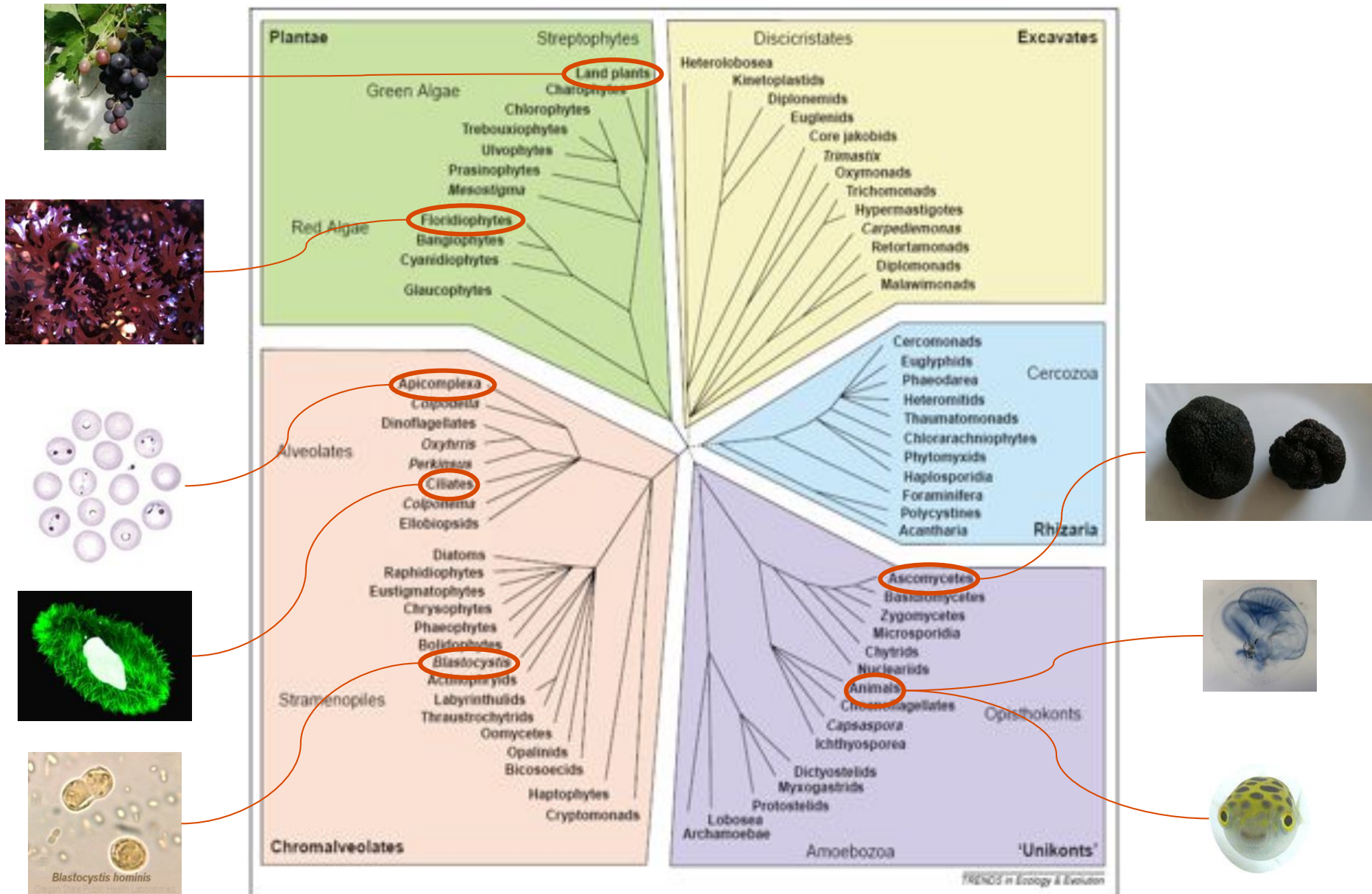


Gaia: at the frontiers of astrometry

Processing massive datasets in genomics

- Created en 1997
- 150 people
- Service for academic collaborative projects
- Support in house R&D research projects
- Large-scale sequencing projects
 - Part of international projects : Human Genome (K14), Arabidopsis, Rice, Medicago, Anopheles etc.
 - Main actor for Tetraodon, Oikopleura, Grape, or Paramecium, Truffle Banana, Coffee, Cacao, wheat 3b chromosome.
- Realized different fungi genome projects (Botrytis, Tuber) and many prokaryote sequencing projects
- Long standing experience in prokaryote and eukaryote genome annotation

