# 1. Introduction to Bioinformatics

Ricardo Gonzalo[1] ,Mireia Ferrer ,Álex Sánchez[1,2]

1 Unitat d'Estadística i Bioinformàtica (UEB) VHIR

2 Departament de Genètica Microbiologia i EStadística

*We are drowning in information and starved for knowledge*

John Naisbitt

*Who on efficient work is bent, Must choose the fittest instrument*.

Goehthe (Fausto)

# What is Bioinformatics?

# A (first) definition

*Bioinformatics is the application of computer technology to the management of biological information.*

*Computers are used to gather, store, analyze and integrate biological information.*
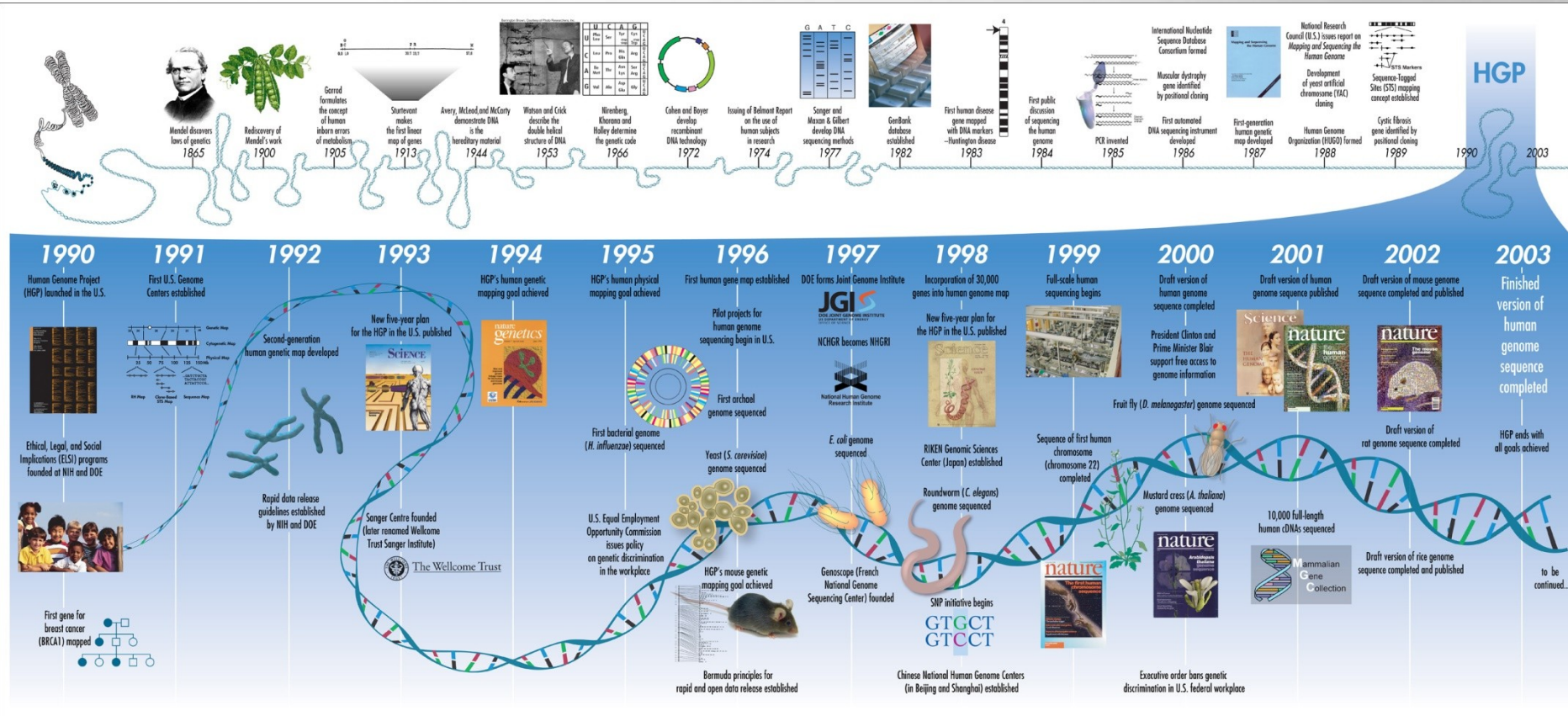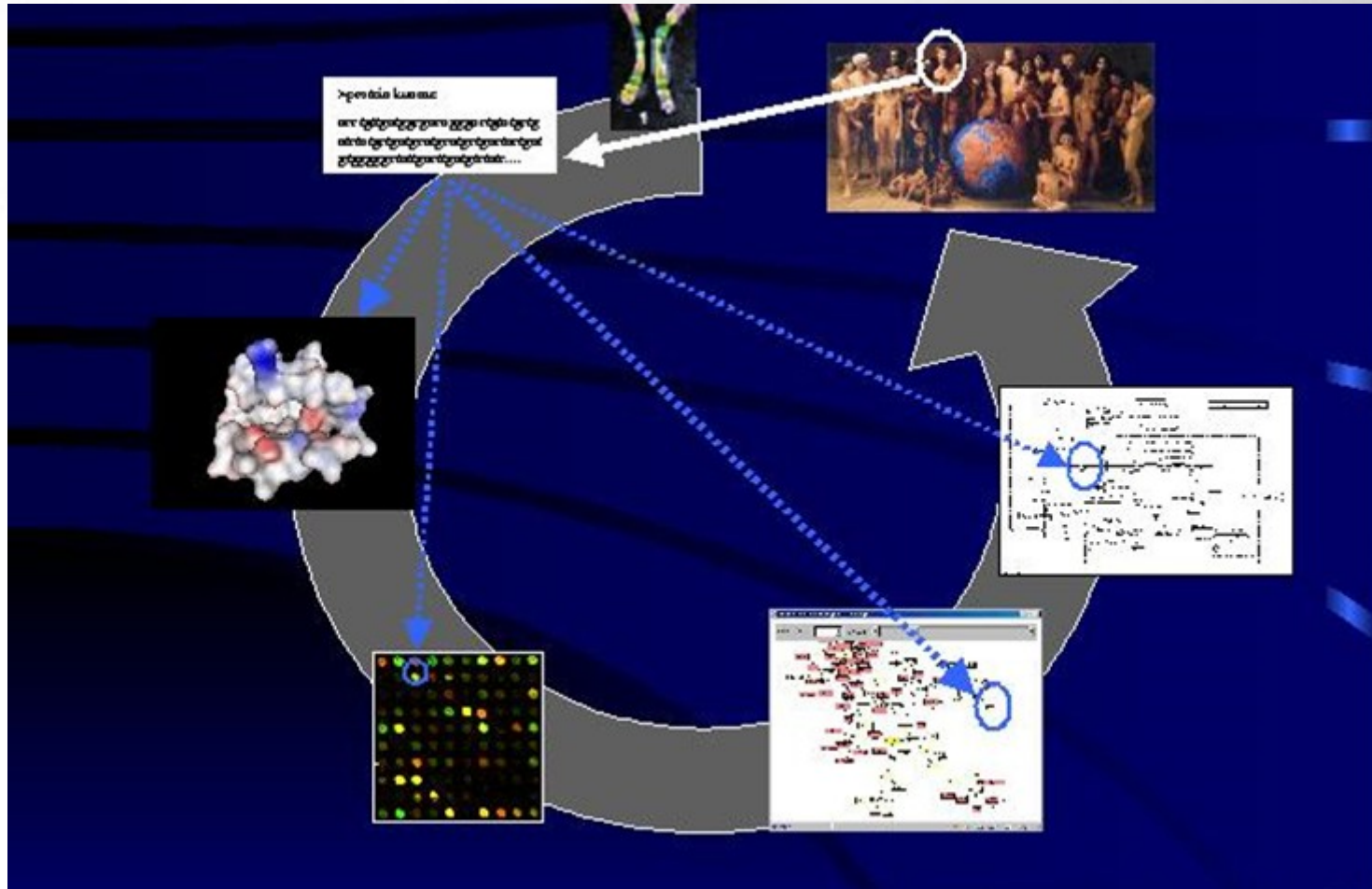
# A historical approach

