

Tree edit distance analysis

Serge Dulucq (LABRI)

Hélène Touzet (LIFL)

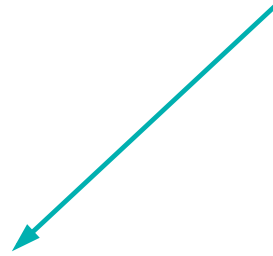
CPM 2003

Tree Edit Distance Problem

- ▷ A pair of ordered rooted trees (A, B)
- ▷ Costs for edit operations
 - `sub`: substituting a node
 - `ins`: inserting a node
 - `del`: deleting a node
- ▷ **Distance**: minimal cost to transform A into B

String edit distance I

A L G O R I T H M
F L O W E R



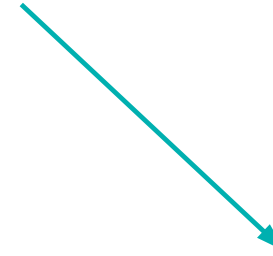
substitution

A	LGORITHM
F	LOWER



insertion

-	ALGORITHM
F	LOWER



deletion

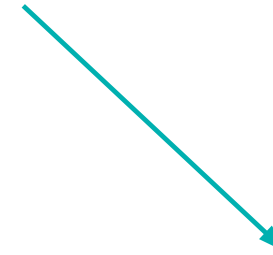
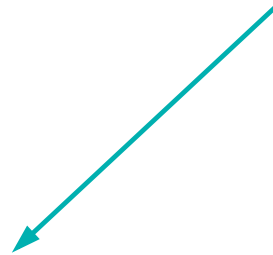
A	LGORITHM
-	FLOWER

Leftmost decomposition

Comparison of all pairs of suffixes

String edit distance II

A L G O R I T H M
F L O W E R



substitution

insertion

deletion

ALGORITHM

M
R

FLOWE

ALGORITHM

-
R

FLOWE

ALGORITHM

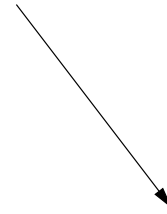
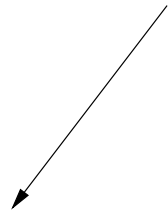
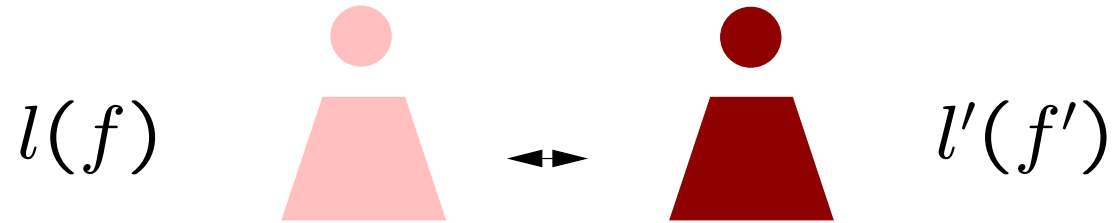
M
-

FLOWER

Rightmost decomposition

Comparison of all pairs of prefixes

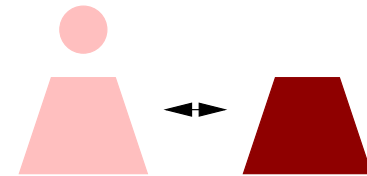
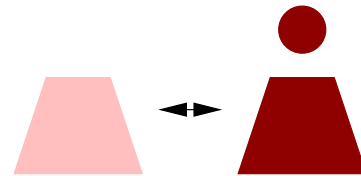
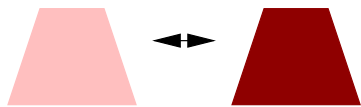
Tree edit distance



Substitution of l into l'

Deletion of l

Insertion of l'

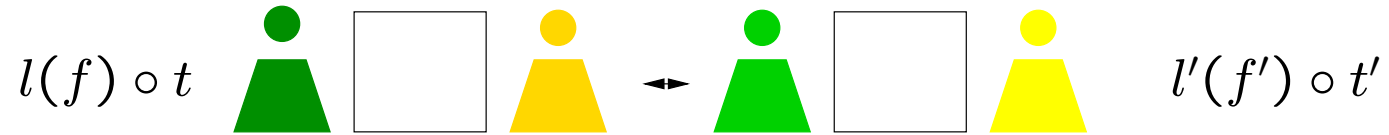


$$\text{Distance}(f, f') + \text{sub}(l, l')$$

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

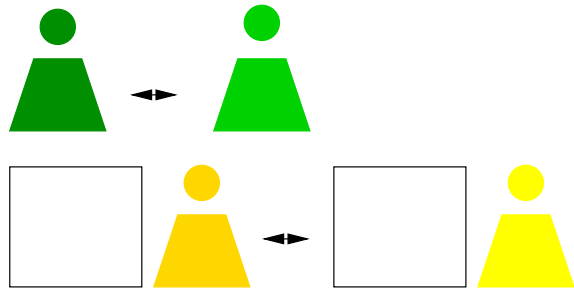
$$\text{Distance}(l(f), f') + \text{ins}(l')$$

Forest edit distance I



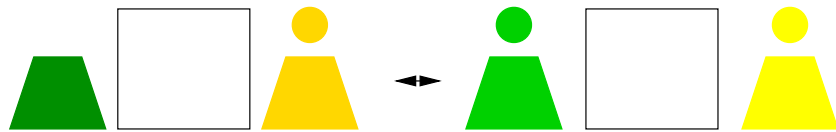
Leftmost decomposition

Substitution of l into l'



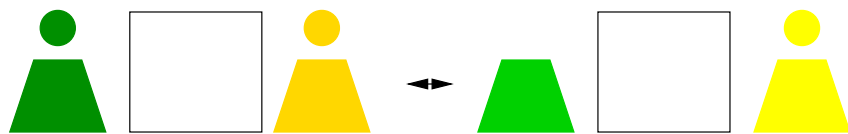
$$\text{Distance}(l(f), l'(f')) + \text{Distance}(t, t')$$

