

# Algorithmics results for RNA design with sequence constraints

Vincent LE GALLIC

Alain DENISE

Yann PONTY

LRI

25/11/2015

- 1 Introduction
- 2 Problem formalisation
- 3 Zhou et al. approach
- 4 Contributions
  - Complexity of generation with mandatory patterns
  - Dynamic programming to enforce one occurrence of a pattern
  - Speed enhancement
- 5 Open questions

# RNA structure

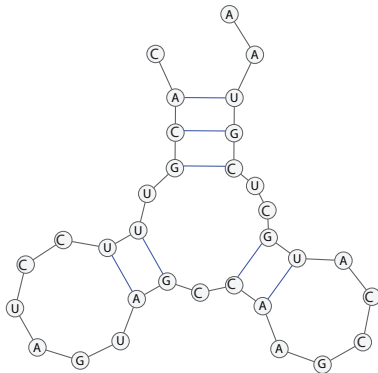
RNA = sequence of  
A, U, G, C

secondary structure =  
hydrogen bonds

A – U, G – C

(Watson-Crick)

G – U (Wobble)



# RNA structure

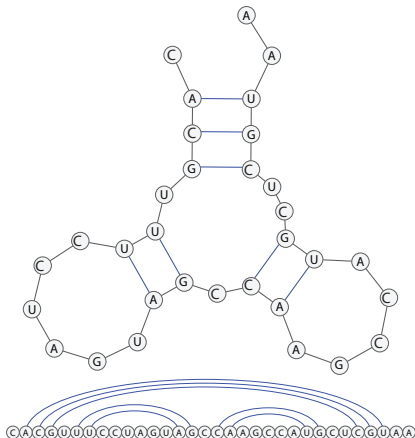
RNA = sequence of  
A, U, G, C

secondary structure =  
hydrogen bonds

A – U, G – C

(Watson-Crick)

G – U (Wobble)



# RNA structure

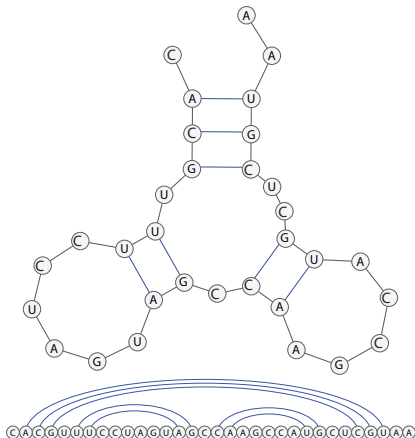
RNA = sequence of  
A, U, G, C

secondary structure =  
hydrogen bonds

A – U, G – C

(Watson-Crick)

G – U (Wobble)



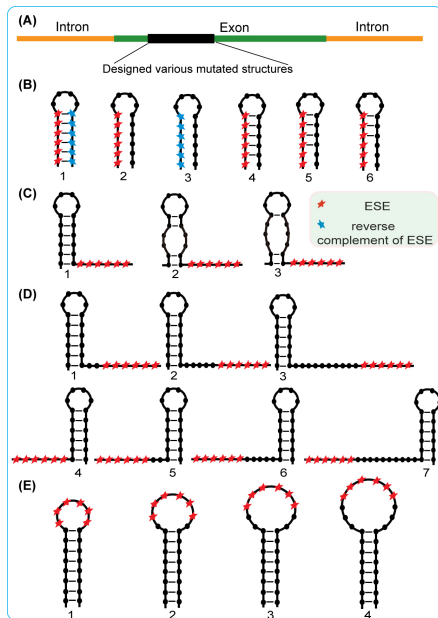
Energy model (Turner) : sequence  $w$ , folded in structure  $S$   
 $\rightarrow E(w, S)$

**Hypothesis** :  $w$  folds into  $S^*$  minimizing  $E(w, S)$

# Biological motivation : Exon Splicing Enhancer

In mRNA : binding site, target for splicing

Q: Influence of structural context?



