

# XML

Yannick Prié  
UFR Informatique – Université Lyon 1  
UE2.2 – Master SIB M1 – 2004-2005

# Objectifs des trois cours

- Etre capable de comprendre des documents XML et des DTD
- Etre capable de construire des documents XML et des DTD
- Découverte de quelques DTD « importantes »

# Un document XML

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE livre SYSTEM "E:\prie\Enseignement\2004-2005\Master
SIB\SIB2.2-bloc2-XML\Fichier CM XML\exemple-intro.dtd">
<livre id="561" nbpages="190" titre="La compagnie des spectres">
  <auteur>
    <nom>Salvayre</nom>
    <prenom>Lydie</prenom>
  </auteur>
  <format type="poche">
    <mesure type="largeur" unite="cm">11</mesure>
    <mesure type="longueur" unite="cm">19</mesure>
    <mesure type="hauteur" unite="mm">10</mesure>
  </format>
</livre>
```

# La DTD correspondante

```
<!ELEMENT livre (auteur, format)>
<!ATTLIST livre
  id CDATA #REQUIRED
  nbpages CDATA #REQUIRED
  titre CDATA #REQUIRED >
<!ELEMENT auteur (nom, prenom)>
<!ELEMENT format (mesure+)>
<!ATTLIST format
  type CDATA #REQUIRED >
<!ELEMENT mesure (#PCDATA)>
<!ATTLIST mesure
  type (hauteur | largeur | longueur) #REQUIRED
  unite (cm | mm | in) #REQUIRED >
<!ELEMENT nom (#PCDATA)>
<!ELEMENT prenom (#PCDATA)>
```

# Plan

- Documents XML
  - Syntaxe XML et documents bien formés
- Types de documents XML
  - DTD et documents valides
  - Introduction à XML-Schema
- Le monde XML
  - Quelques normes liés à XML
  - Quelques DTD importantes

# Plan

- **Documents XML**
  - Syntaxe XML et documents bien formés
- Types de documents XML
  - DTD et documents valides
  - Introduction à XML-Schema
- Le monde XML
  - Quelques normes liés à XML
  - Quelques DTD importantes

## Qu'y a-t-il dans un document XML ?

- o Prologue
  - Déclaration XML
  - Déclarations de DTD
    - o Instructions pour les processeurs XML
    - o Instructions de traitement
    - o Instructions pour applications externes
- o Arbre des éléments
  - Éléments
    - o Balises XML pour le marquage
    - o Contenu
      - texte
      - autres éléments
  - Attributs des éléments
    - o Information associées aux éléments
- o Commentaires

## Déclaration XML

- o Syntaxe générale  

```
<?xml version="1.0" [encoding = "encodage"] [standalone="yes|no »] ?>
```
- o Est une des informations de traitement
- o Indique
  - Conformité du document à une version de la norme XML
    - o version="1.0"
  - Jeu de caractères utilisé dans le document
    - o encoding = "UTF-8"
  - Présence ou non de références externes
    - o standalone="yes"

## Instructions de traitement

- o Informations nécessaires à une application externe
- o Format :
  - `<?NomApplication paramètres ?>`
- o Exemples
  - Déclaration de feuille de style à utiliser
    - o `<?xml-stylesheet href="fichier.xml" type="text/xsl"?>`
  - Déclaration XML de début de fichier
    - o `<?xml version='1.0' ?>`

## Éléments : règles de base

- o Un nom d'élément
  - commence par une lettre ou souligné
  - contient des lettres, chiffres, et "-", ".", ":", " \_"
  - peut posséder un nom de domaine
    - o domaine:nom\_element
    - o Ex. : xsl:template
- o Les noms d'éléments dépendent de la casse
  - `<nom_element> ≠ <nom_Element>`
- o Balises
  - de début : `<nom_element>`
  - de fin : `</nom_element>`
- o Les éléments peuvent être vides
  - pas de contenu
  - `<element_vide />`
  - Ex :
    - o ``

## Arbre des éléments

- o Un seul élément racine qui contient tous les autres
- o Pas d'intersections entre éléments
  - Mauvais : `<nom1><nom2>...</nom1></nom2>`
  - Bon : `<nom1><nom2>...</nom2></nom1>`
- o Blancs ou retours chariot en général non significatifs
  - `<section><p> ... </p></section>`
  - `<section><p> ... </p></section>`
- o Les éléments sont ordonnés

## Caractères spéciaux

- o Ces caractères ont une signification spéciale pour les outils XML
- o Il faut les écrire différemment
  - `<`      `&lt;`;
  - `>`      `&gt;`;
  - `&`      `&amp;`;
  - `'`      `&apos;`;
  - `"`      `&quot;`;

## Attributs associés aux élément : règles de base

- Dans les balises ouvrantes
  - `<el att1="valeur1" att2="valeur2">`
- Les noms d'attributs dépendent de la casse
  - `<el att1="valeur1" Att1="valeur2">`
- Valeurs d'attributs entourées
  - par des guillemets (") ou des apostrophes (')
- Les attributs sont non-ordonnés

## Attributs

- Les valeurs peuvent être
  - des données textuelles
    - `value="N'importe quoi"`
  - des *tokens* (noms XML) simples
    - `value = "blue"`
  - des ensemble *de tokens*
    - `value = "red green blue"`
- Possibilité d'énumérer les valeurs possibles et de mettre des valeurs par défaut (voir DTD)

## Attributs de type ID et IDREF(S)

- Permettent des relations non hiérarchiques entre éléments
  - ID : identificateur unique dans le document XML
  - IDREF : référence à un élément ayant un attribut de type ID
  - IDREFS : références à *des* éléments ayant un attribut de type ID
- Exemple

