



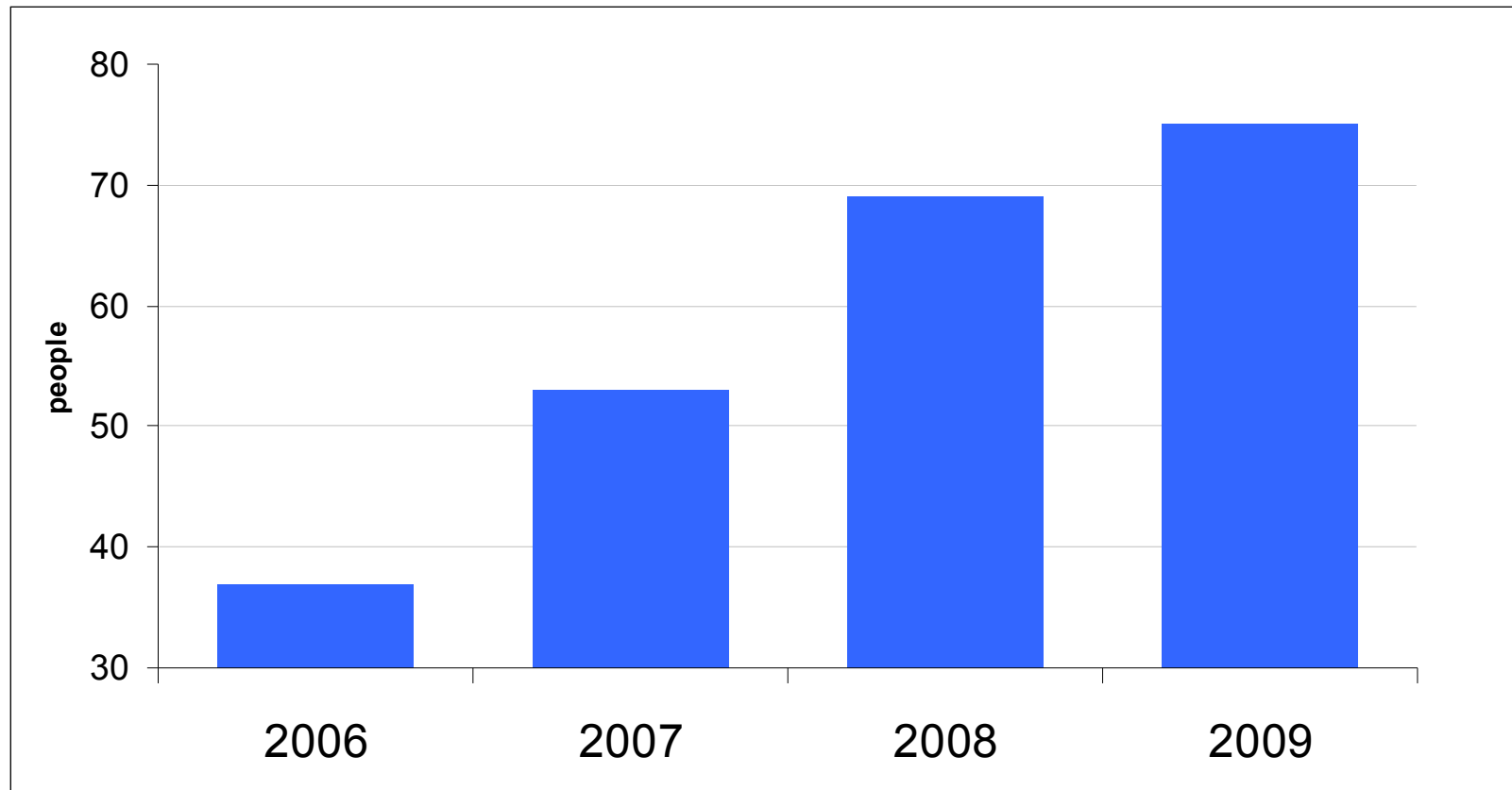
**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación

