*Introdu*...

**Bioinformatics**

# Introduction to Bioinformatics

Dr. Taysir Hassan A. Soliman
Associate Professor,
Head of Information System Department,
Faculty of Computers and Information
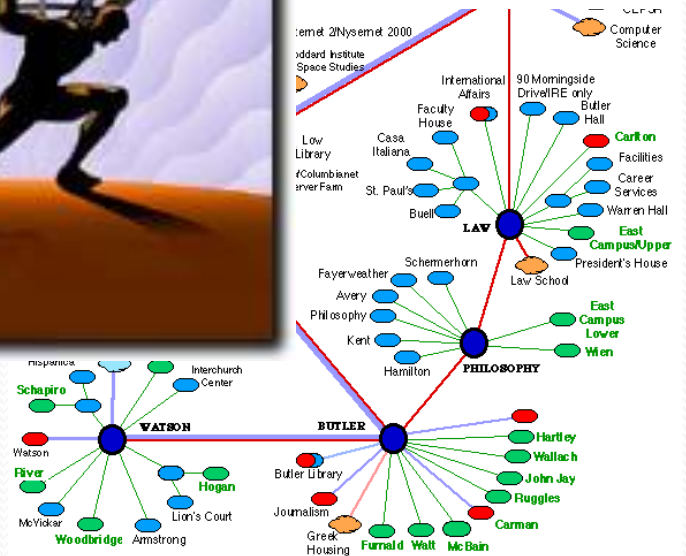Assiut University, Egypt
taysirhs2@gmail.com

# Outline

- Introduction to Bioinformatics
- Bioinformatics Applications
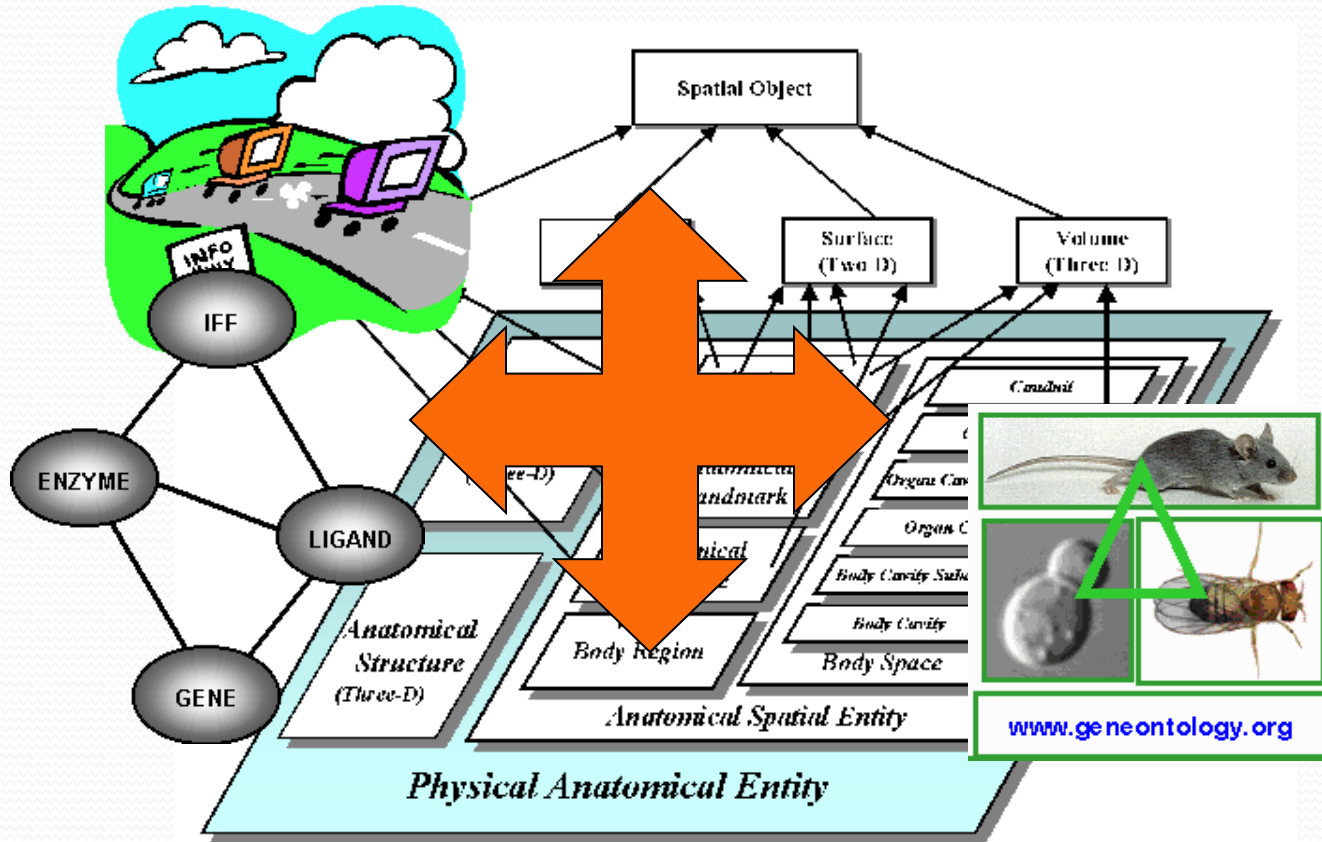- Bioinformatics databases
- Sequence Alignment

# Science then, then and now

A vast amount, rapidly generated related but highly distributed and semantically unconnected information

# Science then, then and now

# How much Computing skills?

- Bioinformatics can be seen as a tool that the biologist needs to use - like PCR
- Or should biologists be able to write their own programs and build databases?
  - it is a big advantage to be able to design exactly the tool that you want
  - this may be the wave of the future

    *"Two months in the lab can easily save an afternoon on the computer."*
    —Alan Bleasby, 1997

- Q: Is this school going to train "bioinformatics professionals" or biologists with informatics skills?
- A: Both!

