# (Positive) Design of RiboNucleic Acids

Yann Ponty[*,•,†]

[*] Centre National de la Recherche Scientifique (CNRS)

[•] LIX, Ecole Polytechnique

[†] AMIBio team, Inria Saclay

`http://goo.gl/mejsFh`

# Yann Ponty

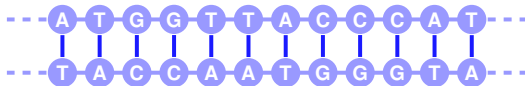PhD in Computer Science, Université Paris-Sud (France)

- CNRS research scientist
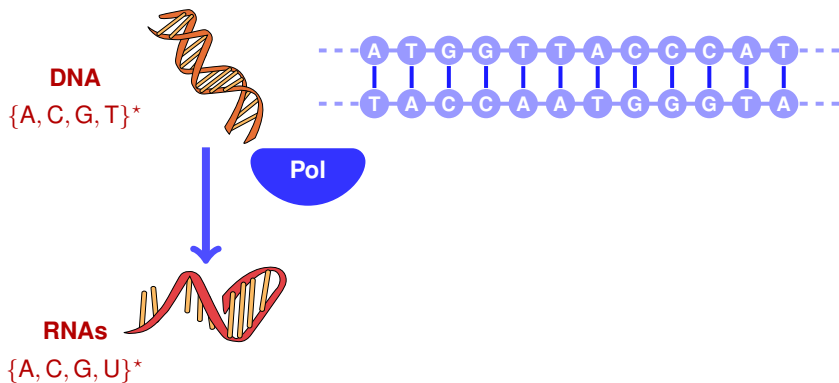- Faculty at LIX, Computer Science department of Ecole Polytechnique

- Group leader – AMIBio team (Ecole Polytechnique and Inria Saclay)
- Postdoc experience in RNA Computational Biology (Boston, Paris) and Discrete Mathematics (Paris)
- Extended sabbatical at Simon Fraser University (Vancouver, Canada)

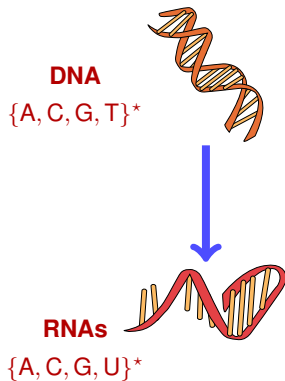# Fundamental *dogma* of molecular biology

**DNA**

$\{A, C, G, T\}^{\star}$



```
A T G G T T A C C C A T
| | | | | | | | | | | |
T A C C A A T G G G T A
```

# Fundamental *dogma* of molecular biology

**DNA**
{A, C, G, T}*



**Pol**

**RNAs**
{A, C, G, U}*

# Fundamental *dogma* of molecular biology

**DNA**
$\{A, C, G, T\}^\star$

**RNAs**
$\{A, C, G, U\}^\star$



A T G G T T A C C C A T

Pol

A

C C A A T G G G T A

# Fundamental *dogma* of molecular biology

**DNA**

$\{A, C, G, T\}^{\star}$

**RNAs**

$\{A, C, G, U\}^{\star}$

A T G G T T A C C C A T

Pol

C A A T G G G T A

A U

# Fundamental *dogma* of molecular biology



**DNA**
$\{A, C, G, T\}^\star$

**RNAs**
$\{A, C, G, U\}^\star$

# Fundamental *dogma* of molecular biology

**DNA**
{A, C, G, T}*

**RNAs**
{A, C, G, U}*

A T G G T T A C C C A T

T A C C A A T G G G **Pol**

A U G G U U A C C C A U

# Fundamental *dogma* of molecular biology



**DNA**
$\{A, C, G, T\}^\star$

**RNAs**
$\{A, C, G, U\}^\star$

# Fundamental *dogma* of molecular biology

**DNA**
$\{A, C, G, T\}^\star$

A T G G T T A C C C A T

T A C C A A T G G G T A

**RNAs**
$\{A, C, G, U\}^\star$

A U G G U U A C C C A U

**Ribosome**

**Proteins**
$\{Ala, Arg, \ldots, Val\}^\star$

$20^+$ **Amino acids**

# Fundamental *dogma* of molecular biology



**DNA**
$\{A, C, G, T\}^\star$

**RNAs**
$\{A, C, G, U\}^\star$

**Ribosome**

**Met**

**Proteins**
$\{Ala, Arg, \ldots, Val\}^\star$

$20^+$ **Amino acids**

# Fundamental *dogma* of molecular biology



**DNA**
$\{A, C, G, T\}^\star$

**RNAs**
$\{A, C, G, U\}^\star$

**Proteins**
$\{Ala, Arg, \ldots, Val\}^\star$

$20^+$ **Amino acids**

# Fundamental *dogma* of molecular biology



**DNA**
$\{A, C, G, T\}^{\star}$

**RNAs**
$\{A, C, G, U\}^{\star}$
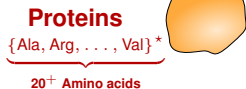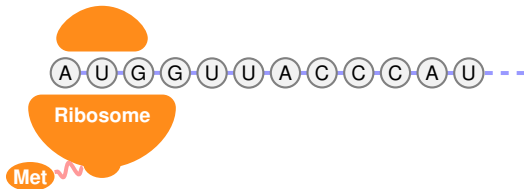
**Proteins**
$\{Ala, Arg, \ldots, Val\}^{\star}$

**$20^{+}$ Amino acids**

# Fundamental *dogma* of molecular biology



**DNA**
{A, C, G, T}*

**RNAs**
{A, C, G, U}*

**Proteins**
{Ala, Arg, . . . , Val}*

$20^{+}$ **Amino acids**

A T G G T T A C C C A T
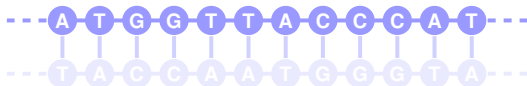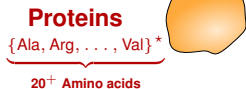T A C C A A T G G G T A

A U G G U U A C C C A U

**Ribosome**

Met Val Thr His

# Fundamental *dogma* of molecular biology

**DNA**
{A, C, G, T}*



A T G G T T A C C C A T

T A C C A A T G G G T A

**RNAs**
{A, C, G, U}*



A U G G U U A C C C A U

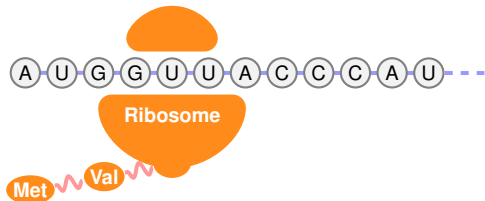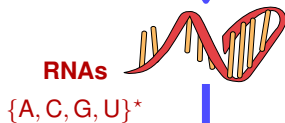**Proteins**
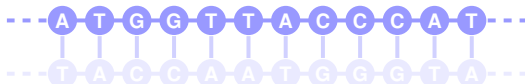{Ala, Arg, . . . , Val}*

20+ Amino acids



Met Val Thr His Ile Leu His Asn

# Fundamental *dogma* of molecular biology

**THE CODE**
**(genes)**
**DNA**
{A, C, G, T}*

A–T–G–G–T–T–A–C–C–C–A–T
T–A–C–C–A–A–T–G–G–G–T–A

**RNAs**
{A, C, G, U}*

A–U–G–G–U–U–A–C–C–C–A–U

**Proteins**
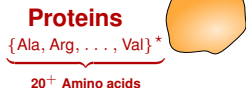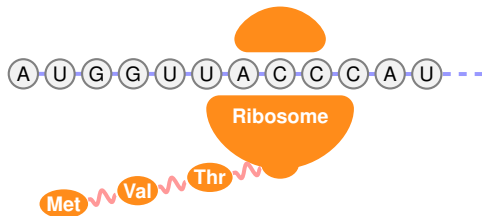{Ala, Arg, . . . , Val}*
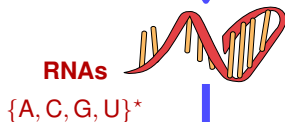
**20⁺ Amino acids**

Met – Val – Thr – His – Ile – Leu – His – Asn

# Fundamental *dogma* of molecular biology

**THE CODE**
**(genes)**

**DNA**
$\{A, C, G, T\}^{\star}$

A T G G T T A C C C A T
T A C C A A T G G G T A

**RNAs**
$\{A, C, G, U\}^{\star}$

A U G G U U A C C C A U

**THE MACHINE**
**(enzymes)**

**Proteins**
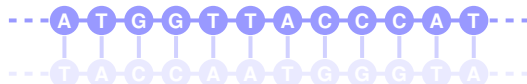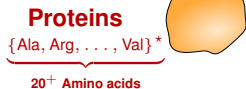$\{Ala, Arg, \ldots, Val\}^{\star}$
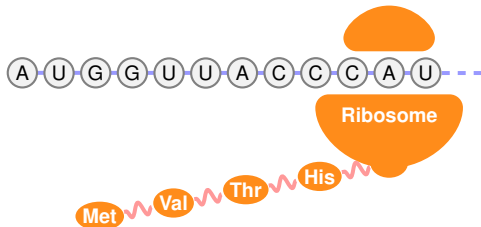
$20^{+}$ **Amino acids**

Met Val Thr His Ile Leu His Asn

# Fundamental *dogma* of molecular biology

**THE CODE**
**(genes)**

**DNA**

$\{A, C, G, T\}^{\star}$

A–T–G–G–T–T–A–C–C–C–A–T

T–A–C–C–A–A–T–G–G–G–T–A

**MEH. . .**

**RNAs**

$\{A, C, G, U\}^{\star}$

A–U–G–G–U–U–A–C–C–C–A–U

**THE MACHINE**
**(enzymes)**

**Proteins**

$\{Ala, Arg, \ldots, Val\}^{\star}$

**20$^{+}$ Amino acids**

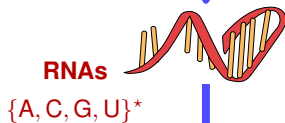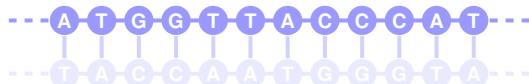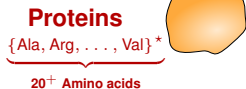Met – Val – Thr – His – Ile – Leu – His – Asn

# Fundamental *dogma* of molecular biology



**DNA**

Transcription
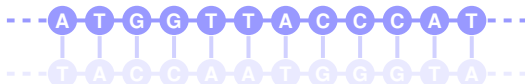
**RNA**

Translation

**Proteins**

# Fundamental *dogma* of molecular biology (v2.0)



DNA

Transfer

Transcription

Maturation

Carrier

RNA

Participates

Regulation

Translation

Synthesis

**RNA functions**
- Messenger
- Translation
- Regulation
- Enzyme
- Catalytic
- …

Proteins

# Fundamental *dogma* of molecular biology



#RNA Functional Families (RFAM DB)

RNA functions
- Messenger
- Translation
- Regulation
- Enzyme
- Catalytic
- …

Proteins

# RNA world: Resolving the *chicken vs egg* paradox at the origin of life...



A gene big enough to specify an enzyme would be too big to replicate accurately without the aid of an enzyme of the very kind that it is trying to specify. So the system *apparently cannot get started*.

[···] This is the RNA World. To see how plausible it is, we need to look at why proteins are good at being enzymes but bad at being replicators; at why DNA is good at replicating but bad at being an enzyme; and finally why *RNA might just be good enough at both roles to break out of the Catch-22.*

**R. Dawkins**. *The Ancestor's Tale: A Pilgrimage to the Dawn of Evolution*

# RNA world: Resolving the *chicken vs egg* paradox at the origin of life...



A gene big enough to specify an enzyme would be too big to replicate accurately without the aid of an enzyme of the very kind that it is trying to specify. So the system *apparently cannot get started*.

[· · · ] This is the RNA World. To see how plausible it is, we need to look at why proteins are good at being enzymes but bad at being replicators; at why DNA is good at replicating but bad at being an enzyme; and finally why *RNA might just be good enough at both roles to break out of the Catch-22.*

**R. Dawkins**. *The Ancestor's Tale: A Pilgrimage to the Dawn of Evolution*

# RNA folding

RNA are **single-stranded** and **fold** on themselves, establishing **complex 3D structures** that are **essential to their function(s)**.

RNA structures are stabilized by **base-pairs**, each mediated by **hydrogen bonds**.



G/C

U/A

U/G

Watson/Crick base-pairs

Wobble base-pair

**Canonical base-pairs**

# RNA Design

**RNA** = Linear Polymer = Sequence in $\{A, C, G, U\}^\star$



```
UUAGGCGGCCACAGC
GGUGGGGUUGCCUCC
CGUACCCAUCCCGAA
CACGGAAGAUAAGCC
CACCAGCGUUCCGGG
GAGUACUGGAGUGCG
CGAGCCUCUGGGAAA
CCCGGUUCGCCGCCA

CC
```

Primary Structure          Secondary Structure          Structure Tertiaire

5s rRNA (PDBID: 1K73:B)

# Evolution of RNAs

**Homologous** genes = **Functionally** equivalent, within or across organisms
Usually well-captured by **sequence similarity** in proteins, binding sites...

**Problem:** Many classes of non-(protein) coding RNAs (ncRNAs) poorly
conserved at the sequence level **but** adopt a **conserved structure!**



RFAM Bacterial RNAse-P class B Alignment

## RNA Design

**RNA** = Linear Polymer = Sequence in $\{A, C, G, U\}^\star$



UUAGGCGGCCACAGC
GGUGGGGUUGCCUCC
CGUACCCAUCCCGAA
CACGGAAGAUAAGCC
CACCAGCGUUCCGGG
GAGUACUGGAGUGCG
CGAGCCUCUGGGAAA
CCCGGUUCGCCGCCA

CC

Primary Structure        Secondary Structure        Structure Tertiaire

5s rRNA (PDBID: 1K73:B)

# RNA Design



**RNA** = Linear Polymer = Sequence in $\{A, C, G, U\}^\star$

```
UUAGGCGGCCACAGC
GGUGGGGUUGCCUCC
CGUACCCAUCCCGAA
CACGGAAGAUAAGCC
CACCAGCGUUCCGGG
GAGUACUGGAGU
CGAGCCUCU
CCCGGUUCGCCGCC

CC
```

**Structure Prediction**

**RNA Design**

Primary Structure        Secondary Structure        Structure Tertiaire

5s rRNA (PDBID: 1K73:B)
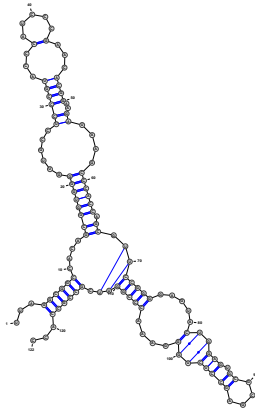
# Why we design RNAs

- **To create building blocks for synthetic systems**
  Rationally-designed RNAs increase orthogonality

- To assess the significance of observed phenomenon
  Random models should include every established characters…
  …including adoption of a single structure

- To test/push our understanding of how RNA folds
  Misfolding RNAs reveal gaps in our energy models and descriptors for the
  conformational spaces

- To help search for homologous sequences
  Incomplete covariance models hindered by limited training sets
  Design can be used to generalize existing alignments

- To fuel RNA-based therapeutics
  Sequence-based (siRNA, synthetic genes), but structure matters

- To perform controlled experiments

# Why we design RNAs

- **To create building blocks for synthetic systems**
  Rationally-designed RNAs increase orthogonality

- **To assess the significance of observed phenomenon**
  Random models should include every established characters...
  ...including adoption of a single structure

- To test/push our understanding of how RNA folds
  Misfolding RNAs reveal gaps in our energy models and descriptors for the
  conformational spaces

- To help search for homologous sequences
  Incomplete covariance models hindered by limited training sets
  Design can be used to generalize existing alignments

- To fuel RNA-based therapeutics
  Sequence-based (siRNA, synthetic genes), but structure matters

- To perform controlled experiments

# Why we design RNAs

- **To create building blocks for synthetic systems**
  Rationally-designed RNAs increase orthogonality

- **To assess the significance of observed phenomenon**
  Random models should include every established characters...
  ...including adoption of a single structure

- **To test/push our understanding of how RNA folds**
  Misfolding RNAs reveal gaps in our energy models and descriptors for the conformational spaces

- To help search for homologous sequences
  Incomplete covariance models hindered by limited training sets
  Design can be used to generalize existing alignments

