

Learning the language of T cell receptor repertoires

Quantitative Immunology

February 2016

Benny Chain

[Innate2Adaptive](#)

[Immunology with numbers](#)

UCL

- **John Shawe-Taylor**
- Yuxin Sun

- Niclas Thomas
- James Heather
- Theres Matjeka
- Mattia Cinelli
- Katharine Best
- Mazlina Ismail
- Connor Husovsky
- Kroopa Joshi

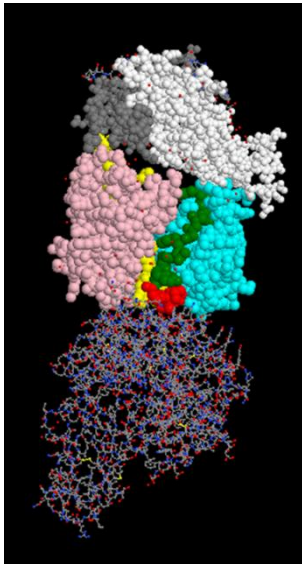
Weizmann Institute

- **Nir Friedman**

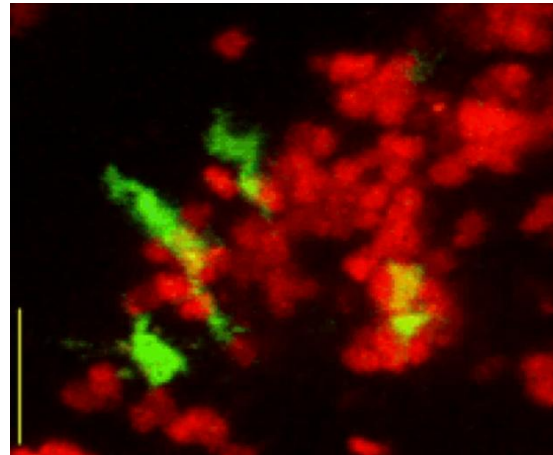
- Hila Gal
- Wilfred Ndifon
- **Shlomit Reich-Zeliger**
- Eric Shifrut
- Michal Marks

Studying the repertoire is an approach to understanding how T cells recognise antigen ...

At an individual receptor level ?



At a population level ?



Thymus

Periphery

Positive selection

Self-antigen driven proliferation

Cytokine driven proliferation

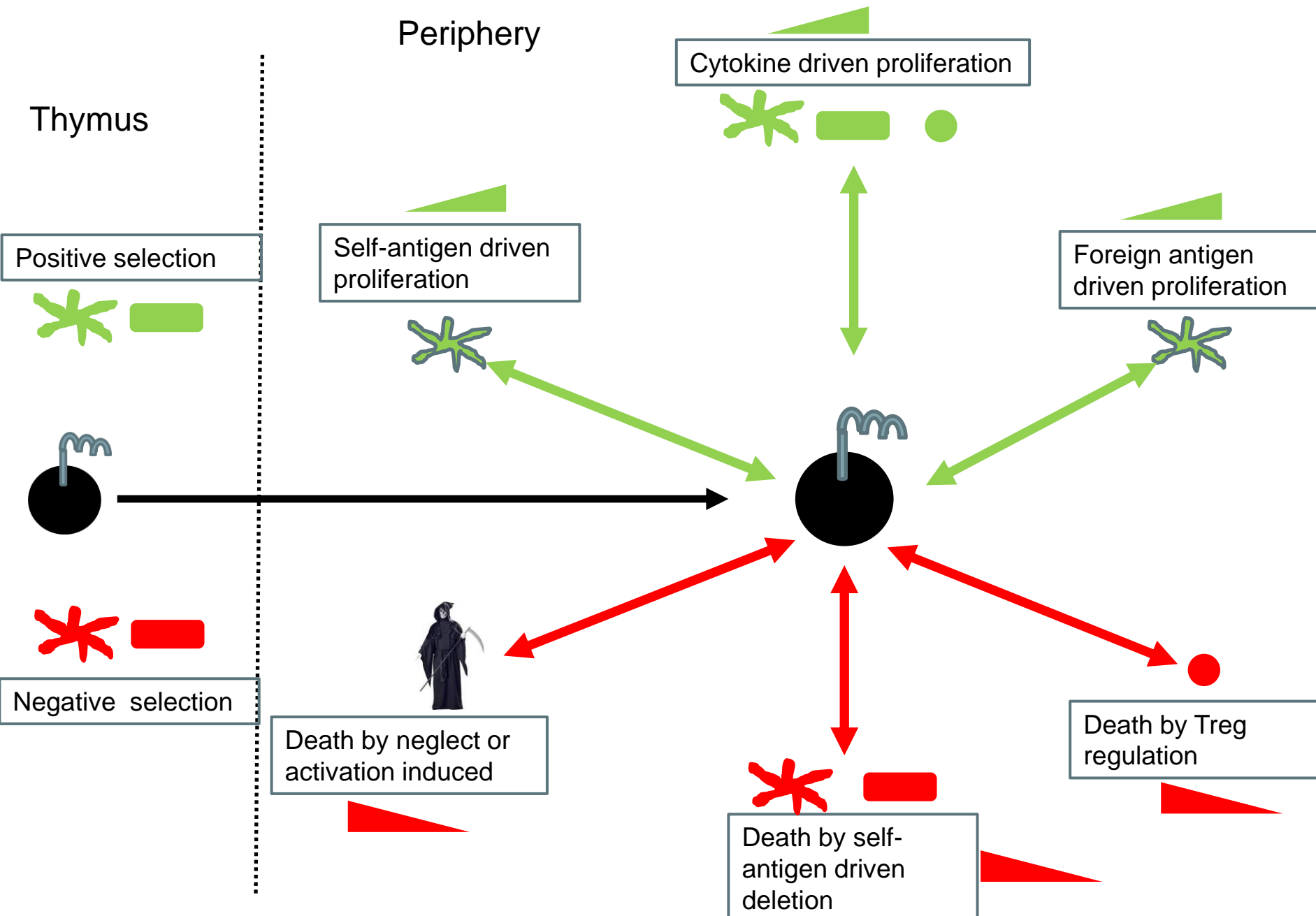
Foreign antigen driven proliferation

Negative selection

Death by neglect or activation induced

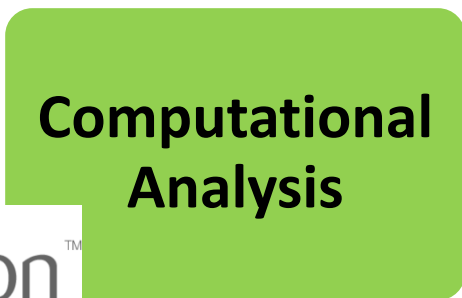
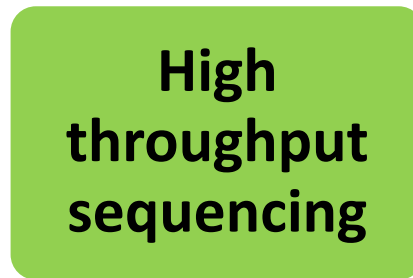
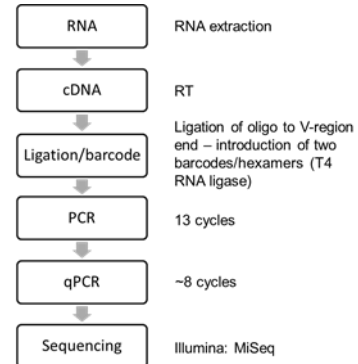
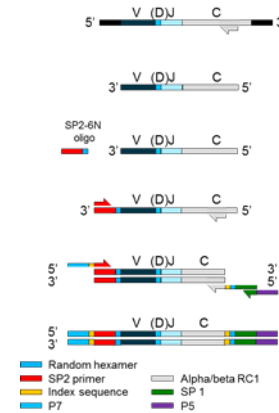
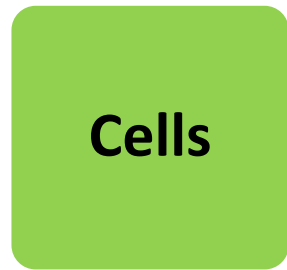
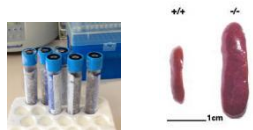
Death by self-antigen driven deletion

