

Intrinsic disorder proteins

Peter Tompa

Institute of Enzymology
Hungarian Academy of Sciences
Budapest, Hungary

Keith Dunker

CCBB & IUSM
Indiana University
Indianapolis, USA

Why do we study disorder proteins?

- 1) Gap in knowledge (1600 vs. thousands of IUPs)
- 2) Structural genomics initiatives
- 3) Bioinformatics studies
- 4) Single protein studies

Analysis of Signaling Interactions

- Examined each interaction on Pawson's website.
- Almost all of the interactions involved ordered regions binding to disordered partners.
- **Conclusion:** if Pawson's examples are typical, then a very significant proportion of protein-protein signaling interactions use disordered regions.

Parallel Paradigms

Catalysis

AA seq → 3-D Structure → Function

Signaling

AA seq → Disordered
Ensemble → Function

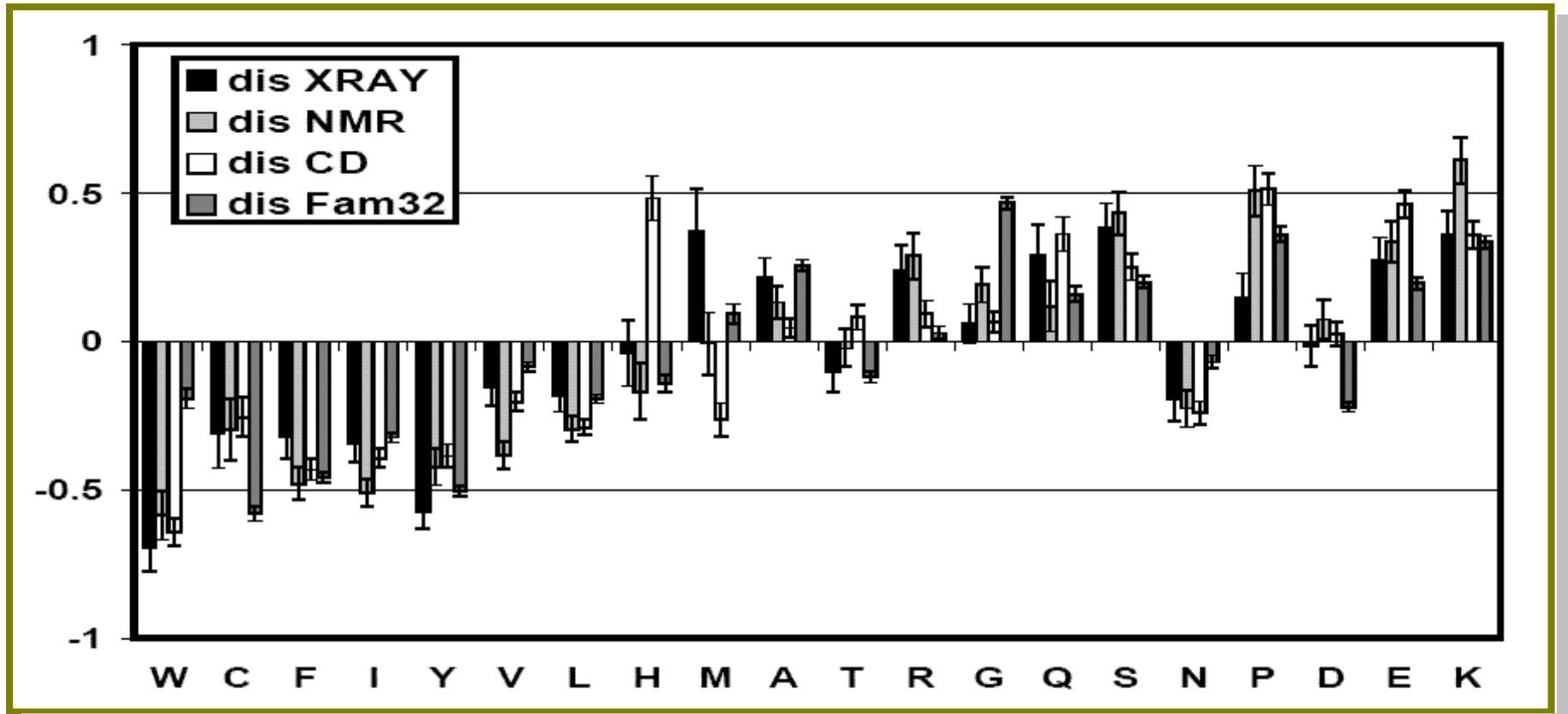
Disorder and Function

Category	Change	Examples	Descriptions
Molecular Recognition	D → O	113	Inter- and Intra-protein, ssDNA, dsDNA, tRNA, rRNA, mRNA, nRNA, bilayers, ligands, co-factors, metals
Protein Modification	Variable	36	Acetylation, fatty acylation, glycosylation, methylation, phosphorylation, ADP-ribosylation, ubiquitination, proteolytic digestion
Entropic Chains	Variable	17	Linkers, spacers, bristles, clocks, springs, detergents, self-transport

Basic approaches to predict disorder

- 1) Machine learning
- 2) Structural approach

The unusual AA composition of IDPs



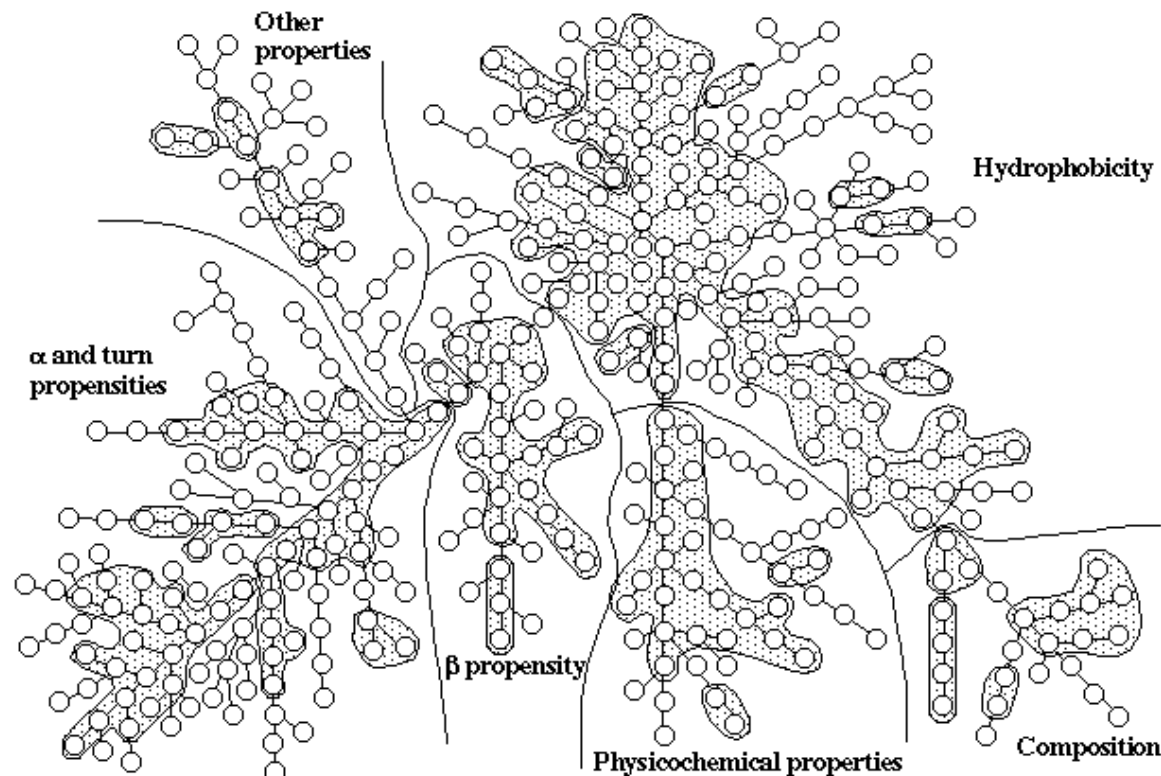
order-promoting

disorder-promoting

AA feature space: AAindex database

<http://www.genome.jp/aaindex>

A number is associated with every amino acid, which quantitatively describes how characteristic the given feature is to the AA (has 566 different scales at present)



Based on AA compositions, two things you might not want to do...

- 1) low complexity regions**
- 2) regions w/o secondary structure**

Low complexity regions

Sequence databases contain a lot of regions in which only a few amino acids occur (simple), or some dominate (biased), this can be described by an entropy function - Wootton (1993)

Shannon`s
entropy

$$K_2 = - \sum_{i=1}^N \frac{n_i}{L} \left(\log_2 \frac{n_i}{L} \right)$$

$L > 20$; window size,
 N alphabet size (20),
 n_i : az i . aminosav száma

...may correspond to disorder

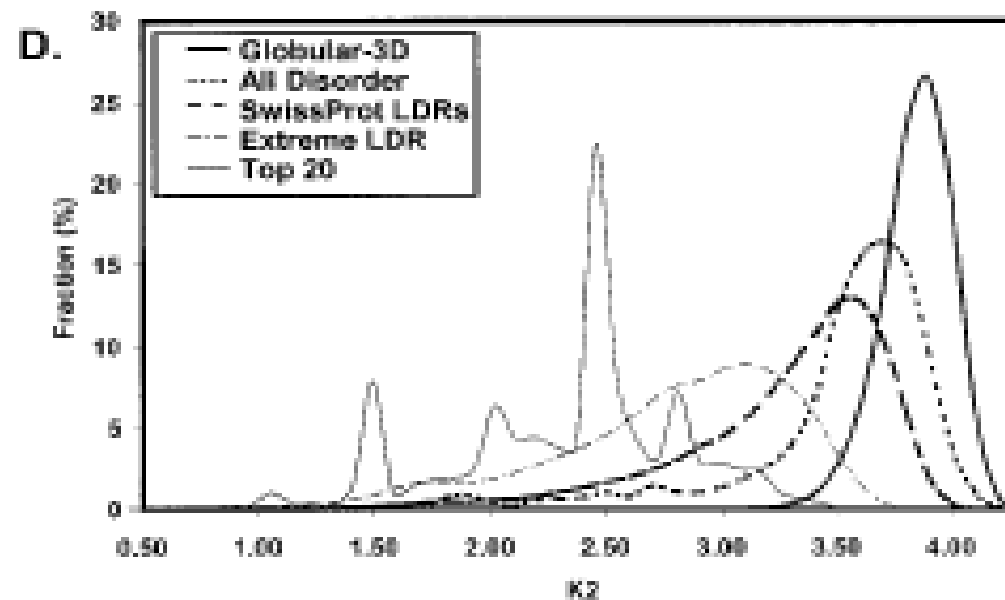
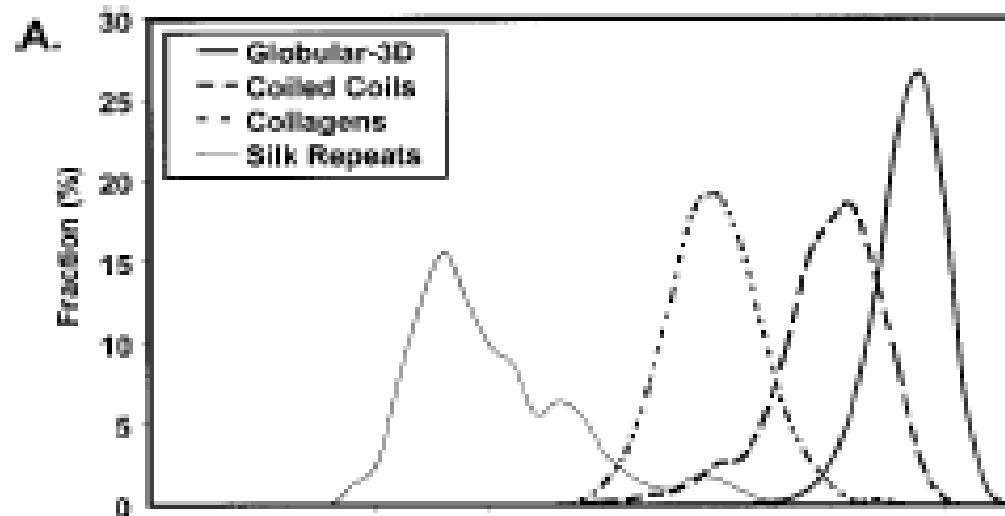
>SRY_MOUSE TRANSCRIPTIONAL ACTIVATOR

