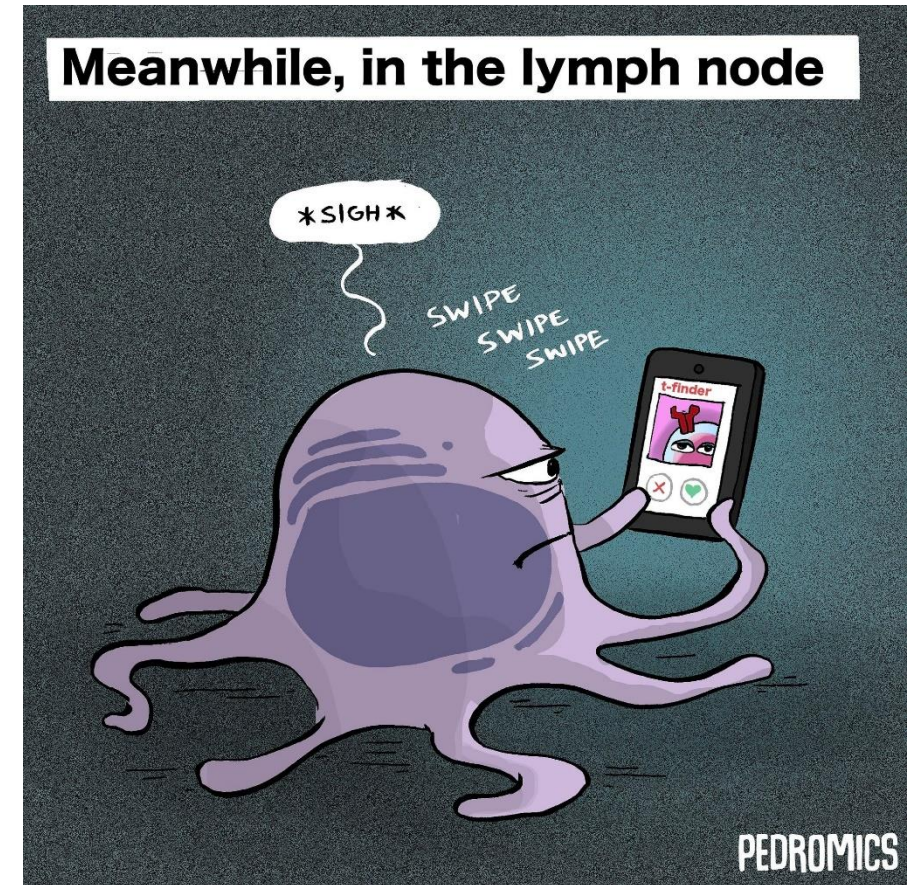
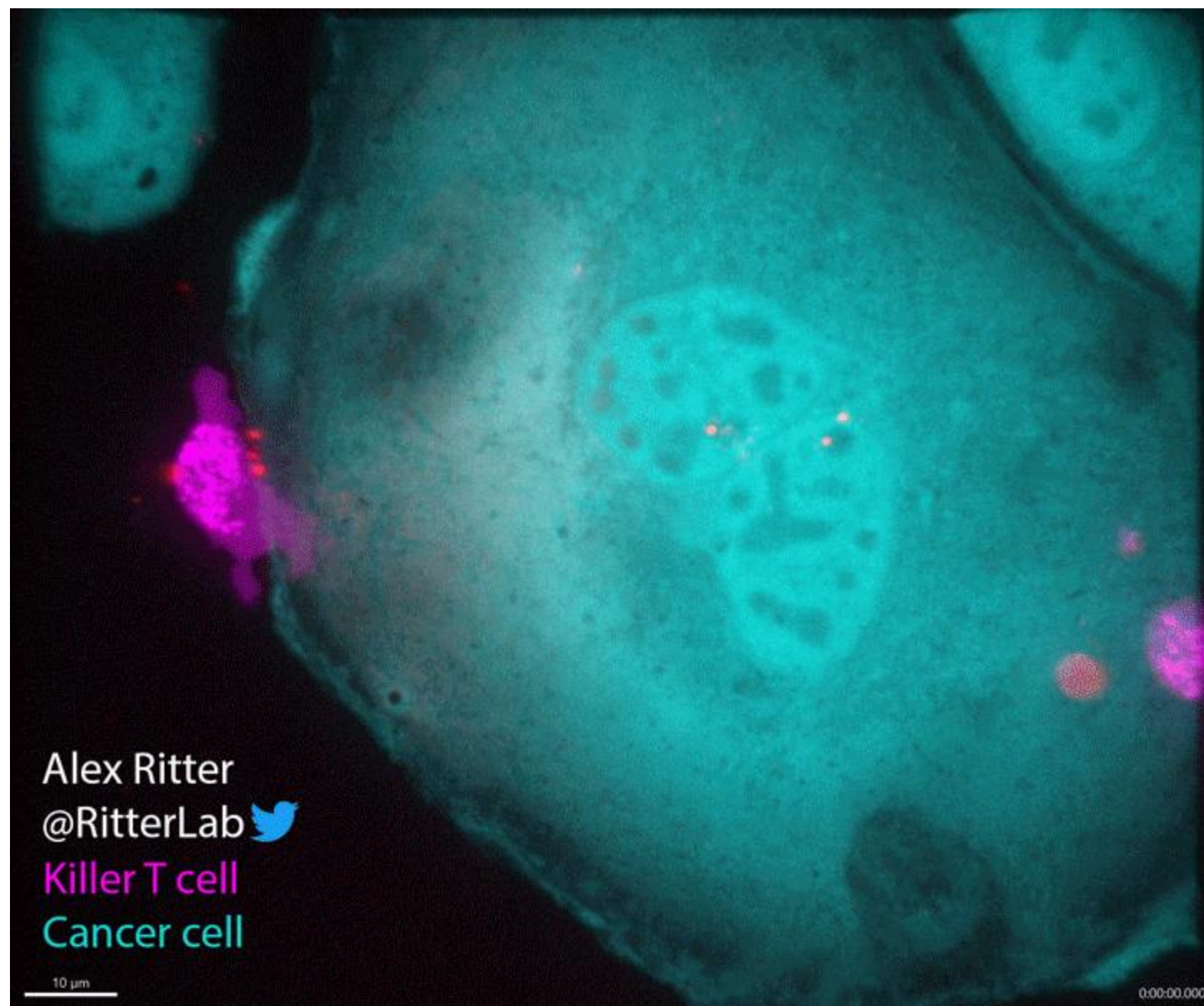


Kavli Institute for Theoretical Physics
October 11th, 2024

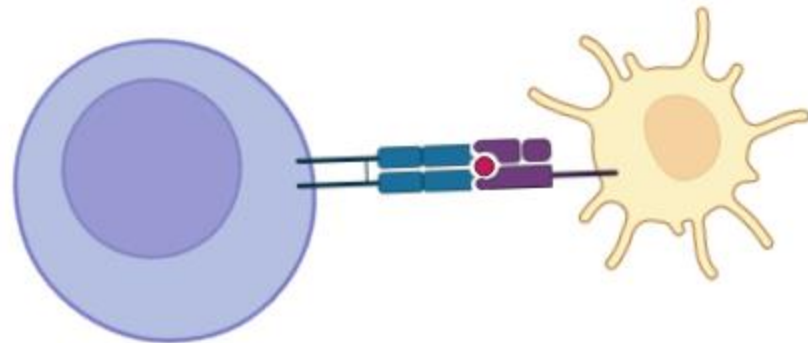
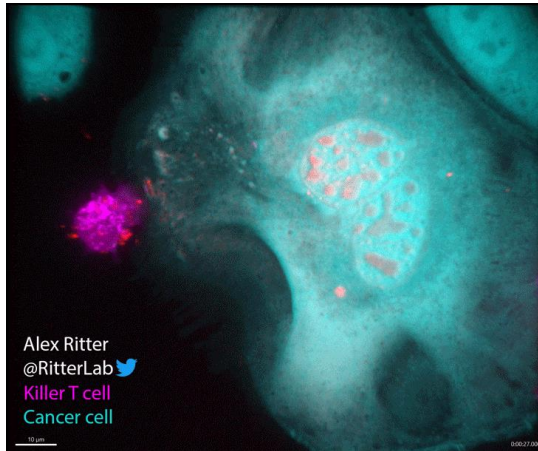
Statistical insights into the immune receptor code





Alex Ritter
@RitterLab 
Killer T cell
Cancer cell

How do lymphocytes know what to attack?

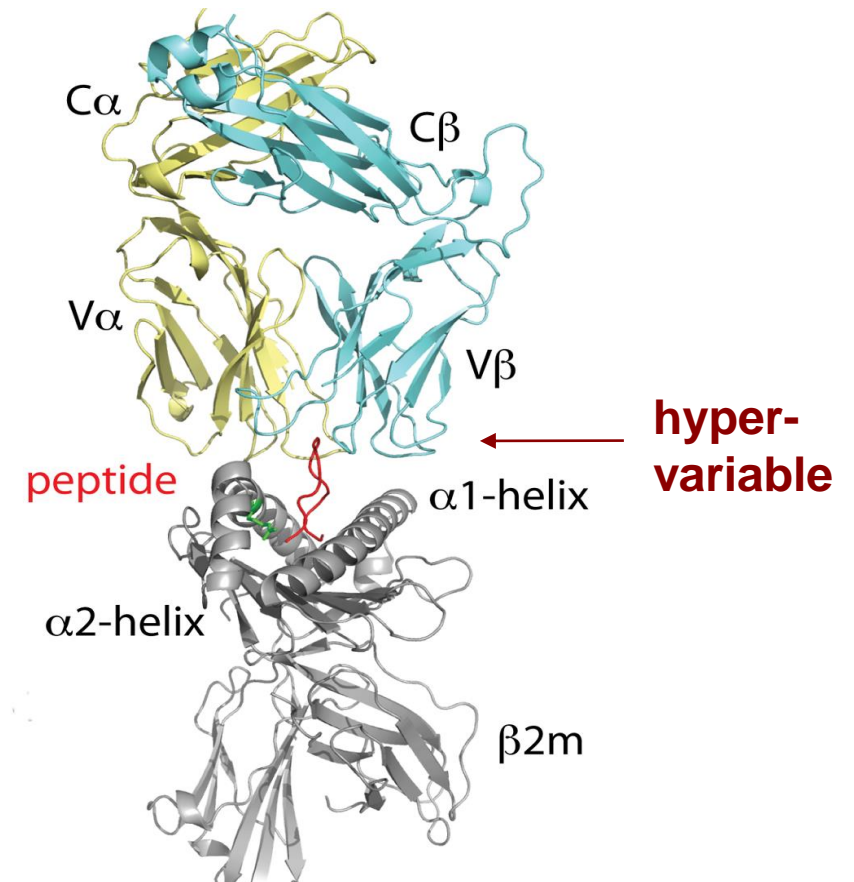


Heroes of today's talk: T cell receptors

T cell receptor (TCR)



Peptide presented on target cell



Outline

1. **Why reading the immune receptor code is hard and why it matters**
Sureshchandra, Henderson, ..., ATM*, Saligrama*, Wagar* bioRxiv 2024
2. **The state of *in silico* TCR-pMHC binding prediction**
Nielsen ... ATM ... Barton Immunoinformatics 2024
Nagano, Pyo, ..., Chain*, ATM* arXiv 2024
3. **Statistical physics of sequence-function maps**
ATM, Callan PNAS 2023
4. **Insights into the biophysical determinants TCR specificity**
Henderson, Nagano, Milighetti, ATM PNAS 2024
Nagano, Pyo, ..., Chain*, ATM* arXiv 2024

Which question would you ask if we could predict antigen-specificity from antibody/TCR sequence?