- To fuel RNA-based therapeutics
  Sequence-based (siRNA, synthetic genes), but structure matters

- To perform controlled experiments

## Why we design RNAs

- To create building blocks for synthetic systems
  Rationally-designed RNAs increase orthogonality

- To assess the significance of observed phenomenon
  Random models should include every established characters. . .
  . . . including adoption of a single structure

- To test/push our understanding of how RNA folds
  Misfolding RNAs reveal gaps in our energy models and descriptors for the
  conformational spaces

- To help search for homologous sequences
  Incomplete covariance models hindered by limited training sets
  Design can be used to generalize existing alignments

- To fuel RNA-based therapeutics
  Sequence-based (siRNA, synthetic genes), but structure matters

- To perform controlled experiments

## Why we design RNAs

- **To create building blocks for synthetic systems**
  Rationally-designed RNAs increase orthogonality

- **To assess the significance of observed phenomenon**
  Random models should include every established characters. . .
  . . . including adoption of a single structure

- **To test/push our understanding of how RNA folds**
  Misfolding RNAs reveal gaps in our energy models and descriptors for the
  conformational spaces

- **To help search for homologous sequences**
  Incomplete covariance models hindered by limited training sets
  Design can be used to generalize existing alignments

- **To fuel RNA-based therapeutics**
  Sequence-based (siRNA, synthetic genes), but structure matters

- **To perform controlled experiments**

# Why we design RNAs

- **To create building blocks for synthetic systems**
  Rationally-designed RNAs increase orthogonality

- **To assess the significance of observed phenomenon**
  Random models should include every established characters...
  ...including adoption of a single structure

- **To test/push our understanding of how RNA folds**
  Misfolding RNAs reveal gaps in our energy models and descriptors for the conformational spaces

- **To help search for homologous sequences**
  Incomplete covariance models hindered by limited training sets
  Design can be used to generalize existing alignments

- **To fuel RNA-based therapeutics**
  Sequence-based (siRNA, synthetic genes), but structure matters

- **To perform controlled experiments**

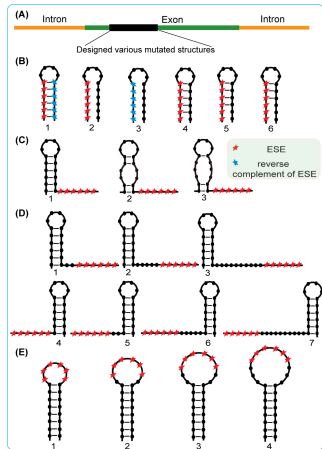# Controlled experiments through RNA design

**Motivation:** Quantifying the impact of structure $S$ on efficacy of a single Exon Splicing Enhancers (ESE):

- Presence of ESE motif $E$;
- Different structures $S_1$, $S_2$...;
- Avoid library of ($\sim$1500!) documented ESEs motifs.

> **Objectives.** Design RNA which:
> 1. Folds into a prescribed structure;
> 2. Features/avoids motifs.
> 3. Control GC%, Boltz. prob.....

Structural context of ESE motif in transcript was shown to affect its functionality. [Liu *et al*, FEBS Lett. 2010]

# Design objectives

## Positive structural design
Optimize **affinity** of designs towards target structure(s)
**Examples:** Most stable sequence for given fold...

## Negative structural design
Limit affinity of designs towards **alternative structures**
**Examples:** Lowest free-energy, High Boltzmann probability/Low entropy...

## Additional constraints:
- **Forbid** motif list to appear **anywhere** in design
- **Force** motif list to appear **each at least once**
- **Limit** available alternatives at certain positions
- **Control** overall composition (GC-content)

**Outline**

- I. Single Structure Design (IncaRNAtion)

- II. Constrained Design using Formal Languages
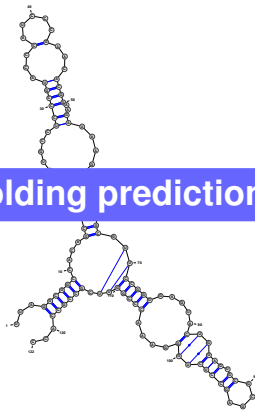
- III. Multiple Structures

# I. Inverse Folding

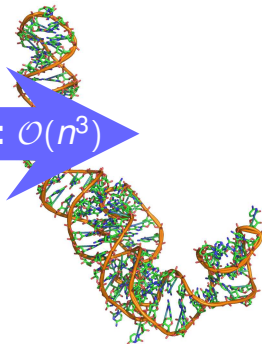**Designing a given structure**

# RNA sequence and structure(s)

**RNA** = Linear Polymer = Sequence in $\{A, C, G, U\}^\star$



**MFE folding prediction:** $\mathcal{O}(n^3)$

```
UUAGGCGGCCACAGC
GGUGGGGUUGCCUCC
CGUACCCAUCCCGAA
CACGGAAGAUAAGCC
CACCAGCGUUCCGGG
GAGUACUGGAGUGCG
CGAGCCUCUGGGAAA
CCCGGUUCGCCGCCA

CC
```

Primary Structure        Secondary Structure        Tertiary Structure

5s rRNA (PDBID: 1K73:B)

# Crossing interactions

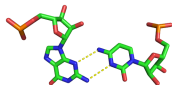Excluded from the secondary structure:

- **Non-canonical base-pairs:**
  Any base-pair **other than** {(A-U), (C-G), (G-U)}
  **OR** interacting in a non-standard way (WC/WC-Cis) [Leontis Westhof, RNA 2001].
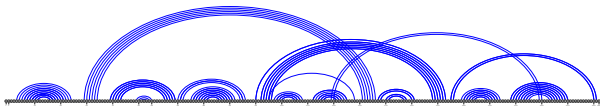


Canonical CG base-pair (WC/WC-Cis)

Non-canonical base-pair (Sugar/WC-Trans)

- **(Pseudo?)knots:** Crossing sets of nested stable base-pairs


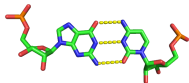
Group I Ribozyme (PDBID: 1Y0Q:A)

# Crossing interactions

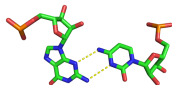Excluded from the secondary structure:

- **Non-canonical base-pairs:**
  Any base-pair **other than** {(A-U), (C-G), (G-U)}
  **OR** interacting in a non-standard way (WC/WC-Cis) **[Leontis Westhof, RNA 2001]**.
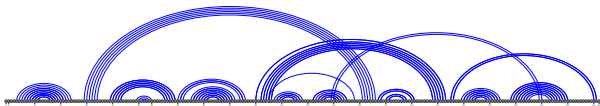


Canonical CG base-pair (WC/WC-Cis)     Non-canonical base-pair (Sugar/WC-Trans)

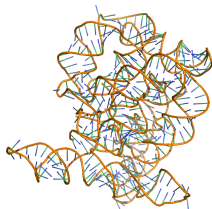- **(Pseudo?)knots:** Crossing sets of nested stable base-pairs



Group I Ribozyme (PDBID: 1Y0Q:A)

# Crossing interactions

Excluded from the secondary structure:

- **Non-canonical base-pairs:**
  Any
  **OR**
  **200**



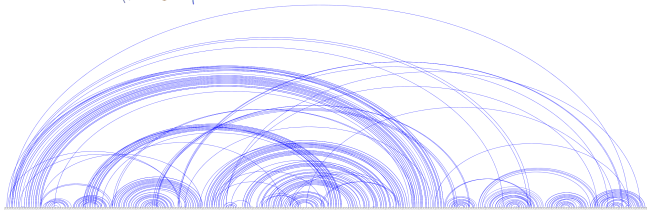**Crossing** interactions **do exist**!

**Example:** Group II Intron (PDB ID: 3IGI)

**But** are **hard** to predict
[Lyngsoe-ICALP'04]
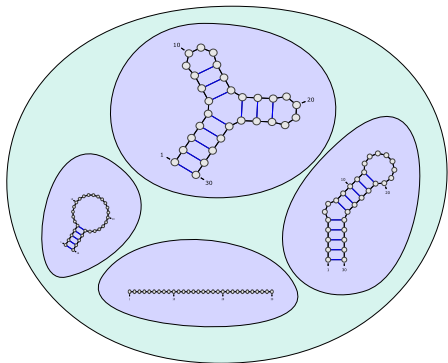[Sheikh Backofen Ponty, CPM'12]

- **(Ps**

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
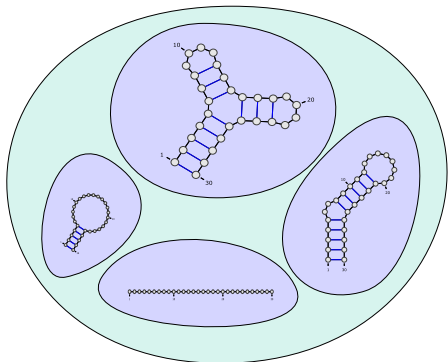- **2010s–????** Embracing the kinetics of RNA folding



mRNA half-life: $\sim$7h
(Mouse [Sharova2009])

$T \to \infty$

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
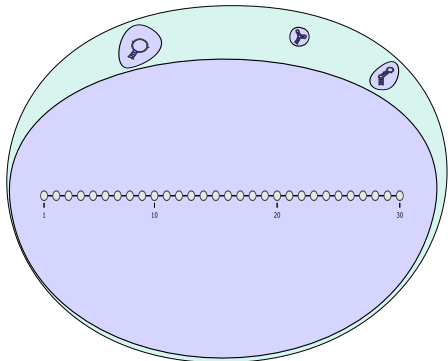- **2010s–????** Embracing the kinetics of RNA folding



mRNA half-life: ∼7h
(Mouse [Sharova2009])

$T \to \infty$

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
- **2010s–????** Embracing the kinetics of RNA folding



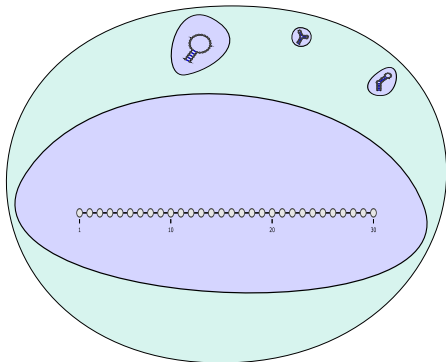mRNA half-life: $\sim$7h
(Mouse [Sharova2009])

$T = 0$h

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
- **2010s–????** Embracing the kinetics of RNA folding



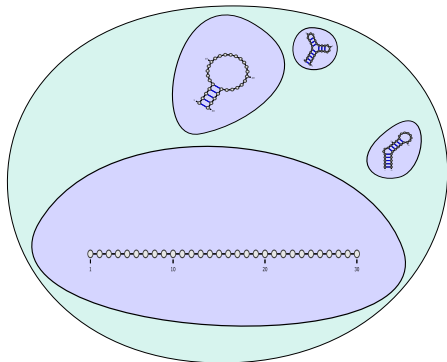mRNA half-life: $\sim$7h
(Mouse [Sharova2009])

$T = 1$h

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
- **2010s–????** Embracing the kinetics of RNA folding



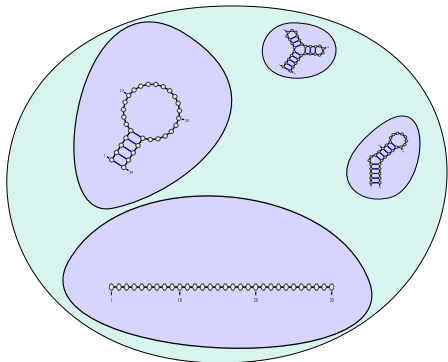mRNA half-life: ∼7h
(Mouse [Sharova2009])

$T = 2h$

1st International Computational Biology workshop

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
- **2010s–????** Embracing the kinetics of RNA folding



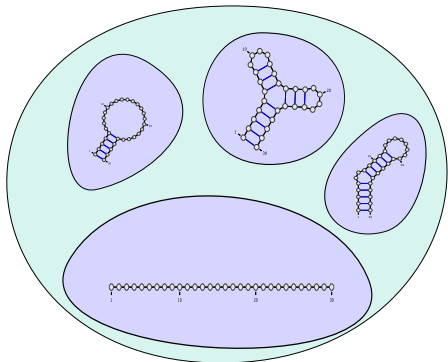mRNA half-life: $\sim$7h
(Mouse [Sharova2009])

$T = 5$h

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
- **2010s–????** Embracing the kinetics of RNA folding



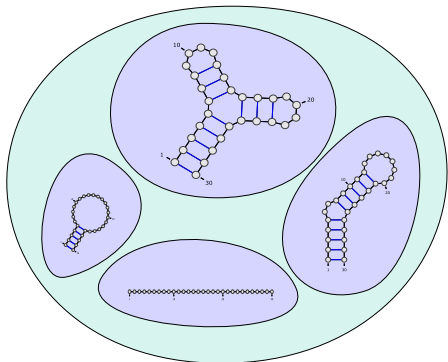mRNA half-life: $\sim$7h
(Mouse [Sharova2009])

$T = 10$h

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
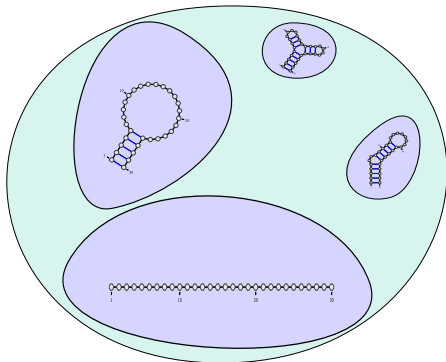- **2010s–????** Embracing the kinetics of RNA folding



mRNA half-life: $\sim$7h
(Mouse [Sharova2009])

$T \to \infty$

# Thermodynamics vs Kinetics

## Paradigms for RNA structure prediction

- **1978–1990s** Most probable structure = Minimal Free-Energy (MFE)
- **1990s–2010s** Functional structure(s) = Boltzmann ensemble (partition function)
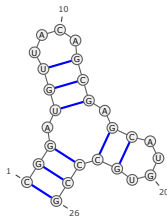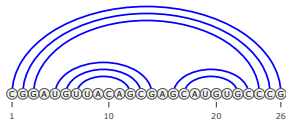- **2010s–????** Embracing the kinetics of RNA folding
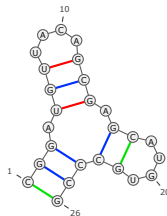


mRNA half-life: $\sim$7h
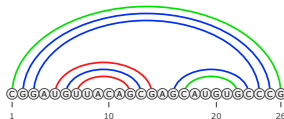(Mouse [Sharova2009])

$T = 10$h

# Problem statement



- **RNA structure** $S$**:** Non-crossing base-pairs for positions in sequence $w$
- **Motifs:** Sequence/structure features (e.g. Base-pairs, Stacks, Loops... )
- **Energy model:**
    **Motif** $\rightarrow$ Free-energy contribution $\Delta(\cdot) \in \mathbb{R}^{-} \cup \{+\infty\}$
    **Free-Energy** $E_w(S)$**:** Sum over (independently contributing) motifs in $S$
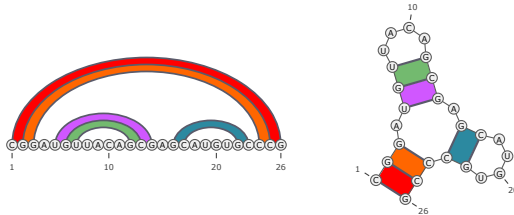
# Problem statement



- **RNA structure** $S$**:** Non-crossing base-pairs for positions in sequence $w$
- **Motifs:** Sequence/structure features (e.g. Base-pairs, Stacks, Loops. . . )
- **Energy model:**
  **Motif** $\to$ Free-energy contribution $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
  **Free-Energy** $E_w(S)$**:** Sum over (independently contributing) motifs in $S$
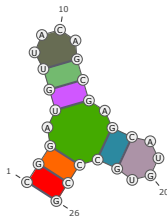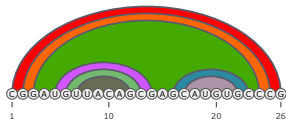
# Problem statement



- **RNA structure** $S$**:** Non-crossing base-pairs for positions in sequence $w$
- **Motifs:** Sequence/structure features (e.g. Base-pairs, Stacks, Loops...)
- **Energy model:**
  - **Motif** $\rightarrow$ Free-energy contribution $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
  - **Free-Energy** $E_w(S)$**:** Sum over (independently contributing) motifs in $S$
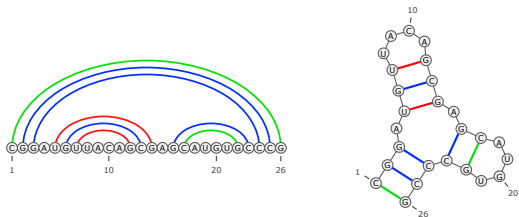
# Problem statement



