

# Langages à balises : une introduction

Yannick Prié  
UFR Informatique – Université Lyon 1

UE2.2 – Master SIB M1 – 2004-2005

## Objectif généraux du cours

- Comprendre les grands principes de la représentation de données et de documents numériques à l'aide d'un langage à balises.
- Découvrir XML, son histoire et son fonctionnement
- Définir des langages basés sur XML à l'aide de DTD
- S'initier à la transformation de documents en utilisant XSL et un moteur XSLT
- Apprendre les bases de XHTML pour la génération de pages web

CM1 : Introduction aux langages à balises – Yannick Prié  
UE2.2 – Master SIB M1 – 2004-2005 : Représentation des données et des connaissances

2

## Objectifs de ce cours introductif

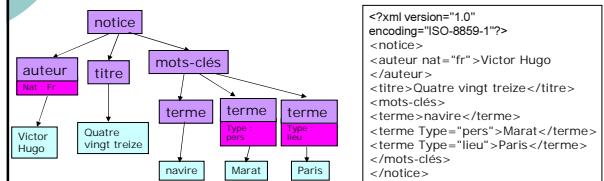
- Introduction aux langages à balise et à leurs principes
  - arbres
  - grammaires
  - langages à balises
- Histoire de ces langages
- Présentation de la galaxie XML et de la suite du cours

CM1 : Introduction aux langages à balises – Yannick Prié  
UE2.2 – Master SIB M1 – 2004-2005 : Représentation des données et des connaissances

3

## Idee générale

- Représenter de l'information dans des structures arborescentes
- Coder ces structures dans des fichiers, qui pourront être échangés

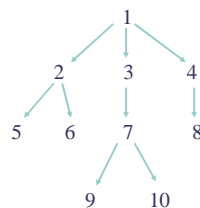


CM1 : Introduction aux langages à balises – Yannick Prié  
UE2.2 – Master SIB M1 – 2004-2005 : Représentation des données et des connaissances

4

## Parler des arbres

- Arbre
- Noeud
  - noeuds fils et pères
- Racine
- Feuille
- Chemin
  - suite de noeud dont chaque élément a pour père le précédent),
- Branche
  - chemin se terminant sur une feuille
- Ancêtres et descendants
- Taille d'un arbre
  - nombre de noeuds
- Profondeur d'un noeud



CM1 : Introduction aux langages à balises – Yannick Prié  
UE2.2 – Master SIB M1 – 2004-2005 : Représentation des données et des connaissances

5

## Les arbres sont partout !



CM1 : Introduction aux langages à balises – Yannick Prié  
UE2.2 – Master SIB M1 – 2004-2005 : Représentation des données et des connaissances

6

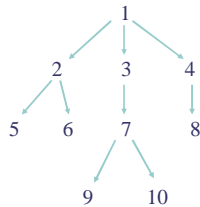
## Parcours d'arbre

### o Largeur d'abord

1  
 → 2 → 3 → 4  
 → 5 → 6 → 7 → 8  
 → 9 → 10

### o Profondeur d'abord

1 → 2 → 5 → 6  
 → 3 → 7 → 9 → 10  
 → 4 → 8



## Exemple de parcours

### o Objectif :

- compter les nœuds entourés

### o Comptage local :

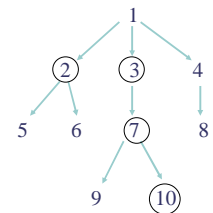
- compter\_localement (nœud)
  - Si il y a un cercle
  - Alors  $n \leftarrow n + 1$
  - Pour tous les nœuds fils  $n_i$ , Compter\_localement ( $n_i$ )

### o Appel général :

- $n \leftarrow 0$
- compter\_localement (nœud 1)
- afficher n

### o Remarques

- Parcours en profondeur d'abord
- Autant de comptages locaux que de nœuds
- Marche sur n'importe quel arbre : on part de la racine et on parcourt tout
- Pas de vision globale de l'arbre



