

Design d'ARN et génération aléatoire

Bonnie Berger* Alain Denise♣ Srinivas Devadas*
Alex Levin* Mieszko Lis* Yann Ponty• Charles O'Donnell*
Vladimir Reinharz‡ Stéphane Vialette† Jérôme Waldispühl‡
Yi Zhang◇ Yu Zhou♣,◇,⊙

*CSail, MIT, USA

♣LRI, Univ. Paris-Sud, Orsay F-91405, France

◇State Key Laboratory of Virology, Wuhan University, China

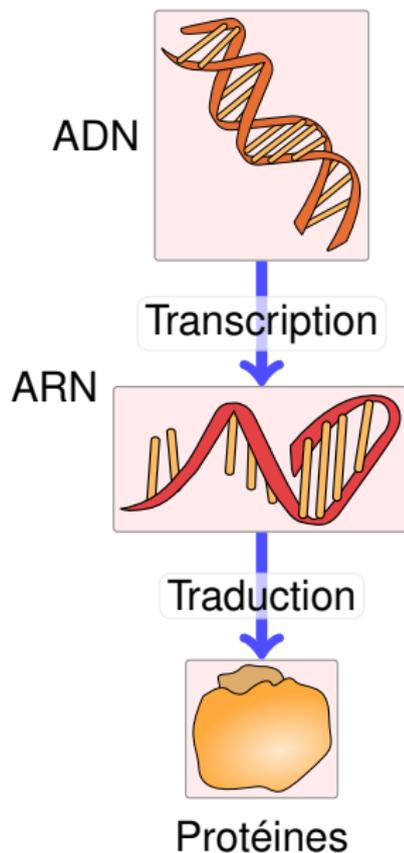
⊙ Cellular and Molecular Medicine, Univ. of California San Diego, USA

• LIX, CNRS/Ecole Polytechnique, France + PIMS, Vancouver, Canada

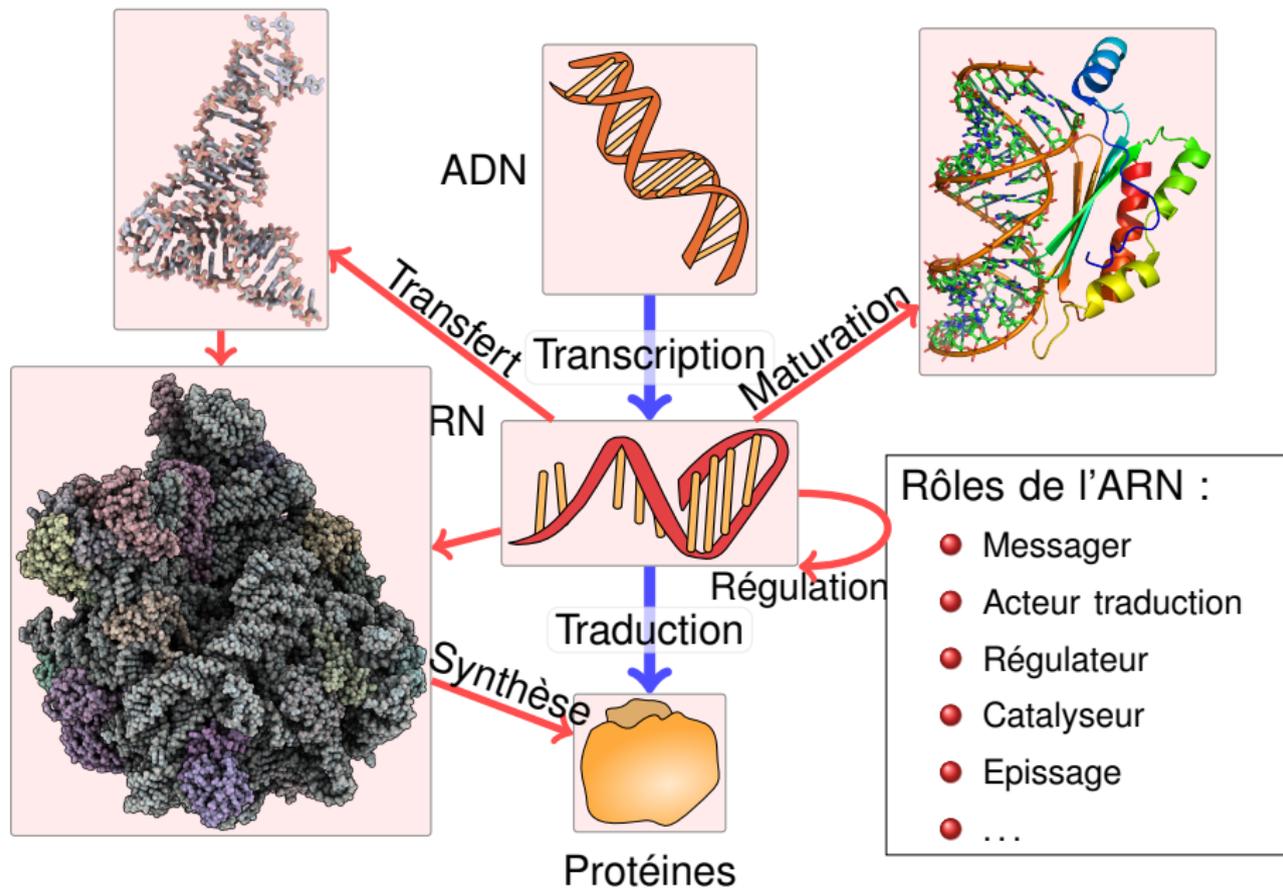
† LIGM, Université Paris-Est, France

‡ Computer Science, Mc Gill University, Canada

Dogme fondamental de la biologie moléculaire



Dogme fondamental de la biologie moléculaire



Motivation d'un design contraint d'ARN

Motivation : Impact de la structure sur les stimulateurs de l'épissage (ESE).

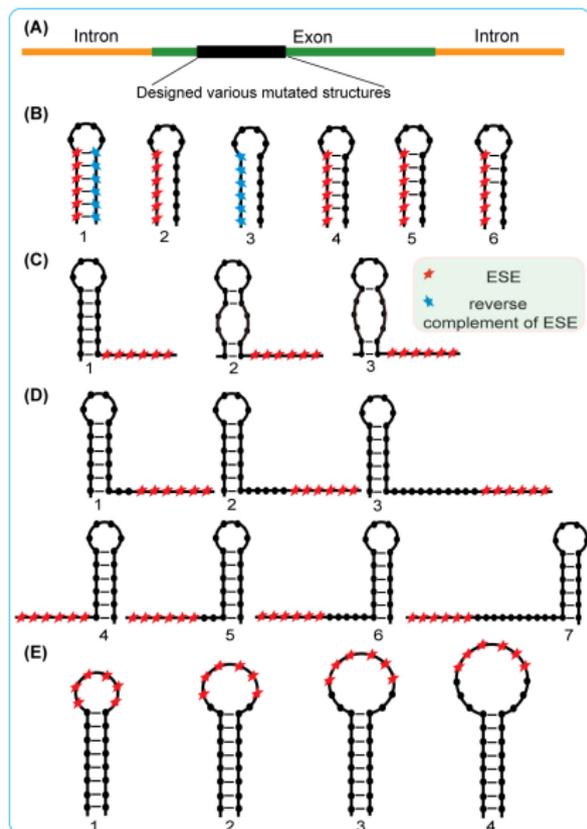
- Structures différentes ;
- Présence localisée d'ESE ;
- Eviter autres ESE (~1500 !)

Validation d'un effet du contexte structurel sur l'efficacité des ESE.

[Liu *et al*, FEBS Lett. 2010]

Objectif : Concevoir séq. d'ARN

- 1 Se repliant correctement
- 2 Présentant/évitant des motifs



I. Repliement inverse

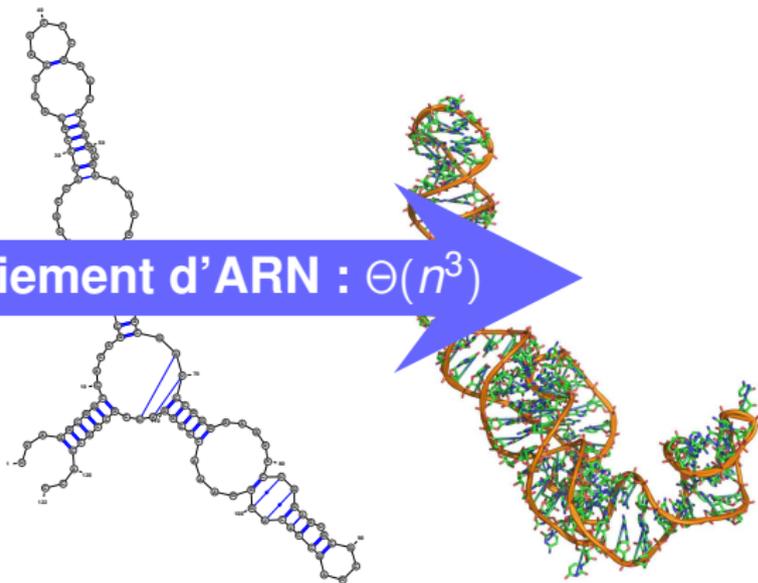
Garantir la structure

Séquence et structure(s) de l'ARN

ARN = Polymère linéaire = Séquence sur $\{A, C, G, U\}^*$

```
UUAGGCGGCCACAGC
GGUGGGGU
CGUACCCAUCCCGAA
CACGGAAGAUAAAGCC
CACCAGCGUUCGGG
GAGUACUGGAGUGCG
CGAGCCUCUGGGAAA
CCCCGUUCGCCCA
CC
```

