

# Recherche d'un ensemble fini de mots

Thierry Lecroq

Université de Rouen  
FRANCE

# Le problème

Localiser toutes les occurrences d'un ensemble fini  $X = \{x_0, x_1, \dots, x_{k-1}\}$  de  $k$  mots dans un texte  $y$  de longueur  $n$ .

Soit  $|X| = |x_0| + |x_1| + \dots + |x_{k-1}|$ .

# Une première solution

Appliquer  $k$  fois un algorithme de recherche exact d'un seul mot

Complexité :  $O(kn)$ .

# Une deuxième solution

En  $O(n)$  sur un alphabet constant.

# Construction de $T(X)$

On commence par construire l'arbre (*trie*)  $T(X) = (Q, q_0, F, \delta)$

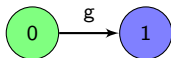
- $Q = \text{Préf}(X)$
- $q_0 = \varepsilon$
- $F = X$
- pour  $u \in \text{Préf}(X)$  et  $a \in A$

$$\delta(u, a) = \begin{cases} ua & \text{si } ua \in \text{Préf}(X) \\ \text{indéfini} & \text{sinon} \end{cases}$$

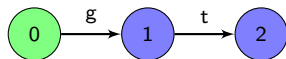
**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

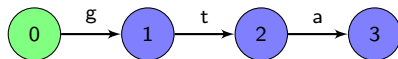


**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

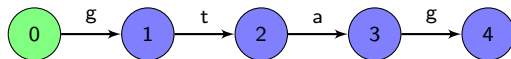




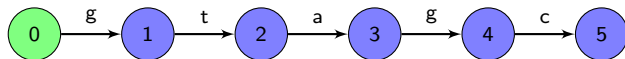
**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



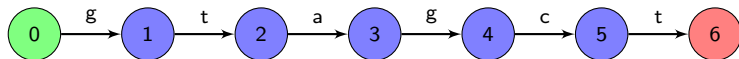
**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



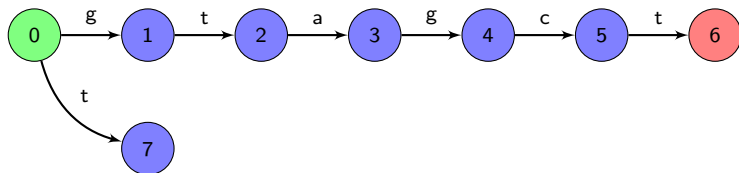
**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



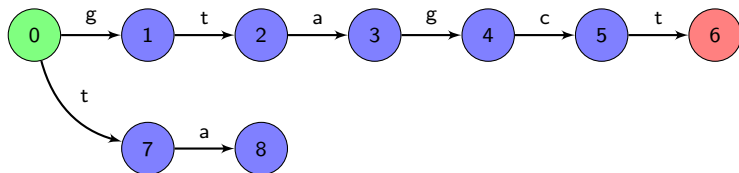
**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



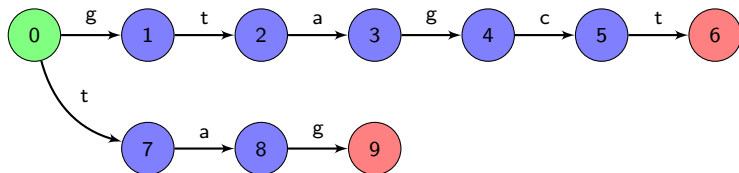
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



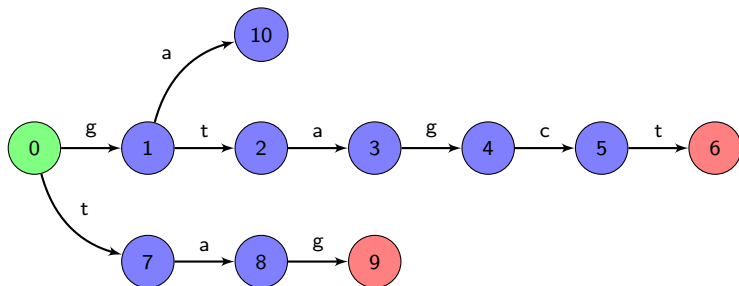
**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



**Exemple**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

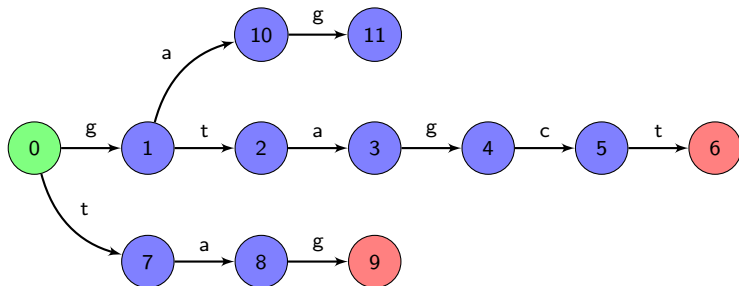


**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

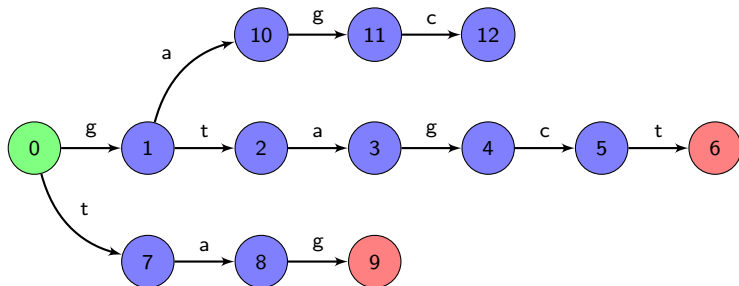




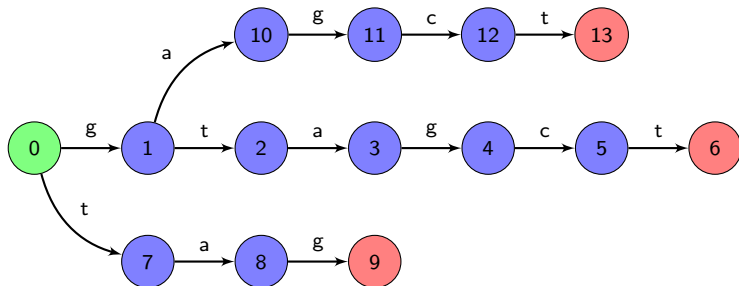
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