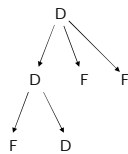
## Notion de grammaire

### o Système formel

- vocabulaire + règles de production
- permet de définir un arbre

### o Exemple

- vocabulaire
  - D (Dossier)
  - F (fichier)
- règle
  - $D \rightarrow (D|F)^*$
  - Avec
    - \* == zéro ou plus
    - | == ou



## Autre exemple

### o Vocabulaire

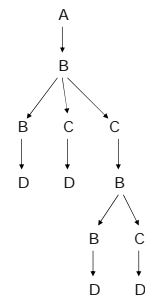
- A, B, C, D

### o Règles

- $A \rightarrow B^+$
- Avec
  - + == 1 ou plus
- $B \rightarrow BC^* | D$
- $C \rightarrow (D | B)$

### o Question

- Quel est l'arbre le plus petit que l'on peut écrire avec cette grammaire ?



## Arbres et séquences d'octets

### o Fichier

- suite d'octets

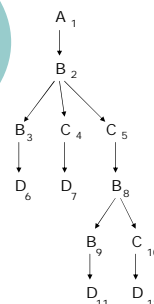
### o Objectif

- représenter un arbre dans un fichier

### o Solution

- décrire l'arbre comme un ensemble d'éléments qui se contiennent les uns les autres.
- représenter les éléments entre deux balises
  - balises ouvrantes
    - on les notera par exemple `<nom>`
  - balises fermantes
    - on les notera par exemple `</nom>`

## Arbres et séquences



### Eléments

A1 ⊂ B2  
 B2 ⊂ B3 C4 C5  
 B3 ⊂ D6  
 C4 ⊂ D7  
 C5 ⊂ B8 C10  
 B8 ⊂ B9  
 B9 ⊂ D11  
 C10 ⊂ D12

### Eléments et balises

```
<A>
<B>
<B>
<D></D>
</B>
<C>
<D></D>
</C>
<C>
<B>
<B><D></D></B>
<C><D></D></C>
</B>
</C>
</A>
```

## Langages à balises

- o Tous les langages ayant pour objectif de représenter de l'information en utilisant des balises
- o Définis par
  - vocabulaire
    - o noms des éléments
  - grammaire
    - o mode d'organisation des éléments
      - des éléments en contiennent d'autres
  - + attributs des éléments
    - o un peu plus de structure (voir cours XML)
- o Une description
  - ensemble d'éléments organisés dans un fichier
  - contenus terminaux (texte)

## Familles de langages à balises

- o Décrire une notice bibliographique
  - notice, titre, auteur, mots-clés, terme, résumé, ...
- o Décrire un poème :
  - poème, quatrain, tercet, vers, ...

```
<poeme type="sonnet">
<quatrain>
<vers>Je vis, je meurs : je me brûle
et me noie.</vers>
<vers>J'ai chaud extrême en
endurant froidure : </vers>
<vers> ... </vers>
<vers> ... </vers>
</quatrain>
...
</poeme>
```

```
<notice>
<auteur nat="fr">Victor Hugo
</auteur>
<titre>Quatre vingt treize</titre>
<mots-clés>
<terme>navire</terme>
<terme Type="pers">Marat</terme>
<terme Type="lieu">Paris</terme>
</mots-clés>
</notice>
```

- vocabulaires différents
- grammaires différentes
- mais *même manière d'exprimer les descriptions*

## Notion de métalangage

- o Langage avec lequel on peut définir d'autres langages
- o Pour les langages à balises
  - langage exprimant la manière dont on peut décrire une famille de langages à balise
    - o comment exprimer les éléments ?
    - o comment organiser les éléments ?
- o Exemples de métalangages
  - SGML
    - o permet de définir : TEI, HTML, ...
  - XML
    - o permet de définir : SVG, TEI, XHTML, ...

