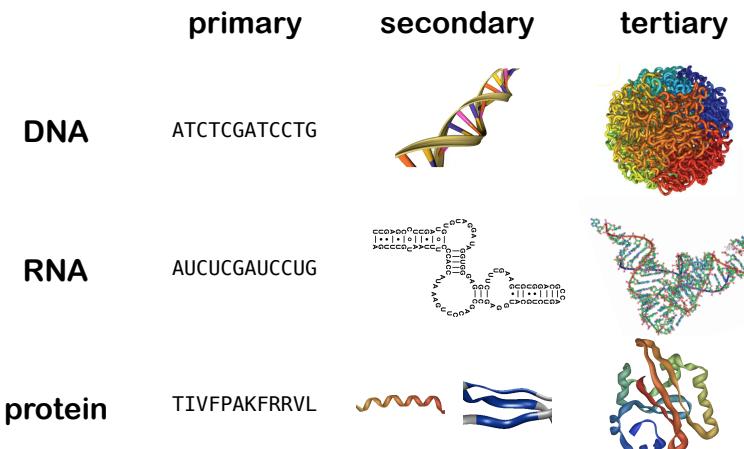


Bioinformatik
Wintersemester 2019 / 2020, Uni Potsdam
Stefanie Hartmann

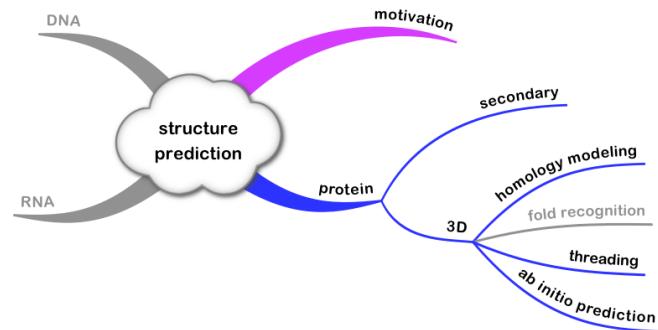
Protein structure prediction Dec 20, 2019

Sequence, structure, function



bioinf WS19 • lec10 • S.Hartmann

Overview



bioinf WS19 • lec10 • S.Hartmann

Determination of sequence & structure

DNA sequence

- Sanger sequencing, NGS technologies
- genomes, protein-coding genes, non-coding RNAs

protein sequence

- mass spectrometry, Edman degradation, ...
- computational translation of DNA sequence

protein 3D structure

- mostly X-ray crystallography, NMR spectroscopy, cryo-electron microscopy

bioinf WS19 • lec10 • S.Hartmann

Determined sequences & structures

NCBI	SRA	3.6	Petabytes
	GenBank	210,000,000	sequences
	WGS	722,000,000	sequences
UniProt	SwissProt	560,000	sequences
	TrEMBL	137,000,000	sequences
PDB	<i>Homo sapiens</i>	50,000	structures
	<i>Escherichia coli</i>	8,200	structures
	<i>Mus musculus</i>	7,000	structures
	<i>S. cerevisiae</i>	4,000	structures
	...		
	Total	140,000	structures

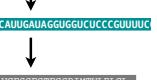
bioinf WS19 • lec10 • S.Hartmann

Prediction of structure

- many more sequences than structures available
 - experimental determination of structure:
difficult, time-consuming, expensive
- early bioinformatics endeavors
 - infer a structure based on sequence information
 - intrinsic information
 - comparative approaches
- the only practical solution?

bioinf WS19 • lec10 • S.Hartmann

Protein: sequence, structure, function

CTTAAGGAAATGCTATTGAGGTGTTTCCAGCGCAATTGCTGCTTACCAAAAGCTTCAGGTGATATCTTATCACATACTCTGTGGTTCTTGCCAGCTTATGTTGCAATGTTGATTATCCTCTT

