

## Comptage et design multiple d'ARN

Stefan Hammer<sup>†</sup>

Yann Ponty<sup>\*</sup>

Wei Wang<sup>\*,•</sup>

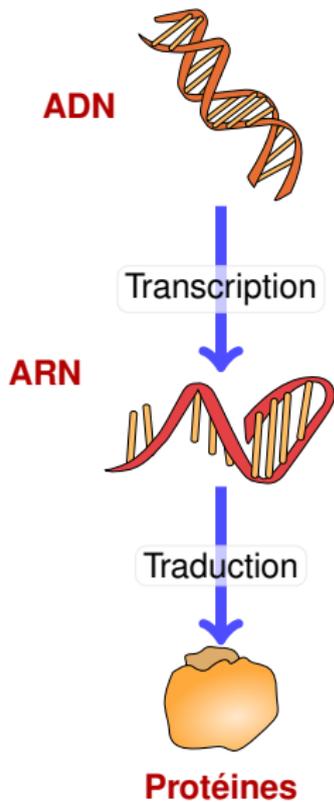
Sebastian Will<sup>†</sup>

\* LIX, CNRS/Ecole Polytechnique

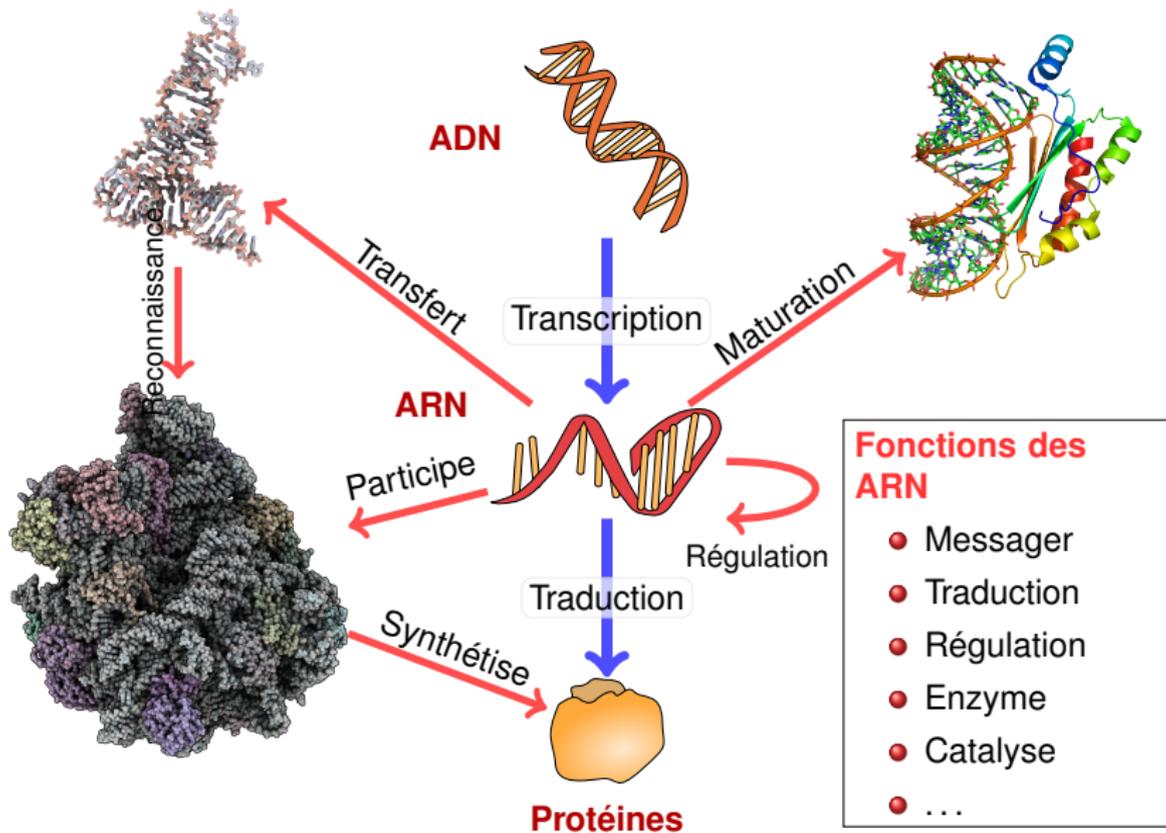
• LRI, Université Paris-Sud

† TBI, University of Vienna

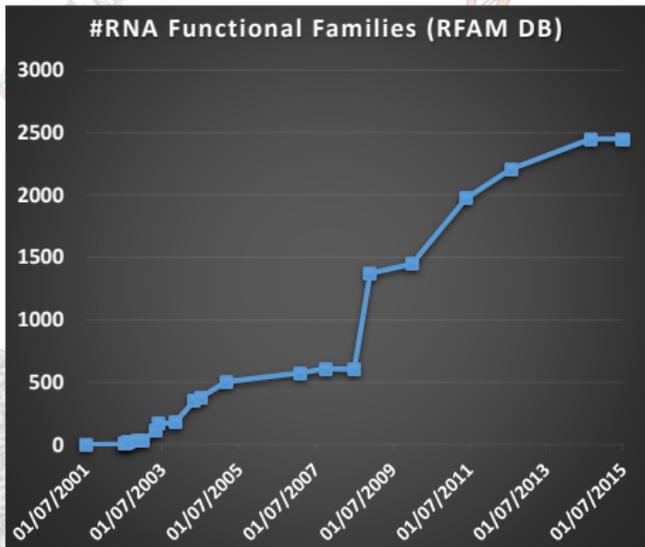
# **Dogme fondamental de la biologie moléculaire**



# Dogme fondamental de la biologie moléculaire (v2.0)



# Dogme fondamental de la biologie moléculaire



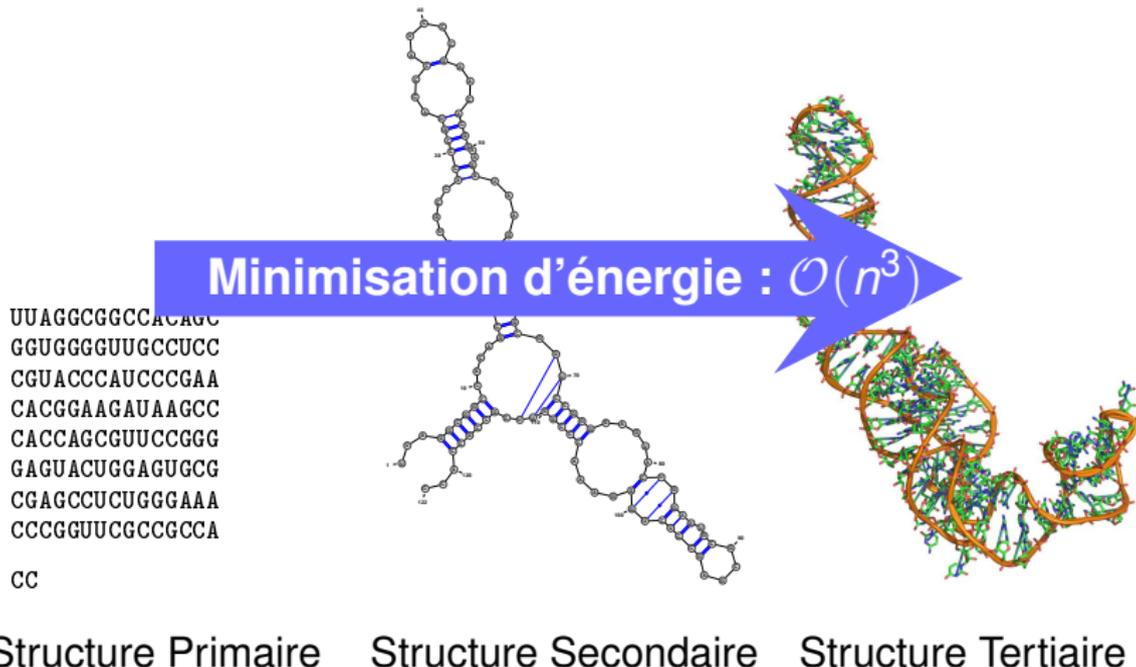
## Fonctions des ARN

- Messenger
- Traduction
- Régulation
- Enzyme
- Catalyse
- ...

Protéines

# Séquence et structure(s) des ARN

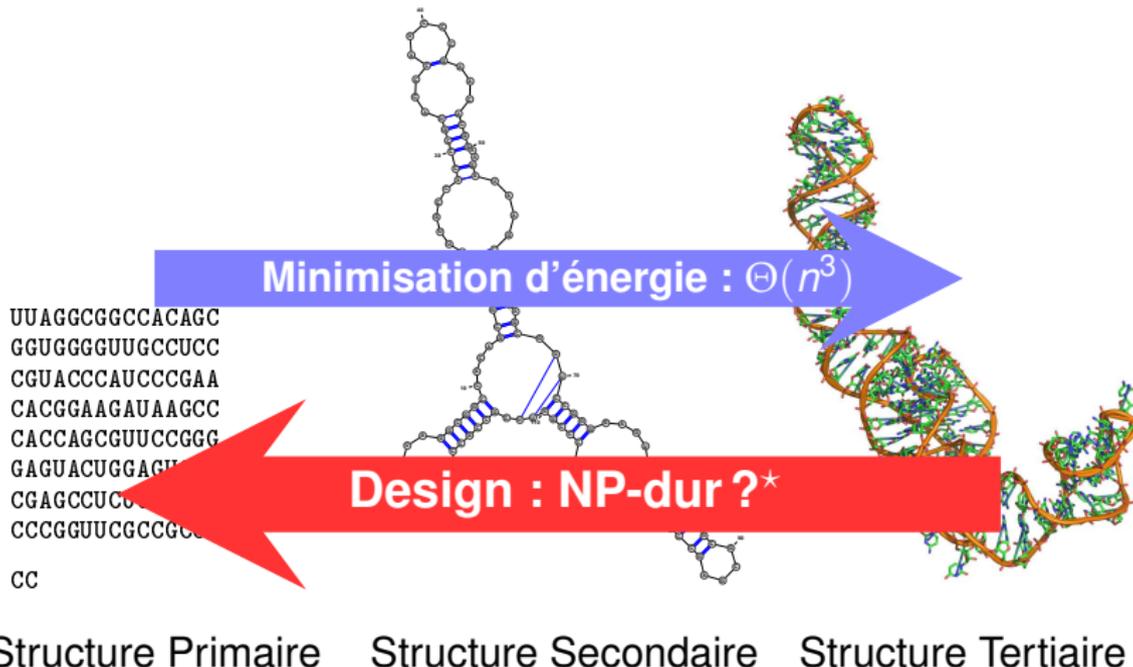
**ARN** = Polymère linéaire = Séquence sur {A, C, G, U}<sup>\*</sup>



ARNr 5s (PDBID : 1K73 :B)

# Séquence et structure(s) des ARN

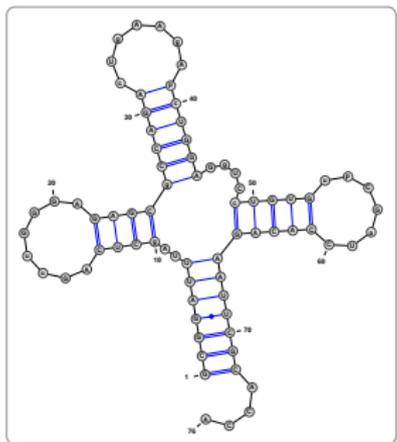
**ARN** = Polymère linéaire = Séquence sur {A, C, G, U}<sup>\*</sup>



ARNr 5s (PDBID : 1K73 :B)

<sup>\*</sup>Probablement ... (cf [Bonnet/Rzężewski/Sikora, RECOMB'18])

