 36). Here we leverage this public fingerprint of TCRs to identify HLA associated sequences, allowing us to impute HLA types with extremely high accuracy and opening a new window into functional characteristics of HLAs.

We identify HLA specific TCR using only the TCR β repertoire, even though the specificity of pHLA is to TCR $\alpha\beta$. Enhanced sequences elicit an immune response only in subjects with the appropriate shared restricting HLA and pathogenic exposure history (Figure 7D), meaning, enhanced sequences are responding to the same pHLAs in different subjects. This result implies that an observed TCR β shared among multiple subjects in a given HLA context is very likely paired with one or a very small set of compatible TCR α s. In this respect, TCR β repertoires are not

unique; analysis of TCR α repertoires would lead to similar results and conclusions as presented here⁶. Thus, our ability to identify HLA specific TCRs using only a single chain of the TCR heterodimer is not a characteristic of the TCR β s but rather a consequence of our procedure for identifying enhanced sequences and the relevant underlying biological processes.

A key finding of this work is that TCRs are typically specific to HLA allotypes (two field resolution) and to class II heterodimers encoded by both the α and β chains, although there are some notable exceptions of TCR β s that are specific to HLA groups (one field resolution) or to either the α or β subunits of HLA-DP and -DQ. While most HLAs elicit a strong and diverse public TCR response, others elicit little or no response, consistent with these HLAs deriving from incompatible α and β subunits and providing further support for the widespread prevalence of non-functional class II heterodimeric pairs resulting from trans-complementation. Furthermore, we find that the breadth of an HLA-specific TCR β response is larger among individuals expressing two (or more) copies of the HLA, suggesting a dose-dependent effect of antigenic exposure on the diversity of expanded T-cell clones. These insights highlight the exquisite specificity of public TCR β s and demonstrate the potential of population-level TCR analysis for probing functional aspects of the immune system.

We show that class I and class II HLA-associated TCR β s are found on CD8⁺ and CD4⁺ memory T cells, respectively, which is consistent with their publicity resulting from clonal expansion in response to antigenic-exposure. The public nature of these TCR β s suggests that they are likely specific to peptides derived from common pathogens, vaccines and conserved endogenous-antigens. Consistent with this hypothesis, ~20% of subjects in our training sample are covid positive and we identify a consequently large fraction of SARS-CoV-2-specific TCR β s as HLA associated. Notably, this overlap provides probable pathogenic and HLA assignments to these TCR β s, as demonstrated by the profound enrichment of these TCR β s only among SARS-CoV-2-positive individuals expressing the appropriate restricting HLA. Thus, while only a subset of HLA-associated TCRs are observed in any given individual, the particular TCR β subset observed reflects that individual's history of exposure to many common antigens.

Our results imply that the vast majority of HLA-associated TCR β s identified in this study likely derive from common antigens, making HLA-association a critical step in decoding the human T-cell repertoire. Moreover, the high imputation accuracy of our HLA models allows us to statistically HLA type all repertoires ever sequenced, thereby expanding the effective size of HLA-typed and TCR-sequenced cohorts by several orders of magnitude and further facilitating decoding efforts. The TCR repertoire is a Rosetta Stone of the human immune system, providing a rich source of information for characterizing both the genetic background and exposure history of individuals at a population scale (70, 71). Mapping TCRs to HLAs and disease exposures and imputing HLA and disease exposure from TCRs represent important steps

toward decoding the immunological history of individuals using immunosequencing.

Materials and methods

Human samples

TCR β and HLA sequence data from human samples used for these studies were aggregated from several independent study collections described below. All necessary patient/participant consent has been obtained for each study and the appropriate institutional forms have been archived.

1. Whole blood samples from DLS (Discovery Life Sciences, Huntsville, AL) were collected under Protocol DLS13 for collection of clinical samples.
2. PBMC used for the sorted repertoire experiments were collected and processed by Bloodworks Northwest (Seattle, WA). Volunteer donors were consented and collected under the Bloodworks Research Donor Collection Protocol BT001.
3. PBMC were obtained from the Fred Hutchinson Cancer Research Center Research Cell Bank biorepository of healthy bone marrow donors. The sample collection protocol was approved and supervised by the Fred Hutchinson Cancer Research Center Institutional Review Board (IRB) (35).
4. Blood collected from human subjects were approved by the IRBs of Johns Hopkins University and Stanford University. All participants provided written informed consent prior to enrollment (38).
5. Blood collected for the ImmuneRACE Study has been approved by the Western IRB (WIRB) (reference number 1-1281891-1). The trial has received appropriate ethical approval from WIRB as described (72).
6. Procedures for the INCOV study were approved by the IRBs at Providence St. Joseph Health with IRB study number STUDY2020000175, the WIRB with IRB study number 20170658, and the University of Washington with IRB study numbers STUDY00000959 and STUDY00002929 (73).
7. Human samples from the Virology Research Clinic at the University of Washington was collected under an IRB-approved study (NCT04338360) (53).
8. Blood from the Institute of Medical, Pharmaceutical, and Health Sciences, Kanazawa University (KU) was collected with informed consent as described by documents approved by the KU review board (document number 585-2).

TCR sequencing

We use the ImmunoSEQ assay developed by Adaptive Biotechnologies to measure the TCR β repertoire. ImmunoSEQ is

⁶ Our group has explored such analyses in disease contexts and found consistent results using TCR α or TCR β repertoires.

TABLE 2 Model performance of class I HLAs.

HLA		Train	Train CV	Holdout
	N Cases	AUCROC	AUCROC	AUCROC
A*01:01	763	0.99	0.99	1.00
A*02:01	1314	1.00	1.00	1.00
A*02:02	74	1.00	1.00	1.00
A*02:05	75	1.00	1.00	0.99
A*02:06	37	1.00	0.97	1.00
A*03:01	703	1.00	1.00	0.99
A*11:01	354	1.00	1.00	0.99
A*23:01	231	1.00	1.00	1.00
A*23:17	31	1.00	1.00	1.00
A*24:02	495	1.00	0.99	0.99
A*25:01	107	1.00	1.00	1.00
A*26:01	144	1.00	1.00	0.99
A*29:02	196	1.00	1.00	1.00
A*30:01	174	1.00	1.00	1.00
A*30:02	133	1.00	1.00	1.00
A*31:01	177	1.00	1.00	1.00
A*32:01	179	1.00	1.00	1.00
A*33:01	74	1.00	1.00	0.99
A*33:03	142	1.00	1.00	1.00
A*34:02	60	1.00	1.00	0.99
A*36:01	58	1.00	1.00	1.00
A*66:01	40	0.99	0.98	1.00
A*68:01	204	1.00	1.00	1.00
A*68:02	160	1.00	1.00	1.00
A*74:01	111	1.00	0.99	1.00
B*07:02	660	1.00	1.00	1.00
B*08:01	559	1.00	0.99	1.00
B*13:02	125	1.00	1.00	1.00
B*14:01	56	1.00	0.99	1.00
B*14:02	175	1.00	1.00	1.00
B*15:01	291	1.00	0.99	1.00
B*15:03	104	1.00	1.00	1.00
B*15:10	72	1.00	1.00	1.00
B*15:16	36	1.00	1.00	1.00
B*18:01	261	1.00	1.00	1.00
B*27:05	169	1.00	1.00	0.99
B*35:01	390	1.00	1.00	1.00

(Continued)