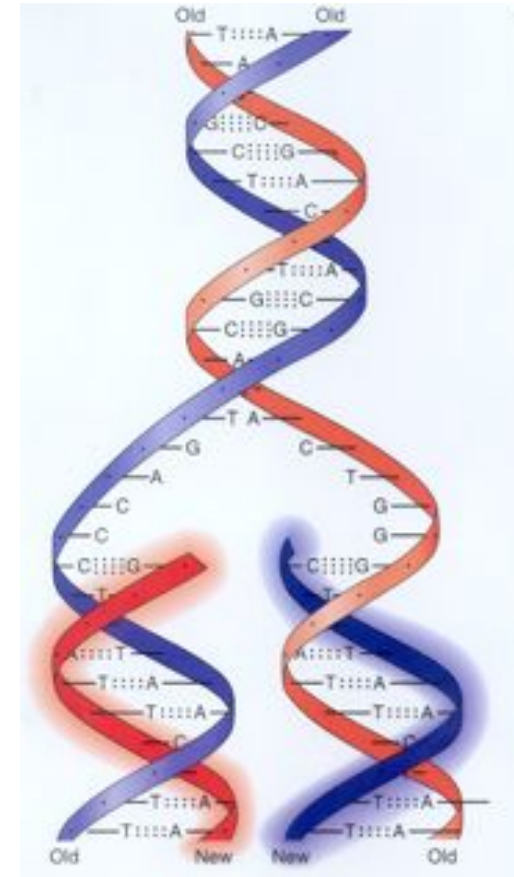
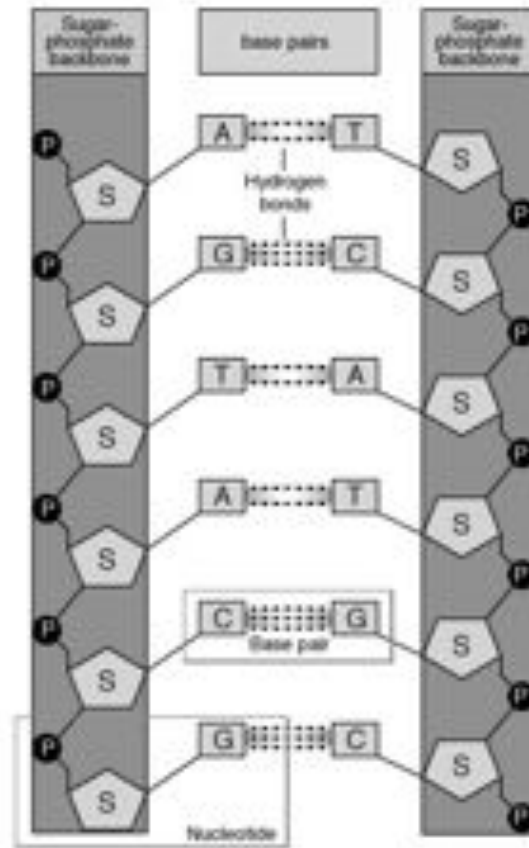




What is genomics ?

- Genomics is “the study of functions and interactions of all the genes in the genome, including their interactions with environmental factors.”
- A genome is “all the DNA contained in an organism or a cell, which includes both the chromosomes within the nucleus and the DNA in mitochondria... all our genes together.”

DNA : the molecule of life

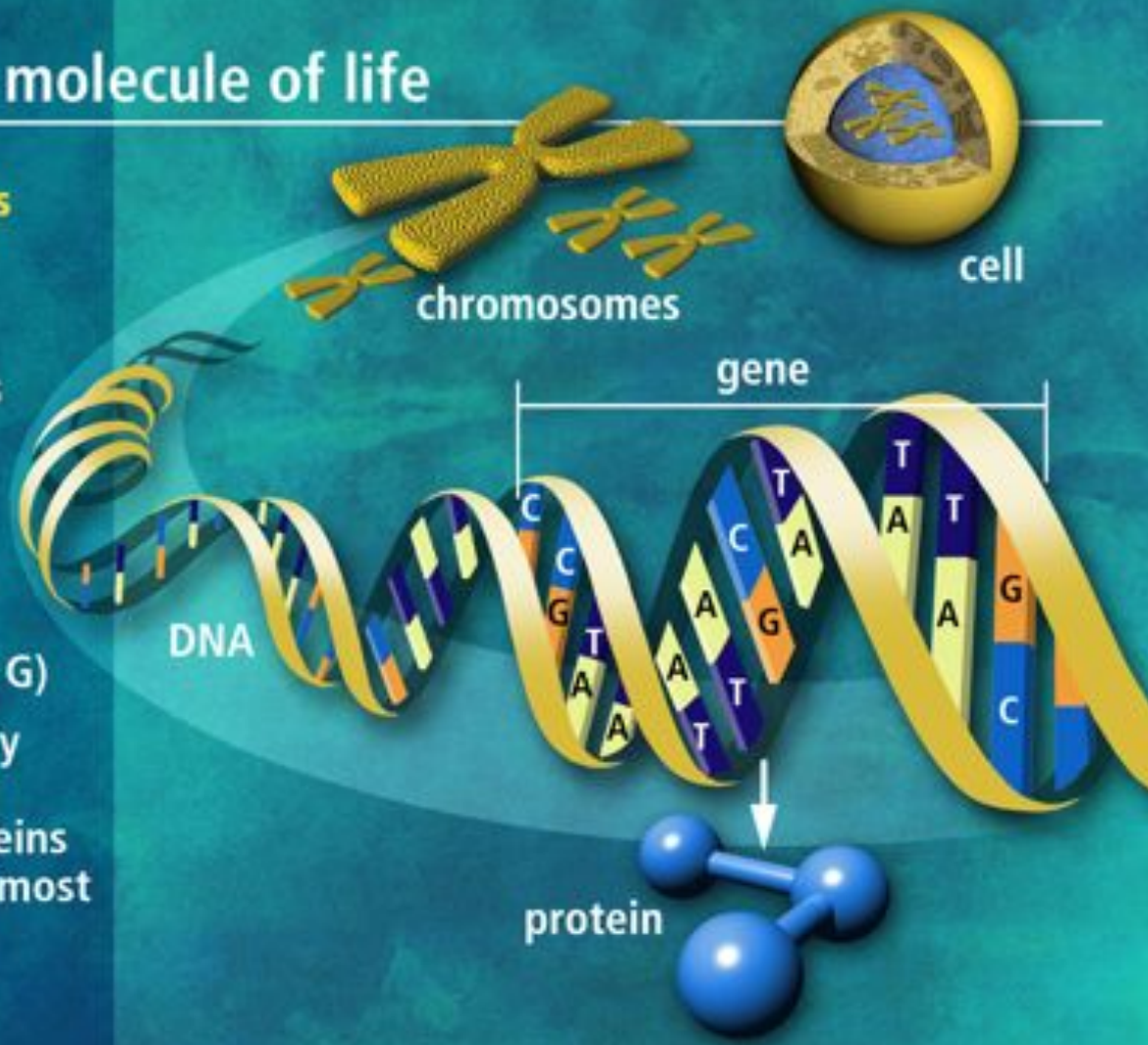


DNA the molecule of life

Trillions of cells

Each cell:

