# Disorder in CAID-2

**Damiano Piovesan**

*University of Padova - Italy*

CASP15, 10-13 December 2022, Antalya
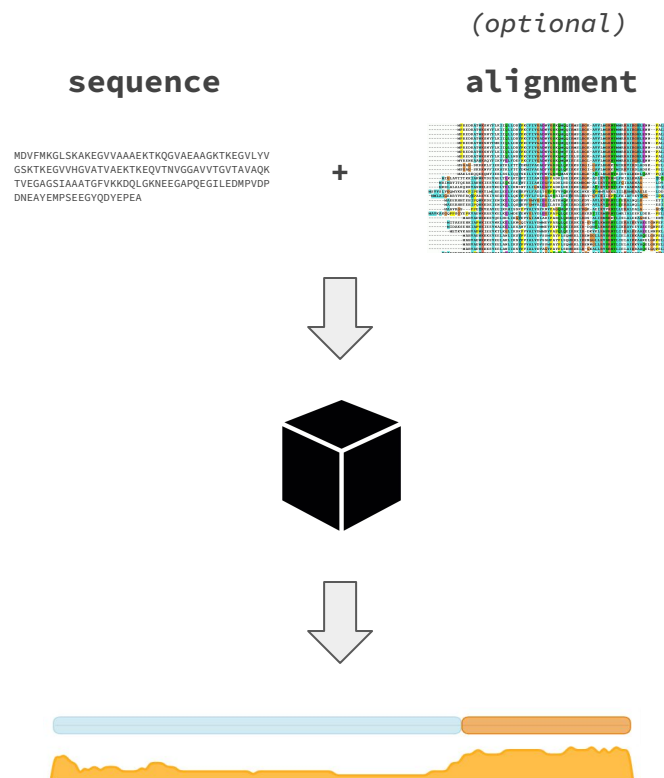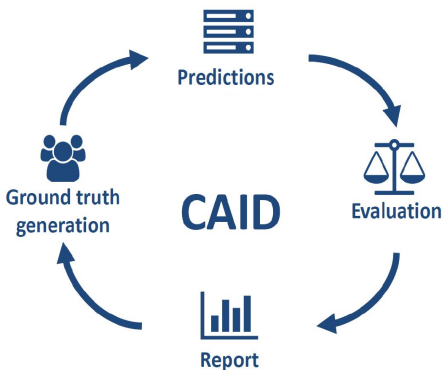
# The challenge

Prediction categories

- **Disorder** - Disordered regions

- **Binding** - Binding residues inside disordered regions

- **Nucleic Acid binding** - Residues inside disordered regions that bind DNA/RNA molecules

- **Linker** - Entropic chains

sequence

alignment

MDVFMKGLSKAKEGVVAAAEKTKQGVAEAAGKTKEGVLYV
GSKTKEGVVHGVATVAEKTKEQVTNVGGAVVTGVTAVAQK
TVEGAGSIAAATGFVKKDQLGKNEEGAPQEGILEDMPVDP
DNEAYEMPSEEGYQDYEPEA

+

2

# The CAID cycle



Predictions → Evaluation → Report → Ground truth generation (CAID cycle)

Marco Necci [1,50], Damiano Piovesan [1,50], CAID Predictors*, DisProt Curators* and Silvio C. E. Tosatto [1]

- **Ground truth generation**
  - Literature curation (DisProt)

- **Prediction**
  - Execution of stand-alone software (containers)

- **Assessment**
  - Accuracy & Technical evaluation
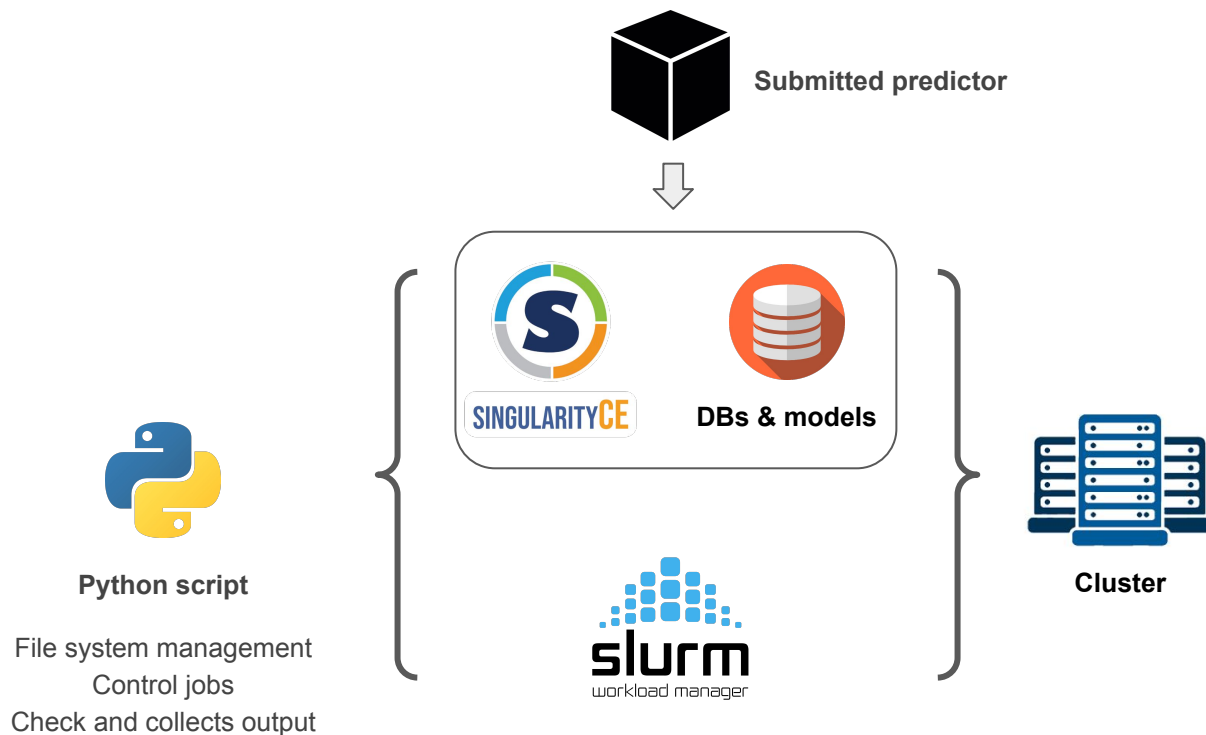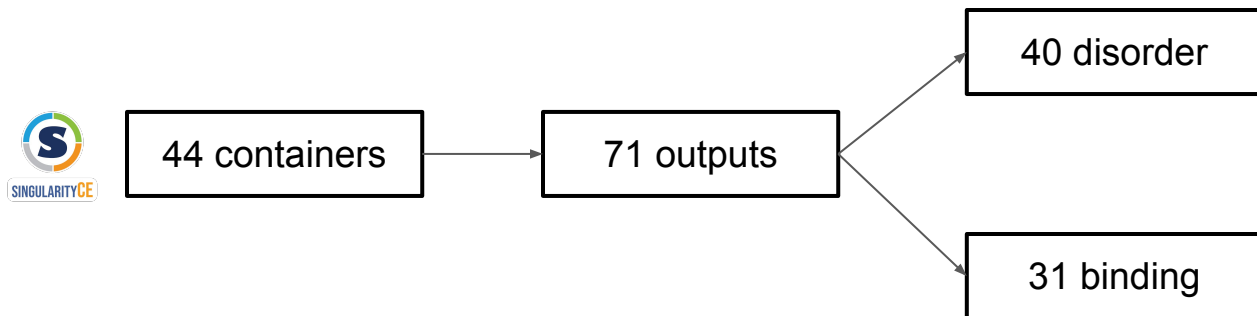
- **Report**
  - CAID & CASP conferences

# Execution pipeline

**Submitted predictor**

**SINGULARITYCE**

**DBs & models**

**Cluster**

**Python script**

File system management
Control jobs
Check and collects output

**slurm**
workload manager

# Methods



| Container size (Gb) | # containers |
|---|---|
| < 0.5 | 14 |
| > 0.5 | 9 |
| > 1 | 13 |
| > 2 | 5 |
| > 3 | 3 |

- **22** packages are **new**, i.e. not tested in CAID 1

- **18** packages also available as **web servers**

# Average execution time per protein (containers)

- Additional execution time is included
  - BLAST
  - …
- Image instantiation time not included
- Some containers include both fast and slow predictors
  - MorfChibi / MorfChibi light
  - SPOT-Disorder / SPOT-Disorder-single
  - …



6

# CAID-2

# Ground truth

# DisProt

Manually curated annotations of from the literature

Community

- 40 curators
- 30 laboratories
- 20 countries
- 2 reviewers

*disprot.org*

# Disorder annotations in DisProt



disorder **ID** DP00018r036 **Curator** Federica Quaglia

**Fragment** 22-105
**Method** nuclear magnetic resonance spectroscopy evidence used in manual assertion
**Reference** p27 binds cyclin-CDK complexes through a sequential mechanism involving binding-induced protein folding. *Lacy ER, Filippov I, Lewis WS, Otieno S, Xiao L, Weiss S, Hengst L, Kriwacki RW.* Nat Struct Mol Biol, 2004

disorder **ID** DP00018r022 **Curator** Federica Quaglia

**Fragment** 22-97
**Method** nuclear magnetic resonance spectroscopy evidence used in manual assertion
**Reference** Functional consequences of preorganized helical structure in the intrinsically disordered cell-cycle inhibitor p27(Kip1). *Bienkiewicz EA, Adkins JN, Lumb KJ.* Biochemistry, 2002

disorder **ID** DP00018r021 **Curator** Federica Quaglia

**Fragment** 97-197
**Method** X-ray crystallography-based structural model with missing residue coordinates used in manual assertion
**Reference** Structural basis of divergent cyclin-dependent kinase activation by Spy1/RINGO proteins. *McGrath DA, Fifield BA, Marceau AH, Tripathi S, Porter LA, Rubin SM.* EMBO J, 2017

disorder **ID** DP00018r011 **Curator** Federica Quaglia

**Fragment** 1-198
**Method** far-UV circular dichroism evidence used in manual assertion
**Reference** Functional consequences of preorganized helical structure in the intrinsically disordered cell-cycle inhibitor p27(Kip1). *Bienkiewicz EA, Adkins JN, Lumb KJ.* Biochemistry, 2002