- The term appeared in the 70's
- It became popular/important with the development of the human genome project
- Bioinformatics is entering a big data era that will foster new possibilities.
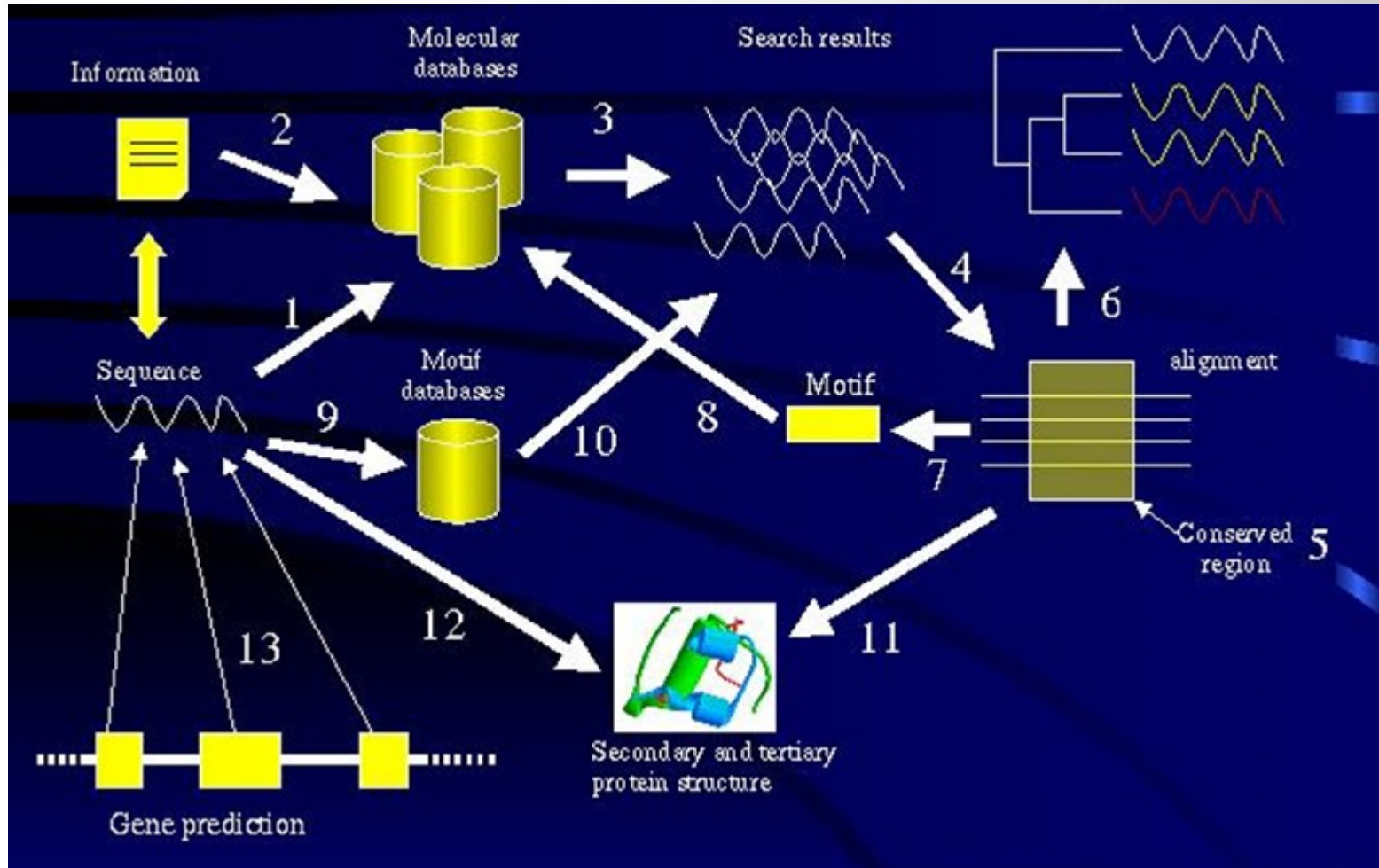
# The Human Genome Project

