

Statistical Machine Learning  
Methods for Bioinformatics  
**IV. Neural Network & Deep  
Learning Applications in  
Bioinformatics**

Jianlin Cheng, PhD

Department of Electrical Engineering & Computer Science  
University of Missouri, Columbia

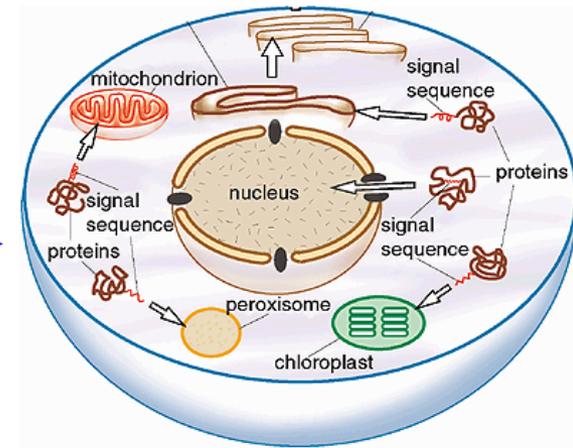
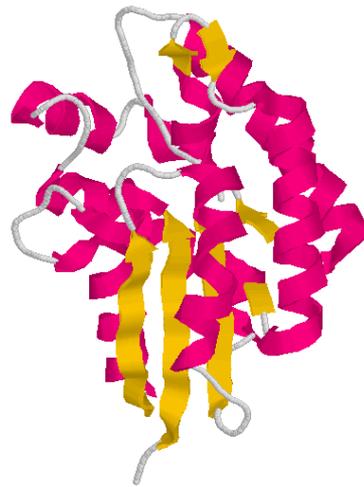
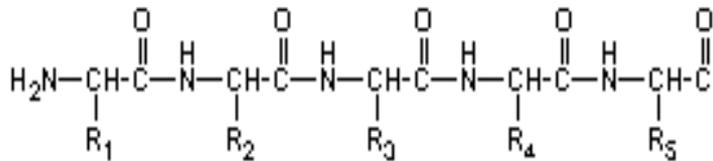
2018

# Neural Network & Deep Learning Application in Bioinformatics

- Neural network is one of the most widely used methods in bioinformatics.
- Deep learning is the most popular method in bioinformatics
- It is used in gene structure prediction, protein structure prediction, gene expression data analysis, ... Almost anywhere when you need to do classification.
- Here we specifically focus on applying neural networks to protein structure prediction (**secondary structure**, solvent accessibility, disorder region, **contact map**, **fold classification**).

# Sequence, Structure and Function

AGCWY.....



**Cell**

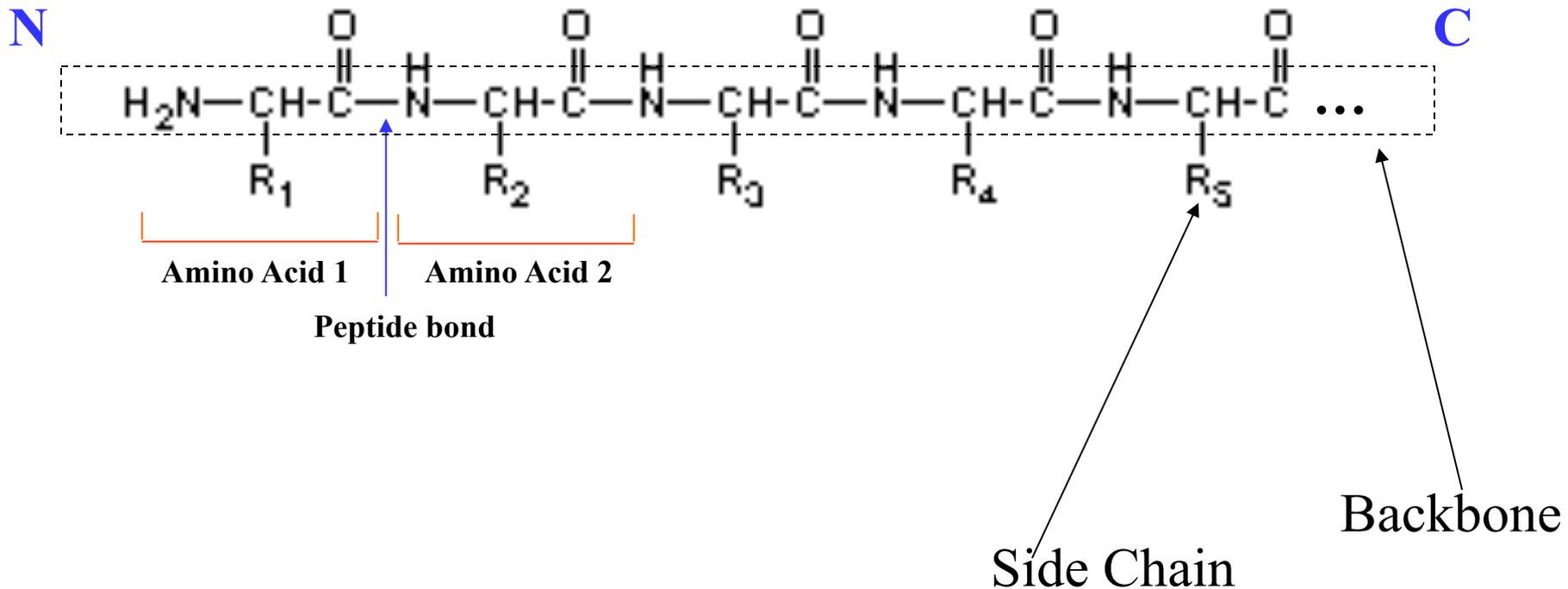
# Protein Sequence – Primary Structure

- The first protein was sequenced by Frederick Sanger in 1953.
- Twice Nobel Laureate (1958, 1980) (other: Curie, Pauling, Bardeen).
- Determined the amino acid sequence of insulin and proved proteins have specific primary structure.

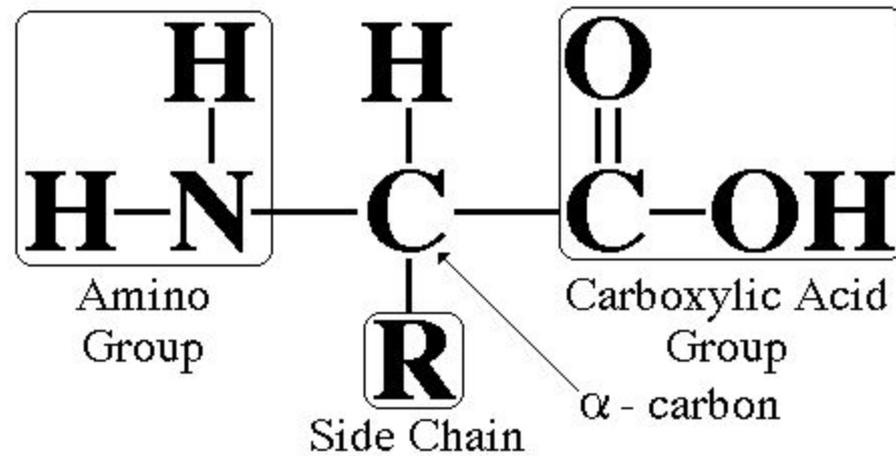


# Protein Sequence

A directional sequence of amino acids/residues



# Amino Acid Structure



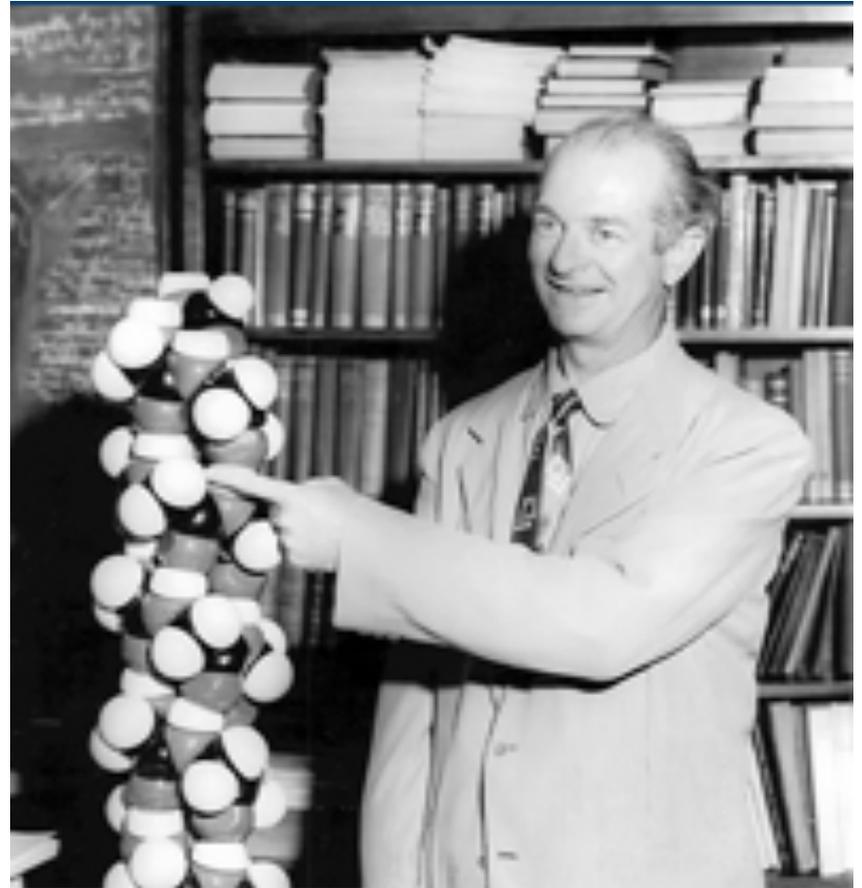
# Amino Acids

| Amino acid    | Abbrev. | Side chain   | Hydrophobic | Polar | Charged  | Small | Tiny | Aromatic or Aliphatic | van der Waals volume | Codon                        | Occurrence in proteins (%) |
|---------------|---------|--|-------------|-------|----------|-------|------|-----------------------|----------------------|------------------------------|----------------------------|
| Alanine       | Ala, A  | -CH <sub>3</sub>   | X           | -     | -        | X     | X    | -                     | 67                   | GCU, GCC, GCA, GCG           | 7.8                        |
| Cysteine      | Cys, C  | -CH <sub>2</sub> SH  | X           | -     | -        | X     | -    | -                     | 86                   | UGU, UGC                     | 1.9                        |
| Aspartate     | Asp, D  | -CH <sub>2</sub> COOH  | -           | X     | negative | X     | -    | -                     | 91                   | GAU, GAC                     | 5.3                        |
| Glutamate     | Glu, E  | -CH <sub>2</sub> CH <sub>2</sub> COOH                          | -           | X     | negative | -     | -    | -                     | 109                  | GAA, GAG                     | 6.3                        |
| Phenylalanine | Phe, F  | -CH <sub>2</sub> C <sub>6</sub> H <sub>5</sub>                 | X           | -     | -        | -     | -    | Aromatic              | 135                  | UUU, UUC                     | 3.9                        |
| Glycine       | Gly, G  | -H   | X           | -     | -        | X     | X    | -                     | 48                   | GGU, GGC, GGA, GGG           | 7.2                        |
| Histidine     | His, H  | -CH <sub>2</sub> -C <sub>3</sub> H <sub>3</sub> N <sub>2</sub> | -           | X     | positive | -     | -    | Aromatic              | 118                  | CAU, CAC                     | 2.3                        |
| Isoleucine    | Ile, I  | -CH(CH <sub>3</sub> )CH <sub>2</sub> CH <sub>3</sub>           | X           | -     | -        | -     | -    | Aliphatic             | 124                  | AUU, AUC, AUA                | 5.3                        |
| Lysine        | Lys, K  | -(CH <sub>2</sub> ) <sub>4</sub> NH <sub>2</sub>               | -           | X     | positive | -     | -    | -                     | 135                  | AAA, AAG                     | 5.9                        |
| Leucine       | Leu, L  | -CH <sub>2</sub> CH(CH <sub>3</sub> ) <sub>2</sub>             | X           | -     | -        | -     | -    | Aliphatic             | 124                  | UUA, UUG, CUU, CUC, CUA, CUG | 9.1                        |
| Methionine    | Met, M  | -CH <sub>2</sub> CH <sub>2</sub> SCH <sub>3</sub>              | X           | -     | -        | -     | -    | -                     | 124                  | AUG                          | 2.3                        |
| Asparagine    | Asn, N  | -CH <sub>2</sub> CONH <sub>2</sub>                             | -           | X     | -        | X     | -    | -                     | 96                   | AAU, AAC                     | 4.3                        |
| Proline       | Pro, P  | -CH <sub>2</sub> CH <sub>2</sub> CH <sub>2</sub> -             | X           | -     | -        | X     | -    | -                     | 90                   | CCU, CCC, CCA, CCG           | 5.2                        |
| Glutamine     | Gln, Q  | -CH <sub>2</sub> CH <sub>2</sub> CONH <sub>2</sub>             | -           | X     | -        | -     | -    | -                     | 114                  | CAA, CAG                     | 4.2                        |
| Arginine      | Arg, R  | -(CH <sub>2</sub> ) <sub>3</sub> NH-C(NH)<br>NH <sub>2</sub>   | -           | X     | positive | -     | -    | -                     | 148                  | CGU, CGC, CGA, CGG, AGA, AGG | 5.1                        |
| Serine        | Ser, S  | -CH <sub>2</sub> OH  | -           | X     | -        | X     | X    | -                     | 73                   | UCU, UCC, UCA, UCG, AGU, AGC | 6.8                        |
| Threonine     | Thr, T  | -CH(OH)CH <sub>3</sub>   | X           | X     | -        | X     | -    | -                     | 93                   | ACU, ACC, ACA, ACG           | 5.9                        |
| Valine        | Val, V  | -CH(CH <sub>3</sub> ) <sub>2</sub>                             | X           | -     | -        | X     | -    | Aliphatic             | 105                  | GUU, GUC, GUA, GUG           | 6.6                        |
| Tryptophan    | Trp, W  | -CH <sub>2</sub> C <sub>8</sub> H <sub>6</sub> N               | X           | -     | -        | -     | -    | Aromatic              | 163                  | UGG                          | 1.4                        |
| Tyrosine      | Tyr, Y  | -CH <sub>2</sub> -C <sub>6</sub> H <sub>4</sub> OH             | X           | X     | -        | -     | -    | Aromatic              | 141                  | UAU, UAC                     | 3.2                        |

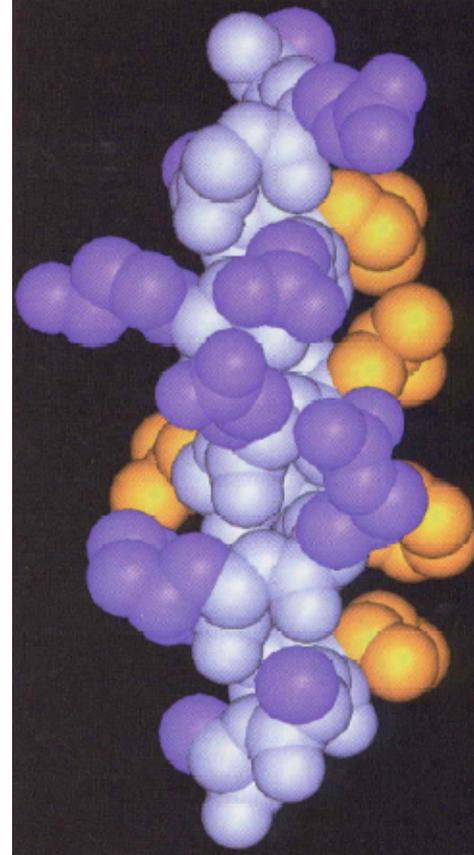
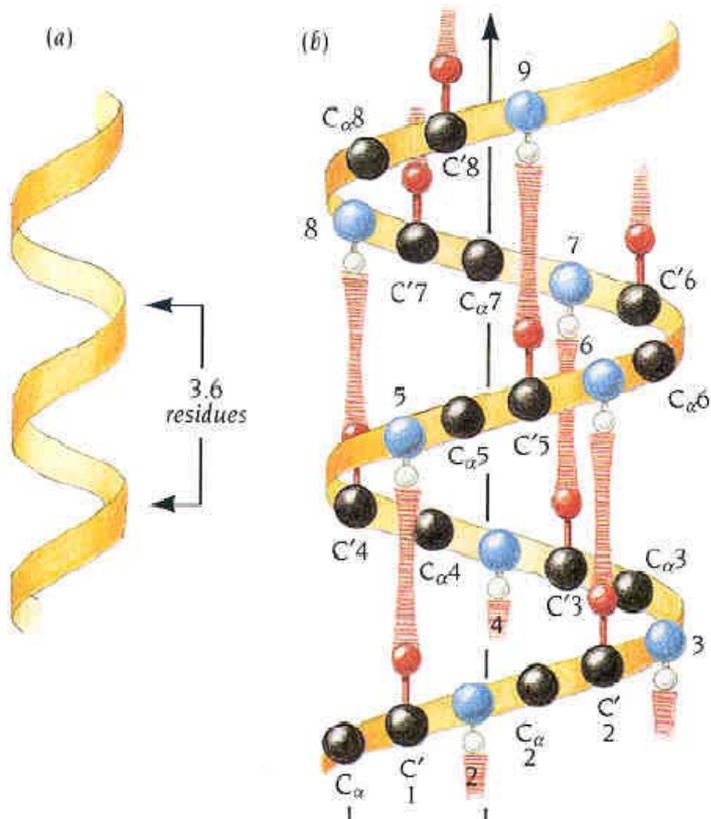
↑  
Hydrophilic

# Protein Secondary Structure

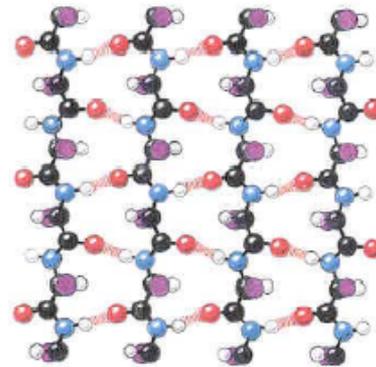
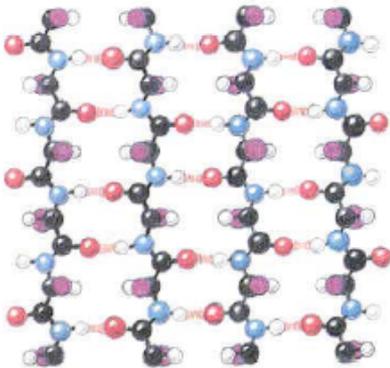
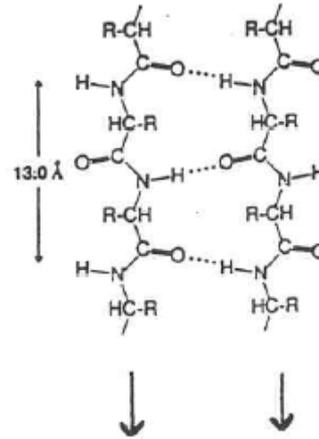
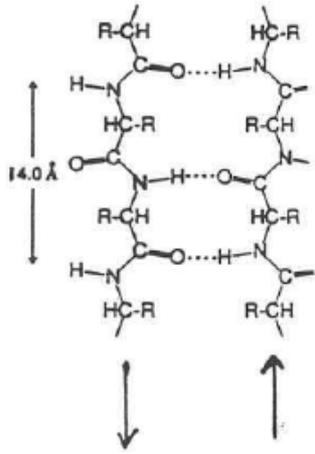
- Determined by hydrogen bond patterns
- 3-Class categories: alpha-helix, beta-sheet, loop (or coil)
- First deduced by Linus Pauling et al.



