

Machine learning models for antigen immunogenicity and T-cell recognition

Barbara Bravi

Department of Mathematics, Imperial College London (UK)

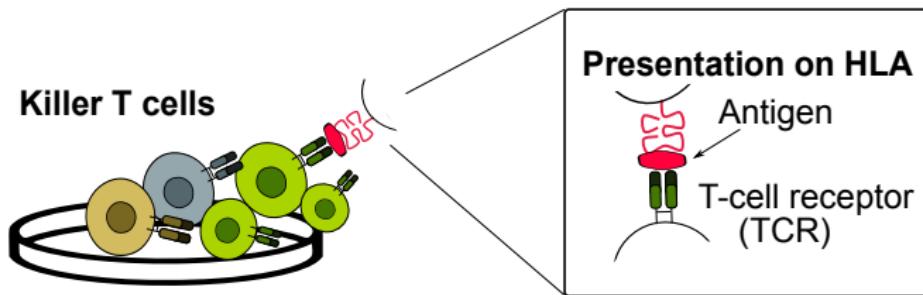


Joint work with: A. Di Gioacchino, J. Fernandez-De-Cossio-Diaz, S. Cocco, R. Monasson,
T. Mora, A.M. Walczak (ENS Paris)

BEvAS, EPFL 17/04/2023

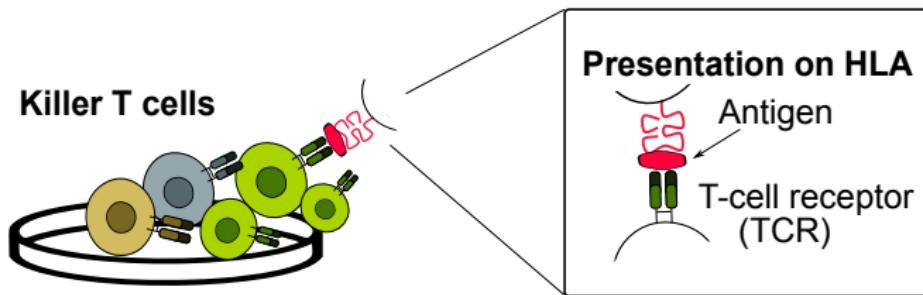
Objective:

Machine learning approach that is able to extract biologically interpretable features on antigen immunogenicity & T-cell epitope specificity



Objective:

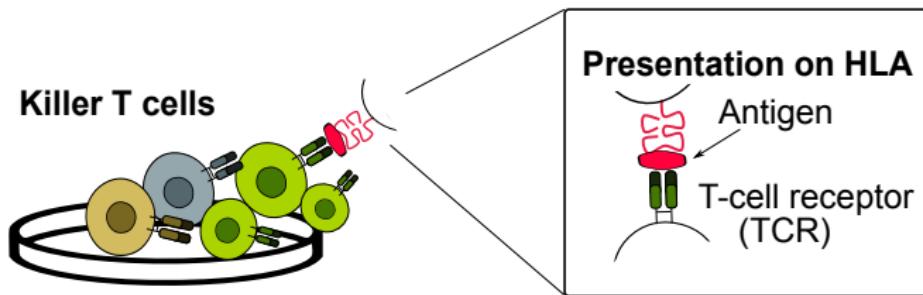
Machine learning approach that is able to extract biologically interpretable features on antigen immunogenicity & T-cell epitope specificity



Only a fraction of HLA-presented antigens are immunogenic (promote a T cell response).
Immunogenicity prediction: key in neoantigen discovery, low success rate (Wells et al. 2020)

Objective:

Machine learning approach that is able to extract biologically interpretable features on antigen immunogenicity & T-cell epitope specificity



Only a fraction of HLA-presented antigens are immunogenic (promote a T cell response).
Immunogenicity prediction: key in neoantigen discovery, low success rate (Wells et al. 2020)

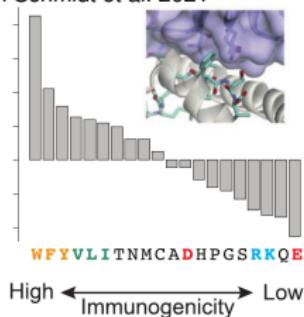
pMHC epitope elicits the response only of specific small subsets of TCRs, recent advances in prediction but insight into molecular properties is still challenging
(Gielis et al. 2019, Montemurro et al. 2021, Weber et al. 2021 and others)

Enrichment in distinctive patterns

Antigen immunogenicity and epitope-specificity of T cell receptors result from physico-chemical constraints on sequence composition

Enrichment in aromatic, hydrophobic residues in immunogenic peptides

From Schmidt et al. 2021



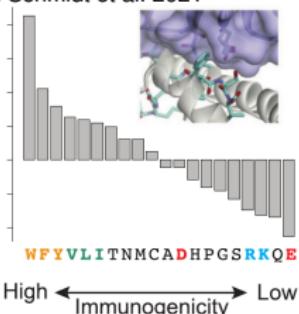
see also: Calis et al. 2013, Chowell et al. 2015

Enrichment in distinctive patterns

Antigen immunogenicity and epitope-specificity of T cell receptors result from physico-chemical constraints on sequence composition

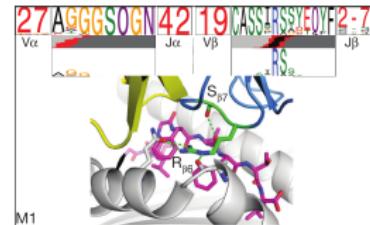
Enrichment in aromatic, hydrophobic residues in immunogenic peptides

From Schmidt et al. 2021



see also: Calis et al. 2013, Chowell et al. 2015

Convergence in sequence motifs in epitope-specific receptors



From Dash et al. 2017

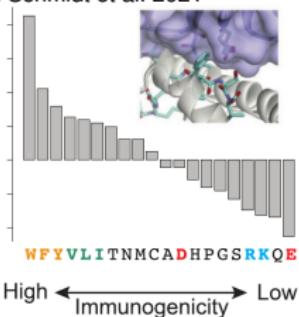
see also: Glanville et al. 2017, Meysman et al. 2019
Pogorelyy et al. 2019, Mayer-Blackwell et al. 2021 etc

Enrichment in distinctive patterns

Antigen immunogenicity and epitope-specificity of T cell receptors result from physico-chemical constraints on sequence composition

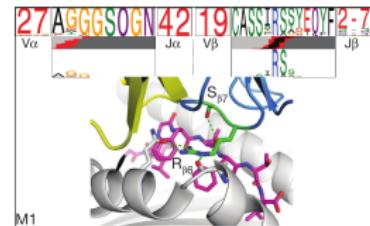
Enrichment in aromatic, hydrophobic residues in immunogenic peptides

From Schmidt et al. 2021



see also: Calis et al. 2013, Chowell et al. 2015

Convergence in sequence motifs in epitope-specific receptors



From Dash et al. 2017

see also: Glanville et al. 2017, Meysman et al. 2019
Pogorelyy et al. 2019, Mayer-Blackwell et al. 2021 etc

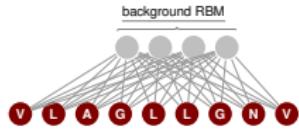
How to disentangle pattern enrichment from baseline constraints?
(ensuring e.g. in antigens high binding affinity to HLA)

Our approach:

Machine learning approach known as ‘transfer learning’
within the model known as Restricted Boltzmann Machines¹

Based on the pre-print: B. Bravi, A. Di Gioacchino, J. Fernandez-de-Cossio-Diaz, A.M. Walczak, T. Mora, S. Cocco, R. Monasson, *Learning the differences: a transfer-learning approach to predict antigen immunogenicity and T-cell receptor specificity*, Biorxiv 2022.12.06.519259v1 (2022)

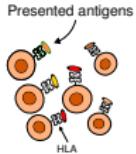
¹ (Smolensky 1986, Hinton 2002; Biophysical modelling: Tubiana, Cocco and Monasson 2019, Shimagaki and Weigt 2019)



HLA-A*02:01 presented antigens

VLAGLLGNV
GILGFVFTL
FLCLFLLPSL
SLQVELAHM
ALYGVWPLL
ALAESIRPL

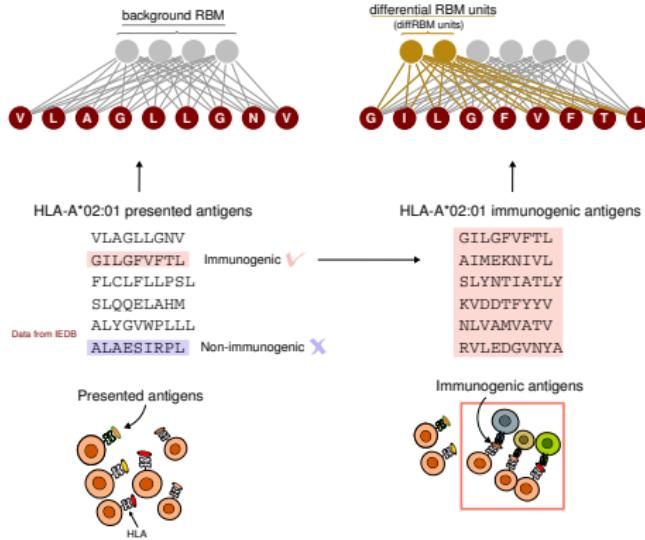
Data from IEDB



Model of presentation

(RBM-MHC, Bravi et al. 2021)

captures background constraints
(binding affinity to HLA)

