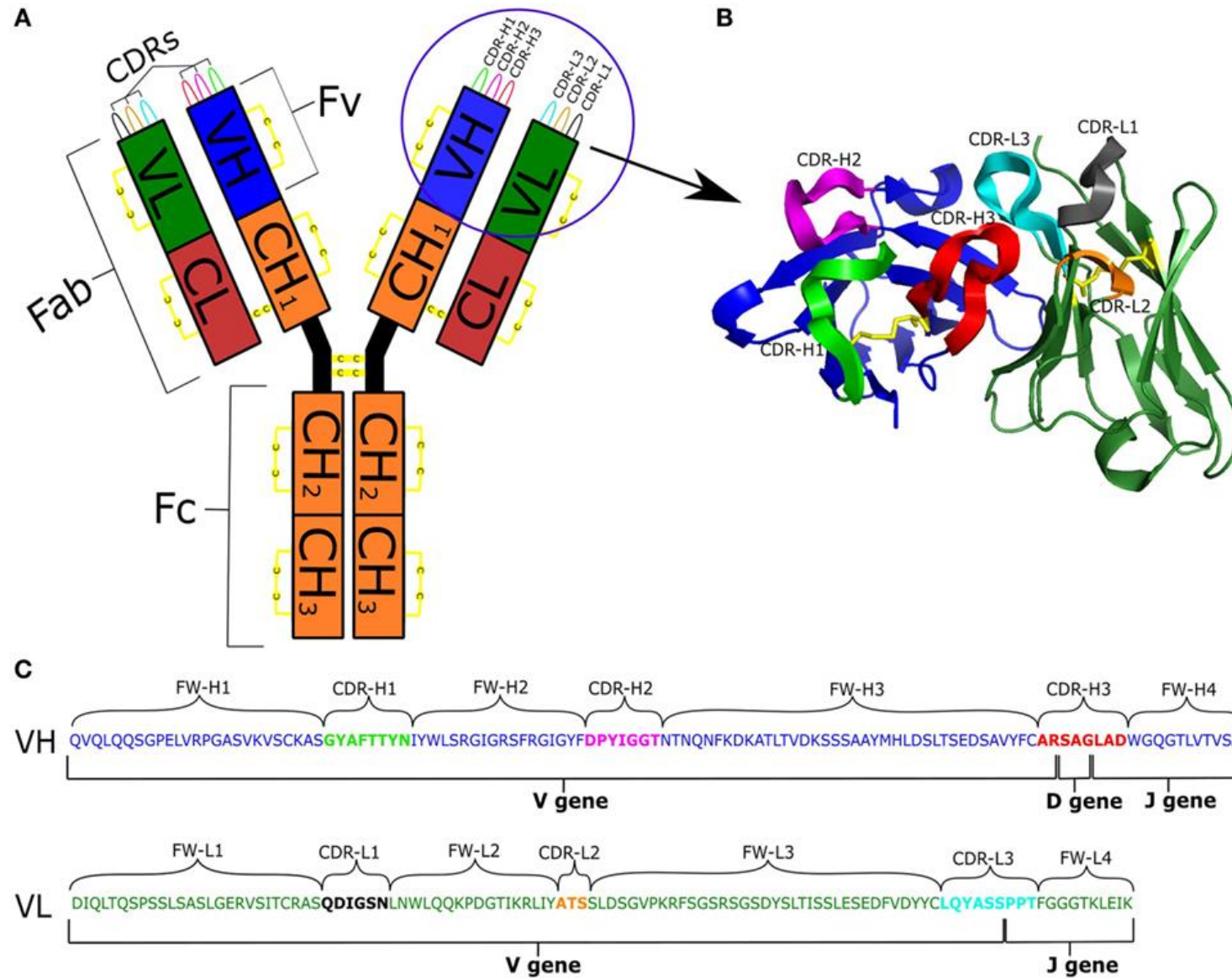


**Using structural information to aid
in-silico therapeutic design in the
era of immunoglobulin repertoire
sequencing**

Charlotte Deane
Department of Statistics
Oxford University

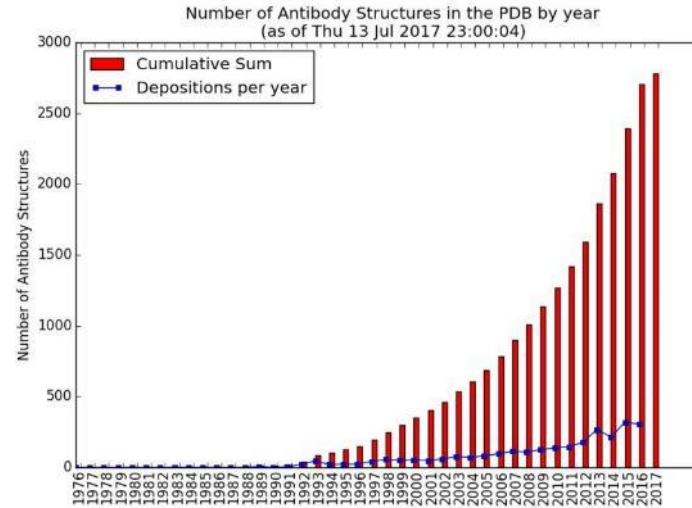
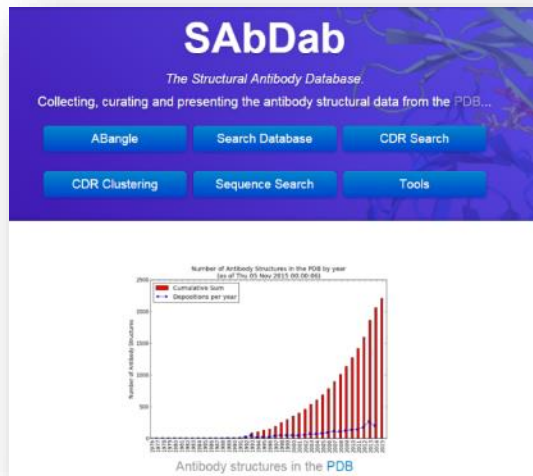
Antibody structure



Structural Antibody Database



Public data repository



Antibody structures in the PDB

Fully automated updating collection of all publicly available antibody structure data

Collect, curate and present structural data consistently.

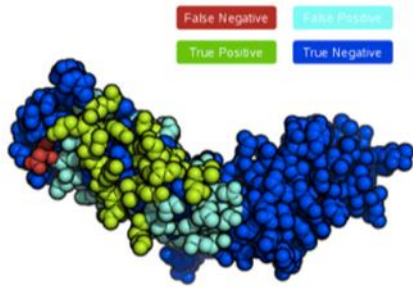
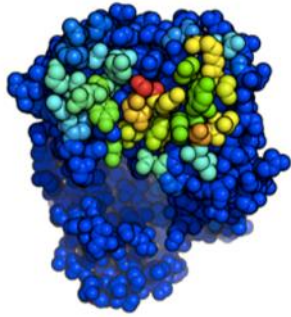
As of March 22nd 2019

3495 antibody structures

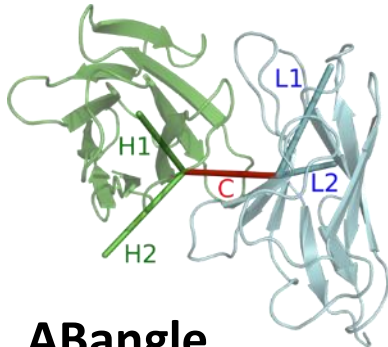
~2398 antibody antigen complexes

<http://opig.stats.ox.ac.uk/webapps/sabdab>

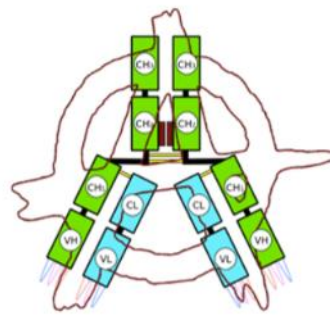
Dunbar *et al.* (2014) *NAR*



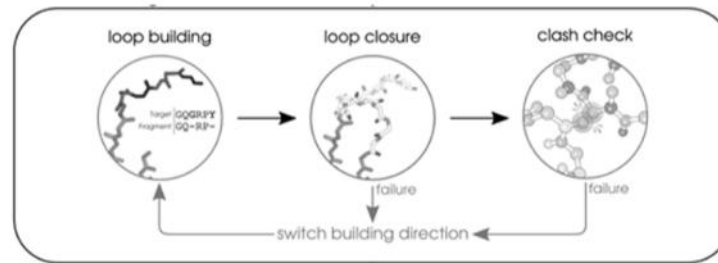
i-Patch, EpiPred:
Krawczyk et al.,
2013; 2014