- **Basic human immunology**

- ✓ Dynamics of immune response and memory at the level of single clones?
- ✓ How stochastic or biased are cell fate decisions taken by antigen-specific populations *in vivo*?
- ✓ ...?

- **Vaccine design and protein engineering**

- ✓ What clones and functional states are associated with protection from disease?
- ✓ Can we use the epitopes they target to elicit such responses using vaccines?
- ✓ Rational protein engineering of TCR T-cells / antibodies ?
- ✓ ...?

Heroes of today's talk (cont.)

Yuta Nagano



Martina Milighetti



Andrew Pyo



James Henderson



UCL

Benny Chain

John Shawe-Taylor

Sankalan Bhattacharyya

Rishika Saxena

Touchchai Chotisorayuth

Ursule Demaël

Curtis Callan (Princeton)

Ned Wingreen (Princeton)

Lisa Wagar (UC Irvine)

Naresha Saligrama (WUSTL)

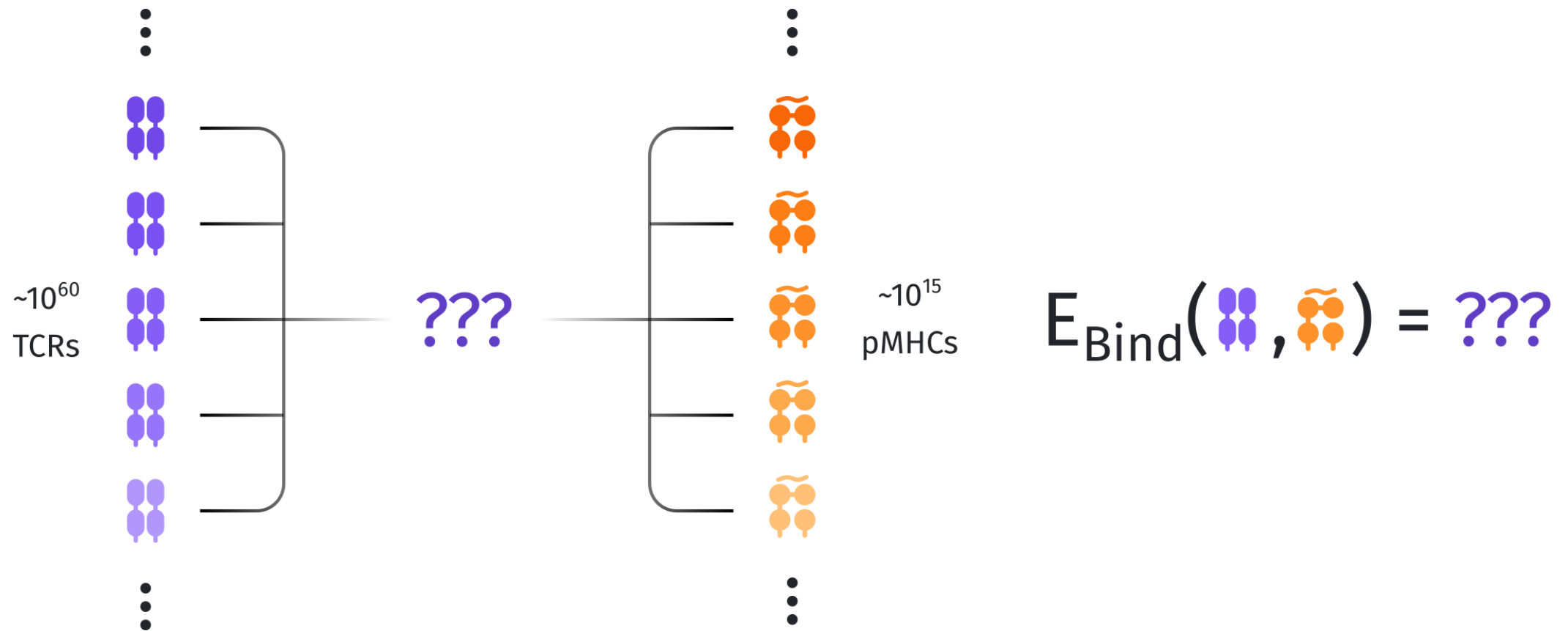
leap^w



Why reading the immune receptor code is hard and why it matters

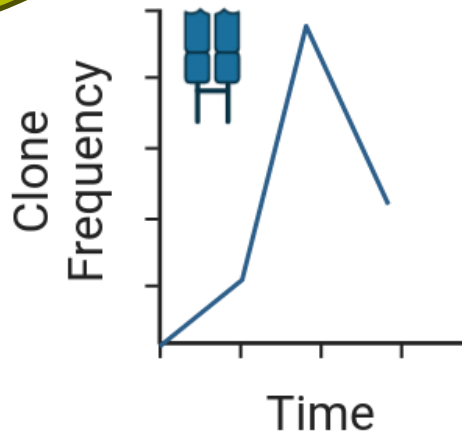
Sureshchandra, Henderson, ..., ATM*, Saligrama*, Wagar*, bioRxiv 2024

Receptors and ligands are immensely diverse

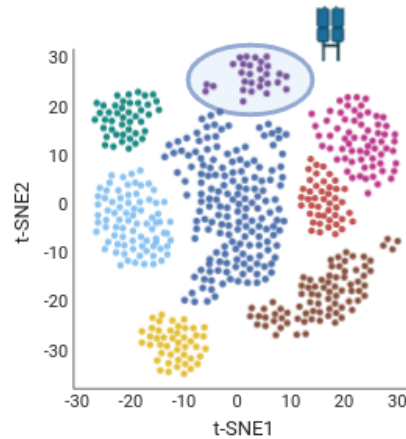


What does this mean for TCRseq?

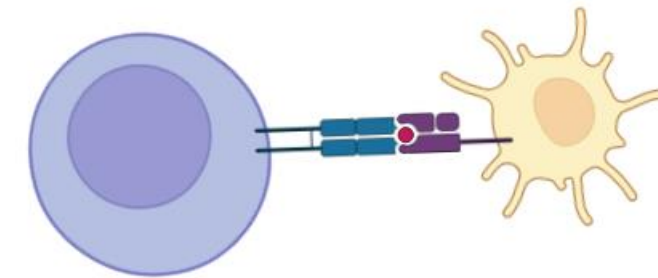
TCRs make great barcodes



Lineage Tracing

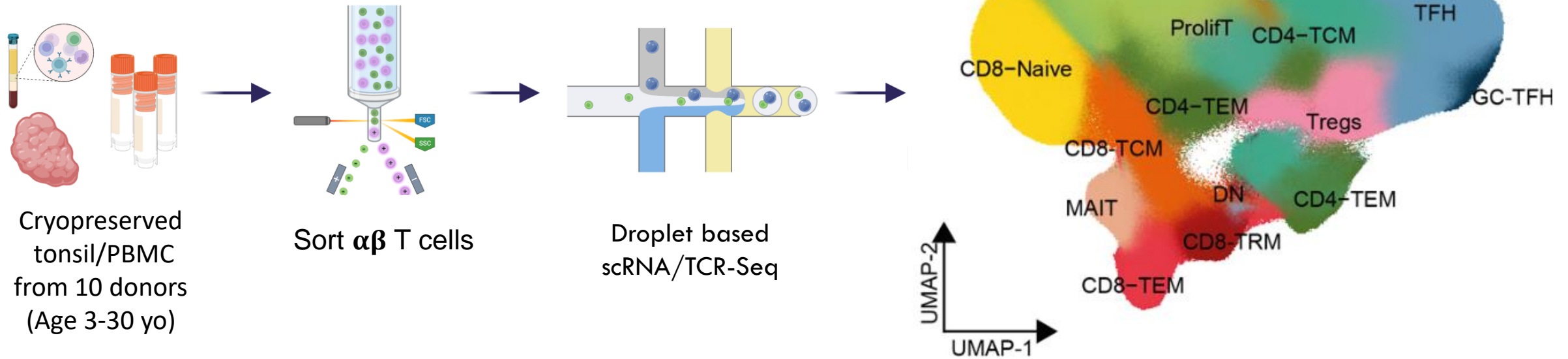


Annotating antigens is hard

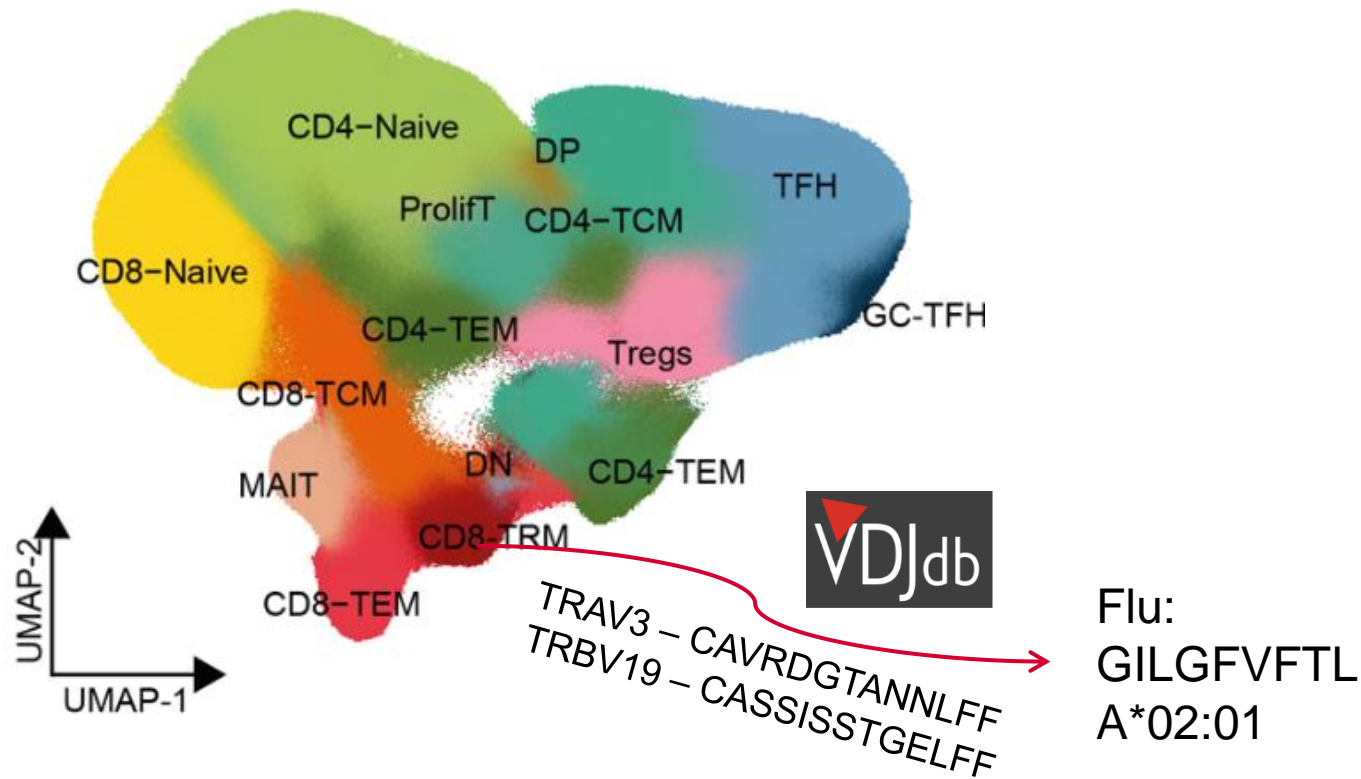


Antigen Deconvolution

DeepTCR study: Profiling the human blood/tonsil repertoire of 5.6 million single T cells



