

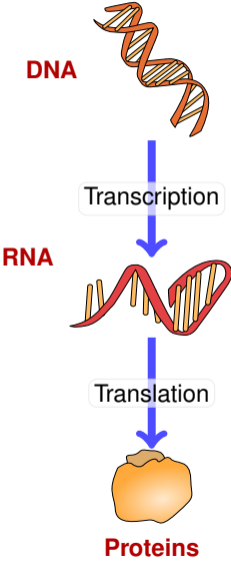


RNA design: Hard(?) but exemplarily combinatorial

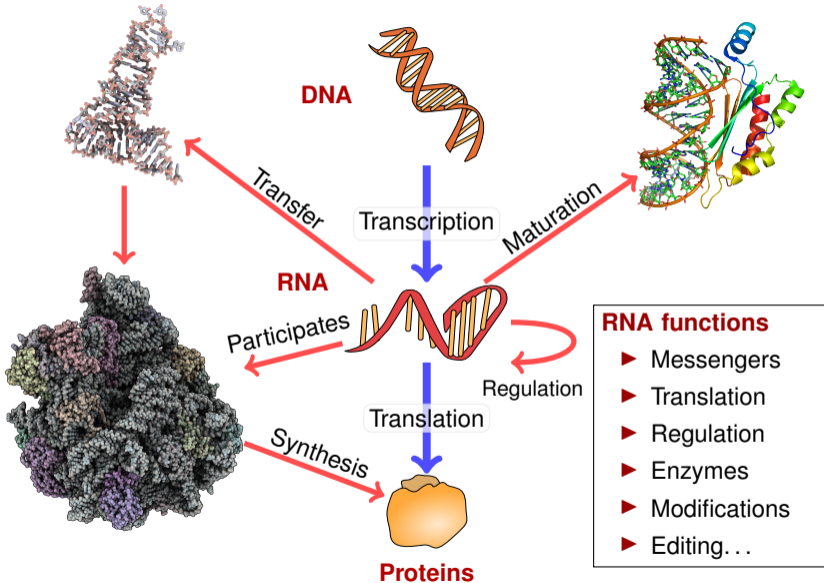
Yann Ponty

LIX, CNRS/Ecole Polytechnique

Fundamental dogma of molecular biology



Fundamental dogma of molecular biology (v2.0)



RiboNucleic Acids (RNA) in Human biology/health: Friend **and** Foe!

RiboNucleic Acids (RNAs)



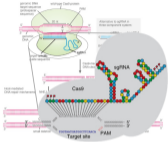
Encodes proteins
mRNA Vaccines
COVID-19, Malaria (Zika, CMV, Cancers?)

RiboNucleic Acids (RNA) in Human biology/health: Friend **and** Foe!

Targeting system for DNA Editing

CRISPR therapies

Sickle-cell anemia, β -thalassamia, Leber congenital amaurosis (LCA), cancers...



Hendel et al, 2015; Agrotis & Ketteler, 2015

RiboNucleic Acids (RNAs)



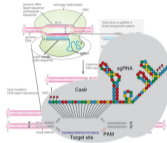
Encodes proteins
mRNA Vaccines
COVID-19, Malaria (Zika, CMV, Cancers?)

RiboNucleic Acids (RNA) in Human biology/health: Friend **and** Foe!

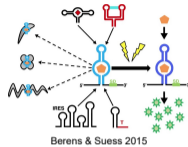
Targeting system for DNA Editing

CRISPR therapies

Sickle-cell anemia, β -thalassaemia, Leber congenital amaurosis (LCA), cancers...



Hendel et al, 2015; Agrotis & Ketteler, 2015



Berens & Suess 2015

Sensor of metabolites

Riboswitches

RiboNucleic Acids (RNAs)



Encodes proteins

mRNA Vaccines

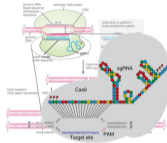
COVID-19, Malaria (Zika, CMV, Cancers?)

RiboNucleic Acids (RNA) in Human biology/health: Friend **and** Foe!

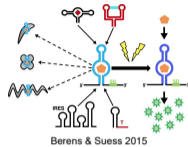
Targeting system for DNA Editing

CRISPR therapies

Sickle-cell anemia, β -thalassaemia, Leber congenital amaurosis (LCA), cancers...



Hendel et al, 2015; Agrotis & Ketteler, 2015

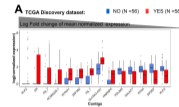


Sensor of metabolites
Riboswitches

Quantitative expression

Transcriptomic signatures

Cancer diagnosis/prognosis/relapse...



[NGuyen et al, 2021]

RiboNucleic Acids (RNAs)



Encodes proteins

mRNA Vaccines

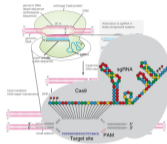
COVID-19, Malaria (Zika, CMV, Cancers?)

RiboNucleic Acids (RNA) in Human biology/health: Friend **and** Foe!

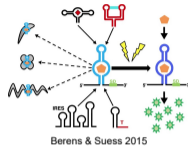
Targeting system for DNA Editing

CRISPR therapies

Sickle-cell anemia, β -thalassaemia, Leber congenital amaurosis (LCA), cancers...



Hendel et al, 2015; Agrotis & Ketteler, 2015



Berens & Suess 2015

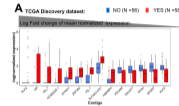
Sensor of metabolites

Riboswitches

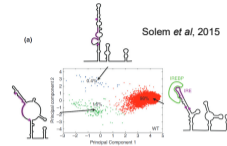
Quantitative expression

Transcriptomic signatures

Cancer diagnosis/prognosis/relapse...



[NGuyen et al, 2021]



Solem et al, 2015

Non-coding mutations

lncRNAs, miRNAs, structure-associated (RiboSnitches)

β -thalassaemia, duchenne muscular dystrophy, Cystic fibrosis, Rett syndrome...

RiboNucleic Acids (RNAs)



Encodes proteins

mRNA Vaccines

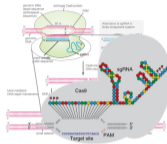
COVID-19, Malaria (Zika, CMV, Cancers?)

RiboNucleic Acids (RNA) in Human biology/health: Friend and Foe!

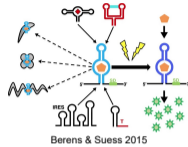
Targeting system for DNA Editing

CRISPR therapies

Sickle-cell anemia, β -thalassaemia, Leber congenital amaurosis (LCA), cancers...



Hendel et al, 2015; Agrotis & Ketteler, 2015



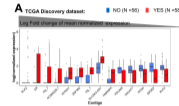
Berens & Suess 2015

Sensor of metabolites
Riboswitches

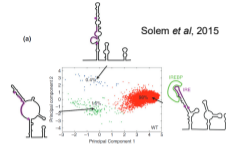
Quantitative expression

Transcriptomic signatures

Cancer diagnosis/prognosis/relapse...



[NGuyen et al, 2021]



Solem et al, 2015

Non-coding mutations

lncRNAs, miRNAs, structure-associated (RiboSnitches)

β -thalassaemia, duchenne muscular dystrophy,
Cystic fibrosis, Rett syndrome...

RiboNucleic Acids (RNAs)



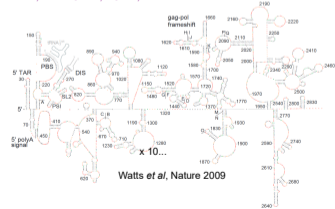
Encodes proteins

mRNA Vaccines

COVID-19, Malaria (Zika, CMV, Cancers?)

Genomic material for Human pathogens

HIV-1, SARS-CoV 2, HCoVs, MERS

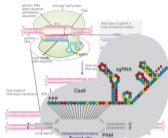


RiboNucleic Acids (RNA) in Human biology/health: Friend and Foe!

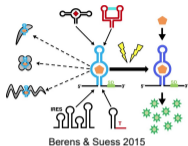
Targeting system for DNA Editing

CRISPR therapies

Sickle-cell anemia, β -thalassaemia, Leber congenital amaurosis (LCA), cancers...



Hendel et al, 2015; Agrotis & Ketteler, 2015

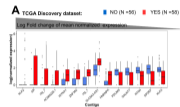


Sensor of metabolites
Riboswitches

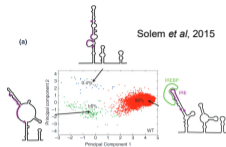
Quantitative expression

Transcriptomic signatures

Cancer diagnosis/prognosis/relapse...



[NGuyen et al, 2021]



Non-coding mutations

lncRNAs, miRNAs, structure-associated (RiboSnitches)

β -thalassaemia, duchenne muscular dystrophy,
Cystic fibrosis, Rett syndrome...