The Power of Computation in Life Sciences

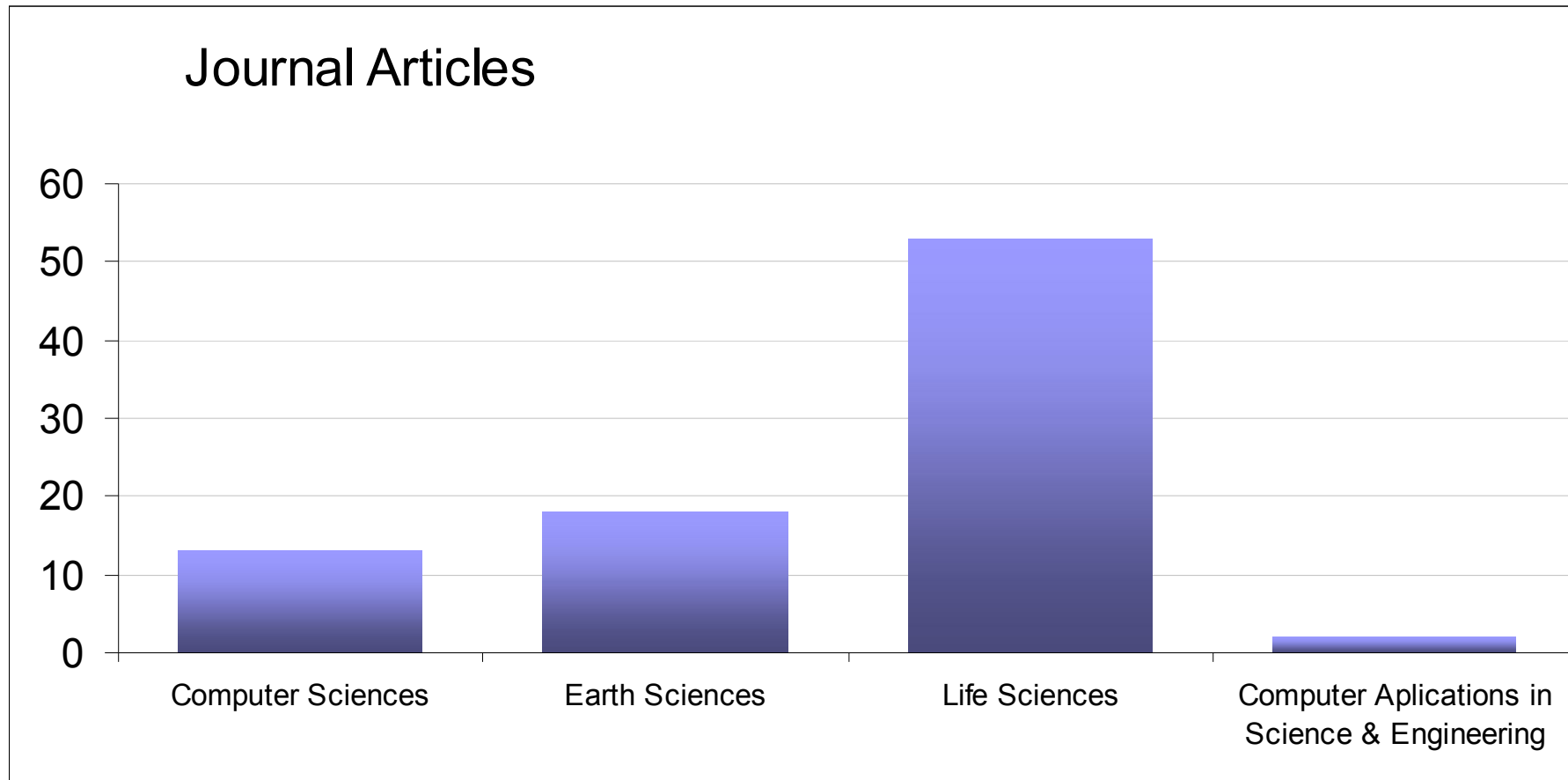
Ramon Goñi

Life Sciences Department





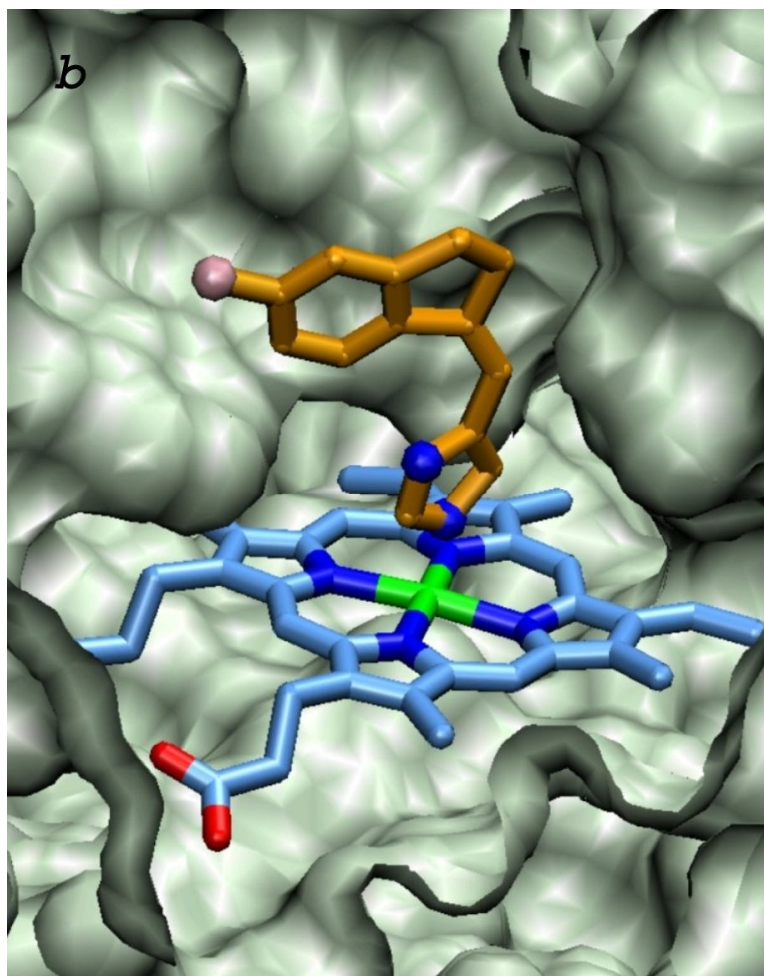
Life Science Research



Research Lines

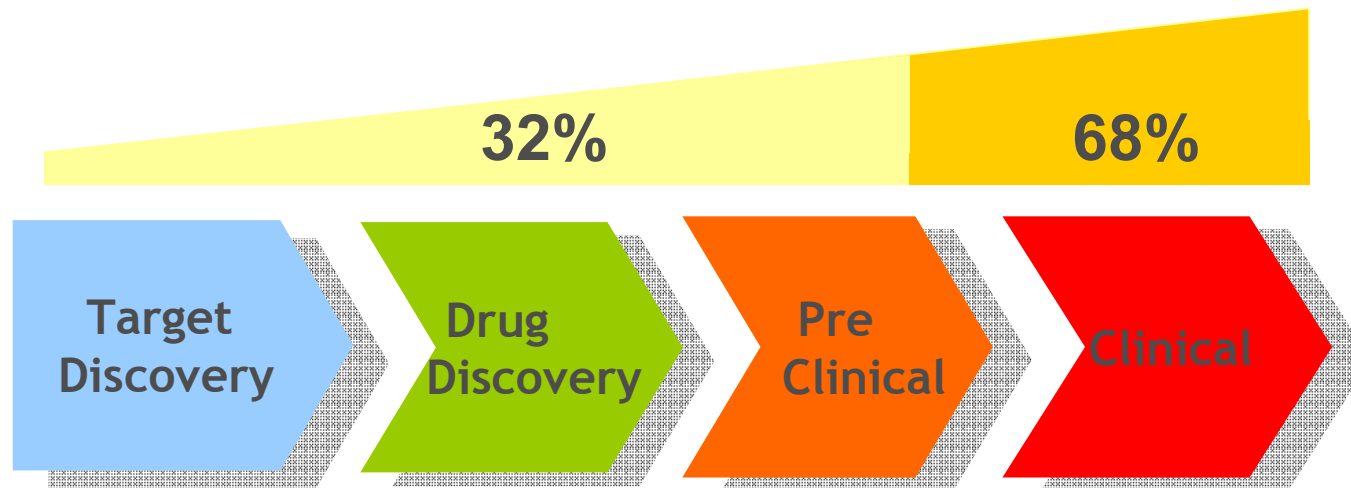
- Computer simulation (Drug Discovery)
 - Molecular Dynamics
 - Protein-Protein Interaction
 - Protein-Ligand Docking
- BioSupercomputing
 - Computational Biology under GPU & CELL
 - User-friendly computing access Web-Services
- Data analysis and Data Management (Target Discovery)
 - Next Generation Sequencing
 - Genomics & DNA structure

Computer Simulation & Drug Discovery



Drug Development

Cost: \$1.2B / drug

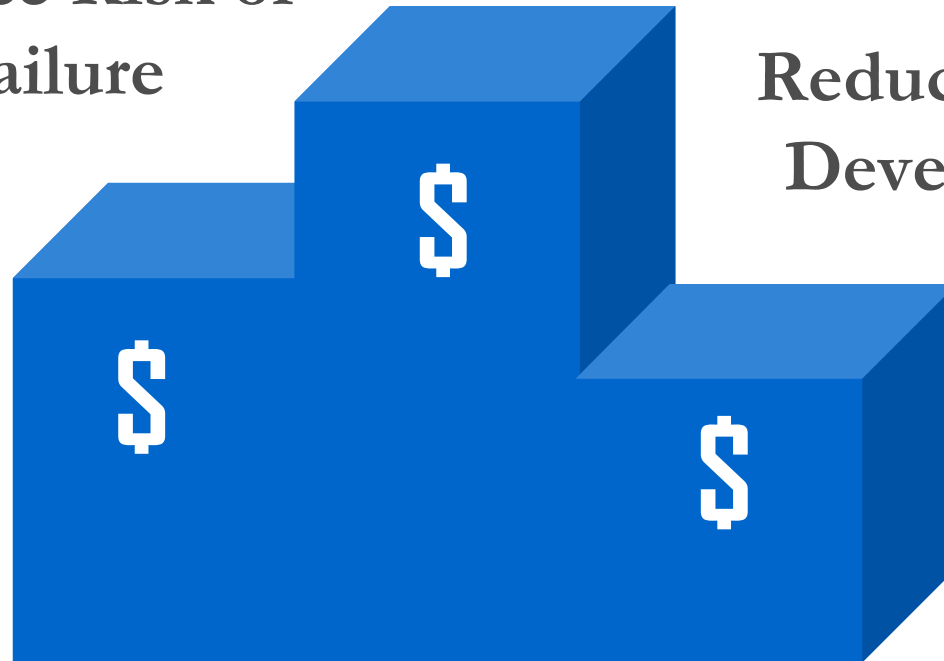


Computer-Based Drug Design

Improve Time
to Market

Reduce Risk of
Failure

Reduce Cost of
Development



Computer Simulation

- Why and when we use it
 - To validate a known model
 - As a cost-effective alternative
 - As the only realistic approach to solve a problem
- The structure of bio-molecules are hardly modeled.
The dynamics through experiments are only available for small molecules.
- There are different methods with different levels of complexity and realism

Molecular Simulation

PRECISION



Quantum Mechanics

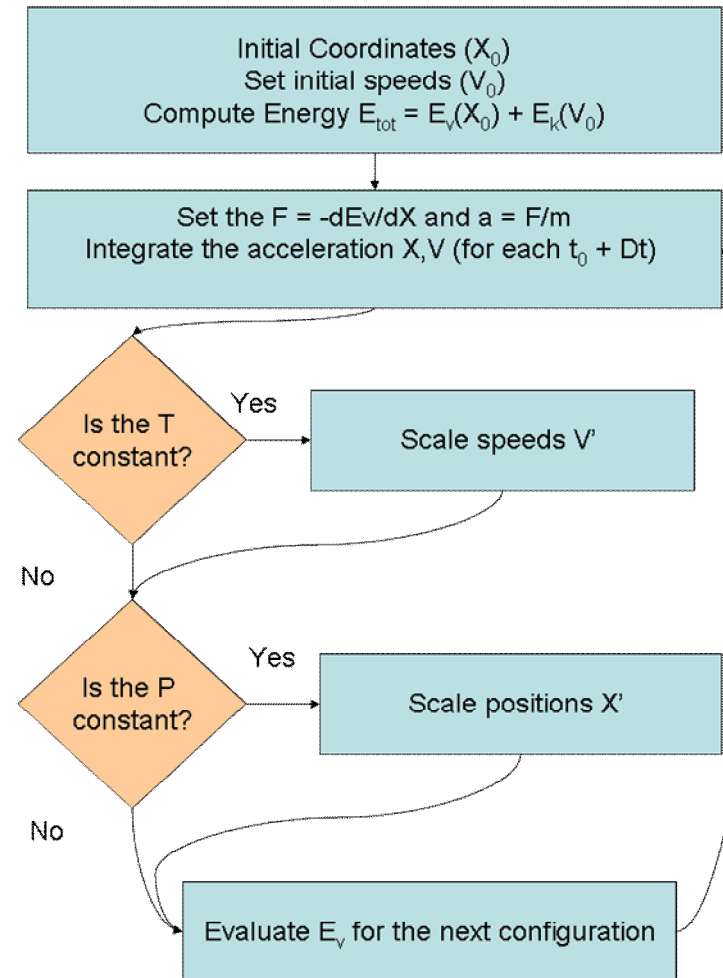
Molecular Dynamics (Newton motions)

Coarse Grained – MD (Pseudo-atoms)

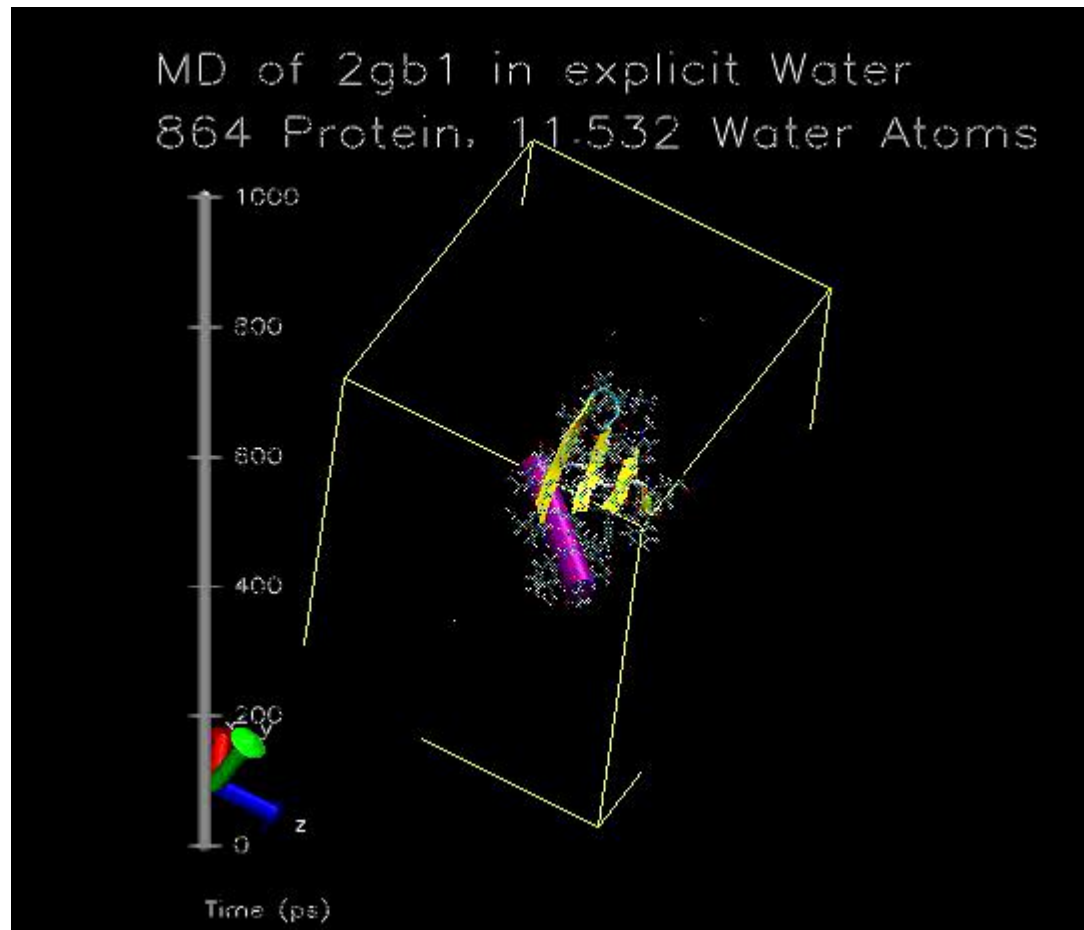
SPEED

Molecular Dynamics

- Atoms and molecules are allowed to interact for a period of time by approximations of known physics.



