

Next-generation epitope prediction using mass spectrometry and integrative genomics

POSTER LB-179

Michael S. Rooney^{1,8,9}, Jennifer G. Abelin¹, Derin B. Keskin^{1,3,4,6}, Siranush Sarkizova^{1,2}, Christina Hartigan¹, Wandi Zhang³, John Sidney⁷, Jonathan Stevens⁵, William Lane⁵, Guang Lan Zhang^{3,6,10}, Karl R. Clauser¹, Nir Hacohen^{1,3,11}, Steven A. Carr¹, Catherine J. Wu^{1,3,4,6}

¹Broad Institute of MIT and Harvard, Cambridge, MA, USA; ²Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA; ³Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA; ⁴Department of Medicine, Brigham and Women's Hospital, Boston, MA, USA; ⁵Tissue Typing Laboratory, Brigham and Women's Hospital, Boston, MA, USA; ⁶Harvard Medical School, Boston, MA, USA; ⁷La Jolla Institute for Allergy and Immunology, La Jolla, CA; ⁸Harvard/MIT Division of Health Sciences and Technology, Cambridge, Massachusetts, USA; ⁹Neon Therapeutics, Cambridge, MA, USA; ¹⁰Department of Computer Science, Metropolitan College, Boston University, Boston, MA, USA; ¹¹Center for Cancer Immunology, Massachusetts General Hospital



Abstract

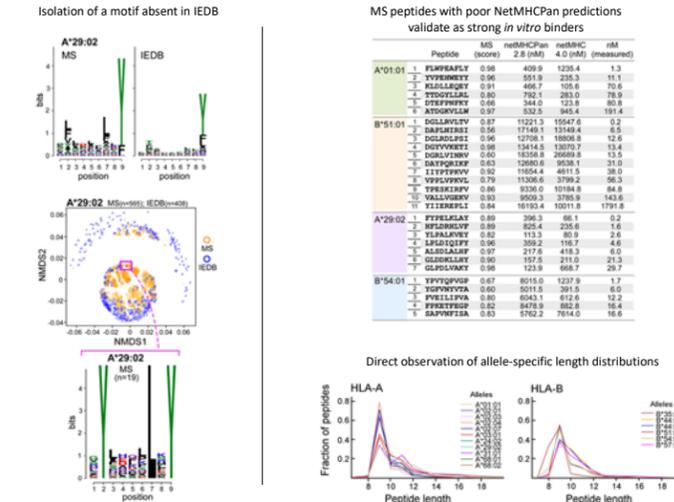
Neoantigens are somatically mutated protein sequences presented on MHC Class I or II molecules. Mounting evidence suggests neoantigens potentiate anti-tumor immune responses during checkpoint blockade, and there is great interest in directly targeting neoantigens therapeutically with vaccine. Such vaccines must rely on personalized *in silico* epitope prediction given the multitude of MHC alleles in the human population (each with its own peptide binding repertoire) and the very small degree of somatic mutation overlap among cancer patients. The current prediction paradigm is based almost entirely on machine learning algorithms (e.g. NetMHC) trained on *in vitro* p:MHC binding affinity assays.

Mass spectrometry (MS), which can be used to directly identify endogenously processed and presented peptides, is an orthogonal approach to define the rules of the MHC ligandome. To date, MS studies have identified large numbers of MHC-bound peptides and produced general observations, such as bias toward highly expressed proteins, but have not yielded new predictors. The main bottleneck to prediction has been the use of cell lines expressing the natural complement of **six distinct MHC Class I molecules** (2 HLA-A alleles, 2 HLA-B alleles, and 2 HLA-C alleles). In this context, peptides can only be assigned to one of the six alleles using pre-existing predictors; many peptides cannot be confidently assigned to any allele. This convolution drastically reduces the amount of forward learning that can be achieved.

Here we modify the MS approach to focus on an **MHC-null cell line** transfected with a **single allele** of our choice. We iteratively applied rapid, high-resolution liquid chromatography mass spectrometry (LC-MS/MS) to 16 HLA alleles identifying >24,000 MHC-bound peptides. Using these data, we were able to learn peptide-binding and proteasomal cleavage motifs *de novo*. In addition, we uncovered new motifs absent in the public database IEDB (which would have been impossible using the multi-allelic MS approach) and validated these with traditional *in vitro* binding assays. Furthermore, we systematically assess variables hypothesized to predict presentation. Notably, we find that transcript expression is highly predictive of presentation and that source protein localization makes minimal incremental contribution.

