

# Cycle decomposition of small RNA configuration space

JING QIN<sup>1,2</sup> AND PETER F. STADLER<sup>1-5</sup>

<sup>1</sup>Max Planck Institute for Mathematics in the Sciences, Inselstraße 22, D-04103 Leipzig, Germany

<sup>2</sup>Bioinformatics Group, Department of Computer Science, and Interdisciplinary Center for Bioinformatics, University of Leipzig, Härtelstrasse 16-18, D-04107 Leipzig, Germany.

<sup>3</sup>RNomics Group, Fraunhofer Institut for Cell Therapy and Immunology, Perlickstraße 1, D-04103 Leipzig, Germany

<sup>4</sup>Inst. f. Theoretical Chemistry, University of Vienna, Währingerstrasse 17, A-1090 Vienna, Austria

<sup>5</sup>The Santa Fe Institute, 1399 Hyde Park Rd., Santa Fe, New Mexico

**Abstract** It is a well known fact that the behavior of simulated annealing algorithms is tightly related to the hierarchical decomposition of their configuration spaces in cycles. We here apply the iterative routine invented by Wentzell and Freidlin [1] to construct the cycle decomposition of small RNA configuration spaces, for instance, hairpins. We furthermore explore the relationships of cycles and the barrier tree of the energy landscape.

## RNA Configuration Space

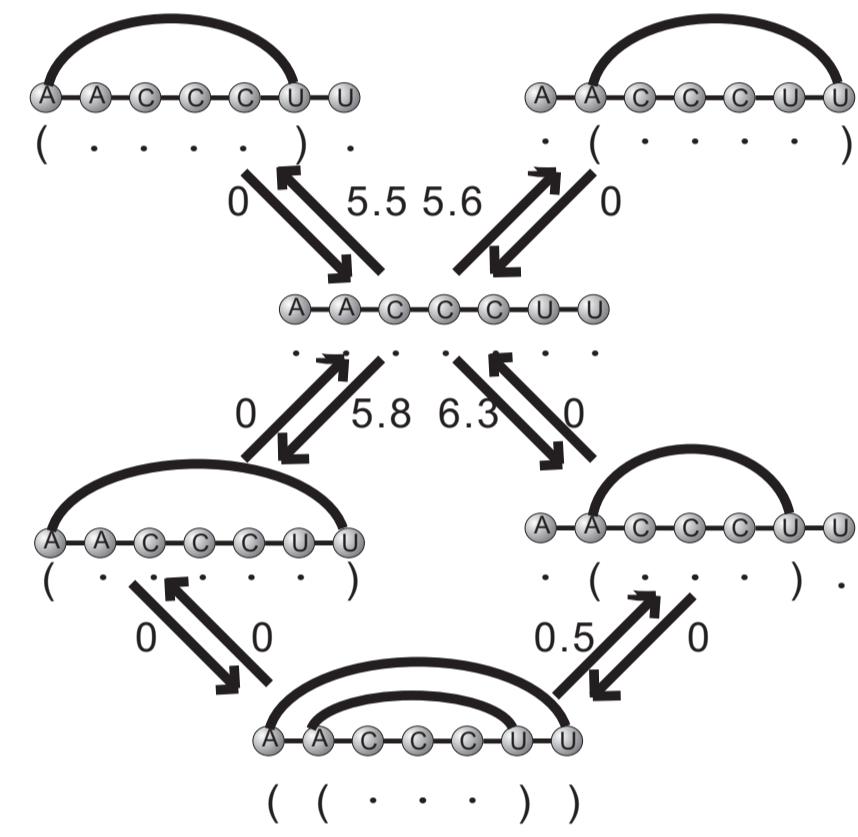
Given an RNA sequence  $R = R_1R_2\dots R_n$ , where  $R_i \in \{A, U, G, C\}$ , a **secondary structure** is a graph on  $n$  vertices with three properties:

- $R_i$  and  $R_{i+1}$  are connected for any  $1 \leq i \leq n-1$ ;
- each vertex  $R_i$  can be paired to at most one other vertex  $R_j$  (exclude the arcs exist in primary structure) if  $R_iR_j \in A = \{AU, UA, GC, CG, GU, UG\}$ ;
- if both  $R_i < R_j$  and  $R_h < R_l$  are paired, then  $i < h < j$  implies that  $i < l < j$ .

