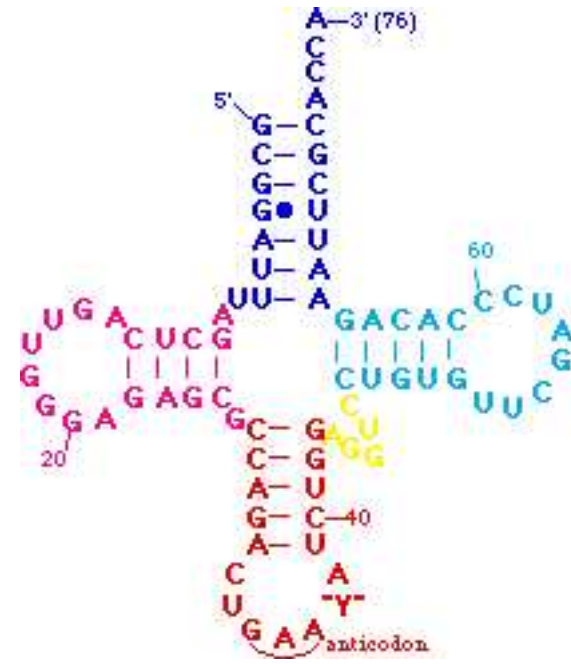
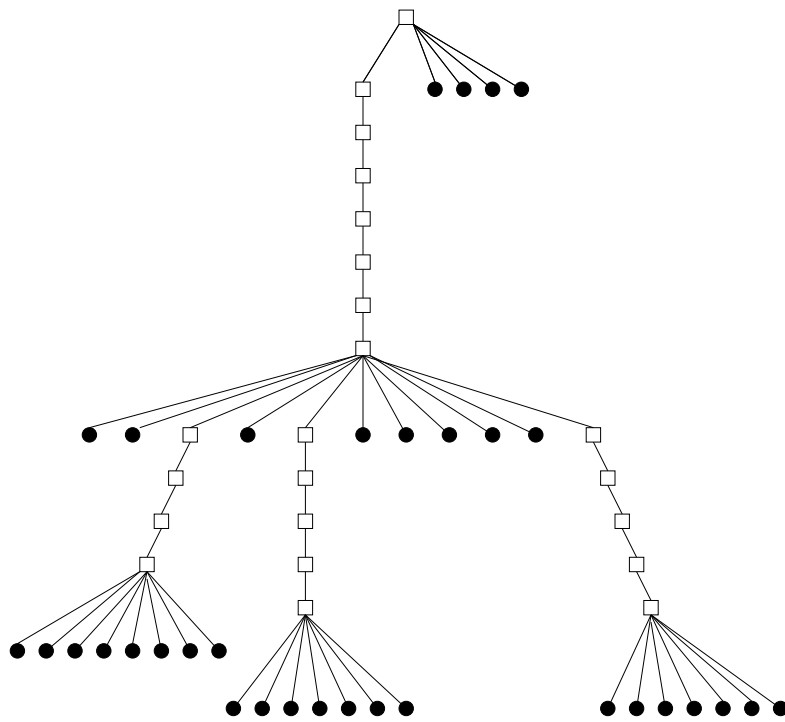


# **Algorithmes pour la comparaison d'ARN : distance d'édition entre arbres**

Serge Dulucq (LABRI)

Hélène Touzet (LIFL)



- : 2 bases appariées (tige)
- : base libre (boucle)

Représentation arborescente pour la structure secondaire de l'ARN de transfert

## Opérations d'édition

- ▷ substitution d'un nœud
- ▷ insertion d'un nœud
- ▷ déléation d'un nœud

# Distance d'édition sur les mots

A L G O R I T H M E  
P A P I L L O N

▷ Décomposition en partant de la gauche

*substitution*

A	LGORITHME
P	APILLON

*insertion*

-	ALGORITHMES
P	APILLON

*délétion*

A	LGORITHME
-	PAPILLON

Appel sur tous les suffixes

▷ Décomposition en partant de la droite

● *substitution*

ALGORITHM	E
PAPILLO	N

● *insertion*

ALGORITHMME	-
PAPILLO	N

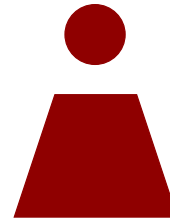
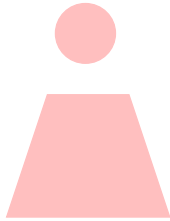
● *délétion*

