# BIOINFORMÁTICA Y BIOLOGÍA COMPUTACIONAL

**MADRID, 4 al 31 de Julio de 2007**

http://www.pdg.cnb.uam.es/cursos/Complutense2007/index.html

## BIOINFORMÁTICA Y BIOLOGÍA COMPUTACIONAL

*DIRECTORES:*

- **Luis Vázquez Martínez.**

  Catedrático Universitario. Dpto. Matemática Aplicada. Facultad de Informática, Universidad Complutense de Madrid.

- **Federico Morán Abad.**

  Profesor Titular del Departamento de Bioquímica y Biología Molecular I, Universidad Complutense de Madrid.

- **Alfonso Valencia Herrera**

  Profesor de Investigación del Consejo Superior de Investigaciones Científicas. Centro Nacional de Investigaciones Oncológicas. Madrid

*FECHAS:*
4 al 31 de Julio de 2007

*HORARIO:*
De 9:00 a 14:00, todos los días lectivos.

*LUGAR DE CELEBRACION* (Clases teóricas y prácticas):
Facultad de Informática de la UCM
Av. Complutense s/n - Ciudad Universitaria.
Aula 2 de Informática 2ª Planta.

*DURACION*:
100 horas

# The Race to Cash In On the Genetic Code

## As DNA Gives Up Secrets, Pursuers Turn Entrepreneurs

### Rival Strategies

The leading genomics companies have contrasting business models: two aim to develop drugs (■); two aim to sell data about genes (■).

### INCYTE PHARMACEUTICALS

**HEADQUARTERS** Palo Alto, Calif.

**BUSINESS** Provides data bases of gene sequences and related analytical software by subscription. Manufactures and sells drug discovery tools.

**ACCOMPLISHMENTS** More than 25 major pharmaceutical and biotechnology subscribers. First genome company to report earnings from operations.

**MARKET CAPITALIZATION** $713 million

'98 REVENUE $135 million

'98 OPERATING INCOME $12 million

### MILLENNIUM PHARMACEUTICALS

**HEADQUARTERS** Cambridge, Mass.

**BUSINESS** Provides genetic targets for partners' drug discovery and agricultural programs. Develops small-molecule drugs; proteins and antibodies; gene-based diagnostic tests.

**ACCOMPLISHMENTS** Partnerships with Bayer, Monsanto, Pfizer, Eli Lilly, American Home Products, Astra and Roche.

**MARKET CAPITALIZATION** $2.15 billion

'98 REVENUE $134 million

'98 OPERATING INCOME −$8 million

### HUMAN GENOME SCIENCES

**HEADQUARTERS** Rockville, Md.

**BUSINESS** Provides genetic targets for partners' drug discovery programs. Develops proprietary proteins, antibodies and gene therapy drugs.

**ACCOMPLISHMENTS** Established first major genomic partnership, receiving $125 million from SmithKline Beecham. Has three drugs in clinical trials.

**MARKET CAPITALIZATION** $1.55 billion

'98 REVENUE $30 million

'98 OPERATING INCOME −$32 million

### CELERA GENOMICS GROUP

**HEADQUARTERS** Rockville, Md.

**BUSINESS** Provides genomic and biological information. Plans to offer complete genomes of human and other species, plus analytical software.

