

# Méthodes formelles pour XML : automates d'arbres

Pierre-Alain Reynier

<http://www.lif.univ-mrs.fr/~preynier/XML/>

# Plan du cours

## 1 - Motivations

## 2 - Automates

### 2.1 sur mots

### 2.2 sur arbres d'arité fixe

### 2.3 sur arbres d'arité variable

# I - Motivations

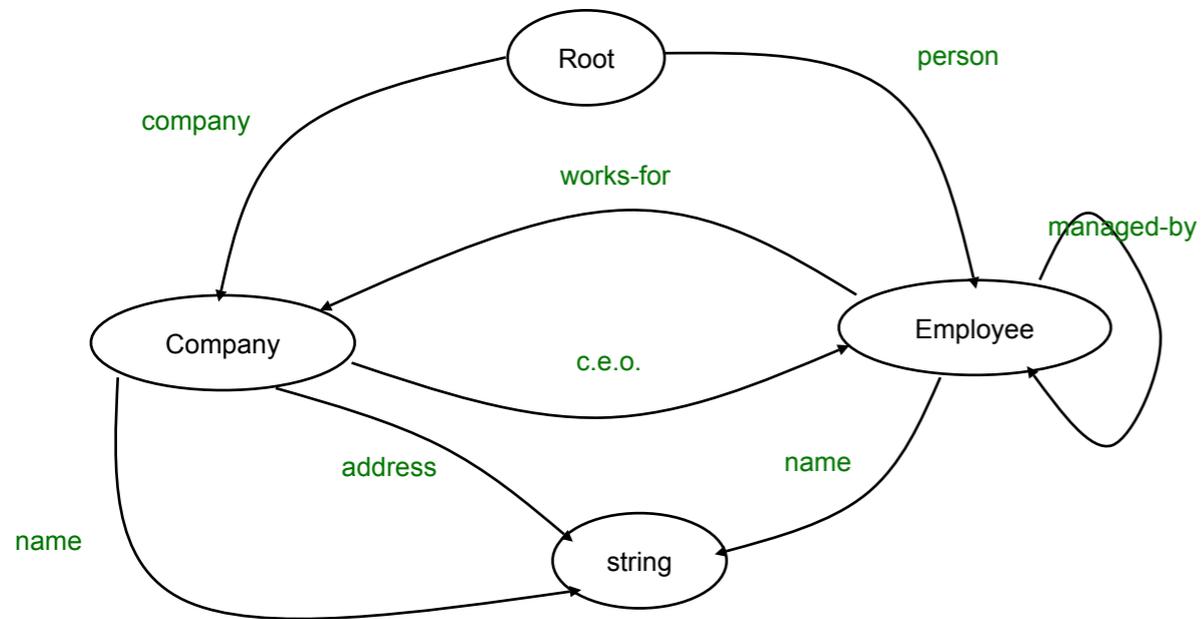
# Objectifs

- Les manipulations de documents XML doivent être sûres, et efficaces
- Besoin de méthodes permettant d'étudier
  - les types des documents XML produits,
  - les requêtes XPath,
  - les transformations XSLT, XQuery...
- Comprendre la puissance expressive et la complexité de ces opérations
- Identifier des fragments présentant un compromis intéressant

# XML Typing

- Type : DTD, Schéma XML, Relax NG...
- Simplifie l'écriture de programmes XML  
(améliore l'interopérabilité entre programmes)
- Améliore le stockage et les performances
- Facilite les requêtes
- Simplifie la protection des données  
(rejeter les mises à jour non conformes)

# Amélioration du stockage



Lower-bound schema

## Company

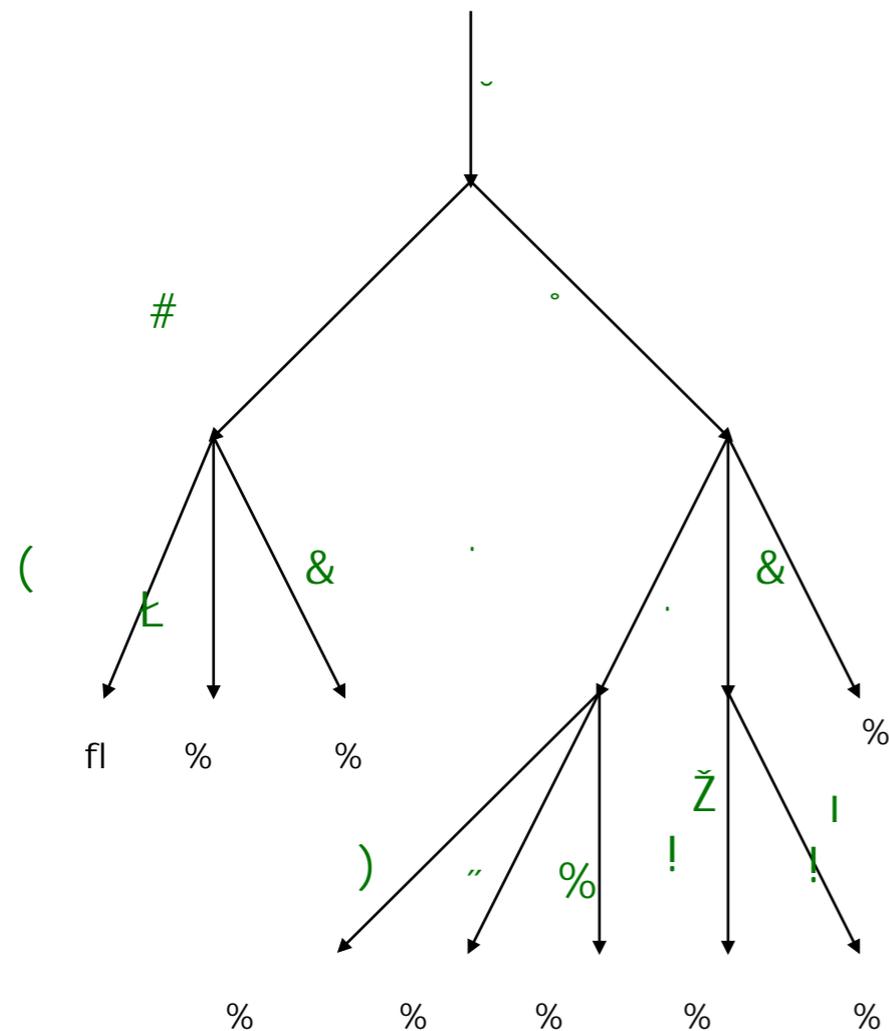
o i d	n a m e	a d d r e s s	c . e . o .
...	...	...	...
...	...	...	...