ALGORITHM	E
PAPILLON	-

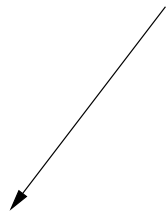
Appel sur tous les préfixes

# Distance d'édition sur les arbres

$l(f)$



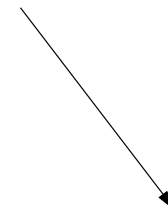
$l'(f')$



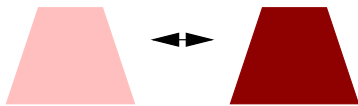
Substitution



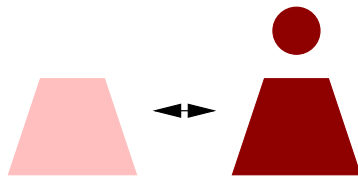
Délétion



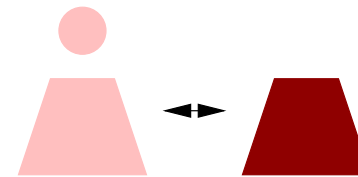
Insertion



$\text{Distance}(f, f')$   
+  
 $\text{sub}(l, l')$

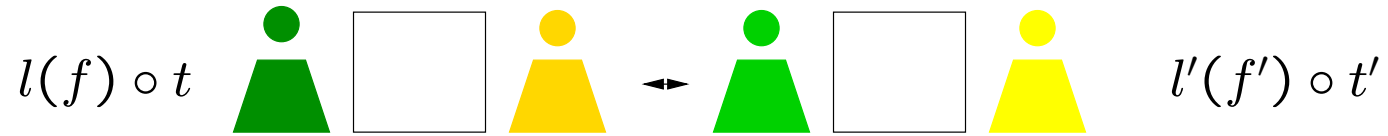


$\text{Distance}(f, l'(f'))$   
+  
 $\text{del}(l)$



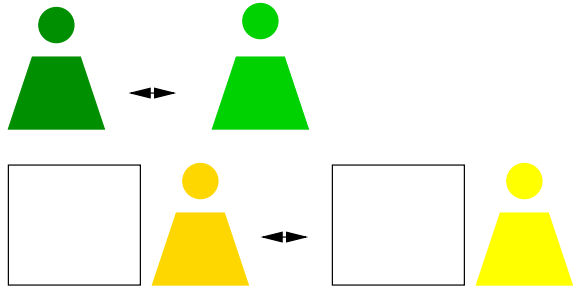
$\text{Distance}(l(f), f')$   
+  
 $\text{ins}(l')$

# Distance d'édition sur les forêts I



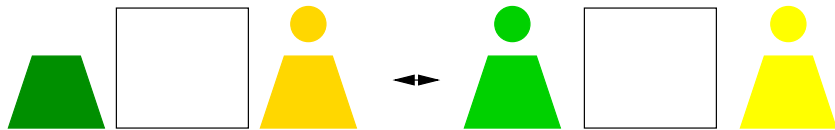
## Décomposition gauche

Substitution



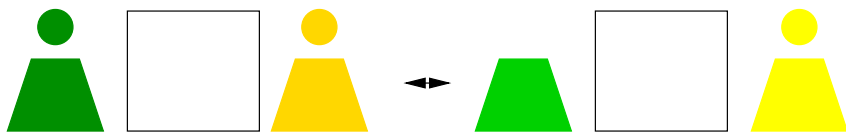
$$\text{Distance}(l(f), l'(f')) + \text{Distance}(t, t')$$

Délétion



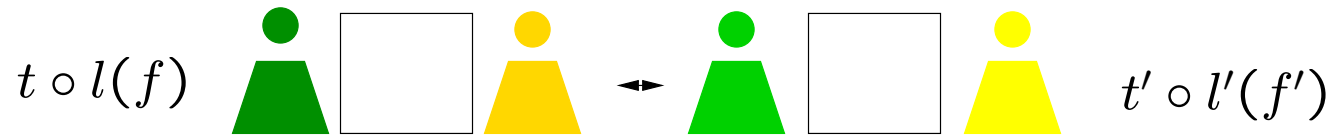
$$\text{Distance}(f \circ t, l'(f') \circ t') + \text{del}(l)$$

Insertion



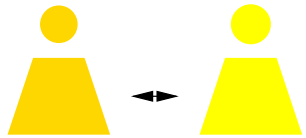
$$\text{Distance}(l(f) \circ t, f' \circ t') + \text{ins}(l')$$

# Distance d'édition sur les forêts II



## Décomposition droite

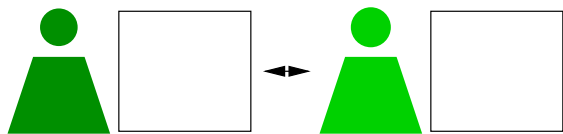
Substitution



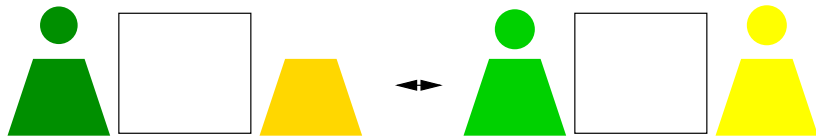
Distance( $l(f), l'(f')$ )

+

Distance( $t, t'$ )

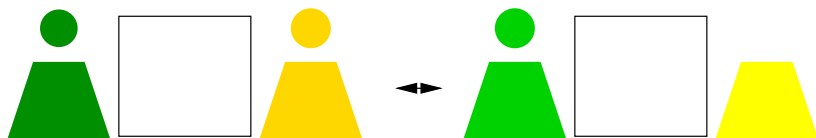


Délétion



Distance( $t \circ f, l'(f') \circ l'(f')$ ) + del( $l$ )