TABLE 2 Continued

HLA		Train	Train CV	Holdout
	N Cases	AUCROC	AUCROC	AUCROC
B*35:02	53	1.00	0.98	0.96
B*35:03	91	1.00	1.00	1.00
B*37:01	65	1.00	1.00	1.00
B*38:01	85	1.00	0.99	0.99
B*39:01	71	1.00	1.00	1.00
B*40:01	290	1.00	1.00	1.00
B*40:02	80	1.00	1.00	0.99
B*42:01	94	1.00	1.00	1.00
B*44:02	405	1.00	1.00	0.99
B*44:03	333	1.00	1.00	1.00
B*45:01	114	1.00	1.00	1.00
B*49:01	123	1.00	1.00	1.00
B*50:01	78	1.00	0.99	1.00
B*51:01	269	1.00	1.00	0.99
B*52:01	74	1.00	1.00	1.00
B*53:01	199	1.00	1.00	1.00
B*55:01	95	1.00	1.00	1.00
B*56:01	32	0.99	0.96	0.83
B*57:01	166	1.00	0.99	0.99
B*57:03	54	1.00	0.99	1.00
B*58:01	109	1.00	1.00	1.00
B*58:02	80	1.00	0.99	0.99
C*01:02	208	0.99	0.95	0.96
C*02:02	209	0.99	0.97	0.94
C*02:10	119	1.00	0.98	0.97
C*03:02	59	1.00	0.95	1.00
C*03:03	282	0.96	0.89	0.89
C*03:04	480	0.98	0.95	0.96
C*04:01	829	0.99	0.99	0.98
C*05:01	446	0.99	0.96	0.95
C*06:02	559	1.00	1.00	0.99
C*07:01	763	1.00	0.98	0.99
C*07:02	757	1.00	0.96	0.97
C*07:04	73	1.00	0.66	0.68
C*07:18	81	0.98	0.96	0.93
C*08:01	30	1.00	0.95	0.87
C*08:02	246	1.00	0.99	0.99

(Continued)

TABLE 2 Continued

HLA	N Cases	Train	Train CV	Holdout
		AUCROC	AUCROC	AUCROC
C*12:02	55	1.00	0.92	0.94
C*12:03	246	1.00	1.00	0.98
C*14:02	75	1.00	0.98	0.99
C*15:02	114	0.99	0.84	0.93
C*15:05	43	0.99	0.93	0.98
C*16:01	320	1.00	0.99	0.98
C*17:01	132	1.00	0.94	0.94

The HLA allele and number of subjects with the given HLA are listed in columns 1 and 2, respectively. Columns 3–5 indicate model performance when the trained model is evaluated on the train, train CV and holdout data sets, respectively.

a multiplexed PCR-based method targeting rearranged TCR β sequences (45). The assay uses a panel of primers specific to all functional TRBV and TRBJ gene segments to amplify the full V(D)J region, including part of the TRBV segment, the entire diversity (D) and a portion of the J segment. The sequenced portions of the V and J segments are sufficient to determine gene usage. This region also includes the complete hypervariable complementarity-determining region 3 (CDR3), which accounts for most of the TCR diversity and is the key determinant of TCR antigen specificity. The reverse primer binds near the 3' end of the J region, adjacent to the constant region and may extend into the 5' end of the TCR. PCR primers contain universal priming sites on the 3'-end and Illumina sequencing adaptors on the 5'-end. This design enables robust amplification across a wide range of rearrangements while controlling for bias. High-throughput sequencing of these amplicons provides accurate quantification of unique TCR β clonotypes.

Flow cytometry and cell sorting $3 - 5 \times 10^7$ PBMC were stained using a cocktail of antibodies that included CD3 (clone UCHT1, Biolegend), CD4 (clone OKT4, Biolegend), CD8 (clone RPA-T8, Biolegend), CD45RA(clone HI100, BD Bioscience), CCR7 (clone G043H7, Biolegend), CD95 (clone DX2, Biolegend), and CD28 (clone CD28.2, Biolegend) for 10 minutes at 4 degrees C. PBMC were washed with MACS buffer (Miltenyi Biotec) and then PE+ CD3+ T cells were enriched using anti-PE Microbeads with LS columns (Miltenyi Biotec) following the manufacturer's protocol. The enriched CD3 sample was sorted using a FACSaria Fusion Flow Cytometer (BD Biosciences) to isolate naive CD4 (CD4⁺ CD3⁺ CD45RA⁺ CD95⁻ CD28⁺ CCR7⁺), memory CD4 (CD4⁺ CD3⁺ non-naive CD4), naive CD8 (CD8⁺ CD3⁺ CD45RA⁺ CD95⁻ CD28⁺ CCR7⁺) and memory CD8 (CD8⁺ CD3⁺ non-naive CD8). Between 150,000 and 3×10^6 T cells were sorted and sent in for sequencing for each subset. TCR β sequencing was carried out using the ImmunoSEQ assay at Adaptive Biotechnologies.

Model training

We independently train a model for each HLA using the following procedure:

1. We randomly split 80% and 20% of all typed samples into a training and holdout set.
2. Using typed HLA data, we label all samples as cases, controls or unlabeled for the HLA being modeled.
3. We derive an ES set using the FET applied to cases and controls.
4. We select a subset of ESs which associate to the HLA being modeled via the L1LR method.
5. We derive an individual weight for each sequence in the list of ESs from step 3).
6. We generate two features which are the logarithm of the weighted sum of the convergent recombination count (i.e., number of unique DNA rearrangements mapping to a given CDR3) of enhanced sequences and the total unique rearrangements in the repertoire, respectively.
7. We fit the standard Scikit-Learn logistic regression classifier to the training data and evaluate the model in cross validation to select the best FET p-value cutoff.
8. We train the final model on all the train data and evaluate on the holdout set.

Below we describe training steps in detail using the following definitions: Let H be the set of N_H typed HLA allotypes in our training set and $h_i \in H, i = 1, \dots, N_H$ refer to any single HLA in the set. The training set of h_i is denoted as $t_i \in T$ where T is the set of N_T training samples. We similarly define a test (or holdout) set for each HLA $t'_i \in T'$ where T' is the set of $N_{T'}$ holdout samples. Below we drop the subscript i unless needed for clarity as the procedure is applied to each HLA independently.

Step 1: T and T' are derived from a random 80/20 split of all typed data with $N_T + N_{T'} = 4,144$. **Supplementary Figure S1** shows the age, sex and ethnicity distribution of sample subjects.

Step 2: When building a model for a given HLA h_i , samples are labeled as cases (label = 1) or controls (label = 0) if a subject expresses or does not express h_i , respectively. A subset of controls are unlabeled if a **p-group matched HLA** to h_i is expressed by the subject. P-group matched HLAs share the same antigen binding domain and thus may share T-cell receptor (TCR) specificities.

Step 3: For any HLA h , an initial set of ESs, E , is defined using Fisher's Exact Test (FET) [74] applied to training set t of HLA h . The a range of p-value thresholds are used to identify ESs and the threshold is treated as a hyperparameter in the modeling; it is independently derived for each HLA we model.

Step 4: If HLA h is in linkage disequilibrium (LD) with other HLAs, the enhanced sequences specific to HLAs in LD will be present in the ES set of HLA h because of their co-occurrence in

TABLE 3 Model performance of class II HLAs.

