

[www.bsc.es](http://www.bsc.es)



**Barcelona  
Supercomputing  
Center**

*Centro Nacional de Supercomputación*

**A collaborative environment to produce, share  
and store DNA  
biomolecular simulations**

Ramon Goñi, PhD



3rd EUDAT User Forum

**University in Prague 23-24 April 2014**

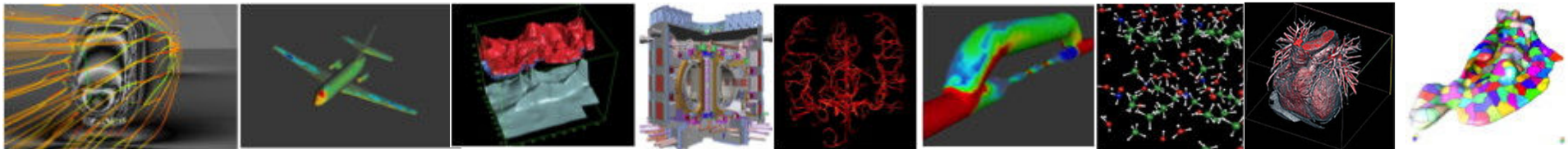
# Computer Simulation

## « Why and when we use it

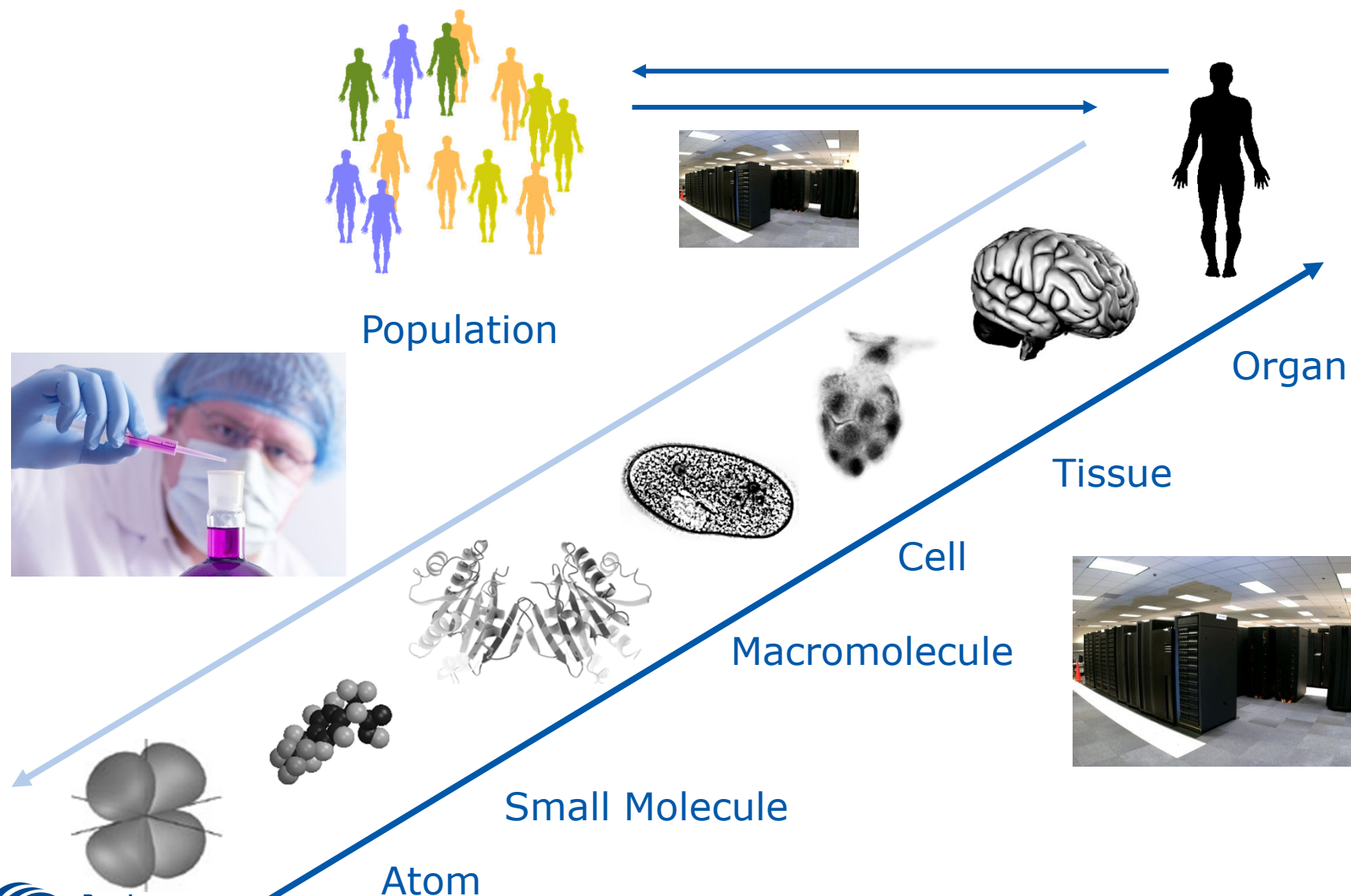
- To validate a known model
- As a cost-effective alternative
- As the only realistic approach to solve a problem