- 46 human chromosomes
- 2 meters of DNA
- 3 billion DNA subunits (the bases: A, T, C, G)
- Approximately 30,000 genes code for proteins that perform most life functions



The diagram illustrates the relationship between different levels of genetic organization. At the top right, a yellow cell contains several yellow chromosomes. Below the cell, several yellow chromosomes are shown in detail. A single chromosome is shown as a long, thin DNA molecule. A section of this DNA molecule is shown as a double helix with nitrogenous bases (A, T, C, G) attached to the sugar-phosphate backbone. A bracket labeled 'gene' spans a portion of this DNA. An arrow points from the gene to a blue ball-and-stick model of a protein.

Y-GG 01-0085

U.S. Department of Energy Human Genome Project ~ www.ornl.gov/hgmis

Annotation of Genomes

Data Distribution

Genome browser

Web site

Submission

Formatting

etc...

Masking

known repeats

Data Collection

cDNAs mapping
Public protein mapping
Public ESTs mapping
Gene models *ab initio* predictions
Repeats *ab initio* detections

Integration

Gene models prediction

Post Annotation Analysis

Proteic domains detection
Paralogs/Orthologs definition
Enzymatic activity detection
Metabolic pathway inference

Homo sapiens Map View build 25

Chromosome: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 X Y

Master: Genes On Sequence Map [Display settings](#)

Total Genes On Chromosome: 988 [13 not localized]

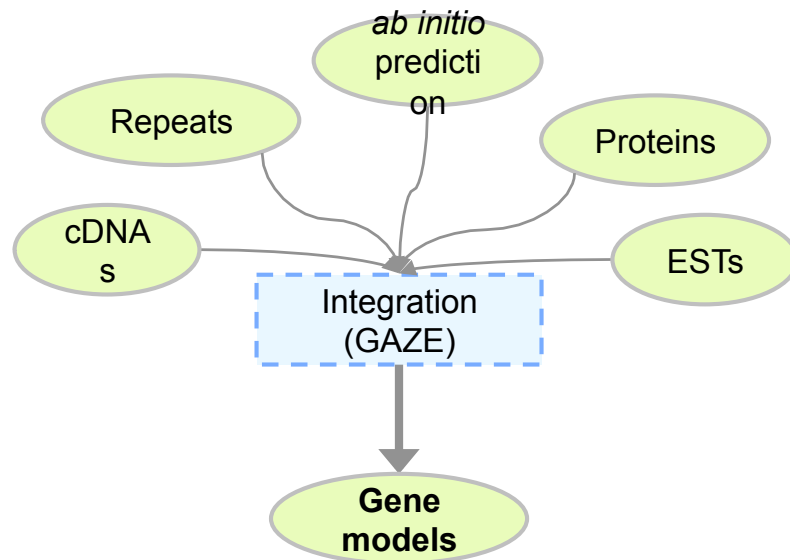
Region Displayed: 97M - 116M bp [Download/View Sequence/Evidence](#)

Genes Labeled: 20 Total Genes in Region: 75

| chr | start | end | transcript | orient | links | cyto. | full name |
|-----|-------|-------|------------|--------|-------|---------------|----------------------------------|
| 11 | 11421 | 11421 | JRKL | + | SV | 11q21 | jerky (mouse) homolog--like |
| 11 | 11421 | 11421 | HSPC048 | + | SV | 11cen-11q12.1 | HSPC048 protein |
| 11 | 11421 | 11421 | KIAA0092 | + | SV | 11cen-q12.1 | KIAA0092 gene product |
| 11 | 11421 | 11421 | PGR | + | SV | 11q22-q23 | progesterone receptor |
| 11 | 11421 | 11421 | YAP1 | + | SV | 11q13 | Yes-associated protein 1, 65 kDa |

The screenshot displays a genome browser interface with a sequence alignment track at the top and a gene model track below. The sequence alignment shows reads from various sources (e.g., ESTs, cDNAs) aligned to a reference sequence. The gene model track shows predicted gene structures with exons and introns. The interface includes a search bar and various navigation options.

This screenshot provides a detailed view of the JRKL gene model. It shows the gene structure with exons and introns, along with associated transcripts and protein domains. The interface includes a search bar and various navigation options. The gene is located on chromosome 11 at position 11421.



