



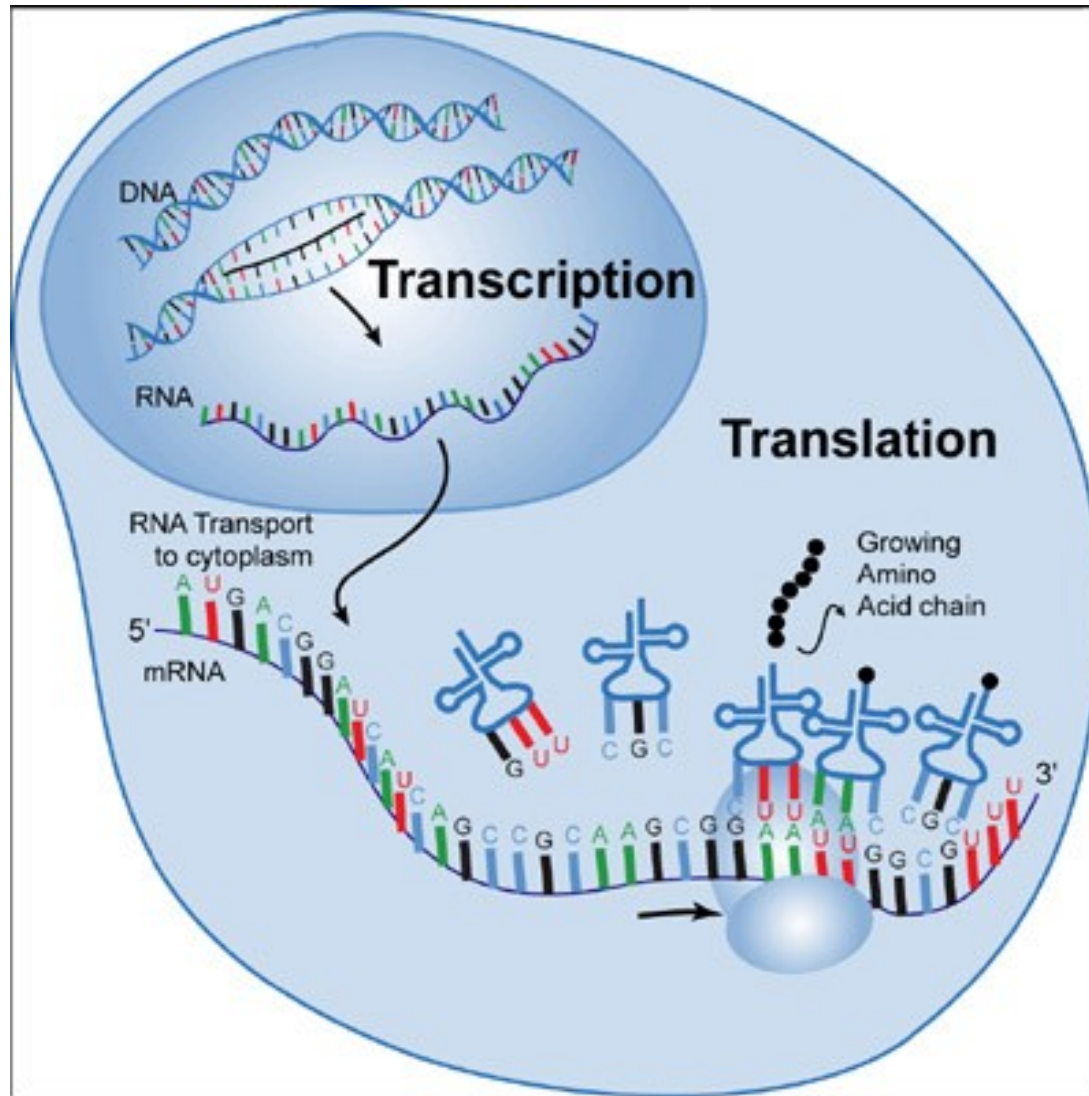
Proteins: Structure & Function

Johannes Starlinger

This Lecture

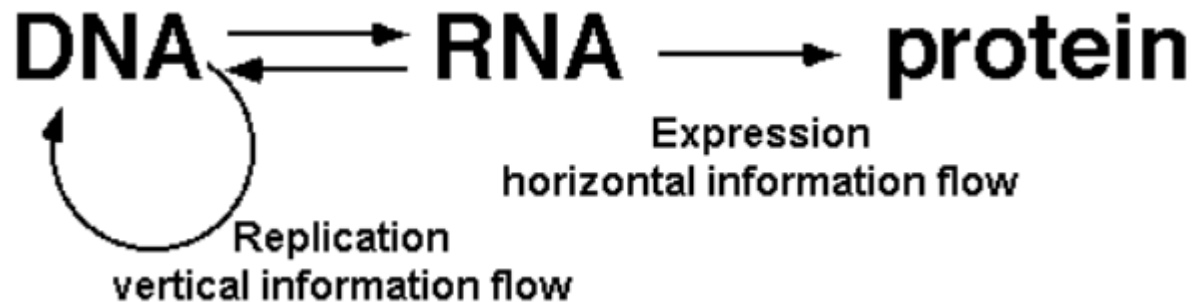
- Proteins
 - Structure
 - Function
 - Databases
- Predicting Protein Secondary Structure
- Many figures from Zvelebil, M. and Baum, J. O. (2008). "Understanding Bioinformatics", Garland Science, Taylor & Francis Group.
- Examples often from O. Kohlbacher, Vorlesung Strukturvorhersage, WS 2004/2005,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

