

Finding Non-direct RNA refolding paths (in silico)

Ronny Lorenz

`ronny@tbi.univie.ac.at`

Institute for Theoretical Chemistry
University of Vienna

Bled, Slovenia, February 16, 2009

1 Basics

- RNA secondary structure energy landscape and neighborhood
- obtaining direct refolding paths
- taking detours

1 Basics

- RNA secondary structure energy landscape and neighborhood
- obtaining direct refolding paths
- taking detours

2 finding meshpoint structures

- Non-deterministic Method(s)
- Deterministic Method

- 1 Basics
 - RNA secondary structure energy landscape and neighborhood
 - obtaining direct refolding paths
 - taking detours
- 2 finding meshpoint structures
 - Non-deterministic Method(s)
 - Deterministic Method
- 3 Results
 - Artificial RNA switches
 - add A-riboswitch

Secondary structure energy landscape

$$\text{Energy Landscape} = \{\mathcal{X}, \mathcal{H}, f\}$$

\mathcal{X} ... set of configurations

\mathcal{H} ... topological structure on \mathcal{X}

f ... fitness/energy function with $f : \mathcal{X} \rightarrow \mathbb{R}$

Secondary structure energy landscape

$$\text{Energy Landscape} = \{\mathcal{X}, \mathcal{H}, f\}$$

\mathcal{X} ... set of configurations (*RNA secondary structures*)

\mathcal{H} ... topological structure on \mathcal{X}

f ... fitness/energy function with $f : \mathcal{X} \rightarrow \mathbb{R}$

Secondary structure energy landscape

$$\text{Energy Landscape} = \{\mathcal{X}, \mathcal{H}, f\}$$

\mathcal{X} ... set of configurations (*RNA secondary structures*)

\mathcal{H} ... topological structure on \mathcal{X} (*insert/delete bp*)

f ... fitness/energy function with $f : \mathcal{X} \rightarrow \mathbb{R}$

Secondary structure energy landscape

$$\text{Energy Landscape} = \{\mathcal{X}, \mathcal{H}, f\}$$

\mathcal{X} ... set of configurations (*RNA secondary structures*)

\mathcal{H} ... topological structure on \mathcal{X} (*insert/delete bp*)

f ... fitness/energy function with $f : \mathcal{X} \rightarrow \mathbb{R}$ (*free energy*)

Folding path p between two secondary structures s_1 and s_2

sequence of secondary structures where:

- $p[0] = s_1$
- $p[n] = s_2$
- $\forall p[i] : p[i + 1]$ is reachable from $p[i]$ by applying one single move from the moveset

direct paths:

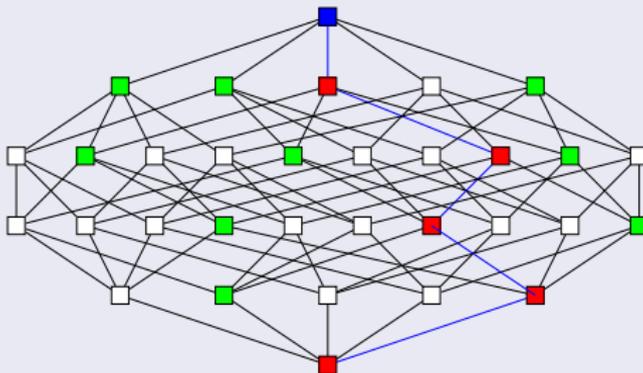
- $n = 1 + d_{BP}(s_1, s_2)$

Energy barrier between two secondary structures along a folding path

- $B(p) = \max_{0 \leq i \leq n} E(p[i]) - E(p[0])$
- finding the lowest energy barrier for any path is too costly in terms of computational effort
- heuristics are applied to obtain a good estimate of the lowest barrier between s_1 and s_2