```
<société codes_services="A001 A003">
  <service code="A001">
    <employé code="E205" code_service="A001"> Jean Dupont </employé>
    <employé code="E206" code_service="A001"> Frédéric Marc </employé>
    <employé code="E207" code_service="A001"> Fabrice Detterne
  </employé>
  <employé code="H107" code_service="A003"> Angélique Millet
  </employé>
</service>
<service code="A003">
  <employé code="A115" code_service="A003"> Isabelle Mascot
  </employé>
</service>
</société>
```

## Commentaires

- Les commentaires ne sont pas considérés comme faisant partie du document XML.
  - `<!-- Un commentaire -->`
- Pas de '--' dans un commentaire !
- Un commentaire ne peut pas se trouver dans une autre déclaration

## Au bilan : dans un document XML

- Prologue
  - en-tête
  - déclaration de DTD (*pas encore vu*)
  - instructions de traitement
- Eléments
  - attributs
  - contenus
- Commentaires

## Plan

- Documents XML
  - Syntaxe XML et documents bien formés
- **Types de documents XML**
  - DTD et documents valides
  - Introduction à XML-Schema
- Le monde XML
  - Quelques normes liés à XML
  - Quelques DTD importantes

## Traiter automatiquement un document XML

- Parser
  - Outil qui lit un document XML et construit l'arbre des éléments en mémoire
- Vérifier qu'un document répond bien à la syntaxe XML
  - Document *bien formé*
  - Possibilité de l'utiliser en tant que tel
    - ex. : le présenter à l'utilisateur
- Vérifier en plus qu'un document suit bien la grammaire définie dans une DTD
  - Document *valide*

## Document Type Definition

- Définir le type de document XML voulu
  - décrire comment construire un document XML qui lui corresponde (grammaire)
- Permet de
  - valider un document XML (parser validant)
    - vérifier que tous les éléments sont présents et corrects
    - vérifier que les noms d'attributs et leurs valeurs sont corrects
  - transmettre cette connaissance à d'autres
    - ils pourront définir leurs propres documents XML dans le même cadre
    - d'où possibilité de standardisation et d'échanges

## DTD

- Une déclaration
  - contenant la définition formelle de la structure autorisée,
  - qui décrit donc
    - quels noms sont utilisés pour les types d'éléments
    - comment ces types d'éléments s'organisent
      - ordre
      - hiérarchie
    - les attributs des éléments
    - des entités analysables ou non
    - des notations pour les types de données binaires
- Liaison DTD / document XML
  - Soit la DTD est dans le document XML (inline)
  - Soit le document XML réfère à la DTD avec une URI (la DTD est dans un fichier externe)

## Déclarations

- Instructions pour le processeur XML
- Format : `<! ... >` ou `<! ... [<! ... >]>`
  - **Document type** - `<!DOCTYPE ... >`
  - **Character data** - `<![CDATA[ ... ]>`
  - **Entities** - `<!ENTITY ... >`
  - **Notation** - `<!NOTATION ... >`
  - **Element** - `<!ELEMENT ... >`
  - **Attributes** - `<!ATTLIST ... >`
  - `<![INCLUDE [...]]>` et `<![IGNORE [...]]>`

## Déclaration *Document Type*

- Identifie le nom de l'élément racine du document
  - `<!DOCTYPE My_XML_Doc>`
- Permet aussi de rajouter des définitions d'entités et des DTD
  - `<!DOCTYPE My_XML_Doc [ ... ] >`  
`<My_XML_Doc>`  
...  
`</My_XML_Doc>`

déclaration de la DTD

## Déclaration *Character Data*

- Dans les occasions pour lesquelles le texte doit contenir des caractères qui ne doivent pas être interprétés
- Deux textes équivalents
  - `Press &lt;&lt;&lt;ENTER&gt;&gt;&gt;`
  - `<![CDATA[Press <<<ENTER>>>]]>`

## DTD et document XML

Valide, contraint

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE livre SYSTEM "exemple-intro.dtd">
<livre id="561" nbpages="190"
  titre="La compagnie des spectres">
  <auteur>
  <nom>Salvayre</nom>
  <prenom>Lydie</prenom>
  </auteur>
  <format type="poche">
  <mesure type="largeur" unite="cm">11</mesure>
  <mesure type="longueur" unite="cm">19</mesure>
  <mesure type="hauteur" unite="mm">10</mesure>
  </format>
</livre>
```

```
<!ELEMENT livre (auteur, format)>
<!ATTLIST livre
  id CDATA #REQUIRED
  nbpages CDATA #REQUIRED
  titre CDATA #REQUIRED >
<!ELEMENT auteur (nom, prenom)>
<!ELEMENT format (mesure+)>
<!ATTLIST format
  type CDATA #REQUIRED >
<!ELEMENT mesure (#PCDATA)>
<!ATTLIST mesure
  type (hauteur | largeur |
  longueur) #REQUIRED
  unite (cm | mm | in)
  #REQUIRED >
<!ELEMENT nom (#PCDATA)>
<!ELEMENT prenom (#PCDATA)>
```