# Représentations de la structure secondaire



Graphe outer-planar

Hamiltonien,  $\Delta(G) \leq 3$ , 2-connexe\*

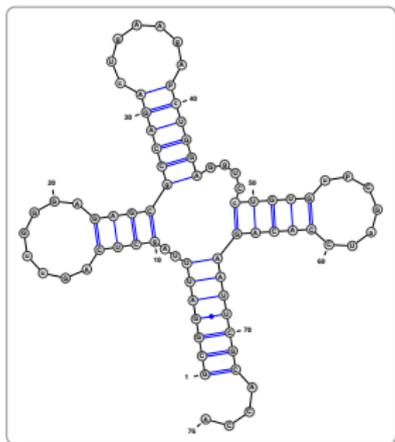
## Aide à l'intuition

Différentes représentations

Même objet combinatoire

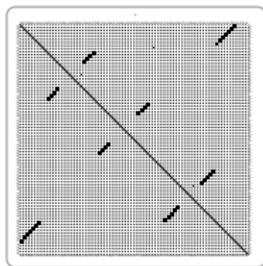
\*Contraintes stériques, pseudo-noeuds

# Représentations de la structure secondaire



Graphe outer-planar

Hamiltonien,  $\Delta(G) \leq 3$ , 2-connecté\*



Dot plots

Matrice d'adjacence\*

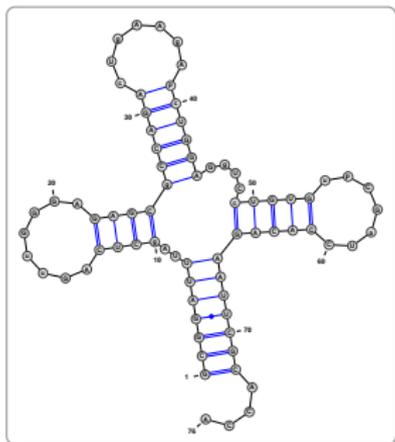
## Aide à l'intuition

Différentes représentations

Même objet combinatoire

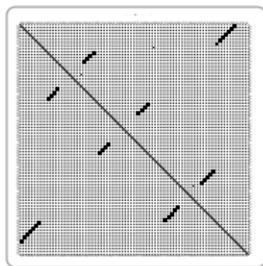
\*Contraintes stériques, pseudo-noeuds

# Représentations de la structure secondaire



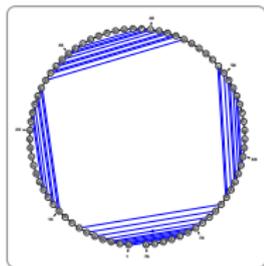
Graphe outer-planar

Hamiltonien,  $\Delta(G) \leq 3$ , 2-connexe\*



Dot plots

Matrice d'adjacence\*



Graphe de cordes\*

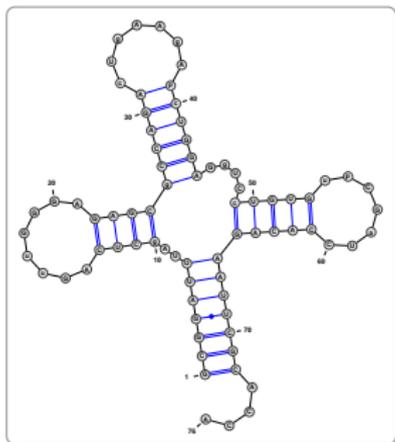
## Aide à l'intuition

Différentes représentations

Même objet combinatoire

\*Contraintes stériques, pseudo-noeuds

# Représentations de la structure secondaire

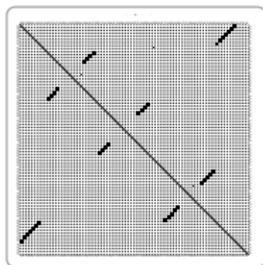


Graphe outer-planar

Hamiltonien,  $\Delta(G) \leq 3$ , 2-connexe\*

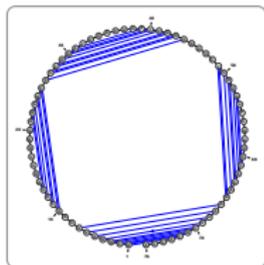
(((((((...(((.....))))))((((.....))))))....((((.....)))))))))....

Mots de Motzkin\*



Dot plots

Matrice d'adjacence\*



Graphe de cordes\*

## Aide à l'intuition

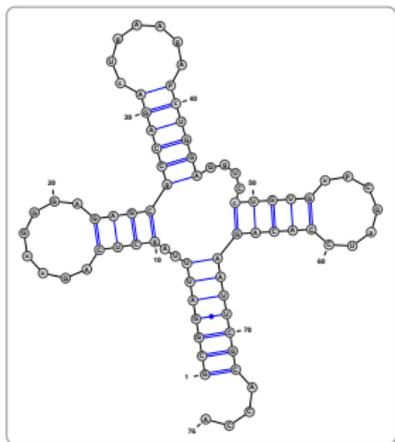
Différentes représentations

Même objet combinatoire

\*Contraintes stériques, pseudo-noeuds



# Représentations de la structure secondaire

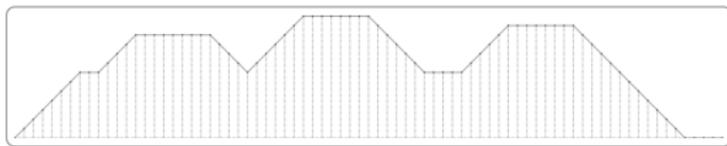


Graphe outer-planar

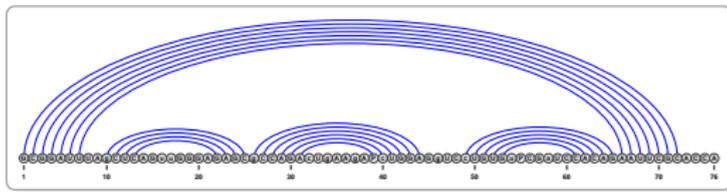
Hamiltonien,  $\Delta(G) \leq 3$ , 2-connecté\*

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

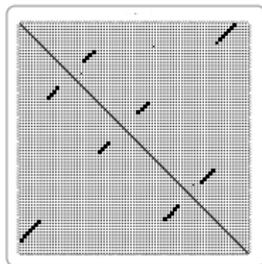
Mots de Motzkin\*



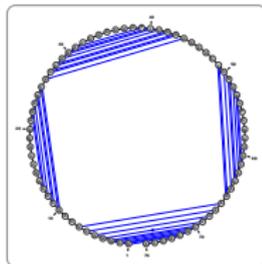
Excursions\* sur  $S = \{+1, -1, 0\}$



Séquence arc-annotées\*



Dot plots  
Matrice d'adjacence\*



Graphe de cordes\*

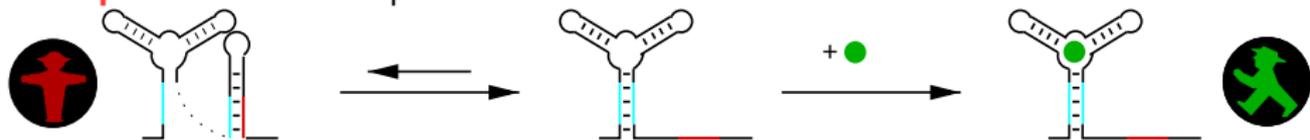
## Aide à l'intuition

Différentes représentations  
Même objet combinatoire

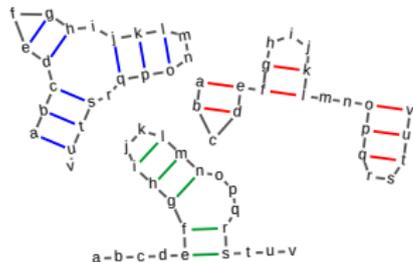
\* Contraintes stériques, pseudo-noeuds

# Design multiple d'ARN : Motivation

**Exemple :** Un *riboswitch* pour le contrôle de la traduction



Plusieurs structures cibles → *Design multiple d'ARN*



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

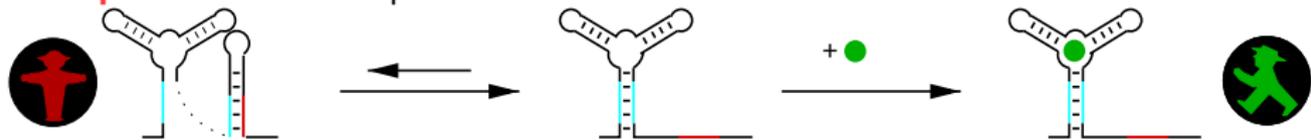
**Objectif :** Engendrer **aléatoirement** des séquences sous contraintes

- 1 **Validité** pour toutes les structures cibles
- 2 **Stabilité** (énergie-libre faible, comparable ...) des structures cibles
- 3 **Composition contrainte** : (contenu prescrit en GC), motifs exclus/forcés ...

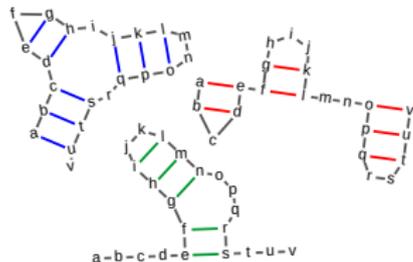
**Méthode réc. (prog. dyn.) :** Compter puis engendrer séquences **valides** uniformes

# Design multiple d'ARN : Motivation

**Exemple :** Un *riboswitch* pour le contrôle de la traduction



**Plusieurs structures cibles** → *Design multiple d'ARN*



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

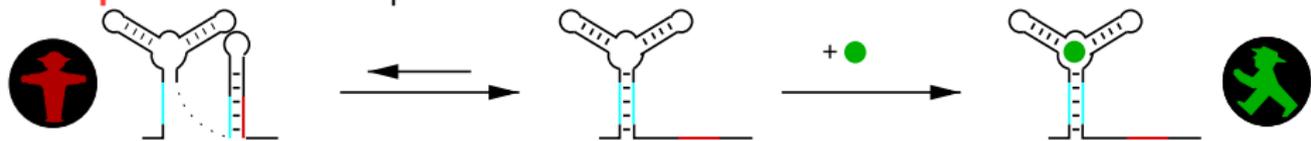
**Objectif :** Engendrer **aléatoirement** des séquences sous contraintes

- 1 **Validité** pour toutes les structures cibles
- 2 **Stabilité** (énergie-libre faible, comparable ...) des structures cibles
- 3 **Composition contrainte** : (contenu prescrit en GC), motifs exclus/forcés ...

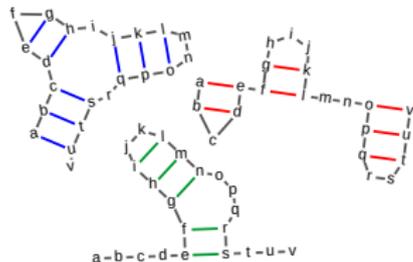
**Méthode réc. (prog. dyn.) :** Compter puis engendrer séquences **valides** uniformes

# Design multiple d'ARN : Motivation

**Exemple :** Un *riboswitch* pour le contrôle de la traduction



Plusieurs structures cibles → *Design multiple d'ARN*



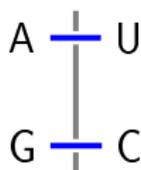
abcdefghijklmnopqrstuv  
((( ( ( ( ( . ) ) . ( ( ( . ) ) ) ) ) ) ) .  
( ( . ) ) ( ( . . . ) ) . . ( ( ( . . . ) ) )  
..... ( ( ( ( ( . . . ) ) ) ) . . . ) . . .

**Objectif :** Engendrer **aléatoirement** des séquences sous contraintes

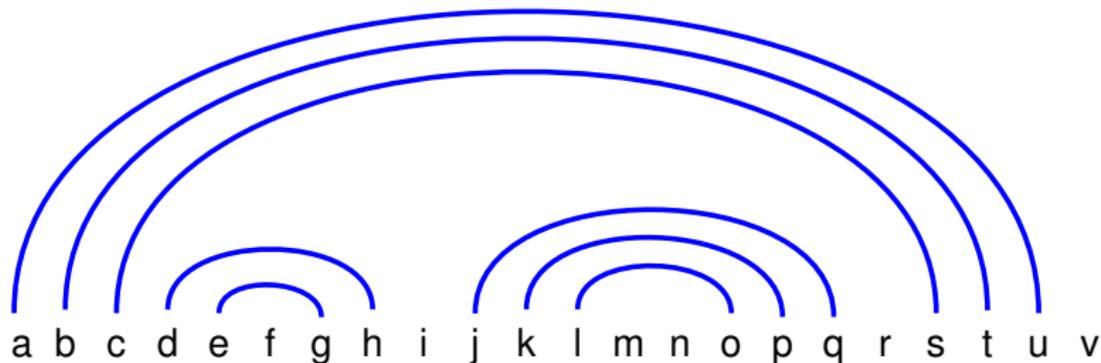
- 1 **Validité** pour toutes les structures cibles
- 2 **Stabilité** (energie-libre faible, comparable ...) des structures cibles
- 3 **Composition contrainte :** (contenu prescrit en GC), motifs exclus/forcés ...

**Méthode réc. (prog. dyn.) :** Compter puis engendrer séquences **valides** uniformes

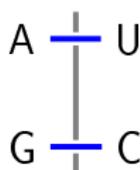
## Compter les séquences valides : Watson-Crick + structure unique



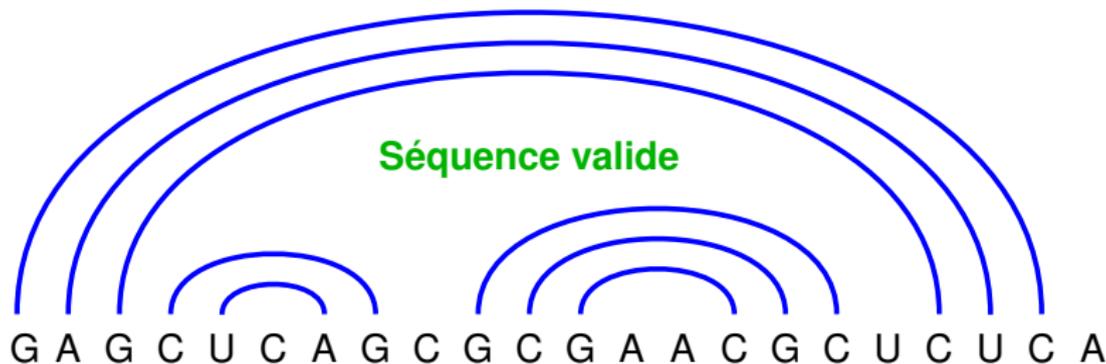
**Paires de Bases (PB) valides = Watson-Crick** seulement



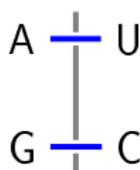
## Compter les séquences valides : Watson-Crick + structure unique



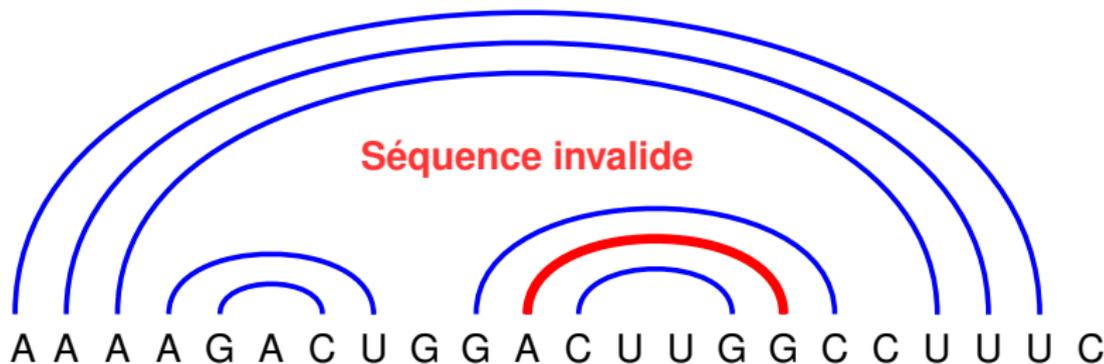
**Paires de Bases (PB) valides = Watson-Crick** seulement