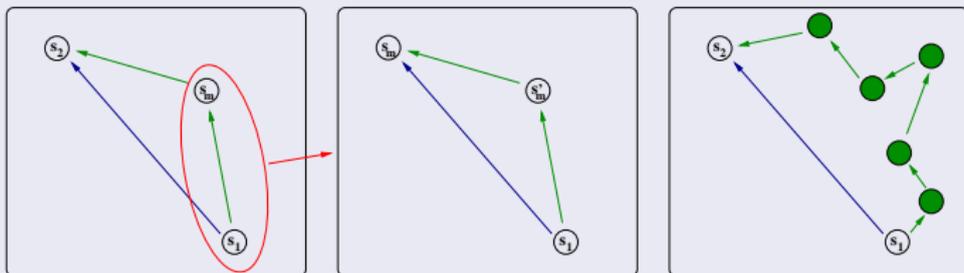
Morgan Higgs et al. approach

- 1 let the current structure A be s_1
- 2 find (i,j) in s_2 which has least incompatible pairs in A
- 3 remove incompatible pairs in A and insert (i,j)
- 4 if $A \neq s_2$ go to (2)



Approach with detours

- 1 let $A = s_1$ and $B = s_2$
- 2 generate best direct folding path p_d between A and B
- 3 find x other structures s_{x_m} (meshpoint x)
- 4 generate best direct paths p_{x_1} from A to s_m and p_{x_2} from s_m to B
- 5 get complete paths $p_{x_c} = p_{x_1} + p_{x_2}$
- 6 accept if $B(p_{x_c}) < B(p_d)$
- 7 if deeper iteration required continue from (2) to level p_{x_1} and p_{x_2}
- 8 otherwise return p_{x_c} with lowest barrier



So, how to obtain meshpoint structures s_{x_m} ?

Non-deterministic meshpoint generation

Monte Carlo with Metropolis rule

- take s_b that exhibits the barrier along the known path between s_1 and s_2
- start a random walk from s_b
- probability to move from structure s_i to one of its neighbors s_j :

$$p(s_i \rightarrow s_j) = \frac{1}{N} \begin{cases} \exp^{-\beta \cdot (E(s_i) - E(s_j))} & \text{if } E(s_i) < E(s_j) \\ 1 & \text{else} \end{cases} \quad (1)$$

- probability to reject the next move at all p_r :

$$p_r = 1 - \sum_j p(s_i \rightarrow s_j) \quad (2)$$

- rejectionless selection of neighbor s_j with probability

$$p_{rl}(s_i \rightarrow s_j) = \frac{1}{\sum_j p(s_i \rightarrow s_j)} \cdot p(s_i \rightarrow s_j) \quad (3)$$

Move penalties

- walks backward on trajectory can be penalized
- penalty contribution of $2 \cdot kT$ added to free energy of previously visited neighbor

Stop criterion

- if all neighbors show heigher free energy than current state, MC stops with probability

$$p_{\text{stop}} = 1 - e^{-\beta \cdot (E(s_i) - \min_j E(s_j))} \quad (4)$$

Simulated annealing

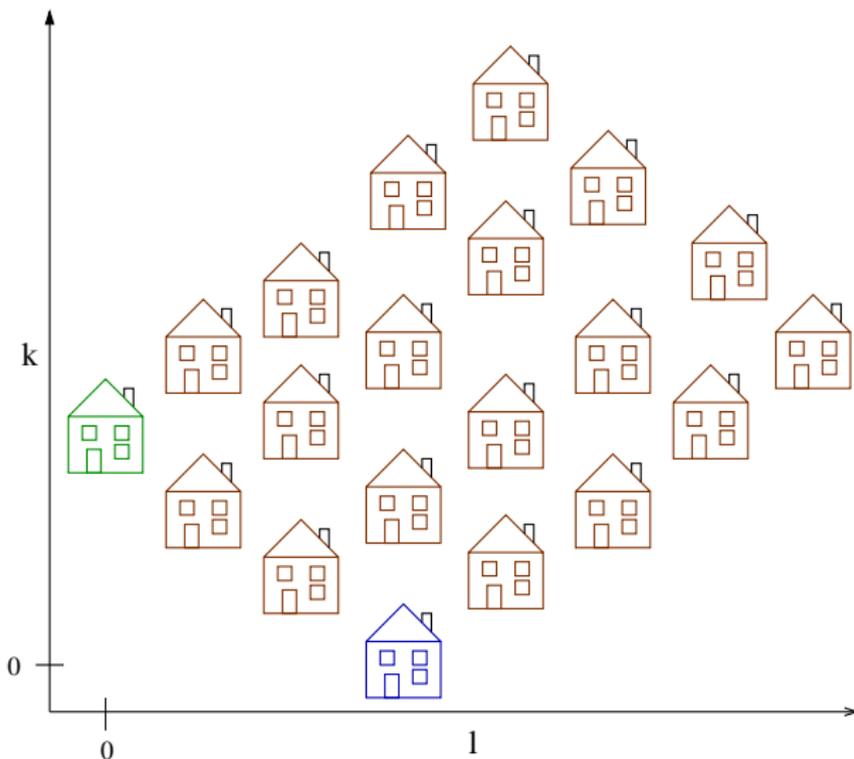
- given a start and stop temperatures t_{start} , t_{stop} and cooling rate r
- acceptance probabilites of next states is altered according to falling temperature of the system