## « 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

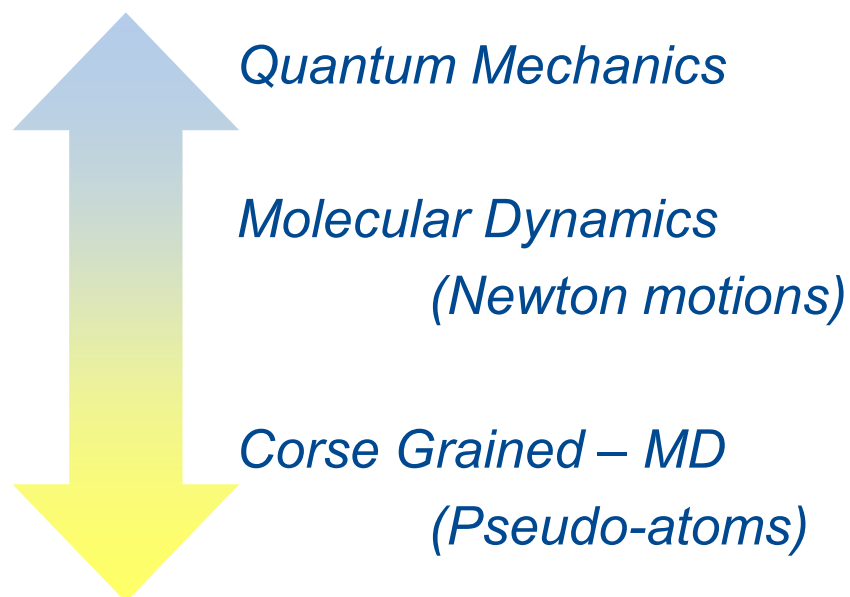


# Life Sciences and Health

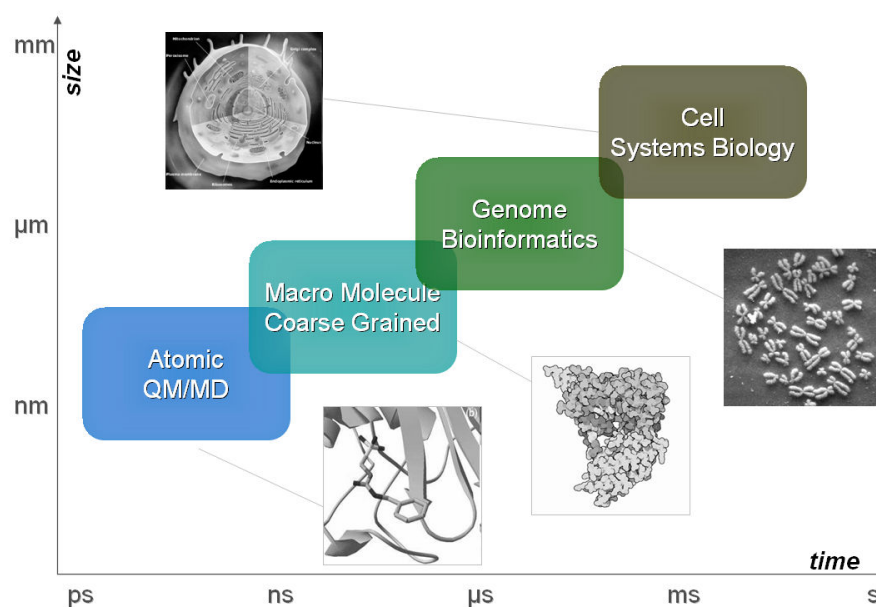


# Molecular Simulation

PRECISION



SPEED





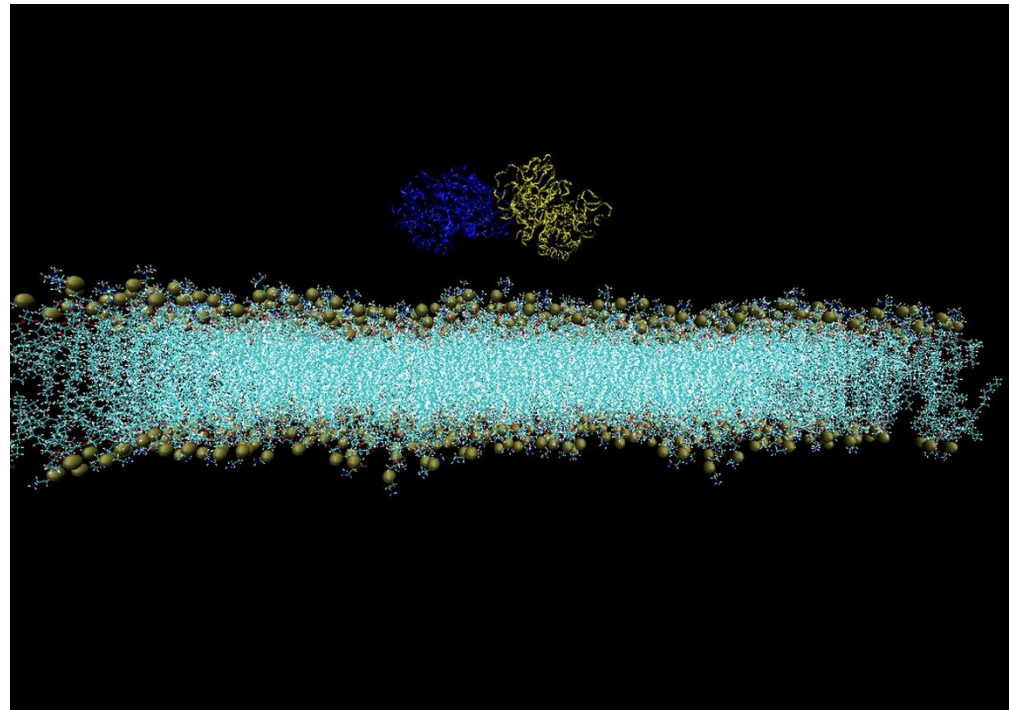
# Current limitations in MD

## « Size of the system

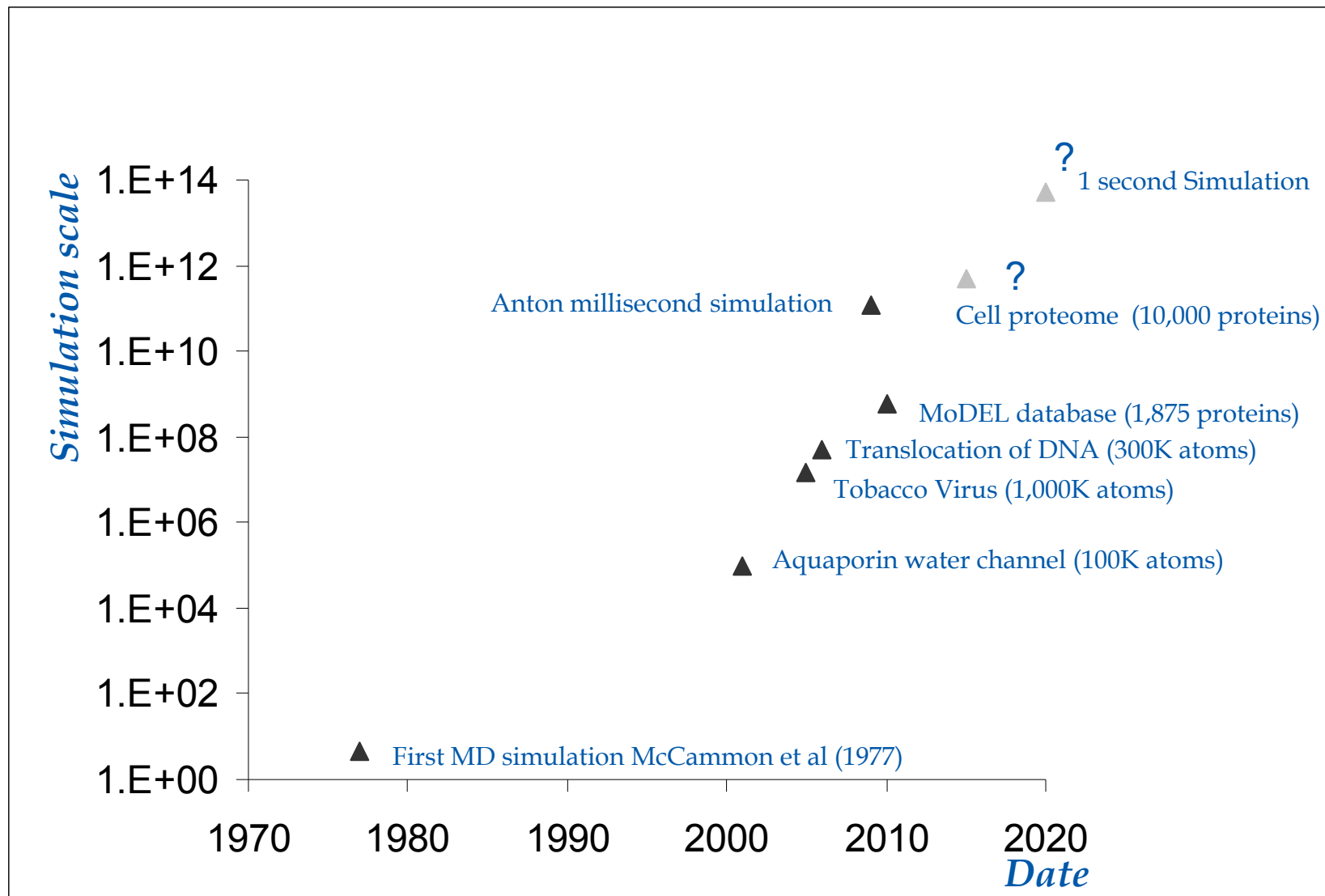
- Typically:  $10^4 - 10^5$  particles
- Flagship:  $10^7$

## « Simulation length ( $10^4$ particles )

- Typically:  $10^1 - 10^2$  ns
- Using HPC:  $\mu$ s
- Using Anthon: ms



# Trend in MD



# Global view of Simulation Information

## Metadata

Application. Version, Forcefield, Simulation time, solvent, etc.

## Pre- process

Output of solvation,  
equilibration, etc.

## Trajectory

Wet/Dry

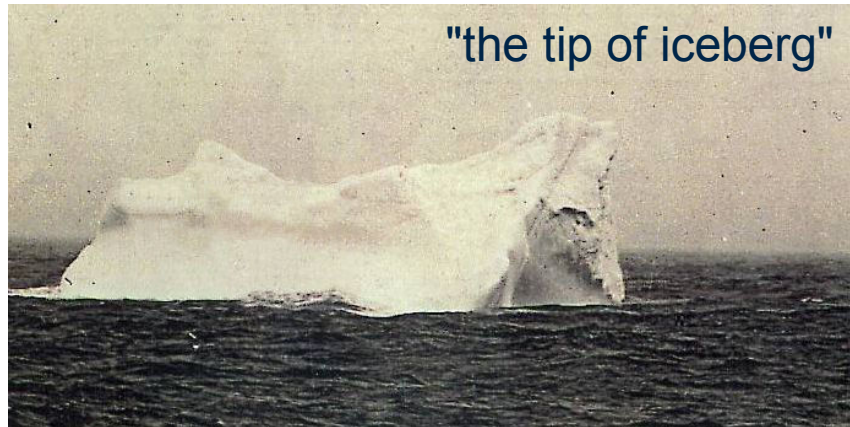
## Post- process

RMSD, B-Factors,  
PCA, Cavities, etc.