# What is Bioinformatics?

- The use of computers to collect, analyze, and interpret biological information at the molecular level.

- A set of software tools for molecular sequence analysis

**YES**

- DNA & protein sequence databases
- Sequence similarity, alignment, & assembly
- Sequence patterns/motifs
- Phylogenetics
- Microarray gene expression data
- Protein structure prediction
- Mapping metabolic and regulatory pathways (graph theory)

**NO**

- patient medical charts, billing, hospital payroll, etc.
- X-ray image analysis

**MAYBE**

- Ontologies
  (biological function, research methods, clinical terminology, etc.)

# Bioinformatics - origins

- Driven by experimental molecular biology
  - lab folks generate the data, then need a way to organize and analyze it
- Grabs methods from many different fields
  - biostatistics, machine learning, data mining, linguistics, etc
- Use computer (algorithms) to gain novel biological knowledge.
- Use biological knowledge to construct algorithms.

# The Biologist in the Age of Information

# Training "computer savvy" scientists

- Know the right tool for the job

- Get the job done with tools available

- Network connection is the lifeline of the scientist

- Jobs change, computers change, projects change, scientists need to be adaptable

# The job of the biologist is changing

- ## As more biological information becomes available ...

  - The biologist will spend more time using computers, building and mining databases

  - The biologist will spend more time on data analysis (and less doing lab biochemistry)

  - Biology will become a more quantitative science (think how the periodic table and atomic theory affected chemistry)

# Biological Data Characteristics

1. Huge  data
2. Heterogeneous distributed data
3. Frequently updated data
4. Defining and representing complex queries are extremely important to the biologist
5. Most biologist will not care or know about the data structure or the schema design
6. Users of biological information often require access to previous versions of existing data.

# I. "Traditional" bioinformatics methods

- Conduct online literature and similarity searches (NCBI Entrez and Blast)
- Use desktop sequence analysis tools
  - restriction digest, PCR primer design, ORF finding
- Assembly of automated sequencing reads

# II. **More advanced stuff**

- Multiple alignment

- Phylogenetic trees

- Motif/domain analysis of proteins

  (Pfam, Blocks, ProDom)

- Motif/domain analysis of DNA

  (promoters, transcription factors, intron splice sites)

- Genefinding in genome data

  combining data from ORFs, promoters, and cDNA homology

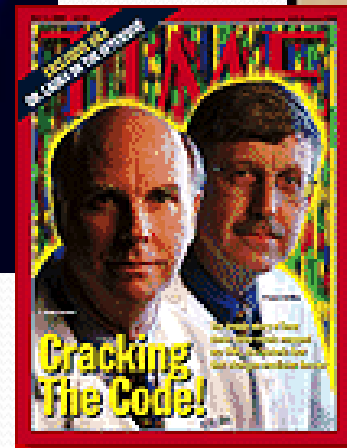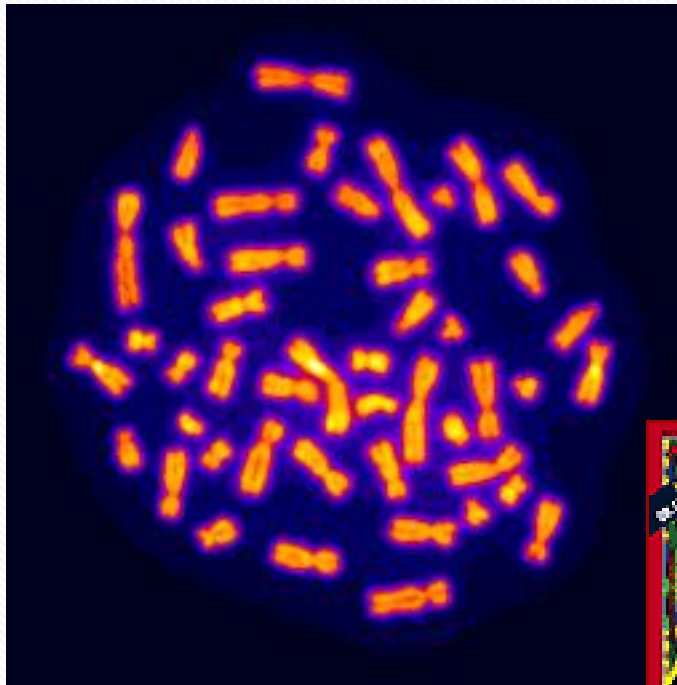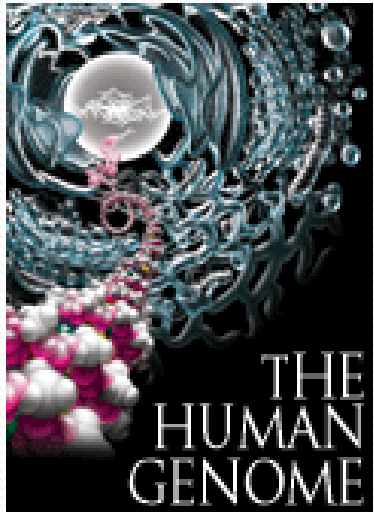# III. Genome scale data analysis

- Handling large amounts of data
  - Create an experiment or lab database
  - use traditional bioinformatics tools on different data
  - scripting languages (simple programming tools, Perl)
- Microarray gene expression analysis
  - differential expression and classification/prediction
  - clustering, principle components
  - functional genomics - pathways, ontology classification
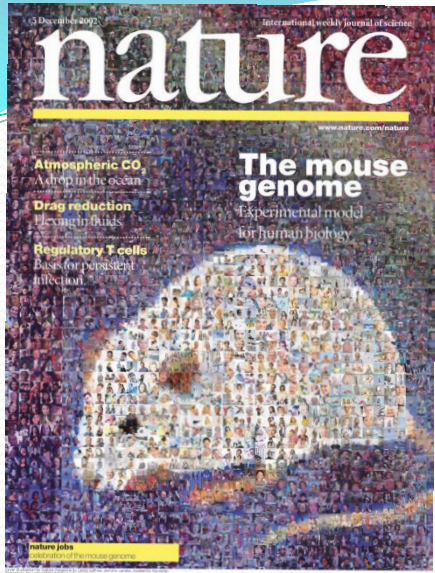- Genome-wide SNP or genome tiling analysis

