

(NP-hard?) RNA Design in linear time and space! (most of the time)

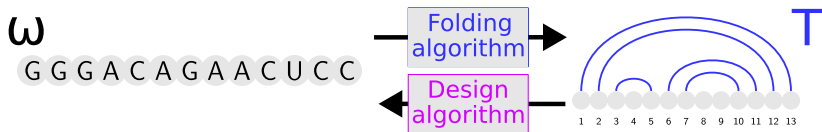
→ e.g. for structures without isolated stacks or base pairs

Théo Boury¹, Laurent Bulteau², Yann Ponty¹

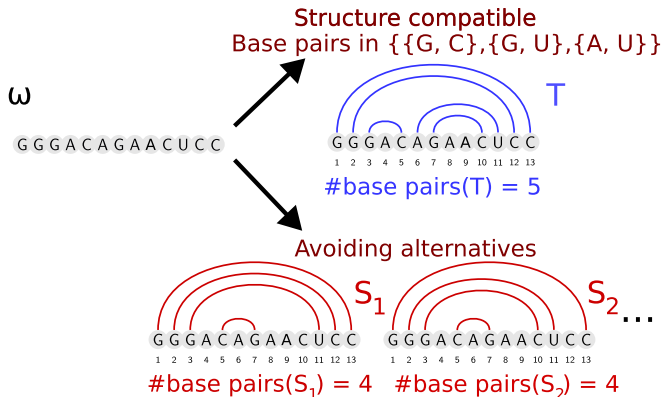
1, Laboratoire d'Informatique de l'Ecole Polytechnique (CNRS/LIX; UMR 7161), Institut Polytechnique de Paris, France

2, Laboratoire d'Informatique Gaspard Monge (CNRS/LIGM; UMR 8049), Université Gustave Eiffel, France

RNA 2D folding vs RNA structural design



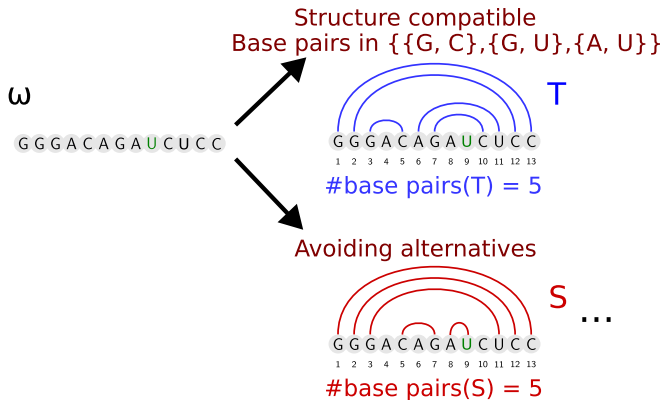
Inverse folding (IF): Formal definition



Goal: Find ω such that T is unique+optimal+valid fold for ω

$\forall S \neq T, S$ comp. with $\omega, \#BasePairs(S) < \#BasePairs(T)$

Inverse folding (IF): Formal definition



Goal: Find ω such that T is unique+optimal+valid fold for ω

$\forall S \neq T, S$ comp. with $\omega, \#BasePairs(S) < \#BasePairs(T)$

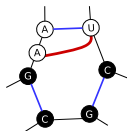
Minimal undesignable structures

[Halès et al, 2017]

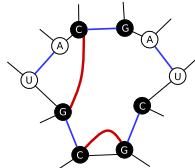
Undesignable structure



m3o motif

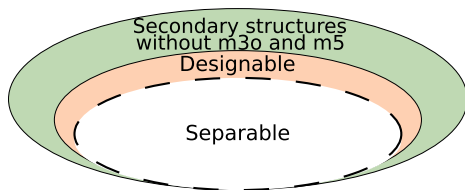


m5 motif



- ▶ Any target featuring occurrence of **m3o** or **m5** is not designable
- ▶ Inverse folding (+ minor constraints) is NP-hard [Bonnet et al, 2018]
→ Infinite (+ exp. growth) list of min undesignable motifs
(unless P=NP)
- ▶ Decision version of Inverse Folding **not reducible** to pattern matching

Design complexity



Complexity of Inverse Folding

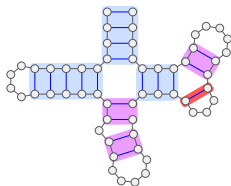
- Linear time and space solvable
- Almost-certainly NP-complete

Inverse Folding efficiently solved for non-bonsai structures

h_{min} = Minimum #base pairs in an helix

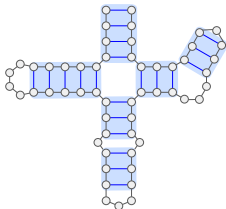
Theorem: Inverse folding (IF) is solved in linear time and space for secondary structures with $h_{min} = 3$. [This talk]

Structure with $h_{min} = 1$






Presumably NP-complete to design...