- **RNA structure** $S$**:** Non-crossing base-pairs for positions in sequence *w*
- **Motifs:** Sequence/structure features (e.g. Base-pairs, Stacks, Loops...)
- **Energy model:**
    - **Motif** $\rightarrow$ Free-energy contribution $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
    - **Free-Energy** $E_w(S)$**:** Sum over (independently contributing) motifs in $S$
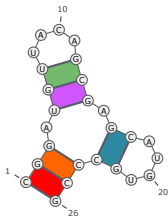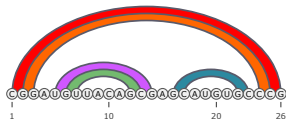
# Problem statement



- **RNA structure** $S$**:** Non-crossing base-pairs for positions in sequence $w$
- **Motifs:** Sequence/structure features (e.g. Base-pairs, Stacks, Loops...)
- **Energy model:**
  - **Motif** $\rightarrow$ Free-energy contribution $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
  - **Free-Energy** $E_w(S)$**:** Sum over (independently contributing) motifs in $S$

$$E_S = 2 \cdot \Delta \begin{pmatrix} \text{U} \\ \text{G} \end{pmatrix} + 4 \cdot \Delta \begin{pmatrix} \text{G} \\ \text{C} \end{pmatrix} + 2 \cdot \Delta \begin{pmatrix} \text{C} \\ \text{G} \end{pmatrix}$$
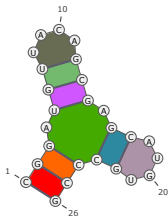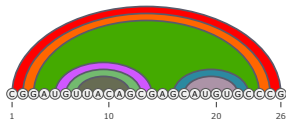
# Problem statement



- **RNA structure** $S$**:** Non-crossing base-pairs for positions in sequence $w$
- **Motifs:** Sequence/structure features (e.g. Base-pairs, Stacks, Loops...)
- **Energy model:**
  - **Motif** $\rightarrow$ Free-energy contribution $\Delta(\cdot) \in \mathbb{R}^{-} \cup \{+\infty\}$
  - **Free-Energy** $E_w(S)$**:** Sum over (independently contributing) motifs in $S$

$$E_S = \Delta \left( \begin{smallmatrix} C & G \\ G & C \end{smallmatrix} \right) + \Delta \left( \begin{smallmatrix} G & G \\ C & C \end{smallmatrix} \right) + \Delta \left( \begin{smallmatrix} U & G \\ G & C \end{smallmatrix} \right) + \Delta \left( \begin{smallmatrix} U & G \\ G & C \end{smallmatrix} \right) + \Delta \left( \begin{smallmatrix} U & G \\ G & C \end{smallmatrix} \right)$$
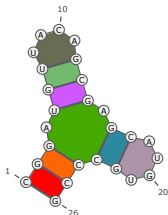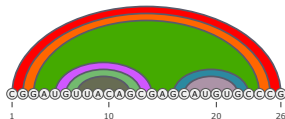
## Problem statement



- **RNA structure** $S$**:** Non-crossing base-pairs for positions in sequence $w$
- **Motifs:** Sequence/structure features (e.g. Base-pairs, Stacks, Loops...)
- **Energy model:**
  **Motif** $\rightarrow$ Free-energy contribution $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
  **Free-Energy** $E_w(S)$**:** Sum over (independently contributing) motifs in $S$

$$E_S = \Delta\left(\text{\scriptsize■}\right) + \Delta\left(\text{\scriptsize■}\right) + \Delta\left(\text{\scriptsize■}\right) + \Delta\left(\text{\scriptsize■}\right) + \Delta\left(\text{\scriptsize■}\right)$$

$$+ \Delta\left(\text{\scriptsize■}\right) + \Delta\left(\text{\scriptsize■}\right) + \Delta\left(\text{\scriptsize■}\right)$$

# Problem statement



- **RNA structure** $S$**:** Non-crossing base-pairs for positions in sequence $w$
- **Motifs:** Sequence/structure features (e.g. Base-pairs, Stacks, Loops...)
- **Energy model:**
  **Motif** $\rightarrow$ Free-energy contribution $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
  **Free-Energy** $E_w(S)$**:** Sum over (independently contributing) motifs in $S$

**Definition (**MFE-PREDICT$(E)$ **problem)**

**Input:** RNA sequence $w \in \{A, C, G, U\}^*$
**Output:** Secondary struct. $S^*$ with Minimal Free-Energy (MFE) $E_w(S^*)$
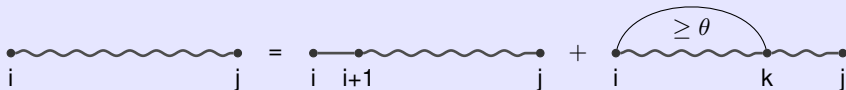
Problem solved **exactly** in $O(n^3)$ time.
**[Nussinov Jacobson, PNAS 1980] [Zuker Stiegler, NAR 1981]**. . . .

# Dynamic programming (DP) for RNA folding

**Theorem ([Nussinov and Jacobson(1980)])**

*Max #base-pairs/min energy structure computed in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$ time/memory*



$E_{i,k}$: Free-energy contribution of base-pair $(i,k)$.     $(-1/+\infty$ or $\Delta G(s_i \stackrel{?}{\equiv} s_k))$

$N_{i,j}$ : Max #base-pairs over interval $[i,j]$

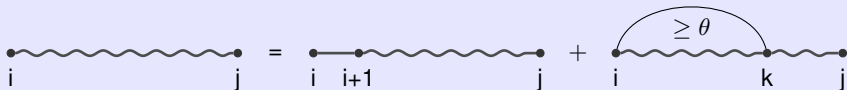$$N_{i,t} = 0, \quad \forall t \in [i, i+\theta]$$

$$N_{i,j} = \min \begin{cases} N_{i+1,j} & \{i \text{ unpaired}\} \\ \min_{k=i+\theta+1}^{j} E_{i,k} + N_{i+1,k-1} + N_{k+1,j} & \{i \text{ paired to } k\} \end{cases}$$

# Dynamic programming (DP) for RNA folding

**Theorem ([Nussinov and Jacobson(1980)])**

*Max #base-pairs/min energy structure computed in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$ time/memory*



$E_{i,k}$: Free-energy contribution of base-pair $(i,k)$.   $(-1/+\infty$ or $\Delta G(s_i \overset{?}{\equiv} s_k))$

$C_{i,j}$ : Number of secondary structures compatible with interval $[i,j]$
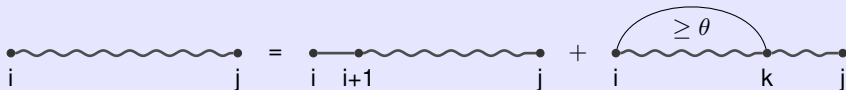
$$C_{i,t} = 1, \quad \forall t \in [i, i+\theta]$$

$$C_{i,j} = \sum \begin{cases} C_{i+1,j} & \{i \text{ unpaired}\} \\ \sum_{k=i+\theta+1}^{j} \mathbb{1}_{\text{comp.}(i,k)} \times C_{i+1,k-1} \times C_{k+1,j} & \{i \text{ paired to } k\} \end{cases}$$

# Dynamic programming (DP) for RNA folding

## Theorem ([Nussinov and Jacobson(1980)])

*Max #base-pairs/min energy structure computed in $\mathcal{O}(n^3)/\mathcal{O}(n^2)$ time/memory*



$E_{i,k}$: Free-energy contribution of base-pair $(i, k)$.     ($-1/+\infty$ or $\Delta G(s_i \overset{?}{\equiv} s_k)$)

$\mathcal{Z}_{i,j} = \sum_{\substack{S \text{ comp.} \\ \text{with } w_{[i,j]}}} e^{\frac{-E_W(S)}{RT}}$ = Partition function for compatible structs within $[i,j]$

$$\mathcal{Z}_{i,t} = 1, \quad \forall t \in [i, i+\theta]$$

$$\mathcal{Z}_{i,j} = \sum \begin{cases} \mathcal{Z}_{i+1,j} & \{i \text{ unpaired}\} \\ \sum_{k=i+\theta+1}^{j} e^{\frac{-E_{i,k}}{RT}} \times \mathcal{Z}_{i+1,k-1} \times \mathcal{Z}_{k+1,j} & \{i \text{ paired to } k\} \end{cases}$$

# Dynamic programming (DP) for RNA folding

**Many extensions:**

- Nearest-neighbor/Turner energy model **[Zuker1981]**
- Comparative folding **[Sankoff1985]**
- Equilibrium base-pairing probabilities **[McCaskill1990]**
- Moments of additive features **[Miklos2005,Ponty2011]**
- $\Delta$ kcal.mol$^{-1}$ suboptimal structures of MFE **[Wuchty1999]**
- Basic crossing structures **[Rivas1999]**...
- Exact sampling in Boltzmann distr. **[Ding2003,Ponty2008]**
- Moments of additive features **[Miklos2005,Ponty2011]**
- Maximum expected accuracy structure **[Do2006]**
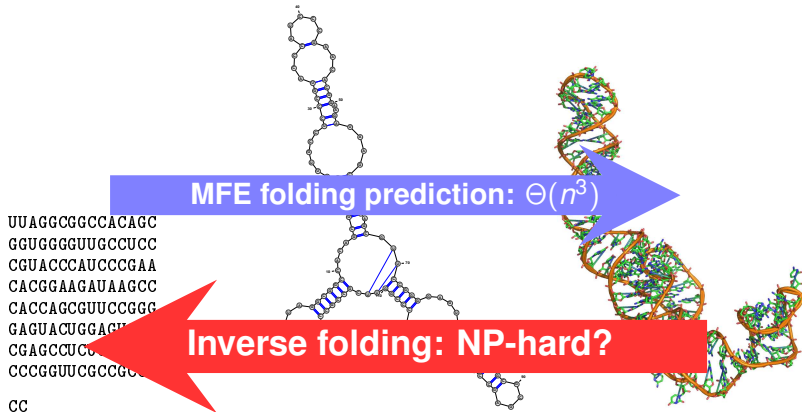- Distance-classified partitioning of Boltzmann ens. **[E.Freyhult2007a]**

**Made possible by:**

- **Completeness**/**Unambiguity** of decomposition
  $\exists$ energy-preserving bijection between **derivations of DP scheme** and **search space**
- Objective function **additive** with respect to DP scheme

# RNA inverse folding

**RNA** = Linear Polymer = Sequence in $\{A, C, G, U\}^\star$



**MFE folding prediction:** $\Theta(n^3)$

```
UUAGGCGGCCACAGC
GGUGGGGUUGCCUCC
CGUACCCAUCCCGAA
CACGGAAGAUAAGCC
CACCAGCGUUCCGGG
GAGUACUGGAGU
CGAGCCUC
CCCGGUUCGCCGCC

CC
```

**Inverse folding: NP-hard?**

Primary Structure          Secondary Structure          Structure Tertiaire

5s rRNA (PDBID: 1K73:B)

# RNA Inverse Folding

## Definition (INVERSE-FOLDING($E$) problem)
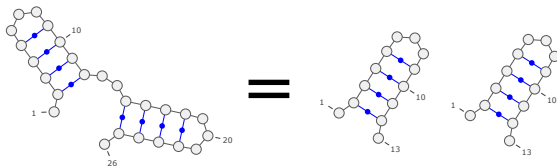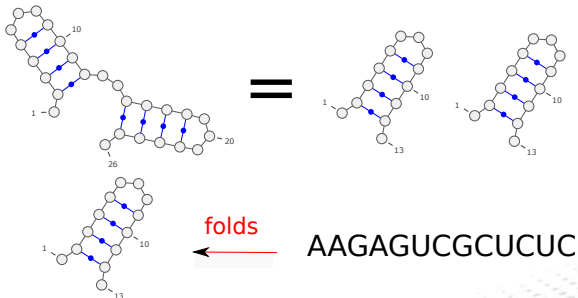
**Input:** Secondary structure $S$ + Energy distance $\Delta > 0$.
**Output:** RNA sequence $w \in \Sigma^\star$ such that:

$$\forall S' \in \mathcal{S}|w| \setminus \{S\} : E_{w,S'} \geq E w, S + \Delta$$

or $\varnothing$ if no such sequence exists.

**Difficult problem:** No **obvious** DP decomposition

- Existing algorithms: Heuristics or Exponential-time
- Complexity of problem unknown (despite **[Schnall Levin** *et al*, **ICML'08]**)
  **Reason:** Non locality, no theoretical frameworks, too many parameters...

# RNA Inverse Folding

**Example:**

# RNA Inverse Folding

**Definition (INVERSE-FOLDING($E$) problem)**
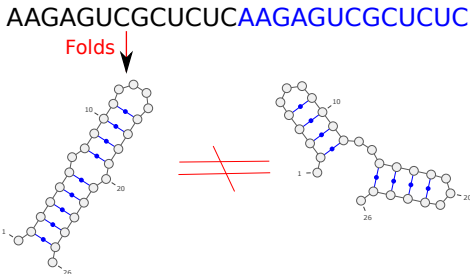
**Input:** Secondary structure $S$ + Energy distance $\Delta > 0$.
**Output:** RNA sequence $w \in \Sigma^\star$ such that:

$$\forall S' \in \mathcal{S}|w| \setminus \{S\} : \ E_{w,S'} \geq E_{w,S} + \Delta$$

or $\varnothing$ if no such sequence exists.

**Example:**

# RNA Inverse Folding

**Definition (INVERSE-FOLDING($E$) problem)**

**Input:** Secondary structure $S$ + Energy distance $\Delta > 0$.
**Output:** RNA sequence $w \in \Sigma^\star$ such that:

$$\forall S' \in \mathcal{S}|w| \setminus \{S\} : \ E_{w,S'} \geq E_{w,S} + \Delta$$

or $\varnothing$ if no such sequence exists.

**Example:**



folds

AAGAGUCGCUCUC

# RNA Inverse Folding

**Input:** Secondary structure $S$ + Energy distance $\Delta > 0$.
**Output:** RNA sequence $w \in \Sigma^\star$ such that:

$$\forall S' \in \mathcal{S}|w| \setminus \{S\} : \ E_{w,S'} \geq E w, S + \Delta$$

or $\varnothing$ if no such sequence exists.

**Example:**

AAGAGUCGCUCUCAAGAGUCGCUCUC

# Existing approaches for negative design

Based on local search...

- RNAInverse - TBI Vienna
- Info-RNA - Backofen@Freiburg
- RNA-SSD - Condon@UBC
- NUPack - Pierce@Caltech

...bio-inspired algorithms...