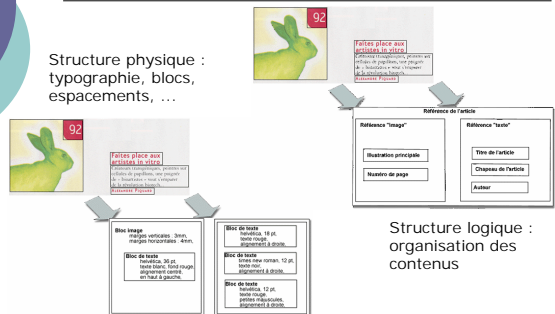
## Dans la suite

- o Petite histoire des langages à balises et des applications liées
  - de SGML à XML

## Représentation de documents

- o Document numérique
  - manipulations et gestion par des ordinateurs
- o Document structuré
  - séparation de la structure physique et de la structure logique
    - o séparation forme / contenu
- o D'où possibilité
  - de manipuler la structure logique des documents
  - d'accéder au texte des différentes parties logiques des documents
  - de générer plusieurs structures physiques à partir d'une structure logique

## Structures logique / physique



## Balisage de texte

- o Idée
  - marquer des zones des textes pour les qualifier
    - o les balises ouvrantes et fermantes délimitent les éléments de description
    - o la structure logique est un arbre « ajouté » au texte

```
<p>Il est de tradition de présenter un langage de programmation à l'aide d'un premier exemple comme : <eg> CHAR*20 GRTG GRTG = 'BONJOUR TOUT LE MONDE' PRINT *, GRTG END </eg></p>
<p>Dans cet exemple, on commence par déclarer la variable <ident>GRTG</ident>, dans la ligne <kw>CHAR*20 GRTG</kw>, qui identifie <ident>GRTG</ident> comme formée de 20 octets de type <kw>CHAR</kw>. On affecte alors à cette variable la valeur <mentioned>BONJOUR TOUT LE MONDE</mentioned>. Suivent alors l'ordre d'impression <kw>PRINT</kw> et l'instruction finale <kw>END</kw>.</p>
```

p : servira à la mise en page  
 eg, kw, mentioned : seront mis en évidence dans la structure physique  
 kw, mentioned : utilisés pour construire un index etc.

## SGML

- o Objectif : représenter l'information contenue dans un document indépendamment
  - des systèmes utilisés pour la saisie et le traitement
  - de la forme physique qu'il sera amené à prendre (papier, CD-ROM, web...)
  - des langues et des alphabets, latins ou non
  - des applications
- o Naissance chez IBM (années soixante)
  - GML
  - Gestion de la documentation technique
- o Normalisation 1986 ISO-8879
  - une dizaine d'années de travail
- o Utilisation
  - Description des documents dans les grosses organisations
    - o Complexité des langages
    - o Lourdeur et cherté des outils (chaîne de traitement)
    - o Journal Officiel, grosses entreprises/documentations, éditeurs...
  - Echange des documents

## SGML : principes

- o Métalangage
  - permet de décrire des modèles (grammaires)
- o Notion de DTD
  - Document Type Definition
  - Permet de décrire un modèle
    - o un type de document
- o Un document SGML
  - Est une instance du type de document
  - Doit être conforme à la DTD associée

## SGML : exemple

(d'après <http://www.cavi.univ-paris3.fr/ilpga/ilpga/tal/>)

### Instance

```
<!DOCTYPE memo SYSTEM "memo.dtd">
<memo statut="conf">
<auteur>Serge Fleury</auteur>
<dest>
<nom>André Salem</nom>
<nom>Pollet Samvelian</nom>
</dest>
< sujet>Cours SLFE6</ sujet>
< corps>
< par>Veuillez noter que le cours SLFE6 sur les documents électronique aura bel et bien lieu au mois de mai 2002</ par>
< par>S'il y avait des changements de votre côté, veuillez m'en aviser dans les plus brefs délais.</ par>
</ corps>
</ memo>
```

