

Bioinformatics for crystallographers

Dan Rigden



UNIVERSITY OF
LIVERPOOL

Preamble

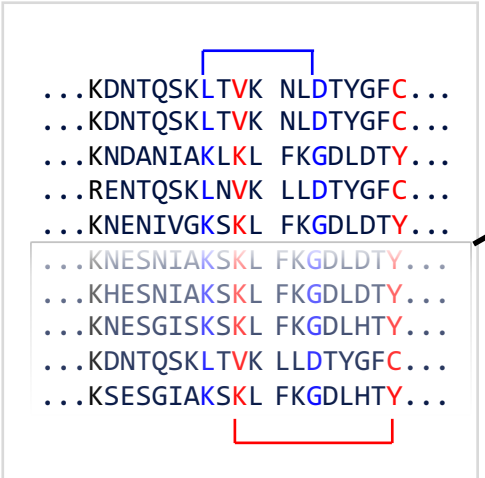
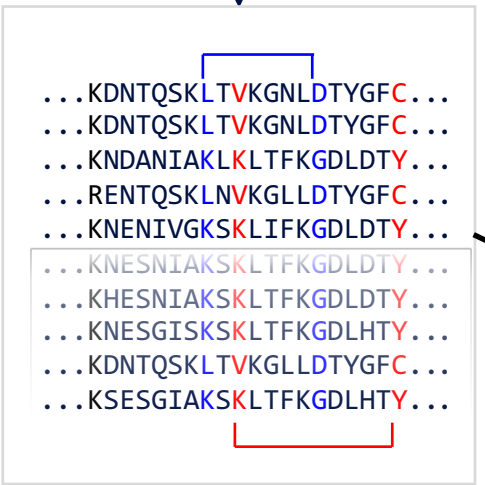
- Can't cover all bioinformatics!
- Prof Garratt will cover structure-based function inference
- I will focus on bioinformatics (= prediction)
 - Using newer/less well-known data
 - Majoring on easily available servers/predictions
 - Related to predicting domain composition and interactions. Mainly relevant to **construct design** and **MR** (all MR, not just AMPLE!)
- Plus some other bits and pieces

A new source of data: further uses of predicted contacts

- Structures unknown
 - Making better *ab initio* models for MR (AMPLE)
 - Predicting domain boundaries
- Structures known
 - Predicting how proteins interact
 - Validating the content of your crystal structure
 - (Predicting functional sites, highlighting conformational states)

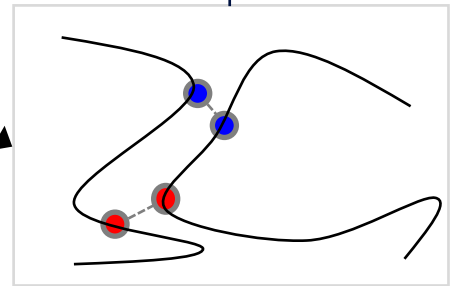
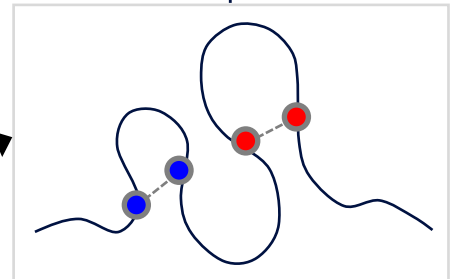
```
>target_seq (s)
GYEYLYRRSTIGNSLVDALDTLISDG
RIEASLAMRVLETFDKVVVAETLKDNT
QSKLTVKGNLDYGFCDVWTFIVKN
CQVTVEDQSVISVDKLRIVACNSKKS
```

Multiple Sequence Alignment of homologous sequences

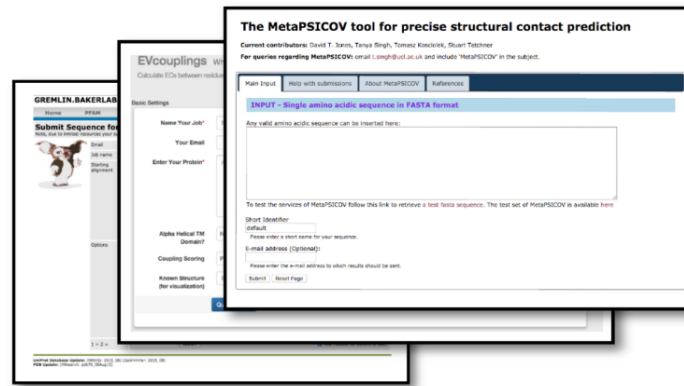


Direct Coupling Analysis of predicted contacts

Applications...




```
>seq
FASDGITF
DRSLFFGH
```

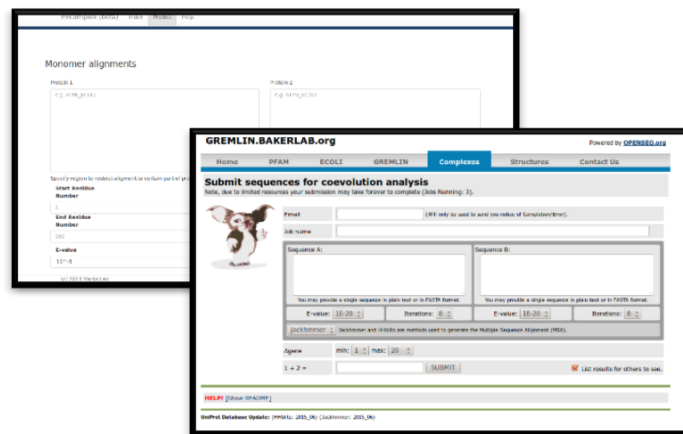


Intramolecular contacts

Intermolecular contacts if a homooligomer

```
>seq1
FASDGITF
DRSLFFGH
```

```
>seq2
NMHKL SDF
PLSERWAQ
```



Intermolecular contacts [currently only reliable for bacterial proteins in operons]

Thinking about your target ...

Which part to express for crystallisation?

Which parts of your crystallised protein might enable phasing by MR, or experimentally?

Recognising folded domains in your sequence

- How novel is your protein target? Recognising distant homology might make it less (or more!) interesting
- You might get extra ideas about its function to guide lab experiments, co-crystallisation, phasing (eg metal-binding sites)
- You might find the whole protein will not express/stay soluble/crystallise etc and want to deal with only part. You might well design construct to exclude disordered regions anyway.
- You might want to parse your protein into domains to explore different MR strategies

Recognising domains by homology with PDB, SCOP, CATH

- (PSI-)BLAST against the PDB might do it

[Edit and Resubmit](#) [Save Search Strategies](#) [Formatting options](#) [Download](#) [YouTube](#) [How to read this page](#) [Blast report description](#)

splq5i8r5| (257 letters)

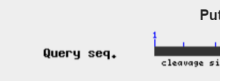
RID [FRP1NV47014](#) (Expires on 03-31 21:44 pm)

Query ID [gii74681646|splQ](#)
 Description Trypsin-like serine
 Molecule type amino acid
 Query Length 257

Other reports: [Search Summary](#)

Graphic Summary

Show Conserved Domains

Query seq. 

Specific hits
Superfamilies
Multi-domains

install-jalview.exe

[Download](#) [GenPept](#) [Graphics](#)

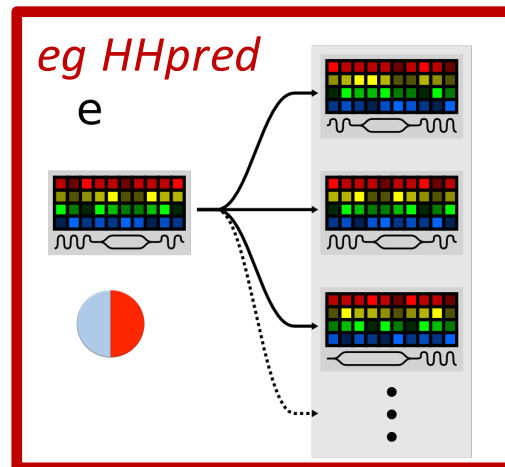
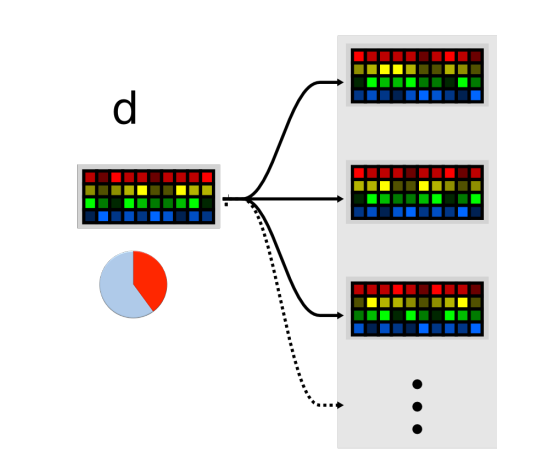
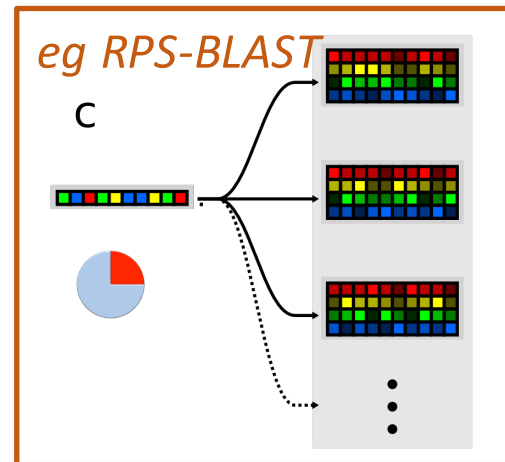
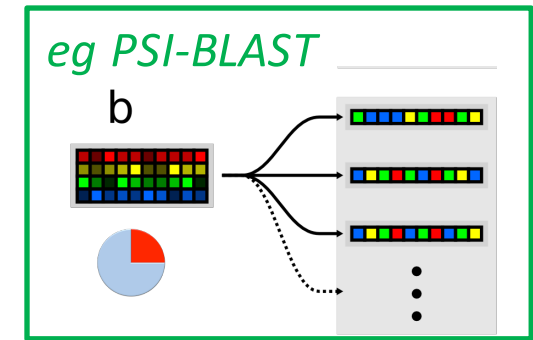
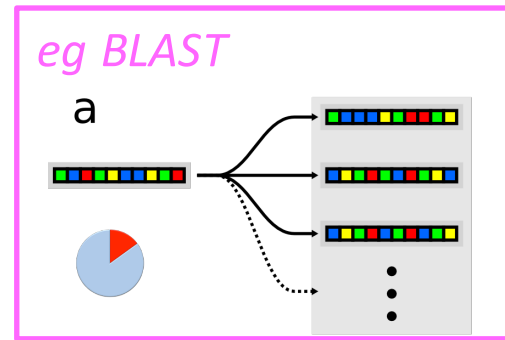
Chain A, Human Hepsin Protease In Complex With The Fab Fragment Of An Inhibitory Antibody
 Sequence ID: [pdb|3T2N|A](#) Length: 372 Number of Matches: 1
[See 1 more title\(s\)](#)

Range 1: 84 to 355 [GenPept](#) [Graphics](#) [Next Match](#) [Previous Match](#)

	Score	Expect	Method	Identities	Positives	Gaps
	142 bits(357)	6e-39	Compositional matrix adjust.	99/281(35%)	140/281(49%)	38/281(13%)
Query 1						
Sbjct 84						
Query 53						
Sbjct 143						
Query 102						
Sbjct 200						
Query 158						
Sbjct 255						
Query 216						
Sbjct 315						

Recognising domains by homology with PDB, SCOP, CATH

- Harder cases require a more sensitive tool. I recommend HHpred. Used by MrBUMP to find homologues to use as search models



HHpred tips and warnings

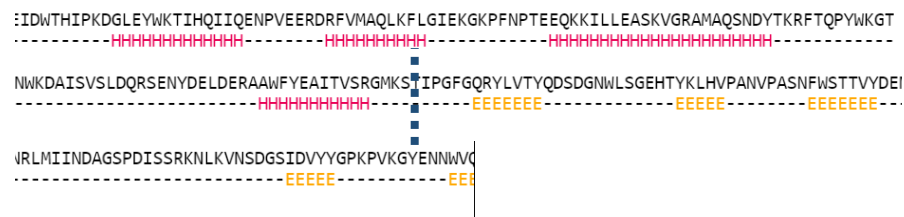
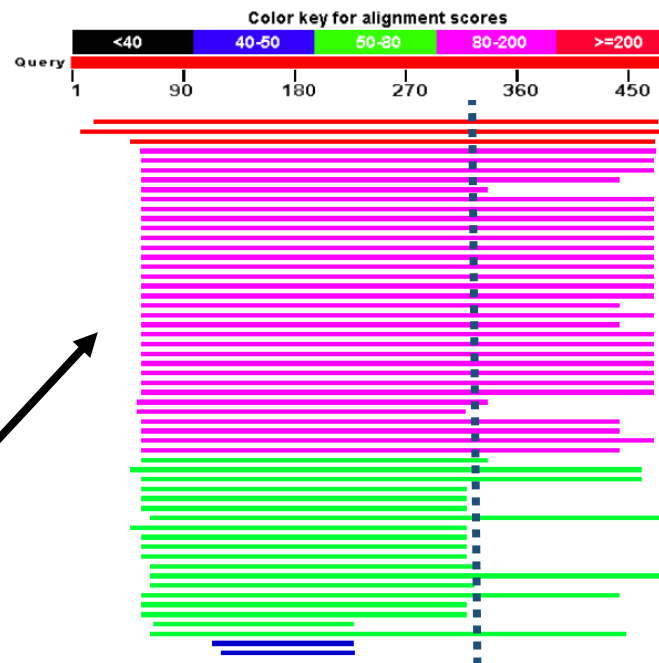
- **Probability** (0-100) is generally a good guide....
- ...but statistics can mislead for **unusual** protein sequences eg coiled-coil, low-complexity, Cys-rich
- Consider if the match makes biological **sense!**
- Look at the matched region
 - Is it a complete domain/structure?
 - If partial, could it reasonably fold?
- Can make reasonable homology models too
- If your query contains multiple domains, '**zooming in**' on particular regions can improve scores and show results not previously seen since
 - Score contains an element favouring similar lengths
 - Only 100 results are shown: can easily get this number for a common domain, making other results invisible!