 ↓
 GAGAUAGGAAACAAUCAUUGAUAGGUUCUCCGUUU
 ↓
 TIVFFPAKFRRLVLFSSSFPSGRIMTWLCL

primary structure (sequence)
 • linear sequence of amino acids

secondary structure
 • hydrogen bonds between backbone NH and CO groups
 • α -helix, β -strand (sheet), loop (coil, turn)



tertiary structure
 • folded secondary structure elements, compact & organized



quaternary structure

bioinf WS19 • lec10 • S.Hartmann

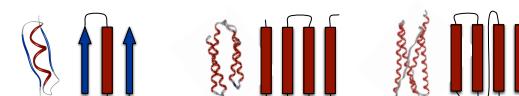
Protein: sequence, structure, function

domain

- functional / structural element (e.g., Calcium-binding domain)

super-secondary structure

- specific geometric arrangements & combinations of α -helices and β -sheets



fold

- general arrangement and connectivity of secondary structure elements in a protein

bioinf WS19 • lec10 • S.Hartmann

example: CATH

- protein structures from PDB: split up into domains

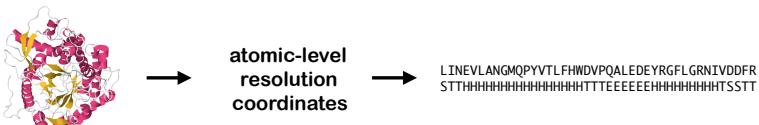
C	class	architectures with similar content of secondary structure (e.g., mainly alpha or beta, mixed, etc)
A	architecture	topologies that share a roughly similar spatial arrangement of secondary structures
H	homologous superfamily	domains that share a clear common ancestor
T	topology	homologous superfamilies that share the same fold; no clear evidence for evolutionary relationship

bioinf WS19 • lec10 • S.Hartmann

Protein secondary structure

1. assign secondary structure from a given 3D structure

- use 3D coordinates for automatic assignment
 - most commonly used approaches are based on hydrogen bonding patterns: DSSP, STRIDE



2. predict secondary structure from a primary sequence

- different approaches & programs exist!



bioinf WS19 · lec10 · S.Hartmann

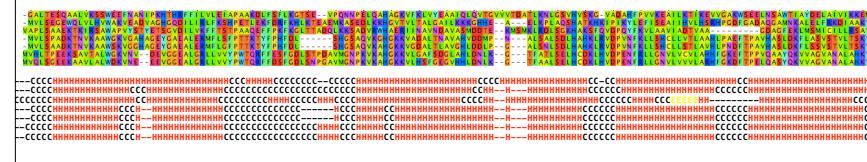
Protein secondary structure

tertiary (or quaternary) structure

- function!

secondary structure

- simpler description of structure
 - important for classifying proteins
 - useful for comparing proteins / protein structure
 - important for predicting 3D structures



bioinf WS19 • lec10 • S.Hartmann

Protein secondary structure prediction

- analyze & learn from known protein structures
 - more data: better prediction results

First generation methods

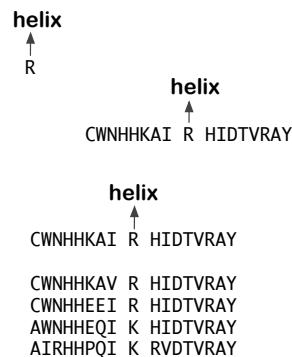
- statistical methods
 - single-residue propensities

Second generation methods

- statistical methods
 - consider adjacent residues

Third generation methods

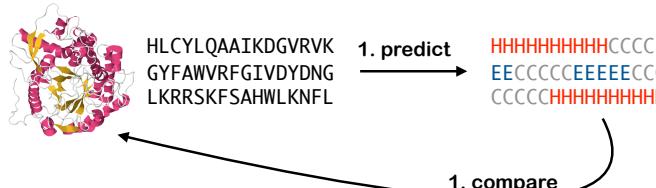
- use evolutionary information
 - use neural networks, HMMs



bioinf WS19 : lec10 : S.Hartmann

Prediction accuracy

- unknown for proteins without a solved 3D structure!
- evaluate accuracy by using protein sequences for which a 3D structure is available