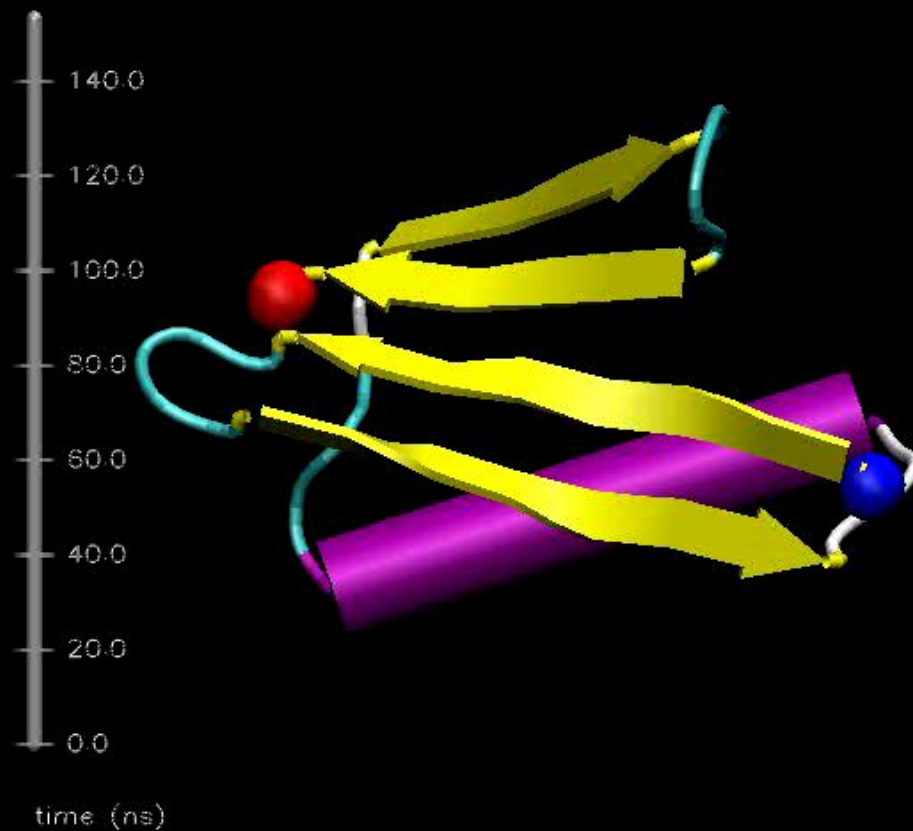
Molecular Dynamics of Solvated Protein



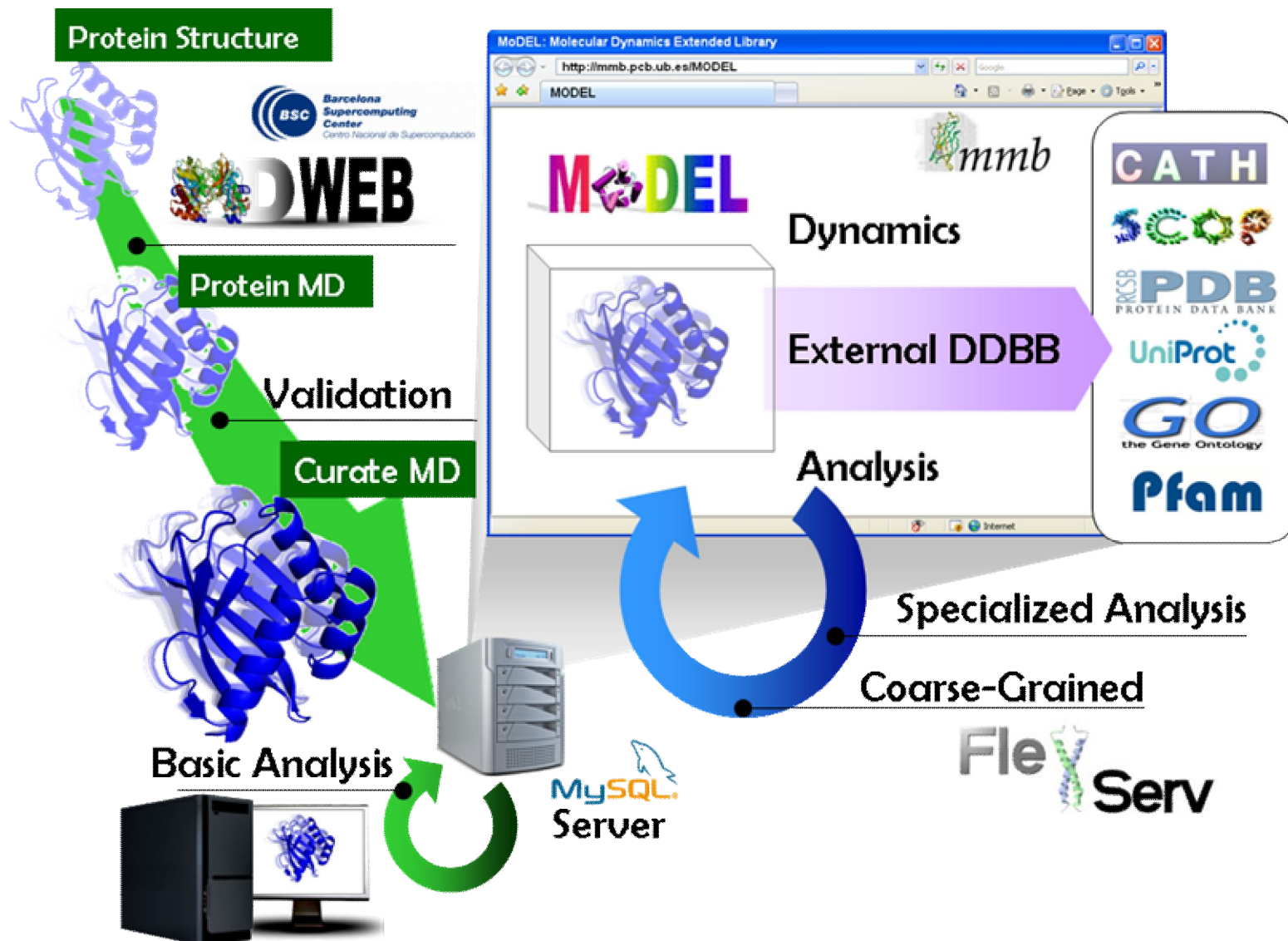
- Snapshot every femtosecond (10^{-15}s)
- System of 10^4 atoms
- 10 operations per atom pair-mate
- Using 16 processors we are able to simulate more than 10 nanoseconds (10^{-9}s) per day