## Employee

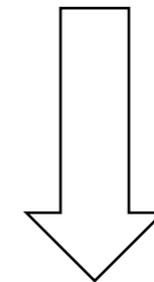
o i d	n a m e	m a n a g e d - b y	w o r k s - f o r
...	...	...	...
...	...	...	...

Store rest in overflow graph

# Amélioration des performances



```
select X.title  
from Bib._ X  
where X.*.zip = "12345"
```



```
select X.title  
from Bib.book X  
where X.address.zip = "12345"
```

# Manipulations de base

- **Validation**  
vérifier qu'un document XML est conforme à un type donné (DTD, Schéma XML...)
- **Navigation/Extraction**  
Sélectionner un ensemble de noeuds dans un document XML (XPath)
- **Transformation**  
Construire à partir d'un doc. XML un nouveau document XML (XSLT, XQuery..)

# Analyse dynamique/statique

- Contexte :
  - Un serveur reçoit un ou plusieurs docs XML respectant des contraintes de types
  - Il applique une transformation à ces données et produit un nouveau doc XML
  - Le document produit est renvoyé à un autre service et doit respecter un nouveau type
- Deux approches :
  - Vérification dynamique : après que les données soit produites
  - Vérification statique : vérification du programme produisant les données

# Problèmes d'analyse statique

- **Satisfaisabilité :**  
Étant donné une requête  $q$ , existe-t-il un doc. XML  $t$  tel que  $q(t)$  est non vide ?
- **Inclusion :**  
Étant donné deux requêtes  $q_1$  et  $q_2$ , est-ce que pour tout doc XML  $t$ , on a  $q_1(t) \subseteq q_2(t)$  ?
- **Equivalence :**  
Étant donné deux requêtes  $q_1$  et  $q_2$ , est-ce que pour tout doc XML  $t$ , on a  $q_1(t) = q_2(t)$  ?

# Problèmes d'analyse statique

- **Inclusion de type :**

Étant donné deux types  $d_1$ ,  $d_2$ , est-ce que tout doc. XML de type  $d_1$  est de type  $d_2$  ?

- **Vérification de type :**

Étant donné deux types  $d_1$ ,  $d_2$ , et une transformation  $T$ , est-ce que pour tout doc. XML  $t$  de type  $d_1$ ,  $T(t)$  est de type  $d_2$  ?

- **Inférence de type :**

Étant donné un type  $d$  et une transformation  $T$ , calculer le type de  $\{T(t) \mid t \text{ de type } d\}$

# Vérification de types

- Vérifié par :
  - Editeurs XML : vérifie que les données sont conformes à leur type
  - Echangeurs de XML (Services Web)
    - serveurs en envoyant les données
    - clients en les recevant
- Mieux : de façon statique  
assurer que l'application produit un document du bon type

# Vérification statique

- Entrée : type  $T$  et code d'une fonction  $f$   
( $f$  appartient à XQuery, XPATH, XSLT...)
- Vérification de  $T'$  :  
Est-ce que  $\forall d \models T, f(d) \models T'$  ?
- Inférence de types :  
Calculer le plus petit  $T'$  tel que  $\forall d \models T, f(d) \models T'$
- Rapidement indécidable à cause des jointures (comparaison d'éléments)

# Example

```
for $p in doc("parts.xml")//part[color="red"]
return <part>
    <name>$p/name</name>
    <desc>$p/desc</desc>
</part>
```

Type résultant :

(part (name (string) desc (any) )\*)

Si le type de parts.xml//part/desc est string :

(part (name (string) desc (string) )\*)

# Difficulté

for  $X$  in Input,  $Y$  in Input do { print (  $\langle b \rangle$  ) }

Input:  $\langle a \rangle \langle a \rangle \langle a \rangle$

Result:  $\langle b \rangle \langle b \rangle$

Problème:  $\{ b^i \mid i=n^2 \text{ for } n \geq 0 \}$  ne peut pas être décrit en XML Schema

Il n'y a pas de "meilleure approximation" :

–  $b^*$

–  $\varepsilon + b^2 b^*$

–  $\varepsilon + b^2 + b^4 b^*$

–  $\varepsilon + b^2 + b^4 + b^9 b^*$

– ...

# Modèle de données

- Un choix naturel : les arbres
  - reflètent la structure des doc XML
  - modèle de données sous-jacent à XML basé sur les arbres
- Limitations :
  - ne modélisent pas ID, IDREF
  - parfois des extensions sont nécessaires
- Dans la suite, nous considérerons cette vision arborescente des doc XML

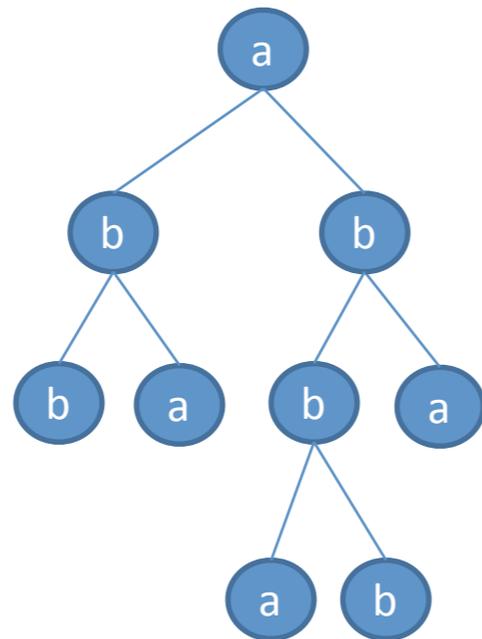
# Aperçu sur les aut. d'arbres

- Théorie riche pour **les mots** : Automates finis
- Etendue en une théorie riche pour **les arbres d'arité fixe** : Automates d'arbres
  - Théorèmes fondamentaux bien connus
  - Algorithmes et complexité des problèmes de décision connus
- Extension aux **arbres d'arité variable** ?  
Oui!