Document XML

Est validé par,  
Est contraint par

## Mettre en place une DTD

- Utiliser les composants de XML...
  - Entités, éléments, déclarations, instructions de traitements, listes d'attributs, etc.
- ... dans des DTD pour spécifier les règles permettant de valider des documents XML
  - Définir un modèle (type) de document de façon formelle
- Une DTD décrit
  - Quels noms peuvent être utilisés pour les types d'éléments
  - L'ordre dans lesquels ceux-ci peuvent apparaître
  - La hiérarchie documentaire
  - Les noms et les types des attributs d'éléments

## Déclaration de DTD

- La DTD est stockée
  - soit dans le fichier XML
  - soit dans un fichier extérieur
  - soit dans les deux
- Une DTD interne peut écraser ou ajouter des ENTITY ou des ATTLIST à des définitions de DTD externes
- Une DTD est composée de déclarations
  - ELEMENT – définitions d'éléments
  - ATTLIST – définitions d'attributs
  - ENTITY – définitions d'entités
  - NOTATION – définitions de notations

## Déclarations d'éléments

- Définir un élément et son contenu
  - `<!ELEMENT name (#PCDATA)>`  
⇒ `<name> ...du texte... </name>`
- Un élément vide n'a pas de contenu
  - `<!ELEMENT name EMPTY>`  
⇒ `<name/>`
- Si on autorise les fils
  - Quelconques :  
`<!ELEMENT name ANY>`
  - Spécifiés :  
`<!ELEMENT person (name, e-mail*)>`

## Spécification des fils (grammaire)

- Définir le contenu des éléments  
`<!ELEMENT person (name, e-mail*)>`
- ...et définir une hiérarchie d'éléments  
`<!ELEMENT name (fname, surname)>`  
`<!ELEMENT fname (#PCDATA)>`  
`<!ELEMENT surname (#PCDATA)>`  
`<!ELEMENT e-mail (#PCDATA)>`
- Organisation des sous-éléments
  - Connecteur de séquence ', ' : (A, B, C) [puis]
  - Connecteur de choix '| ' : (A | B | C) [ou]

## Indicateurs de quantité

- Contraintes sur les éléments des DTD
 

<b>A?</b>	Possible	[0..1]
<b>A+</b>	1 fois et plus	[1..*]
<b>A*</b>	0 ou plus	[0..*]
- Exemples  
(A, B)+  
(A, B?) | C+\*

## Déclaration d'attributs

- Les attributs sont associés aux types d'éléments
- Déclarés dans une déclaration ATTLIST
  - `<!ELEMENT element ... >`
  - `<!ATTLIST element ... >`
  - Il faut ensuite définir
    - le nom de l'attribut
    - le type de l'attribut
    - sa valeur par défaut

## Noms et types d'attributs

- Noms d'attributs
  - `<!ATTLIST elem name type default>`
  - `<!ATTLIST elem first_attr ... second_attr ... third_attr ... >`
- Types d'attributs
 

CDATA	ID
NMTOKEN	IDREF
NMTOKENS	IDREFS
ENTITY	NOTATION
ENTITIES	name group

## Types d'attributs (1)

- CDATA
    - Chaîne de caractères
    - `<!ATTLIST person name CDATA ... >`
    - `<person name = "Tom Jones">`
  - NMTOKEN
    - Token unique
    - `<!ATTLIST mug color NMTOKEN ... >`
    - `<mug color="red">`
  - NMTOKENS
    - Multiples tokens
    - `<!ATTLIST temp values NMTOKENS ... >`
    - `<temp values="12 15 34">`
- Joue sur la manière dont le parser interprète l'attribut

## Types d'attributs (2)

- ENTITY
  - L'attribut est une référence d'entité
  - `<!ATTLIST person photo ENTITY ... >`
  - `<person photo="MyPic">`
- ENTITIES
  - Plusieurs références d'entités
  - `<!ATTLIST album photos ENTITIES ... >`
  - `<album photos="pic1 pic2">`
- ID
  - Identificateur unique
  - `<!ATTLIST person id ID ... >`
  - `<person id = "P09567">`
- IDREF
  - Référence à un ID
  - `<!ATTLIST person father IDREF ... >`
  - `<person father="P09567">`