### DTD (memo.dtd)

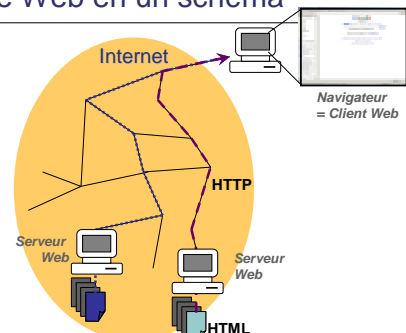
```
<!- DTD utilisable pour baliser les memos en SGML ->
<ELEMENT memo -- ((auteur & (date?) & sujet & dest & (cc?)), corps) -->
<!ATTLIST memo statut (conf | pub) #REQUIRED
<ELEMENT (dest | cc) -- (nom+) -->
<ELEMENT corps -- (par*) -->
<ELEMENT (auteur | date | sujet | nom | par) -- (#PCDATA) -->
```

## HTML

- o 1991 – CERN – Tim Berners Lee
- o Basé sur
  - Principes de l'hypertexte
  - Client/serveur sur IP
- o Principes
  - Des serveurs peuvent fournir des documents hypertextes
  - Les documents seront décrits en suivant une DTD SGML → HTML (HyperText Markup Language)
  - Les liens sont décrits avec leur cible (URL)
  - Un client (navigateur)
    - o permet de lire (présenter) les documents HTML
    - o charge un nouveau document quand on active un lien
  - Protocole d'échange : HTTP (HyperText Transfert Protocol)

## Le Web en un schéma

(d'après cours PCI – Web <http://pci.univ-lyon1.fr>)



## HTML : notion d'URL

- o Uniform Ressource Locator
  - + permet d'identifier une ressource sur le réseau
- o Une ressource peut être
  - une page Web
  - une image (seule ou utilisée dans une page Web)
  - un programme
  - un fichier à télécharger...
- o Une URL indique
  - un protocole (langage de communication entre deux programmes sur deux machines)
    - o FTP (File Transfert Protocol),
    - o HTTP (HyperText Transfert Protocol)...
  - l'adresse d'un serveur
  - un chemin dans l'arborescence des fichiers
- o Forme générale : **protocole://adresse/chemin**
- o Exemples
  - <http://www.univ-lyon1.fr/>
  - <http://www10.univ-lyon1.fr/~yprie/Enseignement/SIB/SIB-UE3-bloc4/CM4.6-7.pdf>



## Une première remarque : URLs et URIs

- o Une URL indique
  - une ressource
  - sur une machine
  - accessible par un protocole
- o Généralisation
  - URI (Uniform Ressource Identifier)
  - Identifier une ressource
    - o disponible sur internet : URL
    - o simplement en lui donnant un nom (URName)
      - urn:ietf:rfc:2396
      - <http://yannick.prie.org/mescollegues/Lionel.Medini>

## Une deuxième remarque : sur la normalisation

- o Norme industrielle
  - Référentiel publié par un organisme officiel (ISO, AFNOR...)
  - En anglais : *standard*
- o Standard
  - Référentiel publié par une entité privée
  - Si diffusion large : *standard de fait*
- o Consortium
  - Ensemble d'entreprises, de centres de recherche, de particuliers qui s'allient pour définir des normes et standards sur tout et n'importe quoi
  - Gain : fournir les outils au moment où le référentiel est publié
    - o JPEG (Joint Picture Expert Group) → norme ISO
    - o MPEG (Moving Picture Expert Group) → norme ISO
    - o W3C (World Wide Web Consortium) → standards
    - o ...

## Pourquoi XML ?

- o Objectif
  - représenter et échanger des données et des documents sur le web
- o SGML
  - un peu vieux
  - trop complexe
- o HTML
  - trop basique
    - o document = en-tête + corps
  - mélange logique / présentation
    - o balise b = bold (mise en gras) :  
<bold>Attention !</b>
    - o bonne approche
      - <important>Attention !</important>
      - présenter la chaîne de caractères importante avec une mise en forme particulière (italique, rouge, gras, etc.)

