# Protein Structure Prediction and Using AlphaFold
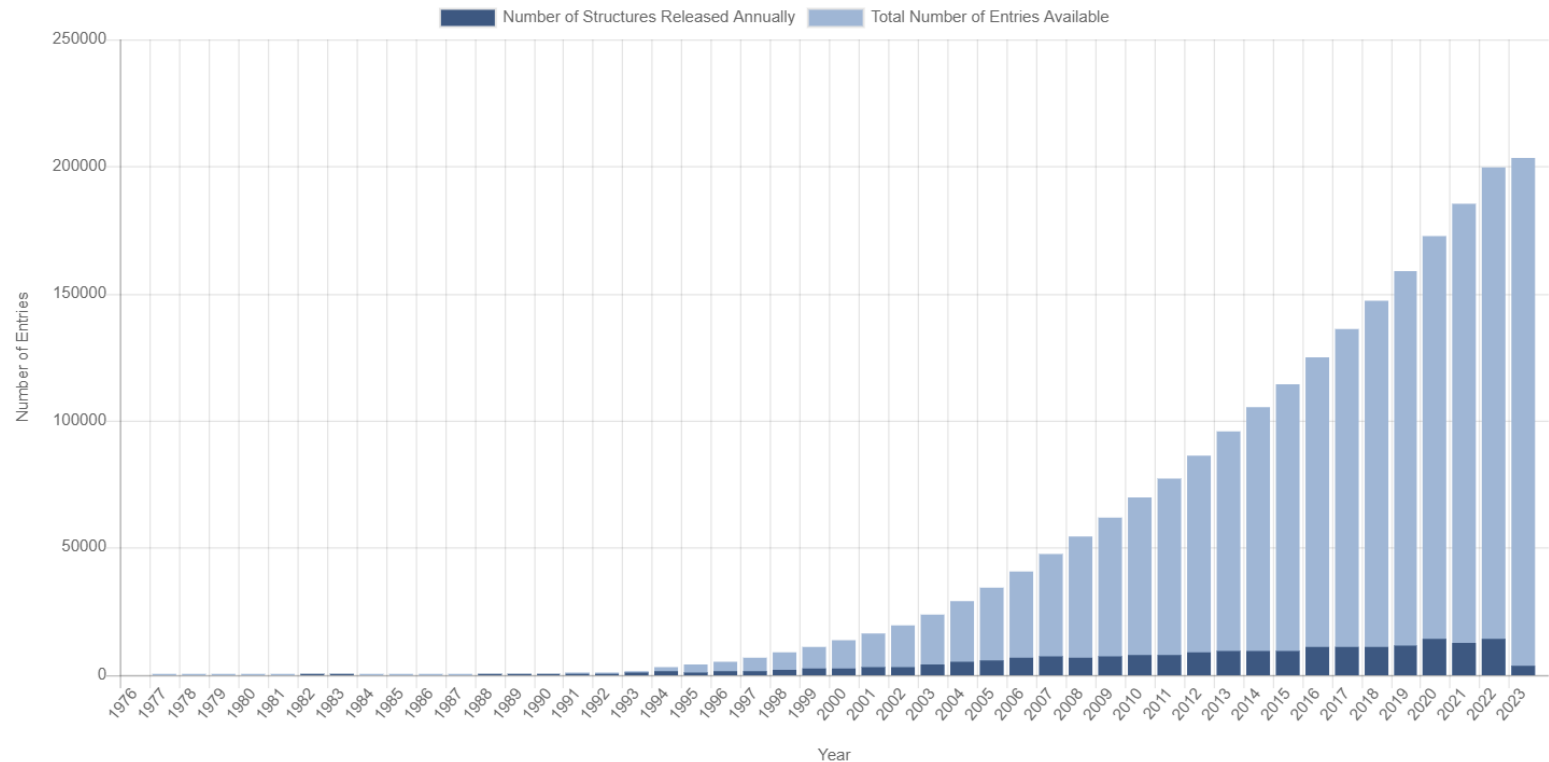
Day 10: Friday 31st March

# 3D Structure Prediction

- Yesterday we learnt about experimental approaches to determine protein structure:
  - NMR
  - X-ray crystallography
  - Electron microscopy (cryo-EM)

- Today we learn about in silico approaches to predict structures:
  - Homology Modelling
  - Artificial Intelligence and AlphaFold
  - What AlphaFold can and can't do (yet)
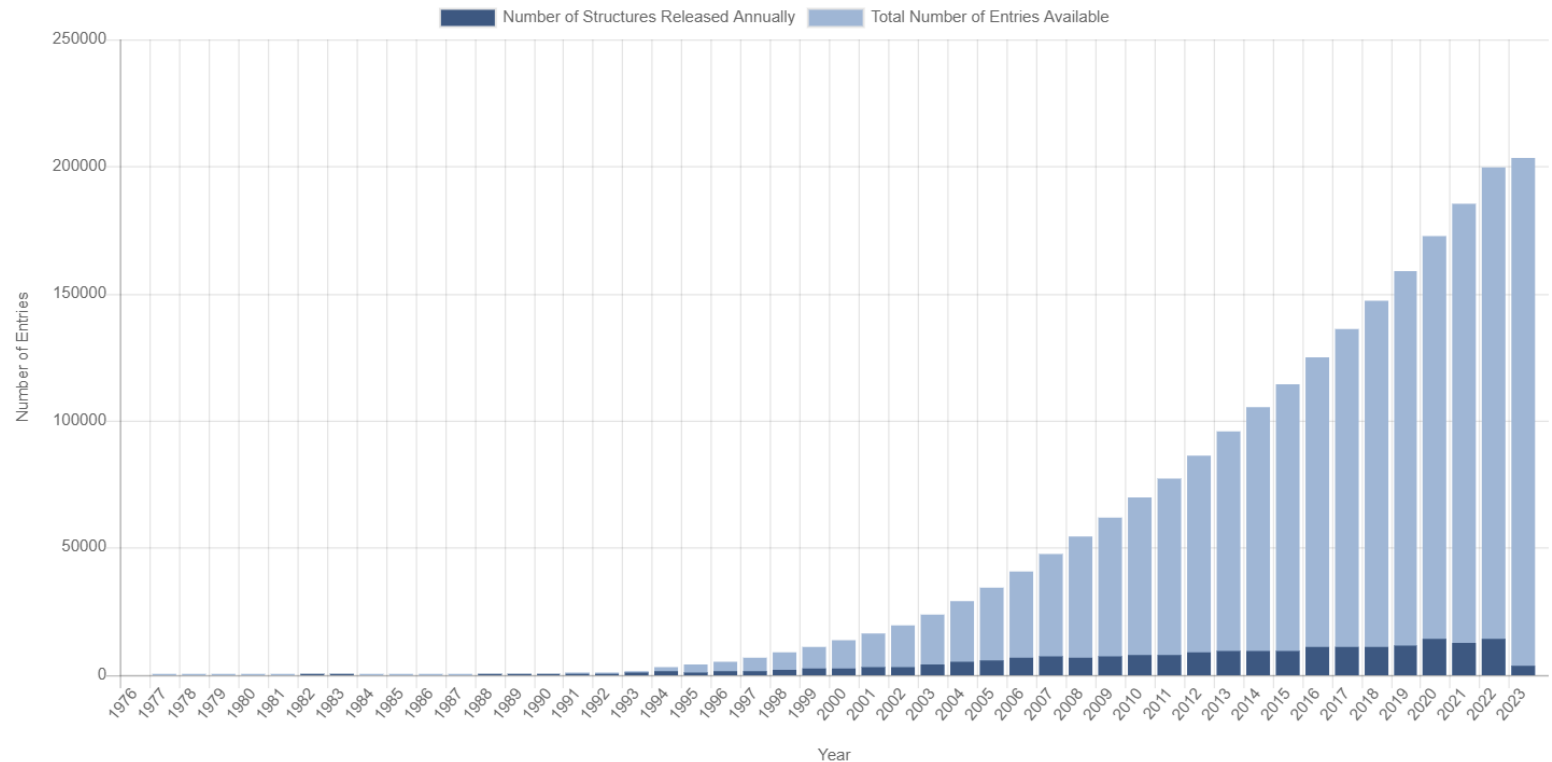
UNIVERSITY OF
CAMBRIDGE

# Experimental Structures in the PDB

- Enormous and growing number of structures that have been experimentally determined