## External links

Links to of Protein, Structure, Small Molecule, Domain ,etc.

# Big Data Challenge



*Photo of the iceberg which was probably rammed by the RMS Titanic.*

## **European Exascale Software Initiative**

*Efficient data management, quick and fast interaction with the computer and flexibility in access to computer resources are in many fields of life sciences at least as important as total theoretical peak power.*

Only a few part of simulation data is visible

- Simulations unpublished
- Failed simulations necessary to generate the good one
- Simulations "lost" in system files
- Simulation repeated in fragmented in laboratories systems
- Errors identified too late (lack of interactivity)
- Backups, copies to share, temporal copies to transfer files
- ... and the cost of maintaining information in disk systems (power)



# First attempts

- Two initial attempts to provide software infrastructure to build biosimulation databases
  - BioSimGrid (2004)
  - P-Found (2006)
- They both shared the philosophy of having a central repository of trajectories that would allow obtaining a comprehensive view of biomolecular structure.
- At the time when these projects were started, computer power was still limited,

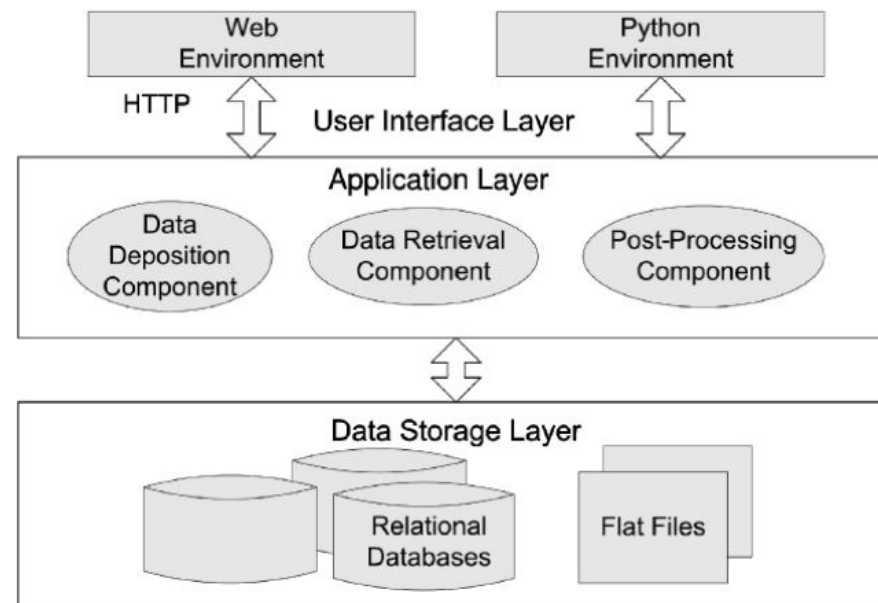
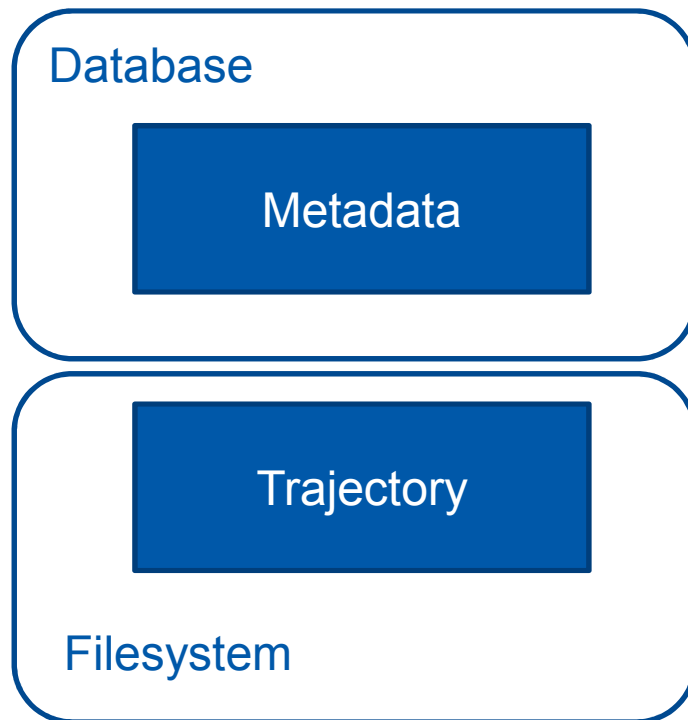
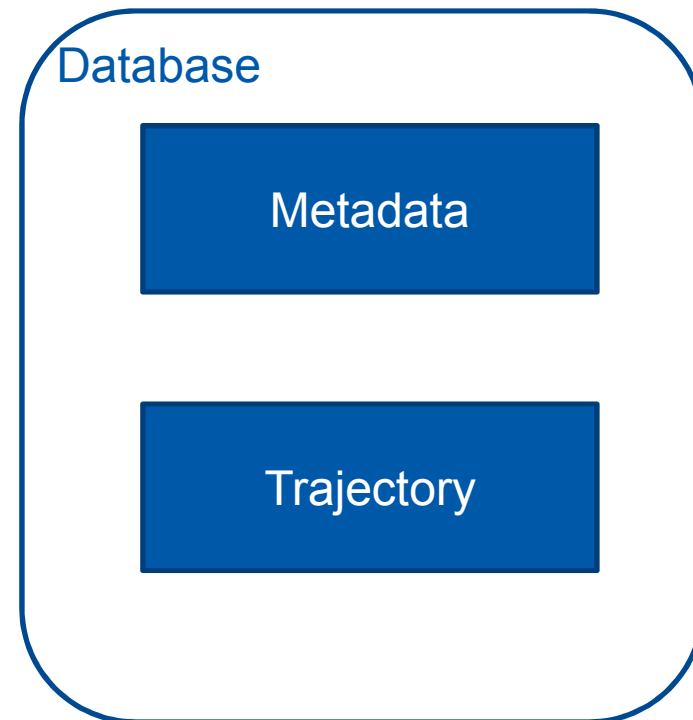


Fig. 1 The architecture of BioSimGrid.