Recognising matches to domains in **sequence** databases, Pfam, Smart etc

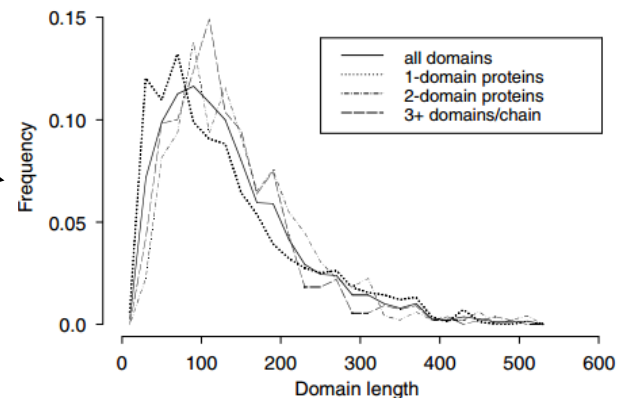
- HHpred is also an excellent, sensitive way to search against these. The same rules for interpretation apply. Matches, even distant ones, can obviously shed light on function
- However, for **structural** purposes need to remember that
 - Many Pfam families don't have structures so that domain limits are much less precise.
 - Many examples where a structure redefines previous Pfam sequence-only domain boundaries
 - Large Pfam entries especially DUFs often turn out to have multiple structural domains
 - Pfam entries for repeats sometimes contain multiple copies

And if the domains can't be matched?

- *ab initio* approaches to domain boundary identification
 - BLAST matches in sequence databases. Domains are often found in different combinations
 - Secondary structure pattern. PSI-PRED or Jpred4 (faster server)
 - Domain guess by size (server defunct)
 - Contact predictions

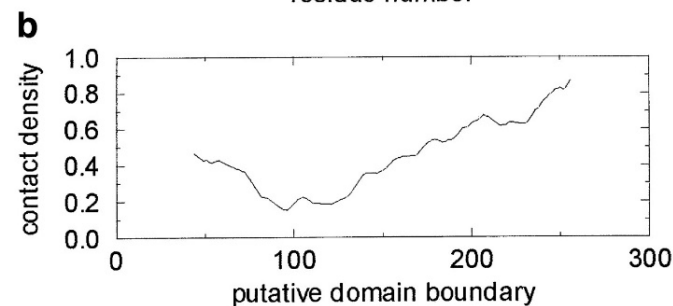
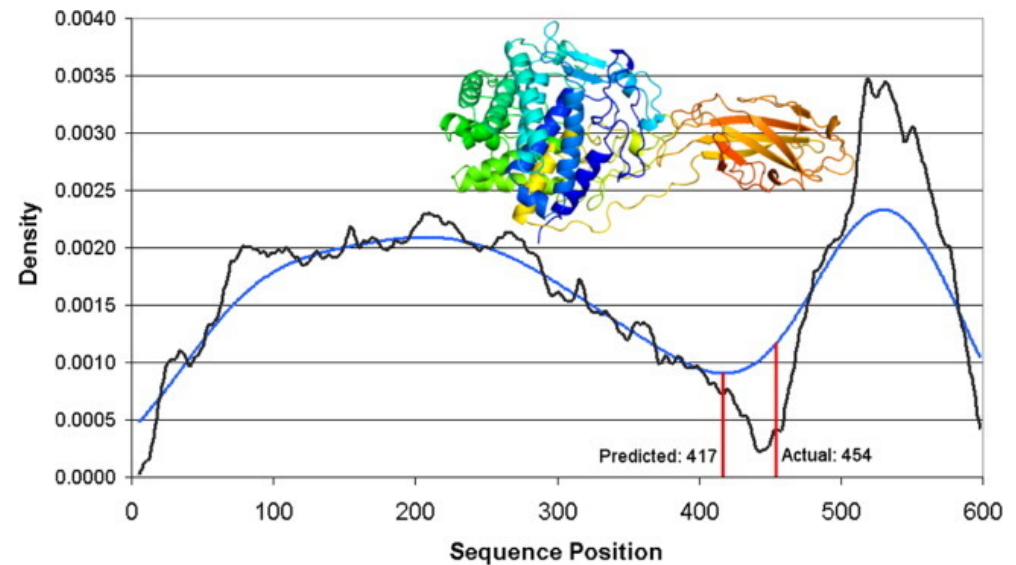
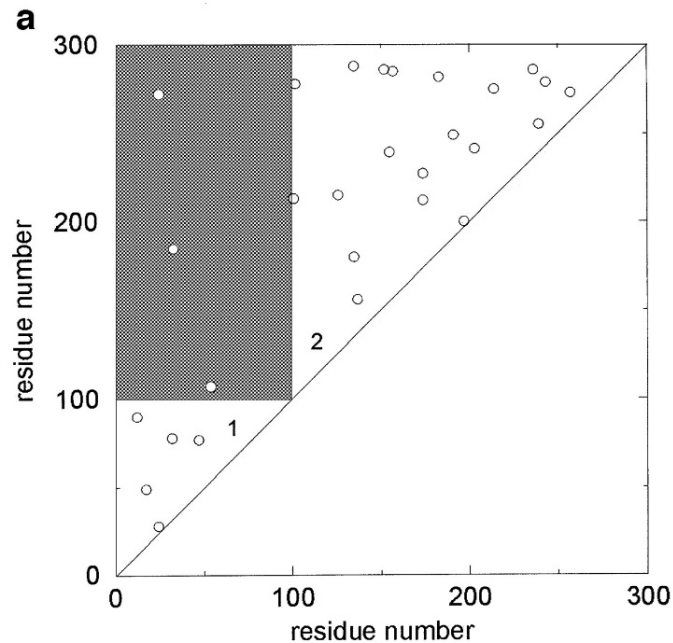


Size of protein domains



Predicted contacts for defining domains

- Domain boundaries required for individual expression *in vitro* and helpful for fold recognition and modelling *in silico*

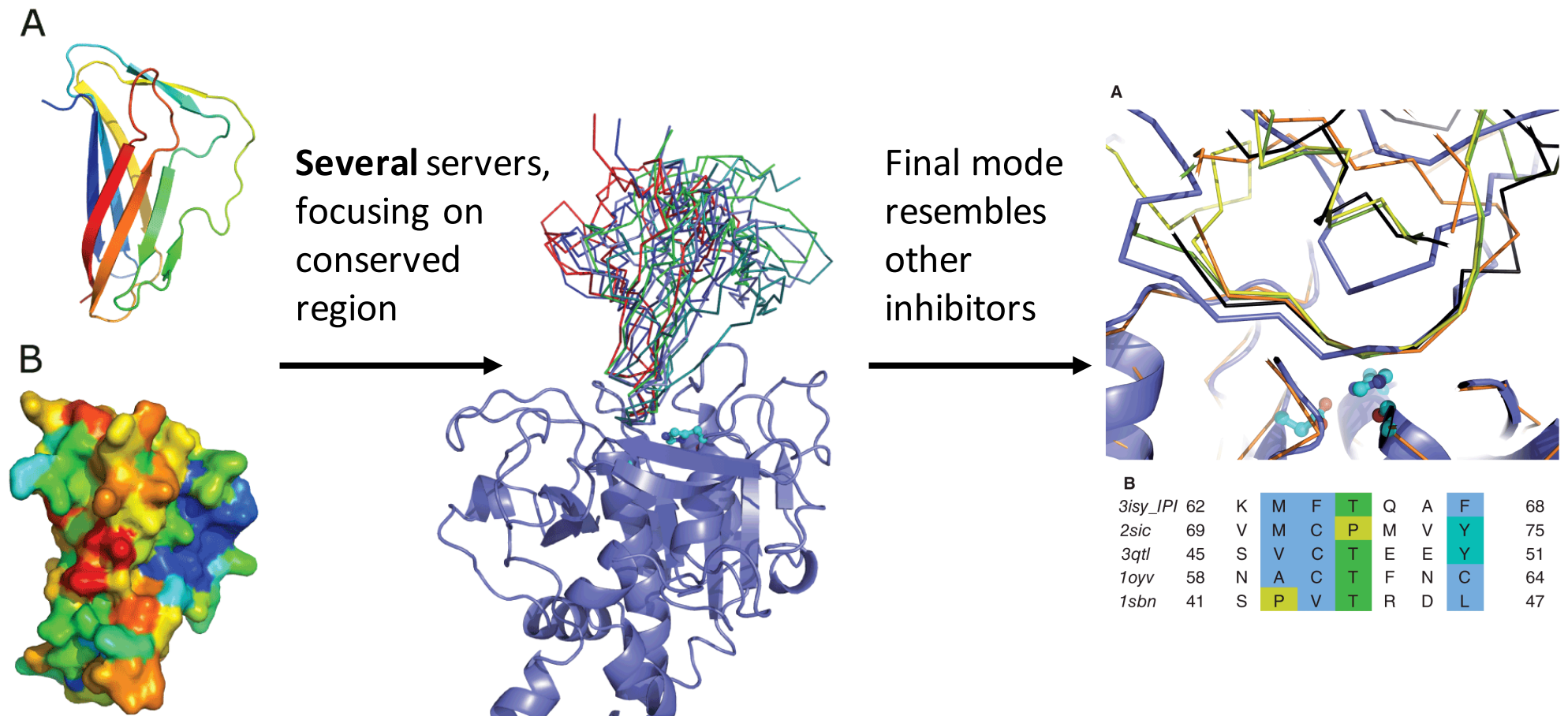


With today's contact predictions, it is now about the best method and works on multiple domains too

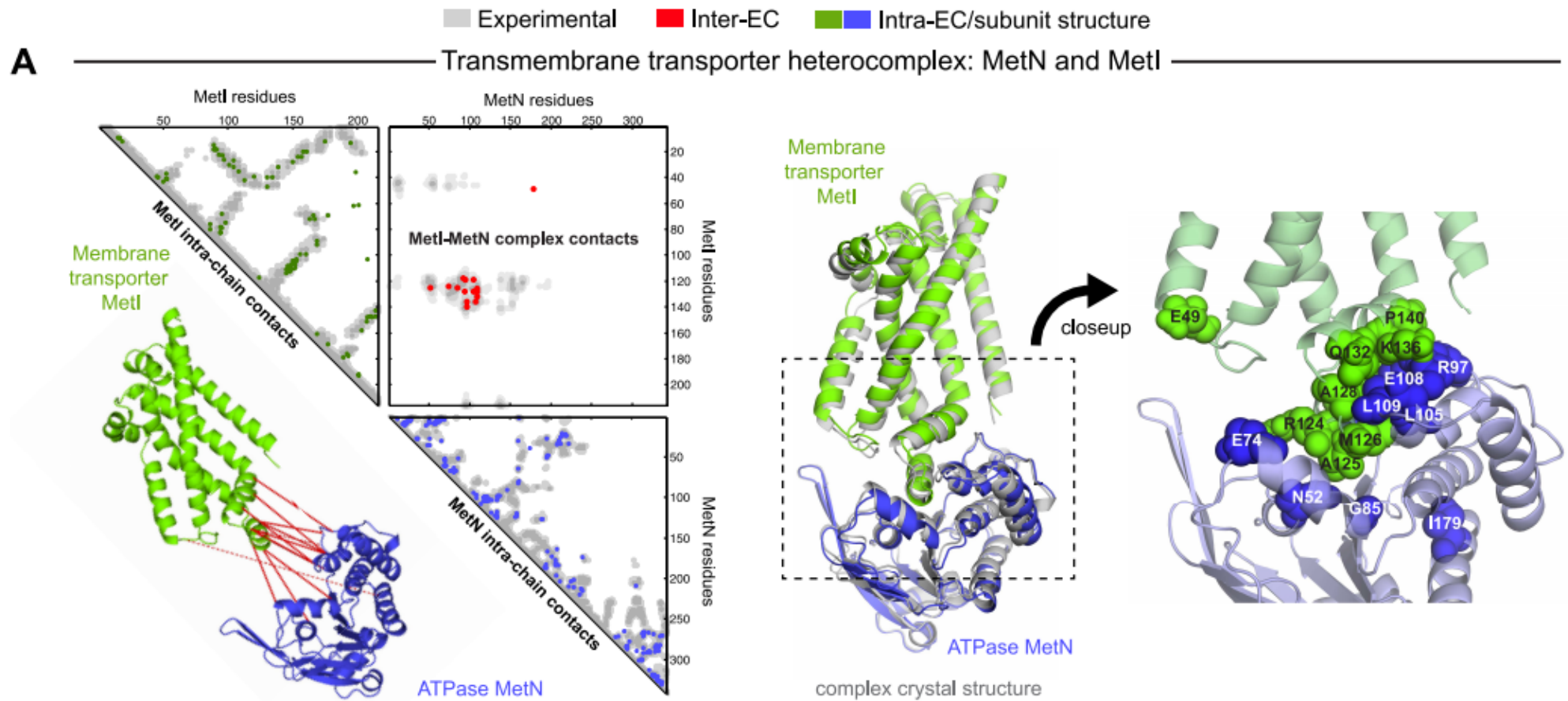
Predicting protein-protein interactions

- Relevant to MR eg proteins A and B are cocrystallised but neither alone solves. An accurately predicted complex, being larger, might solve
- Many methods predict complexes based on steric complementarity plus other scoring functions
- Recommendable servers include
 - ClusPro, the best performing docking method
 - Haddock, which has a good server with different modes
 - Each allows inclusion of other information eg predicted interface residues
 - Symmetric docking at ROSIE server

