

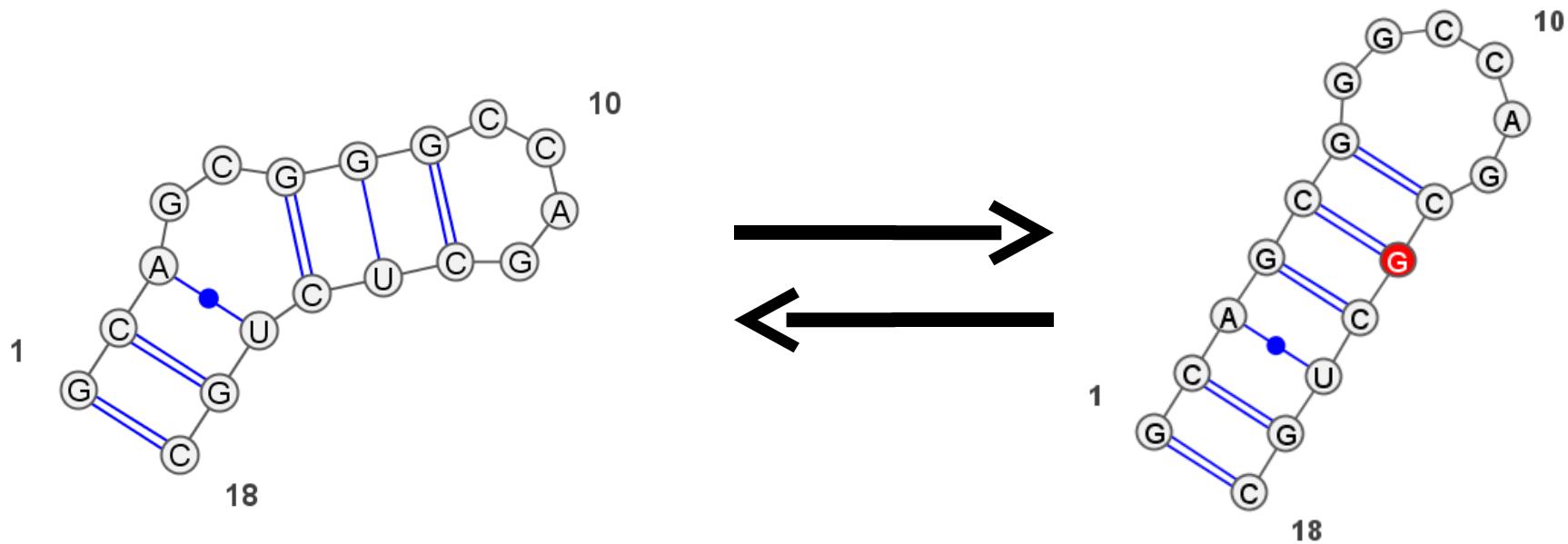
Journées MAGNUM – Partie 1

**An unbiased adaptive sampling
algorithm for the exploration of
RNA mutational landscapes under
evolutionary pressure**

Jérôme Waldspühl, PhD
School of Computer Science,
McGill Centre for Bioinformatics,
McGill University, Canada

Yann Ponty, PhD
Laboratoire d'informatique (LIX),
École Polytechnique, France

RNAmutants: Algorithms to explore the RNA mutational landscape



Understanding how mutations influence RNA secondary structures **AND** how structures influence mutations (Waldripühl et al., PLoS Comp Bio, 2008).

Sampling k-mutants

Seed
↓

CAGUGAUUGCAGUGCGAUGC
..((.((((((...)))))))

(-1.20)

CAGUGAUUGCAGUGCGAU**C**
..(.((((((...)))))))

(-3.40)

CAGUGAUUGCAGUGCG**G**UGC
((.((....)).)).....

(-0.30)

CAGUGAU**C**GCAGUGCGAUGC
.....(((.((...))))..

(-3.10)

uAGc**GccgGgAGacCGgcGC**
..((((((.((...)))))))

(-18.00)

CccUGgccGCA**agGCCAgGg**
((((((.((...))))))))

(-20.40)

CcGUGgccGC**gagGCCAcGg**
((((((.((...))))))))

(-19.10)

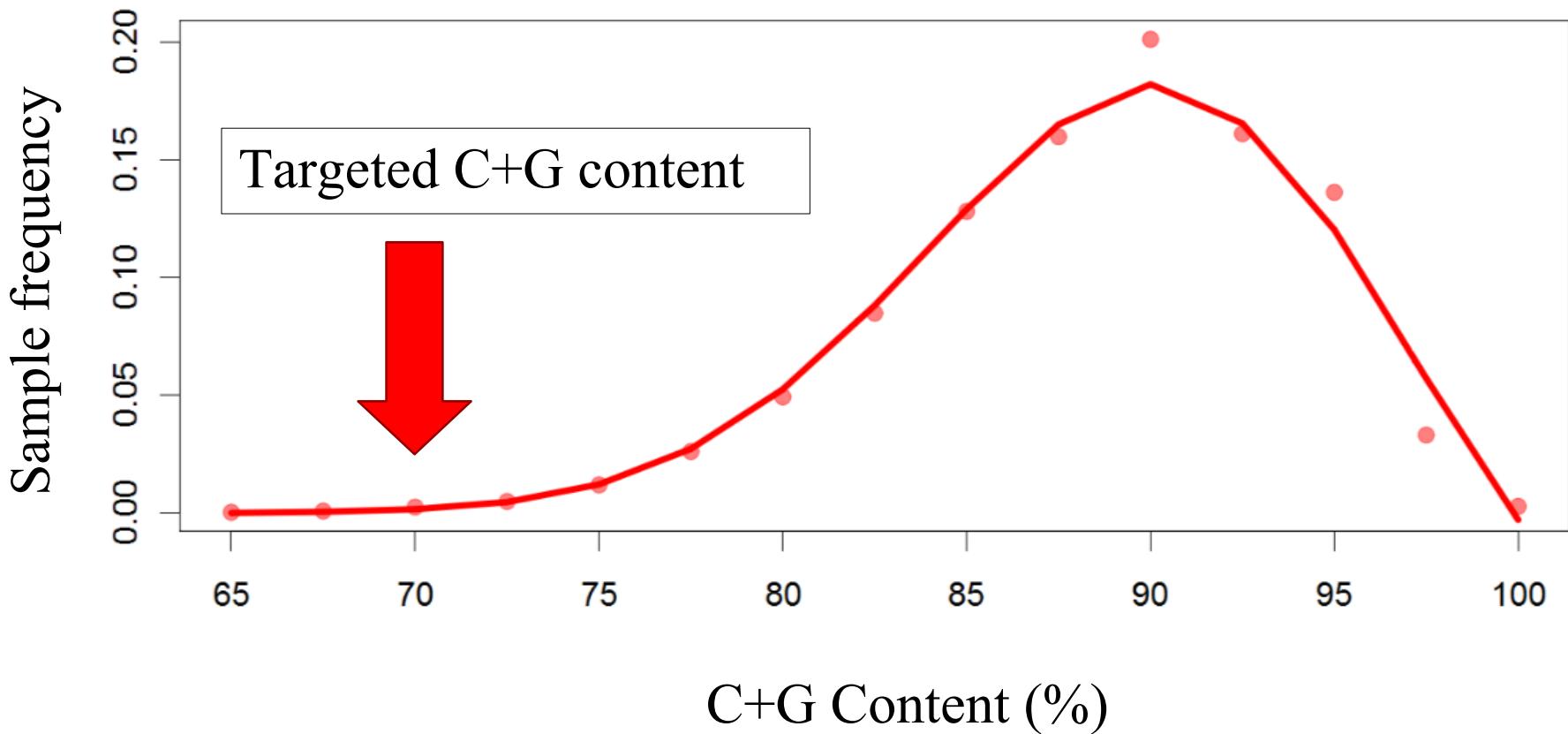
Classic: 0 mutation

RNAmutants: 1 mutation

RNAmutants: 10 mutations

Sample k mutations increasing the folding energy
Consequence: Increase in C+G content

Objectives



How to efficiently sample sequences at arbitrary C+G contents ... without introducing a bias!