Repliement d'ARN : $\Theta(n^3)$

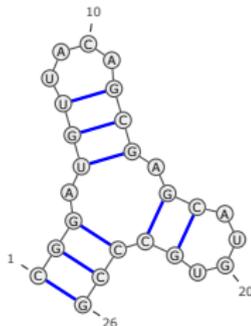
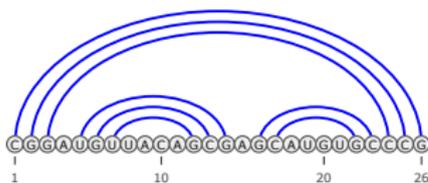


Structure Primaire

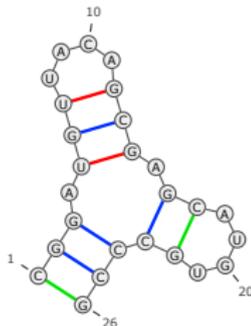
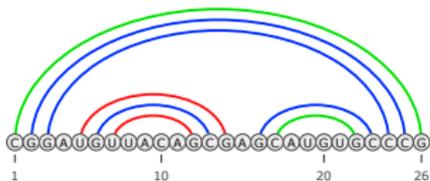
Structure Secondaire

Structure Tertiaire

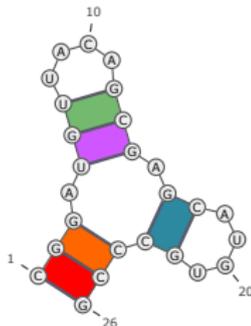
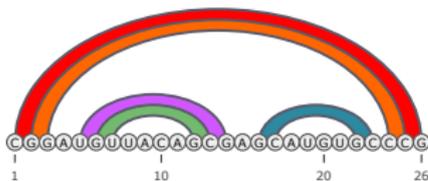
5s rRNA (PDBID : 1K73 :B)



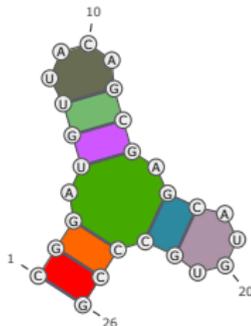
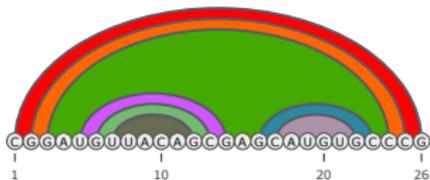
- **Structure d'ARN S : Couplages sans croisement** d'un ARN w
- **Motifs** séquence/structure (ex. Paires, Empilements, Boucles ...)
- **Modèle d'énergie :**
 - Motif** \rightarrow Contribution à l'énergie libre $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
 - Energie libre $E_w(S)$:** Somme Δ sur les motifs de S



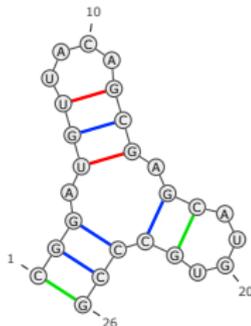
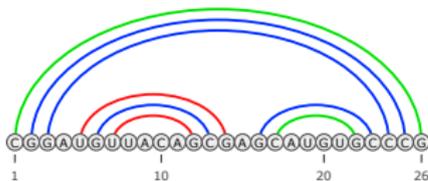
- **Structure d'ARN S** : Couplages **sans croisement** d'un ARN w
- **Motifs** séquence/structure (ex. Paires, Empilements, Boucles ...)
- **Modèle d'énergie** :
 Motif \rightarrow Contribution à l'énergie libre $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
 Energie libre $E_w(S)$: Somme Δ sur les motifs de S



- **Structure d'ARN S** : Couplages **sans croisement** d'un ARN w
- **Motifs** séquence/structure (ex. Paires, Empilements, Boucles ...)
- **Modèle d'énergie** :
 Motif \rightarrow Contribution à l'énergie libre $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
 Energie libre $E_w(S)$: Somme Δ sur les motifs de S

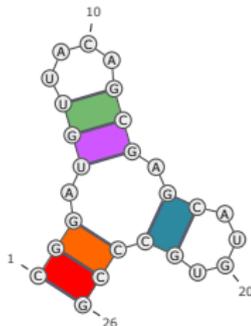
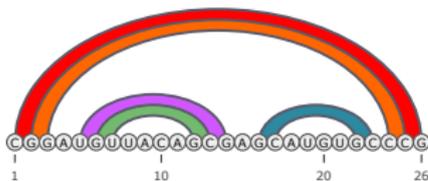


- **Structure d'ARN S** : Couplages **sans croisement** d'un ARN w
- **Motifs** séquence/structure (ex. Paires, Empilements, Boucles ...)
- **Modèle d'énergie** :
 - Motif** \rightarrow Contribution à l'énergie libre $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
 - Energie libre $E_w(S)$** : Somme Δ sur les motifs de S



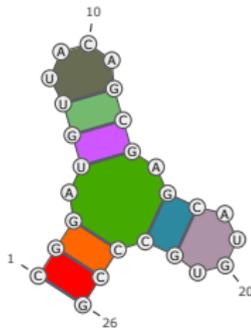
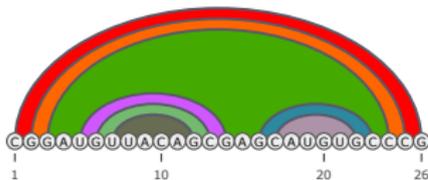
- **Structure d'ARN S** : Couplages **sans croisement** d'un ARN w
- **Motifs** séquence/structure (ex. Paires, Empilements, Boucles ...)
- **Modèle d'énergie** :
Motif \rightarrow Contribution à l'énergie libre $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
Energie libre $E_w(S)$: Somme Δ sur les motifs de S

$$E_S = 2 \cdot \Delta \left(\begin{array}{c} \textcircled{U} \\ | \\ \textcircled{G} \end{array} \right) + 4 \cdot \Delta \left(\begin{array}{c} \textcircled{G} \\ | \\ \textcircled{C} \end{array} \right) + 2 \cdot \Delta \left(\begin{array}{c} \textcircled{C} \\ | \\ \textcircled{G} \end{array} \right)$$



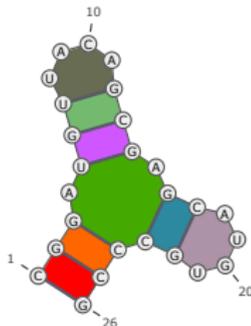
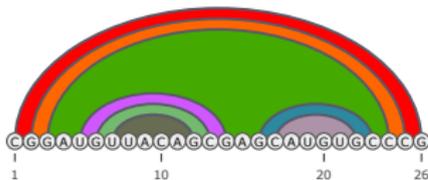
- **Structure d'ARN S** : Couplages **sans croisement** d'un ARN w
- **Motifs** séquence/structure (ex. Paires, Empilements, Boucles ...)
- **Modèle d'énergie** :
Motif \rightarrow Contribution à l'énergie libre $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
Energie libre $E_w(S)$: Somme Δ sur les motifs de S

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



- **Structure d'ARN S** : Couplages **sans croisement** d'un ARN w
- **Motifs** séquence/structure (ex. Paires, Empilements, Boucles ...)
- **Modèle d'énergie** :
 - Motif** \rightarrow Contribution à l'énergie libre $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
 - Energie libre $E_w(S)$** : Somme Δ sur les motifs de S

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



- **Structure d'ARN S** : Couplages **sans croisement** d'un ARN w
- **Motifs** séquence/structure (ex. Paires, Empilements, Boucles ...)
- **Modèle d'énergie** :
 - Motif** \rightarrow Contribution à l'énergie libre $\Delta(\cdot) \in \mathbb{R}^- \cup \{+\infty\}$
 - Energie libre $E_w(S)$** : Somme Δ sur les motifs de S

Definition (Problème REPLIEMENT(E))

Entrée : Séquence d'ARN $w \in \{A, C, G, U\}^*$.

Sortie : Couplage **sans croisement** $S^* = \operatorname{argmin}_S E_w(S)$

\Rightarrow Algorithme(s) en $\Theta(n^3)$ par **programmation dynamique**



Récurrance sur l'énergie minimale d'un repliement :

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

$$N_{i,j} = \min \begin{cases} N_{i+1,j} & (i \text{ non apparié}) \\ \min_{k=i+\theta+1}^j E_{i,k} + N_{i+1,k-1} + N_{k+1,j} & (i \text{ comp. avec } k) \end{cases}$$

Schémas similaires pour des modèles d'énergie plus complexes.

