# About the importance of HPC for Life Sciences

Oswaldo Trelles
University of Malaga, Spain
ortrelles@uma.es

# Computer Architecture Dept.

## Main Research Lines

- VLSI (ASIC integration)

- 'Parallel' Compilers

- Applications (sparse matrices $\rightarrow$ I-O bounded)

Bioinformatics: Computer sciences as applied to biological data



**BITLAB**: Bioinformatics and Information Technologies Laboratory

Basic & applied Research

**www.bitlab-es.com**

O.Trelles, PhD, 2014

# High Performance Computing applied to Life Sciences

**RISC**    **JKU**    **UMA**

## Improving open source software for high performance computing in Biology

### Bingos
(Bioinformatics *next generation* open software)

• **Problem**: new high throughput technologies in several areas of life sciences produce enormous amounts of data. A bottleneck in our ability to process and analyse the data is becoming apparent

• **Solution**: This Action aims to increase communication between bioinformatics, HPC and Open Source communities for adapting / developing HPC capable software tools

**www.bitlab-es.com**

Bingos Project
- Objectives
- Deliverables
- Contact
- Background
- Cientific Program
- Participants

links
- SHARE
- ELIXIR
- CCP4
- GALAXY
- ENFIN
- EGEE
- BIOSAPIENS
- Mancoosi
- FLOSS
- BioJava
- COST Action IC0702
- EMBRACE
- INB

News and meetings
- News
- Meetings

Expand all - Collapse all

USER: ____
PASSWORD: ____
Login!

Team-work makes the difference
www.bitlab.es

RISC, Linz 2010

# Targeting Big-Data problems in BI

1995:   1 US$ per base (3.000M US$ the full human genome)
2000:   1 Mbp ≈ 10.000 US $ [1]
2008:   Full human genome (3,2 Gbp) in 6 weeks, and ≈ $60,000 [2]
        Predicted: US$1,000 genome in next 3 years.
2009:   (October) nanopore DNA sequencing [3, 4]
2011:   (Mar. 2011) 0,5 US$ per Mbp [5]
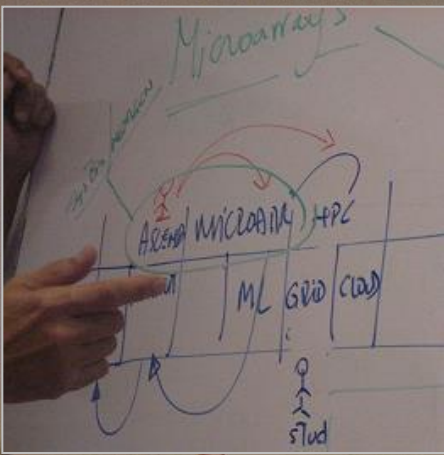2012:   (Feb. 2012) size: USB memory stick / $900 [6]



**Allow individual to get his or her genome sequenced, thus truly ushering in the era of genetics based personalized medicine.**

[1] http://www.genome.gov/11006943
[2] http://www.technologyreview.com/Biotech/20640/
[3] http://nextbigfuture.com/2009/10/ibm-targeting-100-dollar-genome.html
[4] http://www.nature.com/embor/journal/v8/n10/full/7401070.html
[5] http://singularityhub.com/2011/03/05/costs-of-dna-sequencing-falling-fast-look-at-these-graphs/
[6] http://www.nanoporetech.com/news/press-releases/view/39

O.Trelles, PhD, 2014

# The first draft…



O.Trelles, PhD, 2014

Andahuaylas, Apurímac   dreaming...

G. Helies, PhD, 2014

# Mr.SymBioMath



**Dissemination**

Software Applications

Grid + Cloud & HPC Infrastructure

Soft. Clients

Interfaces, visualization & **Data analysis**

Data collections

Models and Base soft: libraries

| | | | |
|---|---|---|---|
| RISC (Cloud) | HCH + UMA | JKU + RISC+ LRZ | LRZ + ITG + ICG |
| UMA (clients) | LNCC | UMA+ ITG | HCH+ LNCC (users) |
| IBM (Big Data) | PUBLIC | IBM + ICG+ BAOBAD | |

O.Trelles, PhD, 2014

# Mr.SymBioMath

THE COMPUTER ARCHITECTURE DEPARTMENT AT - THE - - UNIVERSITY - - OF MALAGA
| | | | | | | | | | | | | | | | | | | | . | | | |                    |        |        | | | | | | | | | |        | | | | | |
DIE COMPUTER ARCHITEKTURE ABTEILUNG - AN - - DER - UNIVERSITAT VON MALAGA
        | | . . | | | . | | |                    | |        | | | | | | | | |        |        | | | | | |
EL DEPARTAMENTO DE ARQUITECTURA DE COMPUTADORES DE LA UNIVERSIDAD - DE MALAGA

# Visualization & Interpretation



Low complexity zones (repetitions of the same symbol in both sequences

Palindromes Inverted diagonals

Dortplots for DNA sequences can be noisy since there are only 4 symbols (each symbol in one sequence will match qith the 25% of the symbols in the other sequence

To avoid noise, instead of compare pair of symbols an sliding windows is used, and a minimal threshold or stringency level is used to assign a real match (e.g W=10, T=6).

Residue deletion in the vertical sequence or insertion in the horizontal

Residue deletion in the horizontal or insertion in the vertical sequences

Repetitions: a zone in the horizontal sequence is similar to more than one in the vertical

O.Trelles, PhD, 2014

H. sapiens  P. troglodytes  M. mulatta  C. familiaris  R. novergicus  M. musculus  B. taurus

H. sapiens

P. troglodytes

M. mulatta

C. familiaris

R. novergicus

M. musculus

B. taurus

O. Trelles, PhD, 2014

# Comparative Genomics
## Detection (& sequence) of Evolution Events



O.Trelles, PhD, 2014

# Introduction

## Survey on biology and bioinformatics

O.Trelles, PhD, 2014

# From genes to pathways



O.Trelles, PhD, 2014

# Cells and organisms

All living things are made of cells.

Prokaryotic & Eukaryotic



Some typical cells
animal cell — cell membrane
plant and animal cells are eukaryotic cells
centriole — vacuole
centrosome — ribosomes
plant cell
plasma membrane — vacuole
chloroplast
endoplasmic reticulum
ribosomes
mitochondrion
nucleus
nucleolus
chromosomes
bacteria cell (bacillus type)
Golgi complex
cytoplasm
cell wall
chromosome
plasmodesma
ribosomes
cell wall
plasma membrane
capsule
pili
mesosome
flagella
bacterial cells are prokaryotic cells
© 2007 Encyclopædia Britannica, Inc.

O.Trelles, PhD, 2014
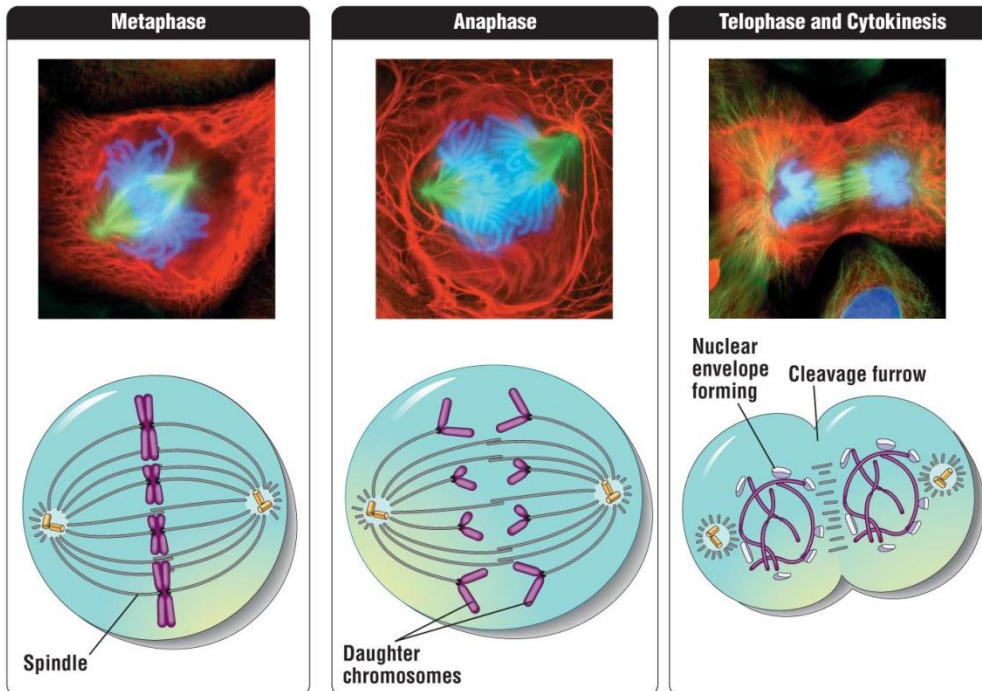
# Cells reproduction



Cell reproduction is the process of a cell splitting and becoming two similar cells.

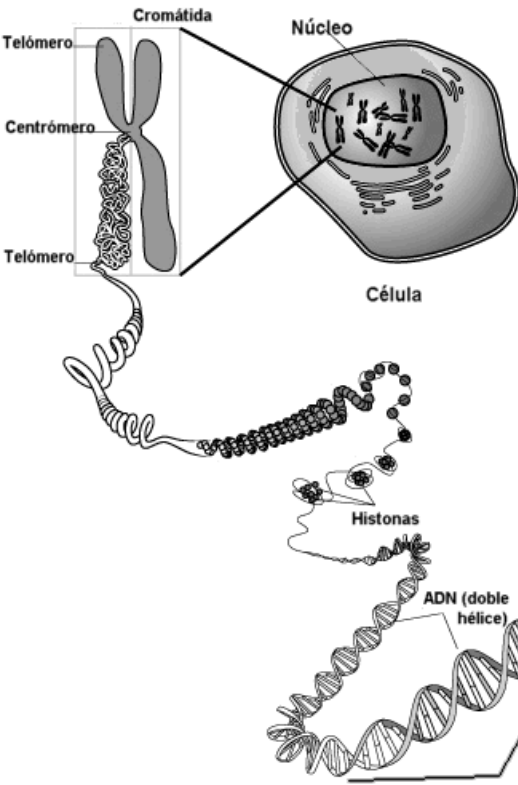Prokaryotes by **binary fission**

Eukaryotic cells reproduce using either **mitosis** (2) or **meiosis**. (4)

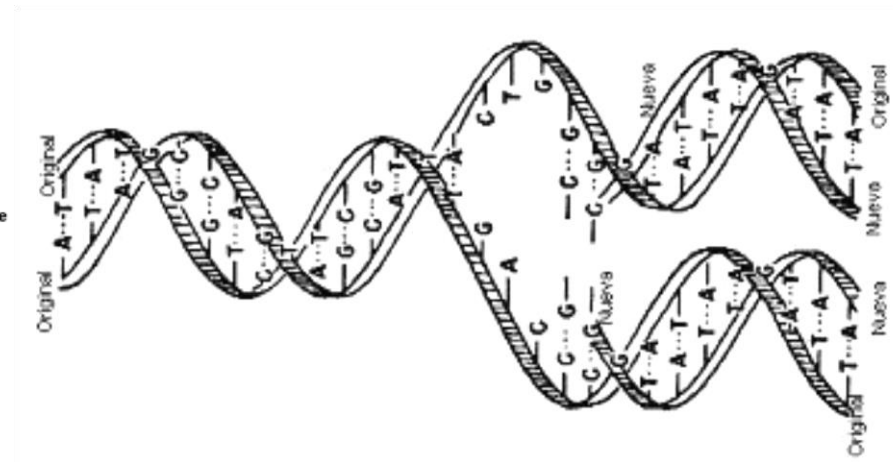daughter cells have identical genetic composition, except for spontaneous **mutations**.