# Alpha-Helix



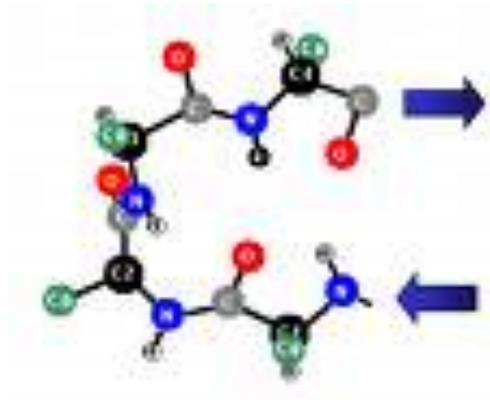
# Beta-Sheet



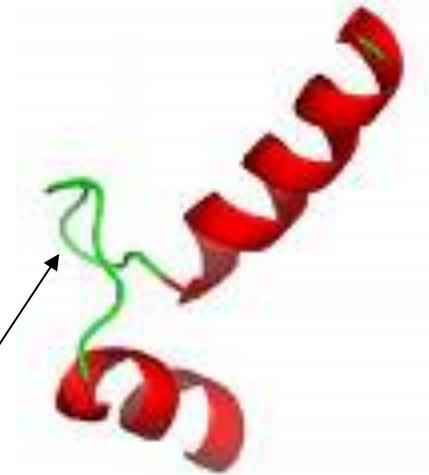
Anti-Parallel

Parallel

# Non-Repetitive Secondary Structure



Beta-Turn



Loop

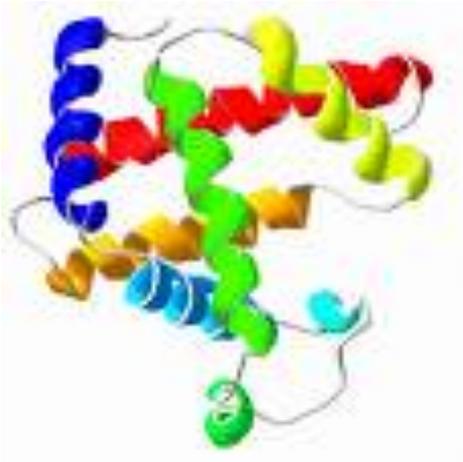
# Tertiary Structure

- John Kendrew et al.,  
Myoglobin
- Max Perutz et al.,  
Haemoglobin
- 1962 Nobel Prize in  
Chemistry

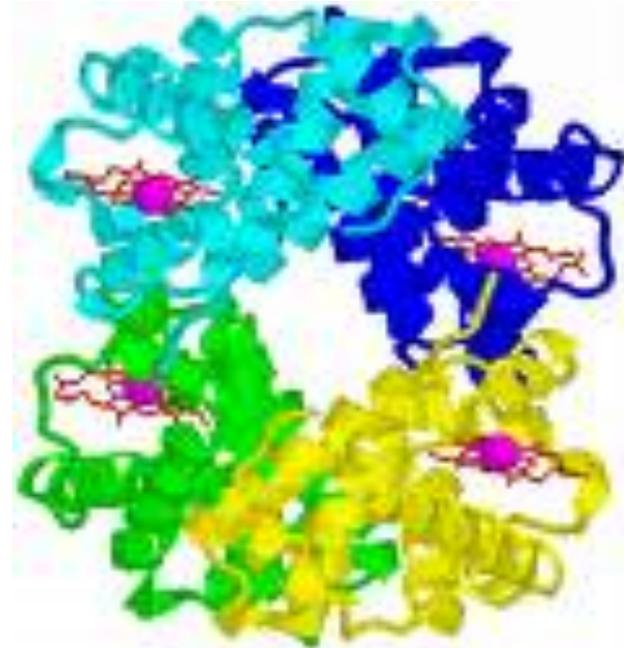


Perutz

Kendrew



myoglobin



haemoglobin

# Quaternary Structure: Complex



G-Protein Complex

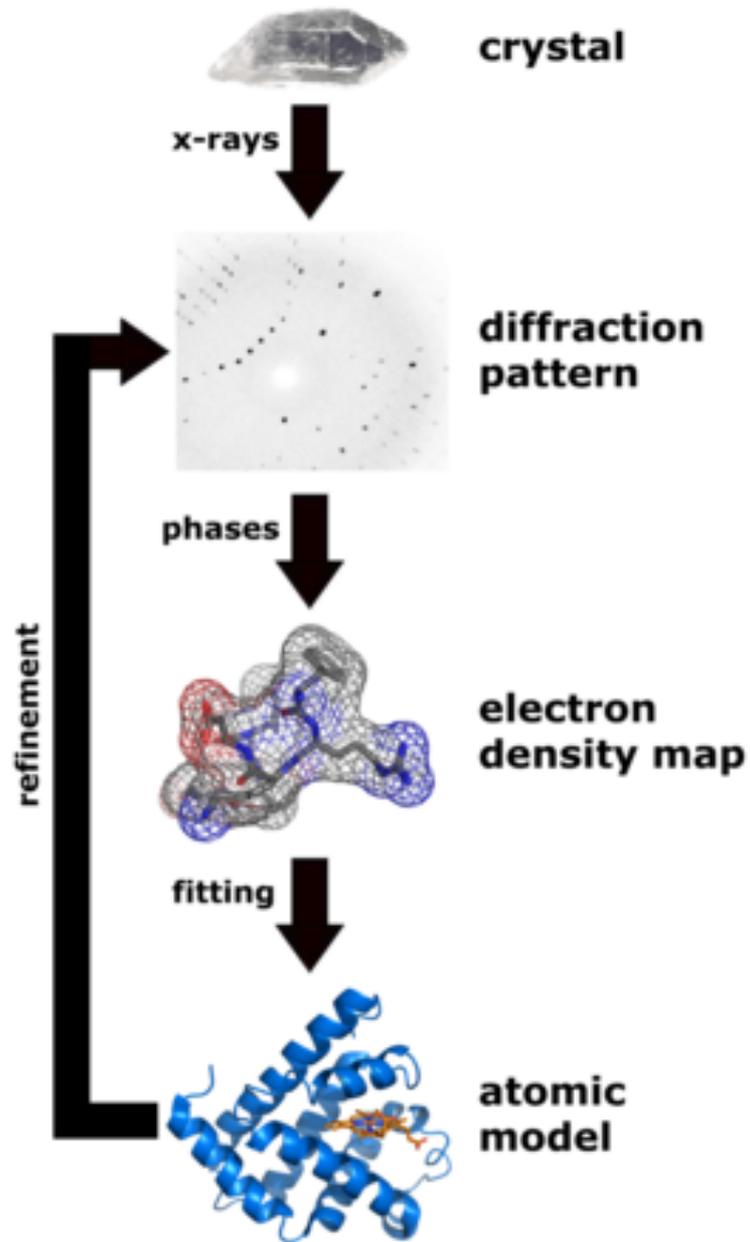
# Anfinsen's Folding Experiment

- Structure is uniquely determined by protein sequence
- Protein function is determined by protein structure



# Protein Structure Determination

- X-ray crystallography
- Nuclear Magnetic Resonance (NMR) Spectroscopy
- X-ray: any size, accurate (1-3 Angstrom ( $10^{-10}$  m)), sometime hard to grow crystal
- NMR: small to medium size, moderate accuracy, structure in solution





[Pacific Northwest National Laboratory](#)'s high magnetic field (800 MHz, 18.8 T) NMR spectrometer being loaded with a sample.

**[Wikipedia, the free encyclopedia](#)**

# Storage in Protein Data Bank

- Home
- Tutorial About This Site
- Getting Started
- Download Files
- Deposit and Validate
- Structural Genomics
- Dictionaries & File Formats
- Software Tools
- Educational Resources
- BioSync
- General Information
- Acknowledgements
- Frequently Asked Questions
- Known Problems
- Report Bugs/Comments

## Welcome to the RCSB PDB

The **RCSB PDB** provides a variety of tools and resources for studying the structures of biological macromolecules and their relationships to sequence, function, and disease.

The RCSB is a member of the **wwPDB** whose mission is to ensure that the PDB archive remains an international resource with uniform data.

This site offers tools for browsing, searching, and reporting that utilize the data resulting from ongoing efforts to create a more consistent and comprehensive archive.

Information about compatible browsers can be found [here](#).

A **narrated tutorial** illustrates how to search, navigate, browse, generate reports and visualize structures using this **NEW SITE**. [This requires the Macromedia Flash player download.]

Comments? [info@rcsb.org](mailto:info@rcsb.org)

### Molecule of the Month: AAA+ Proteases



How would you make a protein cutting machine that would be safe to use inside a cell? Digestive proteases like trypsin and pepsin are small and efficient—they diffuse up to proteins and start cutting. This would never work inside a cell. The cell needs to have more control, so that only obsolete or damaged proteins are destroyed. The

### NEWS

- Complete News
- Newsletter
- Discussion Forum

29-August-2006  
**New RCSB PDB Flyer Available in Print and Online**

Two new brochures are available for RCSB PDB users: The General Information trifold & The Easy Steps for Structure Deposition.



Search database

RCSB PDB : Structure Explorer - Mozilla Firefox

File Edit View Go Bookmarks Tools Help

http://www.rcsb.org/pdb/navbsearch.do?newSearch=yes&isAuthorSearch=no&radioSet=All&inputQuickSearch=1vjg&image.x=0&image.y=0&image=Search

Google pdb

**RCSB PDB**  
PROTEIN DATA BANK

A MEMBER OF THE **PDDB**

An Information Portal to Biological Macromolecular Structures

As of Tuesday Oct 10, 2006 there are 39323 Structures | PDB Statistics

Contact Us | Help | Print Page

PDB ID or keyword Author  SEARCH | Advanced Search

Home Search **Structure** Queries Structure Summary Biology & Chemistry Materials & Methods Sequence Details Geometry

1VJG

[Download Files](#)  
[FASTA Sequence](#)  
[Display Files](#)  
[Display Molecule](#)  
[Structural Reports](#)  
[Structure Analysis](#)  
[Help](#)

**1VJG**

**Title** Crystal structure of putative lipase from the G-D-S-L family from Nostoc sp. at 2.01 Å resolution

**Authors** Joint Center for Structural Genomics (JCSG)

**Primary Citation** Joint Center for Structural Genomics (JCSG) Crystal structure of putative lipase from the G-D-S-L family from Nostoc sp. at 2.01 Å resolution. *To be published*

**History** Deposition 2004-02-19 Release 2004-03-16

**Experimental Method** Type X-RAY DIFFRACTION Data [ EDS ]

| Parameters | Resolution[Å] | R-Value      | R-Free | Space Group          |
|------------|---------------|--------------|--------|----------------------|
|            | 2.01          | 0.175 (obs.) | 0.218  | P 3 <sub>2</sub> 2 1 |

| Unit Cell | Length [Å] | a     | 56.19 | b    | 56.19 | c     | 129.32 |
|-----------|------------|-------|-------|------|-------|-------|--------|
|           | Angles [°] | alpha | 90.00 | beta | 90.00 | gamma | 120.00 |

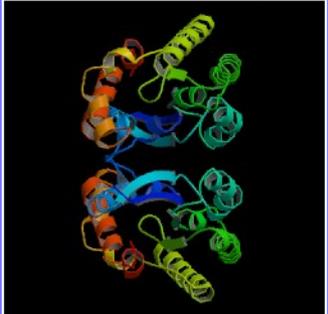
**Molecular Description Asymmetric Unit** Polymer: 1 Molecule: putative lipase from the G-D-S-L family Chains: A

**Functional Class** Structural Genomics Unknown Function

**Source** Polymer: 1 Scientific Name: **Nostoc sp. pcc 7120** Common Name: **Bacteria** Expression system: **Nostoc sp. pcc 7120**

**Images and Visualization**

Biological Molecule



**Display Options**

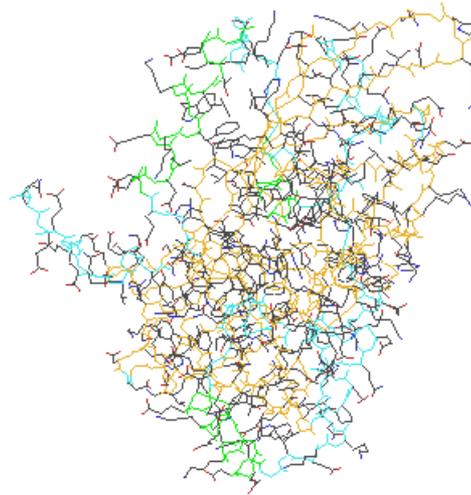
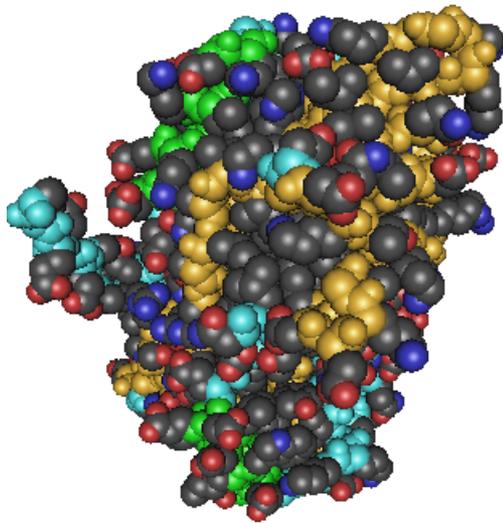
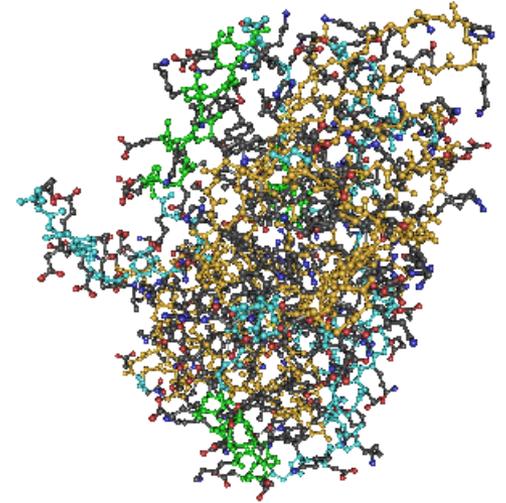
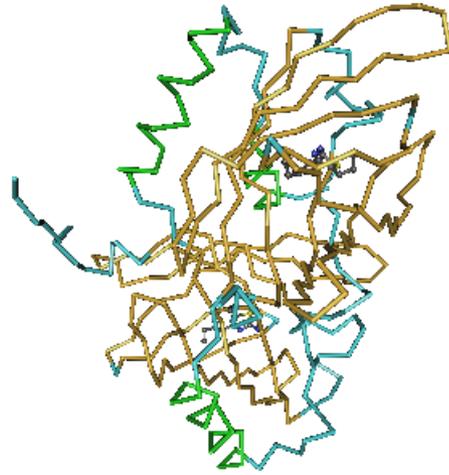
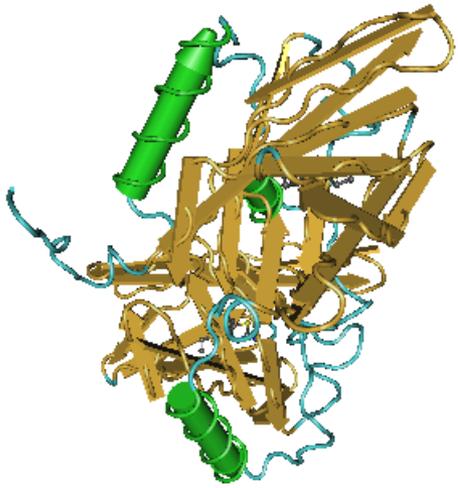
- KiNG
- Jmol
- WebMol
- Protein Workshop
- QuickPDB
- All Images

Done

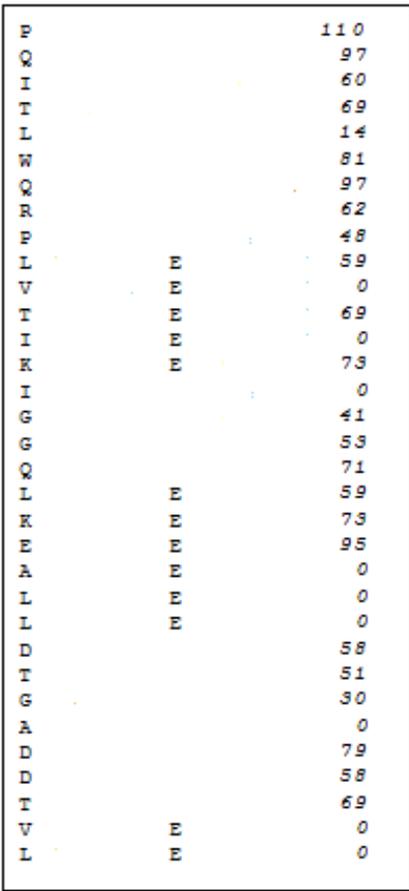
Start | Inboxes - Outlook Express | CAP5937 | slides13 | slides1 | RCSB PDB : Structure ... | 10:42 AM Monday

