# The biological problem

- Extremely important
- Complex
- Diffuse and stochastic nature
- Bad defined problems
- Cover a wide scale (time and space)
- Typically defined with poor data
- Exponential growth (↑↑ Moore law)
- Cost of obtaining data ↓↓
- Cost of processing data ↑↑

# The dual nature of computers

- Ordenador: machine to manage data.

- Computador: machine to do maths.



The biological problem is diverse in nature.
Input of biologists should be important
when selecting architectures

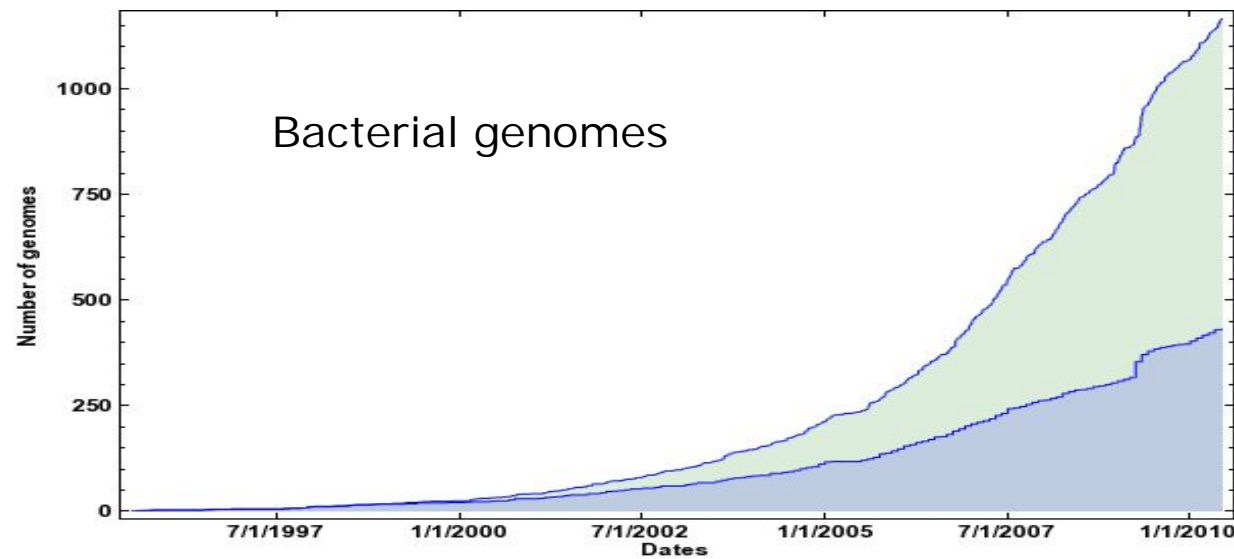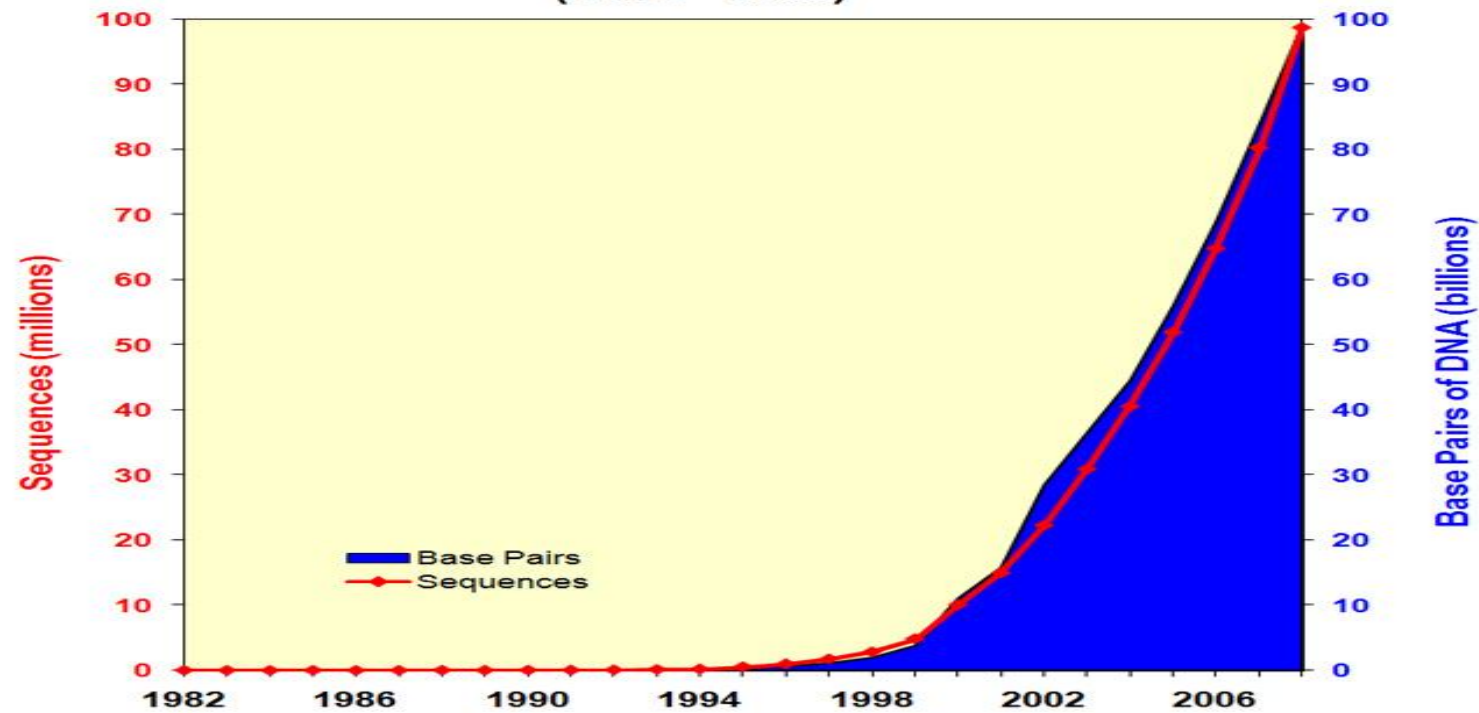# The dual nature of computers

- Ordenador: machine
  to manage data.
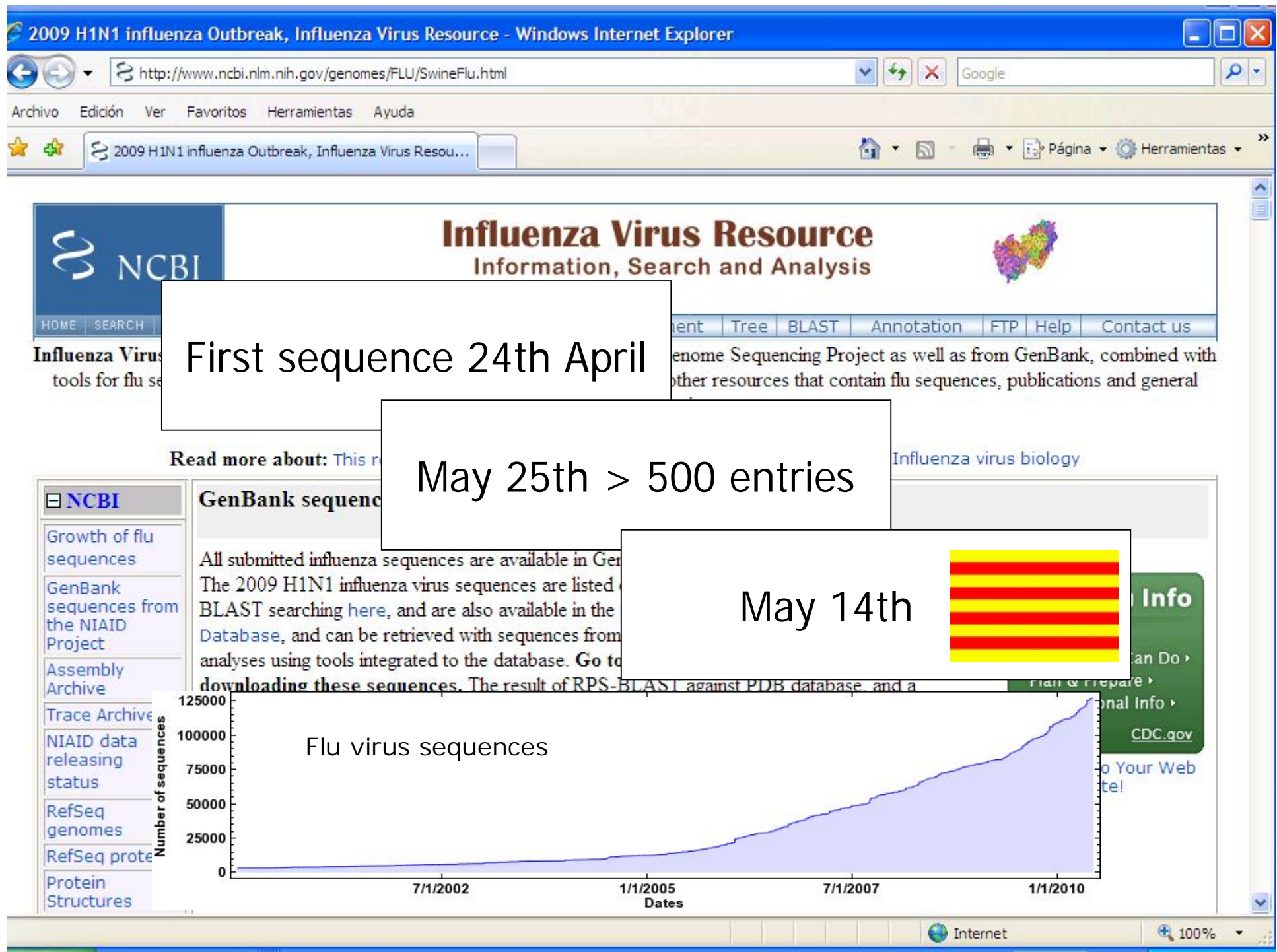
- Computador:
  machine to do maths.

