



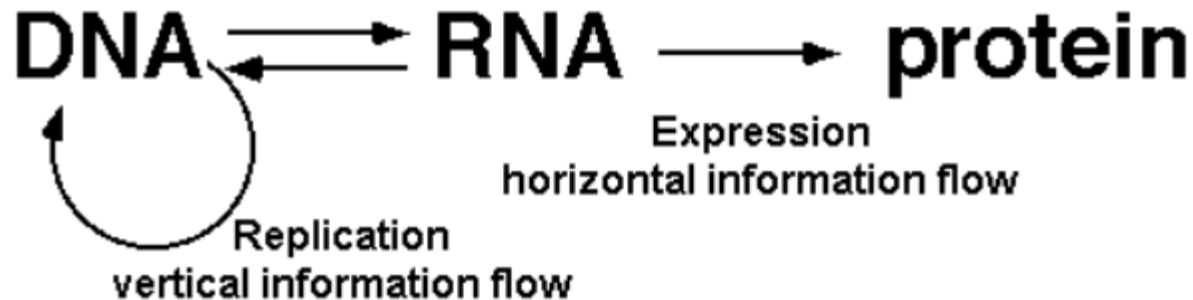
Proteins: Structure & Function

Ulf Leser

This Lecture

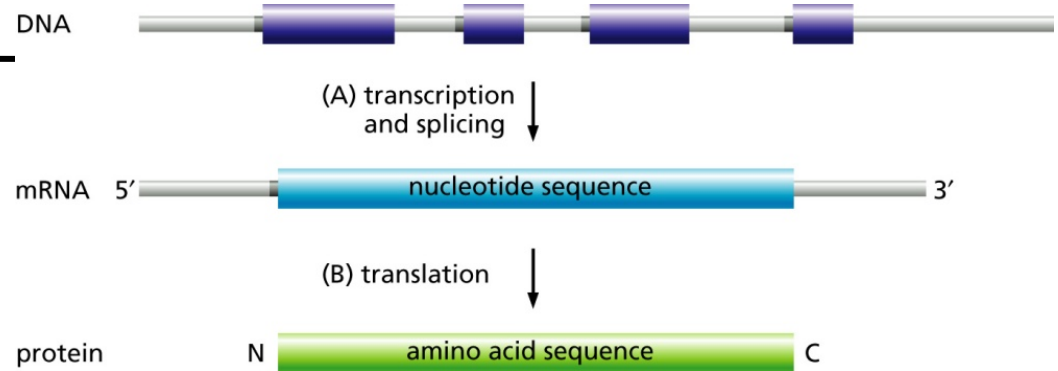
- Proteins
 - Structure
 - Function
 - Databases
- Predicting Protein Secondary Structure
- Many figures from Zvelebil, M. and Baum, J. O. (2008). "Understanding Bioinformatics", Garland Science, Taylor & Francis Group.
- Examples often from O. Kohlbacher, Vorlesung Strukturvorhersage, WS 2004/2005, Universität Tübingen

Central Dogma of Molecular Biology



	U		C		A		G	
U	UUU Phenyl- alanine UUA Leucine UUG		UCU Serine UCC UCA UCG		UAU Tyrosine UAC UAA Stop codon UAG Stop codon		UGU Cysteine UGC UGA Stop codon UGG Tryptophan	U C A G
C	CUU Leucine CUC CUA CUG		CCU Proline CCC CCA CCG		CAU Histidine CAC CAA Glutamine CAG		CGU Arginine CGC CGA CGG	U C A G
A	AUU Isoleucine AUC AUA Methionine; initiation codon AUG		ACU Threonine ACC ACA ACG		AAU Asparagine AAC AAA Lysine AAG		AGU Serine AGC AGA Arginine AGG	U C A G
G	GUU Valine GUC GUA GUG		GCU Alanine GCC GCA GCG		GAU Aspartic acid GAC GAA Glutamic acid GAG		GGU Glycine GGC GGA GGG	U C A G

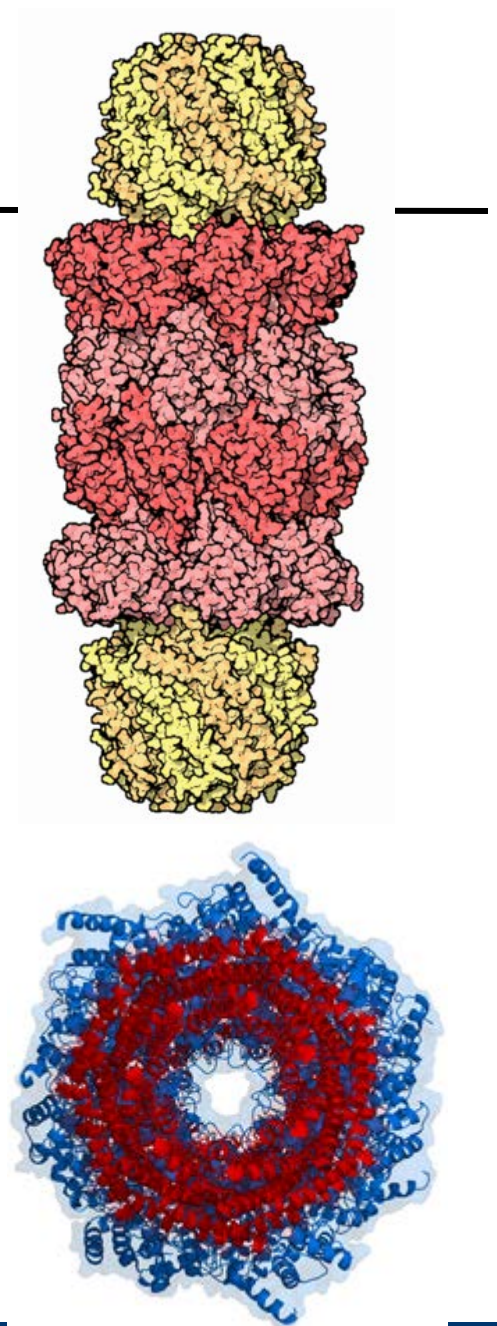
Details



- Alternative Splicing
 - “One gene – one protein” is wrong
 - Exons may be spliced out from the mRNA
 - Human: at least **6 times more unique proteins** than genes
- Post-translational modifications
 - (De-)Phosphorylation, glycolysation, cleavage of signal peptides, ...
 - Human: At least **5 times more protein forms** than proteins
- Complexes: Proteins **physically group together** to perform specific function