# A Genome Revolution in Biology and Medicine

- We are in the midst of a "Golden Era" of biology
- The Human Genome Project has produced a huge storehouse of data that will be used to change every aspect of biological research and medicine
- The revolution is about treating biology as an information science, not about specific biochemical technologies.
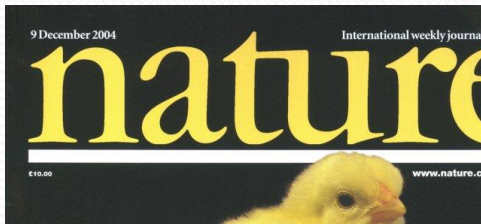
# Genome Projects

The Human Genome sequence is complete approximately 3.2 billion base pairs

# More Genomes

# Bioinformatics …
## A breakthrough towards …
## Various Applications



Clinical Outcomes Data

Targeted Pharmaceuticals

Pathology

Cancer Treatments

Bench

Bedside

Proteomic Data

Treatment and Care

Laboratory Data

Genetic Data

Surgical & Radiotherapeutic Technologies

Structural Genomics

Interactomics

Microarray Gene Expression

Querying

Indexing

Data Mining

Visualization

Genomics

Metabolomics

Ontologies & Specialized Databases

**Various Disciplines**
**But**
▪**How to reach each?**
▪ **How to integrate?**
▪ **What about after integration?**

# So,...



RESOURCES

a **Reality.**

# All the Genes

- Any human gene can now be found in the genome by similarity searching with over 99.9% certainty.
- However, the sequence still has gaps
- Still can't identify pseudogenes, false genes with certainty
  - This will improve as more sequence data accumulates
- We are getting close to a complete list of human genes and proteins
  - Needed as a starting point for gene expression, pattern finding, and systems biology

# Raw Genome Data:

# bioinformatics Databases

# Bioinformatics Challenges

## The huge dataset

- Lots of new sequences being added
  - automated sequencers
  - genome sequencing
  - EST sequencing
  - environmental/metagenomic sequencing

- GenBank has over 100 **Billion** bases and is doubling every year!!
  - problem of exponential growth
  - how can computers keep up?
  - hard drives are cheaper, but processor speeds are not keeping up

**100 Gigabases**

GenBank and its collaborating databases, the European Molecular Biology Laboratory and the DNA Data Bank of Japan, have reached a milestone of 100 billion bases from over 165,000 organisms. See the press release or find more information on GenBank.

Growth of the International Nucleotide Sequence Database Collaboration

Base Pairs contributed by   GenBank® —■   EMBL —■   DDBJ —■

# DNA Sequencing capability has grown exponentially

DNA sequences in GenBank

Doubling time = 18 months

# DNA Sequencing & Assembly

- Automated Sequencers

- ~500 bp reads must be assembled into complete genes & genomes

- faster sequencing relies on better software

# Next Generation Sequencing

# Genomics Technologies

Next-Generation DNA sequencing •

Automated annotation of sequences •

DNA microarrays •

gene expression (measure RNA levels) •

single nucleotide polymorphisms (SNPs) •

ChIP-chip, genomic tiling, etc •

Proteomics

Protein-protein

# Biological Information

mRNA Expression

Protein 3-D Structure

Protein 2-D gel

Genome sequence

The Cell

Mass Spec.

# New Types of <u>Big</u> Biological Data

- Microarrays - gene expression

- Networks of protein-protein interactions

# Microarray Data Analysis

- Linkage between gene expression data and gene sequence/function/metabolic pathways databases

- Discovery of common sequences in co-regulated genes

- Meta-studies using data from multiple experiments

# Impact on Bioinformatics

- Genomics produces high-throughput, high-quality data, and bioinformatics provides the analysis and interpretation of these massive data sets.

- It is impossible to separate genomics laboratory technologies from the computational tools required for data analysis.

# Example of Biological Database Formats

| Database | Data Format | Website | Access |
|----------|-------------|---------|--------|
| IntAct | Flat File | http://www.ebi.ac.uk/intact | ftp://ftp.ebi.ac.uk/pub/databases/intact/current |
| IntEnz | XML | http://www.ebi.ac.uk/intenz/ | ftp://ftp.ebi.ac.uk/pub/databases/intenz/ |
| Pfam | Flat File | http://www.sanger.ac.uk/Software/Pfam/ | ftp://ftp.sanger.ac.uk/pub/databases/Pfam/ |
| UniProt | Fasta, Flat File | http://www.expasy.ch/ | ftp://ftp.expasy.org/ |
| KEGG | XML | http://www.genome.ad.jp/kegg/ | Web Services, ftp://ftp.genome.jp/pub/kegg/ |
| PDB | pdb Flat file, mmCIF, XML | http://www.rcsb.org/pdb | Web Services, ftp://ftp.wwpdb.org/ |

# Sequence Similarity

# Sequence Alignment

- Definition: Procedure for comparing two or more sequences by searching for a series of individual characters or character patterns that are *in the same order* in the sequences
  - **Pair-wise alignment**: compare two sequences
  - **Multiple sequence alignment**: compare more than two sequences

# Bioinformatics

**Stuart M. Brown, Ph.D.**
**NYU School of Medicine**

The next step is obviously to locate all of the genes and describe their functions. This will probably take another 15-20 years!

# Similarity Searching the Databanks

- What is similar to my sequence?