Prediction accuracy

HLCYLQAAIKDGVRVKGYFAWVRFGIVDYDNGLKRRSKFSAHWLKNFL
pred1 HHHHHHHHHHCCCCEEECCCCCCCCCCCCHHHHHHHHHHHHH 9/50
pred2 HHHHHHHHHHCCCCHHHHCCCCCCCCCCCCCCCCHHHHHHHHHHH 9/50

measures of accuracy

- overall accuracy (Q_3)
- accuracy for each state
- overlap between observed vs. predicted structures (SOV)

current best predictors: ~82% Q_3 accuracy

bioinf WS19 • lec10 • S.Hartmann

3D protein structure

- helps to understand the protein's function
- helps to design site-directed mutations & drugs
- can explain binding specificity, antigenic property
- ...

experimental determination of 3D structure

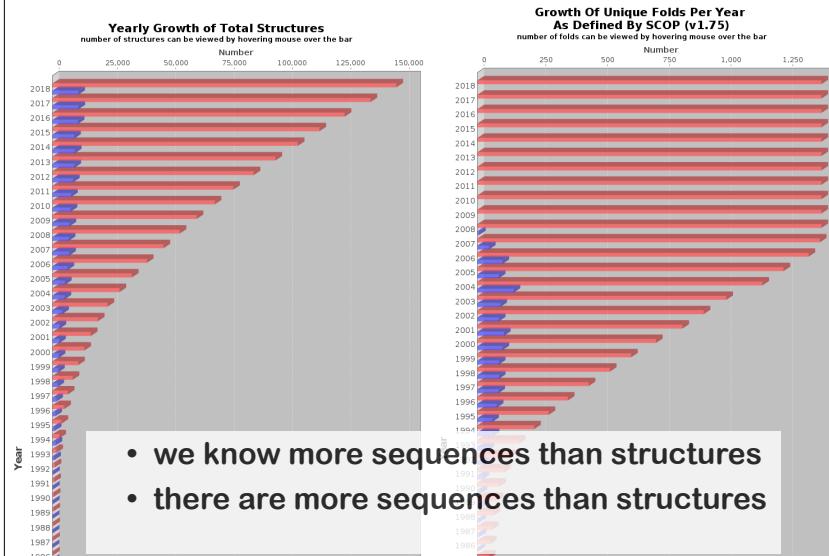
- expensive! time consuming!
- doesn't always work

computational prediction of 3D structure

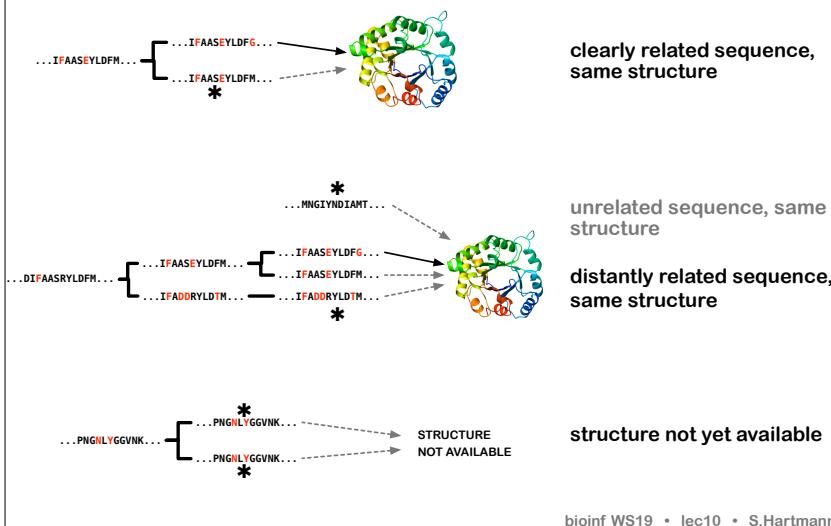
- can provide useful results
- doesn't always work

bioinf WS19 • lec10 • S.Hartmann

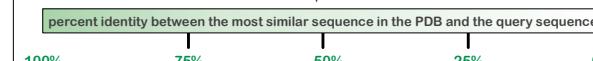
PDB content growth



Approaches to structure prediction



protein sequence; 3D structure unknown



detection of sequence homology

no (easy) detection of sequence homology

no (detectable) sequence homology, similar structures

no sequence homology, no similar structures

homology modeling

fold recognition / threading

ab initio approach



3D protein model

3D structure prediction

is a sufficiently similar sequence found in PDB?

