

Estimation of low-energy refolding paths

Michael Wolfinger

Institute for Theoretical Chemistry
University Vienna

February 21, 2006

Outline

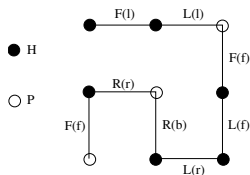
- 1 Lattice Proteins
- 2 Conformation space
- 3 Energy landscapes
- 4 Refolding paths

The HP-model

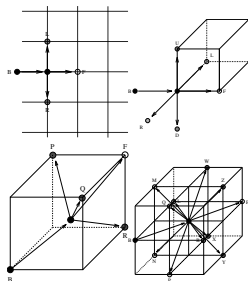
Suggested by Dill, Chan and Lau in the late 1980ies. In this *simplified model*, a conformation is a *self-avoiding walk (SAW)* on a given lattice in 2 or 3 dimensions. Each bond is a straight line, bond angles have a few discrete values. The 20 letter alphabet of amino acids (monomers) is reduced to a two letter alphabet, namely **H** and **P**. H represents **hydrophobic** monomers, P represents **hydrophilic** or *polar* monomers.

Advantages:

- lattice-independent folding algorithms
- simple energy function
- hydrophobicity can be reasonably modeled



FRLLFLF



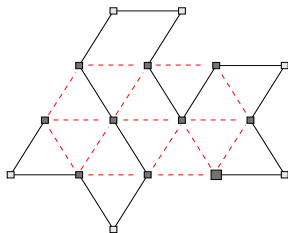
Contact Potentials

Generally, the energy function for a sequence with n residues $\mathfrak{S} = \mathfrak{s}_1\mathfrak{s}_2\dots\mathfrak{s}_n$ with $\mathfrak{s}_i \in \mathcal{A} = \{a_1, a_2, \dots, a_b\}$, the alphabet of b residues, and an overall configuration $x = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ on a lattice \mathcal{L} can be written as the sum of pair potentials

$$E(\mathfrak{S}, x) = \sum_{\substack{i < j-1 \\ |\mathbf{x}_i - \mathbf{x}_j| = 1}} \Psi[\mathfrak{s}_i, \mathfrak{s}_j].$$

Lattice proteins

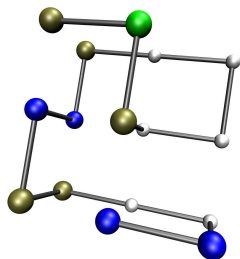
$$\mathcal{S} = \text{HPRPHHHPRPHHHPRPH} \quad n = 16$$



$$E = -15$$

	<i>H</i>	<i>P</i>
<i>H</i>	-1	0
<i>P</i>	0	0

$$\mathcal{S} = \text{NNHHPPNNPHHHHPXP} \quad n = 16$$

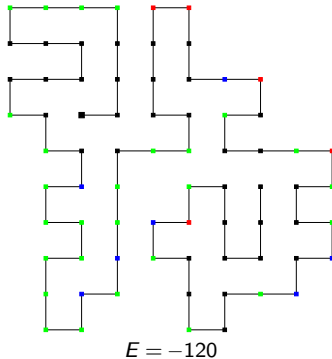


$$E = -16$$

	<i>H</i>	<i>P</i>	<i>N</i>	<i>X</i>
<i>H</i>	-4	0	0	0
<i>P</i>	0	1	-1	0
<i>N</i>	0	-1	1	0
<i>X</i>	0	0	0	0

Lattice proteins - interaction scheme II

$\mathcal{S} = \text{HHHHNNNNHHHHHHHHNHPNNNNNNNPNPNHNNHHHHXXHHPXHNHHNXNHHNPHPNHHNHHNPXNHHHHHH}$
 $n = 74$



	<i>H</i>	<i>P</i>	<i>N</i>	<i>X</i>
<i>H</i>	-4	0	0	0
<i>P</i>	0	0	-1	0
<i>N</i>	0	-1	0	0
<i>X</i>	0	0	0	0

Folding landscape - energy landscape

The energy landscape of a biopolymer molecule is a complex surface of the (free) energy versus the conformational degrees of freedom.