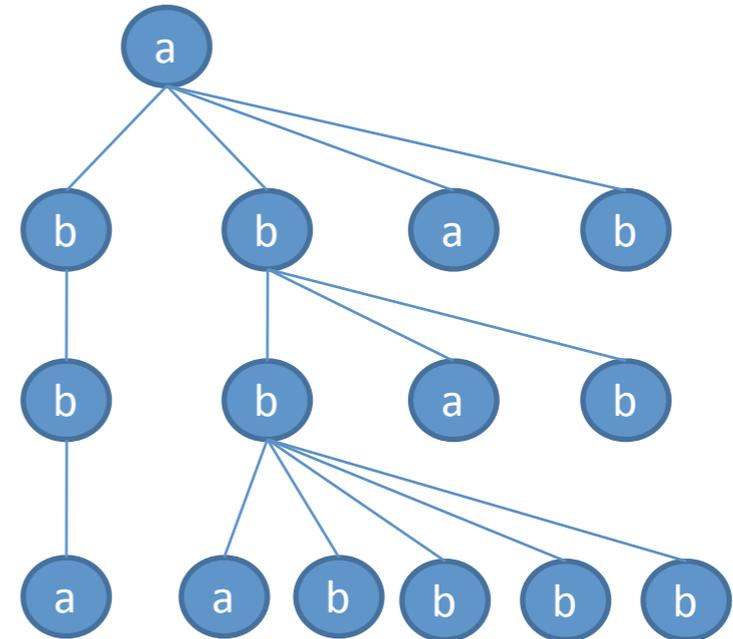
# Des mots aux arbres



Mots :  
Automates  
finis



Arbres binaires :  
Automates d'arbres  
d'arité 2

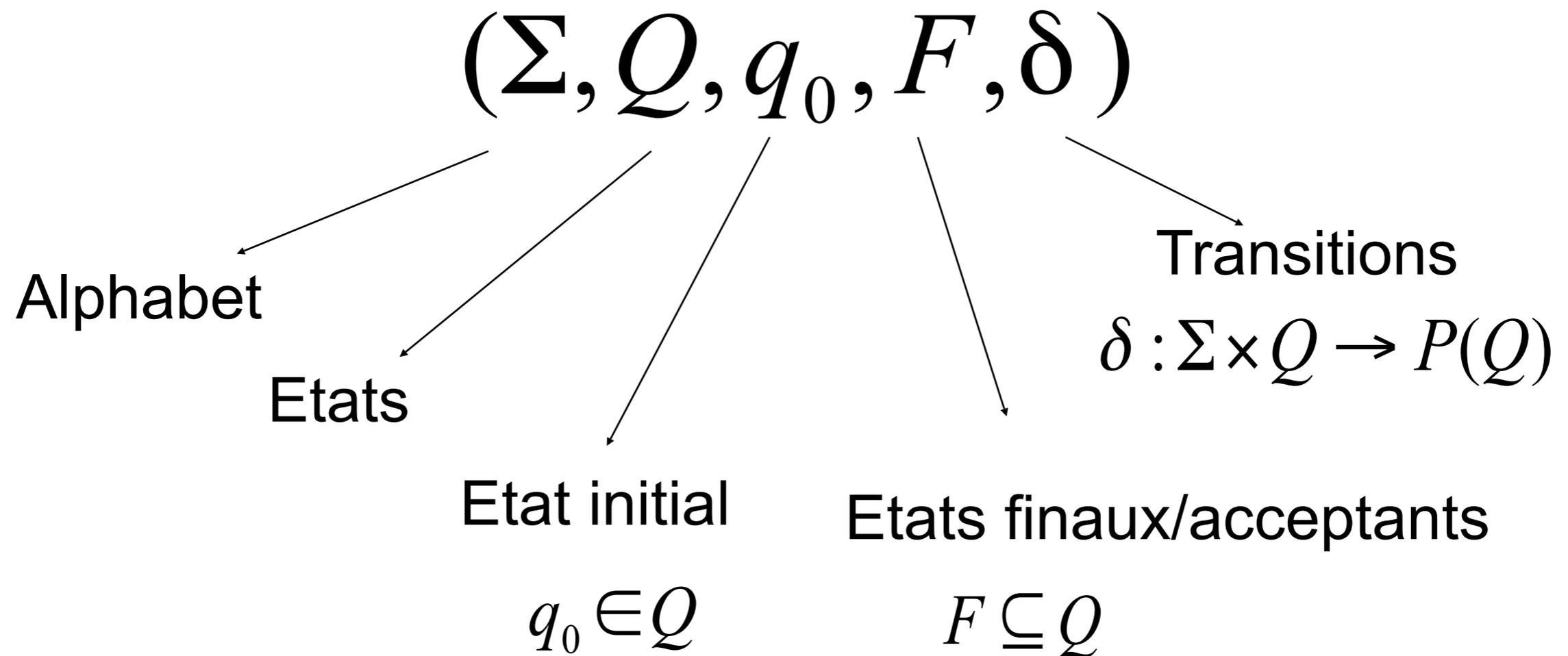


Arbres quelconques :  
Automates d'arbres  
d'arité variable

# **II - Automates**

## **- sur mots -**

# Automates finis sur mots



# Automates non déterministes

$$\Sigma = \{a, b\}$$

$$Q = \{q_0, q_1, q_2, q_3\}$$

$$F = \{q_2\}$$

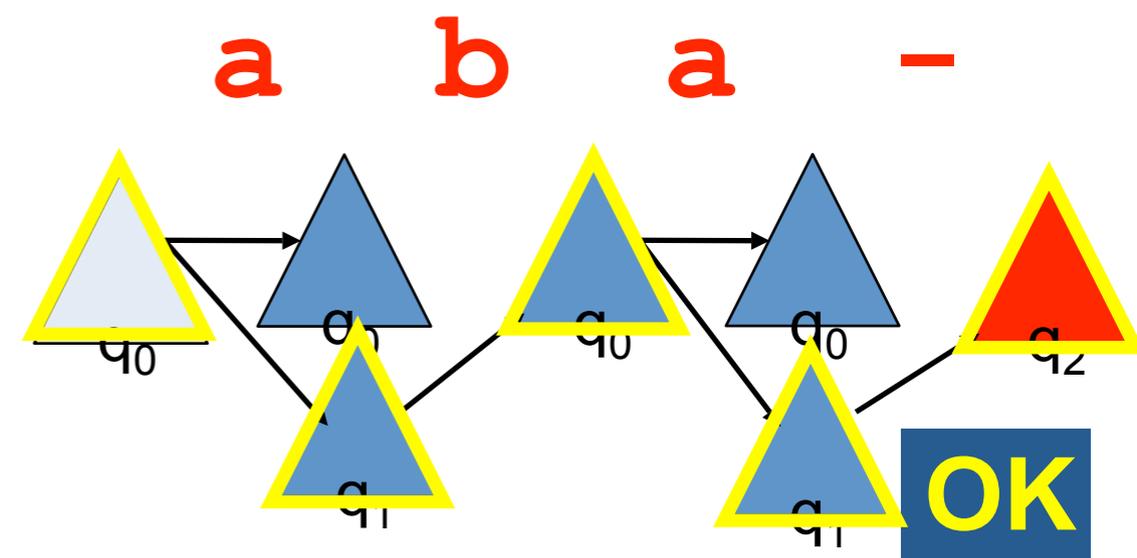
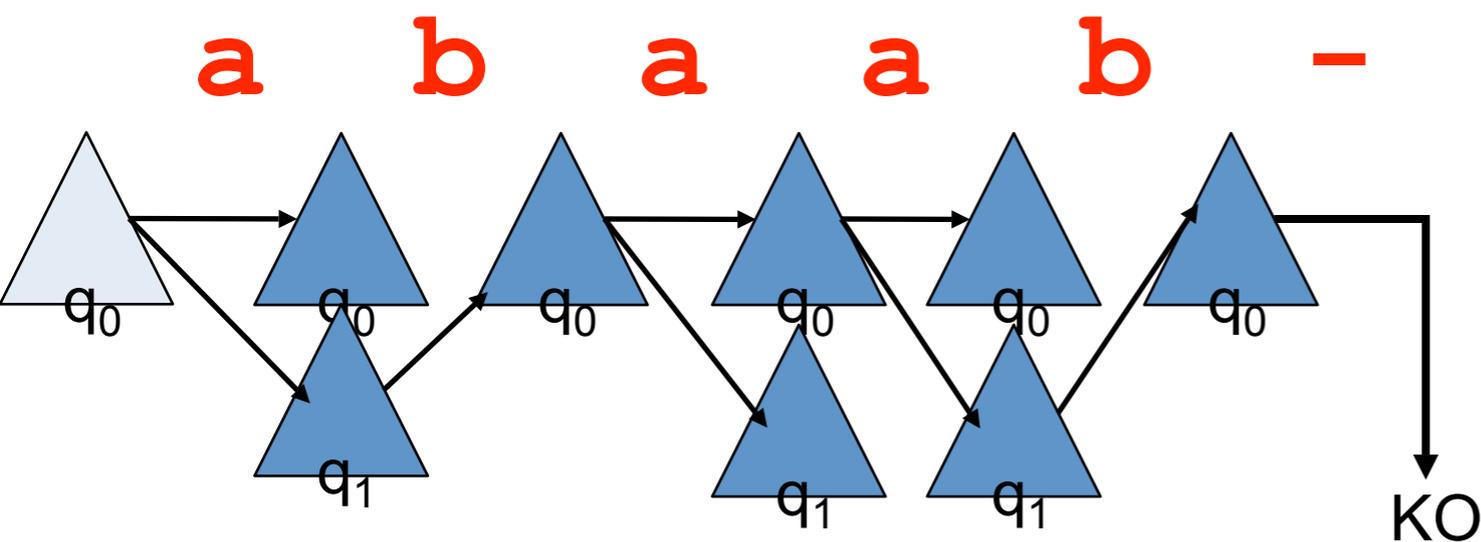
$$\delta(a, q_0) = \{q_0, q_1\}$$

$$\delta(b, q_1) = \{q_0\}$$

$$\delta(-, q_1) = \{q_2\}$$