Linking T cell function to antigen specificity



Only 24 TCRs could be mapped this way

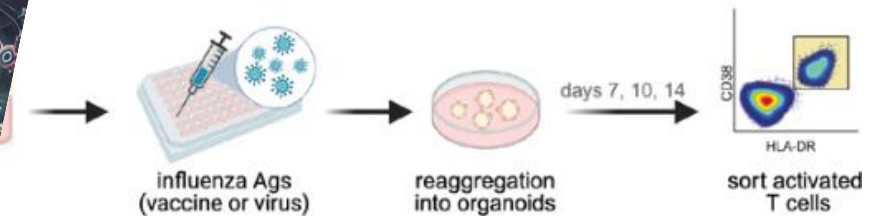
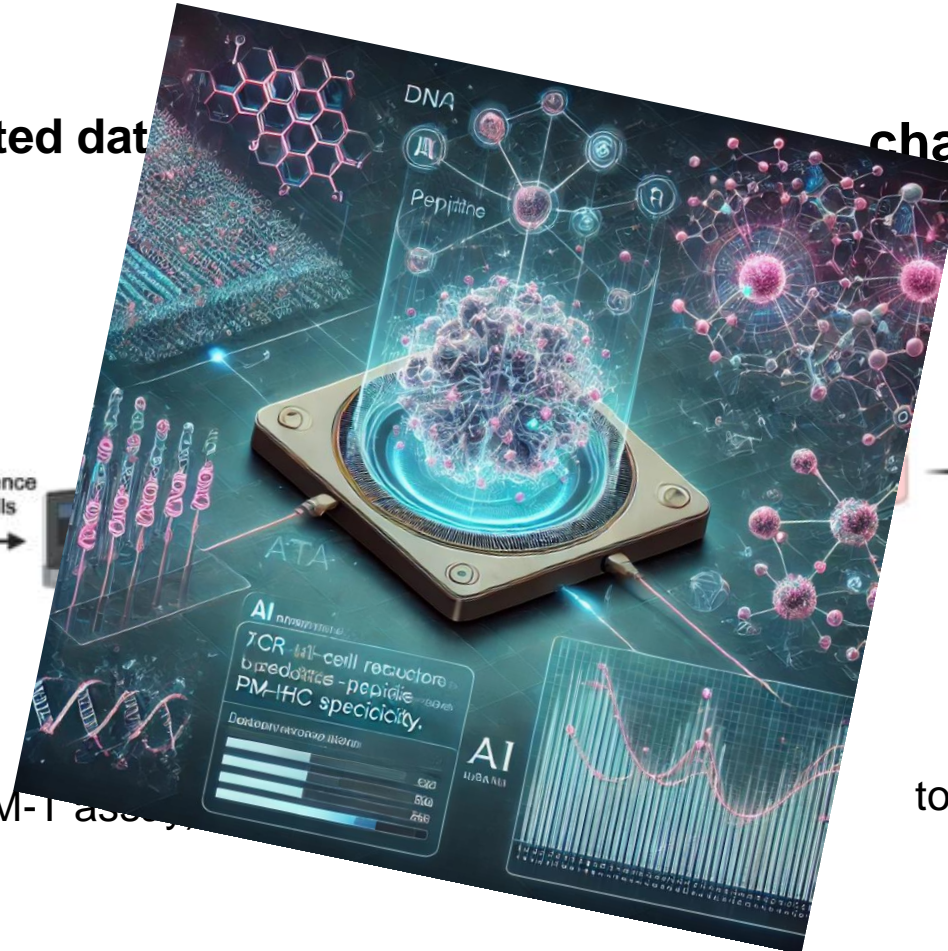
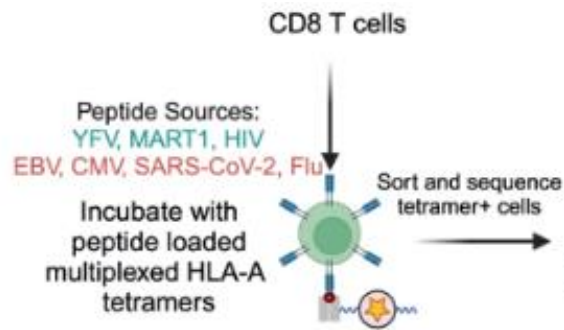
>99.99% of TCRs are orphan

How can we bridge the annotation gap?

>99.99% of TCRs are orphan

generate more annotated data

challenge experiments



antigen-specific sort (BEAM-Flow)

tonsil organoids challenged with flu

The state of *in silico* TCR-pMHC binding prediction

Nielsen ... ATM ... Barton Immunoinformatics 2024

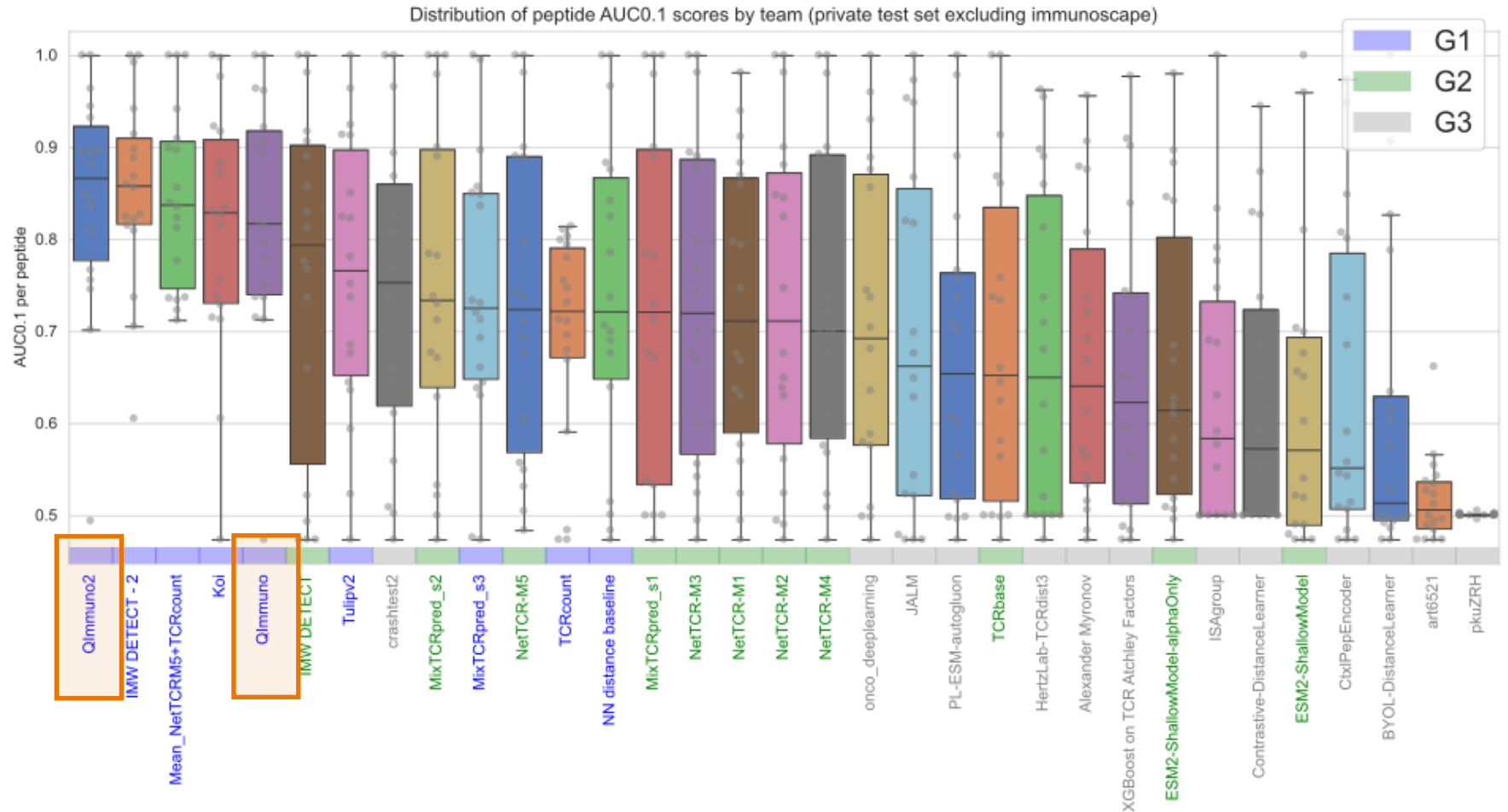
Nagano, Pyo, Milighetti, Henderson, Shawe-Taylor, Chain*, ATM* arXiv 2024

AI to the rescue? Predict TCR-pMHC binding from sequence

Immrep23
benchmark

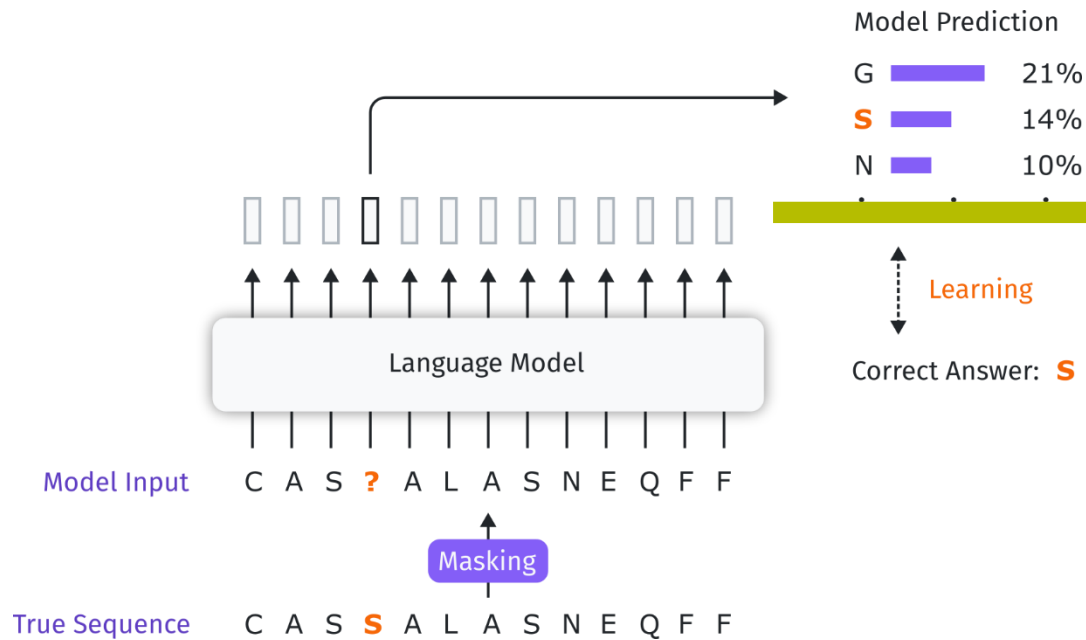
For each TCR
determine most
likely target among
20 pMHCs