HLA	Train		Train CV		Holdout	In LD
	N Cases	AUCROC	AUCROC	AUCROC	AUCROC	
DPA1*01:03+DPB1*01:01	387	1.00	0.99	0.99	0.99	No
DPA1*01:03+DPB1*02:01	796	1.00	0.99	0.99	1.00	Yes
DPA1*01:03+DPB1*02:02	39	1.00	0.97	0.97	0.97	Yes
DPA1*01:03+DPB1*03:01	415	0.99	0.99	0.99	0.99	Yes
DPA1*01:03+DPB1*04:01	1618	1.00	0.99	0.99	0.99	Yes
DPA1*01:03+DPB1*04:02	571	1.00	1.00	1.00	1.00	Yes
DPA1*01:03+DPB1*05:01	117	1.00	1.00	1.00	0.99	No
DPA1*01:03+DPB1*06:01	85	1.00	0.99	0.99	1.00	Yes
DPA1*01:03+DPB1*104:01	101	1.00	1.00	1.00	1.00	Yes
DPA1*01:03+DPB1*105:01	66	0.99	0.93	0.93	1.00	No
DPA1*01:03+DPB1*10:01	70	1.00	0.98	0.98	0.95	No
DPA1*01:03+DPB1*11:01	92	1.00	0.96	0.96	0.99	No
DPA1*01:03+DPB1*13:01	85	1.00	1.00	1.00	1.00	No
DPA1*01:03+DPB1*14:01	55	1.00	0.99	0.99	1.00	No
DPA1*01:03+DPB1*17:01	78	1.00	0.97	0.97	0.99	No
DPA1*01:03+DPB1*18:01	113	1.00	1.00	1.00	1.00	Yes
DPA1*02:01+DPB1*01:01	461	1.00	0.99	0.99	0.99	Yes
DPA1*02:01+DPB1*02:01	172	0.92	0.57	0.57	0.58	No
DPA1*02:01+DPB1*03:01	71	1.00	0.99	0.99	0.95	No
DPA1*02:01+DPB1*04:01	324	1.00	0.57	0.57	0.61	No
DPA1*02:01+DPB1*04:02	96	1.00	0.54	0.54	0.57	No
DPA1*02:01+DPB1*05:01	36	1.00	0.94	0.94	1.00	No
DPA1*02:01+DPB1*09:01	39	0.99	0.94	0.94	0.99	Yes
DPA1*02:01+DPB1*105:01	52	1.00	0.80	0.80	0.79	No
DPA1*02:01+DPB1*10:01	84	1.00	0.99	0.99	0.98	Yes
DPA1*02:01+DPB1*11:01	145	1.00	0.99	0.99	1.00	Yes
DPA1*02:01+DPB1*131:01	39	1.00	1.00	1.00	1.00	Yes
DPA1*02:01+DPB1*13:01	139	1.00	1.00	1.00	1.00	Yes
DPA1*02:01+DPB1*14:01	71	1.00	1.00	1.00	1.00	Yes
DPA1*02:01+DPB1*17:01	143	1.00	1.00	1.00	1.00	Yes
DPA1*02:01+DPB1*18:01	35	1.00	0.70	0.70	0.73	No
DPA1*02:02+DPB1*01:01	268	1.00	1.00	1.00	0.99	Yes
DPA1*02:02+DPB1*02:01	63	1.00	0.61	0.61	0.56	No
DPA1*02:02+DPB1*04:01	78	0.97	0.80	0.80	0.84	No
DPA1*02:02+DPB1*05:01	102	1.00	1.00	1.00	0.95	Yes
DPA1*02:06+DPB1*05:01	34	1.00	0.96	0.96	1.00	Yes
DPA1*03:01+DPB1*01:01	54	1.00	0.78	0.78	0.91	Yes

(Continued)

TABLE 3 Continued

HLA	Train		Train CV		Holdout	In LD
	N Cases	AUCROC	AUCROC	AUCROC	AUCROC	
DPA1*03:01+DPB1*105:01	140	1.00	1.00	1.00	1.00	Yes
DQA1*01:01+DQB1*02:01	65	1.00	0.80	0.83	0.83	No
DQA1*01:01+DQB1*02:02	78	1.00	0.81	0.72	0.72	No
DQA1*01:01+DQB1*03:01	122	1.00	0.78	0.75	0.75	No
DQA1*01:01+DQB1*03:02	53	1.00	0.76	0.68	0.68	No
DQA1*01:01+DQB1*05:01	602	1.00	0.99	1.00	1.00	Yes
DQA1*01:01+DQB1*06:02	81	1.00	0.91	0.98	0.98	No
DQA1*01:01+DQB1*06:03	32	1.00	0.94	0.98	0.98	No
DQA1*01:02+DQB1*02:01	161	1.00	0.81	0.85	0.85	No
DQA1*01:02+DQB1*02:02	181	1.00	0.80	0.77	0.77	No
DQA1*01:02+DQB1*03:01	214	1.00	0.72	0.73	0.73	No
DQA1*01:02+DQB1*03:02	102	1.00	0.69	0.76	0.76	No
DQA1*01:02+DQB1*03:03	41	0.99	0.54	0.55	0.55	No
DQA1*01:02+DQB1*03:19	59	1.00	0.86	0.91	0.91	No
DQA1*01:02+DQB1*04:02	63	1.00	0.81	0.65	0.65	No
DQA1*01:02+DQB1*05:01	221	1.00	0.99	0.99	0.99	No
DQA1*01:02+DQB1*05:02	149	1.00	1.00	1.00	1.00	Yes
DQA1*01:02+DQB1*05:03	33	1.00	0.85	0.90	0.90	No
DQA1*01:02+DQB1*06:02	830	1.00	1.00	1.00	1.00	Yes
DQA1*01:02+DQB1*06:03	70	1.00	0.97	0.99	0.99	No
DQA1*01:02+DQB1*06:04	184	1.00	0.99	1.00	1.00	Yes
DQA1*01:02+DQB1*06:09	135	1.00	1.00	1.00	1.00	Yes
DQA1*01:03+DQB1*02:01	46	1.00	0.62	0.76	0.76	No
DQA1*01:03+DQB1*02:02	43	1.00	0.64	0.59	0.59	No
DQA1*01:03+DQB1*03:01	62	1.00	0.68	0.69	0.69	No
DQA1*01:03+DQB1*03:02	33	1.00	0.62	0.56	0.56	No
DQA1*01:03+DQB1*05:01	55	0.99	0.92	0.91	0.91	No
DQA1*01:03+DQB1*06:01	54	1.00	1.00	1.00	1.00	Yes
DQA1*01:03+DQB1*06:02	53	1.00	1.00	0.95	0.95	No
DQA1*01:03+DQB1*06:03	299	1.00	1.00	0.99	0.99	Yes
DQA1*01:04+DQB1*05:03	122	1.00	1.00	1.00	1.00	Yes
DQA1*01:05+DQB1*05:01	145	1.00	1.00	1.00	1.00	Yes
DQA1*01:05+DQB1*06:02	30	1.00	0.81	0.74	0.74	No
DQA1*02:01+DQB1*02:01	74	1.00	1.00	1.00	1.00	No
DQA1*02:01+DQB1*02:02	602	1.00	1.00	1.00	1.00	Yes
DQA1*02:01+DQB1*03:01	131	1.00	0.99	0.98	0.98	No
DQA1*02:01+DQB1*03:02	64	1.00	0.97	0.97	0.97	No

(Continued)