- Searching gets harder as the databases get bigger - and quality degrades

- Tools: BLAST and FASTA = time saving heuristics (approximate)

- Statistics + informed judgement of the biologist

# BLAST Algorithm



**(1)** For the query, find the list of high scoring words of length w

Query Sequence of length L

Maximum of L-w+1 words (typically w = 3 for proteins)

For each word from the query sequence find the list of words that will score at least T when scored using a pair-score matrix (e.g. PAM 250).

**(2)** Compare the word list to the database and identify exact matches

Word List

Database Sequences

Exact matches of words from word list

>gb|BE588357.1|BE588357 194087 BARC 5BOV Bos taurus cDNA 5'.
 Length = 369
 Score =  272 bits (137), Expect = 4e-71
 Identities = 258/297 (86%), Gaps = 1/297 (0%)
 Strand = Plus / Plus


Query: 17   aggatccaacgtcgctccagctgctcttgacgactccacagataccccgaagccatggca 76
            ||||||||||||||| | ||| | ||| || ||| | ||||  ||||| ||||||||
Sbjct: 1    aggatccaacgtcgctgcggctacccttaaccact-cgcagaccccccgcagccatggcc 59

.

Query: 77   agcaagggcttgcaggacctgaagcaacaggtggaggggaccgcccaggaagccgtgtca 136
            ||||||||||||||||||||||||| | || |||||||||| | ||||||||||| ||| ||
Sbjct: 60   agcaagggcttgcaggacctgaagaagcaagtggaggggggcggcccaggaagcggtgaca 119

.

Query: 137  gcggccggagcggcagctcagcaagtggtggaccaggccacagaggcggggcagaaagcc 196
             ||||||| | || | |||||||||||||| |||||||||| || |||||||||||||
Sbjct: 120  tcggccggaacagcggttcagcaagtggtggatcaggccacagaagcagggcagaaagcc 179

.

Query: 197  atggaccagctggccaagaccacccaggaaaccatcgacaagactgctaaccaggcctct 256
            |||||||| | |||||||| |||||||||||||||||||| ||||||||||||||||||||
Sbjct: 180  atggaccaggttgccaagactacccaggaaaccatcgaccagactgctaaccaggcctct 239

.

Query: 257  gacaccttctctgggattgggaaaaaattcggcctcctgaaatgacagcagggagac 313
            || || ||||| ||  |||||||||| | |||||||||||||||||||| ||||||||
Sbjct: 240  gagactttctcgggtttttgggaaaaaacttggcctcctgaaatgacagaagggagac 296

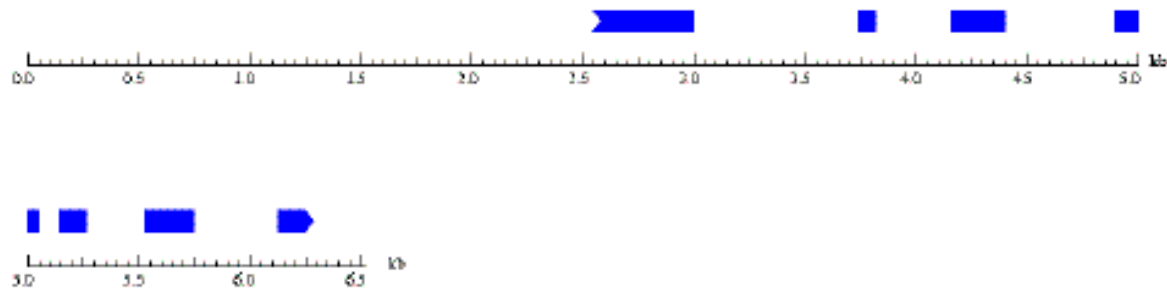# Finding Genes in genome Sequence is Not Easy

- About 1% of human DNA encodes functional genes.

- Genes are interspersed among long stretches of non-coding DNA.

- Repeats, pseudo-genes, and introns confound matters

# Pattern Finding Tools

- It is possible to use DNA sequence patterns to predict genes:
  - promoters
  - translational start and stop codes (ORFs)
  - intron splice sites
  - codon bias

- Can also use similarity to known genes/ESTs

GENSCAN predicted genes in sequence HSKER101

# Alignment

- Alignment is the basis for finding similarity

- Pairwise alignment = dynamic programming

- Multiple alignment: protein families and functional domains

- Multiple alignment is "impossible" for lots of sequences

- Another heuristic - progressive pairwise alignment

# Sample Multiple Alignment

# Structure- Function Relationships

- Can we predict the function of protein molecules from their sequence?

## sequence > structure > function

- Conserved functional domains = motifs

- Prediction of some simple 3-D structures ($\alpha$-helix, $\beta$-sheet, membrane spanning, etc.)

# Protein domains
## (from ProDom database)

# Main Method for Pairwise Alignment

- Word or $k$-tuple methods (FASTA and BLAST)

# Sample Multiple Alignment

# Examples

## I

"Once upon a midnight dreary, while I **ponder**ed, weak and weary,
Over many a quaint and curious volume of **forg**otten **lore**,
While I nodded, nearly **napping**, sudd**enly** there **came** a **tapping**,
As of some one ge**ntly rapping, rapping at my chamber door**.
"'Tis some visitor," I muttered, "**tapping at my chamber door-**
Only this, **and nothing more**."

## IV

"Presently my soul grew **stronger**; hesitating then no **longer**,
"Sir," said I, "or Madam, truly your **forg**iveness I imp**lore**;
But the fact is I was **napping**, and so g**ently** you **came** r**apping**,
And so fai**ntly** you came **tapping, tapping at my chamber door**,
That I scarce was sure I heard you"- here I opened wide the door;-
Darkness there, **and nothing more**. "

# Examples (Cont...)

...I pondered ...                       (I)

...stronger...