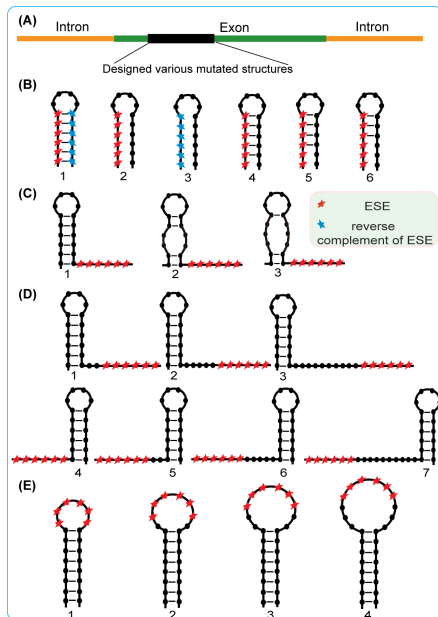
# Biological motivation : Exon Splicing Enhancer

In mRNA : binding site, target for splicing

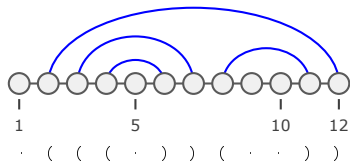
Q: Influence of structural context?

given structure  $S$   
 $ESE_i = AGAACU$   
required at specific position

$ESE_{j \neq i}$   
forbidden



one structure = a well-parenthesized word on  $\{ (, ), \bullet, \}$





$E(w, S)$  = free energy of  $w$  folded into  $S$

## Problem : **Folding**

**Data:** Sequence  $w$

**Result:**  $\text{MFE}(w)$ , set of structures minimizing  $E(w, S)$

## Problem : **Inverse folding** or **Negative design**

**Data:** Secondary structure  $S$

**Result:** Sequence  $w$  such that  $\text{MFE}(s) = \{S\}$  if such a sequence exists,  $\emptyset$  otherwise.

Unknown complexity (NP-hard?).

- Need to include constraints :
  - **Forbidden patterns.** Example : "ESE<sub>*j*</sub> mustn't appear"
  - **Mandatory patterns.** Example : "CAAU must appear at least once"
  - **Positional constraints.** Example : "At position 27 only A or G allowed" or "positions 3-9 ESE<sub>*j*</sub>"
- Other objective : finding sequences which folding is "near" to  $S$

## Problem : Design with constraints

### Data:

- structure  $S$  ( $|S| = n$ )
- set of forbidden patterns  $\mathcal{F}$
- set of mandatory patterns  $\mathcal{M}$
- positional constraints  $PC = \{(i, C_i) / i \in \llbracket 1, n \rrbracket, C_i \subseteq \Sigma\}$

### Result: Sequence $w$ such that

- $w$  is compatible with  $S$
- $\forall f \in \mathcal{F}, f \not\preceq w$
- $\forall m \in \mathcal{M}, m \preceq w$
- $\forall i \in \llbracket 1, n \rrbracket, w_i \in C_i$

or  $\emptyset$  if no sequence fulfill these constraints

( $u \preceq v$  means  $u$  factor of  $v$ )

## Stochastic search :

- RNAInverse

## ... + divide and conquer :

- RNA-SSD
- INFO-RNA
- NUPack

## Random generation :

- RNAesign
- IncaRNAtion

## Genetic algorithms :

- FRNAKenstein
- RNAExInv

## Exact algorithms :

- RNAiFold
- CO4

## Stochastic search :

- RNAInverse

## ... + divide and conquer :

- RNA-SSD
- INFO-RNA
- NUPack

## Random generation :

- RNAensign
- IncaRNAtion

## Genetic algorithms :

- FRNAKenstein
- RNAExInv

## Exact algorithms :

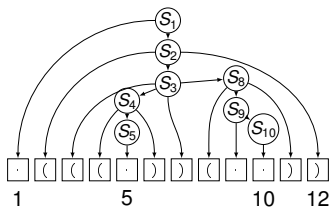
- RNAiFold
- CO4

Only one handles forbidden patterns : NUPack

Zhou *et al.* (2013) introduce CFGRNAD :

- Define  $\mathcal{L}$  = sequences compatible with  $S$  and the constraints
- Random generation in  $\mathcal{L}$  (dynamic programming)
  - counting
  - generating
- Fold each candidate to check if it folds into  $S$

Goal : optimize the counting and generating steps of this method



$$S_1 \rightarrow \bullet S_2$$

$$S_4 \rightarrow (S_5)$$

$$S_8 \rightarrow (S_9)$$

$$S_2 \rightarrow (S_3)$$

$$S_5 \rightarrow \bullet$$

$$S_9 \rightarrow \bullet S_{10}$$

$$S_3 \rightarrow (S_4)S_8$$

$$S_{10} \rightarrow \bullet$$

Full development of  $\mathcal{G}_S$  :

•  $\longrightarrow$  base A, U, G, C

(...)  $\longrightarrow$  pair {A, U}, {G, C}, {G, U}

$$S_1 \rightarrow aS_2 \mid uS_2 \mid gS_2 \mid cS_2$$

$$S_2 \rightarrow aS_3u \mid uS_3a \mid gS_3c \mid cS_3g \mid gS_3u \mid uS_3g$$

...

$$S_5 \rightarrow a \mid u \mid g \mid c$$

...

