

Discovering the general architecture of protein families with OverProt

Adam Midlik, Ivana Hutařová Vařeková,
Jan Hutař, Aliaksei Charehneu,
Radka Svobodová, Karel Berka

ENBIK 2022, 13 June 2022, Němčice

MUNI
SCI

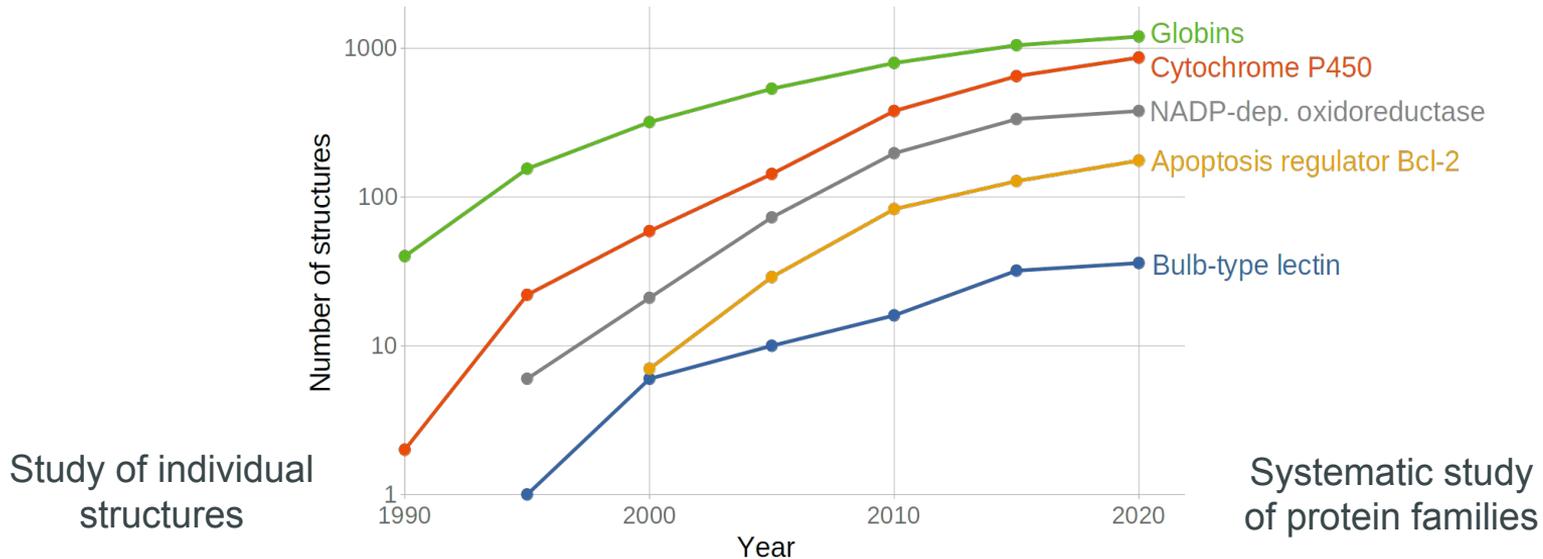
National Centre
for Biomolecular
Research



CEITEC

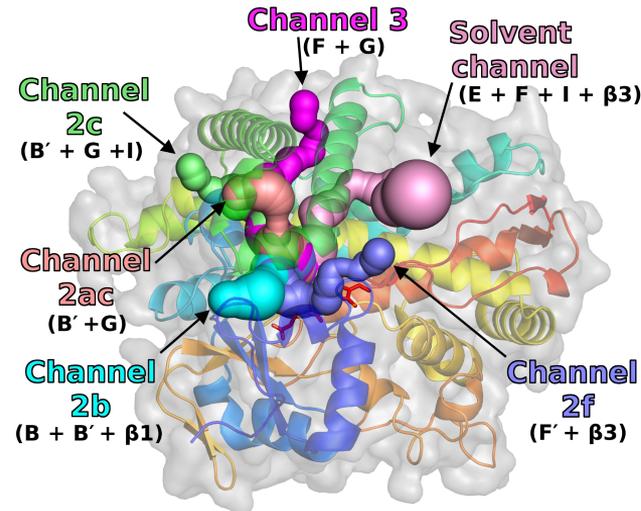
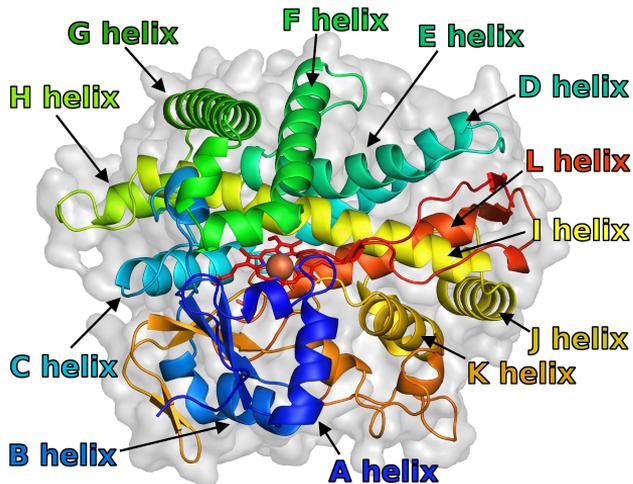
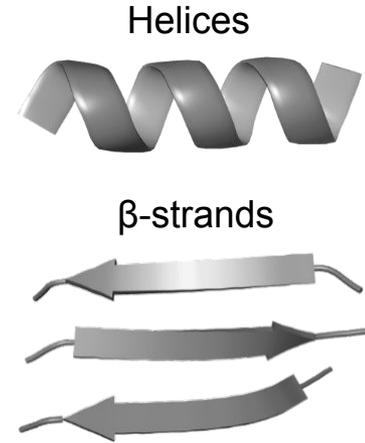
Introduction

