

Course Summary

Protein Prediction I

Summer Term 2017

Based on the course 'Protein Prediction I for Computer Scientists (IN2322)' by the Chair of Bioinformatics, Technical University of Munich.

Disclaimer: This is an unofficial summary without any guarantee for correctness. It was created by students to improve their understanding of the subject and aid the learning process.
It should not in anyway serve as a replacement for visiting the lectures.
Prof. Rost really makes it worthwhile to sit in his lectures, so go there.

Table of Contents

Introduction	1.1
1. Lectures	1.2
1.1 Introduction: Bioinformatics	1.2.1
1.2 Introduction: Structure	1.2.2
1.3 Alignments 1	1.2.3
1.4 Alignments 2	1.2.4
1.5 Comparative Modeling	1.2.5
1.6 Secondary Structure Prediction	1.2.6
1.7 Secondary Structure Prediction 2	1.2.7
1.8 Secondary Structure Prediction 3	1.2.8
1.9 Membrane Structure Prediction	1.2.9
1.10 TMSEG	1.2.10
1.11 Beta Membrane and Accessibility	1.2.11
2. Exercises	1.3
2.1 Introduction	1.3.1
2.2 Biological Background	1.3.2
2.3 Protein Structures	1.3.3
2.4 Alignments	1.3.4
2.5 Resources for BioInformatics	1.3.5
2.6 Secondary Structure Prediction	1.3.6
2.7 Homology Modeling	1.3.7
2.8 Wrap Up	1.3.8
3. Exam Questions	1.4
3.1 Lecture Questions	1.4.1
3.2 Exercise Questions	1.4.2
3.3 Question Catalogue	1.4.3

Protein Prediction I

Course Summary, Summer Term 2017

tl;dr: This purpose of this document is to collaboratively create a both concise and detailed course summary of the *Protein Prediction I* Lecture from 2017 Summer Term at TUM.

To learn as effective as possible, I would like to encourage everyone to engage in the discussion evolving around the content of this document. If you have questions or challenges what someone else wrote please do so in a **constructive way**. We are all new to the subject of Protein Prediction and mistakes happen. Let's learn from them together!

Official Lecture Resources

Lecture Homepage: <https://www.rostlab.org/teaching/ss17/pp1cs>

Lecture Wiki: https://i12r-studfilesrv.informatik.tu-muenchen.de/sose17/pp4cs1/index.php/Main_Page

Youtube Channel: <https://www.youtube.com/channel/UCU6j8BG4RbEtTgyIZJ6Vpow>

Getting Started

This document is set up a **Gitbook** and hosted on **Github**. When you read this, you were already granted access to the repository so the first step is done.

The easiest way to start contributing is to download **Gitbook Editor** (available for Mac, Linux, Windows) from [here](#).

Before you add / change anything, please read through the Contribution Guide.

Contribution Guide

Tell others what you work on | Write meaningful commit messages | Push often | Use American English

Why is there a contribution guide? I think it is in everyone's best interest to keep this summary as easy to understand as possible for everyone. This guideline should help to maintain consistency across the entire document.

Each section may contain a short additional information on how to format things specific to that section. Please have a look there as well.

1. Adding new content

1.1 Adding minor updates

If you add minor updates, like the answer to a single question, you can do this on the `develop` branch directly. Make sure your commit has a meaningful message.

1.2 Adding major updates

If you add major updates, like several related changes (e.g. an entire lecture summary), go along as follows:

1. Add a new **issue** on Github, describing what you are working on
2. Create a `feature/<issue-name>` branch and add your changes
3. Open a pull-request to merge back into `develop` and add the other contributors as reviewers
4. Once the pull request is merged, delete your feature branch and close the issue by referencing the merge commit

Why so complicated? This way the issues reflect new changes and are transparent for all contributors.

2. Challenging existing content

If you find obvious mistakes (typos, clearly wrong statements) just change them directly.

If you are challenging statements, answers to questions etc. which might not be trivial to understand go along as follows:

1. Open a new **issue** on github.
2. Reference the statement in question you consider to be wrong
3. Provide an explanation why you think it is wrong
4. Provide your correct solution.

3. Adding new contributors

The purpose of this document is to foster collaborative learning - hence to make this as inclusive as possible. This being said, too many collaborators would probably lead to chaos. If you know other students personally, you want to add to the project shoot me a message and we will figure it out.

1. Lectures

1.1 Introduction: Bioinformatics

02.05.2017 | [Slides](#) | [Lecture Recording](#)

1. Definitions

Computational Biology: Biology Replacing experiments by computers (including neurobiology, image processing)

Bioinformatics: anything that has to do with storing and using the information about bio-sequences

2. Biology Introduction

Central to biology is the question: *How does life work?*

Question: What is common to life?

DNA, Protein, RNA

Question: How many bacteria do we carry around?

About 2 kilos. Humans carry around more bacterial DNA than human DNA.

Question: Which elements make up life?

- 65.0 % - O, Oxygen
- 18.6 % - C, Carbon
- 9.7 % - H, Hydrogen
- 3.2 % - N, Nitrogen
- 1.8 % - Ca, Calcium
- 1.0 % - P, Phosphorus

Question: What is life? Can you define it?

Descriptive definitions of life:

- Homeostasis (regulation of internal environment to maintain constant state)
- Organization (Unit: Cells)
- Metabolism
- Growth
- Adaptation
- Response to stimuli
- Reproduction

Question: Are viruses life?

Strictly speaking NO. Viruses on their own cannot replicate and thus are not alive. However, one could say that viruses are alive / represent life once they infected a cell and replicate.

2.1 Organisms

Different Type of Cells:

Prokaryotic Cells: Mainly found in bacteria and archaea.

- no nucleus
- usually unicellular
- no cell organells

Eukaryotic Cells: Found in animals and plants

- nucleus
- usually mulicellular
- cell organelles

Note: *The density within cells can be described as almost solid.*

Note: Different organisms use the same amino acids for proteins. However, they differ in their codon usage (which RNA triplets are translated into which amino acid).

Questions

Question: What is the smallest building block of life that can replicate?

cells

Question: How many different cells are in a typical human?

200

Question: What are the parts of the cell called?

organelles

Question: Which part of the cell is called the "powerhouse"?

mitochondria

Question: What part of a plant is involved with photosynthesis?

chloroplast

Question: What is mitosis?

cell division

Question: Who first used the term cell?

Robert Hooke

Question: How many elements are found in amounts larger than trace amounts (0.01%) in our bodies?

11

Question: When communities of living things interact with non living things they are called ... ?

ecosystem

Question: The most common molecule in the human body is ... ?

Water: H₂O

Question: What do bacteria have in common?

Single Cells

2.2 Genes

Question: What is DNA made out of?

DNA is a linear polymer out of 4 bases / nucleotides. DNA exists in cells mainly as a two-stranded structure called the double helix. Each of the bases has a complementary base.

- G: Guanine => Cytosine
- A: Adenine => Thymine
- T: Thymine => Adenine
- C: Cytosine => Guanine

Question: What is RNA made out of?

RNA is a single stranded linear polymer out of 4 bases / nucleotides.

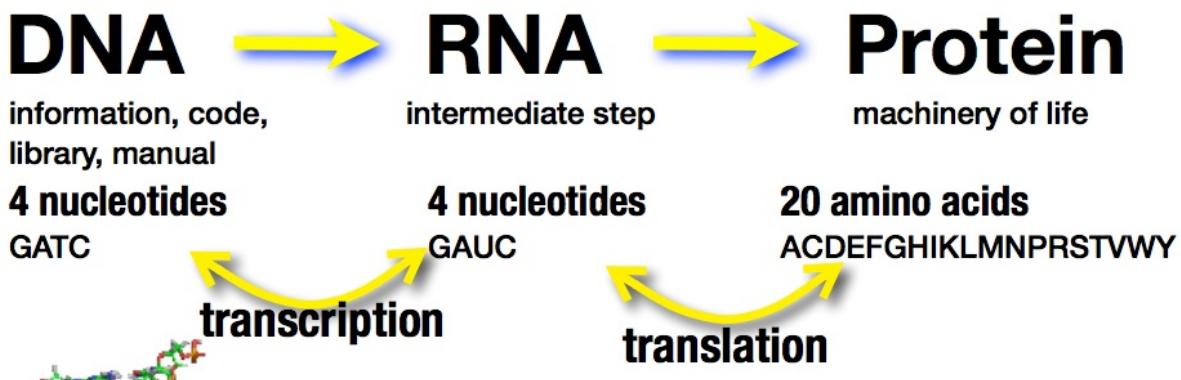
- G: Guanine
- A: Adenine
- U: Uracil
- C: Cytosine

Question: What is a gene?

A gene is a region of DNA, which contains all information for the creation of an entire RNA strand (= protein).

2.3 Central Dogma

Central dogma of molecular biology



DNA: Stores genetical information in a 4 letter alphabet. Double stranded helix.

RNA: Working copy of DNA, needed to produce a protein. (Oversimplification). Single stranded. Different types of RNA.

Protein: Composed of 20 letter alphabet (amino acids). Machinery of Life: Proteins do the work in our body.

Transcription: Process of turning a part of the DNA (a gene) into RNA.

Translation: Process of turning a RNA strand into a protein by a Ribosome. Each amino acid is encoded as a RNA nucleotides triplet.

In rare cases it is also possible that RNA translate to either RNA or DNA

Note: SEQUENCE leads to STRUCTURE leads to FUNCTION. Always!

3. Protein Introduction

Question: How many proteins does a typical human have?

Between 20.000 and 25.000 different kinds of proteins.

Question: What are functions of proteins?

- Defense (e.g. antibodies)
- Structure (e.g. collagen)
- Enzymes (metabolism, catabolism)
- Communication / Signaling (e.g. insulin)
- Ligand binding / Transport (e.g. hemoglobin)
- Storage (e.g. ferritin)

Question: How many residues long are typical proteins?

Between 35 and 30.000 residues. The median is around 400.

Question: Do proteins consist of units?

Proteins are built up of several domains. Most proteins have more than 2 domains.

Question: How many proteins are known?

About 85 millions sequences are known. However, the 3D structure (experimentally determined) of only 120.000 proteins is known.

Question: Is this gap (known sequences vs known 3D structure) expected to increase?

Yes, the gap is expected to increase. The amount of new sequences has increased drastically (far faster than Moore's Law) in the past. This is expected to continue. Advances in experimentally determining protein 3D structure could only improve marginally, but today experimentally determining the 3D structure of a proteins still costs about 100 000 EUR.

1.2 Introduction: Structure

04.05.2017 | [Slides](#) | [Lecture Recording](#)

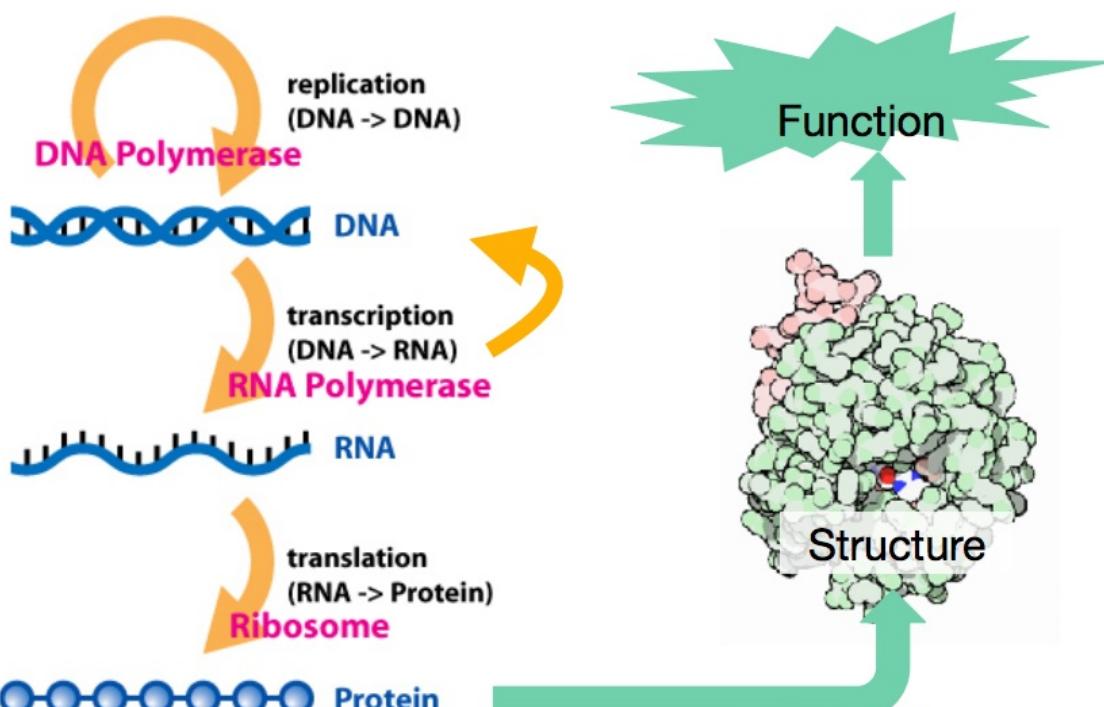
1. Recap

Common to Life: DNA / Cells

Proteins: Machinery of Life - they do everything that needs to be done

- about **85 Million** known protein sequences
- about 120 000 known 3D structures of proteins in PDB
- between 20.000 and 25.000 proteins in a typical human
- protein length (in amino acids): 35 - 30.000, with a median around 400

Central Dogma / Informationflow: DNA → RNA → Proteins



Translation: Proteins are made up of amino acids (20 different kinds). Each amino acid is encoded by a nucleotide triplet (codon) of DNA / RNA.