Data integration

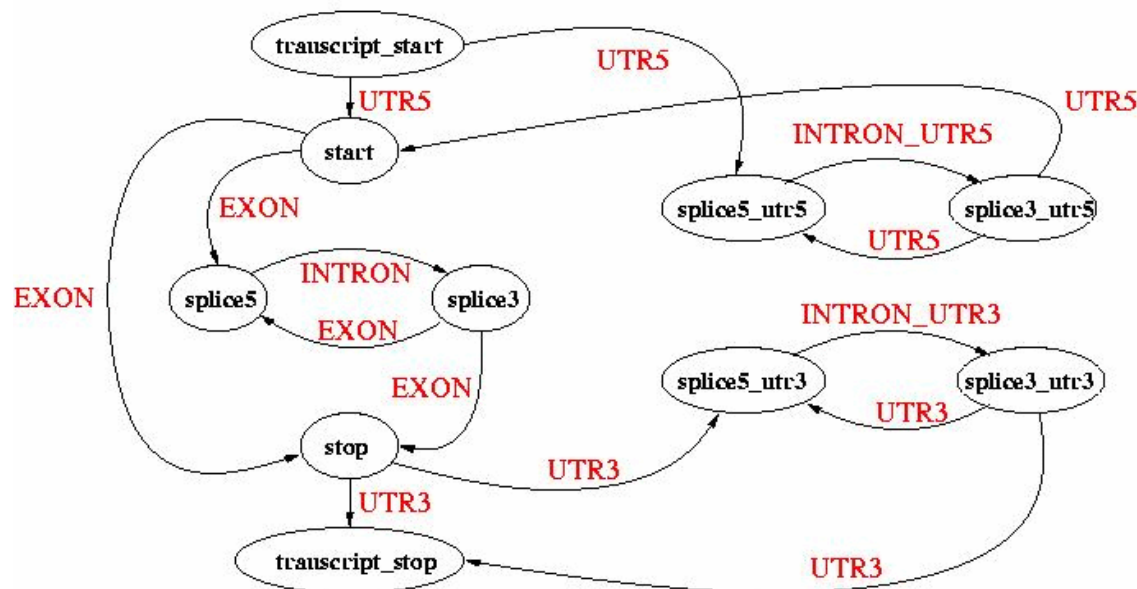
Objective : Build gene models using whole dataset (expression products, proteins, ab initio (machine learning) predictors).

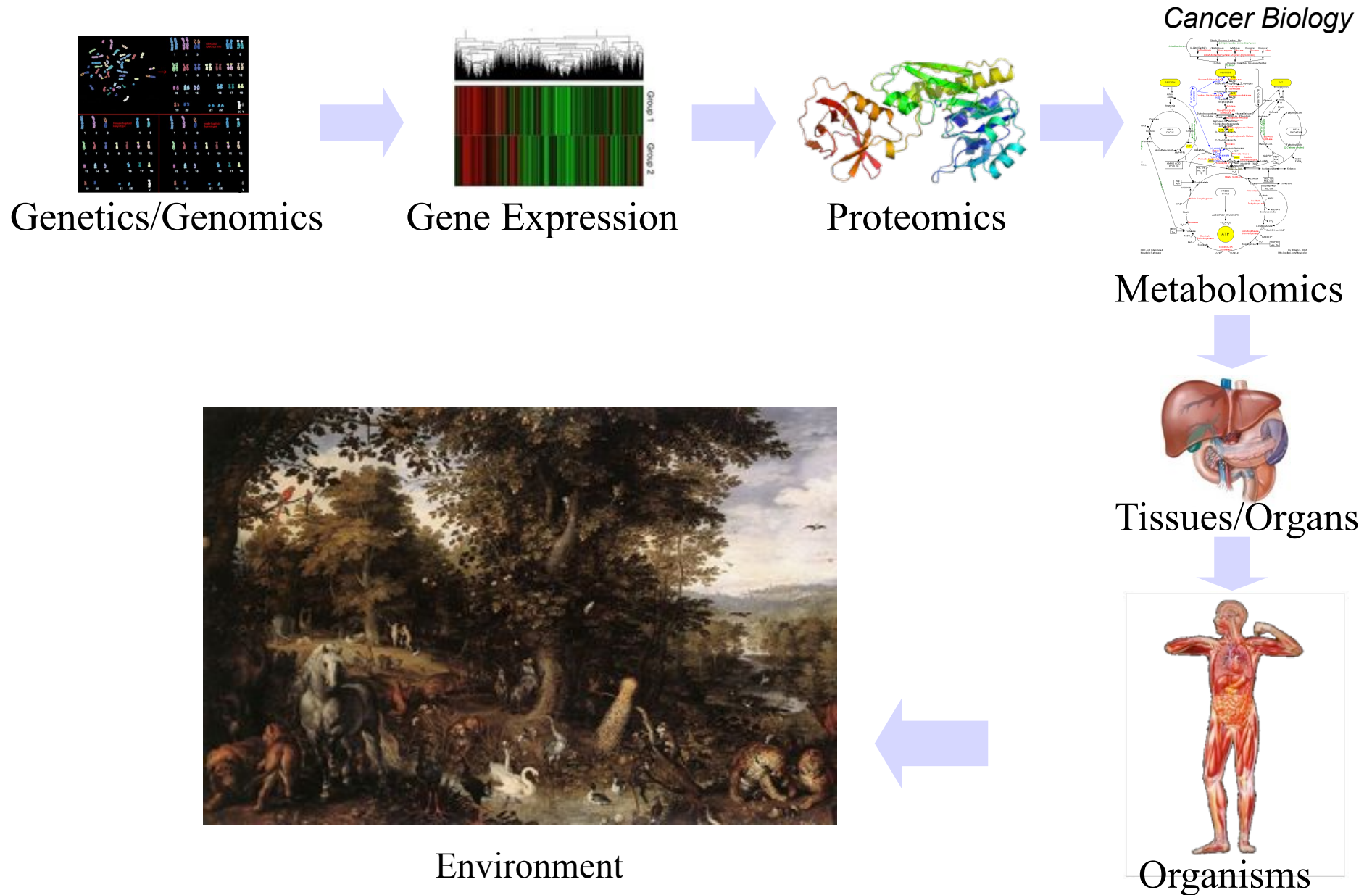
Gene model represented as an automat

2 type of elements

segments (exons, intron,...)