B. subtilis IPI docking to protease



Using predicted contacts to help predict complexes



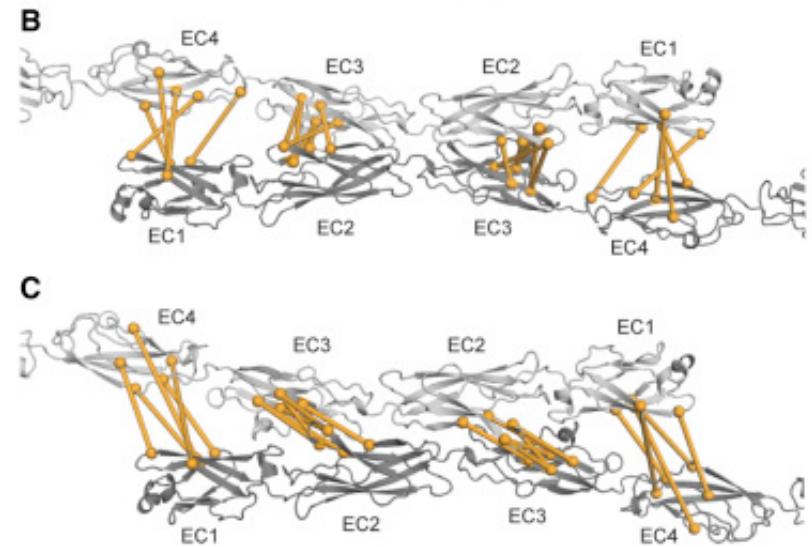
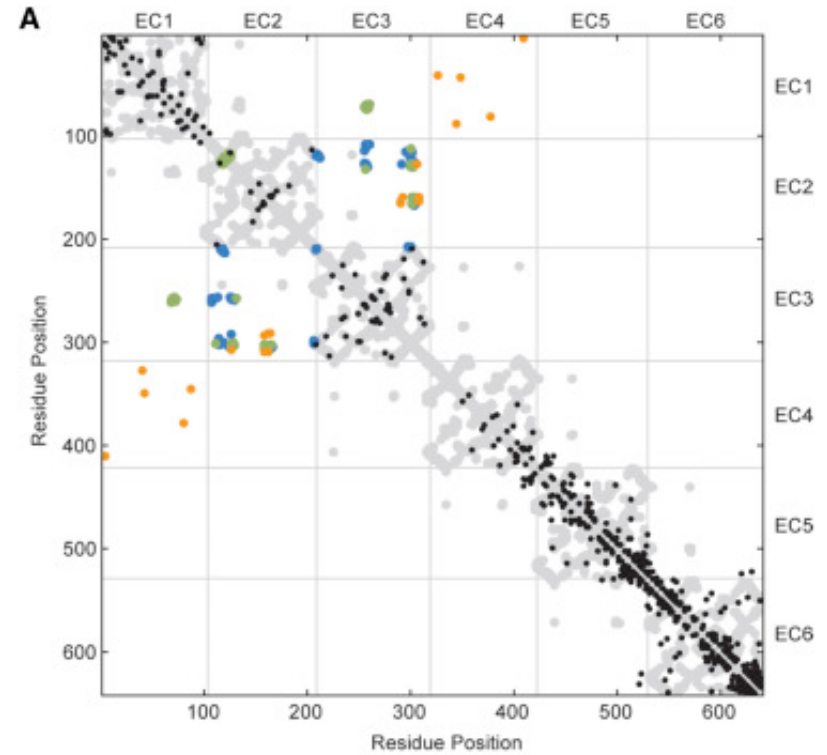
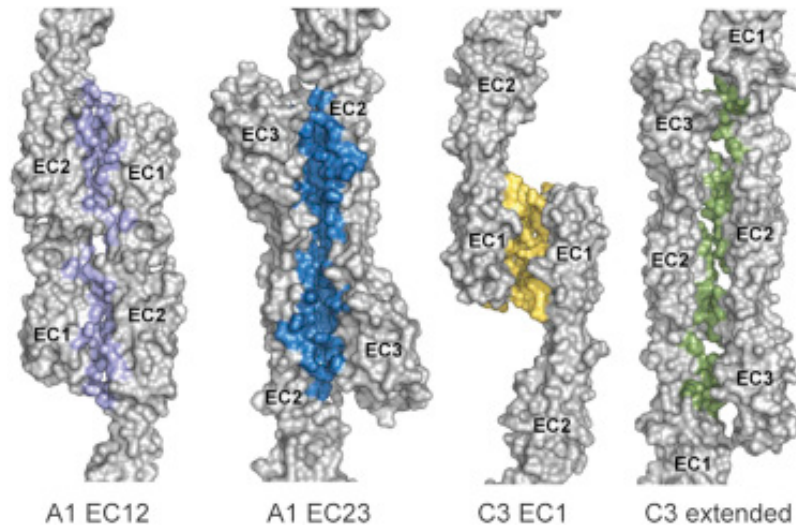
... and once you have your
crystal structure...

What is the biologically relevant quaternary structure?

Where are the functional/catalytic sites?

Validating crystal structure contents

- PISA is an excellent general method, but contact predictions help in some cases
- Crystal showed various ways in which protocadherins could interact
- Predicted contacts supported two of the four modes

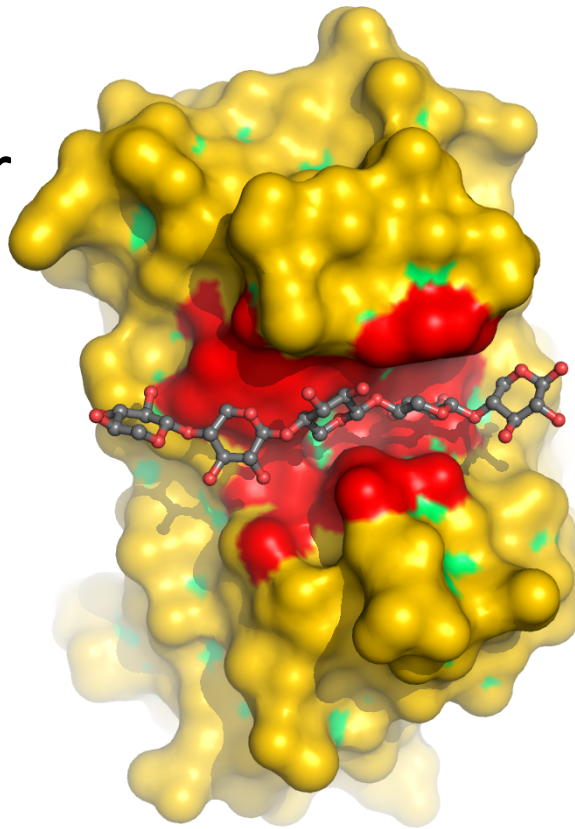


Some lesser-known structure-based function annotation methods

- Finding functional sites is based on their being different somehow to the rest of the protein surface. Important generic methods are based on
 - Shape
 - Electrostatics
 - Physico-chemical characteristics
 - Evolutionary conservation (Consurf)
- Less well-known but valuable characteristics are
 - Statistics of surface atom 'triangles' (STP)
 - Probe interaction energetics (ISMBlab)
 - Rigidity and geometry (EXIA2)
 - Predicted pKa values (THEMATICS/POOL)

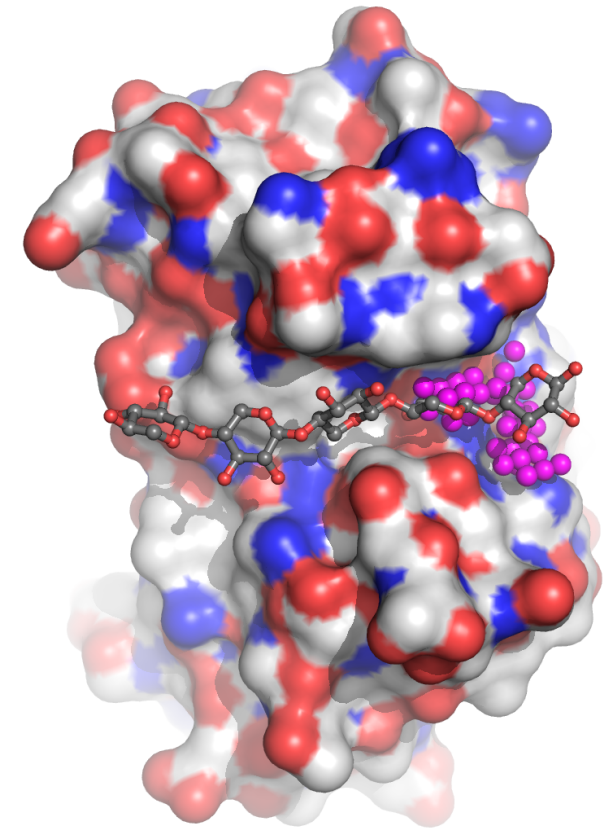
Different probes for different binding sites

- Hydroxyl group can be used to probe for carbohydrate binding sites
- Phosphate oxygen used for binding sites of phosphorylated ligands



ISMBlab

ismblab.genomics.sinica.edu.tw

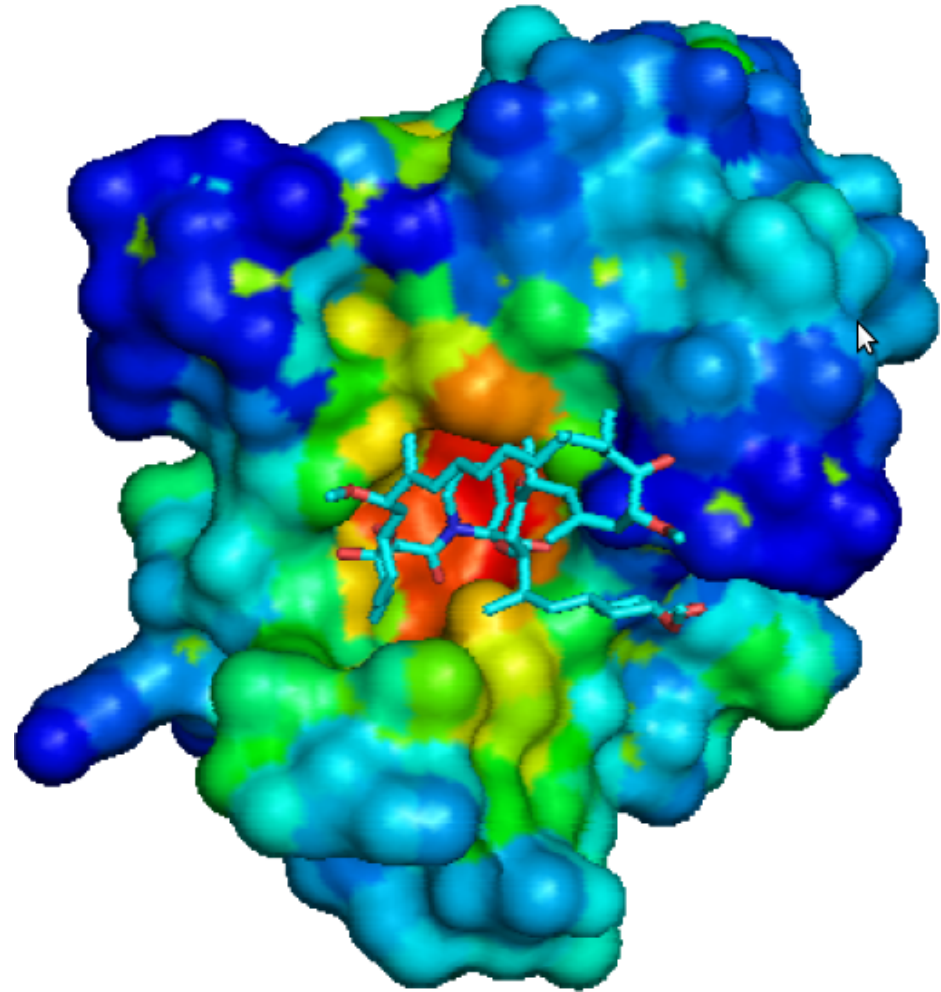


SiteHound

scbx.mssm.edu/sitehound

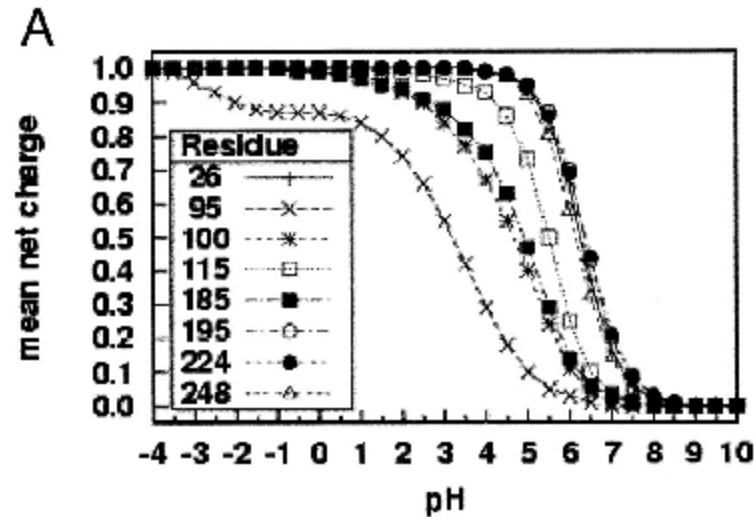
Binding sites from statistics

- STP (surface triplet propensities)
- 13 atom types \rightarrow 455 triplets
- Distribution in binding vs non-binding sites varies
- Designed for small molecules, works on PPIs, including flat surfaces

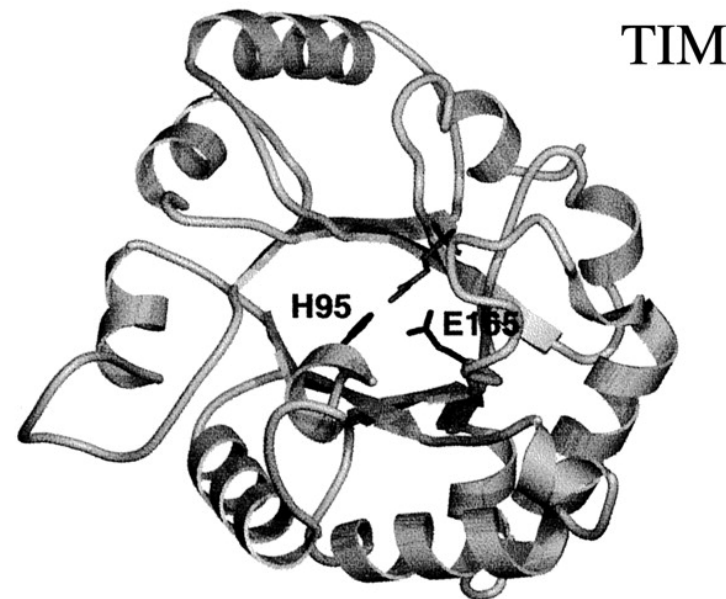


Theoretical microscopic titration

- Computer analysis of a reliable protein structure can predict pKa values for acids and bases. Residues with perturbed pKa values are possible catalytic residues, especially if clustered.