Given an RNA sequence, the **configuration space (Landscape)** can be viewed as a directed graph  $G(V, E)$ . In which, the vertex set  $V$  is formed by all the possible secondary structures with respect to the RNA sequence and we say  $i \rightarrow j \in E$  if  $j$  can be obtained from  $i$  by either adding or removing an arc in  $i$ . Let  $F(i)$  denote the free energy of the secondary structure  $i$ . The weight of a directed edge  $i \rightarrow j$ , denoted by  $w(i \rightarrow j)$ , is defined by

$$w(i \rightarrow j) = \begin{cases} 0 & \text{if } F(i) > F(j); \\ F(j) - F(i) & \text{otherwise.} \end{cases} \quad (1)$$

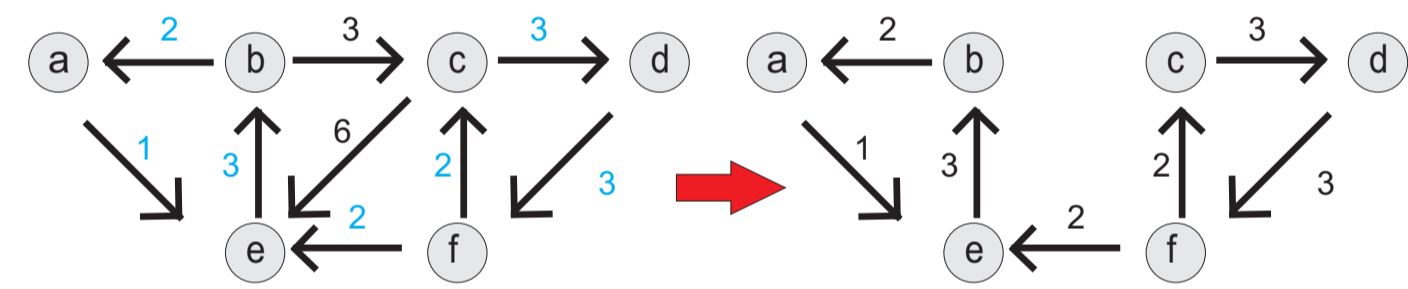
**Example: A toy configuration space for AACCCUU**



## Exit Graph and Strongly Connected Component

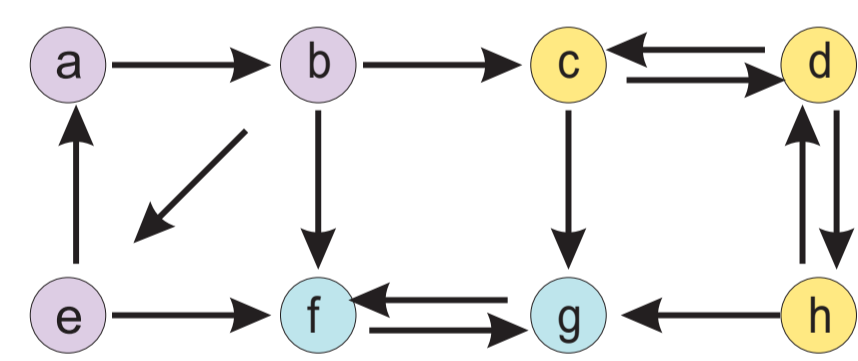
Given a directed graph  $G = (V, E)$ , the **exit graph** of  $G$ , denoted by  $H = (V_H, E_H)$ , is a subgraph of  $G$ , such that  $V_H = V$  and  $E_H = \{i \rightarrow j | C(i, j) = \inf_{i \rightarrow k \in E} C(i, k)\}$ .

**Example: Construction of exit graph**



A directed graph is called **strongly connected** if there is a path from each vertex in the graph to every other vertex. The **strongly connected components** of a directed graph  $G$  are its maximal strongly connected subgraphs with respect to set inclusion.