ABangle
Dunbar et al. 2013

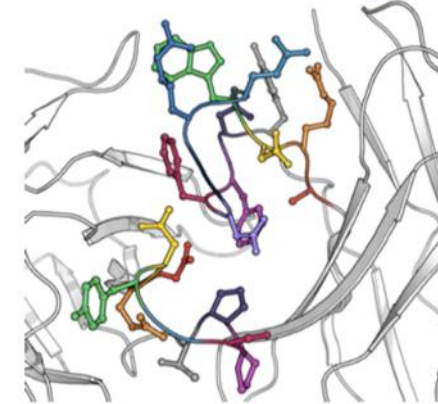
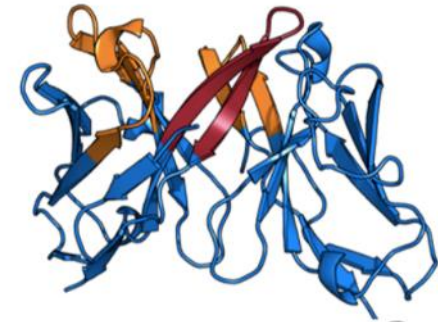
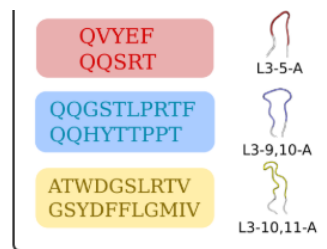


ANARCI: Dunbar & Deane, 2016.

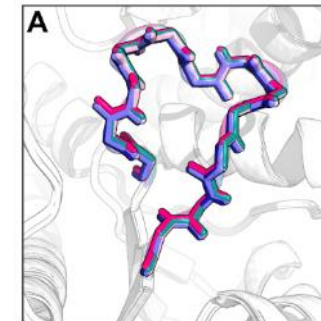


Sphinx: Marks et al., 2017.

SCALOP
Wong et al.

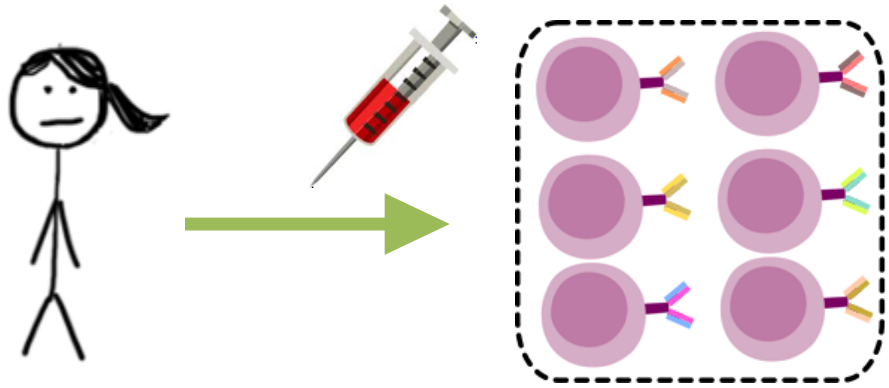


ABodyBuilder, PEARs:
Leem et al., 2016; 2018



FREAD
Choi et al. 2009

Antibody next-generation sequencing (Ig-seq)



```
GGTCCCTGAGACTCTCCTGTGCAGCCTCT  
GGATTCACCTTTGATGATTATGCCATGCAC  
TGGGTCCGCCAAGCTCCAGGGAAGGGCT  
GGAGTGGGTCTCAGGTAAGTTGGAGTA  
GTAGTCCATAGGCTATGTGGACTCTGTGA  
AGGGCCGATTCAACATCTCCAGAGACAAC  
GCCAAGAACTCCCTGTATCTGCAAATGAAC  
AGTCTGAGAGTTGAGGACACGGCCTTATAT  
TACTGTGCAAAAGATGTTCTTAGCCGCAGC  
TGGCGATATCTTGACCCCTGGGGCCATGGA  
ACCCTGGTCACCGTCTCCTCAGCATCCCCG  
ACCAGCCCCAAGGTCTTCC
```



- Snapshots between 10^4 and 10^7 antibody sequences
- Theoretical antibody repertoire $>10^{10}$
- Naïve human antibody repertoire
- Pre and post immunisation datasets
- Sequences repertoires from different species

Observed Antibody Space: OAS

- The first organized collection of a large body of antibody NGS outputs
- 60 studies covering ~ 1 billion antibody sequences across diverse immune states, organisms and individuals.
- We have sorted, cleaned, annotated, translated and numbered these sequences and made the data available
- OAS will facilitate data mining of immune repertoires for improved understanding of the immune system and development of better biotherapeutics.
- The data is all made freely available at antibodymap.org

by Oxford University

Observed Antibody Space

Database of antibody sequences from next generation sequencing analysis of immunoglobulin gene repertoires (Ig-seq). If you would like to read more on this work and/or cite it please see this [pre-print](#).

Bulk data download

Download the entire OAS database [here](#) (nucleotides are [here](#)). This contains >600m redundant antibody sequences from 53 studies. Please refer to the [documentation](#) to learn of the data formats.

Online OAS Search

Specify parameters to return a specific subset of OAS.

NB: the search fields are not exclusive, so if you pick a combination of fields that does not exist in our database, it will yield no results.

