# Statistical learning approaches to modelling T cell response at the molecular level

## Barbara Bravi

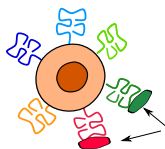Department of Mathematics, Imperial College London

Joint work with: S. Cocco, R. Monasson, T. Mora, A. Walczak (ENS Paris)
**Mini-colloque: Information et Biologie, SFP, 3 July 2023**
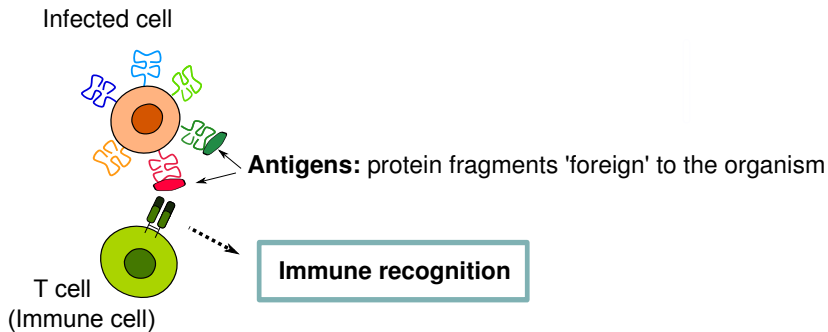
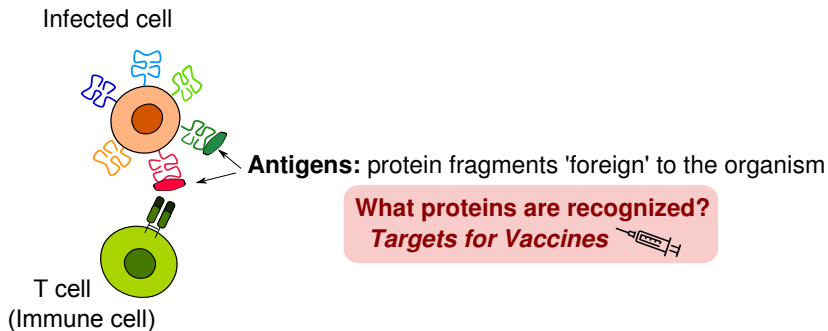**Imperial College**
London

# T cell response

Infected cell



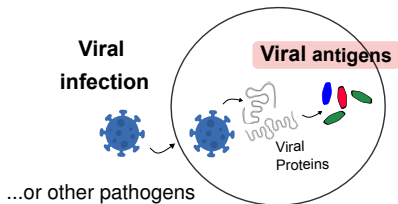**Antigens:** protein fragments 'foreign' to the organism

# T cell response



Infected cell

**Antigens:** protein fragments 'foreign' to the organism

T cell
(Immune cell)

**Immune recognition**

# T cell response



Infected cell

**Antigens:** protein fragments 'foreign' to the organism

**What proteins are recognized?**
*Targets for Vaccines*

T cell
(Immune cell)

# T cell response



Infected cell

**Antigens:** protein fragments 'foreign' to the organism

**What proteins are recognized?**
*Targets for Vaccines*

T cell
(Immune cell)

**Viral infection**

**Viral antigens**

Viral Proteins

...or other pathogens

# T cell response



Infected cell
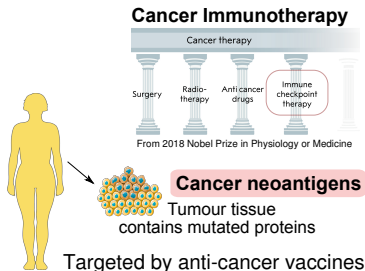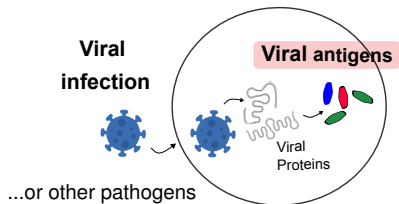
**Antigens:** protein fragments 'foreign' to the organism

**What proteins are recognized?**
*Targets for Vaccines*

T cell
(Immune cell)