...of **forg**ott**en**--- - ---**lore**         (II)

your **forg**iv-**en**ess I im**p**lore

...**napping** sud -    **den-ly** there **came** a t**apping,** (III)

...**napping** and so ge**ntly** you-- **came** - r**apping**

# Examples (Cont...)

As of some one gently --- ---- rapping rapping at my chamber door (IV)
An d- so-- --f aintly you came tapping tapping at my chamber door

...     I muttered **tapping at my chamber door**     (IV')
... came tapping **tapping at my chamber door**

...**rapping** rapping at my chamber door     (IV'')
...tapping tapping at my chamber door
...-------- tapping at my chamber door

# Why do sequence alignments?

- To find out whether homologs of this gene (protein) are already available, and if they are, what is known about them

- To find whether two (or more) genes or proteins are evolutionarily related to each other

- To find structurally or functionally similar regions within proteins

# Origin of similar genes

- Similar genes arise by **gene duplication**
- Copy of a gene inserted next to the original
- Two copies mutate independently
- Each can take on separate functions
- All or part can be transferred from one part of genome to another

# Example sequence alignment

- Task: align **"abcdef"** with **"abdgf"**
- Write second sequence below the first

  **abcdef**

  **abdgf**

- Move sequences to give maximum match between them
- Show characters that match using vertical bar

# Example sequence alignment

**abcdef**

||

**abdgf**

- Insert gap between **b** and **d** on lower sequence to allow **d** and **f** to align

# Example sequence alignment

**abcdef**

|| |

**ab-dgf**

# Example sequence alignment

**abcdef**

|| | |

**ab-dgf**

- Note **e** and **g** don't match

# An alignment of two sequences t and s must satisfy:

- All symbols (residues) in the two sequences have to be in the alignment, and in the same order they appear in the sequences
- We can align one symbol from one sequence with one from the another
- A symbol can be aligned with a blank ('-')

- Two blanks cannot be aligned
- t: c g g g t a t c c a a
- s: c c c t a g g t c c c a

- t: c g g g t a - - t - c c a a
- s: c c c - t a g g t c c c - a

# Matching Similarity vs. Identity

- Alignments can be based on finding only identical characters, or (more commonly) can be based on finding *similar* characters

- More on how to define *similarity* later

# Global vs. Local Alignment

- We distinguish
  - **Global** alignment algorithms which optimize *overall* alignment between two sequences
  - **Local** alignment algorithms which seek only relatively *conserved* pieces of sequence
    - Alignment stops at the ends of regions of strong similarity
    - Favors finding conserved patterns in otherwise different pairs of sequences

# Global vs. Local Alignment

- Global

    **LGPSSKQTGKGS-SRIWDN**

    **| | ||| | |**

    **LN-ITKSAGKGAIMRLGDA**

- Local

    **--------GKG--------**

    **|||**

    **--------GKG--------**

# Global vs. Local Alignment

- Global

  **LGPSSKQTGKGS-SRIWDN**

  | | ||| | |

  **LN-ITKSAGKGAIMRLGDA**

- Local

  **-------TGKG---------**

  |||

  **-------AGKG---------**

# Sequence **FASTA** Format

- In the process of writing a similarity searching program (in 1985), William Pearson designed a simple text format for DNA and protein sequences

- The FASTA format is now universal for all databases and software that handles DNA and protein sequences

**One header line, starts with > with a [return] at end**

All other characters are part of sequence.
Most software ignores spaces, carriage returns.
Some ignores numbers

```
>URO1 uro1.seq  Length: 2018  November 9, 2000 11:50  Type: N  Check: 3854  ..
CGCAGAAAGAGGAGGCGCTTGCCTTCAGCTTGTGGGAAATCCCGAAGATGGCCAAAGAC
A
ACTCAACTGTTCGTTGCTTCCAGGGCCTGCTGATTTTTGGAAATGTGATTATTGGTTGTT
GCGGCATTGCCCTGACTGCGGAGTGCATCTTCTTTGTATCTGACCAACACAGCCTCTACC
CACTGCTTGAAGCCACCGACAACGATGACATCTATGGGGCTGCCTGGATCGGCATATTTG
TGGGCATCTGCCTCTTCTGCCTGTCTGTTCTAGGCATTGTAGGCATCATGAAGTCCAGCA
GGAAAATTCTTCTGGCGTATTTCATTCTGATGTTTATAGTATATGCCTTTGAAGTGGCAT
CTTGTATCACAGCAGCAACACAACAAGACTTTTTCACACCCAACCTCTTCCTGAAGCAGA
TGCTAGAGAGGTACCAAAACAACAGCCCTCCAAACAATGATGACCAGTGGAAAAACAATG
```
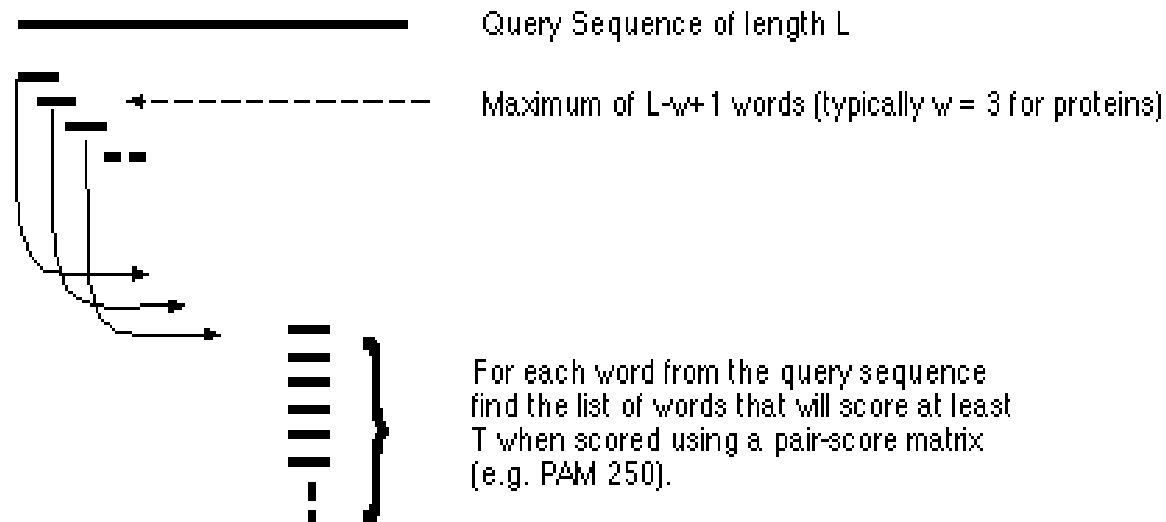
# Multi-Sequence FASTA file