## Compter les séquences valides : Watson-Crick + structure unique



Paires de Bases (PB) valides = Watson-Crick seulement

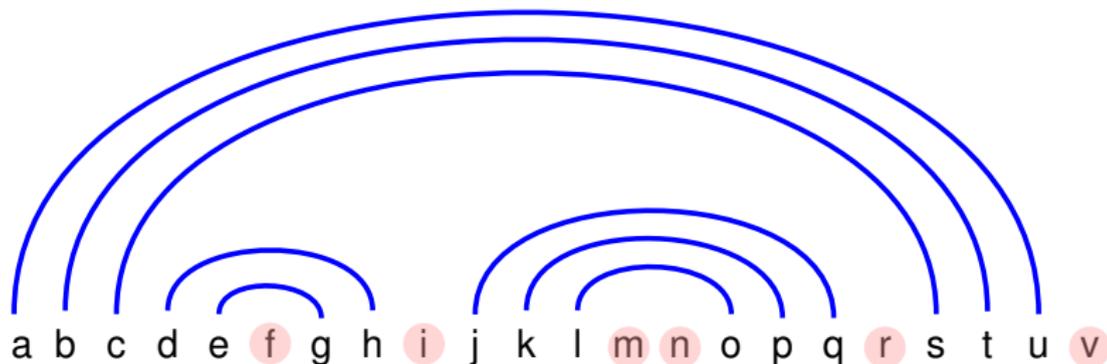


**Question :** Combien de séquences **valides** ?

## Compter les séquences valides : Watson-Crick + structure unique

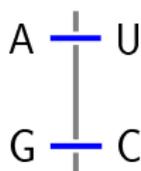


Paires de Bases (PB) valides = Watson-Crick seulement

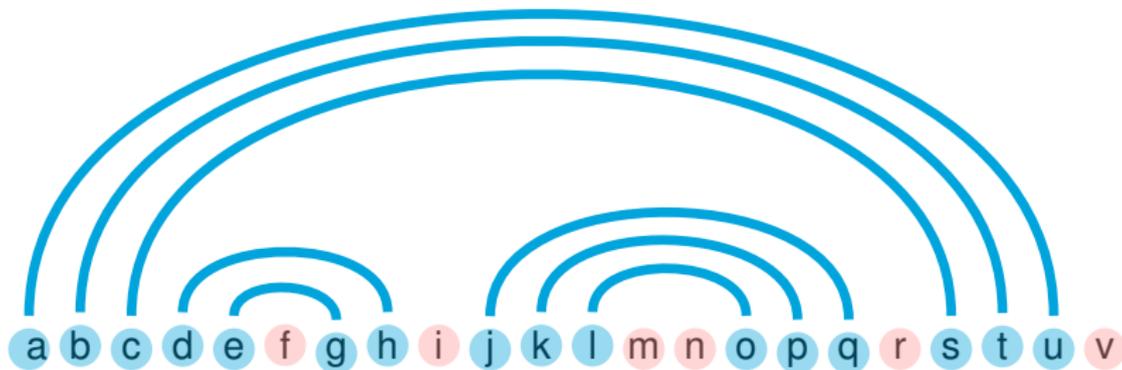


**Question :** Combien de séquences **valides** ?

## Compter les séquences valides : Watson-Crick + structure unique

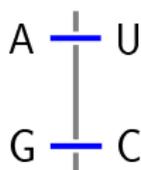


**Paires de Bases (PB) valides = Watson-Crick** seulement

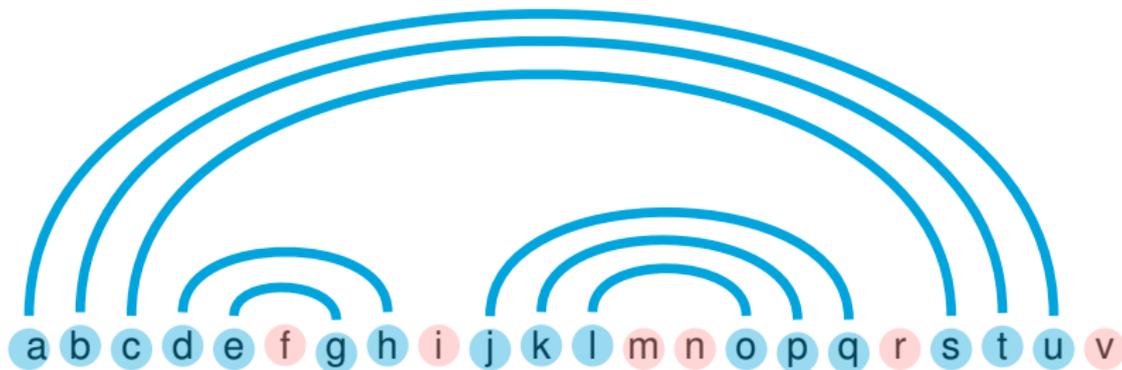


**Question :** Combien de séquences **valides** ?

## Compter les séquences valides : Watson-Crick + structure unique



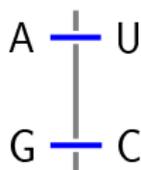
**Paires de Bases (PB) valides = Watson-Crick** seulement



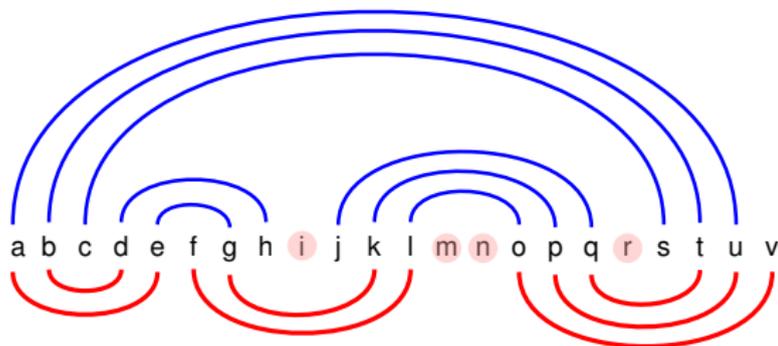
**Question :** Combien de séquences **valides** ?

**Réponse :**  $4^{\# \text{Paires}} \times 4^{\# \text{Non-appariées}} \rightarrow 268\,435\,456$

# Compter les séquences valides : Watson-Crick + 2 structures



Paires de Bases (PB) valides = Watson-Crick seulement



**Graphe de dépendance :**

Cycles + Chemins

i m n r

g-e-a-u h j-q

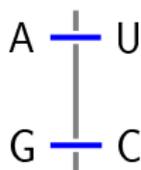
k-p d-b-t

f-l-o-v c-s

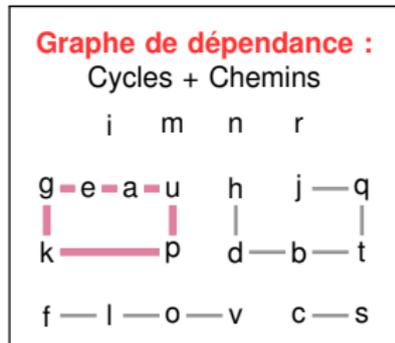
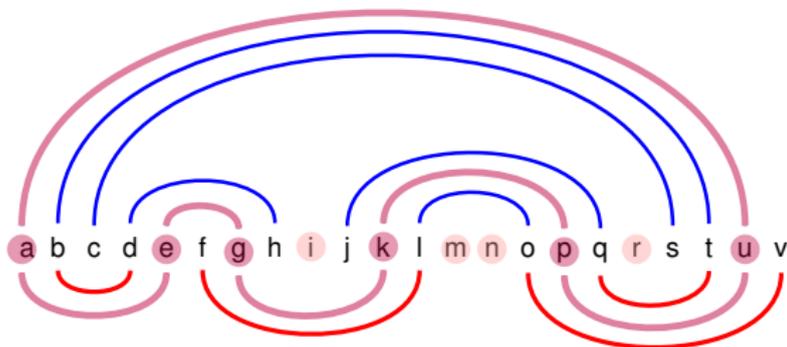
**Question :** Combien de séquences **valides** ?

**Réponse :**  $\neq \emptyset$  (graphes de dépendance ET des paires valides **bipartis**)

# Compter les séquences valides : Watson-Crick + 2 structures



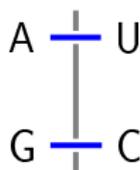
Paires de Bases (PB) valides = Watson-Crick seulement



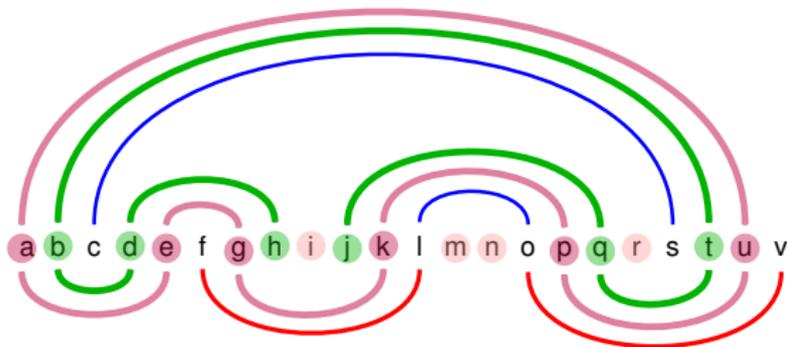
**Question :** Combien de séquences **valides** ?

**Réponse :**  $\neq \emptyset$  (graphes de dépendance ET des paires valides **bipartis**)

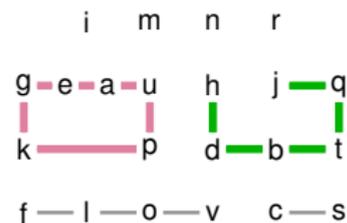
# Compter les séquences valides : Watson-Crick + 2 structures



Paires de Bases (PB) valides = Watson-Crick seulement



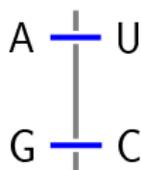
**Graphe de dépendance :**  
Cycles + Chemins



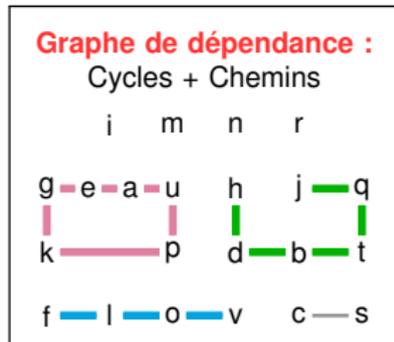
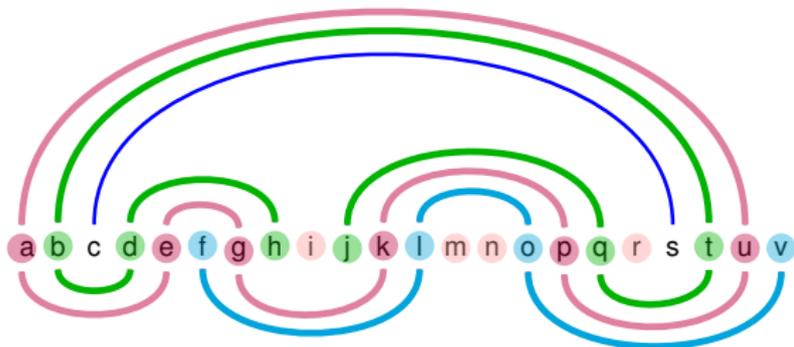
**Question :** Combien de séquences **valides** ?

**Réponse :**  $\neq \emptyset$  (graphes de dépendance ET des paires valides **bipartis**)

# Compter les séquences valides : Watson-Crick + 2 structures



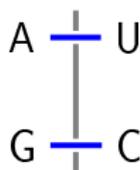
Paires de Bases (PB) valides = Watson-Crick seulement



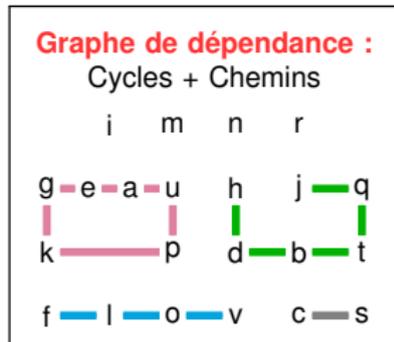
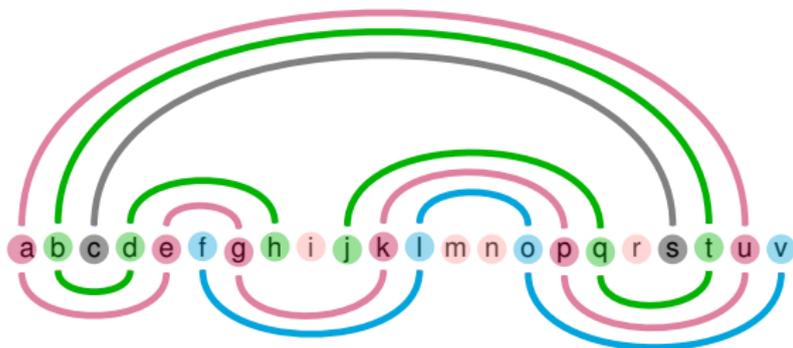
**Question :** Combien de séquences **valides** ?