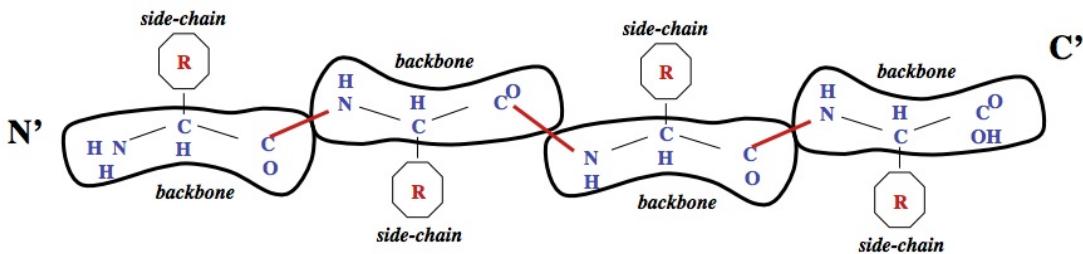
2. Proteins and Domains

```
# It is important to realize that every representation of a protein (sequence, image, ...)
# is
# only a representation of reality.
```

2.1 Amino Acids

Proteins are built up out of a chain of amino acids. These amino acids are joined into a **linear polypeptide chain**, a protein. Each protein is therefore a combination of the **20 different types of amino acids**.

polypeptide chain



- Each **residue** (amino acid) in this chain has a backbone and a side chain
- Different amino acids have **different side-chains**
- Each amino acids has **the same backbone**, along which they are chained
- Proteins are always chained up from the **N-Terminus** to the **C-Terminus** in a condensation reaction (a H₂O molecule is released)

Side Chains

Amino acids only differ in their side chains. These side chains determine the chemical properties of the respective amino acid. There are the following *features* an amino acid can have:

- polar (hydrophilic, likes water)
- non-polar (hydrophobic, avoids water)
- acidic (negatively charged)
- basic (positively charged)

Question: How many different amino acids are there?

20

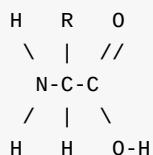
Question: How do amino acids differ? What do they have in common?

Different amino acids have different side-chains, which influence the chemical features of the respective amino acid. All of them share the same backbone.

Question: In which different feature groups can you categorize amino acids?

polar, non-polar, acidic (negatively charged), basic (positively charged)

Question: Draw the basic chemical structure of an amino acid.



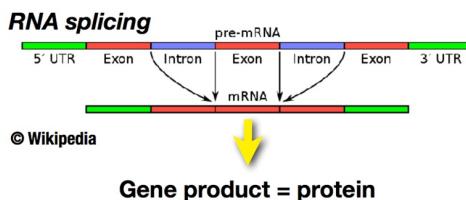
Question: How are amino acids linked together to form a protein?

In the translation process, a **Ribosome** translates a **mRNA** strand to a protein, by decoding the RNA triplets into amino acids and then linking the amino acids by peptide bonds. They chaining **ALWAYS** happens from the **N-Terminus** to the **C-Terminus** releasing an H₂O molecule as part of the reaction.

2.2 Protein Structure

What is a gene?

A gene is a region of DNA, which contains all information for the creation of an entire RNA strand. (= protein)



UTR: Untranslated region (leader sequence, header sequence)

Exon: Part of a gene that will encode a part of the final mature RNA (and thus protein)

Intron: Part of a gene that will be removed by **RNA Splicing** before the protein is translated

2.3 Domains

Definition: If I took a sequence out of a protein, string it up and put it into solvent, it adopts a unique 3D structure on its own.

Proteins are built out of several such substructures. The question is **Can we guess domains from sequence?** By aligning and comparing proteins with known 3D structure, it is possible to find common, overlapping domains that adopt the same 3D structure across different proteins.

Question: What is the definition of a 'domain'?

A domain is a protein sequence, which when put into solvent adopts a unique 3D structure on its own.

Question: How many domains does a protein have?

- 61% of proteins in the PDB are single domain
- 28% of proteins in the PDB are in 62 proteomes

Problem: This is a biased view on proteins. The 3D structure of Single-Domain-Proteins is easier to experimentally determine, so more Single-Domain-Proteins have been analyzed.

Question: Can domains overlap?

Yes, it can happen. However, it is not what is typically observed .

3. 3D Comparisons

There are many different methods for 3D alignments. The point is that comparing 3D structures is highly non trivial and ultimately comes back to the intuition about comparing 3D objects.

Question: How can we compare 3D structures?

One solution would be to align the corresponding residues of both sequences / 3D structures and take the **Root Mean Square Deviation**. (If one pair lies very far apart, it will result in an extremely low score)

$$RMSD(A, B) = \sum_{i=0}^n (r_i^a - r_i^b)^2$$

If the score is below a certain threshold, it is a match, otherwise it is not.

Question: How can align and compare the structure of 2 proteins?

- 1) Find the corresponding points (residues that match in 3D)
- 2) Find Superposition independent of domain movements and calculate score (e.g. RMSD)

Question: Why is global protein comparison most of the time impossible?

The definition of protein enforces a per residue comparison (no scaling). Hence only proteins of the (almost) the same length can be compared globally. Since proteins are between 35 and 30.000 residues long, global comparison does not make sense in most of the cases.

Question: What is the difference between global and local alignment?

In **global alignment** two structures / sequences are compared from beginning to end (compare the whole thing).

In **local alignment** however, subunits (domains) of the proteins are aligned. (Problem: What is a valid unit? Where to cut?)

Question: How to decide what is a valid unit for local comparison of 2 proteins?

(I couldn't identify a valid answer in the lecture recording)

Question: Which comparison not using cartesian RMSD could be used for comparison?

2D distance map: difference of differences. Only information about the chirality (mirror image) is lost.

1.3 Alignments 1

11.05.2017 | [Slides](#) | [Lecture Recording](#)

1. Recap

Sequence leads to Structure leads to Function

Question: Why compare 3D shapes, when we are after function? Why not compare function?

Because ...

- we cannot compare function directly
- structure is related to function
- we CAN compare 3D structures
- sometimes: similar structure -> similar function

Question: How do we get protein 3D shapes?

- primarily by experiment (most accurate)
- computational biology (most inferences)

Question: How much does it cost to experimentally determine the 3D shape of a protein?

Today it costs on average about 100 000 \$ per protein.

2. Tree of Life

- **All life is related (common ancestor)**
- 3 sections of tree of life
 - prokaryotes
 - (unicellular) bacteria
 - archaea
 - eukaryotes (plants, animals, ...)

Homology: Here (in the context of genes), it describes proteins originating from a common ancestor.

Definition of Species: We are talking about two different species, once they cannot produce fertile offspring together. (Example Bonobo and Chimpanzee)

Question: What are the 3 sections found in the tree of life?

bacteria, archaea, eukaryotes

Question: What does Homology stand for?

Here (in the context of genes), it describes proteins originating from a common ancestor. It is also frequently used to describe 'similar structure' for genes / proteins.

3. Pairwise Sequence Comparison

Correct alignment: We need an objective function

- simplest objective function: percentage of letters which are identical
- more complicated functions describing a match

BUT: the match score itself ignores, what we are after - biological similarity in function

Alignment

To find the optimal superposition of two sequence, it is first necessary to define what 'optimal' means.

Global Alignment:

- Align all residues from the beginning to the end
- Needleman-Wunsch

Local Alignment:

- Best match for locally aligned regions
- Smith-Waterman

How do we align 2 sequences?

Basically brute force: Visually (moving around), computationally (dynamic programming)

Dynamic Programming Algorithm: See [Exercise 2.4 Alignments](#)

Gap insertion penalty: Each wildcard (gap) used when aligning 2 sequences has a certain cost.

- Linear gap penalty: N gaps cost $N \cdot x$
- Affine gap penalty: opening gaps become more expensive
 - Gap open: cost $10x$
 - Gap extension (elongation): costs x

Local vs Global Alignment: What is better?

Question: What is better? High sequence identity of a short (local) sequence, or worse sequence identity when matching a longer sequence? How can we decide?

Compile the probability of randomly matching a sequence considering the background distribution. The result of this would be a substitution matrix such as BLOSUM62.

Question: Is identity the best way to match two sequences?

Not necessarily: What we really find is similar biological function. Some amino acids might have similar biophysical features and could be swapped without any significant influence on the structure of the protein. Such matches should also be considered 'positive'.

Building a scoring matrix based on evolutionary conserved residues does optimize the algorithm. (e.g. BLOSUM62)

Question: What is the biological assumption behind an insertion when comparing sequences?

Through evolutionary changes in the DNA (e.g. a point mutation) a new bump (= amino acid(s)) was introduced. Implicitly it is also assumed similar structure -> similar function.

Question: Why do linear gap penalties not model the reality of related genes / proteins well?

With a linear gap penalty (N gaps cost N*x) equally distributed gaps would be as expensive as clustered gaps. Biologically, gaps clustered to blocks, are however far more likely to occur, while the protein maintains similar structure / function.

It is more realistic to use an **Affine gap penalty** with higher costs for opening a new gap.

BLOSUM62

BLOSUM (Scoring Matrix)

BLOcks of amino acid SUbstitution Matrices

Align only conserved regions

compile log-odd ratios

$$S_{i,j} = \log \frac{p_i \cdot M_{i,j}}{p_i \cdot p_j} = \log \frac{M_{i,j}}{p_j} = \log \frac{\text{observed frequency}}{\text{expected frequency}}$$

BLOSUM n =threshold at $n\%$ pairwise sequence identity

Today many more substitution matrices exist.

Interactive Tool to practice dynamic programming: <http://melolab.org/sat>

Question: Does dynamic programming give the best solution?

Yes, dynamic programming produces one optimal solution. (There could be others, though)

Question: What are issues with dynamic programming?

- Time used: $O(n^2)$
 - Especially a problem, when comparing one protein against the entire database.
- How to choose parameters?
 - Gap penalties
 - substitution matrix

Question: How can we speed up the alignment of sequences?

1) Hashing (fast and dirty). e.g. BLAST

Question: How does BLAST (Basic Local Alignment Search Tool) work?

1. Start with indexed (hashed) seeds (words of size = 3) and find matching proteins
2. Extend matching 'words' into both directions
3. Begin dynamic programming from these strong local hits

4. Multiple Sequence Comparison

1.4 Alignments 2

16.05.2017 | [Slides](#) | [Lecture Recording](#)

1. Recap

3D Comparison: There is a way to look at proteins in 3 dimensions

- 1D - secondary structure prediction
- 2D - distance map
- 3D - real 3D coordinates

Dynamic Programming

- **Global** (Needleman-Wunsch) or **Local** (Smith-Waterman)
- With or without **Gaps**
 - Linear Gap Penalty: Opening and extending a gap have the same cost
 - Affine Gap Penalty: Opening a new gap is more expensive than extending an existing one
- **Scoring Matrices**
 - For each pair of residues you can read off, how much you gain by aligning these 2 residues

BLAST: Basic Local Alignment Search Tool

- Dynamic Programming is slow for large scale comparisons
- Speed search up by hashing words (seed = 3 amino acids / residues)
- After matching a word, try to extend the match by dynamic programming
- **Major Challenge:** Get the statistics right
 - *How significant is a match* against the background probability entire database

Question: What is the major challenge of BLAST?

Getting the statistics right: BLAST needs to know, *how significant a match is*, by comparing it against the background probability of the entire database.

2. Pairwise Alignment Accuracy

PSI: Percentage Sequence Identity

Zones:

- **Daylight-Zone:** PSI, where it can be assumed that from a similar sequence follows similar structure
- **Twilight-Zone:** PSI, where it is not possible to infer similar structure from similar sequence (the signal fades)
- **Midnight-Zone:** PSI, where sequence similarity does not tell anything about structure similarity

(down to random changes)

```
# Note: The Midnight-Zone is, where most proteins of similar structure sit
```

Going deeper into the Twilight-Zone, the following results are to be expected:

1. True Positives go up in absolute numbers
2. False positive increase (drastically) in absolute numbers

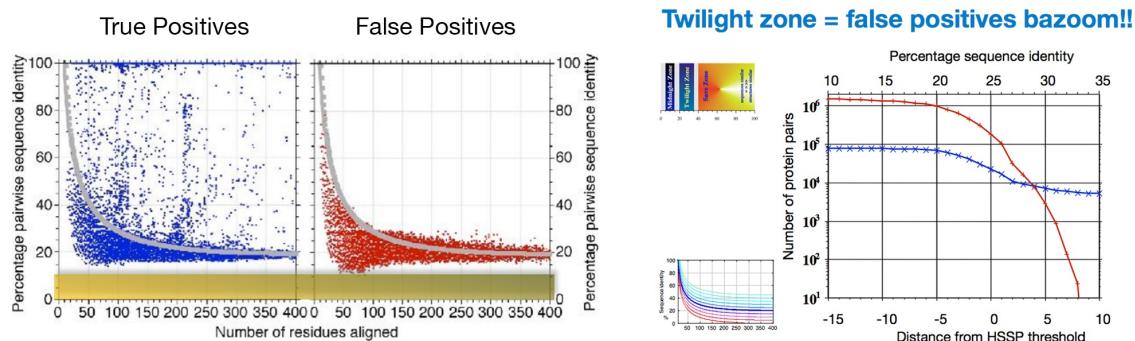
2.1 HSSP Curve

HSSP: Homology-derived Secondary Structure of Proteins

How to get the curve?