**Metaphase**

**Anaphase**

**Telophase and Cytokinesis**

Nuclear envelope forming

Cleavage furrow

Spindle

Daughter chromosomes

O.Trelles, PhD, 2014

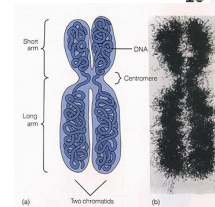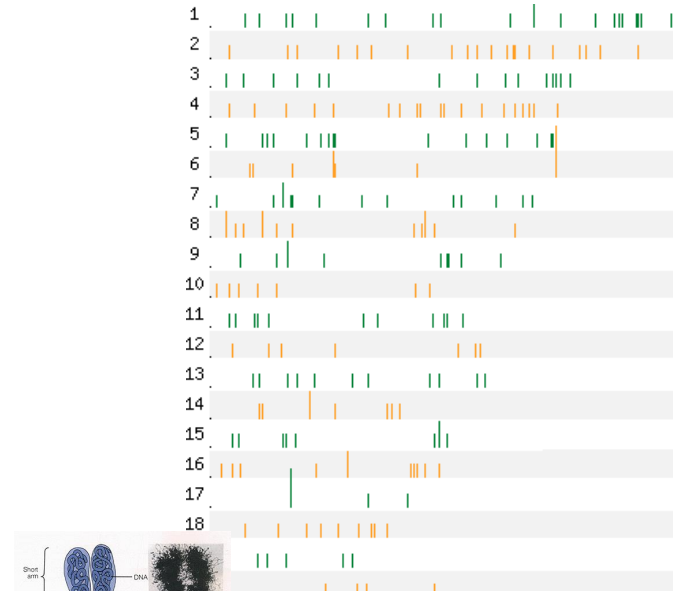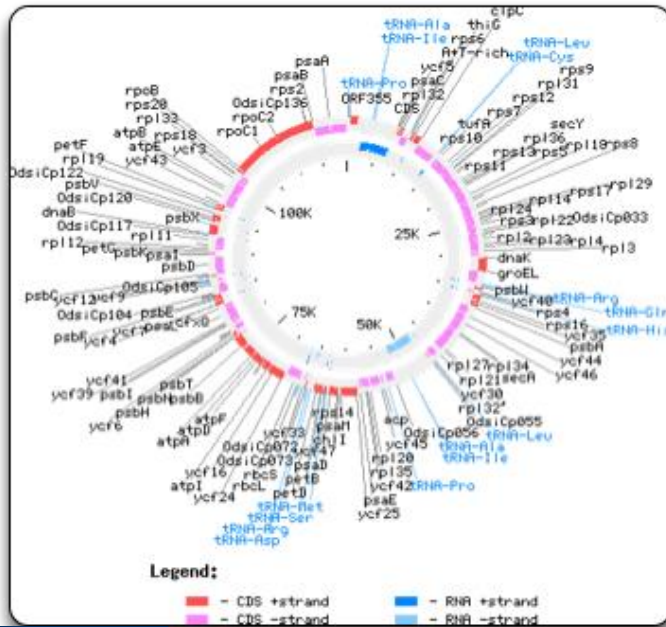# The DNA carries the hereditary information

DNA conforms a long linear polymer using 4 different molecules or monomers: Adenine, Cytosine, Guanine y Thymine (A, C, G and T) also called nucleotides or bases.

# Chromosomes and Genes
## DNA is organized in chromosomes
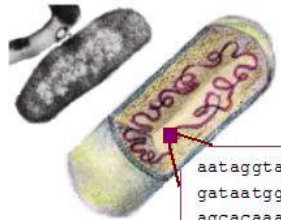## Genes carry out the instructions to synthetize proteins



Mouse genome: 20 pairs of chromosomes

Circular chromosome of *"Odeontella"* with 119,704 base pairs / 174 genes

source: http://chloroplast.ocean.washington.edu/chloroplast_files/images/odontella_genome.png
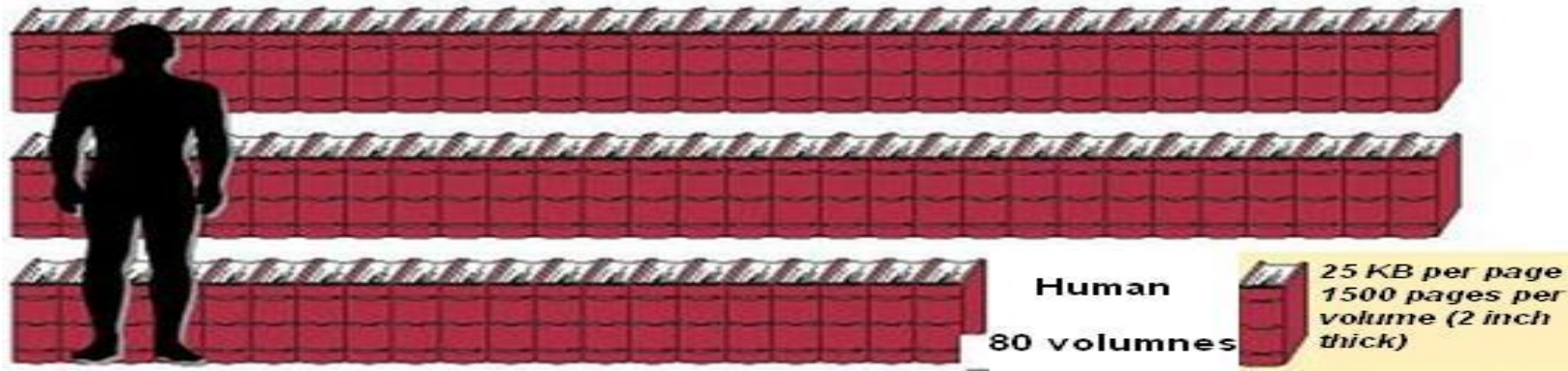
O.Trelles, PhD, 2014

# Genome size

The genome is replicated in each cell
Size: from few thousands of bases in bacteria (viruses?)
To about some GBp (basepairs)