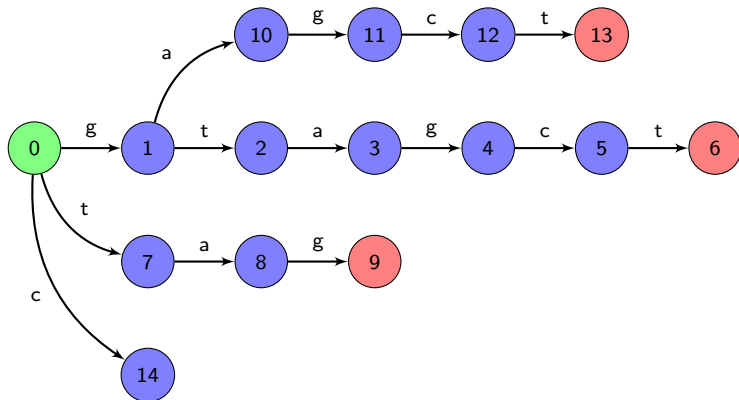
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



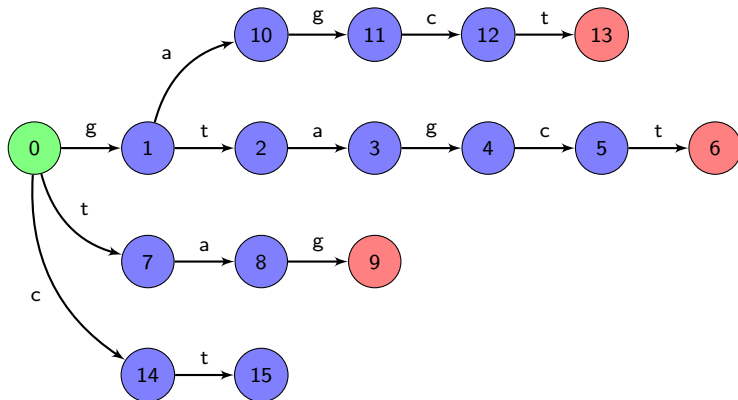
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



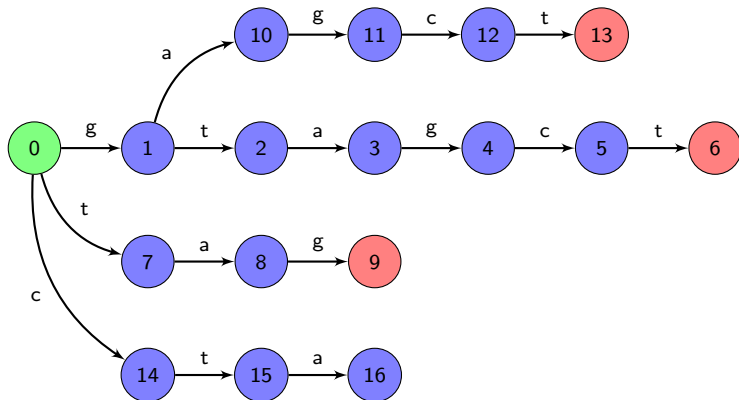
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



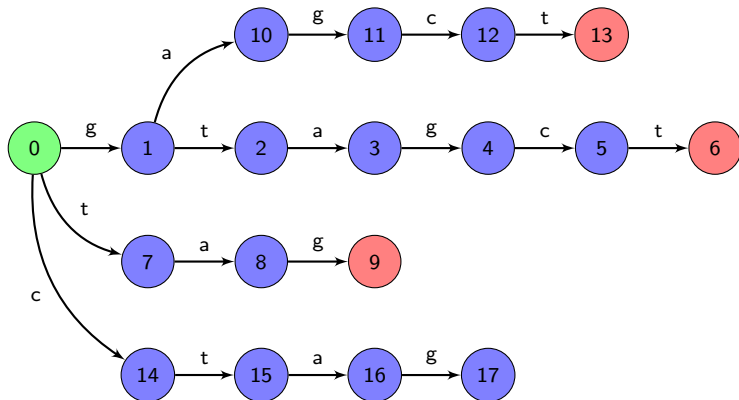
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



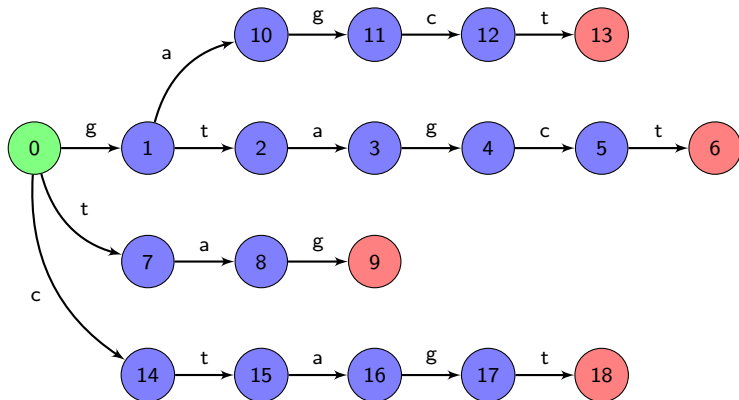
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