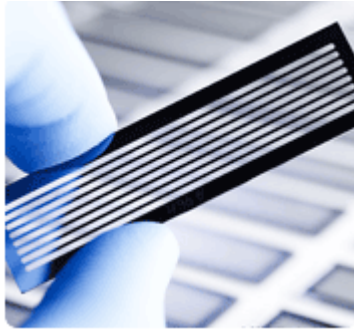In biology capability does not always translate into FLOPs

# The computer as data manager

- Data growths exponentially
- Data management more complex
- Data integration a crucial issue
- Processing of data very costly

- Computers are now the limiting step in many biological problems
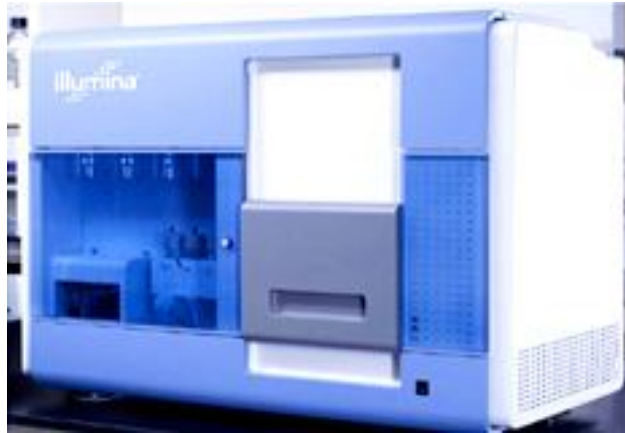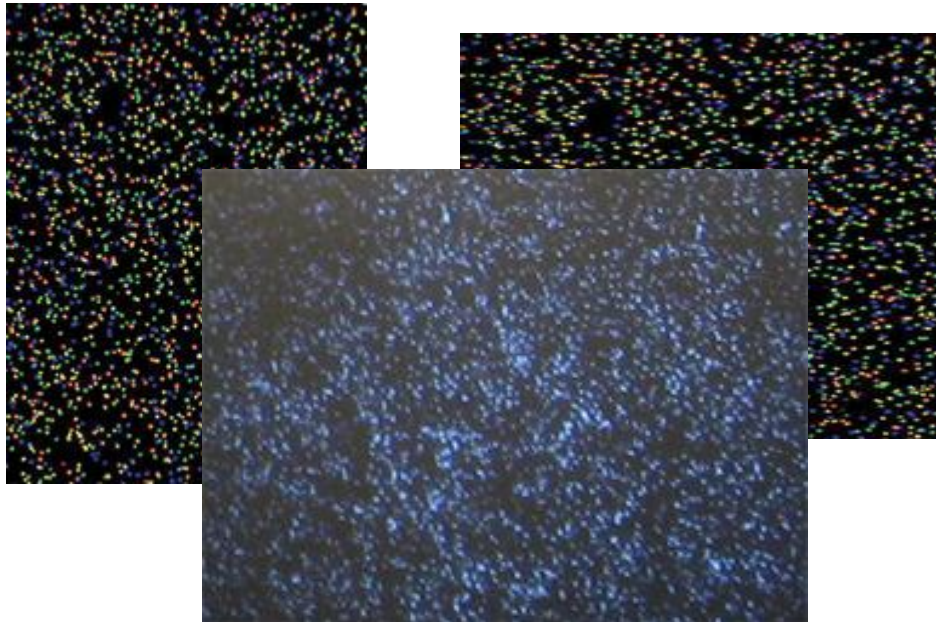- Time is a major issue in Bio-research

Growth of GenBank (1982 - 2008)

Bacterial genomes

http://www.ncbi.nlm.nih.gov/genomes/FLU/SwineFlu.html

Google

Archivo   Edición   Ver   Favoritos   Herramientas   Ayuda

2009 H1N1 influenza Outbreak, Influenza Virus Resou...

Página   Herramientas

**Influenza Virus Resource**
Information, Search and Analysis

NCBI

HOME   SEARCH   ...ment   Tree   BLAST   Annotation   FTP   Help   Contact us

**Influenza Virus** ...enome Sequencing Project as well as from GenBank, combined with
tools for flu se... ...other resources that contain flu sequences, publications and general

## First sequence 24th April

Read more about: This r... ...Influenza virus biology

## May 25th > 500 entries

NCBI

**GenBank sequenc...**

Growth of flu sequences

GenBank sequences from the NIAID Project

Assembly Archive

Trace Archive

NIAID data releasing status

RefSeq genomes

RefSeq prote...

Protein Structures

All submitted influenza sequences are available in Ge...
The 2009 H1N1 influenza virus sequences are listed ...
BLAST searching here, and are also available in the ...
Database, and can be retrieved with sequences from ...
analyses using tools integrated to the database. **Go t**o...
downloading these sequences. The result of RPS-BLAST against PDB database, and a

## May 14th

...Info

...Can Do

...nal Info

CDC.gov

...o Your Web
...te!

### Flu virus sequences

125000
100000
75000
50000
25000
0

Number of sequences

7/1/2002          1/1/2005          7/1/2007          1/1/2010
**Dates**

Internet          100%

# Next generation sequencing

- Every experiment 2 Tb of data
- Every machine 2 experiments a week
- A medium sized center 10 machines

A sequencing center generates around
2 Petabytes data a year!

# Massive (human) genomic projects

# Massive (human) genomic projects



50 cancers

25000 cancer genomes

# Structures, chemical, literature,...



3-Hydroxy-DL-kynurenine

melatonin

**PubMed Literature**

Chemical banks in the order of $10^7$ compounds

SNPs responsible
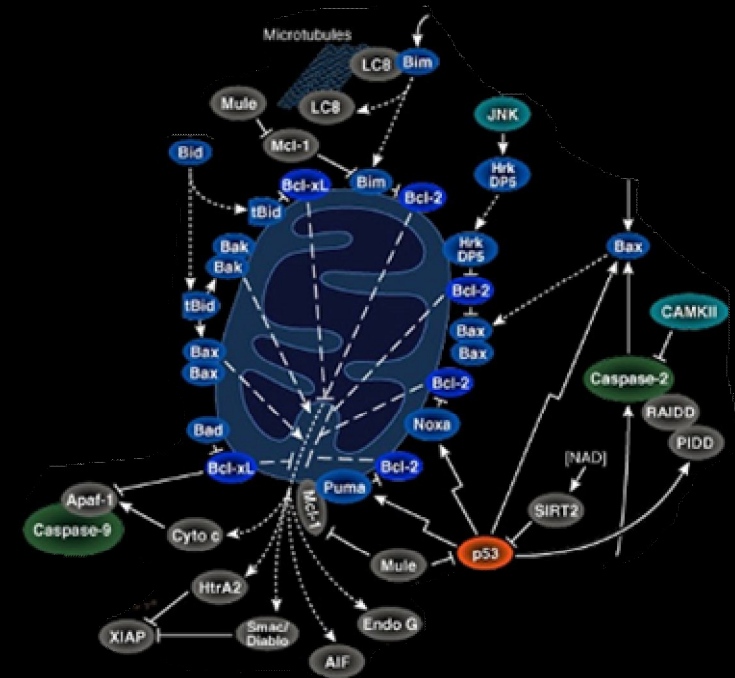90% interspecie variability

Responsible of
monogenic
pathologies

# Modeling complex diseases using genomic data
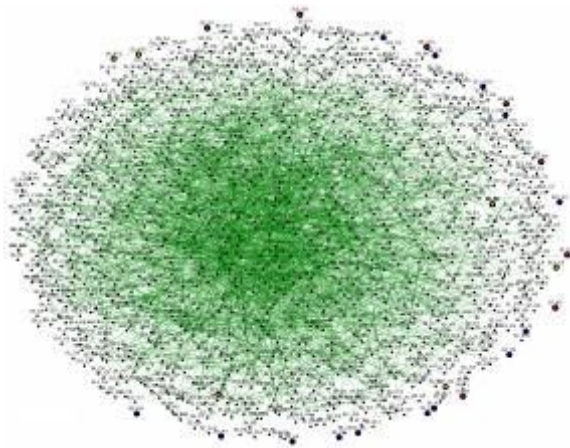
# Ex. the genomic approach to medicine





- Genome healthy vs pathological

- Trace changes

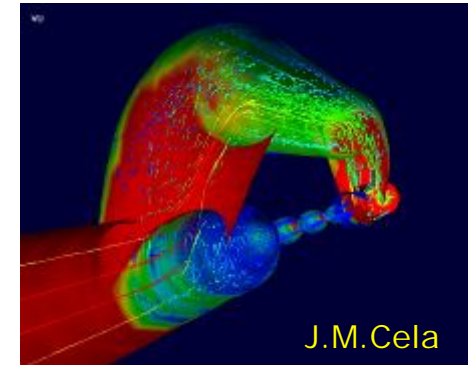- Remove noisse

- Correlate changes with pathology

Every step is computationally demanding
the last simply impossible (for current computers)
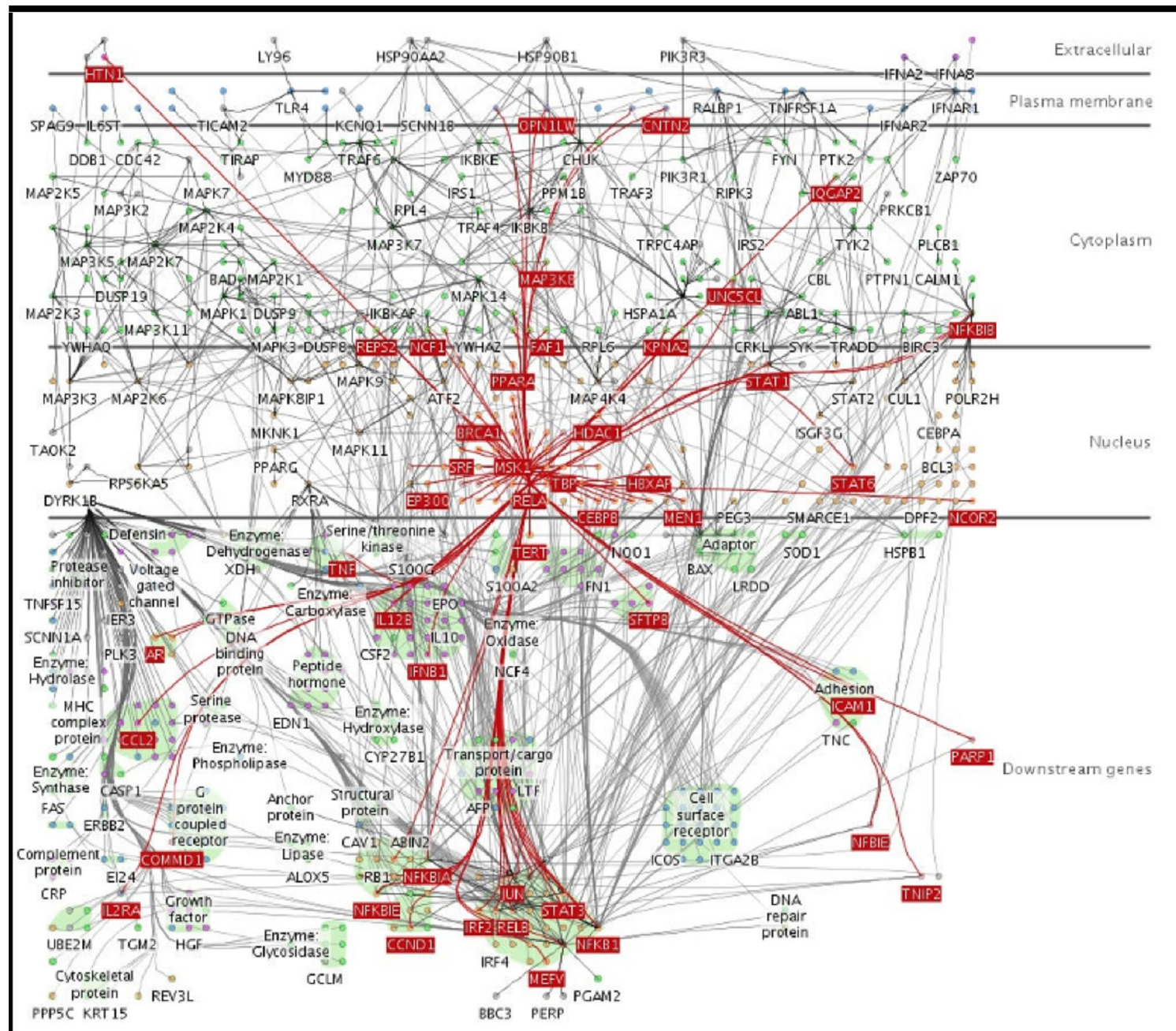
# The computer as a calculator