**Viral infection**

**Viral antigens**

Viral Proteins

...or other pathogens

**Cancer Immunotherapy**

Cancer therapy

Surgery | Radio-therapy | Anti cancer drugs | Immune checkpoint therapy

From 2018 Nobel Prize in Physiology or Medicine

**Cancer neoantigens**
Tumour tissue contains mutated proteins

Targeted by anti-cancer vaccines

# Mechanisms: protein-protein binding



Infected cell

HLA protein

T cell
(Immune cell)

**Antigen - HLA protein binding: 'Presentation'**

# Mechanisms: protein-protein binding



Infected cell

HLA protein

**Antigen - HLA protein binding: 'Presentation'**

**Receptor - antigen binding: recognition**

T cell
(Immune cell)

# Mechanisms: protein-protein binding



Infected cell

HLA protein

T cell
(Immune cell)

**Antigen - HLA protein binding: 'Presentation'**

**Receptor - antigen binding: recognition**

**Aim: build models of immune interactions
from available protein data**

# The statistical learning approach
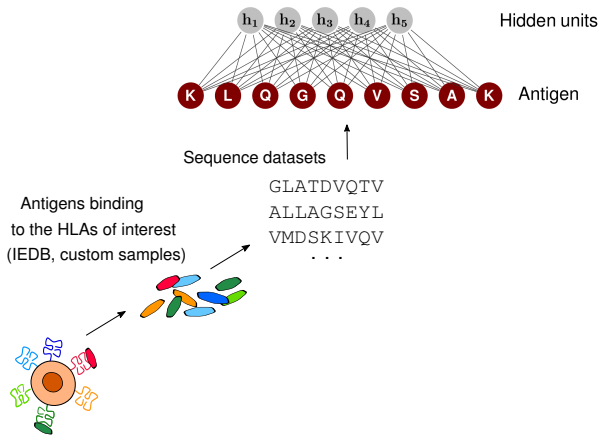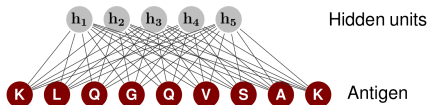


**Restricted Boltzmann Machines (RBMs)**
(Smolensky 1986, Hinton 2002, Tubiana et al. 2019)

# Predicting antigen presentation



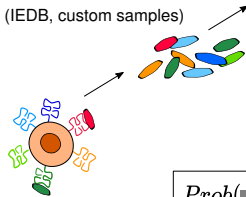**Restricted Boltzmann Machines (RBMs)**
(Smolensky 1986, Hinton 2002, Tubiana et al. 2019)

Hidden units

Antigen

Sequence datasets

GLATDVQTV
ALLAGSEYL
VMDSKIVQV
. . .

Antigens binding
to the HLAs of interest
(IEDB, custom samples)

# Predicting antigen presentation



**Restricted Boltzmann Machines (RBMs)**
(Smolensky 1986, Hinton 2002, Tubiana et al. 2019)

Hidden units

Antigen

Sequence datasets

GLATDVQTV
ALLAGSEYL
VMDSKIVQV
· · ·

Antigens binding
to the HLAs of interest
(IEDB, custom samples)

**Learn from sequence data
a probability of presentation**

$Prob(\blacksquare > \blacksquare) = 0.973$

**RBM discriminates
presented antigens**

Presented antigens

Generic peptides

Presentation probability $log_{10}P(\mathbf{S})$

# Predicting antigen presentation

## Restricted Boltzmann Machines (RBMs)
(Smolensky 1986, Hinton 2002, Tubiana et al. 2019)



Hidden units
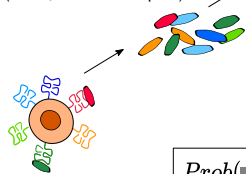
Antigen

Sequence datasets

```
GLATDVQTV
ALLAGSEYL
VMDSKIVQV
. . .
```

Antigens binding
to the HLAs of interest
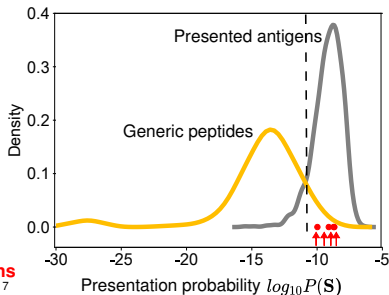(IEDB, custom samples)