## Objectifs XML

- o XML doit être facilement utilisable sur le Web
- o XML doit supporter une grande variété d'applications
- o XML doit être compatible avec SGML
- o Il doit être facile d'écrire des programmes qui traitent des documents XML
- o Le nombre d'options doit être réduit au minimum, idéalement à zéro
- o Les documents XML doivent être lisibles et raisonnablement clairs
- o La conception de XML doit être menée rapidement
- o La description de XML doit être formelle et concise
- o Les documents XML doivent être faciles à créer
- o La concision du balisage XML est d'une importance minime

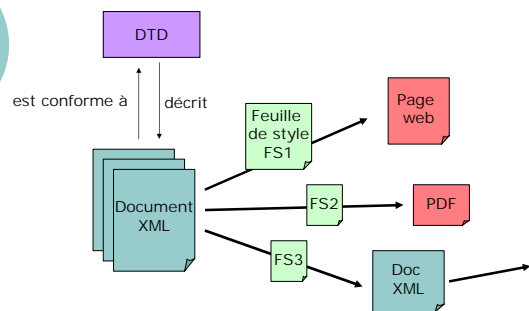
## XML = métalangage

- o Permet de décrire des types de documents
  - DTD, Schémas
- o Permet de définir des instances
  - Documents XML
    - o Répondant à un type de document
      - Classique
    - o Simplement bien construits
      - Nouveau
- o Les instances peuvent décrire
  - des documents (texte balisé)
    - o Classique, cf. SGML
  - des données structurées quelconques
    - o Nouveau !

## Principe général XML

- o DTDs, Schéma
  - Comment décrire les données et les documents ?
- o Documents XML
  - Les données et les documents eux-mêmes, dans des fichiers
- o Feuilles de style
  - Manière de présenter les données et les documents
- o Remarque
  - On ne sait plus trop bien où sont les données, et les documents !

## Schéma récapitulatif

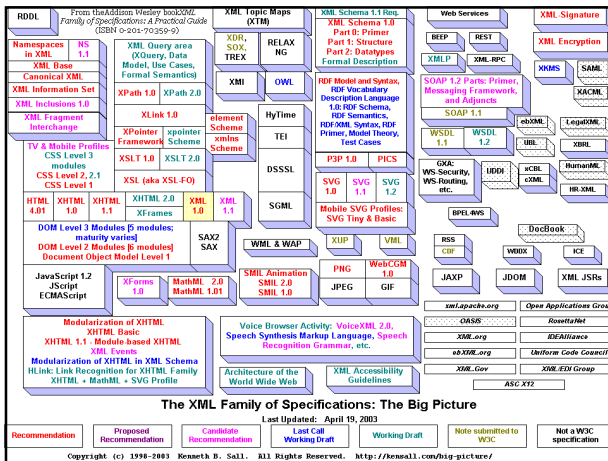


## Troisième remarque : Intégration de XML dans les SI

- o Stockage de données
  - Simples fichiers (ex. configuration)
  - Bases de données semi-structurées (requêtes, etc.)
  - Bases de données documentaires
    - o Documents XML
    - o Documents XHTML (web)
- o Echange de données
  - D'une base de données vers une autre (format d'échange)
  - Serveur vers un navigateur : données + feuille de style
- o Remarque :
  - Circulation de flux XML sur un réseau :
    - o Utilisation de l'arbre entier (le document)
    - o Utilisation à la volée pour les très gros documents (exemple : BIM)

## Différents langages plus ou moins standards liés à XML

- o DTD / Schémas pour décrire
  - Données
  - Documents
- o Normalisation à différents niveaux
  - W3C
  - ISO
  - Organismes liés à un domaine
  - ...



## Et d'autres encore !

### o Suite du module :

- XML
  - o 3 CM / 2 TP
- XPATH, XSL
  - o 2 CM / 1 TP
- (X)HTML / CSS
  - o 2 CM / 2 TP
- Projet