DNA→**Transcription**→RNA→**Translation**→Protein



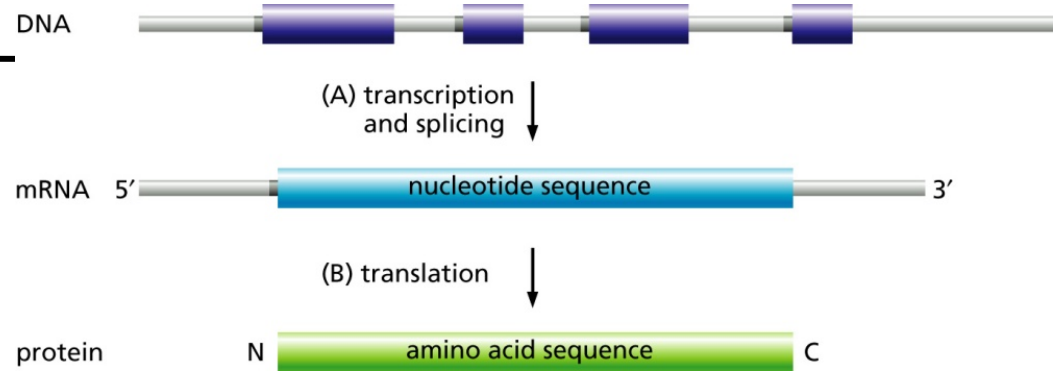
From:
<http://www.tokresource.org>

Central Dogma of Molecular Biology



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

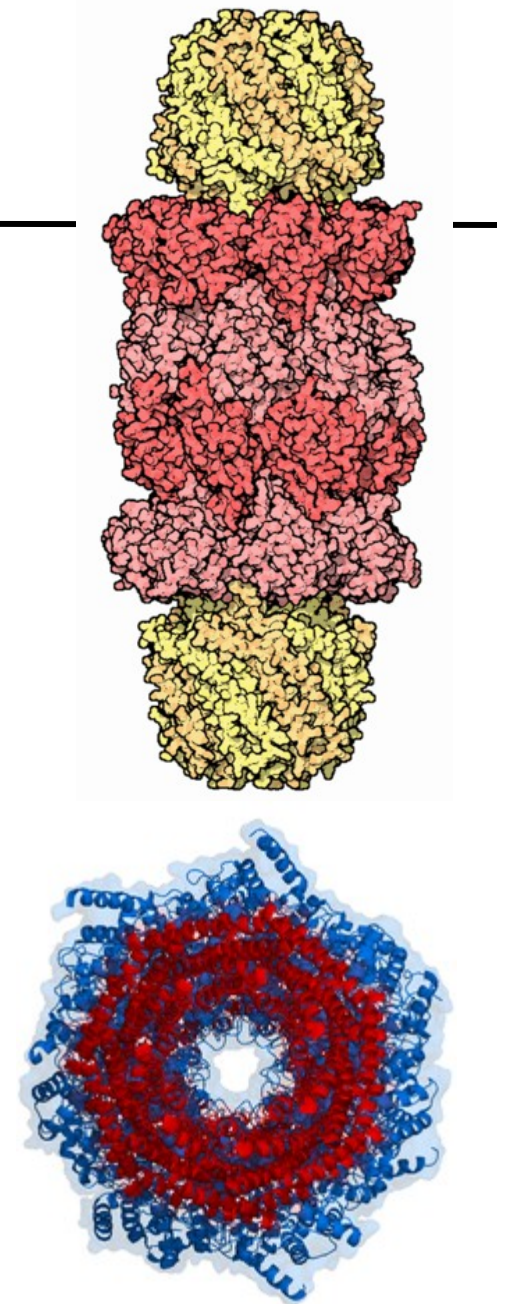
Details



- Alternative Splicing
 - “One gene – one protein” is wrong
 - Exons may be spliced out from the mRNA
 - Human: ~6 times more different proteins than genes
- Post-translational modifications
 - (De-)Phosphorylation, glycolysation, cleavage of signals, ...
 - Rough estimates: 100K proteins, 500K protein forms
- Complexes: Proteins physically group together to perform specific function

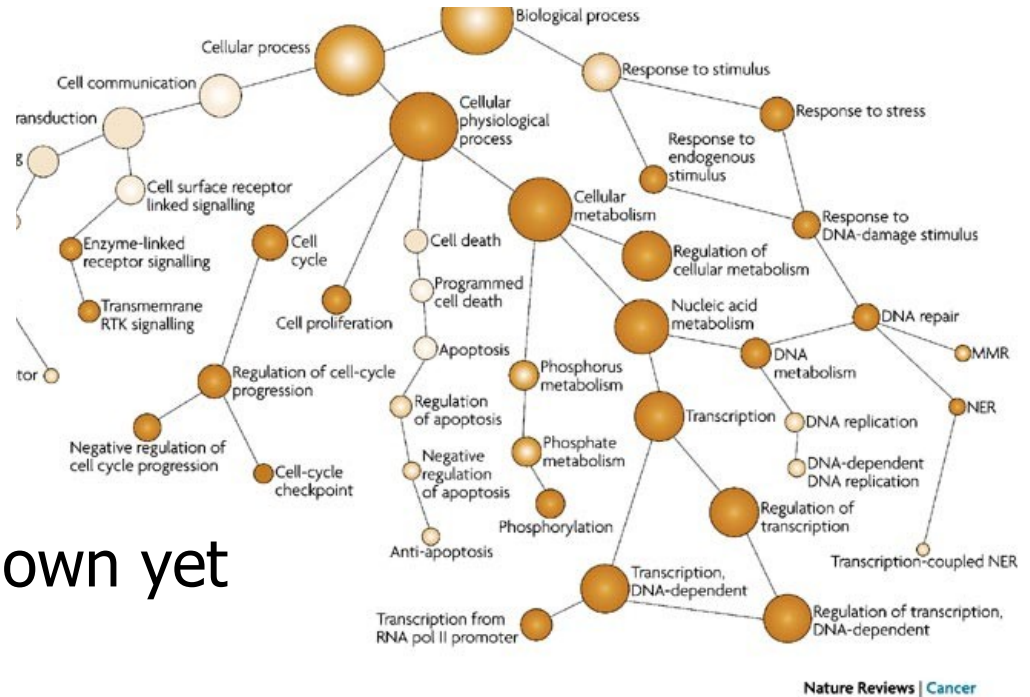
Example: Proteasome

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

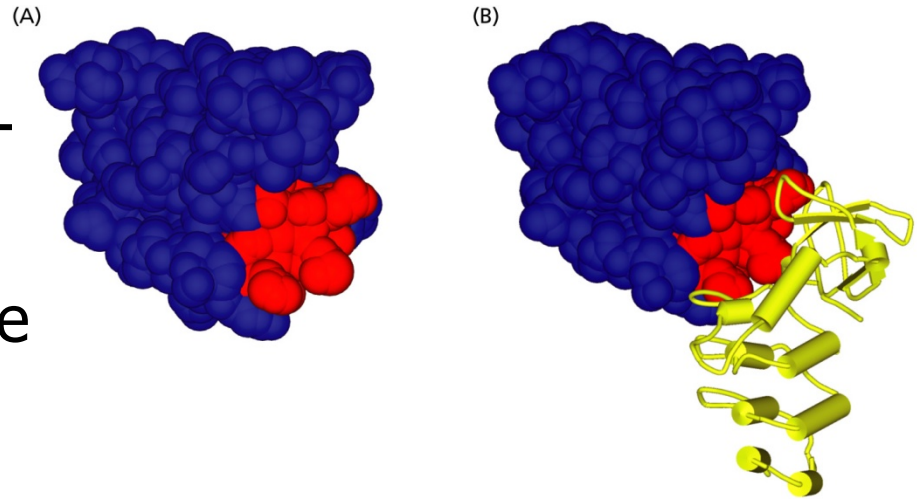


Protein Function

- Proteins perform essentially everything that makes an organism alive
 - Metabolism
 - Signal processing
 - Gene regulation
 - Cell cycle
 - ...
 - For $\sim 1/3$ of all human genes, no function is known yet
 - Describing function
 - Gene Ontology: 3 branches, >30.000 concepts
 - Cellular component, biological process, molecular function
 - Used world-wide to describe gene/protein function
-