## Types d'attributs (3)

- IDREFS
  - Référence à plusieurs ID
  - `<!ATTLIST person children IDREFS ... >`
  - `<person children="P09567 P09677">`
- Name group
  - Liste restreinte de valeurs possibles
  - `<!ATTLIST point coord (X|Y|Z) ... >`
  - `<point coord="X">`
- NOTATION
  - Décrit des données non XML
  - `<!NOTATION jpg SYSTEM "JPEG Image" >`
  - `<!NOTATION gif PUBLIC "-//ISBN 0-7923-9432-1::Graphic Notation//NOTATION CompuServer Graphic Interchange Format//EN">`
  - `<!ATTLIST image format NOTATION (jpg|gif) ... >`
  - `<image format="gif">`

## Valeurs d'attributs par défaut

- Quatre types
  - #REQUIRED** doit être spécifié
  - #IMPLIED** peut être spécifié
  - "default"** valeur par défaut si non spécifié
  - #FIXED** une seule valeur autorisée

<code>&lt;ATTLIST tag</code>	name	type	default>
<code>&lt;!ATTLIST citoyen</code>	parents	IDREFS	<b>#REQUIRED</b>
	id	ID	<b>#IMPLIED</b>
	sex	(m f)	"f"
	adress	CDATA	<b>#IMPLIED</b>
	nat	CDATA	<b>#FIXED "Fr"&gt;</b>

## Définition interne de DTD

### o Dans la déclaration DOCTYPE

```
<?xml version="1.0" standalone="yes" ?>
<!DOCTYPE racine [
  <!-- ici la DTD -->
  <! ... >
  <! ... >
]>
<!-- ici le reste du fichier XML -->
<racine>
...
</racine>
```

## Définition externe privée de DTD

### o Référence à la DTD externe par un chemin dans la déclaration DOCTYPE

```
<?xml version="1.0" standalone="no" ?>
<!DOCTYPE racine
  SYSTEM "dossiers/MyDoc.dtd" [
  <!-- déclarations supplémentaires -->
  <! ... >
  <! ... >
]>
<!-- ici le reste du fichier XML -->
```

DTD externe privée

- Les déclarations spécifiques au document restent définies de façon interne

## Définition externe privée de DTD

### o Référence à la DTD externe par une URL dans la déclaration DOCTYPE

```
<?xml version="1.0" standalone="no" ?>
<!DOCTYPE racine
  SYSTEM "http://.../MyDoc.dtd" [
  <!-- déclarations supplémentaires -->
  <! ... >
  <! ... >
]>
<!-- ici le reste du fichier XML -->
```

DTD externe privée

- Les déclarations spécifiques au document restent définies de façon interne

## Définition externe publique de DTD

### o Utilisation du mot-clé PUBLIC

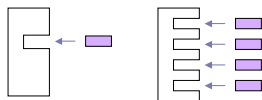
```
<!DOCTYPE racine
  PUBLIC "identifiant public" "url" ?>
```

### o Exemple

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN"
  "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">
<html>
.....
</html>
```

## Les entités

- o Sont des alias associant un nom à des « unités d'information »
- o Les entités spécifiques au document sont décrites dans sa DTD interne
- o Les entités plus générales sont décrites dans des DTD externes
- o Chaque entité
  - o est identifiée par un nom
  - o est définie par une déclaration d'entité
  - o est utilisée en appelant une référence d'entité



## Utilisation des entités

### o Quand l'information

- Est utilisée dans plusieurs endroits
  - o Ex. déclaration légale, caractère spécial
- Est une partie d'un document qui doit être tronçonné pour rester gérable
  - o Ex. livre : 1 fichier + n chapitres : n fichiers
- Est conforme à un format de donnée différent de XML
  - o Ex. : image JPEG

## Types d'entités

- Entités internes
  - générales
    - utilisées dans les documents XML
  - paramètre
    - utilisées dans les déclarations dans les DTD
- Entités externes
  - générales
  - paramètres
- Entités analysables
- Entités non analysables
- Entités caractères
  - déjà vues

## Entités générales internes

- ```
<!ENTITY nom "chaîne de remplacement">
```
- Entités analysables utilisées uniquement dans le document
  - Référence : `&nom_entité;`
  - Exemple
    - Déclaration dans la DTD

```
<!ENTITY PCI "Permis de conduire informatique">
```
    - Utilisation

```
<p>Le cours du PCI (&PCI;) se compose de...</p>
```

## Entités générales externes

- ```
<!ENTITY nom SYSTEM "URI" >
```
- Permet de construire un document XML à partir de plusieurs autres documents
  - Référence : `&nom_entité;`
  - Exemple
    - Déclaration dans la DTD

```
<!ENTITY doc SYSTEM "http://toto.org/doc.xml" >
```
    - Utilisation

```
<aide> &doc; </aide>
```

## Entités paramètres internes

- ```
<!ENTITY % nom "caractères de remplacement" >
```
- Entités analysables uniquement utilisées dans les DTD
  - Référence dans la DTD : `(%nom_entité;)` (parenthèses conseillées)
  - Exemples
    - Déclarations DTD

```
<ENTITY % tout "ANY" >
<ENTITY % common "(paralist|table)">
```
    - Utilisations dans la DTD