**ACCOMPLISHMENTS** Partnerships with Novartis, Pharmacia & Upjohn and Amgen.

**MARKET CAPITALIZATION** $650 million

'98 REVENUE* $6 million

'98 OPERATING INCOME* −$16 million
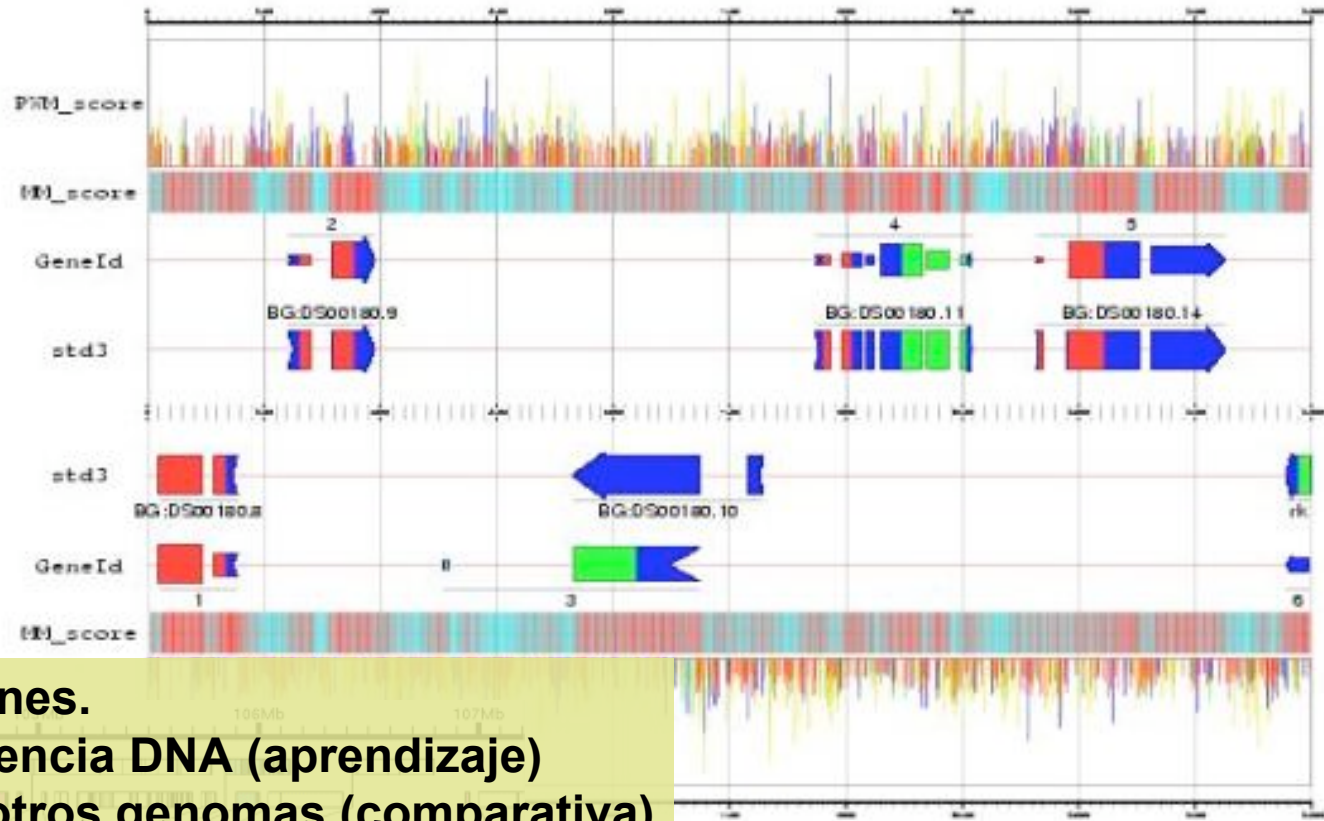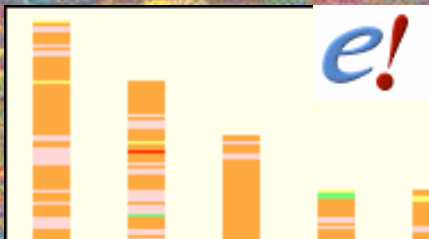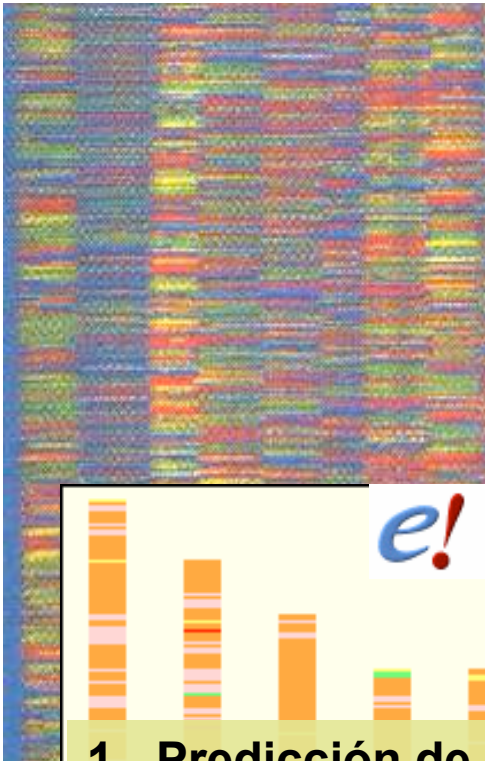
*Six months ended Dec. 31
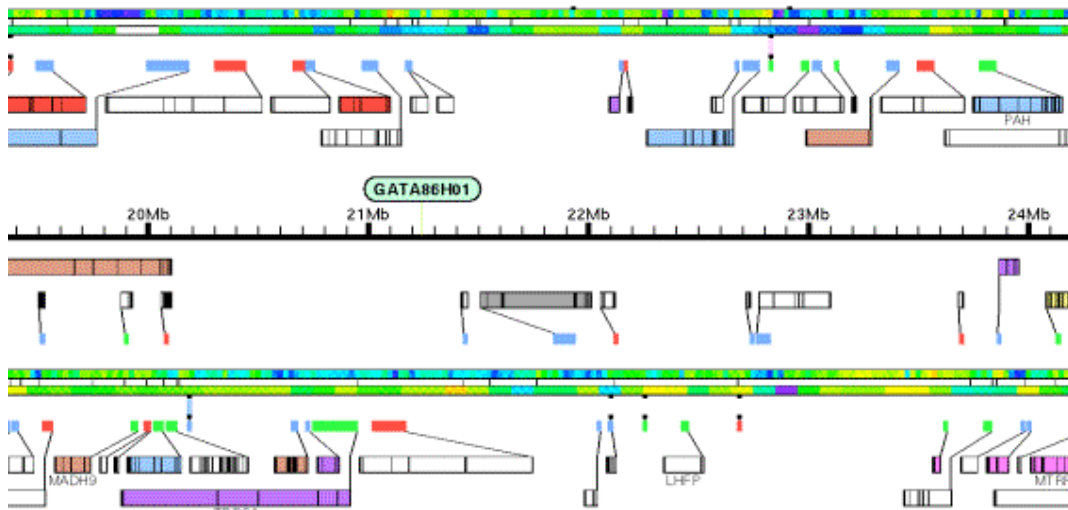
The New York Times/
Illustration by Chris Gould