- RNAFBinv - Barash@Ben Gurion
- FRNAKenstein - Hein@Oxford
- AntaRNA - Backofen@Freiburg
- ERD - Ganjtabesh@Tehran

...exact approaches...
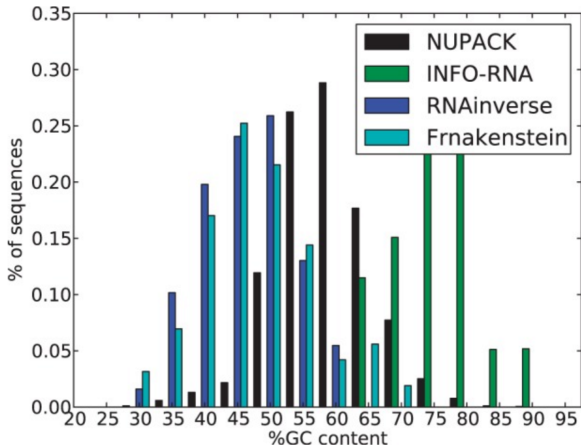
- RNAIFold - Clote@Boston College
- CO4 - Will@Leipzig

## Typical issues:

- Naive initialization strategies
- Poor coverage of sequence space:
  Local search remain *confined* near initial sequence
- GC-rich produced sequences

  ⇒ **Global sampling [Levin *et al*, NAR 12]**

# Existing approaches for negative design

Based on local search. . .
- RNAInverse - TBI Vienna
- Info-RNA - Backofen@Freiburg
- RNA-SSD - Condon@UBC
- NUPack - Pierce@Caltech

. . . bio-inspired algorithms. . .
- RNAFBinv - Barash@Ben Gurion
- FRNAKenstein - Hein@Oxford
- AntaRNA - Backofen@Freiburg
- ERD - Ganjtabesh@Tehran

. . . exact approaches. . .
- RNAIFold - Clote@Boston College
- CO4 - Will@Leipzig

## Typical issues:
- Naive initialization strategies
- Poor coverage of sequence space:
  Local search remain *confined* near initial sequence
- GC-rich produced sequences

  $\Rightarrow$ **Global sampling** [Levin *et al*, NAR 12]

## The case for a control of GC-content



High GC-content suspected to induce **kinetic traps**

# Global sampling [Levin *et al*, NAR 12]

## Target structure $S$

- **Boltzmann distribution** based on **affinity** towards $S$
- **Random generation** from **Boltzmann Distribution**
- **Fold** sampled sequences and **compare** to target

Boltzmann factor:
$$\mathcal{B}_w(S) := e^{\frac{-E_W(S)}{RT}}$$

Pseudo-Partition Function:
$$\mathcal{Z}(S) = \sum_{w \in \Sigma^*} \mathcal{B}_w(S)$$

Boltzmann probability:
$$p(s) := \frac{\mathcal{B}_w(S)}{\mathcal{Z}}$$

Explore sequence space
Structure fixed

# IncaRNAtion [Reinharz *et al*, Bioinformatics 2013]
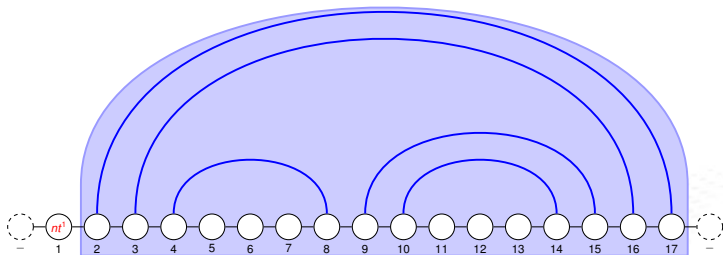


Explore sequence space
Structure fixed

$$\sum_{a'}$$

[Position $i$ unpaired]

# IncaRNAtion [Reinharz *et al*, Bioinformatics 2013]

Explore sequence space
Structure fixed



[Paired ends + Stacking pairs]

# IncaRNAtion [Reinharz *et al*, Bioinformatics 2013]

Explore sequence space
Structure fixed



[Position $i$ paired with $k$]

# IncaRNAtion [Reinharz *et al*, Bioinformatics 2013]



Explore sequence space
Structure fixed

$\sum\limits_{a'}$   [Position $i$ unpaired]

$\sum\limits_{a',b'}$   [Paired ends + Stacking pairs]

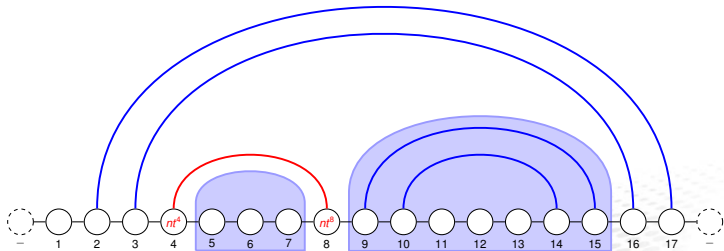$\sum\limits_{a',b'}$   [Position $i$ paired with $k$]

(Only 1 case applies)

# "Global" Stochastic Backtrack

Sequence:

# "Global" Stochastic Backtrack
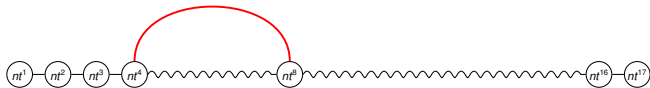
Sequence:    $nt^1$

## "Global" Stochastic Backtrack



Sequence:

$Z = \sum_{\sigma}$ weight of structure

# "Global" Stochastic Backtrack



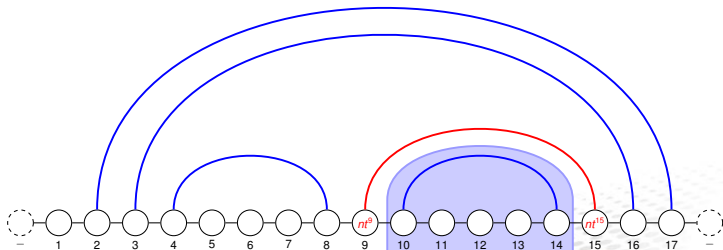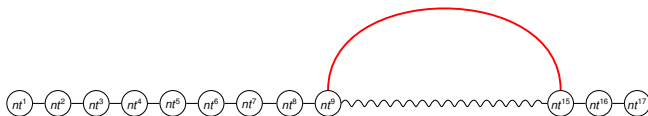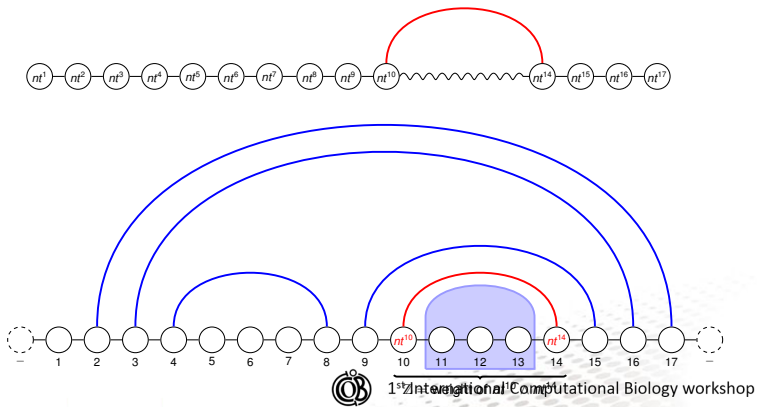Sequence:

# "Global" Stochastic Backtrack

Sequence:



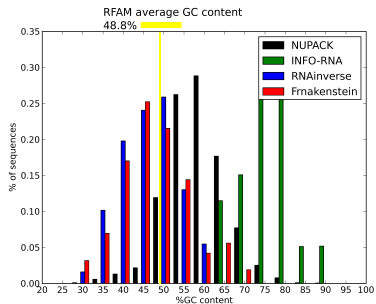$\mathbb{Z}$ = weight of $nt^4 \wedge nt^8$
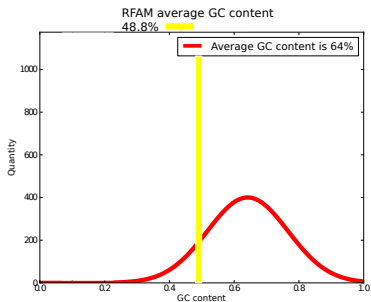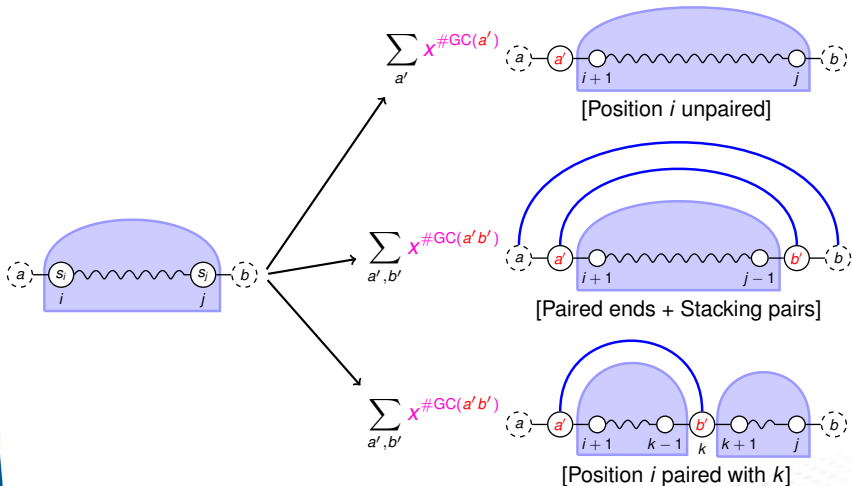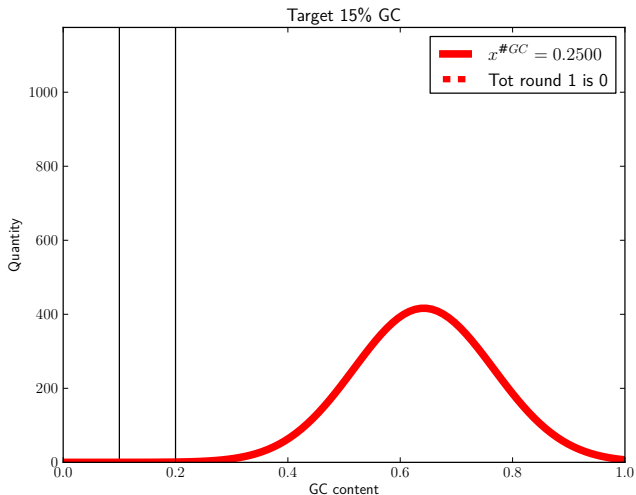
# "Global" Stochastic Backtrack

Sequence:

# "Global" Stochastic Backtrack



Sequence:

# GC-content bias

# Weighted DP Recursions



$$\sum_{a'} X^{\#\mathrm{GC}(a')}$$

[Position $i$ unpaired]

$$\sum_{a',b'} X^{\#\mathrm{GC}(a'b')}$$

[Paired ends + Stacking pairs]

$$\sum_{a',b'} X^{\#\mathrm{GC}(a'b')}$$
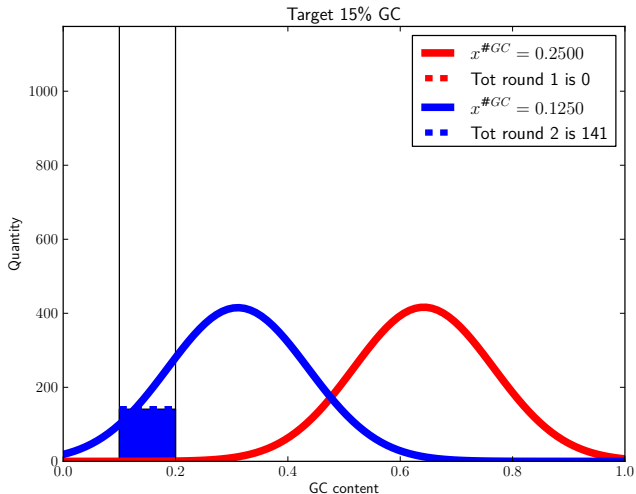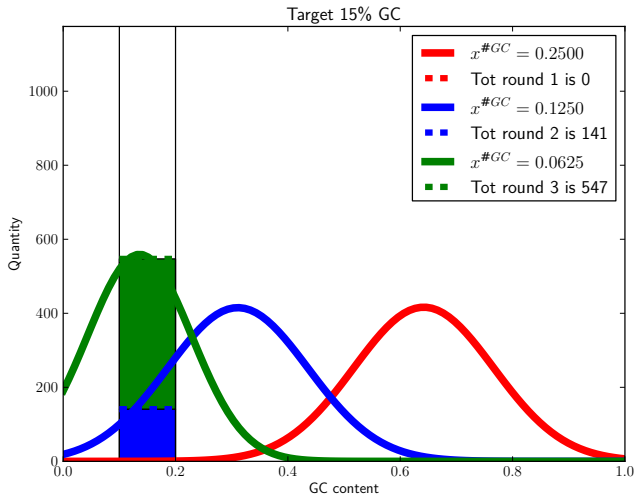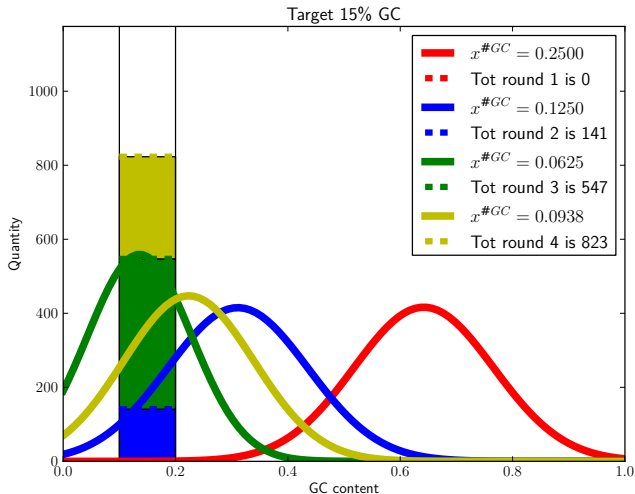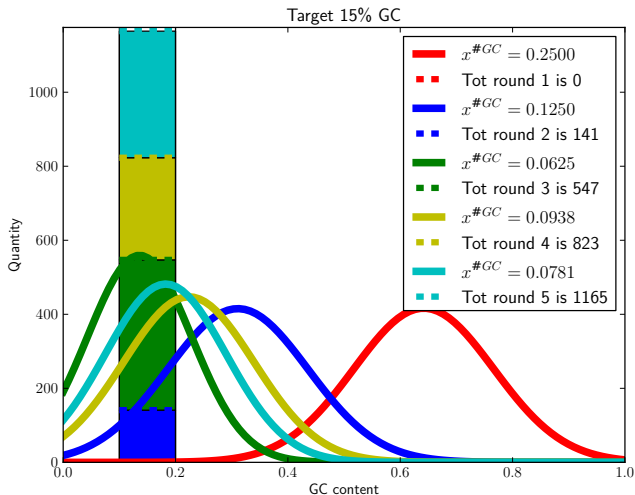
[Position $i$ paired with $k$]

# Incarnation NT distribution: Bissection scheme



[Waldispühl and Ponty, RECOMB, 2011]

# Incarnation NT distribution: Bissection scheme



[Waldispühl and Ponty, RECOMB, 2011]

## Incarnation NT distribution: Bissection scheme



[Waldispühl and Ponty, RECOMB, 2011]

# Incarnation NT distribution: Bissection scheme



Target 15% GC

Legend:
- $x^{\#GC} = 0.2500$
- Tot round 1 is 0
- $x^{\#GC} = 0.1250$
- Tot round 2 is 141
- $x^{\#GC} = 0.0625$
- Tot round 3 is 547
- $x^{\#GC} = 0.0938$
- Tot round 4 is 823

Axis labels: Quantity (y-axis), GC content (x-axis)
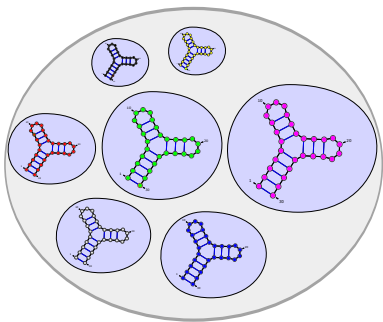
[Waldispühl and Ponty, RECOMB, 2011]

# Incarnation NT distribution: Bissection scheme



[Waldispühl and Ponty, RECOMB, 2011]
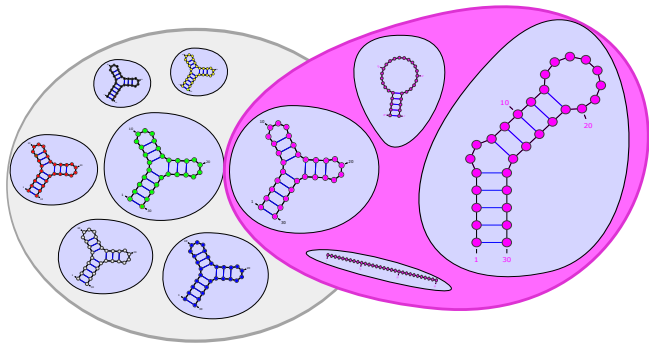
# Limits of the approach



**Heuristic:** Strong affinity is neither sufficient, nor necessary, **but** . . .

- Strong **empirical** correlation affinity/success of design [Levin et al, NAR 2012]
- **Linear** time-complexity [Reinharz Ponty Waldispühl, ISMB/ECCB'13]
- **Composition** control [Bodini Ponty, AofA'10] [Reinharz et al, ISMB/ECCB'13]
- **Complementary** with local search approaches [Reinharz et al, ISMB/ECCB'13]
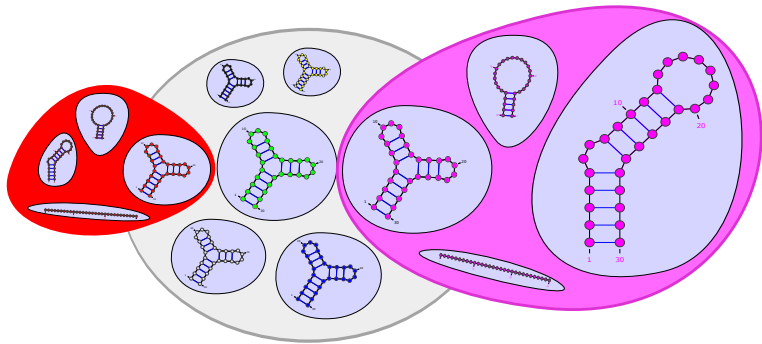
# Limits of the approach