Deletion of l



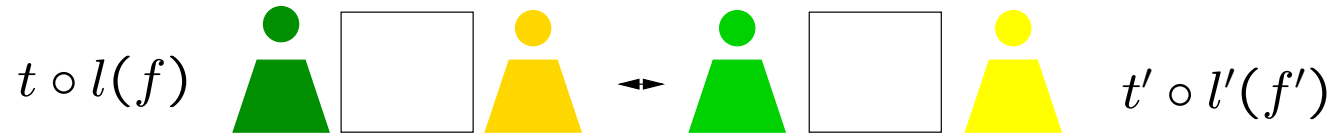
$$\text{Distance}(f \circ t, l'(f') \circ t') + \text{del}(l)$$

Insertion of l'



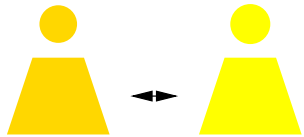
$$\text{Distance}(l(f) \circ t, f' \circ t') + \text{ins}(l')$$

Forest edit distance II



Rightmost decomposition

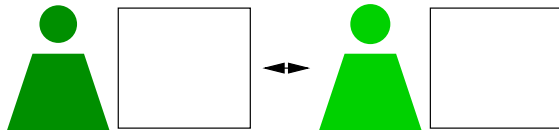
Substitution of l into l'



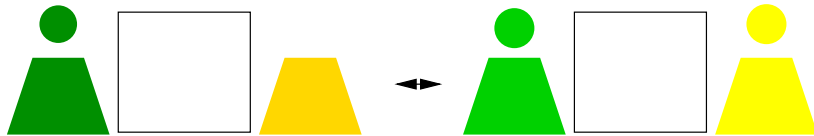
$\text{Distance}(l(f), l'(f'))$

+

$\text{Distance}(t, t')$

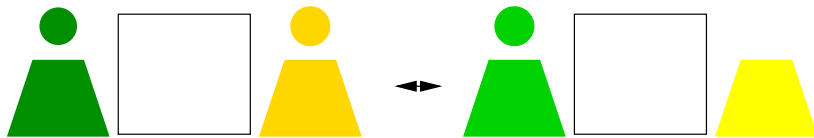


Deletion of l



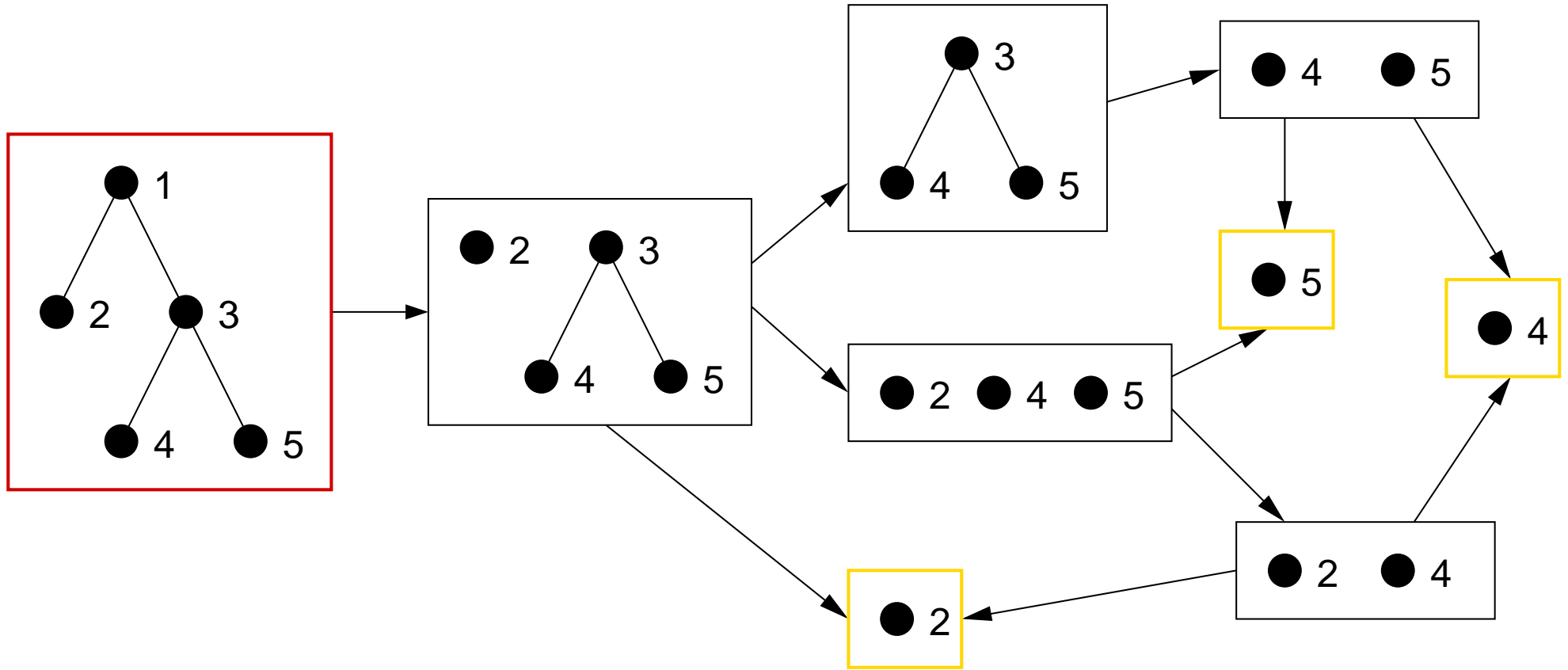
$\text{Distance}(t \circ f, t' \circ l'(f')) + \text{del}(l)$

Insertion of l'