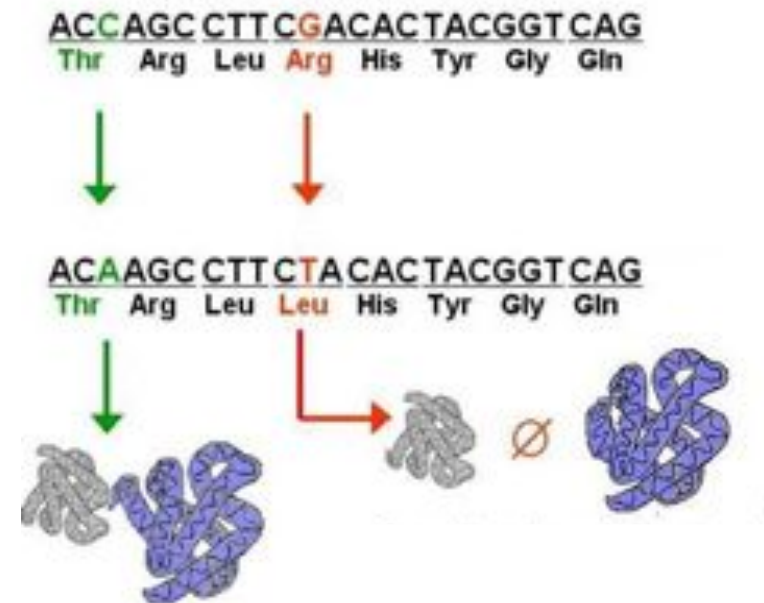
signals (start, GT/AG, stop,...)





Why study genomics ?

- From a public health perspective, genomics is the study of the gene-environment-host interaction that leads to disease — or disease prevention — in populations.
- Rare diseases
- Single gene disorders
- Public health activities
 - Newborn screening
 - Reproductive health
 - Genetic services





Contemporary Public Health Genetics



- Common diseases
- Multiple genes
- Gene/environment interactions
- Public health activities/implications
 - Chronic diseases
 - Infectious diseases
 - Environmental health
 - Epidemiology

- Genetic testing
 - To detect mutations
 - For disease diagnosis and prognosis
 - For the prediction of disease risk in individuals or families
- Several hundred genetic tests are in use.
 - Rare genetic disorders (muscular dystrophies, cystic fibrosis)
 - Complex conditions (breast, ovarian, and colon cancers)
- Pharmacogenomics
 - The development of drugs tailored to specific subpopulations based on genes
 - Pharmacogenomics has the potential to:
 - Decrease side effects of drugs
 - Increase drug effectiveness
 - Make drug development faster and less costly

- 1953 The double helical structure of DNA
- 1977 **Frederick Sanger (MRC) develops methods for sequencing DNA**
- 1980 First method to map the entire human genome based on RFLPs
- 1984 Complete DNA sequence of the Epstein-Barr virus, 170 kb
- 1985 Kary Mullis develops PCR , a technique to replicate vast amounts of DNA
- 1986 **The first automated DNA sequencing machine**
- 1987 Human Genome Initiative begins
- 1992 First physical map of chromosome 21 (D. Cohen)
First genetic maps of human (J. Weissenbach)
- 1995 The first sequence of a complete genome, *Haemophilus influenzae* , 1.8 Mb (C. Venter)
- 1996 NIH funds six groups to attempt large-scale sequencing of the human genome
- 1997 Genoscope is setup in France
- 1998 3700 capillary sequencing machines
Celera company is funded and declares that it will sequence the human genome within 3 years
New objectives of the HGP of creating a "working draft" of the human genome by 2001
Completion date for the finished draft from 2005 to 2003.
- 2000 Complete genome of the first plant, *Arabidopsis thaliana* 125 Mb

The HGP consortium publishes its working draft in *Nature* (15 February), and Celera publishes its draft in *Science* (16 February).



1995 : The human genome was considered a 10 years world wide project

2005 Second generation of DNA sequencing machines

A draft of the human genome could be obtained in less than 3 months in one genome center

2011 : Third generation of DNA sequencing machines

=> A human genome in 15'