Number of lattice protein structures

$$c_n \sim \mu^n \cdot n^{\gamma-1}$$

problem is NP-hard

In the RNA case

$$c_n \sim 1.86^n \cdot n^{-\frac{3}{2}}$$

dynamic programming algorithms available

dim	Lattice Type	μ	γ
2	SQ	2.63820	1.34275
	TRI	4.15076	1.343
	HEX	1.84777	1.345
3	SC	4.68391	1.161
	BCC	6.53036	1.161
	FCC	10.0364	1.162

Formally, three things are needed to construct an energy landscape:

- A set X of configurations
- a symmetric neighborhood relation $\mathfrak{N} : X \times X$
- an energy function $f : X \rightarrow \mathbf{R}$

Folding landscape - energy landscape

The energy landscape of a biopolymer molecule is a complex surface of the (free) energy versus the conformational degrees of freedom.

Number of lattice protein structures

$$c_n \sim \mu^n \cdot n^{\gamma-1}$$

problem is NP-hard

In the RNA case

$$c_n \sim 1.86^n \cdot n^{-\frac{3}{2}}$$

dynamic programming algorithms available

dim	Lattice Type	μ	γ
2	SQ	2.63820	1.34275
	TRI	4.15076	1.343
	HEX	1.84777	1.345
3	SC	4.68391	1.161
	BCC	6.53036	1.161
	FCC	10.0364	1.162

Formally, three things are needed to construct an energy landscape:

- A set X of configurations
- a symmetric neighborhood relation $\mathfrak{N} : X \times X$
- an energy function $f : X \rightarrow \mathbf{R}$

Folding landscape - energy landscape

The energy landscape of a biopolymer molecule is a complex surface of the (free) energy versus the conformational degrees of freedom.

Number of lattice protein structures

$$c_n \sim \mu^n \cdot n^{\gamma-1}$$

problem is NP-hard

In the RNA case

$$c_n \sim 1.86^n \cdot n^{-\frac{3}{2}}$$

dynamic programming algorithms available

dim	Lattice Type	μ	γ
2	SQ	2.63820	1.34275
	TRI	4.15076	1.343
	HEX	1.84777	1.345
3	SC	4.68391	1.161
	BCC	6.53036	1.161
	FCC	10.0364	1.162

Formally, three things are needed to construct an energy landscape:

- A set X of configurations
- a symmetric neighborhood relation $\mathfrak{N} : X \times X$
- an energy function $f : X \rightarrow \mathbf{R}$

Folding landscape - energy landscape

The energy landscape of a biopolymer molecule is a complex surface of the (free) energy versus the conformational degrees of freedom.

Number of lattice protein structures

$$c_n \sim \mu^n \cdot n^{\gamma-1}$$

problem is NP-hard

In the RNA case

$$c_n \sim 1.86^n \cdot n^{-\frac{3}{2}}$$

dynamic programming algorithms available

dim	Lattice Type	μ	γ
2	SQ	2.63820	1.34275
	TRI	4.15076	1.343
	HEX	1.84777	1.345
3	SC	4.68391	1.161
	BCC	6.53036	1.161
	FCC	10.0364	1.162

Formally, three things are needed to construct an energy landscape:

- A set X of configurations
- a symmetric neighborhood relation $\mathfrak{N} : X \times X$
- an energy function $f : X \rightarrow \mathbb{R}$

Folding landscape - energy landscape

The energy landscape of a biopolymer molecule is a complex surface of the (free) energy versus the conformational degrees of freedom.