RiboNucleic Acids (RNAs)



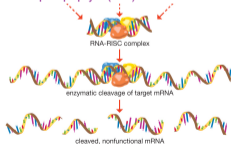
Regulation of gene expression

RNAi therapies (FDA approved)

Primary hyperoxaluria type 1 (PH1),

Hereditary transthyretin amyloidosis (ATTRv),

Acute hepatic porphyria (AHP)



Encyclopaedia Britannica, Inc 2013



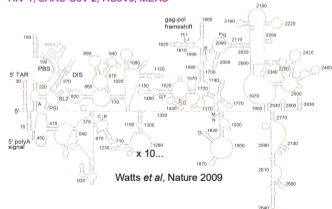
Encodes proteins

mRNA Vaccines

COVID-19, Malaria (Zika, CMV, Cancers?)

Genomic material for Human pathogens

HIV-1, SARS-CoV 2, HCoVs, MERS



RiboNucleic Acids (RNA) in Human biology/health: Friend **and** Foe!

Targeting system for DNA Editing

CRISPR therapies

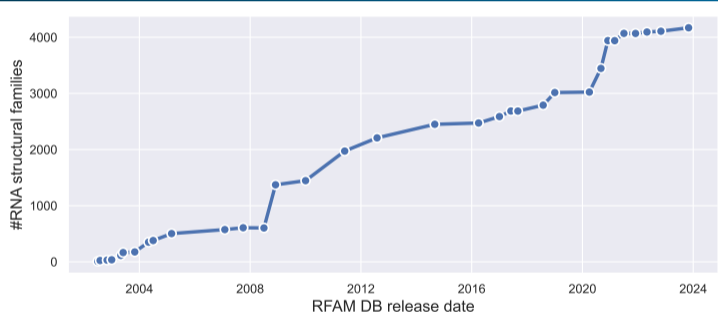
Sickle-cell anemia, β -thalassaemia, Leber congenital amaurosis (β CA), cancer



Quantitative expression

Transcriptomic signatures
Cancer diagnosis/prognosis/relapse...

Solem et al, 2015



RNA functional diversity is (largely) enabled by deep structural diversity

Regul
RNA t
Primar
Heredi
Acute t

devant, [https://doi.org/10.1093/nar/nkz1071](#)
Encyclopaedia Britannica, Inc. 2013

mRNA Vaccines
COVID-19, Malaria (Zika, CMV, Cancers?)

Watts et al, Nature 2009

RiboNucleic Acids (RNA) in Human biology/health: Friend and Foe!

Targeting system for DNA Editing

CRISPR therapies
Sickle-cell anemia, β -thalassaemia, Leber congenital amaurosis (LCA), cancers...

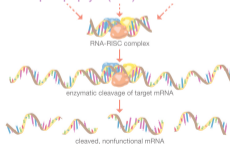


Hendel et al, 2015; Agrotis & Ketteler, 2015

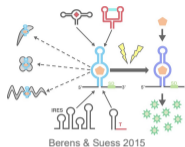
Rational design

Regulation of gene expression

RNAi therapies (FDA approved)
Primary hyperoxaluria type 1 (PH1),
Hereditary transthyretin amyloidosis (ATTRv),
Acute hepatic porphyria (AHP)



Encyclopaedia Britannica, Inc 2013

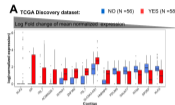


Berens & Suess 2015

Sensor of metabolites
Riboswitches

Quantitative expression

Transcriptomic signatures
Cancer diagnosis/prognosis/relapse...

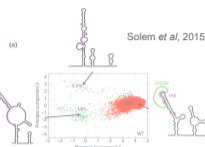


[NGuyen et al, 2021]

RiboNucleic Acids (RNAs)



Encodes proteins
mRNA Vaccines
COVID-19, Malaria (Zika, CMV, Cancers?)



Solem et al, 2015

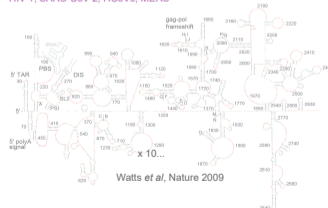
Non-coding mutations

lncRNAs, miRNAs, structure-associated (RiboSnitches)
 β -thalassaemia, duchenne muscular dystrophy,
Cystic fibrosis, Rett syndrome...

(2D) Structure Modeling

Genomic material for Human pathogens

HIV-1, SARS-CoV 2, HCoV, MERS



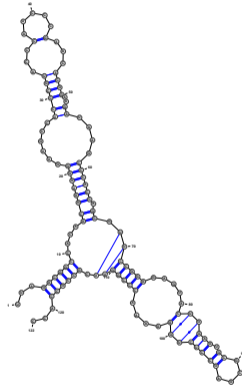
Watts et al, Nature 2009

RNA structure(s)

RNA = Linear Polymer = Nucleotides sequence $w \in \{A, C, G, U\}^*$

```
UUAGGCGGCCACAGC
GGUGGGGUUGCCUCC
CGUACCCAUCCGAA
CACGGAAGAUAGCC
CACCAGCGUCCGGG
GAGUACUGGAGUGCG
CGAGCCUCUGGAAA
CCCGGUUCGCCCA
CC
```

Primary struct.



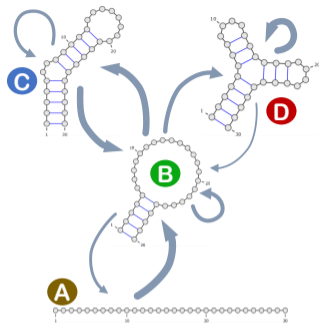
Secondary (2D) struct.



Tertiary (\approx 3D) struct.

Source: 5s rRNA (PDBID: 1K73:B)

Paradigms in RNA structural bioinformatics



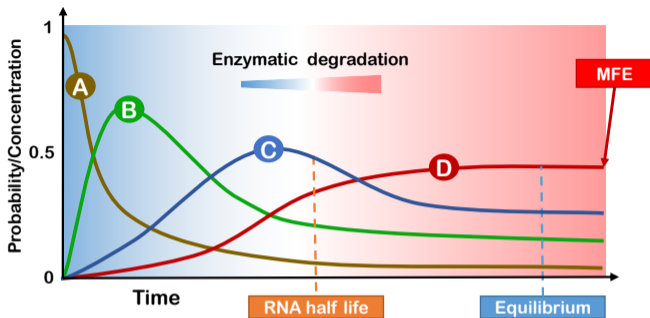
A – Kinetic Landscape

Continuous-time Markov chain

Given **free-energy** $E : \{A, C, G, U\}^* \times \mathcal{S} \rightarrow \mathbb{R}$, at the Boltzmann equilibrium one has:

$$\mathbb{P}(S | w) = e^{-E(w,S)/RT} / \mathcal{Z}(S) \quad (\mathcal{Z} \text{ partition function})$$

- ▶ **Minimum Free-Energy (MFE)**: Relevant structure = Most stable/probable
- ▶ **Partition function**: Equilibrium properties (stationary distribution)
- ▶ **Kinetics**: Finite-time dynamics of concentrations/probabilities



B – Evolution of concentrations

O(99) reasons to perform a rational design of structural RNAs

1. To stress test our understanding of how RNA folds
Misfolded RNAs reveal gaps in our energy models and conformational descriptions
2. To create building blocks for synthetic systems
Rationally-designed RNAs increase orthogonality
3. To assess the significance of observed phenomenon
Random models should include all established traits, including adopting a well-defined structure
4. To help search for homologous sequences (remote homology)
Include designed/unseen homologs in multiple sequence alignments (e.g. cov. models)
5. To perform controlled experiments
Test statistical support of theories (w/o confirmation bias)
6. To fuel RNA-based therapeutics
Sequence-based (siRNA, synthetic genes), but structure also matters

O(99) reasons to perform a rational design of structural RNAs

1. To stress test our understanding of how RNA folds
Misfolded RNAs reveal gaps in our energy models and conformational descriptions
2. To create building blocks for synthetic systems
Rationally-designed RNAs increase orthogonality
3. To assess the significance of observed phenomenon
Random models should include all established traits, including adopting a well-defined structure
4. To help search for homologous sequences (remote homology)
Include designed/unseen homologs in multiple sequence alignments (e.g. cov. models)
5. To perform controlled experiments
Test statistical support of theories (w/o confirmation bias)
6. To fuel RNA-based therapeutics
Sequence-based (siRNA, synthetic genes), but structure also matters

O(99) reasons to perform a rational design of structural RNAs

1. To stress test our understanding of how RNA folds
Misfolded RNAs reveal gaps in our energy models and conformational descriptions
2. To create building blocks for synthetic systems
Rationally-designed RNAs increase orthogonality
3. To assess the significance of observed phenomenon
Random models should include all established traits, including adopting a well-defined structure
4. To help search for homologous sequences (remote homology)
Include designed/unseen homologs in multiple sequence alignments (e.g. cov. models)
5. To perform controlled experiments
Test statistical support of theories (w/o confirmation bias)
6. To fuel RNA-based therapeutics
Sequence-based (siRNA, synthetic genes), but structure also matters