Applied Biosystems
ABI 3730XL

1000 euros / Mbase



Roche / 454
Genome Sequencer FLX

100 euros / Mbase

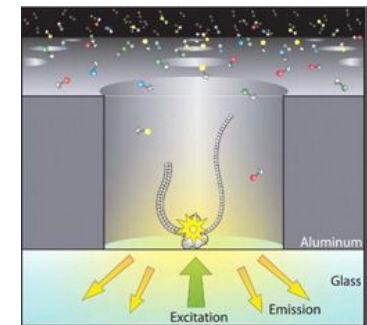


Illumina / Solexa
Genetic Analyzer

10 euros / Mbase

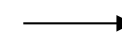


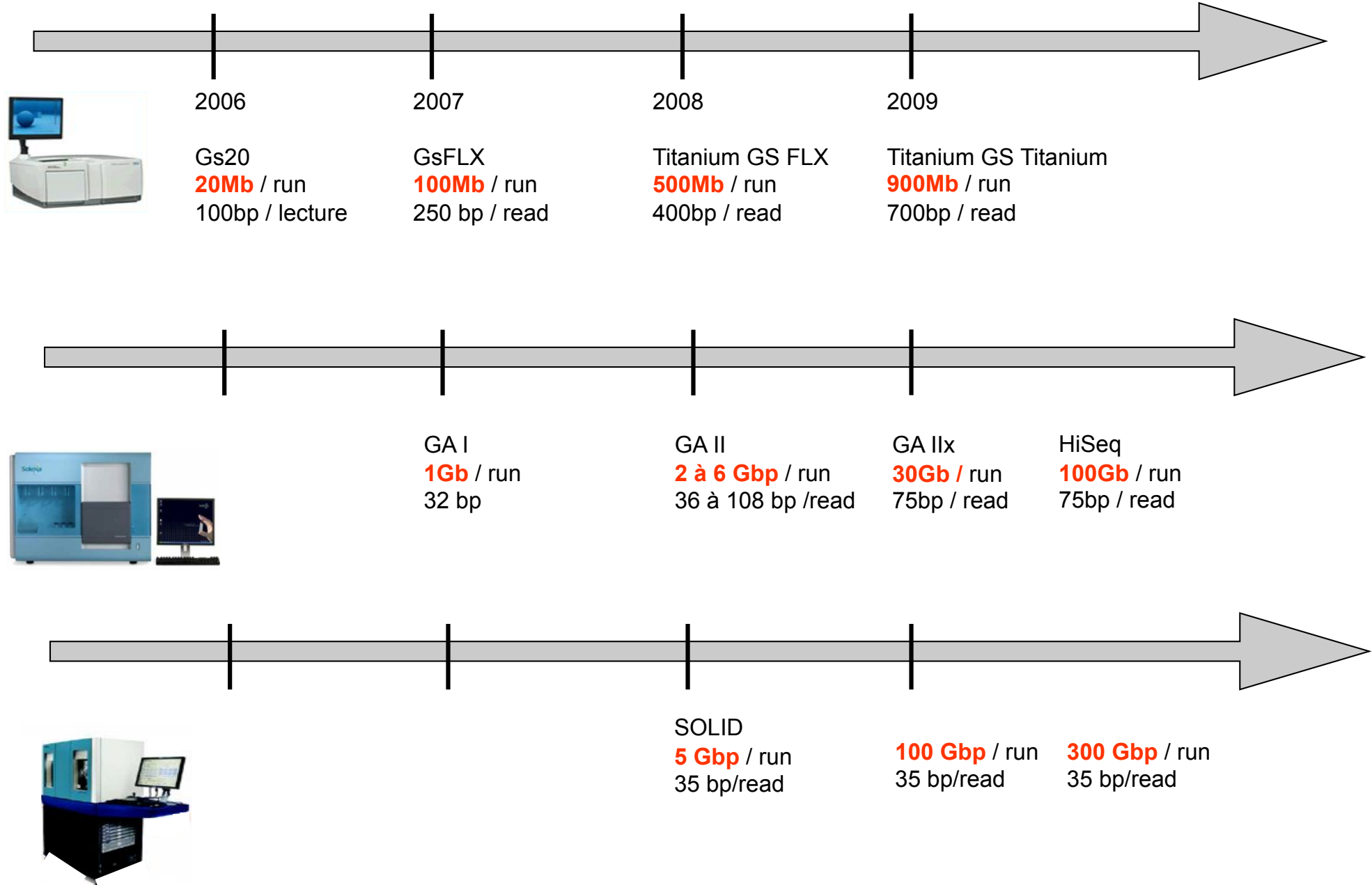