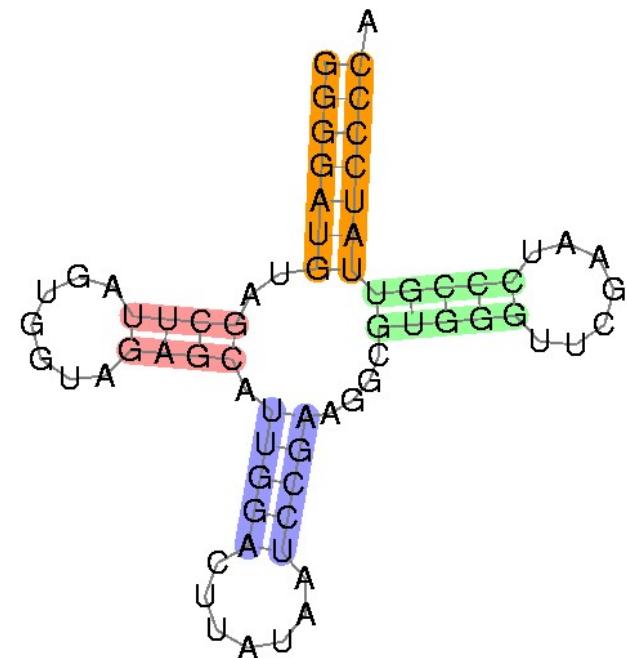
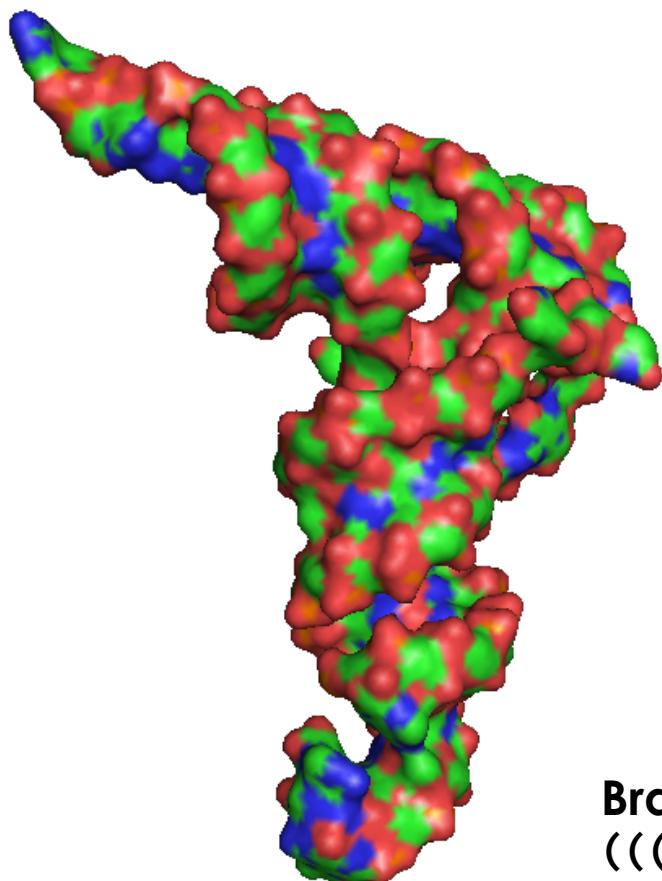
- **Background: RNAmutants in a nutshell**
Algorithms to sample RNA secondary structures and mutations.
- **Our approach: Adaptive sampling**
Uniformly shifting the distribution of samples.
- **Results: Evolutionary studies**
Insights on the evolutionary pressure stemming from an optimization of the thermodynamical stability.

- **Background: RNAmutants in a nutshell**
Algorithms to sample RNA secondary structures and mutations.
- Our approach: Adaptive sampling
Uniformly shifting the distribution of samples.
- Results: Evolutionary studies
Insights on the evolutionary pressure stemming from an optimization of the thermodynamical stability.

RNA secondary structure



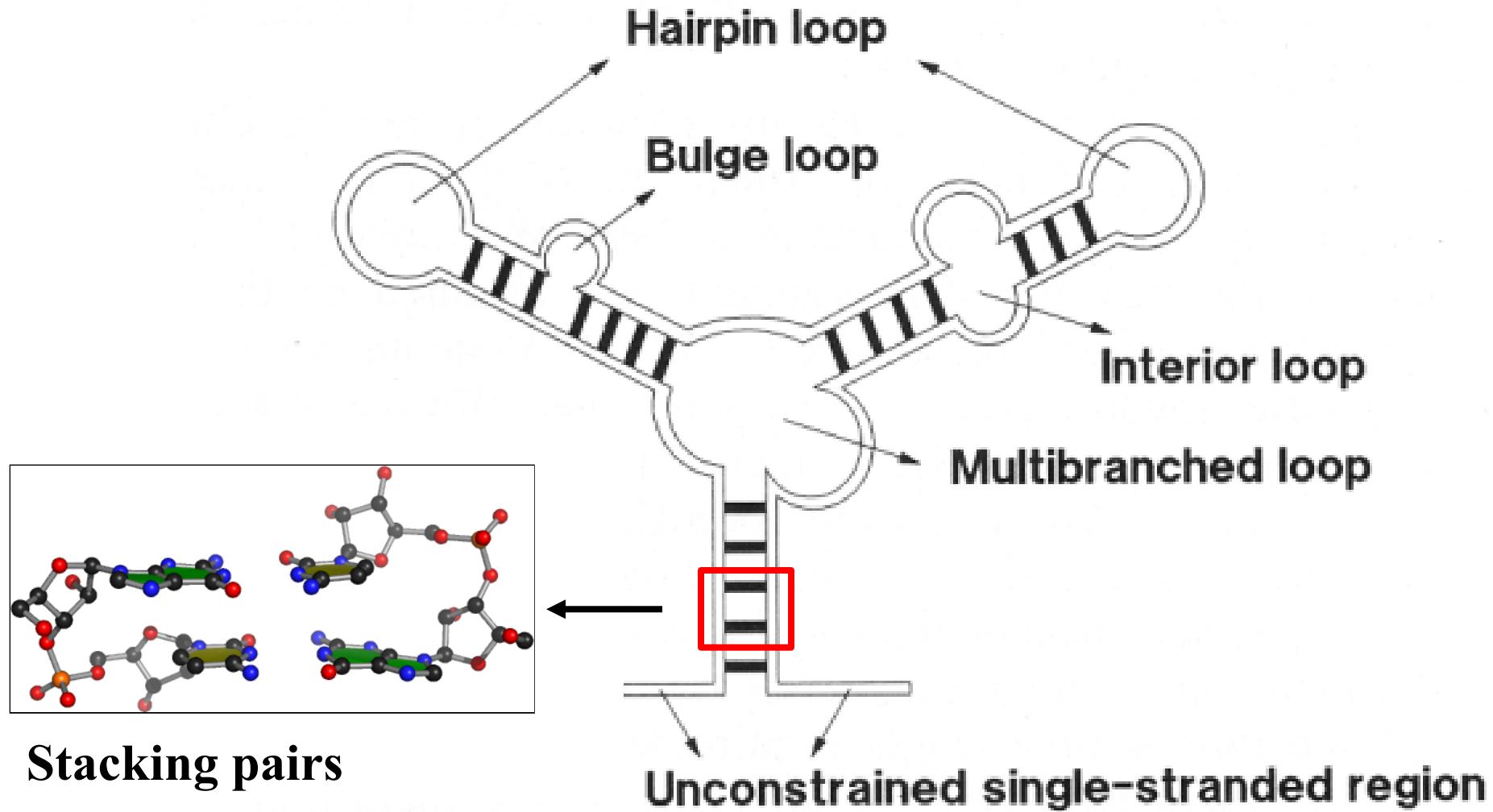
The **secondary structure** is the ensemble of base-pairs in the structure.