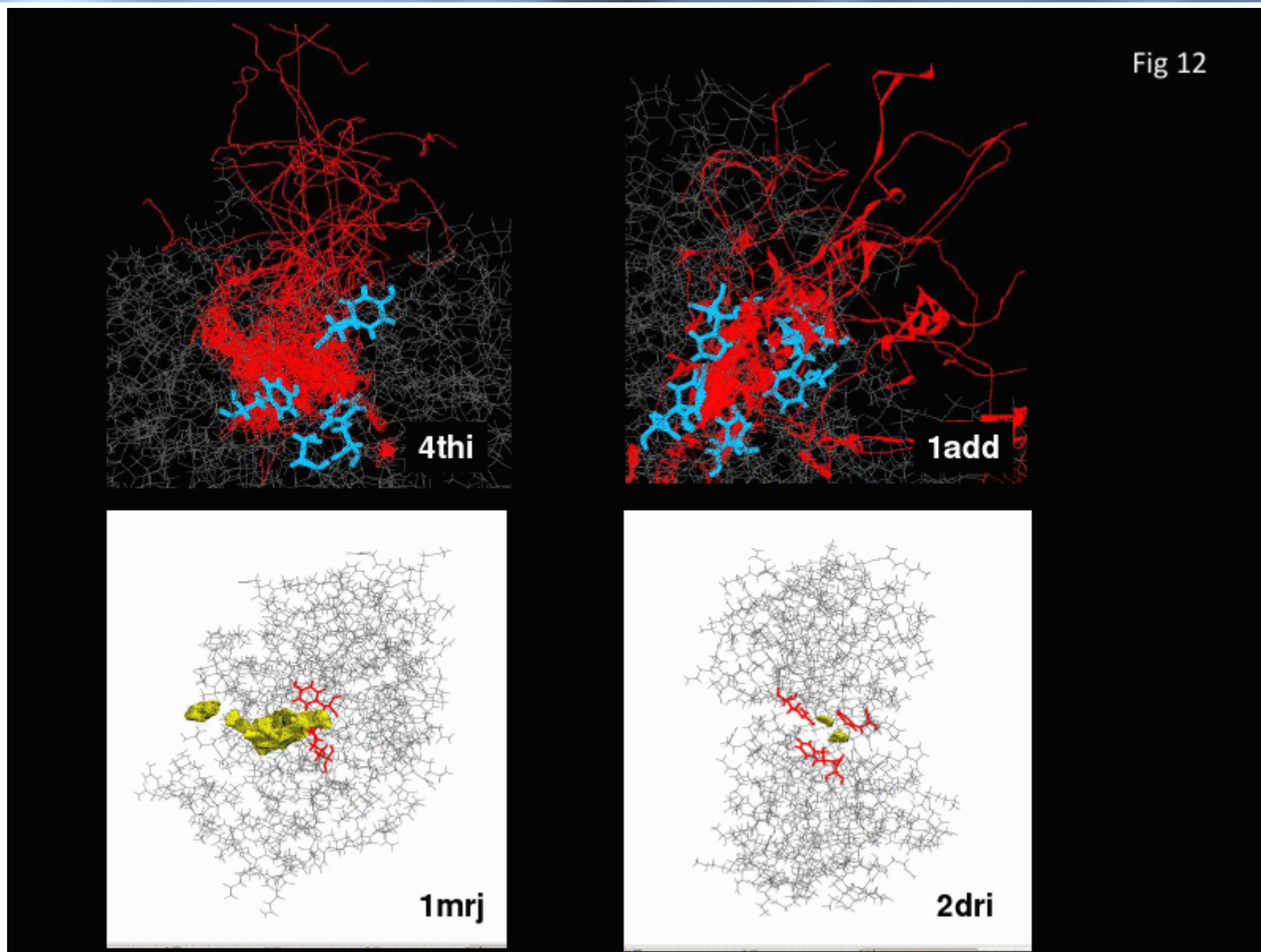
MD Scale

MD Simulation of 2gb1 in Vacuum
all HIS, GLU, and ASP protonated

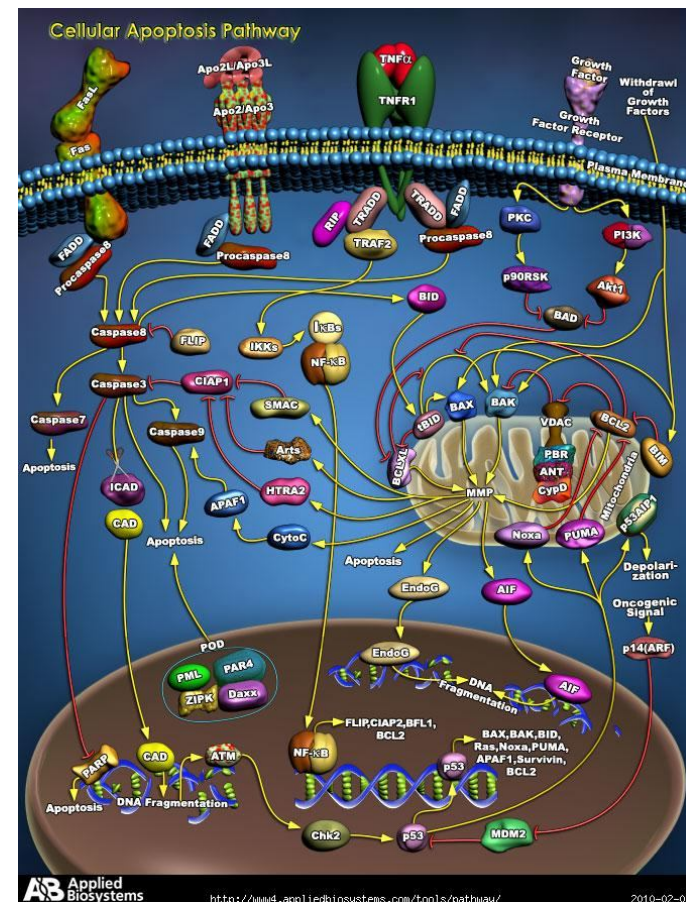
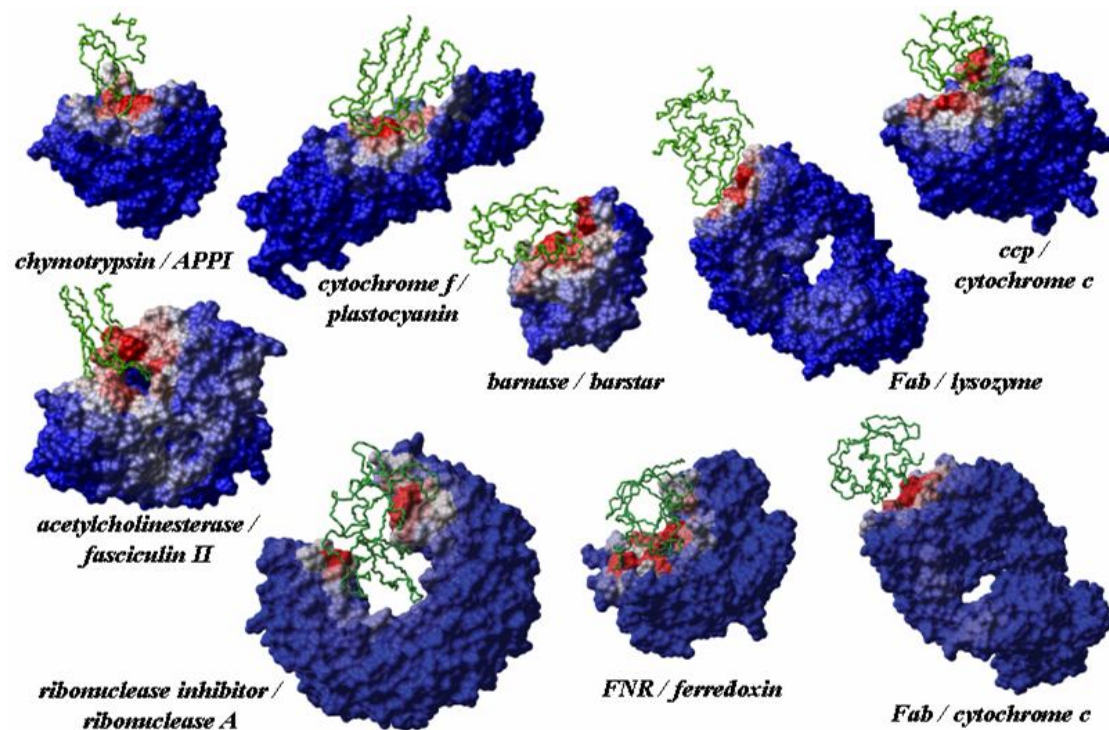


- In the last years the improvement of MD simulation was mainly through software-optimization
- Today the millisecond scale ($10^{-3}s$) is reached using specific hardware (512-processor)
- Atoms and molecules are allowed to interact for a period of time by approximations of known physics



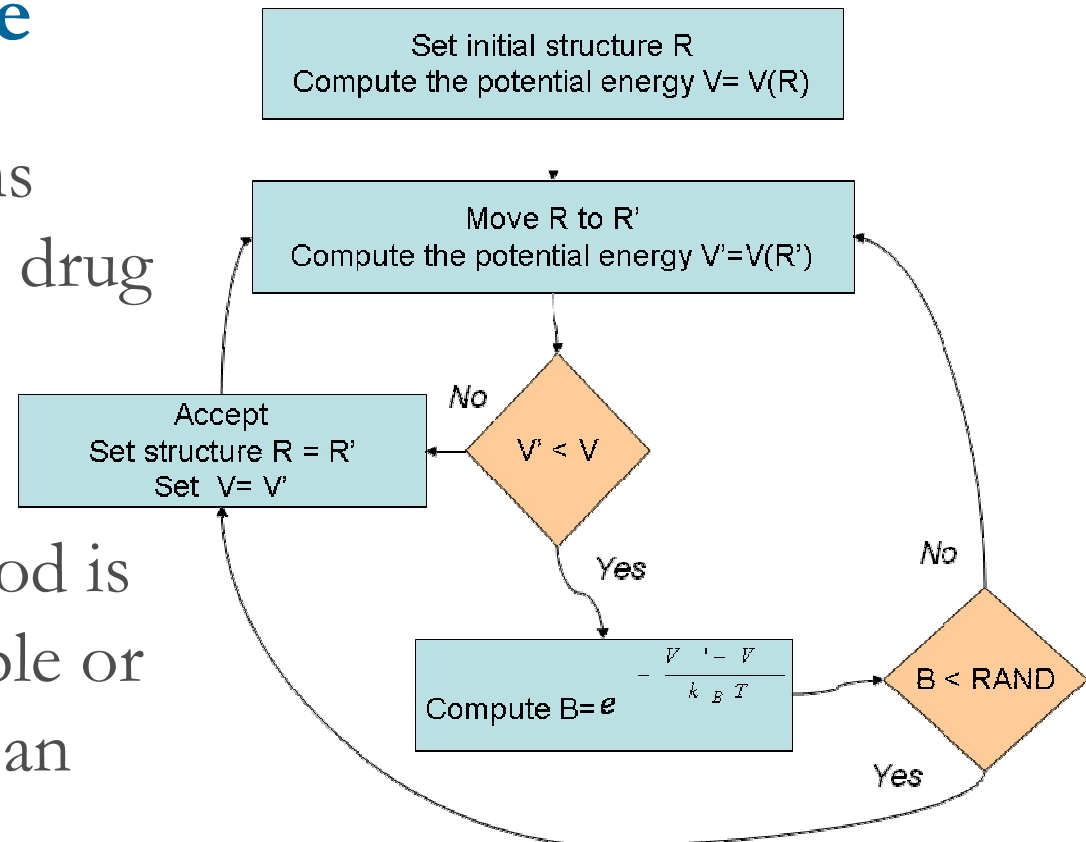


Protein Interaction



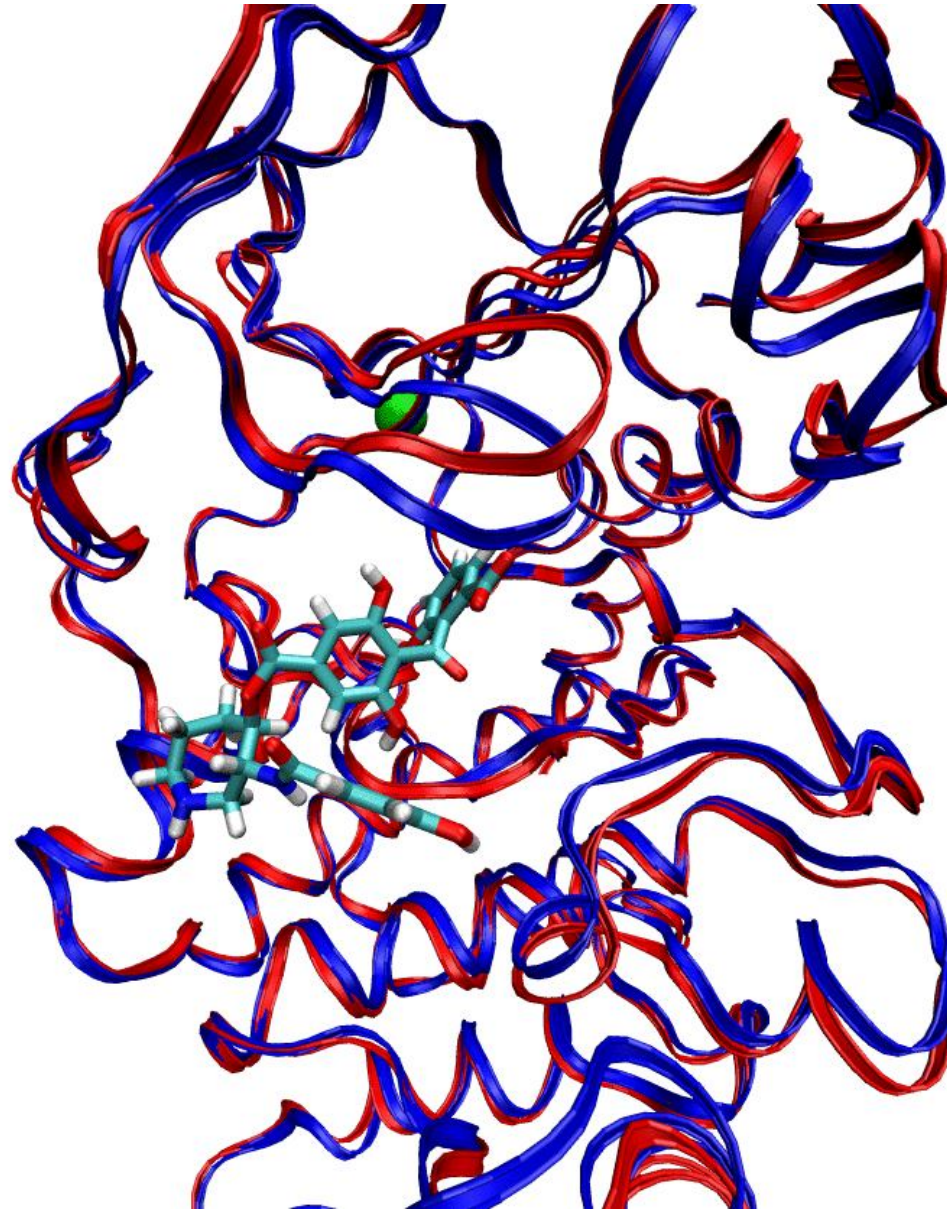
Simulate and Explore

- Explore how proteins interact. Explore how a drug binds its target.
- The Monte Carlo method is used when it is unfeasible or impossible to compute an exact result with a deterministic algorithm.