$$\delta(-, q_2) = \{q_3\}$$

$$\delta(-, q_3) = \{q_3\}$$



# Rappels

- Déterministe :

- pas de transition  $\varepsilon$

$$\delta(-, q) = \{q_0\}$$

- pas de définition à choix comme

$$\delta(a, q_0) = \{q_0, q_1\}$$

- Déterminisation :

- possible de construire un équivalent déterministe

- les états sont des sous-ensembles de  $Q$

- potentielle explosion exponentielle

- Minimisation

- Limitations :  $\{a^n b^n, n \in \mathbb{N}\}$

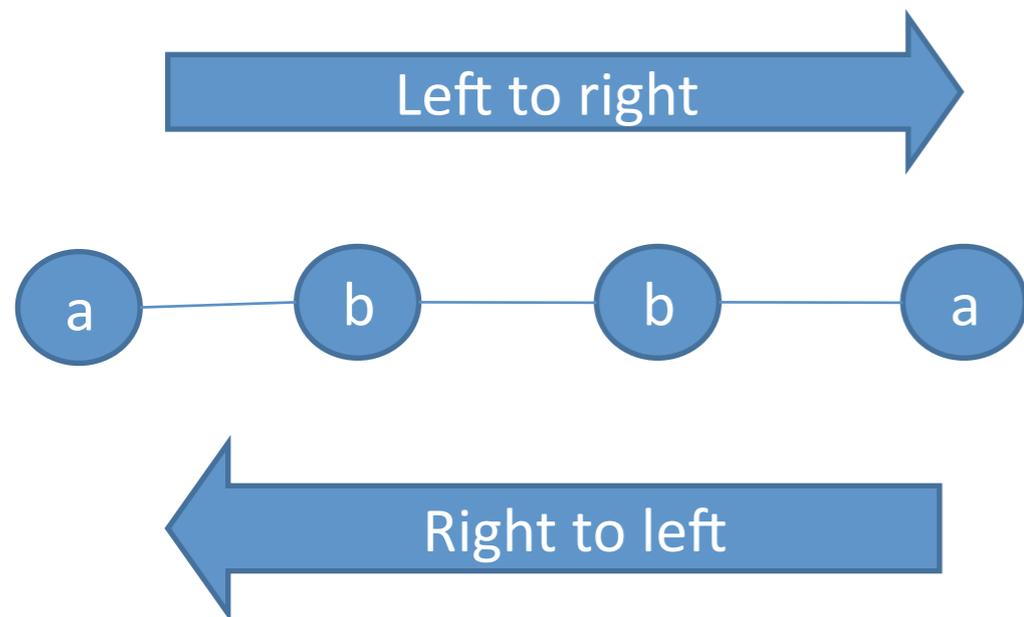
(plus faible que les langages hors-contexte)

- Outil très important (analyse lexicale)