Bracket notation:

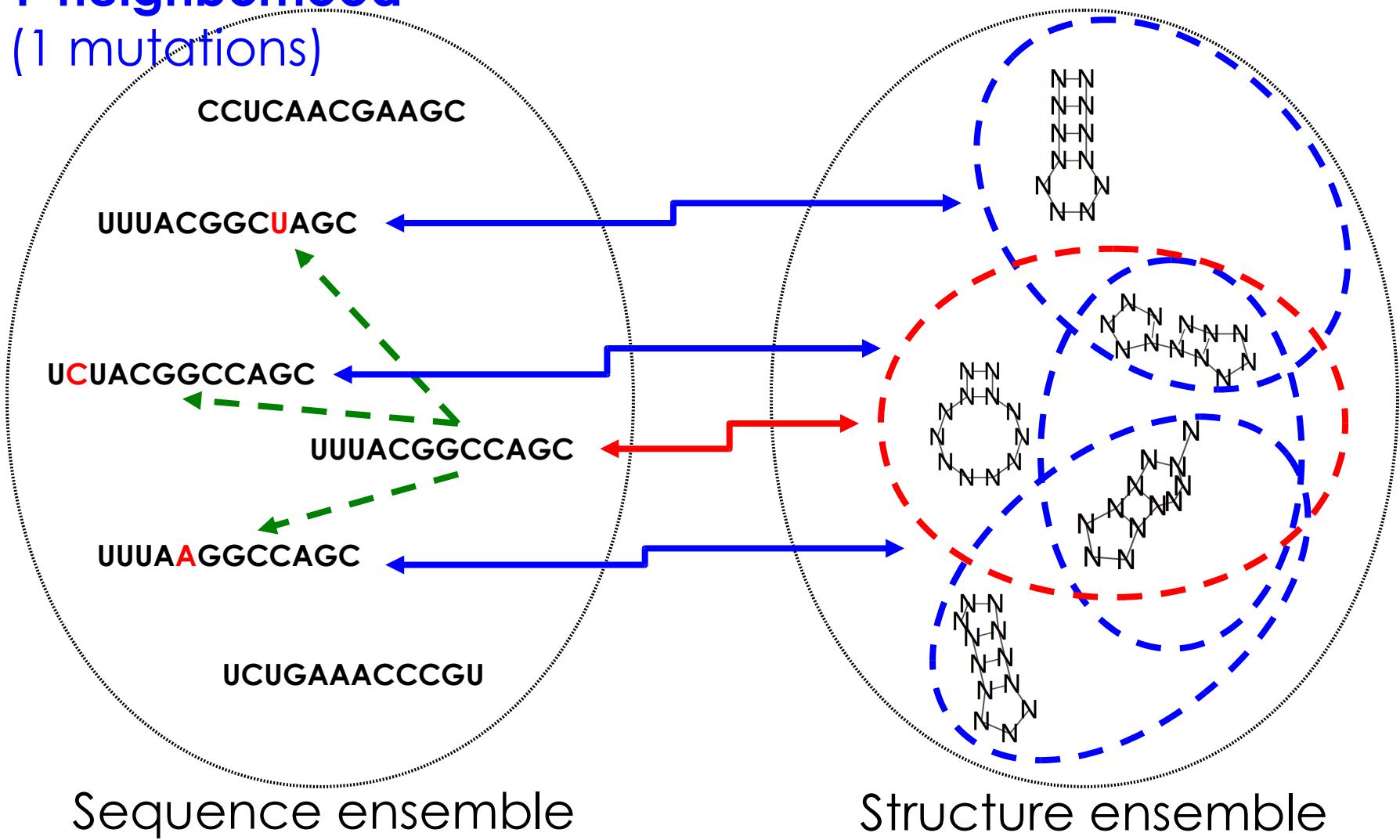
$((((((((....))))..(((....))))..(((....))))..))))))$

Loop decomposition



Parameterization of the mutational landscape

1-neighborhood
(1 mutations)



Classical Recursions (Zuker & Stiegler, McCaskill)

$$\bullet \overbrace{z}^{z} \bullet = \bullet \overbrace{z}^{z} \bullet \overbrace{\quad}^{r-1} \bullet \overbrace{z^B}^{r} \bullet \overbrace{\quad}^{j} + \bullet \overbrace{z}^{z} \bullet \overbrace{\quad}^{j-1} \bullet \overbrace{z}^{j}$$

$$\bullet \overbrace{z^B}^{z^B} \bullet = \bullet \overbrace{\quad}^{i} \bullet \overbrace{i+1}^{i+1} \bullet \overbrace{\quad}^{j-1} \bullet \overbrace{j}^{j} + \bullet \overbrace{\quad}^{i} \bullet \overbrace{r}^{r} \bullet \overbrace{\quad}^{s} \bullet \overbrace{j}^{j} + \bullet \overbrace{\quad}^{i+1} \bullet \overbrace{i+1}^{r-1} \bullet \overbrace{r}^{r} \bullet \overbrace{\quad}^{j-1} \bullet \overbrace{j}^{j}$$

$$\bullet \overbrace{z^{M1}}^{z^{M1}} \bullet = \bullet \overbrace{z^B}^{z^B} \bullet \overbrace{\quad}^{r} \bullet \overbrace{\quad}^{j}$$

$$\bullet \overbrace{z^M}^{z^M} \bullet = \bullet \overbrace{\quad}^{i} \bullet \overbrace{r}^{r} \bullet \overbrace{\quad}^{j} \bullet \overbrace{z^{M1}}^{z^{M1}} + \bullet \overbrace{z^M}^{z^M} \bullet \overbrace{\quad}^{r-1} \bullet \overbrace{r}^{r} \bullet \overbrace{\quad}^{j}$$

Enumerate all secondary structures

RNA mutants Generalize Classical Algorithms

$$\begin{aligned}
 x \xrightarrow{\mathcal{Z}} y &= \text{Diagram } 1 + \text{Diagram } 2 \\
 &\quad \text{Diagram } 1: \text{Diagram } 3 + \text{Diagram } 4 \\
 &\quad \text{Diagram } 2: \text{Diagram } 5 + \text{Diagram } 6 + \text{Diagram } 7 \\
 x \xrightarrow{\mathcal{Z}^{M1}} y &= \text{Diagram } 8 \\
 x \xrightarrow{\mathcal{Z}^M} y &= \text{Diagram } 9 + \text{Diagram } 10
 \end{aligned}$$

Diagrams illustrating RNA secondary structure components and their relationships:

- Diagram 1:** A horizontal bar from x to y labeled \mathcal{Z} .
- Diagram 2:** A horizontal bar from x to y labeled \mathcal{Z}^B , with a red dashed arc above it connecting positions i to j .
- Diagram 3:** A horizontal bar from x to y labeled \mathcal{Z} , with a red dashed arc below it connecting positions $r-1$ to r .
- Diagram 4:** A horizontal bar from x to y labeled \mathcal{Z} , with a red dashed arc below it connecting positions i to $j-1$.
- Diagram 5:** A horizontal bar from x to y labeled \mathcal{Z}^B , with a red dashed arc below it connecting positions i to j . It also features wavy lines between u and v , and between v and y .
- Diagram 6:** A horizontal bar from x to y labeled \mathcal{Z}^B , with a red dashed arc below it connecting positions i to s . It also features wavy lines between u and v , and between v and y .
- Diagram 7:** A horizontal bar from x to y labeled \mathcal{Z}^M , with a red dashed arc below it connecting positions i to j . It also features wavy lines between u and v , and between v and w , and between w and z .
- Diagram 8:** A horizontal bar from x to y labeled \mathcal{Z}^{M1} , with a red dashed arc below it connecting positions i to r .
- Diagram 9:** A horizontal bar from x to y labeled \mathcal{Z}^M , with a red dashed arc below it connecting positions i to r .
- Diagram 10:** A horizontal bar from x to y labeled \mathcal{Z}^{M1} , with a red dashed arc below it connecting positions i to r .

Enumerate all secondary structures over all mutants

Our approach



RNAmutants

- Explore the complete mutation landscape.
- Polynomial time and space algorithm.
- Compute the partition function for all sequences:

RNAmutants:

$$Z = \sum_s \sum_S \exp\left(-\frac{E(s, S)}{RT}\right)$$

Single sequence:

$$Z(s) = \sum_S \exp(b \times E(s, S))$$

- Sample by backtracking in dynamic prog. Tables.
=> Weighted random generation (recursive method)

(Waldspuhl et al., PLoS Comp Bio, 2008)

Sampling k-mutants

Seed
↓

CAGUGAUUGCAGUGCGAUGC
..((.((((((...)))))))

(-1.20)

CAGUGAUUGCAGUGCGAU**C**
..(.((((((...)))))))

(-3.40)

CAGUGAUUGCAGUGCG**g**UGC
((.((....)).)).....

(-0.30)

CAGUGAU**C**GCAGUGCGAUGC
.....(((.((...))))..

(-3.10)

uAG**c**GccgGgAGac**CG**gc**GC**
..((((((.((...)))))))

(-18.00)

CccUGgccGCA**ag**GC**c**Ag**Gg**
(((((.((...)))))))

(-20.40)

CcGUGgccGC**gag**GC**c**Ac**Gg**
(((((.((...)))))))

(-19.10)

Classic: 0 mutation

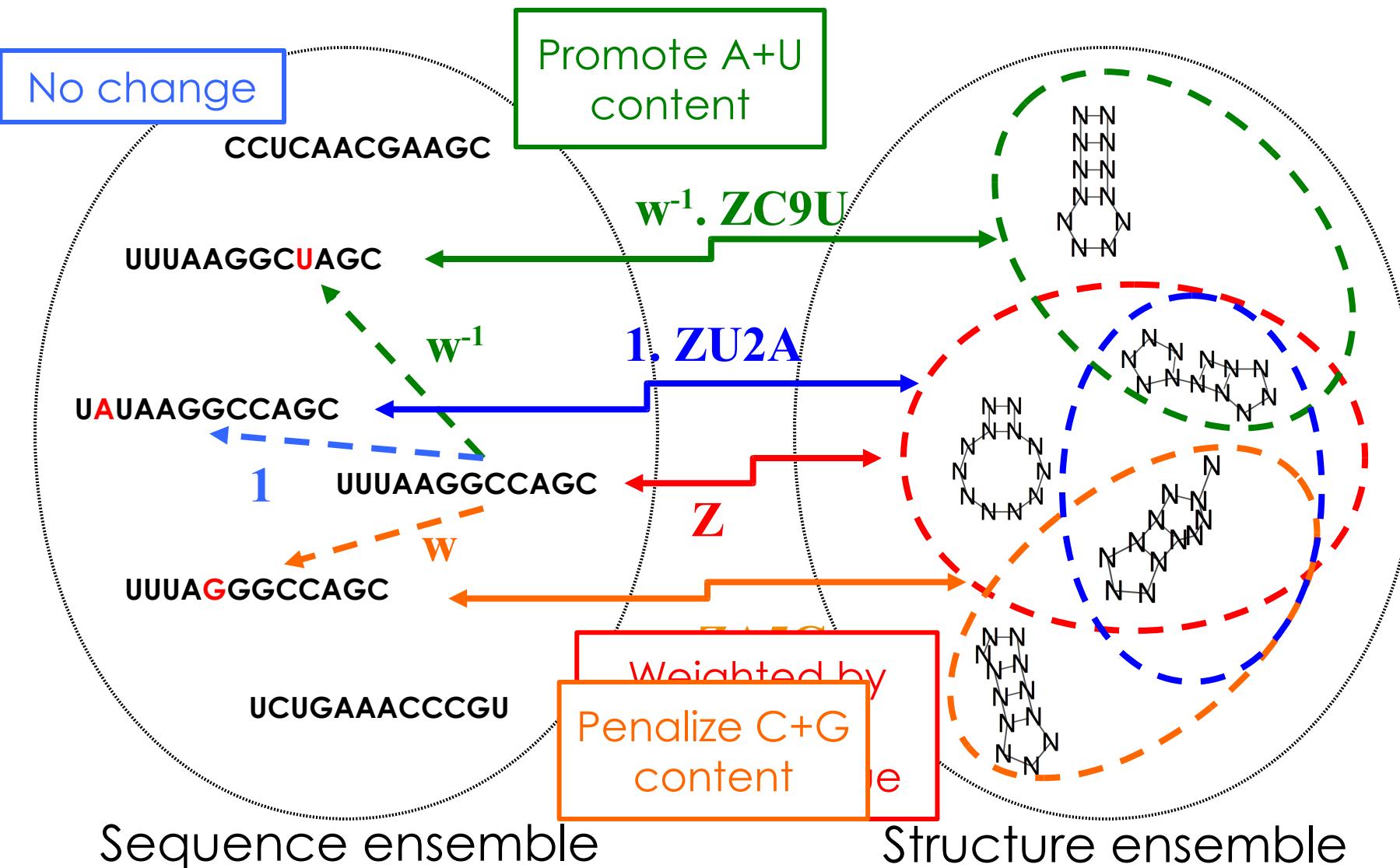
RNAmutants: 1 mutation

RNAmutants: 10 mutations

Sample k mutations increasing the folding energy
Consequence: Increase in C+G content

- Background: RNA mutants in a nutshell
Algorithms to sample RNA secondary structures and mutations.
- **Our approach: Adaptive sampling**
Uniformly shifting the distribution of samples.
- Results: Evolutionary studies
Insights on the evolutionary pressure stemming from an optimization of the thermodynamical stability.