# Strategies

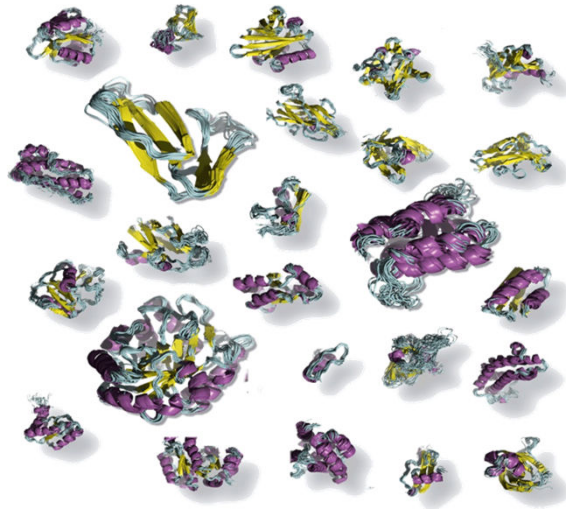


Heavy data file transfer



Poor Scalability  
Analysis tools not ready

# First repositories (2010)



## MySQL

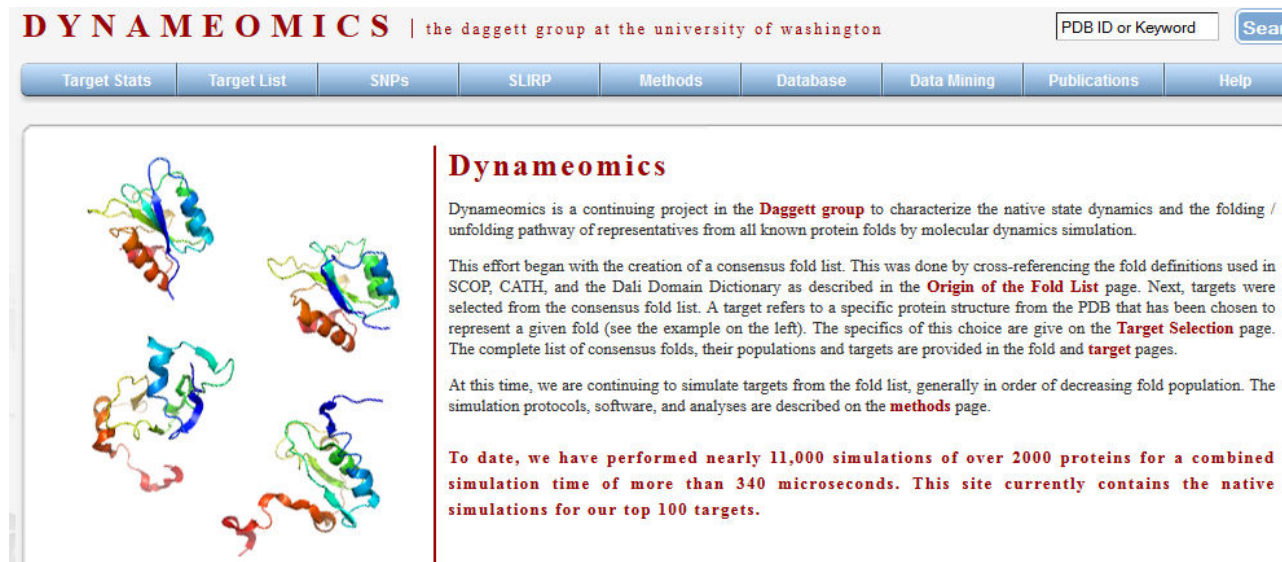
MODEL (Molecular Dynamics Extended Library)

[mmb.pcb.ub.es/MODEL](http://mmb.pcb.ub.es/MODEL)

- more than 1,800 entries (~20 Tb)
- covers around 40% of PDB structures, 8% of UniProtKB sequences, 29% of Human UniProtKB sequences and 33% of DrugBank proteins

## OLAP (SQL)

11,000 simulations  
of over 2,000  
proteins



**DYNAMEO MICS** | the daggett group at the university of washington

PDB ID or Keyword

Target Stats | Target List | SNPs | SLIRP | Methods | Database | Data Mining | Publications | Help

### Dyneameomics

Dyneameomics is a continuing project in the **Daggett group** to characterize the native state dynamics and the folding / unfolding pathway of representatives from all known protein folds by molecular dynamics simulation.

This effort began with the creation of a consensus fold list. This was done by cross-referencing the fold definitions used in SCOP, CATH, and the Dali Domain Dictionary as described in the **Origin of the Fold List** page. Next, targets were selected from the consensus fold list. A target refers to a specific protein structure from the PDB that has been chosen to represent a given fold (see the example on the left). The specifics of this choice are give on the **Target Selection** page. The complete list of consensus folds, their populations and targets are provided in the fold and **target** pages.

At this time, we are continuing to simulate targets from the fold list, generally in order of decreasing fold population. The simulation protocols, software, and analyses are described on the **methods** page.

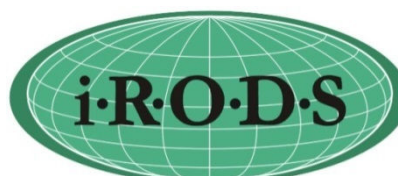
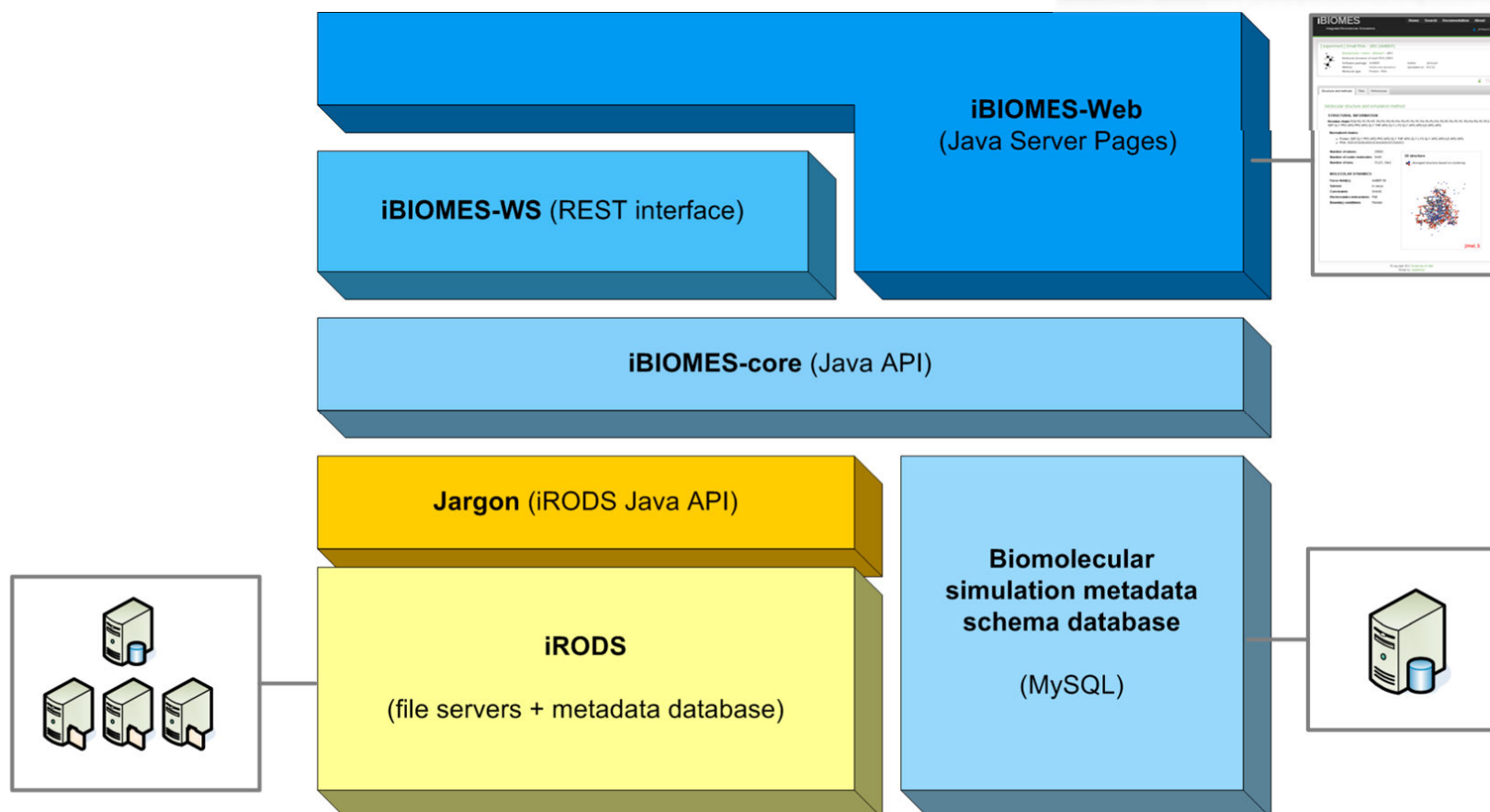
**To date, we have performed nearly 11,000 simulations of over 2000 proteins for a combined simulation time of more than 340 microseconds. This site currently contains the native simulations for our top 100 targets.**