1-143 MEGHVKRPMNAFMVWSRGERHKLAQQNPSM
QNT EISKQLGCRWKS L TEAEKRPF FQEAQR
LKILHREKYPNYKYQPHRRAKVSQ RSGILQ
PAVASTKLYNLLQWDRNP HAITYRQDWSRA
AHLYSKNQQS FYWQPVDIPTGHL

qqqqqqqqqqqqfhn hhqqqqqfydh hqqqq 144-366
qqqqqqqqqfhd hhqqkqqfhd hhqqqqqqfh
dhhhhhqqqqfhd hhqqqqqqfhd hqqqqqqq
qqqqqqfhd hhqqkqqfhd hhhhqqqqqqfhd
hqqqqqqfhd hqqqqqh qfhd hpqqkqqfhd
hpqqqqqqfhd hhhhqqqqkqqfhd hhqqkqq
fhd hhqqkqqfhd hhqqqqqqfhd hhqqqqqq
qqqqqqqqqfhdqq

367-395
LTYLLTADITGEHTPYQEHLSTALWLAVS

The relationship of low complexity and disorder



PROTEINS: Structure, Function, and Genetics 42:38–48 (2001)

red Protein

Garner,^{2†} Celeste J. Brown,² and A. Keith Dunker^{2*}
State University, Pullman, Washington
an, Washington

Based on AA compositions, two things you might not want to do...

- 1) low complexity regions**
- 2) regions w/o secondary structure**

...but there are globular proteins without Secondary structure and IDPs with secondary structure

doi:10.1016/S0022-2836(02)00736-2 available online at <http://www.idealibrary.com> on IDEAL[®]

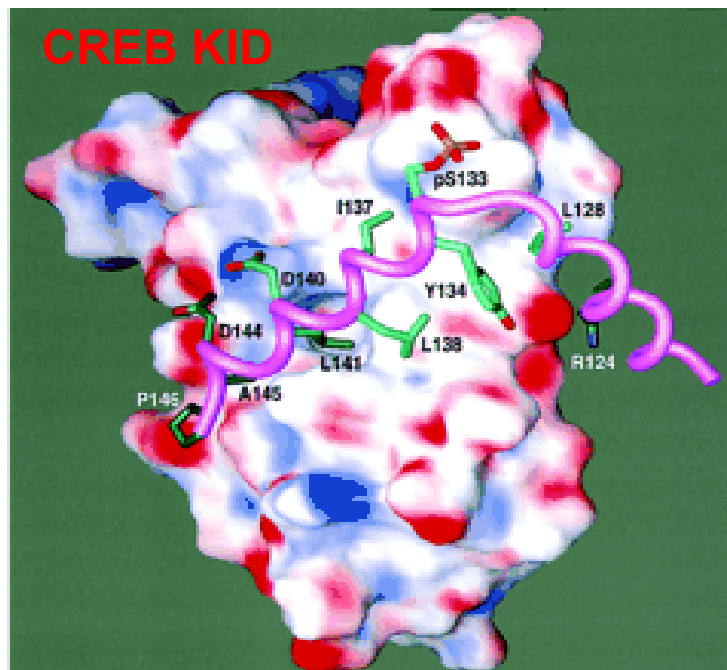
J. Mol. Biol. (2002) 322, 53–64

JMB



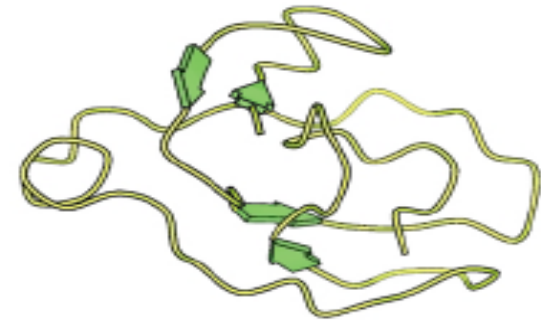
Loopy Proteins Appear Conserved in Evolution