Our approach: Weighting mutations



Weighting recursive equations

$$\begin{array}{c} z \\ \hline x \quad y \\ i \quad j \end{array} = \begin{array}{c} z \\ \hline x \quad u \quad v \quad y \\ i \quad r-1 \quad r \quad j \end{array} + \begin{array}{c} z \\ \hline x \quad u \quad y \\ i \quad j-1 \quad j \end{array} \times W(j,y)$$

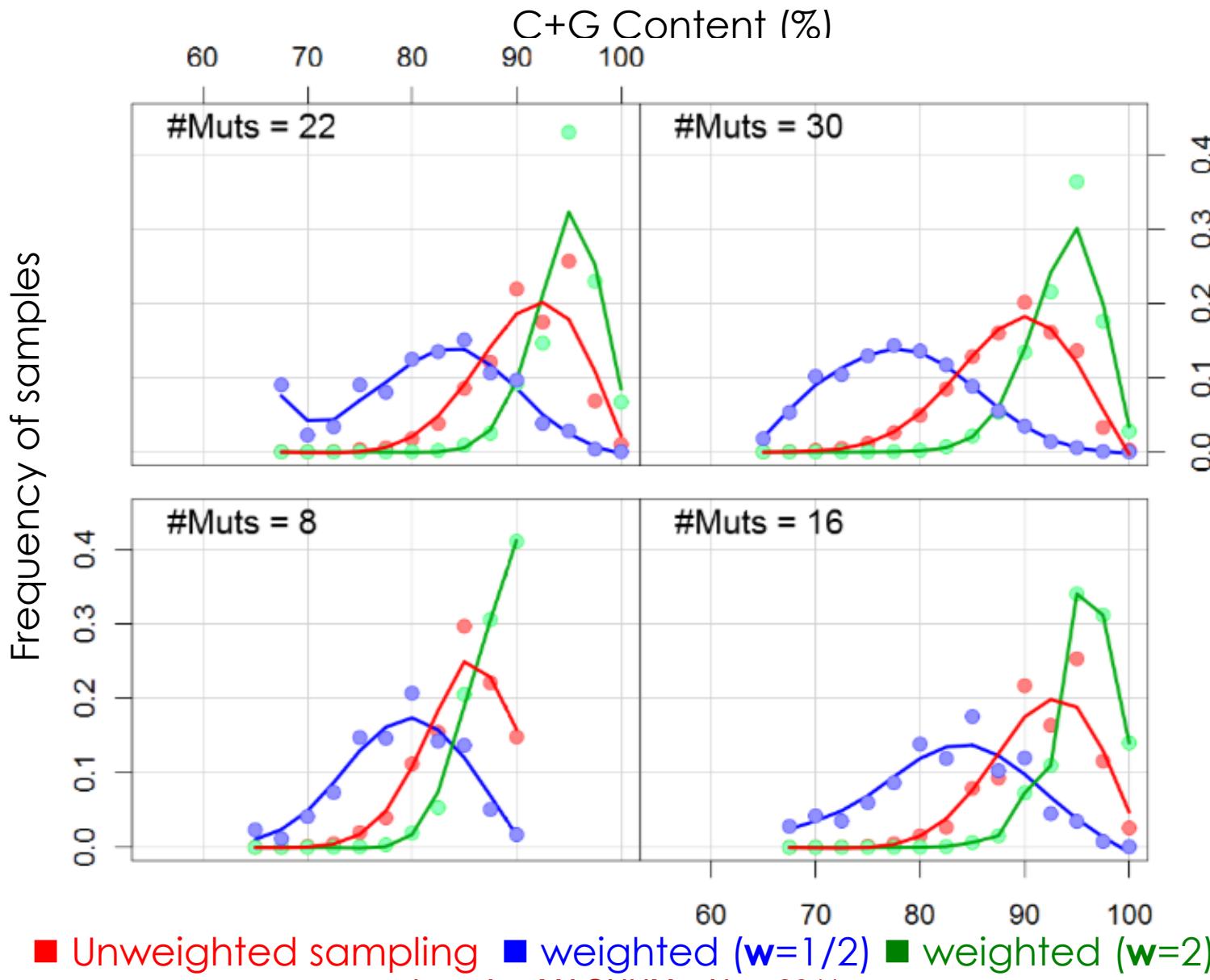
$$\begin{array}{c} z^B \\ \hline x \quad y \\ i \quad j \end{array} = \left(\begin{array}{c} z^B \\ \hline x \quad u \quad v \quad y \\ i \quad i+1 \quad j-1 \quad j \end{array} + \begin{array}{c} z^B \\ \hline x \quad u \quad v \quad y \\ i \quad r \quad s \quad j \end{array} + \begin{array}{c} z^M \\ \hline x \quad u \quad v \quad w \quad z \quad y \\ i \quad i+1 \quad r-1 \quad r \quad j-1 \quad j \end{array} \right) \times W(i,x) \times W(j,y)$$

$$\begin{array}{c} z^{M1} \\ \hline x \quad y \\ i \quad j \end{array} = \begin{array}{c} z^B \\ \hline x \quad u \quad v \quad y \\ i \quad r \quad j \end{array}$$

$$W(i,x) = \begin{cases} w & \text{If } A,U \rightarrow C,G \\ w^{-1} & \text{If } C,G \rightarrow A,U \\ 1 & \text{Otherwise} \end{cases}$$

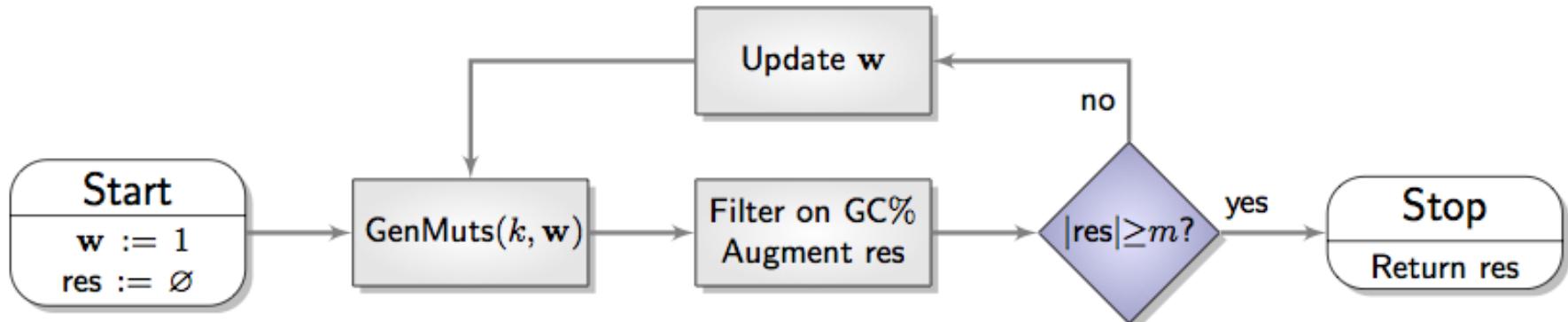
$$\begin{array}{c} z^M \\ \hline x \quad y \\ i \quad j \end{array} = \begin{array}{c} z^{M1} \\ \hline x \quad u \quad v \quad y \\ i \quad r \quad j \end{array} + \begin{array}{c} z^M \\ \hline x \quad u \quad v \quad y \\ i \quad r-1 \quad r \quad j \end{array}$$

Effect of weighted sampling



An Adaptive Sampling approach

Idea: Figure out a suitable weight on the fly.



- Keep all samples at targeted C+G% and reject others.
- Update w at each iteration using a bisection strategy.
- Stop when enough samples have been stored.