## DisProt — DP00018 - Cyclin-dependent kinase inhibitor 1B

Browse Ontology Release notes Download Help About Biocuration

**Organism** Homo sapiens **Gene** CDKN1B (KIP1, p27) **Sequence length** 198 **Disorder content** 100%
**Homologous entries** DP01128 (50%)
**Cross references** UniProtKB:P46527, MobiDB:P46527, FuzDB: FC00036, AlphaFold: P46527, UniRef50:P46527
**Dataset(s)** Autophagy-related proteins   Cancer-related proteins
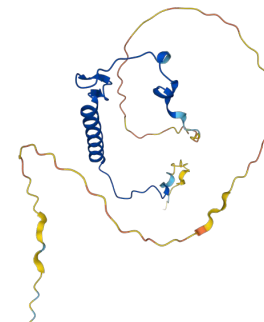**Last update** 2022-06-14

Download
Entry history

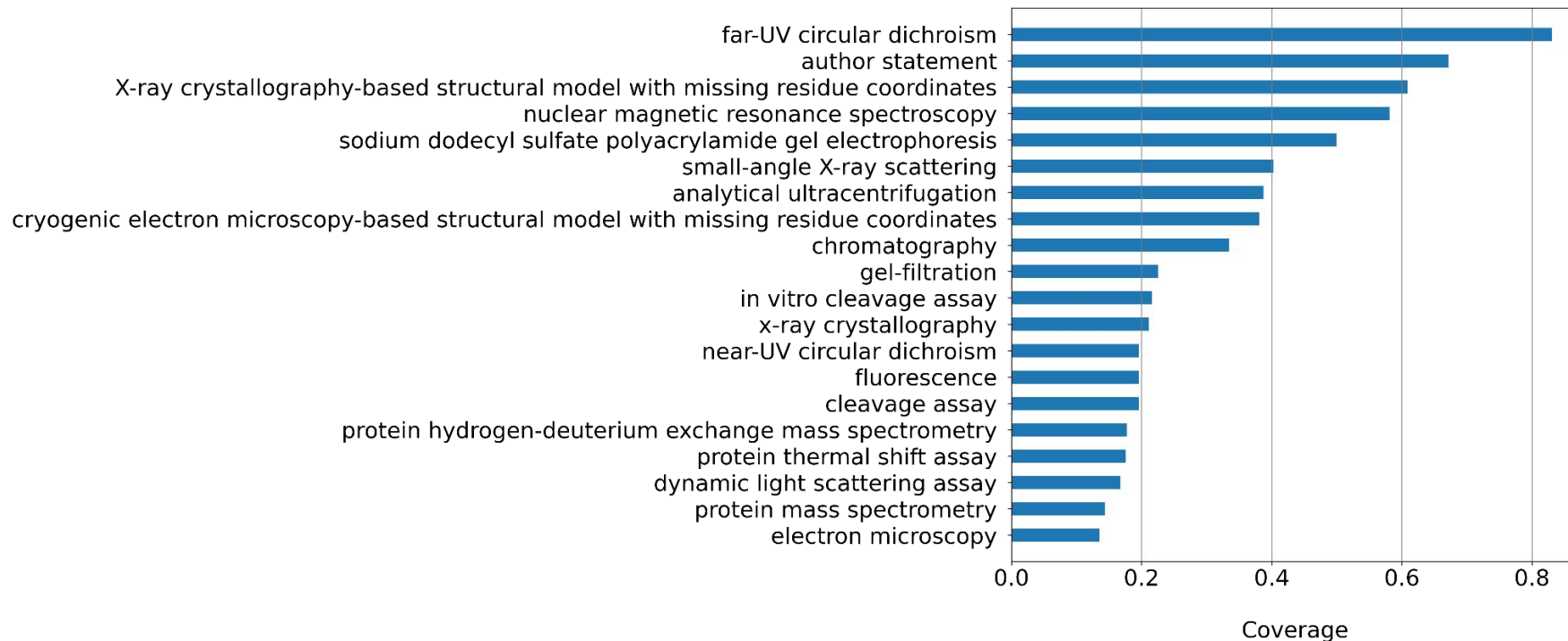Collapse Feature-Viewer   Expand Feature-Viewer   Toggle sequence viewer

DisProt consensus
Structural state
disorder
DP00018r037
DP00018r036
DP00018r023
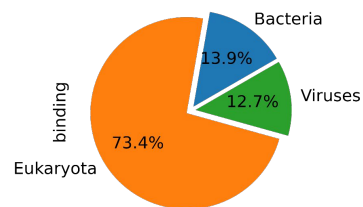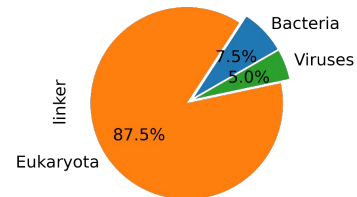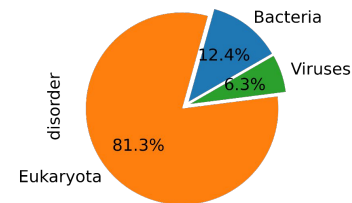DP00018r022
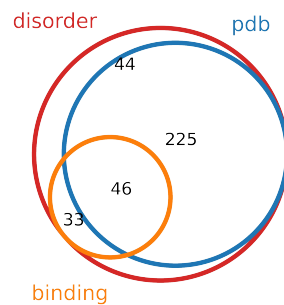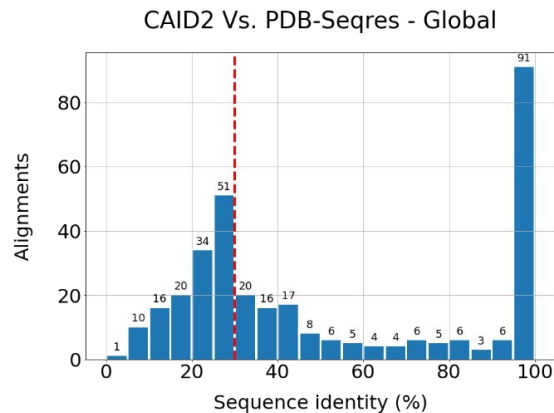DP00018r021
DP00018r020
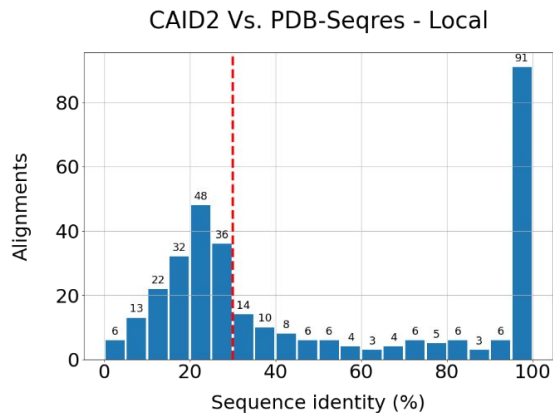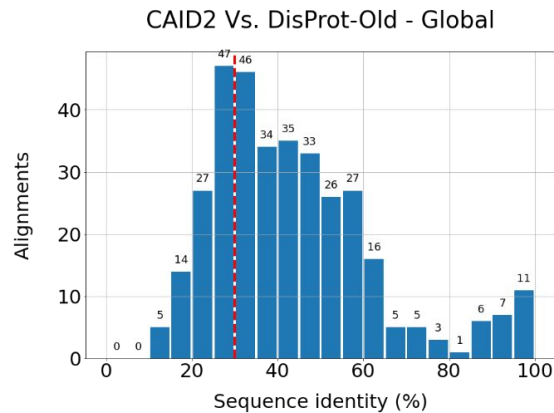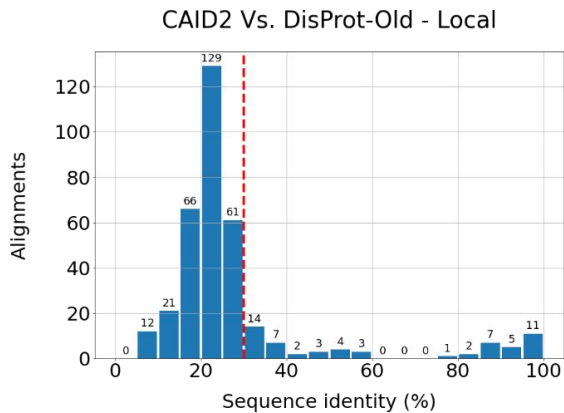DP00018r011
DP00018r008

