

A large peptidome dataset improves HLA class I epitope prediction across most of the human population

Siranush Sarkizova^{1,2,13}, Susan Klaeger¹³, Phuong M. Le³, Letitia W. Li³, Giacomo Oliveira³, Hasmik Keshishian², Christina R. Hartigan², Wandu Zhang³, David A. Braun^{2,3,4,5}, Keith L. Ligon^{2,4,6,7}, Pavan Bachireddy^{2,3,5}, Ioannis K. Zervantonakis⁸, Jennifer M. Rosenbluth⁸, Tamara Ouspenskaia², Travis Law², Sune Justesen⁹, Jonathan Stevens¹⁰, William J. Lane^{4,10}, Thomas Eisenhaure², Guang Lan Zhang^{3,4,11}, Karl R. Clauser², Nir Hacohen^{2,3,12*}, Steven A. Carr^{2*}, Catherine J. Wu^{2,3,4,5*} and Derin B. Keskin^{2,3,4,5,11*}

Prediction of HLA epitopes is important for the development of cancer immunotherapies and vaccines. However, current prediction algorithms have limited predictive power, in part because they were not trained on high-quality epitope datasets covering a broad range of HLA alleles. To enable prediction of endogenous HLA class I-associated peptides across a large fraction of the human population, we used mass spectrometry to profile >185,000 peptides eluted from 95 HLA-A, -B, -C and -G mono-allelic cell lines. We identified canonical peptide motifs per HLA allele, unique and shared binding submotifs across alleles and distinct motifs associated with different peptide lengths. By integrating these data with transcript abundance and peptide processing, we developed HLAthena, providing allele-and-length-specific and pan-allele-pan-length prediction models for endogenous peptide presentation. These models predicted endogenous HLA class I-associated ligands with 1.5-fold improvement in positive predictive value compared with existing tools and correctly identified >75% of HLA-bound peptides that were observed experimentally in 11 patient-derived tumor cell lines.

The HLA genes are the most polymorphic across the human population, with more than 16,200 distinct class I alleles as of May of 2019 (refs. ^{1,2}). Short peptides (8–11-mers) bound to the diverse array of HLA class I molecules (HLA-A, -B, -C and -G) arise from intracellular proteins that are cleaved by the proteasome and peptidases before loading and display by surface HLA class I proteins to cytotoxic T-cell lymphocytes. Given the diversity in HLA binding preferences, an important question is whether one can accurately predict if a peptide is presented by a specific HLA allele. The accuracy of computational models that predict binding between peptides and HLA alleles, especially HLA-A and -B alleles, has been improving^{3–7}. In the field of cancer, these tools are now increasingly used in conjunction with next-generation DNA sequencing of tumors to identify immunogenic cancer neoantigens, which arise from tumor-specific somatic mutations. They have accelerated epitope discovery, as they enable experimental efforts to focus on a narrower list of epitopes with good predicted binding. However, even with widely used algorithms such as NetMHCpan^{3,8}, the numbers of falsely discovered binders increase once the predicted binding affinity decreases (that is, half-maximum inhibitory concentration > 100 nM)⁹. Furthermore, while these algorithms are

designed to predict the binding affinity of peptides to individual HLA molecules, the final step of antigen presentation, they do not account for intracellular availability of the peptide precursors or their processing by proteases. Finally, because previous research has focused on the few alleles highly expressed by Caucasian populations, existing algorithms have uneven accuracy in the prediction of epitopes binding to less common alleles in Caucasians, or those highly prevalent in other populations.

Detection and sequencing of HLA-bound peptides by liquid chromatography–tandem mass spectrometry (LC–MS/MS) has the unique advantage that information on endogenously processed and presented peptides from a cell can be directly learned. Our previous proof-of-concept study demonstrated that characterization of HLA-bound peptides eluted from a limited set of cell lines engineered to express single HLA alleles could reveal allele-specific peptide motifs and be used to train predictive algorithms for endogenous allele-specific peptide presentation⁴. Here, we expand our initial dataset of >24,000 peptides from 16 cell lines and identify and characterize 186,464 eluted peptides from 95 HLA-A, -B, -C and -G alleles. We included HLA-G peptidomes because this HLA is implicated in maternal–fetal tolerance and is also upregulated in many

¹Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA. ²Broad Institute of MIT and Harvard, Cambridge, MA, USA.

³Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA. ⁴Harvard Medical School, Boston, MA, USA. ⁵Department of Medicine, Brigham and Women's Hospital, Boston, MA, USA. ⁶Center for Patient Derived Models, Dana-Farber Cancer Institute, Boston, MA, USA.

⁷Division of Neuropathology, Brigham and Women's Hospital, Boston, MA, USA. ⁸Department of Cell Biology, Harvard Medical School, Boston, MA, USA.

⁹Immunitrack, Copenhagen, Denmark. ¹⁰Department of Pathology, Brigham and Women's Hospital, Boston, MA, USA. ¹¹Department of Computer Science, Metropolitan College, Boston University, Boston, MA, USA. ¹²Center for Cancer Immunology, Massachusetts General Hospital, Boston, MA, USA.

¹³These authors contributed equally: Siranush Sarkizova, Susan Klaeger. *e-mail: nhacohen@mg.harvard.edu; scarr@broadinstitute.org; cwu@partners.org; derin_keskin@dfci.harvard.edu

cancers^{10,11}. These data allow us to compare peptide length preferences and the spectrum of distinct and shared submotifs across HLA class I alleles, revealing the diversity and complexity of endogenous HLA ligands. Using this information, we trained allele- and length-specific and pan-allele-pan-length predictors, which identify 1.5-fold more peptides than conventional prediction tools when evaluating ligands directly detected by LC-MS/MS from 11 patient-derived tumor cell lines. The datasets of HLA binding peptides from mono-allelic cells and patient-derived tumors, as well as the prediction models (HLAthena) and interactive web tools, are all made publicly available.

Results

Systematic LC-MS/MS profiling of HLA class I ligands from mono-allelic cell lines. We engineered 79 cell lines expressing a single HLA class I allele by stably transfecting individual HLA-A, -B, -C or -G alleles into the HLA-null B721.221 cell line (Fig. 1a), adding to the 16 lines we previously reported⁴. Surface expression of the alleles was confirmed by flow cytometry (Fig. 1b, Supplementary Fig. 1a and Supplementary Table 1a). Altogether, the collection of 95 cell lines (31 HLA-A, 40 HLA-B, 21 HLA-C and 3 HLA-G) covered at least one allele in 95% of individuals worldwide for each HLA-A, -B and -C alleles^{12–14} (Supplementary Table 1b and Supplementary Note 1).

HLA-bound peptides for each engineered cell line were isolated by HLA immunopurification and analyzed by high-resolution LC-MS/MS, and sequences were identified by a ‘no-enzyme’ specificity database search at 1% false discovery rate (FDR) (Supplementary Data 1). We identified a median of 1,860 peptides per allele (range, 692–4,033), for a total of 186,464 peptides, after excluding nonspecifically bound peptides (Supplementary Table 1c–e, Supplementary Fig. 1b and see Methods). Most observed modifications, found in 12% of identified peptides, could be explained by sample processing artifacts such as methionine oxidation; adding carbamidomethylation of cysteine into the sample processing workflow recovered more cysteine-containing peptides (Supplementary Fig. 1c,d). HLA-bound peptides mapped to 10,649 human genes (≥ 2 peptides per gene), representing 91% of human gene products detected by LC-MS/MS in an extensively fractionated B721.221 proteome (≥ 2 peptides per gene), and 89% of transcribed genes (> 2 transcripts per million from RNA sequencing (RNA-seq)) (Fig. 1c, Supplementary Fig. 1e and Supplementary Note 3). The top 50 proteins with high HLA peptide coverage were large, highly abundant and consistently observed across HLA-A, -B and -C alleles (Supplementary Fig. 1f). Notably, 1,517 genes represented by HLA-bound peptides were not detected in either of the two expression datasets, suggesting that they had very low transcript and protein levels. These peptides had identification metrics comparable to the rest of the dataset, and are thus reliable identifications (Supplementary Fig. 1g). Peptides identified from sets of six alleles matched to patient genotypes¹⁵ amounted to 4,000–5,000 presented genes (Supplementary Fig. 1h). We conclude that all expressed proteins can undergo processing and presentation by HLA class I, a far higher proportion than previously appreciated¹⁶.

Our mono-allelic data nearly doubles the HLA ligands recorded in the Immune Epitope Database (IEDB)¹⁷, which holds 208,885 ligands from 157 human class I alleles (Supplementary Table 1f and see Methods). Peptides for 80 of 95 alleles are available in IEDB; however, 33 of 95 alleles have fewer than 100 known binders, which hinders reliable motif deduction and accurate prediction (Fig. 1d). For the 15 previously uncharacterized alleles, we identified 1,845 peptides on average (range, 693–4,022). We systematically assessed the length distribution, positional entropy, residue frequencies, binding motif and submotif clusters of HLA-bound peptides per allele (Supplementary Fig. 1i–l, Supplementary Note 4 and see Methods) and created an interactive website for data exploration

(<http://HLAthena.tools>). Altogether, these data and tools greatly expand the current knowledge of HLA class I-bound peptides.

Identification of HLA binding motifs and submotifs that are shared across alleles. Since the numbers of peptide identifications per allele were only weakly correlated with surface HLA levels, differential binding potential likely contributes to the variation in peptide numbers (Supplementary Fig. 1m,n). To better understand the basis for differential binding, we compared HLA alleles based on the motifs of their observed ligands and the physicochemical properties of binding pocket residues in the HLA protein. By computing pairwise correlations between the peptide binding motifs of each allele, we found groups of alleles sharing well-defined HLA-A and -B motifs (Fig. 2a, left, and Supplementary Table 2a). As expected, HLA alleles belonging to supertypes such as HLA-A*02 clustered together (Fig. 2a(ii)) as did split antigen serotypes such as HLA-B*54,55,56 and HLA-A*23,24 (Fig. 2a(i,iv))^{3,18,19}. There was minimal motif sharing outside of the dominant groups (mean motif correlation of each HLA-A allele to all other -A alleles, and each HLA-B allele to all other -B alleles was 0.28 and 0.25, respectively; Fig. 2b, left). In contrast, HLA-C motifs were more similar to each other (mean correlation 0.51), thus sharing more overlapping motifs, consistent with previous studies indicating that HLA-C (and HLA-G) alleles are more evolutionarily recent, with less divergence amongst alleles²⁰. The patterns of similarity revealed by binding motifs were mirrored by similarities in the HLA binding clefts, quantified by physicochemical properties of HLA residues in contact with the ligand (Fig. 2a,b, both right, Supplementary Table 2b and see Methods). To assess the agreement of the two approaches, for each allele, we counted the number of neighboring alleles in motif space analysis that were also proximal in pocket space (Supplementary Fig. 2a and see Methods). This correspondence maps the rules of ligand preferences onto HLA protein sequence and serves as the basis for creating pan-allele predictors that rely on transfer learning from characterized to uncharacterized alleles²¹.

To delineate allele similarity at finer granularity, we decomposed each aggregate motif per allele into submotifs by computing inter-peptide distances, projecting them onto two-dimensional space and clustering the peptides, obtaining 1,133 submotifs (≥ 20 peptides per submotif) across the 95 alleles. Distinct motifs were then identified by clustering the allele-specific submotifs, revealing 101 clusters (Supplementary Fig. 2b and see Methods). While half of the clusters containing HLA-A/B submotifs were contributed solely by HLA-A/B alleles, respectively (Fig. 2c), most of the HLA-C submotifs overlapped with submotifs from HLA-A and/or -B alleles (Fig. 2d and Supplementary Fig. 2c), consistent with HLA-A and -B alleles having more divergent structure than the evolutionarily ‘younger’ HLA-C alleles. Submotif overlap across alleles enables selection of a minimal set of epitopes covering an optimal set of alleles and thus individuals.