```
<ELEMENT paragraphe %tout; >
<ELEMENT chapter ((%common;)*, section*)>
<ELEMENT section (%common;)*>
```

## Entités paramètres externes

- ```
<ENTITY % nom SYSTEM "URI" >
```
- Pour construire une DTD complexe à partir d'autres DTD complémentaires
  - Référence dans la DTD : `%nom_entité;`
  - Exemple
    - Déclaration dans la DTD

```
<ENTITY % règles SYSTEM "http://toto.org/regles.dtd" >
```
    - Utilisation dans la DTD

```
%règles;
```

## Entités analysables

- Le texte de remplacement fait partie intégrale du document
  - les données sont analysées correctement par le parser XML
- Déclaration dans la DTD comme ENTITY
- Utilisation avec `&nom;` ou `%nom;`



## Entités non analysables

```
<!ENTITY % nom SYSTEM "URI" NDATA notation >
```

- Pour déclarer un contenu non XML dans un document XML
  - fichier image, audio, etc.
- Référence : `&nom_entité;` uniquement comme attribut de type ENTITY
- Exemple
  - Déclaration DTD

```
<!NOTATION TIFF SYSTEM "format TIFF">
<!ENTITY photo SYSTEM "photo.tif" NDATA TIFF>
<!ELEMENT pic EMPTY>
<!ATTLIST pic name ENTITY #REQUIRED>
```
  - Utilisation dans le document XML

```
<pic name="photo" />
```

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

49

## Déclarations de notations

```
<!NOTATION nom SYSTEM "URI" >
```

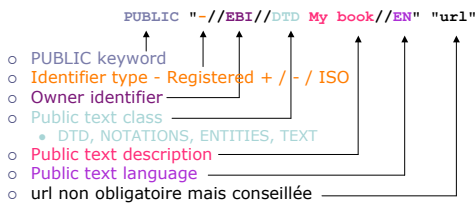
- Pour
  - Identifier par un nom le format des entités non XML externes
  - Définir les formats des données et les applications qui permettent de les traiter
- Exemple

```
<!NOTATION GIF SYSTEM "GIF" >
<!NOTATION GIF89a PUBLIC "-//Compuserve//NOTATION
Graphics Interchange format 89a//EN" >
```

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

50

## Identificateurs publics



Exemples d'utilisation  

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN"
"http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">
<!NOTATION GIF89a PUBLIC "-//Compuserve//NOTATION
Graphics Interchange format 89a//EN" >
```

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

51

## Construire une DTD

- Non trivial : il faut éviter de se tromper
  - Changer une DTD XML a des conséquences sur les documents qui la suivent
- Ressemble à la création d'un schéma de base de données
- Il faut considérer
  - Le problème de la granularité
  - La questions des attributs et des éléments
  - Les limitations inhérentes aux DTD

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

52

## Identifier les données qui nécessitent d'être balisées

- Pour chaque unité d'information, déterminer
  - Peut-on lui donner un nom ?
  - Apparaît-elle tout le temps ?
  - Peut-il y en avoir plusieurs ?
  - Peut-on la décomposer en des unités plus petites ?
  - Y-a-t'il du contenu textuel qui ne change pas ?
  - Comment est-elle associée aux autres unités ?

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

53

## Granularité

- **<PERSON>**  
**<NAME>Jon Smith</NAME>**  
**</PERSON>**
- **<PERSON>**  
**<FORENAME>Jon</FORENAME>**  
**<SURNAME>Smi th</SURNAME>**  
**</PERSON>**

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

54

## Eléments ou attributs ?

- Comment les données doivent-elles être encapsulées ?
  - `<book>`  
  `<title>The Forty-nine Steps</title>`  
  `...`  
  `</book>`
  - `<book title="The Forty-nine Steps">`  
  `...`  
  `</book>`
- Tout dépend de ce que l'on veut faire...
- Il existe des avis tranchés...

## Eléments ou attributs ? (2)

- Séparer le contenu des métadonnées
  - Données qui doivent être imprimées comme du texte comme contenu
  - Métadonnées comme attributs
- Règles générales
  - Si on enlève toutes les balises, le document doit encore être lisible et utilisable
  - S'il y a doute, utiliser un attribut

## Limites des DTD / XML

- XML est seulement une syntaxe
- XML ne porte pas de sémantique
- Uniquement description de structure
- Pas de types
  
- Un des moyens de pallier certains problèmes
  - XML-schema

## DTD

- Une syntaxe de description non-XML, héritée de SGML
  - Oblige à apprendre un langage supplémentaire
  - Ne permet pas de manipuler les DTD avec des outils XML
- Pas assez de contraintes sur les données manipulées
  - Toute données est une chaîne de caractères
  - Impossible de
    - spécifier des types simples
      - ex. entiers, dates, etc.
    - spécifier des cardinalités simples
      - ex. « un ARTICLE aura entre 1 et 4 MOTS-CLE »
    - spécifier des contraintes simples
      - ex. entier positif

## XML-Schema

- Autre manière de spécifier des types de documents XML
- Le schéma est exprimé en XML
- Possibilité de spécifier plus de contraintes sur les données
- Possibilités avancées d'extension des schémas
- On élargit l'approche de gestion documentaire à celle plus générale de gestion de données