# Pre genomics era vision in the lab
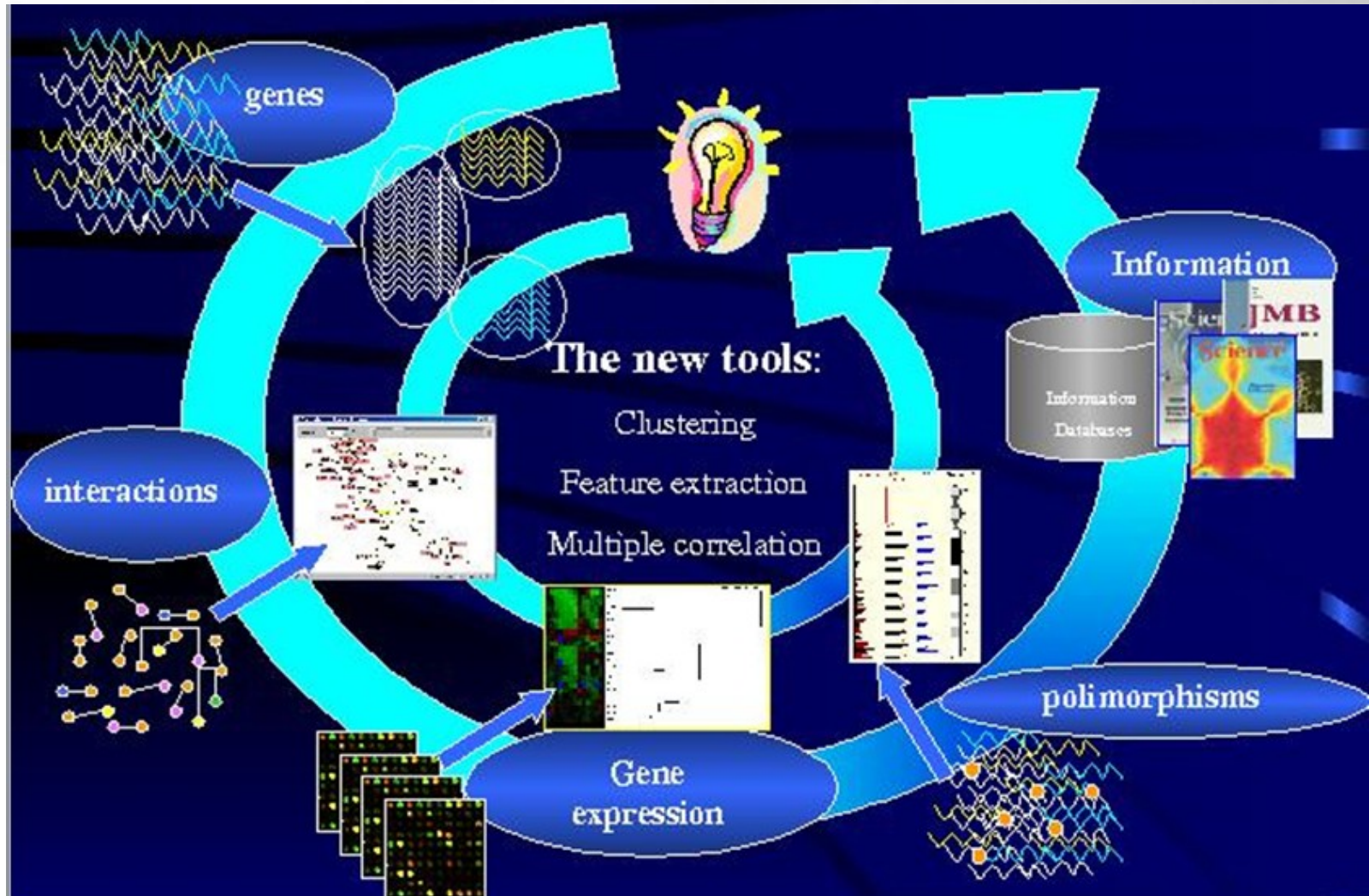


**Adapted from a presentation by J. dopazo**