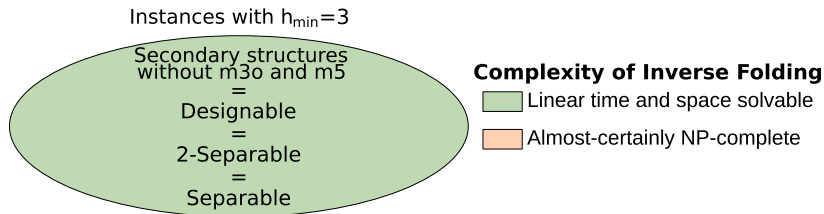
Structure with $h_{min} = 3$



Designed in linear time!

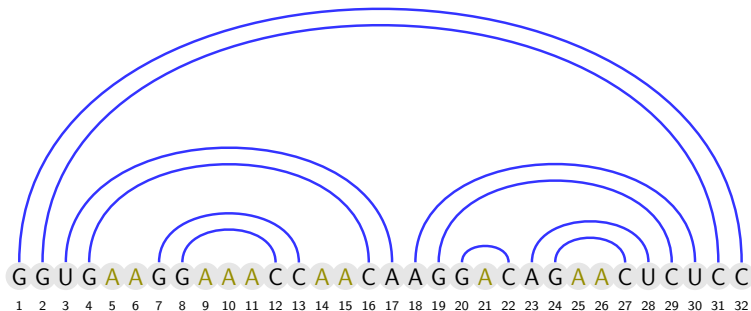
-  Helices of size 3+
-  Isolated stacks (e.g. helices of size 2)
-  Isolated base pairs (e.g. helices of size 1)

Design complexity when $h_{min} = 3$



Technical point: Separable structure

[Halès et al, 2017]



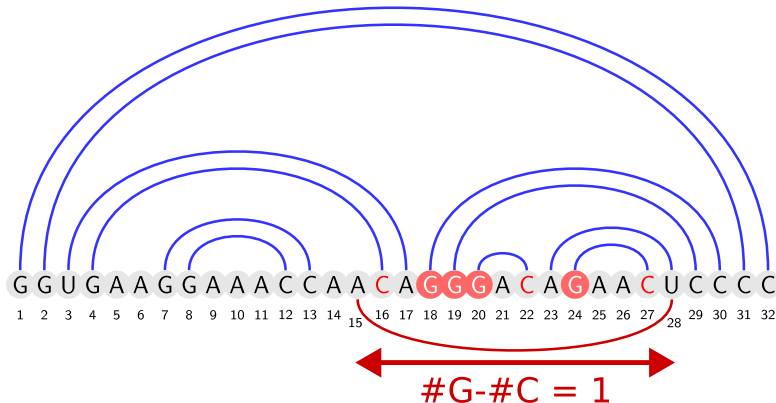
- ▶ Unpaired: A, Paired: G, U, C
- ▶ Separated sequence $\leftrightarrow \forall$ alternative A-U, $\#G - \#C \neq 0$
- ▶ Separable structure $\leftrightarrow \exists \omega$, separated sequence $\rightarrow \omega$, solution for IF

Result:

[This talk]

Deciding if a structure is separable is NP-complete

Technical point: Separable structure



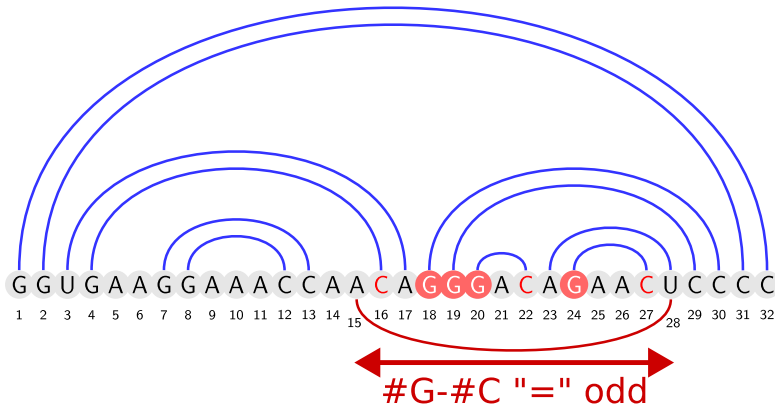
- ▶ Unpaired: A, Paired: G, U, C
- ▶ Separated sequence $\leftrightarrow \forall$ alternative A-U, $\#G - \#C \neq 0$
- ▶ Separable structure $\leftrightarrow \exists \omega$, separated sequence $\rightarrow \omega$, solution for IF