Jinfeng Liu^{1,2,3}, Hepan Tan² and Burkhard Rost^{2,3,4*}

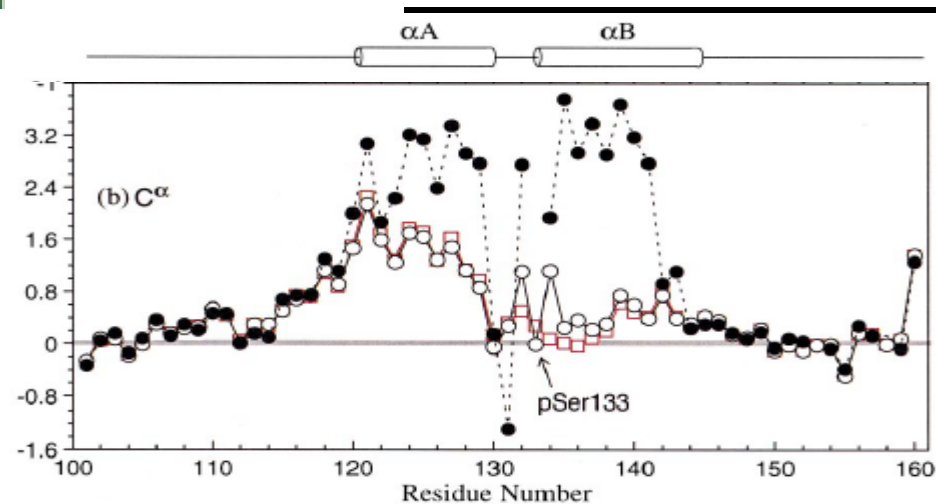


Radhakrishnan (1997) *Cell* 91, 741

(d) Loopy domains



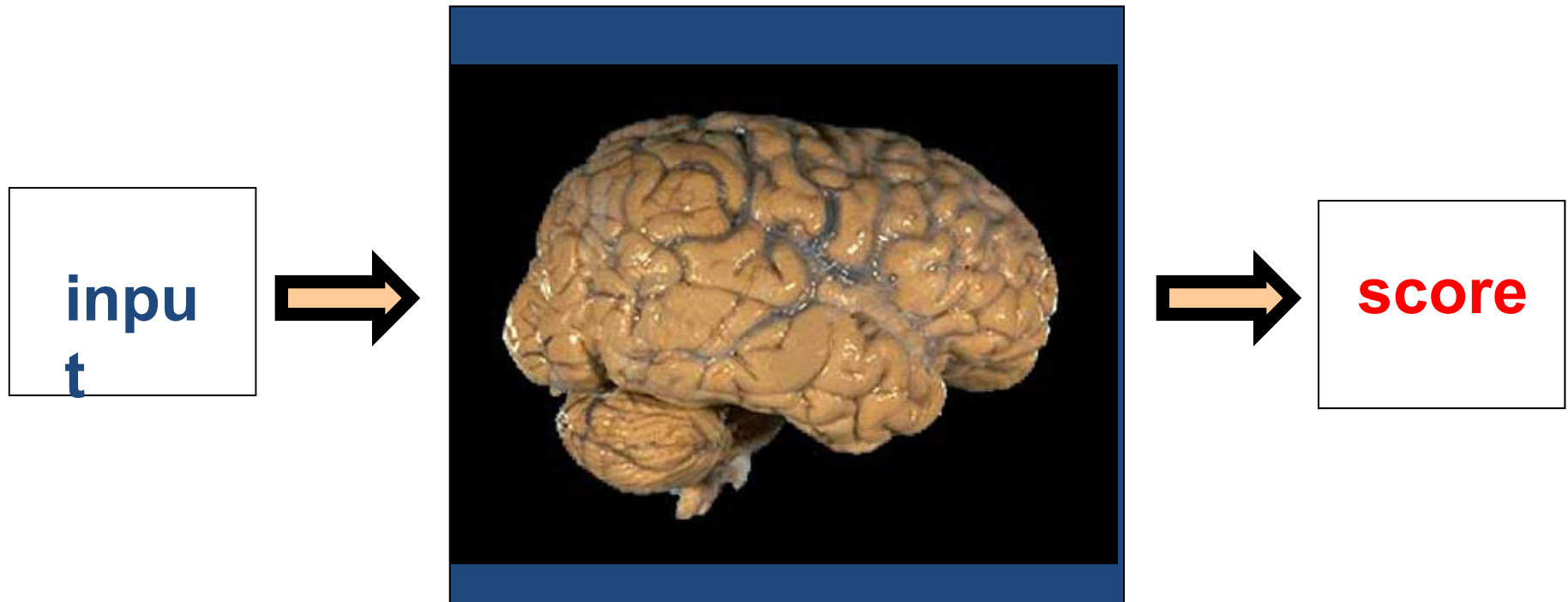
1tbi



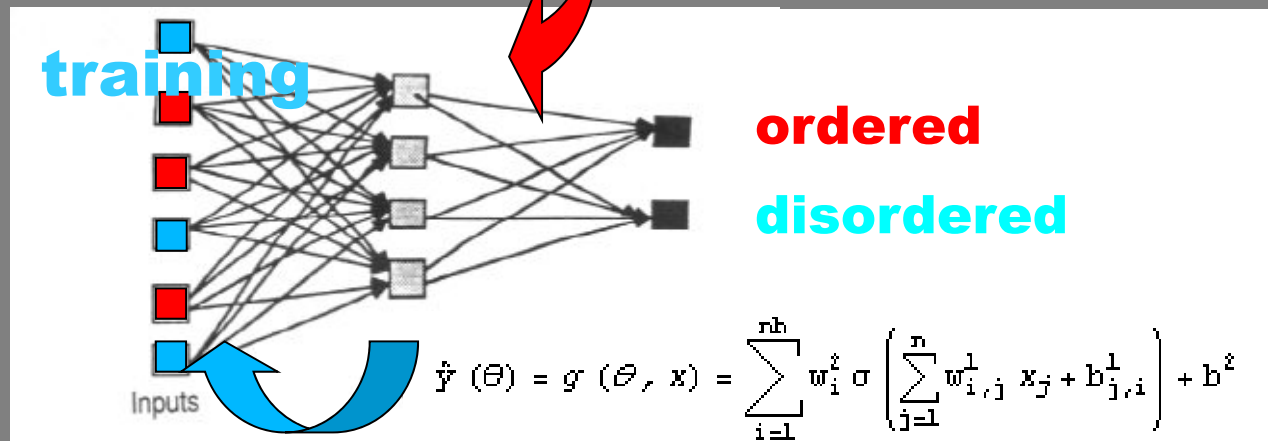
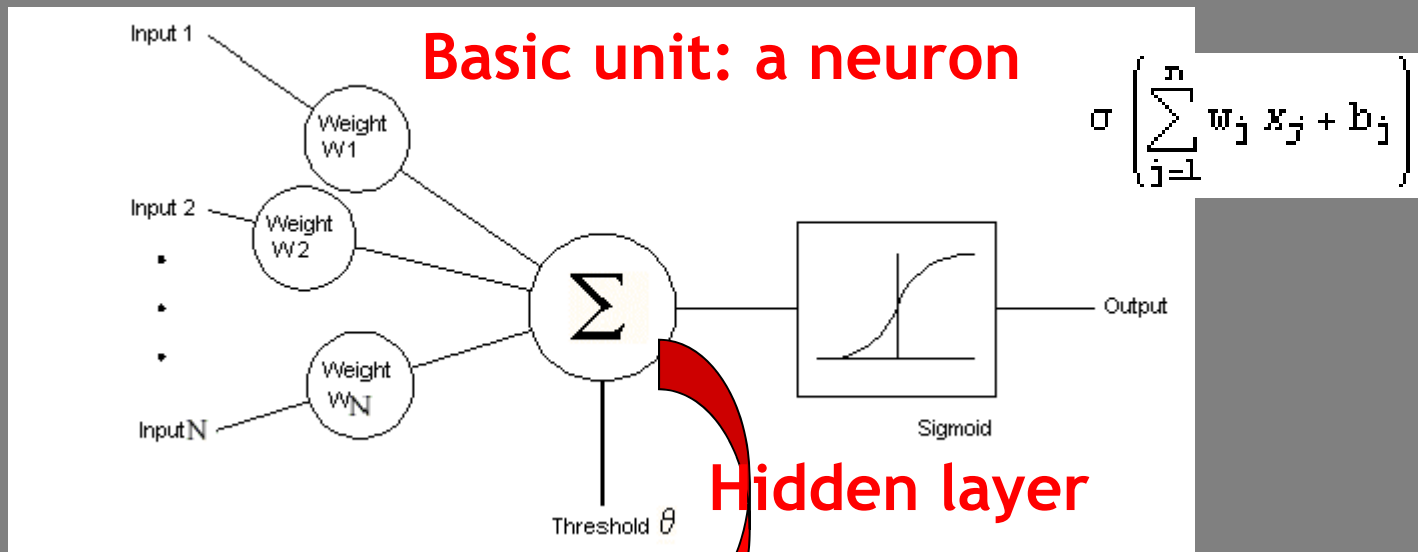
Radhakrishnan (1998) *FEBS Lett.* 430, 317

1) Machine learning

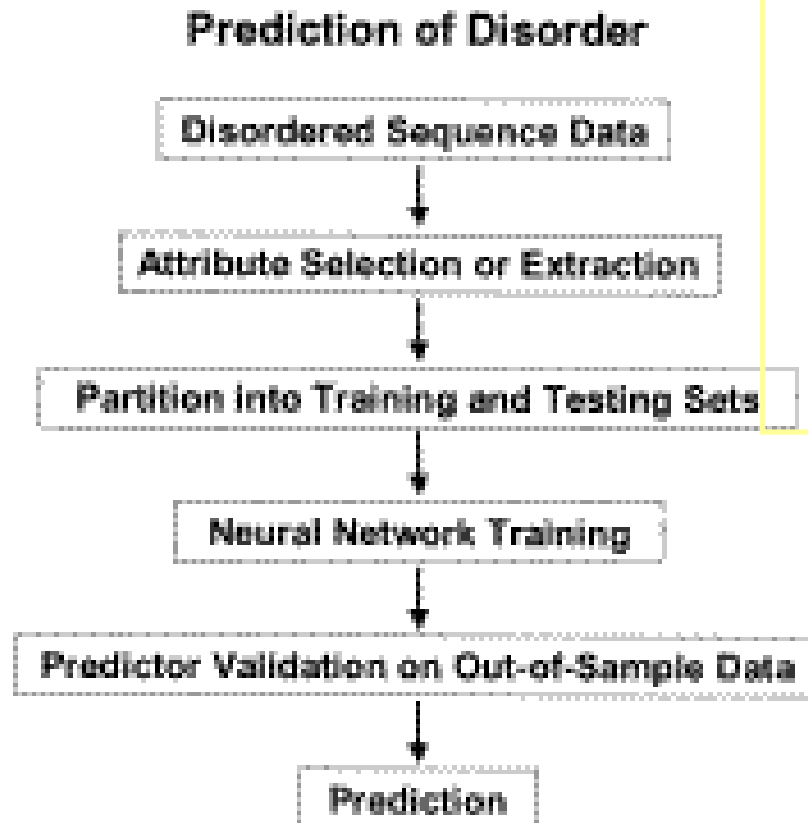
Artificial neural network (NN)



Artificial neural networks



Predictor of naturally disordered regions (PONDR®)



Input

18 amino acids

Hydrophobicity

Sequence complexity

VL2, VL3

VLXT

VSL2

Dunker, 1998

**Predicting Disordered Regions from Amino Acid
Sequence: Common Themes Despite Differing
Structural Characterization**

Ethan Garner ¹

Paul Cannon ²

Pedro Romero ²

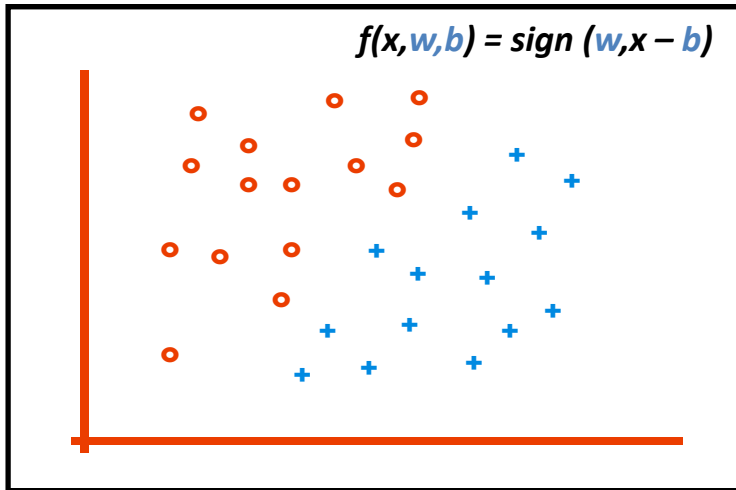
**THOUSANDS OF PROTEINS LIKELY TO HAVE LONG
DISORDERED REGIONS**

PEDRO ROMERO, ZORAN OBRADOVIC
*School of Electrical Engineering and Computer Science
Washington State University, Pullman, WA 99164*

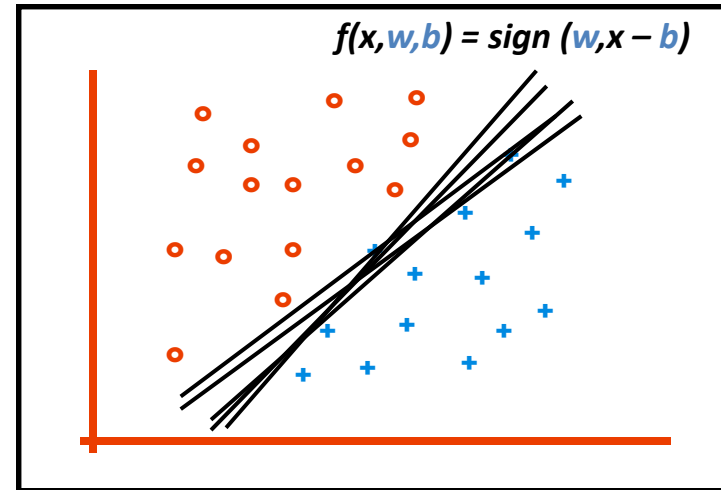
**PROTEIN DISORDER AND THE EVOLUTION OF
MOLECULAR RECOGNITION: THEORY, PREDICTIONS
AND OBSERVATIONS**

A. K. DUNKER, E. GARNER, S. GUILLIOT
dunker@mail.wsu.edu
*Department of Biochemistry & Biophysics,
Washington State University, Pullman, WA 99164-4660*

Support vector machine (SVM)

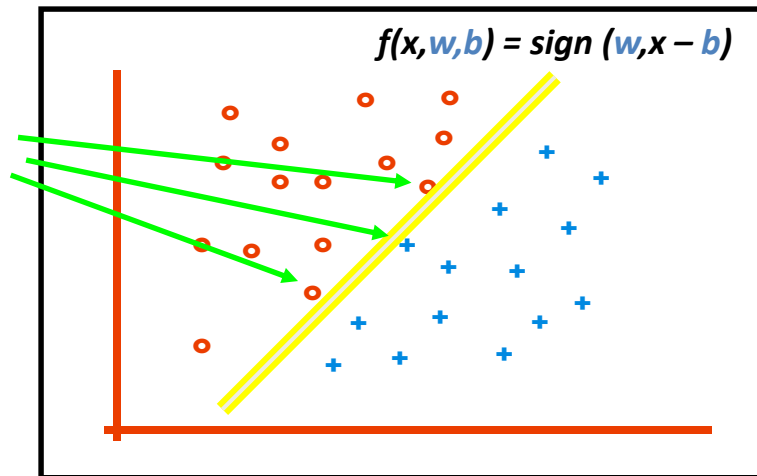


How would you classify this data?



But which is best?

Support vectors
against which the
margin pushes up



This is the simplest
SVM, the linear SVM
(LSVM)

The maximum margin linear
classifier is considered best

2) Structural approach (interaction potential)

The protein non-folding problem

- **Protein folding problem**

How does amino acid sequence determine protein structure ?

