

<presentation>

<title> **Programmation XML:
de la théorie aux outils** </title>

_{24^{ème} Journées Frankofônes des Langages Applicatifs}
<date>3-6 Février 2013, Aussois</date>

<author>**Kim Nguyễn**</author>
<email>kn@lri.fr</email>



Comprendre le monde,
construire l'avenir®

</presentation>

<title>**But de ce cours**</title>

Réciter 1400 pages de spécifications écrites par le **W3C**

<title>But de ce cours</title>

~~Réciter 1400 pages de spécifications écrites par le W3C~~

Donner un bref aperçu d'XML

<title>But de ce cours</title>

~~Réciter 1400 pages de spécifications écrites par le W3C~~

Donner un bref aperçu d'XML

Décrire quelques concepts théoriques

~~Réciter 1400 pages de spécifications écrites par le W3C~~

Donner un bref aperçu d'XML

Décrire quelques concepts théoriques

Présenter des outils issus de la recherche

~~Réciter 1400 pages de spécifications écrites par le W3C~~

Donner un bref aperçu d'XML

Décrire quelques concepts théoriques

Présenter des outils issus de la recherche

⇒ favoriser la réutilisation dans d'autres domaines (typeurs, solveurs SMT , *model checking* ...)

<title>XML?</title>

<description>

<item>XML est un langage de balisage</item>

<item>Standardisé par le W3C</item>

<item>Un document XML est composé</item>

<description>

<item>De texte dont des séquences d'&apostrophe;échappement</item>

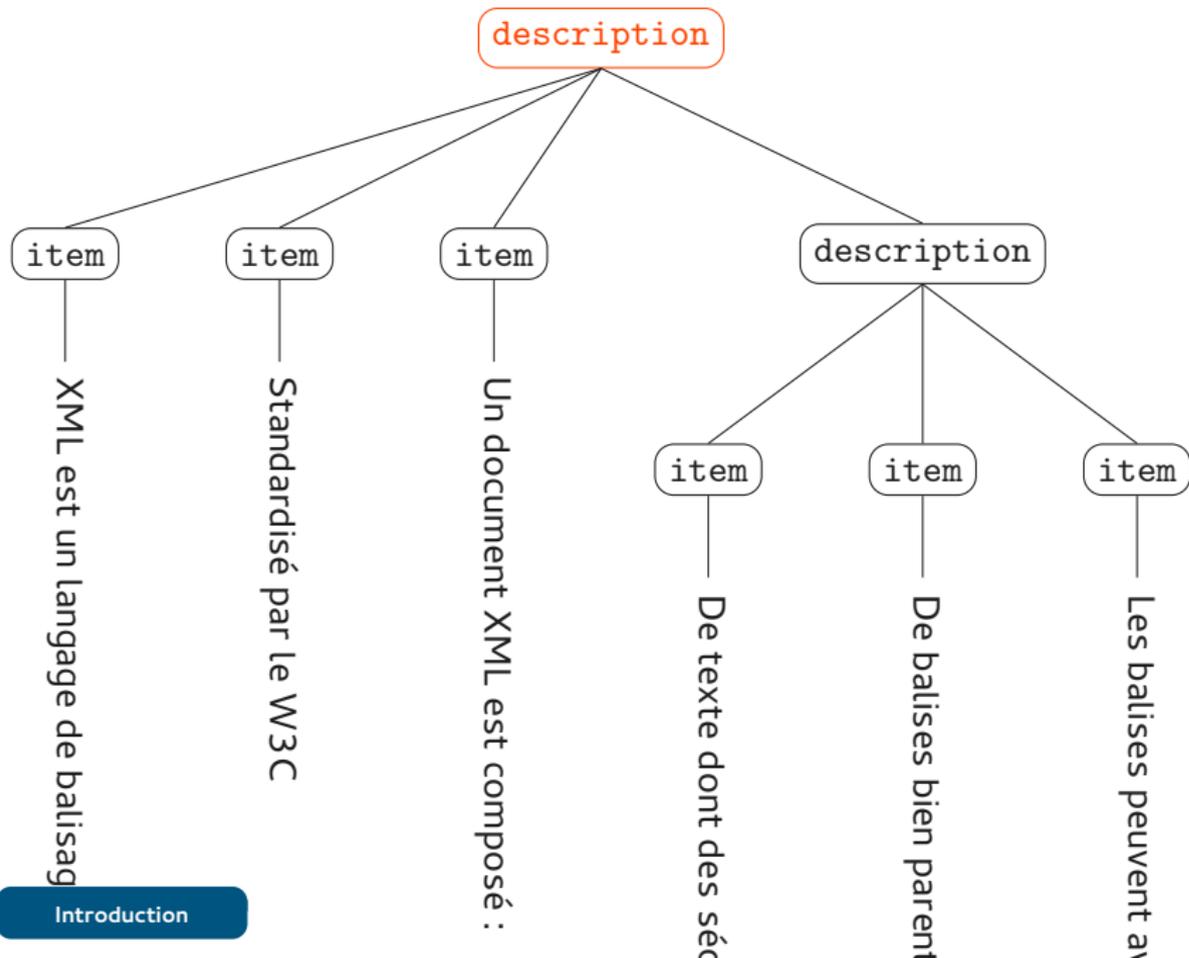
<item>De balises, bien parenthésées (dont une racine)</item>

<item id="id1">Les balises peuvent avoir des attributs</item>

</description>

</description>

<title>Document XML = Arbre</title>



- Données semi-structurées (par oppositions aux tables relationnelles)
- Représentation arborescente de données
- Utilisations multiples :
 - Standards du web XHTML, Atom, RSS, SOAP, SVG, MathML, ...
 - Documents OpenDocument format
 - Bio-informatique UniProt XML, SBML
 - Finance XBRL
 - Linguistique Treebank, PAULA XML, ...

<title>Quels avantages ?</title>

- 1 Standardisé par le W3C
- 2 Cadre commun pour la sérialisation et l'échange de données
- 3 Bibliothèques de *parsing* génériques
- 4 Possibilité de *spécifier* les documents
- 5 Langage(s) standards de manipulation de documents (XSLT, XPath, ...)
- 6 *Mapping* standardisé entre document XML et langages objets (DOM)

<title>Quels inconvénients ?</title>

- 1 Standardisé par le W3C
- 2 Verbeux
- 3 DOM est mal adapté aux langages fonctionnels
- 4 Spécifications très informelles

<title>Au programme</title>

Pratique

1. Validation
de documents

3. CDuce et le
sous-typage
sémantique

4. XPath,
XQuery,
encodage
d'arbres par
intervalles

2. Automates
d'arbres

5. Logiques
pour les arbres

Théorie

1. Validation de documents

On veut pouvoir définir des contraintes sur un document :

- 1 ensemble des noms de balises autorisés
- 2 imbrication et contenu des balises
- 3 (nom et valeurs des attributs)
- 4 (« forme » du texte : date, nombre, ...)

Il y a plusieurs standards pour cela :

- 1 Document Type Definition (DTD)
- 2 (XML Schema)
- 3 (RelaxNG)

<title>Document Type Definition</title>

- 1 Définit le contenu de chaque balise par une expression régulière
- 2 Ne permet pas d'imposer des contraintes sur le texte

Exemple :

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item             (#PCDATA | b | i )* >
<!ELEMENT i                (#PCDATA | b | i )* >
<!ELEMENT b                (#PCDATA | b | i )* >
```

<title>Document Type Definition</title>

- 1 Définit le contenu de chaque balise par une expression régulière
- 2 Ne permet pas d'imposer des contraintes sur le texte

Exemple :

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item             (#PCDATA | b | i )* >
<!ELEMENT i                (#PCDATA | b | i )* >
<!ELEMENT b                (#PCDATA | b | i )* >
```

Contrainte : le document doit être vérifiable en *streaming*

Remarque : on s'autorise une pile proportionnelle à la profondeur de l'arbre

<title>Exemple de validation en streaming</title>

```
<!ELEMENT description      (item | description)+ >  
<!ELEMENT item             (#PCDATA | b | i )*>  
<!ELEMENT i                (#PCDATA | b | i )* >  
<!ELEMENT b                (#PCDATA | b | i )* >
```

<title>Exemple de validation en streaming</title>

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item              (#PCDATA | b | i )*>
<!ELEMENT i                  (#PCDATA | b | i )* >
<!ELEMENT b                  (#PCDATA | b | i )* >
```

Le parseur renvoie la séquence d'évènements :

```
<description>
```

<title>Exemple de validation en streaming</title>

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item             (#PCDATA | b | i )*>
<!ELEMENT i                (#PCDATA | b | i )* >
<!ELEMENT b                (#PCDATA | b | i )* >
```

Le parseur renvoie la séquence d'évènements :

```
<description>
  <item>
```

<title>Exemple de validation en streaming</title>

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item              (#PCDATA | b | i )*>
<!ELEMENT i                 (#PCDATA | b | i )* >
<!ELEMENT b                 (#PCDATA | b | i )* >
```

Le parseur renvoie la séquence d'évènements :

```
<description>
  <item>
    Blah blah
```

<title>Exemple de validation en streaming</title>

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item              (#PCDATA | b | i )*>
<!ELEMENT i                  (#PCDATA | b | i )* >
<!ELEMENT b                  (#PCDATA | b | i )* >
```

Le parseur renvoie la séquence d'évènements :

```
<description>
  <item>
    Blah blah
  <b>
```

<title>Exemple de validation en streaming</title>

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item             (#PCDATA | b | i )*>
<!ELEMENT i                (#PCDATA | b | i )* >
<!ELEMENT b                (#PCDATA | b | i )* >
```

Le parseur renvoie la séquence d'évènements :

```
<description>
  <item>
    Blah blah
    <b>
</description>
```

<title>Exemple de validation en streaming</title>

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item              (#PCDATA | b | i )*>
<!ELEMENT i                 (#PCDATA | b | i )* >
<!ELEMENT b                 (#PCDATA | b | i )* >
<!ELEMENT b                 description+ >
```

Le parseur renvoie la séquence d'évènements :

```
<description>
  <item>
    Blah blah
    <b>
</description>
```

<title>Exemple de validation en streaming</title>

```
<!ELEMENT description      (item | description)+ >
<!ELEMENT item              (#PCDATA | b | i )*, b>
<!ELEMENT i                  (#PCDATA | b | i )* >
<!ELEMENT b                  (#PCDATA | b | i )* >
```

Le parseur renvoie la séquence d'évènements :

```
<description>
  <item>
    Blah blah
    <b>
</description>
```

- 1 Une seule production par nom de balise
- 2 Expression 1-nonambiguë [*Brüggemann-Klein, Wood, 97*]

- 1 Une seule production par nom de balise
- 2 Expression 1-nonambiguë [*Brüggemann-Klein, Wood, 97*]

Exemple :

$(a \mid b)^* a$

- 1 Une seule production par nom de balise
- 2 Expression 1-nonambiguë [*Brüggemann-Klein, Wood, 97*]

Exemple :

$(a \mid b)^* a \rightsquigarrow (a_1 \mid b_1)^* a_2$

- 1 Une seule production par nom de balise
- 2 Expression 1-nonambiguë [*Brüggemann-Klein, Wood, 97*]

Exemple :

$(a \mid b)^* a \rightsquigarrow (a_1 \mid b_1)^* a_2$

Considérons le mot baa :

- 1 on lit b, on sait que c'est b_1
- 2 on lit a, on ne sait pas si c'est a_1 ou a_2

- 1 Une seule production par nom de balise
- 2 Expression 1-nonambiguë [*Brüggemann-Klein, Wood, 97*]

Exemple :

$(a \mid b)^* a \rightsquigarrow (a_1 \mid b_1)^* a_2$

Considérons le mot baa :

- 1 on lit b, on sait que c'est b_1
- 2 on lit a, on ne sait pas si c'est a_1 ou a_2

Remarque : $b_1^* a_1 (b_2^* a_2)^*$ dénote le même langage et est nonambiguë

La 1-nonambiguïté est une propriété syntaxique

Construction simple a partir de l'expression régulière :

- $(a_1 \mid b_1)^* a_2$

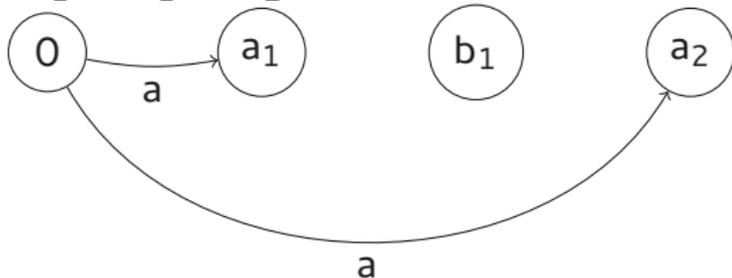


- Automate de Glushkov deterministe \Leftrightarrow expression 1-nonambiguë
- Validation en *streaming* (pile $o(\text{depth}(D))$) [Segoufin, Vianu 2002]