Longitudinal: ?
*


Chain: ?
*

Isotype: ?
*

Age: ?
*

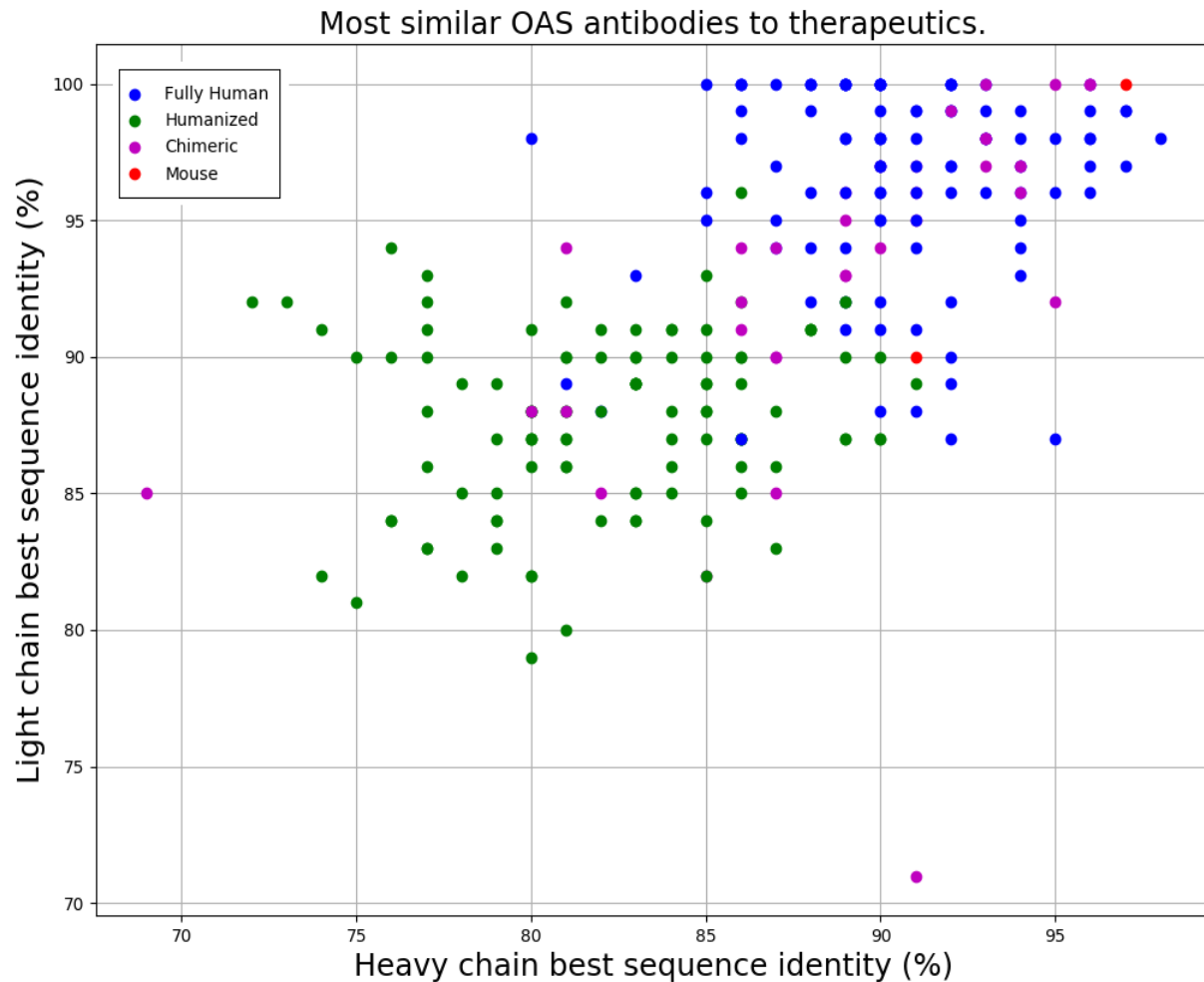
Disease: ?
*

BSource: ?
*



Kovaltsuk *et al.* (2018). *Journal of Immunology*

Looking for Therapeutic Antibodies in OAS



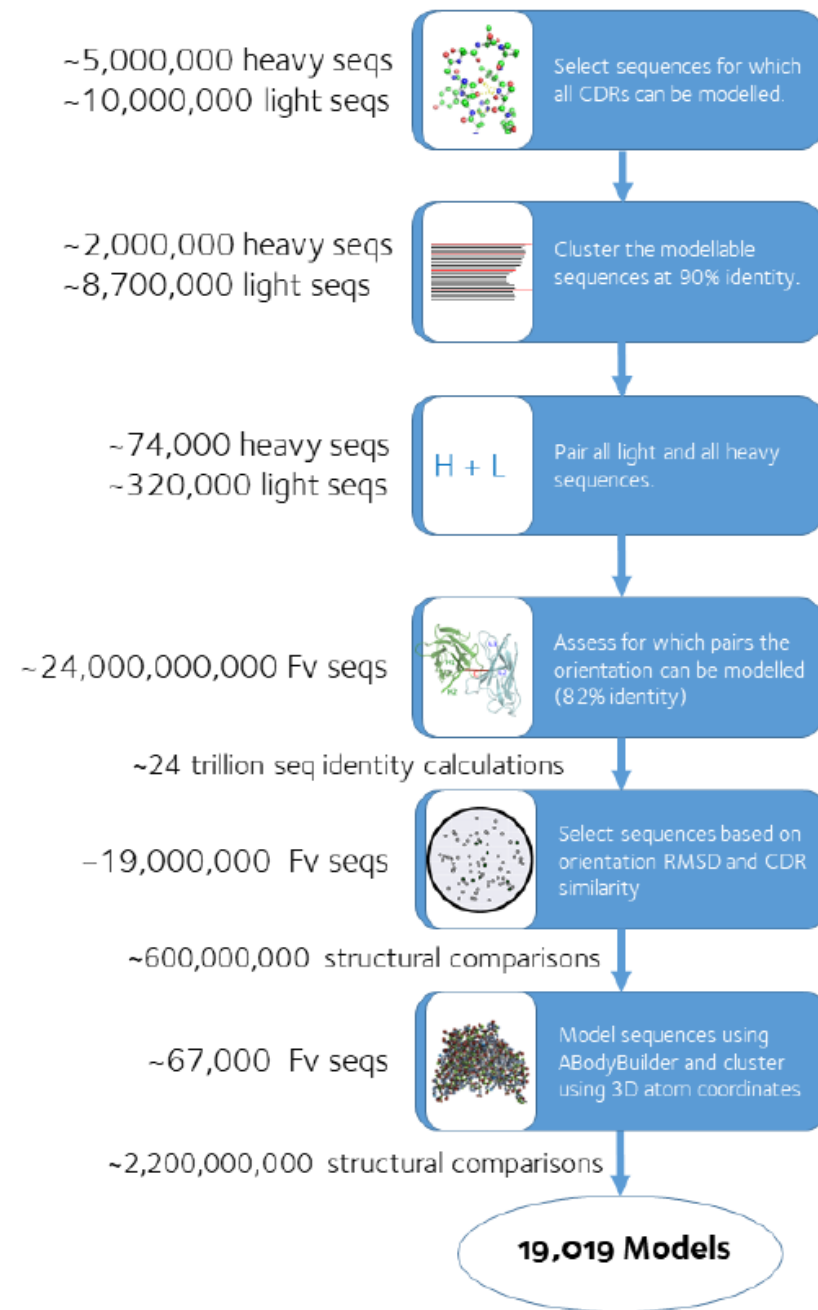
- 242 post phase-1 antibodies
- Unexpected high percentages of sequence overlap with therapeutics
- Many can be found in OAS with sequence identities >95%
- Enfortumab, heavy and light chain have 98% seqID
 - differences H38:N-S, H88:S-Y, L37:G-S, L52:F-L
- 54 have a perfect CDRH3 match
 - 22 of these found in more than one dataset

Building a human antibody model library

- Build a database of models that describe the structural variability of human antibody space
 - Input data is unpaired complete VH and VL sequences from the same sample
- Build a pairing and modelling protocol
 - to capture the sequence and structural diversity
 - within the constraints of modelability
 - and computational expense

A Human Antibody Model Library

- Ig-seq data from circa 500 healthy individual
- The sequences are unpaired naïve and memory IgM molecules sourced from peripheral blood, bone marrow and spleen.
- ~13.5 million unique sequences ~5m heavy variable regions and ~8.5m light variable regions.
- Reduces to ~20,000 structurally diverse antibodies that we can model accurately



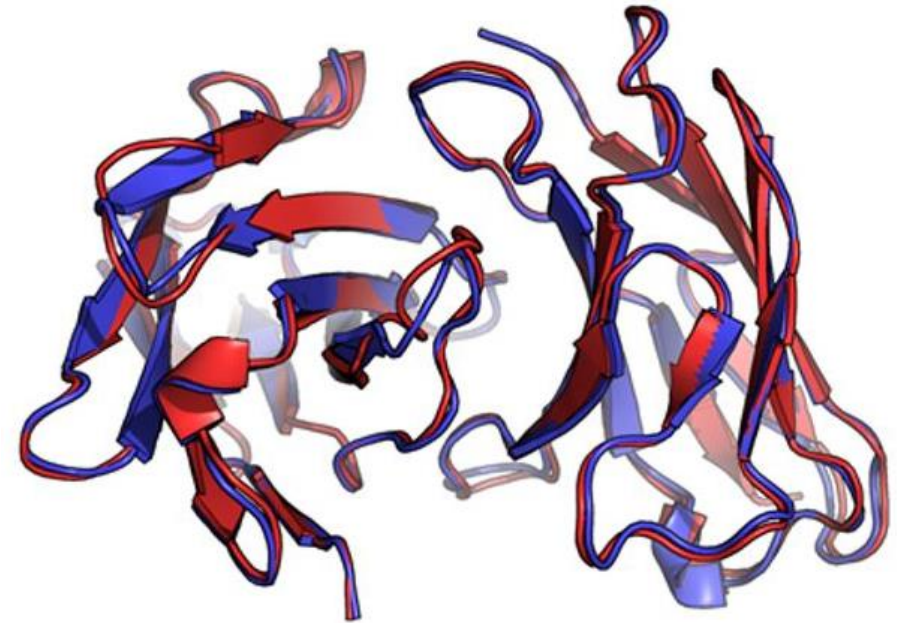
The Therapeutic Antibody Profiler (TAP)*

Five Computational Developability Guidelines

- Therapeutic antibodies must not only bind to their target but must also be free from 'developability issues' such as poor stability or high levels of aggregation.
- TAP is an *in-silico* antibody design analog of the Lipinski's rule of five for small molecules
 - to guide the selection of antibodies with appropriate biophysical properties
- Derive distributions of metrics for clinical stage therapeutics and assume that these indicate the allowed values of these properties.
 - Calculate these metrics on models so can run against potential therapeutics where crystal structures are unavailable
- These metrics don't have to correlate with a particular experiment that tests for developability rather they indicate that a potential therapeutic has outlying values.