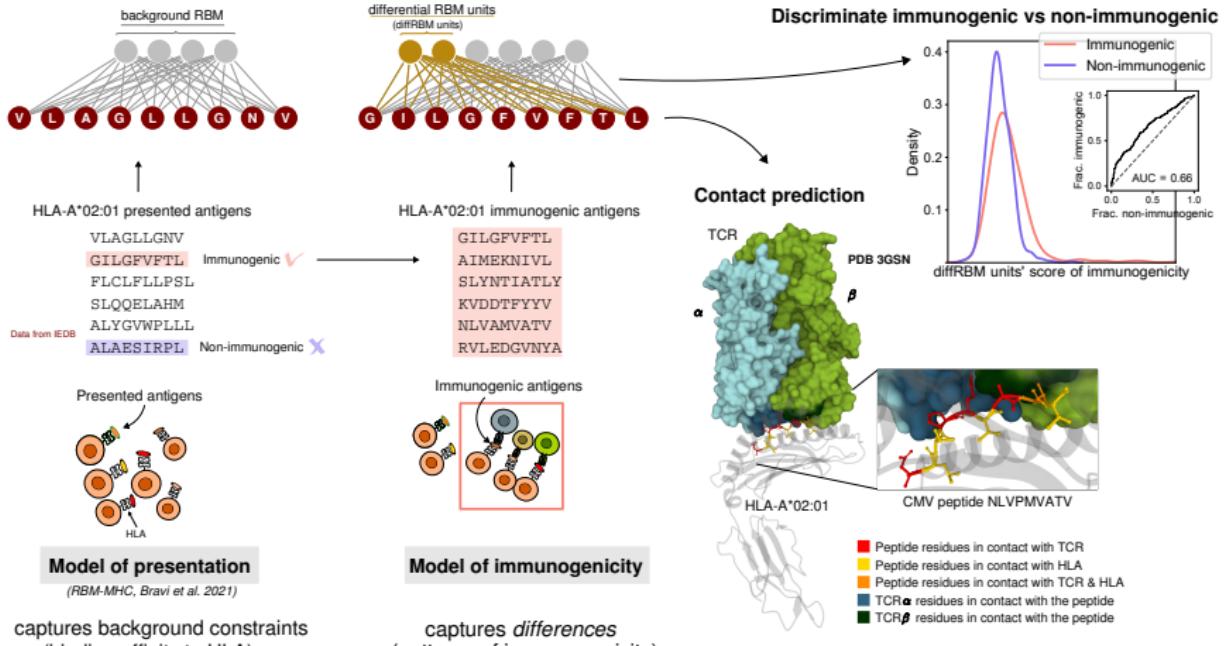


Model of presentation (RBM-MHC, Bravi et al. 2021)

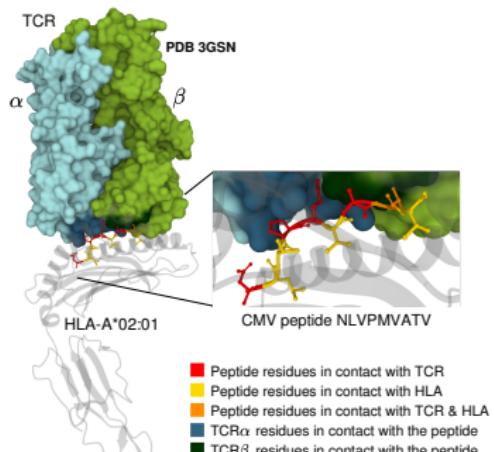
captures background constraints
(binding affinity to HLA)

Model of immunogenicity

captures *differences*
(patterns of immunogenicity)

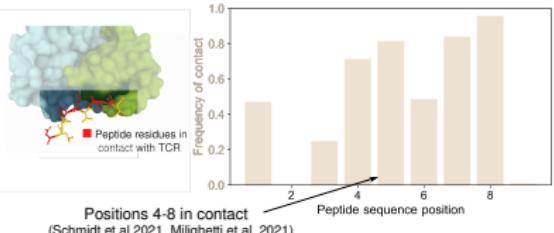


Contact positions in resolved structures:

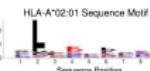
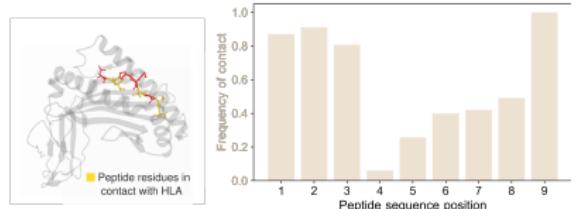


HLA-A*02:01-specific peptides

Peptide-TCR contacts

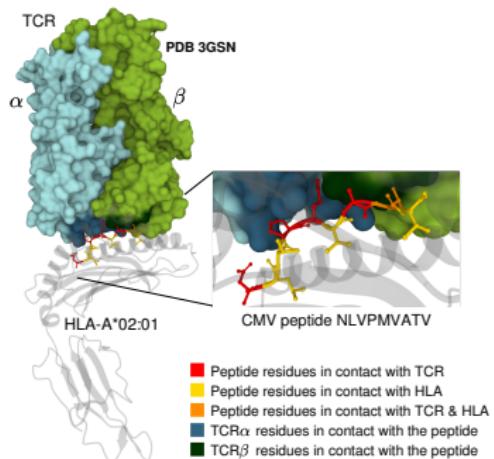


Peptide-HLA contacts



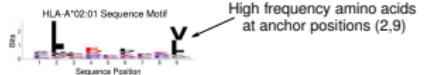
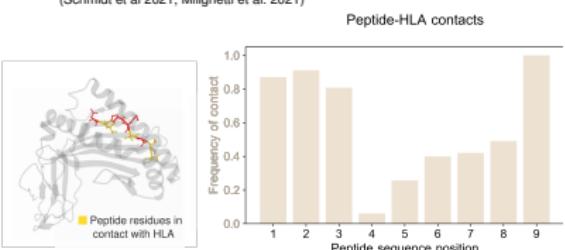
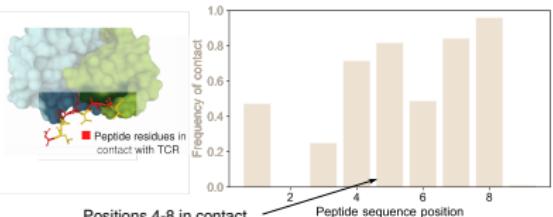
High frequency amino acids at anchor positions (2,9)

Contact positions in resolved structures:



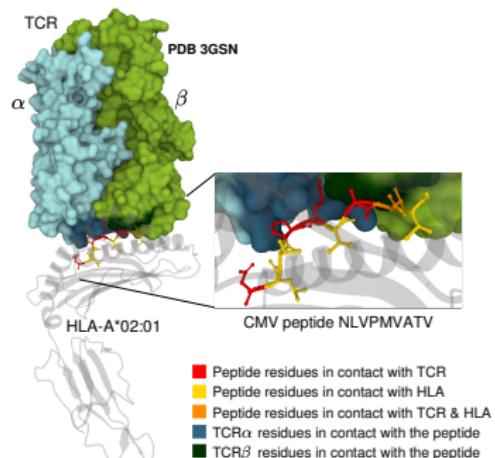
HLA-A*02:01-specific peptides

Peptide-TCR contacts



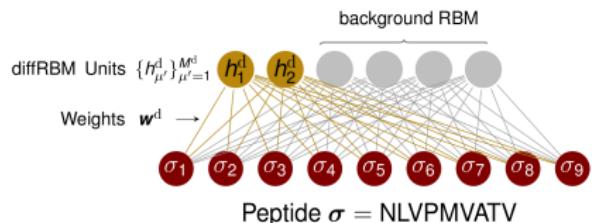
Differences in statistics of immunogenic peptides should reflect contacts

Contact positions in resolved structures:



Model prediction:

diffRBM architecture

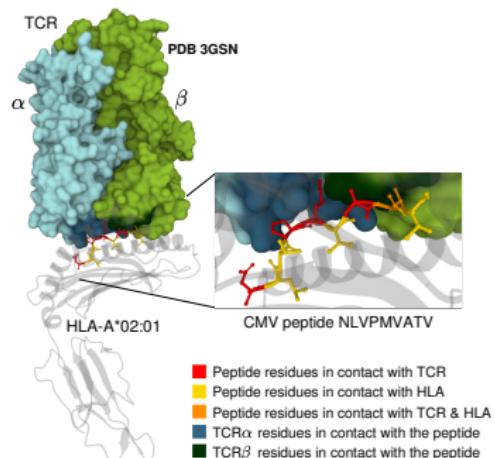


Single-site importance factors

$$T_i(\sigma_i) = \underbrace{g_i^d(\sigma_i)}_{\text{related to amino acid frequency difference between immunogenic and presented}} + \sum_{\mu'=1}^{M^d} \underbrace{w_{i\mu'}^d(\sigma_i) \langle h_{\mu'} | \sigma \rangle}_{\text{captures correlations between positions}}$$

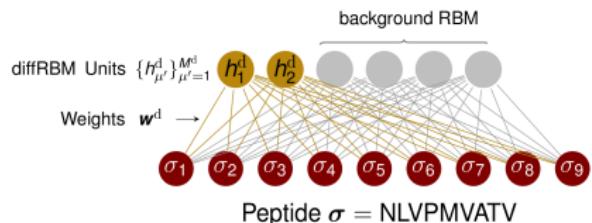
$\langle h_{\mu'} | \sigma \rangle$: from $P(h_{\mu'} | l_{\mu'}(\sigma))$, where $l_{\mu'}(\sigma) = \sum_i w_{i\mu'}^d(\sigma_i)$

Contact positions in resolved structures:



Model prediction:

diffRBM architecture



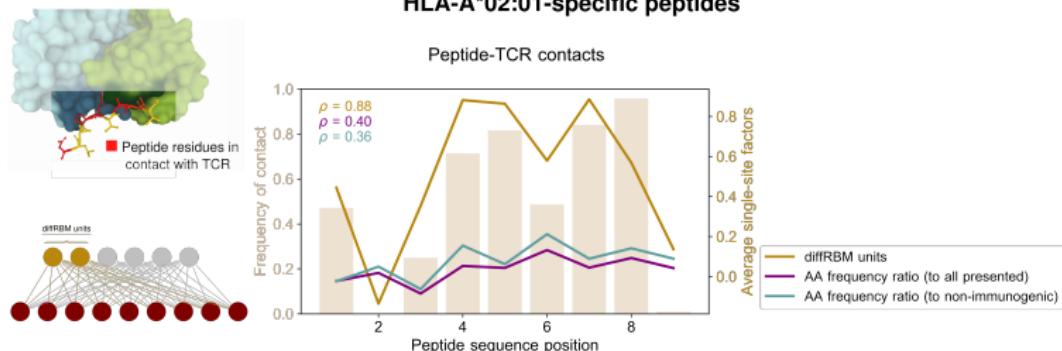
Single-site importance factors

$$T_i(\sigma_i) = \underbrace{g_i^d(\sigma_i)}_{\text{related to amino acid frequency difference between immunogenic and presented}} + \sum_{\mu'=1}^{M^d} \underbrace{w_{i\mu'}^d(\sigma_i) \langle h_{\mu'} | \sigma \rangle}_{\text{captures correlations between positions}}$$