**Heuristic:** Strong affinity is **neither sufficient,** nor necessary, **but** …

- Strong **empirical** correlation affinity/success of design **[Levin et al, NAR 2012]**
- **Linear** time-complexity **[Reinharz Ponty Waldispühl, ISMB/ECCB'13]**
- **Composition** control **[Bodini Ponty, AofA'10] [Reinharz et al, ISMB/ECCB'13]**
- **Complementary** with local search approaches **[Reinharz et al, ISMB/ECCB'13]**

# Limits of the approach



**Heuristic:** Strong affinity is **neither sufficient, nor necessary**, **but** . . .

- Strong **empirical** correlation affinity/success of design [Levin et al, NAR 2012]
- **Linear** time-complexity [Reinharz Ponty Waldispühl, ISMB/ECCB'13]
- **Composition** control [Bodini Ponty, AofA'10] [Reinharz et al, ISMB/ECCB'13]
- **Complementary** with local search approaches [Reinharz et al, ISMB/ECCB'13]
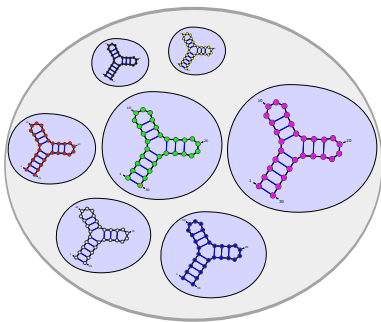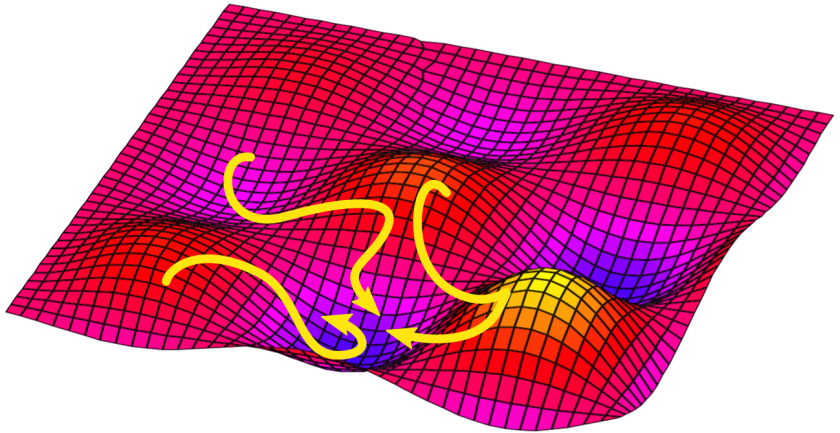
# Limits of the approach



**Heuristic:** Strong affinity is **neither sufficient, nor necessary**, **but** . . .
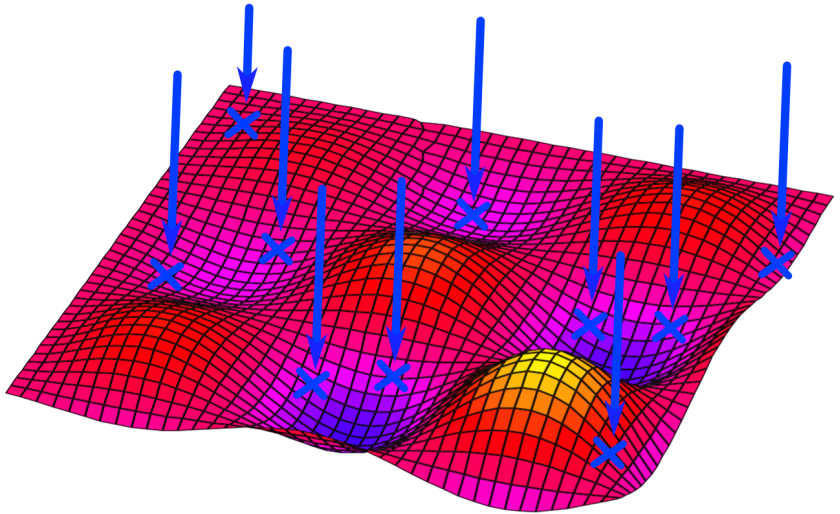
- Strong **empirical** correlation affinity/success of design **[Levin et al, NAR 2012]**
- **Linear** time-complexity **[Reinharz Ponty Waldispühl, ISMB/ECCB'13]**
- **Composition** control **[Bodini Ponty, AofA'10] [Reinharz et al, ISMB/ECCB'13]**
- **Complementary** with local search approaches **[Reinharz et al, ISMB/ECCB'13]**
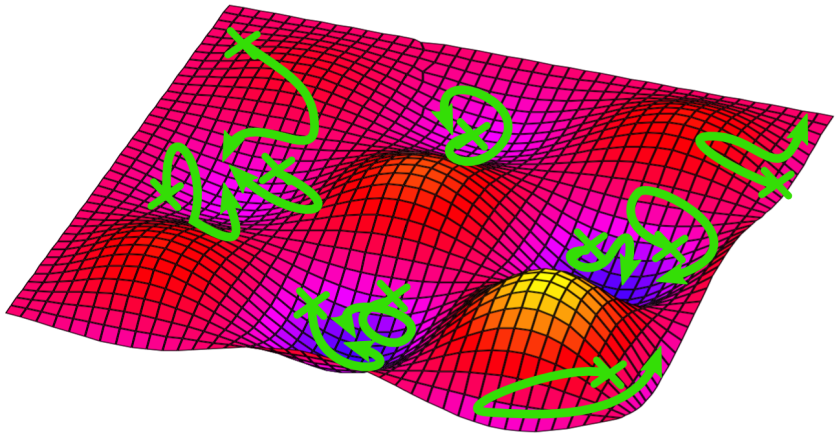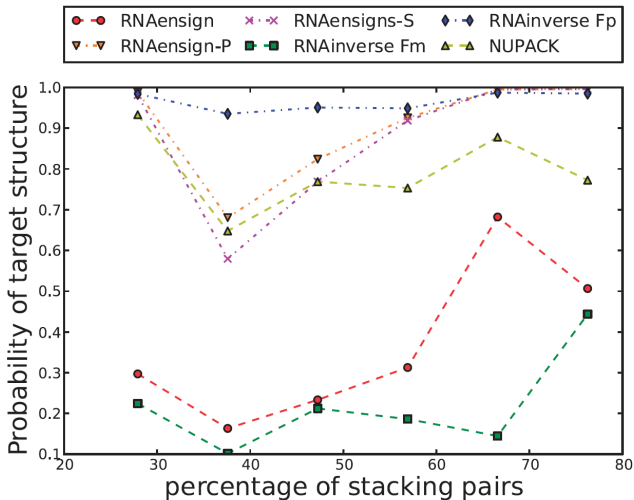
# Local vs Global vs "Glocal"

# Local vs Global vs "Glocal"

# Local vs Global vs "Glocal"

# The success of glocal strategies



Sampling + Optimize creates highly probable design sequences

# II. Constrained design

## Avoiding/forcing motifs

# Existing approaches for negative design

Based on local search...

- RNAInverse - TBI Vienna
- Info-RNA - Backofen@Freiburg
- RNA-SSD - Condon@UBC
- NUPack - Pierce@Caltech

...bio-inspired algorithms...

- RNAFBinv - Barash@Ben Gurion
- FRNAKenstein - Hein@Oxford
- AntaRNA - Backofen@Freiburg

...exact approaches...

- RNAIFold - Clote@Boston College
- CO4 - Will@Leipzig

Few algorithms support avoided/mandatory motifs...

...none guarantees *reasonable* runtime.

**Typical reasons:**

- Deep local minima (Rugged landscape)
- Mandatory motifs ⇒ Late deadends (Branch and Bound)
- Forbidden motifs ⇒ Search space disconnection (Local Search)

# Existing approaches for negative design

Based on local search. . .

- RNAInverse - TBI Vienna
- Info-RNA - Backofen@Freiburg
- RNA-SSD - Condon@UBC
- NUPack - Pierce@Caltech

. . . bio-inspired algorithms. . .

- RNAFBinv - Barash@Ben Gurion
- FRNAKenstein - Hein@Oxford
- AntaRNA - Backofen@Freiburg

. . . exact approaches. . .

- RNAIFold - Clote@Boston College
- CO4 - Will@Leipzig

Few algorithms support avoided/mandatory motifs. . .
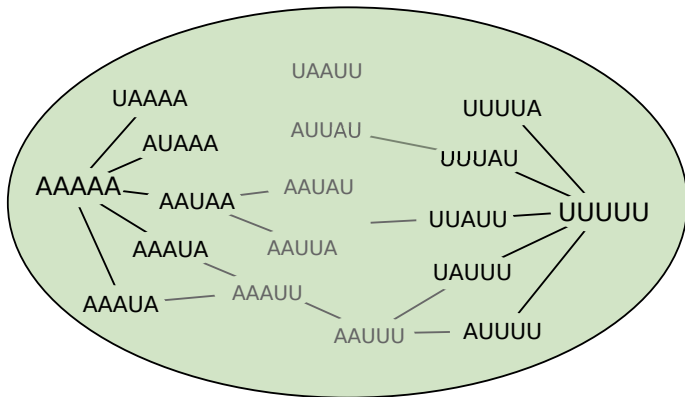
. . . none guarantees *reasonable* runtime.

## Typical reasons:

- Deep local minima (Rugged landscape)
- Mandatory motifs ⇒ Late deadends (Branch and Bound)
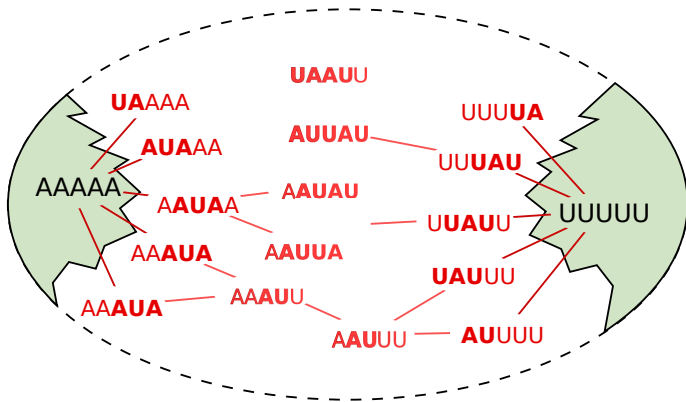- Forbidden motifs ⇒ Search space disconnection (Local Search)

# Problem with local approaches: An example

Simplified vocabulary $\{A, U\}$

# Problem with local approaches: An example

Simplified vocabulary $\{A, U\}$ **+ Forbidden motifs** $\mathcal{F} = \{AU, UA\}$



$\Rightarrow \mathcal{F}$ may **disconnect** search space (holds for **any** move set!)

# Idea

## Use formal language constructs to constrain global sampling

Forced motifs
Avoided motifs $\rightarrow$ Regular language $\mathcal{L}_C \in$ Reg

Structure compatibility
+ Positional constraints $\rightarrow$ **Weighted** Context-Free Lang $\mathcal{L}_S \in$ CFL
+ Energy Model

**Folklore theorem (constructive):** Reg $\cap$ (**W**)CFL $\subseteq$ (**W**)CFL

**Build weighted context-free grammar $\mathcal{G}$ for $\mathcal{L}_C \cap \mathcal{L}_S$**
**+ Random generation**

$\Rightarrow$ **Global sampling under constraints**

## Idea

**Use formal language constructs to constrain global sampling**

Forced motifs
Avoided motifs          $\rightarrow$ Regular language $\mathcal{L}_C \in$ Reg

Structure compatibility
+ Positional constraints   $\rightarrow$ **Weighted** Context-Free Lang $\mathcal{L}_S \in$ CFL
+ Energy Model

**Folklore theorem (constructive):** Reg $\cap$ (**W**)CFL $\subseteq$ (**W**)CFL

**Build weighted context-free grammar $\mathcal{G}$ for $\mathcal{L}_C \cap \mathcal{L}_S$**
**+ Random generation**

$\Rightarrow$ **Global sampling under constraints**

**Idea**


**Use formal language constructs to constrain global sampling**


Forced motifs
Avoided motifs $\quad \rightarrow$ Regular language $\mathcal{L}_C \in$ Reg


Structure compatibility
+ Positional constraints $\quad \rightarrow$ **Weighted** Context-Free Lang $\mathcal{L}_S \in$ CFL
+ Energy Model


**Folklore theorem (constructive):** Reg $\cap$ (**W**)CFL $\subseteq$ (**W**)CFL

Build weighted context-free grammar $\mathcal{G}$ for $\mathcal{L}_C \cap \mathcal{L}_S$
+ Random generation

$\Rightarrow$ Global sampling under constraints

# Idea

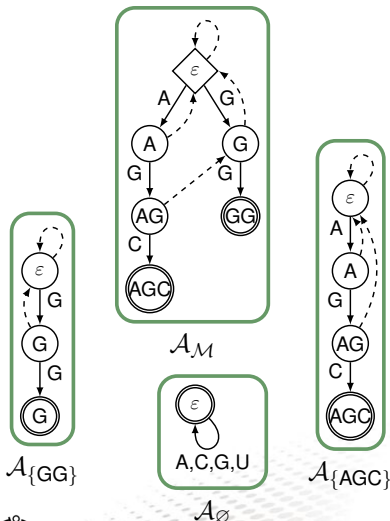## Use formal language constructs to constrain global sampling

Forced motifs
Avoided motifs → Regular language $\mathcal{L}_C \in$ Reg

Structure compatibility
+ Positional constraints → **Weighted** Context-Free Lang $\mathcal{L}_S \in$ CFL
+ Energy Model

**Folklore theorem (constructive):** Reg $\cap$ (**W**)CFL $\subseteq$ (**W**)CFL

**Build weighted context-free grammar $\mathcal{G}$ for $\mathcal{L}_C \cap \mathcal{L}_S$**
**+ Random generation**

$\Rightarrow$ **Global sampling under constraints**

## Building the Finite State Automaton

To force multiple words, **keep track** of generated words:
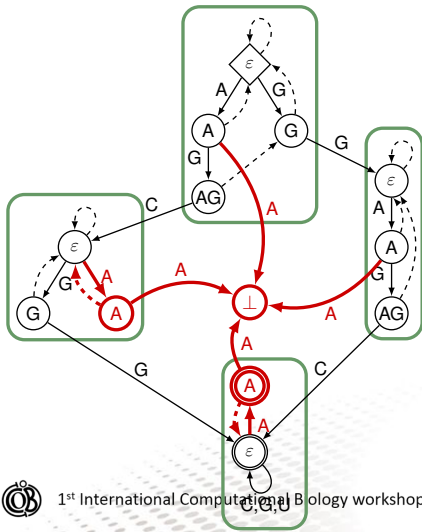
- Create disjunctive automata for each $\mathcal{M}' \subseteq \mathcal{M}$
  - **Reroute** accepting states
  - Accepting state = no **forced word** remaining ($\varepsilon$ in $\mathcal{A}_\varnothing$)
  - Forbidden words can be added to sub-automata

**Example:** $\mathcal{M} = \{AGC, GG\}$



$\mathcal{A}_\mathcal{M}$

$\mathcal{A}_{\{GG\}}$

$\mathcal{A}_\varnothing$

$\mathcal{A}_{\{AGC\}}$

### #States:

$$O\left(2^{|\mathcal{M}|} \cdot \left(\sum_i |f_i| + \sum_j |m_j|\right)\right)$$

## Building the Finite State Automaton

To force multiple words, **keep track** of generated words:

- Create disjunctive automata for each $\mathcal{M}' \subseteq \mathcal{M}$
- **Reroute** accepting states
- Accepting state = no **forced word** remaining ($\varepsilon$ in $\mathcal{A}_\varnothing$)
- Forbidden words can be added to sub-automata

**Example:** $\mathcal{M} = \{AGC, GG\}$



**#States:**

$$O\left(2^{|\mathcal{M}|} \cdot \left(\sum_i |f_i| + \sum_j |m_j|\right)\right)$$