Search protein 1VJG

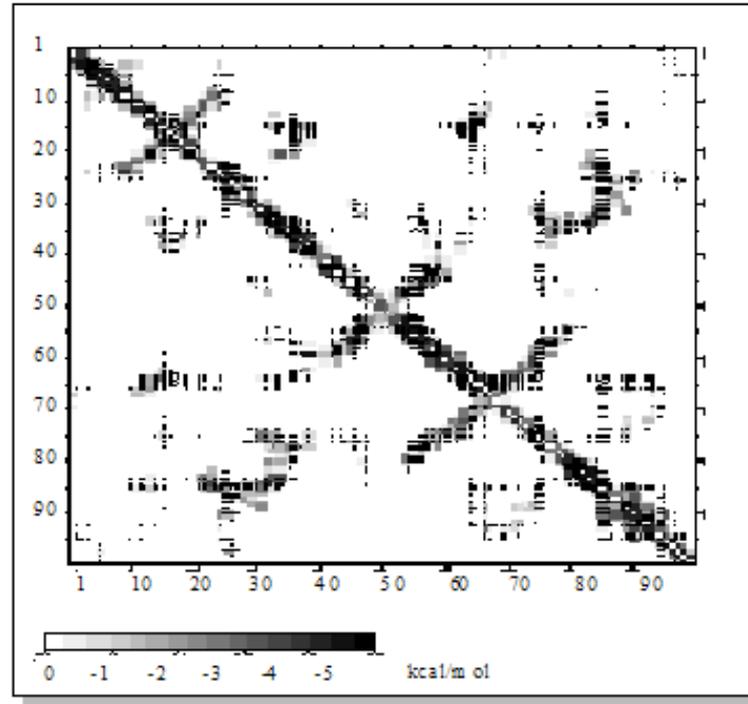




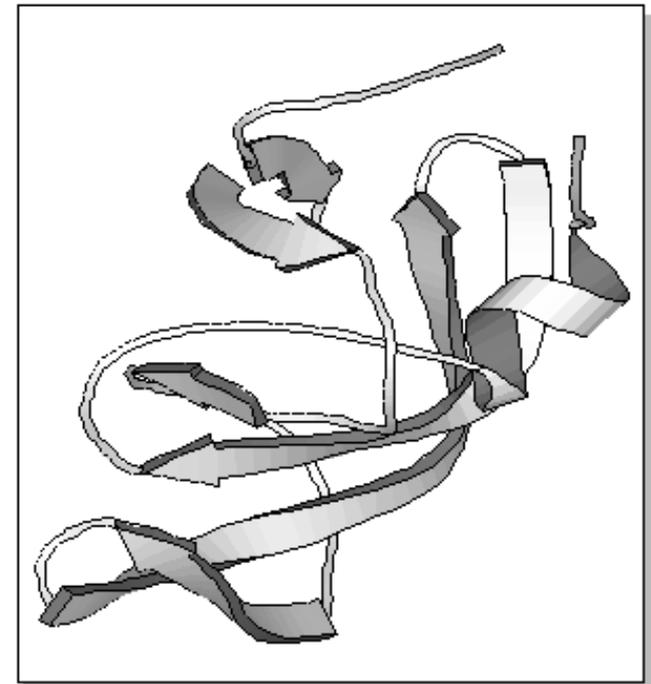
# 1D, 2D, 3D Structure Prediction



1D



2D



3D

# Importance of Computational Modeling

## The Nobel Prize in Chemistry 2013



Photo: A. Mahmoud  
**Martin Karplus**  
Prize share: 1/3



Photo: A. Mahmoud  
**Michael Levitt**  
Prize share: 1/3

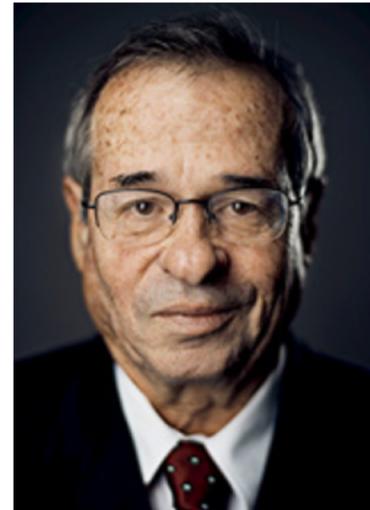
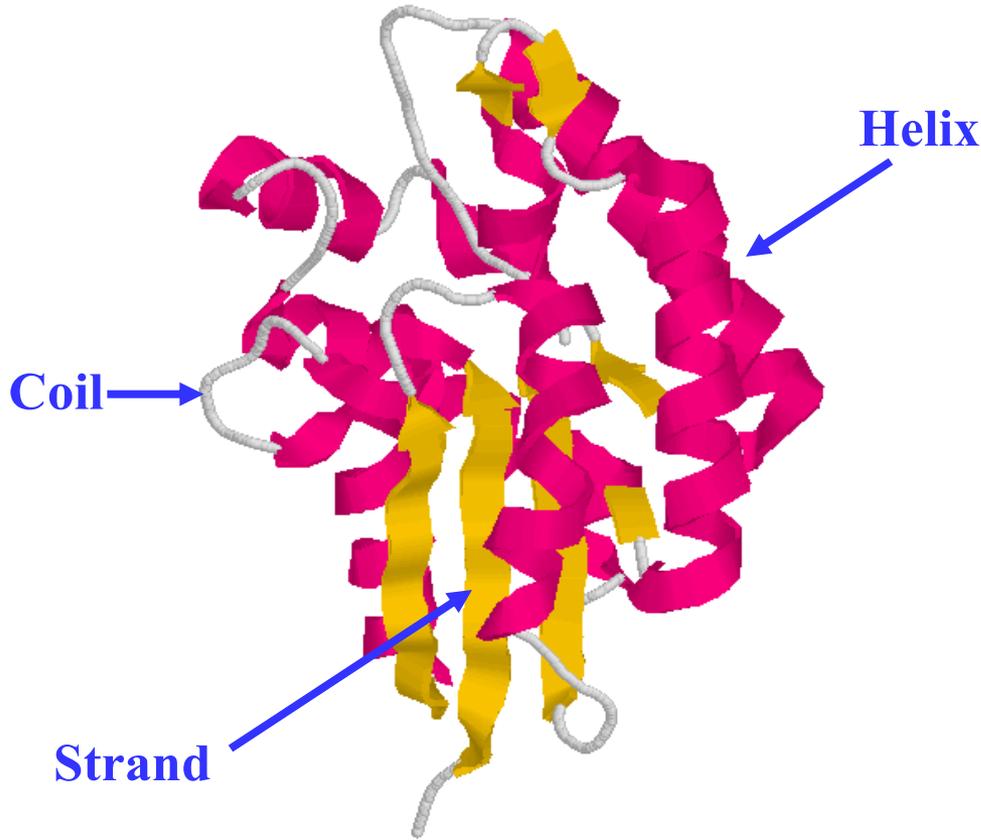


Photo: A. Mahmoud  
**Arieh Warshel**  
Prize share: 1/3

The Nobel Prize in Chemistry 2013 was awarded jointly to Martin Karplus, Michael Levitt and Arieh Warshel *"for the development of multiscale models for complex chemical systems"*.

# 1D: Secondary Structure Prediction



MWLKKFGINLLIGQSV...



**Neural Networks  
+ Alignments**



CCCCHHHHCCCCSSSS...

Accuracy: 78%

# How to Use Neural Network to Predict Secondary Structure

- Create a data set with input sequences ( $x$ ) and output labels (secondary structures)
- Encode the input and output to neural network
- Train neural network on the dataset (training dataset)
- Test on the unseen data (test dataset) to estimate the generalization performance.

# Create a Data Set

- Download proteins from Protein Data Bank
- Select high-resolution protein structures ( $< 2.5$  Angstrom, determined by X-ray crystallography)
- Remove proteins with chain-break (Ca-Ca distance  $> 4$  angstrom)
- Remove redundancy (filter out very similar sequences using BLAST)
- Use DSSP program (Kabsch and Sander, 1983) to assign secondary structure to each residue.

# Train and Test

- Use one data set as training dataset to build neural network model
- Use another data set as test dataset to evaluate the generalization performance of the model
- Sequence similarity any two sequences in test and training dataset is less than 25%.

# Create Inputs and Outputs for Feed-Forward NN for a Single Sequence

Protein Sequence:

MWLKKFGINLLIGQSVQTRSWYYCKRA

SS Sequence:

CCCCHHHHHHEEEEEHHHHEEEEECC



How to encode the input for each position?  
How to encode the output for each position?

# Create Inputs and Outputs for Feed-Forward NN for a Single Sequence

Protein Sequence:

MWLKKFGINLLIGQSVQTRSWYYCKRA

Use 20 inputs of 0s and 1s for each amino acid  
Use 3 inputs to encode the SS alphabet

SS Sequence:

CCCCHHHHHHEEEEEHHHHEEEEECC

100: Helix, 010: Extended strand, 001: Coil

Similarly for 20 different amino acids

# Use a Window to Account for Context

Protein Sequence:

MWLKKFGINLLIGQSVQTRSWYYCKRA

SS Sequence:

CCCCHHHHHHEEEEEHHHHEEEEECC

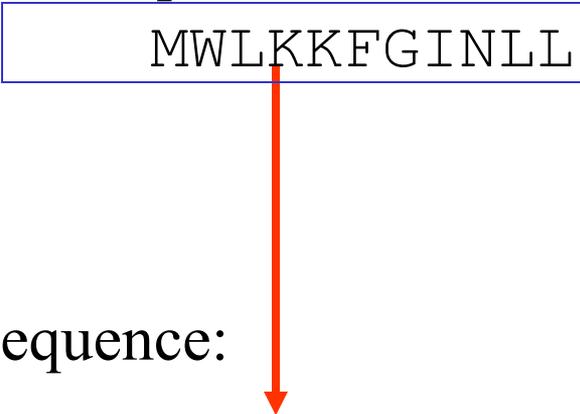


Total number of inputs is window size ( $l$ ) \* 20.  $l$  is a parameter to tune.

# Use an Extra Input to Account for N- and C- Terminal Boundary

Protein Sequence:

MWLKKFGINLLIGQSVQTRSWYYCKRA



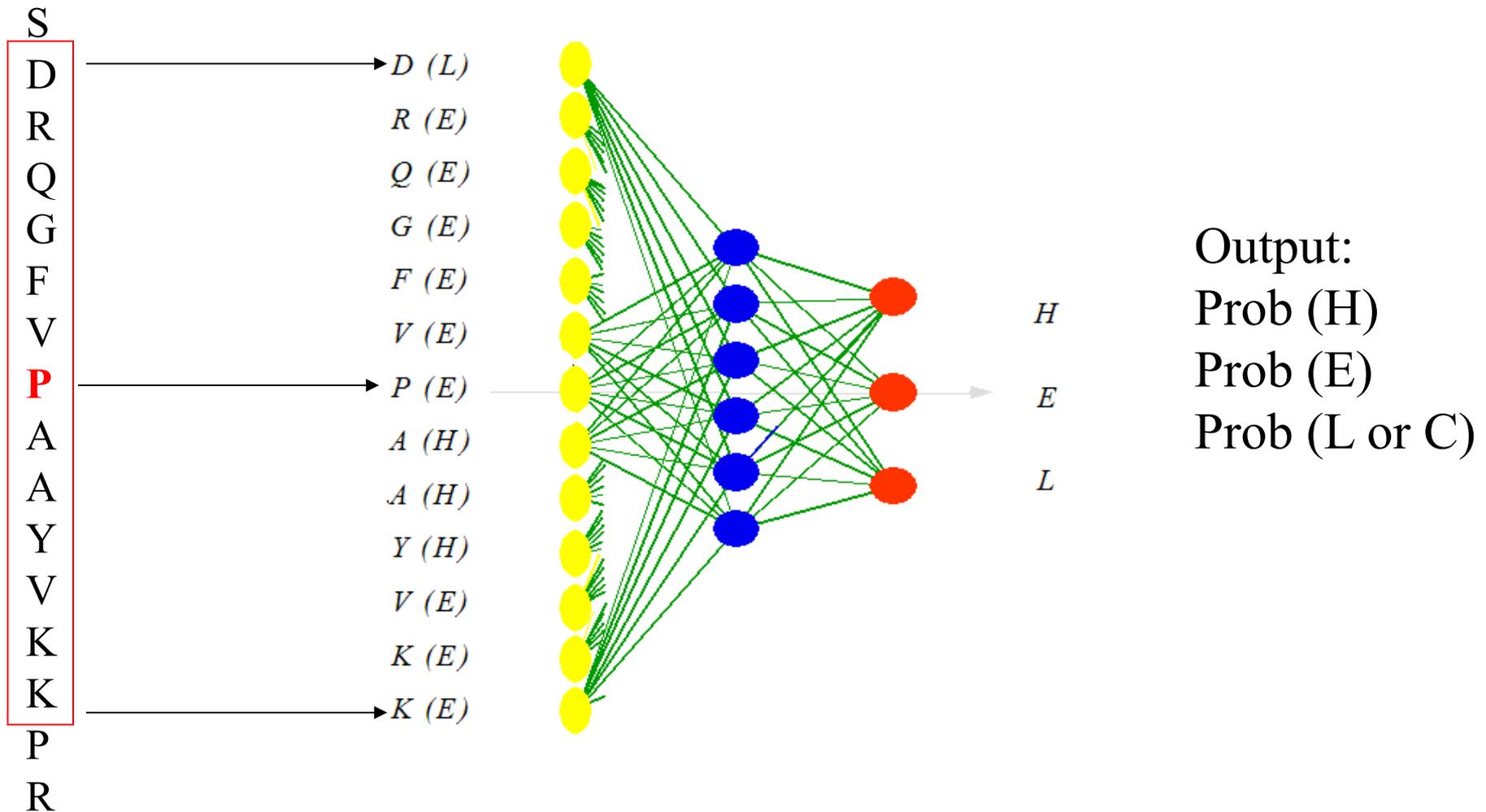
SS Sequence:

CCCCHHHHHHEEEEEHHHHEEEEECC

Add an extra input for each position to indicate if it is out of the boundary of the sequence.

Total number of inputs is window size ( $l$ ) \* 21.  $l$  is a parameter to tune.

# Secondary Structure Prediction (Generation III – Neural Network)



# Evolutionary Information is Important

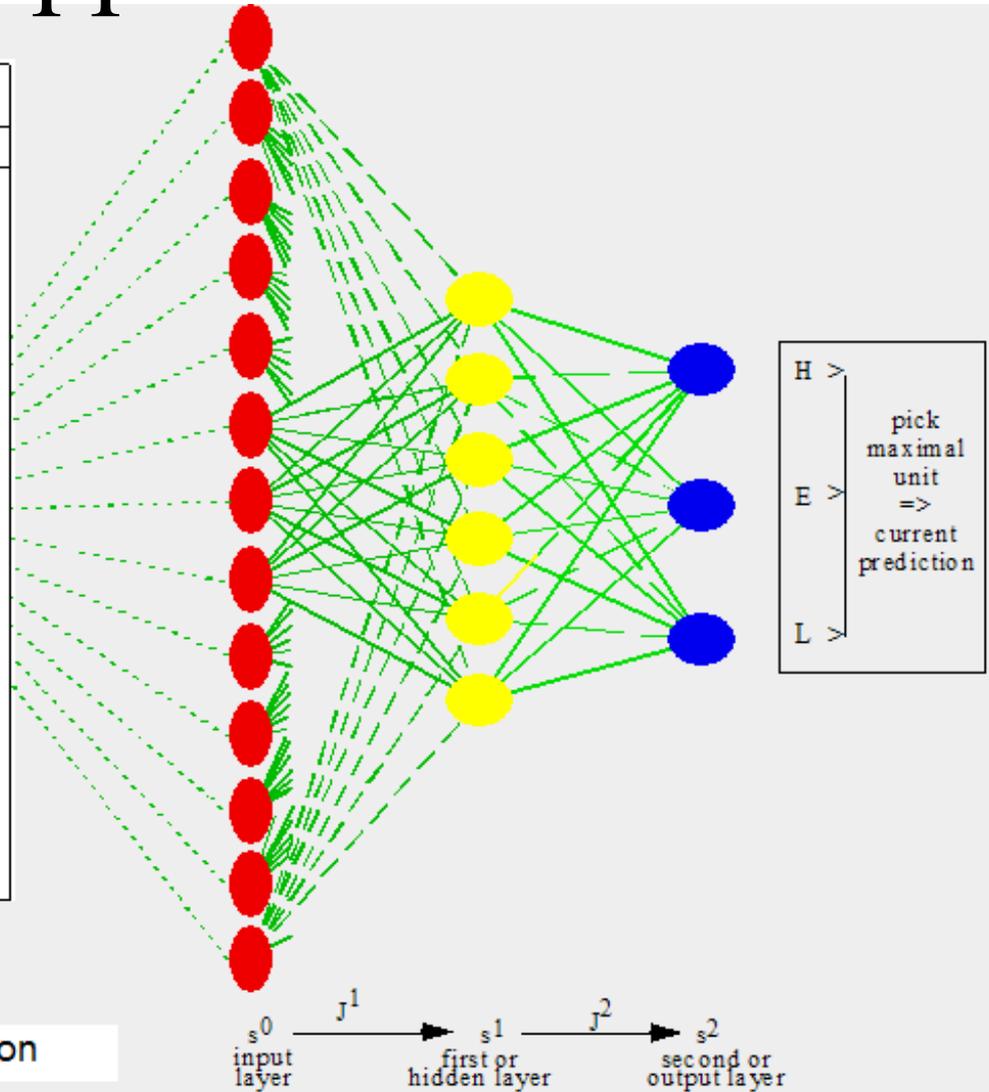
- Single sequence yields accuracy below 70%.
- Use all the sequences in the family of a query sequence can improve accuracy to 78%.
- Structure is more conserved than sequence during evolution. The conservation and variation provides key information for secondary structure prediction.

# How to Find Homologous Sequences and Generate Alignments

- Use PSI-BLAST to search a query sequence against the very large non-redundant protein sequence database (NR database, compiled at NCBI)
- Combine the pairwise alignment between the query sequence and other sequences into a multiple sequence alignment using the query sequence as the center.
-

# PhD Approach

| Protein | Alignments | profile table                       |
|---------|------------|-------------------------------------|
|         |            | GSAPD NTEKQ CVHIR LMYFW             |
| :       | :          | :                                   |
| G       | G G G G    | 5 . . . . .                         |
| Y       | Y Y Y Y    | . . . . . 5 . .                     |
| I       | I I E E    | . . . . . 2 . . . 3 . . . . .       |
| Y       | Y Y Y Y    | . . . . . . . . . . 5 . .           |
| D       | D D D D    | . . . . . 5 . . . . .               |
| P       | P P P P    | . . . . . 5 . . . . .               |
| E       | A E A A    | . . 3 . . . . 2 . . . . .           |
| D       | V V E E    | . . . . . 1 . . 2 . . . 2 . . . . . |
| G       | G G G G    | 5 . . . . . . . . . .               |
| D       | D D D D    | . . . . . 5 . . . . .               |
| P       | P P P P    | . . . . . 5 . . . . .               |
| D       | D T D D    | . . . . . 4 . . 1 . . . . .         |
| D       | N Q N N    | . . . . . 1 3 . . . 1 . . . . .     |
| G       | G N G G    | 4 . . . . . 1 . . . . .             |
| V       | V I V V    | . . . . . . . . . . 4 . 1 . . . . . |
| N       | E P K K    | . . . . . 1 . 1 . 1 2 . . . . .     |
| P       | P P P P    | . . . . . 5 . . . . .               |
| G       | G G G G    | 5 . . . . . . . . . .               |
| T       | T T T T    | . . . . . . 5 . . . . .             |
| D       | E K S A    | . 1 1 . 1 . . 1 1 . . . . .         |
| F       | F F F F    | . . . . . . . . . . . 5 . .         |
| :       | :          | :                                   |



corresponds to 20 input numbers for a position

**Comments: frequency is normalized into probability and sequence needs to be weighted.**

Reference: Rost and Sander. Proteins, 1994.

# PSI-PRED Approach

- PSI-PRED does not use probability matrix instead it uses the another kind of profile: Position Specific Scoring Matrix generated by PSI-BLAST during sequence search.
- The weighting of the sequences is done implicitly by PSI-BLAST.
- The raw PSSM is transformed into values within  $[0,1]$  using sigmoid function.

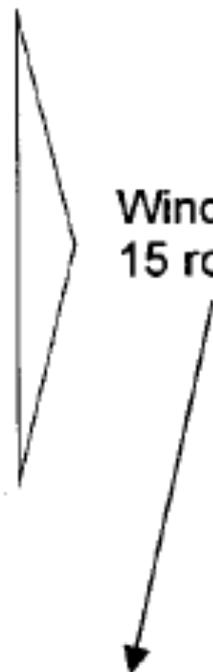
Reference: Jones, Journal of Molecular Biology, 1999.