**Réponse :**  $\neq \emptyset$  (graphes de dépendance ET des paires valides **bipartis**)

# Compter les séquences valides : Watson-Crick + 2 structures



Paires de Bases (PB) valides = Watson-Crick seulement

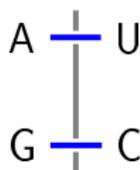


**Question :** Combien de séquences **valides** ?

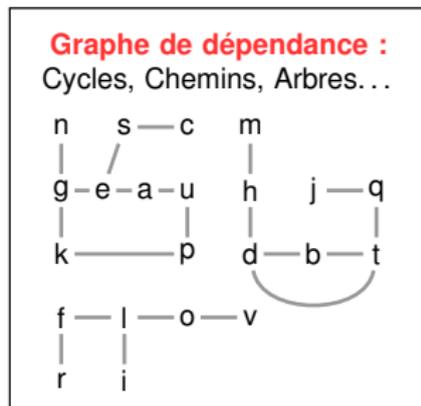
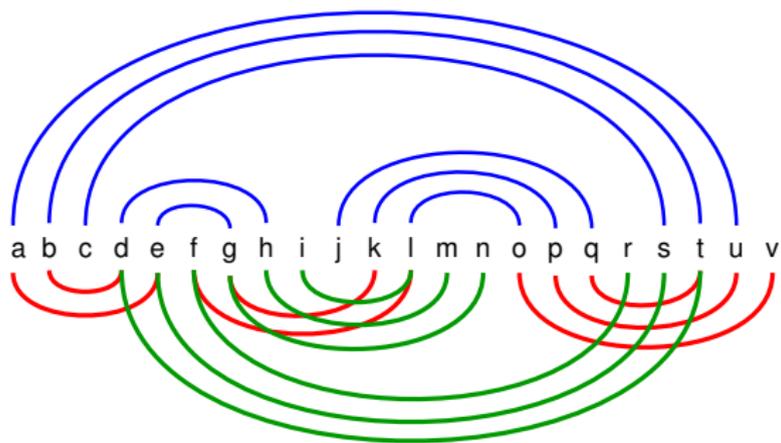
**Réponse :**  $\neq \emptyset$  (graphes de dépendance ET des paires valides **bipartis**)

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

# Compter les séquences valides : Watson-Crick + > 2 structures



**Paires de Bases (PB) valides = Watson-Crick** seulement



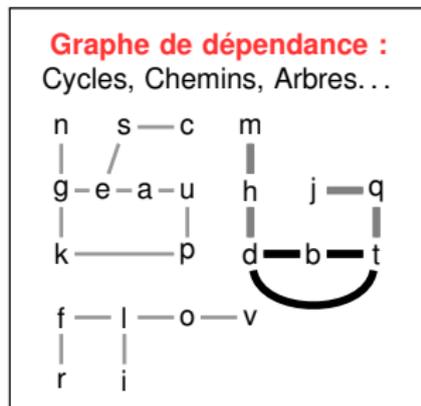
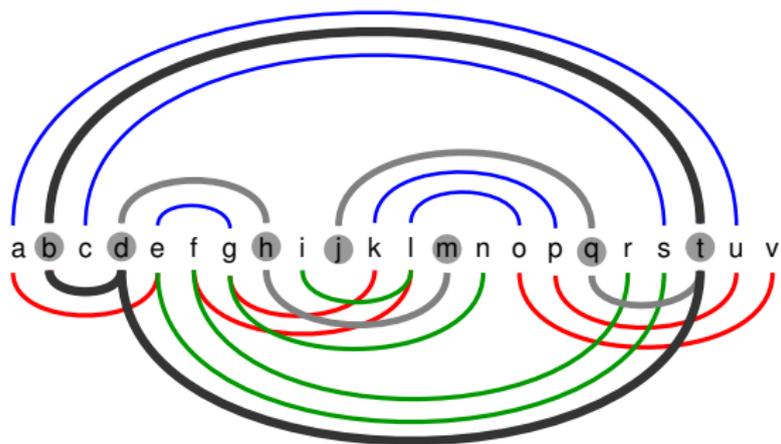
**Question :** Combien de séquences **valides** ?

**Réponse :** Non-biparti  $\rightarrow \emptyset$  ; Biparti  $\rightarrow 4^{\#CCs} = 64$

# Compter les séquences valides : Watson-Crick + > 2 structures



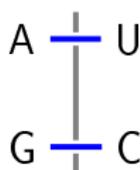
**Paires de Bases (PB) valides = Watson-Crick** seulement



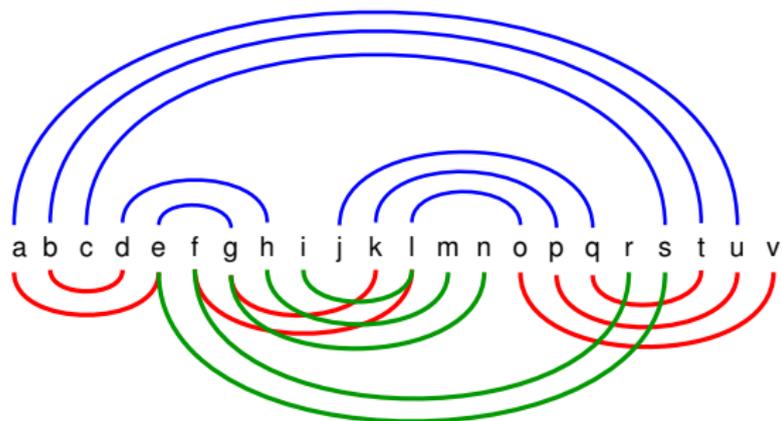
**Question :** Combien de séquences **valides** ?

**Réponse :** Non-biparti  $\rightarrow \emptyset$  ; Biparti  $\rightarrow 4^{\#CCs} = 64$

# Compter les séquences valides : Watson-Crick + > 2 structures

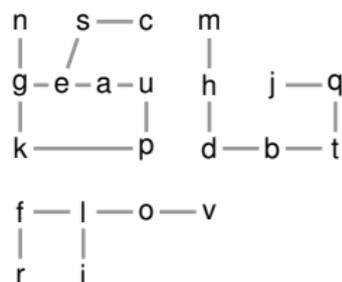


**Paires de Bases (PB) valides = Watson-Crick** seulement



**Graphe de dépendance :**

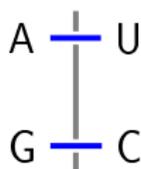
Cycles, Chemins, Arbres...



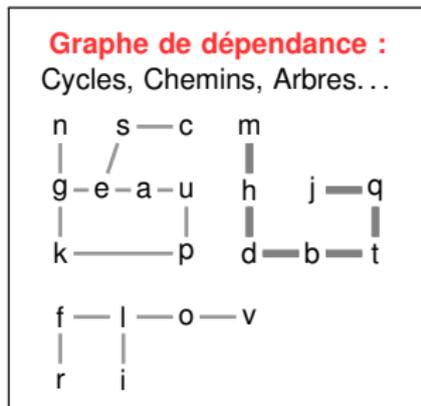
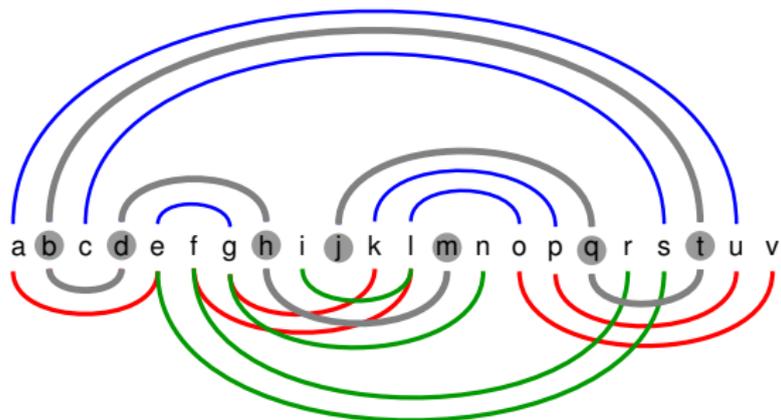
**Question :** Combien de séquences **valides** ?

**Réponse :** Non-biparti  $\rightarrow \emptyset$ ; Biparti  $\rightarrow 4^{\#CCs} = 64$

# Compter les séquences valides : Watson-Crick + > 2 structures



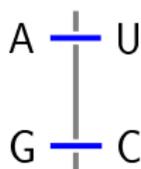
**Paires de Bases (PB) valides = Watson-Crick** seulement



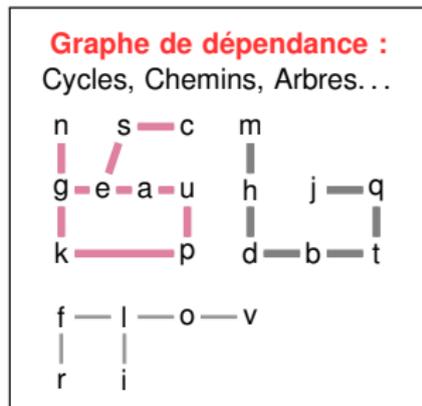
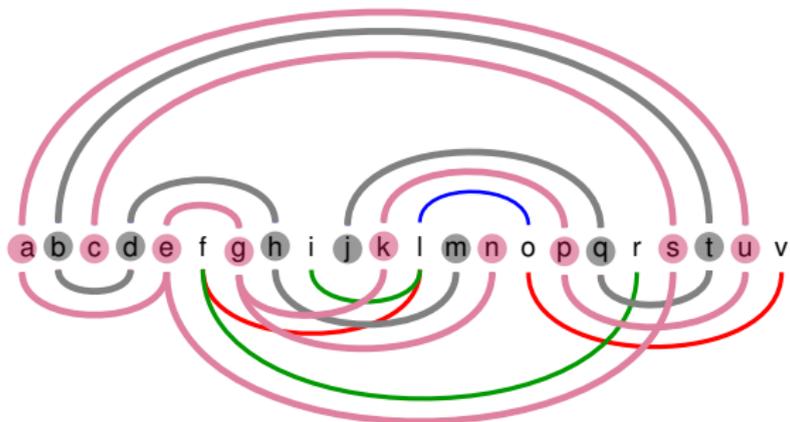
**Question :** Combien de séquences **valides** ?

**Réponse :** Non-biparti  $\rightarrow \emptyset$ ; Biparti  $\rightarrow 4^{\#CCs} = 64$

# Compter les séquences valides : Watson-Crick + > 2 structures



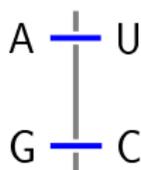
Paires de Bases (PB) valides = Watson-Crick seulement



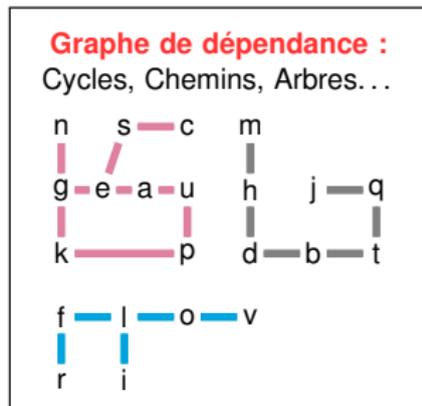
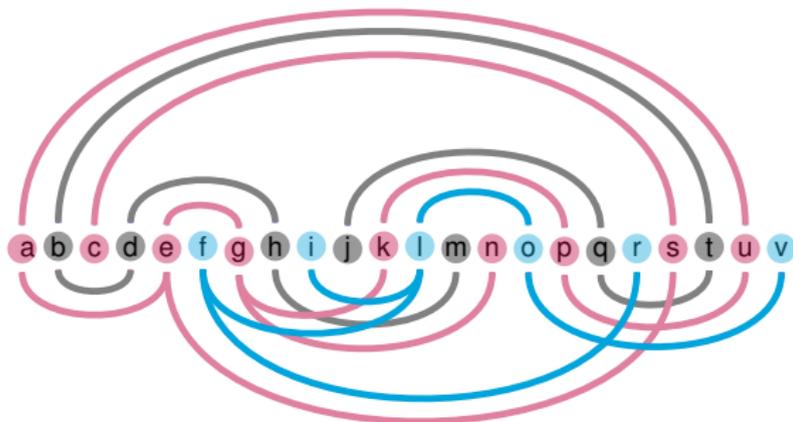
**Question :** Combien de séquences **valides** ?

**Réponse :** Non-biparti  $\rightarrow \emptyset$ ; Biparti  $\rightarrow 4^{\#CCs} = 64$

# Compter les séquences valides : Watson-Crick + > 2 structures



**Paires de Bases (PB) valides = Watson-Crick** seulement



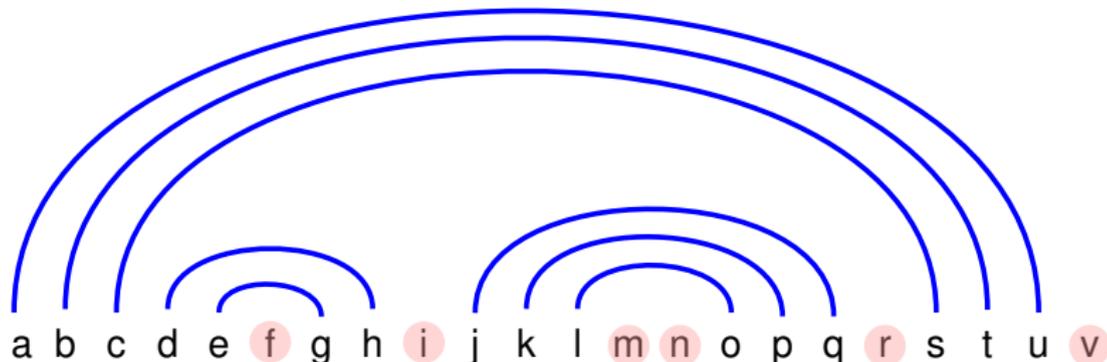
**Question :** Combien de séquences **valides** ?

