



Bioinformatics

Knowledge-representation in molecular biology

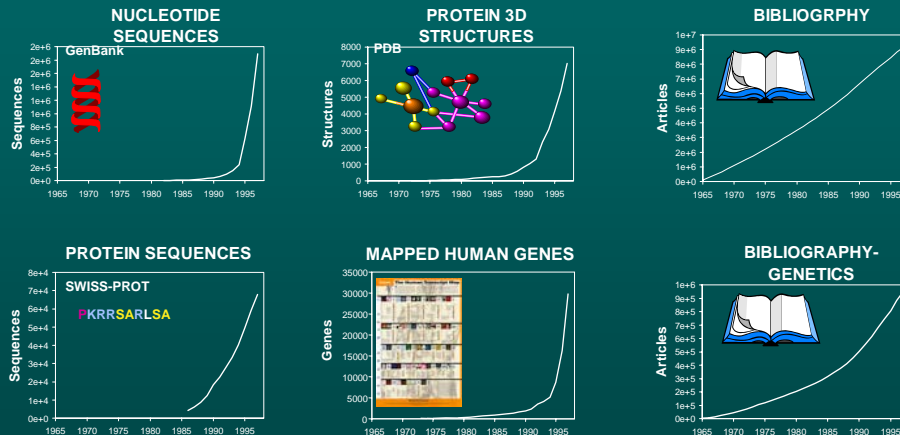
Sándor Pongor

Protein Structure and Bioinformatics, ICGEB, Trieste

An overview of bioinformatics

- History and development
- Models:
 - Sequences,
 - 3D structures
 - Networks
- Similarity and classification:
 - database search,
 - consensus descriptions
- Integrated resources

Representation of biological knowledge



Source: NCBI

Bioinformatics milestones 1

- 1962 - Pauling's theory of molecular evolution
- 1967 - Margaret Dayhoff's Atlas of Protein Sequences
- 1970 - Needleman-Wunsch algorithm
- 1977 - DNA sequencing and software to analyze it (Staden)
- 1981 - The concept of a sequence motif (Doolittle)
- 1982 - Phage lambda genome
- 1983 - Database search (Wilbur-Lipman)
- 1985 - FASTP/FASTN: fast sequence similarity searching
- 1987 - Sequence profiles
- 1987 - EMBL, Genbank, Swiss-Prot databases

Bioinformatics milestones 2

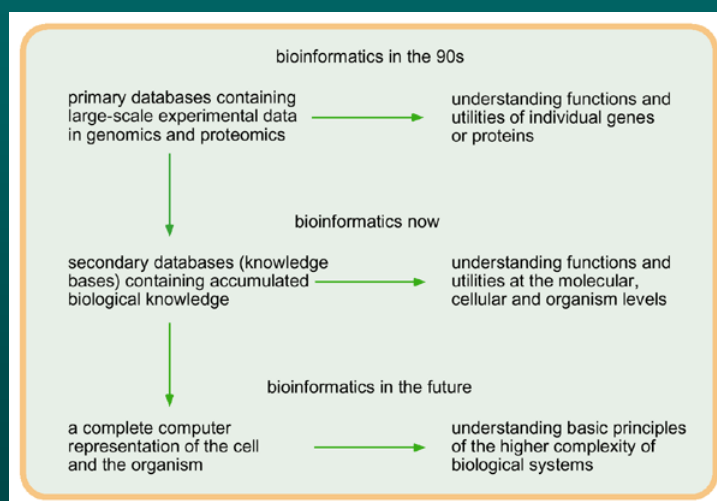
- 1988 - National Center for Biotechnology Information (US)
- 1988 - EMBnet network for database distribution
- 1990 - BLAST: fast sequence similarity searching
- 1991 - EST: expressed sequence tag sequencing
- 1993 - Sanger Centre, Hinxton, UK
- 1994 - EMBL European Bioinformatics Institute, Hinxton, UK
- 1995 - First bacterial genomes
- 1996 - Yeast genome
- 1997 - PSI-BLAST
- 1998 - Worm (multicellular) genome
- 2000+ The rice and human genomes.
- Microarrays

The ingredients

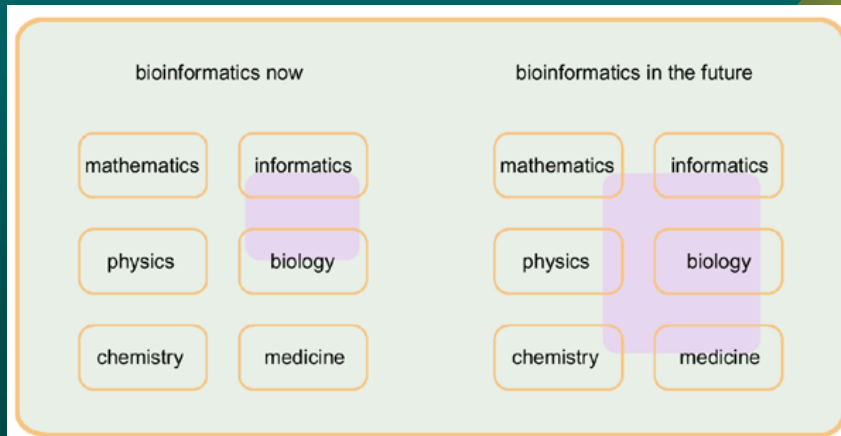
- Data collection techniques (DNA sequencing, protein sequencing, microarrays)
- Theoretical milestones (concepts of DNA structure, protein structure, evolution)
- Algorithms and programs (BLAST, FASTA)
- Databases
- Institutions
- Genomic data



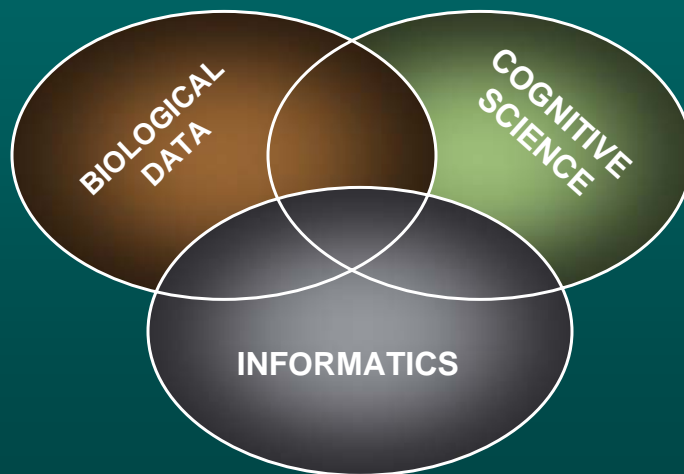
The evolution of bioinformatics



Bioinformatics is an approach to biology...