- Protein structures → help us understand the function
 - Today > 190 000 experimental structures
- Similar protein structures form **protein families**
 - CATH database – 6 631 families
 - Families are getting bigger



Secondary structure elements (SSEs)

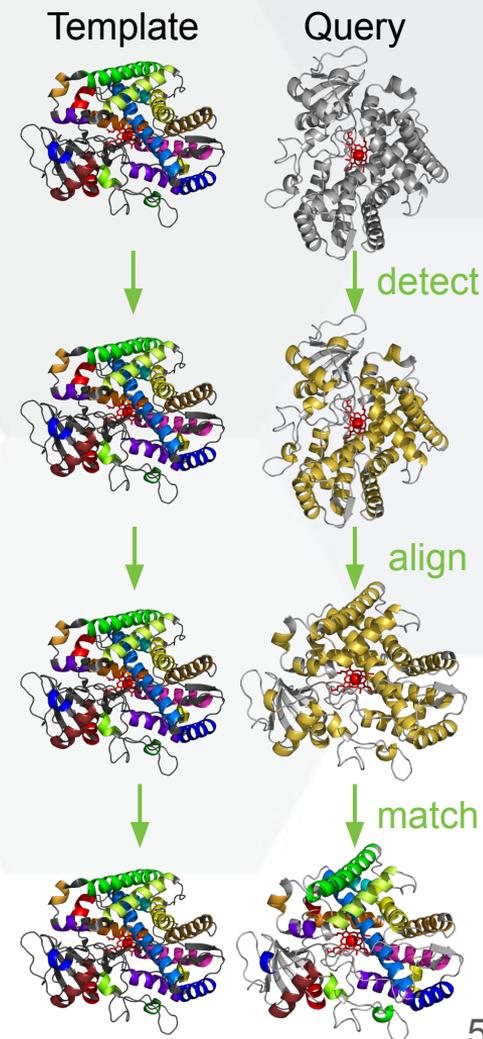
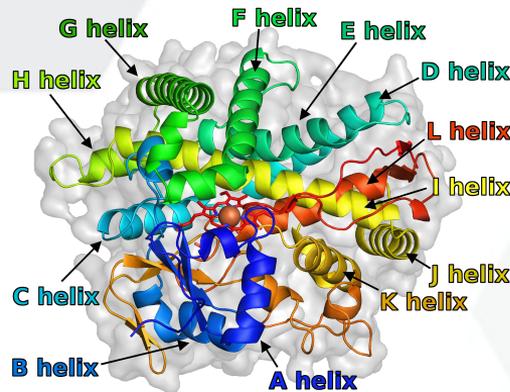
- Characteristic arrangement in each protein family
 - “General architecture” of the family
- Can serve as landmarks
 - Help us orient in the structures
 - Help us locate the key regions (active sites, channels...)



SecStrAnnotator

SSE annotation

- = Labelling of SSEs, consistent throughout the family
- Many protein families have traditional annotations
 - Performed manually :(
- Our solution – **SecStrAnnotator**
 - Annotation based on a provided template



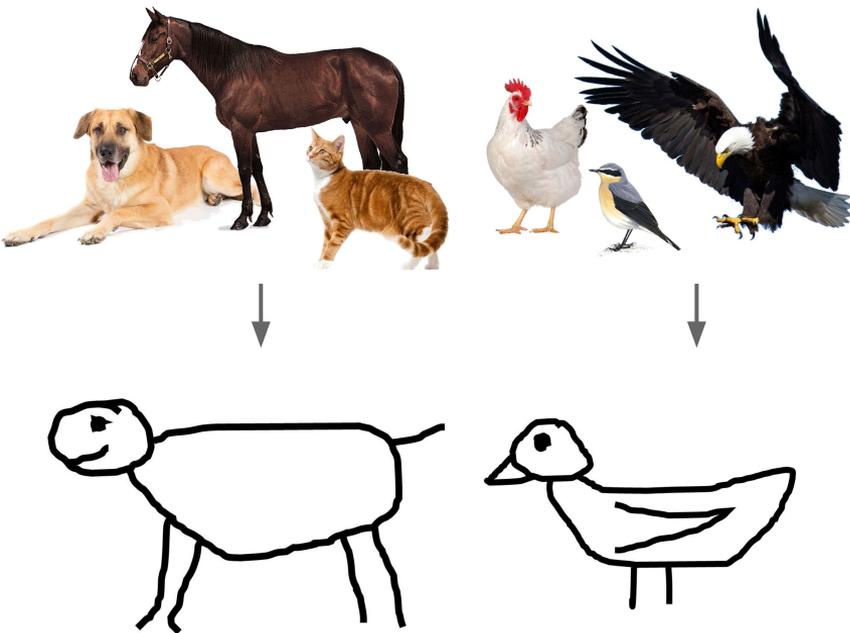
Midlik A. et al. (2019) Automated family-wide annotation of secondary structure elements. In Kister A.E. (ed.) *Protein Supersecondary Structures*, vol. 1958, pp. 47–71.