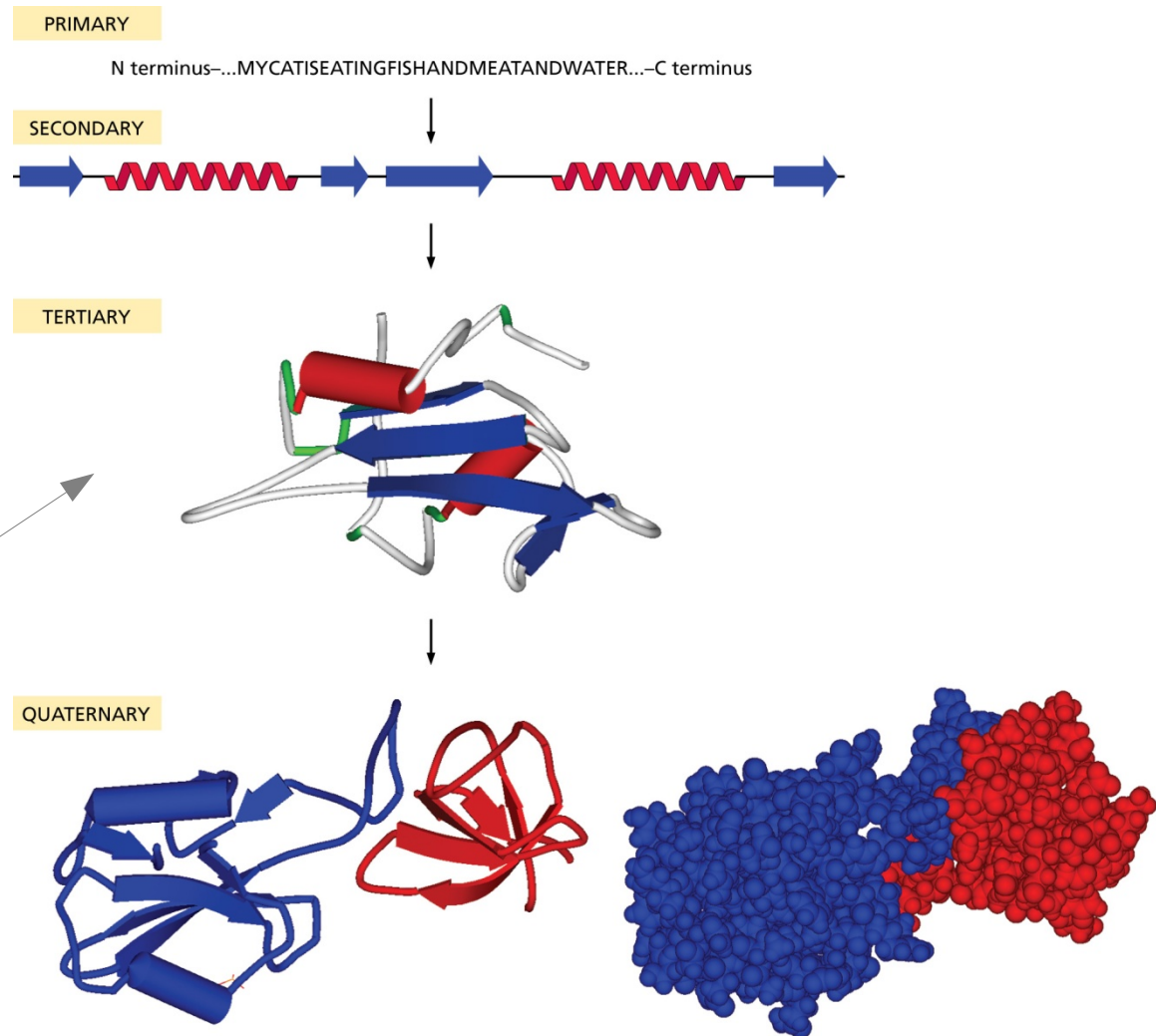
Function and Motifs



- Proteins often have multiple functions
 - Avg. no. of GO terms assigned to a human protein: ~6
- 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**