BIOINFORMATICS



MODELS

Molecular structures

MARTKQTARK
STGGKAPRKQ
LATKAARKSA

Sequences

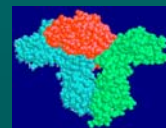
CIPKWNRCGPKMDGVPCCEPYTCTSDYYGNCS



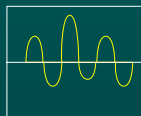
Extended sequences
(e.g. disulphide-topologies)



Domain-cartoons
(sec. str. cartoons)



3D structures

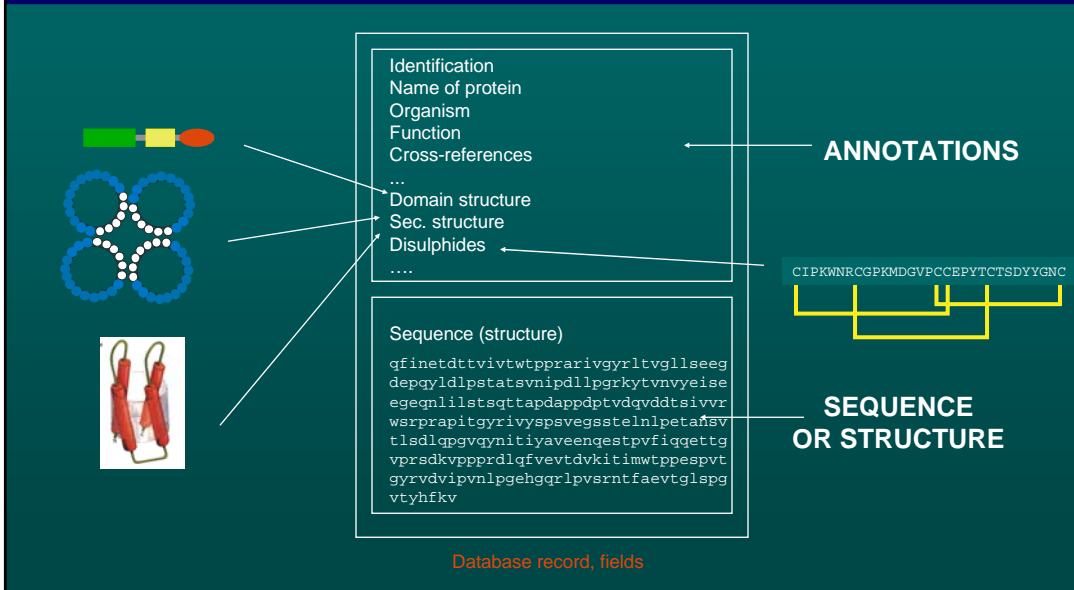


Diagrams (hydrophobicity plots, helical circles)

