# Fragment-Based Protein Structure Prediction and Design
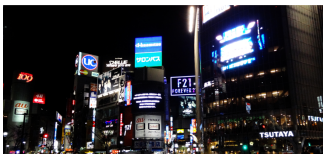
David Simoncini

Mathématiques et Informatique appliquées de Toulouse, INRA

November 19, 2014

**PhD. Thesis, Computer Science**

- University of Nice-Sophia Antipolis, France
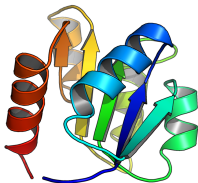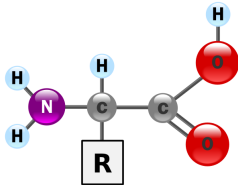- Selective pressure in cellular evolutionary algorithms



**Postdoctoral fellowship**



- Zhang IRU, RIKEN, Japan
- Fragment-based protein structure prediction
- Computational protein design

# What are proteins ?



### Amino acids (and primary structure)

Building blocks of proteins. A sequence of amino acids is called **primary structure** of a protein. Each amino acid in the sequence is called a **residue**.

### Secondary structure

There are two types of recurrent interactions between residues which lead to two regular patterns in the structure of a protein: $\alpha$ **helices** and $\beta$ **sheets**

### Tertiary structure

Oh well, the 3D structure is called **tertiary structure**.

# Role of proteins

## In our body

- Enzymes
- Cell signaling and ligand binding
- Structural proteins

## In our industry

- Biofuel production
- Wastewater treatment

## In our diseases

- Alzheimer
- Parkinson
- Type II diabetes
- Cancer

# Experimental determination of protein structures

## Nucleic Magnetic Resonance spectroscopy

- Establish restraints between atoms (distance, angle, orientation)
- Build models that satisfy the restraints

## X-ray crystallography

- Crystallize a protein
- Shoot it with X-rays and get a diffraction map
- Recover the phase of the structure factors, and solve the structure

All the solved protein structures are centralized in the Protein DataBank (PDB).

# Protein structure prediction

All the information that is needed to determine the tertiary structure of a protein is contained in its sequence.
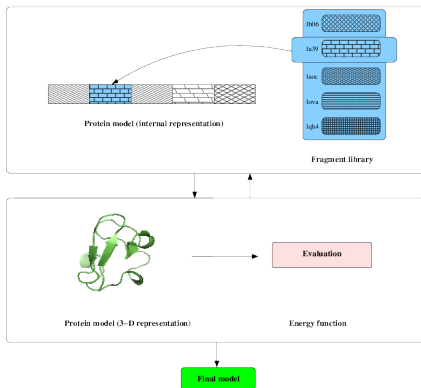
## Comparative modeling

- Identify an homologous protein by sequence alignment
- Generate models with the homologous protein as template

## Ab initio modeling (or free modeling)

Start "from scratch". Fold the models starting from the protein in its extended state, using only its sequence as information
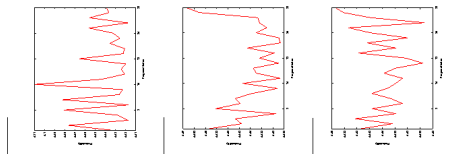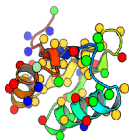
## How does it work ?

- Cut target sequence in small overlapping pieces
- Find fragments of structures in the PDB that match the small pieces
- Build a library of fragments (usually 25 per fragment window)
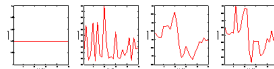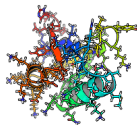- Assemble the fragments and minimize an energy function

9 residues frame

Coarse grained models

Estimation of distribution
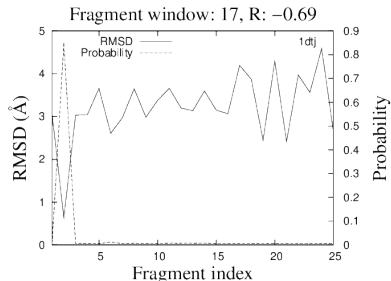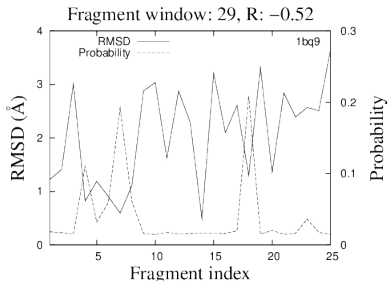
All atom models