Midlik A. et al. (2021) Uncovering of cytochrome P450 anatomy by SecStrAnnotator. *Scientific Reports*, **11**, 12345.

OverProt

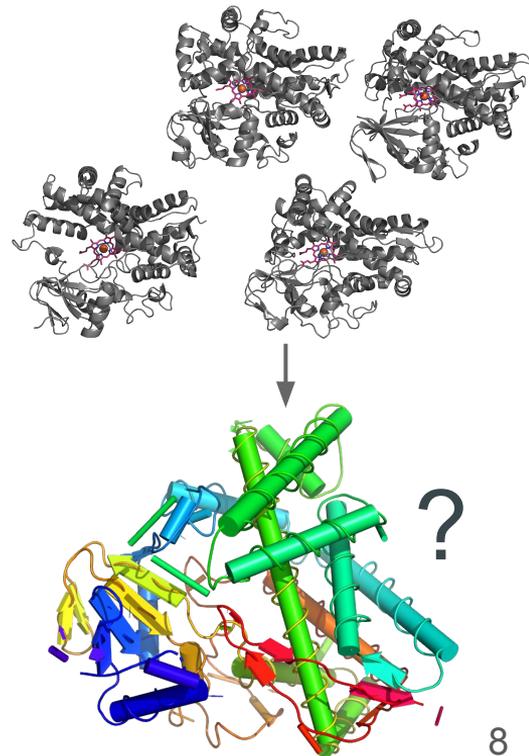
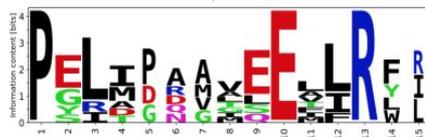
What is the “general architecture” of a family?