# iRODS based



**iBIOMES**

Integrated biomolecular simulations  
julien.thibault@utah.edu

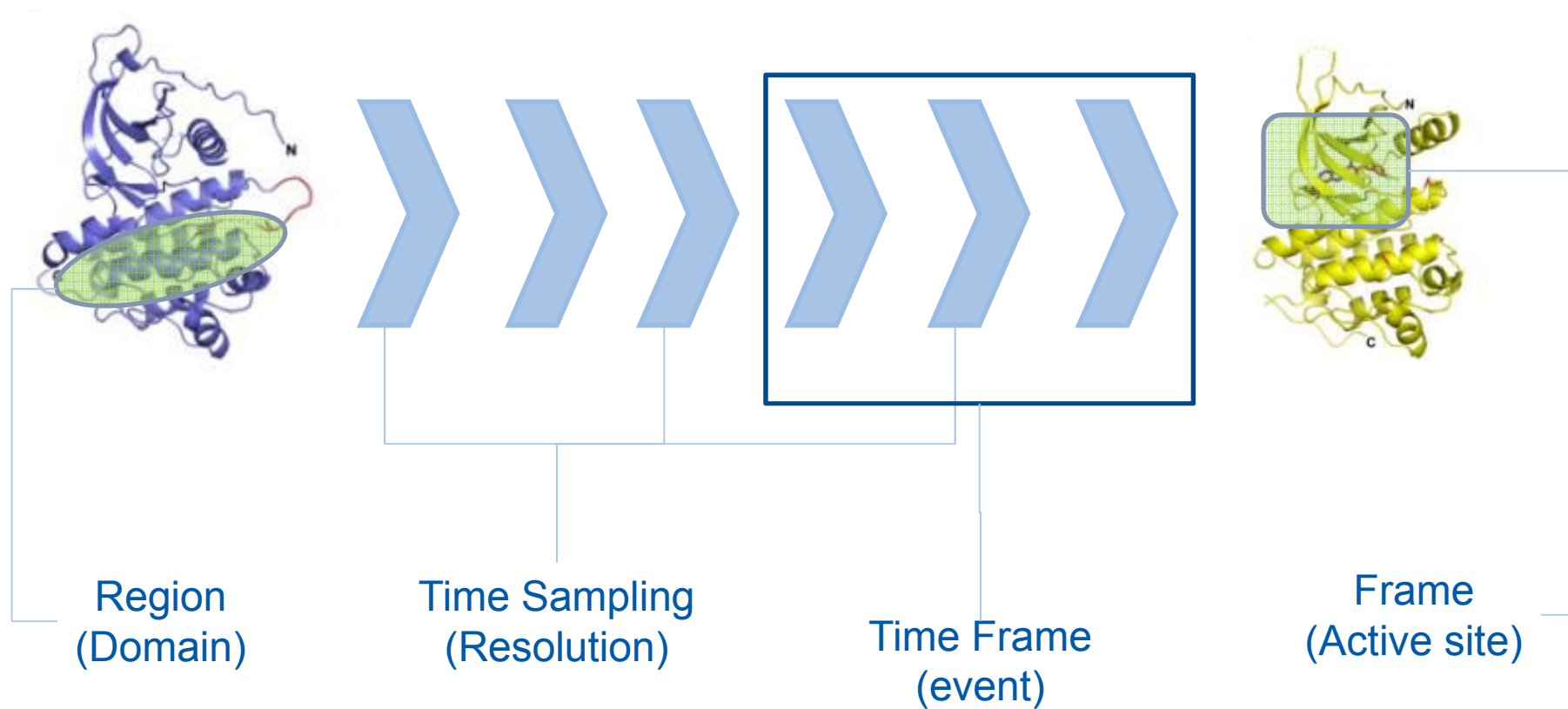


Integrated Rule-Oriented Data System

# Data Analysis

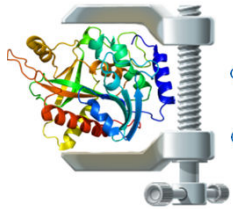
First snapshot

Last snapshot

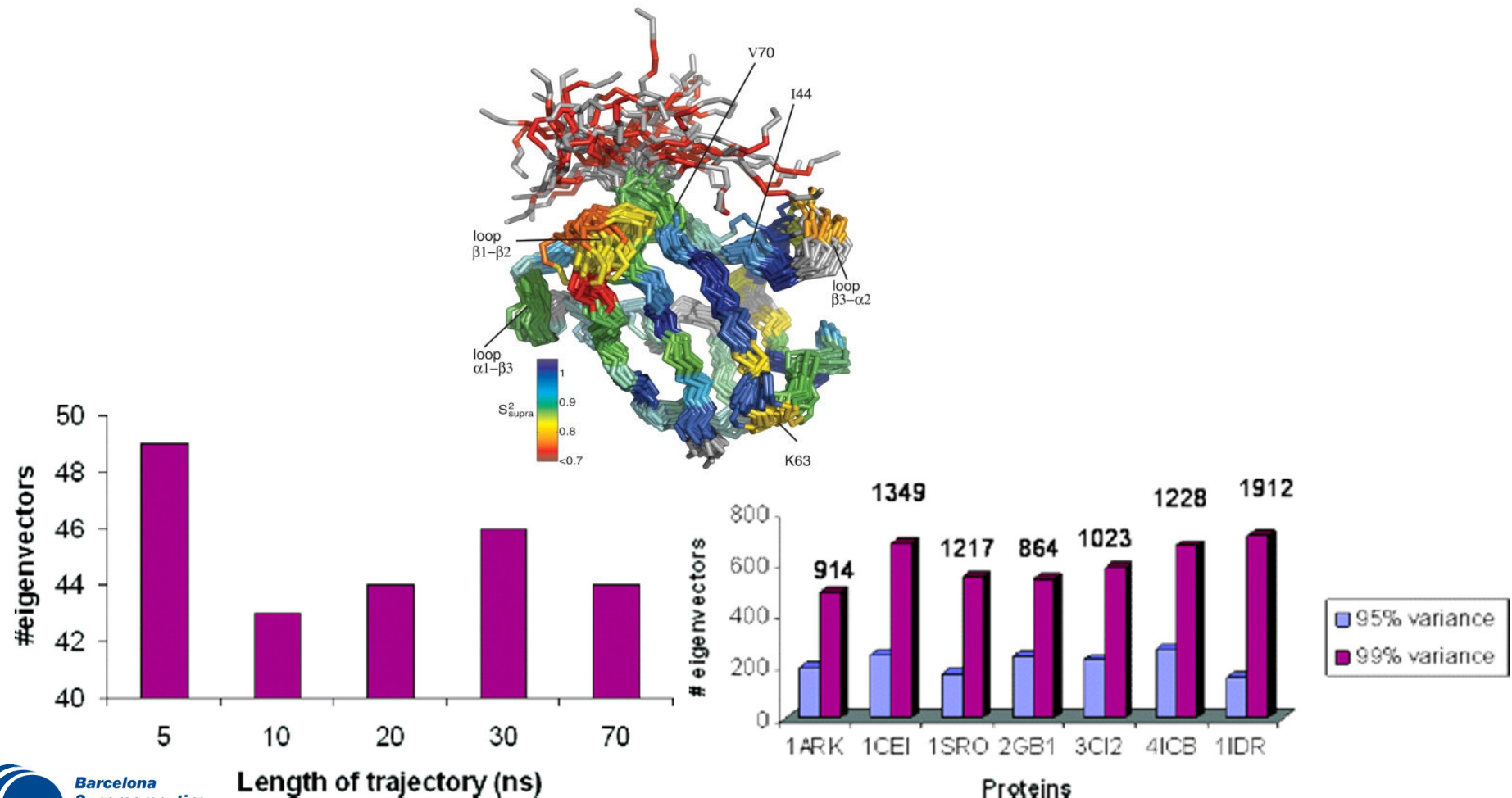




# PCA Compression



- Compression through principal components selection
- T. Meyer et al. J. Chem. Theor. Comp. 2006, 2 251-258



# Cooperative Biomolecular Simulation project

- “ The Ascona B-DNA Consortium (ABC) was set up following discussions during a conference in Ascona (Switzerland) in 2001.
- “ Its aim is to study the effect of base sequence on the structure, dynamics and interactions of DNA using molecular dynamics simulations.
- “ The initial goal of ABC was to generate a database containing structural and dynamic information on all unique tetra nucleotide sequences (2005)
- “ 1Tb of data! (compressed)



# Members

## US

D. Beveridge, Wesleyan U., USA  
T. Bishop, Tulane U., USA  
D. Case, Rutgers U., USA  
T. Cheatham, U. Utah, USA  
A. Pérez, USA  
R. Osman, Mount Sinai, NY, USA