Algorithmic aspects

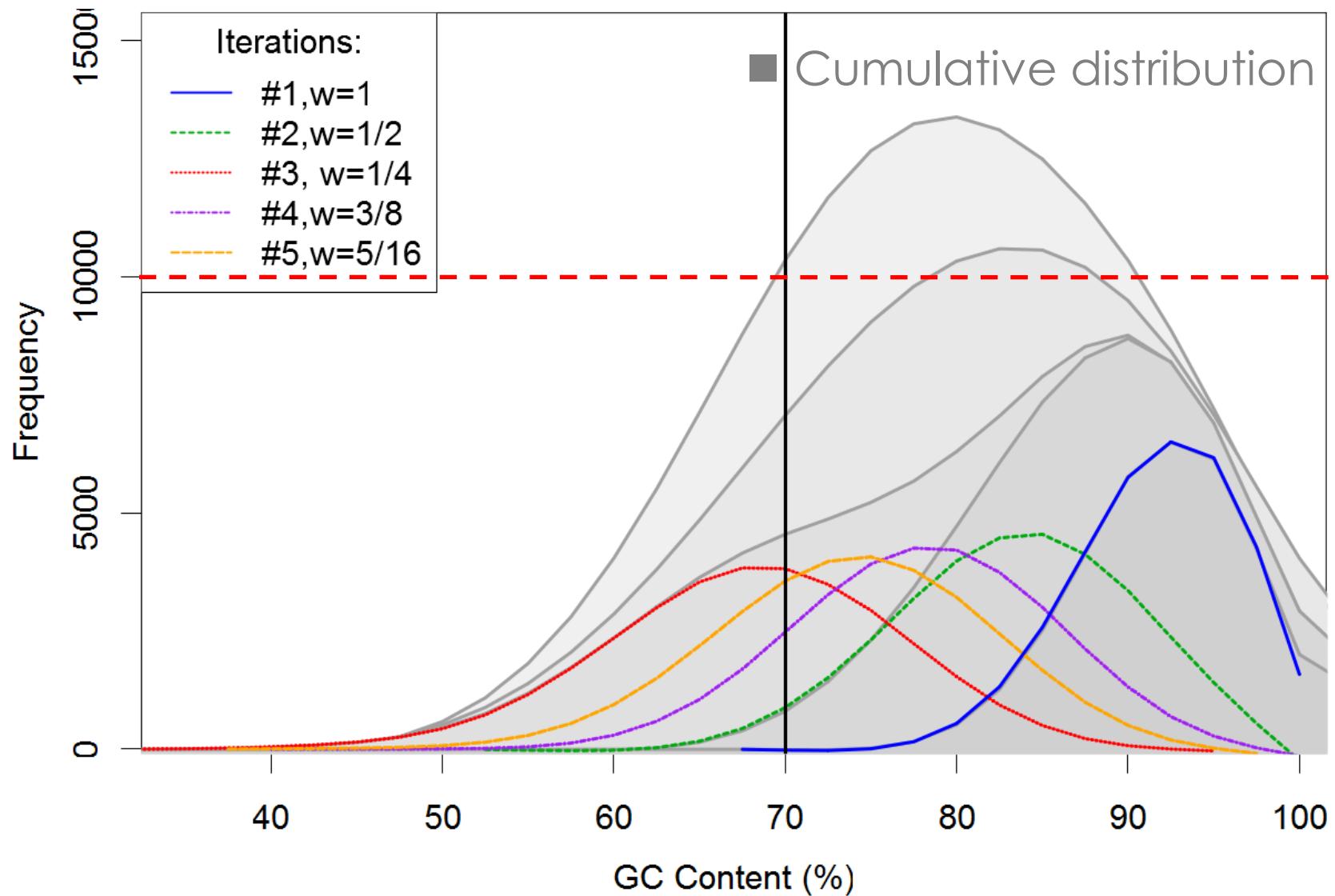
- After rejection, the weighted schema only impacts performance, not the probability distribution.
- Partition function can be written as a polynomial:

$$Z = \sum_{i=0}^n a_i \times w^i$$

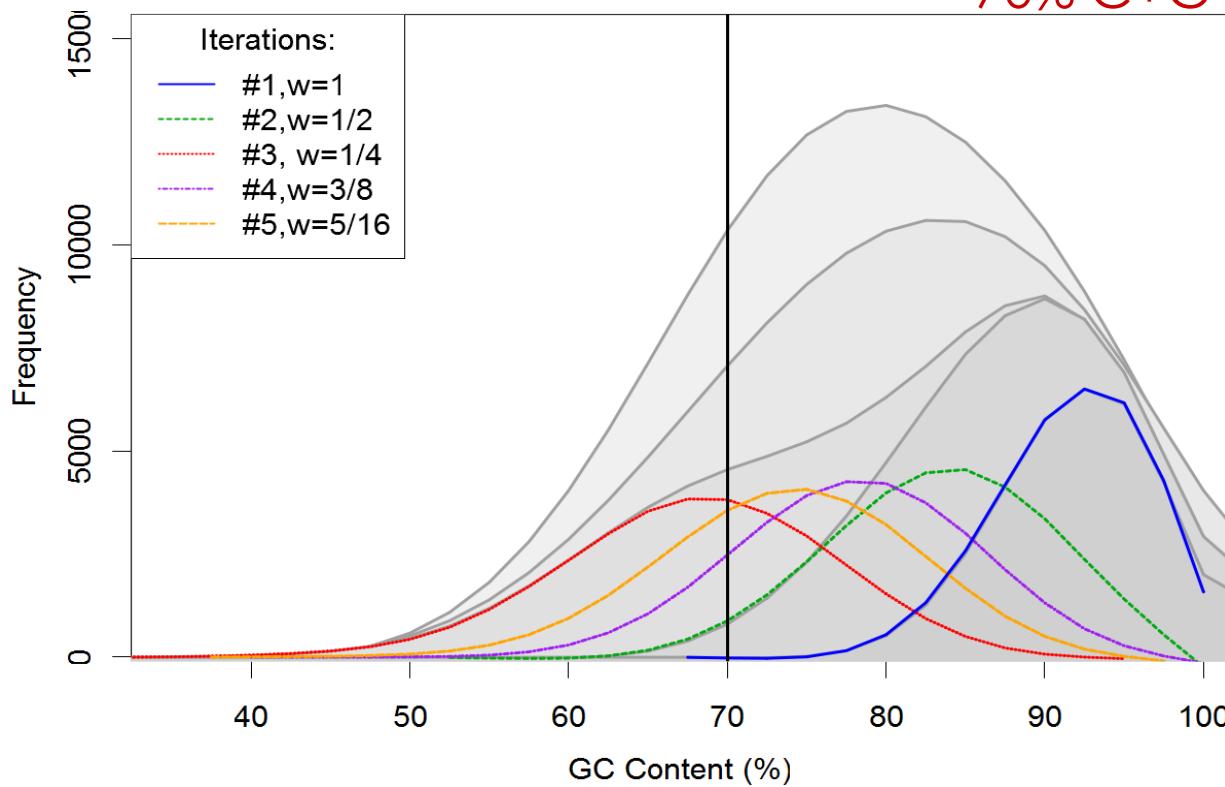
After n iterations one can compute all a_i through interpolation and *invert* the polynomial to compute optimal weight w (Gröbner base).

Remark: In practice, few iterations are necessary

Example: 40 nt., 10000 samples, 30 mutations, 70% C+G content



Example: 40 nt., 10000 samples, 30 mutations,
70% C+G content



Overall, 150 000 samples required.

Naïve rejection: $\sim 30\ 000\ 000$ samples!