1. Get all 3D structures from PDB
2. Remove bias (sequence unique subset)
3. And compare 'all vs redundancy reduced set'
 - i. compare 3D structure (e.g. RSMD)
 - ii. compare sequence

Result: Sequence Conservation of protein structure



Observations

- Over the curve (daylight-zone), many true-positives (sequence identity \rightarrow similar 3D structure)
- Under the curve (upon entering the twilight zone)
 - Explosion of **false positives**
 - but also **significant increase** ($\times 10$) of **true positives** (This is why we want to go into the Twilight Zone!!)

Question: Why is it interesting to find similar proteins out of the Twilight / Midnight Zone?

The Midnight-Zone is, where most proteins of similar structure sit.

Question: Why is it that even with only 40% PSI, we can still assume similar structure? Could we randomly change 60% of the residues in the lab and get a new protein with similar structure?

- These 60% of changed residues happened under evolutionary pressure and are not

- random
 - mutations that did not change structure & function survived (we can observe them today)
 - mutations that did change structure & function most likely did not survive
- Thus randomly changing 60% of residues in a protein, would not result in a similar protein

Question: Why are certain proteins / structures multiple times in the PDB?

- different resolution of 3D structure
- different goals of publication produced (new) 3D structures
 - folding sites
 - binding partners
 - etc ...

3. Multiple Sequence Alignment

Dynamic Programming cannot be done efficiently with multiple sequences.

3.1 Multiple Alignment Hack 1: Iterative Pairwise Dynamic Programming

Progressive 1

1. Align sequences pairwise into consensus sequences
2. Align resulting consensus

Progressive 2

1. Align all sequences pairwise
2. Start with highest matching alignment
3. Iteratively align sequences into consensus sequence

How to find consensus? Different methods, e.g. the first, the more meaningful aa, ...

3.2 Multiple Alignment Hack 1: Map to Tree / Pairwise

ClustalW/ClustalX

- all against all (pairs) by dynamic programming (varying substitution matrices)
- build **phylogenetic tree**
- slow, dynamic programming, for experts

4. Profiles

Profiles profit from relation of 'families'.

Building up a profile, we can see certain amino acids that are more conserved than others. Computationally we can identify **sequence motifs** that describe such a profile as a regular expression.

PSSM (Position Specific Scoring Matrix)

You could also write a **profile** into a substitution matrix: A matrix of numbers with scores for each residue or nucleotide at each position.

Building a PSSM

1. Absolute Frequencies
2. Add pseudo-counts if necessary
3. relative frequency
4. log likelihoods

Question: How are profiles built up? How are the normal noted down? Do we have to know a specific algorithm?

Build up algorithm:

- Take all proteins of PSI over a certain threshold ...
-

Profile Formats:

- Regular Expression
- PSSM (Position Specific Scoring Matrix)

Question: What is a PSSM (Position Specific Scoring Matrix)?

A matrix of numbers with scores for each residue or nucleotide at each position. Built, e.g. by PSI-BLAST.

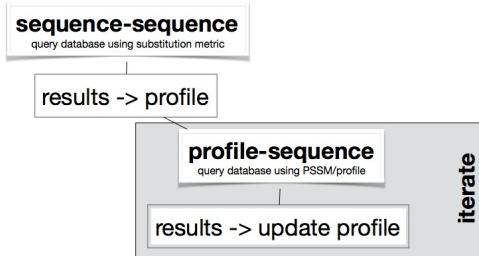
5. PSI-BLAST

Position-Specific Iterative Basic Local Alignment Tool

PSI-BLAST Steps

- 1) **Fast Hashing:** Like BLAST, match 'word'
- 2) **Dynamic Programming Extension between matches:** BLAST + Smith-Waterman
- 3) **Compile Statistics:** EVAL - Expectation Values
- 4) **Collect all pairs and build profile**
- 5) ... compare sequences (profile-sequence) and iterate

Steps involved for profile-based alignments



Question: Which steps are involved in building up a profile with PSI-BLAST?

- 1) **Fast Hashing:** Like BLAST, match 'word'
- 2) **Dynamic Programming Extension between matches:** BLAST + Smith-Waterman
- 3) **Compile Statistics:** EVAL - Expectation Values
- 4) **Collect all pairs and build profile**
- 5) ... compare sequences (profile-sequence) and iterate

Question: Why is PSI-BLAST so fast?

Because it drastically reduces the length of the comparisons with dynamic programming.

6. Hidden Markov Models (HMM)

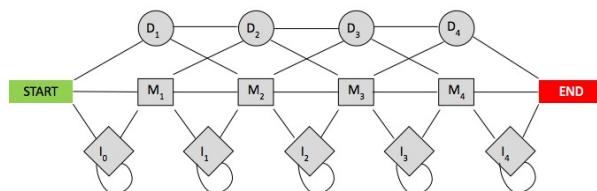
Hidden Markov Model are another method for creating Machine Learning Models. They are a good choice, if the structure of the problem is known beforehand.

- different states
- each state has transition probabilities to the neighboring states / itself

7. HMM for Alignment

Generic Profile HMM for alignment

- Captures matches, insertions, deletions
- Transition and emmission probabilities
- gap penalty handled by variation of transition probabilities
- calculation of probability by multiplying path variables



Entropy in alignment: Consider the residue at position i

- BEFORE any amino acid is aligned, we expect a particular amino acid to have some prior

background probability P_0 with entropy H_0

- AFTER the alignment we consider the same column with a *posterior probability* $P_i + priors \rightarrow H_i$.

$$\text{We expect } H_i = \begin{cases} 0, & \text{if conserved} \\ H_0, & \text{if varied} \end{cases}$$

- $H_i - H_0$ reflects the "bits_saved" by the alignment

With small families (few members, little divergence) the entropy is dominated by priors (= the background noise dominates)

8. Genetic Algorithm for alignment

Independence Assumption NOT needed for genetic algorithm

```
# All algorithms so far assumed *Independence between residues*:
# What happens at position i is independent of what happens at position i+x.

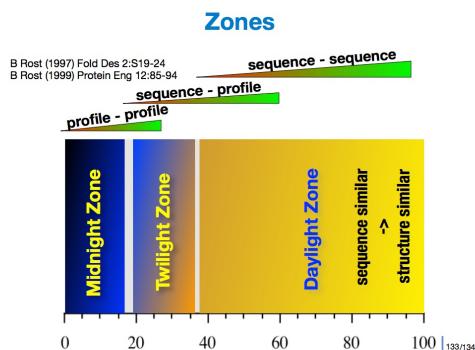
# The genetic algorithm does not make this independence assumption!!
```

- The genetic algorithm works on segments
- through mutations it creates new alignments

T-Coffee: much slower, requires pre-processing, Genetic algorithm

9. Profile-Profile Alignments

Compare Profiles to go even deeper into the **Twilight- / Midnight-Zone**



1.3 Comparative Modelling

18.05.2017 | [Slides](#) | [Lecture Recording](#)

- **1. Recap**

- Profile-Sequence comparisons are more accurate than sequence-sequence alignments
- Profile-Profile alignments gain even more accuracy

Question: How do you build up a family (profile) of sequences?

1. Find proteins of similar structure with BLAST
 2. Build PSSM
 3. build up a set of pairwise alignments
 4. add those over a certain HSSP value to the family
 5. Search with profile-sequence comparison for more distant family members and refine profile
- When building up a profile, start with a high threshold (only very similar sequences are taken), so the profile is not wrong from the beginning

2. Goal of structure prediction

Sequence uniquely determines structure! → Thus, from a sequence it should be possible to predict 3D structure and function

How would you assess prediction performance?

CASP: Critical Assessment of Structure Prediction

- Yearly event
- Submit predictions for structures, which will be experimentally predicted before a deadline
- Compare (after release of experimental structures) how the methods performed

Current State

- Only Homology Modeling is good
- No general prediction of 3D structure from sequence yet
- BUT: Important improvement in many fields

3. Structure by Experiment

Different Methods to determine 3D structure

- 90% - X-Ray Crystallography

- 09% - Nuclear Magnetic Resonance Spectroscopy (NMR)
- 01 % - Cryo Electron Microscope (Cryo-EM)

X-Ray Crystallography

1. **Grow Crystal:** Force the protein to grow a crystal
2. **Observe Diffraction Pattern:** Shoot x-rays onto crystal and observe the diffraction pattern
3. **Compute Electron Density Map**
4. **Fit observations to atomic model**

NMR

1. Protein has to be in similar solution as naturally
2. Massive Magnets required

Cryo-EM

- worse resolution than other methods
- cheaper than other methods
- *Pushing the boundaries of resolution of Cryo-EM is the future*

Question: Which methods to experimentally determine the structure of a protein exist? How much are they used?

Fraction of proteins in the PDB by experimental method:

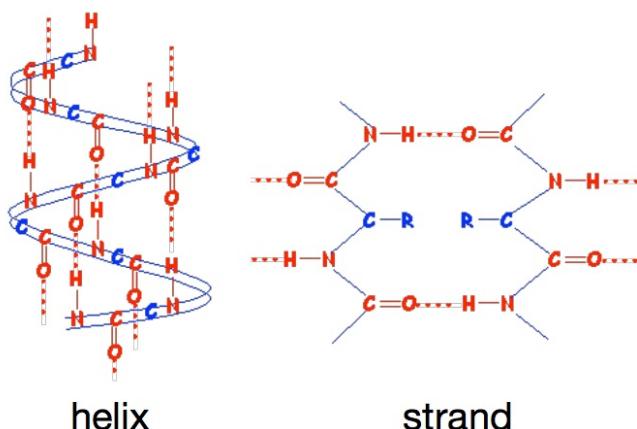
- 90% - X-Ray Crystallography
- 09% - Nuclear Magnetic Resonance Spectroscopy (NMR)
- 01 % - Electron Microscope (EM)

Question: How does X-Ray Crystallography Work

1. **Grow Crystal:** Force the protein to grow a crystal
2. **Observe Diffraction Pattern:** Shoot x-rays onto crystal and observe the diffraction pattern
3. **Compute Electron Density Map**
4. **Fit observations to atomic model**

Hydrogen Bond Formation

Idea: Secondary structure is completely explained by hydrogen bond formation.



Helix: Hydrogen-Bond between residue i and residue $i+4$, which stabilize the helix.

Sheet: Two strands come together to form a sheet by forming hydrogen bonds between them

Question: How to get 1D secondary structure from 3D coordinates?

Two methods were used to annotate 3D coordinates:

- 1) DEFINE, based on geometry (not used anymore)
 - 2) DSSP, based on hydrogen bond pattern (coulomb energy)

4. Comparative Modeling (=Homology Modeling)

Assumption: Sequence uniquely determines structure and therefore, from similar sequence follows similar structure.

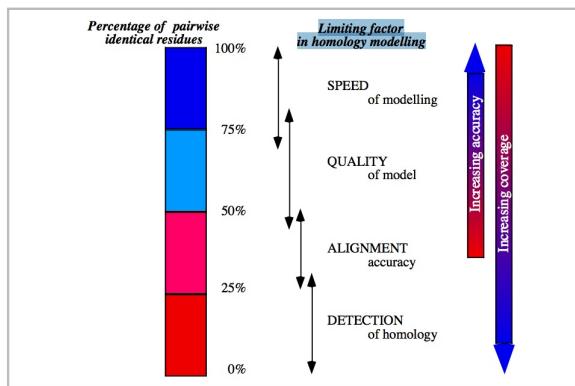
How can we use this to predict 3D structure?

Target: Protein to model

Template: Protein to model from

1. **Identify Template:** Query the PDB for similar sequences to your **Target**
 2. **Align Target / Template:** Select the best match as **Template** and assume the **Target** has the same structure
 3. **Build Model**
 4. **Assess Model**
 5. **Refine Model**

Comparative modeling: quality



Question: How does Homology Modeling (Comparative Modeling) work?

Target: Protein to model

Template: Protein to model from

1. **Identify Template:** Query the PDB for similar sequences to your **Target**
2. **Align Target / Template:** Select the best match as **Template** and assume the **Target** has the same structure
3. **Build Model**
4. **Assess Model**
5. **Refine Model**

Question: Which tradeoff does comparative modeling face? What are the limiting factors based on PSI (Percentage Sequence Identity)?