# PSI-PRED Input

Position-based scoring matrix used

| A  | R  | N  | D  | C  | Q  | E  | G  | H  | I  | L  | K  | M  | F  | P  | S  | T  | W  | Y  | V  |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| -3 | -4 | -4 | -4 | -3 | -4 | -4 | -4 | -2 | -1 | -1 | -4 | -1 | 8  | -5 | -3 | -3 | 0  | 2  | -2 |
| 0  | -1 | -1 | 3  | -4 | 3  | 4  | 1  | -1 | -4 | -4 | 0  | -3 | -4 | -2 | -1 | -2 | -4 | -3 | -3 |
| 0  | -1 | 2  | 1  | -3 | 4  | 0  | -1 | -2 | -4 | -3 | 1  | -2 | -4 | -2 | 2  | 0  | -4 | -3 | -3 |
| -2 | -3 | -4 | -5 | -2 | -3 | -4 | -6 | -4 | 0  | 6  | 0  | 0  | -1 | -4 | -3 | -2 | -4 | -2 | 0  |
| 0  | -3 | -1 | -2 | -3 | 0  | -2 | 4  | -3 | -3 | 0  | -2 | -2 | -4 | -3 | 3  | 1  | -4 | -4 | -3 |
| 0  | 2  | 0  | 4  | -4 | 1  | 2  | 1  | -2 | -4 | -4 | 0  | -3 | -4 | -3 | 1  | -2 | -5 | -4 | -4 |
| -1 | 5  | 3  | -2 | -4 | -1 | -1 | 1  | -2 | -1 | -4 | 1  | -3 | -4 | -3 | 1  | -2 | -5 | -4 | -4 |
| -2 | -3 | -4 | -5 | -3 | -3 | -4 | -5 | -4 | 3  | 4  | -1 | 1  | 2  | -4 | -3 | -2 | -3 | -1 | 0  |
| -2 | 3  | 2  | -2 | -4 | 2  | 1  | -3 | -2 | -3 | -3 | 1  | 1  | -4 | -3 | 2  | 1  | -4 | -3 | -1 |
| 0  | 2  | 3  | 1  | -4 | 0  | 0  | 0  | -2 | -4 | -4 | 1  | -3 | -4 | -3 | 2  | 0  | -5 | -4 | -4 |
| 5  | -3 | -3 | -3 | -2 | -3 | -3 | -2 | -3 | 1  | -2 | -3 | -2 | 1  | -3 | 0  | 1  | -4 | -2 | 0  |
| -1 | -4 | -5 | -5 | -3 | -4 | -4 | -5 | -4 | 3  | 3  | -4 | 2  | 3  | -5 | -3 | -2 | 5  | -1 | 2  |
| 0  | 3  | 3  | 0  | -4 | 3  | 0  | 1  | -2 | -4 | -4 | 1  | -3 | -4 | -3 | 1  | -1 | -4 | -3 | -4 |
| -1 | 0  | 1  | 0  | -4 | 1  | -1 | -1 | -2 | -4 | -3 | 5  | -2 | 0  | -3 | 0  | -2 | -4 | 0  | -3 |
| -2 | -3 | -1 | -5 | -3 | -3 | -4 | -5 | -4 | 3  | 4  | 0  | 4  | 2  | -4 | -3 | -2 | -3 | -2 | 0  |
| 0  | 3  | 0  | -2 | -3 | -1 | 0  | 0  | -2 | 0  | 0  | 1  | 0  | -1 | -3 | 2  | 0  | -4 | -3 | 0  |
| -1 | 1  | 3  | -2 | -4 | 0  | -2 | 4  | -2 | -4 | -4 | 0  | -3 | 0  | -3 | 0  | 0  | -3 | 0  | -4 |

Window of  
15 rows



| A   | R   | N   | D   | C   | Q   | E   | G   | H   | I   | L   | K   | M   | F   | P   | S   | T   | W   | Y   | V   |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.4 | 0.3 | 0.3 | 0.3 | 0.2 | 0.9 | 0.3 | 0.3 | 0.4 | 0.4 | 0.4 | 0.3 | 0.4 | 0.9 | 0.1 | 0.4 | 0.4 | 0.5 | 0.7 | 0.4 |
| 0.3 | 0.2 | 0.3 | 0.8 | 0.4 | 0.3 | 0.7 | 0.1 | 0.6 | 0.2 | 0.4 | 0.3 | 0.5 | 0.2 | 0.1 | 0.4 | 0.8 | 0.2 | 0.3 | 0.2 |
| 0.1 | 0.1 | 0.4 | 0.3 | 0.5 | 0.1 | 0.1 | 0.3 | 0.1 | 0.1 | 0.4 | 0.2 | 0.4 | 0.9 | 0.3 | 0.4 | 0.4 | 0.9 | 0.3 | 0.6 |
| 0.6 | 0.3 | 0.3 | 0.1 | 0.3 | 0.5 | 0.5 | 0.2 | 0.1 | 0.4 | 0.4 | 0.3 | 0.6 | 0.9 | 0.1 | 0.5 | 0.1 | 0.5 | 0.7 | 0.4 |

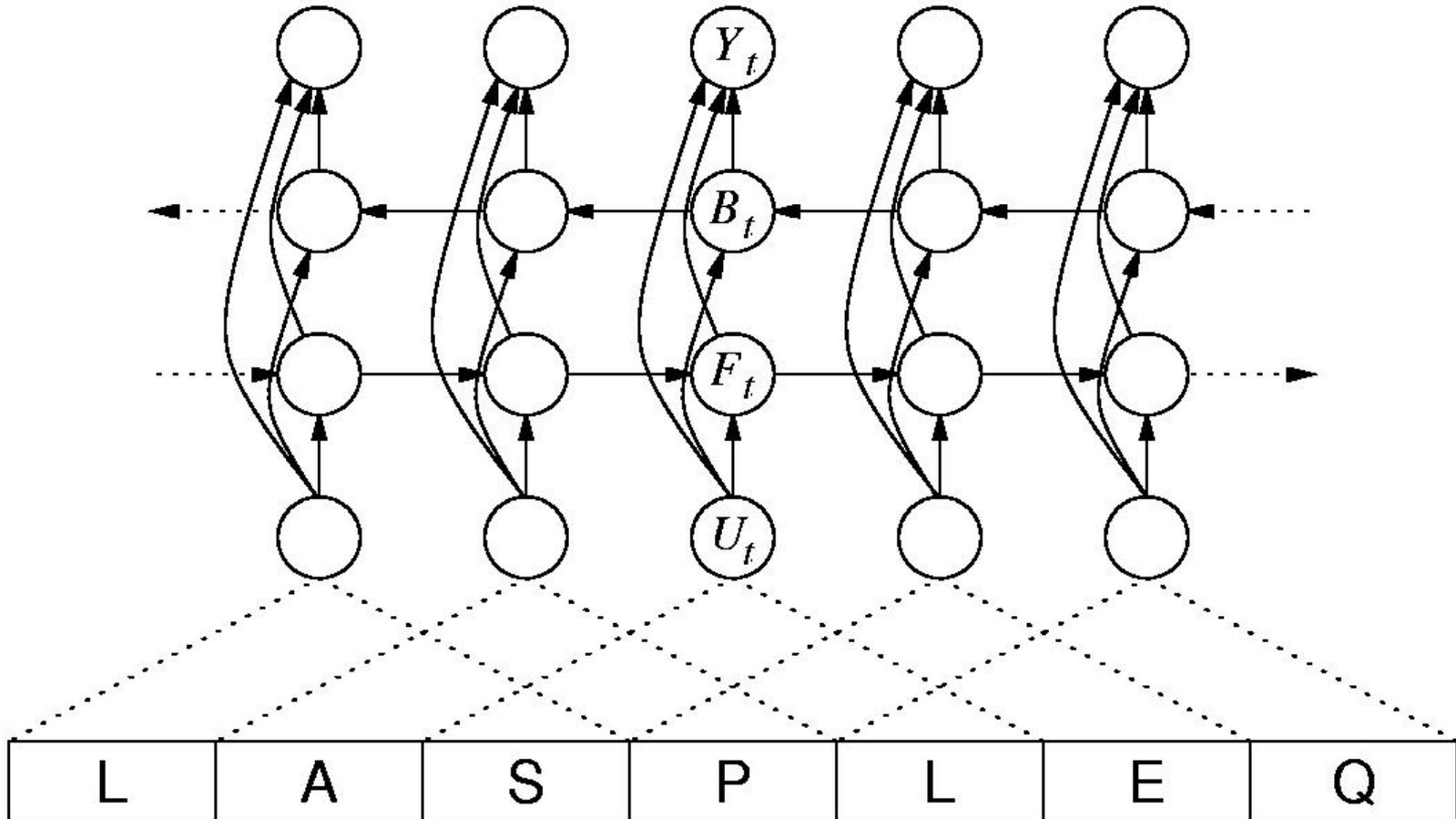
Reference: Jones, Journal of Molecular Biology, 1999.

# SSpro Approach

- SSpro uses probability matrix as inputs
- SSpro uses an information theory approach to weight sequences
- The main novelty of SSpro is to use 1-Dimensional Recurrent Neural Network (1D-RNN)

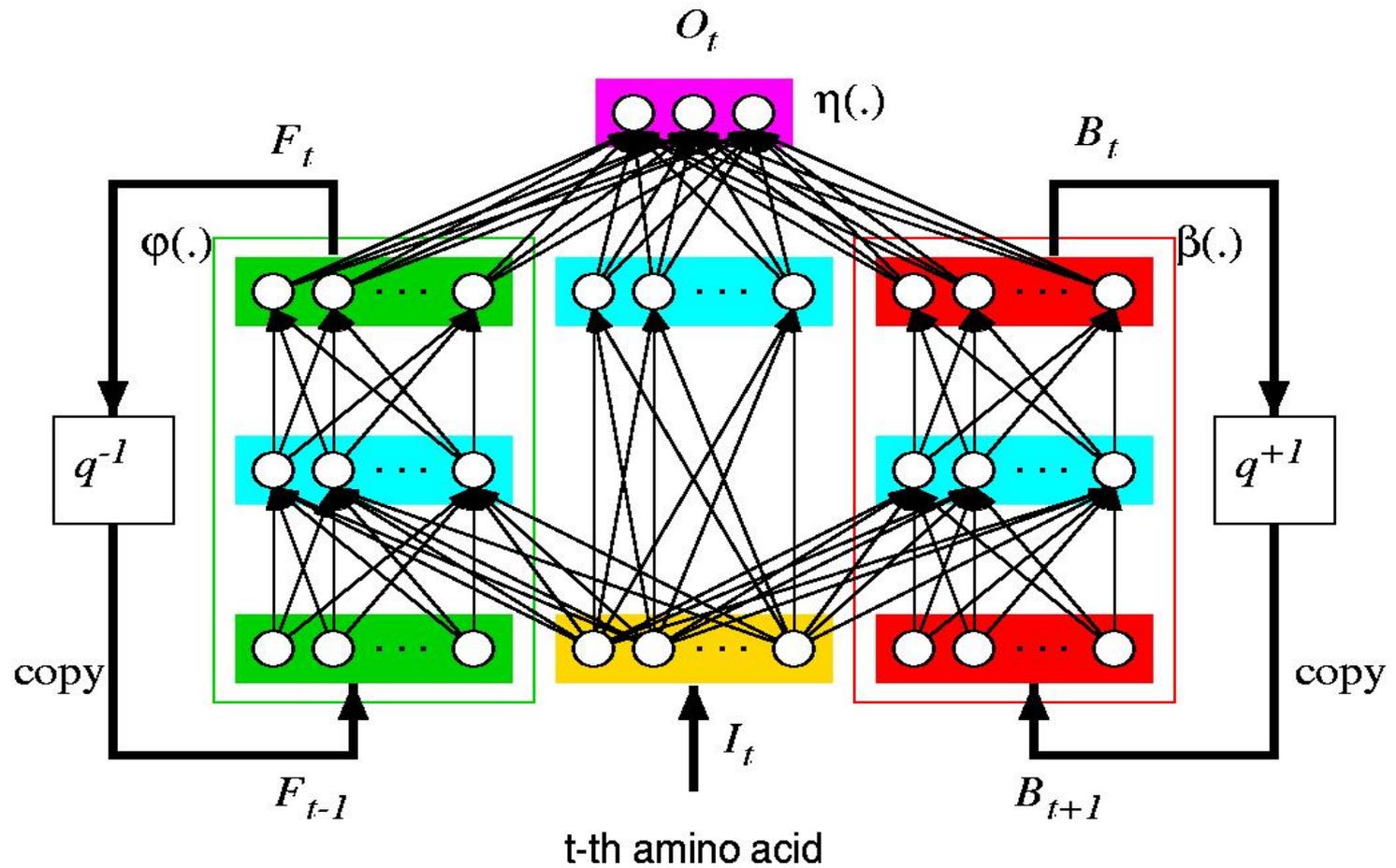
Pollastri et al.. Proteins, 2002.

# Bi-directional Input Output Hidden Markov Model for SS Prediction



Baldi, 2004

# 1D-Recursive Neural Network



Baldi, 2004

# Advantage and Disadvantages of SSpro

- Directly take a sequence with variable length as inputs.
- Hopefully can utilize more information than a fixed-window approach
- More complex, thus harder to implement than feed-forward neural network.

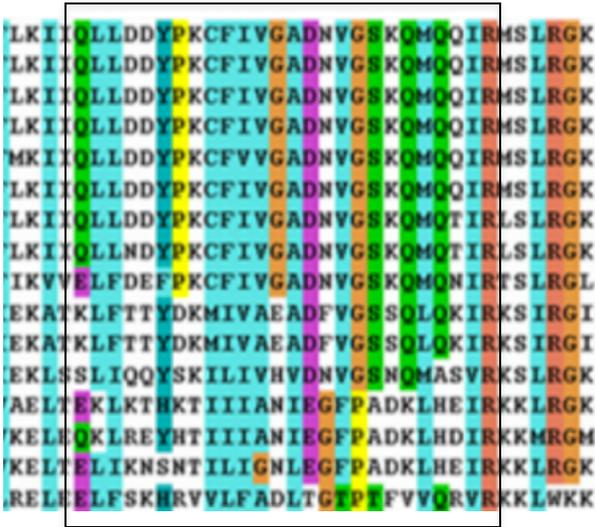
# Second Neural Network to Smooth Output Predictions

- Raw output from one neural network may contain weird predictions such as helix of length 1. But minimum length is 2.
- So use another neural network to smooth output. The inputs are a window of predicted secondary structure. The outputs are the true secondary structures.
- The second neural network makes the predictions more protein-like.

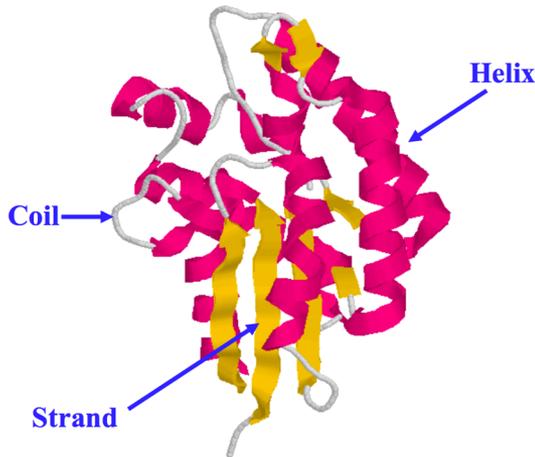
# Deep Learning for Secondary Structure Prediction

M. Spencer, J. Eickholt, J. Cheng. **A Deep Learning Network Approach to ab initio Protein Secondary Structure Prediction.**  
*IEEE Transactions on Computational Biology and Bioinformatics, 2015*

# Deep Learning for Secondary Structure Prediction



Window Size = 19



| Rank | Layers | Nodes <sup>a</sup> | Q <sub>3</sub> (%) | Sov (%) |
|------|--------|--------------------|--------------------|---------|
| 1    | 4      | 600                | 80.24              | 73.07   |
| 2    | 5      | 400                | 79.96              | 72.67   |
| 3    | 3      | 500                | 80.00              | 72.37   |
| 4    | 4      | 300                | 79.84              | 72.48   |
| 5    | 6      | 300                | 79.82              | 72.50   |
| 6    | 3      | 400                | 79.89              | 72.20   |
| 7    | 6      | 400                | 79.87              | 72.15   |
| 8    | 6      | 200                | 79.70              | 72.24   |
| 9    | 5      | 500                | 79.42              | 71.03   |
| 10   | 4      | 400                | 79.35              | 70.83   |
| 11   | 4      | 500                | 79.39              | 70.31   |
| 12   | 3      | 600                | 79.07              | 70.06   |

Spencer et al., *IEEE Transactions on Computational Biology and Bioinformatics*, 2014.