- Importance of simulation increases with:
  - Increase in data on biological systems
  - Better definition of the problem
- Different types of algorithms
- Must be robust to lack of information
- Often set-up conditions are unclear

# Simulation scenario in Life Sciences

- Ecosystem simulation
- Organ simulation


J.M.Cela

- Gene inter-relations (Bayesian logics)
- Cell simulation (systems biology)
- Molecular simulations
  - Structural prediction
  - Docking
  - Atomistic simulation
  - Cell-scale mesoscopic simulations

# Breast cancer interactome

HEME GROUP

STRUCTURE IS THE MOST POWERFUL WAY TO UNDERSTAND BIOLOGY

INHIBITOR

MEMBRANE ANCHORING DOMAIN

# Moving from abstract networks and real cells



Ribosome

Permease

ATP synthase

?

Ftsz

Adhesin

Predicting structure of complexes
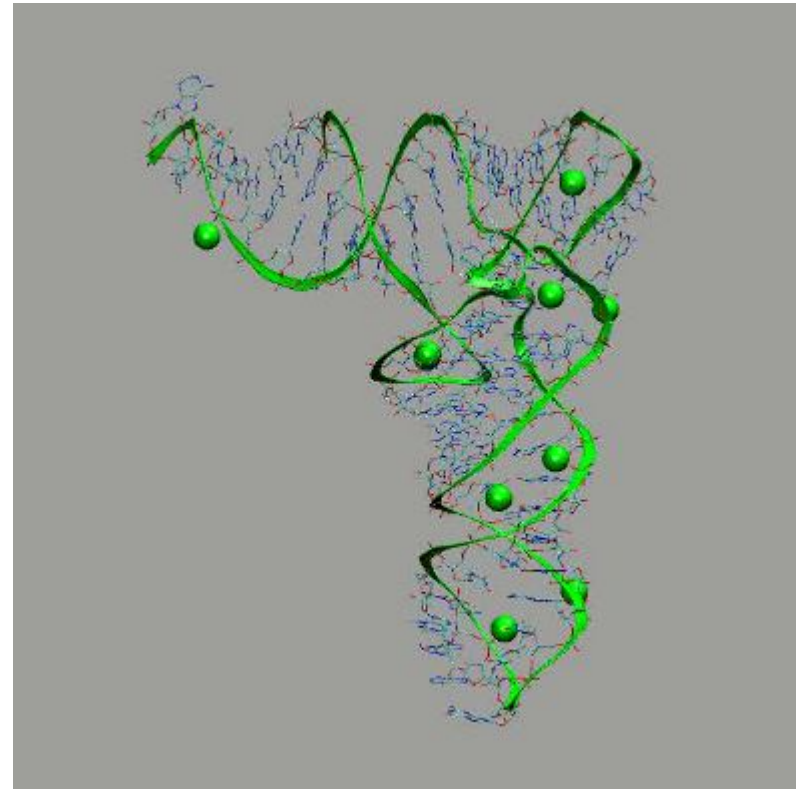
# Structure prediction



$10^8$ sequences vs $10^5$ structures

# Molecular dynamics

$$\vec{f}_i = -\frac{\partial E_i}{\partial \vec{r}_i}$$
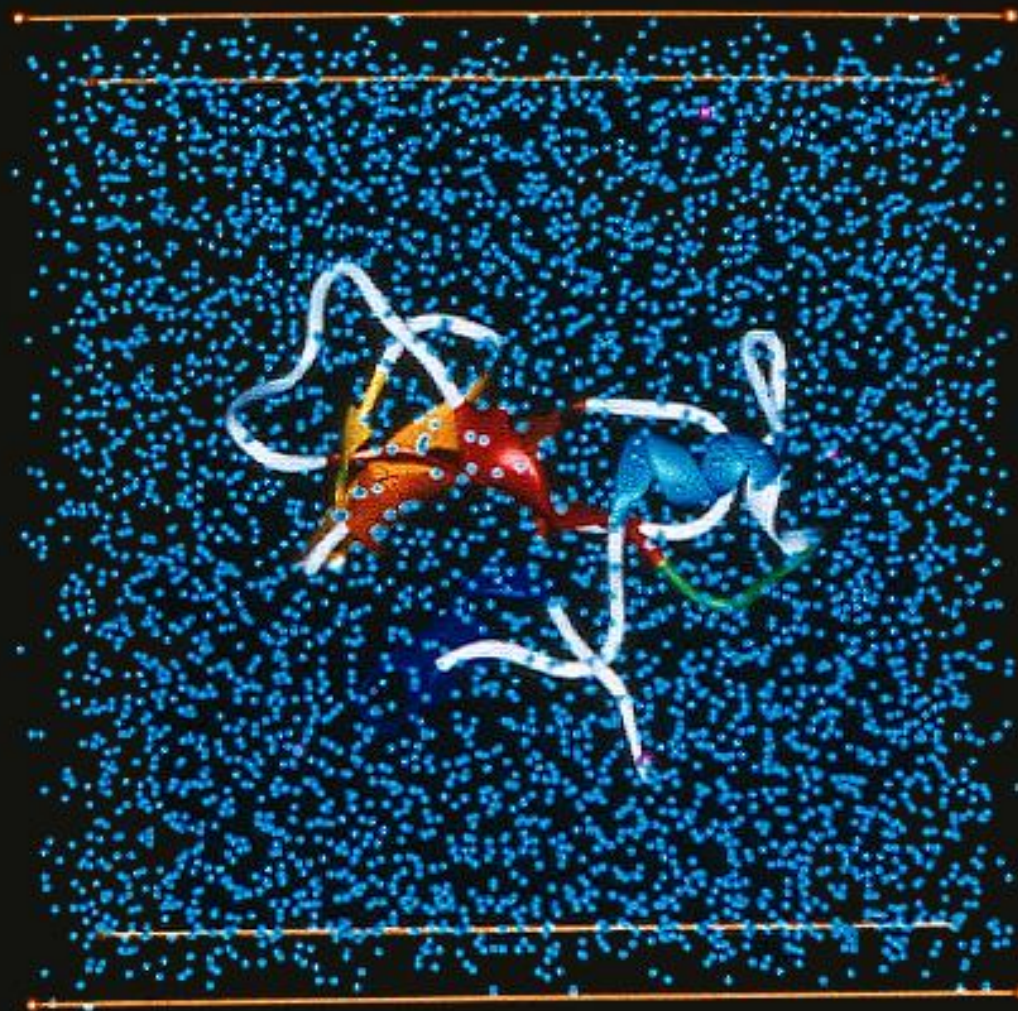
$$\vec{f}_i = m_i \vec{a}_i$$
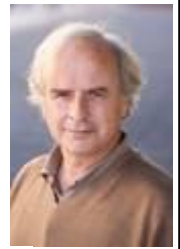
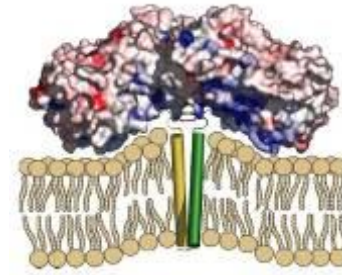$$\vec{v}_i = \int \vec{a}_i \, dt$$
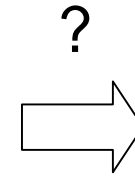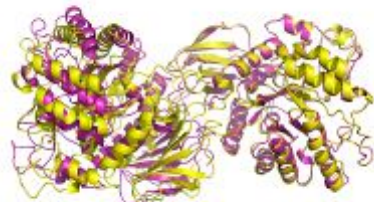
$$\vec{r}_i = \int \vec{v}_i \, dt$$



Integration step 1 fts ($10^{-15}$ seg) $\rightarrow$ 1 mseg = 1 Eur Billion integration steps