**Réponse :** Non-biparti  $\rightarrow \emptyset$  ; Biparti  $\rightarrow 4^{\#CCs} = 64$

## Compter les séquences valides : WC/Wobble + unique structure



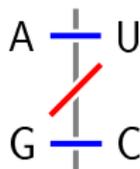
**Paires de Bases (PB) valides** = autorise les paires **Wobble**



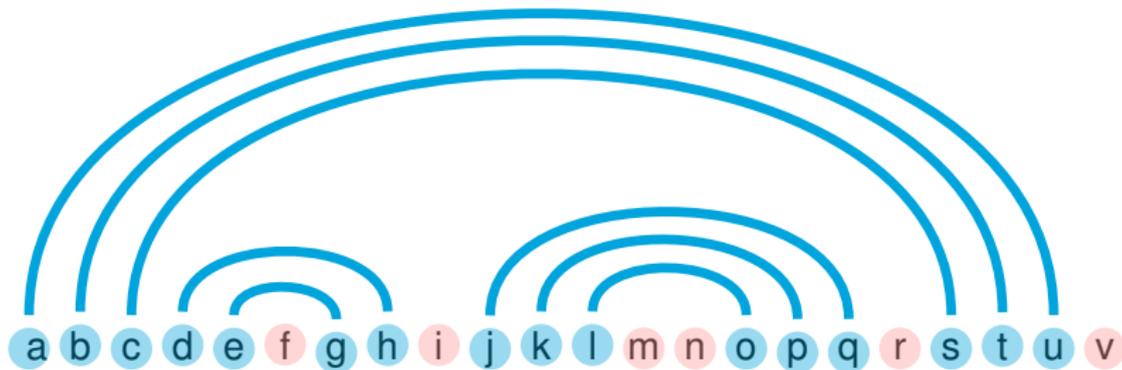
**Question :** Combien de séquences **valides** ?

**Réponse :** 4<sup>#Non-appariées</sup> × 6<sup>#PBs</sup> → 6879707136

## Compter les séquences valides : WC/Wobble + unique structure



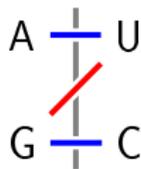
**Paires de Bases (PB) valides** = autorise les paires **Wobble**



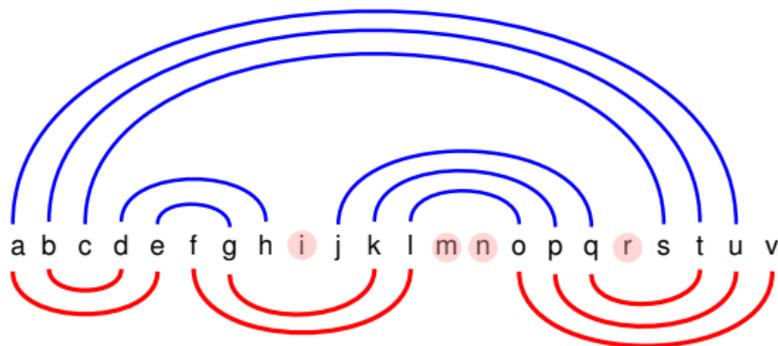
**Question :** Combien de séquences **valides** ?

**Réponse :**  $4^{\# \text{Non-appariées}} \times 6^{\# \text{PBs}} \rightarrow 6\ 879\ 707\ 136$

# Compter les séquences valides : WC/Wobble + Deux structures



**Paires de Bases valides** = autorise les paires **Wobble**



**Graphe de dépendance :**

Cycles + Chemins

i m n r

g-e-a-u h j-q

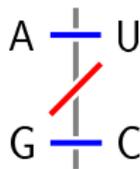
k-p d-b-t

f-l-o-v c-s

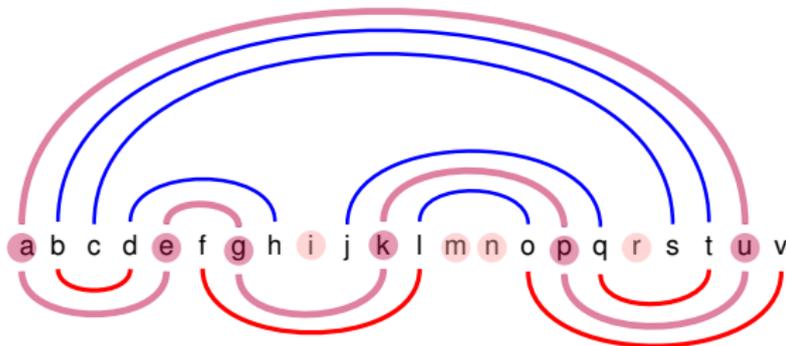
**Question :** Combien de séquences **valides** ?

**Réponses :**  $\neq \emptyset$  ! (graphes de dépendance et paires de bases tous deux **bipartis**)

# Compter les séquences valides : WC/Wobble + Deux structures



**Paires de Bases valides** = autorise les paires **Wobble**



**Graphe de dépendance :**

Cycles + Chemins

i m n r

g-e-a-u h j-q

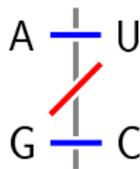
k-p d-b-t

f-l-o-v c-s

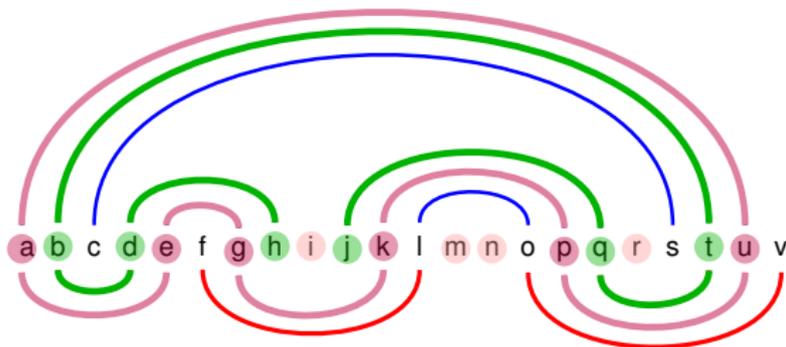
**Question :** Combien de séquences **valides** ?

**Réponses :**  $\neq \emptyset$  ! (graphes de dépendance et paires de bases tous deux **bipartis**)

# Compter les séquences valides : WC/Wobble + Deux structures



**Paires de Bases valides** = autorise les paires **Wobble**



**Graphe de dépendance :**

Cycles + Chemins

i m n r

g-e-a-u h j q

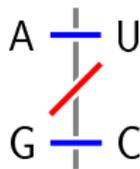
k-p d-b-t

f-l-o-v c-s

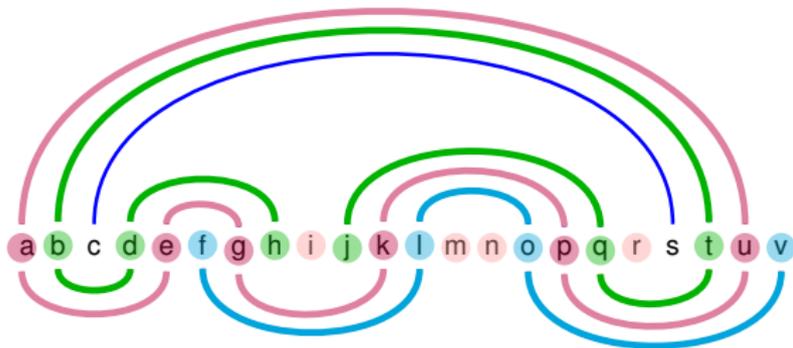
**Question :** Combien de séquences **valides** ?

**Réponses :**  $\neq \emptyset$  ! (graphes de dépendance et paires de bases tous deux **bipartis**)

# Compter les séquences valides : WC/Wobble + Deux structures



**Paires de Bases valides** = autorise les paires **Wobble**



**Graphe de dépendance :**

Cycles + Chemins

i m n r

g - e - a - u    h    j - q

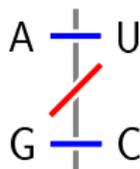
k - p    d - b - t

f - l - o - v    c - s

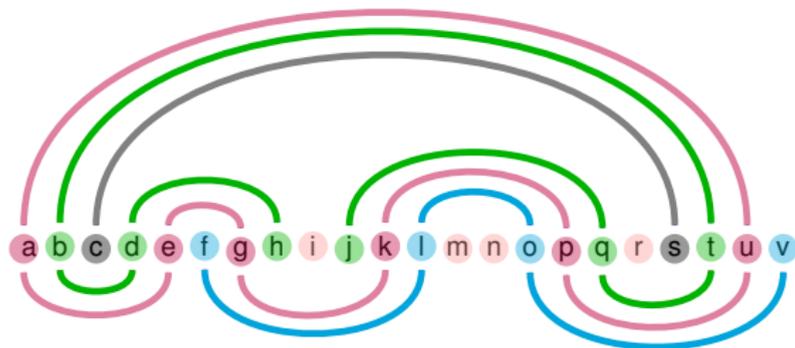
**Question :** Combien de séquences **valides** ?

**Réponses :**  $\neq \emptyset$  ! (graphes de dépendance et paires de bases tous deux **bipartis**)

# Compter les séquences valides : WC/Wobble + Deux structures



**Paires de Bases valides** = autorise les paires **Wobble**



**Graphe de dépendance :**

Cycles + Chemins

i m n r

g - e - a - u h j - q

k - p d - b - t

f - l - o - v c - s

**Question :** Combien de séquences **valides** ?

**Réponses :**  $\neq \emptyset$  ! (graphes de dépendance et paires de bases tous deux **bipartis**)

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

## Compter les séquences valides pour les chemins et cycles

$p(n)$  : Nombre de séquences valides pour un chemin de taille  $n$ .

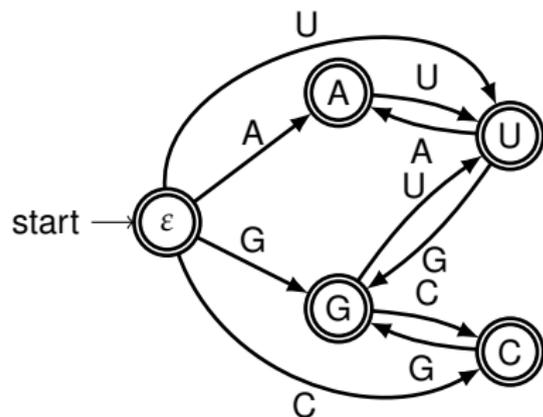
$c(n)$  : Nombre de séquences valides pour un cycle de taille  $n$ .

### Théorème (#Séquences valides pour chemins et cycles)

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

où  $\mathcal{F}_n$  est le  $n$ -ième nombre de Fibonacci.

**Pour les chemins :** Un simple automate ...



**Remarque :** Symétrie  $A \leftrightarrow C/G \leftrightarrow U$

## Compter les séquences valides pour les chemins et cycles

$p(n)$  : Nombre de séquences valides pour un chemin de taille  $n$ .

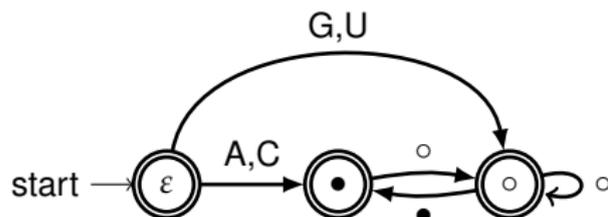
$c(n)$  : Nombre de séquences valides pour un cycle de taille  $n$ .

### Théorème (#Séquences valides pour chemins et cycles)

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

où  $\mathcal{F}_n$  est le  $n$ -ième nombre de Fibonacci.

**Pour les chemins :** Un simple automate ...



**Remarque :** Symétrie  $A \leftrightarrow C/G \leftrightarrow U$

## Compter les séquences valides pour les chemins et cycles

$p(n)$  : Nombre de séquences valides pour un chemin de taille  $n$ .

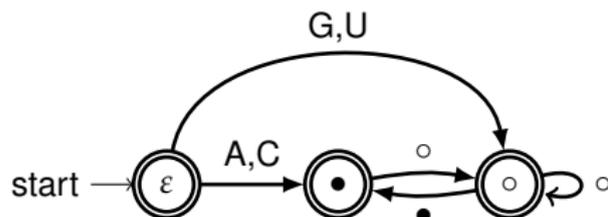
$c(n)$  : Nombre de séquences valides pour un cycle de taille  $n$ .

### Théorème (#Séquences valides pour chemins et cycles)

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

où  $\mathcal{F}_n$  est le  $n$ -ième nombre de Fibonacci.

**Pour les chemins :** Un simple automate ...



**Remarque :** Symétrie  $A \leftrightarrow C/G \leftrightarrow U$

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

# Compter les séquences valides pour les chemins et cycles

$p(n)$  : Nombre de séquences valides pour un chemin de taille  $n$ .

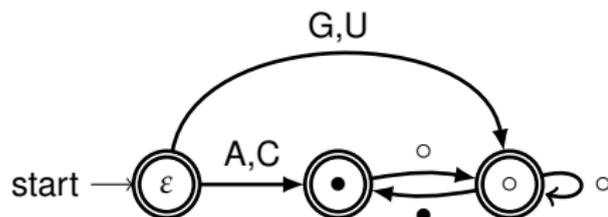
$c(n)$  : Nombre de séquences valides pour un cycle de taille  $n$ .

## Théorème (#Séquences valides pour chemins et cycles)

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

où  $\mathcal{F}_n$  est le  $n$ -ième nombre de Fibonacci.

**Pour les chemins :** Un simple automate ...



**Remarque :** Symétrie  $A \leftrightarrow C/G \leftrightarrow U$

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

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

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

$$= \mathcal{F}(n+2)$$

# Compter les séquences valides pour les chemins et cycles

$p(n)$  : Nombre de séquences valides pour un chemin de taille  $n$ .