O(99) reasons to perform a rational design of structural RNAs

1. To stress test our understanding of how RNA folds
Misfolded RNAs reveal gaps in our energy models and conformational descriptions
2. To create building blocks for synthetic systems
Rationally-designed RNAs increase orthogonality
3. To assess the significance of observed phenomenon
Random models should include all established traits, including adopting a well-defined structure
4. To help search for homologous sequences (remote homology)
Include designed/unseen homologs in multiple sequence alignments (e.g. cov. models)
5. To perform controlled experiments
Test statistical support of theories (w/o confirmation bias)
6. To fuel RNA-based therapeutics
Sequence-based (siRNA, synthetic genes), but structure also matters

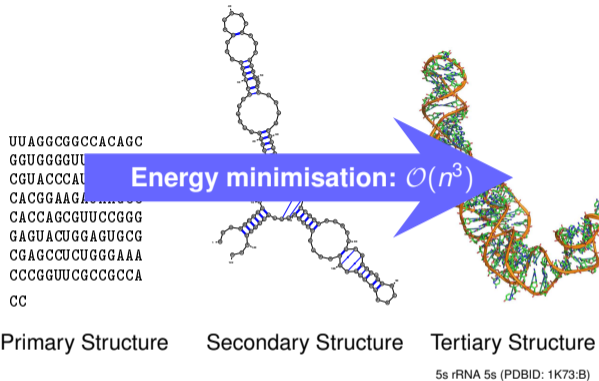
O(99) reasons to perform a rational design of structural RNAs

1. To stress test our understanding of how RNA folds
Misfolded RNAs reveal gaps in our energy models and conformational descriptions
2. To create building blocks for synthetic systems
Rationally-designed RNAs increase orthogonality
3. To assess the significance of observed phenomenon
Random models should include all established traits, including adopting a well-defined structure
4. To help search for homologous sequences (remote homology)
Include designed/unseen homologs in multiple sequence alignments (e.g. cov. models)
5. To perform controlled experiments
Test statistical support of theories (w/o confirmation bias)
6. To fuel RNA-based therapeutics
Sequence-based (siRNA, synthetic genes), but structure also matters

O(99) reasons to perform a rational design of structural RNAs

1. To stress test our understanding of how RNA folds
Misfolded RNAs reveal gaps in our energy models and conformational descriptions
2. To create building blocks for synthetic systems
Rationally-designed RNAs increase orthogonality
3. To assess the significance of observed phenomenon
Random models should include all established traits, including adopting a well-defined structure
4. To help search for homologous sequences (remote homology)
Include designed/unseen homologs in multiple sequence alignments (e.g. cov. models)
5. To perform controlled experiments
Test statistical support of theories (w/o confirmation bias)
6. To fuel RNA-based therapeutics
Sequence-based (siRNA, synthetic genes), but structure also matters

Minimum Free-Energy (MFE) folding



Nussinov's [PNAS, 1980] $\Theta(n^3)$ algorithm finds Min. Free-Energy structure (base-pairs)

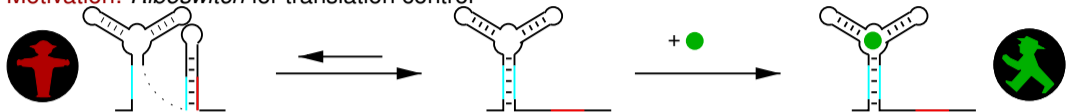
$$\text{MFE}_{i,j} = \min \left\{ \text{MFE}_{i+1,j}; \sum_k E(i,k) + \text{MFE}_{i+1,k-1} + \text{MFE}_{k+1,j} \right\}$$

Trivially adapted into joint OPT of Energy and Codon Adaptation Index for given protein sequence
→ **Yield-optimized mRNA vaccines** [Zhang *et al*, Nature 2023]

Positive multiple design

Positive design for multiple RNA structures

Motivation: *Riboswitch* for translation control



Multiple target structures

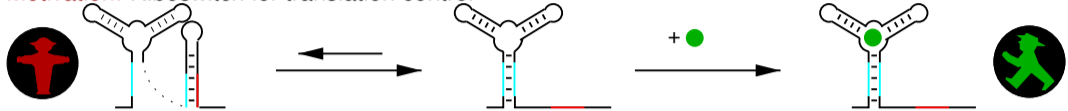
```
abcdefghijklmnopqrstuv  
(((((.)).(((..))).)).).  
((.))((...))..(((..)))  
....((((((..)))...))....
```

Objective: To randomly generate RNA sequences under constraints

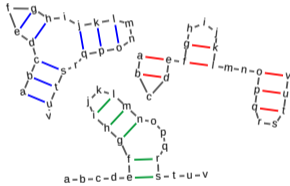
1. Validity for targeted structures wrt base pairing nucleotides
2. Stability (low free-energy, comparable across structures. . .) of target structures
3. Constrained composition: (prescribed G+C content), \pm motifs. . .

Positive design for multiple RNA structures

Motivation: *Riboswitch* for translation control



Multiple target structures



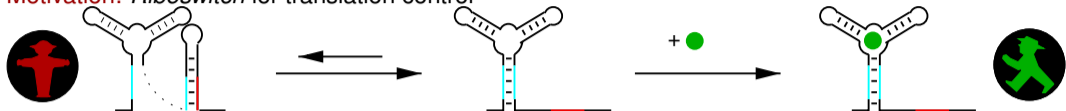
abcdefghijklmnopqrstuv
((((().).(((.)).)).).
((.))((...))..(((.)))
.....(((((.)))....))....

Objective: To randomly generate RNA sequences under constraints

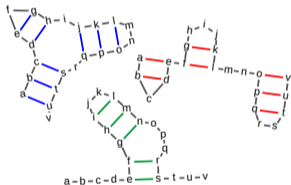
1. Validity for targeted structures wrt base pairing nucleotides
2. Stability (low free-energy, comparable across structures. . .) of target structures
3. Constrained composition: (prescribed G+C content), \pm motifs. . .

Positive design for multiple RNA structures

Motivation: *Riboswitch* for translation control



Multiple target structures

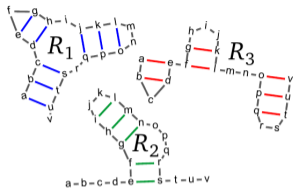


abcdefghijklmnopqrstuv
((((().)).(((.)).)).).
((.))((...))..(((.)))
.....(((((.))).....)

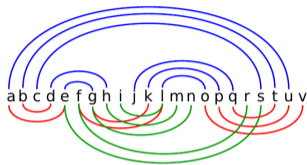
Objective: To **randomly** generate RNA sequences under constraints

1. **Validity** for targeted structures wrt base pairing nucleotides
2. **Stability** (low free-energy, comparable across structures...) of target structures
3. **Constrained composition**: (prescribed G+C content), \pm motifs...

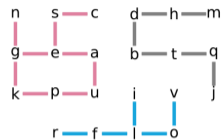
Counting the number of valid sequences



i) Input Structures



ii) Merged Base-Pairs



iii) Compatibility Graph

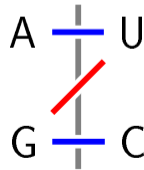
Question: How many valid sequences over $\Sigma^n := \{A, C, G, U\}^n$?

Problem (#ValidSequences)

Input: Secondary structures $\mathcal{R} = \{R_1, \dots, R_k\}$ of length n

Output: Number of valid sequences

$$\#Designs = |\{S \in \Sigma^n \mid \forall (i, j) \in R_\ell, (S_i, S_j) \text{ forms a valid base pair}\}|$$



Valid base pairs

Theorem (Designs \approx Independent sets)

Let G be a **bipartite and connected** dependency graph:

$$\#Designs(G) = 2 \times \#Designs^*(G) = 2 \times \#IndSets(G)$$

$$\Rightarrow \text{For general graphs: } \#Designs(G) = \prod_{cc \in CC(G)} 2 \times \#IndSets(cc) = 2^{|CC(G)|} \times \#IndSets(G)$$

But $\#IndSets(G)$ is **#P-hard** on bipartite graphs (**#BIS**) [Dyer & Greenhill'00]

Theorem (Classic counting complexity)

Counting $\#Designs$ is #P-hard.