P46527

9

# Experimental methods in DisProt

# CAID-2 benchmark proteins

# CAID-2 benchmark sequence identity
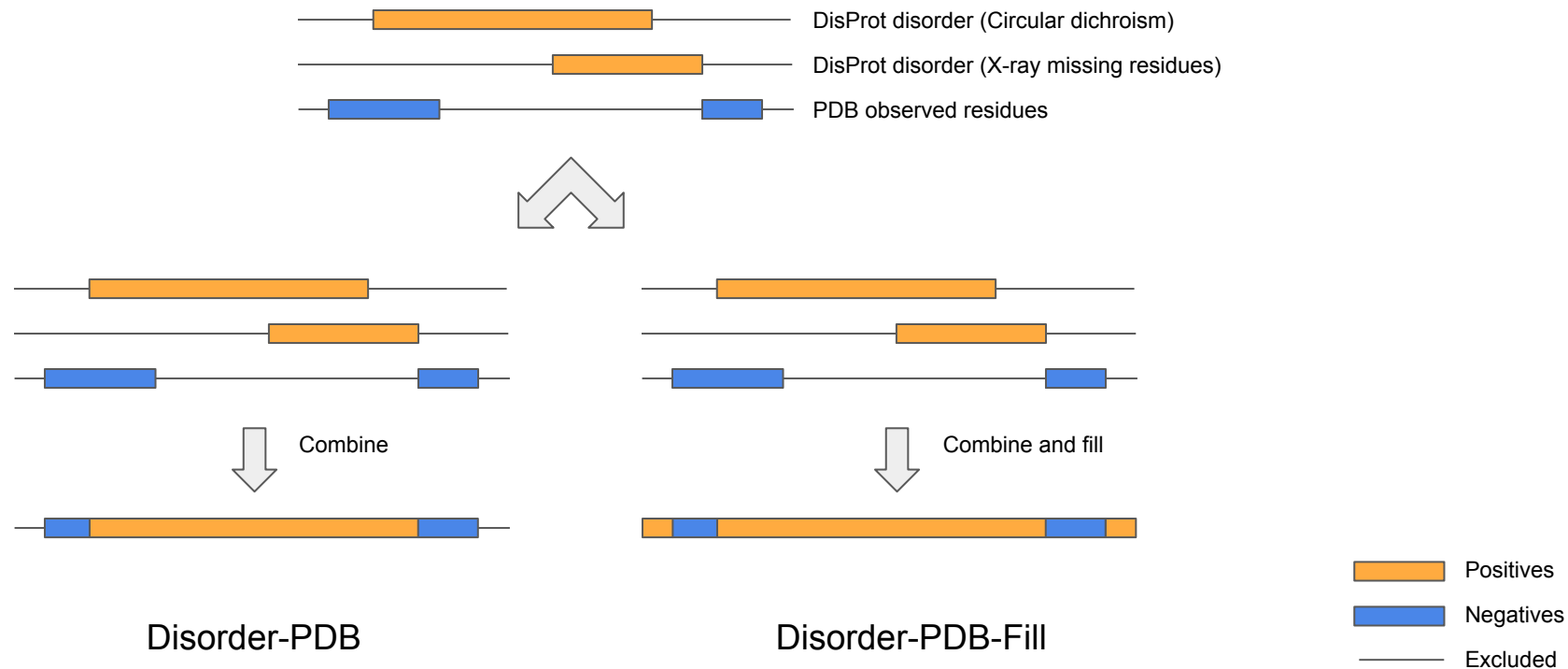
# Reference definition - Disorder



DisProt disorder (Circular dichroism)

DisProt disorder (X-ray missing residues)
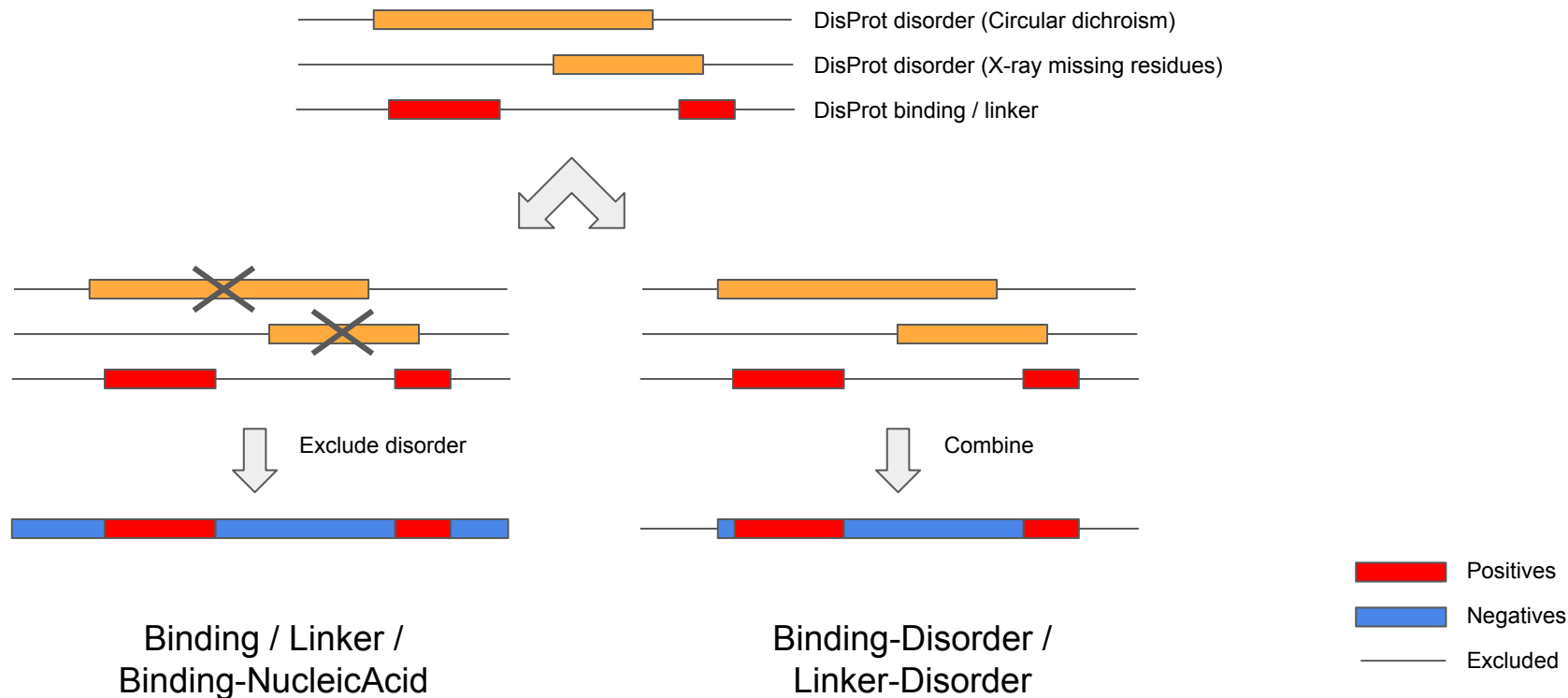
Merge

Remove X-ray annotations

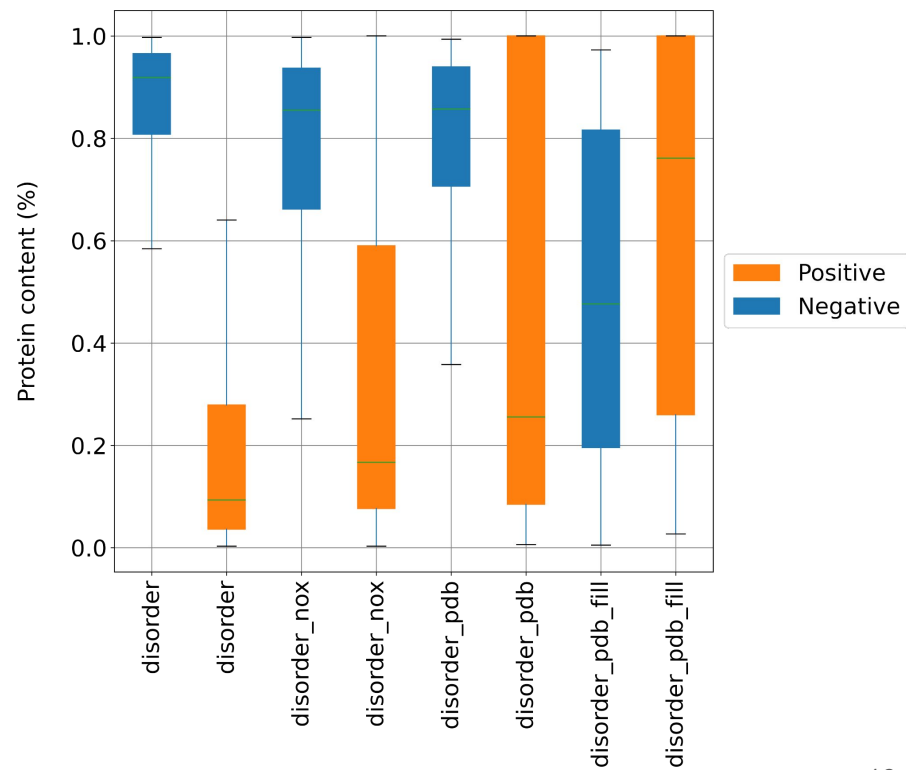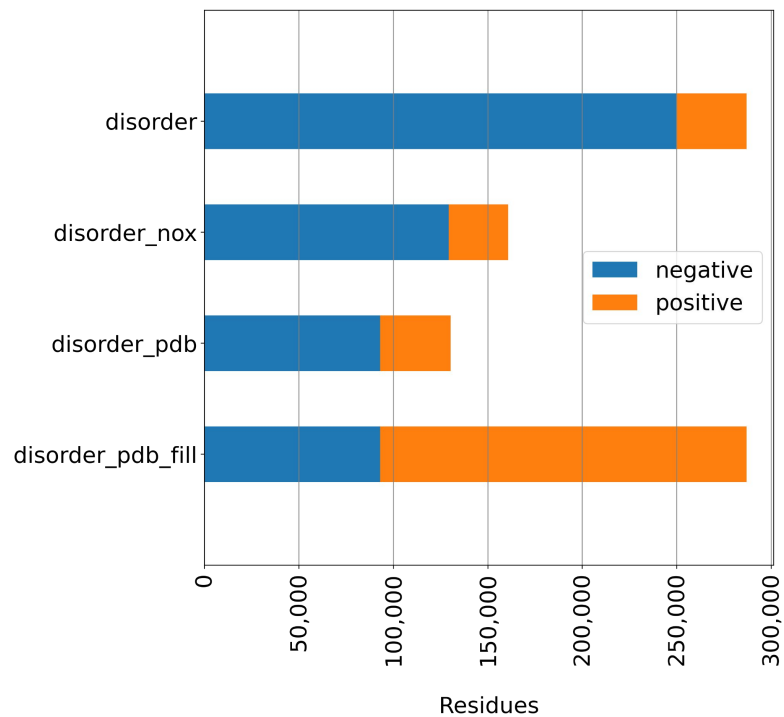Disorder