Definition (Problème REPLIEMENT(E))

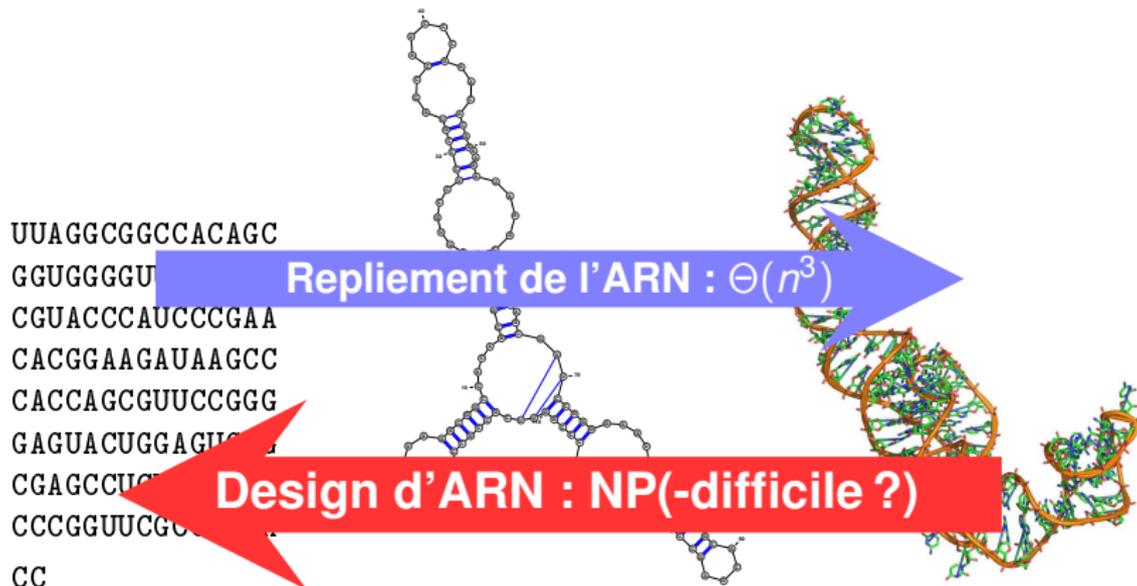
Entrée : Séquence d'ARN $w \in \{A, C, G, U\}^*$.

Sortie : Couplage *sans croisement* $S^* = \operatorname{argmin}_S E_w(S)$

\Rightarrow Algorithme(s) en $\Theta(n^3)$ par **programmation dynamique**

Séquence et structure(s) de l'ARN

ARN = Polymère linéaire = Séquence sur $\{A, C, G, U\}^*$



Structure Primaire

Structure Secondaire

Structure Tertiaire

5s rRNA (PDBID : 1K73 :B)

Repliement inverse d'ARN

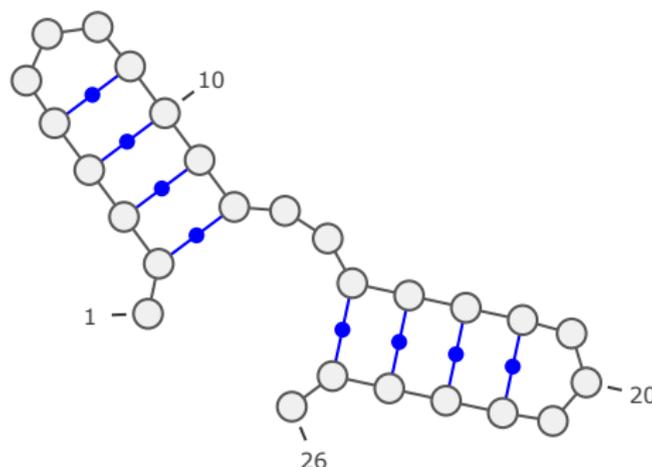
Definition (Problème $\text{DESIGN}(E)$)

Entrée : Couplage **sans croisement** $S \subseteq \mathcal{P}([1, n]^2)$.

Sortie : Séquence w^* telle que $S = \text{argmax}_{S'} E_{w^*}(S')$

Difficile de **diviser pour régner** ...

Exemple :



Repliement inverse d'ARN

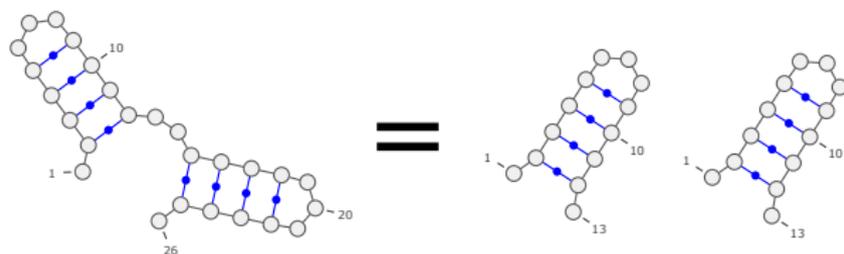
Definition (Problème $\text{DESIGN}(E)$)

Entrée : Couplage **sans croisement** $S \subseteq \mathcal{P}([1, n]^2)$.

Sortie : Séquence w^* telle que $S = \text{argmax}_{S'} E_{w^*}(S')$

Difficile de **diviser pour régner** ...

Exemple :



Repliement inverse d'ARN

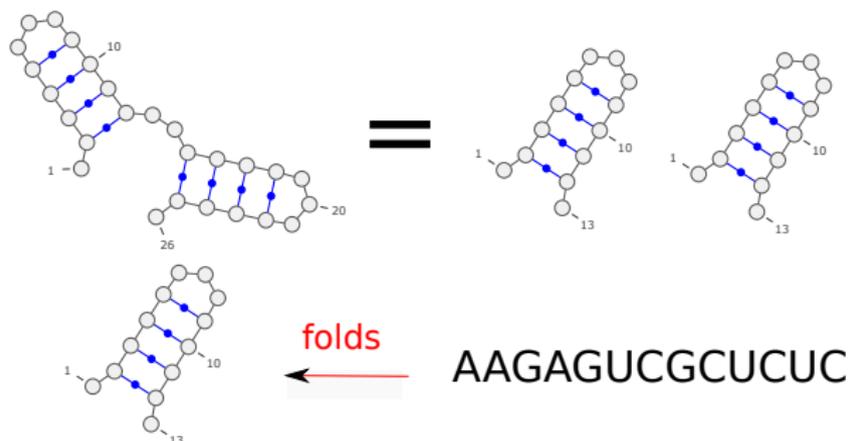
Definition (Problème $\text{DESIGN}(E)$)

Entrée : Couplage **sans croisement** $S \subseteq \mathcal{P}([1, n]^2)$.

Sortie : Séquence w^* telle que $S = \text{argmax}_{S'} E_{w^*}(S')$

Difficile de **diviser pour régner** ...

Exemple :



Repliement inverse d'ARN

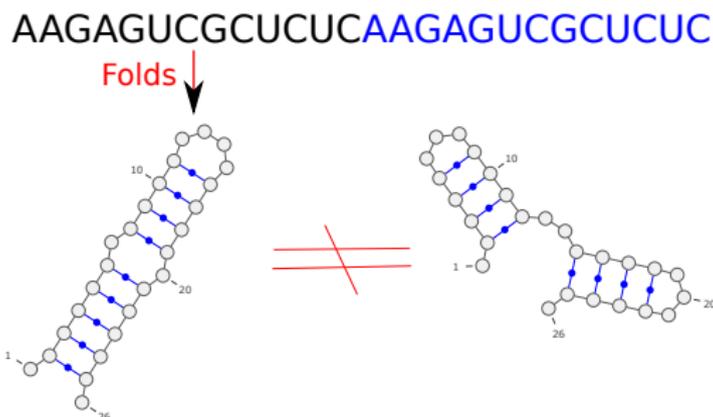
Definition (Problème $\text{DESIGN}(E)$)

Entrée : Couplage **sans croisement** $S \subseteq \mathcal{P}([1, n]^2)$.

Sortie : Séquence w^* telle que $S = \text{argmax}_{S'} E_{w^*}(S')$

Difficile de **diviser pour régner** ...

Exemple :



Repliement inverse d'ARN

Definition (Problème $\text{DESIGN}(E)$)

Entrée : Couplage **sans croisement** $S \subseteq \mathcal{P}([1, n]^2)$.

Sortie : Séquence w^* telle que $S = \text{argmax}_{S'} E_{w^*}(S')$

Difficile de **diviser pour régner** ... mais **complexité inconnue**

⇒ De nombreuses algorithmes basés sur :

... une recherche locale... ... des algorithmes génétiques...