No evidence of
generalization of current
deep learning models

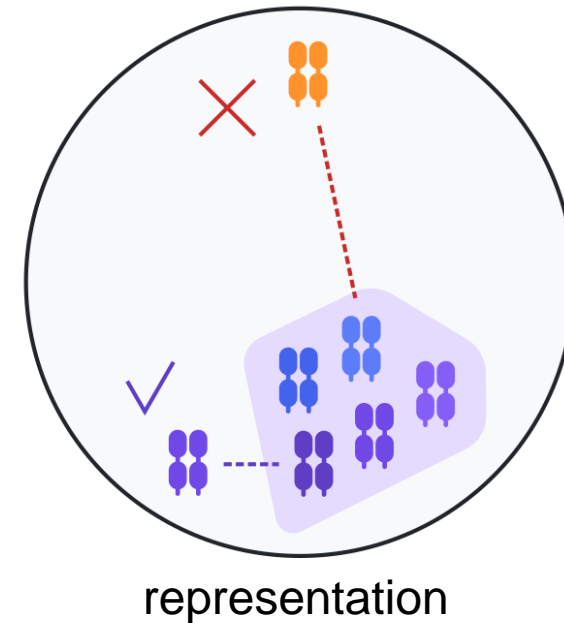


The premise of using large language models

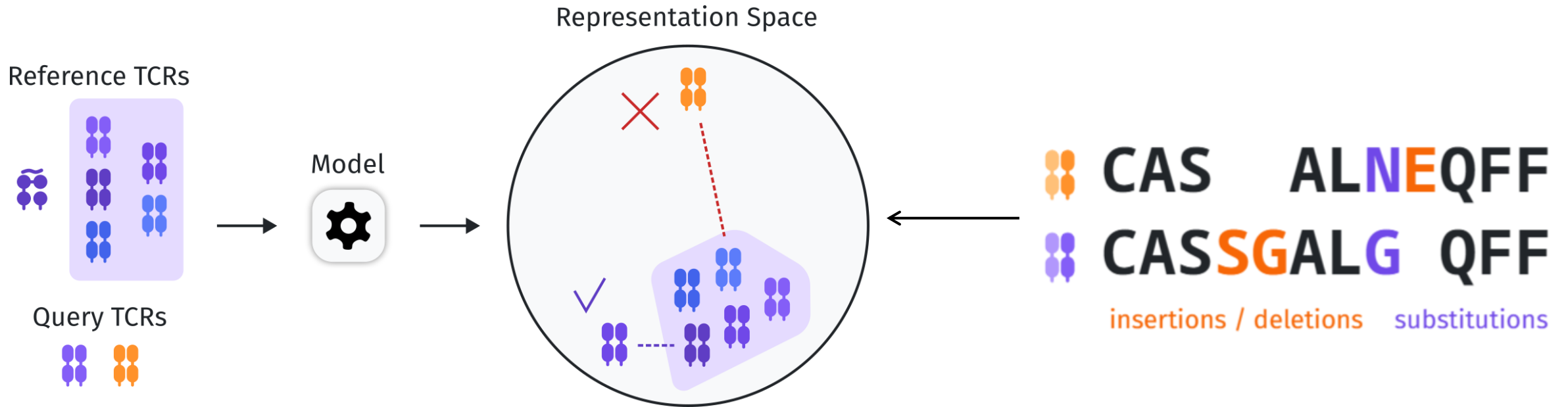
pre-training



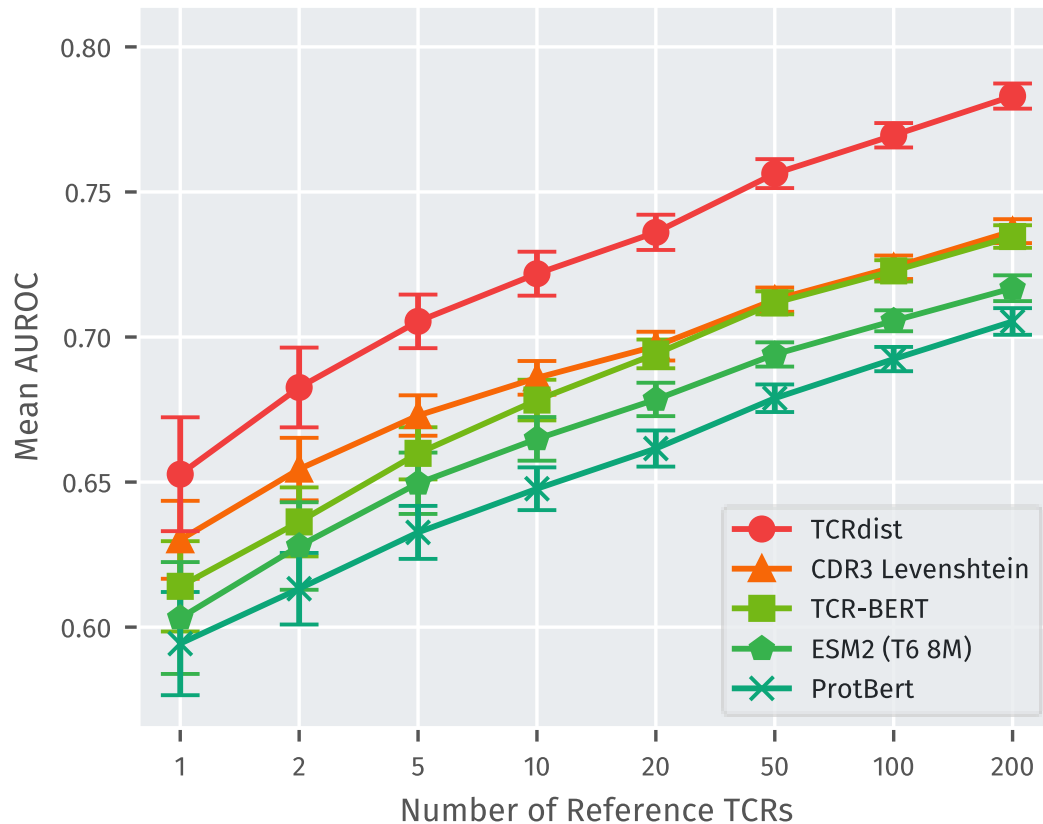
transfer learning



Benchmarking language models against sequence alignment



Language embeddings are inferior to sequence alignment for TCR-pMHC prediction



Nearest neighbour prediction

Positives: binders not used as reference
 Negatives: all remaining TCRs

Data: 6 pMHC with > 300 $\alpha\beta$ TCRs
 Average over pMHC / data splits

A statistical physics view on sequence-function maps

ATM, Callan PNAS 2023

Pyo, Nagano, Milighetti, Henderson, Callan, Chain, Wingreen, ATM (in preparation)

Each ligand defines a selection landscape

Repertoire

$$P(\sigma)$$

CASSWNGPTYEQYF
 CSAPLGGGEQFF
 CADPFRDRGSNQPQHF
 CASSPGQGSYEQYF
 CAISGQGTGEKYQPQHF
 CSARDGTGNGYTF
 CASSTIEGQGGRHTQYF

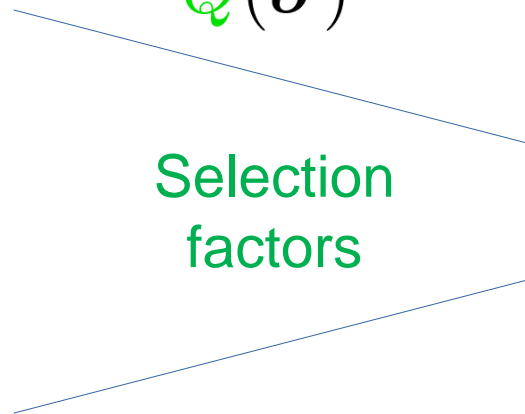
pMHC-specific repertoire

$$Q(\sigma)P(\sigma)$$

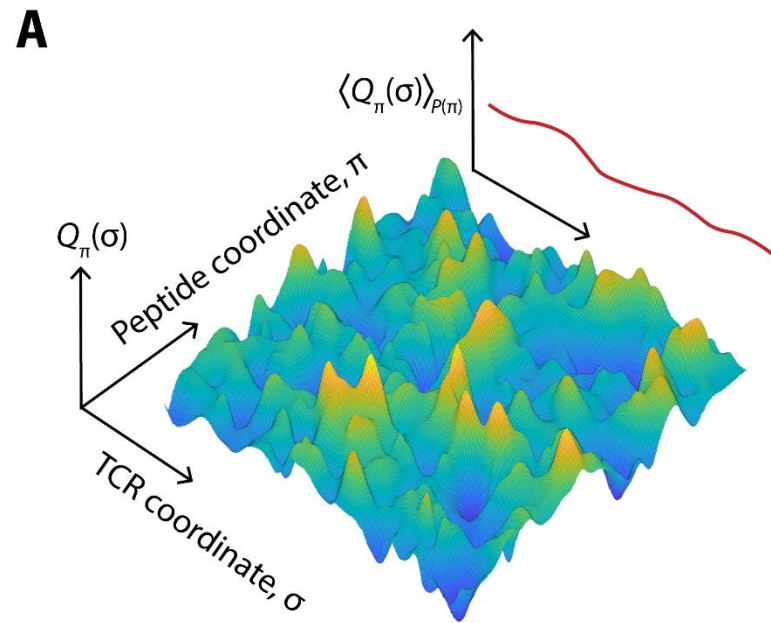
CSARDRVGNGYTF
 CASSPGQGSYEQYF
 CSARDGTGNGYTF
 CSARDGTGNGYTF
 ...

$$Q(\sigma)$$

Selection
factors

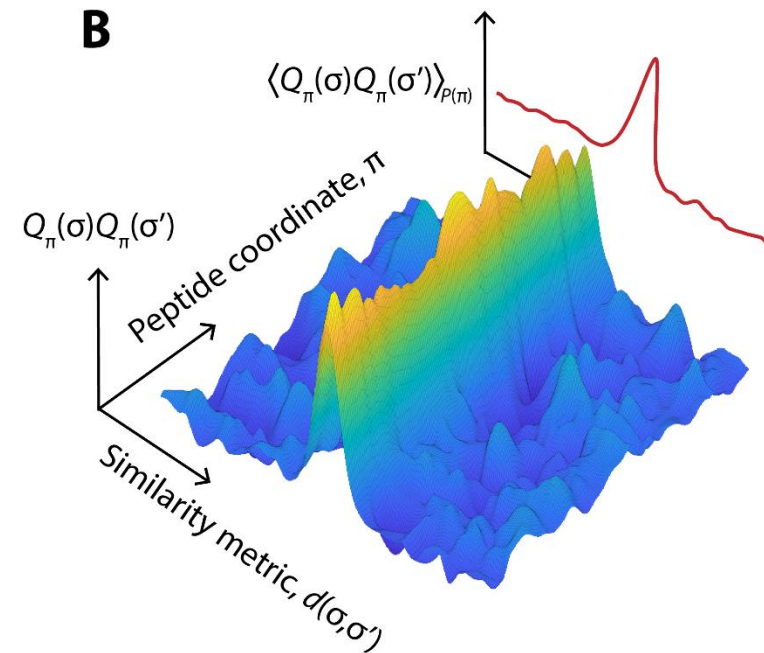


Averaging selection landscapes is destructive



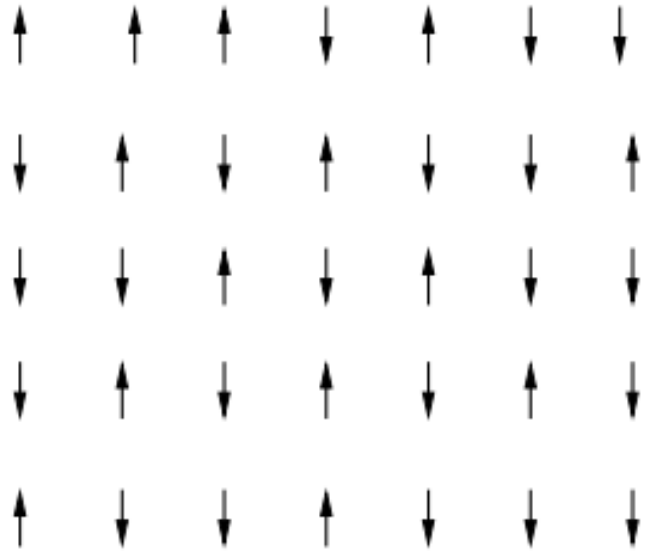
Expect flat marginal distribution

Modest selection relative to Pgen (Elhanati PNAS 2014)



