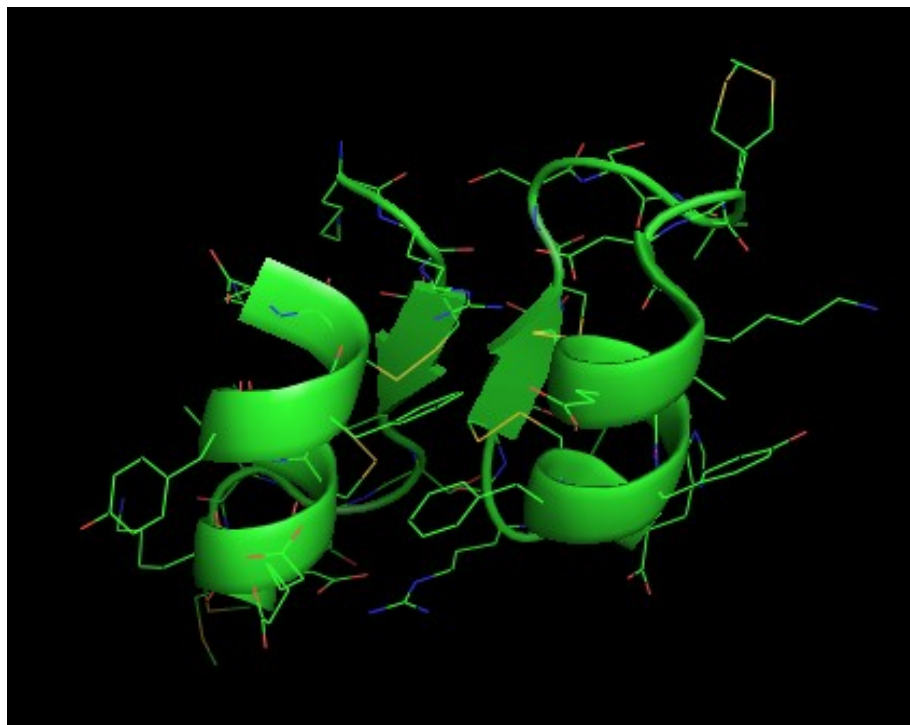


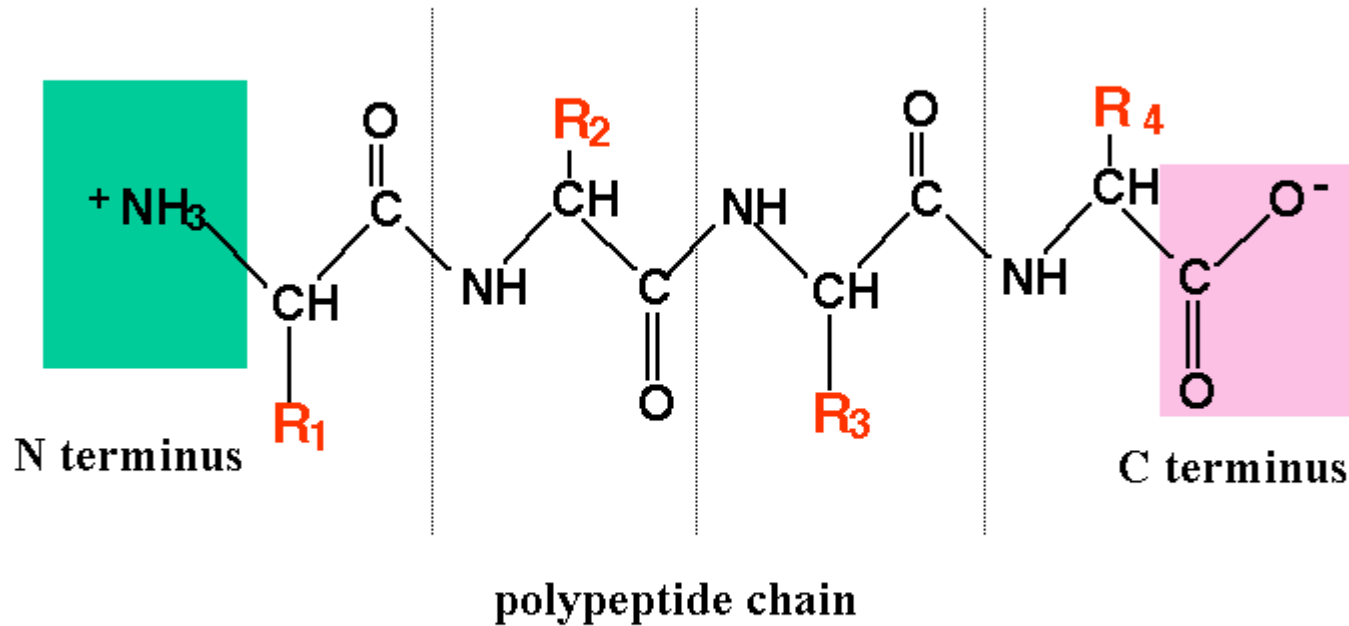
EFFICIENT CONFORMER LIBRARIES TO IMPROVE SIDECCHAIN OPTIMIZATION



CIBM Seminar - 20th September 2011
Sabareesh Subramaniam
Senes Lab, UW Biochemistry

PROTEINS

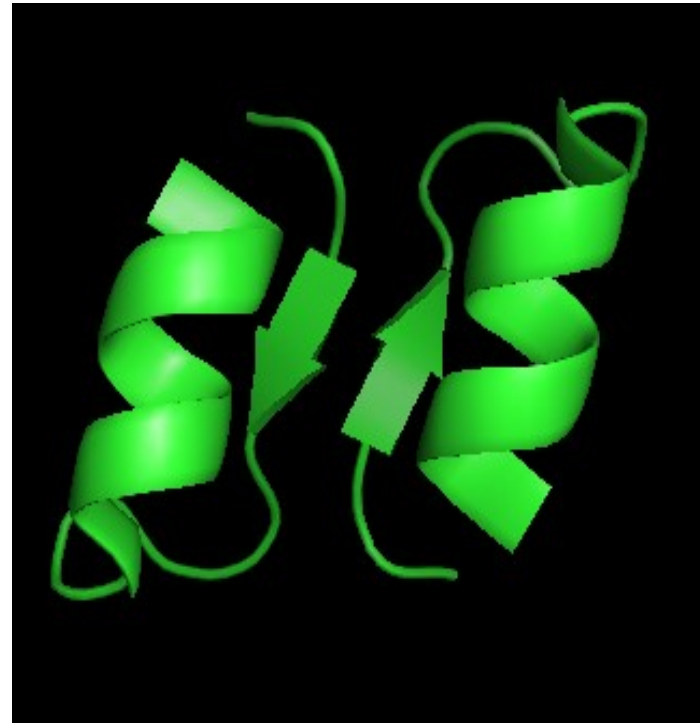
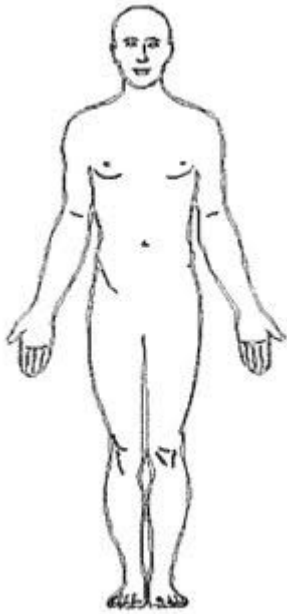
Peptide = chain of amino acids



PROTEINS

- Perform important biological functions
- Structure and function closely related
- Structures VERY helpful to study proteins
- X-Ray Crystallography, NMR to determine structure
- Computational modelling when structure not available

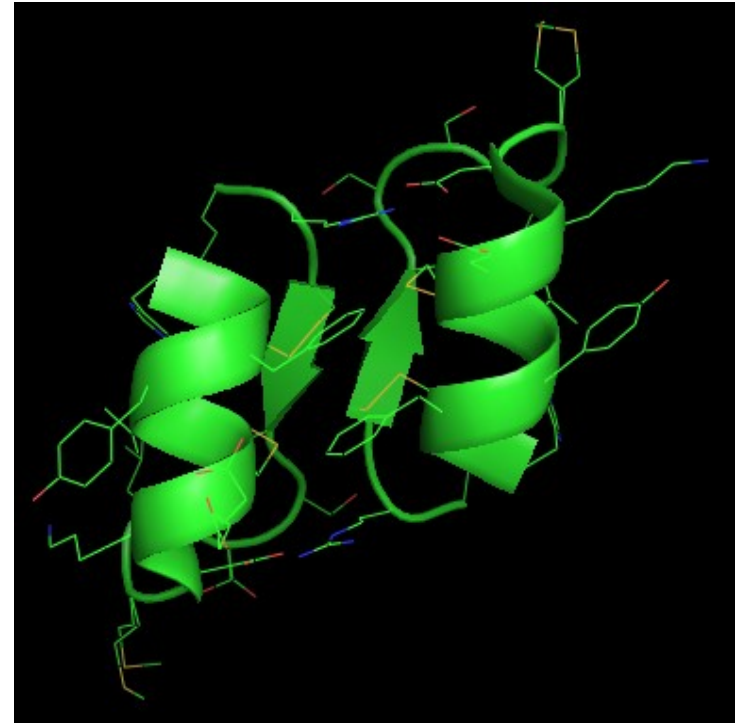
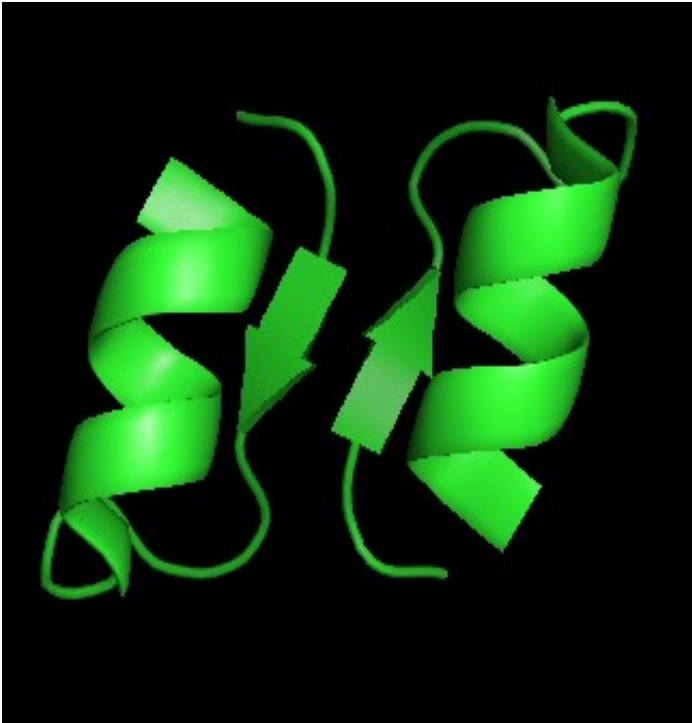
Structure Prediction via Homology Modeling



SIDECCHAIN OPTIMIZATION

Backbone

Add sidechains to achieve
minimum energy configuration

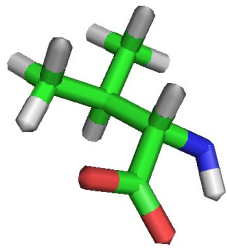


COMPUTATIONAL MODELING

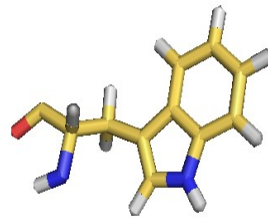
- Structure Prediction
 - Sequence → backbone(from homolog) → **Sidechain optimization**
- Protein Design
- Docking

SIDECCHAIN OPTIMIZATION

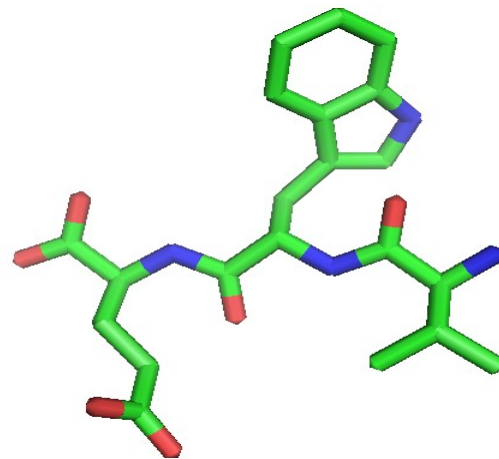
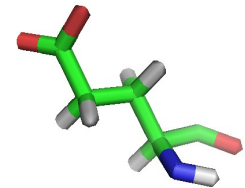
GLU



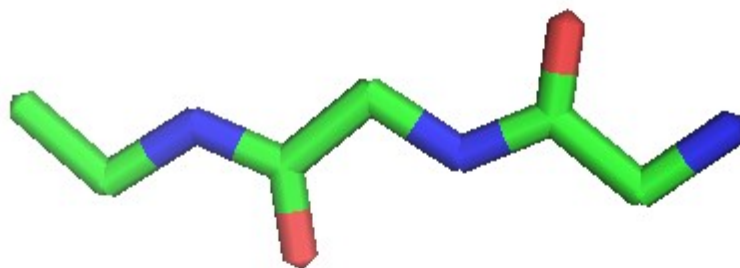
TRP



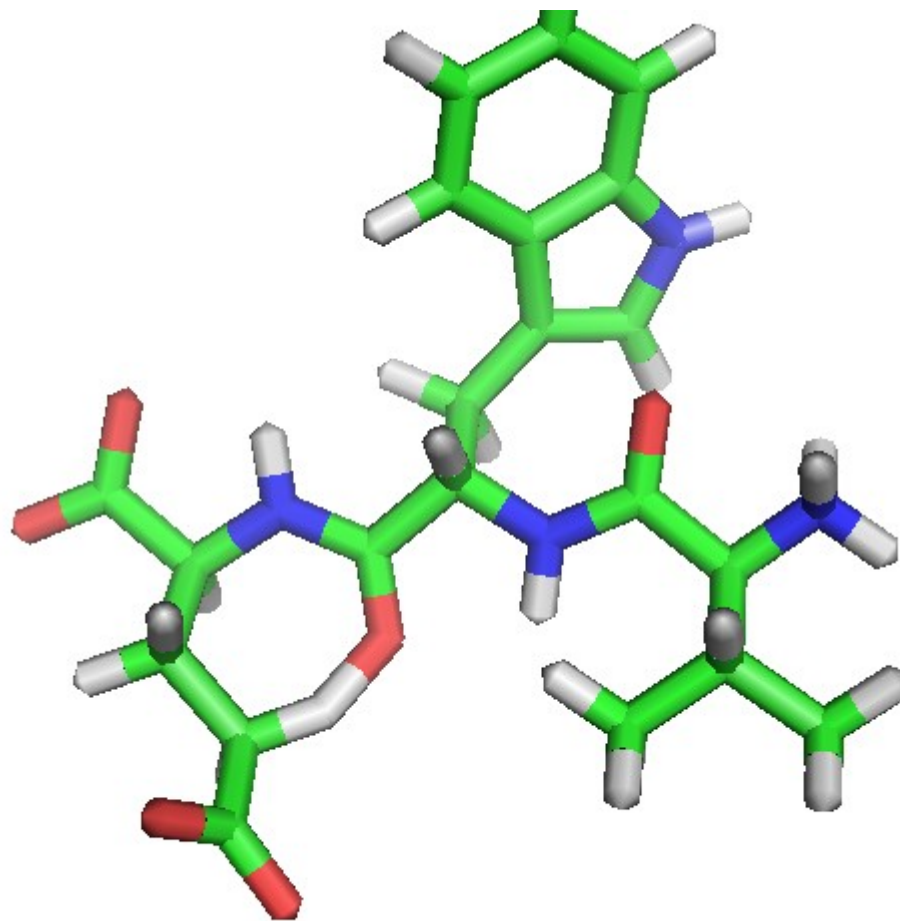
VAL



SIDECCHAIN OPTIMIZATION

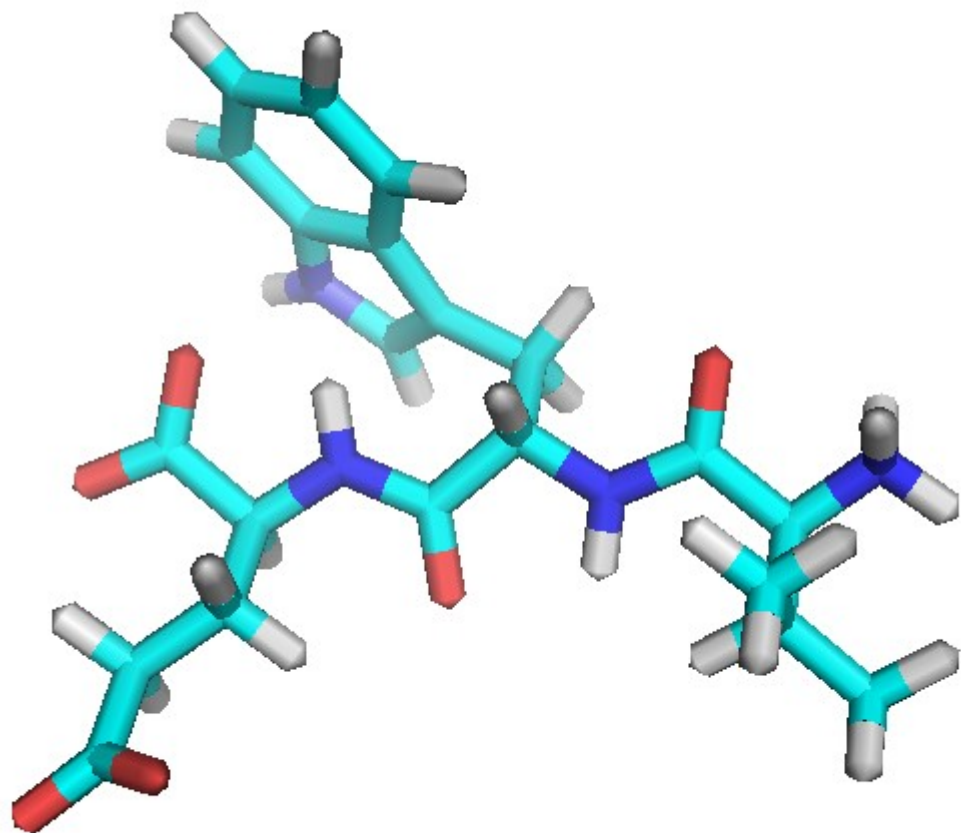


SIDCHAIN OPTIMIZATION



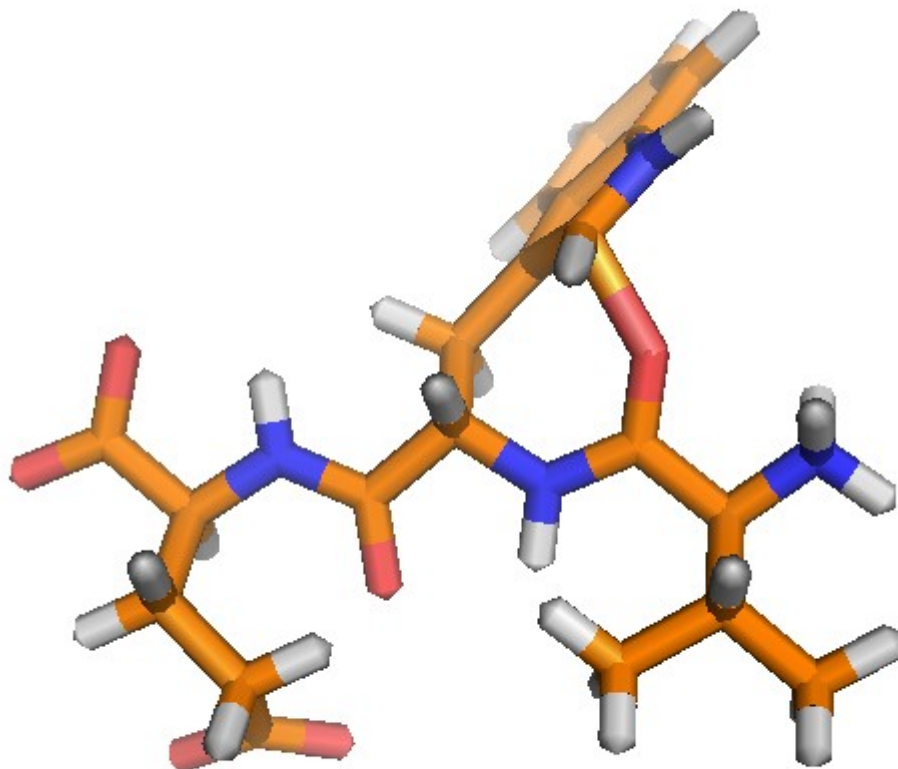
Energy = 40000 Kcal / Mol

SIDECCHAIN OPTIMIZATION



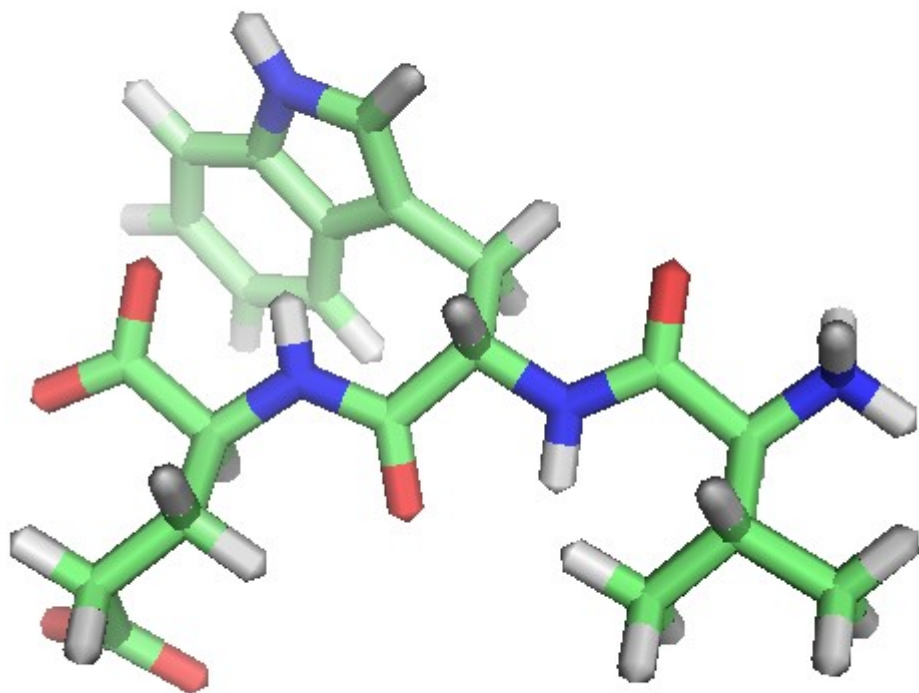
Energy = 60 Kcal / Mol

SIDECCHAIN OPTIMIZATION



Energy = 30000 Kcal / Mol

SIDCHAIN OPTIMIZATION



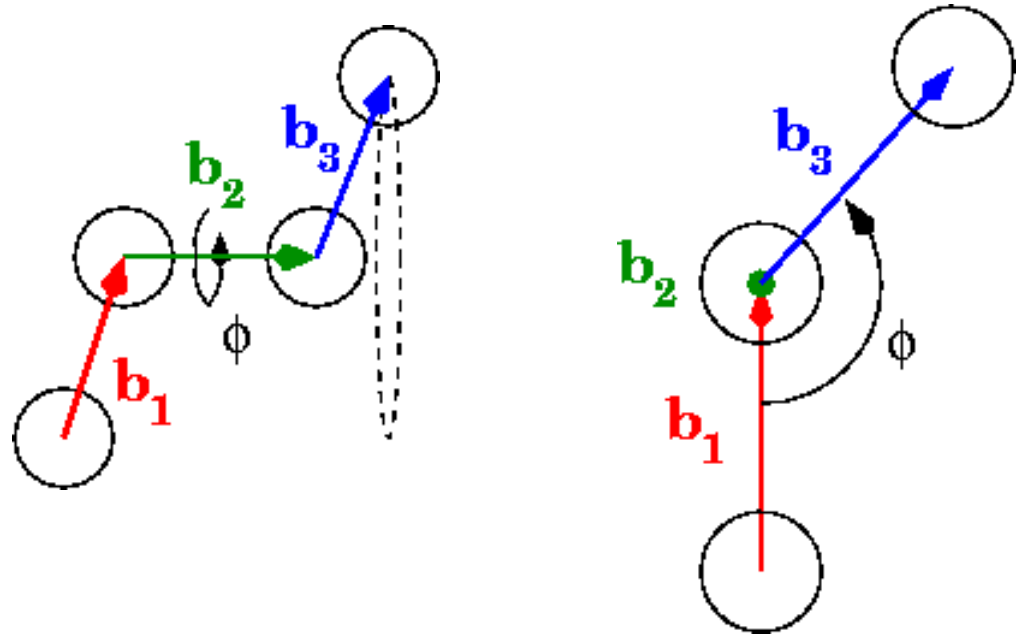
Energy = 70 Kcal / Mol

DEGREES OF FREEDOM

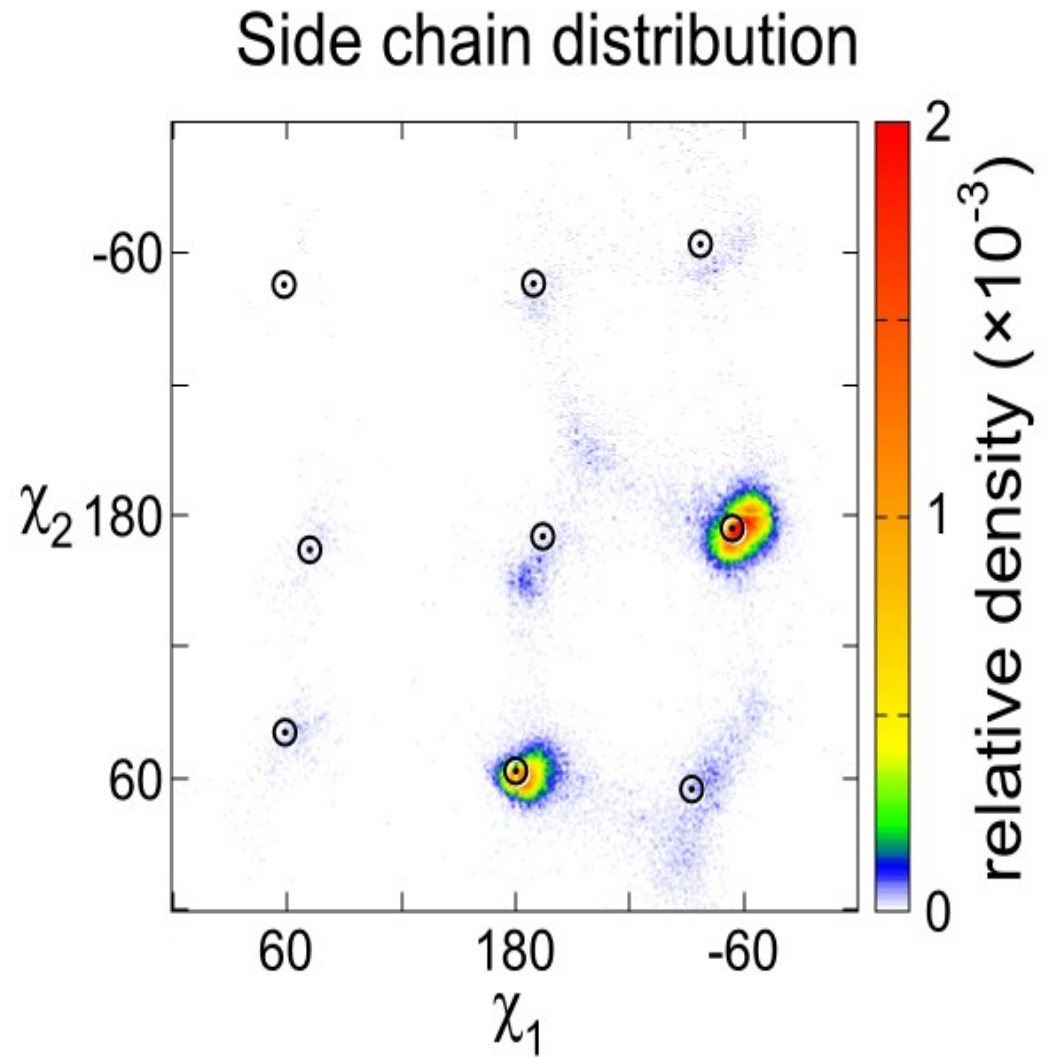
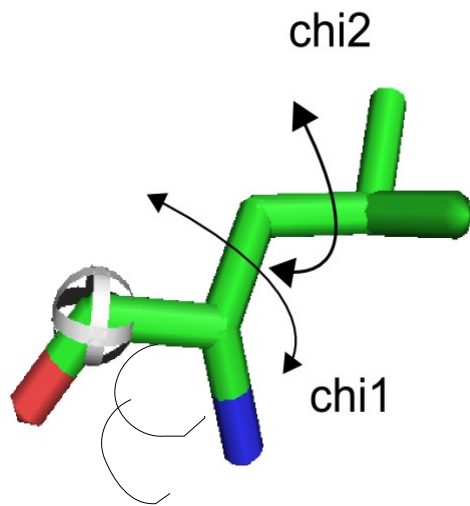
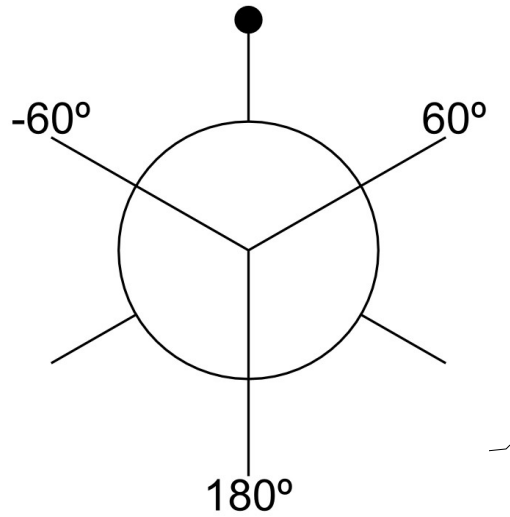
Bond distances

Bond Angles

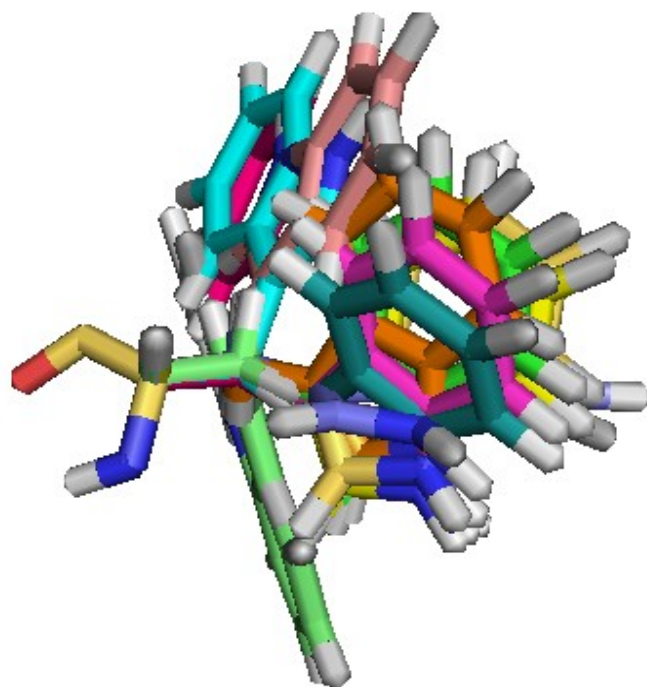
Dihedral or torsional angles



STATISTICS OF DIHEDRAL ANGLES



DISCRETIZED CONFORMATION LIBRARIES

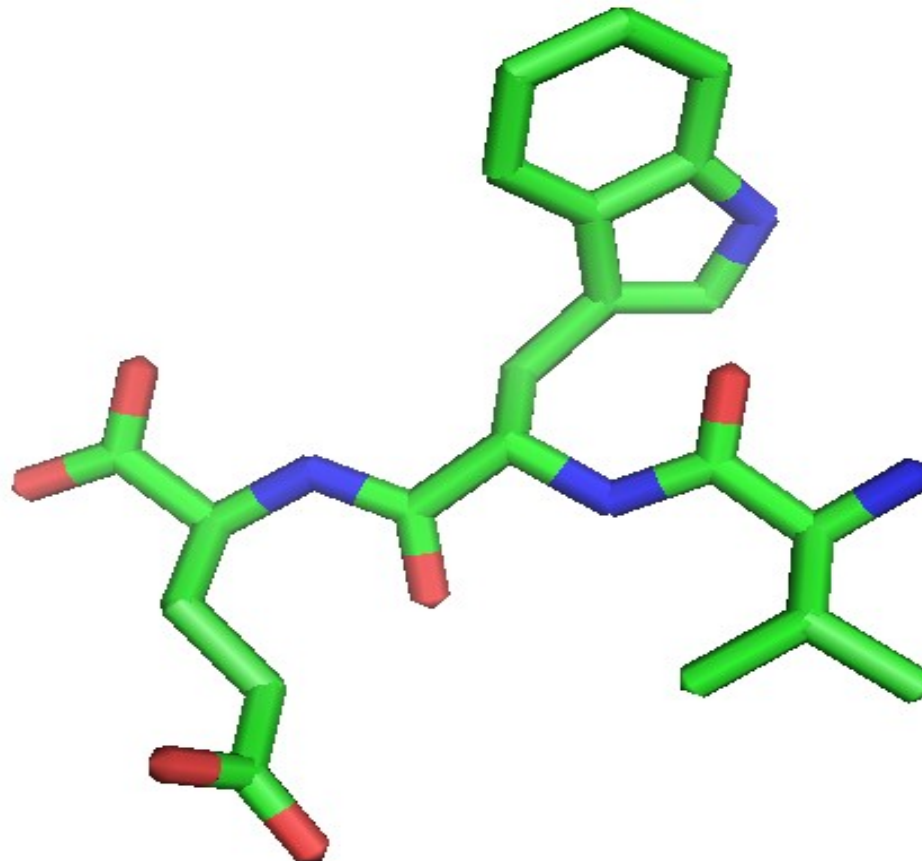


SIDCHAIN OPTIMIZATION

GLU

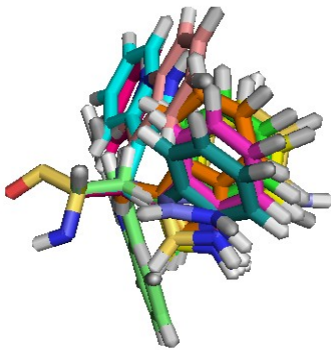
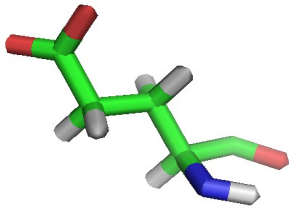
TRP