# Bioinformatic analysis


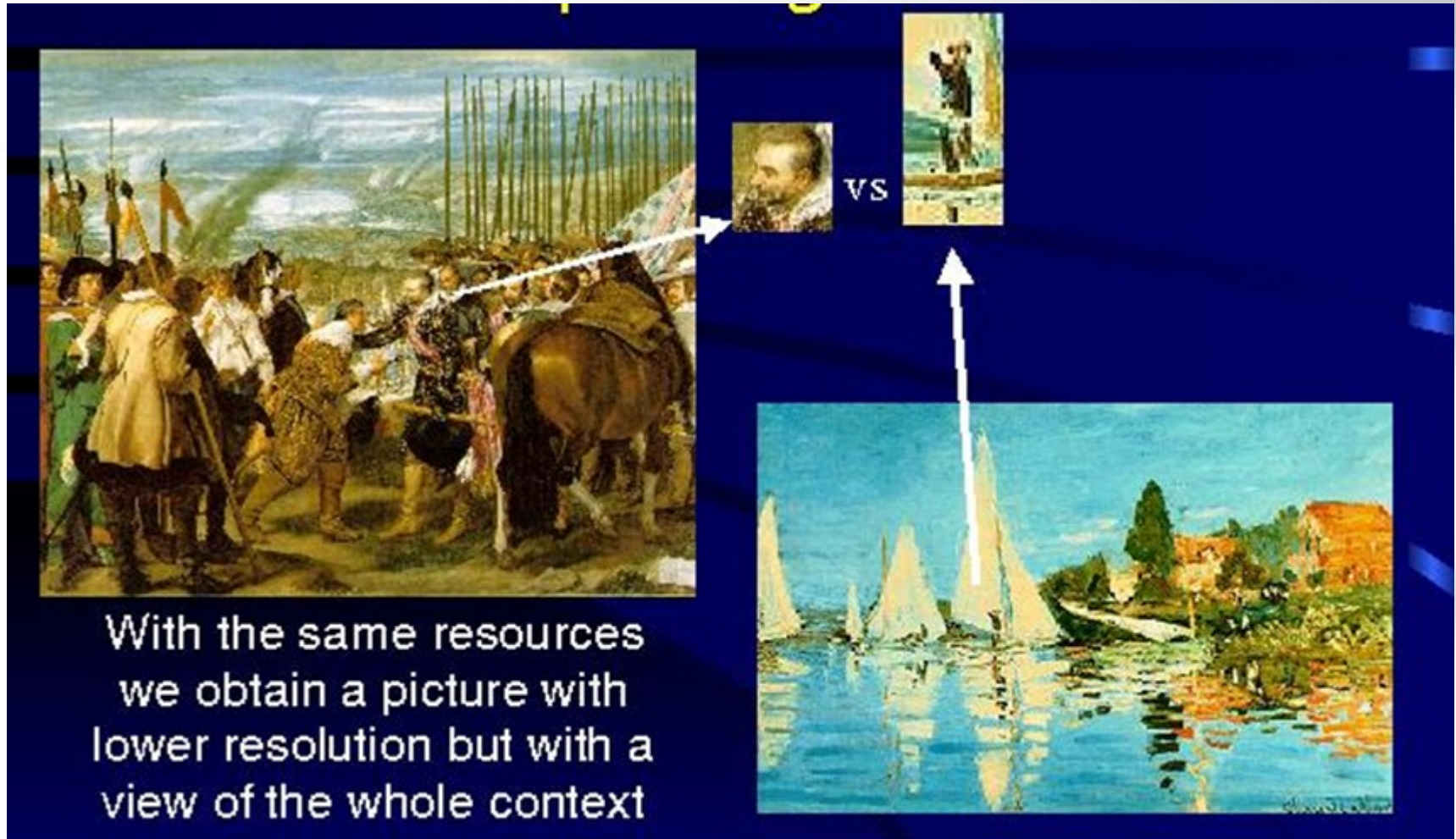
**Adapted from a presentation by J. dopazo**