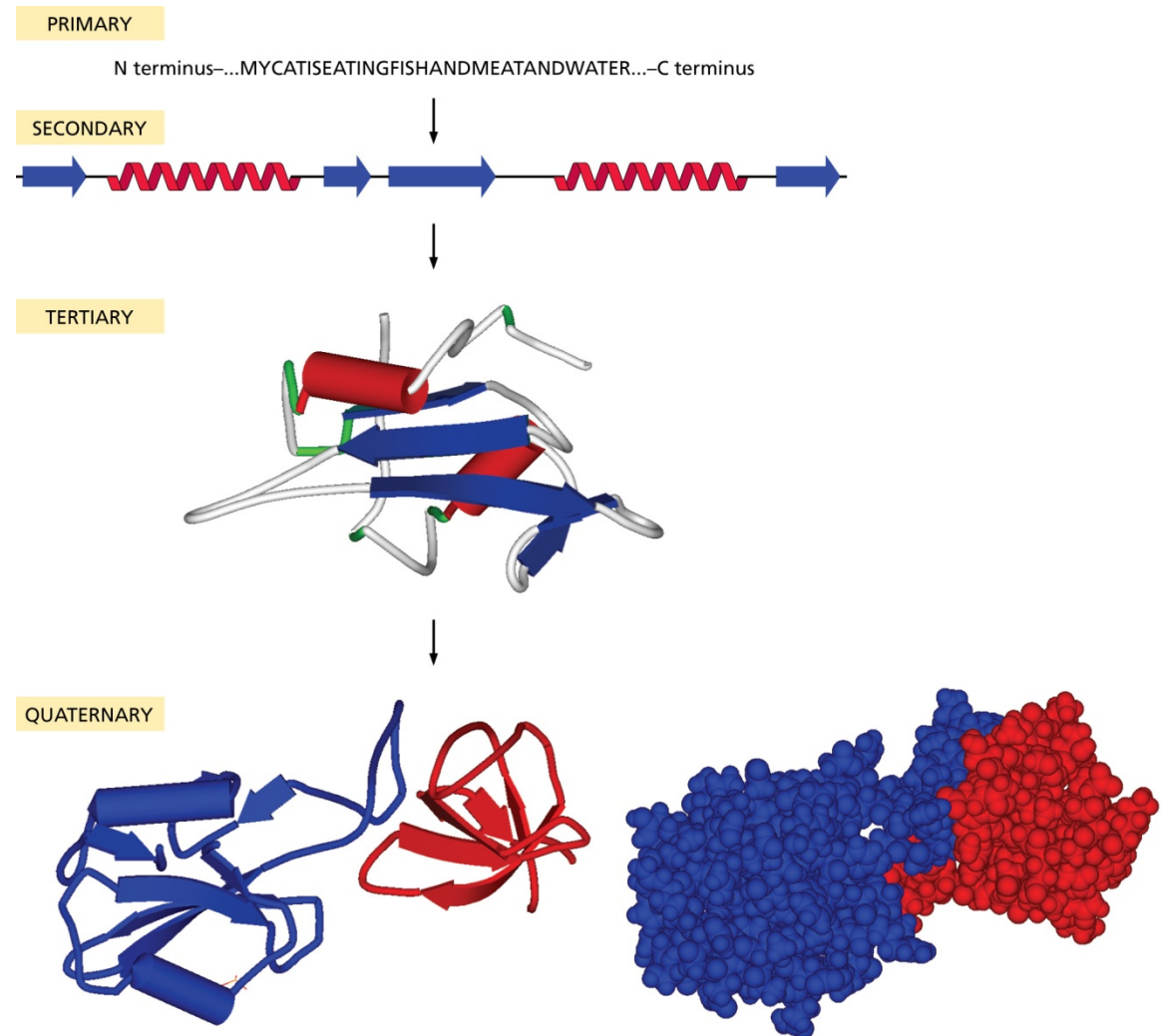
Example: Proteasome

- Function: Breaks (mis-folded, broken, superfluous, ...) proteins into small **peptides for reuse**
- Very large complexes present in all eukaryotes (and more species)
 - >2000 kDa, consists of **dozens of single proteins**
 - Formation of the complex is a complex process only partly understood yet



Protein Structure

- Primary
 - 1D-Seq. of AA
- Secondary
 - 1D-Seq. of “subfolds”
- Tertiary
 - 3D-Structure
- Quaternary
 - Assembled complexes



Protein Function

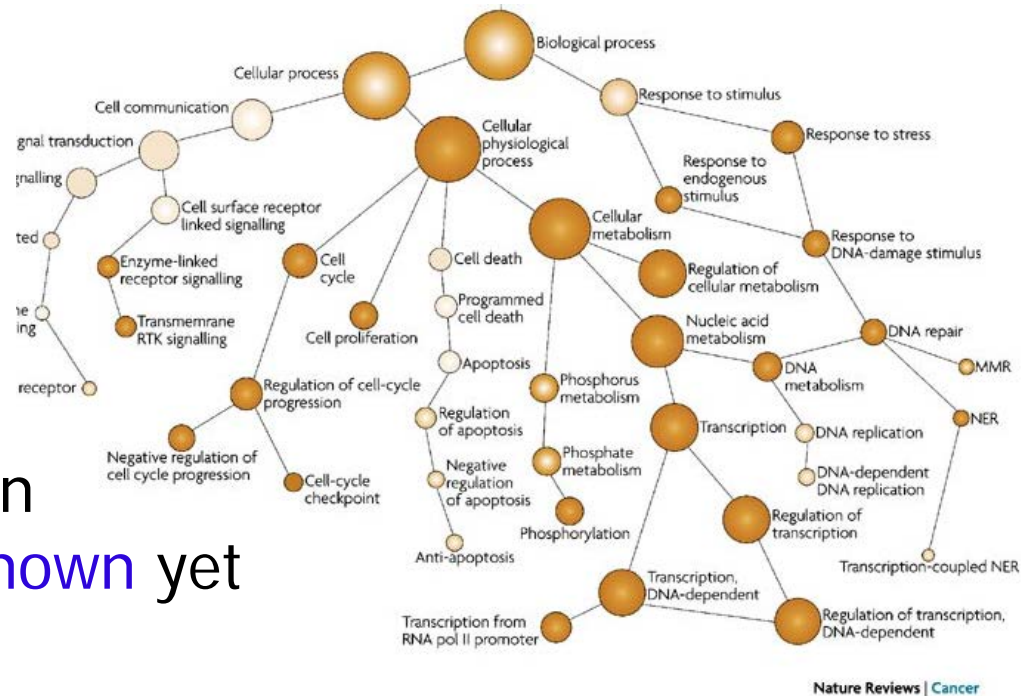
- Proteins perform essentially everything that makes an **organism alive**

- Metabolism
- Signal processing
- Gene regulation
- Cell cycle
- ...

