



# 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 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 challenge 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](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 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 organelles

**Eukaryotic Cells:** Found in animals and plants

- nucleus
- usually multicellular
- 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: H<sub>2</sub>O

**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

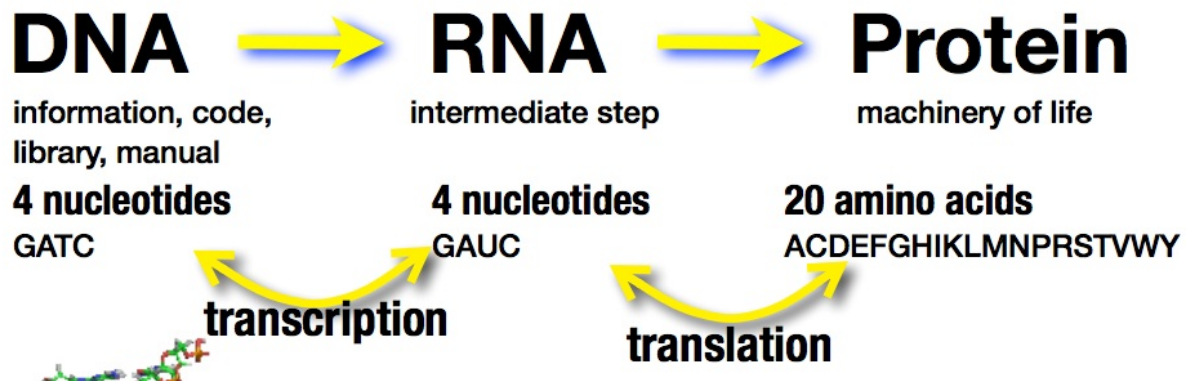
**Question:** What is a gene?

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](#)

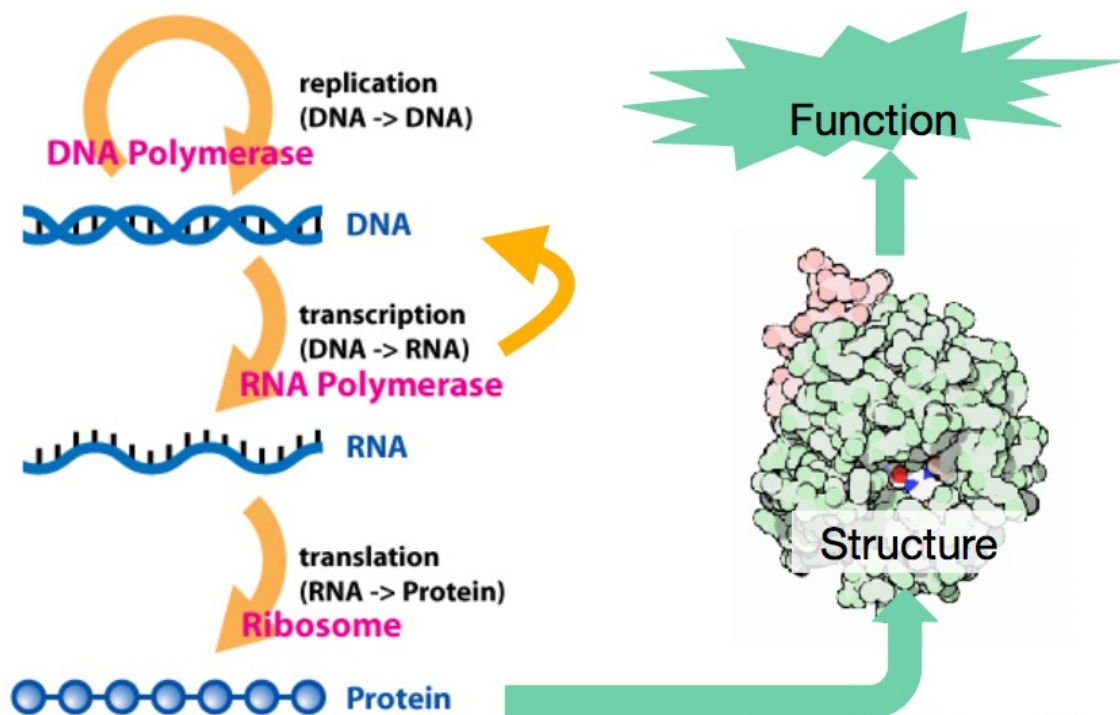
### 1. Recap

**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



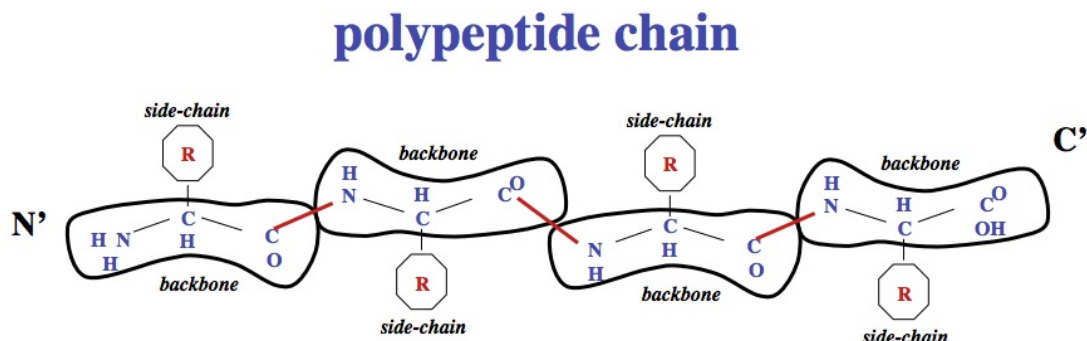
**Translation:** Proteins are made up of amino acids (20 different kinds). Each amino acid is encoded by a nucleotide triplet (codon) of DNA / RNA.