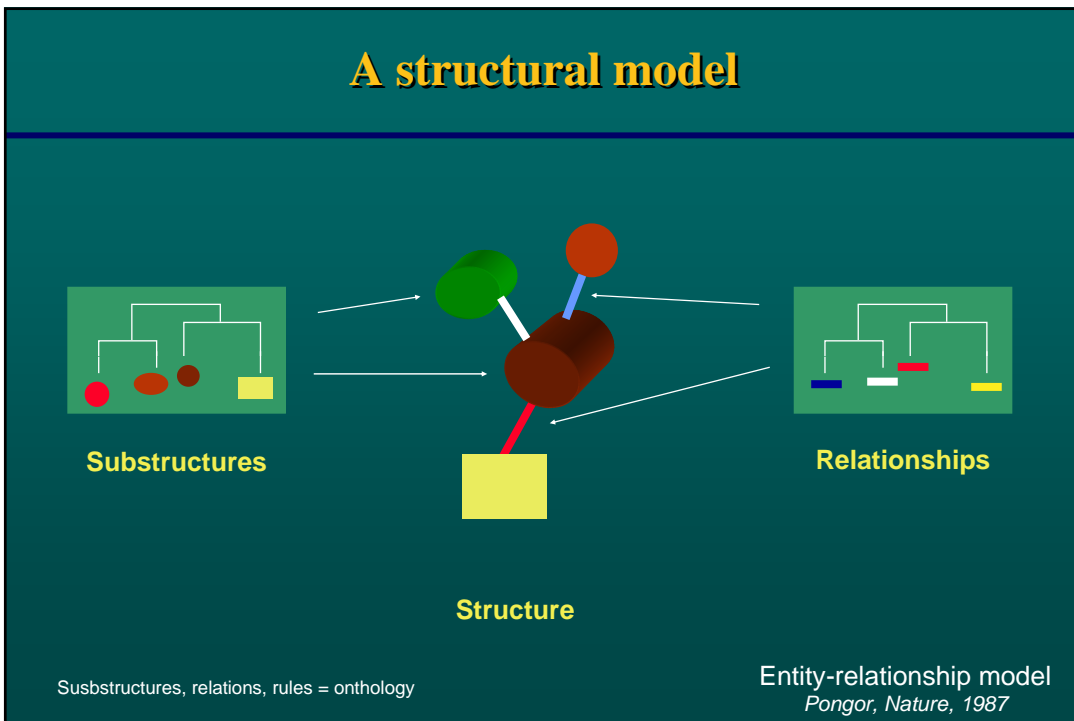


3D cartoons

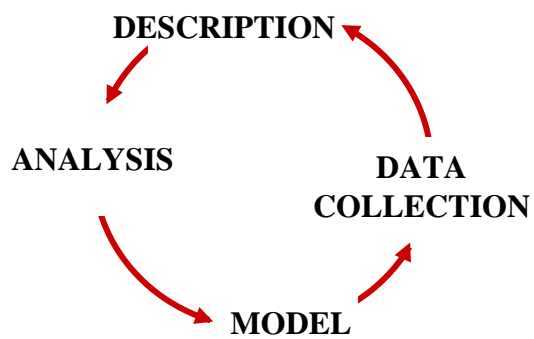
Structures As Database Records



A structural model

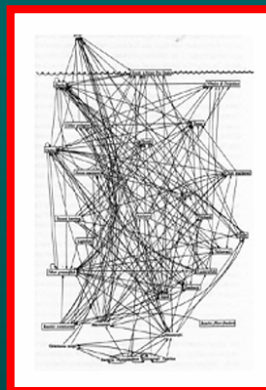
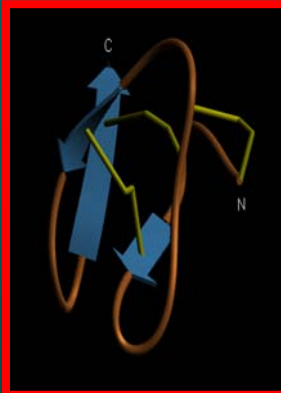


Molecules change



Models

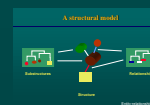
```
tassfvsvwsasdtvsgfrvey  
elseegdepyldlpstatsvni  
pdllpgrkytvnyveiseeqn  
lilstsqttapdappdptvdvd  
dtsivvrwrprapitgyrivys  
psvegsstelnlpetansvtlsd  
lpggvqynitiyaveengestpv  
fiqqettgvprsdkvppprdlqf  
vevtdvkitimwtppesvptgyr  
vdvipvnlpgehgqrlpvsrntf  
aevtglspgvtyhfkvfavnqgr  
eskpltagqatkldaptnlqfin  
etdttvivtwtpprarivgyrlt  
vgltrggqpkqynvgaasqypl  
rnlqgseyavslvavkgnqqsp  
rvtgvfttlqplgsiphyntev  
ettivitwtpaprigfklgvrps  
qggeaprevtsesgsivvsqgtp  
gveyvvtisvlrdgqerdapivk
```



SEQUENCES

3D STRUCTURES

NETWORKS



■ SEQUENCES

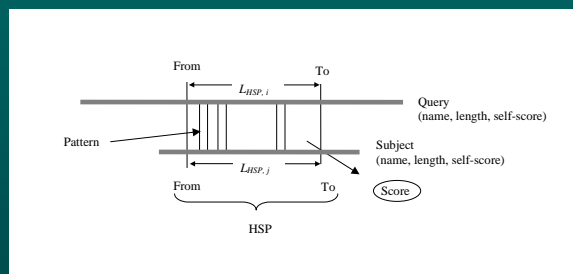
LANGUAGE

Sequences as language

```

qfinetdttvivtwtpprarivgyrl
tvglseegdepyldlpstatsvni
pdllpgrkytvnyveiseeegqnll
stsgttapdapdptvdqddtsivv
rwsrprapitgyrivyspsvegsste
lnlpetansvtlsdlqpgvqynitiy
aveengestpvfiqgettgvprsdkv
ppprdlqfvevtvdkitimwtppesp
vtgyrvdviplnlpgehgqrlpvsrn
tfaevtqlspgvtyhfkvfavnqgre
skpltaqqatkldaptnlqfinetdt
tvivtwtpprarivgyrltvgltrgg
qpkqynvqpaasqyplrnlpqgseya
vslvavkgnqqsprvtgvfttlqplg
siphyntevtttvtvwtppaprigf
klgvrpsgggeaprevtseegsivvs
gltpgveyvytisvlrdgqerdapiv
kkvvtplsppntlhleanpdtgvlv
swersttpditgyrittptngqqyy
sleevvhadqssctfenlspgleynv
svytkddkesvpissfvvsvwsas
dtvsgfrveyelseegdepyldlps
tatsvniplpgrkytvnyveisee
    
```

Alignments



Character strings, computer-languages,
Chomsky et al, etc.

LANGUAGE

The language of bibliographies

Structures As Database Records

- Annotations:
 - Identification
 - Name of protein
 - Organism
 - Function
 - Cross-references
- Sequence (structure):
 - Domain structure
 - Sec. structure
 - Disulphides

Database record fields

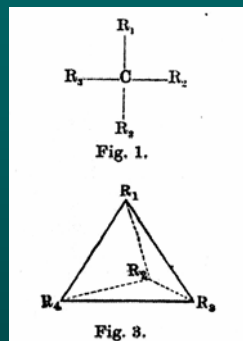
Keyword-collections, ontologies, etc.

■ 3D STRUCTURES

Chimie dans l'espace

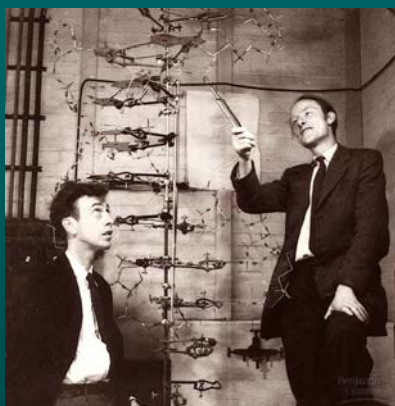
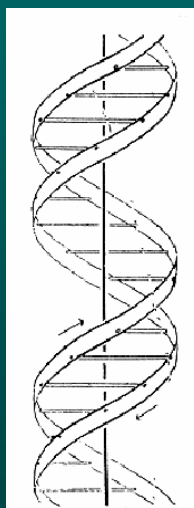


Van 't Hoff
1852-1911



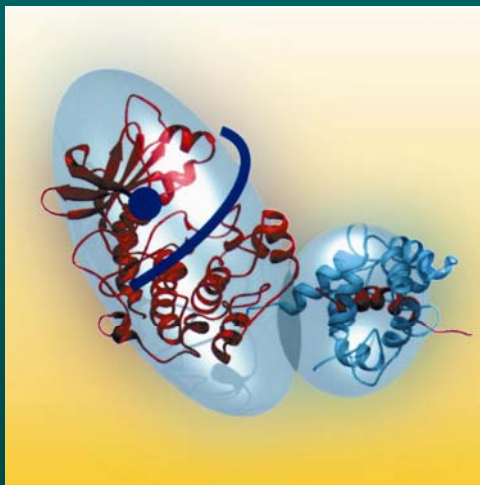
1898

Some molecules are more equal than others...