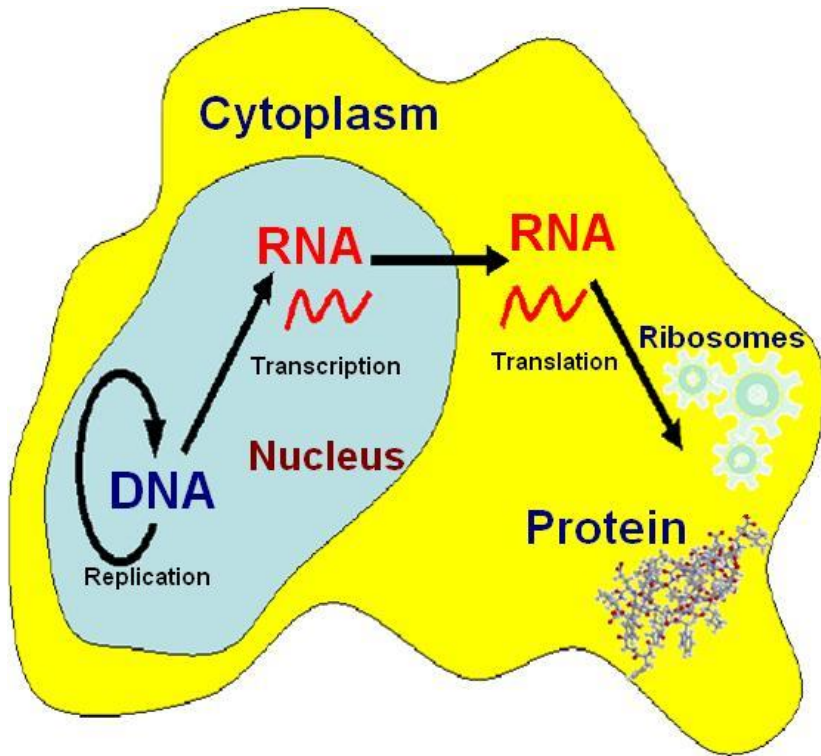```
aataggtaaatctacaacaacacaaaacttaacatcagggcttgctacaatggacaagaa
gataatggtagtaggatgtgatcctaaggctgactcaacaaggttattactaggaggact
agcacaaaaaagtgttcttgatacattaagagaagaaggagatgacgtagatttagattc
aatcttaaagccaggatttagaggtataaaatgtgttgaatcaggcggtccagaaccagg
agttggatgtgcaggaagaggtataataacttcaatcaatatgctagagcaattaggtgc
ttacgaatcagatttagattatgtttttctatgatgtgtattaggt
```

The genome is similar to a recipes book

| Lambda Bacteriophage 2 pages | Escherichia coli Bacteria 200 pages | Saccharomyces cerevisiae 500 pages | Arabidopsis thaliana or C. elegans 3 volumes | Drosophila melanogaster 5 volumes |

Human
80 volumnes

25 KB per page
1500 pages per volume (2 inch thick)

O.Trelles, PhD, 2014

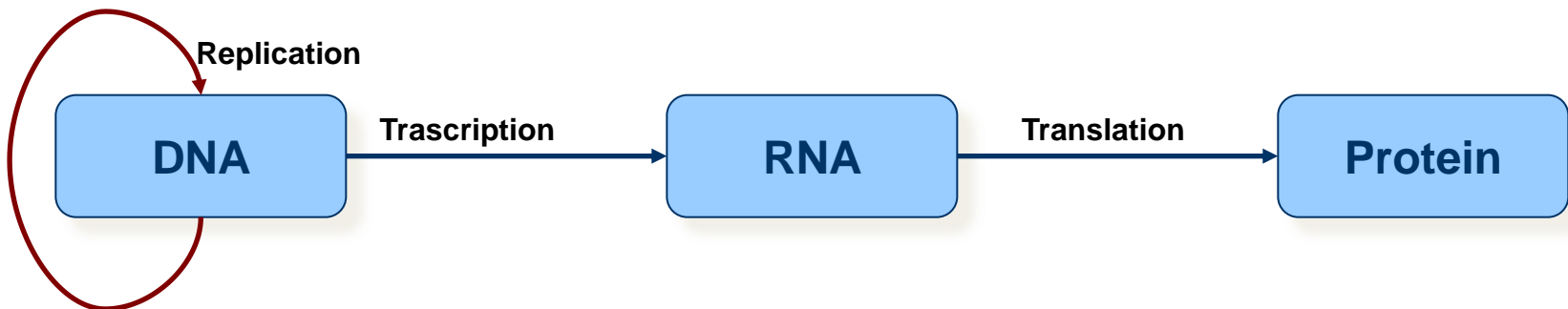# Central Dogma of molecular biology
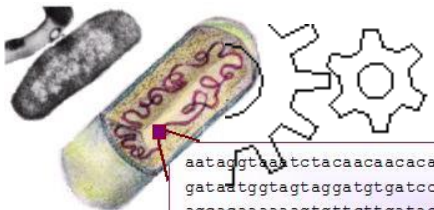


The **central dogma** of molecular biology states:

(1) DNA carries the genetic information of organisms and **replicates** during cell division to allow each daughter cell to contain a full complement of chromosomes.

(2) The genetic information in the DNA is used in a process called **transcription** to produce a complementary one-strand messenger of mRNA

(3) mRNA is interpreted (**translation**) in the ribosomes using the genetic-code to produce a protein.



O.Trelles, PhD, 2014

# *f*rom *G*enes to *P*roteins
## The Genetic Code
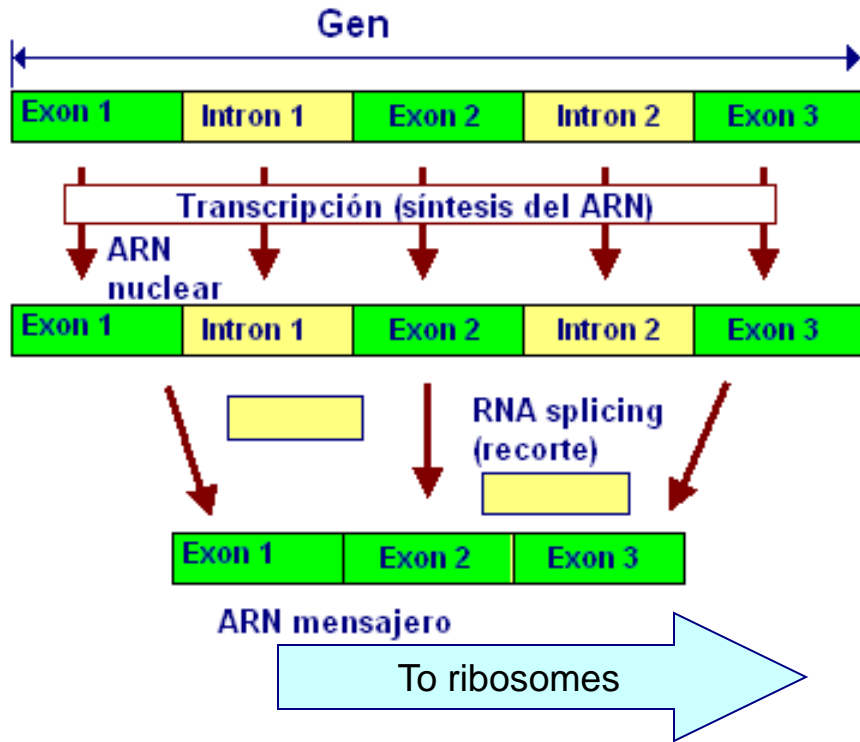


Genes contains the instructions for protein synthesis. That instructions are translated by the cellular machinery using the so called genetic code that translate each consecutive codon (DNA triple) into an specific amino acid

- **Codon**: 3 consecutive bases of DNA

- There are 6 (putative) different ways to read the DNA

- **ORF**: the frame o DNA with not stop codons

O.Trelles, PhD, 2014

# Some details….
## (*i.e.* eukaryotic genes)



Protein synthesis start with a copy of one of the DNA strands into RNA inside the cellular nucleus. This RNA is spliced to remove the introns (mature mRNA).
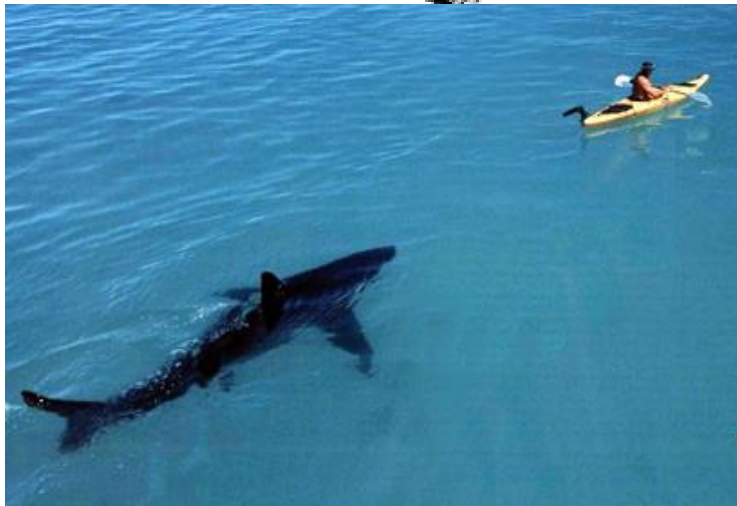
Small signals for starting (donors) of introns and exons and ending points (acceptors) are used to identify the right cutting position, including the stop signals for ending the translation.

Differences between pro and eukaryotic cells

O.Trelles, PhD, 2014

# *P*roteins

DNA: instructions to drive the synthesis of proteins

Proteins: organisms are made of proteins (bones, muscles, nervous…

Protein function is associated to its 3D spatial conformation

DNA is present in all cells

Proteins: each cell produces only those proteins the cell needs

sulphate

heparin

# $\mathcal{P}$roteins levels

## Throught Gene-Expression



Levels of proteins ⟷ Cellular state

Gene levels ≈ Protein levels

Env. stimuly ⟷ Change proteins levels