# Experiments

## Benchmark

- 20 proteins
- 200,000 models generated by protein
- 4 iterations (50,000 models generated at each iteration)
- 200,000 models generated with Rosetta Modeling Suite for comparison
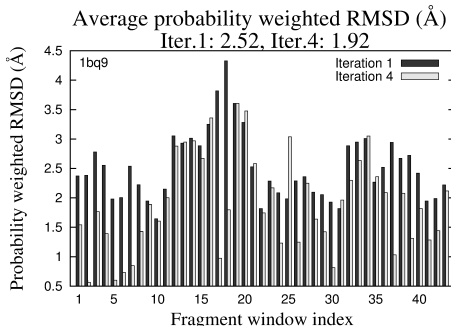
## How do we evaluate

- We analyze the impact of the EDA
- We look at the energy of the models
- We look at how close we are from the native structure (the real structure)
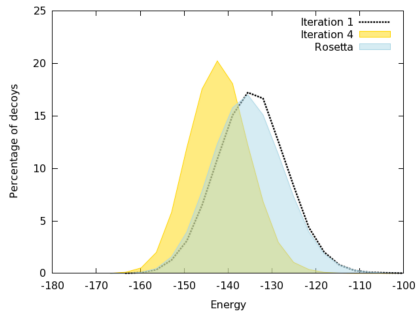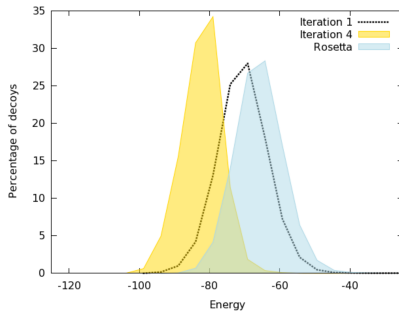
# Probabilities of selecting fragments



One window contains 25 fragments. The Root Mean Square Deviation (RMSD) of each fragment to the native one is plotted as well its probability of being selected.

# Probabilities of selecting fragments



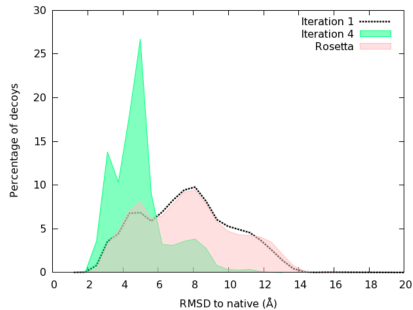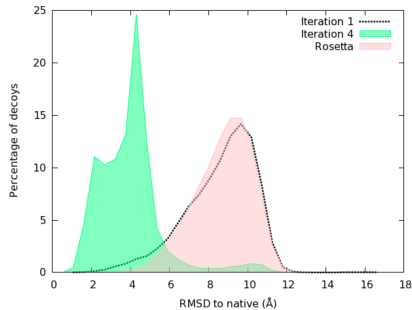Average probability weighted RMSD (Å)
Iter.1: 2.52, Iter.4: 1.92

The RMSD of each fragment to the native one is weighted by its probability of being selected. We plot the average over 25 fragments for each window.

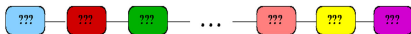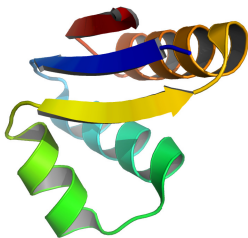# Energy minimization



Fragment-Based Protein Structure Prediction and Design

# RMSD to native structure

# Computational Protein Design

Protein design is the problem of finding an amino acid sequence compatible with a given protein backbone.



### Default method

Minimize a score function that captures residue interactions in the structure. 20 values: the 20 natural residues. Size of the search space: $20^n$ where $n$ is the length of the sequence to predict.

# Extracting knowledge to guide the search



## Our approach

- Establish an interaction profile for our query structure
- Find interaction matches in a database for each sequence position
- Rank matches according to surface overlap