>FBpp0074027 type=protein; loc=X:complement(16159413..16159860,16160061..16160497); ID=FBpp0074027; name=CG12507-PA;
    parent=FBgn0030729,FBtr0074248; dbxref=FlyBase:FBpp0074027,FlyBase_Annotation_IDs:CG12507
    PA,GB_protein:AAF48569.1,GB_protein:AAF48569; MD5=123b97d79d04a06c66e12fa665e6d801; release=r5.1; species=Dmel; length=294;
MRCLMPLLLANCIAANPSFEDPDRSLDMEAKDSSVVDTMGMGMGVLDPTQ
PKQMNYQKPPLGYKDYDYYLGSRRMADPYGADNDLSASSAIKIHGEGNLA
SLNRPVSGVAHKPLPWYGDYSGKLLASAPPMYPSRSYDPYIRRYDRYDEQ
YHRNYPQYFEDMYMHRQRFDPYDSYSPRIPQYPEPYVMYPDRYPDAPPLR
DYPKLRRGYIGEPMAPIDSYSSSKYVSSKQSDLSFPVRNERIVYYAHLPE
IVRTPYDSGSPEDRNSAPYKLNKKKIKNIQRPLANNSTTYKMTL
>FBpp0082232 type=protein; loc=3R:complement(9207109..9207225,9207285..9207431); ID=FBpp0082232; name=mRpS21-PA;
    parent=FBgn0044511,FBtr0082764; dbxref=FlyBase:FBpp0082232,FlyBase_Annotation_IDs:CG32854-
    PA,GB_protein:AAN13563.1,GB_protein:AAN13563; MD5=dcf91821f75ffab320491d124a0d816c; release=r5.1; species=Dmel; length=87;
MRHVQFLARTVLVQNNNVEEACRLLNRVLGKEELLDQFRRTRFYEKPYQV
RRRINFEKCKAIYNEDMNRKIQFVLRKNRAEPFPGCS
>FBpp0091159 type=protein; loc=2R:complement(2511337..2511531,2511594..2511767,2511824..2511979,2512032..2512082); ID=FBpp0091159; name=CG33919-
    PA; parent=FBgn0053919,FBtr0091923; dbxref=FlyBase:FBpp0091159,FlyBase_Annotation_IDs:CG33919-
    PA,GB_protein:AAZ52801.1,GB_protein:AAZ52801; MD5=c91d880b654cd612d7292676f95038c5; release=r5.1; species=Dmel; length=191;
MKLVLVVLLGCCFIGQLTNTQLVYKLKKIECLVNRTRVSNVSCHVKAINW
NLAVVNMDCFMIVPLHNPIIRMQVFTKDYSNQYKPFLVDVKIRICEVIER
RNFIPYGVIMWKLFKRYTNVNHSCPFSGHLIARDGFLDTSLLPPFPQGFY
QVSLVVTDTNSTSTDYVGTMKFFLQAMEHIKSKKTHNLVHN
>FBpp0070770 type=protein; loc=X:join(5584802..5585021,5585925..5586137,5586198..5586342,5586410..5586605); ID=FBpp0070770; name=cv-PA;
    parent=FBgn0000394,FBtr0070804; dbxref=FlyBase:FBpp0070770,FlyBase_Annotation_IDs:CG12410-
    PA,GB_protein:AAF46063.1,GB_protein:AAF46063; MD5=0626ee34a518f248bbdda11a211f9b14; release=r5.1; species=Dmel; length=257;
MEIWRSLTVGTIVLLAIVCFYGTVESCNEVVCASIVSKCMLTQSCKCELK
NCSCCKECLKCLGKNYEECCSCVELCPKPNDTRNSLSKKSHVEDFDGVPE
LFNAVATPDEGDSFGYNWNVFTFQVDFDKYLKGPKLEKDGHYFLRTNDKN
LDEAIQERDNIVTVNCTVIYLDQCVSWNKCRTSCQTTGASSTRWFHDGCC
ECVGSTCINYGVNESRCRKCPESKGELGDELDDPMEEEMQDFGESMGPFD
GPVNNNY
…

# BLAST Algorithm

**(1)** For the query, find the list of high scoring words of length w

Query Sequence of length L

Maximum of L-w+1 words (typically w = 3 for proteins)

For each word from the query sequence find the list of words that will score at least T when scored using a pair-score matrix (e.g. PAM 250).