- These are freely available in the online Protein Data Bank and listed in UniProt
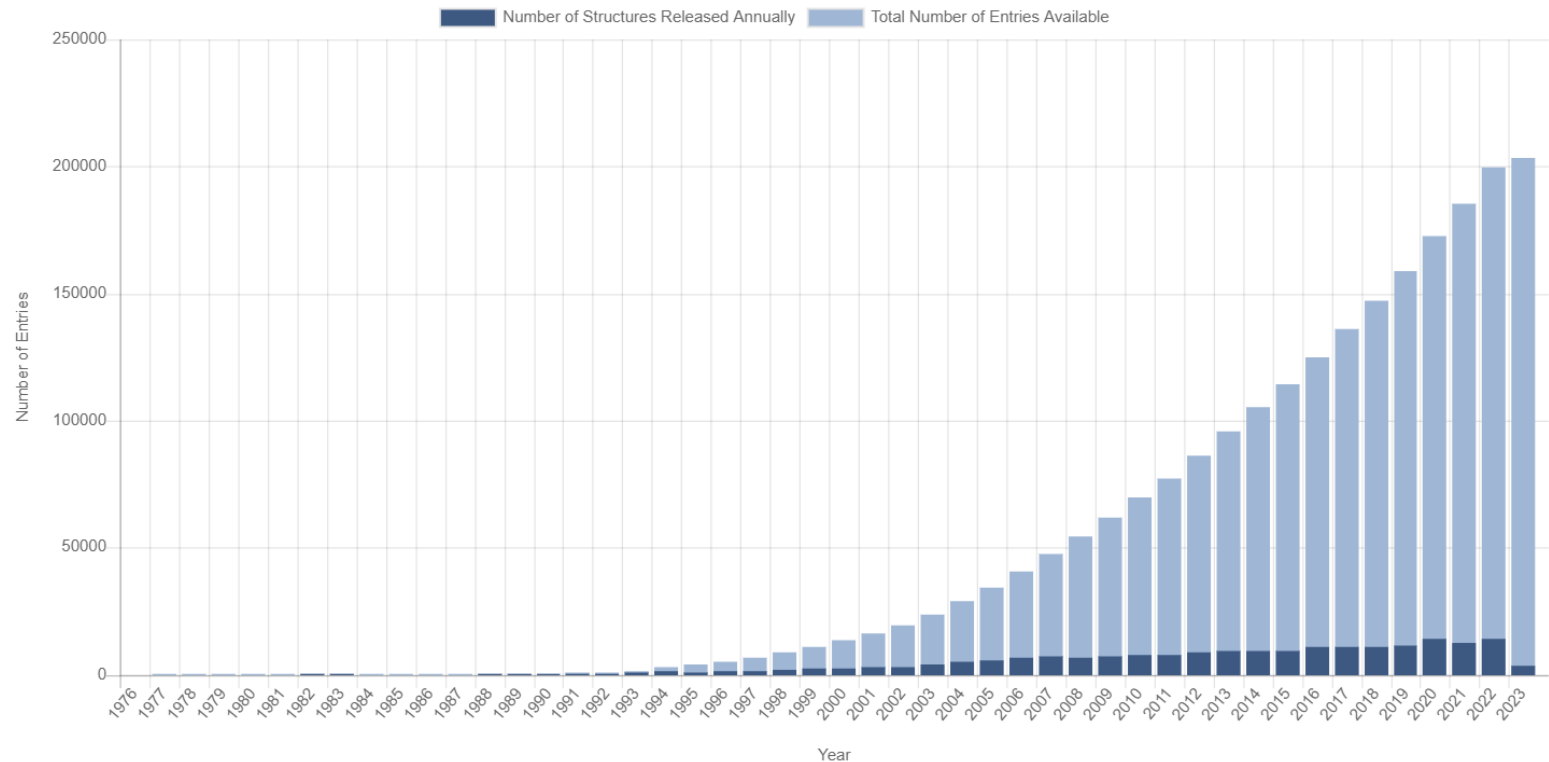
- If an experimental structure exists – USE IT!



Image: www.rcsb.org/stats/growth/growth-released-structures

# Experimental Structures in the PDB

- These experimental structures have been a very rich source of information for structure prediction for decades



Image: www.rcsb.org/stats/growth/growth-released-structures

# Experimental Structures in the PDB

- These experimental structures have been a very rich source of information for structure prediction for decades

- This approach is called Homology Modelling

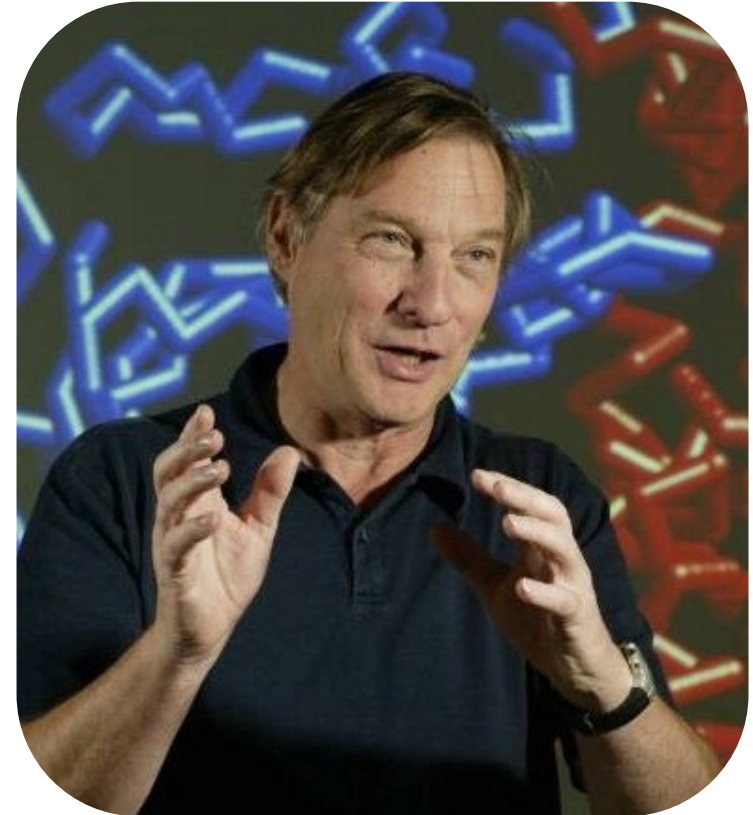Image: www.rcsb.org/stats/growth/growth-released-structures

# Predicting a 3D structure

- Homology modelling has existed for a long time – use a closely related known structure to predict a new one
  - Modeller
  - SwissModel
  - HHPred
  - FFAS
  - SCWRL

# Predicting a 3D structure

- Homology modelling has existed for a long time – use a closely related known structure to predict a new one
    - Modeller
    - SwissModel
    - HHPred
    - FFAS
    - SCWRL

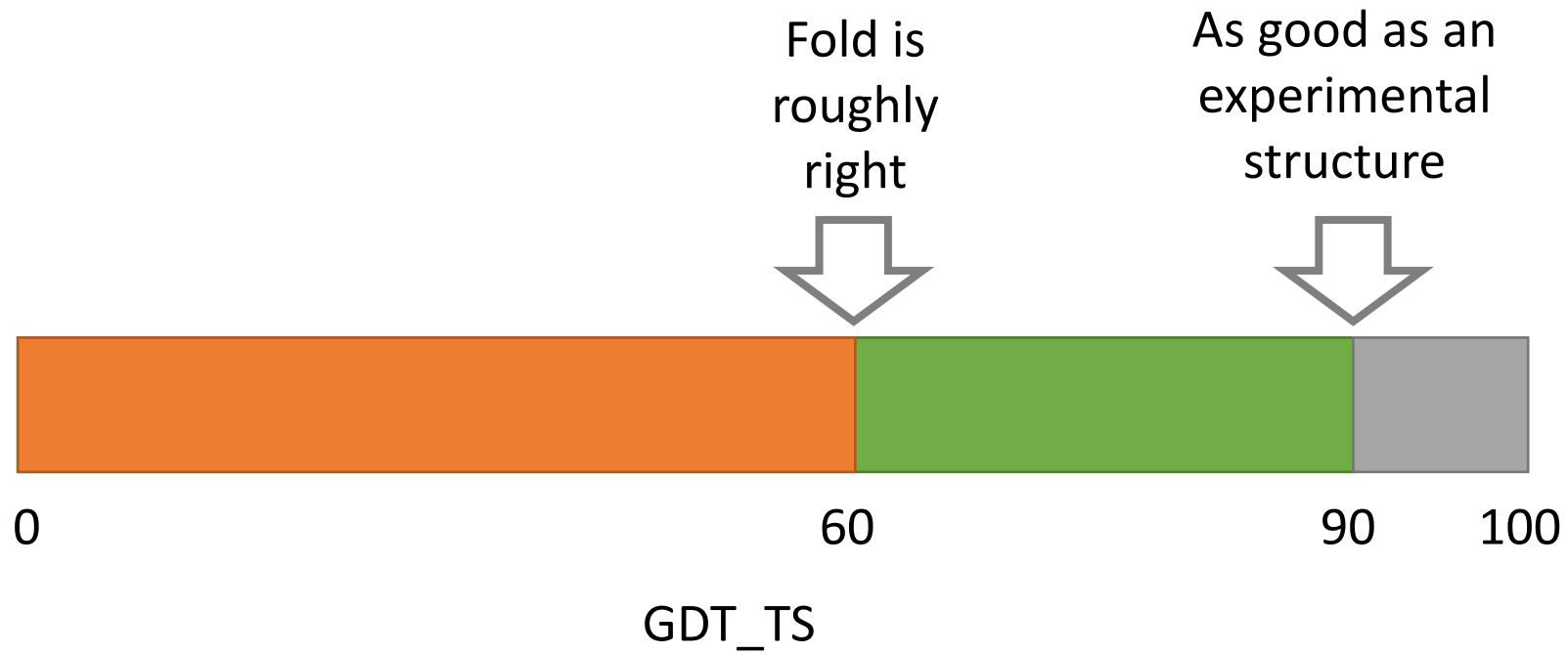- Ab initio modelling has been a huge challenge and actually stimulated an open competition called CASP

UNIVERSITY OF
CAMBRIDGE

# Critical Assessment of protein Structure Prediction (CASP)

- Experiment initiated by John Moult started in 1994

- Independent benchmark of ability to predict novel protein structures

- Targets ranked based on 'difficulty'

- Scored on accuracy and coverage of backbone prediction (GDT_TS)

# How good is the prediction?

Fold is roughly right

As good as an experimental structure

0          60          90   100

GDT_TS

UNIVERSITY OF
CAMBRIDGE

# Predicting structures

- CASP outcomes and scoring

# AlphaFold was a huge step forward

- DeepMind developed a deep learning approach to structure prediction for CASP13 (2016)

- Step-change in quality

# CASP14 and AlphaFold2

- DeepMind completely redesigned their prediction pipeline

- Unparalleled accuracy

- Quality of the prediction was largely decoupled from sequence identity to homologues



UNIVERSITY OF CAMBRIDGE

# How does AlphaFold2 do this?



MSA embedding

Sequence-residue edges

Protein sequence

Genetics search & embed

Embed & outer sum

Residues →
Sequences

Residues →
Sequences

Residues →
Residues

Residues →
Residues

Residue-residue edges

Confidence Score

Structure module

Pairwise distances

3D structure

# How does AlphaFold2 do this?

- The only input needed is a sequence!