Source: *New Yorker*

1.- Predicción de genes.
- señales en la secuencia DNA (aprendizaje)
- comparación con otros genomas (comparativa)
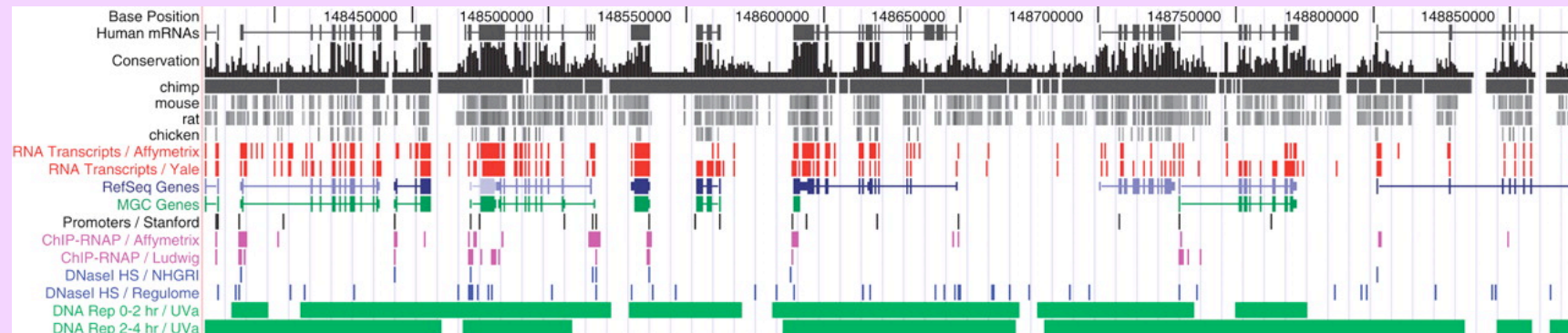
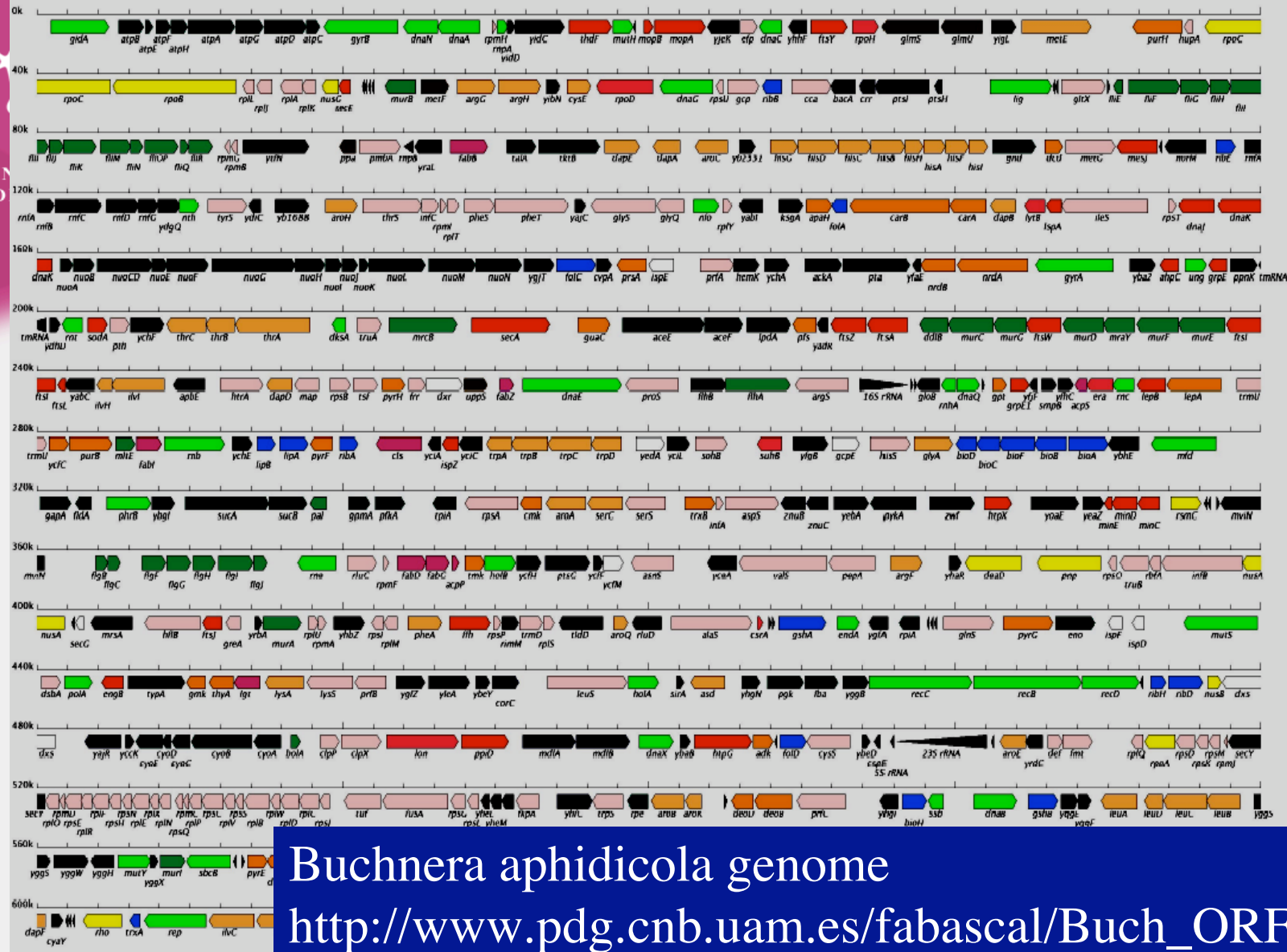Gene Finding / Gene Prediction

# The ENCODE Project

The Encyclopaedia of DNA Elements (ENCODE) is an NIH-backed multi-million dollar project, that brings together an international consortium of scientists in industry and academia with the aim of identifying all the functional elements in the human genome.