**(2)** Compare the word list to the database and identify exact matches

Word List

Database Sequences

Exact matches of words from word list

```
>gb|BE588357.1|BE588357 194087 BARC 5BOV Bos taurus cDNA 5'.
 Length = 369
 Score =  272 bits (137), Expect = 4e-71
 Identities = 258/297 (86%), Gaps = 1/297 (0%)
 Strand = Plus / Plus


Query: 17   aggatccaacgtcgctccagctgctcttgacgactccacagataccccgaagccatggca 76
            |||||||||||||||| | ||| | ||| || ||| | |||| ||||| ||||||||
Sbjct: 1    aggatccaacgtcgctgcggctacccttaaccact-cgcagaccccccgcagccatggcc 59
.
Query: 77   agcaagggcttgcaggacctgaagcaacaggtggaggggaccgcccaggaagccgtgtca 136
            ||||||||||||||||||||||||| | || ||||||||| | |||||||||| ||| ||
Sbjct: 60   agcaagggcttgcaggacctgaagaagcaagtggagggggcggcccaggaagcggtgaca 119
.
Query: 137  gcggccggagcggcagctcagcaagtggtggaccaggccacagaggcggggcagaaagcc 196
             |||||||| | || | |||||||||||||| |||||||||| || |||||||||||||
Sbjct: 120  tcggccggaacagcggttcagcaagtggtggatcaggccacagaagcagggcagaaagcc 179
.
Query: 197  atggaccagctggccaagaccacccaggaaaccatcgacaagactgctaaccaggcctct 256
            ||||||||| | |||||||| ||||||||||||||||||| |||||||||||||||||||
Sbjct: 180  atggaccaggttgccaagactacccaggaaaccatcgaccagactgctaaccaggcctct 239
.
Query: 257  gacaccttctctgggattgggaaaaaattcggcctcctgaaatgacagcagggagac 313
            || || ||||| ||  ||||||||||| |||||||||||||||||| ||||||||
Sbjct: 240  gagactttctcgggttttgggaaaaaacttggcctcctgaaatgacagaagggagac 296
```

# Two classes of widely used protein scoring matrices

PAM = % Accepted Mutations:
1500 changes in 71 groups w/ > 85% similarity

BLOSUM = Blocks Substitution Matrix:
2000 "blocks" from 500 families

```
>mysequence1
atggaggatgatttcatgtgcgatgatgaggaggactacgacctggaatactctga
agatagtaactccgagccaaatgtggatttggaaaatcagtactataattccaaag
cattaaaagaagatgacccaaaagcggcattaagcagtttccaaaaggttttggaa
cttgaaggtgaaaaaggagaatggggatttaaagcactgaaacaaatgattaagat
taacttcaagttgacaaactttccagaaatgatgaatagatataagcagctattga
cctatattcggagtgcagtcacaagaaattattctgaaaaatccattaattctatt
cttgattatatctctacttctaaacagatggatttactgcaggaattctatgaaac
aacactggaagctttgaaagatgctaag
```

Use **Blast** - there are different varieties, depending on what kind of
sequence you have and what kind of sequence you are looking for

| | |
|---|---|
| blastn | Search nucleotide database using a nucleotide query |
| blastp | Search protein database using a protein query |
| blastx | Search protein database using a translated nucleotide query |
| tblastn | Search translated nucleotide database using a protein query |
| tblastx | Search translated nucleotide database using a translated nucleotide query |

# BLAST
*Basic Local Alignment Search Tool*

| Home | Recent Results | Saved Strategies | Help |

▸ NCBI/ BLAST/ blastn suite

| blastn | blastp | blastx | tblastn | tblastx |

## Enter Query Sequence

BLASTN programs search nucleotide databases using a nucleotide query. more...

**Enter accession number(s), gi(s), or FASTA sequence(s)** ❓          Clear

**Query subrange** ❓

From [          ]

To [          ]

**Or, upload file**        [          ] Browse... ❓

**Job Title**        [                                        ]

Enter a descriptive title for your BLAST search ❓

☐ **Align two or more sequences** ❓

## Choose Search Set

**Database**        ◉ Human genomic + transcript   ○ Mouse genomic + transcript   ○ Others (nr etc.):

[ Human genomic plus transcript (Human G+T)    ▾ ] ❓

**Exclude**
Optional        ☐ Models (XM/XP)  ☐ Uncultured/environmental sample sequences

**Entrez Query**
Optional        [                                        ]

Enter an Entrez query to limit search ❓

## Program Selection

**Optimize for**        ◉ Highly similar sequences (megablast)

○ More dissimilar sequences (discontiguous megablast)

○ Somewhat similar sequences (blastn)

▼ Algorithm parameters

## General Parameters

**Max target sequences**   `100 ▼`
Select the maximum number of aligned sequences to display ⑦

**Short queries**   ☑ Automatically adjust parameters for short input sequences ⑦

**Expect threshold**   `10` ⑦

**Word size**   `28 ▼` ⑦
The length of the seed that initiates an alignment. more...

**Max matches in a query range**   `0` ⑦

## Scoring Parameters

**Match/Mismatch Scores**   `1,-2 ▼` ⑦

**Gap Costs**   `Linear ▼` ⑦

## Filters and Masking

**Filter**   ☐ Low complexity regions ⑦
☐ Species-specific repeats for: `Human ▼` ⑦

**Mask**   ☑ Mask for lookup table only ⑦
☐ Mask lower case letters ⑦

**BLAST**   Search **database Human G+T** using **Megablast (Optimize for highly similar sequences)**
☐ Show results in a new window

# How to read a BLAST result

**score** – indicates how similar the query sequence is to the results,
**larger number is better**
BUT: longer sequences lead to higher scores

**e-value** – expectation val
how often would you exp
to find this sequence in th
database randomly (this i
particularly relevant if you
query sequence is short
contains many repeats, e
**smaller number is bette**
Note: 2e-3 = 2*10⁻³ = 0.0

**gene id** of hit (accession number)

description

Legend for links to other resources: **U** UniGene **E** GEO **G** Gene **S** Structure **M** Map Viewer