Change proteins levels ⟵ Change gene levels

⟶ Gene regulation mchanisms:

Changes in protein levels have profund effect in the biology of the organisms (even with phisiological and pathologic effects)
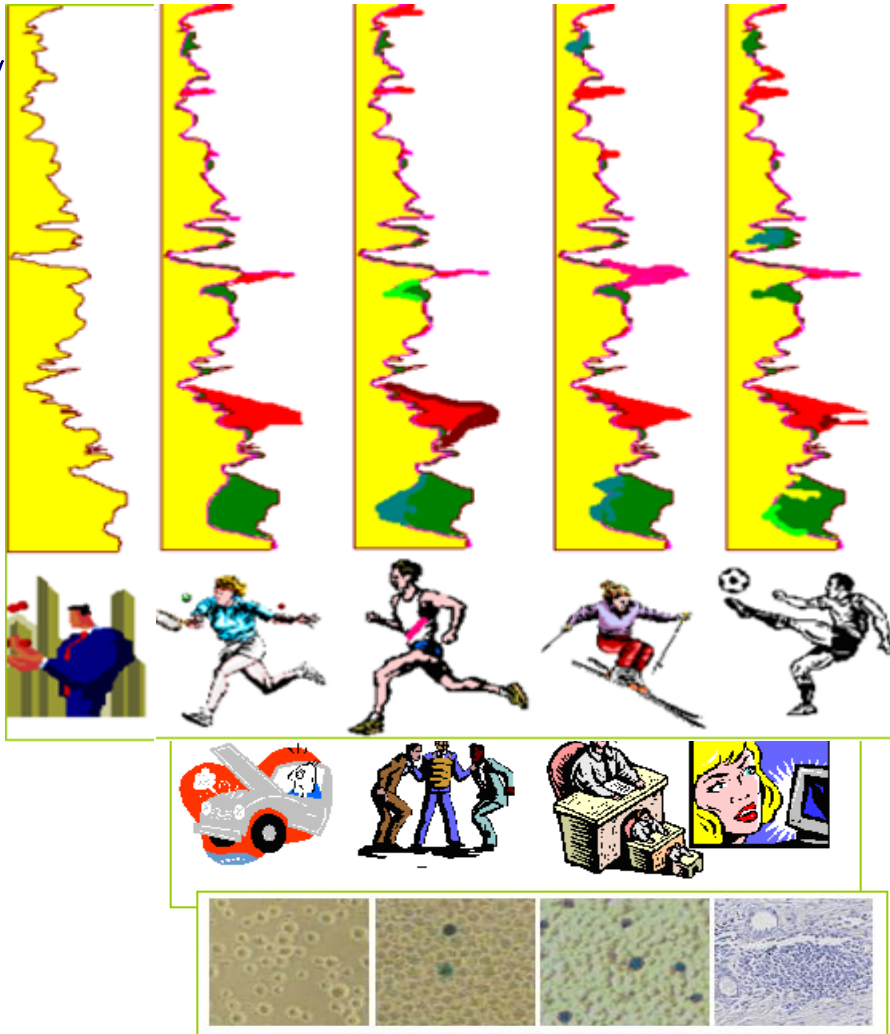
Gene-expression levels are used to determine the response of an organism to a particular event

Different developmental stage, tissues types, clinical conditions, organisms, etc

# Gene Expression

## Quantify the level at which a particular gene is expressed



Gene's catalogue

GE levels are used to determine the response of an organism to a particular event

Simultaneosus anaysis of thousands of gene

Different developmental stage, tissues types, clinical conditions, organisms, etc
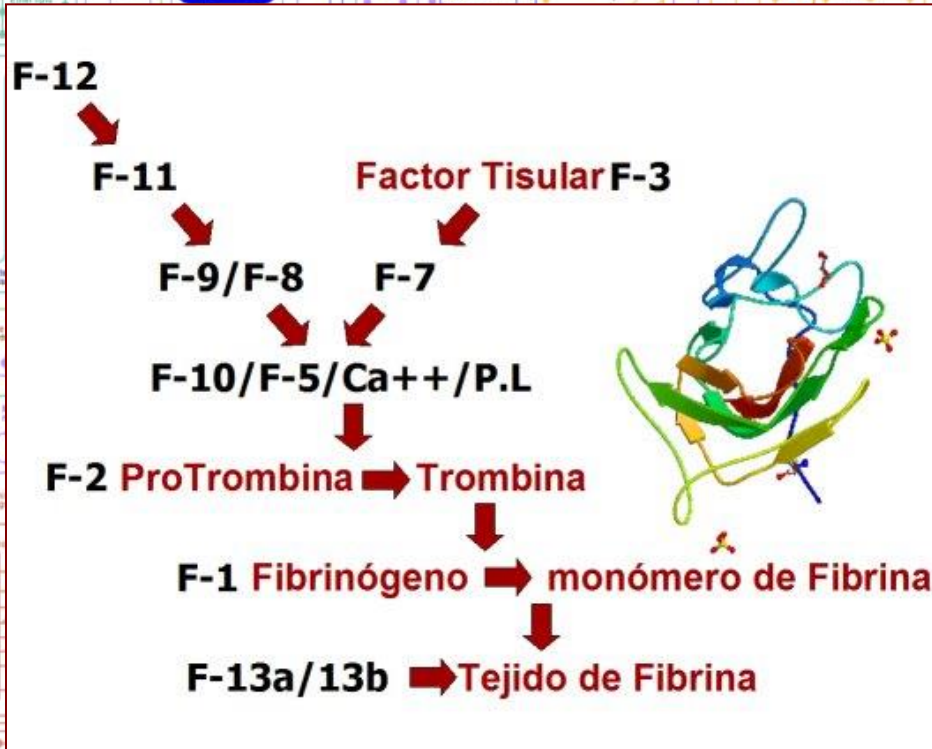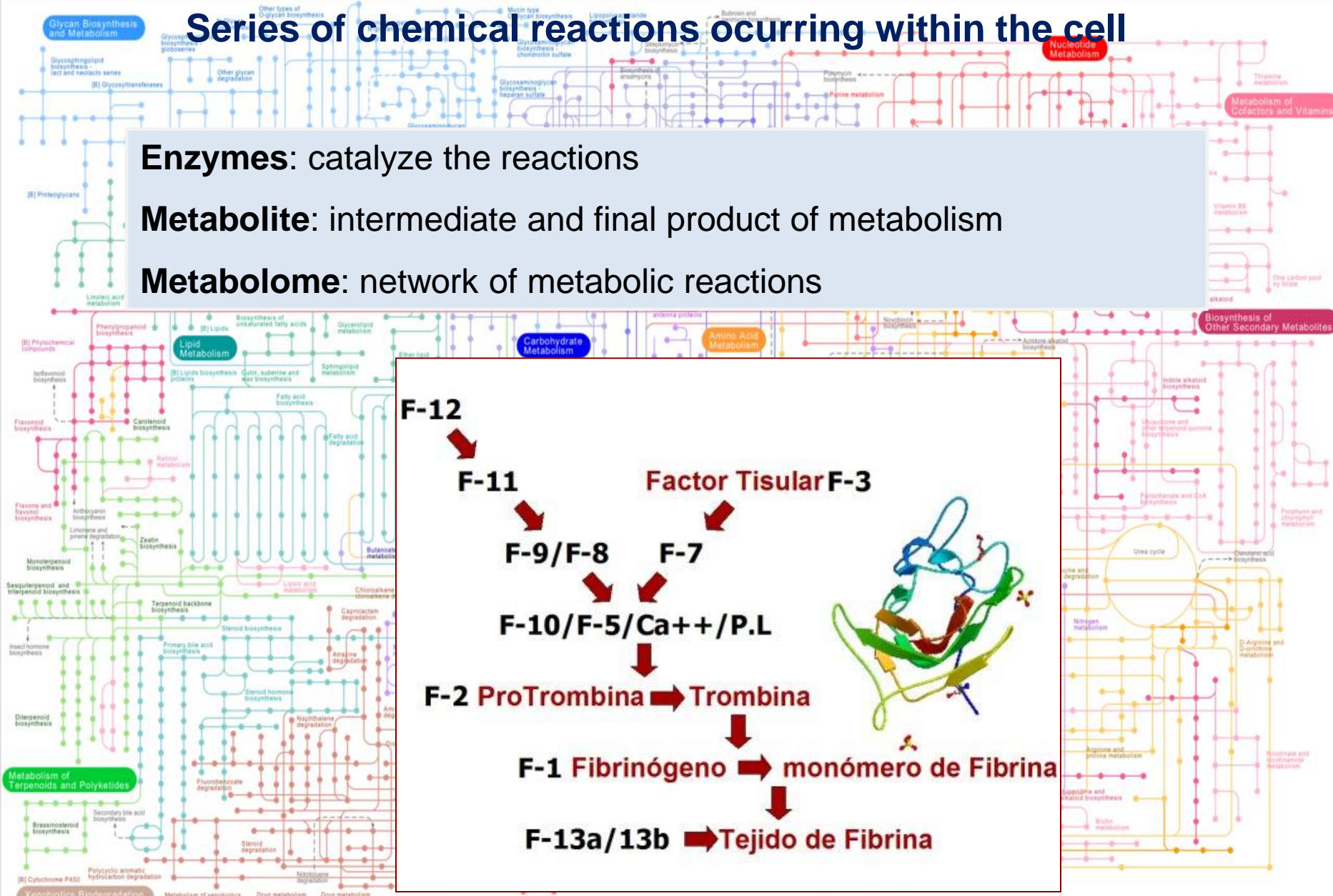
# Metabolic *P*athways

## Series of chemical reactions ocurring within the cell

**Enzymes**: catalyze the reactions

**Metabolite**: intermediate and final product of metabolism

**Metabolome**: network of metabolic reactions

# Bioinformatics



Source: ECCC'02 Web site

## Featuring the application domain

O.Trelles, PhD, 2014

# *B*ioinformatics (Computational biology)

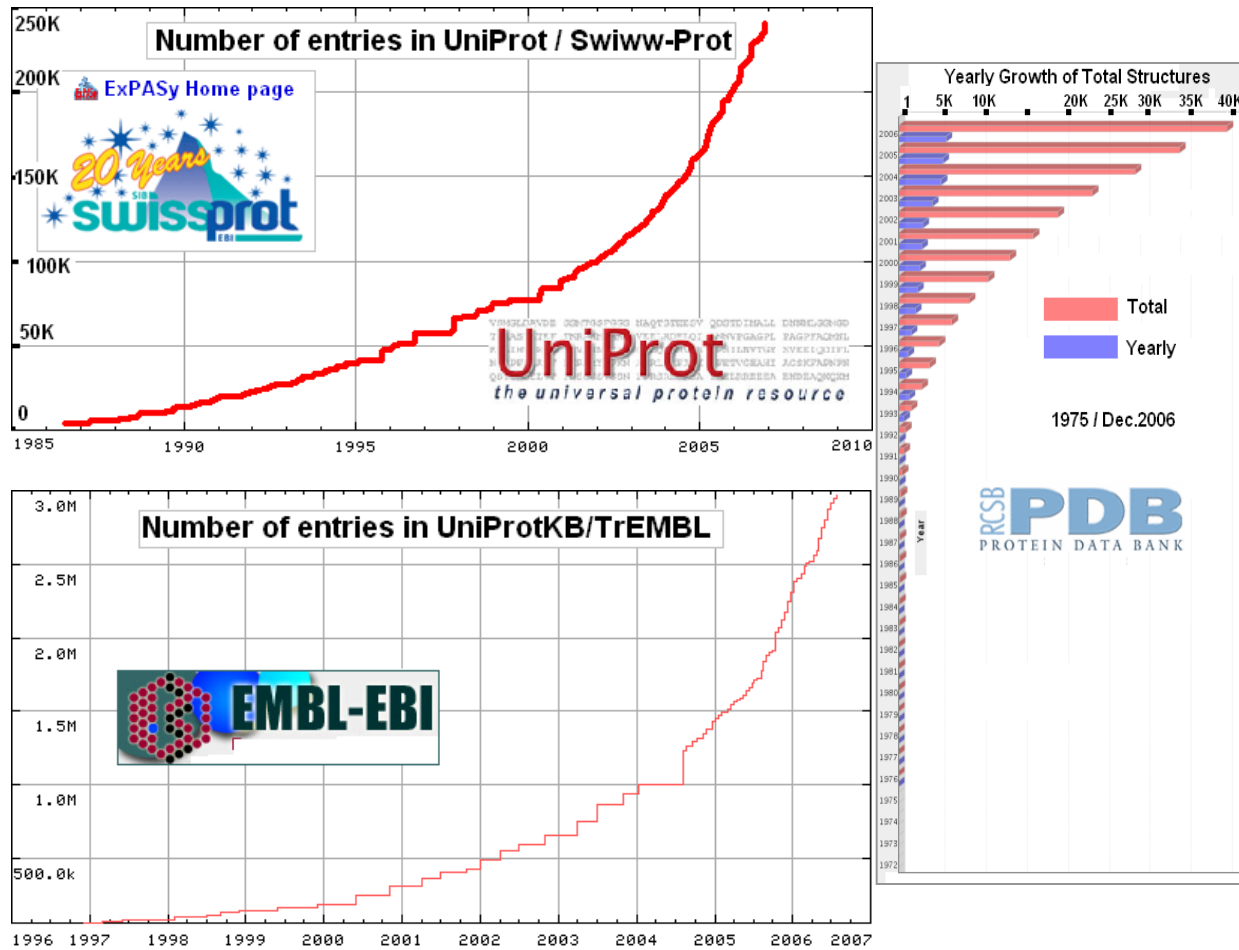## Computer sciences as applied to biological data

Computer sciences, statistics, physiscs, chemistry, IT, …

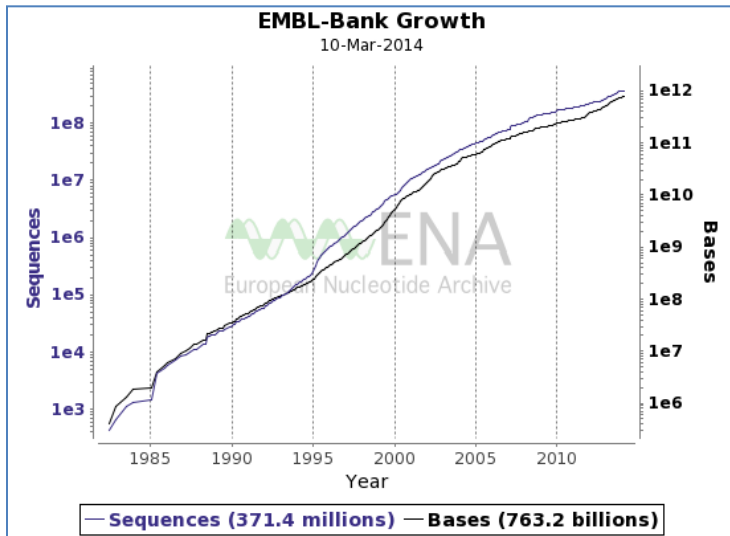Molecular clínical, imaging, population, environmental, ....

# Data production

## Huge data production at different levels
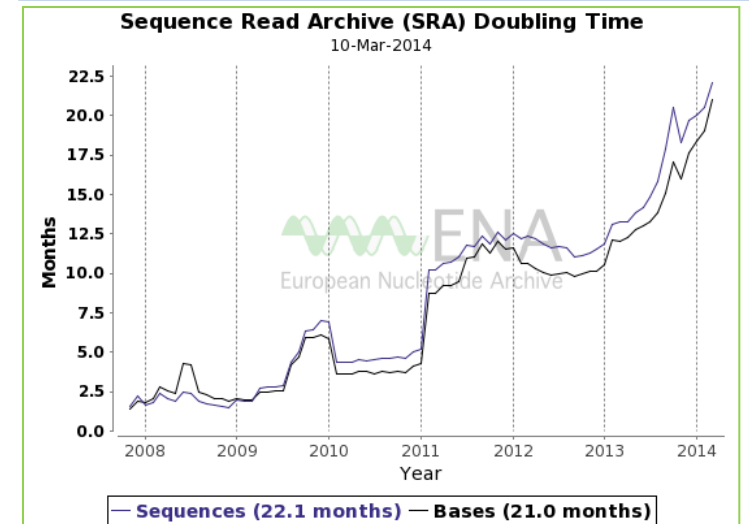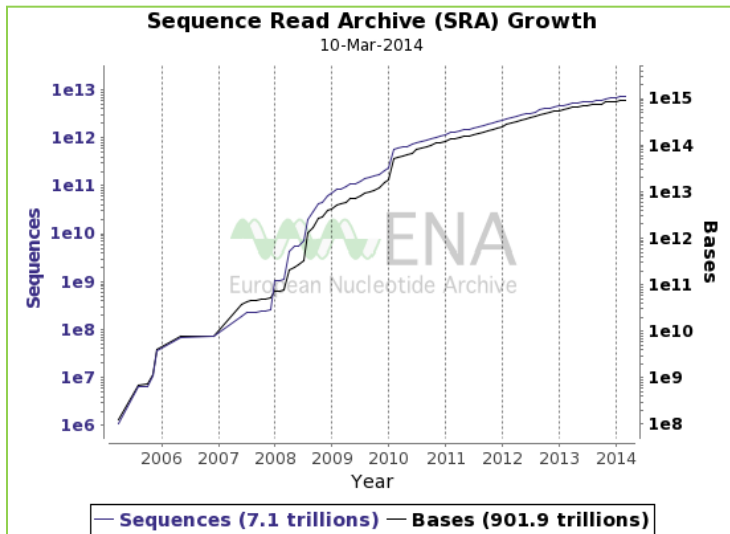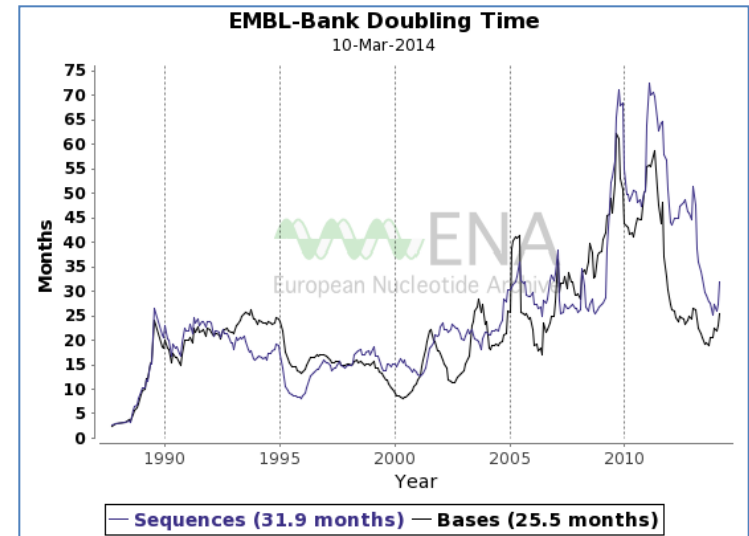


O.Trelles, PhD, 2014

# *f*rom *G*enes to *G*enomes

Assembled/annotated sequence growth

Assembled/annotated sequence doubling time



**EMBL-Bank Growth**
10-Mar-2014

Sequences (371.4 millions) — Bases (763.2 billions)



**EMBL-Bank Doubling Time**
10-Mar-2014

Sequences (31.9 months) — Bases (25.5 months)



**Sequence Read Archive (SRA) Growth**
10-Mar-2014

Sequences (7.1 trillions) — Bases (901.9 trillions)



**Sequence Read Archive (SRA) Doubling Time**
10-Mar-2014

Sequences (22.1 months) — Bases (21.0 months)

READs growth

READs doubling time

ENA statistics: https://www.ebi.ac.uk/ena/about/statistics   O.Trelles, PhD, 2014

# Diverse types of data

Atoms

Proteins

Interactions

Metabolic pathways
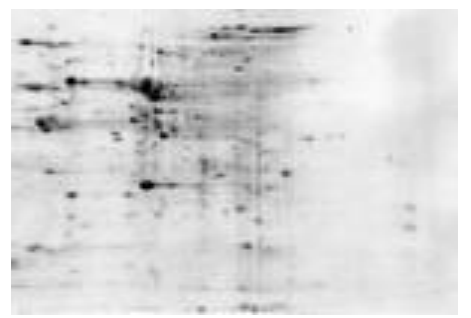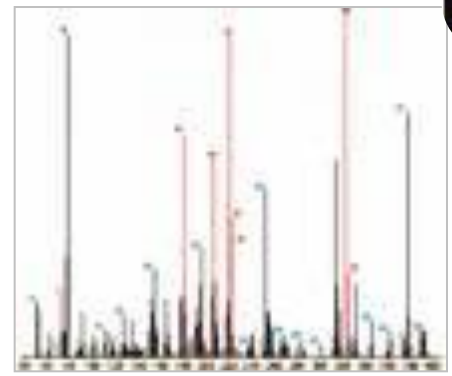
Cells

Organs

Organisms

Populations
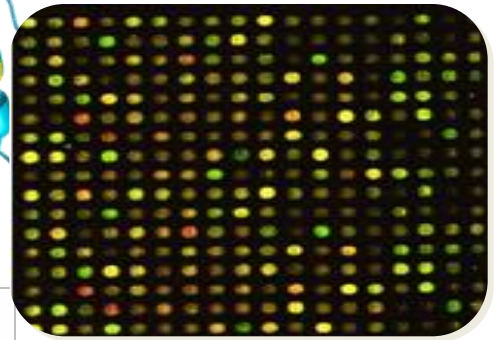
```
> E01306 229 bp  DNA linear
gaattctaac ggtcccgaaa ctctgtgcgg tgctgaactg gttgacgctc tgcagtttgt
ttgcggtgac cgtggttttt attttaacaa acccactggt tatggttctt cttctcgtcg
tgctccccag actcgtactg ttgacgaatg ctgctttcgt tcttgcgacc tgcgtcgtct
ggaaatgta  tgcg tc cc  ga acccgc taaatctgct tagaagctt
```

# Format heterogeneity