The efforts of the ENCODE consortium are <u>focused entirely on DNA aspects of the project</u>, such as transcription sites, TRANSFRAGS, non–protein-coding genes and sequences that mediate chromosome structure and dynamics.



The pilot project is studying 30MB from 44 regions comprising 1% of the genome.

15MB come from regions chosen for their scientific interest, the rest are chosen via a stratified random sampling method.

*by Michael Tress, CNIO*

Buchnera aphidicola genome
http://www.pdg.cnb.uam.es/fabascal/Buch_ORFand_ww

**Proceso común de análisis de secuencias de proteínas**

INSTITUTO NACIONAL DE BIOINFORMÁTICA

**Proteína Problema**

Base de datos de proteínas

**Proteínas similares**

**Grupos de proteínas relacionadas**

**Proteínas homólogas**

ras (H. sapiens)
ras2 (H. sapiens)
ras (M. musculus)
ras (C. elegans)

rab (H. sapiens)
rab (M. musculus)
rab (C. elegans)

**Predicción de función**

UCM, 07
Alfonso Valencia  CNIO

*by  F, Abascal*

**MARVEL: a conserved domain involved in membrane apposition events.**
**Sanchez-Pulido et al., TIBS 2002**



MARVEL novel domain
 four transmembrane-helix architecture
 - myelin and lymphocyte (MAL),
  - physins,
 - gyrins and
 - occludin families.

 cholesterol-rich membrane apposition events
- biogenesis of vesicular transport carriers
- tight junction regulation.

Human diseases: schizophrenia and inflammation.

Alfonso Valencia CNIO

# PREDICCION ESTRUCTURA PROTEINAS

**2.- Predicción de estructura de proteínas.**
**- como problema físico: simulaciones, campos de fuerza**
**- Informático: aprendizaje (secuencia > estructura conocidas),**
**estadísticas, librerías**