Deterministic meshpoint generation

k, l - neighborhood of two RNA secondary structures

- similar to approach of finding optimal k -neighbors (Freyhult et al. 2007)
a.k.a. RNAbor
- optimal k -neighbor is optimal secondary structure s_k with $d_{BP}(s_k, s_{\text{target}}) = k$
- extension to second target structure s_{target_2}
- optimal k, l -neighbor is optimal secondary structure s_{kl} with $d_{BP}(s_{kl}, s_{\text{target}}) = k$ and $d_{BP}(s_{kl}, s_{\text{target}_2}) = l$

k, l - neighborhood of two RNA secondary structures



Finding all optimal k, l neighbors

Definitions

$$d_{BP}(s_1, s_2) = |s_1 \cup s_2| - |s_1 \cap s_2|$$

$$s[i, j] = \{(p, q) \in s : i \leq p < q \leq j\}$$

$$\delta(x, y) = 1, \text{ iff } x = y \text{ (Kronecker function)}$$

$$\lambda_1(i, j, s) = d_{BP}(s[i, j], s[i, j - 1])$$

$$\lambda_2(i, j, s) = d_{BP}(s[i, j], s[i, u] \cup s[u + 1, j])$$

$$\lambda_3(i, j, p, q, s) = d_{BP}(s[i, j], \{(i, j)\} \cup s[p, q])$$

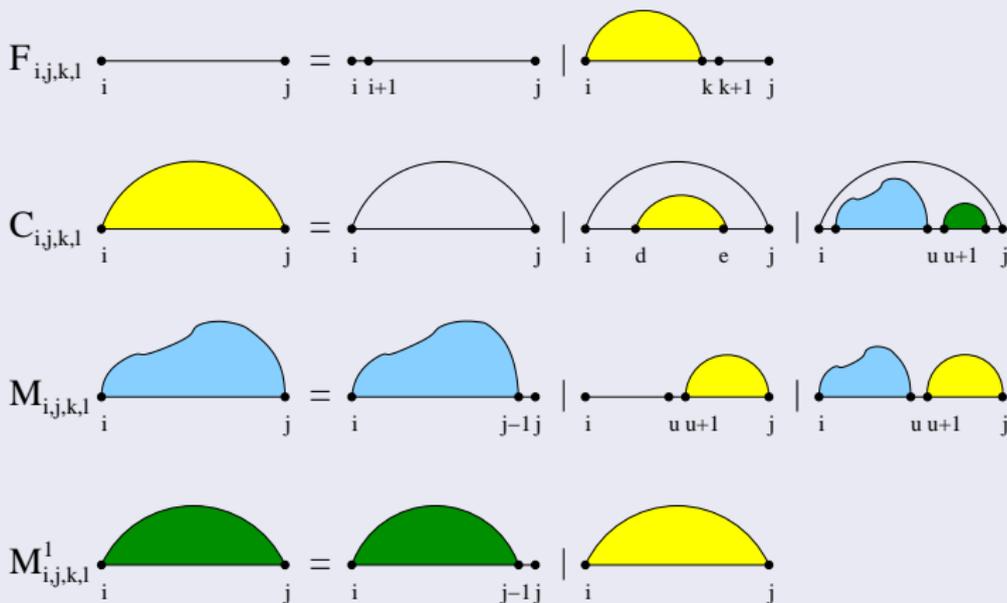
$$\lambda_4(i, j, u, s) = d_{BP}(s[i, j], \{(i, j)\} \cup s[i + 1, u] \cup s[u + 1, j - 1])$$