$\langle h_{\mu'} | \sigma \rangle$: from $P(h_{\mu'} | l_{\mu'}(\sigma))$, where $l_{\mu'}(\sigma) = \sum_i w_{i\mu'}^d(\sigma_i)$

We hypothesize that sites at high $T_i(\sigma_i)$ are potential contacts

Structural interpretation

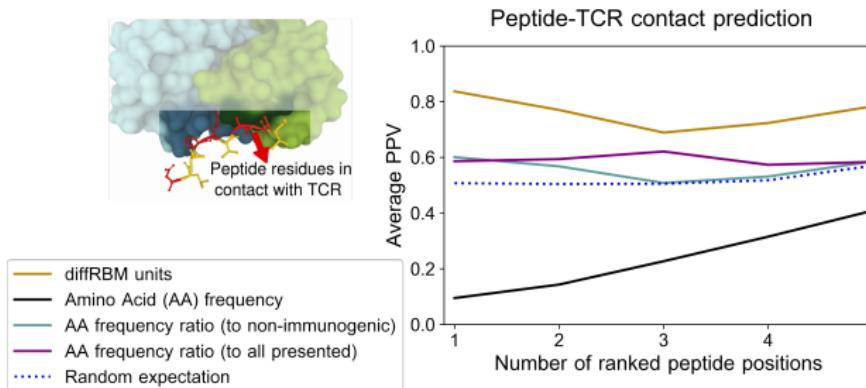


diffRBM identifies positions 4-8 as the most relevant to immunogenicity without restricting a priori the input sequences to a subset of positions

Comparison: independent-site models
based purely amino acid (AA) frequency (see e.g. IEDB tool)

Contact prediction (peptide-TCR)

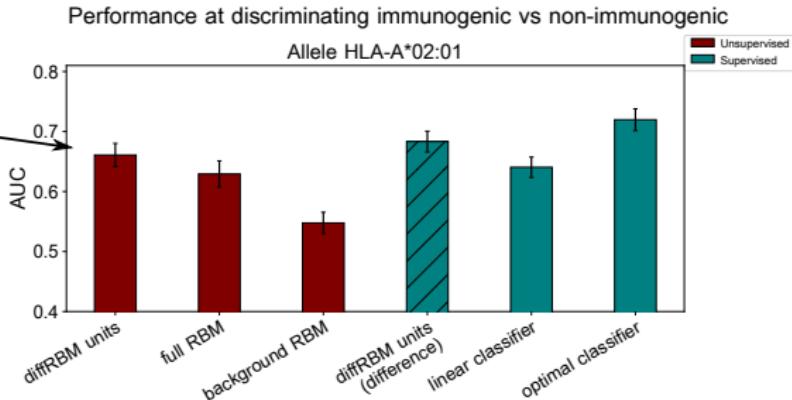
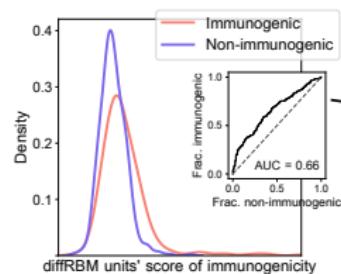
Top ranking positions by $T_i(\sigma_i) \rightarrow$ putative contacts



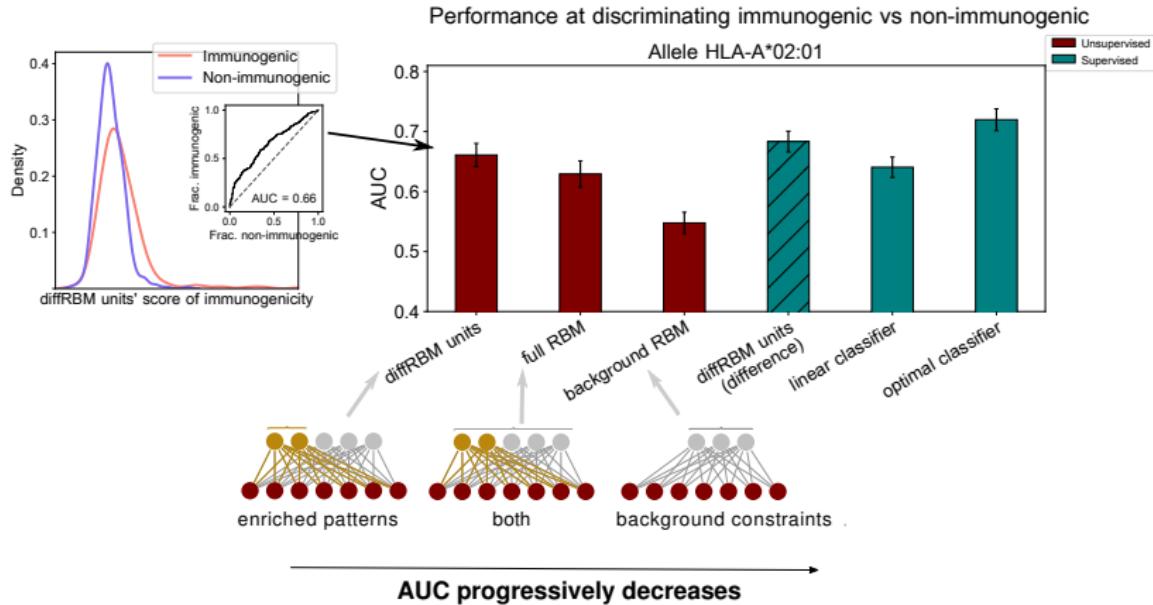
We test the prediction by the Positive Predictive Value (PPV):
fraction of ranked positions corresponding to true contacts

PPV averaged over structures for 3 HLA-I (HLA-A*02:01, HLA-B*35:01, HLA-B*07:02 - in total 46 structures)

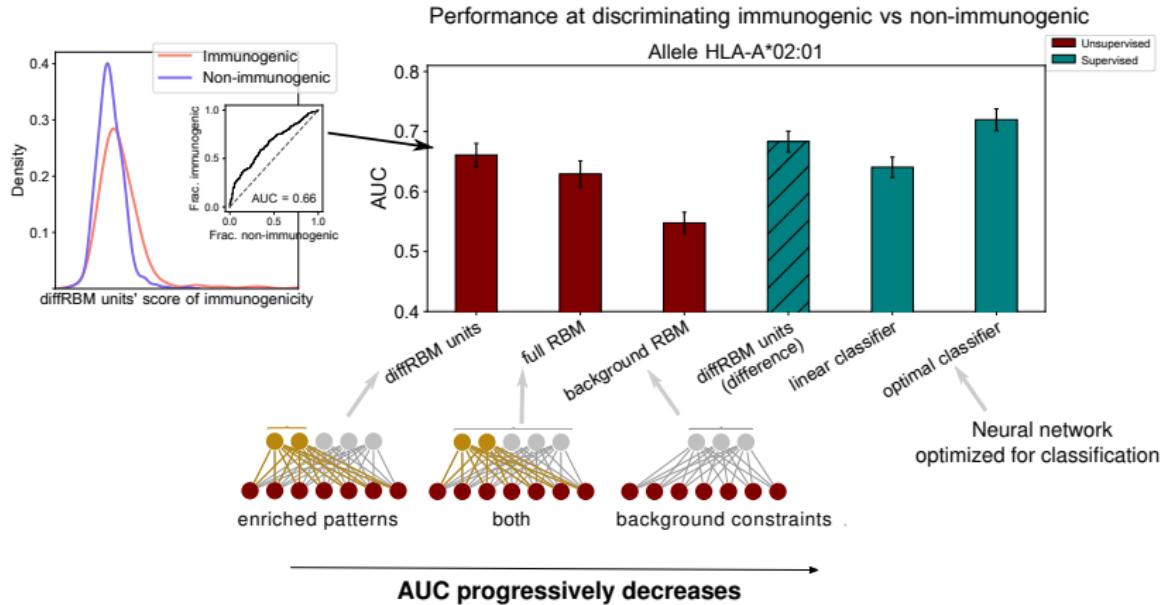
Classifying immunogenic vs non-immunogenic



Classifying immunogenic vs non-immunogenic



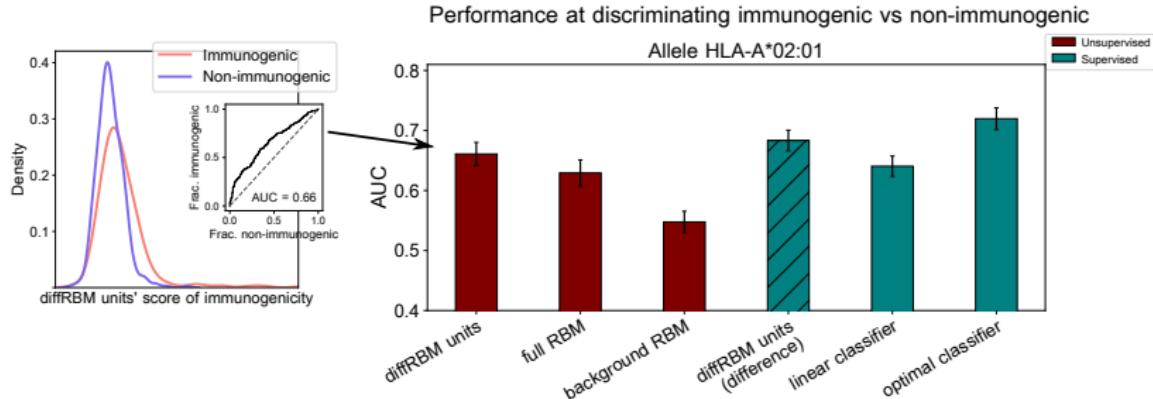
Classifying immunogenic vs non-immunogenic