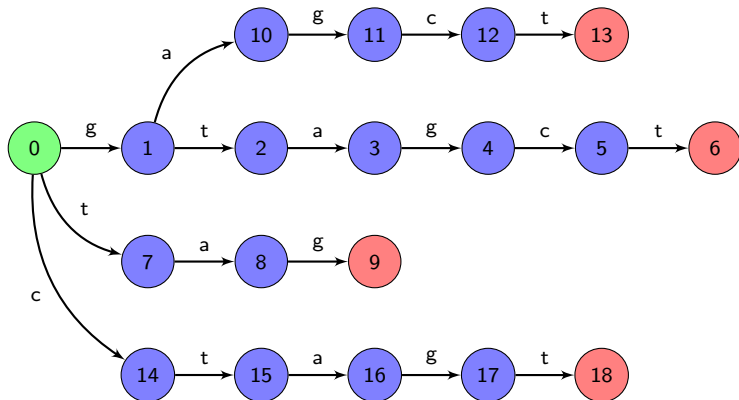




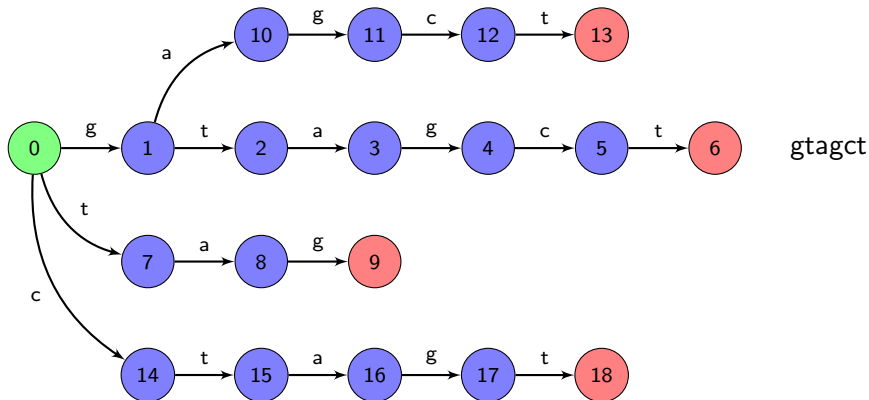
# Suite de la construction

- On associe une fonction de sortie à chaque état terminal :  
 $sortie(x) = \{x\}$  si  $x \in X$
- On crée une boucle sur l'état initial :  $\delta(q_0, a) = q_0$  pour  $a \in A$  et  $a \notin Préf(X)$
- On associe un état suppléant à chaque état :  $sup(q) = u$  où  $u$  est le plus long suffixe propre de  $q$  qui appartient à  $Préf(X)$
- On complète la fonction de sortie : si  $sup(q) = u$  alors  
 $sortie(q) = sortie(q) \cup sortie(u)$

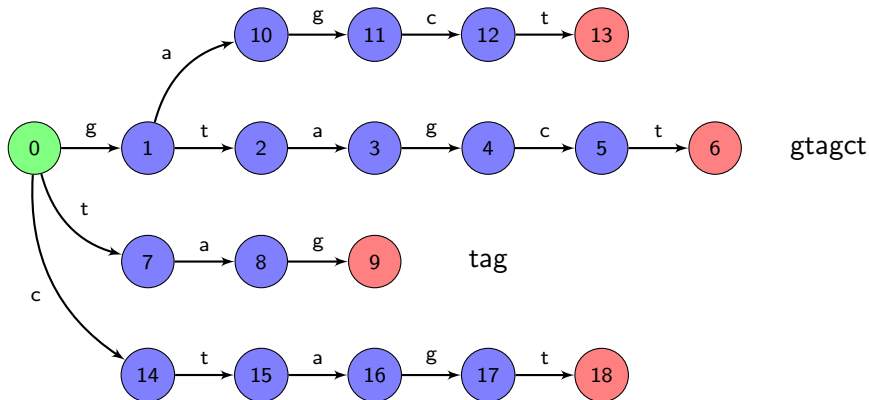
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



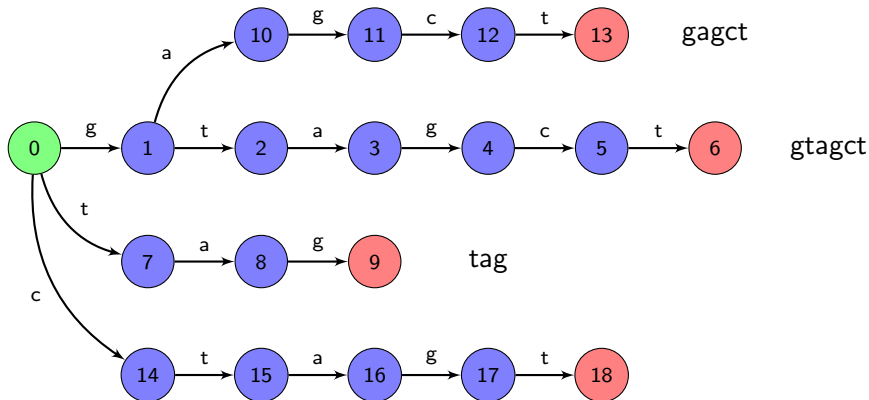
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