No Poly-Time algorithm for $\#Designs(G)$ **unless** $\#P = FP (\Rightarrow P = NP)$

Theorem (Parameterized complexity for treewidth)

$\#Designs$ is Fixed-Parameter Tractable ($O(f(tw) \cdot P(n))$) for the **Treewidth** parameter tw

Theorem (Designs \approx Independent sets)

Let G be a **bipartite and connected** dependency graph:

$$\#Designs(G) = 2 \times \#Designs^*(G) = 2 \times \#IndSets(G)$$

$$\Rightarrow \text{For general graphs: } \#Designs(G) = \prod_{cc \in CC(G)} 2 \times \#IndSets(cc) = 2^{|CC(G)|} \times \#IndSets(G)$$

But $\#IndSets(G)$ is **#P-hard** on bipartite graphs (**#BIS**) [Dyer & Greenhill'00]

Theorem (Classic counting complexity)

Counting $\#Designs$ is **#P-hard**.

No Poly-Time algorithm for $\#Designs(G)$ **unless** $\#P = FP (\Rightarrow P = NP)$

Theorem (Parameterized complexity for treewidth)

$\#Designs$ is Fixed-Parameter Tractable ($O(f(tw) \cdot P(n))$) for the **Treewidth** parameter tw

Theorem (Designs \approx Independent sets)

Let G be a **bipartite and connected** dependency graph:

$$\#Designs(G) = 2 \times \#Designs^*(G) = 2 \times \#IndSets(G)$$

$$\Rightarrow \text{For general graphs: } \#Designs(G) = \prod_{cc \in CC(G)} 2 \times \#IndSets(cc) = 2^{|CC(G)|} \times \#IndSets(G)$$

But $\#IndSets(G)$ is **#P-hard** on bipartite graphs (**#BIS**) [Dyer & Greenhill'00]

Theorem (Classic counting complexity)

Counting $\#Designs$ is **#P-hard**.

No Poly-Time algorithm for $\#Designs(G)$ **unless** $\#P = FP (\Rightarrow P = NP)$

Theorem (Parameterized complexity for treewidth)

$\#Designs$ is Fixed-Parameter Tractable ($O(f(tw) \cdot P(n))$) for the **Treewidth** parameter tw

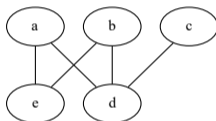
Tree decomposition and treewidth

A **tree decomposition** T for graph $G = (V, E)$:

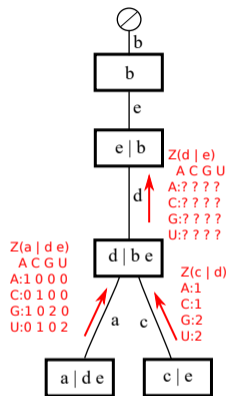
1. Nodes of $T =$ Bags, *i.e.* subsets of V ;
2. Every vertex must be found in ≥ 1 bag;
3. Each edge must be represented in ≥ 1 bag;
4. Nodes featuring any $v \in V$ form a **connected** subtree of T

a b c d e
 (. .) .
 . (())
 ((.))

Target structures



Dependency graph



Tree decomposition

w : **Width** of tree decomposition T ($=\max_{b \in B} |b| - 1$)

Let $b = (v; v_1 \dots) \subseteq V$ a bag of T , and T_b be the subtree rooted at b

$$\# \text{Designs}(T_b | b_2 \leftarrow v_2 \dots) = \sum_{\substack{b_1 \leftarrow v_1 \\ v_1 \in \{A, C, G, U\}}} \prod_{c \text{ child of } b} \# \text{Designs}(T_c | b_1 \leftarrow v_1, b_2 \leftarrow v_2 \dots)$$

\rightarrow **#Designs** (resp. **partition function**) computable in $\Theta(n k 2^w)$ time for k struct. of length n

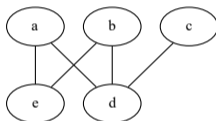
Tree decomposition and treewidth

A **tree decomposition** T for graph $G = (V, E)$:

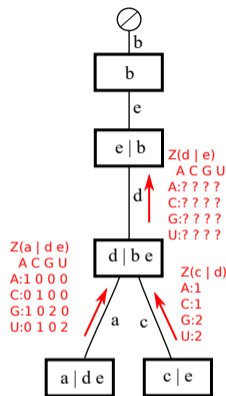
1. Nodes of $T =$ Bags, *i.e.* subsets of V ;
2. Every vertex must be found in ≥ 1 bag;
3. Each edge must be represented in ≥ 1 bag;
4. Nodes featuring any $v \in V$ form a **connected** subtree of T

a	b	c	d	e
(.	.)	.
.	(())
((.))

Target structures



Dependency graph



w : **Width** of tree decomposition T ($= \max_{b \in B} |b| - 1$)

Let $b = (v; v_1 \dots) \subseteq V$ a bag of T , and T_b be the subtree rooted at b **Tree decomposition**

$$\# \text{Designs}(T_b \mid b_2 \leftarrow v_2 \dots) = \sum_{\substack{b_1 \leftarrow v_1 \\ v_1 \in \{A, C, G, U\}}} \prod_{c \text{ child of } b} \# \text{Designs}(T_c \mid b_1 \leftarrow v_1, b_2 \leftarrow v_2 \dots)$$

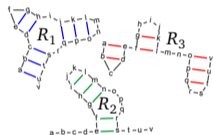
\rightarrow **#Designs** (resp. **partition function**) computable in $\Theta(n k 2^w)$ time for k struct. of length n

Tree decomposition and Boltzmann sampling of sequences

First **count** (FPT), then **stochastic backtrack** $\Theta(n)$, based on Min width (tw) decomposition (FPT)

Theorem

Positive multiple design (unif./Boltzmann distr.) FPT for the treewidth parameter



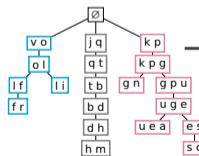
i) Input Structures



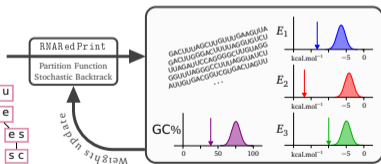
ii) Merged Base-Pairs



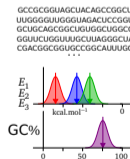
iii) Compatibility Graph



iv) Tree Decomposition



v) Weight Optimization (Adaptive Sampling)



vi) Final Designs

RNARedPrint [Hammer, P, Wang, Will, RECOMB 2018 & BMC Bioinfo 2019]

Infrared, a declarative (weighted) constraint satisfaction framework

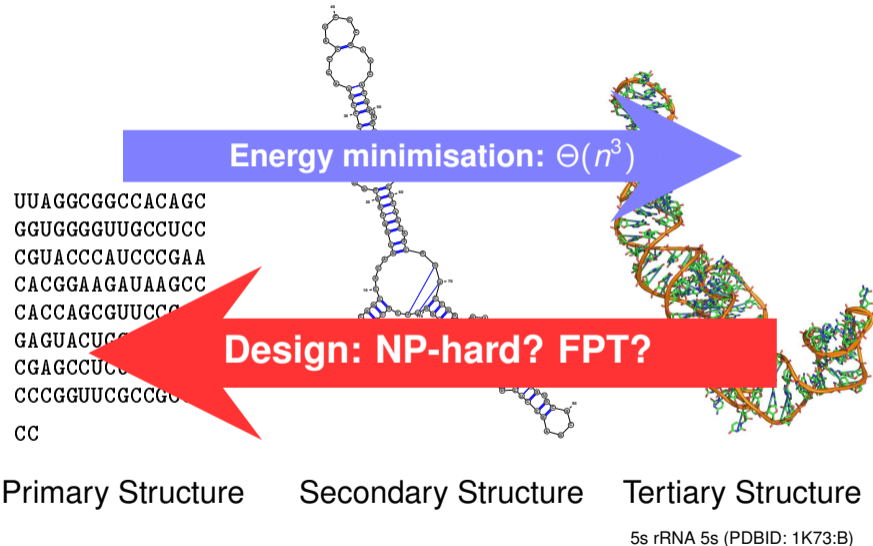
```
import infrared as ir
import infrared.rna as rna
n, bps = len(target), rna.parse(target)
model = ir.Model(n, 4)
model.add_constraints(rna.BPComp(i, j) for (i, j) in bps)
model.add_functions([rna.GCCont(i) for i in range(n)], 'gc')
model.add_functions([rna.BPEnergy(i, j, (i-1, j+1) not in bps)
                    for (i, j) in bps], 'energy')
model.set_feature_weight(-1.5, 'energy')
sampler = ir.Sampler(model)
samples = [sampler.sample() for _ in range(10)]
```

InfraRed [Yao *et al*, *Algorithms Mol Biol* 2024] generalizes **RNARedPrint** beyond RNA design tasks:

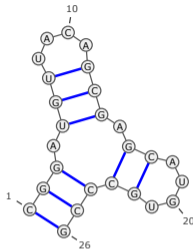
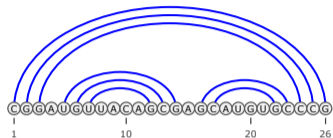
- ▶ Generic solver for sparse/weighted constraints networks, fueled by tree decomposition;
- ▶ Supports: optimization, exact sampling (unif./Boltzmann distr.), integers-value feature targets;
- ▶ Critical sections in C, conveniently interfaced in Python
- ▶ Illustrated on threading, network-based parsimony, alignment...