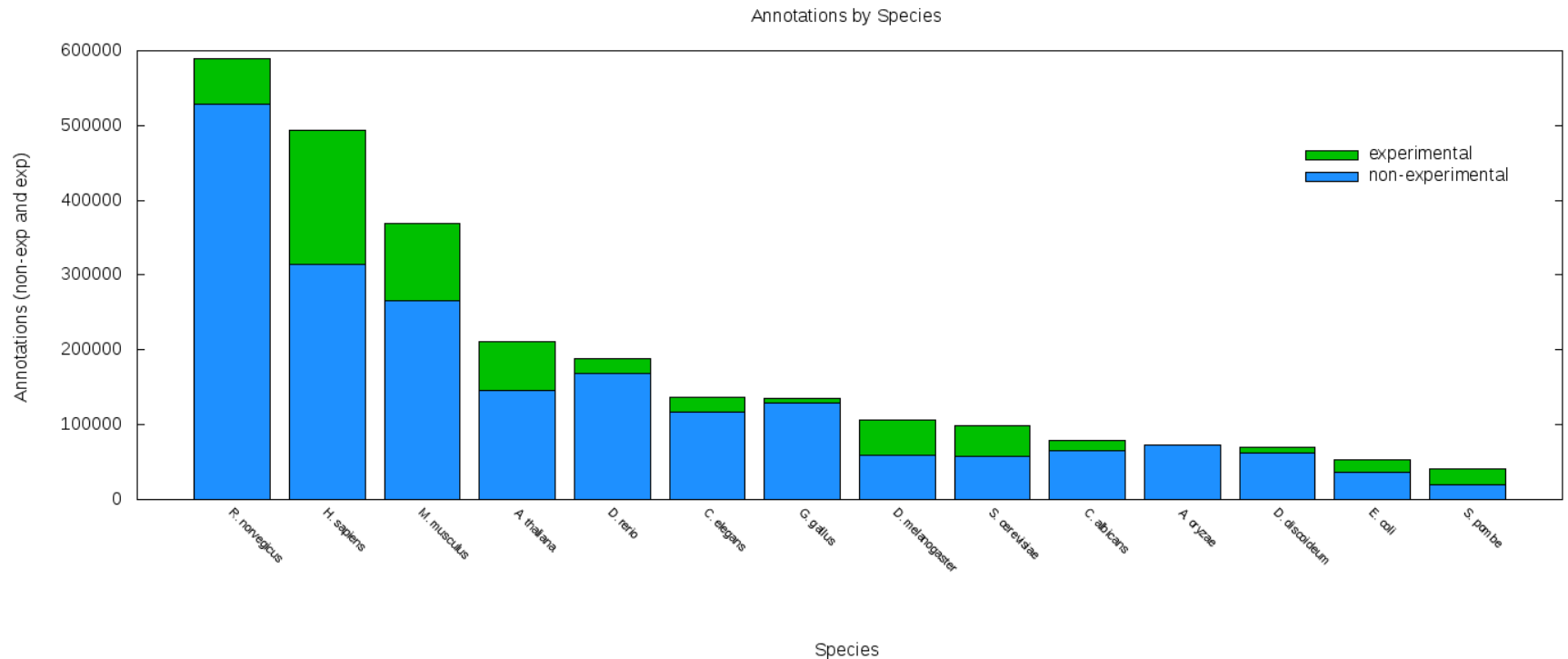
- For ~20% of all human gene, **no function is known** yet

- Describing function

- **Gene Ontology**: 3 branches, >30.000 concepts
- Used world-wide to describe gene/protein function

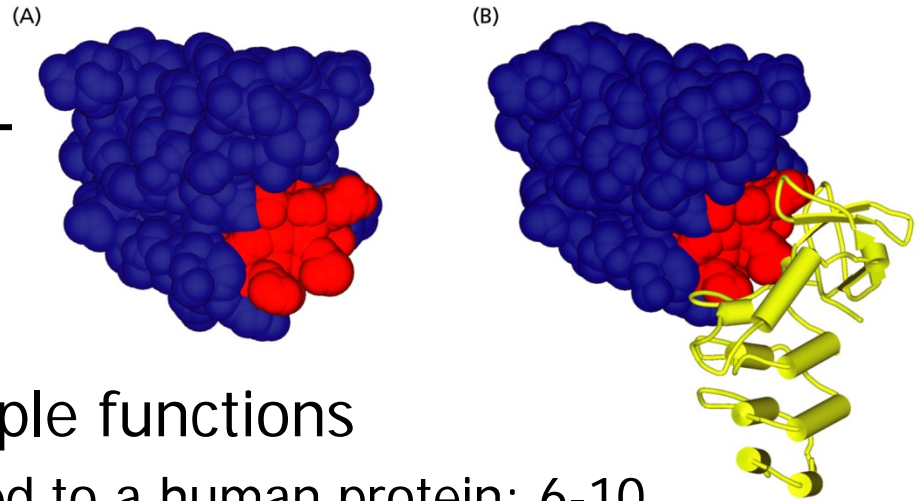


„Known“ Protein Functions



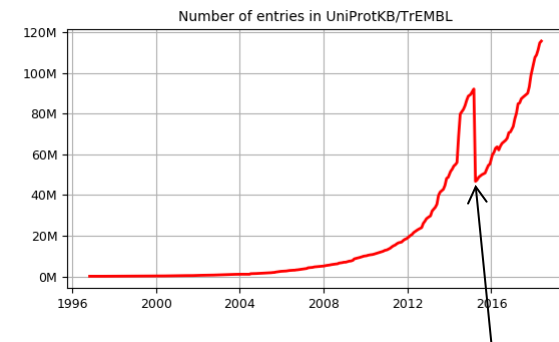
<http://geneontology.org/page/current-go-statistics>, June 2016