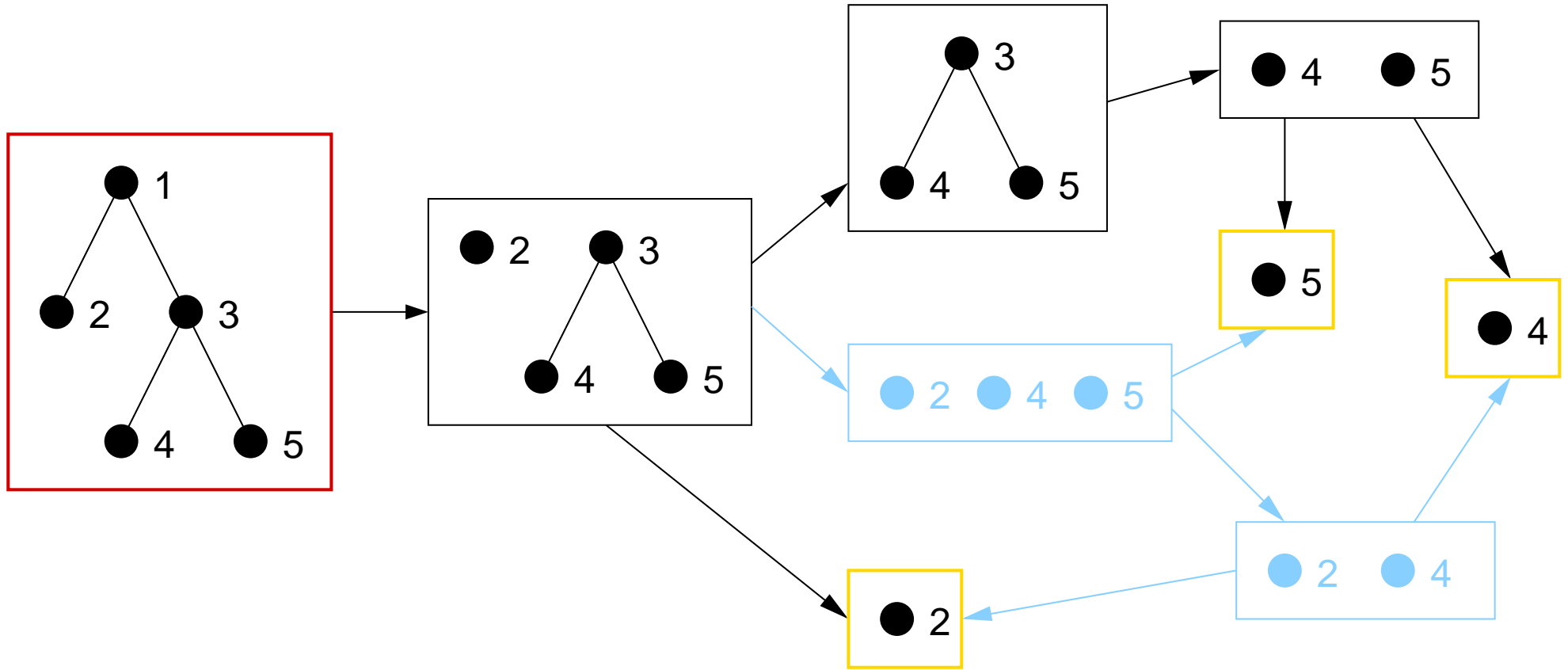
$\text{Distance}(t \circ l(f), t' \circ f') + \text{ins}(l')$

Left/right for trees ?



rightmost decomposition

Left/right for trees ?



leftmost decomposition

Decomposition strategies

▷ Succession of choices **left** or **right**

▷ $S : forest \times forest \rightarrow \{\text{left}, \text{right}\}$

▷ **Zhang & Shasha** (1989) :

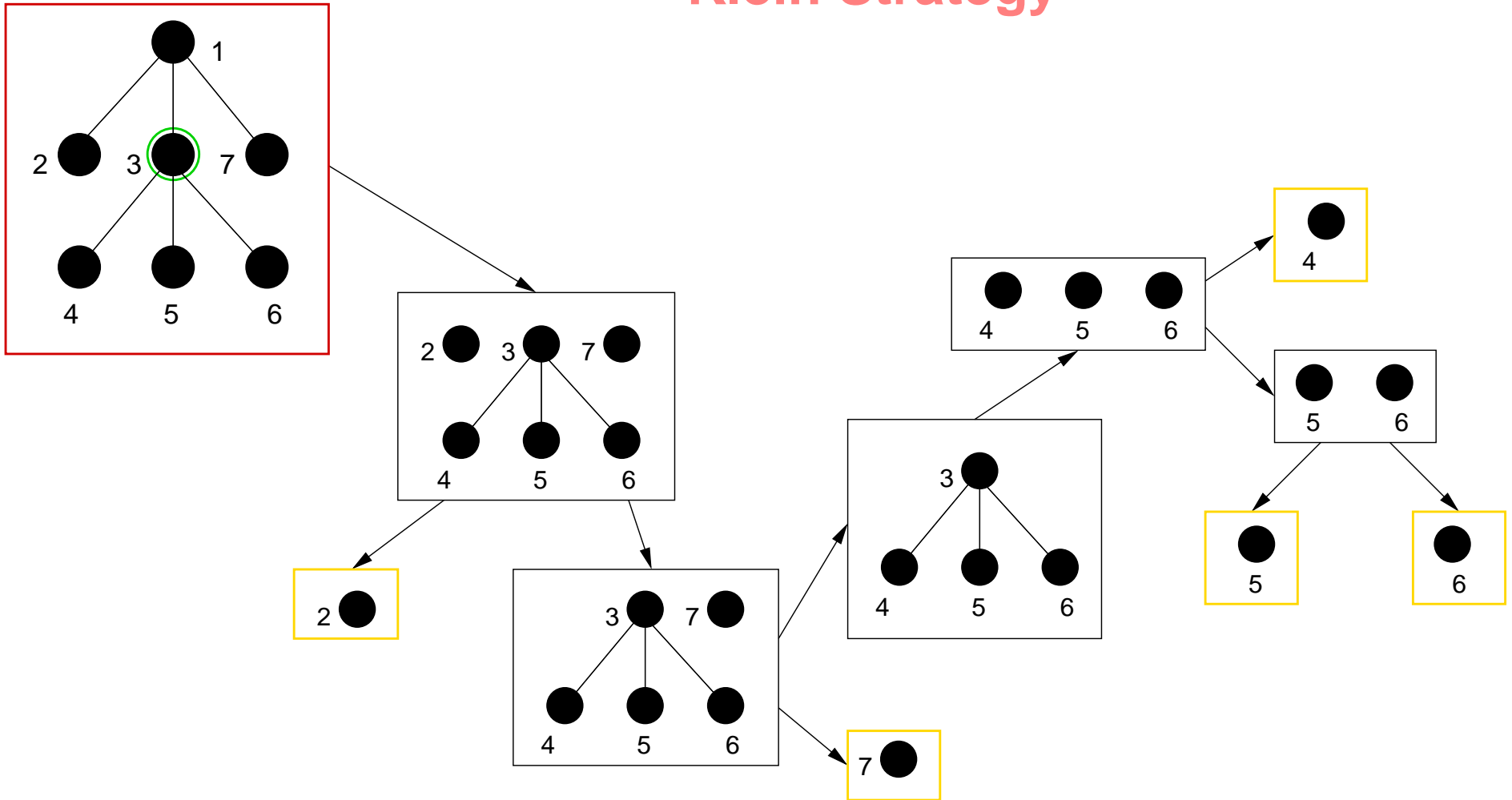
$(f, g) \rightarrow \text{left}$

▷ **Klein** (1998) :