Number of lattice protein structures

$$c_n \sim \mu^n \cdot n^{\gamma-1}$$

problem is NP-hard

In the RNA case

$$c_n \sim 1.86^n \cdot n^{-\frac{3}{2}}$$

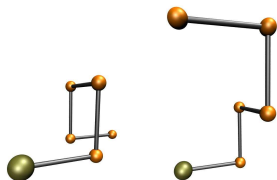
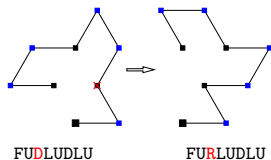
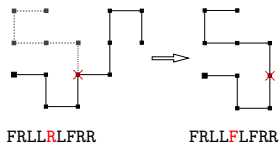
dynamic programming algorithms available

dim	Lattice Type	μ	γ
2	SQ	2.63820	1.34275
	TRI	4.15076	1.343
	HEX	1.84777	1.345
3	SC	4.68391	1.161
	BCC	6.53036	1.161
	FCC	10.0364	1.162

Formally, three things are needed to construct an energy landscape:

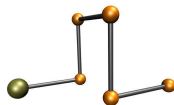
- A set X of configurations
- a symmetric neighborhood relation $\mathfrak{N} : X \times X$
- an energy function $f : X \rightarrow \mathbf{R}$

The move set



FU**U**RR

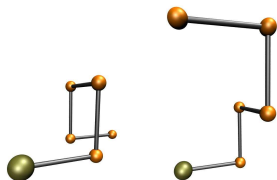
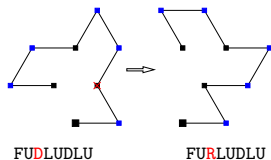
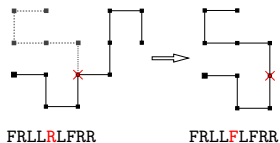
FU**D**RR



FU**D**LR

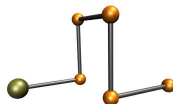
- For each move there must be an inverse move
- Resulting structure must be in X
- Move set must be *ergodic*

The move set



FUURR

FUDRR



FUDLR

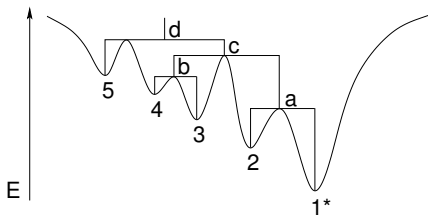
- For each move there must be an inverse move
- Resulting structure must be in X
- Move set must be *ergodic*

Energy barriers and barrier trees

Some topological definitions:

A structure is a

- **local minimum** if its energy is lower than the energy of **all** neighbors
- **local maximum** if its energy is higher than the energy of **all** neighbors
- **saddle point** if there are at least two local minima that can be reached by a downhill walk starting at this point



We further define

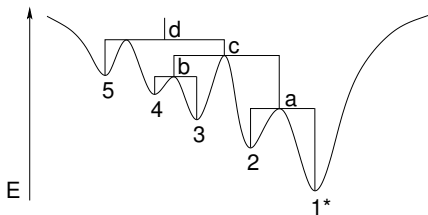
- a **walk** between two conformations x and y as a list of conformations $x = x_1 \dots x_{m+1} = y$ such that $\forall 1 \leq i \leq m : \mathfrak{N}(x_i, x_{i+1})$
- the **lower part** of the energy landscape (written as $X^{\leq \eta}$) as *all* conformations x such that $E(\mathcal{G}, x) \leq \eta$ (with a predefined threshold η).

Energy barriers and barrier trees

Some topological definitions:

A structure is a

- **local minimum** if its energy is lower than the energy of **all** neighbors
- **local maximum** if its energy is higher than the energy of **all** neighbors
- **saddle point** if there are at least two local minima that can be reached by a downhill walk starting at this point



We further define

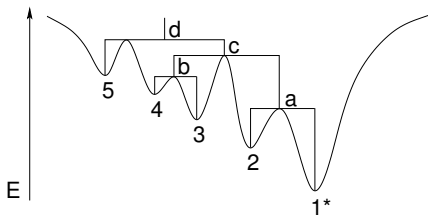
- a **walk** between two conformations x and y as a list of conformations $x = x_1 \dots x_{m+1} = y$ such that $\forall 1 \leq i \leq m : \mathfrak{N}(x_i, x_{i+1})$
- the **lower part** of the energy landscape (written as $X^{\leq \eta}$) as *all* conformations x such that $E(\mathcal{G}, x) \leq \eta$ (with a predefined threshold η).