## Exercice 1 : Proposez plusieurs documents XML valides suivants la DTD suivante

```
<!ENTITY % opt_fields "year?, volume?, pages?, month?, url?, abstract?, note?">
<!ENTITY % req_fields "author, title">
<!ENTITY % key_atts "key ID #REQUIRED">

<!ELEMENT bibtex-file
(article|inproceedings|book|techreport|phdthesis|unpublished|misc)>

<!ELEMENT key (#PCDATA)>
<!ELEMENT author (name+)>
<!ELEMENT name (#PCDATA)>
<!ELEMENT title (#PCDATA)>
<!ELEMENT booktitle (short, long?)>
<!ELEMENT short (#PCDATA)>
<!ELEMENT long (#PCDATA)>
```

```
(suite)
<!ELEMENT year (#PCDATA)>
<!ELEMENT volume (#PCDATA)>
<!ELEMENT pages EMPTY>
<!ATTLIST pages
  first NMTOKEN #REQUIRED
  last NMTOKEN #REQUIRED>
<!ELEMENT month (#PCDATA)>
<!ATTLIST month
  mtype (short | long) "short">
<!ELEMENT url EMPTY>
<!ATTLIST url
  ftype (ps|pdf|html) #REQUIRED
  href CDATA #REQUIRED>
```

```

(suite)
<ELEMENT abstract (#PCDATA)>
<ELEMENT number (#PCDATA)>
<ELEMENT note (#PCDATA)>
<ELEMENT publisher (#PCDATA)>
<ELEMENT institution (#PCDATA)>
<ELEMENT school (#PCDATA)>
<ELEMENT howpublished (#PCDATA)>
<ELEMENT address (#PCDATA)>
<ELEMENT inproceedings (%req_fields;, booktitle, %opt_fields;)>
<ATTLIST inproceedings %key_atts;>
<ELEMENT article (%req_fields;, journal, %opt_fields;)>
<ATTLIST article %key_atts;>
<ELEMENT book (%req_fields;, publisher, %opt_fields;)>
<ATTLIST book %key_atts;>
<ELEMENT techreport (%req_fields;, institution?, number?, %opt_fields;)>
<ATTLIST techreport %key_atts;>
<ELEMENT phdthesis (%req_fields;, school?, %opt_fields;)>
<ATTLIST phdthesis %key_atts;>
<ELEMENT unpublished (%req_fields;, %opt_fields;)>
<ATTLIST unpublished %key_atts;>
<ELEMENT misc (%req_fields;, howpublished, %opt_fields;)>
<ATTLIST misc %key_atts;>

```

## Exercice 2

- Proposez deux DTD permettant de valider le document XML suivant

```

<examen code="coursXML">
<titre>Outils et documents XML</titre>
<date mois="jan" annee="2006"/>
<questions>
<question> <partie>Première partie<partie>Une sous-partie</partie>
/partie</question>
<question> <partie>Deuxième partie<partie> </question>
<question> <partie/> </question>
</questions>
</examen>

```

## Exercice 3

- On veut représenter en XML des données concernant les étudiants en SIB et leurs enseignements
- Proposez une DTD

## Plan

- Documents XML
  - Syntaxe XML et documents bien formés
- Types de documents XML
  - DTD et documents valides
  - Introduction à XML-Schema
- **Le monde XML**
  - Quelques normes liés à XML
  - Quelques DTD importantes

## Standardisation

- XML permet de définir des DTD
  - modèles de documents
  - modèles de représentation de données
- Dès qu'on a un groupe, partage de données/documents
  - nécessité de partager les manières de décrire
  - accord
    - local
    - global → standardisation
- Des standards sous la forme de DTD (ou de schémas),
  - Stricts
  - Qui peuvent être raffinés
    - Les spécialiser avec des DTD internes
    - N'en utiliser que des parties

## Avantages et applications XML

- Avantages
  - Réutilisabilité, partage
  - Pérennité
  - Intégrité
  - Portabilité
- Applications
  - Documents
  - Echange de données
  - Bureautique
  - Web
  - BDD semi-structurées
  - Commerce électronique
  - ...

## Quelques standard XML

- The XML Bookmark Exchange Language (XBEL)
- Open eBook Publication Structure
- SportsML
- NewsML
- XML Book Industry Transaction Standards (XBITS)
- OpenDocument...
- DocBook
- ebXML (electronic Business)
- Universal Description, Discovery & Integration (UDDI)
- Text Encoding and Interchange (TEI)
- XTM (XML Topic Maps)
- ...