## Building the Finite State Automaton

To force multiple words, **keep track** of generated words:

- Create disjunctive automata for each $\mathcal{M}' \subseteq \mathcal{M}$
- **Reroute** accepting states
- Accepting state = no **forced word** remaining ($\varepsilon$ in $\mathcal{A}_\varnothing$)
- Forbidden words can be added to sub-automata

**#States:**

$$O\left(2^{|\mathcal{M}|} \cdot \left(\sum_i |f_i| + \sum_j |m_j|\right)\right)$$

**Example:**

$\mathcal{M} = \{\text{AGC}, \text{GG}\}; \mathcal{F} = \{\text{AA}\}$

# Building the grammar

**Input:** Secondary Structure $S$ + Positional constraints

A **Create Parse Tree** for secondary structure
B **Translate Parse Tree** into single-word grammar
C **Expand** grammar to instantiate compatible base/base-pairs
D **Restrict** to bases/base-pairs allowed at each position

## Building the grammar

**Input:** Secondary Structure $S$ + Positional constraints

**A Create Parse Tree** for secondary structure

B Translate Parse Tree into single-word grammar

C Expand grammar to instantiate compatible base/base-pairs

D Restrict to bases/base-pairs allowed at each position

## Building the grammar

**Input:** Secondary Structure $S$ + Positional constraints

 **A Create Parse Tree** for secondary structure

 **B Translate Parse Tree** into single-word grammar

 C Expand grammar to instantiate compatible base/base-pairs

 D Restrict to bases/base-pairs allowed at each position

$$S_1 \to .S_2 \qquad S_2 \to (S_3) \qquad S_3 \to (S_4)S_8 \qquad S_4 \to (S_5)$$
$$S_5 \to . \qquad S_8 \to (S_9) \qquad S_9 \to .S_{10} \qquad S_{10} \to .$$

## Building the grammar

**Input:** Secondary Structure $S$ + Positional constraints
- **A** **Create Parse Tree** for secondary structure
- **B** **Translate Parse Tree** into single-word grammar
- **C** **Expand** grammar to instantiate compatible base/base-pairs
- **D** **Restrict** to bases/base-pairs allowed at each position

$$V_1 \rightarrow A\,V_2 \mid C\,V_2 \mid G\,V_2 \mid U\,V_2$$

$$V_2 \rightarrow A\,V_3\,U \mid C\,V_3\,G \mid G\,V_3\,C \mid G\,V_3\,U \mid U\,V_3\,A \mid U\,V_3\,G$$

$$V_3 \rightarrow A\,V_4\,U\,V_8 \mid C\,V_4\,G\,V_8 \mid G\,V_4\,C\,V_8 \mid G\,V_4\,U\,V_8 \mid U\,V_4\,A\,V_8 \mid U\,V_4\,G\,V_8$$

$$V_4 \rightarrow A\,V_5\,U \mid C\,V_5\,G \mid G\,V_5\,C \mid G\,V_5\,U \mid U\,V_5\,A \mid U\,V_5\,G$$

$$V_5 \rightarrow A \mid C \mid G \mid U$$

$$V_8 \rightarrow A\,V_9\,U \mid C\,V_9\,G \mid G\,V_9\,C \mid G\,V_9\,U \mid U\,V_9\,A \mid U\,V_9\,G$$

$$V_9 \rightarrow A\,V_{10} \mid C\,V_{10} \mid G\,V_{10} \mid U\,V_{10}$$

$$V_{10} \rightarrow A \mid C \mid G \mid U$$

## Building the grammar

**Input:** Secondary Structure $S$ **+ Positional constraints**

- **A Create Parse Tree** for secondary structure
- **B Translate Parse Tree** into single-word grammar
- **C Expand** grammar to instantiate compatible base/base-pairs
- **D Restrict** to bases/base-pairs allowed at each position

$$V_1 \rightarrow A V_2 \mid C V_2 \mid G V_2 \mid U V_2$$
$$V_2 \rightarrow A V_3 U \mid \cancel{C V_3 G} \mid \cancel{G V_3 C} \mid \cancel{G V_3 U} \mid U V_3 A \mid \cancel{U V_3 G}$$
$$V_3 \rightarrow A V_4 U V_8 \mid \cancel{C V_4 G V_8} \mid \cancel{G V_4 C V_8} \mid \cancel{G V_4 U V_8} \mid U V_4 A V_8 \mid \cancel{U V_4 G V_8}$$
$$V_4 \rightarrow A V_5 U \mid \cancel{C V_5 G} \mid \cancel{G V_5 C} \mid \cancel{G V_5 U} \mid U V_5 A \mid \cancel{U V_5 G}$$
$$V_5 \rightarrow A \mid C \mid G \mid U$$
$$V_8 \rightarrow A V_9 U \mid \cancel{C V_9 G} \mid \cancel{G V_9 C} \mid \cancel{G V_9 U} \mid U V_9 A \mid \cancel{U V_9 G}$$
$$V_9 \rightarrow A V_{10} \mid C V_{10} \mid G V_{10} \mid U V_{10}$$
$$V_{10} \rightarrow A \mid C \mid G \mid U$$

## Random generation

Combine CFG and aut. $\rightarrow$ CFG (Multiplying #Rules by $|Q|^3$)

**GenRGenS** [Ponty Termier Denise, Bioinformatics 2006]**:**
- Precomputes #words for each non-terminal
- Random Generation w.r.t. **weighted distribution**

**Energy models:**
- **Uniform distribution**
- **Nussinov energy model**
- **Stacking-pairs model (Turner 2004)**
  Based on refined, yet similar, grammar

**Overall complexity:** $|S| \cdot 2^{3|\mathcal{M}|} \cdot \left( \sum_i |f_i| + \sum_j |m_j| \right)^3$
- **Linear** on $|S|$
- **Exponential** on $|\mathcal{M}|$, but **NP-Hard** problem

# III. Positive design for multiple structures

# Motivation: Kinetics and riboswitches

# Design objectives

**Positive structural design**
Optimize affinity of designed sequences towards target structure
Or simply ensure their compatibility with one or **several structures**
**Examples:** Most stable sequence for given fold...

**Negative structural design**
Limit affinity of designed sequences towards **alternative structures**
**Examples:** Lowest free-energy, High Boltzmann probability/Low entropy...

**Additional constraints:**

- **Forbid** motif list to appear **anywhere** in design
- **Force** motif list to appear **each at least once**
- **Limit** available alternatives at certain positions
- **Control** overall composition (GC-content)

# Design objectives

## Positive structural design
Optimize affinity of designed sequences towards target structure
Or simply **ensure their compatibility with one or several structures**
**Examples:** Most stable sequence for given fold. . .

## Negative structural design
Limit affinity of designed sequences towards **alternative structures**
**Examples:** Lowest free-energy, High Boltzmann probability/Low entropy. . .

## Additional constraints:

- **Forbid** motif list to appear **anywhere** in design
- **Force** motif list to appear **each at least once**
- **Limit** available alternatives at certain positions
- **Control** overall composition (GC-content)

# Counting compatible RNAs: Watson-Crick + Single structure

A —— U

G —— C

**Compatible Base Pairs** = Only **Watson-Crick** base pairs



a b c d e f g h i j k l m n o p q r s t u v

# Counting compatible RNAs: Watson-Crick + Single structure

A — U

G — C

**Compatible Base Pairs** = Only **Watson-Crick** base pairs

**Compatible Sequence**

G A G C U C A G C G C G A A C G C U C U C A

# Counting compatible RNAs: Watson-Crick + Single structure

A — U

G — C

**Compatible Base Pairs** = Only **Watson-Crick** base pairs

**Incompatible Sequence**

A A A G A C U G G A C U U G G C C U U C

**Question:** How many **Compatible** sequences?

# Counting compatible RNAs: Watson-Crick + Single structure



**Compatible Base Pairs** = Only **Watson-Crick** base pairs



**Question:** How many **Compatible** sequences?

# Counting compatible RNAs: Watson-Crick + Single structure

$$A — U$$

$$G — C$$

**Compatible Base Pairs** = Only **Watson-Crick** base pairs



**Question:** How many **Compatible** sequences?

# Counting compatible RNAs: Watson-Crick + Single structure



$$A \;—\; U$$

$$G \;—\; C$$

**Compatible Base Pairs** = Only **Watson-Crick** base pairs



a b c d e f g h i j k l m n o p q r s t u v

## Question: How many **Compatible** sequences?

**Answer:** $4^{\#BPs} \times 4^{\#Unpaired} \rightarrow 268\,435\,456$

# Counting compatible RNAs: Watson-Crick + Two structures

Compatible Base Pairs = Only Watson-Crick base pairs

Question: How many Compatible sequences?

Answer: $\neq \emptyset$! (both base-pairs and dependency graphs bipartite)

# Counting compatible RNAs: Watson-Crick + Two structures



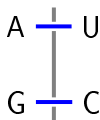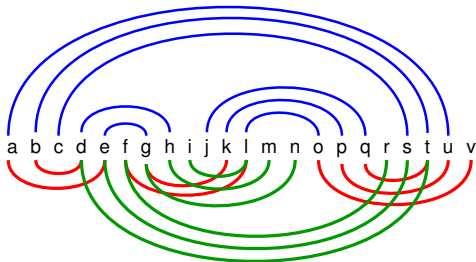**Compatible Base Pairs** = Only **Watson-Crick** base pairs



**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (both base-pairs and dependency graphs **bipartite**)

# Counting compatible RNAs: Watson-Crick + Two structures



**Compatible Base Pairs** = Only **Watson-Crick** base pairs



**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (both base-pairs and dependency graphs **bipartite**)

# Counting compatible RNAs: Watson-Crick + Two structures

A — U

G — C

**Compatible Base Pairs** = Only **Watson-Crick** base pairs

**Dependency graph:**
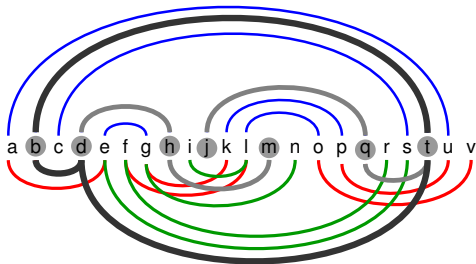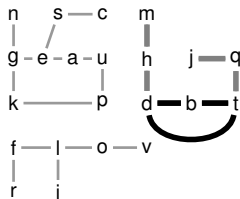Cycles + Paths

**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (both base-pairs and dependency graphs **bipartite**)

# Counting compatible RNAs: Watson-Crick + Two structures



$$A \underline{\quad} U$$

$$G \underline{\quad} C$$

**Compatible Base Pairs** = Only **Watson-Crick** base pairs



**Dependency graph:** Cycles + Paths
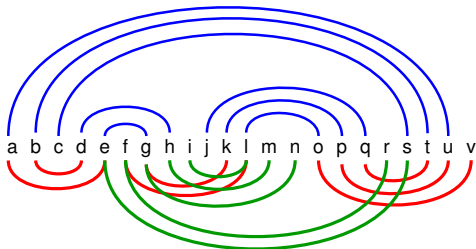
**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (both base-pairs and dependency graphs **bipartite**)

$$4^{\#CCs} \rightarrow 65\,536$$

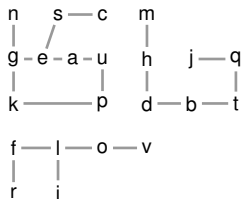# Counting compatible RNAs: Watson-Crick $+ > 2$ structs

**Compatible Base Pairs** = Only **Watson-Crick** base pairs
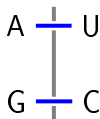
**Dependency graph:**
Cycles, Paths, Trees. . .

**Question:** How many **Compatible** sequences?

**Answer:** Non-bipartite $\rightarrow \varnothing$; Bipartite $\rightarrow 4^{\#CCs} = 64$

# Counting compatible RNAs: Watson-Crick + > 2 structs



A — U

G — C

**Compatible Base Pairs** = Only **Watson-Crick** base pairs

**Dependency graph:**

Cycles, Paths, Trees. . .

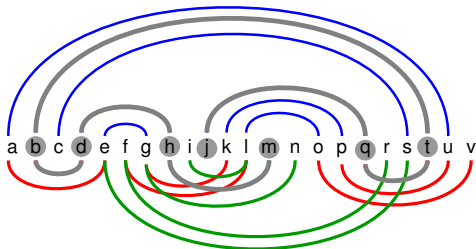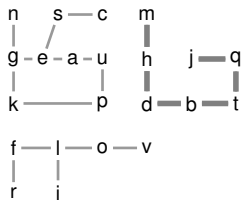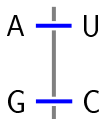**Question:** How many **Compatible** sequences?

**Answer:** Non-bipartite → ∅; Bipartite → $4^{\#CCs} = 64$

# Counting compatible RNAs: Watson-Crick + > 2 structs



A — U

G — C

**Compatible Base Pairs** = Only **Watson-Crick** base pairs

**Dependency graph:**
Cycles, Paths, Trees...
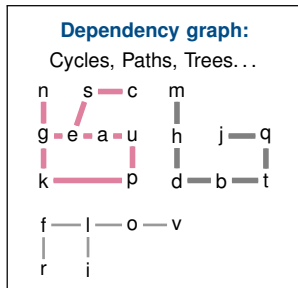
**Question:** How many **Compatible** sequences?

**Answer:** Non-bipartite → ∅; Bipartite → $4^{\#CCs} = 64$

# Counting compatible RNAs: Watson-Crick + > 2 structs



A —— U

G —— C

**Compatible Base Pairs** = Only **Watson-Crick** base pairs



**Dependency graph:**
Cycles, Paths, Trees. . .
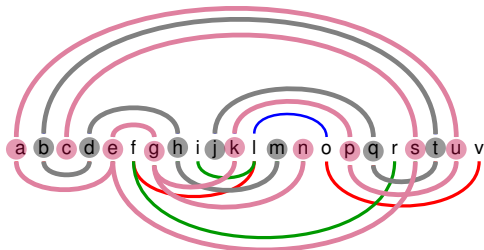
**Question:** How many **Compatible** sequences?

**Answer:** Non-bipartite → ∅; Bipartite → $4^{\#CCs}$ = 64

# Counting compatible RNAs: Watson-Crick + > 2 structs



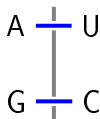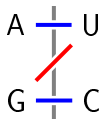**Compatible Base Pairs** = Only **Watson-Crick** base pairs



**Dependency graph:**
Cycles, Paths, Trees...

**Question:** How many **Compatible** sequences?

**Answer:** Non-bipartite → ∅; Bipartite → $4^{\#CCs} = 64$

# Counting compatible RNAs: Watson-Crick + > 2 structs

$$A \text{ —— } U$$
$$G \text{ —— } C$$

**Compatible Base Pairs** = Only **Watson-Crick** base pairs



**Dependency graph:**
Cycles, Paths, Trees...

**Question:** How many **Compatible** sequences?

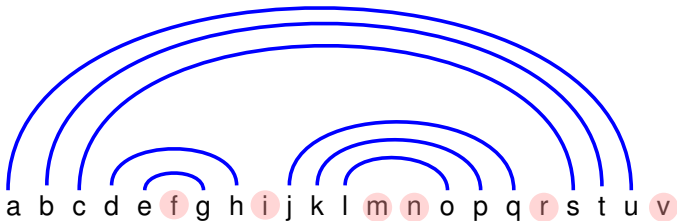**Answer:** Non-bipartite $\to \varnothing$; Bipartite $\to 4^{\#\text{CCs}} = 64$

# Counting compatible RNAs: WC/Wobble + Single struct.

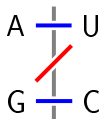

**Compatible Base Pairs** = Include **Wobble** base pairs



a b c d e f g h i j k l m n o p q r s t u v

**Question:** How many **Compatible** sequences?

**Answer:** $4^{\#\text{Unpaired}} \times 6^{\#\text{BPs}} \rightarrow 6\,879\,707\,136$

# Counting compatible RNAs: WC/Wobble + Single struct.



**Compatible Base Pairs** = Include **Wobble** base pairs
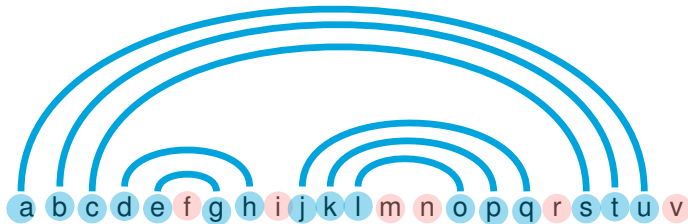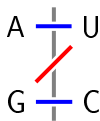


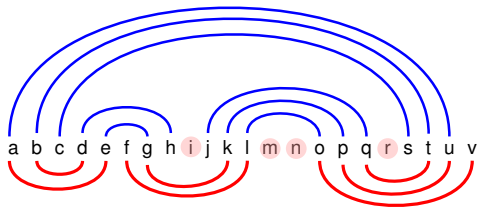**Question:** How many **Compatible** sequences?