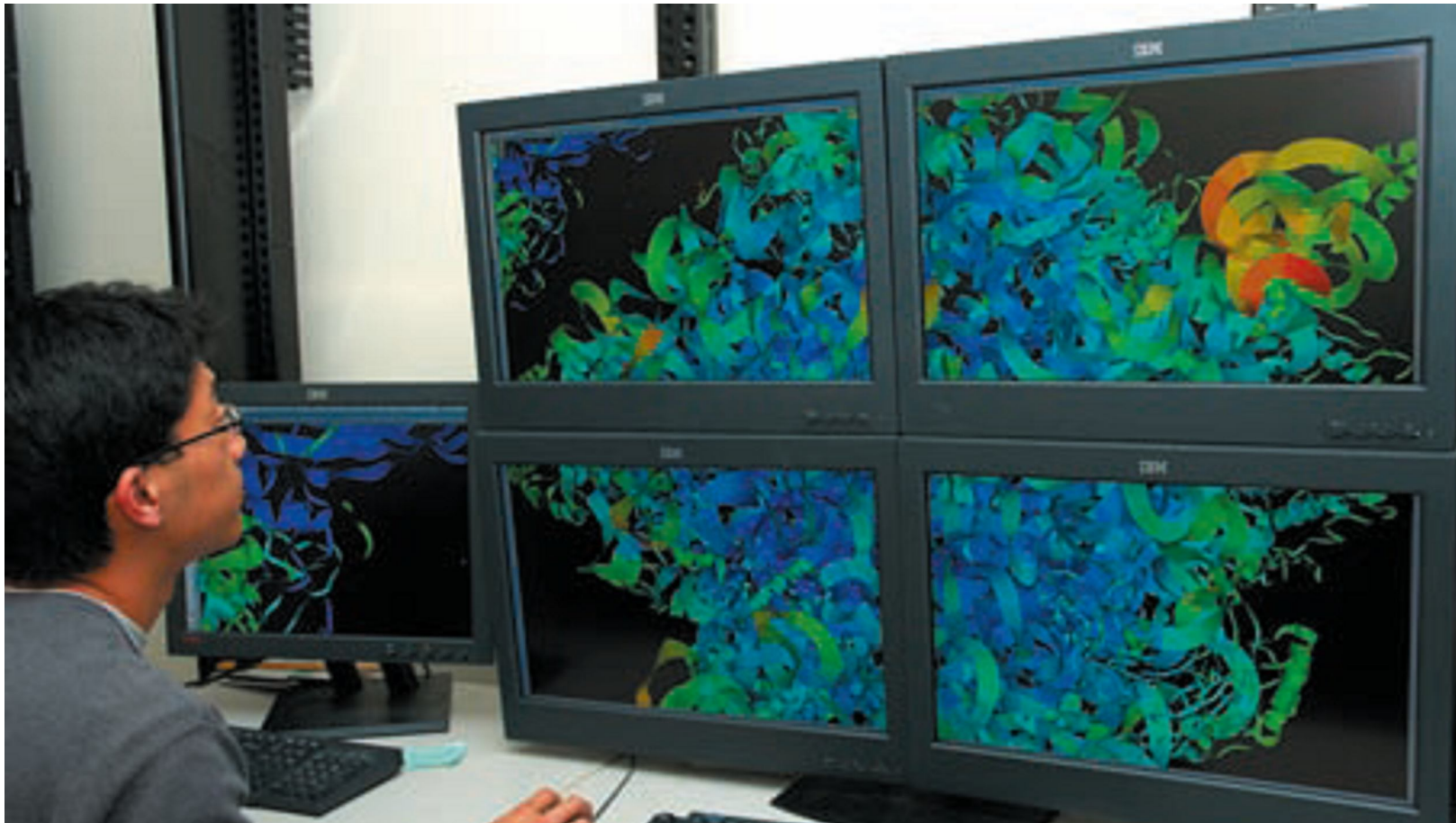


Drug Discovery

- Red: Cristal Protein
- Blue: Modeled Protein (Coarse Grained)
- Drug explores protein binding site using a Monte Carlo algorithm



BioSupercomputing

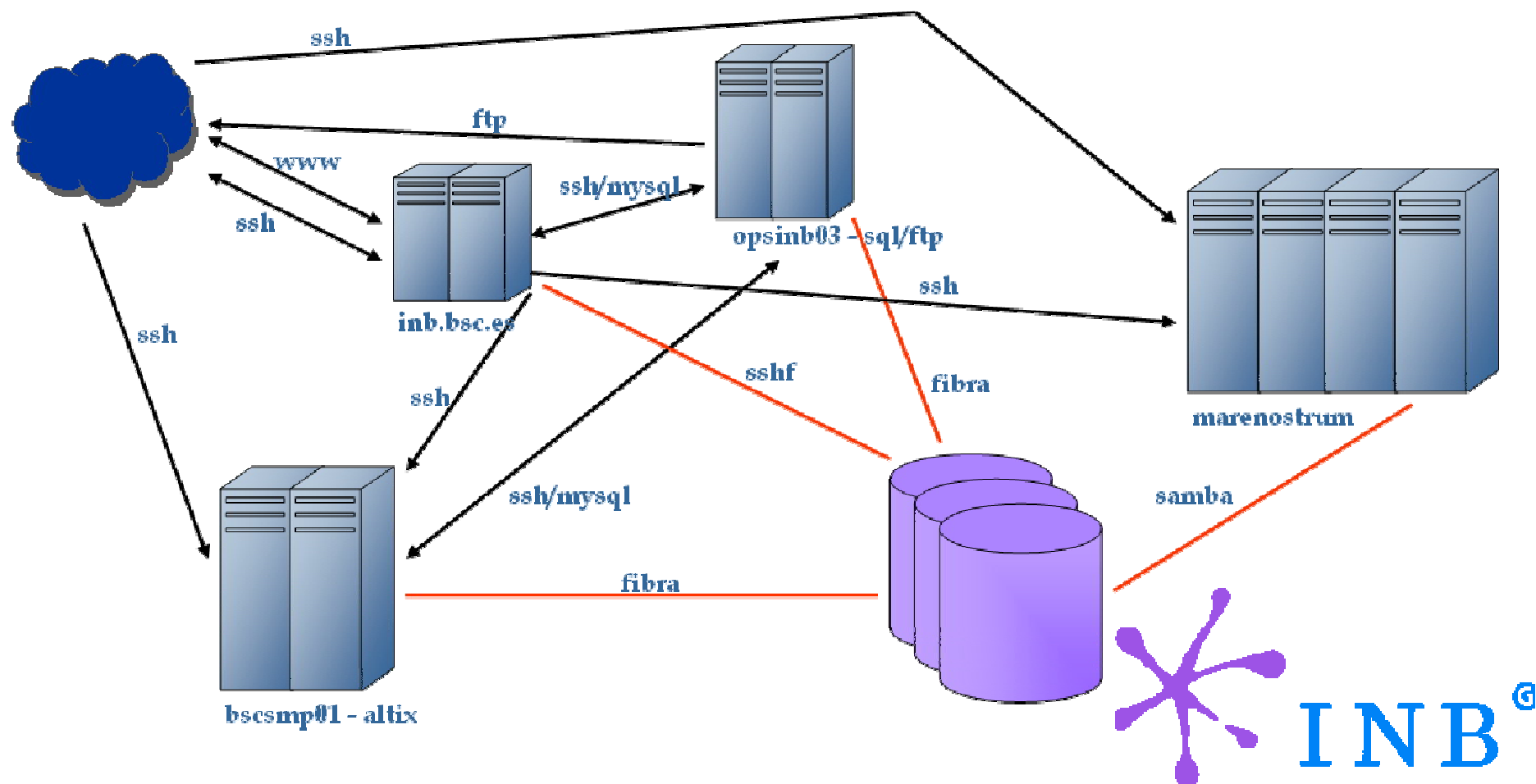


Next Generation of Supercomputers

- We are currently porting our technology to new architecture:
 - **CELL processors:** “Mare Incognito” machine, which aims at providing 10PF in the 2012 time frame . It will be a Supercomputer based on the Cell processor
 - **GPUs:** One of the most time consuming calculations in a typical molecular dynamics simulation is the evaluation of forces between atoms that do not share bonds. The high degree of parallelism and floating point arithmetic capability of GPUs can attain performance levels 20 times that of a single CPU core.



BioSupercomputing open to “BioCommunity”



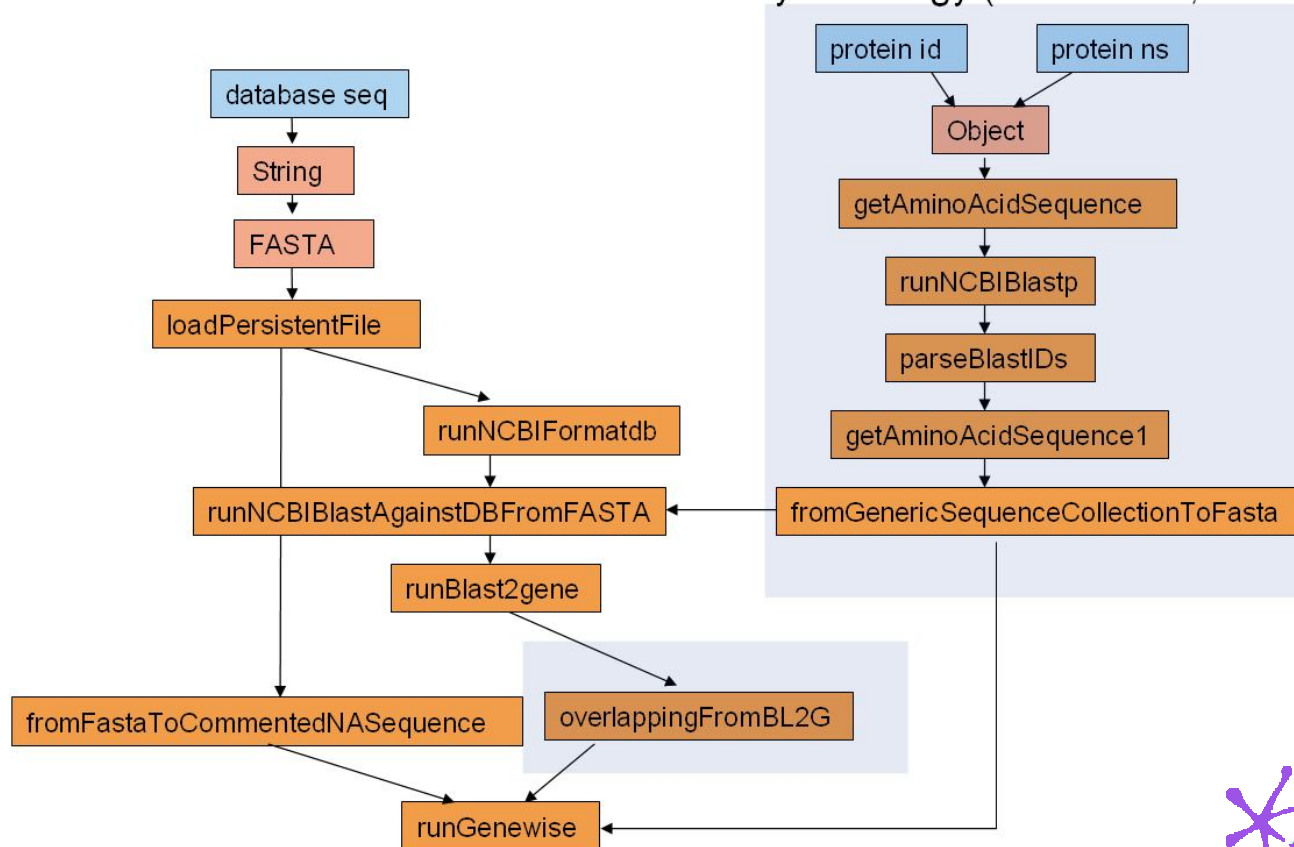
Pre-defined Workflows

ISC 2008

A BioMoby-based workflow for gene detection using sequence homology. Royo, R., López, J., Torrents, D., Gelpi, J.L.