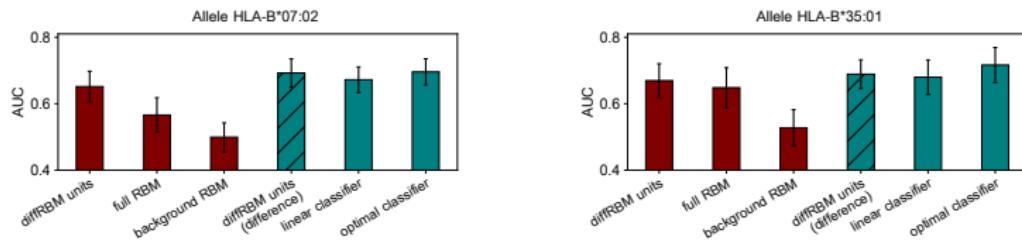
For efficient classification: score of diffRBM (immunogenic) - diffRBM (non-immunogenic)



Classifying immunogenic vs non-immunogenic



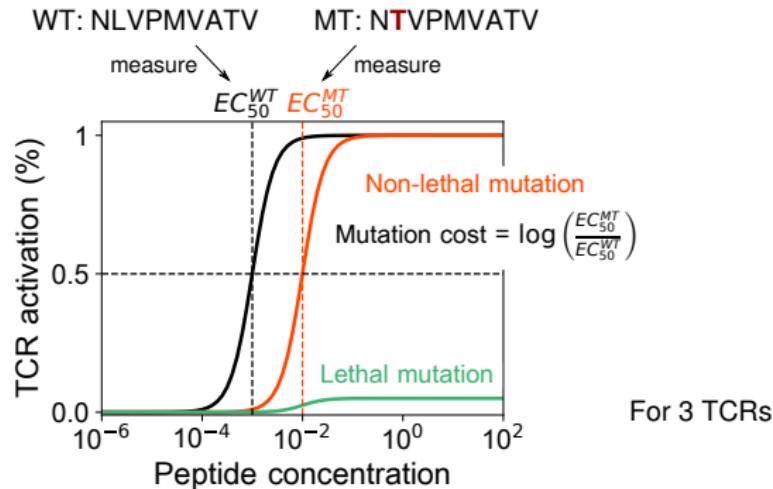
Consistent trend across HLA-specific models



Comparison to TCR avidity assays

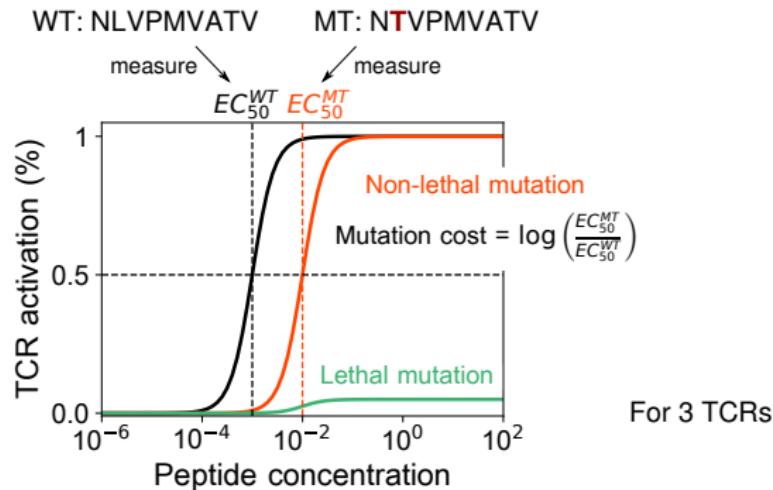
Data on TCR cross-reactivity to NLVPMVATV mutants

From Łuksza et al. Nature 2022



Comparison to TCR avidity assays

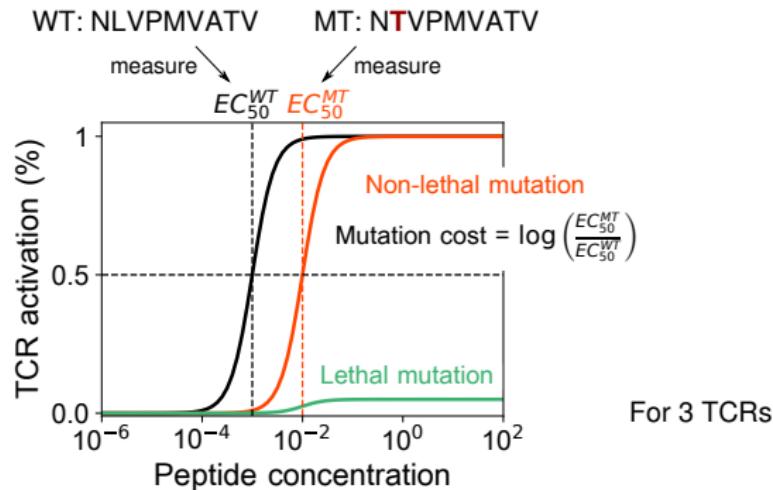
Data on TCR cross-reactivity to NLVPMVATV mutants
From Łuksza et al. Nature 2022



Positive mutation costs → loss in TCR response & decrease in antigen immunogenicity

Comparison to TCR avidity assays

Data on TCR cross-reactivity to NLVPMVATV mutants
From Łuksza et al. Nature 2022



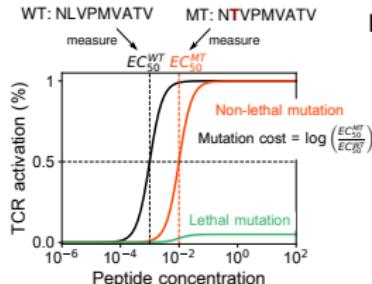
Positive mutation costs → loss in TCR response & decrease in antigen immunogenicity

We distinguish 2 groups of mutations: lethal and non-lethal (for TCR response)

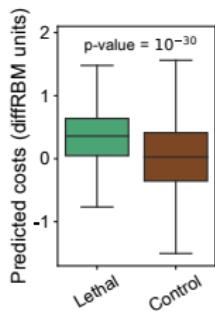
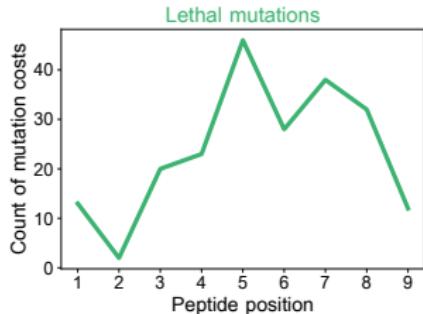
Comparison to TCR avidity assays

Data on TCR cross-reactivity to NLVPMVATV mutants

From Łuksza et al. Nature 2022



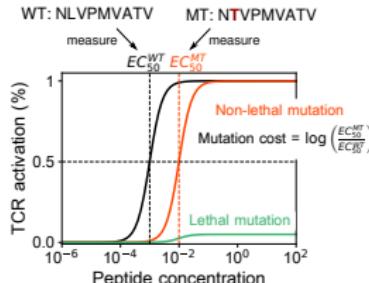
DiffRBM predicts the loss in immunogenicity upon lethal mutations



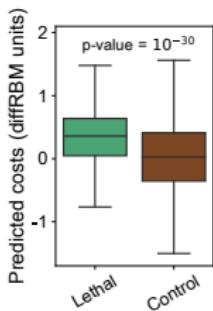
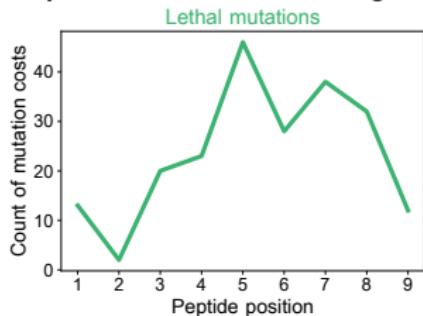
Comparison to TCR avidity assays

Data on TCR cross-reactivity to NLVPMVATV mutants

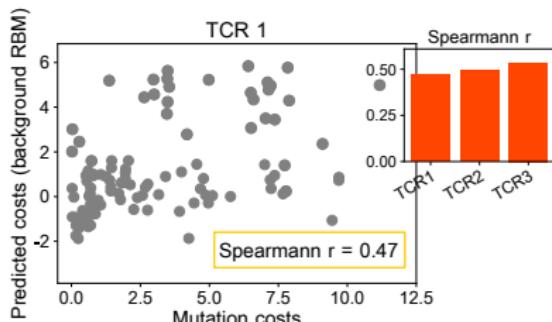
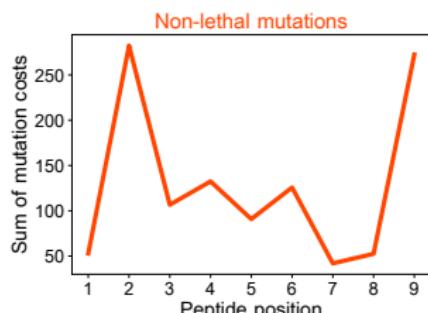
From Łuksza et al. Nature 2022