● RNAInverse

● RNAFBinv

● Info-RNA

● FRNAKenstein

● RNA-SSD

... ou exponentiels exacts

● NUPack

● RNAIFold

● CO4

Notre approche : Génération aléatoire **pondérée**

Échantillonnage global [Levin *et al*, NAR 12]

- **Distribution de Boltzmann (DB)** basée sur l'**affinité** avec S
- **Engendrer aléatoirement** selon **DB**
- **Replier** les séquences engendrées et **comparer** à la cible

Facteur de Boltzmann :

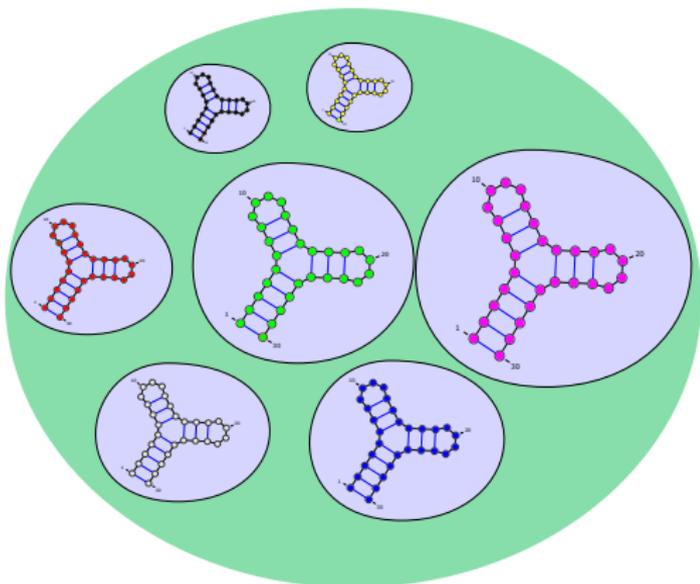
$$\mathcal{B}_w(S) := e^{-\frac{E_w(S)}{RT}}$$

Pseudo-fonction de partition :

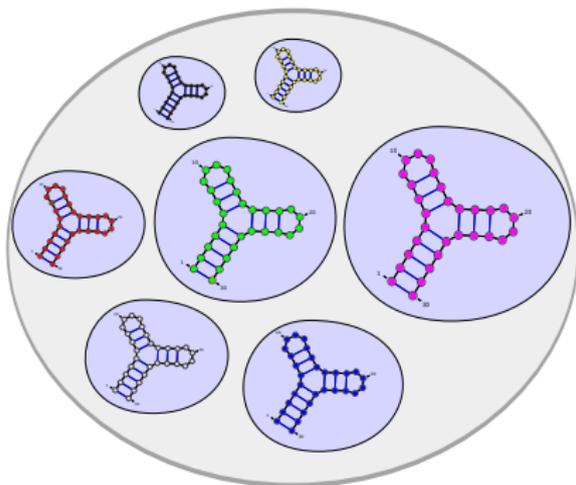
$$\mathcal{Z}(S) = \sum_{w \in \Sigma^*} \mathcal{B}_w(S)$$

Probabilité de Boltzmann :

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



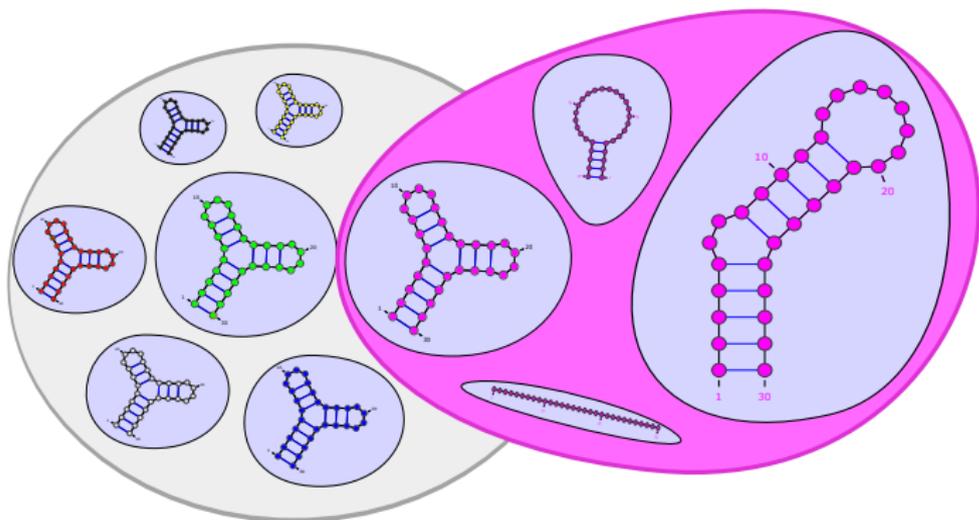
Limites de l'approche



Malheureusement, une forte affinité n'est formellement ni suffisante, ni nécessaire pour être une solution au problème, **mais** ...

- Forte corrélation **empirique** affinité/succès du design [Levin et al 12]
- **Contrôle** de la composition [Bodini P 10] [Reinharz P Waldispühl 10]

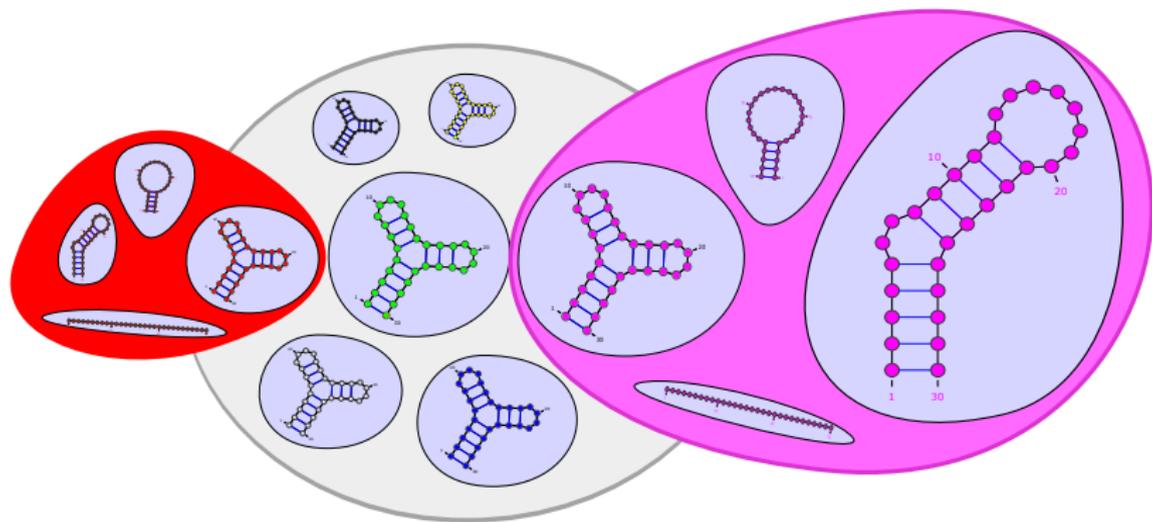
Limites de l'approche



Malheureusement, une forte affinité n'est formellement **ni suffisante**, ni nécessaire pour être une solution au problème, **mais** ...

- Forte corrélation **empirique** affinité/succès du design [Levin et al 12]
- **Contrôle** de la composition [Bodini P 10] [Reinharz P Waldispühl 10]

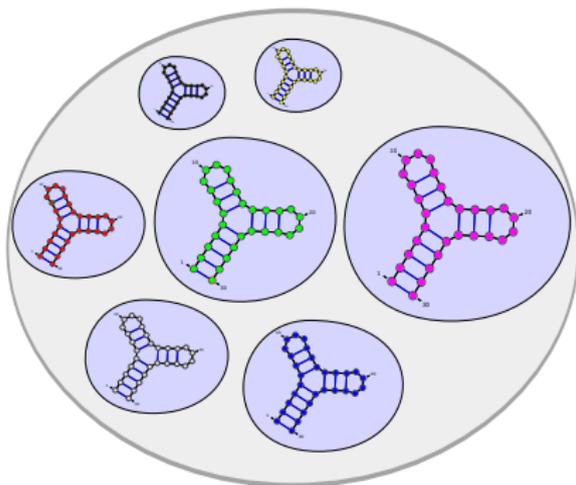
Limites de l'approche



Malheureusement, une forte affinité n'est formellement **ni suffisante**, **ni nécessaire** pour être une solution au problème, **mais** ...

- Forte corrélation **empirique** affinité/succès du design [Levin et al 12]
- **Contrôle** de la composition [Bodini P 10] [Reinharz P Waldispühl 10]