## APAC

B. Jayaram, IIT New Delhi, India  
T. Singh, IIT New Delhi, India

## EUROPE

J. Curuksu, EPFL, Switzerland  
F. Lankas, Prague, Czech Rep.  
C. Laughton, Nottingham, UK  
R. Lavery, IBCP, France  
J. Maddocks, EPFL, Switzerland  
A. Michon, IBCP, France  
M. Orozco, Barcelona, Spain  
D. Petkeviciute, EPFL, Switzerland  
N. Spackova, Brno, Czech Rep.  
J. Sponer, Brno, Czech Rep.  
K. Zakrzewska, IBCP, France

# Results

Biophysical Journal Volume 92 June 2007 3817–3829

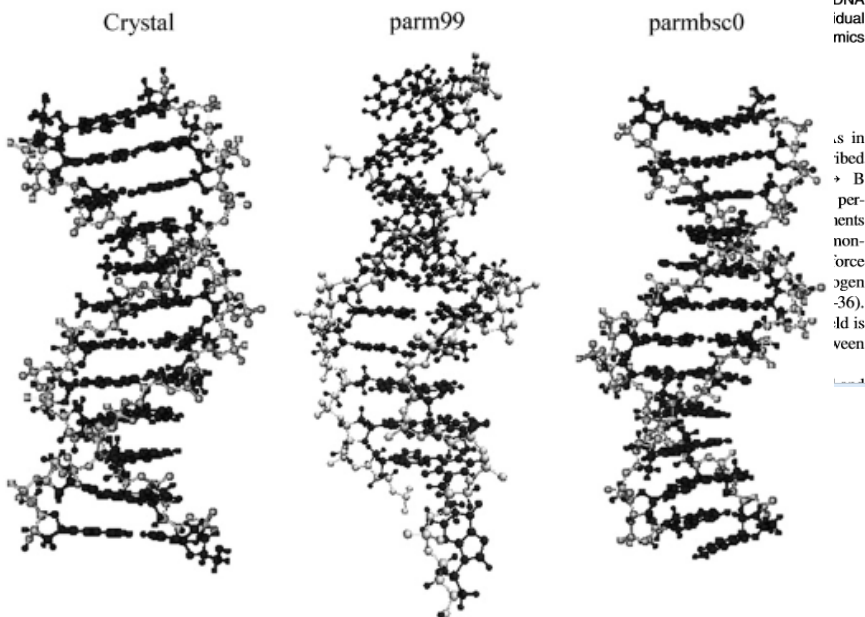
3817

## Refinement of the AMBER Force Field for Nucleic Acids: Improving the Description of $\alpha/\gamma$ Conformers

Alberto Pérez,<sup>\*,†</sup> Iván Marchán,<sup>\*,†</sup> Daniel Svozil,<sup>‡,§</sup> Jiri Sponer,<sup>§,¶</sup> Thomas E. Cheatham III,<sup>||</sup> Charles A. Laughton,<sup>\*\*</sup> and Modesto Orozco<sup>\*,†,††</sup>

<sup>\*</sup>Molecular Modeling and Bioinformatics Unit, Institut de Recerca Biomèdica & Instituto Nacional de Bioinformática, Parc Científic de Barcelona, Barcelona 08028, Spain; <sup>†</sup>Computational Biology Program, Barcelona Supercomputer Centre, Edifici Torre Girona, Barcelona 08028, Spain; <sup>‡</sup>Institute of Organic Chemistry and Biochemistry, Center for Biomolecules and Complex Molecular Systems, Academy of Sciences of the Czech Republic, 166 10 Prague 6, Czech Republic; <sup>§</sup>Institute of Biophysics, Academy of Sciences of the Czech Republic, 612 65 Brno, Czech Republic; <sup>¶</sup>Faculty of Science, Masaryk University, 611 37 Brno, Czech Republic; <sup>||</sup>Departments of Medicinal Chemistry, Pharmaceutical Chemistry and Pharmaceutics and Bioengineering, University of Utah, Salt Lake City, Utah 84112; <sup>\*\*</sup>School of Pharmacy and Centre for Biomolecular Sciences, University of Nottingham, Nottingham NG7 2RD, United Kingdom; and <sup>††</sup>Departament de Bioquímica i Biologia Molecular, Facultat de Biologia, Universitat de Barcelona, Barcelona 08028, Spain

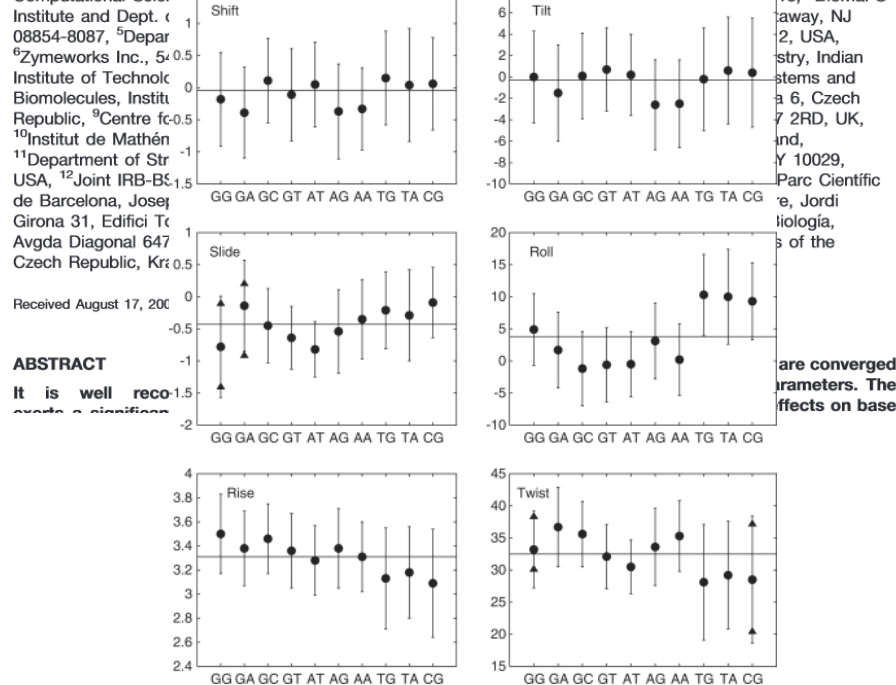
**ABSTRACT** We present here the parmbsc0 force field, a refinement of the AMBER parm99 force field, where emphasis has been made on the correct representation of the  $\alpha/\gamma$  concerted rotation in nucleic acids (NAs). The modified force field corrects overpopulations of the  $\alpha/\gamma = (g+,t)$  backbone that were seen in long (more than 10 ns) simulations with previous AMBER parameter sets (parm94-99). The force field has been derived by fitting to high-level quantum mechanical data and verified by comparison with very high-level quantum mechanical calculations and by a very extensive comparison between simulations



## A systematic molecular dynamics study of nearest-neighbor effects on base pair and base pair step conformations and fluctuations in B-DNA

Richard Lavery<sup>1,\*</sup>, Krystyna Zakrzewska<sup>1</sup>, David Beveridge<sup>2</sup>, Thomas C. Bishop<sup>3</sup>, David A. Case<sup>4</sup>, Thomas Cheatham III<sup>5</sup>, Surjit Dixit<sup>6</sup>, B. Jayaram<sup>7</sup>, Filip Lankas<sup>8</sup>, Charles Laughton<sup>9</sup>, John H. Maddocks<sup>10</sup>, Alexis Michon<sup>1</sup>, Roman Osman<sup>11</sup>, Modesto Orozco<sup>12</sup>, Alberto Perez<sup>12</sup>, Tanya Singh<sup>7</sup>, Nada Spackova<sup>13</sup> and Jiri Sponer<sup>13</sup>