- Task: extract and visualize the characteristic features of a family



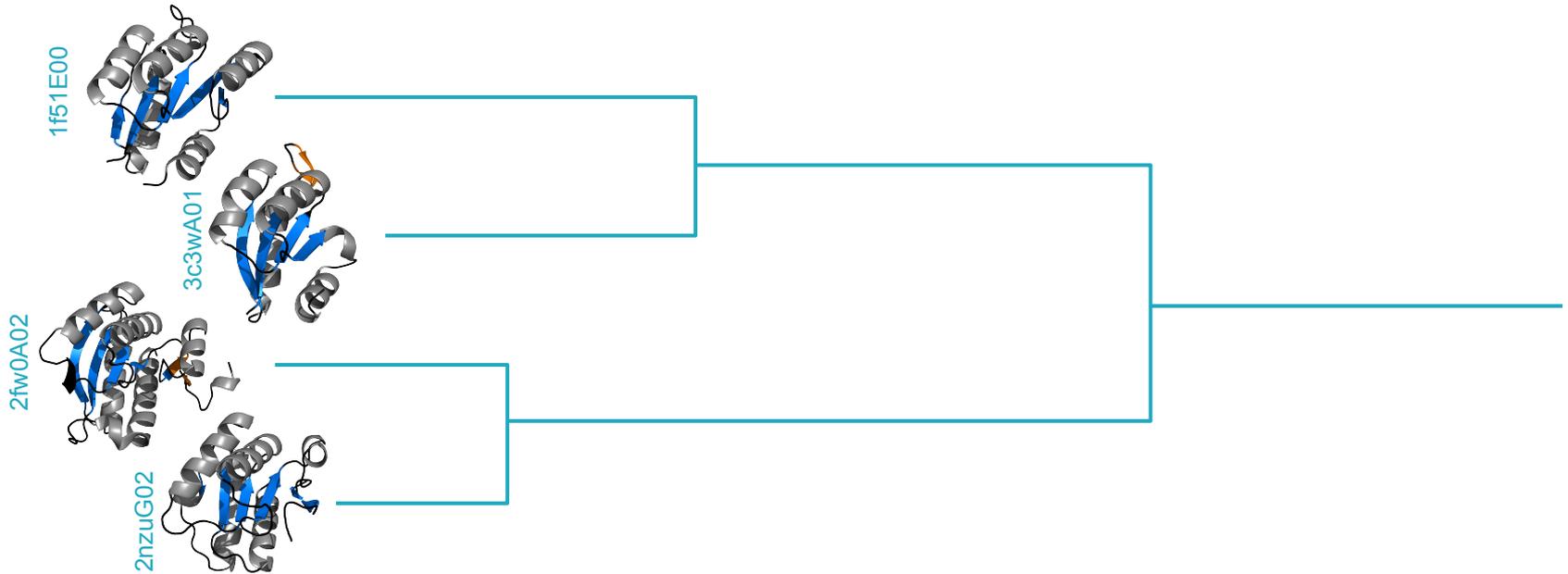
PGLIPAMVSETIRWL
ERIDAATEELLR
PELANAAEEVLRW
PYLDQVMQEILRLI
PELMGEEALRFR

PGLIPAMVSETIRWL
-ERIDAATEELLR--
PELA-NAAEEVLRW-
PYLD-QVMQEILRLI
PELM---GEEALRFR



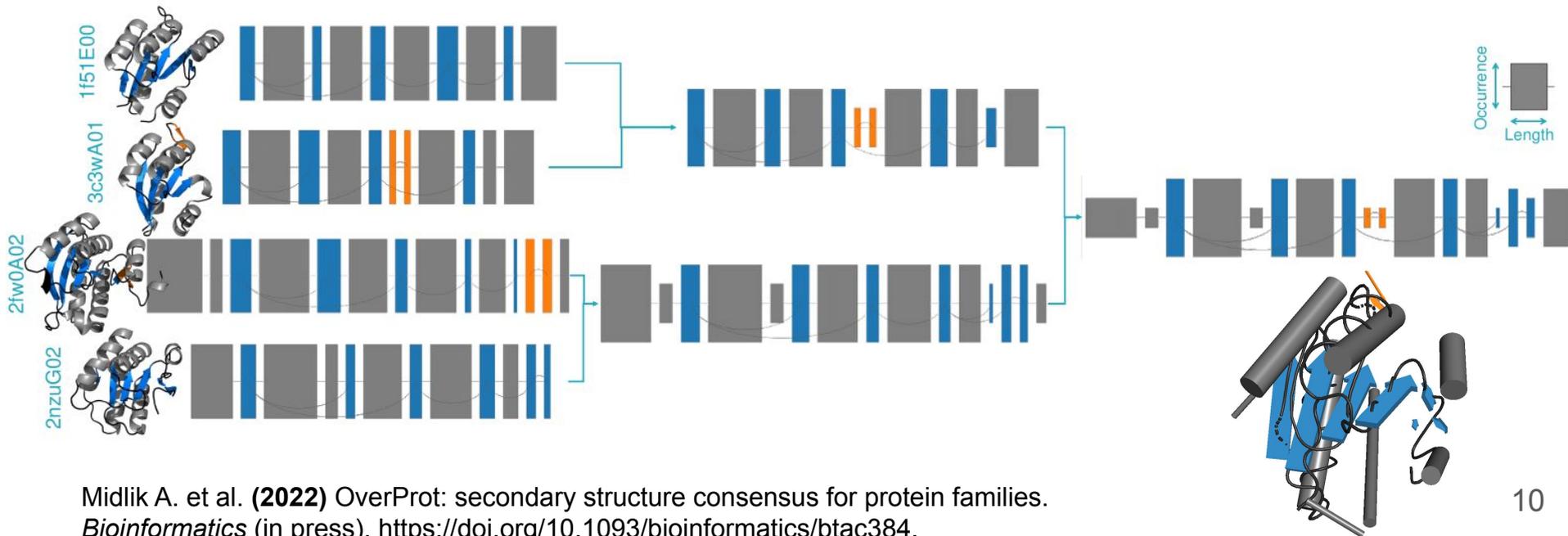
Our solution: OverProt

- Step 1: Align
- Step 2: Build the guide tree (agglomerative clustering of 3D domains)



Our solution: OverProt

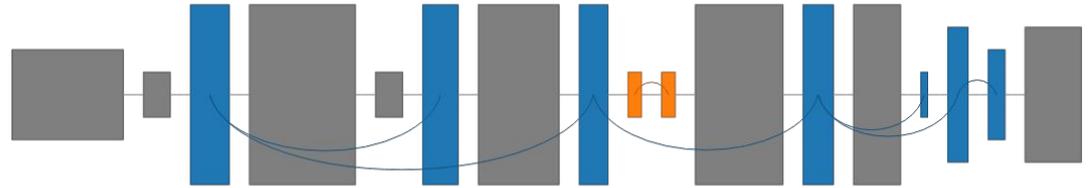
- Step 1: Align
- Step 2: Build the guide tree (agglomerative clustering of 3D domains)
- Step 3: Merge the SSEs through the guide tree



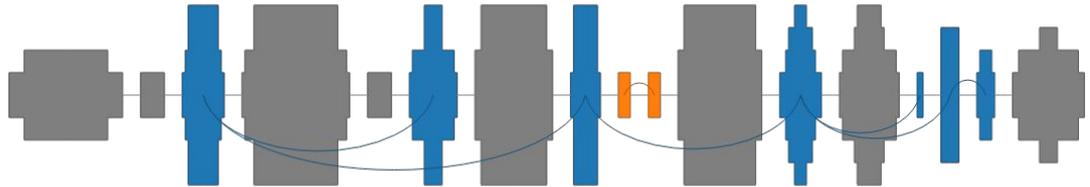
Midlik A. et al. (2022) OverProt: secondary structure consensus for protein families. *Bioinformatics* (in press). <https://doi.org/10.1093/bioinformatics/btac384>.