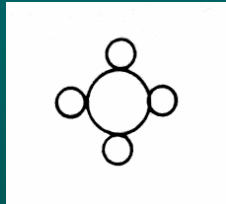
... "This figure is purely diagrammatic. The two ribbons symbolize the the phosphate-sugar chains, and the horizontal rods the pairs of the bases holding the chains together. The vertical line marks the fibre axis"

Protein models

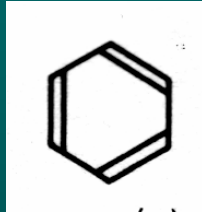


■ NETWORKS

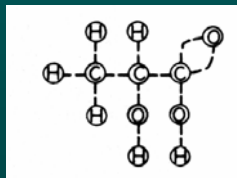
TOPOLOGIES,
GRAPHS **Small molecules – classical graphs**



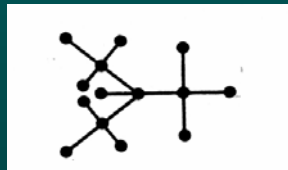
Loschmidt, 1861



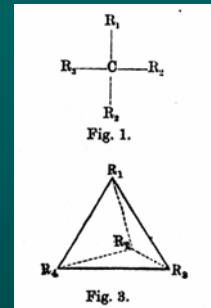
Kekulé, 1865



Crum Brown, 1861



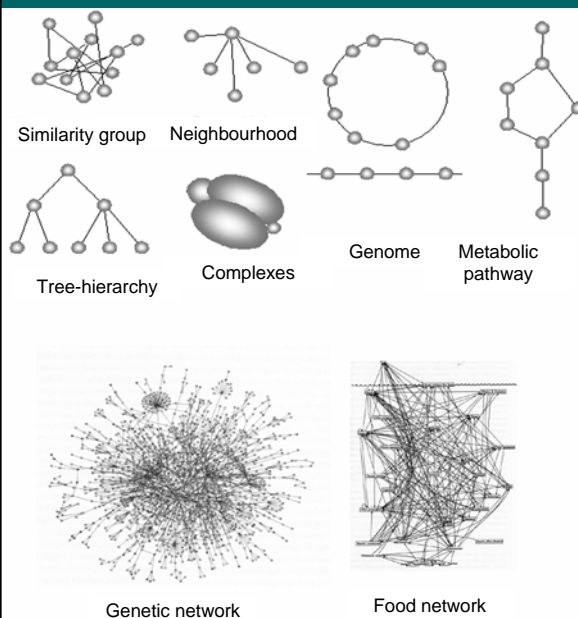
Cayley, 1872



Van't Hoff, 1898

TOPOLOGIES,
GRAPHS

Genomes, assemblies

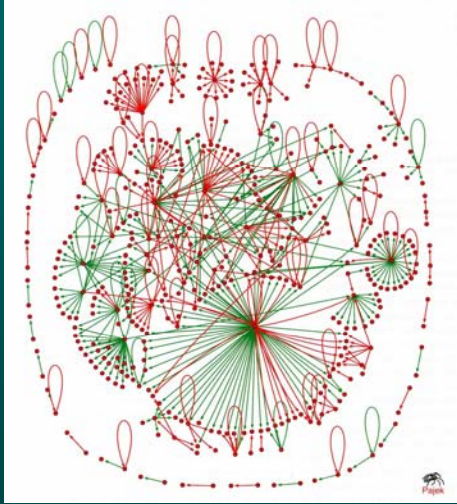


Entity-relationship models
Topological meta-models

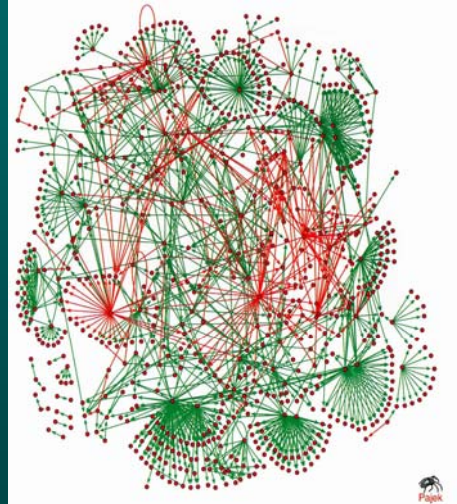
TOPOLOGIES,
GRAPHS

The transcription regulatory networks

+ (up)
- (down)



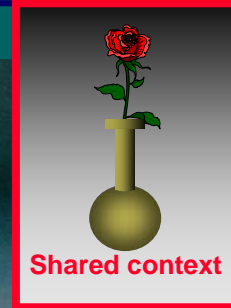
E. coli



S. cerevisiae

SIMILARITY, CLASSIFICATION

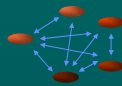
The concept of similarity I



...easier if modular

Multiple Objects

Structural similarity



Similarity groups or neighborhoods

```

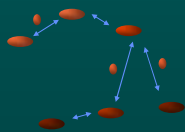
CGPK-MDGVPCCEPY
CGGQNWSGPTCCASG
CSPTSYN---CCR--
CSRLMY---DCCY--
CIPYYL---DCCPEL
    
```

Multiple alignments

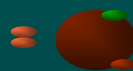


Evolutionary trees

Context (function)