<title>Automate de Glushkov</title>

Construction simple a partir de l'expression régulière :

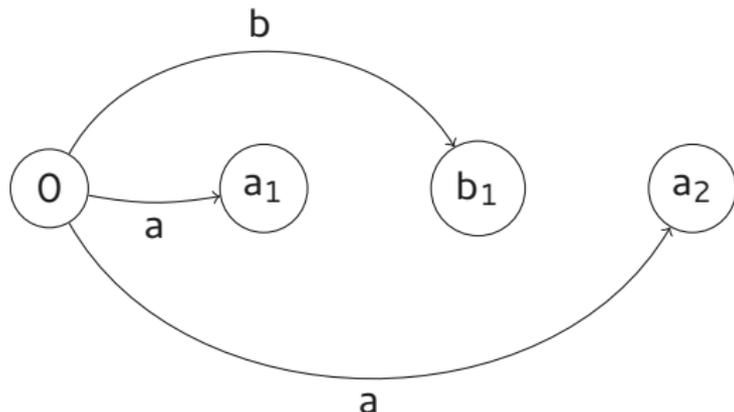
- $(a_1 \mid b_1)^* a_2$



- Automate de Glushkov déterministe \Leftrightarrow expression 1-nonambiguë
- Validation en *streaming* (pile $o(\text{depth}(D))$) [Segoufin, Vianu 2002]

Construction simple a partir de l'expression régulière :

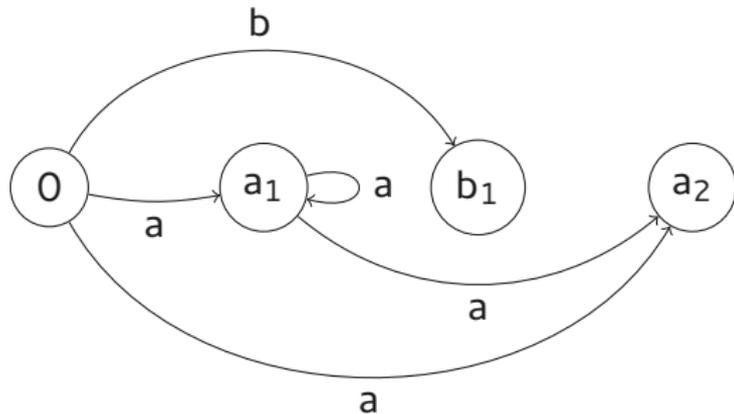
- $(a_1 \mid b_1)^* a_2$



- Automate de Glushkov déterministe \Leftrightarrow expression 1-nonambiguë
- Validation en *streaming* (pile $o(\text{depth}(D))$) [Segoufin, Vianu 2002]

Construction simple a partir de l'expression régulière :

- $(a_1 \mid b_1)^* a_2$

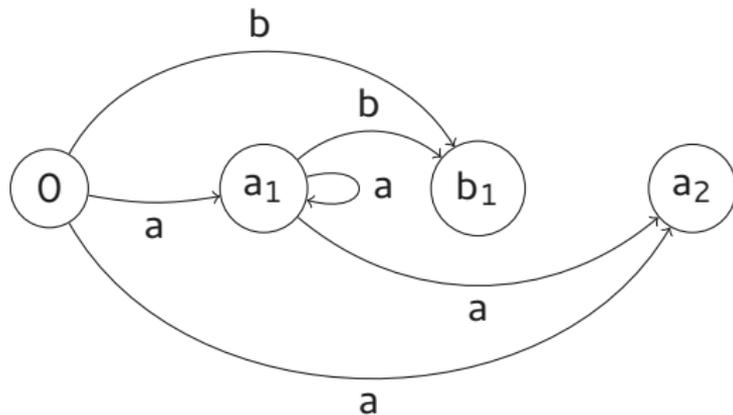


- Automate de Glushkov déterministe \Leftrightarrow expression 1-nonambiguë
- Validation en *streaming* (pile $o(\text{depth}(D))$) [Segoufin, Vianu 2002]

<title>Automate de Glushkov</title>

Construction simple a partir de l'expression régulière :

- $(a_1 \mid b_1)^* a_2$

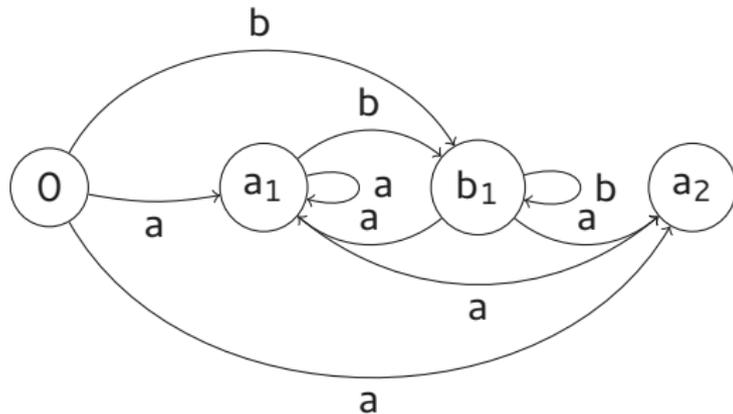


- Automate de Glushkov déterministe \Leftrightarrow expression 1-nonambiguë
- Validation en *streaming* (pile $o(\text{depth}(D))$) [Segoufin, Vianu 2002]

<title>Automate de Glushkov</title>

Construction simple a partir de l'expression régulière :

- $(a_1 \mid b_1)^* a_2$

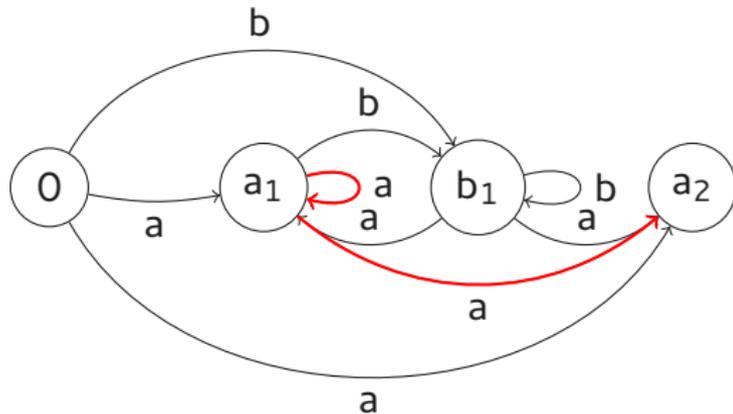


- Automate de Glushkov déterministe \Leftrightarrow expression 1-nonambiguë
- Validation en *streaming* (pile $o(\text{depth}(D))$) [Segoufin, Vianu 2002]

<title>Automate de Glushkov</title>

Construction simple a partir de l'expression régulière :

- $(a_1 \mid b_1)^* a_2$

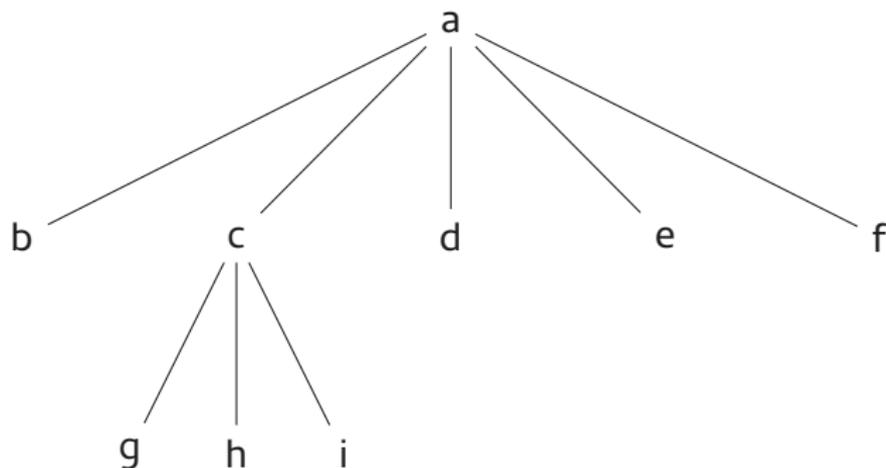


- Automate de Glushkov déterministe \Leftrightarrow expression 1-nonambiguë
- Validation en *streaming* (pile $o(\text{depth}(D))$) [Segoufin, Vianu 2002]

2. Automates d'arbres

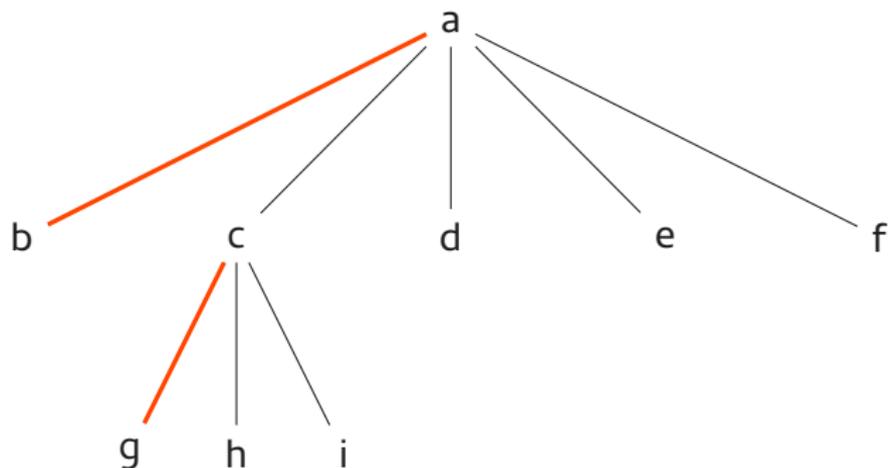
<title>Automates d'arbres (préliminaire)</title>

Encodage d'un arbre n-aire en arbre binaire :



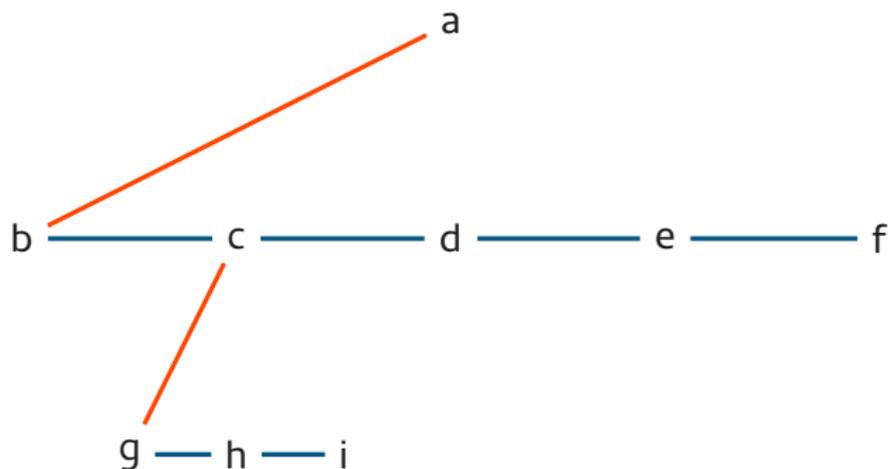
<title>Automates d'arbres (préliminaire)</title>

Encodage d'un arbre n-aire en arbre binaire :