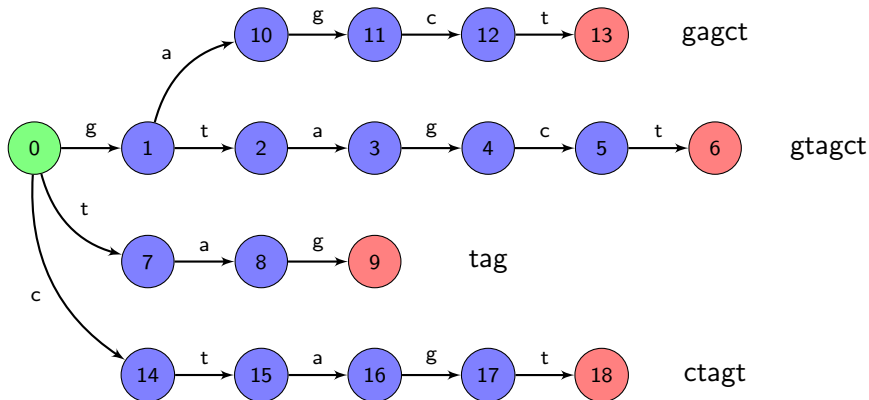
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



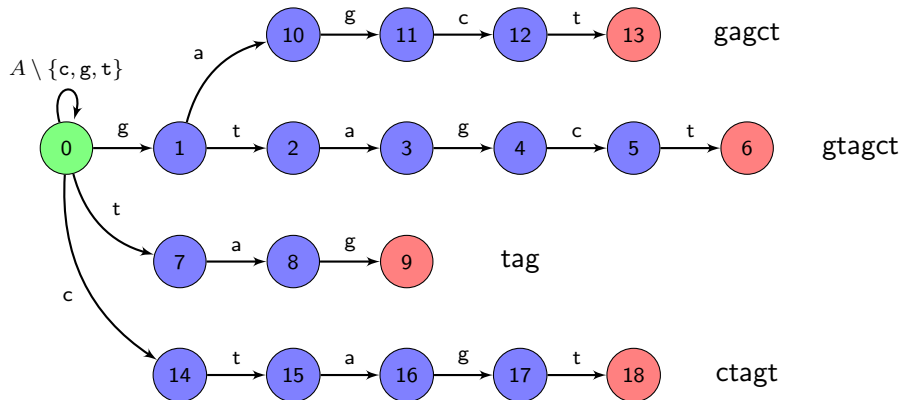
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



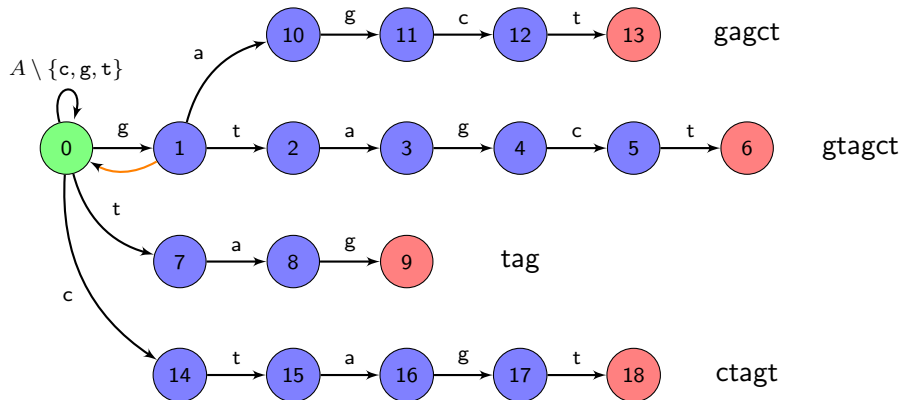
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