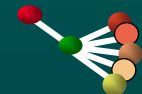
Metabolic pathways



Subunit structures, ligands

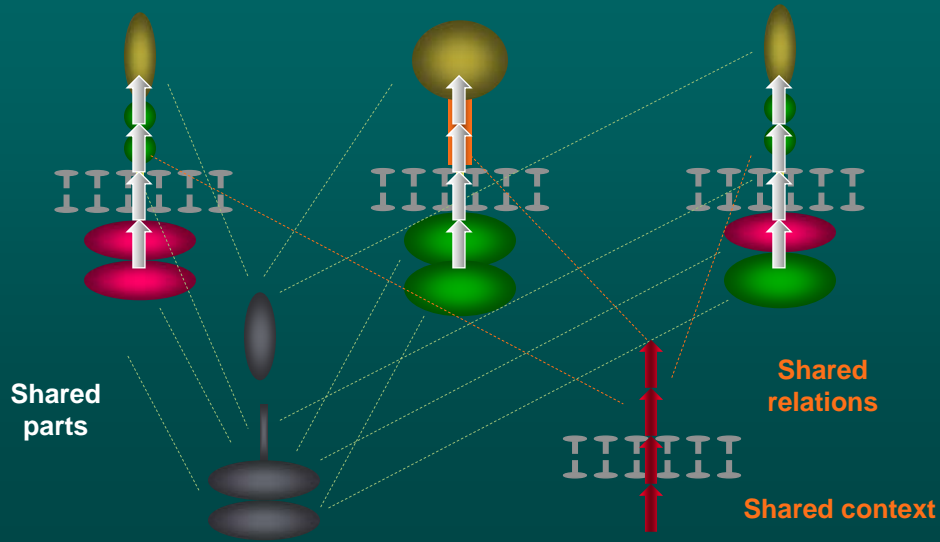


Genomes

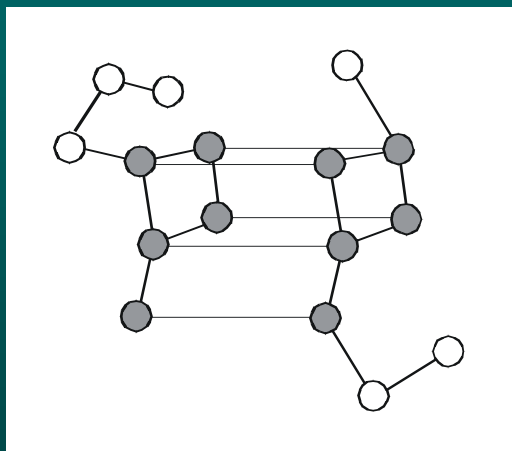


Trajectories

Similarity of molecules



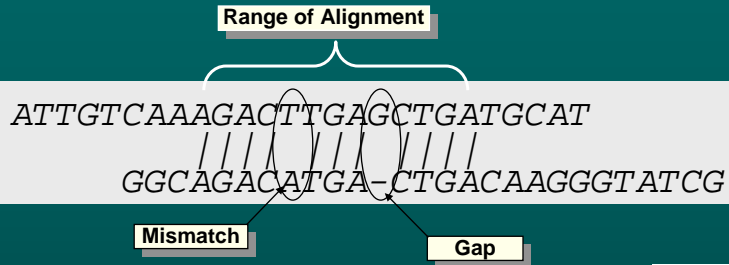
Substructure identity ~ similarity



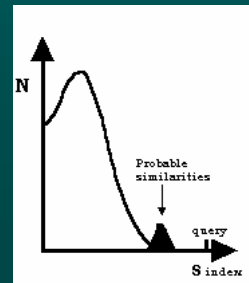
"The similarity of objects can be best described as partial identities of components and relationships

Erich Goldmeier, The similarity of perceived forms, 1936

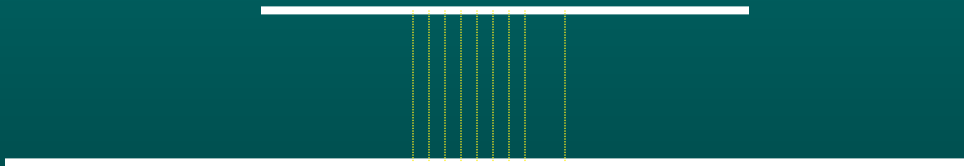
Quantification of sequence similarity



Score ~ sg like an edit distance

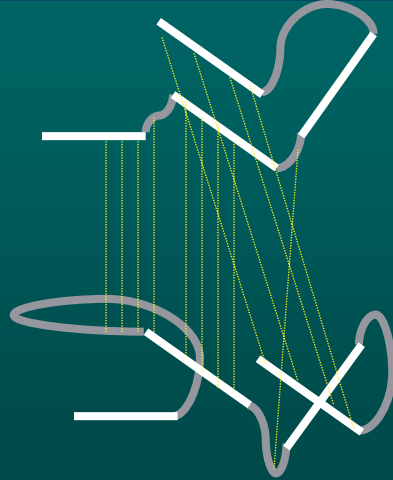


Sequence comparison



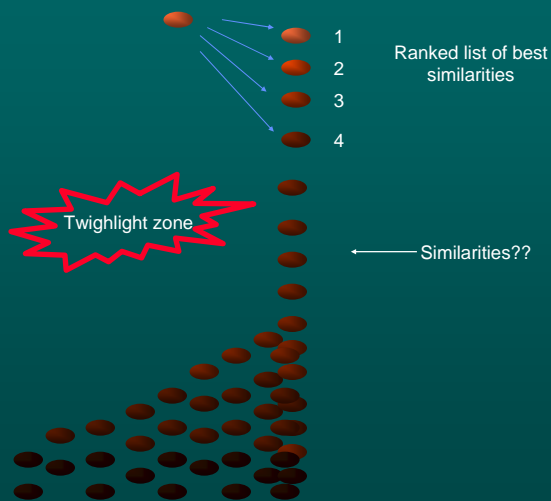
1. Find region of alignment (fast)
2. Calculate similarity score (fast)

3-D comparison



1. Find region of alignment (slow)
2. Calculate similarity score (fast)

Using similarity: Comparing one sequence with a group (database)



SEQUENCE SCORE	DESCRIPTION
SWISSPROT:Q53844 457.36	ALPHA-AMYLASE INHIBITOR AAI_2/95
SWISSPROT:P77977 152.82	CELLULOSE BINDING PROTEIN
SWISSPROT:Q0V1 145.77	EXOGLUCANASE I PRECURSOR
SWISSPROT:Q126 145.66	CELLULOSE (EC 3.2.1.91)