# How does AlphaFold2 do this?

- Powerful neural networks
  - Attention-based neural networks



MSA embedding  Sequence–residue edges

Protein sequence

Genetics search & embed

Embed & outer sum

Residues →

Sequences

Residues →

Residues →

Confidence Score

Structure module

Residues →

Residues →

Pairwise distances

Residue–residue edges

3D structure

UNIVERSITY OF CAMBRIDGE

# How does AlphaFold2 do this?

- Powerful neural networks
  - Attention-based neural networks



- Protein sequences
  - Multiple sequence alignments
  - Interrogate co-evolution of residues

UNIVERSITY OF CAMBRIDGE

Image: www.deepmind.com

# How does AlphaFold2 do this?

- Powerful neural networks
  - Attention-based neural networks



- Protein sequences
  - Multiple sequence alignments
  - Interrogate co-evolution of residues

- Protein structures
  - Identifying residue pairs that should be close to each other
  - Experience of what folded proteins 'look like'

UNIVERSITY OF CAMBRIDGE

Image: www.deepmind.com

# How does AlphaFold2 do this?

- Powerful neural networks
  - Attention-based neural networks



- Very powerful computers

- Protein sequences
  - Multiple sequence alignments
  - Interrogate co-evolution of residues

- Protein structures
  - Identifying residue pairs that should be close to each other
  - Experience of what folded proteins 'look like'

UNIVERSITY OF CAMBRIDGE

Image: www.deepmind.com

Did Quake make AlphaFold happen?

Image: ID software

# Did Quake make AlphaFold happen?

Peak Double Precision (GFLOPs)

GPU
CPU

UNIVERSITY OF CAMBRIDGE

# AlphaFold2 in a nutshell



- Kendrew Lecture 2021:
  https://www.youtube.com/watch?v=jTO6odQNp90

# AF2 predictions of all proteins

- Teamed up with EBI to predict representative set of all known proteins (still ongoing…)

- Results for human and model organisms available already from Uniprot website



Image: UniProt

# How reliable is your AF2 model?

- AF2 will always give you a structure

- But that doesn't mean it is right

UNIVERSITY OF CAMBRIDGE

# How reliable is your AF2 model?

- AF2 will always give you a structure

- But that doesn't mean it is right

- **You have to check** the statistical plots and scores that are also generated

pLDDT scores

PAE plot

UNIVERSITY OF CAMBRIDGE

# pLDDT Scores and Plot

- This is a <u>per residue </u>score on scale of 0 to 100

- Score above 70 is a confident prediction

Model Confidence:

■ Very high (pLDDT > 90)

■ Confident (90 > pLDDT > 70)

■ Low (70 > pLDDT > 50)

■ Very low (pLDDT < 50)

UNIVERSITY OF
CAMBRIDGE

Image: UniProt

# pLDDT Scores and Plot

- This is a <u>per residue</u> score on scale of 0 to 100

- Score above 70 is a confident prediction

- Displayed as a coloured structure



Model Confidence:
- Very high (pLDDT > 90)
- Confident (90 > pLDDT > 70)
- Low (70 > pLDDT > 50)
- Very low (pLDDT < 50)

UNIVERSITY OF CAMBRIDGE

# pLDDT Scores and Plot

- This is a <u>per residue</u> score on scale of 0 to 100

- Score above 70 is a confident prediction

- Displayed as a coloured structure or a plot



Model Confidence:

- Very high (pLDDT > 90)
- Confident (90 > pLDDT > 70)
- Low (70 > pLDDT > 50)
- Very low (pLDDT < 50)

pLDDT scores

UNIVERSITY OF CAMBRIDGE

# PAE plots – Predicted Aligned Error

- This is a measure of confidence relative to other regions of the structure
- Low error is high confidence (blue)
- High error is low confidence (yellow)

# PAE plots – Predicted Aligned Error

- This is a measure of confidence relative to other regions of the structure

- Low error is high confidence (blue)

- High error is low confidence (yellow)

- In this example:
  - Domains 1 and 2 confident relative to each other
  - Domains 1 and 2 not confident relative to domains 4 and 5

# PAE plots – Predicted Aligned Error



Image: Janet Deane CCBY 4.0

# How to judge an AF2 model

- To summarise, a high-confidence per-residue model can be low-confidence overall



pLDDT score

# What about proteins like this?

**Human Afadin**



**Predicted accuracy**
**High → Low**

- If you see something like this, can you learn anything at all?

Image: UniProt

# What about proteins like this?

**Human Afadin**



**Predicted accuracy**

**High → Low**

**Actually quite high confidence these regions are unstructured**

UNIVERSITY OF CAMBRIDGE

Image: UniProt

# Look at the pLDDT plot



Predicted accuracy

**High**

↑

**Low**

UNIVERSITY OF CAMBRIDGE

# Look at the pLDDT plot

- Clearly determine domain boundaries

# Look at the PAE plot

# Look at the PAE plot

# Look at the PAE plot

# Look at the PAE plot

- The predictions are only confident within domains NOT BETWEEN

# In this case:

- This AF2 model is useful for:

  - Determining domain boundaries

  - Fold of individual domains

# In this case:

- This AF2 model is useful for:

  - Determining domain boundaries

  - Fold of individual domains

  - You could then use these individual domains to search using DALI or FoldSeek to find structural homologues that may inform function

# AF2 doesn't know about topology



**Human PTPRK**

UNIVERSITY OF CAMBRIDGE

# AF2 doesn't know about topology



**Intracellular**

**TM**

**Extracellular**

UNIVERSITY OF
CAMBRIDGE

# AF2 doesn't know about topology

**Predicted accuracy**

**High** → **Low**



- Treat membrane-spanning models with caution...