Inverse folding

Minimum Free-Energy (MFE) folding



Energy model



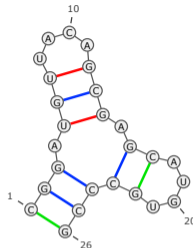
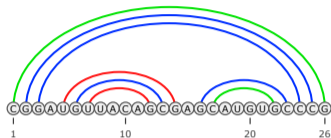
This section: Restriction to **valid** base-pairs = $\{(A, U), (G, C), (G, U)\}$

- ▶ **RNA structure R :** Set of non-crossing base pairs (BPs)
- ▶ **Motifs:** Connected positions + content (e.g. Base Pairs, Stacking, Loops...)
- ▶ **Energy model:**

Motif \rightarrow Free-energy contribution $\Delta(m, a) \in \mathbb{R} \cup \{+\infty\}$, $m \subset [1, n]$, $a \in \Sigma^{|m|}$

Free-energy $E(S, R)$: Sum of energies for motifs in R , given sequence S

Energy model



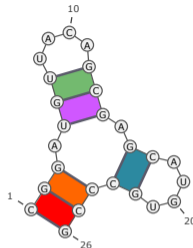
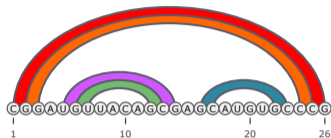
This section: Restriction to **valid** base-pairs = $\{(A, U), (G, C), (G, U)\}$

- ▶ **RNA structure R :** Set of non-crossing base pairs (BPs)
- ▶ **Motifs:** Connected positions + content (e.g. Base Pairs, Stacking, Loops...)
- ▶ **Energy model:**

Motif \rightarrow Free-energy contribution $\Delta(m, a) \in \mathbb{R} \cup \{+\infty\}$, $m \subset [1, n]$, $a \in \Sigma^{|m|}$

Free-energy $E(S, R)$: Sum of energies for motifs in R , given sequence S

Energy model



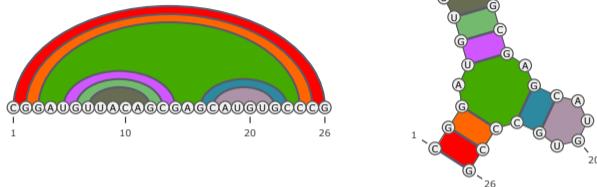
This section: Restriction to **valid** base-pairs = $\{(A, U), (G, C), (G, U)\}$

- ▶ **RNA structure R :** Set of non-crossing base pairs (BPs)
- ▶ **Motifs:** Connected positions + content (e.g. Base Pairs, Stacking, Loops...)
- ▶ **Energy model:**

Motif \rightarrow Free-energy contribution $\Delta(m, a) \in \mathbb{R} \cup \{+\infty\}$, $m \subset [1, n]$, $a \in \Sigma^{|m|}$

Free-energy $E(S, R)$: Sum of energies for motifs in R , given sequence S

Energy model



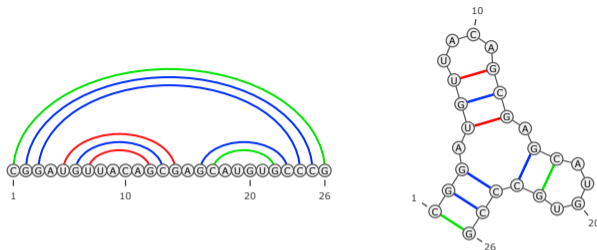
This section: Restriction to **valid** base-pairs = $\{(A, U), (G, C), (G, U)\}$

- ▶ **RNA structure R :** Set of non-crossing base pairs (BPs)
- ▶ **Motifs:** Connected positions + content (e.g. Base Pairs, Stacking, Loops...)
- ▶ **Energy model:**

Motif \rightarrow Free-energy contribution $\Delta(m, a) \in \mathbb{R} \cup \{+\infty\}$, $m \subset [1, n]$, $a \in \Sigma^{|m|}$

Free-energy $E(S, R)$: Sum of energies for motifs in R , given sequence S

Energy model



This section: Restriction to **valid** base-pairs = $\{(A, U), (G, C), (G, U)\}$

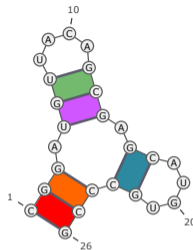
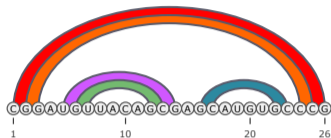
- ▶ **RNA structure R** : Set of non-crossing base pairs (BPs)
- ▶ **Motifs**: Connected positions + content (e.g. Base Pairs, Stacking, Loops...)
- ▶ **Energy model**:

Motif \rightarrow Free-energy contribution $\Delta(m, a) \in \mathbb{R} \cup \{+\infty\}$, $m \subset [1, n]$, $a \in \Sigma^{|m|}$

Free-energy $E(S, R)$: Sum of energies for motifs in R , given sequence S

$$E_R = 2 \cdot \Delta \begin{pmatrix} \text{U} \\ | \\ \text{G} \end{pmatrix} + 4 \cdot \Delta \begin{pmatrix} \text{G} \\ | \\ \text{C} \end{pmatrix} + 2 \cdot \Delta \begin{pmatrix} \text{C} \\ | \\ \text{G} \end{pmatrix}$$

Energy model



This section: Restriction to **valid** base-pairs = $\{(A, U), (G, C), (G, U)\}$

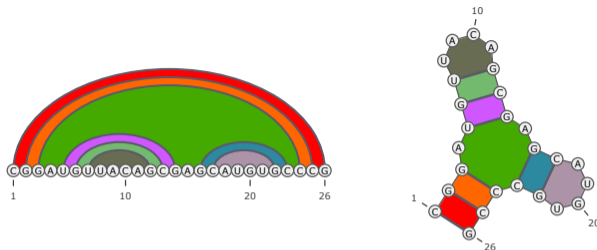
- ▶ **RNA structure R :** Set of non-crossing base pairs (BPs)
- ▶ **Motifs:** Connected positions + content (e.g. Base Pairs, Stacking, Loops...)
- ▶ **Energy model:**

Motif \rightarrow Free-energy contribution $\Delta(m, a) \in \mathbb{R} \cup \{+\infty\}$, $m \subset [1, n]$, $a \in \Sigma^{|m|}$

Free-energy $E(S, R)$: Sum of energies for motifs in R , given sequence S

$$E_R = \Delta \left(\begin{array}{cc} \text{C} & \text{G} \\ | & | \\ \text{G} & \text{C} \end{array} \right) + \Delta \left(\begin{array}{cc} \text{G} & \text{G} \\ | & | \\ \text{C} & \text{C} \end{array} \right) + \Delta \left(\begin{array}{cc} \text{U} & \text{G} \\ | & | \\ \text{G} & \text{C} \end{array} \right) + \Delta \left(\begin{array}{cc} \text{U} & \text{G} \\ | & | \\ \text{G} & \text{C} \end{array} \right) + \Delta \left(\begin{array}{cc} \text{U} & \text{G} \\ | & | \\ \text{G} & \text{C} \end{array} \right)$$

Energy model



This section: Restriction to **valid** base-pairs = $\{(A, U), (G, C), (G, U)\}$

- ▶ **RNA structure R** : Set of non-crossing base pairs (BPs)
- ▶ **Motifs**: Connected positions + content (e.g. Base Pairs, Stacking, Loops...)
- ▶ **Energy model**:

Motif \rightarrow Free-energy contribution $\Delta(m, a) \in \mathbb{R} \cup \{+\infty\}$, $m \subset [1, n]$, $a \in \Sigma^{|m|}$

Free-energy $E(S, R)$: Sum of energies for motifs in R , given sequence S

$$\begin{aligned}
 E_R = & \Delta \left(\begin{array}{c} \text{C} \quad \text{G} \\ | \quad | \\ \text{G} \quad \text{C} \end{array} \right) + \Delta \left(\begin{array}{c} \text{G} \quad \text{G} \\ | \quad | \\ \text{C} \quad \text{C} \end{array} \right) + \Delta \left(\begin{array}{c} \text{U} \quad \text{G} \\ | \quad | \\ \text{G} \quad \text{C} \end{array} \right) + \Delta \left(\begin{array}{c} \text{U} \quad \text{G} \\ | \quad | \\ \text{G} \quad \text{C} \end{array} \right) + \Delta \left(\begin{array}{c} \text{U} \quad \text{G} \\ | \quad | \\ \text{G} \quad \text{C} \end{array} \right) \\
 & + \Delta \left(\begin{array}{c} \text{A} \quad \text{C} \\ | \quad | \\ \text{U} \quad \text{A} \end{array} \right) + \Delta \left(\begin{array}{c} \text{A} \quad \text{C} \\ | \quad | \\ \text{U} \quad \text{A} \end{array} \right) + \Delta \left(\begin{array}{c} \text{C} \quad \text{A} \\ | \quad | \\ \text{G} \quad \text{U} \end{array} \right)
 \end{aligned}$$