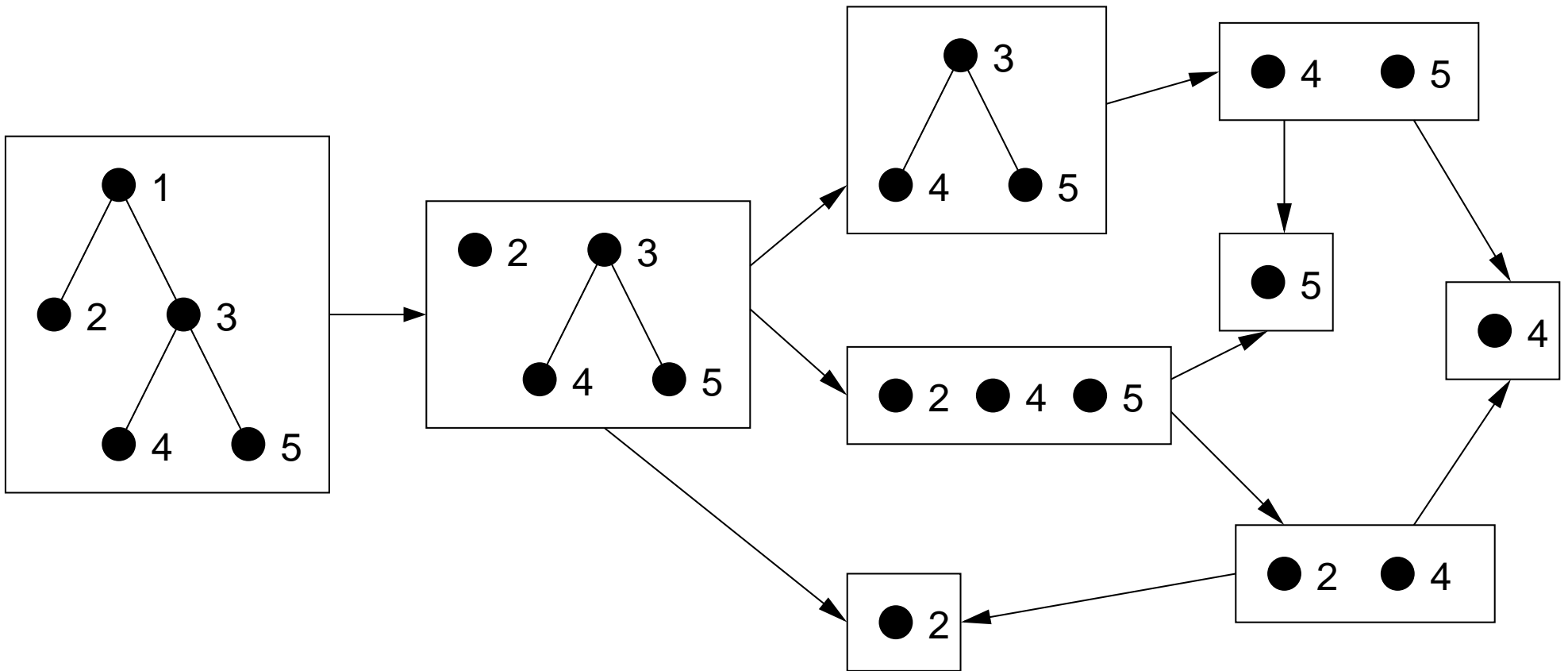
Insertion



Distance( $t \circ l(f), t' \circ t'$ ) + ins( $l'$ )

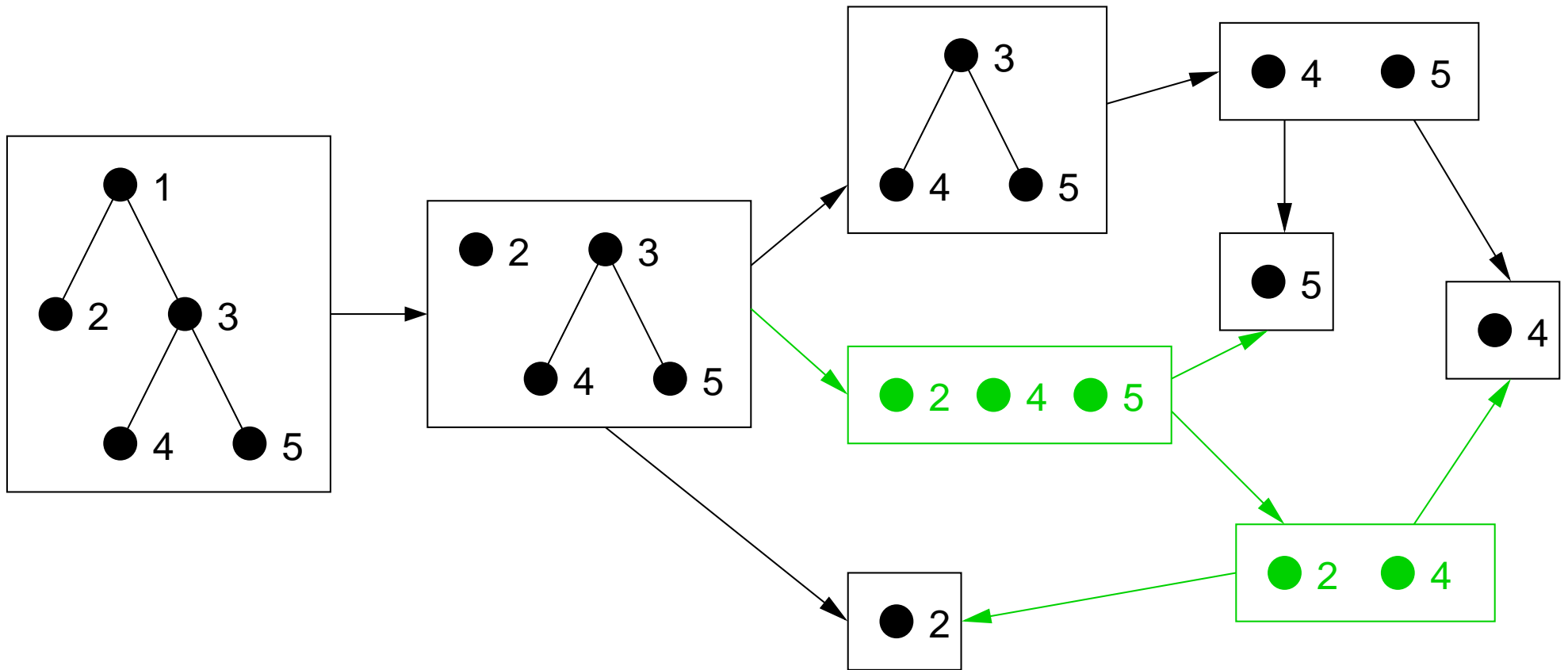


# Droite/gauche pour les arbres ?



décomposition droite

# Droite/gauche pour les arbres ?



décomposition gauche

## Peigne à $n$ branches ( $2n + 1$ nœuds)

▷ Décomposition gauche :

$$3n + 1 \text{ sous-forêts}$$

▷ Décomposition droite :

$$2n^2 + 1 \text{ sous-forêts}$$

# Stratégie de décomposition

▷ succession de choix **droite** ou **gauche**

▷  $S : forêt \times forêt \rightarrow \{droite, gauche\}$

▷ **Zhang & Shasha** (1989) :