Result:

Deciding if a structure is separable is NP-complete

[This talk]

m -separable structure

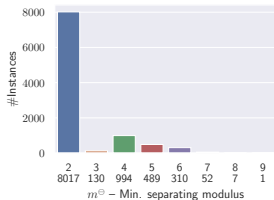


Result:

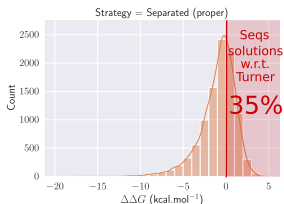
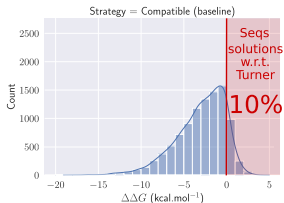
[This talk]

- ▶ Core: Structures with $h_{min} = 3 \rightarrow 2$ -separable
- ▶ In general: m -separated design is $\Theta(n.m.2^m)$ time and $\Theta(m.n)$ space (e.g. Fixed-Parameter Tractable in m)
- ▶ Corollary: Structures with $h_{min} = 3$ solved in $O(n)!$
- ▶ Bonus: Uniform sampling of sequences

Experimental results at a glance



- Structures with $h_{min} \leq 2 \rightarrow$ Mainly designable in practice (around 100 nucleotides), but exponential decay of numbers of solutions



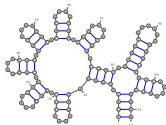
- Sampled solutions from 2-separable structures are more promising according to the Turner energy model than random compatible sequences

Conclusion

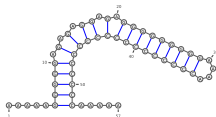
- ▶ $h_{min} = 3 \rightarrow$ IF solved in **linear time** with 2-separable
- ▶ $h_{min} \leq 2 \rightarrow$ IF partially solved, **FPT in m** with m -separable

Ongoing work (Return to the REAL world¹)

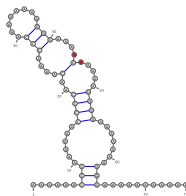
- ▶ Solutions as seed sequences of **RNA_{INVERSE}**
- ▶ Test on Eterna v2 benchmark (100 **artificial** puzzles)
→ almost everyone is **2-separable**
- ▶ One puzzle has a multiloop of degree 25
→ still linear, just with a large constant
- ▶ solved ~80 puzzles, (almost) immediate solutions in half of them



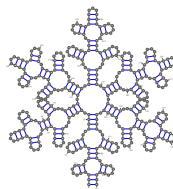
Puzzle #58



Puzzle #27



Puzzle #75



Puzzle #76

¹where Hua-Ting tries to steal an authorship

Thanks to...



Yann Ponty



Search co-chercheur de:
CNRS
ÉCOLE DES PONTS PARISTECH
UNIVERSITÉ GUSTAVE EIFFEL



Laurent Bulteau



Collaborator: Hua Ting Yao (that did the last slide)

And to the other members of the AMIBio team:

Sarah Berkemer

Sebastian Will

Alan Azede

Nan Pan

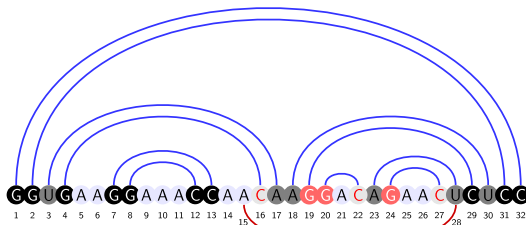
Link to the paper



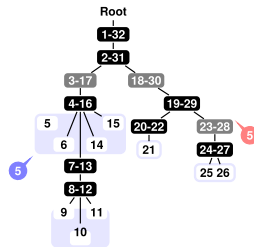
You! ENS DE LYON



Tree formalism



Structure representation



Tree representation

Definition (Levels): Given a tree coloring, the level $L : V(T) \rightarrow \mathbb{Z}$ of a node v is $L(v) := |p|_{\bullet} - |p|_{\circ}$ where p denotes the color vector associated with the node sequence from $\text{parent}(v)$ to Root.