Alternative visualizations

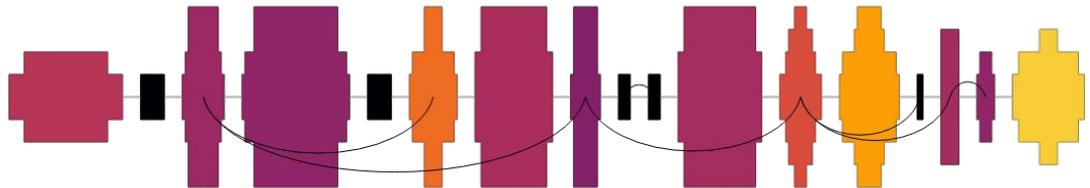
- Classic



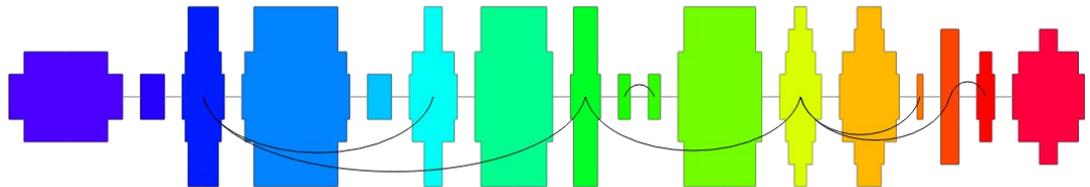
- Length distribution



- 3D variability



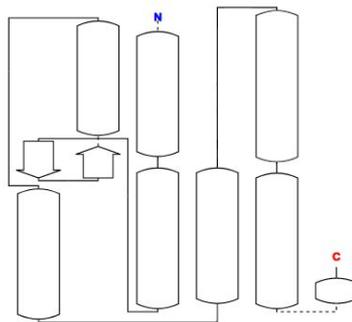
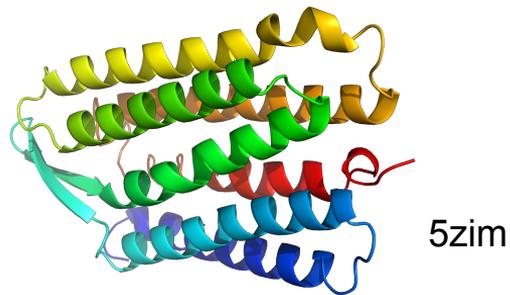
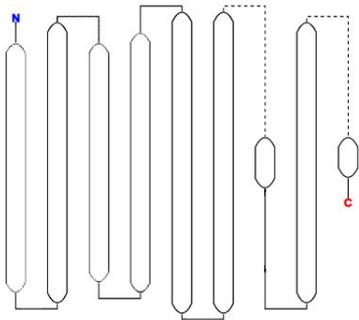
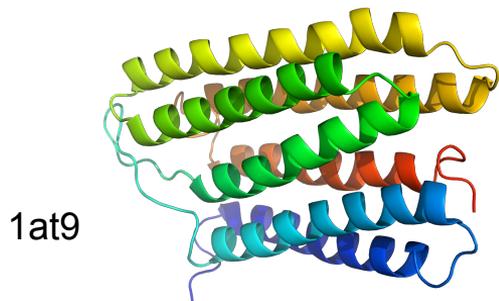
- Chainbow



2DProts

SSE visualization in 2D diagrams

- Existing diagrams in PDBe:
 - Ignore 3D arrangement
 - Ignore protein similarity

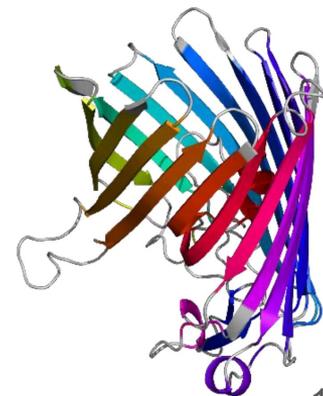
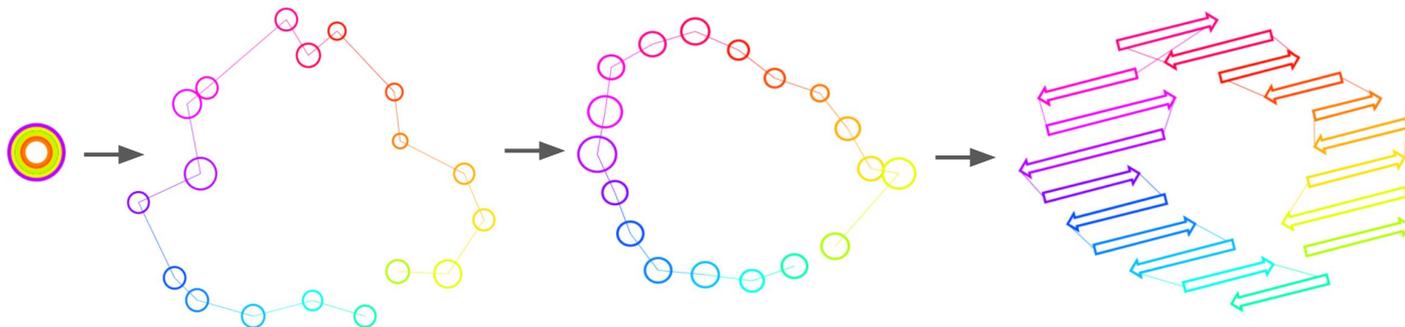


Our solution: 2DProts

- Preparation
 - OverProt creates SSE consensus for the family → annotation template
 - SecStrAnnotator annotates all protein domains in the family