Death by Treg regulation

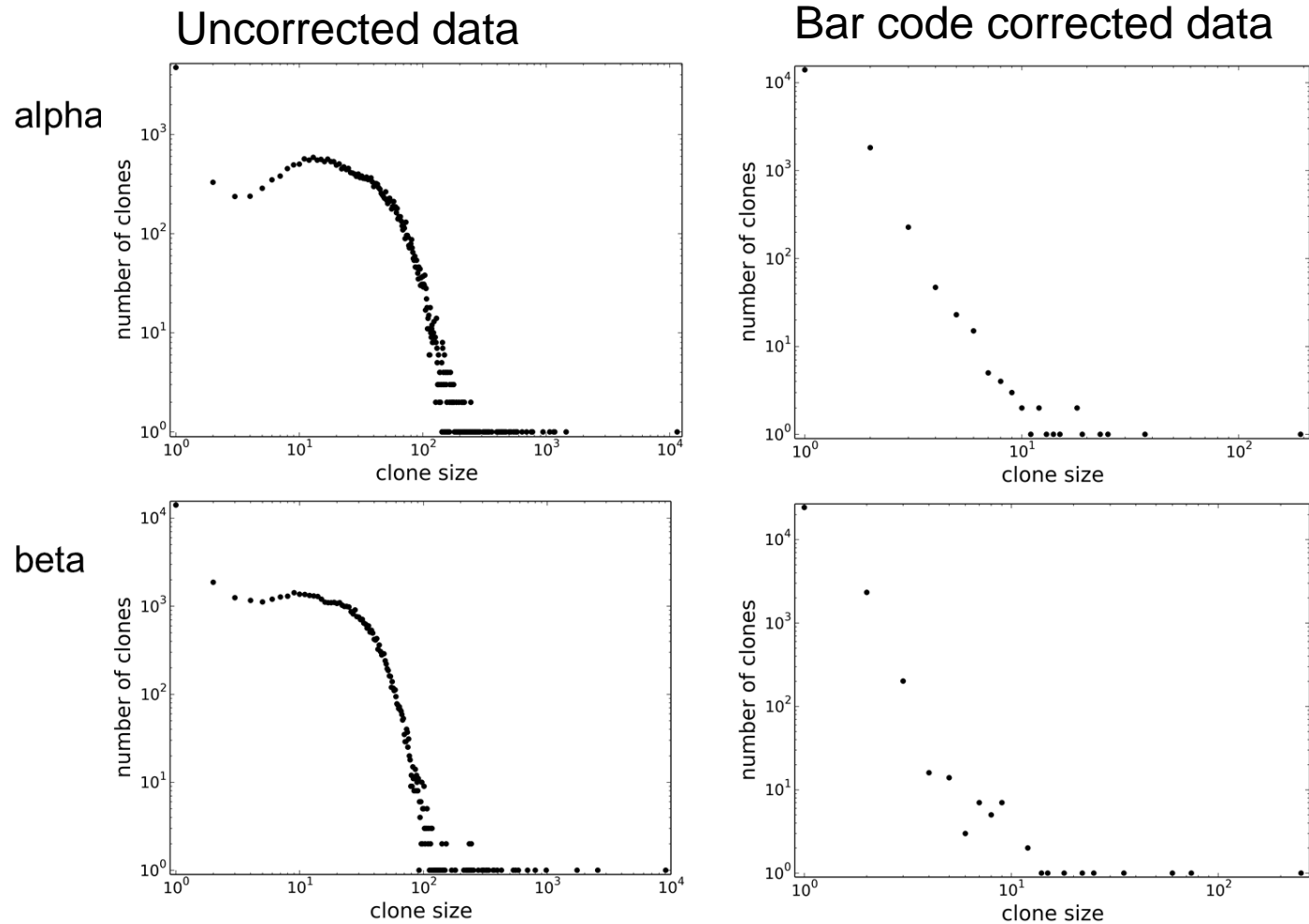


... and may provide biomarkers, tools for diagnosis and stratification , and correlates of protection or pathology for infectious diseases.

The TCR sequencing pipeline



Cord blood sample A



1: Best K, Oakes T, Heather JM, Shawe-Taylor J, Chain B. Computational analysis of stochastic heterogeneity in PCR amplification efficiency revealed by single molecule barcoding. Sci Rep. 2015 Oct 13;5:14629.

What is the T cell repertoire response to antigen ?

Can we recognise antigen specific cells from sequences ?

- There are $>10^{14}$ different alpha and beta chains.
- Each individual has undergone a unique set of stochastic recombination events
- Don't know rules to go from sequence to structure to specificity
- HTS data does not link α and β TCR chains per cell

The experiment

β chain CDR3 TCR repertoire (NOT barcoded) sequenced from CD4 spleen cells.

9 unimmunised mice

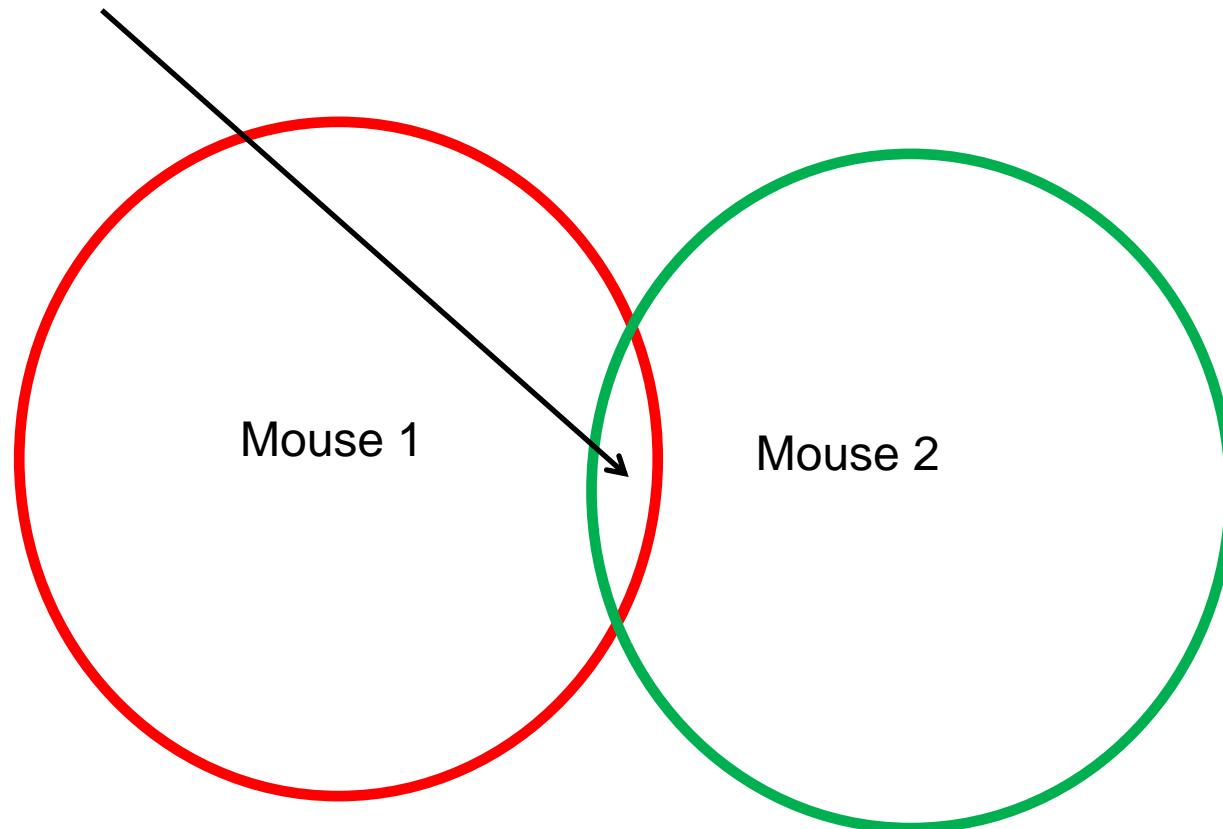
8/9 “early” immunised mice (5,7 and 14 days) immunised with Complete Freund’s Adjuvant (CFA) or ovalbumin with CFA