Length-specific differences in ligand preferences among HLA alleles and loci. Unbiased evaluation of length distributions revealed that 9-mers were dominant for all but two HLA-B alleles (which preferentially bound 8-mers) and that 8-mers were more common for HLA-B/C alleles while 10/11-mers were more common for HLA-A/B (Fig. 3a, Supplementary Table 3a and Table 2a). To systematically identify potential length-specific binding motifs, we compared 8/10/11-mers with 9-mers based on residue frequency and entropy at every peptide position (Fig. 3b and see Methods), and found 26 differences across HLA-A, -B and -C alleles (20 8-mer, two 10-mer and four 11-mer) out of 178 motifs with at least 100 identified 8/10/11-mer peptides (Fig. 3b–d and Supplementary Table 3b). The most notable changes in entropy were at position 5 for 8-mers compared with 9-mers (Supplementary Fig. 3a). This residue position has been implicated in structural changes of certain

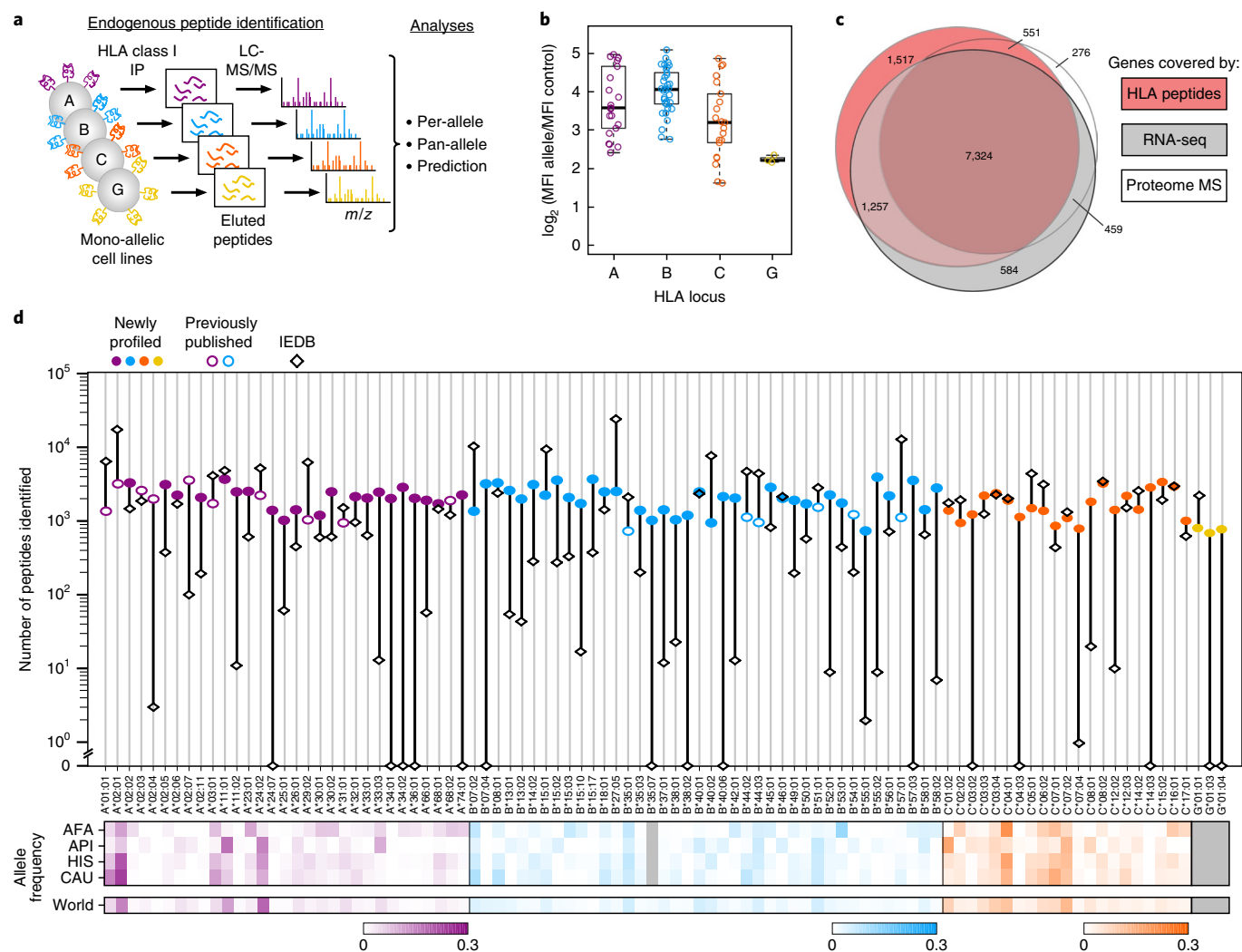


Fig. 1 | Mass spectrometric characterization of peptides eluted from HLA proteins in mono-allelic cell lines. **a**, Schematic of the experimental design: HLA-null B721.221 cells transfected to express a single HLA allele (31 HLA-A, 40 HLA-B, 21 HLA-C and three HLA-G) were subjected to HLA class I immunoprecipitation (IP) with W6/32 antibody from 50–300 million cells per allele followed by identification of eluted peptides by LC-MS/MS, to generate endogenous peptide binding data used to characterize allele-specific or pan-allele binding preferences and train neural network predictors of antigen processing and presentation. **b**, Surface expression of each transfected HLA allele was confirmed by flow cytometric detection against parental cells transfected with an empty vector (MFI, mean fluorescence intensity; $n = 21$ HLA-A, 34 HLA-B, 21 HLA-C, 3 HLA-G biologically independent samples; boxplots depict median intensity, the box contains 25–75% of the data, whiskers extend to lowest and highest values no further than $1.5 \times$ interquartile range; profiles of all lines in Supplementary Fig. 1a). **c**, Overlap of human genes represented by at least two HLA-associated peptides (pink), detected in RNA-seq (transcripts per million > 2 , light gray) or identified in deep proteome analysis (≥ 2 unique peptides, dark gray) of the B721.221 mono-allelic cells lines. **d**, Top, numbers of HLA-bound peptides identified per allele by MS-based profiling (circles: filled, generated data; open, previously reported⁴; diamonds: recorded in IEDB). Bottom, heatmaps of relative median population frequencies per allele across racial groups (AFA, African American; API, Asian or Pacific Islander; HIS, Hispanic; CAU, Caucasian) in the US population¹³ and worldwide. See also Supplementary Fig. 1.

HLA alleles upon binding as it allows embedding of these short peptides in the cleft^{22,23}. Selected peptides were confirmed as strong binders in *in vitro* binding assays, despite poor predicted affinity by NetMHCpan4.0-BA³ (Fig. 3e and Supplementary Fig. 3c). Collectively, these observations motivate a more explicit approach of modeling length-specific HLA binding characteristics.

Peptide-extrinsic properties vary per HLA and length. Since HLA-bound peptides captured from the cell surface reflect the cell-endogenous processes that shape the ligandome, we assessed whether HLA-A, -B, -C and -G ligands of different peptide lengths are preferentially derived from peptides with variable extrinsic properties. We found that HLA-C-bound peptides were biased toward higher expression and hydrophobicity (Fig. 3f, Supplementary

Table 3c and see Methods). HLA-C also showed a preference toward peptides with poorer proteasome cleavability scores, likely driven by the higher frequency of 8-mers, which were observed to have lower cleavage scores (Fig. 3f,g and Supplementary Fig. 3a). These observations agree with previous structural analyses that have reported more shallow HLA-C binding clefts²⁴, as higher abundance and elevated hydrophobicity of cognate peptides could compensate for decreased binding stability. We noted that HLA-G peptides had an even stronger bias toward lower cleavability scores, possibly due to the lack of HLA-G training data for the cleavability predictor, and this may suggest differential protease activity in shaping the HLA-G ligandome²⁵. Other differences with smaller effect sizes were also observed which altogether prompted us to model extrinsic properties per HLA loci in pan-allele predictors.

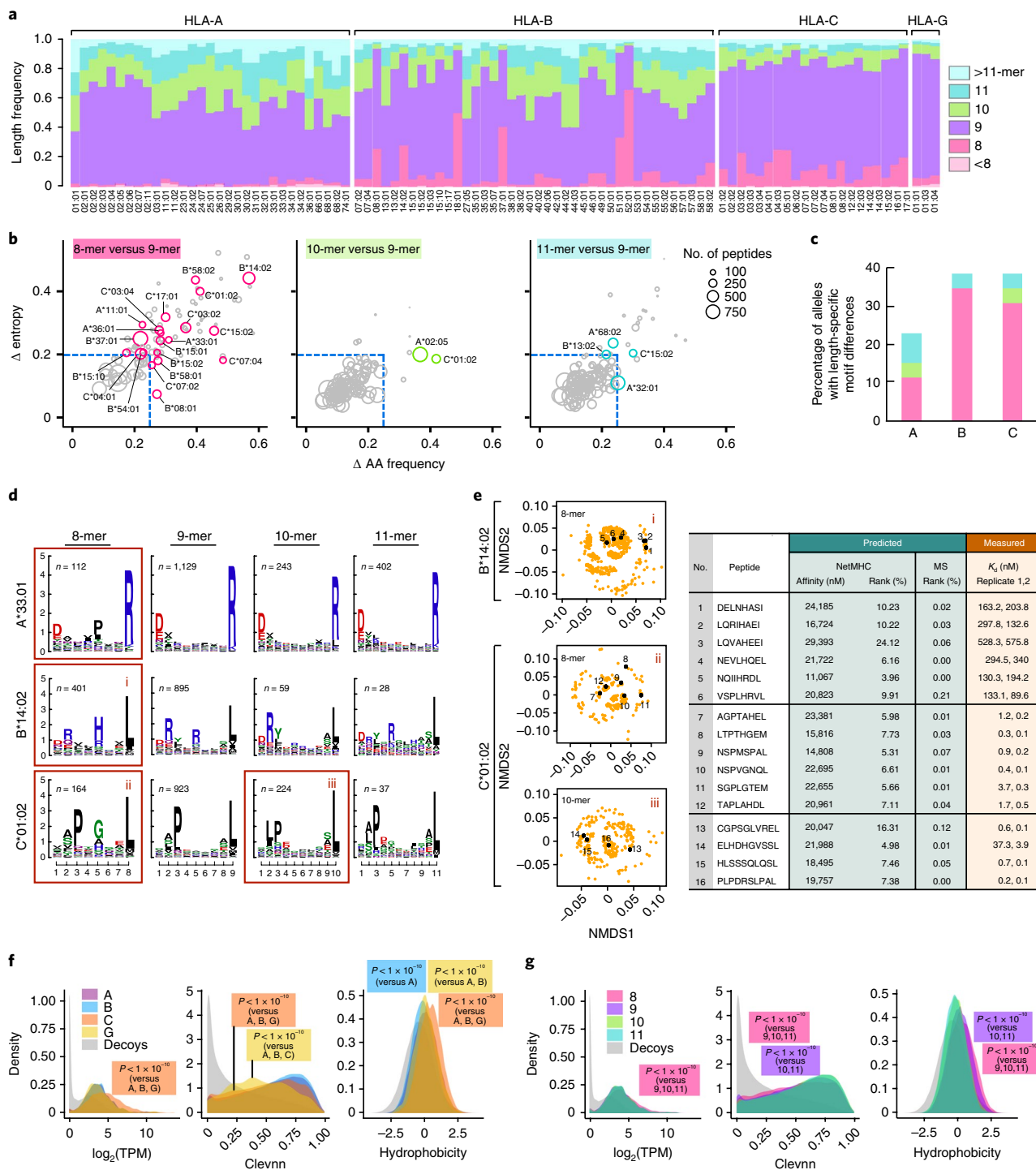


Fig. 3 | Mono-allelic data uncover length-specific HLA binding preferences. **a**, Frequencies of peptide lengths observed across alleles (8, pink; 9, violet; 10, green; 11, cyan). All but two HLA-B alleles preferentially present 9-mers. HLA-A alleles bind longer peptides more frequently than -B and -C alleles, while -B and -C alleles have a higher propensity for short peptides. **b**, The 8-, 10- and 11-mer binding motifs were compared with 9-mer motifs by dropping middle residues (positions 4, 5, 6 or 7 depending on the length) to create pseudo motifs of the same length (8-mers: pseudo 8-mer from 9-mers versus true 8-mer motif; 10- and 11-mers: pseudo 9-mer from 10- and 11-mers versus true 9-mer motif) and selecting the pseudo motif that was most similar to the corresponding true motif. The maximum differences amongst peptide residue positions between the 8-, 10- and 11-mer pseudo motifs and the corresponding true motifs in amino acid frequency (x axis) and entropy (y axis) are shown. Circle size reflects number of peptides; dashed lines indicate cutoff values. Circles in color and labeled denote alleles with >100 peptides with change in amino acid frequency or entropy greater than the selected cutoffs (absolute difference in residue frequency with the true motif of >0.25 or an absolute difference in entropy of >0.2 at any position). **c**, Percentage motif changes within each HLA type colored by length. **d**, Length-dependent log plots for A*33:01, B*14:02 and C*01:02; red boxes outline the changing motifs. **e**, Experimental validation of selected peptides (indicated with black dots on the NMDS plots) by in vitro binding assays compared with their predicted scores by NetMHCpan4.0-BA³ and MS models. **f, g**, Expression, predicted cleavability (cleavnn) and hydrophobicity stratified by HLA loci (**f**, $n = 95$ alleles; 31 HLA-A, 40 HLA-B, 21 HLA-C and 3 HLA-G, 1×10^6 decoys) and peptide length (**g**, $n = 12,970$ 8-mers, 111,898 9-mers, 29,956 10-mers, 18,202 11-mers; all comparisons Welch's two sample t-test, two-sided, provided in Supplementary Table 3c). See also Supplementary Fig. 3. AA, amino acid; NMDS, non-metric multidimensional scaling; TPM, transcripts per million.