# Interaction library

```
BEGIN 7
9 W T 3 L 4 E 5 T 6 R 8 T 9 Y 10 R 11 V 2
9 L I 3 A 4 C 5 I 6 A 8 S 9 S 10 G 11 V 2
9 M Q 3 Y 4 L 5 E 6 L 8 K 9 S 10 I 11 V 2
9 M T 3 W 4 A 5 A 6 E 8 E 9 L 10 V 11 I 2
9 Y T 3 G 4 K 5 Q 6 L 8 E 9 S 10 L 11 A 2
9 L I 3 A 4 R 5 M 6 D 8 S 9 L 10 G 11 I 2
9 I P 3 L 4 L 5 L 6 A 8 E 9 Q 10 I 11 L 2
9 I R 3 P 4 E 5 R 6 K 8 Q 9 L 10 Q 11 A 2
9 L A 3 Q 4 K 5 V 6 A 8 L 9 A 10 Q 11 I 2
9 I I 3 I 4 M 5 K 6 L 8 E 9 K 10 K 11 V 2
9 W A 3 K 4 G 5 F 6 L 8 G 9 F 10 I 11 F 2
9 S S 3 W 4 K 5 A 6 F 8 L 9 R 10 A 11 I 2
9 V T 3 P 4 V 5 R 6 A 8 Q 9 S 10 L 11 I 2
9 Q R 3 N 4 I 5 I 6 L 8 A 9 S 10 I 11 F 2
9 D D 3 K 4 L 5 D 6 C 8 A 9 V 10 I 11 I 2
9 I K 3 T 4 K 5 D 6 V 8 N 9 G 10 L 11 G 2
9 M S 3 E 4 H 5 Q 6 E 8 L 9 L 10 D 11 M 2
9 M T 3 E 4 Q 5 V 6 H 8 L 9 M 10 V 11 V 2
9 I D 3 C 4 K 5 T 6 L 8 K 9 A 10 L 11 L 2
9 K Q 3 L 4 I 5 V 6 E 8 A 9 L 10 V 11 H 2
9 R A 3 C 4 A 5 K 6 V 8 V 9 R 10 D 11 C 2
9 L G 3 L 4 E 5 E 6 Y 8 K 9 L 10 W 11 V 2
9 L A 3 L 4 C 5 A 6 Y 8 L 9 E 10 S 11 L 2
9 A L 3 D 4 K 5 L 6 H 8 L 9 I 10 D 11 I 2
9 C N 3 L 4 A 5 S 6 Y 8 L 9 K 10 Q 11 A 2
END
```

## Method

- Construct a library of interactions
- Randomly select one fragment per window
- Mutate all residues described in one fragment
- Relax and score the structure
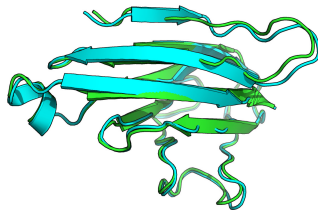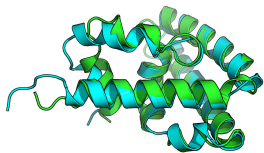- Update the probabilities to select each fragment

## Structural Homology Approach for DESign

- Source code in C++
- Uses Rosetta Macromolecular Modeling Suite to model proteins and score them [a]
- Uses an Estimation of Distribution algorithm for sampling: typically 25 iterations, 50,000 models generated
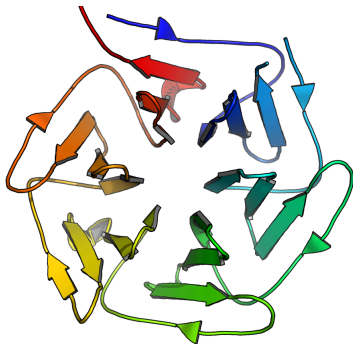- Parallelized with MPI

[a]https://www.rosettacommons.org/

**Protein structure recovery**
*Designer proteins obtained with a prototype of Shades and assessed with I-TASSER.*

## How did Nature create this ?

Hypothesis: by gene duplication and fusion

Problem: no structure with the exactly same sequence in each blade

# Application: symmetrical beta propellers



## Method highlights

- Generate a symmetrical backbone from a single blade with symmetric docking (RosettaDock)
- Derive a phylogenetic tree from the blade sequences to predict possible ancestral blade sequences (FastML)
- Map each sequence to the symmetrical backbone and relax the structure (Rosetta Relax)

# Conclusion

## Protein structure prediction

- Fragment-based approaches still are the best at free modeling
- The problem is still far from being solved
- Protein design can help

## Computational Protein design

- Few positive results
- Very hard to validate any new approach
- Many doors to open

**CPD as a cost function network**

Interactions can be decomposed as a sum of one-body and two-body energy terms.

$$E_S = \sum_{i=1}^{n} E(x_i^{r_1}) + \sum_{i=1}^{n-1} \sum_{j>i} E(x_i^{r_1}, y_j^{r_2})$$