Function and Motifs



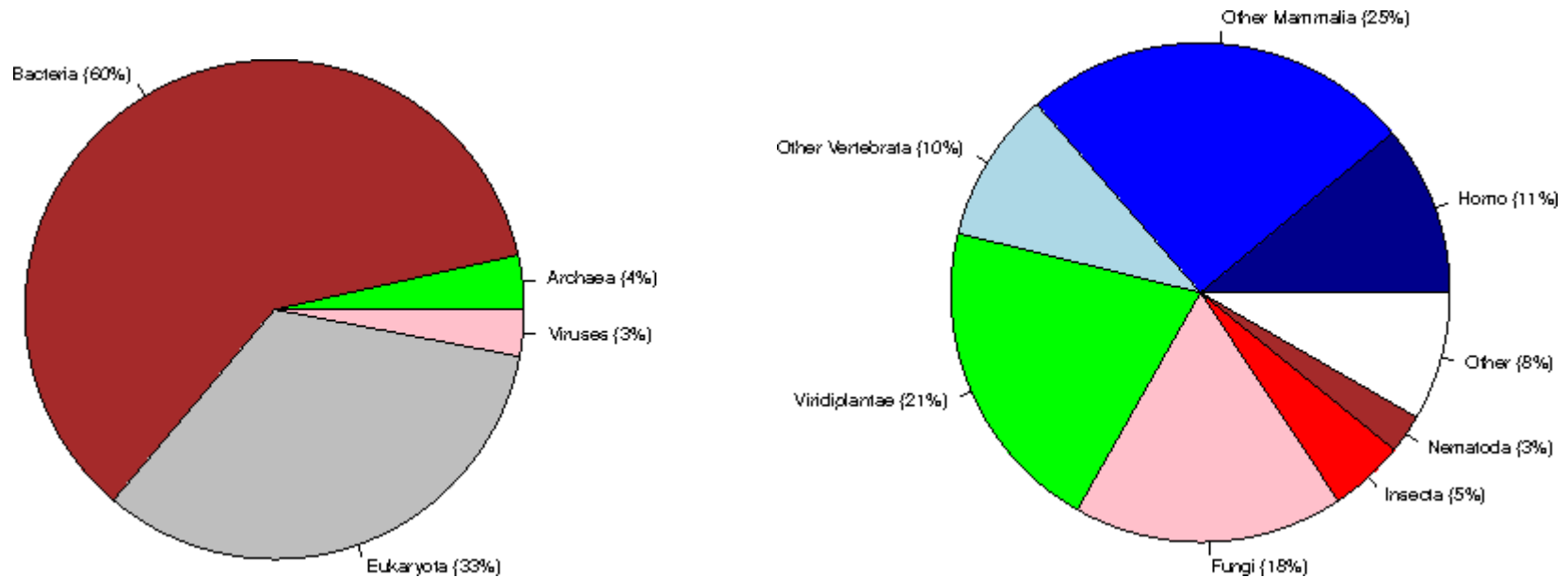
- Proteins usually have multiple functions
 - Avg. n# of GO terms assigned to a human protein: 6-10
- Functions are associated to **motifs or domains**
- There probably exist only 4000-5000 motifs
 - Proteins as assemblies of functional motifs
- Performing a function often requires **binding to another protein** or molecule
 - The binding requires a certain constellation of the protein structure
 - Major target of **pharmacological research**

- “Standard” database for **protein sequences and annotation**
 - Original name: SwissProt
 - Started at the Swiss Institute of Bioinformatics, now mostly EBI
 - Other: PIR, HPRD
- Continuous growth and **curation**
 - >30 „Scientific Database Curators”
 - Quarterly releases
 - **Very rich set of annotations**
- Actually two databases
 - **SwissProt**: Curated, high quality, versioned
 - TrEMBL: Automatic generation from (putative) coding genomic sequences, low quality, redundant, much larger



Def. and removal of
„redundant” sequences

UniProt: Species [<http://www.expasy.org/sprot/relnotes/relstat.html>, June 2016]

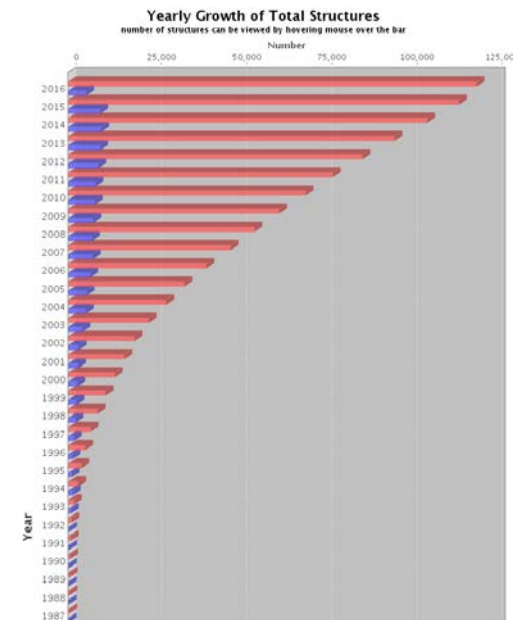


20258	Homo sapiens (Human)
16327	Mus musculus (Mouse)
9842	Arabidopsis thaliana (Mouse-ear cress)
7560	Rattus norvegicus (Rat)
6582	Saccharomyces cerevisiae (Baker's yeast)
5803	Bos taurus (Bovine)

...

PDB – Protein Structure Database

- Oldest protein database, evolved from a book
- Contains experimentally obtained **protein 3D-structures**
 - Plus some DNA, protein-ligand, complexes, ...
 - X-Ray (~75%), NMR (nuclear magnetic resonance, ~23%)
- Costly and **rather slow techniques**
 - Growth much smaller than that of sequence-related DBs
- Many problems with **legacy data** and data formats



<http://www.rcsb.org/pdb/statistics/contentGrowthChart.do?content=total>, June 2016

InterPro

- **Integrated database** of protein signatures, classifications, and motifs
 - Currently ~21.000 signatures
- Associates signatures with function (GO term)
- **InterProScan** – quick identification of signatures in a protein sequence
 - For a fast, first functional annotation



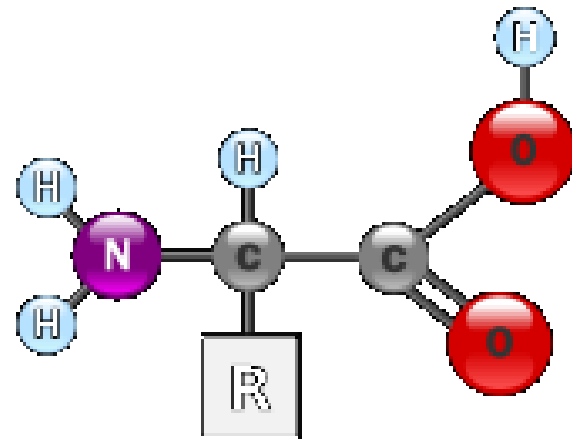
InterPro 30.0

This Lecture

- Introduction
- Predicting Protein Secondary Structure
 - Secondary structure elements
 - Chou-Fasman
 - GOR IV
 - Other methods

Amino Acids (AA)

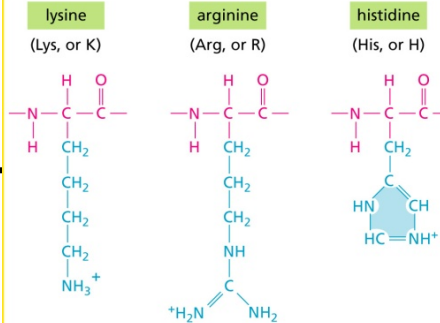
- AA consist of a common core and a **specific residue**
 - Amino group – NH_2
 - Central C_α - Carbon – CH
 - Carboxyl group – COOH