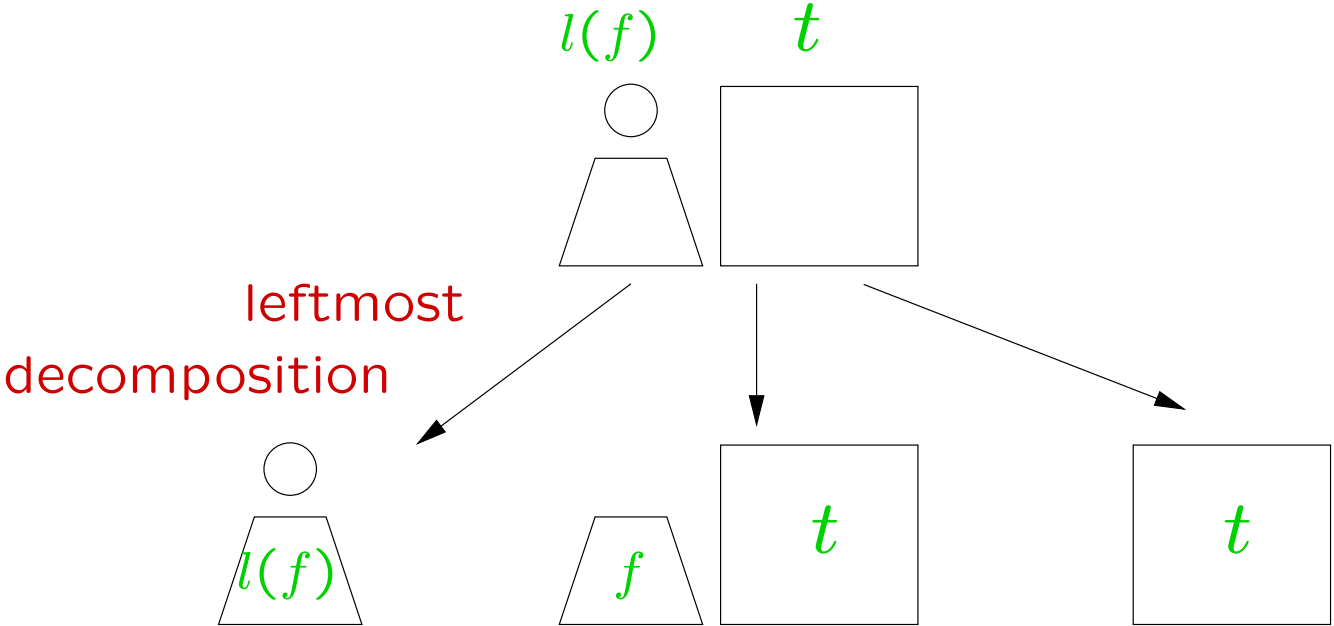
$(f, g) \rightarrow$ **right**, when the first node of f is
the heaviest child