$(f, g) \rightarrow gauche$

▷ **Klein** (1998) :

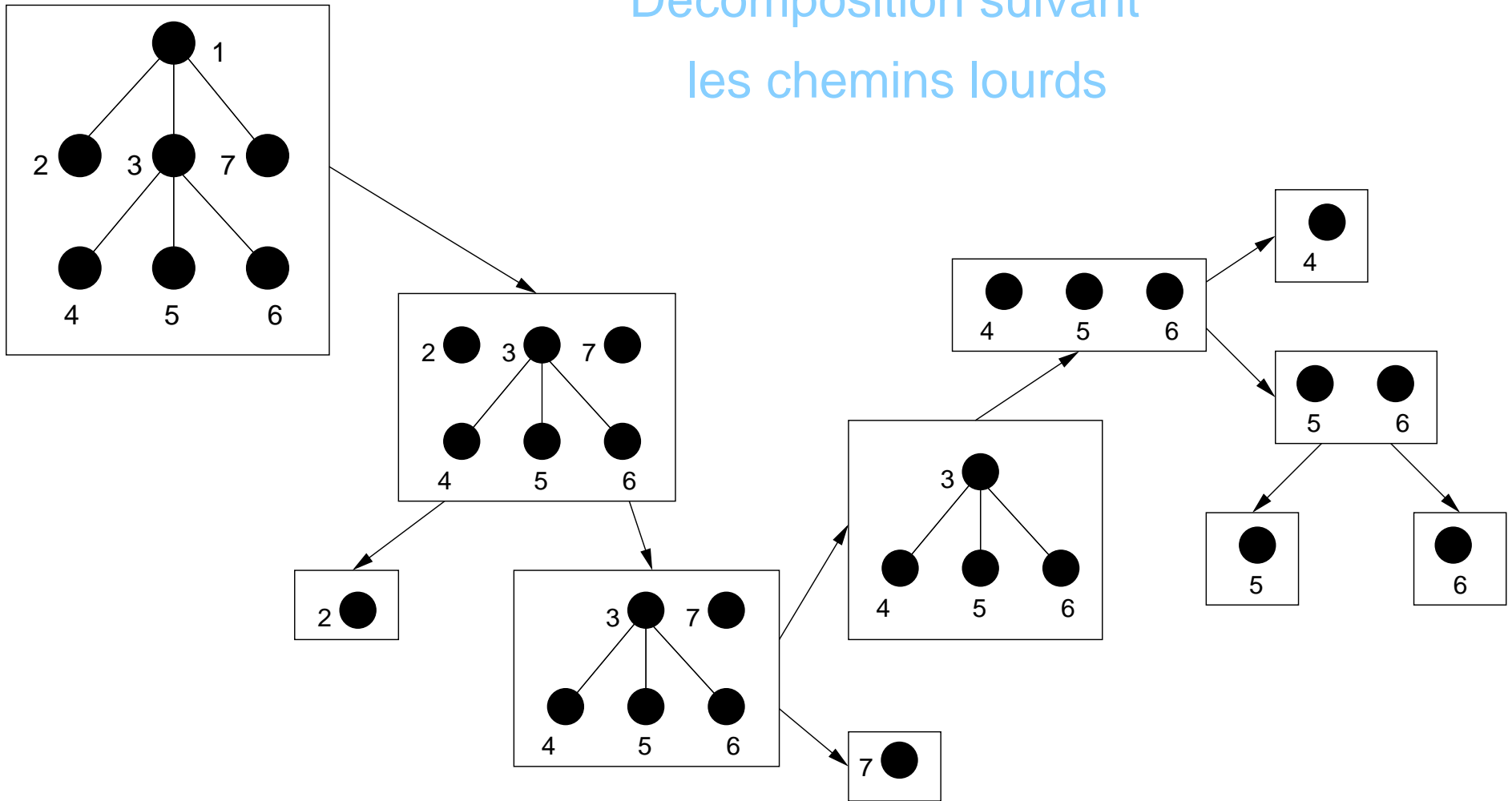
$(f, g) \rightarrow$  **gauche** si le premier nœud de  $f$  est  
sur le chemin lourd