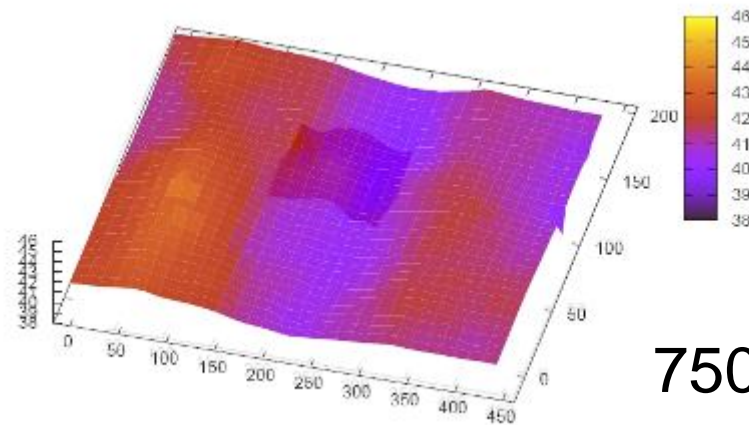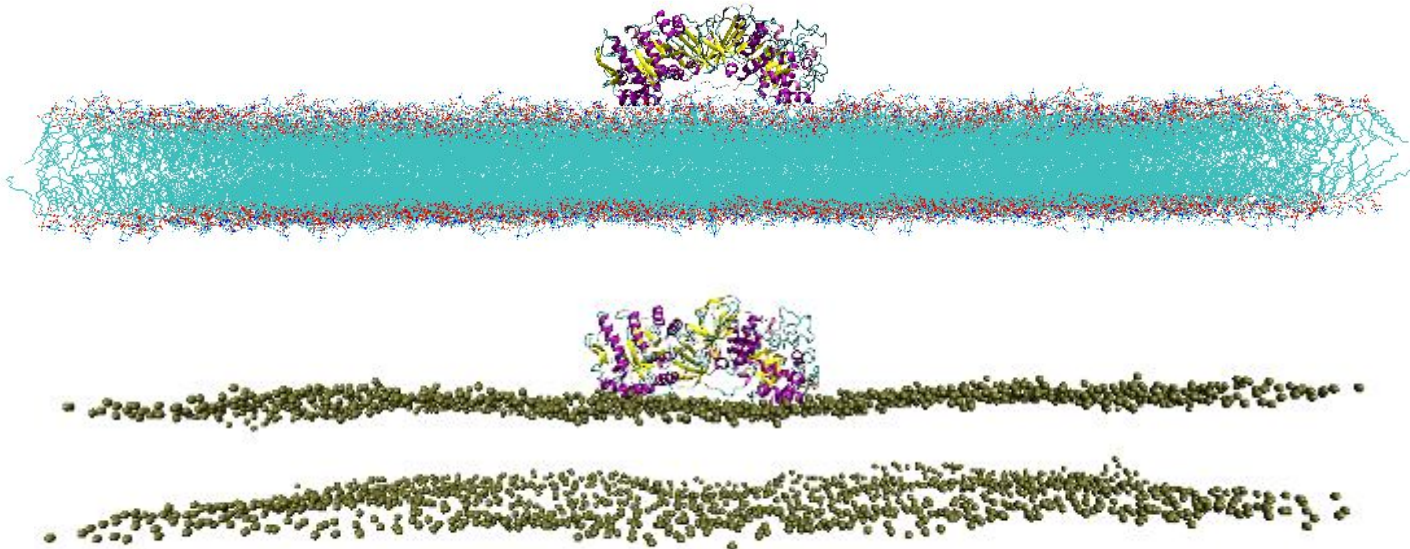Equivalent to follow evolution Nearthental$\rightarrow$ H sapiens with photos every sec
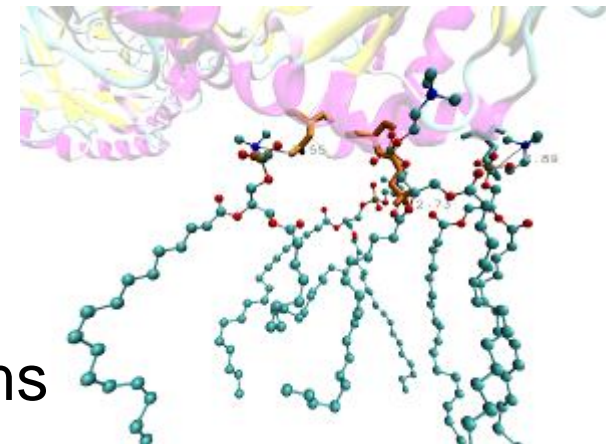
Less than 10000 atoms
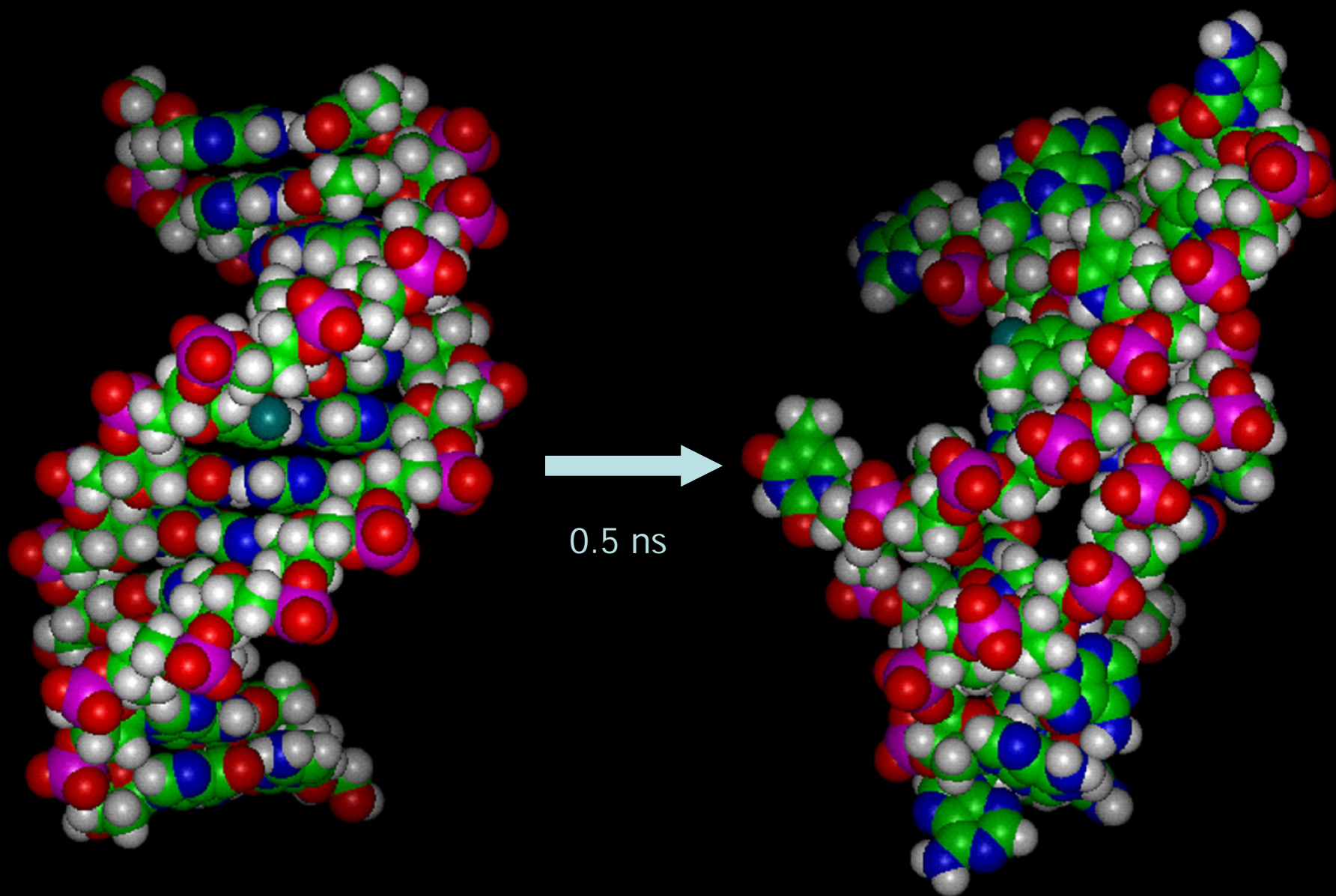
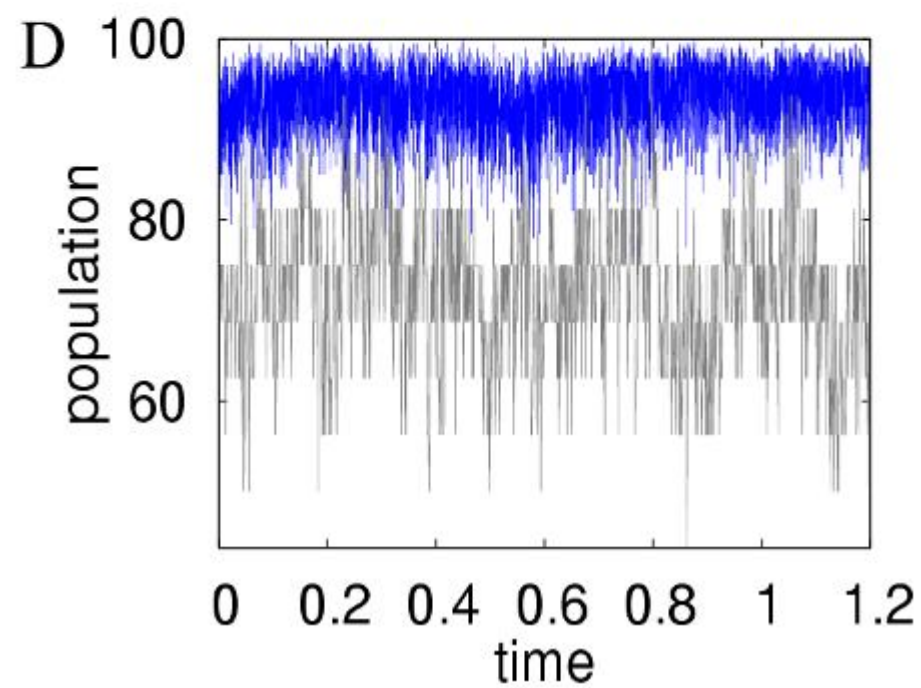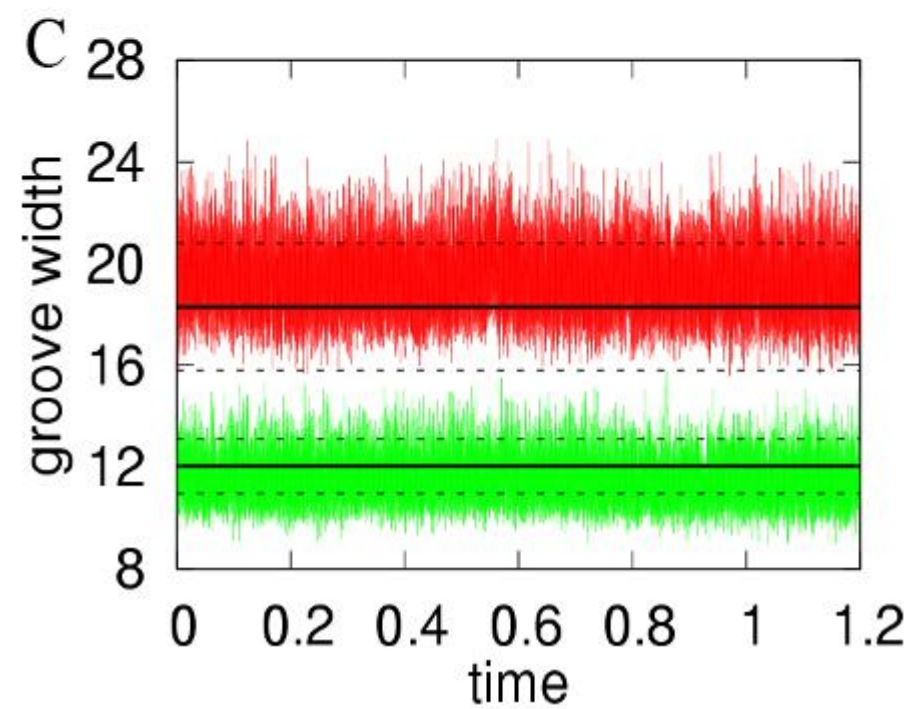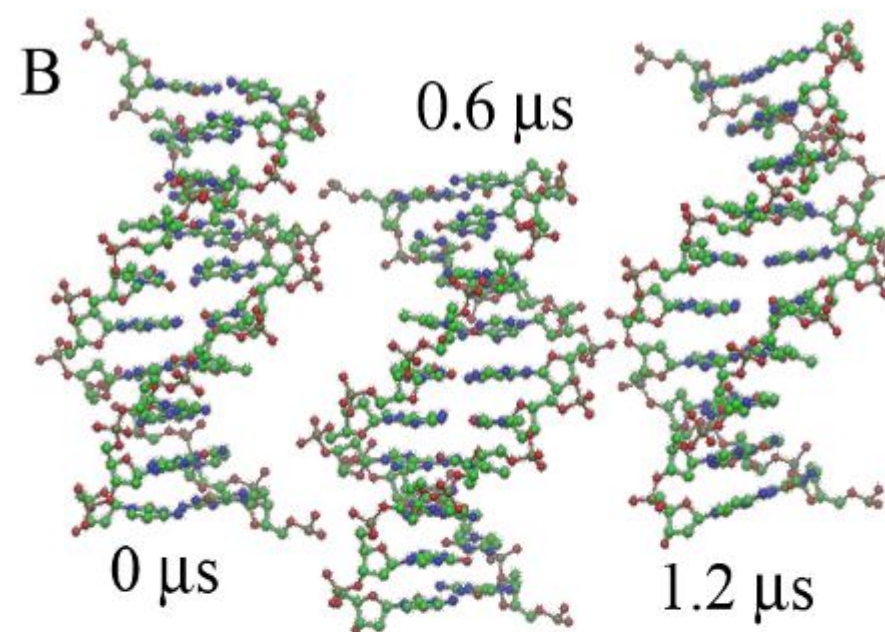Tirado et al., Biochemistry, 36,7313, 1997
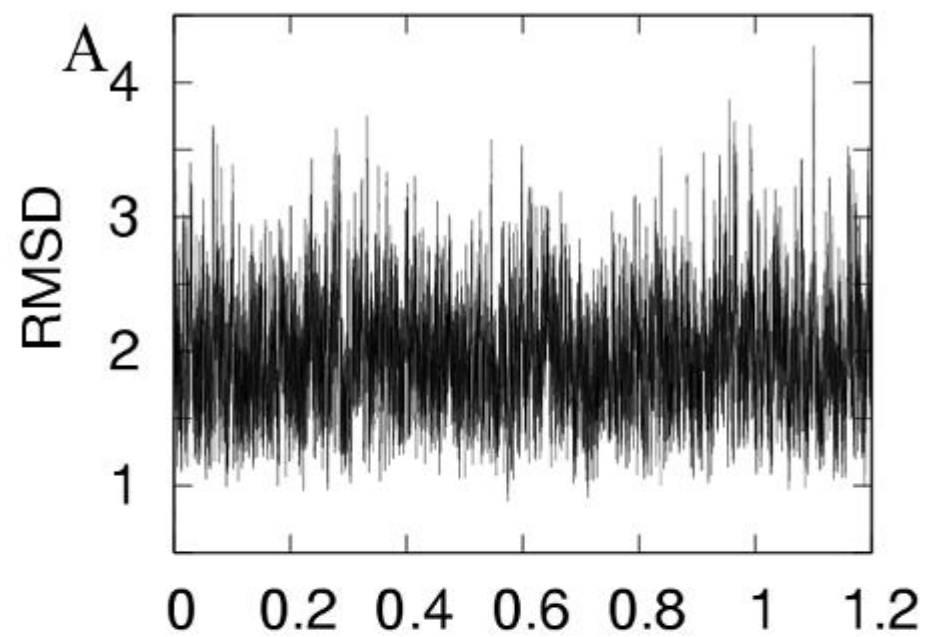
J.Fort, et al  JBC 2007
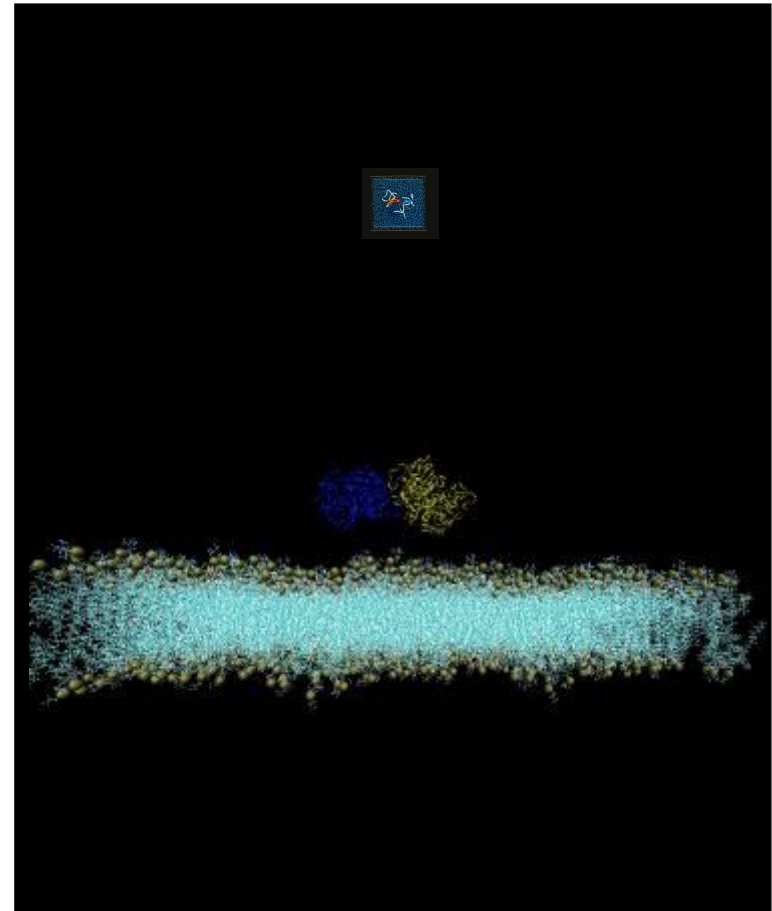
750000 atoms