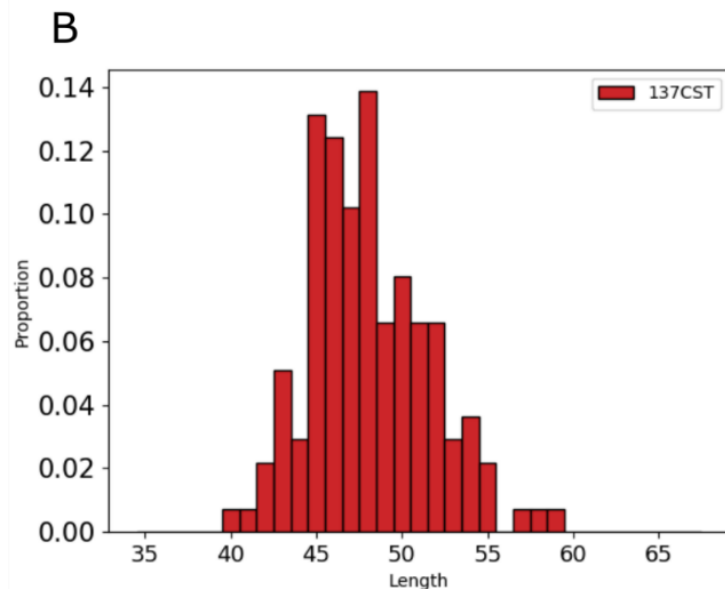
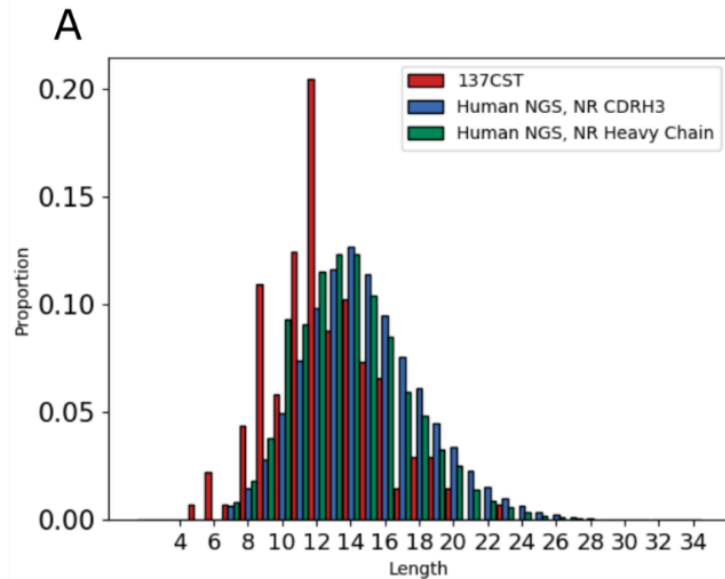
Datasets – structural models

- Models of the variable domain structures of 137 post-Phase I clinical-stage antibody therapeutics (CSTs)*
 - Models are accurate enough for our metrics (tested with the 56 CSTs with known structure)
 - Average RMSD of framework < 1Å
 - Less than 4% of residues are wrongly annotated exposed/buried
- For context we use the ~20,000 models of our larger human antibody model library
 - Would any human antibody make a good therapeutic?



*Jain et al (2017) PNAS

CDRH3 length and total CDR length



CDRH3s of our clinical stage therapeutics tend to be shorter than the human Ig-Seq data

So being a human sequence \neq good therapeutic

Could be for a number of reasons, including **affinity** (reduce entropy loss), or **aggregation**.

Use total CDR length as our metric (across all 6)

To perhaps better capture **binding site shape**

Closely correlates with CDRH3 length ($r = +0.77$, $p < 10^{-28}$)

TAP METRIC 1: TOTAL CDR LENGTH

CDR vicinity Patches of Surface Hydrophobicity (PSH)

$$\sum_{R_1 R_2} \frac{H(R_1, S) H(R_2, S)}{r_{12}^2}$$

R_1, R_2 are different residues

$H(R, S)$ hydrophobicity score

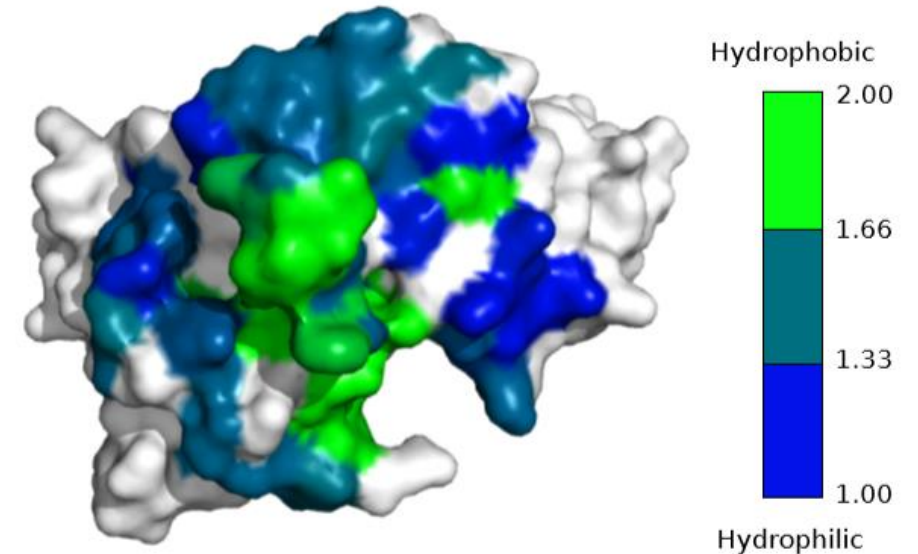
r_{12} – distance between R_1 and R_2

The set of residues considered are the CDR vicinity

All surface-exposed IMGT CDR residues, as well as other surface-exposed residues with a least one heavy atom within a radius of 4Å

TAP METRIC 2: CDR Vicinity PSH

B



Galiximab, with a high CDR Vicinity PSH value, has a large patch of hydrophobicity in its CDRH3 loop.

CDR vicinity Patches of Surface Hydrophobicity (PSH)

$$\sum_{R_1 R_2} \frac{H(R_1, S)H(R_2, S)}{r_{12}^2}$$

R_1, R_2 are different residues

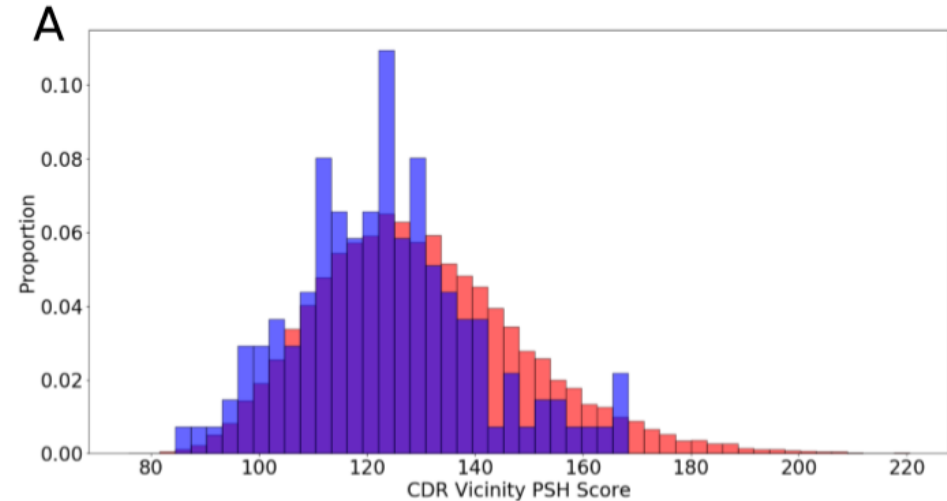
$H(R, S)$ hydrophobicity score

r_{12} – distance between R_1 and R_2

The set of residues considered are the CDR vicinity

All surface-exposed IMGT CDR residues, as well as other surface-exposed residues with a least one heavy atom within a radius of 4Å