Energy barriers and barrier trees

Some topological definitions:

A structure is a

- **local minimum** if its energy is lower than the energy of **all** neighbors
- **local maximum** if its energy is higher than the energy of **all** neighbors
- **saddle point** if there are at least two local minima that can be reached by a downhill walk starting at this point



We further define

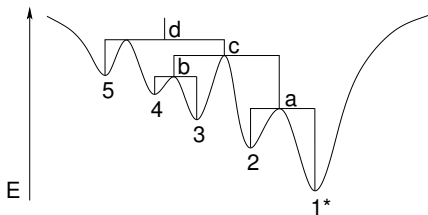
- a **walk** between two conformations x and y as a list of conformations $x = x_1 \dots x_{m+1} = y$ such that $\forall 1 \leq i \leq m : \mathfrak{N}(x_i, x_{i+1})$
- the **lower part** of the energy landscape (written as $X^{\leq \eta}$) as *all* conformations x such that $E(\mathcal{G}, x) \leq \eta$ (with a predefined threshold η).

Energy barriers and barrier trees

Some topological definitions:

A structure is a

- **local minimum** if its energy is lower than the energy of **all** neighbors
- **local maximum** if its energy is higher than the energy of **all** neighbors
- **saddle point** if there are at least two local minima that can be reached by a downhill walk starting at this point



We further define

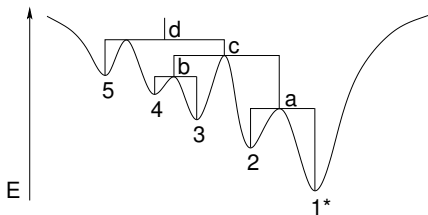
- a **walk** between two conformations x and y as a list of conformations $x = x_1 \dots x_{m+1} = y$ such that $\forall 1 \leq i \leq m : \mathfrak{N}(x_i, x_{i+1})$
- the **lower part** of the energy landscape (written as $X^{\leq \eta}$) as *all* conformations x such that $E(\mathcal{G}, x) \leq \eta$ (with a predefined threshold η).

Energy barriers and barrier trees

Some topological definitions:

A structure is a

- **local minimum** if its energy is lower than the energy of **all** neighbors
- **local maximum** if its energy is higher than the energy of **all** neighbors
- **saddle point** if there are at least two local minima that can be reached by a downhill walk starting at this point



We further define

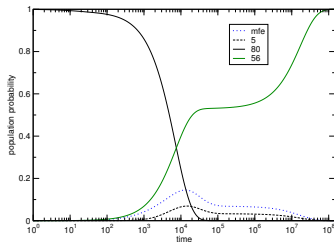
- a **walk** between two conformations x and y as a list of conformations $x = x_1 \dots x_{m+1} = y$ such that $\forall 1 \leq i \leq m : \mathfrak{N}(x_i, x_{i+1})$
- the **lower part** of the energy landscape (written as $X^{\leq \eta}$) as *all* conformations x such that $E(\mathcal{G}, x) \leq \eta$ (with a predefined threshold η).

Information from the barrier tree

- Local minima
- Saddle points
- Barrier heights
- Gradient basins
- Partition functions and free energies of (gradient) basins

This information can be used to approximate the dynamics of biopolymers, i.e. transition rates between different macrostates (basins in the barrier tree)

$$\blacksquare r_{\beta\alpha} = \Gamma_{\beta\alpha} \exp\left(-\frac{(E_{\beta\alpha}^* - G_{\alpha})}{kT}\right)$$



The lower part of the energy landscape

Two conformations x and y are mutually accessible at the level η (written as $x \xleftrightarrow{\eta} y$) if there is a walk from x to y such that all conformations z in the walk satisfy $E(\mathcal{G}, z) \leq \eta$. The *saddle height* $\hat{f}(x, y)$ of x and y is defined by