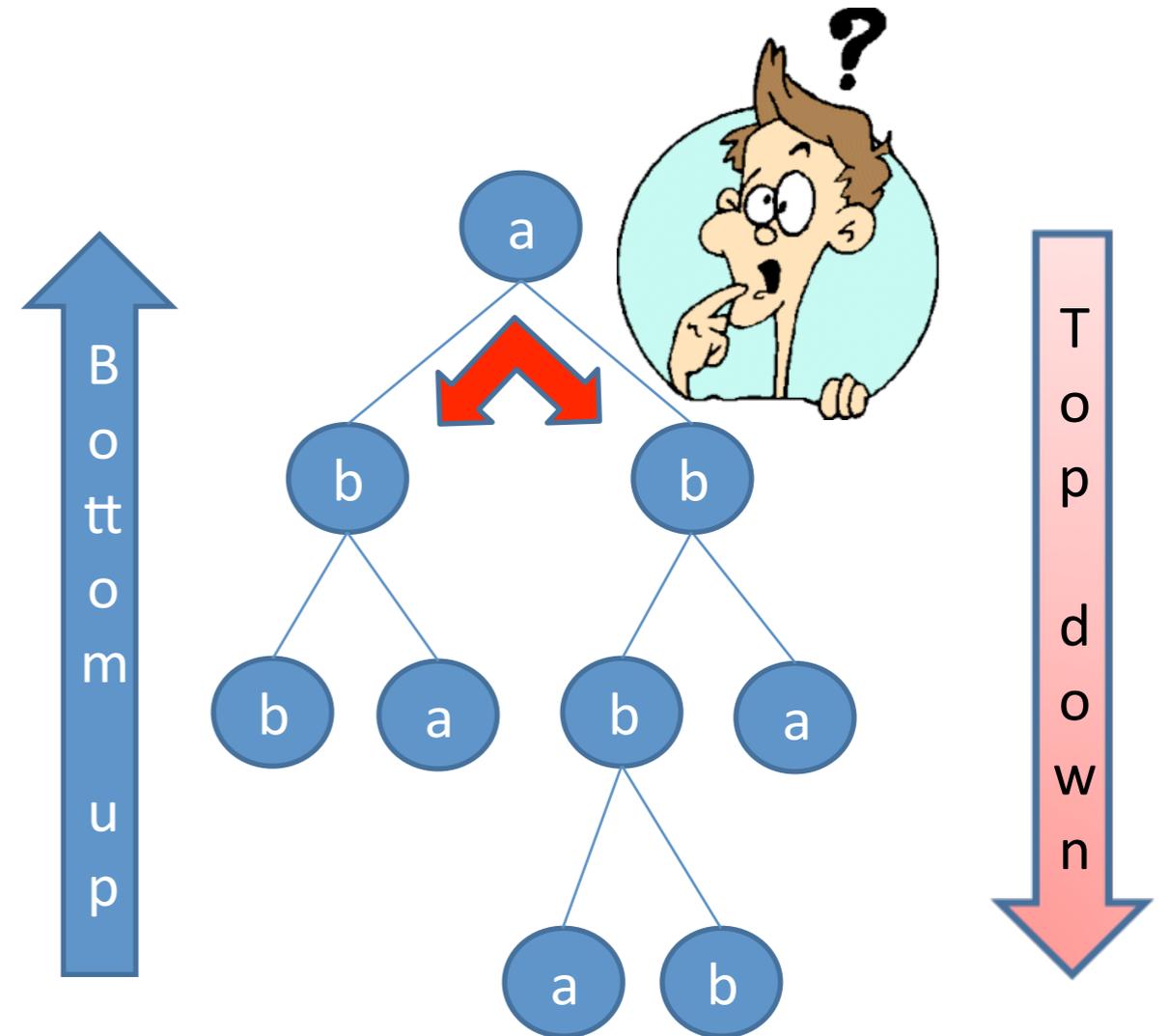
# Rappels (suite)

- $L(A)$  = langage accepté par l'automate  $A$
- Langages réguliers
- Peuvent être décrits par des expressions régulières :  $a(b+c)^*d$
- Fermés par complément :  $\Sigma^* - L(A)$
- Fermés par union, intersection :  $L(A) \cap L(B)$   
produit d'automates avec états  $(s, s')$   $L(A) \cup L(B)$   
où  $s$  (resp.  $s'$ ) est un état de  $A$  (resp.  $A'$ )

# Automates sur mots/ sur arbres



Pas de différence



Différences...

# II - Automates

- sur arbres d'arité fixe -

# Arbres binaires

- Evaluation parallèle

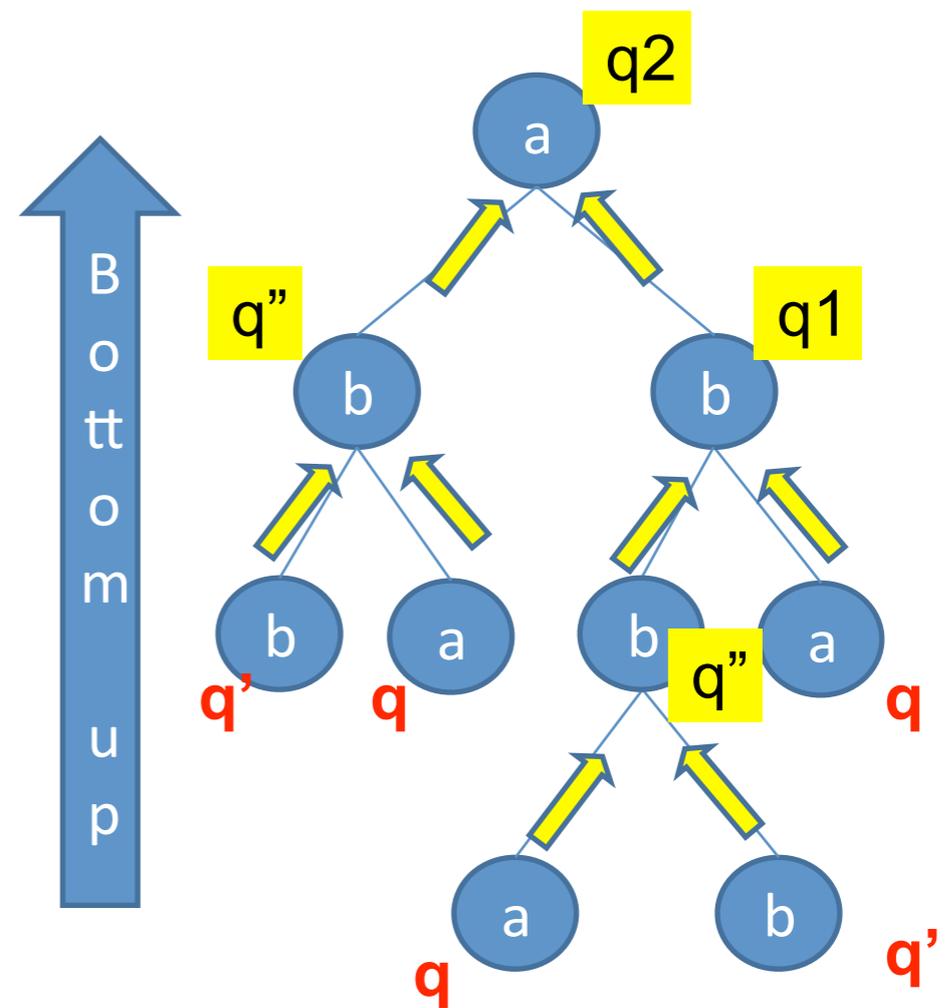
$$(\Sigma, Q, F, \delta)$$

- Pour les feuilles :

$$\delta : \Sigma \rightarrow P(Q)$$

- Pour les autres noeuds :

$$\delta : \Sigma \times Q \times Q \rightarrow P(Q)$$



# Automates d'arbres ascendants

$$\delta(a, q, q') \rightarrow \{r, r'\}$$

- Ascendant : si un noeud, étiqueté par  $a$ , a ses deux fils dans les états  $q$  et  $q'$ , alors ce noeud passe de façon non déterministe à l'état  $r$  ou  $r'$
- Arbre accepté si sa racine est dans un état de  $F$
- Non déterministe si transitions  $\varepsilon$  ou choix :

$$\delta(a, q, q') \rightarrow \{r, r'\} \quad \delta(\cdot\varepsilon, r) \rightarrow \{r'\}$$