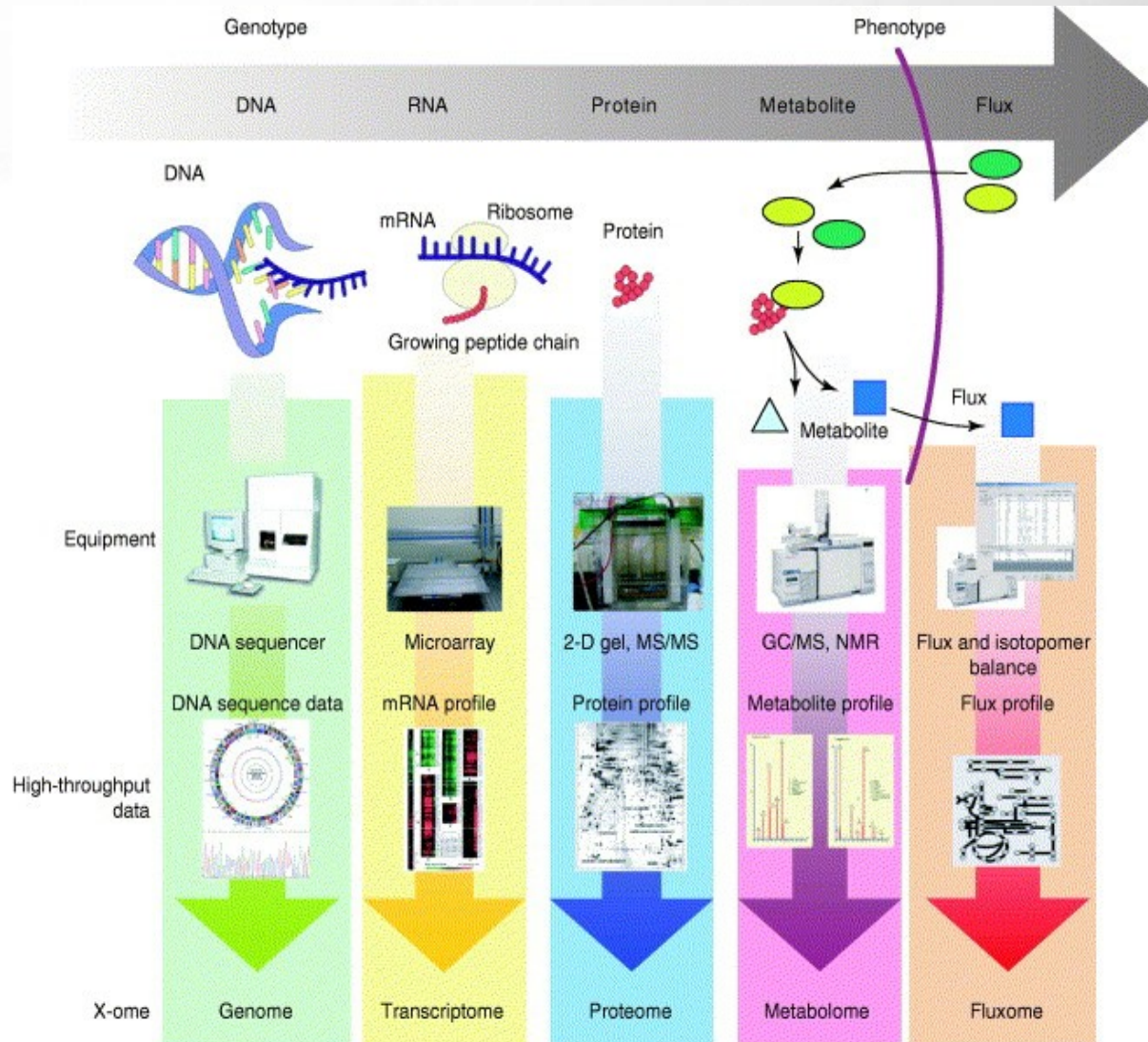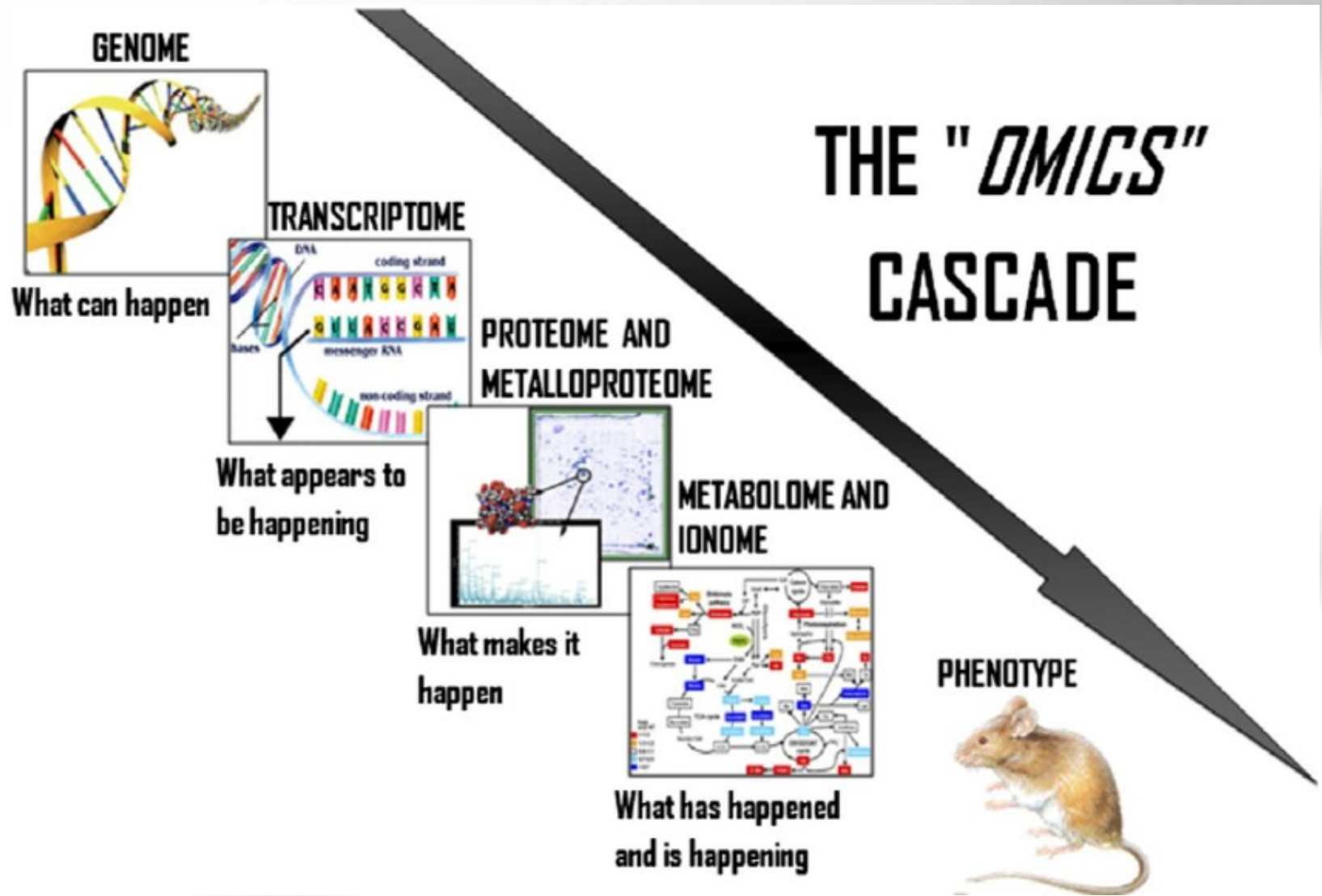# Post-genomic vision