# Deep Learning for Secondary Structure Prediction Project (2<sup>nd</sup> project)

- Training dataset with sequences and secondary structures (1180 sequences) and test dataset (126 sequences). (training data was created by Pollastri et al. and test data was created by Rost and Sander.) ([http://calla.rnet.missouri.edu/cheng\\_courses/mlbioinfo/ss\\_train.txt](http://calla.rnet.missouri.edu/cheng_courses/mlbioinfo/ss_train.txt)(and [ss\\_test.txt](http://calla.rnet.missouri.edu/cheng_courses/mlbioinfo/ss_test.txt))
- Generate multiple alignments using `generate_flatblast.sh` in Pspro 1.2 package ([http://sysbio.rnet.missouri.edu/multicom\\_toolbox/tools.html](http://sysbio.rnet.missouri.edu/multicom_toolbox/tools.html))

# Secondary Structure Prediction Project (continued)

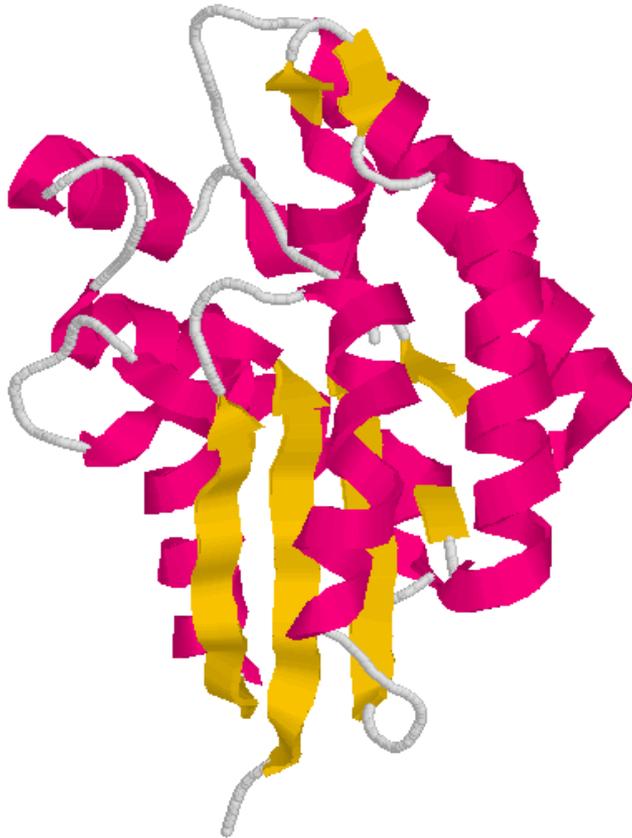
- Generate inputs and outputs (probability matrix or Position Specific Scoring Matrix)
- Develop a deep learning method (e.g. convolutional/recurrent neural network to predict contacts)
- Test the method on test dataset

# References for the Project:

- B. Rost and C. Sander. Combining Evolutionary Information and Neural Networks to Predict Protein Secondary Structure. *Proteins*. 1994.
- D.T. Jones. Protein Secondary Structure Prediction Based on Position-Specific Scoring Matrices. *JMB*. 1999.
- G. Pollastri, D. Przybylski, B. Rost, and P. Baldi. Improving the prediction of protein secondary structure in three and eight classes using recurrent neural networks and profiles. *Proteins*. 2002
- J. Cheng, A. Randall, M. Sweredoski, and P. Baldi. SCRATCH: a protein structure and structural feature prediction server. *Nucleic Acids Research*. 2005.
- Spencer, Matt, Jesse Eickholt, and Jianlin Cheng. "A deep learning network approach to ab initio protein secondary structure prediction." *IEEE/ACM transactions on computational biology and bioinformatics (TCBB)* 12.1 (2015): 103-112
- Wang, Sheng, et al. "Protein secondary structure prediction using deep convolutional neural fields." *Scientific reports* 6 (2016): 18962

# 2D: Contact Map Prediction

## 3D Structure

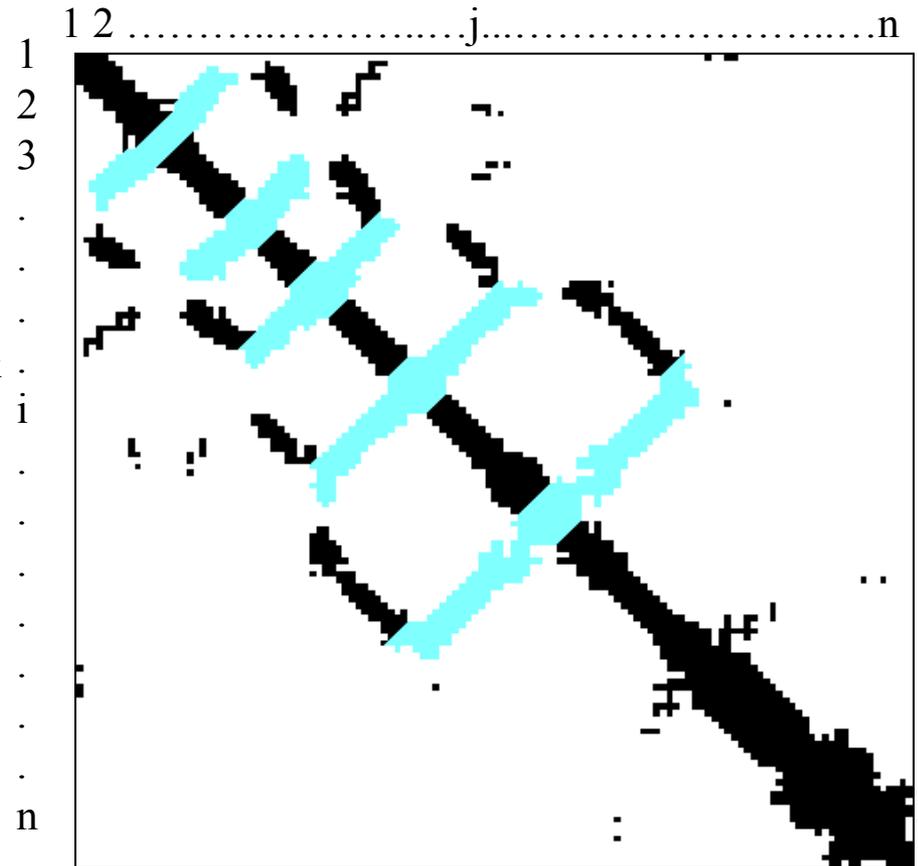


2D-Recursive  
Neural Network



Support Vector  
Machine

## 2D Contact Map



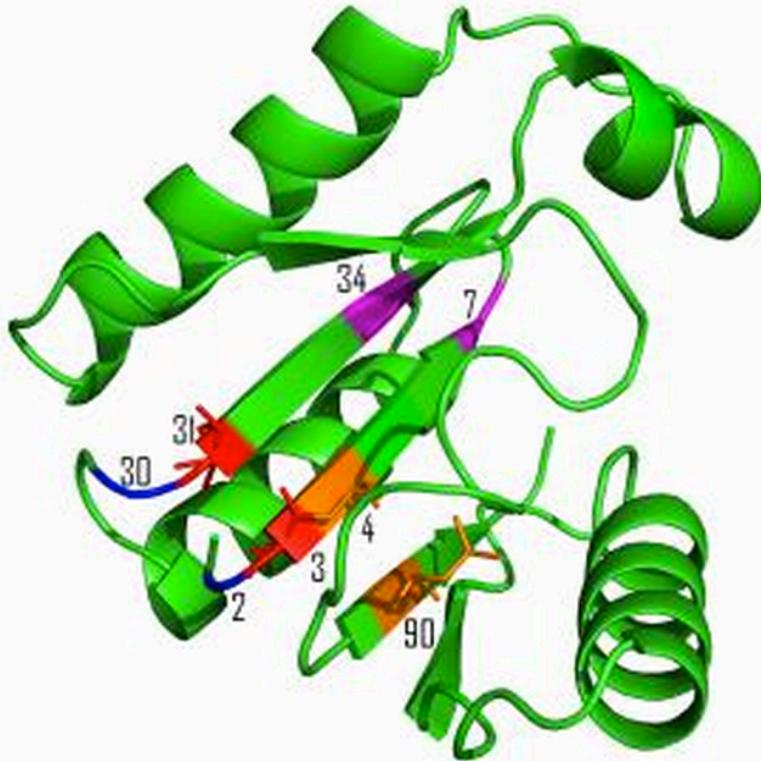
Distance Threshold =  $8\text{\AA}$

Cheng, Randall, Sweredoski, Baldi. *Nucleic Acid Research*, 2005  
Cheng and Baldi. *BMC Bioinformatics*, 2007.

# Residue-Residue Contact Prediction

## 1D Sequence

SDDEVYQYIVSQVKQYGI EPAELLSRKYGDKAKYHLSQRW

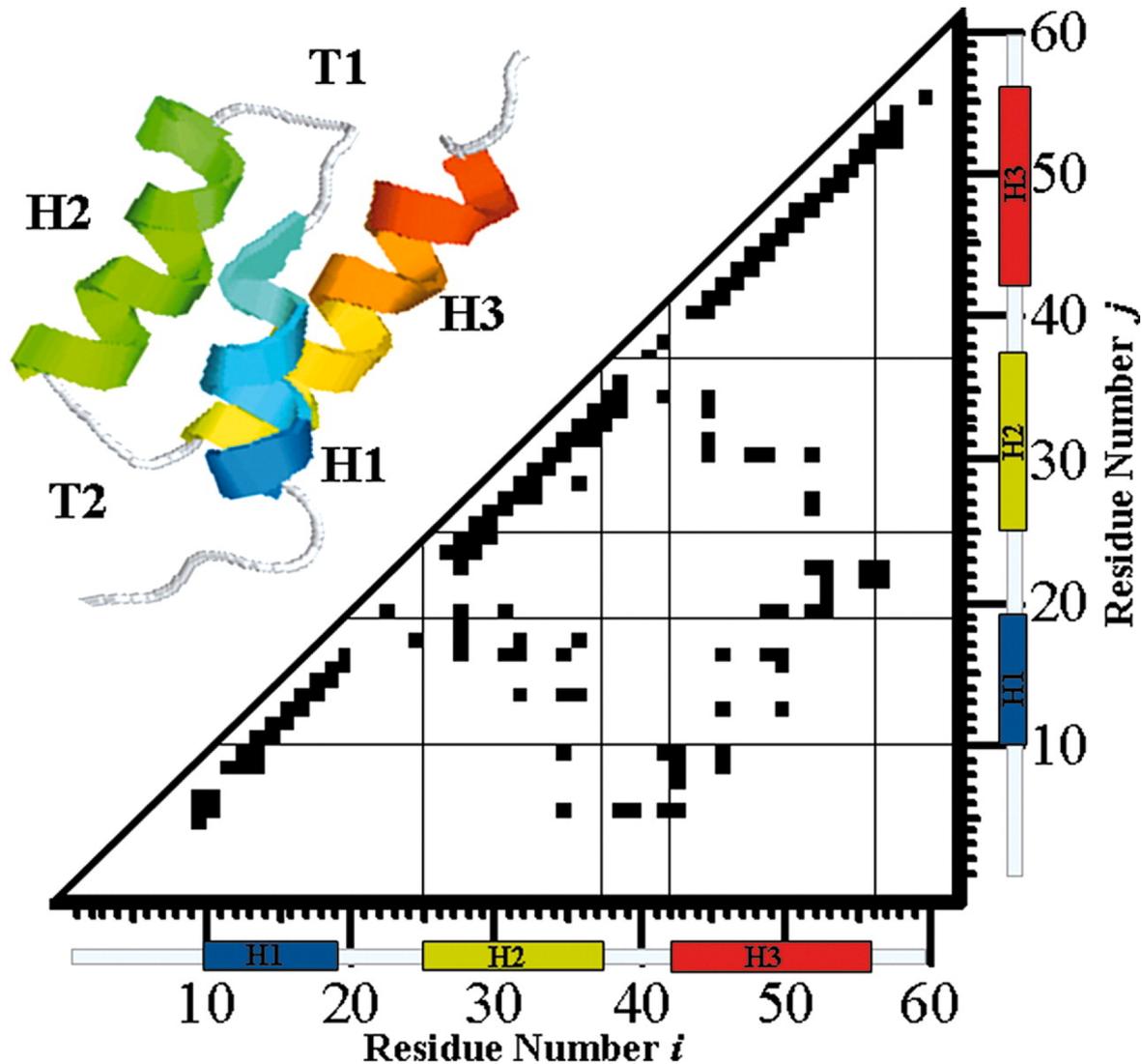


## 3D Structure

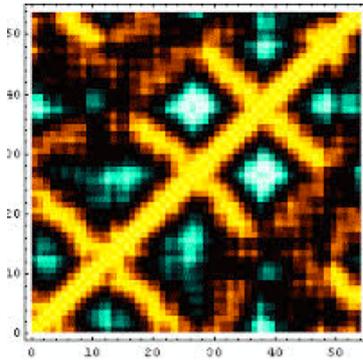
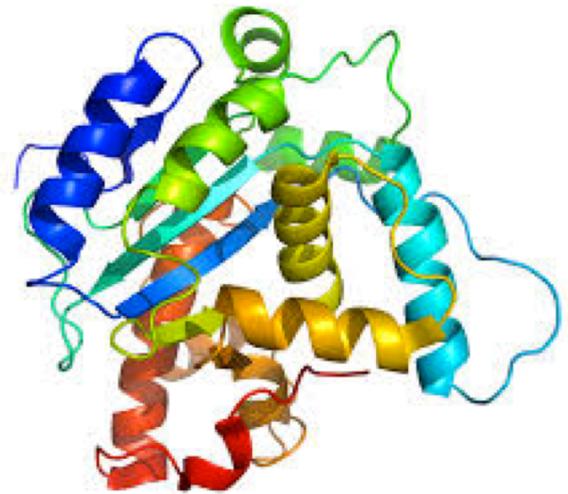
### Objective:

Predict if two residues ( $i, j$ ) are in contact (spatially close), i.e.  $\text{distance}(i, j) < 8 \text{ Angstrom}$

# Visualization of Contact Map



# Contact-Based Structure Prediction



# A Binary Classification Problem

*i* SDDEVYQYI**V**SQVK QYGI EPCSAELLSRKYGD KAK**Y**HLS QRW *j*

Residue identity, secondary structure, solvent accessibility, ...

A Vector of ~400 Features (numbers between 0 and 1)

Probability that **V** and **Y** are in contact?

# Input Features

*i*

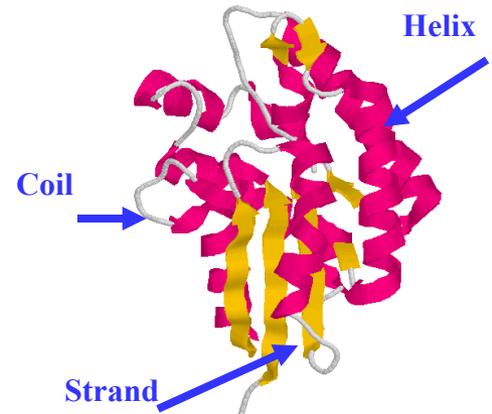
*j*

SDDEVYQYI**V**SQVKQYGI**E**PCSAELLSRKY**G**DKAK**Y**HLSQRW

20 binary numbers

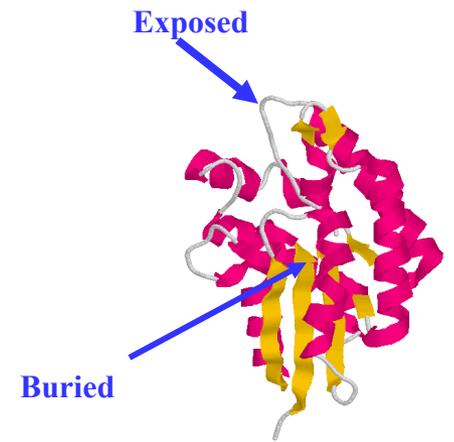
A 10000000000000000000  
 C 01000000000000000000  
 D 00100000000000000000  
 .  
 .  
 .  
 .  
 .  
 .  
 .  
 .  
 Y 000000000000000000001

3 numbers



Helix 100  
 Strand 010  
 Coil 001

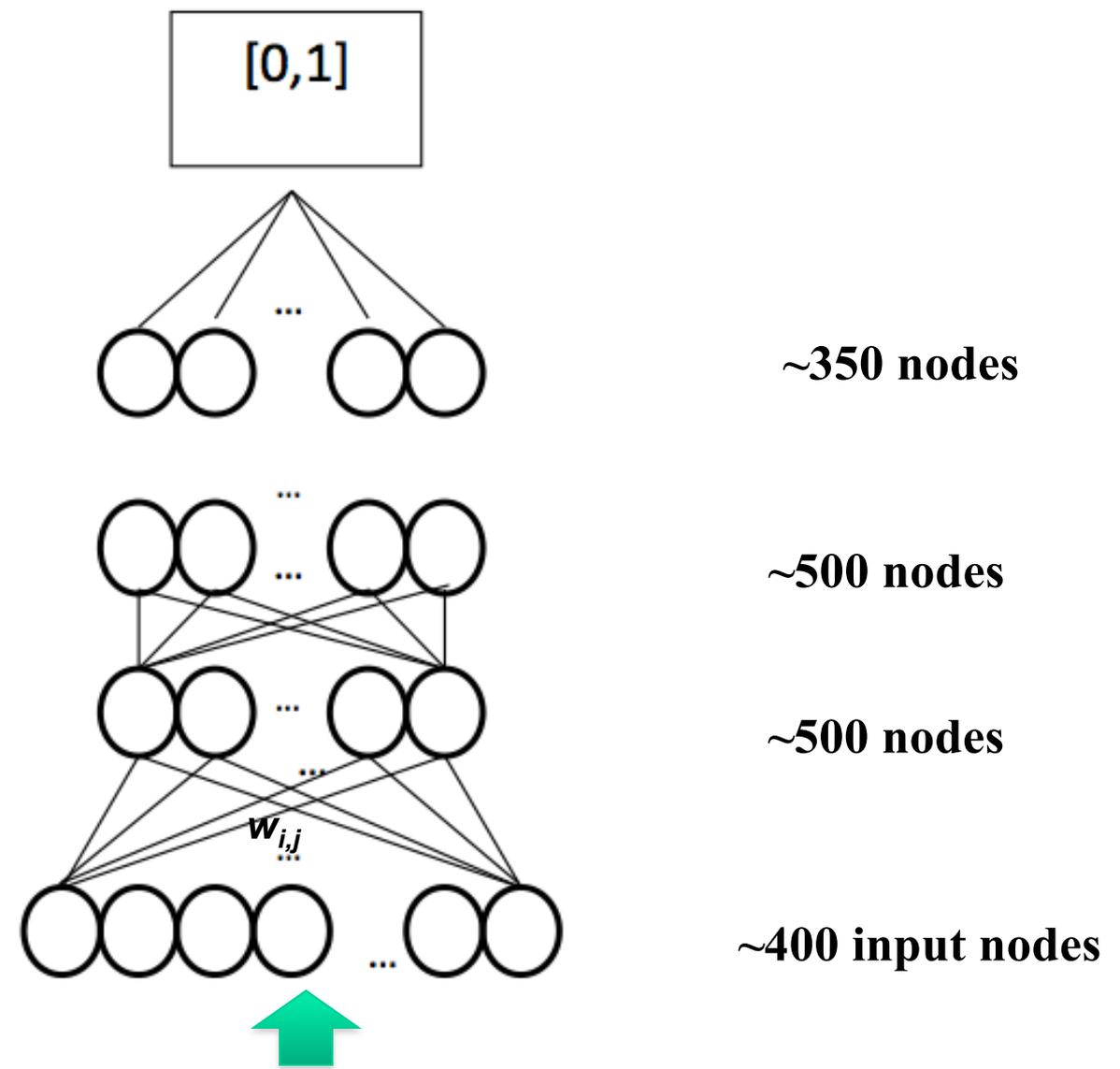
2 numbers



Exposed 10  
 Buried 01

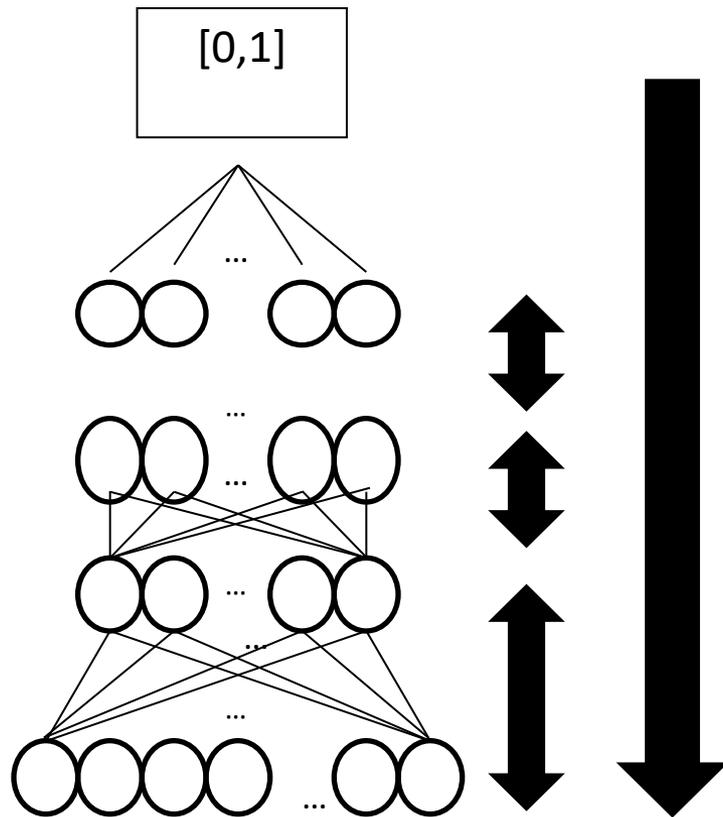
25 \* 18 = 400 features for a pair (*i*, *j*)

# Deep Learning Network Architecture



A Vector of  $\sim 400$  Features (numbers between 0 and 1)