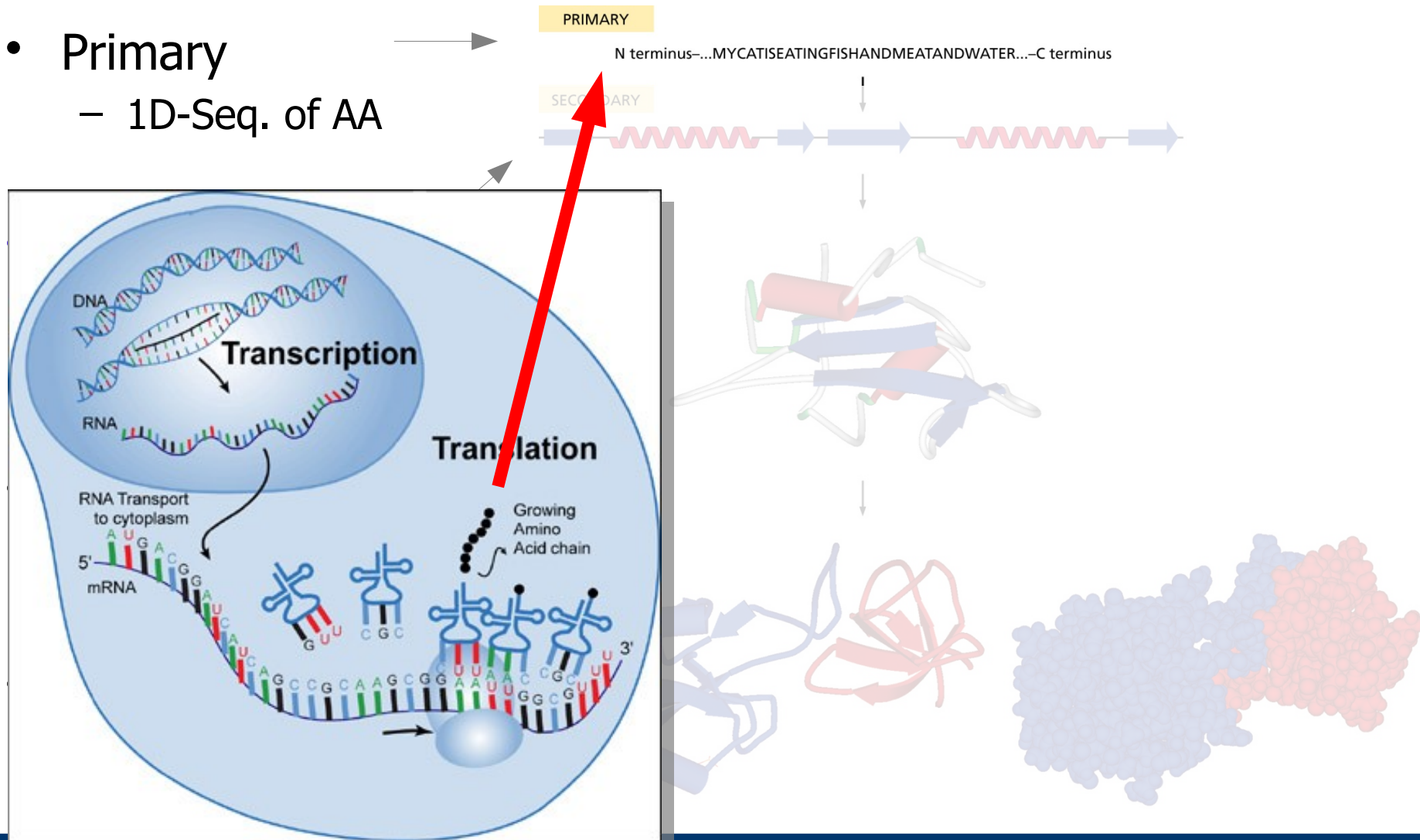
Protein Structure

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

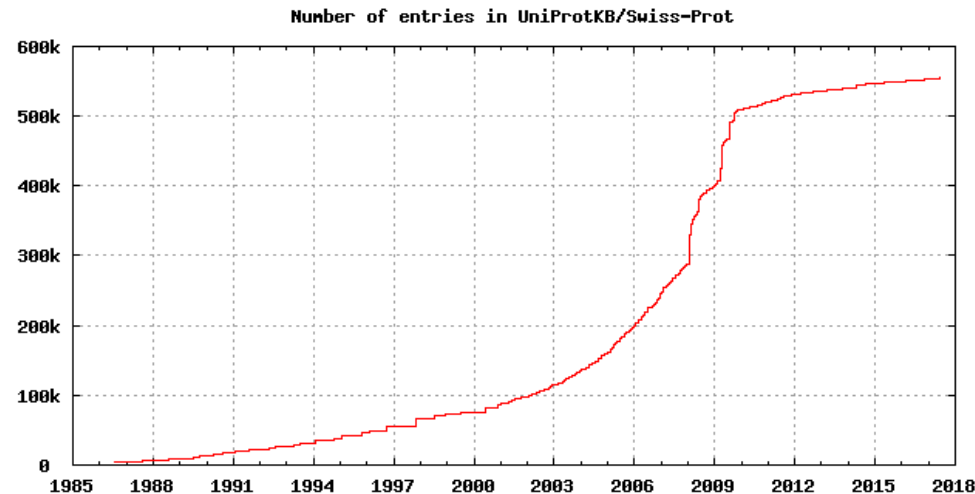


Protein Structure

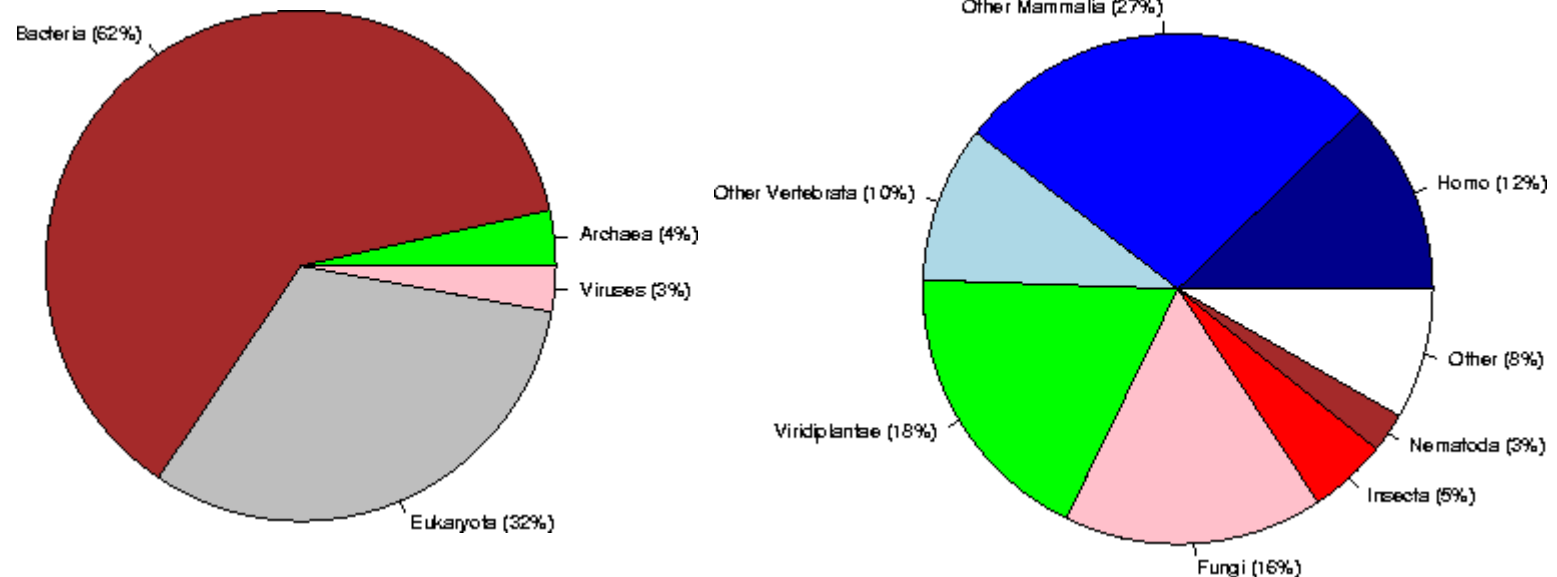
- Primary
 - 1D-Seq. of AA



- “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



UniProt: Species [\[http://www.expasy.org/sprot/relnotes/relstat.html\]](http://www.expasy.org/sprot/relnotes/relstat.html)

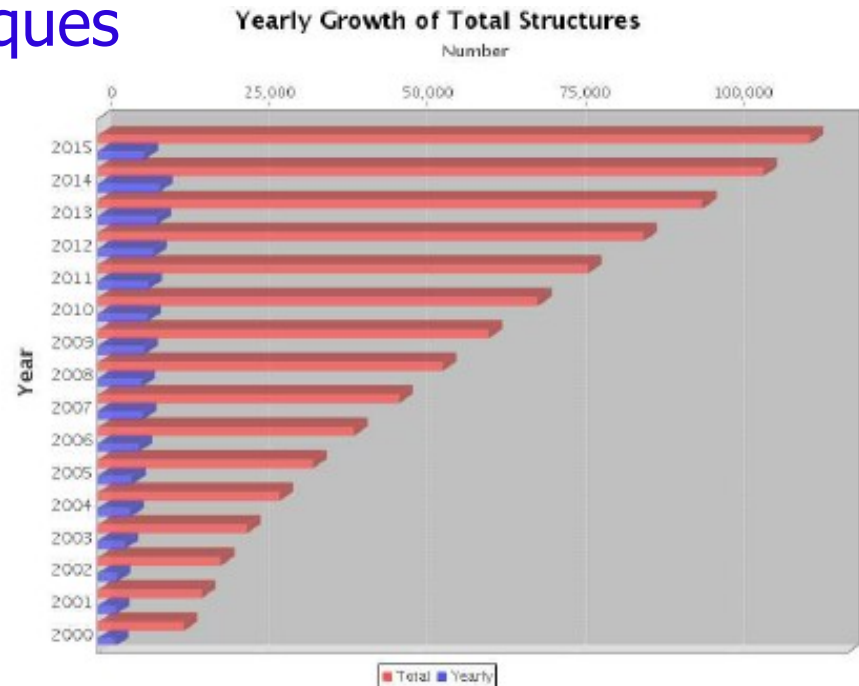


Number	Frequency	Species
1	20205	Homo sapiens (Human)
2	16894	Mus musculus (Mouse)
3	15383	Arabidopsis thaliana (Mouse-ear cress)
4	7991	Rattus norvegicus (Rat)
5	6721	Saccharomyces cerevisiae (Baker's yeast)
6	5999	Bos taurus (Bovine)

...