0.5 ns

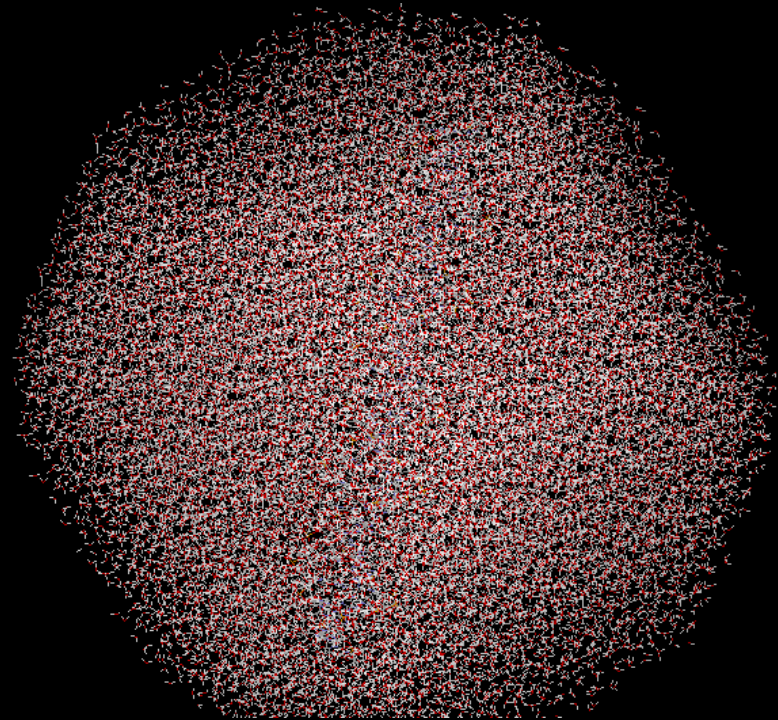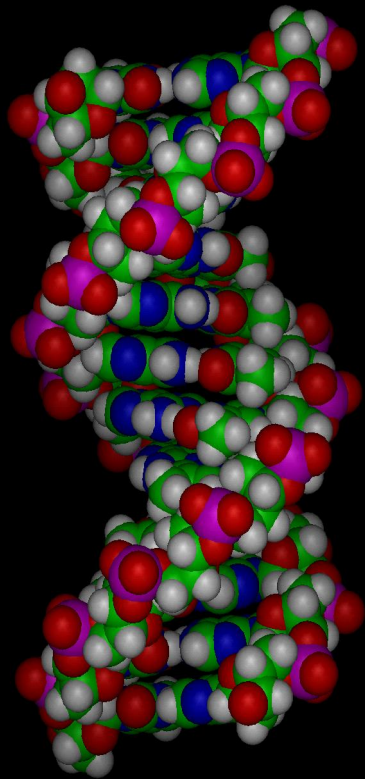CAL, FJL, MO: *J. Phys.Chem* <u>99</u>, 11591-11599 (1995).

# Current limitations in MD

- Size of the system
  - Typically: $10^4$ - $10^5$ particles
  - Flagship: $10^6$

- Simulation length
  - Typically: $10^1$ – $10^2$ ns
  - Using HPC: $\mu s$
  - For $10^4$ particles

Longer and bigger,...