5/6 “late” mice (60 days) with additional boost at 14 days

Isolate CD4 cells, amplify (NO barcoding, multiplex) and sequence on HiSeq

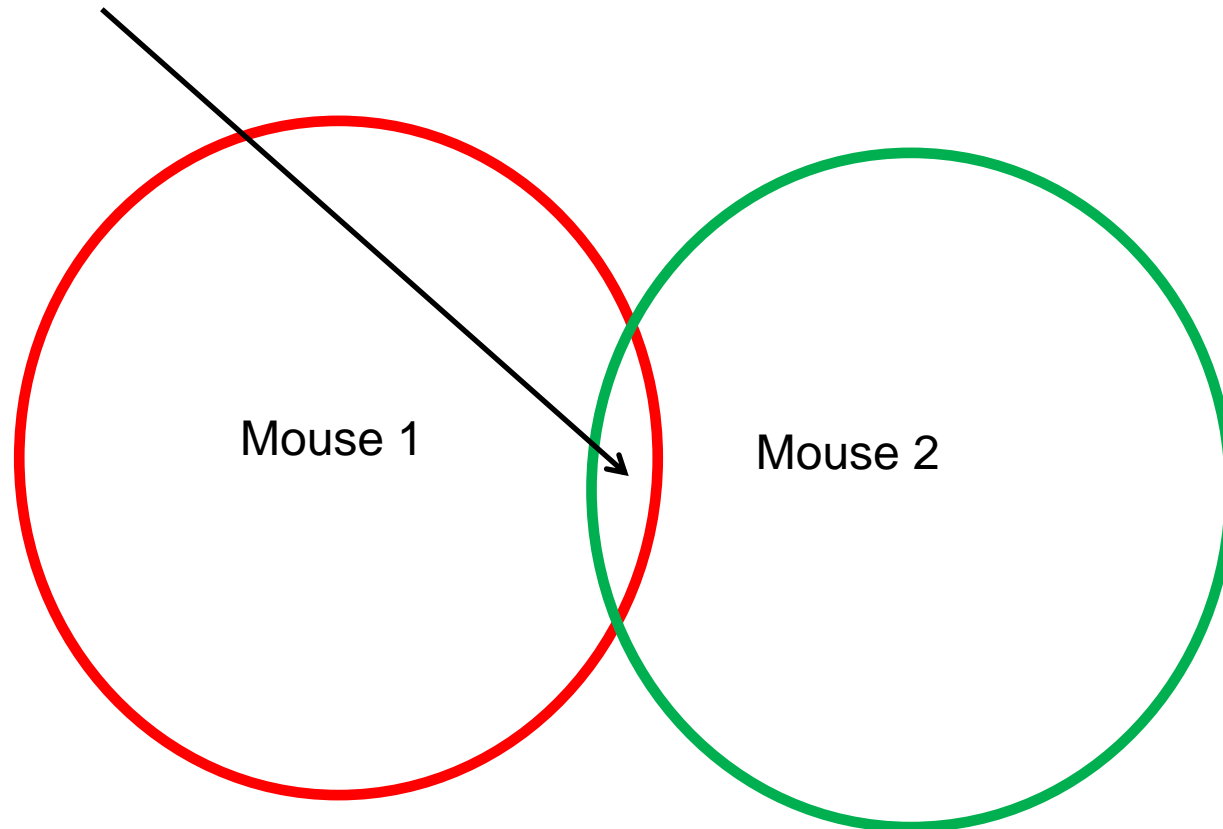
Simplistic approach

The Jaccard index measures overlap

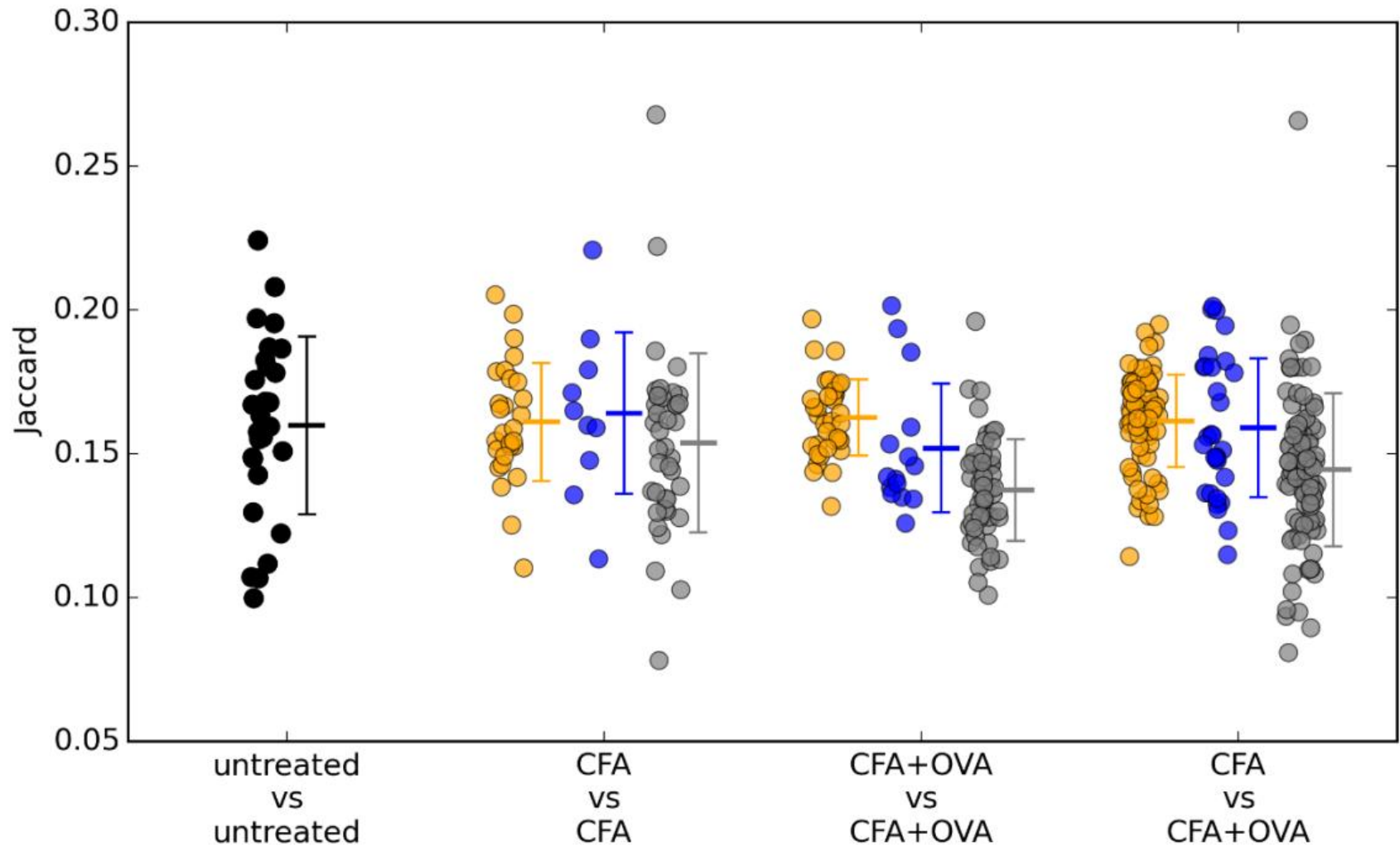


$$J = (M1 \cap M2) / (M1 \cup M2)$$

The Jaccard index measures overlap



Does immunisation drive increased overlap ?

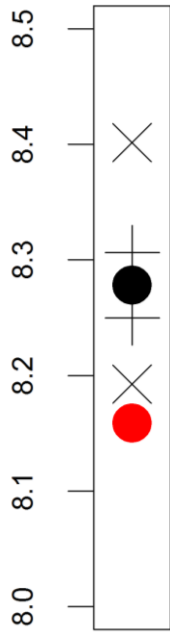


DIVERSITY AND PRIVACY

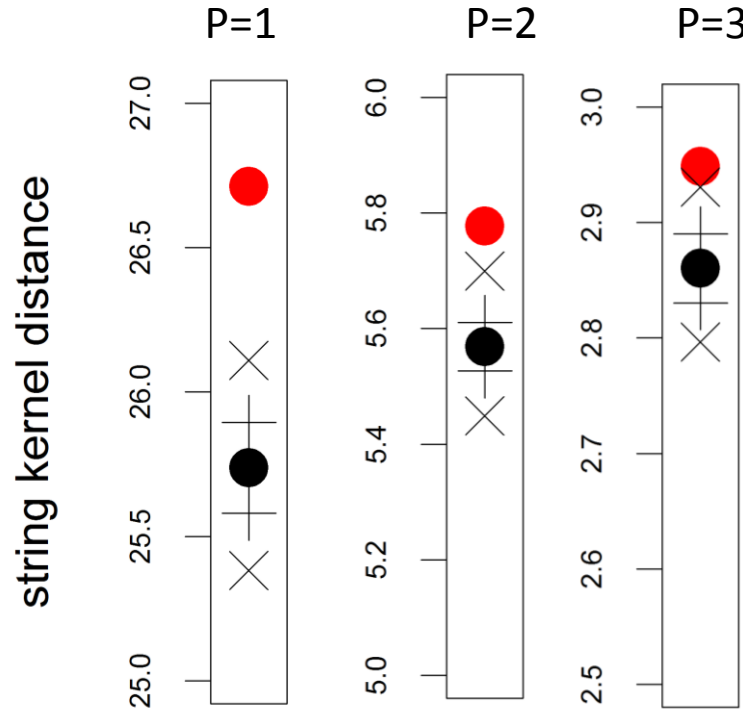
- Focus on early time point
- Measure abundance of each CDR3 in each immunised mice relative to average frequency in unimmunised mice and express as $\text{Log}_2(\text{relative frequency})$ (LRF)
- Identify >2000 sequences with LRF > 6 in OVA immunised mice, but < 2 in CFA only mice
- Plot LRF for each of this set of CDR3s in all mice

SIMILARITY

Are there features of the repertoire which distinguish the OVA responsive repertoire ?



f



What should we use as features for TCR analysis ?