TABLE 3 Continued

HLA	Train		Train CV		Holdout	In LD
	N Cases	AUCROC	AUCROC	AUCROC	AUCROC	
DQA1*02:01+DQB1*03:03	164	1.00	0.99	0.98	0.98	Yes
DQA1*02:01+DQB1*04:02	32	1.00	0.87	0.98	0.98	No
DQA1*02:01+DQB1*05:01	109	1.00	0.77	0.66	0.66	No
DQA1*02:01+DQB1*06:02	110	1.00	0.96	0.97	0.97	No
DQA1*02:01+DQB1*06:03	39	0.90	0.75	0.79	0.79	No
DQA1*02:01+DQB1*06:04	30	1.00	0.78	0.93	0.93	No
DQA1*03:01+DQB1*02:01	59	1.00	0.73	0.77	0.77	No
DQA1*03:01+DQB1*02:02	49	1.00	0.97	0.99	0.99	No
DQA1*03:01+DQB1*03:01	78	1.00	0.99	1.00	1.00	No
DQA1*03:01+DQB1*03:02	442	1.00	0.99	1.00	1.00	Yes
DQA1*03:01+DQB1*05:01	55	1.00	0.64	0.68	0.68	No
DQA1*03:01+DQB1*06:02	55	1.00	0.81	0.68	0.68	No
DQA1*03:02+DQB1*03:03	76	1.00	1.00	1.00	1.00	Yes
DQA1*03:03+DQB1*02:01	52	0.98	0.90	0.99	0.99	No
DQA1*03:03+DQB1*02:02	120	1.00	1.00	1.00	1.00	Yes
DQA1*03:03+DQB1*03:01	306	1.00	1.00	1.00	1.00	Yes
DQA1*03:03+DQB1*03:02	98	1.00	1.00	1.00	1.00	Yes
DQA1*03:03+DQB1*03:03	40	1.00	0.92	0.82	0.82	Yes
DQA1*03:03+DQB1*04:02	30	1.00	0.73	0.71	0.71	No
DQA1*03:03+DQB1*05:01	72	1.00	0.73	0.81	0.81	No
DQA1*03:03+DQB1*06:02	65	1.00	0.77	0.92	0.92	No
DQA1*04:01+DQB1*02:02	38	1.00	0.72	0.79	0.79	No
DQA1*04:01+DQB1*03:01	33	1.00	0.90	0.81	0.81	No
DQA1*04:01+DQB1*03:19	64	1.00	1.00	1.00	1.00	Yes
DQA1*04:01+DQB1*04:02	241	1.00	1.00	1.00	1.00	Yes
DQA1*04:01+DQB1*05:01	38	1.00	0.66	0.57	0.57	No
DQA1*04:01+DQB1*06:02	49	1.00	0.91	0.95	0.95	No
DQA1*05:01+DQB1*02:01	639	1.00	1.00	1.00	1.00	Yes
DQA1*05:01+DQB1*02:02	57	1.00	0.99	1.00	1.00	No
DQA1*05:01+DQB1*03:01	122	1.00	1.00	1.00	1.00	No
DQA1*05:01+DQB1*03:02	63	0.97	0.97	0.99	0.99	No
DQA1*05:01+DQB1*03:03	30	1.00	0.93	0.95	0.95	No
DQA1*05:01+DQB1*04:02	31	1.00	0.69	0.90	0.90	No
DQA1*05:01+DQB1*05:01	80	1.00	0.73	0.72	0.72	No
DQA1*05:01+DQB1*06:02	95	1.00	0.83	0.88	0.88	No
DQA1*05:01+DQB1*06:03	37	1.00	0.70	0.80	0.80	No
DQA1*05:05+DQB1*02:01	84	1.00	1.00	1.00	1.00	No

(Continued)

TABLE 3 Continued

HLA	Train		Train CV		Holdout	In LD
	N Cases	AUCROC	AUCROC	AUCROC	AUCROC	
DQA1*05:05+DQB1*02:02	69	1.00	0.97	1.00	1.00	No
DQA1*05:05+DQB1*03:01	601	1.00	1.00	0.99	0.99	Yes
DQA1*05:05+DQB1*03:02	66	0.99	0.94	0.94	0.94	No
DQA1*05:05+DQB1*03:19	116	1.00	1.00	1.00	1.00	Yes
DQA1*05:05+DQB1*05:01	101	1.00	0.76	0.67	0.67	No
DQA1*05:05+DQB1*06:02	117	1.00	0.84	0.82	0.82	No
DQA1*05:05+DQB1*06:03	41	1.00	0.73	0.77	0.77	No
DQA1*06:01+DQB1*03:01	40	1.00	1.00	1.00	1.00	Yes
DRB1*01:01	442	1.00	0.99	1.00	1.00	Yes
DRB1*01:02	138	1.00	1.00	1.00	1.00	Yes
DRB1*01:03	49	1.00	1.00	1.00	1.00	Yes
DRB1*03:01	655	1.00	1.00	1.00	1.00	Yes
DRB1*03:02	103	1.00	1.00	1.00	1.00	Yes
DRB1*04:01	421	1.00	1.00	1.00	1.00	Yes
DRB1*04:02	46	1.00	0.99	1.00	1.00	Yes
DRB1*04:03	49	1.00	1.00	1.00	1.00	Yes
DRB1*04:04	190	1.00	1.00	1.00	1.00	Yes
DRB1*04:05	76	1.00	1.00	1.00	1.00	Yes
DRB1*04:07	72	1.00	1.00	1.00	1.00	Yes
DRB1*07:01	750	1.00	1.00	1.00	1.00	Yes
DRB1*08:01	109	1.00	1.00	1.00	1.00	Yes
DRB1*08:02	34	1.00	1.00	1.00	1.00	Yes
DRB1*08:04	93	1.00	1.00	1.00	1.00	Yes
DRB1*09:01	126	1.00	1.00	1.00	1.00	Yes
DRB1*10:01	79	1.00	1.00	1.00	1.00	Yes
DRB1*11:01	388	1.00	1.00	1.00	1.00	Yes
DRB1*11:02	85	1.00	1.00	1.00	1.00	Yes
DRB1*11:04	141	1.00	1.00	0.99	0.99	Yes
DRB1*12:01	154	1.00	1.00	0.99	0.99	Yes
DRB1*13:01	364	1.00	1.00	0.99	0.99	Yes
DRB1*13:02	353	1.00	1.00	1.00	1.00	Yes
DRB1*13:03	95	1.00	0.99	1.00	1.00	Yes
DRB1*14:54	129	1.00	1.00	1.00	1.00	Yes
DRB1*15:01	611	1.00	1.00	1.00	1.00	Yes
DRB1*15:02	75	1.00	1.00	1.00	1.00	Yes
DRB1*15:03	202	1.00	1.00	1.00	1.00	Yes
DRB1*16:01	64	1.00	1.00	1.00	1.00	Yes

(Continued)

TABLE 3 Continued

HLA	N Cases	Train	Train CV	Holdout	In LD
		AUCROC	AUCROC	AUCROC	
DRB1*16:02	54	1.00	1.00	1.00	Yes
DRB3*01:01	654	1.00	0.99	0.99	Yes
DRB3*01:62	91	1.00	0.99	1.00	Yes
DRB3*02:02	1005	1.00	1.00	1.00	Yes
DRB3*03:01	374	1.00	1.00	1.00	Yes
DRB4*01:01	361	1.00	1.00	1.00	Yes
DRB4*01:03	1014	0.99	0.98	0.98	Yes
DRB5*01:01	677	1.00	1.00	1.00	Yes
DRB5*01:02	41	1.00	1.00	1.00	Yes
DRB5*02:02	75	1.00	1.00	1.00	Yes

The HLA allele and number of subjects with the given HLA are listed in columns 1 and 2, respectively. Columns 3–5 indicate model performance when the trained model is evaluated on the train, train CV and holdout data sets, respectively. Column 6 indicates whether the heterodimer is observed to be in linkage disequilibrium. HLAs in linkage equilibrium likely result from trans-complementation of the subunits (see main text).

subjects (this is how LD is defined). The goal of our procedure is to identify and remove these sequences from the final set of ESs associated to h . Thus, a fundamental assumption of the model is that for a given train set, an ES may be associated with only a single HLA.