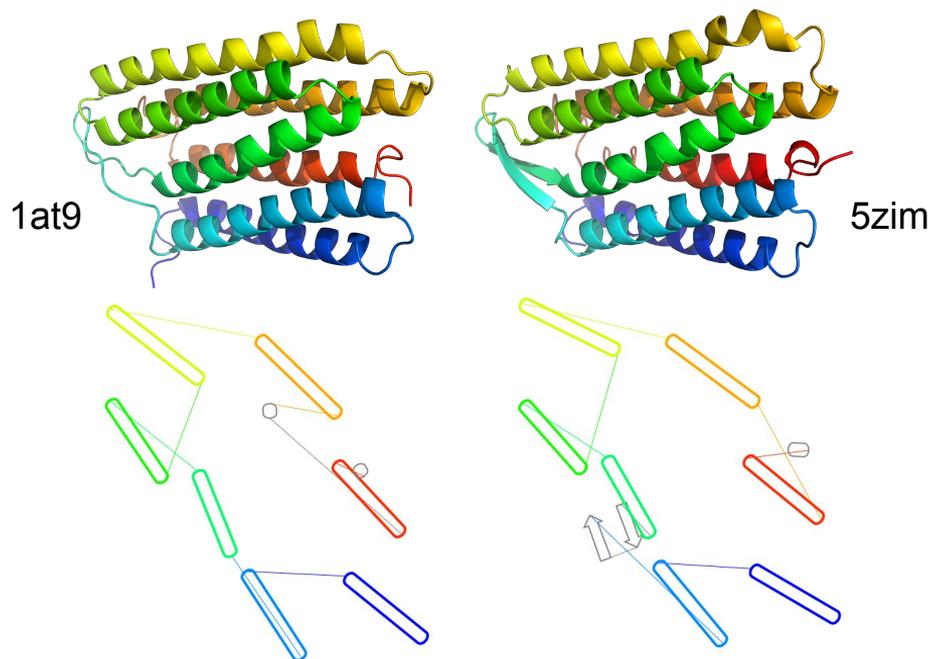
- Objective function $F = \sum_{s_1, s_2 \in S_p} |d_{3D}(s_1, s_2) - d_{2D}(s_1, s_2)| (L(s_1) + L(s_2)) + \sum_{s \in S_p \cap S_0} |r(s) - r_0(s)| L(s) \frac{|S_p|}{20}$
 - Deviation of 2D vs 3D
 - Deviation of 2D vs family template 2D

- Minimize F

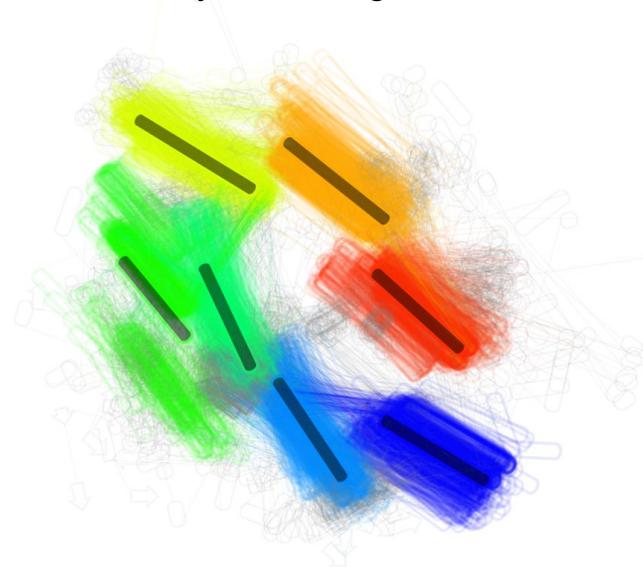


Example family

- 1.20.1070.10 Rhodopsin 7-helix transmembrane proteins



Family multi-diagram



Availability

SecStrAnnotator

Desktop version

gitlab.com/midlik/SecStrAnnot2
(all platforms via .NET 6.0)

Online version

sestra.ncbr.muni.cz
– Cytochrome P450 template
– Custom template
– Detect SSEs (no annotation)

Database

1 family
1 855 domains

Integrations

OverProt

gitlab.com/midlik/overprot
(all platforms via Docker/Podman)

overprot.ncbr.muni.cz

Job calculations
up to 500 protein domains

6 631 families
> 470 000 domains
(2.5 CPU days)

2DProts

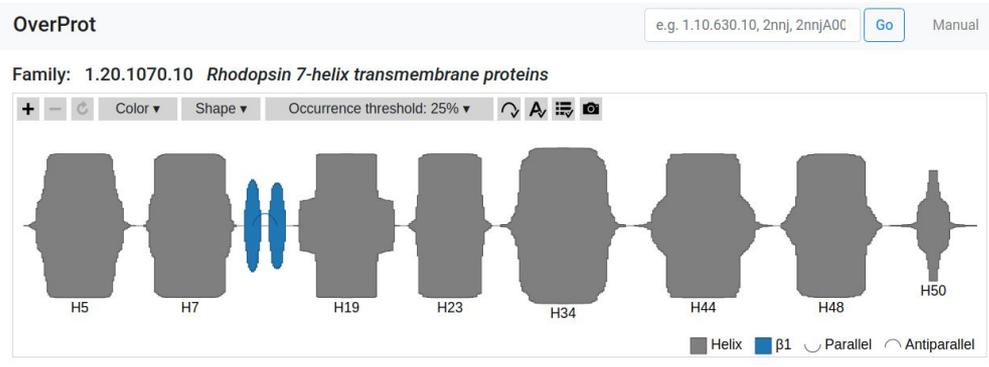
2dprots.ncbr.muni.cz

Job calculations
up to 30 protein domains

6 631 families
> 470 000 domains
(25 CPU days)