Workflow

Gene detection by homology (D. Torrents, BSC.)



User-friendly for non-expert

The screenshot displays the Moby Miner Applet running in a Mozilla Firefox browser window. The address bar shows the URL http://inb.bsc.es/applications/java/mobyminer/moby_miner.html. The browser's menu bar includes Archivo, Editar, Ver, Historial, Marcadores, Herramientas, and Ayuda. The toolbar contains various navigation and utility icons.

The applet interface is divided into several sections:

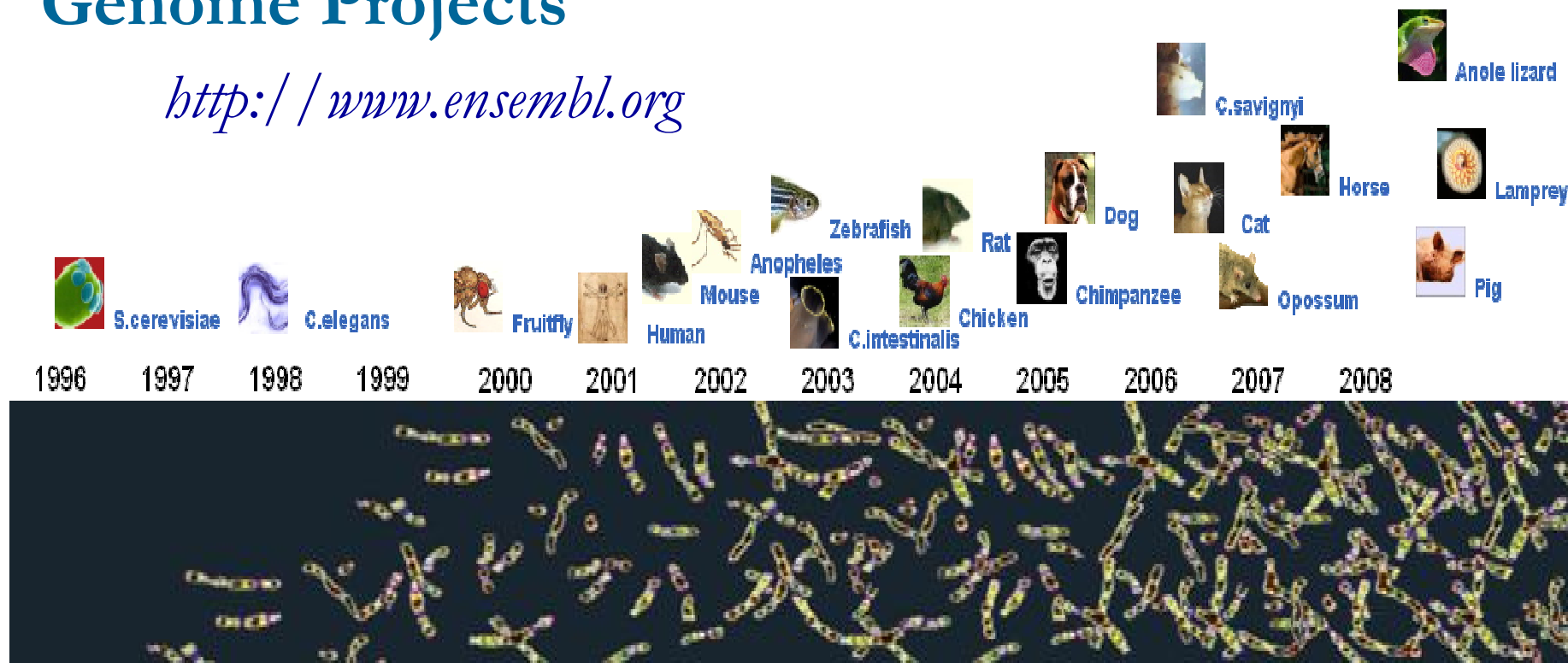
- Search Bar:** A text input field containing '1pio' and a 'search' button.
- Database List:** A scrollable list of databases including:
 - SCOP: Protein families database of alignments and HMMs.
 - Namespace for SCOP database
 - The Universal Protein Knowledgebase.
 - ASLVN9
 - BLAC_STAAU
- Protein Network:** A central node labeled '1pio' is connected to several other nodes:
 - SCOP (2)
 - FSSP (1)
 - PDB (1)
 - DSSP (1)
 - Pfam (1)
 - HSSP (1)
 - UniProt (2)
 - BLAC_STAAU
- Actions:** Two buttons are visible: 'Get the sequence' and 'Get PDB entry'.
- MobyMessage:** A panel on the right side of the applet.
- 3D Visualization:** A Jmol viewer showing a 3D ribbon diagram of a protein structure. A specific residue is highlighted and labeled: '[ASN]203:B.ND2 #3369'.
- Status Bar:** At the bottom left, it says 'Terminado'.

Data Analysis & Data Management



Genome Projects

<http://www.ensembl.org>



<http://www.icgc.org>
<http://www.1000genomes.org>

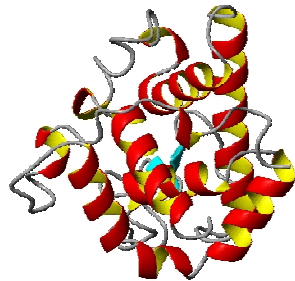


1000

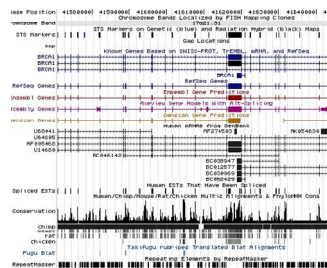


International
Cancer Genome
Consortium

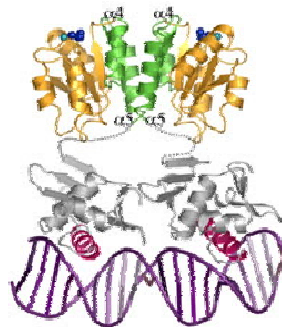
The Genome is the key for Target Discovery



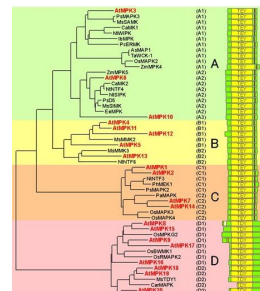
```
DEIGDAAKKLGDASYAFAKEVDWNNGIFLQ
APGKLQPLEALKAI DKMIVMGAAADPKLLK
AAAEAHHKAIGSISGPNVTSRADWDNVNA
ALGRVIASVPENMVMDVYDSVSKITDPKVP
AYMKSLVNGADAEKAYEGFLAFKDVVKKSQ
VTSAA
```



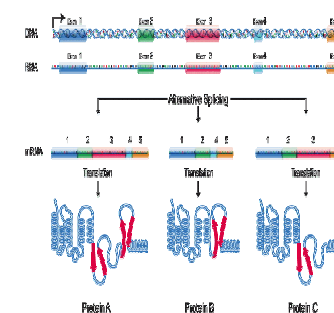
Origin



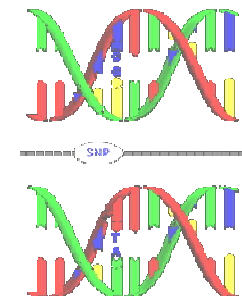
Regulation



Phylogeny



Isoform



Mutations

Data Overload

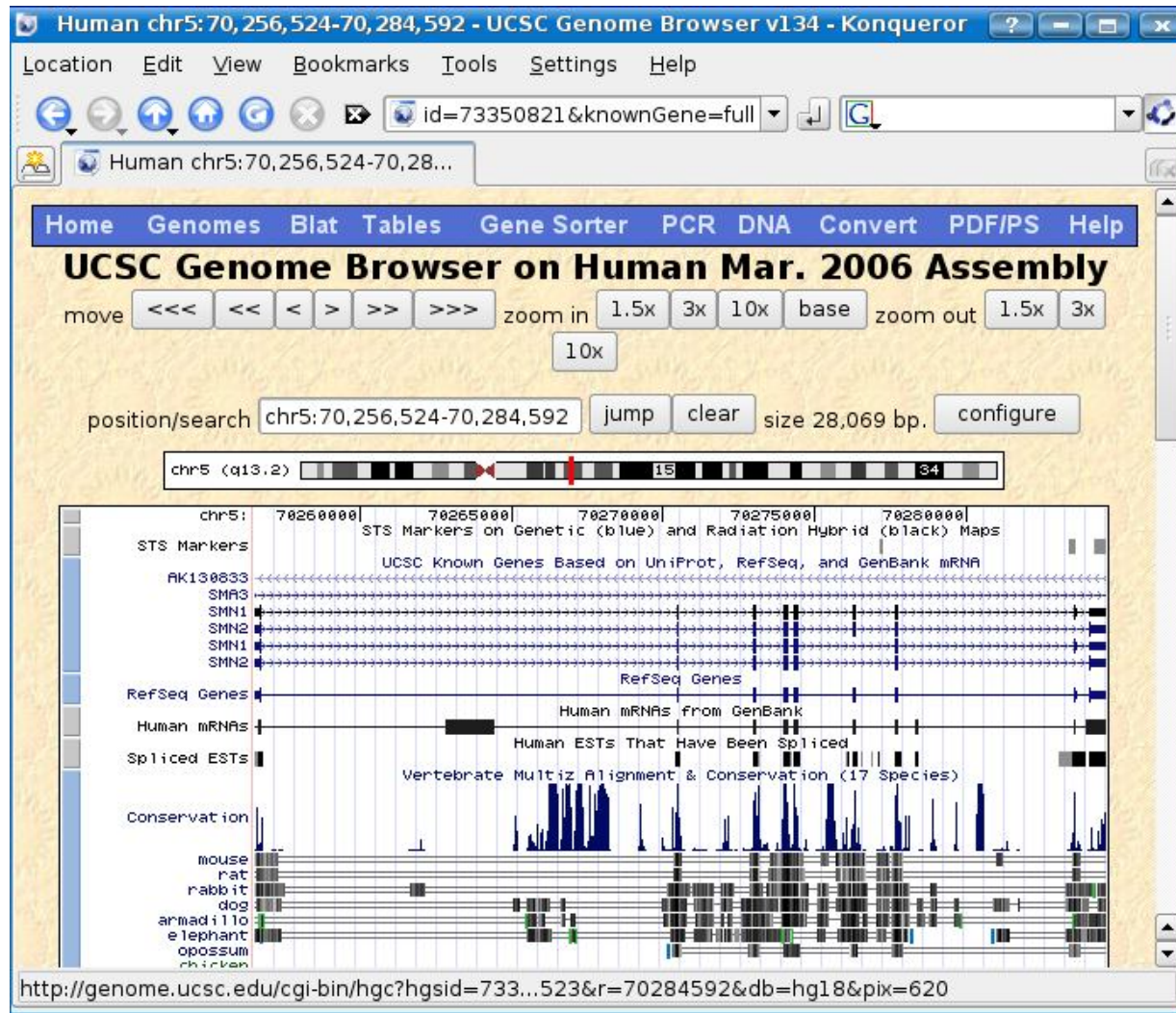
- Next Generation Sequencing platforms can sequence a complete genome fast and cheap. It is estimated that soon this service will cost \$100.
- Approximately 115,200 Tiff formatted files are produced per run, each at about 8 megabytes (MB) in size. This is approximately 1 terabyte (TB) of data, which must be moved from the capture workstation to the analysis resource. A mere 10–20 sequencing runs could overwhelm any storage and archiving system available to individual investigators.
- A 1 GB network is essential within this environment, with 10 GB networks becoming more prevalent.