Tradeoff: Accuracy vs Coverage

Limiting factor in homology modeling:

75% - 100% - Speed of Modeling

50% - 75% - Quality of Model

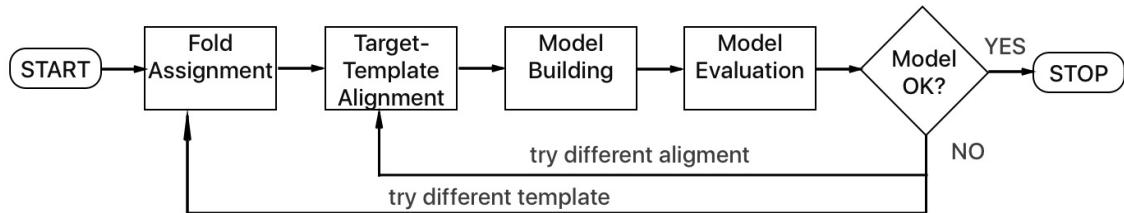
25% - 50% - Alignment Accuracy

0% - 25% - Detection of Homology

5. Comparative Modeling Methods

5.1 MODELLER

Summary: lots of whistles and bells, downloadable, very accurate



Constraint Satisfaction: use a set of objective functions to check whether the model is plausible

- $C_\alpha - C_\alpha$ distance
- Molecular dynamics
- Langevin dynamics
- Rigid bodies
- Rigid molecular dynamics
- ...

Optimization Steps (run repeatedly)

- explore different local minima

Typical Errors

- side chain packing
- misalignment
- wrong template

Pick the right solution:

- DOPE score (Discrete Optimized Protein Energy)
- based on knowledge based pair potentials

Question: How to handle a missing loop in comparative modeling?

- One way would be to find similar loops and compute the average over them.
- Another solution would be to apply molecular dynamics on the loop sequence. (only for short loops)

5.2 SWISS-Model

Summary: automated, increasingly comprehensive and flexible

Underlying 'Philosophy'

- fully automated
- for non-expert users / experimental biologists
- do less, make less mistakes

Original

1. alignment by BLAST / PSI-BLAST
2. copy to coordinates
3. end

Today: More complicated ...

1.6 Secondary Structure Prediction

?? .05.2017 | ??? | ???

1.7 Secondary Structure Prediction 2

01.06.2017 | [Slides](#) | [Lecture Recording](#)

1. Recap

Goal of Structure Prediction: Predict the 3D structure and function from an input sequence.

Proteins:

- are formed by stringing up amino acids in a chain
- amino acids are between 35 to 30.000 residues long
- amino acids form substructures, called domains in order to fold
- in principle: a domain put into solvent (water) folds on its own and adopts a unique 3D structure

Zones: There are different 'zones' by percentage of sequence identity, in which we can identify similar structures

- **Daylight-Zone** (100 % - 40%)
 - Sequence - Sequence Alignment
 - Assumption: sequence similar -> structure similar
- **Twilight-Zone** (40% - 20%)
 - Sequence - Profile Alignment
 - more distant relationships
- **Midnight-Zone** (20% - 0%)
 - Profile - Profile alignment
 - even more distant relationships

Global and Local Alignment:

Question: Relate the terms **Local** and **Global alignment** to the terms **Sequence-Sequence** and **Sequence-Profile**.

Global alignments refers to aligning sequences (proteins) from start to end. Local alignments refers to only aligning parts of the sequences (e.g. 50 residues).

Throughout Sequence-Sequence, Sequence-Profile and Profile-Profile methods both global and local alignment can be used. I practice mostly local alignment is done.

Comparative Modeling:

- Idea:
 - Find a similar sequence with known structure in the PDB (in the daylight-zone)
 - Try to use the known 3D structure of the similar protein to model the structure of the unknown protein
 - fix physical / chemical errors of the predicted 3D structure and find most plausible 3D

structure

- Reliably predicts **over 40 million proteins**
- However, for **most residues comparative modeling cannot be applied**

2. Secondary Structure Prediction

Secondary Structure Prediction happens in 1D, 2D and 3D. The following chapter will mainly be about 1D Secondary Structure Prediction.

DSSP (Define Secondary Structure of Proteins) algorithm: Has 8 states

- **H** = Helix
- **G** = 3_{10} Helix
- **I** = Pi Helix
- **E** = Extended
- **B** = Beta-bridge, single-strand residue
- **S** = bent
- **“”** = loop

Local Sequence determines secondary structure!

- Certain local sequence always form the same secondary structure (α -helix, β -strand, loop).
- Others (penta-peptides) are found in 2 different state, **dependent on their environment**

Question: What would be a simple method to predict secondary structure?

- 1) Take known structure
- 2) Find longest consecutive run of motifs that **ONLY** occur in one of the 3 states: H (Helix), E (Strand), O (Other)
- 3) Check unknown sequence against found motifs

Question: What was the first secondary structure prediction method?

Assuming that a **Proline** would break a helix, the occurrences of proline in a sequence was used to predict helices.

3. 1st Generation Secondary Structure Prediction

Idea: Build a frequency table over all amino acids, how often they occur in the secondary structure states based on the proteins where the structure is known.

Important: Bias reduction, to make the set table representative for future. Remove all proteins in comparative modeling range.

1. find a unique subset of proteins with known 3D structure (PDB)
2. convert 3D to 1D (secondary structure) with DSSP

Question: Where do we get the secondary structures from?

From the DSSP, which defines 8 states in total based on H-bond patterns.

Question: What is the 1st generation of secondary structure prediction based on? What was the accuracy? Was it successful?

- Based on single residues
- Between 50% and 55% accuracy (Q3)
- Clearly better than random - so it can be considered a success

How can we measure the performance of secondary structure prediction?

Q3: three-state per residue accuracy

$$Q3 = \frac{\text{number of correctly predicted residues in states helix, strand, other}}{\text{number of residues in protein}}$$

Question: How can the performance of secondary structure prediction be measured?

One way to do it, would be to calculate the **Q3** accuracy of a method against a test set. The Q3 accuracy is the **number of correctly categorized residues into one of the categories helix, strand, other** divided by the total amount of residues.

4. 2nd Generation Secondary Structure Prediction

Question: How did the second generation of secondary structure prediction improve? Name one algorithm.

Instead of using only single amino acids, it would consider a sliding window of the residues around a center amino acid.

Example: GORIII, with a Q3 accuracy of 55% - 60%

Question: What were problems of secondary structure prediction until 1994?

- the maximum accuracy of predictions was expected to be 65%
- β -sheet prediction was below 40%
- many predicted segments were too short to appear in nature

5. Introduction: Neural Networks

Goal: Use the representation of a set of examples (training set) for which the mapping *input* \rightarrow *output* is known to iteratively refine the weights of the connections between output and input units so that the error is minimized.

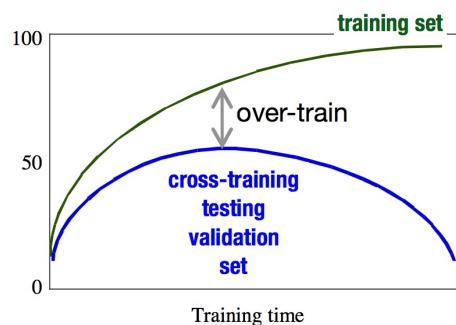
Principles of neural networks

- **Free Variables:** *Connections* $\{J\}$
- **Output:** $out_i = \sum_{j=1}^{N^{in+1}} J_{ij} in_j$
 - in_j value of input unit j
 - out_i value of output unit i

- J_{ij} connection between input unit j and output unit i
- Error: $E = \sum_{i=1}^{N^{out}} (out_i - des_i)^2$
 - out_i value of output unit i
 - des_i secondary structure state observed for central amino acid for output unit j

Training: Change of connections $\{J\}$ such that E decreases (e.g. gradient descent)

Problem: Overtraining - happens if the network becomes too specific to the actual training set and loses accuracy for predicting unknown input. The point when to stop training can be found by using **cross-training, testing, validation sets**.



Cross-Validation: Split your available dataset into 3 sections

1. **Training** (50%): used to train ML algorithm
2. **Cross-Train** (25%): used to find threshold when to stop training and tweak parameters
3. **Testing** (25%): used ONLY to assess performance / accuracy of final ML algorithm

Question: How can the introduction of a new hidden layer in a neural network be described by means of a simple graph?

Each new hidden layer basically introduces a new 'decision line' which can separate datapoints into different categories.

Question: What is cross-validation in the context of Machine Learning and why do we need it?

Cross Validation is a method for estimating the performance of a predictive model (e.g. a neural network). To use it, the available dataset is split in 3 categories, 1) a training set, 2) a cross-training set and 3) a test set.

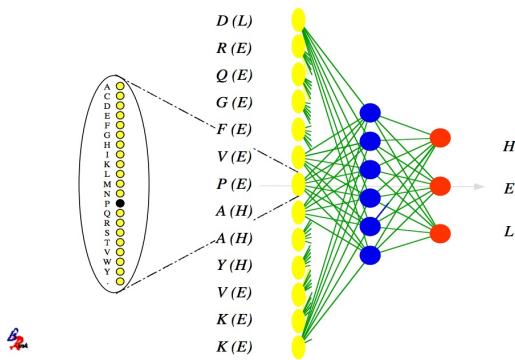
- 1) The training set is used to train the model
- 2) The cross-training set is used to estimate the performance of the model after x training steps
- 3) The test set is used to assess the final performance of the model after training is finished

The cross-training set is needed to decide, when to stop training (when overtraining sets in) and to tweak certain parameters before running against the test-set.

6. Neural Networks for Secondary Structure

Goal: Solve the 3 problems at the time [1] accuracy, [2] strand performance, [3] short segments

Neural Network for secondary structure



Input: $13 * 21$ input units

- 13 ??
- 20 amino acids + 1 spacer

However, the final accuracy was only about **62%**

Balanced Training:

- Helices are overrepresented in the training data
- Choose the training data, so all 3 states (helix, strand, other) are equally represented

Result

- overall accuracy dropped to **60%**
- β -sheet prediction improved from **40%** to around **60%**

Question: Did balanced training improve the Q3 prediction accuracy? Which assumption did it prove wrong?

Balanced training actually decreased the Q3 accuracy. However, it did improve the prediction accuracy for strands significantly, falsifying the hypothesis that strands could not be predicted with local information.

7. PHDSec: Structure to Structure Prediction

We still have the problem of bad segment prediction (too short segments). This is due to the fact that samples from the database are selected at random, losing information about local correlations.

How can we get information about the local correlation (e.g. length of a helix) into the prediction model?

Solution: Add a second Neural Network, which takes the predicted sequences from the first network as input.

BUT: Accuracy was still only $60\% + \varepsilon$

Question: Which problem did PHDSec solve? How did it accomplish it?

By introducing a **Structure-to-Structure** prediction model, PHDSec improved the prediction of too short segment. The Structure-to-Structure network would take structure (helix, strand, other) prediction of a sequence as input and predict segments based on them.

1.8 Secondary Structure Prediction 3

08.06.2017 | [Slides](#) | [Lecture Recording](#)

1. Recap

Goal of Structure Prediction: Predict the 3D structure and function from an input sequence.

Cell: The density of a cell is like solid state, but proteins are still surrounded by water

Relation of Proteins: We can find out about the relation of proteins by comparing their sequence

- direct sequence-sequence comparison only in the daylight-zone
- building up profiles means to 'pick up the implicit evolutionary signal'

Question: Which ways of comparing proteins are there? Why do we need

- Dynamic Programming (Brute Force)
- Hashing (e.g. BLAST)

Question: Why are fast search algorithms such as BLAST needed?

Comparing sequences of length n residues is in $O(n^2)$. For comparing a single pair this is still fine, but comparing one (newly found) protein against all known proteins in the PDB (about 120 000) is impossible. Thus we need 'shortcuts' such as BLAST to speed up the search.

Question: What is the normal approach when you find / analyse a newly found protein?

- 1) Sequence the new protein (if not done yet)
- 2) Run BLAST against the PDB
- 3) Run Dynamic Programming against the results from BLAST

Question: In terms of CPU, is sequence-sequence as fast as sequence-profile?

Question: How can it be that even with only 40% sequence identity we assume / observe similar structure?

The changes in sequences we observe are not random, but follow underlying evolutionary rules. Changes, which affected the structure and thus the function of a protein are simply not likely to survive and thus we do not observe them. Changes, which did not influence the structure / function however, did survive.

Question: Why is protein sequence changing? Why are we mutating?

- Replication Errors (point mutations)

- Radiation
- Viruses

Question: How much do any two unrelated typical humans differ on average?

On average every pair of humans would differ in one amino acid per protein. (Though, changes cluster)

Question: In a structure to structure network, which additional information could be used to improve the prediction?

- E.g. redundant information about the sequence, e.g. parts of it.
- length of protein
- Is the sliding window at the end / mid / start of the protein?

Question: When training a neural network, how do you choose the next training sample from your test set? Why so?

Randomly, to avoid correlations

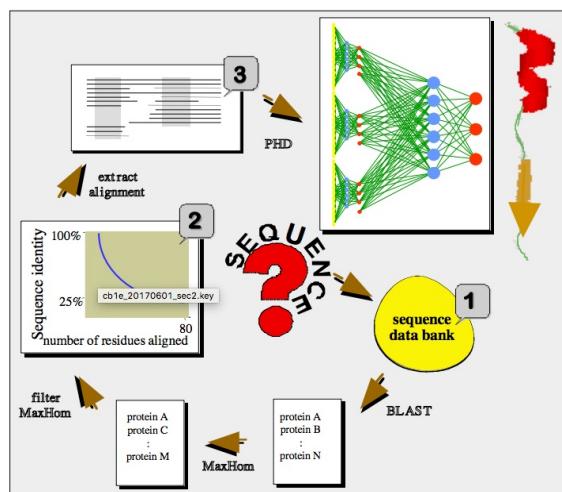
2. 3rd Generation Secondary Structure Prediction

Critical Question: How to improve beyond $60\% + \epsilon$ accuracy?