We enforce the assumption that each TCR β is specific to one HLA by fitting a logistic regression model with L1 regularization and we refer to this as the L1LR method. We independently associate each ES in E with a single HLA. Thus, the assumption that a given TCR β is associated to a single HLA is not globally enforced, i.e., it is possible that due to variations in training data that the same TCR β may be associated to multiple HLAs though this is rare (see [Supplementary Figure S2](#)).

When associating individual ESs to HLAs, we model the presence/absence of any given ES $e_j \in E$ in the training set t as a logistic regression classification. Here j denotes the individual TCR β s in the ES set. For any sample $t_l \in t$, where l denotes individual samples in the training set t , the label for sequence e_j is either 0 or 1 if the sequences is present or absent in t_l , respectively. The feature vector for associating a given ES e_j includes the indicator vector for the typed set of HLAs, i.e. $I_{kl} = 1$ if sample l expresses HLA h_k and $I_{kl} = 0$ if sample l does not express HLA h_k . Here, k indicates the subset of HLAs in H is in LD with HLA h (defined by FET p-value $< 10^{-3}$). Thus, the feature vector is a binary vector indicating whether a subject has or does not have a given HLA which is observed to co-occur with the HLA being modeled. When associating sequences we include the total number of unique rearrangements in each sample repertoire, d_l as a covariate not subject to L1 regularization. We model each ES independently and thus drop the subscript j in the following for clarity. Taken together a single ES e is modeled such that.

$$\text{logit}(\text{label} = 1) = \sum_{k,l} I_{kl} \beta_k + \sum_l \alpha d_l + \lambda \sum_k \beta_k + b. \quad (1)$$

In [Equation 1](#), α and β are free parameters, b is the bias term and λ is the regularization strength applied to β . We find the best-fit parameters by minimizing the log-loss and tune λ to be the smallest value that yields precisely one non-zero coefficient in β_k . If the non-zero coefficient is $\beta_{k=i}$, then the sequence is associated with the HLA being modeled and is retained as part of the final ES set. If the non-zero coefficient is $\beta_{k \neq i}$, the sequences is excluded from the final ES set of HLA h which is denoted as \tilde{E} . We weight each sample by the square root of the convergent recombination count when finding the best-fit parameters and zero count samples are given a weight of unity. The effectiveness of our procedure to resolve LD is demonstrated in [Supplementary Figure S2](#).

Step 5: We derive a per-sequence weight, w_j for each sequence $e_j \in \tilde{E}$ by fitting a two feature logistic regression classifier. Here j refers to individual sequences in \tilde{E} . Each sequence is fit independently. For each sequence the target of the classifier is whether the sample is a case or control of HLA h_i and the features are the convergent recombination count of the sequence and total number of unique rearrangements in the given repertoire sample. Essentially, we build a model to discriminate cases from controls using the count of a single sequence and take the coefficient as a weight. In detail, we standardize the features by subtracting the mean and normalizing to the standard deviation of the features across all samples. We fit the default logistic regression model in the [Scikit-learn](#) library which includes a fixed amount of L2 regularization ($\lambda = 1$). This procedure yields a weight for each sequence in $e_j \in \tilde{E}$.

Step 6: The final classifier model is a two feature model where the features F for subject l are.

$$F_l = \left[\log(1 + \sum_j w_j c_{lj}), \log(1 + N_l) \right]. \quad (2)$$

In [Equation 2](#) c_{lj} is the convergent recombination count of e_j in sample l . These features, F_l , are used in a the standard Scikit-learn

logistic regression to predict the presence/absence of an HLA for a given sample.

Step 7: To model any given HLA, we fit the standard Scikit-Learn logistic regression classifier using features generated in Step 6). We treat the p-value cutoff for identifying ESs via the FET as described in Step 3) as a hyperparameter of each HLA model. We train and evaluate a model using a five-fold cross validation strategy for p-value cutoffs = $[10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}, 10^{-8}]$. We adopt p-value cutoff of the highest performing cross-validation model.

Step 8: We generate our final model by adopting the p-value threshold determined in step 7) and train on the full set of training data. We evaluate the final model on the holdout set (Supplementary Figure S3–S8).

Data availability statement

A subset of the raw data sufficient to reproduce the results of this study is publicly available from Adaptive Biotechnologies at <https://clients.adaptivebiotech.com/pub/Emerson-2017-NatGen>. The remaining data were generated by Adaptive Biotechnologies and are subject to their data sharing policies.

Ethics statement

There are multiple cohorts used in the paper and the IRBs are listed. The studies were conducted in accordance with the local legislation and institutional requirements. The participants provided their written informed consent to participate in this study.

Author contributions

HZ: Conceptualization, Data curation, Investigation, Formal analysis, Methodology, Project administration, Software, Visualization, Writing – original draft, Writing – review & editing. RT: Conceptualization, Methodology, Investigation, Writing – original draft, Writing – review & editing. PE: Conceptualization, Investigation, Writing – original draft, Writing – review & editing. I-TC: Conceptualization, Methodology, Investigation, Writing – original draft, Writing – review & editing. CG: Resources, Data curation, Software, Writing – review & editing. JL: Resources, Data curation, Software, Writing – review & editing. LP: Resources, Data curation, Software, Writing – review & editing. MR: Validation, Writing – review & editing. RE: Validation, Writing – review & editing. HT: Validation, Writing – review & editing. WZ: Project administration, Supervision, Writing – review & editing. JG: Conceptualization, Methodology, Writing – review & editing, Supervision. HR: Conceptualization, Funding acquisition, Resources, Supervision, Writing – review & editing. JC: Conceptualization, Resources, Supervision, Writing – original draft, Writing – review & editing.

Funding

The author(s) declare that no financial support was received for the research and/or publication of this article.

Acknowledgments

We are grateful to Jonathan Duke-Cohan for his careful and engaged editorial oversight of the manuscript. We thank Mary Carrington for discussion and helpful feedback which improved the manuscript and Thomas Snyder for guidance during early stages of manuscript drafting. We also thank the reviewers who provided helpful feedback that improved the clarity of the manuscript.

In memoriam

This paper is dedicated to the memory of Peter Jacob Robert Ebert (1978 – 2023). He marched to his own drumbeat and was unconcerned with popular opinions. His intellectual curiosity, independence and rigor made him an innovative scientific leader and a wonderful collaborator. He was not only an incredibly clever scientist, but also a funny and kind colleague who will be missed by all who had the privilege to know him.

Conflict of interest

HZ, CG, JL, LP, JG, JC were employed by Microsoft. RT, PE, I-TC, MR, RE, WZ, HR were employed by Adaptive Biotechnologies.

The remaining 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.

Any alternative text (alt text) provided alongside figures in this article has been generated by Frontiers with the support of artificial intelligence and reasonable efforts have been made to ensure accuracy, including review by the authors wherever possible. If you identify any issues, please contact us.

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.2025.1603730/full#supplementary-material>

SUPPLEMENTARY FIGURE 1

Demographic distribution of sample population. (A) Age, (B) sex and (C) self-reported ethnicity distribution of sample subjects.

SUPPLEMENTARY FIGURE 2

Effectiveness of L1LR method to resolve LD. Heatmap showing the amount of TCR β sharing before (left) and after (right) applying L1LR method to resolve LD. The color indicates the fractional sharing of TCR β s amongst pairs of HLAs relative to the total number of TCR β s associated with the HLA indicated on the vertical axis. After accounting for LD, we observe significantly less sharing of TCR β s amongst HLA allotypes. The remaining shared TCR β s may be the result of LD that is too strong to resolve using our method or true sharing of TCR β s between multiple HLAs (e.g., as shown in the main text).