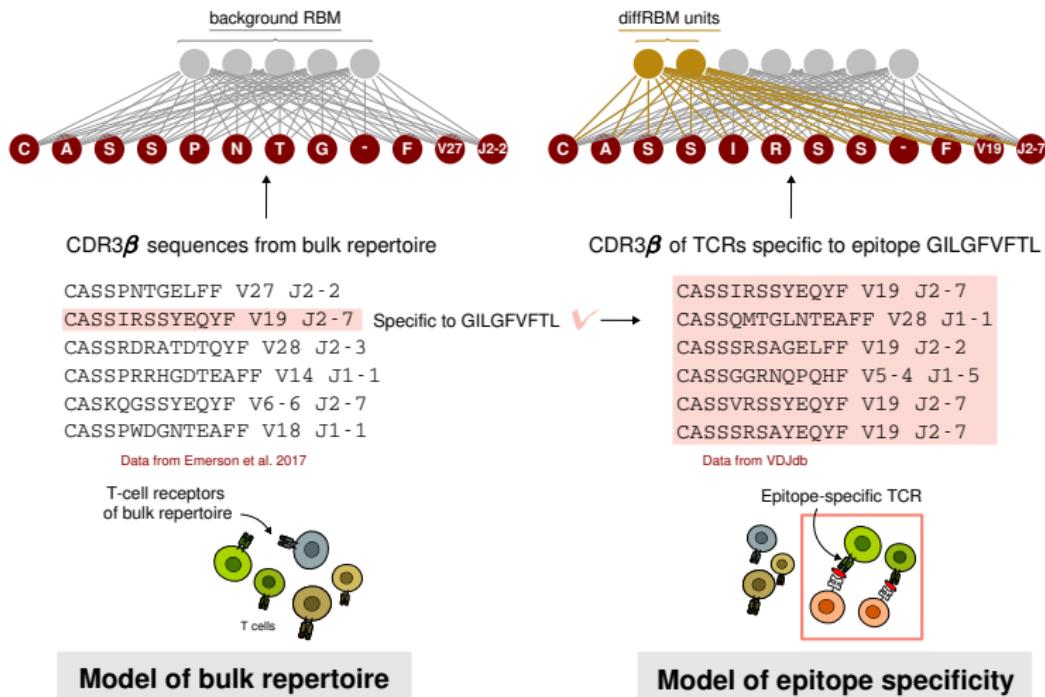
DiffRBM predicts the loss in immunogenicity upon lethal mutations



Non-lethal mutation costs are predicted by the presentation model

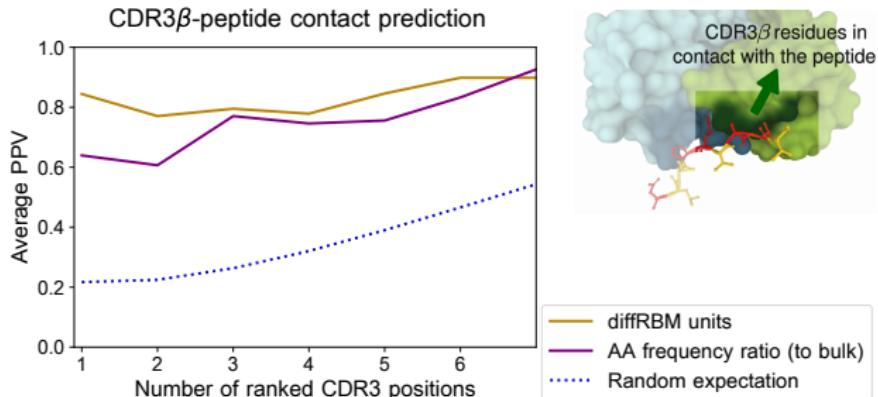


Model of T-cell epitope specificity



diffRBM units: capture antigen-driven convergent features

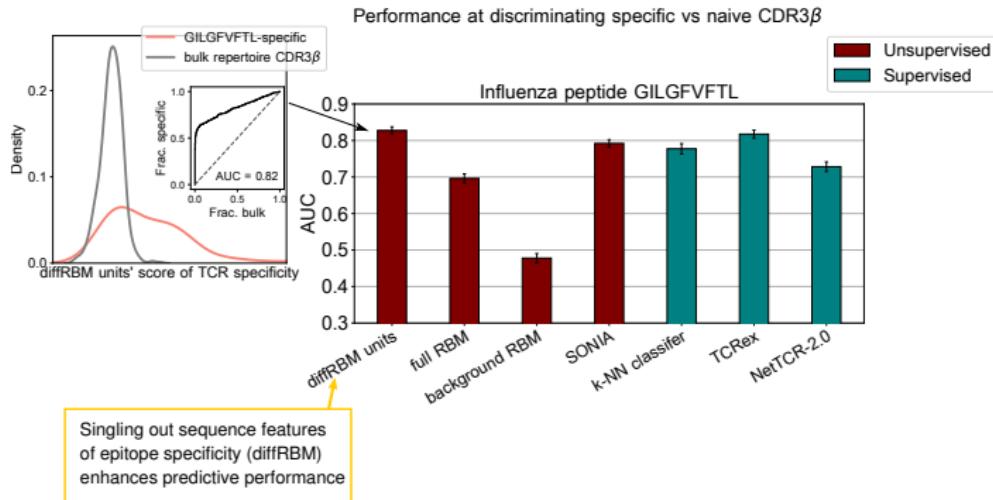
Contact prediction ($\text{CDR3}\beta$ -peptide)



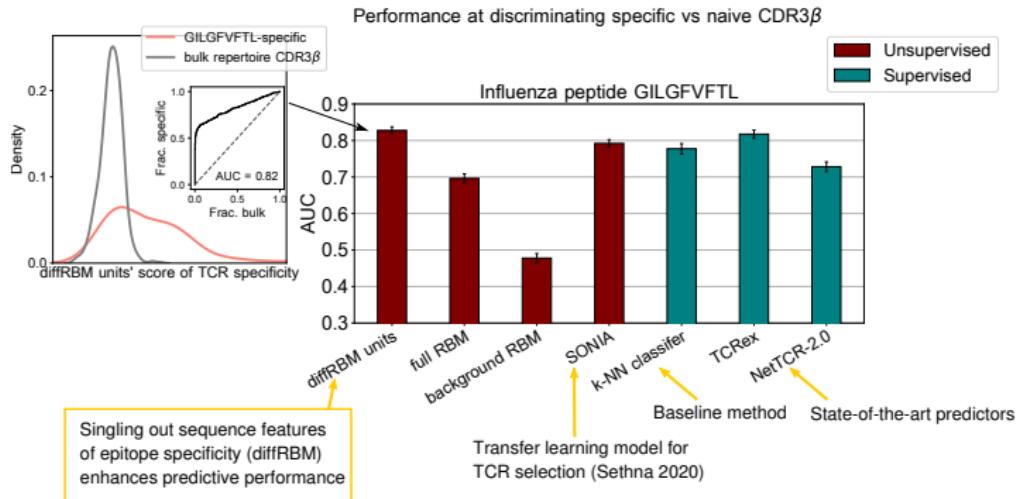
PPV averaged over structures for 4 epitopes (CMV, Influenza, EBV, Sars-CoV-2)

diffRBM performs better than independent models, both higher than random baseline

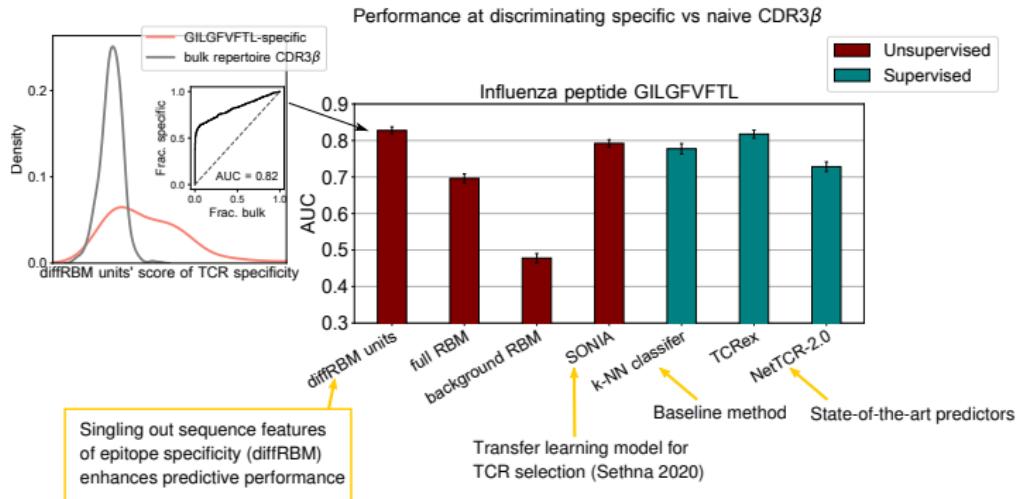
Classifying epitope-specific receptors



Classifying epitope-specific receptors

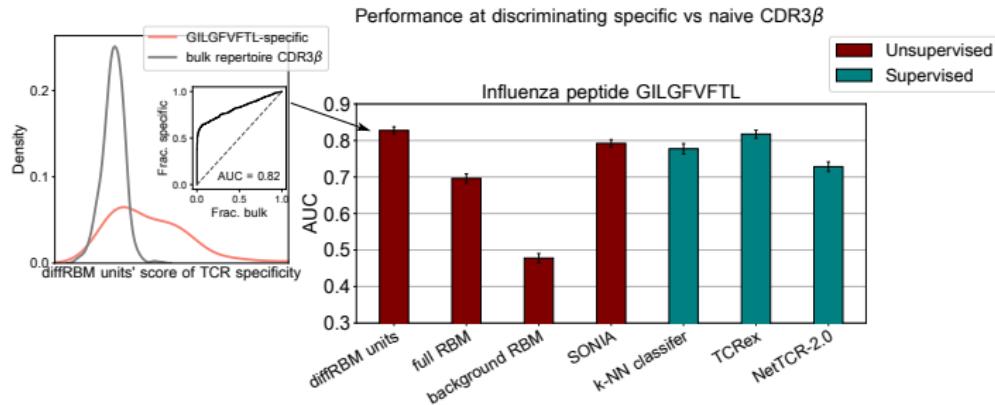


Classifying epitope-specific receptors

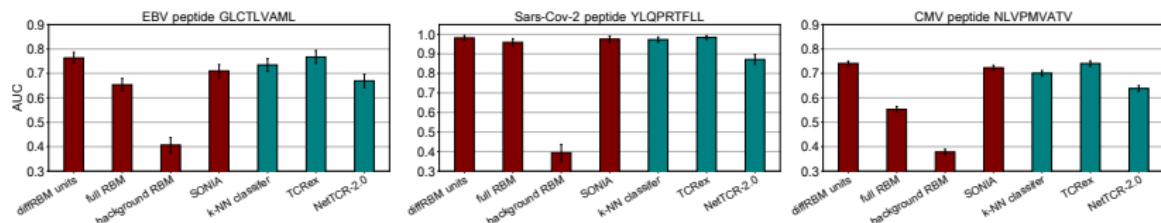


Classifying epitope-specific receptors vs generic non-binders (bulk):
diffRBM reaches the performance of state-of-the-art TCR specificity tools