TAP METRIC 2: CDR Vicinity PSH

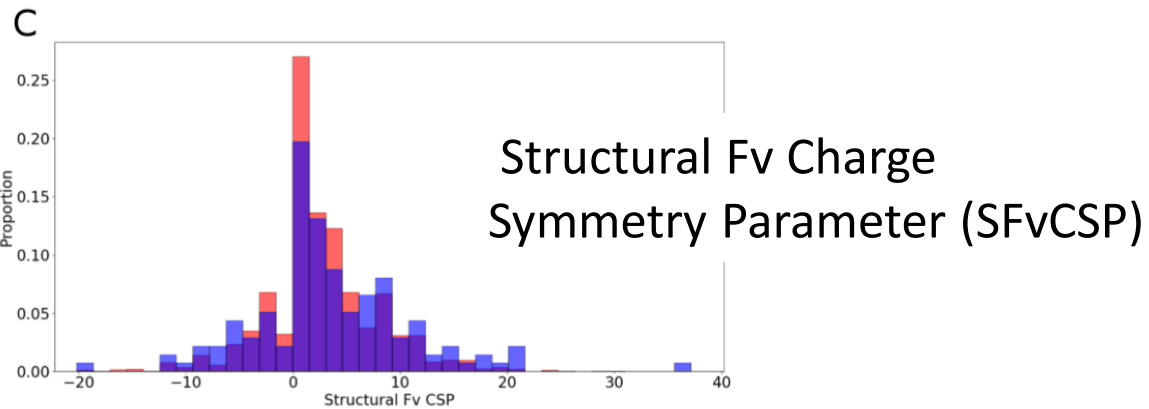
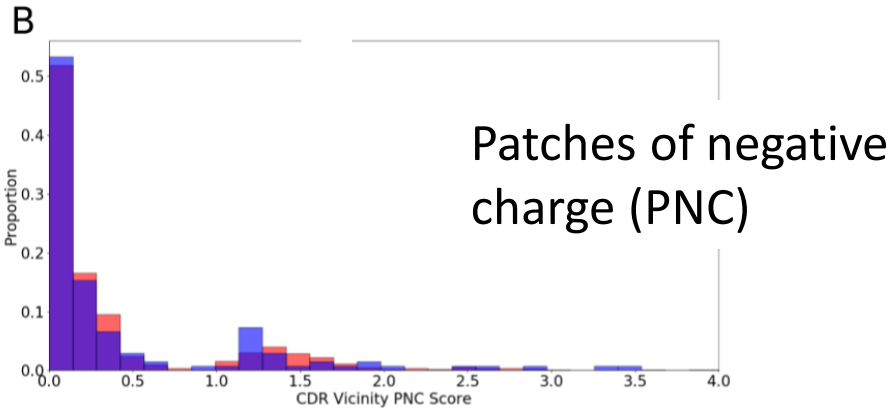
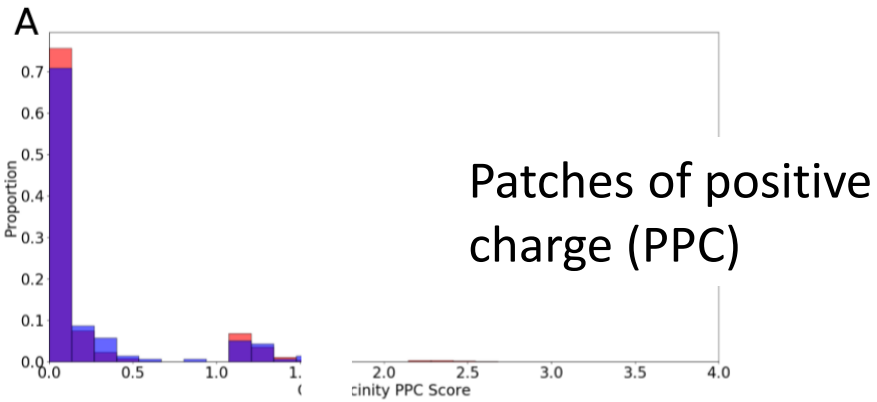


Blue are calculated for the models of the Clinical stage therapeutics

Red are the values for the Human antibody model library

Once again human sequence \neq good therapeutic

Charge



TAP METRIC 3: CDR Vicinity PPC

$$\sum_{R_1 R_2} \frac{|Q(R_1)| |Q(R_2)|}{r_{12}^2}$$

Q(R) charge of residue R

TAP METRIC 4: CDR Vicinity PNC

$$\sum_{R_H} Q_{R_H} \sum_{R_L} Q_{R_L}$$

RH, RL are surface-exposed VH, VL residues respectively.

TAP METRIC 5: SFvCSP

Five selected metrics – amber and red flags

- Total CDR Loop Length

Longer lengths \propto Aggregation, Lower Affinity, Convex Binding Site

- Patches of Surface Hydrophobicity (PSH) in the CDR Vicinity

Higher values \propto Aggregation (Self-association), High Viscosity

- Patches of Positive Charge (PPC) in the CDR Vicinity

Higher values \propto Poor Expression/Stability

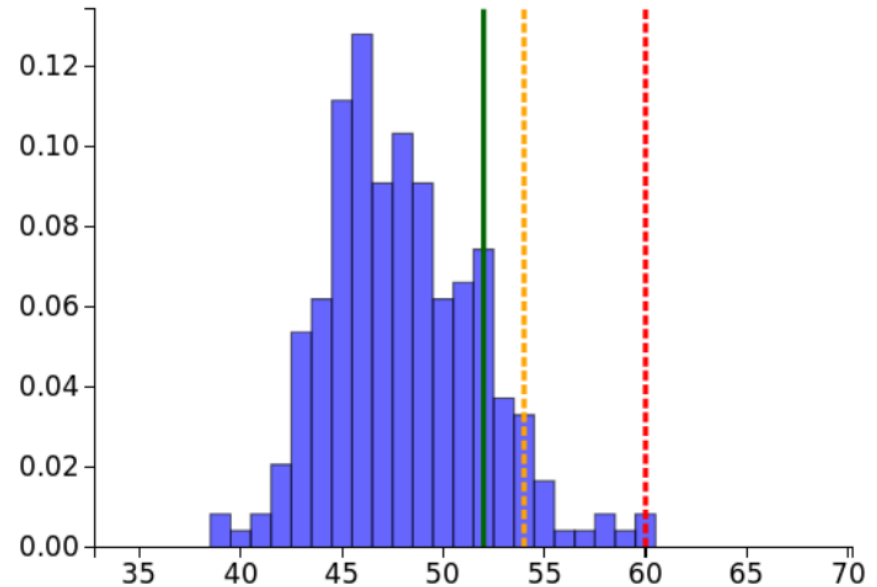
- Patches of Negative Charge (PNC) in the CDR Vicinity

Higher values \propto Poor Expression/Stability

- Structural Fv Charge Symmetry Parameter (SFvCSP)

More negative values \propto Colloidal instability, High Viscosity

Metric 1: CDR Length



Amber within 5% of highest/lowest value seen
Red above any value seen

Validation: Things TAP shouldn't flag

- Tested against 105 extra post-Phase I therapeutics

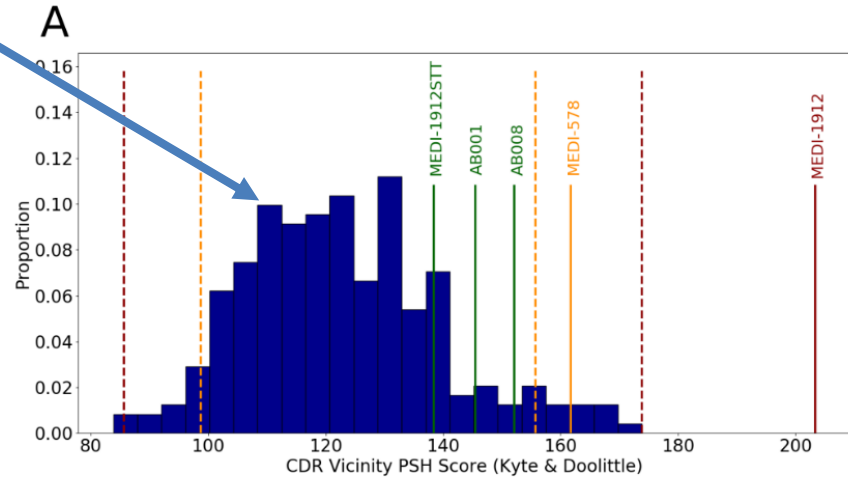
Metric	137 CST Amber Flag Region	Number Amber Flagged	137 CST Red Flag Region	Number Red Flagged
Total CDR Length (L)	$54 < L \leq 59$	6	$L > 59$	2*
PSH, CDR Vicinity (Kyte)	$85.65 \leq \text{PSH} < 98.74$	2	$\text{PSH} < 85.65$	1
	$155.76 < \text{PSH} \leq 171.91$	5	$\text{PSH} < 171.91$	1*
PPC, CDR Vicinity	$1.23 \leq \text{PPC} < 1.51$	1	(> 1.51)	5*
PNC, CDR Vicinity	$1.90 \leq \text{PNC} < 3.50$	4	(> 3.50)	0
SFvCSP	$-39.00 \leq \text{SFvCSP} < -18.00$	1	(< -39.00)	1

*Erenumab flagged for each of these properties

- Low red-flagging rate (8 of 105), implies won't pick out genuine therapeutics as having issues very often.