SUPPLEMENTARY FIGURE 3

HLA-A models. Left column shows the precision-recall curve for the trained model evaluated on the train sample (middle left), train sample in cross-validation (middle right) and on the holdout sample (right) for HLA-A. The green line indicates the calibration threshold.

SUPPLEMENTARY FIGURE 4

HLA-B models. Left column shows the precision-recall curve for the trained model evaluated on the train sample (middle left), train sample in cross-validation (middle right) and on the holdout sample (right) for HLA-B. The green line indicates the calibration threshold.

SUPPLEMENTARY FIGURE 5

HLA-C models. Left column shows the precision-recall curve for the trained model evaluated on the train sample (middle left), train sample in cross-validation (middle right) and on the holdout sample (right) for HLA-C. The green line indicates the calibration threshold.

SUPPLEMENTARY FIGURE 6

HLA-DP models. Left column shows the precision-recall curve for the trained model evaluated on the train sample (middle left), train sample in cross-validation (middle right) and on the holdout sample (right) for HLA-DP. The green line indicates the calibration threshold.

SUPPLEMENTARY FIGURE 7

HLA-DQ models. Left column shows the precision-recall curve for the trained model evaluated on the train sample (middle left), train sample in cross-validation (middle right) and on the holdout sample (right) for HLA-DQ. The green line indicates the calibration threshold.

SUPPLEMENTARY FIGURE 8

HLA-DR models. Left column shows the precision-recall curve for the trained model evaluated on the train sample (middle left), train sample in cross-validation (middle right) and on the holdout sample (right) for HLA-DR. The green line indicates the calibration threshold.

SUPPLEMENTARY FIGURE 9

TCRs specificity apparently independent of synonymous mutations in coding regions of HLAs. Number of TCR β s overlapping in two ES sets of HLAs differing only in their third field designation as a function of the FET p-value threshold used to identify ESs. The solid line shows the overlap in two sets of ESs identified using the three field designation to divide samples into two groups. Error bars are poisson uncertainties. Even if the TCR specificity is identical for HLAs differing only in their third field designation, we expect differences in the ES sets due to statistical fluctuations resulting from sampling. To account for this effect, we compare the observed overlap fraction when segregating samples by third field designation to a null distribution of overlap fractions derived from 100 realizations of two samples which are randomly divided irrespective of the third field designation. The two randomly divided samples in each realization are selected such that the relative size of the two samples matches the size of two samples split by the three field designation. Shaded regions indicate the distribution of the overlap fraction from these 100 realizations. Consistency between the overlap fraction of ESs derived from samples differing in their three field resolution and ones based on random permutation demonstrates that TCR specificity is independent of three field resolution for the two HLA allotypes tested.

SUPPLEMENTARY FIGURE 10

We test the hypothesis that homozygosity at one loci increases breadth at another class matched loci. The distribution of zscores across class matched loci from a pairwise loci comparison. We calculate the zscore of the breadth at one loci in the case where another class matched loci is heterozygous or homozygous. We aggregate all pairwise comparisons for class I and class II separately.

References

- Hughes AL, Yeager M. Natural selection at major histocompatibility complex loci of vertebrates. *Annu Rev Genet.* (1998) 32:415–35. doi: 10.1146/annurev.genet.32.1.415
- Zinkernagel RM, Doherty PC. Restriction of *in vitro* T cell-mediated cytotoxicity in lymphocytic choriomeningitis within a syngeneic or semiallogeneic system. *Nature.* (1974) 248:701–2. doi: 10.1038/248701a0
- Davis MM, Boniface JJ, Reich Z, Lyons D, Hampl J, Arden B, et al. Ligand recognition by (alpha)(beta) T cell receptors. *Annu Rev Immunol.* (1998) 16:523. doi: 10.1146/annurev.immunol.16.1.523
- Martin MP, Carrington M. Immunogenetics of HIV disease. *Immunol Rev.* (2013) 254:245–64. doi: 10.1111/imr.12071
- Schumacher TN, Schreiber RD. Neoantigens in cancer immunotherapy. *Science.* (2015) 348:69–74. doi: 10.1126/science.aaa4971
- Montgomery RA, Tatapudi VS, Leffell MS, Zachary AA. HLA in transplantation. *Nat Rev Nephrol.* (2018) 14:558–70. doi: 10.1038/s41581-018-0039-x
- Schumacher TN, Scheper W, Kvistborg P. Cancer neoantigens. *Annu Rev Immunol.* (2019) 37:173–200. doi: 10.1146/annurev-immunol-042617-053402
- Kovacs AA, Kono N, Wang C-H, Wang D, Frederick T, Operskalski E, et al. Association of HLA genotype with T-cell activation in human immunodeficiency virus (HIV) and HIV/hepatitis C virus-coinfected women. *J Infect Dis.* (2020) 221:1156–66. doi: 10.1093/infdis/jiz589
- Francis JM, Leistritz-Edwards D, Dunn A, Tarr C, Lehman J, Dempsey C, et al. Allelic variation in class I HLA determines CD8+ T cell repertoire shape and cross-reactive memory responses to SARS-CoV-2. *Sci Immunol.* (2021) 7:eabk3070.
- Olafsdottir TA, Bjarnadottir K, Norddahl GL, Halldorsson GH, Melsted P, Gunnarsdottir K, et al. HLA alleles, disease severity, and age associate with T-cell responses following infection with SARS-CoV-2. *Commun Biol.* (2022) 5:914. doi: 10.1038/s42003-022-03893-w
- Robinson J, Halliwell JA, McWilliam H, Lopez R, Parham P, Marsh SG. The imgt/hla database. *Nucleic Acids Res.* (2012) 41:D1222–7. doi: 10.1093/nar/gks949
- Bjorkman PJ, Saper M, Samraoui B, Bennett WS, Strominger Jt, Wiley D. Structure of the human class I histocompatibility antigen, HLA-A2. *Nature.* (1987) 329:506–12. doi: 10.1038/329506a0
- Fremont DH, Matsumura M, Stura EA, Peterson PA, Wilson IA. Crystal structures of two viral peptides in complex with murine MHC class I H-2Kb. *Science.* (1992) 257:919–27. doi: 10.1126/science.1323877
- Babbitt BP, Allen PM, Matsueda G, Haber E, Unanue ER. Binding of immunogenic peptides to Ia histocompatibility molecules. *Nature.* (1985) 317:359–61. doi: 10.1038/317359a0
- Brown JH, Jardetzky TS, Gorga JC, Stern LJ, Urban RG, Strominger JL, et al. Three-dimensional structure of the human class II histocompatibility antigen HLA-DR1. *Nature.* (1993) 364:33–9. doi: 10.1038/364033a0
- Fremont DH, Hendrickson WA, Marrack P, Kappler J. Structures of an MHC class II molecule with covalently bound single peptides. *Science.* (1996) 272:1001–4. doi: 10.1126/science.272.5264.1001
- Marsh SGE, Albert ED, Bodmer WF, Bontrop RE, Dupont B, Erlich HA, et al. An update to HLA nomenclature, 2010. *Bone Marrow Transplant.* (2010) 45:846–8. doi: 10.1038/bmt.2010.79
- Sidney J, Peters B, Frahm N, Brander C, Sette A. HLA class I supertypes: a revised and updated classification. *BMC Immunol.* (2008) 9:1–15. doi: 10.1186/1471-2172-9-1
- Apps R, Qi Y, Carlson JM, Chen H, Gao X, Thomas R, et al. Influence of HLA-C expression level on HIV control. *Science.* (2013) 340:87–91. doi: 10.1126/science.1232685
- Ramsuran V, Naranbhai V, Horowitz A, Qi Y, Martin MP, Yuki Y, et al. Elevated HLA-A expression impairs HIV control through inhibition of NKG2A-expressing cells. *Science.* (2018) 359:86–90. doi: 10.1126/science.aam8825
- Martin MP, Gao X, Lee J-H, Nelson GW, Detels R, Goedert JJ, et al. Epistatic interaction between KIR3DS1 and HLA-B delays the progression to AIDS. *Nat Genet.* (2002) 31:429–34. doi: 10.1038/ng934