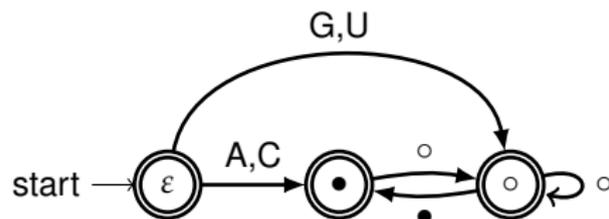
$c(n)$  : Nombre de séquences valides pour un cycle de taille  $n$ .

## Théorème (#Séquences valides pour chemins et cycles)

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

où  $\mathcal{F}_n$  est le  $n$ -ième nombre de Fibonacci.

**Pour les chemins :** Un simple automate ...



**Remarque :** Symétrie  $A \leftrightarrow C/G \leftrightarrow U$

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

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

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

$$= \mathcal{F}(n+2)$$

(Car  $m_{\circ}(0) = 1$  and  $m_{\circ}(1) = 2$ )

# Compter les séquences valides pour les chemins et cycles

$p(n)$  : Nombre de séquences valides pour un chemin de taille  $n$ .

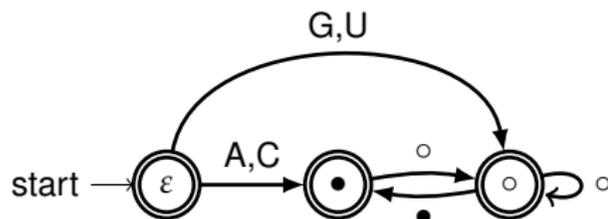
$c(n)$  : Nombre de séquences valides pour un cycle de taille  $n$ .

## Théorème (#Séquences valides pour chemins et cycles)

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

où  $\mathcal{F}_n$  est le  $n$ -ième nombre de Fibonacci.

**Pour les chemins :** Un simple automate ...



**Remarque :** Symétrie  $A \leftrightarrow C / G \leftrightarrow U$

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

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

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

$$= \mathcal{F}(n+2)$$

(Car  $m_{\circ}(0) = 1$  and  $m_{\circ}(1) = 2$ )

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

# Compter les séquences valides pour les chemins et cycles

$p(n)$  : Nombre de séquences valides pour un chemin de taille  $n$ .

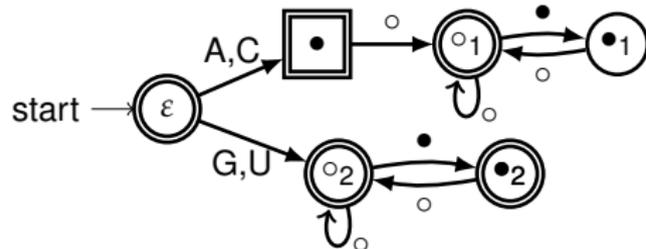
$c(n)$  : Nombre de séquences valides pour un cycle de taille  $n$ .

## Théorème (#Séquences valides pour chemins et cycles)

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

où  $\mathcal{F}_n$  est le  $n$ -ième nombre de Fibonacci.

Pour les cycles : à peine plus compliqué ...



**Remarque :** Symétrie  $A \leftrightarrow C/G \leftrightarrow U$

$$m_{o_2}(n) = \mathcal{F}(n+2)$$

$$m_{o_1}(n) = \mathcal{F}(n+1)$$

(Car  $m_{o_1}(0) = 1$  and  $m_{o_1}(1) = 1$ )

$$\begin{aligned} c(n) &:= m_{\epsilon}(n) = 2 m_{o_1}(n-2) + 2 m_{o_2}(n-1) \\ &= 2(\mathcal{F}(n-1) + \mathcal{F}(n+1)) = 2 \mathcal{F}(n) + 4 \mathcal{F}(n-1) \end{aligned}$$

## Compter les séquences valides pour deux structures

$p(n)$  : Nombre de séquences valides pour un chemin de taille  $n$ .

$c(n)$  : Nombre de séquences valides pour un cycle de taille  $n$ .

### Théorème (#Séquences valides pour chemins et cycles)

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

où  $\mathcal{F}_n$  est le  $n$ -ième nombre de Fibonacci.

$G$  : Graphe de dépendance issu de 2 structures (degré max  $\leq 2$ ).

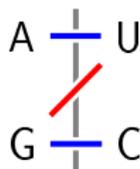
$G$  décomposé en  $\mathcal{P}(G)$  chemins et  $\mathcal{C}(G)$  cycles.

### Théorème (#Séquences valides pour les graphes de 2-structures)

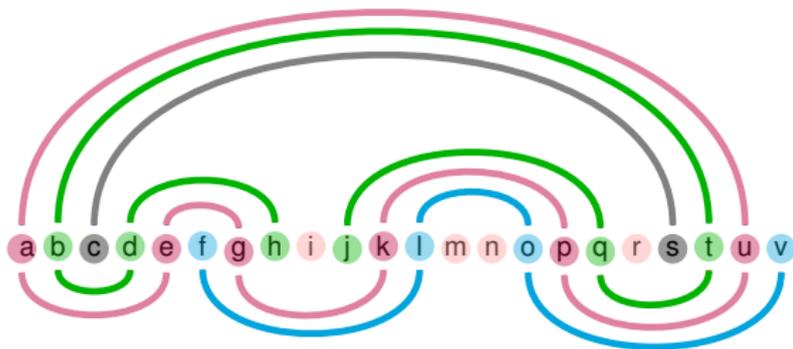
Le nombre  $\#\text{Designs}(G)$  de séquences valides pour  $G$  est

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

# Compter les séquences valides : WC/Wobble + Deux structures



**Paires de Bases valides** = autorise les paires **Wobble**



**Graphe de dépendance :**

Cycles + Chemins

i m n r

g-e-a-u h j-q

k-p d-b-t

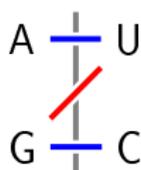
f-l-o-v c-s

**Question :** Combien de séquences **valides** ?

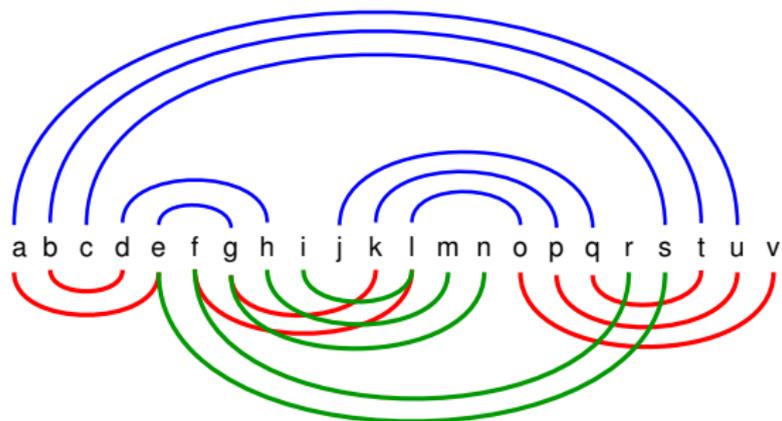
**Réponse :**  $\neq \emptyset$  ! (graphes PB et dépendance **bipartis**)

$$\# \text{Designs}(G) = \prod_{c \in \text{CC}(G)} \# \text{Designs}(cc) = 2\,322\,432$$

# Compter les séquences valides : Watson-Crick + > 2 structures

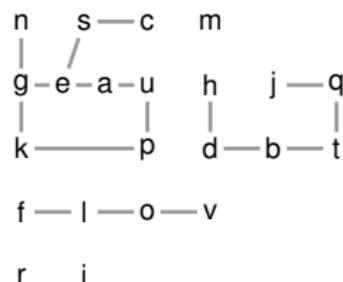


**Paires de Bases valides** = autorise les paires **Wobble**



**Graphe de dépendance :**

Cycles, Chemins, Arbres...



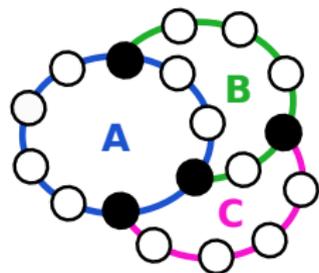
**Question :** Combien de séquences **valides** ?

**Réponse :** Non-biparti  $\rightarrow \emptyset$  ; Biparti  $\rightarrow$  ????

# Etat de l'art

## Approche de Abfalter/Flamm/Stadler 2003 :

- Décomposition *en oreille* du graphe [Whitney 1932]  
*Epluchage du graphe*;



- **Programmation dynamique** : comptage des #chemins valides pour chaque composante, étant donné des valeurs pour les **ancrages** (noeuds noirs) ;
- **Recombinaison** les valeurs pour calculer le nombre de séquences valides.

**Complexité** :  $\Theta(n \cdot 4^A)$  où  $A = \# \text{Max. d'ancrages}$ ,  $A \in \Theta(n)$  dans le cas au pire.

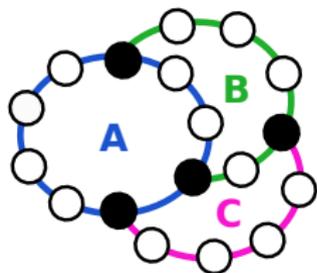
## Quelques réactions :

- Est ce bien *optimal*? Autres algorithmes et paramètres ?
- Quelles extensions possibles ? Génération pondérée (Boltzmann multidim.)
- Est ce bien nécessaire ? **Probablement** car comptage #P-difficile

# Etat de l'art

## Approche de Abfalter/Flamm/Stadler 2003 :

- Décomposition *en oreille* du graphe [Whitney 1932]  
*Epluchage du graphe*;



- **Programmation dynamique** : comptage des #chemins valides pour chaque composante, étant donné des valeurs pour les **ancrages** (noeuds noirs) ;
- **Recombinaison** les valeurs pour calculer le nombre de séquences valides.

**Complexité** :  $\Theta(n \cdot 4^A)$  où  $A = \# \text{Max. d'ancrages}$ ,  $A \in \Theta(n)$  dans le cas au pire.

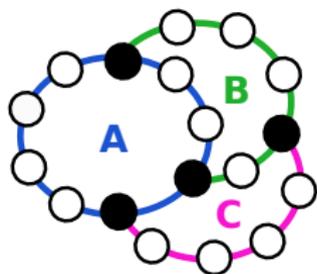
## Quelques réactions :

- Est ce bien *optimal*? Autres algorithmes et paramètres ?
- Quelles extensions possibles ? Génération pondérée (Boltzmann multidim.)
- Est ce bien nécessaire ? **Probablement** car comptage #P-difficile

# Etat de l'art

## Approche de Abfalter/Flamm/Stadler 2003 :

- Décomposition *en oreille* du graphe [Whitney 1932]  
*Epluchage du graphe*;



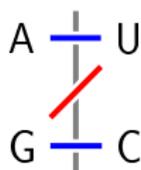
- **Programmation dynamique** : comptage des #chemins valides pour chaque composante, étant donné des valeurs pour les **ancrages** (noeuds noirs) ;
- **Recombinaison** les valeurs pour calculer le nombre de séquences valides.

**Complexité** :  $\Theta(n \cdot 4^A)$  où  $A = \# \text{Max. d'ancrages}$ ,  $A \in \Theta(n)$  dans le cas au pire.

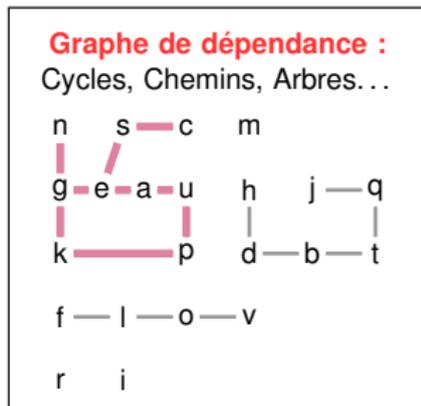
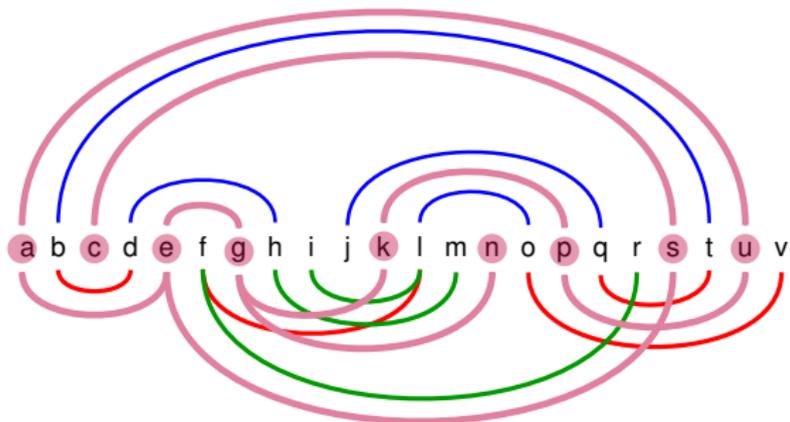
## Quelques réactions :

- Est ce bien *optimal*? Autres algorithmes et paramètres ?
- Quelles extensions possibles ? Génération pondérée (Boltzmann multidim.)
- Est ce bien nécessaire ? **Probablement** car comptage #P-difficile

# Compter les séquences valides : Watson-Crick + > 2 structures



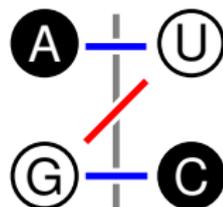
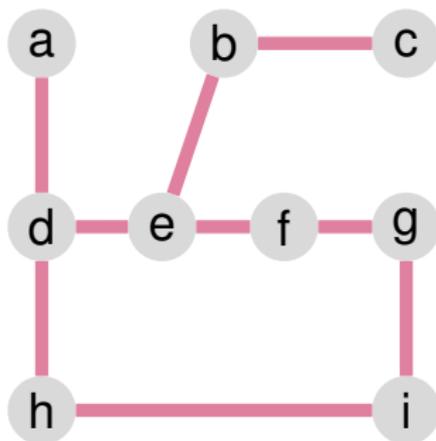
**Paires de Bases valides** = autorise les paires **Wobble**



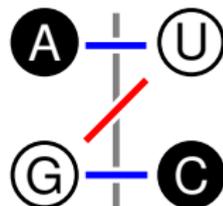
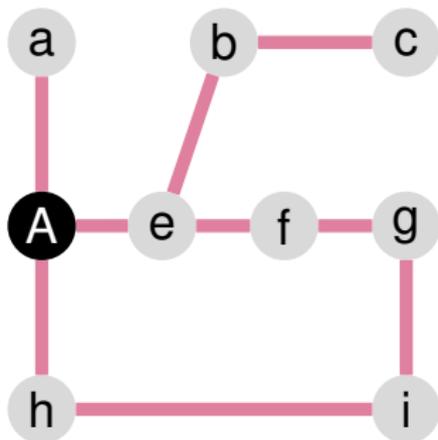
**Question :** Combien de séquences **valides** ?