Intuition:
 similarity-relationships are generalizable

Analogy with spin glasses



Zero magnetization

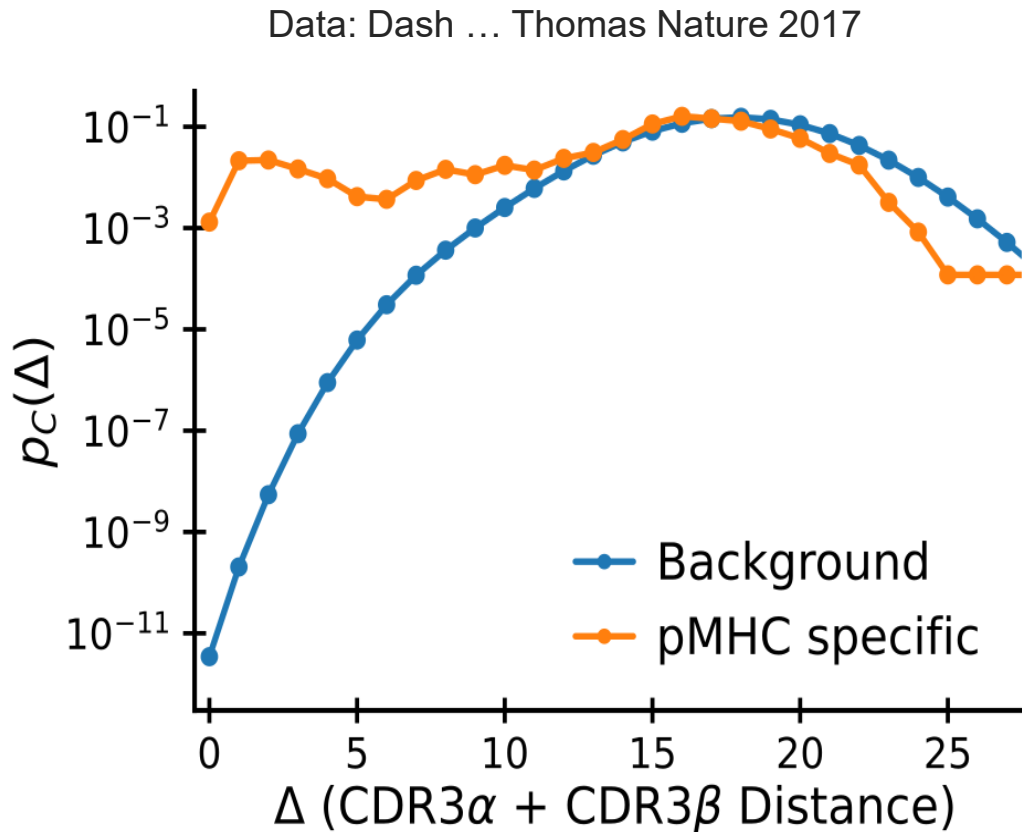
$$m = \langle \sigma_i \rangle = 0$$

Non-trivial order parameter

$$q_{EA} = \lim_{t \rightarrow \infty} \langle \sigma_i(0) \sigma_i(t) \rangle \neq 0$$

$$H = \sum_{\langle i,j \rangle} J_{ij} \sigma_i \sigma_j$$

Convergent TCRs in pMHC-specific repertoires



Histogram of pairwise sequence distances

$$p_C(\Delta) = \sum_{\{\sigma, \sigma'\}} P(\sigma)P(\sigma') \mathbf{I}_{d(\sigma, \sigma')=\Delta}$$

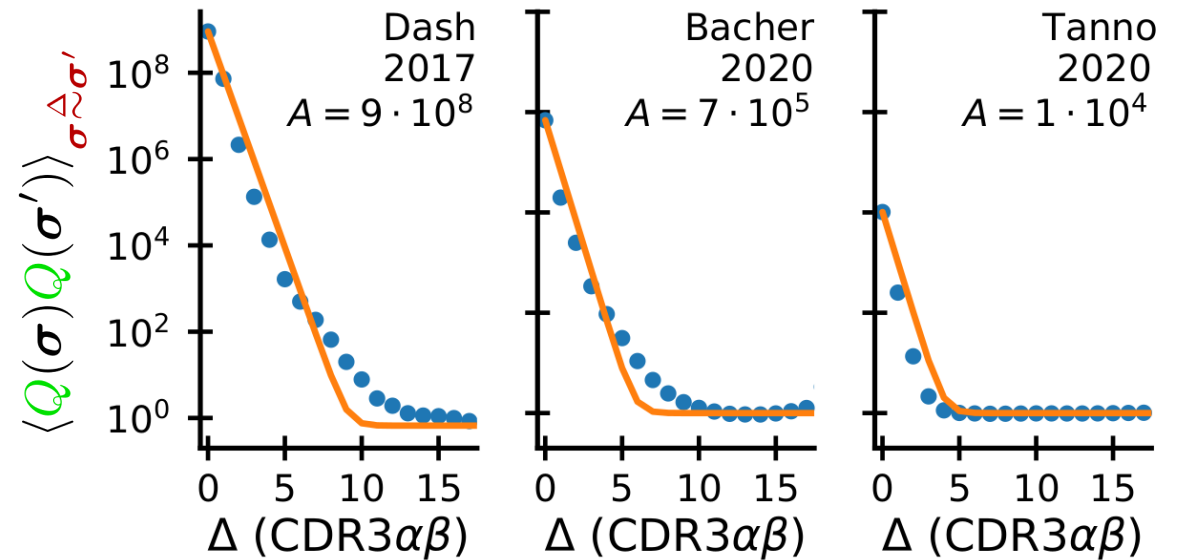
Sequences at distance Δ

Sequence (near-) coincidences are highly enriched

Sequence similarity relates to co-specificity in a generalizable manner

$$\frac{p_C[QP](\Delta)}{p_C[P](\Delta)} = \langle Q(\sigma)Q(\sigma') \rangle_{\sigma \approx \sigma'} \quad \text{Pairs at } d=\Delta$$

Near-coincidence ratios
= how selection co-varies with sequence



Some insights into the biophysical determinants TCR co-specificity

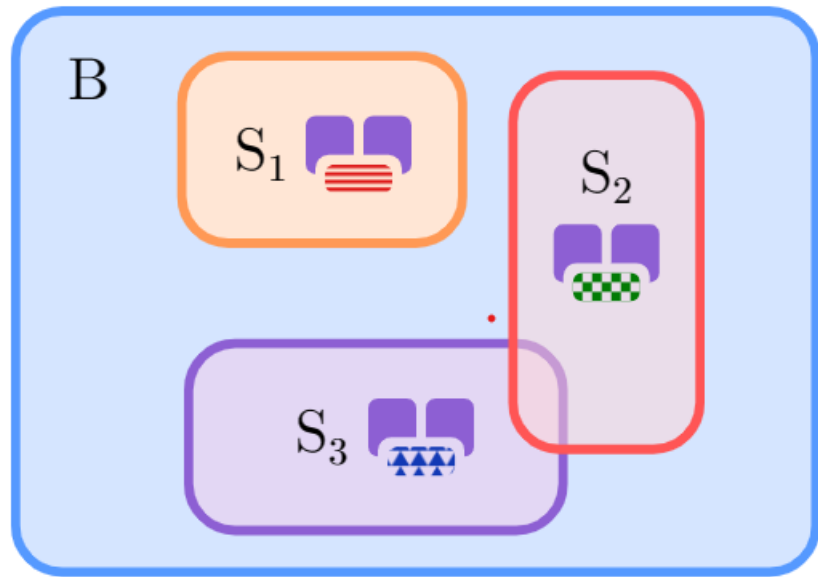
ATM PRE 2024

Henderson, Nagano, Milighetti, ATM PNAS 2024

Nagano, Pyo, Milighetti, Henderson, Shawe-Taylor, Chain*, ATM* arXiv 2024

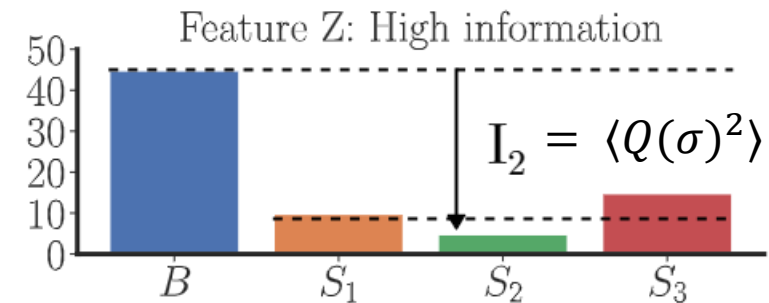
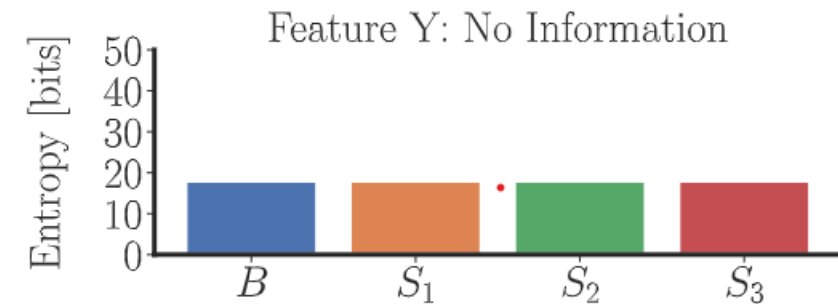
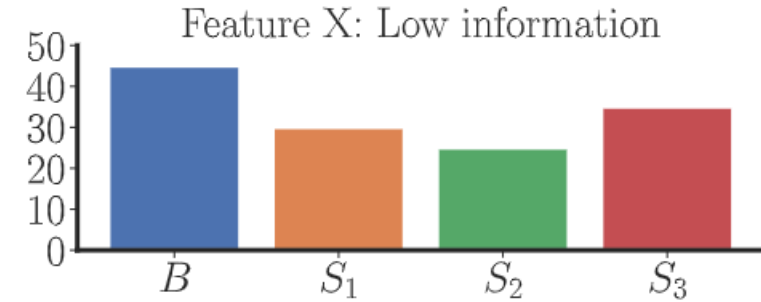
Pyo, Nagano, Milighetti, Henderson, Callan, Chain, Wingreen, ATM (in preparation)