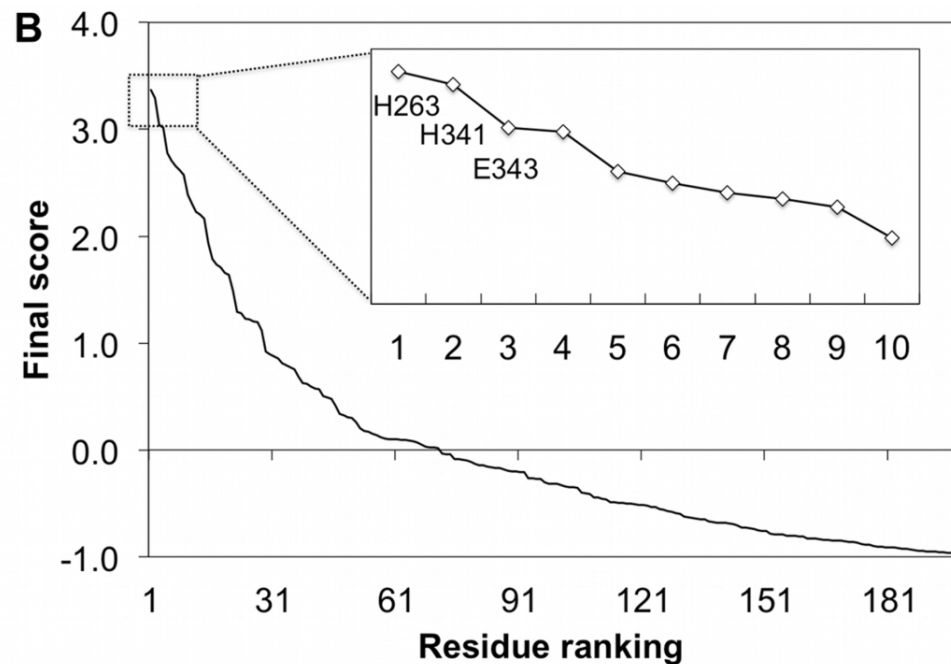
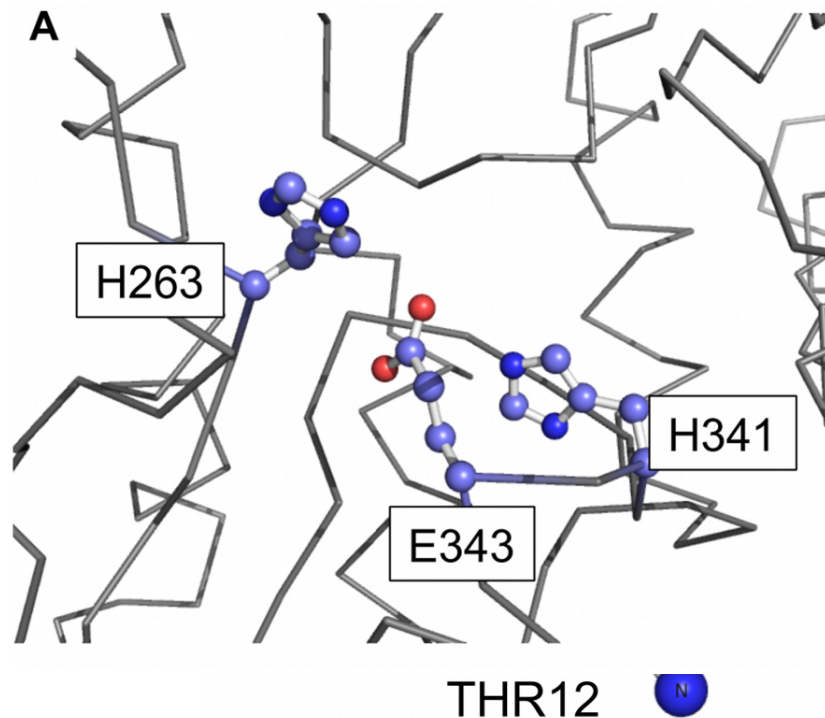
pKa of His95 is atypical compared to other His residues in enzyme



His95 and other residues with atypical pKa cluster at catalytic site

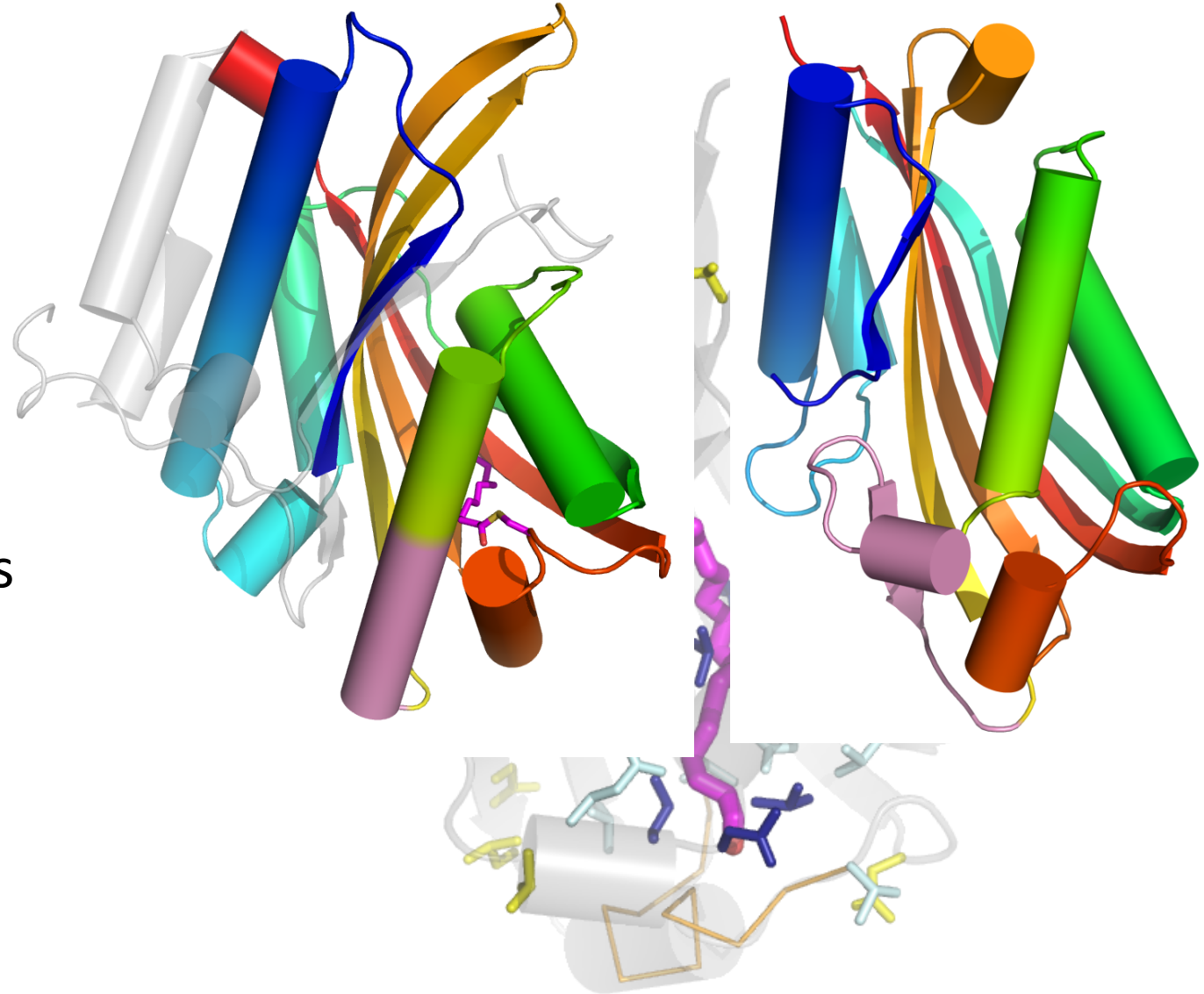
EXIA2: side chain orientation and rigidity

- Find points on the surface with many side chains 'pointing at them'
- Weights these further according to predicted rigidity (measured as number of contacts). Catalytic residues tend to be more packed and so more rigid.



Multiple methods in bioinformatics: Structure comparisons of Evf

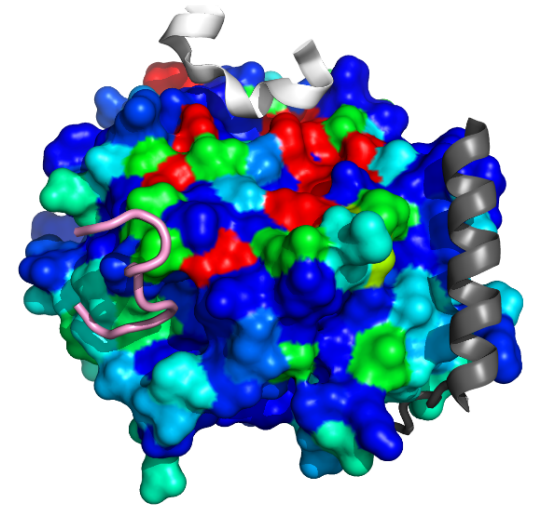
- Reported as novel fold...
- ... but in fact related to *Bacillus* toxin structures
- Both bind to host insect membranes
- Palmitate seen in Evf structure. Matches conserved region of toxins...



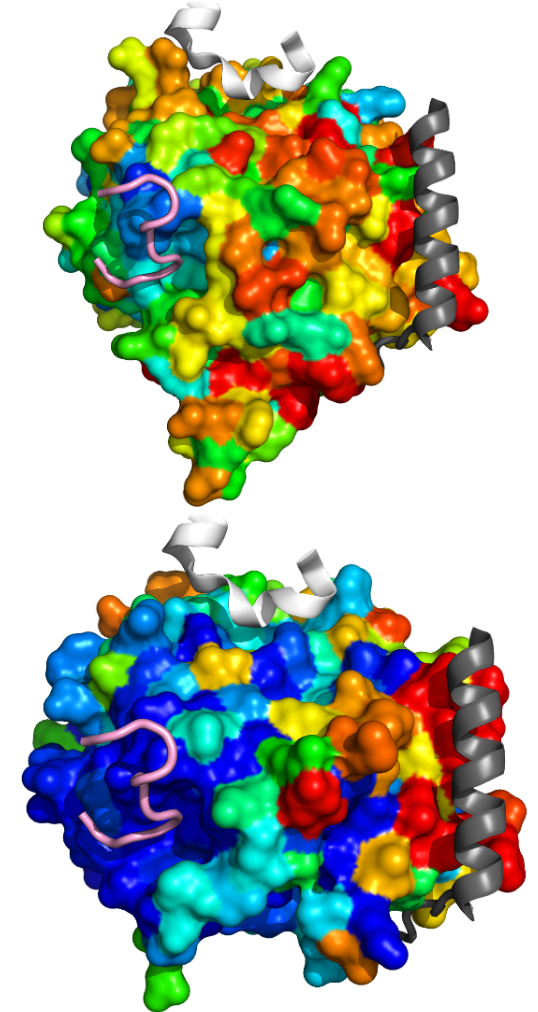
Some servers require thought...