**Learn from sequence data
a probability of presentation**

$Prob(■ > ■) = 0.973$

**RBM discriminates
presented antigens**

**Prediction of cancer neoantigens**
Validated neoantigens from Marty et al. 2017



Presented antigens

Generic peptides

Density

Presentation probability $log_{10}P(\mathbf{S})$

# RBM low-dim. representation



**Restricted Boltzmann Machines (RBMs)**
(Smolensky 1986, Hinton 2002, Tubiana et al. 2019)

Hidden units

Antigen

**Dimensionality Reduction:**

*like PCA but non-linear!*
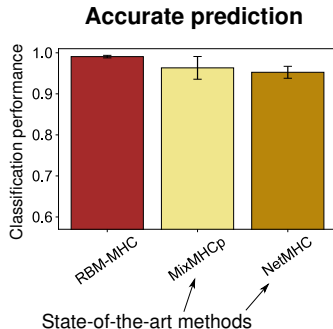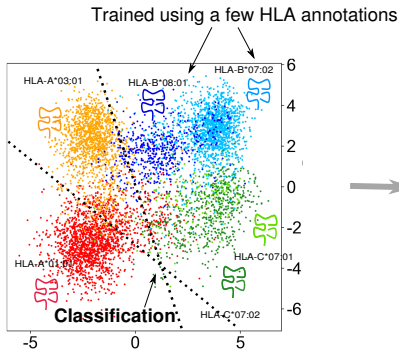
STDPEALHY
AAKKKLQSL
YRAEQINQL
TPRPVTELL
HELGVADRL
LPFKKSLAL

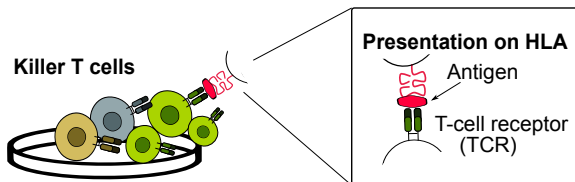**Low-dim. representation**

Clusters =
HLA binding specificity

HLA-A*03:01   HLA-B*08:0T   HLA-B*07:02
HLA-A*01    HLA-C*07 01
HLA-C*07:02

# Prediction of HLA specificity



Trained using a few HLA annotations

**Accurate prediction**

B. Bravi, J. Tubiana, S. Cocco, R. Monasson, T. Mora, A.M. Walczak, *RBM-MHC: a semi-supervised machine-learning method for sample-specific prediction of antigen presentation by HLA-I alleles*, *Cell Systems* (2021)

# Antigen immunogenicity
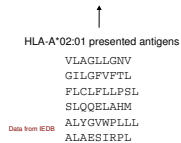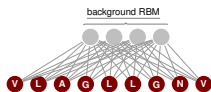
**Presentation alone is not immunogenicity!**



Only a fraction of HLA-presented antigens are immunogenic (promote a T cell response).
**Immunogenicity prediction: still low success rate** (Wells et al. 2020, Buckley et al. 2022)
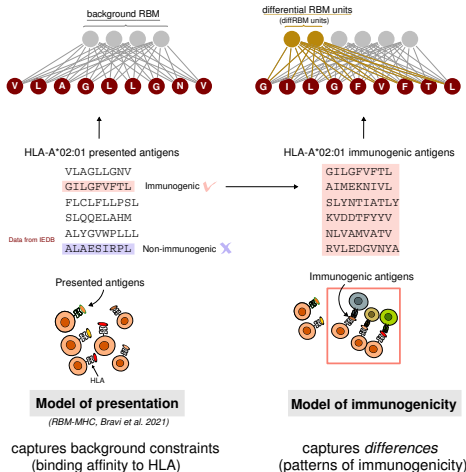
# Transfer learning with 'differential' units



background RBM

**V L A G L L G N V**

HLA-A*02:01 presented antigens

```
VLAGLLGNV
GILGFVFTL
FLCLFLLPSL
SLQQELAHM
ALYGVWPLLL
ALAESIRPL
```

Data from IEDB

Presented antigens

HLA

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

captures background constraints
(binding affinity to HLA)