Disorder-NoX

Positives
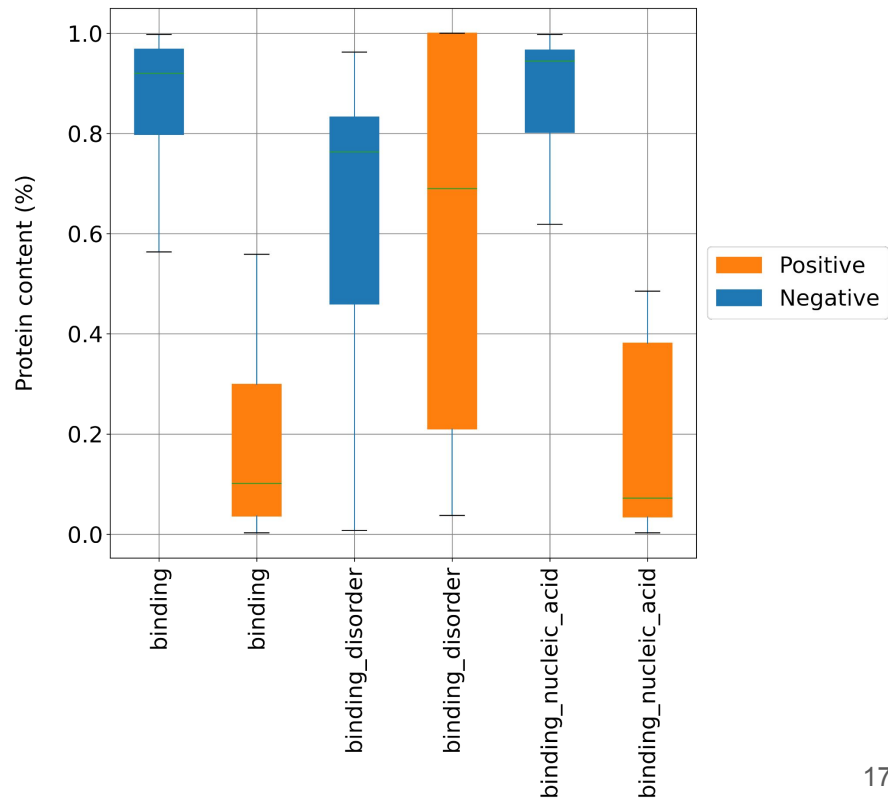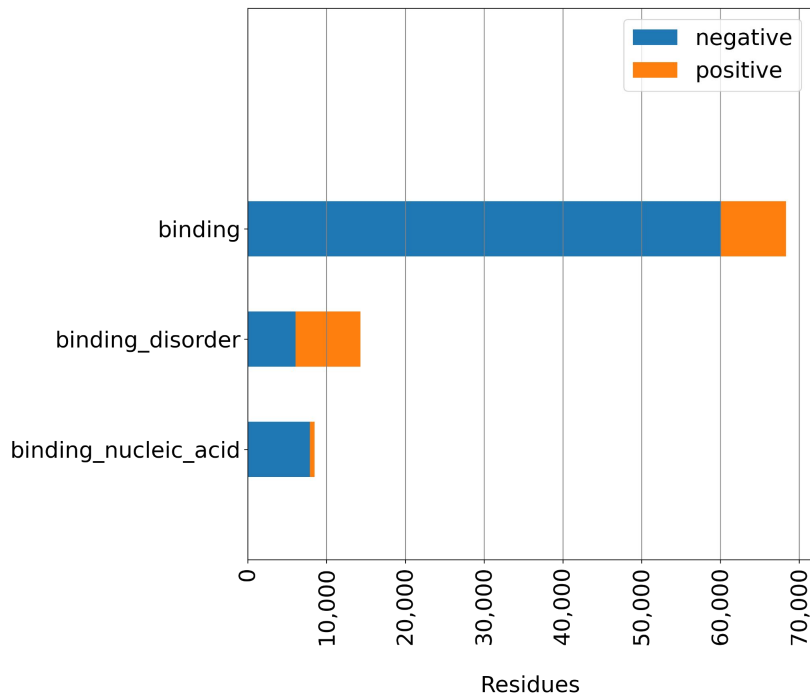
Negatives

# Reference definition - Disorder & PDB
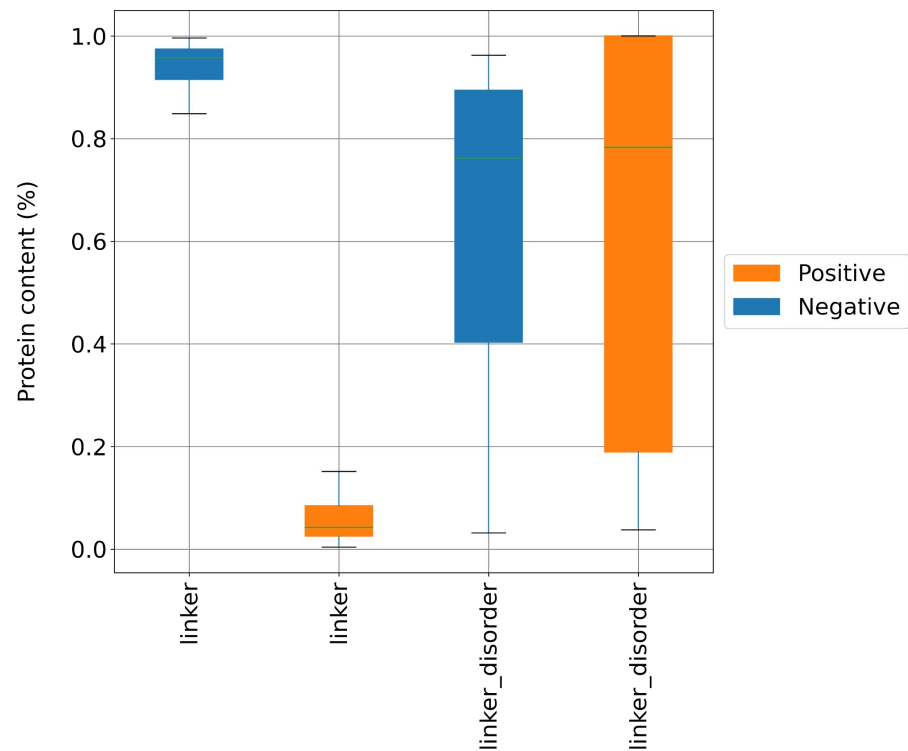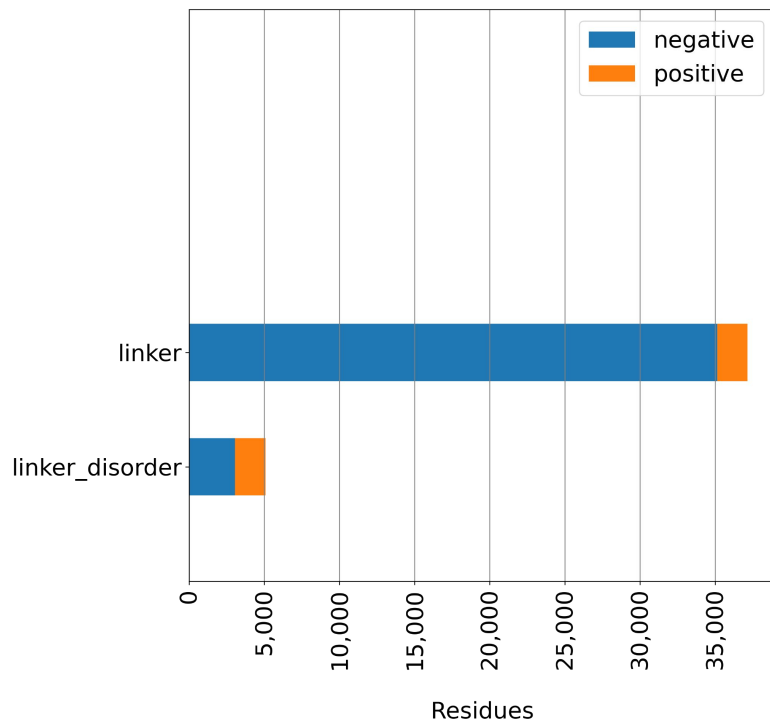
# Reference definition - Binding / Linker

# Class distribution - Disorder

# Class distribution - Binding
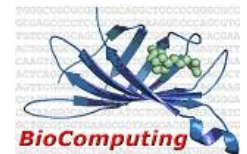
# Class distribution - Linker

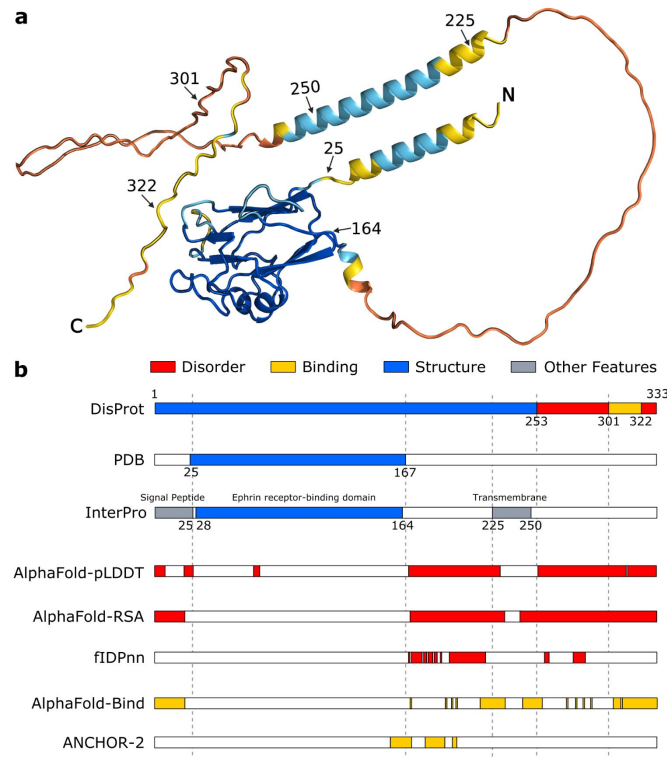# CAID-2

# Assessment

# Assessment

- Predictors **thresholds** are selected in order to optimize the **F1-score** in the considered benchmark

- Statistics are provided both at the **dataset** level or averaged over **targets**