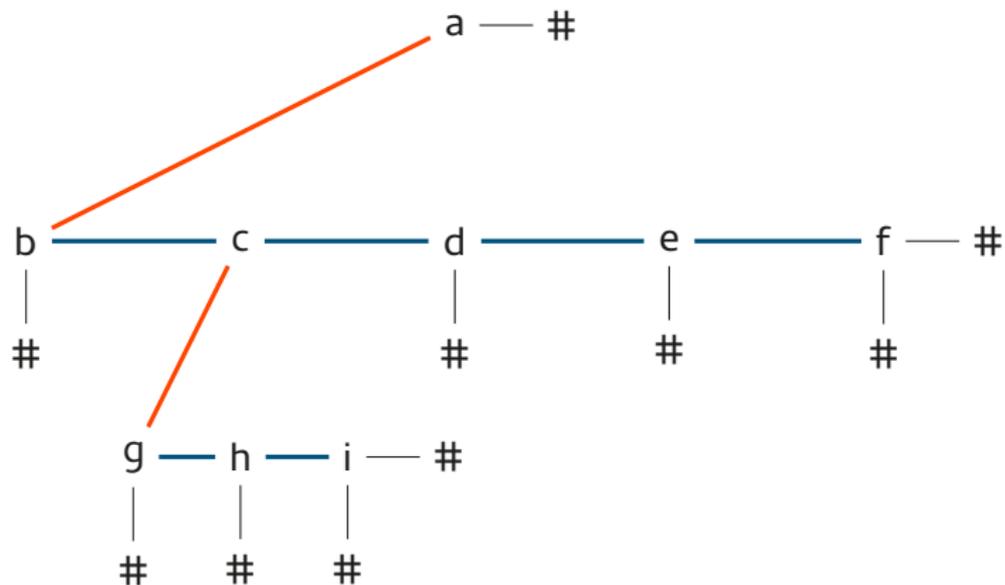
<title>Automates d'arbres (préliminaire)</title>

Encodage d'un arbre n-aire en arbre binaire :



<title>Automates d'arbres (préliminaire)</title>

Encodage d'un arbre n-aire en arbre binaire :



Caractériser des ensembles d'arbres :

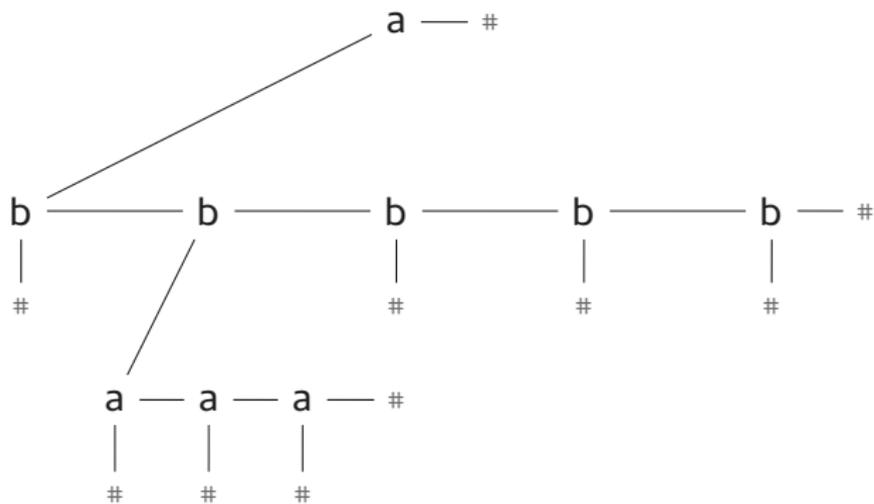
« Arbres ayant des a à profondeur paire et b à profondeur impaire (dans l'arbre n-aire)»

Definition (Automate d'arbres ([Comon et al. 2007]))

Un automate d'arbre est un tuple $(\Sigma, Q, T, B, \delta)$:

- Σ : alphabet (ensemble de symboles et leur arité)
- Q : ensemble d'états
- $T \subseteq Q$: états à la racine (Top)
- $B \subseteq Q$: états aux feuilles (Bottom)
- $\delta \subseteq Q \times \Sigma \times Q \times Q$: relation de transition

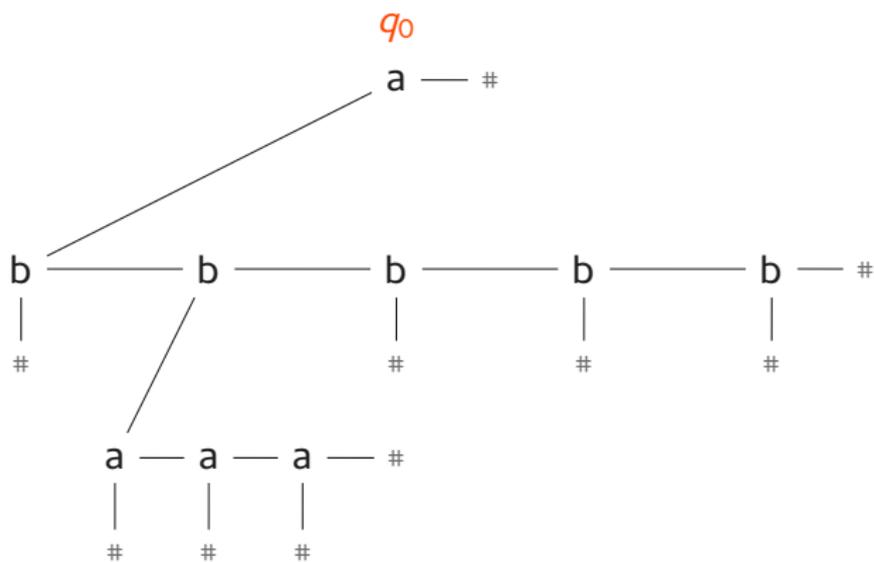
<title>Automate d'arbres </title>



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_\perp\}, \{q_0\}, \{q_0, q_1\},$$

q_0, a	(q_1, q_0)
q_0, b	(q_\perp, q_\perp)
q_1, b	(q_0, q_1)
q_1, a	(q_\perp, q_\perp)
$q_\perp, -$	(q_\perp, q_\perp)

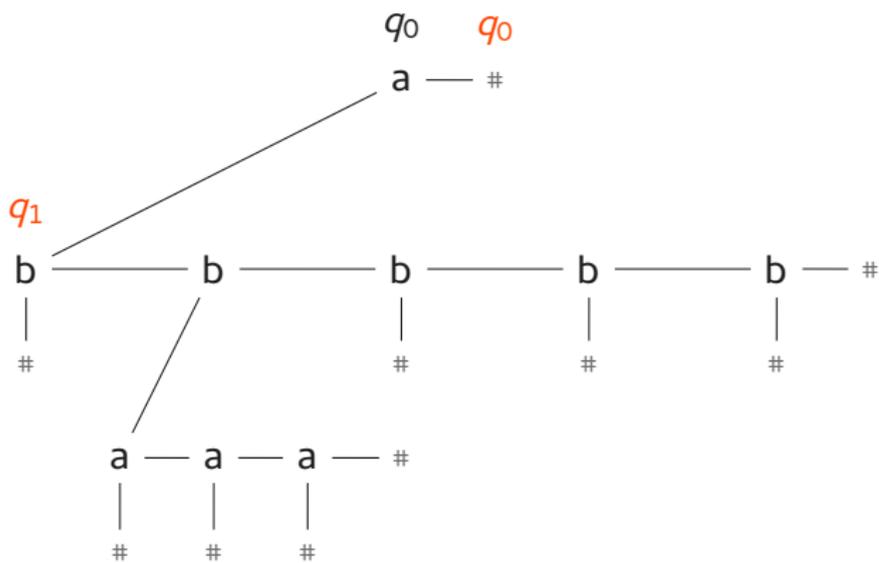
<title>Automate d'arbres </title>



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_{\perp}\}, \{q_0\}, \{q_0, q_1\},$$

q_0, a	(q_1, q_0)
q_0, b	(q_{\perp}, q_{\perp})
q_1, b	(q_0, q_1)
q_1, a	(q_{\perp}, q_{\perp})
$q_{\perp}, -$	(q_{\perp}, q_{\perp})

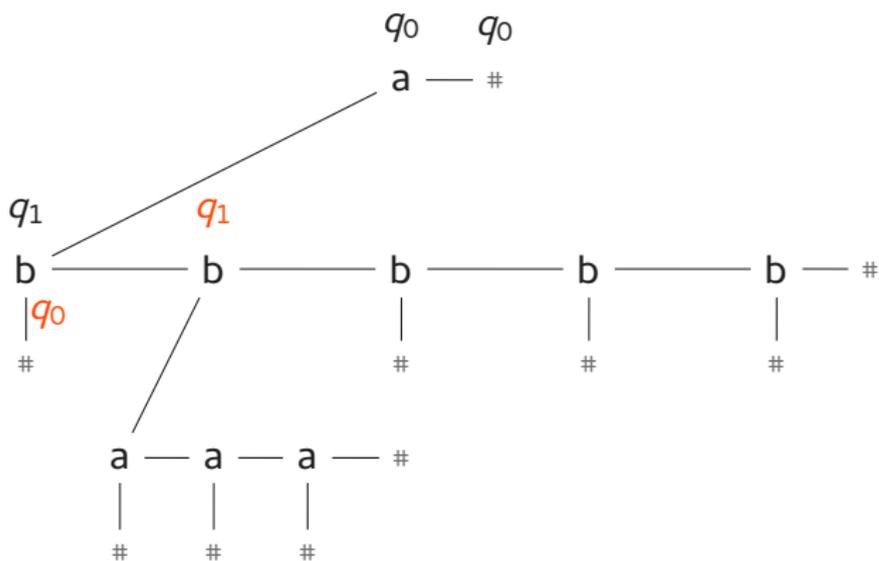
<title>Automate d'arbres </title>



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_\perp\}, \{q_0\}, \{q_0, q_1\},$$

q_0, a	(q_1, q_0)
q_0, b	(q_\perp, q_\perp)
q_1, b	(q_0, q_1)
q_1, a	(q_\perp, q_\perp)
$q_\perp, -$	(q_\perp, q_\perp)

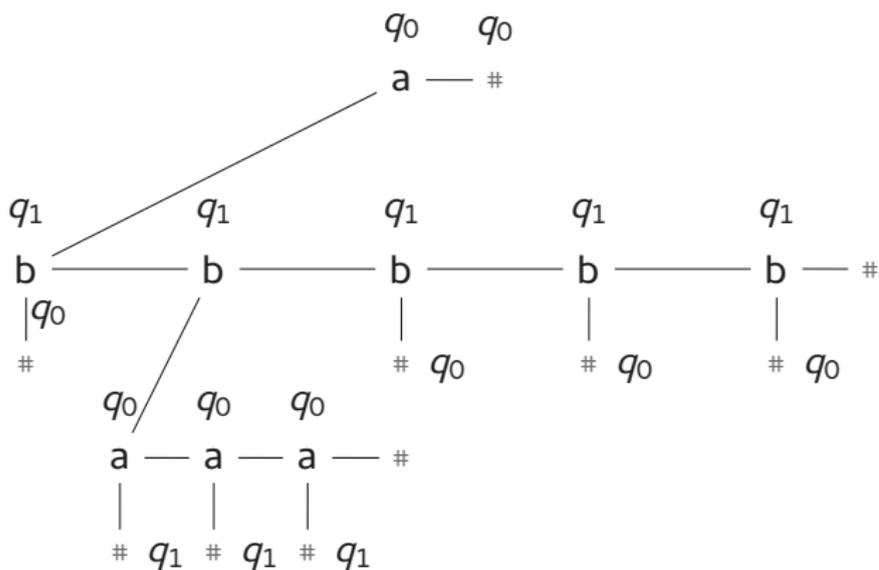
<title>Automate d'arbres </title>



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_\perp\}, \{q_0\}, \{q_0, q_1\},$$

q_0, a	(q_1, q_0)
q_0, b	(q_\perp, q_\perp)
q_1, b	(q_0, q_1)
q_1, a	(q_\perp, q_\perp)
$q_\perp, -$	(q_\perp, q_\perp)

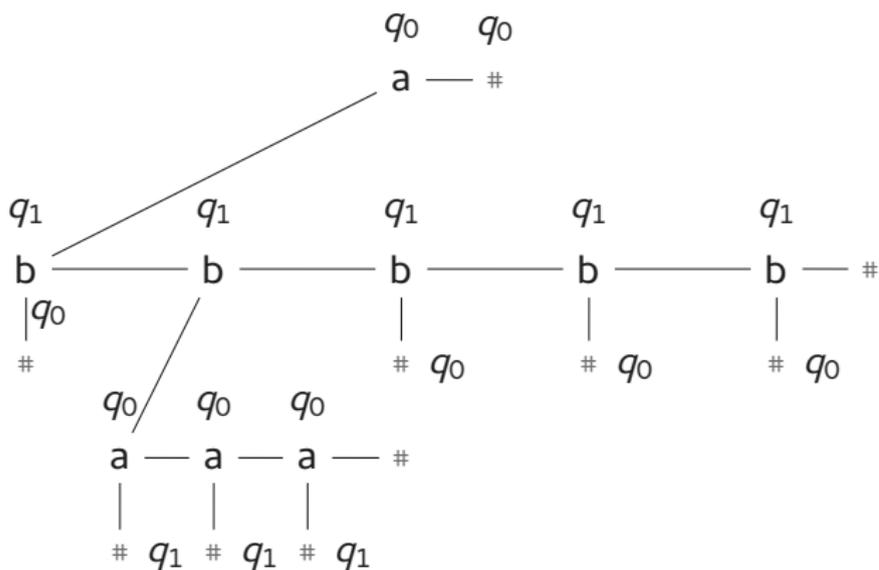
<title>Automate d'arbres </title>



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_{\perp}\}, \{q_0\}, \{q_0, q_1\},$$

q_0, a	(q_1, q_0)
q_0, b	(q_{\perp}, q_{\perp})
q_1, b	(q_0, q_1)
q_1, a	(q_{\perp}, q_{\perp})
$q_{\perp}, -$	(q_{\perp}, q_{\perp})