Definition (INVERSE-FOLDING(E) problem)

Input: Secondary structure R + Energy distance $\Delta > 0$

Output: RNA sequence $S \in \Sigma^*$ such that:

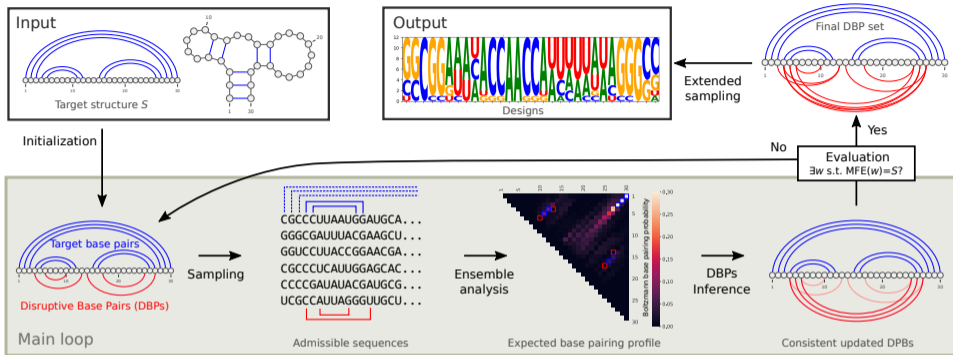
$$\forall R' \in \mathcal{S}_{|S|} \setminus \{R\} : E(S, R') \geq E(S, R) + \Delta$$

or \emptyset if no such sequence exists.

Difficult problem: Probably no **obvious** DP decomposition

- ▶ Informally introduced by [Hofacker *et al* Monatshefte für Chemie/Chemical Monthly 1994]
- ▶ NP-hardness by including energy model in input [Schnall-Levin *et al*, ICML'08]
- ▶ NP-hardness for BP maximization with partial assignment [Bonnet *et al*, RECOMB'18]
- ▶ **Reason(s):** Non locality, no theoretical framework, too many parameters. . .
- ▶ Existing algorithms: Heuristics or Exponential-time

RNA POsitive and Negative Design (RNAPOND) [Yao *et al*, RECOMB 2021]



RNAPond: Human-inspired heuristics based on identification of **Disruptive Base Pairs (DBPs)**

- ▶ Sample sequences compatible with target & avoiding DBPs ← Infrared (NP-hard, FPT on treewidth)
- ▶ Identify and forbid recurrent DPBs
- ▶ Iterate until solution found or treewidth threshold reached (def. $tw \leq 10$)

Close to state of the art heuristics in a crowded field

Existing approaches for negative design

Bio-inspired algorithms...

- ▶ FRNAKenstein - Hein@Oxford
- ▶ AntaRNA - Backofen@Freiburg
- ▶ ERD - Ganjtabesh@Tehran

... exact (exptime) approaches...

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

... based on local search...

- ▶ RNAInverse - TBI Vienna
 - ▶ Info-RNA - Backofen@Freiburg
 - ▶ RNA-SSD - Condon@UBC
 - ▶ (Inca)RNAFBinv - Barash@BGU
 - ▶ NUPack - Pierce@Caltech
- ... or ML/DL (Ribodiffusion...)

Typical issues:

- ▶ Single solution
- ▶ Strong impact of initialization strategy
- ▶ Synthesized sequences do not necessarily fold properly (kinetics)
- ▶ Overly GC-rich sequences
- ▶ Generative ML usually fails to generalize
- ▶ Few options to produce negative results

⇒ **Establish combinatorial foundations!**

Inverse Folding in unitary Base Pair energy model (aka BP maximization)

Definition (INVERSE-FOLDING(BPmax) problem)

Input: Target secondary structure R , i.e. a set of base pairs

Output: RNA sequence $S \in \Sigma^*$ such that:

- ▶ Sequence S valid for structure R
- ▶ If S valid sequence for alt structure $R' \neq R$, then $|R'| < |R|$



Inverse Folding in unitary Base Pair energy model (aka BP maximization)

Definition (INVERSE-FOLDING(BPmax) problem)

Input: Target secondary structure R , *i.e.* a set of base pairs

Output: RNA sequence $S \in \Sigma^*$ such that:

- ▶ Sequence S valid for structure R
- ▶ If S valid sequence for alt structure $R' \neq R$, then $|R'| < |R|$



Designability in simple BP-based energy models [Hales *et al*, CPM'15 & Algorithmica'17]

Partial characterization of **designable** structures

- ▶ **Saturated structures (all positions paired):** Designable \Leftrightarrow Multiloops degrees ≤ 4 (+ $\Theta(n)$ algo.)
- ▶ Designable \Rightarrow Avoid multiloops with *degree* ≥ 5 (m_5), or *degree* ≥ 3 with ≥ 1 *unpaired* ($m_3 \circ$).
Corollary: Fraction of designable structures decreases exponentially with n [Yao *et al*, ACM-BCB'19]

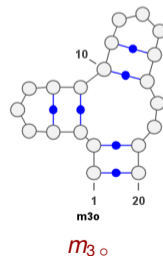
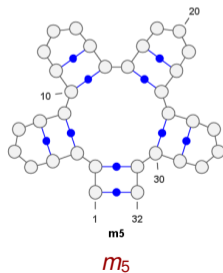
Designability in simple BP-based energy models [Hales *et al*, CPM'15 & Algorithmica'17]

Partial characterization of **designable** structures

- ▶ **Saturated structures (all positions paired)**: Designable \Leftrightarrow Multiloops degrees ≤ 4 (+ $\Theta(n)$ algo.)
- ▶ Designable \Rightarrow Avoid multiloops with **degree ≥ 5 (m_5)**, or **degree ≥ 3 with ≥ 1 unpaired ($m_{3\circ}$)**.

Theorem: Similar motifs exist for any **energy model** and **design criterion**

Corollary: Fraction of designable structures decreases **exponentially** with n [Yao *et al*, ACM-BCB'19]



Designability in simple BP-based energy models [Hales *et al*, CPM'15 & Algorithmica'17]

Partial characterization of **designable** structures

▶ **Saturated structures (all positions paired)**: Designable \Leftrightarrow Multiloops degrees ≤ 4 (+ $\Theta(n)$ algo.)

▶ Designable \Rightarrow Avoid multiloops with **degree ≥ 5** (m_5), or **degree ≥ 3 with ≥ 1 unpaired** ($m_{3\circ}$).

Corollary: Fraction of designable structures decreases **exponentially** with n [Yao *et al*, ACM-BCB'19]

▶ \exists **Separated** coloring for structure \Rightarrow Designable (+ $\Theta(n)$ algo.)

Each base pair \rightarrow one out of 3 colors: ● \rightarrow G · C; ○ \rightarrow C · G; ● \rightarrow A · U or U · A.

Coloring rules: Within each loop, #● ≤ 1 , #○ ≤ 1 , #● ≤ 2 and #● + #○ < 2

Definitions:

▶ **Level** of a base pair = #● - #○ on path to root

▶ Coloring **separated** if ● base pairs and unpaired positions at **different** levels

Idea: Separated sequences (unpaired \rightarrow A) uniquely fold since alt BPs segregate regions with #G \neq #C

Separated Coloring (example)

Base pairs \rightarrow 3 colors:

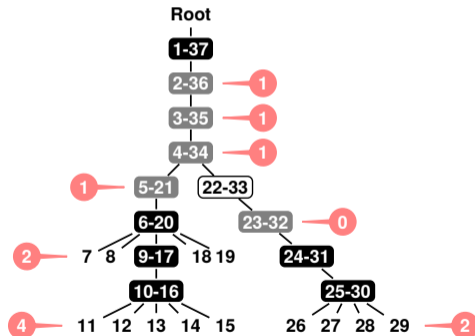
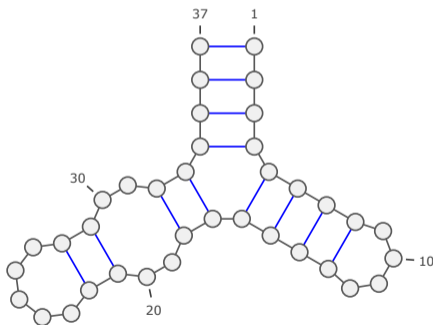
● \rightarrow G · C;

○ \rightarrow C · G;

● \rightarrow A · U or U · A.