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

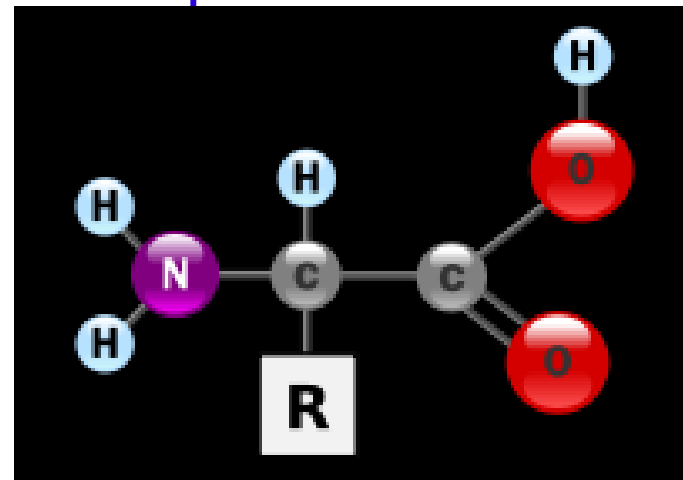


This Lecture

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

Amino Acids

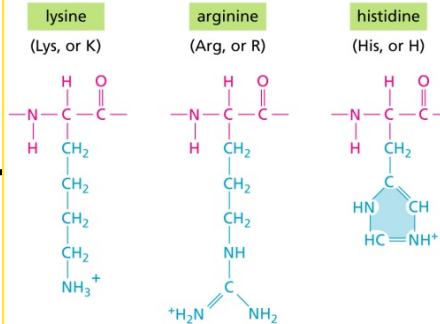
- AA consist of a **common core** and a **specific residue**
 - Amino group – NH_2
 - Central C_α - Carbon – CH
 - Carboxyl group – COOH



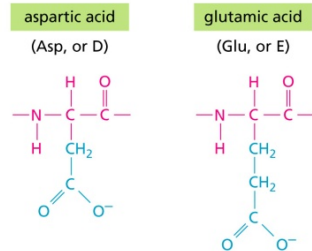
- Residues (side chains) vary greatly between AA
- Residues determine the **specific properties** of a AA

Side Chains

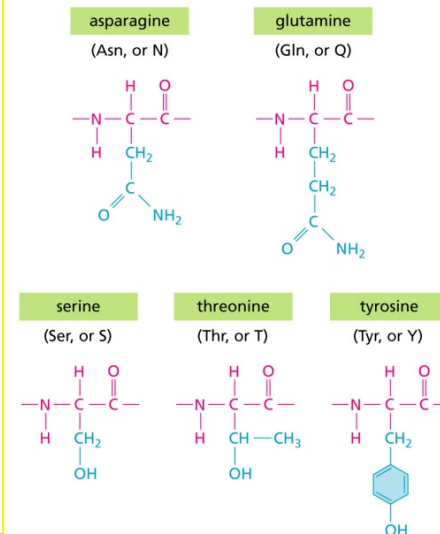
BASIC SIDE CHAINS



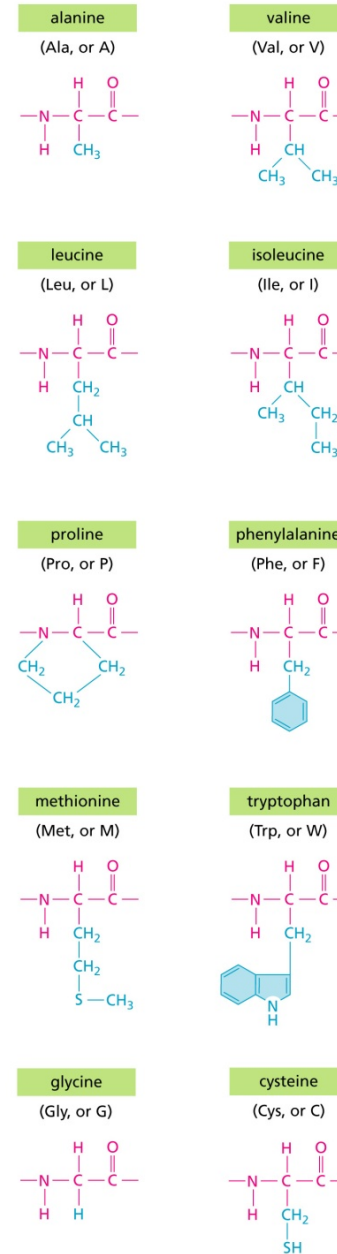
ACIDIC SIDE CHAINS



UNCHARGED POLAR SIDE CHAINS



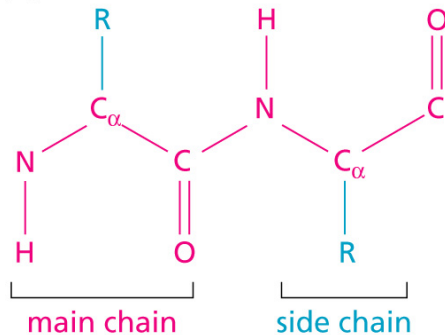
NONPOLAR SIDE CHAINS



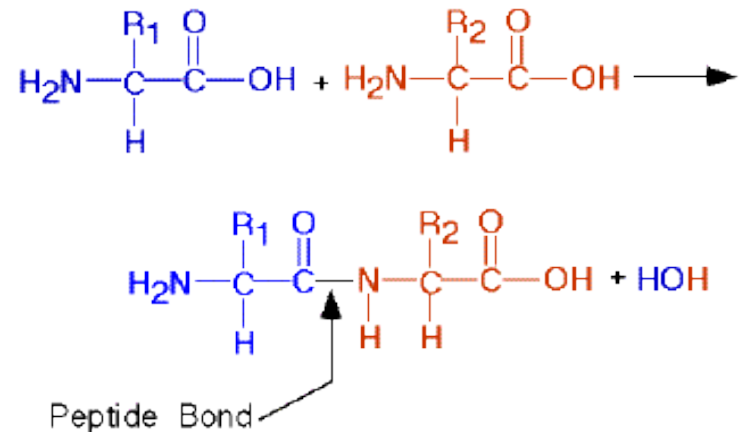
Structure of a Protein

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

(A)