CASP - Asilomar 94, 96, 98, 00, 02 - Italia 04

INB

INSTITUTO NACIONAL DE BIOINFORMÁTICA

MAKEFGIPAAVAGTVLNVVEAGGWVTTIVSILT
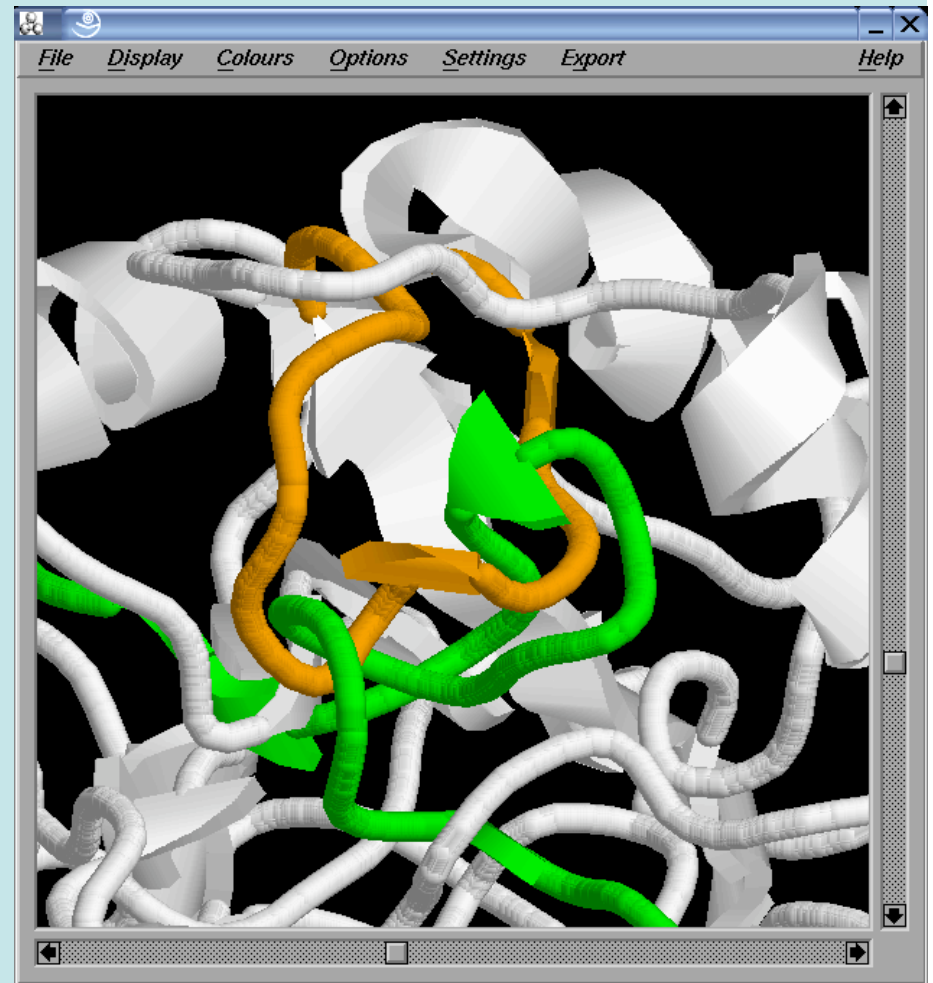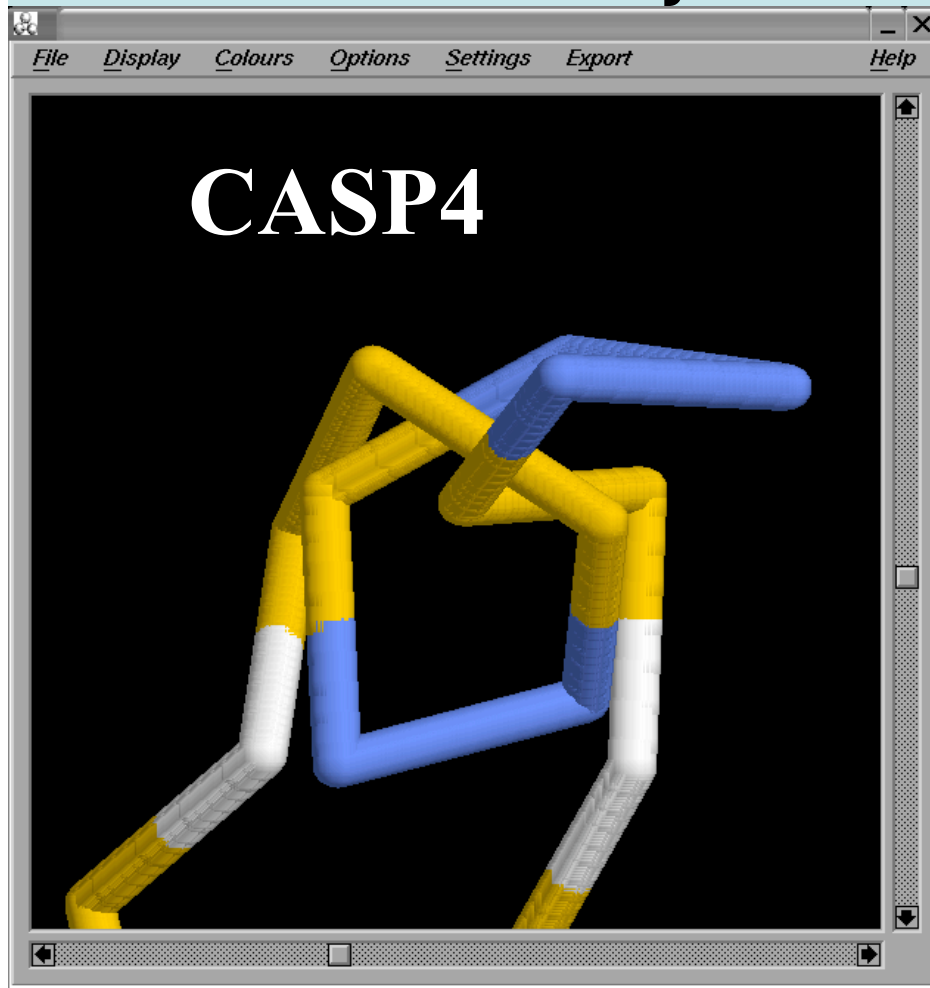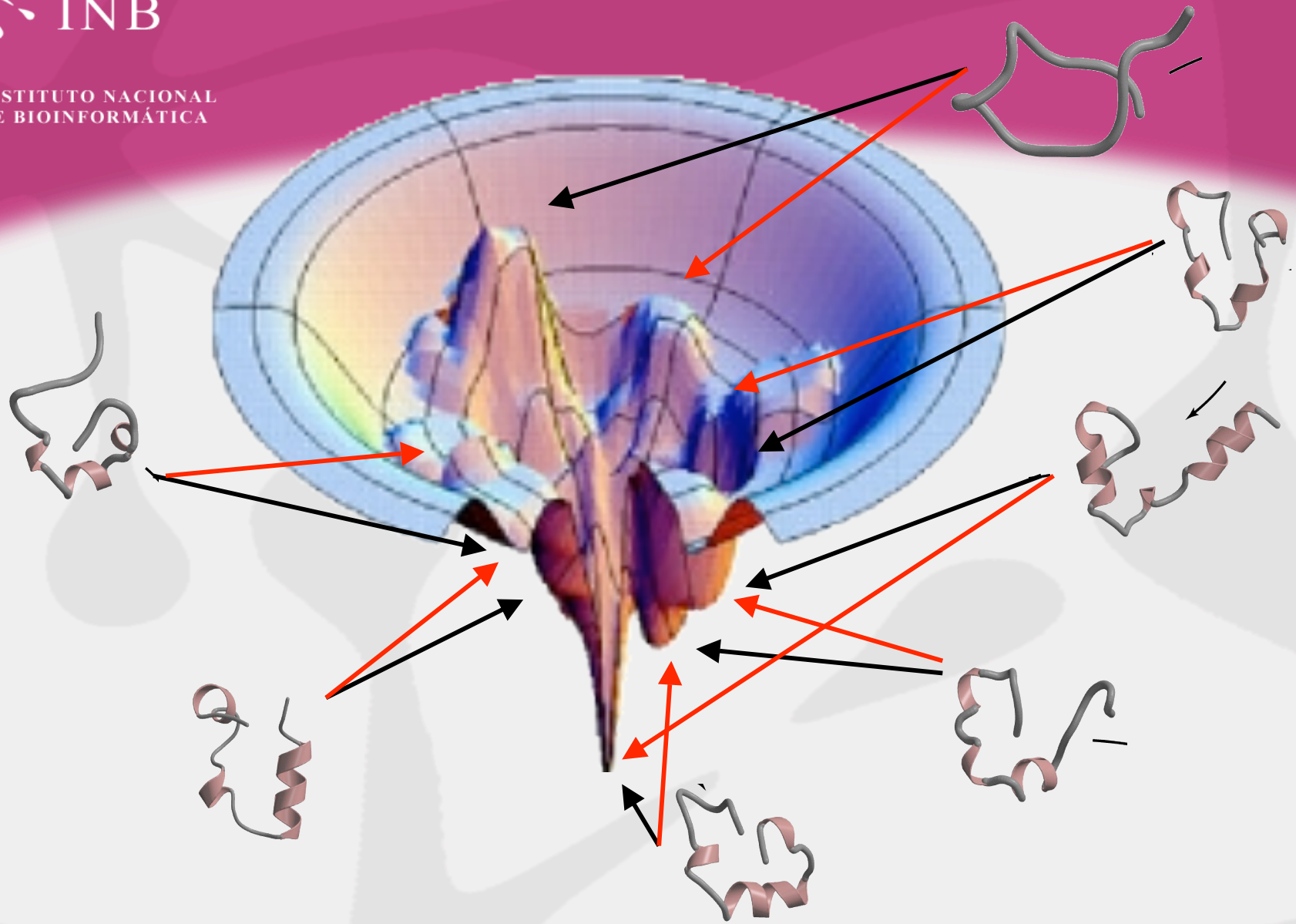AVGSGGLSLLAAAGRESIKAYLKKEI
KKGKRAVIAW