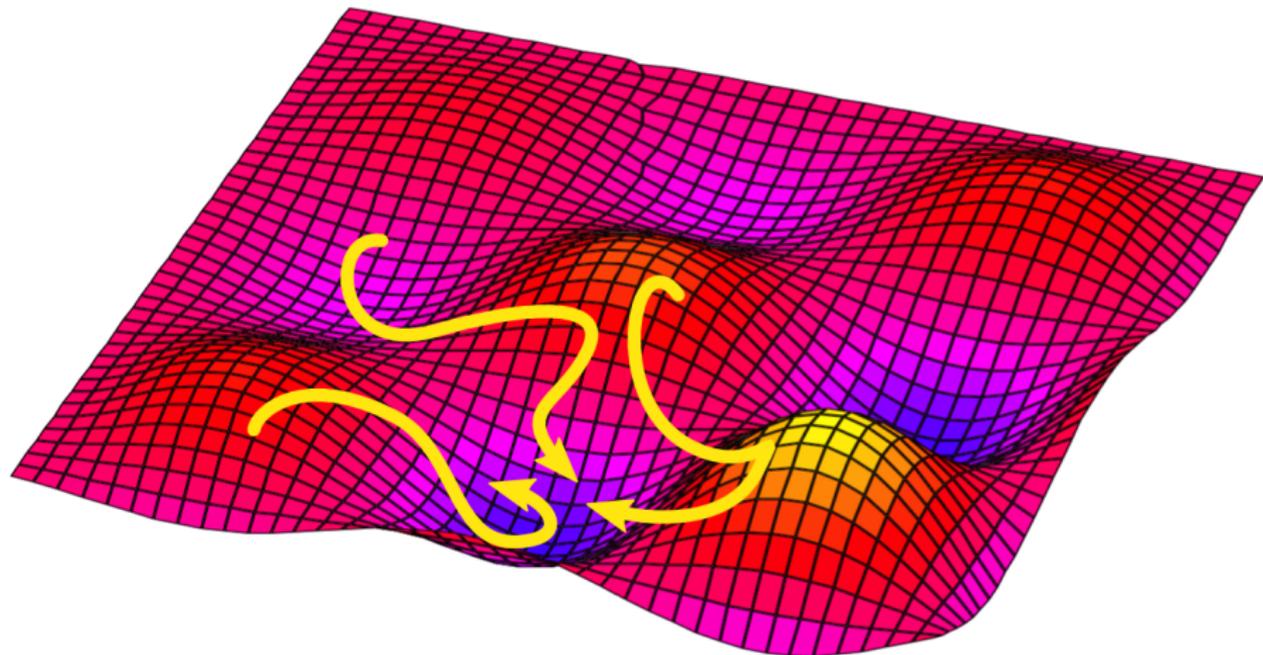
Limites de l'approche



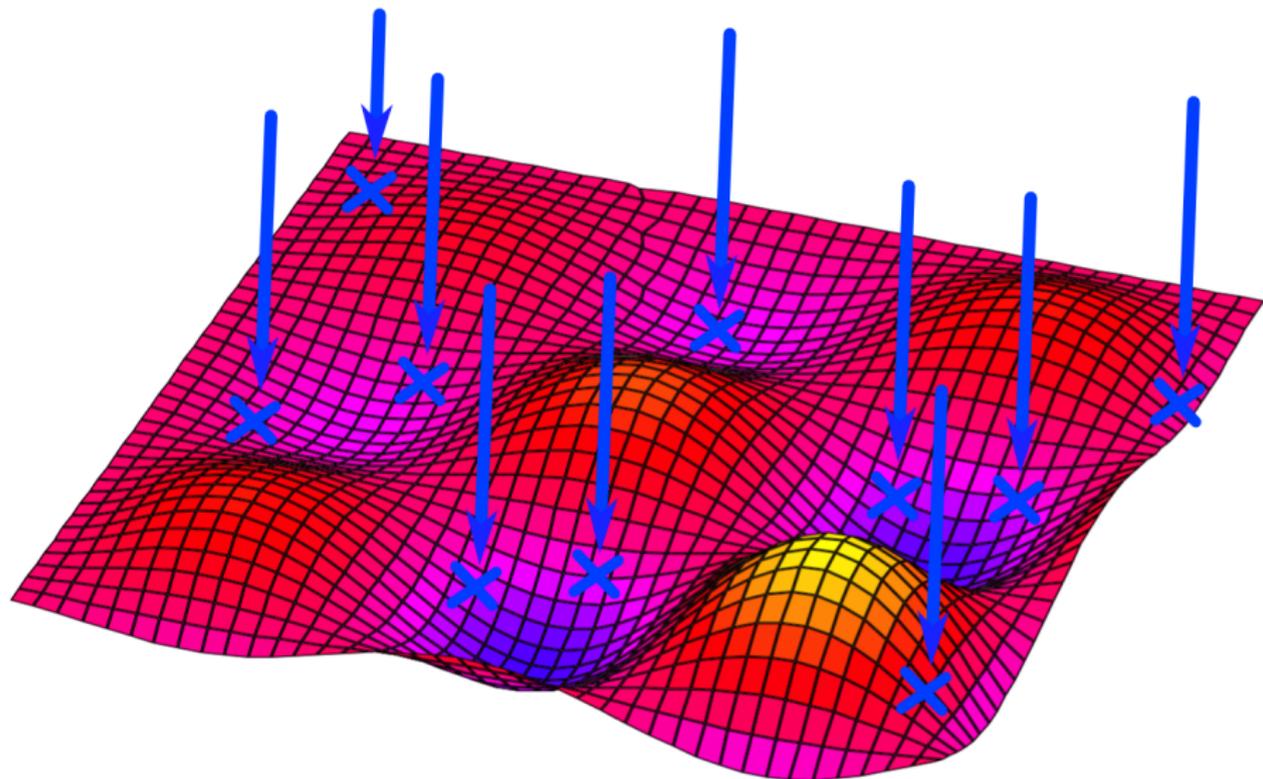
Malheureusement, une forte affinité n'est formellement **ni suffisante**, **ni nécessaire** pour être une solution au problème, **mais** ...

- Forte corrélation **empirique** affinité/succès du design [Levin et al 12]
- **Contrôle** de la composition [Bodini P 10] [Reinharz P Waldispühl 10]

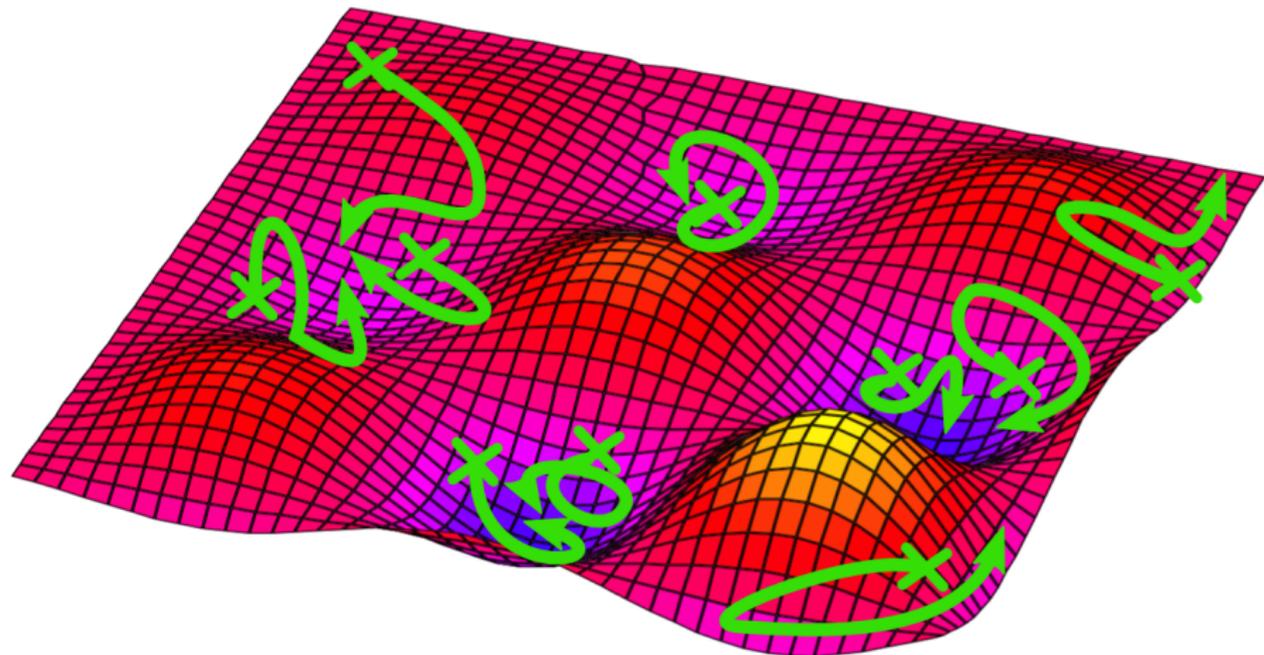
Local vs Global vs "Glocal"



Local vs Global vs "Glocal"



Local vs Global vs "Glocal"



Construire la grammaire (Modèle d'énergie simple)

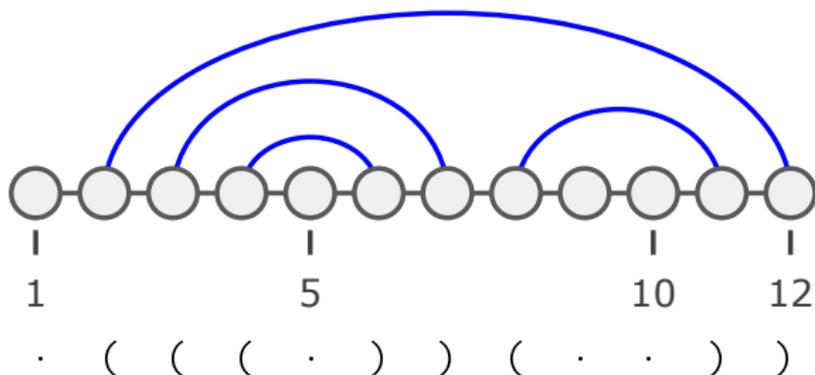
Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (**+ Poids**).