# Validate using experimental techniques

AlphaFold Prediction

# Validate using experimental techniques
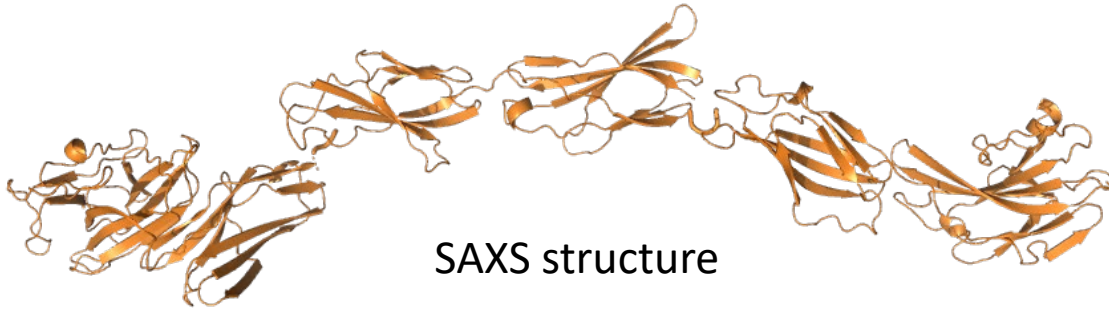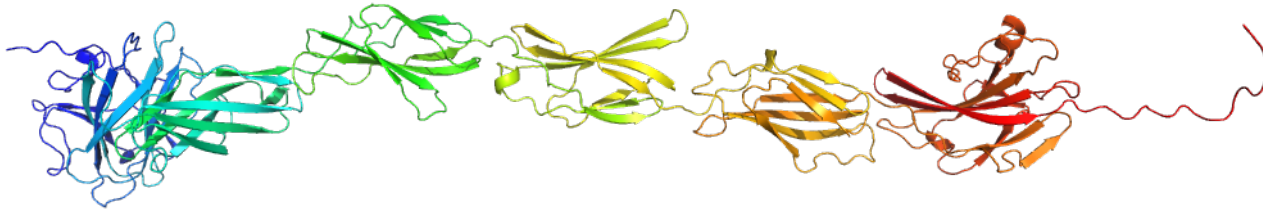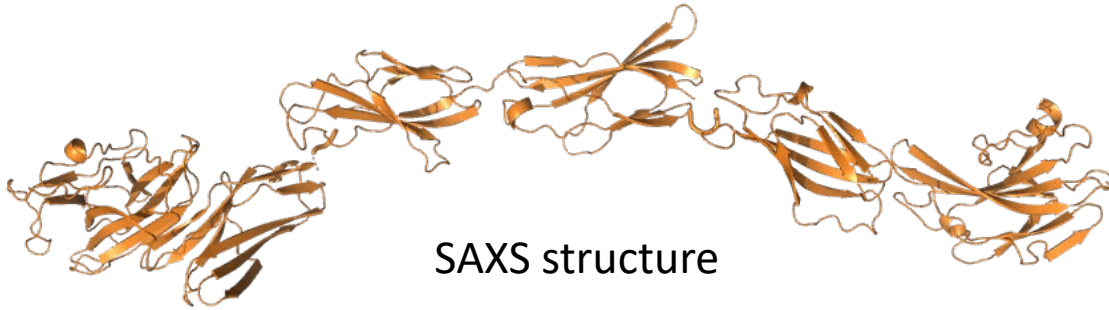
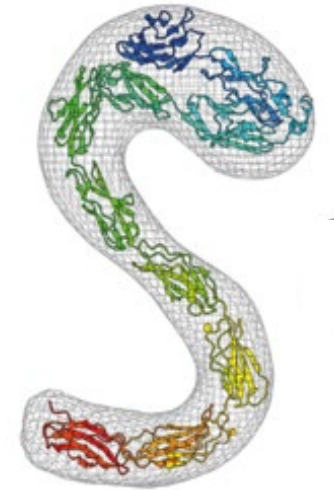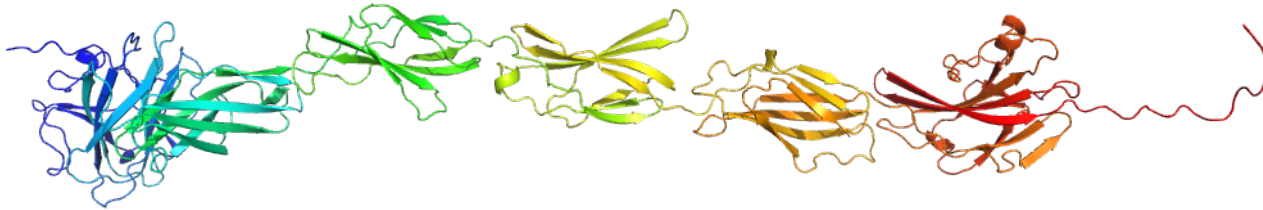AlphaFold Prediction



SAXS structure

# Validate using experimental techniques

AlphaFold Prediction



AlphaFold Prediction



SAXS structure

UNIVERSITY OF
CAMBRIDGE

# Validate using experimental techniques



AlphaFold Prediction

SAXS structure

AlphaFold Prediction

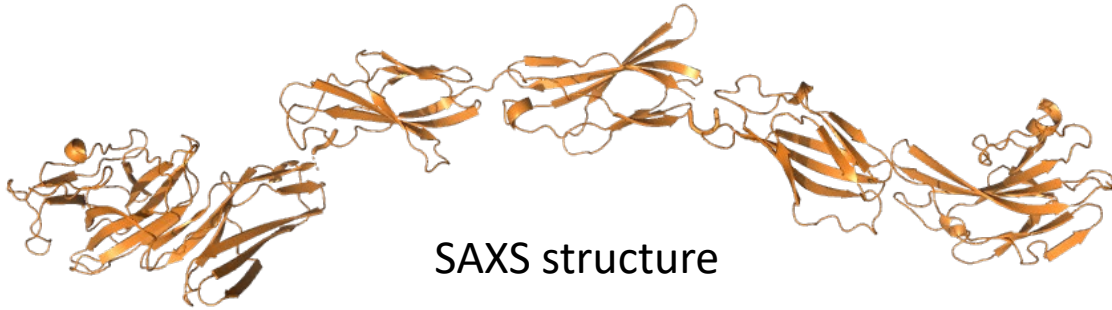EM structure

UNIVERSITY OF CAMBRIDGE

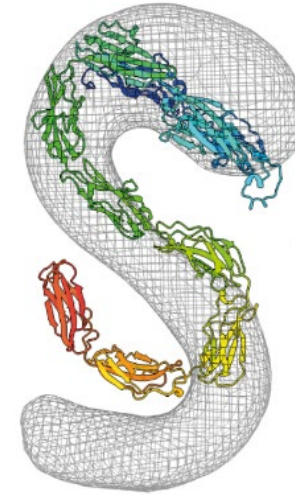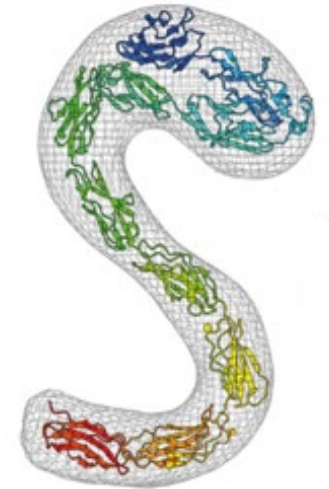# Validate using experimental techniques



AlphaFold Prediction

SAXS structure

AlphaFold Prediction

EM structure

- Importantly, AF2 provided excellent starting models for these experimental approaches

UNIVERSITY OF CAMBRIDGE

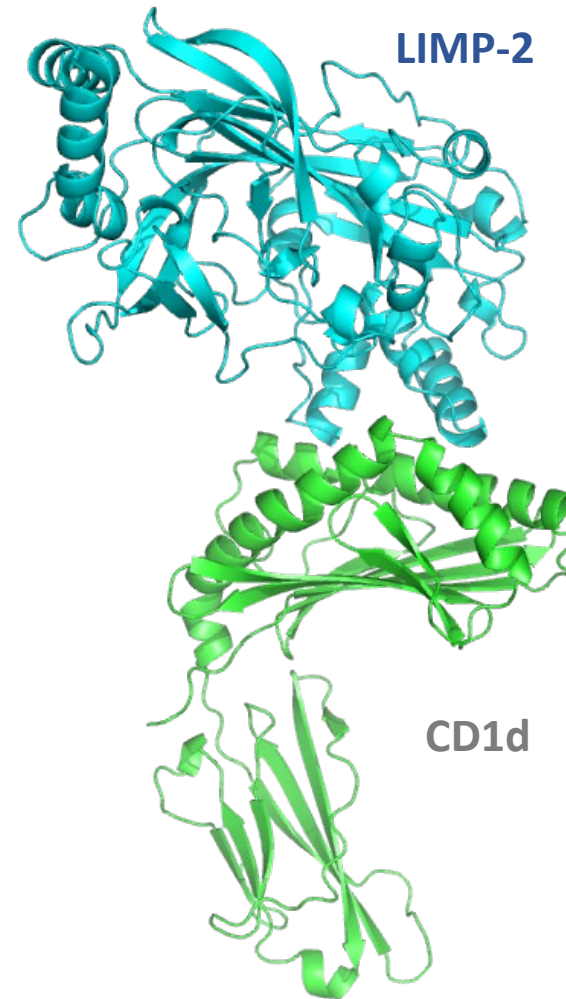# AF2 is pretty good at single proteins…what about complexes?

- AF2 Multimer was developed to try and address this question

- Answer is mixed, again you have to know how to interpret the statistics of the models produced