VAL



COMBINATORIAL SEARCH SPACE (3-D JIGSAW PUZZLE)

GLU

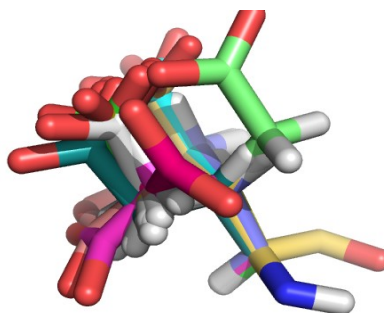
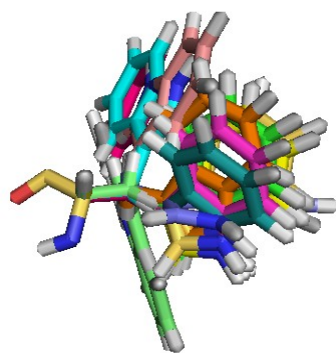
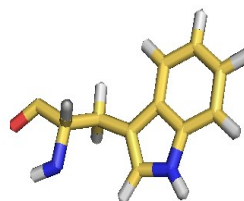
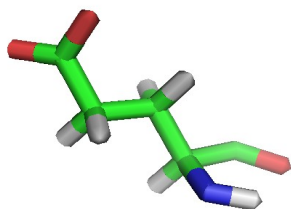


No of conformations to search 36

COMBINATORIAL SEARCH SPACE (3-D JIGSAW PUZZLE)

GLU

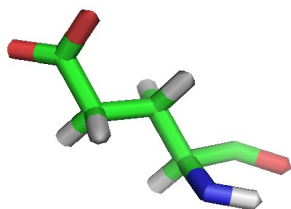
TRP



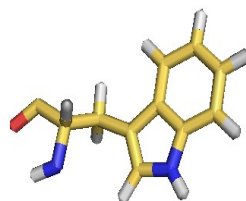
No of conformations to search $36 * 54$

COMBINATORIAL SEARCH SPACE (3-D JIGSAW PUZZLE)

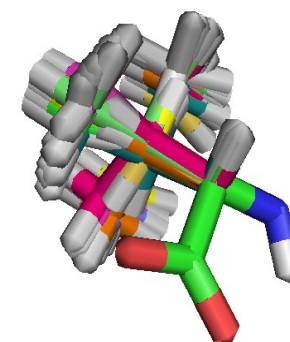
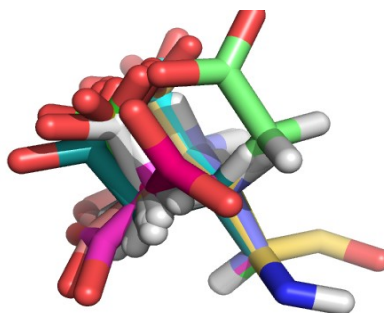
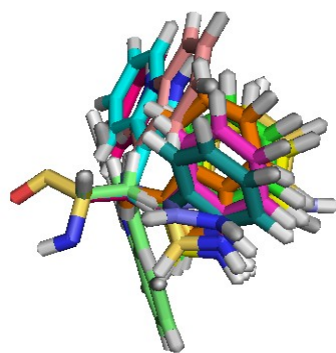
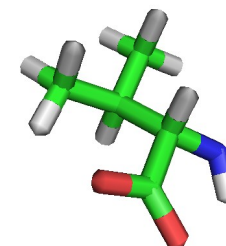
GLU



TRP



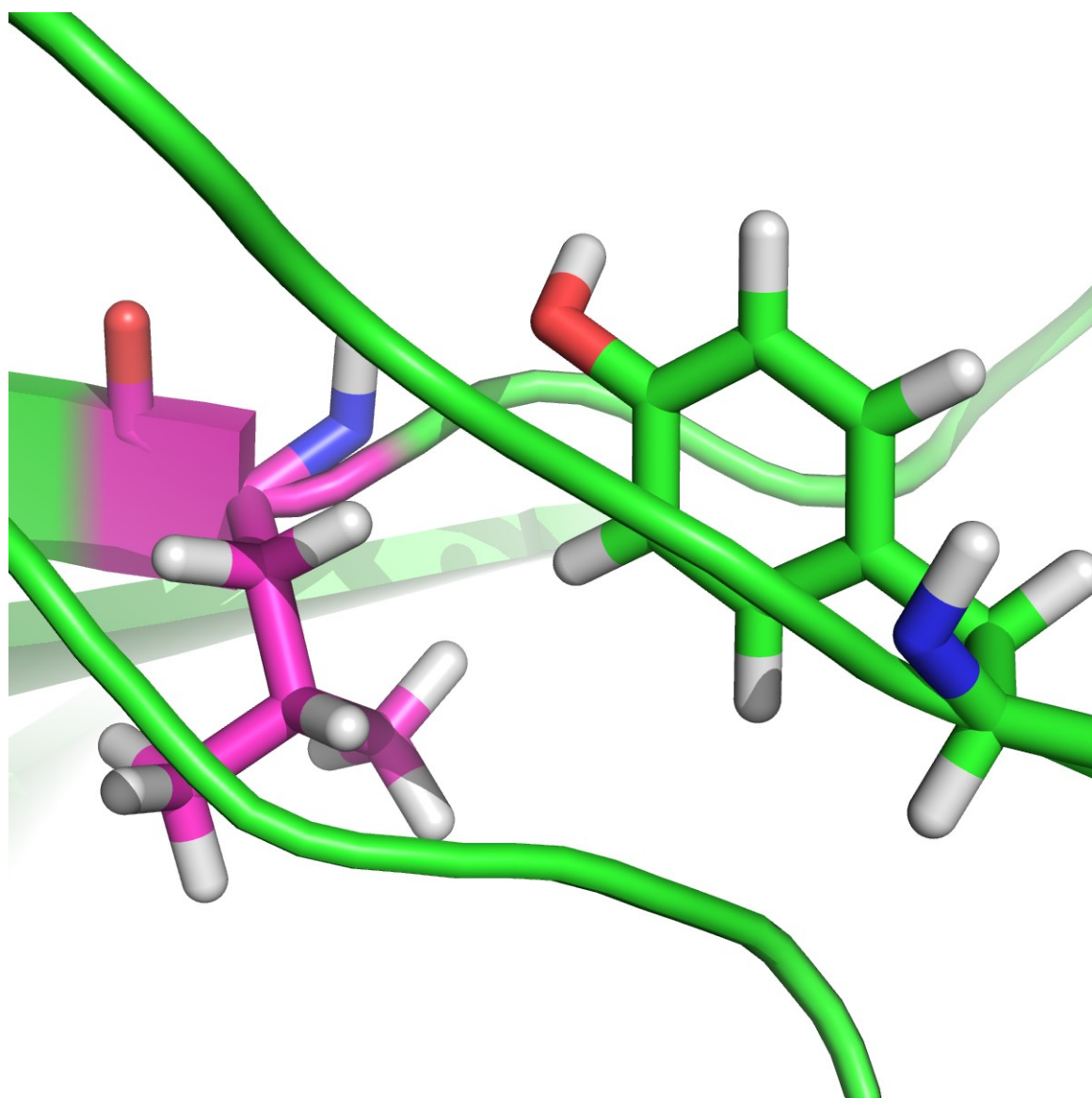
VAL

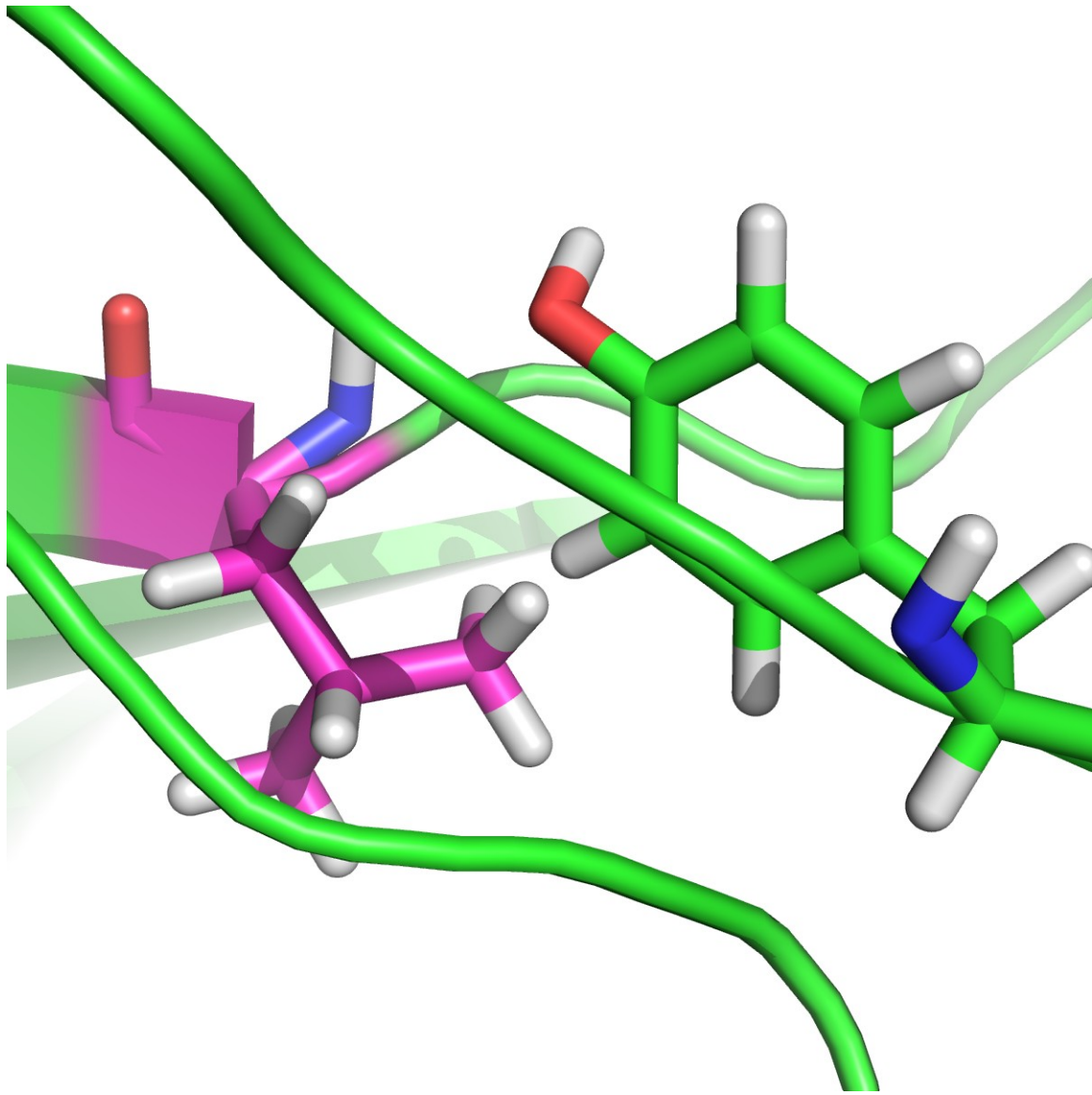


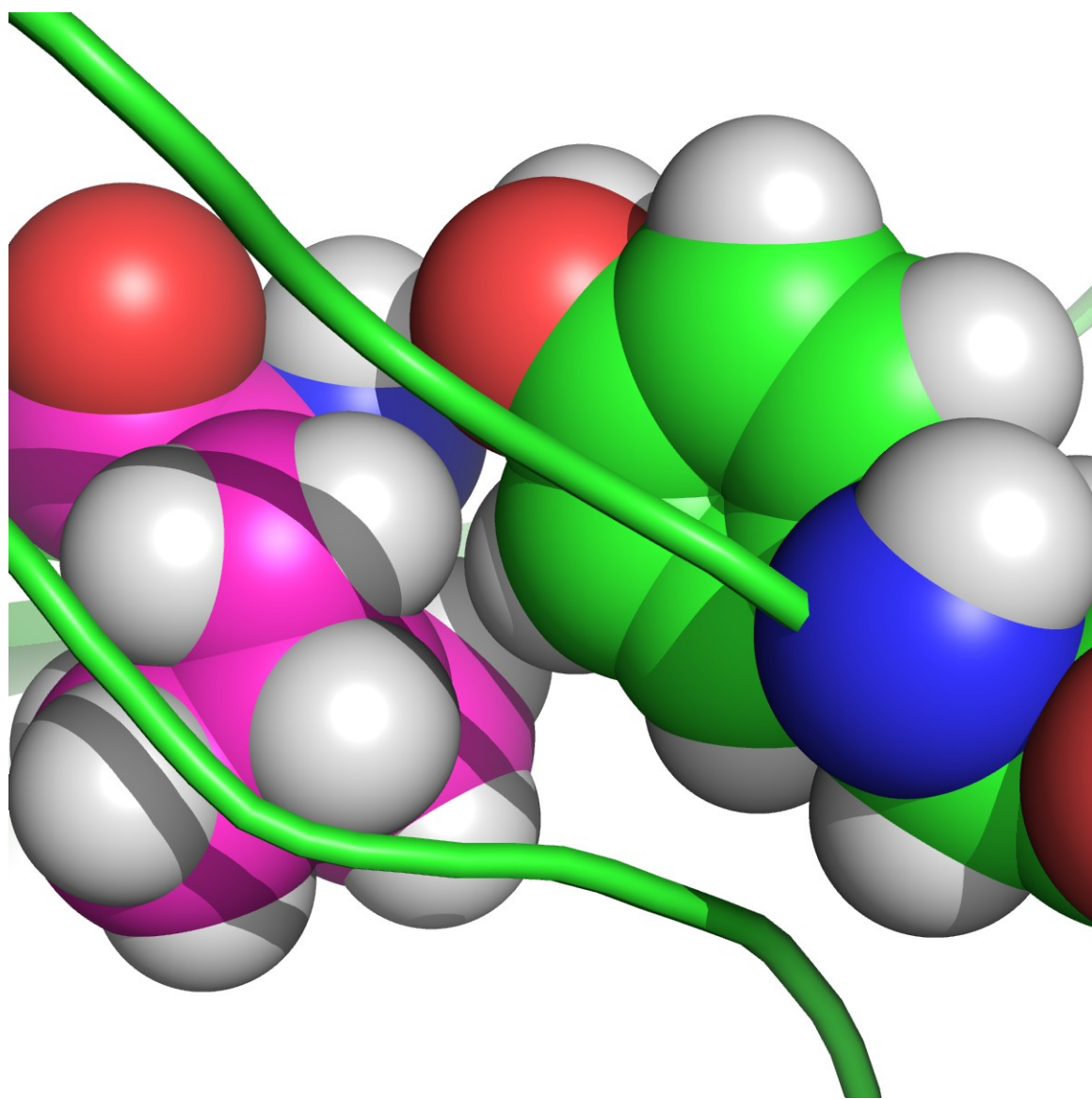
No of conformations to search $36 * 54 * 3 = 5832$
Typically $> 10^{60}$

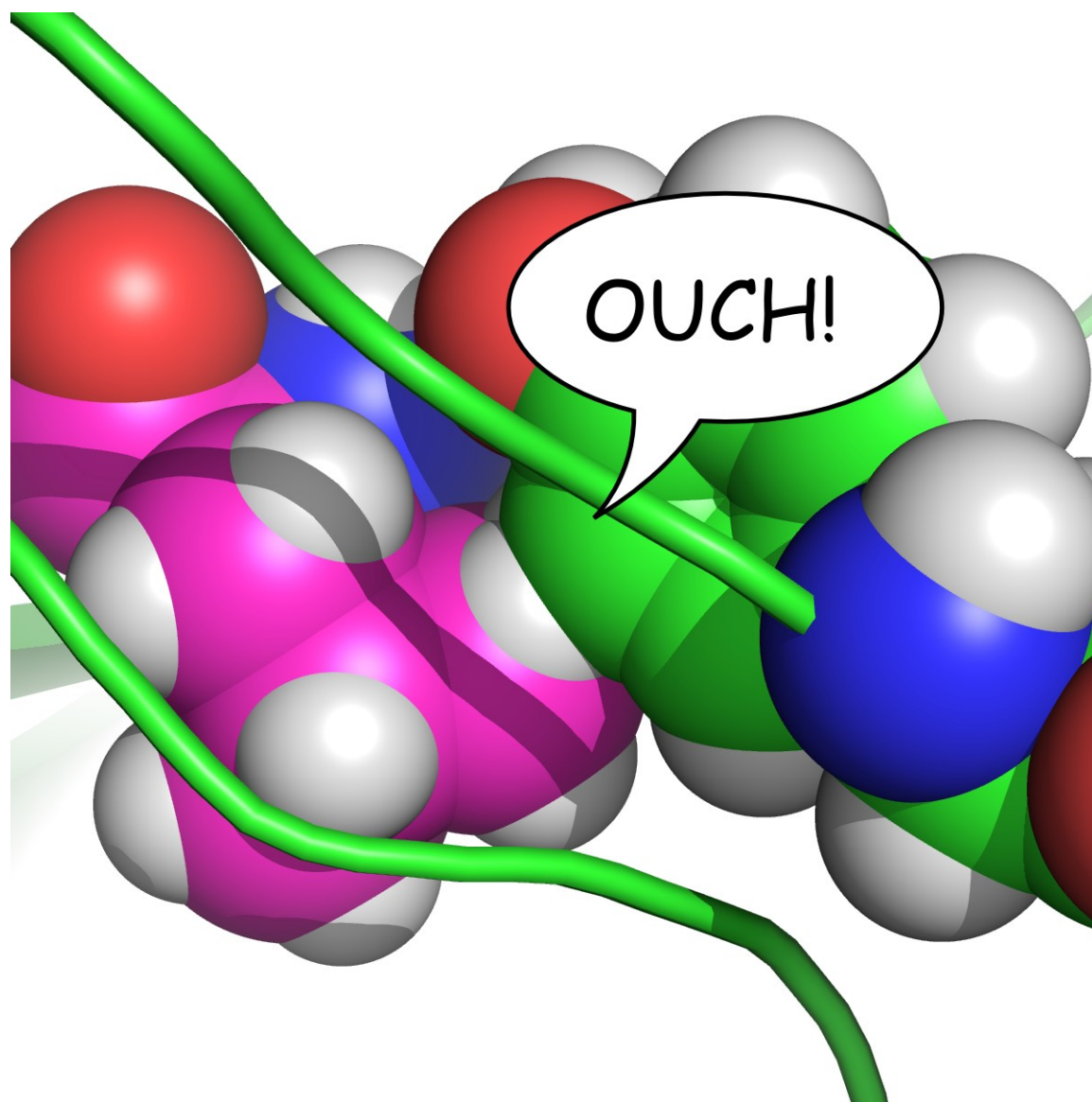
SEARCH ALGORITHMS

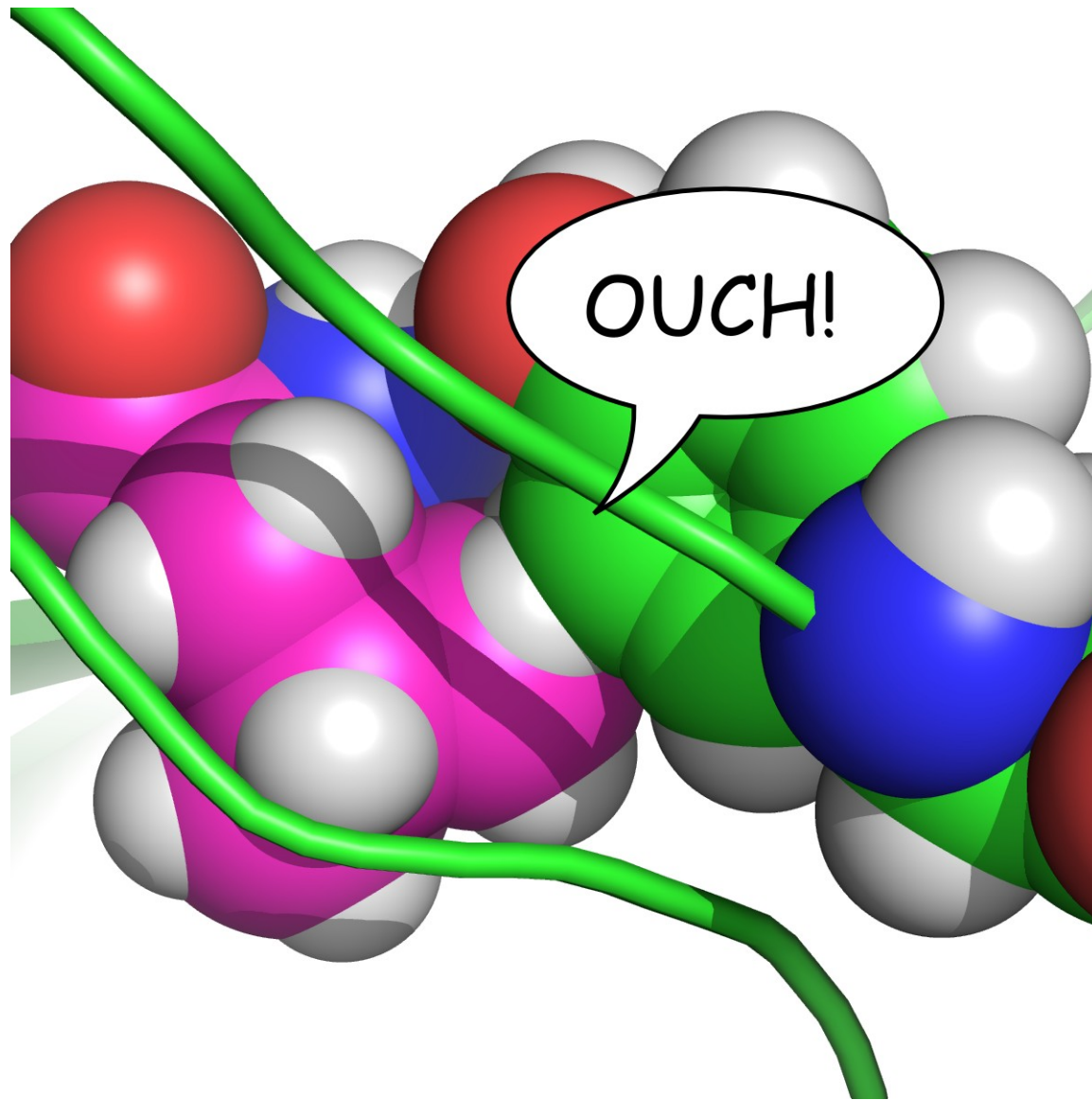
- Dead End Elimination
- Self Consistent Mean Field
- A* search
- Monte Carlo Simulated Annealing
- Graph decomposition







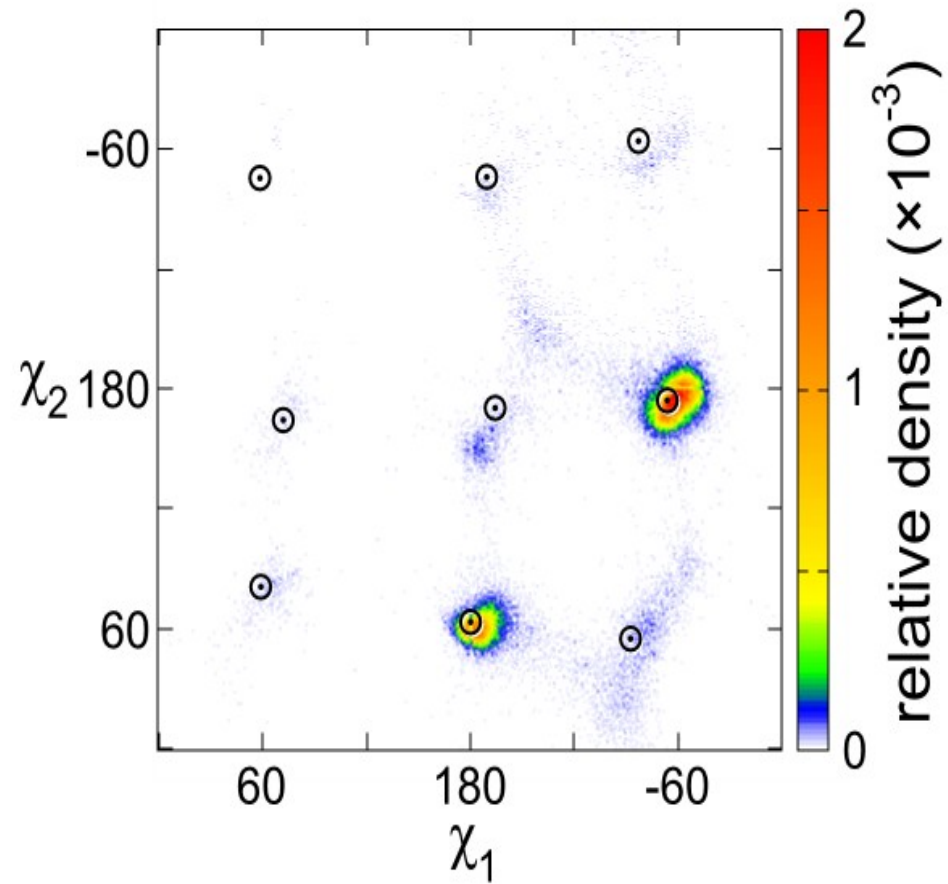




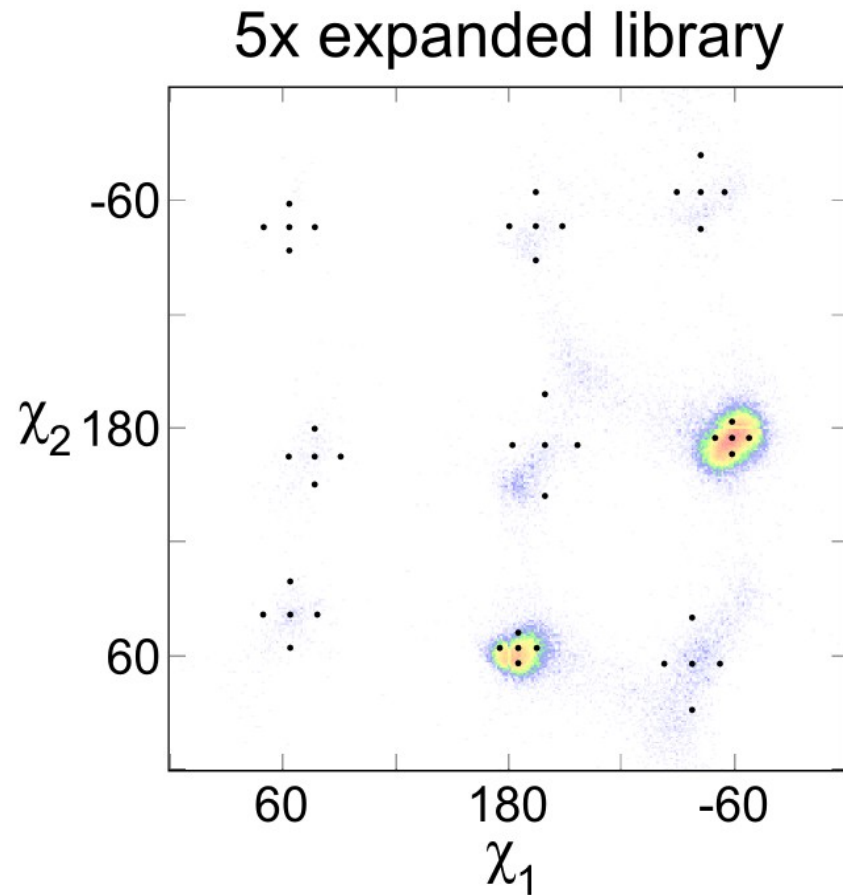
We need more sampling

MORE SAMPLING

Side chain distribution

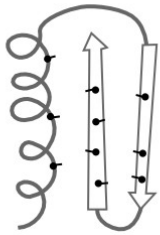


MORE SAMPLING

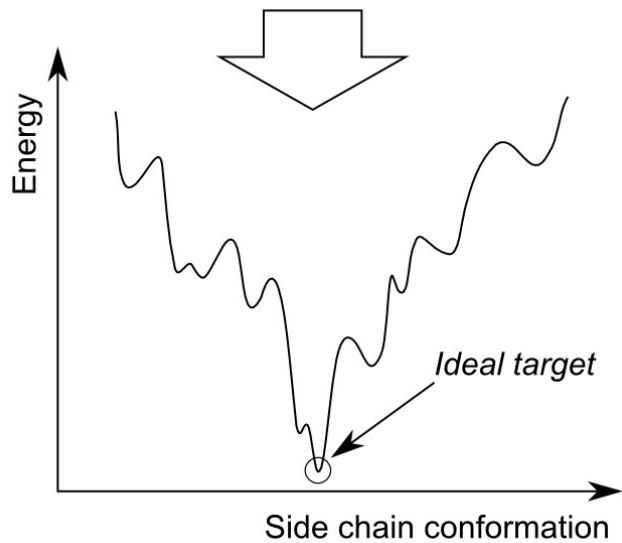


ROTAMER LIBRARY DETERMINES QUALITY OF SOLUTION

a

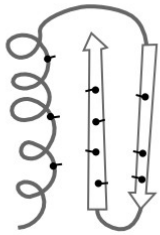


The template and the energy functions define a continuum energy landscape in side chain conformational space

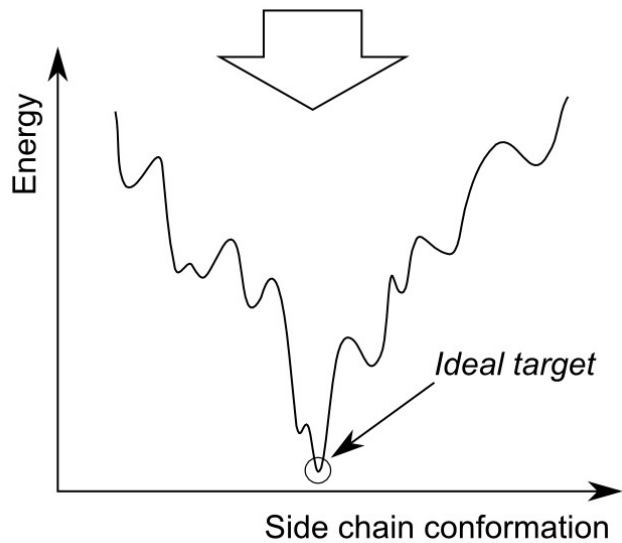


ROTAMER LIBRARY DETERMINES QUALITY OF SOLUTION

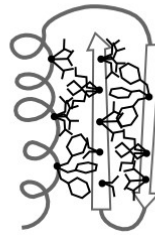
a



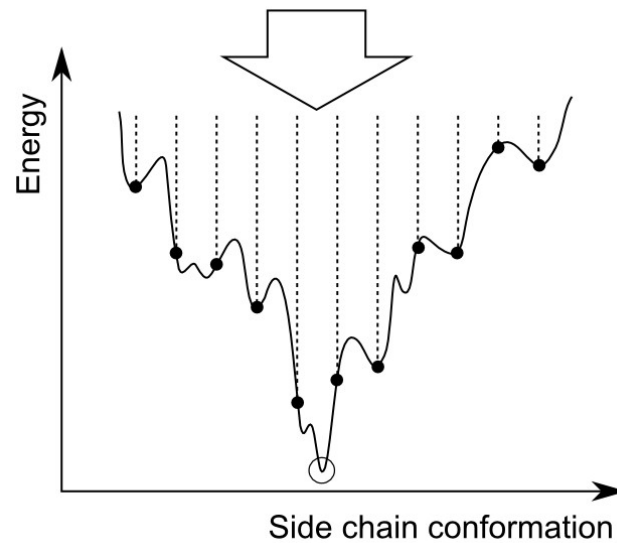
The template and the energy functions define a continuum energy landscape in side chain conformational space