**Answer:** $4^{\#\text{Unpaired}} \times 6^{\#\text{BPs}} \rightarrow 6\,879\,707\,136$

# Counting compatible RNAs: WC/Wobble + Two structures



**Compatible Base Pairs** = Include **Wobble** base pairs



**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (base-pairs and dependency graphs **always bipartite**)

# Counting compatible RNAs: WC/Wobble + Two structures



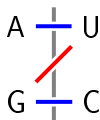**Compatible Base Pairs** = Include **Wobble** base pairs



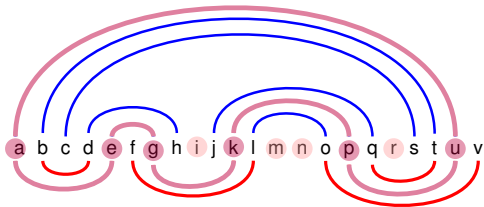**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (base-pairs and dependency graphs **always bipartite**)

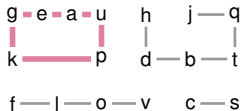# Counting compatible RNAs: WC/Wobble + Two structures



**Compatible Base Pairs** = Include **Wobble** base pairs



**Dependency graph:**
Cycles + Paths

**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (base-pairs and dependency graphs **always bipartite**)

# Counting compatible RNAs: WC/Wobble + Two structures



**Compatible Base Pairs** = Include **Wobble** base pairs



**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (base-pairs and dependency graphs **always bipartite**)

# Counting compatible RNAs: WC/Wobble + Two structures

**Compatible Base Pairs** = Include **Wobble** base pairs
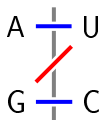
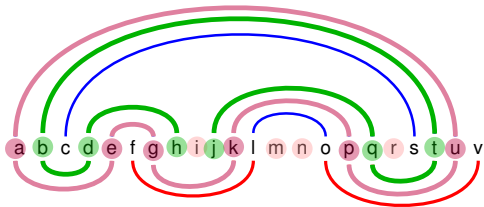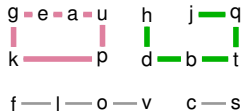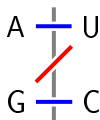**Dependency graph:**
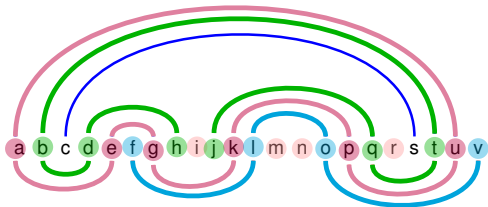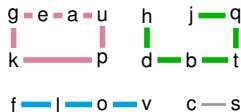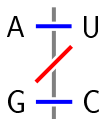Cycles + Paths

**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (base-pairs and dependency graphs **always bipartite**)

$$\#\text{Designs}(G) = \prod_{c \in CC(G)} \#\text{Designs}(cc)$$

# Counting compatible designs for paths and cycles

## Theorem (#Compatible designs for paths and cycles)

*The numbers of compatible designs for paths and cycles of length n are:*

$$p(n) = 2\,\mathcal{F}_{n+2} \qquad and \qquad c(n) = 2\,\mathcal{F}_n + 4\,\mathcal{F}_{n-1}$$

*where $\mathcal{F}_n$: $n^{th}$ **Fibonacci number**, $\mathcal{F}_0 = 0$, $\mathcal{F}_1 = 1$ and $\mathcal{F}_n = \mathcal{F}_{n-1} + \mathcal{F}_{n-2}$.*

**For paths:** A simple DFA generates compatible sequences



**Remark:** A ↔ C/G ↔ U symmetry

# Counting compatible designs for paths and cycles

## Theorem (#Compatible designs for paths and cycles)

*The numbers of compatible designs for* <span style="color:red">*paths*</span> *and* <span style="color:blue">*cycles*</span> *of length n are:*

$$\textcolor{red}{p(n)} = 2\,\mathcal{F}_{n+2} \qquad and \qquad \textcolor{blue}{c(n)} = 2\,\mathcal{F}_n + 4\,\mathcal{F}_{n-1}$$

*where* $\mathcal{F}_n$: $n^{th}$ ***Fibonacci number***, $\mathcal{F}_0 = 0$, $\mathcal{F}_1 = 1$ *and* $\mathcal{F}_n = \mathcal{F}_{n-1} + \mathcal{F}_{n-2}$.

**For paths:** A simple DFA generates compatible sequences



**Remark:** A $\leftrightarrow$ C/G $\leftrightarrow$ U symmetry
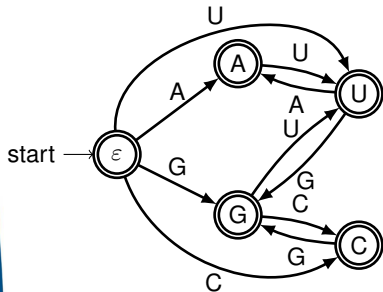
# Counting compatible designs for paths and cycles

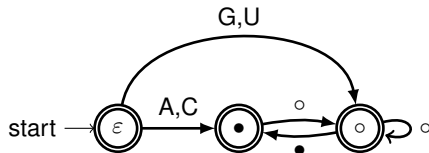**Theorem (#Compatible designs for paths and cycles)**

*The numbers of compatible designs for paths and cycles of length n are:*

$$p(n) = 2\,\mathcal{F}_{n+2} \qquad \text{and} \qquad c(n) = 2\,\mathcal{F}_n + 4\,\mathcal{F}_{n-1}$$

*where $\mathcal{F}_n$: $n^{th}$ **Fibonacci number**, $\mathcal{F}_0 = 0$, $\mathcal{F}_1 = 1$ and $\mathcal{F}_n = \mathcal{F}_{n-1} + \mathcal{F}_{n-2}$.*

**For paths:** A simple DFA generates compatible sequences



**Remark:** A $\leftrightarrow$ C/G $\leftrightarrow$ U symmetry

$$m_\bullet(n) = m_\circ(n-1)$$

# Counting compatible designs for paths and cycles
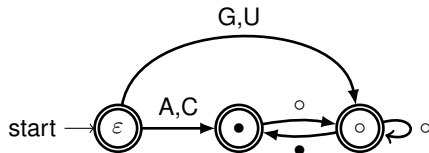
## Theorem (#Compatible designs for paths and cycles)

*The numbers of compatible designs for* <span style="color:red">paths</span> *and* <span style="color:blue">cycles</span> *of length n are:*

$$\textcolor{red}{p(n)} = 2\,\mathcal{F}_{n+2} \qquad \text{and} \qquad \textcolor{blue}{c(n)} = 2\,\mathcal{F}_{n} + 4\,\mathcal{F}_{n-1}$$

*where* $\mathcal{F}_n$: $n^{th}$ **Fibonacci number**, $\mathcal{F}_0 = 0$, $\mathcal{F}_1 = 1$ and $\mathcal{F}_n = \mathcal{F}_{n-1} + \mathcal{F}_{n-2}$.

**For paths:** A simple DFA generates compatible sequences



**Remark:** A $\leftrightarrow$ C / G $\leftrightarrow$ U symmetry

$$m_{\bullet}(n) = m_{\circ}(n-1)$$
$$m_{\circ}(n) = m_{\circ}(n-1) + m_{\bullet}(n-1)$$
$$= m_{\circ}(n-1) + m_{\circ}(n-2)$$
$$= \mathcal{F}(n+2)$$

# Counting compatible designs for paths and cycles

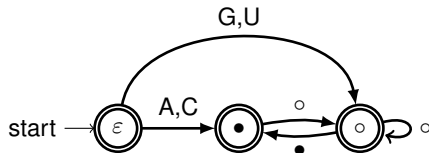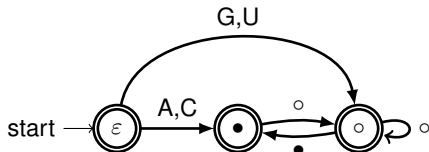## Theorem (#Compatible designs for paths and cycles)

*The numbers of compatible designs for paths and cycles of length n are:*

$$p(n) = 2\,\mathcal{F}_{n+2} \qquad \text{and} \qquad c(n) = 2\,\mathcal{F}_n + 4\,\mathcal{F}_{n-1}$$

*where $\mathcal{F}_n$: $n^{th}$ **Fibonacci number**, $\mathcal{F}_0 = 0$, $\mathcal{F}_1 = 1$ and $\mathcal{F}_n = \mathcal{F}_{n-1} + \mathcal{F}_{n-2}$.*

**For paths:** A simple DFA generates compatible sequences



**Remark:** A $\leftrightarrow$ C/G $\leftrightarrow$ U symmetry

$$m_\bullet(n) = m_\circ(n - 1)$$
$$m_\circ(n) = m_\circ(n - 1) + m_\bullet(n - 1)$$
$$= m_\circ(n - 1) + m_\circ(n - 2)$$
$$= \mathcal{F}(n + 2)$$

(Since $m_\circ(0) = 1$ and $m_\circ(1) = 2$)
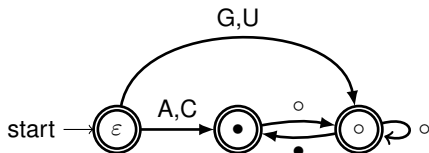
# Counting compatible designs for paths and cycles

## Theorem (#Compatible designs for paths and cycles)

*The numbers of compatible designs for paths and cycles of length n are:*

$$p(n) = 2\,\mathcal{F}_{n+2} \qquad and \qquad c(n) = 2\,\mathcal{F}_n + 4\,\mathcal{F}_{n-1}$$

*where $\mathcal{F}_n$: $n^{th}$ **Fibonacci number**, $\mathcal{F}_0 = 0$, $\mathcal{F}_1 = 1$ and $\mathcal{F}_n = \mathcal{F}_{n-1} + \mathcal{F}_{n-2}$.*

**For paths:** A simple DFA generates compatible sequences



**Remark:** A $\leftrightarrow$ C/G $\leftrightarrow$ U symmetry

$$m_\bullet(n) = m_\circ(n-1)$$
$$m_\circ(n) = m_\circ(n-1) + m_\bullet(n-1)$$
$$= m_\circ(n-1) + m_\circ(n-2)$$
$$= \mathcal{F}(n+2)$$

(Since $m_\circ(0) = 1$ and $m_\circ(1) = 2$)

$$p(n) := m_\varepsilon(n) = 2\,m_\bullet(n-1) + 2\,m_\circ(n-1) = 2(\mathcal{F}(n) + \mathcal{F}(n+1)) =, \mathcal{F}(n+2)$$

# Counting compatible designs for paths and cycles

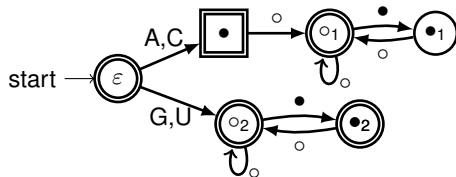## Theorem (#Compatible designs for paths and cycles)

*The numbers of compatible designs for* <span style="color:red">*paths*</span> *and* <span style="color:blue">*cycles*</span> *of length n are:*

$$p(n) = 2\,\mathcal{F}_{n+2} \qquad \text{and} \qquad c(n) = 2\,\mathcal{F}_n + 4\,\mathcal{F}_{n-1}$$

*where* $\mathcal{F}_n$: $n^{th}$ **Fibonacci number**, $\mathcal{F}_0 = 0$, $\mathcal{F}_1 = 1$ *and* $\mathcal{F}_n = \mathcal{F}_{n-1} + \mathcal{F}_{n-2}$.

**For cycle:** A *barely more involved* DFA generates compatible sequences



**Remark:** A $\leftrightarrow$ C/G $\leftrightarrow$ U symmetry

$$m_{\circ_2}(n) = \mathcal{F}(n+2)$$
$$m_{\circ_1}(n) = \mathcal{F}(n+1)$$

(Since $m_{\circ_1}(0) = 1$ and $m_{\circ_1}(1) = 1$)

$$c(n) := m_\varepsilon(n) = 2\,m_{\circ_1}(n-2) + 2\,m_{\circ_2}(n-1)$$
$$= 2(\mathcal{F}(n-1) + \mathcal{F}(n+1)) = 2\,\mathcal{F}(n) + 4\,\mathcal{F}(n-1)$$

# Counting compatible designs for paths and cycles

## Theorem (#Compatible designs for paths and cycles)

*The numbers of compatible designs for <span style="color:red">paths</span> and <span style="color:blue">cycles</span> of length n are:*

$$p(n) = 2\,\mathcal{F}_{n+2} \qquad and \qquad c(n) = 2\,\mathcal{F}_n + 4\,\mathcal{F}_{n-1}$$

*where $\mathcal{F}_n$: $n^{th}$ **Fibonacci number**, $\mathcal{F}_0 = 0$, $\mathcal{F}_1 = 1$ and $\mathcal{F}_n = \mathcal{F}_{n-1} + \mathcal{F}_{n-2}$.*

## Theorem (#Compatible designs for general 2-structures graphs)

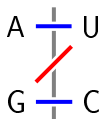*G: dependency graph associated with 2 RNA structures (max deg=2).*
*The number $\#\mathrm{Designs}(G)$ of compatible designs for G is given by*

$$\#\mathrm{Designs}(G) = \prod_{p \in \mathcal{P}(G)} 2\,\mathcal{F}_{|p|+2} \times \prod_{c \in \mathcal{C}(G)} \left(2\,\mathcal{F}_{|c|} + 4\,\mathcal{F}_{|c|-1}\right)$$

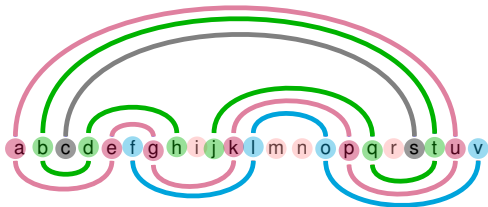*where G decomposes into <span style="color:red">paths</span> $\mathcal{P}(G)$ and <span style="color:blue">cycles</span> $\mathcal{C}(G)$.*

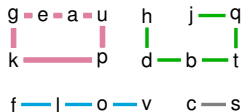# Counting compatible sequences: WC/Wobble + Two structures



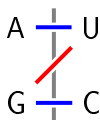**Compatible Base Pairs** = Include **Wobble** base pairs



**Question:** How many **Compatible** sequences?

**Answer:** $\neq \varnothing$! (base-pairs and dependency graphs **always bipartite**)
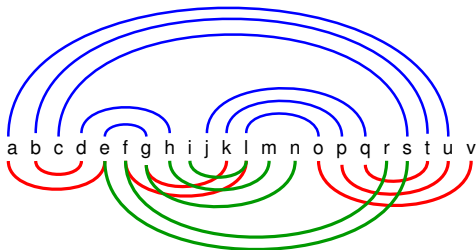
$$\#\text{Designs}(G) = \prod_{c \in CC(G)} \#\text{Designs}(cc) = 2\,322\,432$$

# Counting compatible sequences: Watson-Crick + > 2 structures



Compatible Base Pairs = Include Wobble base pairs



Dependency graph:
Cycles, Paths, Trees...

Question: How many Compatible sequences?

Answer: Non-bipartite → ∅; Bipartite →

1st International Computational Biology workshop

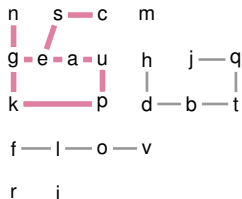# Counting compatible sequences: Watson-Crick + > 2 structures



**Compatible Base Pairs** = Include **Wobble** base pairs
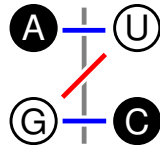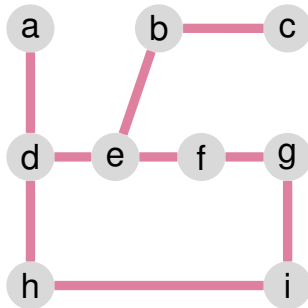


**Dependency graph:**

Cycles, Paths, Trees...
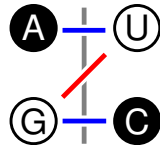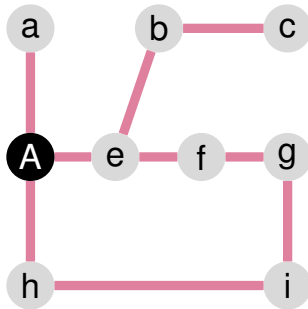
**Question:** How many **Compatible** sequences?