22. Kloverpris HN, Cole DK, Fuller A, Carlson J, Beck K, Schauenburg AJ, et al. A molecular switch in immunodominant HIV-1-specific CD8 T-cell epitopes shapes differential HLA-restricted escape. *Retrovirology*. (2015) 12:1–11. doi: 10.1186/s12977-015-0149-5
23. Illing PT, Pymm P, Croft NP, Hilton HG, Jojic V, Han AS, et al. HLA-B57 micropolymorphism defines the sequence and conformational breadth of the immunopeptidome. *Nat Commun*. (2018) 9:4693. doi: 10.1038/s41467-018-07109-w
24. Carlson JM, Listgarten J, Pfeifer N, Tan V, Kadie C, Walker BD, et al. Widespread impact of HLA restriction on immune control and escape pathways of HIV-1. *J Virol*. (2012) 86:5230–43. doi: 10.1128/JVI.06728-11
25. Hedrick SM, Cohen DI, Nielsen EA, Davis MM. Isolation of Cdna clones encoding T cell-specific membrane-associated proteins. *Nature*. (1984) 308:149–53. doi: 10.1038/308149a0
26. Yanagi Y, Yoshikai Y, Leggett K, Clark SP, Aleksander I, Mak TW. A human T cell-specific cDNA clone encodes a protein having extensive homology to immunoglobulin chains. *Nature*. (1984) 308:145–9. doi: 10.1038/308145a0
27. Fink PJ, Bevan MJ. H-2 antigens of the thymus determinelymphocyte specificity. *J Exp Med*. (1978) 148:766–75. doi: 10.1084/jem.148.3.766
28. Zinkernagel RM, Callahan GN, Klein J, Dennert G. Cytotoxic T cells learn specificity for self H-2 during differentiation in the thymus. *Nature*. (1978) 271:251–3. doi: 10.1038/271251a0
29. Philpott KL, Viney JL, Kay G, Rastan S, Gardiner EM, Chae S, et al. Lymphoid development in mice congenitally lacking T cell receptor $\alpha\beta$ -expressing cells. *Science*. (1992) 256:1448–52. doi: 10.1126/science.1604321
30. Huseby ES, White J, Crawford F, Vass T, Becker D, Pinilla C, et al. How the T cell repertoire becomes peptide and MHC specific. *Cell*. (2005) 122:247–60. doi: 10.1016/j.cell.2005.05.013
31. Kisielow P, Teh HS, Blüthmann H, von Boehmer H. Positive selection of antigen-specific T cells in thymus by restricting MHC molecules. *Nature*. (1988) 335:730–3. doi: 10.1038/335730a0
32. von Boehmer H, Kisielow P, Kishi H, Scott B, Borgulya P, Teh HS. The expression of CD4 and CD8 accessory molecules on mature T cells is not random but correlates with the specificity of the $\alpha\beta$ receptor for antigen. *Immunol Rev*. (1989) 109:143–52. doi: 10.1111/j.1600-065X.1989.tb00023.x
33. Kaye J, Hsu M-L, Sauron M-E, Jameson SC, Gascoigne NR, Hedrick SM. Selective development of CD4+ T cells in transgenic mice expressing a class II MHC-restricted antigen receptor. *Nature*. (1989) 341:746–9. doi: 10.1038/341746a0
34. Burnet FM. A modification of jerne's theory of antibody production using the concept of clonal selection. *CA: A Cancer J Clin*. (1976) 26:119–21. doi: 10.3322/canjclin.26.2.119
35. Emerson RO, DeWitt WS, Vignali M, Gravley J, Hu JK, Osborne EJ, et al. Immunosequencing identifies signatures of cytomegalovirus exposure history and HLA-mediated effects on the T cell repertoire. *Nat Genet*. (2017) 49:659–65. doi: 10.1038/ng.3822
36. DeWitt WSIII, Smith A, Schoch G, Hansen JA, Matsen FAIV, Bradley P. Human T cell receptor occurrence patterns encode immune history, genetic background, and receptor specificity. *Elife*. (2018) 7:e38358. doi: 10.7554/eLife.38358.043
37. Snyder TM, Gittelman RM, Klinger M, May DH, Osborne EJ, Taniguchi R, et al. Magnitude and dynamics of the T-cell response to SARS-CoV-2 infection at both individual and population levels. *MedRxiv*. (2020). doi: 10.1101/2020.07.31.20165647
38. Greissl J, Pesesky M, Dalai SC, Rebman AW, Soloski MJ, Horn EJ, et al. Immunosequencing of the T-cell receptor repertoire reveals signatures specific for diagnosis and characterization of early Lyme disease. *medRxiv*. (2021). doi: 10.1101/2021.07.30.21261353
39. Bradley P, Thomas PG. Using T cell receptor repertoires to understand the principles of adaptive immune recognition. *Annu Rev Immunol*. (2019) 37:547–70. doi: 10.1146/annurev-immunol-042718-041757
40. Cowell LG. The diagnostic, prognostic, and therapeutic potential of adaptive immune receptor repertoire profiling in cancerAdaptive immune receptor repertoire profiling in cancer. *Cancer Res*. (2020) 80:643–54. doi: 10.1158/0008-5472.CAN-19-1457
41. Arstila TP, Casrouge A, Baron V, Even J, Kanellopoulos J, Kourilsky P. A direct estimate of the human $\alpha\beta$ T cell receptor diversity. *Science*. (1999) 286:958–61. doi: 10.1126/science.286.5441.958
42. Sewell AK. Why must T cells be cross-reactive? *Nat Rev Immunol*. (2012) 12:669–77. doi: 10.1038/nri3279
43. Ruiz Ortega M, Pogorely MV, Minervina AA, Thomas PG, Mora T, Walczak AM. Learning predictive signatures of HLA type from T-cell repertoires. *PLoS Comput Biol*. (2025) 21:e1012724. doi: 10.1371/journal.pcbi.1012724
44. ElAbd H, Mahdy AK, Wacker EM, Gretsava M, Ellinghaus D, Franke A. Decoding the restriction of T cell receptors to human leukocyte antigen alleles using statistical learning. *bioRxiv*. (2025), 2025–02. doi: 10.1101/2025.02.06.636910
45. Carlson CS, Emerson RO, Sherwood AM, Desmarais C, Chung M-W, Parsons JM, et al. Using synthetic templates to design an unbiased multiplex PCR assay. *Nat Commun*. (2013) 4:1–9. doi: 10.1038/ncomms3680
46. Robins H. Immunosequencing: applications of immune repertoire deep sequencing. *Curr Opin Immunol*. (2013) 25:646–52. doi: 10.1016/j.coi.2013.09.017
47. Smith AG, Pereira S, Jaramillo A, Stoll ST, Khan FM, Berka N, et al. Comparison of sequence-specific oligonucleotide probe vs next generation sequencing for HLA-A, B, C, DRB1, DRB3/4/5, DQA1, DQB1, DPA1, and DPB1 typing: Toward single-pass high-resolution HLA typing in support of solid organ and hematopoietic cell transplant programs. *Hla*. (2019) 94:296–306. doi: 10.1111/tan.13619
48. Lewontin RC. *The genetic basis of evolutionary change* Vol. 560. . New York: Columbia University Press (1974).
49. Zarnitsyna VI, Evavold BD, Schoettle LN, Blattman JN, Antia R. Estimating the diversity, completeness, and cross-reactivity of the T cell repertoire. *Front Immunol*. (2013) 4:485. doi: 10.3389/fimmu.2013.00485
50. Bentzen AK, Marquard AM, Lyngaa R, Saini SK, Ramskov S, Donia M, et al. Large-scale detection of antigen-specific T cells using peptide-MHC-I multimers labeled with DNA barcodes. *Nat Biotechnol*. (2016) 34:1037–45. doi: 10.1038/nbt.3662
51. 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
52. Gittelman RM, Lavezzo E, Snyder TM, Zahid HJ, Carty CL, Elyanow R, et al. Longitudinal analysis of T cell receptor repertoires reveals shared patterns of antigen-specific response to SARS-CoV-2 infection. *JCI Insight*. (2022) 7(10), p.e151849. doi: 10.1172/jci.insight.151849
53. Elyanow R, Snyder TM, Dalai SC, Gittelman RM, Boonyaratanakornkit J, Wald A, et al. T cell receptor sequencing identifies prior SARS-CoV-2 infection and correlates with neutralizing antibodies and disease severity. *JCI Insight*. (2022) 7(10), p.e150070. doi: 10.1172/jci.insight.150070
54. Macdonald WA, Purcell AW, Mifsud NA, Ely LK, Williams DS, Chang L, et al. A naturally selected dimorphism within the HLA-B44 supertype alters class I structure, peptide repertoire, and T cell recognition. *J Exp Med*. (2003) 198:679–91. doi: 10.1084/jem.20030066
55. Rossjohn J, Gras S, Miles JJ, Turner SJ, Godfrey DI, McCluskey J. T cell antigen receptor recognition of antigen-presenting molecules. *Annu Rev Immunol*. (2015) 33:169–200. doi: 10.1146/annurev-immunol-032414-112334
56. Carrington M, Nelson GW, Martin MP, Kissner T, Vlahov D, Goedert JJ, et al. HLA and HIV-1: heterozygote advantage and B* 35-Cw* 04 disadvantage. *Science*. (1999) 283:1748–52. doi: 10.1126/science.283.5408.1748
57. Chowell D, Morris LGT, Grigg CM, Weber JK, Samstein RM, Makarov V, et al. Patient HLA class I genotype influences cancer response to checkpoint blockade immunotherapy. *Science*. (2018) 359:582–7. doi: 10.1126/science.aao4572
58. Neefjes JJ, Ploegh HL. Allele and locus-specific differences in cell surface expression and the association of HLA class I heavy chain with β 2-microglobulin: differential effects of inhibition of glycosylation on class I subunit association. *Eur J Immunol*. (1988) 18:801–10. doi: 10.1002/eji.1830180522
59. Choo SY. The HLA system: genetics, immunology, clinical testing, and clinical implications. *Yonsei Med J*. (2007) 48:11–23. doi: 10.3349/ymj.2007.48.1.11
60. Braunstein NS, Germain RN. Allele-specific control of Ia molecule surface expression and conformation: implications for a general model of Ia structurefunction relationships. *Proc Natl Acad Sci*. (1987) 84:2921–5. doi: 10.1073/pnas.84.9.2921
61. Kwok WW, Nepom GT. Structural and functional constraints on HLA class II dimers implicated in susceptibility to insulin dependent diabetes mellitus. *Bailliere's Clin Endocrinol Metab*. (1991) 5:375–93. doi: 10.1016/S0950-351X(05)80137-5
62. Tollefsen S, Hotta K, Chen X, Simonsen B, Swaminathan K, Mathews II, et al. Structural and functional studies of trans-encoded HLA-DQ2. 3 (DQA1* 03: 01/ DQB1* 02: 01) protein molecule. *J Biol Chem*. (2012) 287:13611–9. doi: 10.1074/jbc.M111.320374
63. Slatkin M. Linkage disequilibrium—understanding the evolutionary past and mapping the medical future. *Nat Rev Genet*. (2008) 9:477–85. doi: 10.1038/nrg2361
64. Johnson SA, Seale SL, Gittelman RM, Rytlewski JA, Robins HS, Fields PA. Impact of HLA type, age and chronic viral infection on peripheral T-cell receptor sharing between unrelated individuals. *PLoS One*. (2021) 16:e0249484. doi: 10.1371/journal.pone.0249484
65. Wooldridge I, Ekeruche-Makinde J, van den Berg HA, Skowera A, Miles JJ, Tan MP, et al. A single autoimmune T cell receptor recognizes more than a million different peptides. *J Biol Chem*. (2012) 287:1168–77. doi: 10.1074/jbc.M111.289488
66. Robins HS, Campregher PV, Srivastava SK, Wachter A, Turtle CJ, Kahsai O, et al. Comprehensive assessment of T-cell receptor β -chain diversity in $\alpha\beta$ T cells. *Blood*. (2009) 114:4099–107. doi: 10.1182/blood-2009-04-217604
67. Qi Q, Liu Y, Cheng Y, Glanville J, Zhang D, Lee J-Y, et al. Diversity and clonal selection in the human T-cell repertoire. *Proc Natl Acad Sci*. (2014) 111:13139–44. doi: 10.1073/pnas.1409155111