<http://publishing.xml.org/standards/>  
<http://www.oasis-open.org/specs/index.php>

## Quelques spécifications XML (W3C)

- XML Schema
- XLink et XPointer
- XPath
- XSL et XSLT
- XML Query
- Namespaces
- SAX
- DOM
- MathML
- OWL
- RDF
- SMIL
- SOAP
- SVG
- XHTML

Voir <http://www.w3c.org/>

## XPath

- Standard permettant d'identifier et de spécifier toutes données dans un document XML
- Exemples
  - `//toto[@name]`  
→ Tous les élément toto qui ont un attribut name
  - `//tata/descendant::*`  
→ Tous les descendants des éléments tata
- Voir cours Xpath

## XLink

- Objectif
  - Donner la possibilité de liens riches
- XLink
  - XML Linking Specification
  - Liens
    - simples (1:1) et étendus (n:n),
    - typés
    - internes ou externes
- Exemple

```
<site xlink:type="local"
  xlink:href="http://www.xml.fr/FAQ.xml"
  xlink:label="fr"
  xlink:title="Version française"/>
```

## XPointer

- Objectif
  - Pointer précisément dans un document XML
- XPointer
  - XML Extended Pointer Specification
  - Une référence absolue (le document XML)
  - et une référence relative (à l'intérieur du document)
    - expression XPATH
- Exemples
  - `http://www.toto.org/xml/doc.xml#xptr\(/intro/title\)`  
élément title dans élément intro de doc.xml
  - `http://site.fr/page.xml||id\('ref12"\).child\(1,session\)`  
premier élément session enfant de l'élément identifié par ref12, dans le document page.xml

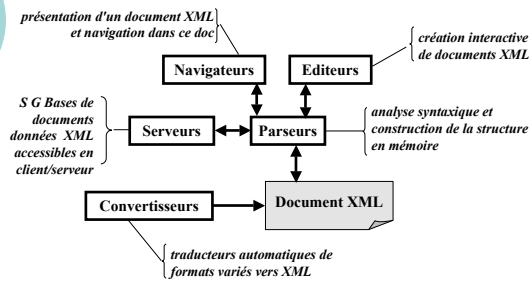
## XSL

- Ensemble d'outils permettant de
  - Transformer les documents XML en d'autres documents XML (changement de format)
  - Visualiser les documents XML sous forme lisible, pour de multiples supports
- eXtensible Style Language
  - XSL-FO (« XSL Formatting Objects »)
    - présenter des informations
  - XSLT (XSL transformation)
    - transformer un arbre XML en un autre arbre XML
- Voir cours XSL-XSLT





## Différents types d'outils XML



## Exemples d'outils

- Parseurs
  - SAX et DOM souvent intégrés directement dans les langages (Java, .NET, etc.)
- Editeurs
  - XML-Spy
  - Cooktop
  - XMetal...
- Navigateurs
  - Firefox, IE, etc.
- Convertisseurs
  - Nombreux outils avec format de sortie textuel
- SGB données/documents XML
  - Évolutions des SGBD classiques
  - SGBD dédiés

## Autres outils

- XHTML / CSS
  - Dreamweaver, NVU...
- XSL
  - Style-vision...
- RDF
  - Outils du web sémantique...
- SMIL
  - Player : REAL ...
- SVG
  - Adobe...

## DocBook

- Objectif
  - Codage de texte de documentation
  - Sorties multi-formats
- <http://www.oasis-open.org/docbook/>
- A voir en TP
  - sdocbook : sous-ensemble de balises

## Open Document

- Objectif
  - codage des documents de suites bureautiques
    - textes, feuilles de calcul, présentations, dessins, formules, bases de données...
    - modèles pour ces documents
- Mis en place par OASIS
  - à partir des formats de OpenOffice (SUN)
- Version 1.0
  - mai 2005
- Utilisé dans Open Office 2.0
- Standard ouvert
  - intérêt pour le partage, la récupération, etc.
  - enjeu politique (Massachusetts, sept 2005)
    - réaction de Microsoft ?

## Conclusion

- XML
  - Norme sortie en 1998
  - Unicode / généricité
  - Documents / données
  - Mondialement adoptée
- Standards et normes
  - Variés : dans tous les domaines nécessitant
    - Pérennité
    - Echange
  - Plus ou moins adaptés et adoptés
  - Questions récurrentes
    - Evolution
    - Interopérabilité

## Annexes : EAD/EAC - TEI

- Pour info
  - quelques transparents sur ces normes utiles aux sciences de la documentation

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

91

## EAD / EAC

- Objectif
  - Normalisation des instruments de recherche et des descriptions de contextes pour les archives
- EAD
  - Encoded Archival Description
- EAC
  - Encoded Archival Context
- Un point d'entrée
  - <http://www.archivesdefrance.culture.gouv.fr/fr/archivistique/DAFlangage.html>
- Il y a des spécialistes à l'ENSSIB

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

92

## Les trois éléments principaux de l'instrument de recherche en XML/EAD

- Sous l'élément racine <ead> :

<eadheader>	en-tête EAD (description bibliographique de l'IR) (obligatoire)
<frontmatter>	préliminaires (page de titre, introduction, préface...)
<archdesc>	description archivistique (obligatoire)

Ce transparent et les suivants sur l'EAD sont extraits de « L'EAD, une DTD pour la rédaction, l'archivage et la diffusion des instruments de recherche archivistiques » F. Clavaud / 5 avril 2004