# Exemple : circuits booléens

$$\delta_1(0) \rightarrow \{q_0\}$$

$$\delta_1(1) \rightarrow \{q_1\}$$

$$\Sigma = \{0, 1, \wedge, \vee\}$$

$$Q = \{q_0, q_1\}$$

$$F = \{q_1\}$$

$$\delta_2(\wedge, q_1, q_1) \rightarrow \{q_1\}$$

$$\delta_2(\wedge, q_0, q_0), \delta_2(\wedge, q_1, q_0), \delta_2(\wedge, q_0, q_1) \rightarrow \{q_0\}$$

$$\delta_2(\vee, q_0, q_0) \rightarrow \{q_0\}$$

$$\delta_2(\vee, q_1, q_1), \delta_2(\vee, q_1, q_0), \delta_2(\vee, q_0, q_1) \rightarrow \{q_1\}$$



# Langages réguliers

Définition :

Un langage d'arbres est dit *régulier* ssi il est accepté par un automate d'arbre ascendant.

# Langages d'arbres réguliers

- Les propositions suivantes sont équivalentes :
  - L est un langage d'arbres régulier
  - L est accepté par un automate d'arbre ascendant non déterministe
  - L est accepté par un automate d'arbre ascendant déterministe
  - L est accepté par un automate d'arbre descendant non déterministe
- Les automates descendants déterministes sont strictement moins expressifs

# Automates descendants

$$\delta(a, q'') \rightarrow \{(q, q')\}$$

- Si un noeud étiqueté  $a$  est dans un état  $q''$ , alors son fils gauche va dans l'état  $q$ , et son fils droit dans l'état  $q'$
- On part de la racine depuis un état initial  $q_0$ .
- L'arbre est accepté si toutes les feuilles sont dans un état final. Pour les feuilles : chgt d'état local.
- Non déterministe si :

$$\delta(a, q'') \rightarrow \{(q, q'), (r, r')\}$$

# Pourquoi déterminisme + descendant est plus faible ?

- Considérons le langage  $L = \{f(a,b), f(b,a)\}$
- $L$  peut être accepté par un automate ascendant  $A$  :  
 $L = L(A)$  (exercice)
- Supposons que  $L$  est accepté par un automate déterministe descendant  $B$  :  $L = L(B)$   
Exercice : montrer que  $B$  accepte  $f(a,a)$

Fait : aucun automate descendant déterministe n'accepte  $L$

# Propriétés arbres d'arité fixe

- Comme les mots, avec une complexité plus élevée
- Déterminisation
- Minimisation
- Fermés par :
  - complément
  - union
  - intersection

# Résumé Arbres d'arité fixe égale à $d$

## **Automates ascendants**

*Forme* :  $(\Sigma, Q, \partial, F)$

*Transitions* :

Feuilles :  $\partial : \Sigma \rightarrow P(Q)$

Noeuds internes :

$\partial : \Sigma \times Q^d \rightarrow P(Q)$

Arbre *accepté* ssi il  
existe un calcul t.q.  
*la racine est dans F*

## **Automates descendants**

*Forme* :  $(\Sigma, Q, \partial, I)$

La racine part de  $I$

*Transitions* :

Noeuds internes :

$\partial : \Sigma \times Q \rightarrow P(Q^d)$

Feuilles :  $\partial : \Sigma \times Q \rightarrow P(Q)$

Arbre *accepté* ssi il existe  
un *calcul terminant*

# Mais...

- Les documents XML sont d'arité variable
- Le type d'arbres que l'on souhaite manipuler :

`book (intro,section*,conclusion)`

## II - Automates

- sur arbres d'arité variable -

# Automates sur arbres d'arité variable : cas ascendant

$$\begin{array}{lll} \delta_2(\wedge, t) \rightarrow \{t\} & \delta_2(\wedge, t, t) \rightarrow \{t\} & \delta_2(\wedge, t, t, t) \rightarrow \{t\} \dots \\ \delta_2(\wedge, f) \rightarrow \{f\} & \delta_2(\wedge, t, f) \rightarrow \{f\} & \delta_2(\wedge, f, t) \rightarrow \{f\} \dots \\ \delta_2(\vee, t) \rightarrow \{t\} & \delta_2(\vee, t, f) \rightarrow \{t\} & \delta_2(\vee, f, t) \rightarrow \{t\} \dots \\ \delta_2(\vee, f) \rightarrow \{f\} & \delta_2(\vee, f, f) \rightarrow \{f\} & \delta_2(\vee, f, f, f) \rightarrow \{f\} \dots \end{array}$$

**Problème : représenter un ensemble infini  
de transitions**

**Solution : utiliser des expressions régulières**

# Automates de haies ascendants

- Règles :  $\delta(a, L(Q)) \rightarrow \{r_1, \dots, r_m\}$
- Feuilles :  $\partial : \Sigma \rightarrow P(Q)$
- Sémantique : si les enfants d'un noeud étiqueté par  $a$  forment un mot dans  $L(Q)$ , alors ce noeud passe dans un état de  $\{r_1, \dots, r_m\}$

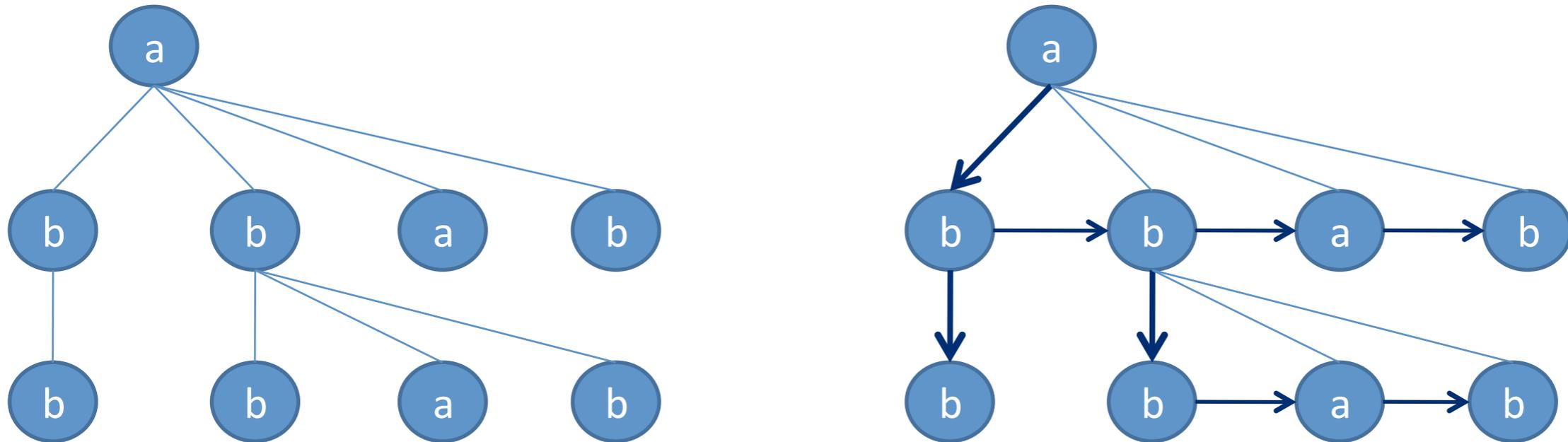
$$\delta_2(\wedge, And1) \rightarrow \{t\} \quad \text{where} \quad And1 = t +$$