- **Protein non-folding problem**

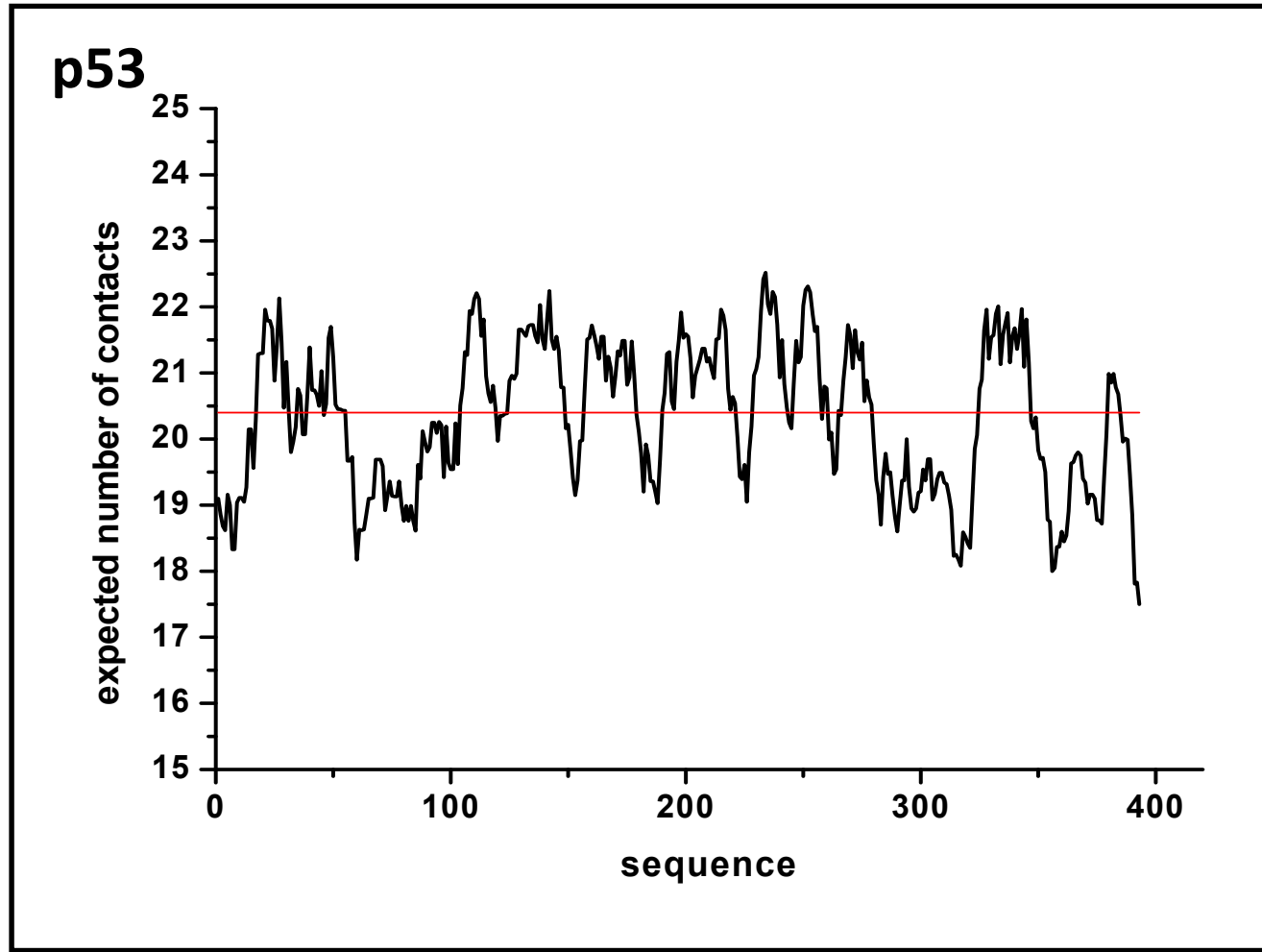
How does amino acid sequence determine the lack of protein structure ?

The protein non-folding problem

- Globular proteins have special sequences that enable the formation of a large number of favorable interactions
- IDPs contain (disorder-promoting) amino acids, which tend to avoid interacting with each other.
- An IDP thus cannot fold into a low-energy conformation.

A simple implementation, FoldUnfold

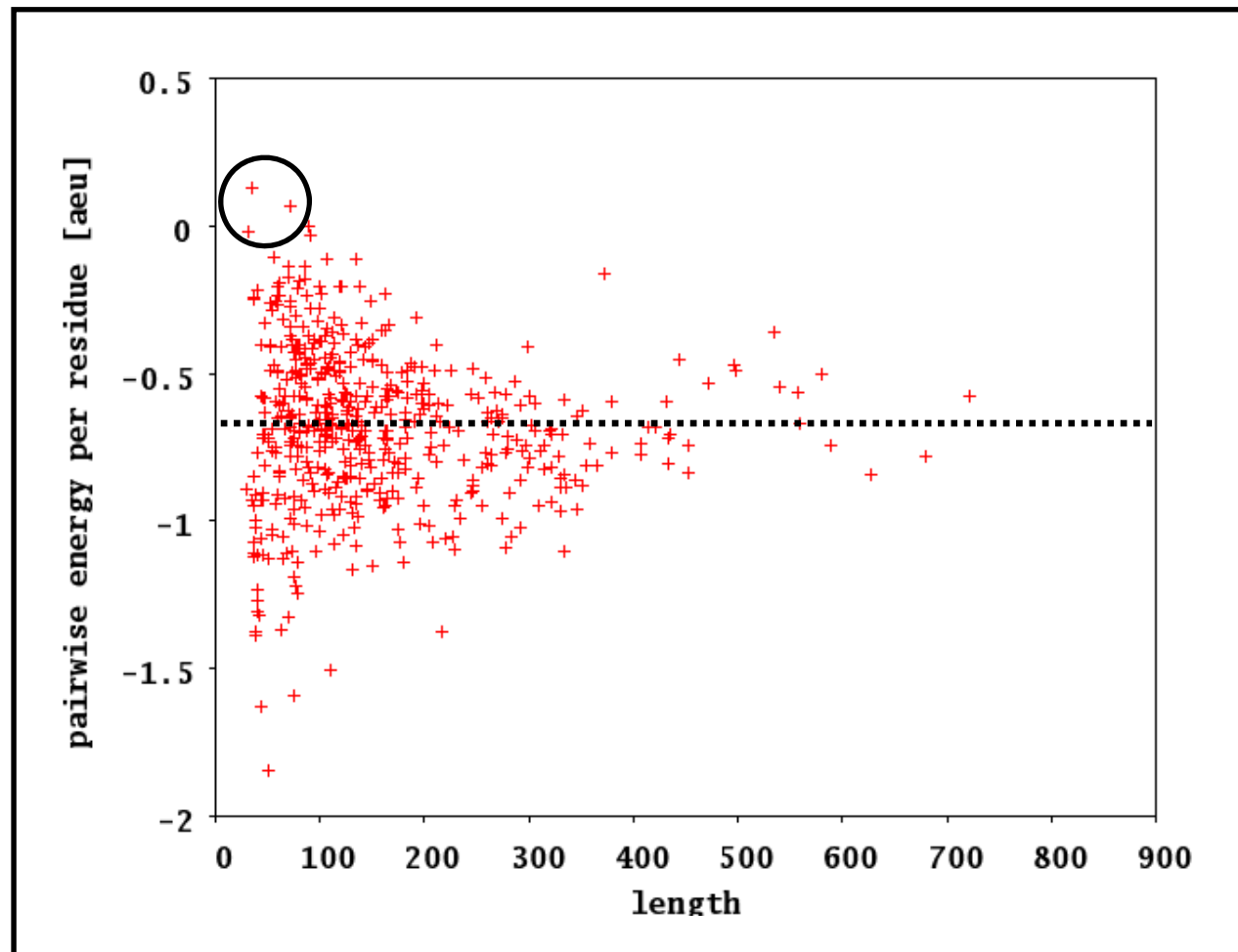
Calculates the contact number of amino acids



Estimating the total pairwise interresidue interaction energy of a sequence: IUPred

- 1) Calculate interresidue interaction energies from structure
- 2) Try to estimate the energy without knowing the structure
- 3) Apply the estimation to sequences w/o structure (e.g. to IDPs, which have no structure)

Interresidue interaction energy calculated for known structures



How to estimate the interresidue interaction energy of a protein of unknown structure or w/o structure?

Structure

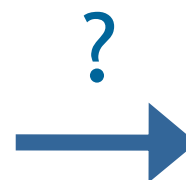
MODEL	1					
ATOM	1	N	MET	A	23	
-4.381						2.191 28.312
ATOM	2	CA	MET	A	23	
-3.305						2.394 27.327
ATOM	3	C	MET	A	23	
-3.706						3.514 26.377
ATOM	4	O	MET	A	23	
-4.867						3.589 25.977



Calculated
energy per
residue

Sequence

```
MKVPPHSIEA EQSVLGGLML
DNERWDDVAE RVVADDFYTR
PHRHIFTEMA RLQESGSPID
LITLAESLER QQLDSVGGF
AYLAELSKNT PSAANISAYA
DIVRERAVVR EMIS
```



Estimated
energy per
residue

Estimation of pairwise interresidue interaction energy from sequence

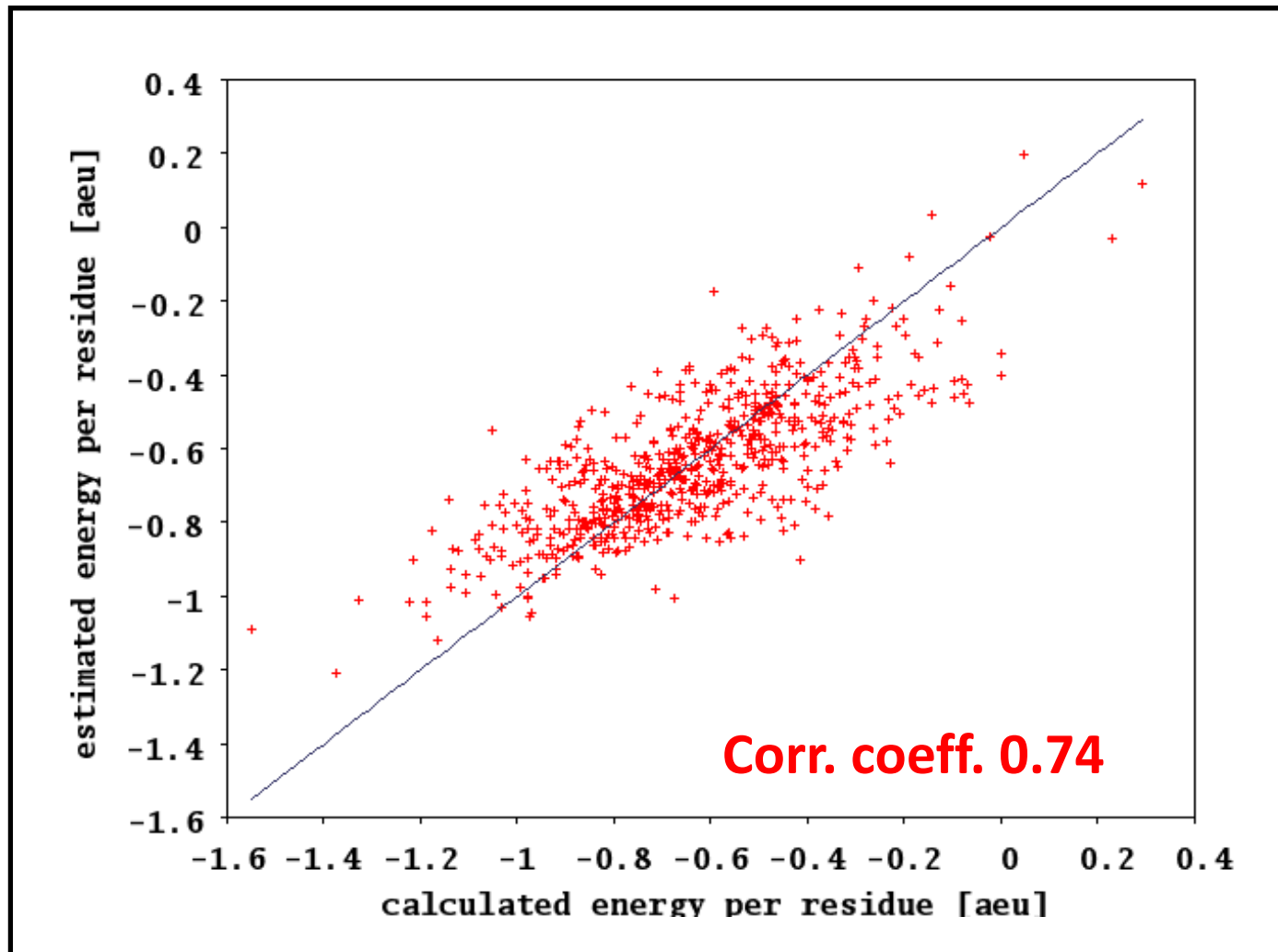
The contribution of individual AAs depends on its potential partners, i.e. its neighborhood. A quadratic formula is needed to take this into consideration.