B. Bravi, A. Di Gioacchino, J. Fernandez-de-Cossio-Diaz, A.M. Walczak, T. Mora, S. Cocco, R. Monasson, pre-print Biorxiv 2022.12.06.519259v1 (2022)

# Transfer learning with 'differential' units



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

Model of immunogenicity

captures background constraints
(binding affinity to HLA)

captures *differences*
(patterns of immunogenicity)

B. Bravi, A. Di Gioacchino, J. Fernandez-de-Cossio-Diaz, A.M. Walczak, T. Mora, S. Cocco, R. Monasson, pre-print Biorxiv 2022.12.06.519259v1 (2022)
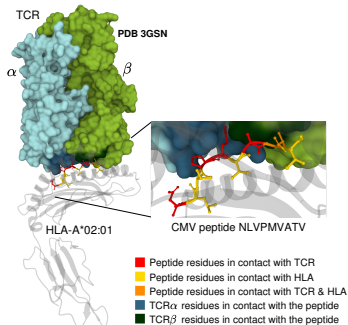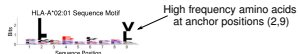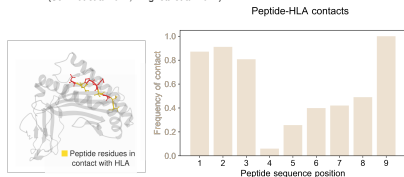
# Transfer learning with 'differential' units
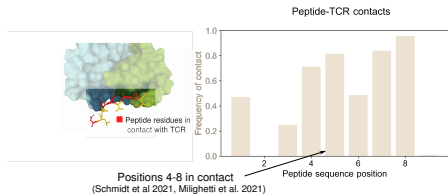
B. Bravi, A. Di Gioacchino, J. Fernandez-de-Cossio-Diaz, A.M. Walczak, T. Mora, S. Cocco, R. Monasson, pre-print Biorxiv 2022.12.06.519259v1 (2022)
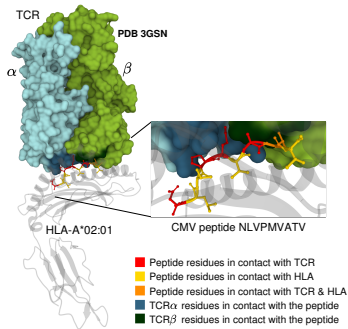
# Contact positions in resolved structures:



**TCR**

**PDB 3GSN**

$\alpha$  $\beta$

**HLA-A*02:01**

CMV peptide NLVPMVATV

- ■ Peptide residues in contact with TCR
- ■ Peptide residues in contact with HLA
- ■ Peptide residues in contact with TCR & HLA
- ■ TCR$\alpha$ residues in contact with the peptide
- ■ TCR$\beta$ residues in contact with the peptide

## HLA-A*02:01-specific peptides

### Peptide-TCR contacts



■ Peptide residues in contact with TCR

Positions 4-8 in contact
(Schmidt et al 2021, Milighetti et al. 2021)

### Peptide-HLA contacts



■ Peptide residues in contact with HLA

HLA-A*02:01 Sequence Motif

High frequency amino acids at anchor positions (2,9)
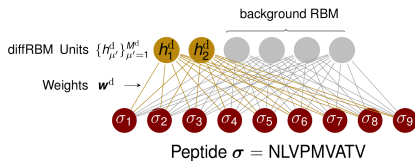
Constraints on statistics of immunogenic antigens should reflect contacts

# Contact positions in resolved structures:



TCR

**PDB 3GSN**

$\alpha$     $\beta$

HLA-A*02:01

CMV peptide NLVPMVATV

- 🟥 Peptide residues in contact with TCR
- 🟨 Peptide residues in contact with HLA
- 🟧 Peptide residues in contact with TCR & HLA
- 🟦 TCR$\alpha$ residues in contact with the peptide
- 🟩 TCR$\beta$ residues in contact with the peptide

# Model prediction:

## diffRBM architecture



background RBM

diffRBM Units $\{h^{\mathrm{d}}_{\mu'}\}^{M^{\mathrm{d}}}_{\mu'=1}$   $h^{\mathrm{d}}_1$   $h^{\mathrm{d}}_2$