Coloring rules: Within each loop, $\# \bullet \leq 1$, $\# \circ \leq 1$, $\# \bullet \leq 2$ and $\# \bullet + \# \circ < 2$

((((((((.....))))))((..(((.....))..))))))



GAAAGUUGGUUUUUCUUCUCAGGUUUUCCUGUUUC

Designability in simple BP-based energy models [Hales *et al*, CPM'15 & Algorithmica'17]

Partial characterization of **designable** structures

▶ **Saturated structures (all positions paired)**: Designable \Leftrightarrow Multiloops degrees ≤ 4 (+ $\Theta(n)$ algo.)

▶ Designable \Rightarrow Avoid multiloops with **degree ≥ 5** (m_5), or **degree ≥ 3 with ≥ 1 unpaired** ($m_{3\circ}$).

Theorem: Similar motifs exist for any **energy model** and **design criterion**

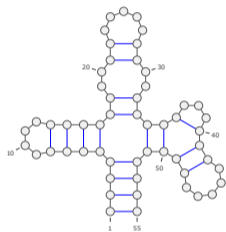
Corollary: Fraction of designable structures decreases **exponentially** with n [Yao *et al*, ACM-BCB'19]

▶ \exists **Separated** coloring for structure \Rightarrow Designable (+ $\Theta(n)$ algo.)

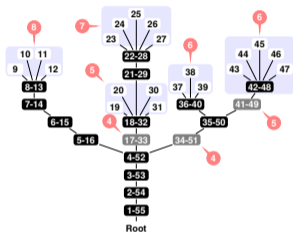
Corollary: **Approximate design** for any structure avoiding m_5 and $m_{3\circ}$ in $\Theta(n)$ time

Idea: Shift unpaired/leaves and  to **odd/even** levels resp. by **adding ≤ 1 BP** in each helix

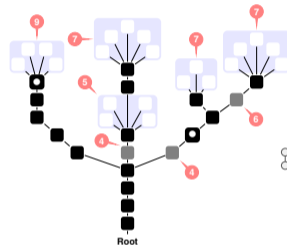
Example of structure-approximating design



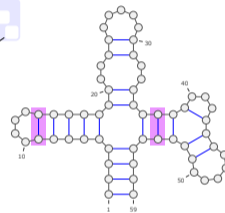
1) Target structure



2) Greedy proper coloring



3) Separated proper coloring



4) Designable structure

Designability in simple BP-based energy models [Hales *et al*, CPM'15 & Algorithmica'17]

Partial characterization of **designable** structures

▶ **Saturated structures (all positions paired)**: Designable \Leftrightarrow Multiloops degrees ≤ 4 (+ $\Theta(n)$ algo.)

▶ Designable \Rightarrow Avoid multiloops with **degree ≥ 5** (m_5), or **degree ≥ 3 with ≥ 1 unpaired** ($m_{3\circ}$).

Theorem: Similar motifs exist for any **energy model** and **design criterion**

Corollary: Fraction of designable structures decreases **exponentially** with n [Yao *et al*, ACM-BCB'19]

▶ \exists **Separated** coloring for structure \Rightarrow Designable (+ $\Theta(n)$ algo.)

Corollary: **Approximate design** for any structure avoiding m_5 and $m_{3\circ}$ in $\Theta(n)$ time

Idea: Shift unpaired/leaves and  to **odd/even** levels resp. by **adding ≤ 1 BP** in each helix

Remark: In structures with **helices of length 3+** offsetting always possible

Theorem (Unpublished!)

Inverse folding solvable in Polytime ($\Theta(n)$) for target structures with 3+ BPs helices

Is inverse folding *really* NP hard?

Designability in simple BP-based energy models [Hales *et al*, CPM'15 & Algorithmica'17]

Partial characterization of **designable** structures

▶ **Saturated structures (all positions paired)**: Designable \Leftrightarrow Multiloops degrees ≤ 4 (+ $\Theta(n)$ algo.)

▶ Designable \Rightarrow Avoid multiloops with **degree ≥ 5** (m_5), or **degree ≥ 3 with ≥ 1 unpaired** ($m_{3\circ}$).

Theorem: Similar motifs exist for any **energy model** and **design criterion**

Corollary: Fraction of designable structures decreases **exponentially** with n [Yao *et al*, ACM-BCB'19]

▶ \exists **Separated** coloring for structure \Rightarrow Designable (+ $\Theta(n)$ algo.)

Corollary: **Approximate design** for any structure avoiding m_5 and $m_{3\circ}$ in $\Theta(n)$ time

Idea: Shift unpaired/leaves and  to **odd/even** levels resp. by **adding ≤ 1 BP** in each helix

Remark: In structures with **helices of length 3+** offsetting always possible

Theorem (Unpublished!)

Inverse folding solvable in Polytime ($\Theta(n)$) for target structures with 3+ BPs helices

Is inverse folding *really* NP hard?

Conclusions

- ▶ RNA design is a **timely** topic for Bio Maths/Computer Sciences with **practical consequences**
- ▶ **Realistic setting**: FPT algo for positive design and general heuristics for inverse folding
- ▶ Structure approximating design: a **tractable** alternative to a (possibly) NP hard model?
- ▶ **Simple BP maximization**: General $\Theta(n)$ algorithm for restricted inputs
- ▶ Forbidden motifs: **Ubiquitous** in DP-based inverse combinatorial optimization
- ▶ RNA structure seemingly harder **in theory** than **in practice**. **Why?**

Extensions and perspectives

- ▶ **Onwards to the bench:** Inspiration from exp. test of designs for SAM riboswitches
- ▶ More complex/realistic energy models (Stacks, Turner's Nearest Neighbors?)
Extended conformational spaces (pseudoknots, non-canonical BPs)
- ▶ Kinetics-aware design (prescribed intermediates, energy barriers. . .)
- ▶ (Parameterized) complexity of general inverse folding?
- ▶ Potential/limitations of Machine Learning towards RNA design:
Can ML learn **negative design strategies** from extent sequences? Novelty/orthogonality?
- ▶ NeuTral networks: **Exponentially** less designable structs (*aka* phenotypes) than initially thought
→ Refine phenotype/genotype studies?

Acknowledgments



Ecole Polytechnique

- ▶ H.T. Yao, B. Marchand, T. Boury
- ▶ S. Will, S. Berkemer
- ▶ M. Régnier, A. Héliou



Univ Gustave Eiffel

- ▶ L. Bulteau



Faculté Pharmacie@Univ Paris Cité

- ▶ B. Sargueil, P. Hardouin



Stat. Physics@ENS Paris

- ▶ J. Fernandez de Cossio Diaz
- ▶ S. Cocco, R. Monasson, J.



Simon Fraser University

- ▶ J. Hales, J. Manuch, L. Stacho
- ▶ C. Chauve



McGill University

- ▶ J. Waldispühl



Université du Québec à Montréal

- ▶ V. Reinharz



University of Vienna

- ▶ S. Hammer, R. Lorenz



Ben Gurion University

- ▶ D. Barash, M. Drory, A. Churkin

Support



Supplementary Slides

Consequences

Corollary (#Approximability for ≤ 5 structures) [Weitz'06]

For ≤ 5 structures (crossings allowed), #Design(G) can be approximated within **any ratio** in **Poly-time** (PTAS)

Corollary (#BIS-hardness for > 5 structures) [Cai, Galanis *et al*'16]

For more than 5 structures (crossings allowed), #Design is **equally as hard** to approximate as general #BIS.

Why crossings/Pseudokots? Because any bipartite graph of max degree Δ can be decomposed in Δ matchings in Poly-Time (Vizing theorem).

Connection between counting and sampling [Jerrum/Valiant/Vazirani'86].

Conjecture (#BIS-hardness of multiple positive design)

Quasi-uniform generation as hard as approximation of general #BIS

\Rightarrow Sampling #P hard?

Consequences

Corollary (#Approximability for ≤ 5 structures) [Weitz'06]

For ≤ 5 structures (crossings allowed), #Design(G) can be approximated within **any ratio** in **Poly-time** (PTAS)

Corollary (#BIS-hardness for > 5 structures) [Cai, Galanis *et al*'16]

For more than 5 structures (crossings allowed), #Design is **equally as hard** to approximate as general #BIS.

Why crossings/Pseudokots? Because any bipartite graph of max degree Δ can be **decomposed** in Δ matchings in **Poly-Time** (Vizing theorem).

Connection between **counting and sampling** [Jerrum/Valiant/Vazirani'86].

Conjecture (#BIS-hardness of multiple positive design)

Quasi-uniform generation as hard as approximation of general #BIS

⇒ **Sampling #P hard?**

Consequences

Corollary (#Approximability for ≤ 5 structures) [Weitz'06]

For ≤ 5 structures (crossings allowed), #Design(G) can be approximated within **any ratio** in **Poly-time** (PTAS)

Corollary (#BIS-hardness for > 5 structures) [Cai, Galanis *et al*'16]

For more than 5 structures (crossings allowed), #Design is **equally as hard** to approximate as general #BIS.