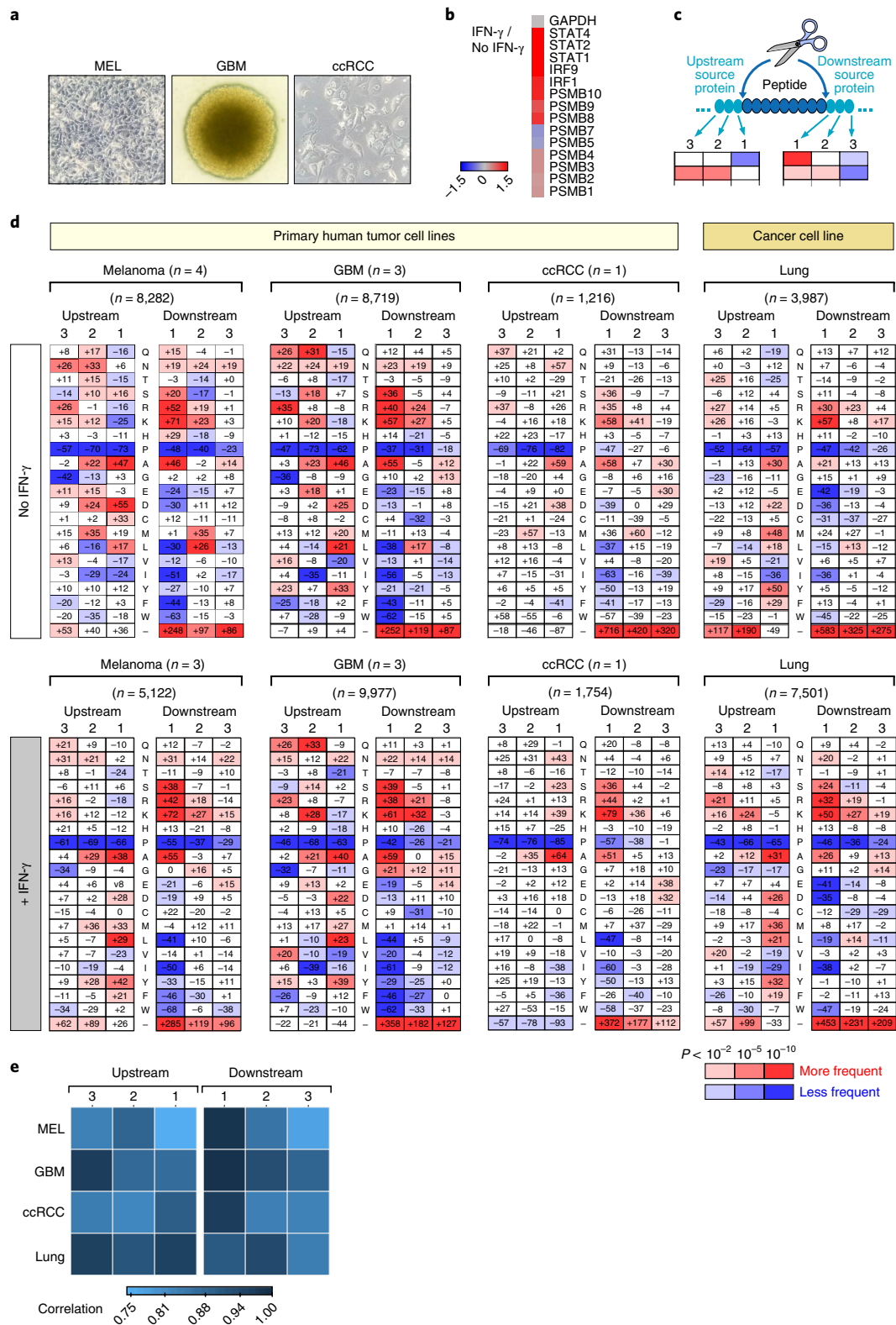


Fig. 4 | Proteasomal and peptidase shaping of the HLA-associated peptidome. a, Three types of primary tumor cell lines (MEL, GBM and ccRCC) used to identify HLA-associated peptidomes. **b**, Changes in relative protein abundance of proteasomal subunits and IFN- γ -inducible genes in patient-derived GBM cells with or without IFN- γ treatment based on MS proteome analysis. **c**, Schematic of cleavage signature analysis. **d**, Peptide processing signatures of HLA ligands presented by primary tumor and cancer cell lines at baseline (top, $n = 4$ MEL, 3 GBM, 1 ccRCC, 1 lung, biologically independent samples) and following IFN- γ treatment (bottom, $n = 3$ MEL, 3 GBM, 1 ccRCC, 1 lung, biologically independent samples), showing overrepresented (red) or underrepresented (blue) amino acid residues or protein N or C termini (indicated by ‘-’) upstream and downstream of the HLA peptide. The number in each cell denotes percentage change over a background decoy set; color intensity indicates significance (see key, chi-squared test, d.f. = 1). **e**, Heatmap of correlations between the processing preferences in untreated and IFN- γ -treated samples at upstream and downstream positions. Signatures for peptides from the IFN- γ -treated cells correlate well with peptides eluted from untreated cells, suggesting minimal to no difference between the two patterns (sample sizes as in **d**; Spearman’s rank correlation). See also Supplementary Fig. 4.

versus a set of decoys drawn from the proteome, controlling for HLA motif biases (Fig. 4c and Supplementary Note 5). As before⁴, in the untreated dataset (Fig. 4d, top), we observed an enrichment for A, K, S and R at downstream positions as well as for peptides derived from protein N or C termini (indicated by ‘-’). Upstream residues R and K were not enriched after removing potential tryptic peptide contaminants. Proline was depleted at both peptide termini in all samples, likely due to steric hindrance. Acidic residues (E, D) and certain hydrophobic residues (I, L, F, W) were under-represented downstream of HLA-associated peptides. Proteasomal signatures of untreated and IFN- γ -treated samples were strongly correlated (Fig. 4d,e; Spearman’s $\rho > 0.76$), suggesting that immunoproteasome activation has minimal impact on the processing of HLA-presented ligands in malignant cells.

Generation and performance of allele-and-length-specific and pan-allele-pan-length models. We previously reported that multivariate models incorporating endogenous HLA presentation descriptors, such as transcript abundance and likelihood of protease cleavage, outperform affinity-trained predictors⁴. With our extended dataset of 95 alleles, we evaluated the predictive contribution of these variables (Supplementary Fig. 4a,b and Supplementary Notes 6 and 7) along with gene presentation bias and translation quantification via ribosome profiling. Presentation bias quantifies the discrepancy between expected and observed number of mass spectrometry (MS)-identified peptides per gene in a given set of samples and can boost recall of lowly expressed peptides⁷, while ribosomal profiling (Ribo-seq) captures actively translated messenger RNA molecules and could provide a more accurate proxy for peptide precursor abundance. To evaluate predictive power, we constructed evaluation datasets with 1:999 ratio of observed binders to random genomic peptides and considered the fraction of true hits scoring within the top 0.1% (that is, positive predictive value (PPV)). The MS models trained on peptide sequence features alone achieved an average PPV across the 95 alleles of 47% (Fig. 5a and Supplementary Table 5a). This PPV strongly correlated with the number of decoys with scores higher than 50% of hits, which fit the binding motifs, suggesting that such decoys could be real binders and hence artificially decrease PPV (Supplementary Fig. 5). Integrating RNA-seq, as a proxy for peptide abundance, boosted PPV to 60%, while protein abundance achieved 54%. Combining RNA-seq with Ribo-seq reached a PPV of 61% (data not shown). Protein presentation bias and cleavability were the next most predictive variables, adding 2.9% and 1.5% to PPV. Based on these results, we trained prediction models that integrate intrinsic peptide features (MSintrinsic, or MSi) with extrinsic properties: cleavability (C), expression (E) and gene presentation bias (B).

The observed length-specific binding motifs (Fig. 3) in conjunction with the high frequency of non-9-mer presentation for some alleles motivated the generation of length-specific binding predictors, trained exclusively on ligands of specific lengths (8-, 9-, 10- or 11-mers), without ‘borrowing’ information from 9-mers (Fig. 5b, left, and see Methods). Model performance was compared against the most recent version of NetMHCpan, 4.0 (ref. ³), which incorporates training data from binding affinity (BA) and MS-sequenced eluted peptides (EL; including our previously published 16 alleles), as well as against MHCflurry⁵ and MixMHCpred2 (ref. ⁶). We found an average improvement in the MS-based sequence-only models (MSi) across lengths of 2.2-, 1.9-, 1.5- and 1.2-fold compared with MHCflurry (for overlapping alleles), NetMHCpan4.0-BA, NetMHCpan4.0-EL and MixMHCpred (for overlapping alleles), respectively (Fig. 5c and Supplementary Table 5a). Models that additionally integrated cleavability (MSiC) added +2.6 percentage points to the PPV achieved by MSi; cleavability and expression (MSiCE) an additional +9.3; and cleavability, expression and gene presentation bias (MSiCEB) a further +1.1 PPV. Length-specific

models for 8-, 10- and 11-mers outperformed the corresponding non-length-specific models currently used (Fig. 5c) with average increases of +15, +18, +26 and +27 PPV for MSi, MSiC, MSiCE and MSiCEB over NetMHCpan4.0-EL, respectively, and +7, +9, +17 and +18 PPV over MixMHCpred. The largest benefits were observed for 8-mer models, which was expected since 8-mer motifs were most different from 9-mer motifs (Fig. 3b, left). Ultimately, MSiCEB achieved 2.7-, 2.4-, 1.8- and 1.5-fold improvements compared with the four benchmark algorithms.

To enable prediction for any HLA allele (beyond our MS dataset), we built a pan-allele-pan-length model (panMSintrinsic or panMSi; Fig. 5b, right). Although the performance of our pan models was on average highly comparable to our nonpan models (mean and median differences of -2% PPV) and improvements over the nonpan models were observed for over 35% of allele-length combinations (Supplementary Table 5a), we also noted several cases with considerable decrease in predictive power. The subpar performance of pan models largely coincided with alleles for which non-9-mer motifs were different from the 9-mer motif (Fig. 3b). Compared with previous algorithms, the largest gains in PPV were observed for poorly characterized alleles (+20% PPV for HLA-C and +38% for HLA-G for MSi against NetMHCpan4.0-EL as the best-scoring pan-allele benchmark algorithm), but gains were also observed for all other alleles (+12% for HLA-A and +14% for HLA-B), even when only the 16 previously profiled alleles were considered (+9% for HLA-A and +5% for HLA-B).