- **Baseline**

  - **Random**

  - **Shuffled dataset** →Class imbalance at the dataset level is preserved

- Assessment **code**

  - **CodeOcean** capsule - https://codeocean.com/capsule/2223745/tree/v1

  - **GitHub** ("v2" branch → CAID2) - https://github.com/BioComputingUP/CAID
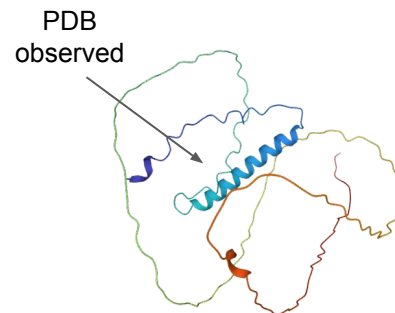
# AlphaFold & disorder

- AlphaFold-**disorder**
  - 1 - pLDDT

- AlphaFold-**RSA**
  - DSSP relative solvent accessibility)

- AlphaFold-**Binding**
  - $$\begin{cases} AlphaFold\_RSA, & AlphaFold\_RSA \leq T \\ T + pLDDT(1 - T), & AlphaFold\_RSA > T \end{cases}$$

Piovesan D, Monzon AM, Tosatto SCE. *Intrinsic protein disorder and conditional folding in AlphaFoldDB.* Protein Sci. 2022. 31(11):e4466
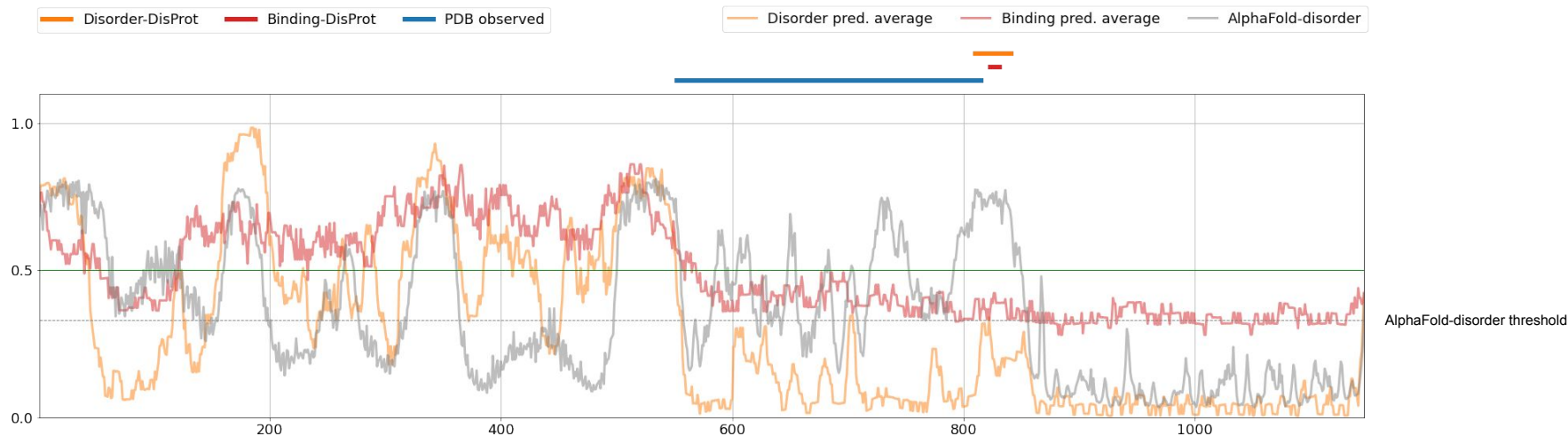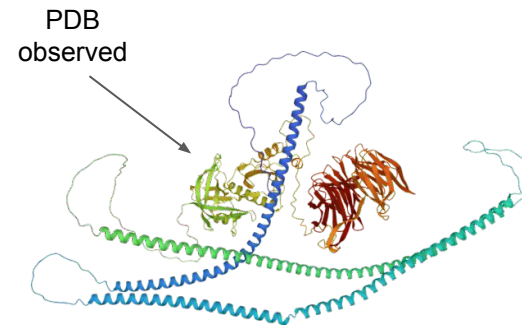


21

# DP02342 - P06837
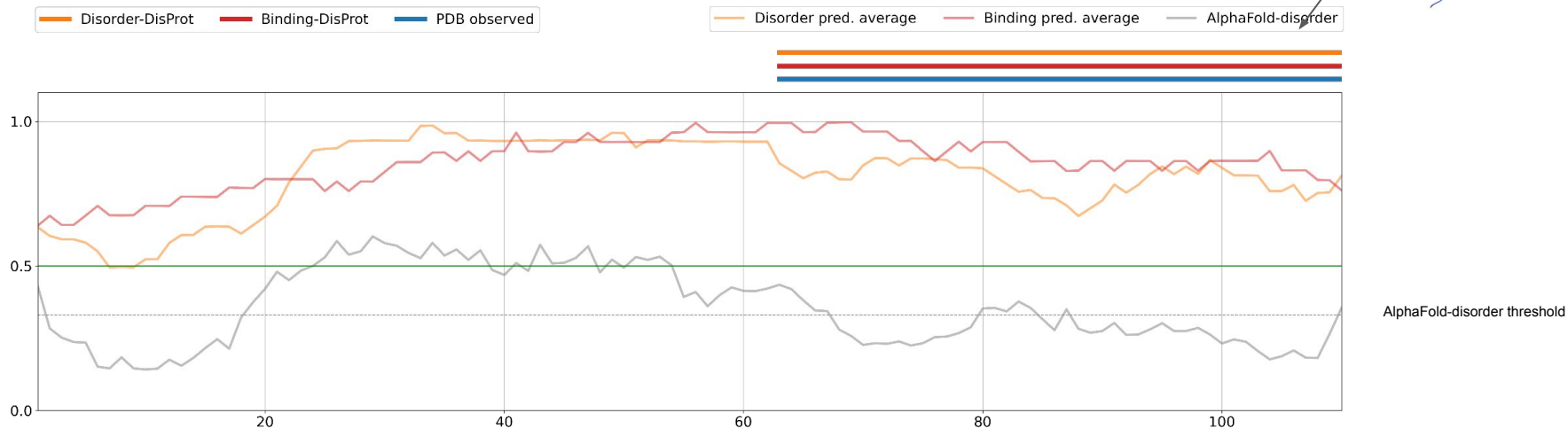
- *Neuromodulin*

- Fully disordered protein
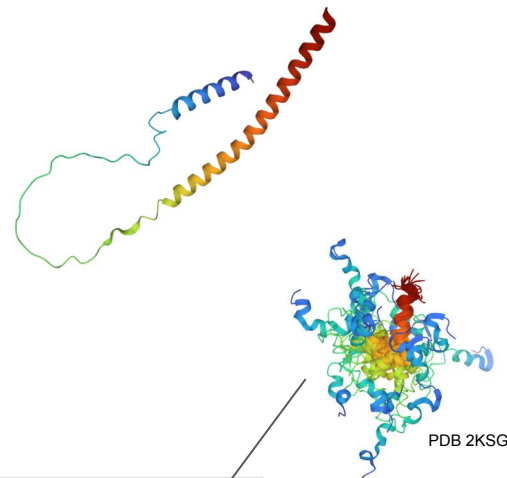
- Average $F_{max}$ ~ 0.96



PDB observed



AlphaFold-disorder threshold

# DP02959 - P42527

- *Myosin heavy chain kinase A*

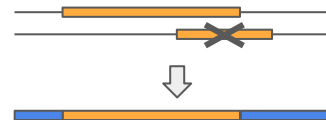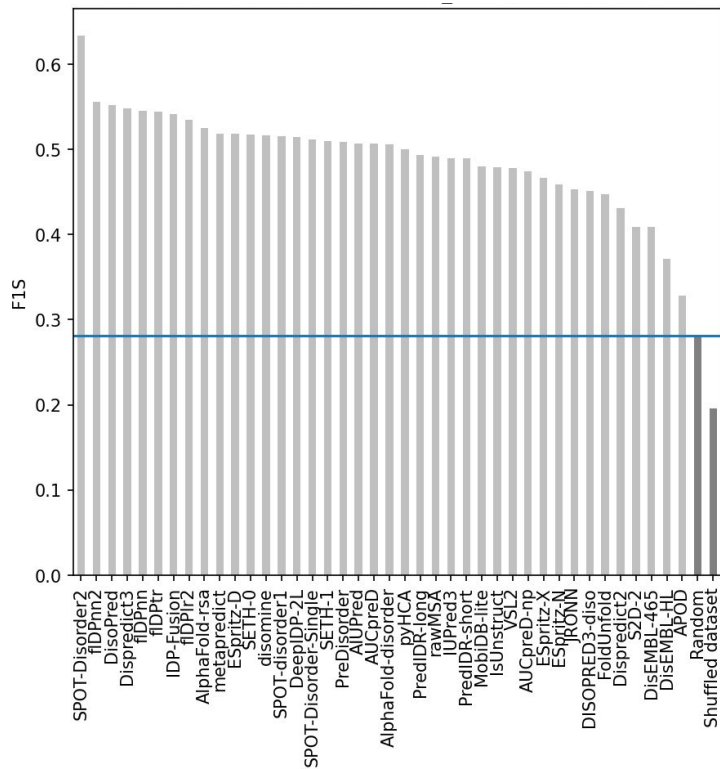- X-ray missing residues evidence

- Average disorder $F_{max}$ ~ 0.1


PDB observed



Legend: Disorder-DisProt | Binding-DisProt | PDB observed | Disorder pred. average | Binding pred. average | AlphaFold-disorder

AlphaFold-disorder threshold

# DP03635 - P81605

- *Dermcidin*

- Average $F_{max}$ ~ 0.74 (Disorder-PDB), ~ 0.45 (Disorder-NoX)