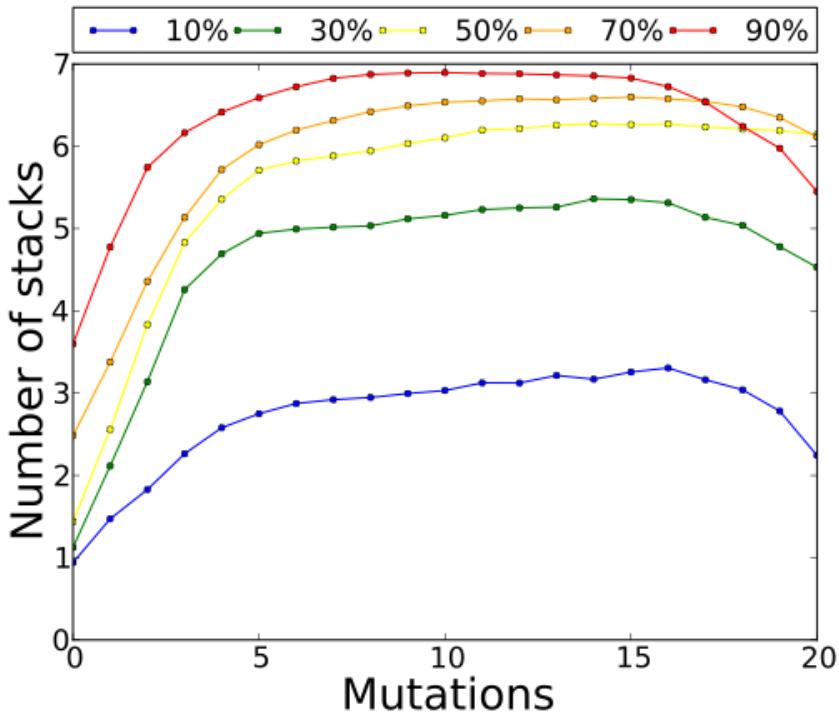
+ Ratio grows exponentially on length (gen. Fun.)

Note: Dynamic programming alternative in $O(n^7)$

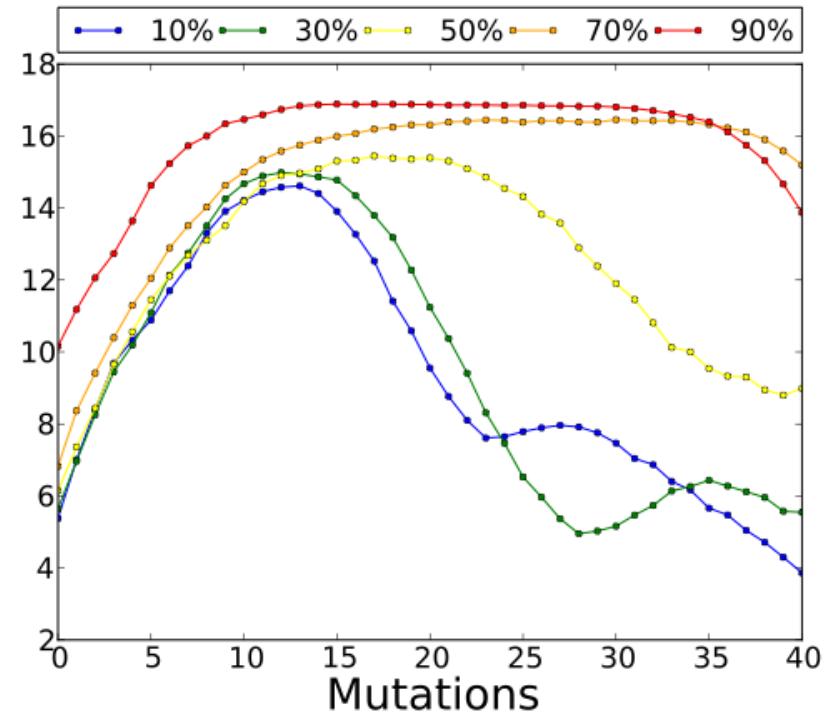
- Background: RNA mutants in a nutshell
Algorithms to sample RNA secondary structures and mutations.
- Our approach: Adaptive sampling
Uniformly shifting the distribution of samples.
- **Results: Evolutionary studies**
Insight into the evolutionary pressure stemming from an optimization of the thermodynamical stability.

Low CG-contents favor structural diversity

Simulation at **fixed G+C content** from random seeds



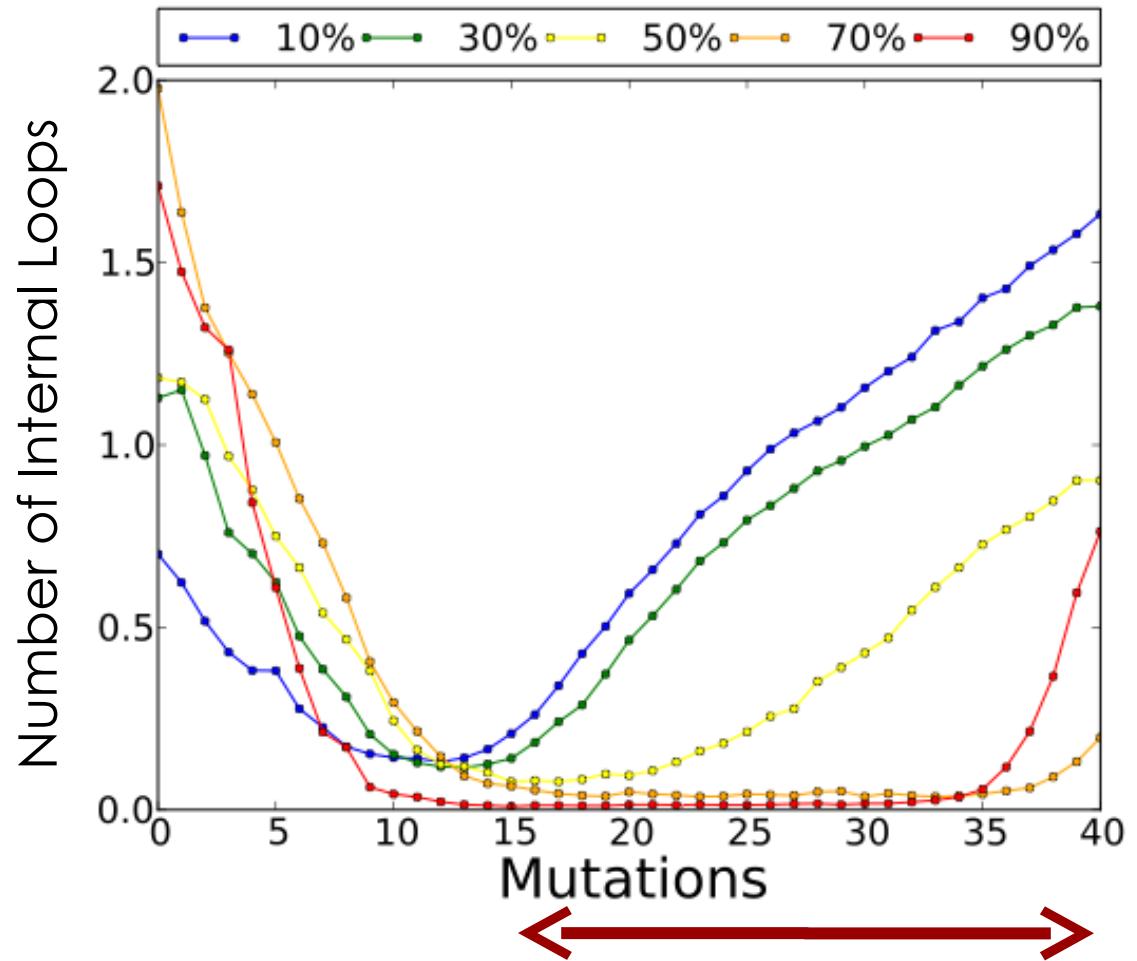
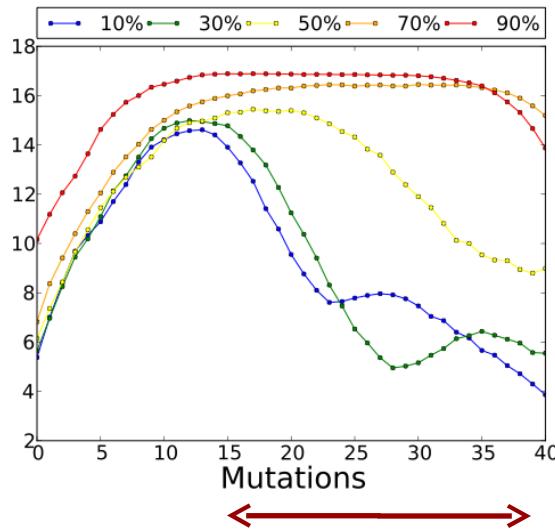
20 nucleotides