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

93

## En-tête EAD <eadheader>

```
<eadheader>
<eadid> Identifiant EAD (du fichier électronique)
<filedesc> Description du fichier
  <titleproper> Titre propre de l'instrument de recherche
  <subtleproper> Sous-titre de l'instrument de recherche
  <author> Auteur de l'instrument de recherche
  <edition> Mention d'édition
  <publicationstmt> Mention de publication
  <seriesstmt> Mention de collection
  <notestmt> Mention de note
<profiledesc> Description du profil
  <creation> Informations sur la création de l'inventaire
  <date> Date de l'inventaire
  <langusage> Langue utilisée
  <desrules> Règles de description (archivistique utilisée)
  <revisiondesc> Descriptions des révisions ; permet de gérer les versions successives de l'instrument de recherche encodé
</eadheader>
```

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

94

## En-tête EAD <eadheader>

```
<eadheader audience="external" findaidstatus="edited-partial-draft" encodinganalog="DC">
  <eadid encodinganalog="Identifier" FRDAFANCH09AP_000000001 </eadid>
  <filedesc>
    <titleproper encodinganalog="title">Etat sommaire des fonds d'archives privées du Centre historique des Archives nationales</titleproper>
    <subtleproper> Séries 317AP à 427AP et ABXXX368 à ABXXX325</subtleproper>
    <author encodinganalog="creator">Instrument de recherche rédigé par Claire Sibille, [...] encodé en XML conformément à la DTD EAD puis converti en HTML par Martin Sévigny (société AJLSM) et Florence Clavaud (service des nouvelles technologies du CHAN)</author>
    <edition>
      <editionstmt> Quatrième édition</editionstmt>
      <publicationstmt>
        <publisher encodinganalog="publisher">Centre historique des Archives nationales de France (CHAN)</publisher>
        <address>
          <addressline>69 rue des Francs-Bourgeois</addressline>
          <addressline>F-75141 PARIS CEDEX 03</addressline>
        </address>
        <date encodinganalog="date">aout 2003</date>
      </publicationstmt>
    </filedesc>
    <profiledesc>
      <creation>Instrument de recherche produit sous Word et converti en XML[...] en <date> novembre 2001</date>. Encodage retu et complété par [...]</creation>
      <langusage>Instrument de recherche rédigé en <language langcode="fr">français</language>. </langusage>
      <revisiondesc>
        <change>
          <date>8 septembre 2003</date>
          <item>Fichier converti en EAD 2002 par programme XSL-T fourni par D. Pihl (SAA) adapté</item>
        </change>
      </revisiondesc>
    </eadheader>
```

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

95

## Préliminaires <frontmatter> : exemple

```
<frontmatter>
  <titlepage>
    <publisher>Centre historique des Archives nationales</publisher>
    <titleproper>Etat sommaire</titleproper> des fonds d'archives privées</titleproper>
    <subtleproper>Séries AP (317AP à 427AP) et AB XX (ABXXX368 à ABXXX325)</subtleproper>
    <author>par Claire SIBILLE</author> conservateur du Patrimoine</author>
    <author>avec la collaboration de Ferry AUDOUS et de Violaine LE NENAON</author>
    <publisher>Direction des Archives de France</publisher>
    <date>aout 2003</date>
  </titlepage>
  <div>
    <head>Introduction</head>
    <p>Paru en 1973, l'état sommaire des fonds de la série AP (archives de personnes et de familles) rédigé par Suzanne d'Huart et Chantal Bonazzi traitait des fonds privés 1AP à 315AP. Depuis, plus de 300 nouveaux fonds d'archives privées sont entrés aux Archives nationales. [...] Ces notices sont également consultables dans la base <extrref href="http://sdx.archivesdefrance.culture.gouv.fr/ap">BORA</extrref></p>
  </div>
  <div>
    <head>Mode d'emploi</head>
  </div>
  <div>
  </div>
</frontmatter>
```

CM2-3-4 : eXtensible Markup Language – Yannick Prié  
UE2.2 – Master SIB M1 – 2005-2006 : Représentation des données et des connaissances

96







## TEI : codage

---

- Construction de la DTD de façon modulaire :
  - jeu de balises «noyau» (*core tag set*) composé d'éléments communs à tous les types de textes (divisions, paragraphes, etc.)
  - des ensembles de balises de base (*base tag sets*) pour chaque type particulier de texte (prose, poésie en vers, etc.)
  - des jeux de balises additionnelles (*additional tag sets*) pour des mécanismes particuliers qui peuvent se superposer à n'importe quel type de texte (liens hypertextuels, etc.).
- Noyau obligatoire, autre éléments facultatifs
- Importance de l'en-tête
  - codage systématique des méta-données de n'importe quel document électronique



## TEI : où en est on ?

---

- Au départ : SGML, puis XML
- Utilisation dans de très nombreux projets
- Reste la question du codage
  - Pour qui code-t'on ?
  - Un chercheur peut-il réutiliser le codage d'un autre ? Ajouter au codage d'un autre ?
  - Possibilité de mixer avec d'autres formats ? (particulièrement docbook)
  - etc.