Predictores

Cristalógrafos

1/3 predicciones correctas

Databases

Algoritmos

Evaluación evaluation

Jueces

UCM, 07
Alfonso Valencia  CNIO

By  D. Devos
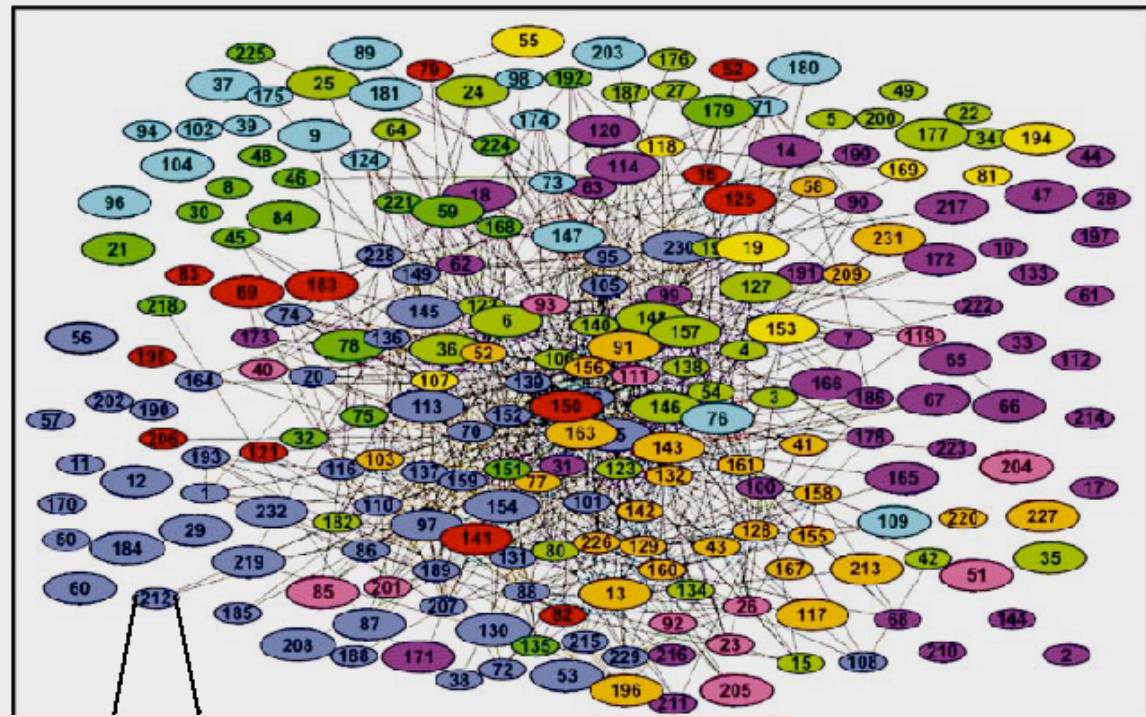
# You think you have seen knots ...