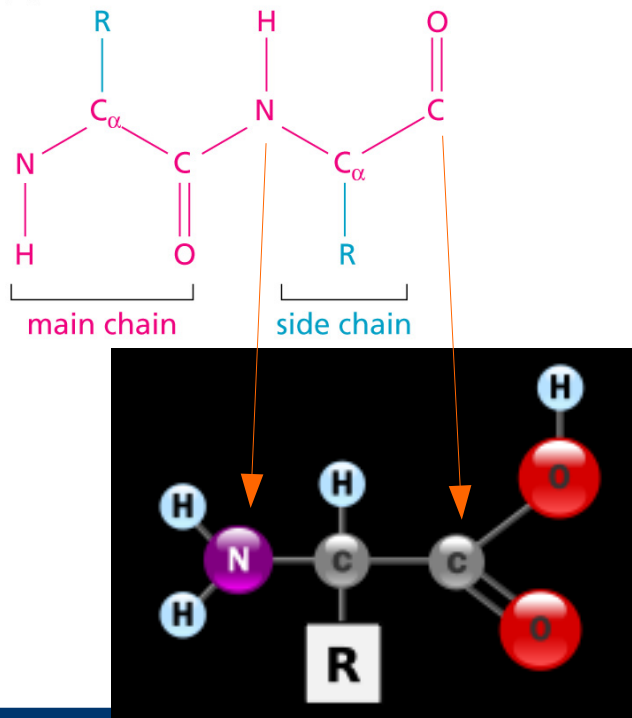
Peptide Bond Formation



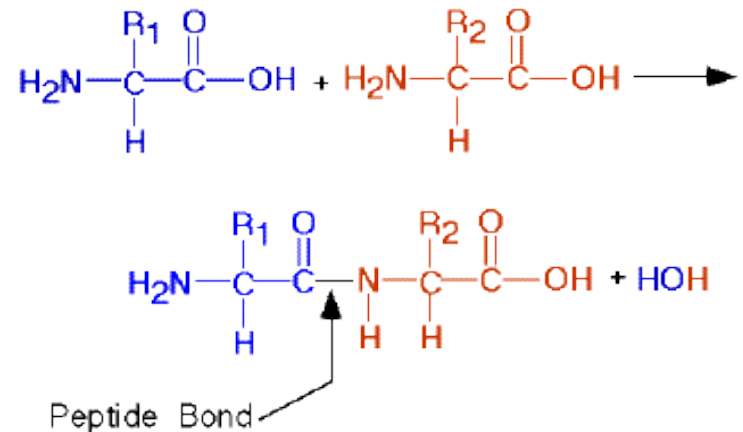
Structure of a Protein

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

(A)

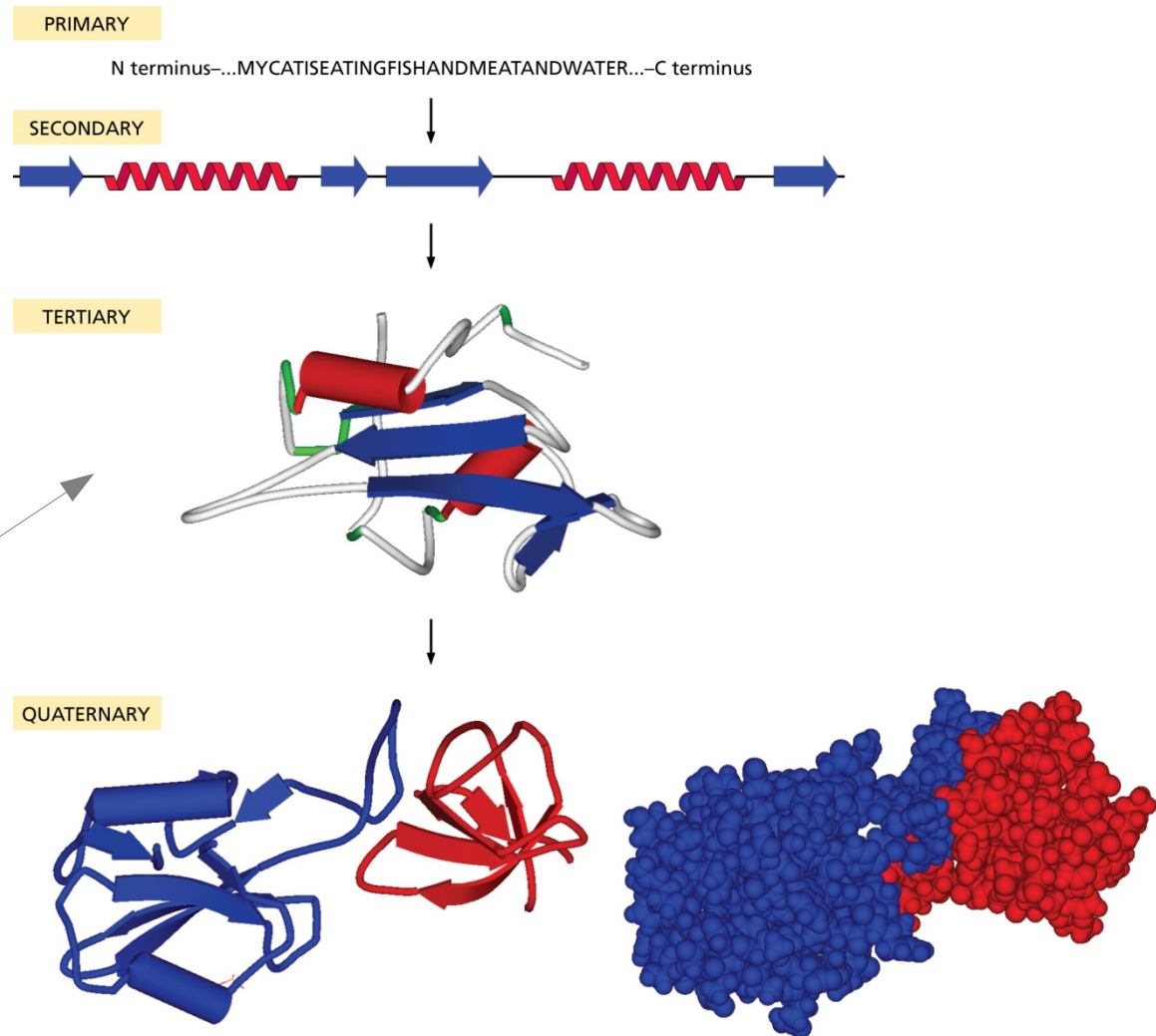


Peptide Bond Formation



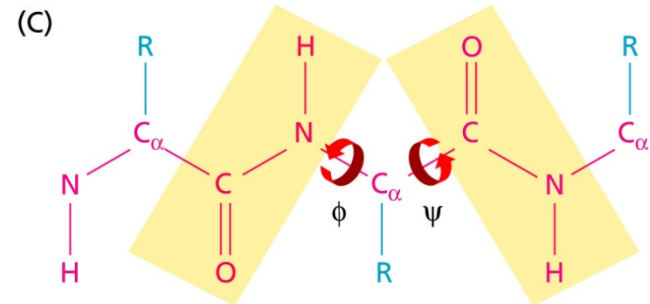
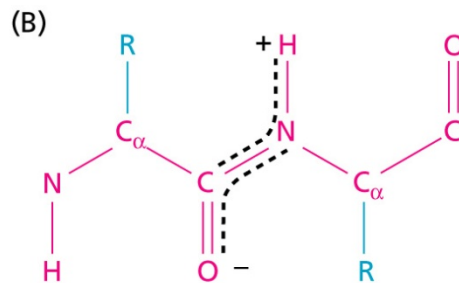
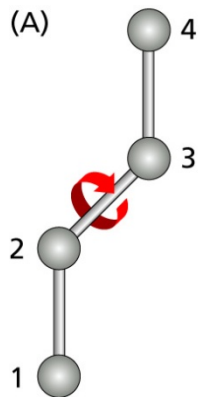
Protein Structure