**droite**, sinon

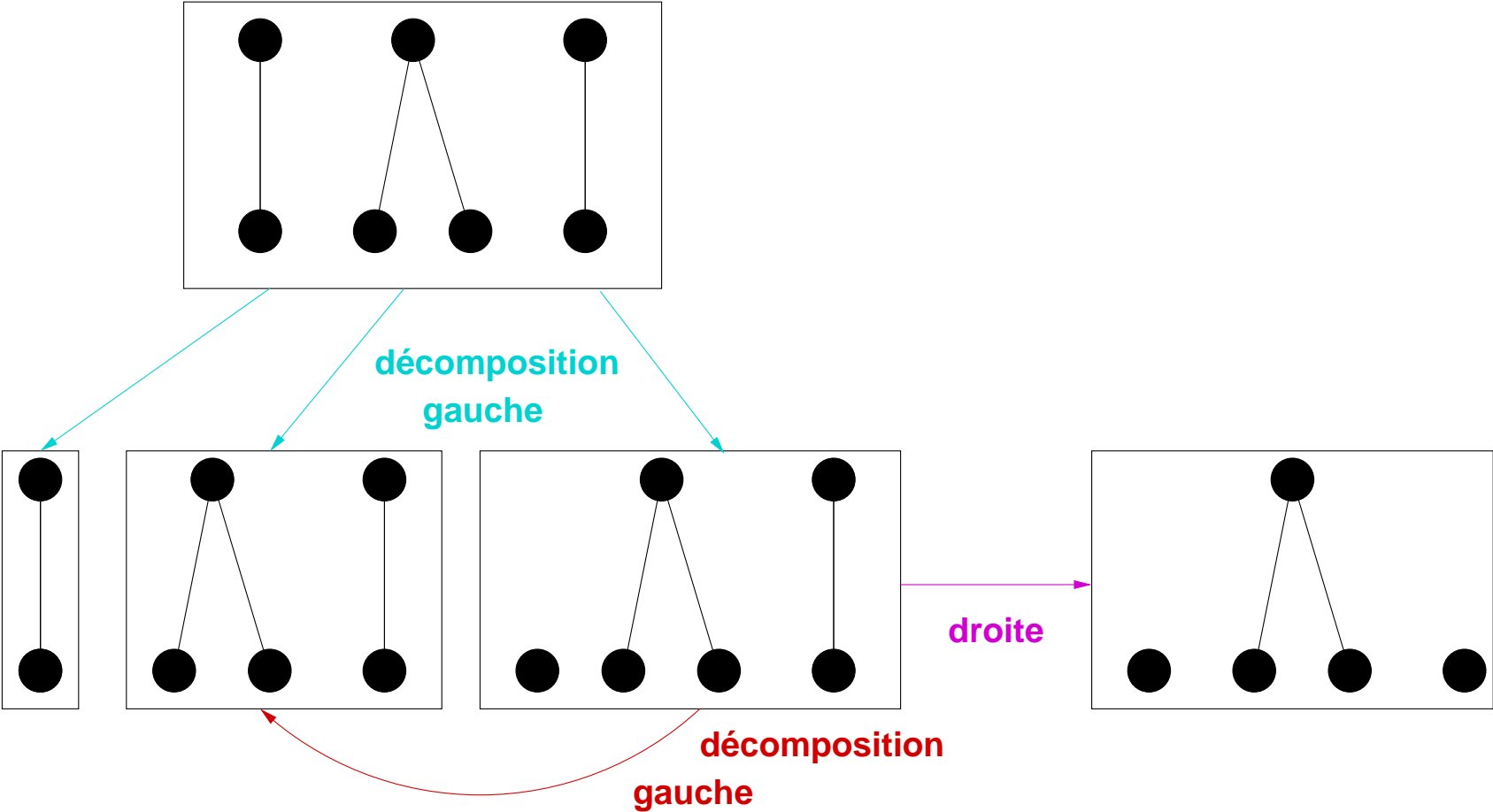


# Stratégie de Klein

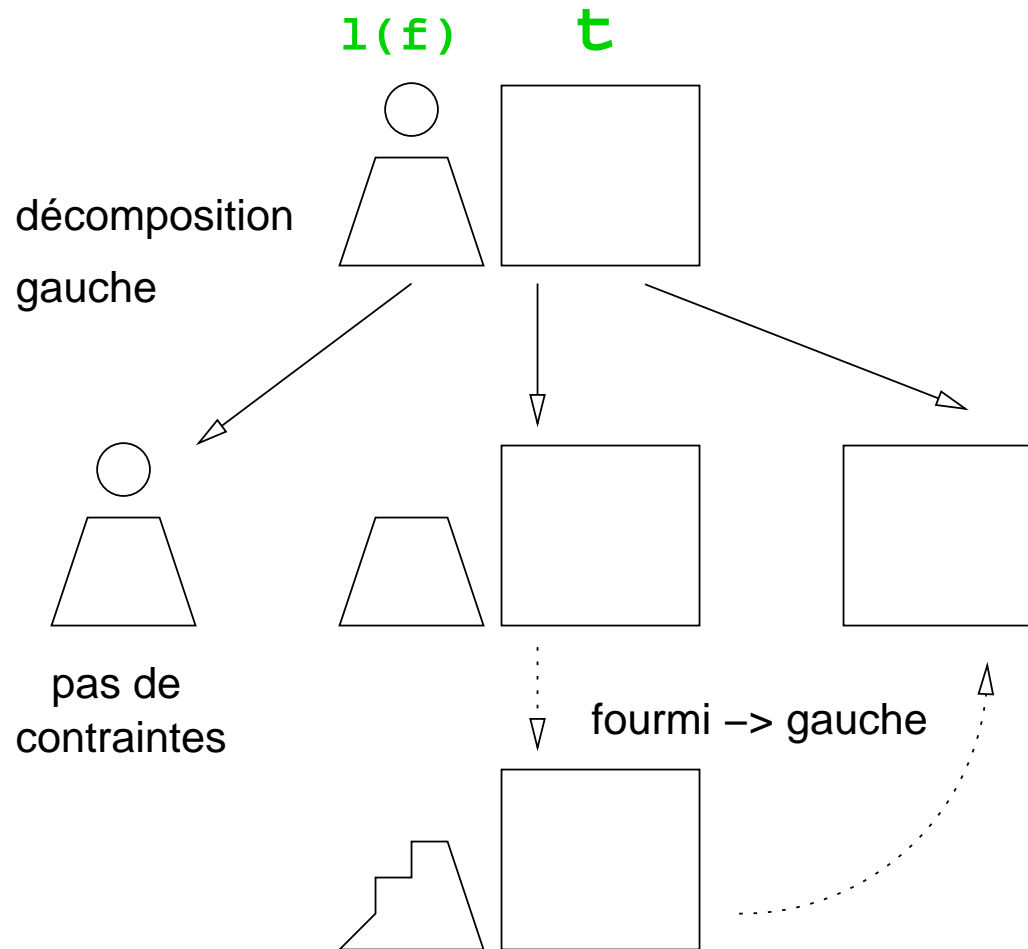
Décomposition suivant  
les chemins lourds