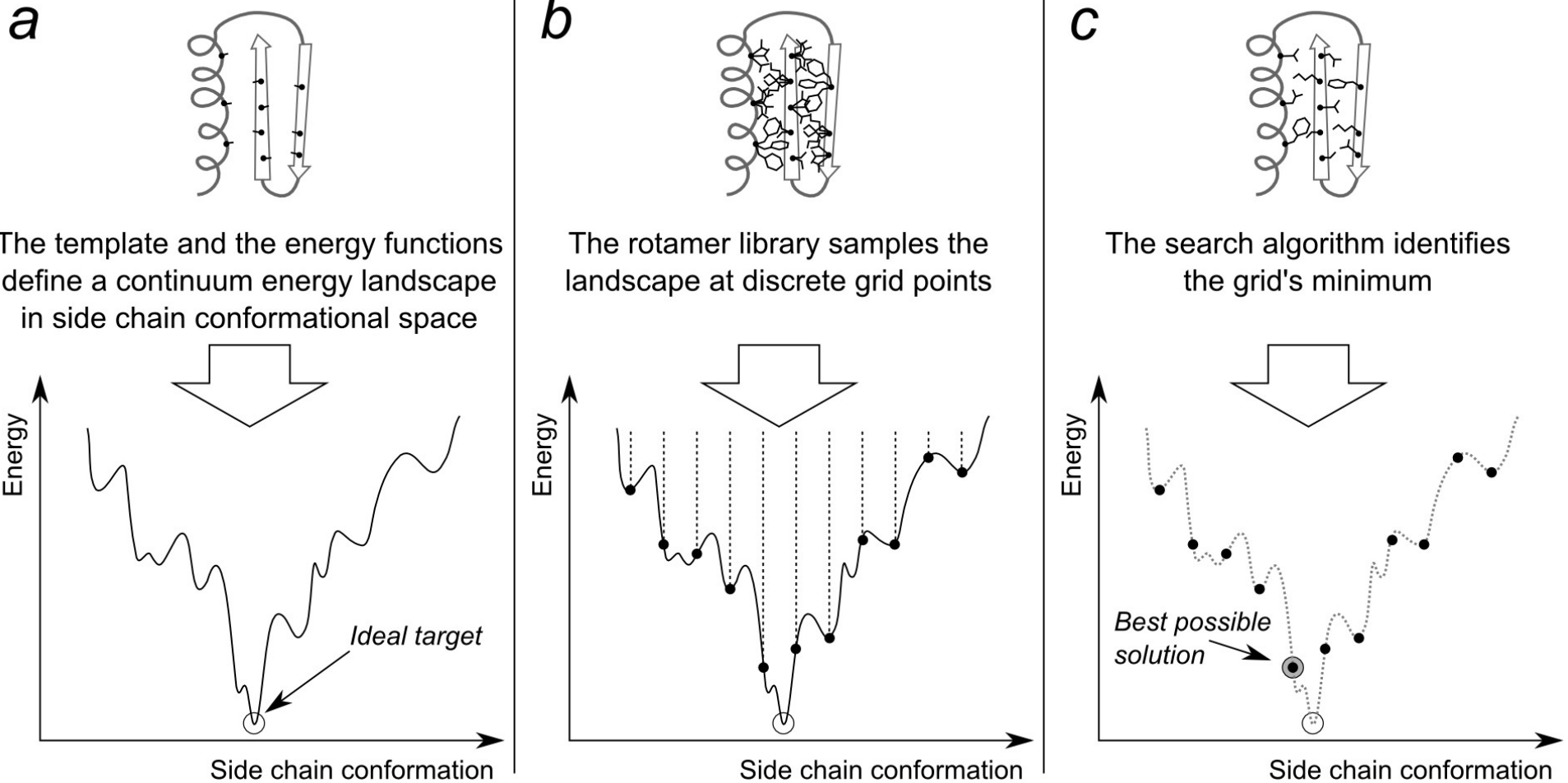
b



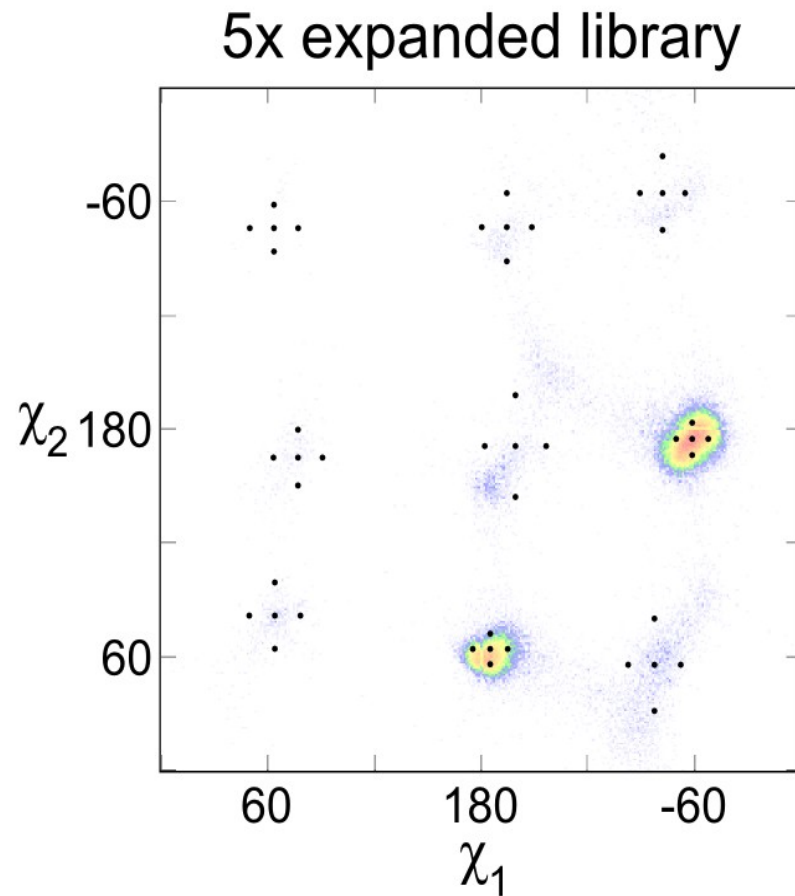
The rotamer library samples the landscape at discrete grid points



ROTAMER LIBRARY DETERMINES QUALITY OF SOLUTION



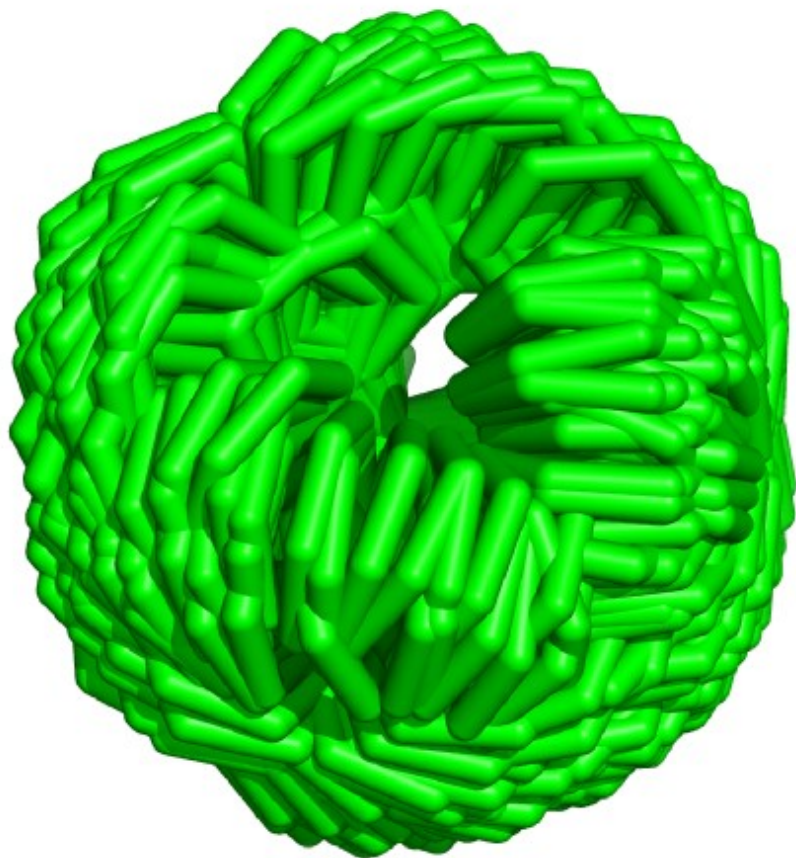
FIXED NUMBER OF CONFORMERS



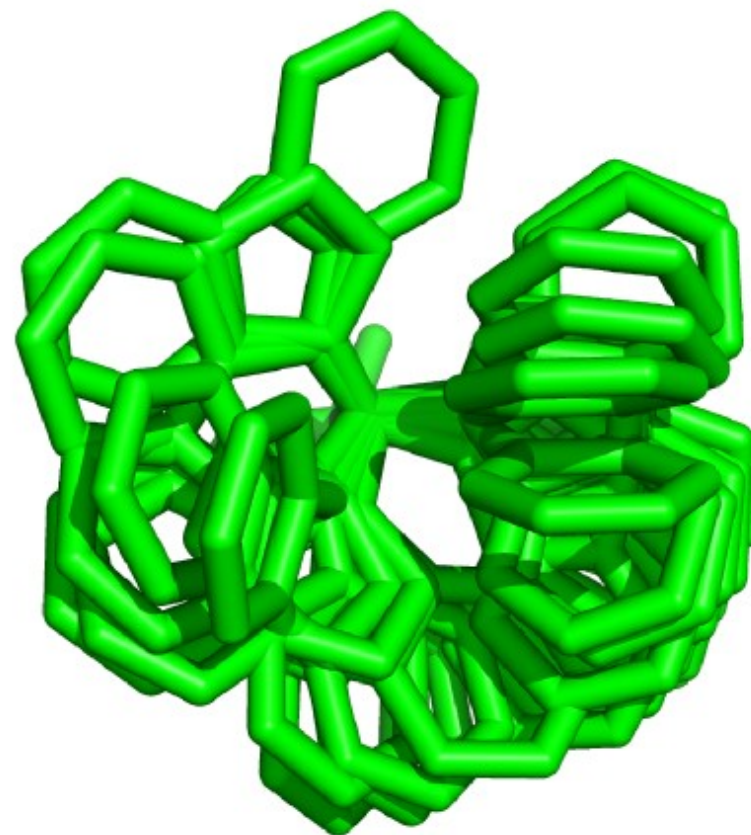
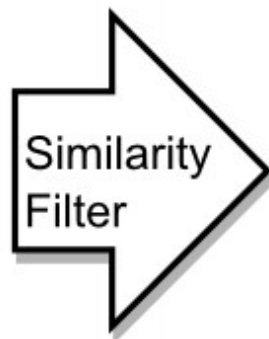
9 \rightarrow 45 rotamers

**No number in
between**

GEOMETRIC FILTERS LEAD TO CONFORMER LIBRARIES



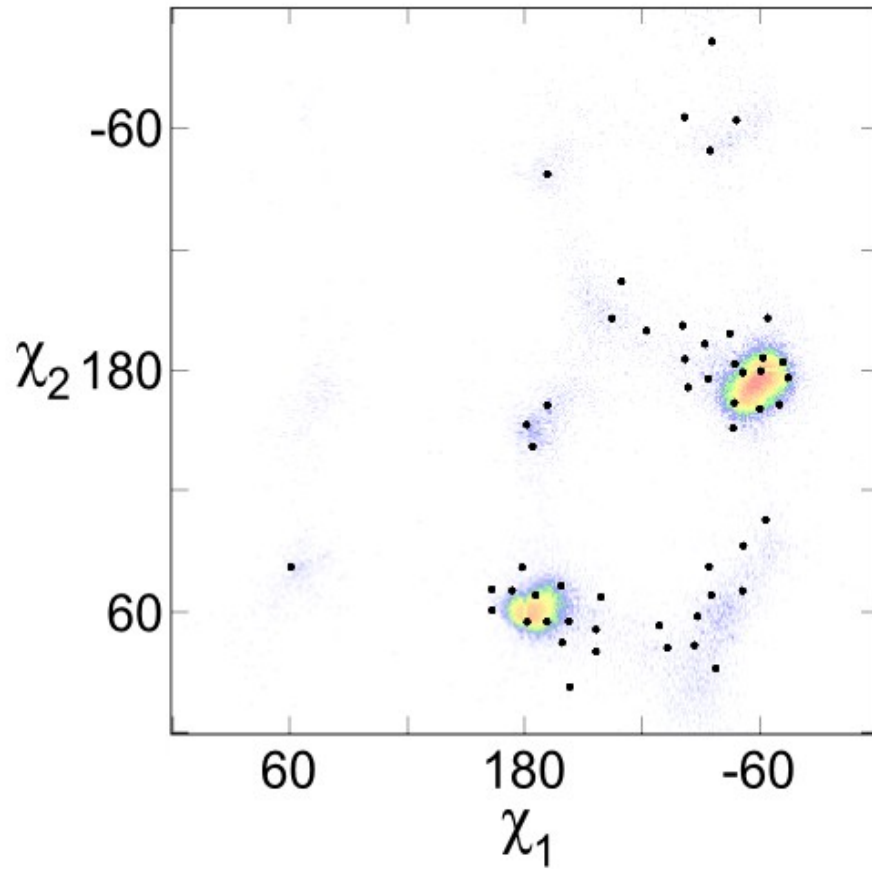
Conformers from high-res PDBs



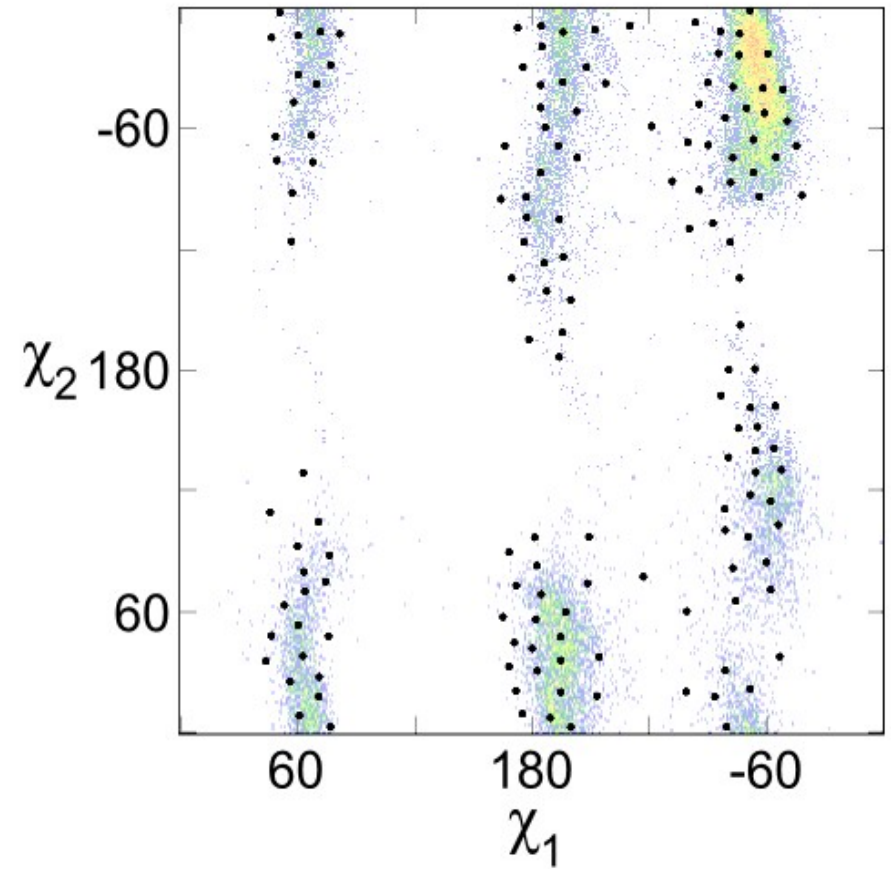
Representative conformers

IGNORES THE NATURAL DISTRIBUTION

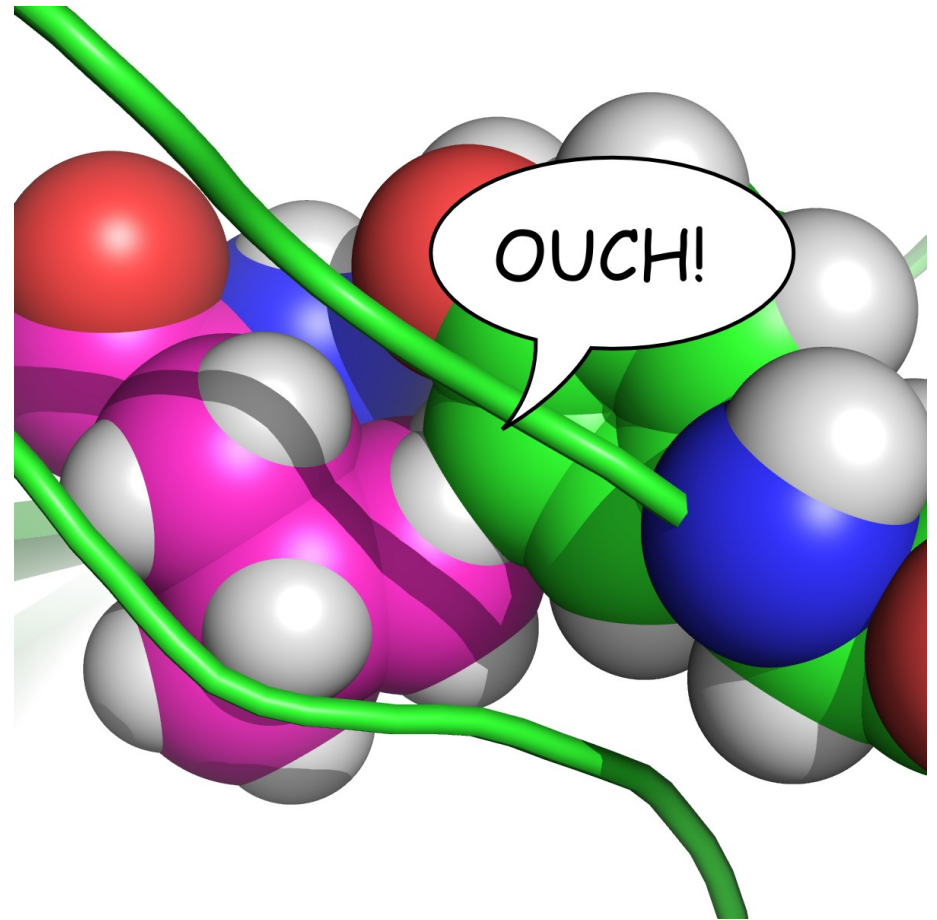
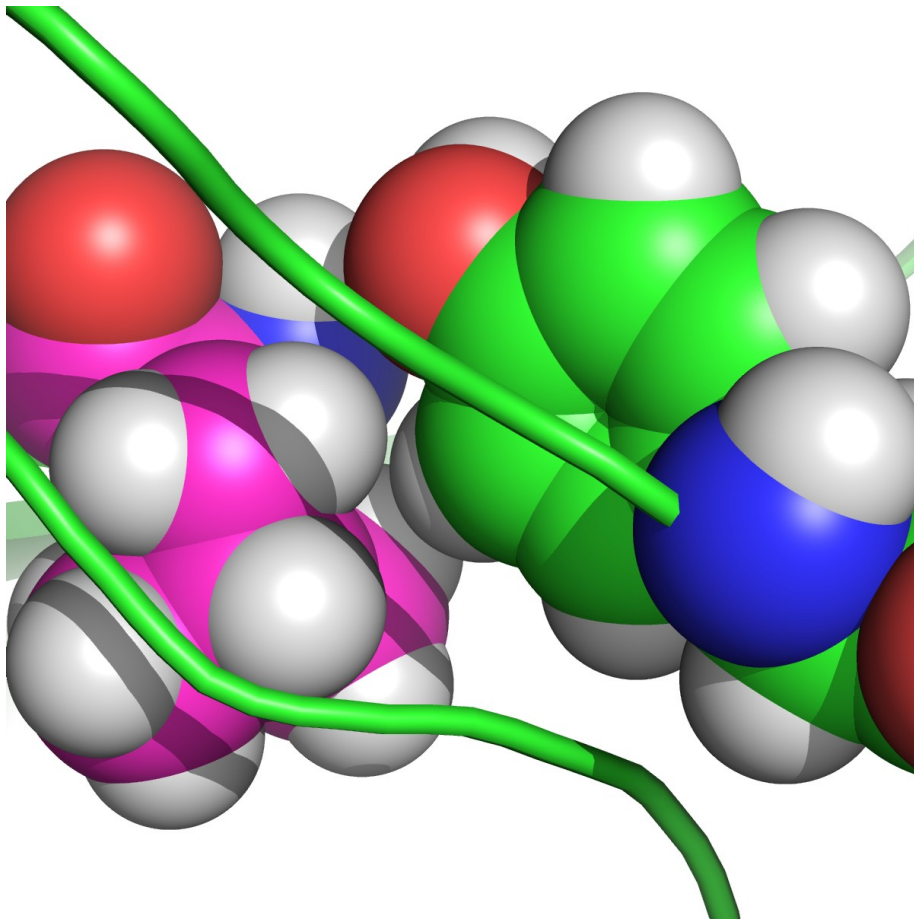
LEU



ASN



- People have been looking for solutions using the statistical distribution in structures
- However, the problem is that sampling is related to the energetics in a way that is difficult to predict



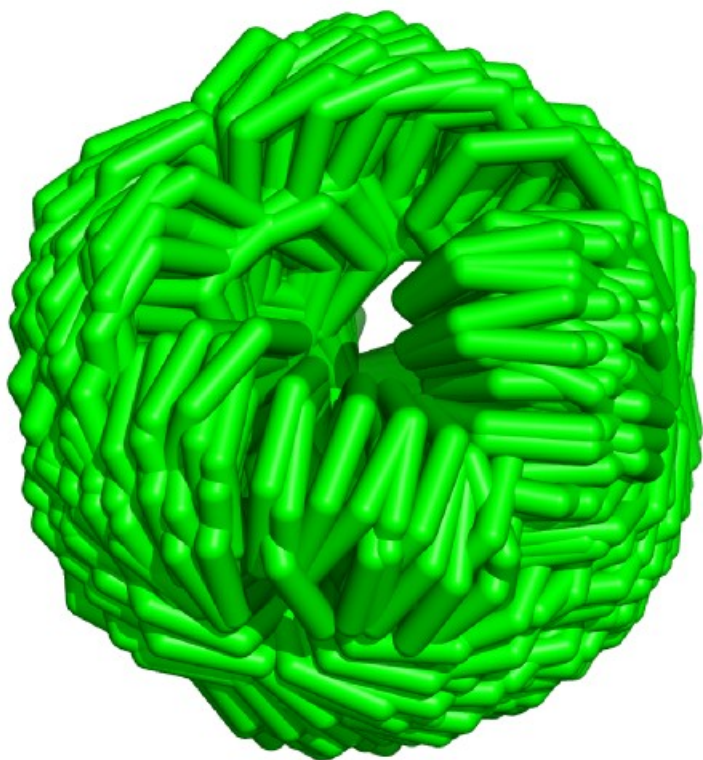
- People have been looking for solutions using the statistical distribution in structures
- However, the problem is that sampling is related to the energies in a way that is difficult to predict
- Solution: use energetics to identify the best sampling strategy for side chain optimization

GOALS

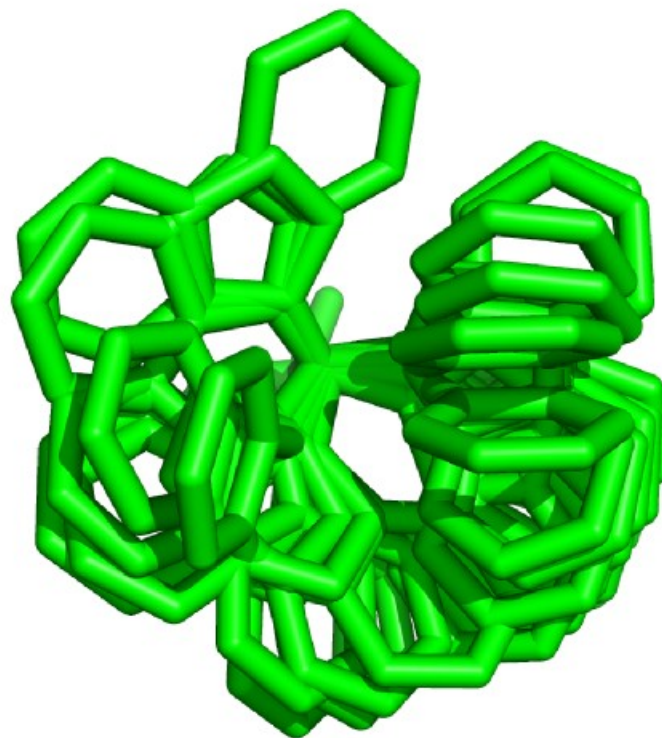
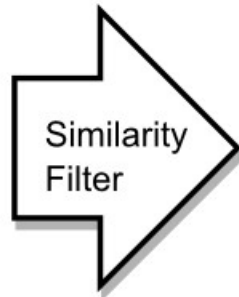
Can we create a library that can outperform existing libraries in terms of speed and/or Energies ?

Can we create a flexible library where the conformers are in some useful order?

Can we sort the conformers instead of extracting a fixed-size subset?

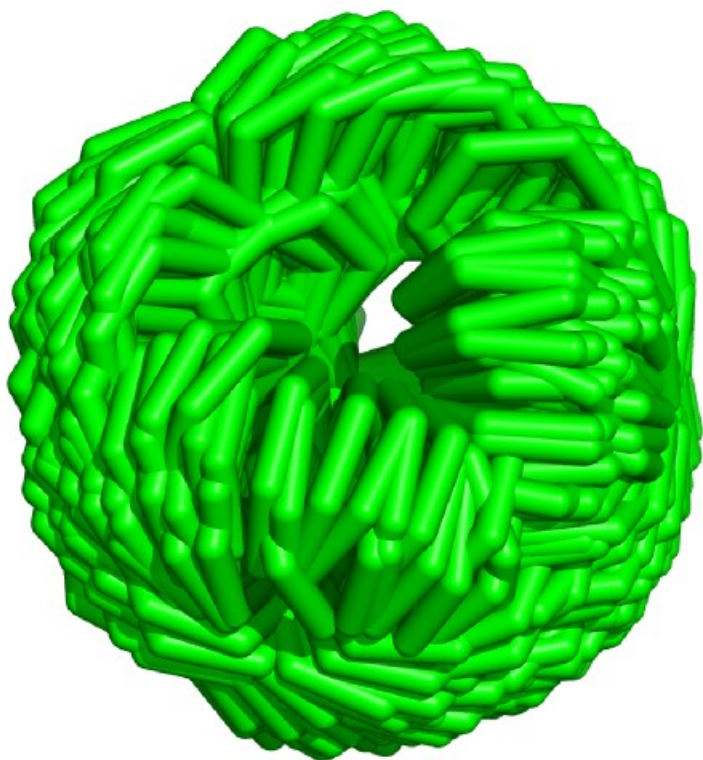


5,000 conformers from high-res PDBs

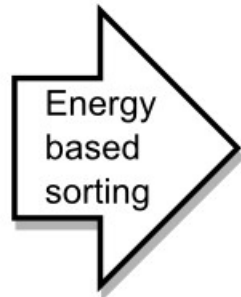


50 representative conformers

Can we sort the conformers instead of extracting a fixed-size subset?