- Two examples: CD1d-LIMP2
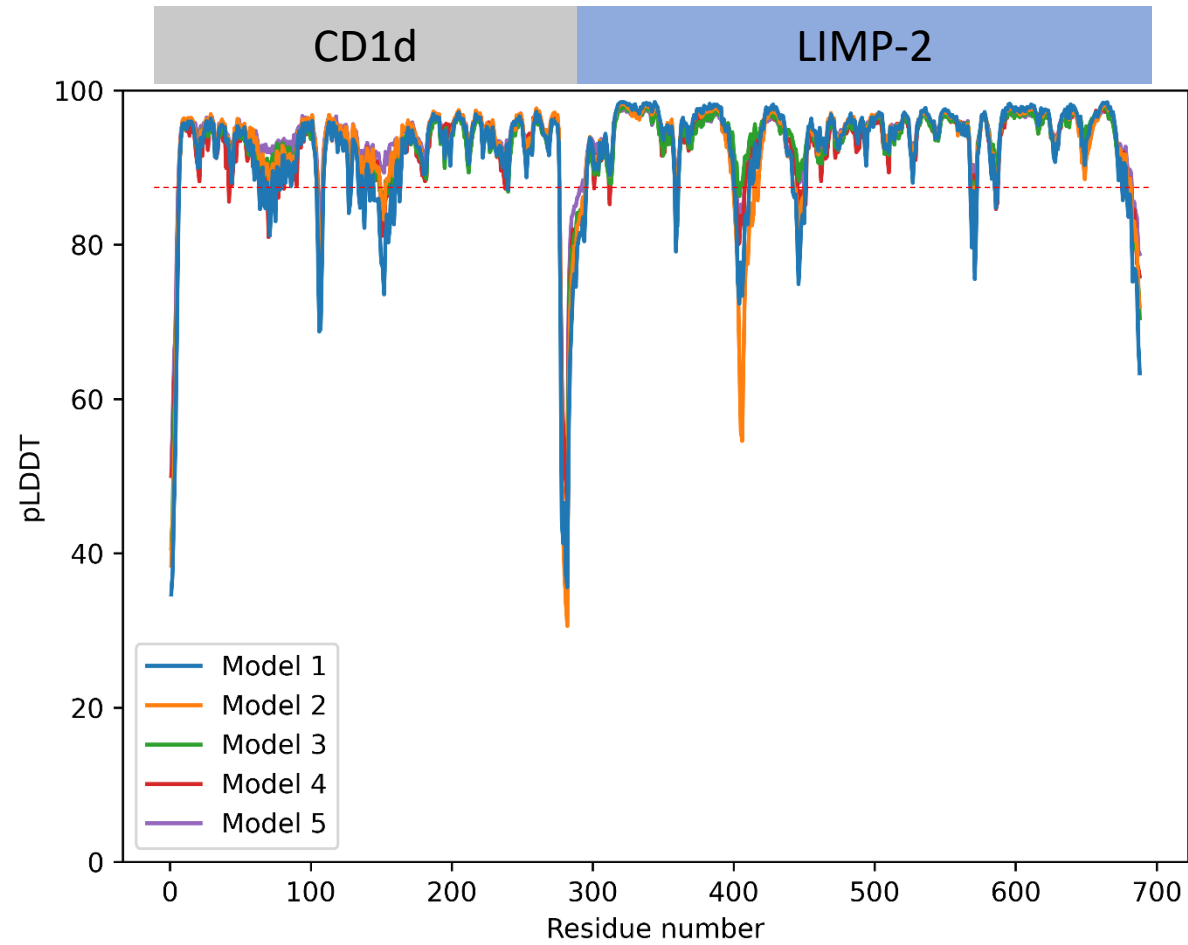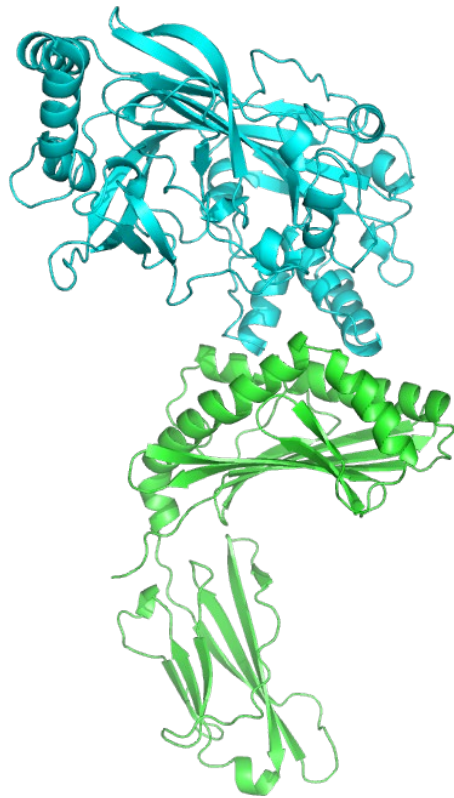
    PTPRK-Afadin

# CD1d-LIMP2

- Lipid binding proteins: CD1d is like MHC-I, LIMP-2 has a lipid tunnel

- From literature they're predicted to interact

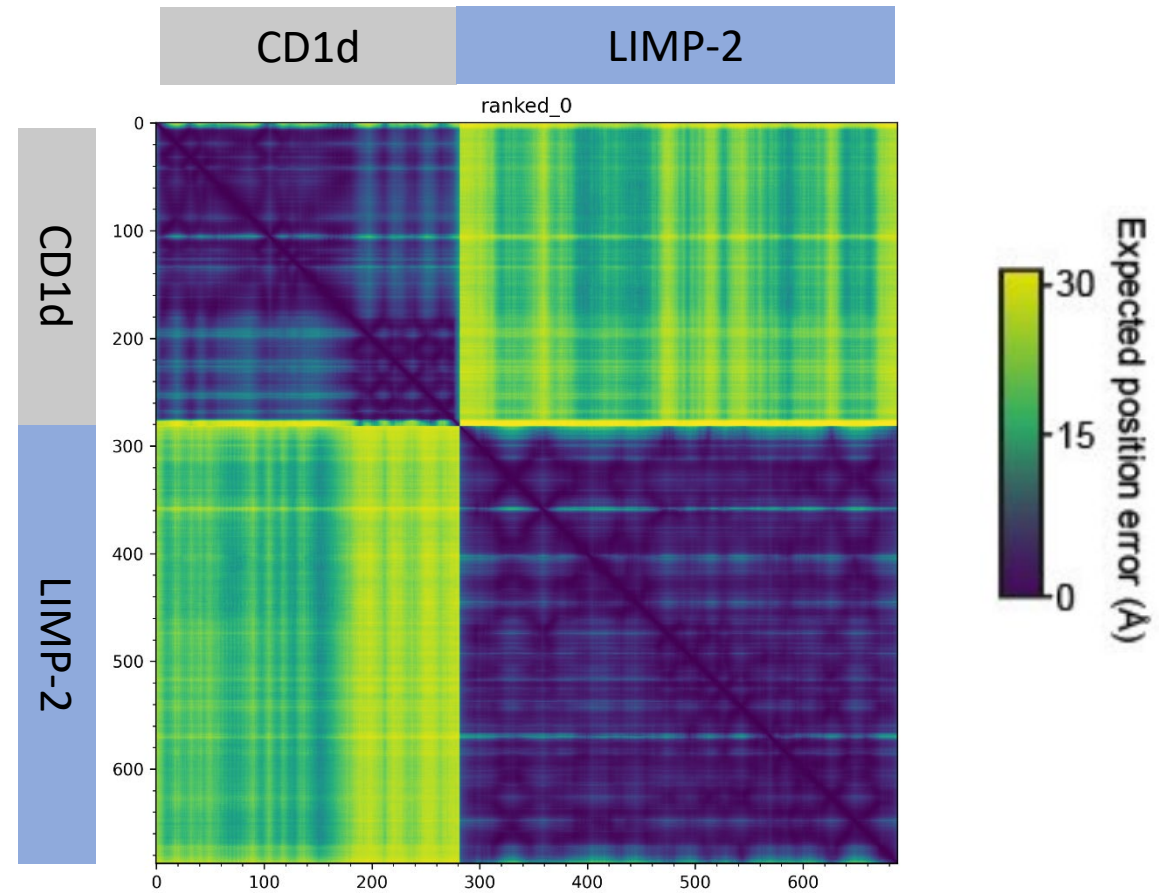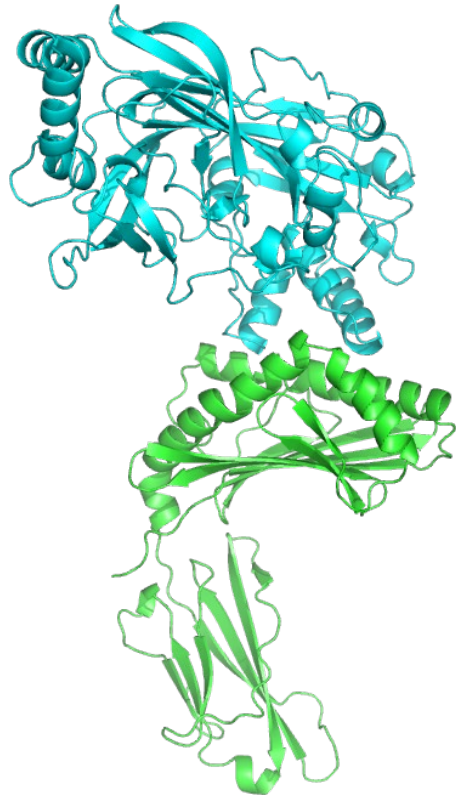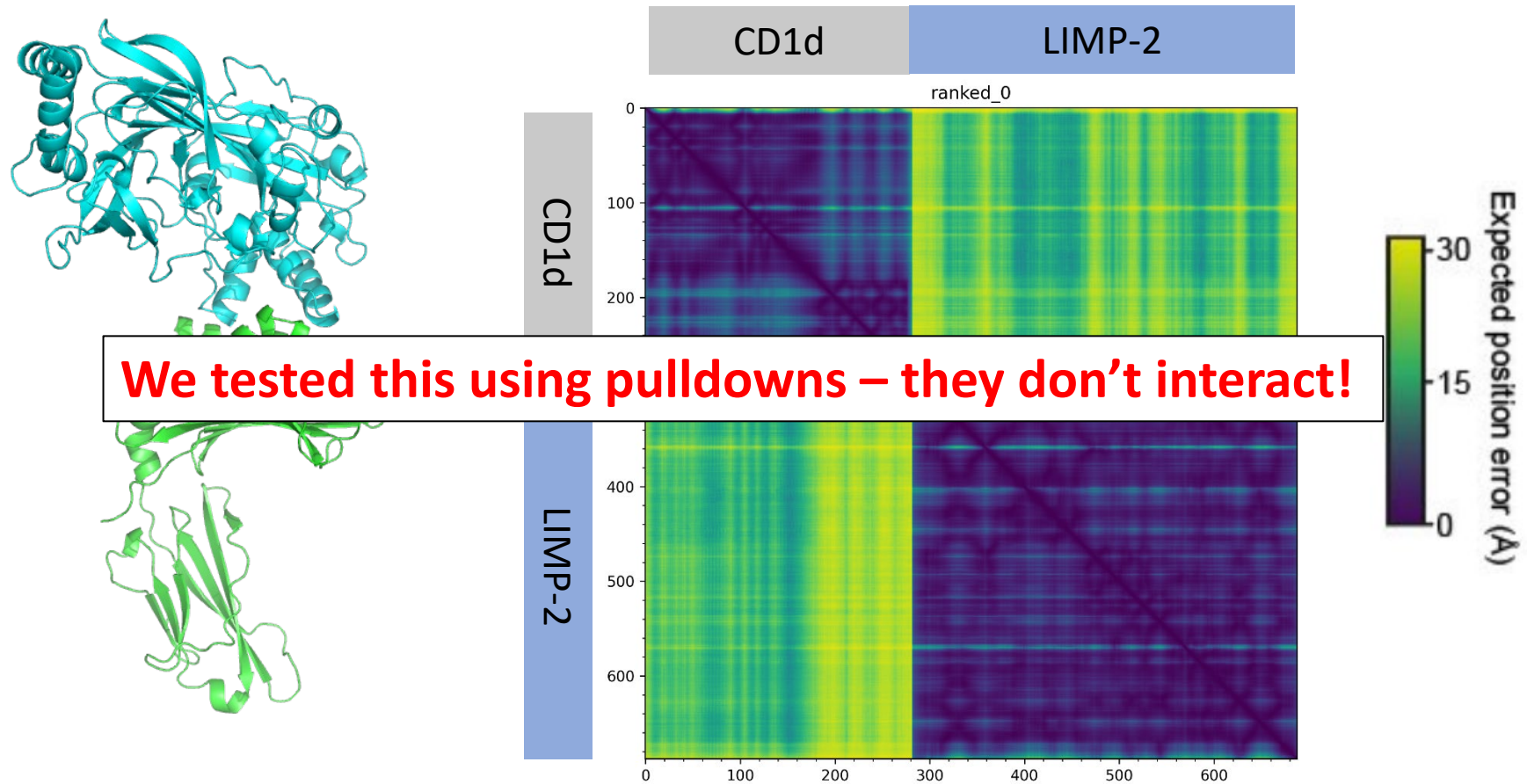- AF2 predicts a consistent complex (all models agree)