# Exemple de stratégie dispendieuse



# Les stratégies fourmis



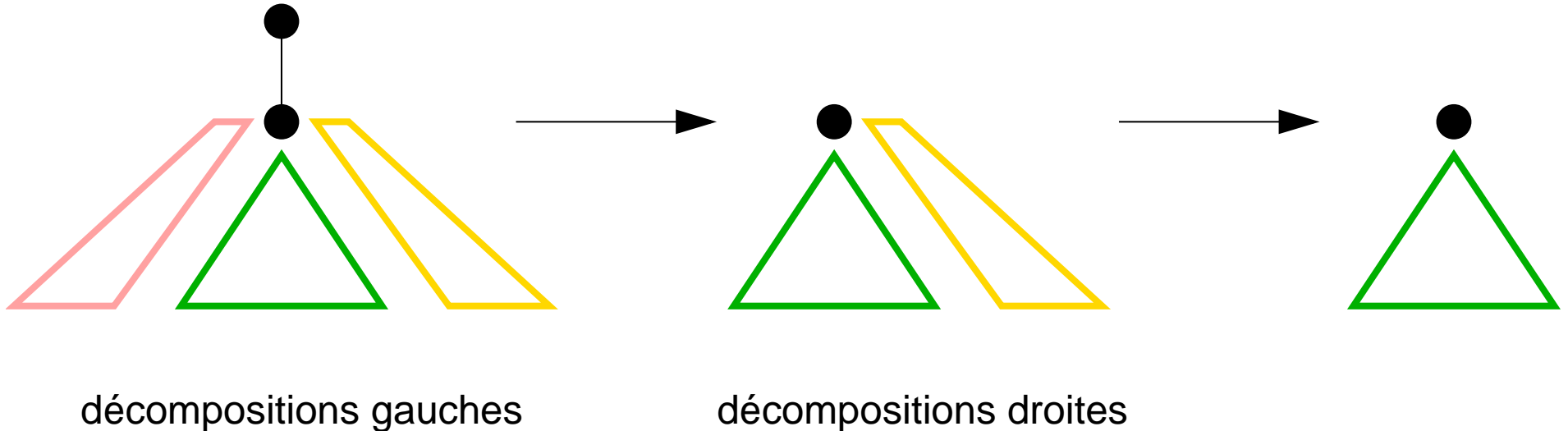
$\#sous\text{-forets}(l(f)t) =$

$\#sous\text{-forets}(l(f)) + |l(f)| + \#sous\text{-forets}(t)$



# Définir une stratégie fourmi à partir d'un recouvrement

- ▷  $\phi(i) \in \{droite, gauche\}$  si  $i$  est un noeud d'arité 0 ou 1
- ▷  $\phi(i)$  est un fils de  $i$  sinon : le fils **favori**



# Recouvrement pour Zhang et Klein

## ▷ Zhang :

- le fils favoris est le fils droit
- la direction est gauche

Zhang est fourni pour les deux arbres

## ▷ Klein :

- le fils favori est le fils de taille maximale
- la direction est gauche

Klein est fourni pour le premier arbre

# Nombre de sous-forêts sur un arbre

$$A = l(A_1 \circ \dots \circ A_n)$$

▷ pour une stratégie quelconque

meilleur des cas

$$|A| - |A_j| + \#\text{sous-forêts}(A_1) + \dots + \#\text{sous-forêts}(A_n) \quad O(n \log(n))$$

$A_j$  est le fils de taille maximale

pire des cas

$$\frac{n(n+3)}{2} - \sum_{i \in A} |A(i)|$$

$$\frac{1}{2} n^2 + \frac{\sqrt{\pi}}{2} n^{\frac{3}{2}} + O(n) \text{ en moyenne}$$

▷ pour un arbre muni d'un recouvrement

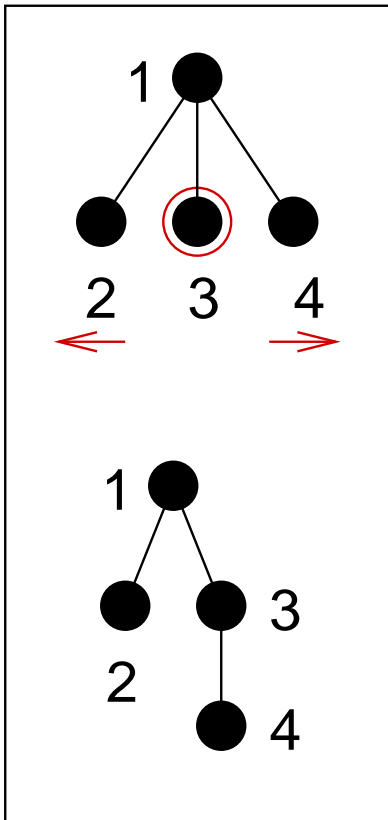
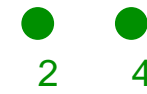
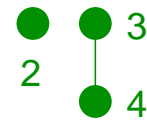
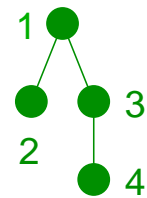
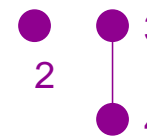
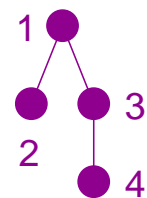
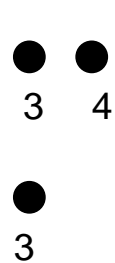
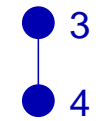
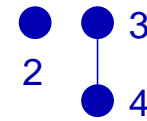
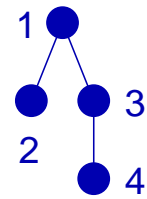
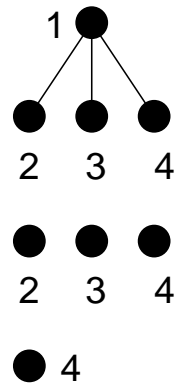
$$\#sous\text{-forêts}(A) = |A| - |A_j| + \#sous\text{-forêts}(A_1) + \dots + \#sous\text{-forêts}(A_n)$$