Data scales as 3N*T
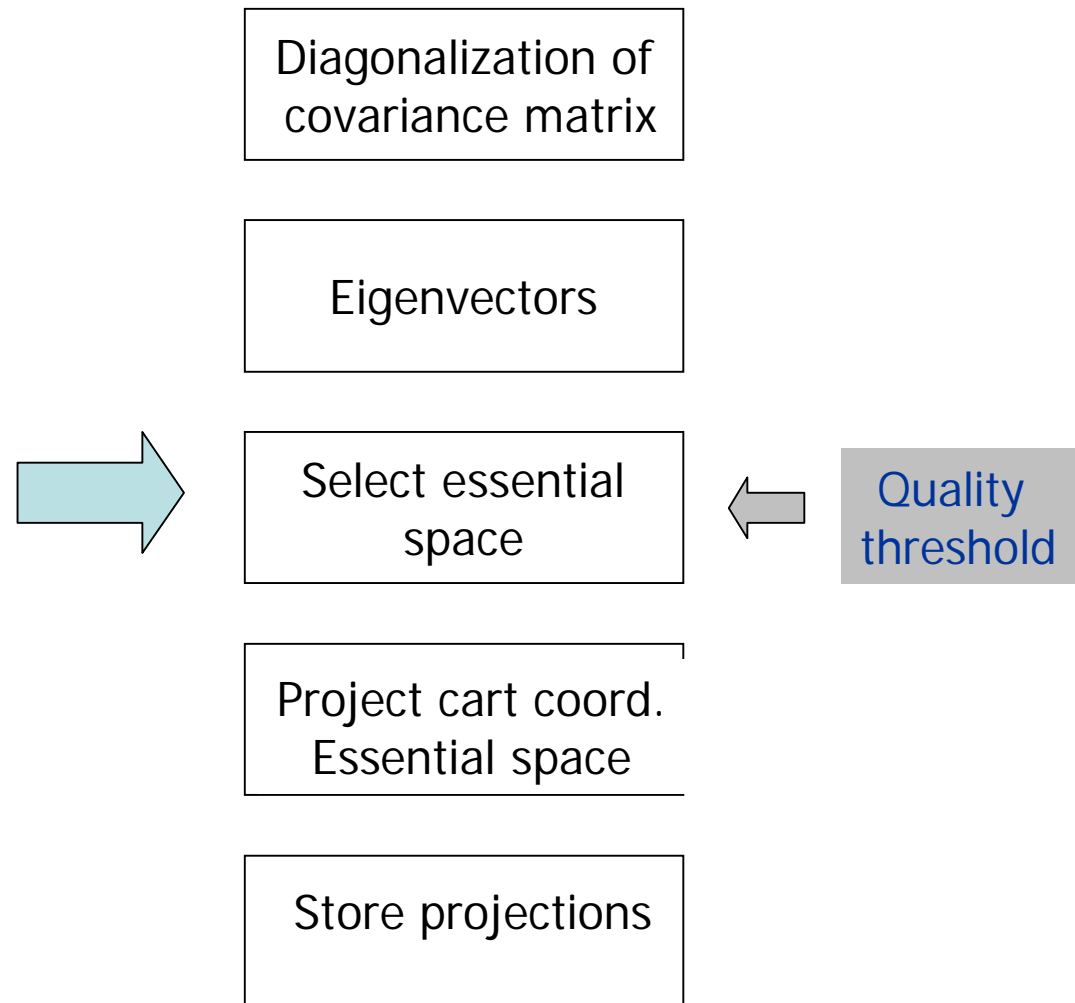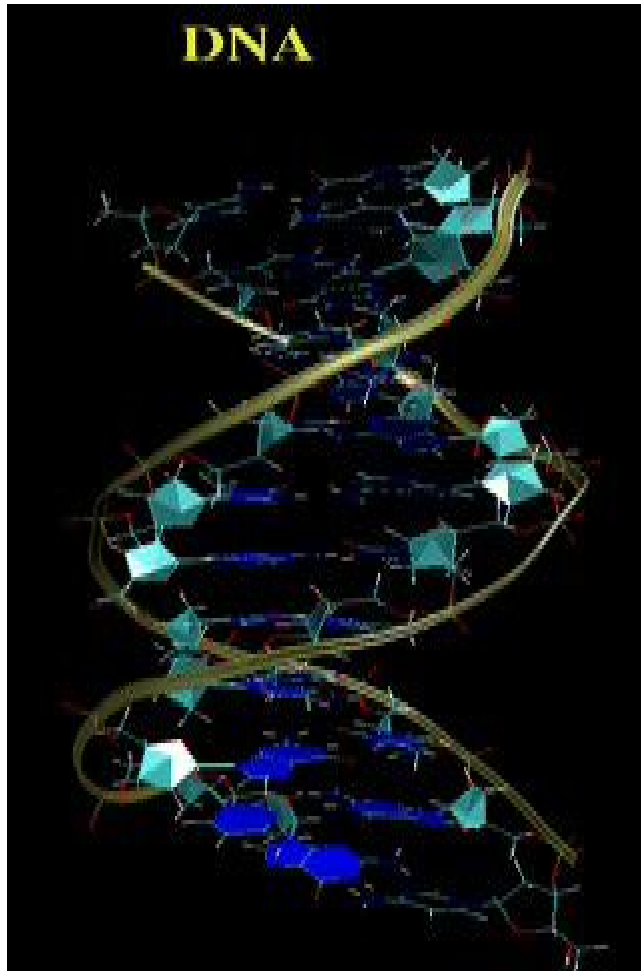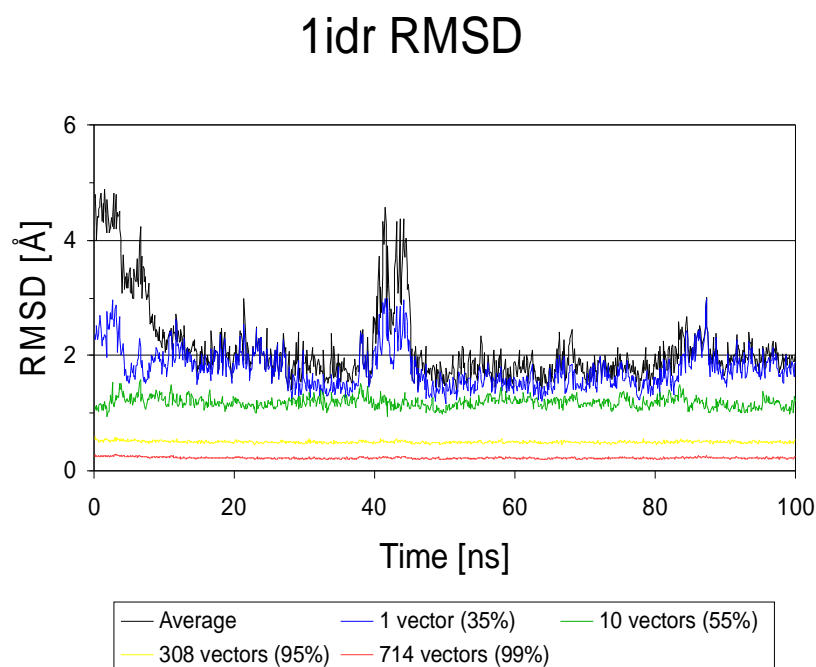N= number of atoms
T= time length
Ex. Scale up $3 \times 10^3$

10 ns

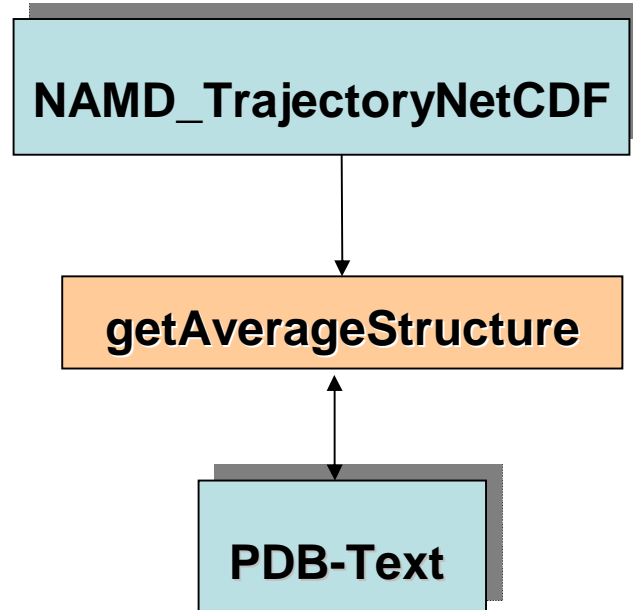Individual trajectories 0.1 Tb

1000 ns

# Compression tools (PCAZIP)



Diagonalization of covariance matrix

Eigenvectors

Select essential space ← Quality threshold

Project cart coord. Essential space

Store projections

Col. C.Laughton

# PCAZIP data reduction proteins

1idr RMSD



| | 95% cutoff | | |
|---|---|---|---|
| Protein | RMSd | File size | |
| 1ark | 0.36 | 8.5 | |
| 1cei | 0.36 | 7.8 | |
| 1sr0 | 0.45 | 6.0 | |
| 2gb1 | 0.29 | 10.0 | |
| 3ci2 | 0.36 | 8.6 | |
| 2icb | 0.33 | 8.8 | |
| 1idr | 0.50 | 5.1 | |

# Web-based MD Trajectory Analysis Toolkit

- From Reference Format NetCDF.