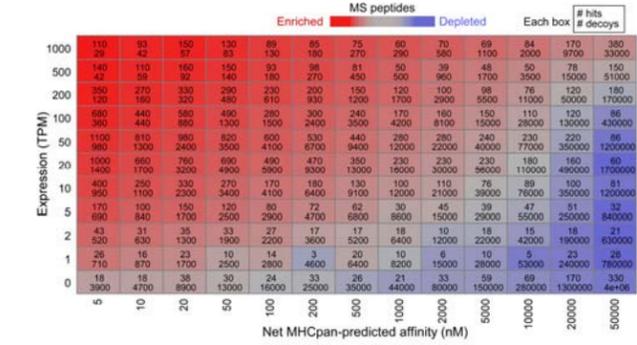
Finally, we trained neural network models on our endogenous peptide dataset to produce integrated predictors that leverage both peptide sequence and auxiliary variables such as gene expression and the likelihood of efficient proteasomal processing. We demonstrate effective prediction in MS data sets collected by independent groups, achieving 2-3x more correct identifications than NetMHC or NetMHCpan. We also show compelling performance in predicting immunogenic HIV epitopes. We conclude that measuring mono-allelic MHC peptides using LC-MS/MS improves the utility of predictive algorithms and provides a rapid and scalable method to generate rules for neoantigen prediction.

Single-allele system identifies novel motifs that are not discoverable using traditional multi-allelic approach

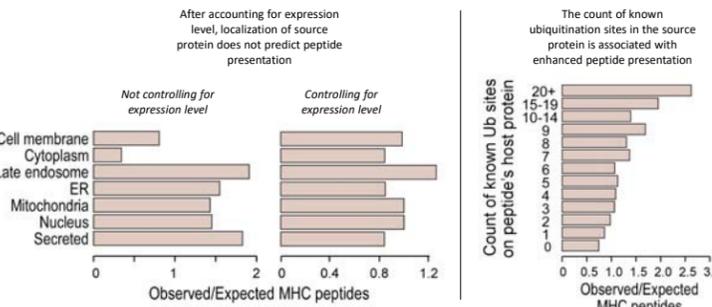


MHC peptides reveal multiplicative relationship between gene expression and binding affinity

- MS peptides and random genomic 9mer decoys were binned according expression (RNA-Seq) and predicted affinity
- Binder:decoy ratio per bin shows that a 10x increase in expression roughly compensates for a 90% reduction in binding strength

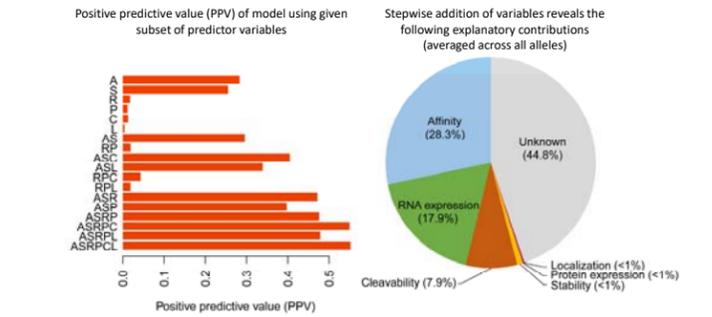


Most cellular localization bias can be explained by expression



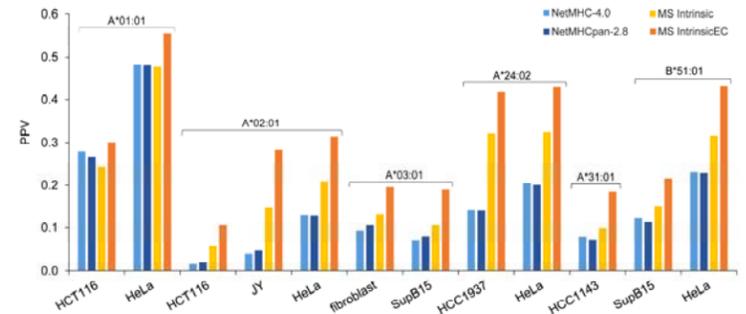
Single-allele data enables development of data-driven integrated predictors; expression and cleavability greatly increase accuracy

- Logistic regression models were developed to discriminate *n* MS-observed peptides from 999n genomic 9mer decoys
- Each model included a different subset of predictors, e.g. $\text{IS}(\text{binder}) \sim \beta_0 + \beta_1 \text{NetMHCscore} + \beta_2 \text{Expression}$
- Model performance was assessed according to the fraction of true MS peptides among the most highly ranked peptides (top 0.1%); this score is called the positive predictive value (PPV)



Novel predictors based on our single-allele MS data outperform NetMHC on external test data

- Bassani-Sternberg *et al* (2015) identified MHC-bound peptides for six human cell lines using mass spectrometry
- We evaluated our algorithms by seeing whether we could distinguish these peptides from decoys (999 decoys per true peptide)
- Our MS-based predictor bests NetMHC for almost every allele we could evaluate
- Adding expression and cleavage to the model ("+ EC") improves prediction further – **2-3x** the performance of NetMHC



Two tables showing HIV epitope counts for HLA-A and HLA-B alleles. The left table is for HLA-A and the right is for HLA-B. Each table has columns for 'Number of HIV epitopes' and 'Number of peptides'.