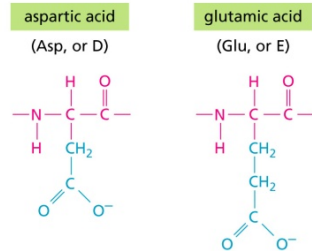
- Residues (side chains) vary greatly between AA
- Residues determine the **specific properties** of a AA

Side Chains

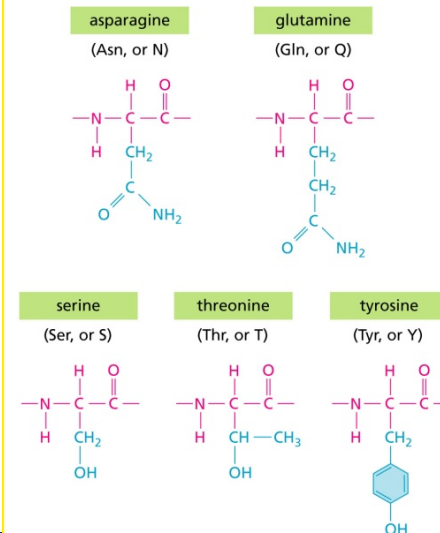
BASIC SIDE CHAINS



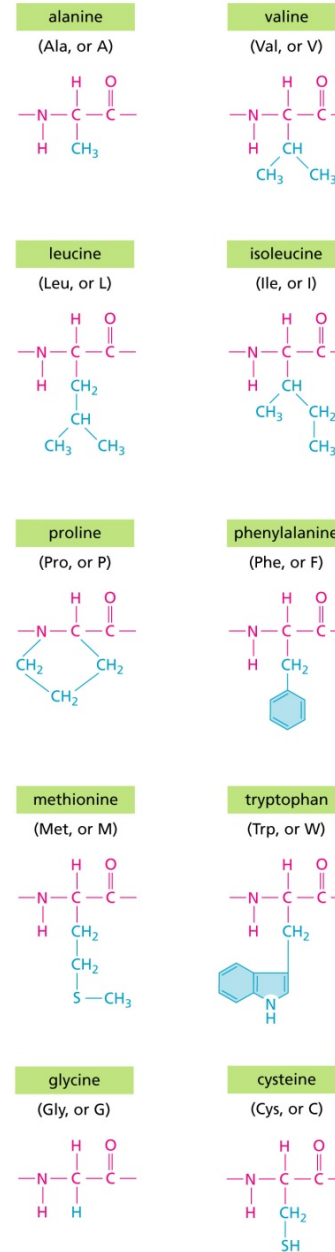
ACIDIC SIDE CHAINS



UNCHARGED POLAR SIDE CHAINS

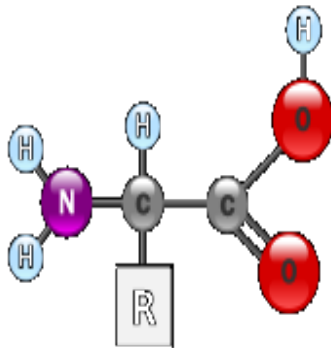
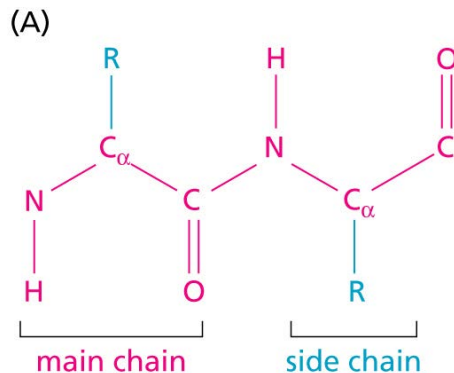


NONPOLAR SIDE CHAINS

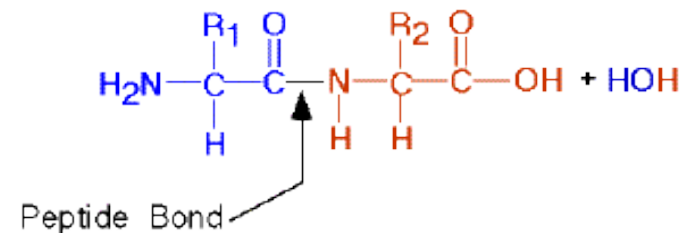
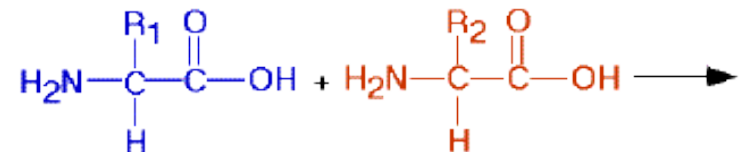


Structure of a Protein

- Concatenation of **cores**: Backbone of AA chain (a protein)
- Covalent **peptide bonds** between carboxyl and amino group (with loss of H_2O)

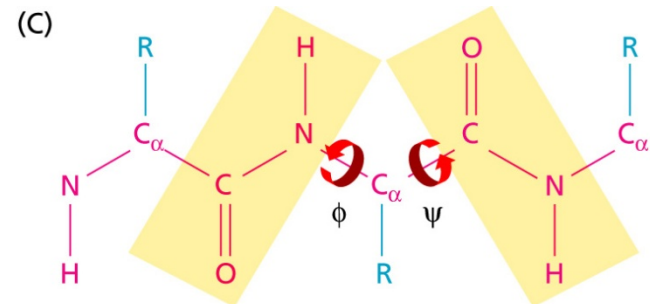
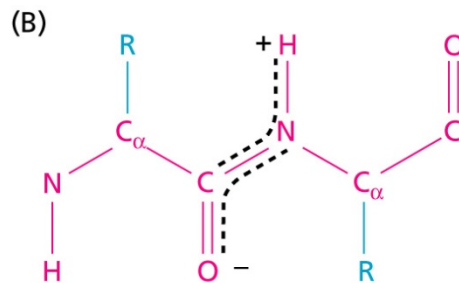
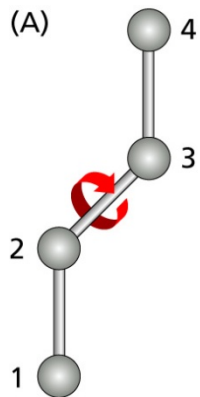