- Consurf maps sequence conservation onto a structure revealing functional sites
- Excellent, general method, but results depend on sequence set chosen for mapping: selecting all or only near relatives gives different results. Either might be more appropriate for you

Mapping 300 homologues mixes different activities so no information on binding sites



But restricting to a single protein family shows only 'pink' site is function in both Diptera and Lepidoptera



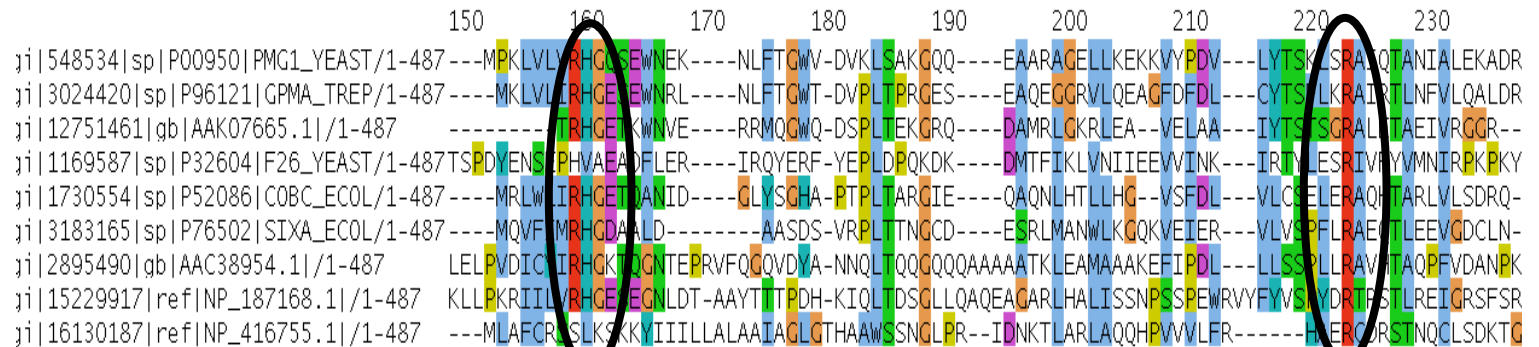
...and finally, you're putting a
manuscript together

Calculating and presenting sequence alignments

Your sequence alignment

- Don't use ClustalW! It's 22 years old! Modern methods like MUSCLE, Probcons and MAFFT are much better

ClustalW misses relatively obvious RHG motif in some of diverse sequence set...



... but MUSCLE gets it



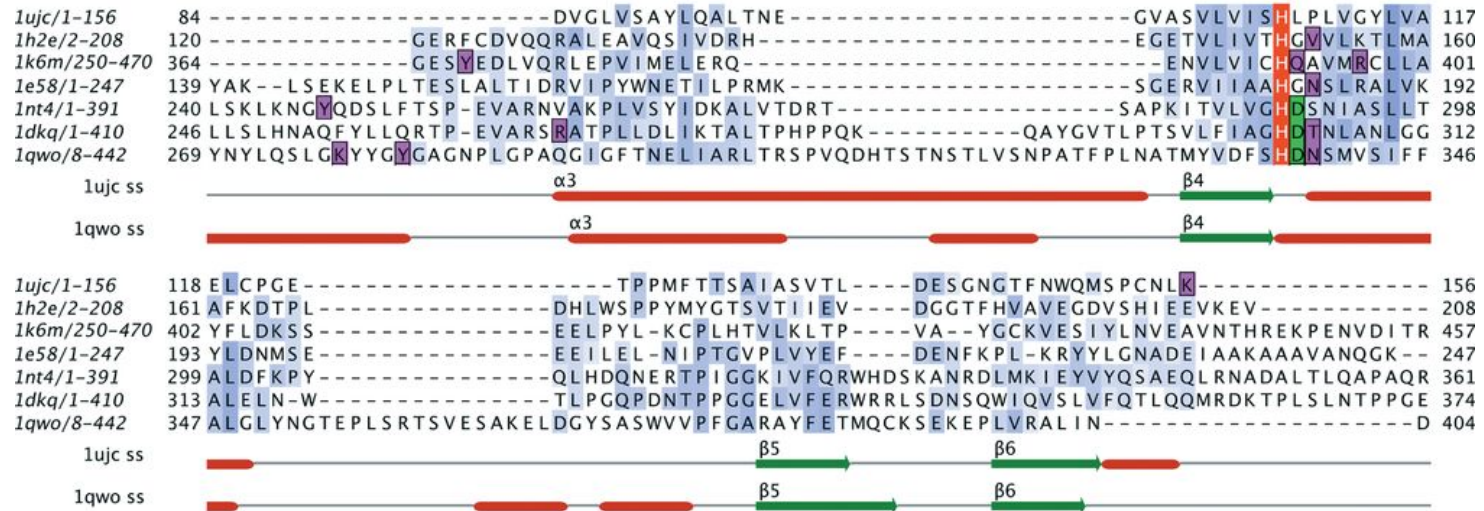
Jalview.org, recommended for sequence alignments

- All these alignment methods and more are available through Jalview on Dundee servers

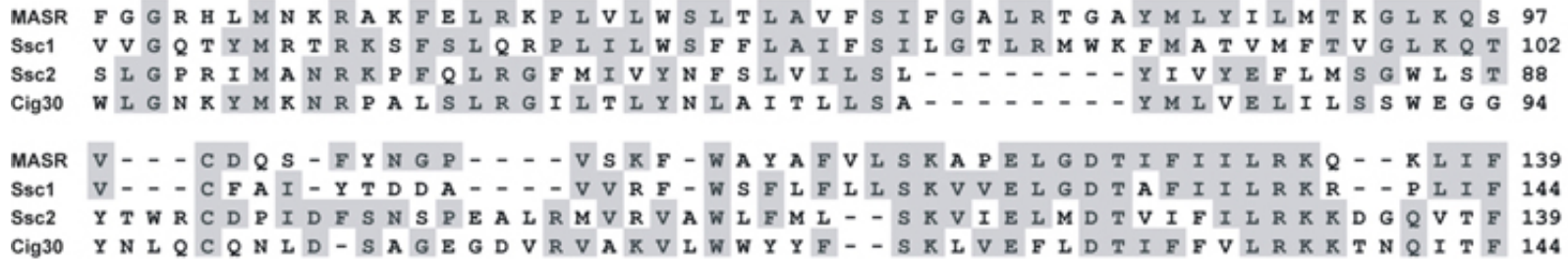
The screenshot displays the Jalview 2.9.0b2 interface. The main window, titled "MAFFT Multiple Sequence Alignment of Retrieved from Uniprot", shows a multiple sequence alignment of FER proteins. The alignment is color-coded by amino acid type. Below the alignment, there are tracks for Secondary Structure, Iron Sulphur Contacts, Conservation, Quality, and Consensus. A menu is open over the alignment, listing various alignment methods: Tcoffee, Probcons, Muscle (selected), MAfft, MSAProbs, and GLProbs. To the right, a "File View" window shows a phylogenetic tree with branches labeled with protein IDs like FER1_PEA, Q7XAG6_TRIPR, FER1_SOLLC, FER1_CAPAA, FER1_SPIOL, FER1_MESCR, FER1_ARATH, FER3_RAPSA, FER2_ARATH, and FER1_MAIZE. Below the tree, a "3D view for FER1_S..." window shows a 3D ribbon diagram of a protein structure with a yellow sphere labeled "THR 89". The bottom of the screen shows a Windows taskbar with icons for Connection, Citrix Receiver, and a terminal window.

Jalview

- Also helps you produce figures like this...



- ... rather than like this



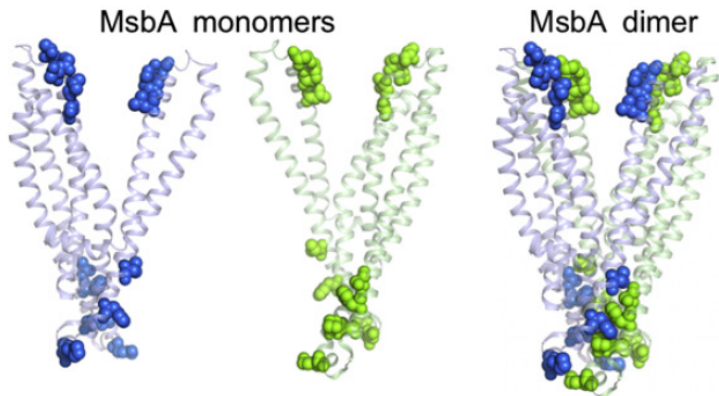
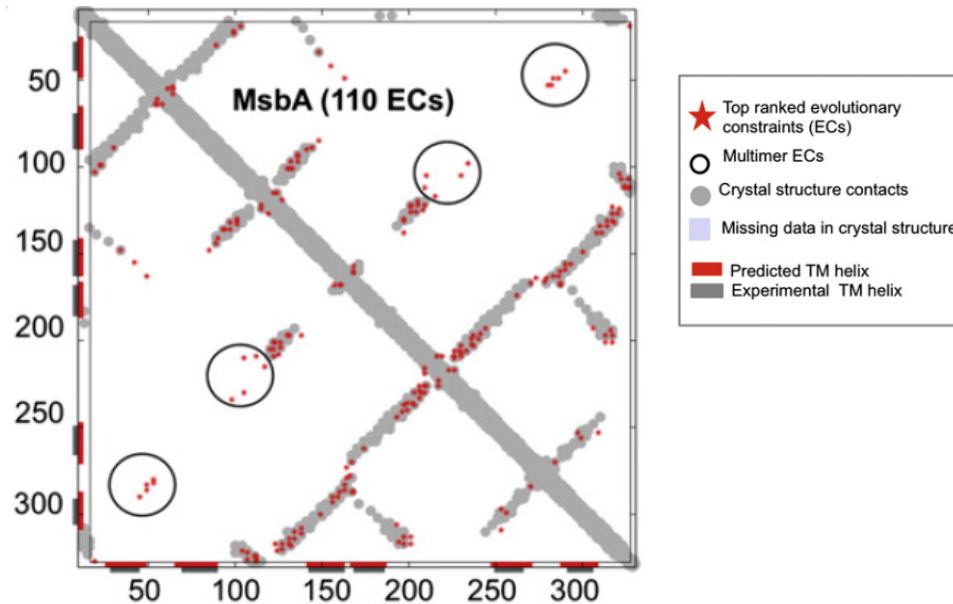
Questions? Feedback?



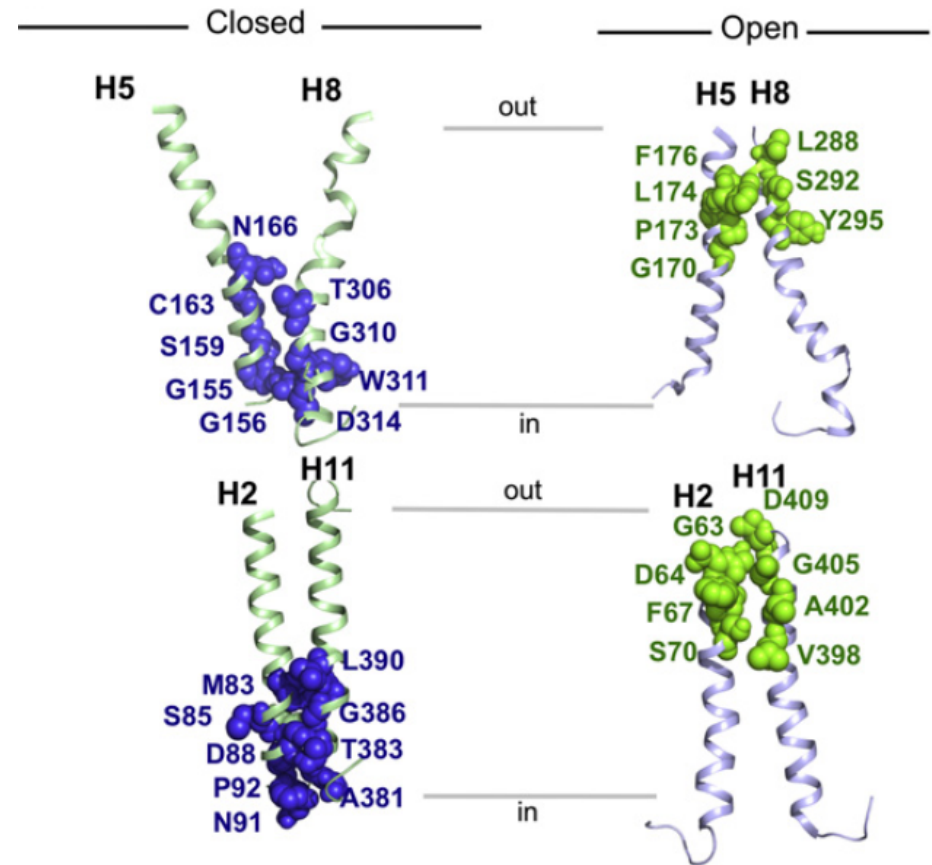
UNIVERSITY OF
LIVERPOOL

Contacts also inform on...