We proposed PPV at 0.1% of the top-scoring peptides as a more suitable metric to evaluate HLA presentation predictors⁴ because of the importance of assessing true positive rates at the realistic 0.1% prevalence of binders (see Methods). A recent study quantified model performance using a different version of PPV, where the fraction of true positive calls out of all predictions necessary to retrieve 40% of the true binders is calculated, and the hits:decoys ratio in the evaluation set was modified from 1:999 to 1:10,000 (ref. ⁷) (that is, ‘PPV at 40% recall’). Due to these differences, the two metrics are not directly comparable. Using PPV at 40% recall, we found that MSi outperforms MHCflurry, NetMHCpan-BA, NetMHCpan-EL and MixMHCpred by 5-, 4-, 2- and 1.3-fold, while MSiCEB achieved 12-, 10-, 6- and 3-fold improvements (Fig. 5d and Supplementary Table 5b). Finally, we observed similar gains in PPV compared with MHCflurry, NetMHCpan and MixMHCpred when evaluating an external dataset of HLA-C and -G binders²⁸, although MixMHCpred performed better on lengths other than 9 with the caveat that this dataset was used in its training (Fig. 5e and Supplementary Table 5c). In summary, we observed 1.5–2.7 \times improvements in PPV (at the top 0.1% of the dataset) compared with existing predictors, which corresponds to 3–12 \times gains at 40% recall.

Motif complexity and motif abundance largely explain PPV variability. We observed that some allele preferences were harder to learn than others, with 9-mer-specific model (MSi) PPV values ranging from 37% to 68%. This variation was not readily explained by the amount of training data available, as model performance plateaus at several hundred peptides⁴ and >375 9-mer peptides were identified for all alleles. Since PPVs after the addition of endogenous features (MSiCEB) were strongly correlated with PPVs achieved with the simple model (Pearson’s correlation = 0.92, $P < 2.2 \times 10^{-16}$), and since we observed differences in the allele-specific motifs and submotifs (Supplementary Figs. 1d and 2b), we posited that PPV variability is driven by differential complexity in the peptide repertoire of each allele. To model complexity, we summed the entropy along each peptide position, to test whether higher information content implies more easily learned motifs. Similarly, we considered all submotifs identified per allele by summing positional entropies over the submotifs, each weighted by the number of supporting peptides,

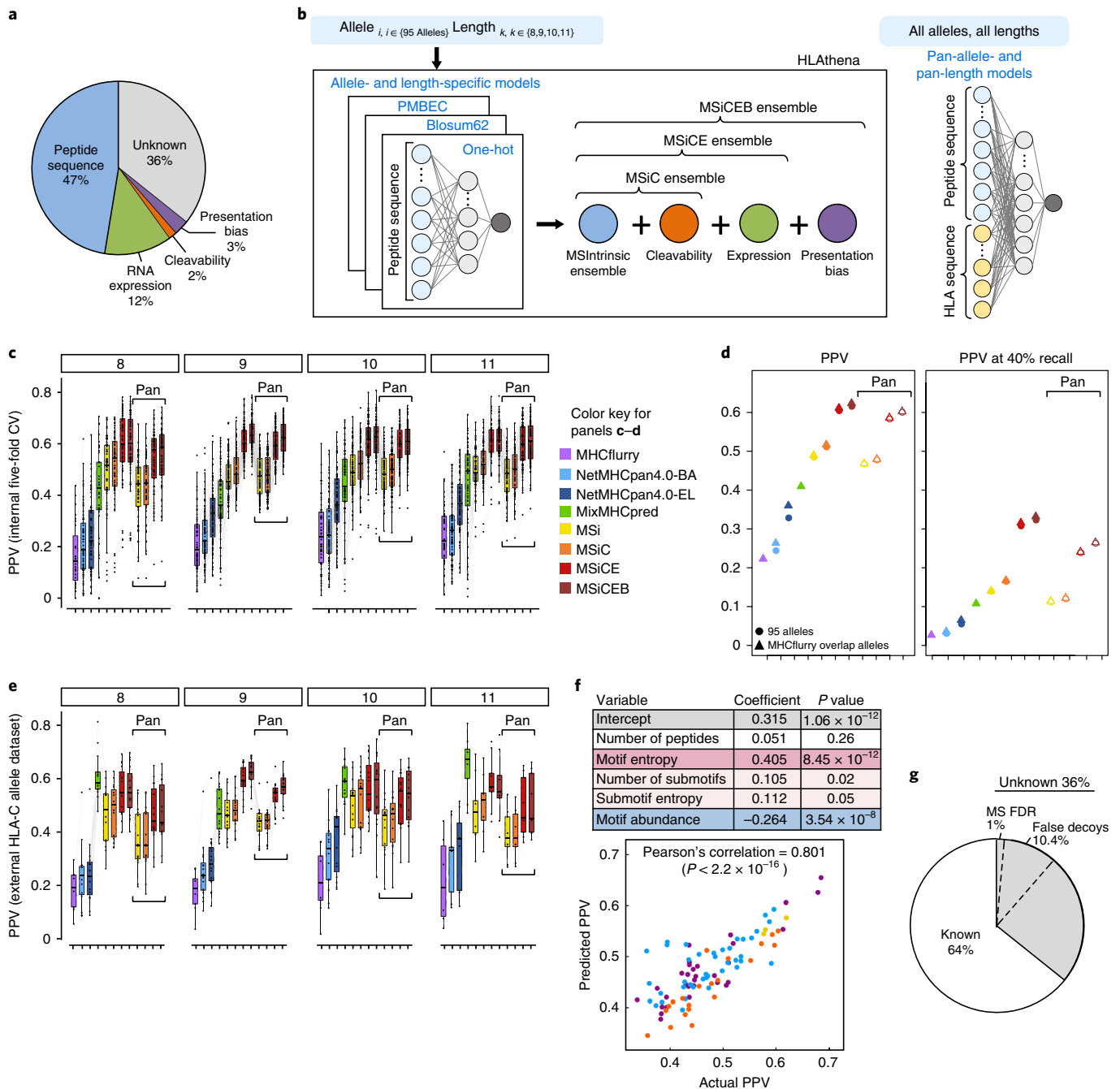


Fig. 5 | Generation and evaluation of allele-and-length-specific and pan-allele-pan-length MS-based models on mono-allelic data. a, Incremental contributions of predictor variables (peptide binding, transcript expression, cleavability and gene presentation bias) to PPV as the most informative variables are added one at a time (analysis performed for 9-mer peptides). **b**, Cartoon schematic of the neural networks used to generate allele-and-length-specific and pan-allele-pan-length predictive models. **c**, Models are evaluated based on their ability to score MS-observed binders in the top 0.1% amongst a 999-fold excess of random decoys (PPV). Shown are five-fold cross-validation (CV) PPVs across each of the $n = 95$ HLA alleles (gray lines) achieved by MHCflurry (available overlapping alleles = 31), NetMHCpan4.0-BA, NetMHCpan4.0-EL, MixMHCpred (available overlapping alleles = 72) and MS-informed models (boxplots depict median PPV, the box contains 25–75% of the data, whiskers reach to lowest and highest values no further than 1.5 × interquartile range). **d**, Average PPVs across the internal dataset of 95 alleles and four lengths for state-of-the-art and MS-based models resulting in a 1.5–2.7-fold improvement in PPV in an evaluation dataset with 1:999 hit:decoy ratio, or 3–12-fold improvement in PPV at 40% recall with 1:10,000 hit:decoy ratio. **e**, Model evaluation as in **d** on an external dataset of HLA-C-presented peptides identified by MS²⁸. **f**, Correlation of actual PPVs achieved by the allele-specific 9-mer MSi models versus PPVs predicted by a multivariate linear regression fit, with variables and their respective effect sizes and significances tabulated ($n = 95$). **g**, The negative contribution of motif abundance to PPV (that is, negative regression coefficient) suggests that ~10.4% of ‘unknown’ PPV (estimated as the average motif abundance scaled by the coefficient) can be attributed to false decoys present in the negative set, which artificially decreases PPV. Similarly, 1% of unexplained PPV could be due to false-positive identifications by MS at the 1% FDR threshold ($n = 95$). See also Supplementary Table 5.

and the number of submotifs normalized by the total number of peptides. Finally, we approximated motif abundance by the natural frequency of amino acids in the human genome underlying each binding motif. When motifs are more likely to occur by chance, more peptides from the random decoy set are real binders, leading to decreased PPV. To assess whether these variables are predictive of PPV, we used a multivariate linear fit, controlling for the size of the training data (see Methods), and found a strong correlation between predicted and actual PPV (Fig. 5f and Supplementary Table 5d). The number of peptides per allele was not predictive, while the entropy of the main motif, and the number and entropy of submotifs, were positively associated with PPV, whereas motif abundance strongly negatively contributed to PPV. Based on the model coefficients, we estimated that motif abundance could be responsible for ~10.4% of the unexplained PPV, although motif abundance likely underestimates the rate of undiscovered binders amongst the decoys. An additional 1% could be due to false-positive MS identification at 1% FDR (Fig. 5g). Overall, these findings suggest that limitations in learning motifs can be in large part attributed to motif complexity and abundance.

Peptides proposed to be derived from proteasomal splicing have poor predicted binding scores. Peptides derived from proteasomally ligated fragments (spliced peptides) have been recently proposed as a major component of the HLA ligandome^{29,30}. Since our collection of mono-allelic data covered the HLA alleles evaluated in those studies, we compared the binding potentials of reported linear and proposed spliced peptide sets using our *de novo* predictors. Consistent with previous analyses^{31,32}, we found that the majority of reported spliced peptides had poor predicted binding: although 81% of canonical linear peptides described by Liepe et al.²⁹ had an HLA binding likelihood score >0.75, only 28% of spliced peptides passed the same threshold (Supplementary Fig. 6a, left, and see Methods). Similar results were obtained for peptides described by Faridi et al.³⁰: 84% linear vs. 36% cis- and 37% trans-spliced (Supplementary Fig. 6a, right). While spliced peptides have been reported to make up 30% of the HLA class I peptidome²⁷, our computational results suggest that no more than 11% (37% of 30%) of presented HLA ligands could be derived from spliced peptides, a number previously shown to be further diminished by factors such as ambiguity in peptide spectral matches and variability in sequence database search strategies^{31,33}.

Leading sensitivity performance of MS-trained integrative algorithms. To evaluate the utility of our predictive models for clinical samples, we assessed their sensitivity to retrieve HLA-bound peptides observed in patient-derived tumor cell lines. To this end, we (1) used LC-MS/MS to identify 51,531 HLA-associated peptides from 11 tumor samples (three chronic lymphocytic leukemia (CLL), one ovarian (OV), three GBM, four MEL) and used external peptide datasets from four MEL³⁴ and 27 OV³⁵ tumors; (2) predicted the likelihood that each observed peptide is presented per HLA allele per sample; and (3) compared the proportions of correctly predicted peptides amongst a large set of random genomic peptides relative to four previous tools. Observed ligands that scored better than 99.9% of random peptides for at least one allele (top 0.1 percentile) were considered correct identifications (Fig. 6a). Our mono-allelic dataset covered 50 of 57 unique HLA alleles found amongst the 42 patient samples used in the evaluation. For covered alleles/lengths, predictions were made with our allele-and-length-specific models, while missing alleles were scored by our pan-allele-pan-length predictors. Across malignancies, we consistently observed a higher proportion of observed peptides predicted by the MS-based models compared with existing algorithms (Fig. 6b, Supplementary Fig. 6 and Supplementary Table 6a,b). At the 0.1 percentile threshold, 26% of observed peptides were recalled by MHCflurry (for MHCflurry

and MixMHCpred we were only able to assess supported alleles), followed by 31%, 46% and 49% predicted by NetMHCpan4.0-BA, NetMHCpan4.0-EL and MixMHCpred, respectively, compared with 56%, 60%, 77% and 78% predicted by MSi, MSiC, MSiCE and MSiCEB, on average across all samples. This constitutes a 1.6-fold improvement in recall.

Allele contribution to peptide presentation varies by individual. Since MS-detected epitopes were assigned to the best-scoring HLA allele(s), this allowed calculation of allele frequency among the presented peptides (Fig. 6c and Supplementary Table 6c). Notably, 3% of peptides on average were uniquely assigned to HLA-C alleles and an additional 6% of peptides were compatible with an HLA-C allele jointly with other alleles, thus suggesting that HLA-C has the potential to harbor neoantigens. In all six samples for which we profiled HLA-associated peptidomes ± IFN- γ treatment, we observed a shift toward HLA-B presentation, consistent with HLA-B having two IFN- γ -inducible promoters elements^{27,36}. We further examined the HLA-B allele combinations of each and found elevated presentation for alleles with both tryptic and chymotryptic C-terminal preferences. This suggests that HLA-B upregulation could be in part responsible for a shift in presentation from tryptic-like to chymotryptic-like peptides that was observed in a cell line with a C-terminal chymotryptic-like HLA-B motif³⁷. Finally, the contribution of each allele to the antigen repertoire varied by patient, suggesting that MS profiling of tumor-presented epitopes can reveal allele-specific utilization and further guide peptide vaccine selection.