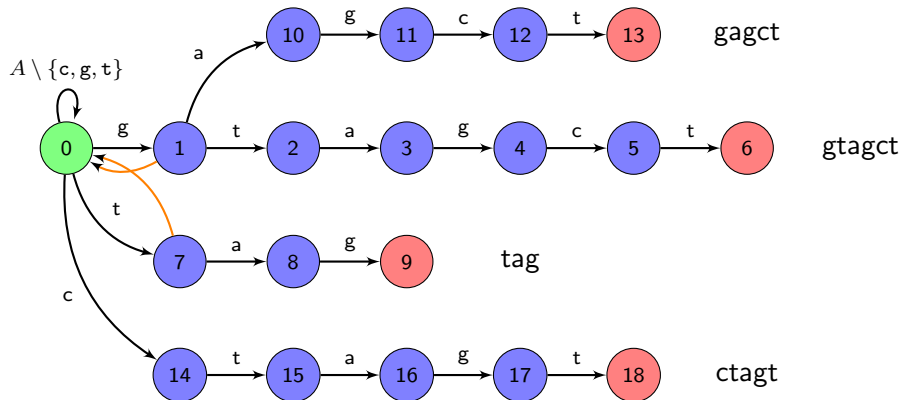


**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

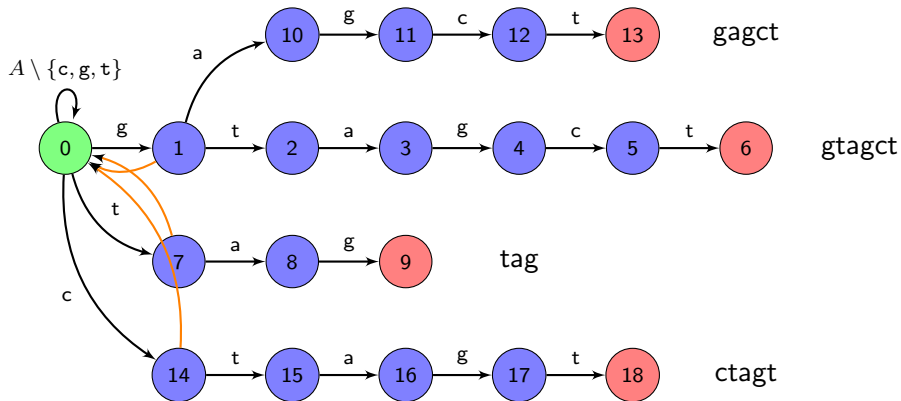




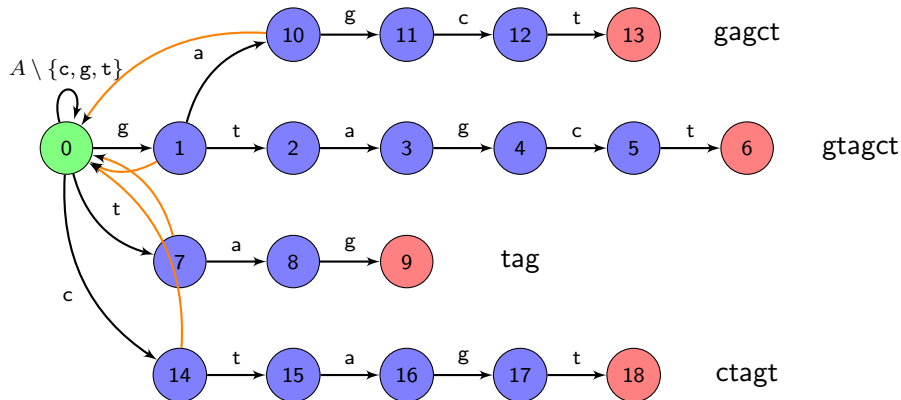
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



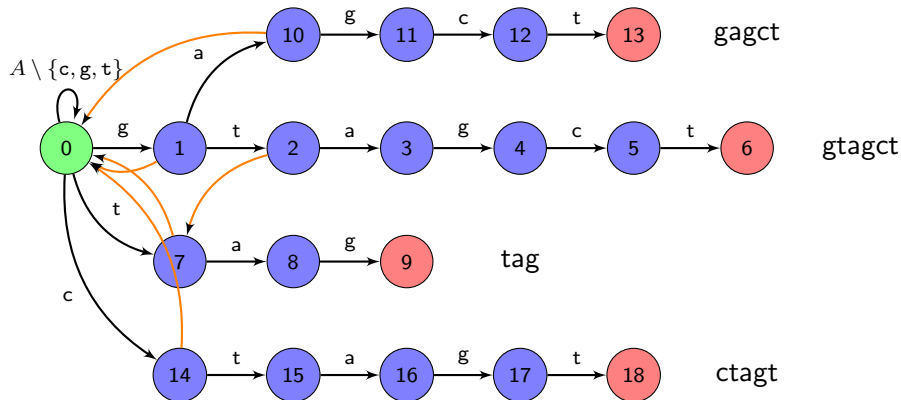
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



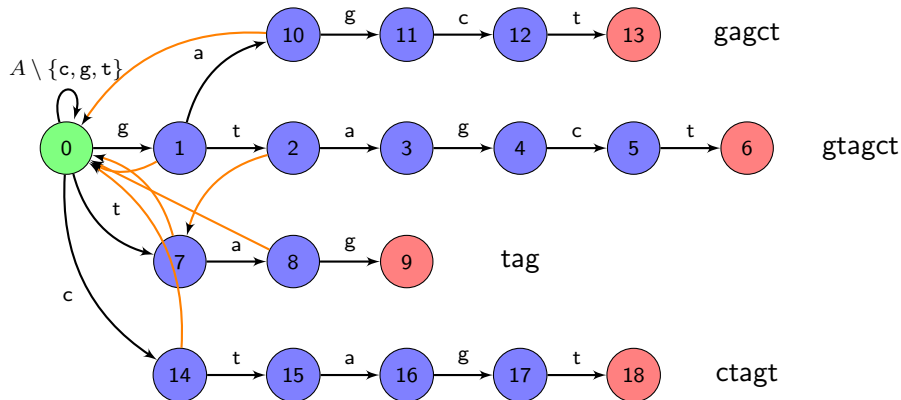
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



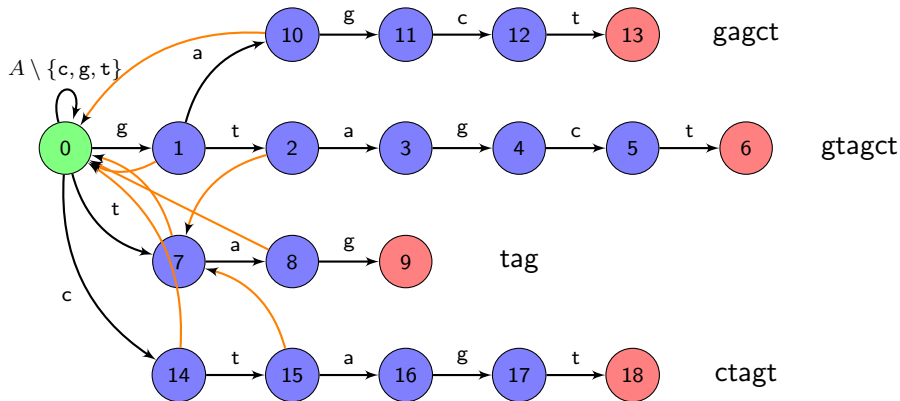
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