*Cellzome, Nature 2002*

## feedforward loop

X → Y → Z

crp → araC → araBAD

## single input module (SIM)

X → Z₁ Z₂ ... Zₙ

argR → argCBH, argD, argE, argF, argI

## dense overlapping regulons (DOR)

X₁ X₂ X₃ ... Xₙ → Z₁ Z₂ Z₃ Z₄ ... Zₘ

carbon utilization DOR · osmotic stress DOR · stationary phase DOR · DNA metabolism DOR · drug and superoxide DOR

maltose · flagella · heat shock

legend:
- ⬭ transcription factor (TF)
- dense overlapping regulons (DOR)
- ☐ single input module (SIM)
- △ coherent feedforward loop
- ▲ incoherent feedforward loop
- ○ single operon
- global TF
- — postive regulation
- — negative regulation
- — dual regulation
- ◇ multi-input module

Shen-Orr, Milo, Mangan & Alon, (2002).
*Network motifs in the transcriptional regulation network of Escherichia coli.* Nature Genet 31, 64-68.

*Milo, Itzkovitz, Kashtan, Levitt, Shen–Orr, Ayzenshtat, Sheffer, Alon, (2004)*
***Superfamilies of Evolved and Designed Networks.***
*Science, 303, 1538–1542*

Triad Significance Profile

The triad significance profile (TSP) of networks from various

UCM, 07
Alfonso Valencia  CNIO

- **14 M publications**
- **10M accessible in Medline**
- **8M abstracts in Medline**

- **Increasing number of accessible journals (full text)**

- **Information organized in web repositories**

Source: MyGRID

UCM, 07
Alfonso Valencia CNIO

# Analisis sistemático de familias de proteinas utilizando tecnología web (biomoby, contexto INB)

12 CPUs = 1.5 min

Human genome (30.000 sequences)

Weekly updates

Other genomes

More complex workflows

UCM, 07
Alfonso Valencia CNI

# A Virtual Institute

**INB**

INSTITUTO NACIONAL
DE BIOINFORMÁTICA

Training
**F.Sanz - UPF**

Computational Node 1
**J. L. Gelpi  BSC**

Protein structure
**M-Orozco – PCB**

Central Node
Protein function
**A.Valencia –CNIO**

Computacional Node 2
**J.M. Carazo – PCM**

Genome Analysis
**R.Guigo – IMIM/CRG**

Company Node
**BIOAlma**
**eBioIntel**
**Applied Biosystems**

Functional Genomics
**J.Dopazo – CIPF**

Genetic Variability
**To be organized**

Integrative Bioinformatics
**O.Trelles – U. Malaga**

UCM, 07
Alfonso Valencia  CNIO

**GenBank Release 144.0 — October 15, 2004**

| Species | Genome size | Bases | Entries |
|---|---|---|---|
| Homo sapiens | 3,400,000,000 | 10,965,381,932 | 8,338,229 |
| Mus musculus | 3,454,200,000 | 6,774,229,516 | 6,071,679 |
| Rattus norvegicus | 2,900,000,000 | 5,644,398,664 | 985,869 |
| Danio rerio | 1,900,000,000 | 1,957,414,191 | 771,409 |
| Zea mays | 5,000,000,000 | 1,455,760,045 | 2,292,596 |
| Oryza sativa | 5,000,000,000 | 779,829,843 | 336,051 |
| Drosophila melanogaster | 180,000,000 | 754,291,835 | 482,102 |
| Bos taurus | 3,651,500,000 | 650,653,065 | 884,892 |
| Gallus gallus | 1,200,000,000 | 605,802,046 | 697,037 |
| Arabidopsis thaliana | 100,000,000 | 584,114,192 | 845,876 |
| Canis familiaris | 3,355,500,000 | 582,919,466 | 1,015,724 |
| Xenopus tropicalis | 3,355,500,000 | 464,627,420 | 560,218 |
| Pan troglodytes | 3,577,500,000 | 439,544,237 | 193,505 |
| Ciona intestinalis | 200,000,000 | 418,098,823 | 693,084 |
| Brassica oleracea | 759,500,000 | 403,897,848 | 595,915 |
| Macaca mulatta | 3,543,000,000 | 372,152,352 | 55,192 |
| Medicago truncatula | 400,000,000 | 327,501,222 | 348,369 |
| Triticum aestivum | 16,978,500,000 | 311,942,146 | 570,595 |
| Xenopus laevis | 3,100,000,000 | 298,427,470 | 455,955 |
| Caenorhabditis elegans | 100,000,000 | 283,634,604 | 309,719 |
| Total | | 43,194,602,655 | 38,941,263 |



UCM 07
Alfonso Valencia, CNIO

# MareNostrum

4.812 IBM PowerPC 970FX processors 2,2 GHz

(2.406 dual 64-bit processor blade nodes).