**Réponse :** Non-biparti  $\rightarrow \emptyset$  ; Biparti  $\rightarrow$  ????

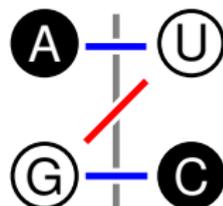
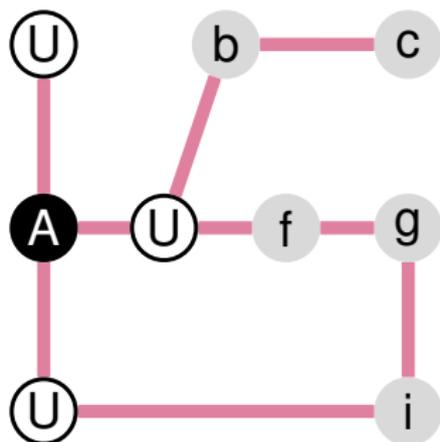
## Bijection entre stables d'un graphe biparti et séquences valides



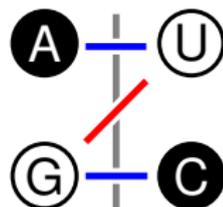
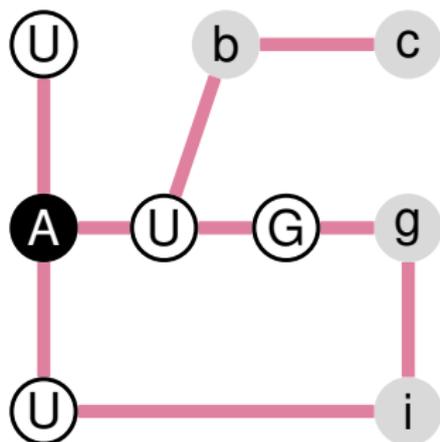
## Bijection entre stables d'un graphe biparti et séquences valides



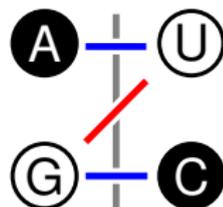
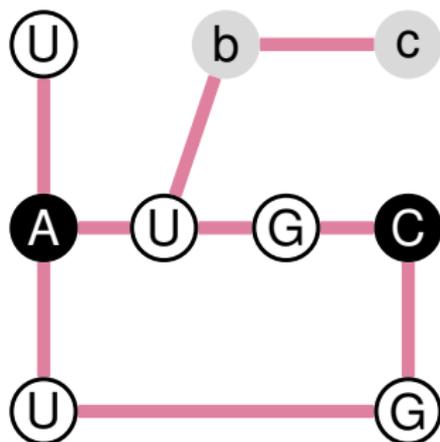
## Bijection entre stables d'un graphe biparti et séquences valides



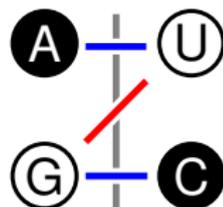
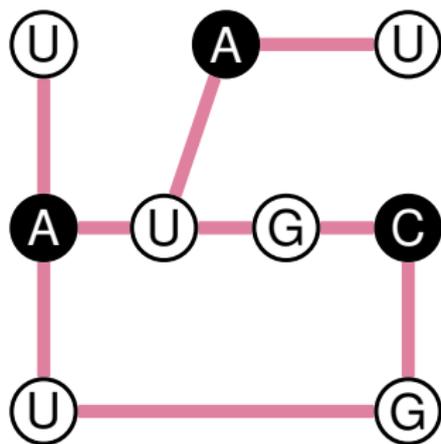
## Bijection entre stables d'un graphe biparti et séquences valides



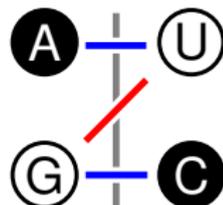
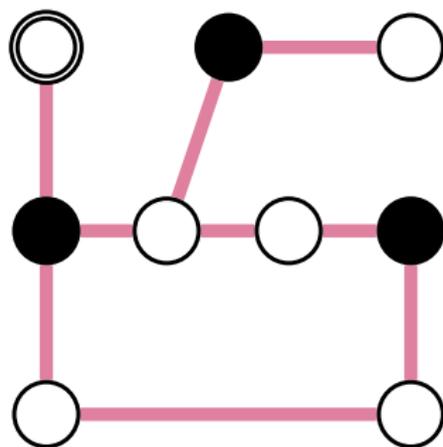
## Bijection entre stables d'un graphe biparti et séquences valides



## Bijection entre stables d'un graphe biparti et séquences valides



## Bijection entre stables d'un graphe biparti et séquences valides

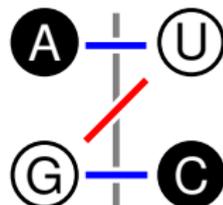
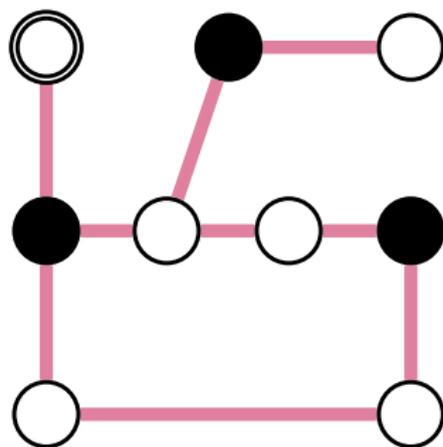


**Remarque :** Lettres noires non-adjacentes dans séquences valides

À symétrie triviale près\* (par ex. position  $\nearrow \in \{U, C\}$ ) :

$$\text{Designs}^*(cc) \subseteq \text{Stables}(cc)$$

## Bijection entre stables d'un graphe biparti et séquences valides



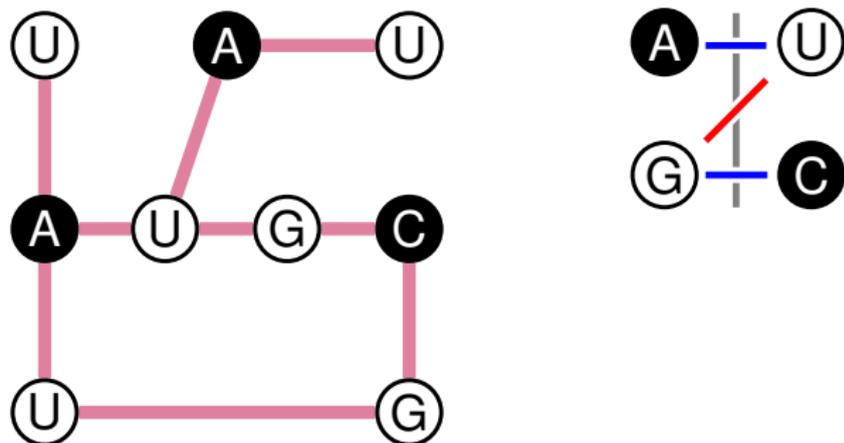
**Remarque :** Lettres noires non-adjacentes dans séquences valides

À symétrie triviale près\* (par ex. position  $\nwarrow \in \{U, C\}$ ) :

$$\text{Designs}^*(cc) \subseteq \text{Stables}(cc)$$

Par ailleurs, Stable (en noir) +  $\nwarrow$  vert.  $\in \{U, C\} \Rightarrow$  Séquence valide

## Bijection entre stables d'un graphe biparti et séquences valides



**Remarque :** Lettres noires non-adjacentes dans séquences valides

À symétrie triviale près\* (par ex. position  $\nwarrow \in \{U, C\}$ ) :

$$\text{Designs}^*(cc) \subseteq \text{Stables}(cc)$$

Par ailleurs, Stable (en noir) +  $\nwarrow$  vert.  $\in \{U, C\} \Rightarrow$  Séquence valide

$\Rightarrow$  Bijection entre  $\text{Designs}^*(cc)$  et  $\text{Stables}(cc)$ .

## Séquences valides et stables d'un graphe biparti

### Théorème (#Séquences valides et stables dans les graphes connexes bipartis)

Soit  $G$  un graphe de dépendance **biparti et connexe** :

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

Pour un graphe de dépendance **biparti**  $G$ , on a :

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

**Mais**  $\#Stables(G)$  est **#P-difficile** sur les graphes bipartis (**#BIS**) [Bubbley&Dyer'01]  
(+ Tout  $G$  est graphe de dépendance de quelque famille de structures)

**Donc**  $\exists$  algo.  $\mathcal{A} \in P$  pour  $\#Designs(G) \Rightarrow \exists$  algo.  $\mathcal{A}' \in P$  pour **#BIS**...

### Théorème

*$\#Designs$  est #P-hard.*

Pas d'algorithme polynomial pour  $\#Designs(G)$  **sauf si**  $\#P = FP (\Rightarrow P = NP)$

## Séquences valides et stables d'un graphe biparti

### Théorème (#Séquences valides et stables dans les graphes connexes bipartis)

Soit  $G$  un graphe de dépendance **biparti et connexe** :

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

Pour un graphe de dépendance **biparti**  $G$ , on a :

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

**Mais**  $\#Stables(G)$  est **#P-difficile** sur les graphes bipartis (**#BIS**) [Bubbly&Dyer'01]  
(+ Tout  $G$  est graphe de dépendance de quelque famille de structures)

**Donc**  $\exists$  algo.  $\mathcal{A} \in P$  pour  $\#Designs(G) \Rightarrow \exists$  algo.  $\mathcal{A}' \in P$  pour **#BIS**...

### Théorème

*#Designs est #P-hard.*

Pas d'algorithme polynomial pour  $\#Designs(G)$  **sauf si**  $\#P = FP (\Rightarrow P = NP)$

## Séquences valides et stables d'un graphe biparti

### Théorème (#Séquences valides et stables dans les graphes connexes bipartis)

Soit  $G$  un graphe de dépendance **biparti et connexe** :

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

Pour un graphe de dépendance **biparti**  $G$ , on a :

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

**Mais**  $\#Stables(G)$  est **#P-difficile** sur les graphes bipartis ( $\#BIS$ ) [Bubbly&Dyer'01]  
(+ Tout  $G$  est graphe de dépendance de quelque famille de structures)

**Donc**  $\exists$  algo.  $\mathcal{A} \in P$  pour  $\#Designs(G) \Rightarrow \exists$  algo.  $\mathcal{A}' \in P$  pour  $\#BIS \dots$

Théorème

*#Designs est #P-hard.*

Pas d'algorithme polynomial pour  $\#Designs(G)$  **sauf si**  $\#P = FP (\Rightarrow P = NP)$

## Séquences valides et stables d'un graphe biparti

### Théorème (#Séquences valides et stables dans les graphes connexes bipartis)

Soit  $G$  un graphe de dépendance **biparti et connexe** :

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

Pour un graphe de dépendance **biparti**  $G$ , on a :

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

**Mais**  $\#Stables(G)$  est **#P-difficile** sur les graphes bipartis ( $\#BIS$ ) [Bubbly&Dyer'01]  
(+ Tout  $G$  est graphe de dépendance de quelque famille de structures)

**Donc**  $\exists$  algo.  $\mathcal{A} \in P$  pour  $\#Designs(G) \Rightarrow \exists$  algo.  $\mathcal{A}' \in P$  pour  $\#BIS$ ...

### Théorème

*$\#Designs$  est #P-hard.*

Pas d'algorithme polynomial pour  $\#Designs(G)$  **sauf si**  $\#P = FP$  ( $\Rightarrow P = NP$ )

# Conséquences

## Corollaire (#Approximabilité pour $\leq 5$ structures) [Weitz'06]

Pour tout graphe  $G$  construit à partir de  $\leq 5$  structures (avec croisements),  $\#Design(G)$  peut être approché dans **n'importe quel ratio** en **temps polynomial** (PTAS)

## Corollaire ??? (#BIS-difficulté pour $> 5$ structures) [Cal, Galanis *et al*'16]

À partir de 5 structures (avec croisements), approcher  $\#Design$  est aussi **aussi difficile que** d'approcher  $\#BIS$  sans contrainte.

**Pourquoi des croisements ?** Parce que tout graphe biparti de degré max  $\Delta$  peut être décomposé en  $\Delta$  couplages en **temps polynomial** (Théorème de Vizing).

Enfin, forte connexion entre **comptage** et **échantillonnage** [Jerrum/Valiant/Vazirani'86].

## Conjecture (#BIS-difficulté de l'échantillonnage)

**Génération quasi-uniforme** aussi difficile qu'approcher  $\#BIS$  en général

⇒ **Génération aléatoire uniforme**  $\#P$  difficile ?