5,000 conformers from high-res PDBs

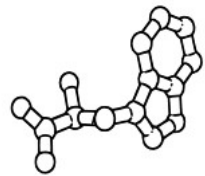


1.	_____
2.	_____
3.	_____
4.	_____
5.	_____
6.	_____
7.	_____
8.	_____
9.	_____
10.	_____
.....	_____
.....	_____
4999.	_____
5000.	_____

Sorted list of 5,000 conformers

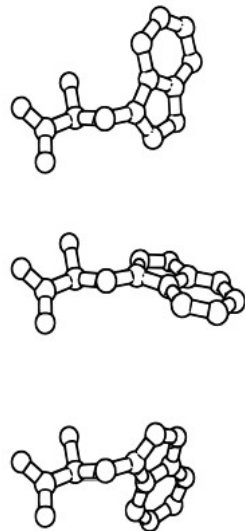
- Use energetics to identify the best sampling strategy for side chain optimization

Conformers

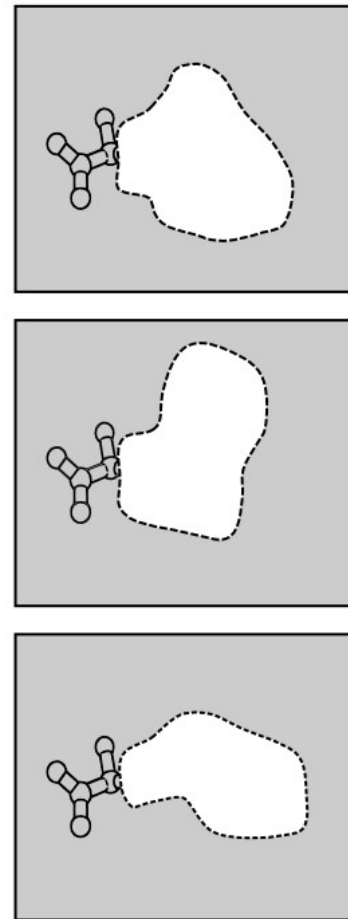


- Use energetics to identify the best sampling strategy for side chain optimization

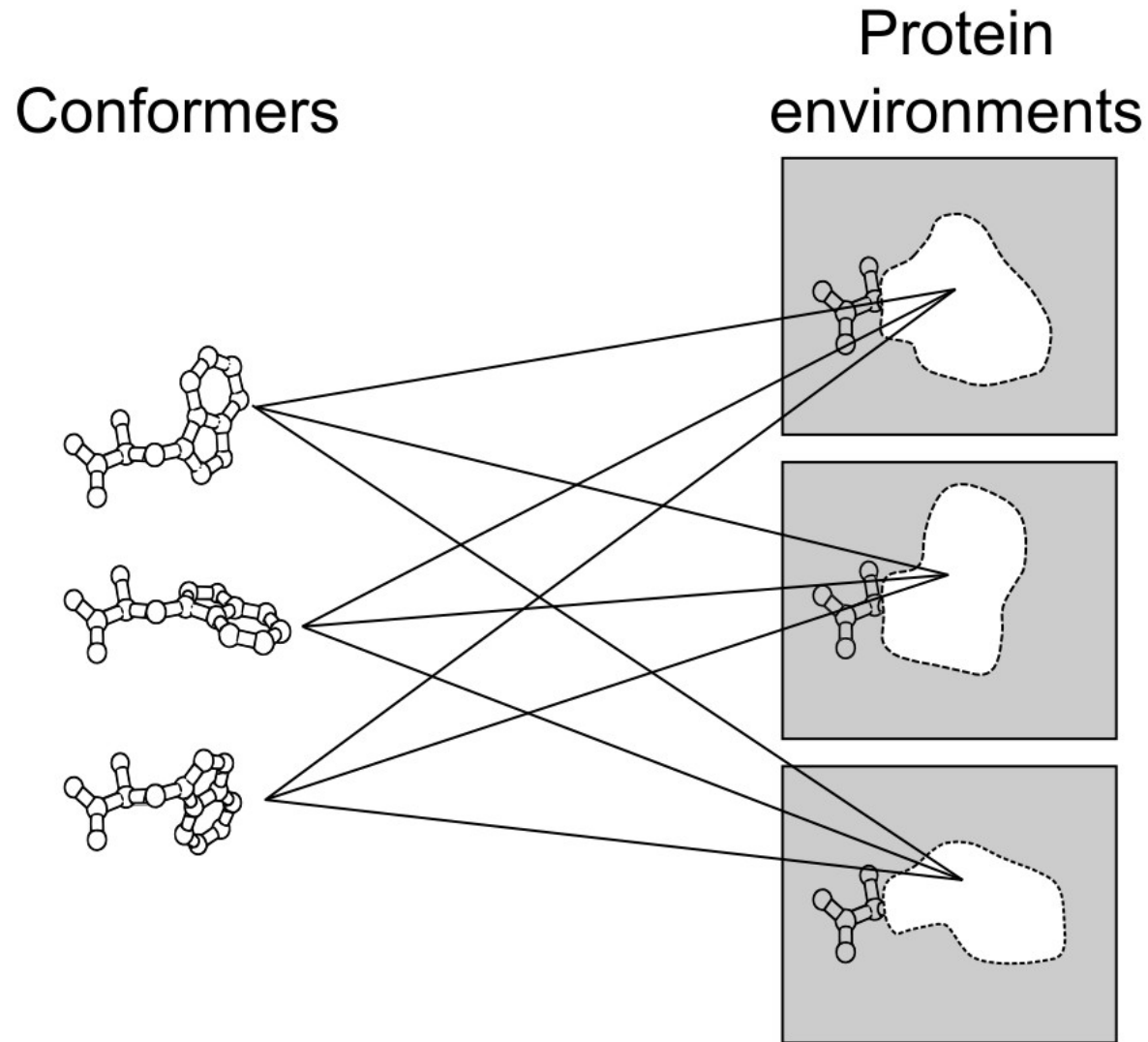
Conformers



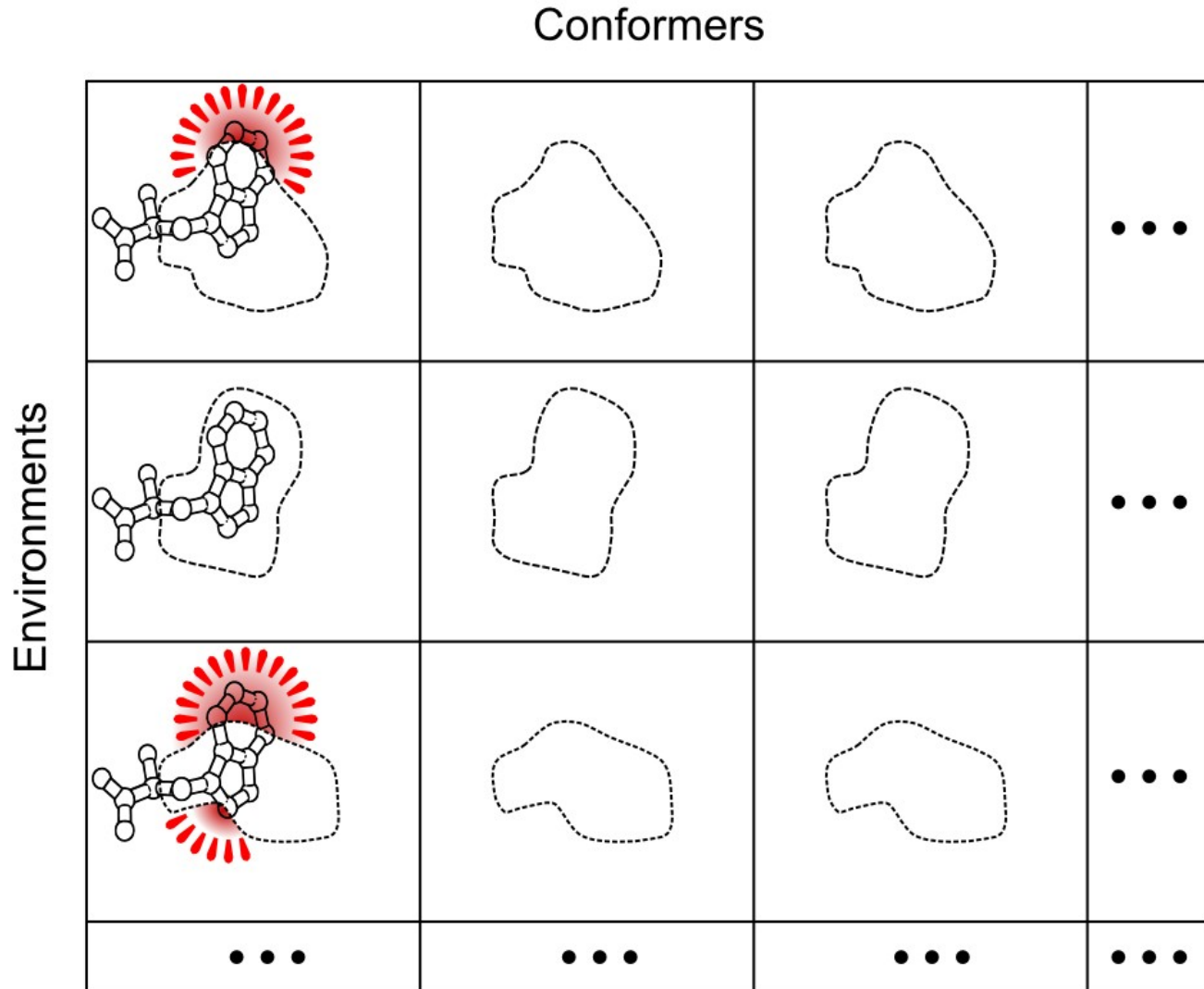
Protein environments



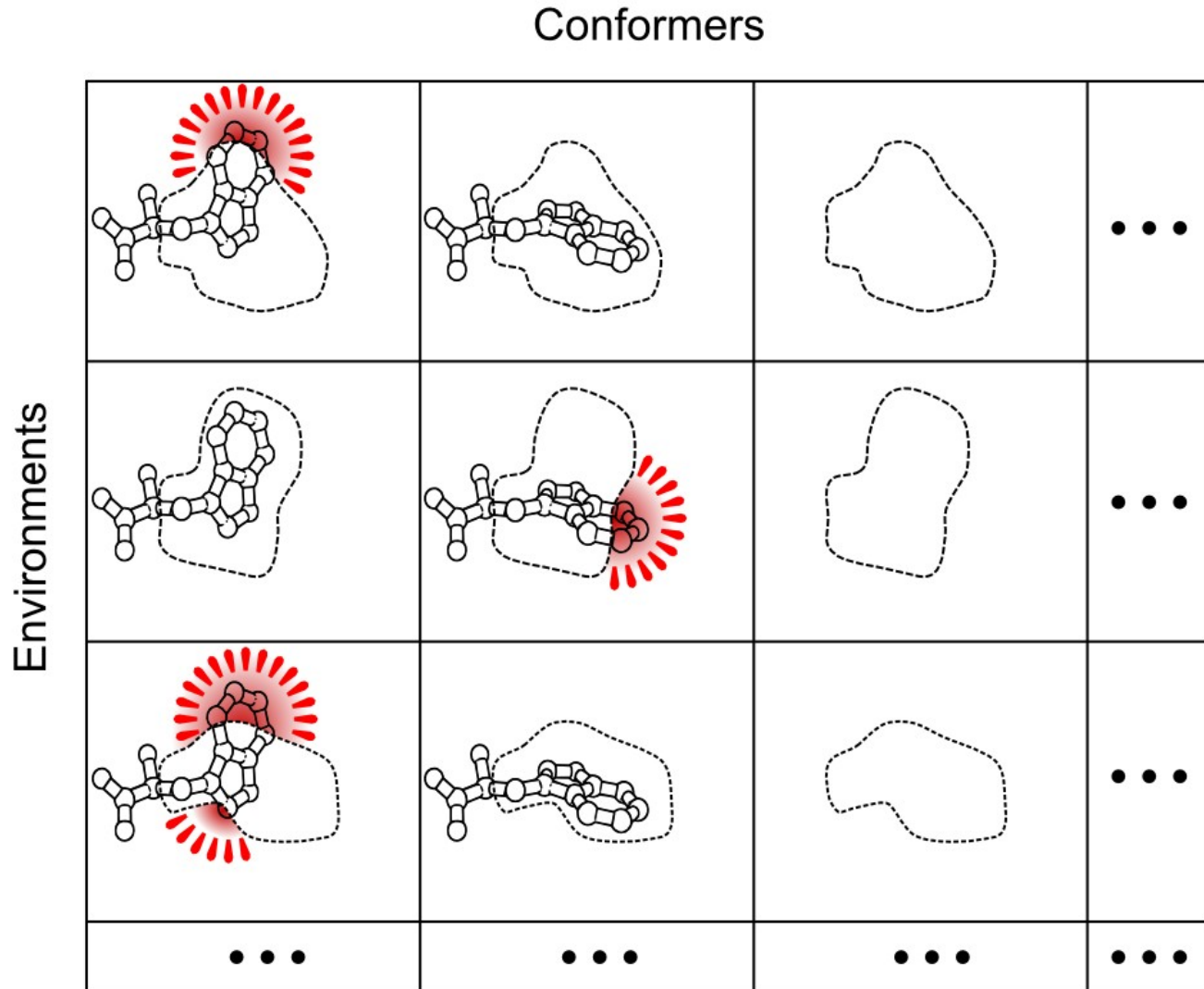
- Use energetics to identify the best sampling strategy for side chain optimization



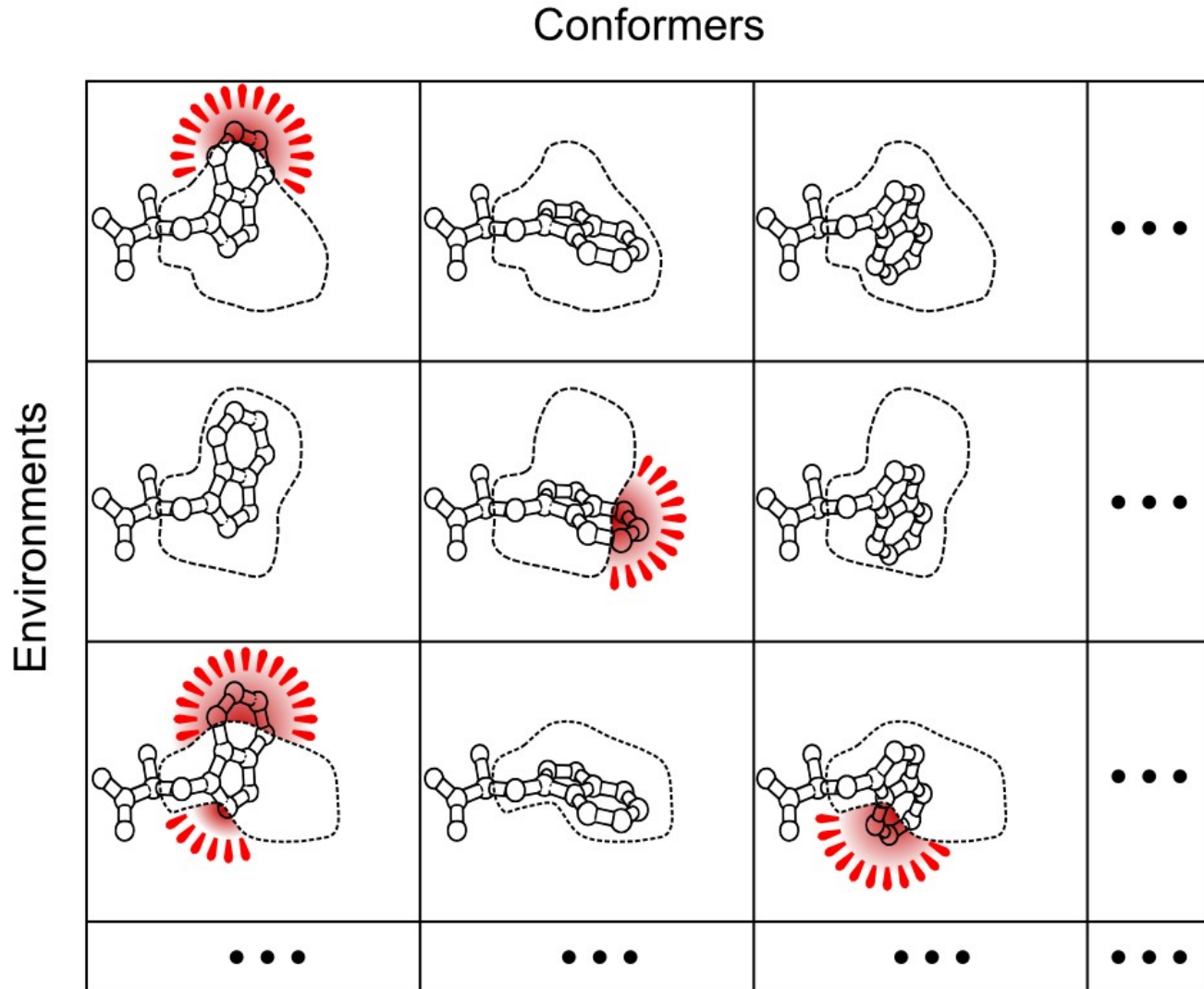
- Use energetics to identify the best sampling strategy for side chain optimization



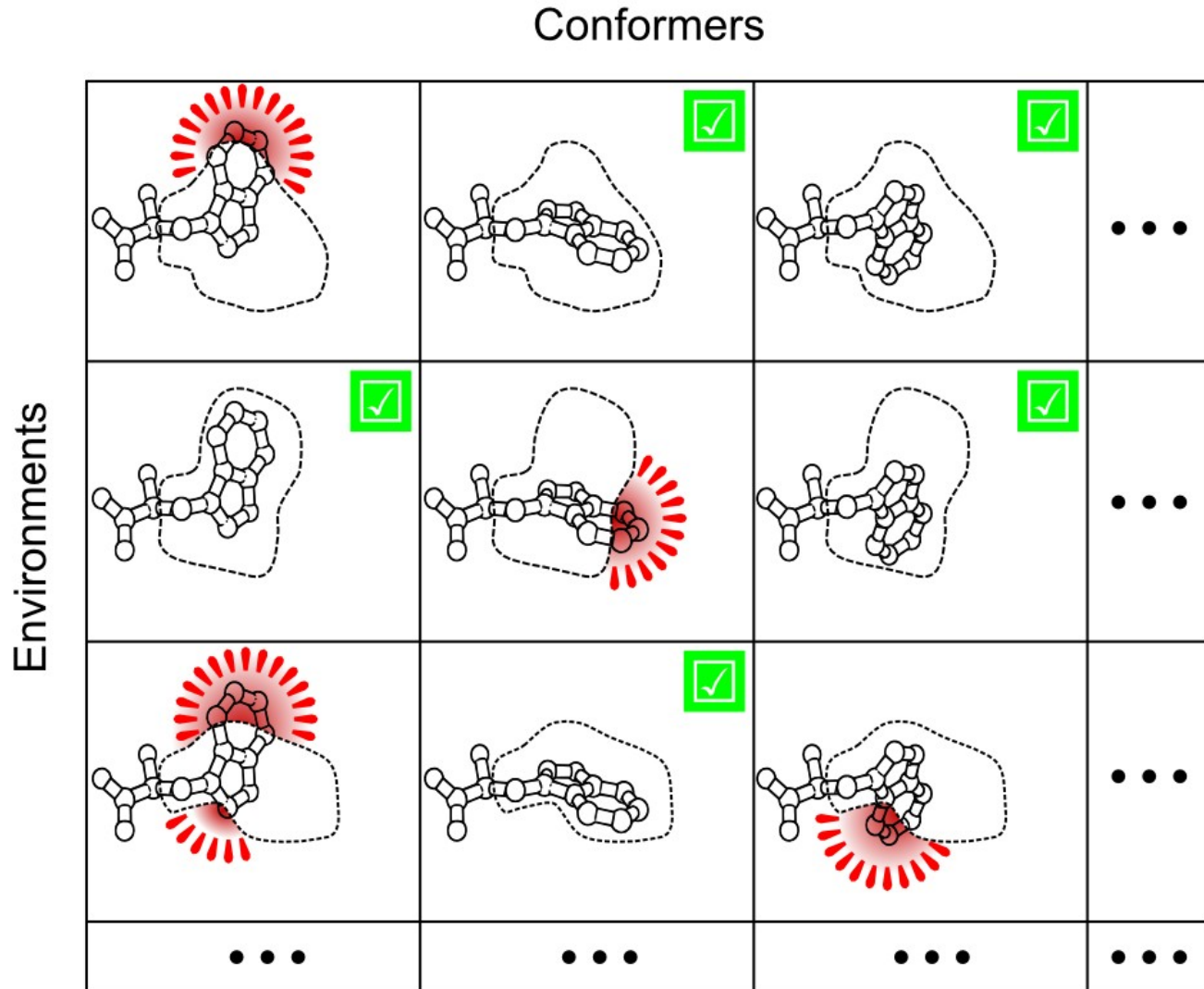
- Use energetics to identify the best sampling strategy for side chain optimization



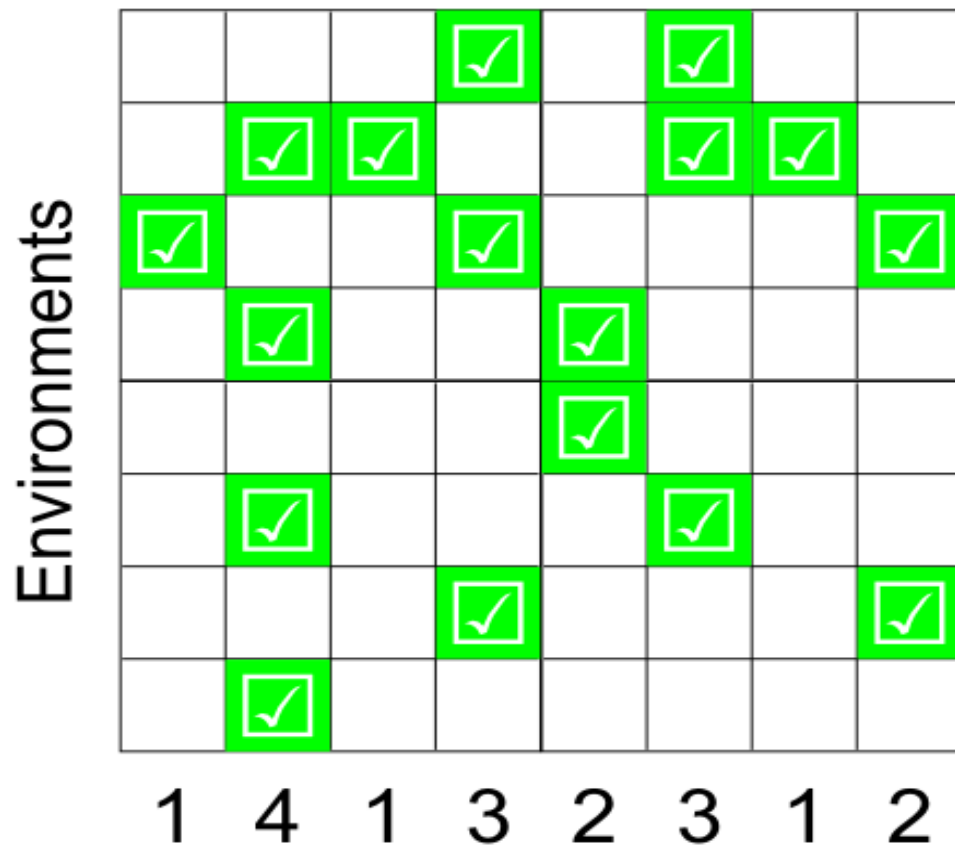
- Use energetics to identify the best sampling strategy for side chain optimization



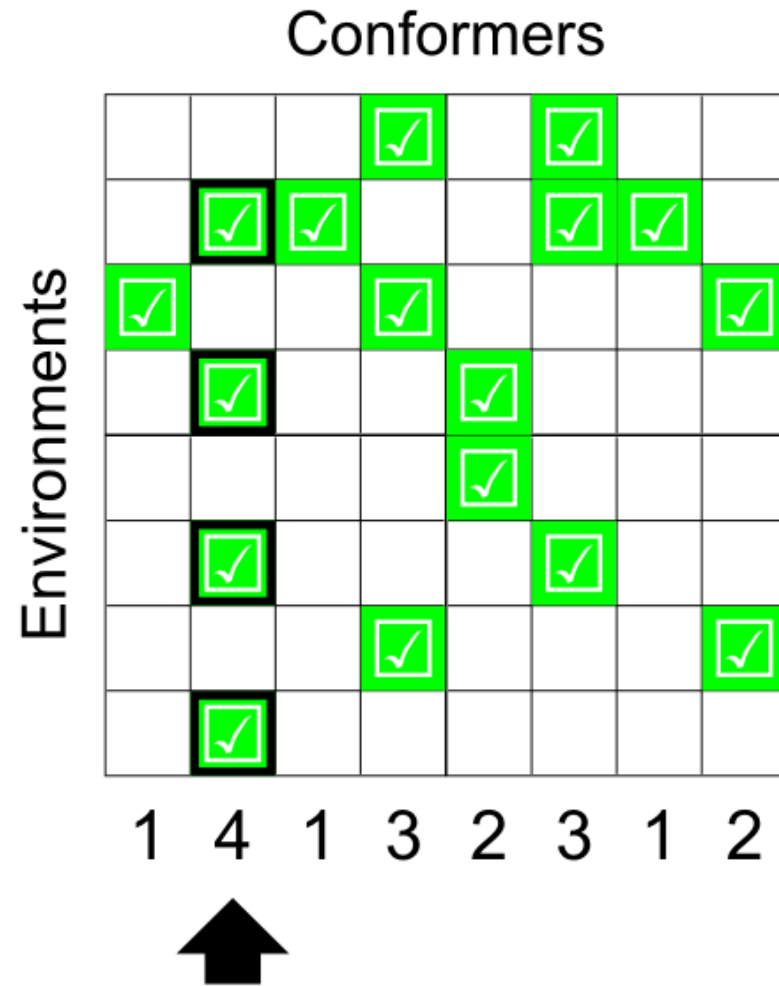
- Use energetics to identify the best sampling strategy for side chain optimization



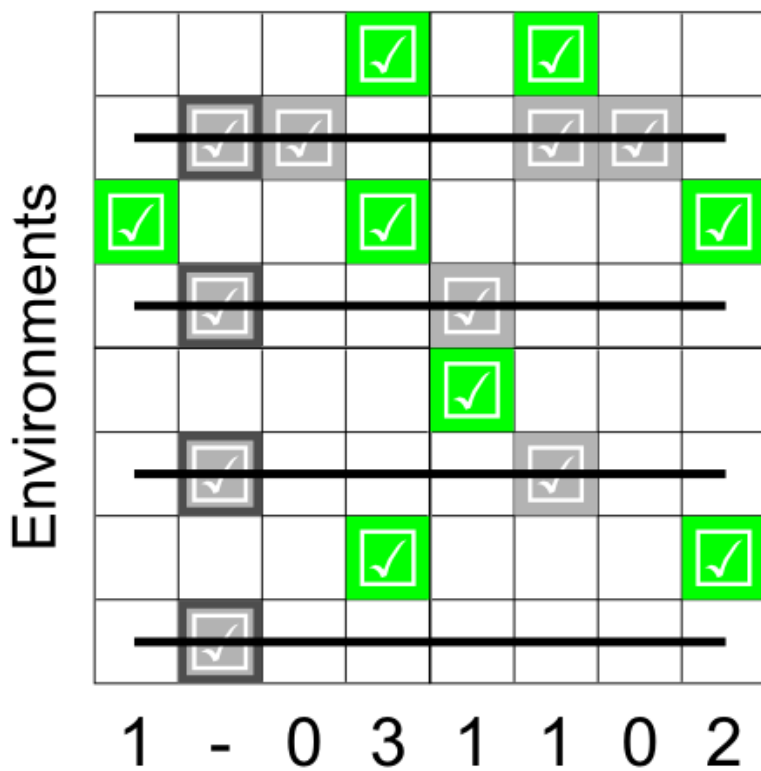
Conformers



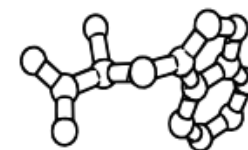
THE first conformer



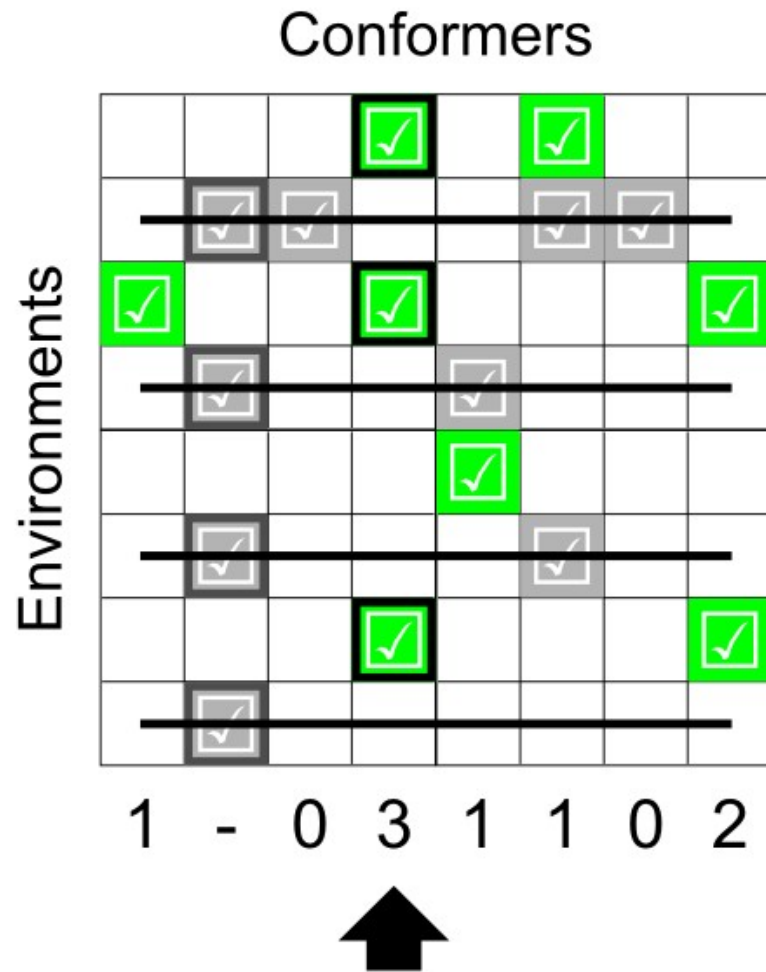
Conformers



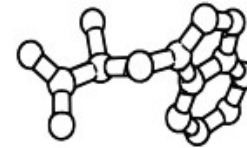
1st



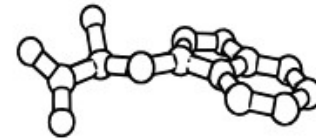
THE second conformer



1st

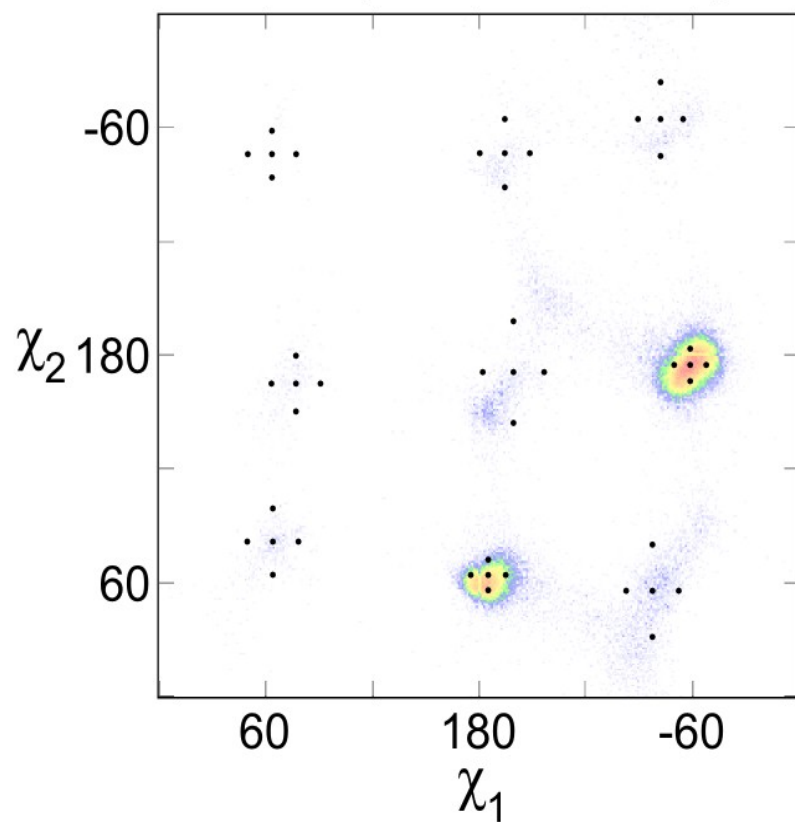


2nd

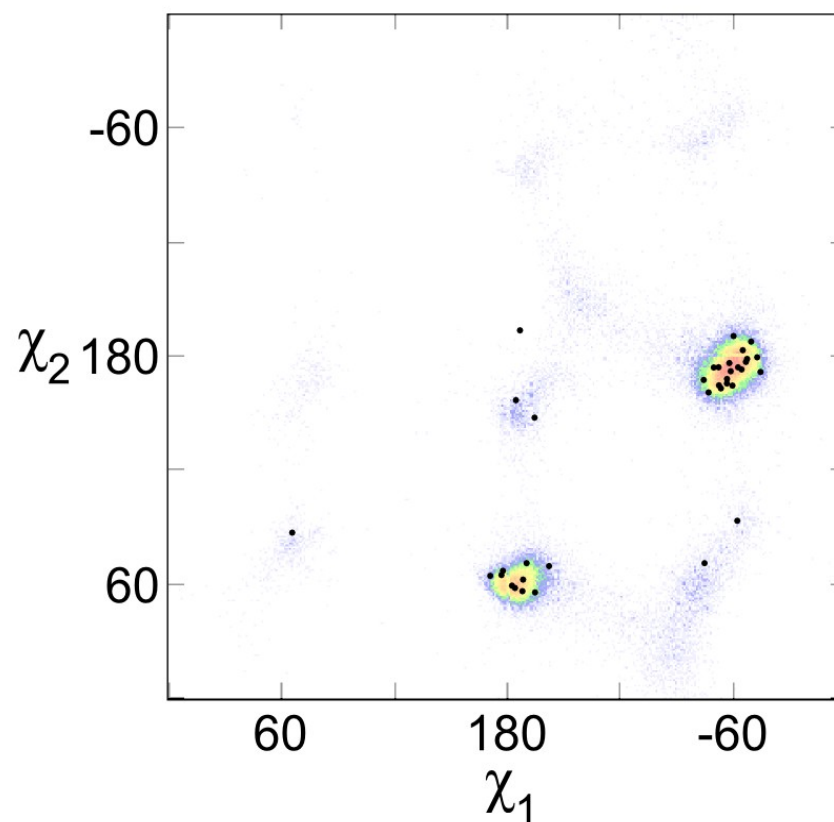


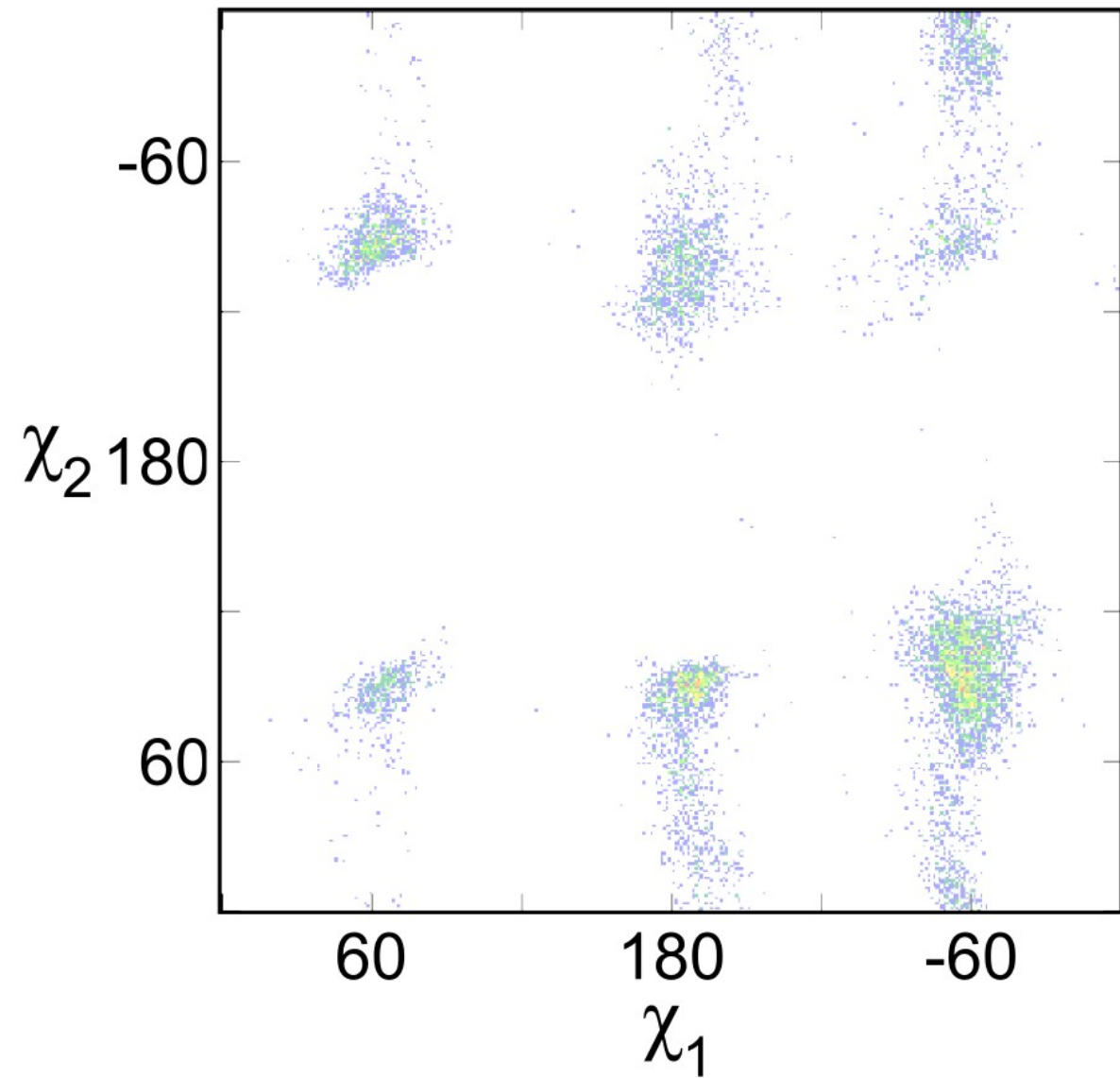
CONFORMERS IN PROPORTION TO DISTRIBUTION