68. Murugan A, Mora T, Walczak AM, Callan CG. Statistical inference of the generation probability of T-cell receptors from sequence repertoires. *Proc Natl Acad Sci*. (2012) 109:16161–6. doi: 10.1073/pnas.1212755109
69. Yewdell JW. Confronting complexity: real-world immunodominance in antiviral CD8+ T cell responses. *Immunity*. (2006) 25:533–43. doi: 10.1016/j.immuni.2006.09.005
70. May DH, Woodhouse S, Zahid HJ, Elyanow R, Doroschak K, Noakes MT, et al. Identifying immune signatures of common exposures through co-occurrence of T-cell receptors in tens of thousands of donors. *BioRxiv*. (2024), 2024–03. doi: 10.1101/2024.03.26.583354
71. Chapfuwa P, Demirel I, Pisani L, Zazo J, Portugaly E, Zahid HJ, et al. , in: *In The Thirteenth International Conference on Learning Representations*, .
72. Dines JN, Manley TJ, Svejnova E, Simmons HM, Taniguchi R, Klinger M, et al. The immunerace study: A prospective multicohort study of immune response action to covid-19 events with the immunecode™ open access database. *medRxiv*. (2020). doi: 10.1101/2020.08.17.20175158
73. Su Y, Yuan D, Chen DG, Ng RH, Wang K, Choi J, et al. Multiple early factors anticipate post-acute COVID-19 sequelae. *Cell*. (2022) 185:881–95. doi: 10.1016/j.cell.2022.01.014
74. Fisher RA. On the interpretation of χ^2 from contingency tables, and the calculation of P. *J R Stat Soc*. (1922) 85:87–94. doi: 10.2307/2340521
75. Surh CD, Sprent J. Homeostasis of naive and memory T cells. *Immunity*. (2008) 29:848–62. doi: 10.1016/j.immuni.2008.11.002