$$E(\text{estimated}) / L = (n_A \quad n_C \quad \cdots \quad n_Y) \begin{pmatrix} P_{AA} & P_{AC} & \cdots & P_{AY} \\ P_{CA} & P_{CC} & & \\ \vdots & & \ddots & \\ P_{YA} & \cdots & \cdots & P_{YY} \end{pmatrix} \begin{pmatrix} n_A \\ n_C \\ \vdots \\ n_Y \end{pmatrix}$$

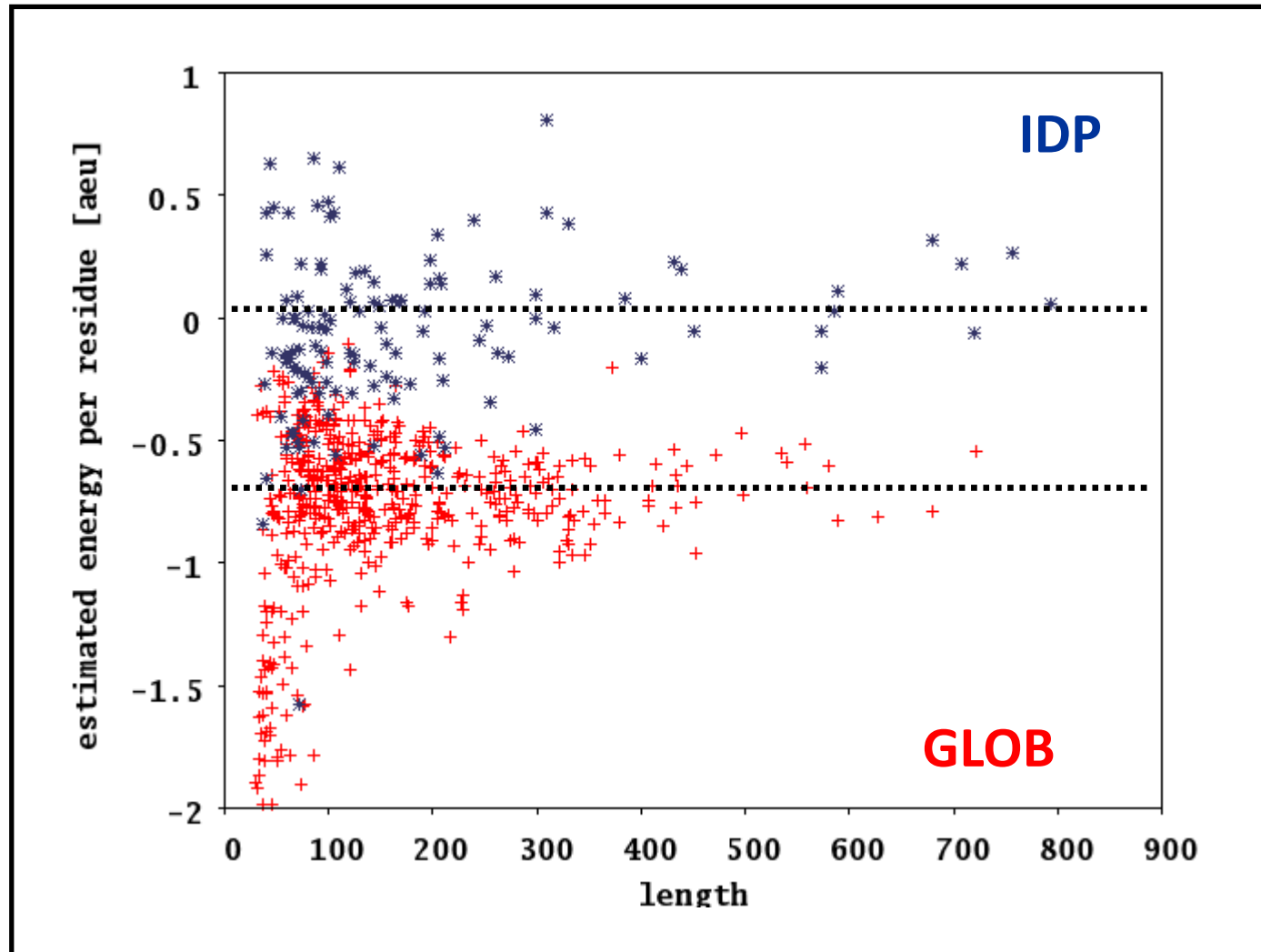
The relationship between AA composition and energy is given by an optimized 20x20 energy predictor matrix, P_{ij}

$$P_{ij} : \sum_{\substack{\text{Globular} \\ \text{proteins} \\ k}} (E_k(\text{calc}) - E_k(\text{est}))^2 \rightarrow \min$$

Correlation of calculated and estimated interaction energies



The estimated energy for globular proteins and IDPs



Making it position-specific: the IUPred algorithm



The screenshot shows the IUPred web interface within a Mozilla browser window. The browser's title bar reads "IUPred - Mozilla". The address bar is empty. The main content area features the IUPred logo at the top left, followed by the heading "Prediction of Intrinsically Unstructured Proteins". Below this heading is a sidebar with a list of links: "IUPs", "Theory", "How to use", "IUPred", "Related links", and "Comments". The main text area contains a paragraph explaining the algorithm: "Intrinsically unstructured/disordered proteins have no single well-defined tertiary structure in their native, functional state. Our server recognizes such regions from the amino acid sequence based on the estimated pairwise energy content. The underlying assumption is that globular proteins are composed of amino acids which have the potential to form a large number of favorable interactions, whereas intrinsically unstructured proteins (IUPs) adopt no stable structure because their amino acid composition does not allow sufficient stable interactions to form." Below the text is a form with several input fields and buttons. The "Title:" field is empty. The "Sequence:" field is a large text area. The "Prediction type:" section has three radio buttons: "long disorder" (selected), "short disorder" (with a note "(e.g. missing residues of X-ray structures)"), and "structured regions". The "Output type:" section has two radio buttons: "raw data only" (selected) and "generate plot". Below the "generate plot" option is a dropdown menu set to "500" and the text "plot window size". At the bottom of the form are two buttons: "SUBMIT" and "CLEAR". The browser's status bar at the bottom shows the Mozilla logo and some icons.

IUPred

Prediction of Intrinsically Unstructured Proteins

- IUPs
- Theory
- How to use
- IUPred
- Related links
- Comments

Intrinsically unstructured/disordered proteins have no single well-defined tertiary structure in their native, functional state. Our server recognizes such regions from the amino acid sequence based on the estimated pairwise energy content. The underlying assumption is that globular proteins are composed of amino acids which have the potential to form a large number of favorable interactions, whereas intrinsically unstructured proteins (IUPs) adopt no stable structure because their amino acid composition does not allow sufficient stable interactions to form.

Title:

Sequence:

Prediction type:

- ☒ long disorder
- ☐ short disorder (e.g. missing residues of X-ray structures)
- ☐ structured regions

Output type:

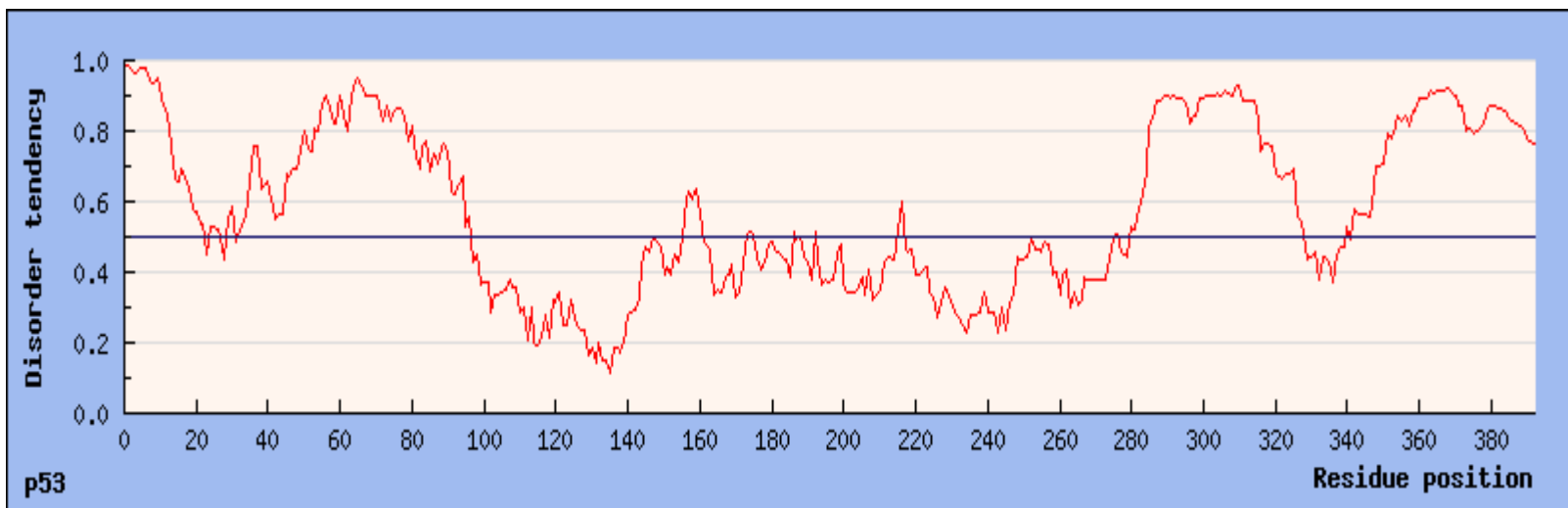
- ☒ raw data only
- ☐ generate plot

500 plot window size

SUBMIT CLEAR

IUPred: <http://iupred.enzim.hu>

IUPred, p53



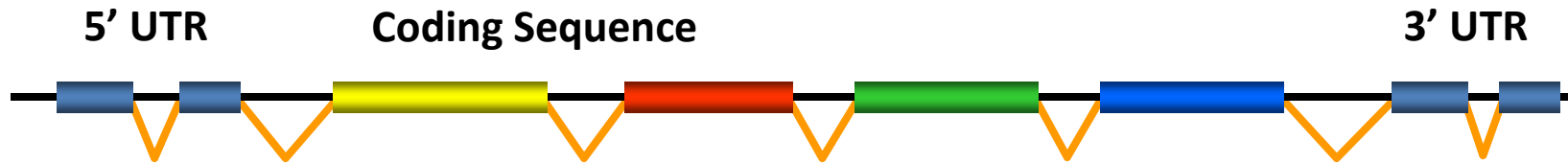
www.disprot.org

<u>DisEMBLTM</u>	Intrinsic Protein Disorder Prediction
<u>DISOPRED2</u>	Disorder Prediction Server
<u>DRIPPRED</u>	Web based predictor for disordered regions in proteins
<u>FoldIndex©</u>	Estimate the fold probability of a protein
<u>GlobPlot 2</u>	Intrinsic Protein Disorder, Domain & Globularity Prediction
<u>IUPred</u>	Prediction of Intrinsically Unstructured Proteins
<u>PONDR[®]</u>	Predictors of Natural Disordered Regions
<u>PreLink</u>	Prediction of unfolded segments in a protein sequence based on amino acid composition
<u>RONN</u>	Regional Order Neural Network
<u>VL2</u>	DisProt Predictor of Intrinsically Disordered Regions
<u>VL3, VL3H, VL3E</u>	DisProt Predictor of Intrinsically Disordered Regions