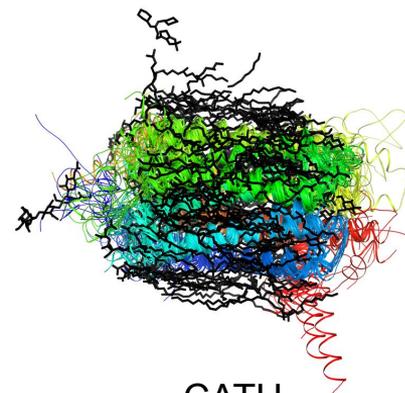
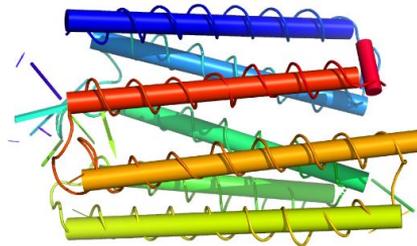
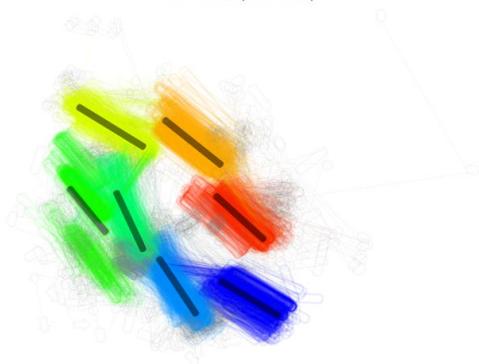
CATH
cathdb.info

All together – overprot.ncbr.muni.cz



2D view (2DProts)

3D view (MAPSCI + OverProt)



CATH
superimposition

Family info

PDB entries: 412 ([List](#)) [ⓘ](#)

Domains: 606 ([List](#)) [ⓘ](#)

Included domains: 412 ([List](#)) [ⓘ](#)

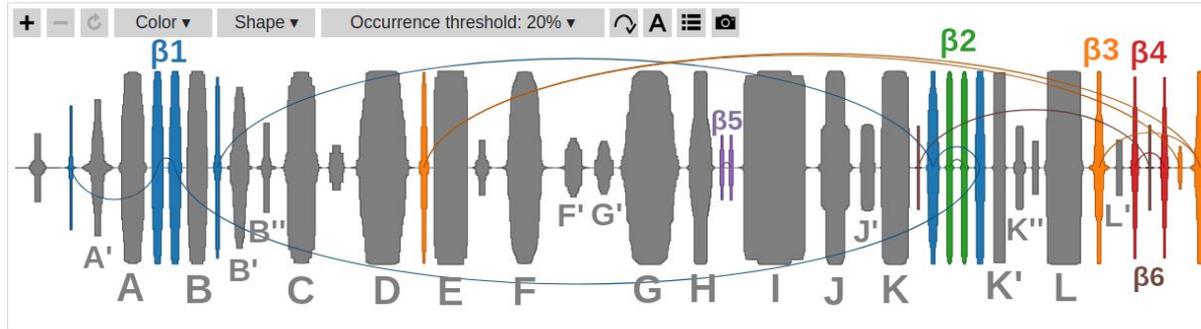
Example domain: 5zimA00 [ⓘ](#)

External links: [CATH](#) [2DProts](#)

Download: [results.zip](#)

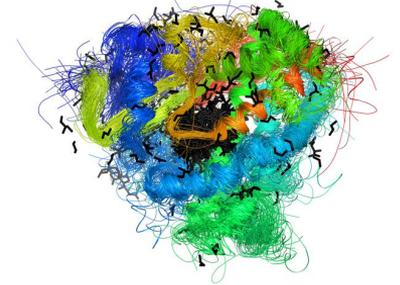
All together – overprot.ncbr.muni.cz

Family: 1.10.630.10 Cytochrome P450

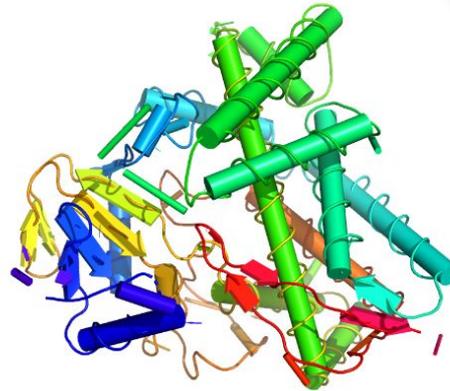
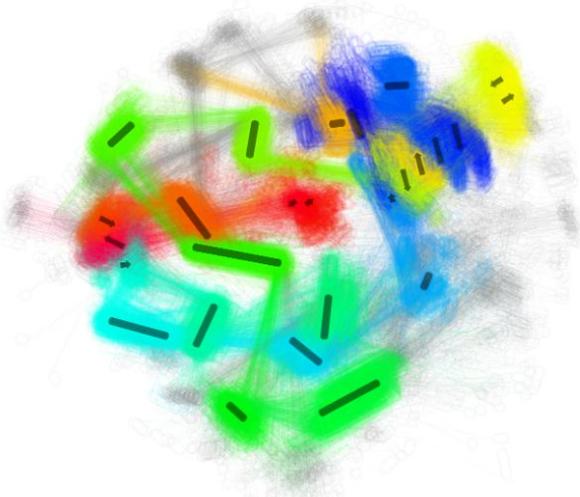


2D view (2DProts)

3D view (MAPSCI + OverProt)