$A_j$  est le fils favori

**Conséquence 1:** nombre exact de sous-forêts pour Zhang (forêts droites)

**Conséquence 2:** Klein est optimal pour l'arbre muni du recouvrement

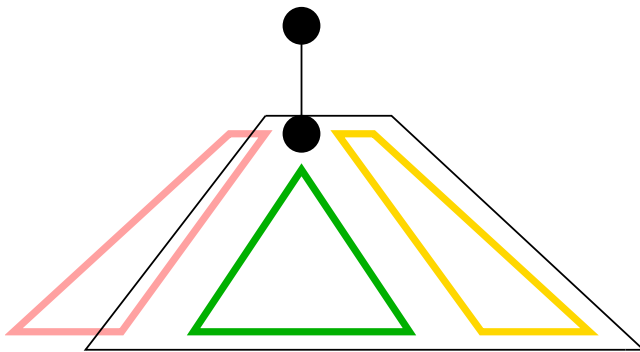
# Que se passe-t-il pour l'autre arbre ?



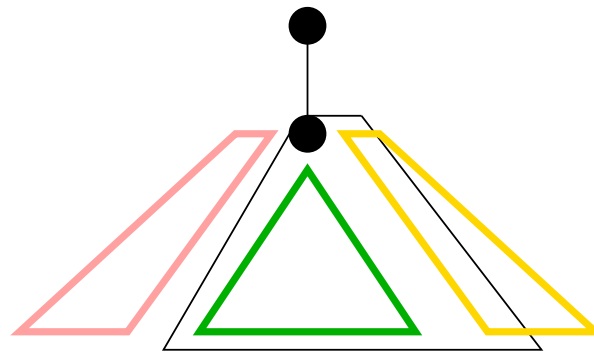
# Que se passe-t-il pour l'autre arbre ?

Il a exactement trois types de forêts pour l'arbre directeur

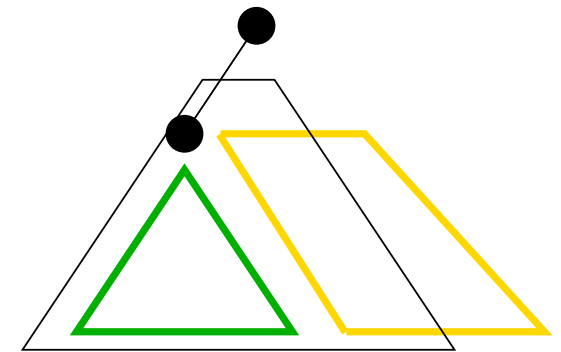
- ▷ celles qui sont comparées avec toutes les forêts gauches de  $B$
- ▷ celles qui sont comparées avec toutes les forêts droites de  $B$
- ▷ celles qui sont comparées avec toutes les forêts de  $B$



toutes les forets droites



toutes les forets



toutes les forets gauches

Les fils favoris et les fils uniques héritent des forêts de leur parent.

- ▷ **Free:** les nœuds qui n'héritent de rien
- ▷ **Left :** les nœuds qui héritent des forêts gauches
- ▷ **Right :** les nœuds qui héritent des forêts droites
- ▷ **All:** les nœuds qui héritent de toutes les forêts

1.  $A$  est un seul nœud de direction droite

$$\begin{aligned}\text{Free}(A) &= \text{Left}(A) = \#\text{left}(B) \\ \text{All}(A) &= \text{Right}(A) = \#\text{special}(B)\end{aligned}$$

2.  $A$  est un seul nœud de direction gauche

$$\begin{aligned}\text{Free}(A) &= \text{Right}(A) = \#\text{right}(B) \\ \text{All}(A) &= \text{Left}(A) = \#\text{special}(B)\end{aligned}$$

3.  $A = l(A')$  et la direction de  $l$  est droite

$$\begin{aligned}\text{Free}(A) &= \text{Left}(A) = \#\text{left}(B) + \text{Right}(A') \\ \text{All}(A) &= \text{Right}(A) = \#\text{special}(B) + \text{All}(A')\end{aligned}$$

4.  $A = l(A')$  et la direction de  $l$  est gauche

$$\begin{aligned}\text{Free}(A) &= \text{Right}(A) = \#\text{right}(B) + \text{Left}(A') \\ \text{All}(A) &= \text{Left}(A) = \#\text{special}(B) + \text{All}(A')\end{aligned}$$



5.  $A = l(A_1 \circ \dots \circ A_n)$  et le fils favori est  $A_1$

$$\begin{aligned} \text{Free}(A) &= \text{Left}(A) = \sum_{i>1} \text{Free}(A_i) + \text{Left}(A_1) + \#\text{left}(B)(|A| - |A_1|) \\ \text{All}(A) &= \text{Right}(A) = \sum_{i>1} \text{Free}(A_i) + \text{All}(A_1) + \#\text{special}(B)(|A| - |A_1|) \end{aligned}$$

6.  $A = l(A_1 \circ \dots \circ A_n)$  et le fils favori est  $A_1$

$$\begin{aligned} \text{Free}(A) &= \text{Right}(A) = \sum_{i<n} \text{Free}(A_i) + \text{Right}(A_n) + \#\text{right}(B)(|A| - |A_n|) \\ \text{All}(A) &= \text{Left}(A) = \sum_{i<n} \text{Free}(A_i) + \text{All}(A_n) + \#\text{special}(B)(|A| - |A_n|) \end{aligned}$$