+ Génération Pondérée : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]



Construire la grammaire (Modèle d'énergie simple)

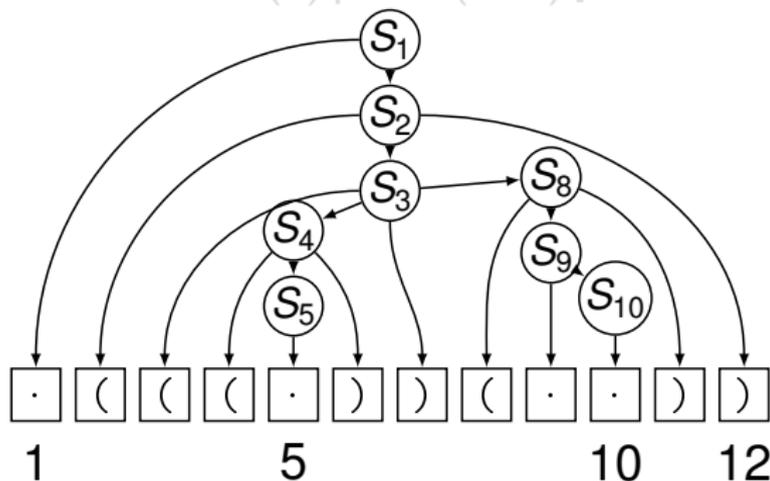
Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (**+ Poids**).

+ **Génération Pondérée** : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]



Construire la grammaire (Modèle d'énergie simple)

Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (+ **Poids**).

+ **Génération Pondérée** : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]

$$\begin{array}{llll} S_1 \rightarrow \cdot S_2 & S_2 \rightarrow (S_3) & S_3 \rightarrow (S_4) S_8 & S_4 \rightarrow (S_5) \\ S_5 \rightarrow \cdot & S_8 \rightarrow (S_9) & S_9 \rightarrow \cdot S_{10} & S_{10} \rightarrow \cdot \end{array}$$

Construire la grammaire (Modèle d'énergie simple)

Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (**+ Poids**).

+ Génération Pondérée : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]

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

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

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

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

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

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

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

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

Construire la grammaire (Modèle d'énergie simple)

Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (**+ Poids**).

+ Génération Pondérée : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]

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

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

Réurrences + Génération

$$v_2(n) := v_3(n-2) + v_3(n-2) + v_3(n-2) + v_3(n-2) + v_3(n-2) + v_3(n-2)$$

$$v_3(n) := \dots$$

Construire la grammaire (Modèle d'énergie simple)

Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (**+ Poids**).

+ Génération Pondérée : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]

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

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

Récurrances + Génération

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

$$V_2 := V_3 + V_3 + V_3 + V_3 + V_3 + V_3$$

$$V_3 := \dots$$

Construire la grammaire (Modèle d'énergie simple)

Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (**+ Poids**).

+ Génération Pondérée : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]

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

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

Récurrentes + Génération

$$V_2 := V_3 \cdot e^{-\Delta_{AU}/RT} + V_3 \cdot e^{-\Delta_{CG}/RT} + V_3 \cdot e^{-\Delta_{GC}/RT} + V_3 \cdot e^{-\Delta_{GU}/RT} \\ + V_3 \cdot e^{-\Delta_{UA}/RT} + V_3 \cdot e^{-\Delta_{UG}/RT}$$

$$V_3 := \dots$$

+ Pondérations

+ Composition contrôlée (Boltzmann multidim. [Bodini P 10])

Rem. : Grammaires réalistes plus complexes (effets *Markoviens*) ...

Construire la grammaire (Modèle d'énergie simple)

Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (**+ Poids**).

+ Génération Pondérée : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]

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

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

Récurr $p = \frac{v_3 \cdot e^{-\Delta_{CG}/RT}}{v_2}$ **ération**

$$v_2 := v_3 \cdot e^{-\Delta_{AU}/RT} + v_3 \cdot e^{-\Delta_{CG}/RT} + v_3 \cdot e^{-\Delta_{GC}/RT} + v_3 \cdot e^{-\Delta_{GU}/RT} \\ + v_3 \cdot e^{-\Delta_{UA}/RT} + v_3 \cdot e^{-\Delta_{UG}/RT}$$

$$v_3 := \dots \mid C V_{10} \mid G V_{10} \mid U V_{10}$$

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

+ Pondérations

+ Composition contrôlée (Boltzmann multidim. [Bodini P 10])

Rem. : Grammaires réalistes plus complexes (effets *Markoviens*) ...

Construire la grammaire (Modèle d'énergie simple)

Entrée : Structure secondaire S

A Construire l'arbre de syntaxe abstraite (ASA) de S ;

B Traduire l'ASA en une grammaire ;

C Étendre la grammaire pour engendrer les lettres (**+ Poids**).

+ Génération Pondérée : $\Theta^*(n)$ puis $\Theta(k \cdot n)$ [Denise P Termier 10]

$$V_1 \rightarrow A V_2 | C V_1 | G V_1 | U V_2$$
$$V_2 \rightarrow A V_3 | C | G V_3 U | U V_3 A | U V_3 G$$

$p = \frac{v_3 \cdot e^{-\Delta_{GC}/RT}}{v_2}$

Réurrences + Génération

$$v_2 := v_3 \cdot e^{-\Delta_{AU}/RT} + v_3 \cdot e^{-\Delta_{CG}/RT} + v_3 \cdot e^{-\Delta_{GC}/RT} + v_3 \cdot e^{-\Delta_{GU}/RT}$$
$$+ v_3 \cdot e^{-\Delta_{UA}/RT} + v_3 \cdot e^{-\Delta_{UG}/RT}$$
$$v_3 := \dots$$

+ Pondérations

+ Composition contrôlée (Boltzmann multidim. [Bodini P 10])

Rem. : Grammaires réalistes plus complexes (effets *Markoviens*) ...

II. Design sous contraintes

Éviter des motifs

But du design

Contraintes supplémentaires :

- **Interdire** des motifs d'apparaître **où que ce soit** ;
- **Obliger** des motifs à apparaître **au moins une fois** ;
- **Limiter** les bases disponibles à certaines positions.

Peu d'algorithmes tolèrent des motifs interdits/obligés ...

... aucun ne garantit des temps d'exécutions *raisonnables*.

Raisons typiques :

- Motifs obligés ⇒ Echecs tardifs (Branch and Bound)
- Motifs interdits ⇒ Espace d'états déconnecté (Recherche Locale)

But du design

Contraintes supplémentaires :

- **Interdire** des motifs d'apparaître **où que ce soit** ;
- **Obliger** des motifs à apparaître **au moins une fois** ;
- **Limiter** les bases disponibles à certaines positions.

Peu d'algorithmes tolèrent des motifs interdits/obligés ...

... aucun ne garantit des temps d'exécutions *raisonnables*.

Raisons typiques :

- Motifs obligés \Rightarrow Echecs tardifs (Branch and Bound)
- Motifs interdits \Rightarrow Espace d'états déconnecté (Recherche Locale)

Utiliser la clôture des langages hors-contexte

Motifs obligés
Motifs interdits

→ Langage Rationnel $\mathcal{L}_C \in \text{Reg}$

Compatibilité structure
+ Contraintes position
+ **Modèle d'énergie**

→ Langage Hors-Contextes $\mathcal{L}_S \in \text{LHC}$
(Pondérés)

Théorème Folklore (constructif) : $\text{Reg} \cap \text{LHC}(\mathbf{P}) \subseteq \text{LHC}(\mathbf{P})$

**Construire une grammaire hors-contexte \mathcal{G} pour $\mathcal{L}_C \cap \mathcal{L}_S$
+ Génération aléatoire
⇒ Échantillonnage global sous contrainte**

Utiliser la clôture des langages hors-contexte

Motifs obligés
Motifs interdits

→ Langage Rationnel $\mathcal{L}_C \in \text{Reg}$

Compatibilité structure
+ Contraintes position
+ **Modèle d'énergie**

→ Langage Hors-Contextes $\mathcal{L}_S \in \text{LHC}$
(Pondérés)

Théorème Folklore (constructif) : $\text{Reg} \cap \text{LHC}(\mathbf{P}) \subseteq \text{LHC}(\mathbf{P})$

**Construire une grammaire hors-contexte \mathcal{G} pour $\mathcal{L}_C \cap \mathcal{L}_S$
+ Génération aléatoire
⇒ Échantillonnage global sous contrainte**

Utiliser la clôture des langages hors-contexte

Motifs obligés
Motifs interdits

→ Langage Rationnel $\mathcal{L}_C \in \text{Reg}$

Compatibilité structure
+ Contraintes position
+ **Modèle d'énergie**

→ Langage Hors-Contextes $\mathcal{L}_S \in \text{LHC}$
(Pondérés)

Théorème Folklore (constructif) : $\text{Reg} \cap \text{LHC}(\mathbf{P}) \subseteq \text{LHC}(\mathbf{P})$

Construire une grammaire hors-contexte \mathcal{G} pour $\mathcal{L}_C \cap \mathcal{L}_S$
+ Génération aléatoire
⇒ Échantillonnage global sous contrainte

Utiliser la clôture des langages hors-contexte

Motifs obligés
Motifs interdits

→ Langage Rationnel $\mathcal{L}_C \in \text{Reg}$

Compatibilité structure
+ Contraintes position
+ **Modèle d'énergie**

→ Langage Hors-Contextes $\mathcal{L}_S \in \text{LHC}$
(Pondérés)