Recognized language = all sequences compatible with  $S$



Full development of  $\mathcal{G}_S$  :

•  $\rightarrow$  base A, U, G, C

(...)  $\rightarrow$  pair {A, U}, {G, C}, {G, U}

$$S_1 \rightarrow aS_2 \mid uS_2 \mid gS_2 \mid cS_2$$

$$S_2 \rightarrow aS_3u \mid \cancel{uS_3a} \mid gS_3c \mid cS_3g \mid gS_3u \mid \cancel{uS_3g}$$

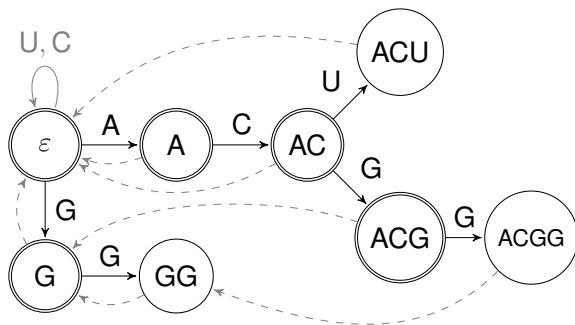
...

$$S_5 \rightarrow a \mid u \mid g \mid c$$

...

Recognized language = all sequences compatible with  $S$   
with  $C_2 = \{A, G, C\}$

# Forbidden words : Aho-Corasick automaton



Automaton  $\mathcal{A}_{\mathcal{F}}$  for  $\mathcal{F} = \{ACU, ACGG, GG\}$   
pictured with *failure transitions*

# Mandatory patterns

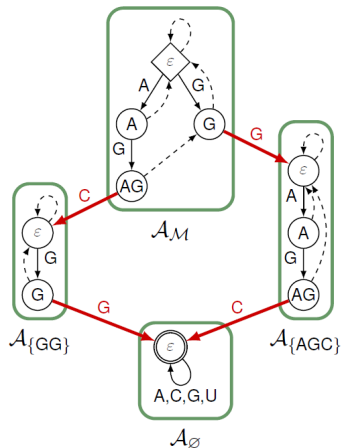


Figure 1: Automaton  $\mathcal{A}_M$  for  $M = \{AGC, GG\}$

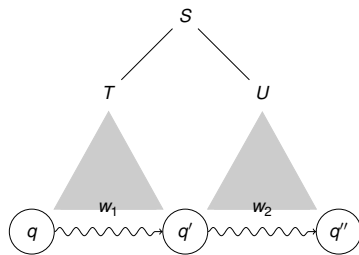
$$\mathcal{A}_{\mathcal{F},\mathcal{M}} = \mathcal{A}_{\mathcal{F}} \times \mathcal{A}_{\mathcal{M}}$$

$$|Q_{\mathcal{F},\mathcal{M}}| \in O\left(2^{|\mathcal{M}|} \left(\sum_{f \in \mathcal{F}} |f| + \sum_{m \in \mathcal{M}} |m|\right)\right)$$

- structural + positional constraints  $\leftrightarrow \mathcal{G}_{\mathcal{S}}$
- forbidden and mandatory patterns  $\leftrightarrow \mathcal{A}_{\mathcal{F},\mathcal{M}}$

context-free  $\cap$  regular = context-free

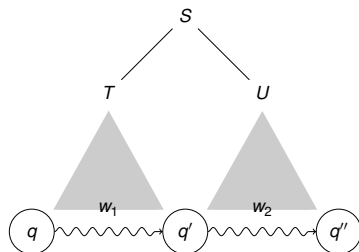
$$\begin{aligned} S &\rightarrow TU \\ &\downarrow \\ S_{q,q''} &\rightarrow T_{q,q'} U_{q',q''} \end{aligned}$$



$$|\mathcal{G}_{\mathcal{G} \times \mathcal{A}}| = |\mathcal{G}| \cdot |\mathcal{A}|^3$$

context-free  $\cap$  regular = context-free

$$\begin{aligned} S &\rightarrow TU \\ \downarrow \\ S_{q,q''} &\rightarrow T_{q,q'} U_{q',q''} \end{aligned}$$



$$|\mathcal{G}_{\mathcal{G} \times \mathcal{A}}| = |\mathcal{G}| \cdot |\mathcal{A}|^3 = n \cdot |\mathcal{Q}_{\mathcal{F}, \mathcal{M}}|^3$$