## Conséquences

### Corollaire (#Approximabilité pour $\leq 5$ structures) [Weitz'06]

Pour tout graphe  $G$  construit à partir de  $\leq 5$  structures (avec croisements), #Design( $G$ ) peut être approché dans **n'importe quel ratio** en **temps polynomial** (PTAS)

### Corollaire ??? (#BIS-difficulté pour $> 5$ structures) [Cai, Galanis *et al*'16]

À partir de 5 structures (avec croisements), approcher #Design est aussi **aussi difficile que** d'approcher #BIS sans contrainte.

**Pourquoi des croisements ?** Parce que tout graphe biparti de degré max  $\Delta$  peut être **décomposé** en  $\Delta$  couplages en **temps polynomial** (Théorème de Vizing).

Enfin, forte connexion entre **comptage** et **échantillonnage** [Jerrum/Valiant/Vazirani'86].

### Conjecture (#BIS-difficulté de l'échantillonnage)

**Génération quasi-uniforme** aussi difficile qu'approcher #BIS en général

⇒ **Génération aléatoire uniforme** #P difficile ?

# Conséquences

## Corollaire (#Approximabilité pour $\leq 5$ structures) [Weitz'06]

Pour tout graphe  $G$  construit à partir de  $\leq 5$  structures (avec croisements),  $\#Design(G)$  peut être approché dans **n'importe quel ratio** en **temps polynomial** (PTAS)

## Corollaire ??? (#BIS-difficulté pour $> 5$ structures) [Cai, Galanis *et al*'16]

À partir de 5 structures (avec croisements), approcher  $\#Design$  est aussi **aussi difficile que** d'approcher  $\#BIS$  sans contrainte.

**Pourquoi des croisements ?** Parce que tout graphe biparti de degré max  $\Delta$  peut être **décomposé** en  $\Delta$  couplages en **temps polynomial** (Théorème de Vizing).

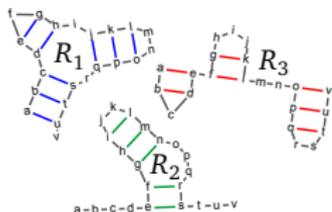
Enfin, forte connexion entre **comptage** et **échantillonnage** [Jerrum/Valiant/Vazirani'86].

## Conjecture (#BIS-difficulté de l'échantillonnage)

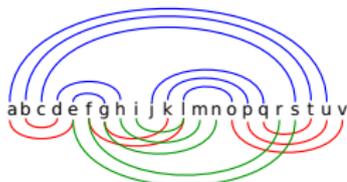
**Génération quasi-uniforme** aussi difficile qu'approcher  $\#BIS$  en général

⇒ **Génération aléatoire uniforme**  $\#P$  difficile ?

# Décomposition arborescente et génération de Boltzmann



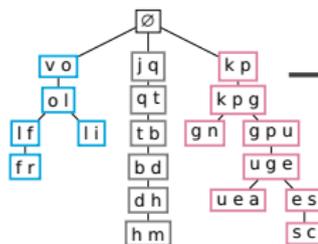
i) Input Structures



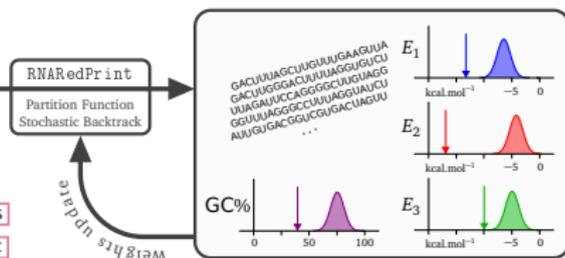
ii) Merged Base-Pairs



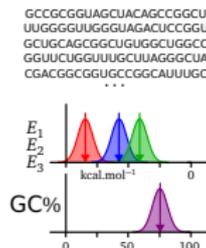
iii) Compatibility Graph



iv) Tree Decomposition



v) Weight Optimization (Adaptive Sampling)



vi) Final Designs

[Hammer/P/Wang/Will, RECOMB'18]

- **Algorithme** de complexité paramétrée sur la **largeur arborescente**
- Génération de Boltzmann multidimensionnelle pour contrôler l'énergie, GC%...

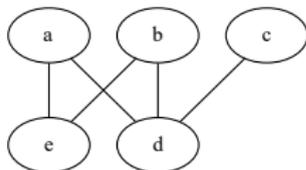
# Décomposition arborescente et comptage/génération

**Décomposition arborescente**  $T$  pour graphe  $G = (V, E) =$

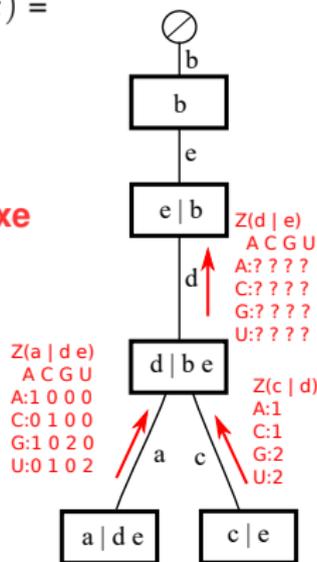
- 1 Noeuds  $B$  tels que  $\forall b \in B, b \subseteq V$
- 2  $\forall v \in V, \exists b \in B$  tel que  $v \in b$
- 3  $\forall (v, v') \in E, \exists b \in B$  tel que  $\{v, v'\} \subseteq b$
- 4 Arêtes telles que sous arbre induit par  $v \in V$  **connexe**

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

**Structures cibles**



**Graphe de dépendance**



**Décomp. arborescente**

$b = \{b_1, b_2 \dots\}$  : noeud de  $D$

$T_b$  : sous-arbre enraciné en  $b$

$w$  : **Largeur** de la décomposition  $D$

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

**Complexité**  $\Theta(nmk + nk2^w)$  pour **génération uniforme** de  $m$  séq. ( $k$  struct.)

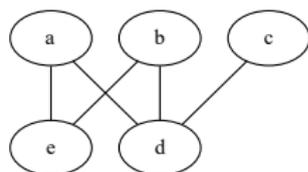
# Décomposition arborescente et comptage/génération

**Décomposition arborescente**  $T$  pour graphe  $G = (V, E) =$

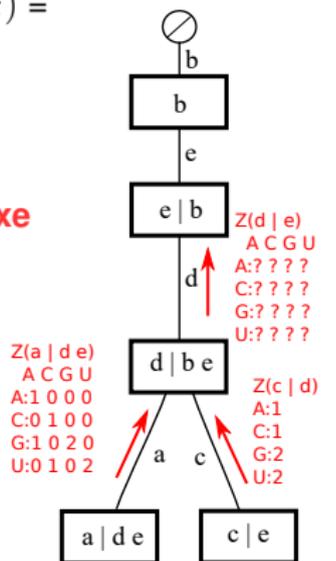
- 1 Noeuds  $B$  tels que  $\forall b \in B, b \subseteq V$
- 2  $\forall v \in V, \exists b \in B$  tel que  $v \in b$
- 3  $\forall (v, v') \in E, \exists b \in B$  tel que  $\{v, v'\} \subseteq B$
- 4 Arêtes telles que sous arbre induit par  $v \in V$  **connexe**

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

Structures cibles



Graphe de dépendance



Décomp. arborescente

$b = \{b_1, b_2 \dots\}$  : noeud de  $D$

$T_b$  : sous-arbre enraciné en  $b$

$w$  : **Largeur** de la décomposition  $D$

$$\mathcal{Z}(T_b | b_2 \leftarrow v_2 \dots) = \sum_{\substack{b_1 \leftarrow v_1 \\ v_1 \in \{A, C, G, U\}}} \prod_{c \text{ fils de } b} \mathcal{Z}(T_c | b_1 \leftarrow v_1, b_2 \leftarrow v_2 \dots)$$

**Complexité**  $\Theta(nmk + nk2^{w+\#CC})$  pour **génération uniforme** de  $m$  séq. ( $k$  struct.)

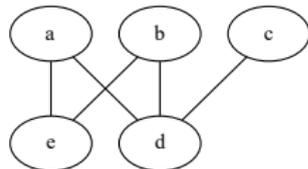
# Décomposition arborescente et comptage/génération

**Décomposition arborescente**  $T$  pour graphe  $G = (V, E) =$

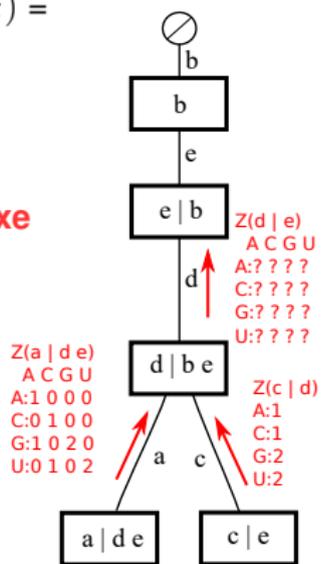
- 1 Noeuds  $B$  tels que  $\forall b \in B, b \subseteq V$
- 2  $\forall v \in V, \exists b \in B$  tel que  $v \in b$
- 3  $\forall (v, v') \in E, \exists b \in B$  tel que  $\{v, v'\} \subseteq b$
- 4 Arêtes telles que sous arbre induit par  $v \in V$  **connexe**

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

Structures cibles



Graphe de dépendance



Décomp. arborescente

$b = \{b_1, b_2 \dots\}$  : noeud de  $D$

$T_b$  : sous-arbre enraciné en  $b$

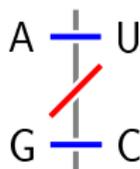
$w$  : **Largeur** de la décomposition  $D$

$$\mathcal{Z}(T_b | b_2 \leftarrow v_2 \dots) = \sum_{b_1 \leftarrow v_1} \prod_{i=1}^k x_i^{E(b, v_1, \dots)} \prod_{c \text{ fils de } b} \mathcal{Z}(T_c | b_1 \leftarrow v_1, b_2 \leftarrow v_2 \dots)$$

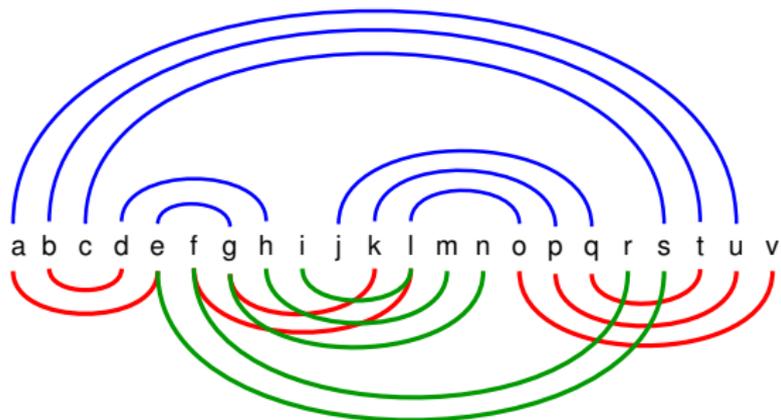
$v_1 \in \{A, C, G, U\}$

**Complexité**  $\Theta(nmk + nk2^w)$  pour **génération pondérée** de  $m$  séq. ( $k$  struct.)

## Compter les séquences valides : Watson-Crick + > 2 structures



**Paires de Bases valides** = autorise les paires **Wobble**



**Graphe de dépendance :**

Cycles, Chemins, Arbres...

n s—c m

g—e—a—u h j—q

| | | |  
k—p d—b—t

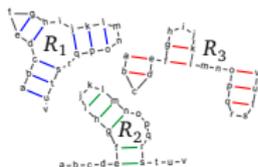
f—l—o—v

r i

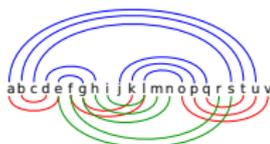
**Question :** Combien de séquences **valides** ?

**Réponse :** Biparti  $\rightarrow \prod_{cc \in CC(G)} 2 \times \#Stables(cc) = 496\,672$

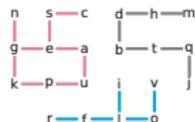
# Décomposition arborescente et génération de Boltzmann



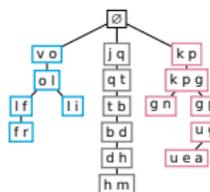
i) Input Structures



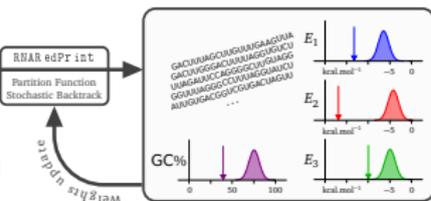
ii) Merged Base-Pairs



iii) Compatibility Graph

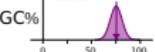
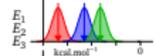


iv) Tree Decomposition



v) Weight Optimization (Adaptive Sampling)

```
GCCGCGGUAGCUACAGCCGGCU
UUGGGGUUGGGUAGUCUCGGGU
GGUCACGGCGCUUGGGUUGGCC
GGUUUCGUUUUCUUAAGGCUCA
CGAUCGCGGUUGCCGCAUUUGC
...
```



vi) Final Designs

[Hammer/P/Wang/Will, RECOMB'18]

- **Algorithme** de complexité paramétrée sur la **largeur arborescente**
- **Génération de Boltzmann multidim.** [Bodini/P. 2011] pour contrôler l'énergie ...

**Entrée :** Energie visées  $(E_\ell^*)_{\ell=1}^k$ , poids  $(x_\ell)_{\ell=1}^k$  tels que  $\mathbb{E}(E(w, S_\ell)) = E_\ell^*, \forall \ell$  :

$$\mathbb{P}(w \mid x_1 \dots x_k) \sim \prod_{\ell=1}^k x_\ell^{E(w, S_\ell)} + \text{Rejet efficace}$$

**A faire :** Itération numérique pour déterminer  $(x_\ell)_{\ell=1}^k$  (cf exposé Sergey)

**Merci !**