left, otherwise

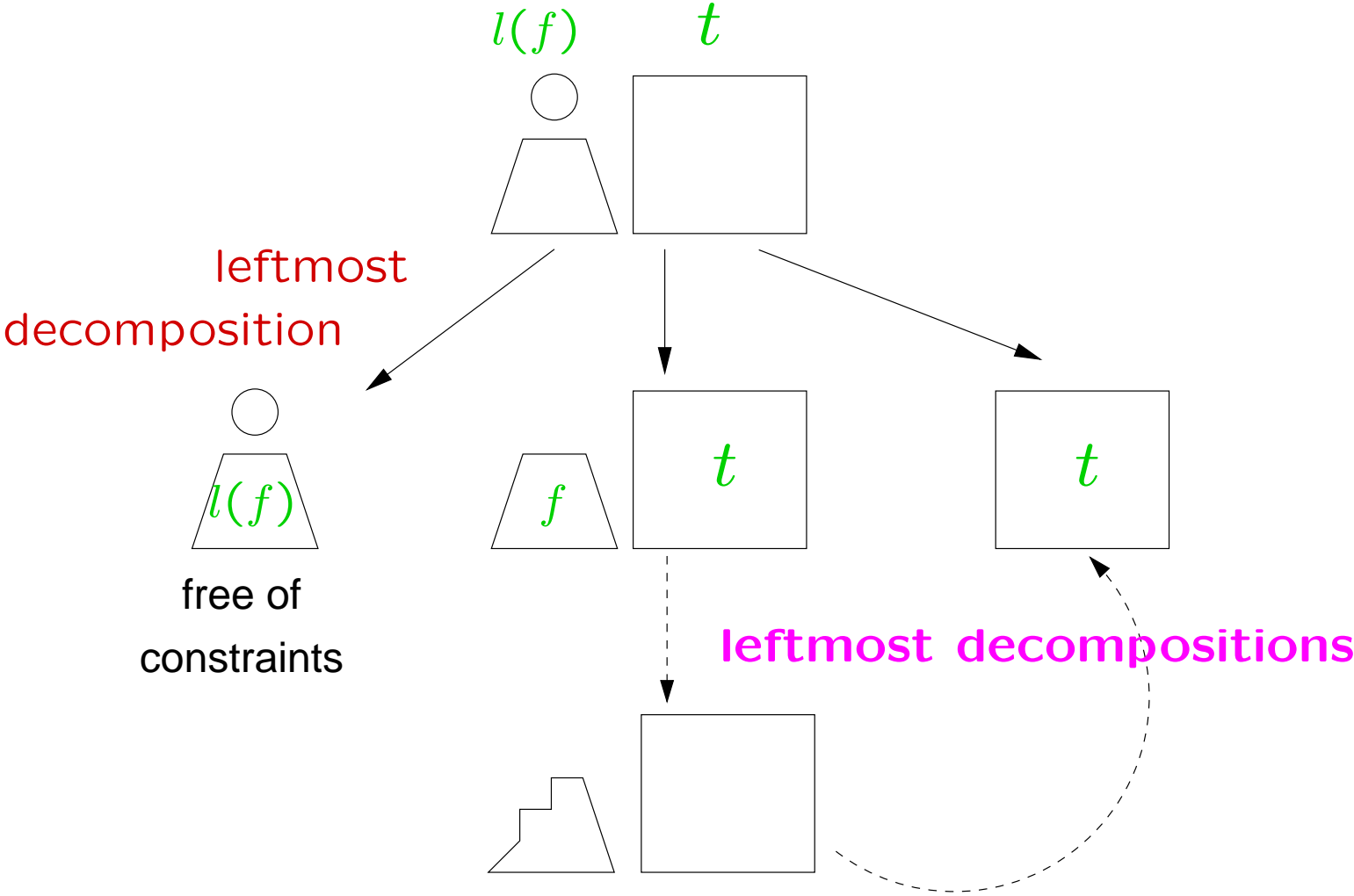
Klein Strategy



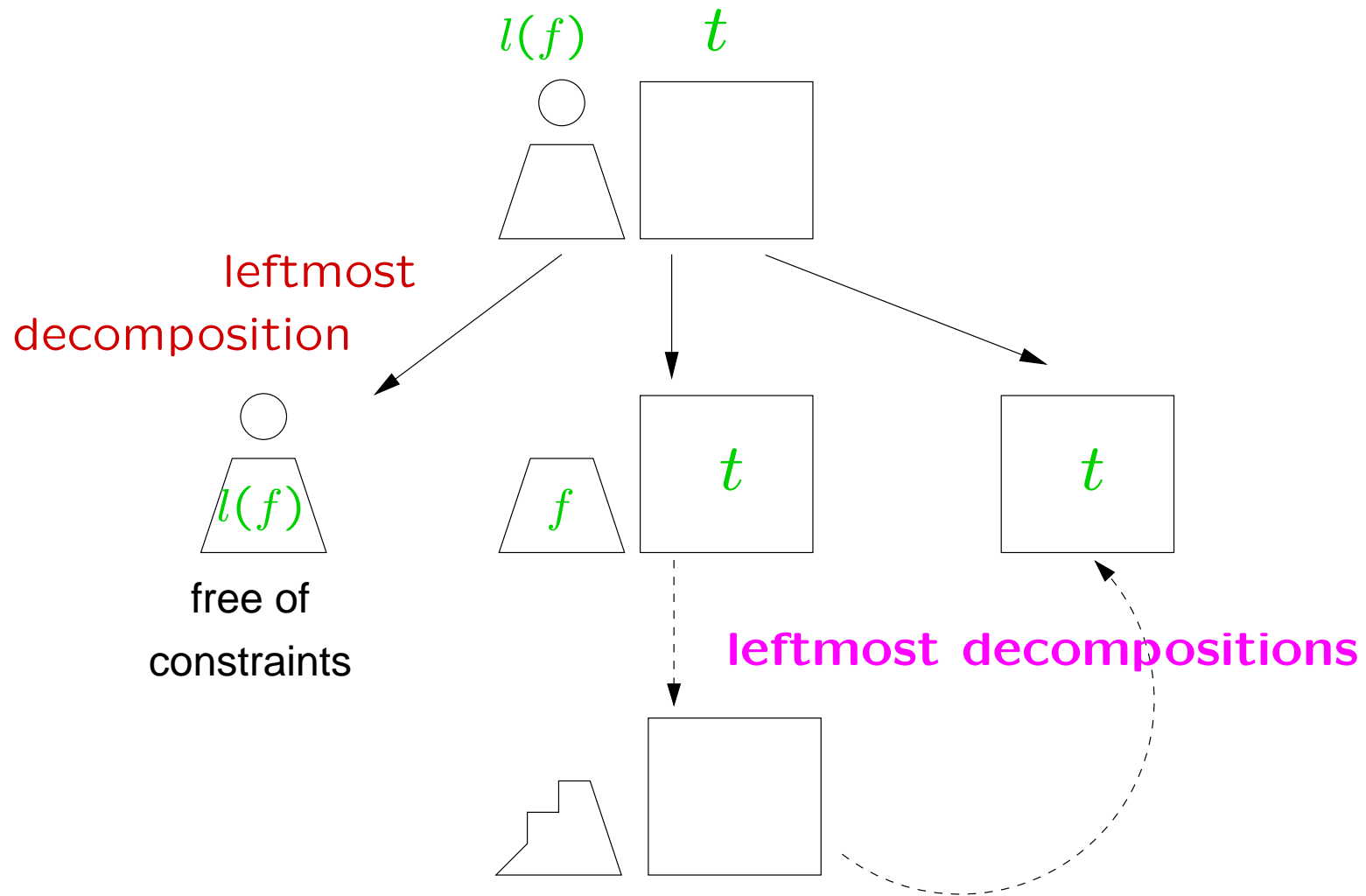
How to built up an economical strategy ?



How to built up an economical strategy ?



How to built up an economical strategy ?