- Short amino acid motifs as features

CAISLAKDRAYNEQFF CAALRATGGNNKLTFGQG CASSLRATGNLTFGQG ...

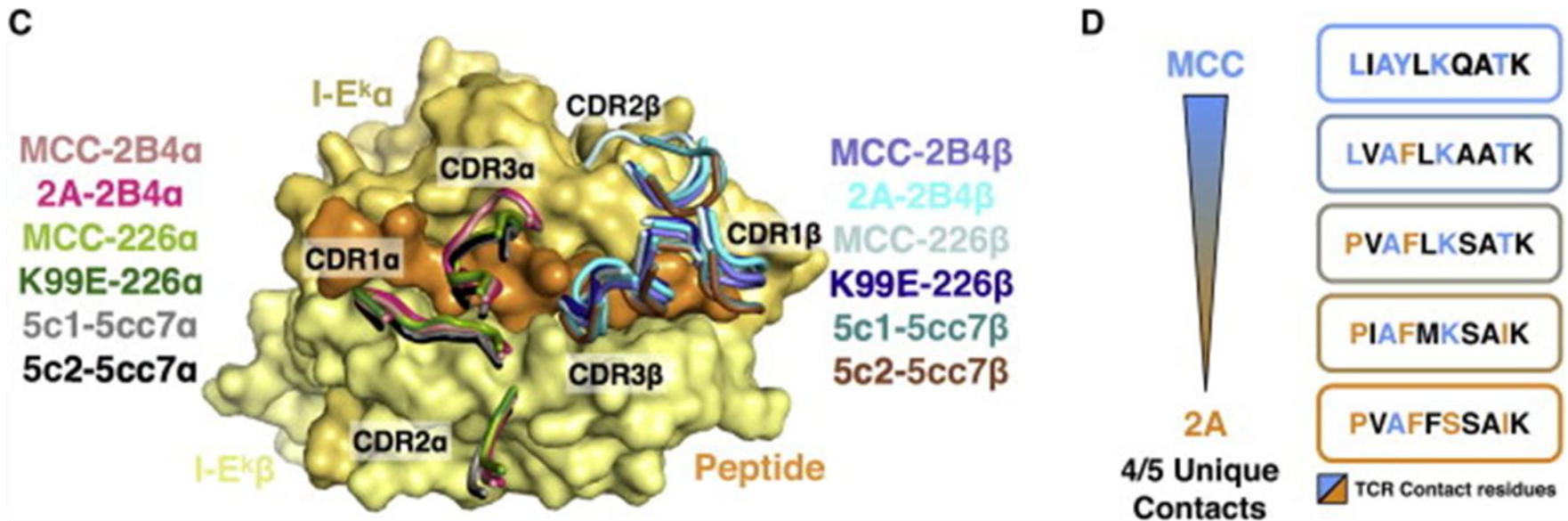


CAI AIS ISL SLA LAK AKD

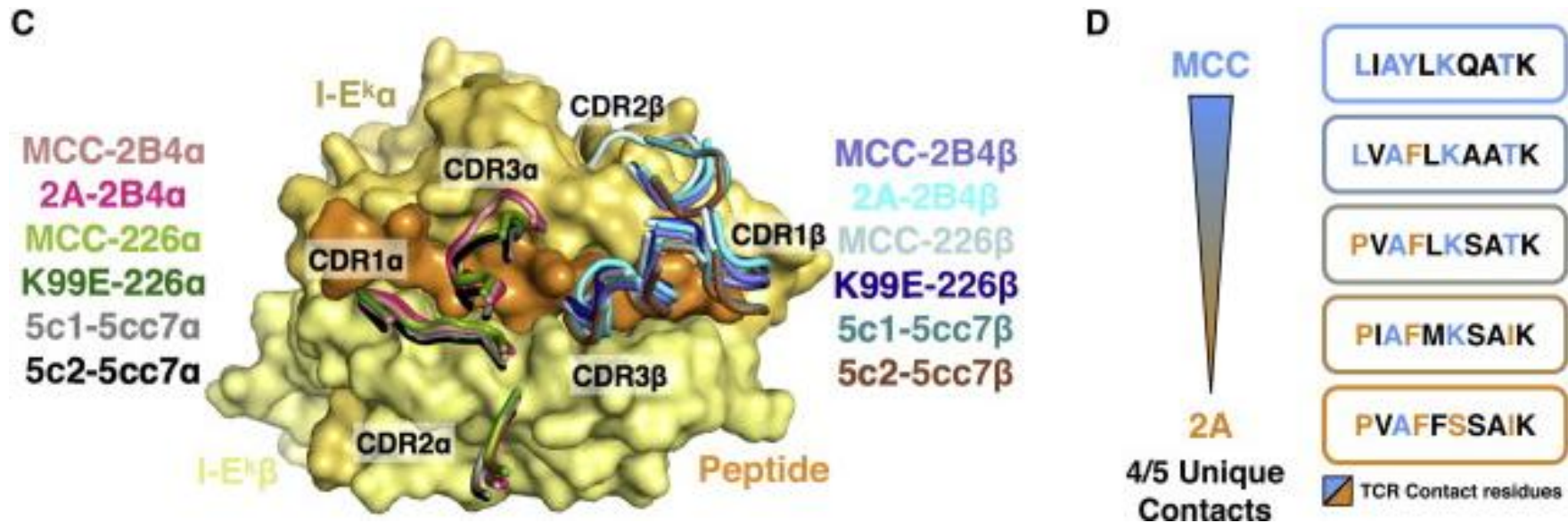




Some *post facto* motivation for this approach



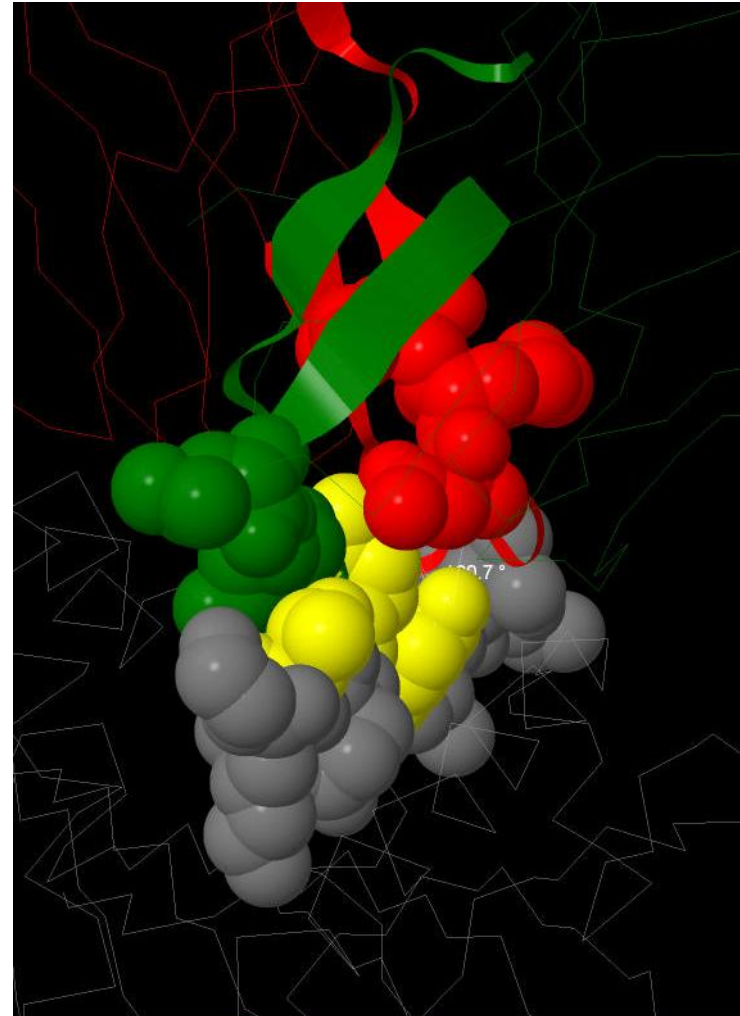
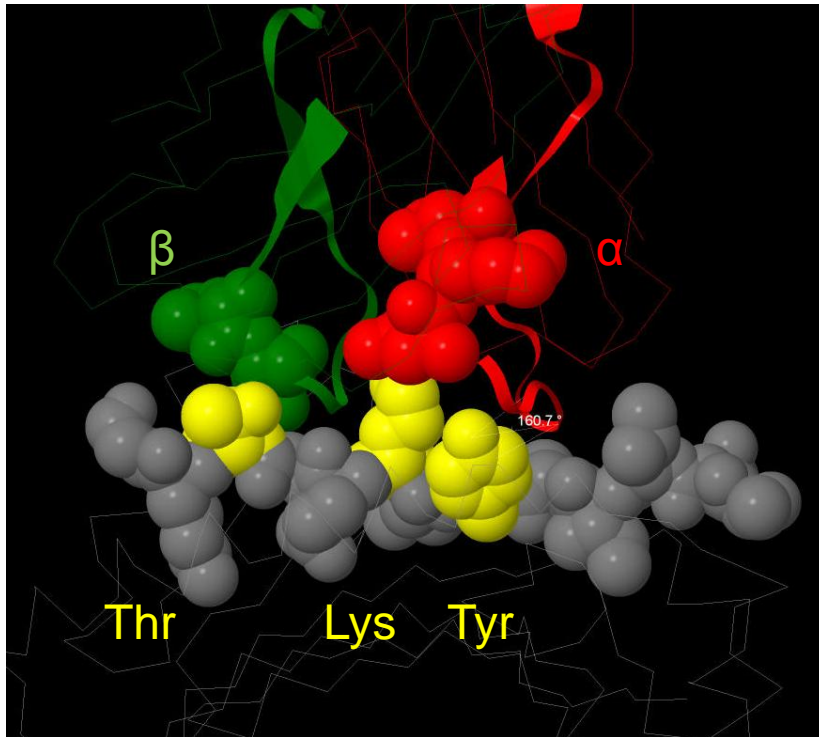
Deconstructing the peptide-MHC specificity of T cell recognition. *Birnbaum ME, Mendoza JL, Sethi DK, Dong S, Glanville J, Dobbins J, Ozkan E, Davis MM, Wucherpfennig KW, Garcia KC. Cell. 2014 May 22;157(5):1073-87*