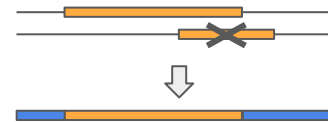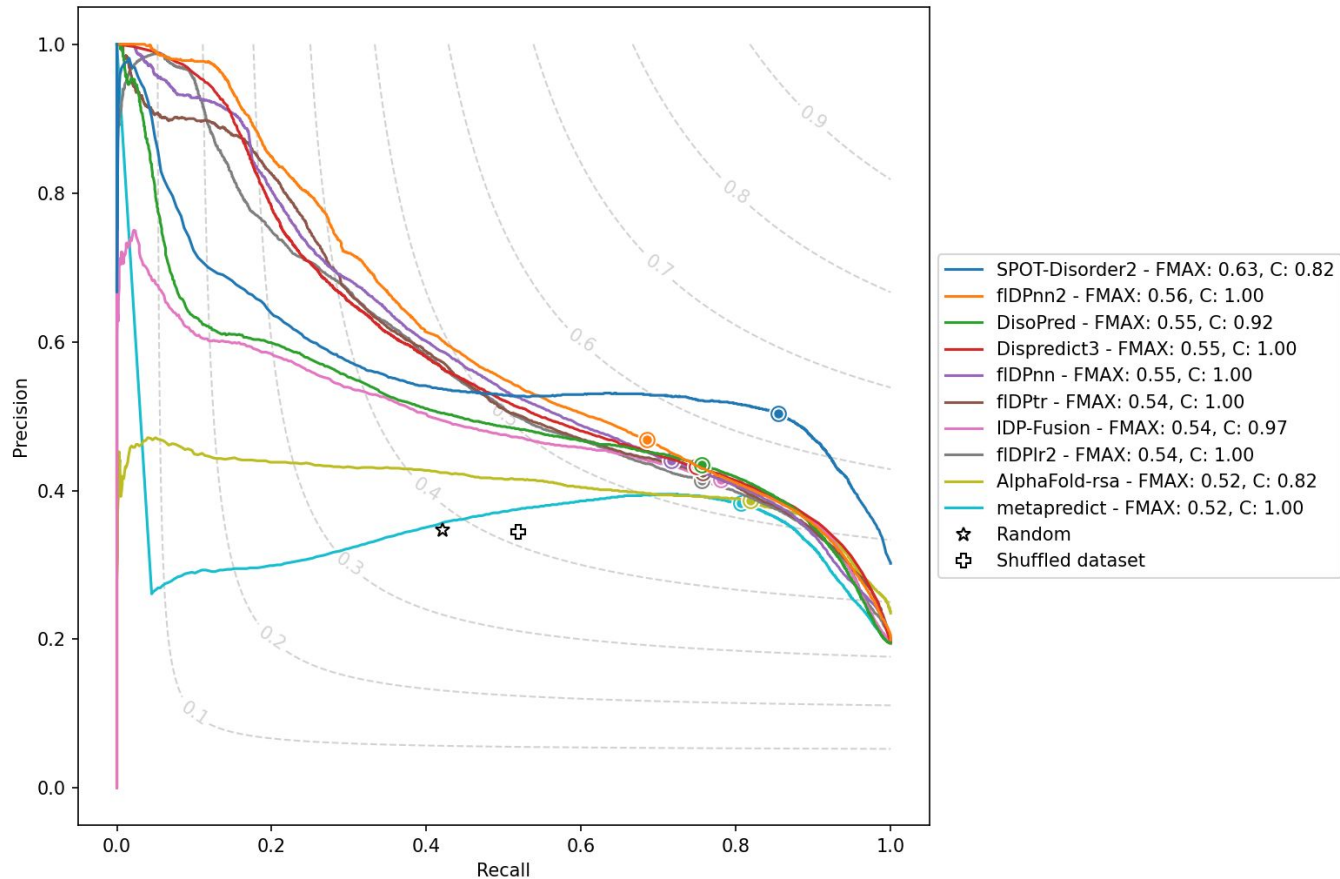- Circular dichroism and NMR evidence



PDB 2KSG



Disorder-DisProt  Binding-DisProt  PDB observed    Disorder pred. average    Binding pred. average    AlphaFold-disorder

AlphaFold-disorder threshold
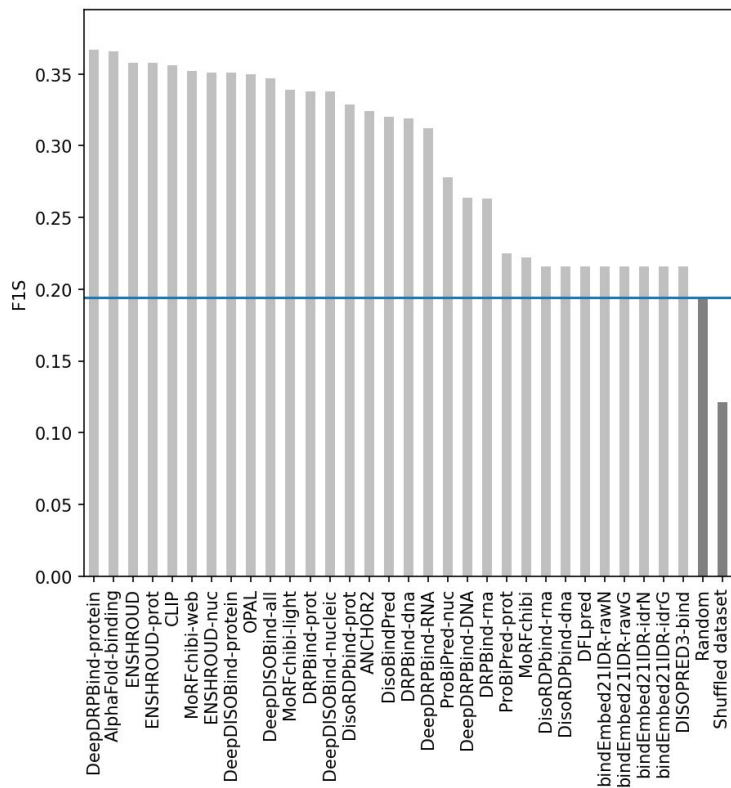
# **Disorder-NoX** (no X-ray)
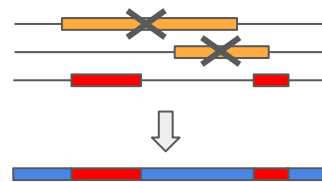
Dataset evaluation
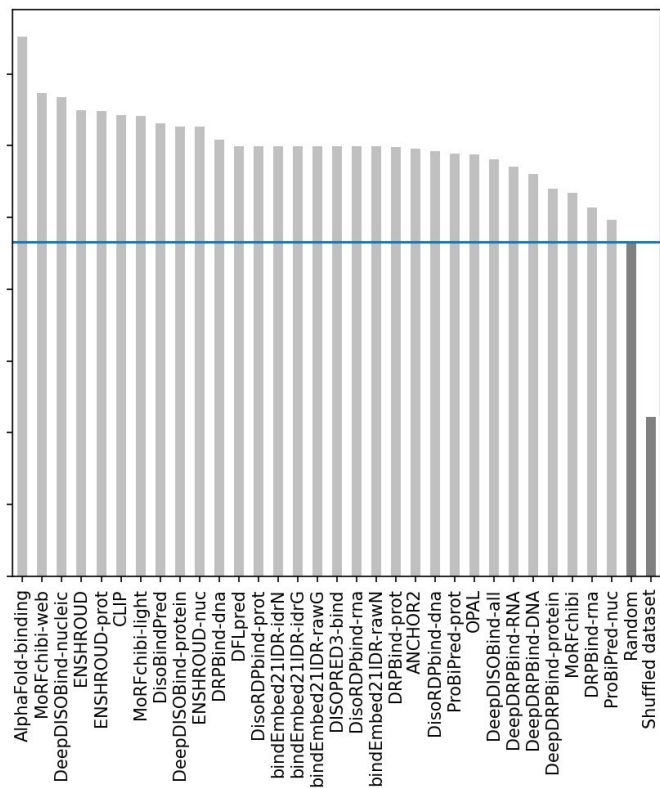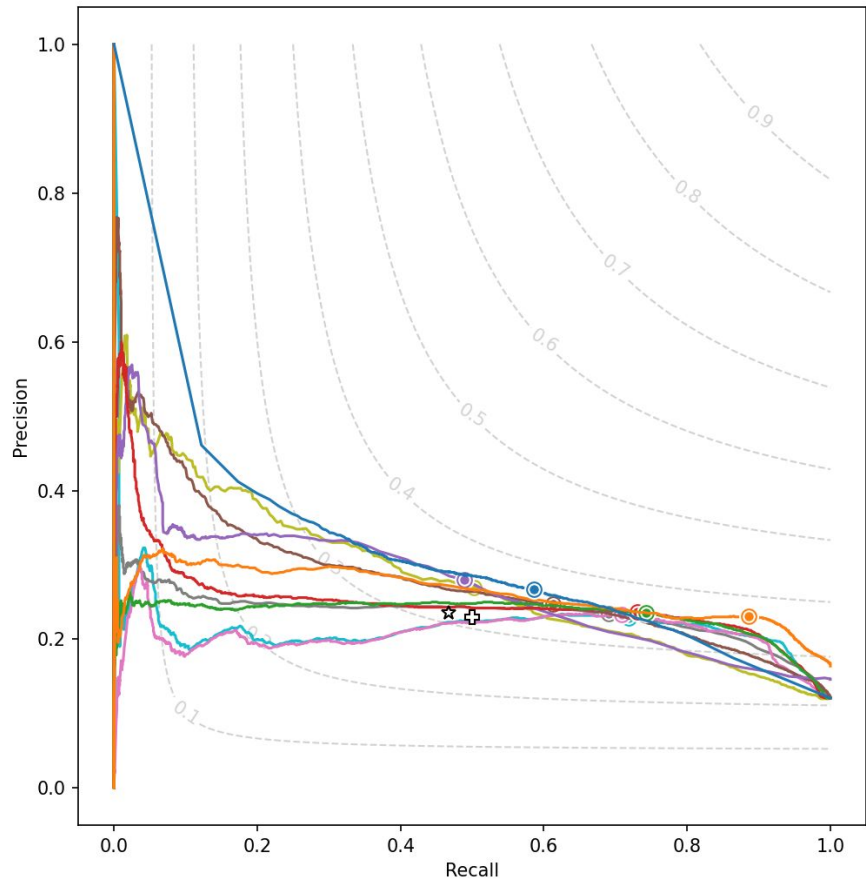
Target average