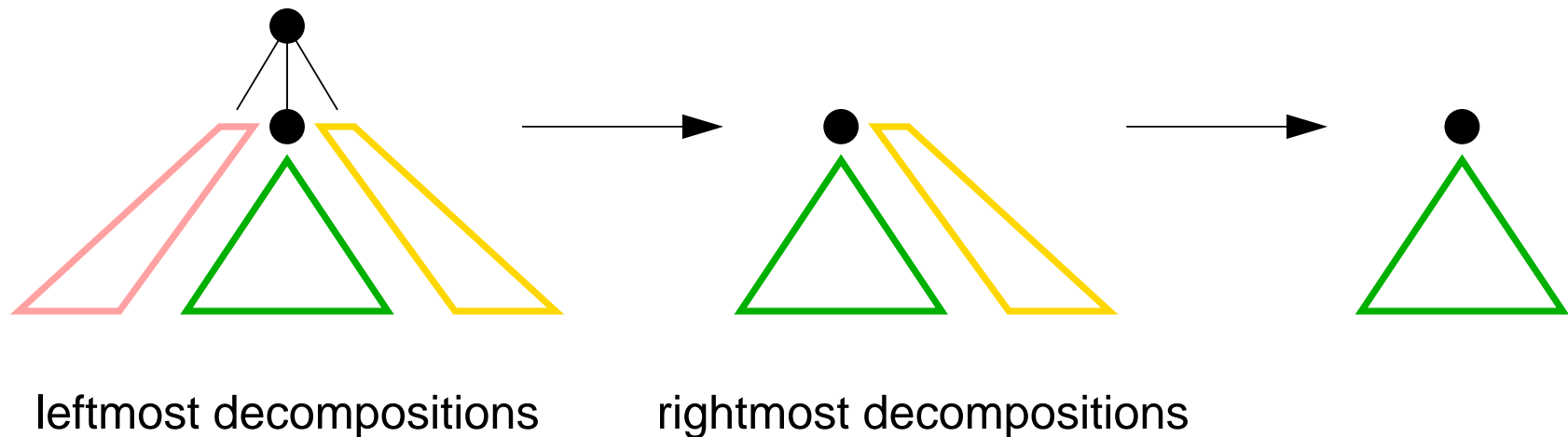


$$\#subforest(l(f) \circ t) = \#subforest(l(f)) + |l(f)| + \#subforest(t)$$

Cover strategies are economical strategies

For a tree A , define a **cover** ϕ for A as

- ▷ $\phi(i) \in \{left, right\}$ if the degree of i is 0 or 1 : **direction**
- ▷ $\phi(i)$ is a child of i : **favorite** child



Zhang & Shasha and Klein are cover strategies.

Number of subforests for one tree

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

▷ Lower bound (no assumption on the strategy)

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

A_j is the heaviest child

▷ Upper bound (no assumption on the strategy)

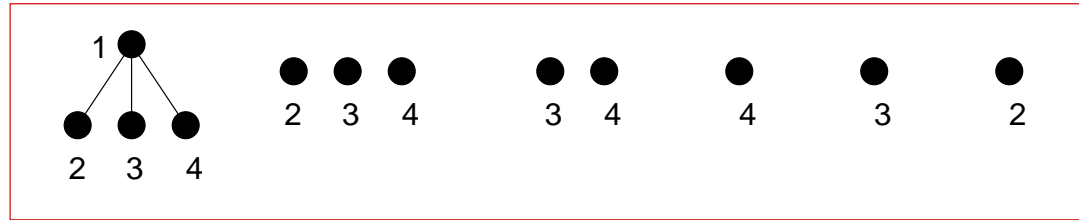
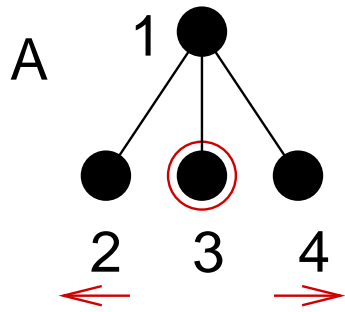
$$\#\text{subforest}(A) \leq \frac{n(n+3)}{2} - \sum_{i \in A} |A(i)| \quad \frac{1}{2} n^2 + \frac{\sqrt{\pi}}{2} n^{\frac{3}{2}} + O(n) \text{ in average}$$

▷ Exact number for a cover strategy

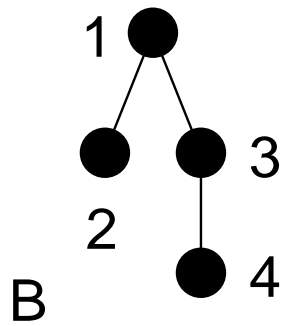
$$\#\text{subforest}(A) = |A| - |A_j| + \#\text{subforest}(A_1) + \dots + \#\text{subforest}(A_n)$$

A_j is the favorite child

Example



$$4 - 1 + 1 + 1 + 1 = 6 \text{ subforests for } A$$



What happens to the other tree B ?

There are only three possibilities for the subforests of the cover tree A

- ▷ being compared with all leftmost forests of B
- ▷ being compared with all rightmost forests of B
- ▷ being compared with all forests of B

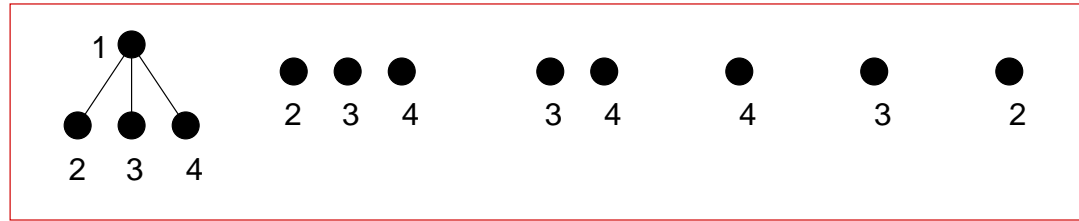
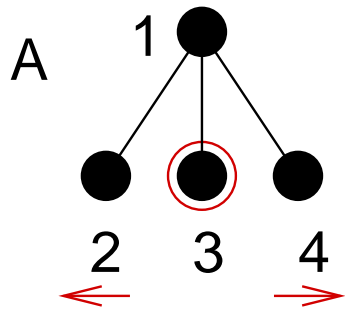
$\#left(B)$: number of leftmost subforests of B

$\#right(B)$: number of rightmost subforests of B

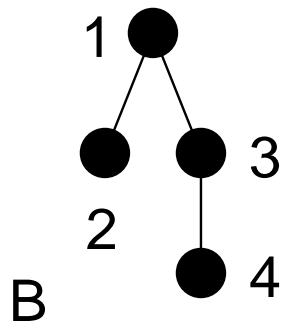
$\#special(B)$: number of subforests of B

$\#left(B)$, $\#right(B)$, $\#special(B)$ are known

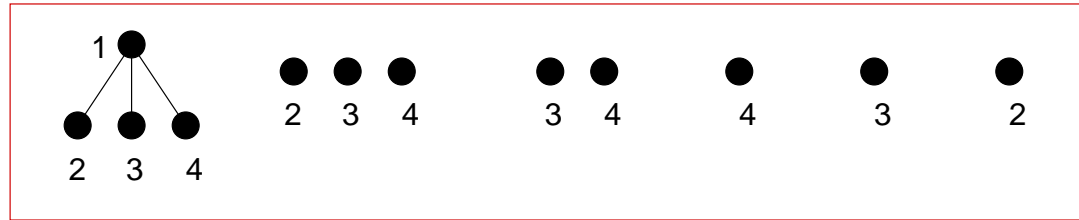
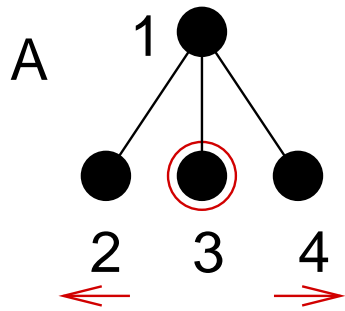
Example (continued)



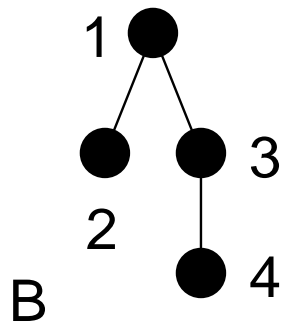
$$4 - 1 + 1 + 1 + 1 = 6 \text{ subforests for } A$$



Example (continued)



$4 - 1 + 1 + 1 + 1 = 6$ subforests for A



$$\#right(B) = 4 - 2 + 1 + 2 = 5$$

$$\#left(B) = 4 - 1 + 1 + 2 = 6$$

$$\#special(B) = 4 * 7/2 - 4 - 1 - 2 - 1 = 6$$

The favorite children inherit subforests from their parent.