Our Challenge is ...the Tomato



Genome Annotation



DNA sequence and DNA structure

Genomic Sequence

SEQUENCE1 :

TGCACGTAGCTAAAAA ██████████ GACGG ██████████ GCGTCGAGCAGCGAG

SEQUENCE2 :

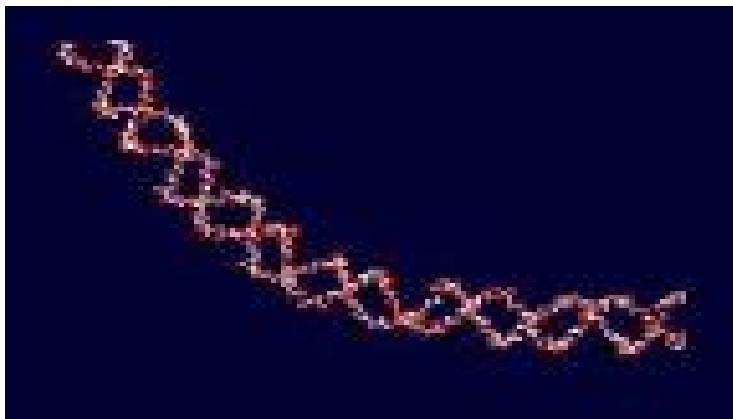
CGACTAGCACTACGTACATCGACATGCGACATCATCGACCTCGACAGTCGACGTCAACGACGAGTCAGCACGTGTAGTCGACAGTGAGCGGCAG

BLAST (sequence alignment): ID 49%

Sequence1	TGCACGTAGCTAAAAAGAGGAGGGAGGAGAG-GA----GAGGGATGACGGAAGAGGAGGGAGAGAGAG-AGGCGTC--GAGCAGCGAG
	. .
Sequence2	CGACTAGCACTACG-TACATCGACATGCGACATCATCGACCTCGACAGTCGACGTCAACGACGAGTCAGCACGTGTAGTCGACAGTGAGCGGC-AG

DNA Molecule

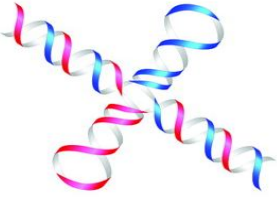

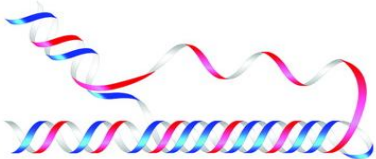

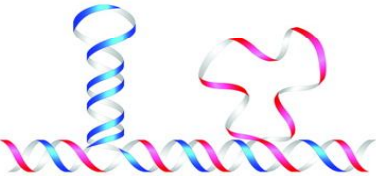

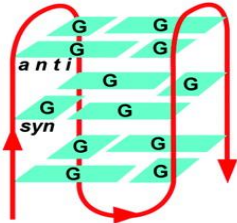

Sequence1



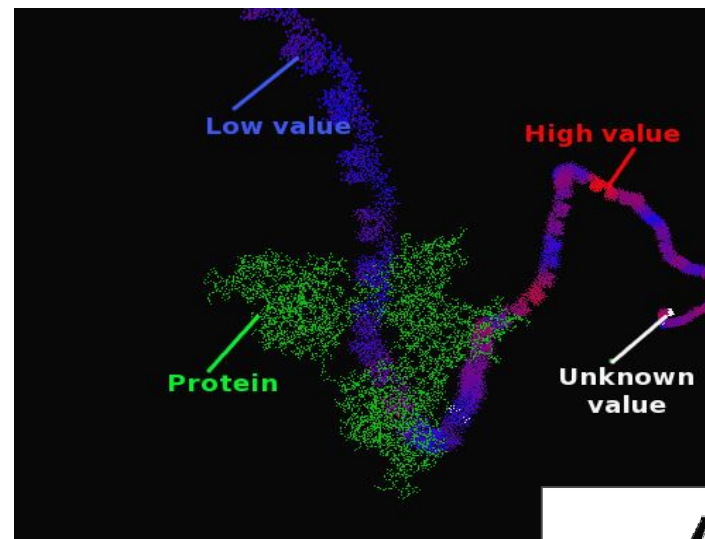
Sequence2



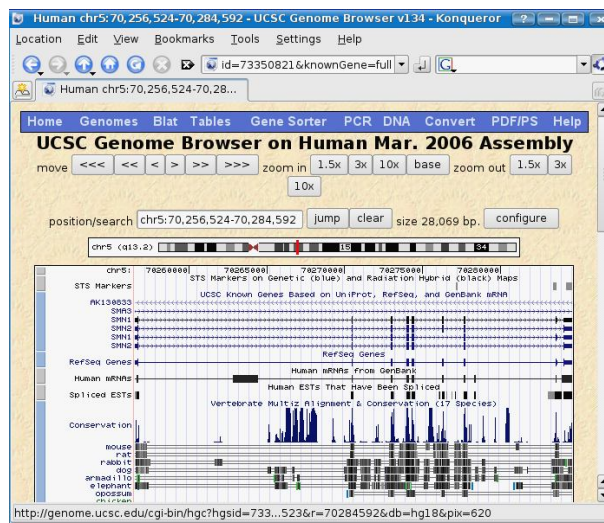
Unusual Conformations of DNA

Name	Conformation	General Seq. Requirements	Sequence
Cruciform		Inverted Repeats	 TCGGTACCGA AGCCATGGCT
Triplex		(R•Y) _n Mirror Repeats	 AAGAGG GGAGAA TTCTCC CCTCTT
Slipped (Hairpin) Structure		Direct Repeats	 TCGGTTTCGGT AGCCAAGCCA
Tetraplex		Oligo (G) _n Tracts	AG ₃ (T ₂ AG ₃) ₃ single strand
Left-handed Z - DNA		(YR•YR) _n	B - Z Junctions CGCGTGCGTGTG GCGCACGCACAC

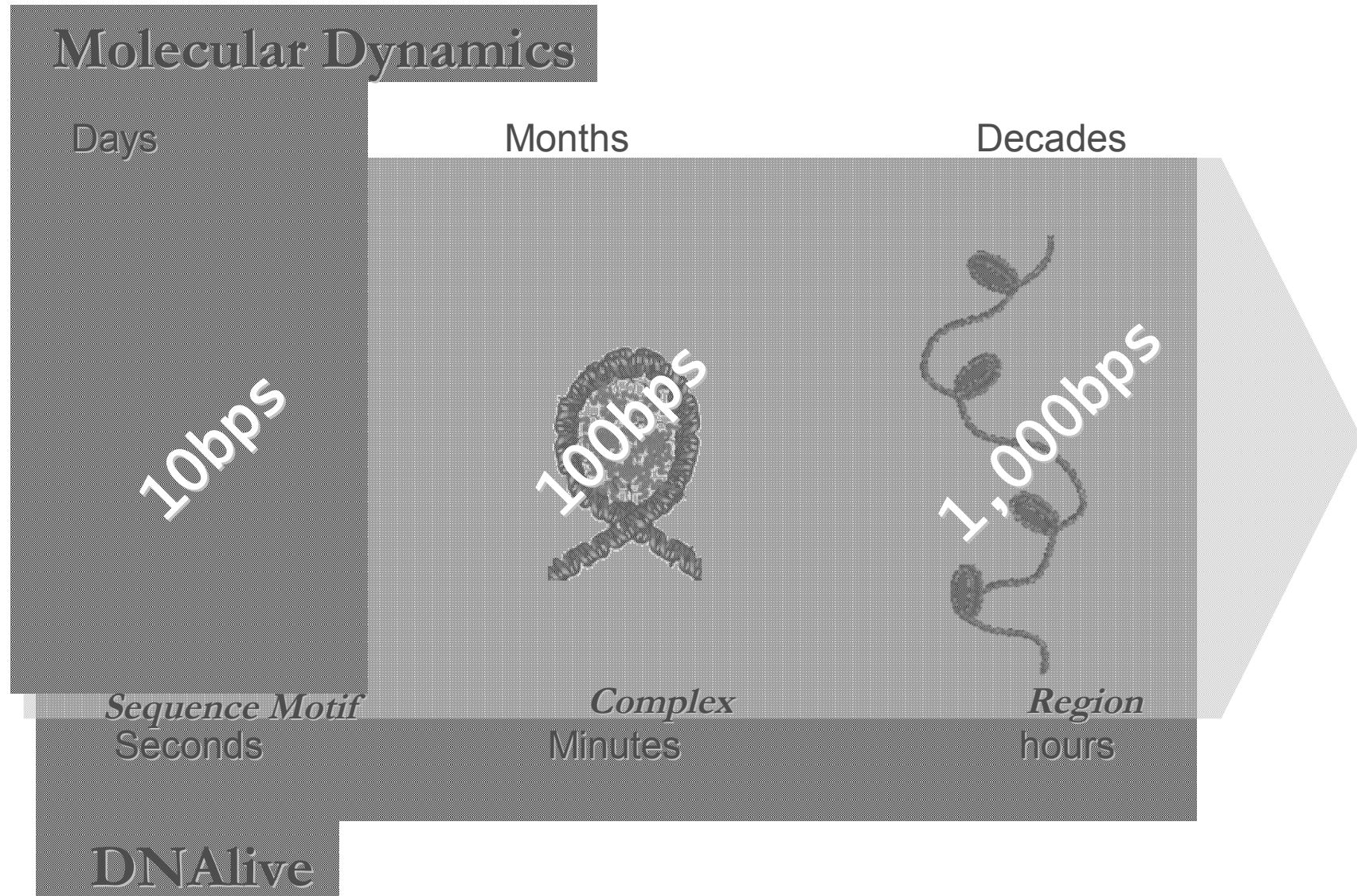
Bacolla et al.
 (2004)
J. Biol. Chem



GENOME IN 3D !