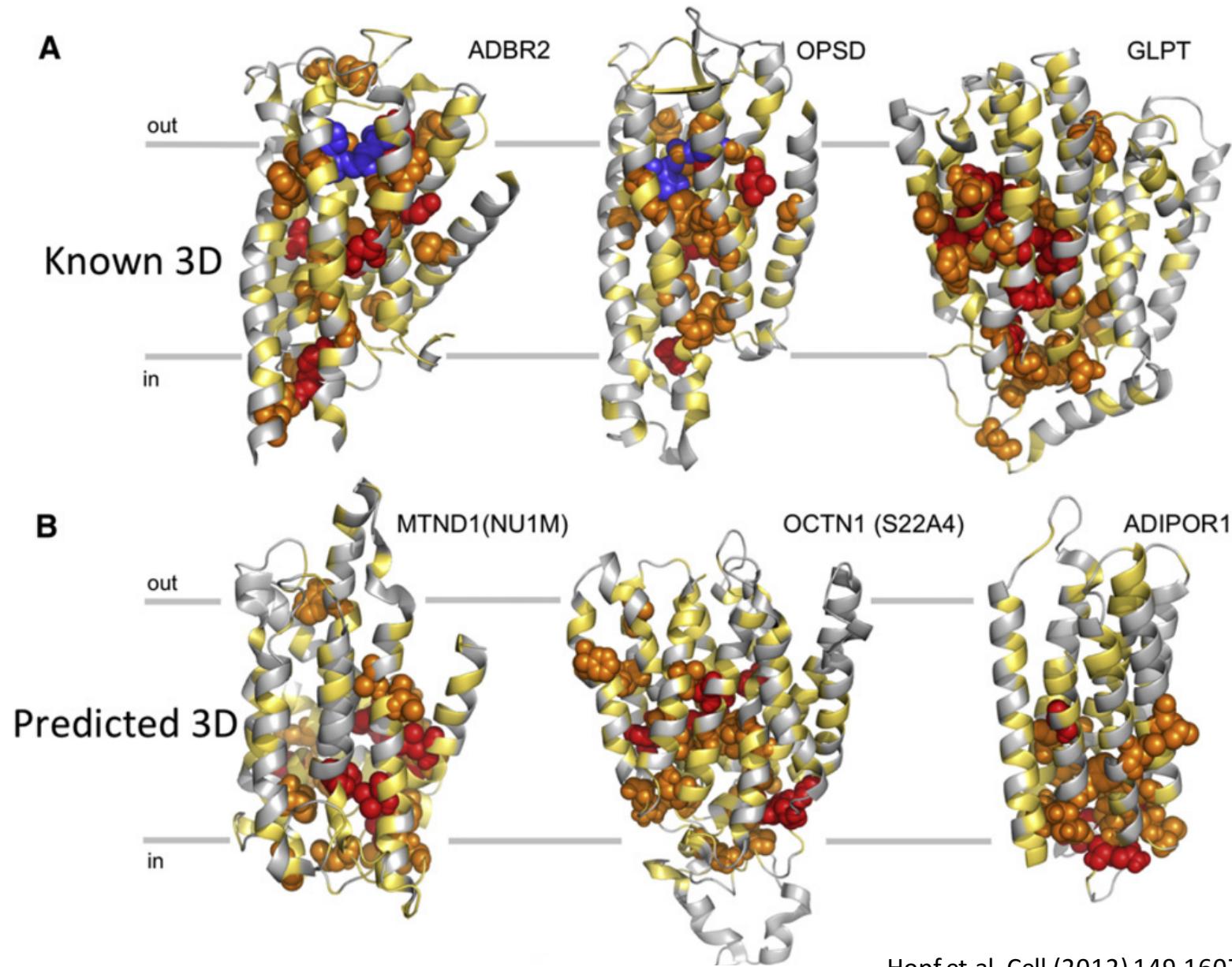
- Oligomerisation



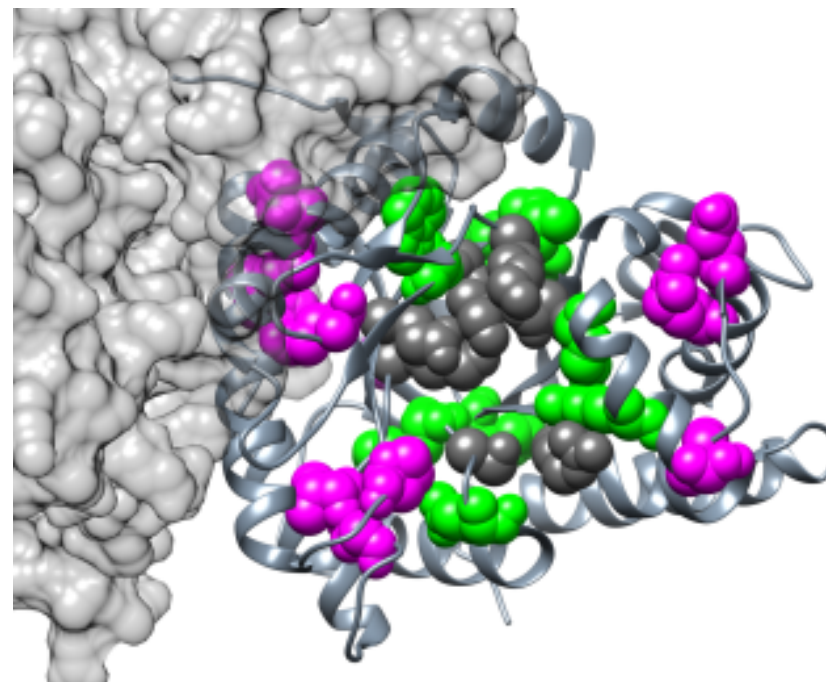
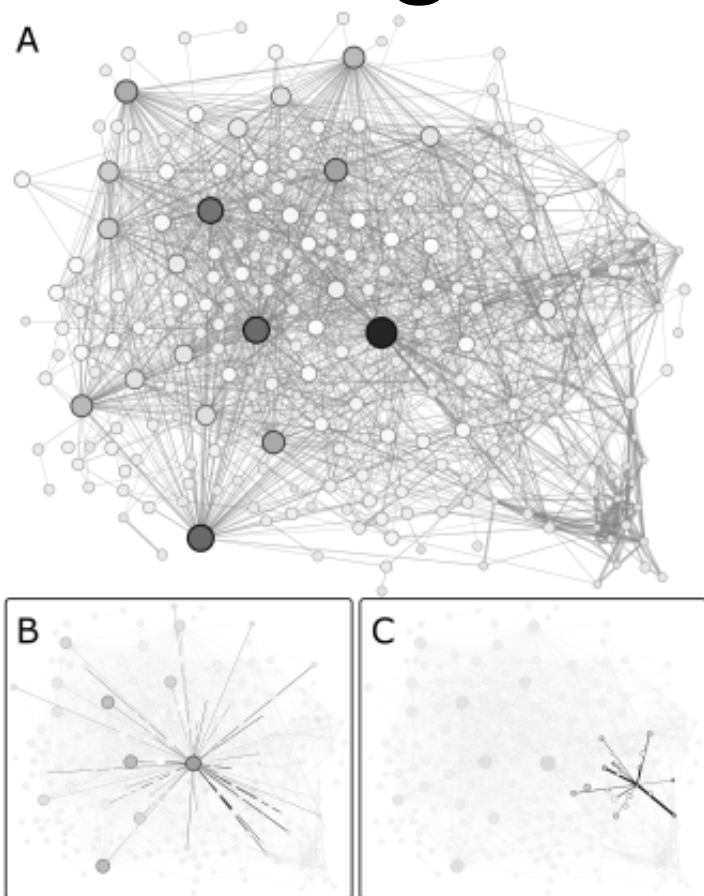
- Conformational change



Predicting functional sites



Predicting functional sites



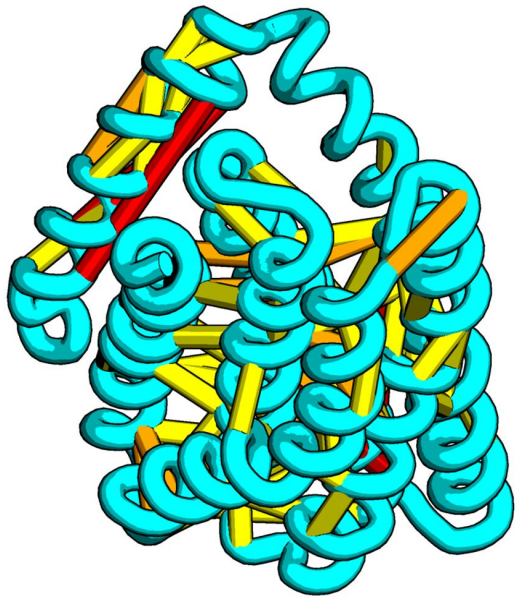
- Considering network structure of contact prediction map gives scores....

- ... that pick out important sites more convincingly
- Grey, conserved catalytic (no covariance signal possible!)
- Green, neighbours of catalytic site; Magenta include interface residues

You don't know the structure

Predicted contacts for folding

A



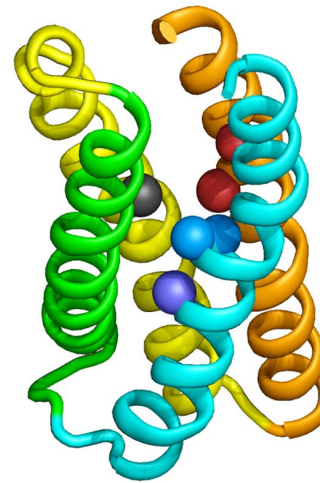
Model showing **satisfied** and **unsatisfied** contacts

B

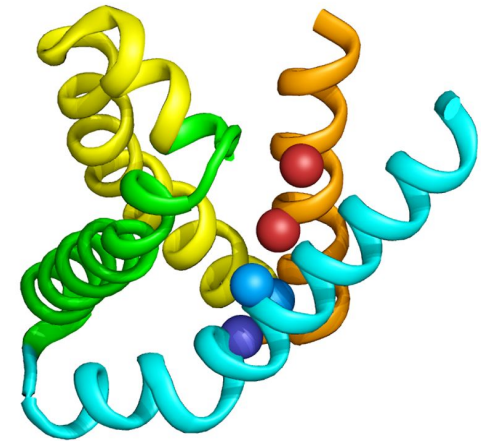


Substructure of model with conserved motifs

C



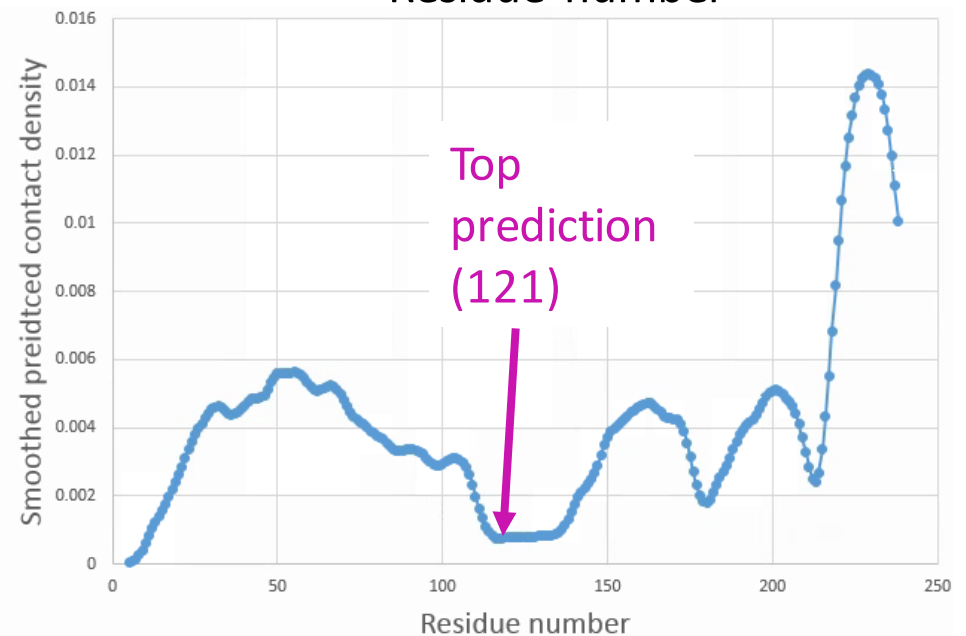
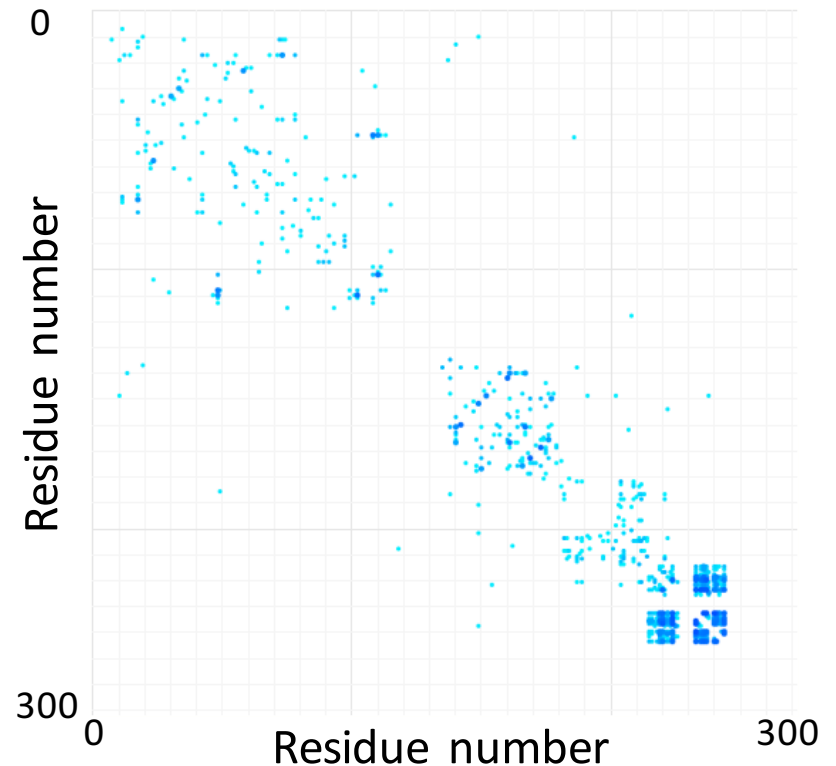
D



Crystal structure with similar conserved motifs

An example from the original paper

- AL1 geminivirus protein of ~250 residues
- “... LM1 of this profile lay at residue **132** ... the depth of LM1 corresponds to an average error of approximately 19 residues. ... the domain definition agrees very well with the functionally defined AL1 origin DNA-binding site domain from residues **1–116**.”
- Domain now structurally defined as ~7-118