The variance of selection factors is a measure of information

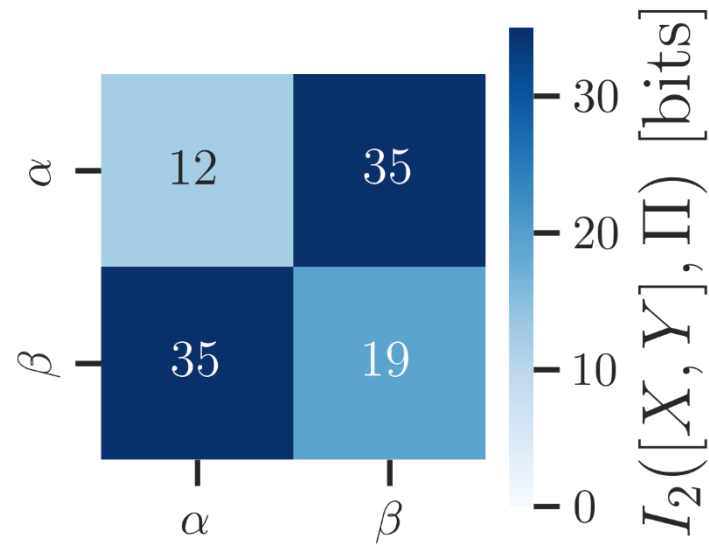


sequencing

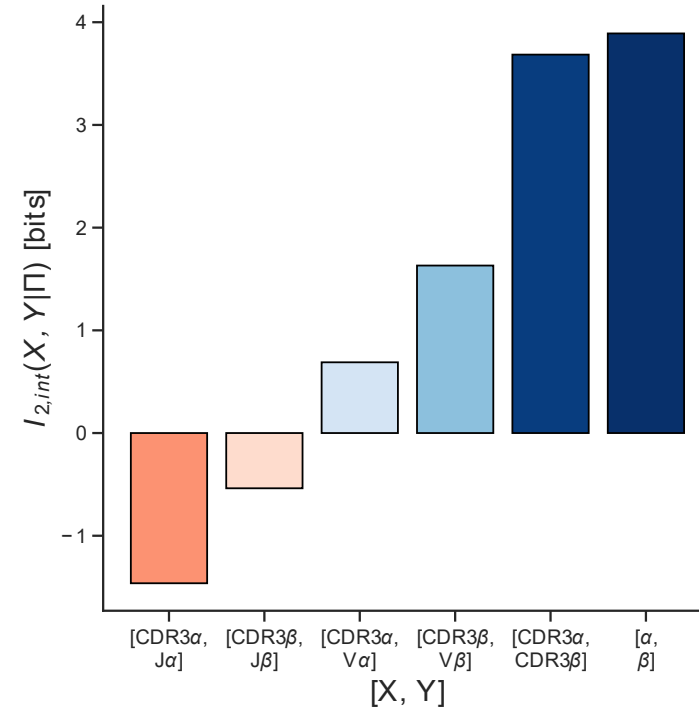
$$H_2[P(X)] = -\log p_C[P(X)]$$



Mapping TCR specificity information

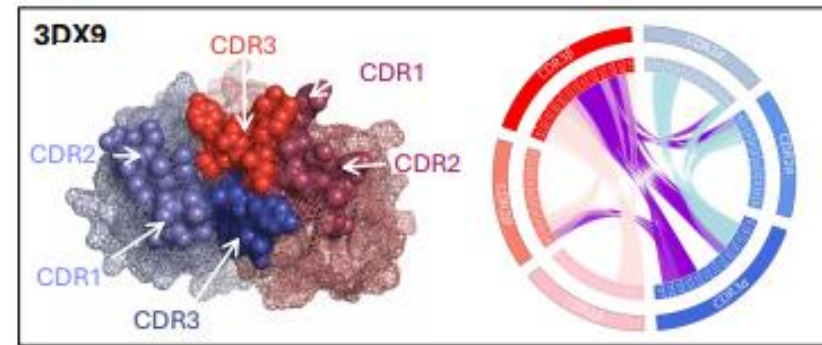
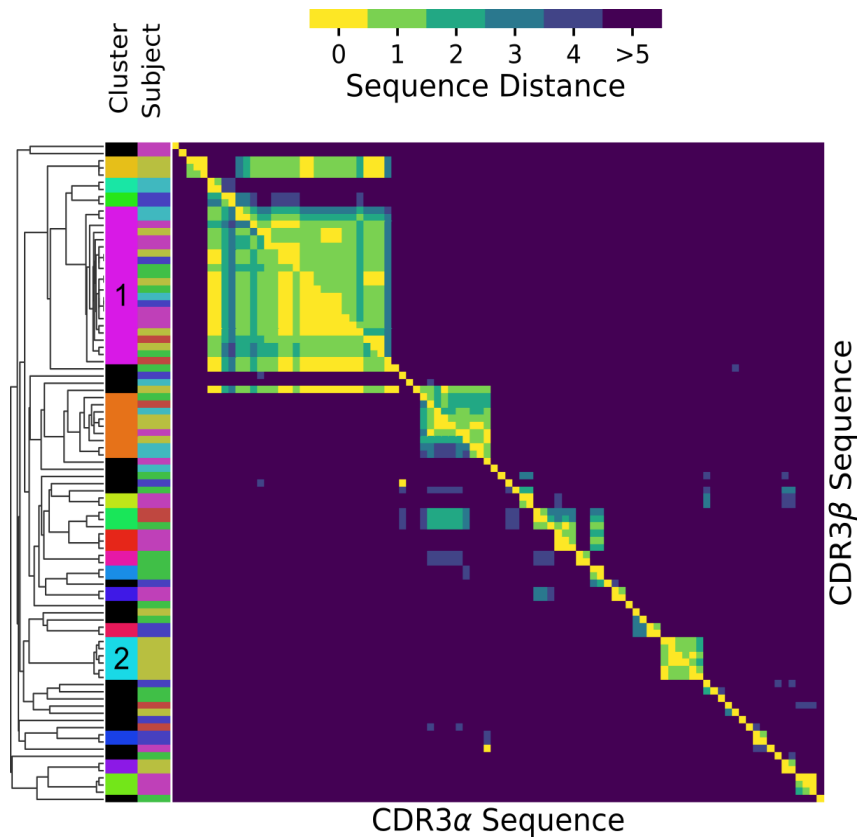


$\beta > \alpha$
but: both chains matter

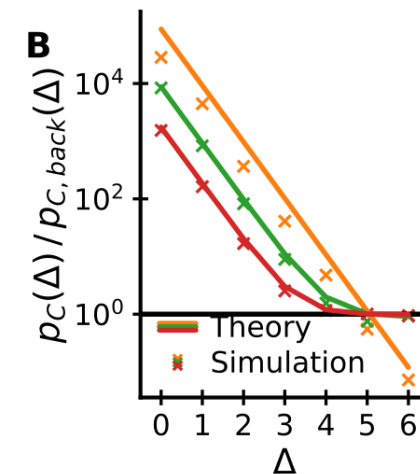
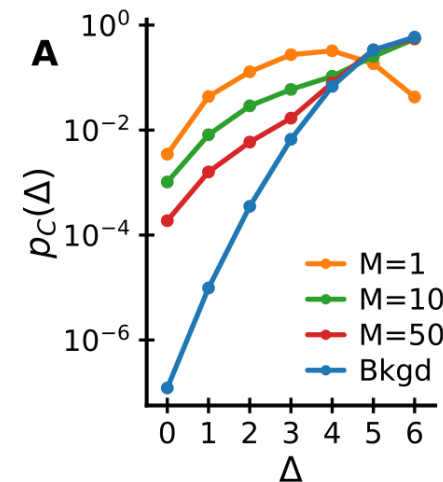


Chains are not independent

α / β chains have constrained pairing



Milighetti ... ATM ...
Chain bioRxiv 2024

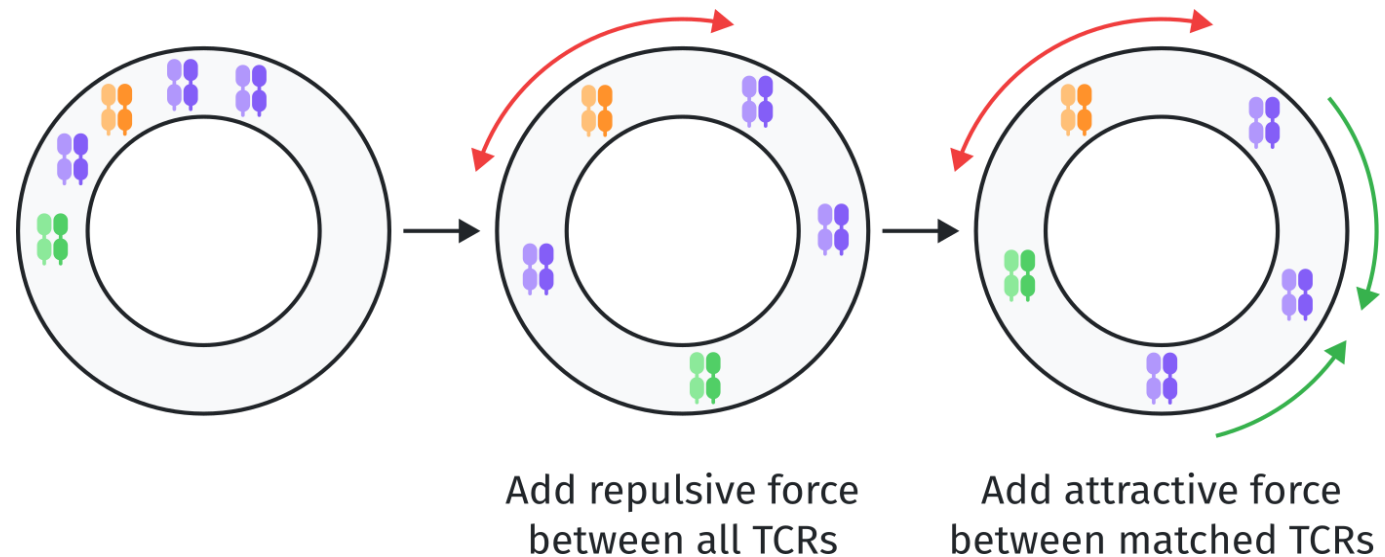


Discover biophysical determinants of co-specificity by contrastive learning

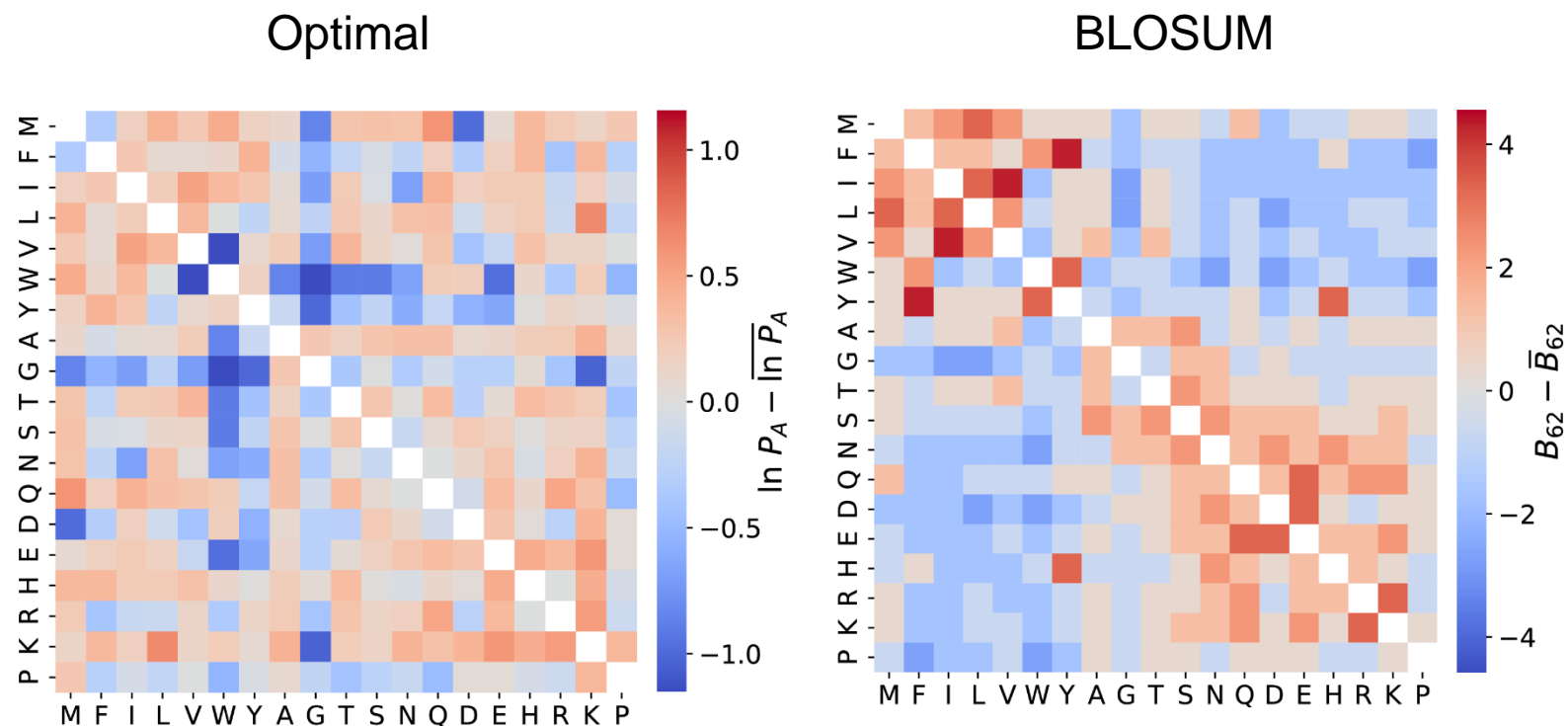
$$Q(\sigma, \sigma') \propto e^{-d(\sigma, \sigma')}$$