Classifying epitope-specific receptors



Consistent trend across epitopes



Summary

- Transfer learning: diffRBM parameters capture characteristic differences of immunogenic peptides and epitope-specific receptors

Summary

- Transfer learning: diffRBM parameters capture characteristic differences of immunogenic peptides and epitope-specific receptors
- They allow us to estimate peptide-CDR3 β contacts

Summary

- Transfer learning: diffRBM parameters capture characteristic differences of immunogenic peptides and epitope-specific receptors
- They allow us to estimate peptide-CDR3 β contacts
- Probabilistic scores distinguish immunogenic vs non-immunogenic peptides, epitope-specific vs generic receptors, with performance comparable to classifiers

Summary

- Transfer learning: diffRBM parameters capture characteristic differences of immunogenic peptides and epitope-specific receptors
- They allow us to estimate peptide-CDR3 β contacts
- Probabilistic scores distinguish immunogenic vs non-immunogenic peptides, epitope-specific vs generic receptors, with performance comparable to classifiers
- Applications in vaccine design, TCR engineering, cancer neoantigen discovery, study of viral evolution and immunoediting in cancer

Summary

- Transfer learning: diffRBM parameters capture characteristic differences of immunogenic peptides and epitope-specific receptors
- They allow us to estimate peptide-CDR3 β contacts
- Probabilistic scores distinguish immunogenic vs non-immunogenic peptides, epitope-specific vs generic receptors, with performance comparable to classifiers
- Applications in vaccine design, TCR engineering, cancer neoantigen discovery, study of viral evolution and immunoediting in cancer
- Broader domain of application: distinctive sequence features that are selected upon (directed evolution, etc.)

ACKNOWLEDGEMENTS

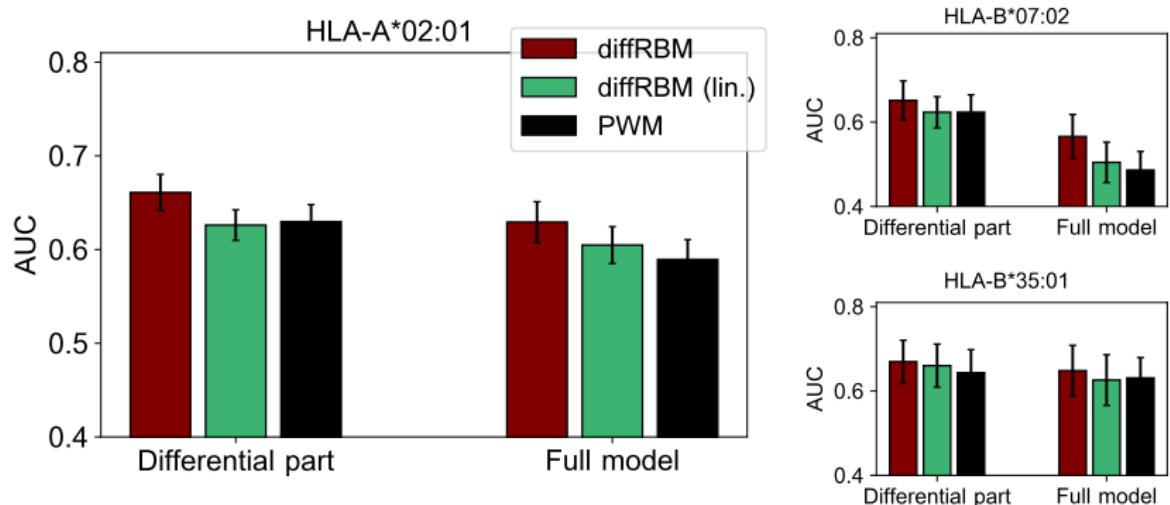
Statistical Physics and Inference for Biology Laboratory, ENS Paris

A. Di Gioacchino, J. Fernandez-De-Cossio-Diaz, S. Cocco, R. Monasson,
T. Mora, A.M. Walczak

Thank you for your attention!

Comparison of differential models (Immunogenicity)

Performance at discriminating immunogenic vs non-immunogenic



Comparison to PRIME and IEDB tool

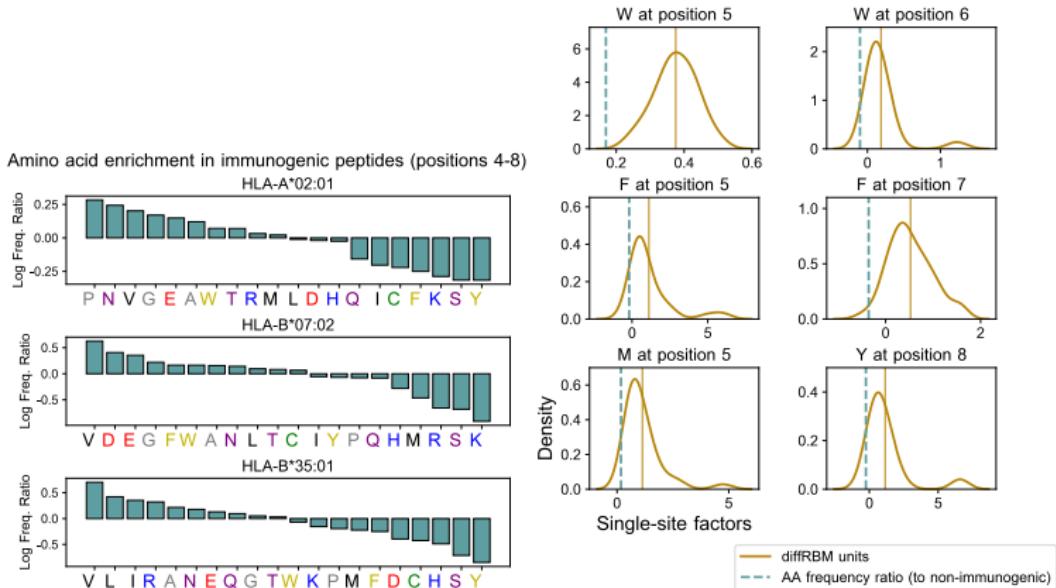
top RBM: AUC = 0.66 (HLA-A*02:01), AUC = 0.65 (HLA-B*07:02), AUC = 0.67 (HLA-B*35:01)

PRIME (Schmidt et al. 2021): AUC = 0.56 (HLA-A*02:01), AUC = 0.52 (HLA-B*07:02), AUC = 0.58 (HLA-B*35:01)

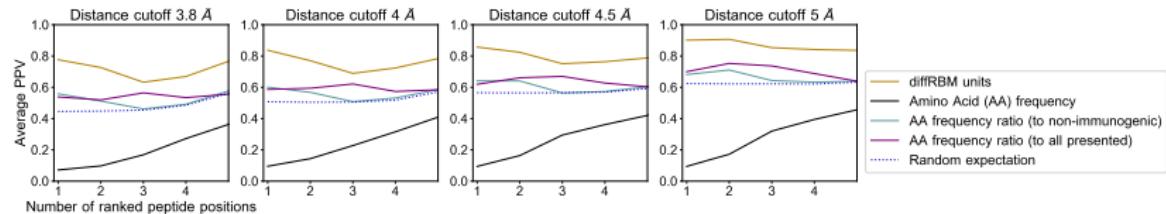
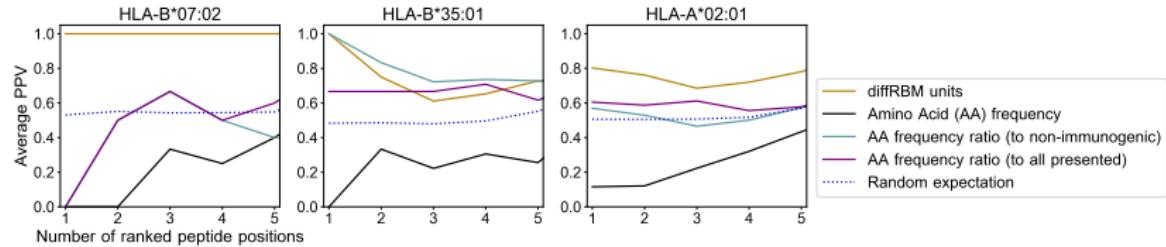
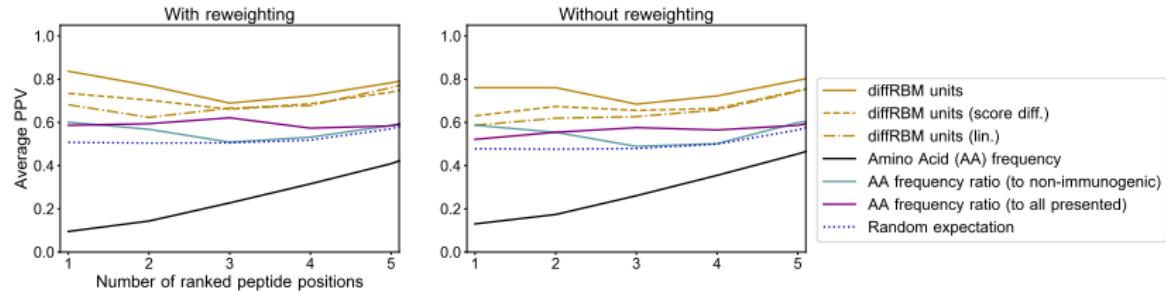
IEDB tool (Calis et al. 2013): AUC = 0.53 (HLA-A*02:01), AUC = 0.60 (HLA-B*07:02), AUC = 0.57 (HLA-B*35:01)