Bacterial competence and ComEC



MBiol Project

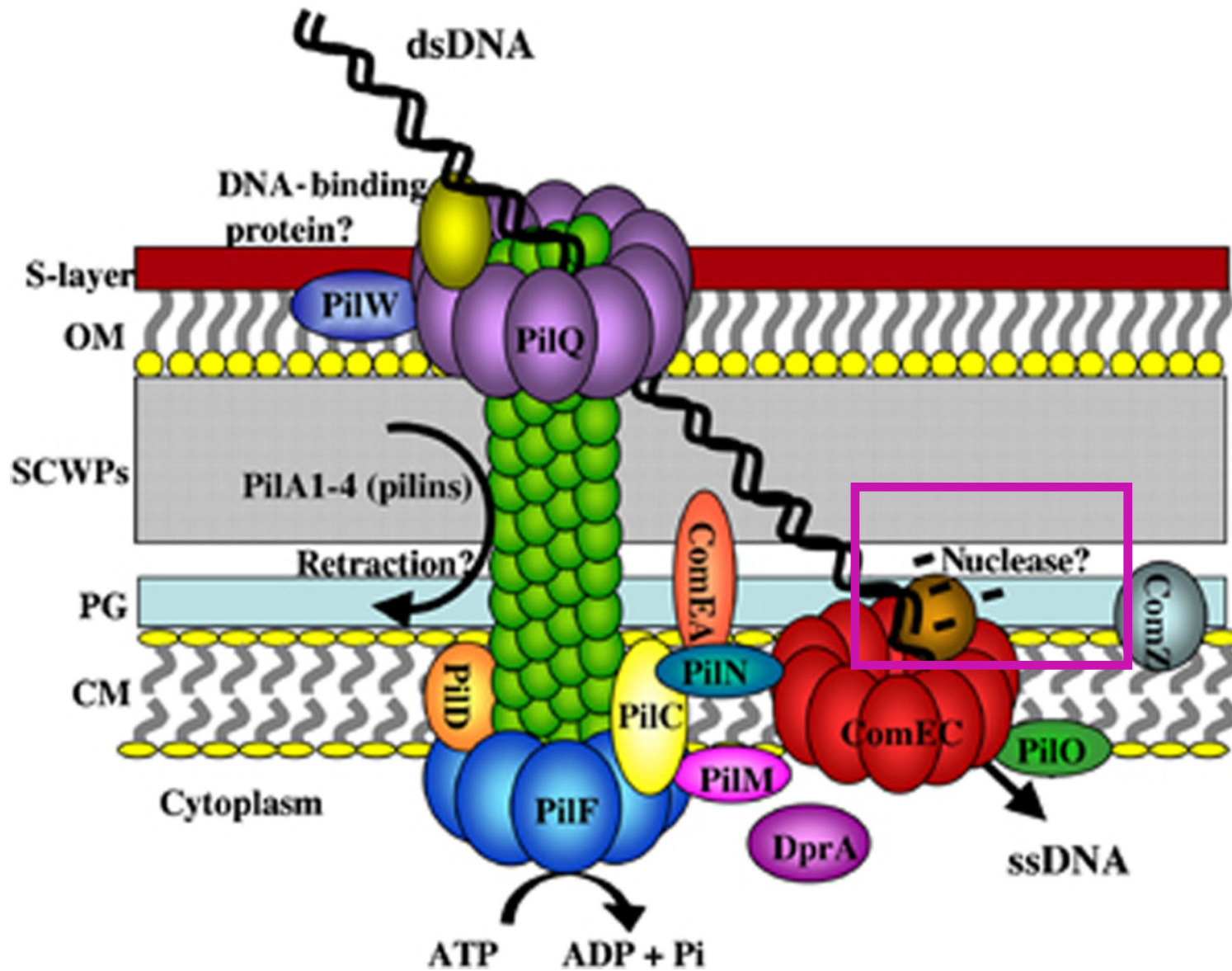


BBSRC Research
Experience Placement

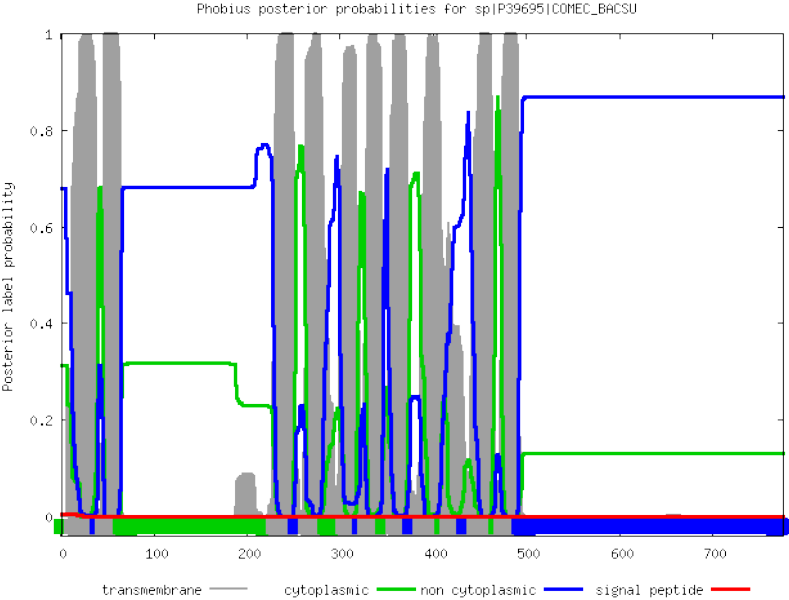
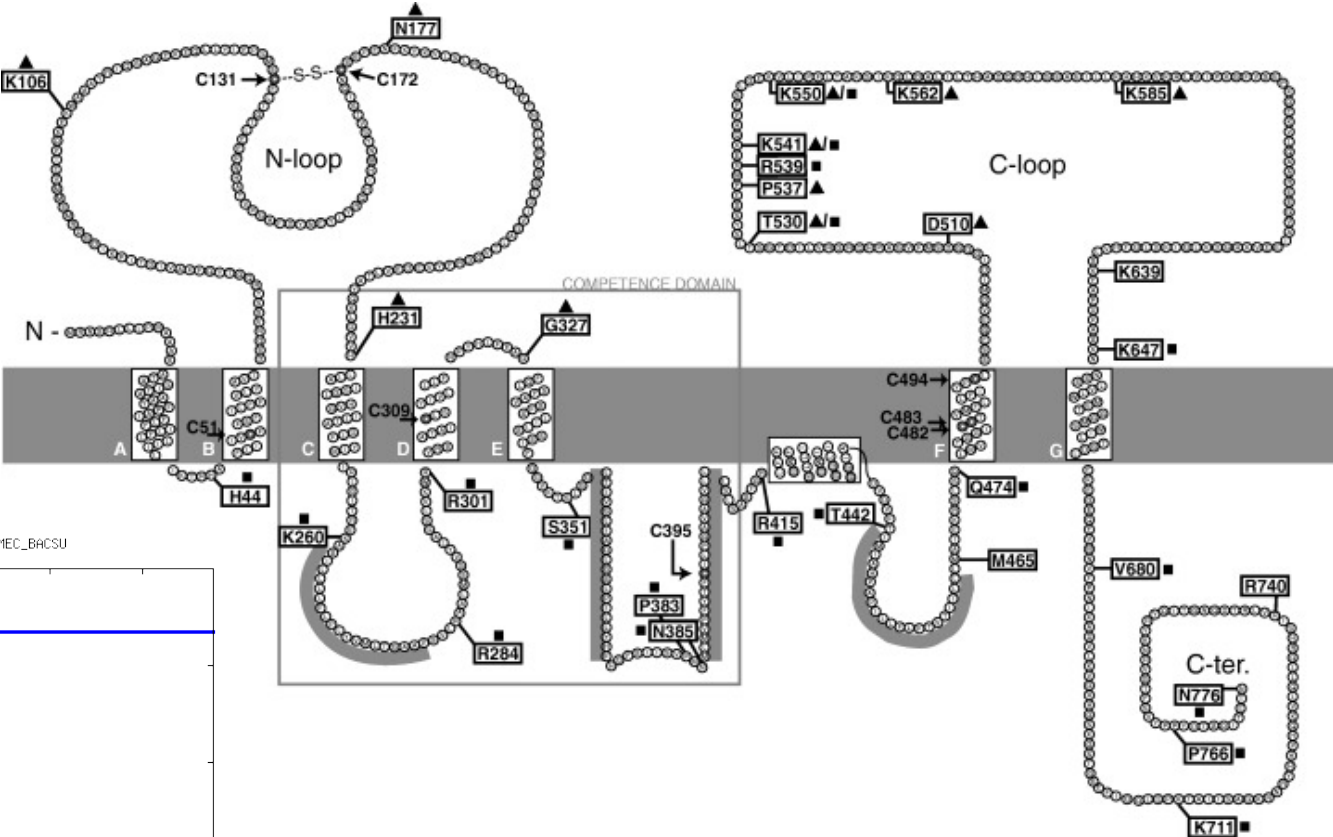
Bacterial competence

- The innate ability of some bacteria to take up extracellular DNA
- Proteins involved vary between species eg Gram +ve vs -ve
- Many poorly understood
- Bacterial competence involved in antibiotic resistance spread and pathogenicity of some bugs
- One of the proteins most strongly linked to competence is ComEC

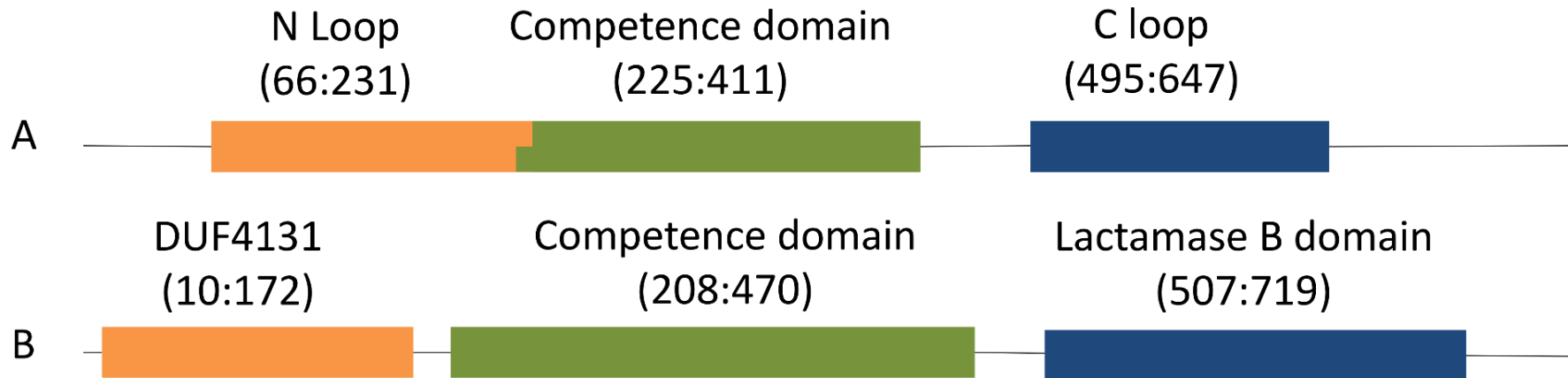
ComEC function



ComEC structure



ComEC domain structure



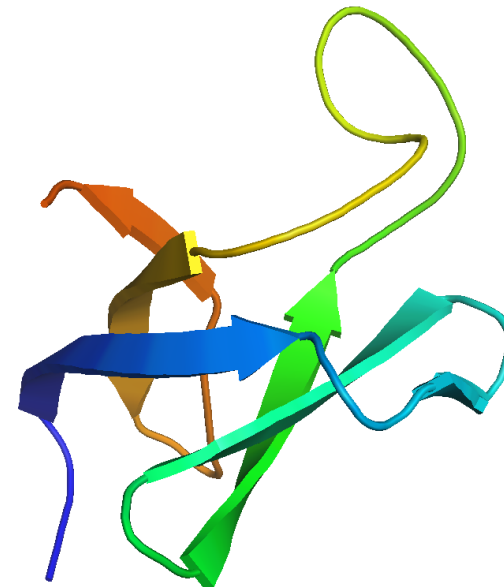
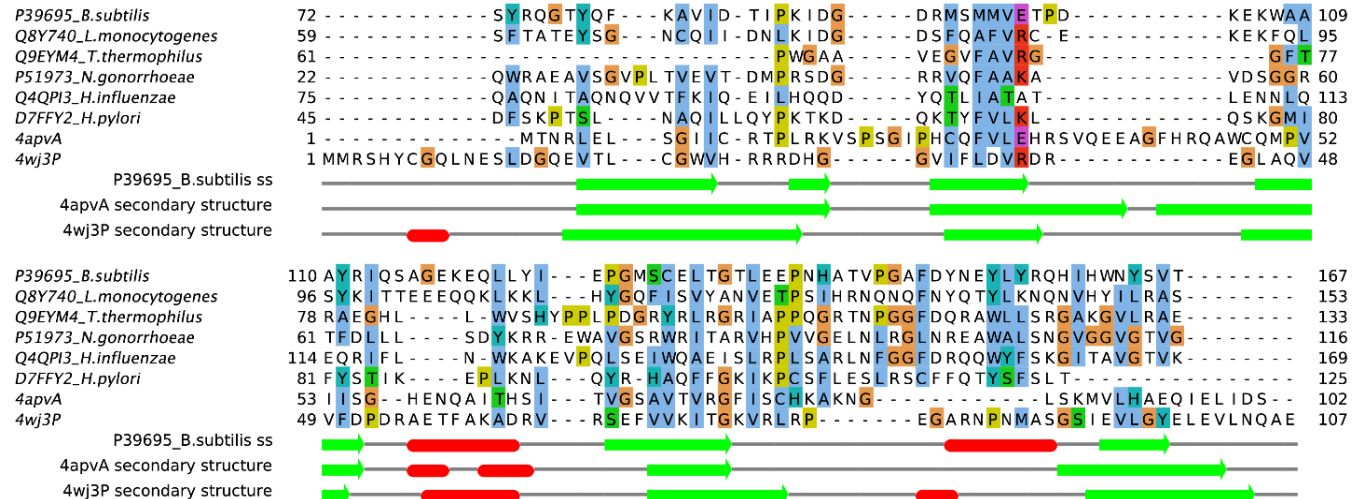
ComEC lactamase-like domain is predicted to be a DNase

- Positively charged
- Predicted as DNA binding by structure-based predictors
- Accommodates DNA duplex neatly
- We think the mystery nuclease activity is encoded within ComEC!