Evolution improves prediction: An **evolutionary profile** averaged built up over several species implicitly captures the history of an individual protein.

PHD: Neural Network and Evolutionary Information

- Build up the family (profile) for the protein and add it to the input of the network
- Each amino acid in the input now has a probability on how often it occurs in the family



Additional Input for PHDSec network

- family (profile)
- percentage of each amino acid in protein
- length of protein

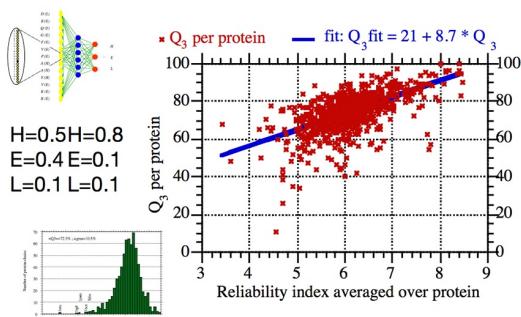
- distance: center, N-term
- distance: center, C-term

Jury decision improve accuracy: All of these input features are fed to different Networks, resulting in many independent predictors (**Jury**). All of these networks add their own 'white noise' to the prediction. The average over all the predictors is better than every single one.

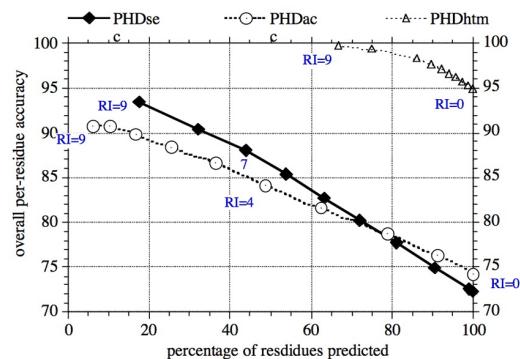
The final accuracy (on average) of ProfSec is about **72%**.

Prediction of correctly predicted residues: In addition, **ProfSec can give an estimation on the strength of the prediction** for each protein. (By counting the 'stars')

Stronger predictions more accurate!



Correct prediction of correctly predicted residues



Global Information improves ProfSec's per protein prediction.

	Q3 (per residue)	Q4 (per protein)
Only Sliding window (local)	72%	70%
Local & Global	72%	75%

Question: How does ProfSec overcome the 60% accuracy hurdle in secondary structure prediction?

ProfSec uses evolutionary information - the family of the protein - as additional input. Furthermore, other relevant input data (e.g length of protein, distribution of amino acids, ...) are used to build up several different networks that independently predict the secondary structure. Together this jury of networks achieves a more accurate prediction than they would on their own.

Question: How would you build up a family for a protein?

1. Search the PDB for proteins in comparative modeling range. (Assumption: same sequence, same 3D structure, same secondary structure)
2. Use profile to search in twilight-zone for potential proteins of that family (possibly verify whether the found protein is plausible to have similar 3D structure) and add to family (recompute profile)

Question: How do you get from a sequence to a secondary structure prediction with PHD?

1. Use BLAST to find potentially similar proteins in sequence data bank

2. For the resulting proteins calculate the sequence identity (homology) with dynamic programming
3. Filter all proteins, which are below a threshold of sequence identity (only take those "over the curve")
4. Extract the profile by aligning the remaining proteins
5. Predict the secondary structure with the sequence and its family as input

Question: Which accuracy does ProfSec achieve on average? What are additional advantages of other secondary structure prediction methods?

ProfSec achieves a Q3 accuracy of about 72% on average. Additionally it can also predict the strength of the prediction.

Question: Does adding global information improve ProfSec prediction?

Yes it does. While the Q3 accuracy (per residue) is not improved, the Q4 accuracy (per protein) does improve.

3. Proper comparison of methods

For a meaningful comparison the methods should

- use the same (meaningful) measure (e.g. Q3)
- use the same dataset
- split training / testing
 - there must not be an overlap between sets
- is the difference (in accuracy) significant (= difference > standard error)
- was the test set not used for making decisions?

1.9 Membrane Structure Prediction

13.06.2017 | [Slides](#) | [Lecture Recording](#)

1. Introduction Membrane

Requirements for Cell Membrane

- separate the content of the cell from its surroundings
- control traffic into and out of the cell
 - - keep malicious things out
 - let good things in
- must be a dynamic structure

Main components of the cell membrane

1. Carbohydrates
2. Cholesterol
3. Phospholipids
4. Proteins

Phospholipids

- form the barrier that separates the inside of a cell from the outside
- phospholipids are arranged in a **lipid bilayer**
 - **Inside:** fatty acid tails (non-polar, hydrophobic)
 - **Outside:** phosphate group (polar, hydrophilic)

Membrane Proteins

- Provide several functions to the cell
 - help to be recognized by immune cells
 - transport proteins control substance flow in and out of the cell
 - receptor proteins bind hormones, which can change cell function
 - provide structural stability
- Membrane proteins can (easily) shift around laterally

Membrane Proteins are especially important for drug targets

Note

Similar to the membrane, proteins also tend to have a **hydrophobic core**.

Membrane proteins however, tend to have a **hydrophobic outside** and an **hydrophilic core**.

Question: What are the requirements of a cell membrane?

- separate the content of the cell from its surroundings
- control traffic into and out of the cell
 - keep malicious things out
 - let good things in
- must be a dynamic structure

Question: What are the 4 main structural components of the cell membrane?

Carbohydrates, Cholesterol, Phospholipids, Proteins

Question: What is the cell membrane mainly made out of?

The cell membrane is a so called **lipid bilayer** of **phospholipids**. Phospholipids have a non-polar, hydrophobic tail (membrane center) and a polar, hydrophilic head (outside of membrane).

Question: What are functions of membrane proteins?

- help to be recognized by immune cells
- transport proteins control substance flow in and out of the cell
- receptor proteins bind hormones, which can change cell function
- provide structural stability

Question: Can membrane proteins easily move around?

It depends:

- Membrane proteins can easily move laterally
- But it is hard to move into / out of the lipid bilayer

2. Introduction Transmembrane Helix (TMH)

Although membrane proteins are especially interesting for drug targets, there are only limited 3D structures to be found in the PDB. Mostly, because it is extremely difficult to put them into a crystal in their 'natural' membrane environment.

There are essentially 2 important questions when it comes to TMH prediction:

- How many helices go through the membrane?
- In which direction do they go through the membrane? (topology)

Question: Why are there so few membrane proteins in the PDB?

It is particularly difficult to experimentally determine the structure of membrane proteins due to the special environment they naturally occur (the membrane).

Question: What are the key questions TMH prediction tries to answer?

- How many helices go through the membrane?

- In which direction do they go through the membrane? (topology)

3. TMH Prediction

Question: Why could be a plausible reason why PHDSec failed for predicting transmembrane helices?

Unlike 'normal' proteins, transmembrane proteins have an hydrophobic outside and a hydrophilic inside.

What could be a strategy to come around that? **Build up a hydrophobicity index.**

- There are different hydrophobicity scales, optimized for different problems

Identifying hydrophobic regions

- Whenever the hydrophobicity is over a certain threshold, consider it a membrane helices
- except the hydrophobic residues over (the lower) threshold are not long enough for a TMH (20 residues)

Identifying Topology: What is inside/outside of a TMH?

- **Positive Inside Rule:** Looking at the parts which connect TMHs within on protein, they look different depending on which side of the membrane they are: **There is a excess of positively charged residues on the inside.**

Question: How should we choose the threshold for the hydrophobic regions?

1. Predict the hydrophobicity for the protein
2. Assign a positive inside-out
3. choose the threshold to **optimize the inside out difference**

Question: What is the Positive Inside Rule and what is it used for?

The positive Inside Rule is used to find the topology of transmembrane proteins: The loops connecting TMHs on the inside of the cell membrane have an **excess of positively charged residues.**

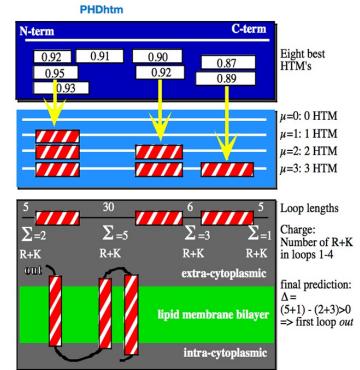
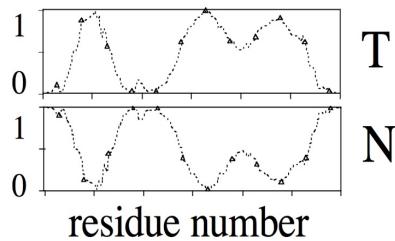
PHDhtm: Membrane Helix Prediction

Predict: **Membrane Helix or Not Membrane Helix**

Using the Sequence-to-Structure -> Structure-Structure approach now led to helices twice the length of observed TMH. Based on the number of TMH to expect based on the hydrophobicity a dynamic programming approach can be applied on the *NN energy*.

```
# Note: I didn't really get how this helps in reducing the problem of too long helices??
```

Dynamic programming on NN 'energy'



TMHMM: Membrane Helix Prediction

- Based on a Hidden Markov Model

4. When are TMHs correctly predicted?

Today: at maximum ± 5 residues overlap

1.10 TMSEG

20.06.2017 | [Slides](#) | [Lecture Recording](#)

1. Introduction TMSEG

1.1 Rational: Why another predictor

- More Data available
- Less expensive machine learning (more computing power available)
- Improve runtime

1.2 Dataset

- 166 membrane protein sequences after redundancy reduction
- Data curated and linked from several databases (PDB, OPM, ...)
- 1441 proteins from the SignalP Training Set
 - 1142 soluble (after RR)
 - 199 membrane (after RR)
- **Split Dataset into 4 subsets**
 - each set maintaining distribution of TMPs, SPs and sequence length
 - use 3 sets for **cross-validation**
 - use 1 set for final independent evaluation (**blind set**)

2. TMSEG Prediction

Intro: Classification Trees and Random Forest

Classification trees

- Given N training samples and M input features find the best recursive partitioning to predict the class labels in the leaf nodes
- Approaches differentiate algorithm: splitting, pruning, balancing, ...

Random Forest

How does it work?

- ensemble method: grow T trees for a forest
- for M input features choose $m < M$
- for each $t \in T$
 - Select N training samples with replacement from all N samples
 - At every split, choose m random features. Use the best split among those to build the tree

- The final prediction uses the prediction of all trees

Advantages

- fast
- robust against overtraining
- no black box
- intuitive to interpret
- good performance

2.1 Step 1 - Feature Sets

Initial Prediction

- Random Forest ($T = 100, m = 9$)
- Sliding Window of 19 residues ($w = 19$)
- 3 scores for each residue (0 - 1000)
 - Signal Peptide (often mistaken for TMHs)
 - TMH
 - Soluble

Feature Set

- **Global Features**
 - Global amino acid composition
 - Protein length
- **Local Features**
 - PSSM Score (Position Specific Scoring Matrix)
 - Distance to N- / C- Terminus
 - Average hydrophobicity (Kyte-Doolittle)
 - percentage of hydrophobic residues (in window size $w = 9$)
 - percentage of negative / positive charged residues (in window size $w = 9$)
 - percentage of polar residues (in window size $w = 9$)

2.2 Step 2 - Empirical Filter

- smooth score with median filter ($x = y$)
- Adjust scores to avoid overprediction
 - soluble ≈ -185
 - TMH ≈ -60
- Assign each residue the state with the highest score
- Remove signal peptides with <4 residues
- Remove TMHs with <7 residues

2.3 Step 3 - Refine TMH prediction

- **Neural Network** (25 hidden nodes)

- Input: TMH segments of variable segments of variable length
- Features:
 - Amino acid composition
 - Average hydrophobicity
 - percentage of hydrophobic residues
 - percentage of charged residues
 - segment length
- Split long TMHs (≥ 35 residues) into 2 shorter TMHs (≥ 17 residues)
- Adjust TMH endpoints by up to ± 3 residues

2.4 Step 4 - Topology Prediction

- Random Forest ($T = 100, m = 7$)
- Assign soluble segments to side 1 or 2
- Features
 - Amino acid composition
 - percentage of positive charged residues
 - percentage of absolute difference of positive charged residues on side 1 vs side 2
- Only consider residues close to TMHs
 - 15 residues nest to TMHs and 8 residues into TMHs
- Predict topology of N-Terminus and extrapolate
- if a SP is predicted, the residues after the SP are always 'outside' (SP = Signal Peptide)

Question: What are advantages of using a Random Forest?

- Fast
- robust against overtraining
- no black box
- Intuitive to interpret
- good performance

3. TMSEG Performance measures

```
# Note: Per residue measures are often misleading!
#       => better score TMH segments
```

Whole Protein Scores: Q_{ok} and Q_{top} => What is a correctly predicted TMH?

- **Strict Criteria**
 - Endpoint deviation ≤ 5 residues
 - Overlap at (observed / predicted) at least 50%

How can we measure the performance on predicting Transmembrane Helices?