- YES: use homology (comparative) modeling
- model unknown structure based on known structure



assumption

- different but homologous sequences fold into the same or highly similar structures

Homology modeling

1. Template recognition

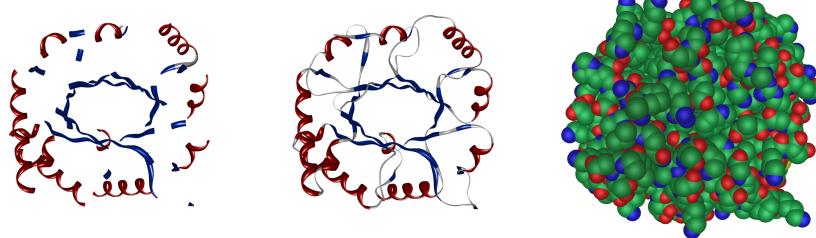
- use BLAST to compare query sequence to all sequences with a structure in PDB
- identify homologous, highly similar sequence(s)

2. Alignment

- optimal alignment of query sequence and homologous sequences with structures
- structure-aware alignment

seq	...VPLVKHLADLSKS K TSPYVL P V P FLNV L NGGSHAG G ALALQEFMIAPTGA-KTFEAELRI...
1tim	...QEVHEKLRGLKTHVSDAVAV--QSRIIYGGSVTGGNCVDGFLVGGASLKPEFVDIINA...
P00925	...VPLVQH L ADLSKS K TSPYVL P V P FLNV L NGGSHAG G ALALQEFMIAPTGA-KTFEAEMRI...
P30575	...IPLYKHIANISNAKKGFVLP P V P QNV L NGGSHAG G ALAFQEFMIAPTGV-STFSEALRI...

Homology modeling



3. Backbone generation

- copy coordinates from template for conserved regions

4. Loop building

- based on existing structures or modeled

5. Side chain modeling

bioinf WS19 • lec10 • S.Hartmann

Homology modeling

6. Model refinement

- side chains, backbone, energy calculations, etc

7. Model verification

- check bond angles, torsion angles, bond lengths
- check for distribution of polar/apolar residues

8. Model evaluation

```
NAMOLGSLCAMLLIGFALRNTNAVRTDPPSHCPVL  
NRSSEFSLVPGFIFGFTASAAQYQVGAEGEGRGFS  
IWDAYTHHNPERIKDRNSGDIAIDQYHRYKEDVGI  
MKNGLDSRSLISWSRSLLPNGKLSGGVNKEGIEY  
YNNLNLTELLRNKGITPPVTLWVQVQVQVQVQVQVQVQV  
LSPRIVHUYKDYTELCPKFCGRKIHWTIPLPYA  
VSHHGAYAIGHAPGRCSDWEACLGDDSAIEPYLVLT  
HNQLLAHASTVKVVKDQYQASQNGVIGITVSHWI  
EPASKSKEDIDAAASRYLDFMFQWMSPLTIGDYPH  
SMRHLGERLPVQEOKSLLLNGSDFPFIGLNYYSA  
RGSASPLNIVPPCIIALLNDTMTDYYHHUCLQAA  
GIDEFNPKXLSLEELNDTMTDYYHHUCLQAA  
IKGGRVVKYGFANSVLNFENWSGYTVRFGINYWD  
YDNGLKRRSKFSAHWLNFKLNKNSKSKEIRVRVD  
DNARDTKAGYEI
```



Oligo-State	Ligands	GMQE	QMEAN4
MONOMER (matching prediction)	None	0.74	-0.78 kJ
Template	Seq Identity Coverage Description		
4a3y:1:A	57.00%		RAUCAFRICINE-O-BETA-D-GLUCOSIDASE

bioinf WS19 • lec10 • S.Hartmann

How well does it work?