- 1 Introduction
- 2 Problem formalisation
- 3 Zhou et al. approach
- 4 Contributions**
  - Complexity of generation with mandatory patterns
  - Dynamic programming to enforce one occurrence of a pattern
  - Speed enhancement
- 5 Open questions

Complexity in  $O(2^{|\mathcal{M}|})$

Awaited: it's impossible to do "better"



Complexity in  $O(2^{|\mathcal{M}|})$

Awaited: it's impossible to do "better"

## Superstring Problem

**Data:** Set  $\{u_1, \dots, u_k\}$  of words on  $\Sigma$  and  $\ell \in \mathbb{N}$

**Result :** "True" if it exists  $u$  such that  $\forall i \in \llbracket 1, k \rrbracket, u_i \preceq u$ , et  $|u| \leq \ell$ , "False" otherwise

This problem is NP-complete (cf Maier *et al.* (1977))

# Superstring NP-hard $\Rightarrow$ Design NP-hard

- $S = \bullet \bullet \dots \bullet$  of length  $\ell$
- $\mathcal{F} = \emptyset$
- $\mathcal{M} = \{u_1, \dots, u_k\}$

Then **Superstring**  $\prec$  **pre-Design**

$\rightsquigarrow$  no hope **Design with constraints**  $\in \text{Poly}(|\mathcal{M}|)$

Weakness of Zhou *et al.* approach:

How to enforce a pattern  $u_{mf}$  *only* at position  $i_{u_{mf}}$  and nowhere else?

- use  $\mathcal{G}_S$  ?

$$\mathcal{A}_{\mathcal{F},\mathcal{M}} \text{ forbids } u_{mf} \rightsquigarrow \mathcal{L}_{\mathcal{G}_S} \cap \mathcal{L}_{\mathcal{A}_{\mathcal{F},\mathcal{M}}} = \emptyset$$

- use  $\mathcal{A}_{\mathcal{F},\mathcal{M}}$  ?

$$|\mathcal{Q}_{\mathcal{F},\mathcal{M}}| \in \Omega(n) \rightsquigarrow \text{complexity in } O(n|\mathcal{Q}_{\mathcal{F},\mathcal{M}}|^3) = O(n^4)$$

# Pattern once and only once

Weakness of Zhou *et al.* approach:

How to enforce a pattern  $u_{mf}$  *only* at position  $i_{u_{mf}}$  and nowhere else?

- use  $\mathcal{G}_S$  ?

$$\mathcal{A}_{\mathcal{F},\mathcal{M}} \text{ forbids } u_{mf} \rightsquigarrow \mathcal{L}_{\mathcal{G}_S} \cap \mathcal{L}_{\mathcal{A}_{\mathcal{F},\mathcal{M}}} = \emptyset$$

- use  $\mathcal{A}_{\mathcal{F},\mathcal{M}}$  ?

$$|\mathcal{Q}_{\mathcal{F},\mathcal{M}}| \in \Omega(n) \rightsquigarrow \text{complexity in } O(n|\mathcal{Q}_{\mathcal{F},\mathcal{M}}|^3) = O(n^4)$$

$\implies$  method in  $O(n|\mathcal{Q}_{\mathcal{F},\mathcal{M}}|^3)$  where  $Q$  depends on  $\mathcal{F}$  and  $\mathcal{M}$ , but not  $n$

$C(i, j, q, q'')$  = number of words  $w$  compatible with  $S_{[i,j]}$  such that  $q \xrightarrow{w} q'' \in \delta_{\mathcal{F}, \mathcal{M}}^*$

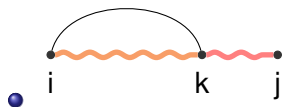
(((( ( ... ) ))))  
          *i*                              *j*  
          *q*                              *q''*

$C(i, j, q, q'')$  needed for random generation

Computing  $C(i, j, q, q'')$  :



- choose  $a \in \Sigma$
- compute  $C(i+1, j, q_{\text{after}}, q'')$  where  $q \xrightarrow{a} q_{\text{after}} \in \delta_{\mathcal{F}, \mathcal{M}}$



- choose  $a, b \in \Sigma$
- choose  $q' \in Q$
- compute  $C(i+1, k-1, q_{a1}, q')$  where  $q \xrightarrow{a} q_{a1} \in \delta_{\mathcal{F}, \mathcal{M}}$
- compute  $C(k+1, j, q_{a2}, q'')$  where  $q' \xrightarrow{b} q_{a2} \in \delta_{\mathcal{F}, \mathcal{M}}$