Applied Biosystems
SOLiD

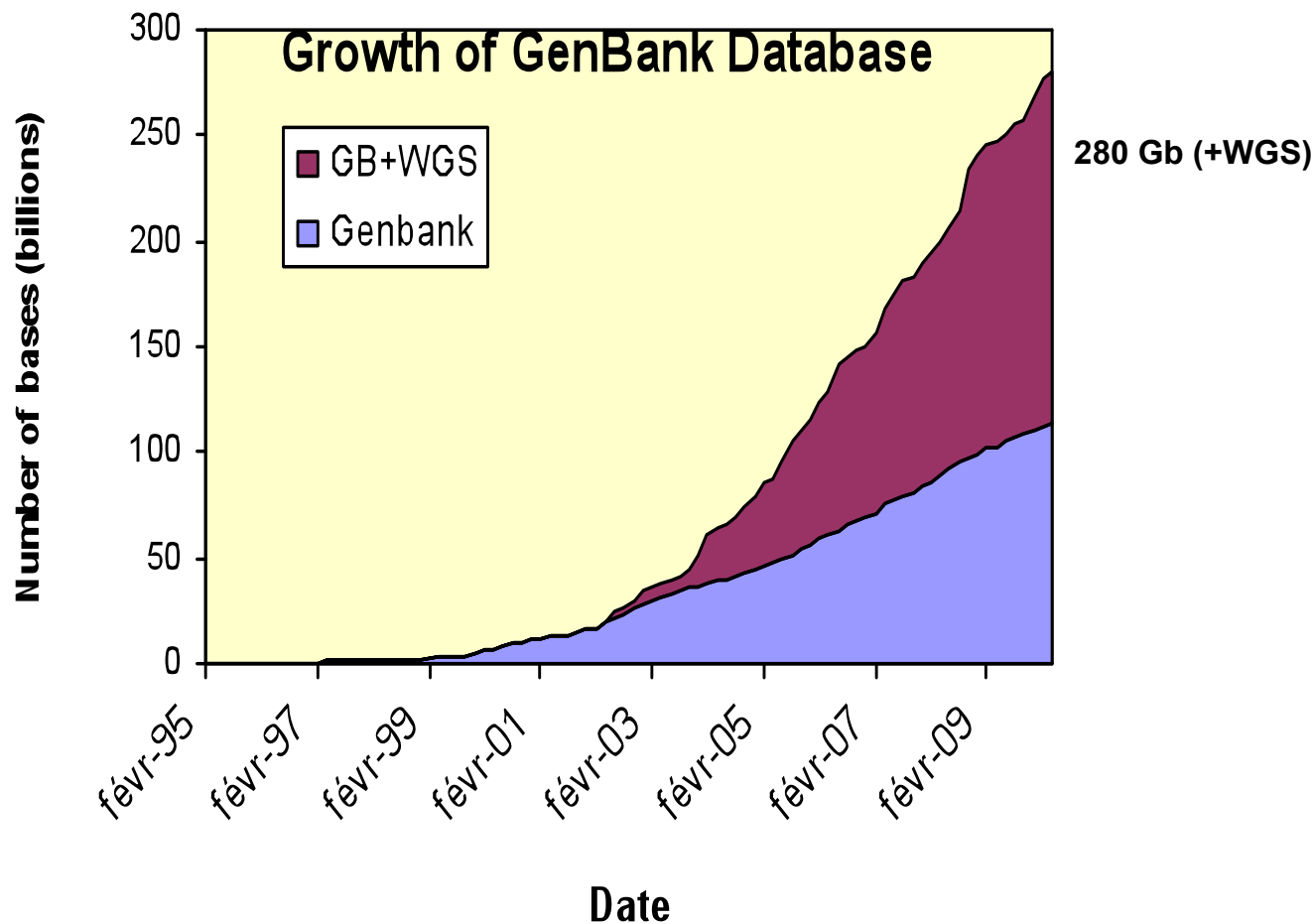


Pacific biosciences

<1 euros / Mbase







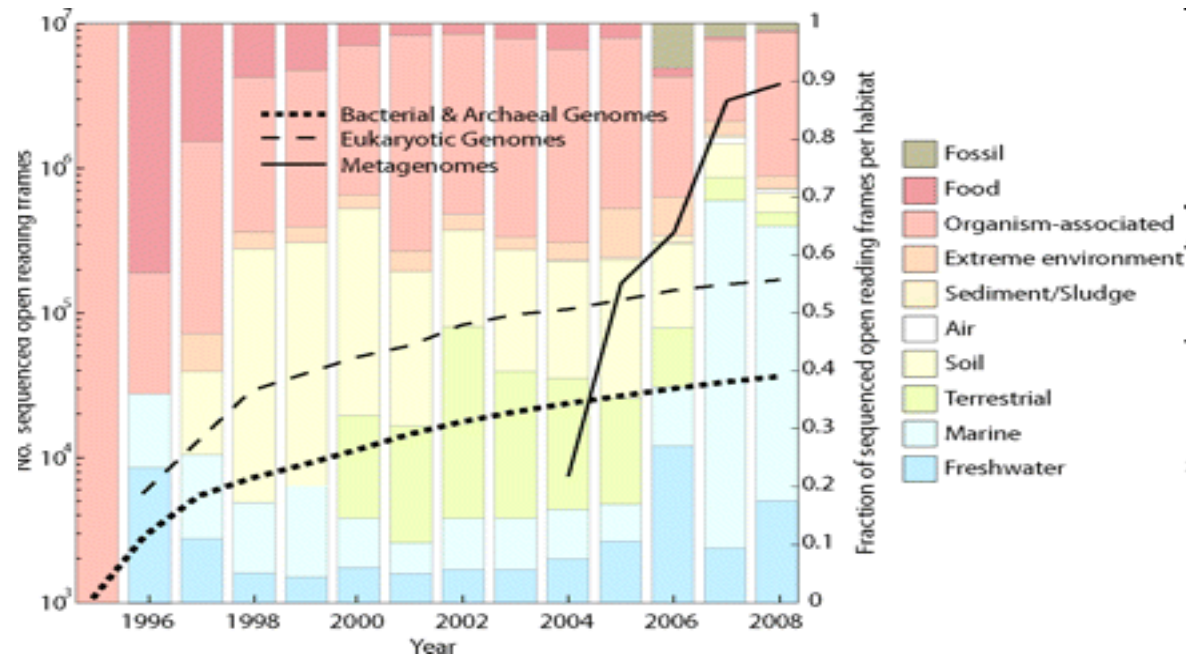
We project in 3-5 years: 100x increase in sequencing volume

Fundamental computing capabilities should increase:

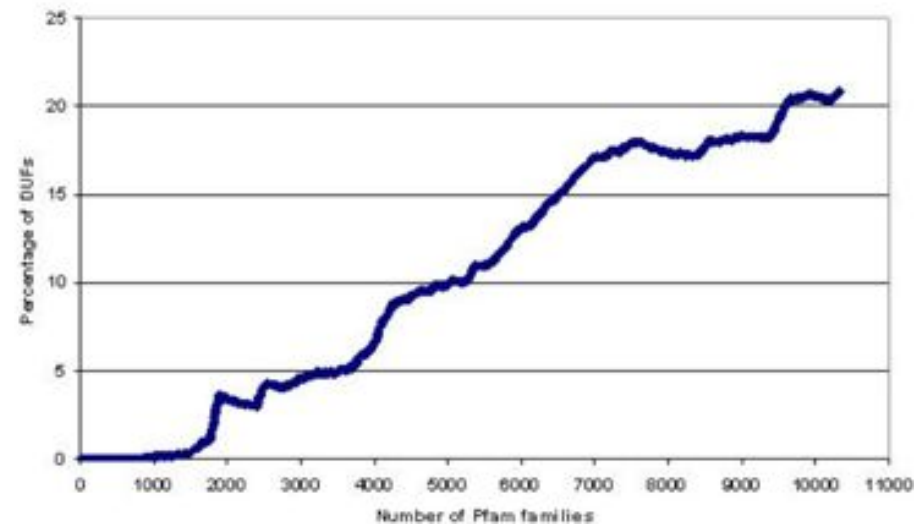
7-10x in 5 years

50-100x in 10 years

Trends in the increase of genomics data



Domains of unknown function (DUFs) in PFAM database

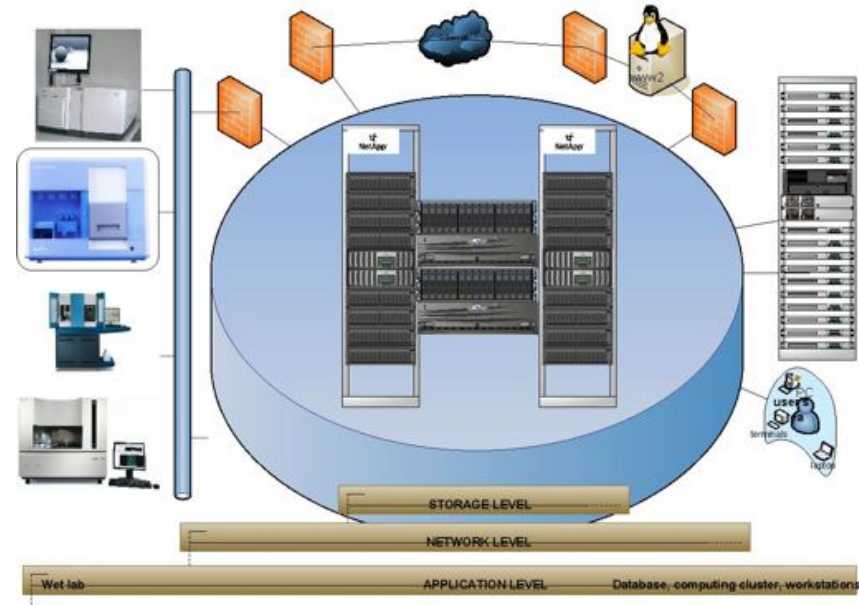


- Computational facilities :

- Genoscope :

- Main compute cluster

- ~40 Sun type Sunfire x4100, 330 cores,
 - > 2TB memory.
 - Storage : ~50 TB, increasing.



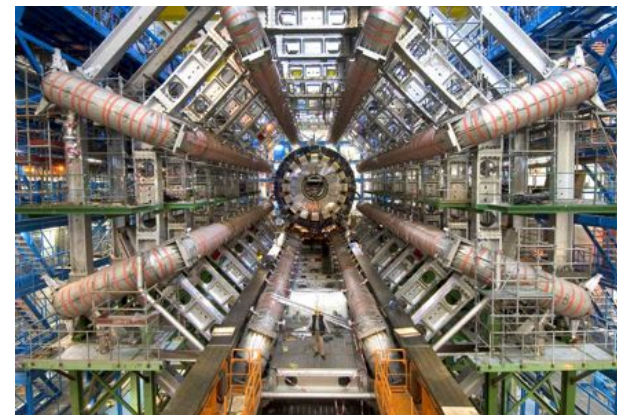
- Genci / CCRT (CEA – Bruyères-le-Château)

- Bull itanium : 47.7 Tflops, ~4000 cores
 - 23To memory
 - Storage : 420To.

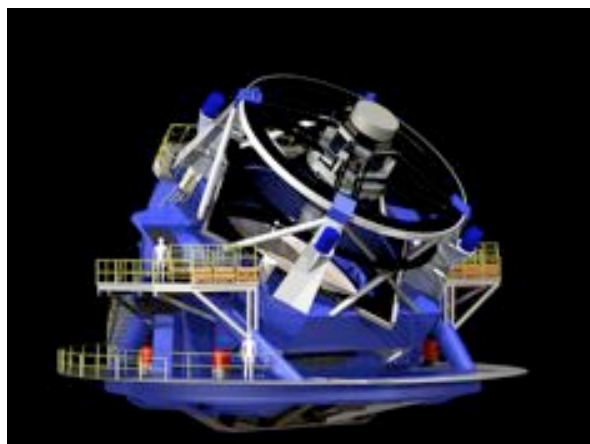


CERN
Large Hadron Collider (LHC)

~10 PB/year at start
~1000 PB in ~10 years



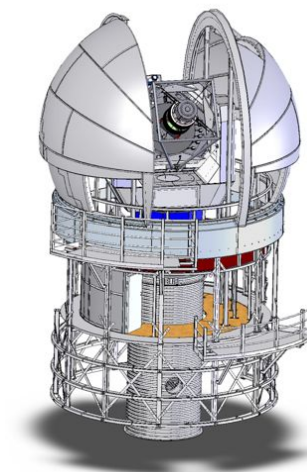
<http://www.cern.ch>



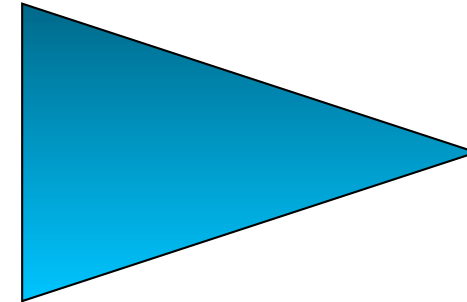
Large Synoptic Survey Telescope (LSST)
NSF, DOE, and private donors

~5-10 PB/year at start in 2012
~100 PB by 2025

Pan-STARRS (Haleakala, Hawaii)
US Air Force
now: 800 TB/year
soon: 4 PB/year

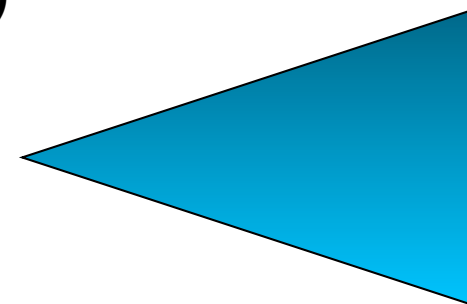


**CERN LHC Atlas detector
generates 10^5 more data than is
stored**



Biology stores * 1000 than detected (Genoscope)

**A need to filter the data at every stop along the
way using strategy appropriate to a particular
experiment/analysis**



Current discussion on data normalization

Thank you