Discussion

We demonstrate the superior performance of HLA class I predictors trained on large-scale data of peptides eluted from cellular HLA proteins, consistent with growing appreciation of MS-derived datasets as a basis for epitope prediction algorithms^{3,5-7}. Using an optimized experimental workflow, we eluted peptides from immunoprecipitated HLA proteins and used high-performance MS followed by a refined database searching approach to build a large dataset of HLA ligands eluted from single HLA-expressing cell lines. The resulting collection of >185,000 peptides from 95 alleles greatly expands available knowledge of the human HLA-associated peptidome¹⁷, such that at least 95% of individuals worldwide have at least one of their A, B and C alleles covered. The data are publicly available, thus providing a valuable resource for researchers. To facilitate access, we have implemented a web-based tool for data visualization, interactive exploration and prediction.

Our large dataset enables more comprehensive insights into the rules of peptide presentation by HLA-A, -B, -C and -G alleles, each of which impacted our model design. First, we ascertained that peptide presentation does access the entire proteome for potential sources of antigen, in contrast to previous reports¹⁶ that relied on a small number of peptides. Second, our analysis revealed 101 binding submotifs, many of which were shared amongst the 95 HLA alleles. We observed strong similarity in physicochemical features of the HLA-C alleles along with their greater promiscuity in binding peptides, compared with the more divergent HLA-A and -B alleles. Moreover, HLA-C alleles only rarely had unique submotif clusters that were not also shared with HLA-A and -B alleles, consistent with their recent evolutionary history²⁰. We speculate that this may increase competition with HLA-A and -B alleles for peptides and may explain our observation that HLA-C epitopes originate from more highly expressed genes. Third, we detected not only differences in length distribution, but also that ~10% of alleles displayed length-based epitope preferences. Altogether, the detailed knowledge gained from our extensive dataset enabled us to generate allele-and-length-specific and pan-allele-pan-length prediction models that we demonstrate to outperform state-of-the-art

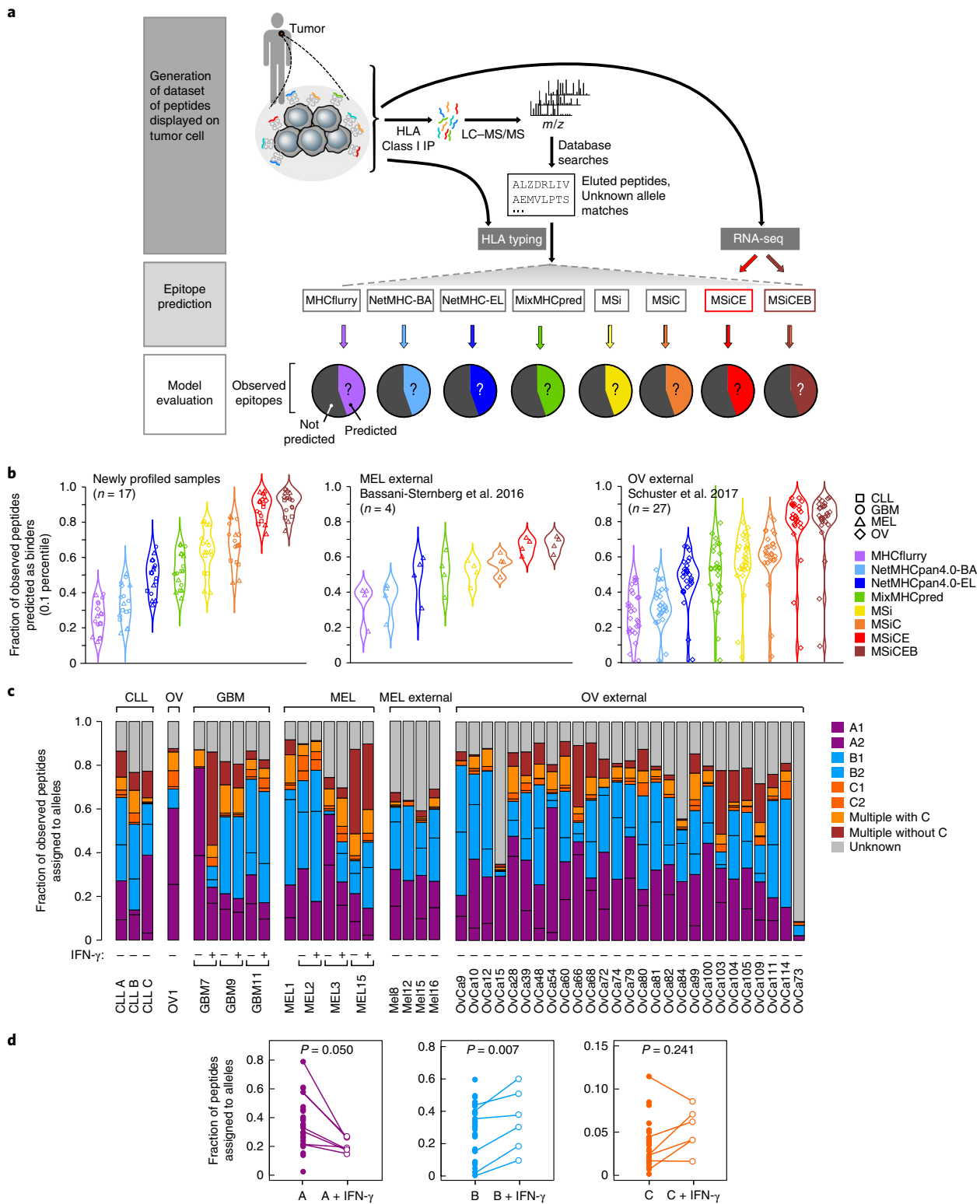


Fig. 6 | Integrative MS-informed models more accurately predict peptides directly observed on primary tumor cells. a, Schematic of data generation and model evaluation: peptides displayed on primary tumor specimens are isolated and sequenced by MS, the HLA alleles of the patient sample are clinically typed and matched RNA-seq data are generated; each observed epitope is evaluated for binding against each of the unique HLA alleles in the sample, and predictions that are better than 0.1% scores within a large decoy set are considered binders and assigned to the corresponding allele; the performance of eight algorithms is evaluated as the fraction of observed binders that are successfully predicted as binders. **b**, MS-based predictor ranks MS-detected peptides better than NetMHCpan and MHCflurry. Internal data on 11 patients ($n = 3$ CLL (squares), 1 OV (diamonds), 3 GBM (circles), 4 MEL (triangles), all biologically independent) and external data ($n = 4$ MEL and 27 OV, biologically independent samples)^{34,35}. **c**, Peptides were assigned to alleles in each sample based on the best-scoring peptide-allele combination. Allele contribution to peptide presentation varies per tumor, per IFN- γ treatment and per individual. **d**, Fraction of peptides contributing per allele type \pm IFN- γ ($n = 6$, biologically independent samples). Peptide presentation on HLA-B increases with IFN- γ stimulation (Wilcoxon signed-rank tests). See also Supplementary Fig. 5.

algorithms, especially for understudied alleles or those with length-specific preferences. Fourth, while IFN- γ signaling broadly modulates gene expression and thus alters the genes HLA ligands derive from, we did not observe prominent differences in cleavage preferences across primary tumor cells of various lineages when exposed to IFN- γ . This is consistent with the expression of both constitutive and immunoproteasome subunits in cancers and supports the application of a unified cleavability predictor.

The improved performance of the HLAthena models can be attributed to several factors: (1) our models are trained exclusively on MS data of eluted peptides from mono-allelic cell lines; (2) we integrate several critical endogenous features, such as peptide cleavage and gene expression; (3) our rich dataset reliably captures not only allele- but also length-specific motifs, widely covering the space of HLA binding preferences; and (4) we preferentially predict with allele-and-length-specific models for their demonstrated accuracy over pan-allele-pan-length predictors which are employed only for uncharacterized alleles.

Despite improvements in epitope prediction, we recognize that further innovations are required if we are to achieve near-perfect accuracy. We offer evidence that allele complexity and motif abundance may partially drive the observed variability in prediction power across alleles. The former implies a benefit in obtaining even larger training datasets, while the latter necessitates techniques to determine nonbinders at large scale to collect reliable true negative datasets against which to evaluate model performance. Other innovations that could boost prediction fidelity include increased LC-MS/MS instrument sensitivity and better de novo HLA peptide identification methods. For predicting peptide presentation in tumor cells, better accuracy can be achieved by taking into account the uneven allele utilization and weighting predictions accordingly. While we include expression as a variable in our predictors, it is important to note that RNA-seq of tumors may not be representative of all clones and usually includes nonmalignant cells that obfuscate tumor gene expression. Finally, we emphasize that our models do not predict whether the HLA-presented peptides can interact with the T-cell receptors in an individual, a problem that remains unsolved.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41587-019-0322-9>.

Received: 28 May 2019; Accepted: 24 October 2019;

Published online: 16 December 2019

References

- Lefranc, M.-P. et al. IMGT*, the international ImMunoGeneTics information system* 25 years on. *Nucleic Acids Res.* **43**, D413–D422 (2015).
- Robinson, J. et al. The IPD and IMGT/HLA database: allele variant databases. *Nucleic Acids Res.* **43**, D423–D431 (2015).
- Jurtz, V. et al. NetMHCpan-4.0: improved peptide–MHC class I interaction predictions integrating eluted ligand and peptide binding affinity data. *J. Immunol.* **199**, 3360–3368 (2017).
- Abelin, J. G. et al. Mass spectrometry profiling of HLA-associated peptidomes in mono-allelic cells enables more accurate epitope prediction. *Immunity* **46**, 315–326 (2017).
- O'Donnell, T. J. et al. MHCflurry: open-source class I MHC binding affinity prediction. *Cell Syst.* **7**, 129–132.e4 (2018).
- Gfeller, D. et al. The length distribution and multiple specificity of naturally presented HLA-I ligands. *J. Immunol.* **201**, 3705–3716 (2018).
- Bulik-Sullivan, B. et al. Deep learning using tumor HLA peptide mass spectrometry datasets improves neoantigen identification. *Nat. Biotechnol.* **37**, 55–63 (2018).
- Nielsen, M. & Andreatta, M. NetMHCpan-3.0: improved prediction of binding to MHC class I molecules integrating information from multiple receptor and peptide length datasets. *Genome Med.* **8**, 33 (2016).