Peptide Bond Formation

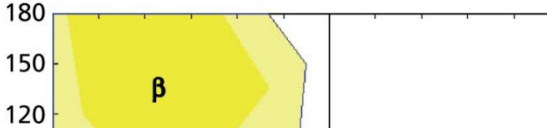


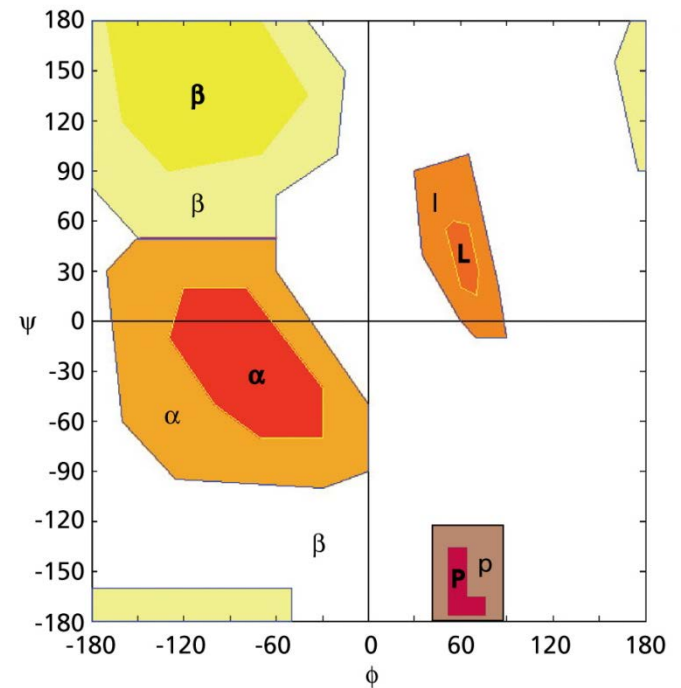
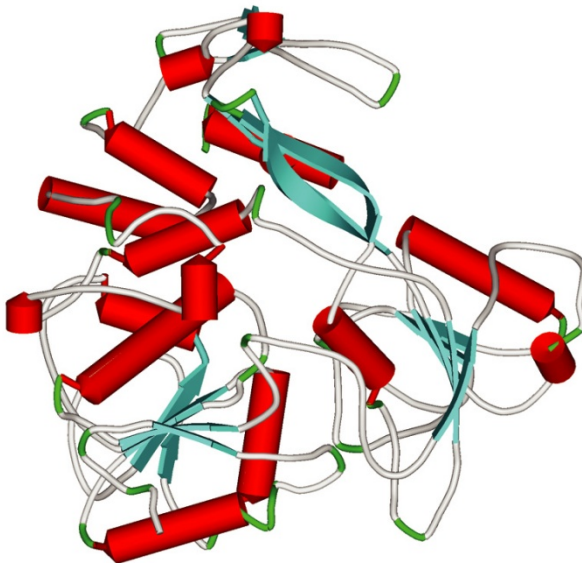
Flexibility

- In principle, every chemical bond can rotate freely
 - Would allow arbitrary backbone structures
- In real proteins, things are more restricted
 - Peptide bond (B) is “flat” – almost no torsion possible
 - Flexibility only in the C_α -flanking bonds ϕ and ψ

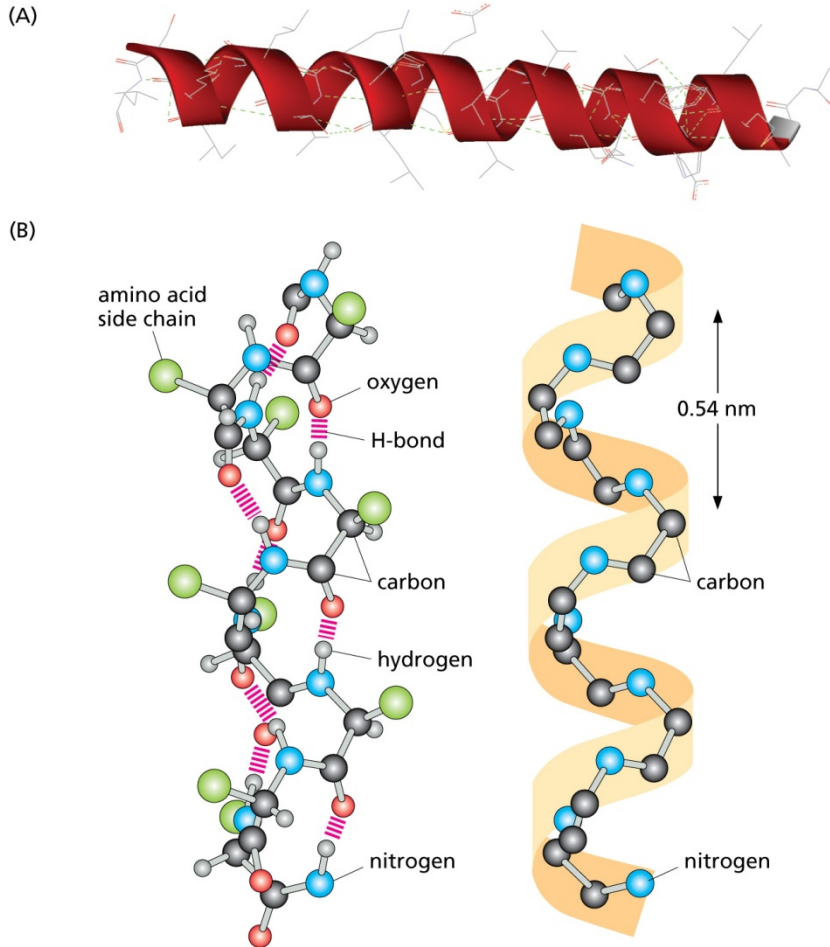


Ramachandran Plots

- Combinations of ϕ and ψ are **highly constrained**
 - Due to chemical properties of the backbone / side chains
 - Two combinations are favored: **α -helixes** and **β -sheets**
 - More detailed classifications exist
 - Angles lead to specific structures
 - **Secondary structure**
- 