- seq identity limiting factor
- <25% template recognition
- 25-50% alignment
- 50-75% loop modeling
- >75% side chain modeling



→

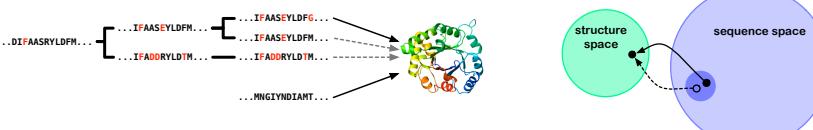
Oligo-State	Ligands	GMQE	QMEAN4
MONOMER (matching prediction)	None	0.74	-0.78 kJ
Template	Seq Identity Coverage Description		
4a3y:1:A	57.00%		RAUCAFRICINE-O-BETA-D-GLUCOSIDASE

bioinf WS19 • lec10 • S.Hartmann

3D structure prediction

is a sufficiently similar sequence found in PDB?

- YES: use homology (comparative) modeling
- NO
 - can I find homologous sequences with more sensitive methods? (sequence-based fold recognition)
 - can I find a similar structure? (threading)

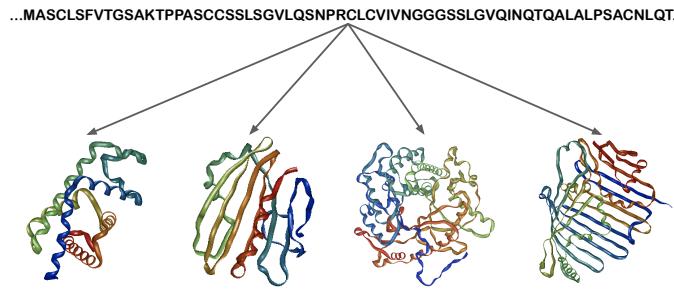


assumption

- structure space is smaller than sequence space
- physical constraints on protein structure

bioinf WS19 • lec10 • S.Hartmann

Structure-based approaches: threading



compute compatibility scores

- free energy of the sequence “in” structure A, B, ...?

bioinf WS19 • lec10 • S.Hartmann

Structure-based approaches

threading

- align (fit) a sequence to a structure
- compute compatibility between sequence & structure
- if a sequence is highly compatible with a structure, it might fold into that structure

for a given structure

- relate atomic coordinates to potential energy
- use mathematical equations & parameters
 - potential energy functions (force field)
 - based on empirical data

bioinf WS19 • lec10 • S.Hartmann

Fold library for threading

fit a sequence to all structures in PDB?

- no!
- one protein with different ligands or mutations
- many proteins with highly similar folds

use a library of non-redundant folds

- CATH topology
- SCOP fold

►	1	Mainly Alpha	5 Architectures, 397 Folds, 907 Superfamilies, 48121 Domains
►	2	Mainly Beta	20 Architectures, 241 Folds, 547 Superfamilies, 5994 Domains
►	3	Alpha/Beta	14 Architectures, 626 Folds, 1158 Superfamilies, 12572 Domains
►	3.10	Roll	58 Folds, 101 Superfamilies, 9748 Domains
►	3.15	Super Roll	3 Folds, 3 Superfamilies, 5 Domains
►	3.20	Alpha-Beta Barrel	18 Folds, 46 Superfamilies, 10515 Domains
►	3.20.10	D-Amino Acid Aminotransferase; Chain A, domain 2	1 Superfamilies, 131 Domains
►	3.20.14	L-Fucose Isomerase; Chain A, domain 3	1 Superfamilies, 15 Domains
►	3.20.16	Serine Protease, Human Cytomegalovirus Protease; Chain A, families	47 Domains
►	3.20.19	Aconitase; domain 4	1 Superfamilies, 38 Domains
►	3.20.20	TIM Barrel	29 Superfamilies, 10050 Domains
►	3.20.70	Anaerobic Ribonucleotide-triphosphate Reductase Large Chain families	35 Domains
►	3.20.80	Multidrug-efflux Transporter 1 Regulator Bmr; Chain A	1 Superfamilies, 28 Domains
►	3.20.90	Tubby Protein; Chain A	1 Superfamilies, 7 Domains

3D structure prediction

is a sufficiently similar sequence found in PDB?