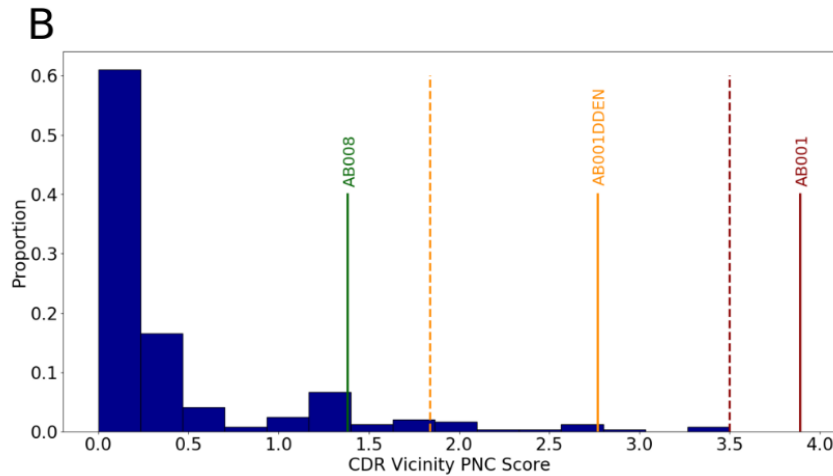
Validation: Things TAP should flag

242CSTs



Anti-NGF antibody (1)
MEDI-578: Minor aggregator;
MEDI-1912: Extreme aggregator;
MEDI-1912STT: No aggregation.

TAP predicts a similar behaviour from just the sequence and its ABodyBuilder model



Anti-IL13 candidate (2)
AB008: Good expression;
AB001: Poor expression (7x lower than AB008);
AB001DDEN: Good expression.

TAP flags AB001 for further consideration.

(1) Dobson et al (2016) Sci Reports

(2) Popovic et al (2017) Protein Eng Des Selec

Splitting by Species Origin

Table S8. 242 CST TAP values split by species origin.

TAP Metric	101 Human ($\mu \pm \sigma$)	108 Humanized ($\mu \pm \sigma$)	30 Chimeric ($\mu \pm \sigma$)	3 Mouse ($\mu \pm \sigma$)
Total CDR Length	48.68 \pm 4.09	47.80 \pm 3.42	46.77 \pm 3.55	46.33 \pm 1.25
PSH	127.76 \pm 18.56	120.90 \pm 14.20	115.73 \pm 15.58	117.26 \pm 9.44
PPC	0.29 \pm 0.58	0.20 \pm 0.36	0.26 \pm 0.55	0.05 \pm 0.06
PNC	0.34 \pm 0.56	0.50 \pm 0.75	0.30 \pm 0.63	0.50 \pm 0.50
SFvCSP	4.06 \pm 7.44	3.13 \pm 7.80	3.29 \pm 5.99	7.58 \pm 6.75

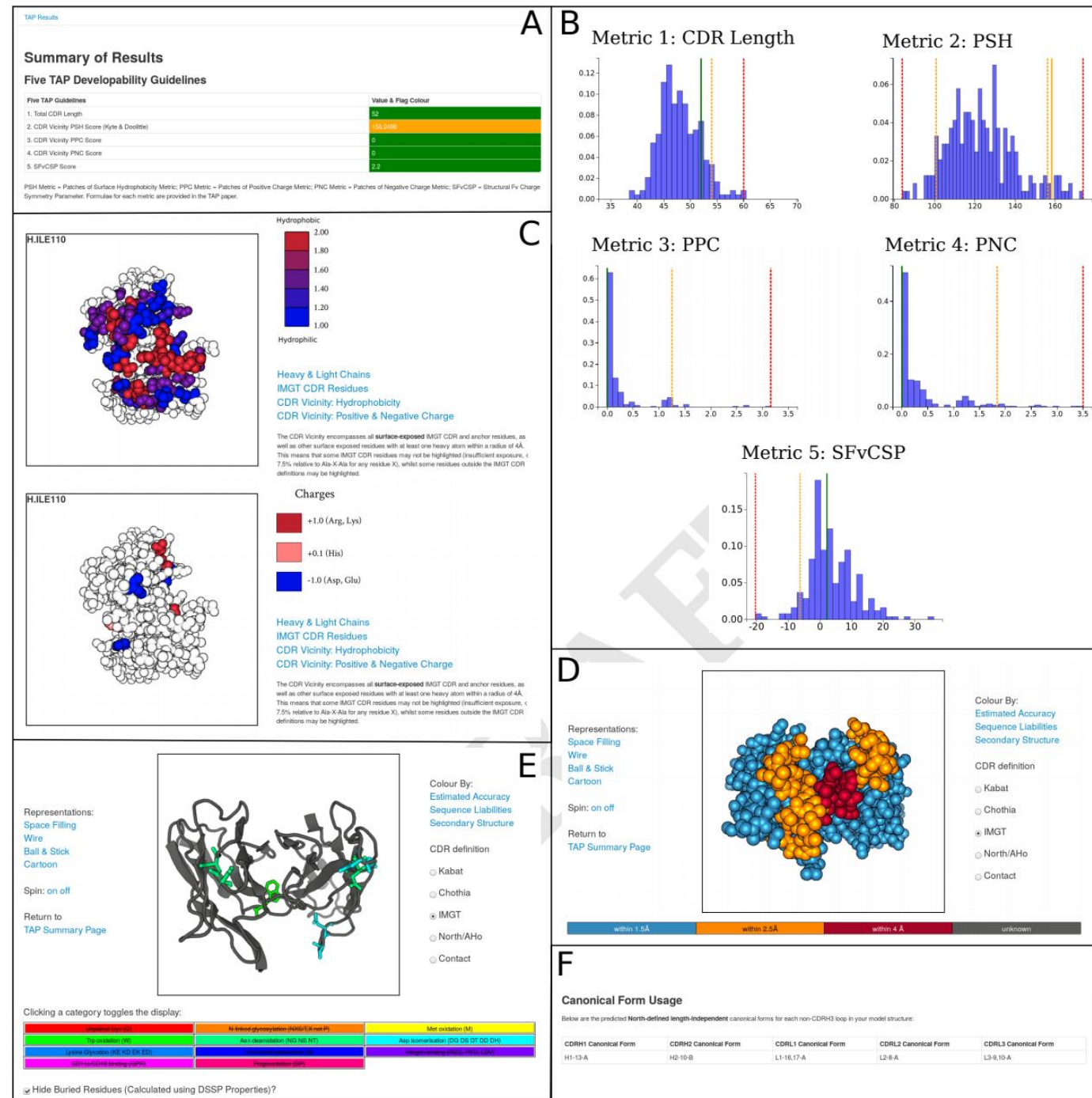
- Appears that the more human mAbs have larger patches of hydrophobicity than mouse mAbs
- We also split by clinical progression (P2, P3, Approved) and drug campaign status (active/discontinued) but found no significant differences in TAP metric values.

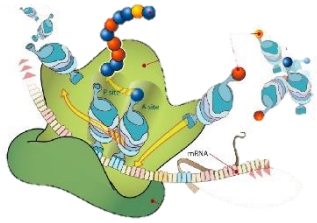
TAP website

Distributions will become more accurate as model quality improves (more up-to-date version of SAbDab) or as more therapeutics come to Phase-II of clinical trials.