**Adapted from a presentation by J. dopazo**

# The (first) paradigm shift



With the same resources we obtain a picture with lower resolution but with a view of the whole context

**Adapted from a presentation by J. dopazo**

# Omics technologies



TRENDS in Biotechnology

# The "Omics cascade"

# New paradigm shift: *Integromics*



H NMR metabolites

"NGS-Sequences

Affy Transcriptome

LC-MS proteomicss

"Non-omic" markers

# "Nextgen" sequencing revolution

# New paradigm shift: *Precision Medicine*

# Bioinformatics domains

- Information management
  - Databases, databanks
  - Algorithms and tools for database querying and searching
- Information Modelling
  - Protein structure characterization
- Analysis & Interpretation of results
  - Genome sequencing and analysis
  - Comparative genomics
  - Transcriptomics and gene expression
  - Metagenomics
  - Proteomics, metabolomics, …
- Biological system modeling

# Information management

# Analysis and interpretation

A G A G T T C T G C T C G
A G G G T T A T G C G C G

# Biological system modeling

# Integrative bioinformatics

# In summary...



**Data**
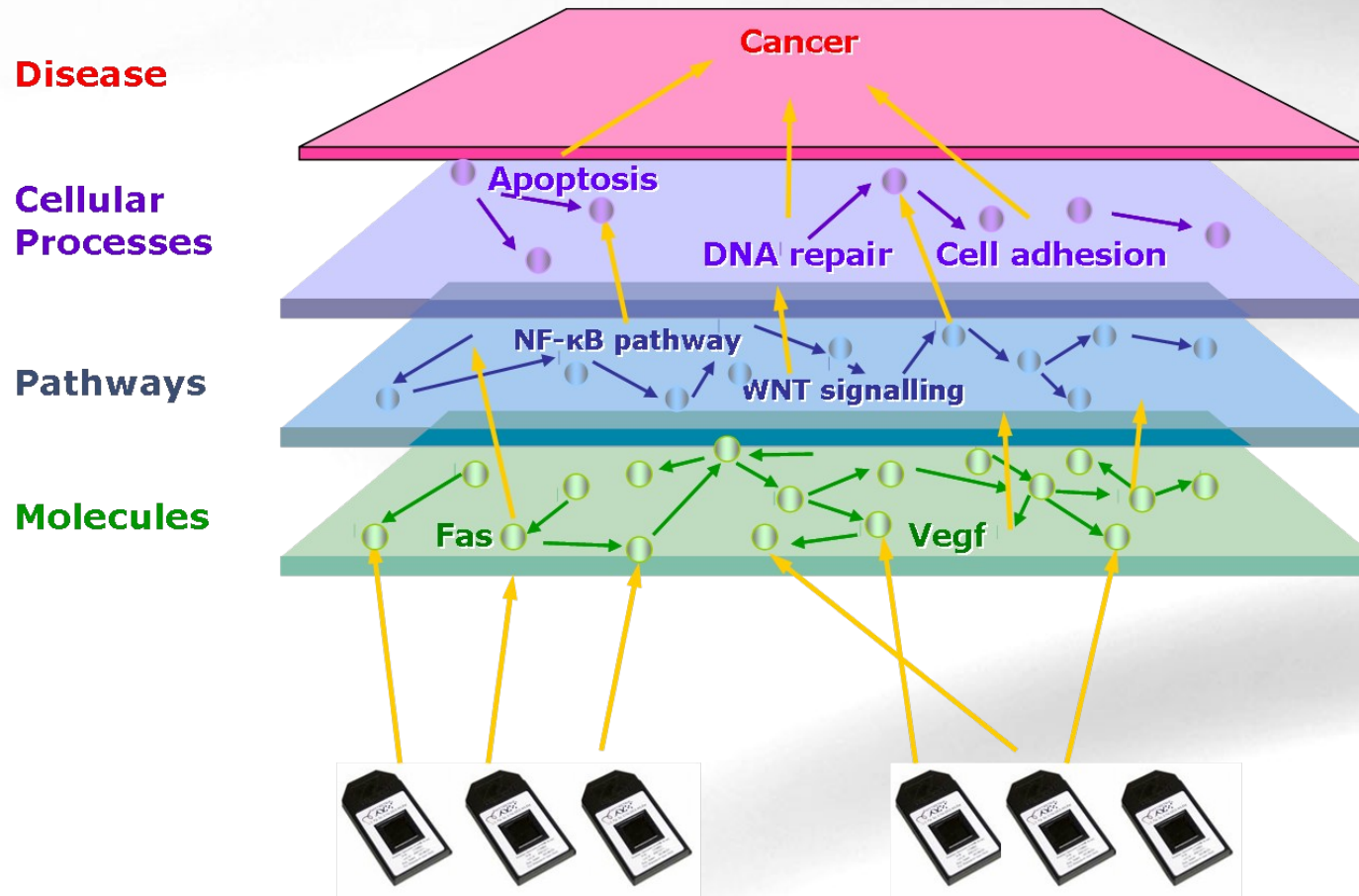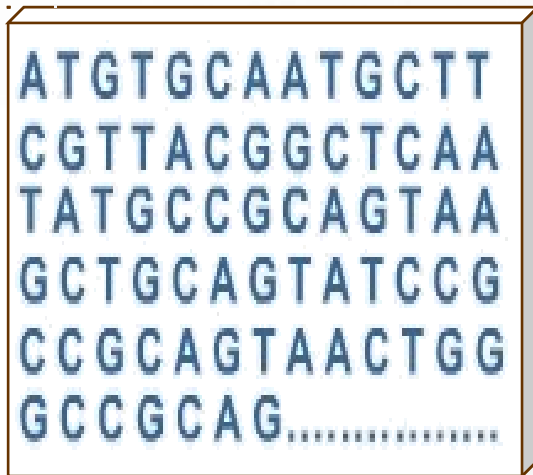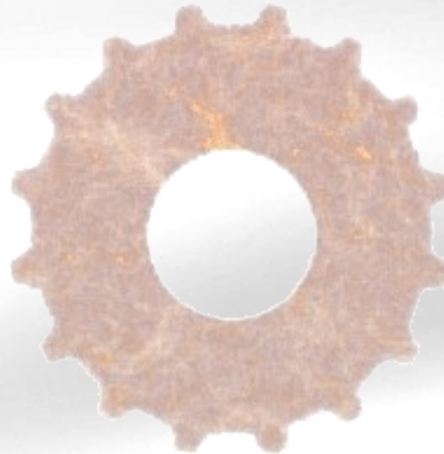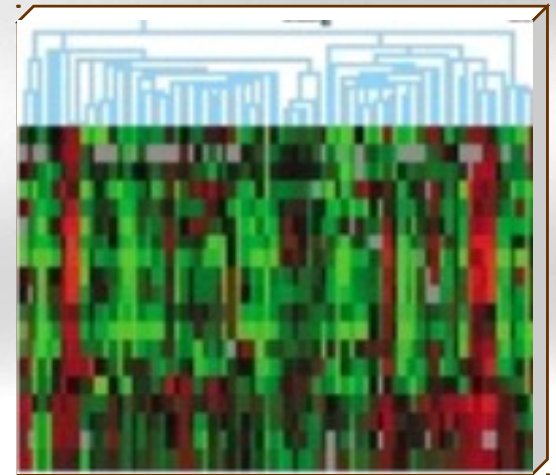
**Bioinformatics methods and resources**