5x expanded library

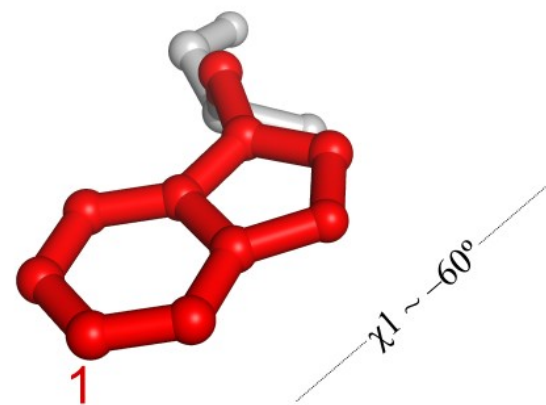
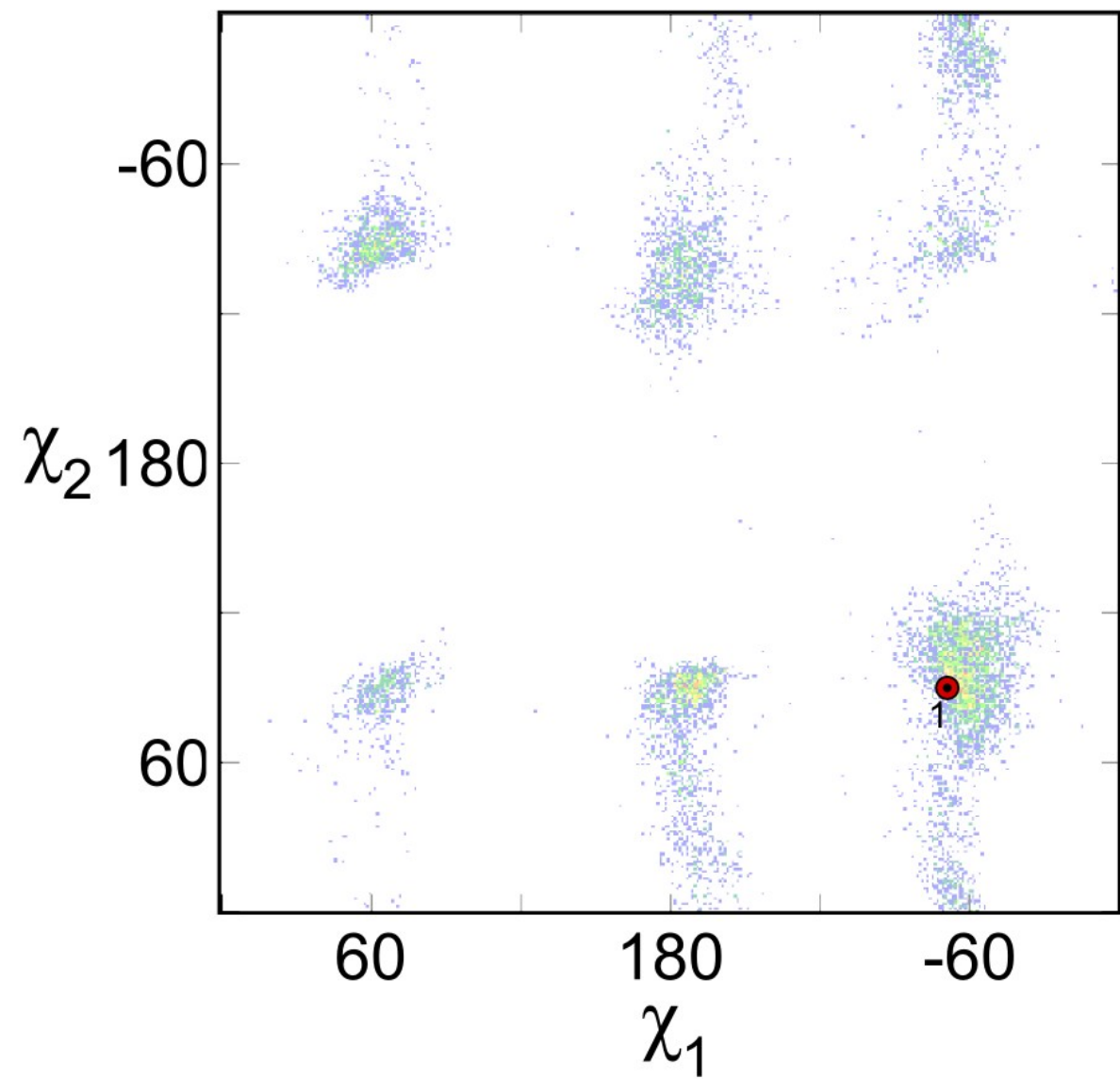


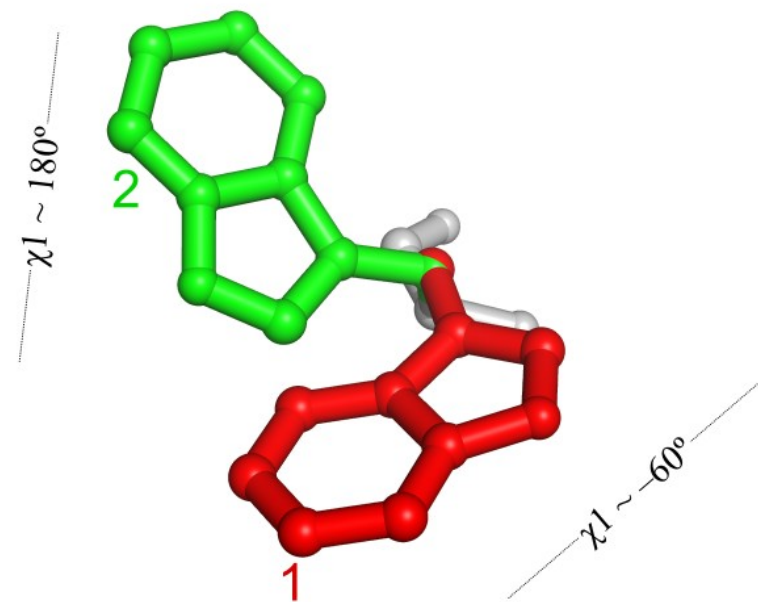
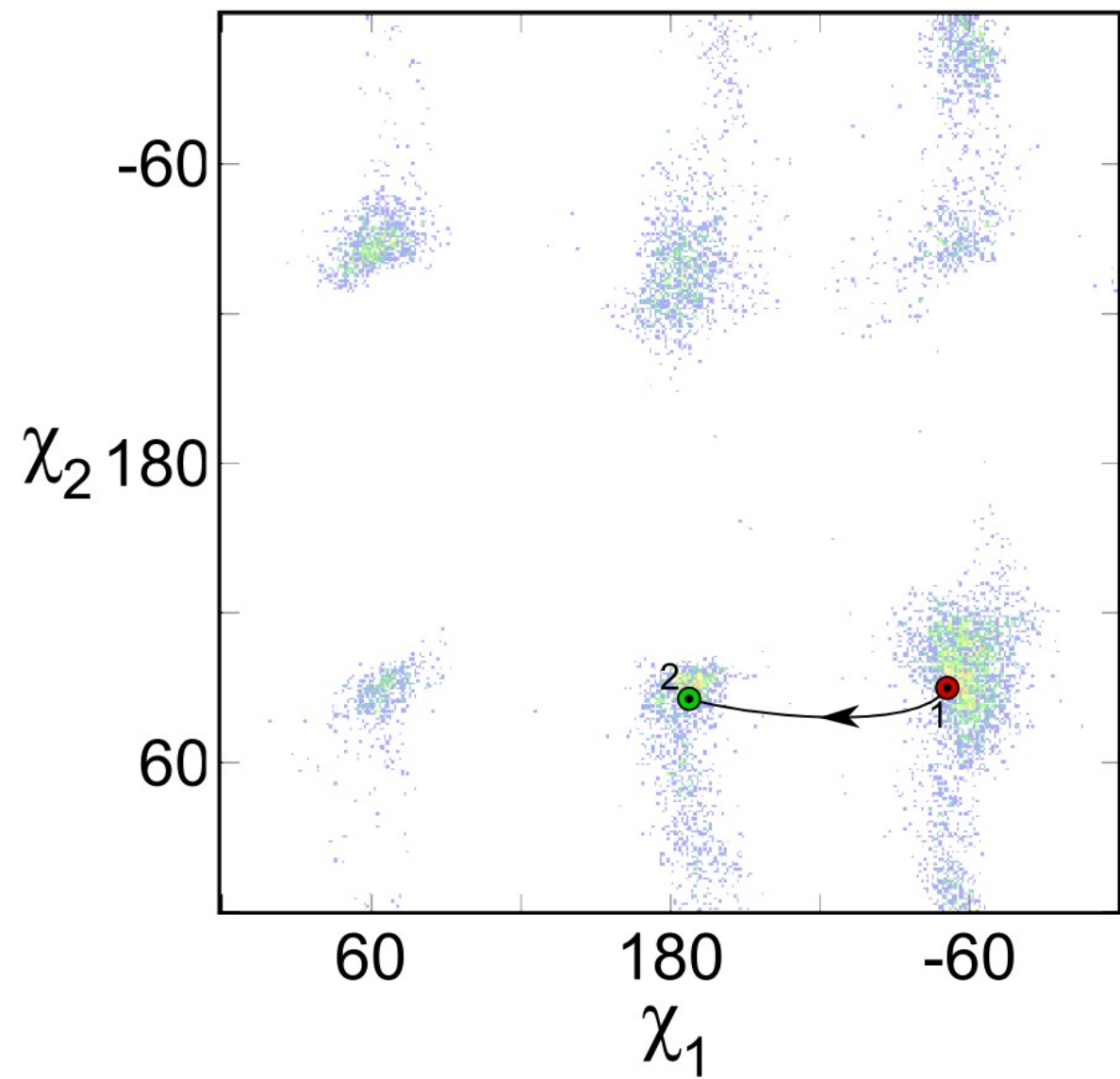
Energy-Based Library

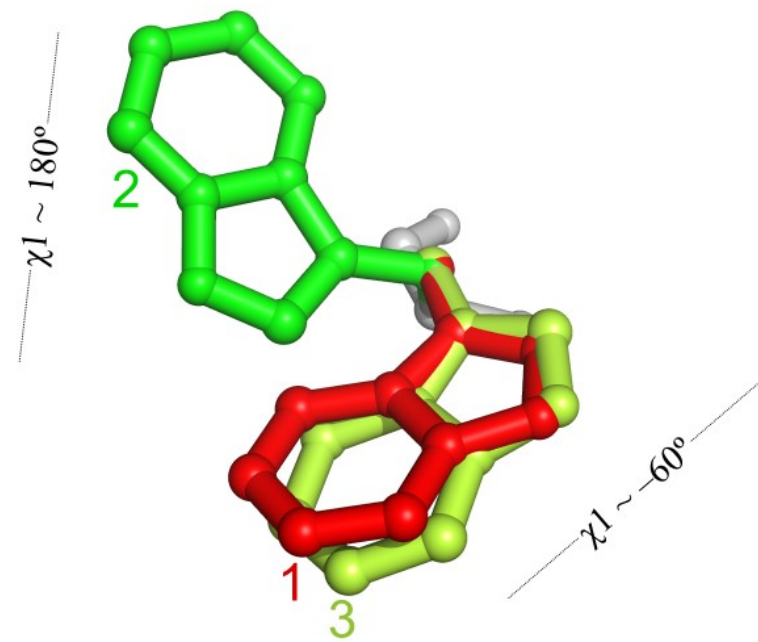
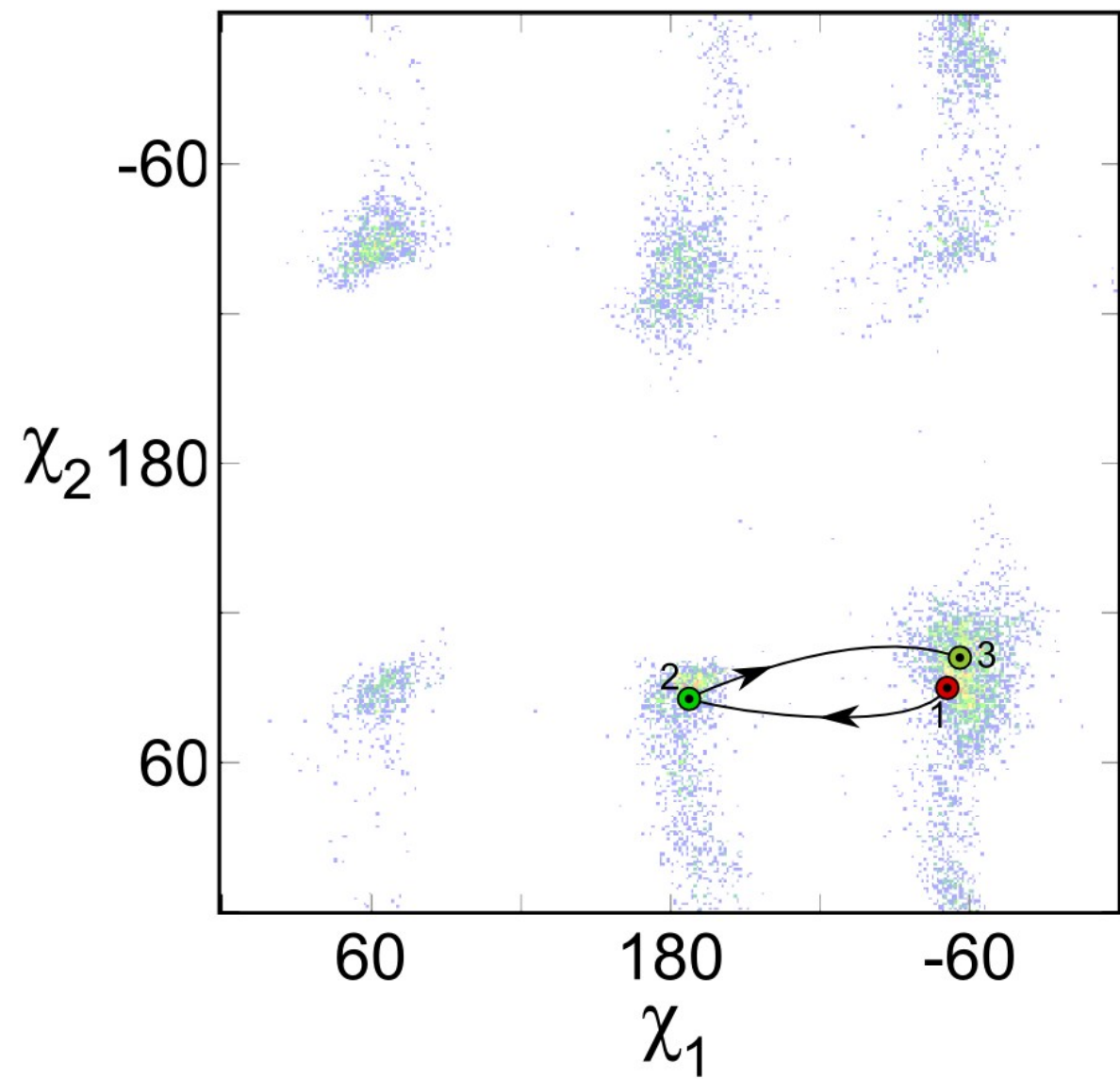


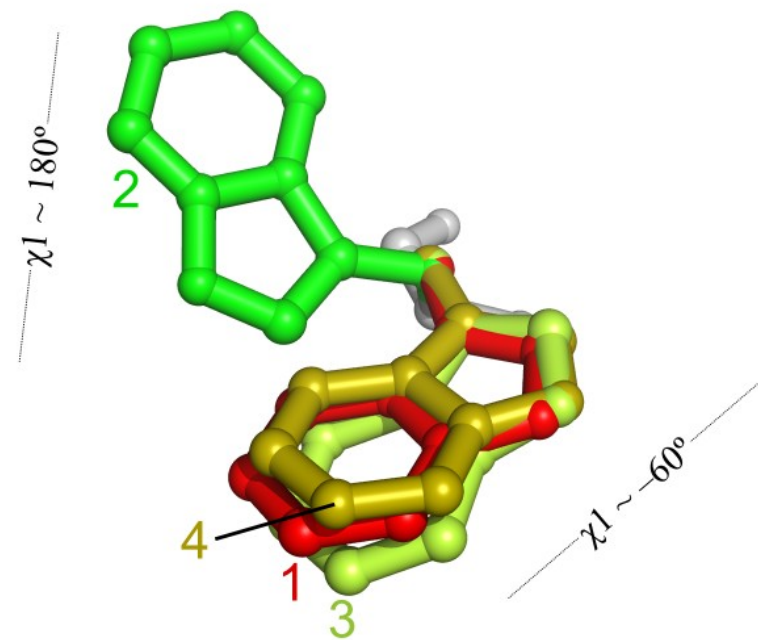
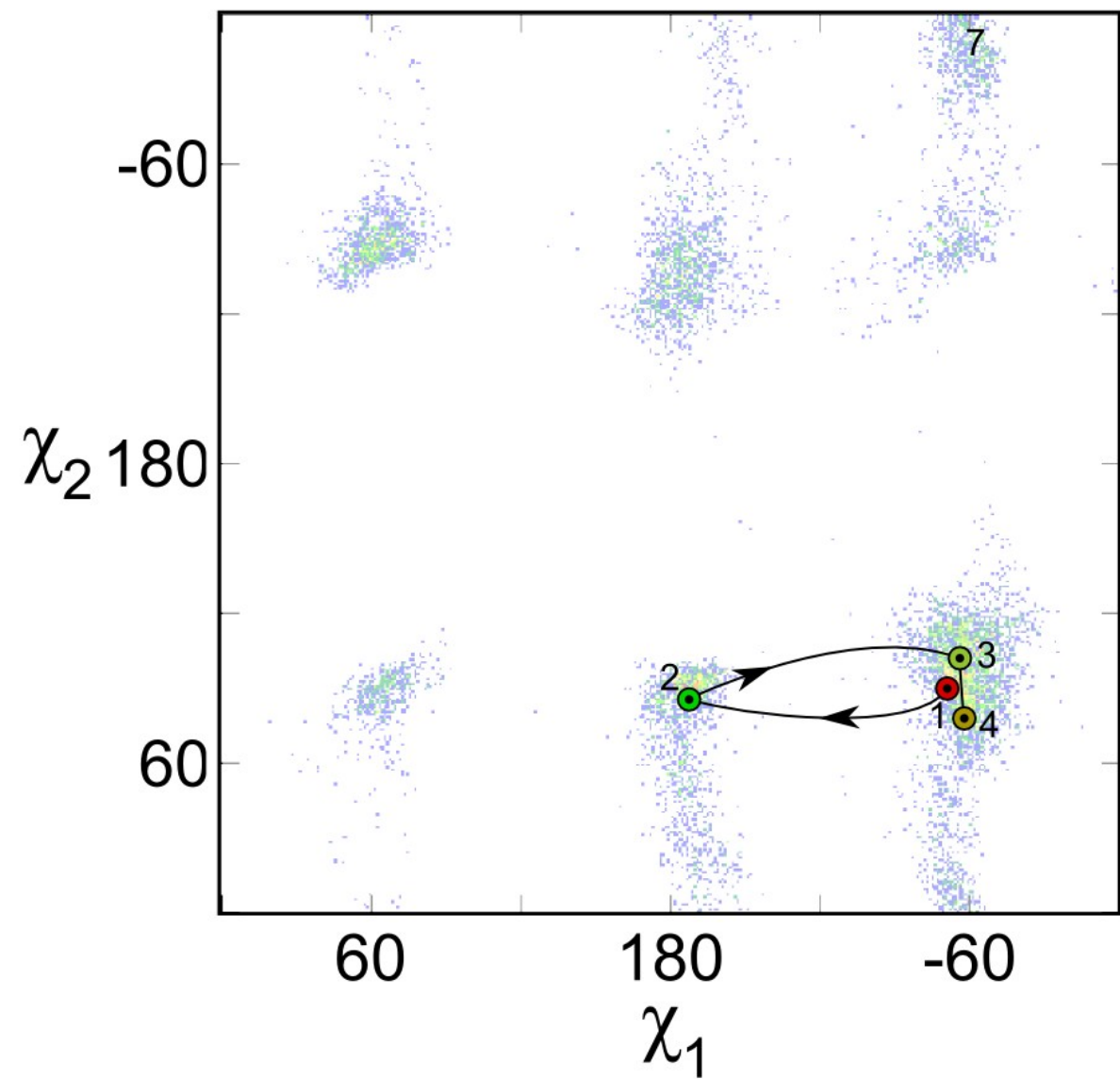


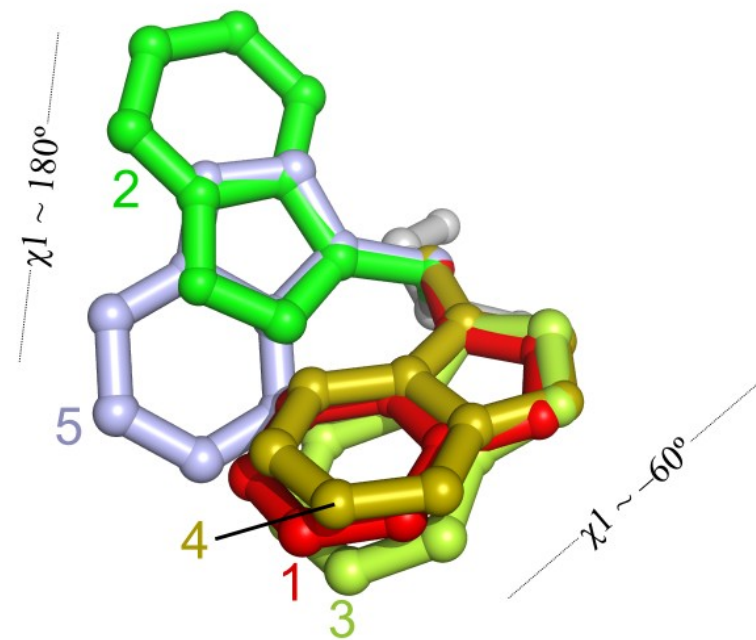
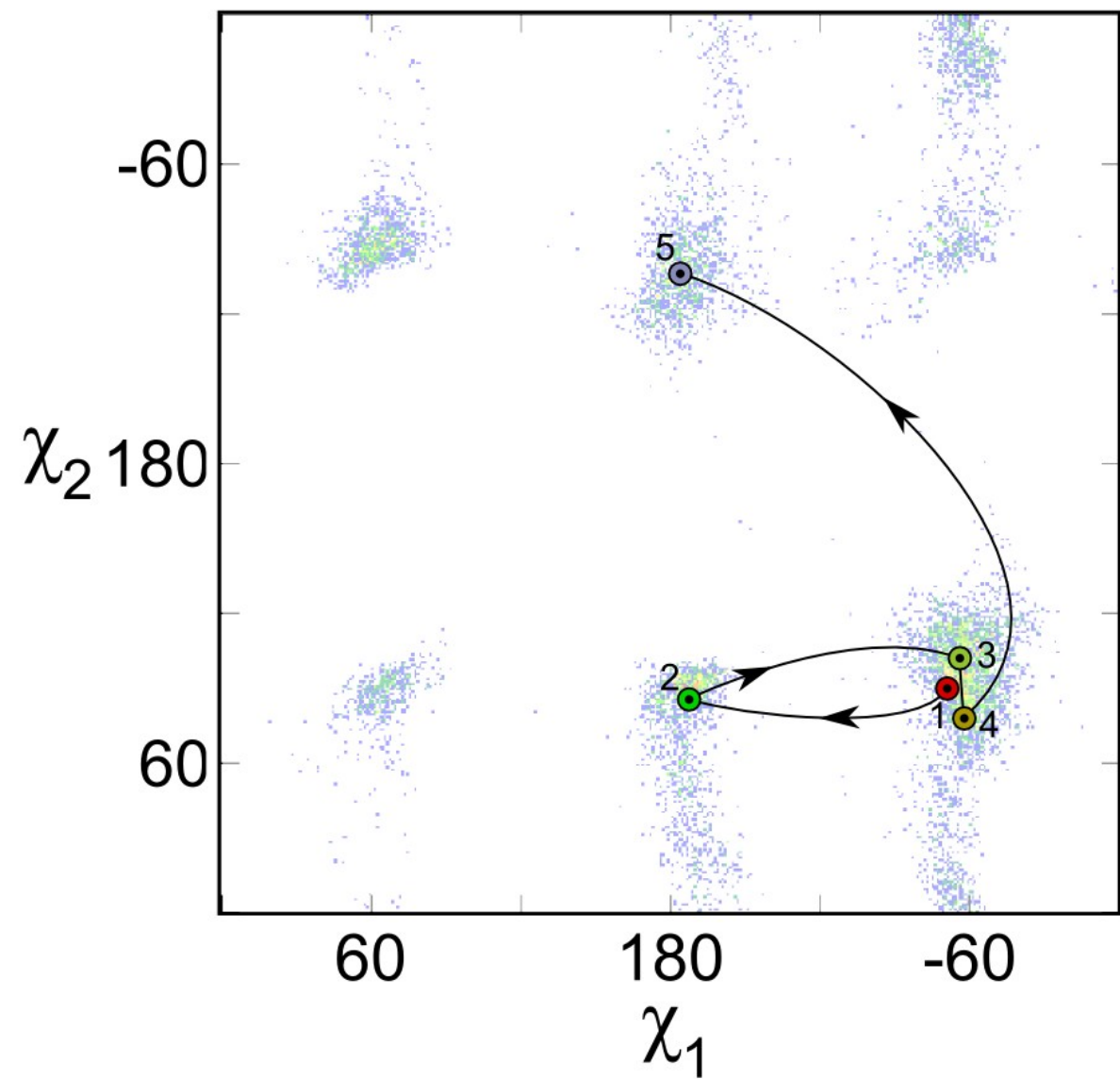
**A walk in
Trp space**

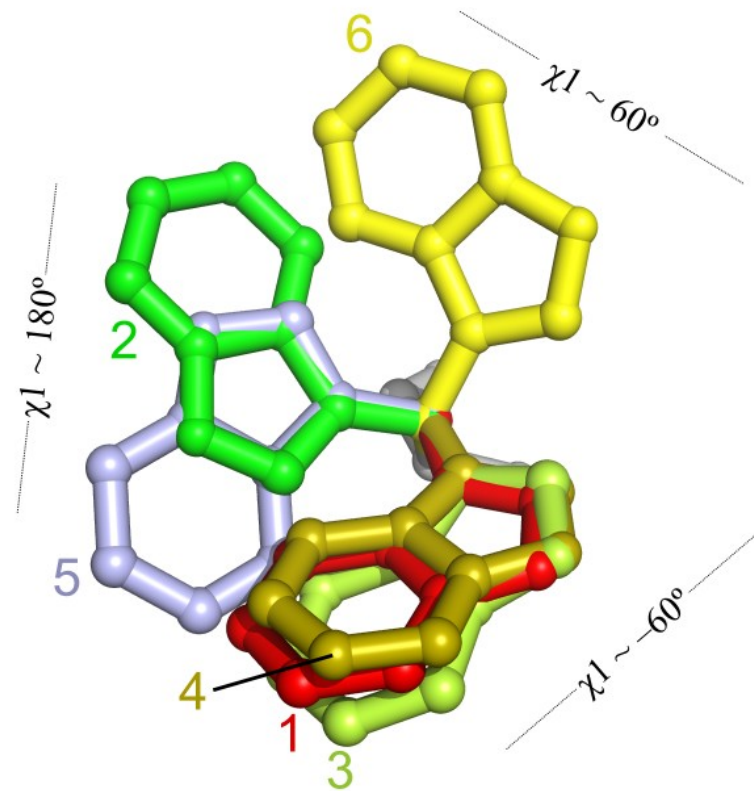
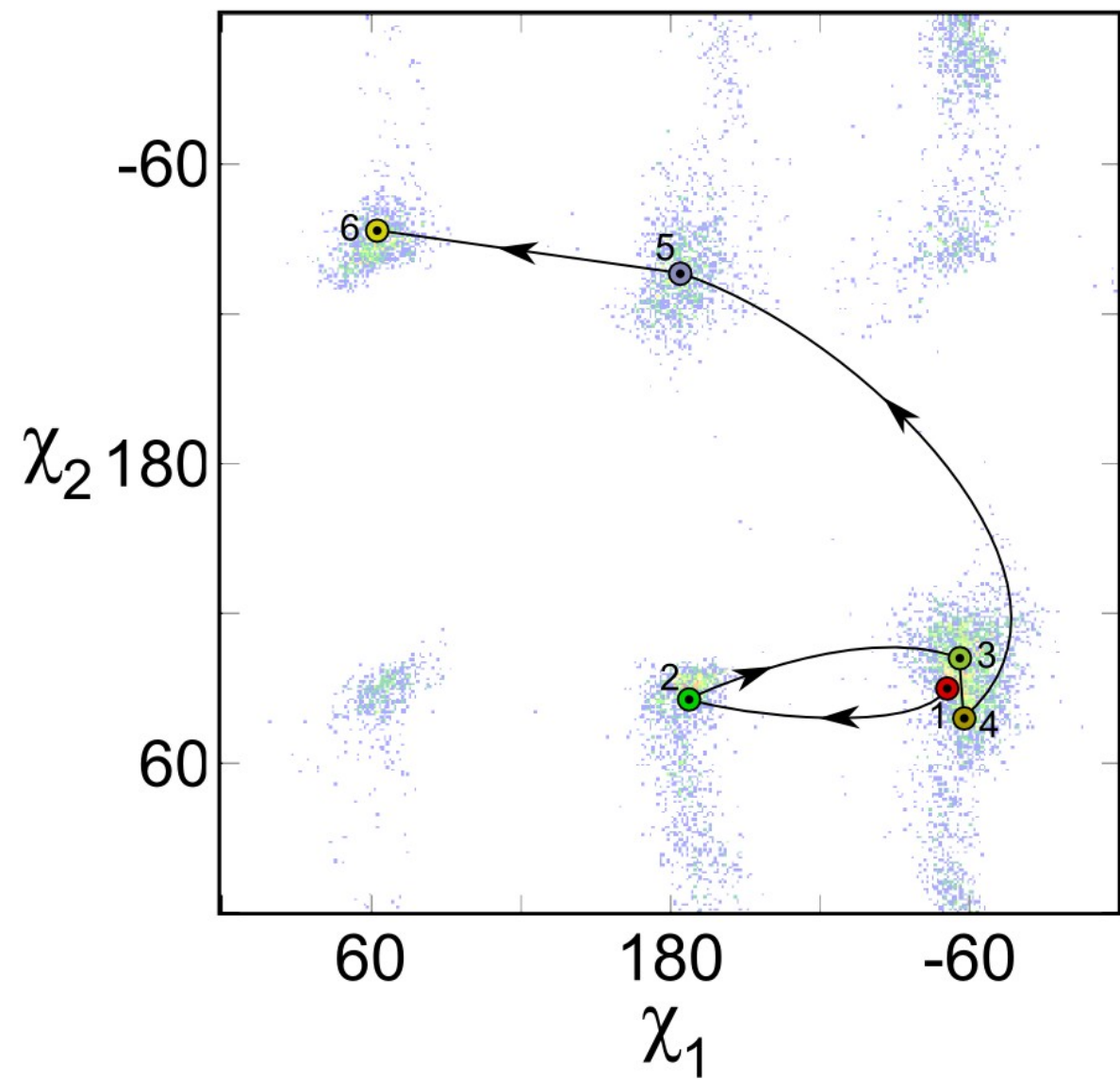