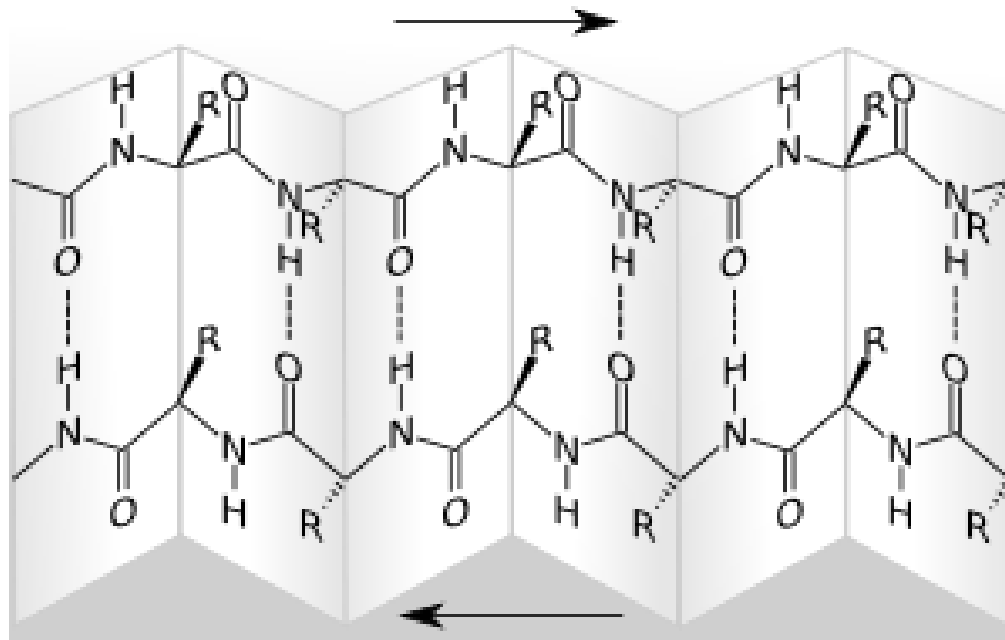
α -Helix



- Sequence of angles forming a regularly structured **helix**
- Additional bonds between amino and carboxyl groups
 - Very **stable structure**
- May have two orientations
 - Most are right-handed
- 3.4 AA per twist
- Often short, sometimes very long

β -Sheet

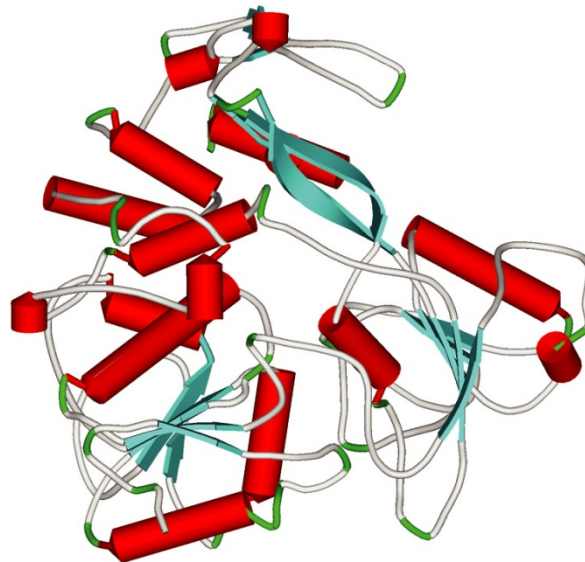
- Two linear and **parallel stretches** (β -strands)
- Strands are bound together by hydrogen bonds
- Can be parallel or anti-parallel (wrt. N/C terminus)



Quelle: Wikipedia

Other Substructures

- α -helices and β -sheets cover 50-80% of most proteins
- Other parts are called **loops or coils**
 - Usually not very important for the structure of the protein
 - But **very important for its function**
 - Often exposed on the surface; important for binding to other molecules



Importance of Secondary Structure Prediction (SSP)

- Secondary structure elements (SSE) are vital for the overall structure of a protein
- Often evolutionary well conserved
- SSE can be used to classify proteins
 - Such classes are highly correlated with function
- SSE gives important clues to protein structure
- SSP much simpler than 3D structure prediction
 - And 3D structure prediction can benefit a lot from a good SSP

Predicting Secondary Structure

- SSP: Given a protein sequence, **assign each AA** in the sequence to one of the **three classes** Helix (H), Strand (E), or Coil (-)

```
KVYGRCELAAAMKRLGLDNYRGYSLGNWVCAAKFESNFNTHATNRNTD
GSTDYGILQINSRWWCNDGRTPGSKNLCNIPCSALLSSDITASVNCAK
KIASGGNGMNAWVAWRNRCKGTDVHAWIRGCRL
```



```
KVYGRCELAAAMKRLGLDNYRGYSLGNWVCAAKFESNFNTHATNRNTD
-----HHHHHHHHHH-----EEEE-----HHHHHHHH--
GSTDYGILQINSRWWCNDGRTPGSKNLCNIPCSALLSSDITASVNCAK
----EEEEEEEEEEEEEEEEEEEE-----HHHHHH
KIASGGNGMNAWVAWRNRCKGTDVHAWIRGCRL
HHH-----EEE-----EEEE----
```

Classification

- **Classification**: Classify each AA into one of three classes
- Classification is a **fundamental problem**
 - Classify the readout of a microarray as diseased / healthy
 - Classify a subsequence of a genome as coding / non-coding
 - Classify an email as spam / no spam
- Many **different techniques**: Naïve Bayes, Regression, Decision Trees, SVMs, Neural Networks, ...
 - **Classification function** learned from properties of known objects
 - Often use same representation (feature vectors) of objects – methods exchangeable
- The following is a rather unsystematic approach
 - But simple to explain and classical for this application

This Lecture

- Introduction
- Predicting Protein Secondary Structure
 - Secondary structure elements
 - Chou-Fasman
 - Other methods

Chou-Fasman Algorithm

Chou & Fasman (1974). Prediction of protein conformation. Biochemistry 13

- Observation: Different AA favor different folds
 - Different AA are more or less often in H, E, C
 - Different AA are more or less often within, starting, or ending a stretch of H, E, C
- Chou-Fasman algorithm (rough idea)
 - Classifies each AA into E or H; unclassified AA are assigned C
 - Compute a score for the probability of any AA to be E (H)
 - Basis: Relative frequencies in a set of sequences with known SSE
 - In principle, assigns each AA its most frequent class
 - Add several heuristic constraints like minimal length of stretches or tendency to build ungapped sequences
 - Rather CCEEEEECC, not CCECEECCECE

Details [sketch, some heuristics omitted]

- Let $f_{j,k}$ be the relative frequency of observing AA j in class k
- Let f_k be the average over all 20 $f_{j,k}$ values
- Compute the propensity $P_{j,k}$ of AA j to be part of class k as

$$P_{j,k} = f_{j,k} / f_k$$

- Using $P_{j,k}$, classify each AA j for every class k into
 - Strong, normal, weak builder ($H_\alpha, h_\alpha, I_\alpha, H_\beta, h_\beta, I_\beta$)
 - Tendency to build a SS-element
 - Strong, weak breaker ($B_\alpha, b_\alpha, B_\beta, b_\beta$)
 - Tendency to stop a SS-element
 - Indifferent (i_α, i_β)