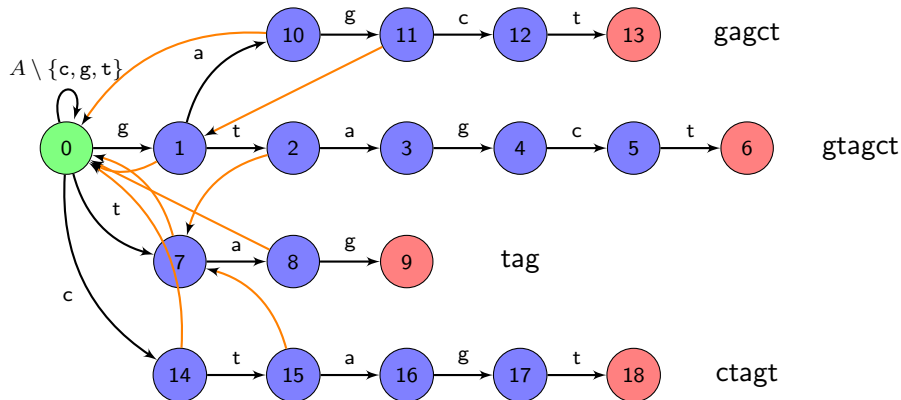
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



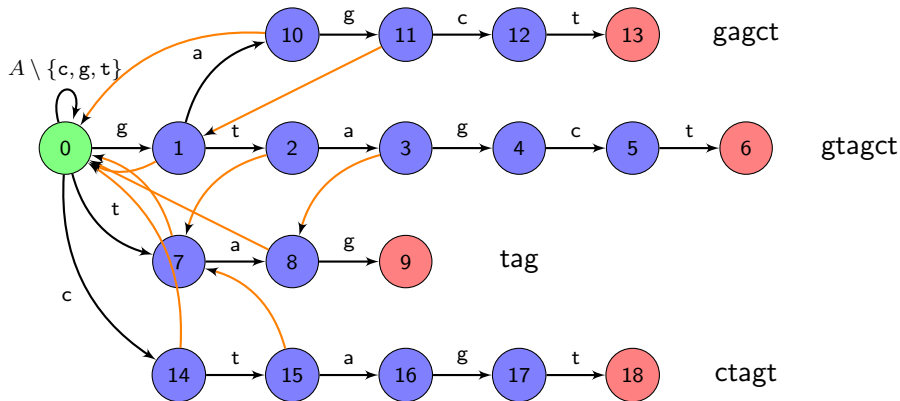
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

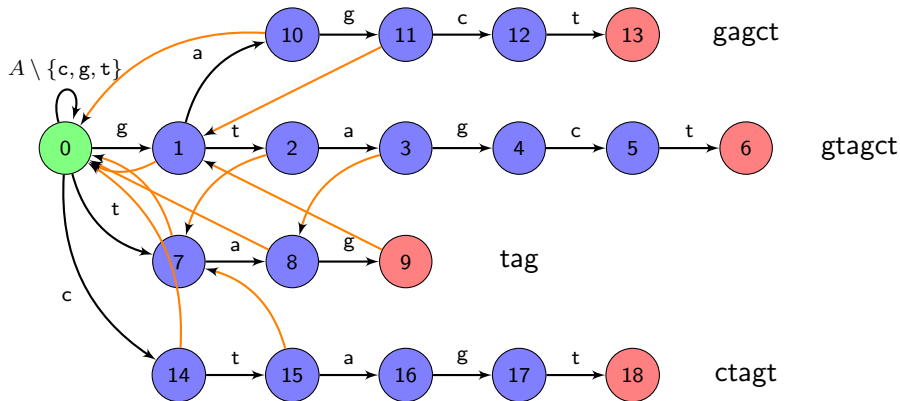


**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

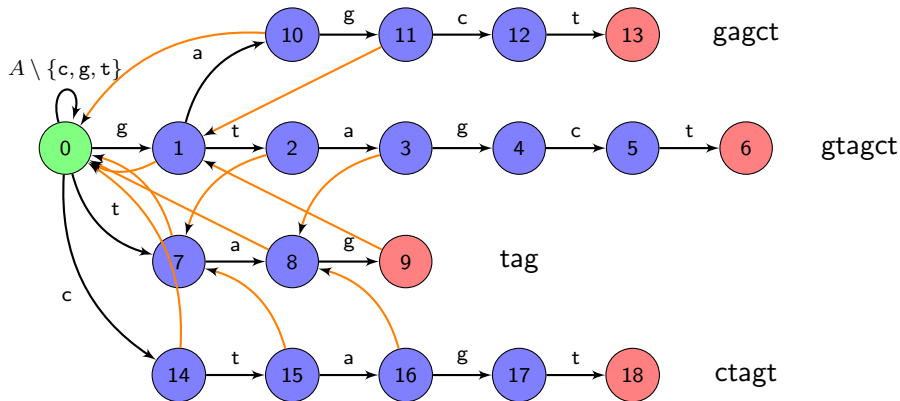




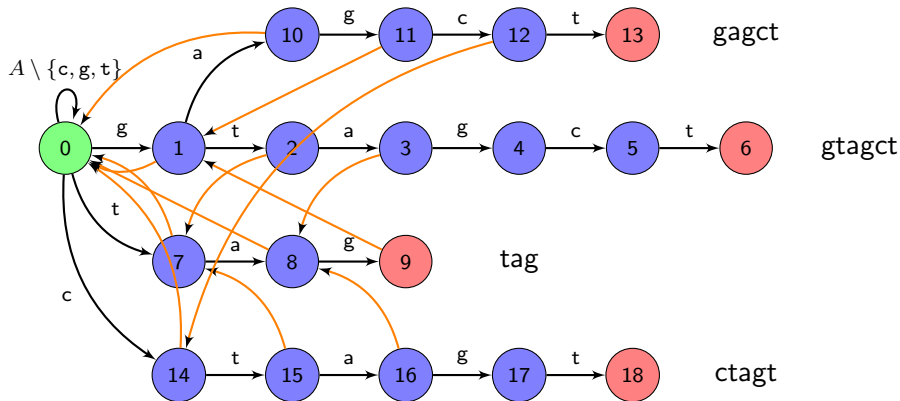
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