$$\lambda_5(i, j, u, s) = d_{BP}(s[i, j], s[u, j])$$

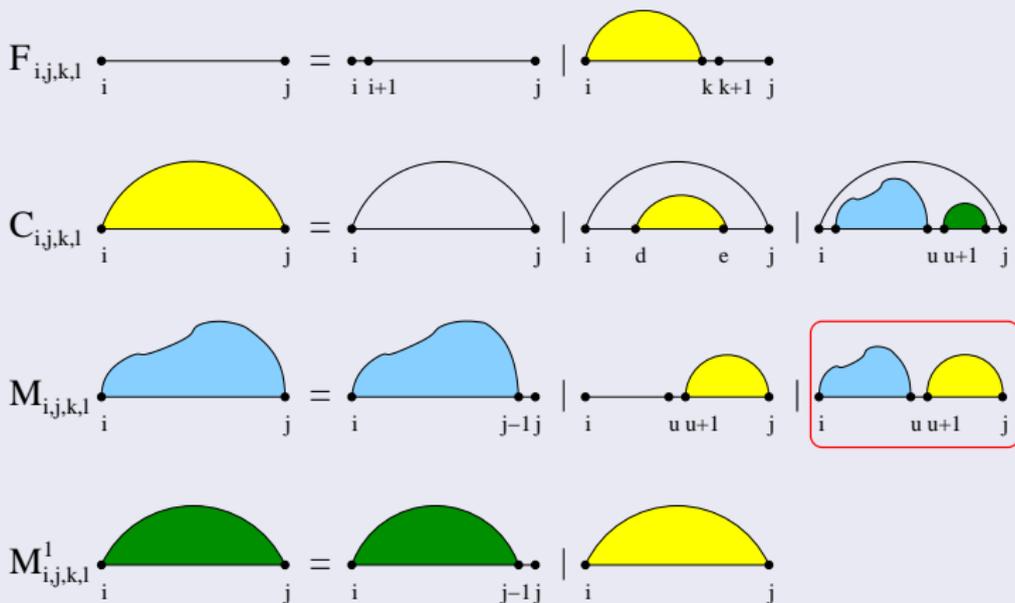
MFE Algorithm for finding all optimal k, l -neighbors

$$\begin{aligned}
 F_{i,j,k,l} &= \min \left\{ \begin{array}{l} F_{i+1,j,k-\lambda_1(i,j,s_1),l-\lambda_1(i,j,s_2)}, \\ \min_{i < u \leq j} \omega_1 + \hat{\omega}_1 = k - \lambda_2(i,j,u,s_1) \min_{\omega_2 + \hat{\omega}_2 = l - \lambda_2(i,j,u,s_2)} C_{i,u,\omega_1,\omega_2} + F_{u+1,j,\hat{\omega}_1,\hat{\omega}_2} \end{array} \right\} \\
 C_{i,j,k,l} &= \min \left\{ \begin{array}{l} \delta(d_{\text{BP}}(s_1[i,j], \{(i,j)\})) \cdot \delta(d_{\text{BP}}(s_2[i,j], \{(i,j)\})) \cdot \mathcal{H}(i,j), \\ \min_{i < p < q < j} \left\{ C_{p,q,k-\lambda_3(i,j,p,q,s_1),l-\lambda_3(i,j,p,q,s_2)} + \mathcal{I}(i,j,p,q) \right\}, \\ \min_{i < u < j} \omega_1 + \hat{\omega}_1 = k - \lambda_4(i,j,u,s_1) \min_{\omega_2 + \hat{\omega}_2 = l - \lambda_4(i,j,u,s_2)} \left\{ M_{i+1,u,\omega_1,\omega_2} + M_{u+1,j-1,\hat{\omega}_1,\hat{\omega}_2}^1 + a \right\} \end{array} \right\} \\
 M_{i,j,k,l} &= \min \left\{ \begin{array}{l} \min_{i \leq u < j} \left\{ (u-i) \cdot c + C_{u,j,k-\lambda_5(i,j,u,s_1),l-\lambda_5(i,j,u,s_2)} + b \right\}, \\ \min_{i < u < j} \omega_1 + \hat{\omega}_1 = k - \lambda_2(i,j,u,s_1) \min_{\omega_2 + \hat{\omega}_2 = l - \lambda_2(i,j,u,s_2)} \left\{ M_{i,u,\omega_1,\omega_2} + C_{u+1,j,\hat{\omega}_1,\hat{\omega}_2} + b \right\}, \\ M_{i,j-1,k-\lambda_1(i,j,s_1),l-\lambda_1(i,j,s_2)} + c \end{array} \right\} \\
 M_{i,j,k,l}^1 &= \min \left\{ \begin{array}{l} M_{i,j-1,k-\lambda_1(i,j,s_1),l-\lambda_1(i,j,s_2)}^1 + c, \\ C_{i,j,k,l} + b \end{array} \right\}
 \end{aligned}$$