<sup>1</sup>Institut de Biologie et Chimie des Protéines, CNRS UMR 5086/Université de Lyon, 7 passage du Vercors, 69367 Lyon, France; <sup>2</sup>Department of Chemistry, Wesleyan University, Middletown, CT 06459; <sup>3</sup>Center for Computational Science and Institute and Dept. of Chemistry, 08854-8087, Middletown, CT 06459; <sup>4</sup>Department of Chemistry, Wesleyan University, Middletown, CT 06459; <sup>5</sup>Center for Computational Science and Institute and Dept. of Chemistry, 08854-8087, Middletown, CT 06459; <sup>6</sup>Zymeworks Inc., 5000 Zymeworks Drive, San Diego, CA 92121; <sup>7</sup>Department of Chemistry, University of Illinois at Chicago, Chicago, IL 60607; <sup>8</sup>Department of Chemistry, University of Illinois at Chicago, Chicago, IL 60607; <sup>9</sup>Department of Chemistry, University of Illinois at Chicago, Chicago, IL 60607; <sup>10</sup>Institut de Mathématiques et de Physique, Université de Strasbourg, Strasbourg, France; <sup>11</sup>Department of Chemistry, University of Illinois at Chicago, Chicago, IL 60607; <sup>12</sup>Joint IRB-BSC, Institut de Recerca Biomèdica & Instituto Nacional de Bioinformática, Parc Científic de Barcelona, Josep Girona 31, Edifici Torre Girona, Barcelona 08028, Spain; <sup>13</sup>Department of Chemistry, University of Illinois at Chicago, Chicago, IL 60607



Received August 17, 2006

### ABSTRACT

It is well recognized that nearest-neighbor effects on base pair and base pair step conformations and fluctuations in B-DNA are converged parameters. The effects on base



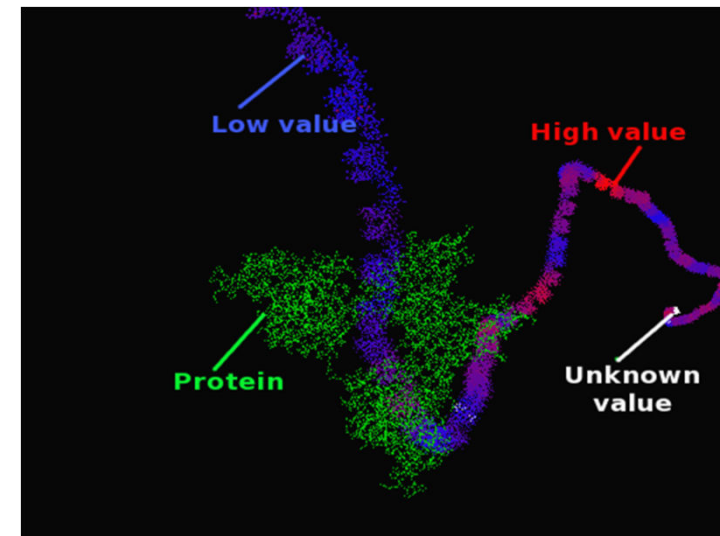
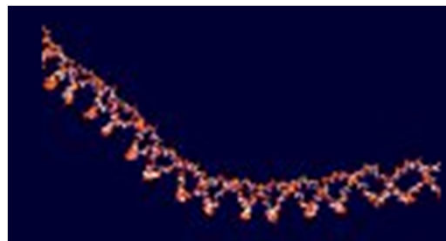
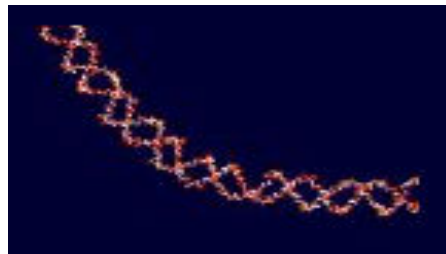
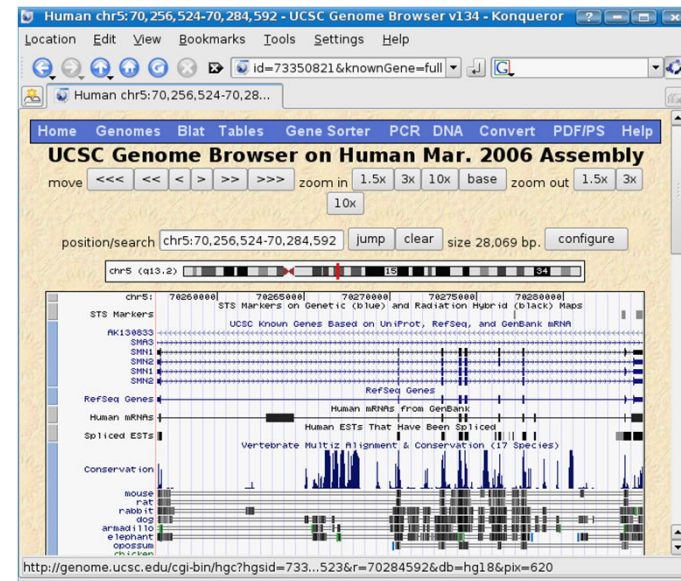
# DNA sequence and DNA structure

```

<3> 701~724      AGCCTTGTATCCGTATC-TTTCAA-----
<4> 132~154      AGCCTTGTATCCGTATC-TTTCA-----
<8> 264~288      -GCCTTGTATCCGTATC-TTTCAACG--
<7> 70~95        --CCTTGTATCCGTATC-TTTCAACGTG
<7> 344~366      --CCTTGTATCCGTATC-TTTCAAC--
<0> 304~326      ---CTTGTATCCGTATCTTTTCAAC---
<2> 482~497      ---CTTGTATCCGTATC-T-----
<2> 522~537      ----TTGTATCCGTATC-TT-----
<5> 440~461      -----GTCATCCGTATC-TTTCAACGTG
<6> 16~37         -----GTCATCCGTATC-TTTCAACGTG
<1> 910~925       -----CATCCGTATC-TTTCAA-----
<6> 986~1000      -----CATCCGTATC-TTTCA-----
<1> 1051~1069     -----ATCCGTATC-TTTCAACGTG
    
```

```

<1> 28~68         AACAAAGCA-A-ACTTTTATCCATGGTCGTTGTTACAGAGGGGTC
<4> 333~373      AACAAAGCA-A-ACTTTTATCCATGGTCGTTGTTACAGAGGGGTC
<8> 154~193      AACAAAGCA-A-ACTTTTATCCATGGTCGTTGTTACAGAGGGGTC
<6> 615~647      AACAAAGCAGA-ACTTTTATCCATGGTCGTTGTTAC-----
<4> 502~533      AACAAAGCA-ACCCTTTTATCCATGGTCGTTGTTA-----
<1> 844~872      AACAAAGCA-A-ACTTTTATCCATGGTCGTTGTTAC-----
<8> 194~220      -----A-ACTTTTATCCATGGTCGTTGTTACAGA-----
<5> 451~480      -----CTTTCA-ACGTGGTCGTTGTTACAAAGGGGTC
    
```





# Results

7220-7230 Nucleic Acids Research, 2013, Vol. 41, No. 15  
doi:10.1093/nar/gkt511

Published online 12 June 2013

Deniz et al. BMC Genomics 2011, 12:489  
http://www.biomedcentral.com/1471-2164/12/489



## Unravelling the hidden DNA structural/physical code provides novel insights on promoter location

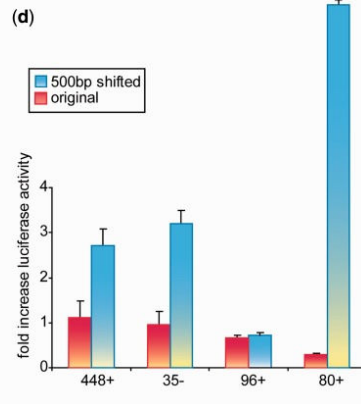
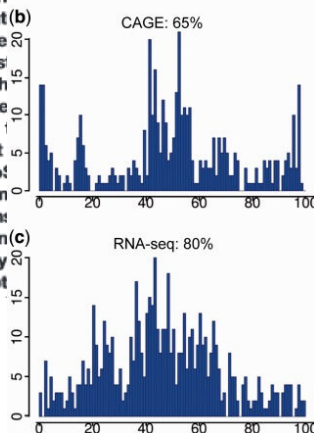
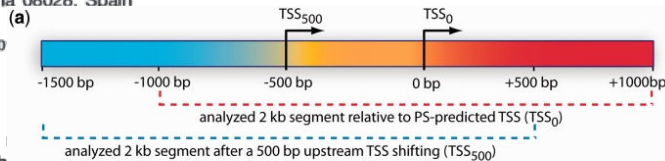
Elisa Durán<sup>1,2</sup>, Sarah Djebali<sup>3</sup>, Santi González<sup>2,4</sup>, Oscar Flores<sup>1,2</sup>, Josep Maria Mercader<sup>2,4</sup>, Roderic Guigó<sup>3</sup>, David Torrents<sup>2,4</sup>, Montserrat Soler-López<sup>1,2</sup> and Modesto Orozco<sup>1,2,4,5,\*</sup>