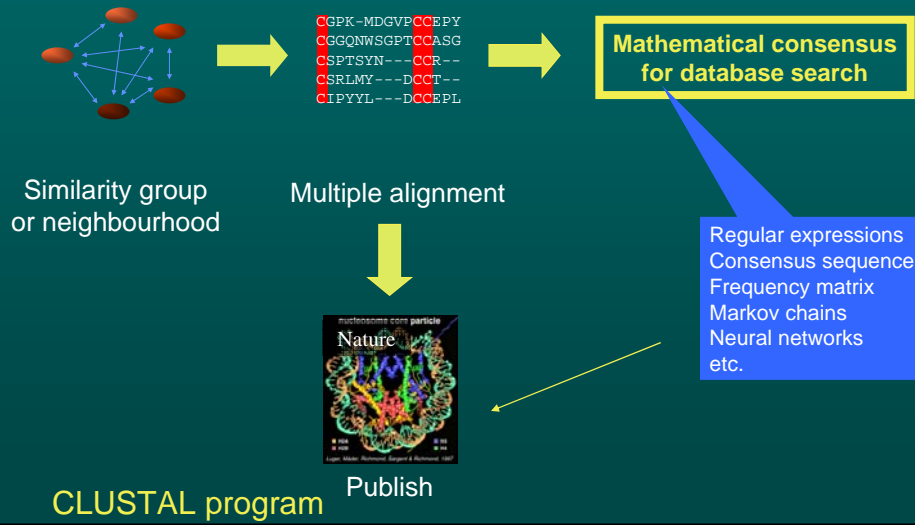
```

| 1 001 5000 0
| 2 001 0900 0
| 3 001 1000 0
| 4 001 010 0
| 5 001 950 0
| 6 001 500 0
| 7 001 090 0
| 8 001 030 0
| 9 001 010 0
|10 001 000 0
|11 001 090 0
|12 001 000 0
|13 001 001 0
|14 001 001 0
|15 001 001 0
|16 001 001 0
|17 001 001 0
|18 001 001 0
|19 001 001 0
|20 001 001 0
|21 001 001 0
|22 001 001 0
|23 001 001 0
|24 001 001 0
|25 001 001 0
|26 001 001 0
|27 001 001 0
|28 001 001 0
|29 001 001 0
|30 001 001 0
|31 001 001 0
|32 001 001 0
|33 001 001 0
|34 001 001 0
|35 001 001 0
|36 001 001 0
|37 001 001 0
|38 001 001 0
|39 001 001 0
|40 001 001 0
|41 001 001 0
|42 001 001 0
|43 001 001 0
|44 001 001 0
|45 001 001 0
|46 001 001 0
|47 001 001 0
|48 001 001 0
|49 001 001 0
|50 001 001 0
|51 001 001 0
|52 001 001 0
|53 001 001 0
|54 001 001 0
|55 001 001 0
|56 001 001 0
|57 001 001 0
|58 001 001 0
|59 001 001 0
|60 001 001 0
|61 001 001 0
|62 001 001 0
|63 001 001 0
|64 001 001 0
|65 001 001 0
|66 001 001 0
|67 001 001 0
|68 001 001 0
|69 001 001 0
|70 001 001 0
|71 001 001 0
|72 001 001 0
|73 001 001 0
|74 001 001 0
|75 001 001 0
|76 001 001 0
|77 001 001 0
|78 001 001 0
|79 001 001 0
|80 001 001 0
|81 001 001 0
|82 001 001 0
|83 001 001 0
|84 001 001 0
|85 001 001 0
|86 001 001 0
|87 001 001 0
|88 001 001 0
|89 001 001 0
|90 001 001 0
|91 001 001 0
|92 001 001 0
|93 001 001 0
|94 001 001 0
|95 001 001 0
|96 001 001 0
|97 001 001 0
|98 001 001 0
|99 001 001 0
|100 001 001 0

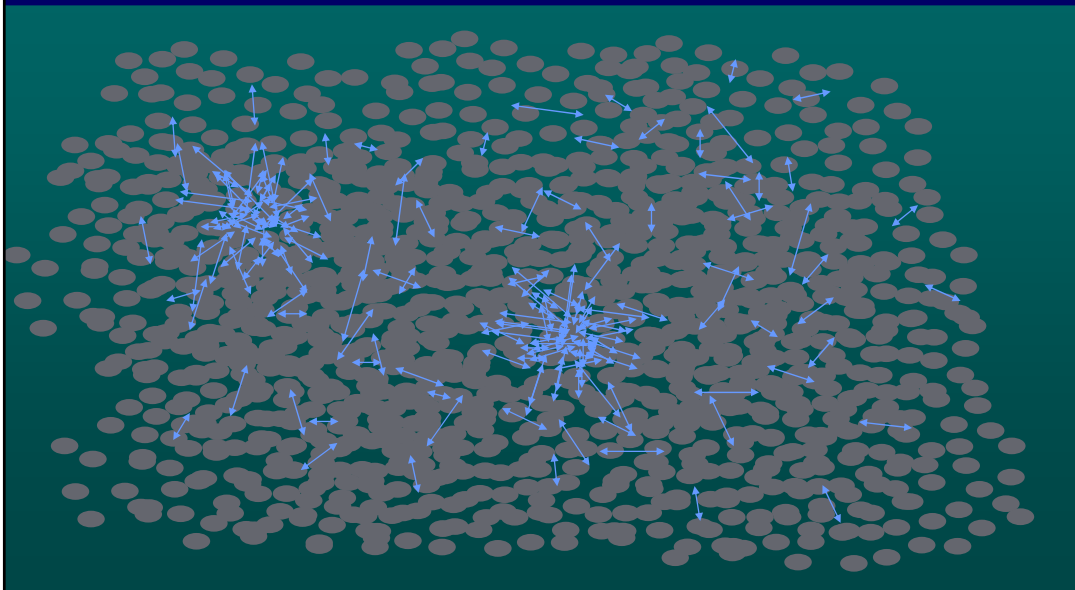
```