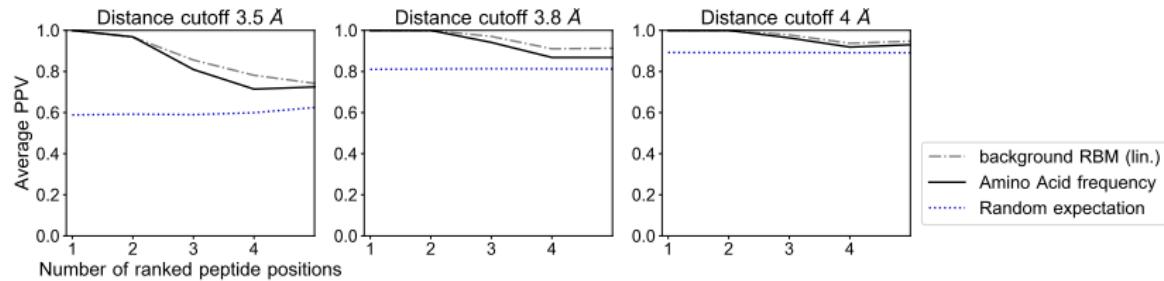
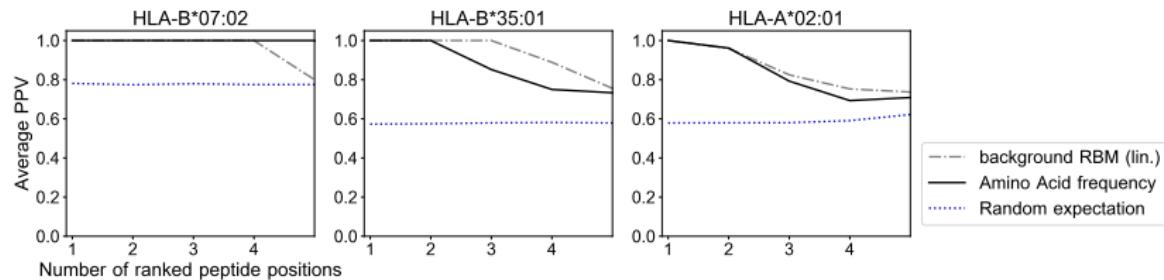
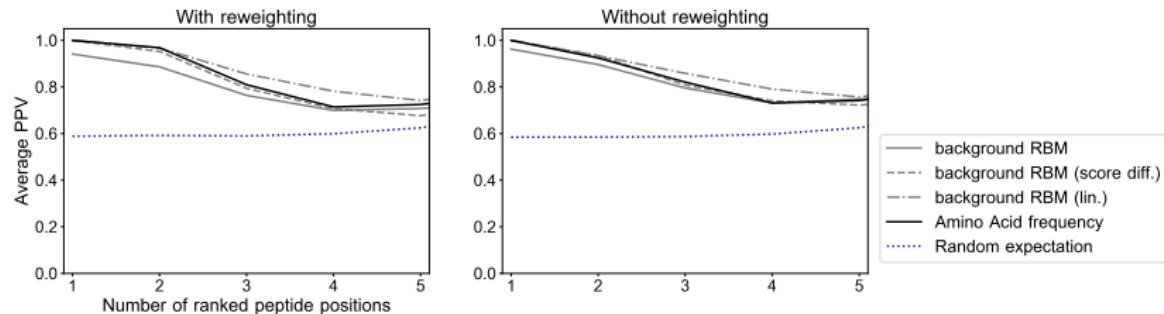
(Note: different training set)

Residues' contribution to immunogenicity

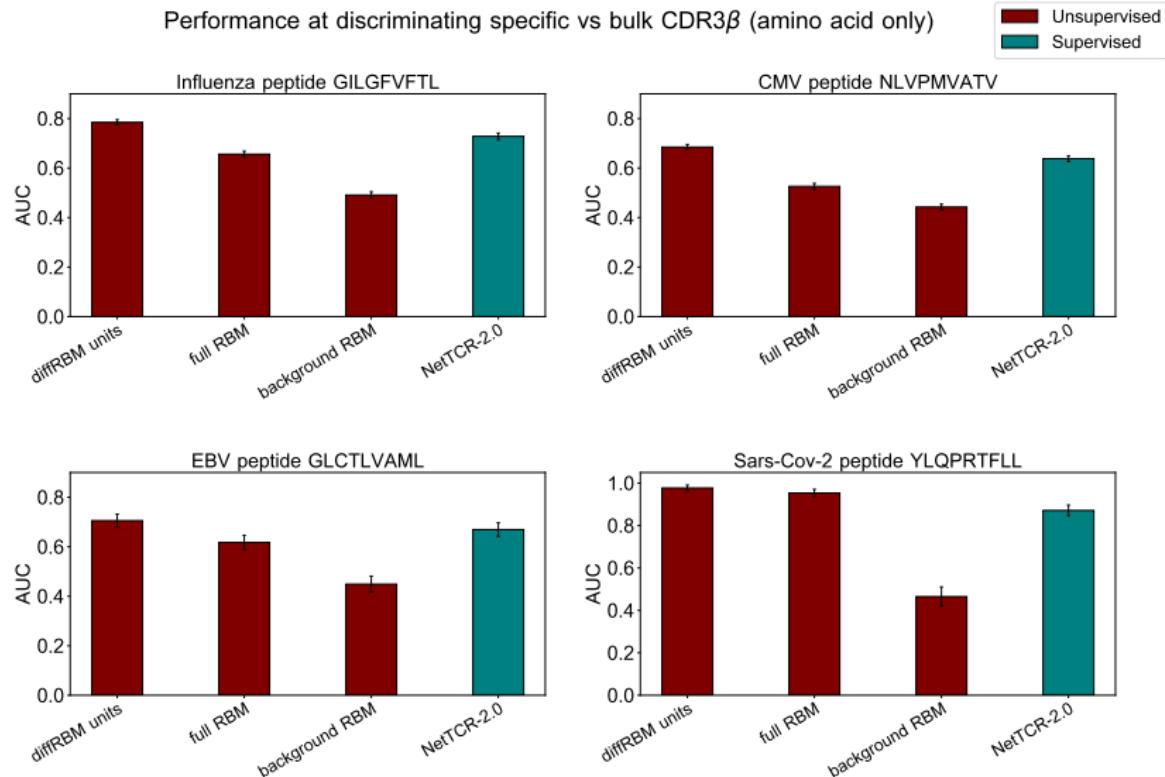


W at position 5 and 6, F at position 5 and 7 (Schmidt et al. 2021), M at position 5 (Luksza et al. 2022), Y at position 8 (Piepenbrink et al. 2013)



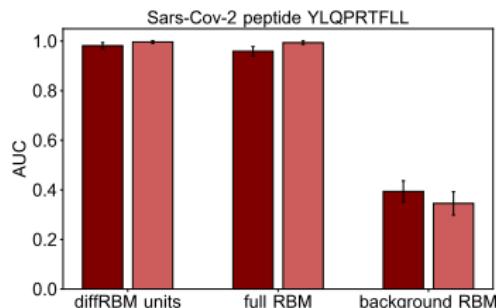
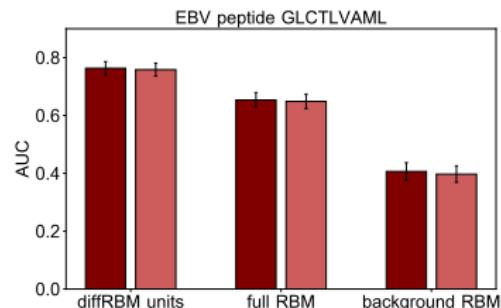
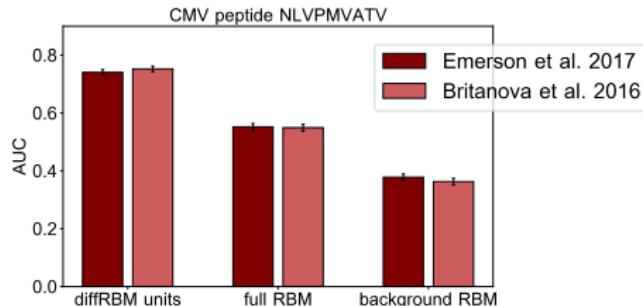
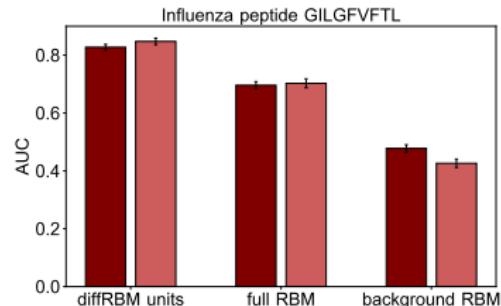