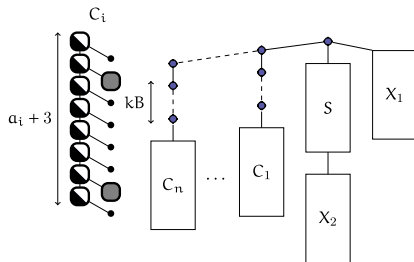
Decide separability is NP-complete

Problem 1 (INTERVAL PACKING):

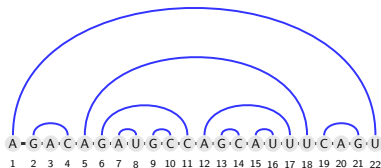
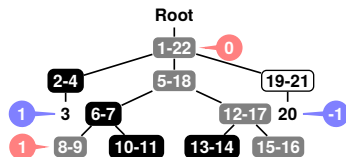
Input: set of distinct integers $A = \{a_1, \dots, a_n\}$, integers k and B

Output: function x from A to intervals of $\llbracket 0, kB - 1 \rrbracket$ such that:

- ▶ $x(a_i)$ is an interval of size a_i
- ▶ $x(a_i)$ and $x(a_j)$ are disjoint for $i \neq j$
- ▶ $x(a_i)$ does not contain both $jB - 1$ and jB for any i, j .

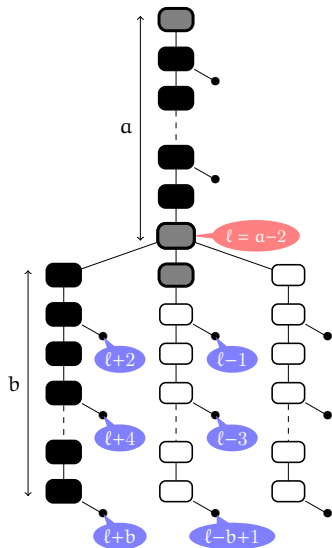
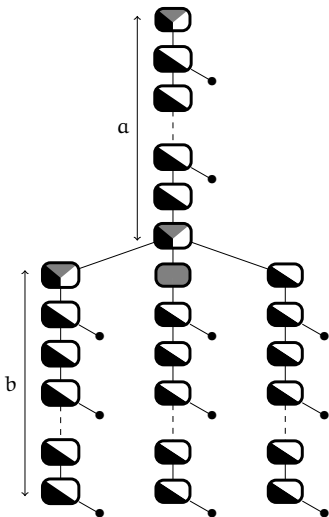


Which instances are non-separable but designable?



- Harder to find with helices of size 2 or more. (currently more than 1000 nucleotides long)

Core widget of the designable non-separable instances with helices of size 2



Modulo-separability

Definition ((Modulo) m -separability): Let m be an integer. A coloring $Color$ is m -separated (or separated with modulus m) for a target secondary structure T , if and only if

$$\{Lv(v) \bmod m \mid Color(v) = \bullet\} \cap \{Lv(v) \bmod m \mid v \text{ is a leaf}\} = \emptyset$$

- Modulo separability coincides with separability with $m \geq \frac{n}{2}$

Problem 2 (MODULO SEPARABILITY):

Input: A tree T (with no m_3 or m_5 motif), a modulus $m \in \mathbb{N}$

Output: A coloring of T that is m -separated, or \perp if no such coloring exists.

Dynamic programming scheme for modulo separability

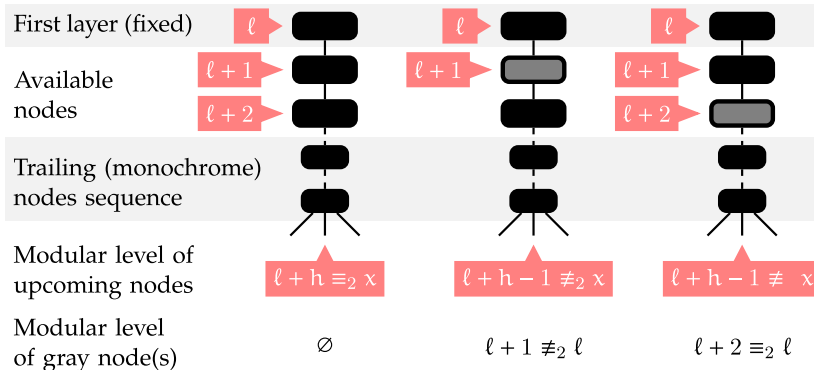
$$d_{v \rightarrow c, \ell}^{\xi_L} = \begin{cases} \text{False} & \text{if } \ell \in \xi_L \wedge c = \bullet \\ & \text{or } \ell' \notin \xi_L, \text{ and} \\ & \exists \text{ leaf in children}(v) \\ \text{True} & \text{if children}(v) = \emptyset \\ \bigvee_{\substack{c' \text{ "valid" coloring of} \\ \text{children}(v) \text{ given } v \rightarrow c}} \bigwedge_{v' \in \text{children}(v)} d_{v' \rightarrow c'(v'), \ell'}^{\xi_L} & \text{otherwise.} \end{cases}$$

with $\ell' := \ell + \delta(c) \bmod m$

- ▶ $d_{v \rightarrow c, \ell}^{\xi_L}$: existence of a valid assignment for a subtree of T rooted at internal node v , with v occurring at level ℓ , and being assigned a prior color c .
- ▶ ξ_L : Leaves levels (thus $\llbracket 0, m \rrbracket \setminus \xi_L$ are \bullet levels.)
- ▶ δ : level increment induced by a color c