BLAST program

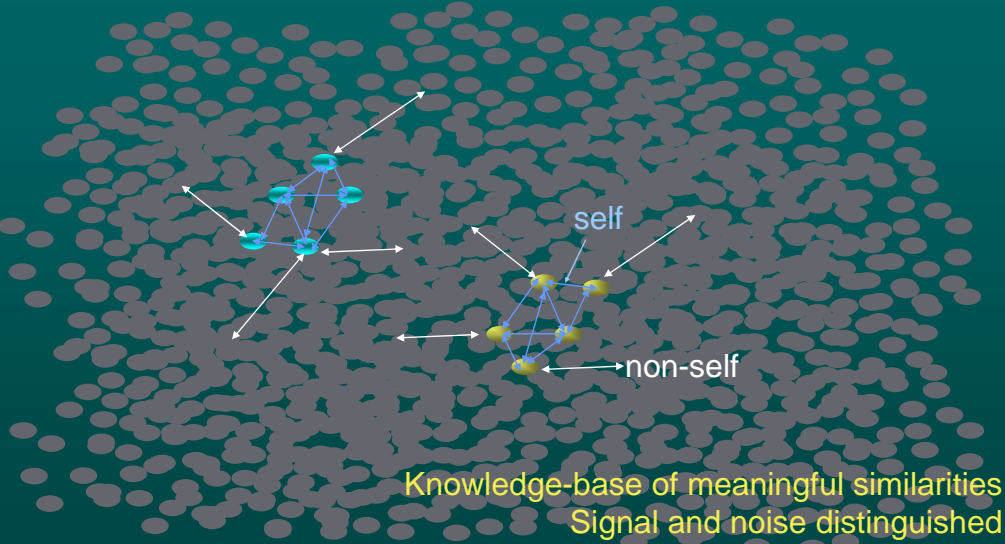
Using similarity: comparing a group with itself



The database as network of similarities: A memory network model

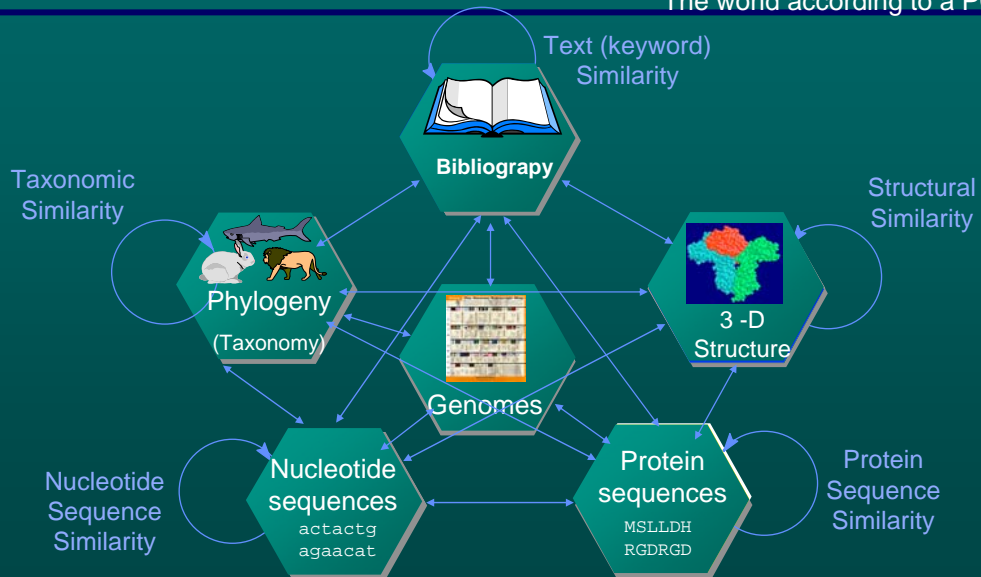


Das Wohltemperierte Database: Similarity network as a knowledge representation



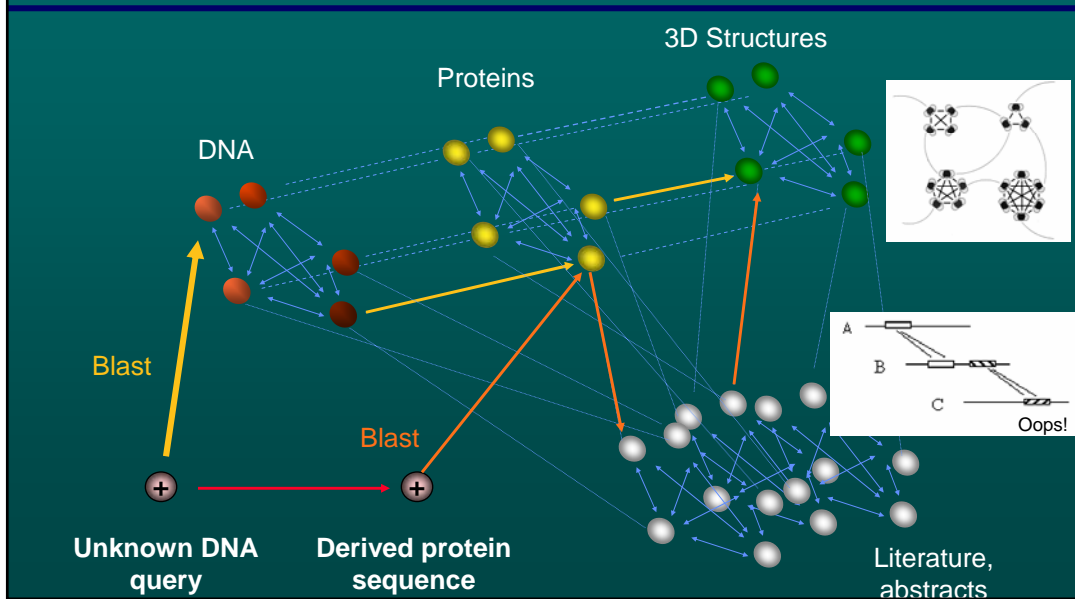
Biological knowledge as a network of data

The world according to a PC...