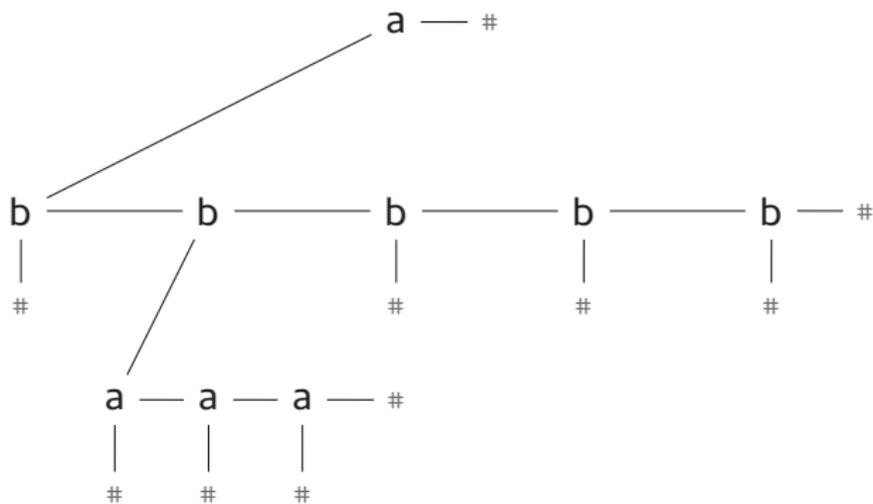
<title>Automate d'arbres descendant</title>



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_{\perp}\}, \{q_0\}, \{q_0, q_1\},$$

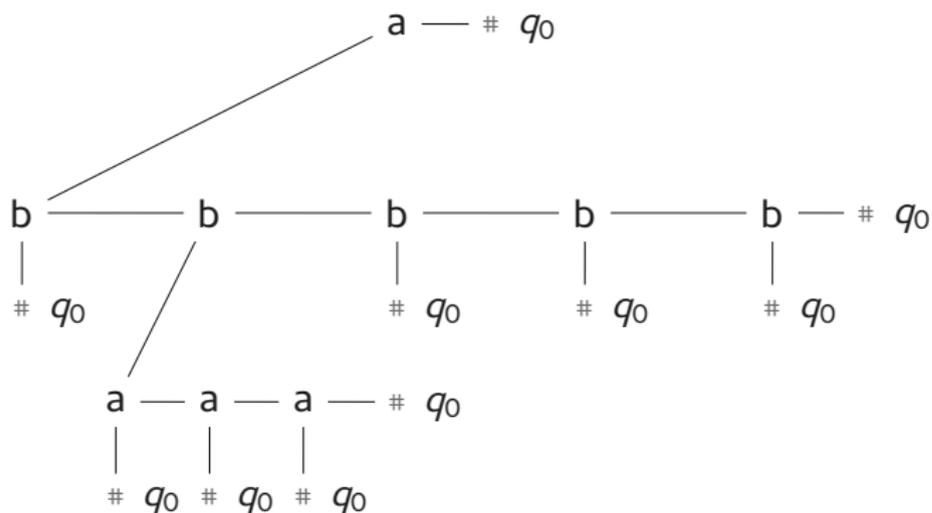
q_0, a	(q_1, q_0)
q_0, b	(q_{\perp}, q_{\perp})
q_1, b	(q_0, q_1)
q_1, a	(q_{\perp}, q_{\perp})
$q_{\perp}, -$	(q_{\perp}, q_{\perp})

<title>Automate d'arbres </title>



$$A_2 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_2, q_\perp\}, \{q_1\}, \{q_0\}, \begin{array}{ll} (q_0|q_2, q_0|q_1, a) & q_1 \\ (q_0|q_1, q_0|q_2, b) & q_2 \\ \dots & q_\perp \end{array}$$

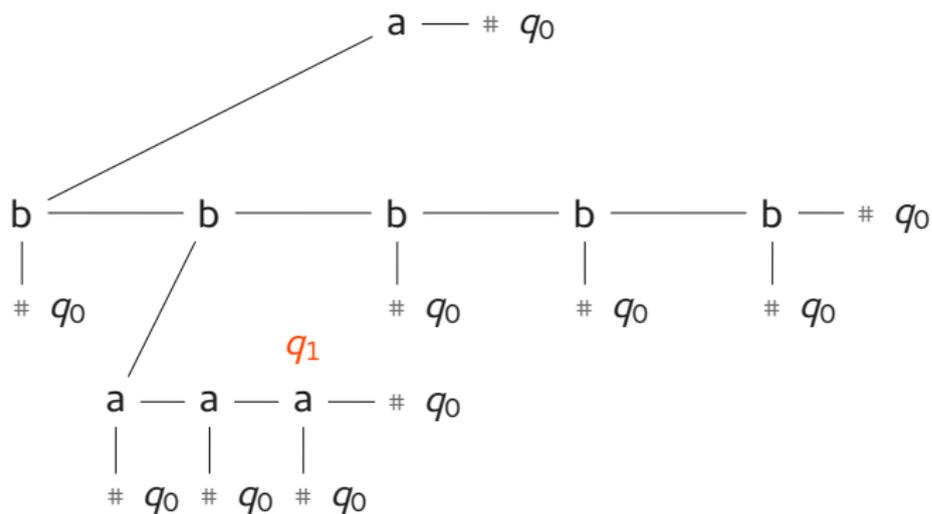
<title>Automate d'arbres </title>



$$A_2 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_2, q_\perp\}, \{q_1\}, \{q_0\}, \dots$$

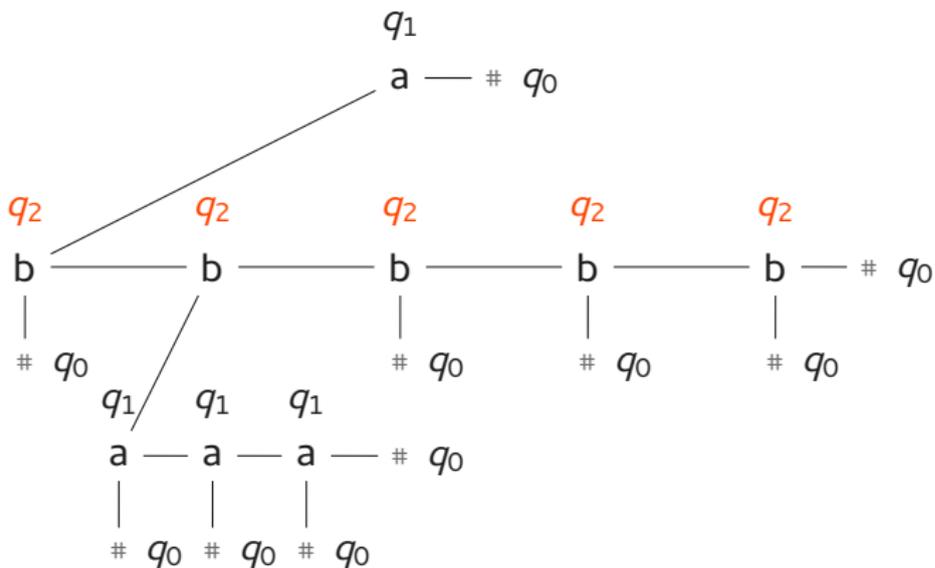
$(q_0 q_2, q_0 q_1, a)$	q_1
$(q_0 q_1, q_0 q_2, b)$	q_2
\dots	q_\perp

<title>Automate d'arbres </title>



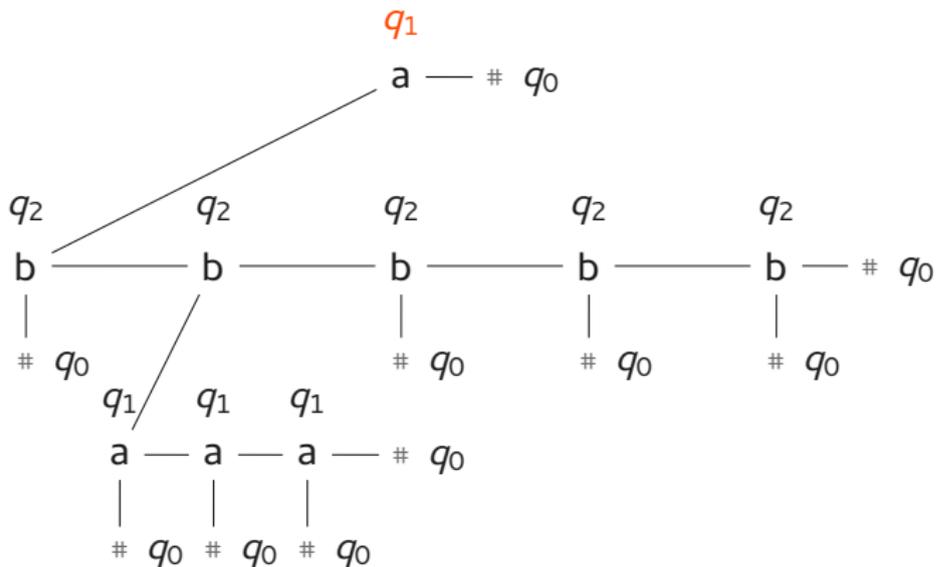
$$A_2 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_2, q_{\perp}\}, \{q_1\}, \{q_0\}, \begin{matrix} (q_0|q_2, q_0|q_1, a) & q_1 \\ (q_0|q_1, q_0|q_2, b) & q_2 \\ \dots & q_{\perp} \end{matrix}$$

<title>Automate d'arbres </title>



$$A_2 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_2, q_{\perp}\}, \{q_1\}, \{q_0\}, \begin{matrix} (q_0|q_2, q_0|q_1, a) & q_1 \\ (q_0|q_1, q_0|q_2, b) & q_2 \\ \dots & q_{\perp} \end{matrix}$$

<title>Automate d'arbres ascendant</title>



$$A_2 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_2, q_{\perp}\}, \{q_1\}, \{q_0\}, \dots$$

$(q_0 q_2, q_0 q_1, a)$	q_1
$(q_0 q_1, q_0 q_2, b)$	q_2
\dots	q_{\perp}

<title>Automates d'arbres</title>

- déterministes (D) / non-déterministes (N)
- top-down (TD) / bottom-up (BU)

<title>Automates d'arbres</title>

- déterministes (D) / non-déterministes (N)
- top-down (TD) / bottom-up (BU)

Quelques propriétés :

- 1 D.TD. \subsetneq D.BU. = N.TD. = N.BU. = Langages Réguliers d'arbres
Exercice : trouver un contre exemple pour D.TD. \subsetneq D.BU.
- 2 N.BU. : opérations d'union, intersection, complémentaire
- 3 N.BU. peuvent être déterminisés (explosion du nombre d'états)
- 4 Appartenance :
 - déterministes $O(|T|)$, en *streaming* pour les D.TD.
 - non-det. $O(|\delta| \times |T|)$ ou $O(2^{|\mathcal{Q}|} + |T|)$
- 5 Inclusion, vide, équivalence :
 - déterministes décidable en temps polynomial
 - non-det. EXPTIME-complet

<title>Digression récursive</title>

```
(* Arbres binaires *)
```

```
type t = Nil | Node of string * t * t
```

```
(* Appartenance pour D. TD. *)
```

```
let rec td_mem t q =  
  match t with  
  Nil -> q ∈ B  
| Node (l, t1, t2) ->  
  let q1, q2 =  $\delta(q,l)$  in  
  if td_mem t1 q1 then td_mem t2 q2  
  else false
```

<title>Digression récursive</title>

```
(* Arbres binaires *)
```

```
type t = Nil | Node of string * t * t
```