Théorème Folklore (constructif) : $\text{Reg} \cap \text{LHC}(\mathbf{P}) \subseteq \text{LHC}(\mathbf{P})$

**Construire une grammaire hors-contexte \mathcal{G} pour $\mathcal{L}_C \cap \mathcal{L}_S$
+ Génération aléatoire
⇒ Échantillonnage global sous contrainte**

Contraintes positionnelles

Entrée : Structure secondaire S

- A Construire l'arbre de syntaxe abstraite (ASA)** de S ;
- B Traduire l'ASA** en une grammaire ;
- C Étendre** la grammaire pour engendrer les lettres (**+ Poids**).
- D Limiter** les lettres autorisées à certaines positions

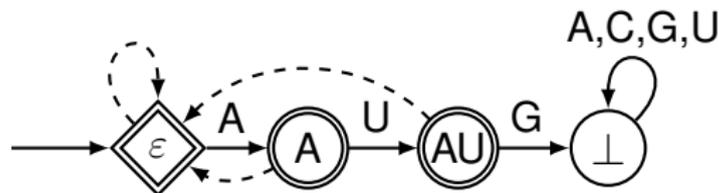
$$\begin{aligned}V_1 &\rightarrow A V_2 \mid C V_2 \mid G V_2 \mid U V_2 \\V_2 &\rightarrow A V_3 U \mid \underline{C V_3 G} \mid \underline{G V_3 C} \mid \underline{G V_3 U} \mid U V_3 A \mid \underline{U V_3 G} \\V_3 &\rightarrow A V_4 U V_8 \mid \underline{C V_4 G V_8} \mid \underline{G V_4 C V_8} \mid \underline{G V_4 U V_8} \mid U V_4 A V_8 \mid \underline{U V_4 G V_8} \\V_4 &\rightarrow A V_5 U \mid \underline{C V_5 G} \mid \underline{G V_5 C} \mid \underline{G V_5 U} \mid U V_5 A \mid \underline{U V_5 G} \\V_5 &\rightarrow A \mid C \mid G \mid U \\V_8 &\rightarrow A V_9 U \mid \underline{C V_9 G} \mid \underline{G V_9 C} \mid \underline{G V_9 U} \mid U V_9 A \mid \underline{U V_9 G} \\V_9 &\rightarrow A V_{10} \mid C V_{10} \mid G V_{10} \mid U V_{10} \\V_{10} &\rightarrow A \mid C \mid G \mid U\end{aligned}$$

Construire un automate fini déterministe

Motifs interdits $\mathcal{F} = \{f_1, f_2 \dots\}$ Motifs obligés $\mathcal{M} = \{m_1, m_2 \dots\}$

- Éviter un **motif unique** \Rightarrow Automate linéaire
- Éviter **plusieurs motifs** \Rightarrow Construction **Aho-Corasick**
- Obliger un **motif** \Rightarrow Automate linéaire (Complémenté)
- Obliger **plusieurs motifs en disjonction** \Rightarrow Aho-Corasick

Exemple : Éviter AUG ($\mathcal{L}_{\overline{\text{AUG}}}$)

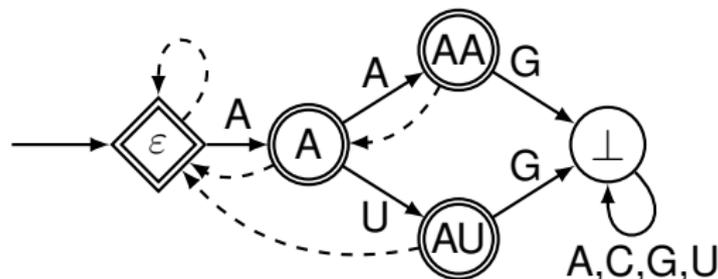


Construire un automate fini déterministe

Motifs interdits $\mathcal{F} = \{f_1, f_2 \dots\}$ Motifs obligés $\mathcal{M} = \{m_1, m_2 \dots\}$

- Éviter un **motif unique** \Rightarrow Automate linéaire
- Éviter **plusieurs motifs** \Rightarrow Construction **Aho-Corasick**
- Obliger un **motif** \Rightarrow Automate linéaire (Complémenté)
- Obliger **plusieurs motifs en disjonction** \Rightarrow Aho-Corasick

Exemple : Éviter AUG + AAG ($\mathcal{L}_{\overline{\text{AUG}}} \cap \mathcal{L}_{\overline{\text{AAG}}}$)

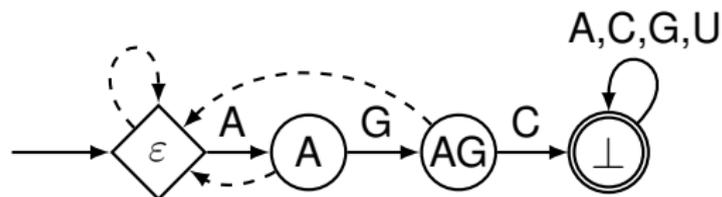


Construire un automate fini déterministe

Motifs interdits $\mathcal{F} = \{f_1, f_2 \dots\}$ Motifs obligés $\mathcal{M} = \{m_1, m_2 \dots\}$

- Éviter un **motif unique** \Rightarrow Automate linéaire
- Éviter **plusieurs motifs** \Rightarrow Construction **Aho-Corasick**
- Obliger un **motif** \Rightarrow Automate linéaire (Complémenté)
- Obliger **plusieurs motifs en disjonction** \Rightarrow Aho-Corasick

Exemple : Obliger AGC (\mathcal{L}_{AGC})

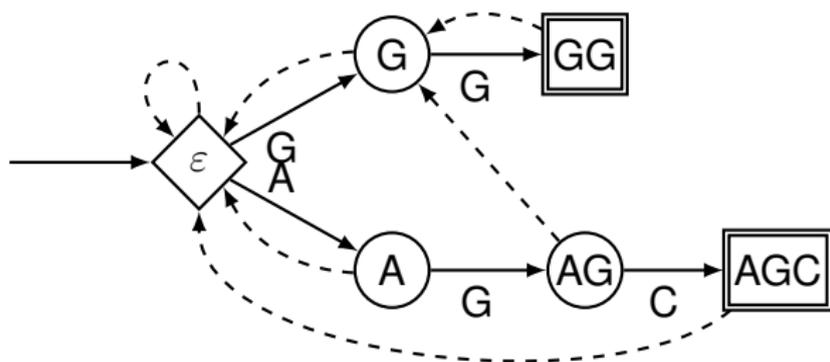


Construire un automate fini déterministe

Motifs interdits $\mathcal{F} = \{f_1, f_2 \dots\}$ Motifs obligés $\mathcal{M} = \{m_1, m_2 \dots\}$

- Éviter un **motif unique** \Rightarrow Automate linéaire
- Éviter **plusieurs motifs** \Rightarrow Construction **Aho-Corasick**
- Obliger un **motif** \Rightarrow Automate linéaire (Complémenté)
- Obliger **plusieurs motifs en disjonction** \Rightarrow Aho-Corasick

Exemple : Obliger AGC **ou** GG ($\mathcal{L}_{AGC} \cup \mathcal{L}_{GG}$)



Construire un automate fini déterministe

Obliger l'apparition des obligés

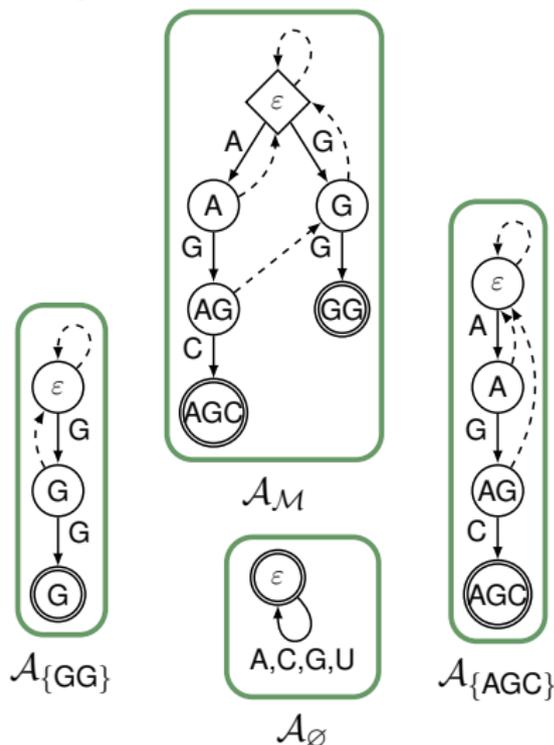
→ **retenir ceux à engendrer** :

- Créer l'automate disjonctif $\mathcal{A}_{\mathcal{M}'}$ pour tout $\mathcal{M}' \subseteq \mathcal{M}$
- **Rediriger** états acceptants
- État acceptant = **Tout motif obligé** est généré ($\varepsilon \in \mathcal{A}_{\emptyset}$)
- Ajouter motifs interdits aux sous-automates

#États :