- YES: use homology (comparative) modeling
- NO
 - was fold recognition or threading successful?
 - NO: use ab initio prediction



assumption & approach

- the native fold of a protein is the one with the lowest free energy
- (in theory) evaluate free energy of all possible conformations for a given sequence

bioinf WS19 • lec10 • S.Hartmann

Ab initio structure prediction

thermodynamic hypothesis

- proteins fold by following physical laws
- native conformation: global free energy minimum

minimizing free energy

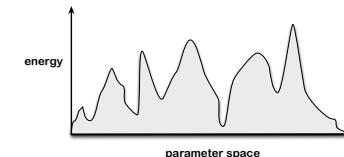
- stereochemically allowed / preferred positions and angles of backbone and side chains
- buried hydrophobic surface
- hydrogen bonds between polar residues and other polar residues (or backbone or water)
- electrostatic & Van der Waals interactions
- ...

bioinf WS19 • lec10 • S.Hartmann

Ab initio structure prediction

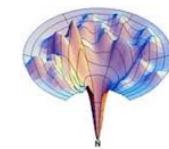
KSARPEYLQKLTYDIVLPADMP

thermodynamic & physicochemical theory

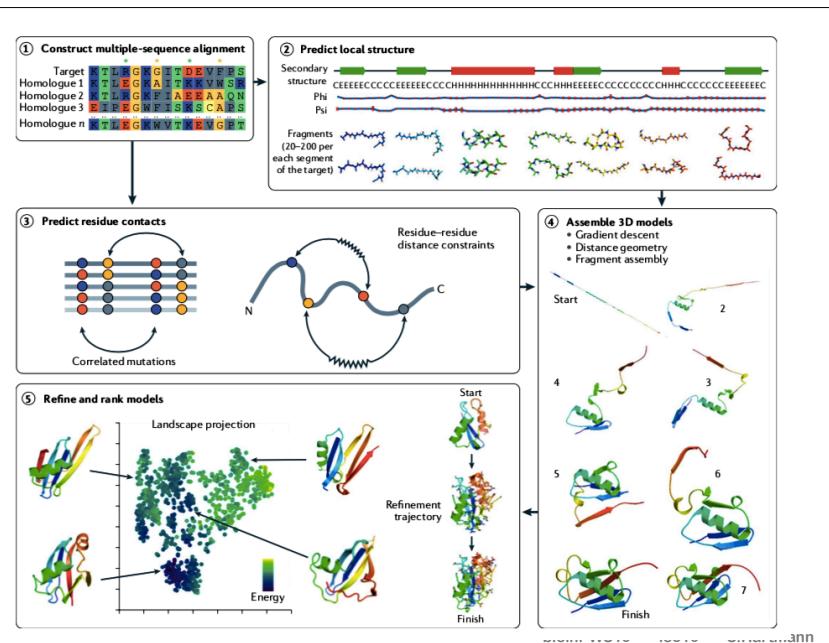


challenges

- number of possible conformations
- computing power
- local minima



bioinf WS19 • lec10 • S.Hartmann



How well does it work?

CASP: Critical Assessment of protein Structure Prediction

- every two years, 1994 - current (CASP 13)
- structures about to be solved experimentally
- their sequences are made available
- predictions are submitted ...
...and then compared to the solved structures

assessment of:

- structure prediction, structure refinement, ...

<http://www.predictioncenter.org>

bioinf WS19 • lec10 • S.Hartmann

How well does it work?

CASP13: progress!

- increased amount of experimental data
- sequence data for huge protein families
- improved hardware & algorithms
- use of deep learning (deep neural networks)
 - machine learning approach
 - use of training data (sequences & their solved structures)
 - “learn”
 - predict structure of new sequence

bioinf WS19 • lec10 • S.Hartmann

How well does it work?