40 nucleotides

■ 10% ■ 30% ■ 50% ■ 70% ■ 90%

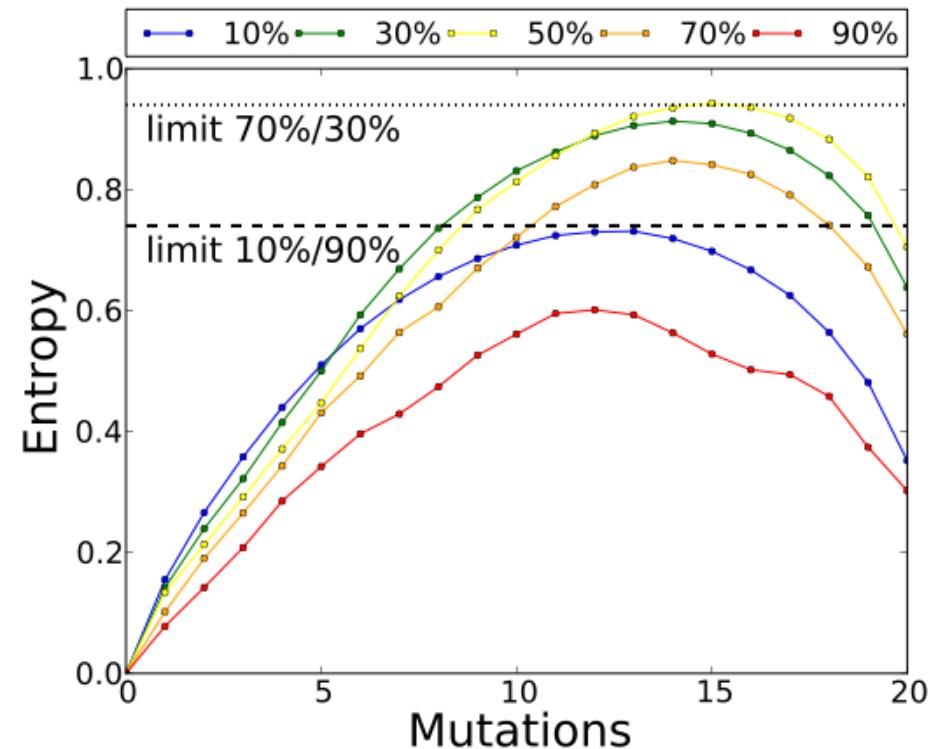
Low C+G contents favor internal loop insertion



■ 10% ■ 30% ■ 50% ■ 70% ■ 90%

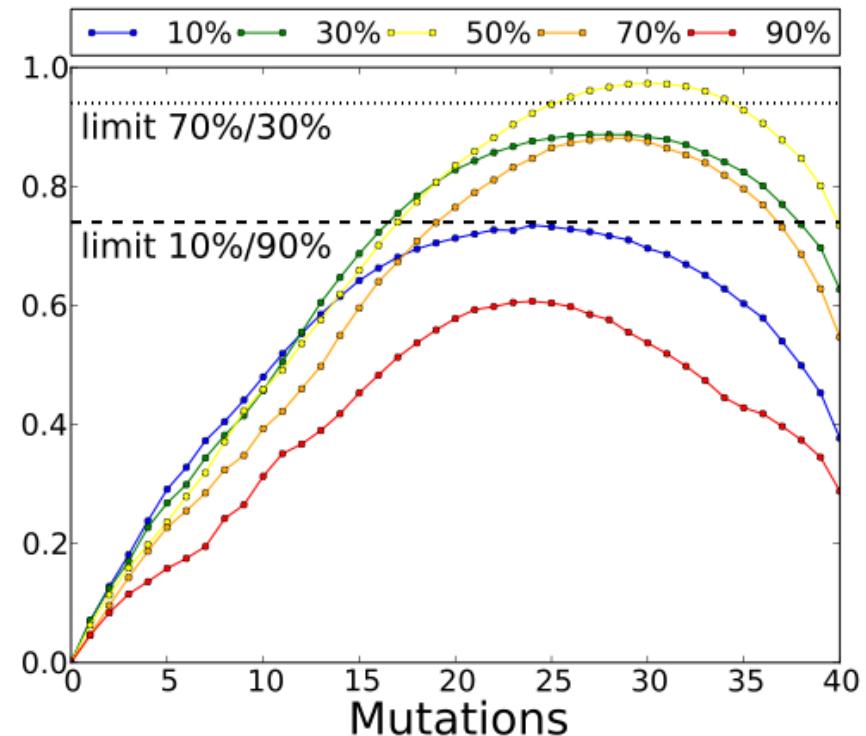
High G+C-contents reduce evolutionary accessibility

Simulation at **fixed G+C content** from random seeds



20 nucleotides

■ 10% ■ 30% ■ 50% ■ 70% ■ 90%



40 nucleotides

- More studies of Sequence-Structure maps.
- Applications to RNA design.
- Same techniques can be applied to other parameters (e.g. number of base pairs).
- Can be generalized to multiple parameters.

Acknowledgments



Jérôme Waldispühl
Assistant Prof. McGill, Montreal Canada.

INRIA

- Philippe Flajolet

MIT

- Bonnie Berger
- Srinivas Devadas
- Mieszko Lis
- Alex Levin
- Charles W. O'Donnell

Google Inc.

- Behshad Behzadi

Ecole Polytechnique

- Jean-Marc Steyaert

Boston College

- Peter Clote

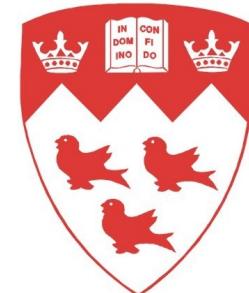


University of Paris 6

- Olivier Bodini

University of Paris 11

- Alain Denise



Would you like to know more? ■

- O. Bodini, and Y. Ponty
Multi-dimensional Boltzmann Sampling of Languages,
Proceedings of AOFA'10, 49--64, 2010
- J. Waldspühl, S. Devadas, B. Berger and P. Clote,
Efficient Algorithms for Probing the RNA Mutation Landscape,
Plos Computational Biology, 4(8):e1000124, 2008.

ht tp://csb.cs.mcgill.ca/RNAnutants