<sup>1</sup>Institute for Research in Biomedicine (IRB Barcelona), Barcelona 08028, Spain, <sup>2</sup>Joint IRB-BSC Research Program on Computational Biology, Barcelona 08028, Spain, <sup>3</sup>Bioinformatics and Genomics Group, Center for Genomic Regulation and Universitat Pompeu Fabra, Barcelona 08003, Spain, <sup>4</sup>Barcelona Supercomputing Center, Barcelona 08034, Spain and <sup>5</sup>Department of Biochemistry and Molecular Biology, University of Barcelona, Barcelona 08028, Spain

Received October 23, 2012

### ABSTRACT

Although protein promoter regions have been considered as critical regulatory elements, the location of promoter start sites (TSSs), as well as the performance of a comprehensive promoter sequence analysis, is still defined by distinct methods. In our previous representative sample, we subjected to extensive analyses. Interestingly, we found a specific sequence motif



motifs are still considered as transcription start sites (TSSs) still assumed as the only signals for gene expression. However, we have found that the city of unusual motifs at the TSS is located

### RESEARCH ARTICLE

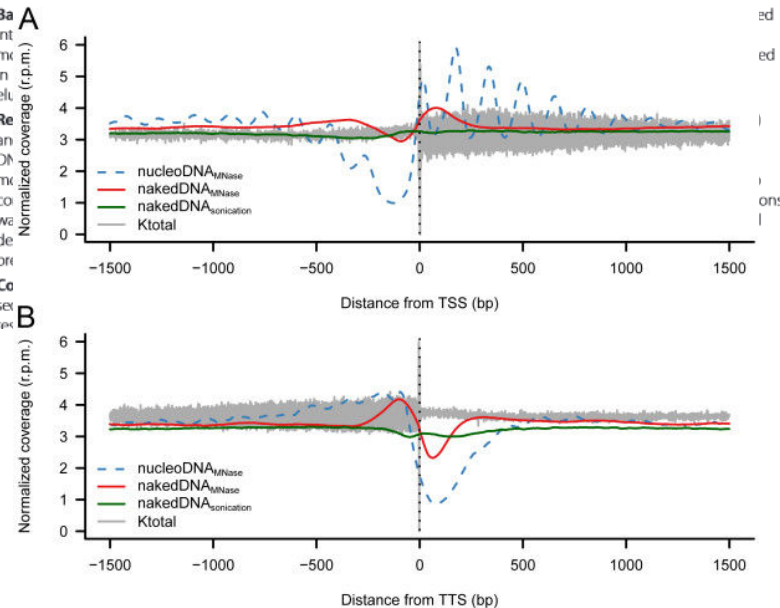
### Open Access

## Physical properties of naked DNA influence nucleosome positioning and correlate with transcription start and termination sites in yeast

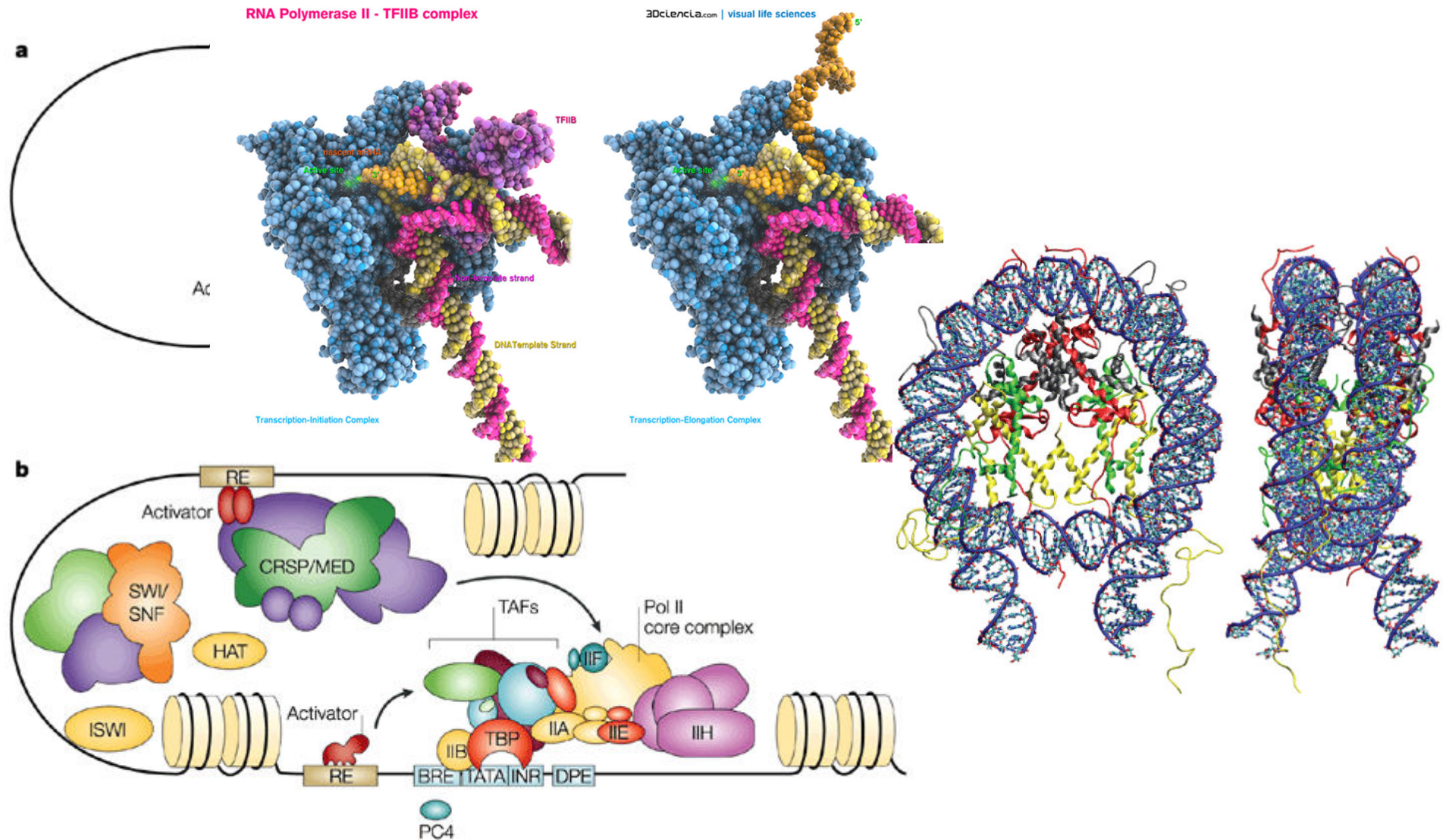
Özgen Deniz<sup>1†</sup>, Oscar Flores<sup>1†</sup>, Federica Battistini<sup>1</sup>, Alberto Pérez<sup>2</sup>, Montserrat Soler-López<sup>1</sup> and Modesto Orozco<sup>1,3,4\*</sup>

### Abstract

**Background:** The physical properties of naked DNA influence nucleosome positioning and correlate with transcription start and termination sites in yeast.



# Next Steps (ABC 2014)



# ABC Meeting 2014: New set up

- Big trajectory files
  - Trajectory Format
  - Compression & I/O efficiency



- Data standardization
  - Data distribution
  - Efficient interoperability

- « Data Management
  - Distributed environment
  - Scalability





# ABC & EUDAT



Home Services & Support EUDAT Communities EUDAT Events Working Groups News & Publications



Search 

- “ Consortium members are data producers, we need a partner to store results and guarantee open access in the long term. /B2SAFE/
- “ Data production is computer intensive (HPC), data analysis too. The consortium is not centralized in a single region. We are looking for a solution that optimizes data transfer.
- “ Simulation data generates heavy trajectory files /B2STAGE/, but also a large set of small unstructured files with metadata, pre-processing information and analysis results /B2SHARE/.

# In the long term

- Build an integrated HPC-Big data solution for MD (not only DNA)
- ... and integration with Life Sciences databases ('omics, etc.)