```
LOCUS       E01306      229 bp    DNA       linear   PAT 04-NOV-2005
DEFINITION  DNA encoding human insulin-like growth factor I(IGF-I).
ACCESSION   E01306
VERSION     E01306.1  GI:2169565
KEYWORDS    JP 1987190088-A/1.
SOURCE      synthetic construct
  ORGANISM  synthetic construct
            other sequences; artificial sequences.
REFERENCE   1  (bases 1 to 229)
  AUTHORS   Raasu,A., Toomasu,M., Berun,N. and Majiasu,U.
  TITLE     METHOD FOR TRANSPORTING GENE PRODUCT TO MEDIUM PROPAGATING GRAM
            NEGATIVE BACTERIA
  JOURNAL   Patent: JP 1987190088-A 1 20-AUG-1987;
            KABIGEN AB
COMMENT     OS   Artificial gene
            OC   Artificial sequence; Genes.
            OS   Homo sapiens
            PN   JP 1987190088-A/1
            PD   20-AUG-1987
            CC   strandedness: Single;
            CC   topology: Linear;
            CC   hypothetical: No;
            CC   anti-sense: No;
            FH   Key    Location/Qualifiers
            FT         /product='human insuline-Like growth factor I
            FT   CDS            >2..223
FEATURES          Location/Qualifiers
     source          1..229
                     /organism="synthetic construct"
                     /mol_type="unassigned DNA"
                     /db_xref="taxon:32630"
ORIGIN
        1 gaattctaac ggtcccgaaa ctctgtgcgg tgctgaactg gttgacgctc tgcagtttgt
       61 ttgcggtgac cgtggttttt attttaacaa acccactggt tatggttctt cttctcgtcg
      121 tgctccccag actggtattg ttgacgaatg ctgctttcgt tcttgcgacc tgcgtcgtct
      181 ggaaatgtat tgcgctcccc tgaaacccgc taaatctgct tagaagctt
//
```