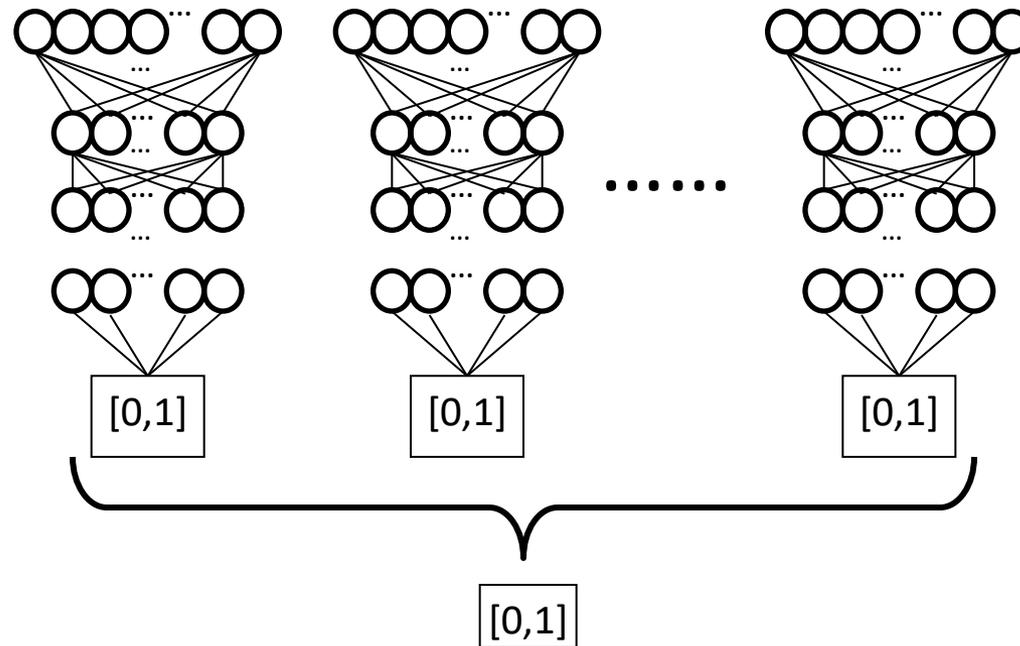
# Training a Deep Network



**1239 Proteins for Training  
Residue Pairs ( $|i-j| \geq 6$ );  
Millions of Residue Pairs**



# Boosted Ensembles for Contact Prediction

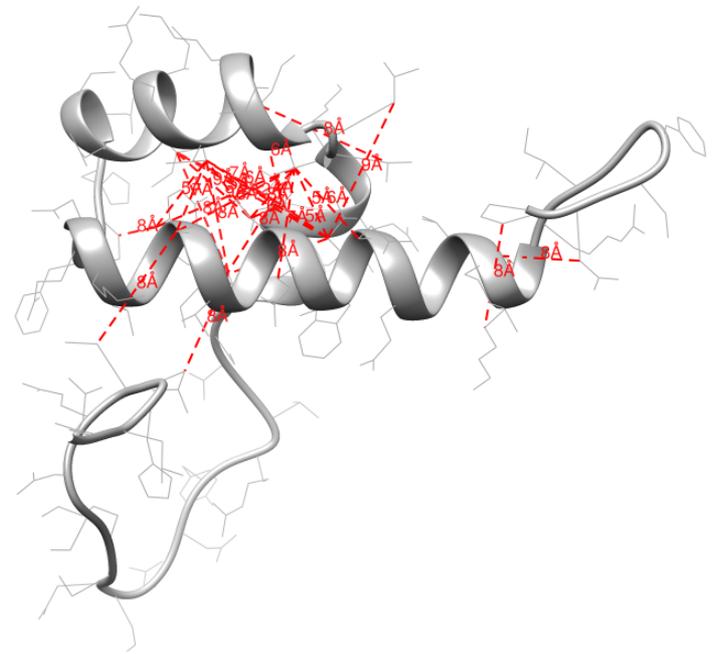


**Final output of ensemble  
is the weighted sum of  
individual DN outputs.**

# Results on Test Data Set (196 Proteins)

| Metric                                   | Acc. L/5 | Acc. L/5<br>(one shift) |
|--|----------|-------------------------|
| Short Range<br>( $6 \leq  i-j  < 12$ )   | 0.51     | 0.79                    |
| Medium Range<br>( $12 \leq  i-j  < 24$ ) | 0.38     | 0.65                    |
| Long Range<br>( $ i-j  \geq 24$ )        | 0.34     | 0.55                    |

## An Example:

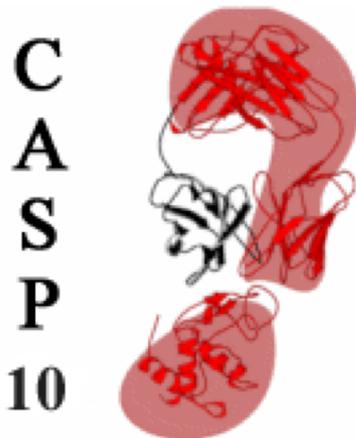


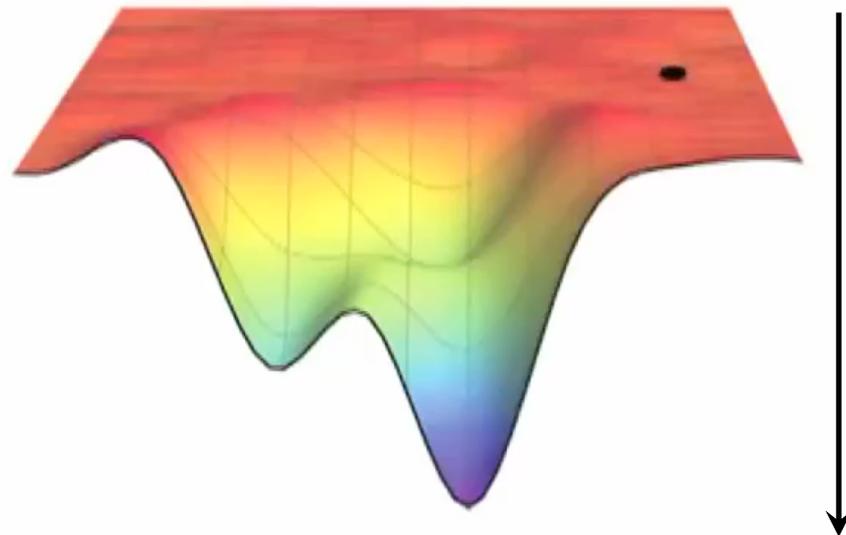
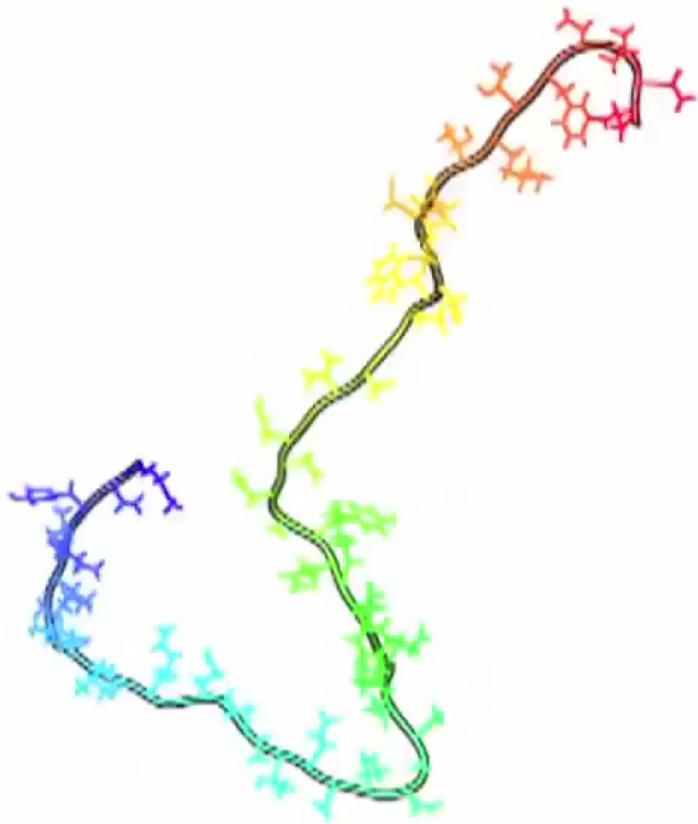
# Blind Test on CASP10 Targets

**Exact match (96 proteins, long-range contacts)**

| Method       | Acc. L/5    |
|--------------|-------------|
| <b>DNcon</b> | <b>0.30</b> |
| SVMcon       | 0.19        |

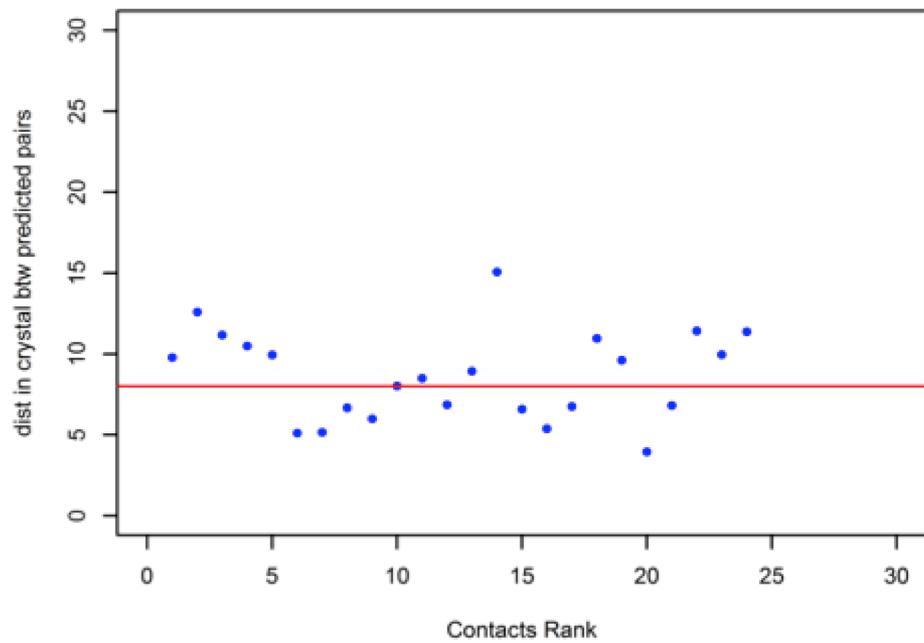
→ 9-fold better than random



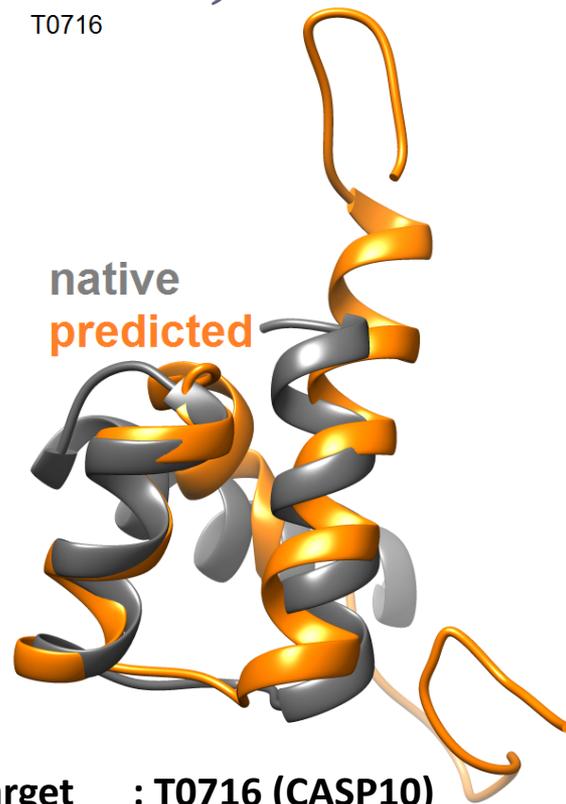


**Contact-based energy landscape**

# 3D Reconstruction from Predicted Contacts (CASP Target T0716)

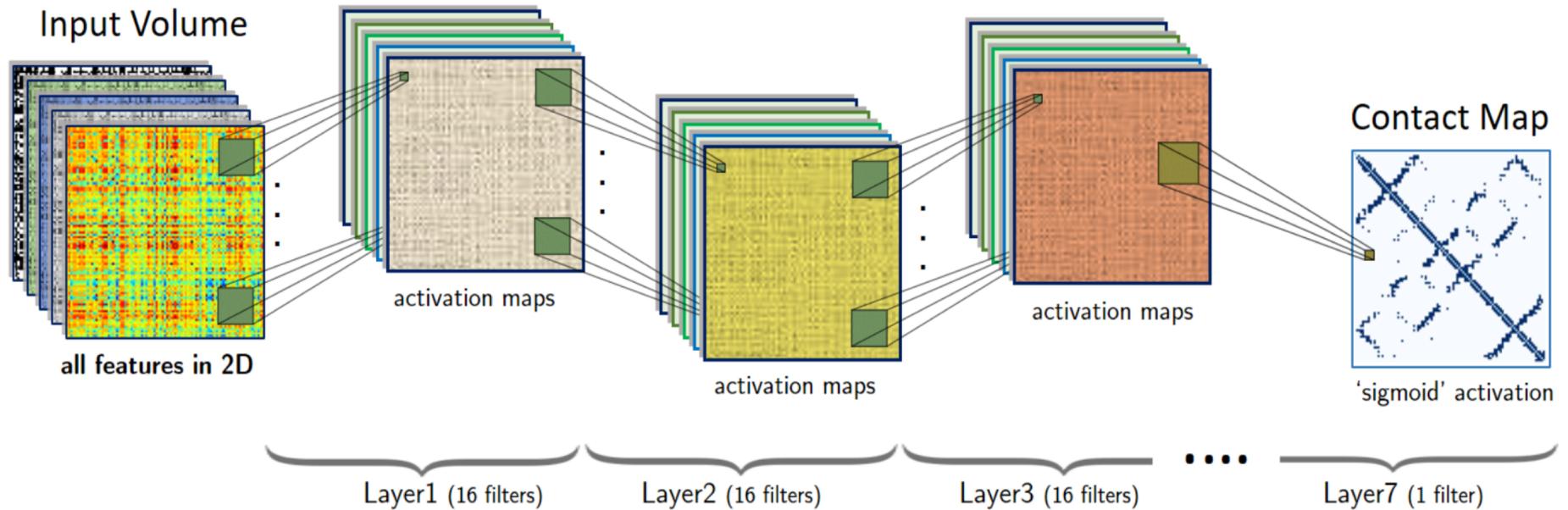


T0716



**Target : T0716 (CASP10)**  
**Length : 71**  
**RMSD : 4.3A**  
**GDT-TS : 0.58**  
**Contacts : DNcon (filtered and selected 0.4L)**

# 2D Convolutional Network for Contact Prediction



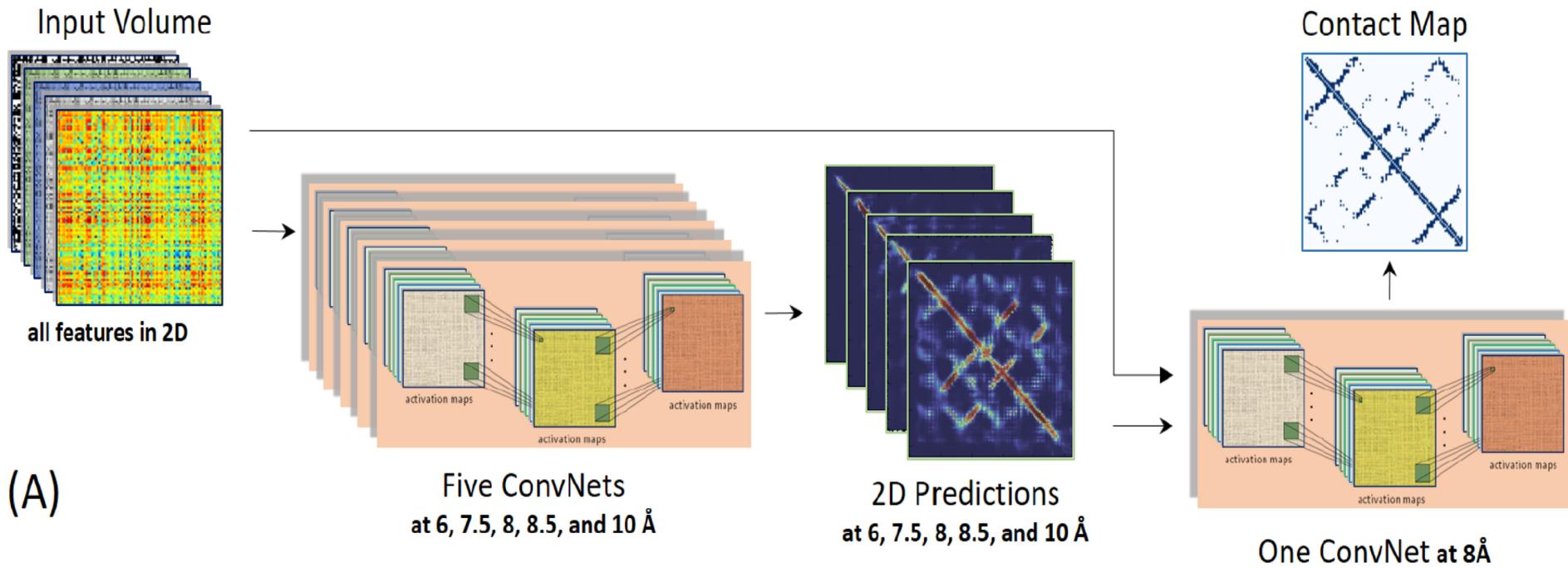
$i$   $j$

ALTLHRDRFTTARRTAPIPQLQCLGGSAGCPAHIPEIVQCRNKGWDGEFDVQWECKAELDT  
 VLTLHRGRYTTARRTAAPVQLQCIGGSAGCS-DIPEVVQCYNRGWDGYDVQWQCKADLEN  
 TITLYADRYTNARRSAPVQLKCI GGNAGCHAMVPQVVQCHNRGWDGLDVQWECRVDMND  
 AITLYADRYTNARRSAPVQLKCI GGSAGCHTMVPQVVQCHNRGWDGEFDVQWECKVDMND  
 VLTLYRGRYTTARRSSPVPQLQCIGGSAGCGSFPEVVQCYNRGSDGIDAQWECKADMND  
 VLTLYKGYTTARRSSAVPQLQCVGGSAGCGSFPEVVQCKNKGWDGVDVAQWECKTDMND  
 VLTLYRGLYTTARRSSPVPQLQCVGGSAGCHAFVPEVVQCQNKGWDGMDIQWECRTDMND  
 TLTLYRGRYTTARRSSPVPQLRCVGGGAGCQAFVPEVVQCQNRGWDGVDVQWECKTDMND  
 ALTLYKNRYTTARRASPVLPQLQCVGGSAGCQAFVPEVVQCQNKGWDGVDVQWECRTDMND  
 VLTLYKGRYTTARRSSPVLQLQCAGGTAGCGSFVPEVVQCYNRGSDGIDTQWECKADMND  
 AITLHKGMTTGRRVSPVIFQLKCVGG-SAKGAFTPKVVQCANQGFDSVDVQWRCADLPH