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

Exemple : $\mathcal{M} = \{\text{AGC}, \text{GG}\}$



Construire un automate fini déterministe

Obliger l'apparition des obligés

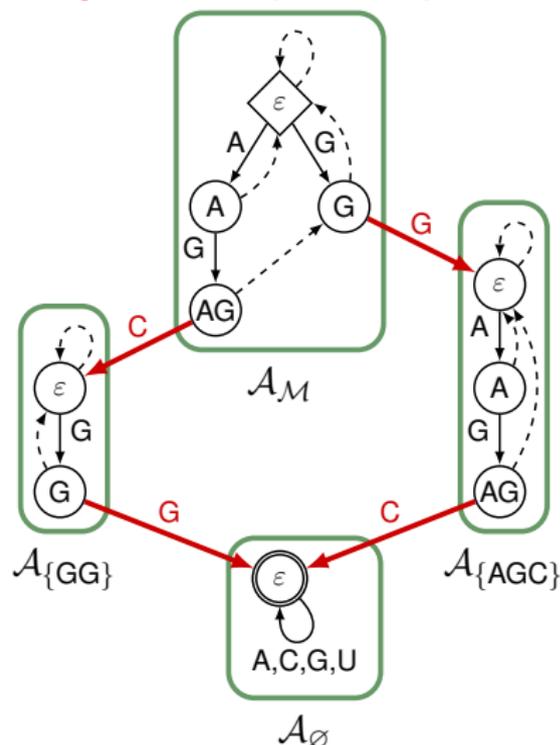
→ **retenir ceux à engendrer** :

- Créer l'automate disjonctif $\mathcal{A}_{\mathcal{M}'}$ pour tout $\mathcal{M}' \subseteq \mathcal{M}$
- **Rediriger** états acceptants
- État acceptant = **Tout motif obligé** est généré ($\varepsilon \in \mathcal{A}_{\emptyset}$)
- Ajouter motifs interdits aux sous-automates

#États :

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

Exemple : $\mathcal{M} = \{AGC, GG\}$



Construire un automate fini déterministe

Obliger l'apparition des obligés

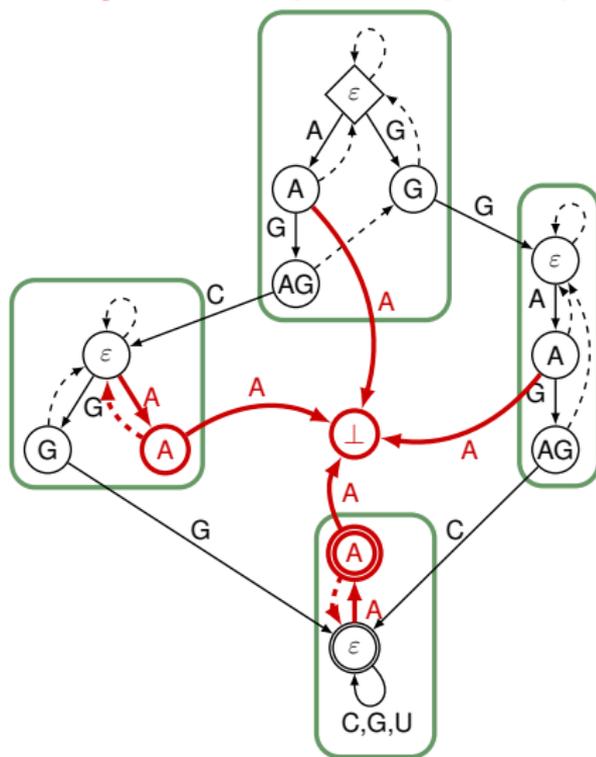
→ **retenir ceux à engendrer** :

- Créer l'automate disjonctif $\mathcal{A}_{\mathcal{M}'}$ pour tout $\mathcal{M}' \subseteq \mathcal{M}$
- **Rediriger** états acceptants
- État acceptant = **Tout motif obligé** est généré ($\varepsilon \in \mathcal{A}_{\emptyset}$)
- Ajouter motifs interdits aux sous-automates

#États :

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

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



Combiner une grammaire \mathcal{G} et un automate \mathcal{A}

Idée (Folklore) :

- **Simuler l'exécution** de \mathcal{A} lors des dérivations de \mathcal{G} ;
- Tout non-terminal $V_{q \rightarrow q'}$ **s'engage** à **relier** deux états de \mathcal{A} ;
- Les règles terminales **vérifient** le respect de l'engagement ;
- Certaines règles deviennent **improductives** \rightarrow nettoyage.

$$V \rightarrow t V' \quad \Longrightarrow \quad V_{q \rightarrow q'} \rightarrow t V'_{\delta(q,t) \rightarrow q'}$$

$$V \rightarrow t V' t' \quad \Longrightarrow \quad V_{q \rightarrow q'} \rightarrow t V'_{\delta(q,t) \rightarrow q''} t' \text{ t.q. } \delta(q'', t') = q'$$

$$V \rightarrow t V' t' V'' \quad \Longrightarrow \quad V_{q \rightarrow q'} \rightarrow t V'_{\delta(q,t) \rightarrow q''} t' V''_{\delta(q'', t') \rightarrow q'}$$

$$V \rightarrow t \quad \Longrightarrow \quad \begin{cases} V_{q \rightarrow q'} \rightarrow t & \text{if } \delta(q, t) = q' \\ \emptyset & \text{otherwise} \end{cases}$$

$\delta(\cdot, \cdot)$: Fonction de transition de \mathcal{A}

Complexité

- **Automate** : Temps \approx #États $|Q| \in \mathcal{O}\left(2^{|\mathcal{M}|} \cdot (\sum_i |f_i| + \sum_j |m_j|)\right)$
- **Grammaire initiale** : #Règles \approx #Non Terminaux $\in \Theta(n)$
- **Grammaire contrainte** :
 - #Non Terminaux $|N| \in \mathcal{O}(n \cdot |Q|^2)$
 - #Règles $|R| \in \mathcal{O}(n \cdot |Q|^3)$
- **Génération aléatoire**
 - Précalcul** : $\Theta^*(|R|)$ temps/ $\Theta(|N|)$ mémoire
 - Génération** : $\Theta(k \cdot n \cdot |Q|)$ pour k séquences
 - Rem.** : Longueur des mots produits déterminée par NT.

L'algorithme entier a **complexité linéaire sur la taille**
mais la **taille totale des motifs** influence le précalcul.

Conclusion/Perspectives

Fait :

- Algorithme en $\Theta(n)$ (+ctes raisonnables) pour le design **positif**
- Outils :
 - Méthode récursive [Flajolet Van Cutsem Zimmermann 94, Denise P Termier 10]
 - Boltzmann multidimensionnel [Bodini P 10]
- Impliqué dans une expérience biologique **contrôlée**
- Aléa pondéré **compétitif/complémentaire** avec recherche locale

Conclusion/Perspectives

Fait :

- Algorithme en $\Theta(n)$ (+ctes raisonnables) pour le design **positif**
- Outils :
 - Méthode récursive [Flajolet Van Cutsem Zimmermann 94, Denise P Termier 10]
 - Boltzmann multidimensionnel [Bodini P 10]
- Impliqué dans une expérience biologique **contrôlée**
- Aléa pondéré **compétitif/complémentaire** avec recherche locale

À faire/comprendre :

- Meilleures constructions pour produits automates/grammaires.
En particulier, non-terminal \rightarrow taille \Rightarrow **gain asymptotique ?**
- Plus petits automates **sans minimiser** ? [Ait Mous Bassino Nicaud 12]
- Génération **sans redondance** [Lorenz P 12]
- Etude de complexité (Obstacles : optim. inverse, dénaturation . . .)

Remerciements

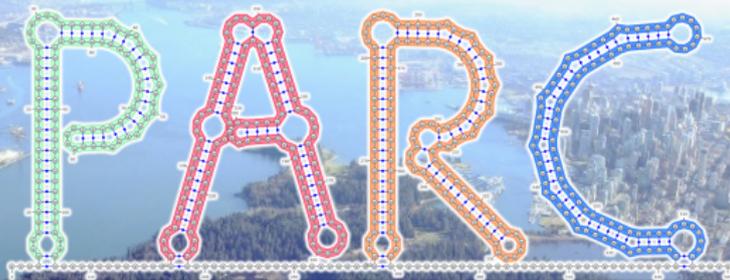
● Co-auteurs

- Yu Zhou (UCSD)
- Stéphane Vialette (LIGM)
- Jérôme Waldispühl
+ Vlad Reinharz (McGill)
- Yi Zhang (Wuhan, China)
- Alain Denise (LRI/IGM)



Merci de votre attention !





PIMS Analytic RNA Combinatorics

April 15 - 16 2014

Simon Fraser University, Burnaby

Speakers

Ralf Bundschuh (Ohio State)

Anne Condon (UBC)

Christine Heitsch (Georgia Tech)

Jan Manuch (UBC/SFU)

Markus Nebel (TU Kaiserslautern)

Yann Ponty (CNRS/Polytechnique)

Current RNA bioinformatics strongly relies on discrete abstractions and decomposition schemes for the conformation of nucleic acids. This workshop will bring together scientists from bioinformatics, mathematics and computer science in an attempt to bridge the gap between the enumeration and (asymptotic) analysis of models for RNA structure, and the design of RNA algorithms.

<http://tinyurl.com/parc-2014> contact: mmishna@sfu.ca yann.ponty@lix.polytechnique.fr