$$\delta_2(\wedge, And0) \rightarrow \{f\} \quad \text{where} \quad And0 = (t + f)^* f (t + f)^*$$

$$\delta_2(\vee, Or1) \rightarrow \{t\} \quad \text{where} \quad Or1 = (t + f)^* t (t + f)^*$$

$$\delta_2(\vee, Or0) \rightarrow \{f\} \quad \text{where} \quad Or0 = f +$$

# Construction à partir du cas borné



Arbre de degré 2 : au plus 2 fils  
Premier fils - Prochain frère

F: codage dans un arbre d'arité 2 :

- F est une bijection

$F^{-1}$ : décodage

# Construction à partir du cas borné (2)

- Pour tout automate de haies  $A$ , il existe un automate d'arbres d'arité 2 acceptant  $F(L(A))$
- Pour tout automate d'arbre  $A$  d'arité 2, il existe un automate de haies acceptant  $F^{-1}(L(A))$
- Les deux sont faciles à construire

Conséquence : les langages de haies sont clos par union, intersection, complément.

# Déterminisation

- Déf. de déterministe dans le cas ascendant d'arité variable :

$\forall \alpha \in \Sigma, \forall w \in Q^*$ , there exists a unique rule  $\delta(\alpha, L)$   
such that  $w \in L$ .

- Tout automate d'arbre ascendant d'arité variable peut être déterminisé
- Peut-on utiliser l'encodage "Premier fils - Prochain suivant" ?  
Non ! il ne préserve pas le déterminisme

# Automates descendants ?

- On part de la racine et on descend
- On lit l'étiquette du noeud, son état, et on donne un état à chaque fils, selon une expression régulière
- Au lieu de donner un tuple d'états (cas d'arité fixe), on donne une expression régulière sur les états.

# Cas descendant (suite)

*Forme* :  $(\Sigma, Q, \partial, l)$  ; La racine part de  $l$

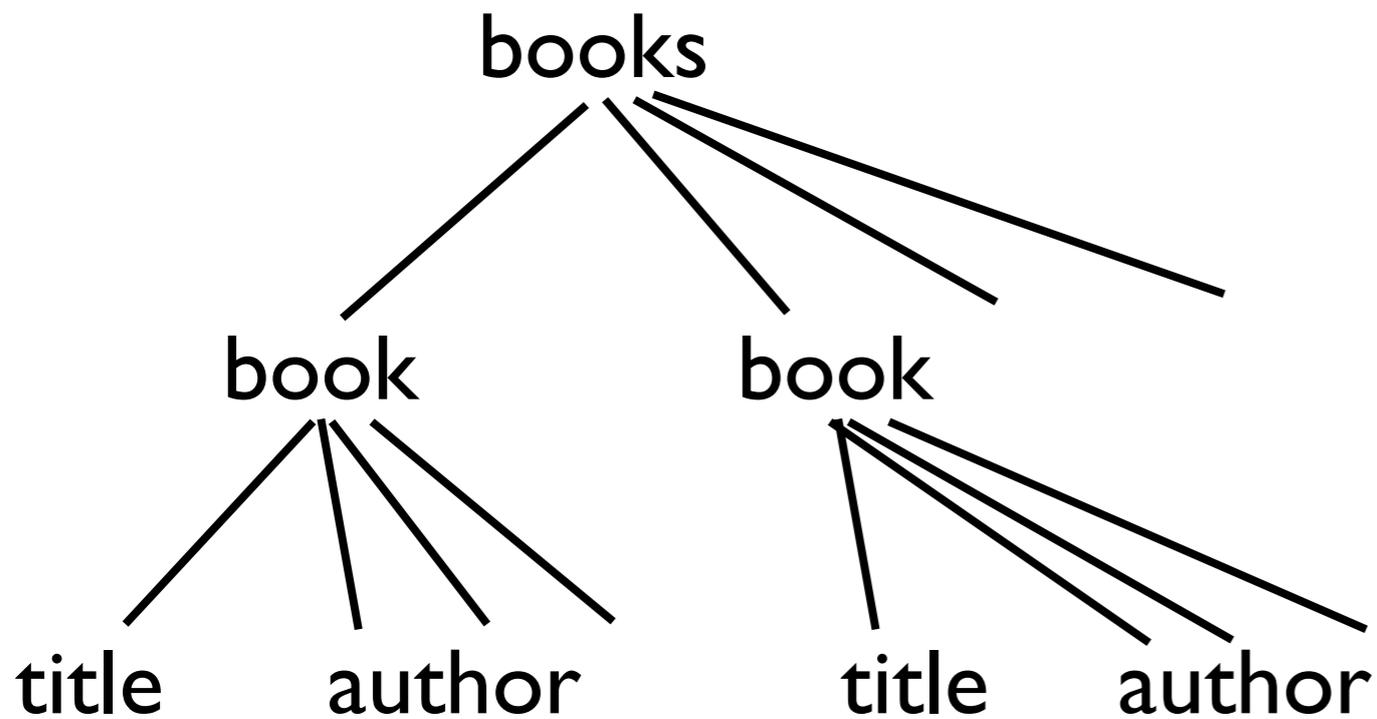
*Transitions* : Noeuds internes :  $\partial : \Sigma \times Q \rightarrow R$   
où  $R$  est dans  $\text{Reg}(Q)$

*Arbre accepté* ssi il existe un *calcul terminant*

L'automate est *déterministe* s'il propose au plus une expression rationnelle.

*Remarque* : Cette expression est toujours représentable par un automate fini déterministe.