Weights $\boldsymbol{w}^{\mathrm{d}} \rightarrow$

$\sigma_1 \; \sigma_2 \; \sigma_3 \; \sigma_4 \; \sigma_5 \; \sigma_6 \; \sigma_7 \; \sigma_8 \; \sigma_9$

Peptide $\boldsymbol{\sigma} = $ NLVPMVATV

## Single-site importance factors

$$T_i(\sigma_i) = \underbrace{g^{\mathrm{d}}_i(\sigma_i)}_{} + \sum_{\mu'=1}^{M^{\mathrm{d}}} \underbrace{w^{\mathrm{d}}_{i\mu'}(\sigma_i)\langle h_{\mu'}|\boldsymbol{\sigma}\rangle}_{}$$

related to amino acid frequency difference between immunogenic and presented

captures correlations between positions

$\langle h_{\mu'}|\boldsymbol{\sigma}\rangle$: from $P(h_{\mu'}|I_{\mu'}(\boldsymbol{\sigma}))$, where $I_{\mu'}(\boldsymbol{\sigma}) = \sum_i w^{\mathrm{d}}_{i\mu'}(\sigma_i)$

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

# Structural interpretation



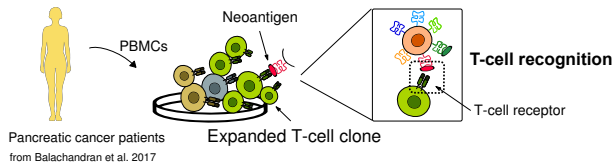**HLA-A*02:01-specific peptides**

Peptide-TCR contacts

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
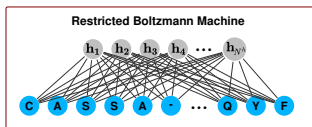
# T cell response specificity
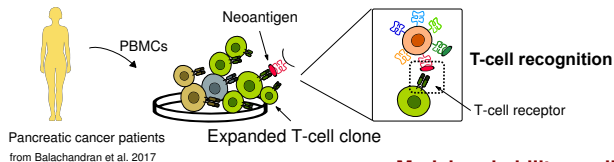


Bravi et al., PLoS Comput. Biol. (2021)

# T cell response specificity



Bravi et al., PLoS Comput. Biol. (2021)

**Model probability predicts specificity of response**

# T cell response specificity



Bravi et al., PLoS Comput. Biol. (2021)

**Uncover sequence patterns determining specificity**

# T cell response specificity

We want a summary metric of these convergence in amino acid patterns
**Dissimilarity Index:**

$$DI = \frac{1}{f} \qquad f = \frac{1}{T} \sum_{i<j} e^{-\left(\frac{d(\sigma_i, \sigma_j)}{\delta}\right)^2}$$

<u>Distance between CDR3 $\sigma_i$ and $\sigma_j$</u>

(CDR3-only version of TCRdist from Dash et al 2017)

# T cell response specificity



Bravi et al., PLoS Comput. Biol. (2021)

We want a summary metric of these convergence in amino acid patterns
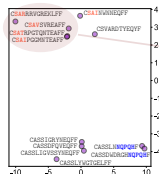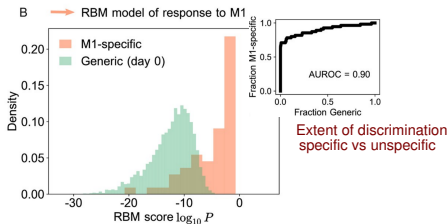**Dissimilarity Index:**

$$DI = \frac{1}{f} \qquad f = \frac{1}{T} \sum_{i<j} e^{-\left(\frac{d(\sigma_i, \sigma_j)}{\delta}\right)^2}$$

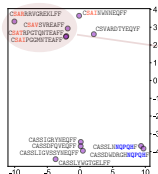Distance between CDR3 $\sigma_i$ and $\sigma_j$

(CDR3-only version of TCRdist from Dash et al 2017)

We consider: neoantigen-specific TCRs from PBMCs,

tetramer-sorted TCRs specific to viral epitopes (e.g. M1)