# Dynamic programming : recurrence relations

- $S_{[i,i-1]} = \varepsilon$

$$C(i, i-1, q, q'') = \begin{cases} 1 & \text{if } q = q'' \\ 0 & \text{otherwise} \end{cases}$$

- $S_{[i,j]} = \bullet S_{[i+j,j]}$

$$C(i, j, q, q'') = \sum_{a \in C_i} \begin{cases} 0 & \text{if } \delta(q, a) \in \mathcal{F} \\ 0 & \text{if } \delta(q, a) = u_{mf} \text{ and } i - (|u_{mf}| - 1) \neq i_{u_{mf}} \\ C(i+1, j, \delta(q, a), q'') & \text{otherwise} \end{cases}$$

- $S_{[i,j]} = (S_{[i+1,k-1]})S_{[k+1,j]}$

$$C(i, j, q, q'') = \sum_{\substack{(a,b) \in C_i \times C_j \\ q' \in Q}} \begin{cases} 0 & \text{if } \delta(q, a) \in \mathcal{F} \text{ ou } \delta(q', b) \in \mathcal{F} \\ 0 & \text{if } \delta(q, a) = u_{mf} \text{ and } i - (|u_{mf}| - 1) \neq i_{u_{mf}} \\ 0 & \text{if } \delta(q', b) = u_{mf} \text{ and } k - (|u_{mf}| - 1) \neq i_{u_{mf}} \\ C(i+1, k-1, \delta(q, a), q') \cdot C(k+1, j, \delta(q', b), q'') & \text{otherwise} \end{cases}$$

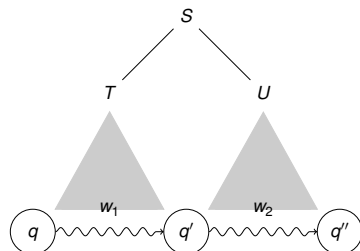
# Speed enhancement

CFG/GRNAD tests

every  $(q, q'') \in Q^2$

every  $q' \in Q$

$\rightsquigarrow O(|Q|^3)$



But:

- accessible/co-accessible ?
- non-terminal  $\rightarrow$  fixed-length words

$\rightsquigarrow$  most  $(q, q', q'')$ , even  $(q, q'')$ , are *irrelevant*

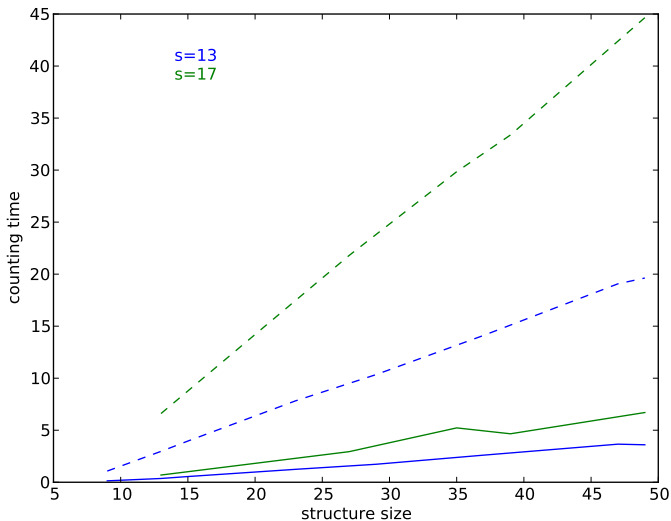


Optimisation : pre-compute triplets  $(q, q', q'')$  such that there is a path  $q \rightarrow^* q' \rightarrow^* q''$  in  $\mathcal{A}_{\mathcal{F}, \mathcal{M}}$ .

$q \rightarrow^{n_1} q' \rightarrow^{n_2} q''$  :

- $q \rightarrow^{n_1+n_2} q'' \in \delta_{\mathcal{F}, \mathcal{M}}^*$
- $q \rightarrow^{n_1} q' \in \delta_{\mathcal{F}, \mathcal{M}}^*$  et  $q' \rightarrow^{n_2} q'' \in \delta_{\mathcal{F}, \mathcal{M}}^*$

# Speed enhancement : preliminary results



Better algorithmic complexity ?

- $\mathcal{A}_{\mathcal{F}, \mathcal{M}}$  optimal *in general*?  
(minimizing is costly, problems with prog. dyn.)
- folding of  $w$  into  $S$  not guaranteed  $\rightsquigarrow$  enforce a *bias* in the random generation
- characterize families of automata with whom complexity in  $O(n \cdot |Q|^p)$ ,  $p < 3$

Thank you for your attention