**NAMD_TrajectoryNetCDF**
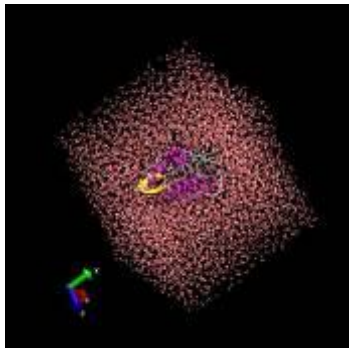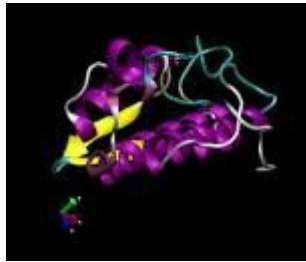
↓

**getAverageStructure** →

↕

**PDB-Text**

- *Bfactor*
- *Bfactor Per Residue*
- *Rms*
- *Rms Per Residue*
- *Superposition*
- *Average Structure*
- *Radius of Gyration*
- *SASA*
- *GaussianRmsd*
- *Esential dynamics*
- *etc.*
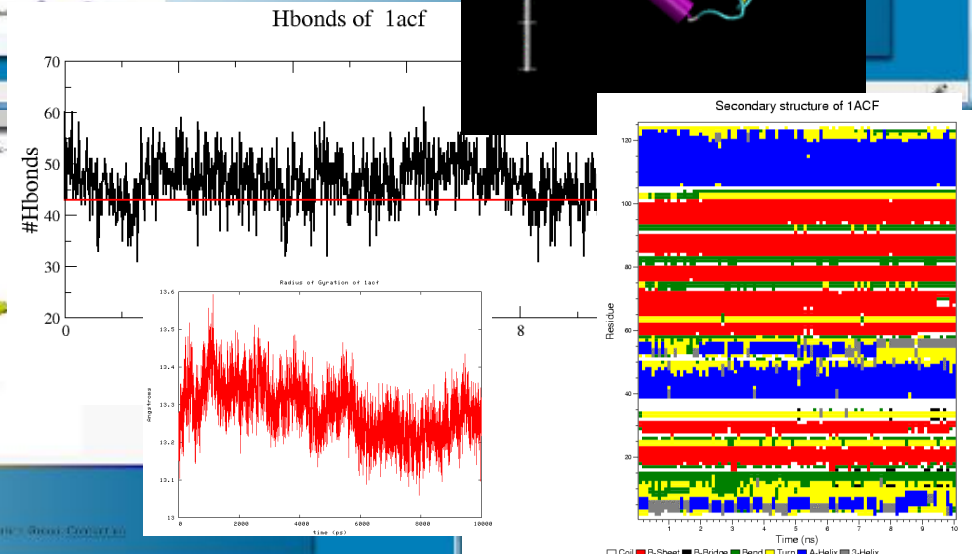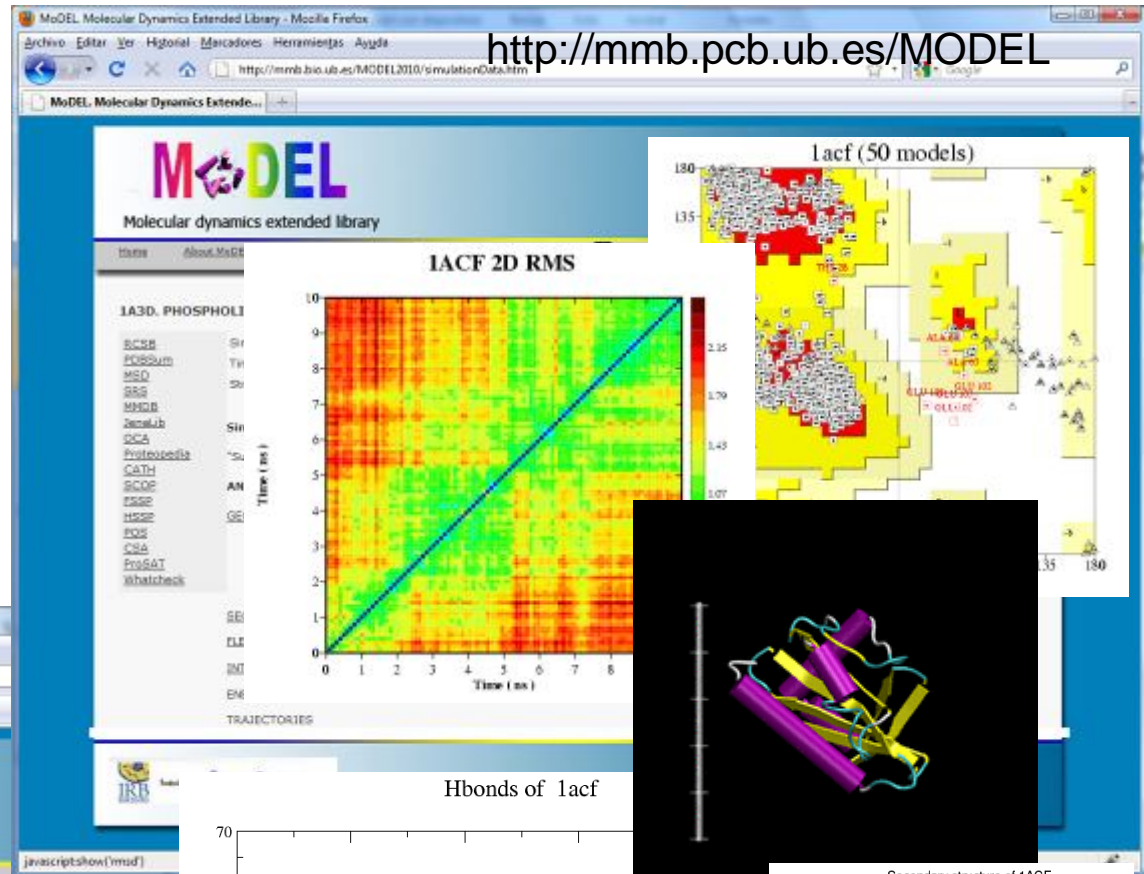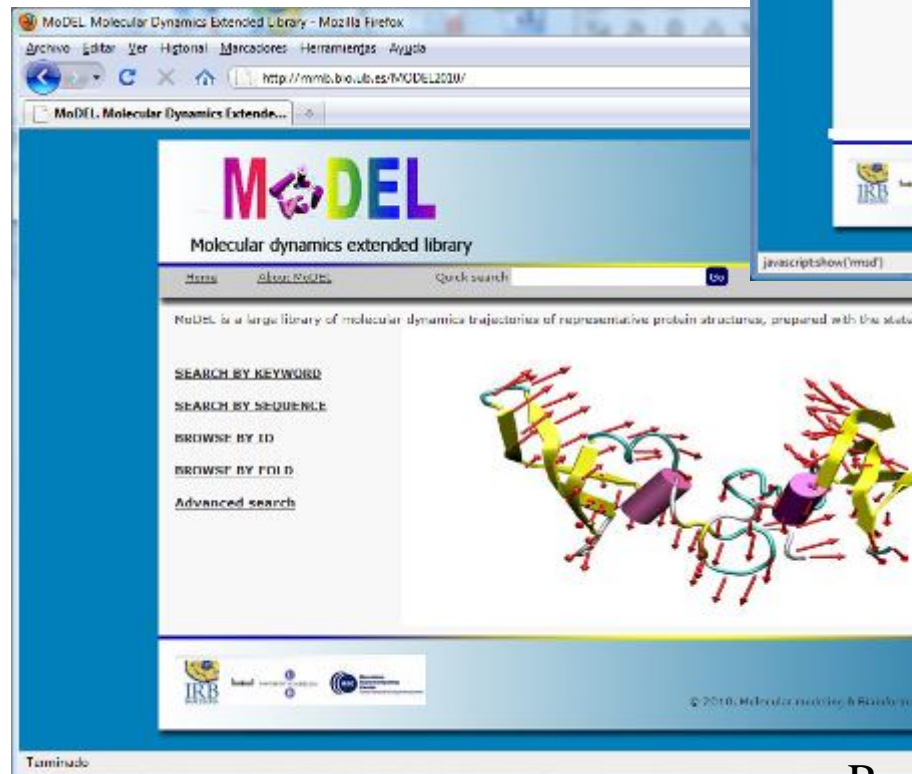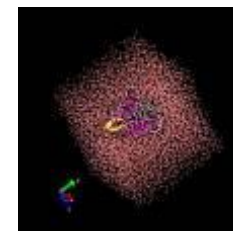
# Molecular Dynamics Workflow
# MDweb



**PDB Cleaning**

↓

**Adding Disulphide Bonds**

↓

**Adding Hydrogens**

↓

**Hydrogen Minimization**

↓

**Protein Minimization**

↓

**Adding Solvent/Ions**

↓

**System Equilibration**



**Inputs**

pdbCode

↓

Object

↓

getStructureFromPDB

↓

cleanPDB

↓

addDisulphideBondsPDB

↓

addHydrogensPDB

↓

runMDFromNAMD_Structure

↓

runMDFromNAMD_Structure1

↓

solvateProteinFromNAMD_Structure

↓

runMDFromNAMD_Structure2

↓

**Outputs**

preparedStructure

http://mmb.pcb.ub.es/MODEL

- 1800 proteins
- 20 Tb of data
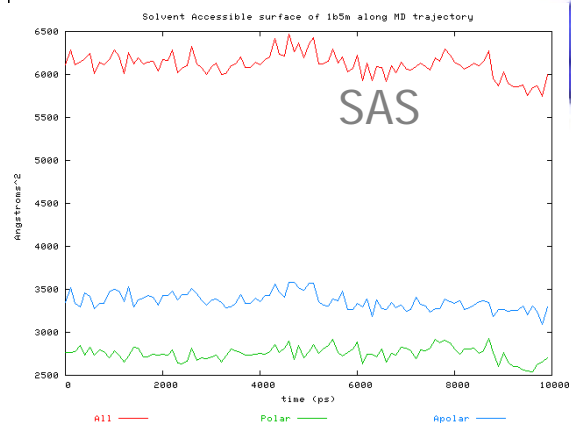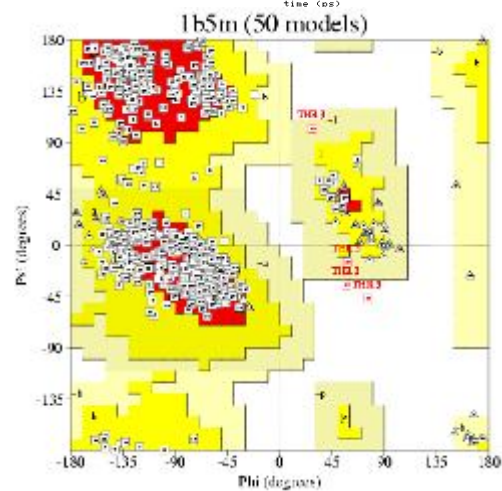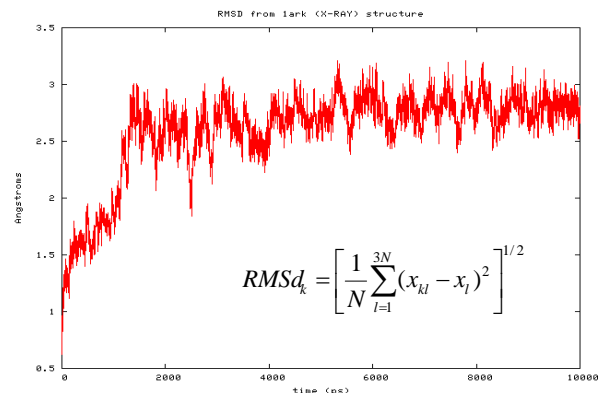- 250.000 residues
- 4.5 million protein atoms
- 19 million water molecules
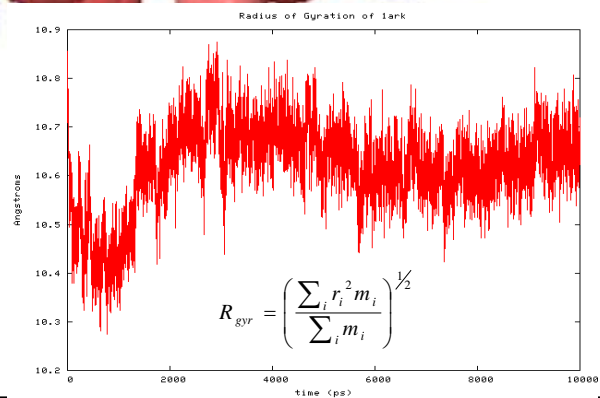
Rueda et al. PNAS 2007; Meyer et al., Structure 2010