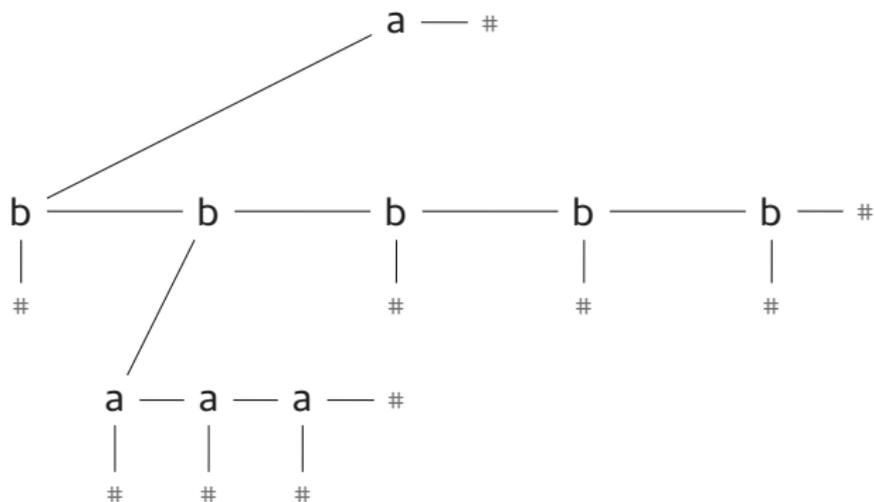
```
(* Appartenance pour D. TD. *)
```

```
let rec td_mem t q =  
  match t with  
  Nil -> ...  
| Node (l, t1, t2) ->  
  ...  
  if td_mem t1 q1 then td_mem t2 q2  
  else false
```

Les fonctions qui ont ce schéma sont *streamables*

<title>Automate avec marquage</title>

On rajoute un ensemble d'états de selection



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_{\perp}\}, \{q_0\}, \{q_0, q_1\}, \{q_1\},$$

$$q_0, a \quad (q_1, q_0)$$

$$q_0, b \quad (q_{\perp}, q_{\perp})$$

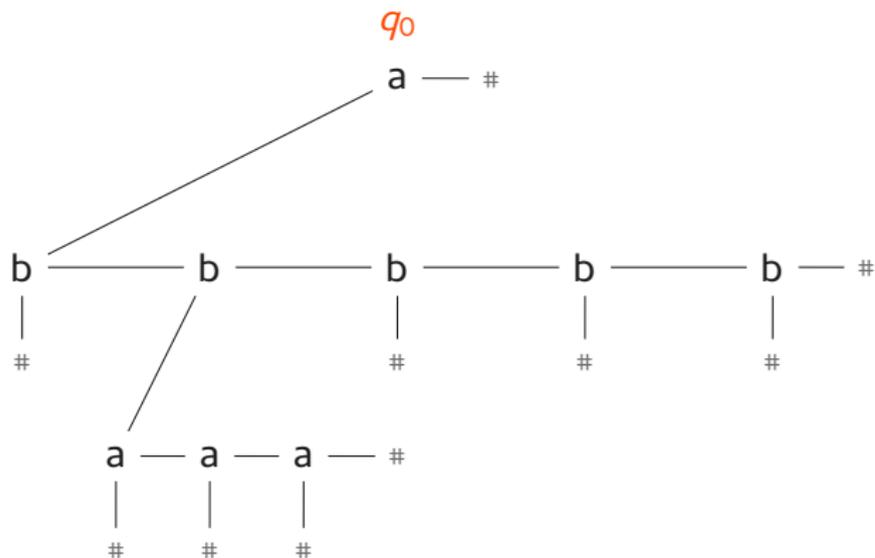
$$q_1, b \quad (q_0, q_1)$$

$$q_1, a \quad (q_{\perp}, q_{\perp})$$

$$q_{\perp}, - \quad (q_{\perp}, q_{\perp})$$

<title>Automate avec marquage</title>

On rajoute un ensemble d'états de selection

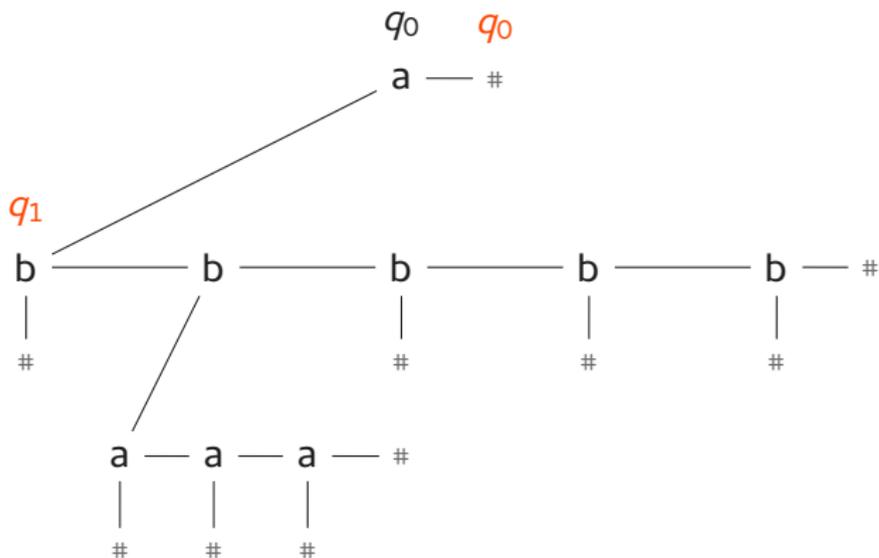


$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_\perp\}, \{q_0\}, \{q_0, q_1\}, \{q_1\},$$

q_0, a	(q_1, q_0)
q_0, b	(q_\perp, q_\perp)
q_1, b	(q_0, q_1)
q_1, a	(q_\perp, q_\perp)
$q_\perp, -$	(q_\perp, q_\perp)

<title>Automate avec marquage</title>

On rajoute un ensemble d'états de selection

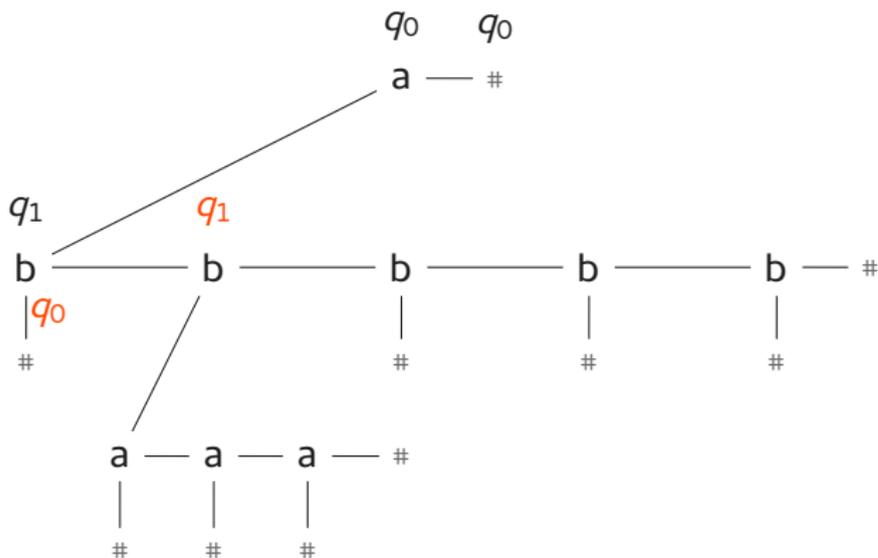


$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_\perp\}, \{q_0\}, \{q_0, q_1\}, \{q_1\},$$

$$\begin{array}{ll} q_0, a & (q_1, q_0) \\ q_0, b & (q_\perp, q_\perp) \\ q_1, b & (q_0, q_1) \\ q_1, a & (q_\perp, q_\perp) \\ q_\perp, - & (q_\perp, q_\perp) \end{array}$$

<title>Automate avec marquage</title>

On rajoute un ensemble d'états de selection



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_\perp\}, \{q_0\}, \{q_0, q_1\}, \{q_1\}, q_0, a \quad (q_1, q_0)$$

$$q_0, b \quad (q_\perp, q_\perp)$$

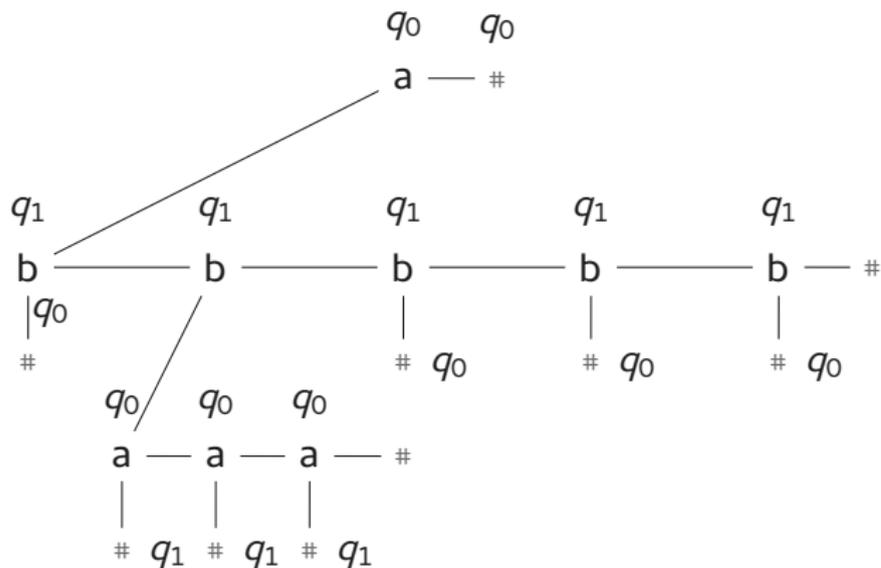
$$q_1, b \quad (q_0, q_1)$$

$$q_1, a \quad (q_\perp, q_\perp)$$

$$q_\perp, - \quad (q_\perp, q_\perp)$$

<title>Automate avec marquage</title>

On rajoute un ensemble d'états de selection



$$A_1 = \{a^2, b^2, \#^0\}, \{q_0, q_1, q_\perp\}, \{q_0\}, \{q_0, q_1\}, \{q_1\},$$

$$\begin{array}{ll} q_0, a & (q_1, q_0) \\ q_0, b & (q_\perp, q_\perp) \\ q_1, b & (q_0, q_1) \\ q_1, a & (q_\perp, q_\perp) \\ q_\perp, - & (q_\perp, q_\perp) \end{array}$$

<title>Automate avec marquage</title>

- Bien définis dans le cas déterministe
- Le cas non déterministe est plus exotique
 - Plusieurs parcours acceptants
 - Nœuds marqués lors d'un parcours acceptant ? lors de tous ?

- Bien définis dans le cas déterministe
- Le cas non déterministe est plus exotique
 - Plusieurs parcours acceptants
 - Nœuds marqués lors d'un parcours acceptant ? lors de tous ?

Digression Unixienne :

- Perl utilise des automates *déterministes* pour les expressions régulières
- Complexité en $O(2^{|R|} + |D|)$ au lieu de $O(|R| \times |D|)$!
- Lié aux *back-references* \1, \2, ...
- Idem pour `lex`, `awk`, `grep`, ...