- Rajasagi, M. et al. Systematic identification of personal tumor-specific neoantigens in chronic lymphocytic leukemia. *Blood* **124**, 453–462 (2014).
- de Kruijff, E. M. et al. HLA-E and HLA-G expression in classical HLA class I-negative tumors is of prognostic value for clinical outcome of early breast cancer patients. *J. Immunol.* **185**, 7452–7459 (2010).
- Zhang, R.-L. et al. Predictive value of different proportion of lesion HLA-G expression in colorectal cancer. *Oncotarget* **8**, 107441–107451 (2017).
- Dawson, D. V., Ozgur, M., Sari, K., Ghanayem, M. & Kostyu, D. D. Ramifications of HLA class I polymorphism and population genetics for vaccine development. *Genet. Epidemiol.* **20**, 87–106 (2001).
- Gragert, L., Madbouly, A., Freeman, J. & Maiers, M. Six-locus high resolution HLA haplotype frequencies derived from mixed-resolution DNA typing for the entire US donor registry. *Hum. Immunol.* **74**, 1313–1320 (2013).
- Solberg, O. D. et al. Balancing selection and heterogeneity across the classical human leukocyte antigen loci: a meta-analytic review of 497 population studies. *Hum. Immunol.* **69**, 443–464 (2008).
- Ott, P. A. et al. An immunogenic personal neoantigen vaccine for patients with melanoma. *Nature* **547**, 217–221 (2017).
- Pearson, H. et al. MHC class I-associated peptides derive from selective regions of the human genome. *J. Clin. Invest.* **126**, 4690–4701 (2016).
- Vita, R. et al. The immune epitope database (IEDB) 3.0. *Nucleic Acids Res.* **43**, D405–D412 (2015).
- Sette, A. & Sidney, J. HLA supertypes and supermotifs: a functional perspective on HLA polymorphism. *Curr. Opin. Immunol.* **10**, 478–482 (1998).
- Robinson, J., Malik, A., Parham, P., Bodmer, J. G. & Marsh, S. G. E. IMGT/HLA database—a sequence database for the human major histocompatibility complex. *Tissue Antigens* **55**, 280–287 (2000).
- Parham, P. & Moffett, A. Variable NK cell receptors and their MHC class I ligands in immunity, reproduction and human evolution. *Nat. Rev. Immunol.* **13**, 133–144 (2013).
- Nielsen, M. et al. NetMHCpan, a method for quantitative predictions of peptide binding to any HLA-A and -B locus protein of known sequence. *PLoS ONE* **2**, e796 (2007).
- Rist, M. J. et al. HLA peptide length preferences control CD8⁺ T cell responses. *J. Immunol.* **191**, 561–571 (2013).
- Maenaka, K. et al. Nonstandard peptide binding revealed by crystal structures of HLA-B*5101 complexed with HIV immunodominant epitopes. *J. Immunol.* **165**, 3260–3267 (2000).
- Kaur, G. et al. Structural and regulatory diversity shape HLA-C protein expression levels. *Nat. Commun.* **8**, 15924 (2017).
- Celik, A. A., Simper, G. S., Hiemisch, W., Blasczyk, R. & Bade-Döding, C. HLA-G peptide preferences change in transformed cells: impact on the binding motif. *Immunogenetics* **70**, 485–494 (2018).
- Keskin, D. B. et al. Neoantigen vaccine generates intratumoral T cell responses in phase Ib glioblastoma trial. *Nature* **565**, 234–239 (2019).
- Javitt, A. et al. Pro-inflammatory cytokines alter the immunopeptidome landscape by modulation of HLA-B expression. *Front. Immunol.* **10**, 141 (2019).
- Di Marco, M. et al. Unveiling the peptide motifs of HLA-C and HLA-G from naturally presented peptides and generation of binding prediction matrices. *J. Immunol.* **199**, 2639–2651 (2017).
- Liepe, J. et al. A large fraction of HLA class I ligands are proteasome-generated spliced peptides. *Science* **354**, 354–358 (2016).
- Faridi, P. et al. A subset of HLA-I peptides are not genomically templated: evidence for cis- and trans-spliced peptide ligands. *Sci. Immunol.* **3**, eaar3947 (2018).
- Mylonas, R. et al. Estimating the contribution of proteasomal spliced peptides to the HLA-I ligandome. *Mol. Cell. Proteom.* **17**, 2347–2357 (2018).
- Rolfs, Z., Solntsev, S. K., Shortreed, M. R., Frey, B. L. & Smith, L. M. Global identification of post-translationally spliced peptides with neo-fusion. *J. Proteome Res.* **18**, 349–358 (2018).
- Rolfs, Z., Müller, M., Shortreed, M. R., Smith, L. M. & Bassani-Sternberg, M. Comment on 'A subset of HLA-I peptides are not genomically templated: evidence for cis- and trans-spliced peptide ligands'. *Sci. Immunol.* **4**, eaaw1622 (2019).
- Bassani-Sternberg, M. et al. Direct identification of clinically relevant neoepitopes presented on native human melanoma tissue by mass spectrometry. *Nat. Commun.* **7**, 13404 (2016).
- Schuster, H. et al. The immunopeptidomic landscape of ovarian carcinomas. *Proc. Natl Acad. Sci. USA* **114**, E9942–E9951 (2017).
- Girdlestone, J. Regulation of HLA class I loci by interferons. *Immunobiology* **193**, 229–237 (1995).
- Chong, C. et al. High-throughput and sensitive immunopeptidomics platform reveals profound interferon-mediated remodeling of the human leukocyte antigen (HLA) ligandome. *Mol. Cell. Proteom.* **17**, 533–548 (2018).

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature America, Inc. 2019

Methods

Generation of HLA-A, -B and -C single-allele cell lines. Single HLA allele-expressing complementary DNA vectors in a pcDNA-3 backbone were ordered from GenScript. The HLA class I-deficient B721.221 cell line was transfected with the HLA allele expression vectors using lipofectamine, as described previously⁴. Cell lines with stable surface HLA expression were generated first through selection using 800 $\mu\text{g ml}^{-1}$ G418 (Thermo Scientific), followed by enrichment of HLA-positive cells through up to two serial rounds of FACS and isolation using a pan-HLA antibody (W6/32; Santa Cruz) on a FACSAria II instrument (BD Biosciences). Priority was given to HLA alleles with lack of binding data in public databases or over 1% frequency in the US organ donor registry populations¹³.

Generation of primary human samples. All human tissues were obtained through Dana-Farber Cancer Institute- or Partners Healthcare-approved Institutional Review Board (IRB) protocols. Conditions for growth and in vitro propagation of MEL and GBM tumor cell lines and of monocyte-derived dendritic cells were described previously^{15,26}. Peripheral blood mononuclear cells (PBMCs) from patients with CLL were enriched for CD19-positive CLL tumor cells and were used in immunoprecipitation followed by mass spectrometry (IP/MS) analysis. Tumor specimens from patients with ccRCC were collected following informed consent for enrollment on a tissue collection research protocol approved by the Dana-Farber/Harvard Cancer Center Institutional Review Board. Surgically resected ccRCC tumor tissue was mechanically dissociated with scalpels, and then enzymatically dissociated using a mixture of collagenase D (Roche), Dispase (STEMCELL Technologies) and DNase I (New England BioLabs) at room temperature, and filtered through a 100- μm cell strainer using the sterile plunger of a syringe. Red blood cells were lysed using ammonium-chloride-potassium buffer (Gibco). The cell suspension was stained for viability (Zombie Aqua; BioLegend), anti-CD45 (BV605; BD Biosciences) and anti-carbonic anhydrase IX (PE; R&D Systems). Viable, CD45-negative, CAIX-positive tumor cells were isolated by FACS (BD FACSAria II cell sorter; BD Biosciences). Cells were cultured in a specialized growth medium consisting of OptiMEM GlutaMax media (Gibco), 5% fetal bovine serum, 1 mM sodium pyruvate (Gibco), 100 units per ml penicillin and streptomycin, 50 $\mu\text{g ml}^{-1}$ gentamicin, 5 $\mu\text{g ml}^{-1}$ insulin (Sigma) and 5 ng ml^{-1} epidermal growth factor (Sigma). Following successive passages, CAIX expression was confirmed by flow cytometry (anti-CAIX, PE-conjugated; R&D Systems) and by immunohistochemical analysis of a cell pellet. Ovarian cancer patient-derived cells were propagated within a xenograft model, which was generated by serial passaging of tumor cells from a patient with advanced ovarian cancer. These cells originated from solid tumor or pleural effusion (3 million cells per mouse) that was injected orthotopically in the abdominal cavity in NOD-SCID mice (8 weeks old, Jackson Laboratory). Tumor growth was monitored weekly by observing mice for signs of abdominal distension. Cells were collected 4 months after initial injection and banked for future experiments. For interferon stimulation, cultured cells were stimulated with 2,000 units per ml of IFN- γ (Peprotech) for 3 d and were used in IP/MS analysis.

For primary tumors and patient cell lines, HLA-peptide complexes were immunoprecipitated from 0.1–0.2 g tissue or up to 50 million cells. Solid tumor samples were dissociated using a tissue homogenizer (Fisher Scientific 150) and HLA complexes were enriched as described in the next section (HLA peptide enrichment and LC-MS/MS analysis).

HLA peptide enrichment and LC-MS/MS analysis. Soluble lysates from up to 50 million single HLA-expressing B721.221 cells and up to 0.2 g from tumor samples were immunoprecipitated with W6/32 antibody (sc-32235, Santa Cruz) as described previously⁴. Iodoacetamide (10 mM) was added to the lysis buffer to alkylate cysteines for 71 alleles (Supplementary Table 1c and Supplementary Data 2). Peptides of up to three immunoprecipitates for single HLA-expressing samples and up to four immunoprecipitates for tumor samples were combined, acid eluted either on StageTips or SepPak cartridges³⁴, and analyzed in technical duplicates using high-resolution LC-MS/MS on a QExactive Plus (QE+), QExactive HF (QE-HF) or Fusion Lumos (Thermo Scientific). For acquisition parameters see Supplementary Note 2.

HLA peptide identification using Spectrum Mill. Mass spectra were interpreted using the Spectrum Mill software package v6.1 pre-release (Agilent Technologies). Tandem MS (MS/MS) spectra were excluded from searching if they did not have a precursor sequence MH+ in the range 600–4,000, had a precursor charge > 5 or had a minimum of <5 detected peaks. Merging of similar spectra with the same precursor m/z acquired in the same chromatographic peak was disabled. Before searches, all MS/MS spectra had to pass the spectral quality filter with a sequence tag length > 2 (that is, minimum of four masses separated by the in-chain masses of three amino acids). MS/MS spectra were searched against a protein sequence database that contained 98,298 entries, including all UCSC Genome Browser genes with hg19 annotation of the genome and its protein-coding transcripts (63,691 entries), common human virus sequences (30,181 entries) and recurrently mutated proteins observed in tumors from 26 tissues (4,167 entries), as well as 259 common laboratory contaminants including proteins present in cell culture media and immunoprecipitation reagents. Mutation files for 26 tumor tissue types

were obtained from the Broad GDAC portal (gdac.broadinstitute.org). Recurrent mutations in the coding region within each of the 26 tumor types (frequency = 3 for stomach adenocarcinoma, uterine corpus endometrial carcinoma; frequency = 5 for adrenocortical carcinoma, pancreatic adenocarcinoma, MEL; frequency = 2 for rest) were included. MS/MS search parameters included: no-enzyme specificity; fixed modification: cysteinylolation of cysteine; variable modifications: carbamidomethylation of cysteine, oxidation of methionine and pyroglutamic acid at peptide N-terminal glutamine; precursor mass tolerance of ± 10 ppm; product mass tolerance of ± 10 ppm; and a minimum matched peak intensity of 30%. Variable modification of carbamidomethylation of cysteine was only used for HLA alleles that included an alkylation step (performed in 2017 or later). Peptide spectrum matches (PSMs) for individual spectra were automatically designated as confidently assigned using the Spectrum Mill autovalidation module to apply target-decoy-based FDR estimation at the PSM level of <1% FDR. Peptide autovalidation was done separately for each HLA allele with an auto thresholds strategy to optimize score and delta Rank1–Rank2 score thresholds separately for each precursor charge state (1 through 4) across all LC-MS/MS runs for an HLA allele. Score threshold determination also required that peptides had a minimum sequence length of 7, and PSMs had a minimum backbone cleavage score (BCS) of 5. BCS is a peptide sequence coverage metric and the BCS threshold enforces a uniformly higher minimum sequence coverage for each PSM, at least four or five residues of unambiguous sequence. The BCS score is a sum after assigning a 1 or 0 between each pair of adjacent amino acids in the sequence (max score is peptide length – 1). To receive a score, cleavage of the peptide backbone must be supported by the presence of a primary ion type for higher energy collisional dissociation (HCD): b, y or internal ion C terminus (that is, if the internal ion is for the sequence PWN then BCS is credited only for the backbone bond after the N). The BCS metric serves to decrease false positives associated with spectra having fragmentation in a limited portion of the peptide that yields multiple ion types. PSMs were consolidated to the peptide level to generate lists of confidently observed peptides for each allele using the Spectrum Mill Protein/Peptide summary module's Peptide-Distinct mode with filtering distinct peptides set to case sensitive. A distinct peptide was the single highest scoring PSM of a peptide detected for each allele. MS/MS spectra for a particular peptide may have been recorded multiple times (for example, as different precursor charge states, from replicate immunoprecipitates, from replicate LC-MS/MS injections). Different modification states observed for a peptide were each reported when containing amino acids configured to allow variable modification; a lowercase letter indicates the variable modification (C-cysteinylated, c-carbamidomethylated). These unfiltered peptide lists are provided as Supplementary Data 1.

MS/MS data from patient-derived cell lines were handled as described in the previous paragraph except that they were searched against the database mentioned above with further inclusion of patient-specific neoantigen sequences^{15,26}. These peptide lists are provided as Supplementary Data 2.