Recall:

$$r_i = \frac{\text{correctly predicted TMHs}}{\text{observed TMHs}}$$

Precision:

$$p_i = \frac{\text{correctly predicted TMHs}}{\text{predicted TMHs}}$$

Q_{ok} :

$$Q_{ok} = \frac{100}{N} \sum_{i=1}^N x_i; x_i = \begin{cases} 1, & \text{if } p_i = r_i = 100\% \\ 0, & \text{else} \end{cases}$$

t_i :

$$t_i = \begin{cases} 100\%, & \text{if topology correct} \\ 0, & \text{else} \end{cases}$$

Q_{top} :

$$Q_{top} = \frac{100}{N} \sum_{i=1}^N y_i; y_i = \begin{cases} 1, & \text{if } t_i = p_i = r_i = 100\% \\ 0, & \text{else} \end{cases}$$

How can we measure the performance on distinguishing soluble proteins from transmembrane proteins?

FPR:

$$FPR = 100 * \frac{\text{incorrectly predicted TMPs}}{\text{soluble proteins}}$$

Sensitivity:

$$\text{Sensitivity} = 100 * \frac{\text{correctly predicted TMPs}}{\text{observed TMPs}}$$

```
# Result: TMSEG has exceptionally low misclassification rates compared to other methods.
#           Additionally, it is strong on topology predictions.
```

4. Future Work

How to get more data?

Check against data published after the release of the method. The data is then unknown by any method.

4.1 Applying TMSEG to other methods

- High modularity (step 1 - 4) of TMSEG allows it to be applied to other methods
- Apparently it can

4.2 Potential extensions

- Re-entrant regions not modeled (too little data)

1.11 Beta Membrane and Accessibility

22.06.2017 | [Slides](#) | [Lecture Recording](#)

Recap

Lipid bilayer (membranes)

- hydrophilic outside,
- hydrophobic inside

Normal surroundings of proteins are **solvent** (hydrophilic, water). Generally, the core of a protein is **hydrophobic**.

Trans Membrane Helices (TMH)

- really small fraction of experimentally known proteins (3D structure)
- but 15% to 25% of all proteins
- 60% of drug targets
- only about 2% of all *unique* structures have membrane helices
- **1D prediction very successful**

Beta Barrels

TMB = Trans Membrane Barrel

- "barrels" formed out of β -sheets connected by hydrogen bonds, which go through the membrane
- looking from the tops they have a hole
- they are pores, letting anything pass that is small enough

Beta Barrel Prediction: PROFtmb

Model Design:

- Hidden Markov Model
- structure based labels (states)
 - inside loop
 - outside loop
 - strand up
 - strand down

How to assess whether this model makes sense?

- Count the different states in the set of proteins, where you know (from experiments where the barrels are)

- Put the observation into the **priors** for the **HMM** (Hidden Markov Model) and train for all the others
- Check the results (per residue) predicted vs observed

Conclusion: Remarkable performance

BUT: Can we distinguish proteins with / without TMB?

Challenges:

- Where do barrel domains start / end?
- Sometimes barrels are built out of several peptide chains (proteins)
- Per Protein Performance: **Accuracy vs Coverage**
 - Where to put the threshold when analysing a new protein?
 - Intuitive / Literature: **Intersection of Accuracy and Coverage**
 - Optimized per Case: E.g. for Master Thesis high accuracy if more important than coverage, as experimental biologists will follow up on only a few of the found proteins in further research

Accessibility

What is it about? Why is this relevant?

- accessibility of residues to water
- outside vs inside

1) Absolute Accessibility: ASA (square Ångstrøm, 1 Å = 0.1 nm)

Long side chains may appear more accessible: *Different amino acids have a different length of their side chain and thus the absolute accessibility per amino acid differs.*

Using absolute accessibility may lead to wrong conclusions.

2) Relative Accessibility: ASA / max ASA

3) "States":

- buried, exposed
- buried, intermediate, exposed

Note: It doesn't matter whether something is 80% or 100% exposed, but it does matter whether something is 0% or 20% exposed. Also, drawing the line where to set the "best" threshold between the states is discussed in academia.

RostLab Approach: Square Root -> Switch from percentage to predicting 10 states

Solvent Accessibility

Accessibility helps in predicting protein function.

- sub cellular localization
- protein-protein interactions
- flexibility / motion from structure

Historically: Prediction by hydrophobicity

- hydrophobic: inside
- hydrophilic: outside

PHDacc: Machine Learning Approach

- 10 output units
- Advantage: No need to decide on threshold beforehand. Threshold can be chosen for future needs.
- Advantage: Mapped to a 2 state system (buried / exposed) each prediction also carries the confidence in the prediction

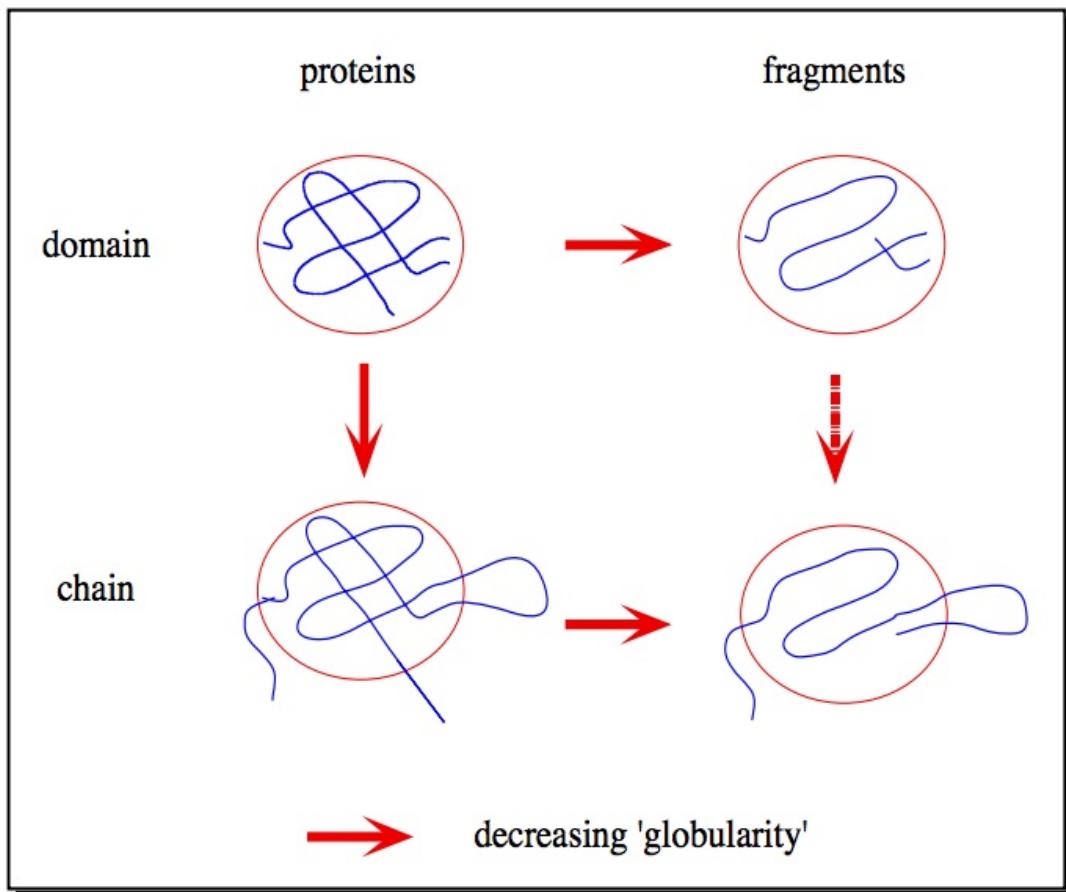
Detailed Prediction Problematic:

ConSurf: Significant gain by evolutionary information (in/out with > 75% accuracy)

More Globular - More likely expressed

Note: I really don't get this slide / content. Anyone an idea, what is meant by that?

- **Domains** are compact structures on their own (= they fold on their own)
- **Question:** How can we see (by a sequence) what we are related to? (Related to what?)
 - Answer: Predict the residues on the surface. (Why???)
 - 1) Take a 2 state model (buried / exposed)
 - 2) Predict the residues which are exposed
 - 3) Check to which of these (see image) the sequence fits best
 - Assumption: Proteins are spheres. (Which is apparently the case in an overwhelming fraction of proteins)



2. Exercises

2.1 Introduction

11.05.2017 | [Slides](#) | [Wiki](#)

2.2 Biological Background

11.05.2017 | [Slides](#) | [Wiki](#)

2.3 Protein Structure

18.05.2017 | [Slides](#) | [Wiki](#)

2.4 Alignments

01.06.2017 | [Slides](#) | [Wiki](#)

2.5 Resources for Bioinformatics

08.06.2017 | [Slides](#) | [Wiki](#)

2.6 Secondary Structure Prediction

22.06.2017 | [Slides](#) | [Wiki](#)

2.7 Homology Modeling

29.05.2017 | [Slides](#) | [Wiki](#)

2.8 Wrap Up

06.07.2017 | Slides | Wiki

3. Exam Questions

This section contains possible exam questions compiled from different sources.

Contribution Guide: Exam Questions

Since all questions here are answered by students, there might be some mistakes in them. Hence a few more words on how to best handle this section.

1. Adding a new question

Just add the question in the respective file. Optimally, you can already provide an answer.

2. Answering a question

To clearly distinguish questions from answers, please put answers in **blockquotes** right under the respective question.

Example:

- How can 1D secondary structure information be used to derive a 3D model?
- It is not possible to derive a 3D model from 1D information. (Trick Question)

3. Updating an answer

If you think an answer does not properly answer a question (e.g. it is wrong or the answer is not sufficient), mark the answer and open a new **issue** on Github to discuss the question and share your improved answer.

(Use the `???` emoji to mark the possibly wrong answer inline)

Example:

- How can 1D secondary structure information be used to derive a 3D model?
- ?? It is not possible to derive a 3D model from 1D information. (Trick Question)

3.1 Lecture Questions

This section contains possible exam questions asked Professor Rost in the lectures he dedicated to answering student questions. They are **highly relevant**, because he will sample exam questions from this pool.

Questions (Thursday, 22nd June)

Question: How can you choose the **e-value** for PSI-BLAST depending on the size of the dataset?

E-value indicates significance of alignment/ hits returned by chance when searching through DB. It depends on the size of dataset and length of query. So higher e-values from large DB aren't always bad (and opposite: smaller e-values from small sample space isn't always good).

Question: You want to develop a new method to predict e-values, how do you prepare your data?

?

Maybe to consider length of sequences, data base size, redundancy of data

Question: What is the regular process when you want to analyse a new sequence?

Go to DB and search your sequence to find out whether homologs of this protein are already available, and if they are, what is known about them.

If we didn't retrieve any significant hits, try to search for motifs, patterns.

Question: What is a structural domain? What is a functional domain and How can we deal with the fact that they can be in different places?

- structural domain: part of sequence with unique 3D structure
- functional domain: part of sequence with unique function
- use local alignment

Questions (Tuesday, 27th June)

Question: How do we predict proteins?

?

- use DB to search for homologs
- use Machine learning and other methods for secondary structure prediction
- use Comparative modelling for 3D structure prediction

Question: Why would someone give you a sequence?

Amount of newly discovered sequences constantly increases. And even for known sequences we still can refine information that we have about them.

Question: How do you run a sequence against the DB?

- Blast search (uses indexing technique, scoring matrix and dynamic programming to find short similar segments)
- Psi-Blast search (at first it uses Blast search, creates profile (PSSM) from highest scoring hits and uses it to replace substitution matrix in a subsequent Blast search, this process can be repeated many times to refine profile)

Question: How do you build a family?

Use non-redundant database. Set a threshold for E-value. Perform similarity search of all the proteins against each other using Blast. Proceed results i.e. set thresholds for including proteins in family (similarity, sequence length). For each family align proteins using ClustalW.

Question: Pairwise/multiline alignment: what can we achieve, what is the risk?

Multiple alignment can be used to find pattern characteristics of specific protein families, build phylogenetic trees, detect homology between new sequence and existing families, help in predicting secondary structure of new sequence.

Sequence alignment can help in: finding conserved regions between the 2 sequences, similarity searches in DB.

In case of Multiple alignment there is a risk of pollution, any errors in the initial alignments cannot be corrected later as new information from other sequences is added.

Question: What is the difference between pairwise and multiple alignment?

Pairwise alignment compares 2 sequences, multiple - 3 and more. (Also look previous question)

Question: Can we predict something we have not observed?

When we predict features we can try to find something that people "want" to see.

Question: Sliding windows introduce information from the sequence environment. Why do we need convolutional network on top of that? Why do we need anything else on top? (*not sure about correctness of question*)

It gives us ability to detect motifs whenever it's in sequence window. It exploits spatially local correlation.

Question: How do we prepare data to predict B-value? (*not sure about correctness of question*)

?

Training set: non-redundant set of high resolution protein structures. Network is trained on properties that can be obtained from primary a.a. sequence: secondary structures and solvent accessibility. We can also use evolutionary profile and global information about sequence.

(There is an answer in the last video lecture, which is not uploaded yet)

Questions (Thursday, 29th June)

Question: With a matching *profile-profile* comparison what can you say about the two families?

The assumption is that matching profiles share a similar / same structure and function.

Question: When I build a profile of a family: Do they share the same structure? Should I verify that they do? How do I do that?

The very assumption is that the proteins of one family share the same structure and function. When iteratively refining the profile with proteins retrieved by *profile-sequence* of *_profile-profile _comparison* (from the twilight- / midnight-zone), it can make sense to double check the new proteins with secondary structure prediction to avoid adding false-positives.