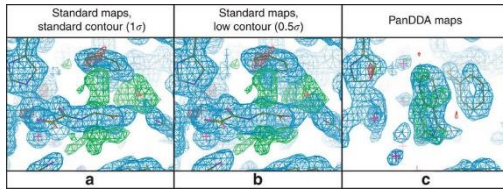
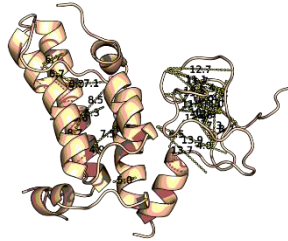
<http://opig.stats.ox.ac.uk/resources>

Raybould et al (2019) PNAS

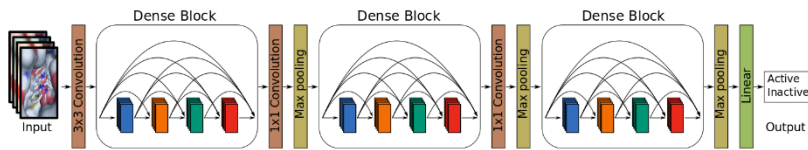




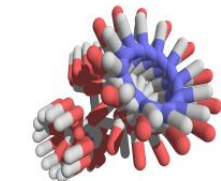
Biologically inspired protein structure prediction



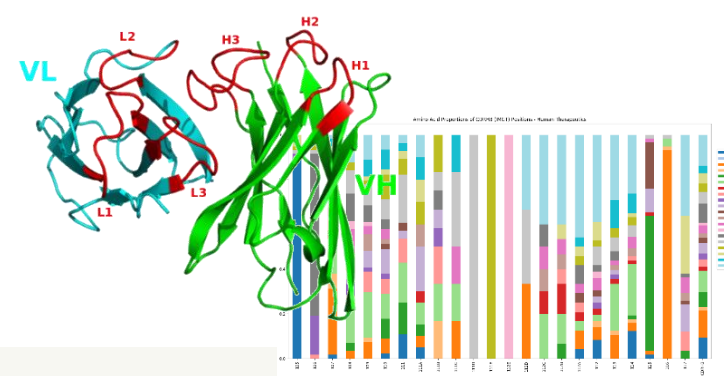
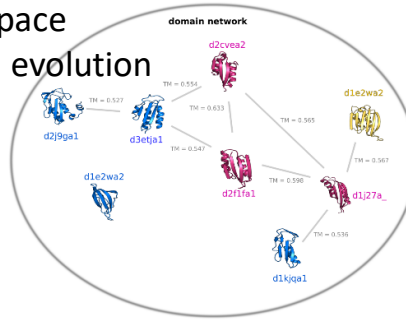
Developing software for Crystallography



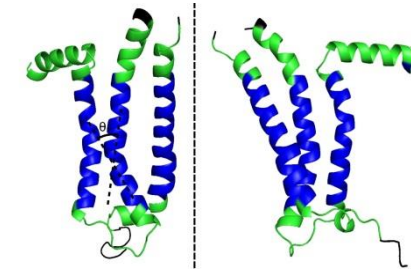
Small molecule drug discovery



Fold space & structural evolution



Immunoinformatics

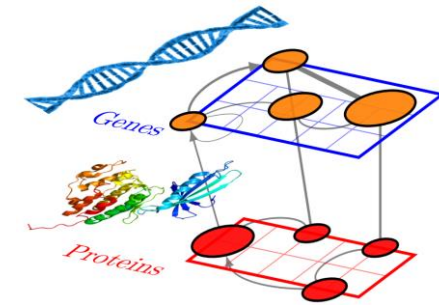
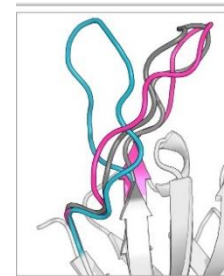


Membrane proteins

The Oxford Protein Informatics Group

“more than just antibodies”

Loop structure Conformations /Prediction

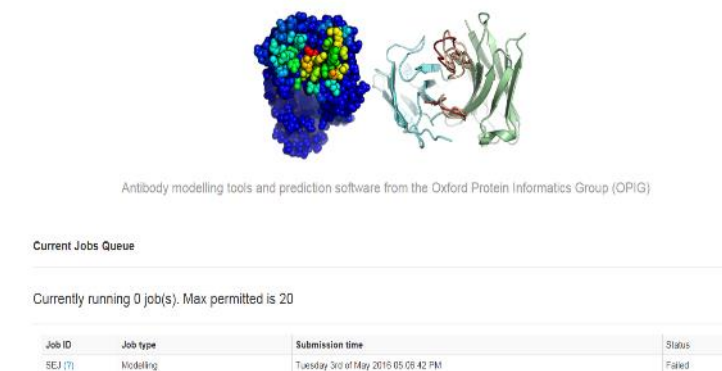
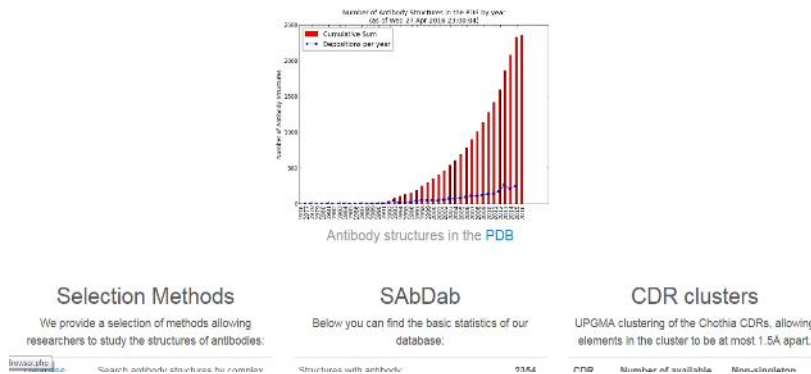
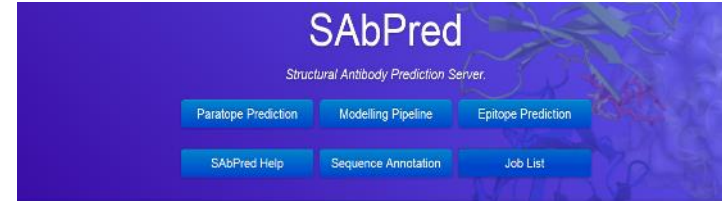
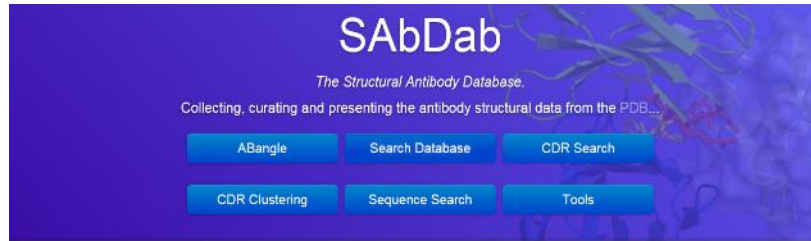


Networks

ACKNOWLEDGEMENTS



SAbDab & SAbPred – a computational antibody design platform



Antibody modelling tools and prediction software from the Oxford Protein Informatics Group (OPIG)

Current Jobs Queue

Currently running 0 job(s). Max permitted is 20

Job ID	Job type	Submission time	Status
SEJ (7)	Modelling	Tuesday 3rd of May 2015 05:09:42 PM	Failed