LIMP-2

CD1d

UNIVERSITY OF CAMBRIDGE

# CD1d-LIMP2

- pLDDT plot

# CD1d-LIMP2

- PAE plot

# CD1d-LIMP2
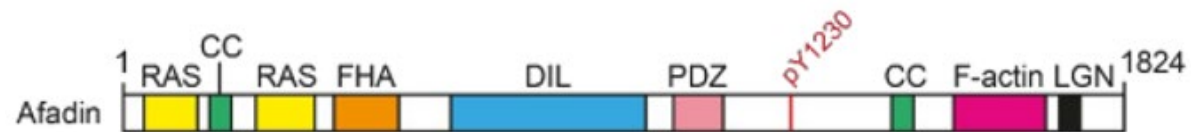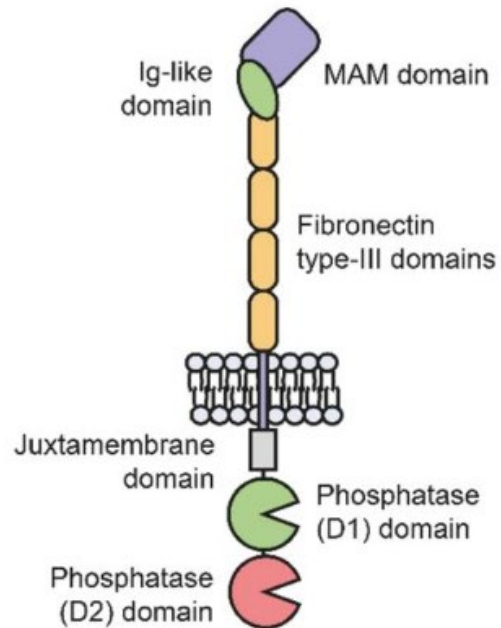
- PAE plot



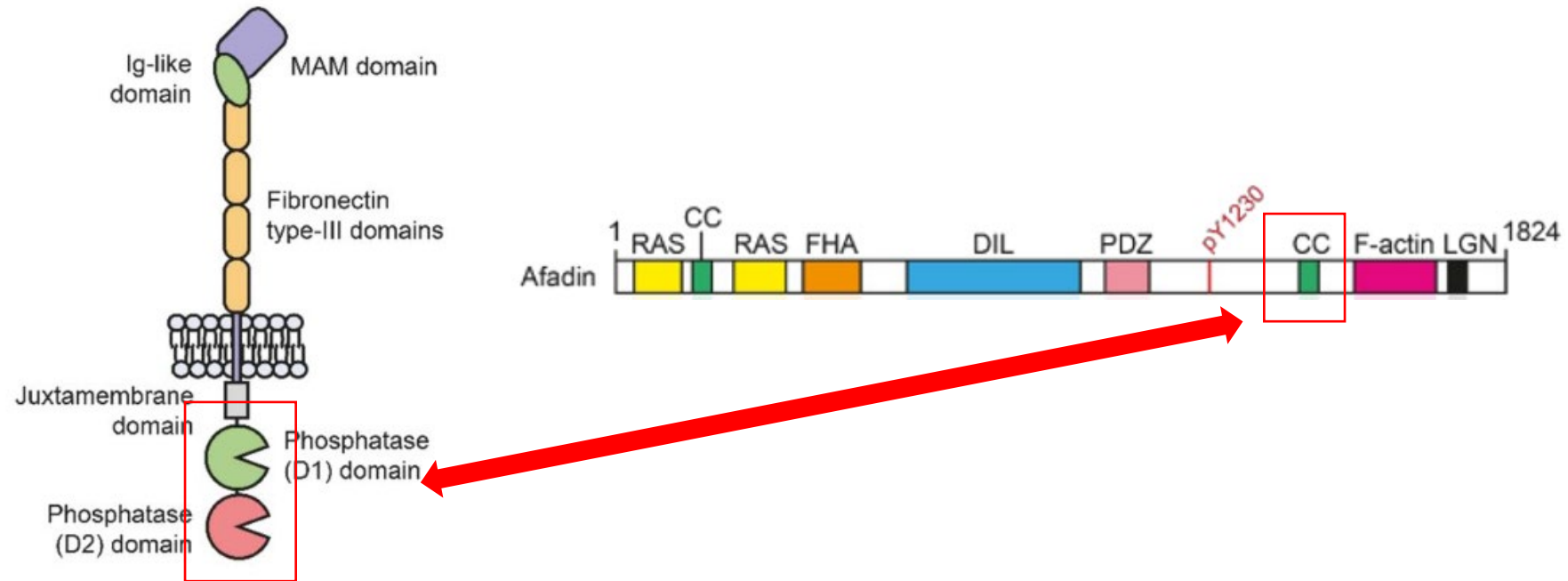**We tested this using pulldowns – they don't interact!**

# PTPRK-Afadin

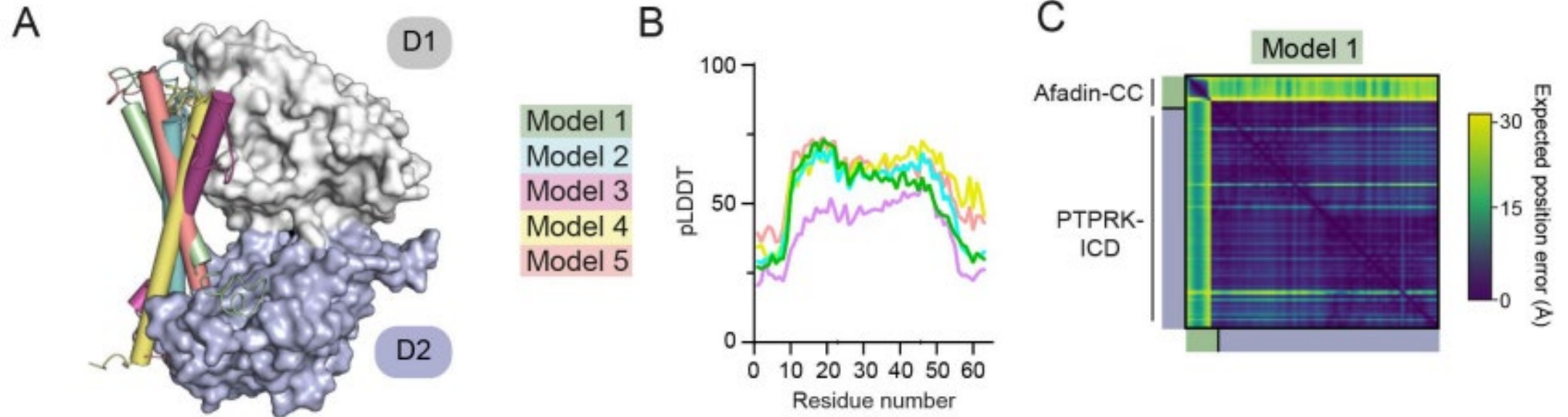- We knew that PTPRK binds Afadin but these are both BIG proteins