Cyt C peptide IAYLKQATKG

2B4 α CAALRATGGN**NKL**TFGQG
226 α CAAEPSSG--**QKL**VFGQG
5CC7 α CAAE-AS-NT**NKV**VFGTG

2B4 β CAS**SLN**WSQDTQYFGPG
226 β CAS**SLN**NANSDYTFGSG
5CC7 B CAS**SLN**WSQDTQYFGPG



2B4 TCR binding cytochrome C peptide showing conserved residues

Selecting features

Beyond simple counting : a Markov process probability model for sequence generation

Consider a (protein) sequence (pseq) of length l as made up of a series of (contiguous) aa strings of length k .

Define $P_{u \rightarrow x}$ as the transition probability of moving from string u (length $k-1$) to character x .

We might consider the protein as a Markov process where the likelihood of generating the protein is given by :

$$\mathcal{L}(\text{pseq}) = \prod_{j=k}^l P(\text{seq}_{j-k+1:j-1} \rightarrow \text{seq}_j) \text{ given a set of parameters } \Theta = P(u \rightarrow x) \text{ for all } u \text{ and } x.$$

Given this generative model the Fisher score is defined as :

$$\begin{aligned} \text{Fisher score} &= \frac{\partial \log (\mathcal{L}(\text{seq}))}{\partial \Theta} \\ &= \frac{\partial \log \prod_{j=k}^l P(\text{seq}_{j-k+1:j-1} \rightarrow \text{seq}_j)}{\partial P(u \rightarrow x)} \quad \forall ux \in \Theta \end{aligned}$$

This is a vector indexed by all possible ux (k -mers) whose elements can be approximated by

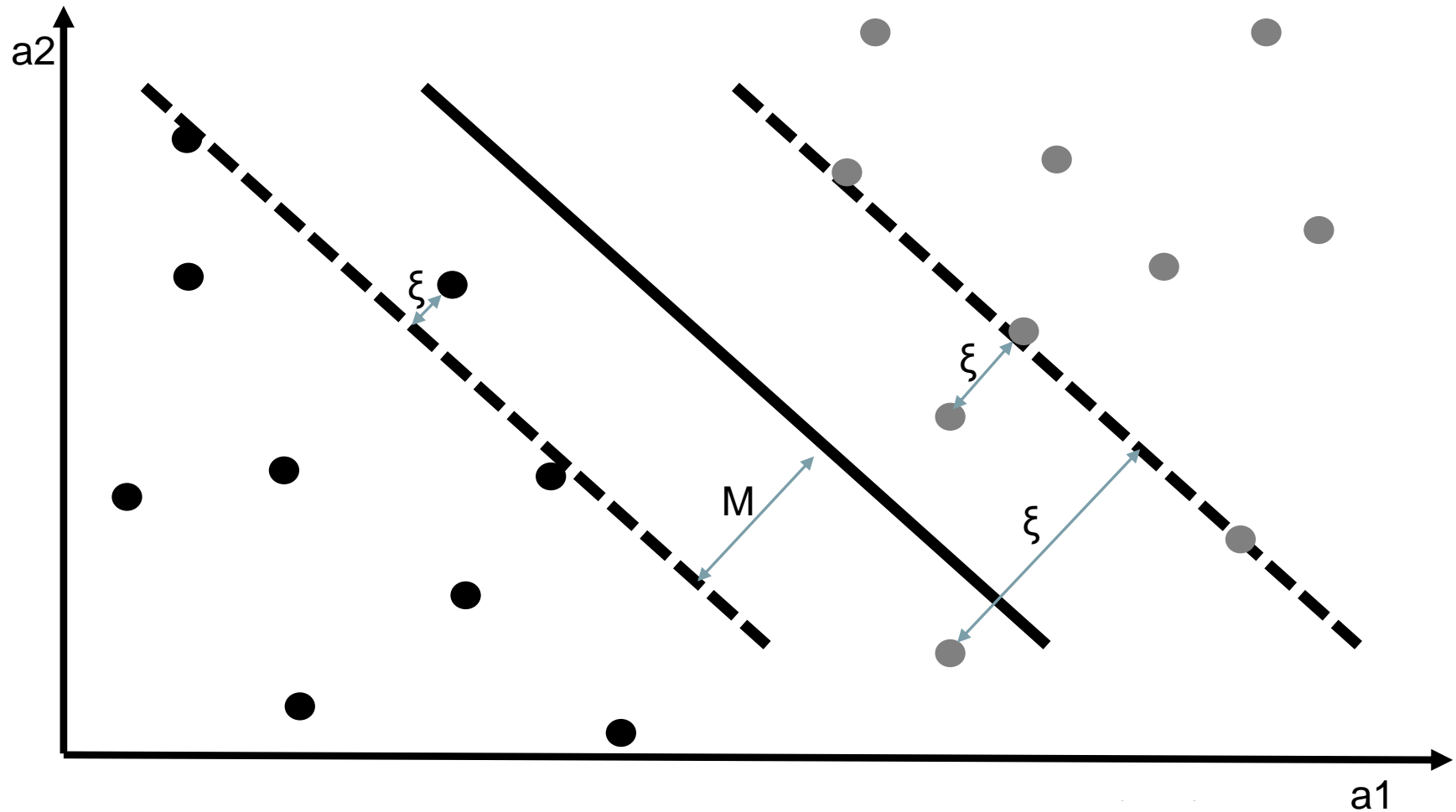
$$\frac{\text{Term frequency } (v)}{P(u \rightarrow x)}$$

Saunders, C., Vinokourov, A. & Shawe-taylor, J. S. String kernels, fisher kernels and finite state automata. in *Adv. Neural Inf. Process. Syst.* 633–640 (2002).

Practically, this corresponds to counting k-mers (e.g. triplets) and weighting according to some observed transition probability (or function of this probability)

Conceptually, the Fisher score can be considered as focusing attention on those features with the largest gradient in parameter space : i.e. those where a small change in parameter will have a large effect on the likelihood. This approach has proved fruitful in the context of image recognition.