Why crossings/Pseudokots? Because any bipartite graph of max degree Δ can be **decomposed** in Δ matchings in **Poly-Time** (Vizing theorem).

Connection between **counting** and **sampling** [Jerrum/Valiant/Vazirani'86].

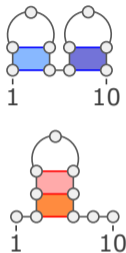
Conjecture (#BIS-hardness of multiple positive design)

Quasi-uniform generation as hard as approximation of general #BIS

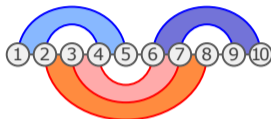
\Rightarrow **Sampling** #P hard?

Our problem for general free-energy models

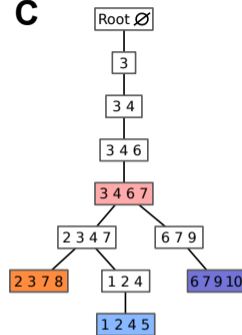
A



B



C



Question: Which partition function for **valid sequences**

Problem (PFDesigns)

Input: Structures $\mathcal{R} = \{R_1, \dots, R_k\}$ of length n + Weight (x_1, \dots, x_k)

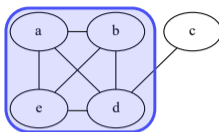
Output: Partition function

$$\mathcal{Z} = \sum_{\substack{S \in \Sigma^n \\ S \text{ valid for } \mathcal{R}}} \prod_{i=1}^k x_i^{E(S, R_i)}$$

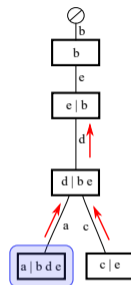
Counting/sampling, the Boltzmann-Gibbs way

a	b	c	d	e
(.	.)	.
.	(())
((.))

Target Structures



Dependency Hypergraph



Tree Decomposition

$b = \{b_1, b_2 \dots\}$: node of D
 T_b : subtree rooted at b
 w : **Width** of treedecomposition D

$$\mathcal{Z}(T_b | b_2 \leftarrow v_2 \dots) = \sum_{\substack{b_1 \leftarrow v_1 \\ v_1 \in \{A, C, G, U\}}} \prod_{i=1}^k x_i^{\sum_{E \in b} E(b, v_1, \dots)} \prod_{c \text{ child of } b} \mathcal{Z}(T_c | b_1 \leftarrow v_1, \dots)$$

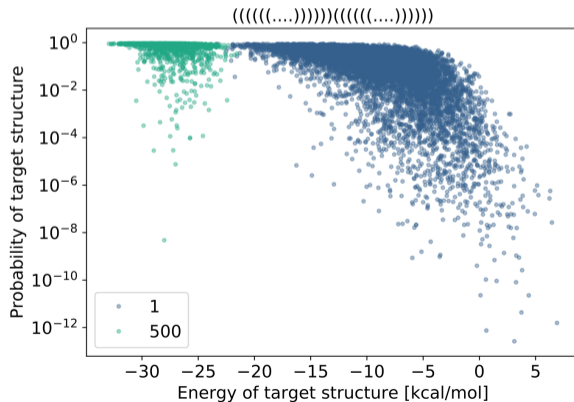
Complexity: $\Theta(nmk + nk2^{w+\#CC})$ for sampling in Boltzmann-Gibbs distrib.

Practical impact of Boltzmann-Gibbs sampling

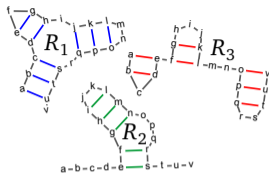
Boltzmann probability of **structure** R , pour une séquence S :

$$\mathbb{P}(R | S) = \frac{e^{-\frac{E(S,R)}{\beta T}}}{Z_S} \quad Z_S := \sum_{R'} e^{-\frac{E(S,R')}{\beta T}}$$

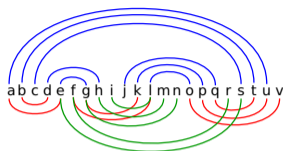
Objectif classique du design négatif (\rightarrow spécificité)



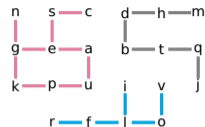
RNAredPrint: a flexible method for (positive) design



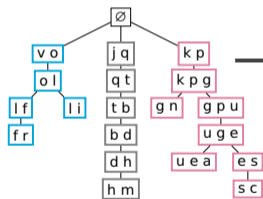
i) Input Structures



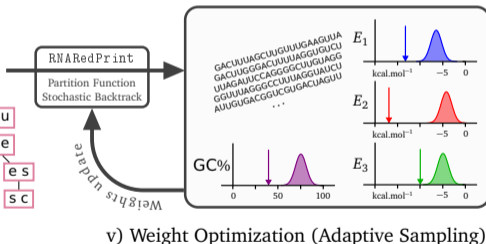
ii) Merged Base-Pairs



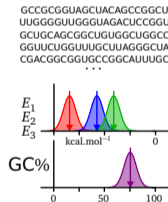
iii) Compatibility Graph



iv) Tree Decomposition



v) Weight Optimization (Adaptive Sampling)



vi) Final Designs

[Hammer/P/Wang/Will, RECOMB'18 + BMC Bioinfo 2019]

- ▶ Fixed Parameter Tractable algorithm based on tree width
- ▶ Uniform or Boltzmann-Gibbs sampling, to favor diversity and stability
- ▶ Multidimensional Boltzmann sampling for controlling free-energy, GC%...

<https://github.com/yannponty/RNAredPrint>

Multidimensional Boltzmann sampling

Multidimensional Boltzmann sampling [Bodini, P, DMTCS 2011]

Input: Targeted free-energies $(E_\ell^*)_{\ell=1}^k$, weights $(x_\ell)_{\ell=1}^k$ such that $\mathbb{E}(E(w, S_\ell)) = E_\ell^*, \forall \ell$:

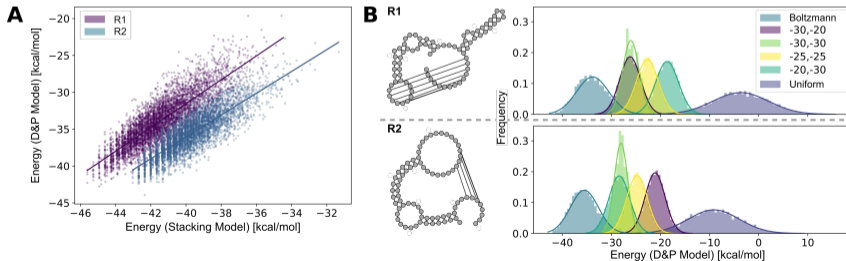
$$\mathbb{P}(w \mid x_1 \cdots x_k) \sim \prod_{\ell=1}^k x_\ell^{E(w, S_\ell)} + \text{Efficient rejection} \rightarrow \mathcal{O}(n^{k/2}) \text{ exact} / \mathcal{O}(\alpha^k) \text{ approx.}$$

Empirical efficiency for additive *concentrated* constraints (GC%, dinucleotides ...)

→ Partial functions → Hyper-edges, aka cliques¹

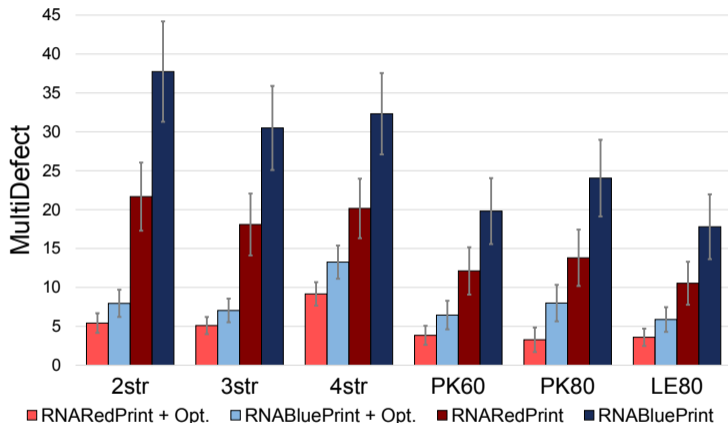


General framework for integer-valued constraints; Concentration tests.



¹But tree width ↗

Strangely enough, it actually works!



$$\text{MultiDefect}(S, R_1 \dots R_k) := \frac{\sum_{\ell=1}^k E(S, R_\ell) - \text{EFE}(S)}{k} + \frac{\sum_{1 \leq \ell < j \leq k} |E(S, R_\ell) - E(S, R_j)|}{2 \binom{k}{2}}$$