Filtering of MS-identified peptides. The list of LC-MS/MS-identified peptides was filtered to remove potential contaminants in the following ways: (1) peptides observed in negative controls runs (blank beads and blank immunoprecipitates); (2) peptides originating from the following species: 'STRSG', 'HEVBR', 'ANGIO432', 'ANGIO394', 'ANGIO785', 'ANGIO530', 'ACHLY', 'PIG', 'ANGIO523', 'RABIT', 'STAAU', 'CHICK', 'Pierce-iRT', 'SOYBN', 'ARMRU' and 'SHEEP' as common laboratory contaminants including proteins present in immunoprecipitation reagents (note that 'BOVINE' peptides derived from cell culture media were not excluded as they appear to have undergone processing and presentation and exhibit anchor residue motifs consistent with the human peptides observed for each allele); (3) peptides that were also identified in a typically digested full-proteome Jurkat sample; (4) peptides for which both the preceding and C-terminal amino acids were tryptic residues (R or K); (5) all possible leader peptides of lengths 8–11 from HLA-A, -B, -C and -G (first exon, $n = 410$) as they are likely to be presented by HLA-E; (6) peptides with negative *deltaFw/RevScore* as likely falling in the 1% false-positive MS identifications; (7) peptides identified for 20 or more of the 95 alleles ($n = 168$); and (8) peptides identified as potential C*01:02 contaminants in other alleles due to residual C*01:02 expression in B721.221 ($n = 383$). These peptides were identified by scoring all peptides with the allele-specific C*01:02 model and selecting those with predicted likelihood binding score > 0.95 that were also outliers for the allele (mean distance to the nearest 10 peptides > 90 percentile). A summary of counts of removed peptides is provided in Supplementary Table 1d. The filtered peptide lists are provided in Supplementary Table 1e (mono-allelic cell lines) and Supplementary Table 6a (patient cell lines).

IEDB data access and preparation. A curated set of previously identified HLA class I ligands was downloaded from IEDB at http://www.iedb.org/downloader.php?file_name=doc/mhc_ligand_full.zip (accessed on 14 June 2018)¹⁷ (related to Fig. 1). Records were filtered to MHC allele class = I, Epitope Object Type = Linear peptide and Allele Name consistent with human HLA class I nomenclature with four-digit typing (that is, regex: "HLA-[ABCG]*[0–9]{2};[0–9]{2}\$"). Peptides with quantitative measurements in units other than nM were removed and so were the following three assay types due to detected inconsistency between predicted (NetMHC 3.0) and actual affinity: 'purified MHC/direct/radioactivity/dissociation

constant KD , 'purified MHC/direct/fluorescence/half maximal effective concentration (EC_{50})' and 'cellular MHC/direct/fluorescence/half maximal effective concentration (EC_{50})'. A peptide was considered a binder if it had a quantitative affinity of <500 nM or qualitative label of 'Positive', 'Positive-High', 'Positive-Intermediate' or 'Positive-Low'. In cases where multiple records are available for the same {peptide, allele} pair, we either took the mean affinity or removed the peptide when the difference between the maximum and minimum log-transformed affinities ($1 - \log(nM)/\log(50,000)$) was >0.2. Similarly, peptides with multiple qualitative records were removed if the same number of positive and negative labels were found or kept otherwise. Our previously published data for 16 HLA-A and -B alleles were removed from the analysis of IEDB counts (PubMedID=28228285).

Allele similarity analyses. To assess which alleles are similar to each other, we considered similarity according to the observed binding motifs (peptide space) as well as similarity according to the HLA binding grooves (HLA binding pocket or HLA protein space) (related to Fig. 2). Similarity in peptide space was evaluated by tabulating the frequency of each of the 20 amino acids at each position along the peptide sequence (1 through 9) per allele, forming a vector of size $20 \times 9 = 180$. The pairwise correlations of these frequency vectors were used to quantify similarity (Fig. 2a).

To evaluate similarity in HLA binding pocket space, HLA protein sequences were downloaded from IMGT, the international ImMunoGeneTics information system <http://www.imgt.org> (http://hla.alleles.org/alleles/text_index.html, accessed 5 May 2018), and aligned. From the full HLA protein sequences, we selected positions that are in contact with the peptide (within a distance of 2) or positions that are most frequently mutated across alleles to represent the binding pocket: {7, 9, 13, 24, 31, 45, 59, 62, 63, 65, 66, 67, 69, 70, 71, 73, 74, 76, 77, 80, 81, 84, 95, 97, 99, 110, 114, 116, 118, 138, 143, 147, 150, 152, 156, 158, 159, 163, 167, 171} (Fig. 2b). The residue at each position of the binding pocket was featured by its amino acid physical properties encoded as ten Kidera Factors (available from R package Peptides v2.4, data(AAdata))³⁸ and three principal components derived from a dimensionality reduction of a large set of physicochemical properties³⁹. The full binding pocket was represented by the concatenated list of positions and allele similarity was assessed by Euclidean distance.

Given the two approaches to evaluating allele similarity (motif space and pocket space), we assessed how well they agree by identifying the closest neighbors for each allele in motif space and the closest neighbors for each allele in pocket space and counting how many of the former are also found in the latter (Supplementary Fig. 2a). The closest neighbors in motif space were considered to be alleles with correlation greater than 97.5% of all pairwise correlations. Analogously, closest neighbor alleles in pocket space were considered to be alleles within distance less than 97.5% of all pairwise distances.

Analysis of submotifs across alleles. Grouping the 9-mer peptides identified for each allele into submotifs (related to Fig. 2) (see *Peptide distance visualization and sub-clustering of binding motifs* section in Supplementary Note 4), identified 1,133 submotifs across the 95 alleles supported by at least 20 peptides (Fig. 2d and Supplementary Fig. 2c). To determine whether any of those submotifs are shared by two or more alleles, each submotif was represented as a vector of amino acid frequencies per peptide position (analogously to main motif representation in allele similarity analysis), projected onto two dimensions (umap() function from R package umap v0.2)⁴⁰ and clustered (dbscan() function from dbscan R package). This approach identified 101 distinct clusters of submotifs with 1–22 alleles participating in each (Fig. 2c,d and Supplementary Fig. 2b,c).

Evaluation and validation of length-dependent motif differences. To compare 8-mer motifs with 9-mer motifs, we generated pseudo 8-mer motifs from 9-mers by dropping middle residues (positions 4, 5 or 6) (related to Fig. 3). To compare 10- and 11-mer motifs with 9-mers, we generated pseudo 9-mer motifs by dropping middle residues from 10-mers (positions 5 or 6) and 11-mers (positions 5 and 6, or 6 and 7). The pseudo motif that was most similar to the true motif was used to evaluate the change in frequency and entropy at every peptide position. We considered 8-, 10- or 11-mer motifs with at least 100 identified peptides, that had an absolute difference in residue frequency with the true motif of >0.25 or an absolute difference in entropy of >0.2 at any position, as different. For example, if proline is observed at position 5 in 5% of peptides associated with motif1 but 40% of peptides associated with motif2, the absolute difference in frequency is $|0.05 - 0.40| = 0.35 > 0.25$, thus deeming motif1 and motif2 different.

To experimentally validate the observed length-specific motifs, we selected 18 peptides from three alleles representing four length-specific motifs that were predicted to be strong binders by our algorithm (MSi) but weak binders according to NetMHCpan4.0 and tested them for binding in *in vitro* binding assays. Peptide affinity measurements were performed at Immunitrack, Copenhagen, Denmark, as previously described⁴¹.

PPV versus PPV at 40% recall. To evaluate predictive power, we constructed datasets consisting of the observed allele- and length-specific binders in our MS data (n) along with $999 \times n$ random decoys from the human proteome and considered the fraction of correctly predicted binders in the top 0.1% of the dataset (that is, $PPV = \frac{\text{true positive calls}}{\text{all positive calls}} = \frac{\text{true positive calls}}{n}$). We advocated for the PPV evaluation

metric of Abelin et al.⁴, over the commonly used area under the receiver operating characteristic curve (AUC), because it is better suited for the HLA presentation prediction problem space where a relatively small number of true binders need to be identified amongst an excess of nonbinders. Each HLA allele is expected to present a repertoire of approximately ~10,000 peptides^{17,42–45} among the 1.1×10^7 9-mer peptides in the proteome, meaning that approximately only 1 of 1,000 peptides (0.1%) gets presented.

The definition of PPV described above is equivalent to PPV at recall equal to PPV% since the number of positive calls equals the number of true positives. A different version of this metric used recently to quantify algorithm performance is PPV at recall equals 40%; that is,

$$PPV_{40\% \text{ Recall}} = \frac{\text{true positive calls}}{\text{all positive calls}} = \frac{\text{true positive calls}}{\text{positive calls when 40\% of then true positives have been called}}^7$$

In addition, Bulik-Sullivan et al.⁷ used a dataset with a ratio of 1 hit to 10,000 decoys (h:d ratio) for the evaluation of a single-allele dataset, rather than a ratio of 1:999 used by us, which reduced PPV. Due to these differences, model performance evaluated with PPV is not comparable to model performance evaluated with $PPV_{40\% \text{ recall}, h:d=1:10,000}$ (Fig. 5d).

Development of integrative HLA binding prediction models. *Overview of model training procedure.* Machine learning models were created to predict the likelihood that a given peptide will be endogenously presented by a given HLA class I molecule. The positive training set consisted of MS-identified peptides from the set of 95 B721.221 mono-allelic cell lines (hits). The negative set consisted of random peptides drawn from the human proteome that did not overlap with the hits (decoys). Note that one decoy set was used for training and a separate nonoverlapping decoy set was used for evaluation. The number of positive and negative training examples was balanced by sampling $10 \times \# \text{hits}$ decoys, and sampling the hits ten times with replacement (this was done after splitting the data into folds to ensure that each hit was only present in one unique fold). Training was carried out in a standard fivefold cross-validation procedure: training data were split into five equal parts; each part was left out one at a time and a model was trained on the remaining four parts, obtaining five models, each of which was evaluated on its corresponding left out set. The fivefold cross-validation training was repeated three times with different model initialization random seeds. The final predictions score for each peptide was the average of the three initializations. Neural network models were trained using the Theano and Keras Python libraries.

Allele- and length-specific models. To build models that are both allele- and length-specific, only the MS hits identified for that particular allele and length were used to train the model. This was done for each of the 95 alleles and for lengths 8, 9, 10 and 11 where at least 40 peptides were identified. The MSi neural network models were fully connected with one hidden layer of size 50 and tanh activation. Training was carried out in batches of size 30, for ten epochs, and early stopping determined by evaluating on a 20% hold-out partition of the training set to avoid overfitting. Three different models were trained with three different encodings of the peptide sequence: (1) one-hot (also known as binary or dummy) encoding; (2) similarity encoding using the blosum62 matrix; and (3) similarity encoding based on the PMBEC matrix⁴⁶. In addition to the peptide sequence-encoding features, the MSi models included the following features:

- (1) Amino acid properties—each peptide residue was represented by the first three principal components derived from a dimensionality reduction of a large set of physicochemical properties³⁹.
- (2) Peptide-level characteristic—eight peptide-level features were computed with the following R package peptides: 'boman', 'hmoment', 'hydrophobicity', 'helixbend', 'sidechain', 'xstr', 'partspec', 'pkc'.

The logit-transformed output scores from MSi models were used as input features to train logistic regression models that integrate endogenous signals (Supplementary Notes 6 and 7): MSiC models were trained with two features (MSi scores and cleavability), MSiCE were trained with three input features (MSi scores, cleavability and expression) and MSiCEB were trained with four input features (MSi scores, cleavability, expression and presentation bias). Note that despite expression having a larger predictive contribution over cleavability, the cleavability feature is incorporated first to the integrative models (that is, MSiC, MSiCE instead of MSiE, MSiEC) to allow for samples that lack expression data to use the cleavability since upstream and downstream peptide context sequences are readily available from the source protein.

Pan-allele-pan-length models. Pan models were trained similarly with some differences. A single panMSi model was trained with all 8–11-mer peptides identified across the 95 alleles. The panMSi neural networks had additional input features to describe the binding pocket of the HLA protein—each binding pocket residue was represented with ten Kidera factors and three principal components (see Allele similarity analyses). The size of the hidden layer was 250, batch size was 5,000 and the hidden layer activation function was rectified linear unit. Training proceeded for 15 epochs with early stopping determined by lack of improvement in four consecutive epochs.