```
ID   E01306; SV 1; linear; unassigned DNA; PAT; SYN; 229 BP.
AC   E01306;
DT   07-OCT-1997 (Rel. 52, Created)
DT   09-NOV-2005 (Rel. 85, Last updated, Version 3)
DE   DNA encoding human insulin-like growth factor I(IGF-I).
KW   JP 1987190088-A/1.
OS   synthetic construct
OC   other sequences; artificial sequences.
RA   Raasu A., Toomasu M., Berun N., Majiasu U.;
RT   "METHOD FOR TRANSPORTING GENE PRODUCT TO MEDIUM PROPAGATING GRAM
RT   NEGATIVE BACTERIA";
RL   Patent number JP1987190088-A/1, 20-AUG-1987.
RL   KABIGEN AB.
CC   OS   Artificial gene
CC   OC   Artificial sequence; Genes.
CC   OS   Homo sapiens
CC   CC   strandedness: Single;
CC   CC   topology: Linear;
CC   CC   hypothetical: No;
CC   CC   anti-sense: No;
CC   FH   Key     Location/Qualifiers
CC   FT   mat_peptide     11..220
CC   FT   CDS       >2..223
CC   FT          /product="human insulin-like growth factor I"
FH   Key          Location/Qualifiers
FT   source       1..229
FT                /organism="synthetic construct"
FT                /mol_type="unassigned DNA"
FT                /db_xref="taxon:32630"
SQ   Sequence 229 BP; 40 A; 57 C; 55 G; 77 T; 0 other;
     gaattctaac ggtcccgaaa ctctgtgcgg tgctgaactg gttgacgctc tgcagtttgt    60
     ttgcggtgac cgtggttttt attttaacaa acccactggt tatggttctt cttctcgtcg   120
     tgctccccag actggtattg ttgacgaatg ctgctttcgt tcttgcgacc tgcgtcgtct   180
     ggaaatgtat tgcgctcccc tgaaacccgc taaatctgct tagaagctt             229
//
```

The DNA encoding human insulin-like growth factor I(IGF-I) available at GenBank: E01306.1 http://www.ncbi.nlm.nih.gov/

The same insulin (E01306) sequence at EBI www.ebi.ac.uk
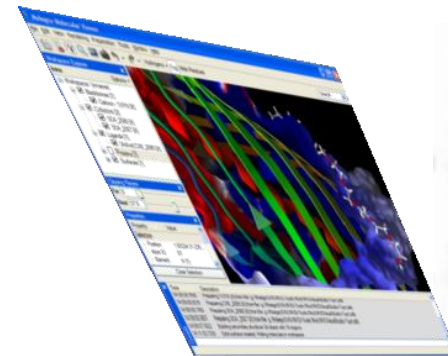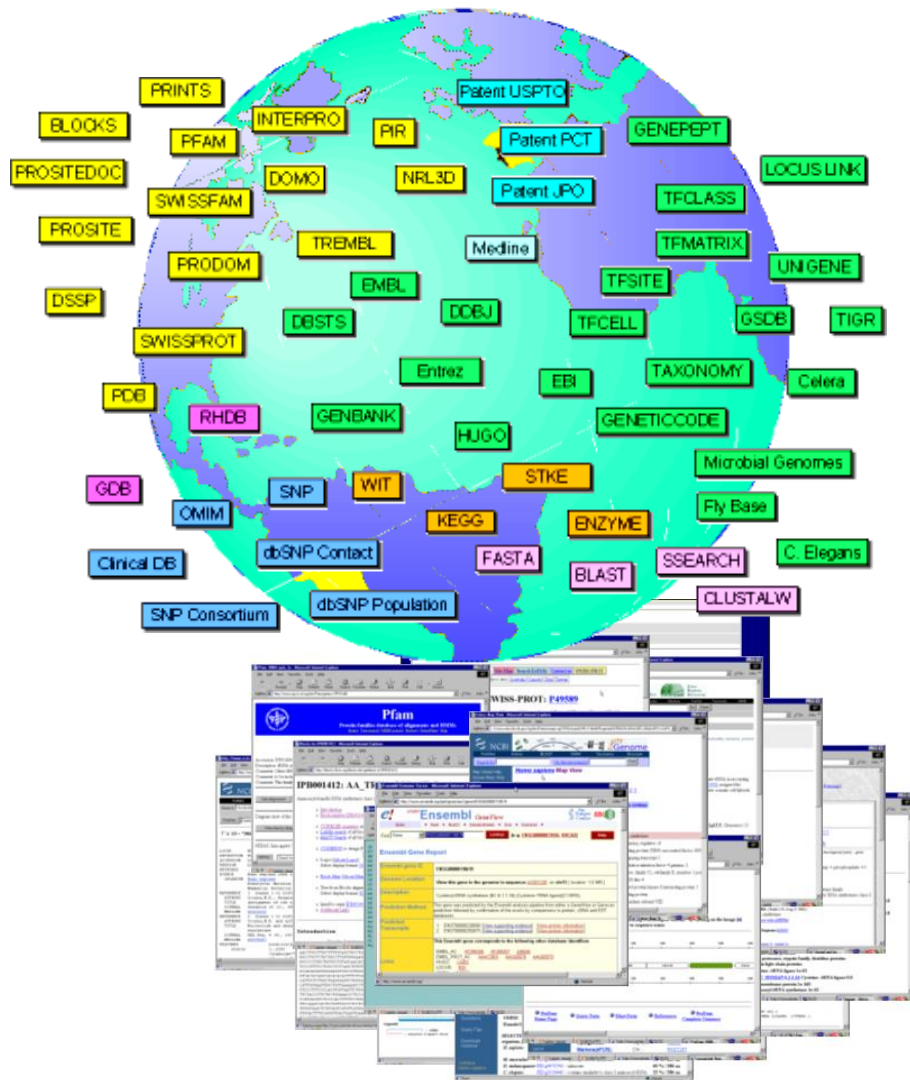
(in both text-boxes some lines has been removed)

# Dispersion of data sources



More than 1000 biological DB collections

Workflows: the usual way to work

See: [1] Infobiogen: Catalog of DBs:

http://www.infobiogen.fr/services/dbcat

Bioinformatics: a web-based domain

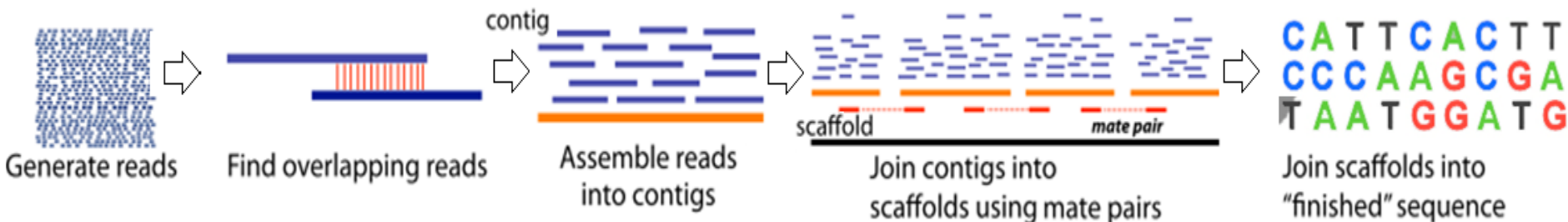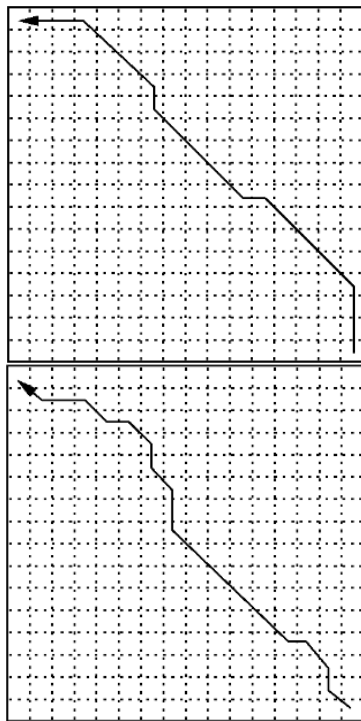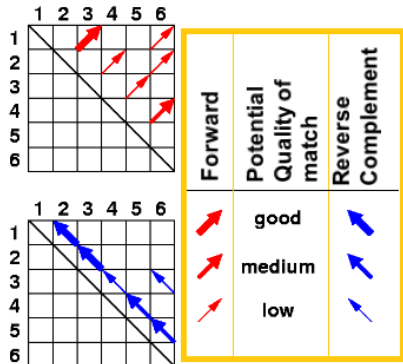O.Trelles, PhD, 2014

**Applications**

O.Trelles, PhD, 2014

# DNA Sequencing (*n*NGS) & *A*ssembly
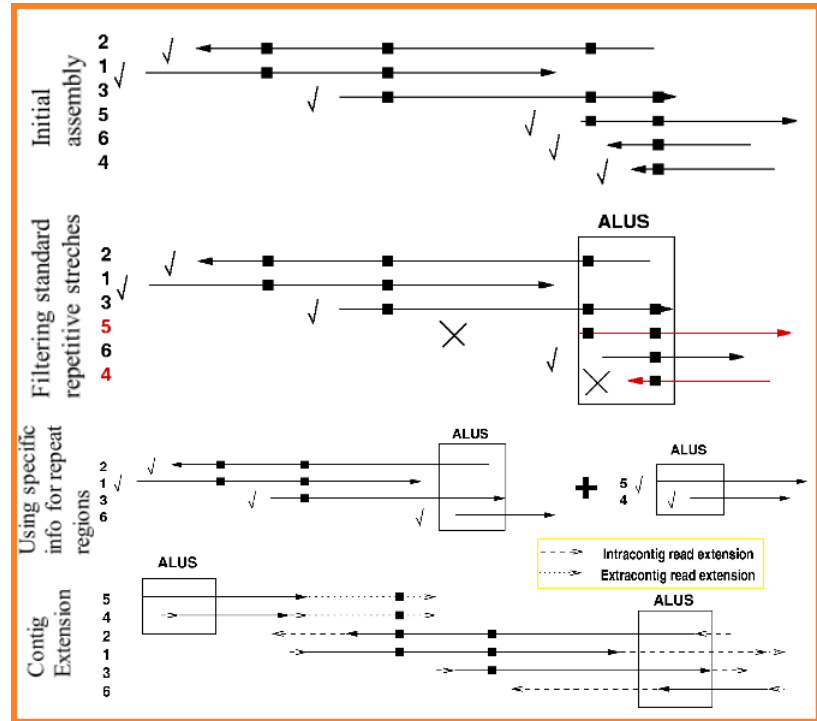
## (>> $10^9$ sequence reads / 36bp to 1kb)

| Field / algorithm | Input | Data volume | | Processing features: computational load, memory access pattern |
|---|---|---|---|---|
| | | In | Out | |
| **Genomics** | | | | |
| **1.1 Next Generation Sequencing (NGS)** | | | | |
| Data acquisition. Image processing | Chromatograms | 300 GB | 4 GB | Image processing. Light and regular pattern |
| Quality filters | Short sequences and by-residue quality value | 4 GB | 4 GB | Sequential processing. Light and regular pattern |
| By-homology clustering of fragments (de novo) | Short sequences | 4 GB | 4 GB | All-All. Out-of-memory. New algorithms |
| By-homology mapping of fragments (mapping) | Short sequences | | | Huge mapping space. High irregular load |
| Assembly contigs from clusters (overlapp) | Group of sequences | 4 GB | 4 GB | All-All. Out-of-memory. New algorithms |
| Copy Number Variations (CNV) | Group of sequences | 4 GB | 10 MB | All-All for each group + MSA. Irregular with data dependencies |
| Single Nucleotide Polimorphism (SNP) | Group of sequences | 4 GB | 10 MB | All-All for each group + MSA. Irregular with data dependencies |



Generate reads → Find overlapping reads → Assemble reads into contigs → Join contigs into scaffolds using mate pairs → Join scaffolds into "finished" sequence
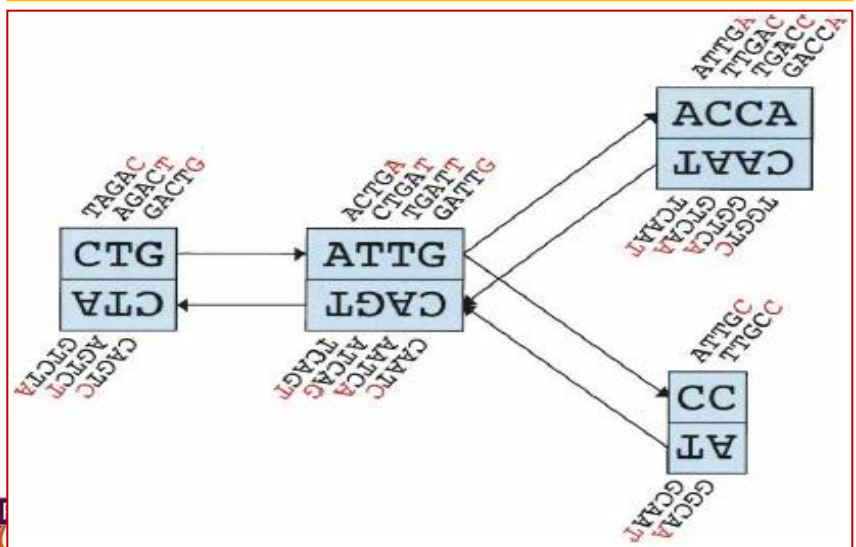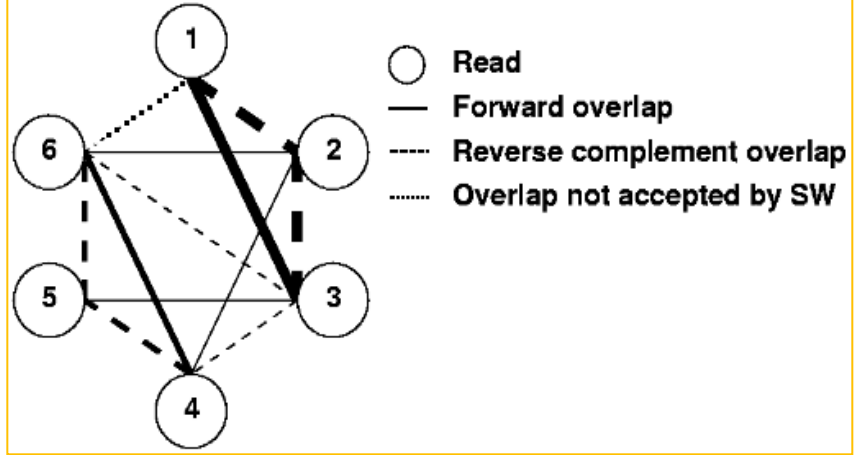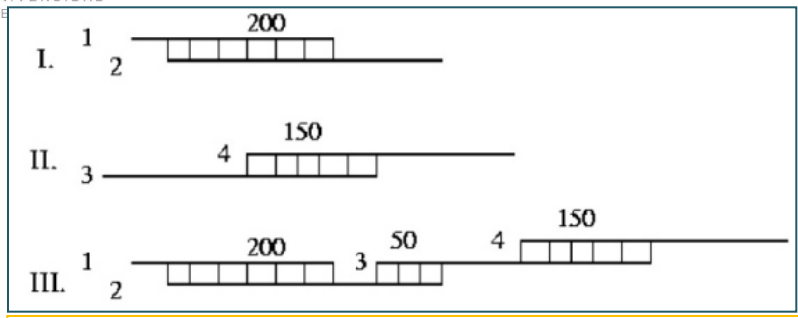
# General concepts for NGS assembly *A*lgorithms (2)



Accepted match
Expected Score : 196
Computed Score: 180
Goodness : 92%
Weight : 1518117

| All-vs-All + reversed complement | Dynamic Programming bounded Gaps | Built-up Contigs and extensions |

# Assembly Algorithms: Data Management



**Greedy** assembly: progressive joint of overlapping fragments.

**Overlap Layout consensus**:
reads are nodes and overlaps are edges. Identify a Hamiltonian path through the graph that contains all the nodes

**Eulerian path** approaches breaks up each read into their overlapping k-mers. Each k-mer is and edge connecting two nodes of its k-1 prefix and suffix respectively. The assembly solution is a path in the graph that uses all the edges - an Eulerian path.

(see also: Bruijin graphs (Velvet)
http://en.wikipedia.org/wiki/De_Bruijn_graph

Align-layout-consensus - Mapping of reads over a related genome (or reference)

# **S**equence **A**nalysis & **P**hylogeny

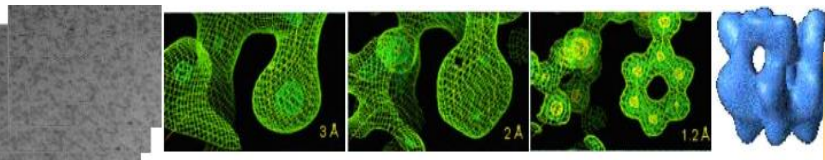| Field / algorithm | Input | Data volume | | Processing features: computational load, memory access pattern |
| --- | --- | --- | --- | --- |
| | | In | Out | |
| **Genomics** | | | | |
| **1.2 Sequence analisis and large scale phylogeny** | | | | |
| Gene identification | Large sequences, full genomes | 1 GB | 10 MB | E/S (local) busquedas intensivas por semejanza, tareas independientes, livianas |
| Searching by homology (Blast, Fasta, Dyn.Prog.) | Query and Sequences DB | 4 GB | 1 MB | E/S (local) intensiva, tareas independientes, livianas |
| Genome scale comparissons (dotplots) | 2 Genomes | 6 GB | 200 MB | Gran demanda de memoria (alg. fuera de memoria) |
| Comparative genomics | Tens of genomes | 30 GB | 4 GB | Gran demanda de memoria y de E/S, nuevos algoritmos |
| Multiple Sequence Alignments (MSA) | Groups of sequences | 10 MB | 1 MB | Todos vs. Todos + resolución de arbol de alineamiento (irregular, dependencias) y diferentes tipos de tareas |
| Phylogeny (by parsimony) | Groups of sequences | 10 MB | 1 MB | Todos vs. Todos + resolución de arbol de alineamiento (irregular, dependencias) y diferentes tipos de tareas |
| Phylogeny (maximum likelihood) | Groups of sequences | 10 MB | 1 MB | Patrón irregular y dependencias de datos. Tareas pesadas |



PhD, 2014

# K-mers numbers

| K | number of combinations in DNA | | | Number of combinations in Proteins | Aprox. |
|---|---|---|---|---|---|
| 1 | 4 | | | 20 | |
| 2 | 16 | | | 400 | |
| 3 | 64 | | | 8.000 | 8 KB |
| 4 | 256 | | | 160.000 | |
| 5 | 1.024 | 1 KB | | 3.200.000 | 3 MB |
| 6 | 4.096 | | | 64.000.000 | |
| 7 | 16.384 | | | 1.280.000.000 | 1,2 GB |
| 8 | 65.536 | | | 25.600.000.000 | |
| 9 | 262.144 | | | 512.000.000.000 | |
| 10 | 1.048.576 | 1 MB | | 10.240.000.000.000 | 10 TB |
| 11 | 4.194.304 | | | 204.800.000.000.000 | |
| 12 | 16.777.216 | | | 4.096.000.000.000.000 | 4 PB |
| 13 | 67.108.864 | | | 81.920.000.000.000.000 | |
| 14 | 268.435.456 | | | 1.638.400.000.000.000.000 | 1,6 EXA |
| 15 | 1.073.741.824 | 1 GB | | 32.768.000.000.000.000.000 | |
| 16 | 4.294.967.296 | | | 655.360.000.000.000.000.000 | |
| 17 | 17.179.869.184 | | | 13.107.200.000.000.000.000.000 | 13 Zetta |
| 18 | 68.719.476.736 | | | 262.144.000.000.000.000.000.000 | |
| 19 | 274.877.906.944 | | | 5.242.880.000.000.000.000.000.000 | |
| 20 | 1.099.511.627.776 | 1 TB | | 104.857.600.000.000.000.000.000.000 | 100 YottaB |
| 25 | 1.125.899.906.842.620 | 1 PETA | | | |
| 30 | 1.152.921.504.606.850.000 | 1 EXA | | | |
| 32 | 18.446.744.073.709.600.000 | | | | |

## Computational space reduction
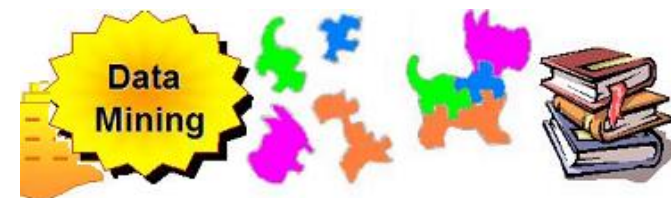
O.Trelles, PhD, 2014

# Structural Analysis: Proteins

| Field / algorithm | Input | Data volume | | Processing features: computational load, memory access pattern |
|---|---|---|---|---|
| | | In | Out | |
| **Proteomics** | | | | |
| **1.2 Sequence analisis and large scale phylogeny** | | | | |
| Structure prediction | 3D protein BD and query seq. | 100 MB | 1 MB | Diversos tipos de tareas (búsquedas BD, comparación de estructuras, refinamiento de topologías), diferente carga |
| By-structural-homology database searching | 3D protein BD and query seq. | 100 MB | 4 MB | E/S, tareas homogéneas, pesadas |
| Structural pattern matching | 3D protein BD and query seq. | 100 MB | 4 MB | E/S, tareas heterogeneas, livianas (conteo) |
| Dynamic protein folding | Query sequence | 4 MB | 4 MB | Heavy tasks with data dependencies |
| Molecular inteactions and docking | Query sequence | 4 MB | 4 MB | Heavy tasks with data dependencies |