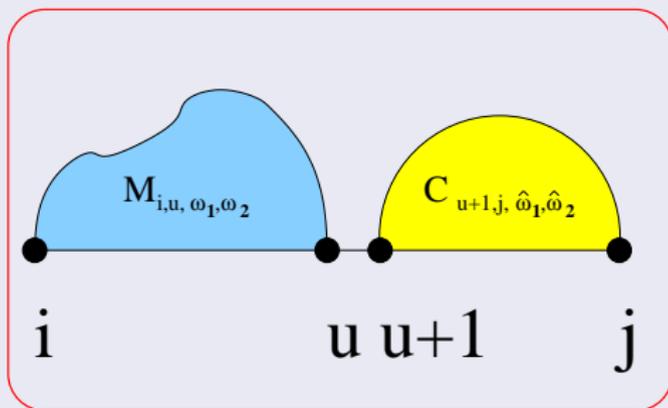
Recursion scheme for finding k, l -neighbors



Recursion scheme for finding k, l -neighbors



Recursion scheme for finding k, l -neighbors



$$\omega_1 + \hat{\omega}_1 = k - \lambda_2(i, j, u, s_1) = k - d_{\text{BP}}(s_1[i, j], s_1[i, u] \cup s_1[u + 1, j])$$

$$\omega_2 + \hat{\omega}_2 = l - \lambda_2(i, j, u, s_2) = l - d_{\text{BP}}(s_2[i, j], s_2[i, u] \cup s_2[u + 1, j])$$

Improving runtime properties of the algorithm

Time complexity

- finding all optimal k, l -neighbors with $k \leq d_1$ and $l \leq d_2$:
 $\mathcal{O}(d_1 \cdot d_2 \cdot n^3)$
- $d_x \leq MM(seq) + |s_x|$
- $MM(seq) < n/2 - 1$ and $|s_x| < n/2 - 1$
- $\mathcal{O}(n^5)$

Memory complexity

- finding all optimal k, l -neighbors with $k \leq d_1$ and $l \leq d_2$:
 $\mathcal{O}(m \cdot n^2)$, $m = \frac{1}{2} \cdot d_1 \cdot d_2$
- checkerboard like memory fill pattern
- $\mathcal{O}(n^4)$

Improving runtime properties of the algorithm

Observations

straight implementation results:

- more than 3.5h runtime for sequence of length $100nt$
- more than 30GB RAM for sequences with length $> 200nt$

Improving runtime properties of the algorithm

Observations

straight implementation results:

- more than 3.5h runtime for sequence of length $100nt$
- more than 30GB RAM for sequences with length $> 200nt$

BUT: Energy arrays are filled by less than 2%

Improving runtime properties of the algorithm

Observations

straight implementation results:

- more than 3.5h runtime for sequence of length $100nt$
- more than 30GB RAM for sequences with length $> 200nt$

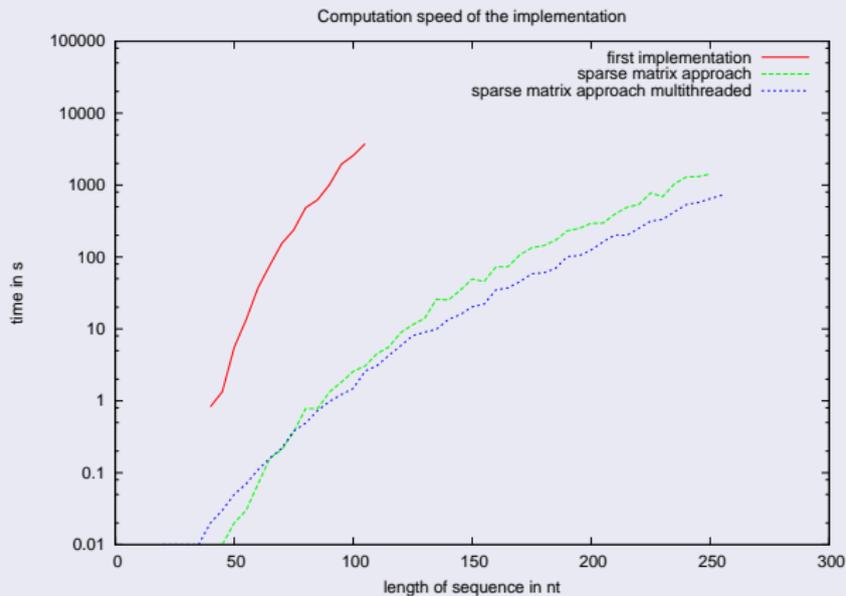
BUT: Energy arrays are filled by less than 2%