3. \mathbb{C} Duce et le sous-typage sémantique

[Frisch et al. 2008]

- Langage fonctionnel, d'ordre supérieur
- Adapté à la manipulation de documents XML
- Surcharge/polymorphisme ad-hac
- *Pattern-matching* efficace

<title>C Duce : un exemple</title>

```
include "xhtml.cd"
type t_phone = <phone>[ ('0'--'9')+ ]
type t_email = <email>[ (Char\ '@')+ '@' (Char\ '@')+ ]
type t_name = <name>[ <first>[Char+] <last>[Char+] ]
type t_contact = <contact>[ t_name t_email* t_phone* ]
type t_addressbook = <addressbook>[ t_contact* ]

let to_li (x : t_email|t_phone) : t_li =
  match x with
  | <email>s -> <li>("courriel: " @ s)
  | <_ >(s & [ '0' ('6' | '7') _* ]) -> <li>("portable : " @ s)
  | <_ >s -> <li>("téléphone: " @ s)
```

<title>C Duce : un exemple</title>

```
include "xhtml.cd"
type t_phone = <phone>[ ('0'--'9')+ ]
type t_email = <email>[ (Char\ '@')+ '@' (Char\ '@')+ ]
type t_name = <name>[ <first>[Char+] <last>[Char+] ]
type t_contact = <contact>[ t_name t_email* t_phone* ]
type t_addressbook = <addressbook>[ t_contact* ]

let to_li (x : t_email|t_phone) : t_li =
  match x with
  | <email>s -> <li>("courriel: " @ s)
  | <_ >(s & [ '0' ('6' | '7') _* ]) -> <li>("portable : " @ s)
  | <_ >s -> <li>("téléphone: " @ s)
```

■ Constructeurs

<title>C Duce : un exemple</title>

```
include "xhtml.cd"
type t_phone = <phone>[ ('0'--'9')+ ]
type t_email = <email>[ (Char\ '@')+ '@' (Char\ '@')+ ]
type t_name = <name>[ <first>[Char+] <last>[Char+] ]
type t_contact = <contact>[ t_name t_email* t_phone* ]
type t_addressbook = <addressbook>[ t_contact* ]

let to_li (x : t_email|t_phone) : t_li =
  match x with
  | <email>s -> <li>("courriel: " @ s)
  | <_ >(s & [ '0' ('6' | '7') *_ ]) -> <li>("portable : " @ s)
  | <_ >s -> <li>("téléphone: " @ s)
```

- Constructeurs
- Expressions régulières de types

<title>C Duce : un exemple</title>

```
include "xhtml.cd"
type t_phone = <phone>[ ('0'--'9')+ ]
type t_email = <email>[ (Char\ '@')+ '@' (Char\ '@')+ ]
type t_name = <name>[ <first>[Char+] <last>[Char+] ]
type t_contact = <contact>[ t_name t_email* t_phone* ]
type t_addressbook = <addressbook>[ t_contact* ]

let to_li (x : t_email|t_phone) : t_li =
  match x with
  | <email>s -> <li>("courriel: " @ s)
  | <_ >(s & [ '0' ('6' | '7') _* ]) -> <li>("portable : " @ s)
  | <_ >s -> <li>("téléphone: " @ s)
```

- Constructeurs
- Expressions régulières de types
- Union, intersection, singleton, différence

<title>C Duce : algèbre de type </title>

$t ::=$	$\text{int} \mid \text{string} \mid \dots$	(types de base)
	$\text{'nil} \mid 42 \mid \dots$	(types singletons)
	(t, t)	(produit)
	$t \rightarrow t$	(flèche)
	$t \mid t$	(union)
	$t \& t$	(intersection)
	$\neg t$	(négation)
	empty	(bottom)
	any	(top)
	$\mu T. t$	(types récursifs)
	T	(variable de récursion)

Expressions régulières :

$$(\text{int}, \mu T. \text{'nil} \mid (\text{int}, T)) \equiv [\text{Int}^+]$$

<title>Problème de l'approche syntaxique</title>

$$\frac{\Gamma \vdash e_1 : t \rightarrow s \quad \Gamma \vdash e_2 : u \quad u \leq t}{\Gamma \vdash e_1 e_2 : s}$$

Difficile de définir $u \leq t$ syntaxiquement :

- On veut prendre en compte les propriétés de distributivité, commutativité de $\&$ et $|$
- Beaucoup de cas à gérer
- Difficile d'avoir un algorithme syntaxique complet

Definition (Sous-typage sémantique)

$$s \leq t \Leftrightarrow \llbracket s \rrbracket \subseteq \llbracket t \rrbracket$$

$\llbracket - \rrbracket$: interprétation ensembliste des types

Équivalence **sémantique** des types :

- $(\text{int}, \text{int}) \mid (2, 4) \equiv (\text{int}, \text{int})$
- $[\text{int int}^* \text{ bool? } \text{ bool}^+] \equiv [\text{int}^+ \text{ bool}^+]$
- $(t_1 \mid t_2) \& (t_1 \mid t_3) \equiv t_1 \mid t_2 \& t_3$

Il faut pouvoir décider $\llbracket s \rrbracket \subseteq \llbracket t \rrbracket \Leftrightarrow \llbracket s \rrbracket \cap \overline{\llbracket t \rrbracket} = \emptyset$.

<title>Implantation par automates d'arbres</title>

[int*] $\equiv \mu T. \text{'nil'|(int, T)}$

<title>Implantation par automates d'arbres</title>

$[\text{int}^*] \equiv \mu T. 'nil|(int, T)$

$T \equiv 'nil|(int, T)$

<title>Implantation par automates d'arbres</title>

$$[\text{int}^*] \equiv \mu T. 'nil | (\text{int}, T)$$
$$T \equiv 'nil | (\text{int}, T)$$
$$T \equiv 'nil$$
$$T \equiv (\text{int}, T)$$

<title>Implantation par automates d'arbres</title>

[int*] $\equiv \mu T. 'nil|(int, T)$

$T \equiv 'nil|(int, T)$

$T \equiv 'nil$

$T \equiv (int, T)$

$q_T, 'nil \quad q_{\#}, q_{\#}$

$q_T, \times \quad q_{int}, q_T$

<title>Implantation par automates d'arbres</title>

$$[\text{int}^*] \equiv \mu T. 'nil | (\text{int}, T)$$

$$T \equiv 'nil | (\text{int}, T)$$

$$T \equiv 'nil$$

$$T \equiv (\text{int}, T)$$

$$q_T, 'nil \quad q_{\#}, q_{\#}$$

$$q_T, \times \quad q_{\text{int}}, q_T$$

Types récursifs \equiv Automates d'arbres !

Union, intersection, complément, test du vide,...

<title>Implantation par automates d'arbres</title>

On souhaite décider $\llbracket t \rrbracket \subseteq \emptyset$. Intuition :

- On met t en forme normale disjonctive :

$$t \equiv \left(\bigvee_{\text{any}_{\text{basic}}} \bigwedge t_i \wedge \bigwedge \neg t_j \right) \vee \left(\bigvee_{\text{any}_x} \bigwedge t_i \wedge \bigwedge \neg t_j \right) \vee \dots$$

- On utilise des automates d'arbres pour décider le vide des types récurifs
- Représentation ad-hoc pour chaque type de base (ex : les entiers sont représentés comme des ensemble d'intervalles maximaux)

<title>Implantation des motifs</title>

```
match v with  
<foo>[ x&t_bar *_ y ] -> ...
```

<title>Implantation des motifs</title>

```
match v with  
<foo>[ x&t_bar *_ y ] -> ...
```

- Syntactiquement, les motifs sont des types avec des variables de capture

```
match v with  
<foo>[ x&t_bar *_ y ] -> ...
```

- Syntactiquement, les motifs sont des types avec des variables de capture
- On utilise des automates avec marquage
- Théorème (typage exact du filtrage)
Pour tout motif p , on peut calculer l'ensemble $\{p\}$ des valeurs qui ne font pas échouer le filtrage. De plus, $\{p\}$ est un type.
- On peut déterminer de manière exacte l'exhaustivité ou la redondance lors du filtrage :
 - $t \subseteq \{p_1\} \vee \dots \vee \{p_n\}$: filtrage redondant
 - $\{p_1\} \vee \dots \vee \{p_n\} \subseteq t$: filtrage non exhaustif (et production d'un contre-exemple)

4. XPath, XQuery, encodage d'arbres par intervalles

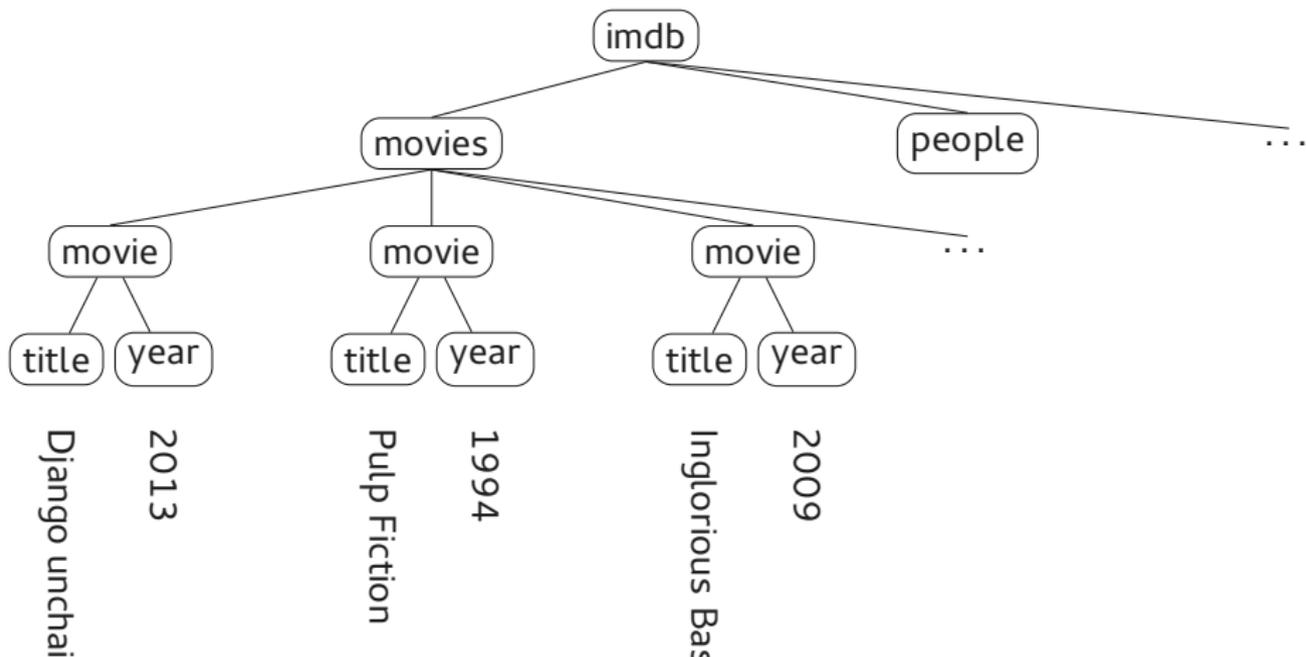
<title>XPath</title>

Langage de selection de nœuds dans un arbre (W3C).

Exemple de requête :

```
/descendant::movie[ child::year >= 2000 ]/child::title
```

```
/descendant::*[ self::year >= 2000 ]/parent::*[child::title
```



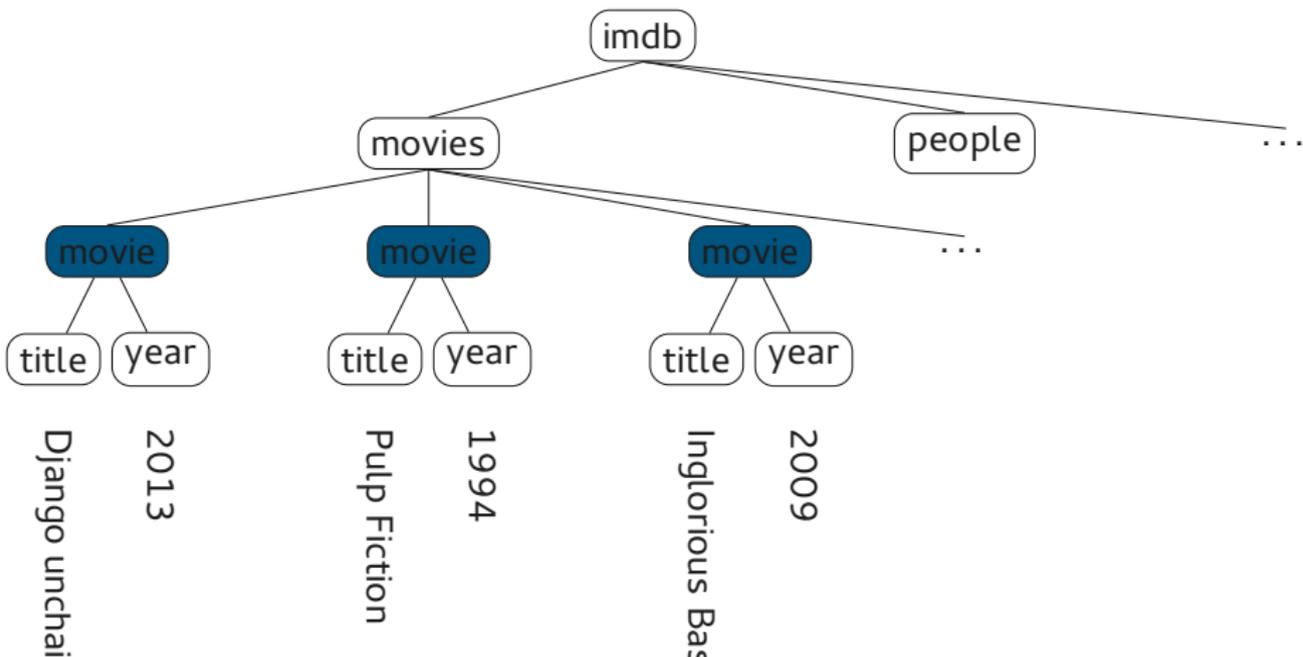
<title>XPath</title>

Langage de selection de nœuds dans un arbre (W3C).