B → RBM model of response to M1

Extent of discrimination
specific vs unspecific

# T cell response specificity



Bravi et al., PLoS Comput. Biol. (2021)

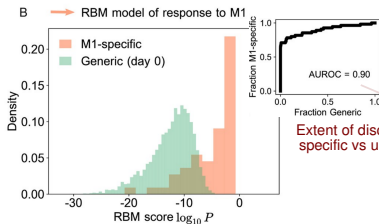We want a summary metric of these convergence in amino acid patterns

**Dissimilarity Index:**

$$DI = \frac{1}{f} \qquad f = \frac{1}{T} \sum_{i<j} e^{-\left(\frac{d(\sigma_i, \sigma_j)}{\delta}\right)^2}$$
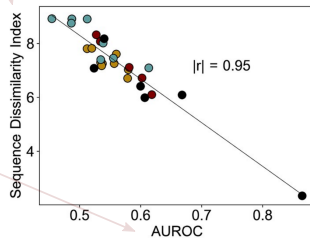
Distance between CDR3 $\sigma_i$ and $\sigma_j$

(CDR3-only version of TCRdist from Dash et al 2017)

We consider: neoantigen-specific TCRs from PBMCs,

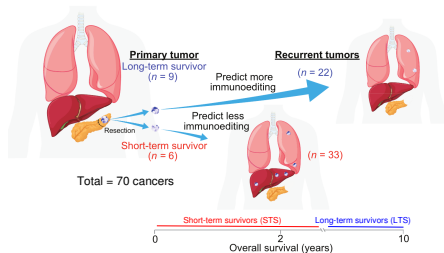tetramer-sorted TCRs specific to viral epitopes (e.g. M1)



$|r| = 0.95$

Extent of discrimination
specific vs unspecific

# Immunoediting in pancreatic cancer



**Large-scale study on immunoediting in PDAC**
conducted by Vinod Balachandran, Benjamin Greenbaum,
Marta Łuksza, Zachary Sethna, Luis Rojas, Jayon Lihm and others

Primary tumor
Long-term survivor
(n = 9)

Recurrent tumors
(n = 22)

Predict more
immunoediting

Resection

Predict less
immunoediting

Short-term survivor
(n = 6)

(n = 33)

Total = 70 cancers

Short-term survivors (STS)          Long-term survivors (LTS)
0                    2                    10
Overall survival (years)

Łuksza* , Sethna*, Rojas*, Lihm, Bravi et al., Nature (2022)

# Immunoediting in pancreatic cancer

**Large-scale study on immunoediting in PDAC**
  conducted by Vinod Balachandran, Benjamin Greenbaum,
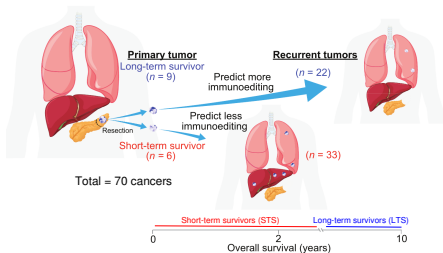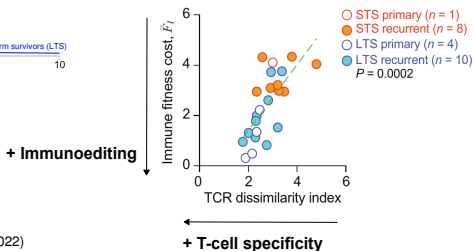Marta Łuksza, Zachary Sethna, Luis Rojas, Jayon Lihm and others



Our T-cell specificity analysis provides
additional statistical evidence of immunoediting

Łuksza* , Sethna*, Rojas*, Lihm, Bravi et al., Nature (2022)

## Summary

**Statistical learning approach based on Restricted Boltzmann Machines**

- Amino acid patterns $\rightarrow$ scores of molecular specificity (Antigen presentation, immunogenicity, T cell response)

# Summary

**Statistical learning approach based on Restricted Boltzmann Machines**

- Amino acid patterns $\rightarrow$ scores of molecular specificity (Antigen presentation, immunogenicity, T cell response)

- Transfer-learning approach extracts biologically interpretable features on immunogenicity

# ACKNOWLEDGEMENTS

Thank you for your attention!