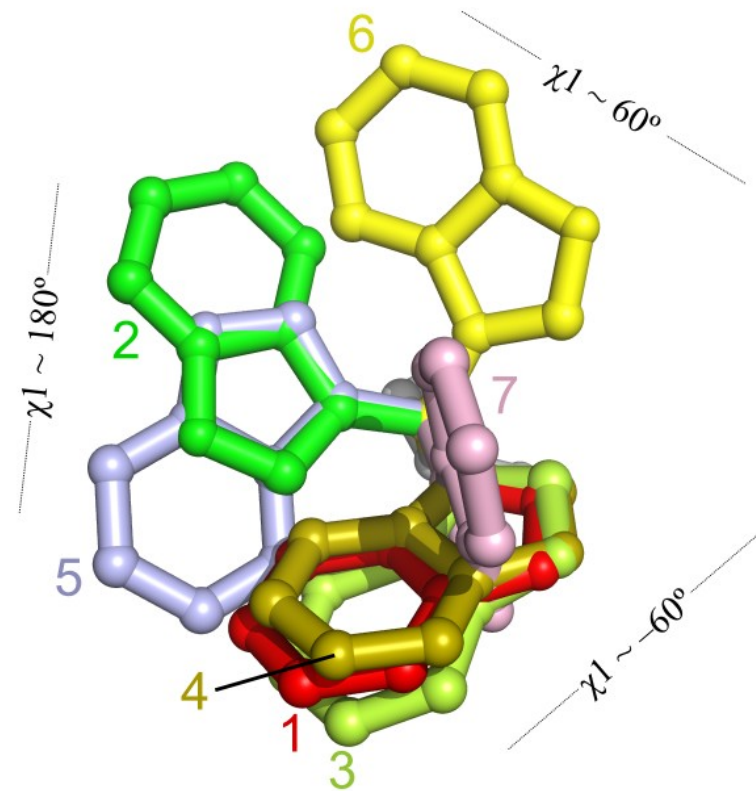
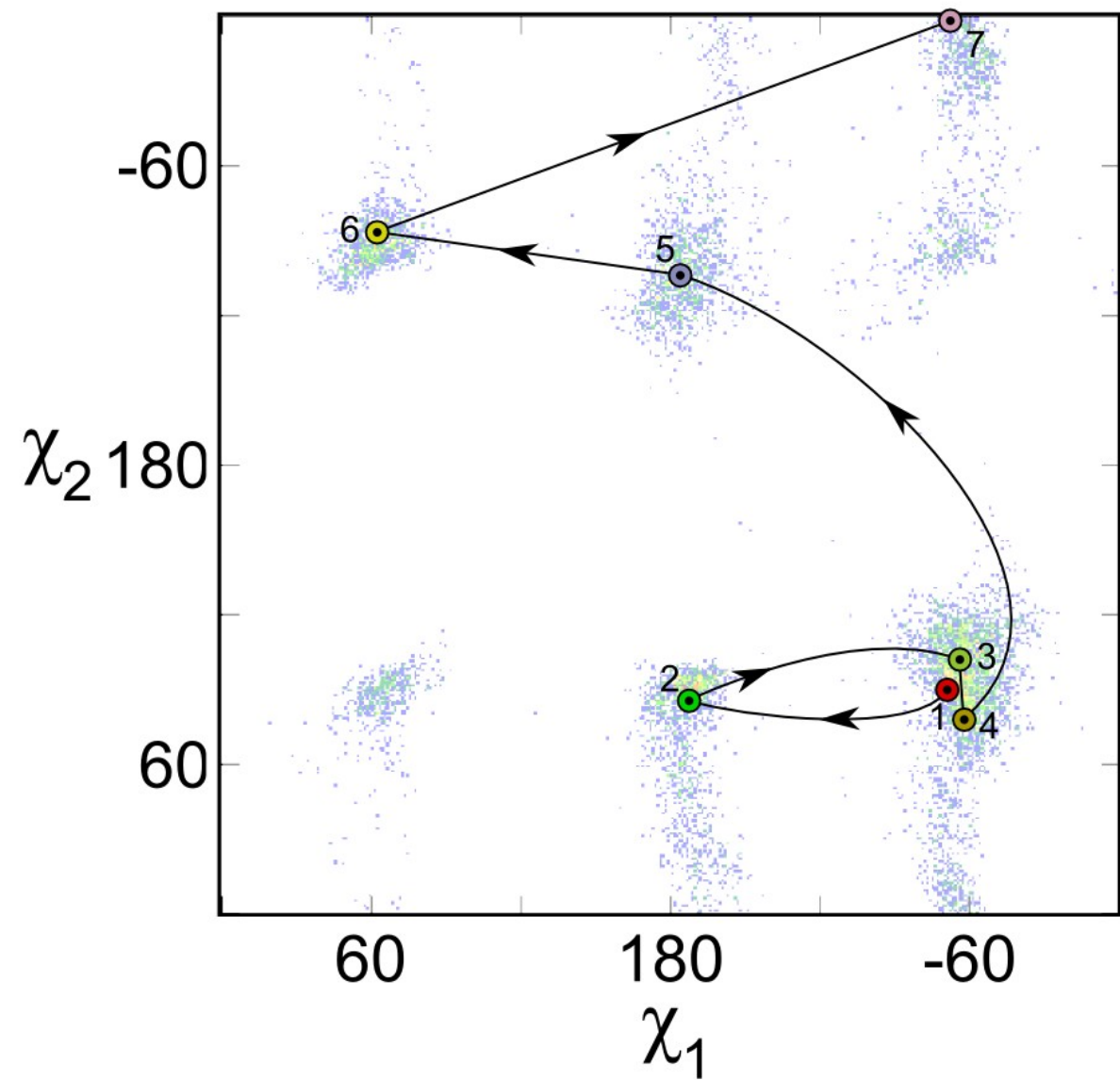


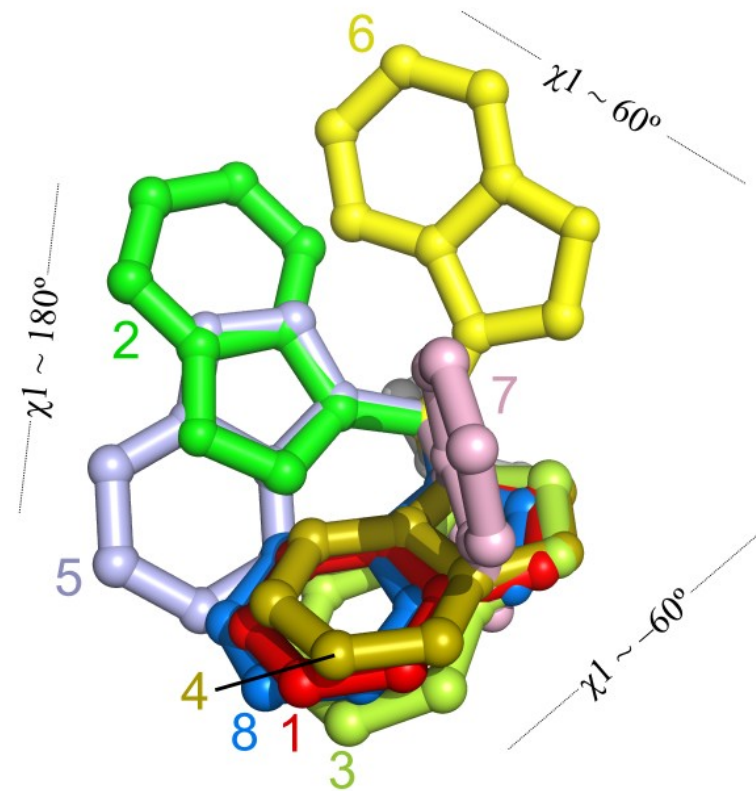
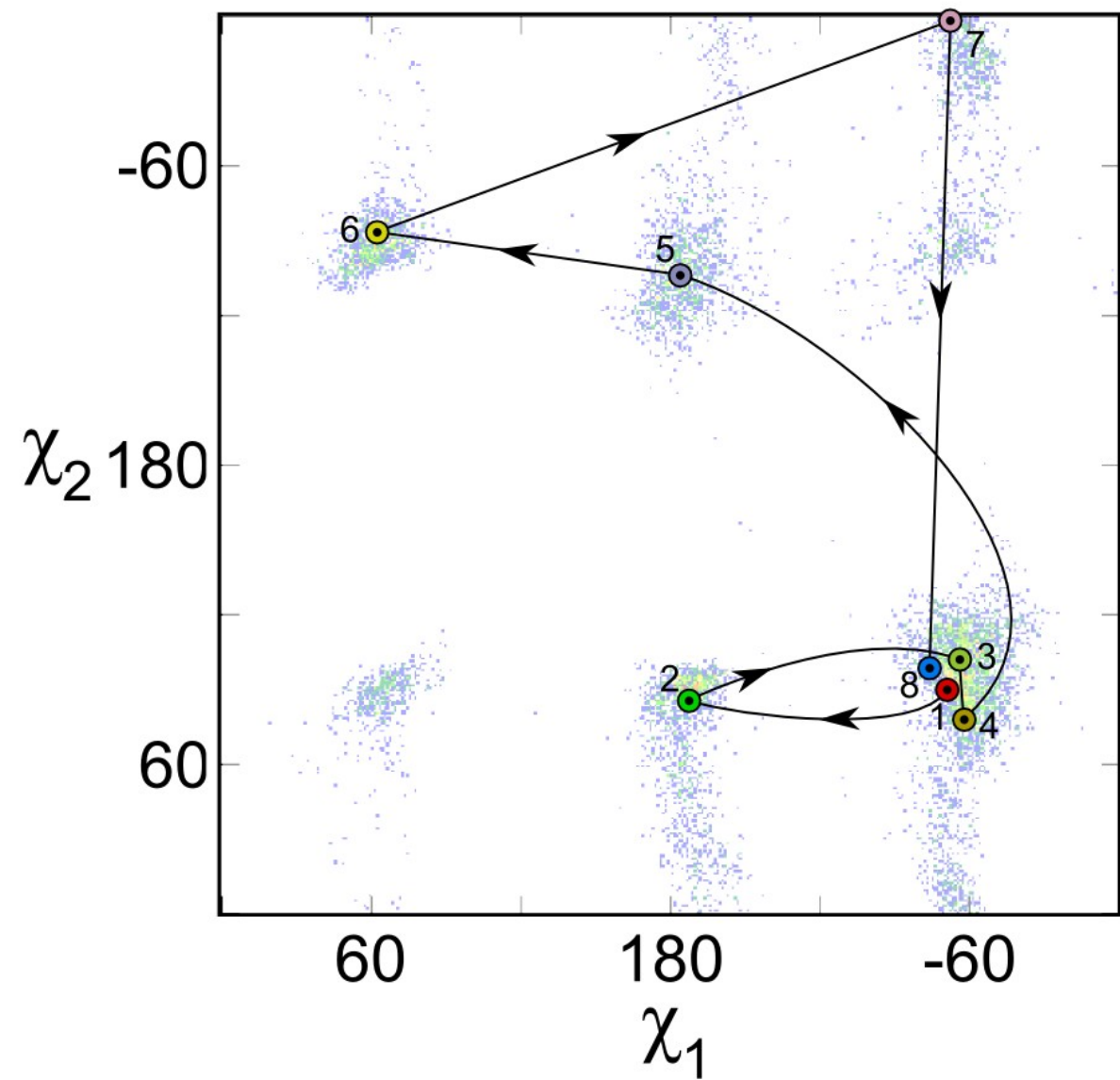


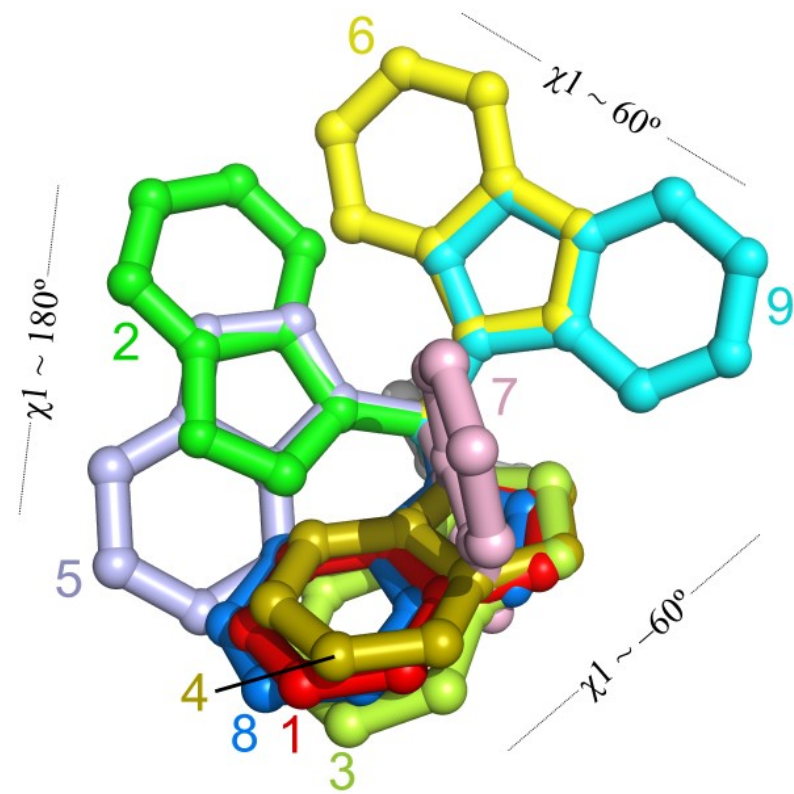
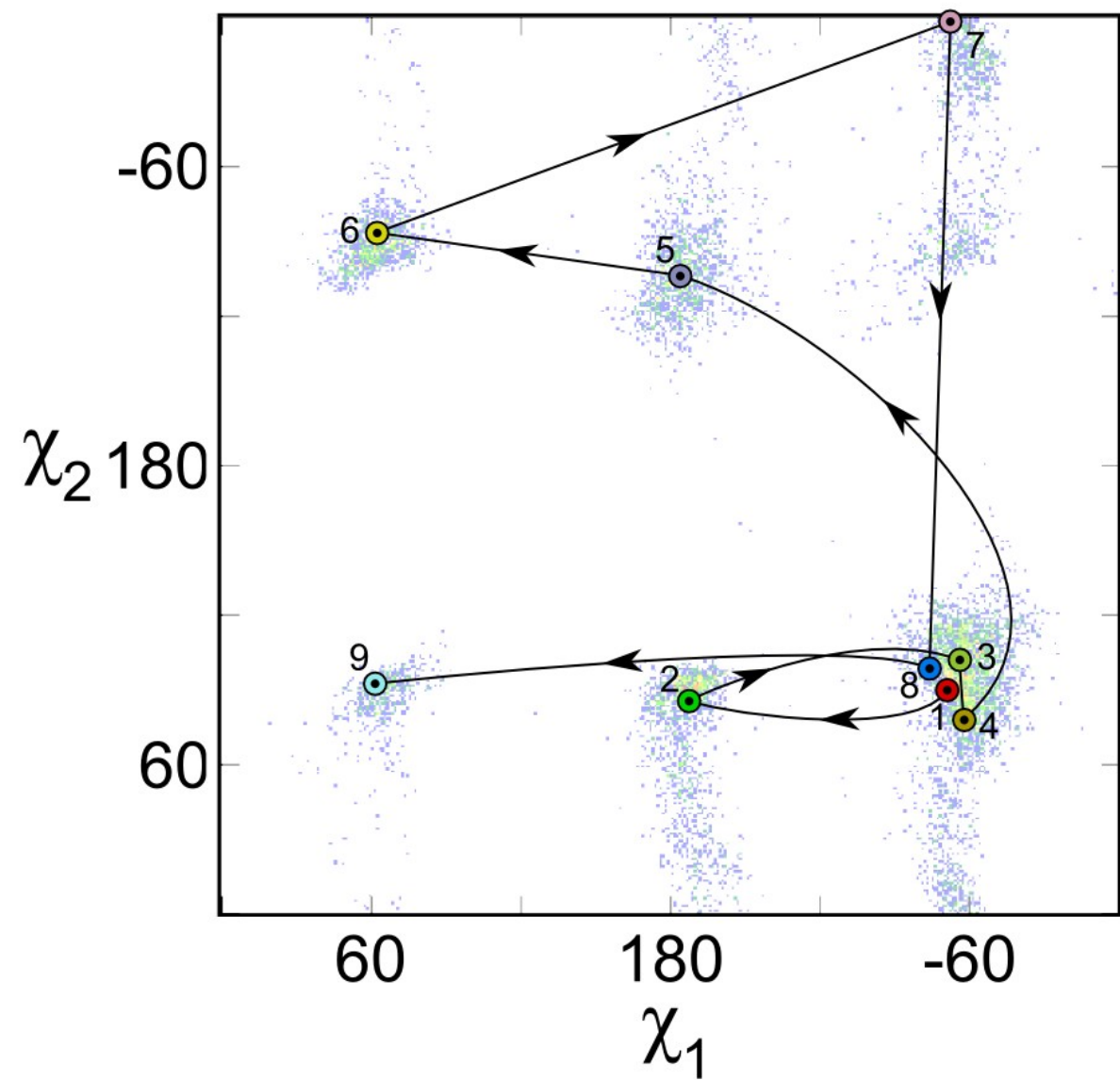


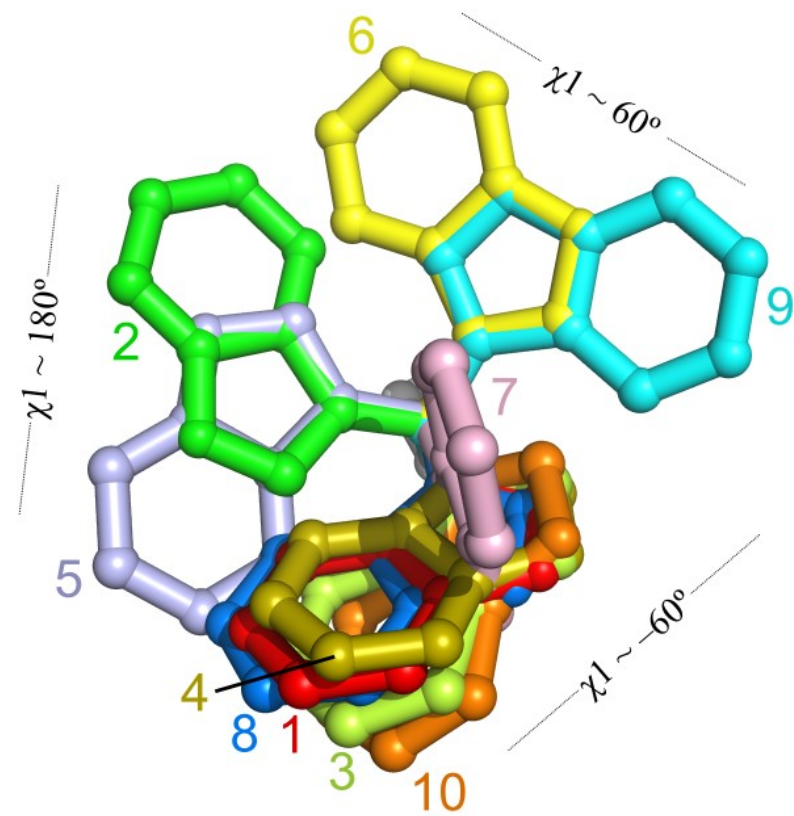
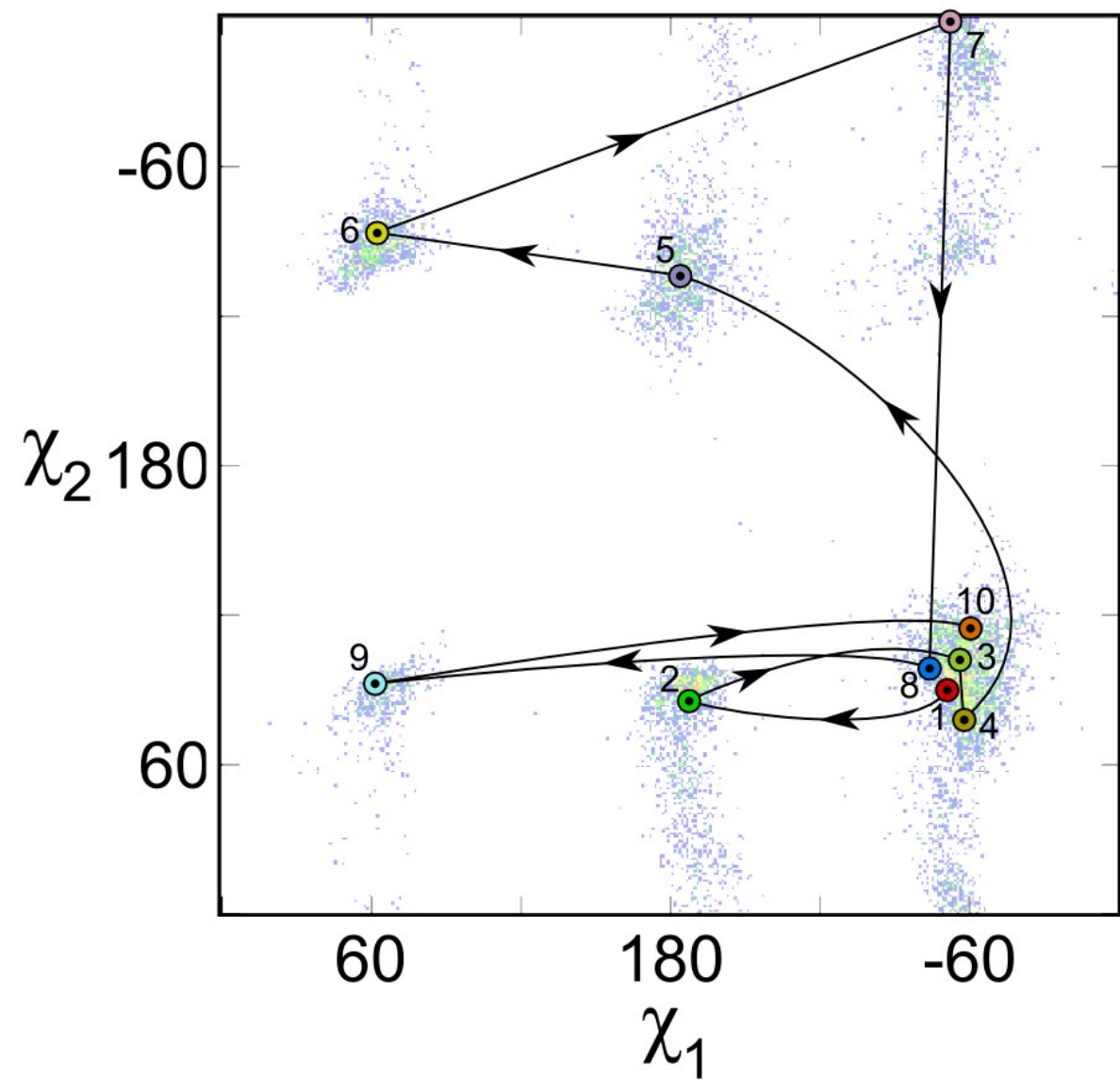


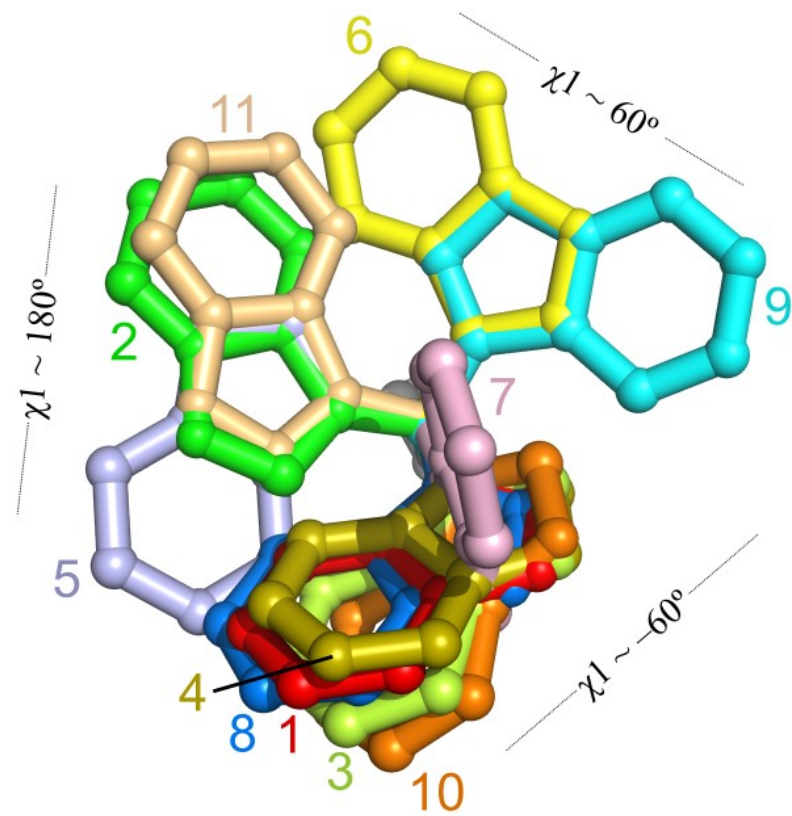
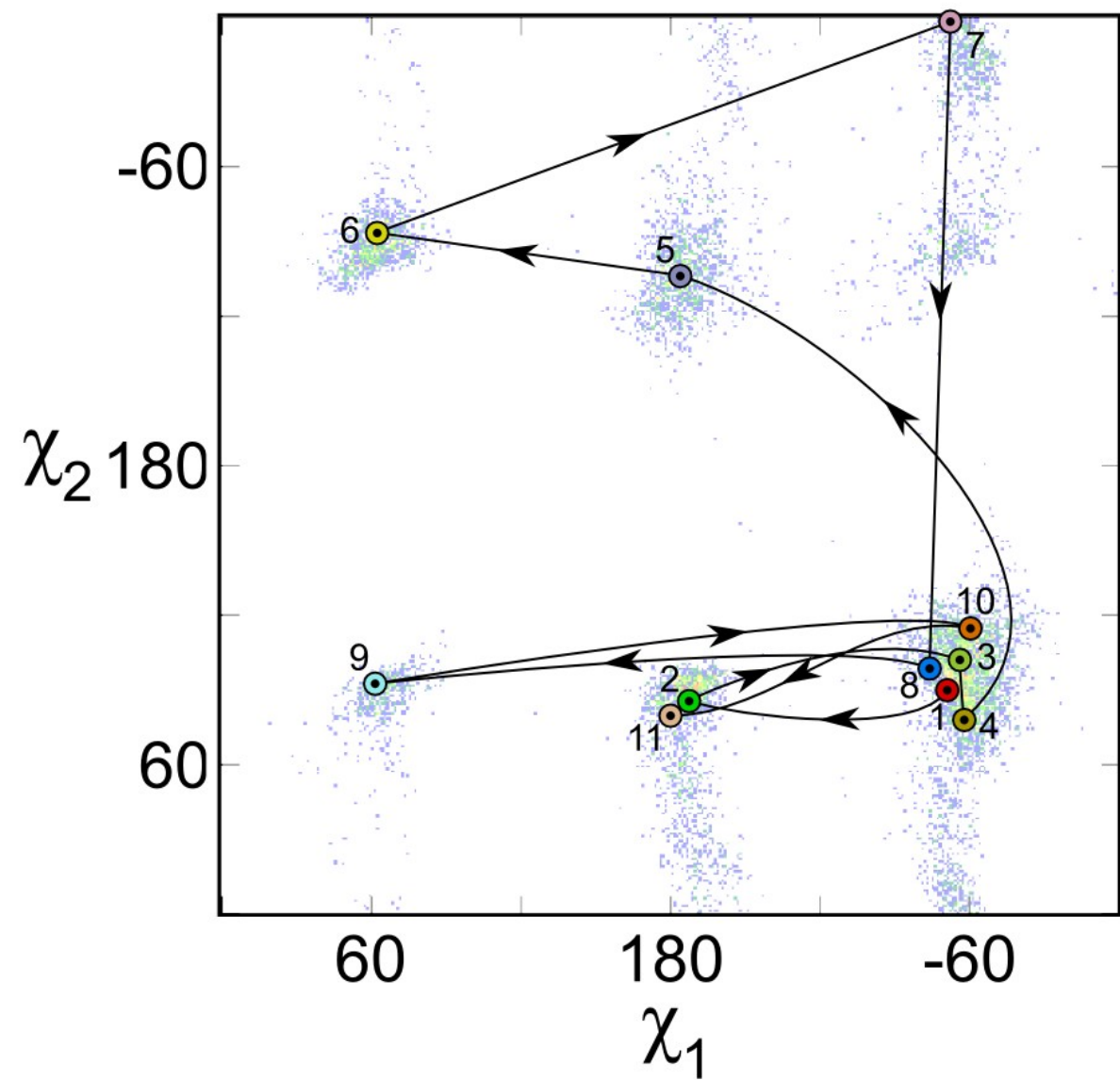


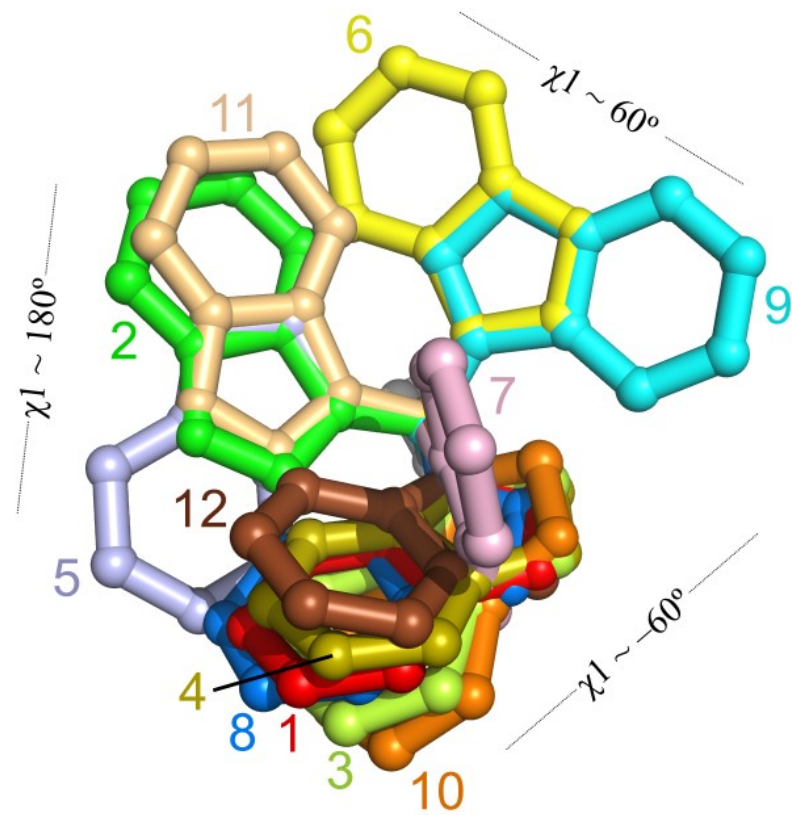
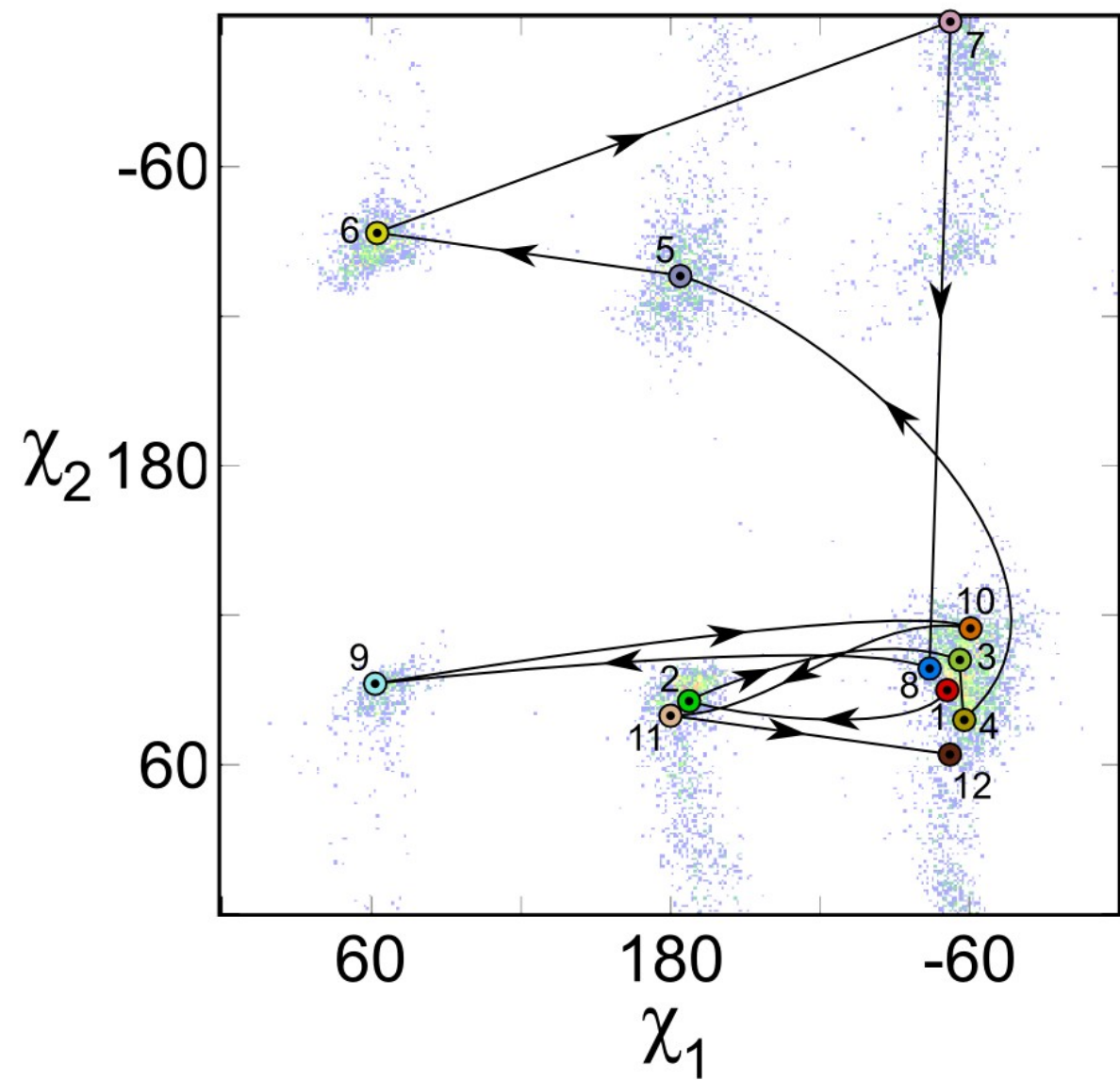