4 kinds of nodes :

- ▷ **Free:** nodes that do not receive anything
- ▷ **Left :** nodes that inherit leftmost forests of B
- ▷ **Right :** nodes that inherit rightmost forests of B
- ▷ **All:** nodes that inherit all subforests of B

The status of a node depends of the direction and of the heritage.

1. A is reduced to a node with direction right

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

2. A is reduced to a node with direction left

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

3. $A = l(A')$ and the direction of l is right

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

4. $A = l(A')$ and the direction of l is left

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

5. $A = l(A_1 \circ \dots \circ A_n)$ and the favorite child is A_1 ?

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

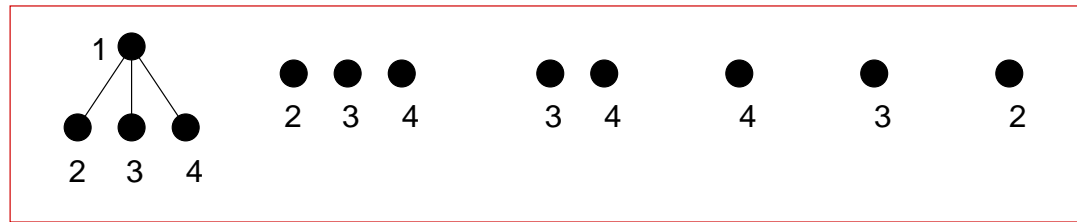
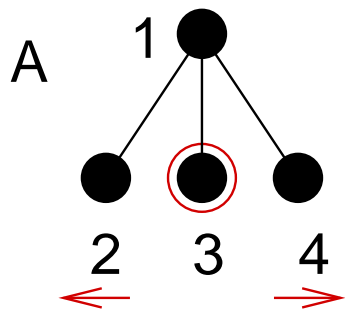
6. $A = l(A_1 \circ \dots \circ A_n)$ and the favorite child is A_n ?

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

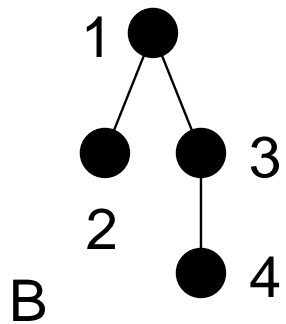
7. otherwise: let A_j ($1 < j < n$) be the favorite child

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

Example (end)



$4 - 1 + 1 + 1 + 1 = 6$ subforests for A

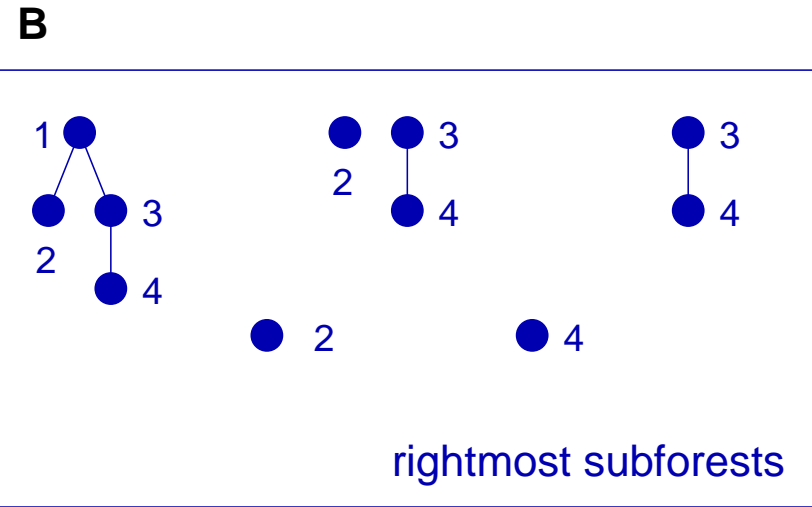
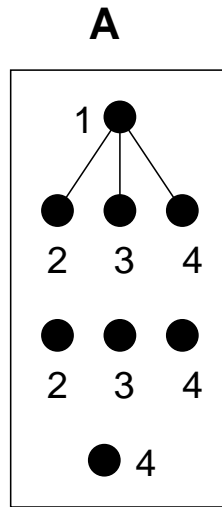
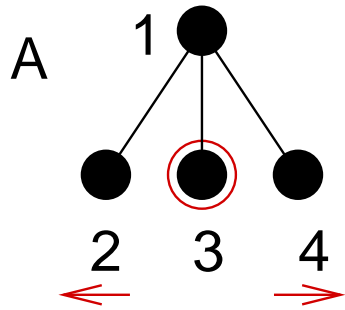


$$\#right(B) = 4 - 2 + 1 + 2 = 5$$

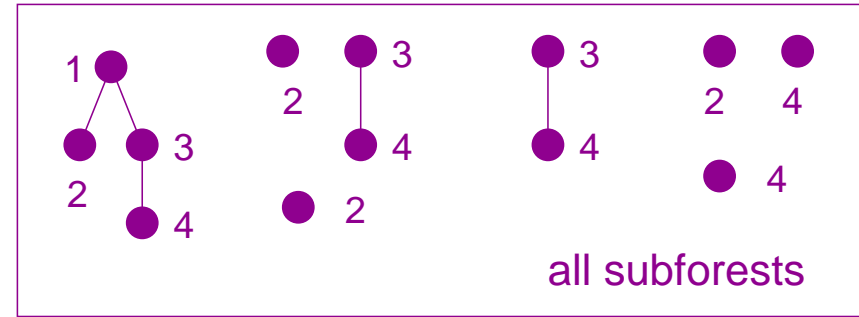
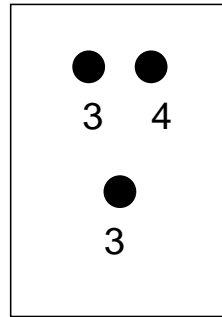
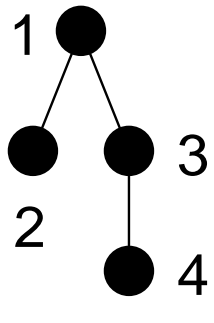
$$\#left(B) = 4 - 1 + 1 + 2 = 6$$

$$\#special(B) = 4 * 7/2 - 4 - 1 - 2 - 1 = 6$$

Example (end)