**Knowledge**

# How does one do bioinformatics?

http://biomedicalcomputationreview.org/content/landscape-bioinformatics-education

# "Doing bioinformatics"

- Bioinformatics analyses
  - Database searching/querying
  - Sequence analysis, Omics data analysis
  - Systems biology
- Can be done differently
  - From console-based systems
    - Using scripts (perl/python/R) for automating processes
    - Doing data analysis with R
  - Or working with graphical/web interfaces to do (almost) the same things
- Each user has a different preferred approach

# What does a bioinformatician know?

Must have "good background" in

- Some biological discipline
    - Molecular biology, biochemistry, evolution …
- Computer science
    - Operating systems: Linux
    - Programming languages: Python, R, Perl
    - Databases SQL
    - Web development: HTML, PHP, …
- Some "quantitative" science
    - Mathematics, Physics, Statistics

Ideally 1, ½, ⅓  from the previous three!!

# A typical "bioinformatics user"

**Welch L, Lewitter F, Schwartz R, Brooksbank C, Radivojac P, et al. (2014) Bioinformatics Curriculum Guidelines: Toward a Definition of Core Competencies. PLoS Comput Biol 10(3): e1003496. doi:10.1371/journal.pcbi.1003496**
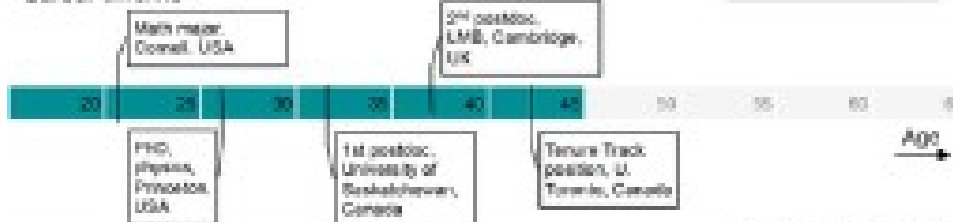**http://journals.plos.org/ploscompbiol/article?id=info:doi/10.1371/journal.pcbi.1003496**

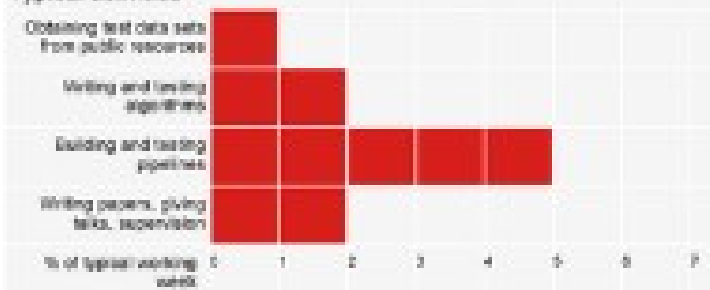# A typical "bioinformatics scientist"
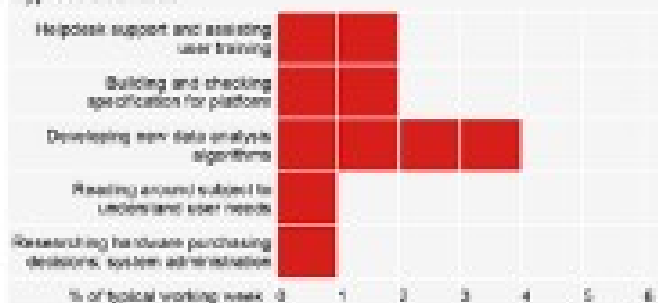
# A typical "bioinformatics engineer"

**Welch L, Lewitter F, Schwartz R, Brooksbank C, Radivojac P, et al. (2014) Bioinformatics Curriculum Guidelines: Toward a Definition of Core Competencies. PLoS Comput Biol 10(3): e1003496. doi:10.1371/journal.pcbi.1003496**
**http://journals.plos.org/ploscompbiol/article?id=info:doi/10.1371/journal.pcbi.1003496**

# In summary, bioinformatics ...

- Was born with
  - Development of new technologies
  - Its application for generating –increasingly huge- of big masses of biological data.

- *Has become now an interdisciplinar science encompassing all aspects of the Acquisition, Processing, Distribution, Analysis, Integration and Interpretation of biological information.*