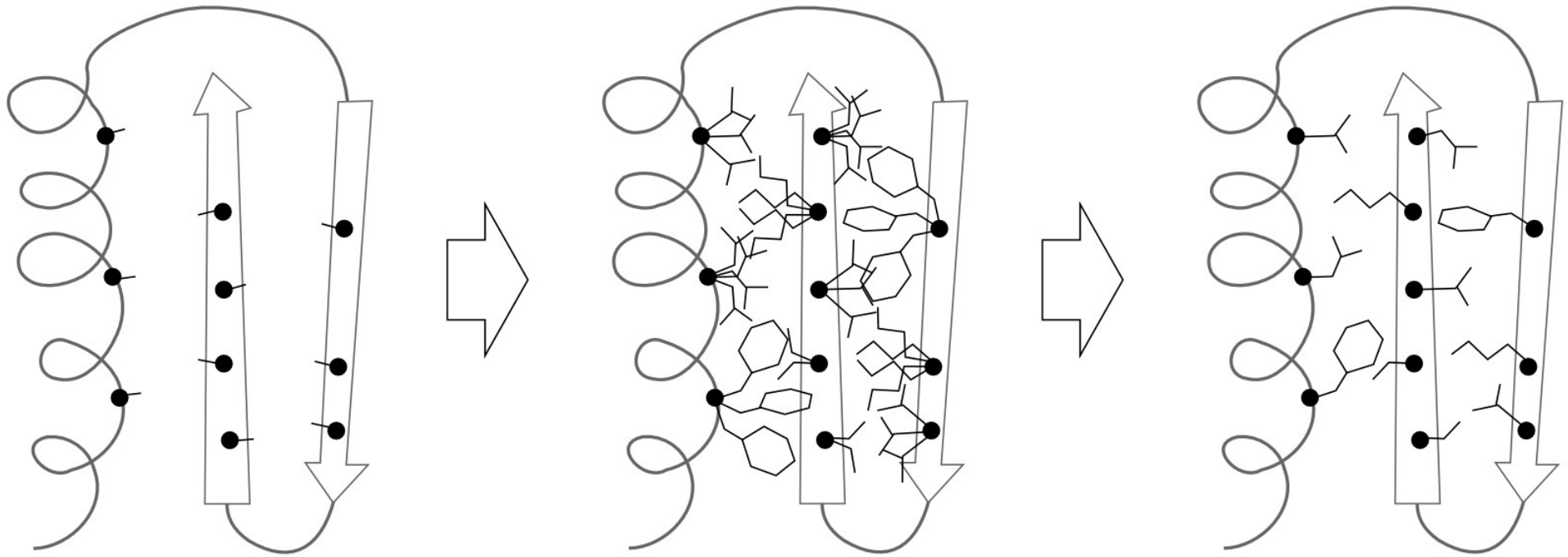








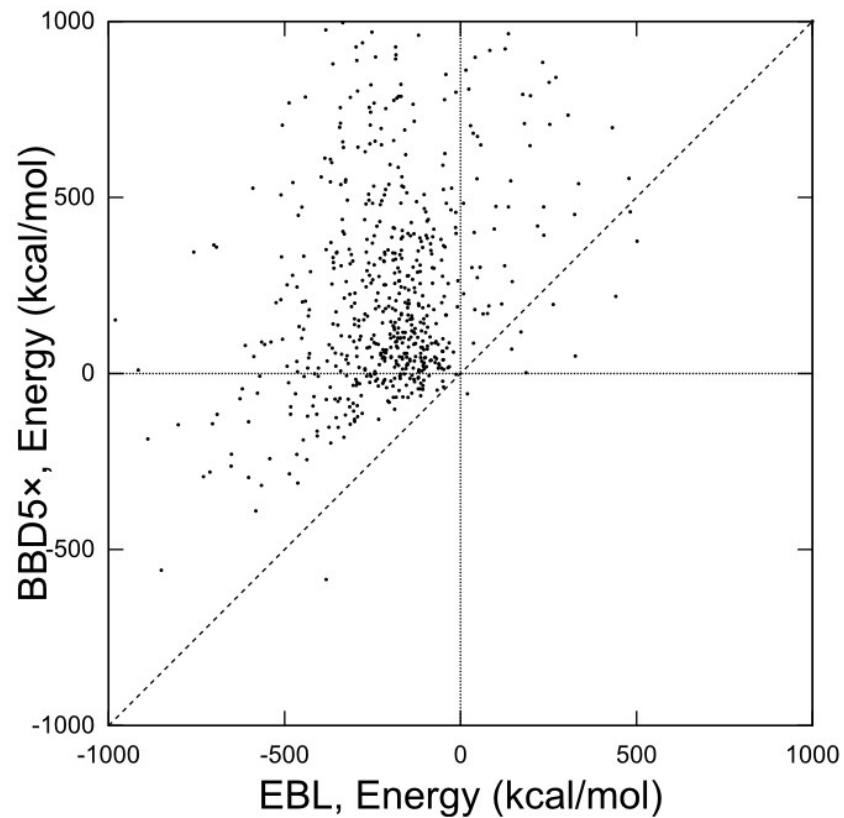




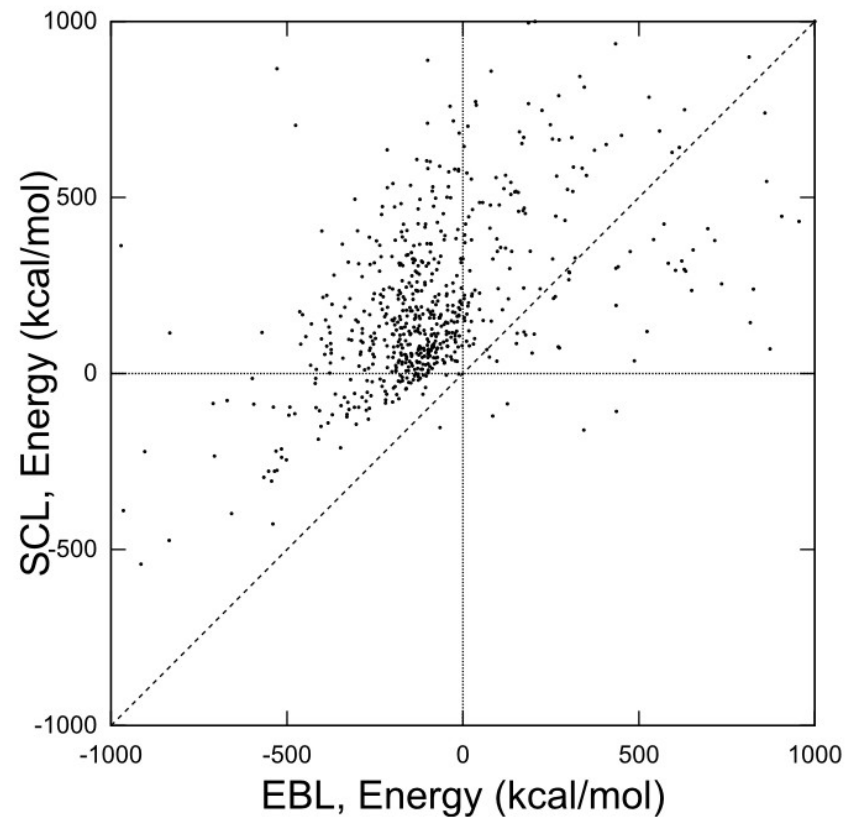
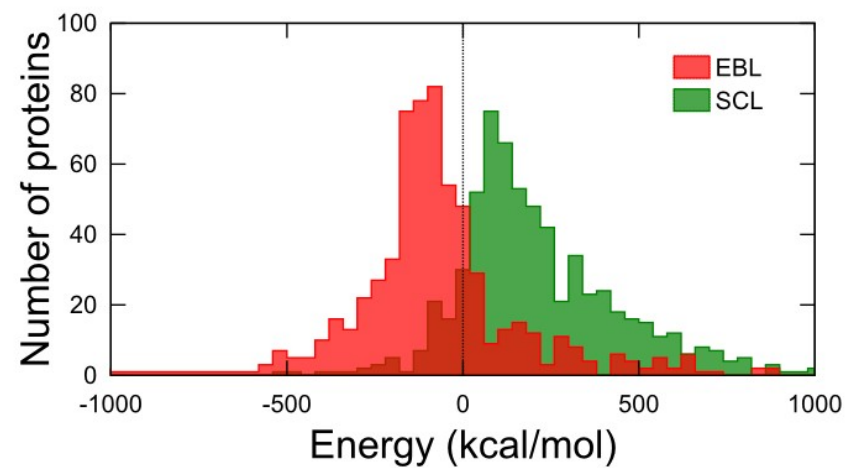
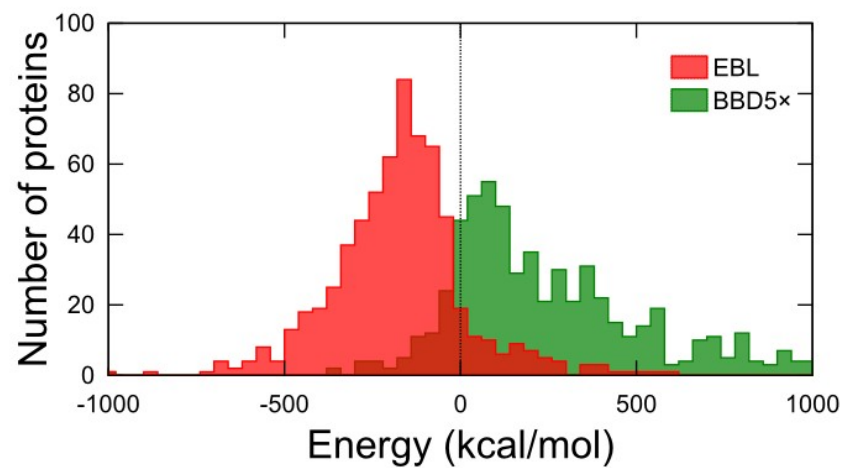
721 complete protein repacks.
The lower the energy, the better.

a

EBL vs BBD5x



EBL vs SCL

**b**

Flexible High-Performance Conformer Library

- Conformers chosen using the same criterion as the optimization algorithm - Energy
- The new library is a sorted list of conformations
- Unprecedented flexibility – the first 'n' conformers is probably the best set of 'n' conformers



of the side chain conformation library



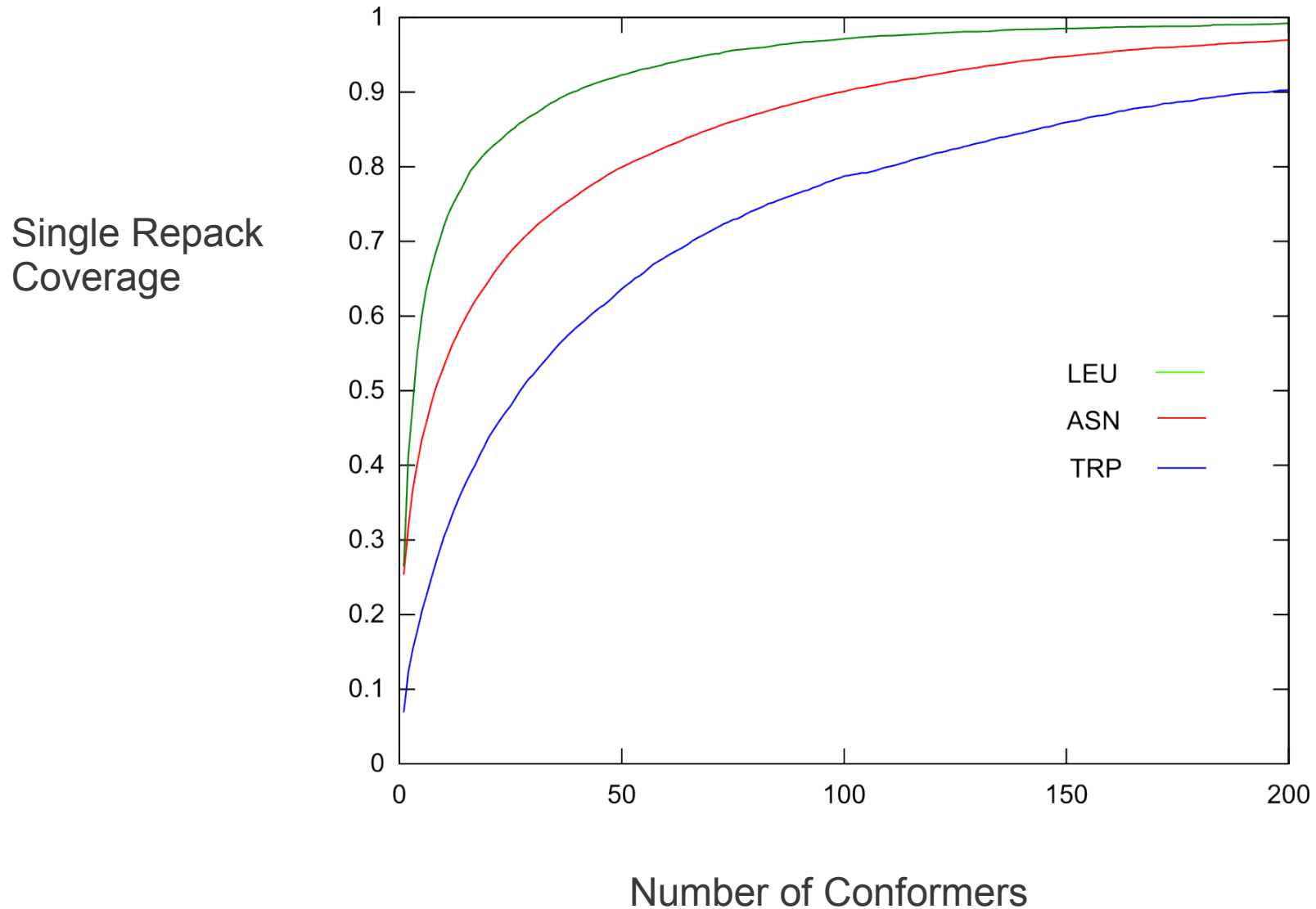
side chain



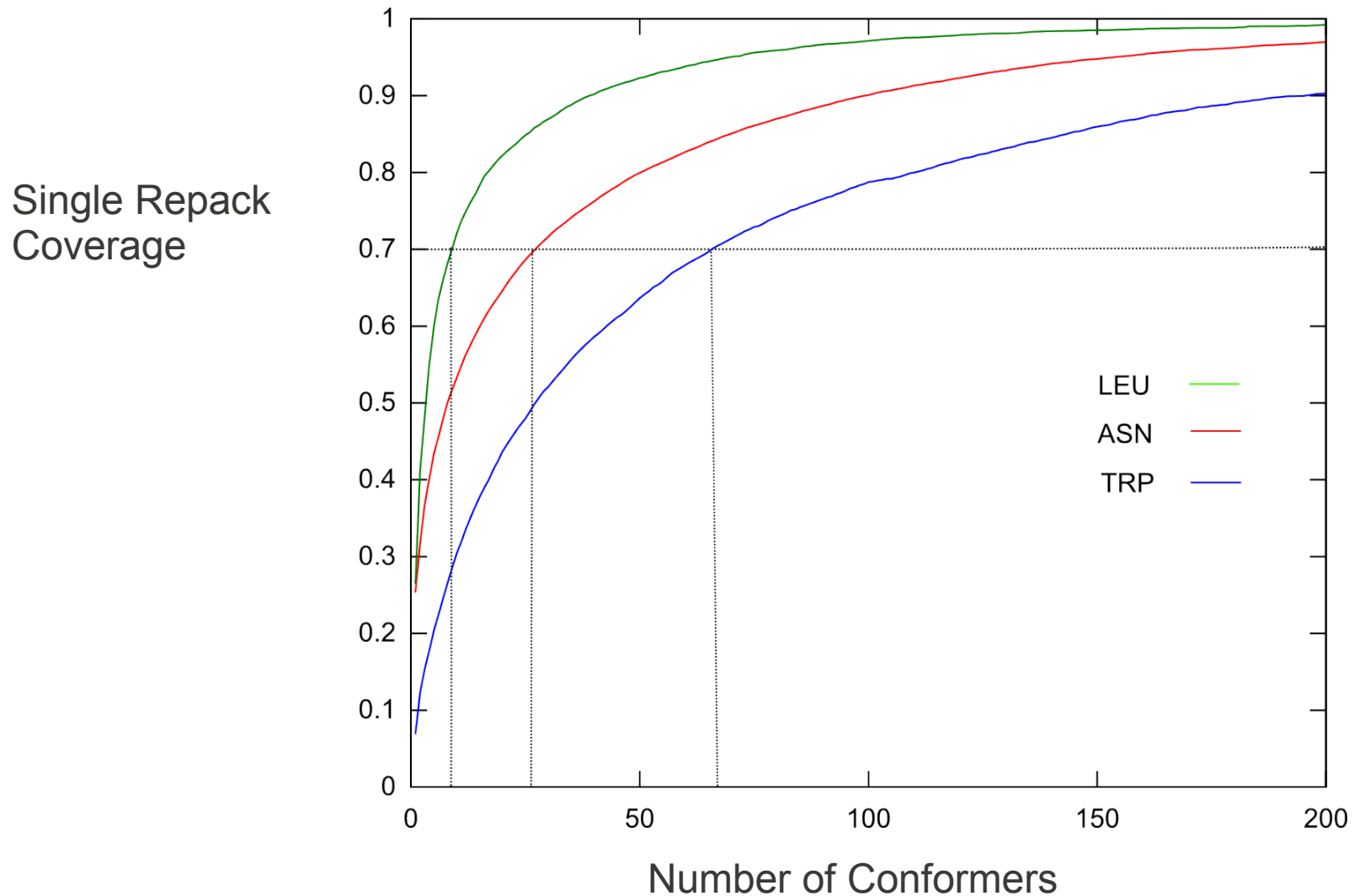
Optimization before the optimization
exceeds optimistic expectation

Alessandro Senes
IPiB Retreat 2011

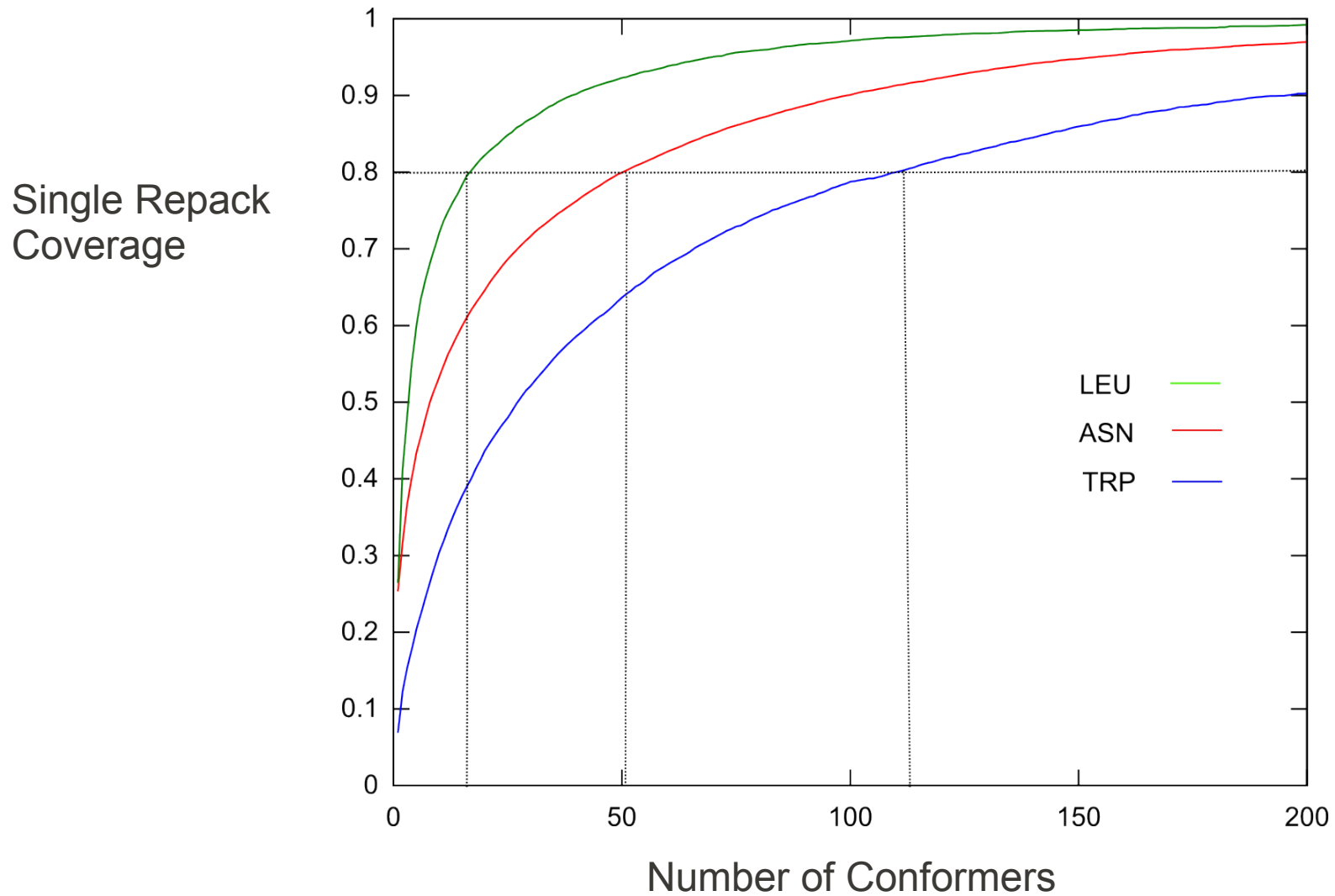
NUMBER OF CONFORMERS FOR EACH AMINO ACID TYPE



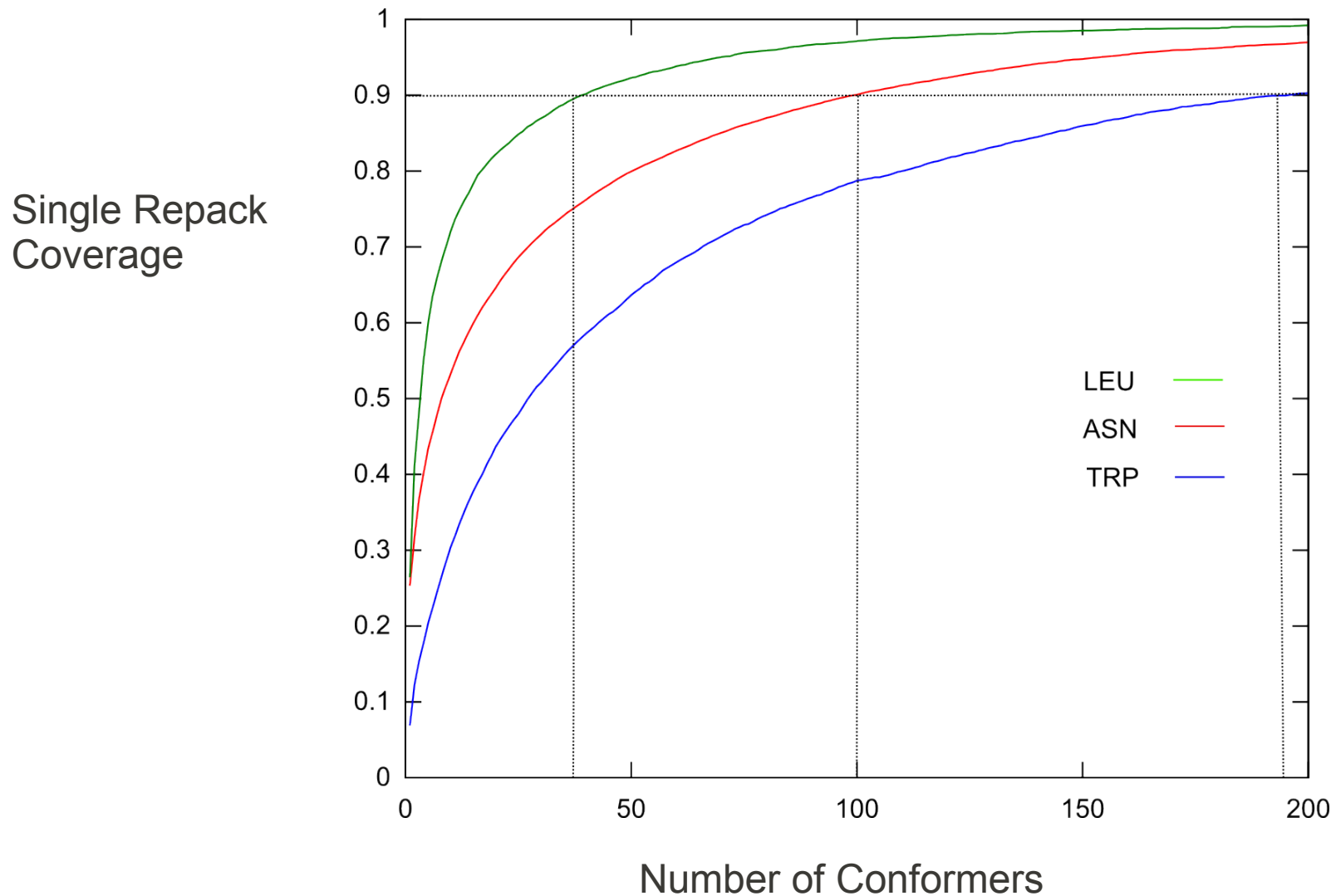
NUMBER OF CONFORMERS FOR EACH AMINO ACID TYPE



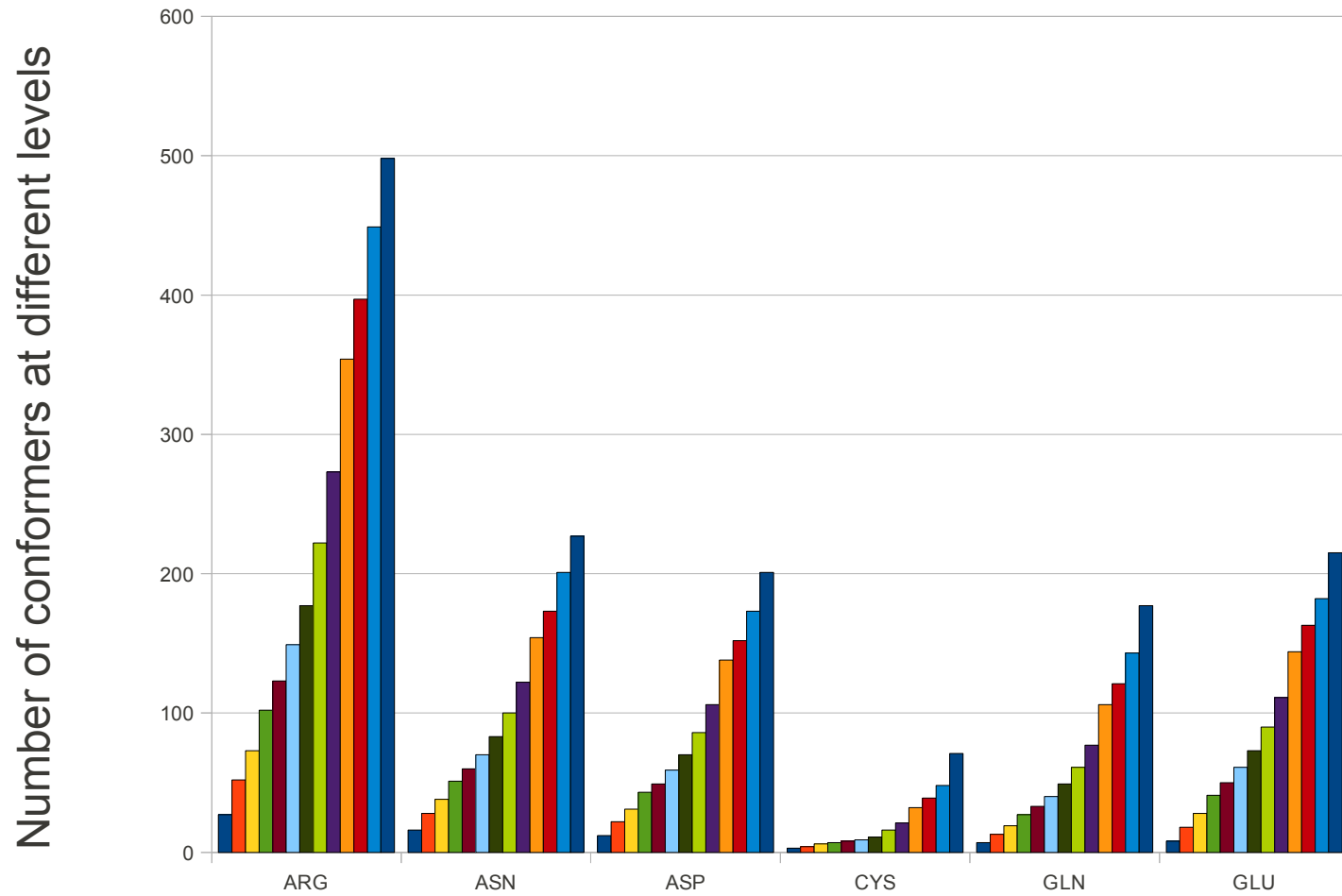
NUMBER OF CONFORMERS FOR EACH AMINO ACID TYPE