Improvements

- operating with sparse matrices to lower runtime and memory requirements
- parallelization of recursive algorithm (parallel computation of elements in diagonals of energy matrices)

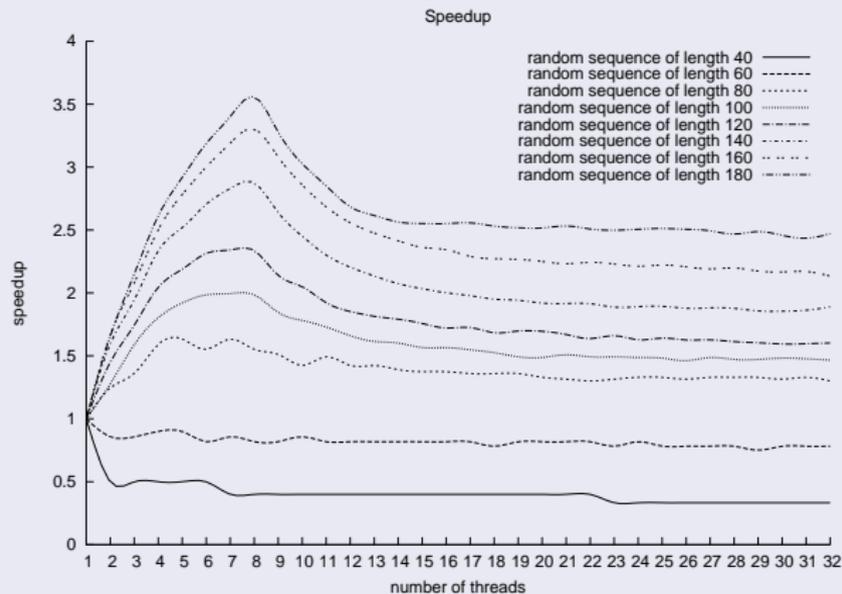
Improving runtime properties of the algorithm