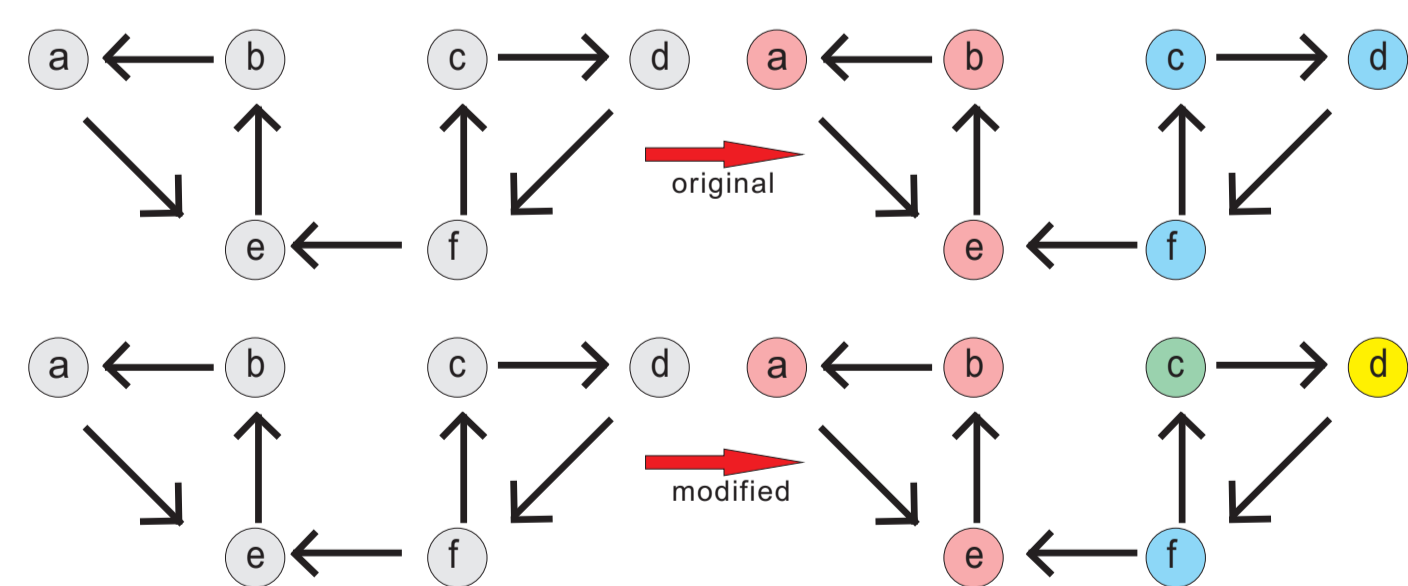
**Example: Graph with strongly connected components colored in different colors**



## Tarjan-V: slight modification of Tarjan's Algorithm

After the exit graph  $H = (V_H, E_H)$  of  $G$  is constructed, Tarjan's algorithm can be used to search for the strongly connected components in  $H$  with complexity  $O(|V_H| + |E_H|)$ . But here, a slight modification is we have to distinguish the "leaf" strongly connected component from the others.

**Example: Difference between original and modified Tarjan's algorithm**



## Cycle decomposition Procedure

Given the RNA sequence, the cycle decomposition in cycles is realized in an iterative way:

- Initialize  $C^0(V^0, E^0)$  which is identified with the original RNA configuration space;
- Assume all the information of  $C^k(V^k, E^k)$  has been derived. We construct  $C^{k+1}(V^{k+1}, E^{k+1})$  as follows. Firstly, we build the exit graph of  $C^k$  and derive the components according to the Tarjan-V procedure. By contracting the vertices of  $C^k$  within the same component, we obtain the vertex-set of  $C^{k+1}$ , i.e.  $V^{k+1}$ . We say the  $C^k$ -vertices  $\{v_1^k, v_2^k, \dots\}$  contracted to a single  $V^{k+1}$ -vertex  $m$  are the son-structures of  $m$ , denoted by  $m \downarrow v_i^k$ . To define  $E^{k+1}$ , given two  $V^{k+1}$ -vertices  $p$  and  $q$ ,  $p$  are directed to  $q$  if and only if there at least exist one pair of  $V^k$ -vertices  $v_1$  and  $v_2$  such that  $p \downarrow v_1, q \downarrow v_2$  and  $v_1 \rightarrow v_2 \in E^k$ . The weight of the directed edge  $p \rightarrow q$  is assigned by Eqn. (2).

$$w^{k+1}(p \rightarrow q) = \inf\{H_m^{k+1}(p) + w^k(v_1, v_2) - H_c^k(v_2)\}. \quad (2)$$

In which,  $H_c^k(v)$  and  $H_m^k(v)$  denote the so-called **escape energy** and **mixing energy** respectively given as follows:

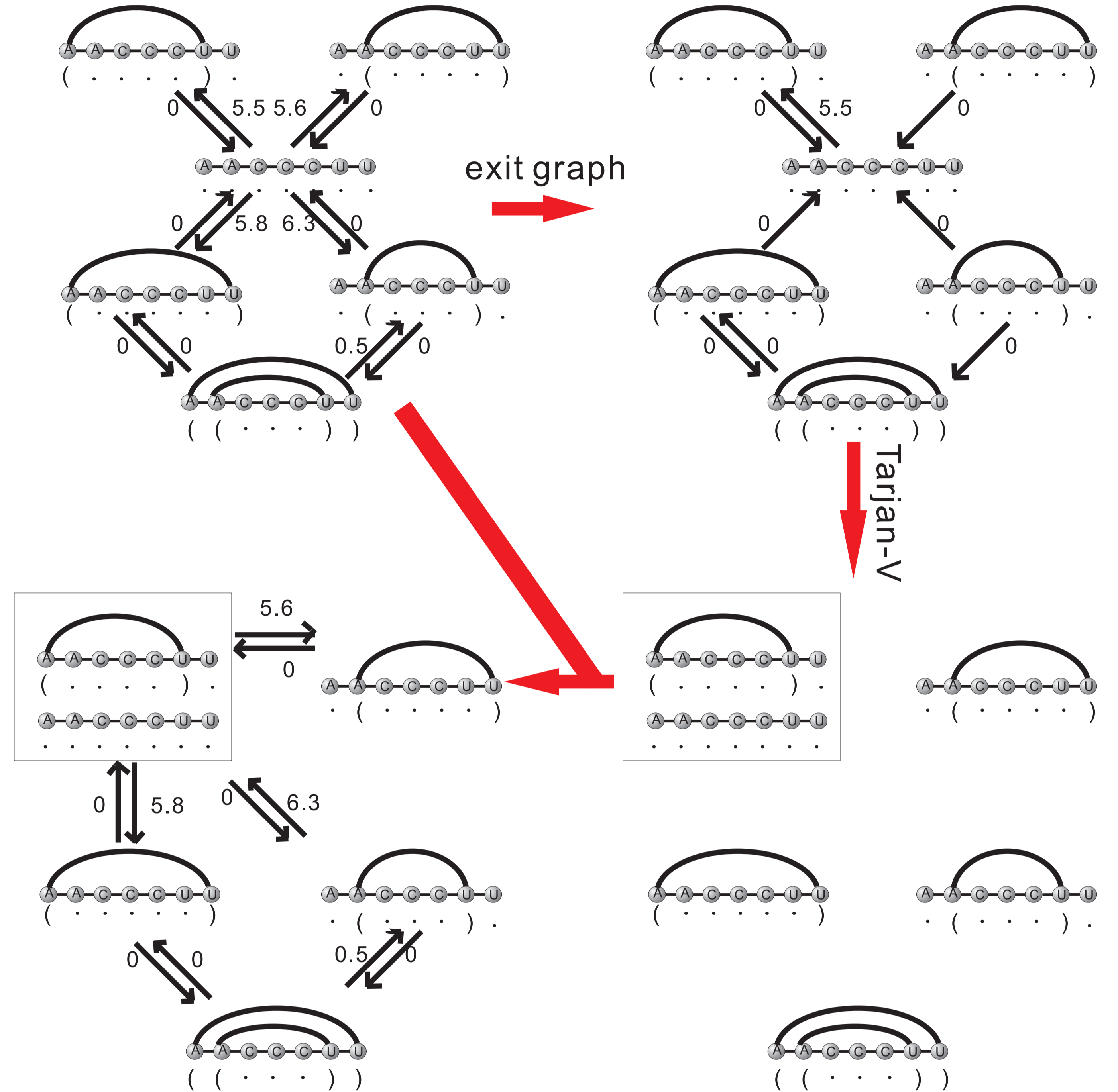
$$H_c^k(v) = \inf\{w^k(v \rightarrow w) | v \rightarrow w \in E^k, v, w \in V^k\} \quad (3)$$

$$H_m^k(v) = \sup\{H_c^{k-1}(w) | v \downarrow w, w \in V^{k-1}\}. \quad (4)$$

- The whole procedure terminates in case of  $|V^{k+1}| = 1$ , iterates otherwise.