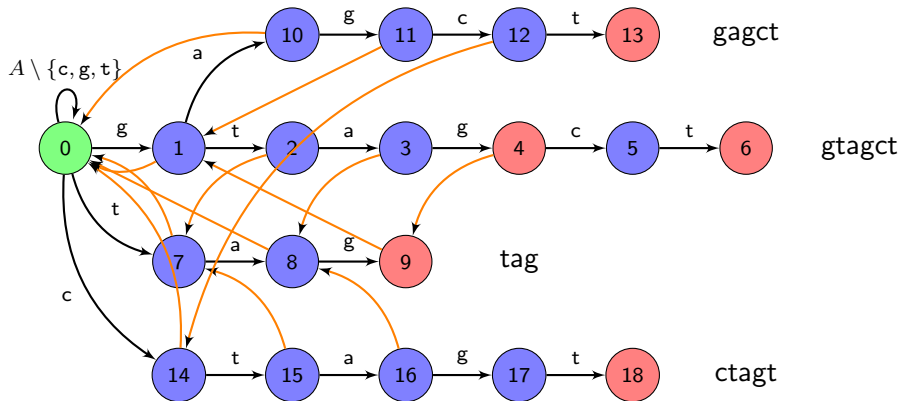
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



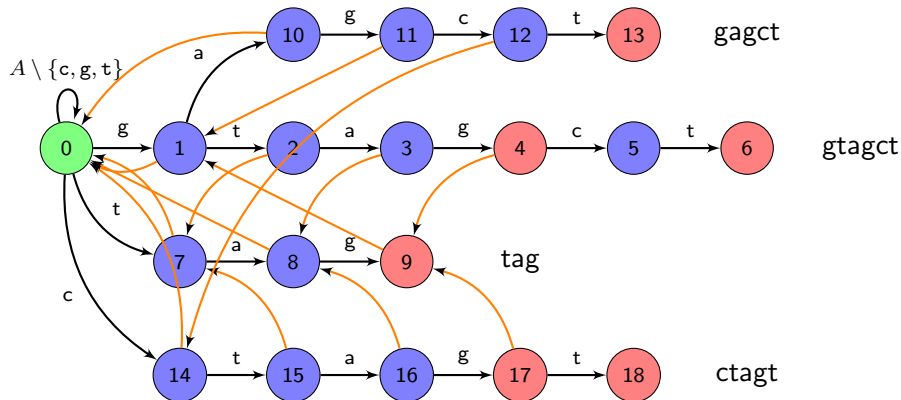
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



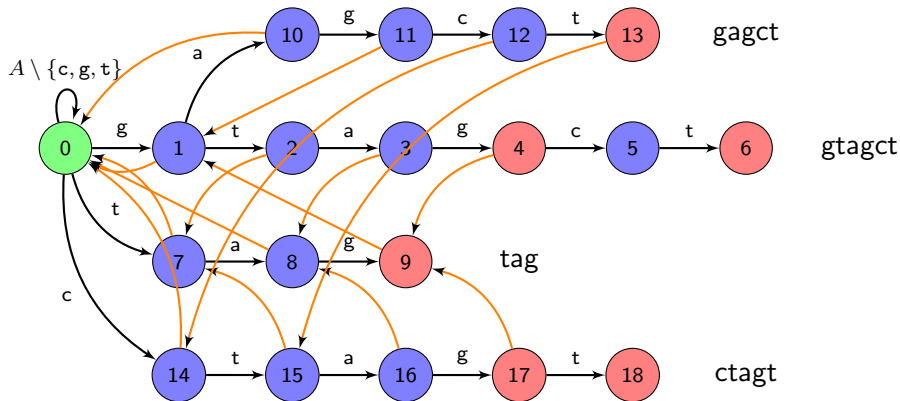
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



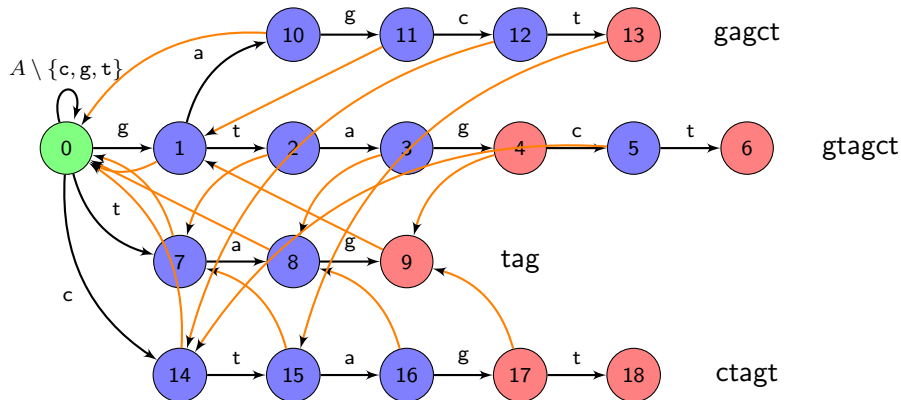
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