Sequences producing significant alignments:
(Click headers to sort columns)

| Accession | Description | Max score | Total score | Query coverage | E value | Max iden |
|---|---|---|---|---|---|---|
| NM_004236.2 | Homo sapiens COP9 constitutive photomorphogenic homolog subunit 2 (Arabidopsis) (COPS2), mRNA | 457 | 457 | 100% | 1e-125 | 100% |
| XM_001166766.1 | PREDICTED: Pan troglodytes similar to COP9 complex subunit 2, transcript variant 1 (LOC453417), mRN. | 457 | 457 | 100% | 1e-125 | 100% |
| XM_510388.2 | PREDICTED: Pan troglodytes similar to COP9 complex subunit 2, transcript variant 2 (LOC453417), mRN. | 457 | 457 | 100% | 1e-125 | 100% |
| BC012629.1 | Homo sapiens COP9 constitutive photomorphogenic homolog subunit 2 (Arabidopsis), mRNA (cDNA clone | 457 | 457 | 100% | 1e-125 | 100% |
| AK222590.1 | Homo sapiens mRNA for COP9 constitutive photomorphogenic homolog subunit 2 variant, clone: CAS050 | 457 | 457 | 100% | 1e-125 | 100% |
| AB209799.1 | Homo sapiens mRNA for COP9 constitutive photomorphogenic homolog subunit 2 variant protein | 457 | 457 | 100% | 1e-125 | 100% |
| AF212227.1 | Homo sapiens TRIP15-ISO mRNA, complete cds | 457 | 457 | 100% | 1e-125 | 100% |
| CR614722.1 | full-length cDNA clone CS0DI070YA02 of Placenta Cot 25-normalized of Homo sapiens (human) | 457 | 457 | 100% | 1e-125 | 100% |
| CR601131.1 | full-length cDNA clone CS0DA011YC02 of Neuroblastoma of Homo sapiens (human) | 457 | 457 | 100% | 1e-125 | 100% |
| AF084260.1 | Homo sapiens signalosome subunit 2 (SGN2) mRNA, complete cds | 457 | 457 | 100% | 1e-125 | 100% |
| AF100762.1 | Homo sapiens thyroid receptor interactor trip15 mRNA, complete cds | 457 | 457 | 100% | 1e-125 | 100% |
| AF120268.1 | Homo sapiens ALIEN (ALIEN) mRNA, complete cds | 457 | 457 | 100% | 1e-125 | 100% |
| L40388.2 | Homo sapiens thyroid receptor interactor (TRIP15) mRNA, 5' end of cds | 457 | 457 | 100% | 1e-125 | 100% |

# An individual "BLAST hit" in more detail

accession number  + description

>ref|NM_004236.2| UG Homo sapiens COP9 constitutive photomorphogenic homolog subuni
2 (Arabidopsis) (COPS2), mRNA
Length=1947

score, e-value,
identical nucleotides,
gaps, orientation →

```
 Score =  457 bits (247),  Expect = 1e-125
 Identities = 247/247 (100%), Gaps = 0/247 (0%)
 Strand=Plus/Plus

Query  1    ATGGAGGATGATTTCATGTGCGATGATGAGGAGGACTACGACCTGGAATACTCTGAAGAT  60
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  21   ATGGAGGATGATTTCATGTGCGATGATGAGGAGGACTACGACCTGGAATACTCTGAAGAT  80
```

your sequence →
blast hit →
```
Query  61   AGTAACTCCGAGCCAAATGTGGATTTGGAAAATCAGTACTATAATTCCAAAGCATTAAAA  120
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  81   AGTAACTCCGAGCCAAATGTGGATTTGGAAAATCAGTACTATAATTCCAAAGCATTAAAA  140

Query  121  GAAGATGACCCAAAAGCGGCATTAAGCAGTTTCCAAAAGGTTTTGGAACTTGAAGGTGAA  180
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  141  GAAGATGACCCAAAAGCGGCATTAAGCAGTTTCCAAAAGGTTTTGGAACTTGAAGGTGAA  200

Query  181  AAAGGAGAATGGGGATTTAAAGCACTGAAACAAATGATTAAGATTAACTTCAAGTTGACA  240
            ||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct  201  AAAGGAGAATGGGGATTTAAAGCACTGAAACAAATGATTAAGATTAACTTCAAGTTGACA  260

Query  241  AACTTTC    247
            |||||||
Sbjct  261  AACTTTC    267
```

This is a perfect match!

# References

[1]   T. Etzold, A. Ulyanow, and P. Argos, "SRS: Information Retrieval System for Molecular Biology Data Banks," *Methods Enzymol.*, vol. 226, pp. 114–128, 1996.

[2] W. Fujibuchi, S. Goto, H. Migimatsu, I. Uchiyama, A. Ogiwara, M. Kanehisa, "DBGET/LinkDB: An Integrated Database Retrieval System," *Proc. Pacific Symp. Biocomputing*, pp. 681–692, 1998.

[3] C. Ramu, "SIR: A Simple Indexing and Retrieval System for Biological Flat File Databases," *Bioinformatics*, vol. 17, no. 8, pp. 756–758, 2001.

[4] L. M. Haas, P. Schwarz, M. Kodali, E.  Kotlar, J. E. Rice, and W. C. Swope, "DiscoveryLink: A System for Integrated Access to Life Sciences Data Sources," *IBM Syst. J.*, vol. 40, pp. 489–510, 2001.

[5] S. Chawathe, H. Garcia-Molina, J. Hammer, K. Ireland, Y. Papakonstantinou, J. Ullman, *et al*, "The TSIMMIS Project: Integration of Heterogeneous Information Sources," *Proceedings of IPSJ conference*, Tokyo, Japan; pp. 7–18, 1994.

[6] A. Freier, R. Hofestadt, M. Lange, U. Scholz, A. Stephanik, "BioDataServer: a SQL-based Service for the Online Integration of Life Science Data," *Silico Biol*, 2(2):37–57, 2002.

[7] R. D. Stevens, P. Baker, S. Bechhofer,  G. Jacoby, A. N.,  Ng, W. Paton, C. A.  Goble, and A. Brass, "TAMBIS: Transparent Access to Multiple Bioinformatics Information Sources," *Bioinformatics*, vol. 16, no. 2, pp. 184–186, 2000.

[8] J. Köhler, S. Philippi, and M. Lange, "SEMEDA – Ontology Based Integration of Biological Databases," *Bioinformatics*, vol. 19, no. 18, pp.2420–2427, 2003.

Thank YOU