Exemple de requête :

```
/descendant::movie[ child::year >= 2000 ]/child::title
```

```
/descendant::*[ self::year >= 2000 ]/parent::*[child::title
```



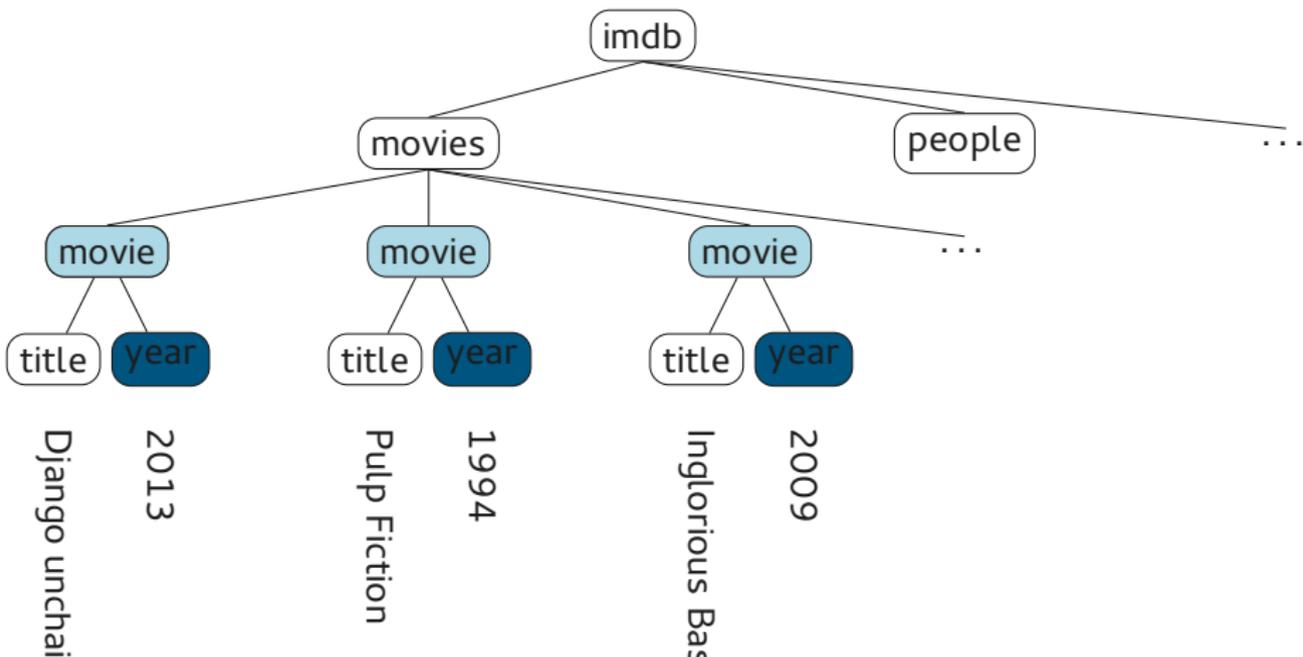
<title>XPath</title>

Langage de selection de nœuds dans un arbre (W3C).

Exemple de requête :

```
/descendant::movie[ child::year >= 2000 ]/child::title
```

```
/descendant::*[ self::year >= 2000 ]/parent::*[child::title
```



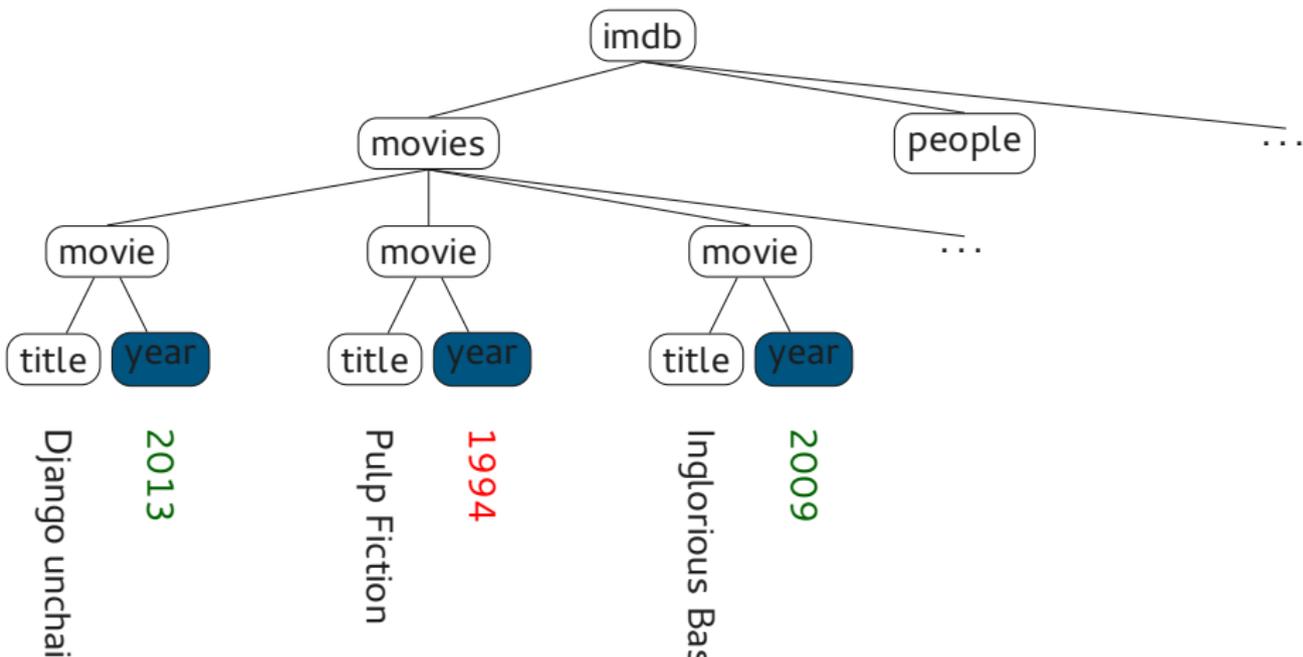
<title>XPath</title>

Langage de selection de nœuds dans un arbre (W3C).

Exemple de requête :

```
/descendant::movie[ child::year >= 2000 ]/child::title
```

```
/descendant::*[ self::year >= 2000 ]/parent::*[child::title
```



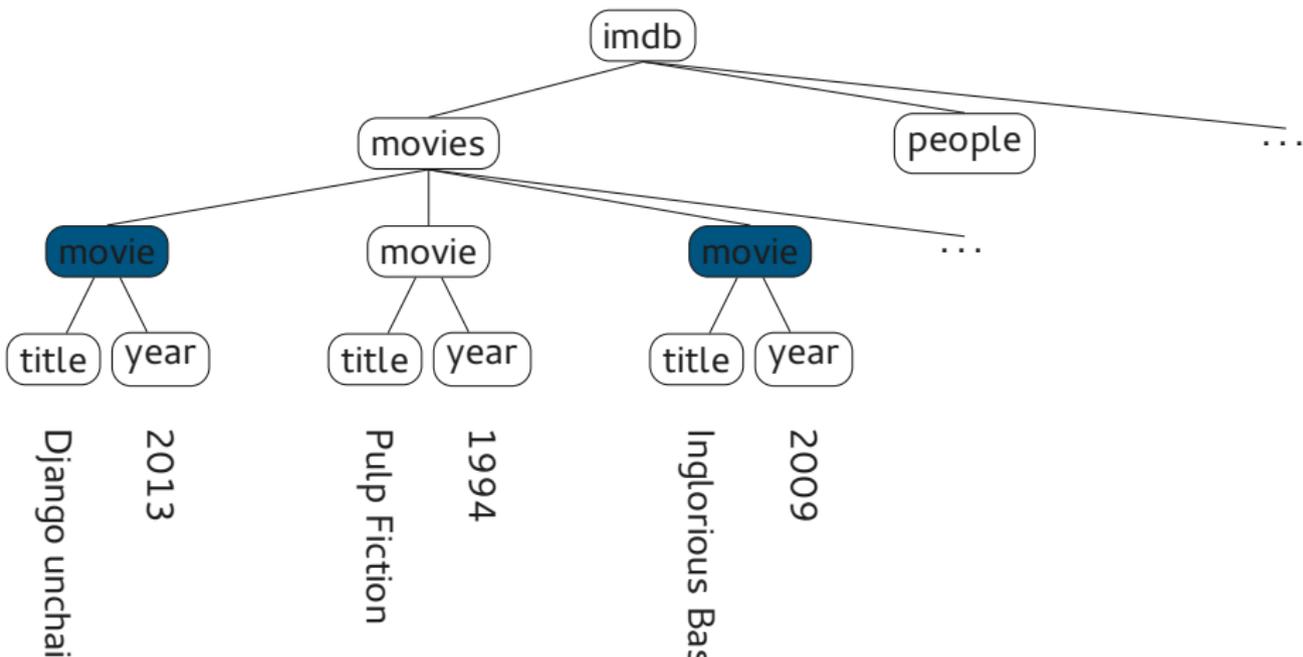
<title>XPath</title>

Langage de selection de nœuds dans un arbre (W3C).

Exemple de requête :

```
/descendant::movie[ child::year >= 2000 ]/child::title
```

```
/descendant::*[ self::year >= 2000 ]/parent::*[child::title
```



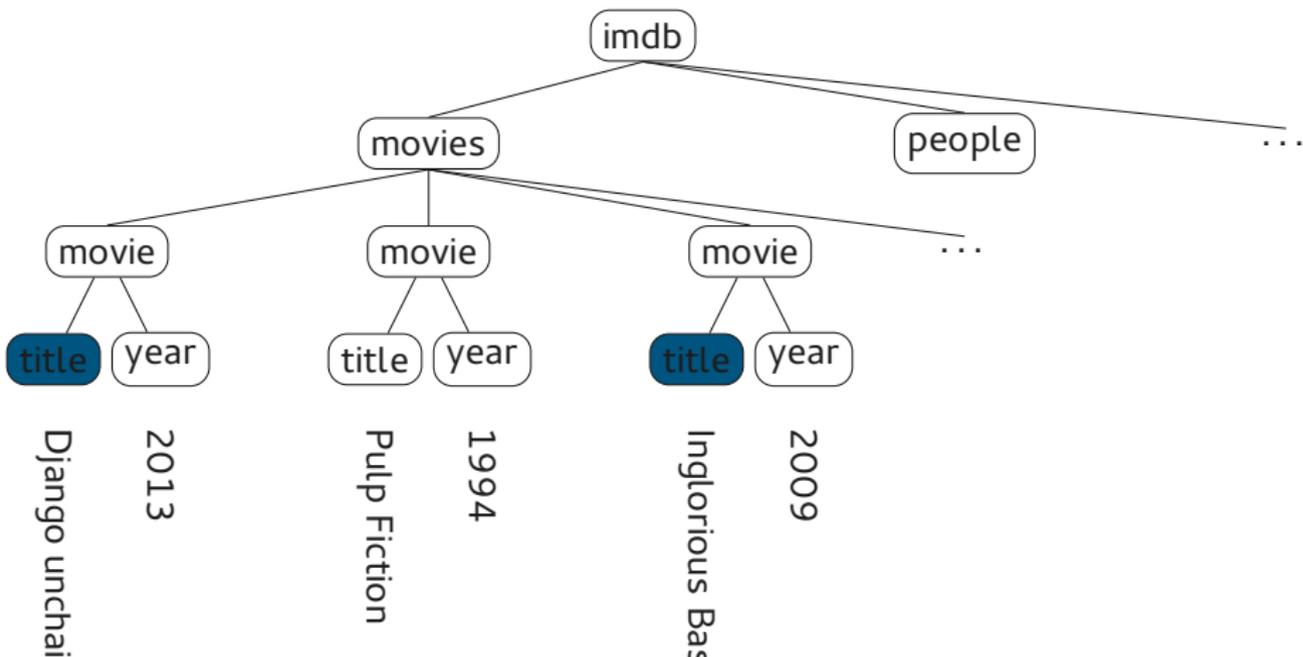
<title>XPath</title>

Langage de selection de nœuds dans un arbre (W3C).