### 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**.



- 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 H<sub>2</sub>O 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

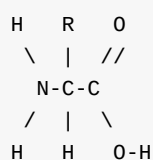
**Question:** How do amino acids differ? What do they have in common?

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.



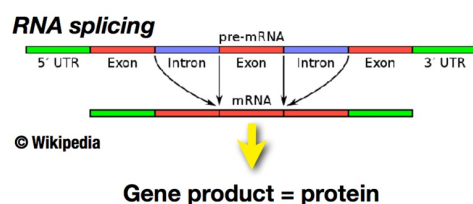
**Question:** How are amino acids linked together to form a protein?

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 H<sub>2</sub>O 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 know 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)

$$\text{RMSD}(A,B) = \text{SQRT}(\sum (a_i - b_i)^2)$$

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 Bono an 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 'postive'.



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  $N \cdot x$ ) 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.

## BLOSUM62

### BLOSUM (Scoring Matrix)

- ❑ **BLO**cks of amino acid **S**ubstitution **M**atrices
- 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](#)

---

## 1.3 Comparative Modelling

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

---

## 1.6 Secondary Structure Prediction

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

---

## 1.7 Secondary Structure Prediction 2

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

---

## 1.8 Secondary Structure Prediction 3

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

---

## 1.9 Membrane Structure Prediction 2

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

---



## 1.10 TMSEG

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

---

# 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  $\beta$ -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 Ångström, 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

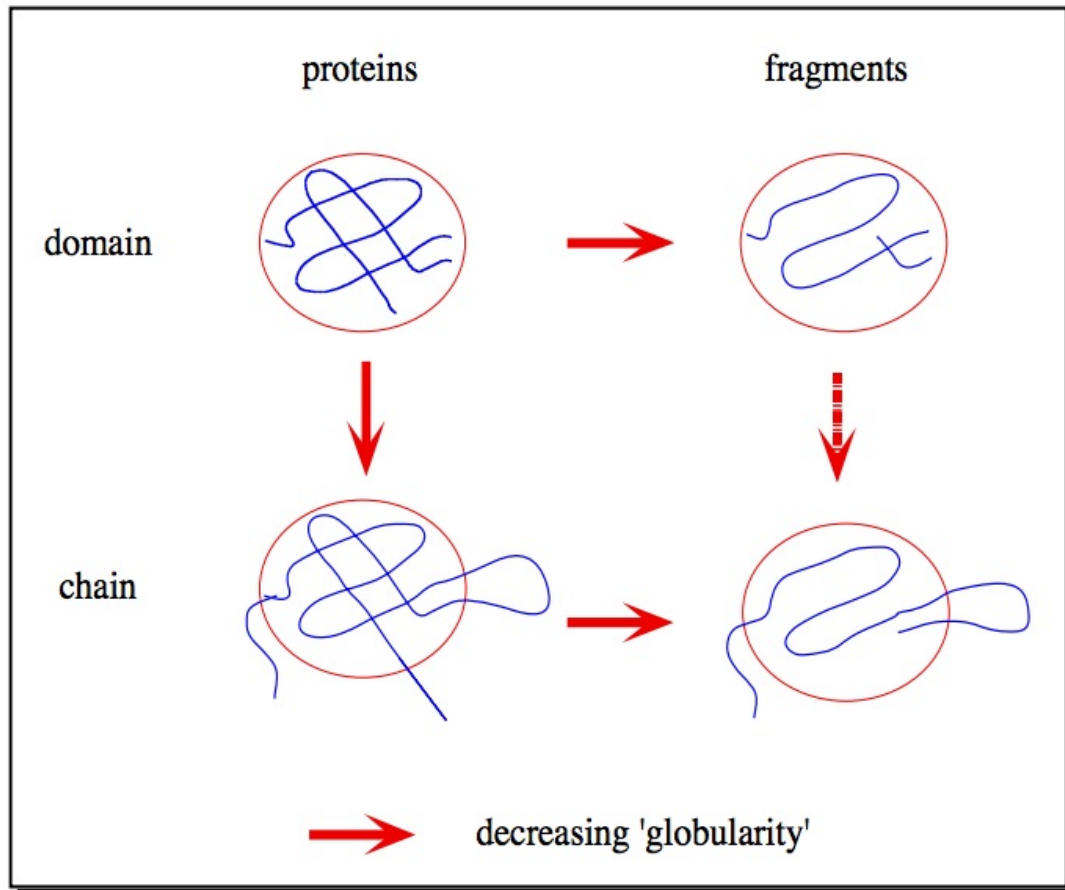
**Detailed Prediction Problematic:**

**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: Protein 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?

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

**Question:** What is a structural domain? What is a functional domain?

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

### Questions (Tuesday, 27th June)

### 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?

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

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



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

**Question:** How can you predict structure in the [a] daylight- [b] twilight- [c] midnight-zone?  
**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  $i+1$ .  
(Short alignment  $i$  and  $i+1$  are independent).

Method w/o that?

**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.

**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?

**Question:** How do you compare proteins of different length?

**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?

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

## 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 year's 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 ( $\alpha$ -helix,  $\beta$ -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 speeds 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: H<sub>2</sub>O

**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

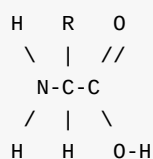
**Question:** How do amino acids differ? What do they have in common?

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)

$$\text{RMSD}(A,B) = \text{SQRT}( \sum (a_i - b_i)^2 )$$

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  $N \cdot x$ ) 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:** 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