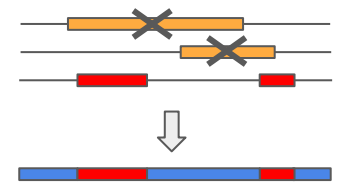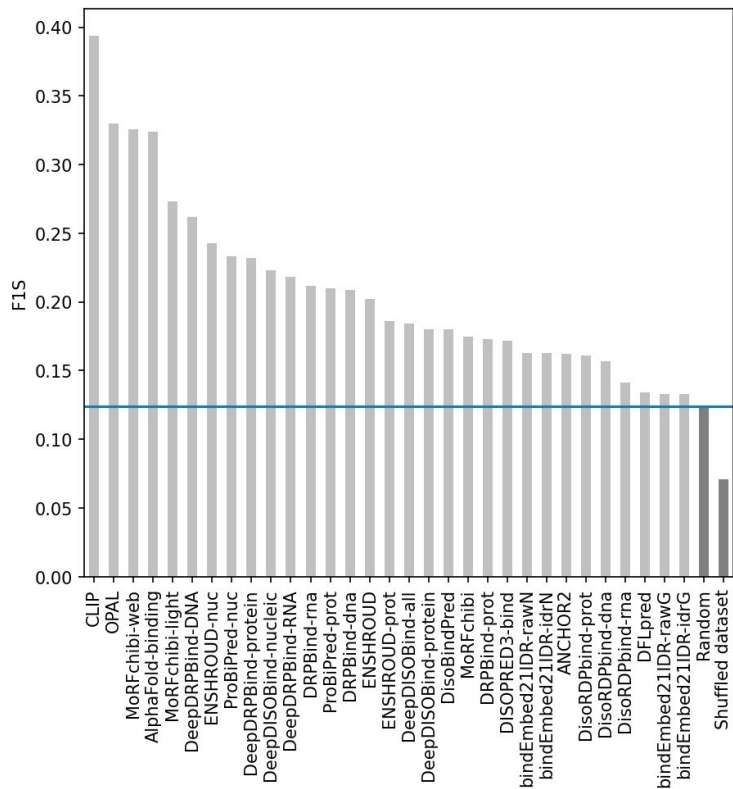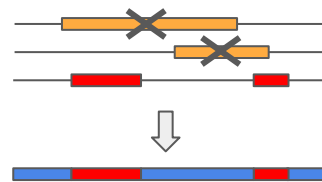# **Disorder-NoX** (no X-ray)
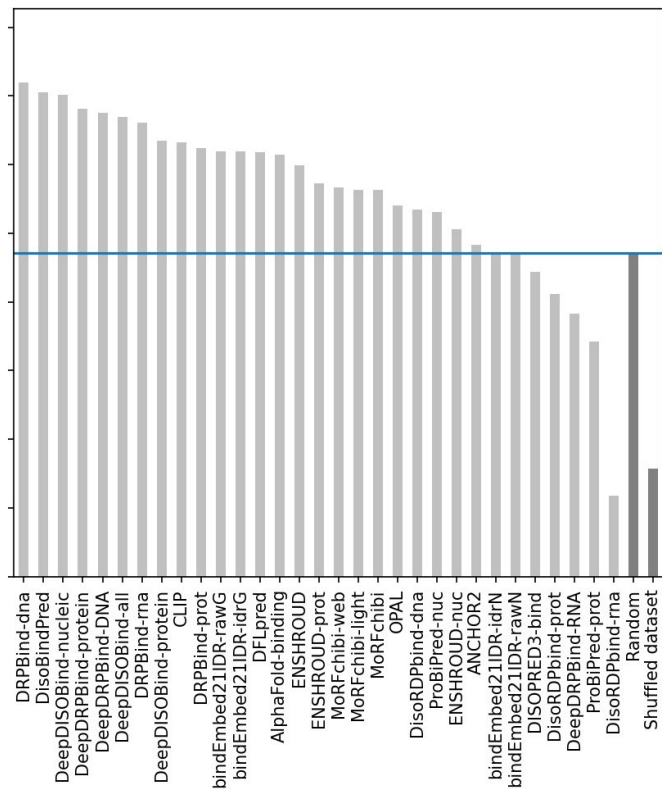
# Binding



Dataset evaluation

Target average

# Binding

# Binding-NucleicAcid
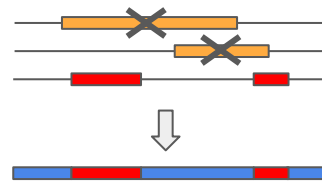


Dataset evaluation

Target average

# Binding-NucleicAcid



CLIP - FMAX: 0.39, C: 0.89
OPAL - FMAX: 0.33, C: 1.00
MoRFchibi-web - FMAX: 0.33, C: 1.00
AlphaFold-binding - FMAX: 0.32, C: 0.78
MoRFchibi-light - FMAX: 0.27, C: 1.00
DeepDRPBind-DNA - FMAX: 0.26, C: 1.00
ENSHROUD-nuc - FMAX: 0.24, C: 1.00
ProBiPred-nuc - FMAX: 0.23, C: 1.00
DeepDRPBind-protein - FMAX: 0.23, C: 1.00
DeepDISOBind-nucleic - FMAX: 0.22, C: 1.00
☆ Random
✚ Shuffled dataset
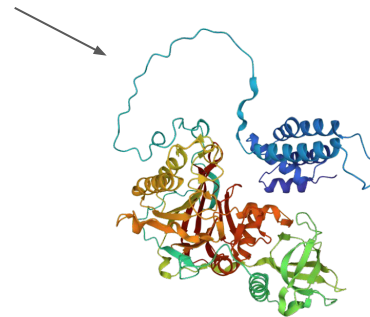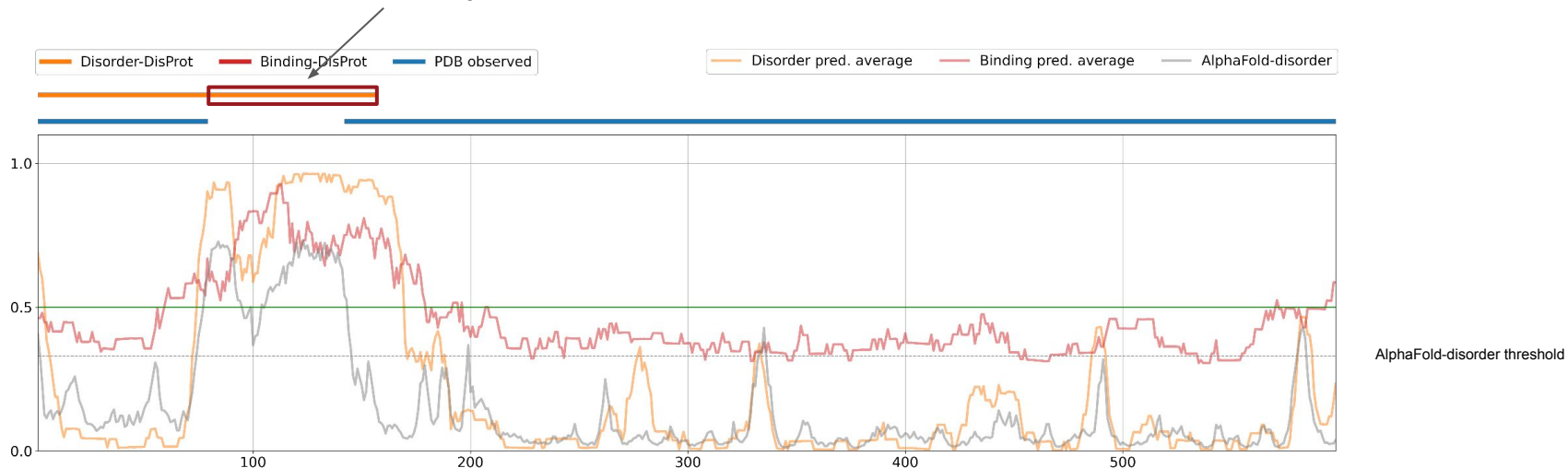
# DP02759 - Q14181