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



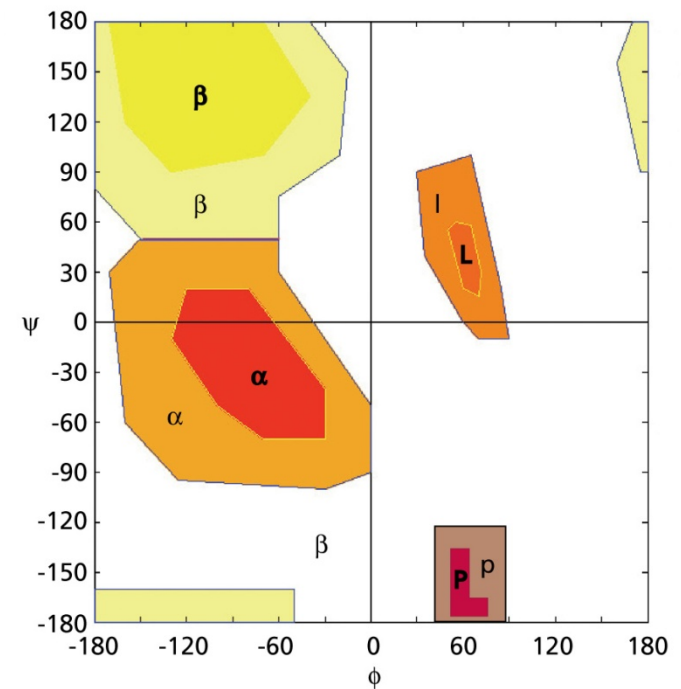
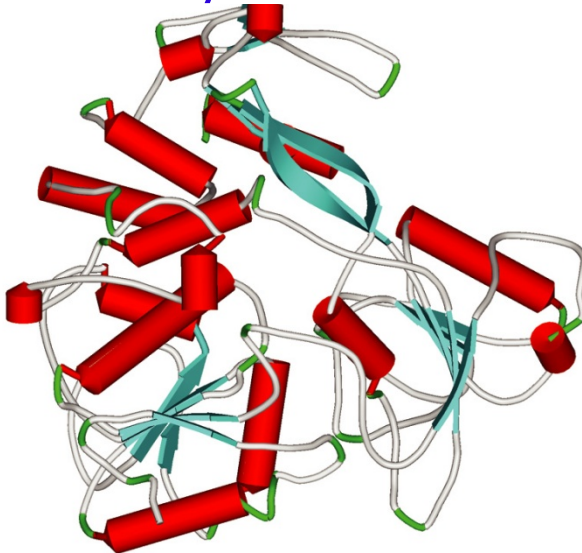
Flexibility