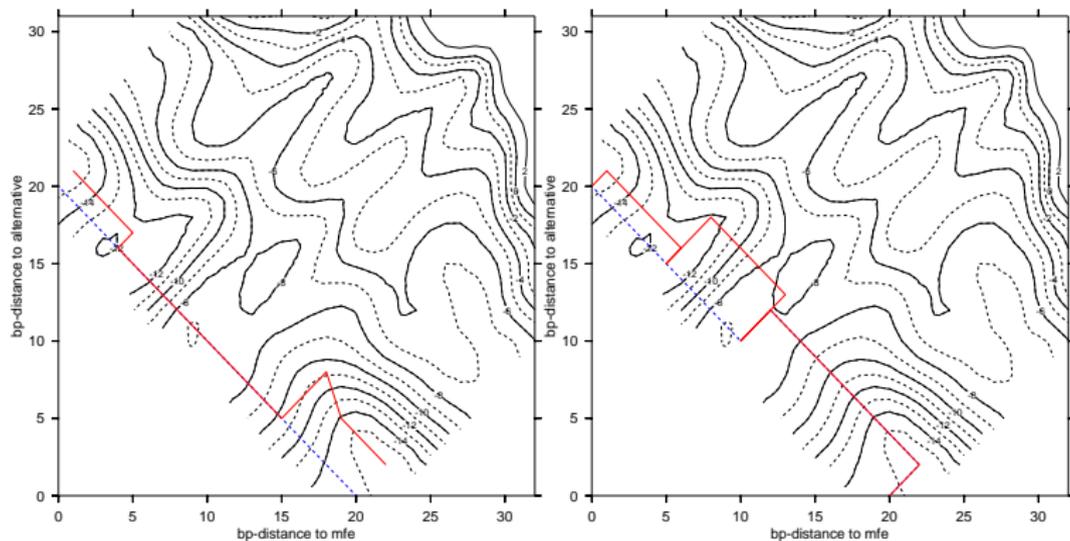
Stats



Improving runtime properties of the algorithm

Stats



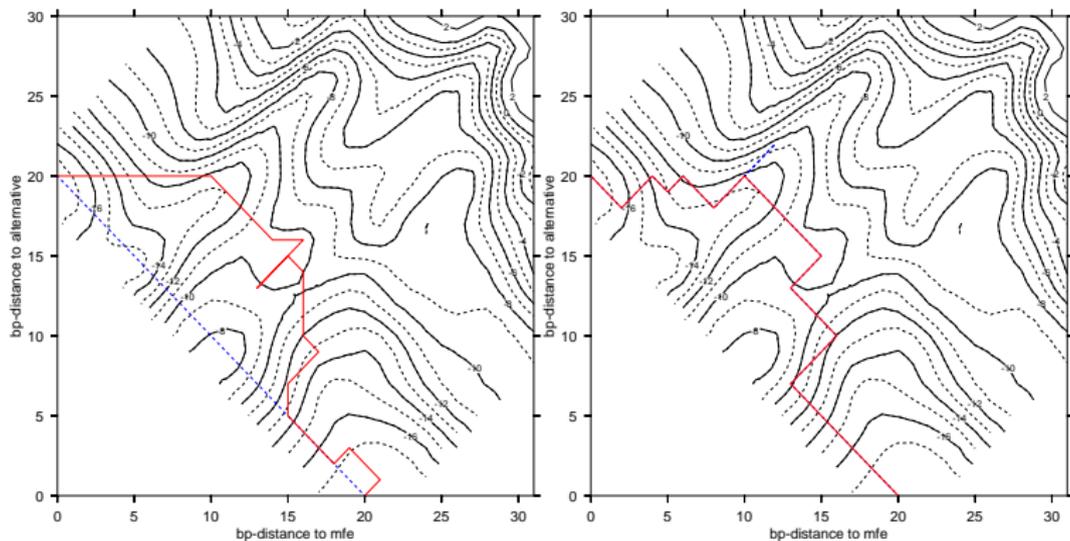


- GGCUGUAUGGCAGCUGCCUCGUUAAGAGGUGAUUACUAUGUAGUC
 ((((((.....))))))(((.....))))(((.....)))) -15.40 (*MFE structure*)
 ((((((.....))))))(((.....)))).....)))))) -14.70 (*alternative structure*)

- Height field contours: Gibbs free energy of the corresponding k , l -neighborhood

- Left side: Trajectories of
 (**barrier tree generated path**) 9.6 kcal/mol (−5.8 kcal/mol)
 (**direct folding path**) 10.5 kcal/mol (−4.9 kcal/mol)

- Right side: Trajectories of
 (**Monte carlo method (SA)**) 9.8 kcal/mol (−5.8 kcal/mol)
 (**Deterministic method**) 9.6 kcal/mol (−5.8 kcal/mol)



- GGGCGGGUUCGCCUCCGCUAAAUGCGGAAGAUAUUUGUGUCU
 ((((((.....))))))(((.....)))(((.....))) -18.20 (*MFE structure*)
 ((((((.....)))))).....))))) -17.70 (*alternative structure*)

- Height field contours: Gibbs free energy of the corresponding k , l -neighborhood

- Left side: Trajectories of
 (**barrier tree generated path**) 10.7 kcal/mol (−7.5 kcal/mol)
 (**direct folding path**) 13.33 kcal/mol (−4.87 kcal/mol)

- Right side: Trajectories of
 (**Monte Carlo method (SA)**) 10.7 kcal/mol (−7.5 kcal/mol)
 (**Deterministic method**) 10.7 kcal/mol (−7.5 kcal/mol)

Thanks to:

Stephan Bernhart

Ilenia Boria

Christoph Flamm

Andreas Gruber

Christian Höner

Ivo Hofacker

Andrea Tanzer

Caroline Thurner

Alexander Ullrich

The audience for listening