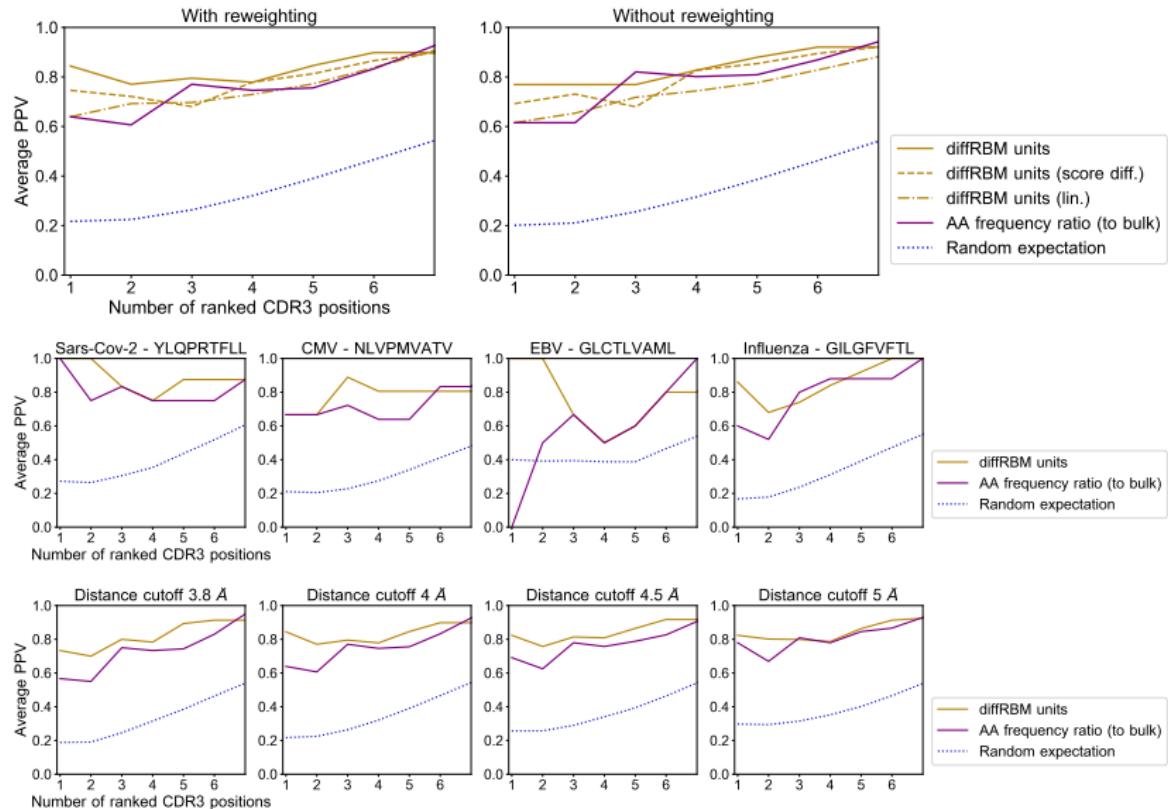
TCR specificity



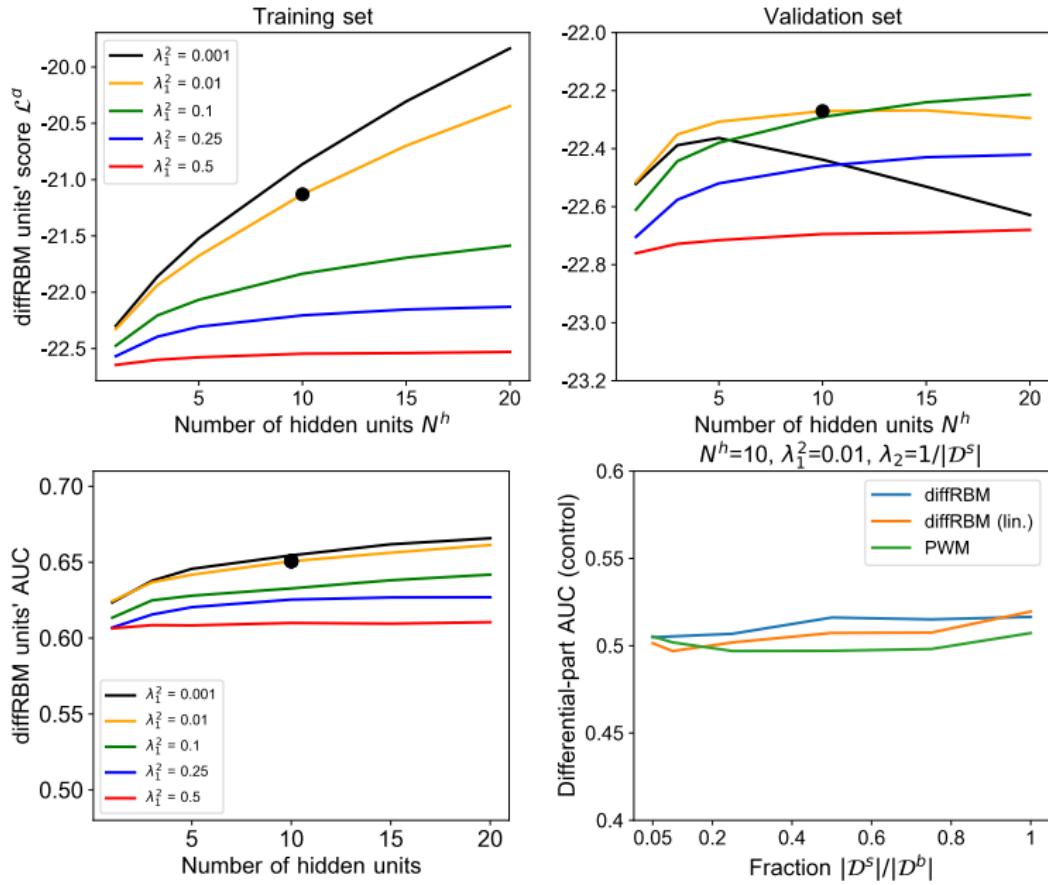
Robustness to choice of naive TCR repertoire

Performance at discriminating specific vs bulk CDR3 β (different background datasets)

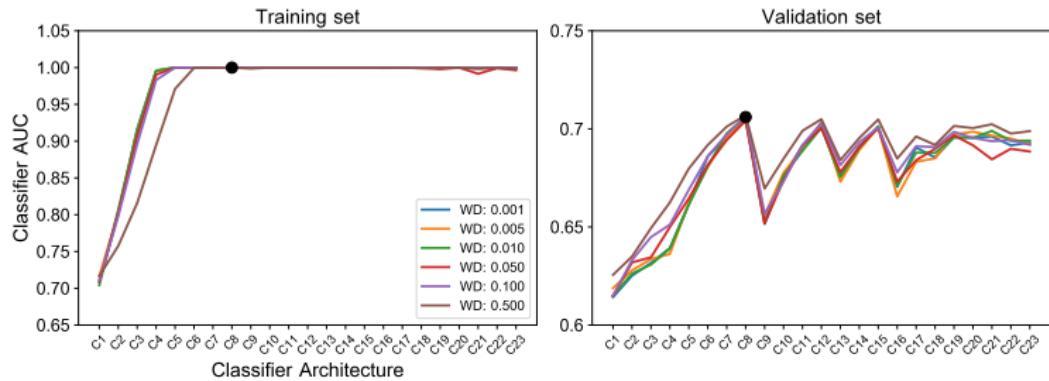




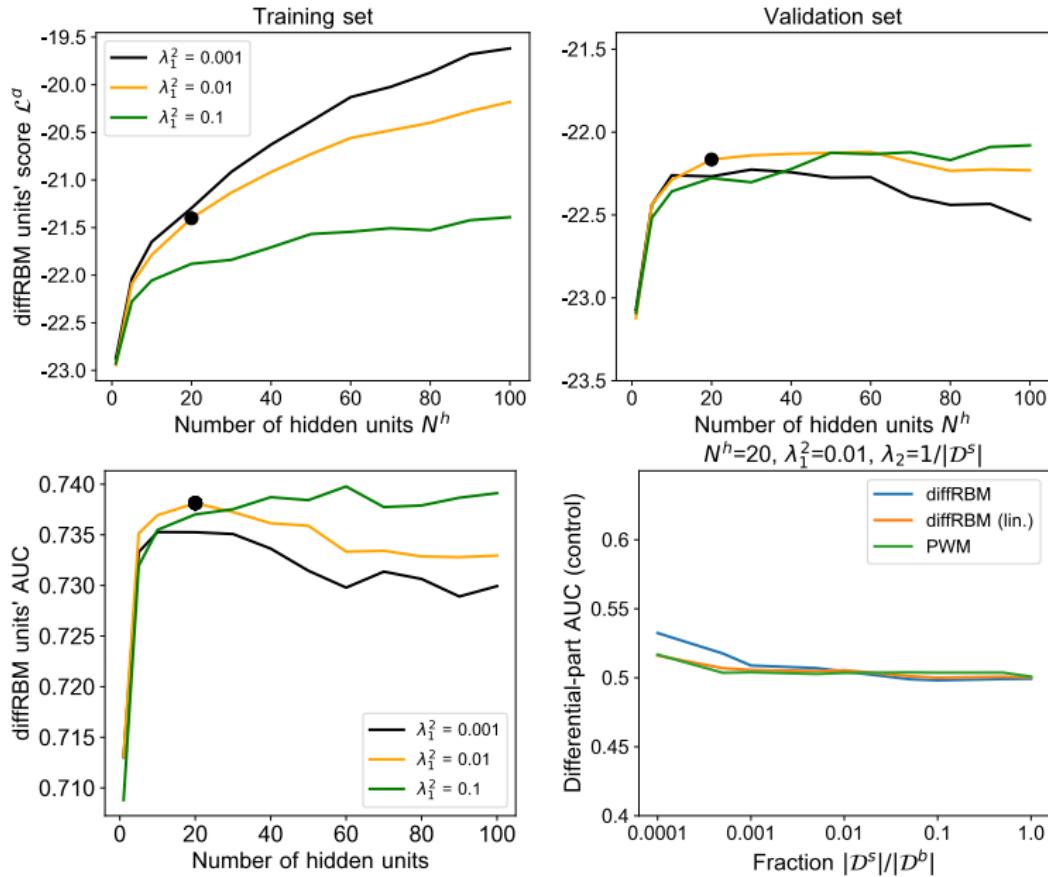
Hyper-parametric search



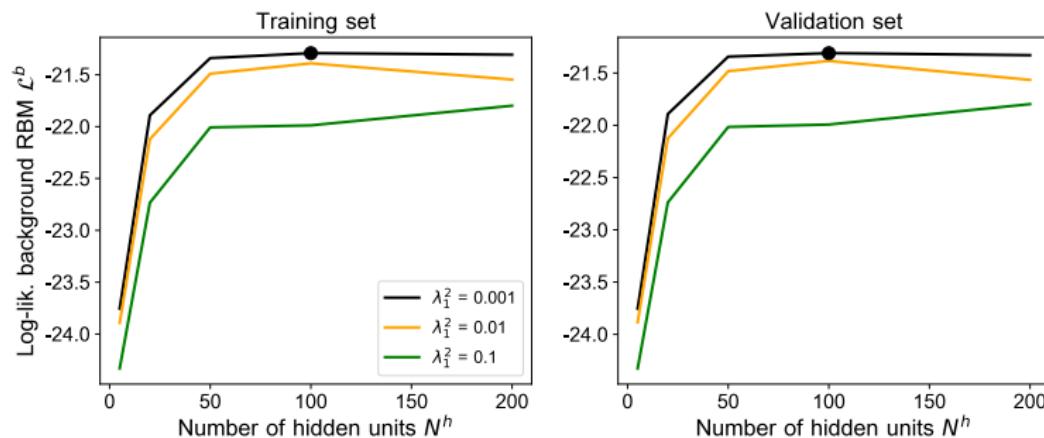
Hyper-parametric search



Hyper-parametric search



Hyper-parametric search



Hyper-parametric search