```
HEADER     HORMONE                          08-OCT-96   2HIU
TITLE      NMR STRUCTURE OF HUMAN INSULIN IN 20% ACETIC ACID,
COMPND     MOLECULE: INSULIN;
COMPND     3 BIOLOGICAL_UNIT: HETERODIMER
SOURCE     ORGANISM_SCIENTIFIC: HOMO SAPIENS;
KEYWDS     INSULIN, HORMONE, GLUCOSE METABOLISM
AUTHOR     Q.X.HUA,S.N.GOZANI,R.E.CHANCE,J.A.HOFFMANN,B
MODEL        1
ATOM      1  N   GLY A   1      -6.132   6.735   1.016  1.00   0.00
ATOM      2  CA  GLY A   1      -4.686   6.753   1.376  1.00   0.00
ATOM      3  C   GLY A   1      -3.864   6.149   0.235  1.00   0.00
```

Sequence homology | Structural similarity | Physico-chemical properties | Building-up the model | Fine-tuning

Medium/large protein: 3-4 weeks of CPU per 1 nanosecond of simulation
[128 or 256 cores]    : 2-3 ns per day.
Biological processes : rank [micro to millisecond time scale]

O.Trelles, PhD, 2014

# Gene-expression

| Field / algorithm | Input | Data volume | | Processing features: computational load, memory access pattern |
| --- | --- | --- | --- | --- |
| | | In | Out | |
| **Transcriptomics** | | | | |
| **1.4 Gene-expression analysis** | | | | |
| Data acquisition. Image processing | 100 Exp. & 6 M samples | 10 GB | 10 GB | Image processing. Light and regular pattern |
| Data Quality and normalization | 100 Exp. & 6 M samples | 10 GB | 10 GB | Image processing. Heavy, regular pattern |
| Clustering of gene-expression profiles | 100 Exp. & 6 M samples | 10 GB | 10 GB | Out-of-memory, data dependencies, lighted tasks |
| Marker genes identification | 100 Exp. & 6 M samples | 10 GB | 10 GB | Heavy I/O, out-of-memory, light tasks |

**Diferential Expression**

**Clustering**

**Clasification**

**Data Mining**

**KDD: Association studies**

O.Trelles, PhD, 2014

# *I*llustrative use cases

# HPC: The basic model

## DB-searching Applications

- High number of tasks
- Heterogeneity (tasks & CPU-power)
- Network overload
- Scheduling / distribution overload
- Task synchronization
- Fault tolerance
- Portability

## The model

- Task parallel (coarse grained)
- Dynamic load balancing
- Network optimization (message size)
- Minimize number of messages
- Buffering (speculative scheduling)
- Check-points
- SM, DM, D&SM architectures

**Database**

**Master**

**Slaves**

GSS & modGSS

| Iter | GSS R[i] | GSS X[i] | mod-GSS R[i] | mod-GSS X[i] |
|---|---|---|---|---|
| 1 | 100 | 25 | | 1 |
| 2 | 75 | 18 | | 1 |
| 3 | 57 | 14 | | 1 |
| 4 | 43 | 10 | | 1 |
| 5 | 33 | 8 | | 1 |
| 6 | 25 | 6 | | 1 |
| 7 | 19 | 4 | | 2 |
| 8 | 15 | 3 | | 2 |
| 9 | 12 | 3 | | 3 |
| 10 | 9 | 2 | | 4 |
| 11 | 7 | 1 | | 5 |
| 12 | 6 | 1 | | 7 |
| 13 | 5 | 1 | | 9 |
| 14 | 4 | 1 | | 12 |
| 15 | 3 | 1 | 50 | 12 |
| 16 | 2 | 1 | 38 | 9 |
| 17 | 1 | 1 | 29 | 7 |
| 18 | | | 22 | 5 |
| 19 | | | 17 | 4 |
| 20 | | | 13 | 3 |
| 21 | | | 10 | 2 |
| 22 | | | 8 | 2 |
| 23 | | | 6 | 1 |
| 24 | | | 5 | 1 |
| 25 | | | 4 | 1 |
| 26 | | | 3 | 1 |
| 27 | | | 2 | 1 |
| 28 | | | 1 | 1 |

# Sequence DBsrch with Dynamic programming
## Coarse grained

**Master**