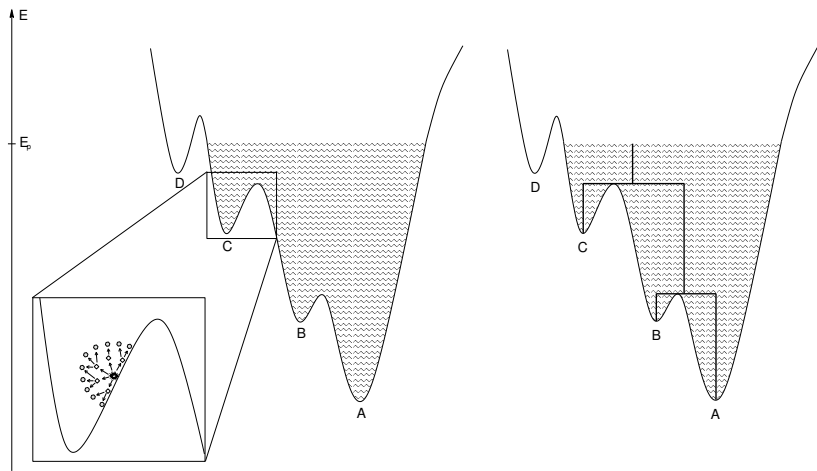
$$\hat{f}(x, y) = \min\{\eta \mid x \xleftrightarrow{\eta} y\}$$

Given the set of all local minima $X_{\min}^{\leq \eta}$ below threshold η , the lower energy part $X^{\leq \eta}$ of the energy landscape can alternatively be written as

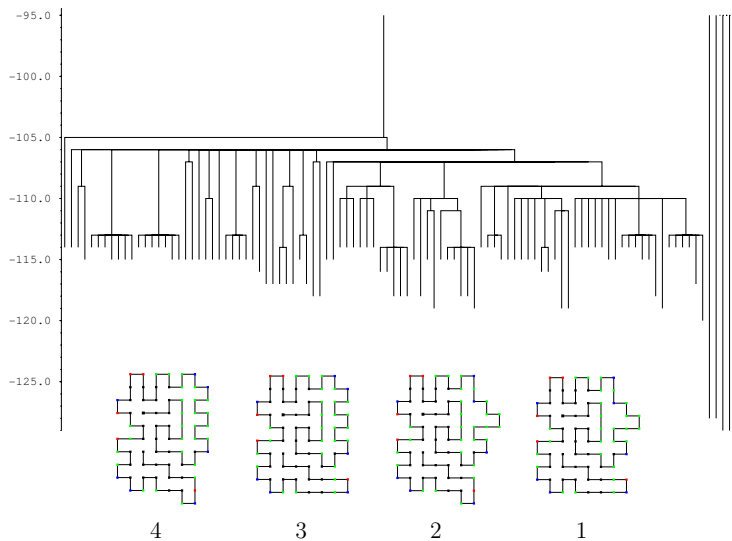
$$X^{\leq \eta} = \{y \mid \exists x \in X_{\min}^{\leq \eta} : \hat{f}(x, y) \leq \eta\}$$

Given a restricted set of low-energy conformations, X_{init} , and a reasonable value for η , the lower part of the energy landscape can be calculated.