**Answer:** Non-bipartite $\rightarrow \varnothing$; Bipartite $\rightarrow \displaystyle\prod_{cc \in CC(G)} 2 \times \#IS(cc)$

**Bijection between Independent Sets and Valid Designs**

# Bijection between Independent Sets and Valid Designs

# Bijection between Independent Sets and Valid Designs

# Bijection between Independent Sets and Valid Designs
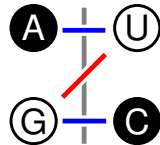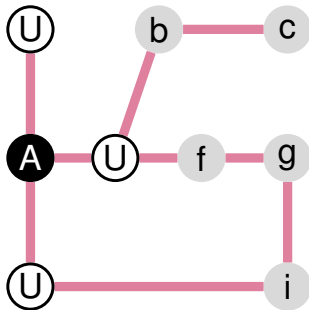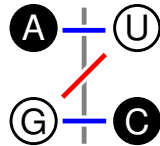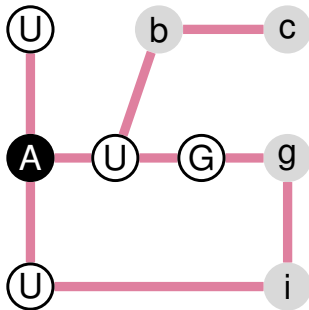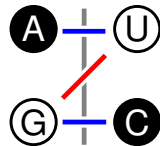
# Bijection between Independent Sets and Valid Designs

## Bijection between Independent Sets and Valid Designs



**Remark:** No adjacent **black letters** in compatible designs

Up to trivial symmetry* (*e.g.* top-left position $\in \{G, A\}$):

$$\text{Designs}^\star(cc) \subseteq \text{IndependentSets}(cc)$$

## Bijection between Independent Sets and Valid Designs



**Remark:** No adjacent **black letters** in compatible designs

Up to trivial symmetry* (*e.g.* top-left position $\in \{G, A\}$):

$$\text{Designs}^\star(cc) \subseteq \text{IndependentSets}(cc)$$

Also, IS (black) + $\nwarrow$ vert. $\in \{G, A\} \Rightarrow$ **Unique** compatible design
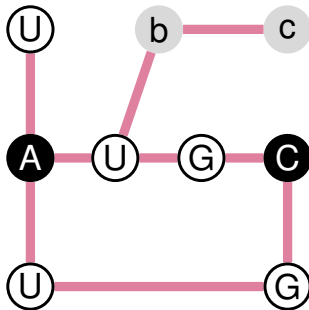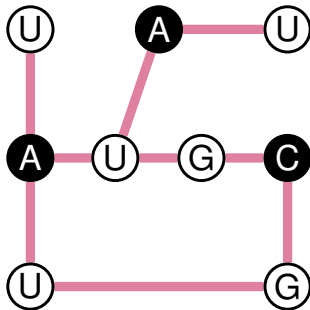
# Bijection between Independent Sets and Valid Designs



**Remark:** No adjacent **black letters** in compatible designs

Up to trivial symmetry$^\star$ (*e.g.* top-left position $\in \{G, A\}$):

$$\text{Designs}^\star(cc) \subseteq \text{IndependentSets}(cc)$$

Also, IS (black) + $\nwarrow$ vert. $\in \{G, A\} \Rightarrow$ **Unique** compatible design

$\Rightarrow$ Bijection between Designs$^\star$(cc) and IndependentSets(cc).

# Valid designs and independent sets

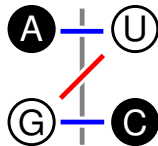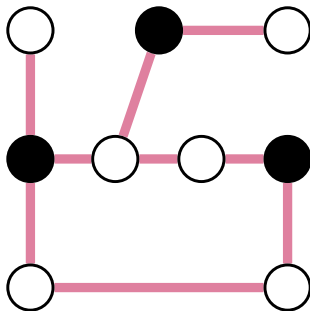**Theorem (#Valid design for bipartite connected dependency graphs)**

*Let G be a **bipartite connected** dependency graph, one has:*

$$\#\text{Designs}(G) = 2 \times \#\text{Designs}^{\star}(G) = 2 \times \#\text{IS}(G)$$

For a **bipartite** dependency graph $G$ we get:

$$\#Designs(G) = \prod_{cc \in CC(G)} 2 \times \#IS(cc) = 2^{|CC(G)|} \times \#IS(G)$$

**But** $\#IS(G)$ is #P-hard on bipartite graphs **[Bubbley&Dyer'01]**
(+ Any $G$ is a dependency graph)

Algorithm $\mathcal{A} \in P$ for $\#Designs(G) \rightarrow$ Algorithm $\mathcal{A}' \in P$ for $\#BIS\dots$

**Theorem**

*#Designs is #P-hard.*

No polynomial algorithm for #Designs(G) unless $\#P = FP$ ($\Rightarrow P = NP$)

# Valid designs and independent sets

**Theorem (#Valid design for bipartite connected dependency graphs)**

*Let G be a **bipartite connected** dependency graph, one has:*

$$\#\text{Designs}(G) = 2 \times \#\text{Designs}^\star(G) = 2 \times \#\text{IS}(G)$$

For a **bipartite** dependency graph *G* we get:

$$\#Designs(G) = \prod_{cc \in CC(G)} 2 \times \#IS(cc) = 2^{|CC(G)|} \times \#IS(G)$$

**But** $\#IS(G)$ is #P-hard on bipartite graphs **[Bubbley&Dyer'01]**
(+ Any $G$ is a dependency graph)

Algorithm $\mathcal{A} \in P$ for $\#Designs(G) \rightarrow$ Algorithm $\mathcal{A}' \in P$ for $\#BIS\ldots$

**Theorem**

*#Designs is #P-hard.*

No polynomial algorithm for #Designs(G) unless $\#P = FP\ (\Rightarrow P = NP)$

# Valid designs and independent sets

**Theorem (#Valid design for bipartite connected dependency graphs)**

*Let $G$ be a **bipartite connected** dependency graph, one has:*

$$\#\text{Designs}(G) = 2 \times \#\text{Designs}^\star(G) = 2 \times \#\text{IS}(G)$$

For a **bipartite** dependency graph $G$ we get:

$$\#Designs(G) = \prod_{cc \in CC(G)} 2 \times \#IS(cc) = 2^{|CC(G)|} \times \#IS(G)$$

**But** $\#IS(G)$ is #P-hard on bipartite graphs **[Bubbley&Dyer'01]**
(+ Any $G$ is a dependency graph)

Algorithm $\mathcal{A} \in P$ for $\#Designs(G) \rightarrow$ Algorithm $\mathcal{A}' \in P$ for $\#BIS$...

**Theorem**

*#Designs is #P-hard.*

No polynomial algorithm for #Designs(G) unless $\#P = FP \; (\Rightarrow P = NP)$

## Valid designs and independent sets

**Theorem (#Valid design for bipartite connected dependency graphs)**

*Let G be a **bipartite connected** dependency graph, one has:*

$$\#\text{Designs(G)} = 2 \times \#\text{Designs}^{\star}(G) = 2 \times \#\text{IS(G)}$$

For a **bipartite** dependency graph $G$ we get:

$$\#Designs(G) = \prod_{cc \in CC(G)} 2 \times \#IS(cc) = 2^{|CC(G)|} \times \#IS(G)$$

**But** $\#IS(G)$ is #P-hard on bipartite graphs **[Bubbley&Dyer'01]**
(+ Any $G$ is a dependency graph)

Algorithm $\mathcal{A} \in P$ for $\#Designs(G) \rightarrow$ Algorithm $\mathcal{A}' \in P$ for $\#BIS\ldots$

**Theorem**

*#Designs is #P-hard.*

No polynomial algorithm for #Designs(G) unless $\#P = FP (\Rightarrow P = NP)$

# Consequences

**Corollary (#Approximability for $\leq 5$ structures)** [Weitz'06]

For any $G$ built from $\leq 5$ pseudoknotted structures, #Design($G$) can be approximated within **any ratio** in **polynomial time** (PTAS)

**Corollary (#BIS hardness for $> 5$ struct.)** [Cai, Galanis, Goldberg, Jerrum, McQuillan'16]

Beyond 5 **pseudoknotted** structures, approximating #Design becomes **as hard as** approximating #BIS without any constraint.

**Why pseudoknotted?** Because any bipartite graph of max degree $\Delta$ can be **decomposed** into $\Delta$ matchings **in polynomial time** (Vizing's theorem).

Lastly, connection between **counting** and **sampling** [Jerrum, Valiant, Vazirani'86].

**Conjecture (#BIS hardness of sampling)**

Generating comp. sequences **(almost) uniformly** for general input is **#BIS-hard**.

# Consequences

**Corollary (#Approximability for $\leq 5$ structures)** [Weitz'06]

For any $G$ built from $\leq 5$ pseudoknotted structures, #Design($G$) can be approximated within **any ratio** in **polynomial time** (PTAS)

**Corollary (#BIS hardness for $> 5$ struct.)** [Cai, Galanis, Goldberg, Jerrum, McQuillan'16]

Beyond 5 **pseudoknotted** structures, approximating #Design becomes **as hard as** approximating #BIS without any constraint.

**Why pseudoknotted?** Because any bipartite graph of max degree $\Delta$ can be **decomposed** into $\Delta$ matchings **in polynomial time** (Vizing's theorem).
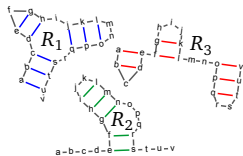
Lastly, connection between **counting** and **sampling** [Jerrum, Valiant, Vazirani'86].

**Conjecture (#BIS hardness of sampling)**

Generating comp. sequences **(almost) uniformly** for general input is **#BIS-hard**.

## Consequences

**Corollary (#Approximability for $\leq 5$ structures)** [Weitz'06]

For any $G$ built from $\leq 5$ pseudoknotted structures, #Design($G$) can be approximated within **any ratio** in **polynomial time** (PTAS)

**Corollary (#BIS hardness for $> 5$ struct.)** [Cai, Galanis, Goldberg, Jerrum, McQuillan'16]

Beyond 5 **pseudoknotted** structures, approximating #Design becomes **as hard as** approximating #BIS without any constraint.

**Why pseudoknotted?** Because any bipartite graph of max degree $\Delta$ can be **decomposed** into $\Delta$ matchings **in polynomial time** (Vizing's theorem).

Lastly, connection between **counting** and **sampling** [Jerrum, Valiant, Vazirani'86].

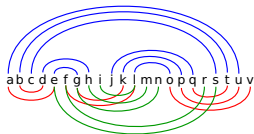**Conjecture (#BIS hardness of sampling)**

Generating comp. sequences **(almost) uniformly** for general input is **#BIS-hard**.
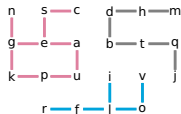
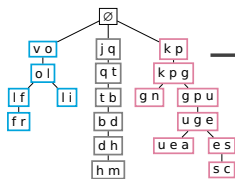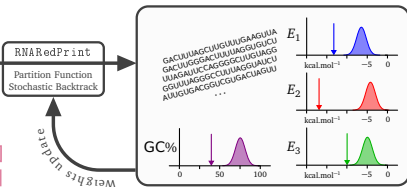# Perspectives: FPT and Boltzmann sampling algorithms
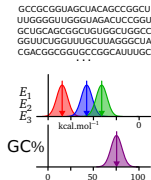


i) Input Structures

ii) Merged Base-Pairs

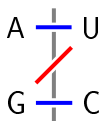iii) Compatibility Graph

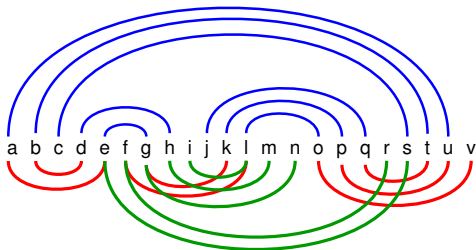iv) Tree Decomposition

v) Weight Optimization (Adaptive Sampling)

vi) Final Designs

- **FPT algorithm** for counting based on **tree decomposition**
- Multidimensional Boltzmann sampling to control energies, GC...

# Counting compatible sequences: Watson-Crick + $> 2$ structures



**Compatible Base Pairs** = Include **Wobble** base pairs



**Dependency graph:**

Cycles, Paths, Trees...

n    s — c    m

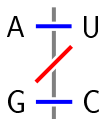g — e — a — u    h      j — q
|           |                  |
k ————— p    d — b — t

f — l — o — v

r    i

**Question:** How many **Compatible** sequences?

**Answer:** Bipartite $\rightarrow$
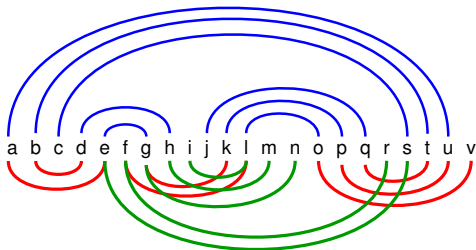
# Counting compatible sequences: Watson-Crick + > 2 structures



**Compatible Base Pairs** = Include **Wobble** base pairs



Dependency graph:
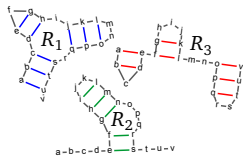Cycles, Paths, Trees...

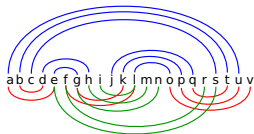**Question:** How many **Compatible** sequences?

**Answer:** Bipartite $\rightarrow \prod_{cc \in CC(G)} 2 \times \#IS(cc) = 496\,672$
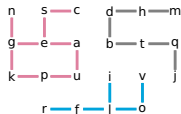
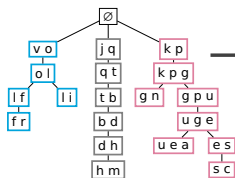# Perspectives: FPT and Boltzmann sampling algorithms
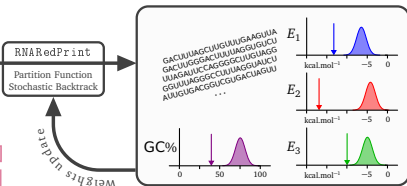


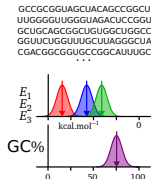i) Input Structures      ii) Merged Base-Pairs      iii) Compatibility Graph

iv) Tree Decomposition      v) Weight Optimization (Adaptive Sampling)      vi) Final Designs

- **FPT algorithm** for counting based on **tree decomposition**
- **Multidimensional Boltzmann sampling** to control energies, GC…

## Conclusions

- **RNA** is **cool!**

- **RNA design** is one of the current challenge of RNA bioinformatics with far-reaching consequences for drug design, synthetic biology...

- Practical use-cases require **expressive and modular constraints**

- Future methods: **kinetics**, **interactions**, **multiple structures**, **pseudoknots**...

- **RNA inverse folding** is the combinatorial core of design. It remains **largely unsolved**, and opens **new lines of research** in Comp. Sci.

# Collaborators

## University McGill 🇨🇦
Vladimir Reinharz
Jérôme Waldispühl

## MIT 🇺🇸
Bonnie Berger
Srinivas Devadas
Alex Levin
Mieszko Lis
Charles O'Donnell

## LRI – Univ. Paris Sud 🇫🇷
Alain Denise
Vincent Le Gallic

## Wuhan University 🇨🇳
Yi Zhang
Yu Zhou

## LIGM – Marne la Vallée 🇫🇷
Stéphane Vialette

## LIX – Ecole Polytechnique 🇫🇷
Alice Héliou
Mireille Regnier

## Simon Fraser University 🇨🇦
Jozef Hales
Jan Manuch (UBC)
Ladislav Stacho

Cédric Chauve
Julien Courtiel

## TBI Vienna 🇦🇹
Ronnie Lorenz
Andrea Tanzer

**Thanks!**

**Poster submission** & **Registration** open soon...
(+ ISCB travel fellowships for students)

# References I

R. Nussinov and A.B. Jacobson.

Fast algorithm for predicting the secondary structure of single-stranded RNA.
*Proc Natl Acad Sci U S A*, 77:6903–13, 1980.