```
Get Parameters, Initialize
Start_Workers
Get QuerySeq
Broadcast(QuerySeq)
While (!eof or TransitMess) {
    for all Free_Workers {
         (!eof) Get DBseq
        Prepare(Message)
        Send(Message)
        TransitMess++;
     }
    Receive(R_mess)
    TransitMess--;
}
Broadcast(END_mess)
Report_Best_Results
```

**Workers**

```
Start with params
Perform Initializations
Receive (Query_seq)
while (! END_mess) {



    Receive(Message)
    Score=Algorithm(QuerySeq,DBseq,par);

    Send(Results)
}
```

# ClustalW overview: ( Thompson J. *et al*, NAR, 1994, 2003, 2007 )

- **Pairwise (PW) alignment matrix**
  **average alignment calculation spends most of its time
  here easy to parallelize as all  *N\*(N-1)/2*  elements are**

  **independent**



Multiple Sequence Alignment

32,34 %

0,01 %   67,65 %

Score Matrix
Topology
Alignment

**Cross Similarity Matrix**

|      | [0] | [1] | [2] | [3] | [4] | [5] | [6] |
|------|-----|-----|-----|-----|-----|-----|-----|
| [ 0] | –   | –   | –   | –   | –   | –   | –   |
| [ 1] | 82  | –   | –   | –   | –   | –   | –   |
| [ 2] | 52  | 54  | –   | –   | –   | –   | –   |
| [ 3] | 60  | 62  | 86  | –   | –   | –   | –   |
| [ 4] | 22  | 24  | 18  | 24  | –   | –   | –   |
| [ 5] | 26  | 20  | 12  | 16  | 78  | –   | –   |
| [ 6] | 22  | 14  | 10  | 8   | 46  | 48  | –   |

- **Guide tree calculation**
  **Calculation of closest sequences (branch) is a relatively
  light task, that can be solved sequentially.**

- **Progressive alignment**
  **Remaining ~30% of the code can be parallelized at
  this stage by calculating profile scores in parallel,
  and by solving data dependencies. *(N-1) cluster vs cluster*
  alignments must be solved.
  As a result the whole application is ~90% parallel**

  **depending on a size of a problem**

O.Trelles, PhD, 2014

**Shared Memory Parallel Model**

4

# Irregular algorithms: DNAml

Current-best-tree $T_k$ $(L_k)$ [from insertion step]

for i = 1 to n-tasks

    Remove sub-tree i from $T_k$ and produce $T_{k1}$ and $T_{k2}$

    Likelihood evaluation for $T_{k1}$ and $T_{k2}$ $(L_{k1}$ and $L_{k2})$

    Current-best-tree $T_k$ = tree with greater likelihood $(T_k, T_{k1}, T_{k2})$

end for



O.Trelles, PhD, 2014

# Irregular algorithms
## Speculative computing



DNA-ml: Algorithm Run-Time Behaviour

| Evaluations | 5732 |
| Penalties | 105 |
| Percentage | 1.83% |

# Open & Provoking questions

O.Trelles, PhD, 2014

# Scenarios: the real work

|  | **Comparative genomics scenarios (CG)** |
|---|---|
| CG1 | multi-genome comparison on higher mammalians |
| CG2 | Multi-genome comparison and phylogenomics. |
| CG3 | Symbionts study case |
| CG4 | Metagenome analysis |
|  |  |
|  | **Biomedical scenarios (BM)** |
| BM1 | access to summarized information of the clinical DB through mobiles |
| BM2 | Data analysis: Combining protein interaction and pathway data |
| BM3 | Discovering correlations between clinical and molecular patient data |

**Aims**
(1) big-data, HPC, Grid & Cloud, visualization
(2) Security, data sensitivity, data analysis
(3) New statatistic, math & biological models

O.Trelles, PhD, 2014

# The global idea: HSPs out-of-core



O.Trelles, PhD, 2014

# Pairwise sequence/genome comparison Sequence DBsrch with Dynamic programming

$$S_{i,j} = max\ [$$
$$S_{i-1,j-1} + w(x_i,y_j),$$
$$S_{i-1,j} + \alpha_g ,$$
$$S_{i,j-1} + \alpha_g\ ]$$



Extensión de un alineamiento con gaps



Programación Dinámica

```
qNEW-SEQ    --SARGDFLNAA YALFFMRSHN FGHSDVLPVL
            ||||||||||   |||  |||||    |||||||
KNOWN-SEQ  MMSARGDFLN-- YALSLMRSHN DEHSDVLPVL


qNEW-SEQ    --CSLKHVAY WDAYQALIYW IKAMNQQTDTSI
            |||||||||    ||||||| | ||||||||||
KNOWN-SEQ  DVCSLKHVAY -VFQALIYW IKAMNQQTTLDT


qNEW-SEQ    --RPPDDQAF GHHHLPQAMH --SRLYVPS-SK
            |||        |   || |       ||||||| ||
KNOWN-SEQ TIRPPA---- GAFGLPTANT CISRLYVPSMSK
```

# Multi-genome comparison (CG1)



Func. Annot

Big-data
HPC
Modeling
Visualization Data analysis
GUIs

O.Trelles, PhD, 2014

# MG comparison & phylogeny (CG2)



Pairwise genome alignment
DB searchin strategy
→ Big-Data Kmers dictionaries
Workflows
Visualization

**Modeling**
Blocks detection + refinement
Breakpoints identification
EE frequencies
Inter-genome distances
GUis

O.Trelles, PhD, 2014

Provide a HP&CC suite for highly demanding CG studies

# Meta-genomes comparison



Computational space reduction based on the fast identification of matching reads

Trelles O. et al. Computational Space Reduction and Parallelization of a new Clustering Approach for Large Groups of Sequences"; Bioinformatics vol.14 no.5 1998 (pp.439-451)O.Trelles, PhD, 2014

# Human vs Chimpanzee
## (close related organisms)



Fragment distribution by length and similarity

**Different distributions?**

Introns, exons, intergenic… produce the same type of fragments?

# Statistical Significance of HSPs
## New models are needed



random

intron / intron

Exon / exon

Fragments dist

90 80 60 40 30 20 10
784x10^6 Fragments

Length : 50 - 100

Length 101 - 200

length 201 - 300

Length 301 - 400

Length 401 - 500

Length 501 - 1000

# Evolution events in full genomes



O.Trelles, PhD, 2014

# Bio-Medical scenarios
## (Allergies)



Prostaglandin pathway | Amine metabolism | Unknown function

Genomic Data

Expression Profile

Clinical Record

**Pheno-geno Correlation Models**

J. Tr. PhD, 2014

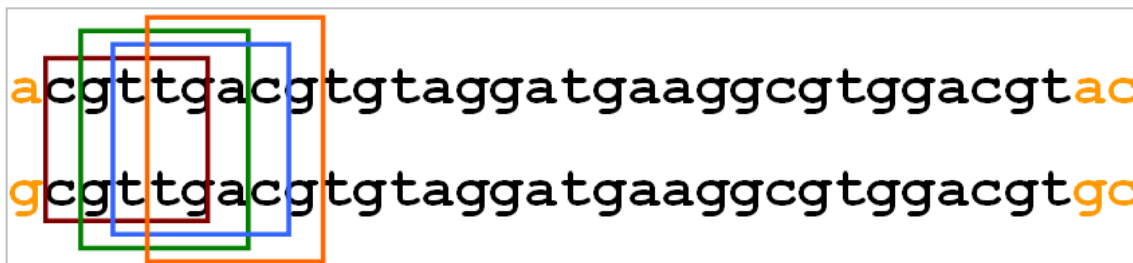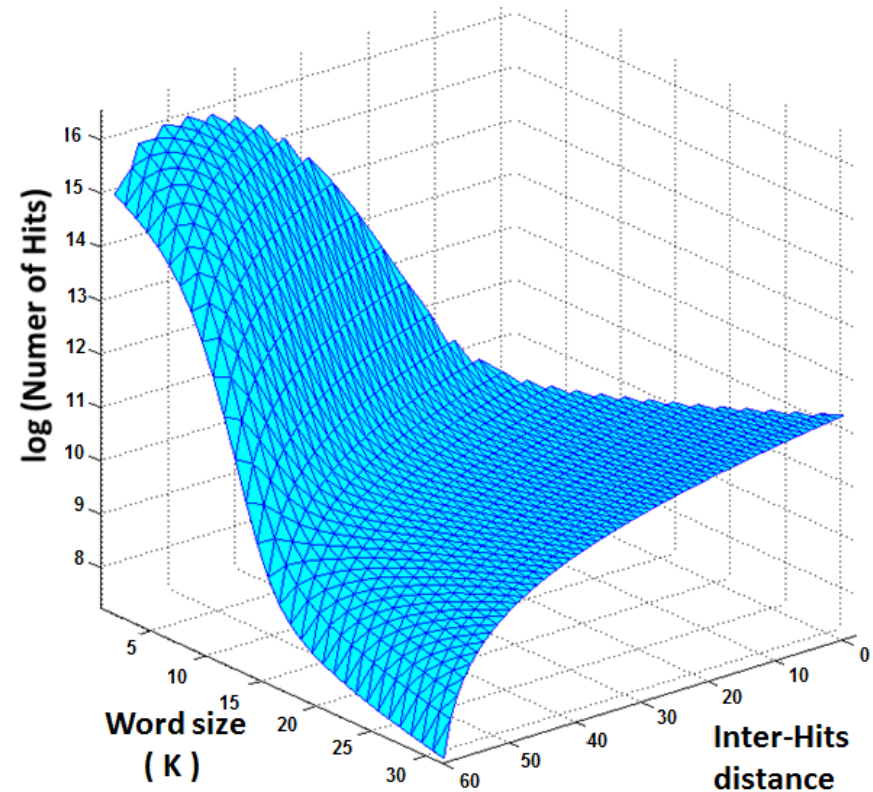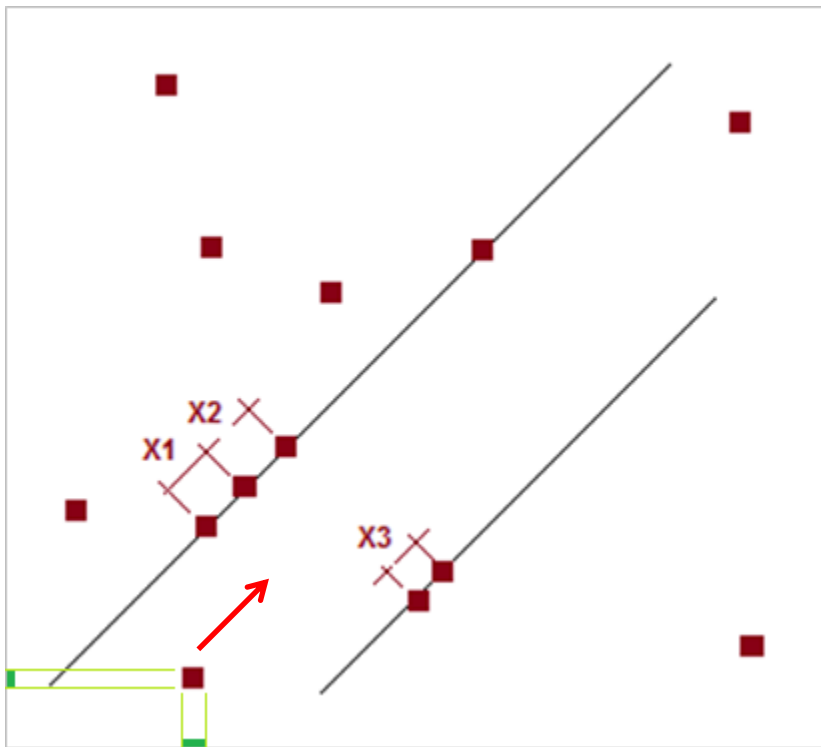# Apps: Patient + Genomic data



GWAS Analysis

O.Trelles, PhD, 2014

# Computational space reduction

**Retrieve compounds in a database that are similar to a query compound**

It's not only a problem of siize but complexity



O.Trelles, PhD, 2014

**Thank you**