## Input Matrix

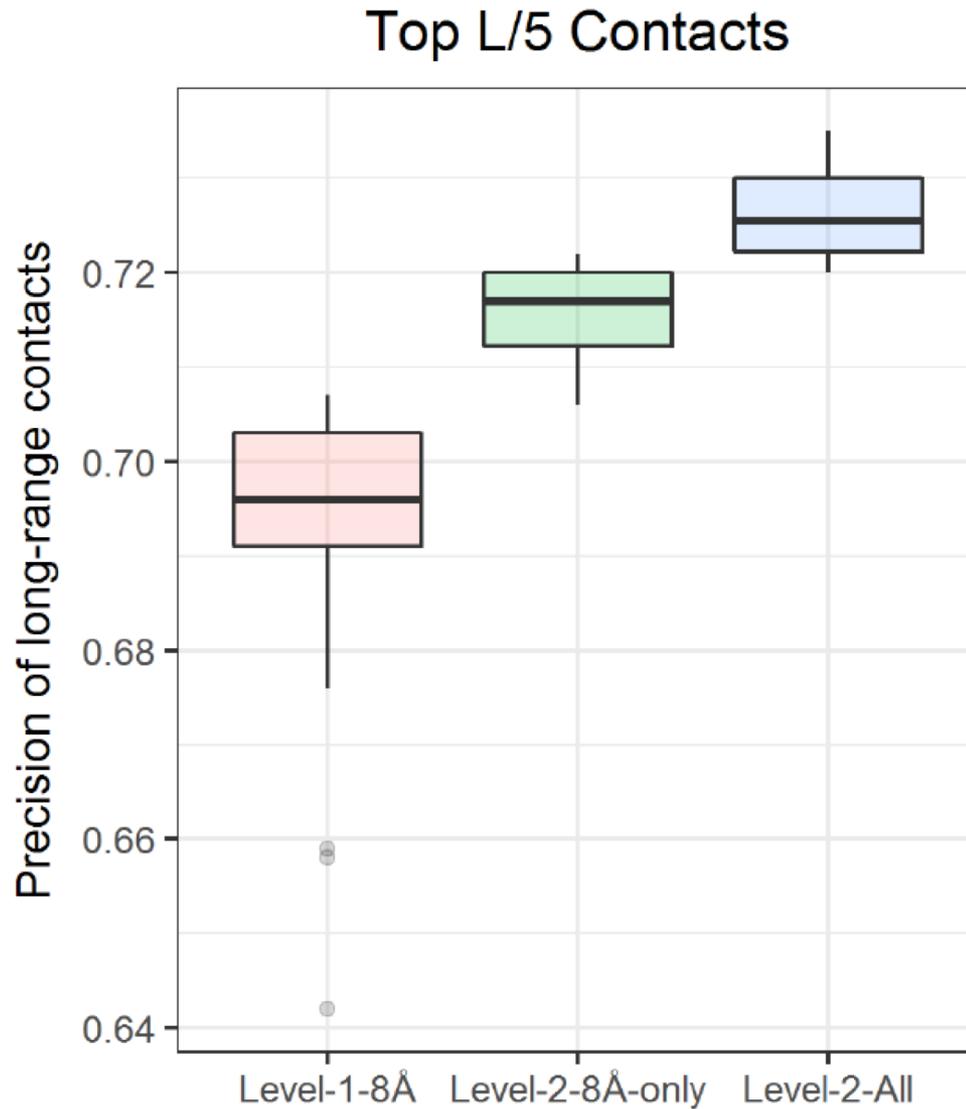
- Co-evolution
- Mutation information
- Secondary structure
- Solvent accessibility
- ...

# Two-Level Deep Convolutional Neural Networks



- **Training on 1426 proteins**
- **196 proteins validation dataset**
- **CASP10 dataset, CASP11 dataset, and CASP12 dataset**
- **Keras and TensorFlow**
- **Tesla K20 Nvidia GPUs**

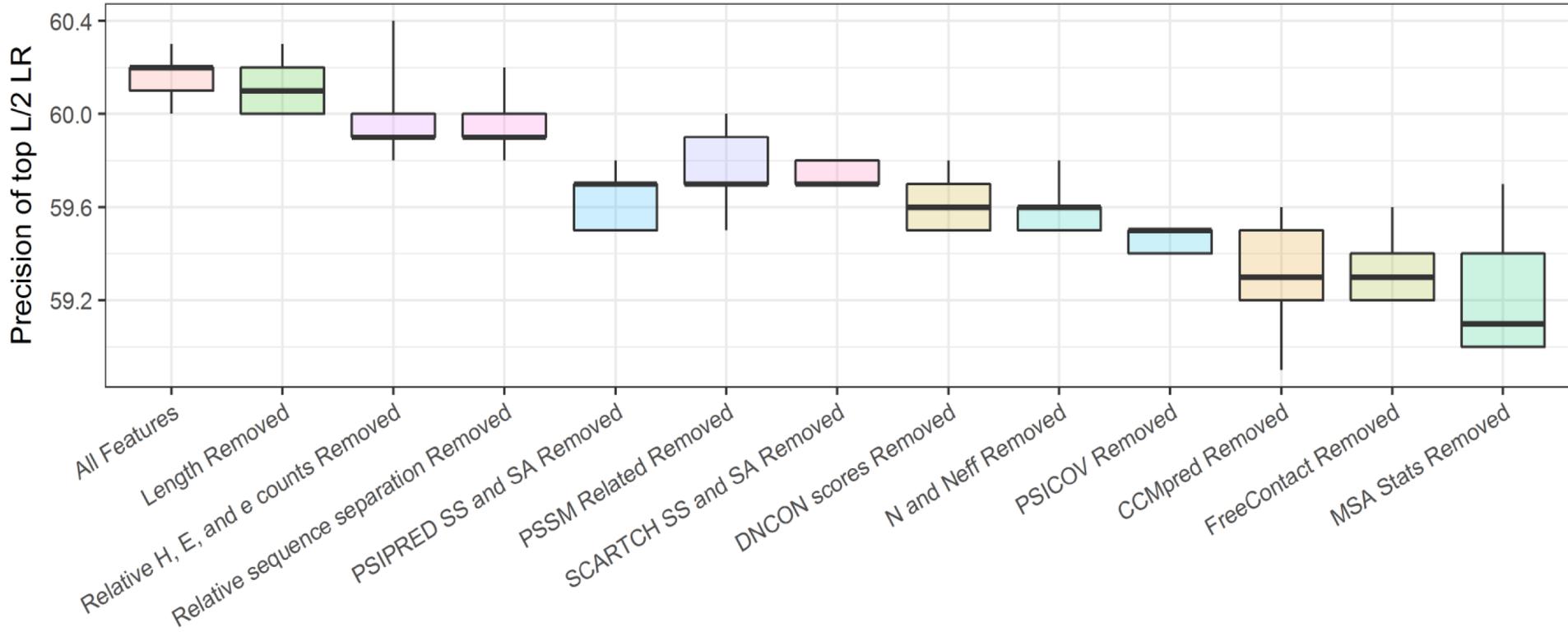
# Impact of Two-Level Network



# Testing on CASP Datasets

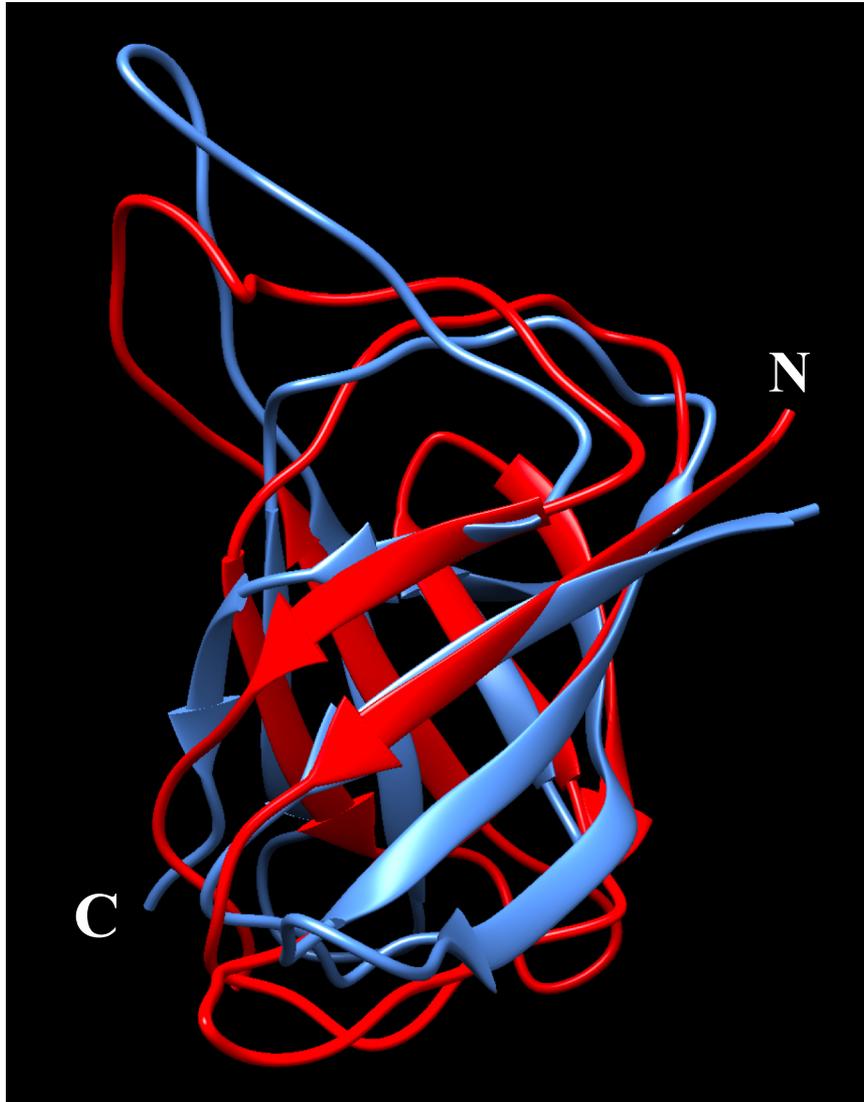
| FM<br>Dataset | Domain<br>Count | Precision of top L/5 long-range contacts (%) |            |        |
|---------------|-----------------|--|------------|--------|
|               |                 | Top CASP Group                               | MetaPSICOV | DNCON2 |
| CASP10        | 15              | 18.1 (DNCON 1.0)                             | 30.6       | 35.0   |
| CASP11        | 30              | 29.7 (CONSIP2)                               | 34.4       | 50.0   |
| CASP12        | 37              | 46.3 (Raptor-X)                              | 42.9       | 53.4   |

# Impact of Multiple Sources of Raw Information





# 3D Reconstruction from Predicted Contacts



RMSD of Core: 2.3 Å

TM-score: 0.67

Modeling Tool: CONFOLD

Blue: true structure

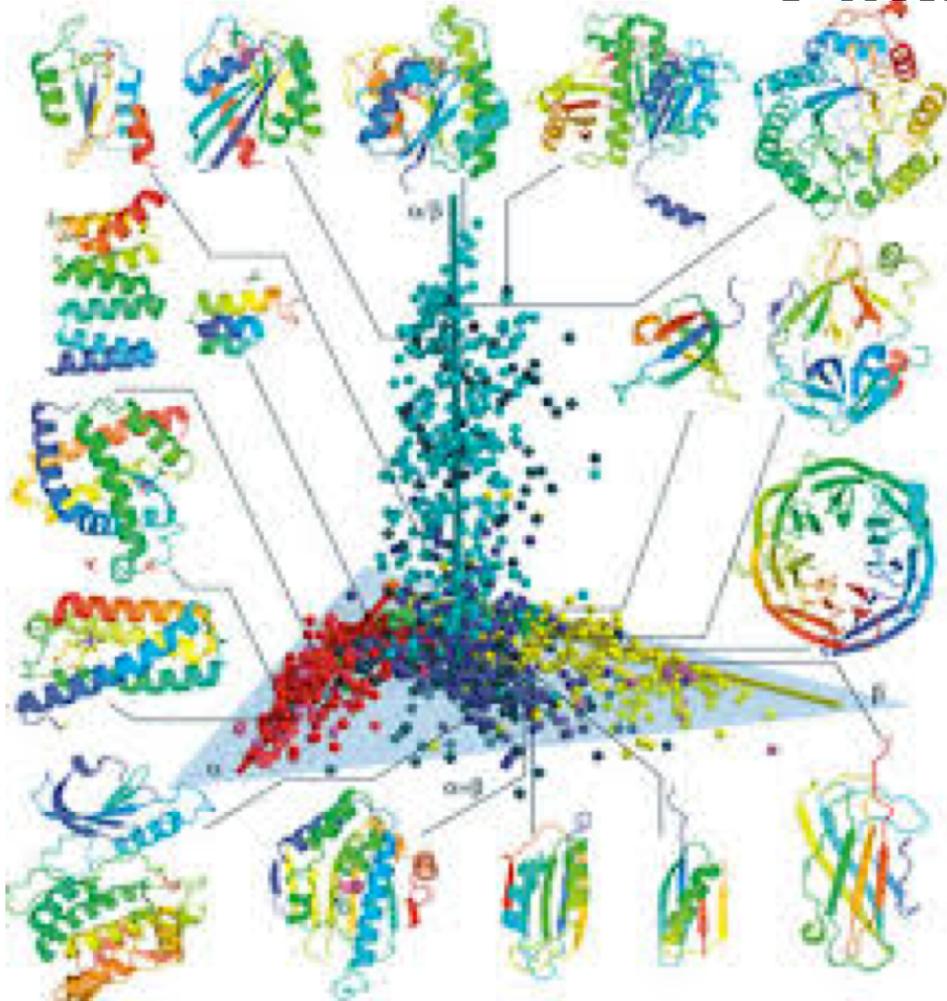
Red: predicted model

# Project 3: Deep Learning for Residue-Residue Contact Prediction

- DNCON2 source code:  
<https://github.com/multicom-toolbox/DNCON2>
- Data set:  
[http://sysbio.rnet.missouri.edu/multicom\\_toolbox/datasets.html](http://sysbio.rnet.missouri.edu/multicom_toolbox/datasets.html)
- Develop a simple deep convolutional network to predict protein contacts

# **1D Deep Convolutional Neural Network for Protein Fold Classification**

# Sequence-Structure Relationship (Millions of Sequences, but Only Thousands of Folds)



**SCOP database: ~1,200 folds**

**A central problem:**

- **Can we directly map protein sequence (language) into folds (semantics)?**

**Similar problems:**

- **Image Recognition**
- **Natural Language Processing**

# Limitation of Existing Methods

- **Sequence alignment (PSI-BLAST, HHSearch): nearest neighbor, indirect**
- **Ab initio folding: computationally intensive, inaccurate, hard to interpret**
- **Machine learning classification ( $f$ : sequence  $\rightarrow$  fold): only dozens of classes, fixed window, hand crafted features**

# Challenges

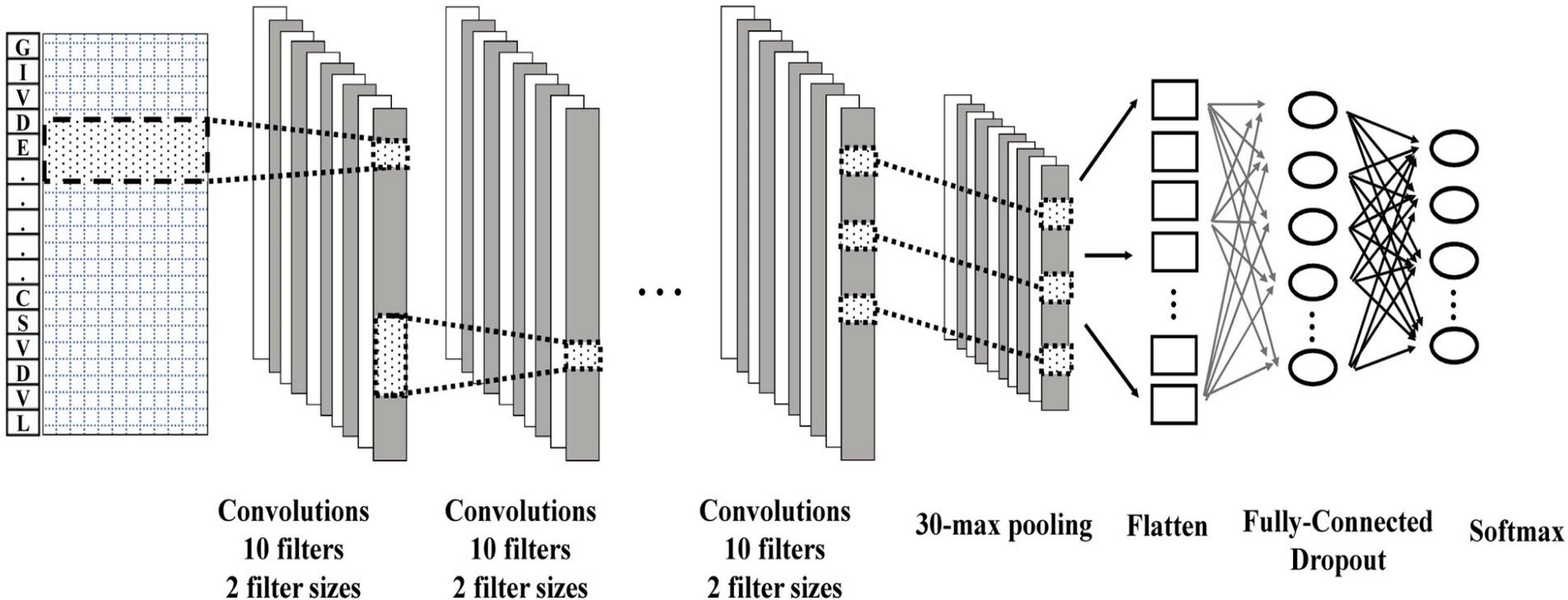
- **A very large number of classes (~1,200 folds)**
- **Natural sequence data without informative features**
- **Varied sequence lengths (from dozens to thousands)**

# **1D Deep Convolutional Neural Network for Fold Classification**

- Classify any protein sequence into thousands of protein folds directly (alignment-free, structure-free)**
- Automatically extract useful features from sequences of any length**
- Convert any sequence into hidden semantic features useful for protein comparison, clustering, and structure prediction**

# Deep 1D-Convolutional Neural Network

INPUT  $L \times 45$     Conv Layer1  $10 \times (L \times 2)$     Conv Layer2  $10 \times (L \times 2)$     ...    Conv Layer10  $10 \times (L \times 2)$     Pooling Layer  $10 \times (30 \times 2)$     Flatten Layer  $(1 \times 600)$     Dense layer  $(1 \times 500)$     Output  $(1 \times 1195)$