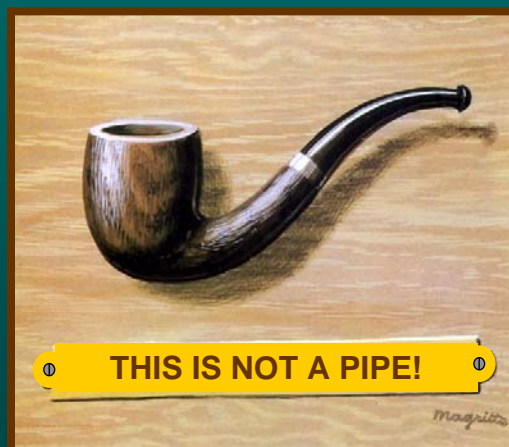


Source: NCBI

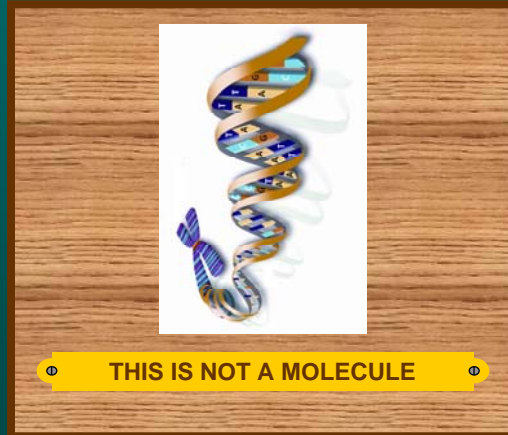
Search on a preprocessed, integrated database: the importance of a good neighbourhood



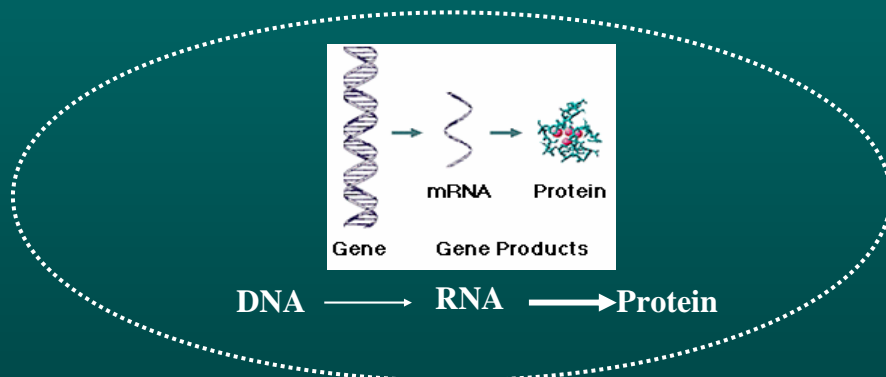
Models are human constructs...



Models are human constructs...

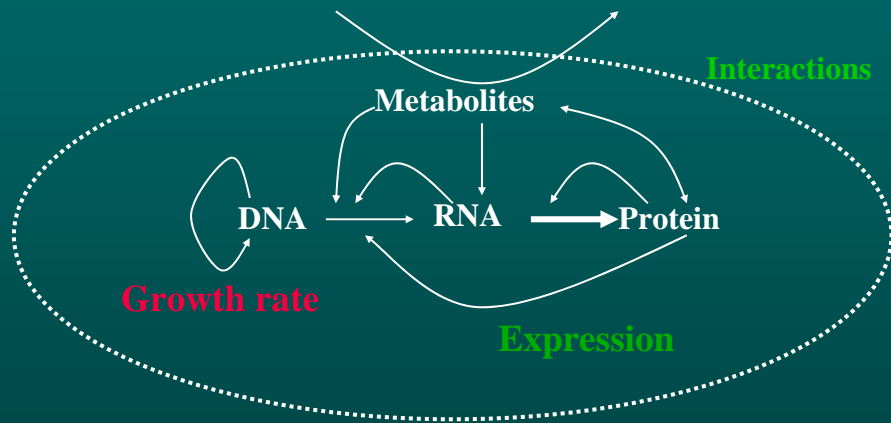


The central dogma:



Dogma, paradigm, mythology

New central dogma: Self-assembly, catalysis, replication, networks



Polymers: Initiate, elongate, terminate, fold, modify, localize, degrade

Evolution + Self assembly, Systems biology

Summary of topics discussed

- History and development
- Models:
 - sequences,
 - 3D structures
 - Networks
- Similarity and classification:
 - database search,
 - consensus descriptions
- Integrated resources, knowledge integration

Summary of the introduction

- Bioinformatics is the science of biological information or rather a computer-based approach to biological problems.
- All kinds of biological data are structures defined with entities and relationships (metabolites, genes, networks).
- Typical tasks: Similarity search, categorization and clustering
- Simultaneous handling of many, complex datatypes

Computer methods in Molecular Biology Trieste June 24 - July 1, 2005

- Sequence database searching, theory and practice (Dave Judge and Jack Leunissen)
- Nucleic acid databases, Medline, Pubmed (David Landsman)
- Protein databases, Swissprot, Prosite (Elisabeth Gasteiger)
- WWW servers, EBI (Kristian Vlahovicek)
- Gene discovery (Luciano Milanese)
- Genome analysis (Martin Bishop)
- Algorithms, domain similarities (Kristian Vlahovicek, Sándor Pongor)

On-line help to this lecture

- Bioinformatics tutorials on-line

<http://www.ebi.ac.uk/2can/home.html>

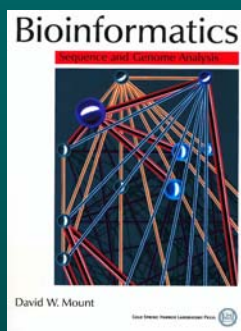
- ICGEBnet

<http://www.icgeb.org/~netsrv/>

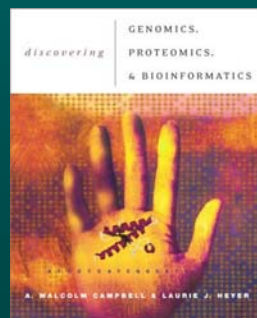
- The Trieste bioinformatics course

<http://www.icgeb.org/~netsrv/netcourse.html>

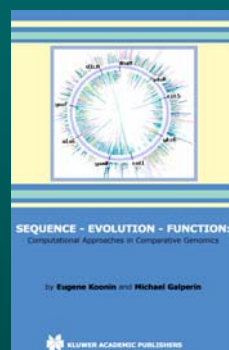
Reading about bioinformatics



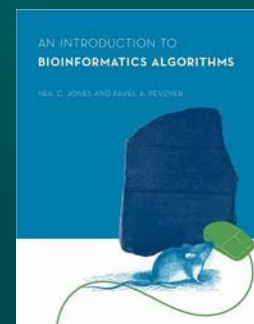
In depth introduction



Genomics research problems



Evolutionary principles



Math principles