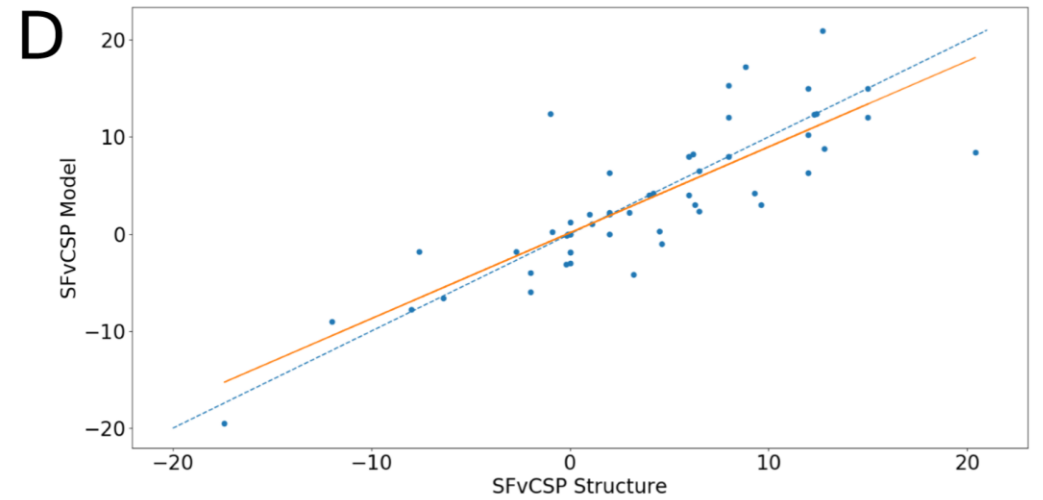
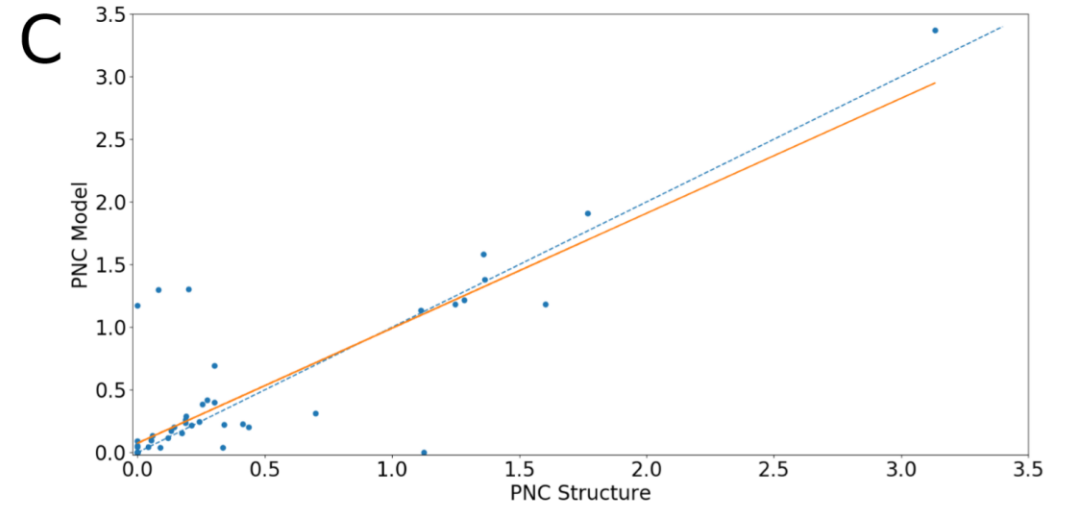
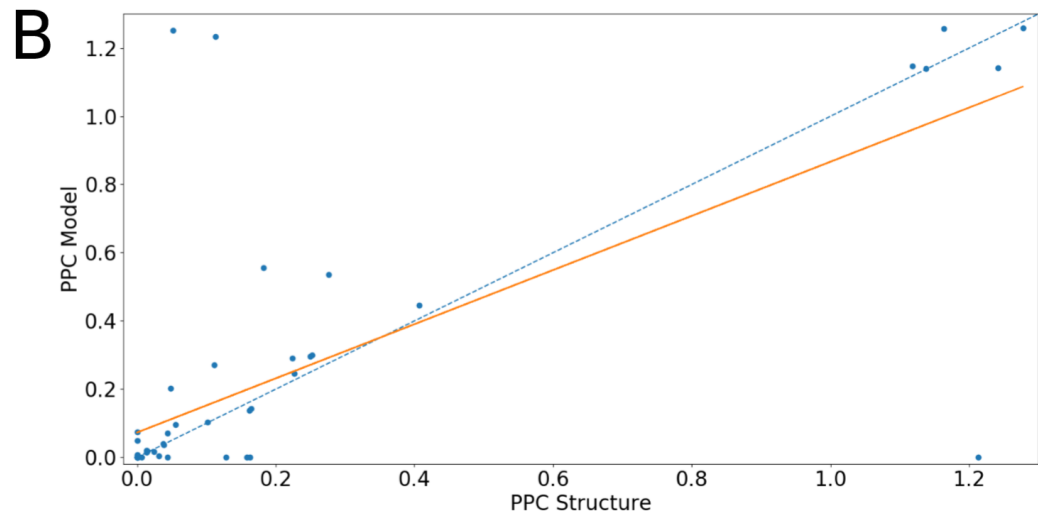
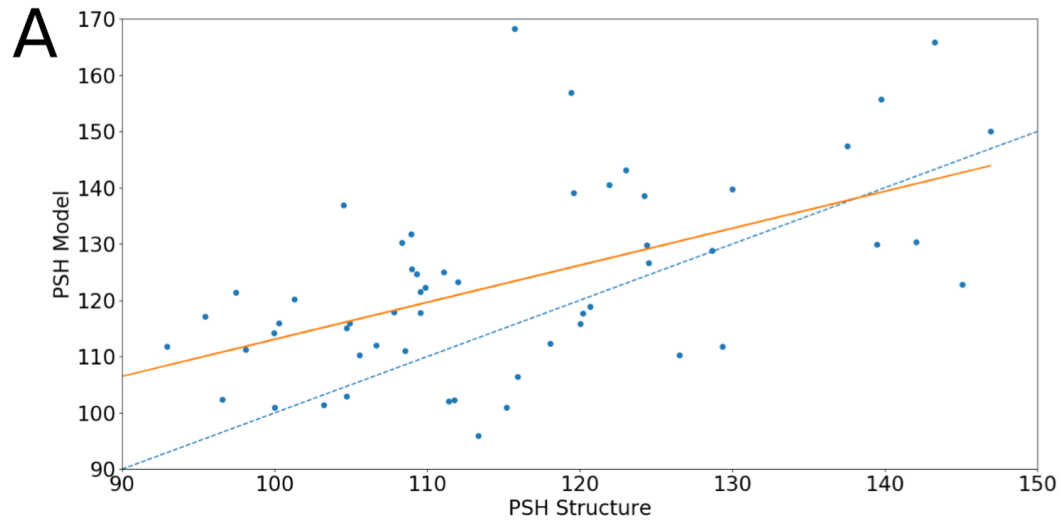
All our software is open and free to use

We offer consulting to create custom install specific for your infrastructure and/or training on the platform.

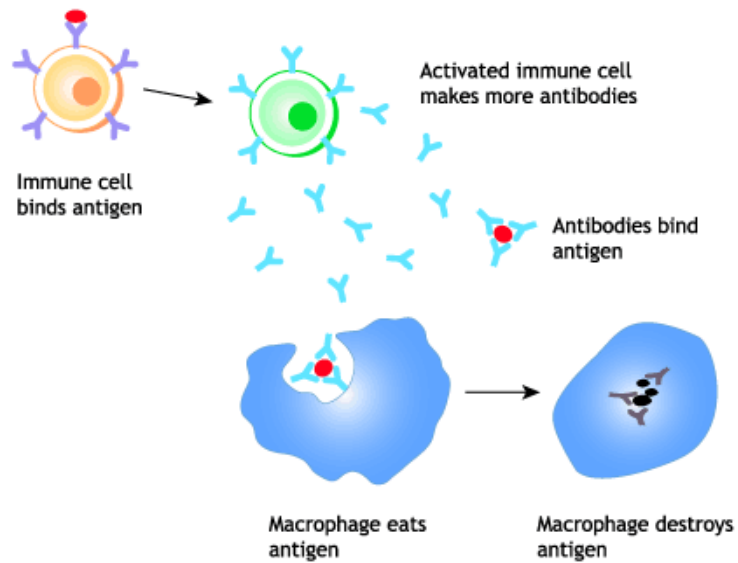
Newly available a VirtualBox virtual machine encompassing SAbDab/SAbPred.

More information contact deane@stats.ox.ac.uk

Values for models vs real structures



Antibodies



- It has been estimated that a typical human is capable of producing more than 10^{10} different antibodies, each capable of binding a distinct epitope
- Recognise and bind to potentially harmful molecules (antigens)
- Either inhibit the antigen themselves or recruit other parts of the immune system to deal with them



- Target specifically and with high affinity
- Can be raised against almost any antigen
- Effective treatment for immune disorders, cancers and arthritis
- Currently >60 approved antibody “drugs”
- Hundreds in late stage development

Beyond antibodies

Nanobodys

Home SABDac SABProt Software Help

H.GLN110 Accuracy: s4A

Annotation options:
Secondary structure
E-estimated accuracy
Sequence liabilities
Solvent E-xposed
Domains

By CDR definition:
 Kabat
 Chothia
 IMGT
 North/VAho
 Contact

Display options:
Spaccfill
Wire
Ball&stick
Cartoon
Spin: on off

Return to
Results Page

within 1Å within 2.5Å within 4 Å unknown

Your model should be shown above (we recommend using google-chrome to load this page; also have WebGL enabled.).
To Alternatively, download your [model](#) (PDB file).

Estimated model accuracy annotations

- Shown below are the img-numbered sequences of the modeled VH and VL domains
- Each position is coloured according to the level of confidence to which it has been modelled. [low]
- For example, a blue coloured position is estimated to be modelled to within 1Å with at least 95% confidence.
- Both the accuracy threshold and the confidence threshold can be changed using the sliders below.

VH domain:

IMGT Target: 123456789101112131415161718192021222324252627282930313233343536373839404142434445464748495051525354555657585960616263646566676869707172737475767778798081828384858687888990919293949596979899100

TCRS