Large Scale DNA Dynamics

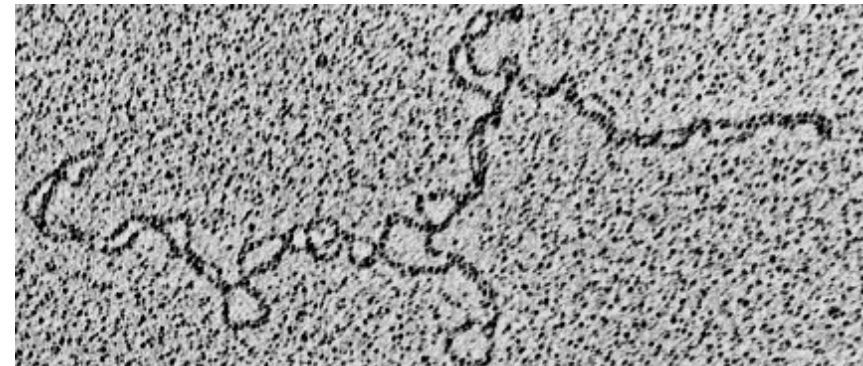


Representing the genomic DNA

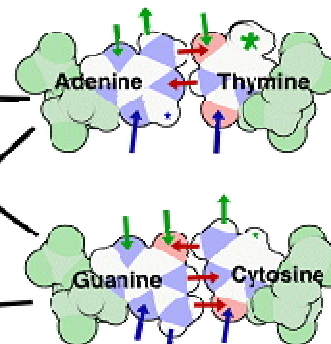


Supercomputing
Resources

Molecular
Dynamics

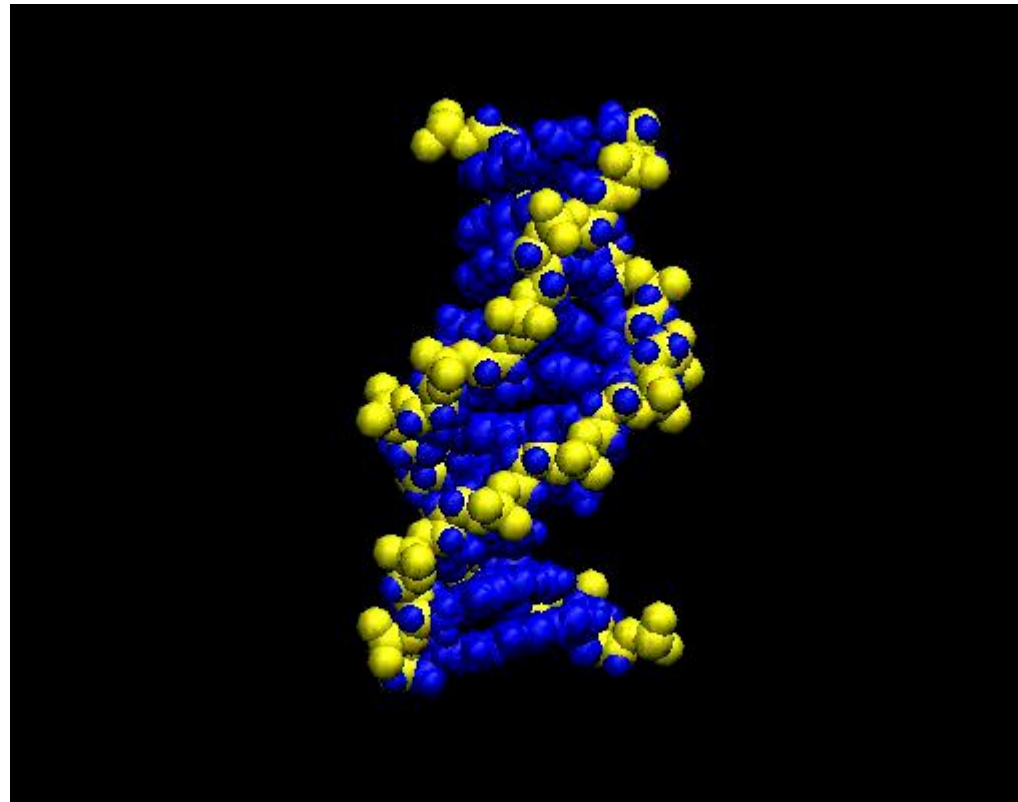


Observed DNA

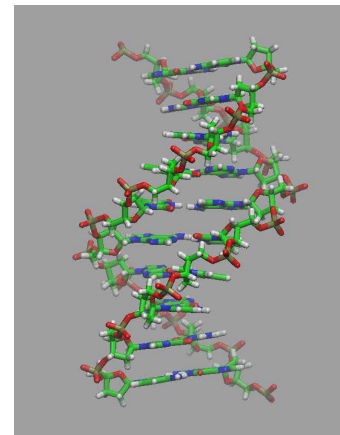
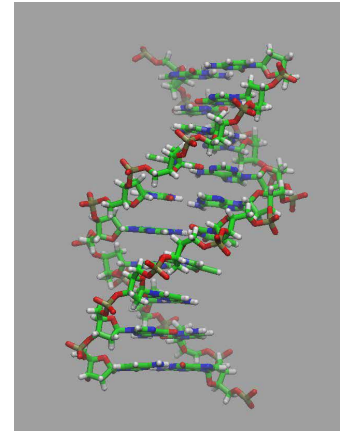
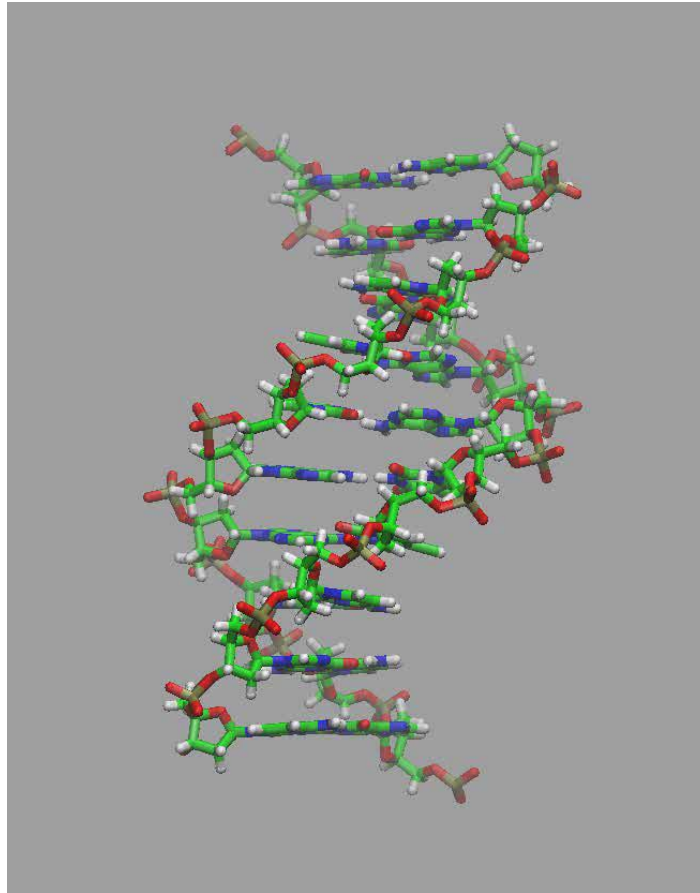


Mesoscopic

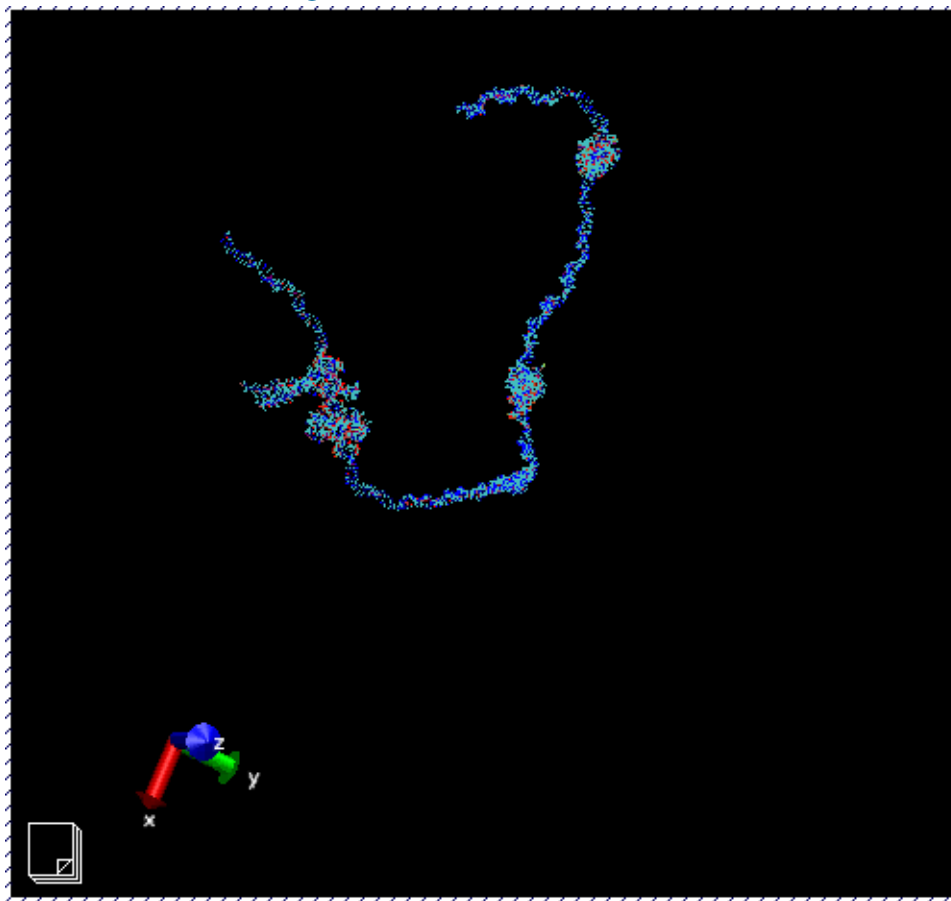
ParmBSC: a Force Field for long DNA Simulations



Study of Principal Component Analysis and Essential Dynamics



Genic DNA Dynamics



$$E = \sum_{i=1}^M \sum_{j=1}^6 K_{ij} \left(\xi_{ij} - \xi_{ij}^0 \right)^2$$

MonteCarlo Simulation (n. steps = 100,000 x number of flexible dinucleotides)
Adjusted to accept 40% of the steps
Energy: Perez et al. (2007) JACS

Summary

- Computer simulation (Drug Discovery)
 - **Molecular Dynamics**
 - **Protein-Protein Interaction**
 - **Protein-Ligand Docking**
- BioSupercomputing
 - **Computational Biology under GPU & CELL**
 - **User-friendly computing access Web-Services**
- Data analysis and Data Management (Target Discovery)
 - **Next Generation Sequencing**
 - **Genomics & DNA structure**

Contact:

Ramon Goñi
Senior Researcher
BSC, Life Sciences
ramon.goni@bsc.es