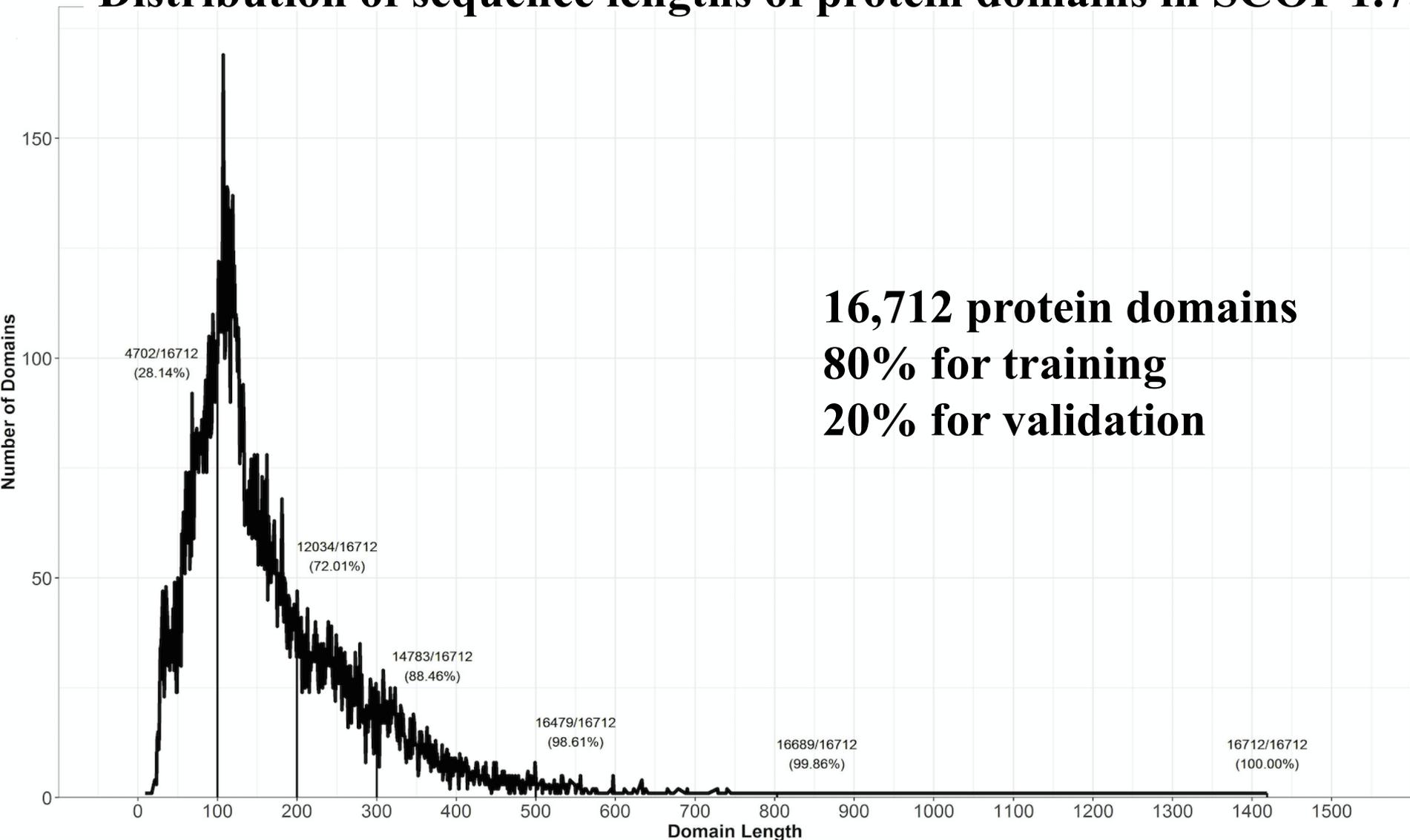


**Rectified Linear Unit (ReLU):**  $f(x) = \max(0, x)$

**Output layer: 1,195 nodes with sigmoid function**

# Training and Validation Data

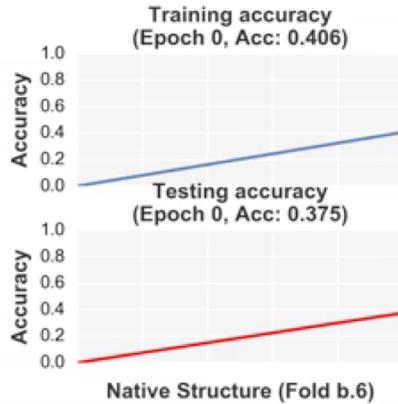
## Distribution of sequence lengths of protein domains in SCOP 1.75



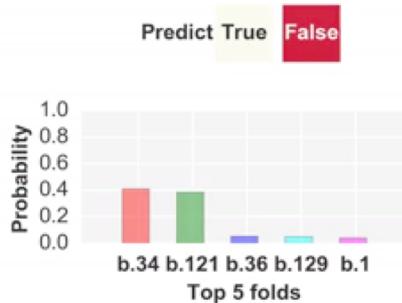
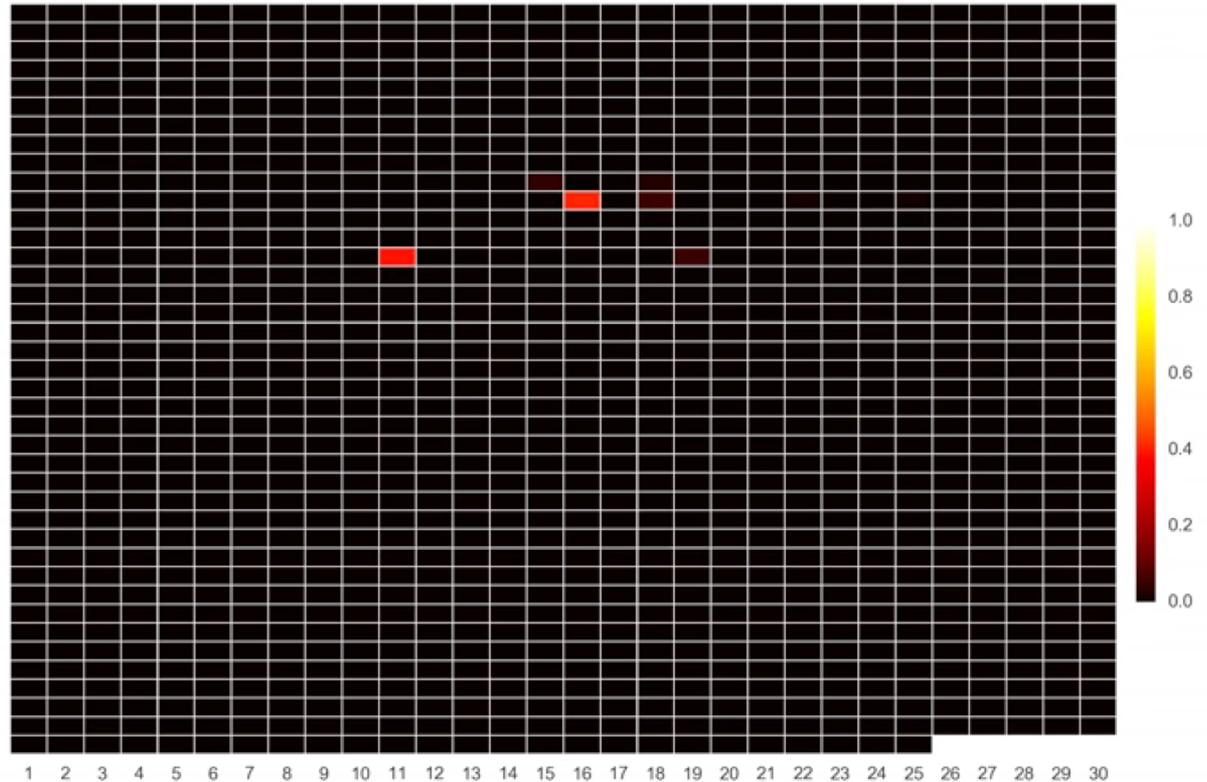
# Demo of Training DCNN



**Protein Sequence**  
ASCETTVTSGDGMTYSTRSISVSPASCAEFTVNFEHKGHMPKTMGMHNWVLAKSADVGDVA  
KEGAHAGADNNFVTPGDKRVI AFTPI IGGGEKTSVKFKVSALS KDEAYTYFCSYPGHFSM  
MRGTLKLEE



1195 Folds in SCOP 1.75



# Classification Accuracy on Validation Data

- Average accuracy of **top 1 prediction** on four validation datasets having <95%, 70%, 40%, and 25% identity with the training dataset is **75%**.
- Average accuracy of **top 5 prediction** is **91%**.

# Accuracy on Hard Template-based Targets of CASP9-12

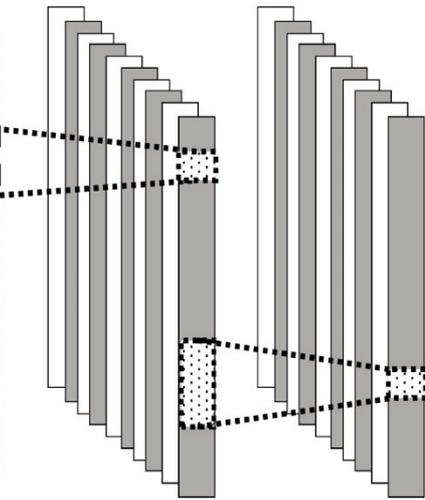
88 targets, <10% similarity with SCOP 1.75

| <b>Method</b>      | <b>Top1</b>   | <b>Top5</b>   | <b>Top10</b>  |
|--------------------|---------------|---------------|---------------|
| <b>DeepSF</b>      | <b>46.59%</b> | <b>73.86%</b> | <b>84.09%</b> |
| <b>HHSearch</b>    | 43.20%        | 61.40%        | 67.00%        |
| <b>DeepSF + HH</b> | <b>59.10%</b> | <b>77.30%</b> | <b>85.20%</b> |
| <b>PSI-BLAST</b>   | 15.90%        | 31.80%        | 47.70%        |

# Can the Hidden Features Represent Proteins?

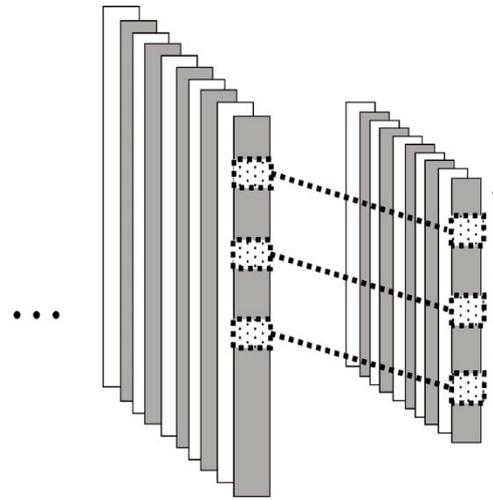
INPUT  $L \times 45$     Conv Layer1  $10 \times (L \times 2)$     Conv Layer2  $10 \times (L \times 2)$     ...    Conv Layer10  $10 \times (L \times 2)$     Pooling Layer  $10 \times (30 \times 2)$     Flatten Layer  $(1 \times 600)$     Dense layer  $(1 \times 500)$     Output  $(1 \times 1195)$

G  
I  
V  
D  
E  
.  
.  
.  
.  
C  
S  
V  
D  
V  
L



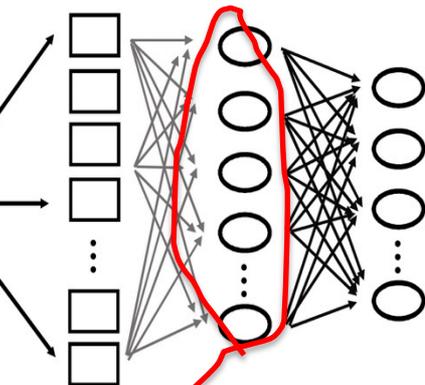
Convolutions  
10 filters  
2 filter sizes

Convolutions  
10 filters  
2 filter sizes



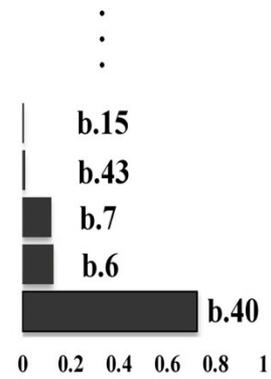
Convolutions  
10 filters  
2 filter sizes

30-max pooling



Flatten

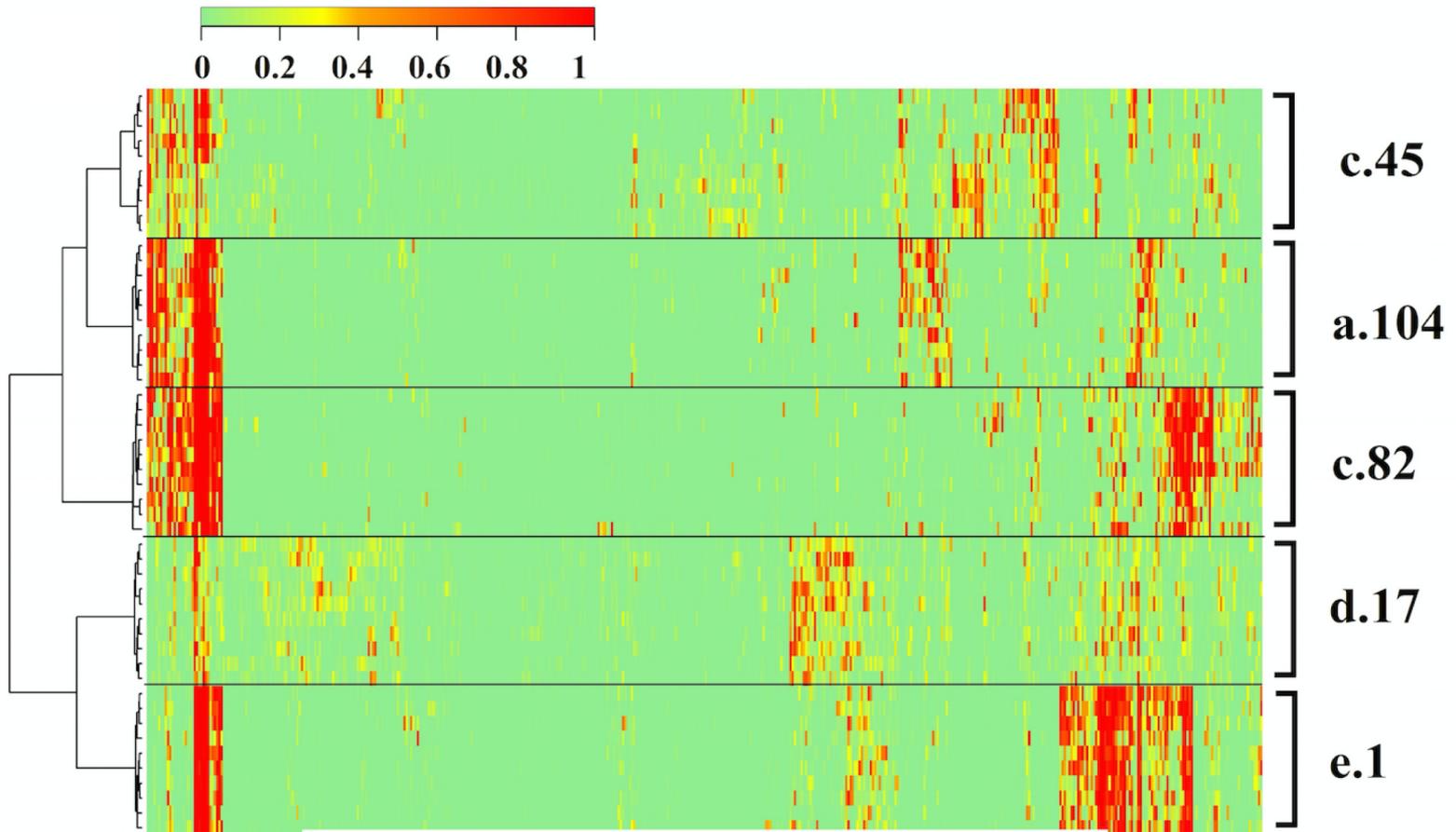
Fully-Connected  
Dropout



Softmax

**Hidden (Semantic) Features**

# Hidden Features Used for Protein Clustering

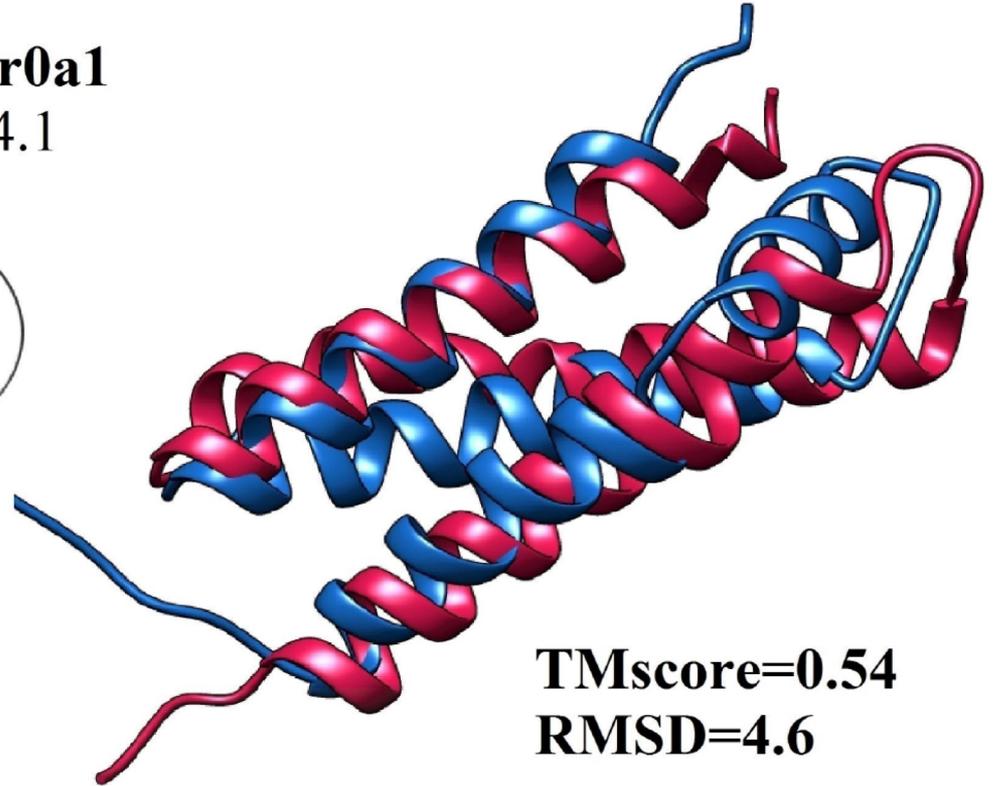
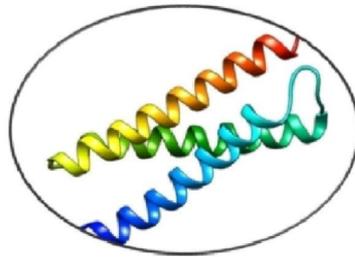


**Fold-specific features for domains  
in fold a.104, c.45, c.82, d.17, e.1**

# Fold Classification-Assisted Structure Prediction

Template: d1wr0a1

SCOP\_id: a.7.14.1



a.2: Long alpha-hairpin

a.4: DNA/RNA-binding 3-helical bundle

a.47: STAT-like

a.24: Four-helical up-and-down bundle

a.7: Spectrin repeat-like

0 0.2 0.4 0.6 0.8 1

TMscore=0.54

RMSD=4.6

T0862-D1 (DeepSF)

**Better than HHSearch model (TMscore = 0.3, RMSD = 8.22)**

# Project 4

- Develop a deep learning method to classify protein sequences into folds
- DeepSF source code and dataset:  
<https://github.com/multicom-toolbox/DeepSF>