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

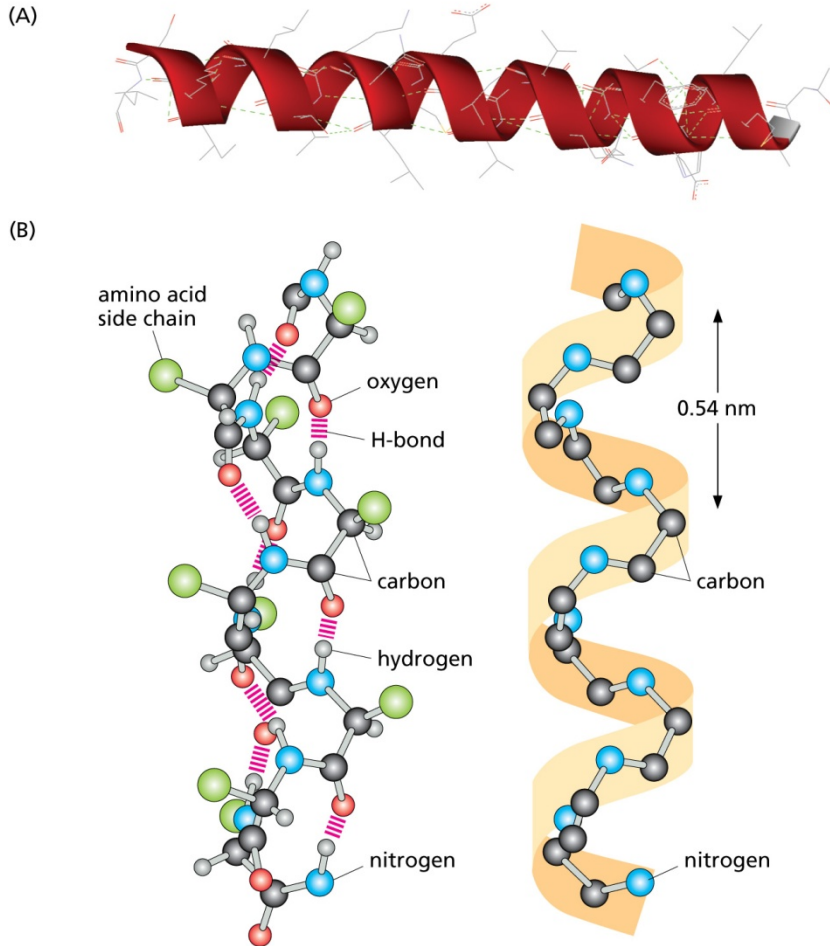


Ramachandran Plots

- 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**

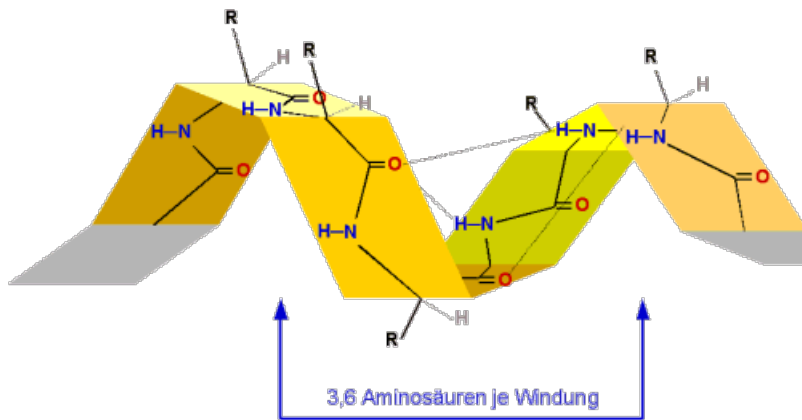


α -Helix



- 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

α -Helix

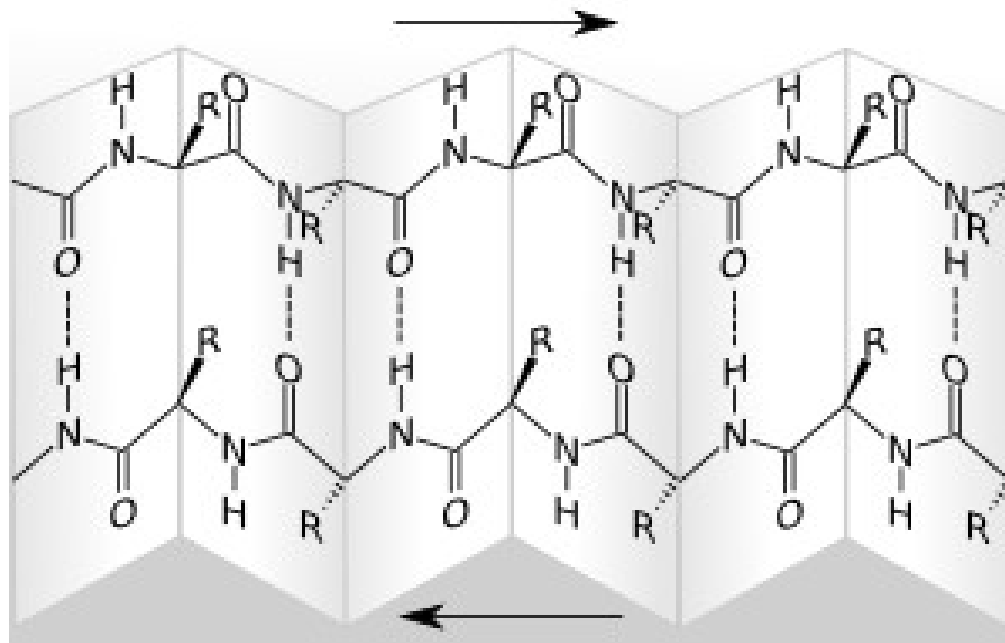


From:
www2.chemie.uni-erlangen.de

- 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

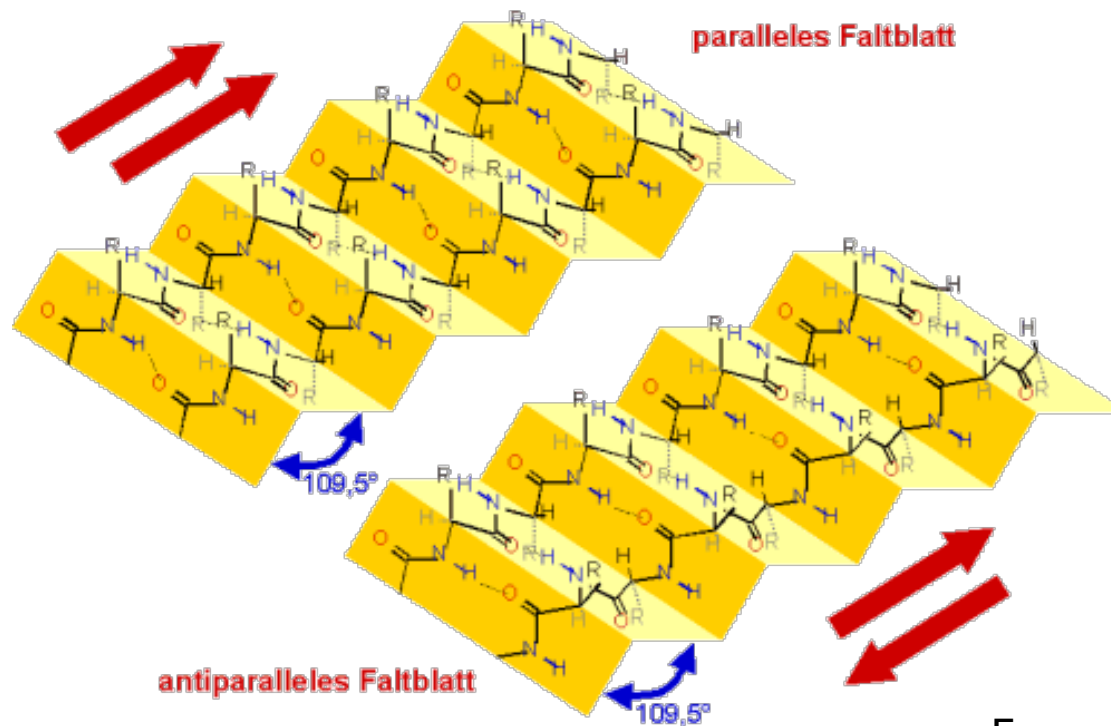
- Two linear and **parallel stretches** (β -strands)
- Strands are bound together by hydrogen bonds
- Can be parallel or anti-parallel (wrt. N/C terminus)



Quelle: Wikipedia

β -Sheet

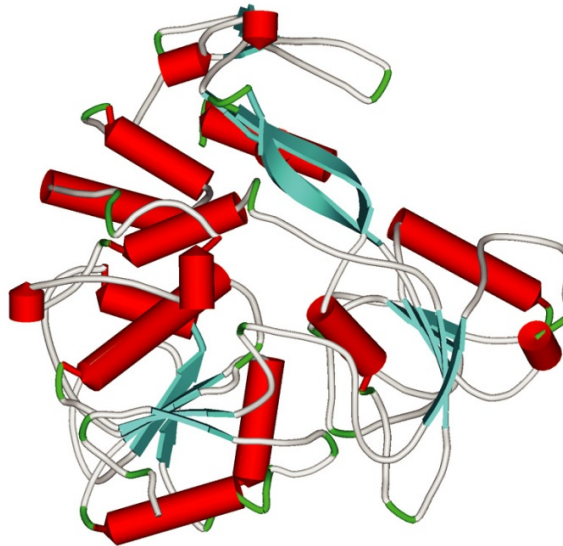
- Two linear and **parallel stretches** (β -strands)
- Strands are bound together by hydrogen bonds
- Can be parallel or anti-parallel (wrt. N/C terminus)



From: chemgapedia.de

Other Substructures

- α -helixes and β -sheets form around 50-80% of a protein
- 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

- Secondary structure elements (SSE) are vital for the overall structure of a protein
- Often **evolutionarily well conserved**
- SSE can be used to classify proteins
 - Such classes are highly associated to function
- Knowing the SSE gives important **clues to protein function**
- Secondary structure prediction (SSP) is **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
-----HHHHHHHH-----EEEE-----
GSTDYGILQINSRWWCNDGRTPGSKNLCNIPCSALLSSDITASVNCAK
----EEEE-----HHHHHH
KIASGGNGMNAWVAWRNRCKGTDVHAWIRGCRL
HHH-----EEE-----

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, ...
 - Based on same principles can be exchanged easily
 - **Classification function** learned from properties of known objects
- 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 constraints about **minimal length** of E (H) stretches
 - Several further heuristics

Some 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 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 ($H_\alpha, h_\alpha, I_\alpha$)
 - Strong, weak breaker (B_α, b_α)
 - Indifferent (i_α)

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	.53	H_α	Met	.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	.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)
 - based on known values from propensity table
- 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
 - If $A \geq 4$, continue; otherwise stop
- Similar method for strands
- **Conflicts** (regions assigned both H and E) are resolved based on aggregated scores

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%
 - E.g. PHD (1994)
- Forth generation: Build ensembles of good methods
 - Accuracy: ~80%
 - E.g. Jpred (1998)

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)