9,6 TB Main Memory
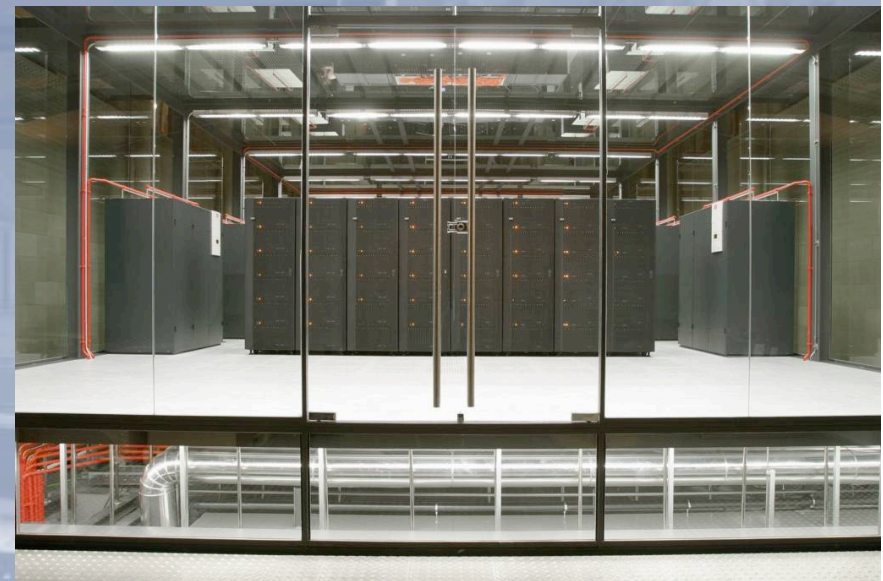 (4GB ECC 333 DDR memory per node).

42,35 Tflops (peak).

140 + 96 TB disk.

3 networks:

Myrinet, Gigabit, 10/100 Ethernet.
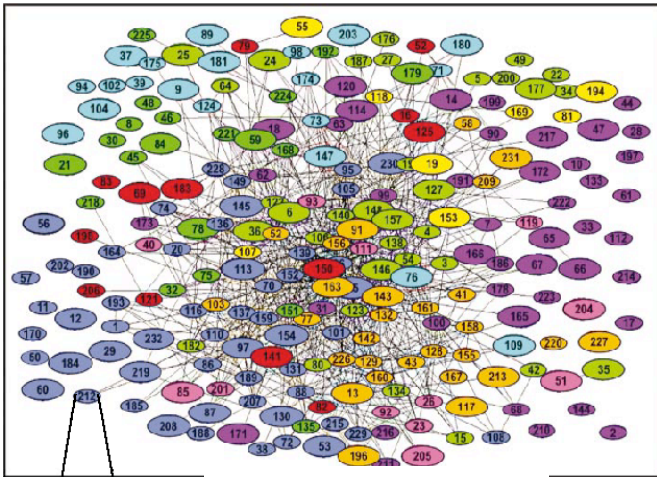
Linux 2.6 cluster (SuSe).

Diskless network support.



**Collaboration with the INB node at the BSC**

INB: Max. 40% capacity of the external use.
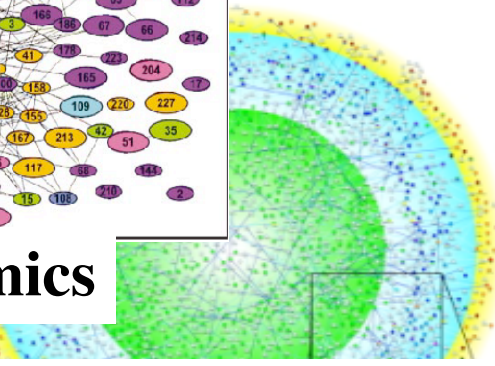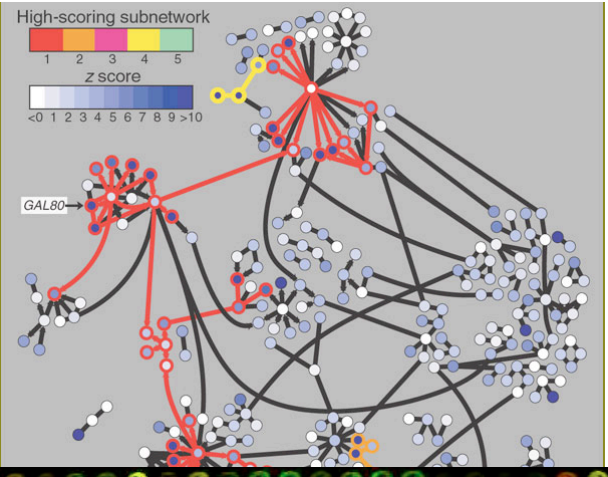Department of Comp. Biol. (M. Orozco director)
Access / parallelization help service

**Proteomics**

**literature**

**New Problems in the postgenomics era**

**PDG** Protein Design Group

www.cnio.es

INB · BIOSAPIENS NETWORK · ECCB

BioCreAtIvE · COMBIO

ENFIN — Enabling Systems Biology

**EMERGENCE COST CA on Synthetic Biology**

www.pdg.cnb.uam.es

www.inba.org

**Curso UCM07 Personasl CNIO / CNB.**

- **Secuencias: L.Sanchez, J.C.Sanchez, F. Abascal, J.M.G-Izarzugaza, A. Rojas**
- **Prediccion genes: J.J.Wesselink** - **Estructuras: D. de Juan, G. Lopez, M. Tress**
- **Text Mining: M.Krallinger** - **Arrays: J.C.Oliveros, G.Gomez**
- **Redes: F.Pazos, I. Cases** - **Web services: J.M. Fernandez, J.M.Rodriguez**
- **Introduccion y Resumen: A. Valencia**

**– BIOCREATIVE II.    biocreative.sourceforge.net**
**– CASP7    predictioncenter.org**

cnio

*Alfonso Valencia*