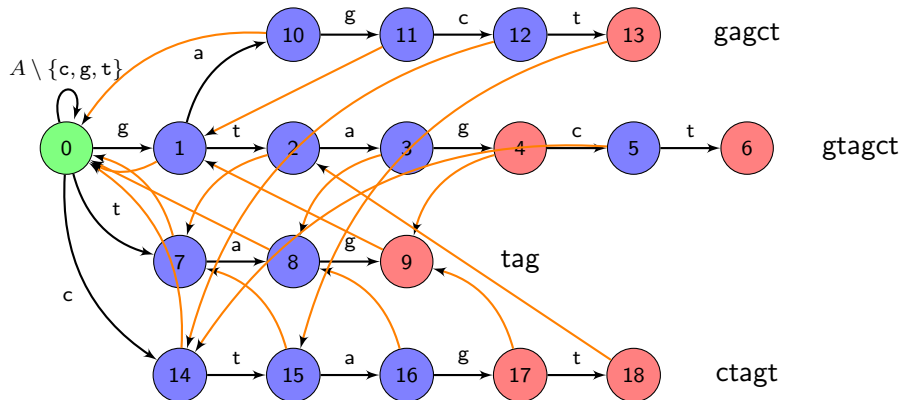
**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$

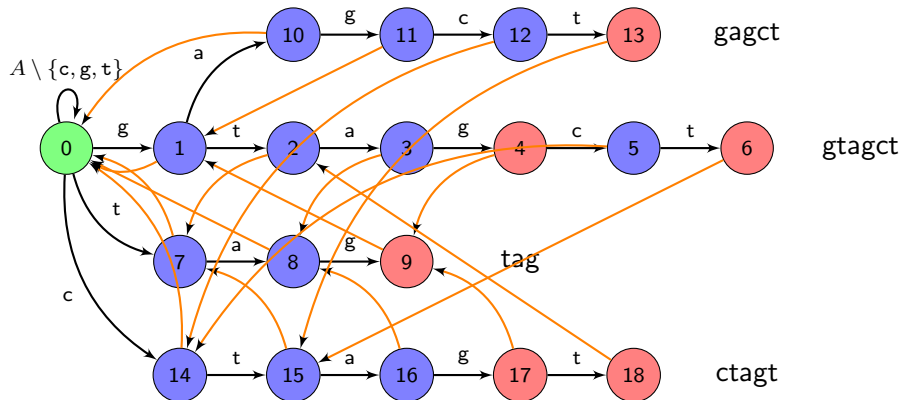


**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$





**Example**  $X = \{\text{gtagct}, \text{tag}, \text{gagct}, \text{ctagt}\}$



## Construction de la fonction de suppléance

La construction de la fonction de suppléance est effectuée par un parcours en largeur de l'arbre  $T(X)$  donc en utilisant une file.

## Pré-AC( $X, k$ )

```
1  créer l'état  $q_0$ 
2  pour  $i \leftarrow 0$  à  $k - 1$  faire
3    ENTRER( $X[i], q_0$ )
4  pour  $a \in A$  faire
5    si  $\delta(q_0, a)$  n'est pas définie alors
6       $\delta(q_0, a) \leftarrow q_0$ 
7  COMPLÉTER( $q_0$ )
8  Retourner  $q_0$ 
```

# Algorithmes

## Entrer( $x, e$ )

```
1   $i \leftarrow 0$ 
2  tantque  $i < |x|$  et  $\delta(e, x[i])$  est définie faire
3     $e \leftarrow \delta(e, x[i])$ 
4     $i \leftarrow i + 1$ 
5  tantque  $i < |x|$  faire
6    créer un état  $s$ 
7     $\delta(e, x[i]) \leftarrow s$ 
8     $e \leftarrow s$ 
9     $i \leftarrow i + 1$ 
10  $\text{sortie}(e) \leftarrow \{x\}$ 
```

# Algorithmes

## Compléter( $e$ )

```
1   $f \leftarrow$  file vide
2   $\ell \leftarrow$  liste des transitions  $(e, a, p)$  telles que  $p \neq e$ 
3  tantque  $\ell$  est non vide faire
4       $(r, a, p) \leftarrow \text{PREMIER}(\ell)$ 
5       $\ell \leftarrow \text{SUIVANT}(\ell)$ 
6       $\text{ENFILER}(f, p)$ 
7       $\text{sup}(p) \leftarrow e$ 
8  tantque  $f$  est non vide faire
9       $r \leftarrow \text{DÉFILER}(f)$ 
10      $\ell \leftarrow$  liste des transitions  $(r, a, p)$ 
11     tantque  $\ell$  est non vide faire
12          $(r, a, p) \leftarrow \text{PREMIER}(\ell)$ 
13          $\ell \leftarrow \text{SUIVANT}(\ell)$ 
14          $\text{ENFILER}(f, p)$ 
15          $s \leftarrow \text{sup}(r)$ 
16         tantque  $\delta(s, a)$  est non définie faire
17              $s \leftarrow \text{sup}(s)$ 
18          $\text{sup}(p) \leftarrow \delta(s, a)$ 
19      $\text{sortie}(e) \leftarrow \text{sortie}(e) \cup \text{sortie}(\text{sup}(p))$ 
```

# Complexité de la phase de prétraitement

La phase de prétraitement s'effectue en temps  $O(|X|)$ .

## Exemple

$y = \text{ctgagtagctag}$

## $AC(X, k, y, n)$

```
1  $e \leftarrow \text{PRE-AC}(X, k)$ 
2 pour  $j \leftarrow 0$  à  $n - 1$  faire
3   tantque  $\delta(e, a)$  est non définie faire
4      $e \leftarrow \text{sup}(e)$ 
5    $e \leftarrow \delta(e, y[j])$ 
6   si  $\text{sortie}(e) \neq \emptyset$  alors
7     reporter une occurrence des éléments de  $\text{sortie}(e)$ 
```



# Complexité

L'algorithme de Aho-Corasick trouve toutes les occurrences d'un ensemble  $X$  de  $k$  mots dans un texte  $y$  de longueur  $n$  en temps  $O(|X| + n)$  (sur un alphabet constant).

# Référence



Efficient string matching : an aid to bibliographic search

A. Aho et M.J. Corasick

*Communications of the ACM* **18**(6) (1975) 333–340.