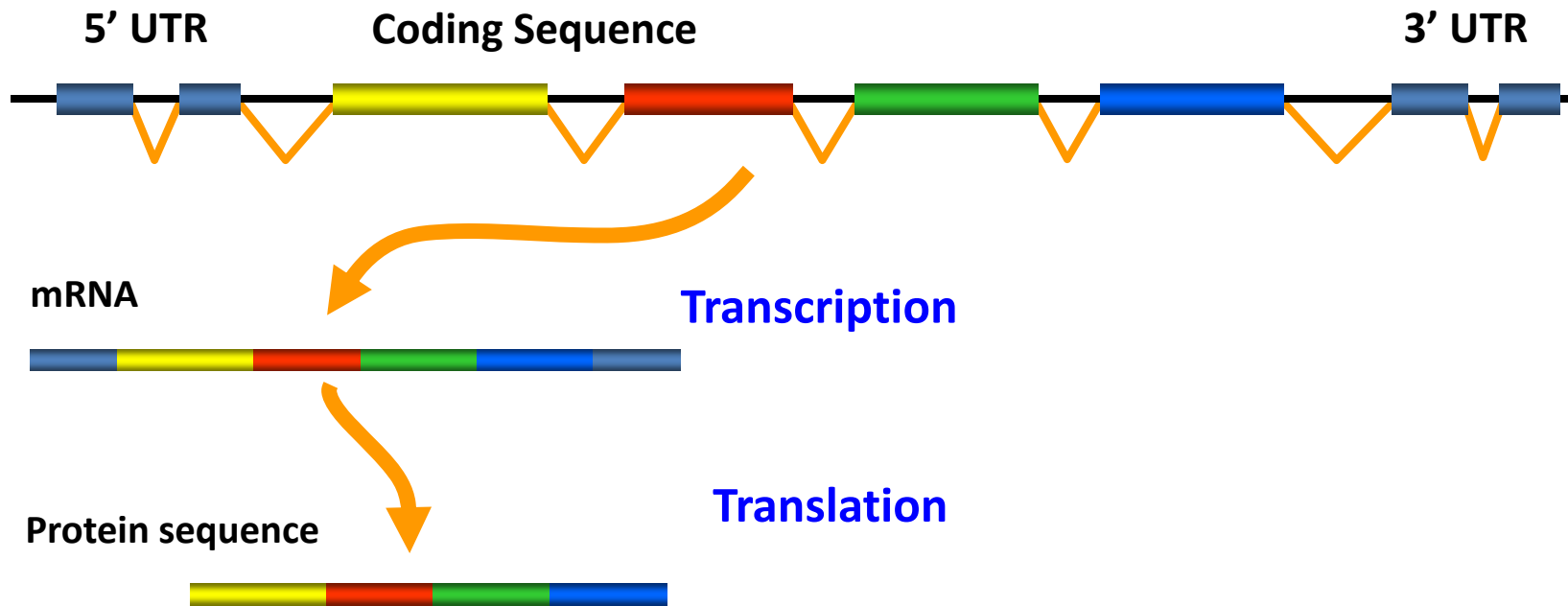
Alternative Splicing and Intrinsic Disorder

- Find proteins with both **ordered** and **disordered** regions.
- Find mRNA alternative splicing information for these proteins and map to the **ordered** and **disordered** regions.
- For alternatively spliced regions of mRNA, do they code for **ordered protein** more often or do they code for **disordered protein** more often?