Optimal d ? Pseudo-likelihood maximization:

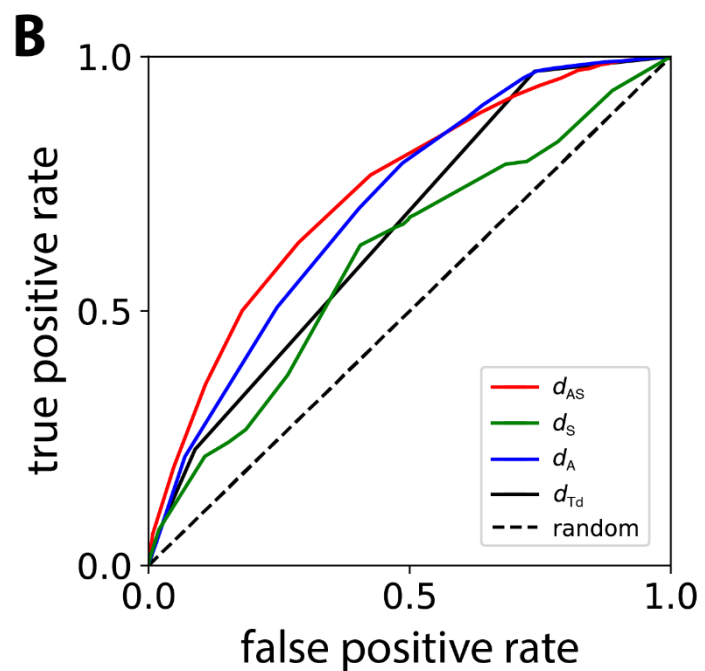
$$\varepsilon(d|\mathcal{P}, \mathcal{U}) = \langle d \rangle_{\mathcal{P}} + \log \langle e^{-d} \rangle_{\mathcal{U}}$$



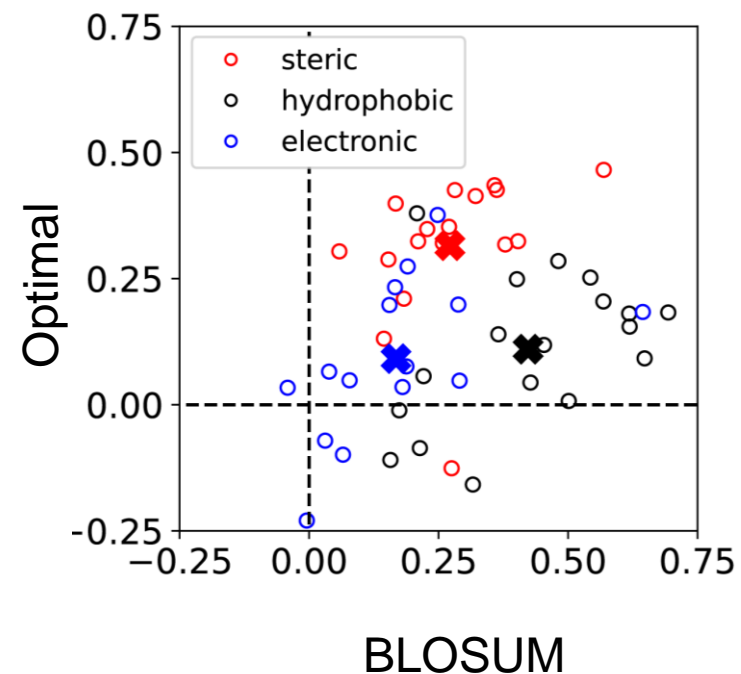
Optimal alignment weights differ from BLOSUM



Learned metric generalizes and is interpretable



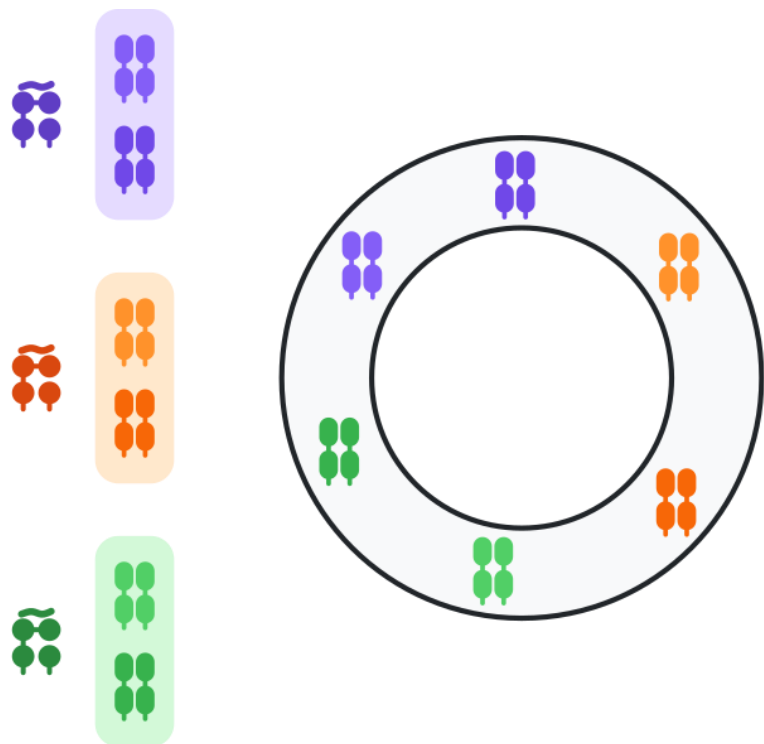
generalization to pMHCs
not seen during training



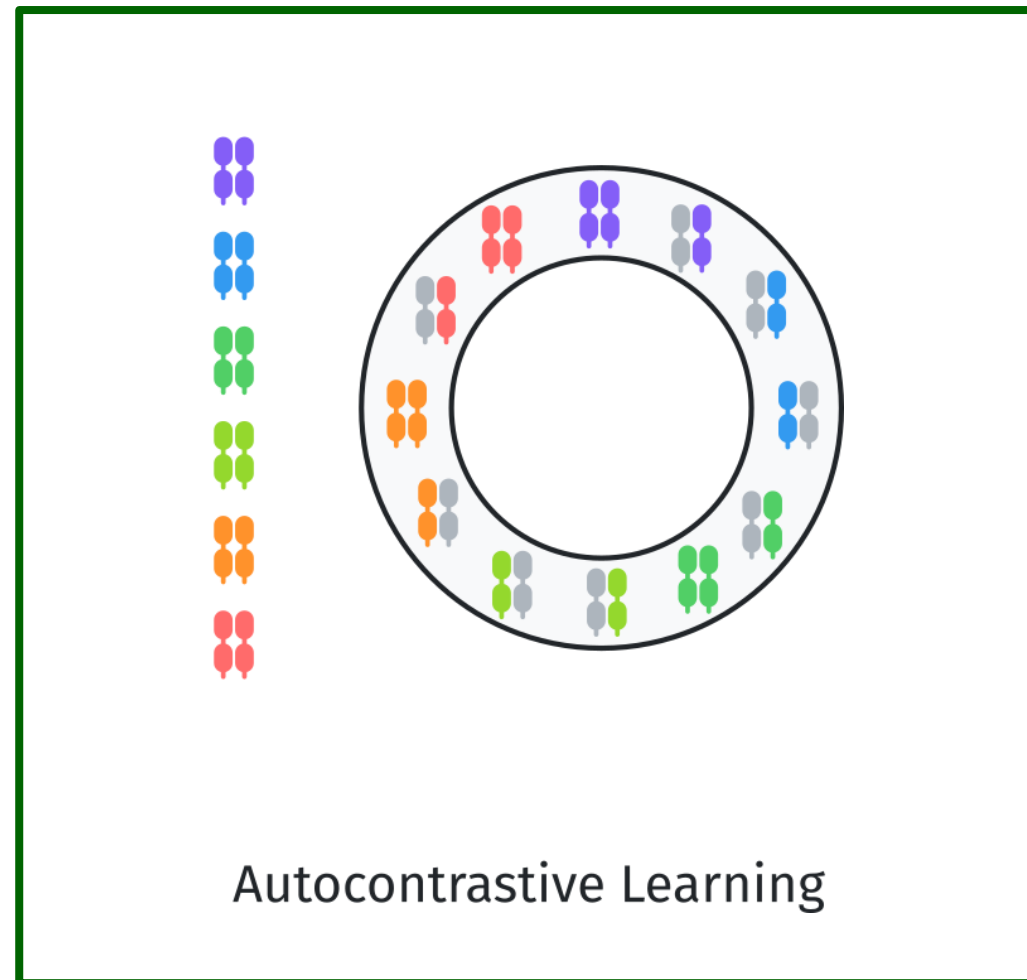
Correlation coefficients of substitution matrix with amino acid characteristics