Combining feature selection and classification



Looking for linear classifiers of the form : $f(\mathbf{x}) = \sum_1^m a_j H_j$

Linear programming boosting

Primal

$$\min_{a, \xi} \sum_{j=1}^m a_j + C \sum_{i=1}^l \xi_i$$

Such that

$$\sum_j^m a_j y_i H_{i,j} \geq 1 - \xi \quad \forall i = 1 \dots l$$

$$a_j \geq 0 \quad \forall j = 1 \dots m$$

$$\xi_i \geq 0 \quad \forall i = 1 \dots l$$

Dual

$$\max_u \sum_{i=1}^l u_i$$

Such that

$$\sum_{i=1}^l u_i y_i H_{ij} \leq 1 \quad \forall j = 1 \dots m$$

$$0 \leq u_i \leq C \quad \forall i = 1 \dots l$$

$$\max_u \sum_{i=1}^l u_i - \frac{1}{2} \sum_{i=1}^l \sum_{k=1}^l u_i u_k y_i y_k \langle H_i \cdot H_k \rangle$$

$$0 \leq u_i \leq C, \forall i = 1 \dots l$$

$$\sum_{i=1}^l u_i y_i = 0$$

Definitions

- l : number of training data points
- m : number dimensions (or weak learners or features)
- H : the set of features associated with each data point
- ξ : slack variables for data points within margin
- a : coefficients defining hyperplane (i.e. feature weights)
- m : sample weights
- C : penalty weights
- y : classification

Support vector machine

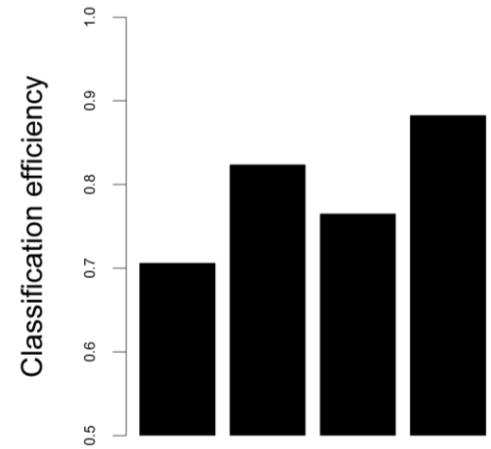
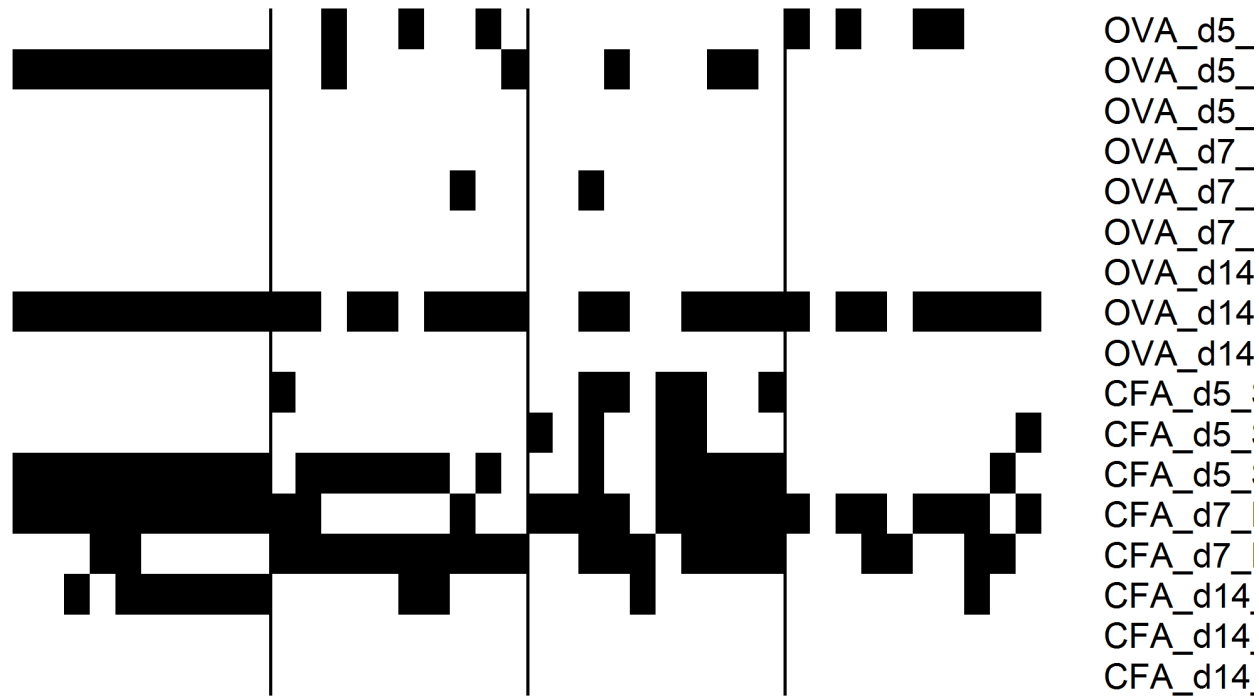
$$\min_{a, \xi} \frac{1}{2} \sum_{i=1}^m a_i^2 + C \sum_{i=1}^l \xi_i$$

Linear Boost
Triplet counts

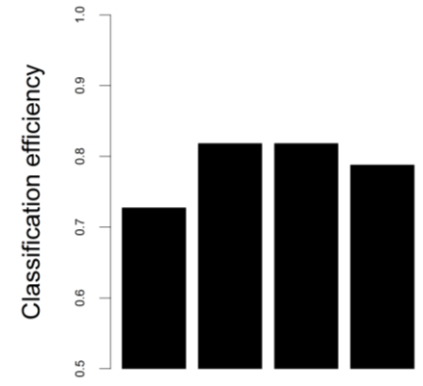
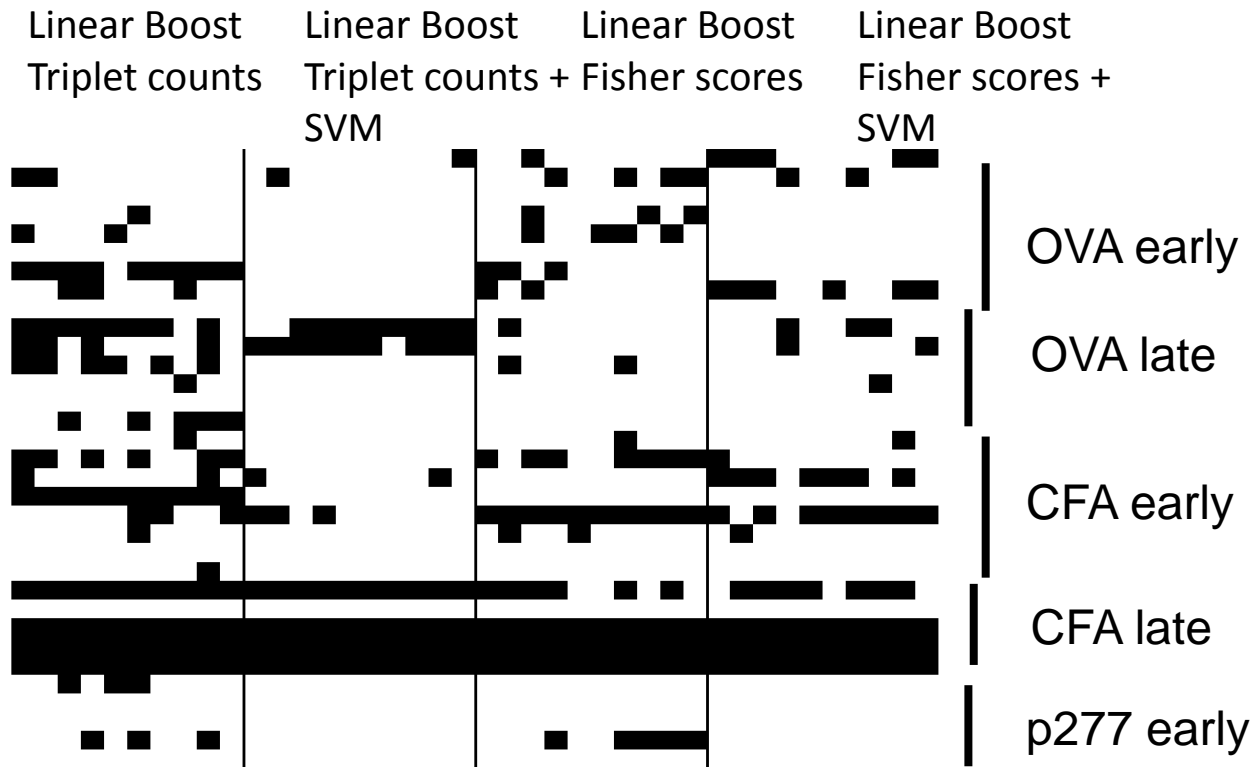
Linear Boost
Triplet counts + Fisher scores
SVM