RMSD from 1ark (X-RAY) structure

$$RMSd_k = \left[\frac{1}{N}\sum_{l=1}^{3N}(x_{kl}-x_l)^2\right]^{1/2}$$
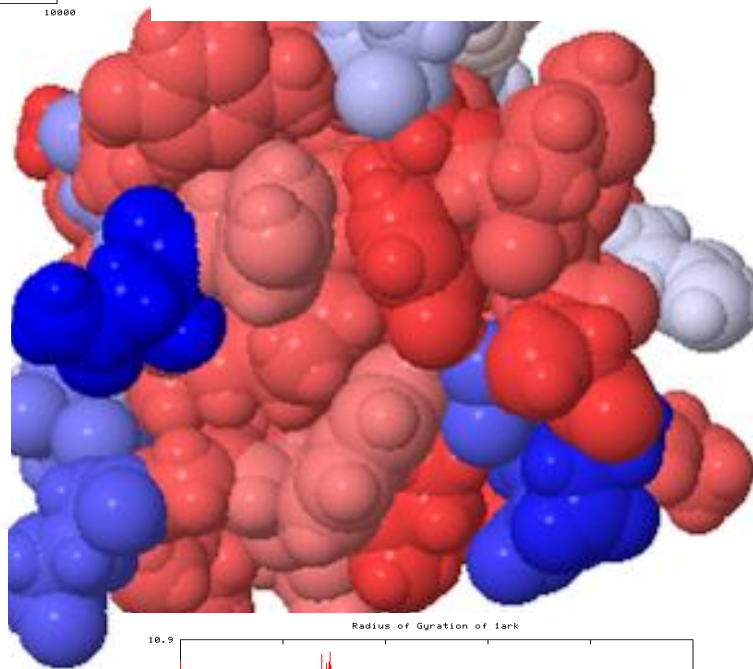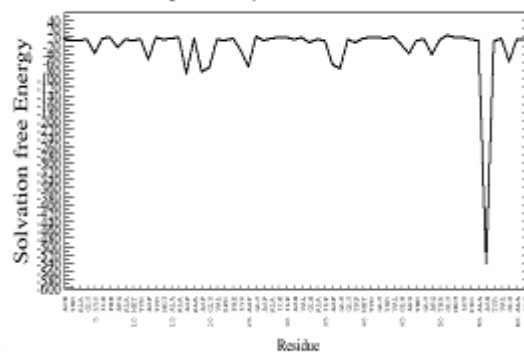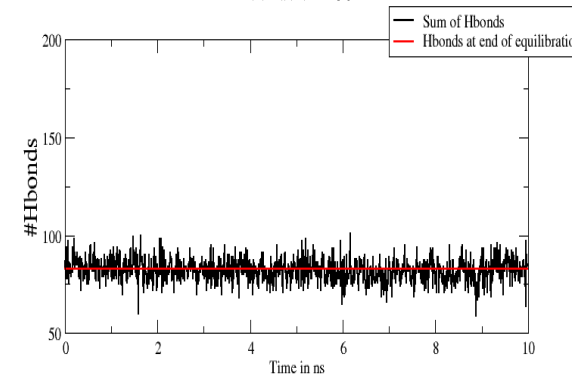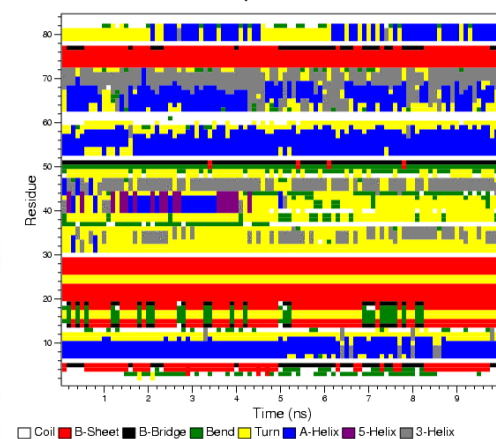
Per Residue Solvation free Energy of 1ark

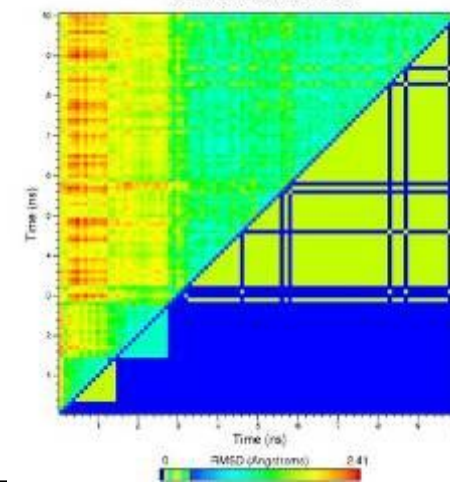Average over 100 snapshots of last nanosecond of simulation
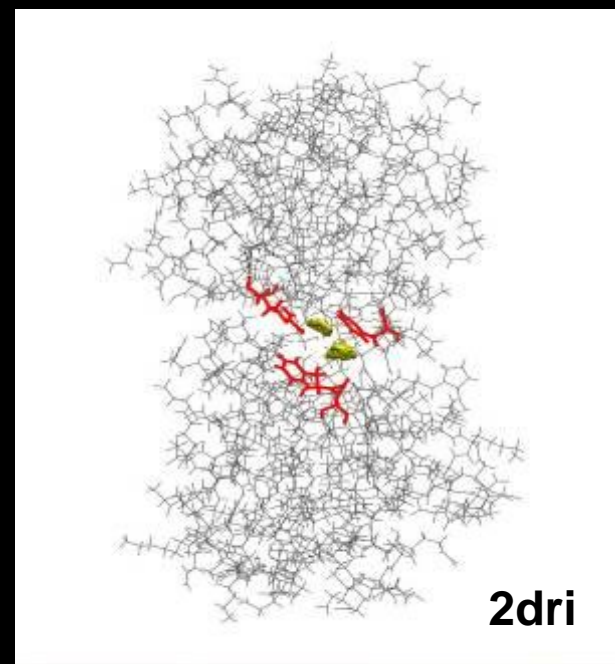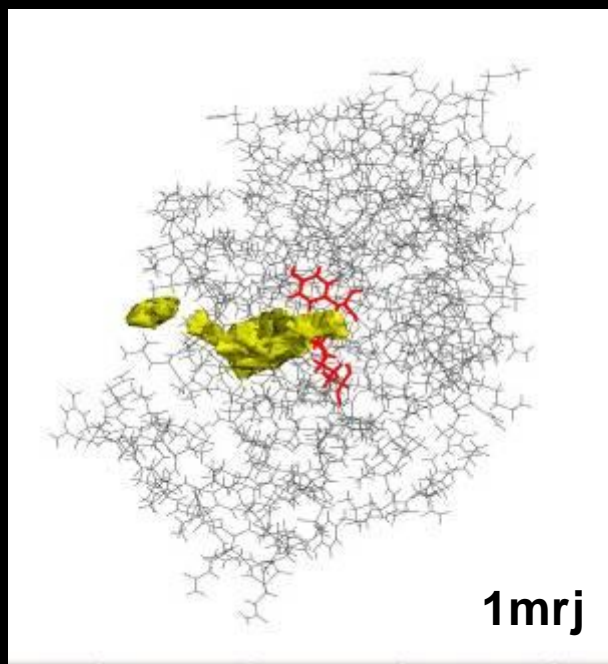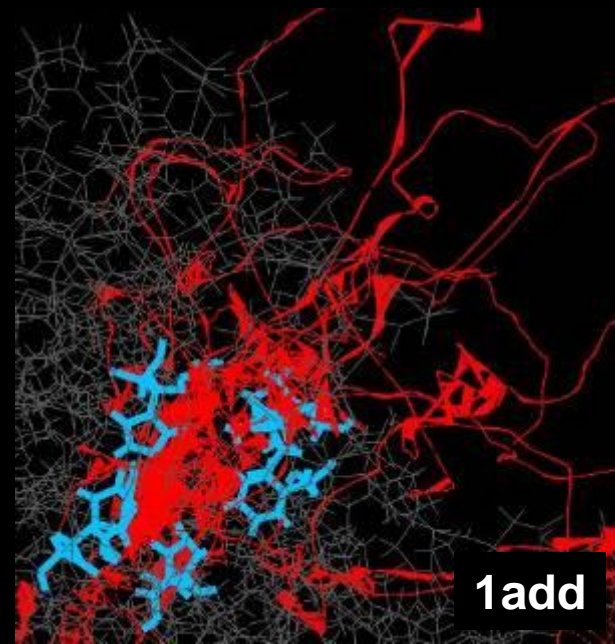
Hbonds of 153l

Sum of Hbonds
Hbonds at end of equilibration

1b5m (50 models)

Secondary structure of 1B5M

Coil B-Sheet B-Bridge Bend Turn A-Helix 5-Helix 3-Helix

Solvent Accessible surface of 1b5m along MD trajectory

SAS

All    Polar    Apolar

Radius of Gyration of 1ark

$$R_{gyr} = \left(\frac{\sum_i r_i^2 m_i}{\sum_i m_i}\right)^{1/2}$$

2D-RMS / Cluster for 1ark

RMSD (Angstroms)    0    2.41

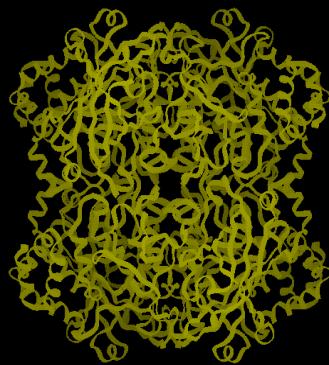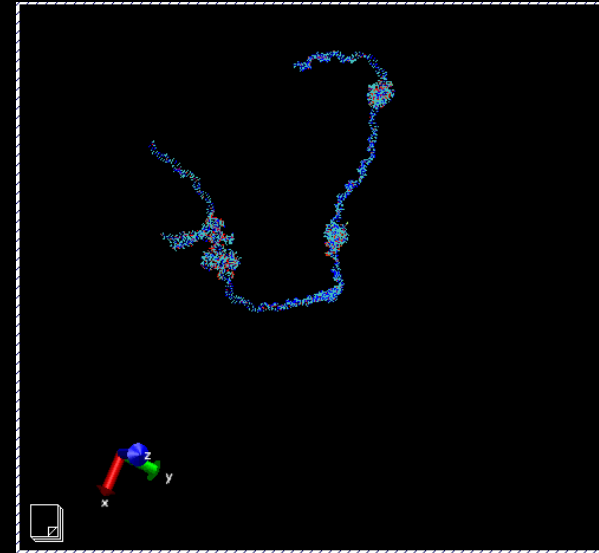4thi

1add

1mrj

2dri
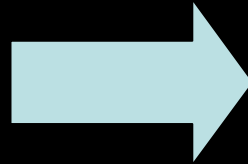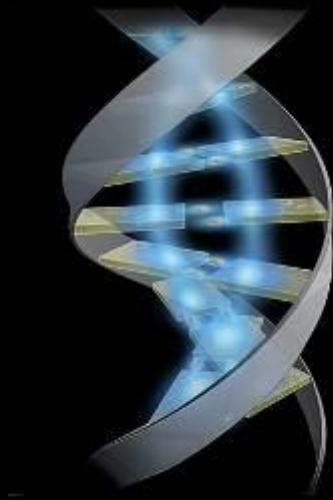
**1ADD**

**1MRJ**

**2DRI**

**4THI**

# Moving atomistic simulations to cell-scale



270 million Hemoglobin molecules