Question: Cross-Validation: What is it? How does it work? Why do we need it?

Is an validation technique for assessing how the results of a statistical analysis will generalize to an independent data set. We need it to estimate how accurately a predictive model will work in practice

N-fold cross-validation: Partition data in 'n' folds. Use 'n-1' for training and 1 for performance assessment. Repeat with different hold out partitions. Average performance.

Question: What is the difference between a BLOSUM matrix and a PSSM (Position Specific Substitution Matrix)?

In case of PSSM a.a. substitution scores are given separately for each position in a protein multiple sequence alignment. And in BLOSUM we have substitution scores for all possible substitution pairs (210 in total) of 20 standard a.a.

Question: What is the most successful method to predict 3D structure?

comparative modelling

Question: What is homology modeling (= comparative modeling) and how does it work? What are the limitations of it?

Align sequence with proteins in PDB and set a threshold. Introduce gaps as loops.

Limitation - no similarities found (templates are unavailable or fragmentary), matching right residues (errors in sequence alignment produce errors in the homology model)

Question: How can you predict structure in the [a] daylight- [b] twilight- [c] midnight-zone?

While using sequence-sequence, sequence-profile and profile-profile alignment respectively.

Question: What is the assumption behind all alignment methods that is incorrect and nevertheless seems to work? Give a method that aligns 2 proteins without that assumption.

The assumption is that the alignment of the residue at position i is independent of the $i+x$. (In short: alignment i and $i+1$ are independent).

The only method that does not rely on this assumption is the **Genetic Algorithm** (e.g. T-

Coffee)

Question: Why do we have so few experimentally confirmed structures in the PDB?

ca. 85 million proteins are known, but only about 120 000 are in the PDB.

Sequencing technology has improved and become a lot cheaper and hence many more proteins were discovered. While there were improvements in 3D structure determination, it still costs at least 100 000 euros to experimentally determine the 3D structure of a protein. It is expected that this divide will further increase.

Question: Why do we need redundancy reduction for machine learning?

The training data for the ML model should be representative for the problem for which it should predict.

Half of known structures in DB have 100% sequence similarity

Question: Say the 3D structure for **N** thousand proteins were known and they serve as input for a method predicting 1D structure. How can you define the value for _sequence-unique _that you have to apply to create an unbiased data set? Why do you need an unbiased dataset?

Using RMSD and putting a threshold

We need an unbiased dataset in order to assess performance of 3D to 1D prediction.

Question: How do you compare proteins of different length?

By using local alignment

Question: What is the significance in using information from protein families (also inferred as evolutionary information) as input to the ML device predicting 3D structure?

- It is additional information for the ML device, which is clearly relevant for the structure
- the profile is a record with information about the 3D reality of the protein

Question: How can I use 1D information to get a 3D structure? What can you do with a 1D structure?

It is impossible to reconstruct a full 3D structure from 1D information. 1D structure can be used for

- optimizing a profile
- predict whether a protein is soluble
- predict whether a protein is a transmembrane protein
- input for further secondary structure prediction
-

Question: What is a 2D contact map (distance map)? How can it be obtained?

Shows distance between all possible a.a. pairs.

By using 3D structure and distance functions (i.e. Voronoi contacts)

Question: Explain the concept between the notation of 1D, 2D, 3D structure. What is in the PDB? What does the DSSP give?

?

1D - secondary structure

2D - contact map

3D - tertiary structure, 3D shape of protein

In PDB there is only 3D structure

DSSP assigns secondary structure according to hydrogen-bond pattern

3.2 Exercise Questions

This section contains possible exam questions asked and answered as part of the exercises.

3.3 Question Catalogue

This section contains possible exam questions sourced from students of previous Protein Prediction I lectures and the lecture recordings.

3.3.1 Exam Structure: 2016ST

We were able to obtain last years exam structure. Let's try to answer all of the concrete questions :-)

Part 1 is mandatory, for the rest choose 3 out of 4.

1. Multiple Choice (5 questions, 10 points)

- Secondary Structure
- RMSD - Protein Similarity
- Hydrogen Bonds (α -helix, β -sheets, long/short bonds)
- 100% sequence identity => same structure? (PIDE)
- Can modern prediction methods correctly predict structure in the midnight zone?
- About "Cryo-Microscope"
- About "X-Rays"

2. Sequence Alignment (10 points)

- Explain each of the following alignment techniques and provide one method for each
 - Sequence - Sequence
 - Sequence - Profile
 - Profile - Profile
- General scoring BLOSUM62 matrix vs. PSSM
- Why is the sequence information valuable?
- How BLAST speed up pairwise alignments?
- Global vs Local alignment

3. Sequence Structure (10 points)

- What data is needed to predict the structure with ML?
- Which tools and db you will use?
- How to prepare data for ML
- Which 2-3 features will help to predict?
- Would you apply method to all protein (query)?
- Which measure would you use to evaluate your method?

4. Protein Structure (10 points)

- Why it is important to know 3D structure?

- Why is it so hard to compare 3D structure?
- Most successful ML algorithm for predicting structure, steps
- Method for experimental structure determination. Short explanation. How many structures are experimentally known?

5. Machine Learning (10 points)

- General definition of Machine Learning
- Cross validation
- What is 'feature'?
- ETP explain, example
- Name and describe one ML method
- Name and describe "sequence" in context of PP
- Discuss how to predict Protein Structure from amino-acid sequence using ML
- Q2: which is better, how to prove your's is better, which value you will publish? (What is Q2)

3.3.2 Lecture 1: Introduction Bioinformatics

Question: What is common to life?

DNA, Protein, RNA

Question: How many bacteria do we carry around?

About 2 kilos. Humans carry around more bacterial DNA than human DNA.

Question: Which elements make up life?

- 65.0 % - O, Oxygen
- 18.6 % - C, Carbon
- 9.7 % - H, Hydrogen
- 3.2 % - N, Nitrogen
- 1.8 % - Ca, Calcium
- 1.0 % - P, Phosphorus

Question: What is life? Can you define it?

There is no holistic definition of life: Descriptive definitions of life are

- Homeostasis (regulation of internal environment to maintain constant state)
- Organization (Unit: Cells)
- Metabolism
- Growth
- Adaptation
- Response to stimuli
- Reproduction

Question: Are viruses life?

Strictly speaking NO. Viruses on their own cannot replicate and thus are not alive. However, one could say that viruses are alive / represent life once they infected a cell and replicate.

Question: What do bacteria have in common?

Single Cells

Question: What are the differences between prokaryotic and eukaryotic cells?

Prokaryotic Cells: mainly found in bacteria and archaea, usually unicellular, no nucleus, no cell organelles

Eukaryotic Cells: Found in animals and plants, usually multicellular, have nucleus, have cell organelles

Question: How can the density of a cell be described?

The state inside a cell is almost solid. We can think of a cell similar to a Christmas day on Time Square: Everything is densely packed, but there is still movement.

Question: What is the smallest building block of life that can replicate?

cells

Question: How many different cells are in a typical human?

200

Question: What are the parts of the cell called?

organelles

Question: Which part of the cell is called the "powerhouse"?

mitochondria

Question: What part of a plant is involved with photosynthesis?

chloroplast

Question: What is mitosis?

cell division

Question: Who first used the term cell?

Robert Hooke

Question: How many elements are found in amounts larger than trace amounts (0.01%) in our bodies?

11

Question: When communities of living things interact with non living things they are called ... ?

ecosystem

Question: The most common molecule in the human body is ... ?

Water: H₂O

Question: What do bacteria have in common?

Single Cells

Question: What is a gene?

A gene is a region of DNA, which contains all information for the creation of an entire RNA strand.

Question: What is DNA made out of?

DNA is a linear polymer out of 4 bases / nucleotides. DNA exists in cells mainly as a two-stranded structure called the double helix. Each of the bases has a complementary base.

- G: Guanine => Cytosine
- A: Adenine => Thymine
- T: Thymine => Adenine
- C: Cytosine => Guanine

Question: What is RNA made out of?

RNA is a single stranded linear polymer out of 4 bases / nucleotides.

- G: Guanine
- A: Adenine
- U: Uracil
- C: Cytosine

Question: Do all organisms use the same amino acids / codons?

Different organisms use the same amino acids for proteins. However, they differ in their codon usage (which RNA triplets are translated into which amino acid).

Question: How many proteins does a typical human have?

Between 20.000 and 25.000 different kinds of proteins.

Question: What are functions of proteins?

- Defense (e.g. antibodies)
- Structure (e.g. collagen)
- Enzymes (metabolism, catabolism)
- Communication / Signaling (e.g. insulin)
- Ligand binding / Transport (e.g. hemoglobin)
- Storage (e.g. ferritin)

Question: How many residues long are typical proteins?

Between 35 and 30.000 residues. The median is around 400.

Question: Do proteins consist of units?

Proteins are built up of several domains. Most proteins have more than 2 domains.

Question: How many proteins are known?

About 85 millions sequences are known. However, the 3D structure (experimentally determined) of only 120.000 proteins is known.

Question: Is this gap (known sequences vs known 3D structure) expected to increase?

Yes, the gap is expected to increase. The amount of new sequences has increased drastically (far faster than Moore's Law) in the past. This is expected to continue. Advances in experimentally determining protein 3D structure could only improve marginally, but today experimentally determining the 3D structure of a proteins still costs about 100 000 EUR.

Question: How much data is produced by one sequencing machine per day?

At full capacity about 5 - 10 terabytes of data per day.

3.3.3 Lecture 2: Introduction Protein Structure

Question: How many different amino acids are there?

20

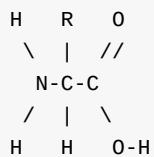
Question: How do amino acids differ? What do they have in common?

Different amino acids have different side-chains, which influence the chemical features of the respective amino acid. All of them share the same backbone.

Question: In which different feature groups can you categorize amino acids?

polar, non-polar, acidic (negatively charged), basic (positively charged)

Question: Draw the basic chemical structure of an amino acid.



Question: How are amino acids linked together to form a protein?

In the translation process, a **Ribosome** translates a **mRNA** strand to a protein, by decoding the RNA triplets into amino acids and then linking the amino acids by peptide bonds. They chaining ALWAYS happens from the **N-Terminus** to the **C-Terminus** releasing an H₂O molecule as part of the reaction.

Question: What is the definition of a 'domain'?

A domain is a protein sequence, which when put into solvent adopts a unique 3D structure on its own.

Question: How many domains does a protein have?

- 61% of proteins in the PDB are single domain
- 28% of proteins in the PDB are in 62 proteomes

Problem: This is a biased view on proteins. The 3D structure of Single-Domain-Proteins is easier to experimentally determine, so more Single-Domain-Proteins have been analyzed.

Question: Can domains overlap?

Yes, it can happen. However, it is not what is typically observed .

Question: How can we compare 3D structures?

One solution would be to align the corresponding residues of both sequences / 3D structures and take the **Root Mean Square Deviation**. (If one pair lies very far apart, it will result in an extremely low score)

$$RMSD(A, B) = \sum_{i=0}^n (r_i^a - r_i^b)^2$$

If the score is below a certain threshold, it is a match, otherwise it is not.

Question: How can align and compare the structure of 2 proteins?

- 1) Find the corresponding points (residues that match in 3D)
- 2) Find Superposition independent of domain movements and calculate score (e.g. RMSD)

Question: Why is global protein comparison most of the time impossible?

The definition of protein enforces a per residue comparison (no scaling). Hence only proteins of the (almost) the same length can be compared globally. Since proteins are between 35 and 30.000 residues long, global comparison does not make sense in most of the cases.

Question: What is the difference between global and local alignment?

In **global alignment** two structures / sequences are compared from beginning to end (compare the whole thing).

In **local alignment** however, subunits (domains) of the proteins are aligned. (Problem: What is a valid unit? Where to cut?)

Question: How to decide what is a valid unit for local comparison of 2 proteins?

(I couldn't identify a valid answer in the lecture recording)

Question: Which comparison not using cartesian RMSD could be used for comparison?

2D distance map: difference of differences. Only information about the chirality (mirror image) is lost.

3.3.4 Lecture 3: Alignments I

Question: Why compare 3D shapes, when we are after function? Why not compare function?

Because ...

- we cannot compare function directly
- structure is related to function
- we CAN compare 3D structures
- sometimes: similar structure -> similar function

Question: How do we get protein 3D shapes?

- primarily by experiment (most accurate)
- computational biology (most inferences)

Question: How much does it cost to experimentally determine the 3D shape of a protein?

Today it costs on average about 100 000 \$ per protein.

Question: What are the 3 sections found in the tree of life?

bacteria, archaea, eukaryotes

Question: What does Homology stand for?

Here (in the context of genes), it describes proteins originating from a common ancestor. It is also frequently used to describe 'similar structure' for genes / proteins.

Question: Why do linear gap penalties not model the reality of related genes / proteins well?

With a linear gap penalty (N gaps cost $N \cdot x$) equally distributed gaps would be as expensive as clustered gaps. Biologically, gaps clustered to blocks, are however far more likely to occur, while the protein maintains similar structure / function. It is more realistic to use an **Affine gap penalty** with higher costs for opening a new gap.

Question: What is better? High sequence identity of a short (local) sequence, or worse sequence identity when matching a longer sequence? How can we decide?