shape matters

Contrastive training variants



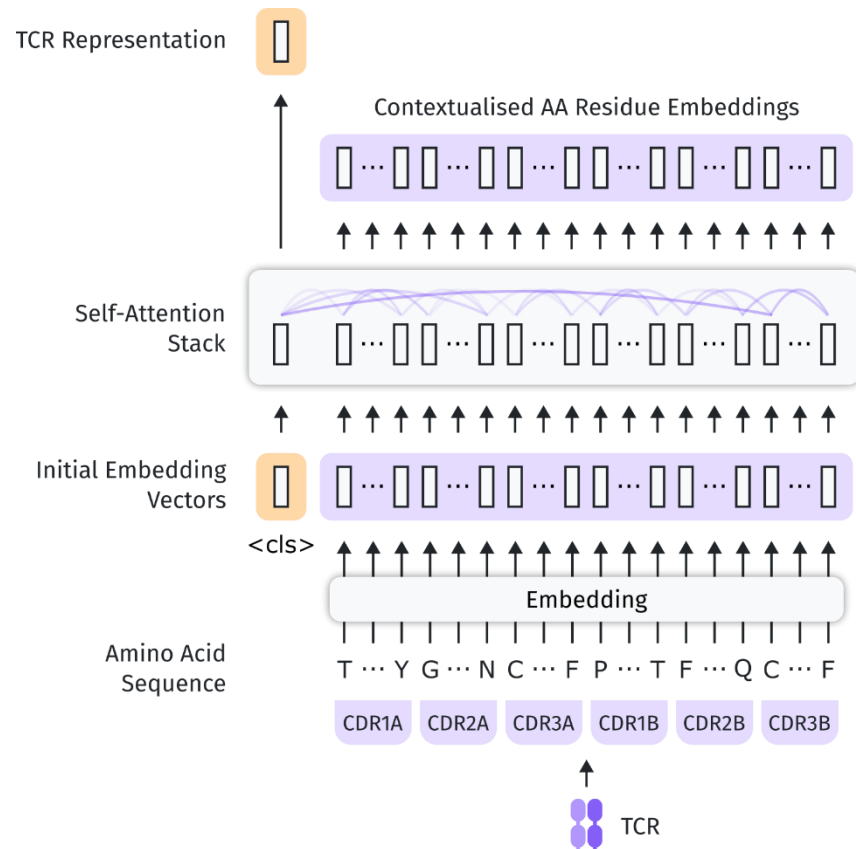
Supervised Contrastive Learning



Autocontrastive Learning

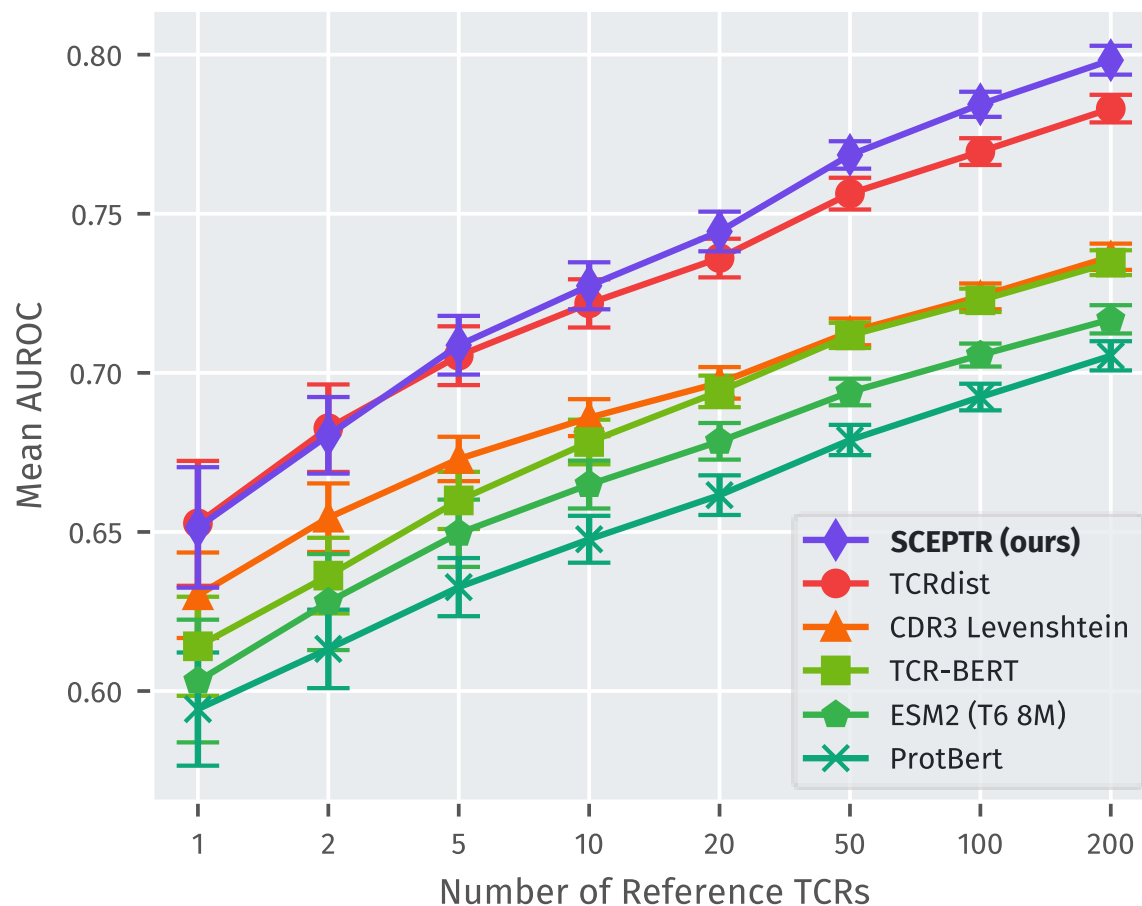
Architecture of SCEPTR

Simple Contrastive Embedding of the Peptide sequence of the T cell Receptor

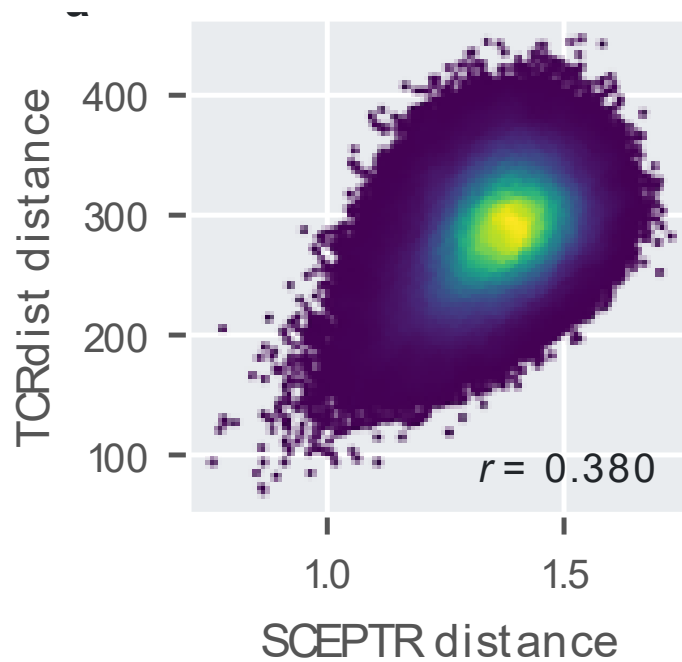


- **Input:** AA sequence of CDRs 1/2/3 of the α and β TCR
- **Output:** Vector representation of the TCR

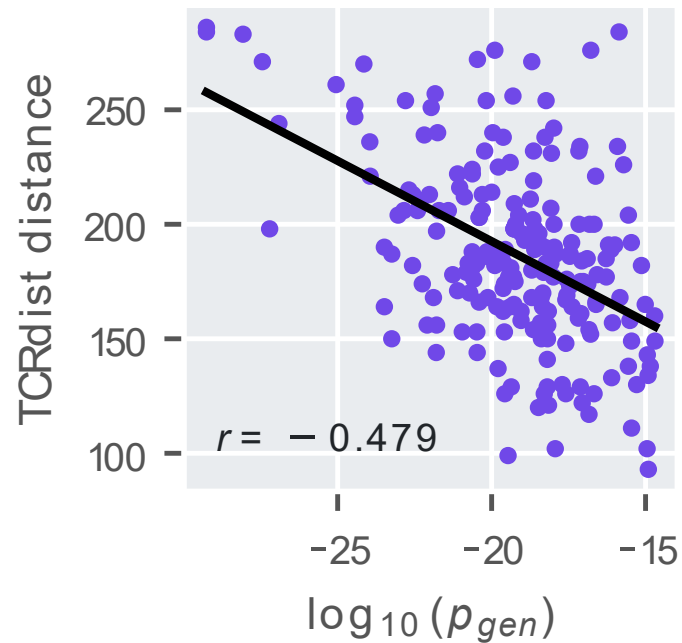
SCEPTR matches alignment accuracy using auto-contrastive pre-training



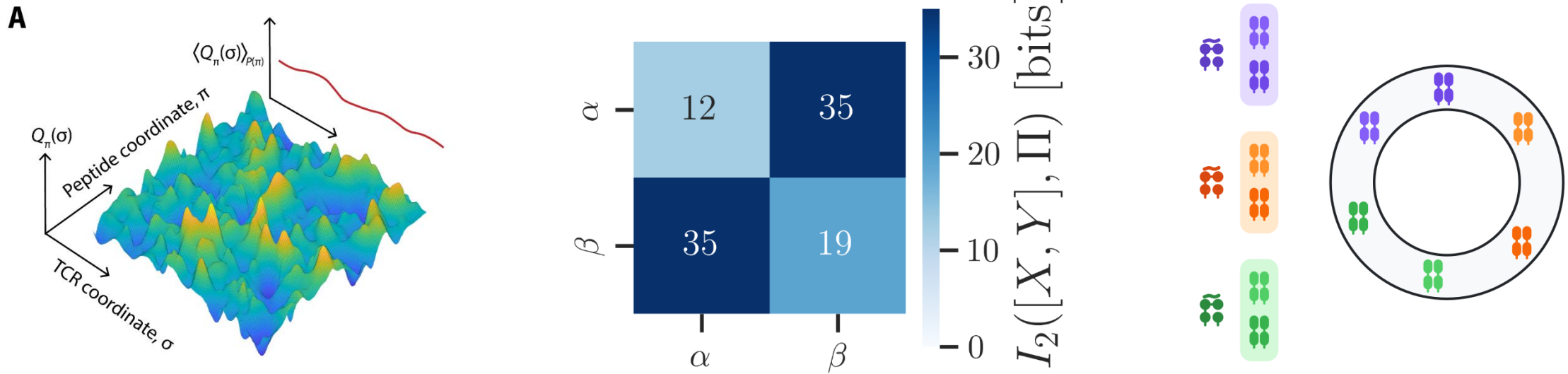
Comparing SCEPTR to TCRdist



Approximates sequence alignment distances

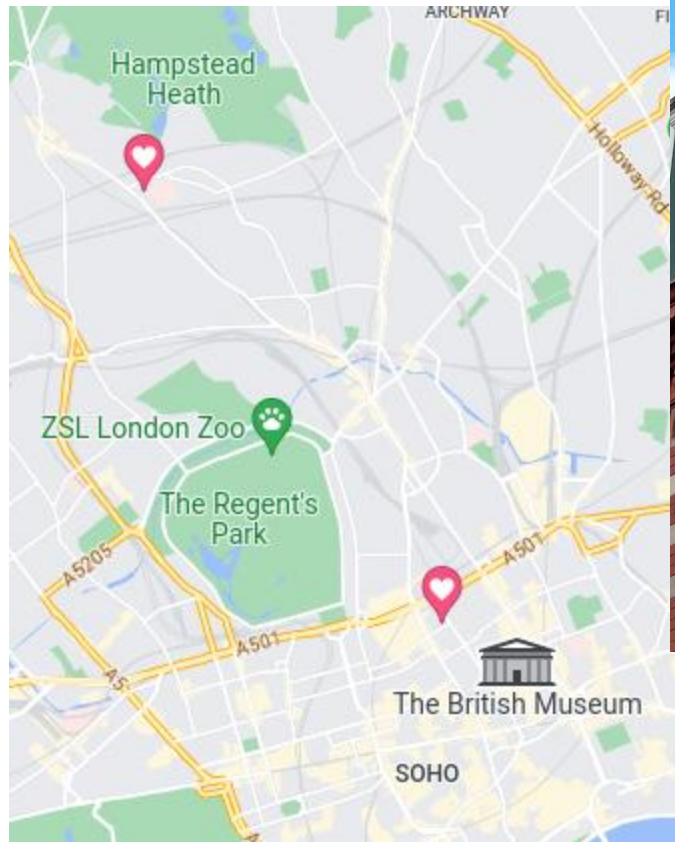


Corrects for recombination biases



Which question would you ask if we could predict antigen-specificity from antibody/TCR sequence?

Come visit!



@andimscience

Interested in joining the group? Reach out!

www.qimmuno.com