Concrete Values

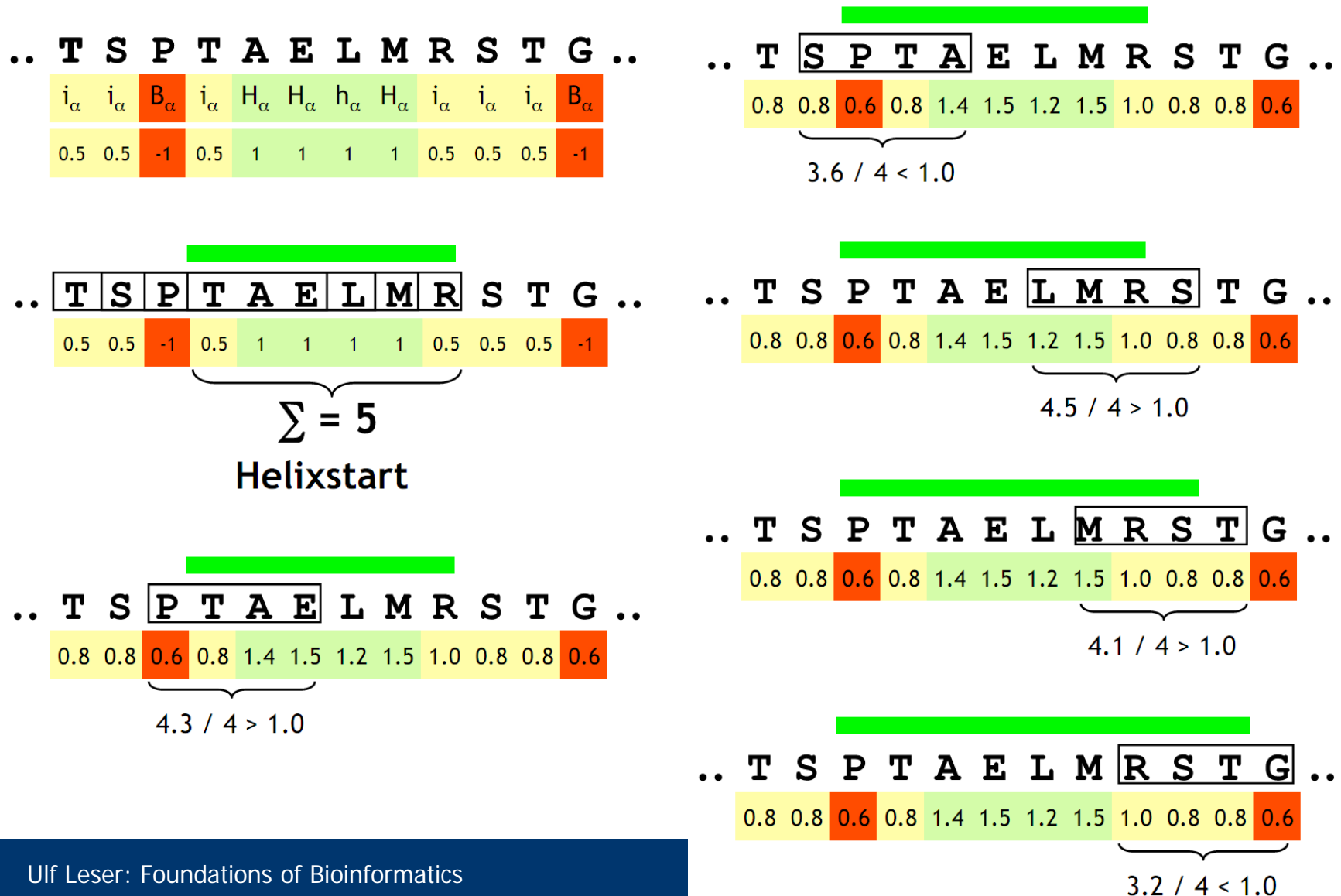
- Originally computed on only 15 proteins (1974)

AS	P_α	Klasse	AS	P_β	Klasse	AS	P_α	Klasse	AS	P_β	Klasse
Glu	1.53	H_α	Met	1.67	H_β	Ile	1.00	I_α	Ala	0.93	I_β
Ala	1.45		Val	1.65		Asp	0.98	i_α	Arg	0.90	i_β
Leu	1.34		Ile	1.60		Thr	0.82		Gly	0.81	
His	1.24	h_α	Cys	1.30	h_β	Ser	0.79		Asp	0.80	
Met	1.20		Tyr	1.29		Arg	0.79		Lys	0.74	b_β
Gln	1.17		Phe	1.28		Cys	0.77		Ser	0.72	
Trp	1.14		Gln	1.23		Asn	0.73	b_α	His	0.71	
Val	1.14		Leu	1.22		Tyr	0.61		Asn	0.65	
Phe	1.12		Thr	1.20		Pro	0.59	B_α	Pro	0.62	
Lys	1.07	I_α	Trp	1.19		Gly	0.53		Glu	0.26	B_β

Algorithm for Helices

- Go through the protein sequence
- Score each AA with 1 (H_α, h_α), 0.5 (I_α, i_α), or -1 (B_α, b_α)
- Find helix cores: subsequences of length 6 with an aggregated AA score ≥ 4
- Starting from the middle of each core, shift a window of length 4 to the left (then to the right)
 - Compute aggregated score A using values $P_{j,k}$ inside the window
 - Not the coarse-grained H/I/B/I values
 - If $A \geq 4$, continue; otherwise stop
- Similar method for strands
- Conflicts (regions assigned both H and E) are resolved based on aggregated scores

Example [Source: O. Kohlbacher, "Strukturvorhersage"]



Performance

- Accuracy app. 50-60%
 - Measured on per-AA correctness
- Prediction is **more accurate in helices** than in strands
 - Because helices build **local bridges** (hydrogen bounds between the turns; each AA binds to the +4 AA)
- General problem
 - Secondary structure is not only a local problem
 - Looking **only at single AAs** is not enough
 - Note: Scores are based on individual AA; aggregation by summation assumes **statistical independence** of pairs, triples ... in a class
- One needs to include the **context of an AA**

This Lecture

- Introduction
- Predicting Protein Secondary Structure
 - Secondary structure elements
 - Chou-Fasman
 - Other methods

Classes of Methods

- First generation: Properties of single AA only
 - Accuracy: 50-60%, e.g. Chou-Fasman (1974)
- Second generation: Include info. **about neighborhood**
 - Accuracy: ~65%, e.g. GOR (1974 – 1987)
- Third generation: Include info. from **homologous seq's**
 - Accuracy: ~70-75%, w.g. PHD (1994)
- Forth generation: Build **ensembles** of good methods
 - Accuracy: ~80%, e.g. Jpred (1998)
- Current performance
 - Jpred 4 (2015): 82% overall, ~90% for certain other properties
 - Spine-X (2012): 84% overall

Further Reading

- Gerhard Steger (2003). "Bioinformatik – Methoden zur Vorhersage von RNA- und Proteinstrukturen", Birkhäuser, chapter 8,10,11,13
- Zvelebil, M. and Baum, J. O. (2008). "Understanding Bioinformatics", Garland Science, Taylor & Francis Group, chapter 2, 11, 12 (partly)