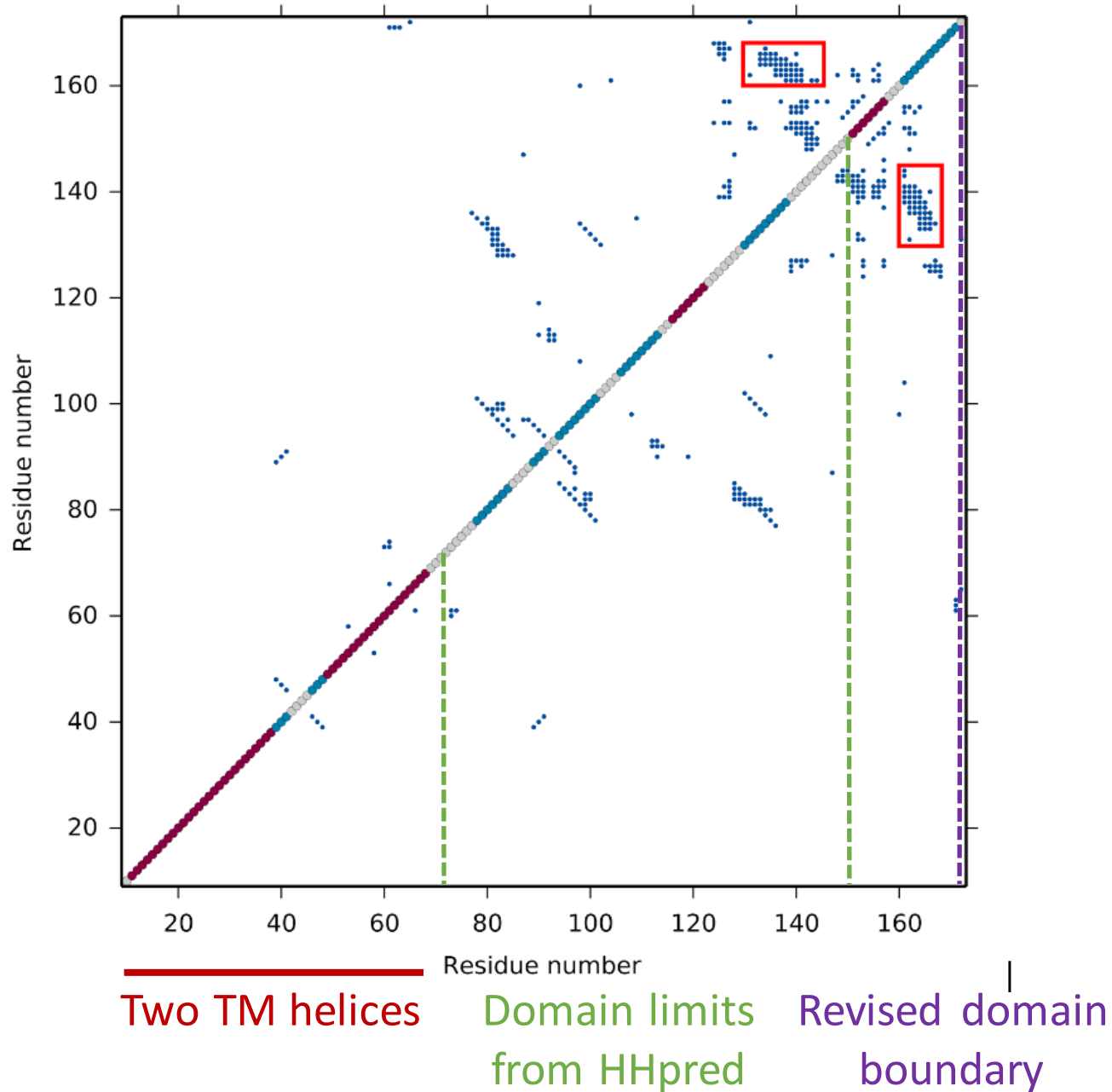


DUF4131 is predicted to be an OB fold

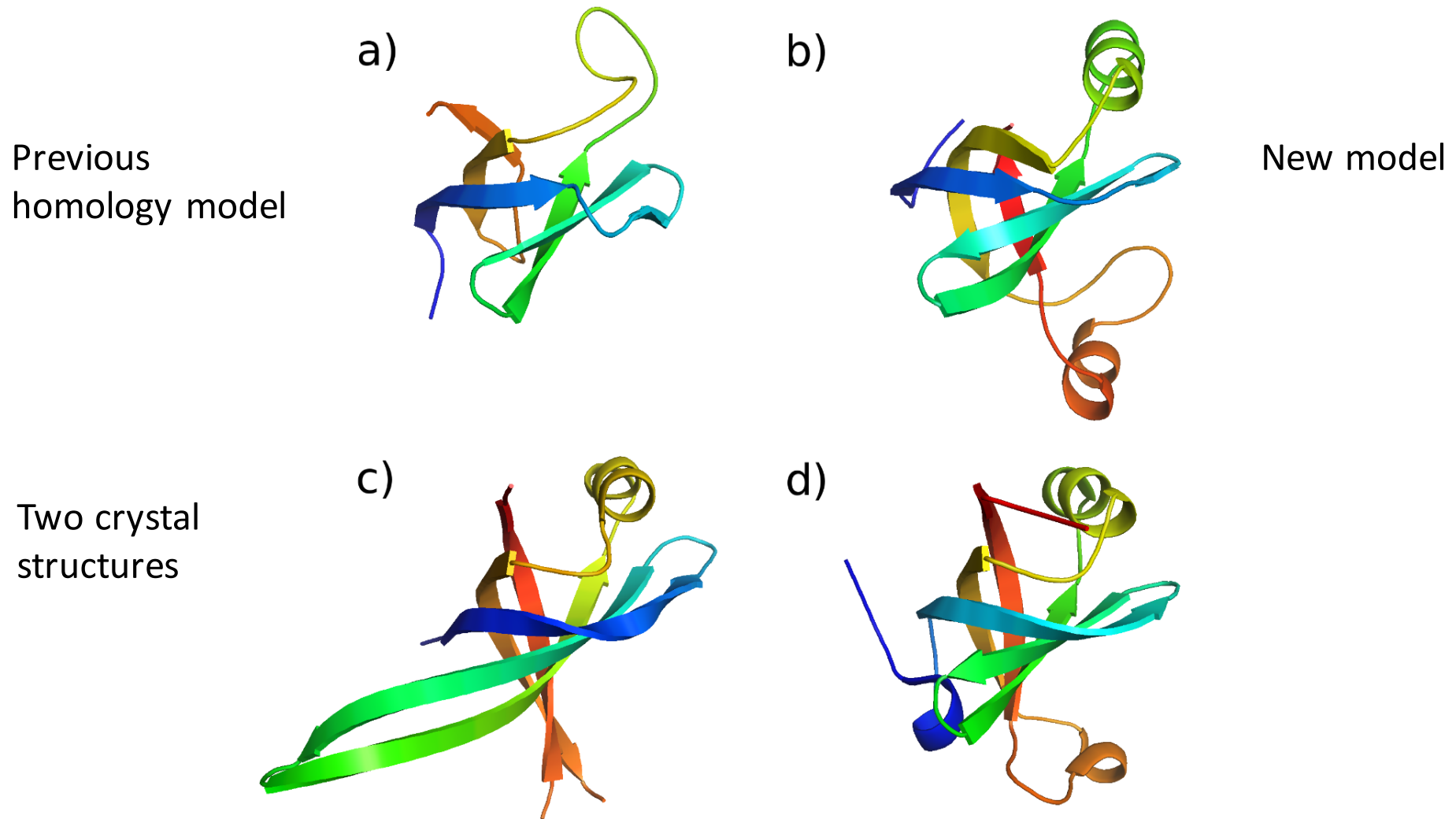
- Pretty clear result by HHpred distant homology detection
- OB folds bind single-stranded nucleic acids or oligosaccharides
- Context suggests former, though structure-based methods do not predict NA binding (trained on dsDNA?)
- But model looks small and not beautifully formed



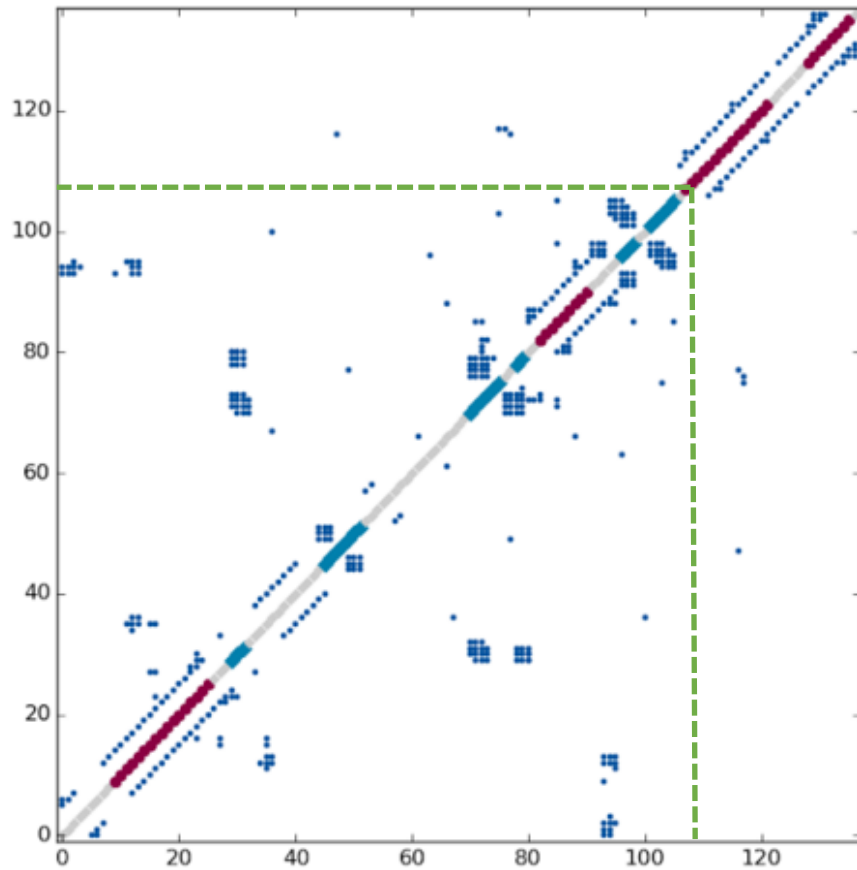
Do we have the right domain boundary?



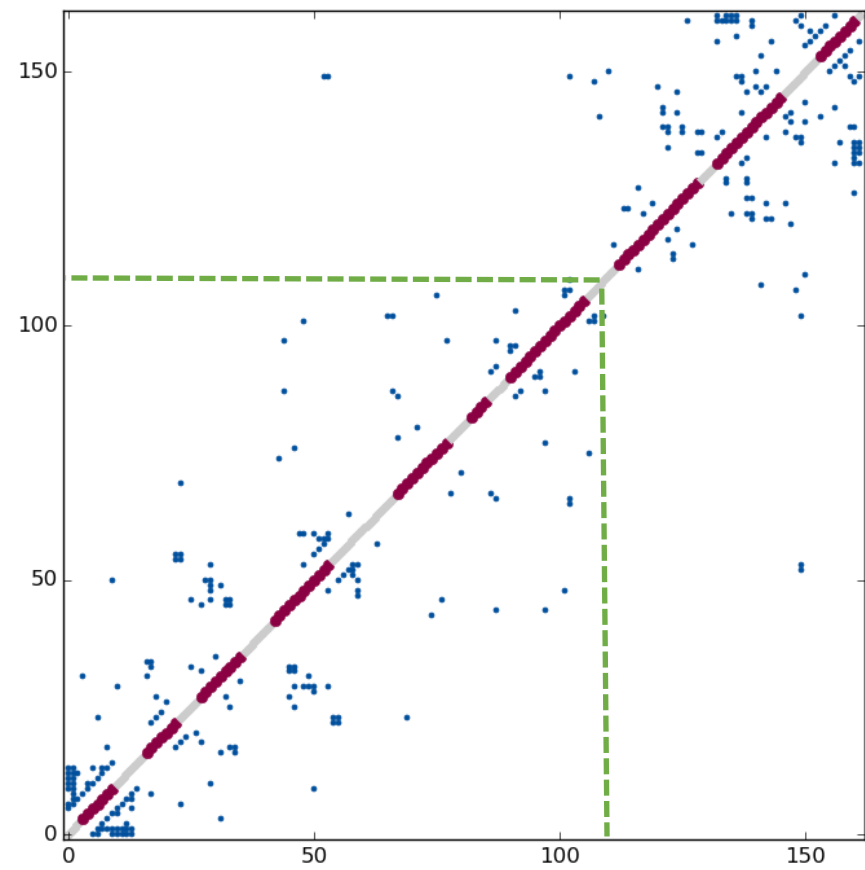
A contact-assisted fragment assembly model from Rosetta looks much better



Many large DUFs seem to contain multiple domains



DUF3670 – separate helices at C-terminus?



DUF4158 – two helical domains?

Limitations and opportunities

- Generally need large number ($>\sim 1000$) of reasonably diverse sequences
- Servers for **single** sequence are available
 - Evcouplings (a day)
 - RaptorX (two days)
 - Gremlin (1-2 hours)
- Only one server for **two** sequences and doesn't work well
- Only one server for **folding** from a sequence, but it's slow and is not the top method
- We have tools for single sequences and folding installed locally (Felix). No local method for two sequences available

Limitations and opportunities

- Get better domain definitions for cloning, bioinformatics
- Predict folds *ab initio* for large families (AMPLE)
- Rank interfaces in crystal structures and docking results

- ? Predict functional sites
- ? Filter true from false positives in Y2H, affinity tagged complexes
 - ??Genome-scale *in silico* interactomes
- ? Supplement incomplete NMR data