Linear Boost
Fisher scores

Linear Boost
Fisher scores + SVM

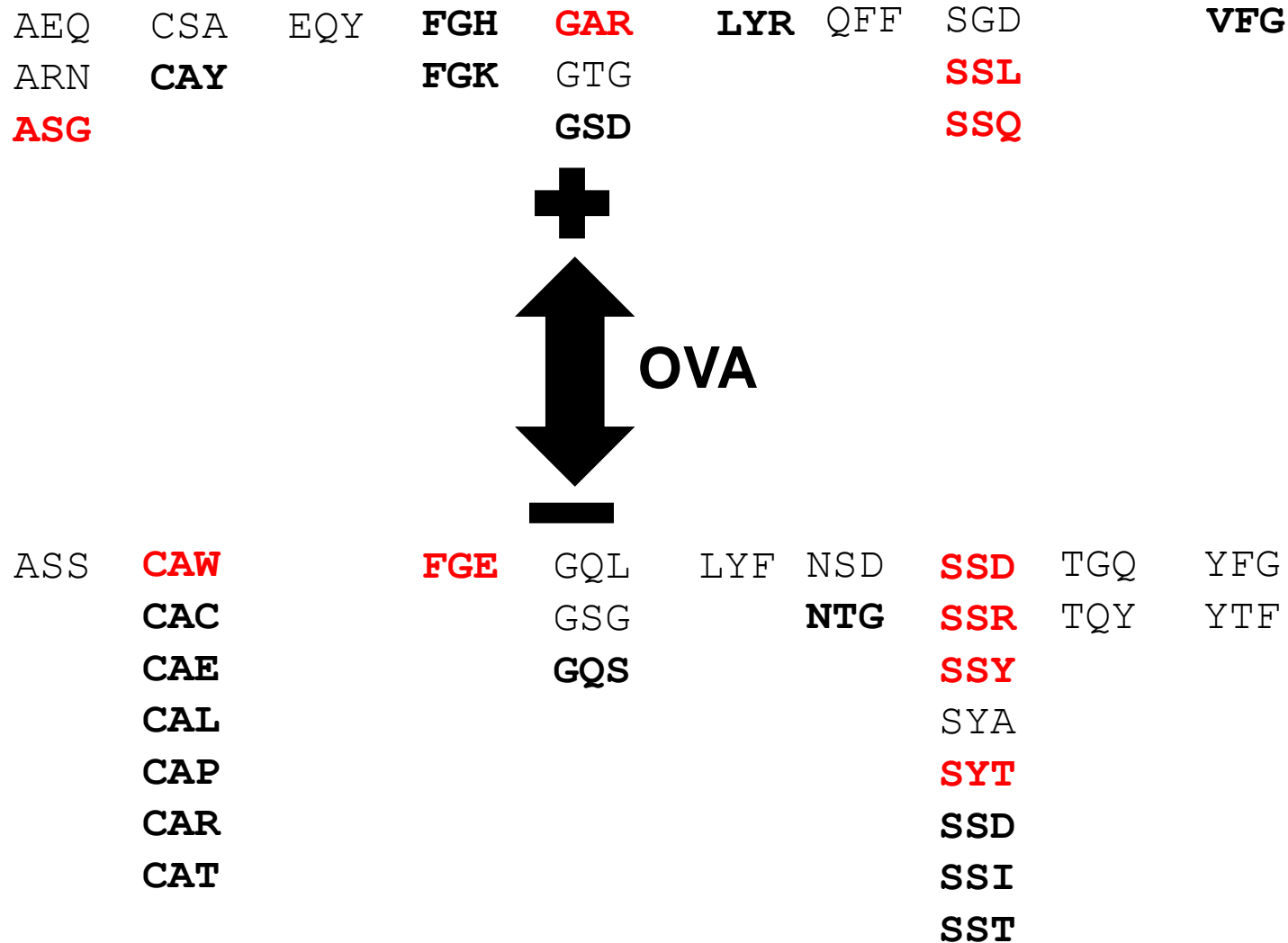


Linear boosting using 50 thousand CDR3s from each mouse, leave-one-out validation

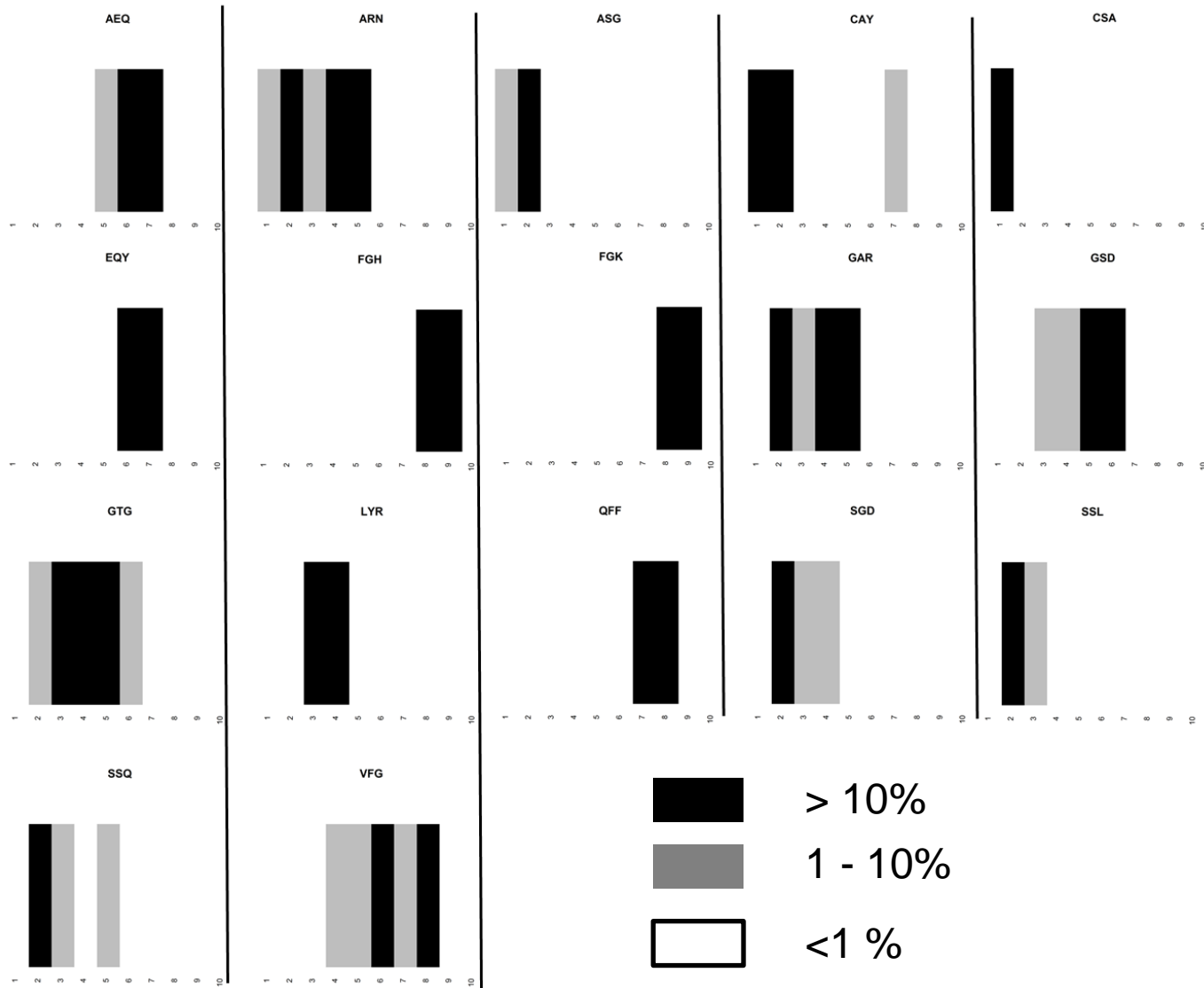


Linear boosting using 50 thousand CDR3s from each mouse, leave-one-out validation