Exemple de requête :

```
/descendant::movie[ child::year >= 2000 ]/child::title
```

```
/descendant::*[ self::year >= 2000 ]/parent::*[child::title
```



Langage de requête et transformation (W3C).

Exemple de programme :

```
let $doc := document("imdb.xml")
return <films>{
  for $i in $doc/descendant::movie[child::year >= 2000]/child::title
  return <titre>{$i}</titre>
}</films>
```

Produit :

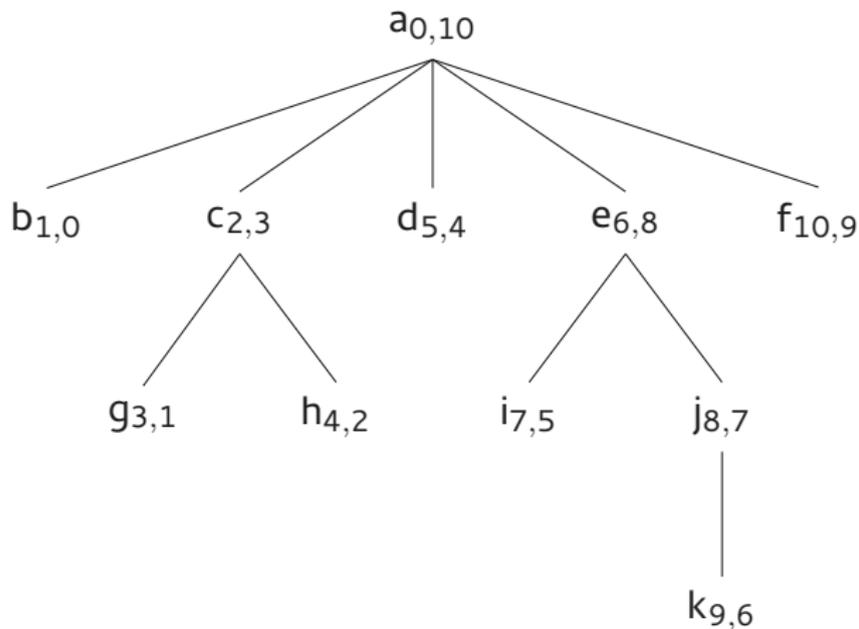
```
<films>
  <titre>Django Unchained</titre>
  <titre>Inglorious Basterds</titre>
</films>
```

<title>Problèmes avec XPath/XQuery</title>

- Le W3C dit quoi faire, pas comment
- Évaluation naïve inefficace
- Pas ou peu d'analyse statique/typage

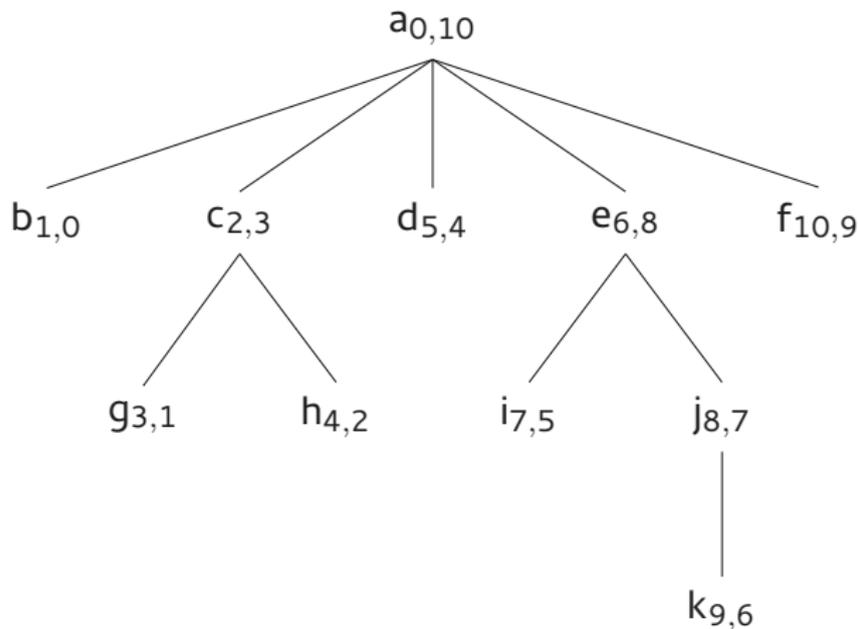
<title>Encodage d'un arbre en intervalles</title>

- 1 Permet de calculer efficacement des chemins XPath
- 2 Encodage compatible avec les BD relationnelles



<title>Encodage d'un arbre en intervalles</title>

- 1 Permet de calculer efficacement des chemins XPath
- 2 Encodage compatible avec les BD relationnelles



$\text{descendant}(x, y) =$
 $x.\text{pre} < y.\text{pre} \wedge$
 $x.\text{post} > y.\text{post}$

$\text{ancestor}(x, y) =$
 $\text{descendant}(y, x)$

$\text{child}(x, y) =$
 $\text{descendant}(x, y) \wedge$
 $\text{level}(x) = \text{level}(y) - 1$

<title>Encodage d'un arbre en intervalles</title>

1 Évaluation d'une requête XPath en $O(|Q| \times |T|)$

2 Traduction immédiate d'XPath en SQL :

```
SELECT DISTINCT z FROM root doc, x doc, y doc, z doc WHERE
  root.pre = 0 AND
  x.pre > root.pre AND x.post < root.post
  AND x.label = "movie" AND
  y.pre > x.pre AND y.post < x.post AND y.level - 1 = x.level
  AND y.label = "year" AND y.value >= 2000
  z.pre > x.pre AND z.post < x.post AND z.level - 1 = x.level
  AND z.label = "title"
ORDER BY z.pre
```

3 Base des évaluateurs XPath/XQuery modernes

4 Variantes avec des indices non contigus pour permettre l'insertion/suppression de nœuds

Étant données deux requêtes XPath R_1 et R_2 on veut savoir si :

$$R_1 \subseteq R_2$$

pour tout document d'entrée

- détecter du code mort (prendre une R_2 vide)
- maintenance de vue/cache
- ...

On peut utiliser l'encodage pre/post pour l'analyse statique :

- Transforme la requête XPath en formule logique du 1^{er} ordre + arithmétique entière
- On donne le tout à un solveur SMT [*Benedikt, Cheney, 2010*]
- On attend (et on espère)

5. Logiques pour les arbres

<title>Logiques d'arbres</title>

Sur quel langage de requêtes sur les arbres peut-on raisonner :

XPath trop verbeux, pas de bonnes propriétés de clotures

Automates d'arbres trop « bas-niveau » (récursion explicite)

Formules logiques oui mais lesquelles ?

Remarque : on veut conserver le pouvoir expressif des langages réguliers (clotures par opérations booléennes, implantations efficaces, langage utilisé pour les types, ...)

<title>Logique du 1^{er} ordre sur les arbres</title>

ϕ	::=	$U(x)$	(relation unaire)
		$B(x, y)$	(relation binaire)
		$\phi \vee \phi$	(disjonction)
		$\exists x.\phi$	(quant. existentielle)
		\top	(vrai)

avec le sucre syntaxique usuel : $\phi \wedge \phi, \forall x.\phi, \phi \Rightarrow \phi, \dots$

Exemple : « tous les nœuds a ont un fils b »

$$\forall x.(\text{label}_a(x) \Rightarrow \exists y.\text{child}(x, y) \wedge \text{label}_b(y))$$

<title>Logique du 1^{er} ordre sur les arbres</title>

ϕ	::=	$U(x)$	(relation unaire)
		$B(x, y)$	(relation binaire)
		$\phi \vee \phi$	(disjonction)
		$\exists x.\phi$	(quant. existentielle)
		\top	(vrai)

avec le sucre syntaxique usuel : $\phi \wedge \phi, \forall x.\phi, \phi \Rightarrow \phi, \dots$

Exemple : « tous les nœuds a ont un fils b »

$$\forall x.(\text{label}_a(x) \Rightarrow \exists y.\text{child}(x, y) \wedge \text{label}_b(y))$$

Limites du 1^{er} ordre ?

« L'ensemble des arbres qui ont un nombre pair de a » n'est pas FO.

<title>Logique du 1^{er} ordre sur les arbres</title>

ϕ	::=	$U(x)$	(relation unaire)
		$B(x, y)$	(relation binaire)
		$\phi \vee \phi$	(disjonction)
		$\exists x.\phi$	(quant. existentielle)
		\top	(vrai)

avec le sucre syntaxique usuel : $\phi \wedge \phi, \forall x.\phi, \phi \Rightarrow \phi, \dots$

Exemple : « tous les nœuds a ont un fils b »

$$\forall x.(\text{label}_a(x) \Rightarrow \exists y.\text{child}(x, y) \wedge \text{label}_b(y))$$

Limites du 1^{er} ordre ?

« L'ensemble des arbres qui ont un nombre pair de a » n'est pas FO.

(Exercice : écrire l'automate d'arbre correspondant)

<title>Logique Monadique du Second ordre</title>

$$\begin{array}{l} \phi ::= \dots \quad (1^{\text{er}} \text{ ordre}) \\ | x \in X \quad (\text{relation unaire variable}) \\ | \forall X.\phi \quad (\text{quant. univ. second-ordre}) \\ | \exists X.\phi \quad (\text{quant. exist. second-ordre}) \end{array}$$

Théorème :

MSO-définissable \equiv langage régulier d'arbre

<title>Logique Monadique du Second ordre</title>

$$\begin{array}{l} \phi ::= \dots \quad (1^{\text{er}} \text{ ordre}) \\ | \quad x \in X \quad (\text{relation unaire variable}) \\ | \quad \forall X.\phi \quad (\text{quant. univ. second-ordre}) \\ | \quad \exists X.\phi \quad (\text{quant. exist. second-ordre}) \end{array}$$

Théorème :

MSO-définissable \equiv langage régulier d'arbre

Problème :

Pour tout entier n , il existe ϕ MSO-définissable telle que

$$|A_\phi| = 2^{2^{\dots^{|\phi|}}} \} n \text{ fois}$$

<title> μ -calcul sur les arbres</title>

$\phi ::=$	\top	(vrai)
	$l \mid \neg l$	(nom de nœud)
	$\phi \wedge \phi$	(conjonction)
	$\phi \vee \phi$	(disjonction)
	$\downarrow_i \phi \mid \uparrow_i \phi \mid \neg \downarrow_i \top \mid \neg \uparrow_i \top$	$i \in \{1, 2\}$ (modalité)
	$\mu X. \phi$	(plus petit point fixe)

Exemple :

$$\uparrow_1 a \wedge b \wedge \mu X. (c \vee \downarrow_2 X)$$

<title> μ -calcul sur les arbres</title>

$\phi ::=$	\top	(vrai)
	$l \mid \neg l$	(nom de nœud)
	$\phi \wedge \phi$	(conjonction)
	$\phi \vee \phi$	(disjonction)
	$\downarrow_i \phi \mid \uparrow_i \phi \mid \neg \downarrow_i \top \mid \neg \uparrow_i \top$	$i \in \{1, 2\}$ (modalité)
	$\mu X. \phi$	(plus petit point fixe)

Exemple :

$$\uparrow_1 a \wedge b \wedge \mu X. (c \vee \downarrow_2 X)$$

- Équivalent aux langages réguliers
- $\text{SAT}(\phi)$ décidable en $O(2^{|\phi|})$ (seulement !)
- Implantation efficace :

<http://wam.inrialpes.fr/web-solver/webinterface.html>

[Genevès et al. 2005, 2011]

Conclusion

Au final, un bon prétexte pour s'intéresser à plein de choses :

- Langages réguliers, automates
- Logiques (de toutes sortes), expressivité, complexité
- Représentations efficaces (encodage d'arbres, BDD, ...)
- Langages de programmation (compilation efficace de requête XML)
- ...

Ça permet aussi de voir que toutes ces choses sont reliées

<title>1999, l'année d'XML*?</title>

* : source, PC Impact, décembre 1998.

XML n'est plus vraiment une techno. à la mode (sauf si on en fait dans le *cloud*). Il reste quand même plein de choses à faire :

- Structure de données efficaces avec mise à jour pour les arbres
- Évaluateurs efficaces basés sur des automates
- Populariser l'analyse statique dans la communauté BD
- Sous-typage sémantique + polymorphisme à la ML
- μ -calcul ou MSO + procédures de décisions
- ...