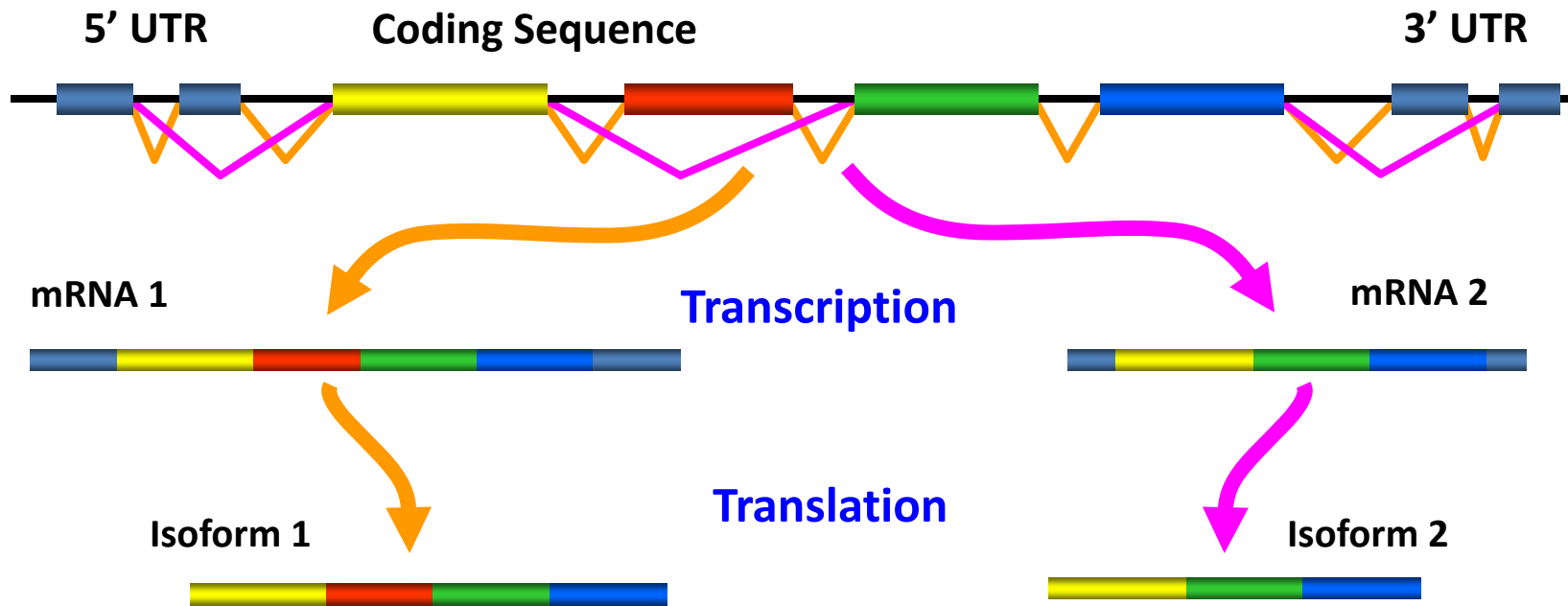
Alternative Splicing



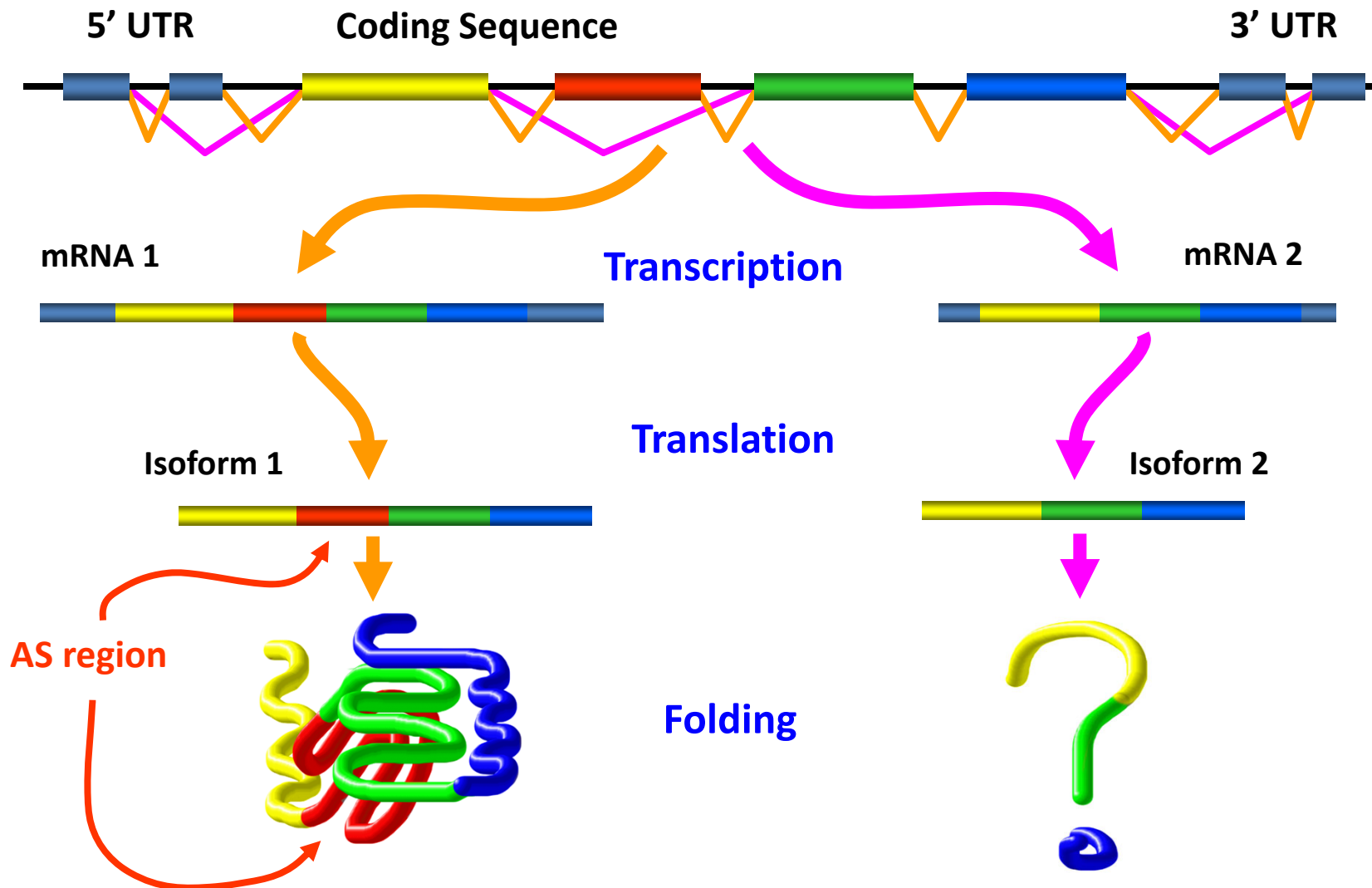
Alternative Splicing



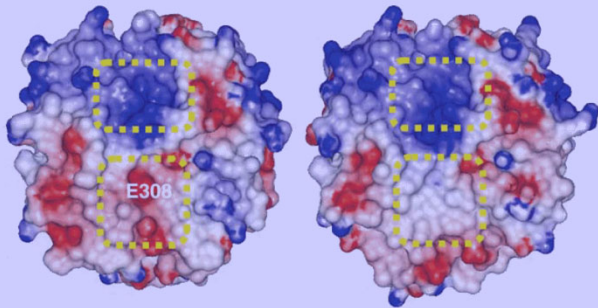
Alternative Splicing



Alternative Splicing



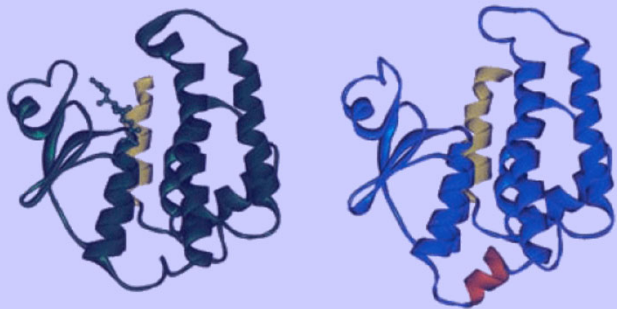
Structural Studies of AS



EDA-A1

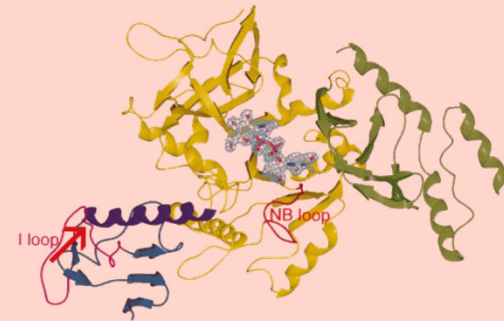
EDA-A2

Structured AS regions

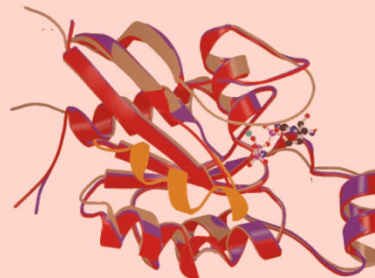


Glutathione S-transferase

Disordered AS regions



Pyrophosphorylase



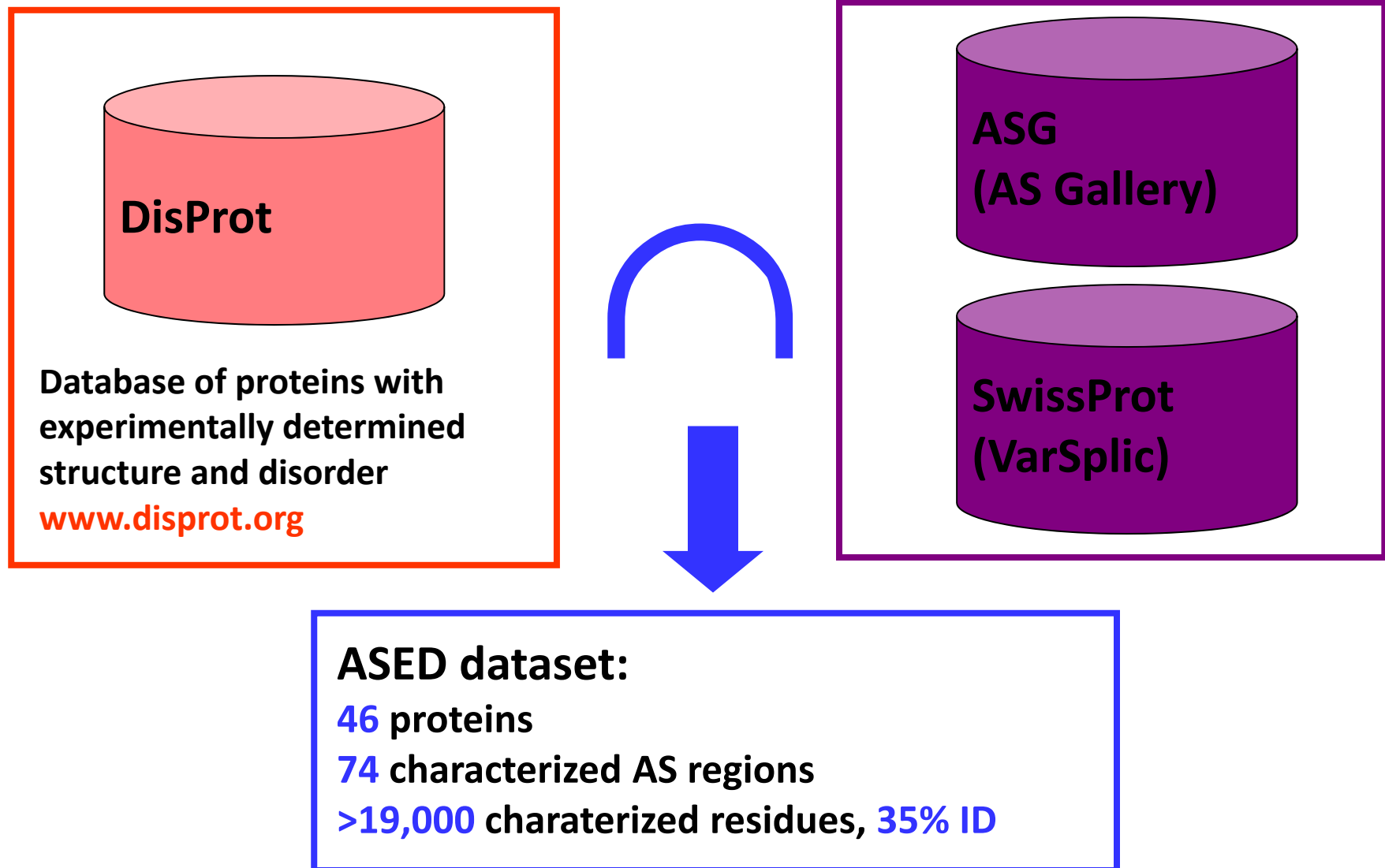
RAC1

Tumor necrosis factor



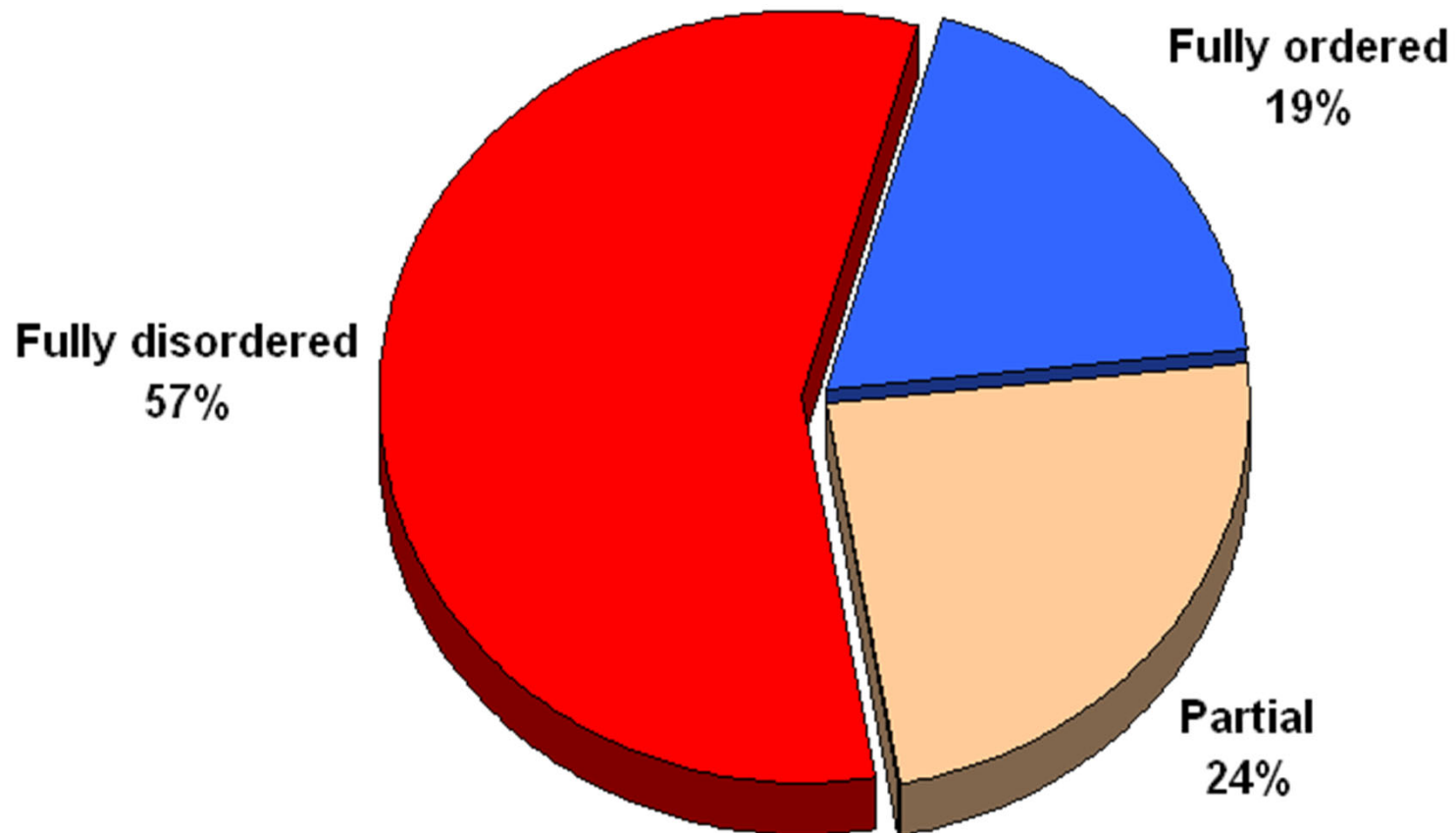
Sulphotransferase

Studying the Relationship ID \leftrightarrow AS

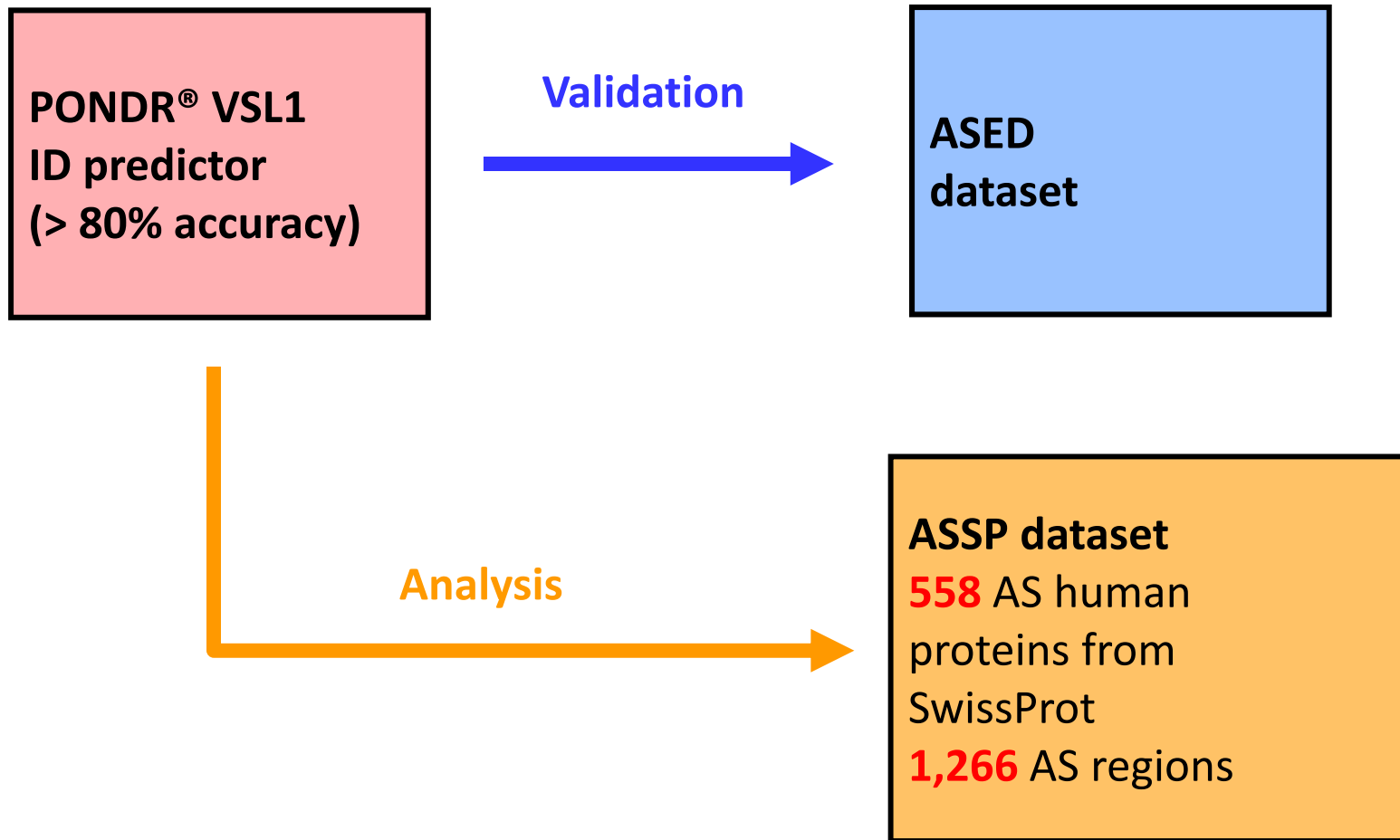


Results on ASED

Distribution of structurally characterized AS regions

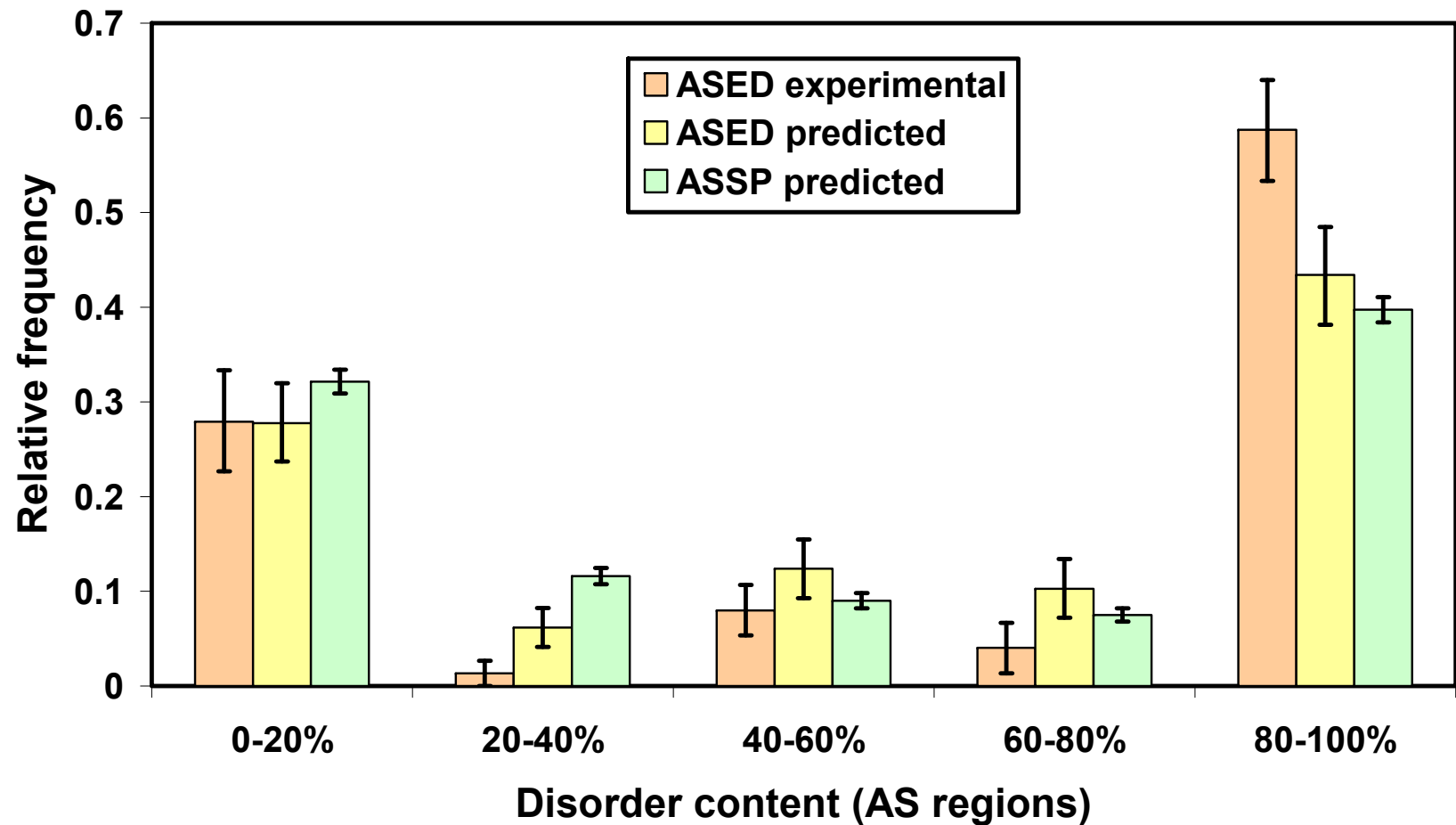


Enlarging the Dataset



Global Results

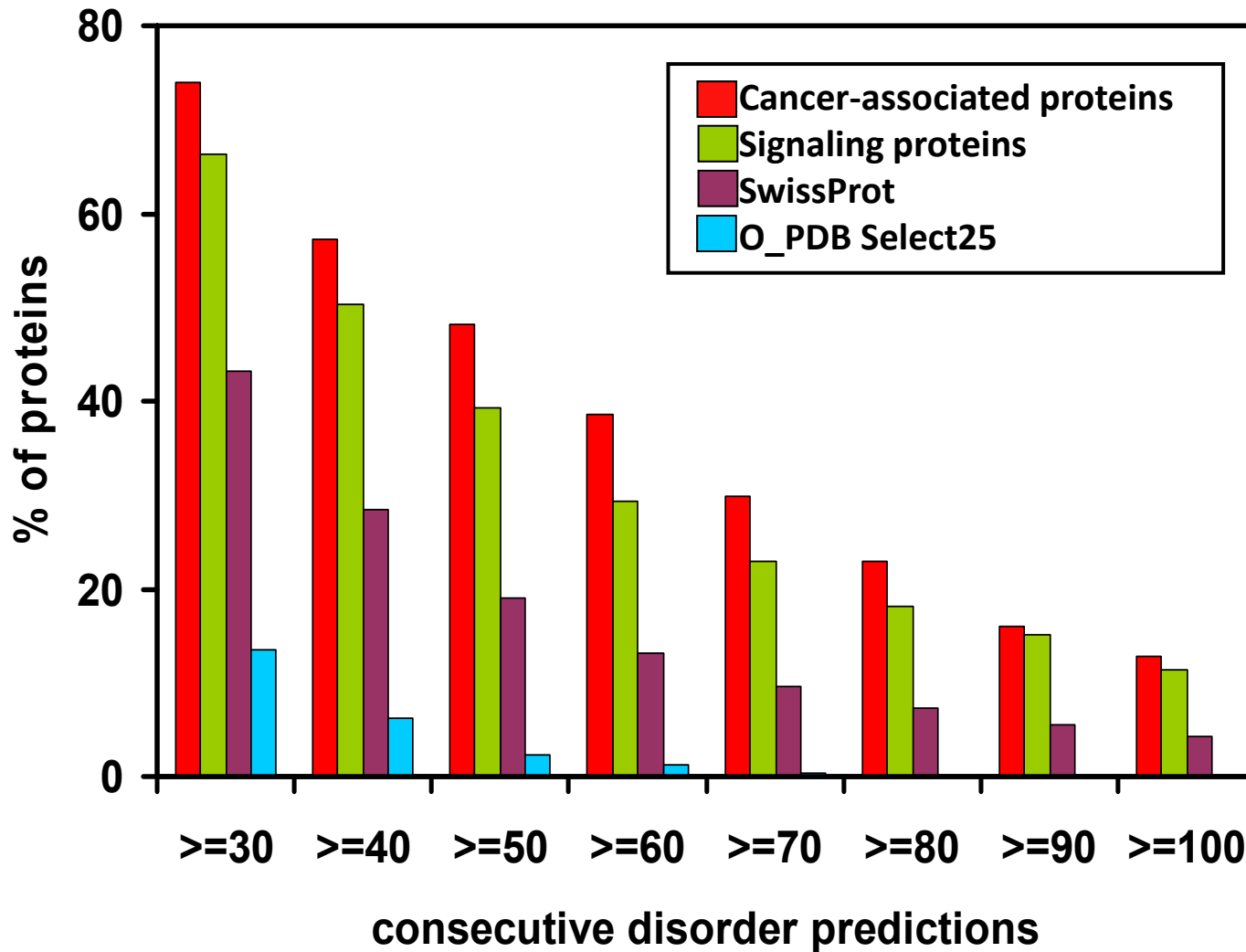
AS regions disorder distributions in ASED and ASSP



Alternative Splicing and Disorder

- **Ordered Proteins:** active site residues non-local in sequence, become associated by protein folding
- **Disordered Proteins and regions:** functional residues localized in sequence
- Functional regions for signaling and regulation are located one after another
- Alternative splicing edits functional sets and thereby leads to regulatory and signaling diversity

Disorder and Cell Signaling



Disorder and Drug Discovery

- The p53-MDM2 interaction is blocked by several drugs; one is in clinical trials and shows promise as an anti-cancer drug.
- The drug molecules bind to the ordered partner, preventing the disordered partner from binding.
- Such interactions are typically weak per unit of surface area, and the interaction surfaces can be small, thus such interactions are ideal drug targets.
- Molecular Kinetics has strategy to find all druggable MoRF-based interactions; bioinformatics indicates that more than one hundred are in cancer-associated proteins.
- Is this approach a new drug discovery pathway?