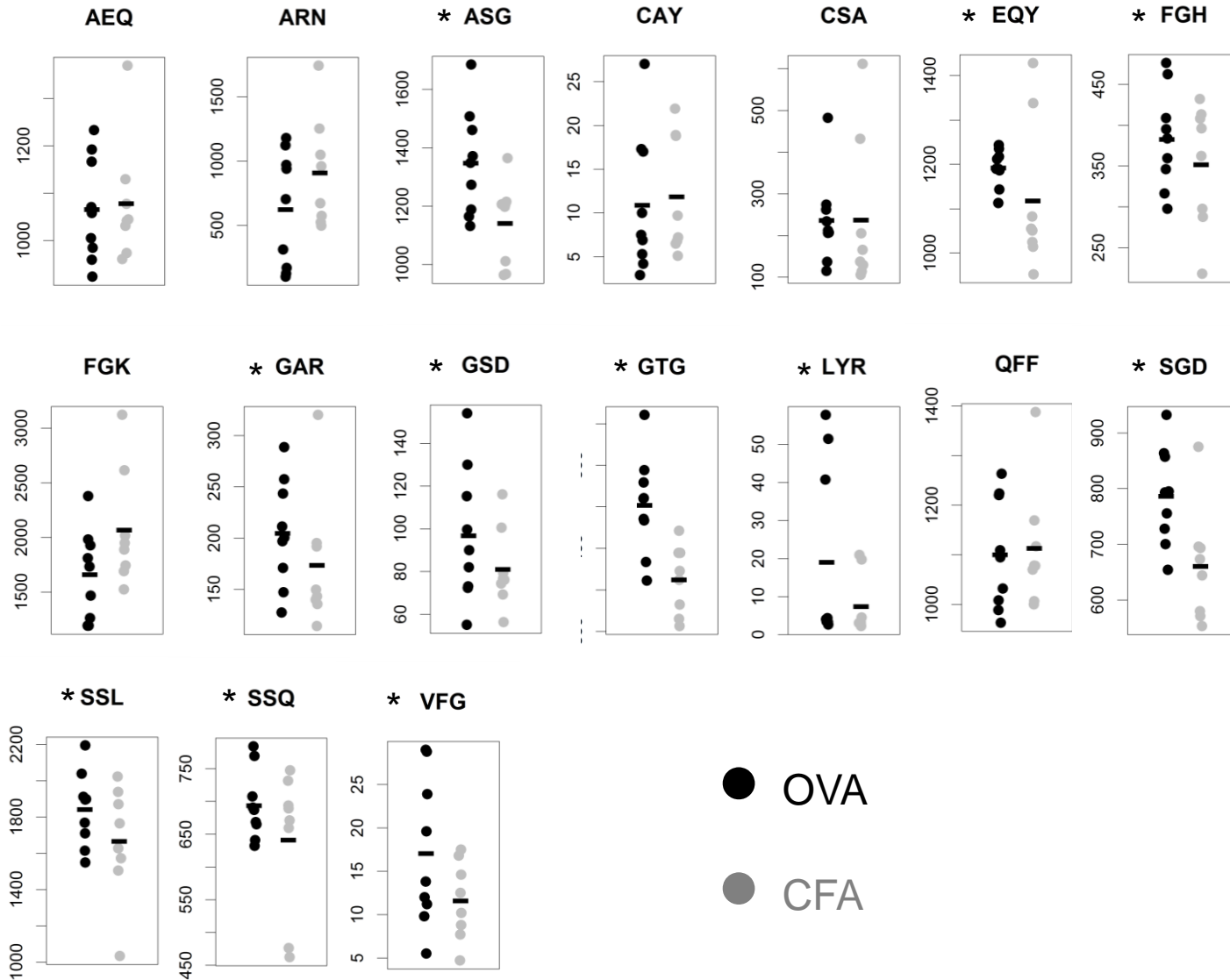
Features selected by LPBoost using counts, Fisher score weighted counts (bold) or both algorithms (red).



Position of OVA selected triplets along CDR3.



Frequency of OVA selected triplets in early immunised repertoires.

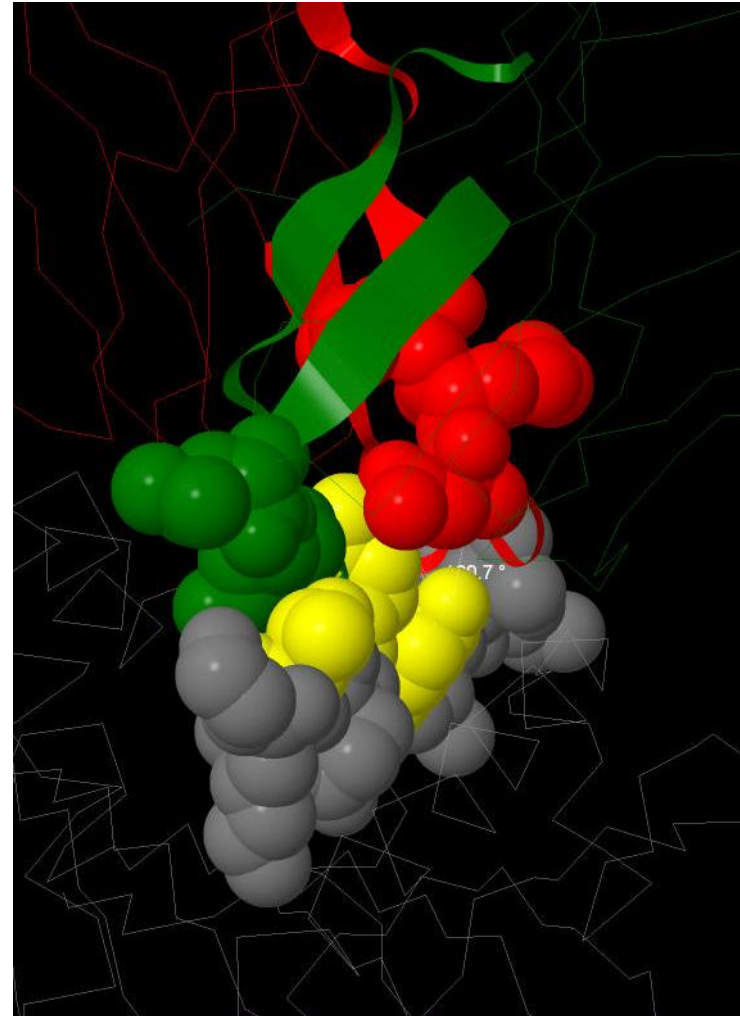
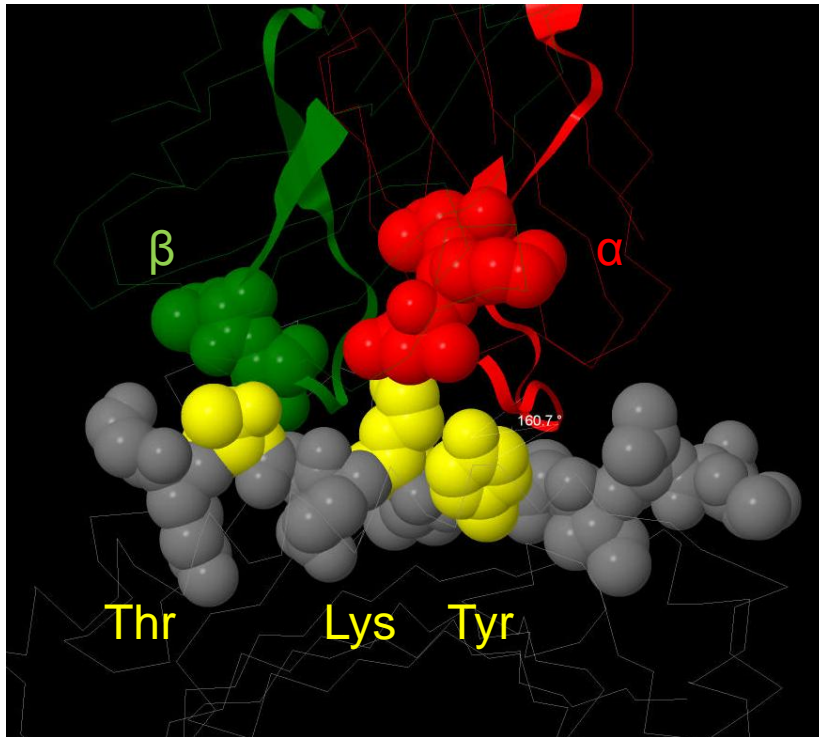


Conclusions

- The OVA response is diverse and predominantly private at the level of CDR3 β
- OVA expanded CDR3 β s have some sequence similarity
- Amino acid triplets provide features which in combination contribute to defining an OVA response
- Each triplet has a well-defined position along the CDR3
- Many selected triplets are found at the ends of the CDR3, within the sequence coded by V or J region genomic

Speculations (i.e. my future research program)

- The antigen response is distributed across many TCRs of differing frequency
- Antigen specificity is an emerging property of a repertoire and is not defined by individual TCR/pMHC interactions
- Individual TCR/pMHC interactions can be modelled as a set of interactions between individual antigen amino acids and a set of conserved short amino acids motifs within the CDRs



2B4 TCR binding cytochrome C peptide showing conserved residues