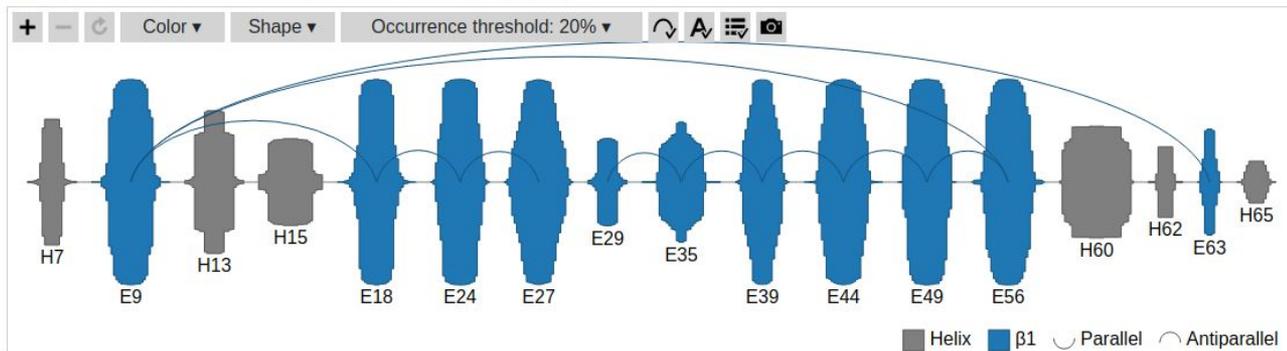


CATH
superimposition



All together – overprot.ncbr.muni.cz

Family: 2.40.128.20 *Calycin beta-barrel core domain*

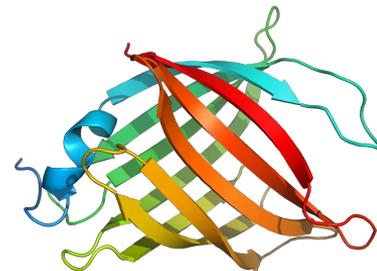


2D view (2DProts)

3D view (MAPSCI + OverProt)



8-stranded subgroup (3k3l)

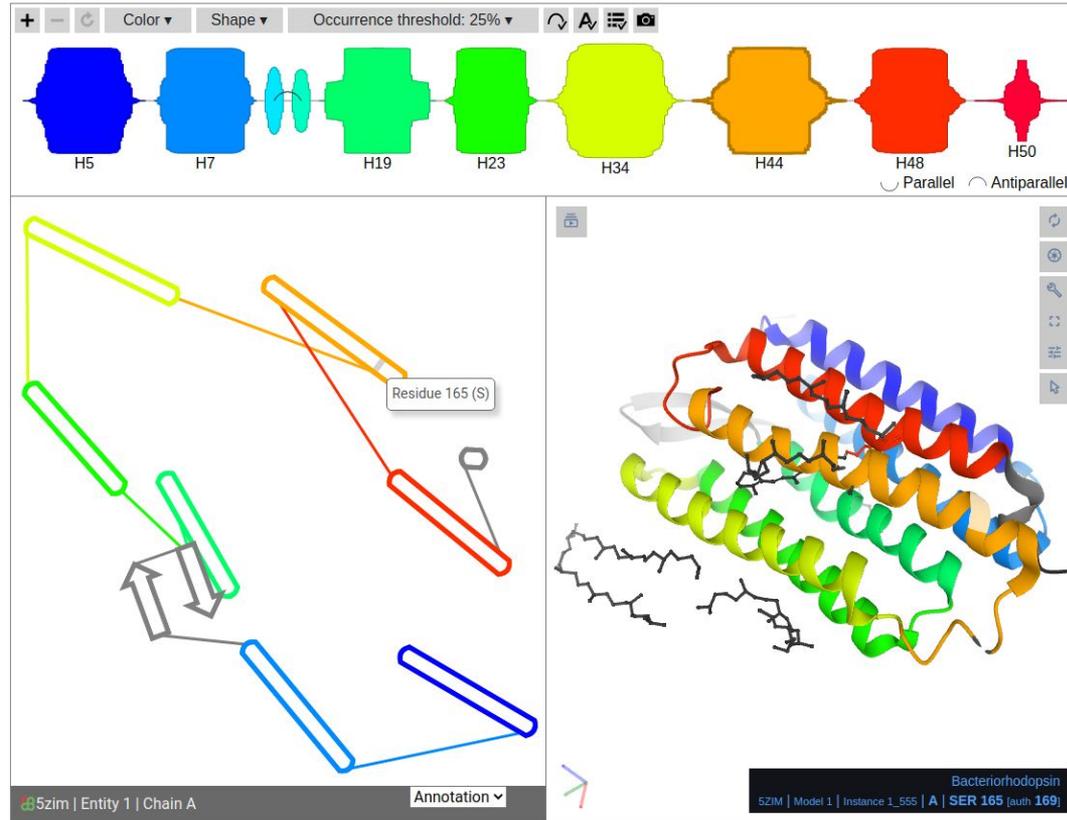


10-stranded subgroup (3wjb)

Integrated view

Family: 1.20.1070.10 *Rhodopsin 7-helix transmembrane proteins*

Domain: 5zimA00



← Family consensus

← Selected domain

Conclusion

- **SecStrAnnotator** – SSE annotation
- **OverProt** – SSE consensus + 1D visualization
- **2DProts** – SSE 2D visualization

Future plans

- **OverProt** – Integration with PDBe-KB, CATH
- **2DProts** – Inclusion of ligands
- Integration of AlphaFold DB via CATH cooperation

Acknowledgement



Central European Institute of Technology
BRNO | CZECH REPUBLIC

Structural bioinformatics group, Masaryk University

- Ivana Hutařová Vařeková, Jan Hutař – **2DProts**
- Aliaksei Chareshneu – **Integrated view**
- Radka Svobodová, Jaroslav Koča

MUNI National Centre
SCI for Biomolecular
Research

MUNI
FI

Palacký University

- Karel Berka



Thank you for your attention



overprot.ncbr.muni.cz



sestra.ncbr.muni.cz



2dprots.ncbr.muni.cz



midlik@mail.muni.cz



CEITEC