Instances with helices of size 3 or more are all separable

$$x := \ell + h \bmod 2$$



Theorem: Secondary structures with helices of size 3 or more are 2-separable (thus designable) in linear time

Inverse folding: Complexity Zoo

- ▶ **NP-hard**, 2008, Schnall-Levin et al ...

- ▶ **Linear**, 2017, Halès et al ...

But only on a subset called "separable instances".

- ▶ **NP-hard**, 2018, Bonnet et al ...

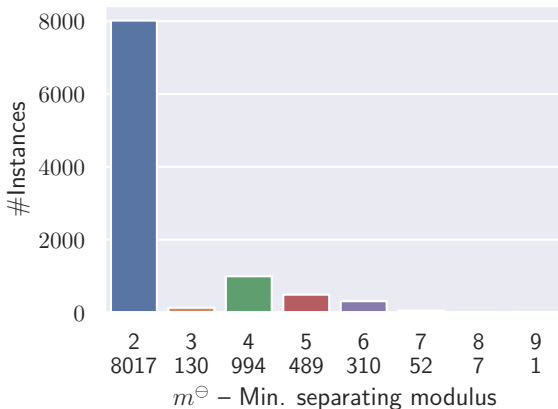
But only an extension with constrained base pairs.

- ▶ Our contribution:

Linear by avoiding isolated base pairs and stacks, 2024, Boury et al.

Beyond helices of size 3: instances with helices of size 2

There is no certainty that these instances are Modulo m -separable!



- Surprisingly enough, all instances containing helices of size 2 were found Modulo m -separable **thus designable**.

Turner energy of designed sequence with helices of size 3

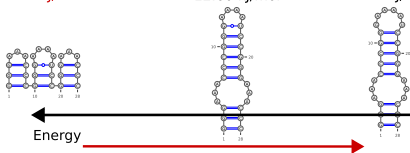
$$\Delta\Delta G(\omega, T) := \Delta G(\omega, \alpha(\omega, T)) - \Delta G(\omega, T)$$

$$\alpha(\omega, T) := \min\{\Delta G(\omega, T') \mid |T' \triangle T| \geq 3\}$$

Target structure T
-4.90 kJ/mol

First alternative structure
-12.60 kJ/mol

Turner structure
-13.0 kJ/mol

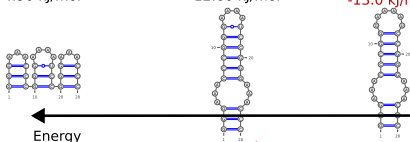


$$\Delta\Delta G(\omega, T) = -8.10 \text{ kJ/mol}$$

Alternative structure
-4.90 kJ/mol

First alternative structure
-12.60 kJ/mol

Target structure T
= Turner structure
-13.0 kJ/mol

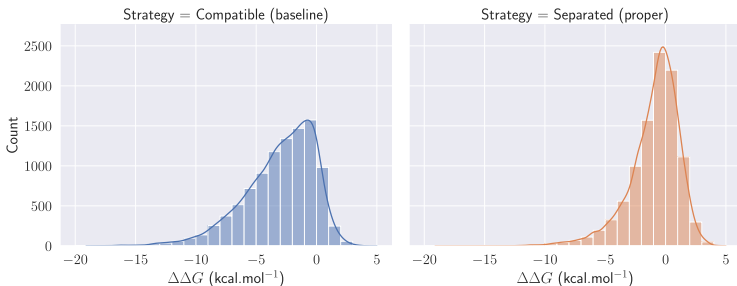


$$\Delta\Delta G(\omega, T) = +0.4 \text{ kJ/mol}$$

Turner energy of designed sequence with helices of size 3

$$\Delta\Delta G(\omega, T) := \Delta G(\omega, C(\omega, T)) - \Delta G(\omega, T)$$

$$C(\omega, T) := \min\{\Delta G(\omega, T') \mid |T' \triangle T| \geq 3\}$$



- Even if guaranteed only in a base pairs model, our sequences represent **better competitor in Turner energy model** than simply compatible sequences