Compile the probability of randomly matching a sequence considering the background distribution. The result of this would be a substitution matrix such as BLOSUM62.

Question: Is identity the best way to match two sequences?

Not necessarily: What we really find is similar biological function. Some amino acids might have similar biophysical features and could be swapped without any significant influence on the structure of the protein. Such matches should also be considered 'positive'.

Building a scoring matrix based on evolutionary conserved residues does optimize the algorithm. (e.g. BLOSUM62)

Question: What is the biological assumption behind an insertion when comparing sequences?

Through evolutionary changes in the DNA (e.g. a point mutation) a new bump (= amino acid(s)) was introduced. Implicitly it is also assumed similar structure -> similar function.

Question: Why do linear gap penalties not model the reality of related genes / proteins well?

With a linear gap penalty (N gaps cost $N \times x$) equally distributed gaps would be as expensive as clustered gaps. Biologically, gaps clustered to blocks, are however far more likely to occur, while the protein maintains similar structure / function. It is more realistic to use an **Affine gap penalty** with higher costs for opening a new gap.

Question: Does dynamic programming give the best solution?

Yes, dynamic programming produces one optimal solution. (There could be others, though)

Question: What are issues with dynamic programming?

- Time used: $O(n^2)$
 - Especially a problem, when comparing one protein against the entire database.
- How to choose parameters?
 - Gap penalties
 - substitution matrix

Question: How can we speed up the alignment of sequences?

1) Hashing (fast and dirty). e.g. BLAST

Question: How does BLAST (Basic Local Alignment Search Tool) work?

1. Start with indexed (hashed) seeds (words of size = 3) and find matching proteins
2. Extend matching 'words' into both directions
3. Begin dynamic programming from these strong local hits

3.3.5 Lecture 4: Alignments II

Question: What is the major challenge of BLAST?

Getting the statistics right: BLAST needs to know, how *significant a match is*, by comparing it against the background probability of the entire database.

Question: Why is it interesting to find similar proteins out of the Twilight / Midnight Zone?

The Midnight-Zone is, where most proteins of similar structure sit.

Question: Why is it that even with only 40% PSI, we can still assume similar structure? Could we randomly change 60% of the residues in the lab and get a new protein with similar structure?

- These 60% of changed residues happened under evolutionary pressure and are not random
 - mutations that did not change structure and function survived (we can observe them)

today)

- mutations that did change structure and function most likely did not survive
- Thus randomly changing 60% of residues in a protein, would not result in a similar protein

Question: Why are certain proteins / structure multiple times in the PDB?

- different resolution of 3D structure
- different goals of publication produced (new) 3D structures
 - folding sites
 - binding partners
 - etc ...

Question: How are profiles built up? How are the normal noted down? Do we have to know a specific algorithm?

Build up algorithm:

- Take all proteins of PSI over a certain threshold ...
-

Profile Formats:

- Regular Expression
- PSSM (Position Specific Scoring Matrix)

Question: What is a PSSM (Position Specific Scoring Matrix)?

A matrix of numbers with scores for each residue or nucleotide at each position. Built, e.g. by PSI-BLAST.

Question: Which steps are involved in building up a profile with PSI-BLAST?

- 1) **Fast Hashing:** Like BLAST, match 'word'
- 2) **Dynamic Programming Extension between matches:** BLAST + Smith-Waterman
- 3) **Compile Statistics:** EVAL - Expectation Values
- 4) **Collect all pairs and build profile**
- 5) ... compare sequences (profile-sequence) and iterate

Question: Why is PSI-BLAST so fast?

Because it drastically reduces the length of the comparisons with dynamic programming.

3.3.6 Lecture 5: Comparative Modeling

Question: How do you build up a family (profile) of sequences?

1. Find proteins of similar structure with BLAST
2. Build PSSM
3. build up a set of pairwise alignments

4. add those over a certain HSSP value to the family
5. Search with profile-sequence comparison for more distant family members and refine profile

Question: Which methods to experimentally determine the structure of a protein exist? How much are they used?

Fraction of proteins in the PDB by experimental method:

- 90% - X-Ray Crystallography
- 09% - Nuclear Magnetic Resonance Spectroscopy (NMR)
- 01 % - Electron Microscope (EM)

Question: How does X-Ray Crystallography Work

1. **Grow Crystal:** Force the protein to grow a crystal
2. **Observe Diffraction Pattern :** Shoot x-rays onto crystal and observe the diffraction pattern
3. **Compute Electron Density Map**
4. **Fit observations to atomic model**

Question: How to get 1D secondary structure from 3D coordinates?

Two methods were used to annotate 3D coordinates:

1) DEFINE, based on geometry (not used anymore) 2) DSSP, based on hydrogen bond pattern (coulomb energy)

Question: How does Homology Modeling (Comparative Modeling) work?

Target: Protein to model

Template: Protein to model from

1. **Identify Template:** Query the PDB for similar sequences to your **Target**
2. **Align Target / Template:** Select the best match as **Template** and assume the **Target** has the same structure
3. **Build Model**
4. **Assess Model**
5. **Refine Model**

Question: Which tradeoff does comparative modeling face? What are the limiting factors based on PSI (Percentage Sequence Identity)?

Tradeoff: Accuracy vs Coverage

Limiting factor in homology modeling:

75% - 100% - Speed of Modeling

50% - 75% - Quality of Model

25% - 50% - Alignment Accuracy

0% - 25% - Detection of Homology

Question: How to handle a missing loop in comparative modeling?

- One way would be to find similar loops and compute the average over them.
- Another solution would be to apply molecular dynamics on the loop sequence. (only for short loops)

3.3.7 Lecture 6: Secondary Structure Prediction 1

3.3.8 Lecture 7: Secondary Structure Prediction 2

Question: Relate the terms **Local** and **Global alignment** to the terms **Sequence-Sequence** and **Sequence-Profile**.

Global alignments refers to aligning sequences (proteins) from start to end. Local alignments refers to only aligning parts of the sequences (e.g. 50 residues).

Throughout Sequence-Sequence, Sequence-Profile and Profile-Profile methods both global and local alignment can be used. I practice mostly local alignment is done.

Question: What would be a simple method to predict secondary structure?

- 1) Take known structure
- 2) Find longest consecutive run of motifs that **ONLY** occur in one of the 3 states: H (Helix), E (Strand), O (Other)
- 3) Check unknown sequence against found motifs

Question: What was the first secondary structure prediction method?

Assuming that a **Proline** would break a helix, the occurrences of proline in a sequence was used to predict helices.

Question: Where do we get the secondary structures from?

From the DSSP, which defines 8 states in total based on H-bond patterns.

Question: What is the 1st generation of secondary structure prediction based on? What was the accuracy? Was it successful?

- Based on single residues
- Between 50% and 55% accuracy (Q3)
- Clearly better than random - so it can be considered a success

Question: How did the second generation of secondary structure prediction improve? Name one algorithm.

Instead of using only single amino acids, it would consider a sliding window of the residues around a center amino acid. **Example:** GORIII, with a Q3 accuracy of 55% - 60%

Question: What were problems of secondary structure prediction until 1994?

- the maximum accuracy of predictions was expected to be 65%

- β -sheet prediction was below 40%
- many predicted segments were too short to appear in nature

Question: How can the performance of secondary structure prediction be measured?

One way to do it, would be to calculate the **Q3** accuracy of a method against a test set. The Q3 accuracy is the **number of correctly categorized residues into one of the categories helix, strand, other divided by the total amount of residues.**

Question: How can the introduction of a new hidden layer in a neural network be described by means of a simple graph?

Each new hidden layer basically introduces a new 'decision line' which can separate datapoints into different categories.

Question: What is cross-validation in the context of Machine Learning and why do we need it?

Cross Validation is a method for estimating the performance of a predictive model (e.g. a neural network). To use it, the available dataset is split in 3 categories, 1) a training set, 2) a cross-training set and 3) a test set.

- 1) The training set is used to train the model
- 2) The cross-training set is used to estimate the performance of the model after x training steps
- 3) The test set is used to assess the final performance of the model after training is finished

The cross-training set is needed to decide, when to stop training (when overtraining sets in) and to tweak certain parameters before running against the test-set.

Question: Did balanced training improve the Q3 prediction accuracy? Which assumption did it prove wrong?

Balanced training actually decreased the Q3 accuracy. However, it did improve the prediction accuracy for strands significantly, falsifying the hypothesis that strands could not be predicted with local information.

Question: Which problem did PHDSec solve? How did it accomplish it?

By introducing a **Structure-to-Structure** prediction model, PHDSec improved the prediction of too short segment. The Structure-to-Structure network would take structure (helix, strand, other) prediction of a sequence as input and predict segments based on them.

3.3.9 Lecture 8: Secondary Structure Prediction 3

Question: Which ways of comparing proteins are there? Why do we need

- Dynamic Programming (Brute Force)
- Hashing (e.g. BLAST)

Question: Why are fast search algorithms such as BLAST needed?

Comparing sequences of length n residues is in $O(n^2)$. For comparing a single pair this is still fine, but comparing one (newly found) protein against all known proteins in the PDB (about 120 000) is impossible. Thus we need 'shortcuts' such as BLAST to speed up the search.

Question: What is the normal approach when you find / analyse a newly found protein?

- 1) Sequence the new protein (if not done yet)
- 2) Run BLAST against the PDB
- 3) Run Dynamic Programming against the results from BLAST

Question: In terms of CPU, is sequence-sequence as fast as sequence profile?

Question: How can it be that even with only 40% sequence identity we assume / observe similar structure?

The changes in sequences we observe are not random, but follow underlying evolutionary rules. Changes, which affected the structure and thus the function of a protein are simply not likely to survive and thus we do not observe them. Changes, which did not influence the structure / function however, did survive.

Question: Why is protein sequence changing? Why are we mutating?

- Replication Errors (point mutations)
- Radiation
- Viruses

Question: How much do any two unrelated typical humans differ on average?

On average every pair of humans would differ in one amino acid per protein. (Though, changes cluster)

Question: In a structure to structure network, which additional information could be used to improve the prediction?

- E.g. redundant information about the sequence, e.g. parts of it.

Question: When training a neural network, how do you choose the next training sample from your test set? Why so?

Randomly, to avoid correlations

Question: How would you build up a family for a protein?

1. Search the PDB for proteins in comparative modeling range. (Assumption: same sequence, same 3D structure, same secondary structure)
2. Use profile to search in twilight-zone for potential proteins of that family (possibly verify whether the found protein is plausible to have similar 3D structure) and add to family (recompute profile)

Question: How do you get from a sequence to a secondary structure prediction with PHD?

1. Use BLAST to find potentially similar proteins in sequence data bank
2. For the resulting proteins calculate the sequence identity (homology) with dynamic programming
3. Filter all proteins, which are below a threshold of sequence identity (only take those "over the curve")
4. Extract the profile by aligning the remaining proteins
5. Predict the secondary structure with the sequence and its family as input

Question: Which accuracy does ProfSec achieve on average? What are additional advantages of other secondary structure prediction methods?

ProfSec achieves a Q3 accuracy of about 72% on average. Additionally it can also predict the strength of the prediction.

Question: Does adding global information improve ProfSec prediction?

Yes it does. While the Q3 accuracy (per residue) is not improved, the Q4 accuracy (per protein) does improve.

3.3.10 Lecture 9: Membrane Structure Prediction

Question: What are the requirements of a cell membrane?

- separate the content of the cell from its surroundings
- control traffic into and out of the cell
 - keep malicious things out
 - let good things in
- must be a dynamic structure

Question: What are the 4 main structural components of the cell membrane?

Carbohydrates, Cholesterol, Phospholipids, Proteins

Question: What is the cell membrane mainly made out of?

The cell membrane is a so called **lipid bilayer of phospholipids**. Phospholipids have a non-polar, hydrophobic tail (membrane center) and a polar, hydrophilic head (outside of membrane).

Question: What are functions of membrane proteins?

- help to be recognized by immune cells
- transport proteins control substance flow in and out of the cell
- receptor proteins bind hormones, which can change cell function
- provide structural stability

Question: Can membrane proteins easily move around?

It depends:

- Membrane proteins can easily move laterally
- But it is hard to move into / out of the lipid bilayer

Question: Why are there so few membrane proteins in the PDB?

It is particularly difficult to experimentally determine the structure of membrane proteins due to the special environment they naturally occur (the membrane).

Question: What are the key questions TMH prediction tries to answer?

- How many helices go through the membrane?
- In which direction do they go through the membrane? (topology)

Question: Why could be a plausible reason why PHDSec failed for predicting transmembrane helices?

Unlike 'normal' proteins, transmembrane proteins have an hydrophobic outside and a hydrophilic inside.

Question: How should we choose the threshold for the hydrophobic regions?

1. Predict the hydrophobicity for the protein
2. Assign a positive inside-out
3. choose the threshold to **optimize the inside out difference**

Question: What is the Positive Inside Rule and what is it used for?

The positive Inside Rule is used to find the topology of transmembrane proteins: The loops connecting TMHs on the inside of the cell membrane have an **excess of positively charged residues**.

3.3.11 Lecture 10: TMSEG

Question: What are advantages of using a Random Forest?

- Fast
- robust against overtraining
- no black box
- Intuitive to interpret
- good performance

3.3.12 Lecture 11: Beta Membrane and Accessibility

3.3.13 Lecture 12: Protein Disorder