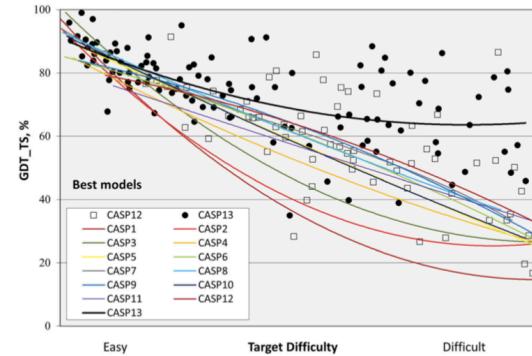


FIGURE 1 Trend lines of backbone accuracy for the best models in each of the 13 CASP experiments. Individual target points are shown for the two most recent experiments. The accuracy metric, GDT_TS, is a multiscale indicator of the closeness of the C α atoms in a model to those in the corresponding experimental structure. Target difficulty is based on sequence and structure similarity to other proteins with known experimental structures (see Reference 5 for details). There is a striking improvement in model accuracy in CASP13 (top black line), particularly for the more difficult targets

A Kryshtafovych et al., 2019. Critical assessment of methods of protein structure prediction (CASP)—Round XIII

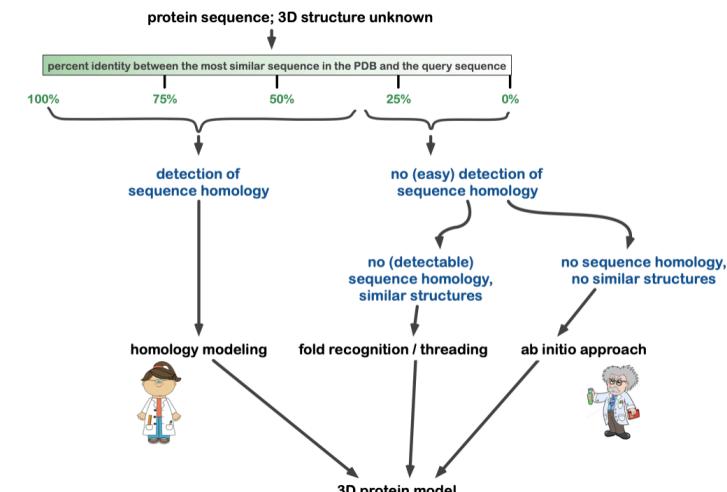
bioinf WS19 • lec10 • S.Hartmann

Terms and concepts

- domain, fold, super-secondary structure, CATH, threading, ab initio
- structure elements of proteins
- ~number of sequences vs. structures
- main approaches to predict 3D protein structure: principles, assumptions, applications, limitations
 - homology modeling
 - threading
 - ab initio prediction

bioinf WS19 • lec10 • S.Hartmann

Today's computer exercise



bioinf WS19 • lec10 • S.Hartmann

BLAST® > blastp suite

Home

Standard Protein BLAST

blastn **blastp** **blastx** **tblastn** **tblastx**

Enter Query Sequence BLASTP programs search protein databases using a protein query.

Enter accession number(s), gi(s), or FASTA sequence(s) (?)

From To

Or, upload file No file selected. (?)

Job Title

Enter a descriptive title for your BLAST search (?)

Align two or more sequences (?)

Choose Search Set

Database Non-redundant protein sequences (nr) (?)

Organism Non-redundant protein sequences (nr) (?)

Optional Reference proteins (refseq_protein) (?)

Exclude Model Organisms (landmark) (?)

Optional UniProtKB/Swiss-Prot(swissprot) (?)

Exclude UniProtKB/Swiss-Prot(swissprot) (?)

Optional Patented protein sequences(pat) (?)

Optional Protein Data Bank proteins(pdb) (?) **Protein Data Bank proteins(pdb)** (?)

Optional Metagenomic proteins(env_nr) (?)

Optional Transcriptome Shotgun Assembly proteins (tsa_nr) (?)

Create custom database

Enter an Entrez query to limit search (?)

Program Selection

Algorithm **blastp (protein-protein BLAST)** (?)

PSI-BLAST (Position-Specific Iterated BLAST) (?)

PHI-BLAST (Pattern Hit Initiated BLAST) (?)

DFAST (Domain Enhanced FAST in Time Accelerated BLAST) (?)

DRAFT WS19 • bioinf WS19 • lec10 • S.Hartmann

3 pools only!

- 1a, 1b
- 2a