7. sinon, le fils favori est  $A_j$  avec  $1 < j < n$

$$\begin{aligned} \text{Free}(A) &= \sum_{i \neq j} \text{Free}(A_i) + \text{All}(A_j) + \#\text{right}(B)(1 + |A_1 \circ \dots \circ A_{j-1}|) \\ &\quad + \#\text{special}(B)|A_j \circ \dots \circ A_n| \\ \text{Right}(A) &= \text{Free}(A) \\ \text{All}(A) &= \text{Left}(A) = \sum_{i \neq j} \text{Free}(A_i) + \text{All}(A_j) + \#\text{special}(B)(|A| - |A_j|) \end{aligned}$$

# Comment en déduire un recouvrement optimal

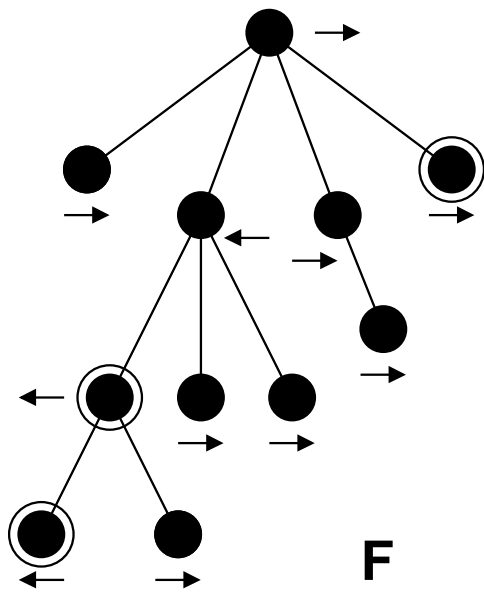
▷ programmation dynamique

▷ 4 tables : *Free*, *All*, *Left*, *Right*

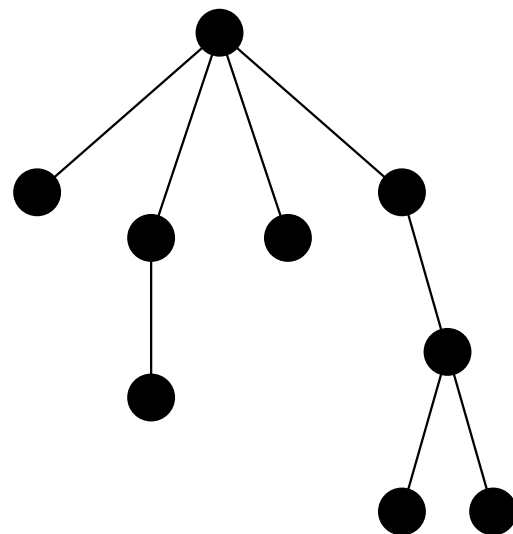
$$\begin{aligned} \text{Free}(A) &= \sum_{i \geq 1} \text{Free}(A_i) \\ &+ \min \left\{ \begin{array}{l} \text{Left}(A_1) - \text{Free}(A_1) + \#\text{left}(B) * (|A| - |A_1|) \\ \text{All}(A_j) - \text{Free}(A_j) + \#\text{special}(B) |A_j \circ \dots \circ A_n| \\ \quad + \#\text{right}(B) (1 + |A_1 \circ \dots \circ A_{j-1}|), \quad 1 < j < n \\ \text{Right}(A_n) - \text{Free}(A_n) + \#\text{right}(B) * (|A| - |A_n|) \end{array} \right. \end{aligned}$$

▷ Coût du prétraitement :  $O(\sum_i \text{arité}(A(i))) + O(|B|) = O(|A|) + O(|B|)$

optimal	:	340
droite	:	405
gauche	:	350
Klein	:	391

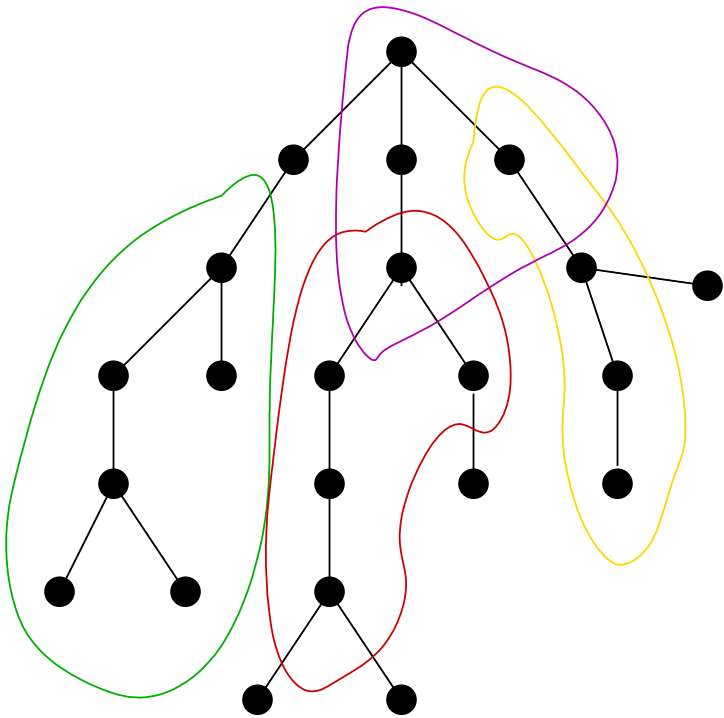


**F**



**G**

## Distance d'édition avec gaps



Distance: NP-dur

Indel sur les sous-structures  
terminales :

$$O(n^2)$$