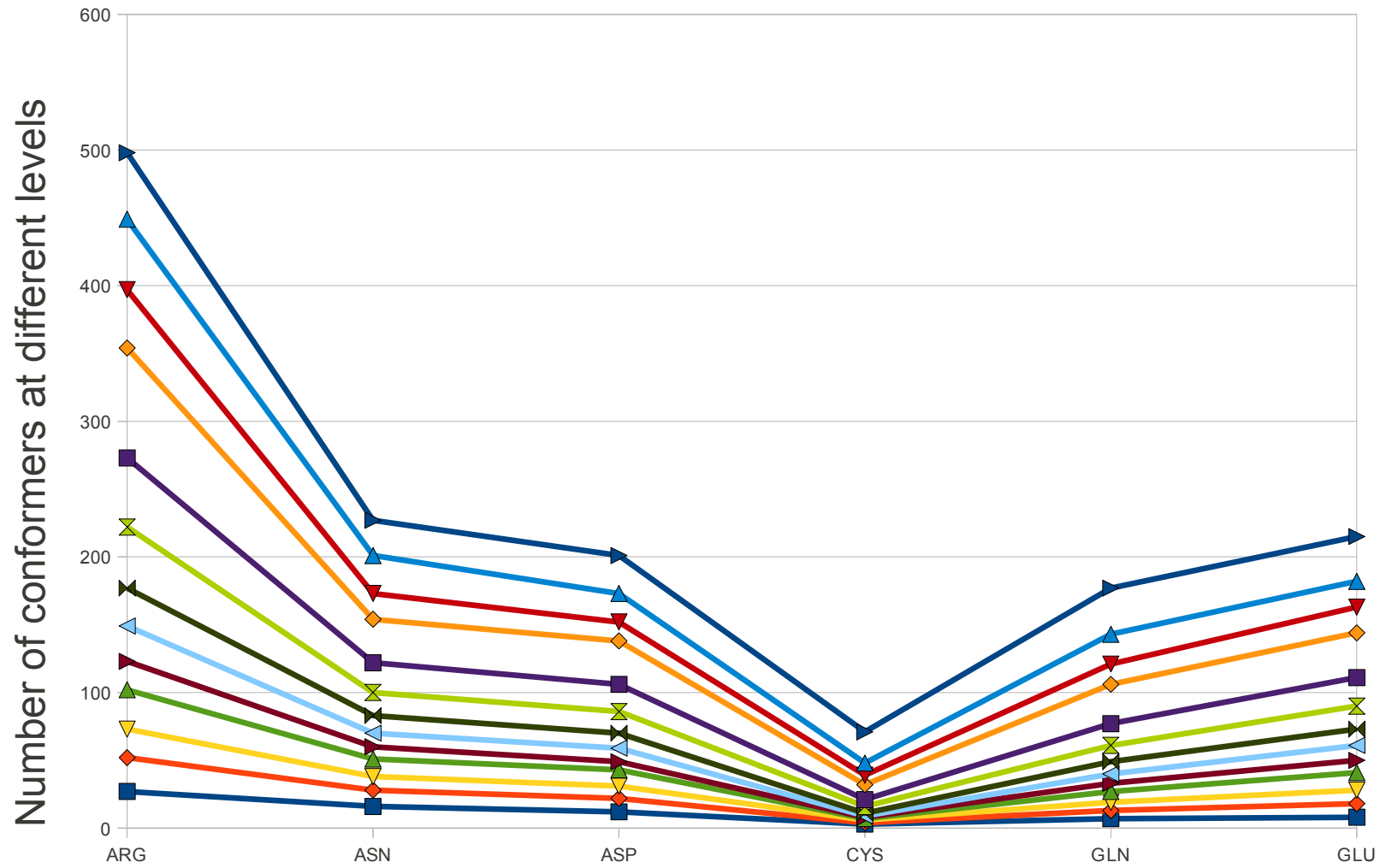
NUMBER OF CONFORMERS FOR EACH AMINO ACID TYPE



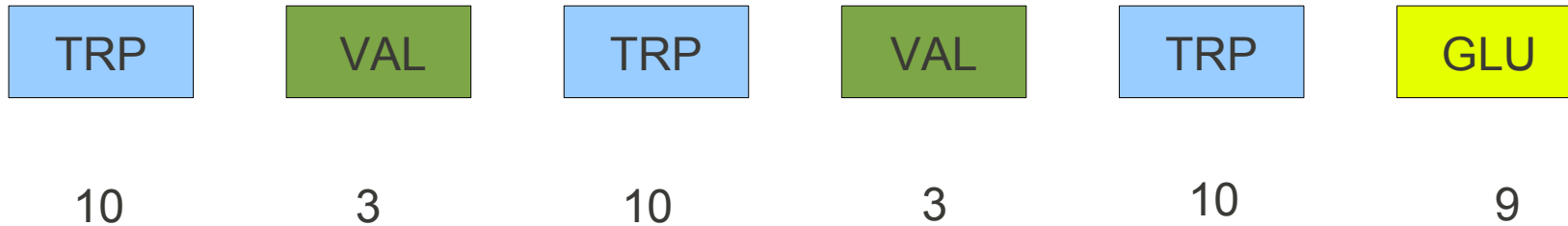
NUMBER OF CONFORMERS FOR EACH AMINO ACID TYPE



NUMBER OF CONFORMERS FOR EACH AMINO ACID TYPE



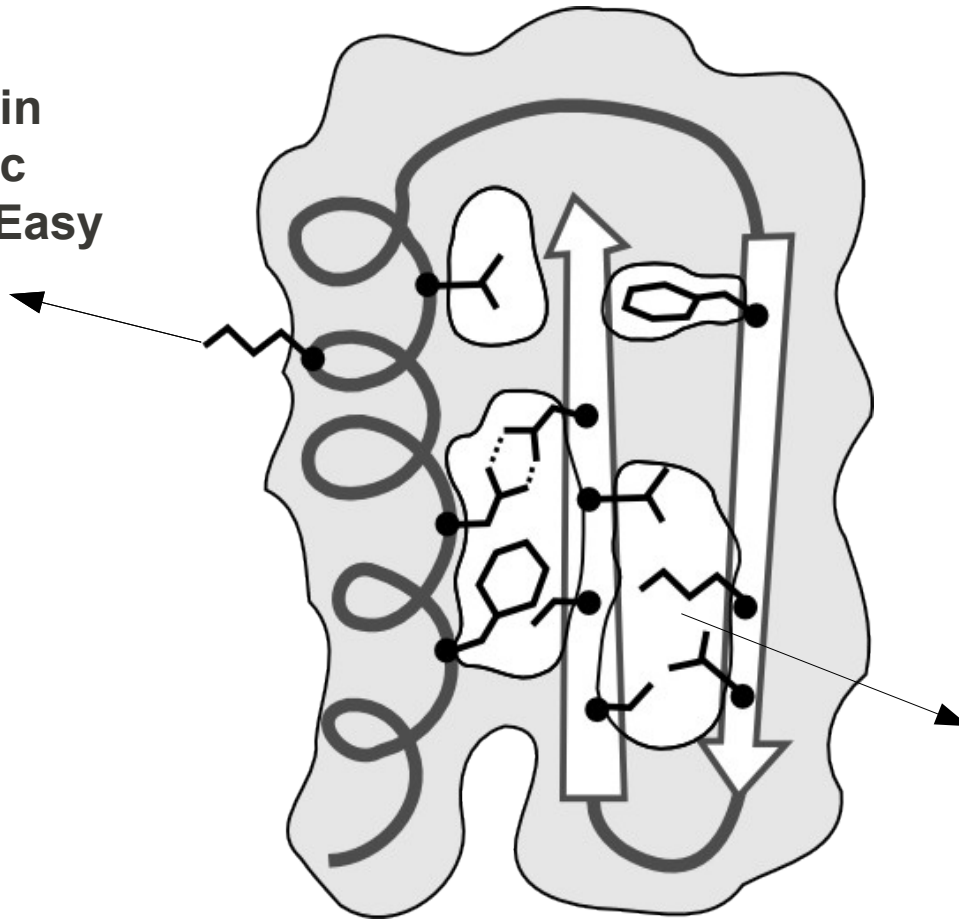
Combinatorial search space



possible conformations = $10 * 3 * 10 * 3 * 10 * 9 = 81000$

Do all positions have the same sampling requirements?

Exposed sidechain
Many isoenergetic
Conformations - Easy



Buried sidechain
Fewer isoenergetic
Conformations - Hard

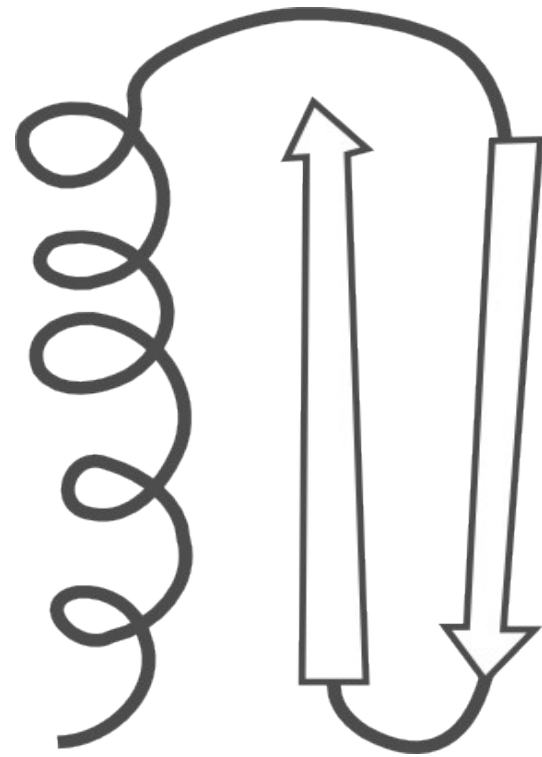
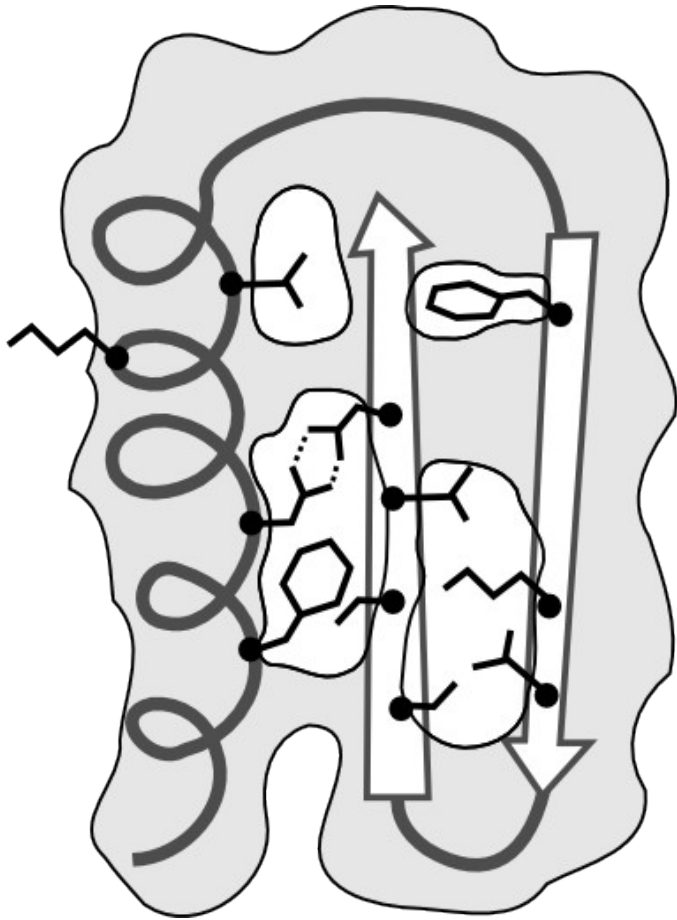
Smaller/Better search space with distributed sampling

TRP	VAL	TRP	VAL	TRP	GLU	
10	3	10	3	10	9	81000
5	2	15	5	2	9	13500 (Faster)
7	3	17	7	3	11	82467 (Better)

By moving sampling from the easy positions to the hard ones, we could be more efficient (fast) and/or achieve better energies

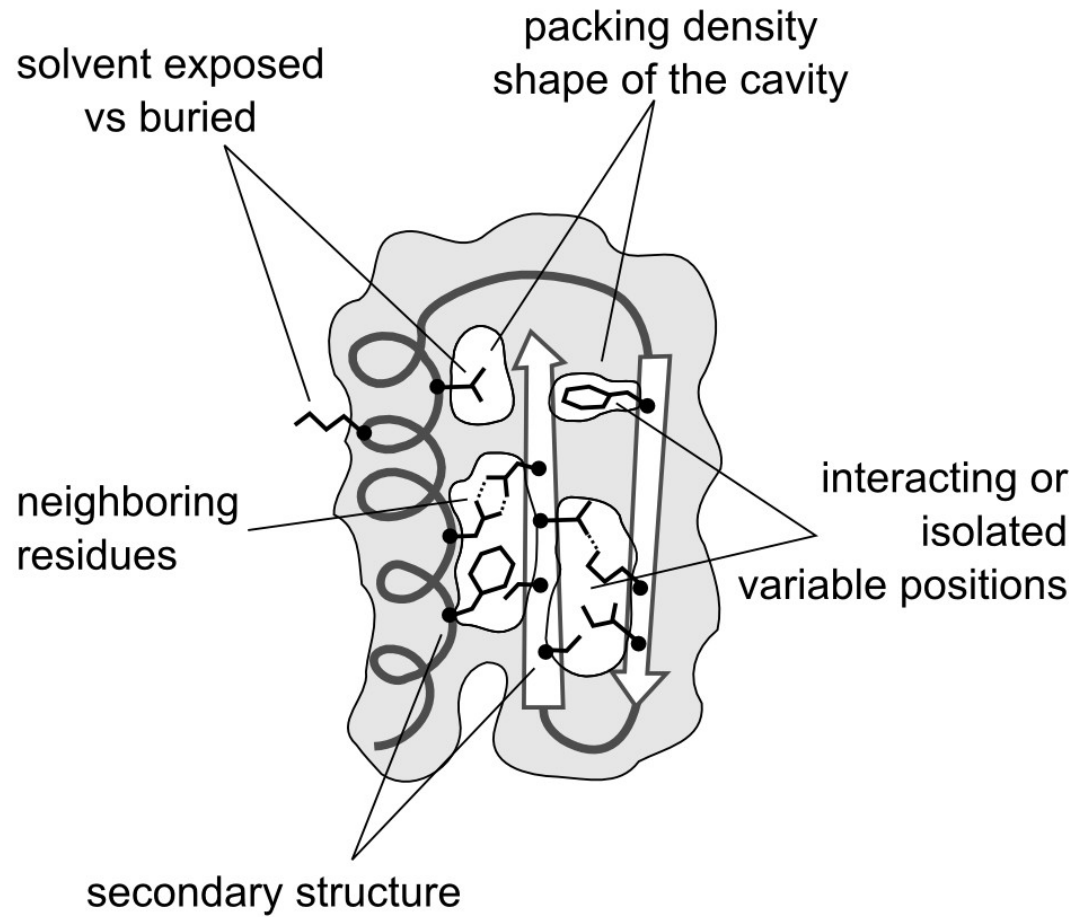
But, can we predict if a position is easy or hard?

SIDECCHAIN OPTIMIZATION



Use Machine Learning

- Information in the backbone - pattern recognition problem?
- Use machine learning to predict requirements based on features of the backbone



Goals

- Identify sampling requirements for each position on backbone
- Reduce run time, find better conformations (lower energy), or both

Issues

- Identify useful features from the backbone structure
- Identify most meaningful labeling strategy to label the dataset
- Devise the best machine learning strategy to predict the label

Label the dataset using the EnergyTable

Conformers

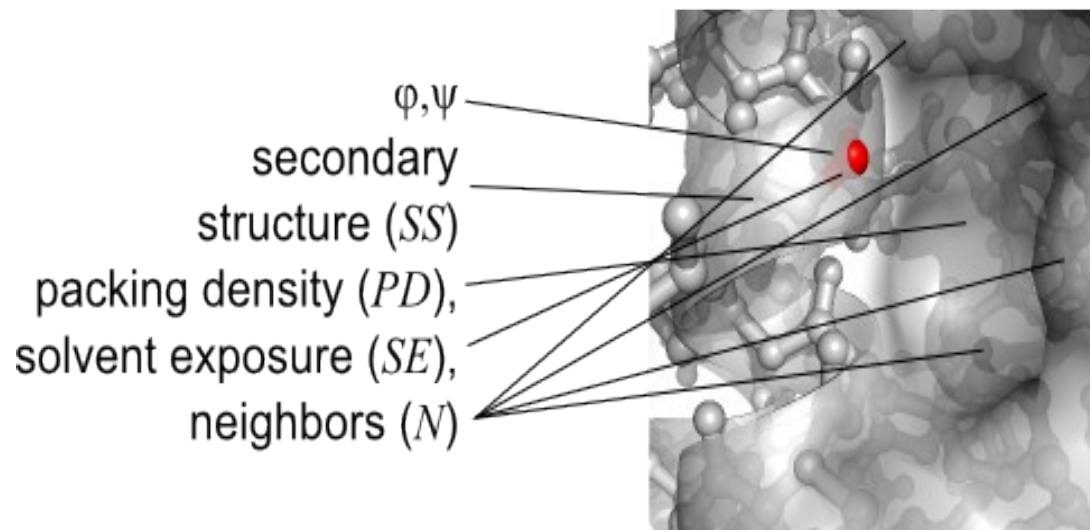
			✓		✓			2
	✓	✓			✓	✓		4
✓			✓				✓	3
	✓			✓				2
				✓				1
	✓				✓			2
			✓				✓	2
	✓							1

Label the dataset using the EnergyTable

Conformers

			✓		✓			M(2)
	✓	✓			✓	✓		E(4)
✓			✓				✓	E(3)
	✓			✓				M(2)
				✓				H(1)
	✓				✓			M(2)
			✓				✓	M(2)
	✓							H(1)

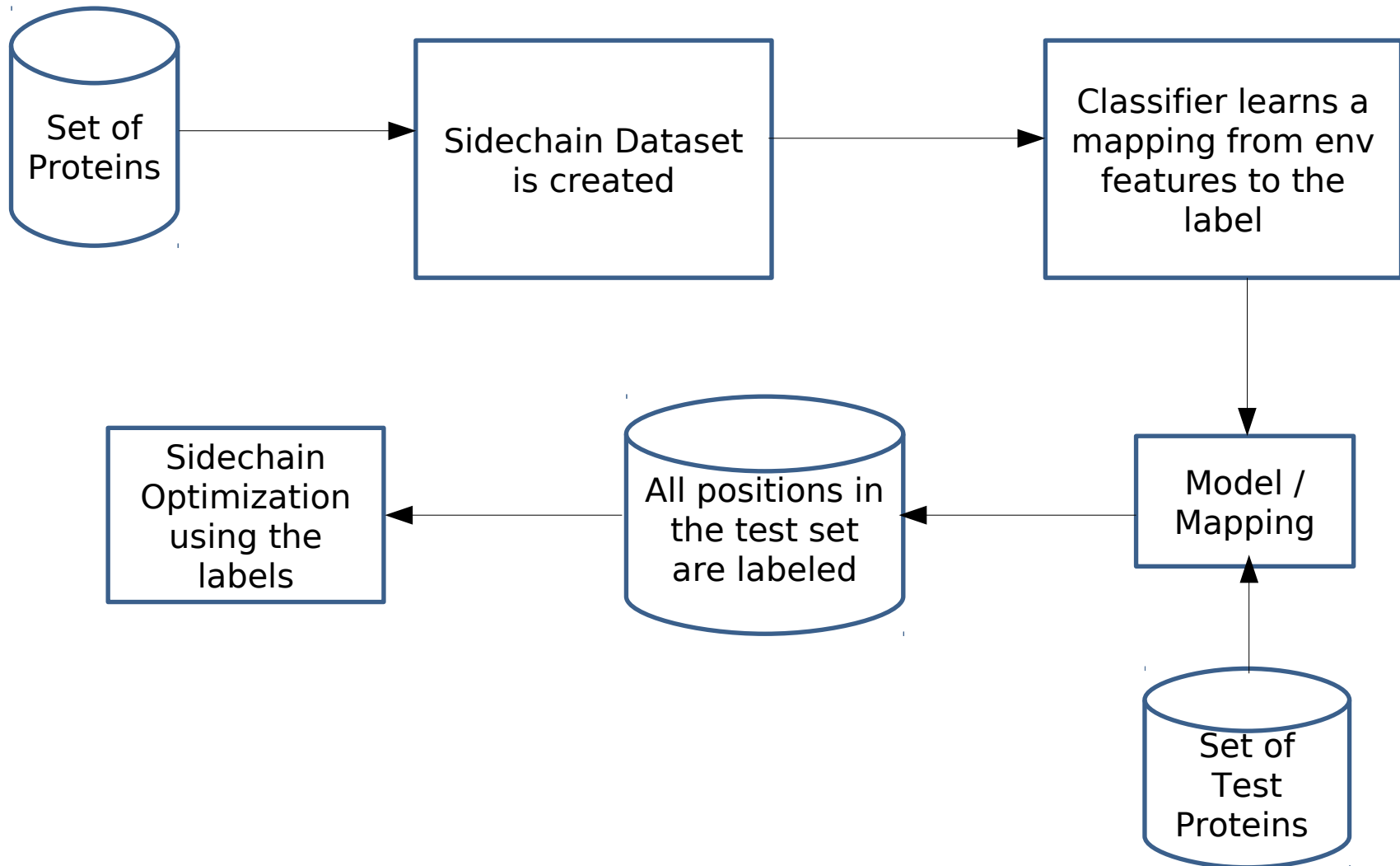
Associate each sidechain in the database
with a feature vector X and a label Y



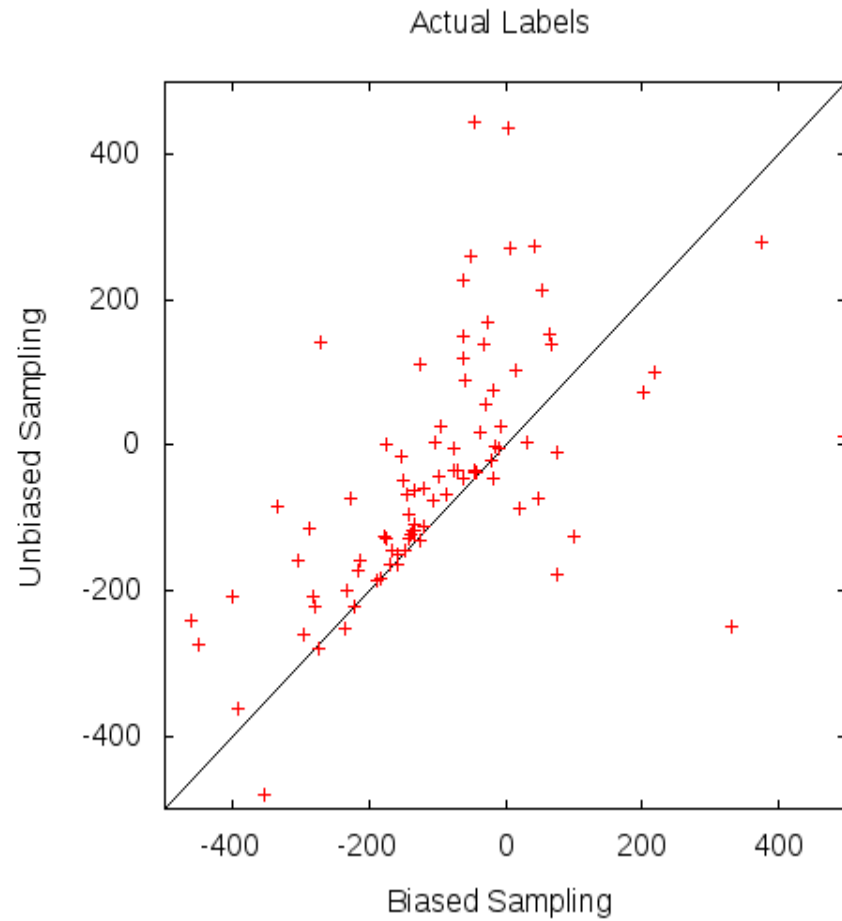
$$X = \{\varphi, \psi, SS, SE, PD, N, \dots\}$$

$$Y = \{\text{Hard, Medium, Easy}\}$$

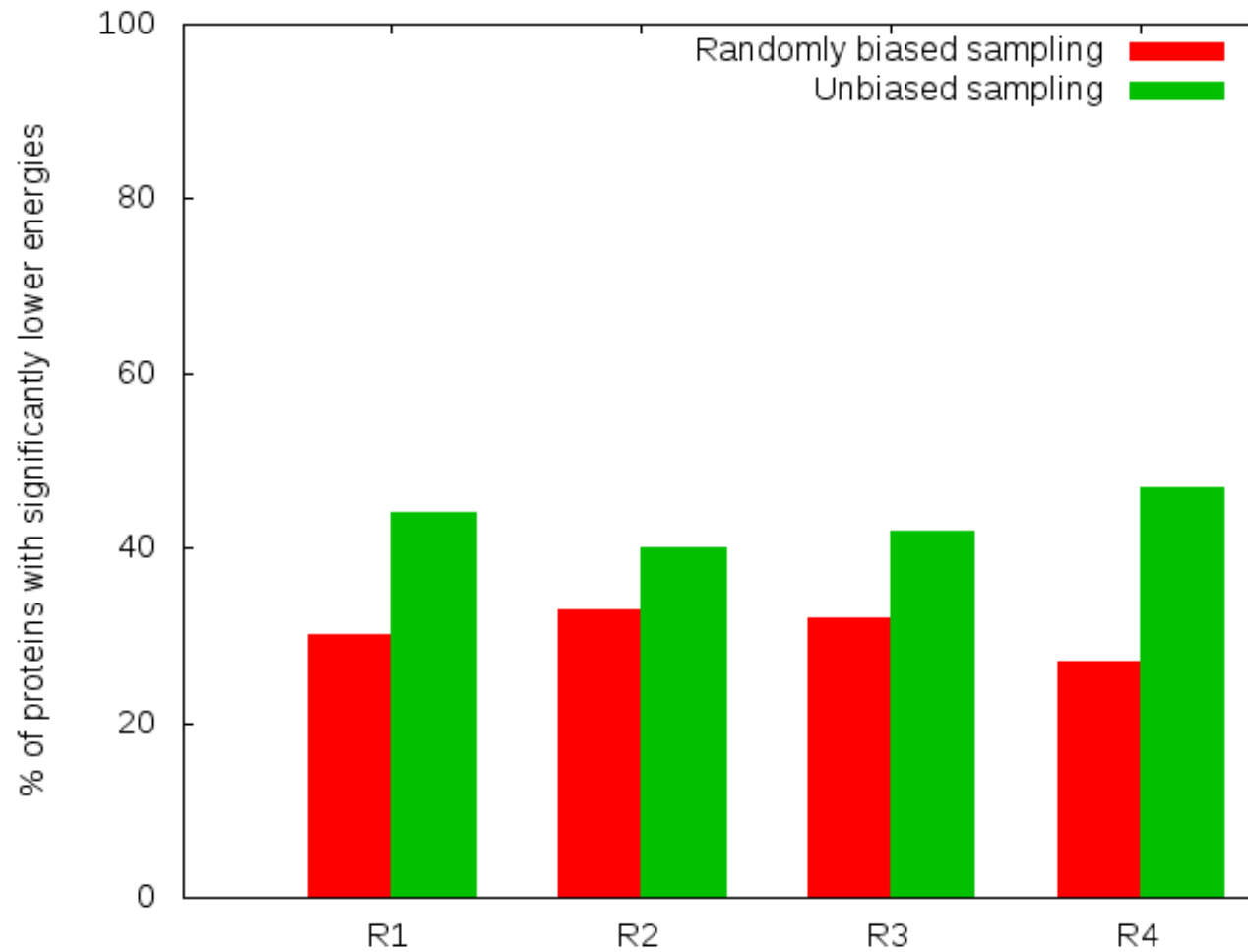
Overall Strategy



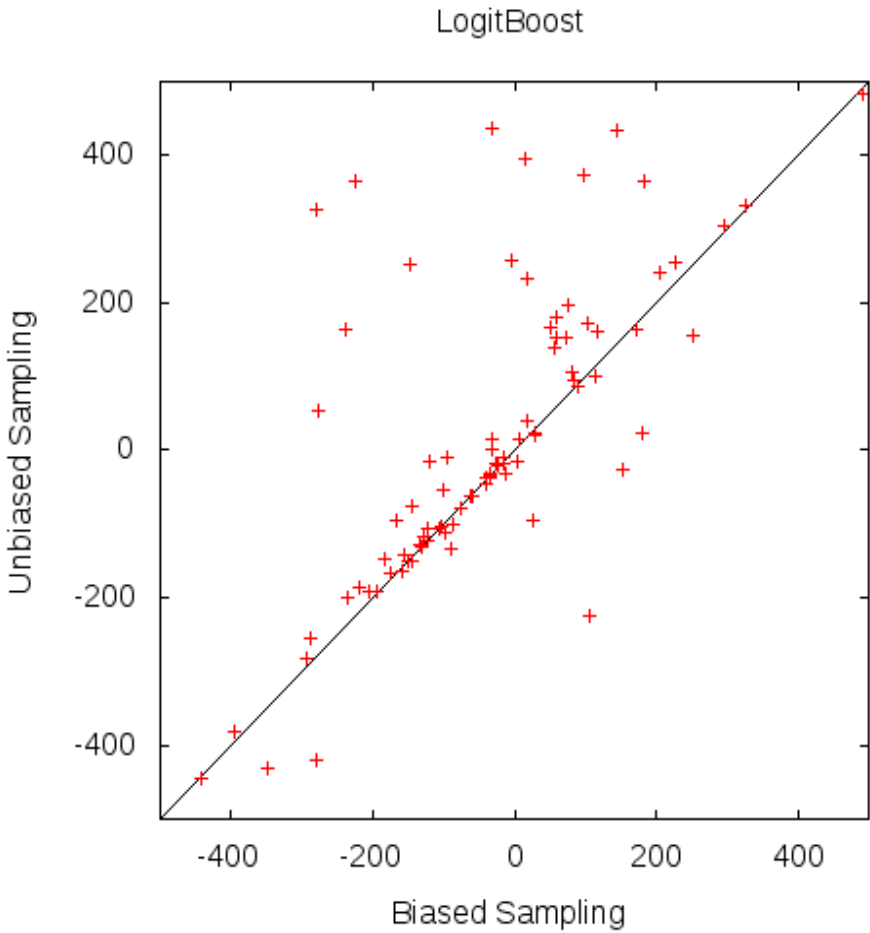
HOW EFFECTIVE ARE THE LABELS?



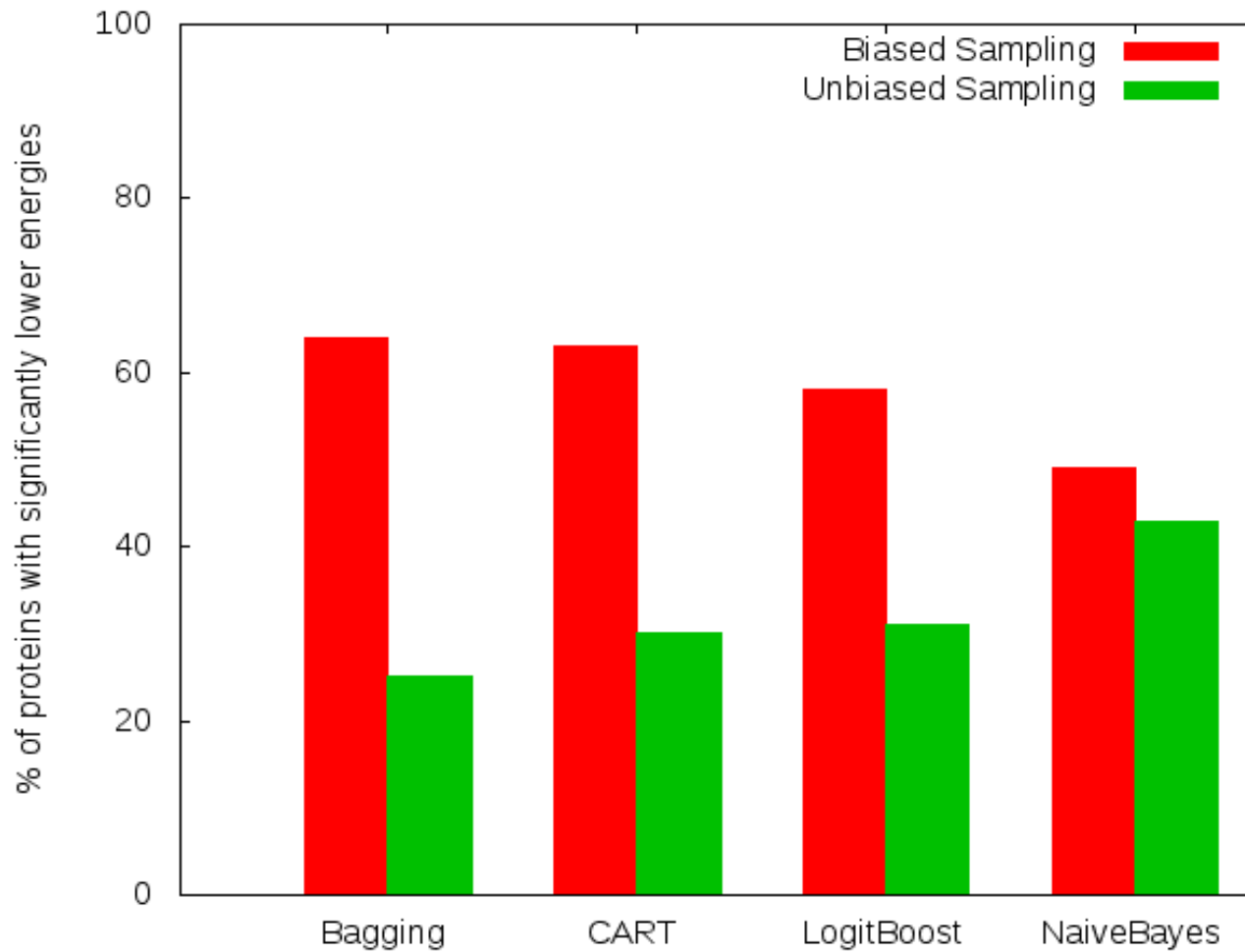
Results – Permuted labels



HOW WELL DOES MACHINE LEARNING WORK?



Results – All Classifiers



Future Research

- More features
- Alternative labeling strategies
- Optimize classifiers performance

Thank You

Alessandro Senes



Sriraam Natarajan
Ambalika Khadria
Loren LaPointe
Ben Mueller



Center of High Throughput Computing
University of Wisconsin-Madison