Conclusions and Next Steps

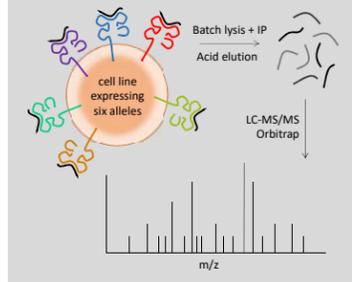
- Single-allele system generates high coverage and enables deep integrative analysis by removing the problem of allele ambiguity
- MS data can be used to learn peptide motifs and cleavage motifs *de novo* and to systematically probe how various factors (protein localization, expression, etc.) impact epitope presentation
- These findings will yield new predictors that will be critical to the success of personalized neoantigen approaches
- Method can be rapidly deployed to expand allele coverage for Class I and can be adapted for Class II MHC predictions

Acknowledgments

We thank Wilfredo F. Garcia-Beltran of the Ragon Institute of MGH, MIT and Harvard for providing the cell lines. The work was supported by the Blavatnik Family Foundation (N.H. and C.W.), the NCI R01 CA155010 (C.W.), NCI CPTAC U24 CA160034 (S.A.C. and K.C.), NHGRI T32 HG002295 (S.S.), and Dana-Farber/Harvard Cancer Center Kidney Cancer SPORC 2P50CA101942-11A1 (D.B.K.). C.W. is a Scholar of the Leukemia and Lymphoma Society.

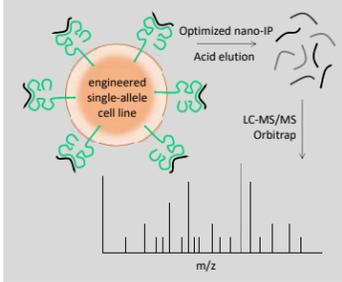
Traditional MS approach

- Cell lines express **six** different class I alleles
- Must rely on known motifs to assign peptides to alleles



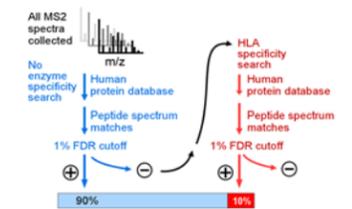
Our approach

- Class I-null B721.221 cells are transfected with **single allele**
- No ambiguity** in allele assignment; higher depth per allele



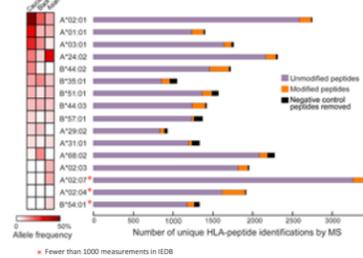
Iterative spectral search

Boosts identifications by 5-40%



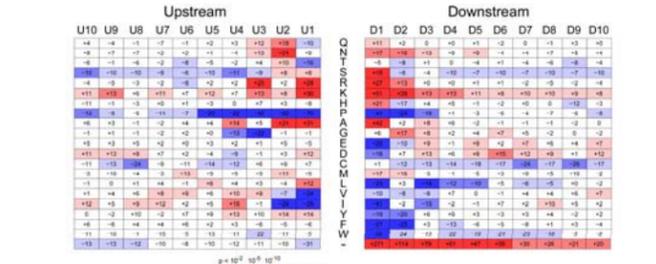
16 alleles, 24,000 peptides

Compare to total size of IEDB – 97,000 peptides

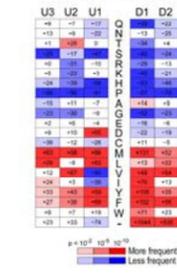


Quantitative analysis of *in vivo*-processed peptide shows tryptic signature, preference for alanine, and strong proline avoidance

- The upstream and downstream sequence of each MS binder was compared to decoy sequences
- Heatmap shows percent change in frequency of given amino acid at given position when comparing binders to decoys
- Signature is preserved in breast, colorectal, and fibroblast cell lines as well as HeLa and PBMC.



Class II MUT3



Imputed from NetChop