# PTPRK-Afadin
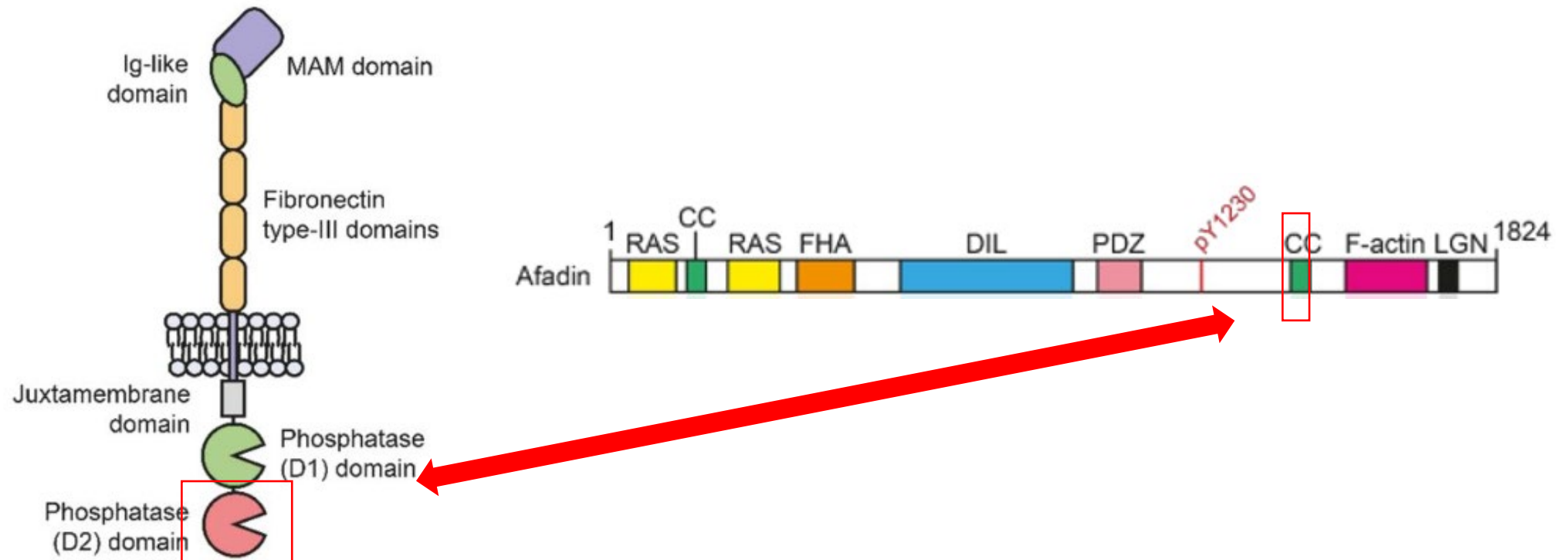
- In the lab, mapped this down to much smaller domains

# PTPRK-Afadin

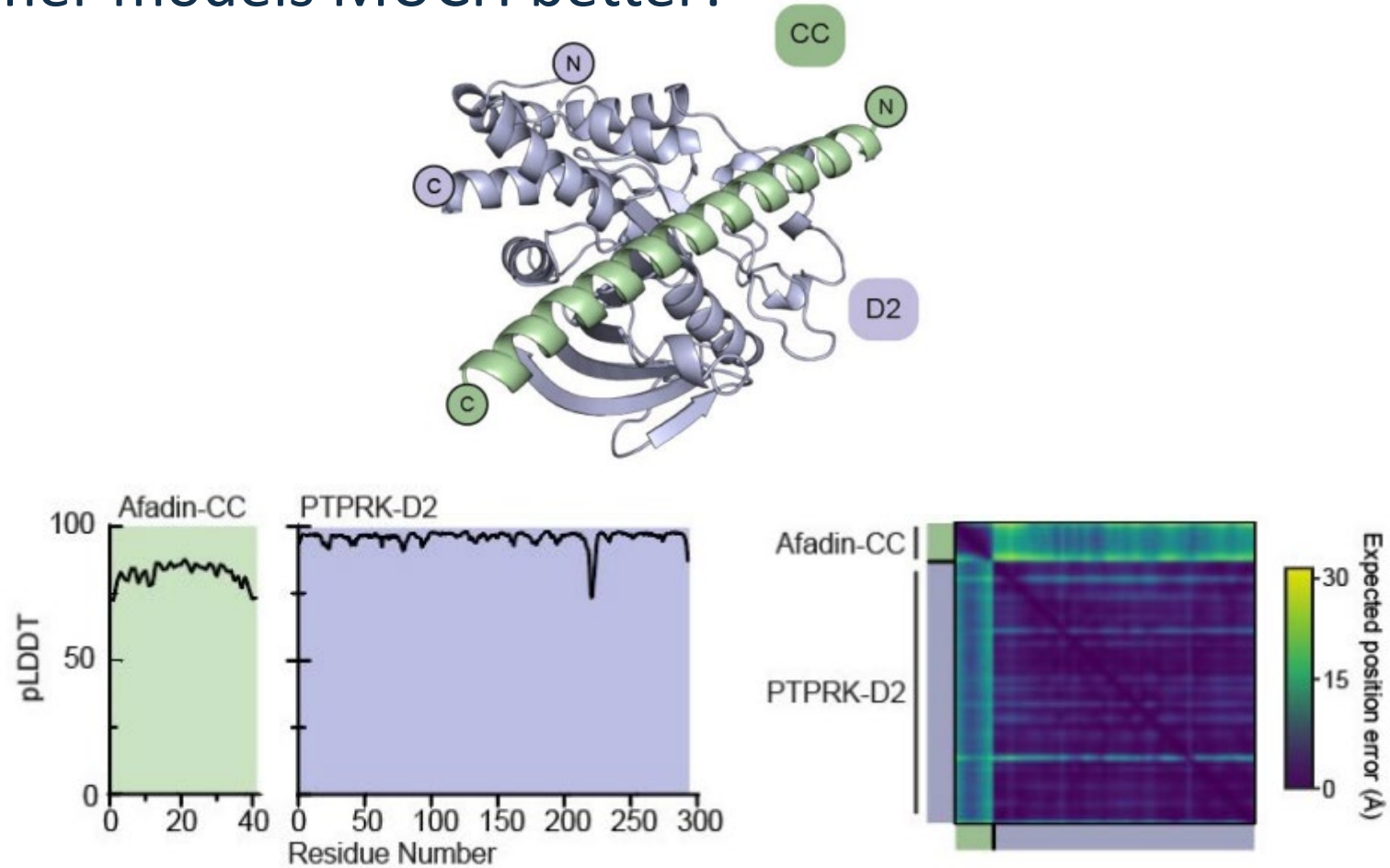- But AF2 Multimer models weren't good

# PTPRK-Afadin

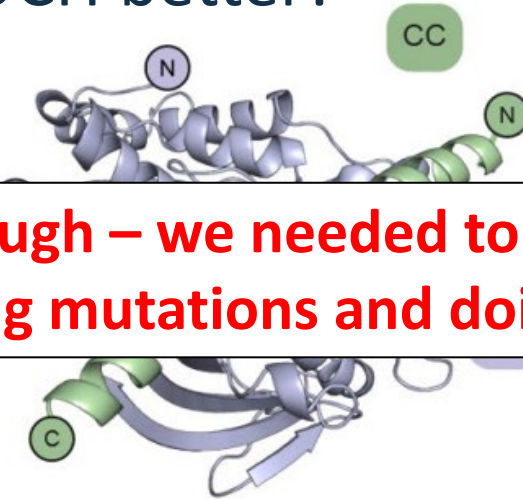- Experimentally mapped it down to smaller pieces
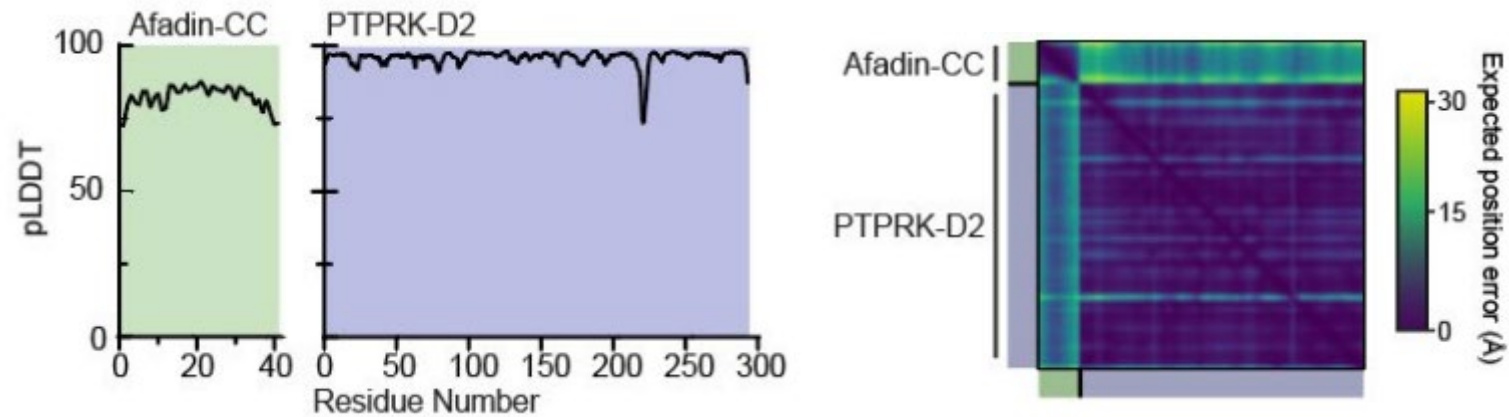
# PTPRK-Afadin

- AF2 Multimer models MUCH better!