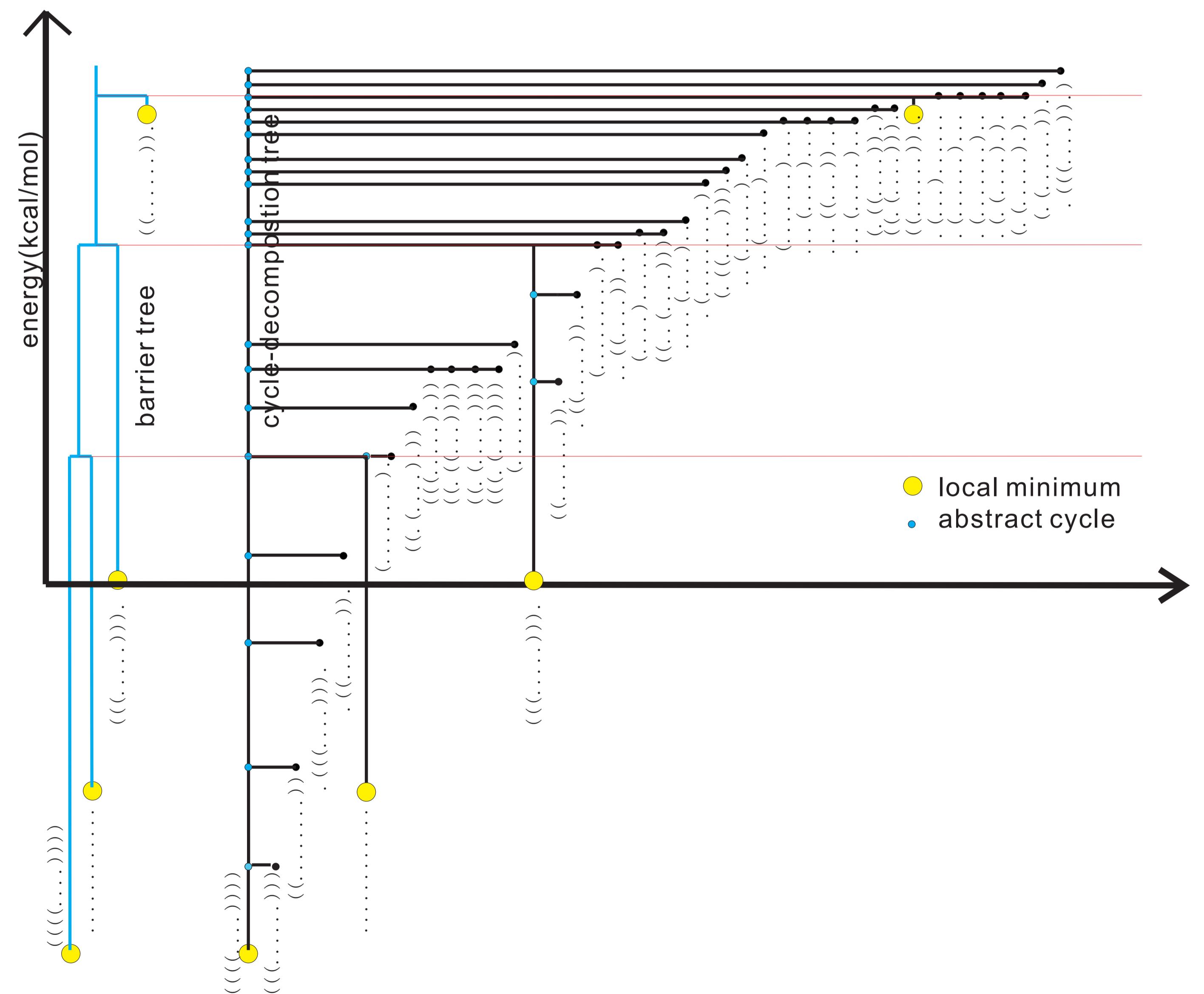
## Cycle decomposition Procedure

**Example: Towards Level-1 from Level-0 for AACCCUU**



## Computational Result

**Example: Barrier tree and cycle decomposition tree for sequence GGAAUAAUCC**



**Remarks:**

- In the worst case, the complexity of the cycle decomposition algorithm is  $O(|V|(|V| + |E|))$ . In which,  $O(|V| + |E|)$  is the complexity of the (modified) Tarjan's algorithm and the other factor,  $|V|$ , caused by (in the worst case) the number of iterations, denoted by  $\sigma$ , need to run before the procedure is terminated. We remark here, the integer  $\sigma$  is a valuable parameter that also reflects the topology of the landscape itself. For RNA configuration space, interestingly, we observe that the bigger the size of the space, the smaller the ratio  $\frac{\sigma}{|V|}$  is.
- By construction, the mixing energy of the cycle (firstly) merging two local minima identified with their saddle height obtained by the RNAbARRIER[2] included in the Vienna package. I.e. the barrier tree can be viewed as a subtree of the cycle decomposition tree. The RNAbARRIER is based on the flooding-algorithm which need more effort in case of the degenerate RNA landscape. Comparing with RNAbARRIER, cycle decomposition by construction avoid the problem may caused by multiple saddle points. Also, we note here the mixing energy is only one of the valuable parameters that can be read from cycle decomposition procedure.

## References

- [1] M.I. Freidlin and A.D. Wentzell, Perturbations of Stochastic Dynamic Systems, Springer-Verlag, 1984.
- [2] C. Flamm, I.L. Hofacker, P.F. Stadler and M.T. Wolfinger, Barrier Trees of Degenerate Landscapes, Z. Phys. Chem., 2002 (216):155-173.