To construct integrative pan models, we considered the four HLA genes (HLA-A, HLA-B, HLA-C and HLA-G) and the four lengths (8, 9, 10 and 11) separately:

linear panMSiC, panMSiCE and panMSiCEB models were trained for all HLA-A alleles and peptides of length 8, all HLA-A alleles and peptides of length 9, all HLA-A alleles and peptides of length 10, and all HLA-A alleles and peptides of length 10, and analogously for HLA-B, -C and -G.

Modeling PPV variability. A linear regression model was trained to predict the achieved PPV (MSi, 9-mer models) given the following variables (related to Fig. 5):

- (1) The total number of 9-mer hits observed for the allele
- (2) The sum of entropies at positions 1 through 9 (main motif entropy)
- (3) The number of identified submotifs for the allele
- (4) The sum of submotif entropies
- (5) Estimated abundance of the binding motif calculated by weighting the frequency of each residue at each peptide position by the natural abundance frequency of each of the 20 amino acids

Evaluation of model performance in multi-allelic samples. To evaluate model performance in multi-allelic patient samples, each tumor-presented peptide was scored for binding to each of the sample-specific HLA alleles (related to Fig. 6). To compare scores for different alleles, each score was converted to percentile rank. To this end, empirical cumulative distribution functions (R package `stats`, function `ecdf()`) were computed for each model (including benchmark algorithms) from the scores of a background set of 1×10^6 random decoys. Each decoy set was constructed such that it contains proportions of 8-, 9-, 10- and 11-mers that are equal to the observed length distribution for the allele (or the HLA gene in the case of pan models). A peptide was considered to be correctly identified as a binder if the predicted binding score for at least one of the alleles in the sample was better than 99.9% of the scores in the corresponding decoy set (0.1 percentile rank; evolution was also performed at different rank thresholds; Supplementary Fig. 5). This approach is very similar to the %Ranks introduced by NetMHC³. A peptide was assigned to the allele(s) for which it had %rank score < 0.1, where assignment to more than one allele was allowed.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

The original mass spectra for 79 of 95 mono-allelic datasets generated for this study, the protein sequence database and tables of peptide spectrum matches for all 95 alleles have been deposited in the public proteomics repository MassIVE (<https://massive.ucsd.edu>) and are accessible at <ftp://massive.ucsd.edu/MSV000084172/>. MS data for the 16 previously published mono-allelic datasets in MassIVE can be downloaded at <ftp://massive.ucsd.edu/MSV000080527>. Datasets for the patient samples are accessible at <ftp://massive.ucsd.edu/MSV000084442/>. B721.221 RNA-seq data for HLA-C (C*04:01, C*07:01) are deposited under GEO: GSE131267. Melanoma RNA-seq data are deposited in dbGaP (https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001451.v1.p1 (ref. ¹⁵)). Glioblastoma bulk RNA-seq data are available through dbGaP (<https://www.ncbi.nlm.nih.gov/gap>) with accession number [phs001519.v1.p1](https://www.ncbi.nlm.nih.gov/gap) (ref. ²⁶). All other data are available from the corresponding authors upon reasonable request.

Code availability

Code used to generate plots characterizing allele-specific preferences (for example, logo plots, entropy plots, peptide projection and clustering, overlap with IEDB data and so on) as well as code to build a sample neural network prediction model is provided as Supplementary Code. The HLATHENA predictors are available to use online for research purposes only at <http://HLATHENA.tools>. For commercial usage inquiries please contact the authors or the Broad Institute.

References

38. Kidera, A., Konishi, Y., Oka, M., Ooi, T. & Scheraga, H. A. Statistical analysis of the physical properties of the 20 naturally occurring amino acids. *J. Protein Chem.* **4**, 23–55 (1985).
39. Bremel, R. D. & Homan, E. J. An integrated approach to epitope analysis I: dimensional reduction, visualization and prediction of MHC binding using amino acid principal components and regression approaches. *Immunome Res.* **6**, 7 (2010).
40. McInnes, L., Healy, J. & Melville, J. UMAP: uniform manifold approximation and projection for dimension reduction. Preprint at arXiv:1802.03426 [stat. ML] (2018).
41. Harndahl, M. et al. Peptide binding to HLA class I molecules: homogenous, high-throughput screening, and affinity assays. *J. Biomol. Screen.* **14**, 173–180 (2009).
42. Bassani-Sternberg, M., Pletscher-Frankild, S., Jensen, L. J. & Mann, M. Mass spectrometry of human leukocyte antigen class I peptidomes reveals strong effects of protein abundance and turnover on antigen presentation. *Mol. Cell. Proteom.* **14**, 658–673 (2015).
43. Hunt, D. F. et al. Characterization of peptides bound to the class I MHC molecule HLA-A2.1 by mass spectrometry. *Science* **255**, 1261–1263 (1992).
44. Rammensee, H. G., Friede, T. & Stevanović, S. MHC ligands and peptide motifs: first listing. *Immunogenetics* **41**, 178–228 (1995).
45. Rammensee, H., Bachmann, J., Emmerich, N. P., Bachor, O. A. & Stevanović, S. SYFPEITHI: database for MHC ligands and peptide motifs. *Immunogenetics* **50**, 213–219 (1999).
46. Kim, Y., Sidney, J., Pinilla, C., Sette, A. & Peters, B. Derivation of an amino acid similarity matrix for peptide:MHC binding and its application as a Bayesian prior. *BMC Bioinformatics* **10**, 394 (2009).

Acknowledgements

We acknowledge technical assistance from K. Pelton, S. Santagata, O. Spiro, L. Elagina, B. Knisbacher, S. Shukla, J. Brugge and A. Appfel. We further express gratitude for constructive input from M. Rooney, J. Abelin and Z. Hu. We acknowledge support from the National Institutes of Health: grant nos. NCI-1R01CA155010-02 (to C.J.W.), NHLBI-5R01HL103532-03 (to C.J.W.), NIH/NCI U24 CA224331 (to C.J.W.), NIH/NCI R21 CA216772-01A1 (to D.B.K.), NCI-SPORE-2P50CA101942-11A1 (to D.B.K.), NHGRI T32HG002295 and NIH/NCI T32CA207021 (to S.S.), NCI 5T32CA009172-41 (to D.A.B.), NIH/NCI U24-CA210986 and NIH/NCI U01 CA214125 (to S.A.C.). This work was supported in part by The G. Harold and Leila Y. Mathers Foundation and the Bridge Project, a partnership between the Koch Institute for Integrative Cancer Research at MIT and the Dana-Farber/Harvard Cancer Center. D.A.B. is supported in part by the John R. Svenson Fellowship. C.J.W. is a scholar of the Leukemia and Lymphoma Society, and is supported in part by the Parker Institute for Cancer Immunotherapy. S.K. is a Cancer Research Institute/Hearst Foundation fellow.

Author contributions

D.B.K., C.J.W., N.H. and S.A.C. directed the overall study design. S.S. performed computational analyses and developed predictive models. S.K., C.R.H., H.K. and K.R.C. generated the MS data and performed data analysis. D.B.K. and G.L.Z. selected the HLA alleles for analysis. D.B.K., P.M.L. and L.W.L. generated the single-HLA allele cell lines and performed data generation. D.B.K., G.O., K.L.L., D.A.B., P.M.L. and L.W.L. developed the patient-derived tumor cell lines. I.K.Z. and J.M.R. generated and provided cells from an ovarian cancer PDX model. P.B. provided CLL samples for analysis. W.Z. provided expert technical assistance. T.E. generated RNA-seq data for mono-allelic cell lines. T.O. and T.L. generated and quantified ribosome profiling data. J.S. and W.J.L. performed HLA typing and validation of all cell lines. S.J. performed HLA binding validation assays. S.S., S.K., N.H., C.J.W. and D.B.K. wrote the manuscript, with contributions from all co-authors.

Competing interests

D.B.K. has previously advised Neon Therapeutics, and owns equity in Aduro Biotech, Agenus Inc., Armata Pharmaceuticals, Biomarin Pharmaceutical Inc., Bristol-Myers Squibb Com., Celldex Therapeutics Inc., Editas Medicine Inc., Exelixis Inc., Gilead Sciences Inc., IMV Inc., Lexicon Pharmaceuticals Inc. and Stemline Therapeutics Inc. D.A.B. has received consulting fees from Octane Global, Defined Health, Dedham Group, Adept Field Solutions, Slingshot Insights, Blueprint Partnership, Charles River Associates, Trinity Group and Insight Strategy, and is a member of the RCC translational medicine advisory board of Bristol-Myers Squibb. K.L.L. owns equity and is a founder of Travera LLC and is an advisor to Bristol-Myers Squibb Com. and Rarecyte. S.A.C. is a member of the scientific advisory boards of Kymera, PTM BioLabs and BioAnalytix and a scientific advisor to Pfizer and Biogen. C.J.W. and N.H. are founders of Neon Therapeutics and members of its scientific advisory board. N.H. is also an advisor for IFM therapeutics. W.J.L. is a member of the scientific advisory board of CareDx. All other authors have no competing interests. Patent applications have been filed on aspects of the described work entitled as follows: 'HLA single allele lines', and 'Methods for identifying neoantigens'.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41587-019-0322-9>.

Correspondence and requests for materials should be addressed to N.H., S.A.C., C.J.W. or D.B.K.

Reprints and permissions information is available at www.nature.com/reprints.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- | n/a | Confirmed |
|-------------------------------------|--|
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The statistical test(s) used AND whether they are one- or two-sided
<i>Only common tests should be described solely by name; describe more complex techniques in the Methods section.</i> |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of all covariates tested |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
<i>Give P values as exact values whenever suitable.</i> |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated |

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection	Mass spectrometry data was collected using Xcalibur software, Thermo Fisher Scientific. RNA-seq data was collected with standard RNA extraction and sequencing procedures as described in the methods.
Data analysis	Mass spectrometry data was analyzed using Spectrum Mill v6.1 pre-Release (Agilent Technologies, Santa Clara, CA) as described in the methods section. RNA data was processed using Bowtie 2 and RSEM as described in the methods section. Analysis of peptide characteristics was performed with common R packages all names in the methods. Models were trained with Theano/Keras python libraries as described in the methods.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The original mass spectra, the protein sequence database, and tables of peptide spectrum matches for the 79 newly described monoallelic datasets have been deposited in the public proteomics repository MassIVE and are accessible at <ftp://MSV000084172@massive.ucsd.edu>. Mass spectrometry data for the 16 previously published mono-allelic datasets can be downloaded from MassIVE under the identifier MassIVE: MSV000080527. Datasets for the patient samples are accessible at <ftp://MSV000084442@massive.ucsd.edu>. B721.221 RNA seq data for HLA-C (C*04:01, C*07:01) is deposited under GEO: GSE131267. All other data are available from the corresponding author upon reasonable request.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	95 HLA alleles; 11 new-profiled patient samples
Data exclusions	Not applicable
Replication	Technical replicate injections
Randomization	Not applicable
Blinding	Not applicable

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involved in the study
<input type="checkbox"/>	<input checked="" type="checkbox"/> Antibodies
<input type="checkbox"/>	<input checked="" type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input type="checkbox"/>	<input checked="" type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data

Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

Antibodies

Antibodies used	W6/32 (sc-32235, Santa Cruz)
Validation	Enrichment of pan-HLA Class I molecules

Eukaryotic cell lines

Policy information about [cell lines](#)

Cell line source(s)	721.221 B cell line was purchased from The international HLA reference Standards (IHWG) biorepository cell bank, Fred Hutchinson Cancer Research Center. Tumor cell lines were derived from MEL and GBM patient resected tumor specimen.
Authentication	721.221 B cell line were authenticated at IHWB
Mycoplasma contamination	All cell lines tested negative for mycoplasma
Commonly misidentified lines (See ICLAC register)	not applicable

Human research participants

Policy information about [studies involving human research participants](#)

Population characteristics	Not applicable
----------------------------	----------------

Recruitment

Describe how participants were recruited. Outline any potential self-selection bias or other biases that may be present and how these are likely to impact results.

Ethics oversight

All human tissues were obtained through DFCI or Partners Healthcare approved IRB protocols.

Note that full information on the approval of the study protocol must also be provided in the manuscript.