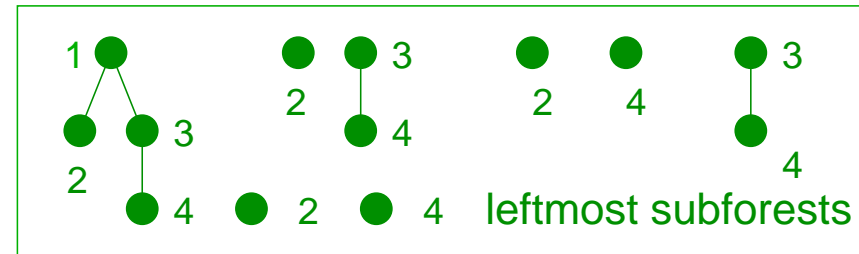
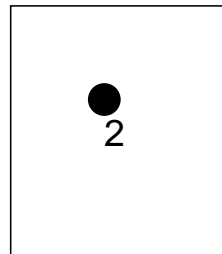


X



X

32 pairs of subforests



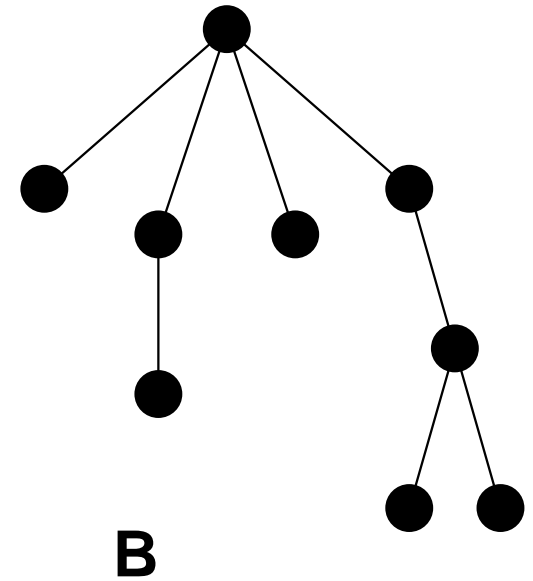
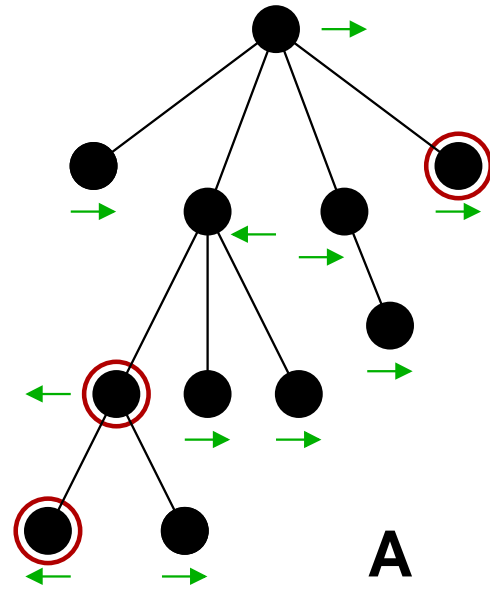
X

How to construct an optimal cover ?

- ▷ Dynamic programming
- ▷ Four tables : *Free*, *All*, *Left*, *Right*

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

- ▷ Preprocessing : $O(\sum_i \text{degree}(A(i))) + O(|B|) = O(|A|) + O(|B|)$



optimal covering : \rightarrow \leftarrow direction

\bigcirc favorite child

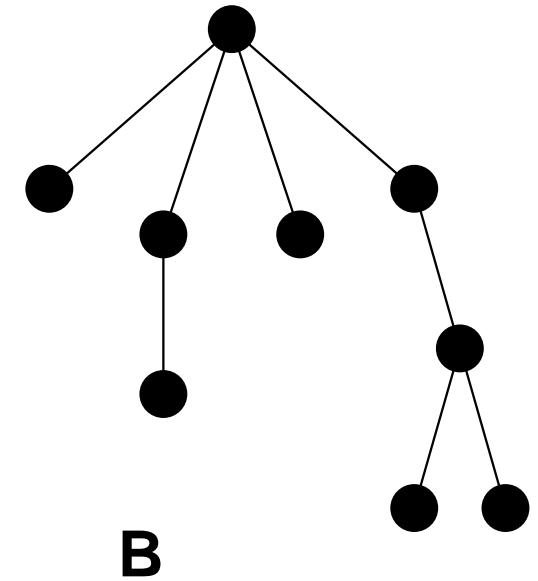
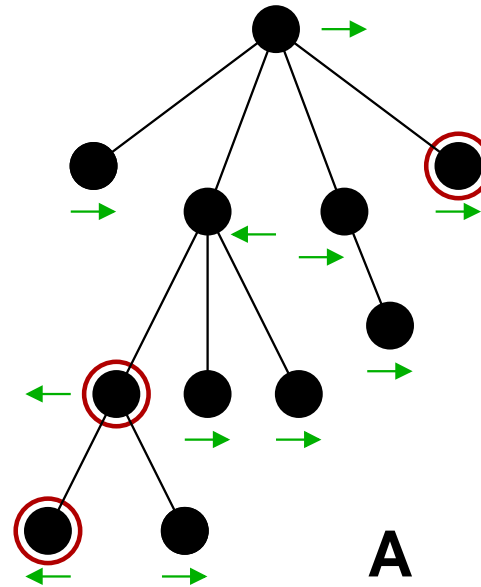
Number of pairs
of subforests

optimal : 340

right : 405

left : 350

Klein : 391



optimal covering : → ← direction

○ favorite child