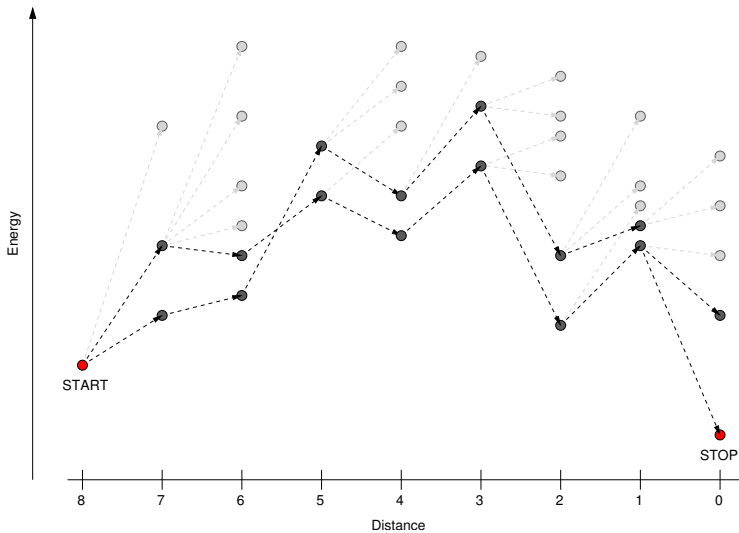
The Flooder approach



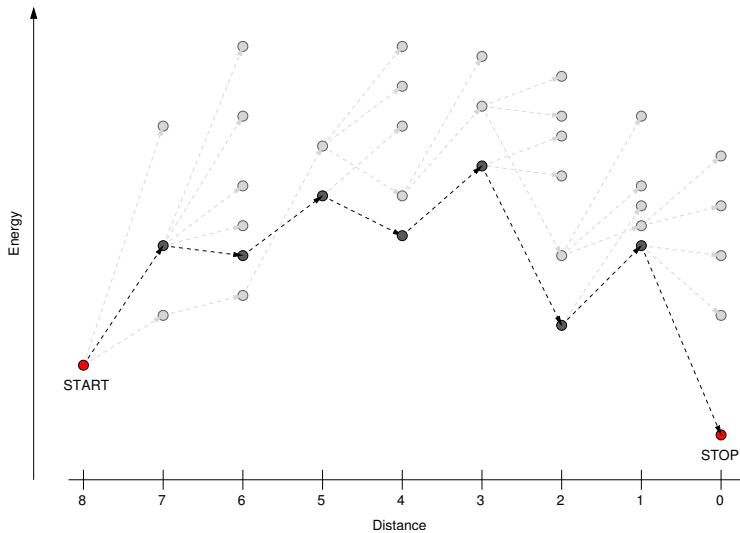
!Connected



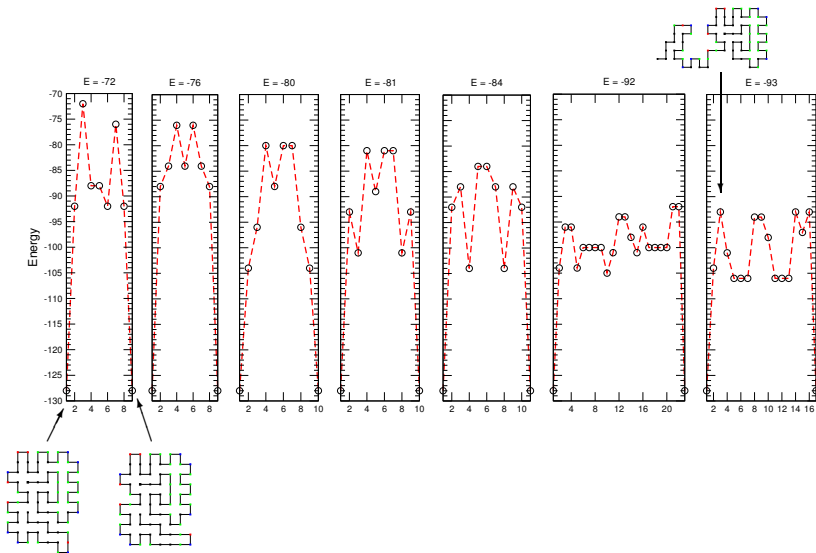
LatticePath - illustration



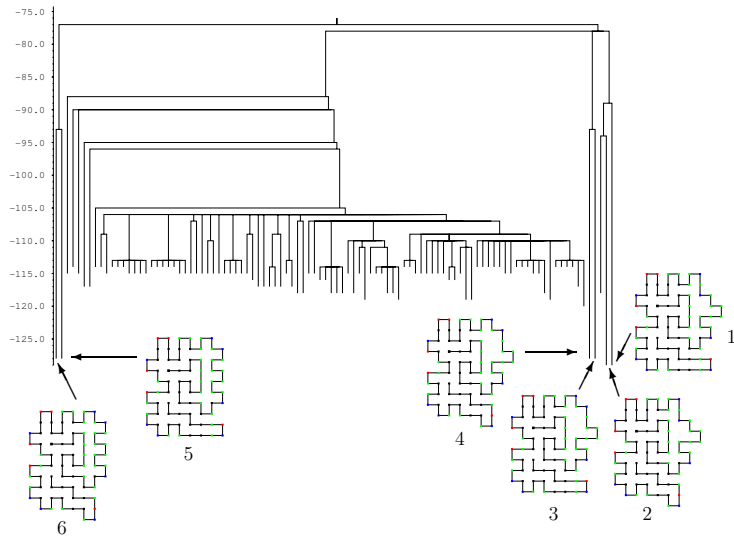
LatticePath - illustration



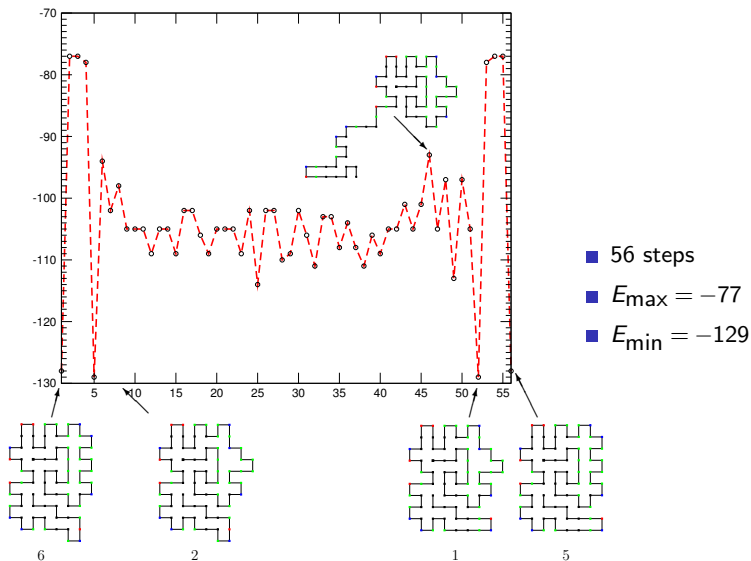
Refolding profiles



Connected!



A longer refolding profile



Conclusion

- **Discrete models** allow a detailed study of the energy surface.
- **Barrier trees** approximate the landscape topology and folding kinetics.
- A **heuristic approach** allows to sample low-energy refolding paths between different structures
- This **newly generated framework** provides a powerful method for further refinement of biopolymer folding landscapes.

Conclusion

- **Discrete models** allow a detailed study of the energy surface.
- **Barrier trees** approximate the landscape topology and folding kinetics.
- A **heuristic approach** allows to sample low-energy refolding paths between different structures
- This **newly generated framework** provides a powerful method for further refinement of biopolymer folding landscapes.

Conclusion

- **Discrete models** allow a detailed study of the energy surface.
- **Barrier trees** approximate the landscape topology and folding kinetics.
- A **heuristic approach** allows to sample low-energy refolding paths between different structures
- This **newly generated framework** provides a powerful method for further refinement of biopolymer folding landscapes.

Conclusion

- **Discrete models** allow a detailed study of the energy surface.
- **Barrier trees** approximate the landscape topology and folding kinetics.
- A **heuristic approach** allows to sample low-energy refolding paths between different structures
- This **newly generated framework** provides a powerful method for further refinement of biopolymer folding landscapes.

Thanks

Sebastian Will

Rolf Backofen

Peter Stadler

Ivo Hofacker

Christoph Flamm

The electric, without whom this would not be possible