# Cas descendant : exemple



- Un élément racine books possédant des fils book
- Un élément book possède un fils title et (au moins un) fils author

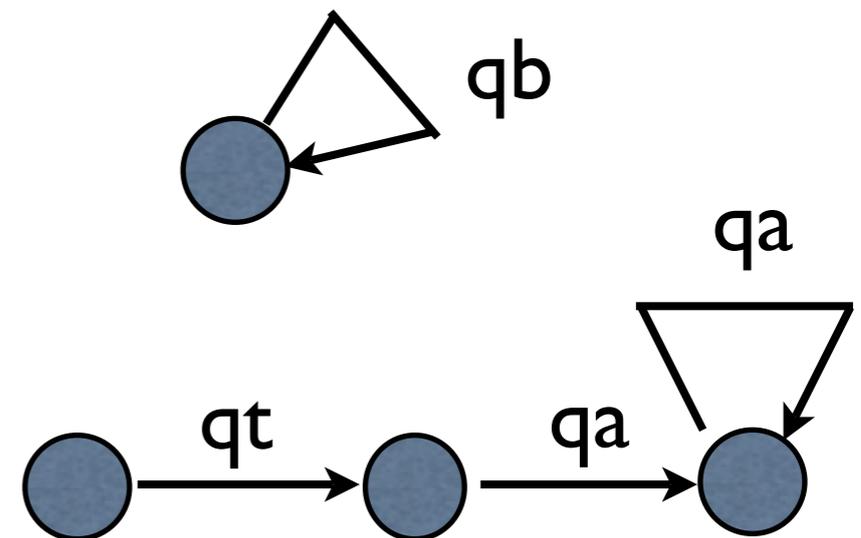
$$\Sigma = \{\text{books, book, title, author}\}$$

$$Q = \{q_0, q_b, q_s, q_t, q_a, q_f\} \quad I = \{q_0\}$$

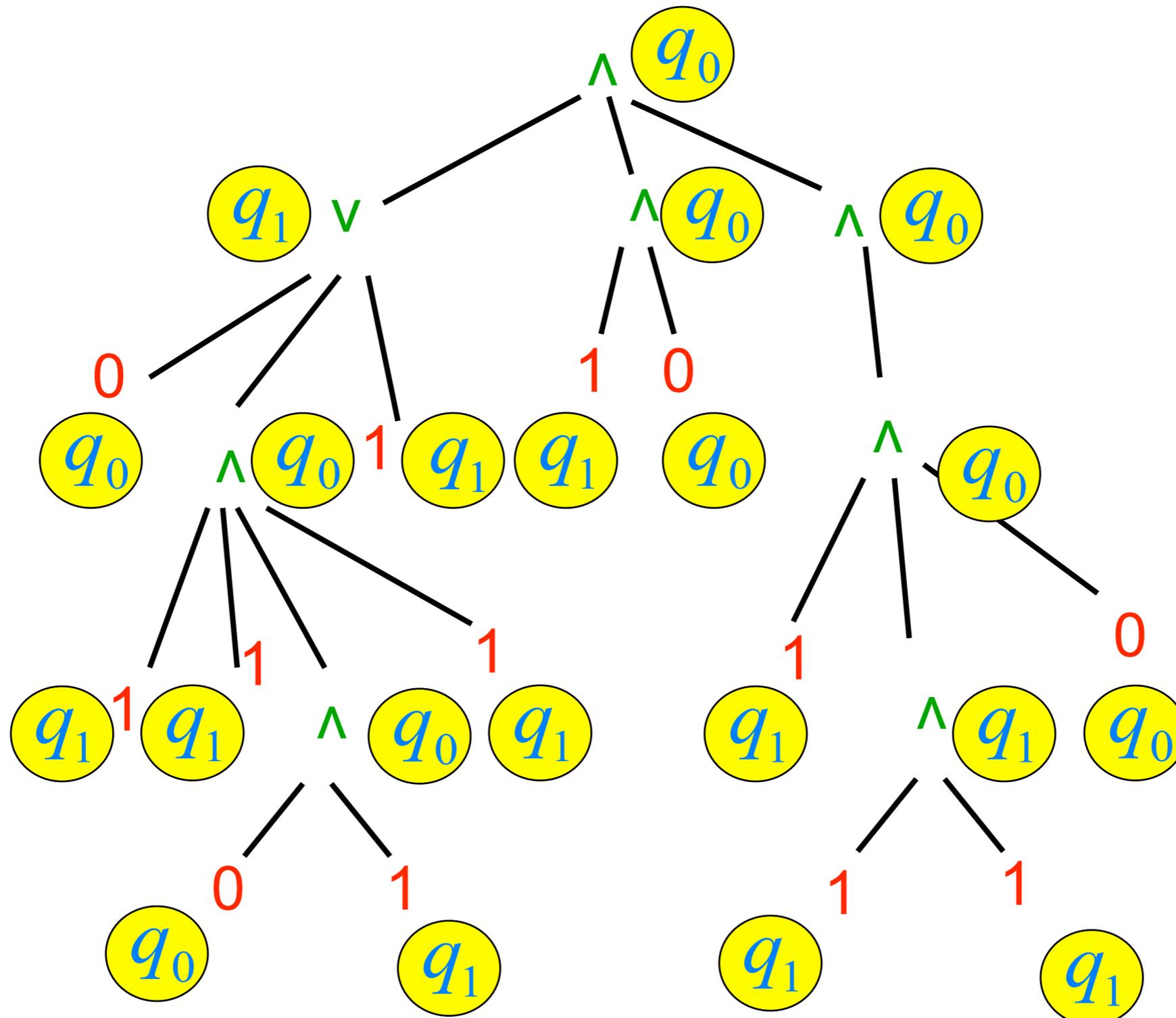
$$\partial(\text{books}, q_0) = \{q_b^*\}$$

$$\partial(\text{book}, q_b) = \{q_t, q_a^+\}$$

$$\partial(\text{title}, q_t) = \partial(\text{author}, q_a) = q_f$$



# Evaluation de circuits booléens



# Résumé Arbres d'arité variable

## **Automates ascendants**

*Forme* :  $(\Sigma, Q, \partial, F)$

*Transitions* :

Feuilles :  $\partial : \Sigma \rightarrow P(Q)$

Noeuds internes :

$\partial : \Sigma \times \text{Reg}(Q) \rightarrow P(Q)$

*Arbre accepté* ssi il existe un calcul t.q. *la racine est dans F*

## **Automates descendants**

*Forme* :  $(\Sigma, Q, \partial, I)$

La racine part de I

*Transitions* :

Noeuds internes :

$\partial : \Sigma \times Q \rightarrow P(\text{Reg}(Q))$

*Arbre accepté* ssi il existe un *calcul terminant*