- *DNA Polymerase Alpha subunit B*

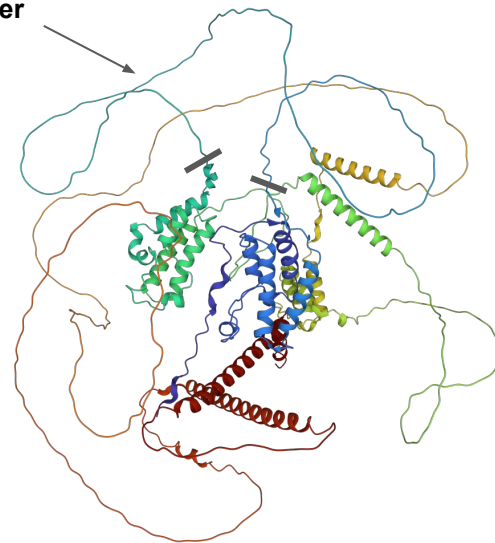- X-ray missing residues evidence

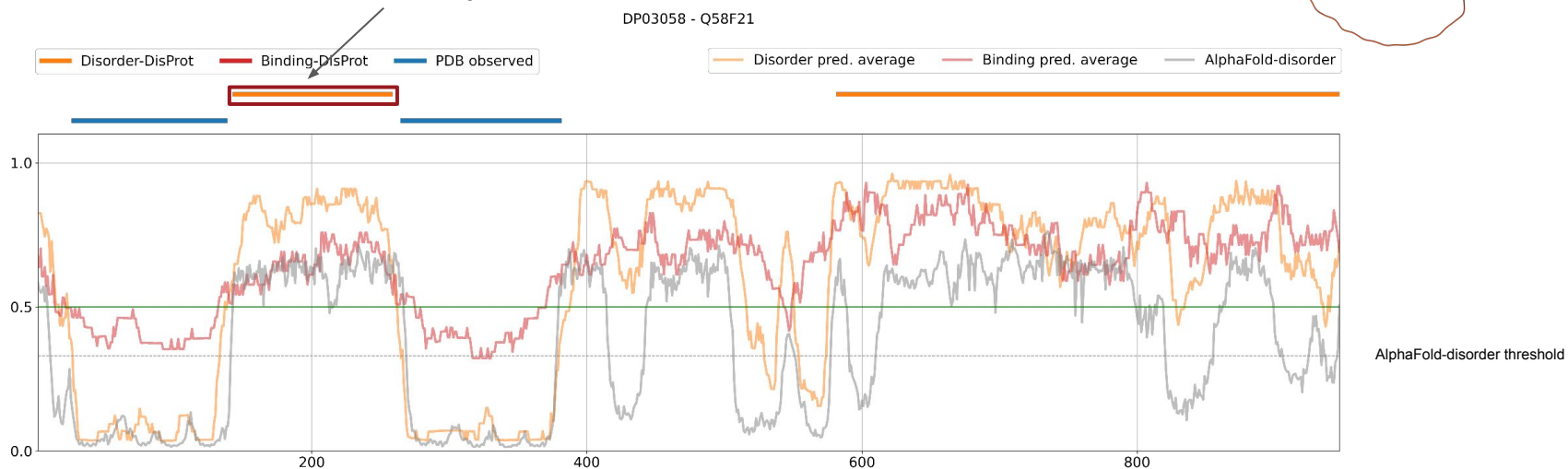- Average $F_{max}$ ~ 0.65



**Linker**

# DP03058 - Q58F21

- *Bromodomain testis-specific protein*

- X-ray missing residues and NMR evidence

- Average $F_{max}$ ~ 0.39





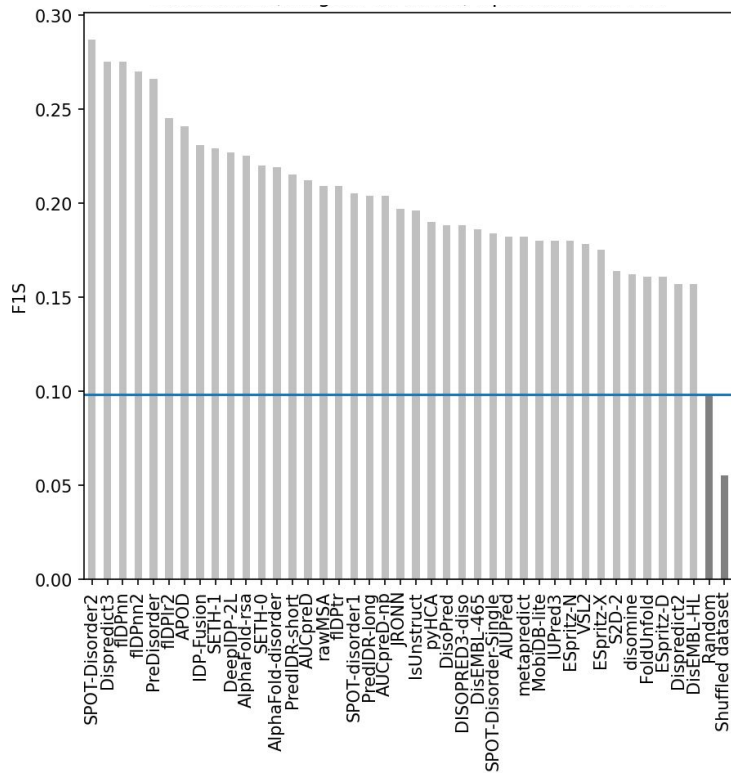DP03058 - Q58F21

**Linker**

Disorder-DisProt    Binding-DisProt    PDB observed      Disorder pred. average    Binding pred. average    AlphaFold-disorder
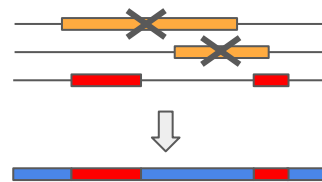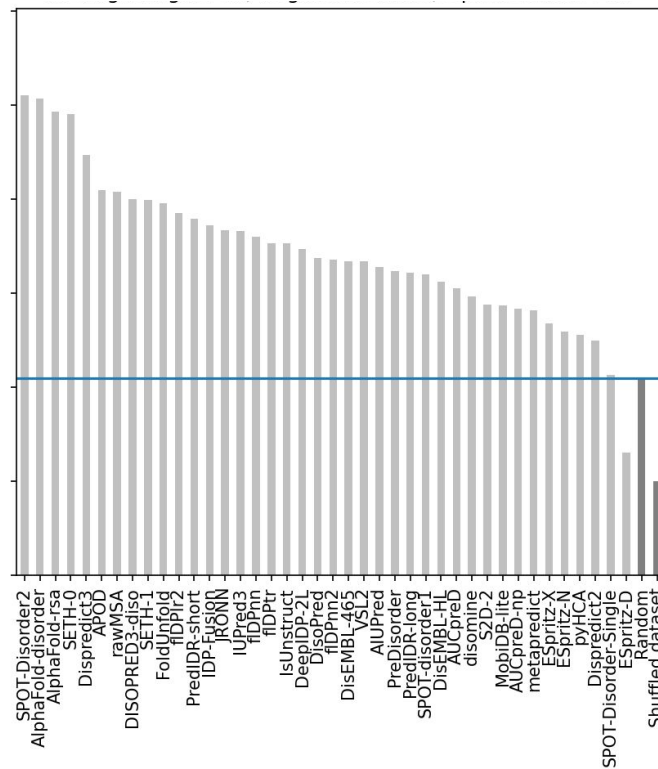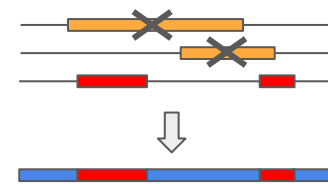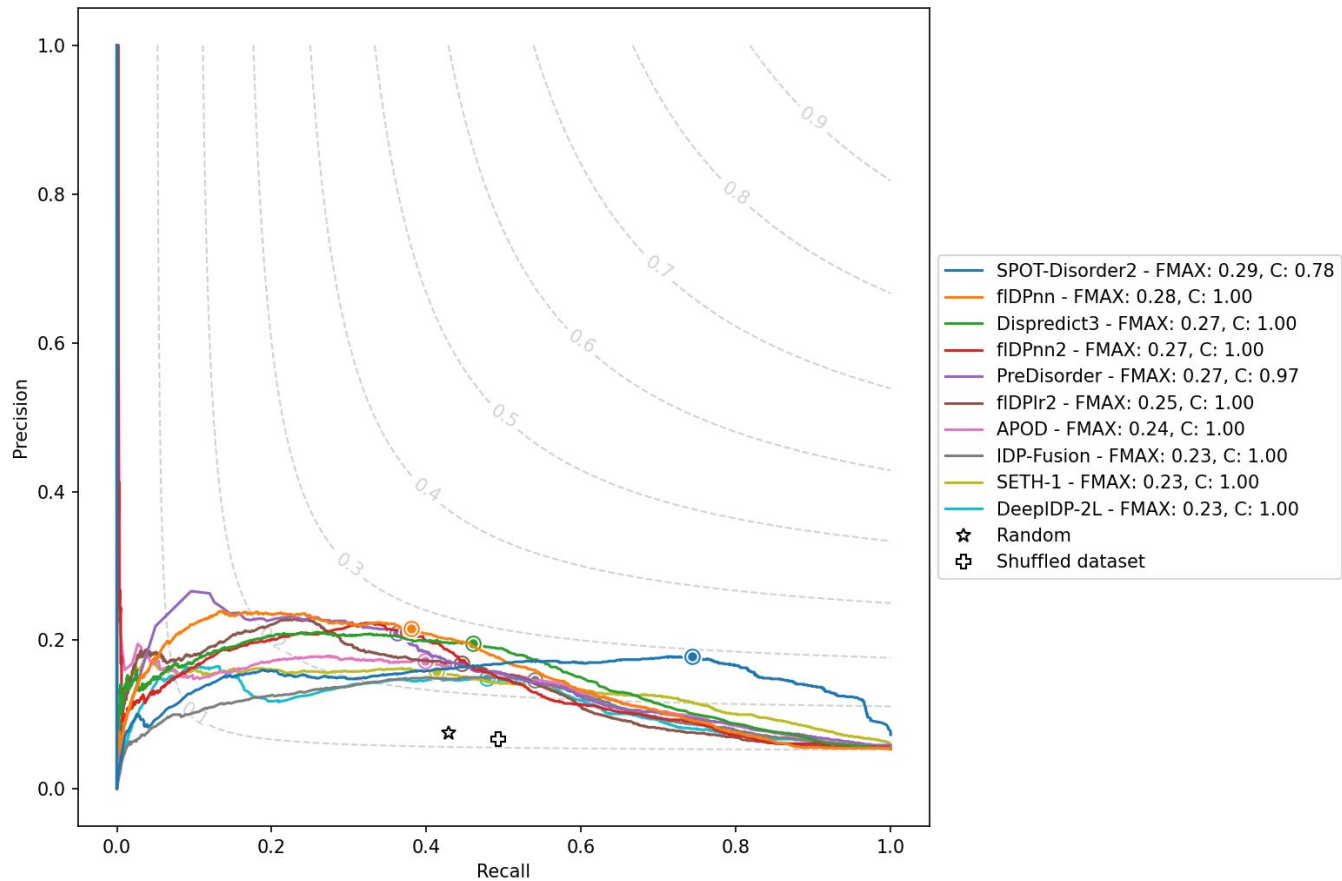
AlphaFold-disorder threshold

32

# Linker



Dataset evaluation

Target average

# Linker

# Acknowledgements
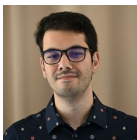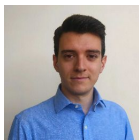
**Organizers**
Silvio CE Tosatto
Damiano Piovesan
Alexander M Monzon

**Curators**
Federica Quaglia
Victoria Nugnes
Maria Cristina Aspromonte

**Assessors**
Alessio Del Conte
Adel Bouhraoua

**BioComputingUP**
University of Padova
Biomedical Sciences

https://biocomputingup.it

@BioComputingUP