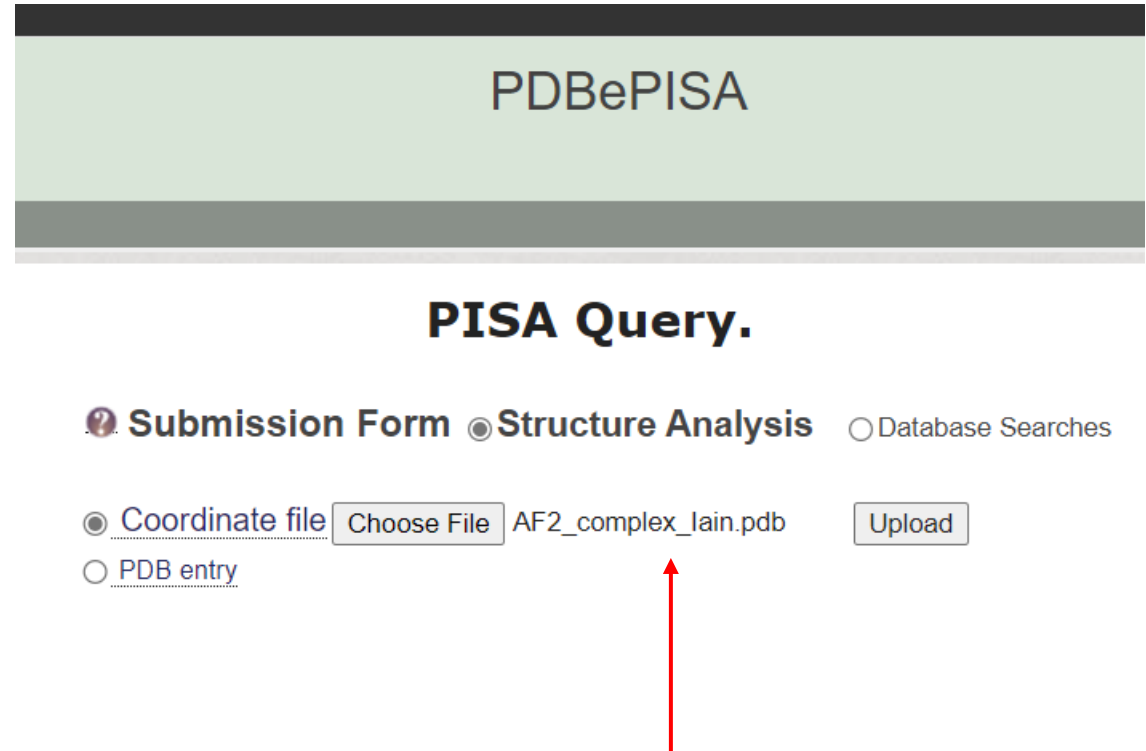# PTPRK-Afadin

- AF2 Multimer models MUCH better!



**This was not enough – we needed to validate the interface by making mutations and doing pulldowns**

# Using PDBePISA with AlphaFold Models

- Upload AF2 model of complex to PDBePISA

# Using PDBePISA with AlphaFold Models

# Using PDBePISA with AlphaFold Models

# Using PDBePISA with AlphaFold Models

**Hydrogen bonds** [XML]

| ## | Structure 1 | Dist. [Å] | Structure 2 |
|----|-------------|-----------|-------------|
| 1 | A:ARG 23[HH11] | 1.77 | B:GLU 179[ OE2] |
| 2 | A:ARG 23[HH21] | 2.15 | B:ASP 118[ O ] |
| 3 | A:ARG 25[HH22] | 2.01 | B:GLU 223[ OE1] |
| 4 | A:ARG 25[HH21] | 2.16 | B:GLU 223[ OE2] |
| 5 | A:GLN 30[HE22] | 1.79 | B:GLU 220[ OE2] |
| 6 | A:LYS 33[ HZ3] | 2.13 | B:GLU 221[ OE1] |
| 7 | A:GLU 22[ OE1] | 2.07 | B:ARG 225[HH11] |
| 8 | A:GLU 22[ OE2] | 1.83 | B:ARG 225[HH22] |

**Salt bridges** [XML]

| ## | Structure 1 | Dist. [Å] | Structure 2 |
|----|-------------|-----------|-------------|
| 1 | A:ARG 23[ NE ] | 3.84 | B:GLU 179[ OE2] |
| 2 | A:ARG 23[ NH1] | 2.73 | B:GLU 179[ OE2] |
| 3 | A:ARG 23[ NH1] | 3.29 | B:GLU 179[ OE1] |
| 4 | A:ARG 23[ NH2] | 3.75 | B:ASP 118[ OD2] |
| 5 | A:ARG 25[ NH2] | 2.90 | B:GLU 223[ OE1] |
| 6 | A:ARG 25[ NH2] | 2.81 | B:GLU 223[ OE2] |
| 7 | A:LYS 26[ NZ ] | 3.45 | B:GLU 223[ OE2] |
| 8 | A:LYS 33[ NZ ] | 2.81 | B:GLU 221[ OE1] |
| 9 | A:GLU 22[ OE1] | 3.07 | B:ARG 225[ NH1] |
| 10 | A:GLU 22[ OE1] | 3.70 | B:ARG 225[ NH2] |
| 11 | A:GLU 22[ OE2] | 3.56 | B:ARG 225[ NH1] |
| 12 | A:GLU 22[ OE2] | 2.82 | B:ARG 225[ NH2] |

No disulfide bonds found

No covalent bonds found

**UNIVERSITY OF CAMBRIDGE**

# Using PDBePISA with AlphaFold Models

# PTPRK-Afadin

- Our pulldowns using mutations based on the AF2 model validated the interface experimentally

# A few caveats

- PDBePISA didn't predict this interface to be significant - but it was!



- AF2 renumbers your residues so they might no longer match the Uniprot entry – you can renumber your model using Coot

```
A:ARG  25[ NH2]    2.81    B:GLU 223[ OE2]
A:LYS  26[ NZ ]    3.45    B:GLU 223[ OE2]
```

# Renumber residues in Coot

# Renumber residues in Coot

# So, what is AlphaFold2 good for?



- Determining the fold of protein domain(s)
  - Identify potential functional homology

- Determining domain boundaries
  - Clone sensible constructs

- Protein:peptide complexes
  - And some protein:protein complexes

# So, what is AlphaFold2 good for?



- Determining the fold of protein domain(s)
  - Identify potential functional homology

- Determining domain boundaries
  - Clone sensible constructs

- Protein:peptide complexes
  - And some protein:protein complexes

- AF2 models should always be:
  - Shown with their statistical plot
  - Tested experimentally

# What AF2 isn't good at (yet!)

- Most protein:protein complexes
  - But gives testable hypotheses

- Predicting surface properties
  - OK but not perfect, interpret with caution

# What AF2 isn't good at (yet!)

- Most protein:protein complexes
  - But gives testable hypotheses
- Predicting surface properties
  - OK but not perfect, interpret with caution
- Predicting ligands (Zn, haem, co-factors, drugs etc)
- Understanding topology, intracellular vs extracellular domains

# What AF2 isn't good at (yet!)

- Most protein:protein complexes
  - But gives testable hypotheses

- Predicting surface properties
  - OK but not perfect, interpret with caution

- Predicting ligands (Zn, haem, co-factors, drugs etc)

- Understanding topology, intracellular vs extracellular domains

- Importantly, AF2 is not designed to test the effect of point mutations
  - Structure predictions rely on multiple sequence alignments and co-evolution
  - To understand point mutations you still need to manually inspect the structure

# What AF2 isn't good at (yet!)

- Most protein:protein complexes
  - But gives testable hypotheses

- Predicting surface properties

> AlphaFold is being constantly developed and expanded
> It is likely several of these limitations will be overcome eventually

- Understanding topology, intracellular vs extracellular domains

- Importantly, AF2 is not designed to test the effect of point mutations
  - Structure predictions rely on multiple sequence alignments and co-evolution
  - To understand point mutations you still need to manually inspect the structure

# Try it yourself

- You can access all the pre-calculated AlphaFold structures by DeepMind/EMBL-EBI:
    - https://alphafold.ebi.ac.uk/

- You can run AF2 yourself via the browser (Google Colab):
    - https://colab.research.google.com/github/deepmind/alphafold/blob/main/notebooks/AlphaFold.ipynb

UNIVERSITY OF CAMBRIDGE

# Try it yourself

- You can access all the pre-calculated AlphaFold structures by DeepMind/EMBL-EBI:
  - https://alphafold.ebi.ac.uk/

- You can run AF2 yourself via the browser (Google Colab):
  - https://colab.research.google.com/github/deepmind/alphafold/blob/main/notebooks/AlphaFold.ipynb

- NOTE: if you want to run locally on your computer you need a very powerful computer (GPU with lots of RAM) and we recommend installing ColabFold not AlphaFold

UNIVERSITY OF CAMBRIDGE